

# Análisis Rigurosamente Matemático del Modelo LiquidFoundationModelOptimized

Asistente de IA

October 31, 2024

## Contents

<b>1</b>	<b>Introducción</b>	<b>3</b>
<b>2</b>	<b>Visión General del Modelo</b>	<b>3</b>
2.1	Flujo de Datos General . . . . .	3
<b>3</b>	<b>Análisis Rigurosamente Matemático de los Componentes Principales</b>	<b>4</b>
3.1	FrequencyHyenaAttention . . . . .	4
3.1.1	Descripción Algorítmica Completa . . . . .	4
3.1.2	Fundamentos Matemáticos Avanzados . . . . .	5
3.1.3	Análisis Matemático Detallado . . . . .	5
3.1.4	Análisis Matemático Rigurosamente Detallado . . . . .	7
3.1.5	Propiedades Matemáticas y Ventajas . . . . .	8
3.1.6	Implicaciones Matemáticas para el Modelo . . . . .	8
3.1.7	Integración de FrequencyHyenaAttention en EnhancedLiquidEmbedding . . . . .	8
3.1.8	Impacto de FrequencyHyenaAttention en el Modelo LiquidFoundationModelOptimized . . . . .	9
3.2	EnhancedLiquidEmbedding . . . . .	11
3.2.1	Descripción Algorítmica . . . . .	11
3.2.2	Fundamentos Matemáticos . . . . .	12
3.2.3	Análisis Detallado . . . . .	13
3.2.4	Implicaciones para el Modelo . . . . .	13
3.3	EnhancedLocalAttentionWithGQA . . . . .	14
3.3.1	Descripción Algorítmica . . . . .	14
3.3.2	Fundamentos Matemáticos . . . . .	14
3.3.3	Análisis Detallado . . . . .	15
3.3.4	Implicaciones para el Modelo . . . . .	15
3.4	MoELayer (Mixture of Experts) . . . . .	16
3.4.1	Descripción Algorítmica . . . . .	16
3.4.2	Fundamentos Matemáticos . . . . .	16
3.4.3	Análisis Detallado . . . . .	17
3.4.4	Implicaciones para el Modelo . . . . .	17
3.5	DeformableConv1d . . . . .	18
3.5.1	Descripción Algorítmica . . . . .	18
3.5.2	Fundamentos Matemáticos . . . . .	18
3.5.3	Análisis Detallado . . . . .	18
3.5.4	Implicaciones para el Modelo . . . . .	19
3.6	OptimizedGatedConvolution . . . . .	20
3.6.1	Descripción Algorítmica . . . . .	20
3.6.2	Fundamentos Matemáticos . . . . .	20
3.6.3	Análisis Detallado . . . . .	20
3.6.4	Implicaciones para el Modelo . . . . .	20
3.7	EnhancedLSTM . . . . .	22
3.7.1	Descripción Algorítmica . . . . .	22

3.7.2	Fundamentos Matemáticos . . . . .	22
3.7.3	Análisis Detallado . . . . .	22
3.7.4	Implicaciones para el Modelo . . . . .	23
3.8	ImprovedTransformerBlock . . . . .	24
3.8.1	Descripción Algorítmica . . . . .	24
3.8.2	Fundamentos Matemáticos . . . . .	25
3.8.3	Análisis Detallado . . . . .	25
3.8.4	Implicaciones para el Modelo . . . . .	25
3.9	BidirectionalEncoder . . . . .	26
3.9.1	Descripción Algorítmica . . . . .	26
3.9.2	Fundamentos Matemáticos . . . . .	26
3.9.3	Análisis Detallado . . . . .	26
3.9.4	Implicaciones para el Modelo . . . . .	26
3.10	LiquidFoundationModelOptimized . . . . .	28
3.10.1	Descripción Algorítmica . . . . .	28
3.10.2	Fundamentos Matemáticos . . . . .	28
3.10.3	Análisis Detallado . . . . .	29
3.10.4	Implicaciones para el Modelo . . . . .	29
<b>4</b>	<b>Consideraciones de Estabilidad Numérica y Eficiencia</b>	<b>30</b>
4.0.1	Análisis Detallado . . . . .	31
4.1	Análisis de Eficiencia y Estabilidad . . . . .	31
<b>5</b>	<b>Conclusión</b>	<b>31</b>
5.1	Puntos Clave del Modelo . . . . .	31
5.2	Resumen de Implicaciones de las Mejoras . . . . .	32

# 1 Introducción

El presente documento realiza un análisis matemático exhaustivo del modelo de aprendizaje profundo denominado **LiquidFoundationModelOptimized**, diseñado para tareas avanzadas de procesamiento de lenguaje natural (NLP). Este modelo integra múltiples arquitecturas y técnicas de vanguardia, incluyendo Transformers mejorados, Mixture-of-Experts (MoE), convoluciones deformables, mecanismos de atención local optimizados y la novedosa **FrequencyHyenaAttention**. Además, incorpora mecanismos avanzados de regularización y estabilización numérica, garantizando eficiencia y robustez durante el entrenamiento y la inferencia.

A lo largo del documento, se profundizará en la estructura del modelo, desglosando cada uno de sus componentes y explicando detalladamente los fundamentos matemáticos que los sustentan. Asimismo, se discutirán las implicaciones de las mejoras implementadas en términos de capacidad de modelado, eficiencia computacional y estabilidad durante el entrenamiento.

## 2 Visión General del Modelo

El **LiquidFoundationModelOptimized** se organiza en dos secciones principales: un **encoder bidireccional** y un **decoder unidireccional**. Ambas secciones emplean embeddings líquidos personalizados y bloques Transformer mejorados que incorporan módulos avanzados como atención local optimizada, MoE, convoluciones deformables y **FrequencyHyenaAttention**. Además, el modelo integra una memoria externa basada en una versión mejorada de LSTM para gestionar información secuencial de manera eficiente.

### 2.1 Flujo de Datos General

#### 1. Entrada del Encoder:

- Los identificadores de tokens de entrada (**encoder\_input\_ids**) se procesan a través del módulo de embeddings líquidos.
- La representación embebida resultante se pasa por múltiples capas de bloques Transformer mejorados en el encoder.

#### 2. Entrada del Decoder:

- Los identificadores de tokens de entrada del decoder (**decoder\_input\_ids**) también se embeben utilizando embeddings líquidos.
- La representación embebida se procesa mediante múltiples capas de bloques Transformer mejorados en el decoder, que incorporan una memoria externa LSTM para manejar dependencias a largo plazo.

#### 3. Generación de Salida:

- La salida final del decoder se proyecta a través de una capa lineal para generar logits que representan las probabilidades de los tokens de salida.

### 3 Análisis Rigurosamente Matemático de los Componentes Principales

#### 3.1 FrequencyHyenaAttention

**Propósito:** Implementar un mecanismo de atención basado en la arquitectura Hyena, optimizado para operar en el dominio de la frecuencia. Este mecanismo busca mejorar la eficiencia computacional y la capacidad de capturar relaciones a largo plazo en secuencias de entrada, aprovechando las propiedades espectrales de las señales.

##### 3.1.1 Descripción Algorítmica Completa

###### 1. Transformada de Fourier Rápida (FFT) en Entradas:

$$\mathbf{X}_{\text{fft}} = \mathcal{F}(\mathbf{X}, \text{dim} = 1)$$

donde  $\mathcal{F}$  representa la Transformada de Fourier Rápida (FFT) aplicada a la dimensión temporal (dim=1) de la entrada  $\mathbf{X} \in \mathbb{R}^{B \times L \times C}$ , con  $B$  el tamaño del batch,  $L$  la longitud de la secuencia y  $C$  la dimensión de embedding.

###### 2. Descomposición en Magnitud y Fase:

$$\mathbf{M} = |\mathbf{X}_{\text{fft}}| \quad \text{y} \quad \Phi = \angle \mathbf{X}_{\text{fft}}$$

donde  $\mathbf{M} \in \mathbb{R}^{B \times L \times C}$  es la magnitud y  $\Phi \in \mathbb{R}^{B \times L \times C}$  es la fase de la representación frecuencial.

###### 3. Proyección Lineal y Normalización por Banda de Frecuencia: Las bandas de frecuencia se definen como potencias de dos, permitiendo una descomposición jerárquica eficiente.

$$\mathbf{H}_b = \text{LayerNorm}(\text{Linear}_b(\mathbf{M}_b)) \quad \forall b \in \{1, 2, \dots, B\}$$

donde  $\mathbf{M}_b$  corresponde a la banda de frecuencia  $b$ , y  $\text{Linear}_b \in \mathbb{R}^{d \times d}$  es una capa lineal específica para cada banda  $b$ .

###### 4. Atención Multi-Cabeza por Banda: Para cada banda de frecuencia $b$ , se aplica una atención multi-cabeza:

$$\mathbf{A}_b = \text{MultiHeadAttention}(\mathbf{H}_b, \mathbf{H}_b, \mathbf{H}_b)$$

donde  $\mathbf{A}_b \in \mathbb{R}^{B \times L \times C}$  es la salida de atención para la banda  $b$ .

###### 5. Modulación con Fase: La salida de la atención se modula con la fase correspondiente para preservar la información espectral:

$$\mathbf{O}_b = \mathbf{A}_b \odot e^{i\Phi_b}$$

donde  $\odot$  denota la multiplicación elemento a elemento y  $e^{i\Phi_b}$  es la representación compleja de la fase.

###### 6. Combinación de Bandas y Proyección Final: Las salidas de todas las bandas se combinan mediante una media ponderada y se proyectan de vuelta al espacio de embedding:

$$\mathbf{O} = \frac{1}{B} \sum_{b=1}^B \mathbf{O}_b$$

$$\mathbf{Y} = \text{LayerNorm}(\text{Linear}(\Re(\mathbf{O})))$$

donde  $\Re(\cdot)$  extrae la parte real de la representación compleja y  $\text{Linear} \in \mathbb{R}^{d \times d}$  es una capa lineal para la proyección final.

### 3.1.2 Fundamentos Matemáticos Avanzados

**1. Transformadas de Fourier y Propiedades Espectrales** La Transformada de Fourier (FT) convierte una señal  $\mathbf{X}$  del dominio temporal al dominio frecuencial:

$$\mathbf{X}_{\text{fft}}[b, k, c] = \sum_{l=0}^{L-1} \mathbf{X}[b, l, c] e^{-i2\pi kl/L}$$

donde  $k$  representa la frecuencia. La FFT es una implementación eficiente de la FT con una complejidad computacional de  $\mathcal{O}(L \log L)$ .

Las propiedades espectrales, como la linealidad y la descomposición en magnitud y fase, permiten manipular y analizar la señal en componentes de frecuencia, facilitando la identificación de patrones repetitivos y relaciones a largo plazo.

**2. Normalización y Proyección por Banda** La normalización de capa (**LayerNorm**) asegura que cada banda de frecuencia tenga una media cero y una varianza unitaria, lo que estabiliza el entrenamiento y mejora la convergencia:

$$\text{LayerNorm}(\mathbf{H}_b) = \frac{\mathbf{H}_b - \mu_b}{\sqrt{\sigma_b^2 + \epsilon}} \gamma_b + \beta_b$$

donde  $\mu_b$  y  $\sigma_b^2$  son la media y varianza de  $\mathbf{H}_b$ , respectivamente, y  $\gamma_b, \beta_b$  son parámetros aprendibles.

