

INTRODUCCIÓN A NLP CON APRENDIZAJE PROFUNDO

Parte 2

Pablo Martínez Olmos, pamartin@ing.uc3m.es

Redes neuronales para secuencias

Motivación: predecir la siguiente palabra de un texto

“This morning I took my cat for a walk.”

given these words

predict the
next word

© MIT 6.S191: Introduction to Deep Learning
introtodeeplearning.com

Motivación: predecir la siguiente palabra de un texto

Idea #1: use a fixed window

“This morning I took my cat **for a walk.**”

given these predict the
two words next word

One-hot feature encoding: tells us what each word is

[1 0 0 0 0 0 1 0 0 0]

for

a



prediction

© MIT 6.S191: Introduction to Deep Learning
introtodeeplearning.com

Motivación: predecir la siguiente palabra de un texto

Idea #1: use a fixed window

“This morning I took my cat **for a walk.**”

given these predict the
two words next word

© MIT 6.S191: Introduction to Deep Learning
introtodeeplearning.com

One-hot feature encoding: tells us what each word is

[1 0 0 0 0 0 1 0 0 0]

for

a



prediction

Problem #1: can't model long-term dependencies

“France is where I grew up, but I now live in Boston. I speak fluent ____.”

Motivación: predecir la siguiente palabra de un texto

“This morning I took my cat for a walk.”

given these words

predict the
next word

© MIT 6.S191: Introduction to Deep Learning
introtodeeplearning.com

Motivación: predecir la siguiente palabra de un texto

Idea #2: use entire sequence as set of counts

“This morning I took my cat for a”



“bag of words”

[0 1 0 0 1 0 0 ... 0 0 1 1 0 0 0 1]



prediction

© MIT 6.S191: Introduction to Deep Learning
introtodeeplearning.com

Motivación: predecir la siguiente palabra de un texto

Idea #2: use entire sequence as set of counts

“This morning I took my cat for a”



© MIT 6.S191: Introduction to Deep Learning
introtodeeplearning.com

“bag of words”

Problem #2: counts don't preserve order



The food was good, not bad at all.

vs.

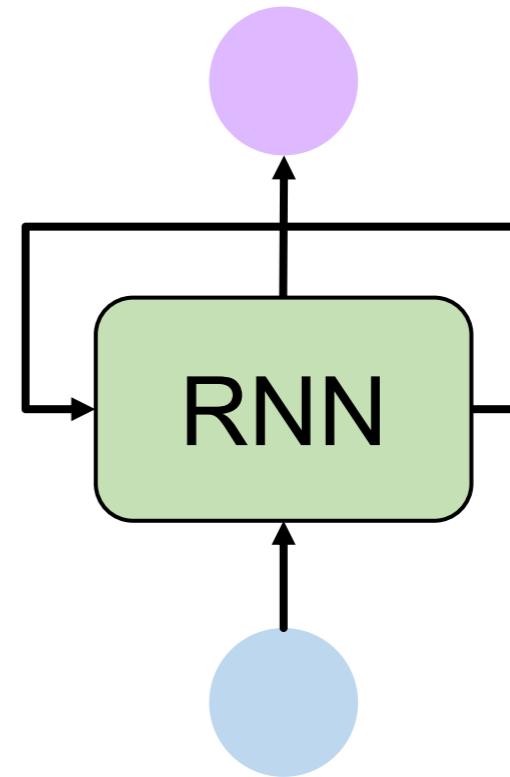
The food was bad, not good at all.



Redes neuronales recurrentes (RNNs)

