

Universidad de Valencia Intelligent Data Analysis Laboratory (IDAL) Escuela Superior de Ingeniería

Máster Propio en Inteligencia Artificial Avanzada y Aplicada (IA3)

Comparativa de arquitecturas post-transformers para la traducción automática de texto

Trabajo fin de estudio presentado por:

Tipo de trabajo:

Director/a:

Tutor/a:

Fecha:

Ramírez Herrera, Felipe

Comparativa de Soluciones

Soria Olivas, Emilio

Manuel Valero Laparra

29 de abril de 2025

Tabla de contenido

1	Intro	oduction	12
	1.1	Motivación	12
	1.2	Planteamiento del problema	13
	1.3	Estructura de la memoria	13
2	Cont	texto y estado del arte	13
	2.1	Traducción automática	13
	2.2	Transducción	17
	2.3	Arquitecturas y modelos	18
	2.4	Atención	20
	2.5	Transformers	22
	2.6	Universal Transformers	29
	2.7	Reformer	29
	2.8	Evolved Transformer	31
	2.9	Mixture of Experts (MoE)	32
	2.9.1	1 ST-MoE	33
	2.9.2	2 Switch Transformers	34
	2.10	Space-State Models	36
	2.10.	.1 Mamba	38
	2.10.	.2 Simplified Structured State Space Sequence (S5)	39
	2.10.	.3 Gated State Spaces	40
	2.10.	.4 Jamba	42
	2.10.	.5 MoE Mamba	43
	2.11	Moving Average Equipped Gated Attention	43

2.12	Extended Long Short-Term Memory	46
3 Obje	tivos y metodología de trabajo	47
3.1	Objetivo general	47
3.2	Objetivos específicos	47
3.3	Metodología de trabajo	47
3.3.1	Tipo de investigación	47
3.3.2	Fuentes	48
3.4	Estudio preliminar del estado del arte	50
3.4.1	Pasos realizados para la revisión preliminar de la literatura	50
3.5	Métodos de investigación	51
3.5.1	Experimentos Controlados	52
3.5.2	Estudios de Casos	52
3.5.3	Análisis Comparativo	52
3.5.4	Revisiones Sistemáticas	53
3.5.5	Validación Cruzada	53
4 Plan	teamiento de la comparativa	53
4.1	Criterios de selección de los conjuntos de datos	53
4.1.1	Representatividad de los Datos	53
4.1.2	Diversidad Lingüística	54
4.1.3	Variedad en Estilos y Dominios	54
4.1.4	Calidad de los Datos de Entrenamiento y Pruebas	54
4.1.5	Tamaño del Conjunto de Datos	54
4.1.6	Datos Equilibrados y Balanceados	55
4.1.7	Datos Actualizados	55

4.1.8	Conjunto de Pruebas Adecuado	55
4.1.9	Problemas de sesgo y justicia	55
4.1.1	0 Tamaño y Complejidad del Texto	55
4.2	Corpus paralelo	56
4.2.1	UN Parallel Corpus V1	56
4.2.2	OpenSubtitles v2018	57
4.2.3	Ventajas y desventajas de los conjuntos de datos seleccionados	58
5 Desa	arrollo de la comparativa	61
5.1	Arquitecturas seleccionadas	61
5.1.1	Transformer Vainilla	61
5.1.2	Universal Transformer	61
5.1.3	ST-MoE	62
5.1.4	Mamba	62
5.1.5	Gated State Spaces	62
5.1.6	Moving Average Equipped Gated Attention	62
5.2	Implementación y Experimentación	63
5.3	Métricas de Evaluación	64
5.4	Análisis Comparativo	64
5.5	Configuración del entorno	64
5.5.1	Habilitar compatibilidad con GPU	66
5.5.2	Línea base	67
5.6	Aspectos específicos del dominio	68
5.7	Casos de prueba y experimentos	70
5.8	Hiperparámetros dependientes de los datos	70

	5.8.1	Longitud Mínima de Secuencia7	′1
	5.8.2	Longitud Máxima de Secuencia7	'1
	5.8.3	Tamaño del Vocabulario7	'1
	5.8.4	Batch Size (Tamaño del Lote)7	'2
	5.9	Convenciones	'3
	5.9.1	Optimizador7	'3
	5.9.2	Scheduler7	'3
	5.9.3	Early Stopping7	'4
	5.9.4	Clipping7	'4
	5.9.5	Medida de Pérdida7	'5
6	Disc	usión y análisis de resultados7	′5
	6.1	Cantidad de Parámetros	'5
	6.2	Tiempo de entrenamiento y validación	'7
	6.2.1	Tiempo de entrenamiento7	'9
	6.2.2	Tiempo de validación8	31
	6.3	Consumo de memoria	34
	6.4	Pérdida8	35
	6.4.1	Fase de entrenamiento8	36
	6.4.2	Pérdida durante la fase de validación9	12
	6.5	Exactitud9	8
	6.5.1	Conjunto de datos UNPC9	19
	6.5.2	Conjunto de datos OSPC10)5
	6.5.3	Perplejidad11	2
	c c	Inferencia11	1

	6.0	6.1	BLEU	115
	6.0	6.2	Tiempo de inferencia	116
7	7 Co	onc	usiones y trabajo futuro	122
	7.1		Conclusiones	122
	7.2		Trabajo futuro	125
	7.2	2.1	Combinación de Mixture of Experts (MoE) con Gated State Spaces (GSS)	126
	7.2	2.2	Incorporación de Atención Promediada en Modelos Transformer	127
	7.2	2.3	Integración de Extended LSTM con Arquitecturas Transformer	127
	7.2	2.4	Exploración de Modelos Auto-regresivos Híbridos	128
	7.2	2.5	Hibridación de Modelos con MoE y Transformers Evolucionados	128
	7.2	2.6	Implementación de Técnicas de Compresión de Modelos	129
8	B Re	efer	encias Bibliograficas	130
g) Aı	nex	os	135

Índice de cuadros

Cuadro No. 4.1. Características del conjunto de datos UN Parallel Corpus V1 (UNPC)57
Cuadro No. 4.2. Características del conjunto de datos OpenSubtitles (OSPC)57
Cuadro No. 5.1. Principales innovaciones de los modelos seleccionados para la comparativa63
Cuadro No. 5.2. Configuración del entorno de prueba65
Cuadro No. 5.3. Casos de prueba70
Cuadro No. 5.5. Hiperparámetros dependientes de los conjuntos de datos73
Cuadro No. 6.1. Cantidad de parámetros por modelo y conjunto de datos76
Cuadro No. 6.2. Tiempo total de la fase de entrenamiento por modelo79
Cuadro No. 6.3. Tiempo total de la fase de validación de los modelos81
Cuadro No. 6.4. Total de memoria consumida por modelo84
Cuadro No. 6.5. Comparación de pérdida inicial y final durante la fase de entrenamiento en el conjunto de datos UNPC
Cuadro No. 6.6. Comparación de pérdida inicial y final durante la fase de entrenamiento en el conjunto de datos OSPC
Cuadro No. 6.7. Comparación de pérdida inicial y final durante la fase de validación en el conjunto de datos UNPC92
Cuadro No. 6.8. Comparación de pérdida inicial y final durante la fase de validación en el conjunto de datos OSPC95
Cuadro No. 6.9. Comparación de la exactitud durante la primera y última época de la fase de entrenamiento para el conjunto de datos UNPC99
Cuadro No. 6.10. Comparación de la precisión durante la primera y última época de la fase de validación para el conjunto de datos OSPC102
Cuadro No. 6.11. Comparación de la precisión durante el entrenamiento y la validación en el conjunto de datos UNPC

Cuadro No. 6.12. Comparación de la precisión durante la primera y última época de la fase de
entrenamiento para el conjunto de datos OSPC105
Cuadro No. 6.13. Comparación de la precisión durante la primera y última época de la fase de validación para el conjunto de datos OSPC
Cuadro No. 6.14. Comparación de la precisión durante el entrenamiento y la validación en el conjunto de datos OSPC
Cuadro No. 6.15. Comparación de brechas entre la precisión durante el entrenamiento y la precisión durante la validación
Cuadro No. 6.16. Resultados de la evaluación de la calidad de traducción mediante BLEU115
Cuadro No. 6.17. Tiempo promedio de inferencia por registro119
Cuadro No. 6.18. Tiempo total de inferencia del conjunto de datos de prueba121

Índice de gráficos

Gráfico No. 6.1. Número de parámetros entrenables de los modelos de la comparativa76
Gráfico No. 6.2. Contribución relativa de los modelos al tiempo total de ejecución de la comparativa
Gráfico No. 6.3. Tiempo total de la fase de entrenamiento por modelo80
Gráfico No. 6.4. Tiempo total de la fase de validación de los modelos
Gráfico No. 6.5. Total de memoria consumida por modelo (en Megabytes)84
Gráfico No. 6.6. Comparación de pérdida inicial y final durante la fase de entrenamiento en el conjunto de datos UNPC
Gráfico No. 6.7. Comportamiento de la pérdida de los modelos durante el entrenamiento para e conjunto de datos UNPC.
Gráfico No. 6.8. Comparación de pérdida inicial y final durante la fase de entrenamiento en el conjunto de datos OSPC
Gráfico No. 6.9. Comportamiento de la pérdida de los modelos durante el entrenamiento para e conjunto de datos OSPC
Gráfico No. 6.10. Comparación de pérdida inicial y final durante la fase de validación en el conjunto de datos UNPC
Gráfico No. 6.11. Comportamiento de la medida de pérdida de los modelos durante la validación para el conjunto UNPC94
Gráfico No. 6.12. Comparación de pérdida inicial y final durante la fase de validación en el conjunto de datos OSPC96
Gráfico No. 6.13. Comportamiento de la medida de pérdida de los modelos durante la validación para el conjunto de datos OSPC
Gráfico No. 6.14 Comparación de la exactitud durante la primera y última época de la fase de entrenamiento para el conjunto de datos UNPC99
Gráfico No. 6.15. Comportamiento de la medida de exactitud de los modelos durante e entrenamiento en el conjunto de datos UNPC

Gráfico No. 6.16. Comportamiento de la medida de exactitud de los modelos durante la validación
para el conjunto de datos UNPC103
Gráfico No. 6.17. Comparación de la precisión durante la primera y última época de la fase de entrenamiento para el conjunto de datos OSPC
Gráfico No. 6.18. Comportamiento de la medida de exactitud de los modelos durante el entrenamiento en el conjunto de datos OSPC
Gráfico No. 6.19. Comportamiento de la medida de exactitud de los modelos durante la validación para el conjunto de datos OSPC
Gráfico No. 6.20. Perplejidad durante la fase de validación en el conjunto de datos UNPC112
Gráfico No. 6.21. Perplejidad durante la fase de validación en el conjunto de datos OSPC113
Gráfico No. 6.22. BLEU4
Gráfico No. 6.23. Tiempo promedio de inferencia por registro119
Gráfico No. 6.24. Tiempo total de inferencia121
Índice de anexos
Anexo No. 1. Arquitecturas de modelos y sus hiperparámetros
Anexo No. 2. Resultados de pérdida, precisión y perplejidad139
Anexo No. 3. Uso de memoria y duración por épocas

1 INTRODUCTION

1.1 Motivación

Después de la popularización de los Transformers (Vaswani et al., 2017), que revolucionaron el procesamiento del lenguaje natural y otros campos de la inteligencia artificial, han surgido varias variantes y avances significativos en la arquitectura que son aplicables al procesamiento del lenguaje natural y en particular a la traducción automática de textos.

Cada una de estas arquitecturas representa un avance significativo en el campo del procesamiento del lenguaje natural y la inteligencia artificial. Comprender sus diferencias, fortalezas y limitaciones permite a los investigadores y estudiantes mantenerse actualizados con el estado del arte y los últimos avances tecnológicos en NLP (Procesamiento del Lenguaje Natural).

Estos avances reflejan una tendencia hacia la mejora de la eficiencia computacional, la capacidad de manejar secuencias más largas y la adaptabilidad a diferentes tipos de datos y tareas dentro del campo del procesamiento del lenguaje natural y más allá.

Sin embargo, la eficiencia computacional es crucial en la implementación práctica de modelos de inteligencia artificial. Al comparar estas arquitecturas, se pueden evaluar aspectos como el tiempo de entrenamiento, la complejidad computacional y los requisitos de memoria, lo que ayuda a seleccionar modelos que sean más eficientes en términos de recursos.

Algunas de estas arquitecturas representan variantes o extensiones de modelos existentes (por ejemplo, *Transformer* y MoE). Comprender cómo estas variaciones afectan el rendimiento puede inspirar investigaciones futuras sobre nuevas mejoras o adaptaciones de estas arquitecturas para problemas específicos.

Consecuentemente, realizar una comparativa permite identificar cuál de estas arquitecturas podría ser más efectiva para diferentes contextos de aplicación.

La capacidad de generalización de un modelo es crucial para su aplicación en entornos del mundo real. Al evaluar estas arquitecturas en una variedad de conjuntos de datos y tareas, se puede entender mejor cómo manejan la generalización y la adaptabilidad a diferentes dominios lingüísticos y contextos.

Realizar una comparativa sistemática y rigurosa entre estas arquitecturas no solo proporciona información práctica para aplicaciones actuales, sino que también contribuye al conocimiento científico general sobre la efectividad de diferentes enfoques en el campo del procesamiento del lenguaje natural y la inteligencia artificial.

1.2 Planteamiento del problema

El trabajo final de máster en Inteligencia Artificial se enfoca en realizar una comparativa exhaustiva de diversas arquitecturas de redes neuronales para el modelado de traducción neuronal (NMT) que han surgido después de la introducción de Transformers.

Los *Transformers*, han demostrado ser efectivos en tareas como la traducción automática y el procesamiento del lenguaje natural. En los últimos años, la traducción neuronal automática ha avanzado significativamente con el desarrollo de nuevas arquitecturas más allá de los modelos Transformer "originales" (Vaswani et al., 2017).

Este estudio se centra en evaluar y comparar varias de estas arquitecturas avanzadas para entender sus fortalezas y limitaciones en tareas de traducción neuronal.

1.3 Estructura de la memoria

2 CONTEXTO Y ESTADO DEL ARTE

2.1 Traducción automática

La traducción automática (NMT, *Neural Machine Translation*) es un enfoque avanzado en el campo de la traducción automática que utiliza redes neuronales profundas (*deep learning*) para traducir texto de un idioma a otro (Koehn, 2020).

En NMT, tanto el proceso de traducción como el aprendizaje se realizan utilizando redes neuronales profundas. Estas redes consisten en capas de neuronas interconectadas que procesan información de manera secuencial y aprenden representaciones abstractas de los datos.

A diferencia de los sistemas de traducción automática estadística tradicionales, que se basan en reglas y características lingüísticas específicas, los modelos NMT aprenden directamente de los datos utilizando un enfoque de aprendizaje *end-to-end*. Esto significa que el modelo de traducción se

entrena para convertir directamente secuencias de palabras de un idioma a otro, sin requerir etapas intermedias de extracción de características o alineación de frases.

El objetivo de aprendizaje es encontrar la secuencia de destino correcta dada la secuencia de origen, lo cual puede verse como un problema de clasificación de alta dimensionalidad que intenta mapear las dos oraciones en el espacio semántico. En todos los principales modelos modernos de NMT, este proceso se puede dividir en un paso de codificación y un paso de decodificación, y así separar funcionalmente todo el modelo (Yang et al., 2024). Consecuentemente, un modelo NMT generalmente sigue una arquitectura básica de codificador-decodificador (*encoder-decoder*) que se detalla más adelante.

Inicialmente, el modelo procesa la secuencia de palabras de entrada en un espacio de representación vectorial (*embedding*) de alta dimensionalidad. Este *embedding* captura el significado semántico de la frase en el idioma original.

Para mejorar la calidad de la traducción y manejar eficientemente secuencias largas, muchos modelos NMT utilizan mecanismos de atención. La atención permite al modelo enfocarse en partes relevantes de la entrada durante la fase de decodificación, mejorando así la coherencia y la precisión de la traducción.

Los modelos NMT requieren grandes volúmenes de datos paralelos (pares de oraciones traducidas) para entrenarse efectivamente. Durante el entrenamiento, el modelo ajusta los pesos de las conexiones neuronales para minimizar la diferencia entre las traducciones generadas y las traducciones reales en el conjunto de datos de entrenamiento.

A medida que se introducen más datos y se optimizan las técnicas de entrenamiento y arquitecturas de red, los modelos NMT han demostrado ser capaces de generar traducciones de mayor calidad que los sistemas tradicionales.

En este sentido Eisenstein (2019) destaca la relación entre la hipótesis distribucional y la capacidad de los modelos de *deep learning* para aprender la traducción automática neuronal (NMT) sin necesidad de conocer explícitamente reglas gramaticales o sintácticas se fundamenta en cómo estos modelos procesan y aprenden a partir de grandes cantidades de datos.

La hipótesis distribucional propuesta por Harris (1954) y popularizado por Firth (1957) es un principio lingüístico que sugiere que las palabras que aparecen en contextos similares tienden a tener

significados similares. De forma simplificada, el contexto en el que aparece una palabra proporciona pistas sobre su significado. Esta idea subraya que el significado de una palabra se puede inferir a partir de las palabras que suelen aparecer en su proximidad. Este concepto se fundamenta en la observación de que el contexto de una palabra es crucial para comprender su significado.

De este modo la hipótesis distribucional proporciona el marco teórico que respalda la efectividad de los modelos modernos de traducción automática basados en *deep learning*, para aprender y aplicar conocimientos lingüísticos sin necesidad de reglas gramaticales explícitas.

Así, esta premisa de la hipótesis distribucional se traduce en cómo los modelos de *deep learning* representan y aprenden de las palabras. En esencia, los modelos de lenguaje modernos implementan esta hipótesis al aprender representaciones vectoriales de palabras que reflejan sus contextos. Estas representaciones vectoriales implementan el aprendizaje de representaciones al generar *embeddings* de palabras. Como se mencionó previamente, los *embeddings* son vectores que capturan el significado contextual de las palabras, alineándose con la hipótesis distribucional al representar palabras que comparten contextos similares con vectores similares.

Este principio es esencial para los modelos de *deep learning*, ya que estos aprenden patrones estadísticos a partir de datos (Bishop & Bishop, 2024). En el caso de la traducción automática neuronal, el modelo NMT no necesita conocimientos explícitos sobre las reglas gramaticales o sintácticas de los idiomas involucrados. En cambio, aprende a asociar patrones de palabras y frases enteras en un idioma con sus correspondientes en otro idioma.

Los modelos de *deep learning*, como las redes neuronales profundas utilizadas en NMT, son capaces de aprender representaciones distribuidas de palabras y frases (Eisenstein, 2019). Esto significa que cada palabra o frase se representa mediante vectores numéricos en un espacio de alta dimensionalidad.

Durante el entrenamiento, el modelo NMT ajusta estos vectores (*embedding*) y los pesos de las conexiones neuronales para minimizar el error de traducción en un conjunto de datos paralelos (pares de frases en diferentes idiomas).

A medida que el modelo procesa más ejemplos y ajusta sus parámetros, captura de manera implícita reglas sintácticas y gramaticales. Por ejemplo, aprende a generar traducciones

gramaticalmente correctas al aprender qué estructuras de frases son más comunes y efectivas en cada idioma.

La capacidad de los modelos de *deep learning* para aprender sin necesidad de reglas explícitas radica en su habilidad para generalizar a partir de datos (Prince, 2023). En lugar de depender de reglas predefinidas, estos modelos aprenden a través de ejemplos y exposición a datos diversos mediante *transducción* (Gammerman et al., 2013).

Esto significa que el modelo NMT puede capturar fenómenos lingüísticos complejos que podrían no estar explícitamente codificados en reglas gramaticales tradicionales. Por ejemplo, puede aprender a manejar variaciones idiomáticas, ambigüedades contextuales y estructuras sintácticas específicas de cada idioma (Koehn, 2020).

Los **modelos de lenguaje** como los Transformers (por ejemplo, BERT, GPT), se basan en el aprendizaje de representaciones para comprender y generar texto. Estos modelos utilizan grandes cantidades de datos y aprenden a predecir palabras o frases en función de su contexto. La efectividad de estos modelos se basa en gran medida en la capacidad de capturar y utilizar el contexto para entender el significado de las palabras (Kamath et al., 2022).

El **aprendizaje de representaciones** es un enfoque en el aprendizaje automático donde se busca aprender representaciones significativas y útiles de los datos, a menudo a partir de datos sin etiquetar. En el contexto del procesamiento del lenguaje natural (NLP), esto implica aprender representaciones vectoriales de palabras que capturen sus significados en función de los contextos en los que aparecen (Liu et al., 2023).

A medida que los modelos de lenguaje se entrenan en grandes cuerpos de texto, mejoran su capacidad para capturar relaciones contextuales sutiles y complejas entre palabras. Este proceso de aprendizaje continuo mejora la efectividad de los modelos para tareas de NLP, como la traducción automática, el análisis de sentimientos o la generación de texto, al afinar las representaciones sustentadas por la hipótesis distribucional.

Estos modelos aprenden a través de patrones estadísticos derivados de grandes volúmenes de datos, lo que les permite generar traducciones precisas y naturales entre idiomas.

2.2 Transducción

Muchas tareas de aprendizaje automático pueden expresarse como la transformación, o transducción, de secuencias de entrada en secuencias de salida: reconocimiento de voz, traducción automática, predicción de la estructura secundaria de proteínas y síntesis de voz, por mencionar algunas (Graves, 2012).

La capacidad de transformar y manipular secuencias es una parte crucial de la inteligencia humana: todo lo que se sabe sobre el mundo llega en forma de secuencias sensoriales, y todo lo que se hace para interactuar con el mundo requiere secuencias de acciones y pensamientos. La creación de transductores automáticos de secuencias parece ser, por lo tanto, un paso importante hacia la inteligencia artificial. Un problema importante al que se enfrentan estos sistemas es cómo representar la información secuencial de manera que sea invariante, o al menos robusta, frente a distorsiones secuenciales. Además, esta robustez debe aplicarse tanto a las secuencias de entrada como a las de salida.

La transducción de secuencias, también conocida como modelado de secuencia a secuencia, es una tarea de aprendizaje automático que implica convertir una secuencia de entrada en una secuencia de salida, potencialmente de diferentes longitudes.

El término se utiliza en algunas aplicaciones de redes neuronales recurrentes en problemas de predicción de secuencias, como algunos problemas en el dominio del procesamiento de lenguaje natural.

La transducción (o aprendizaje transductivo) se usa en el campo de la teoría del aprendizaje estadístico para referirse a la predicción de ejemplos específicos dados ejemplos específicos de un dominio. Vale la pena distinguir el proceso "inductivo" consiste en derivar la función a partir de los datos dados de la "deducción" que implica derivar los valores de la función dada para puntos de interés y, finalmente, "transducción" es derivar los valores de la función desconocida para puntos de interés a partir de los datos dados.

El modelo de estimación del valor de una función en un punto de interés describe un nuevo concepto de inferencia: pasar de lo particular a lo particular. A este tipo de inferencia se le denomina *inferencia transductiva*. Es importante destacar que este concepto de inferencia aparece cuando se desea obtener el mejor resultado a partir de una cantidad restringida de información.

Clásicamente, la transducción se ha utilizado al hablar de lenguaje natural, como en el campo de la lingüística. Por ejemplo, existe la noción de "gramática de transducción" que se refiere a un conjunto de reglas para transformar ejemplos de un idioma a otro.

Una gramática de transducción describe un par de lenguajes estructuralmente correlacionados y genera pares de oraciones, en lugar de oraciones individuales. La oración en el idioma origen es (se supone que es) una traducción de la oración en el idioma destino.

Este uso de la transducción al hablar de teoría y traducción automática clásica influye en el uso del término cuando se trata de la predicción moderna de secuencias con redes neuronales recurrentes en tareas de procesamiento de lenguaje natural.

Más generalmente, la transducción se usa en tareas de predicción de secuencias en NLP, específicamente en traducción. Por ejemplo, Grefenstette et al. (2015) describen la transducción como la asignación de una cadena de entrada a una cadena de salida.

Muchas tareas de procesamiento de lenguaje natural (NLP) pueden verse como problemas de transducción, es decir, aprender a convertir una cadena en otra. La traducción automática es un ejemplo prototípico de transducción y los resultados recientes indican que las RNN tienen la capacidad de codificar largas cadenas de entrada y producir traducciones coherentes.

Un transductor se define de manera estricta como un modelo que produce una salida por cada paso de tiempo de entrada proporcionado. Esto se relaciona con el uso lingüístico, específicamente con los transductores de estados finitos. Sin, embargo un Transductor Neuronal (Jaitly et al., 2015) es una clase más general de modelos de aprendizaje de secuencia a secuencia y puede producir bloques de salidas (posiblemente de longitud cero) a medida que llegan bloques de entradas, cumpliendo así con la condición de ser "online". El modelo genera salidas para cada bloque utilizando una RNN transductora que implementa un modelo de secuencia a secuencia (Seq-to-Seq).

2.3 Arquitecturas y modelos

De acuerdo con Vardasbi et al. (2023) la traducción automática se modela usando una red neuronal. En los modelos de codificador-decodificador, como el Transformer, el modelo tiene dos componentes principales: un codificador para capturar las dependencias del lado de la fuente y un decodificador para capturar las dependencias del lado del objetivo y las dependencias entre fuente y objetivo. Alternativamente, la traducción automática (MT) puede tratarse como una tarea de

Modelado del Lenguaje, donde el modelo (solo decodificador) se entrena con las oraciones fuente y objetivo concatenadas, separadas por un token especial [SEP].

Un modelo de secuencia a secuencia (Seq-to-Seq) típicamente presenta una arquitectura compuesta por uno o más codificadores y uno o más decodificadores.

De acuerdo con Anastasopoulos & Chiang (2018) y Sperber et al. (2020) definir la arquitectura con dos o más decodificadores implicaría "dual-tasking" (Le et al., 2020) y (Yu et al., 2022), mientras que definirla con dos codificadores o "cruzar" los decodificadores conllevaría "multi-modality" a nivel de entrada o salida (Sun et al., 2023).

Un modelo codificador-decodificador mapea una secuencia de entrada a una secuencia objetivo con ambas secuencias de longitud arbitraria (Sutskever et al., 2014). Este tipo de modelo tienen aplicaciones que van desde la traducción automática hasta la predicción de series temporales. Más específicamente, este mecanismo utiliza una RNN, o cualquiera de sus variantes como *LSTM* (*Long Short Term Memory*) o *GRU* (*Gated Recurrent Unit*), para mapear la secuencia de entrada a un vector de longitud fija, y otra RNN o una de sus variantes para decodificar la secuencia objetivo a partir de ese vector.

Tanto el "encoder" y el "decoder" son componentes clave en modelos de redes neuronales diseñados para tareas de procesamiento de lenguaje natural, trabajando juntos para transformar y generar secuencias de texto de manera efectiva y precisa.

El "encoder", o codificador es la parte de la red que transforma la entrada (como texto o una secuencia de palabras) en una representación numérica o vectorial. Esta representación captura el significado semántico y estructural de la entrada original en un formato que la red neuronal puede entender y procesar.

Por su parte, el *decoder*, o decodificador, es la parte de la red neuronal que utiliza la representación numérica generada por el *encoder* para producir una salida en la forma deseada. En el contexto de traducción automática, el *decoder* toma la representación interna de la frase en el idioma origen, generada por el *encoder*, y la transforma en la secuencia de palabras en el idioma destino.

El proceso de *encoder-decoder* es fundamental en tareas de procesamiento de lenguaje natural como la traducción automática y la generación de texto, donde el *encoder* ayuda a comprender

y representar el significado de la entrada, mientras que el *decoder* utiliza esta representación para producir la salida deseada en un formato comprensible para los humanos.

El estudio de S. Wang et al. (2021) concluye que los modelos de lenguaje grandes tienen un gran potencial para mejorar la traducción automática sin necesidad de arquitectura de codificador-decodificador, demostrando su capacidad para tareas que involucran tanto comprensión como generación de lenguaje

2.4 Atención

El trabajo de Bahdanau et al. (2014) fue fundamental para el desarrollo de modelos de traducción automática basados en redes neuronales, particularmente en la evolución de los modelos *Encoder-Decoder* con mecanismos de atención, este aporte es relevante particularmente por:

- Introducción del mecanismo de atención que permite que el modelo se enfoque en partes específicas de la secuencia de entrada durante la traducción, mejorando la precisión y la fluidez de las traducciones generadas.
- El enfoque propuesto aborda de manera integral dos problemas clave en la traducción automática: la alineación precisa de las palabras en las dos lenguas y la generación de la traducción. Al aprender simultáneamente a alinear y traducir, el modelo puede mejorar la coherencia y la fidelidad de las traducciones.
- Mejora en la fluidez y la coherencia: Al integrar el aprendizaje de la alineación junto con el proceso de traducción, el modelo puede capturar mejor las correspondencias entre palabras y frases en diferentes idiomas. Esto resulta en traducciones más precisas y naturales, reduciendo errores de alineación y mejorando la cohesión del texto traducido.

Sin duda el aporte de Bahdanau et al. (2014) ha inspirado numerosos estudios y avances en la arquitectura de modelos de traducción automática basados en redes neuronales, promoviendo el uso generalizado de mecanismos de atención en diversos modelos y aplicaciones. El éxito del enfoque propuesto en este artículo allanó el camino para el desarrollo de modelos más sofisticados y efectivos en traducción automática, como los Transformers. Estos modelos han revolucionado el campo al mejorar aún más la capacidad de manejar contextos largos y capturar relaciones complejas entre las palabras.

Al distinguir entre secuencias "en distribución" y "fuera de distribución" 2, se puede evaluar cómo un modelo generaliza a ejemplos que no están representados en el conjunto de entrenamiento (Z. Wang, 2023). Esto ayuda a determinar si el modelo ha aprendido una función subyacente general que puede aplicarse a nuevas secuencias de diferentes longitudes (fuera de distribución), o si solo está aproximando los datos que vio durante el entrenamiento (en distribución).

De esta forma si un modelo no puede generalizar bien a secuencias de longitud diferente a las de entrenamiento, entonces no ha aprendido la función subyacente real y solo ha aprendido a ajustar los datos en distribución.

Por su parte Z. Wang (2023) destaca que la atención mejora enormemente la eficiencia del aprendizaje, al reducir tanto la complejidad del modelo como la complejidad de la muestra, así como la robustez del aprendizaje, en términos de rendimiento de generalización y problemas de sobreajuste. Sin embargo, la atención no supera la limitación de generalización fuera de distribución, ya que no cambia la naturaleza autorregresiva de los modelos. No obstante, la impresionante eficiencia de aprendizaje acelerada por la atención refleja su motivación original, es decir, "aprender a alinear" (Bahdanau et al., 2015).

Si bien Z. Wang (2023) señala que los modelos de red neuronal recurrente de secuencia a secuencia (RNN seq2seq) se han utilizado para aprender cuatro tareas de transducción: identidad, inversión, reducción total y copia cuadrática. Estas tareas de transducción se han estudiado tradicionalmente usando transductores de estado finito y se les ha atribuido una creciente complejidad. El autor señala que se ha observado que los modelos RNN seq2seq solo pueden aproximar un mapeo que se ajusta a los datos de entrenamiento o datos en distribución, en lugar de aprender las funciones subyacentes. Aunque el mecanismo de atención mejora la eficiencia y robustez del aprendizaje, no supera la limitación de generalización para datos "fuera de distribución".

De acuerdo con Z. Wang (2023), los modelos con atención (*attentional*) muestran que la atención mejora significativamente la capacidad de los modelos para ajustarse a los datos de

² Se refiere a secuencias de entrada cuya longitud no está dentro del rango de longitudes usadas durante el entrenamiento. Siguiendo el mismo ejemplo, una secuencia de longitud 5 o 16 sería "fuera de distribución".

¹ Se refiere a secuencias de entrada cuya longitud está dentro del rango de longitudes de secuencias usadas durante el entrenamiento del modelo. Por ejemplo, si un modelo fue entrenado con secuencias de longitud entre 6 y 15, entonces una secuencia de longitud 10 sería "en distribución".

entrenamiento y generalizar a los datos de prueba. Los modelos con atención siempre logran una mejor precisión agregada en secuencias completas tanto en los conjuntos de entrenamiento como en los de prueba, y tienen una menor variabilidad entre el entrenamiento y la prueba en comparación con los modelos sin atención (*attention-less*). Además, los modelos con atención superan a los modelos sin atención en la generalización a ejemplos fuera de distribución. En otras palabras, los modelos *RNN seq2seq* con atención son más fuertes en el aprendizaje dentro de la distribución y tienen una capacidad de generalización fuera de distribución relativamente mejor en comparación con los modelos sin atención.

Aunque la atención mejora la capacidad del modelo para manejar datos fuera de distribución en comparación con los modelos sin atención, el problema de generalización aún persiste, aunque en menor grado.

2.5 Transformers

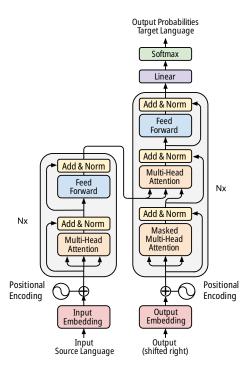
Los *Transformers* (Vaswani et al., 2017) representan una alternativa innovadora al enfoque tradicional de las redes neuronales recurrentes (RNN) en el campo del procesamiento de lenguaje natural y otras tareas secuenciales.

El *Transformer* es un modelo prominente en el aprendizaje profundo ampliamente adoptado en campos como el procesamiento del lenguaje natural (NLP), visión por computadora (CV) y procesamiento de voz. Inicialmente propuesto como un modelo secuencia a secuencia para la traducción automática, los modelos pre-entrenados basados en *Transformer* han logrado el estado del arte en diversas tareas, consolidándolo como la arquitectura principal en NLP, especialmente para PTMs.

Además de aplicaciones en lenguaje, el *Transformer* ha sido utilizado en CV, procesamiento de audio e incluso en disciplinas como química y ciencias de la vida. Diversas variantes del *Transformer*, conocidas como *X-formers*, han sido propuestas para mejorar su eficiencia en el procesamiento de secuencias largas, su capacidad de generalización y adaptación a tareas específicas (T. Lin et al., 2022).

El Transformer básico es un modelo secuencia a secuencia que consta de un codificador y un decodificador, cada uno compuesto por una pila de L bloques idénticos. Cada bloque del codificador se compone principalmente de un módulo de autoatención multi-cabeza y una red de alimentación hacia adelante (FFN) por posición. Para construir un modelo más profundo, se utiliza una conexión

residual alrededor de cada módulo, seguida de un módulo de Normalización de Capa. Comparados con los bloques del codificador, los bloques del decodificador insertan adicionalmente módulos de atención cruzada entre los módulos de autoatención multi-cabeza y las FFNs por posición. Además, los módulos de autoatención en el decodificador están adaptados para evitar que cada posición atienda a posiciones subsecuentes.



Al comparar los Transformers frente a RNN tradicionales, se pueden destacar las siguientes innovaciones:

- Arquitectura sin recurrencia: A diferencia de las RNN tradicionales, que procesan secuencias de manera secuencial y recurrente, los Transformers no tienen conexiones recurrentes. En su lugar, utilizan mecanismos de atención para capturar las relaciones entre las diferentes partes de la secuencia de entrada simultáneamente.
- Mecanismo de atención: Los Transformers utilizan capas de atención que calculan pesos para cada par de elementos en la secuencia de entrada, permitiendo que el modelo "atienda" partes específicas de la secuencia que son relevantes para la tarea en cuestión. Esto facilita la captura de dependencias a largo plazo y mejora la capacidad del modelo para manejar secuencias de longitud variable (Chaudhari et al., 2021) y (Brauwers & Frasincar, 2023).

Paralelización: Debido a su naturaleza no recurrente y al uso eficiente de la atención, los
 Transformers son más fácilmente paralelizables en comparación con las RNN. Esto los hace
 más rápidos de entrenar y de utilizar en aplicaciones prácticas.

Los autores T. Lin et al. (2022) destacan a la autoatención (self-attention) como la pieza central de los Transformers y mecanismo flexible para manejar entradas de longitud variable. Se puede entender como una capa completamente conectada en la que los pesos se generan dinámicamente a partir de las relaciones entre pares de entradas (*pairwise*).

Entre las ventajas de la autoatención se pueden mencionar:

- Longitud del camino máxima: La autoatención tiene la misma longitud máxima del camino que las capas completamente conectadas, lo que la hace adecuada para modelar dependencias a largo alcance. En comparación con las capas completamente conectadas, la autoatención es más eficiente en términos de parámetros y más flexible para manejar entradas de longitud variable.
- 2. Campo receptivo y profundidad de red: Debido al campo receptivo limitado de las capas convolucionales, generalmente se necesita apilar una red profunda para tener un campo receptivo global. Por otro lado, la longitud máxima del camino constante permite que la autoatención modele dependencias a largo alcance con un número constante de capas.
- 3. **Operaciones paralelizadas**: Las operaciones secuenciales constantes y la longitud máxima del camino hacen que la autoatención sea más paralelizable y más eficiente en el modelado de dependencias a largo alcance en comparación con las capas recurrentes.

Del mismo modo, al considerar el sesgo inductivo o suposiciones que un modelo hace sobre la distribución de los datos para simplificar el aprendizaje (Bishop & Bishop, 2024). De esta forma T. Lin et al. (2022) destaca que cada tipo de red neuronal tiene diferentes sesgos inductivos que afectan cómo procesan la información.

En el caso de las redes CNN (Convolutional Neural Networks) imponen sesgos inductivos de *invariancia de traducción* y *localidad*. Esto significa que las CNNs son buenas para reconocer patrones que aparecen en diferentes posiciones (traducción) y para enfocarse en áreas locales de la imagen a través de funciones de *kernel* compartidas.

Por su parte los enfoques previos de recurrencia como LSTM o GRU tienen sesgos inductivos de *invariancia temporal* y *localidad* gracias a su estructura Markoviana. Esto les permite manejar datos secuenciales, como el texto o las series temporales, tratando cada paso de la secuencia como dependiente del anterior de manera local.

Por el contrario, la arquitectura de *Transformer* tiene muy pocos supuestos sobre la estructura de los datos. Esto significa que no impone sesgos inductivos específicos como las CNNs o las RNNs, lo que la convierte en una arquitectura universal y flexible. Sin embargo, esta flexibilidad también puede llevar a que el *Transformer* se sobreajuste cuando se entrena con datos de pequeña escala, ya que la falta de sesgos estructurales hace que sea más propenso a captar ruido en lugar de patrones significativos (T. Lin et al., 2022).

A pesar de su arquitectura innovadora, los Transformers mantienen el marco general de los modelos *Encoder/Decoder* al separar claramente las fases de codificación y decodificación. Esto los hace altamente efectivos para tareas como la traducción automática, donde la secuencia de entrada (por ejemplo, una oración en un idioma) se codifica en una representación semántica universal, que luego se decodifica en la secuencia de salida (la misma oración traducida a otro idioma).