La proyección lineal (**Linear** <sub>$b$</sub> ) *transforma la representación frecuencial para cada banda, permitiendo que cada una aporte*

**3. Atención Multi-Cabeza en el Dominio Frecuencial** La atención multi-cabeza se define como:

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax} \left( \frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d_k}} \right) \mathbf{V}$$

donde  $\mathbf{Q}, \mathbf{K}, \mathbf{V}$  son las matrices de consultas, claves y valores, respectivamente, y  $d_k$  es la dimensión de las consultas.

En el dominio frecuencial, esta operación permite que cada cabeza de atención enfoque diferentes componentes espectrales, capturando relaciones complejas y de largo alcance sin incurrir en la complejidad cuadrática típica de la atención estándar.

**4. Modulación con Fase** La modulación con la fase preserva la información de fase de la señal original, que es crucial para mantener la coherencia temporal y la estructura de la señal:

$$\mathbf{O}_b = \mathbf{A}_b \odot e^{i\Phi_b}$$

Esta operación garantiza que la salida de la atención conserve las propiedades espectrales originales, permitiendo una reconstrucción precisa durante la Transformada Inversa de Fourier (IFFT).

**5. Reconstrucción y Proyección Final** La combinación de las salidas de todas las bandas mediante una media ponderada:

$$\mathbf{O} = \frac{1}{B} \sum_{b=1}^B \mathbf{O}_b$$

reduce la dimensionalidad de la representación frecuencial combinada y garantiza que la información de todas las bandas de frecuencia contribuya equitativamente a la salida final.

La proyección lineal y la normalización de capa transforman la representación combinada de vuelta al espacio de embedding original, asegurando que las características aprendidas en el dominio frecuencial sean integradas de manera coherente con el resto de la arquitectura del modelo:

$$\mathbf{Y} = \text{LayerNorm}(\text{Linear}(\Re(\mathbf{O})))$$

### 3.1.3 Análisis Matemático Detallado

Para comprender completamente el funcionamiento de **FrequencyHyenaAttention**, consideremos cada paso con mayor detalle:

**1. Transformada de Fourier y Representación Frecuencial** La Transformada de Fourier de una señal discreta  $\mathbf{X} \in \mathbb{R}^{B \times L \times C}$  se define como:

$$\mathbf{X}_{\text{fft}}[b, k, c] = \sum_{l=0}^{L-1} \mathbf{X}[b, l, c] e^{-i2\pi kl/L}, \quad k = 0, 1, \dots, L-1$$

donde  $i$  es la unidad imaginaria. La FFT convierte la señal temporal  $\mathbf{X}$  en una representación frecuencial  $\mathbf{X}_{\text{fft}}$ , permitiendo el análisis y procesamiento en el dominio de la frecuencia.

**2. Descomposición en Magnitud y Fase** La señal frecuencial  $\mathbf{X}_{\text{fft}}$  se descompone en magnitud  $\mathbf{M}$  y fase  $\Phi$ :

$$\begin{aligned} \mathbf{M}[b, k, c] &= |\mathbf{X}_{\text{fft}}[b, k, c]| = \sqrt{\Re(\mathbf{X}_{\text{fft}}[b, k, c])^2 + \Im(\mathbf{X}_{\text{fft}}[b, k, c])^2} \\ \Phi[b, k, c] &= \angle \mathbf{X}_{\text{fft}}[b, k, c] = \arctan \left( \frac{\Im(\mathbf{X}_{\text{fft}}[b, k, c])}{\Re(\mathbf{X}_{\text{fft}}[b, k, c])} \right) \end{aligned}$$

donde  $\Re(\cdot)$  e  $\Im(\cdot)$  representan las partes real e imaginaria, respectivamente.

**3. Proyección y Normalización por Banda** Cada banda de frecuencia  $b$  se procesa individualmente:

$$\mathbf{H}_b = \text{LayerNorm}(\mathbf{M}_b \mathbf{W}_b + \mathbf{b}_b)$$

donde  $\mathbf{W}_b \in \mathbb{R}^{d \times d}$  y  $\mathbf{b}_b \in \mathbb{R}^d$  son los parámetros de la capa lineal para la banda  $b$ .

La Layer Normalization se define como:

$$\text{LayerNorm}(\mathbf{H}_b) = \frac{\mathbf{H}_b - \mu_b}{\sqrt{\sigma_b^2 + \epsilon}} \gamma_b + \beta_b$$

con  $\mu_b$  y  $\sigma_b^2$  la media y varianza de  $\mathbf{H}_b$  a lo largo de la dimensión de características, y  $\gamma_b, \beta_b$  parámetros aprendibles.

**4. Atención Multi-Cabeza en Cada Banda** Para cada banda  $b$ , se aplica atención multi-cabeza:

$$\mathbf{A}_b = \text{softmax} \left( \frac{\mathbf{Q}_b \mathbf{K}_b^\top}{\sqrt{d_k}} \right) \mathbf{V}_b$$

donde:

$$\mathbf{Q}_b = \mathbf{H}_b \mathbf{W}_Q, \quad \mathbf{K}_b = \mathbf{H}_b \mathbf{W}_K, \quad \mathbf{V}_b = \mathbf{H}_b \mathbf{W}_V$$

y  $\mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V \in \mathbb{R}^{d \times d_k}$  son matrices de proyección aprendibles.

La atención multi-cabeza permite que diferentes cabezas se enfoquen en distintos subespacios de características, capturando diversas relaciones dentro de cada banda de frecuencia.

**5. Modulaci3n con Fase** La salida de atenci3n  $\mathbf{A}_b$  se modula con la fase  $\Phi_b$ :

$$\mathbf{O}_b = \mathbf{A}_b \odot e^{i\Phi_b}$$

Esta modulaci3n asegura que la estructura temporal original de la se1al se preserve, manteniendo la coherencia espectral necesaria para una reconstrucci3n precisa mediante IFFT.

**6. Combinaci3n de Bandas y Proyecci3n Final** Las salidas de todas las bandas se combinan:

$$\mathbf{O} = \frac{1}{B} \sum_{b=1}^B \mathbf{O}_b$$

y se proyectan de vuelta al espacio de embedding:

$$\mathbf{Y} = \text{LayerNorm}(\mathbf{O} \mathbf{W}_O + \mathbf{b}_O)$$

donde  $\mathbf{W}_O \in \mathbb{R}^{d \times d}$  y  $\mathbf{b}_O \in \mathbb{R}^d$  son parámetros aprendibles.

### 3.1.4 Análisis Matemático Rigurosamente Detallado

Para comprender completamente el funcionamiento de `FrequencyHyenaAttention`, es esencial analizar cada componente desde una perspectiva matemática profunda.

**1. Transformada de Fourier y Representación en el Dominio de la Frecuencia** La Transformada de Fourier de una señal discreta permite descomponerla en sus componentes frecuenciales:

$$\mathbf{X}_{\text{fft}}[b, k, c] = \sum_{l=0}^{L-1} \mathbf{X}[b, l, c] e^{-i2\pi kl/L}, \quad k = 0, 1, \dots, L-1$$

Esta transformación es lineal y reversible, permitiendo que la información temporal se represente en términos de amplitud y fase de diferentes frecuencias.

**2. Descomposición en Magnitud y Fase** La magnitud  $\mathbf{M}$  y fase  $\Phi$  de la señal frecuencial son fundamentales para reconstruir la señal original:

$$\mathbf{X}_{\text{fft}}[b, k, c] = \mathbf{M}[b, k, c] e^{i\Phi[b, k, c]}$$

La magnitud captura la amplitud de las frecuencias, mientras que la fase codifica la información de desplazamiento temporal.

**3. Proyección y Normalización por Banda** Cada banda de frecuencia se proyecta linealmente y se normaliza:

$$\mathbf{H}_b = \text{LayerNorm}(\mathbf{M}_b \mathbf{W}_b + \mathbf{b}_b)$$

Esto transforma la representación magnitudinal de cada banda en un espacio de características más adecuado para la atención, asegurando que las diferentes bandas contribuyan de manera equilibrada.

**4. Atención Multi-Cabeza en el Dominio Frecuencial** La atención multi-cabeza permite que cada cabeza de atención se enfoque en diferentes aspectos de la representación frecuencial:

$$\mathbf{A}_b = \text{softmax} \left( \frac{\mathbf{Q}_b \mathbf{K}_b^\top}{\sqrt{d_k}} \right) \mathbf{V}_b$$

donde:

$$\mathbf{Q}_b = \mathbf{H}_b \mathbf{W}_Q, \quad \mathbf{K}_b = \mathbf{H}_b \mathbf{W}_K, \quad \mathbf{V}_b = \mathbf{H}_b \mathbf{W}_V$$

La normalización por  $\sqrt{d_k}$  estabiliza las puntuaciones de atención, previniendo gradientes demasiado grandes.

**5. Modulación con Fase** La modulación de la salida de atención con la fase preserva la información temporal original:

$$\mathbf{O}_b = \mathbf{A}_b \odot e^{i\Phi_b}$$

Esta operación garantiza que la reconstrucción mediante IFFT mantenga la coherencia temporal y la integridad espectral de la señal.

**6. Combinación de Bandas y Proyección Final** La combinación de todas las bandas se realiza mediante una media ponderada:

$$\mathbf{O} = \frac{1}{B} \sum_{b=1}^B \mathbf{O}_b$$

Posteriormente, una capa lineal y normalización ajusta la representación combinada de vuelta al espacio de embedding:

$$\mathbf{Y} = \text{LayerNorm}(\Re(\mathbf{O}) \mathbf{W}_O + \mathbf{b}_O)$$

La parte real  $\Re(\mathbf{O})$  se extrae para asegurar que la señal reconstruida sea una representación válida en el espacio real.

### 3.1.5 Propiedades Matemáticas y Ventajas

**1. Linealidad y Reversibilidad** La FFT es una transformación lineal y reversible, lo que garantiza que no se pierda información durante la conversión al dominio frecuencial:

$$\mathcal{F}^{-1}(\mathcal{F}(\mathbf{X})) = \mathbf{X}$$

Esto es crucial para asegurar que la reconstrucción mediante IFFT sea precisa.

**2. Reducción de la Complejidad Computacional** La atención multi-cabeza tradicional tiene una complejidad de  $\mathcal{O}(L^2)$ , donde  $L$  es la longitud de la secuencia. Al operar en el dominio frecuencial y procesar bandas específicas, la complejidad se reduce significativamente:

$$\text{Complejidad Total} = \mathcal{O}(B \cdot L \cdot C \cdot \log L)$$

Esto permite manejar secuencias más largas de manera eficiente, manteniendo una alta calidad en las representaciones aprendidas.

**3. Captura de Relaciones a Largo Plazo** En el dominio frecuencial, las relaciones a largo plazo se manifiestan como componentes de baja frecuencia. La atención multi-cabeza en este dominio puede identificar y reforzar estas relaciones sin la necesidad de almacenar todas las interacciones individuales, mejorando la capacidad del modelo para capturar dependencias de largo alcance.

**4. Preservación de la Integridad Temporal** La modulación con fase asegura que la reconstrucción mediante IFFT mantenga la coherencia temporal de la señal original. Sin esta modulación, la reconstrucción podría perder información crucial sobre la estructura temporal, afectando negativamente la calidad de las representaciones aprendidas.