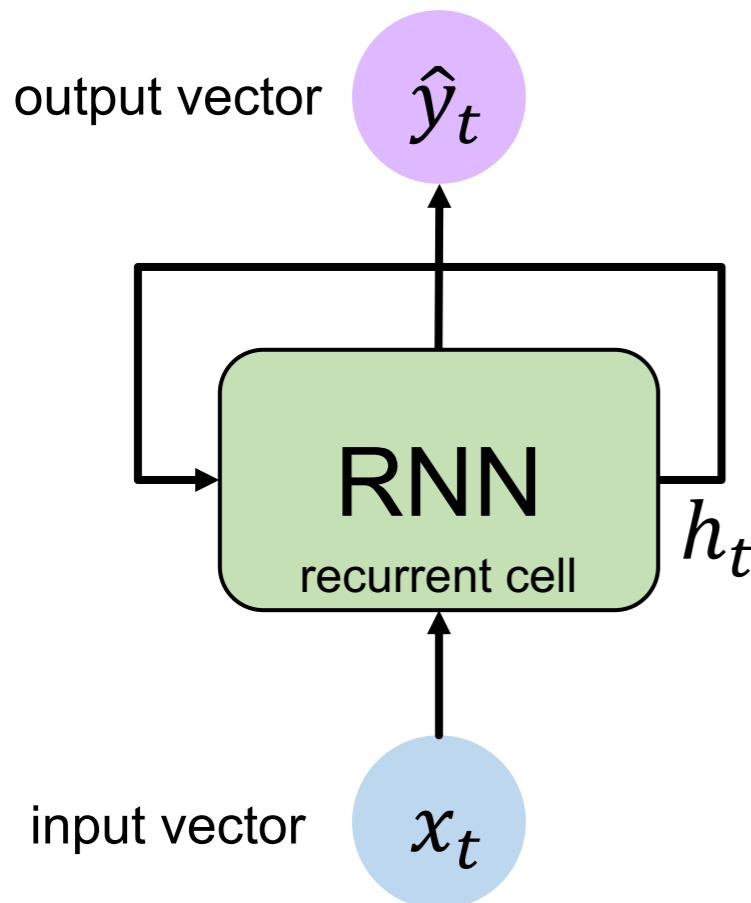
Para modelar secuencias, debemos:

1. Tratar secuencias de **longitud variable**
2. Capturar **dependencias** temporales
3. Mantener información del **orden**



© MIT 6.S191: Introduction to Deep Learning
introtodeeplearning.com

Una RNN básica



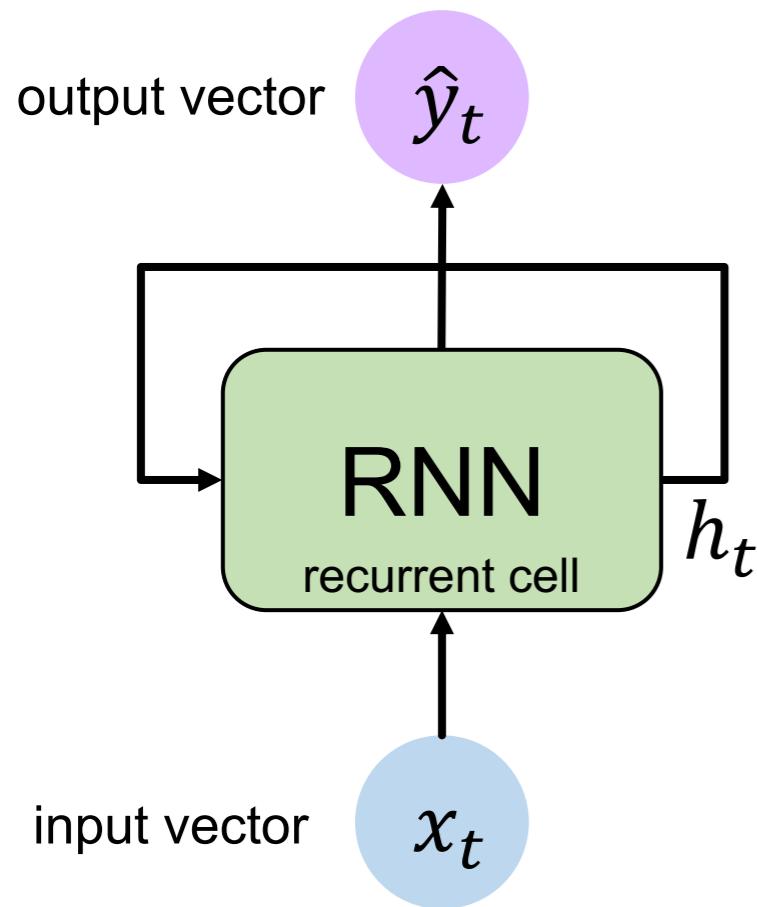
Misma **recurrencia** a lo largo del tiempo:

$$h_t = f_W(h_{t-1}, x_t)$$

La misma función y **parámetros se reutilizan** a lo largo del tiempo.