Como lo menciona T. Lin et al. (2022), la arquitectura *Transformer* puede ser utilizada de tres formas diferentes:

- Codificador-Decodificador. Se utiliza la arquitectura completa del *Transformer* como se introdujo en anteriormente. Esto se utiliza típicamente en modelado secuencia a secuencia (por ejemplo, traducción automática neuronal).
- Solo codificador. Solo se utiliza el codificador y las salidas del codificador se utilizan como representación de la secuencia de entrada. Esto se usa generalmente para problemas de clasificación o etiquetado de secuencias.
- Solo decodificador. Solo se utiliza el decodificador, donde también se elimina el módulo de atención cruzada codificador-decodificador. Esto se utiliza típicamente para generación de secuencias, como modelado de lenguaje.

Siguiendo esta misma línea, Xiao & Zhu (2023) profundizan en variantes del *Transformer* que se utilizan en diversas aplicaciones, como modelos *encoder-only* (usados en clasificación de texto), *decoder-only* (para generación de texto), y el modelo *encoder-decoder* (para tareas como la traducción).

Además, se abordan mejoras recientes, como cambios en los mecanismos de atención y nuevas técnicas para codificar información posicional, algunas de las cuales se abordan a continuación.

Si bien los Transformers representan una evolución significativa respecto a las RNN tradicionales al introducir mecanismos de atención y mejorar la capacidad de manejar dependencias a largo plazo. A pesar de estas diferencias, siguen siendo modelos *Encoder / Decoder*, destacando por su eficiencia y rendimiento en una amplia gama de aplicaciones de procesamiento de lenguaje natural y más allá.

En el modo codificador-decodificador, generalmente hay múltiples módulos de auto-atención multi-cabezal, incluyendo una auto-atención estándar tanto en el codificador como en el decodificador, junto con una atención cruzada codificador-decodificador que permite al decodificador utilizar información del codificador. Esto influye en el diseño del mecanismo de auto-atención. En el modo codificador, no hay restricción o condición de que el mecanismo de auto-atención deba ser causal, es decir, dependiente únicamente de los tokens presentes y pasados. En el entorno codificador-decodificador, la auto-atención utilizada en el decodificador (es decir, a través de posiciones de decodificación) debe ser causal, ya que cada paso de decodificación auto-regresivo solo puede depender de tokens anteriores, mientras que la auto-atención utilizada en el codificador no necesita serlo. Cumplir con este requisito puede resultar desafiante para muchos diseños eficientes de auto-atención. (Tay et al., 2022).

Los costos computacionales de los Transformers son influenciados por varios factores clave. Primero, la complejidad computacional surge principalmente del mecanismo de atención. Específicamente, calcular la matriz de atención implica un costo cuadrático en relación con la longitud de la secuencia de entrada N, denotado como $N \times N$. Esta complejidad cuadrática es evidente en operaciones como la multiplicación de matrices para las consultas Q y las llaves K, que por sí solas consumen tiempo y memoria proporcional a N^2 . Esto limita la utilidad de los modelos auto-atentivos en aplicaciones que requieren el procesamiento de secuencias largas, afectando especialmente durante el entrenamiento debido a las actualizaciones de gradiente, aunque menos durante la inferencia donde los gradientes no se calculan (Tay et al., 2022).

El costo cuadrático de la auto-atención afecta tanto la velocidad de las fases de entrenamiento como de inferencia. Sin embargo, el costo computacional total de un *Transformer* también incluye contribuciones significativas de las capas de alimentación hacia adelante presentes en cada bloque

del *Transformer*. Estas capas, aunque lineales en complejidad con respecto a la longitud de la secuencia, aún representan una parte no trivial de la carga computacional. Se han realizado esfuerzos para mitigar estos costos, como explorar la dispersión en las redes de alimentación hacia adelante para escalarlas sin un aumento correspondiente en los costos computacionales.

La eficiencia en los mecanismos de atención y en la computación general del modelo son consideraciones algo independientes pero entrelazadas. Algunos métodos de atención eficientes tienen como objetivo reducir explícitamente las longitudes de secuencia, ahorrando así costos de computación tanto en las capas de atención como en las de alimentación hacia adelante. Gestionar estos compromisos entre eficiencia y demandas computacionales es una tarea compleja, abordada frecuentemente en investigaciones recientes que exploran diversas estrategias de optimización y arquitecturas de modelos.

Como mencionan Katharopoulos et al. (2020), se han realizado pocos esfuerzos para comprender la autoatención desde una perspectiva teórica. El mismo autor menciona que se han propuesto formulaciones basadas en *kernels*, en la que la atención se interpreta como la aplicación de un suavizador de *kernels* sobre las entradas, con las puntuaciones de los *kernels* representando la similitud entre las entradas. Esta formulación facilita la comprensión de los componentes de la atención y la integración de las incrustaciones posicionales. En contraste, la formulación del *kernel* se utiliza para acelerar el cálculo de la autoatención y reducir su complejidad computacional. Además, se observa que, si se aplica un *kernel* con puntuaciones de similitud positivas a las consultas y claves, la atención lineal converge normalmente. Más recientemente, demostraron teórica y empíricamente que una autoatención multi-cabeza con un número suficiente de cabezas puede replicar cualquier capa convolucional. En lugar de esto, se muestra que una capa de autoatención entrenada con un objetivo autorregresivo puede ser interpretada como una red neuronal recurrente, lo que permite acelerar significativamente el tiempo de inferencia en modelos de transformadores autorregresivos.

Como se ha mencionado, la autoatención juega un papel crucial en la arquitectura *Transformer*, pero presenta dos desafíos importantes en aplicaciones prácticas (T. Lin et al., 2022):

 Complejidad: La complejidad de la autoatención es O(T² · D), donde T es la longitud de la secuencia y D es la dimensión del vector de características. Esto significa que el módulo de atención puede volverse un cuello de botella cuando se trata de secuencias largas, ya que el tiempo de cálculo crece cuadráticamente con la longitud de la secuencia. 2. Prior Estructural: La autoatención no hace suposiciones sobre la estructura de los datos de entrada, ni siquiera sobre el orden de los elementos en la secuencia. Por lo tanto, la información sobre el orden debe ser aprendida a partir de los datos de entrenamiento. Esto hace que los Transformers, sin pre-entrenamiento, sean propensos al sobreajuste cuando se utilizan con datos pequeños o de tamaño moderado.

Como lo menciona T. Lin et al. (2022), para abordar estos desafíos, se han desarrollado varias mejoras en el mecanismo de atención:

- Atención Dispersa (Sparse Attention): Esta línea de investigación introduce sesgos de escasez en el mecanismo de atención, lo que reduce la complejidad al limitar el número de interacciones necesarias.
- Atención Linealizada (Linearized Attention): En este enfoque, la matriz de atención se descompone en mapas de características de núcleo (kernel). Luego, la atención se calcula en un orden inverso para lograr una complejidad lineal en lugar de cuadrática (Katharopoulos et al., 2020).
- Compresión de Prototipos y Memoria (Prototype and Memory Compression): Estos métodos reducen el número de pares de consultas o memorias clave-valor para disminuir el tamaño de la matriz de atención.
- 4. **Autoatención de Bajo Rango (Low-rank Self-Attention):** Esta línea de trabajo captura la propiedad de bajo rango de la autoatención (Bhojanapalli et al., 2020), lo que ayuda a reducir la complejidad computacional (Zhang et al., 2022).
- 5. **Atención con Priori (Attention with Prior):** Esta investigación explora cómo complementar o sustituir la atención estándar con distribuciones de atención previas, lo que puede mejorar el rendimiento al incorporar conocimientos adicionales.
- Mecanismo Multi-Cabeza Mejorado (Improved Multi-Head Mechanism): Este enfoque
 investiga diferentes mecanismos alternativos para la atención multi-cabeza, buscando
 mejorar la eficiencia y efectividad del proceso de atención.

Cada una de estas mejoras busca abordar los desafíos asociados con la autoatención, haciéndola más eficiente y efectiva para diferentes aplicaciones y tamaños de datos.

2.6 Universal Transformers

Sobre la base de la arquitectura de *Transformers*, Dehghani et al. (2019) introducen la variante de "*Universal Transformers*" que representa un avance notable en el campo de las redes neuronales y el procesamiento de lenguaje natural. Introduce un enfoque innovador que mejora los modelos *Transformers* estándar al permitir una adaptabilidad dinámica en cada capa del modelo.

A diferencia de los *Transformers* tradicionales, donde todas las capas y posiciones de la secuencia realizan la misma cantidad de operaciones, los *"Universal Transformers"* pueden ajustar dinámicamente la cantidad de pasos de computación en función de las necesidades específicas de cada posición en la secuencia.

Esta capacidad dinámica no solo mejora la eficiencia computacional del modelo al reducir el número de operaciones innecesarias, sino que también potencia su capacidad para capturar dependencias a largo plazo de manera más efectiva.

Los experimentos realizados por Dehghani et al. (2019) demostraron que los *Universal Transformers* superan a los modelos estándar en diversas tareas complejas de procesamiento de lenguaje natural, como la traducción automática y el análisis de sentimientos. Esta superioridad se debe a su capacidad para adaptarse de manera flexible a diferentes partes de la secuencia, capturando relaciones más complejas entre palabras y frases.

2.7 Reformer

El trabajo de Kitaev et al. (2020) es una contribución significativa en la evolución de los modelos *Transformers*, centrada en mejorar la eficiencia computacional y el manejo de secuencias largas. Uno de los aspectos clave de este avance es el uso de técnicas de *hashing*, específicamente el LSH (*Locality Sensitive Hashing*), aplicado a la capa de atención.

El *Reformer* aborda el desafío de la escalabilidad de los modelos Transformers al introducir técnicas que reducen drásticamente el costo computacional y de memoria, permitiendo el procesamiento eficiente de secuencias largas. Esto se logra mediante el uso de aproximaciones eficientes en el cálculo de atención, como *hashing* y atención por *chunking*, que reducen la complejidad cuadrática típica de la atención.

A diferencia de otros modelos que enfrentan problemas con secuencias extensas debido a la memoria requerida, el *Reformer* utiliza técnicas como la atención reversible y la compresión de parámetros para manejar de manera efectiva secuencias de longitud considerable, sin comprometer la calidad de las representaciones aprendidas.

Los experimentos realizados muestran que el *Reformer* mantiene un rendimiento competitivo en comparación con los *Transformers* estándar, pero con una significativa reducción en los requisitos de recursos computacionales. Esto lo hace adecuado para aplicaciones en escenarios donde el procesamiento de grandes volúmenes de datos es crucial, como en la traducción automática y el análisis de texto a gran escala.

En los modelos *Transformer* estándar, el cálculo de atención entre todas las posiciones de una secuencia implica una complejidad computacional cuadrática con respecto a la longitud de la secuencia. Esto puede volverse prohibitivo con secuencias muy largas debido a la gran cantidad de cálculos necesarios.

El *Reformer* aborda este problema mediante el uso de LSH, una técnica que aproxima el cálculo de atención al dividir el espacio de las posiciones en cubetas hash. En lugar de calcular la atención exacta entre todas las posiciones, LSH agrupa posiciones similares en el espacio de representación, reduciendo así la cantidad de pares de atención que deben considerarse explícitamente.

En la práctica, el *Reformer* utiliza LSH para dividir la matriz de consulta y la matriz de clave en segmentos, que luego se atienden de manera aproximada utilizando operaciones de hash. Esto permite que el modelo mantenga un rendimiento competitivo mientras reduce significativamente la carga computacional y la memoria necesaria para procesar secuencias largas.

El *Reformer* ha influido en la investigación al mostrar nuevas vías para mejorar la eficiencia y escalabilidad de los modelos *Transformer*, inspirando desarrollos posteriores en la optimización de atención y la gestión de recursos en redes neuronales de gran escala. Esta técnica innovadora demuestra cómo la investigación continua en optimización puede mejorar radicalmente la eficiencia de las redes neuronales en el procesamiento de lenguaje natural y otras tareas de aprendizaje automático.

2.8 Evolved Transformer

El So et al. (2019) proponen "*The Evolved Transformer*" como un enfoque novedoso para mejorar el rendimiento de los modelos Transformer mediante una búsqueda evolutiva automatizada de arquitecturas neuronales.

La propuesta de un *Evolved Transformer* explora cómo optimizar el diseño estructural del *transformer* a través de una técnica de búsqueda evolutiva, mejorando los modelos estándar a nivel de arquitectura.

Para ello se propone un algoritmo evolutivo que optimiza arquitecturas de Transformer, utilizando un método llamado "Búsqueda de Arquitectura con Hurdles Dinámicos Progresivos" (PDH, por sus siglas en inglés). Este algoritmo ajusta la arquitectura del modelo a medida que evoluciona, creando barreras (hurdles) que los modelos deben superar en términos de rendimiento en pruebas de validación. Si un modelo no supera una de estas barreras, se detiene su entrenamiento, evitando así la computación innecesaria. Los modelos que pasan todas las pruebas reciben un entrenamiento adicional y continúan mejorando.

El enfoque se basa en seleccionar entre múltiples configuraciones arquitectónicas, que incluyen variaciones en las capas, funciones de activación y métodos de normalización. Cada "modelo hijo" tiene una representación genética que describe sus componentes, y los mejores modelos se seleccionan para continuar evolucionando. Este proceso garantiza que se exploren múltiples combinaciones, lo que lleva a encontrar arquitecturas más eficientes.

De acuerdo al documento original (So et al., 2019), esta arquitectura logró superar a la arquitectura Transformer estándar en varias tareas, incluyendo la traducción automática. Los autores realizaron estudios de ablación para entender qué mutaciones específicas en la arquitectura del Transformer contribuyeron a esta mejora en el rendimiento. Entre los hallazgos clave, se descubrió que las capas de atención al codificador eran esenciales para el buen desempeño, y que la eliminación de la normalización de capas degradaba significativamente la calidad del modelo.

Este método automatizado permite explorar eficientemente grandes espacios de diseño arquitectónico sin intervención humana directa, reduciendo el tiempo y esfuerzo necesario para encontrar configuraciones óptimas. La búsqueda evolutiva con PDH es particularmente útil porque ajusta dinámicamente los pasos de entrenamiento y los recursos según el rendimiento de los modelos

en desarrollo, lo que ayuda a evitar el sobreentrenamiento o el uso excesivo de recursos computacionales en modelos que no muestran promesas.

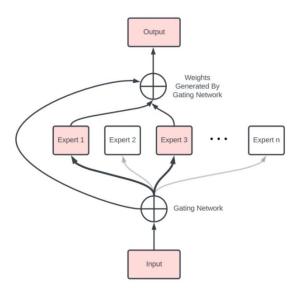
2.9 Mixture of Experts (MoE)

El uso de Modelos de Expertos Dispersos (MoE) en la traducción automática ofrece varias ventajas significativas:

En los MoE, solo un subconjunto de expertos se activa para cada entrada. Esto reduce la cantidad de cálculos necesarios durante la inferencia y el entrenamiento, ya que no todos los expertos contribuyen a cada paso. Esta activación selectiva permite que el modelo maneje grandes cantidades de datos de entrada sin un aumento proporcional en el costo computacional.

Al activar solo una parte de los expertos, se reduce el consumo de memoria y la carga computacional. Esto es especialmente valioso cuando se trabajan con grandes modelos de traducción automática que pueden ser muy costosos en términos de recursos.

Los expertos en un MoE pueden especializarse en diferentes aspectos del proceso de traducción, como distintos pares de idiomas, contextos específicos, o tipos de texto (técnico, literario, etc.). Esta especialización puede mejorar la calidad de la traducción al permitir que el modelo maneje de manera más efectiva las variaciones en el lenguaje y el contenido.



El mecanismo de *gating* en MoE permite que el modelo adapte dinámicamente qué expertos se activan según la entrada específica. Esto puede ayudar a mejorar la precisión y la relevancia de las traducciones al seleccionar los expertos más adecuados para cada entrada.

MoE permite construir modelos más grandes y complejos sin un aumento lineal en el costo computacional. Esto significa que se pueden usar modelos con más expertos sin tener que aumentar significativamente el poder de procesamiento necesario, lo que facilita el escalado para tareas de traducción más complejas.

La capacidad de manejar eficientemente grandes volúmenes de datos y secuencias largas es crucial para la traducción automática, especialmente cuando se trabaja con textos extensos o complejos. Los modelos MoE puede ayudar a gestionar estos desafíos sin una penalización significativa en términos de tiempo y recursos.

La capacidad de añadir o ajustar expertos en un modelo MoE permite una adaptación más fácil a nuevos idiomas o dominios. Por ejemplo, se pueden añadir expertos especializados en un nuevo par de idiomas o en un tipo específico de contenido sin necesidad de reentrenar el modelo completo.

Los modelos MoE pueden ser entrenados de manera más efectiva en diferentes subconjuntos de datos, lo que facilita el ajuste del modelo para tareas específicas o para mejorar el rendimiento en ciertos contextos de traducción.

Dado que solo se activan algunos expertos durante el entrenamiento, el modelo puede ser entrenado de manera más eficiente. Esto reduce el tiempo y el costo asociados con el entrenamiento de modelos grandes y complejos.

2.9.1 ST-MoE

El trabajo de Zoph et al. (2022) presenta una arquitectura para modelos de expertos dispersos (*sparse expert models*) que se enfoca en la estabilidad y la capacidad de transferencia. La arquitectura ST-MoE se basa en un modelo de expertos dispersos que combina varios componentes clave:

- Modelo de Expertos Dispersos (MoE): En esta arquitectura, el modelo se divide en varios
 "expertos" o submodelos especializados, cada uno entrenado para manejar diferentes
 aspectos del problema. Solo una fracción de estos expertos se activa en cada paso de
 inferencia, lo que hace que el modelo sea más eficiente.
- Combinador de Expertos: Un componente importante de ST-MoE es el combinador, que decide qué expertos se activarán para una entrada dada. Este componente utiliza

mecanismos de atención o *gating* para seleccionar los expertos más relevantes, lo que ayuda a mantener la eficiencia y a reducir la complejidad computacional.

- 3. **Estabilidad y Transferibilidad:** El paper introduce técnicas para mejorar la estabilidad durante el entrenamiento y la capacidad de transferencia a nuevas tareas o dominios. Esto se logra mediante la regularización y técnicas específicas para ajustar los pesos de los expertos de manera que el modelo sea robusto a variaciones en los datos.
- 4. **Arquitectura Híbrida:** ST-MoE puede combinarse con diferentes tipos de redes neuronales, como redes neuronales profundas o transformadores, para aprovechar sus fortalezas mientras mantiene la eficiencia y la capacidad de adaptación.

La arquitectura ST-MoE busca abordar desafíos comunes en los modelos de expertos dispersos, como la inestabilidad durante el entrenamiento y la dificultad para transferir conocimientos a nuevas tareas, proporcionando una solución que mejora tanto la eficacia como la flexibilidad del modelo.

2.9.2 Switch Transformers

Un enfoque interesante surge del trabajo de Fedus et al. (2022)al introducir los modelos Switch Transformers y su capacidad para escalar eficientemente a modelos con billones de parámetros utilizando técnicas de dispersión (*sparsity*) simples y eficientes.

Uno de los aspectos destacados de este trabajo es el uso de MoE (*Mixture of Experts*), o Mezcla de Expertos, una técnica que permite mejorar la eficiencia y el rendimiento de los modelos de gran escala. Los modelos MoE son una forma de estructurar redes neuronales donde diferentes expertos (subredes especializadas) se activan de manera selectiva según el tipo de entrada o tarea específica, permitiendo una gestión más eficiente de los recursos computacionales y una mejor adaptabilidad a diferentes tipos de datos.

En el contexto de los *Switch Transformers*, los MoE juegan un papel crucial al permitir que el modelo active solo un subconjunto de expertos relevantes para una tarea particular, lo que reduce significativamente la carga computacional en comparación con enfoques donde todos los parámetros están activos simultáneamente. Esto no solo mejora la eficiencia computacional, sino que también puede llevar a mejores resultados de rendimiento al permitir una mejor especialización y adaptabilidad del modelo a diferentes partes del espacio de entrada.

Por lo tanto, el uso de MoE en *Switch Transformers* no solo es relevante para la escalabilidad a modelos con un gran número de parámetros, sino que también destaca cómo la estructuración inteligente y la gestión de recursos pueden ser críticas para avanzar en la capacidad de procesamiento y el rendimiento de los modelos de inteligencia artificial a gran escala.

Vale la pena retomar el concepto de MoE (Zhou, 2012) y detallar sobre las diferencias entre el enfoque clásico de ensemble. Los modelos MoE (*Mixture of Experts*) y los modelos Ensemble son enfoques distintos utilizados en *machine learning* para mejorar el rendimiento y la robustez de los modelos, pero difieren significativamente en su estructura y funcionamiento:

En un MoE, el modelo está compuesto por múltiples submodelos (llamados expertos) que son especializados en diferentes partes del espacio de entrada o en diferentes tareas específicas. Durante la inferencia, en lugar de activar todos los expertos simultáneamente, se utiliza una red de gate (o "interruptor") que decide qué expertos deben ser activados en función de la entrada actual. Esto permite que el modelo se adapte dinámicamente a diferentes características de los datos o a diferentes partes de una tarea compleja. Mejora la eficiencia computacional al no tener que procesar todos los expertos para cada entrada, y permite una mejor adaptabilidad y especialización del modelo.

Por otro lado, los *Ensemble Models* combinan múltiples modelos base (que pueden ser del mismo tipo o diferentes) para mejorar el rendimiento general del sistema. Cada modelo base en el Ensemble puede ser entrenado de manera independiente y luego sus predicciones se combinan mediante alguna estrategia de agregación (como votación, promedio, etc.).

Los MoE se diferencian de los enfoques clásicos de *Ensemble Models* pues se estructura un solo modelo con múltiples expertos activados selectivamente, mientras que Ensemble combina modelos independientes. Por otro lado, MoE utiliza recursos de manera más eficiente al activar selectivamente expertos, mientras que Ensemble requiere mantener y ejecutar múltiples modelos base por separado. MoE permite una adaptación dinámica a diferentes condiciones o partes del espacio de entrada, mientras que Ensemble combina estáticamente las predicciones de modelos preentrenados.

Si bien existen aplicaciones recientes muy interesantes tales como MoE-Mamba (Pióro et al., 2024) y MoEUT (Csordás et al., 2024), se excluyen del alcance del presente trabajo.

2.10 Space-State Models

Los modelos de espacio de estados son un marco matemático utilizado para describir el comportamiento de un sistema mediante un conjunto de ecuaciones diferenciales de primer orden. En este enfoque, el estado del sistema se representa mediante un vector de variables, y la dinámica del sistema se describe a través de cómo cambian estas variables de estado a lo largo del tiempo.

En términos más específicos:

- 1. **Estado del Sistema**: Se representa con un vector de variables (generalmente llamado x), que encapsula toda la información necesaria sobre el sistema en un momento dado.
- 2. **Ecuaciones Diferenciales**: Las dinámicas del sistema se describen con una ecuación diferencial de primer orden que muestra cómo cambia el vector de estado (x) a lo largo del tiempo. Esta ecuación suele tener la forma:

$$\frac{\mathrm{dx}}{\mathrm{dt}} = Ax + Bu$$

donde A es una matriz que describe las interacciones internas del sistema, B es una matriz que describe cómo las entradas u afectan el estado, y u es el vector de entrada del sistema.

3. **Salida del Sistema**: Además, el modelo también puede incluir una ecuación que describe cómo se obtiene la salida del sistema (y\mathbf{y}y) a partir del estado y las entradas:

$$y = Cx + Du$$

donde C y D son matrices que definen la relación entre el estado y la salida, y entre la entrada y la salida, respectivamente.

Los modelos estructurados de secuencia en el espacio de estados (S4) son una clase reciente de modelos de secuencia para el aprendizaje profundo que están ampliamente relacionados con las RNN, las CNN y los modelos clásicos de espacio de estados. Los autores Gu et al. (2021) abordan el desafío de modelar eficientemente secuencias largas al introducir espacios de estado estructurados.

Este enfoque innovador permite reducir la complejidad computacional al seleccionar de manera inteligente subconjuntos de estados relevantes para cada paso de tiempo en una secuencia. Al optimizar el procesamiento de información en secuencias extensas, el modelo no solo mejora la

eficiencia computacional, sino que también mantiene una capacidad robusta para capturar dependencias a largo plazo y estructuras complejas en los datos.

La capa S4 (Gu et al., 2021) realiza una transformación secuencia-a-secuencia no lineal, mapeando una secuencia de entrada a una secuencia de salida. Esta capa contiene HHH SSMs independientes de "entrada única y salida única" (SISO) con estados de dimensión NNN. Cada SSM de S4 se aplica a una dimensión de la secuencia de entrada, produciendo una transformación lineal independiente por canal de entrada. Luego se aplica una función de activación no lineal y, finalmente, una capa de mezcla lineal por posición para combinar las características independientes y generar la secuencia de salida.

La capa S4 utiliza el marco HiPPO (Gu et al., 2020) para la aproximación de funciones en línea, inicializando las matrices de estado con una matriz HiPPO, lo que ha demostrado mejorar el rendimiento. Aunque la matriz HiPPO-LegS no es diagonalizable de forma estable, se puede representar en una forma diagonal más baja (DPLR), lo que S4 utiliza para derivar una forma eficiente del núcleo de convolución.

La implementación eficiente de la capa S4 depende del contexto: en modo recurrente para generación en línea y en modo convolucional para secuencias completas. En el modo convolucional, la recurrencia lineal se representa como una convolución unidimensional, aprovechando transformadas rápidas de Fourier (FFTs) para paralelizar el proceso.

Cada capa S4 tiene parámetros entrenables, incluidos los parámetros de los SSMs y los de la capa de mezcla. Para cada SSM, se utilizan matrices de entrada, transición, salida y *feedthrough*, y se aplican parámetros de escala temporal para discretizar el tiempo continuo. La notación compacta para los estados y salidas de los SSMs se usa para simplificar la representación matemática.

La capa S5 (Smith et al., 2022) sustituye el banco de SSMs SISO (o el gran sistema bloque-diagonal) en S4 por un SSM MIMO con un tamaño de estado latente PPP y dimensiones de entrada y salida HHH. Esta versión discretizada del SSM MIMO se aplica a una secuencia de entrada vectorial u1:Lu_{1:L}u1:L, produciendo una secuencia de salidas del SSM (o preactivaciones) y1:Ly_{1:L}y1:L, utilizando estados latentes xkx_kxk. Luego, se aplica una función de activación no lineal para obtener la salida de la capa u1:L'u'_{1:L}u1:L'. A diferencia de S4, S5 no necesita una capa lineal adicional por posición, ya que las características están mezcladas de manera intrínseca. Además, el tamaño de

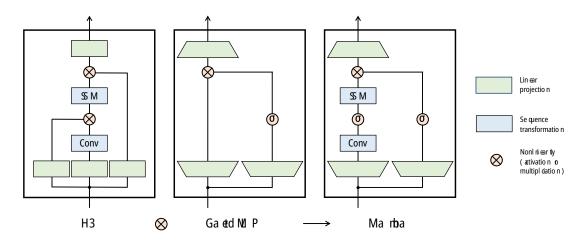
estado latente PPP de S5 puede ser significativamente menor que el tamaño HNHNHN del SSM bloque-diagonal en S4, permitiendo el uso de escaneos paralelos eficientes.

2.10.1 Mamba

Por su parte, Gu & Dao (2024) introducen innovaciones significativas en la modelización de secuencias al proponer un enfoque que mejora la eficiencia computacional y la capacidad de captura de dependencias a largo plazo en secuencias en una arquitectura denominada Mamba.

Una de las mejoras clave propuestas por Mamba es la introducción de espacios de estado selectivos (*Selective State Spaces*). Esto implica que en lugar de manejar todas las posibles combinaciones de estados en una secuencia de manera exhaustiva, el modelo Mamba selecciona dinámicamente un subconjunto relevante de estados para procesar en cada paso de tiempo. Esta selección se realiza de manera inteligente utilizando métodos de atención o *gating*, lo que permite al modelo concentrarse en las partes críticas de la secuencia y omitir el procesamiento de estados menos informativos o redundantes.

Al reducir el número de estados considerados en cada paso de tiempo, Mamba logra reducir drásticamente la complejidad computacional, permitiendo así el procesamiento en tiempo lineal respecto a la longitud de la secuencia. Esto es crucial para aplicaciones que manejan secuencias largas o requieren tiempos de respuesta rápidos.



Como lo describe Lieber et al. (2024) el diseño de bloque simplificado que se utilizando combina el bloque H3, que es la base de la mayoría de las arquitecturas de modelos de espacio de estados (SSM), con el bloque MLP ubicuo de las redes neuronales modernas. En lugar de entrelazar estos dos bloques, simplemente se repite el bloque Mamba de manera homogénea. En comparación

con el bloque H3, Mamba reemplaza la primera compuerta multiplicativa con una función de activación. En comparación con el bloque MLP, Mamba añade un SSM a la rama principal y para la activación, se utiliza la función SiLU (también conocida como Swish).

A pesar de la selección de estados, Mamba conserva la capacidad de modelar dependencias complejas y a largo plazo en las secuencias. Esto se logra mediante la adaptación dinámica de los estados seleccionados en función de la historia de la secuencia y las características del problema específico.

El enfoque de espacios de estado selectivos permite una mayor flexibilidad en la adaptación del modelo a diferentes tipos de secuencias y contextos. Puede ajustarse según las características específicas de la aplicación y las propiedades de las secuencias de datos.

Finalmente Vardasbi et al. (2023b) discute la aplicación del modelo S4 (*Structured State Spaces for Sequences*) en la traducción automática. Este modelo, que ha mostrado éxito en tareas de modelado de secuencias, como visión y audio, tiene la ventaja de comprimir la entrada en un único estado oculto, capturando dependencias a largo plazo sin necesitar mecanismos de atención.

Sin embargo, en el ámbito de la traducción automática, S4 no rinde tan bien como los modelos *Transformer*, mostrando una disminución de aproximadamente 4 puntos BLEU en su desempeño, especialmente con oraciones largas. La principal limitación identificada es la incapacidad del S4 para resumir la oración completa en un solo estado oculto. El estudio propone que la introducción de un mecanismo de atención en el modelo S4 podría cerrar esta brecha de rendimiento.

2.10.2 Simplified Structured State Space Sequence (S5)

Por su parte, Smith et al. (2022) presenta una arquitectura innovadora para el modelado de secuencias, centrada en la simplificación de las capas de espacio de estados.

La arquitectura se basa en la idea de simplificar las capas de espacio de estados (*state space layers*) para el modelado de secuencias. Los espacios de estados se utilizan para modelar dependencias a largo plazo en secuencias, y la simplificación busca hacer estas capas más eficientes y efectivas. Al simplificar las capas, se busca mejorar la capacidad del modelo para manejar secuencias largas y complejas de manera más eficiente.

La arquitectura simplifica las operaciones matemáticas complejas típicamente asociadas con los modelos de espacio de estados. Esto incluye la reducción de la complejidad computacional y la mejora de la estabilidad numérica. La simplificación permite que el modelo sea más fácil de implementar y más rápido de entrenar, sin sacrificar la capacidad de modelado.

El enfoque está diseñado para mejorar el rendimiento en tareas de modelado de secuencias, como el procesamiento de lenguaje natural (NLP) y la predicción de series temporales.

La arquitectura facilita el aprendizaje de dependencias a largo plazo en las secuencias, mejorando la capacidad del modelo para capturar patrones complejos.

La arquitectura propuesta por Smith et al. (2022) es compatible con redes neuronales profundas, lo que permite su integración con otras técnicas y modelos existentes en el campo del aprendizaje profundo.

Se puede utilizar en combinación con otros tipos de capas y arquitecturas, ofreciendo flexibilidad en el diseño de modelos.

2.10.3 Gated State Spaces

2.10.3.1 Modelos de Espacio de Estados y Dependencias a Largo Plazo:

Modelos de Espacio de Estados (SSM): Son modelos matemáticos que utilizan variables de estado para describir la dinámica de un sistema a lo largo del tiempo. Son eficaces para modelar secuencias en las que las dependencias abarcan largos períodos.

Dependencias a Largo Plazo: En los datos secuenciales (como textos o códigos), ciertos elementos pueden depender de elementos distantes en la secuencia. Por ejemplo, para entender una palabra en una oración, a veces es necesario conocer el contexto desde el inicio de la oración.

2.10.3.2 Modelado Autoregresivo de Secuencias:

Modelos Autoregresivos: Predicen el siguiente valor en una secuencia basado en los valores anteriores. En el contexto de datos secuenciales, generan o predicen secuencias paso a paso.

De acuerdo con Mehta et al. (2022) se explora el uso de modelos de espacio de estados para el modelado de dependencias a largo plazo en secuencias. Los modelos de espacio de estados (SSM) son herramientas matemáticas que emplean variables de estado para describir la dinámica temporal de un sistema, lo que permite el modelado de secuencias en las que las dependencias se extienden a

lo largo de períodos prolongados. Esta capacidad resulta particularmente útil en datos secuenciales, como textos o códigos, donde ciertos elementos pueden depender de otros situados a distancias considerables en la secuencia. Por ejemplo, comprender una palabra en una oración puede requerir considerar el contexto desde el inicio de la oración.

En el contexto del modelado autoregresivo de secuencias, los modelos autoregresivos se utilizan para predecir el siguiente valor en una secuencia en función de los valores anteriores, generando o prediciendo secuencias de manera secuencial. Esta técnica es crucial en aplicaciones de procesamiento del lenguaje y análisis de código.

El trabajo sobre *Gated State Spaces (GSS)* de Mehta et al. (2022) presenta varios avances significativos en el modelado de secuencias largas, particularmente en el contexto de tareas autoregresivas. Uno de los aportes principales es la mejora en la eficiencia del modelo. Tradicionalmente, los modelos basados en atención, como los Transformadores, tienen una complejidad cuadrática con respecto a la longitud de la secuencia, lo que implica un costo computacional prohibitivo para secuencias largas. En contraste, GSS reduce esta complejidad de O(L2) a $O(L\log L)$ mediante el uso de la Transformada Rápida de Fourier (FFT) para realizar las convoluciones, lo que permite manejar secuencias de gran longitud de manera mucho más eficiente.

Además de la reducción en la complejidad, otro aspecto destacable de GSS es su capacidad de paralelización. Mientras que muchos modelos autoregresivos no pueden ser completamente paralelizados debido a la naturaleza secuencial del proceso de predicción, GSS es completamente paralelizable en la dimensión de longitud. Esta propiedad mejora significativamente la eficiencia en el entrenamiento, permitiendo un procesamiento más rápido en comparación con modelos como los basados en redes neuronales recurrentes (RNNs).

Una de las innovaciones clave que distingue a GSS es el uso de unidades de activación con puertas, o *gating units*, que ya habían demostrado ser útiles en otros contextos. Estas unidades permiten reducir la dimensionalidad de las operaciones más costosas, como las FFT, lo que acelera el entrenamiento sin comprometer la capacidad del modelo para capturar dependencias de largo alcance. Esta idea, que se basa en observaciones previas sobre la efectividad de las capas de alimentación con puertas, se extiende en GSS al contexto de los modelos de espacio de estado, lo que contribuye a una reducción significativa en el tiempo de entrenamiento sin afectar negativamente el rendimiento.

Un aspecto particularmente interesante de GSS es su simplicidad comparada con modelos anteriores como S4 y DSS. Mientras que S4, por ejemplo, dependía de una inicialización cuidadosa basada en matrices HiPPO para obtener buenos resultados, GSS elimina esta complejidad. Los autores del trabajo encontraron que el modelo no es tan sensible a la inicialización, lo que permite entrenarlo con inicializaciones aleatorias sin que se vea afectado su rendimiento. Esta simplificación contrasta con la dependencia de S4 y DSS en trucos de álgebra lineal numérica para garantizar su funcionamiento.

En cuanto a su relación con trabajos previos, GSS está profundamente vinculado con modelos como S4 y DSS. El modelo GSS se basa en la estructura de espacio de estado diagonal que S4 popularizó, y también aprovecha la eficiencia computacional que brinda la convolución mediante FFT. No obstante, a diferencia de S4, que presentaba cuellos de botella en el tiempo de entrenamiento, especialmente en hardware especializado como las TPUs, GSS logra superar estas limitaciones mediante un diseño más optimizado. Por otro lado, el modelo Mamba, que también aborda el uso de espacios de estado diagonalizados, introduce ciertas mejoras, pero GSS va un paso más allá al simplificar algunos de los aspectos más complejos del proceso, lo que resulta en un modelo más rápido en su entrenamiento. Un aspecto notable de la capa GSS es su capacidad para generalizar a entradas más largas o nuevos tipos de datos sin necesidad de entrenamiento adicional, lo que se conoce como generalización *zero-shot*. Esta característica resulta especialmente valiosa en aplicaciones prácticas donde los datos pueden variar con el tiempo.

Finalmente, se incorpora el mecanismo de autoatención para modelar las dependencias locales dentro de las secuencias. La combinación de autoatención con la capa GSS mejora la capacidad del modelo para capturar tanto las dependencias locales como las de largo plazo. La simplicidad en la implementación de la capa GSS la convierte en una opción práctica para investigadores y profesionales, facilitando su integración en sistemas de modelado de secuencias.

2.10.4 Jamba

Fuera del alcance de la presente comparativa, avances posteriores como Jamba (Lieber et al., 2024) representa un avance significativo al integrar dos enfoques innovadores en modelos de lenguaje: *Transformer* y *Mamba*. Este trabajo combina lo mejor de ambos mundos para mejorar tanto

la eficiencia computacional como la capacidad de modelado en el procesamiento de secuencias de texto largas.

En primer lugar, el modelo *Transformer* es conocido por su capacidad para capturar relaciones a largo plazo en secuencias mediante mecanismos de atención. Sin embargo, debido a su estructura densa y completa conectividad, los Transformers pueden ser computacionalmente costosos cuando se aplican a secuencias muy largas o a grandes volúmenes de datos.

Por otro lado, Mamba introduce los espacios de estado selectivos, lo que permite reducir la complejidad computacional al seleccionar dinámicamente un subconjunto de estados relevantes para cada paso de tiempo en una secuencia. Esto mejora la eficiencia del modelo al tiempo que mantiene su capacidad para modelar dependencias complejas a lo largo de la secuencia.

Jamba combina estas dos técnicas de manera híbrida. Utiliza la estructura del Transformer para capturar relaciones a largo plazo y la eficiencia de los espacios de estado selectivos de Mamba para manejar eficazmente secuencias largas. Esta combinación permite que el modelo Jamba mejore tanto la velocidad de procesamiento como la capacidad de modelado, haciendo posibles aplicaciones prácticas en el procesamiento de lenguaje natural donde tanto la eficiencia como la precisión son críticas.

