



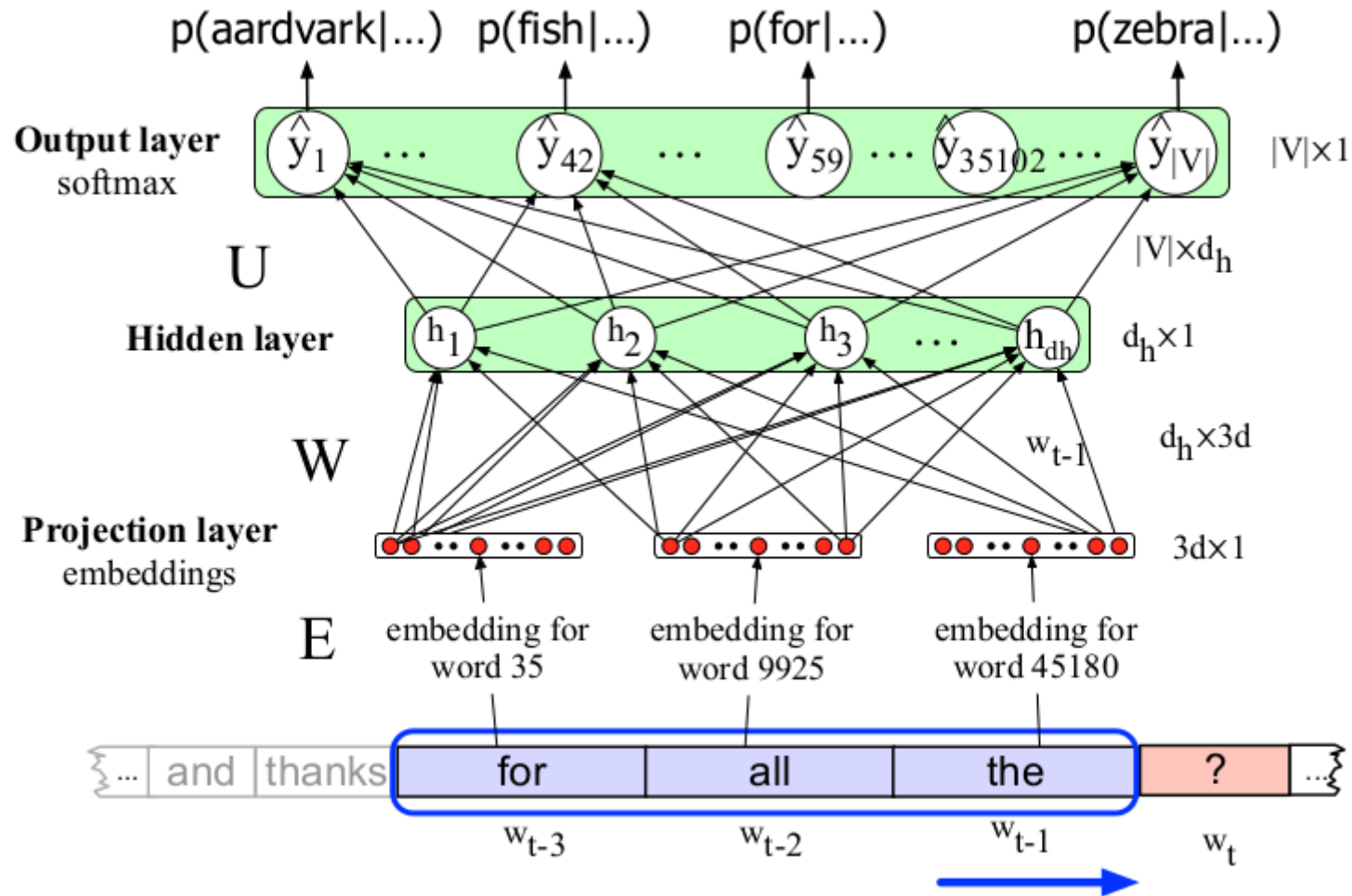
IIC 3800 Tópicos en CC NLP

<https://github.com/marcelomendoza/IIC3800>

- MODELOS DE LENGUAJE -

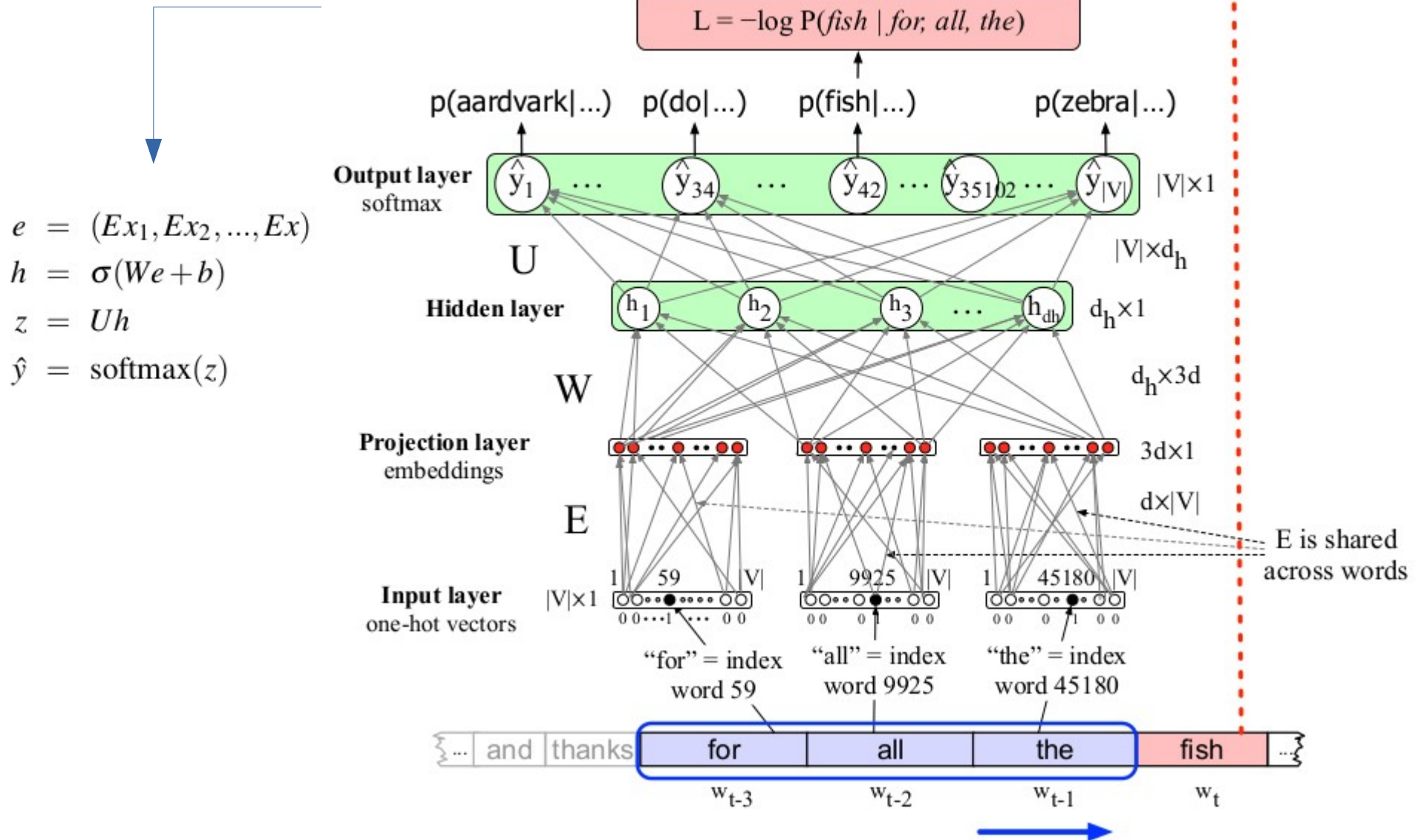
