

Redes Neuronales para PLN

Parte 2

Jurafsky and Martin 3rd edition. Caps 7,9,11. <https://web.stanford.edu/~jurafsky/slp3/>

FFN o MLP (multi-layer perceptron)

(Multi-layer fully-connected feed forward Network)

Cada **capa** => secuencia de **neuronas artificiales**
Neurona - recibe vector, devuelve escalar
- pesos ajustables, func. activación

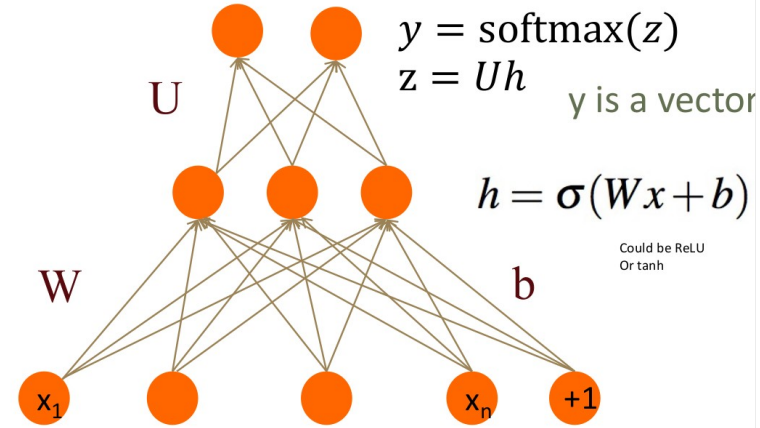
Un MLP es una secuencia de capas apiladas
Entrada → capas ocultas → capa de salida
Sucesivas transformaciones de la entrada

... es **composición de funciones diferenciables**

Regla de la cadena => gradiente de los pesos por capa (**backpropagation**)

Función objetivo (MAE, MSE, Cross-entropy, contrastive, triplet)

Se usan técnicas basadas en el gradiente y mini-batches (adam, rmsprop, sgd)



Algunos conceptos/prácticas generales

Corpus (train, validation/development, test)

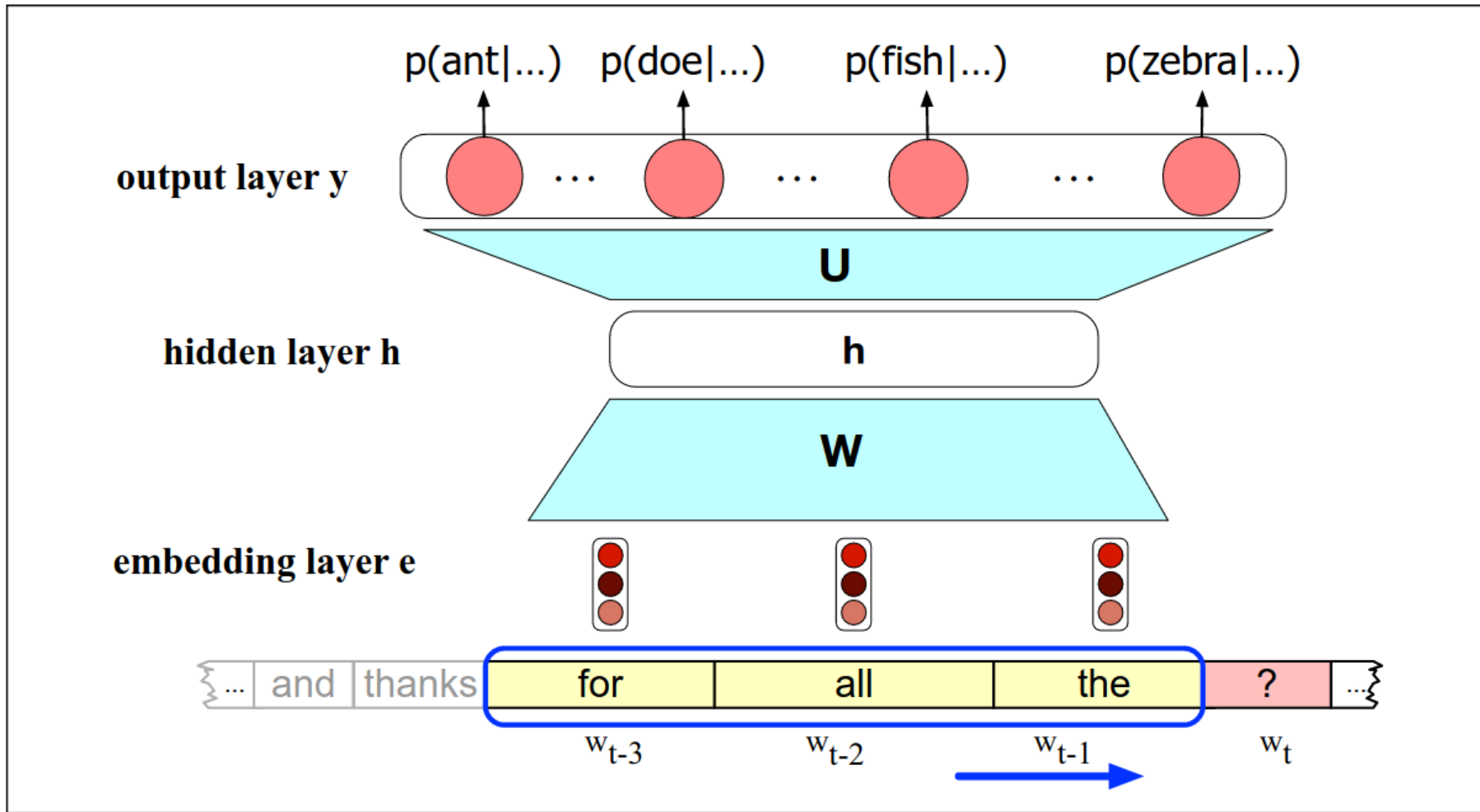
Early Stopping (using validation)

Regularizations (dropout, noise, normalization)

Hyperparameter Search (manual, grid, random)

Embedding layer (look-up/embedding matrix)

Y si la entrada es una secuencia?





Seq2Vec

Ej. Text classification, Análisis de Sentimiento, Detección de Humor, Hatespeech,

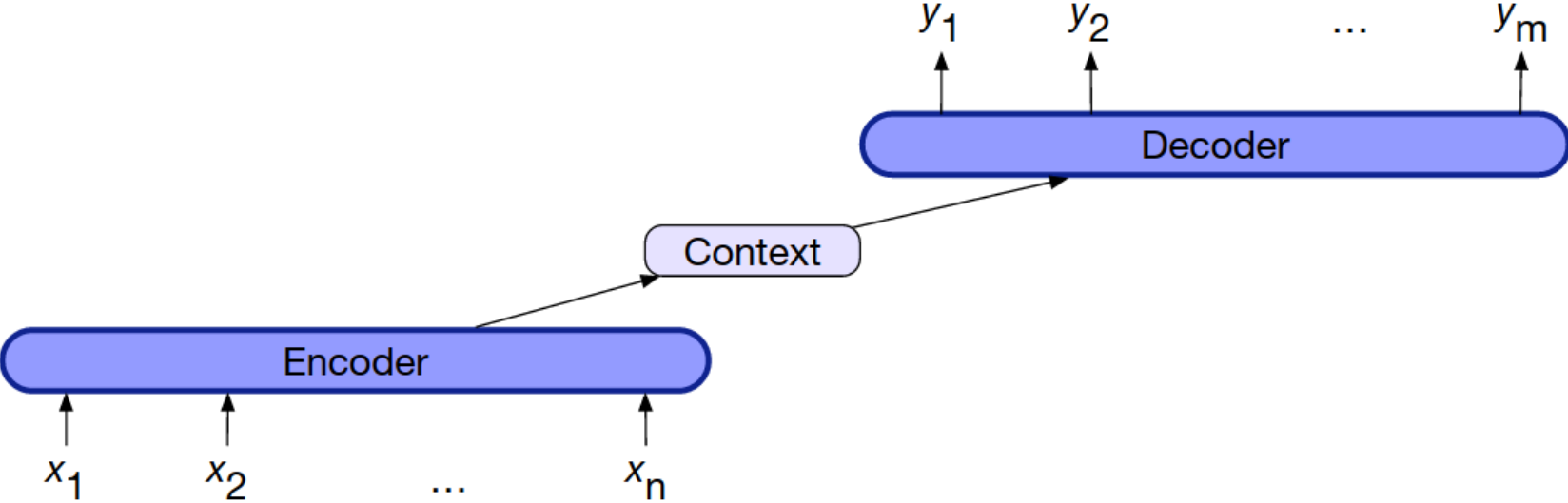
Vec2Seq

Ej. Generación de texto, Image captioning

Seq2Seq (encoder-decoder)

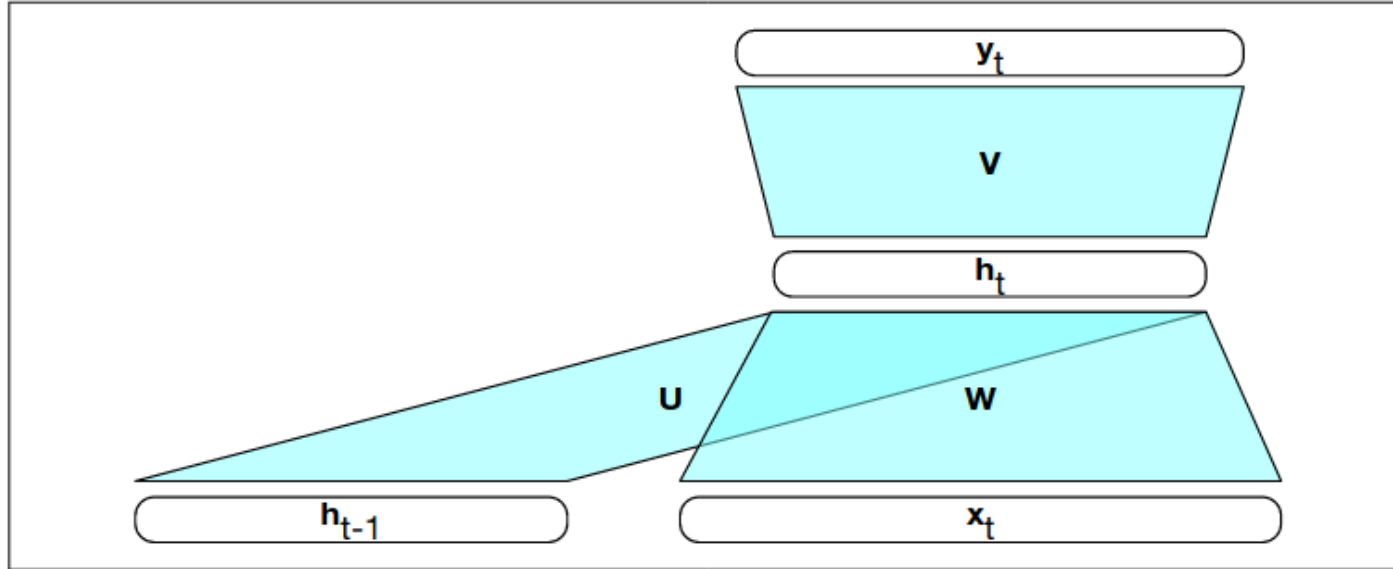
Ej. Traductores automáticos, Resumen abstractivo, Question-answering, Generative chatbots,

