# Transformers y atención como bloque universal

Self-attention (Q,K,V), multi-head, conexiones residuales + LayerNorm, positional encodings y variantes encoder/decoder.

# Temario

- Self-attention: consultas (Q), claves (K) y valores (V); atención token<>token y complejidad $O(T^2)$.
- Multi-head attention, normalización (LayerNorm) y conexiones residuales (skip).
- Arquitectura encoder/decoder y variantes (encoder-only, decoder-only para LLM).
- Positional encodings y variantes (seno/coseno, RoPE, ALiBi, etc.).

$$Attention\ (Q, K, V)\ =\ softmax\left(\frac{Q \cdot k^T}{\sqrt{d_k}}\right) \cdot V$$

$$y = \frac{e^{x_i}}{\sum_{i=1}^{i=n} e^{x_i}}$$

$$y = \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix} = softmax(x) = f\left( \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \right) = \begin{bmatrix} \dfrac{e^{x_1}}{e^{x_1} + e^{x_2} + e^{x_3}} \\ \dfrac{e^{x_2}}{e^{x_1} + e^{x_2} + e^{x_3}} \\ \dfrac{e^{x_3}}{e^{x_1} + e^{x_2} + e^{x_3}} \end{bmatrix}$$

$$python: y = \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix} = softmax(x) = f\left( \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \right) = \frac{[x_1 \quad x_2 \quad x_3]}{e^{x_1} + e^{x_2} + e^{x_3}} = \frac{e^x}{sum(e^x)}$$