Background: Modelos de lenguaje neuronales

Feed-forward (modelo de lenguaje):



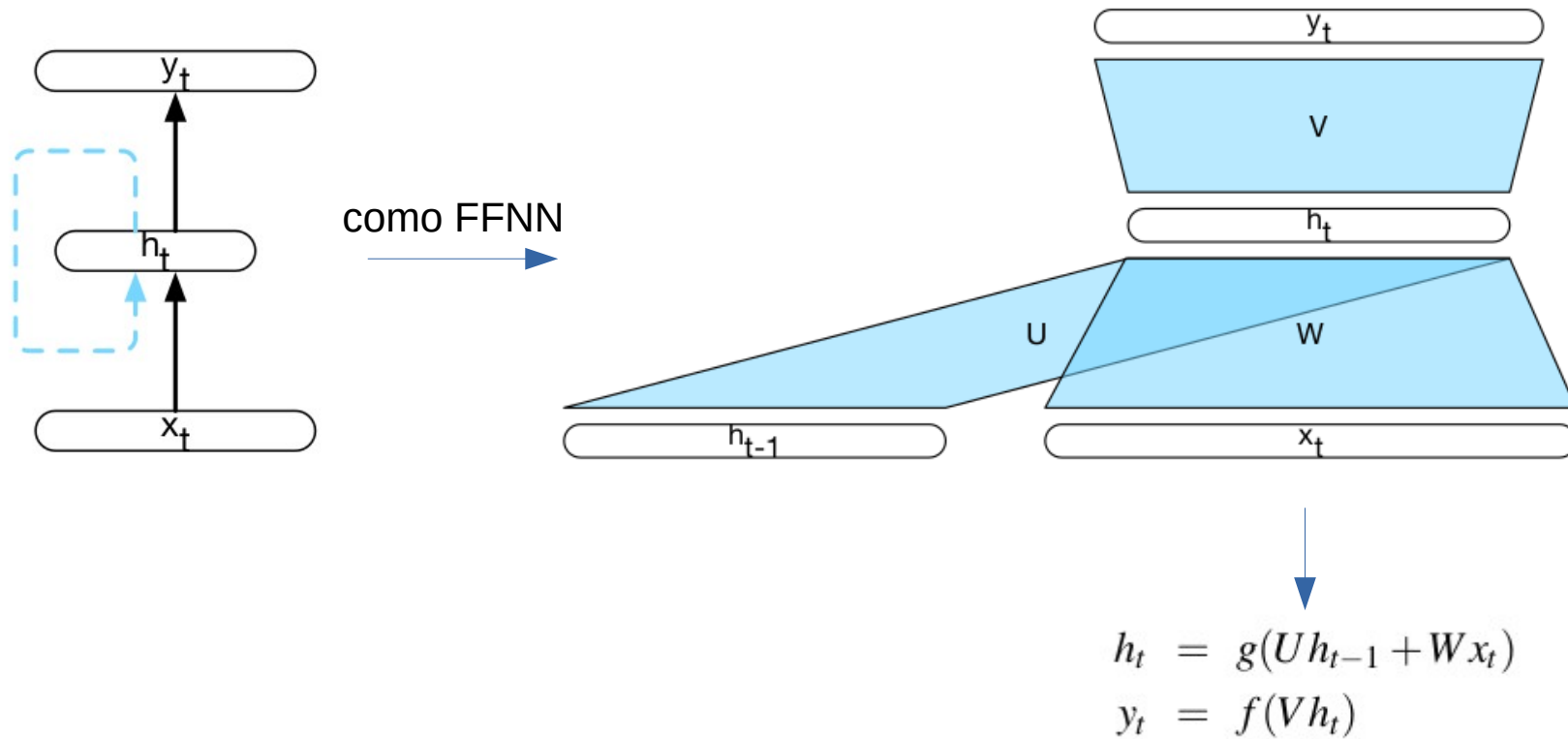
Background: Modelos de lenguaje neuronales