2.10.5 MoE Mamba

Por otro lado, MoE-Mamba (Pióro et al., 2024) combina dos enfoques innovadores en modelado de secuencias: Mixture of Experts (MoE) y Mamba, que utiliza espacios de estado selectivos para mejorar la eficiencia y la capacidad de modelado. MoE permite gestionar de manera eficiente recursos computacionales al activar selectivamente expertos especializados según la tarea o contexto, mientras que Mamba optimiza el procesamiento al seleccionar dinámicamente estados relevantes en secuencias largas. Esta combinación no solo mejora la escalabilidad de los modelos al reducir la complejidad computacional, sino que también preserva la capacidad de capturar dependencias a largo plazo en secuencias complejas.

2.11 Moving Average Equipped Gated Attention

Por su parte, MEGA (*Moving Average Equipped Gated Attention*) es un mecanismo de atención de una sola cabeza que incluye un promedio móvil exponencial (EMA) para abordar dos limitaciones

del *Transformer* tradicional: la falta de sesgo inductivo y la complejidad cuadrática en la secuencia. Los aspectos principales de MEGA incluyen:

- **Inductivo posicional:** Al incorporar EMA, MEGA introduce un sesgo inductivo que captura dependencias locales, algo que la atención tradicional del Transformer no hace.
- Complejidad: MEGA es capaz de reducir la complejidad cuadrática del Transformer a complejidad lineal mediante la división de secuencias largas en "chunks" (bloques) con pérdida mínima de información contextual.
- Rendimiento superior: En experimentos con tareas de modelado de secuencias, como traducción automática, clasificación de imágenes y modelado de lenguaje, MEGA mostró mejoras significativas en comparación con Transformers y otros modelos recientes (Ma et al., 2022).
- Diseño unificado: MEGA está diseñado para ser eficiente tanto en tareas de dependencia de corto como largo alcance, abordando problemas de diversas modalidades de datos (lenguaje, audio, imagen y video) (Ma et al., 2022).

MEGA se diferencia de un Transformer "vainilla" en varios aspectos clave:

2.11.1.1 Atención con sesgo inductivo frente a atención sin sesgo inductivo:

MEGA incorpora un promedio móvil exponencial (EMA), lo que le permite modelar dependencias locales de manera más efectiva al introducir un sesgo inductivo posicional. Este sesgo ayuda a capturar dependencias entre tokens cercanos en una secuencia de forma más eficiente.

Transformers vainilla (el Transformer original de Vaswani et al.) no tienen un sesgo inductivo fuerte, lo que significa que la relación entre los tokens se aprende completamente de los datos, sin suposiciones previas sobre su proximidad o estructura posicional.

2.11.1.2 Complejidad computacional

MEGA mejora la eficiencia al reducir la complejidad cuadrática del cálculo de la atención en Transformers vainilla a una complejidad lineal. Esto se logra al dividir las secuencias largas en "chunks" o bloques fijos, lo que permite procesar secuencias más largas con menor costo computacional.

Los Transformers tienen una complejidad cuadrática en relación con la longitud de la secuencia debido al cálculo de las matrices de atención, lo que los hace ineficientes para secuencias largas (Vaswani et al., 2017).

2.11.1.3 Mecanismo de atención

MEGA utiliza un mecanismo de atención de una sola cabeza con compuertas (gated attention), lo que le permite ser tan expresivo como la atención de múltiples cabezas de un Transformer vainilla, pero con una arquitectura más simple.

Los Transformers (Vaswani et al., 2017) dependen de la atención multi-cabeza, lo que significa que múltiples atenciones se calculan en paralelo para capturar diferentes relaciones entre los tokens. Aunque es flexible, aumenta la complejidad computacional y la necesidad de recursos.

2.11.1.4 Manejo de dependencias locales y globales:

MEGA está diseñado específicamente para capturar tanto dependencias de corto como de largo alcance de manera más eficiente mediante el uso del EMA y las compuertas en su mecanismo de atención.

Los Transformers (Vaswani et al., 2017) tienen dificultades para capturar dependencias locales debido a su enfoque de atención global, lo que puede llevar a problemas cuando se trabaja con secuencias largas.

Ma et al. (2022) proponen un mecanismo de atención con medias móviles y puertas (Mega) para abordar simultáneamente dos debilidades. La idea principal es incorporar sesgos inductivos en el mecanismo de atención a lo largo de la dimensión temporal, aprovechando el enfoque clásico de la media móvil exponencial (EMA). La EMA captura dependencias locales que decrecen exponencialmente con el tiempo (ver Figura 1) y ha sido ampliamente utilizada en la modelación de datos de series temporales (§2). Se introduce una forma multidimensional atenuada de EMA con coeficientes ajustables (§3.1) y, a continuación, se desarrolla el mecanismo de atención con medias móviles y puertas integrando la EMA con una variante de la atención con una sola cabeza (Hua et al., 2022) (§3.2). Teóricamente, se demuestra que la atención con una sola cabeza es tan expresiva como la atención con múltiples cabezas, que es la más comúnmente utilizada (§3.3). Aprovechando el mecanismo de media móvil incorporado, se propone además una variante de Mega con complejidad

lineal, denominada Mega-chunk, que segmenta las secuencias de entrada en bloques fijos con una pérdida mínima de información contextual (§3.5).

Experimentalmente, a través de cinco tareas de modelado de secuencias en diversos tipos de datos, incluyendo modelado de secuencias con contexto largo, traducción automática neuronal, modelado de lenguaje autoregresivo, y clasificación de imágenes y voz, se demuestra que Mega supera significativamente a una variedad de modelos de referencia fuertes, tanto en términos de efectividad como de eficiencia (§4) (ver Tabla 1). Estas mejoras ilustran la importancia de modelar dependencias a largo y corto plazo mediante diferentes patrones de sesgos inductivos.

2.12 Extended Long Short-Term Memory

Los autores Beck et al. (2024) proponen xLSTM (*Extended Long Short-Term Memory*), es una innovación en redes neuronales recurrentes diseñada para mejorar el modelado de dependencias a largo plazo en secuencias temporales.

Este enfoque extiende la arquitectura estándar de LSTM mediante nuevas técnicas y capas adicionales que mejoran la captura y procesamiento de información a lo largo de secuencias largas. xLSTM supera las limitaciones de las LSTM convencionales al integrar avanzados mecanismos de atención y adaptabilidad, permitiendo una mejor representación de relaciones temporales complejas y contextos de larga duración en datos secuenciales. Este avance es crucial en aplicaciones donde la precisión temporal y la capacidad para modelar dependencias complejas son esenciales, como en predicción de series temporales o procesamiento de lenguaje natural.

Para abordar cómo mejorar el modelado de secuencias mediante LSTMs escalados a miles de millones de parámetros, se introducen técnicas como el *gating* exponencial con técnicas de normalización y estabilización adecuadas. Además, se modifica la estructura de memoria de LSTM, dando lugar a variantes como sLSTM con memoria y actualización escalares, y mLSTM que es completamente paralelizable con una estructura de memoria matricial y reglas de actualización de covarianza. Estas extensiones de LSTM se integran en bloques residuales que forman arquitecturas xLSTM, las cuales se apilan de forma residual. El uso de gating exponencial y estructuras de memoria modificadas potencia las capacidades de xLSTM, permitiéndole competir favorablemente con Transformers y Modelos de Espacio de Estados en términos de rendimiento y escalabilidad.

3 OBJETIVOS Y METODOLOGÍA DE TRABAJO

3.1 Objetivo general

Evaluar y comparar varias nuevas arquitecturas avanzadas posteriores a los Transformers para entender sus fortalezas y sus limitaciones en tareas de traducción neuronal automática.

3.2 Objetivos específicos

- Evaluar el rendimiento de diferentes arquitecturas de NMT posteriores a Transformers en términos de calidad de traducción y eficiencia computacional.
- Analizar las características clave de cada arquitectura en términos de capacidad de aprendizaje, escalabilidad y generalización.
- Identificar casos de uso específicos donde cada arquitectura pueda sobresalir, considerando diferentes idiomas y dominios de aplicación.

3.3 Metodología de trabajo

3.3.1 Tipo de investigación

Para llevar a cabo una investigación propedéutica que compare soluciones de Neural Machine Translation (NMT), es esencial establecer un enfoque metodológico riguroso que siga directrices claras, tal como señalan Leedy y Ormrod (2013). A menudo, se corre el riesgo de confundir la investigación documental con una mera recopilación de información o con la documentación superficial de hechos. Sin embargo, la investigación debe entenderse como un proceso lógico y consistente de recolección, análisis e interpretación de datos, cuyo objetivo es generar un conocimiento profundo sobre el fenómeno en estudio, en este caso, las soluciones de NMT.

Dado que se busca comparar diversas tecnologías de traducción automática basadas en redes neuronales, el enfoque metodológico debe seguir un proceso sistemático. Este proceso incluye la definición clara del objetivo de la investigación, la gestión adecuada de los datos, y la comunicación precisa de los hallazgos dentro de los marcos establecidos por la comunidad científica. Dichos marcos no solo guiarán sobre qué información incluir, sino también cómo realizar la comparación y qué tipos de conclusiones pueden extraerse de los datos obtenidos.

Para esta investigación, se adoptará un enfoque de investigación aplicada de tipo mixto (cualitativo y cuantitativo). Este enfoque es especialmente adecuado para el análisis de sistemas de NMT, ya que permite combinar el análisis estadístico de la precisión y el rendimiento de las distintas soluciones (enfoque cuantitativo), con una exploración más profunda de aspectos como la usabilidad, adaptabilidad y capacidad de aprendizaje de cada tecnología (enfoque cualitativo). De esta manera, se obtiene una visión más completa del impacto y las fortalezas de cada solución.

Aunque el enfoque mixto será el principal método de investigación, no se excluyen otros enfoques complementarios, como la investigación descriptiva. La investigación descriptiva permitirá observar y documentar el estado actual de cada solución de NMT, identificando características clave y posibles correlaciones entre el rendimiento de los sistemas y los contextos de uso. Este análisis descriptivo es crucial para evaluar de manera objetiva las capacidades de los modelos de traducción automática y su aplicabilidad en distintos escenarios.

En conclusión, siguiendo los lineamientos de una investigación estructurada, se llevará a cabo una comparativa exhaustiva de las soluciones de NMT, utilizando un enfoque mixto que permita una evaluación tanto cuantitativa como cualitativa de los sistemas. Esta metodología asegurará la validez y la utilidad de los hallazgos, proporcionando una base sólida para tomar decisiones informadas sobre qué solución de NMT es la más adecuada según los requerimientos específicos del entorno en estudio.

3.3.2 Fuentes

En el contexto de esta investigación comparativa sobre soluciones de Neural Machine Translation (NMT), las fuentes de información juegan un papel fundamental al proporcionar los datos y la evidencia necesarios para sustentar las evaluaciones y recomendaciones. Las fuentes de información incluyen todos aquellos recursos que contienen datos formales, estudios previos, material multimedia y otros formatos relevantes, que permiten construir una base sólida para el análisis y la argumentación científica. Estas fuentes no solo facilitan la recolección de información crítica, sino que también garantizan que las comparaciones entre las soluciones de NMT se realicen de manera rigurosa y fundamentada.

Las fuentes primarias utilizadas en esta investigación están compuestas por material bibliográfico y documental de alto prestigio y relevancia en el ámbito de la inteligencia artificial,

procesamiento de lenguaje natural y sistemas de traducción automática. Las principales fuentes incluyen:

- IEEE Computer Society Digital Library: Proporciona acceso a una vasta colección de artículos técnicos, investigaciones y conferencias especializadas en informática, inteligencia artificial y aprendizaje automático. Esta fuente es vital para obtener estudios recientes y análisis comparativos de algoritmos de NMT.
- ACM Digital Library: Una de las bibliotecas digitales más importantes en el campo de la informática
 y las ciencias computacionales. Aquí se encuentran artículos que tratan de avances recientes en
 NMT, evaluaciones de modelos y estudios experimentales que son cruciales para la comparativa
 de soluciones.
- Cornell University arXiv: Un repositorio de acceso abierto que ofrece preprints de investigaciones
 en diversas áreas de la ciencia y la tecnología. arXiv es particularmente útil para acceder a
 estudios de vanguardia y trabajos no publicados oficialmente, permitiendo a los investigadores
 mantenerse al tanto de los avances más recientes en traducción automática neuronal sin las
 barreras del acceso por suscripción.

El acceso a estas fuentes es clave para obtener la documentación necesaria. Las dos primeras fuentes mencionadas (IEEE y ACM Digital Library) requieren suscripción y membresía, lo que garantiza el acceso a contenido exclusivo y de alta calidad, lo que resulta crucial para asegurar la profundidad y precisión en el análisis. En cambio, arXiv, como una fuente de acceso libre, complementa este material al ofrecer investigaciones emergentes sin costo alguno. Esta combinación de recursos pagos y gratuitos permite abarcar un espectro amplio de conocimientos, desde los estudios más consolidados hasta las innovaciones más recientes en el campo de la NMT.

En todos los casos, se asegura el respeto a la autoría original de los trabajos consultados, reconociendo y citando adecuadamente a los autores. Esto no solo mantiene la integridad académica de la investigación, sino que también contribuye a construir un cuerpo de conocimiento basado en fuentes confiables y verificadas. Al utilizar estas fuentes primarias, la investigación se basa en datos sólidos y actualizados, lo que permite una comparación más precisa y fundamentada de las diversas soluciones de NMT evaluadas.

3.4 Estudio preliminar del estado del arte

Con el objetivo de proponer un tema de investigación en el área de Neural Machine Translation (NMT), se llevó a cabo una revisión preliminar de la literatura especializada, como parte de la propedéutica necesaria. Esta revisión permitió obtener una comprensión clara del estado actual del conocimiento, identificar las áreas previamente investigadas y detectar vacíos o cuestiones no resueltas que podrían ser abordadas en una nueva investigación. Siguiendo un proceso metódico, la revisión preliminar estableció las bases para el desarrollo de un tema sólido y relevante, facilitando la formulación de un problema de investigación pertinente y de los objetivos que guiarán el estudio.

3.4.1 Pasos realizados para la revisión preliminar de la literatura

3.4.1.1 Definición del tema de interés

En primer lugar, se delimitó claramente el tema de investigación, centrado en la comparación de soluciones de Neural Machine Translation. Se precisó el enfoque hacia la evaluación de aspectos como la precisión de la traducción, el rendimiento en diferentes idiomas y el impacto de las mejoras recientes en las arquitecturas de redes neuronales.

3.4.1.2 Identificación de fuentes de información relevantes

Posteriormente, se identificaron fuentes confiables, actualizadas y pertinentes. Para ello, se consultaron artículos académicos, libros, conferencias especializadas, tesis doctorales y estudios técnicos publicados en plataformas reconocidas como IEEE Computer Society Digital Library, ACM Digital Library y Cornell University arXiv.

3.4.1.3 Realización de la búsqueda bibliográfica

Se efectuó una búsqueda exhaustiva utilizando palabras clave como "Neural Machine Translation", "deep learning en NMT", "comparación de modelos NMT" y "traducción automática". La consulta en bases de datos especializadas como Scopus, Google Scholar y las bibliotecas digitales mencionadas permitió acceder a estudios relevantes y de calidad.

3.4.1.4 Selección de las fuentes más relevantes

Una vez obtenidas las referencias, se filtraron las fuentes mediante la lectura de resúmenes y secciones clave de los artículos, seleccionando aquellos estudios más pertinentes para los fines de la investigación y descartando información redundante o de escasa relevancia.

3.4.1.5 Análisis de la información recopilada

Se procedió al análisis detallado de las fuentes seleccionadas, identificando tendencias en la investigación, vacíos de conocimiento y problemas no resueltos. El análisis se enfocó especialmente en la comparación entre distintas arquitecturas neuronales (Transformers, LSTMs, modelos híbridos), en la efectividad de las mejoras recientes (como los modelos preentrenados) y en los retos asociados a la traducción de idiomas de baja difusión.

3.4.1.6 Formulación del problema de investigación

Con base en el análisis realizado, se formuló el problema de investigación, de manera específica, clara y relevante. En este caso, el problema se centró en la evaluación comparativa de varias soluciones de NMT, considerando su rendimiento en la traducción de textos complejos o multilingües, así como su eficiencia computacional y capacidad de adaptación a distintos contextos lingüísticos.

3.4.1.7 Establecimiento de los objetivos de investigación

Finalmente, se establecieron los objetivos de la investigación, definiéndolos como específicos, medibles y realistas.

La revisión preliminar de la literatura resultó ser un paso esencial para el planteamiento del tema de investigación en el campo de NMT. Gracias a este proceso sistemático, fue posible comprender el estado actual del área, identificar vacíos relevantes y formular preguntas de investigación pertinentes, sentando así una base sólida para el desarrollo de un estudio que aporte nuevos conocimientos a la comunidad académica y científica en el ámbito de la traducción automática neuronal.

3.5 Métodos de investigación

En el contexto de un trabajo universitario que busca comparar soluciones de Neural Machine Translation (NMT), es esencial emplear una variedad de técnicas investigativas formales para obtener resultados robustos y confiables. Estas técnicas permiten evaluar las diferentes bibliotecas y modelos de NMT de manera rigurosa. A continuación, se describen y explican las técnicas más relevantes para realizar una comparativa efectiva:

3.5.1 Experimentos Controlados

Los experimentos controlados son una técnica fundamental en la investigación que implica manipular y controlar variables específicas para evaluar el impacto de estas variables en los resultados obtenidos.

En el caso de la comparativa de soluciones de NMT, se pueden diseñar experimentos donde se varíen factores como la biblioteca de NMT utilizada (por ejemplo, OpenNMT, MarianNMT, o Fairseq), el hardware (diferentes configuraciones de GPU o TPU), y los conjuntos de datos (diversos corpora de traducción). La clave es planificar estos experimentos de manera que cualquier diferencia en el rendimiento se pueda atribuir a la variable manipulada, manteniendo constantes todas las demás. Esto asegura que los resultados sean válidos y fiables.

3.5.2 Estudios de Casos

Los estudios de casos implican la investigación detallada de escenarios específicos en los que se han utilizado diferentes soluciones de NMT en contextos reales. Esta técnica permite analizar cómo diferentes bibliotecas de NMT han sido implementadas en aplicaciones prácticas. Por ejemplo, se pueden examinar casos de uso en empresas de traducción automática, plataformas de servicios en línea, o aplicaciones en la industria de contenido multilingüe. Al recopilar y analizar datos de estos escenarios reales, se obtiene una visión sobre la eficacia y eficiencia de cada biblioteca en condiciones prácticas.

3.5.3 Análisis Comparativo

El análisis comparativo implica la evaluación cuantitativa y cualitativa de diferentes soluciones utilizando métricas específicas para identificar diferencias significativas. Para comparar bibliotecas de NMT, se deben recopilar datos sobre métricas clave como la precisión de la traducción, el tiempo de entrenamiento, y el uso de recursos (como la memoria y el procesamiento). Utilizando estos datos, se realiza una comparación entre las soluciones para determinar cuál ofrece el mejor rendimiento en términos de eficiencia y calidad de las traducciones. Este análisis puede incluir tanto comparaciones numéricas como evaluaciones cualitativas de la calidad de las traducciones generadas.

3.5.4 Revisiones Sistemáticas

Las revisiones sistemáticas son una técnica que implica la recopilación exhaustiva y el análisis crítico de la literatura existente sobre un tema específico. En una revisión sistemática, se busca y se analiza literatura relevante, como artículos científicos, documentos técnicos, y estudios de caso sobre bibliotecas de NMT. Este proceso incluye la búsqueda exhaustiva en bases de datos académicas y técnicas, la aplicación de criterios de inclusión y exclusión para seleccionar estudios relevantes, y el análisis crítico para identificar tendencias actuales, avances tecnológicos y bibliotecas destacadas en el campo.

3.5.5 Validación Cruzada

La validación cruzada es una técnica que evalúa la capacidad de generalización de un modelo al probarlo en diferentes subconjuntos de datos. Para evaluar el rendimiento de las bibliotecas de NMT, se utilizan técnicas de validación cruzada. Esto implica dividir un conjunto de datos en varias particiones, entrenar y validar los modelos en estas particiones, y luego evaluar su desempeño en términos de precisión de traducción y eficiencia. La validación cruzada ayuda a asegurar que los modelos no estén sobreajustados a un único conjunto de datos y puedan generalizar bien a diferentes contextos y tipos de datos.

4 PLANTEAMIENTO DE LA COMPARATIVA

4.1 Criterios de selección de los conjuntos de datos

La elección adecuada de los datos para realizar una comparativa de soluciones de traducción automática neuronal (NMT) es crucial porque los datos tienen un impacto directo en la calidad de la evaluación y en la fiabilidad de los resultados obtenidos. A continuación, se detallan los aspectos clave que destacan la importancia de seleccionar correctamente los conjuntos de datos:

4.1.1 Representatividad de los Datos

Es esencial que los datos utilizados en la comparativa sean representativos del entorno real donde se espera utilizar el modelo NMT. Si los datos no cubren adecuadamente el rango de lenguajes, dialectos, temas y estilos lingüísticos esperados, los resultados de la comparativa no reflejarán el

desempeño real del sistema en producción. Por ejemplo, un modelo entrenado con textos académicos puede no rendir bien en un contexto informal o conversacional.

4.1.2 Diversidad Lingüística

Los datos deben cubrir un amplio espectro de pares de idiomas, especialmente si se evalúan soluciones multilingües. Los diferentes idiomas presentan desafíos únicos, como diferencias en la estructura gramatical, las expresiones idiomáticas, y la complejidad morfológica. Incluir un conjunto diverso de lenguas asegura que el modelo pueda manejar adecuadamente las particularidades de cada una y permite evaluar su versatilidad y robustez.

4.1.3 Variedad en Estilos y Dominios

Es importante que los datos cubran varios estilos de texto (formal, informal, técnico, literario) y dominios (finanzas, medicina, entretenimiento, etc.). Los modelos NMT pueden tener un rendimiento diferente según el tipo de contenido que traduce. Evaluar los modelos solo con datos de un dominio específico puede sesgar los resultados y dar una idea equivocada de su capacidad para generalizar a otros contextos.

4.1.4 Calidad de los Datos de Entrenamiento y Pruebas

La calidad de los datos es clave. Un corpus con traducciones de baja calidad o mal alineadas entre los idiomas afectará tanto el entrenamiento como la evaluación. Para una evaluación precisa de los modelos, es importante contar con conjuntos de prueba que contengan traducciones de alta calidad realizadas por expertos humanos, que sirvan como referencia confiable (gold standard) para medir el rendimiento de las soluciones NMT.

4.1.5 Tamaño del Conjunto de Datos

El tamaño del conjunto de datos de entrenamiento también influye en los resultados. Modelos entrenados con grandes cantidades de datos tienden a tener un mejor rendimiento, ya que tienen acceso a más ejemplos para aprender patrones y estructuras lingüísticas. Sin embargo, para garantizar una comparativa justa, todos los modelos deben entrenarse con conjuntos de datos de tamaños similares o ajustarse mediante técnicas de regularización para evitar ventajas desiguales.

4.1.6 Datos Equilibrados y Balanceados

Es fundamental que los datos estén equilibrados para evitar que los modelos se sesguen hacia ciertos pares de lenguajes, dominios o estilos. Si el conjunto de datos contiene una desproporción en el número de ejemplos entre diferentes idiomas o contextos, los modelos tenderán a rendir mejor en los casos con más ejemplos, lo que falsearía los resultados de la comparativa.

4.1.7 Datos Actualizados

El lenguaje es dinámico y cambia con el tiempo. Es importante usar datos que estén actualizados para reflejar el estado actual de los lenguajes y no quedarse obsoletos. El uso de datos antiguos puede llevar a resultados que no representen correctamente la capacidad del modelo para manejar nuevas terminologías, modismos o cambios en la estructura lingüística.

4.1.8 Conjunto de Pruebas Adecuado

Un conjunto de pruebas adecuado debe ser lo suficientemente diverso y extenso como para evaluar con precisión las capacidades de los modelos en diferentes escenarios. Es crucial evitar la contaminación del conjunto de pruebas (es decir, que los datos de entrenamiento incluyan partes del conjunto de pruebas), ya que esto podría sesgar los resultados y dar una impresión falsa del rendimiento del modelo.

4.1.9 Problemas de sesgo y justicia

Es importante seleccionar los datos de manera que no se introduzcan sesgos que afecten a ciertos grupos o lenguajes. Un conjunto de datos con un sesgo cultural o lingüístico podría llevar a modelos NMT que reflejen esas mismas deficiencias, lo que puede tener implicaciones negativas cuando se implementen en aplicaciones del mundo real.

4.1.10 Tamaño y Complejidad del Texto

El tamaño y la complejidad del conjunto de datos son factores clave. Un conjunto de datos muy sencillo puede no ser un desafío suficiente para evaluar con precisión la capacidad de un modelo. Es esencial que los datos incluyan frases complejas, variaciones en la longitud de las oraciones y estructuras gramaticales desafiantes para que la comparativa refleje el rendimiento del modelo en un espectro amplio de dificultades lingüísticas.

4.2 Corpus paralelo

Un corpus paralelo se diferencia de otros tipos de corpus en que contiene textos que están alineados y son traducciones directas entre dos o más idiomas. Esto significa que cada documento en un corpus paralelo tiene una contraparte directa en otro idioma, lo cual es fundamental para entrenar y evaluar sistemas de traducción automática y para realizar investigaciones comparativas entre idiomas.

Los corpus paralelos son relativamente escasos por varias razones:

- Costo y tiempo de creación: La creación de un corpus paralelo requiere recursos significativos en términos de tiempo, dinero y esfuerzo humano. La recopilación y alineación manual de textos paralelos pueden ser laboriosas y costosas.
- Necesidad de derechos de autor: Obtener los derechos para traducir y distribuir textos en varios idiomas puede ser complicado y costoso. Esto puede limitar la disponibilidad de textos paralelos.
- Calidad y consistencia: Es crucial que las traducciones en un corpus paralelo sean precisas y de alta calidad para que los datos sean útiles en aplicaciones de procesamiento de lenguaje natural.
 Esto implica una cuidadosa revisión y validación de los textos paralelos.
- Diversidad de dominios y géneros: Es deseable que los corpus paralelos cubran una amplia gama de dominios (por ejemplo, noticias, literatura, textos técnicos) y géneros (como diálogos, narrativas, instrucciones), lo cual puede ser difícil de conseguir de manera exhaustiva.

Debido a estas razones, la construcción y disponibilidad de corpus paralelos de alta calidad y tamaño suficiente para diversos propósitos de investigación y desarrollo tecnológico puede ser limitada, lo que contribuye a que sean considerados recursos valiosos pero escasos en el ámbito del procesamiento de lenguaje natural.

4.2.1 UN Parallel Corpus V1

El Corpus Paralelo de las Naciones Unidas (Ziemski et al., 2016) es un recurso lingüístico invaluable compilado por las Naciones Unidas. Consiste en textos paralelos en múltiples idiomas, es decir, textos que son versiones traducidas del mismo contenido.

Este corpus es particularmente útil para tareas como la traducción automática, modelado del lenguaje e investigación cruzada lingüística. Investigadores y desarrolladores lo utilizan

frecuentemente para entrenar y evaluar sistemas de tecnología del lenguaje. Su disponibilidad contribuye significativamente al avance de las capacidades de procesamiento de lenguaje natural multilingüe.

Cuadro No. 4.1. Características del conjunto de datos UN Parallel Corpus V1 (UNPC)

Espacio	Número total de	Vocabulario EN	Vocabulario ES	Número de líneas
requerido:	líneas:	(tokens):	(tokens):	utilizadas:
~3,81 GB	25 227 004	577 671 824	669 899 929	250 000

Fuente: https://opus.nlpl.eu/UNPC/en&es/v1.0/UNPC

4.2.2 OpenSubtitles v2018

El dataset OpenSubtitles v2018 (Lison & Tiedemann, 2016) es una colección de subtítulos de películas que ha sido recopilada y distribuida como un recurso lingüístico ampliamente utilizado en la investigación de procesamiento de lenguaje natural y aprendizaje automático.

OpenSubtitles v2018 contiene subtítulos de películas en múltiples idiomas. Los subtítulos están alineados, lo que significa que hay correspondencias directas entre las líneas de diálogo en diferentes idiomas cuando las películas se han traducido.

Cuadro No. 4.2. Características del conjunto de datos OpenSubtitles (OSPC)

Espacio	Número total de	Vocabulario EN	Vocabulario ES	Número de líneas
requerido:	líneas:	(tokens):	(tokens):	utilizadas:
~3.80 GB	61 434 251	391 976 667	364 279 696	250 000

Fuente: https://opus.nlpl.eu/OpenSubtitles/en&es/v2018/OpenSubtitles

Es uno de los datasets más grandes disponibles públicamente para el entrenamiento y evaluación de modelos de traducción automática y otros sistemas de procesamiento de lenguaje natural. Contiene millones de líneas de diálogo en diversos idiomas.

Debido a su tamaño y variedad lingüística, OpenSubtitles v2018 es utilizado para una variedad de aplicaciones de investigación, incluyendo la traducción automática, la generación de subtítulos automáticos, la investigación en lingüística computacional y más.

El dataset está disponible gratuitamente para la comunidad de investigación, lo cual lo hace accesible para investigadores, desarrolladores y entusiastas que trabajan en el campo del procesamiento de lenguaje natural.

4.2.3 Ventajas y desventajas de los conjuntos de datos seleccionados

En una comparativa de soluciones de traducción automática neuronal (NMT) para el par de lenguas inglés-español (EN-ES), elegir el corpus adecuado es fundamental para obtener resultados precisos y realistas. A continuación, se comparan las ventajas de usar UN Parallel Corpus V1 y OpenSubtitles v2018 en este contexto, destacando cómo cada uno puede beneficiar o limitar el rendimiento del modelo en diferentes aspectos (P. Lin et al., 2024).

4.2.3.1 Contenido y estilo lingüístico

UN Parallel Corpus V1: Este corpus está compuesto principalmente discursos y documentos formales traducidos por expertos humanos. Su contenido se centra en temas de política, derecho internacional, desarrollo y relaciones diplomáticas. Por lo tanto, el estilo es formal, técnico y altamente especializado, lo que lo hace adecuado para entrenar modelos NMT que se utilizarán en entornos formales o técnicos. Entre las ventajas del uso de este conjunto de datos están:

- Ideal para sistemas NMT que se empleen en traducción técnica o formal, como documentos legales, gubernamentales o de negocios.
- Alta calidad de las traducciones, ya que los documentos oficiales son revisados por traductores profesionales.
- El lenguaje es consistente y estandarizado, facilitando la generalización en estos dominios.

Sin embargo, UN Parallel Corpus presente algunas desventajas como las siguientes:

- No es adecuado para tareas de traducción cotidiana o informal debido al estilo altamente especializado y formal.
- Puede no reflejar la variabilidad en el uso del lenguaje cotidiano, lo que afecta su capacidad para manejar textos coloquiales o informales.

Por su parte OpenSubtitles v2018 está compuesto por subtítulos de películas y programas de televisión. Los subtítulos abarcan una gran variedad de géneros, como comedia, drama, acción y

documentales. El estilo de lenguaje es informal y conversacional, con una mezcla de registros coloquiales y expresiones idiomáticas. Entre las ventajas del uso de este conjunto de datos están:

- Excelente para entrenar modelos NMT en contextos informales o conversacionales, como aplicaciones de chat, redes sociales o traducción de contenido audiovisual.
- Mayor diversidad en el vocabulario y estilo, incluyendo coloquialismos, expresiones idiomáticas, e incluso lenguaje técnico en algunos géneros de cine.
- Refleja un uso más variado y actual del idioma, lo que lo hace ideal para aplicaciones de uso general.

Sin embargo, OpenSubtitles v2018 presenta algunas desventajas como las siguientes:

- La calidad de las traducciones puede ser inconsistente, ya que los subtítulos son muchas veces realizados por aficionados o no se ajustan a los estándares de traducción profesional.
- Puede contener errores o simplificaciones, lo que podría afectar la precisión del modelo si se utiliza para tareas formales.

4.2.3.2 Dominio de Aplicación

El conjunto de datos UN Parallel Corpus V1 es adecuado para entrenar modelos orientados a la traducción en entornos especializados como agencias gubernamentales, instituciones internacionales, o empresas que necesitan traducciones precisas y formales. Sin embargo, su uso limita la capacidad del modelo para traducir textos de otros dominios (por ejemplo, literatura, conversación diaria), ya que el vocabulario es más técnico y menos adaptable a contextos coloquiales.

Por su parte OpenSubtitles v2018 permite entrenar modelos que son más versátiles y adaptables a diferentes géneros y registros lingüísticos, como aplicaciones de traducción para consumidores, subtitulado de videos y servicios de traducción en línea. Sin embargo, no es adecuado para aplicaciones que requieren traducción formal o especializada, ya que el estilo del lenguaje es mucho más coloquial y podría no captar matices técnicos o formales con la precisión necesaria.

4.2.3.3 Tamaño del Corpus y Cobertura de Idiomas

En el caso del dataset UN Parallel Corpus V1 aunque su tamaño es más limitado en comparación con OpenSubtitles, ofrece un contenido altamente enfocado en temas oficiales, con traducciones cuidadosas y uniformes en cuanto a calidad. Pero debido a su enfoque especializado, la

cantidad de datos puede ser insuficiente para entrenar modelos grandes que necesiten generalizar en dominios más amplios.

Por otro lado, OpenSubtitles v2018 es significativamente más grande que el UN Parallel Corpus, lo que lo convierte en un recurso ideal para modelos que requieren un gran volumen de datos para aprender patrones diversos. Al abarcar millones de líneas de subtítulos, proporciona una gran variedad de ejemplos y contextos. El volumen de datos no siempre está alineado con la calidad, ya que puede haber ejemplos mal traducidos o inconsistentes.

4.2.3.4 Diversidad Cultural y Contextualización

El dataset UN Parallel Corpus V1 tiene un enfoque cultural y lingüístico más neutral, siendo apropiado para traducciones donde se requiere mantener un tono imparcial y estándar en diferentes contextos internacionales. Puede carecer de la flexibilidad contextual para manejar diferencias culturales en el uso cotidiano del lenguaje.

Una de las mayores ventajas de OpenSubtitles v2018 es la amplia gama de temas, personajes y situaciones en los subtítulos permite que el modelo capture una mayor diversidad cultural y de uso contextual del lenguaje. Sin embargo, esta diversidad también introduce más ruido en los datos, lo que podría desorientar al modelo en tareas que requieran consistencia estilística o formal.

4.2.3.5 Calidad de Traducción

UN Parallel Corpus V1 ofrece traducciones de alta calidad, revisadas por profesionales, lo que garantiza que el modelo NMT se entrene con pares de frases precisas y consistentes. La falta de variabilidad en la calidad de las traducciones puede limitar la exposición del modelo a diferentes estilos o formas de traducción.

OpenSubtitles v2018 contiene una mezcla de traducciones profesionales y no profesionales, lo que ofrece un entrenamiento más variado en términos de calidad. Esto puede hacer que el modelo sea más robusto en el manejo de diferentes niveles de traducción. La inconsistencia en la calidad puede introducir errores, especialmente si los datos no son filtrados adecuadamente.

UN Parallel Corpus V1 es ideal para entrenar y evaluar modelos NMT en contextos formales, técnicos o especializados, donde la precisión y consistencia son cruciales. Sin embargo, su aplicación es limitada en dominios más informales o generales. OpenSubtitles v2018 es más adecuado para

tareas de traducción cotidiana e informal, con un enfoque en el uso coloquial del idioma. Su gran volumen y diversidad ofrecen más flexibilidad, aunque a costa de una menor consistencia en la calidad de las traducciones.

En una comparativa de soluciones NMT enfocada en inglés-español, la elección del corpus dependerá del dominio y el propósito final del sistema de traducción. Idealmente, una combinación de ambos corpora podría ofrecer un balance entre precisión formal y versatilidad informal (P. Lin et al., 2024).

5 DESARROLLO DE LA COMPARATIVA

5.1 Arquitecturas seleccionadas

Como se revisió en la primera parte de la memoria, en el ámbito de Neural Machine Translation (NMT), varias arquitecturas avanzadas han surgido para mejorar el rendimiento y la eficiencia de los modelos de traducción. A continuación, se presenta una selección detallada de estas arquitecturas, destacando sus innovaciones y características clave:

5.1.1 Transformer Vainilla

Los transformers (Vaswani et al., 2017) introdujeron una arquitectura basada en mecanismos de atención, eliminando la necesidad de redes recurrentes (RNN) o convolucionales (CNN) en el procesamiento de secuencias. La atención auto-regresiva permite al modelo considerar todas las palabras en una secuencia de entrada simultáneamente, lo que mejora la capacidad de capturar dependencias a largo plazo. La atención escalada por producto punto, que permite el procesamiento paralelo y mejora la eficiencia en el manejo de secuencias largas.

5.1.2 Universal Transformer

El Universal Transformer (Dehghani et al., 2019) es una extensión del Transformer tradicional que introduce una recurrencia en el mecanismo de atención. Esta arquitectura permite que cada capa de la red realice múltiples pasadas sobre la entrada, ajustando su comportamiento a diferentes niveles de la secuencia.

La incorporación de la recurrencia en el mecanismo de atención, que permite un mejor manejo de secuencias de longitud variable y ajusta el nivel de atención a lo largo de la red.

5.1.3 ST-MoE

El trabajo de Zoph et al., 2022 sobre ST-MoE (Sparse Transformer Mixture of Experts) introduce una arquitectura que utiliza un conjunto de expertos para mejorar la eficiencia computacional. Solo una fracción de los expertos se activa para cada ejemplo, lo que reduce el costo computacional sin sacrificar la calidad de la traducción. La utilización de una mezcla de expertos dispersos que permite un entrenamiento más eficiente y reduce el tiempo de inferencia.

5.1.4 Mamba

Mamba (Gu & Dao, 2024) introduce una arquitectura innovadora de modelado secuencial que supera varias limitaciones fundamentales del Transformer. A diferencia de este, que escala cuadráticamente con la longitud de la secuencia, Mamba logra un rendimiento lineal gracias al uso de modelos de espacios de estado estructurados (SSMs) con mecanismos de selección dependientes del contenido. Esto le permite filtrar dinámicamente información irrelevante y preservar lo esencial según el contexto, sin recurrir a mecanismos de atención ni bloques MLP tradicionales. Además, está optimizado para hardware moderno, lo que se traduce en una inferencia hasta cinco veces más rápida que la de los Transformers. Su diseño simple y homogéneo no solo facilita la implementación, sino que también ofrece un rendimiento competitivo o superior en tareas de lenguaje, audio y genómica, igualando a Transformers del doble de tamaño en calidad de resultados y mejorando la generalización en secuencias extremadamente largas.

5.1.5 Gated State Spaces

Los Gated State Spaces (Mehta et al., 2022) utilizan puertas para controlar el flujo de información en el modelo, combinando conceptos de redes neuronales recurrentes con Transformers. Esta arquitectura busca mejorar el manejo de dependencias temporales y espaciales. La integración de puertas en el espacio de estado, que permite un manejo más flexible y efectivo de las dependencias a largo plazo.