Encoder-decoder



Recurrent Neural Networks

Simple Recurrent Network (Elman, 1990)



$$\mathbf{h}_t = g(\mathbf{U}\mathbf{h}_{t-1} + \mathbf{W}\mathbf{x}_t)$$
$$\mathbf{y}_t = f(\mathbf{V}\mathbf{h}_t)$$

Simple Recurrent Network (Elman, 1990)

function FORWARDRNN(\mathbf{x} , *network*) **returns** output sequence \mathbf{y}

$\mathbf{h}^0 \leftarrow 0$

for $i \leftarrow 1$ to LENGTH(\mathbf{x}) **do**

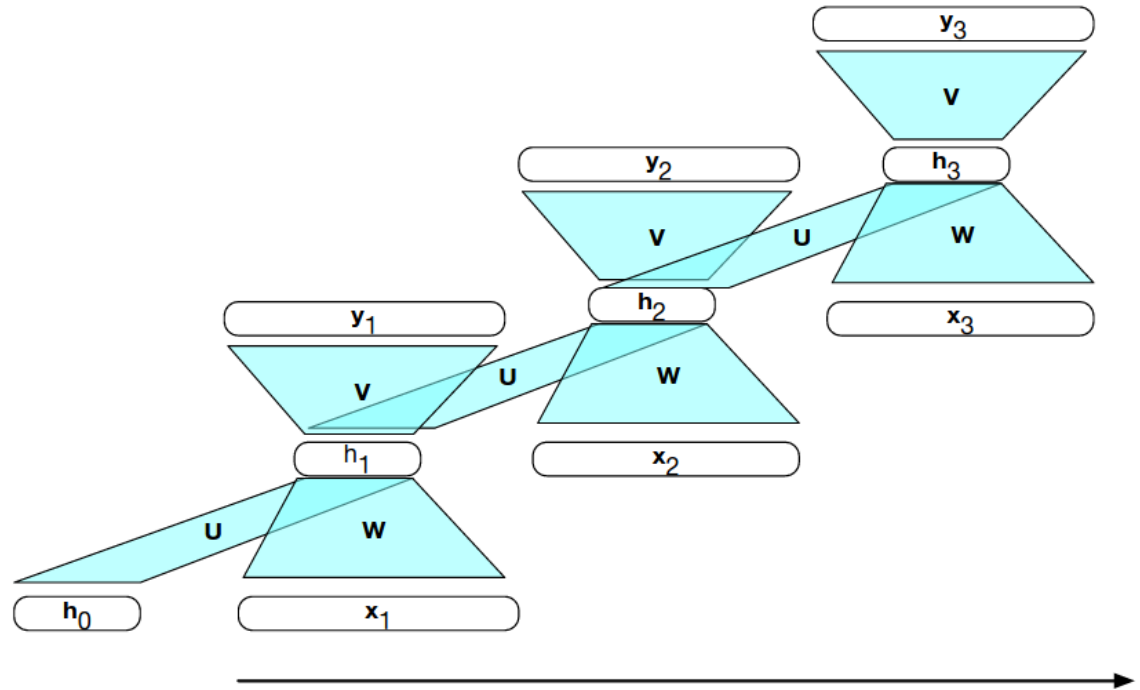
$\mathbf{h}_i \leftarrow g(\mathbf{U}\mathbf{h}_{i-1} + \mathbf{W}\mathbf{x}_i)$

$\mathbf{y}_i \leftarrow f(\mathbf{V}\mathbf{h}_i)$

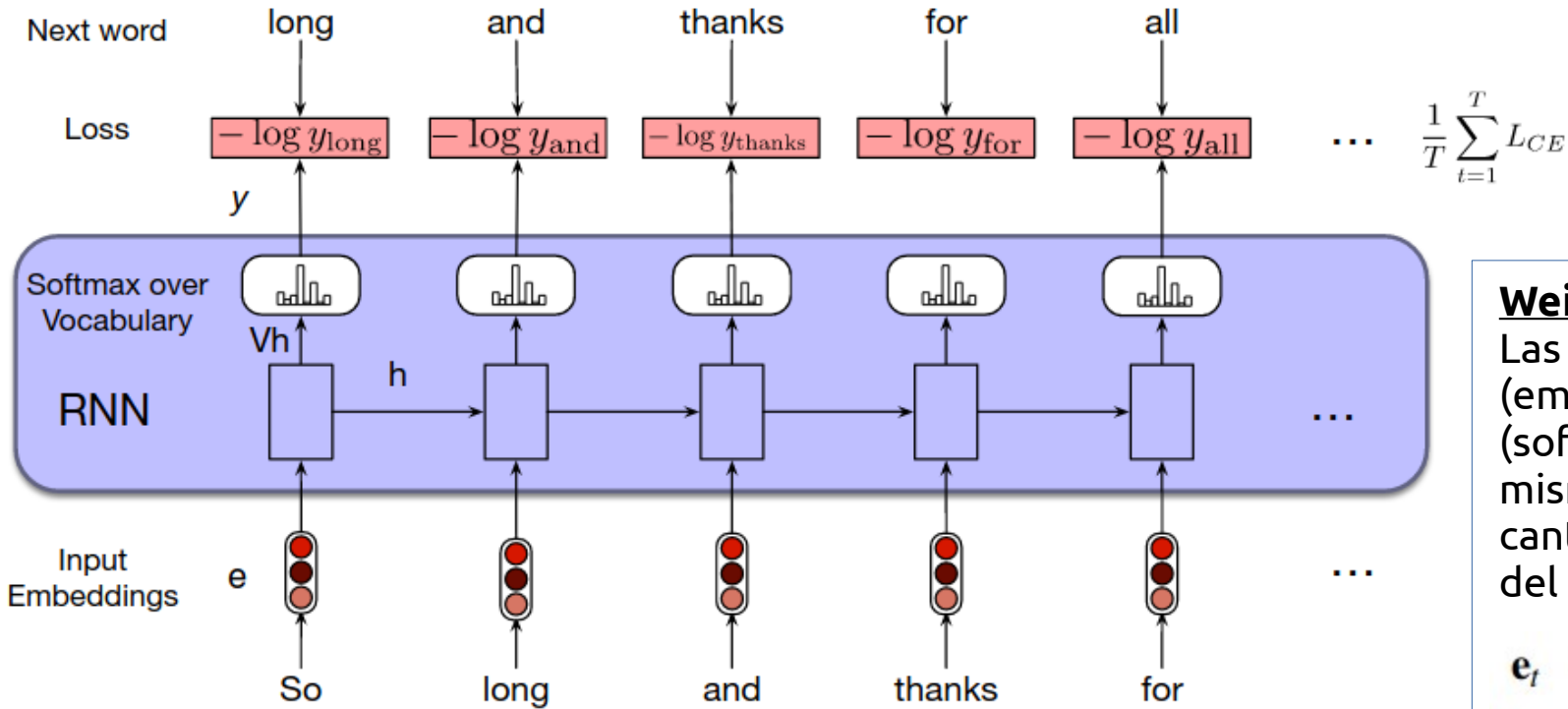
return \mathbf{y}

Entrenamiento

Cada paso es una composición
(Backpropagation Through Time, BPTT)



RNN as Language Models (Mikolov, 2010)



Weights tying

Las matrices E (embeddings) y V (softmax) pueden ser la misma, reduciendo la cantidad de parámetros del modelo.

$$\mathbf{e}_t = \mathbf{E}\mathbf{x}_t$$

$$\mathbf{h}_t = g(\mathbf{U}\mathbf{h}_{t-1} + \mathbf{W}\mathbf{e}_t)$$

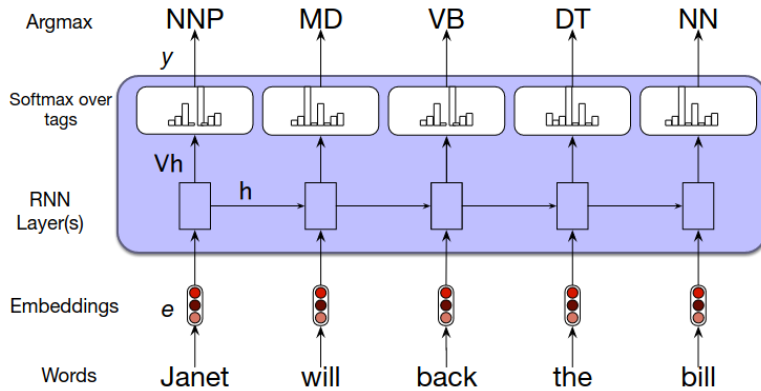
$$\mathbf{y}_t = \text{softmax}(\mathbf{E}^{\text{intercal}} \mathbf{h}_t)$$

Teacher forcing

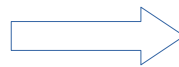
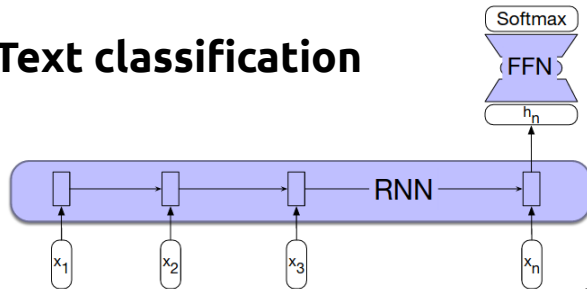
Al entrenar cada paso se usa la secuencia correcta de entrada, se ignoran las predicciones del modelo.