El **número de parámetros no crece con la longitud** de la secuencia.

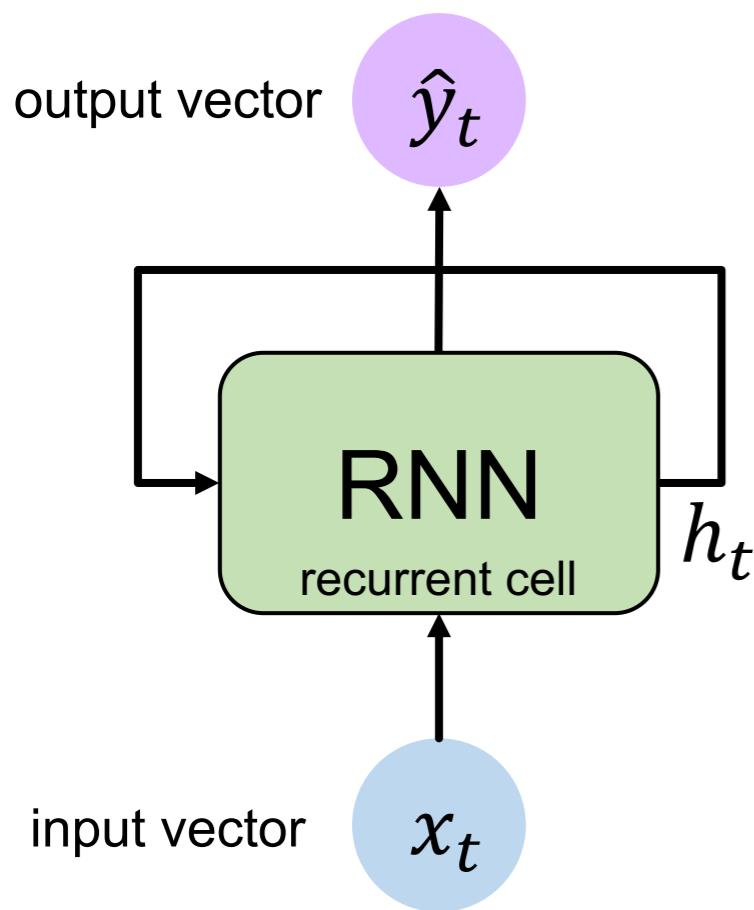
Una RNN básica



$$h_t = \tanh (\mathbf{W}_{hh} h_{t-1} + \mathbf{W}_{xh} x_t)$$