$$Q = X \cdot W^Q, \; K = X \cdot W^V, \; V = X \cdot W^V$$

$$time = \mathcal{O}(T^2 + d) \quad space = \mathcal{O}(T^2 + Td)$$

# Codificación posicional seno/coseno



$$PE_{(pos,2i)} = \sin\left(pos/1000^{2i/d}\right)$$

$$PE_{(pos,2i+1)} = \cos\left(pos/1000^{2i/d}\right)$$

## RoPE, ALiBi y embeddings posicionales aprendidos

- Objetivo: inyectar información de orden sin cambiar el bloque de atención (mismo Q,K,V).
  - RoPE (Rotary): rota Q y K según posición -> atención depende de distancias relativas; muy común en LLM.
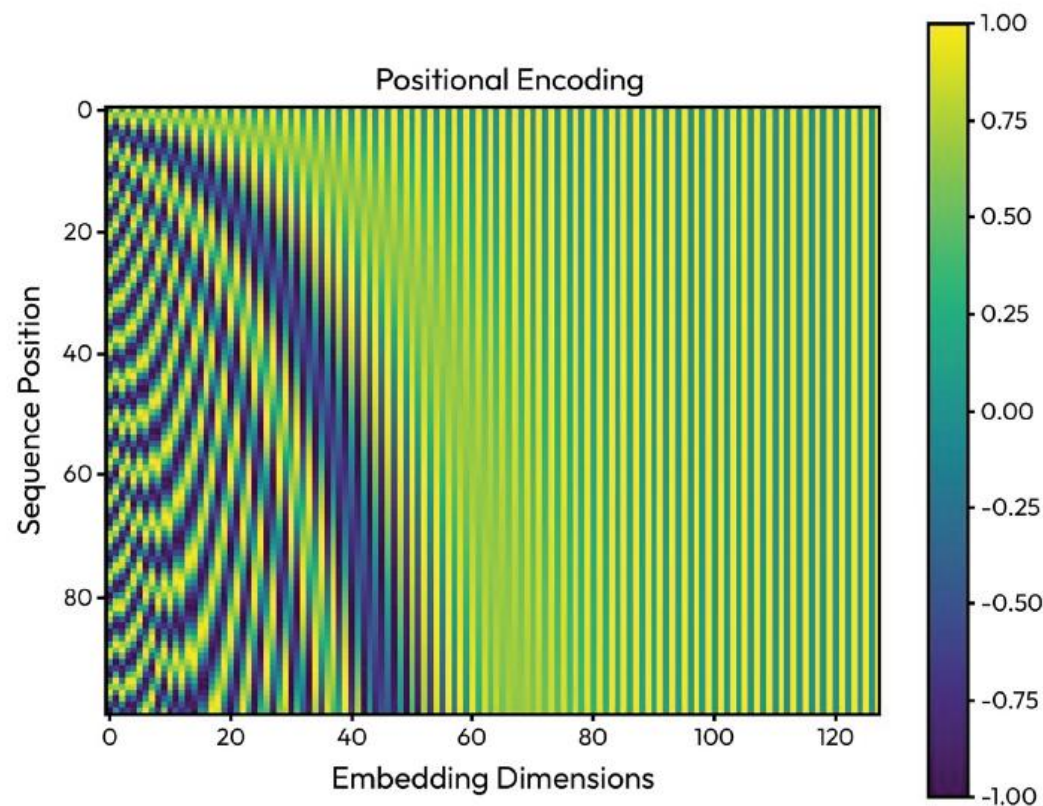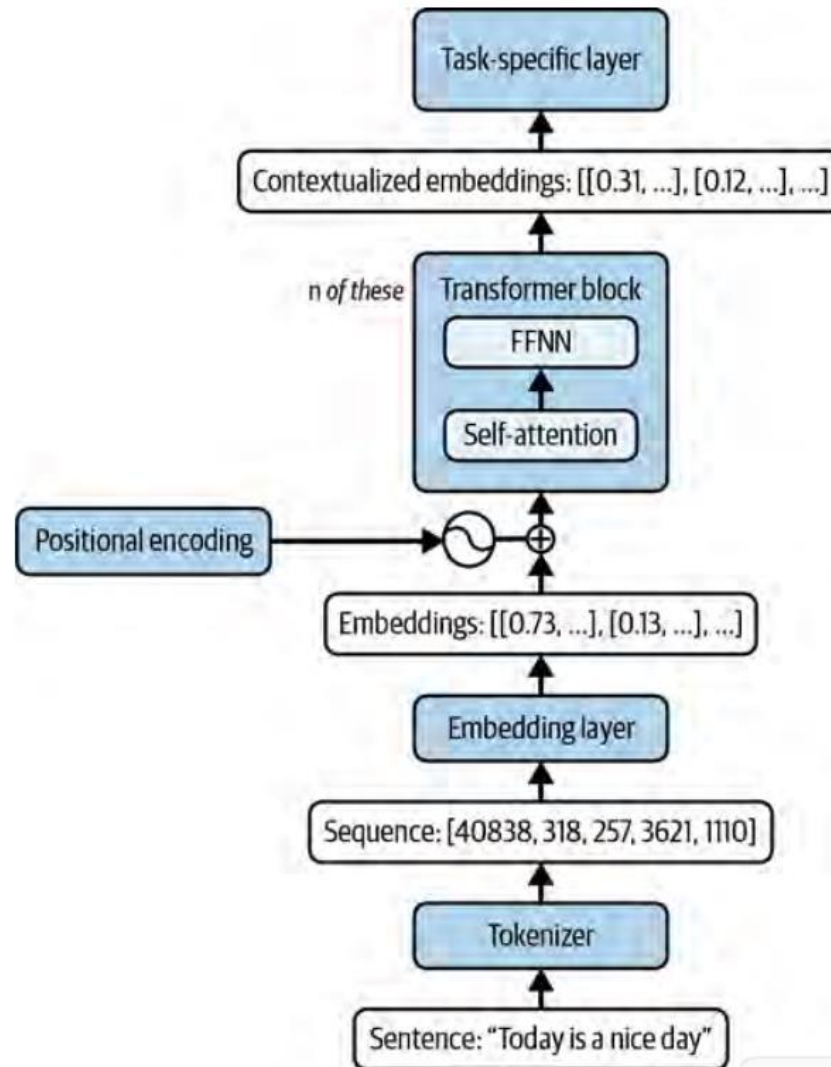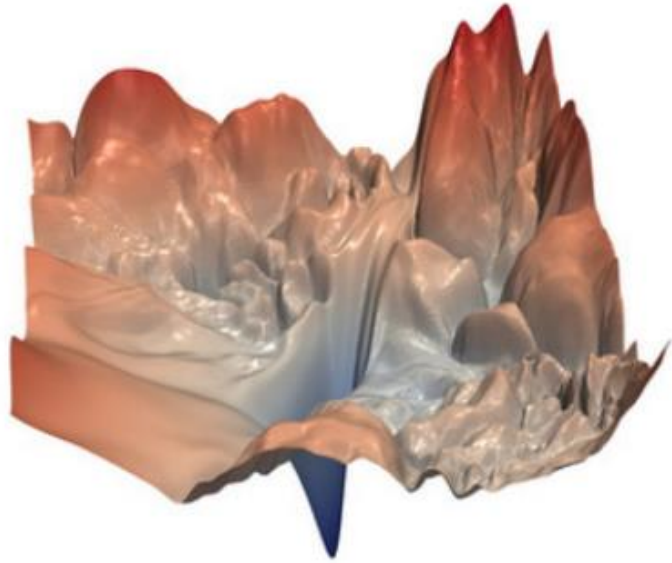  - ALiBi: sesgo lineal en logits de atención por distancia -> favorece relaciones locales y escala a contextos largos.
  - Embeddings posicionales aprendidos: vector por posición (absoluto) -> simple/efectivo, pero extrapola peor fuera del rango entrenado.