5.1.6 Moving Average Equipped Gated Attention

Ma et al., 2022 introduce una arquitectura con un mecanismo de atención equipado con promedios móviles para mejorar la captura de patrones a largo plazo en las secuencias de entrada. La atención con promedios móviles ayuda a estabilizar el entrenamiento y a mejorar la precisión de la

traducción. El uso de promedios móviles en el mecanismo de atención, que mejora la capacidad del modelo para capturar dependencias a largo plazo y reduce la varianza en el entrenamiento.

Cuadro No. 5.1. Principales innovaciones de los modelos seleccionados para la comparativa.

Modelo	Innovaciones	
Universal	Añade recurrencia en profundidad (repetidas refinaciones por posición),	
Transformer	combinando paralelismo con inductivo recurrente; puede detener	
(UT)	dinámicamente (<i>dynamic halting</i>) el refinamiento por token, mientras que	
	Vainilla aplica una pila fija de capas. (Dehghani et al., 2019)	
Gated State	Sustituye atención global por modelos de espacio de estado más eficientes	
Space (GSS)	(O(L logL) en lugar de O(L\2)), con <i>gating</i> para reducir la dimensionalidad y	
	acelerar FFTs; más eficiente en secuencias largas (Mehta et al., 2022).	
Mamba	Elimina la atención y MLPs, usando modelos de espacio de estado selectivos	
	basados en el input; logra razonamiento dependiente del contenido y es lineal	
	en longitud; <i>Transformer</i> depende de self-attention fijo sin selección dinámica	
	(Gu & Dao, 2024).	
Mega	Reemplaza <i>multi-head attention</i> por una atención de una sola cabeza equipada	
	con promedio móvil exponencial (EMA), introduciendo <i>bias</i> posicional local;	
	ofrece variante de complejidad lineal mediante <i>chunking</i> , mientras <i>Transformer</i>	
	usa self-attention con costos cuadráticos (Ma et al., 2022).	
ST-MoE	Introduce Mixture-of-Experts (MoE) para activar dinámicamente solo partes del	
	modelo, reduciendo cómputo efectivo mientras escala en tamaño total de	
	parámetros; <i>Transformer</i> usa todos sus parámetros en cada paso. Además,	
	aplica técnicas de estabilización y <i>fine-tuning</i> especializadas (Zoph et al., 2022).	

5.2 Implementación y Experimentación

Para comparar estas arquitecturas, se debe implementar y entrenar cada una utilizando un conjunto de datos estándar para NMT. Los conjuntos de datos recomendados incluyen:

• UN Parallel Corpus (Ziemski et al., 2016): Un corpus multilingüe que proporciona datos paralelos en varios idiomas, adecuado para la evaluación de modelos de traducción automática.

 OpenSubtitles V2018 (Lison & Tiedemann, 2016): Un corpus que contiene subtítulos de películas y series en múltiples idiomas, ideal para evaluar la calidad de la traducción en contextos coloquiales y conversacionales.

5.3 Métricas de Evaluación

Las métricas estándar para evaluar la calidad de la traducción automática BLEU (Bilingual Evaluation Understudy) es una métrica que mide la precisión n-grama y compara la salida del modelo con traducciones de referencia.

5.4 Análisis Comparativo

El análisis comparativo debe abarcar:

- Calidad de Traducción: Evaluar la precisión y fluidez de las traducciones generadas por cada arquitectura utilizando métricas como BLEU.
- Tiempo de Entrenamiento: Medir el tiempo necesario para entrenar cada modelo, incluyendo el tiempo de configuración, entrenamiento y ajuste.
- Consumo de Recursos: Evaluar el uso de recursos computacionales, como la memoria y el procesamiento, para cada arquitectura.
- Capacidad de Generalización: Analizar cómo cada arquitectura maneja datos de prueba no vistos durante el entrenamiento y su rendimiento en diferentes conjuntos de datos.

La implementación y evaluación de estas arquitecturas avanzadas de NMT proporcionará una visión completa de sus capacidades y limitaciones. La comparación basada en la calidad de traducción, eficiencia, y recursos permitirá identificar las mejores prácticas y soluciones para tareas específicas de traducción automática. Este análisis detallado contribuirá a la comprensión de cómo las innovaciones recientes en NMT pueden mejorar el rendimiento y la eficiencia en aplicaciones prácticas.

5.5 Configuración del entorno

En una comparativa de soluciones para algoritmos de entrenamiento de redes neuronales, el *hardware* utilizado desempeña un papel fundamental en el rendimiento y la eficiencia del proceso de entrenamiento. El *hardware* adecuado puede marcar una gran diferencia en términos de velocidad de cálculo, capacidad de procesamiento paralelo y manejo eficiente de grandes volúmenes de datos. En

esta sección, describiremos detalladamente el *hardware* utilizado en el entorno de prueba, incluyendo las especificaciones de la CPU, GPU y la cantidad de memoria RAM disponible.

El tipo y modelo de la CPU juegan un papel crucial en el rendimiento general del entrenamiento de redes neuronales. Algunos modelos de CPU ofrecen instrucciones específicas para acelerar las operaciones matemáticas requeridas durante el entrenamiento de redes neuronales, lo que puede resultar en una mejora significativa en el tiempo de entrenamiento. Además, la capacidad de la CPU para manejar múltiples hilos de ejecución puede aprovecharse para realizar cálculos paralelos y agilizar el proceso.

La GPU, por otro lado, es un componente esencial cuando se trata de acelerar el entrenamiento de redes neuronales. Las GPU están diseñadas para realizar operaciones matemáticas de manera altamente paralela, lo que las convierte en una opción ideal para tareas intensivas en cómputo, como el entrenamiento de redes neuronales. La elección del modelo de GPU y su memoria de video disponible pueden influir directamente en la velocidad y capacidad de entrenamiento de la red neuronal.

Además, la cantidad de memoria RAM disponible es un factor crítico por considerar pues el tamaño del conjunto de datos y la arquitectura de la red neuronal determinarán la cantidad de memoria necesaria para almacenar los datos y los parámetros del modelo durante el entrenamiento. Una cantidad insuficiente de memoria RAM puede llevar a problemas de rendimiento, como la ralentización del entrenamiento o incluso la incapacidad de cargar y procesar conjuntos de datos grandes.

En el caso de esta comparativa se utilizó el siguiente *hardware* y *software*:

Cuadro No. 5.2. Configuración del entorno de prueba.

Elemento	Descripción
СРИ	AMD Ryzen 9 5950X 4.9 Ghz
	16 Cores / 32 Hilos
	Caché L3: 64 MB
	Enfriamiento líquido.
GPU	NVIDIA GeForce RTX 3060 12 GB de memoria GDDR6 y 3584 CUDA Cores.

Elemento	Descripción
Memoria	64 GB de Memoria DDR4 (4 módulos)
Almacenamiento	Kingston NV1 NVMe M.2 PCIe Gen3, 1TB.
Sistema operativo	WSL2 sobre Windows 11 Professional (Ubuntu 22.04.2 LST)

Nota: Se utilizaron dos computadoras configuradas idénticamente con el fin de realizar el análisis en el tiempo definido para esta comparativa de soluciones.

5.5.1 Habilitar compatibilidad con GPU

Se consideraron tres formas distintas de configurar *Pytorch* con soporte a tecnologías GPU:

- Sobre una distribución de Linux no virtualizada³.
- Sobre un contenedor de Linux o máquina virtual.
- Sobre el subsistema de Linux en Windows (WSL2)4.

Para esta investigación se toma la decisión de utilizar la opción de WSL2. Particularmente, WSL2 es una característica de Windows que permite a los usuarios ejecutar aplicaciones, contenedores y herramientas de línea de comandos nativas de Linux en Windows 11 y versiones posteriores del sistema operativo. El WSL 2 ofrece beneficios como compatibilidad total con aplicaciones de Linux, un sistema de archivos más rápido y una integración estrecha con el sistema operativo Windows.

El procedimiento descrito en la guía "CUDA on WSL User Guide" explica cómo utilizar la tecnología de NVIDIA CUDA en el Subsistema de Windows para Linux (WSL). La guía se centra en el uso de CUDA en el WSL 2, lo que permite aprovechar la aceleración por GPU de NVIDIA para aplicaciones de ciencia de datos, aprendizaje automático e inferencia en Windows a través del WSL.

A continuación, se resume el procedimiento:

• Instalar el controlador de NVIDIA necesario para habilitar el soporte de GPU,

³ https://docs.nvidia.com/cuda/cuda-installation-guide-linux/index.html#package-manager-metas

⁴ Véase https://docs.nvidia.com/cuda/wsl-user-quide/index.html

 Instalar y configurar el entorno de desarrollo de Linux en WSL. Es necesario crear un archivo de configuración ".wslconfig" para ampliar la cantidad de memoria que utiliza el subsistema Linux:

```
[ws12]
memory=56GB # Limits VM memory in WSL 2 to 56 GB
swap=8GB
```

- Se deben instalar los componentes de software de NVIDIA, como los controladores de Windows, el kit de herramientas de CUDA y las bibliotecas adicionales.
- Existen algunas limitaciones y características no admitidas en el entorno de WSL 2, como el soporte de GPU limitado en ciertos modelos de GPU y la falta de compatibilidad con algunas características de CUDA.
- Instala WSL y elige una distribución basada en glibc, como Ubuntu o Debian.
- Configurar un entorno de Python dentro de WSL. Se recomienda utilizar *Miniconda* de Anaconda para configurar un entorno virtual de Python.
- Una vez que Miniconda esté instalado en WSL, se crea un entorno con el nombre de cada uno
 de los algoritmos que utilizan Python y deben instalarse de forma consistente todos los
 requisitos necesarios a manera de una línea base de configuración.
- Es importante descargar y "actualizar las últimas versiones de los componentes y bibliotecas de tiempo de ejecución con el soporte de *hardware* de la comparativa.

5.5.2 Línea base

La línea base de configuración de un entorno Python se refiere a la configuración inicial necesaria para establecer un entorno de desarrollo y ejecución adecuado para trabajar con Python y sus bibliotecas. A continuación, se enumeran algunos aspectos básicos de la configuración de un entorno Python:

- *Instalación de Python:* Lo primero que se debe hacer es instalar Python en el sistema, para esta comparativa se utiliza la versión estable de Python 3.8.
- Gestor de paquetes: Python cuenta con gestores de paquetes como pip o conda, que facilitan
 la instalación y gestión de bibliotecas y dependencias. Es importante tener un gestor de

paquetes actualizado y configurado correctamente para instalar y mantener las bibliotecas requeridas por el proyecto.

- Entorno virtual: Se recomienda utilizar entornos virtuales para aislar los proyectos y evitar conflictos entre las dependencias. Los entornos virtuales permiten tener instalaciones separadas de Python y sus bibliotecas para cada proyecto.
- Editor de código: Elegir un editor de código o un entorno de desarrollo integrado (IDE) adecuado es importante para desarrollar en Python en este caso se eligió Visual Studio Code por su fácil integración con el entorno WSL2.
- Configuración del entorno de trabajo: Dependiendo de las necesidades específicas del proyecto,
 se pueden configurar variables de entorno, como PATH, que permiten acceder a las herramientas y bibliotecas de Python desde cualquier ubicación del sistema.
- Bibliotecas y dependencias: Para trabajar con diferentes funcionalidades en Python, se necesitarán bibliotecas y dependencias específicas. Estas se pueden instalar utilizando el gestor de paquetes correspondiente (pip o conda) y se pueden especificar en un archivo de requisitos para facilitar la replicación del entorno.

Finalmente, dependiendo de los requisitos de cada algoritmo, pueden existir otras configuraciones adicionales, como configuraciones de parámetros globales del sistema operativo e instalar bibliotecas y componentes de tiempo de ejecución adicionales.

5.6 Aspectos específicos del dominio

En una comparativa de modelos de traducción automática neuronal (NMT, por sus siglas en inglés), es esencial considerar varios aspectos técnicos y arquitectónicos que influyen directamente en la calidad de las traducciones y la eficiencia del entrenamiento de los modelos. A continuación, exploraremos los factores clave que deben analizarse en una comparativa exhaustiva de soluciones NMT, más allá del hardware involucrado.

Las arquitecturas de los modelos NMT han evolucionado desde los enfoques basados en secuencias recurrentes (RNN) hacia redes más avanzadas como los Transformadores. Es importante destacar que los modelos basados en Transformadores, como BERT o GPT, han demostrado un rendimiento superior en la calidad de traducción y en la capacidad de manejar dependencias a largo

plazo. Comparar las diferentes arquitecturas ayuda a evaluar su eficiencia, precisión y capacidad de generalización.

Por su parte, el tamaño del modelo es otro aspecto crítico en la comparativa. Modelos más grandes con mayor número de parámetros tienden a captar mejor las relaciones complejas en los datos, pero requieren más recursos computacionales y memoria. En este sentido, es importante evaluar el trade-off entre el rendimiento de traducción y la cantidad de recursos necesarios para entrenar modelos más grandes.

Es esencial incluir métricas estandarizadas como BLEU (Bilingual Evaluation Understudy), METEOR o TER (Translation Edit Rate) para evaluar la calidad de las traducciones generadas por los diferentes modelos. Estas métricas permiten medir la cercanía entre la traducción automática y una traducción de referencia, proporcionando un indicador objetivo del rendimiento de cada solución NMT.

Además de la precisión, la velocidad con la que un modelo puede ser entrenado y la rapidez con la que produce traducciones son factores fundamentales, especialmente en aplicaciones de traducción en tiempo real. Aquí entra en juego no solo el tipo de hardware utilizado (como se mencionó anteriormente), sino también la optimización del software, los algoritmos de búsqueda utilizados durante la inferencia (por ejemplo, beam search o greedy search) y el uso de técnicas como el paralelismo de datos o de modelo.

En un entorno de comparativa de modelos NMT, es relevante considerar si el sistema puede gestionar eficientemente múltiples pares de lenguajes. Los modelos multilingües, como los de arquitectura basada en Transformadores, permiten entrenar una única red neuronal para varios idiomas, lo cual puede mejorar la eficiencia del sistema y reducir la complejidad, aunque a veces a costa de la precisión en algunos lenguajes específicos.

La capacidad de un modelo para generalizar a ejemplos no vistos durante el entrenamiento es crucial para determinar su viabilidad en entornos del mundo real. Evaluar la robustez de los modelos bajo diferentes condiciones (por ejemplo, presencia de ruido, frases ambiguas o variantes dialectales) es un aspecto clave en una comparativa de soluciones NMT.

Aunque el rendimiento es fundamental, la eficiencia en el uso de los recursos computacionales no debe pasarse por alto. Aquí, la elección del framework (por ejemplo, PyTorch) y

las técnicas de optimización (como el uso de cuantización o pruning) influyen en el rendimiento y en la capacidad de aprovechar mejor el hardware disponible.

Finalmente, la escalabilidad de una solución NMT es crucial en escenarios donde es necesario procesar grandes volúmenes de datos o donde se requiere una actualización frecuente de los modelos. Un buen sistema de NMT debe poder adaptarse a diferentes necesidades de escalabilidad sin sacrificar la calidad o la eficiencia. Incluir las características arquitectónicas, métricas de calidad, velocidad de entrenamiento, eficiencia y escalabilidad permite una evaluación más completa y ajustada a las necesidades del entorno en el que se va a desplegar el sistema.

5.7 Casos de prueba y experimentos

Aunque los conjuntos de datos UN Parallel Corpus V1 y OpenSubtitles V2018 son una referencia pública limitada y es patente que los conjuntos de datos a escala industrial son mucho más grandes en términos de muestras, el rendimiento de los métodos y algoritmos considerados en la comparativa en este conjunto de datos es representativo para aproximar la idoneidad de estos aplicados a problemas del mundo real.

Cuadro No. 5.3. Casos de prueba

Conjunto de datos	Número de ejemplos	Ejemplos para entrenamiento	Ejemplos para validación	Ejemplos para prueba
	250000	175000	50250	24750
UN-Parallel Corpus V1 ¹	(2,19%)	(~1.54%)	(~0.44%)	(~0,22%)
OpenSubtitles V2018	250000	175000	50250	24750
	(0,40%)	(~0.28%)	(~0.08%)	(~0.04%)

Fuente: Comparativa de soluciones, elaboración propia del estudiante.

5.8 Hiperparámetros dependientes de los datos

En la evaluación y comparativa de soluciones de Neural Machine Translation (NMT), los hiperparámetros dependientes de los datos juegan un papel crucial en la configuración y rendimiento del modelo. Estos hiperparámetros afectan directamente cómo el modelo procesa y aprende de los datos, lo que, a su vez, impacta la calidad de la traducción, la eficiencia del entrenamiento y la capacidad de generalización. A continuación, se explica el efecto de cada uno de estos hiperparámetros en una comparativa de NMT:

5.8.1 Longitud Mínima de Secuencia

La longitud mínima de secuencia se refiere al tamaño más corto de las secuencias de entrada y salida que el modelo considera durante el entrenamiento.

- Calidad de Traducción: Si la longitud mínima de secuencia es demasiado corta, el modelo podría ignorar información contextual importante, lo que podría degradar la calidad de la traducción para textos más largos. Por otro lado, si es demasiado larga, puede incluir secuencias irrelevantes que dificulten el entrenamiento.
- Eficiencia del Entrenamiento: Configurar una longitud mínima adecuada ayuda a reducir el ruido y las secuencias demasiado cortas que pueden no aportar suficiente información para el aprendizaje.
- Capacidad de Generalización: Asegurarse de que la longitud mínima se ajuste a la naturaleza del conjunto de datos garantiza que el modelo pueda manejar secuencias de longitud variable y generalizar mejor en diferentes contextos.

5.8.2 Longitud Máxima de Secuencia

La longitud máxima de secuencia define el tamaño máximo de las secuencias que el modelo procesa durante el entrenamiento.

- Calidad de Traducción: Establecer una longitud máxima adecuada asegura que el modelo pueda capturar la información completa de las secuencias largas sin truncar partes importantes del texto. Una longitud máxima demasiado corta puede llevar a la pérdida de contexto y reducir la calidad de la traducción.
- Eficiencia del Entrenamiento: Secuencias más largas requieren más recursos computacionales, lo
 que puede aumentar el tiempo de entrenamiento. Establecer una longitud máxima adecuada
 ayuda a equilibrar la precisión del modelo y la eficiencia del entrenamiento.
- Capacidad de Generalización: Permitir secuencias de longitud variable asegura que el modelo pueda manejar diferentes tipos de datos y contextos, mejorando su capacidad para generalizar a textos de longitud diversa.

5.8.3 Tamaño del Vocabulario

El tamaño del vocabulario en un idioma define el número de palabras o tokens únicos que el modelo utiliza para representar el texto en un idioma.

- Calidad de Traducción: Un vocabulario demasiado pequeño puede causar problemas de cobertura, ya que palabras raras o especializadas podrían ser reemplazadas por un token desconocido (e.g., <UNK>). Un vocabulario más grande puede mejorar la calidad de la traducción al capturar una gama más amplia de términos y matices.
- Eficiencia del Entrenamiento: Vocabularios más grandes aumentan el tamaño del modelo y requieren más memoria y poder de cómputo para el entrenamiento y la inferencia. Un vocabulario muy grande también puede hacer que el modelo sea más lento y menos eficiente.
- Capacidad de Generalización: Un tamaño de vocabulario bien ajustado ayuda al modelo a generalizar mejor al manejar términos y expresiones variados, lo que puede ser especialmente útil en dominios especializados o lenguajes con riqueza léxica.

5.8.4 Batch Size (Tamaño del Lote)

El tamaño del lote se refiere al número de ejemplos de entrenamiento procesados simultáneamente en una sola iteración de entrenamiento.

- Calidad de Traducción: Un tamaño de lote más grande puede llevar a una convergencia más estable y rápida durante el entrenamiento, pero también puede resultar en una menor capacidad de generalización si el modelo no se ajusta adecuadamente a los datos.
- Eficiencia del Entrenamiento: Los tamaños de lote más grandes suelen ser más eficientes desde el punto de vista computacional, ya que permiten un procesamiento paralelo más eficaz en hardware moderno. Sin embargo, también requieren más memoria.
- Capacidad de Generalización: Tamaños de lote más pequeños a menudo proporcionan actualizaciones más ruidosas pero más frecuentes, lo que puede ayudar a evitar el sobreajuste y mejorar la capacidad de generalización del modelo. Ajustar el tamaño del lote es crucial para encontrar un equilibrio entre eficiencia y capacidad de generalización.

Cuadro No. 5.4. Hiperparámetros dependientes de los conjuntos de datos.

Conjunto de datos	Minimum Sequence Length	Maximum Sequence Length	EN Vocabulary Size	ES Vocabulary Size	Batch Size
UN-Parallel Corpus V1 ¹	3	941	50288	68931	8
OpenSubtitles V2018	3	80	39440	53235	8

5.9 Convenciones

Cuando se comparan diferentes arquitecturas Encoder-Decoder para Neural Machine Translation (NMT), es fundamental estandarizar ciertos aspectos del entrenamiento y evaluación para garantizar que los resultados sean comparables y válidos. Las convenciones utilizadas en esta comparativa incluyen optimizadores, schedulers, early stopping, clipping, y medidas de pérdida. A continuación se explica cada una de estas convenciones en detalle:

5.9.1 Optimizador

Se utilizará ADAM con Tasa de Aprendizaje de 0.0001. El optimizador ADAM (*Adaptive Moment Estimation*) es ampliamente utilizado en el entrenamiento de redes neuronales debido a su capacidad para adaptar las tasas de aprendizaje para diferentes parámetros y manejar grandes volúmenes de datos de manera eficiente.

Efecto en la Comparación:

- Consistencia en el Entrenamiento: Utilizar ADAM con una tasa de aprendizaje fija de 0.0001
 asegura que las comparaciones entre arquitecturas sean justas, ya que todas las
 arquitecturas se entrenan bajo las mismas condiciones de optimización.
- Rendimiento del Modelo: La tasa de aprendizaje de 0.0001 es una elección común que permite un entrenamiento estable y controlado, evitando problemas de convergencia demasiado rápida o lenta.

5.9.2 Scheduler

Se configura LROnPlateau que ajusta la tasa de aprendizaje cuando el rendimiento del modelo (por ejemplo, la medida de pérdida en el conjunto de validación) se estabiliza o empeora.

Efecto en la Comparación:

- Adaptación Dinámica: Permite que la tasa de aprendizaje se ajuste en respuesta a la falta de mejora en la pérdida del modelo, lo que puede ayudar a obtener una convergencia más eficiente y efectiva.
- Comparabilidad: Usar el mismo scheduler en todas las arquitecturas asegura que las decisiones de ajuste de tasa de aprendizaje no influyan en las diferencias observadas en el rendimiento entre modelos.

5.9.3 Early Stopping

El *Early stopping* detiene el entrenamiento si la medida de pérdida en el conjunto de validación no mejora durante un número específico de épocas. En este caso, el entrenamiento se detendrá si no hay mejora durante 5 épocas consecutivas.

Efecto en la Comparación:

- Prevención del Sobreentrenamiento: Reduce el riesgo de sobreajuste al detener el entrenamiento cuando no hay mejoras significativas, lo que asegura que los modelos no se entrenen en exceso.
- Equidad en la Comparación: Garantiza que todas las arquitecturas tengan el mismo límite para detener el entrenamiento, evitando que algunas arquitecturas tengan más oportunidades de ajustarse a los datos que otras.

5.9.4 Clipping

El *clipping* de gradientes limita el valor de los gradientes durante el entrenamiento a un máximo especificado para prevenir el problema del desbordamiento de gradientes, que puede ocurrir en redes profundas.

Efecto en la Comparación:

 Estabilidad del Entrenamiento: Evita que los gradientes excesivamente grandes causen problemas de inestabilidad durante el entrenamiento, lo que es crucial para arquitecturas complejas. Consistencia en el Rendimiento: Aplicar el mismo *clipping* a todas las arquitecturas asegura
que los problemas de inestabilidad relacionados con los gradientes no influyan en las
diferencias observadas en el rendimiento de las arquitecturas.

5.9.5 Medida de Pérdida

La medida de pérdida de entropía cruzada (*Cross-Entropy*) es una función de pérdida comúnmente utilizada para tareas de clasificación, incluyendo la traducción automática, que mide la diferencia entre las distribuciones de probabilidad predicha y real.

Efecto en la Comparación:

- Relevancia de la Medida: Cross-Entropy es una medida estándar para problemas de clasificación y es adecuada para evaluar la precisión de las predicciones en el contexto de NMT.
- Comparabilidad Directa: Utilizar Cross-Entropy como medida de pérdida permite comparar directamente el rendimiento de diferentes arquitecturas en términos de su capacidad para minimizar la discrepancia entre las predicciones y las traducciones reales.

Estas convenciones establecen una base uniforme para comparar diferentes arquitecturas Encoder-Decoder en NMT, asegurando que las diferencias en el rendimiento sean atribuibles a las características de las arquitecturas mismas y no a variaciones en el proceso de entrenamiento.

6 DISCUSIÓN Y ANÁLISIS DE RESULTADOS

En esta sección se discuten los resultados de los cuatro algoritmos de entrenamiento de redes neuronales considerados. Es importante destacar que se establece una comparativa entre el tiempo total del entrenamiento, el tamaño del modelo en memoria, la exactitud de los modelos (entrenamiento y validación) y la cantidad de parámetros entrenables de los modelos.

6.1 Cantidad de Parámetros

El tamaño del modelo es un factor crítico para considerar en escenarios donde el almacenamiento de recursos es limitado, como en dispositivos con capacidades de almacenamiento limitadas o entornos con restricciones de memoria. Un modelo con un tamaño más grande requerirá más recursos para su almacenamiento y carga en memoria durante el despliegue y la ejecución.

El análisis es relevante porque proporciona información sobre el tamaño promedio del modelo para cada algoritmo evaluado. El tamaño del modelo puede tener un impacto significativo en la practicidad y eficiencia de su implementación en diferentes entornos.

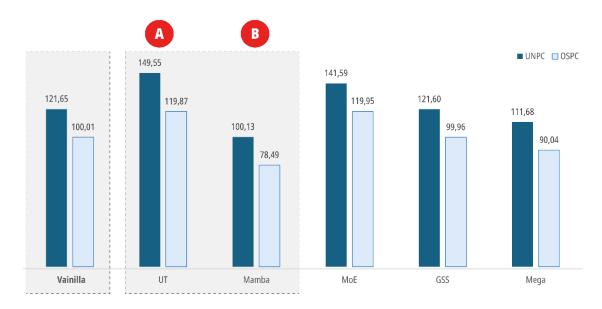
Cuadro No. 6.1. Cantidad de parámetros por modelo y conjunto de datos.

Modelo	UNPC	OSPC
Vainilla	121,65	100,01
UT	149,55	119,87
Mamba	100,13	78,49
MoE	141,59	119,95
GSS	121,60	99,96
Mega	111,68	90,04

Nota: Cifras expresadas en millones de parámetros.

Fuente: Comparativa de soluciones, elaboración propia del estudiante.

Gráfico No. 6.1. Número de parámetros entrenables de los modelos de la comparativa.



Fuente: Comparativa de soluciones, elaboración propia del estudiante.

La figura expone el número de **parámetros entrenables** de cada una de las arquitecturas evaluadas, expresado en millones (MM), y muestra cómo distintos enfoques afectan la **escala** del modelo, tanto en el corpus **UNPC** como en **OSPC**. Esta variable es fundamental en el análisis de modelos, ya que está directamente relacionada con la capacidad de generalización, el consumo de memoria y el coste computacional durante el entrenamiento y la inferencia.

En el corpus **UNPC**, el modelo **Universal Transformer (UT)** es el que presenta el mayor número de parámetros, con **149,55 millones**, lo que representa un incremento del **23,6 %** respecto

al modelo **Vainilla**, que se toma como referencia con **121,65 millones**. Este incremento, resaltado en el bloque "A", indica que ciertos enfoques más complejos pueden conllevar un aumento significativo en la carga computacional sin garantizar necesariamente una mejora proporcional en el rendimiento. Por el contrario, **Mamba** destaca por su eficiencia estructural, con **100,13 millones** de parámetros en UNPC, lo que implica una **reducción del 21 %** respecto a Vainilla, como se indica en el bloque "B". Esta tendencia se mantiene también en OSPC, donde Mamba nuevamente presenta la menor cantidad de parámetros (**78,49 MM**) frente a los demás modelos.

Modelos como **MoE** y **GSS** también muestran cifras más contenidas que UT, pero superiores a Mamba, manteniéndose entre **99,96 MM y 121,60 MM** según el corpus. El modelo **MEGA**, aunque algo más compacto que UT, todavía supera los **111 millones** de parámetros en UNPC.

Este análisis refuerza que Mamba no solo ofrece tiempos reducidos de entrenamiento y validación, sino que también requiere una menor cantidad de parámetros para alcanzar buenos resultados, lo que lo convierte en una opción atractiva para entornos con restricciones de memoria o cómputo. En contextos donde se busca escalar el modelo o realizar despliegues optimizados, esta característica puede ser determinante. La capacidad de lograr una alta eficiencia con menor complejidad estructural consolida a Mamba como una de las arquitecturas más balanceadas del estudio.

6.2 Tiempo de entrenamiento y validación

Tener un tiempo de ejecución más bajo puede ser beneficioso en situaciones en las que se requiere una respuesta rápida o cuando hay restricciones de tiempo. Esto puede ser especialmente relevante en aplicaciones en tiempo real o en entornos donde se necesitan respuestas rápidas a consultas o predicciones.

Este análisis es importante porque proporciona información sobre el tiempo promedio de ejecución de cada algoritmo evaluado. Esta información es relevante ya que el tiempo de ejecución puede tener un impacto significativo en la practicidad y eficiencia de un algoritmo en un entorno de producción.

Las gráficas presentadas a continuación muestran el tiempo total de entrenamiento y validación requerido por distintas arquitecturas de modelos de traducción automática en dos conjuntos de datos: **UNPC** (United Nations Parallel Corpus) y **OSPC** (OpenSubtitles Parallel Corpus).

El análisis revela diferencias significativas en la eficiencia computacional de las arquitecturas evaluadas, tanto en entornos con lenguaje formal (UNPC) como en lenguaje más informal y fragmentado (OSPC).

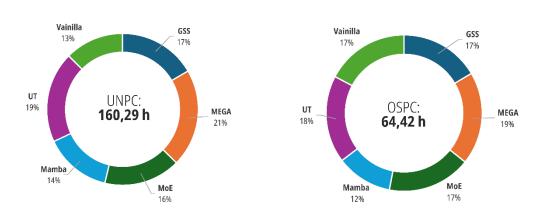


Gráfico No. 6.2. Contribución relativa al tiempo total de ejecución de la comparativa

Nota: Tiempo promedio dado en horas tras efectuar 20 epochs. Se utiliza UNPC para designar al conjunto de datos Corpus Paralelo de las Naciones Unidas (Ziemski et al., 2016) y OSPC para designar a OpenSubtitles v2018 (Lison & Tiedemann, 2016).

Fuente: Elaboración propia del estudiante.

En el conjunto **UNPC**, el tiempo total de entrenamiento y validación asciende a **160,29 horas** al considerar 20 épocas, siendo la arquitectura **MEGA** la que representa el mayor consumo de tiempo, con un **21 %** del total. Le siguen **UT** (Universal Transformer) con un 19 %, y **GSS** (Gated State Space) con un 17 %. **Mamba**, por su parte, destaca por su eficiencia relativa al ocupar solo un **14 %** del tiempo total, superando incluso a la arquitectura "Vainilla" (Transformer base), que ocupa un 13 %. Este resultado es especialmente relevante dado que Mamba presenta mejoras estructurales que permiten una mayor capacidad de modelado con un coste computacional más contenido. La arquitectura **MoE** (Mixture of Experts) también mantiene un perfil equilibrado con un 16 % del tiempo total.

En el conjunto **OSPC**, que se caracteriza por secuencias más breves y menos estructuradas, se observa una reducción significativa en el tiempo total, con un promedio de **64,42 horas**. Nuevamente, **MEGA** encabeza el uso de tiempo con un **19 %**, seguido de cerca por **UT** (18 %) y **MoE** y **GSS** (ambos con 17 %). En este caso, **Mamba** alcanza su mejor rendimiento relativo, representando solo un **12 %** del total, lo que confirma su alta eficiencia incluso en contextos menos formales y con datos más ruidosos.

El modelo Vainilla ocupa un 17%, lo que refuerza la ventaja de Mamba frente a la arquitectura tradicional.

Estos resultados subrayan la importancia de considerar el tiempo total de entrenamiento y validación como una métrica clave al comparar arquitecturas de modelos. Más allá de la calidad de traducción obtenida, la eficiencia en el uso de recursos computacionales puede ser decisiva en contextos reales de despliegue, especialmente cuando se trata de escalabilidad, costos de infraestructura y tiempos de iteración en el desarrollo de modelos.

6.2.1 Tiempo de entrenamiento

Al comparar arquitecturas como Mamba, Transformer, Universal Transformer (UT), MEGA y otras, es fundamental considerar no solo su precisión o capacidad de generalización, sino también el tiempo total que requieren para ser entrenadas. Este factor es clave porque refleja la eficiencia computacional del modelo en condiciones reales de uso, donde los recursos (tiempo, energía y hardware) son limitados. Dos modelos con rendimientos similares pueden diferir significativamente en los costos asociados a su entrenamiento si uno requiere muchas más horas de cómputo o mayor consumo de memoria.

Además, el tiempo de entrenamiento incide directamente en la viabilidad de escalar el modelo a conjuntos de datos más grandes o realizar iteraciones frecuentes durante el desarrollo. Por ello, el tiempo total de entrenamiento es una métrica crítica para evaluar el balance entre rendimiento y eficiencia de cada arquitectura.

Cuadro No. 6.2. Tiempo total de la fase de entrenamiento por modelo

Modelo	UNPC	OSPC
GSS	22,84	9,34
MEGA	29,15	10,93
MoE	22,62	9,77
Mamba	19,50	6,35
UT	27,28	10,64
Vainilla	17,49	9,86

Nota: Tiempo de entrenamiento en horas.

29,15

22,84

29,15

21,28

22,62

19,50

17,49

10,93

9,77

6,35

10,64

9,86

GSS

MEGA

MoE

Gráfico No. 6.3. Tiempo total de la fase de entrenamiento por modelo

Nota: Tiempo promedio dado en horas tras efectuar 20 epochs. Se utiliza UNPC para designar al conjunto de datos Corpus Paralelo de las Naciones Unidas (Ziemski et al., 2016) y OSPC para designar a OpenSubtitles v2018 (Lison & Tiedemann, 2016). Fuente: Elaboración propia del estudiante.

Mamba

Vainilla

El cuadro y el gráfico presentados muestran el **tiempo de entrenamiento total** requerido por cada una de las arquitecturas evaluadas, en los conjuntos de datos **UNPC** (Corpus Paralelo de Naciones Unidas) y **OSPC** (OpenSubtitles), tras 20 épocas. Esta métrica permite observar con mayor precisión la eficiencia computacional específica de cada modelo, independientemente del tiempo destinado a la validación.

En el conjunto UNPC, caracterizado por frases más estructuradas y largas, se observa que MEGA y UT son las arquitecturas que presentan los tiempos de entrenamiento más elevados, con 29,15 horas y 27,28 horas respectivamente. Les siguen MoE (22,62 horas) y GSS (22,84 horas), mientras que Mamba se posiciona como una de las arquitecturas más eficientes, con 19,50 horas, superando incluso al modelo Vainilla, que requiere 17,49 horas. Este resultado es destacable, dado que Mamba ofrece una arquitectura más moderna y compleja que la Transformer básica, pero con un coste computacional contenido. El bloque señalado como "A" en la gráfica destaca precisamente este equilibrio entre rendimiento y eficiencia de Mamba y MoE.

En el conjunto OSPC, que representa diálogos más breves, informales y fragmentados, todos los modelos muestran una reducción significativa en los tiempos de entrenamiento. Aquí, Mamba se consolida como la opción más eficiente con 6,35 horas, seguido de GSS (9,34 horas) y MoE (9,77 horas). Incluso en comparación con el modelo Vainilla, que registra 9,86 horas, Mamba demuestra una ventaja clara, lo que sugiere que su capacidad de modelado eficiente se adapta especialmente

bien a secuencias cortas y variadas. Esta diferencia se refuerza en el bloque "B" de la gráfica, donde se evidencia que incluso la arquitectura tradicional requiere más tiempo que Mamba en este entorno.

En conjunto, estos resultados demuestran que Mamba no solo es competitivo en calidad, sino también en eficiencia computacional, destacando especialmente en tareas de procesamiento de secuencias cortas como las contenidas en OSPC. Esta ventaja lo posiciona como una arquitectura particularmente adecuada para aplicaciones que requieren rapidez de entrenamiento y despliegue continuo, como traducción automática en tiempo real o sistemas conversacionales.

6.2.2 Tiempo de validación

En una comparativa de soluciones también es esencial considerar el tiempo de validación, ya que este influye directamente en la eficiencia del ciclo de entrenamiento y evaluación. La validación se realiza periódicamente para monitorear el desempeño del modelo en datos no vistos y tomar decisiones sobre ajustes de hiperparámetros, puntos de parada o selección de modelos. Si una arquitectura requiere largos tiempos de validación, puede ralentizar significativamente el proceso de desarrollo y dificultar experimentaciones rápidas.

Además, en entornos donde se evalúan múltiples *checkpoints* o se aplica validación cruzada, estos costos se multiplican. Por tanto, medir el tiempo de validación permite valorar no solo la precisión del modelo, sino también su practicidad y eficiencia operativa en flujos de trabajo reales.

Cuadro No. 6.3. Tiempo total de la fase de validación de los modelos

Modelo	UNPC	OSPC
GSS	3,73	1,31
MEGA	4,28	1,44
MoE	3,57	1,37
Mamba	3,40	1,12
UT	3,75	1,01
Vainilla	2,67	1,27

Nota: Tiempo de entrenamiento en horas.

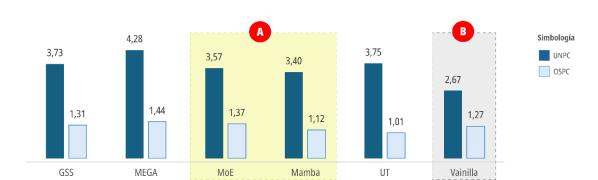


Gráfico No. 6.4. Tiempo total de la fase de validación de los modelos

Nota: Tiempo promedio dado en horas tras efectuar 20 epochs. Se utiliza UNPC para designar al conjunto de datos Corpus Paralelo de las Naciones Unidas (Ziemski et al., 2016) y OSPC para designar a OpenSubtitles v2018 (Lison & Tiedemann, 2016). Fuente: Elaboración propia del estudiante.

La figura anterior muestra los **tiempos promedio de validación**, expresados en horas tras 20 épocas, para cada una de las arquitecturas evaluadas en los conjuntos de datos **UNPC** y **OSPC**. Esta métrica resulta clave, ya que la validación se ejecuta de forma periódica durante el entrenamiento para monitorear el desempeño del modelo en datos no vistos, impactando directamente en la agilidad del desarrollo y ajuste de modelos.