**5. Flexibilidad y Adaptabilidad** Al procesar diferentes bandas de frecuencia de manera independiente, `FrequencyHyenaAttention` puede adaptarse a diferentes tipos de patrones y relaciones en las secuencias de entrada. Esto mejora la versatilidad del modelo en diversas tareas de NLP, permitiendo una mejor captura de relaciones semánticas complejas.

### 3.1.6 Implicaciones Matemáticas para el Modelo

- **Eficiencia y Escalabilidad:** La reducción de la complejidad computacional mediante el uso de FFT y procesamiento por bandas permite que el modelo escale a secuencias de mayor longitud sin un aumento proporcional en el uso de recursos, facilitando su aplicación en tareas que requieren el manejo de grandes volúmenes de datos.
- **Mejora en la Captura de Relaciones Contextuales:** La capacidad de identificar y procesar relaciones a largo plazo en el dominio frecuencial mejora la comprensión contextual del modelo, permitiéndole capturar dependencias que se extienden a lo largo de toda la secuencia.
- **Estabilidad y Precisión Numérica:** La normalización por capa y la proyección lineal garantizan que las representaciones frecuenciales se mantengan estables y bien distribuidas, lo que facilita el entrenamiento y previene problemas de desvanecimiento o explosión de gradientes.
- **Integración con Otros Componentes del Modelo:** La salida de `FrequencyHyenaAttention` se integra de manera cohesiva con el resto de la arquitectura Transformer, permitiendo una fluida propagación de la información a través de las capas del modelo y asegurando que las representaciones aprendidas sean coherentes y efectivas.

### 3.1.7 Integración de FrequencyHyenaAttention en EnhancedLiquidEmbedding

**1. Función dentro de EnhancedLiquidEmbedding** Dentro de la arquitectura `EnhancedLiquidEmbedding`, `FrequencyHyenaAttention` juega un papel crucial en el enriquecimiento de las representaciones embebidas de los tokens. El proceso completo se puede describir de la siguiente manera:



1. **Embeddings Iniciales:** Los identificadores de tokens se convierten en vectores densos mediante matrices de embedding  $\mathbf{E}_{\text{token}}$ , y se suman a los embeddings de posición  $\mathbf{E}_{\text{pos}}$  para incorporar información posicional:

$$\mathbf{X} = \mathbf{E}_{\text{token}}(\mathbf{x}) + \mathbf{E}_{\text{pos}}(\mathbf{p})$$

2. **Capas Convolucionales:** Se aplican dos capas convolucionales 1D con activaciones GELU y normalización de capa para capturar patrones locales en la secuencia:

$$\mathbf{X}_1 = \text{GELU}(\text{Conv}_1(\mathbf{X}))$$

$$\mathbf{X}_2 = \text{GELU}(\text{Conv}_2(\mathbf{X}_1))$$

$$\mathbf{X}_3 = \text{LayerNorm}(\mathbf{X}_2)$$

3. **FrequencyHyenaAttention:** La representación convolucionalmente procesada se pasa a través de FrequencyHyenaAttention para capturar relaciones a largo plazo en el dominio frecuencial:

$$\mathbf{Y} = \text{FrequencyHyenaAttention}(\mathbf{X}_3)$$

4. **Compresión Adaptativa:** Basado en la complejidad de la señal, se aplica una compresión adaptativa en el dominio frecuencial, permitiendo que el modelo reduzca la dimensionalidad de la representación sin perder información relevante.

5. **Reconstrucción y Proyección Final:** La representación comprimida se reconstruye mediante IFFT y se proyecta de vuelta al espacio de embedding original:

$$\mathbf{X}_{\text{ifft}} = \mathcal{F}^{-1}(\mathbf{Y}, \text{dim} = 1)$$

$$\mathbf{X}_{\text{proj}} = \text{LayerNorm}(\mathbf{W}\mathbf{X}_{\text{ifft}} + \mathbf{b})$$

**2. Impacto en EnhancedLiquidEmbedding** La inclusión de FrequencyHyenaAttention dentro de EnhancedLiquidEmbedding tiene múltiples efectos beneficiosos:

- **Captura de Dependencias a Largo Plazo:** Al operar en el dominio frecuencial, el mecanismo de atención permite que el embedding capture relaciones que se extienden a lo largo de toda la secuencia, mejorando la representación contextual de cada token.
- **Eficiencia Computacional:** La reducción de la complejidad de la atención mediante la FFT permite manejar secuencias más largas sin un aumento significativo en el costo computacional, haciendo que el embedding sea más escalable.
- **Preservación de Información Espectral:** La modulación con fase garantiza que la información espectral esencial se mantenga intacta, lo que es crucial para la reconstrucción precisa de la señal original y para mantener la coherencia temporal en las representaciones embebidas.
- **Integración Sólida con Convoluciones:** La combinación de convoluciones locales con atención en el dominio frecuencial permite que el embedding capture tanto patrones locales como relaciones globales, resultando en representaciones más ricas y completas.
- **Regularización Implícita:** La normalización de capa y las proyecciones lineales actúan como mecanismos de regularización, estabilizando las activaciones y facilitando un entrenamiento más eficiente y robusto.

### 3.1.8 Impacto de FrequencyHyenaAttention en el Modelo LiquidFoundationModelOptimized

**1. Función dentro del Modelo Completo** En el contexto del modelo LiquidFoundationModelOptimized, FrequencyHyenaAttention desempeña un rol integral en las siguientes etapas:

1. **Procesamiento de Embeddings:**

$$(\mathbf{X}, \mathcal{L}_{\text{recon}}) = \text{EnhancedLiquidEmbedding}(\mathbf{x})$$

Aquí, FrequencyHyenaAttention mejora las representaciones embebidas al capturar relaciones a largo plazo y reducir la dimensionalidad de la secuencia mediante la compresión adaptativa.

## 2. Propagación a través de las Capas Transformer:

$$(\mathbf{X}, \mathcal{L}_{\text{total\_entropy}}) = \sum_{l=1}^K \text{ImprovedTransformerBlock}_l(\mathbf{X})$$

Las representaciones mejoradas por **FrequencyHyenaAttention** se procesan a través de múltiples bloques Transformer mejorados, donde cada bloque puede aprovechar las dependencias a largo plazo capturadas anteriormente.

## 3. Generación de Salida:

$$\mathbf{Y} = \text{LayerNorm}(\mathbf{H}_d)\mathbf{W}_{\text{out}} + \mathbf{b}_{\text{out}}$$

La salida final del decoder se proyecta a través de una capa lineal, donde las representaciones enriquecidas por **FrequencyHyenaAttention** contribuyen a generar logits de alta calidad para la predicción de tokens.

**2. Beneficios para el Modelo Completo** La integración de **FrequencyHyenaAttention** en **LiquidFoundationModel** aporta múltiples ventajas:

- **Mejora de la Representación Contextual:** La capacidad de capturar relaciones a largo plazo en el dominio frecuencial enriquece las representaciones contextuales, lo que es fundamental para tareas de generación y comprensión de lenguaje natural.
- **Escalabilidad y Manejo de Secuencias Largas:** Al reducir la complejidad computacional de la atención y permitir la compresión adaptativa, **FrequencyHyenaAttention** facilita que el modelo maneje secuencias extremadamente largas sin comprometer la eficiencia o el rendimiento.
- **Optimización de Recursos Computacionales:** La eficiencia mejorada en el procesamiento frecuencial reduce el uso de memoria y acelera las operaciones de atención, permitiendo entrenar y desplegar modelos más grandes y profundos de manera práctica.
- **Preservación de Información Espectral para Reconstrucción Precisa:** La modulación con fase asegura que la información espectral relevante se mantenga intacta, lo que es crucial para tareas que requieren una reconstrucción precisa de la señal o una comprensión detallada de la estructura temporal.
- **Facilitación de la Generalización:** Las representaciones contextuales más ricas y completas permiten que el modelo generalice mejor a datos no vistos, mejorando su desempeño en diversas tareas y dominios de NLP.
- **Compatibilidad con Otros Componentes Avanzados:** **FrequencyHyenaAttention** se integra de manera cohesiva con otros componentes avanzados del modelo, como MoE, convoluciones deformables y memoria externa basada en LSTM, potenciando la capacidad global del modelo para manejar tareas complejas y variadas.

**3. Impacto en la Calidad de las Representaciones** Las representaciones generadas por **FrequencyHyenaAttention** son más robustas y ricas en información contextual debido a la capacidad de capturar y procesar relaciones a largo plazo en el dominio frecuencial. Esto se traduce en una mayor precisión y coherencia en las tareas de generación de lenguaje, traducción y comprensión de texto, ya que el modelo puede mantener una visión holística de la secuencia de entrada.

**4. Aceleración del Entrenamiento y la Inferencia** La eficiencia computacional mejorada permite que **LiquidFoundationModelOptimized** se entrene más rápidamente y realice inferencias de manera más ágil. Esto es particularmente beneficioso en entornos de producción donde el tiempo de respuesta y el uso eficiente de recursos son críticos.

**5. Robustez y Estabilidad del Modelo** La combinación de normalización de capa, proyecciones lineales y compresión adaptativa dentro de **FrequencyHyenaAttention** contribuye a una mayor estabilidad numérica y robustez durante el entrenamiento. Esto previene problemas como el desvanecimiento o explosión de gradientes, asegurando que el modelo pueda aprender de manera efectiva sin interrupciones.

## 3.2 EnhancedLiquidEmbedding

**Propósito:** Convertir identificadores de tokens en representaciones vectoriales densas que incorporan información posicional y aplicar técnicas de compresión adaptativa para manejar secuencias largas de manera eficiente.

### 3.2.1 Descripción Algorítmica

#### 1. Embeddings Iniciales:

- Sea  $\mathbf{x} \in \mathbb{N}^{B \times L}$  la matriz de identificadores de tokens, donde  $B$  es el tamaño del batch y  $L$  la longitud de la secuencia.
- Los embeddings de tokens se representan como  $\mathbf{E}_{\text{token}} \in \mathbb{R}^{V \times d}$ , donde  $V$  es el tamaño del vocabulario y  $d$  la dimensión de embedding.
- Los embeddings de posición son  $\mathbf{E}_{\text{pos}} \in \mathbb{R}^{L_{\text{max}} \times d}$ , con  $L_{\text{max}}$  la longitud máxima de secuencia.

#### 2. Combinación de Embeddings:

$$\mathbf{X} = \mathbf{E}_{\text{token}}(\mathbf{x}) + \mathbf{E}_{\text{pos}}(\mathbf{p})$$

donde  $\mathbf{p}$  es el vector de posiciones.

#### 3. Capas Convolucionales:

$$\mathbf{X}_1 = \text{GELU}(\text{Conv}_1(\mathbf{X}))$$

$$\mathbf{X}_2 = \text{GELU}(\text{Conv}_2(\mathbf{X}_1))$$

donde  $\text{Conv}_i$  representa la capa convolucional  $i$ .