  Regla práctica: si esperas extrapolar a contextos más largos, RoPE/ALiBi suelen ser más robustos que embeddings absolutos aprendidos.
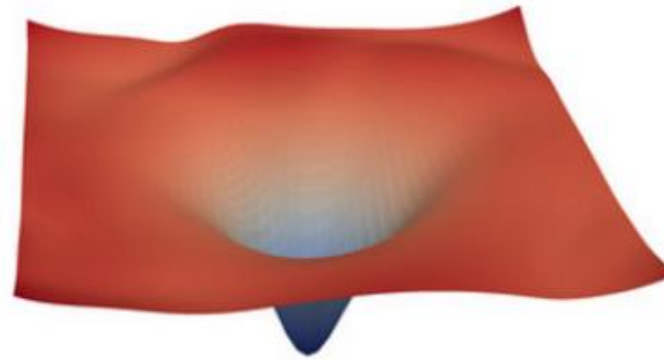
Loss landscape without residual

Loss landscape with residual

# LayerNorm: normalizar y re-escalar

$$\mu = \frac{1}{d}\sum_{i=1}^{d} x_i \qquad \sigma = \sqrt{\frac{1}{d}\sum_{i=1}^{d}(x_i - \mu)^2}$$

$$\hat{x} = \frac{(x - \mu)}{\sigma}$$
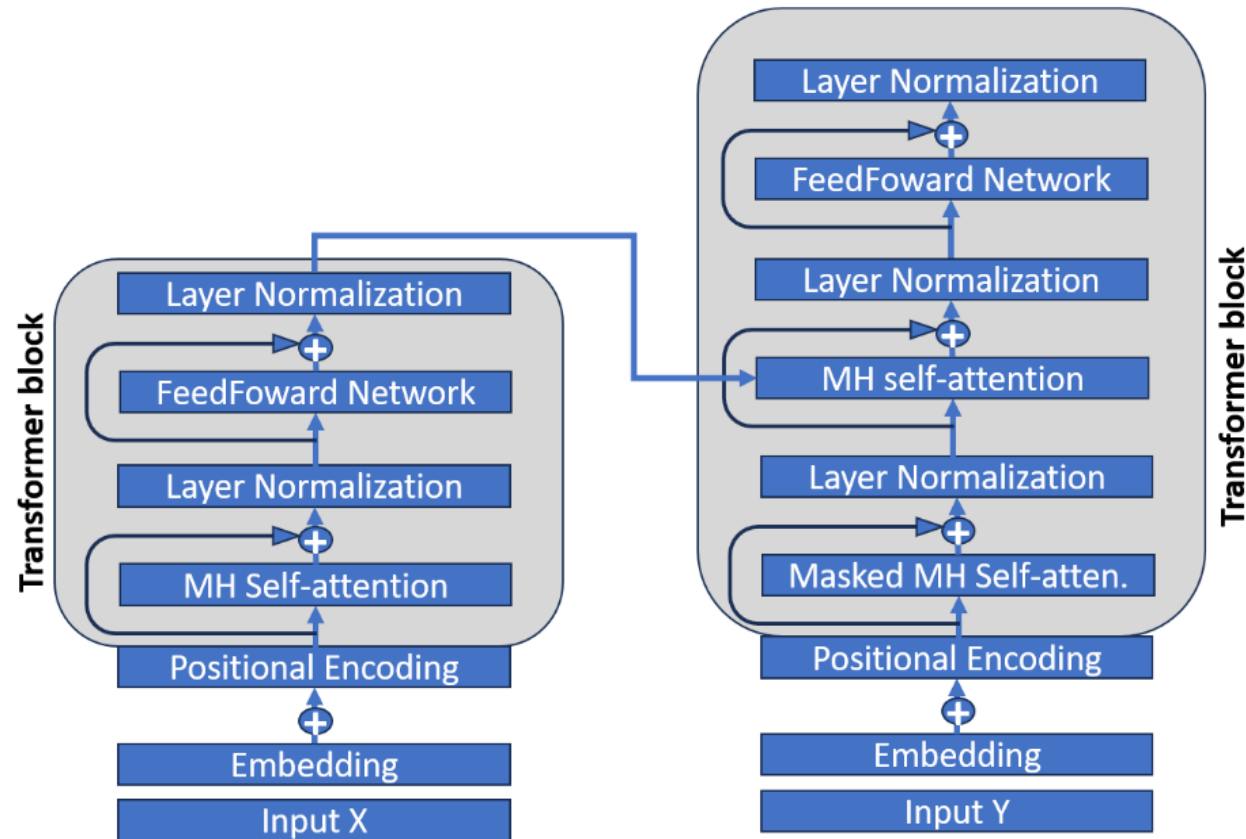
$$LayerNormalization = \gamma\hat{x} + \beta$$

# Pre-Norm en Transformers modernos

$$H = LayerNorm(X + MultiHeadSelfAttention(X))$$
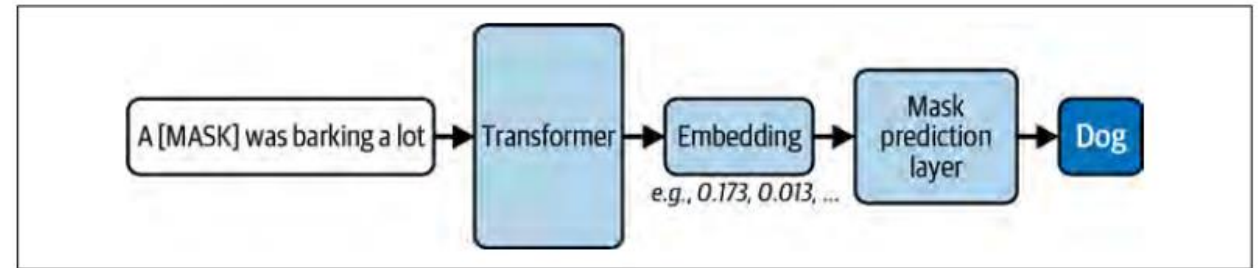
$$H = LayerNorm(H + FFN(H))$$