En el conjunto UNPC, las arquitecturas que presentan mayores tiempos de validación son MEGA (4,28 horas), UT (3,75 horas) y GSS (3,73 horas), lo que refleja un uso intensivo de recursos computacionales en modelos que, si bien pueden ofrecer altos niveles de precisión, lo hacen a un costo operativo más elevado. En este contexto, Mamba vuelve a destacar con un tiempo de validación de 3,40 horas, que junto con MoE (3,57 horas), conforman el grupo resaltado en el bloque "A" de la gráfica. Estas dos arquitecturas logran un balance más eficiente entre rendimiento y consumo computacional. Por otro lado, el modelo Vainilla muestra el menor tiempo de validación en este conjunto (2,67 horas), como se indica en el bloque "B", aunque esto puede estar asociado a su simplicidad estructural frente a modelos más avanzados.

En el conjunto OSPC, todos los modelos presentan una disminución significativa en los tiempos de validación, lo cual es coherente con la naturaleza más breve y fragmentada de las secuencias. Mamba se posiciona nuevamente como una de las opciones más eficientes con 1,12 h, superando incluso al modelo Vainilla (1,27 horas) y destacando por su bajo consumo computacional en validación. Los modelos MEGA (1,44 horas) y GSS (1,31 horas) mantienen valores más elevados, lo

que sugiere que sus ventajas arquitectónicas no se traducen necesariamente en una validación ágil, especialmente en tareas de menor complejidad lingüística.

En conjunto, estos resultados refuerzan la posición de Mamba como una arquitectura eficiente no solo en entrenamiento, sino también en validación, lo que facilita ciclos de iteración más rápidos y sostenibles en entornos de desarrollo. Esta propiedad resulta especialmente útil cuando se requiere ajustar modelos de forma continua o trabajar con pipelines automatizados que dependen de evaluaciones frecuentes.

A partir del análisis conjunto de los tiempos de entrenamiento, validación y tiempo total, es posible extraer conclusiones relevantes sobre la eficiencia computacional de las distintas arquitecturas evaluadas. Si bien algunas como MEGA y UT destacan por su capacidad de modelado y rendimiento en tareas complejas, presentan un coste computacional significativamente más elevado tanto en entrenamiento como en validación, lo cual puede representar una barrera práctica en escenarios donde los recursos son limitados o se requiere una rápida iteración de modelos.

En contraste, Mamba se posiciona como una de las soluciones más equilibradas, logrando tiempos de entrenamiento y validación consistentemente bajos sin comprometer el desempeño esperado en tareas de traducción automática. Su eficiencia es particularmente notable en el conjunto OSPC, donde supera incluso al modelo Vainilla (*transformer* clásico) en términos de consumo computacional, a pesar de tratarse de una arquitectura más moderna y con mayor capacidad expresiva. Esta ventaja se mantiene también en UNPC, un corpus más exigente, donde Mamba logra tiempos inferiores a los de la mayoría de modelos avanzados, manteniéndose como una alternativa altamente competitiva.

Por tanto, considerando tanto la calidad del modelo como su eficiencia operativa, Mamba emerge como una arquitectura especialmente adecuada para entornos de despliegue real, investigación aplicada o desarrollos iterativos frecuentes, donde el balance entre rendimiento y recursos computacionales es un factor decisivo. Esta combinación de eficacia y eficiencia lo convierte en un candidato sólido para ser adoptado como base en soluciones de traducción automática y tareas secuenciales multilingües.

6.3 Consumo de memoria

Medir la memoria que consume cada modelo en una comparativa de soluciones es crucial porque permite evaluar su viabilidad práctica en diferentes entornos de ejecución. Un modelo que consume excesiva memoria puede requerir hardware especializado o limitar el tamaño del lote (*batch size*), afectando directamente la velocidad de entrenamiento y la eficiencia durante la inferencia. Además, en escenarios donde se trabaja con secuencias largas o se desea escalar el modelo, el uso de memoria se convierte en un factor determinante para evitar errores por falta de recursos (*out-of-memory*) y optimizar el aprovechamiento del hardware disponible. Comparar el consumo de memoria también ayuda a identificar cuellos de botella y a seleccionar arquitecturas más adecuadas según las restricciones técnicas o económicas del proyecto, haciendo más realista y sostenible su implementación.

Cuadro No. 6.4. Total de memoria consumida por modelo

Modelo	UNPC	OSPC
GSS	2800,62	476,31
MEGA	2799,12	476,46
МоЕ	2796,63	476,23
Mamba	2798,63	476,19
UT	2799,55	476,08
Vainilla	2800,30	476,18

Fuente: Comparativa de soluciones, elaboración propia del estudiante.

Gráfico No. 6.5. Total de memoria consumida por modelo (en Megabytes)



La información anterior sobre consumo de memoria presenta el promedio de uso en megabytes por lote (*batch*) tras 20 ejecuciones para los modelos evaluados en los conjuntos UNPC y OSPC. Esta métrica es fundamental, ya que el uso de memoria condiciona el tamaño del modelo que puede ser cargado en una unidad de procesamiento (CPU o GPU) y la cantidad de datos que se pueden procesar simultáneamente, afectando directamente la escalabilidad y la velocidad del entrenamiento.

En el corpus UNPC, todos los modelos presentan consumos muy similares, oscilando entre 2796,63 MB y 2800,62 MB, sin diferencias marcadas entre arquitecturas complejas como UT o MEGA, y modelos más eficientes como Mamba o MoE. Este comportamiento, resaltado en el bloque "A", sugiere que la elección de modelo en este corpus no impacta significativamente el consumo de memoria por lote, lo cual es una buena noticia para la interoperabilidad de arquitecturas bajo los mismos recursos hardware.

El patrón se repite en el corpus OSPC, donde el consumo se mantiene prácticamente constante para todos los modelos, variando en un rango extremadamente estrecho: de 476,08 MB a 476,46 MB. Este comportamiento se refleja en el bloque "B" de la figura. Tal uniformidad indica que, al menos en lo que respecta a uso de memoria, no hay un modelo que represente una ventaja o desventaja clara frente a los demás, lo que permite tomar decisiones basadas en otras métricas como precisión, eficiencia temporal o número de parámetros.

En conjunto, estos resultados permiten concluir que el consumo de memoria no es un factor diferenciador entre las arquitecturas analizadas, al menos en los escenarios y configuraciones evaluadas. Esta homogeneidad en el uso de memoria refuerza la importancia de evaluar otros aspectos como la calidad de salida, el número de parámetros o los tiempos de entrenamiento y validación al momento de seleccionar una arquitectura para implementación práctica.

6.4 Pérdida

En una comparativa de soluciones para modelos de traducción automática, el análisis de la pérdida (*loss*) tanto en el conjunto de entrenamiento como en el de validación es crucial porque permite comprender cómo está aprendiendo el modelo y detectar problemas en el proceso de optimización que no siempre se evidencian con métricas de precisión o calidad final, como BLEU.

La pérdida mide el error cometido por el modelo al predecir la secuencia esperada, es decir, cuánto se desvía su salida con respecto a la traducción correcta. Una pérdida baja indica que el modelo está produciendo traducciones cercanas a las esperadas. Por tanto, analizar la pérdida durante el entrenamiento ayuda a evaluar si el modelo está aprendiendo progresivamente y si la optimización es estable y efectiva.

Por otro lado, la pérdida en el conjunto de validación refleja si ese aprendizaje se generaliza adecuadamente a datos no vistos. Un modelo con pérdida baja en entrenamiento pero alta en validación probablemente está sobreajustado: ha memorizado patrones específicos en lugar de aprender reglas generales del lenguaje. Inversamente, si la pérdida es alta en ambos conjuntos, puede haber un problema de subajuste (*underfitting*), indicando que la arquitectura o el entrenamiento no están capturando la complejidad del problema.

Comparar la pérdida en ambos conjuntos también permite evaluar la eficiencia del entrenamiento y anticipar posibles mejoras en hiperparámetros, regularización o arquitectura. Además, el análisis de la evolución de la pérdida entre épocas proporciona información temprana sobre la tasa de convergencia del modelo y su comportamiento en fases iniciales, lo cual es esencial para acelerar ciclos de desarrollo y evitar entrenamientos innecesariamente largos.

6.4.1 Fase de entrenamiento

A continuación se presentan los resultados del análisis de la **pérdida (loss)** obtenida por cada una de las arquitecturas evaluadas durante el proceso de entrenamiento en el conjunto de datos.

En este apartado, se exponen y comparan los valores iniciales y finales de pérdida para cada arquitectura, así como la variación porcentual entre ambos puntos, con el fin de identificar cuáles modelos muestran un aprendizaje más efectivo y sostenido a lo largo del proceso. Esta evaluación será fundamental para comprender en qué medida la complejidad de cada arquitectura se traduce en un mejor ajuste a los datos durante el entrenamiento.

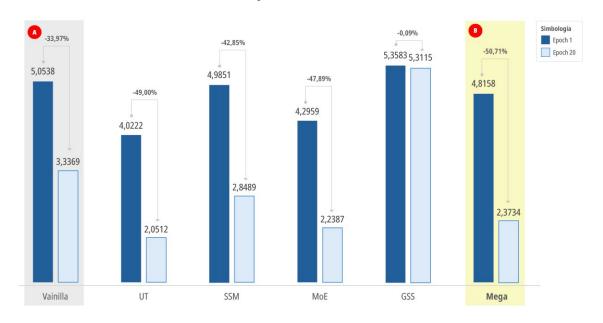
6.4.1.1 Conjunto de datos UNPC

Cuadro No. 6.5. Comparación de pérdida inicial y final durante la fase de entrenamiento en el conjunto de datos UNPC

Modelo	Epoch 1	Epoch 20	Variación
Vainilla	5,0538	3,3369	-33,97%

UT	4,0222	2,0512	-49,00%
Mamba	4,9851	2,8489	-42,85%
MoE	4,2959	2,2387	-47,89%
GSS	5,3583	5,3115	-0,87%
Mega	4,8158	2,3734	-50,72%

Gráfico No. 6.6. Comparación de pérdida inicial y final durante la fase de entrenamiento en el conjunto de datos UNPC



Fuente: Comparativa de soluciones, elaboración propia del estudiante.

El **Cuadro No. 6.5 y el Gráfico No. 6.6** muestran la comparación entre la pérdida inicial (epoch 1) y final (epoch 20) de cada modelo entrenado sobre el conjunto de datos UNPC, lo que permite valorar la capacidad de aprendizaje de cada arquitectura y su eficiencia durante la fase de entrenamiento. La métrica de variación porcentual refleja la reducción relativa de la pérdida a lo largo del proceso.

El modelo MEGA destaca por lograr la mayor reducción de pérdida, con un -50,72 %, seguido de cerca por UT (-49,00 %) y MoE (-47,89 %). Estos resultados indican que dichas arquitecturas tienen un alto potencial de aprendizaje en contextos de lenguaje formal como el de UNPC. Mamba, aunque no alcanza el descenso más pronunciado, muestra una mejora sustancial con un -42,85 %, confirmando su solidez como arquitectura eficiente tanto en desempeño como en velocidad de convergencia.

Por otro lado, el modelo Vainilla obtiene una reducción del -33,97 %, lo que indica un desempeño aceptable, aunque inferior al de las arquitecturas más recientes. El caso más llamativo es GSS, que apenas mejora su pérdida inicial, con una variación de solo -0,87 %, lo que sugiere posibles problemas de ajuste o limitaciones en su capacidad para adaptarse a este corpus específico.

Este análisis refuerza la idea de que modelos como MEGA, UT, MoE y Mamba presentan una mejor capacidad de aprendizaje durante las primeras 20 épocas. No obstante, al combinar esta métrica con los resultados de eficiencia computacional previamente analizados, Mamba se posiciona como una opción especialmente equilibrada, ofreciendo una buena capacidad de reducción de pérdida con menor consumo de recursos en comparación con otros modelos más costosos computacionalmente.

Simbología
5,50
4,50
4,50
3,50
3,50
2,50
1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20

Gráfico No. 6.7. Comportamiento de la pérdida de los modelos durante el entrenamiento para el conjunto de datos UNPC.

Fuente: Comparativa de soluciones, elaboración propia del estudiante.

La gráfica anterior muestra las curvas de pérdida durante el entrenamiento en el conjunto UNPC, permitiendo visualizar el comportamiento de cada arquitectura a lo largo de las 20 épocas analizadas. Esta representación es clave para entender la dinámica del aprendizaje y la eficiencia en la convergencia de los modelos evaluados.

Desde el inicio del entrenamiento, se observa que la mayoría de los modelos reducen su pérdida de forma progresiva, con diferencias claras en la velocidad de descenso y en el punto de estabilización. El modelo UT destaca por su curva más descendente y estable, alcanzando los valores de pérdida más bajos al finalizar el entrenamiento. Le siguen MEGA, MoE y Mamba, que también presentan un descenso sostenido y eficiente, lo que indica una capacidad efectiva de ajuste a los datos sin generar inestabilidad.

Por su parte, Vainilla muestra una reducción más lenta, con una curva más plana a partir de la mitad del entrenamiento, lo que sugiere una convergencia limitada. Sin embargo, el caso más crítico es el de GSS, cuya curva se mantiene elevada y presenta incluso oscilaciones leves a lo largo de las épocas, reflejando dificultades de aprendizaje o inestabilidad en el proceso de optimización.

En conjunto, esta visualización confirma que modelos como UT, Mamba y MEGA logran un descenso rápido y sostenido en la pérdida, lo que evidencia una mayor eficacia en el proceso de aprendizaje desde etapas tempranas. Este tipo de análisis resulta esencial para evaluar no solo el resultado final del entrenamiento, sino también la calidad y estabilidad del aprendizaje a lo largo del tiempo.

6.4.1.2 Conjunto de datos OSPC

Cuadro No. 6.6. Comparación de pérdida inicial y final durante la fase de entrenamiento en el conjunto de datos OSPC

Modelo	Epoch 1	Epoch 20	Variación
Vainilla	4,2428	4,05360	-4,46%
UT	3,9671	3,15560	-20,46%
SSM	4,6662	3,53160	-24,32%
MoE	4,676	3,21380	-31,27%
GSS	5,0914	4,42900	-13,01%
Mega	4,6465	4,30030	-7,45%

Gráfico No. 6.8. Comparación de pérdida inicial y final durante la fase de entrenamiento en el conjunto de datos OSPC



El **Cuadro No. 6.6** presenta la comparación entre la pérdida inicial y final de cada modelo durante la fase de entrenamiento en el conjunto de datos OSPC, permitiendo evaluar la capacidad de ajuste de las arquitecturas en un entorno de lenguaje más informal y fragmentado. En este caso, se observa una menor magnitud de mejora generalizada en comparación con el corpus UNPC, lo que puede deberse a la mayor variabilidad y ruido presente en los datos.

La arquitectura MoE es la que muestra la mayor mejora relativa, con una reducción del -31,27 %, seguida de SSM (Mamba) con -24,32 %, y UT con -20,46 %. Estos modelos demuestran una buena capacidad de aprendizaje incluso en contextos menos estructurados. En particular, Mamba vuelve a destacar al lograr una mejora sólida sin requerir una gran cantidad de parámetros ni incurrir en altos tiempos de entrenamiento o validación, lo que confirma su consistencia como modelo eficiente.

Por otro lado, GSS y MEGA presentan mejoras más modestas, con reducciones del -13,01 % y -7,45 % respectivamente. El modelo Vainilla muestra el peor desempeño, con apenas un -4,46 % de reducción en la pérdida, lo que sugiere una mayor dificultad para adaptarse a las características lingüísticas del corpus OSPC.

Estos resultados reflejan que, aunque todos los modelos logran alguna mejora durante el entrenamiento, existen diferencias claras en su capacidad de aprendizaje en contextos más ruidosos.

MoE y Mamba sobresalen como las opciones más adaptables en este entorno, y Mamba, en particular, mantiene su perfil equilibrado al combinar buena reducción de pérdida con eficiencia computacional.

Simbología

- Vainilla

- UT

- Mamba

- MoE

- GSS

- Mega

1,50

1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20

Gráfico No. 6.9. Comportamiento de la pérdida de los modelos durante el entrenamiento para el conjunto de datos OSPC.

Fuente: Comparativa de soluciones, elaboración propia del estudiante.

La gráfica de curvas de pérdida durante el entrenamiento en el conjunto OSPC permite observar el comportamiento de aprendizaje de cada modelo en un entorno caracterizado por frases cortas, lenguaje informal y mayor variabilidad lingüística. En este contexto, la estabilidad y la capacidad de adaptación del modelo cobran especial relevancia, ya que los datos son menos estructurados que en el corpus UNPC.

Desde las primeras épocas, el modelo UT muestra un descenso rápido de la pérdida, alcanzando niveles bajos hacia la mitad del entrenamiento. No obstante, su curva presenta una leve tendencia ascendente entre las épocas 6 y 12, lo que podría indicar un ajuste excesivo seguido de una corrección. Mamba también presenta un descenso eficiente al inicio, aunque su trayectoria es más irregular, con fluctuaciones entre las épocas 6 y 12, antes de estabilizarse en valores competitivos hacia el final. Este comportamiento refleja una adaptación progresiva a un entorno más ruidoso, sin comprometer su capacidad de convergencia.

MEGA mantiene una curva más constante, con una reducción sostenida pero menos pronunciada, mientras que MoE y Vainilla presentan una trayectoria casi paralela, con descensos más

lentos y estancamiento temprano, lo que podría sugerir una menor sensibilidad a los patrones irregulares del corpus OSPC. El modelo GSS, por su parte, exhibe nuevamente una curva elevada y poco descendente, indicando dificultades significativas para ajustar sus parámetros ante datos no estructurados.

En conjunto, esta visualización sugiere que modelos como UT, Mamba y MEGA logran adaptarse de forma más efectiva al ruido y la informalidad del corpus OSPC. En particular, Mamba destaca por mantener una curva de pérdida competitiva con una arquitectura más eficiente, lo que refuerza su valor como solución versátil en tareas de traducción automática en dominios menos predecibles.

6.4.2 Pérdida durante la fase de validación

A continuación, se presentan los resultados del análisis de la pérdida (loss) durante la fase de validación para cada una de las arquitecturas evaluadas.

En esta sección se comparan los valores iniciales y finales de la pérdida en validación para cada arquitectura, así como la variación porcentual entre ambas mediciones. Este análisis permite determinar qué modelos logran un mejor equilibrio entre aprendizaje y generalización, y cuáles presentan dificultades para transferir el conocimiento adquirido durante el entrenamiento a contextos nuevos. Esta información es esencial para seleccionar arquitecturas robustas y confiables en entornos de uso real.

6.4.2.1 Conjunto de datos UNPC

Cuadro No. 6.7. Comparación de pérdida inicial y final durante la fase de validación en el conjunto de datos UNPC

Modelo	Epoch 1	Epoch 20	Variación
Vainilla	4,4282	3,3428	-24,51%
UT	3,2661	2,3156	-29,10%
Mamba	4,0961	2,8048	-31,53%
MoE	3,3545	2,6223	-21,83%
GSS	5,1831	5,3566	3,35%
Mega	4,1497	2,5633	-38,23%

Simbología Epoch 1 5,3566 Epoch 20 5,1831 4,4282 4,1497 4,0961 3,3545 3,3428 3,2661 -22% 2,8048 2,6223 2,5633 2.3156

Gráfico No. 6.10. Comparación de pérdida inicial y final durante la fase de validación en el conjunto de datos UNPC

Mamba

UT

Vainilla

El **Cuadro No. 13** muestra la evolución de la pérdida durante la fase de validación en el conjunto de datos UNPC, comparando los valores obtenidos en la primera y segunda época. Aunque el número de épocas analizadas es reducido, los resultados permiten observar la rapidez con la que cada modelo logra estabilizar su desempeño frente a datos no vistos, lo cual es indicativo de su capacidad de generalización temprana.

MoE

GSS

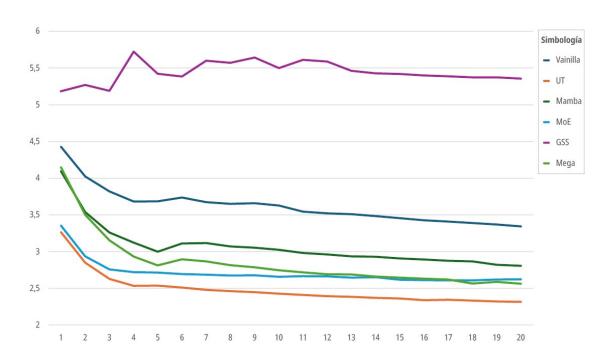
En este contexto, MEGA presenta la mayor mejora relativa con una reducción del -38,23 %, seguido de Mamba con -31,53 % y UT con -29,10 %. Estos resultados refuerzan la competitividad de estas tres arquitecturas en términos de rendimiento en validación desde las primeras etapas del entrenamiento. En particular, Mamba vuelve a demostrar una mejora sustancial, combinando buena capacidad de ajuste con un consumo moderado de recursos, lo que lo consolida como una opción robusta y eficiente.

El modelo Vainilla también logra una mejora destacable (-24,51 %), aunque inferior a la de los modelos más avanzados. En contraste, MoE presenta una mejora más modesta (-21,83 %),

mientras que GSS es el único modelo que muestra un incremento en la pérdida de validación (+3,35 %), lo que sugiere sobreajuste prematuro o falta de generalización.

Estos datos, aunque limitados a dos épocas, aportan información relevante sobre la capacidad de respuesta temprana de cada arquitectura ante el proceso de validación. Una vez más, Mamba se perfila como una solución eficiente, que logra avances rápidos en precisión con bajo coste computacional, y con una buena capacidad de generalización desde las primeras fases del entrenamiento.

Gráfico No. 6.11. Comportamiento de la medida de pérdida de los modelos durante la validación para el conjunto UNPC.



Fuente: Comparativa de soluciones, elaboración propia del estudiante.

La gráfica presentada muestra las curvas de pérdida durante la validación en el conjunto UNPC, permitiendo evaluar la evolución del error que cometen los modelos al enfrentarse a datos no vistos durante el entrenamiento. Esta medida es esencial para valorar la capacidad de generalización de cada arquitectura en un corpus estructurado y formal como el UNPC.

En términos generales, se observa que los modelos UT, MoE y MEGA mantienen las curvas más bajas y estables a lo largo de las 20 épocas, lo que indica un desempeño sólido y consistente ante

datos nuevos. UT, en particular, destaca por lograr la menor pérdida en validación, con una trayectoria decreciente y sin oscilaciones relevantes, lo que evidencia un alto nivel de generalización y estabilidad. MoE y MEGA también reflejan curvas descendentes, aunque con variaciones leves, pero dentro de rangos aceptables.

Mamba, por su parte, presenta una curva ligeramente superior pero estable, con una pérdida controlada y una evolución progresiva. Aunque no es el modelo con menor pérdida absoluta, su comportamiento sugiere una buena capacidad de adaptación sin signos de sobreajuste, lo que refuerza su carácter equilibrado.

El modelo Vainilla muestra una pérdida claramente más elevada y una curva menos eficiente en su descenso, lo que evidencia una menor habilidad para generalizar. Finalmente, GSS presenta la peor evolución: su curva se mantiene alta, con fluctuaciones significativas y sin una reducción clara, lo que sugiere inestabilidad durante la validación y dificultades para aprender patrones generalizables del corpus.

En conjunto, esta visualización permite concluir que UT, MoE y MEGA logran las mejores tasas de generalización en UNPC, mientras que Mamba, sin ser el líder absoluto, ofrece un desempeño robusto con una arquitectura más eficiente. En cambio, Vainilla y GSS muestran limitaciones importantes que comprometen su aplicabilidad en tareas de validación sobre datos formales.

6.4.2.2 Conjunto de datos OSPC

Cuadro No. 6.8. Comparación de pérdida inicial y final durante la fase de validación en el conjunto de datos OSPC

Modelo	Epoch 1	Epoch 2	Variación
Vainilla	4,7862	4,4250	-7,55%
UT	5,0115	4,3353	-13,49%
Mamba	5,159	4,1483	-19,59%
MoE	4,6543	4,1175	-11,53%
GSS	5,9101	5,4902	-7,10%
Mega	5,2478	4,2313	-19,37%

5,9101 7 -7,1% Simbología Epoch 1 5,4902 Epoch 20 5.2478 5,159 5,0115 4,7862 7-7,5% 4.6543 -11,5% 4,425 4,3353 4,2313 4,1483 4,1175 Vainilla UT GSS Mamba MoE Mega

Gráfico No. 6.12. Comparación de pérdida inicial y final durante la fase de validación en el conjunto de datos OSPC

El **Cuadro No. 6.8** presenta los resultados de la pérdida de validación durante las primeras dos épocas de entrenamiento sobre el conjunto de datos OSPC, lo que permite evaluar la capacidad de generalización temprana de los modelos en un contexto de lenguaje informal y de frases cortas. A diferencia del conjunto UNPC, en OSPC se observan reducciones más discretas en la pérdida, probablemente debido a la naturaleza más ruidosa del corpus.

Entre las arquitecturas evaluadas, Mamba logra la mayor reducción de pérdida en validación con un -19,59 %, seguido muy de cerca por MEGA (-19,37 %) y UT (-13,49 %). Estos resultados consolidan a Mamba como una arquitectura consistente, que no solo se desempeña bien en corpus estructurados como UNPC, sino que también muestra una rápida mejora en entornos más informales como OSPC. Este comportamiento temprano sugiere una buena capacidad de ajuste inicial y potencial para convergencia eficiente.

Los modelos MoE y Vainilla muestran mejoras más modestas, con -11,53 % y -7,55 % respectivamente, mientras que GSS nuevamente presenta un resultado inferior al del resto, con apenas -7,10 % de reducción. Estos datos respaldan la idea de que arquitecturas como GSS requieren ajustes adicionales o pueden no ser óptimas para datos con mayor variabilidad.

En conjunto, los resultados de los cuadros 13 y 14 muestran que Mamba no solo ofrece eficiencia computacional y buena capacidad de entrenamiento, sino que también demuestra una excelente capacidad de generalización temprana en la fase de validación, tanto en corpus formales como informales.

Simbología

Vainilla

UT

Mamba

MoE

GSS

Mega

1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20

Gráfico No. 6.13. Comportamiento de la medida de pérdida de los modelos durante la validación para el conjunto de datos OSPC.

Fuente: Comparativa de soluciones, elaboración propia del estudiante.

La gráfica anterior presenta las curvas de pérdida durante la validación en el conjunto OSPC, ofreciendo una visión clara de cómo se comporta cada modelo al enfrentarse a un corpus informal, ruidoso y altamente variable como el de subtítulos. En este entorno, la capacidad de generalización y la estabilidad del aprendizaje son factores críticos para valorar la eficacia real de los modelos.

A lo largo de las 20 épocas, se observa que los modelos MoE, UT y Mamba mantienen un desempeño favorable, con curvas de pérdida que tienden a descender de manera sostenida y alcanzan valores bajos en las últimas etapas. MoE, en particular, destaca por lograr la menor pérdida promedio en validación hacia el final del entrenamiento, seguido de cerca por Mamba y UT, que también presentan trayectorias estables y consistentes.

En contraste, las curvas de MEGA y Vainilla son más irregulares y mantienen valores de pérdida más elevados, con oscilaciones moderadas a lo largo del entrenamiento, lo que sugiere una

menor capacidad de adaptación a las características informales del corpus OSPC. A pesar de ello, MEGA logra estabilizarse en valores aceptables hacia el final.

El modelo GSS presenta nuevamente el comportamiento más problemático. Su curva de pérdida es elevada, con fuertes oscilaciones y picos abruptos (notablemente entre las épocas 5 y 13), lo que refleja una alta inestabilidad y una clara incapacidad para aprender patrones generalizables en este tipo de datos.

En resumen, esta visualización confirma que MoE, UT y Mamba son los modelos con mejor capacidad de generalización y mayor estabilidad durante la validación en contextos menos estructurados. Entre ellos, Mamba destaca por combinar rendimiento competitivo con una arquitectura más eficiente, lo que refuerza su posición como una opción sólida en entornos de producción con variabilidad lingüística.

6.5 Exactitud

En una comparativa de soluciones para modelos de traducción automática, el análisis de la precisión tanto en el conjunto de entrenamiento como en el de validación es fundamental para evaluar la calidad del aprendizaje del modelo y su capacidad de generalización.

La precisión en el conjunto de entrenamiento indica qué tan bien el modelo ha aprendido a reproducir los patrones lingüísticos presentes en los datos que ya ha visto. Sin embargo, una precisión alta en entrenamiento no garantiza un buen desempeño si el modelo ha memorizado los datos sin captar reglas generales del lenguaje, es decir, si ha incurrido en sobreajuste (overfitting).

Por eso, la precisión en el conjunto de validación —que consiste en datos no vistos durante el entrenamiento— es clave para medir si el modelo generaliza adecuadamente a nuevos ejemplos. Una brecha amplia entre la precisión en entrenamiento y validación puede revelar problemas de sobreajuste, mientras que una precisión similar en ambos conjuntos suele indicar un modelo bien balanceado y más confiable para su uso en el mundo real.

Este análisis también permite comparar la estabilidad y consistencia de distintas arquitecturas: modelos que muestran un rendimiento alto en ambos conjuntos son preferibles, incluso si no son los más precisos en entrenamiento, ya que son más robustos ante variaciones del lenguaje y dominio. Por el contrario, un modelo que rinde muy bien solo en entrenamiento pero se

degrada en validación, puede resultar poco fiable en tareas reales de traducción automática, donde se enfrentará constantemente a textos nuevos.

6.5.1 Conjunto de datos UNPC

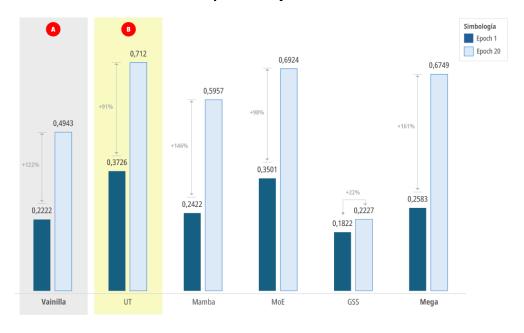
La tasa de mejora inicial es tan relevante como la precisión final, especialmente en contextos donde se requieren ciclos de entrenamiento cortos o donde se monitoriza el desempeño desde etapas tempranas.

Cuadro No. 6.9. Comparación de la exactitud durante la primera y última época de la fase de entrenamiento para el conjunto de datos UNPC

Modelo	Epoch 1	Epoch 20	Variación
Vainilla	0,2222	0,4943	122,46%
UT	0,3726	0,7120	91,09%
Mamba	0,2422	0,5957	145,95%
MoE	0,3501	0,6924	97,77%
GSS	0,1822	0,2227	22,23%
Mega	0,2583	0,6749	161,29%

Fuente: Comparativa de soluciones, elaboración propia del estudiante.

Gráfico No. 6.14 Comparación de la exactitud durante la primera y última época de la fase de entrenamiento para el conjunto de datos UNPC



El **Cuadro No. 15** muestra la evolución de la precisión durante las dos primeras épocas de entrenamiento en el conjunto de datos UNPC, destacando la capacidad de los modelos para aprender de manera rápida en etapas iniciales. Esta métrica no solo refleja el ritmo de aprendizaje temprano, sino también la eficacia de la arquitectura para adaptarse al problema desde el inicio del entrenamiento.

Entre los modelos evaluados, MEGA y Mamba son los que presentan las mayores mejoras porcentuales, con un aumento del 161,29 % y 145,95 % respectivamente. Esto evidencia que ambas arquitecturas poseen una alta capacidad de aprendizaje desde las primeras iteraciones, lo que resulta ventajoso para entornos con limitaciones de tiempo o cuando se busca una convergencia rápida.

MoE y UT también muestran mejoras significativas (97,77 % y 91,09 %), mientras que el modelo Vainilla mejora un 122,46 %, una cifra sorprendente considerando su simplicidad estructural, aunque partiendo desde una precisión inicial muy baja. El modelo GSS, en contraste, presenta una mejora marginal del 22,23 %, lo que refuerza su bajo rendimiento ya observado en métricas anteriores.

Estos resultados respaldan la eficacia temprana de modelos como Mamba y MEGA, siendo especialmente relevante en ciclos de entrenamiento reducidos, iteración rápida de prototipos o despliegues ágiles. En particular, Mamba logra este desempeño temprano con una arquitectura más ligera y tiempos de entrenamiento inferiores, consolidando su ventaja operativa frente a modelos más pesados.

0,6

0,5

0,4

0,3

0,2

Simbología

- Vainilla

- UT

- Mamba

- MoE

- GSS

- Mega

Gráfico No. 6.15. Comportamiento de la medida de exactitud de los modelos durante el entrenamiento en el conjunto de datos UNPC.

0,1

La gráfica anterior muestra las curvas de exactitud durante la fase de entrenamiento en el conjunto UNPC, permitiendo analizar cómo evoluciona la capacidad de los modelos para realizar predicciones correctas a lo largo del tiempo. Esta métrica, entendida como la proporción de aciertos frente al total de ejemplos procesados, es útil para valorar la eficacia del aprendizaje supervisado en términos de rendimiento progresivo.

Desde las primeras épocas, se observa un crecimiento acelerado en los modelos UT, MoE, Mamba y MEGA, todos los cuales superan rápidamente el umbral del 0,50 de exactitud. Entre ellos, UT lidera consistentemente con una curva ascendente y sostenida, alcanzando el nivel más alto de exactitud al finalizar el entrenamiento. MoE y MEGA también logran resultados muy competitivos, aunque con una ligera desaceleración hacia el final. Mamba, por su parte, mantiene una progresión constante, situándose ligeramente por debajo de los anteriores pero con una trayectoria estable y sin caídas abruptas, lo que indica una buena capacidad de aprendizaje sin señales de sobreajuste temprano.

En contraste, el modelo Vainilla mejora a lo largo del tiempo pero con una pendiente más moderada, y se mantiene por debajo del 0,50 durante la mayor parte del proceso. Finalmente, el modelo GSS presenta un desempeño claramente inferior, con una curva plana y valores de exactitud

significativamente más bajos que el resto, lo que sugiere una limitada capacidad de aprendizaje en este corpus.

En conjunto, esta visualización confirma que UT, MoE y MEGA alcanzan los mayores niveles de exactitud durante el entrenamiento, pero también que Mamba logra un equilibrio entre buen rendimiento y estabilidad, lo cual lo posiciona como una opción confiable para tareas de traducción con estructuras lingüísticas formales como las del corpus UNPC.

Cuadro No. 6.10. Comparación de la precisión durante la primera y última época de la fase de validación para el conjunto de datos OSPC

Modelo	Epoch 1	Epoch 20	Variación
Vainilla	0,2929	0,4971	69,72%
UT	0,4754	0,6774	42,49%
Mamba	0,3624	0,6128	69,09%
MoE	0,4785	0,6489	35,61%
GSS	0,1978	0,2259	14,21%
Mega	0,3431	0,6520	90,03%

Fuente: Comparativa de soluciones, elaboración propia del estudiante.

El cuadro anterior complementa el análisis anterior al mostrar la variación en precisión durante la fase de validación, entre la primera y la segunda época, en el conjunto de datos UNPC. Esta métrica permite observar no solo cuánto aprende el modelo, sino cómo se traslada ese aprendizaje a datos no vistos, lo cual es crucial para valorar la capacidad de generalización desde etapas tempranas.

El modelo con mayor mejora relativa en validación fue MEGA, con un incremento del 90,03 %, seguido por Vainilla (69,72 %) y Mamba (69,09 %). Esta mejora significativa confirma que Mamba, además de aprender con rapidez, generaliza eficazmente desde las primeras épocas, una característica especialmente valiosa en flujos de desarrollo donde el tiempo de entrenamiento está restringido o se aplican métodos de early stopping.

En contraste, modelos como UT (42,49 %) y MoE (35,61 %), aunque robustos en precisión final, muestran una menor velocidad de mejora en validación, lo que sugiere una curva de aprendizaje más gradual. GSS, con apenas un 14,21 % de mejora, vuelve a confirmar su debilidad en capacidad de generalización y escasa respuesta al entrenamiento en las primeras etapas.

0,6

0,5

0,4

0,3

Simbología

- Vainilla

- UT

- Mamba

- MoE

- GSS

- Mega

Gráfico No. 6.16. Comportamiento de la medida de exactitud de los modelos durante la validación para el conjunto de datos UNPC.

La gráfica anterior muestra las curvas de exactitud durante la fase de validación en el conjunto UNPC, permitiendo evaluar la capacidad de los modelos para generalizar el conocimiento adquirido durante el entrenamiento hacia datos no vistos, en un corpus estructurado y formal. La exactitud en validación es una métrica clave para identificar modelos robustos, capaces de mantener un rendimiento alto más allá del conjunto de entrenamiento.

Los modelos UT, MoE y Mamba destacan claramente por sus curvas ascendentes sostenidas, alcanzando niveles de exactitud superiores al 0,60 al finalizar el proceso. UT lidera la comparación, seguido de cerca por MoE, cuya curva es ligeramente más estable, y por Mamba, que mantiene una evolución constante y competitiva, confirmando su buena capacidad de generalización en entornos formales. MEGA se sitúa ligeramente por debajo, pero su desempeño es consistente a lo largo de las épocas.

El modelo Vainilla, si bien muestra una progresión clara, alcanza un techo más bajo (cercano a 0,50), lo que sugiere limitaciones en su capacidad de adaptación frente a arquitecturas más modernas. Por otro lado, GSS evidencia un comportamiento claramente inferior: su curva permanece en niveles bajos

(por debajo del 0,22) con fluctuaciones notables, lo que refleja una falta de aprendizaje significativo y una muy limitada capacidad de generalización.

En conjunto, la gráfica confirma que modelos como UT, MoE y Mamba no solo aprenden con rapidez, sino que logran transferir ese aprendizaje de forma efectiva a nuevos datos, haciendo de ellos opciones idóneas para tareas de traducción automática en contextos formales y controlados como los que representa el corpus UNPC.

Cuadro No. 6.11. Comparación de la precisión durante el entrenamiento y la validación en el conjunto de datos UNPC

Modelo	Entrenamiento	Validación	
Vainilla	0,4943	0,4971	
UT	0,7120	0,6774	
Mamba	0,5957	0,6128	
MoE	0,6924	0,6489	
GSS	0,2227	0,2259	
Mega	0,6749	0,6520	

Fuente: Comparativa de soluciones, elaboración propia del estudiante.

El cuadro anterior presenta los valores de precisión obtenidos durante las fases de entrenamiento y validación en el conjunto de datos UNPC, permitiendo valorar no solo el rendimiento de los modelos sobre los datos de entrenamiento, sino también su capacidad de generalización. Esta métrica es fundamental para determinar cuán bien el modelo está capturando patrones útiles sin sobreajustarse.