Feed-forward (modelo de lenguaje):



Background: Modelos de lenguaje neuronales

Red recurrente:



Típicamente, en NLP: $y_t = \text{softmax}(Vh_t)$

Background: Modelos de lenguaje neuronales

Red recurrente:

output distribution

$$\hat{y}^{(t)} = \text{softmax} \left(U h^{(t)} + b_2 \right) \in \mathbb{R}^{|V|}$$

hidden states

$$h^{(t)} = \sigma \left(W_h h^{(t-1)} + W_e e^{(t)} + b_1 \right)$$

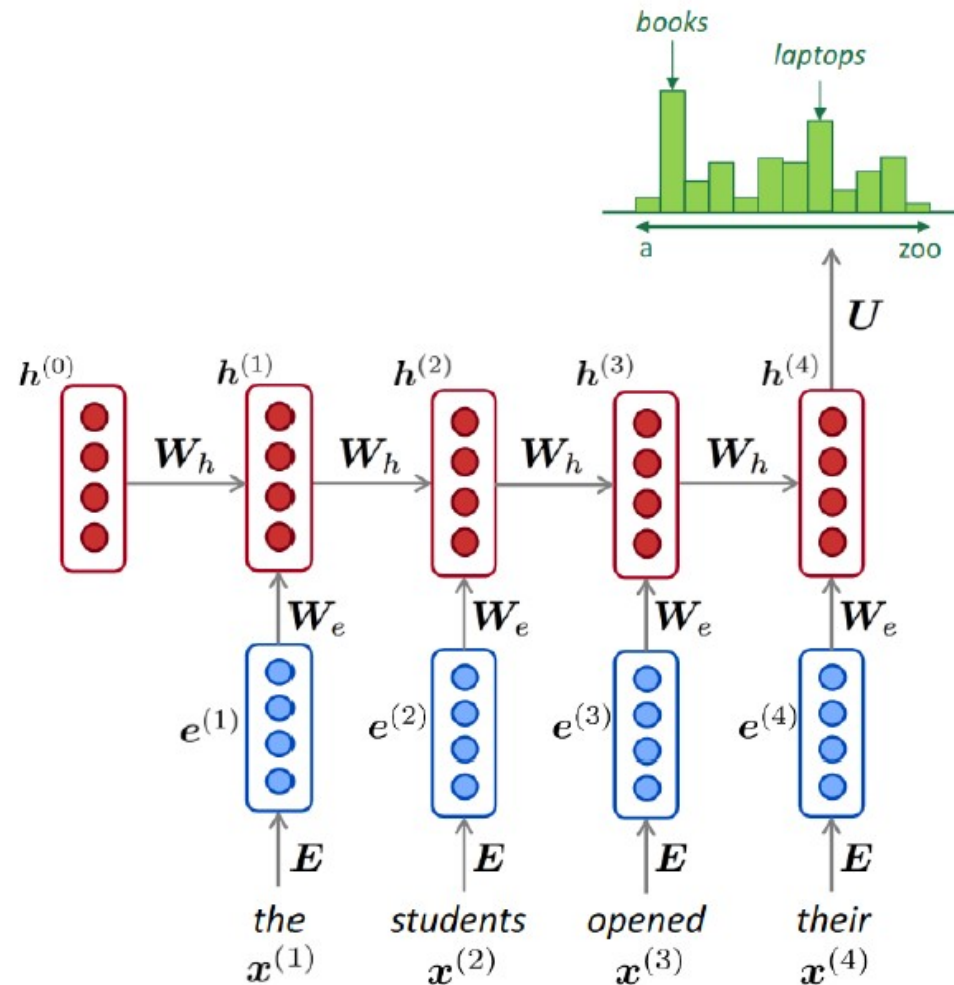
$h^{(0)}$ is the initial hidden state

word embeddings

$$e^{(t)} = E x^{(t)}$$

words / one-hot vectors

$$x^{(t)} \in \mathbb{R}^{|V|}$$



Bengio, Y., Ducharme, R., Vincent, P., Jauvin, C. A Neural Probabilistic Language Model. *Journal of Machine Learning Research* 3:1137-1155, 2003.