$$\hat{y}_t = \mathbf{W}_{hy} h_t$$

Una RNN básica

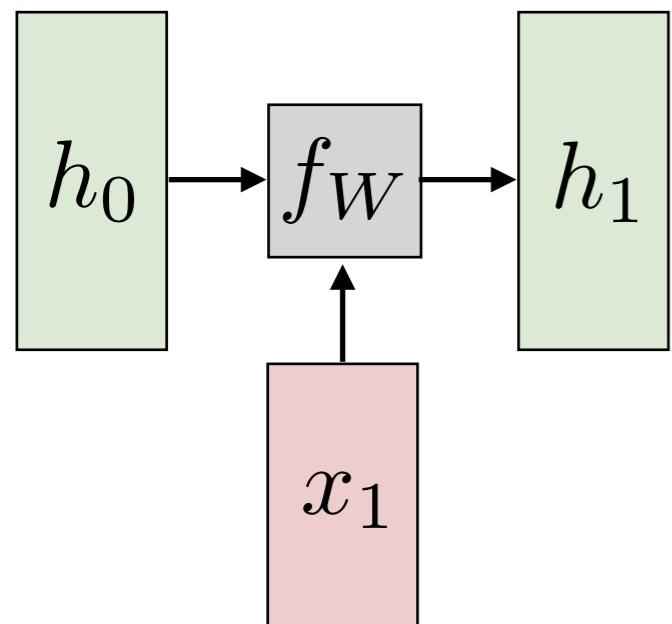


$$h_t = \tanh (\mathbf{W}_{hh} h_{t-1} + \mathbf{W}_{xh} x_t)$$

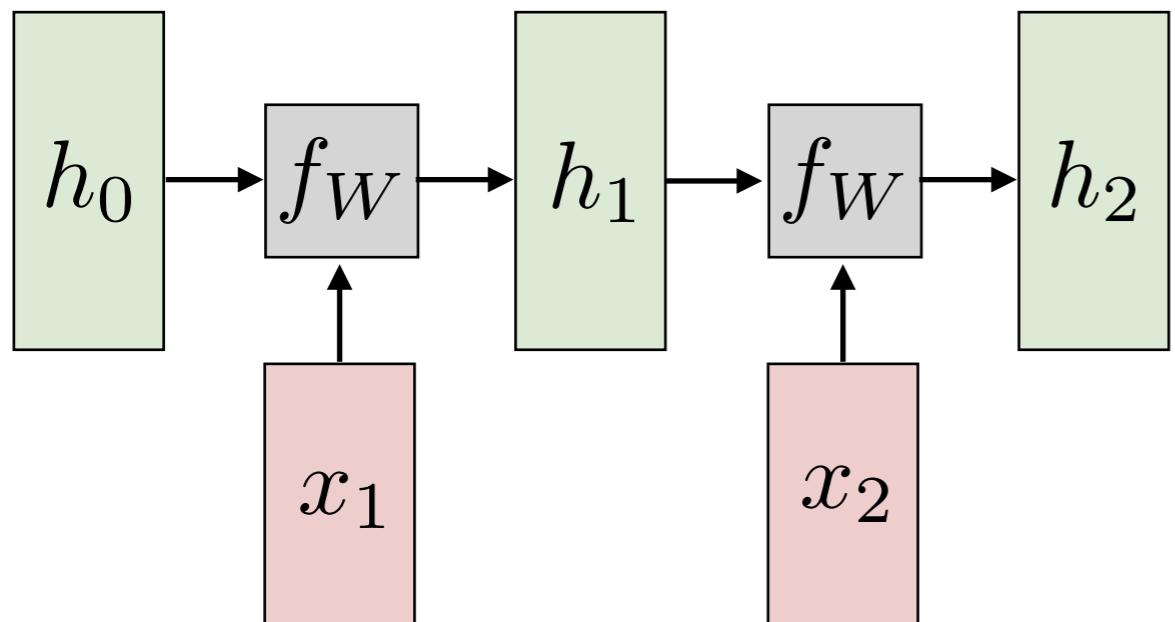
$$\hat{y}_t = \mathbf{W}_{hy} h_t$$