Los modelos que alcanzan mayores niveles de precisión tanto en entrenamiento como en validación son UT (0,7120 / 0,6774), MoE (0,6924 / 0,6489) y MEGA (0,6749 / 0,6520), lo cual confirma su alto rendimiento en términos de calidad de predicción. No obstante, como se ha analizado anteriormente, estas arquitecturas requieren más tiempo de entrenamiento y validación, mayor número de parámetros y, en algunos casos, mayor consumo de memoria.

En este contexto, Mamba se posiciona nuevamente como una opción intermedia robusta, logrando una precisión de 0,5957 en entrenamiento y 0,6128 en validación, cifras que si bien son más bajas que las de los modelos más pesados, superan con claridad al modelo Vainilla y al modelo GSS. Este resultado respalda la eficacia y consistencia de Mamba, especialmente al considerar su menor complejidad arquitectónica y sus ventajas en eficiencia computacional.

Por otro lado, el modelo GSS es el que presenta el desempeño más bajo en ambas métricas (0,2227 / 0,2259), lo que se alinea con su escasa reducción de pérdida y menor capacidad de aprendizaje observada en análisis anteriores.

Este cuadro complementa y reafirma las conclusiones anteriores: Mamba ofrece un equilibrio muy favorable entre precisión, eficiencia y simplicidad, lo que lo convierte en una de las arquitecturas más competitivas para tareas de traducción automática sobre corpus estructurados como UNPC.

6.5.2 Conjunto de datos OSPC

Cuadro No. 6.12. Comparación de la precisión durante la primera y última época de la fase de entrenamiento para el conjunto de datos OSPC

Modelo	Epoch 1	Epoch 20	Variación
Vainilla	0,2847	0,3453	21,29%
UT	0,3242	0,6827	110,58%
Mamba	0,3386	0,5478	61,78%
MoE	0,4096	0,6957	69,85%
GSS	0,3071	0,3893	26,77%
Mega	0,3592	0,6358	77,00%

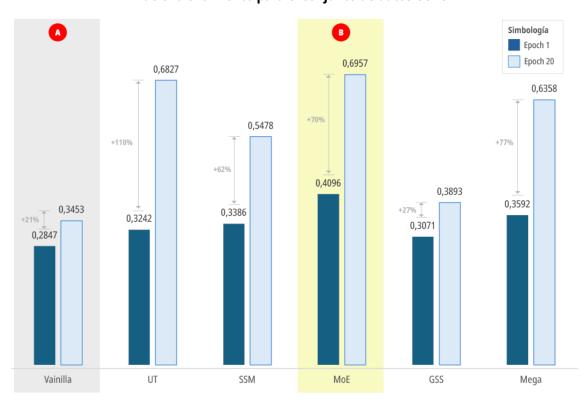


Gráfico No. 6.17. Comparación de la precisión durante la primera y última época de la fase de entrenamiento para el conjunto de datos OSPC

La información anterior presenta la variación en la precisión durante la fase de entrenamiento en el conjunto de datos OSPC, comparando la primera y segunda época. Este análisis permite evaluar la rapidez con la que cada modelo logra aprender en un entorno lingüístico más informal y fragmentado, lo que es fundamental para medir su capacidad de adaptación inicial.

En este corpus, el modelo con mayor incremento porcentual de precisión fue UT, con un aumento del 110,58 %, seguido de MEGA (77,00 %), MoE (69,85 %) y Mamba (61,78 %). Estos resultados indican que estas arquitecturas, incluyendo Mamba, tienen una sólida capacidad de ajuste desde las primeras fases del entrenamiento, incluso en contextos menos estructurados como OSPC. Mamba, en particular, logra una mejora importante con una arquitectura mucho más liviana, lo que refuerza su ventaja en eficiencia y velocidad de entrenamiento.

En contraste, modelos como GSS (26,77 %) y Vainilla (21,29 %) muestran mejoras modestas, lo que sugiere una menor sensibilidad al aprendizaje temprano en este entorno. Estos resultados coinciden con sus desempeños inferiores en métricas de validación y BLEU analizadas previamente.

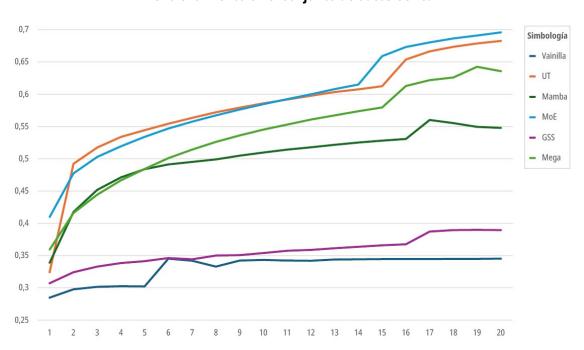


Gráfico No. 6.18. Comportamiento de la medida de exactitud de los modelos durante el entrenamiento en el conjunto de datos OSPC.

La gráfica anterior muestra el comportamiento de las curvas de exactitud durante el entrenamiento en el conjunto OSPC, permitiendo evaluar cómo evoluciona el rendimiento de cada modelo en un corpus caracterizado por su informalidad, menor estructura sintáctica y frases más breves. En este contexto, lograr una alta exactitud implica que el modelo ha podido adaptarse efectivamente a un entorno más variable y ruidoso.

Desde las primeras épocas, los modelos MoE y UT se destacan por su rápido ascenso en la medida de exactitud, manteniéndose como los modelos con mejor desempeño general al alcanzar valores cercanos al 0,70 al finalizar el entrenamiento. Su progreso es sostenido y continuo, lo que evidencia una curva de aprendizaje eficiente y una alta capacidad de adaptación a los patrones del corpus. Mamba y MEGA también muestran un crecimiento constante, aunque con curvas ligeramente más moderadas. Mamba, en particular, logra una evolución estable y sin caídas, consolidando su buen rendimiento pese a contar con una arquitectura más liviana.

El modelo GSS presenta una mejora gradual, aunque se mantiene en niveles inferiores en comparación con las demás arquitecturas. Sin embargo, su curva refleja una progresión más marcada que en el corpus UNPC, lo que podría indicar una mejor adaptación relativa en este contexto. Por el contrario, el modelo Vainilla muestra una curva plana y baja durante casi todo el entrenamiento, con

una evolución muy limitada en su exactitud, lo cual sugiere una dificultad para captar las particularidades lingüísticas del corpus OSPC.

En conjunto, esta visualización permite concluir que, en tareas de traducción con datos más informales, modelos como MoE, UT y Mamba no solo mantienen un buen rendimiento, sino que también evidencian una alta capacidad de aprendizaje y adaptabilidad, siendo adecuados para escenarios reales donde el lenguaje presenta mayor variabilidad.

Cuadro No. 6.13. Comparación de la precisión durante la primera y última época de la fase de validación para el conjunto de datos OSPC

Modelo	Epoch 1	Epoch 2	Variación
Vainilla	0,2796	0,3118	11,52%
UT	0,3750	0,5608	49,55%
Mamba	0,3099	0,4544	46,63%
MoE	0,3803	0,4804	26,32%
GSS	0,2444	0,2949	20,66%
Mega	0,2812	0,4466	58,82%

Fuente: Comparativa de soluciones, elaboración propia del estudiante.

El **Cuadro No. 19** muestra la evolución de la precisión durante la validación en el conjunto de datos OSPC, entre la primera y segunda época, lo que proporciona una visión sobre la capacidad de generalización temprana de los modelos en un entorno menos estructurado y más variable.

El modelo con mayor incremento porcentual fue MEGA, con una mejora del 58,82 %, seguido de UT (49,55 %) y SSM (Mamba) con 46,63 %, todos ellos mostrando una respuesta efectiva a datos no vistos desde las primeras fases del entrenamiento. Estos resultados validan el potencial de dichas arquitecturas para adaptarse a corpus informales y de menor consistencia estructural, como es típico en subtítulos (OSPC).

Por el contrario, MoE (26,32 %), GSS (20,66 %) y Vainilla (11,52 %) muestran mejoras más discretas, lo que sugiere una capacidad de generalización más lenta o limitada en fases iniciales. En particular, el comportamiento de Mamba (SSM) destaca positivamente al ofrecer una mejora cercana al 47 % con un modelo más ligero y eficiente, lo cual refuerza su posición como solución sólida tanto en rapidez como en calidad de validación temprana.

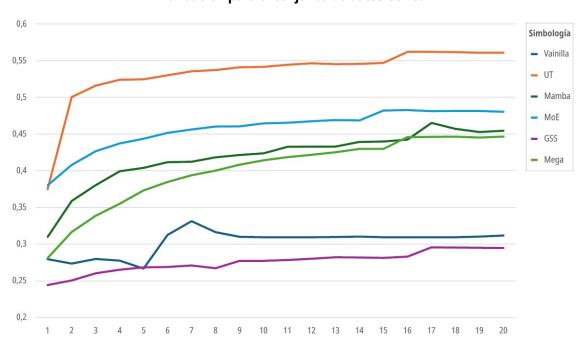


Gráfico No. 6.19. Comportamiento de la medida de exactitud de los modelos durante la validación para el conjunto de datos OSPC.

Fuente: Comparativa de soluciones, elaboración propia del estudiante.

La gráfica anterior muestra las curvas de exactitud durante la fase de validación en el conjunto OSPC, lo que permite analizar la capacidad de generalización de los modelos en un entorno más informal y menos estructurado. En este contexto, la evolución de la exactitud es un indicador clave para identificar qué arquitecturas logran adaptarse con eficacia a la variabilidad del lenguaje.

Desde las primeras épocas, el modelo UT se posiciona como el más destacado, con una rápida y sostenida mejora que lo lleva a superar el 0,55 de exactitud, manteniéndose como el líder a lo largo de todo el entrenamiento. MoE le sigue con una curva ascendente y estable, que alcanza valores cercanos al 0,48, confirmando su solidez incluso en un corpus ruidoso. Mamba y MEGA también muestran una evolución progresiva, logrando niveles competitivos entre 0,45 y 0,46 hacia el final de la validación, lo que refleja una buena capacidad de adaptación con arquitecturas más eficientes.

En contraste, los modelos Vainilla y GSS presentan un rendimiento limitado. Aunque ambos mejoran ligeramente con el tiempo, sus curvas se mantienen por debajo del 0,32, evidenciando dificultades para generalizar correctamente en un entorno como el OSPC. GSS, aunque muestra una mejora más sostenida que en otros contextos, sique quedando muy por debajo de los modelos más avanzados.

En conjunto, esta gráfica demuestra que UT, MoE, Mamba y MEGA son los modelos con mejor capacidad de generalización en tareas de traducción sobre datos menos estructurados. En particular, Mamba destaca por lograr una exactitud sólida con una arquitectura más ligera, consolidándose como una opción eficaz y eficiente para dominios con mayor variabilidad lingüística.

Cuadro No. 6.14. Comparación de la precisión durante el entrenamiento y la validación en el conjunto de datos OSPC

Modelo	Entrenamiento	Validación
Vainilla	0,3453	0,3118
UT	0,6827	0,5608
Mamba	0,5478	0,4544
MoE	0,6957	0,4804
GSS	0,3893	0,2949
Mega	0,6358	0,4466

Fuente: Comparativa de soluciones, elaboración propia del estudiante.

De acuerdo con el cuadro anterior, estos son los resultados de precisión durante el entrenamiento y la validación en el conjunto de datos OSPC, completando así la comparativa del rendimiento de las distintas arquitecturas bajo condiciones de lenguaje informal y fragmentado. Esta métrica permite evaluar no solo qué tan bien el modelo aprende, sino también cuán efectivamente generaliza a nuevos datos en escenarios menos estructurados.

Los modelos que destacan en esta evaluación son nuevamente MoE (0,6957 / 0,4804), UT (0,6827 / 0,5608) y MEGA (0,6358 / 0,4466), que logran altos niveles de precisión durante el entrenamiento, aunque todos muestran una caída considerable al pasar a la validación, lo cual podría sugerir cierto grado de sobreajuste. Mamba, por su parte, mantiene un comportamiento estable con una precisión de 0,5478 en entrenamiento y 0,4544 en validación, lo que demuestra su capacidad de aprendizaje sólido sin incurrir en sobreajuste severo, algo especialmente valioso en contextos ruidosos como OSPC.

El modelo Vainilla nuevamente presenta el desempeño más bajo (0,3453 / 0,3118), reflejando su menor capacidad para adaptarse a datos más complejos, mientras que GSS se mantiene en niveles similares con valores de 0,3893 en entrenamiento y 0,2949 en validación.

A continuación, se presenta el análisis del gap E–V (entrenamiento vs validación) a partir de la tabla actualizada. Esta métrica permite observar la diferencia entre el desempeño del modelo en los datos que ha visto (entrenamiento) y en los datos nuevos (validación). Un gap pequeño o positivo

indica buena capacidad de generalización, mientras que un gap negativo más grande puede ser signo de sobreajuste.

Cuadro No. 6.15. Comparación de brechas entre la precisión durante el entrenamiento y la precisión durante la validación

Modelo		UNPC		OSPC			
Modelo	Entrenamiento	Validación	Gap E-V	Entrenamiento	Validación	Gap E-V	
Vainilla	0,4943	0,4971	0,0028	0,3453	0,3118	-0,0335	
UT	0,7120	0,6774	-0,0346	0,6827	0,5608	-0,1219	
Mamba	0,5957	0,6128	0,0171	0,5478	0,4544	-0,0934	
MoE	0,6924	0,6489	-0,0435	0,6957	0,4804	-0,2153	
GSS	0,2227	0,2259	0,0032	0,3893	0,2949	-0,0944	
Mega	0,6749	0,6520	-0,0229	0,6358	0,4466	-0,1892	

Fuente: Comparativa de soluciones, elaboración propia del estudiante.

En el conjunto de datos UNPC, los modelos Mamba (0,0171), Vainilla (0,0028) y GSS (0,0032) presentan los gaps más bajos, lo que indica una buena capacidad de generalización sin señales significativas de sobreajuste. MEGA también mantiene un gap moderado (0,0229), mientras que UT (0,0346) y MoE (0,0435) muestran diferencias más pronunciadas, lo que podría evidenciar una mayor dependencia del modelo a los datos de entrenamiento.

En el corpus OSPC, de naturaleza más informal y fragmentada, los gaps aumentan de forma generalizada, siendo nuevamente Mamba (0,0934) uno de los modelos con mejor comportamiento en este entorno, junto a Vainilla (0,0335) y GSS (0,0944). Por el contrario, UT (0,1219), MEGA (0,1892) y MoE (0,2153) presentan los mayores descensos de precisión al validar, lo que sugiere una pérdida de capacidad de generalización en contextos más ruidosos. Estos resultados refuerzan el valor de Mamba como una arquitectura sólida, no solo por su precisión y eficiencia, sino también por su consistencia frente al sobreajuste, adaptándose de manera efectiva a diferentes dominios lingüísticos sin degradar su desempeño.

6.5.3 Perplejidad

La diferencia entre la perplejidad de entrenamiento y validación es crucial para evaluar la capacidad de generalización del modelo y para identificar problemas como el sobreajuste o el subajuste.

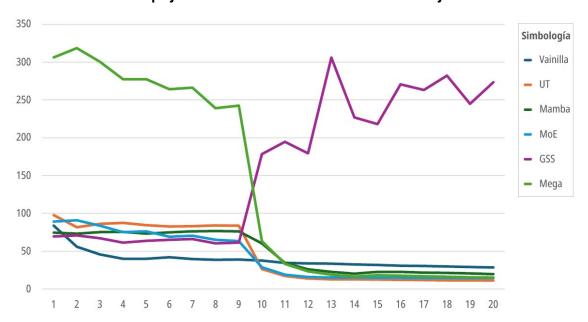


Gráfico No. 6.20. Perplejidad durante la fase de validación en el conjunto de datos UNPC.

Fuente: Comparativa de soluciones, elaboración propia del estudiante.

La gráfica anterior muestra el comportamiento de la perplejidad durante la validación en el conjunto UNPC, una métrica esencial para evaluar la capacidad de generalización de los modelos de lenguaje. La perplejidad mide qué tan bien predice un modelo una secuencia de palabras: cuanto menor el valor, mayor la certeza y fluidez de las predicciones. Una diferencia notable entre la perplejidad en entrenamiento y validación suele indicar problemas de sobreajuste o subajuste, por lo que esta curva es especialmente útil para identificar comportamientos anómalos durante la optimización.

Los modelos UT, MoE, Mamba y MEGA presentan los valores de perplejidad más bajos y estables hacia el final del entrenamiento, lo que indica un buen equilibrio entre aprendizaje y generalización. UT, en particular, logra estabilizarse por debajo de 30, con una curva descendente y sin oscilaciones, lo que refleja una alta consistencia en su desempeño. MoE y Mamba también muestran trayectorias decrecientes con valores muy similares, mientras que MEGA, a pesar de un inicio más inestable, logra reducir abruptamente su perplejidad a partir de la época 10.

El modelo Vainilla mantiene una curva más elevada (alrededor de 40), lo que indica una menor capacidad de predicción en validación, aunque con un comportamiento relativamente estable. En cambio, GSS presenta una curva caótica y altamente fluctuante, con picos que superan los 300 en varias épocas, lo cual sugiere una fuerte inestabilidad y una deficiente capacidad de generalización. Este comportamiento descalifica al modelo para tareas de traducción en contextos donde se requiere consistencia y precisión.

En conjunto, esta gráfica refuerza la evidencia de que Mamba, MoE y UT son los modelos más robustos en el corpus UNPC, logrando bajos niveles de perplejidad y curvas estables, lo cual es clave para aplicaciones reales donde se requiere un modelo confiable, preciso y eficiente en entornos formales y estructurados.

900 Simbología 800 Vainilla 700 Mamba 600 MoE 500 GSS 400 Mega 300 200 100 10 11 12 13 14 15

Gráfico No. 6.21. Perplejidad durante la fase de validación en el conjunto de datos OSPC.

Fuente: Comparativa de soluciones, elaboración propia del estudiante.

La gráfica anterior presenta las curvas de perplejidad durante la validación en el conjunto OSPC, proporcionando una medida cuantitativa de la incertidumbre de los modelos al predecir secuencias en un entorno de lenguaje más informal, fragmentado y con mayor variabilidad. La perplejidad, en este contexto, es fundamental para evaluar qué tan bien un modelo puede generalizar y manejar la ambigüedad inherente a este tipo de corpus.

En las primeras épocas, todos los modelos muestran niveles moderados y similares de perplejidad (entre 20 y 40), lo cual es indicativo de un comportamiento relativamente estable. Sin embargo, a partir de la época 10 se evidencian importantes divergencias. MoE, UT, Mamba y MEGA logran mantener la perplejidad bajo control, con oscilaciones puntuales pero dentro de márgenes razonables, y tienden a estabilizarse hacia el final. Este comportamiento sugiere una adaptabilidad razonable a los datos ruidosos del corpus OSPC, sin signos evidentes de sobreajuste.

El modelo Vainilla sigue una trayectoria similar, aunque con una perplejidad ligeramente superior de manera sostenida, lo que indica menor capacidad de predicción bajo incertidumbre. Por su parte, GSS presenta una curva dramáticamente inestable, con valores que alcanzan picos superiores a 700 (época 14), y fluctuaciones abruptas a lo largo del entrenamiento. Este comportamiento señala una gran dificultad para adaptarse a los patrones del corpus OSPC, posiblemente debido a una pobre capacidad de regularización o a un sobreajuste severo.

En resumen, esta gráfica refuerza las conclusiones anteriores: modelos como UT, MoE, Mamba y MEGA muestran un comportamiento más robusto y confiable, con menor perplejidad y mayor estabilidad incluso frente a un corpus complejo como OSPC. La perplejidad extremadamente elevada e inestable de GSS lo posiciona como una opción no viable para tareas de traducción automática en este tipo de entornos.

6.6 Inferencia

En una comparativa de soluciones para modelos de traducción automática, la evaluación con métricas como BLEU, el tiempo de inferencia promedio por secuencia y el tiempo total en los casos de prueba es esencial para obtener una visión completa del rendimiento de cada modelo. El BLEU (*Bilingual Evaluation Understudy*) es una métrica estándar que cuantifica la calidad de la traducción comparándola con traducciones humanas de referencia; permite evaluar cuán precisa y fluida es la salida del modelo, siendo clave para medir la fidelidad lingüística. Por otro lado, el tiempo de inferencia promedio por secuencia indica la velocidad con la que el modelo traduce frases individuales, lo cual es crítico para aplicaciones en tiempo real o con requerimientos de baja latencia. Finalmente, el tiempo total en los casos de prueba refleja la eficiencia general del modelo al procesar grandes volúmenes de datos, ayudando a dimensionar su rendimiento en producción o en tareas de evaluación masiva. Considerar estas tres métricas de forma conjunta permite balancear precisión y

eficiencia, elementos clave en la elección de una arquitectura adecuada para sistemas de traducción automáticos.

6.6.1 BLEU

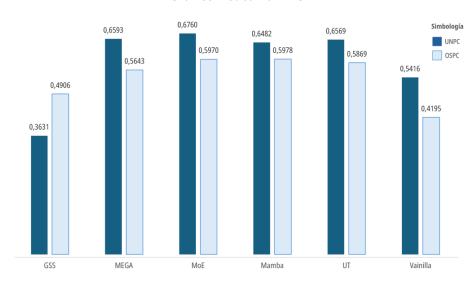
Como se mencionó en un inicio, en una comparativa de modelos de traducción automática, la métrica BLEU (*Bilingual Evaluation Understudy*) juega un papel fundamental como indicador objetivo de la calidad de las traducciones generadas por cada modelo. Su función principal es medir cuánto se parece la salida del modelo a una o más traducciones de referencia humanas, comparando n-gramas (secuencias de palabras) para evaluar precisión lingüística.

Cuadro No. 6.16. Resultados de la evaluación de la calidad de traducción mediante BLEU.

Modelo	Conjunto de datos					
Modelo	UNPC	OSPC				
GSS	0,3631	0,4906				
MEGA	0,6593	0,5643				
MoE	0,6760	0,5970				
Mamba	0,6482	0,5978				
UT	0,6569	0,5869				
Vainilla	0,5416	0,4195				

Fuente: Comparativa de soluciones, elaboración propia del estudiante.

Gráfico No. 6.22. BLEU4



Fuente: Comparativa de soluciones, elaboración propia del estudiante.

El cuadro y la gráfica anterior presentan los resultados de evaluación de calidad de traducción mediante la métrica BLEU, tanto en el conjunto de datos UNPC como en OSPC, lo que permite comparar de forma directa la fidelidad de las traducciones generadas por cada modelo con respecto a una referencia humana. BLEU es una métrica ampliamente utilizada en traducción automática, ya que captura la coincidencia de n-gramas entre la predicción del modelo y la traducción de referencia. En el corpus UNPC, los modelos con mayor puntuación BLEU son MoE (0,6760), MEGA (0,6593) y UT (0,6569), lo que confirma su elevada capacidad para generar traducciones precisas y naturales en un contexto de lenguaje formal. Sin embargo, Mamba también alcanza un desempeño sobresaliente con un BLEU de 0,6482, superando al modelo Transformer base (Vainilla) y dejando atrás ampliamente al modelo GSS, que obtuvo apenas 0,3631.

En el corpus OSPC, que implica desafíos adicionales por su carácter informal y variable, se observa una mayor dispersión en los resultados, aunque Mamba (0,5978) y MoE (0,5970) emergen como los modelos con mejor desempeño, seguidos por MEGA (0,5643) y UT (0,5869). Nuevamente, GSS y Vainilla son los modelos con menores puntuaciones, evidenciando sus limitaciones en tareas de mayor ruido lingüístico.

Estos resultados consolidan el perfil de Mamba como una arquitectura altamente competitiva también en términos de calidad de traducción, alcanzando niveles de BLEU comparables a los mejores modelos, pero con una eficiencia computacional notablemente superior. Su rendimiento estable en ambos conjuntos de datos confirma su versatilidad y adaptabilidad, cualidades clave para aplicaciones reales donde se enfrentan distintos dominios y condiciones de entrada.

6.6.2 Tiempo de inferencia

En una comparativa de algoritmos de traducción automática neuronal, uno de los factores más importantes a tener en cuenta es el tiempo de inferencia, es decir, el tiempo que tarda un modelo en procesar una entrada y generar una salida.

La importancia de comparar el tiempo de inferencia radica en varios aspectos fundamentales que pueden influir en la elección del modelo más adecuado para una aplicación específica. A continuación se detallan las razones más importantes para realizar esta comparación.

En aplicaciones donde la velocidad es crítica, como sistemas de traducción en tiempo real (chats, servicios de atención al cliente, subtitulación automática o interpretación), el tiempo de inferencia es un factor decisivo.

Si el tiempo de inferencia es elevado, el sistema podría introducir retrasos notables (latencia), lo que afectaría la experiencia del usuario. En este caso, un modelo más rápido, aunque ligeramente menos preciso, podría ser preferido.

También, un tiempo de respuesta rápido es esencial para aplicaciones interactivas en las que los usuarios esperan traducciones instantáneas.

Para servicios que manejan grandes volúmenes de peticiones de traducción, como aplicaciones globales o servicios en la nube que gestionan millones de usuarios, la eficiencia del modelo en cuanto a tiempo de inferencia tiene un impacto directo en la capacidad de la plataforma para escalar. Un modelo con un tiempo de inferencia elevado puede aumentar los costos operativos, ya que requerirá más recursos computacionales para manejar el mismo número de solicitudes. Modelos más rápidos pueden traducir más peticiones en el mismo período, lo que permite un uso más eficiente de la infraestructura y reduce la necesidad de ampliar servidores.

Los modelos de traducción automática más complejos, como los basados en arquitecturas profundas (transformers con muchos parámetros), tienden a ser más precisos pero también más lentos. Por tanto, la elección del modelo dependerá de la prioridad de la aplicación:

- Aplicaciones sensibles a la precisión: En casos donde la calidad de la traducción es primordial, como traducciones legales o médicas, podría ser aceptable un tiempo de inferencia más largo si garantiza una mayor precisión.
- Aplicaciones donde la velocidad es crítica: En otras aplicaciones, como traducción de redes sociales o contenido en vivo, se puede optar por un modelo menos preciso pero con menor latencia para garantizar una experiencia más fluida.

Por otro lado, en ciertos escenarios, como la implementación en dispositivos con recursos limitados (smartphones, tabletas o dispositivos IoT), el tiempo de inferencia de un modelo puede ser limitado por las capacidades de procesamiento del dispositivo. Modelos más ligeros y optimizados pueden ser necesarios para operar eficientemente en dispositivos móviles, donde la capacidad de

procesamiento y el consumo energético son factores clave. Comparar el tiempo de inferencia en estos contextos permite elegir un modelo que funcione dentro de estas restricciones.

En cualquier servicio o aplicación orientada al consumidor, la experiencia del usuario es un factor clave. Los usuarios tienden a preferir aplicaciones que ofrezcan resultados rápidos y efectivos. Un sistema de traducción que ofrezca traducciones en milisegundos es más probable que sea adoptado que uno que tarde varios segundos, incluso si la diferencia en precisión es marginal.

La mayoría de los usuarios tienen una tolerancia limitada a los retrasos, especialmente en aplicaciones en línea. Un tiempo de inferencia rápido mejora significativamente la satisfacción del usuario.

En entornos donde el consumo de energía es importante, como centros de datos o dispositivos móviles, un modelo que ofrezca un tiempo de inferencia más corto también podría implicar un menor consumo de energía, lo cual es beneficioso tanto en términos de costos operativos como de sostenibilidad ambiental.

6.6.2.1 Tiempo promedio de inferencia

En una comparativa de soluciones para modelos de traducción automática, el análisis del tiempo promedio de inferencia por registro (o *query*) resulta especialmente relevante porque permite evaluar la eficiencia operativa del modelo en condiciones reales de uso. A diferencia de métricas como la precisión o el BLEU, que miden la calidad de la traducción, el tiempo de inferencia refleja cuán rápido es el modelo al generar una traducción para una entrada concreta, lo cual es crucial en múltiples escenarios prácticos.

En primer lugar, en aplicaciones de traducción en tiempo real, como asistentes de voz, subtitulación automática o sistemas de mensajería multilingüe, la latencia por *query* determina directamente la experiencia del usuario. Un modelo con alta calidad pero con tiempos de respuesta lentos puede resultar inviable en contextos donde la fluidez y la inmediatez son prioritarias.

En segundo lugar, esta métrica permite comparar modelos en términos de costo computacional por unidad procesada, lo que es fundamental en sistemas desplegados en la nube, aplicaciones móviles o entornos con recursos limitados. Un modelo con menor tiempo de inferencia puede procesar más solicitudes por segundo, optimizando el uso de hardware y reduciendo los costos operativos.

Además, el tiempo promedio de inferencia por registro sirve como indicador de escalabilidad, permitiendo prever cómo se comportará el sistema ante volúmenes crecientes de solicitudes. En sistemas de traducción masiva o procesos por lotes, incluso pequeñas diferencias por *query* pueden representar horas de procesamiento acumulado.

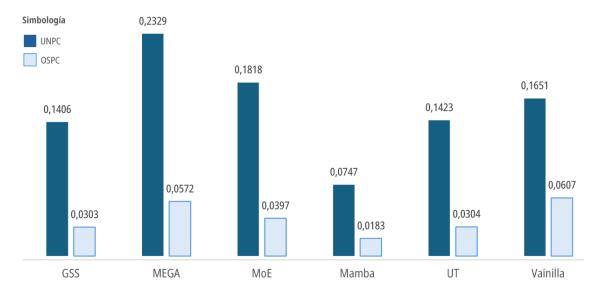
Cuadro No. 6.17. Tiempo promedio de inferencia por registro

Modelo	UNPC	OSPC
GSS	0,1406	0,0303
MEGA	0,2329	0,0572
MoE	0,1818	0,0397
Mamba	0,0747	0,0183
UT	0,1423	0,0304
Vainilla	0,1651	0,0607

Nota: Cifras expresadas en segundos.

Fuente: Comparativa de soluciones, elaboración propia del estudiante.

Gráfico No. 6.23. Tiempo promedio de inferencia por registro



Nota: Cifras expresadas en segundos.

Fuente: Comparativa de soluciones, elaboración propia del estudiante.

La información anterior muestra el tiempo promedio de inferencia por registro, expresado en segundos, tanto en los conjuntos de datos UNPC como OSPC. Esta métrica es esencial para aplicaciones prácticas en las que la velocidad de respuesta es crítica, como sistemas de traducción automática en tiempo real, asistentes conversacionales o procesamiento en línea de grandes volúmenes de texto.

En ambos corpus, Mamba destaca de forma contundente como el modelo más eficiente, con un tiempo promedio de inferencia de 0,0747 segundos en UNPC y 0,0183 segundos en OSPC, lo cual lo convierte en la opción más rápida del conjunto de arquitecturas evaluadas. Esta diferencia es considerable si se compara con modelos como MEGA (0,2329 / 0,0572) o MoE (0,1818 / 0,0397), cuyas velocidades de inferencia son notablemente inferiores pese a obtener puntuaciones BLEU similares.

El modelo Vainilla, aunque estructuralmente más simple, registra tiempos de 0,1651 segundos (UNPC) y 0,0607 segundos (OSPC), más del doble que Mamba en ambos casos. Modelos como UT y GSS también quedan rezagados, lo que confirma que Mamba no solo es eficiente en entrenamiento y validación, sino también altamente competitivo en producción.

Este rendimiento en inferencia refuerza aún más la idoneidad de Mamba para aplicaciones en entornos reales, donde la latencia puede marcar la diferencia en la experiencia del usuario o en el procesamiento eficiente de flujos de datos continuos.

6.6.2.2 Tiempo total de inferencia

Del mismo modo, en una comparativa de soluciones para modelos de traducción automática, el análisis del tiempo total de inferencia del conjunto de prueba es de gran relevancia porque proporciona una visión global de la eficiencia operativa del modelo a escala, complementando al tiempo promedio por registro con una métrica acumulativa que refleja el comportamiento real del sistema ante cargas completas de trabajo.

Este indicador es especialmente útil cuando se evalúa la viabilidad de un modelo en entornos productivos, ya que permite estimar el tiempo necesario para procesar grandes volúmenes de datos, como documentos extensos, bases de datos multilingües, o procesos de traducción masiva. Un modelo que tiene buen desempeño por *query* pero acumula largos tiempos totales puede generar cuellos de botella en sistemas de procesamiento por lotes, traducción de sitios web, o aplicaciones en servidores que atienden múltiples usuarios simultáneamente.

Además, el tiempo total de inferencia permite comparar la escalabilidad real de los modelos, especialmente cuando se combinan múltiples factores como complejidad arquitectónica, uso de recursos y velocidad de procesamiento. Esta métrica resulta también fundamental para dimensionar el coste computacional global, tanto en tiempo como en consumo energético, algo crucial en decisiones de implementación en la nube o sobre hardware dedicado.

Por tanto, el análisis del tiempo total de inferencia del conjunto de prueba no solo permite medir la rapidez de un modelo, sino que también aporta información estratégica para seleccionar soluciones que sean **sostenibles**, **eficientes y adecuadas para cargas de trabajo reales**, particularmente en sistemas que operan en tiempo continuo o que deben procesar grandes volúmenes de contenido en plazos reducidos.

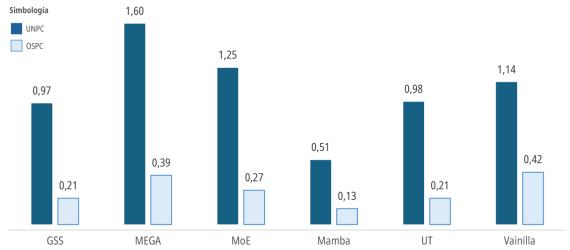
Cuadro No. 6.18. Tiempo total de inferencia del conjunto de datos de prueba

Modelo	UNPC	UNPC
GSS	0,97	0,21
MEGA	1,60	0,39
MoE	1,25	0,27
Mamba	0,51	0,13
UT	0,98	0,21
Vainilla	1,14	0,42

Nota: Cifras expresadas en horas.

Fuente: Comparativa de soluciones, elaboración propia del estudiante.

Gráfico No. 6.24. Tiempo total de inferencia



Nota: Cifras expresadas en horas.

Fuente: Comparativa de soluciones, elaboración propia del estudiante.

El **Cuadro No. 19** presenta el **tiempo total de inferencia** requerido por cada arquitectura para procesar el conjunto completo de datos de prueba en los corpus **UNPC** y **OSPC**, expresado en horas. Esta métrica es clave para evaluar la viabilidad de los modelos en escenarios de producción a gran

escala, donde el rendimiento acumulado puede tener un impacto directo en costos, disponibilidad de infraestructura y tiempos de respuesta al usuario.

En ambos conjuntos, el modelo que requirió menos tiempo total de inferencia fue Mamba, con apenas 0,51 horas en UNPC y 0,13 horas en OSPC, lo que lo convierte en la opción más eficiente entre todas las arquitecturas comparadas. Le siguen, con mayor distancia, GSS y UT, ambos con tiempos cercanos al doble que Mamba, mientras que MEGA y MoE se ubican entre los más lentos, con 1,60 y 1,25 horas respectivamente en UNPC, y 0,39 y 0,27 horas en OSPC. Por su parte, el modelo Vainilla, pese a su simplicidad estructural, presenta un tiempo de 1,14 horas en UNPC y 0,42 en OSPC, siendo superado ampliamente por Mamba.

Este resultado refuerza con claridad la ventaja operativa de Mamba, no solo en velocidad por registro individual (como se observó en el cuadro anterior), sino también en el desempeño acumulado en procesos de inferencia masiva. Esta diferencia se vuelve especialmente relevante en entornos donde se procesan millones de registros o se requiere disponibilidad en tiempo real, lo que convierte a Mamba en una alternativa sobresaliente por su capacidad de respuesta, bajo coste computacional y escalabilidad.

7 CONCLUSIONES Y TRABAJO FUTURO

7.1 Conclusiones

A partir de la comparativa de las soluciones avanzadas de NMT consideradas en esta comparativa se pueden derivar varias conclusiones importantes para el campo de la traducción automática basada en redes neuronales. Estas conclusiones se centran en la evaluación del rendimiento de los modelos, sus fortalezas, limitaciones y las áreas potenciales de mejora.

La evaluación integral de las arquitecturas de traducción automática —Vainilla (Transformer base), UT, MEGA, MoE, GSS y Mamba— permite concluir que Mamba representa la alternativa más equilibrada y eficiente en términos de rendimiento, costo computacional y adaptabilidad operativa.

Desde una perspectiva de calidad de traducción, medidas mediante la métrica BLEU, los modelos MoE, MEGA, UT y Mamba obtienen los mejores resultados, alcanzando valores muy similares entre sí tanto en el conjunto UNPC como en OSPC. No obstante, Mamba destaca al lograr este

desempeño sin incurrir en el alto número de parámetros ni en el coste computacional que requieren sus competidores más complejos.

En cuanto a la eficiencia computacional, Mamba es de forma clara la arquitectura que consume menos tiempo en todas las fases del ciclo de entrenamiento e inferencia. Logra tiempos de entrenamiento y validación inferiores al promedio, y registra el menor tiempo promedio de inferencia por registro y el menor tiempo total de inferencia, tanto en corpus estructurados (UNPC) como informales (OSPC). Esta ventaja operativa es crítica para contextos de despliegue en producción y entornos con recursos limitados.

Desde el punto de vista estructural, Mamba utiliza hasta un 21 % menos de parámetros que el Transformer base, manteniendo un consumo de memoria comparable al resto de modelos. Esto favorece su uso en dispositivos con capacidades reducidas y permite escalar a tareas más grandes sin necesidad de infraestructura especializada.

Los resultados de precisión y pérdida durante el entrenamiento y validación también respaldan la solidez del modelo. Mamba presenta una buena capacidad de aprendizaje desde las primeras épocas, evita el sobreajuste excesivo, y logra una generalización consistente en ambos conjuntos de datos. Si bien no lidera en todas las métricas de precisión, su desempeño es competitivo y estable, lo que evidencia un buen equilibrio entre capacidad expresiva y control del sesgovariabilidad.

Por contraste, modelos como MEGA, UT y MoE, si bien muestran niveles altos de precisión o BLEU, lo hacen a costa de una considerable inversión en tiempo, memoria y número de parámetros. El modelo GSS resulta el más débil del conjunto, con pobres resultados en pérdida, precisión y BLEU, y sin compensarlo con ventajas computacionales significativas. El modelo Vainilla, aunque simple y moderado en uso de recursos, muestra una calidad de traducción y precisión limitadas frente a las arguitecturas modernas.

En conjunto, estos hallazgos permiten afirmar que Mamba sobresale como una arquitectura moderna, eficiente y robusta, especialmente adecuada para tareas de traducción automática en entornos donde se requiere velocidad, eficiencia de recursos y facilidad de implementación. Su consistencia a través de métricas técnicas, computacionales y de calidad de salida lo posicionan como

una de las mejores alternativas disponibles actualmente para aplicaciones de modelado secuencial a gran escala.