#### 4. Normalización de Capa y Padding:

$$\mathbf{X}_3 = \text{LayerNorm}(\mathbf{X}_2)$$

$$L_{\text{padded}} = 2^{\lceil \log_2 L \rceil}$$

$$\mathbf{X}_{\text{padded}} = \text{Pad}(\mathbf{X}_3, L_{\text{padded}})$$

#### 5. Transformada de Fourier Rápida (FFT) y Compresión Adaptativa:

$$\mathbf{X}_{\text{fft}} = \text{FFT}(\mathbf{X}_{\text{padded}}, \text{dim} = 1)$$

$$\text{magnitude} = |\mathbf{X}_{\text{fft}}|$$

$$\text{freq\_energy} = \frac{1}{Bd} \sum_{b=1}^B \sum_{c=1}^d \sum_{k=\lfloor 0.8L \rfloor}^{L-1} |\mathbf{X}_{\text{fft}}[b, k, c]|^2$$

$$\text{local\_var\_norm} = \frac{1}{BdL} \sum_{b=1}^B \sum_{c=1}^d \sum_{l=1}^L |\mathbf{X}[b, l, c] - \mathbf{X}[b, l-1, c]|^2$$

$$\text{complexity} = 0.7 \cdot \text{freq\_energy} + 0.3 \cdot \text{local\_var\_norm}$$

$$\text{compression\_ratio} = \text{clamp}(\text{min\_compression\_ratio} + (\text{base\_compression\_ratio} - \text{min\_compression\_ratio}) \cdot \text{complexity})$$

$$N_b = \text{clamp}(\text{compression\_ratio} \times L, 1, L)$$

$$\mathbf{X}_{\text{fft\_compressed}}[b, :, N_b, c] = \mathbf{X}_{\text{fft}}[b, :, N_b, c] \quad \forall b, c$$

#### 6. Reconstrucción mediante IFFT y Pérdida de Reconstrucción:

$$\mathbf{X}_{\text{ifft}} = \text{IFFT}(\mathbf{X}_{\text{fft\_compressed}}, \text{n} = L, \text{dim} = 1)$$

$$\mathbf{X}_{\text{proj}} = \text{LayerNorm}(\mathbf{W}\mathbf{X}_{\text{ifft}} + \mathbf{b})$$

$$\mathcal{L}_{\text{recon}} = \frac{\sum_{b=1}^B \sum_{c=1}^d \sum_{l=1}^L \mathbb{I}_{\text{mask}} \cdot |\mathbf{X}_{\text{proj}}[b, l, c] - \mathbf{X}_{\text{recon}}[b, l, c]|}{\sum \mathbb{I}_{\text{mask}} + \epsilon}$$

donde  $\mathbb{I}_{\text{mask}}$  es una máscara que indica las posiciones válidas tras la compresión.

### 3.2.2 Fundamentos Matemáticos

- **Embeddings:** Los embeddings transforman identificadores discretos en representaciones vectoriales continuas. Matemáticamente, esto se realiza mediante matrices de embedding  $\mathbf{E}_{\text{token}}$  y  $\mathbf{E}_{\text{pos}}$ , donde cada fila de estas matrices corresponde a un vector de embedding para un token o una posición específica. Esta transformación permite que el modelo capture relaciones semánticas entre tokens a través de las distancias en el espacio vectorial.

- **Convoluciones 1D:** Las convoluciones 1D se definen por la operación:

$$(\text{Conv}_k * \mathbf{X})[b, l, c'] = \sum_{i=1}^d \sum_{m=-k}^k \mathbf{W}[c', i, m] \cdot \mathbf{X}[b, l + m, i]$$

donde  $\mathbf{W}$  son los pesos aprendibles del filtro de convolución. Estas operaciones capturan patrones locales en la secuencia, permitiendo que el modelo identifique características espaciales o temporales relevantes.

- **Layer Normalization:** La normalización de capa estabiliza las activaciones al normalizar cada muestra de manera independiente:

$$\text{LayerNorm}(\mathbf{x}) = \frac{\mathbf{x} - \mu}{\sqrt{\sigma^2 + \epsilon}} \cdot \gamma + \beta$$

donde  $\mu$  y  $\sigma^2$  son la media y varianza de las activaciones, respectivamente, y  $\gamma, \beta$  son parámetros aprendibles que permiten la escala y desplazamiento de las activaciones normalizadas.

- **FFT e IFFT:** Las transformadas de Fourier rápida (FFT) y su inversa (IFFT) son transformadas lineales que convierten señales del dominio del tiempo al dominio de la frecuencia y viceversa. En este contexto, se utilizan para comprimir la representación de las secuencias, reteniendo solo las componentes frecuenciales más significativas. Matemáticamente, la FFT de una señal  $\mathbf{X}$  en la dimensión de la secuencia se define como:

$$\mathbf{X}_{\text{fft}}[b, k, c] = \sum_{l=0}^{L-1} \mathbf{X}[b, l, c] e^{-i2\pi kl/L}$$

donde  $i$  es la unidad imaginaria, y  $k$  es la frecuencia.

- **Compresión Adaptativa:** La compresión adaptativa ajusta el grado de compresión basado en la complejidad de la secuencia. La complejidad se mide como una combinación ponderada de la energía frecuencial y la variación local de la señal en el dominio del tiempo:

$$\text{complexity} = 0.7 \cdot \text{freq\_energy} + 0.3 \cdot \text{local\_var\_norm}$$

El ratio de compresión `compression_ratio` se calcula como:

$$\text{compression\_ratio} = \text{clamp}(\text{min\_compression\_ratio} + (\text{base\_compression\_ratio} - \text{min\_compression\_ratio}) \cdot \text{complexity})$$

Esto permite que el modelo mantenga un balance óptimo entre eficiencia computacional y preservación de información relevante.

- **Pérdida de Reconstrucción:** La pérdida de reconstrucción  $\mathcal{L}_{\text{recon}}$  mide la discrepancia entre la señal reconstruida  $\mathbf{X}_{\text{proj}}$  y la señal objetivo  $\mathbf{X}_{\text{recon}}$ , aplicada únicamente en las posiciones válidas tras la compresión:

$$\mathcal{L}_{\text{recon}} = \frac{\sum_{b=1}^B \sum_{c=1}^d \sum_{l=1}^L \mathbb{I}_{\text{mask}} \cdot |\mathbf{X}_{\text{proj}}[b, l, c] - \mathbf{X}_{\text{recon}}[b, l, c]|}{\sum \mathbb{I}_{\text{mask}} + \epsilon}$$

donde  $\epsilon$  es un término pequeño añadido para evitar divisiones por cero.

### 3.2.3 Análisis Detallado

La capa `EnhancedLiquidEmbedding` es una capa de embedding mejorada que combina convoluciones, atención en el dominio de frecuencia y compresión adaptativa. A continuación, se detallan sus componentes y operaciones matemáticas:

1. **Embeddings Iniciales:** Los embeddings de tokens y posición se combinan y normalizan para formar una representación embebida robusta de la entrada.
2. **Capas Convolucionales:** Se aplican dos capas convolucionales con activación GELU y normalización de capa para capturar patrones locales en la secuencia.
3. **Transformada Rápida de Fourier (FFT):** La representación embebida se transforma al dominio de la frecuencia mediante FFT, facilitando la compresión adaptativa basada en la complejidad de la señal.
4. **Compresión Adaptativa:** Basado en la complejidad de la señal, se determina un ratio de compresión que ajusta dinámicamente la longitud de la representación en frecuencia, creando una máscara que indica las posiciones seleccionadas.
5. **Reconstrucción mediante IFFT:** La representación comprimida se transforma de vuelta al dominio del tiempo mediante IFFT, seguida de una proyección y normalización adicional.
6. **Cálculo de Pérdida:** Se calcula una pérdida de reconstrucción que penaliza la discrepancia entre la salida y el objetivo, considerando únicamente las posiciones válidas tras la compresión.

### 3.2.4 Implicaciones para el Modelo

- **Eficiencia Computacional:** La compresión adaptativa reduce la cantidad de datos a procesar en las siguientes capas, disminuyendo la carga computacional y permitiendo manejar secuencias más largas sin un incremento proporcional en el consumo de recursos.
- **Preservación de Información Relevante:** Al ajustar el ratio de compresión según la complejidad de la secuencia, el modelo puede mantener una representación rica en información para secuencias complejas mientras reduce la redundancia en secuencias más simples.
- **Estabilidad durante el Entrenamiento:** La inclusión de términos de regularización como la pérdida de reconstrucción y la normalización de capa contribuye a una mejor convergencia y estabilidad numérica, mitigando problemas como el desvanecimiento o explosión de gradientes.
- **Capacidad de Generalización:** Al incorporar convoluciones y transformadas de Fourier, el modelo puede capturar patrones tanto locales como globales en los datos, mejorando su capacidad para generalizar a diferentes contextos y tareas dentro del dominio de NLP.

### 3.3 EnhancedLocalAttentionWithGQA

**Propósito:** Implementar un mecanismo de atención local mejorado que optimiza la eficiencia y precisión mediante grupos de atención y ventanas deslizantes.

#### 3.3.1 Descripción Algorítmica

##### 1. Proyección Lineal de Q, K, V:

$$\mathbf{Q} = \mathbf{XW}_Q, \quad \mathbf{K} = \mathbf{XW}_K, \quad \mathbf{V} = \mathbf{XW}_V$$

donde  $\mathbf{W}_Q \in \mathbb{R}^{d \times d}$ ,  $\mathbf{W}_K \in \mathbb{R}^{d \times \frac{d}{g}}$ ,  $\mathbf{W}_V \in \mathbb{R}^{d \times \frac{d}{g}}$ , y  $g$  es el número de grupos de  $K$  y  $V$ .

##### 2. División en Grupos:

$$\mathbf{K}_g = \text{reshape}(\mathbf{K}, [B, L, g, \frac{d}{g}])$$

$$\mathbf{V}_g = \text{reshape}(\mathbf{V}, [B, L, g, \frac{d}{g}])$$

$$\mathbf{K}'_g = \text{repeat}(\mathbf{K}_g, [1, 1, h/g, 1])$$

$$\mathbf{V}'_g = \text{repeat}(\mathbf{V}_g, [1, 1, h/g, 1])$$

donde  $h$  es el número de cabezas de atención.

##### 3. Ventanas Deslizantes:

$$w_i = \{l \mid s_i \leq l < e_i\}$$

donde  $s_i$  y  $e_i$  son los índices de inicio y fin de la ventana  $i$ .

##### 4. Atención Flash en Ventanas:

$$\mathbf{A}_i = \text{FlashAttention}(\mathbf{Q}_{w_i}, \mathbf{K}_{g,w_i}, \mathbf{V}_{g,w_i})$$

$$\mathbf{A}_i = \text{softmax} \left( \frac{\mathbf{Q}_{w_i} \mathbf{K}_{g,w_i}^\top}{\sqrt{d_k}} \right) \mathbf{V}_{g,w_i}$$

donde  $d_k = \frac{d}{h}$ .

##### 5. Concatenación y Proyección Final:

$$\mathbf{O} = \text{concat}(\mathbf{A}_1, \mathbf{A}_2, \dots, \mathbf{A}_W) \mathbf{W}_O$$

donde  $\mathbf{W}_O \in \mathbb{R}^{d \times d}$ .

#### 3.3.2 Fundamentos Matemáticos

- **Atención Multi-Cabeza:** La atención multi-cabeza se define como:

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax} \left( \frac{\mathbf{QK}^\top}{\sqrt{d_k}} \right) \mathbf{V}$$

donde  $d_k$  es la dimensión de las proyecciones de las consultas  $Q$ . Este mecanismo permite que el modelo enfoque en diferentes aspectos de la secuencia simultáneamente, capturando diversas relaciones semánticas.