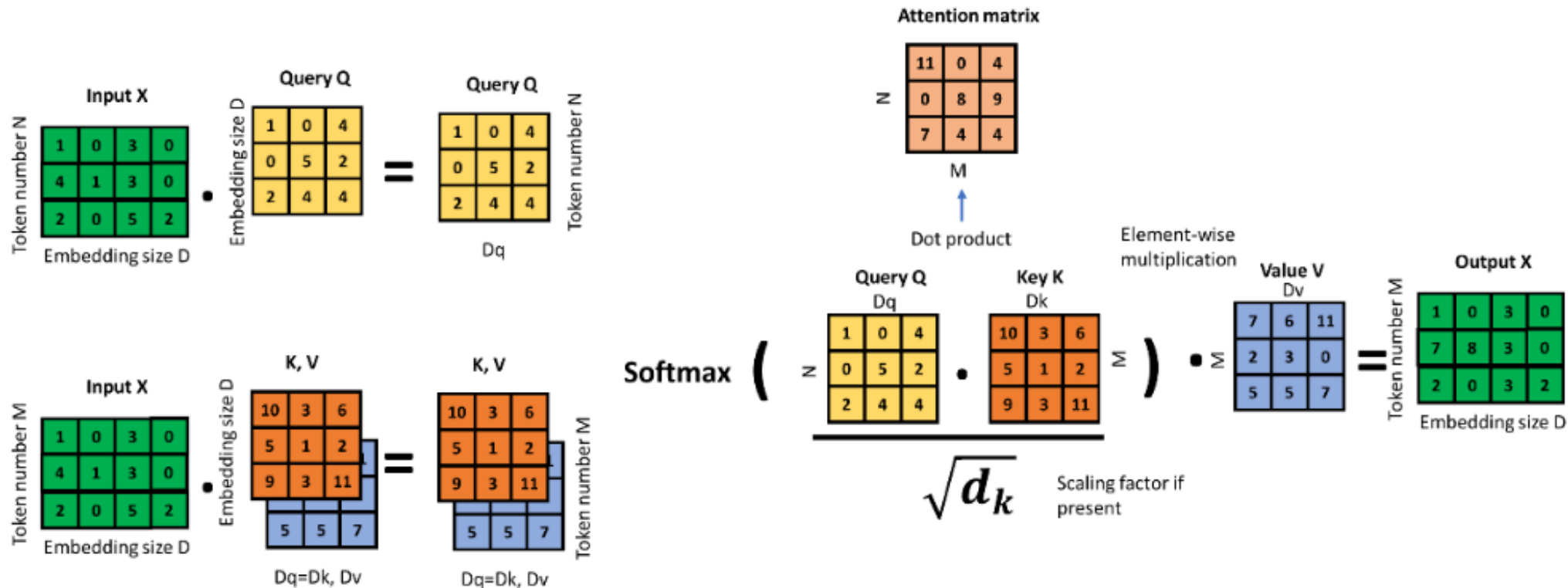# Encoder vs Decoder (self, masked, cross-attention)

# Variantes: encoder-only (MLM) vs decoder-only (LLM)

- Encoder-only: predicción de [MASK] (BERT).
- Decoder-only: modelo autoregresivo (next-token) con máscara causal.
- Encoder/decoder: cross-attention para condicionamiento (traducción, T5, etc.).