Combinar con una MLP para predicción final

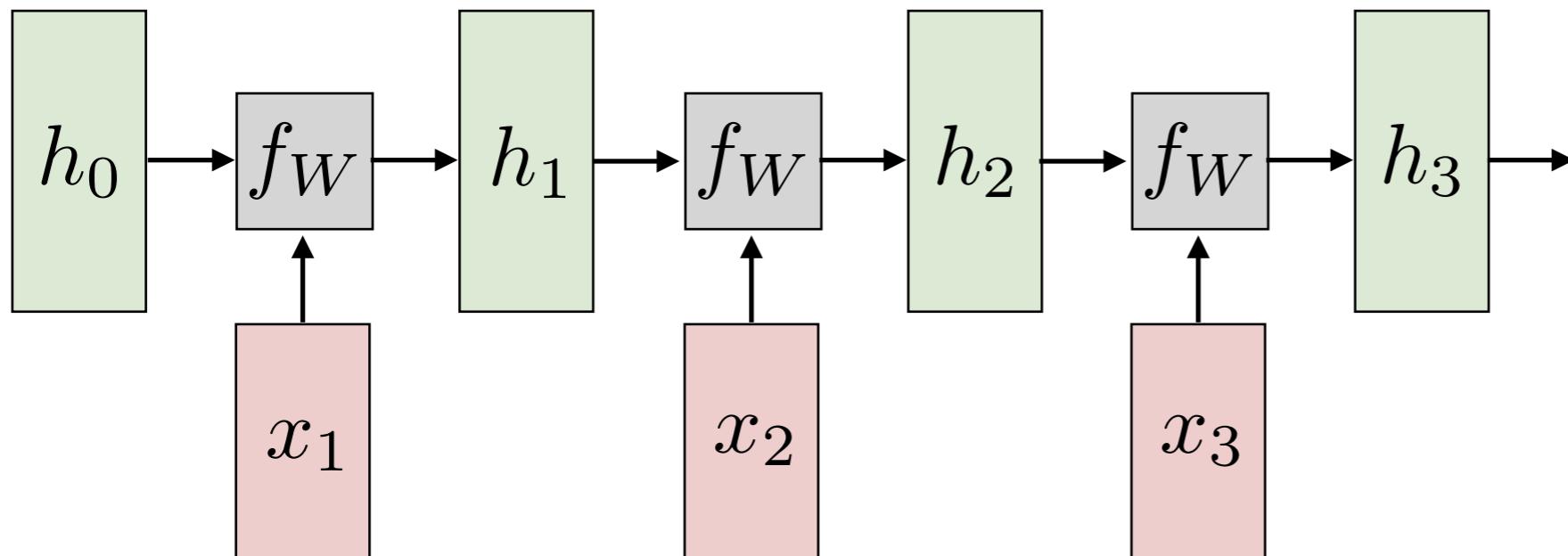
Una RNN básica



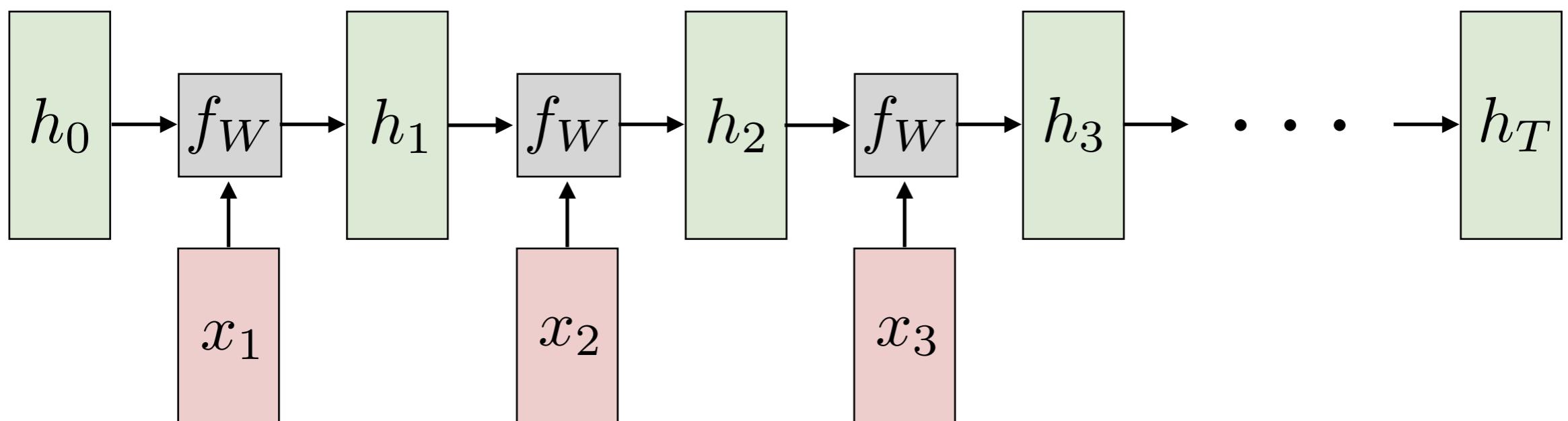
Una RNN básica



Una RNN básica

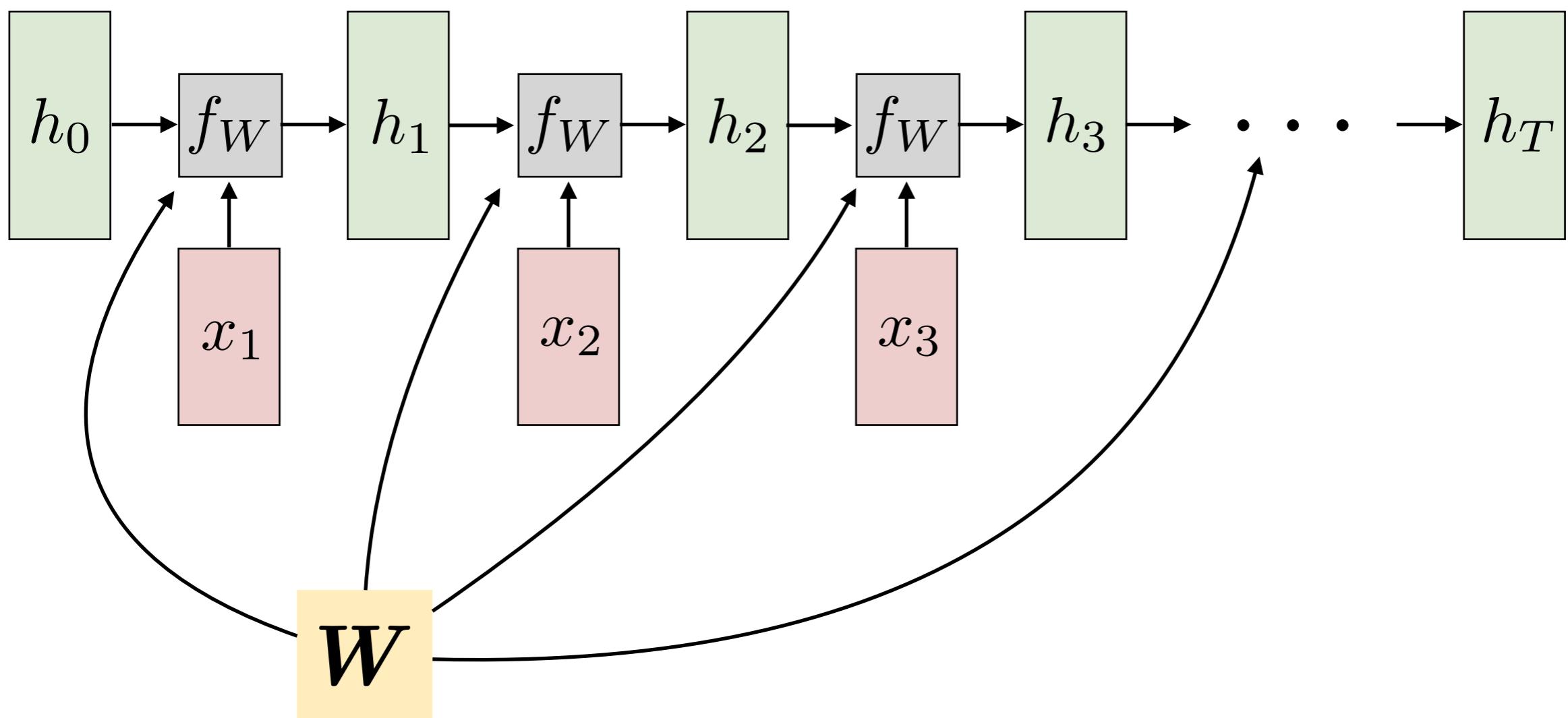