- **Flash Attention:** Flash Attention es una implementación optimizada de la atención estándar que reduce la complejidad computacional y el uso de memoria. Matemáticamente, optimiza la operación de atención reduciendo la complejidad de  $\mathcal{O}(L^2)$  a  $\mathcal{O}(L \cdot W)$ , donde  $W$  es el tamaño de la ventana. Esto se logra mediante técnicas avanzadas de paralelización y aprovechamiento de la memoria caché.

- **Ventanas Deslizantes:** La atención basada en ventanas deslizantes restringe el ámbito de atención a regiones locales de la secuencia. Esto se formaliza como la división de la secuencia en ventanas  $w_i$  que se superponen parcialmente, permitiendo que cada ventana deslice sobre la secuencia con un solapamiento definido. Esta técnica reduce la complejidad de la atención y mantiene la capacidad de capturar relaciones locales y de largo alcance mediante múltiples capas de atención.
- **Grupos de Atención (GQA - Grouped Query Attention):** El agrupamiento de las proyecciones de  $K$  y  $V$  en  $\mathbf{K}_g$  y  $\mathbf{V}_g$  respectivamente, permite manejar de manera eficiente las interacciones entre múltiples cabezas de atención. Esto se logra dividiendo las proyecciones en  $g$  grupos, lo que facilita una mayor diversidad en las representaciones de atención sin incrementar significativamente la complejidad computacional.

### 3.3.3 Análisis Detallado

La implementación de atención local mejorada con GQA introduce una eficiencia computacional significativa al reducir la complejidad de la atención tradicional. Al dividir las proyecciones de  $K$  y  $V$  en grupos y emplear ventanas deslizantes, el modelo puede manejar secuencias de mayor longitud sin un incremento proporcional en el uso de memoria o tiempo de cómputo.

Además, Flash Attention optimiza la operación de atención al reducir la redundancia en cálculos y aprovechar mejor la arquitectura de hardware moderna, lo que se traduce en una aceleración del proceso de atención sin sacrificar la precisión.

### 3.3.4 Implicaciones para el Modelo

- **Escalabilidad:** La reducción de la complejidad computacional permite que el modelo escale a secuencias más largas, aumentando su aplicabilidad en tareas que requieren el procesamiento de grandes volúmenes de texto o información.
- **Eficiencia en el Uso de Memoria:** Al limitar el ámbito de atención a ventanas locales y agrupar las proyecciones de  $K$  y  $V$ , se optimiza el uso de la memoria, permitiendo entrenar modelos más profundos o con mayor capacidad sin exceder las limitaciones de hardware.
- **Diversidad de Atención:** El agrupamiento de las proyecciones de atención facilita que diferentes cabezas enfoquen en distintos aspectos de la secuencia, mejorando la capacidad del modelo para capturar relaciones semánticas complejas y variadas.
- **Velocidad de Entrenamiento e Inferencia:** Las optimizaciones introducidas por Flash Attention y las ventanas deslizantes reducen el tiempo de cómputo requerido para cada paso de atención, acelerando tanto el entrenamiento como la inferencia del modelo.

### 3.4 MoELayer (Mixture of Experts)

**Propósito:** Introducir un mecanismo de enrutamiento adaptativo que dirige diferentes partes de la entrada a expertos especializados, mejorando la capacidad del modelo para manejar diversas tareas o contextos.

#### 3.4.1 Descripción Algorítmica

##### 1. Capa de Gateo:

$$\mathbf{z} = \mathbf{x}\mathbf{W}_g + \mathbf{b}_g$$

donde  $\mathbf{W}_g \in \mathbb{R}^{d \times m}$  y  $m$  es el número de expertos.

$$\mathbf{p} = \text{softmax}(\mathbf{z})$$

##### 2. Selección de Expertos:

$$\text{TopK}(\mathbf{p}) = \{e_1, e_2, \dots, e_k\}$$

Si `dynamic_k` está activado:

$$k = \max(1, \min(m, \lfloor k_{\text{base}} \times (1 + \text{complexity}) \rfloor))$$

donde la complejidad se define como:

$$\text{complexity} = H(\mathbf{p}) = - \sum_{i=1}^m p_i \log p_i$$

##### 3. Enrutamiento y Cálculo de Salidas:

$$\mathbf{y}_j = \mathbf{W}_j \mathbf{x} + \mathbf{b}_j \quad \forall e_j \in \{e_1, e_2, \dots, e_k\}$$

$$\mathbf{y} = \sum_{j=1}^k p_j \mathbf{y}_j$$

##### 4. Regularización y Penalización:

$$\mathcal{L}_{\text{entropy}} = - \sum_{i=1}^m p_i \log p_i$$

$$\mathcal{L}_{\text{overuse}} = \sum_{j=1}^m \text{ReLU}\left(\frac{n_j}{N} - \theta\right)$$

$$\mathcal{L}_{\text{total}} = \lambda_{\text{entropy}} \mathcal{L}_{\text{entropy}} + \lambda_{\text{overuse}} \mathcal{L}_{\text{overuse}}$$

#### 3.4.2 Fundamentos Matemáticos

- **Mixture of Experts (MoE):** La arquitectura MoE combina múltiples submodelos, denominados expertos, donde cada uno se especializa en diferentes aspectos de la tarea. Matemáticamente, esto se representa como una combinación ponderada de las salidas de los expertos:

$$\mathbf{y} = \sum_{j=1}^k p_j \mathbf{y}_j$$

donde  $p_j$  son las probabilidades asignadas por la capa de gateo a cada experto seleccionado.

- **Softmax Gate:** La función softmax aplicada en la capa de gateo asegura que las probabilidades  $\mathbf{p}$  sumen a uno, facilitando una distribución probabilística sobre los expertos. Esto es crucial para permitir que el modelo sea completamente diferenciable y permita el entrenamiento end-to-end mediante backpropagation.



- **Entropía de la Distribución de Gateo:** La entropía  $H(\mathbf{p})$  mide la incertidumbre en la distribución de probabilidades. Al maximizar esta entropía, se promueve una asignación más uniforme de las entradas a los diferentes expertos, evitando que un subconjunto reducido de expertos domine el procesamiento.
- **Penalización por Uso Excesivo:** La penalización  $\mathcal{L}_{\text{overuse}}$  impone restricciones para que el uso de cada experto se mantenga dentro de límites razonables. Esto se formaliza mediante la función ReLU que activa la penalización cuando el ratio de uso de un experto  $\frac{n_j}{N}$  excede un umbral  $\theta$ :

$$\mathcal{L}_{\text{overuse}} = \sum_{j=1}^m \text{ReLU} \left( \frac{n_j}{N} - \theta \right)$$

### 3.4.3 Análisis Detallado

El mecanismo MoE introduce una flexibilidad significativa en el modelo al permitir que diferentes expertos se especialicen en distintos aspectos de los datos de entrada. Este enfoque permite que el modelo escale su capacidad de manera eficiente, ya que solo un subconjunto de expertos se activa para cada entrada, reduciendo el costo computacional en comparación con un modelo monolítico de igual capacidad.

La capa de gateo y la selección dinámica de expertos aseguran que la asignación de expertos sea adaptativa y dependiente del contenido de la entrada. La regularización mediante la entropía y la penalización por uso excesivo promueve una distribución equilibrada del trabajo entre los expertos, mejorando la capacidad de generalización del modelo y evitando el sobreajuste de ciertos expertos a subconjuntos específicos de datos.

### 3.4.4 Implicaciones para el Modelo

- **Aumento de la Capacidad del Modelo:** Al incorporar múltiples expertos, el modelo puede aprender representaciones más ricas y especializadas, mejorando su desempeño en tareas complejas y variadas.
- **Eficiencia Computacional:** Solo un subconjunto de expertos se activa para cada entrada, lo que permite que el modelo escale su capacidad sin un incremento proporcional en el costo computacional y el consumo de memoria.
- **Mejora de la Generalización:** La distribución equilibrada de la carga de trabajo entre los expertos previene el sobreajuste y fomenta una mejor generalización a datos no vistos.
- **Flexibilidad Adaptativa:** La selección dinámica de expertos permite que el modelo se adapte a la complejidad de la entrada, optimizando el uso de recursos según las necesidades específicas de cada tarea.

### 3.5 DeformableConv1d

**Propósito:** Implementar una convolución 1D deformable que adapta dinámicamente los desplazamientos de los filtros, permitiendo una mayor flexibilidad en la captura de patrones locales variables.

#### 3.5.1 Descripción Algorítmica

1. **Generación de Offsets:**

$$\Delta = \text{Conv}_{\text{offset}}(\mathbf{x})$$

donde  $\mathbf{x} \in \mathbb{R}^{N \times C \times L}$  y  $\Delta \in \mathbb{R}^{N \times 2k \times L_{\text{out}}}$ .

2. **Aplicación de Offsets:**

$$\Delta = \Delta.\text{view}(N, k, 2, L_{\text{out}})$$

$$\Delta = \Delta.\text{permute}(0, 3, 1, 2)$$

3. **Interpolación Bilineal:**

$$\mathbf{g} = \mathbf{G} + \Delta$$

donde  $\mathbf{G}$  es la grilla base de posiciones.

$$\mathbf{x}_{\text{deform}} = (1 - \alpha)\mathbf{x}_{\text{left}} + \alpha\mathbf{x}_{\text{right}}$$

donde  $\alpha$  es la parte fraccionaria de  $\mathbf{g}$ .

4. **Convolución Principal:**

$$\mathbf{y} = \text{Conv}_{\text{main}}(\mathbf{x}_{\text{deform}})$$

#### 3.5.2 Fundamentos Matemáticos

- **Convoluciones Deformables:** Las convoluciones deformables permiten que las posiciones de muestreo en la convolución sean dinámicas y aprendidas, adaptándose a patrones de entrada complejos. Matemáticamente, la salida de una convolución deformable se define como:

$$\mathbf{y}[n, c, l] = \sum_{m=1}^k \mathbf{w}[c, m] \cdot \mathbf{x}[n, c, l \cdot s + m + \Delta_{n, m, l}]$$

donde  $s$  es el stride y  $\Delta$  son los offsets aprendidos que determinan los desplazamientos de los filtros.

- **Interpolación Bilineal:** La interpolación bilineal permite la extracción de características en posiciones no enteras mediante una combinación lineal de valores adyacentes. Matemáticamente, para una posición deformada  $g$ , la señal interpolada  $\mathbf{x}_{\text{deform}}$  se calcula como:

$$\mathbf{x}_{\text{deform}} = (1 - \alpha)\mathbf{x}_{\text{left}} + \alpha\mathbf{x}_{\text{right}}$$

donde  $\alpha = g - \lfloor g \rfloor$  es la parte fraccionaria de  $g$ .

- **Reorganización de Offsets:** La manipulación de tensores para ajustar los desplazamientos según las dimensiones de entrada y salida garantiza que los offsets se apliquen correctamente a través de la secuencia. Esto se logra mediante operaciones de reshaping y permutación que alinean los offsets con las posiciones de la grilla base.

#### 3.5.3 Análisis Detallado

La capacidad de adaptar dinámicamente los desplazamientos de los filtros en las convoluciones permite que el modelo capture patrones locales que pueden variar en su ubicación dentro de la secuencia. Esta flexibilidad es especialmente útil en tareas de NLP donde las dependencias entre palabras pueden no ser estrictamente locales o pueden variar en distancia.