Background: Modelos de lenguaje neuronales

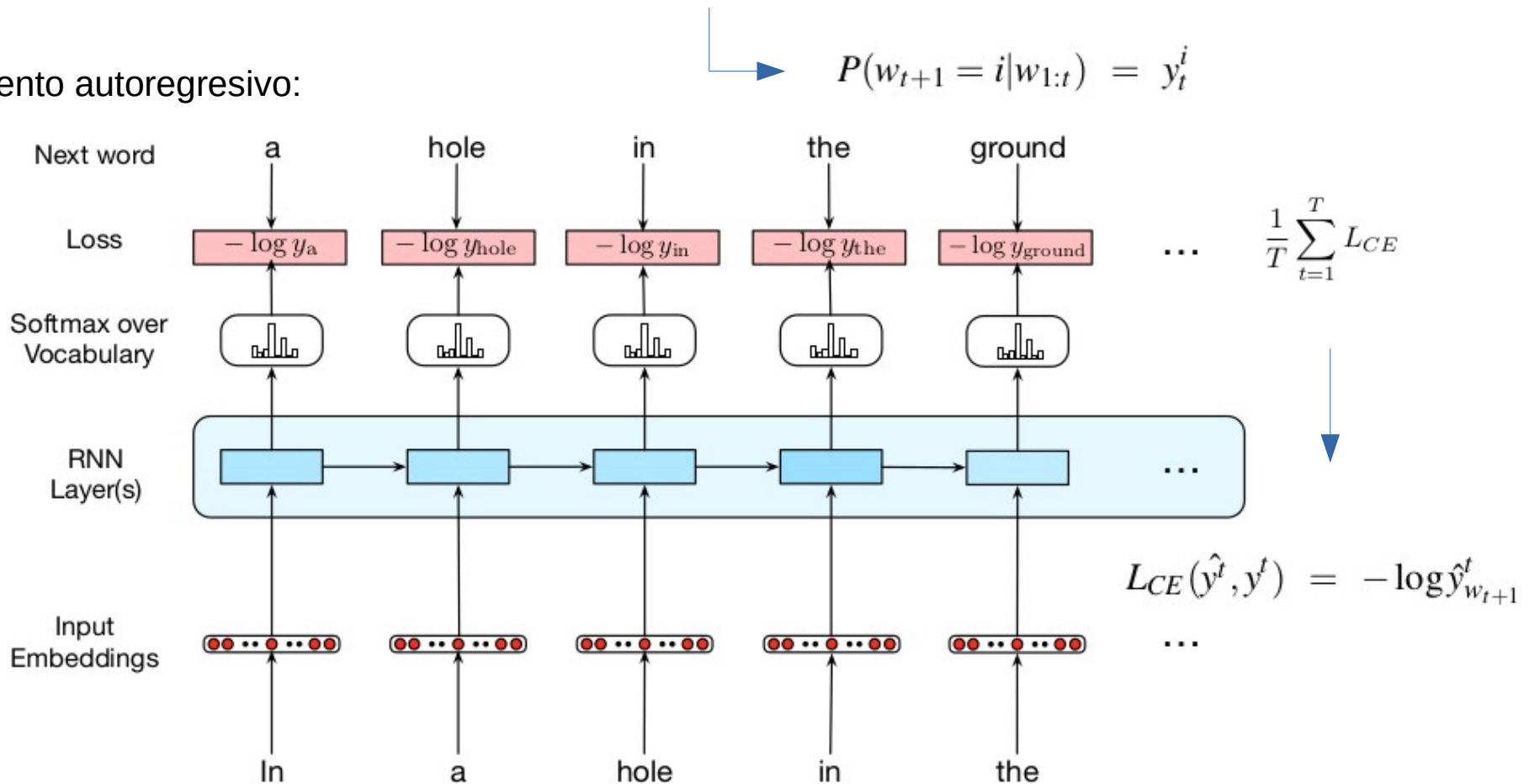
Red recurrente:

$$e_t = E^T x_t$$

$$h_t = g(Uh_{t-1} + We_t)$$

$$y_t = \text{softmax}(Vh_t)$$

Entrenamiento autoregresivo:



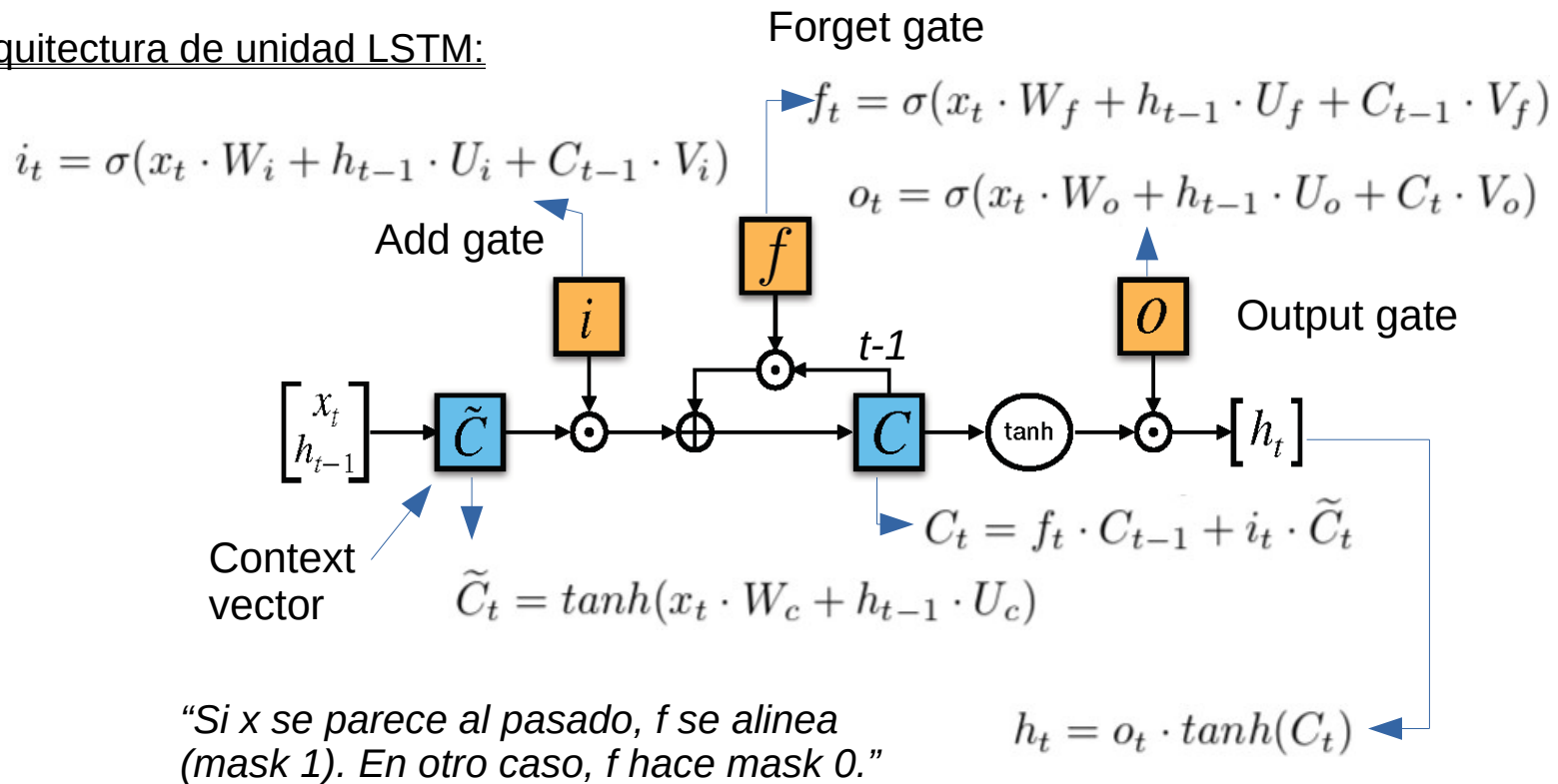
- DEPENDENCIA DEL CONTEXTO -

Extensión a LSTM

Long-short term memory (LSTM): dividen el problema en dos sub-problemas: a) remover información innecesaria, b) agregar información necesaria para resolver una tarea.

Las LSTM hacen lo anterior agregando unidades de control de flujo de información. Cada unidad consiste de una capa FF seguida por una sigmoid (activación). Combinando la salida de la sigmoid con productos, se obtiene un efecto similar al de una máscara binaria.

Arquitectura de unidad LSTM:



Dependencia del contexto

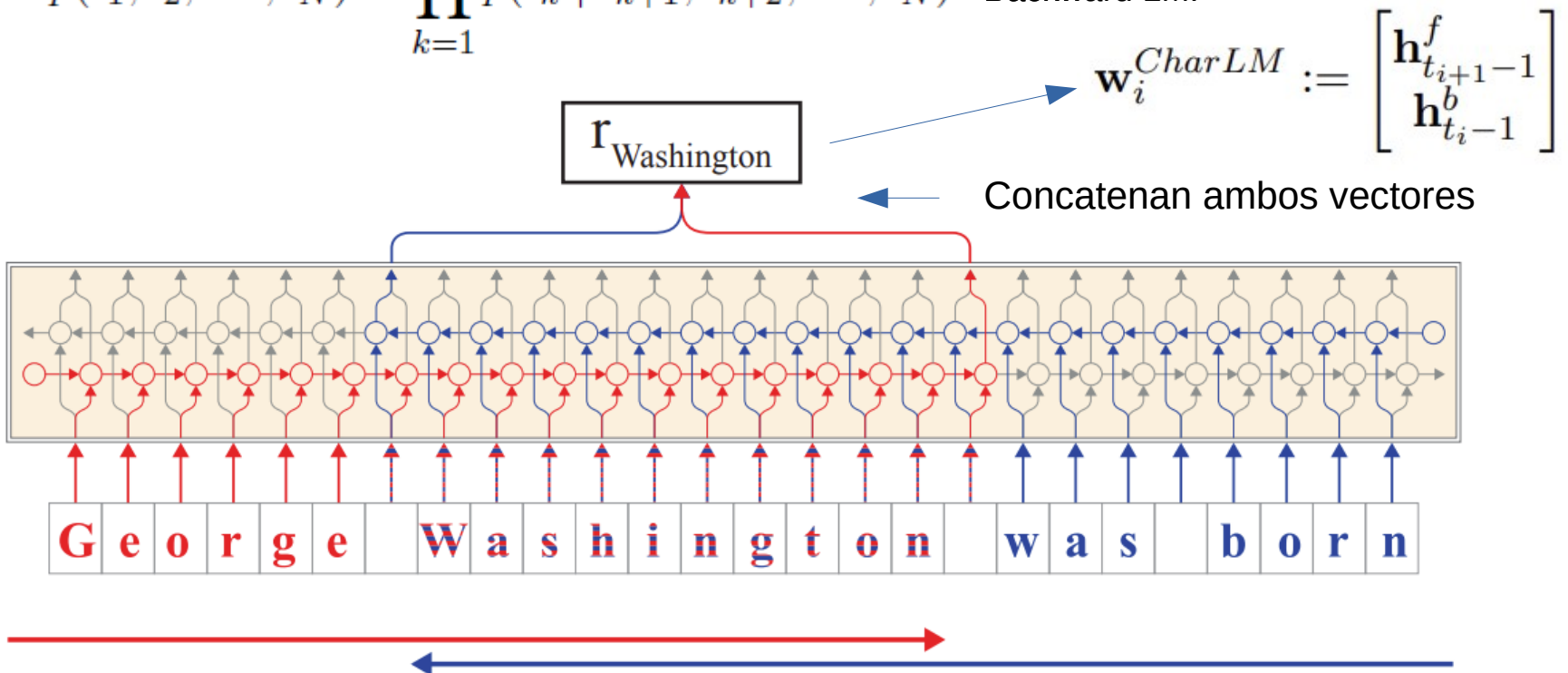
FLAIR

Dada una secuencia de n símbolos, se calculan dos modelos de lenguajes:

LSTM

$$p(t_1, t_2, \dots, t_N) = \prod_{k=1}^N p(t_k \mid t_1, t_2, \dots, t_{k-1}). \quad \text{Forward LM.}$$

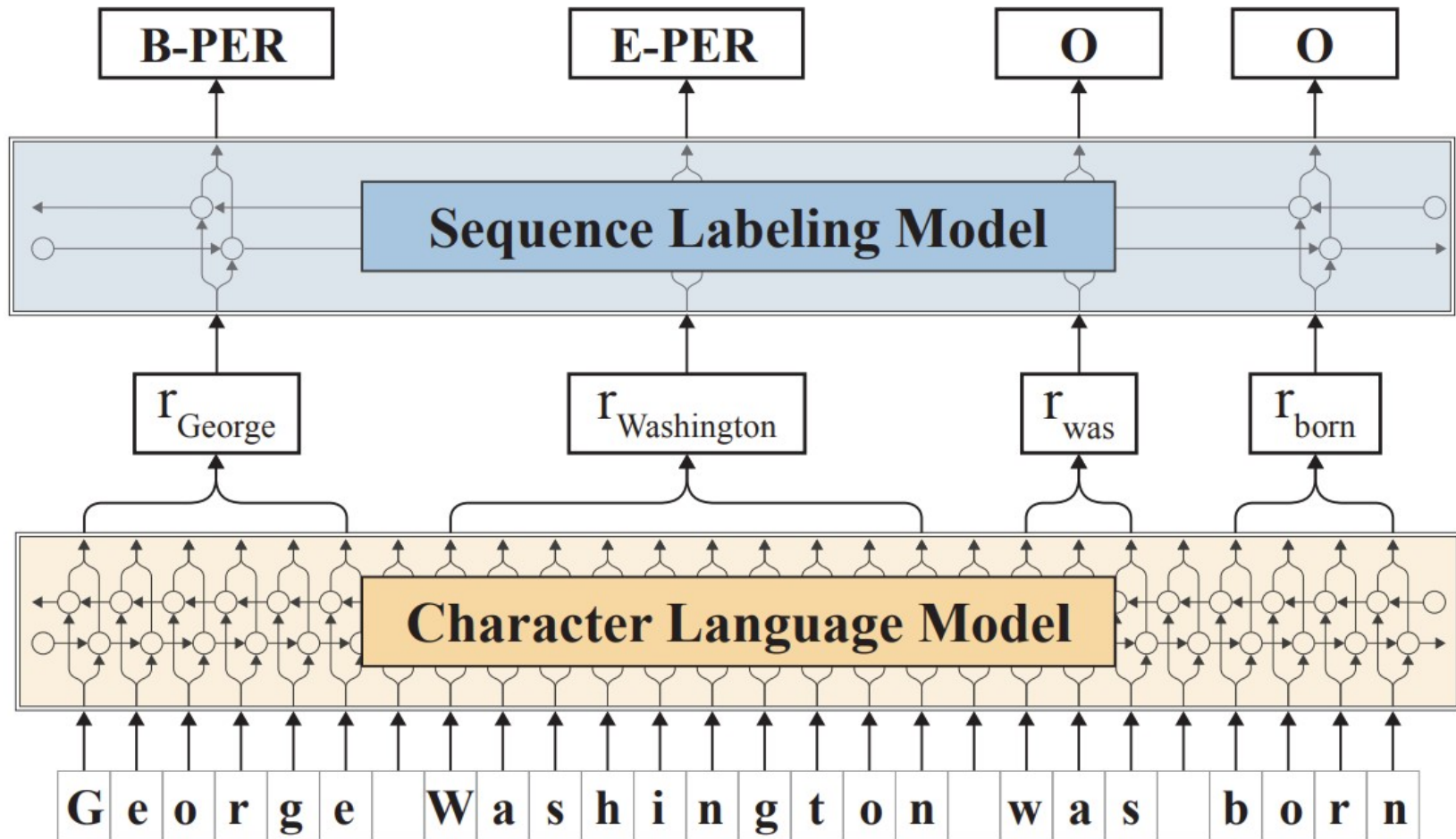
$$p(t_1, t_2, \dots, t_N) = \prod_{k=1}^N p(t_k \mid t_{k+1}, t_{k+2}, \dots, t_N). \quad \text{Backward LM.}$$