RNN for Other NLP Tasks

Sequence labeling



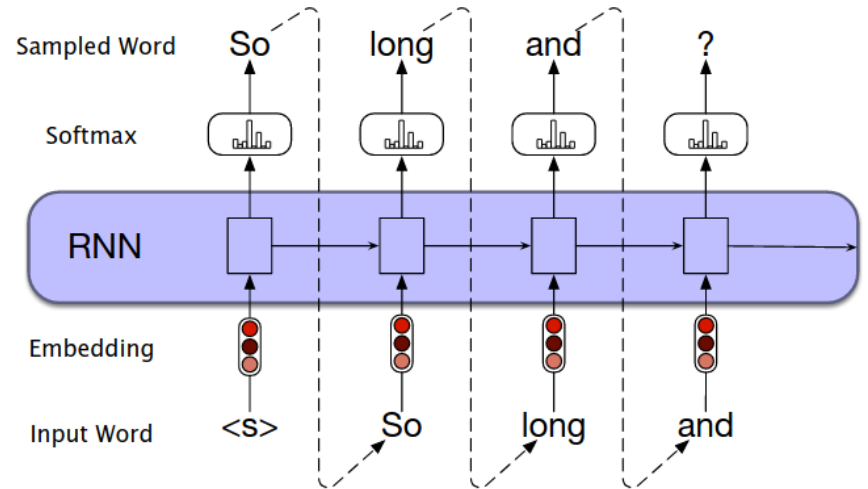
Text classification



$$\mathbf{h}_{mean} = \frac{1}{n} \sum_{i=1}^n \mathbf{h}_i$$

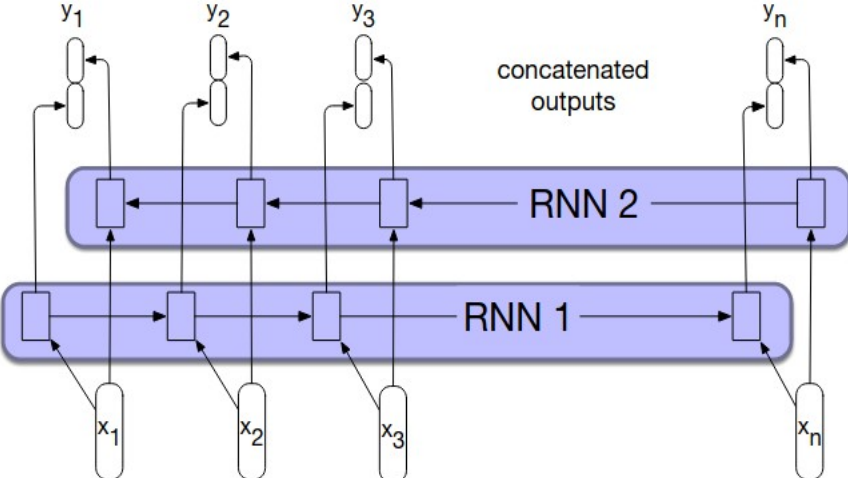
..o usar una combinación de todos los estados ocultos

Text generation

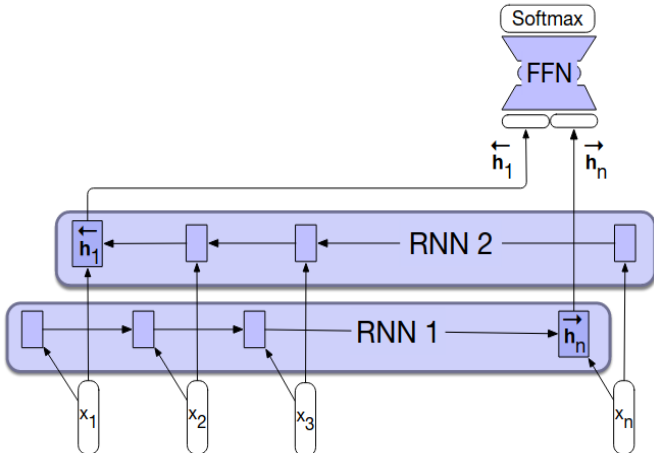


Bidirectional RNNs

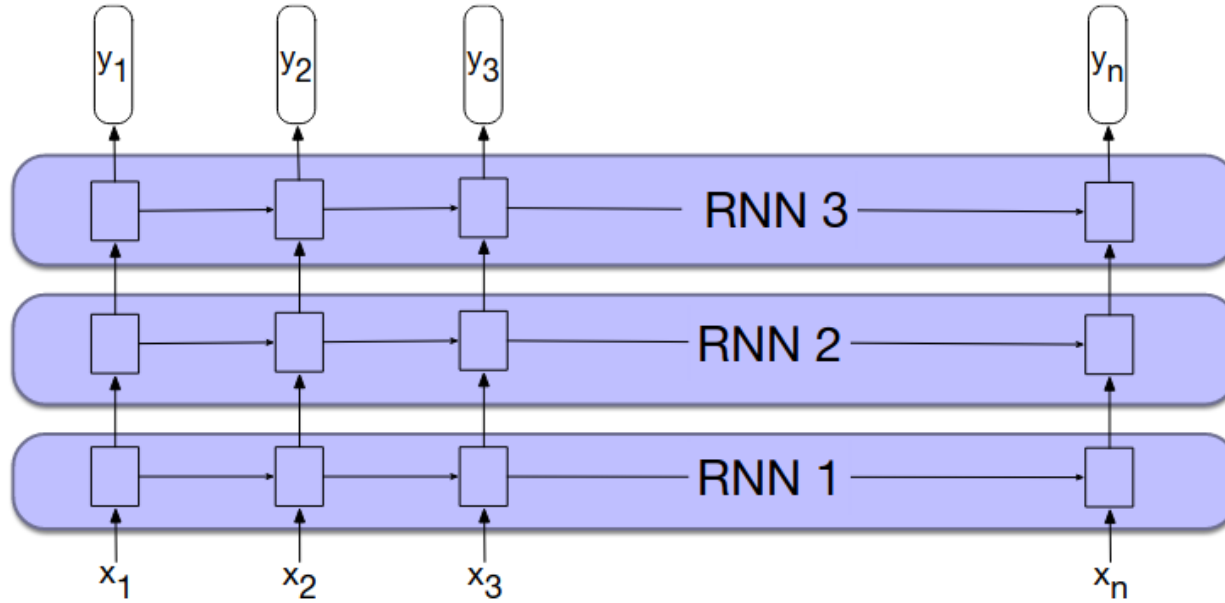
Sequence labeling



Sequence classification



Stacked RNNs



RNNs Vanishing Gradient Problem

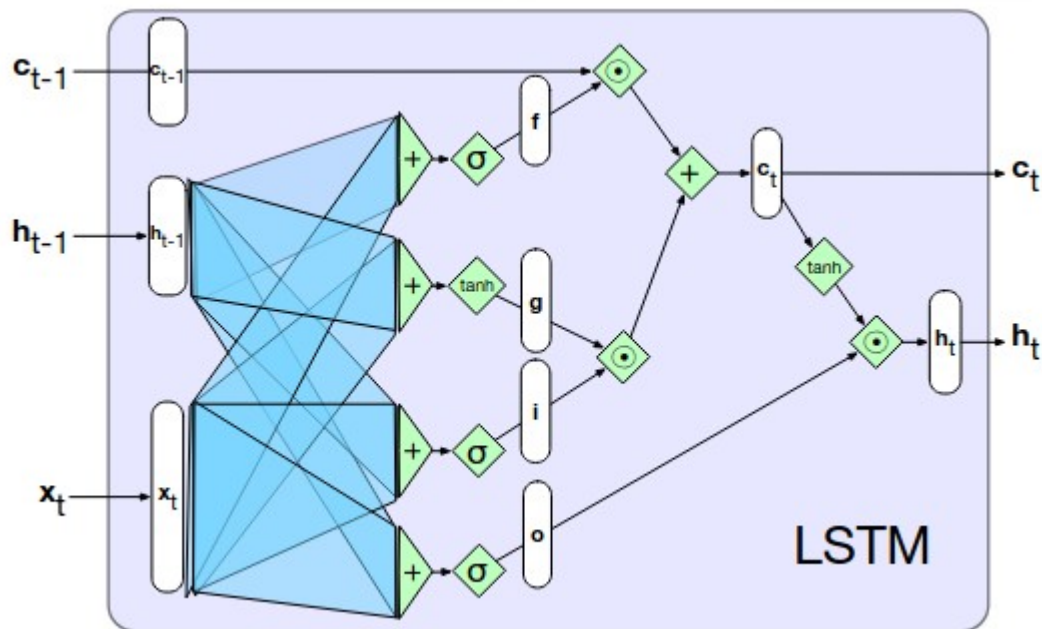
Las RNNs tienden a degradar su funcionamiento con la información más alejada (ej. dependencias de largo alcance)

Vanishing Gradients Problem

Los *steps* en una RNN (ej. cada token de la secuencia) resulta en multiplicaciones de gradientes que eventualmente se van a 0, dependiendo de su cercanía a 0 y el largo de la secuencia

Long short-term memory (Hochreiter and Schmidhuber, 1997)