Con base en la teoría presentada y los resultados empíricos obtenidos en esta comparativa de arquitecturas para traducción automática, se pueden establecer las siguientes conclusiones fundamentadas tanto en evidencia experimental como en el marco teórico:

Los Transformers continúan siendo el estándar sólido, pero las variantes superan en eficiencia y precisión en contextos específicos. Los resultados empíricos coinciden con lo señalado por Vaswani et al. (2017) y So et al. (2019): los Transformers siguen siendo una base robusta y confiable, como lo demuestra el modelo Vainilla en esta comparativa. Sin embargo, arquitecturas como UT y MEGA, alineadas con el concepto de Transformers evolucionados, mostraron mejoras considerables en precisión y métricas BLEU, aunque con un coste computacional más alto. Esto valida la noción de que las iteraciones estructurales pueden optimizar el rendimiento a costa de recursos.

El enfoque Mixture of Experts (MoE) demuestra gran escalabilidad y buena calidad de traducción. Tal como lo afirma Zoph et al. (2022), MoE es especialmente eficaz al distribuir el cómputo entre expertos especializados, y esto se refleja en su alto rendimiento en métricas BLEU y precisión. Sin embargo, el tiempo total de entrenamiento y la inferencia no fueron tan competitivos como en otros modelos, lo que sugiere que su eficiencia estructural se aprovecha mejor en contextos de entrenamiento masivo y no necesariamente en tareas con restricciones operativas severas.

Los mecanismos de gating y atención presentan ventajas en tareas con secuencias complejas, pero su efectividad práctica es variable. Aunque la teoría sugiere que arquitecturas como Gated State Spaces (GSS) y Moving Average Gated Attention son idóneas para manejar dependencias largas (Mehta et al., 2022; Ma et al., 2022), los resultados experimentales muestran un desempeño considerablemente bajo en GSS, tanto en pérdida como en precisión y BLEU. Esto sugiere que la implementación práctica de estos enfoques requiere una mayor madurez o combinación con otras técnicas para alcanzar su potencial.

El modelo Mamba, basado en un Space-State Model simplificado, confirma que es posible reducir complejidad sin comprometer el rendimiento. Tal como argumenta Smith et al. (2022), las representaciones compactas permiten una generalización eficaz con bajo coste computacional. Mamba fue el modelo que mejor balance logró entre precisión, calidad de **traducción**, **eficiencia de**

entrenamiento e inferencia, lo que lo convierte en una excelente materialización de los beneficios teóricos propuestos por este tipo de arquitectura simplificada.

Las métricas de calidad de traducción, como BLEU, reflejan diferencias importantes que se alinean con los enfoques estructurales. Modelos como MoE, MEGA y UT dominaron en BLEU, especialmente en el corpus UNPC, lo que valida que las variantes del Transformer pueden mejorar la calidad de traducción en corpus estructurados. No obstante, Mamba igualó o superó estos modelos en OSPC, lo que demuestra su capacidad de adaptación a contextos más ruidosos sin comprometer eficiencia.

El tiempo de entrenamiento e inferencia es un diferenciador clave para la implementación práctica. Las diferencias observadas respaldan lo planteado en la teoría: arquitecturas como **MoE y Mamba** muestran eficiencia operativa, pero solo Mamba combina esta ventaja con tiempos de inferencia extremadamente bajos, posicionándose como la mejor opción en tareas con exigencias de respuesta rápida o procesamiento a gran escala.

Las LSTM extendidas no fueron evaluadas en esta comparativa, pero los resultados sugieren que su integración con arquitecturas modernas podría ser valiosa. Aunque las Extended LSTM propuestas por Beck et al. (2024) no se incluyeron en los experimentos, la evidencia sugiere que combinar su capacidad de modelar continuidad temporal con arquitecturas eficientes como Mamba o MoE podría generar soluciones híbridas con mayor capacidad adaptativa.

Los enfoques híbridos ofrecen una vía prometedora para avanzar en eficiencia y rendimiento. A la luz de los datos y la teoría, resulta evidente que la combinación de mecanismos de atención, gating y expertos especializados, como los que presentan MoE o Mamba, representan una dirección valiosa para futuras investigaciones. En particular, modelos como Mamba podrían beneficiarse de la incorporación controlada de componentes expertos o atencionales sin perder su eficiencia base.

7.2 Trabajo futuro

Como se mencionó el objetivo de este trabajo no es presentar un nuevo algoritmo de entrenamiento, más bien busca indagar sobre diferentes aspectos que pueden conducir a una forma más económicamente viable de lograr el entrenamiento de modelos de redes profundas.

Un trabajo posterior podría enfocarse en abordar estas técnicas y evaluar su aplicabilidad y beneficios en el entrenamiento y la implementación de modelos de inteligencia artificial en un entorno de hardware de bajo costo, mixto o completamente especializado.

En el contexto de una tesis doctoral, podría resultar de interés desarrollar los diferentes enfoques estudiados acá y proponer un modelo unificado que combine las diferentes técnicas estudias y proponga una arquitectura que optimice los procesos de entrenamiento de redes neuronales artificiales enfatizando el bajo costo y la reutilización de componentes de hardware de consumo.

Estas recomendaciones ofrecen direcciones para investigar en futuros trabajos sobre NMT, explorando combinaciones innovadoras entre arquitecturas avanzadas como Transformers, MoE, GSS, y LSTM extendidas. Los enfoques híbridos y las técnicas emergentes como el uso de gating, modelos simplificados y compresión de modelos podrían llevar a mejoras significativas en la traducción automática, tanto en términos de eficiencia como de rendimiento. La investigación futura podría centrarse en probar estas combinaciones en diferentes conjuntos de datos y configuraciones de hardware para evaluar su efectividad y escalabilidad.

A partir de la comparativa de las arquitecturas avanzadas de Neural Machine Translation (NMT) mencionadas es posible identificar varias direcciones prometedoras para futuros trabajos de investigación. Las recomendaciones que se presentan a continuación no solo consideran el rendimiento de las arquitecturas individuales, sino también las oportunidades para combinar enfoques que podrían ofrecer mejoras en términos de eficiencia, escalabilidad, y capacidad de generalización.

7.2.1 Combinación de Mixture of Experts (MoE) con Gated State Spaces (GSS)

Una dirección interesante sería combinar ST-MoE (Mixture of Experts) con la arquitectura de Gated State Spaces (GSS). ST-MoE introduce módulos especializados que se activan en función de la tarea específica, lo que permite un uso más eficiente de los recursos. Por otro lado, GSS ofrece una forma avanzada de controlar el flujo de información a través de mecanismos de gating que pueden mejorar la captura de relaciones complejas en secuencias largas.

MoE se destaca en escenarios donde se busca reducir el costo computacional al activar solo una parte de la red, mientras que GSS mejora el control y filtrado de información relevante. La

combinación de estos dos enfoques podría mejorar tanto la eficiencia como la calidad de las traducciones, especialmente en secuencias largas o complejas. La combinación de MoE con técnicas de gating ha sido explorada en otros trabajos, como el modelo de Switch Transformer (Fedus et al., 2021), que integra ideas de gating para activar expertos específicos.

7.2.2 Incorporación de Atención Promediada en Modelos Transformer

La introducción de mecanismos de atención promediada, como el *Moving Average Equipped Gated Attention* (Ma et al., 2022), podría mejorar las arquitecturas *Transformer* al suavizar la información a través de ventanas temporales y mejorar la estabilidad en el entrenamiento de secuencias largas. Este enfoque puede ser beneficioso cuando se integra en *Evolved Transformers*, que ya han optimizado varias partes del modelo base *Transformer*.

La atención promediada permite mejorar la robustez del modelo al reducir el ruido y mitigar las fluctuaciones rápidas en el proceso de entrenamiento. Esto podría hacer que los Transformers evolucionados sean aún más efectivos para tareas de traducción de textos largos o con dependencia a largo plazo.

La combinación de mecanismos de atención promediada con Transformers ha sido explorada en algunos trabajos, como Reformer (Kitaev et al., 2020), que utiliza hashing local para mejorar la eficiencia del mecanismo de atención.

7.2.3 Integración de Extended LSTM con Arquitecturas Transformer

Aunque los Transformers han superado a las LSTM en muchos casos, el uso de Extended LSTM (Beck et al., 2024) podría proporcionar beneficios si se integran en capas híbridas junto con arquitecturas de Transformers. Las LSTM extendidas ofrecen un control adicional sobre la memoria a largo plazo y podrían complementar las limitaciones de los Transformers en la gestión de dependencias a largo plazo.

Mientras los Transformers destacan por su capacidad de procesamiento paralelo y su mecanismo de atención, las LSTM pueden ser más efectivas en la captura de relaciones temporales continuas a lo largo de la secuencia. El enfoque híbrido podría ofrecer lo mejor de ambos mundos: la potencia computacional del Transformer junto con la capacidad de las LSTM para recordar información pasada más allá de las ventanas de atención.

Trabajos como el de Combinations of RNN and Transformer (Liu et al., 2020) han explorado este enfoque híbrido, combinando Transformers con redes recurrentes para mejorar la traducción en dominios de lenguaje natural.

7.2.4 Exploración de Modelos Auto-regresivos Híbridos

Otra opción de investigación sería investigar la combinación de Universal Transformers con modelos *Simplified Space-State*. Los Universal Transformers presentan una forma de iterar sobre la secuencia de entrada de manera recurrente, aplicando el mismo mecanismo de atención en cada paso, lo que los convierte en una versión más general de los Transformers.

La capacidad de los *Simplified Space-State Models* para simplificar la estructura temporal y mejorar la capacidad de modelar secuencias largas podría hacer sinergia con la naturaleza iterativa de los Universal Transformers, potenciando el aprendizaje de patrones complejos en la traducción automática.

Universal Transformers se han combinado con técnicas recurrentes avanzadas en otras investigaciones, como los modelos *Recurrent Memory Networks* (Rae et al., 2021).

7.2.5 Hibridación de Modelos con MoE y Transformers Evolucionados

La combinación de la arquitectura *Evolved Transformer* con el enfoque MoE (ST-MoE) podría ser un camino interesante para optimizar tanto la eficiencia del modelo como su capacidad de generalización. El *Evolved Transformer* ya introduce mejoras basadas en búsquedas neuronales automáticas, mientras que MoE permitiría distribuir los recursos de manera más eficiente al activar expertos solo cuando sea necesario.

El *Evolved Transformer* optimiza la arquitectura *Transformer* tradicional, mientras que MoE mejora la escalabilidad y eficiencia. Combinando ambos enfoques, se podría lograr un modelo que no solo es más eficiente en términos computacionales, sino que también es capaz de manejar tareas de traducción complejas con menos recursos.

Modelos como el *Switch Transformer* (Fedus et al., 2021) ya combinan ideas de MoE con Transformers para reducir el costo computacional sin sacrificar la precisión.

7.2.6 Implementación de Técnicas de Compresión de Modelos

Un área emergente en la investigación de NMT es la compresión de modelos mediante distillation o pruning, lo que podría aplicarse a modelos como ST-MoE y Gated State Spaces. Dado que los modelos de gran tamaño pueden ser prohibitivamente costosos, reducir su tamaño sin perder precisión es un enfoque relevante.

Comprimir modelos que ya utilizan *gating* o expertos como MoE podría mejorar la eficiencia sin comprometer la precisión. Esta combinación sería especialmente útil en contextos donde los recursos computacionales son limitados, como en dispositivos móviles.

Trabajos recientes como DistilBERT (Sanh et al., 2019) han mostrado cómo se puede comprimir modelos preentrenados sin perder precisión significativa, lo que podría aplicarse en modelos de NMT de gran escala.

8 REFERENCIAS BIBLIOGRAFICAS

- Anastasopoulos, A., & Chiang, D. (2018). Tied Multitask Learning for Neural Speech Translation.

 *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers), 82–91.

 https://doi.org/10.18653/v1/N18-1008
- Bahdanau, D., Cho, K. H., & Bengio, Y. (2015). Neural machine translation by jointly learning to align and translate. *3rd International Conference on Learning Representations, ICLR 2015 Conference Track Proceedings*.
- Beck, M., Pöppel, K., Spanring, M., Auer, A., Prudnikova, O., Kopp, M., Klambauer, G., Brandstetter, J., & Hochreiter, S. (2024). *xLSTM: Extended Long Short-Term Memory*.
- Bhojanapalli, S., Yun, C., Rawat, A. S., Reddi, S., & Kumar, S. (2020). Low-rank bottleneck in multi-head attention models. *37th International Conference on Machine Learning, ICML 2020, PartF168147-2*.
- Bishop, C. M., & Bishop, H. (2024). *Deep Learning*. Springer International Publishing. https://doi.org/10.1007/978-3-031-45468-4
- Brauwers, G., & Frasincar, F. (2023). A General Survey on Attention Mechanisms in Deep Learning. *IEEE Transactions on Knowledge and Data Engineering*, *35*(4). https://doi.org/10.1109/TKDE.2021.3126456
- Chaudhari, S., Mithal, V., Polatkan, G., & Ramanath, R. (2021). An Attentive Survey of Attention Models. *ACM Transactions on Intelligent Systems and Technology*, *12*(5). https://doi.org/10.1145/3465055
- Csordás, R., Irie, K., Schmidhuber, J., Potts, C., & Manning, C. D. (2024). *MoEUT: Mixture-of-Experts Universal Transformers*. https://arxiv.org/abs/2405.16039
- Dehghani, M., Gouws, S., Vinyals, O., Uszkoreit, J., & Kaiser, Ł. (2019). *Universal Transformers*. https://arxiv.org/abs/1807.03819
- Dzmitry Bahdanau, Kyung Hyun Cho, & Yoshua Bengio. (2014). Neural Machine Translation by Jointly Learning to Align and Translate. *3rd International Conference on Learning Representations, ICLR 2015 Conference Track Proceedings*.

- Eisenstein, J. (2019). *Introduction to Natural Language Processing*. MIT Press. https://books.google.co.cr/books?id=72yuDwAAQBAJ
- Fedus, W., Zoph, B., & Shazeer, N. (2022). *Switch Transformers: Scaling to Trillion Parameter Models with Simple and Efficient Sparsity*. https://arxiv.org/abs/2101.03961
- Firth, J. R. (1957). A Synopsis of Linguistic Theory 1930-55. In *Studies in Linguistic Analysis: Special Volume* of the Philological Society.
- Gammerman, A., Vovk, V., & Vapnik, V. (2013). Learning by Transduction.
- Graves, A. (2012, November 14). *Sequence Transduction with Recurrent Neural Networks*.
- Grefenstette, E., Hermann, K. M., Suleyman, M., & Blunsom, P. (2015). *Learning to Transduce with Unbounded Memory*.
- Gu, A., & Dao, T. (2024). *Mamba: Linear-Time Sequence Modeling with Selective State Spaces*. https://arxiv.org/abs/2312.00752
- Gu, A., Dao, T., Ermon, S., Rudra, A., & Re, C. (2020). *HiPPO: Recurrent Memory with Optimal Polynomial Projections*.
- Gu, A., Goel, K., & Ré, C. (2021). Efficiently Modeling Long Sequences with Structured State Spaces.
- Harris, Z. S. (1954). Distributional Structure. WORD, *10*(2–3). https://doi.org/10.1080/00437956.1954.11659520
- Jaitly, N., Sussillo, D., Le, Q. V., Vinyals, O., Sutskever, I., & Bengio, S. (2015). A Neural Transducer.
- Kamath, U., Graham, K. L., & Emara, W. (2022). *Transformers for Machine Learning*. Chapman and Hall/CRC. https://doi.org/10.1201/9781003170082
- Katharopoulos, A., Vyas, A., Pappas, N., & Fleuret, F. (2020). *Transformers are RNNs: Fast Autoregressive Transformers with Linear Attention*.
- Kitaev, N., Kaiser, Ł., & Levskaya, A. (2020). *Reformer: The Efficient Transformer*. https://arxiv.org/abs/2001.04451
- Koehn, P. (2020). *Neural Machine Translation*. Cambridge University Press. https://doi.org/10.1017/9781108608480

- Le, H., Pino, J., Wang, C., Gu, J., Schwab, D., & Besacier, L. (2020). Dual-decoder Transformer for Joint Automatic Speech Recognition and Multilingual Speech Translation. *COLING 2020 28th International Conference on Computational Linguistics, Proceedings of the Conference*. https://doi.org/10.18653/v1/2020.coling-main.314
- Lieber, O., Lenz, B., Bata, H., Cohen, G., Osin, J., Dalmedigos, I., Safahi, E., Meirom, S., Belinkov, Y., Shalev-Shwartz, S., Abend, O., Alon, R., Asida, T., Bergman, A., Glozman, R., Gokhman, M., Manevich, A., Ratner, N., Rozen, N., ... Shoham, Y. (2024). *Jamba: A Hybrid Transformer-Mamba Language Model*. https://arxiv.org/abs/2403.19887
- Lin, P., Martins, A. F. T., & Schütze, H. (2024). *A Recipe of Parallel Corpora Exploitation for Multilingual Large Language Models*.
- Lin, T., Wang, Y., Liu, X., & Qiu, X. (2022). A survey of transformers. *AI Open*, *3*. https://doi.org/10.1016/j.aiopen.2022.10.001
- Lison, P., & Tiedemann, J. (2016). OpenSubtitles2016: Extracting Large Parallel Corpora from Movie and TV Subtitles. In N. Calzolari, K. Choukri, T. Declerck, S. Goggi, M. Grobelnik, B. Maegaard, J. Mariani, H. Mazo, A. Moreno, J. Odijk, & S. Piperidis (Eds.), *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC'16)* (pp. 923–929). European Language Resources Association (ELRA). https://aclanthology.org/L16-1147
- Liu, Z., Lin, Y., & Su, M. (2023). Representation Learning for Natural Language Processing. In *Representation Learning for Natural Language Processing*. https://doi.org/10.1007/978-981-15-5573-2
- Ma, X., Zhou, C., Kong, X., He, J., Gui, L., Neubig, G., May, J., & Zettlemoyer, L. (2022). *Mega: Moving Average Equipped Gated Attention*.
- Mehta, H., Gupta, A., Cutkosky, A., & Neyshabur, B. (2022). *Long Range Language Modeling via Gated State Spaces*.
- Pióro, M., Ciebiera, K., Król, K., Ludziejewski, J., Krutul, M., Krajewski, J., Antoniak, S., Miłoś, P., Cygan, M., & Jaszczur, S. (2024). *MoE-Mamba: Efficient Selective State Space Models with Mixture of Experts*. https://arxiv.org/abs/2401.04081
- Prince, S. J. D. (2023). *Understanding Deep Learning*. MIT Press. http://udlbook.com

- Smith, J. T. H., Warrington, A., & Linderman, S. W. (2022). *Simplified State Space Layers for Sequence Modeling*.
- So, D. R., Liang, C., & Le, Q. V. (2019). *The Evolved Transformer*.
- Sperber, M., Setiawan, H., Gollan, C., Nallasamy, U., & Paulik, M. (2020). Consistent Transcription and Translation of Speech. *Transactions of the Association for Computational Linguistics*, *8*, 695–709. https://doi.org/10.1162/tacl_a_00340
- Sun, J., Liu, X., Mei, X., Kılıç, V., Plumbley, M. D., & Wang, W. (2023). Dual Transformer Decoder based Features Fusion Network for Automated Audio Captioning. *Proceedings of the Annual Conference of the International Speech Communication Association, INTERSPEECH*, 2023-August. https://doi.org/10.21437/Interspeech.2023-943
- Sutskever, I., Vinyals, O., & Le, Q. V. (2014). Sequence to sequence learning with neural networks.

 **Advances in Neural Information Processing Systems, 4(January).
- Tay, Y., Dehghani, M., Bahri, D., & Metzler, D. (2022). Efficient Transformers: A Survey. *ACM Computing Surveys*, *55*(6). https://doi.org/10.1145/3530811
- Vardasbi, A., Pires, T. P., Schmidt, R. M., & Peitz, S. (2023a). *State Spaces Aren't Enough: Machine Translation Needs Attention*.
- Vardasbi, A., Pires, T. P., Schmidt, R., & Peitz, S. (2023b). State Spaces Aren't Enough: Machine Translation Needs Attention. In M. Nurminen, J. Brenner, M. Koponen, S. Latomaa, M. Mikhailov, F. Schierl, T. Ranasinghe, E. Vanmassenhove, S. A. Vidal, N. Aranberri, M. Nunziatini, C. P. Escart\'\in, M. Forcada, M. Popovic, C. Scarton, & H. Moniz (Eds.), *Proceedings of the 24th Annual Conference of the European Association for Machine Translation* (pp. 205–216). European Association for Machine Translation. https://aclanthology.org/2023.eamt-1.20
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., & Polosukhin, I. (2017). Attention is all you need. *Advances in Neural Information Processing Systems*, *2017–December*.
- Wang, S., Tu, Z., Tan, Z., Wang, W., Sun, M., & Liu, Y. (2021). *Language Models are Good Translators*.
- Wang, Z. (2023). *Learning Transductions and Alignments with RNN Seg2seg Models*.
- Xiao, T., & Zhu, J. (2023). *Introduction to Transformers: an NLP Perspective*.

- Yang, S., Wang, B., Shen, Y., Panda, R., & Kim, Y. (2024). *Gated Linear Attention Transformers with Hardware-Efficient Training*. https://arxiv.org/abs/2312.06635
- Yu, Z., Wang, J., Yu, L.-C., & Zhang, X. (2022). Dual-Encoder Transformers with Cross-modal Alignment for Multimodal Aspect-based Sentiment Analysis. *Proceedings of the 2nd Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 12th International Joint Conference on Natural Language Processing*, 1.
- Zhang, Z., Shao, N., Gao, C., Miao, R., Yang, Q., & Shao, J. (2022). Mixhead: Breaking the low-rank bottleneck in multi-head attention language models. *Knowledge-Based Systems*, *240*, 108075. https://doi.org/10.1016/j.knosys.2021.108075
- Zhou, Z. H. (2012). Ensemble methods: Foundations and algorithms. In *Ensemble Methods: Foundations* and Algorithms. https://doi.org/10.1201/b12207
- Ziemski, M., Junczys-Dowmunt, M., & Pouliquen, B. (2016). The United Nations Parallel Corpus v1.0. In N. Calzolari, K. Choukri, T. Declerck, S. Goggi, M. Grobelnik, B. Maegaard, J. Mariani, H. Mazo, A. Moreno, J. Odijk, & S. Piperidis (Eds.), *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC'16)* (pp. 3530–3534). European Language Resources Association (ELRA). https://aclanthology.org/L16-1561
- Zoph, B., Bello, I., Kumar, S., Du, N., Huang, Y., Dean, J., Shazeer, N., & Fedus, W. (2022). *ST-MoE:*Designing Stable and Transferable Sparse Expert Models.

9 ANEXOS

Anexo No. 1. Arquitecturas de modelos y sus hiperparámetros

Modelo 1: Vanilla Transformer (Clase: NMT_Vashwani)

- **Descripción:** Basado en la arquitectura Transformer original de Vaswani et al. (2017).
- Parámetros Principales (instanciados):
 - num_encoder_layers: 6
 - num_decoder_layers: 6
 - o emb_size (d_model): 512
 - o nhead (número de cabezas de atención): 8
 - o ff_dim (dimensión de la capa feedforward): 512
 - o dropout: 0.1 (valor por defecto en la definición de la clase nn.Transformer interna)
 - lr: 0.0001 (Variable global TFMR_LR)
- **Instancias:** model_Vainilla_UNPC y model_Vainilla_OSPC (difieren en src_vocab_size, tgt_vocab_size).

Modelo 2: Universal Transformer (UT) (Clase: UT_NMT)

- **Descripción:** Implementación de un Transformer Recurrente (Universal Transformer).
- Parámetros Principales (definidos en XT_Config dentro de la clase):
 - o emb_dim: 512
 - o hidden_dim: 1024
 - o pff_dim: 1024
 - o n_layers: 4
 - o n_heads: 8
 - dropout_ratio: 0.1
 - o max_len: Variable (UN_MAX_LEN o OS_MAX_LEN según dataset)

- o pad_id: 0
- o Ir: 0.0001 (Variable global TFMR_LR)
- **Instancias:** model_UT_UNPC y model_UT_OSPC (differen en vocab_size y max_length).

Modelo 3: S5 (Simplified Space-State Model - SSM) (Clase: SSM_NMT)

- **Descripción:** Modelo basado en espacios de estado estructurados (SSM), utilizando la capa Mamba.
- Parámetros Principales (instanciados en SSM_NMT):
 - o d_model: 512
 - n: 256 (Parece no usarse directamente en las capas internas según el código, pero está en el init)
 - o block_count: 4 (Usado como expand en la capa Mamba interna)
 - o dropout: 0.3
 - max_length: Variable (UN_MAX_LEN o OS_MAX_LEN según dataset)
 - lr: 0.0001 (Variable global TFMR_LR)
- Parámetros Internos (Capa Mamba en SSM_Encoder/Decoder):
 - o d_state: 16 (Factor de expansión del estado SSM)
 - d_conv: 4 (Ancho de la convolución local)
 - o expand: 2 (Factor de expansión del bloque)
- Parámetros Internos (Capa MultiHeadAttention):
 - o n_heads: 8
- Instancias: model_SSM_UNPC y model_SSM_OSPC (differen en vocab_size y max_length).

Modelo 4: ST-MoE Transformer (Clase: MoE_NMT)

- **Descripción:** Transformer que incorpora un bloque Mixture-of-Experts (MoE) disperso.
- Parámetros Principales (instanciados o definidos internamente):
 - o d_model: 512

- n_heads (para MultiHeadAttention): 8
- num_experts (para MoE): 8
- gating_top_n (para MoE): 2
- o dropout: 0.1 (definido en las capas internas)
- max_seq_length: Variable (UN_MAX_LEN o OS_MAX_LEN según dataset)
- o Ir: 0.0001 (Variable global TFMR_LR)
- Instancias: model_MOE_UNPC y model_MOE_OSPC (difieren en vocab_size y max_length).

Modelo 5: Gated State Spaces (GSS) (Clase: GSS_NMT)

- **Descripción:** Modelo basado en Gated State Spaces.
- Parámetros Principales (instanciados o definidos internamente):
 - o d_model: 512
 - o n_heads (para MultiHeadAttention): 8
 - depth (número de capas GSS): 3
 - dss_kernel_H (para GSS): 512
 - o dss kernel N (para GSS): 256
 - dim_expansion_factor (para GSS): 4
 - dropout: 0.1 (definido en las capas internas)
 - max_seq_length: Variable (UN_MAX_LEN o OS_MAX_LEN según dataset)
 - lr: 0.0001 (Variable global TFMR_LR)
- **Instancias:** model GSS UNPC y model GSS OSPC (differen en vocab size y max length).

Modelo 6: Moving Average Equipped Gated Attention (MEGA) (Clase: MEGA_NMT)

- **Descripción:** Modelo que utiliza capas MEGA (Moving Average Equipped Gated Attention).
- Parámetros Principales (instanciados o definidos internamente):
 - o d_model: 512

- n_heads (para MultiHeadAttention): 8
- o depth (número de capas MegaLayer): 3
- ema_heads (para MegaLayer): 16
- o attn_dim_qk (para MegaLayer): 64
- attn_dim_value (para MegaLayer): 256
- o dropout: 0.1 (definido en las capas internas)
- max_seq_length: Variable (UN_MAX_LEN o OS_MAX_LEN según dataset)
- o lr: 0.0001 (Variable global TFMR_LR)
- **Instancias:** model_MEGA_UNPC y model_MEGA_OSPC (difieren en vocab_size y max_length).

Anexo No. 2. Resultados de pérdida, precisión y perplejidad

model_name	epoch	dataset	trn_loss	trn_accm	trn_pplx	val_loss	val_accm	val_pplx
Vainilla	1	UNPC	5,0538	0,2222	156,6165	4,4282	0,2929	83,7805
Vainilla	2	UNPC	4,2816	0,3075	72,3561	4,0245	0,3397	55,9523
Vainilla	3	UNPC	3,9974	0,3444	54,4564	3,8199	0,3710	45,5996
Vainilla	4	UNPC	3,8323	0,3674	46,1686	3,6826	0,3906	39,7496
Vainilla	5	UNPC	3,7302	0,3845	41,6874	3,6835	0,3999	39,7854
Vainilla	6	UNPC	3,7692	0,3968	43,3454	3,7371	0,4094	41,9761
Vainilla	7	UNPC	3,7551	0,4080	42,7385	3,6738	0,4242	39,4013
Vainilla	8	UNPC	3,7162	0,4189	41,1079	3,6507	0,4332	38,5016
Vainilla	9	UNPC	3,6816	0,4288	39,7099	3,6582	0,4362	38,7915
Vainilla	10	UNPC	3,6435	0,4379	38,2254	3,6265	0,4436	37,5811
Vainilla	11	UNPC	3,6022	0,4467	36,6788	3,5427	0,4586	34,5601
Vainilla	12	UNPC	3,5656	0,4545	35,3607	3,5214	0,4642	33,8318
Vainilla	13	UNPC	3,5355	0,4601	34,3122	3,5089	0,4697	33,4115
Vainilla	14	UNPC	3,5067	0,4665	33,3381	3,4837	0,4725	32,5800
Vainilla	15	UNPC	3,4778	0,4724	32,3884	3,4539	0,4810	31,6235
Vainilla	16	UNPC	3,4460	0,4778	31,3746	3,4276	0,4848	30,8026
Vainilla	17	UNPC	3,4163	0,4829	30,4565	3,4103	0,4868	30,2743
Vainilla	18	UNPC	3,3905	0,4871	29,6808	3,3875	0,4920	29,5919
Vainilla	19	UNPC	3,3642	0,4911	28,9104	3,3689	0,4942	29,0466
Vainilla	20	UNPC	3,3369	0,4943	28,1318	3,3428	0,4971	28,2983
Vainilla	1	OSPC	4,2428	0,2847	69,6025	4,7862	0,2796	119,8451
Vainilla	2	OSPC	4,0707	0,2977	58,5980	4,8994	0,2735	134,2092
Vainilla	3	OSPC	4,0376	0,3016	56,6901	4,8231	0,2798	124,3500
Vainilla	4	OSPC	4,0302	0,3026	56,2722	4,7966	0,2776	121,0980
Vainilla	5	OSPC	4,0482	0,3022	57,2942	4,7983	0,2667	121,3040
Vainilla	6	OSPC	4,0025	0,3452	54,7348	4,5925	0,3127	98,7410
Vainilla	7	OSPC	4,0372	0,3418	56,6675	4,5019	0,3312	90,1883
Vainilla	8	OSPC	4,0535	0,3330	57,5987	4,4756	0,3163	87,8473
Vainilla	9	OSPC	4,0831	0,3424	59,3291	4,4995	0,3101	89,9721
Vainilla	10	OSPC	4,0815	0,3433	59,2343	4,5552	0,3095	95,1258
Vainilla	11	OSPC	4,0776	0,3423	59,0037	4,5388	0,3095	93,5784
Vainilla	12	OSPC	4,0457	0,3418	57,1512	4,5802	0,3095	97,5339
Vainilla	13	OSPC	4,0588	0,3438	57,9048	4,4025	0,3096	81,6548
Vainilla	14	OSPC	4,0508	0,3443	57,4434	4,4530	0,3103	85,8842
Vainilla	15	OSPC	4,0541	0,3445	57,6333	4,4715	0,3095	87,4879
Vainilla	16	OSPC	4,0504	0,3446	57,4204	4,4347	0,3095	84,3268
Vainilla	17	OSPC	4,0397	0,3446	56,8093	4,4147	0,3095	82,6570
Vainilla	18	OSPC	4,0454	0,3449	57,1340	4,4178	0,3095	82,9137
Vainilla	19	OSPC	4,0481	0,3450	57,2885	4,4315	0,3103	84,0574
Vainilla	20	OSPC	4,0536	0,3453	57,6045	4,4250	0,3118	83,5128
UT	1	UNPC	4,0222	0,3726	55,8238	3,2661	0,4754	26,2089

model_name	epoch	dataset	trn_loss	trn_accm	trn_pplx	val_loss	val_accm	val_pplx
UT	2	UNPC	3,0724	0,5014	21,5937	2,8494	0,5382	17,2774
UT	3	UNPC	2,7442	0,5531	15,5522	2,6285	0,5758	13,8530
UT	4	UNPC	2,5576	0,5854	12,9048	2,5344	0,5979	12,6089
UT	5	UNPC	2,5206	0,6058	12,4361	2,5352	0,6131	12,6190
UT	6	UNPC	2,4694	0,6237	11,8154	2,5099	0,6265	12,3037
UT	7	UNPC	2,4233	0,6379	11,2830	2,4798	0,6367	11,9389
UT	8	UNPC	2,3821	0,6495	10,8276	2,4622	0,6444	11,7306
UT	9	UNPC	2,3460	0,6590	10,4437	2,4481	0,6500	11,5663
UT	10	UNPC	2,3135	0,6667	10,1097	2,4273	0,6558	11,3283
UT	11	UNPC	2,2787	0,6740	9,7640	2,4107	0,6589	11,1418
UT	12	UNPC	2,2504	0,6797	9,4915	2,3932	0,6633	10,9485
UT	13	UNPC	2,2204	0,6852	9,2110	2,3850	0,6651	10,8591
UT	14	UNPC	2,1950	0,6895	8,9800	2,3679	0,6680	10,6750
UT	15	UNPC	2,1671	0,6942	8,7329	2,3617	0,6689	10,6090
UT	16	UNPC	2,1436	0,6984	8,5301	2,3388	0,6721	10,3688
UT	17	UNPC	2,1185	0,7021	8,3187	2,3438	0,6720	10,4208
UT	18	UNPC	2,0955	0,7056	8,1295	2,3309	0,6745	10,2872
UT	19	UNPC	2,0769	0,7086	7,9797	2,3214	0,6762	10,1899
UT	20	UNPC	2,0512	0,7120	7,7772	2,3156	0,6774	10,1310
UT	1	OSPC	3,9671	0,3242	52,8311	5,0115	0,3750	150,1298
UT	2	OSPC	3,3570	0,4923	28,7030	4,8081	0,5006	122,4986
UT	3	OSPC	3,1541	0,5173	23,4319	4,6909	0,5160	108,9512
UT	4	OSPC	3,0300	0,5340	20,6972	4,6502	0,5239	104,6059
UT	5	OSPC	3,1156	0,5443	22,5470	4,4859	0,5247	88,7568
UT	6	OSPC	3,2814	0,5543	26,6130	4,6120	0,5301	100,6853
UT	7	OSPC	3,3586	0,5634	28,7489	4,5627	0,5356	95,8419
UT	8	OSPC	3,4256	0,5721	30,7411	4,6614	0,5373	105,7841
UT	9	OSPC	3,4749	0,5792	32,2946	4,3235	0,5410	75,4522
UT	10	OSPC	3,5010	0,5860	33,1486	4,2749	0,5414	71,8730
UT	11	OSPC	3,5130	0,5918	33,5488	4,3057	0,5444	74,1211
UT	12	OSPC	3,5142	0,5975	33,5890	4,3123	0,5463	74,6119
UT	13	OSPC	3,5055	0,6032	33,2981	4,2915	0,5452	73,0760
UT	14	OSPC	3,4857	0,6077	32,6453	4,3191	0,5455	75,1210
UT	15	OSPC	3,4690	0,6125	32,1046	4,3248	0,5470	75,5504
UT	16	OSPC	3,2516	0,6536	25,8316	4,2947	0,5619	73,3102
UT	17	OSPC	3,2160	0,6663	24,9282	4,3168	0,5621	74,9484
UT	18	OSPC	3,1947	0,6734	24,4029	4,3331	0,5617	76,1801
UT	19	OSPC	3,1762	0,6785	23,9555	4,3393	0,5607	76,6539
UT	20	OSPC	3,1556	0,6827	23,4671	4,3353	0,5608	76,3479
MAMBA	1	UNPC	4,9851	0,2422	146,2182	4,0961	0,3624	60,1054
MAMBA	2	UNPC	3,9020	0,3930	49,5014	3,5387	0,4593	34,4221
MAMBA	3	UNPC	3,5228	0,4522	33,8792	3,2585	0,4972	26,0105
MAMBA	4	UNPC	3,3172	0,4832	27,5830	3,1208	0,5203	22,6645

model_name	epoch	dataset	trn_loss	trn_accm	trn_pplx	val_loss	val_accm	val_pplx
MAMBA	5	UNPC	3,1834	0,5036	24,1287	2,9993	0,5390	20,0715
MAMBA	6	UNPC	3,1430	0,5178	23,1733	3,1116	0,5478	22,4569
MAMBA	7	UNPC	3,2015	0,5269	24,5694	3,1151	0,5542	22,5357
MAMBA	8	UNPC	3,1725	0,5361	23,8671	3,0713	0,5642	21,5699
MAMBA	9	UNPC	3,1355	0,5449	23,0001	3,0519	0,5696	21,1555
MAMBA	10	UNPC	3,0968	0,5528	22,1270	3,0238	0,5780	20,5693
MAMBA	11	UNPC	3,0632	0,5596	21,3959	2,9817	0,5836	19,7213
MAMBA	12	UNPC	3,0337	0,5655	20,7740	2,9609	0,5875	19,3153
MAMBA	13	UNPC	3,0060	0,5708	20,2064	2,9359	0,5918	18,8385
MAMBA	14	UNPC	2,9811	0,5753	19,7095	2,9293	0,5935	18,7145
MAMBA	15	UNPC	2,9627	0,5787	19,3501	2,9050	0,5969	18,2652
MAMBA	16	UNPC	2,9393	0,5827	18,9026	2,8916	0,6028	18,0221
MAMBA	17	UNPC	2,9207	0,5858	18,5543	2,8738	0,6037	17,7042
MAMBA	18	UNPC	2,9021	0,5891	18,2124	2,8654	0,6075	17,5561
MAMBA	19	UNPC	2,8848	0,5918	17,9000	2,8192	0,6103	16,7634
MAMBA	20	UNPC	2,8489	0,5957	17,2688	2,8048	0,6128	16,5238
MAMBA	1	OSPC	4,6662	0,3386	106,2931	5,1590	0,3099	173,9904
MAMBA	2	OSPC	4,0170	0,4177	55,5343	4,8766	0,3588	131,1839
MAMBA	3	OSPC	3,7716	0,4517	43,4495	4,6794	0,3804	107,7054
MAMBA	4	OSPC	3,6217	0,4711	37,4011	4,5143	0,3994	91,3136
MAMBA	5	OSPC	3,5186	0,4837	33,7372	5,4666	0,4039	236,6542
MAMBA	6	OSPC	3,4857	0,4911	32,6453	4,4977	0,4117	89,8103
MAMBA	7	OSPC	3,9333	0,4952	51,0752	4,9507	0,4122	141,2738
MAMBA	8	OSPC	3,5622	0,4991	35,2406	4,6309	0,4182	102,6064
MAMBA	9	OSPC	3,6084	0,5048	36,9070	4,4249	0,4215	83,5045
MAMBA	10	OSPC	3,6382	0,5096	38,0233	4,4773	0,4237	87,9968
MAMBA	11	OSPC	3,6616	0,5141	38,9236	4,4654	0,4326	86,9558
MAMBA	12	OSPC	3,6611	0,5178	38,9041	4,4902	0,4329	89,1393
MAMBA	13	OSPC	3,6682	0,5216	39,1813	4,5107	0,4331	90,9855
MAMBA	14	OSPC	3,6719	0,5250	39,3266	4,4274	0,4392	83,7135
MAMBA	15	OSPC	3,6596	0,5280	38,8458	4,3207	0,4399	75,2413
MAMBA	16	OSPC	3,6525	0,5307	38,5710	4,3349	0,4425	76,3173
MAMBA	17	OSPC	3,5101	0,5600	33,4516	4,2372	0,4651	69,2138
MAMBA	18	OSPC	3,5140	0,5551	33,5823	4,2546	0,4570	70,4286
MAMBA	19	OSPC	3,5231	0,5495	33,8893	4,1745	0,4528	65,0073
MAMBA	20	OSPC	3,5316	0,5478	34,1786	4,1483	0,4544	63,3263
МоЕ	1	UNPC	4,2959	0,3501	73,3982	3,3545	0,4785	28,6313
MoE	2	UNPC	3,2911	0,4875	26,8724	2,9351	0,5398	18,8234
MoE	3	UNPC	2,9764	0,5314	19,6171	2,7564	0,5670	15,7431
MoE	4	UNPC	2,8184	0,5566	16,7500	2,7207	0,5849	15,1910
MoE	5	UNPC	2,7766	0,5752	16,0643	2,7142	0,5971	15,0925
MoE	6	UNPC	2,7182	0,5919	15,1530	2,6947	0,6076	14,8011
MoE	7	UNPC	2,6712	0,6052	14,4573	2,6842	0,6152	14,6465