Una RNN básica



Una RNN básica

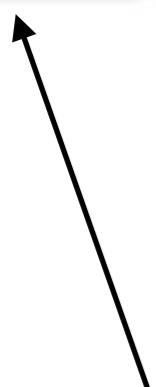
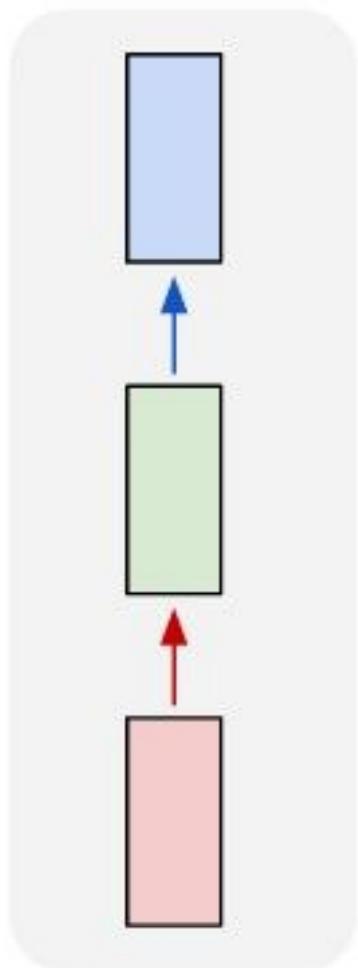
Reutilización de los mismos parámetros para procesar toda la secuencia