LSTM es un tipo de capa recurrente con un vector de contexto (c_t) que puede ser “propagado”. El diseño se basa en **compuertas** con capas feedforward, activaciones sigmoid y producto element-wise (Hadamard) (funciona como una máscara binaria)



$$f_t = \sigma(U_f h_{t-1} + W_f x_t)$$

$$k_t = c_{t-1} \odot f_t$$

$$g_t = \tanh(U_g h_{t-1} + W_g x_t)$$

$$i_t = \sigma(U_i h_{t-1} + W_i x_t)$$

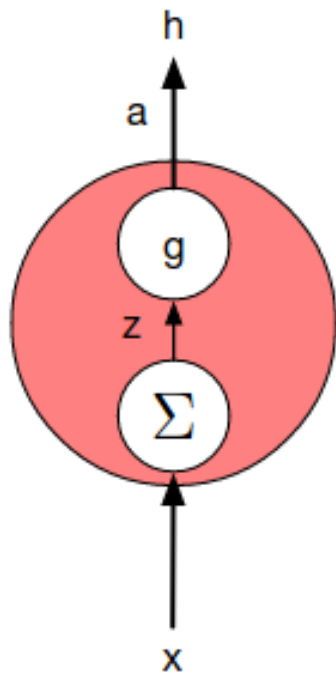
$$j_t = g_t \odot i_t$$

$$c_t = j_t + k_t$$

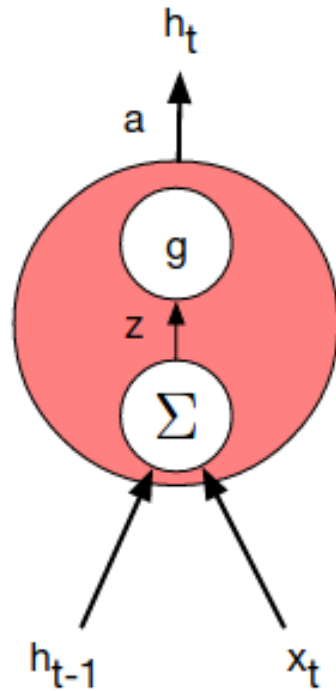
$$o_t = \sigma(U_o h_{t-1} + W_o x_t)$$

$$h_t = o_t \odot \tanh(c_t)$$

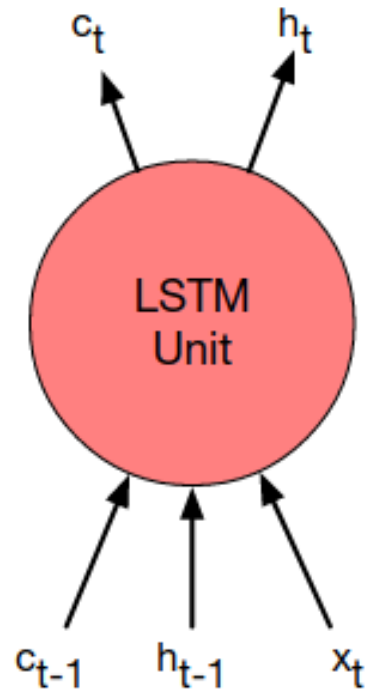
Este tipo de formulación de la capa recurrente no es única, por ejemplo una GRU (Gated Recurrent Unit) es un tipo de capa recurrente con ideas similares, comparable a LSTM.



(a)



(b)



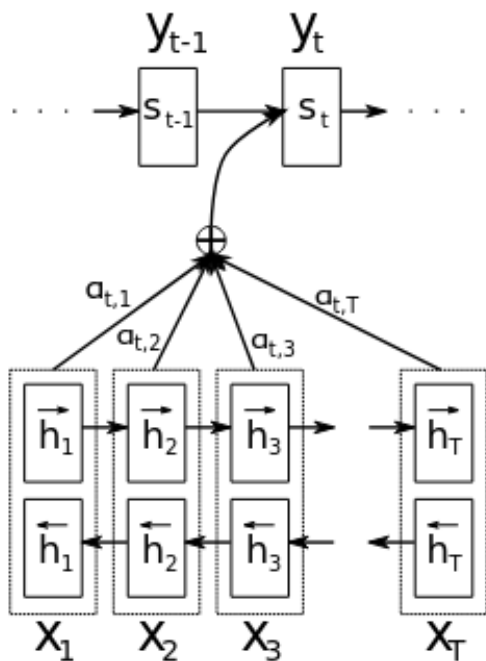
(c)

LSTMs y GRUs mejoran el manejo de información más distante pero no resuelven el problema

RNN + Attention (soft-search)

Bahdanau, D., Cho, K., & Bengio, Y. (2014)

Planteado para RNN encoder-decoder. Al generar, en cada paso, se considera como contexto un promedio ponderado (usando softmax y ffn) de los estados ocultos (anotaciones) de la entrada.



$$p(y_i | y_1, \dots, y_{i-1}, \mathbf{x}) = g(y_{i-1}, s_i, c_i),$$

where s_i is an RNN hidden state for time i , computed by

$$s_i = f(s_{i-1}, y_{i-1}, c_i).$$

$$c_i = \sum_{j=1}^{T_x} \alpha_{ij} h_j.$$

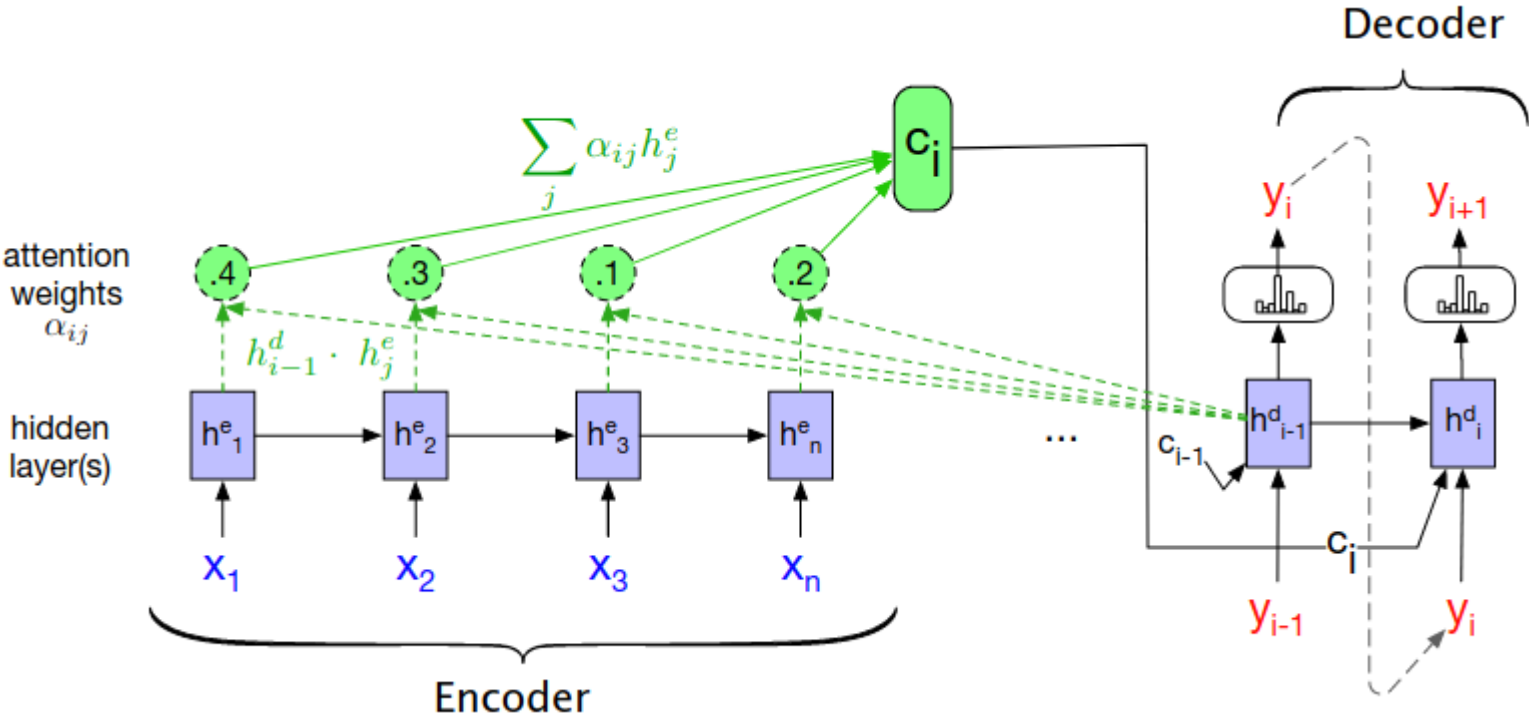
The weight α_{ij} of each annotation h_j is computed by

$$\alpha_{ij} = \frac{\exp(e_{ij})}{\sum_{k=1}^{T_x} \exp(e_{ik})},$$

where

$$e_{ij} = a(s_{i-1}, h_j)$$

RNN + Attention



ELMo (representaciones contextualizadas)

Peters, Matthew E. et al. "Deep Contextualized Word Representations." NAACL (2018).

ULMFiT (fine-tuning)

Howard, Jeremy and Sebastian Ruder. "Fine-tuned Language Models for Text Classification." ArXiv abs/1801.06146 (2018): n. pag.

Transformer

Transformer

La naturaleza recurrente de las RNNs dificulta el procesamiento paralelo.