La interpolación bilineal asegura que las transiciones entre diferentes patrones sean suaves y continuas, evitando discontinuidades que podrían afectar negativamente la capacidad de aprendizaje del modelo.

### 3.5.4 Implicaciones para el Modelo

- **Flexibilidad en la Captura de Patrones:** Las convoluciones deformables permiten que el modelo se adapte a diferentes estructuras de datos, mejorando su capacidad para identificar patrones que no son fijos en su ubicación.
- **Mejora de la Representación de Características:** Al permitir desplazamientos dinámicos, el modelo puede capturar relaciones más complejas y sutiles entre los elementos de la secuencia, enriqueciendo las representaciones aprendidas.
- **Aumento de la Robustez:** La capacidad de adaptarse a variaciones en los datos de entrada mejora la robustez del modelo frente a diferentes tipos de secuencias y contextos.
- **Optimización Computacional:** Aunque las convoluciones deformables agregan una complejidad adicional, la implementación optimizada asegura que el impacto en el rendimiento computacional sea mínimo, manteniendo la eficiencia del modelo.

### 3.6 OptimizedGatedConvolution

**Propósito:** Combinar convoluciones deformables con mecanismos de puertas para controlar el flujo de información, mejorando la capacidad del modelo para manejar dependencias complejas.

#### 3.6.1 Descripción Algorítmica

1. **Aplicación de Convolución Deformable:**

$$\mathbf{y} = \text{DeformConv1d}(\mathbf{x}) = [\mathbf{y}_{\text{main}}, \mathbf{y}_{\text{gate}}]$$

2. **Activaciones No Lineales:**

$$\mathbf{y}'_{\text{main}} = \text{GELU}(\mathbf{y}_{\text{main}})$$

$$\mathbf{y}'_{\text{gate}} = \sigma(\mathbf{y}_{\text{gate}})$$

3. **Mecanismo de Puerta:**

$$\mathbf{y}_{\text{gated}} = \mathbf{y}'_{\text{main}} \otimes \mathbf{y}'_{\text{gate}}$$

4. **Normalización y Dropout:**

$$\mathbf{y}_{\text{norm}} = \text{LayerNorm}(\mathbf{y}_{\text{gated}})$$

$$\mathbf{y}_{\text{drop}} = \text{Dropout}(\mathbf{y}_{\text{norm}})$$

#### 3.6.2 Fundamentos Matemáticos

- **Mecanismo de Puerta (Gated Mechanism):** El mecanismo de puertas controla el flujo de información, permitiendo que ciertas características sean retenidas o descartadas. Matemáticamente, se define como:

$$\mathbf{y}_{\text{gated}} = \mathbf{y}'_{\text{main}} \otimes \sigma(\mathbf{y}_{\text{gate}})$$

donde  $\sigma$  es la función sigmoide que actúa como una puerta de activación, y  $\otimes$  denota la multiplicación elemento a elemento.

- **Activaciones No Lineales:** Las funciones de activación no lineales, como GELU y la sigmoide, introducen no linealidades que permiten al modelo aprender representaciones más complejas. GELU se define como:

$$\text{GELU}(x) = x\Phi(x)$$

donde  $\Phi(x)$  es la función de distribución acumulativa de la distribución normal estándar.

- **Normalización de Capa:** La normalización de capa aplicada después del mecanismo de puertas asegura que las activaciones estén dentro de un rango estable, facilitando la convergencia durante el entrenamiento y mejorando la estabilidad numérica.

#### 3.6.3 Análisis Detallado

La combinación de convoluciones deformables optimizadas con mecanismos de puertas introduce un control dinámico sobre las características extraídas por las convoluciones. Este control permite al modelo priorizar o suprimir ciertas características según la relevancia para la tarea específica, mejorando así la capacidad de modelado y la eficiencia del aprendizaje.

La aplicación de dropout después de la normalización contribuye a la regularización del modelo, previniendo el sobreajuste y promoviendo la generalización a datos no vistos.

#### 3.6.4 Implicaciones para el Modelo

- **Control Dinámico del Flujo de Información:** El mecanismo de puertas permite que el modelo ajuste dinámicamente qué información es relevante en cada contexto, mejorando la capacidad de capturar dependencias complejas y contextos variados.
- **Mejora de la Representación de Características:** Al controlar qué características son retenidas o descartadas, el modelo puede enfocarse en aspectos más significativos de la secuencia, enriqueciendo las representaciones aprendidas.

- **Regularización y Estabilidad:** La normalización de capa y el uso de dropout contribuyen a una mayor estabilidad durante el entrenamiento y previenen el sobreajuste, asegurando que el modelo generalice mejor a datos no vistos.
- **Aumento de la Flexibilidad:** La combinación de técnicas permite que el modelo se adapte a una amplia variedad de patrones de datos, mejorando su versatilidad en diferentes tareas de NLP.

### 3.7 EnhancedLSTM

**Propósito:** Introducir una versión mejorada de la arquitectura LSTM que integra capas adicionales y mecanismos de regularización para capturar dependencias secuenciales complejas.

#### 3.7.1 Descripción Algorítmica

1. **Capa LSTM Estándar:**

$$\mathbf{h}_t^{(l)}, \mathbf{c}_t^{(l)} = \text{LSTM}^{(l)}(\mathbf{x}_t^{(l)}, \mathbf{h}_{t-1}^{(l)}, \mathbf{c}_{t-1}^{(l)})$$

2. **Capa de Salida Adicional:**

$$\mathbf{y}_t^{(l)} = \text{GELU}(\mathbf{W}_o \mathbf{h}_t^{(l)} + \mathbf{b}_o)$$

$$\mathbf{y}_t^{(l)} = \mathbf{W}'_o \mathbf{y}_t^{(l)} + \mathbf{b}'_o$$

3. **Conexión Residual:**

$$\mathbf{y}_t^{(l)} = \mathbf{y}_t^{(l)} + \mathbf{x}_t^{(l)}$$

4. **Dropout:**

$$\mathbf{y}_t^{(l)} = \text{Dropout}(\mathbf{y}_t^{(l)})$$

#### 3.7.2 Fundamentos Matemáticos

- **Arquitectura LSTM:** Las LSTM (Long Short-Term Memory) son una variante de las redes neuronales recurrentes diseñadas para manejar dependencias a largo plazo en secuencias. Las ecuaciones de actualización de las puertas de entrada, olvido y salida se definen como:

$$\mathbf{f}_t = \sigma(\mathbf{W}_f \mathbf{x}_t + \mathbf{U}_f \mathbf{h}_{t-1} + \mathbf{b}_f)$$

$$\mathbf{i}_t = \sigma(\mathbf{W}_i \mathbf{x}_t + \mathbf{U}_i \mathbf{h}_{t-1} + \mathbf{b}_i)$$

$$\mathbf{o}_t = \sigma(\mathbf{W}_o \mathbf{x}_t + \mathbf{U}_o \mathbf{h}_{t-1} + \mathbf{b}_o)$$

$$\mathbf{g}_t = \tanh(\mathbf{W}_c \mathbf{x}_t + \mathbf{U}_c \mathbf{h}_{t-1} + \mathbf{b}_c)$$

$$\mathbf{c}_t = \mathbf{f}_t \odot \mathbf{c}_{t-1} + \mathbf{i}_t \odot \mathbf{g}_t$$

$$\mathbf{h}_t = \mathbf{o}_t \odot \tanh(\mathbf{c}_t)$$

donde  $\sigma$  es la función sigmoide,  $\tanh$  es la función tangente hiperbólica, y  $\odot$  denota la multiplicación elemento a elemento.

- **Red Feedforward y Conexión Residual:** La inclusión de una red feedforward y conexiones residuales mejora la capacidad de modelado y facilita el flujo de gradientes a través de la red profunda, estabilizando el entrenamiento y permitiendo que el modelo aprenda representaciones más complejas.
- **Regularización mediante Dropout:** El dropout es una técnica de regularización que previene el sobreajuste al desactivar aleatoriamente neuronas durante el entrenamiento. Matemáticamente, esto se implementa como:

$$\mathbf{x}_{\text{drop}} = \mathbf{x} \otimes \mathbf{m}$$

donde  $\mathbf{m} \sim \text{Bernoulli}(p)$  es una máscara binaria que determina qué neuronas están activas.

#### 3.7.3 Análisis Detallado

La integración de una capa de salida adicional con activaciones no lineales y una conexión residual permite que la información fluya de manera más eficiente a través de la red, facilitando el aprendizaje de representaciones complejas y mejorando la capacidad del modelo para capturar dependencias a largo plazo en las secuencias de entrada.

La aplicación de dropout después de la capa LSTM y la capa de salida adicional introduce una regularización adicional que ayuda a prevenir el sobreajuste, promoviendo una mejor generalización del modelo a datos no vistos.

#### 3.7.4 Implicaciones para el Modelo

- **Captura de Dependencias a Largo Plazo:** La arquitectura LSTM mejorada permite al modelo mantener y actualizar estados ocultos que capturan información de largo alcance en las secuencias, mejorando la coherencia y relevancia de las representaciones generadas.
- **Flujo de Gradientes Mejorado:** Las conexiones residuales facilitan el flujo de gradientes durante el entrenamiento, mitigando problemas como el desvanecimiento de gradientes y permitiendo que el modelo aprenda de manera más eficiente.
- **Regularización Avanzada:** El uso combinado de dropout y regularización por entropía en la capa MoE contribuye a un mejor equilibrio entre capacidad de modelado y prevención del sobreajuste, mejorando la capacidad de generalización del modelo.
- **Flexibilidad y Robustez:** La capacidad de manejar dependencias secuenciales complejas y la inclusión de mecanismos de regularización avanzados hacen que el modelo sea más robusto frente a diferentes tipos de datos y tareas de NLP.

### 3.8 ImprovedTransformerBlock

**Propósito:** Implementar un bloque Transformer mejorado que combina atención local optimizada, convoluciones deformables, MoE, y una red feedforward con mecanismos de regularización avanzados.

#### 3.8.1 Descripción Algorítmica

1. **Bloque de Atención con Conexión Residual:**

$$\mathbf{X}' = \text{LayerNorm}(\mathbf{X})$$

$$\mathbf{A} = \text{EnhancedLocalAttentionWithGQA}(\mathbf{X}')$$

$$\mathbf{X} = \mathbf{X} + \text{Dropout}(\mathbf{A})$$

2. **Bloque de Convolución con Conexión Residual:**

$$\mathbf{X}' = \text{LayerNorm}(\mathbf{X})$$

$$\mathbf{C} = \text{OptimizedGatedConvolution}(\mathbf{X}')$$

$$\mathbf{X} = \mathbf{X} + \text{Dropout}(\mathbf{C})$$

3. **Bloque MoE con Conexión Residual:**

$$\mathbf{X}' = \text{LayerNorm}(\mathbf{X})$$

$$\mathbf{M}, \mathcal{L}_{\text{MoE}} = \text{MoELayer}(\mathbf{X}')$$

$$\mathbf{X} = \mathbf{X} + \text{Dropout}(\mathbf{M})$$

4. **Bloque Feedforward con Conexión Residual:**

$$\mathbf{F} = \text{LayerNorm}(\mathbf{X})$$

$$\mathbf{F} = \text{Linear}_1(\mathbf{F})$$

$$\mathbf{F} = \text{LayerNorm}(\mathbf{F})$$

$$\mathbf{F} = \text{GELU}(\mathbf{F})$$

$$\mathbf{F} = \text{Dropout}(\mathbf{F})$$

$$\mathbf{F} = \text{Linear}_2(\mathbf{F})$$

$$\mathbf{F} = \text{Dropout}(\mathbf{F})$$

$$\mathbf{F} = \text{Clamp}(\mathbf{F}, -100, 100)$$

$$\mathbf{F} = \mathbf{F} + \mathbf{X}$$

$$\mathbf{F} = \text{nan\_to\_num}(\mathbf{F})$$

$$\mathbf{X} = \mathbf{F}$$

5. **Regularización y Clipping de Gradientes:**

$$\|\nabla\theta\|_2 \leq 1.0$$

donde  $\theta$  son los parámetros del bloque.