# Máscara causal (no mirar el futuro)

# Modelo de lenguaje autoregresivo

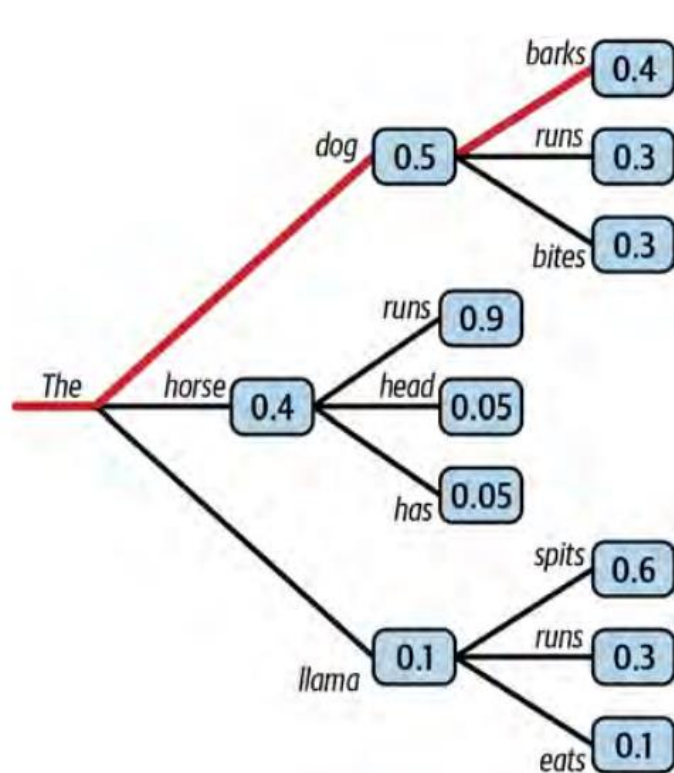$$P(w|h) = P(w_n|w_{1:n-1}) = \prod_{i=1}^{n} P(w_i|w_{1:i-1})$$

$$L_{CE} = -\sum_{w \in V} y_t[w] \log \hat{y}_t[w]$$
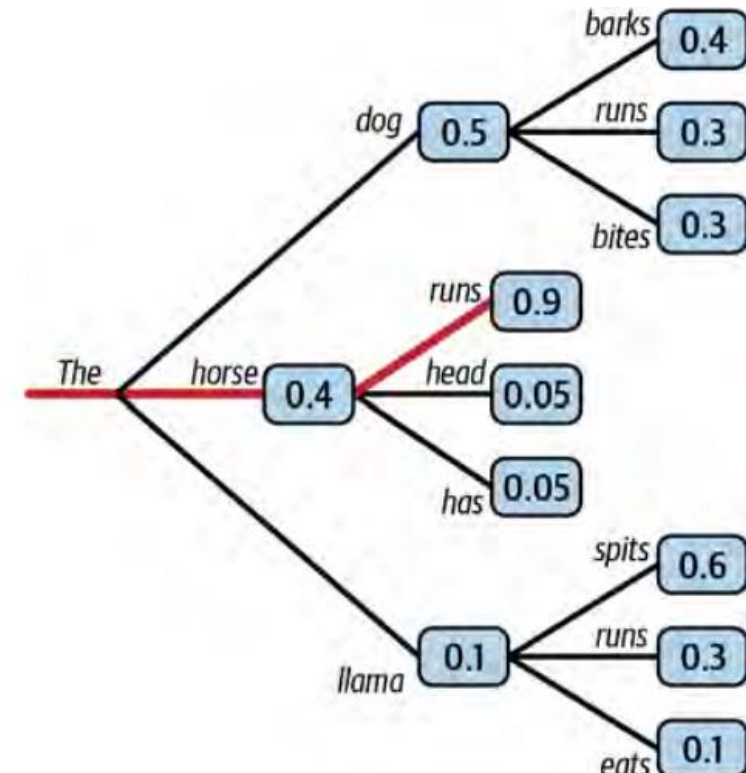
# Decoding: temperatura, top-k y top-p

- Temperatura (T): controla entropía (más alto = más diversidad).
- Top-k: restringe a los k tokens más probables.
- Top-p (nucleus): restringe al menor conjunto con probabilidad acumulada ≥ p.

# Decoding: Greedy y beam search



Decodificación greedy

Decodificación beam-search

# Entrenamiento y Optimización

- AdamW (weight decay desacoplado).
- LR: warmup corto (1-5%) -> cosine decay.
- Gradient clipping (norm ≤ 1.0) para evitar spikes.
- Label smoothing ($\varepsilon \approx 0.1$) mejora calibración.
- Mixed precision (BF16/FP16) para acelerar y ahorrar memoria.