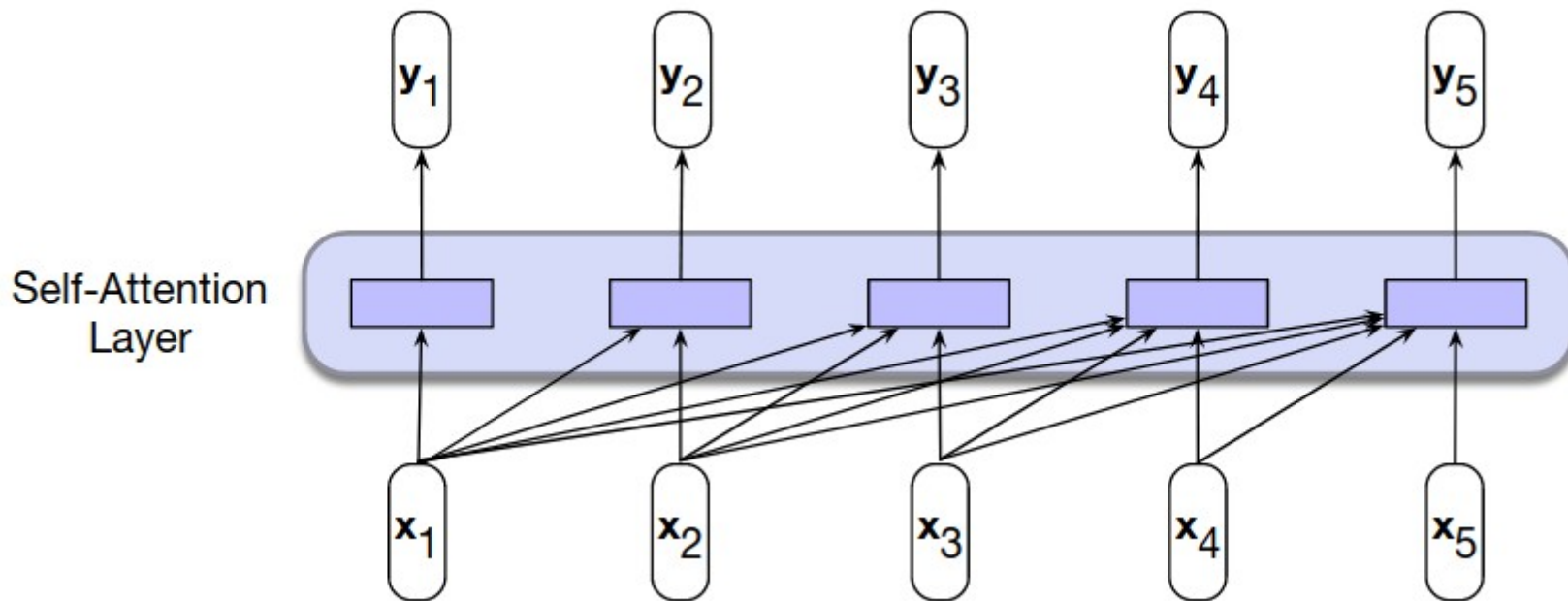
Reúne varias ideas exitosas

- self-attention
- positional embeddings
- contextualized embeddings
- fine-tuning
- masked language modeling
- next sentence prediction

El transformer puede procesar todas las entradas en paralelo.

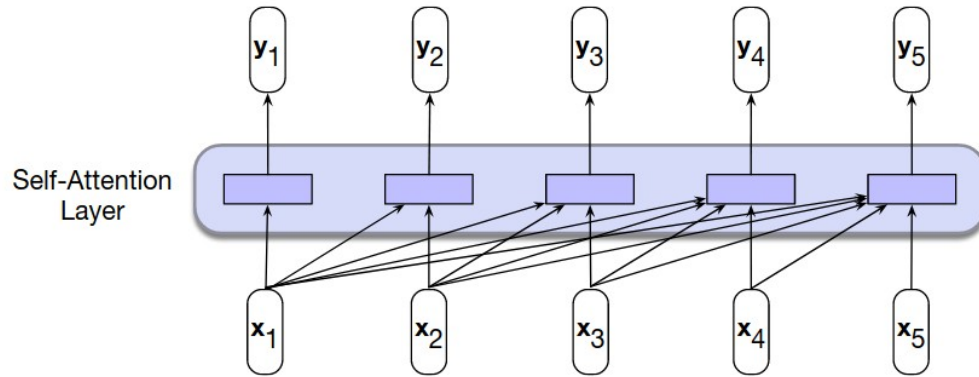
Transformer: Self-Attention

Transformer: Self-Attention



A diferencia de una RNN, cada paso es independiente (se puede paralelizar).

Transformer: Self-Attention

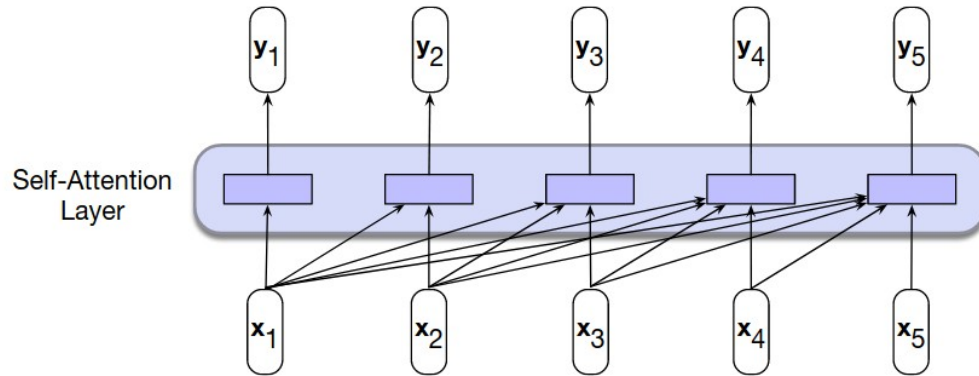


$$y_i = \sum_{j \leq i} \alpha_{ij} x_j$$

$$\begin{aligned} \alpha_{ij} &= \text{softmax}(\text{score}(x_i, x_j)) \quad \forall j \leq i \\ &= \frac{\exp(\text{score}(x_i, x_j))}{\sum_{k=1}^i \exp(\text{score}(x_i, x_k))} \quad \forall j \leq i \end{aligned}$$

$$\text{score}(x_i, x_j) = ?$$

Transformer: Self-Attention



Se plantea *self-attention* con 3 proyecciones de cada x_i (**query, key y value**)

$$\mathbf{q}_i = \mathbf{W}^Q \mathbf{x}_i; \quad \mathbf{k}_i = \mathbf{W}^K \mathbf{x}_i; \quad \mathbf{v}_i = \mathbf{W}^V \mathbf{x}_i$$

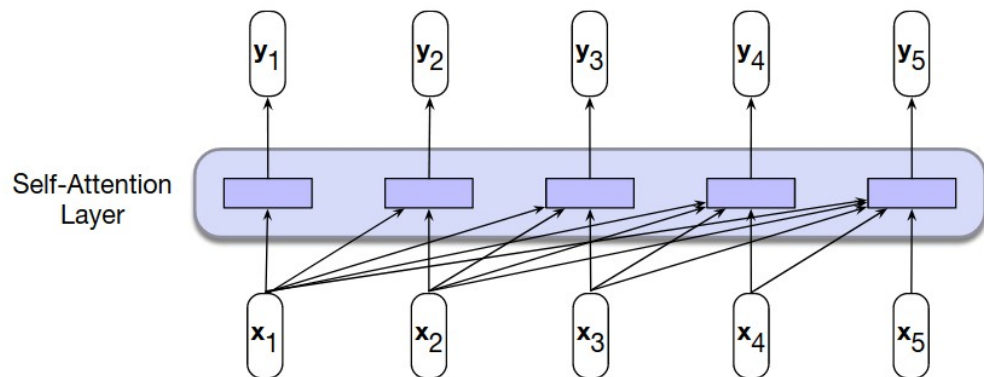
Los vectores \mathbf{q}_i y \mathbf{k}_j computan **score** de j para i

$$\text{score}(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{q}_i \cdot \mathbf{k}_j$$

Con softmax se suman los \mathbf{v}_j

$$\mathbf{y}_i = \sum_{j \leq i} \alpha_{ij} \mathbf{v}_j$$

Transformer: Self-Attention



Se plantea *self-attention* con 3 proyecciones de cada x_i (**query, key y value**)

$$\mathbf{q}_i = \mathbf{W}^Q \mathbf{x}_i; \quad \mathbf{k}_i = \mathbf{W}^K \mathbf{x}_i; \quad \mathbf{v}_i = \mathbf{W}^V \mathbf{x}_i$$

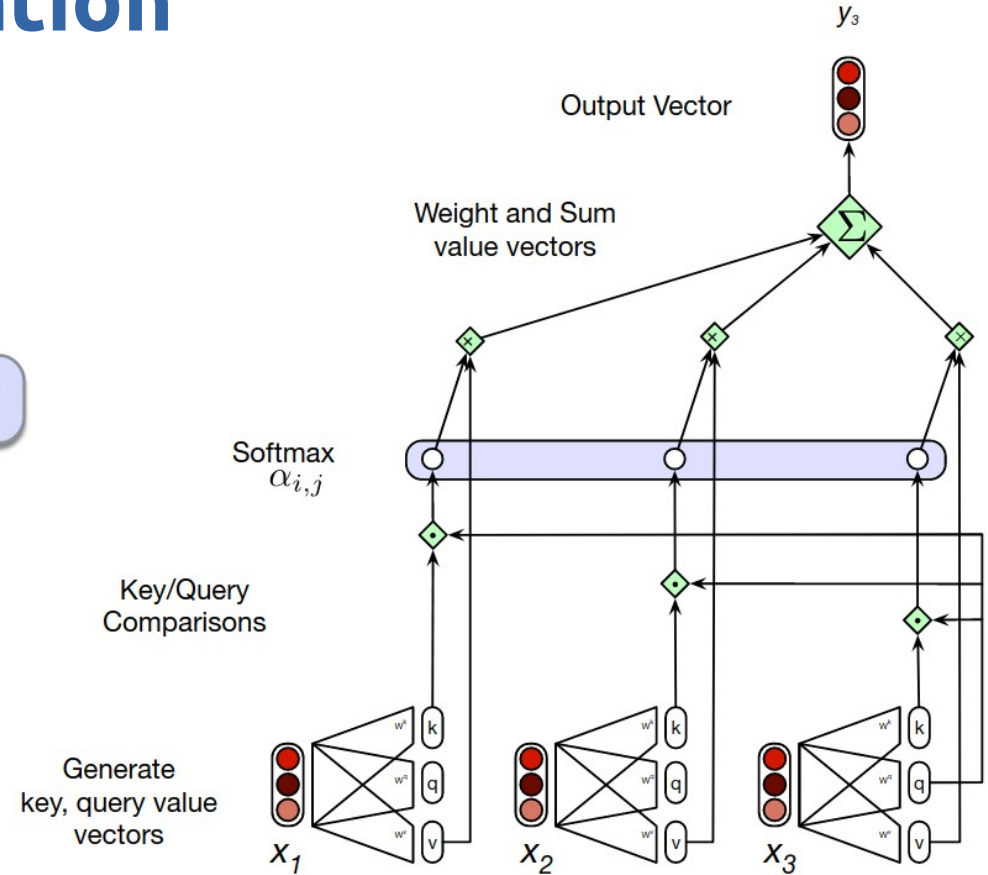
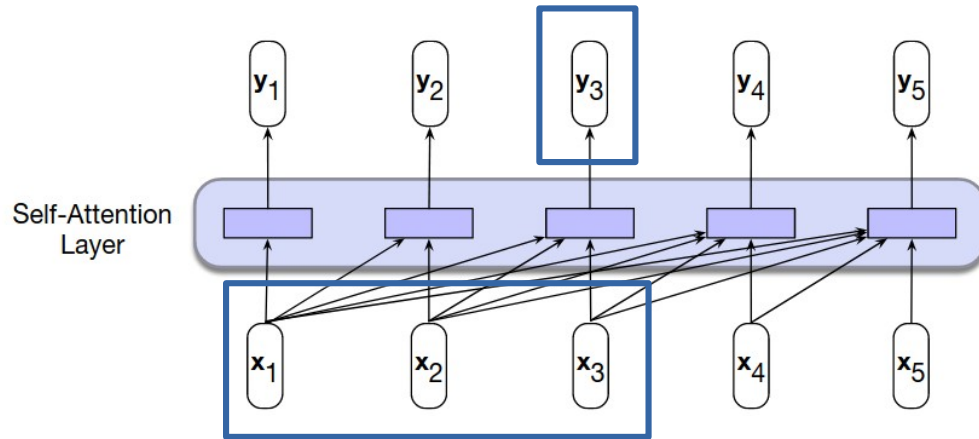
Los vectores \mathbf{q}_i y \mathbf{k}_j computan **score** de j para i