### 3.8.2 Fundamentos Matemáticos

- **Atención Multi-Cabeza Mejorada:** La atención multi-cabeza mejorada combina mecanismos de atención local optimizada con MoE, aumentando la capacidad de modelado y eficiencia del bloque Transformer. Matemáticamente, se integra la atención multi-cabeza dentro del bloque Transformer siguiendo la arquitectura pre-LN:

$$\text{MultiHeadAttention}(\mathbf{X}) = \text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V})\mathbf{W}_O$$

donde  $\mathbf{W}_O$  es la matriz de pesos de salida.

- **Convoluciones Deformables Optimizada:** Las convoluciones deformables optimizadas permiten una mayor flexibilidad en la captura de patrones locales, adaptándose dinámicamente a la entrada. Esto se formaliza mediante la aplicación de desplazamientos aprendidos en las posiciones de muestreo de la convolución.
- **Mixture of Experts Integrado:** La integración de MoE dentro del bloque Transformer permite que diferentes expertos se especialicen en diferentes aspectos de la entrada, mejorando la capacidad del modelo para manejar heterogeneidad y diversidades en los datos de entrada.
- **Red Feedforward con Clipping y Normalización:** La red feedforward modificada incorpora capas de normalización y clipping para mantener la estabilidad numérica y prevenir la explosión de gradientes. Matemáticamente, esto se representa mediante:

$$\mathbf{F} = \text{Clamp}(\mathbf{F}, -100, 100)$$

asegurando que las activaciones no excedan un rango predefinido.

- **Clipping de Gradientes:** El clipping de gradientes limita la norma de los gradientes a un valor máximo, evitando que los gradientes excesivamente grandes afecten negativamente el proceso de actualización de pesos:

$$\|\nabla\theta\|_2 \leq 1.0$$

### 3.8.3 Análisis Detallado

El `ImprovedTransformerBlock` es una fusión avanzada de múltiples componentes que optimizan tanto la capacidad de modelado como la eficiencia computacional del modelo. La combinación de atención local mejorada, convoluciones deformables optimizadas y MoE permite al bloque Transformer capturar relaciones complejas y variadas en las secuencias de entrada, adaptándose dinámicamente a diferentes contextos y patrones.

La inclusión de una red feedforward robusta con mecanismos de clipping y normalización asegura que las activaciones permanezcan dentro de rangos estables, mejorando la estabilidad del entrenamiento y facilitando la convergencia del modelo.

### 3.8.4 Implicaciones para el Modelo

- **Aumento de la Capacidad de Modelado:** La integración de múltiples mecanismos de atención y convolución permite que el bloque Transformer capture una gama más amplia de relaciones semánticas y patrones en los datos, mejorando el rendimiento en tareas complejas.
- **Eficiencia Computacional Mejorada:** La optimización de los mecanismos de atención y convolución, junto con el uso de MoE, reduce la carga computacional y el consumo de memoria, permitiendo entrenar modelos más profundos y con mayor capacidad sin incurrir en costos computacionales prohibitivos.
- **Estabilidad Numérica:** Los mecanismos de clipping y normalización aseguran que las activaciones y los gradientes permanezcan dentro de rangos manejables, evitando problemas comunes como el desvanecimiento o explosión de gradientes, y mejorando la robustez del modelo durante el entrenamiento.
- **Generalización y Robustez:** La combinación de técnicas avanzadas de regularización y la capacidad de capturar relaciones complejas contribuyen a una mejor generalización del modelo a datos no vistos, aumentando su aplicabilidad en diferentes dominios y tareas de NLP.

### 3.9 BidirectionalEncoder

**Propósito:** Construir un encoder bidireccional que procesa secuencias de entrada en ambas direcciones, capturando contextos pasados y futuros simultáneamente.

#### 3.9.1 Descripción Algorítmica

1. **Embeddings Líquidos:**

$$(\mathbf{X}, \mathcal{L}_{\text{recon}}) = \text{EnhancedLiquidEmbedding}(\mathbf{x})$$

Aplicación de dropout:

$$\mathbf{X} = \text{Dropout}(\mathbf{X})$$

2. **Capas Transformer Mejoradas:**

$$\mathcal{L}_{\text{total\_entropy}} = 0$$

$$\forall l \in \{1, \dots, K\} :$$

$$(\mathbf{X}, \mathcal{L}_l) = \text{ImprovedTransformerBlock}(\mathbf{X})$$

$$\mathcal{L}_{\text{total\_entropy}} += \mathcal{L}_l$$

3. **Normalización Final:**

$$\mathbf{X} = \text{LayerNorm}(\mathbf{X})$$

#### 3.9.2 Fundamentos Matemáticos

- **Encoder Bidireccional:** Un encoder bidireccional procesa la secuencia de entrada en ambas direcciones, hacia adelante y hacia atrás, permitiendo que cada token capture contexto tanto del pasado como del futuro. Matemáticamente, esto se logra concatenando las representaciones obtenidas de ambas direcciones:

$$\mathbf{h}_t = [\mathbf{h}_t^{\text{forward}}; \mathbf{h}_t^{\text{backward}}]$$

- **Integración de Embeddings Líquidos y Transformer Mejorado:** La combinación de embeddings líquidos con bloques Transformer mejorados permite al modelo capturar representaciones ricas y contextualmente informadas. Los embeddings líquidos proporcionan una base robusta que se enriquece a través de múltiples capas de atención y convolución, optimizadas para manejar secuencias largas y complejas.

#### 3.9.3 Análisis Detallado

El encoder bidireccional es fundamental para capturar el contexto completo de cada token en la secuencia. Al procesar la secuencia en ambas direcciones, el encoder puede entender mejor las relaciones semánticas y las dependencias a largo plazo entre los tokens, lo que es crucial para tareas como la traducción automática, el resumen de textos y la comprensión de lenguaje natural.

La acumulación de pérdidas de entropía durante el procesamiento de las capas Transformer mejoradas contribuye a la regularización del modelo, promoviendo una distribución equilibrada en el uso de los expertos y evitando el sobreajuste.

#### 3.9.4 Implicaciones para el Modelo

- **Captura Completa del Contexto:** La bidireccionalidad permite que el encoder capture información contextual tanto del pasado como del futuro, mejorando la comprensión semántica de la secuencia de entrada.
- **Mejora en la Representación de Contexto:** La integración de embeddings líquidos y bloques Transformer mejorados facilita la creación de representaciones contextuales profundas, enriqueciendo la capacidad del modelo para entender y generar lenguaje natural.
- **Regularización y Estabilidad:** La acumulación de pérdidas de entropía y el uso de LayerNorm contribuyen a una mayor estabilidad durante el entrenamiento y una mejor generalización del modelo.

- **Eficiencia en el Manejo de Secuencias Largas:** Las técnicas avanzadas de compresión y atención local optimizada permiten que el encoder maneje secuencias de gran longitud de manera eficiente, manteniendo la calidad de la representación sin incurrir en costos computacionales prohibitivos.

### 3.10 LiquidFoundationModelOptimized

**Propósito:** Integrar todos los componentes anteriores en un modelo completo que puede ser utilizado para tareas de generación de lenguaje, traducción, o cualquier otra aplicación de NLP que requiera comprensión y generación de secuencias.

#### 3.10.1 Descripción Algorítmica

1. **Encoder:**

$$(\mathbf{H}_e, \mathcal{L}_{\text{recon\_enc}}, \mathcal{L}_{\text{entropy\_enc}}) = \text{BidirectionalEncoder}(\mathbf{E}_{\text{in}})$$

2. **Decoder:**

$$(\mathbf{H}_d, \mathcal{L}_{\text{recon\_dec}}) = \text{EnhancedLiquidEmbedding}(\mathbf{D}_{\text{in}})$$

$$\mathbf{H}_d = \text{Dropout}(\mathbf{H}_d)$$

Inicialización de la memoria LSTM:

$$\mathbf{h}_0, \mathbf{c}_0 = \text{EnhancedLSTM.init\_hidden}(B)$$

$$\mathcal{L}_{\text{total\_entropy\_dec}} = 0$$

$$\forall l \in \{1, \dots, K\} :$$

$$(\mathbf{H}_d, \mathcal{L}_l) = \text{ImprovedTransformerBlock}(\mathbf{H}_d)$$

$$\mathcal{L}_{\text{total\_entropy\_dec}} += \mathcal{L}_l$$

$$(\mathbf{H}_d, \mathbf{h}, \mathbf{c}) = \text{EnhancedLSTM}(\mathbf{H}_d, (\mathbf{h}, \mathbf{c}))$$

3. **Normalización Final y Proyección a Logits:**

$$\mathbf{H}_d = \text{LayerNorm}(\mathbf{H}_d)$$

$$\mathbf{Y} = \mathbf{H}_d \mathbf{W}_{\text{out}} + \mathbf{b}_{\text{out}}$$

4. **Generación de Secuencia:** Método de generación iterativa que, dado un token inicial, genera tokens sucesivos hasta alcanzar una longitud máxima o un token de fin de secuencia (EOS).

#### 3.10.2 Fundamentos Matemáticos

- **Arquitectura Encoder-Decoder:** Basada en la arquitectura Transformer de Vaswani et al. (2017), la arquitectura encoder-decoder permite que el encoder capture la representación contextual de la entrada y el decoder genere la salida basada en esta representación. Matemáticamente, esto se formaliza como:

$$\mathbf{Y} = \text{Decoder}(\text{Encoder}(\mathbf{X}))$$

- **Memoria Externa (EnhancedLSTM):** La memoria externa basada en una versión mejorada de LSTM facilita el manejo de dependencias secuenciales a largo plazo, integrando estados ocultos que capturan información contextual relevante para la generación de secuencias.
- **Generación Iterativa:** El proceso de generación iterativa se define como una secuencia de pasos donde cada token generado se utiliza como entrada para generar el siguiente token. Matemáticamente, esto se puede representar como:

$$\mathbf{y}_t = \text{Decoder}(\mathbf{y}_{<t}, \mathbf{H}_e)$$

donde  $\mathbf{y}_{<t}$  son los tokens generados hasta el momento  $t$ , y  $\mathbf{H}_e$  es la representación contextual proporcionada por el encoder.