Dependencia del contexto

FLAIR

Luego usan los *string embeddings* en NER:



Dependencia del contexto

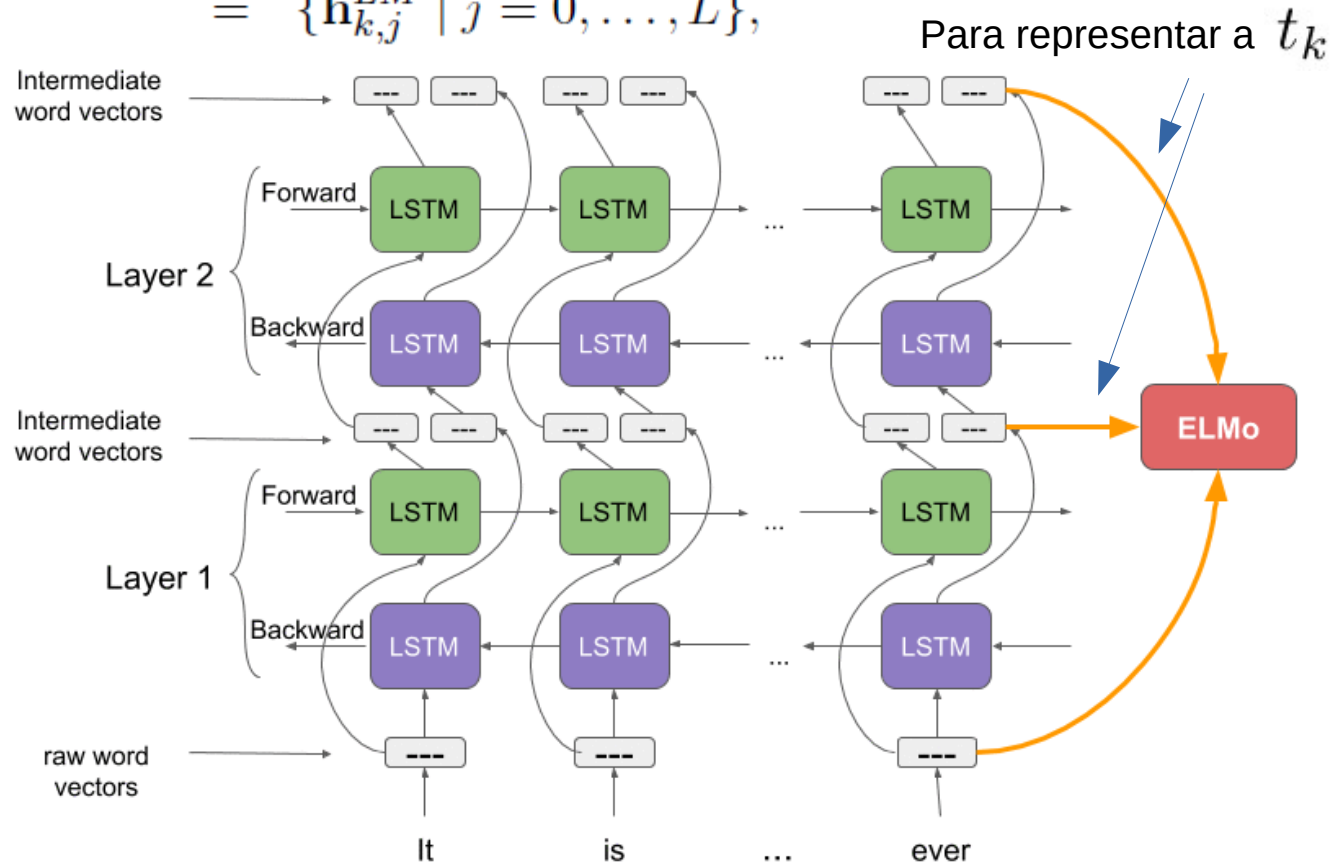
Embeddings basados en modelos de lenguaje (ELMo)

Idea similar a FLAIR pero con multicapa y modelos basados en palabras:

Parámetros de la representación ← Parámetros de las LSTM

$$\begin{aligned} R_k &= \{\mathbf{x}_k^{LM}, \vec{\mathbf{h}}_{k,j}^{LM}, \overleftarrow{\mathbf{h}}_{k,j}^{LM} \mid j = 1, \dots, L\} \\ &= \{\mathbf{h}_{k,j}^{LM} \mid j = 0, \dots, L\}, \end{aligned}$$

Arquitectura:



Dependencia del contexto

Embeddings basados en modelos de lenguaje (ELMo)