$$\text{score}(\mathbf{x}_i, \mathbf{x}_j) = \frac{\mathbf{q}_i \cdot \mathbf{k}_j}{\sqrt{d_k}} \quad \leftarrow \text{para suavizar valores grandes}$$

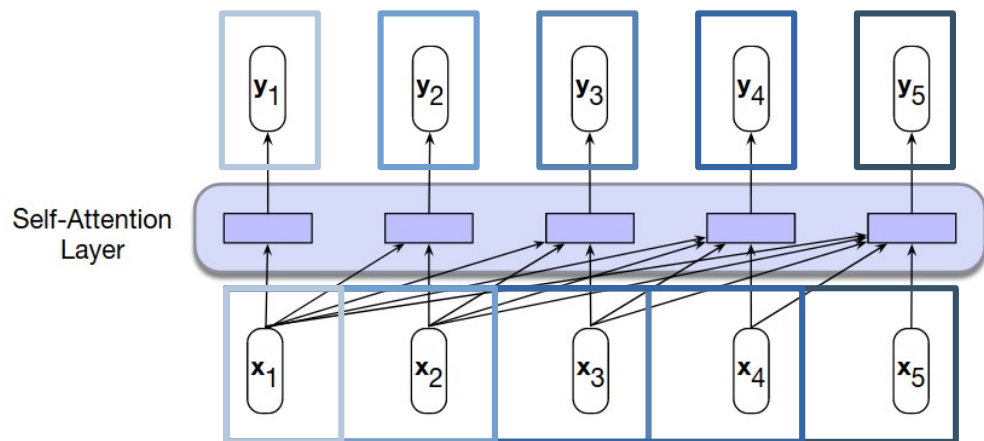
Con softmax se suman los \mathbf{v}_j

$$\mathbf{y}_i = \sum_{j \leq i} \alpha_{ij} \mathbf{v}_j$$

Transformer: Self-Attention



Transformer: Self-Attention



$$q_i = W^Q x_i; \quad k_i = W^K x_i; \quad v_i = W^V x_i$$

Como cada paso es independiente podemos hacerlos en paralelo

$$Q = XW^Q; \quad K = XW^K; \quad V = XW^V$$

... quedando

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

$QK^T \Rightarrow N$

q1•k1	-∞	-∞	-∞	-∞
q2•k1	q2•k2	-∞	-∞	-∞
q3•k1	q3•k2	q3•k3	-∞	-∞
q4•k1	q4•k2	q4•k3	q4•k4	-∞
q5•k1	q5•k2	q5•k3	q5•k4	q5•k5

Se rellena con -∞ (se vuelven 0 con softmax)

Multihead Attention

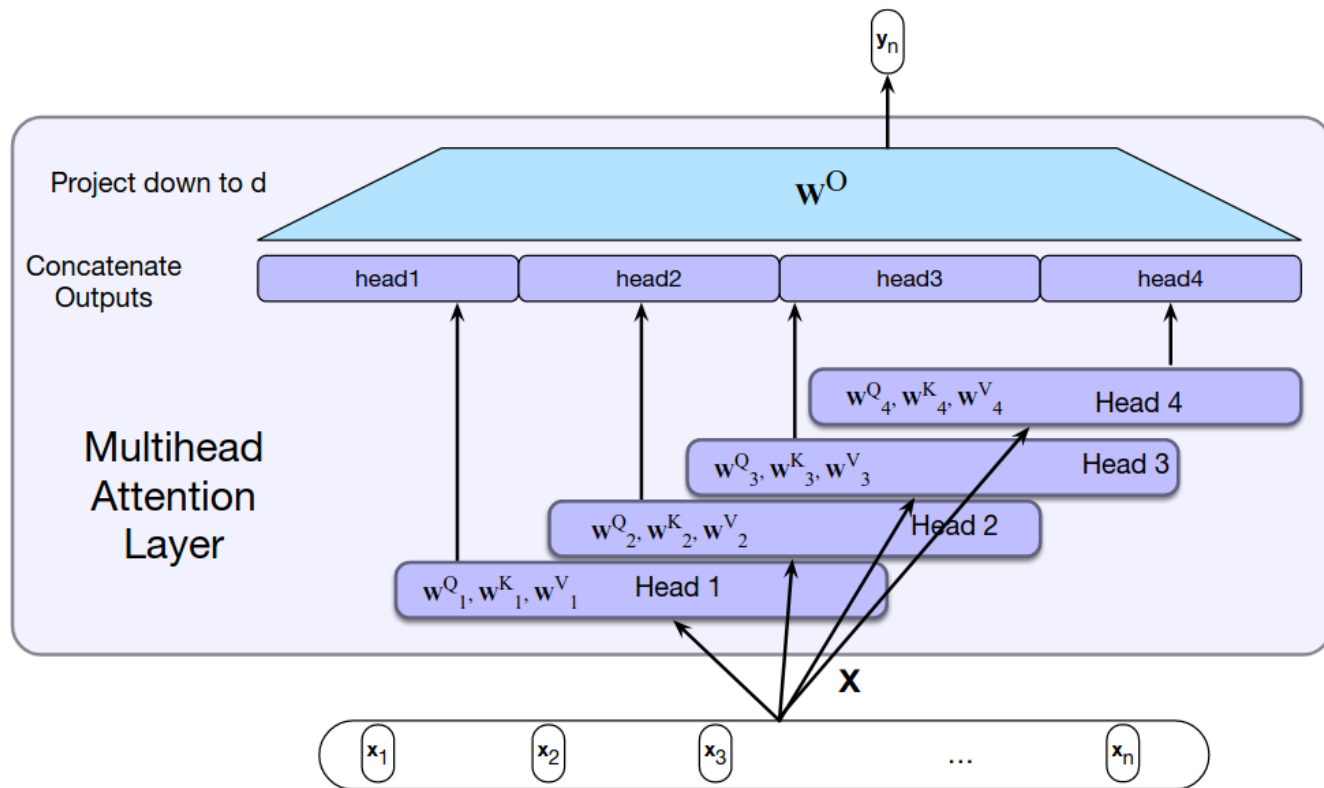
Multihead Attention

$$\mathbf{W}_i^Q \in \mathbb{R}^{d \times d_k}, \mathbf{W}_i^K \in \mathbb{R}^{d \times d_k}, \mathbf{W}_i^V \in \mathbb{R}^{d \times d_v}$$

$$\text{MultiHeadAttn}(\mathbf{X}) = (\text{head}_1 \oplus \text{head}_2 \dots \oplus \text{head}_h) \mathbf{W}^O$$