- **Regularización y Estabilidad:** Técnicas como dropout y clipping de gradientes son esenciales para mantener la estabilidad durante el entrenamiento, especialmente en modelos profundos y complejos. Matemáticamente, esto se implementa como:

$$\text{Dropout}(\mathbf{x}) = \mathbf{x} \otimes \mathbf{m}$$

$$\|\nabla \theta\|_2 \leq 1.0$$

donde  $\mathbf{m} \sim \text{Bernoulli}(p)$  es una máscara de dropout y  $\theta$  son los parámetros del modelo.

### 3.10.3 Análisis Detallado

El modelo `LiquidFoundationModelOptimized` integra de manera cohesiva todos los componentes analizados anteriormente, creando una arquitectura robusta y eficiente para tareas avanzadas de NLP. La estructura encoder-decoder permite una separación clara de responsabilidades: el encoder captura las representaciones contextuales de la entrada, mientras que el decoder genera las salidas basadas en estas representaciones.

La memoria externa basada en `EnhancedLSTM` introduce una capacidad adicional para manejar dependencias a largo plazo, mejorando la coherencia y relevancia de las secuencias generadas. La inclusión de mecanismos avanzados de atención y convolución optimizados, junto con la capa MoE y **FrequencyHyenaAttention**, aumenta significativamente la capacidad del modelo para manejar diversas tareas y contextos de manera eficiente.

### 3.10.4 Implicaciones para el Modelo

- **Versatilidad en Aplicaciones de NLP:** La arquitectura completa permite que el modelo sea utilizado en una amplia gama de aplicaciones de NLP, incluyendo generación de lenguaje, traducción automática, resumen de textos, y más.
- **Eficiencia y Escalabilidad:** La combinación de técnicas de compresión dinámica, atención optimizada y MoE asegura que el modelo pueda escalar a grandes volúmenes de datos y secuencias largas sin incurrir en costos computacionales prohibitivos.
- **Estabilidad y Robustez durante el Entrenamiento:** La implementación de mecanismos de regularización avanzados y técnicas de estabilidad numérica garantizan que el modelo pueda entrenarse de manera eficiente y robusta, minimizando problemas como el desvanecimiento o explosión de gradientes.
- **Capacidad de Adaptación Dinámica:** La inclusión de mecanismos como la selección dinámica de expertos, convoluciones deformables y **FrequencyHyenaAttention** permite que el modelo se adapte de manera flexible a diferentes tipos de datos y contextos, mejorando su desempeño en tareas variadas.
- **Mejora Continua mediante Monitoreo:** La implementación de hooks para monitoreo de activaciones y gradientes asegura que el modelo pueda ser supervisado y ajustado en tiempo real, facilitando la detección y corrección de problemas durante el entrenamiento.

## 4 Consideraciones de Estabilidad Numérica y Eficiencia

El modelo incorpora múltiples mecanismos para garantizar la estabilidad numérica y la eficiencia durante el entrenamiento y la inferencia:

### 1. Verificaciones de Finitud:

- Después de operaciones críticas como atención, convoluciones y capas de gating, se verifica que todas las activaciones sean finitas (`torch.isfinite`). Si se detectan valores no finitos, se lanza una excepción o se reemplazan los valores con ceros, evitando la propagación de errores numéricos.

### 2. Uso de Mixed Precision (autocast):

- Utiliza el mezclado de precisión para aprovechar la eficiencia de los cálculos en punto flotante de menor precisión (FP16), reduciendo el uso de memoria y acelerando el entrenamiento sin sacrificar significativamente la precisión:

$$\mathbf{y} = \text{autocast}(\mathbf{x})$$

### 3. Checkpointing:

- Implementa checkpointing para manejar grandes modelos de manera eficiente en términos de memoria, permitiendo el entrenamiento de secuencias largas o modelos profundos sin exceder la capacidad de la GPU:

$$\mathbf{y} = \text{checkpoint}(\text{func}, \mathbf{x})$$

donde `func` es una función que se ejecuta bajo la estrategia de checkpointing.

### 4. Inicialización de Pesos:

- Aplica esquemas de inicialización avanzados como Xavier y Kaiming para asegurar que los gradientes fluyan adecuadamente durante el entrenamiento, evitando problemas como el desvanecimiento o explosión de gradientes:

$$\mathbf{W} \sim \mathcal{N}\left(0, \frac{2}{n_{\text{in}} + n_{\text{out}}}\right) \quad (\text{Xavier})$$

$$\mathbf{W} \sim \mathcal{N}\left(0, \frac{2}{n_{\text{in}}}\right) \quad (\text{Kaiming})$$

### 5. Dropout y Clipping de Gradientes:

- **Dropout:**

$$\mathbf{x}_{\text{drop}} = \mathbf{x} \otimes \mathbf{m}$$

donde  $\mathbf{m} \sim \text{Bernoulli}(p)$ .

- **Clipping de Gradientes:**

$$\mathbf{g} = \frac{\mathbf{g}}{\max\left(1, \frac{\|\mathbf{g}\|_2}{\text{clip\_value}}\right)}$$

donde  $\mathbf{g}$  es el gradiente.

### 6. Normalización de Capa (LayerNorm):

$$\text{LayerNorm}(\mathbf{x}) = \frac{\mathbf{x} - \mu}{\sqrt{\sigma^2 + \epsilon}} \gamma + \beta$$

#### 4.0.1 Análisis Detallado

- **Verificaciones de Finitud:** La verificación de finitud en cada etapa crítica del modelo asegura que las activaciones y los gradientes no contengan valores infinitos o NaNs, lo que podría interrumpir el proceso de entrenamiento o llevar a un rendimiento inconsistente.
- **Mixed Precision Training:** El uso de precisión mixta (`autocast`) permite una reducción significativa en el uso de memoria y un aumento en la velocidad de cómputo sin comprometer la precisión del modelo. Matemáticamente, esto implica que ciertas operaciones se realizan en FP16, mientras que otras críticas se mantienen en FP32 para preservar la exactitud.
- **Checkpointing:** El checkpointing optimiza el uso de memoria al almacenar solo ciertos puntos de activación durante el forward pass y recalculando las activaciones durante el backward pass según sea necesario. Esto se formaliza como:

$$\mathbf{y} = \text{checkpoint}(\text{func}, \mathbf{x})$$

donde `func` es una función que encapsula una parte del modelo.

- **Inicialización de Pesos:** La inicialización adecuada de los pesos es crucial para asegurar un flujo de gradientes adecuado durante el entrenamiento. Las inicializaciones Xavier y Kaiming están diseñadas para mantener la varianza de las activaciones y gradientes constante a través de las capas, lo que facilita la convergencia del modelo.
- **Dropout y Clipping de Gradientes:** El dropout actúa como una forma de regularización que previene el sobreajuste, mientras que el clipping de gradientes evita que los gradientes excesivamente grandes afecten negativamente el proceso de actualización de pesos. Matemáticamente:

$$\mathbf{x}_{\text{drop}} = \mathbf{x} \otimes \mathbf{m}, \quad \mathbf{g} = \frac{\mathbf{g}}{\max\left(1, \frac{\|\mathbf{g}\|_2}{\text{clip\_value}}\right)}$$

- **LayerNorm:** La normalización de capa estabiliza las activaciones y acelera la convergencia del modelo al reducir la covariate shift interna. Esto se logra mediante la normalización de cada muestra de manera independiente, manteniendo la distribución de las activaciones dentro de un rango manejable.

### 4.1 Análisis de Eficiencia y Estabilidad

Las mejoras implementadas en el modelo `LiquidFoundationModelOptimized` no solo aumentan su capacidad de modelado, sino que también garantizan una eficiencia y estabilidad óptimas durante su entrenamiento y despliegue. La combinación de técnicas de regularización avanzadas, optimización de atención y convoluciones, junto con mecanismos de control de gradientes, asegura que el modelo pueda entrenarse de manera robusta y escalable, manteniendo un rendimiento consistente incluso en entornos de producción de gran escala.

## 5 Conclusión

El `LiquidFoundationModelOptimized` representa una arquitectura integral y altamente avanzada que combina múltiples técnicas modernas de aprendizaje profundo para abordar tareas complejas de procesamiento de lenguaje natural. Su diseño incorpora mecanismos para capturar relaciones contextuales profundas, manejar secuencias largas, y adaptarse dinámicamente a la complejidad de los datos de entrada. Además, se presta especial atención a la estabilidad numérica y la eficiencia computacional, haciéndolo apto para aplicaciones a gran escala en entornos de producción.

### 5.1 Puntos Clave del Modelo

- **Embeddings Líquidos:** Combina embeddings de tokens y posición con convoluciones y técnicas de compresión en el dominio de la frecuencia, permitiendo manejar secuencias de gran longitud de manera eficiente.

- **Atención Local Mejorada:** Optimiza la atención multi-cabeza mediante mecanismos de agrupamiento y ventanas deslizantes, reduciendo la complejidad computacional sin sacrificar la capacidad de capturar relaciones a largo plazo.
- **Mixture of Experts (MoE):** Introduce un enrutamiento adaptativo que dirige las entradas a expertos especializados, mejorando la capacidad del modelo para manejar diversidad en las tareas y contextos.
- **Convoluciones Deformables y Gated Mechanisms:** Proporcionan flexibilidad en la captura de patrones locales y controlan el flujo de información, mejorando la capacidad del modelo para manejar dependencias complejas.
- **FrequencyHyenaAttention:** Integra la arquitectura Hyena optimizada para el dominio de la frecuencia, mejorando la eficiencia y la captura de relaciones a largo plazo en secuencias de entrada.
- **Memoria Externa (EnhancedLSTM):** Facilita el manejo de dependencias secuenciales a largo plazo, mejorando la coherencia y relevancia de las secuencias generadas.
- **Regularización y Estabilidad Numérica:** Incorpora múltiples mecanismos para mantener la estabilidad durante el entrenamiento, prevenir el sobreajuste y optimizar el flujo de gradientes.

## 5.2 Resumen de Implicaciones de las Mejoras

- **Capacidad de Modelado Aumentada:** La integración de MoE, convoluciones deformables y **FrequencyHyenaAttention** mejora significativamente la capacidad del modelo para capturar y procesar información compleja y variada.
- **Eficiencia Computacional:** Las optimizaciones en los mecanismos de atención, convolución y compresión dinámica permiten un manejo más eficiente de recursos, reduciendo el tiempo de entrenamiento y el consumo de memoria.
- **Estabilidad y Robustez:** Los mecanismos de verificación de finitud, junto con técnicas avanzadas de normalización y regularización, garantizan una mayor estabilidad durante el entrenamiento, prevenir el sobreajuste y optimizar el flujo de gradientes.
- **Escalabilidad:** La arquitectura modular y escalable del modelo permite su adaptación a diferentes tamaños de vocabulario, longitudes de secuencia y capacidades de hardware, facilitando su implementación en una variedad de entornos y aplicaciones.
- **Generalización Mejorada:** Las técnicas de regularización y la capacidad de manejar dependencias a largo plazo contribuyen a una mejor generalización del modelo, permitiéndole desempeñarse eficazmente en datos no vistos y en diferentes contextos de aplicación.

En resumen, el modelo `LiquidFoundationModelOptimized` representa una convergencia de múltiples avances en arquitectura de redes neuronales, optimización computacional y técnicas de regularización. Estas mejoras no solo aumentan la capacidad y eficiencia del modelo, sino que también aseguran su estabilidad y robustez, posicionándolo como una solución robusta y eficiente para desafíos avanzados en procesamiento de lenguaje natural.