Vector de ELMo

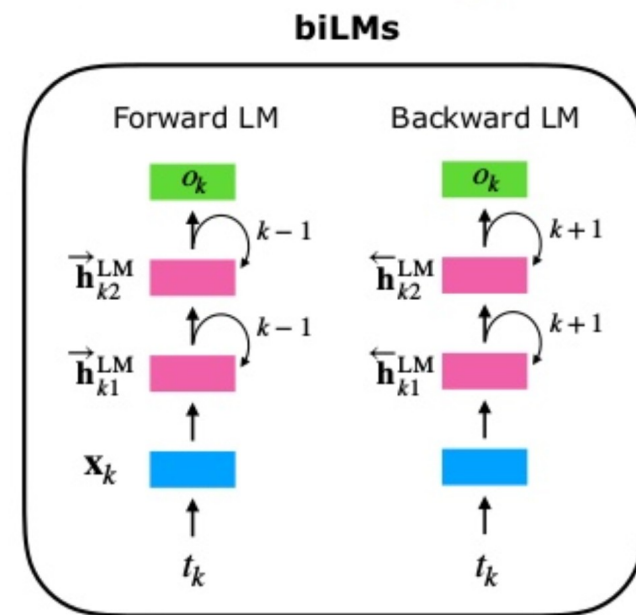
ELMo: Se calculan pesos específicos para *downstream tasks* que producen un vector ad-hoc a la tarea.

$$\mathbf{ELMo}_k^{task} = E(R_k; \Theta^{task}) = \gamma^{task} \sum_{j=0}^L s_j^{task} \mathbf{h}_{k,j}^{LM}$$

El parámetro s^{task} corresponde a pesos de la softmax. γ^{task} es un factor de escala que se introduce para efectos de optimización. En algunos casos se puede hacer normalización a nivel de capas antes de combinar los parámetros $\mathbf{h}_{k,j}^{LM}$.

$$\mathbf{ELMo}_k^{task} = \gamma^{task} \times \sum \left\{ \begin{array}{l} s_2^{task} \times \mathbf{h}_{k2}^{LM} \\ s_1^{task} \times \mathbf{h}_{k1}^{LM} \\ s_0^{task} \times \mathbf{h}_{k0}^{LM} \end{array} \right. \quad \text{Concatenate} \quad \left[\vec{\mathbf{h}}_{kj}^{LM}; \overleftarrow{\mathbf{h}}_{kj}^{LM} \right]$$

($[\mathbf{x}_k; \mathbf{x}_k]$)

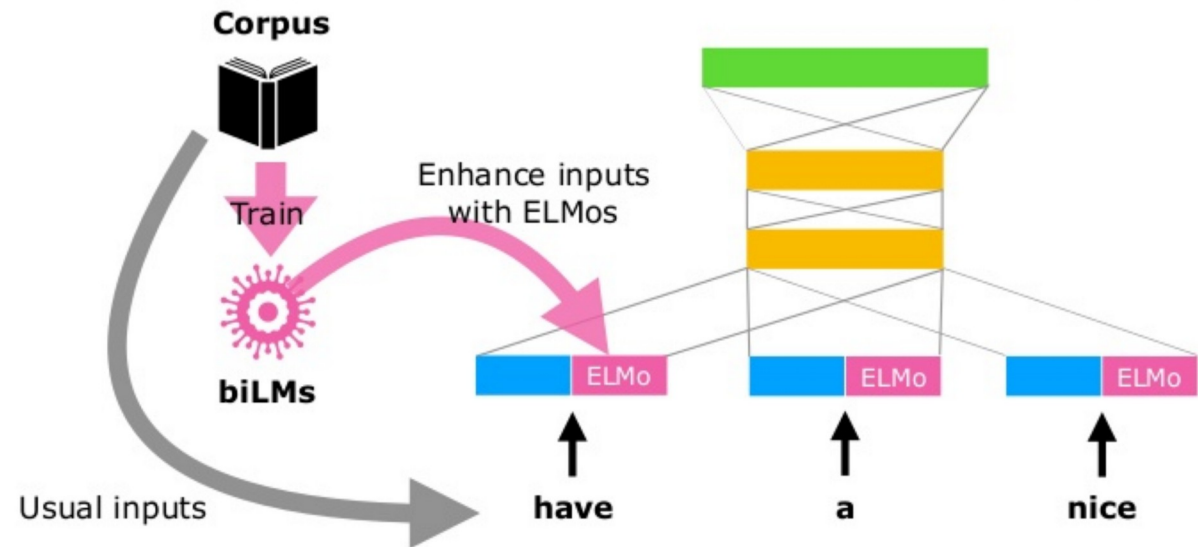


- En aprendizaje no supervisado, los vectores h se promedian.
- Al regularizar, el efecto que se produce es que la agregación aproxima al promedio de vectores h .

Dependencia del contexto

Embeddings basados en modelos de lenguaje (ELMo)

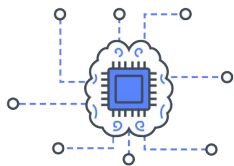
ELMo: se puede combinar con otras representaciones según la tarea específica:



El modelo final de ELMo (pre-entrenado) usa dos capas LSTM con vectores de dimensionalidad 512. Para la capa de representación, ELMo usa char n-grams y una capa de proyección de dimensionalidad 512.

ELMo fue entrenado sobre el **1 Billion Word Benchmark**.

Fine tuning para downstream tasks implica disminución en *perplexity* y mejoras en desempeño.



Ciprian Chelba, Tomas Mikolov, Mike Schuster, Qi Ge, Thorsten Brants, Phillipp Koehn, and Tony Robinson. 2014. One billion word benchmark for measuring progress in statistical language modeling. In INTERSPEECH.