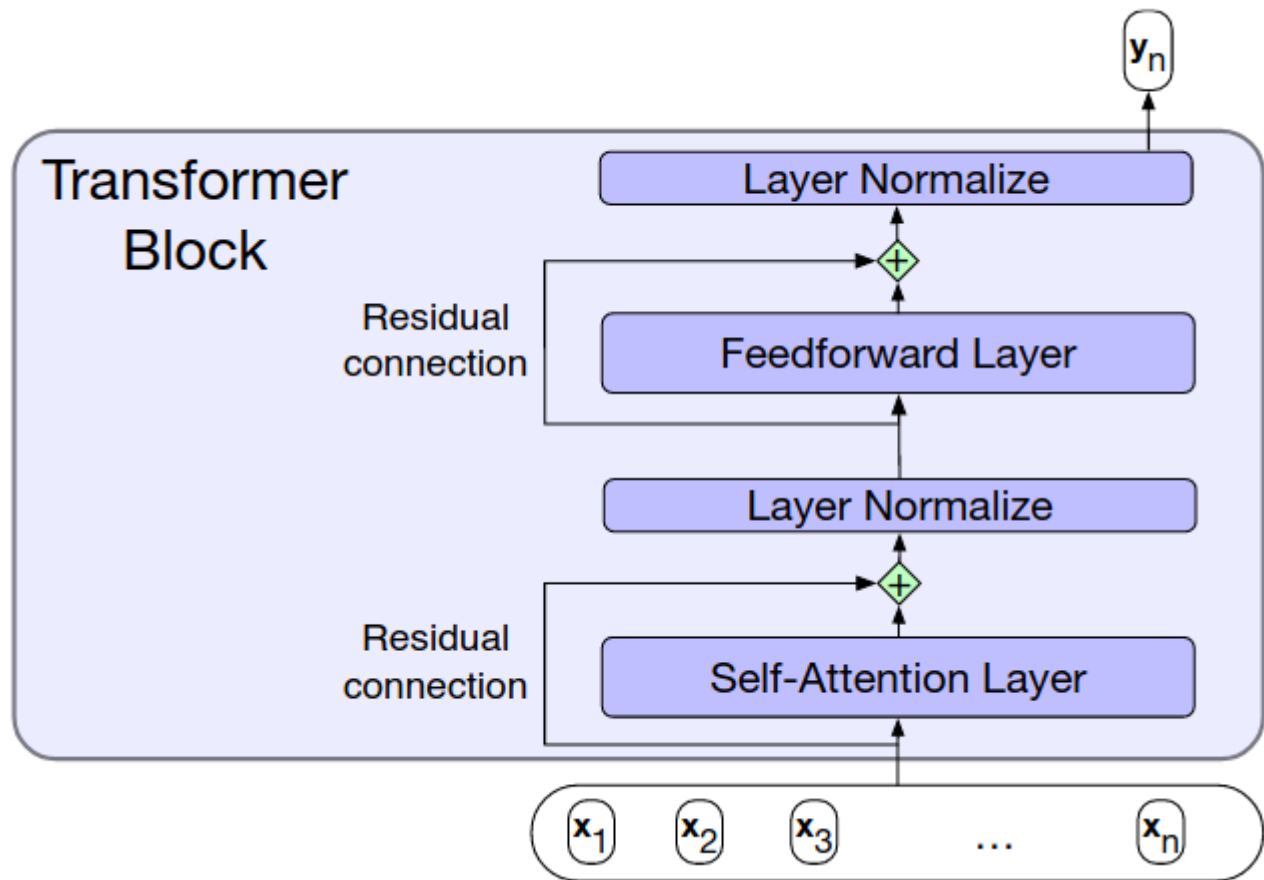
$$\mathbf{Q}_i = \mathbf{X} \mathbf{W}_i^Q; \mathbf{K}_i = \mathbf{X} \mathbf{W}_i^K; \mathbf{V}_i = \mathbf{X} \mathbf{W}_i^V$$

$$\text{head}_i = \text{SelfAttention}(\mathbf{Q}_i, \mathbf{K}_i, \mathbf{V}_i)$$



Transformer Block

Transformer Block



$$\mathbf{z} = \text{LayerNorm}(\mathbf{x} + \text{SelfAttn}(\mathbf{x}))$$
$$\mathbf{y} = \text{LayerNorm}(\mathbf{z} + \text{FFNN}(\mathbf{z}))$$

Layer Normalization

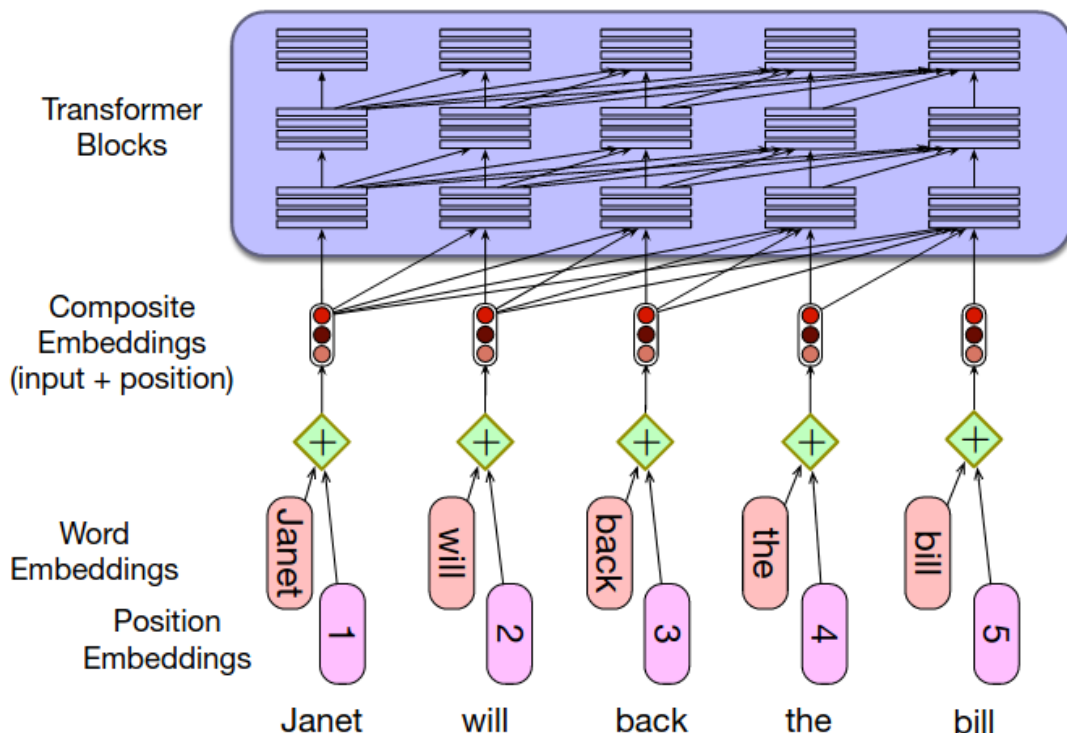
$$\mu = \frac{1}{d_h} \sum_{i=1}^{d_h} x_i$$
$$\sigma = \sqrt{\frac{1}{d_h} \sum_{i=1}^{d_h} (x_i - \mu)^2}$$
$$\hat{\mathbf{x}} = \frac{(\mathbf{x} - \mu)}{\sigma} \quad (\gamma \text{ y } \beta \text{ se aprenden})$$
$$\text{LayerNorm} = \gamma \hat{\mathbf{x}} + \beta$$

Positional Embeddings

Word Order: Positional Embeddings

Hasta este punto no hay noción de orden en la entrada.

Solucion: Agregar a cada x_i un vector que represente su posición (**positional embeddings**)

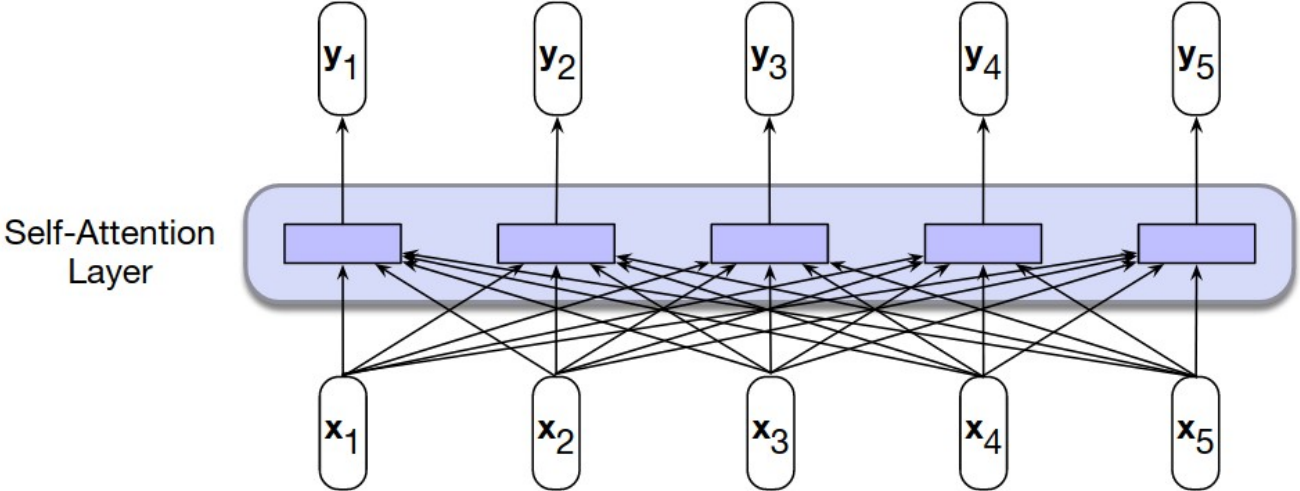


Pueden ser entrenados pero las primeras posiciones van a ser más comunes (por los diferentes largos), dificultando el entrenamiento.

Una alternativa es que sean estáticos. Una combinación de senos y cosenos es usada en el transformer original.

Mejorar los positional embeddings es un área actual de investigación.