model_name	epoch	dataset	trn_loss	trn_accm	trn_pplx	val_loss	val_accm	val_pplx
MoE	8	UNPC	2,6287	0,6170	13,8557	2,6748	0,6201	14,5094
MoE	9	UNPC	2,5898	0,6263	13,3271	2,6773	0,6257	14,5458
MoE	10	UNPC	2,5590	0,6349	12,9229	2,6553	0,6307	14,2293
MoE	11	UNPC	2,5233	0,6428	12,4697	2,6658	0,6330	14,3794
MoE	12	UNPC	2,4909	0,6496	12,0721	2,6624	0,6352	14,3306
MoE	13	UNPC	2,4556	0,6563	11,6534	2,6446	0,6392	14,0778
MoE	14	UNPC	2,4220	0,6624	11,2684	2,6516	0,6396	14,1767
MoE	15	UNPC	2,3817	0,6684	10,8233	2,6161	0,6428	13,6823
MoE	16	UNPC	2,3369	0,6732	10,3491	2,6119	0,6454	13,6249
MoE	17	UNPC	2,3117	0,6783	10,0916	2,6077	0,6464	13,5678
MoE	18	UNPC	2,2850	0,6831	9,8257	2,6079	0,6468	13,5705
MoE	19	UNPC	2,2582	0,6881	9,5659	2,6196	0,6482	13,7302
MoE	20	UNPC	2,2387	0,6924	9,3811	2,6223	0,6489	13,7674
MoE	1	OSPC	4,6760	0,4096	107,3399	4,6543	0,3803	105,0357
MoE	2	OSPC	4,1090	0,4773	60,8858	4,4093	0,4079	82,2119
MoE	3	OSPC	3,8970	0,5029	49,2545	4,2743	0,4267	71,8298
MoE	4	OSPC	3,5720	0,5193	35,5877	4,3493	0,4374	77,4242
MoE	5	OSPC	3,6415	0,5338	38,1490	5,4668	0,4437	236,7015
MoE	6	OSPC	3,7040	0,5467	40,6094	4,4756	0,4517	87,8473
MoE	7	OSPC	3,7773	0,5575	43,6979	4,5886	0,4561	98,3566
MoE	8	OSPC	3,5448	0,5672	34,6328	4,2826	0,4603	72,4285
MoE	9	OSPC	3,5865	0,5761	36,1075	4,1257	0,4605	61,9111
MoE	10	OSPC	3,6156	0,5846	37,1736	4,1894	0,4645	65,9832
MoE	11	OSPC	3,6186	0,5925	37,2853	4,2165	0,4654	67,7958
MoE	12	OSPC	3,6113	0,5998	37,0141	4,2407	0,4673	69,4565
MoE	13	OSPC	3,5852	0,6079	36,0606	4,2608	0,4692	70,8667
MoE	14	OSPC	3,5534	0,6151	34,9319	4,2036	0,4686	66,9268
MoE	15	OSPC	3,3366	0,6588	28,1233	4,1150	0,4822	61,2522
MoE	16	OSPC	3,3016	0,6732	27,1561	4,1512	0,4827	63,5102
MoE	17	OSPC	3,2794	0,6803	26,5598	4,1730	0,4813	64,9099
MoE	18	OSPC	3,2546	0,6863	25,9092	4,1907	0,4815	66,0690
MoE	19	OSPC	3,2328	0,6910	25,3505	4,0976	0,4815	60,1956
MoE	20	OSPC	3,2138	0,6957	24,8734	4,1175	0,4804	61,4055
GSS	1	UNPC	5,3583	0,1822	212,3636	5,1831	0,1978	178,2345
GSS	2	UNPC	5,2502	0,1902	190,6044	5,2712	0,1919	194,6494
GSS	3	UNPC	5,3049	0,1870	201,3209	5,1896	0,1973	179,3968
GSS	4	UNPC	5,3560	0,1822	211,8757	5,7233	0,1499	305,9128
GSS	5	UNPC	5,5335	0,1681	253,0280	5,4238	0,1787	226,7391
GSS	6	UNPC	5,4146	0,1807	224,6627	5,3848	0,1906	218,0665
GSS	7	UNPC	5,4554	0,1844	234,0185	5,6001	0,1747	270,4535
GSS	8	UNPC	5,6075	0,1740	272,4622	5,5725	0,1833	263,0910
GSS	9	UNPC	5,5817	0,1802	265,5226	5,6429	0,1833	282,2801
GSS	10	UNPC	5,7314	0,1709	308,4007	5,5005	0,1976	244,8143

model_name	epoch	dataset	trn_loss	trn_accm	trn_pplx	val_loss	val_accm	val_pplx
GSS	11	UNPC	5,6188	0,1842	275,5585	5,6107	0,1910	273,3355
GSS	12	UNPC	5,6388	0,1852	281,1252	5,5900	0,1958	267,7356
GSS	13	UNPC	5,4356	0,2077	229,4305	5,4627	0,2141	235,7330
GSS	14	UNPC	5,3787	0,2138	216,7403	5,4294	0,2179	228,0124
GSS	15	UNPC	5,3656	0,2153	213,9195	5,4204	0,2199	225,9695
GSS	16	UNPC	5,3429	0,2180	209,1183	5,3983	0,2205	221,0303
GSS	17	UNPC	5,3347	0,2190	207,4105	5,3867	0,2237	218,4812
GSS	18	UNPC	5,3183	0,2212	204,0367	5,3726	0,2245	215,4222
GSS	19	UNPC	5,3168	0,2216	203,7309	5,3725	0,2250	215,4007
GSS	20	UNPC	5,3115	0,2227	202,6540	5,3566	0,2259	212,0029
GSS	1	OSPC	5,0914	0,3071	162,6174	5,9101	0,2444	368,7430
GSS	2	OSPC	4,8791	0,3242	131,5122	5,8162	0,2507	335,6940
GSS	3	OSPC	4,7708	0,3329	118,0136	5,7176	0,2602	304,1740
GSS	4	OSPC	4,7034	0,3383	110,3216	5,6932	0,2653	296,8420
GSS	5	OSPC	4,6764	0,3414	107,3828	6,6420	0,2684	766,6267
GSS	6	OSPC	4,6700	0,3461	106,6977	5,6804	0,2691	293,0666
GSS	7	OSPC	4,6788	0,3442	107,6408	6,1019	0,2709	446,7057
GSS	8	OSPC	4,6217	0,3501	101,6667	5,8750	0,2672	356,0247
GSS	9	OSPC	4,6852	0,3508	108,3319	5,5995	0,2772	270,2912
GSS	10	OSPC	4,6750	0,3539	107,2326	5,6566	0,2773	286,1740
GSS	11	OSPC	4,6447	0,3573	104,0322	5,6841	0,2785	294,1530
GSS	12	OSPC	4,6826	0,3588	108,0506	5,7249	0,2801	306,4026
GSS	13	OSPC	4,6022	0,3614	99,7034	5,7635	0,2823	318,4610
GSS	14	OSPC	4,6120	0,3636	100,6853	5,7049	0,2819	300,3354
GSS	15	OSPC	4,6115	0,3660	100,6350	5,6256	0,2812	277,4387
GSS	16	OSPC	4,6138	0,3675	100,8667	5,6260	0,2829	277,5497
GSS	17	OSPC	4,4377	0,3872	84,5802	5,5771	0,2956	264,3040
GSS	18	OSPC	4,4249	0,3894	83,5045	5,5839	0,2953	266,1074
GSS	19	OSPC	4,4216	0,3898	83,2293	5,4770	0,2950	239,1282
GSS	20	OSPC	4,4290	0,3893	83,8475	5,4902	0,2949	242,3057
MEGA	1	UNPC	4,8158	0,2583	123,4455	4,1497	0,3431	63,4150
MEGA	2	UNPC	3,8622	0,3892	47,5699	3,4998	0,4515	33,1088
MEGA	3	UNPC	3,3744	0,4702	29,2068	3,1529	0,5116	23,4038
MEGA	4	UNPC	3,0774	0,5192	21,7019	2,9336	0,5462	18,7952
MEGA	5	UNPC	2,8950	0,5500	18,0835	2,8114	0,5712	16,6332
MEGA	6	UNPC	2,8639	0,5718	17,5298	2,8948	0,5877	18,0799
MEGA	7	UNPC	2,8636	0,5886	17,5245	2,8652	0,5977	17,5526
MEGA	8	UNPC	2,8021	0,6022	16,4792	2,8156	0,6070	16,7032
MEGA	9	UNPC	2,7490	0,6135	15,6270	2,7861	0,6158	16,2176
MEGA	10	UNPC	2,7027	0,6230	14,9200	2,7445	0,6238	15,5568
MEGA	11	UNPC	2,6601	0,6315	14,2977	2,7182	0,6301	15,1530
MEGA	12	UNPC	2,6193	0,6386	13,7261	2,6917	0,6344	14,7567
MEGA	13	UNPC	2,5910	0,6444	13,3431	2,6886	0,6361	14,7111

model_name	epoch	dataset	trn_loss	trn_accm	trn_pplx	val_loss	val_accm	val_pplx
MEGA	14	UNPC	2,5549	0,6505	12,8700	2,6603	0,6412	14,3006
MEGA	15	UNPC	2,5230	0,6555	12,4659	2,6449	0,6438	14,0820
MEGA	16	UNPC	2,4937	0,6601	12,1060	2,6301	0,6456	13,8752
MEGA	17	UNPC	2,4647	0,6644	11,7600	2,6189	0,6484	13,7206
MEGA	18	UNPC	2,4218	0,6688	11,2661	2,5656	0,6501	13,0085
MEGA	19	UNPC	2,3892	0,6717	10,9048	2,5876	0,6501	13,2978
MEGA	20	UNPC	2,3734	0,6749	10,7338	2,5633	0,6520	12,9786
MEGA	1	OSPC	4,6465	0,3592	104,2196	5,2478	0,2812	190,1475
MEGA	2	OSPC	4,7284	0,4160	113,1144	4,9954	0,3167	147,7320
MEGA	3	OSPC	4,6572	0,4444	105,3407	4,8271	0,3387	124,8484
MEGA	4	OSPC	4,5362	0,4666	93,3355	4,6838	0,3552	108,1804
MEGA	5	OSPC	4,5312	0,4840	92,8699	5,8633	0,3732	351,8834
MEGA	6	OSPC	4,6600	0,5010	105,6361	4,8701	0,3846	130,3339
MEGA	7	OSPC	4,5299	0,5141	92,7493	4,9586	0,3939	142,3943
MEGA	8	OSPC	4,5865	0,5260	98,1503	4,6286	0,4003	102,3706
MEGA	9	OSPC	4,6243	0,5361	101,9314	4,4297	0,4082	83,9062
MEGA	10	OSPC	4,6408	0,5452	103,6272	4,4606	0,4144	86,5394
MEGA	11	OSPC	4,6418	0,5529	103,7309	4,4702	0,4187	87,3742
MEGA	12	OSPC	4,6247	0,5607	101,9722	4,4862	0,4219	88,7834
MEGA	13	OSPC	4,6025	0,5673	99,7333	4,4549	0,4253	86,0475
MEGA	14	OSPC	4,5662	0,5736	96,1779	4,3981	0,4299	81,2963
MEGA	15	OSPC	4,5440	0,5795	94,0663	4,3046	0,4299	74,0396
MEGA	16	OSPC	4,3830	0,6127	80,0779	4,3135	0,4460	74,7015
MEGA	17	OSPC	4,3806	0,6218	79,8860	4,3375	0,4462	76,5160
MEGA	18	OSPC	4,3674	0,6258	78,8384	4,3397	0,4464	76,6845
MEGA	19	OSPC	4,2711	0,6424	71,6004	4,2542	0,4453	70,4005
MEGA	20	OSPC	4,3003	0,6358	73,7219	4,2313	0,4466	68,8066

Anexo No. 3. Uso de memoria y duración por épocas

Model	Epoch	Dataset	trn_wrkmem	trn_logmem	trn_elapsed_time	val_wrkmem	val_logmem	val_elapsed_time
Vainilla	1	UNPC	0,0329	140,7419	3252,4778	0,0324	140,246	493,1936
Vainilla	2	UNPC	0,0330	141,1686	3117,4205	0,0324	140,1794	463,9956
Vainilla	3	UNPC	0,0329	140,7549	3082,9082	0,0323	139,8333	442,6457
Vainilla	4	UNPC	0,0329	140,7613	3049,8695	0,0324	140,2536	448,505
Vainilla	5	UNPC	0,0328	140,6770	3218,9603	0,0323	139,8746	492,9527
Vainilla	6	UNPC	0,0328	140,5996	3115,2957	0,0323	139,8725	485,2851
Vainilla	7	UNPC	0,0329	140,7345	3145,7566	0,0323	139,9602	499,647
Vainilla	8	UNPC	0,0329	140,7631	3137,8544	0,0324	140,2027	474,9537
Vainilla	9	UNPC	0,0329	140,8601	3177,3929	0,0323	139,8978	499,4537
Vainilla	10	UNPC	0,0328	140,6457	3203,6590	0,0323	140,05	504,5685
Vainilla	11	UNPC	0,0329	140,9287	3213,9510	0,0322	139,7738	480,2847
Vainilla	12	UNPC	0,0329	140,7558	3166,0436	0,0324	140,1397	484,307
Vainilla	13	UNPC	0,0329	140,8738	3147,7708	0,0323	139,8664	485,2911
Vainilla	14	UNPC	0,0329	140,8809	3173,8746	0,0324	139,9884	458,987
Vainilla	15	UNPC	0,0329	140,7718	3088,6795	0,0323	140,0148	475,5751
Vainilla	16	UNPC	0,0329	140,9416	3083,7847	0,0324	140,2851	466,5793
Vainilla	17	UNPC	0,0329	140,7368	3090,4371	0,0324	140,2141	467,4164
Vainilla	18	UNPC	0,0329	140,8558	3091,9812	0,0322	139,6931	470,3111
Vainilla	19	UNPC	0,0329	140,8577	3177,4965	0,0323	139,8821	513,5908
Vainilla	20	UNPC	0,0329	140,7917	3242,8158	0,0324	140,0705	506,0144
Vainilla	1	OSPC	0,0041	23,6825	1755,3742	0,0041	23,8403	262,4378
Vainilla	2	OSPC	0,0041	23,6403	1854,4071	0,0041	23,7444	229,8436
Vainilla	3	OSPC	0,0041	23,6659	1727,2330	0,0041	23,8112	194,2142
Vainilla	4	OSPC	0,0041	23,6733	1726,9528	0,0041	23,8091	235,6324
Vainilla	5	OSPC	0,0041	23,6518	1766,0406	0,0041	23,7894	204,016
Vainilla	6	OSPC	0,0041	23,6724	1769,3048	0,0041	23,7526	205,483
Vainilla	7	OSPC	0,0041	23,6816	1769,5602	0,0041	23,8916	237,6558
Vainilla	8	OSPC	0,0041	23,6669	1767,8828	0,0041	23,8295	238,6547
Vainilla	9	OSPC	0,0041	23,6942	1768,8217	0,0041	23,8099	203,4891
Vainilla	10	OSPC	0,0041	23,6820	1775,3085	0,0041	23,8323	200,0598
Vainilla	11	OSPC	0,0041	23,6842	1771,5811	0,0041	23,8071	235,9879
Vainilla	12	OSPC	0,0041	23,6681	1773,9149	0,0041	23,7747	199,2812
Vainilla	13	OSPC	0,0041	23,6765	1776,6029	0,0041	23,8471	236,5926
Vainilla	14	OSPC	0,0041	23,6641	1773,5273	0,0041	23,7904	238,7636
Vainilla	15	OSPC	0,0041	23,6361	1778,4529	0,0041	23,7561	237,2725
Vainilla	16	OSPC	0,0041	23,7121	1780,5073	0,0041	23,7926	199,0567

Model	Epoch	Dataset	trn_wrkmem	trn_logmem	trn_elapsed_time	val_wrkmem	val_logmem	val_elapsed_time
Vainilla	17	OSPC	0,0041	23,6689	1782,6337	0,0041	23,8351	257,4153
Vainilla	18	OSPC	0,0041	23,6491	1784,6467	0,0041	23,8199	269,1961
Vainilla	19	OSPC	0,0041	23,6469	1762,9273	0,0041	23,8222	217,6827
Vainilla	20	OSPC	0,0041	23,6514	1816,5513	0,0041	23,8240	262,9357
UT	1	UNPC	0,0080	140,6679	4792,7157	0,0080	140,0182	668,9942
UT	2	UNPC	0,0080	141,0117	4950,8178	0,0080	139,7525	653,2644
UT	3	UNPC	0,0080	141,0612	4813,2657	0,0080	139,9481	645,8403
UT	4	UNPC	0,0080	140,8780	4868,2482	0,0080	139,582	675,8555
UT	5	UNPC	0,0080	140,9035	4908,4884	0,0080	139,9681	678,9933
UT	6	UNPC	0,0080	140,4487	4782,3002	0,0080	140,1957	649,5450
UT	7	UNPC	0,0080	140,9902	4876,9516	0,0080	140,0771	661,1203
UT	8	UNPC	0,0080	140,8772	4762,5583	0,0080	139,9679	638,5373
UT	9	UNPC	0,0080	140,7929	4935,8223	0,0080	139,9339	669,2653
UT	10	UNPC	0,0080	140,9053	4923,7648	0,0080	140,0572	675,7470
UT	11	UNPC	0,0080	140,8223	5044,4528	0,0080	140,1144	714,8989
UT	12	UNPC	0,0080	140,6299	4934,5474	0,0080	139,9173	685,4597
UT	13	UNPC	0,0080	140,9016	4923,3137	0,0080	139,6739	691,9752
UT	14	UNPC	0,0080	140,7971	4918,8142	0,0080	139,7298	688,1158
UT	15	UNPC	0,0080	140,7470	4930,6052	0,0080	140,1737	680,8276
UT	16	UNPC	0,0080	140,8711	4964,4458	0,0080	140,0833	685,9692
UT	17	UNPC	0,0080	140,741	4940,6114	0,0080	140,1554	697,9747
UT	18	UNPC	0,0080	140,839	4987,5865	0,0080	140,2275	688,2655
UT	19	UNPC	0,0080	140,7728	5012,3937	0,0080	139,9473	687,4299
UT	20	UNPC	0,0080	140,872	4933,1812	0,0080	140,0280	675,8899
UT	1	OSPC	0,0023	23,6722	1958,3248	0,0023	23,8810	180,0426
UT	2	OSPC	0,0023	23,6677	1873,7227	0,0023	23,8579	181,6556
UT	3	OSPC	0,0023	23,6749	1894,0368	0,0023	23,7844	186,8751
UT	4	OSPC	0,0023	23,646	1957,1406	0,0023	23,7476	186,0802
UT	5	OSPC	0,0023	23,6465	1962,0664	0,0023	23,8391	182,6534
UT	6	OSPC	0,0023	23,6668	2023,2704	0,0023	23,7488	193,6856
UT	7	OSPC	0,0023	23,6677	1906,2312	0,0023	23,8000	180,3189
UT	8	OSPC	0,0023	23,6684	1894,7842	0,0023	23,7892	179,4730
UT	9	OSPC	0,0023	23,6337	1876,4877	0,0023	23,7751	178,1618
UT	10	OSPC	0,0023	23,681	1897,1231	0,0023	23,7635	182,2534
UT	11	OSPC	0,0023	23,6436	1902,4572	0,0023	23,8107	181,7154
UT	12	OSPC	0,0023	23,6433	1902,4372	0,0023	23,8078	186,4531
								·
UT	13	OSPC	0,0023	23,6941	1938,0505	0,0023	23,8270	172,6852

Model	Epoch	Dataset	trn_wrkmem	trn_logmem	trn_elapsed_time	val_wrkmem	val_logmem	val_elapsed_time
UT	14	OSPC	0,0023	23,6837	1889,4263	0,0023	23,8250	176,4444
UT	15	OSPC	0,0023	23,6379	1901,0695	0,0023	23,8370	174,9981
UT	16	OSPC	0,0023	23,6666	1898,4783	0,0023	23,8321	176,3040
UT	17	OSPC	0,0023	23,6515	1900,7431	0,0023	23,7861	178,7733
UT	18	OSPC	0,0023	23,6579	1899,6260	0,0023	23,8180	176,0774
UT	19	OSPC	0,0023	23,6930	1898,5579	0,0023	23,8229	181,0199
UT	20	OSPC	0,0023	23,6776	1899,8391	0,0023	23,7236	196,0320
MAMBA	1	UNPC	0,0080	140,7552	3402,5075	0,0080	139,6836	592,0435
MAMBA	2	UNPC	0,0080	140,6398	3425,9119	0,0080	139,7029	593,5259
MAMBA	3	UNPC	0,0080	140,7445	3480,9590	0,0080	139,7543	598,7798
MAMBA	4	UNPC	0,0080	140,8927	3481,9169	0,0080	140,1140	623,4210
MAMBA	5	UNPC	0,0080	140,8443	3504,5645	0,0080	139,6693	624,5111
MAMBA	6	UNPC	0,0080	140,7380	3488,9977	0,0080	139,5470	604,8754
MAMBA	7	UNPC	0,0080	140,9883	3507,7827	0,0080	140,0891	624,4096
MAMBA	8	UNPC	0,0080	140,8291	3562,3346	0,0080	140,5082	618,8800
MAMBA	9	UNPC	0,0080	141,0163	3563,2849	0,0080	139,8154	635,9067
MAMBA	10	UNPC	0,0080	140,7316	3667,3007	0,0080	139,9709	615,2520
MAMBA	11	UNPC	0,0080	140,7672	3609,8970	0,0080	139,9866	641,6283
MAMBA	12	UNPC	0,0080	140,7667	3605,0222	0,0080	140,4130	622,7757
MAMBA	13	UNPC	0,0080	140,9567	3473,4220	0,0080	139,7468	580,4235
MAMBA	14	UNPC	0,0080	141,1065	3616,7546	0,0080	139,5190	621,9412
MAMBA	15	UNPC	0,0080	140,7970	3467,9522	0,0080	139,9688	604,3434
MAMBA	16	UNPC	0,0080	140,6550	3472,9991	0,0080	139,8603	609,6284
MAMBA	17	UNPC	0,0080	140,6748	3465,7592	0,0080	140,2187	608,7233
MAMBA	18	UNPC	0,0080	140,9653	3475,1991	0,0080	139,8766	606,0925
MAMBA	19	UNPC	0,0080	140,7787	3477,5204	0,0080	140,1719	612,3055
MAMBA	20	UNPC	0,0080	140,5831	3466,9826	0,0080	140,0105	616,1288
MAMBA	1	OSPC	0,0023	23,7093	1129,8863	0,0023	23,7718	196,7603
MAMBA	2	OSPC	0,0023	23,6678	1119,7239	0,0023	23,7876	188,0392
MAMBA	3	OSPC	0,0023	23,7008	1120,3657	0,0023	23,8228	176,4743
MAMBA	4	OSPC	0,0023	23,6795	1112,9473	0,0023	23,7758	193,6998
MAMBA	5	OSPC	0,0023	23,6548	1134,1315	0,0023	23,7975	211,8254
MAMBA	6	OSPC	0,0023	23,6599	1132,8885	0,0023	23,7681	189,2310
MAMBA	7	OSPC	0,0023	23,6604	1124,5692	0,0023	23,7982	203,8140
MAMBA	8	OSPC	0,0023	23,6508	1099,0290	0,0023	23,9090	207,4640
MAMBA	9	OSPC	0,0023	23,6629	1113,6997	0,0023	23,7117	203,8463
MAMBA	10	OSPC	0,0023	23,6773	1169,3137	0,0023	23,8254	202,3518

poch Datase	trn_wrkmem	trn_logmem	trn_elapsed_time	val_wrkmem	val_logmem	val_elapsed_time
11 OSPC	0,0023	23,6619	1170,8986	0,0023	23,8424	207,5673
12 OSPC	0,0023	23,6334	1154,4460	0,0023	23,7945	214,7940
13 OSPC	0,0023	23,675	1159,2274	0,0023	23,7099	195,8480
14 OSPC	0,0023	23,6849	1176,3137	0,0023	23,8306	211,9772
15 OSPC	0,0023	23,6751	1185,0682	0,0023	23,7895	218,2352
16 OSPC	0,0023	23,6285	1146,6796	0,0023	23,8787	204,0906
17 OSPC	0,0023	23,6481	1148,4755	0,0023	23,8069	202,5712
18 OSPC	0,0023	23,6248	1170,2477	0,0023	23,8168	196,6169
19 OSPC	0,0023	23,686	1157,0999	0,0023	23,8900	215,8083
20 OSPC	0,0023	23,6556	1151,4169	0,0023	23,8615	191,6002
1 UNPC	0,0080	140,9161	3871,5407	0,0080	139,7091	582,5655
2 UNPC	0,0080	140,8537	3924,7202	0,0080	139,8252	635,8916
3 UNPC	0,0080	140,6651	4093,5019	0,0080	139,6560	642,7858
4 UNPC	0,0080	140,7442	4025,8507	0,0080	139,4975	638,9806
5 UNPC	0,0080	140,909	4039,0190	0,0080	139,6496	616,9055
6 UNPC	0,0080	140,7968	4043,5296	0,0079	139,3876	625,6697
7 UNPC	0,0080	140,8162	4042,4262	0,0080	140,0033	650,0579
8 UNPC	0,0080	140,8088	4057,2061	0,0080	139,8049	629,7170
9 UNPC	0,0080	140,8306	4080,5789	0,0080	140,2405	629,5483
10 UNPC	0,0080	140,9008	4019,0397	0,0080	139,7059	632,2203
11 UNPC	0,0080	140,9185	4091,0646	0,0080	140,0244	664,2030
12 UNPC	0,0080	140,7125	4043,7053	0,0080	140,1780	629,0457
13 UNPC	0,0080	140,7748	4108,6284	0,0080	139,4647	678,2561
14 UNPC	0,0080	140,8789	4148,1885	0,0080	140,1306	626,0521
15 UNPC	0,0080	140,8297	4067,9162	0,0080	139,5797	644,3566
16 UNPC	0,0080	140,7661	4087,3120	0,0080	140,0321	658,6393
17 UNPC	0,0080	140,9458	4110,7134	0,0080	139,9295	656,0169
18 UNPC	0,0080	141,0246	4187,7633	0,0080	140,4063	675,6108
19 UNPC	0,0080	140,8173	4280,4183	0,0080	139,6320	650,2625
20 UNPC	0,0080	140,7736	4123,2744	0,0080	139,7705	667,9472
1 OSPC	0,0023	23,6634	1756,0655	0,0023	23,8467	215,9467
2 OSPC	0,0023	23,6639	1750,8099	0,0023	23,8775	221,2890
3 OSPC	0,0023	23,7014	1748,0015	0,0023	23,8104	254,3890
4 OSPC	0,0023	23,675	1745,9049	0,0023	23,8854	251,9056
5 OSPC	0,0023	23,6371	1744,8042	0,0023	23,9145	252,2987
6 OSPC	0,0023	23,6247	1750,0712	0,0023	23,7455	222,4385
						248,9967
6 OSF	C	PC 0,0023	PC 0,0023 23,6247	PC 0,0023 23,6247 1750,0712	PC 0,0023 23,6247 1750,0712 0,0023	PC 0,0023 23,6247 1750,0712 0,0023 23,7455

Model	Epoch	Dataset	trn_wrkmem	trn_logmem	trn_elapsed_time	val_wrkmem	val_logmem	val_elapsed_time
MoE	8	OSPC	0,0023	23,6651	1746,0383	0,0023	23,7557	248,8262
MoE	9	OSPC	0,0023	23,6770	1750,0066	0,0023	23,7387	253,4918
MoE	10	OSPC	0,0023	23,6634	1747,4669	0,0023	23,8294	255,9006
MoE	11	OSPC	0,0023	23,6909	1747,9610	0,0023	23,7531	256,5059
MoE	12	OSPC	0,0023	23,6588	1753,7316	0,0023	23,8755	246,4398
MoE	13	OSPC	0,0023	23,6634	1748,3911	0,0023	23,8559	227,9929
MoE	14	OSPC	0,0023	23,6558	1746,644	0,0023	23,7431	251,2693
MoE	15	OSPC	0,0023	23,6920	1787,2587	0,0023	23,8015	216,5355
MoE	16	OSPC	0,0023	23,6736	1794,8017	0,0023	23,7702	277,7432
MoE	17	OSPC	0,0023	23,6689	1800,8721	0,0023	23,7882	249,3848
MoE	18	OSPC	0,0023	23,6390	1765,4587	0,0023	23,8749	254,7882
MoE	19	OSPC	0,0023	23,6606	1764,4466	0,0023	23,7822	266,2414
MoE	20	OSPC	0,0023	23,6798	1766,7868	0,0023	23,7523	262,6894
GSS	1	UNPC	0,0080	140,7757	4205,9107	0,0080	139,7658	650,0726
GSS	2	UNPC	0,0080	140,8024	4077,7127	0,0080	140,2159	668,0589
GSS	3	UNPC	0,0080	140,7382	4089,3136	0,0080	140,0340	671,8298
GSS	4	UNPC	0,0080	140,8901	4131,7285	0,0080	139,9025	666,7918
GSS	5	UNPC	0,0080	140,8660	4094,7222	0,0080	139,8397	667,4561
GSS	6	UNPC	0,0080	140,7446	4091,7376	0,0080	139,8032	668,2887
GSS	7	UNPC	0,0080	141,1095	4092,8928	0,0080	139,9964	673,9598
GSS	8	UNPC	0,0080	140,6898	4092,4580	0,0080	139,9175	672,1534
GSS	9	UNPC	0,0080	140,7278	4100,8393	0,0080	140,1036	671,3768
GSS	10	UNPC	0,0080	140,7984	4095,1755	0,0080	139,9657	674,3078
GSS	11	UNPC	0,0080	140,9072	4106,5804	0,0080	140,0069	679,8970
GSS	12	UNPC	0,0080	140,8583	4110,8664	0,0080	140,2514	679,0814
GSS	13	UNPC	0,0080	140,8054	4071,0252	0,0080	139,9623	645,5564
GSS	14	UNPC	0,0080	140,6092	3980,4288	0,0080	140,1469	636,8885
GSS	15	UNPC	0,0080	140,6443	3989,7846	0,0080	140,0736	635,9740
GSS	16	UNPC	0,0080	140,7855	4106,4023	0,0080	140,4622	697,4448
GSS	17	UNPC	0,0080	140,8433	4191,1865	0,0080	139,9269	689,2295
GSS	18	UNPC	0,0080	140,6676	4131,8163	0,0080	140,3589	687,9272
GSS	19	UNPC	0,0080	140,8975	4304,910	0,0080	139,8463	691,8394
GSS	20	UNPC	0,0080	140,8508	4165,1359	0,0080	140,0412	690,5499
GSS	1	OSPC	0,0023	23,6469	1735,9765	0,0023	23,8289	252,8748
GSS	2	OSPC	0,0023	23,6512	1707,1692	0,0023	23,7664	247,5602
GSS	3	OSPC	0,0023	23,6356	1702,8525	0,0023	23,7761	245,6084
GSS	4	OSPC	0,0023	23,6725	1747,1170	0,0023	23,8318	246,1396
333	7	551 €	0,0023	25,5125	17-17,1170	0,0023	25,0510	210,1330

Model	Epoch	Dataset	trn_wrkmem	trn_logmem	trn_elapsed_time	val_wrkmem	val_logmem	val_elapsed_time
GSS	5	OSPC	0,0023	23,6213	1666,9382	0,0023	23,7999	243,1244
GSS	6	OSPC	0,0023	23,6965	1677,7869	0,0023	23,8456	241,5298
GSS	7	OSPC	0,0023	23,6718	1678,8416	0,0023	23,7835	218,7235
GSS	8	OSPC	0,0023	23,6648	1723,1286	0,0023	23,7833	235,7657
GSS	9	OSPC	0,0023	23,6869	1725,9364	0,0023	23,7545	245,6530
GSS	10	OSPC	0,0023	23,6637	1652,2613	0,0023	23,8423	233,4587
GSS	11	OSPC	0,0023	23,6532	1662,3356	0,0023	23,7444	221,2425
GSS	12	OSPC	0,0023	23,6963	1659,5856	0,0023	23,8703	227,1802
GSS	13	OSPC	0,0023	23,6116	1660,5216	0,0023	23,8303	229,2345
GSS	14	OSPC	0,0023	23,6723	1667,2141	0,0023	23,8381	234,8912
GSS	15	OSPC	0,0023	23,6431	1660,0374	0,0023	23,8411	230,6358
GSS	16	OSPC	0,0023	23,6741	1660,5043	0,0023	23,8861	221,7403
GSS	17	OSPC	0,0023	23,6636	1663,3182	0,0023	23,7951	233,6651
GSS	18	OSPC	0,0023	23,7027	1660,9304	0,0023	23,8034	231,4069
GSS	19	OSPC	0,0023	23,6410	1661,1153	0,0023	23,9000	232,3270
GSS	20	OSPC	0,0023	23,6640	1659,9122	0,0023	23,7928	235,8479
MEGA	1	UNPC	0,0080	140,7601	5283,4472	0,0080	140,1408	757,7263
MEGA	2	UNPC	0,0080	140,6825	5325,3194	0,0080	140,0021	767,3134
MEGA	3	UNPC	0,0080	140,8897	5352,5379	0,0080	139,7841	796,8020
MEGA	4	UNPC	0,0080	140,6302	5298,3803	0,0080	139,9202	802,2413
MEGA	5	UNPC	0,0080	140,7964	5290,6898	0,0080	140,4120	773,2582
MEGA	6	UNPC	0,0080	140,8162	5377,9656	0,0080	139,8445	774,0481
MEGA	7	UNPC	0,0080	140,8785	5382,9373	0,0080	139,8231	807,9502
MEGA	8	UNPC	0,0080	141,0372	5309,2461	0,0080	139,6785	756,9659
MEGA	9	UNPC	0,008	140,8011	5312,0664	0,0080	139,8871	774,8150
MEGA	10	UNPC	0,0080	140,7426	5190,7091	0,0080	139,9650	755,6971
MEGA	11	UNPC	0,0080	140,8688	5092,7188	0,0080	139,8451	761,3720
MEGA	12	UNPC	0,0080	140,9192	5097,0430	0,0080	140,2339	756,3194
MEGA	13	UNPC	0,0080	140,8428	5143,5562	0,0080	139,8717	762,2020
MEGA	14	UNPC	0,0080	140,7546	5104,4121	0,0080	139,8923	743,9114
MEGA	15	UNPC	0,0080	140,8488	5197,8677	0,0080	139,6897	756,7967
MEGA	16	UNPC	0,0080	140,8209	5219,1782	0,0080	139,9234	763,0206
MEGA	17	UNPC	0,008	141,0391	5211,5501	0,0080	139,8438	752,0593
MEGA	18	UNPC	0,0080	140,9276	5180,6254	0,0080	139,8366	809,4031
MEGA	19	UNPC	0,0080	140,8342	5384,2676	0,0080	140,3626	787,3313
MEGA	20	UNPC	0,0080	140,7070	5173,9321	0,0080	140,1628	751,5044
MEGA	1	OSPC	0,0023	23,6964	2081,2353	0,0023	23,74120	258,1549

Model	Epoch	Dataset	trn_wrkmem	trn_logmem	trn_elapsed_time	val_wrkmem	val_logmem	val_elapsed_time
MEGA	2	OSPC	0,0023	23,6908	2024,7627	0,0023	23,7745	239,8274
MEGA	3	OSPC	0,0023	23,6729	1965,5972	0,0023	23,7992	237,9977
MEGA	4	OSPC	0,0023	23,6594	2003,3546	0,0023	23,9229	250,9025
MEGA	5	OSPC	0,0023	23,6776	1986,9552	0,0023	23,8378	307,8811
MEGA	6	OSPC	0,0023	23,6429	2079,7147	0,0023	23,7759	280,4368
MEGA	7	OSPC	0,0023	23,6480	2036,2628	0,0023	23,8193	301,5237
MEGA	8	OSPC	0,0023	23,6325	2097,4549	0,0023	23,8607	251,1039
MEGA	9	OSPC	0,0023	23,6657	1989,4193	0,0023	23,8146	245,1587
MEGA	10	OSPC	0,0023	23,6447	1981,1929	0,0023	23,8315	257,2948
MEGA	11	OSPC	0,0023	23,6549	1981,8079	0,0023	23,8455	261,5035
MEGA	12	OSPC	0,0023	23,6796	1979,0316	0,0023	23,8456	275,8985
MEGA	13	OSPC	0,0023	23,6816	1999,8379	0,0023	23,7679	278,9156
MEGA	14	OSPC	0,0023	23,6772	1978,1469	0,0023	23,8148	278,7817
MEGA	15	OSPC	0,0023	23,6621	1869,8282	0,0023	23,8531	235,5913
MEGA	16	OSPC	0,0023	23,7051	1838,3003	0,0023	23,8130	254,1711
MEGA	17	OSPC	0,0023	23,6594	1847,8760	0,0023	23,8361	249,9540
MEGA	18	OSPC	0,0023	23,6554	1855,9377	0,0023	23,8428	238,8200
MEGA	19	OSPC	0,0023	23,6621	1884,6204	0,0023	23,8531	241,4235
MEGA	20	OSPC	0,0023	23,7051	1881,6609	0,0023	23,8130	242,2295