Bidirectional Transformer



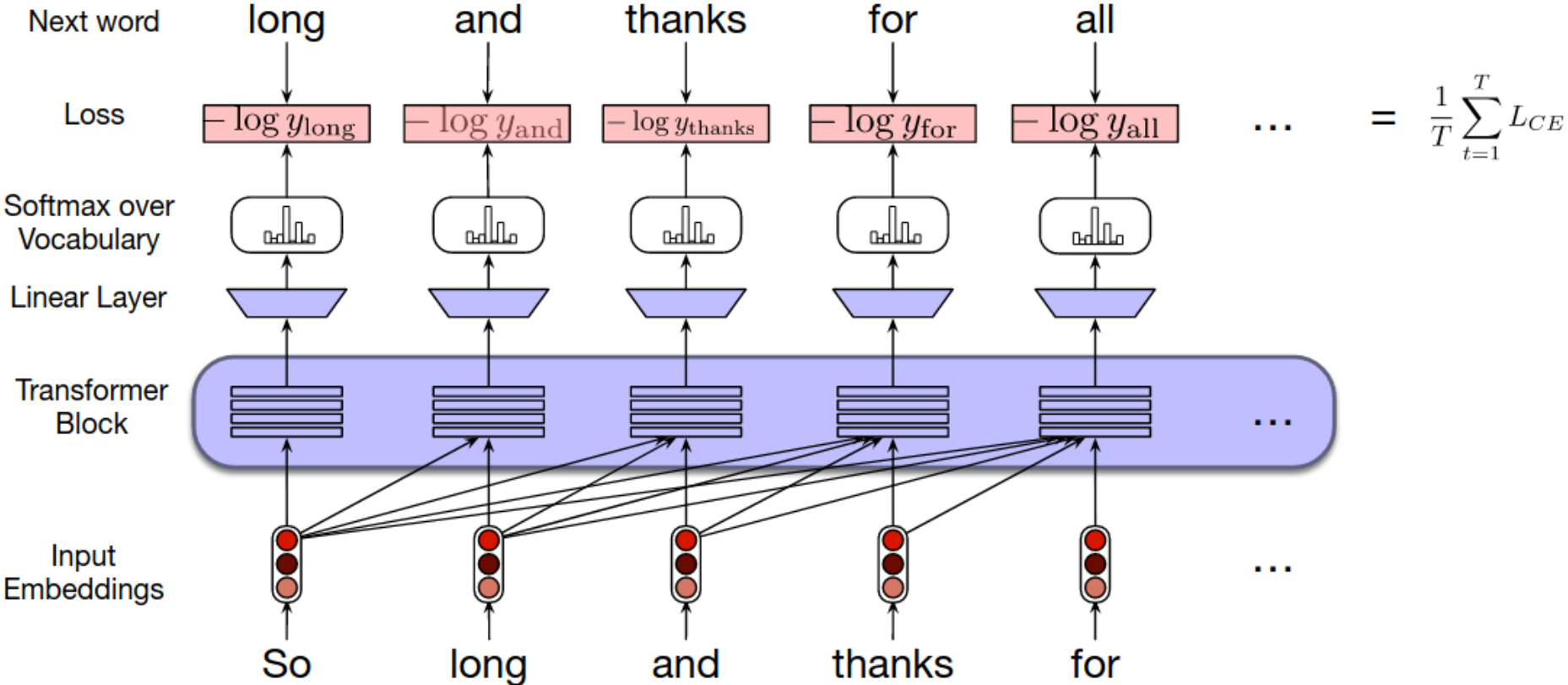
N

$q1 \cdot k1$	$q1 \cdot k2$	$q1 \cdot k3$	$q1 \cdot k4$	$q1 \cdot k5$
$q2 \cdot k1$	$q2 \cdot k2$	$q2 \cdot k3$	$q2 \cdot k4$	$q2 \cdot k5$
$q3 \cdot k1$	$q3 \cdot k2$	$q3 \cdot k3$	$q3 \cdot k4$	$q3 \cdot k5$
$q4 \cdot k1$	$q4 \cdot k2$	$q4 \cdot k3$	$q4 \cdot k4$	$q4 \cdot k5$
$q5 \cdot k1$	$q5 \cdot k2$	$q5 \cdot k3$	$q5 \cdot k4$	$q5 \cdot k5$

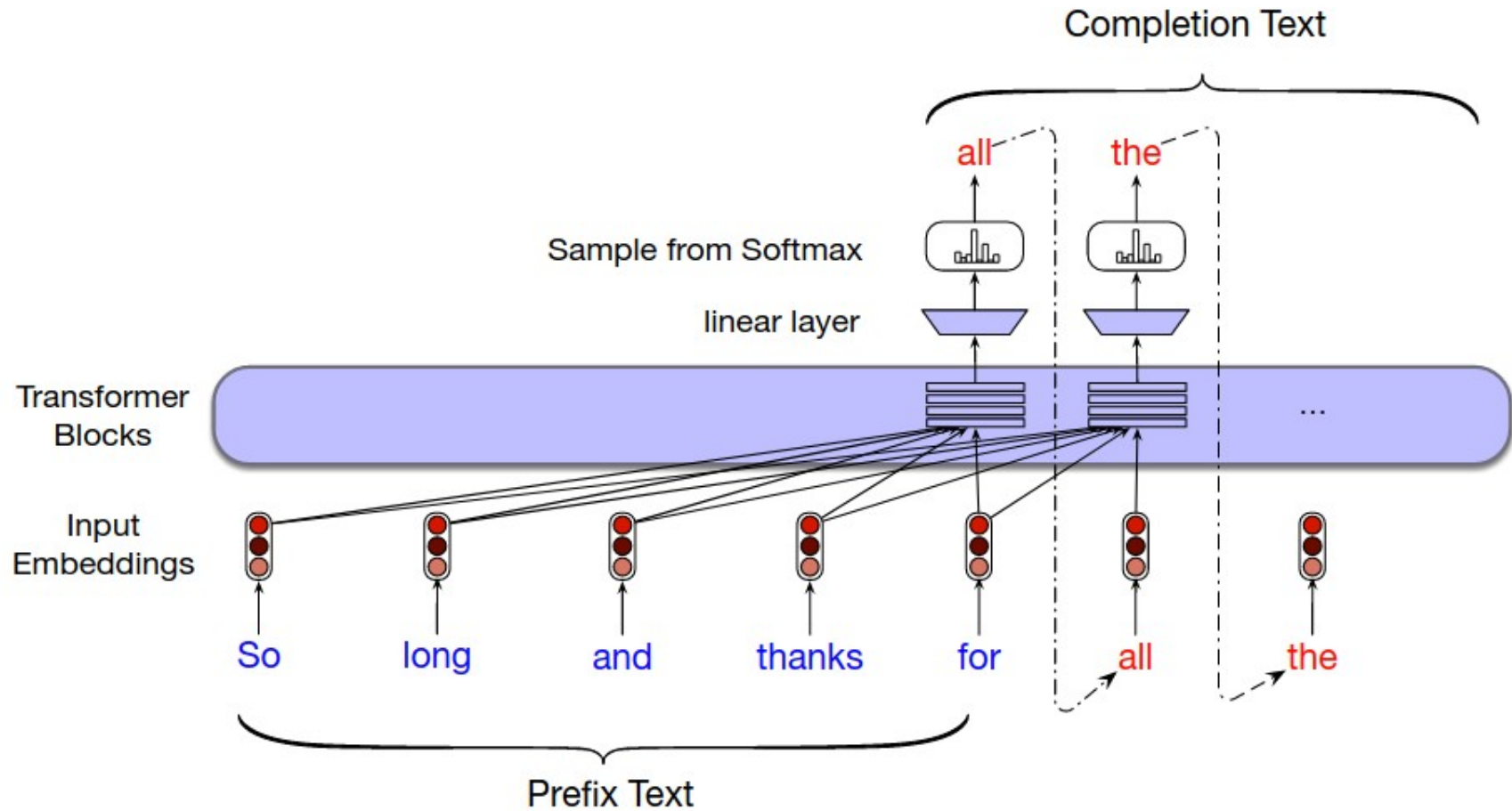
N

¿Cómo los usamos?

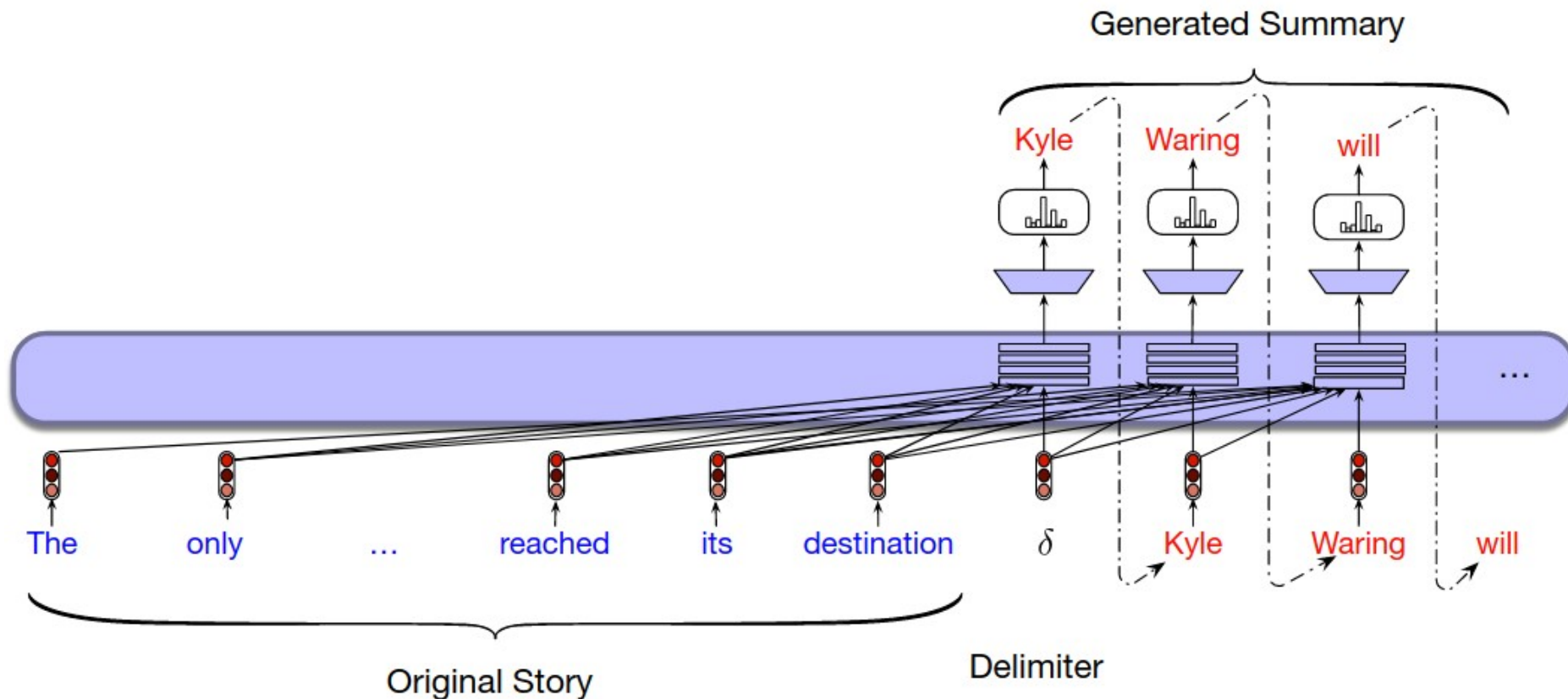
Language Model con Transformer



Generar con Transformers



Conditioned Generation with Transformers



Pretraining y fine-tuning

Se preentrena un transformer como un modelo de lenguaje con un corpus grande (**pretraining**)

Luego, se agrega una capa feedforward que es ajustada con un corpus anotado (**fine-tuning**)