Variedad de estructuras utilizando RNNs

Variedad de estructuras utilizando RNNs

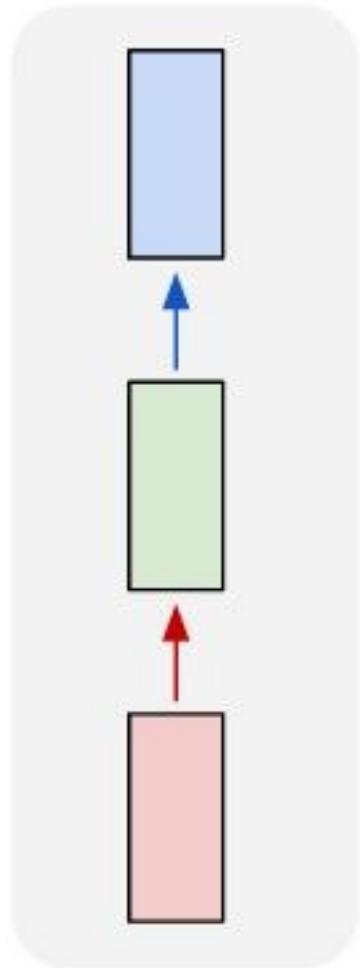
one to one



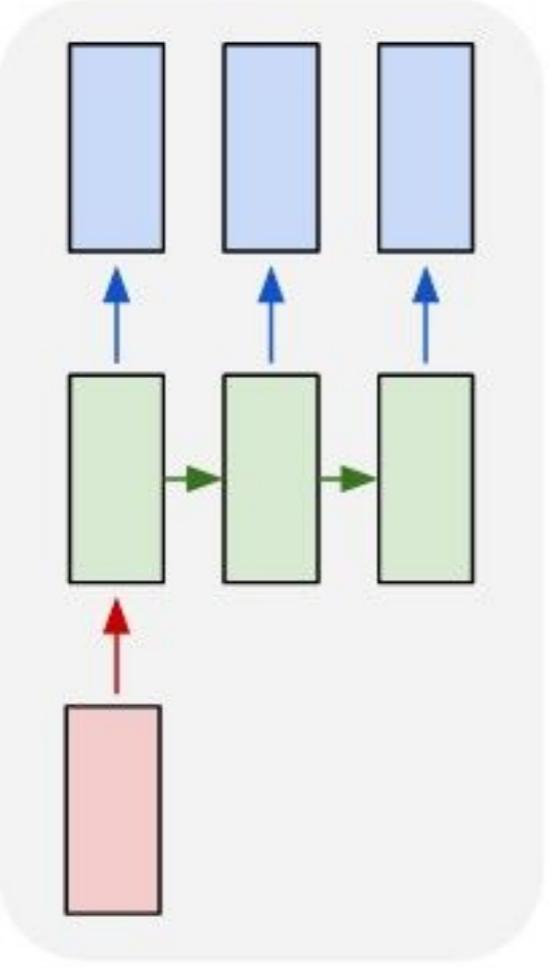
MLP

Variedad de estructuras utilizando RNNs

one to one



one to many

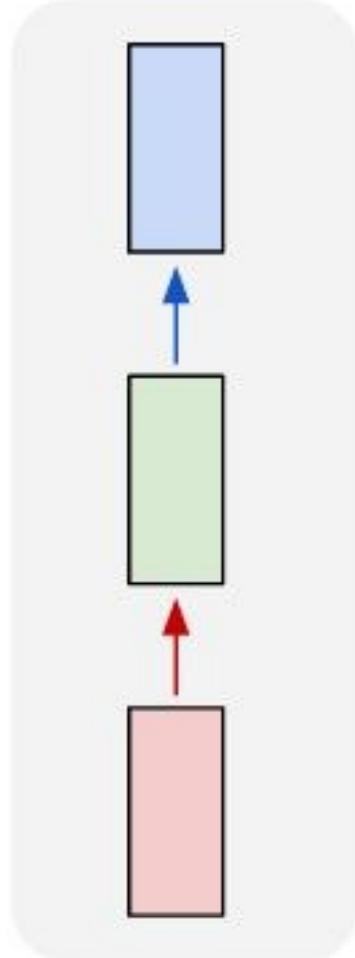


e.g. **Image Captioning**
image -> sequence of words

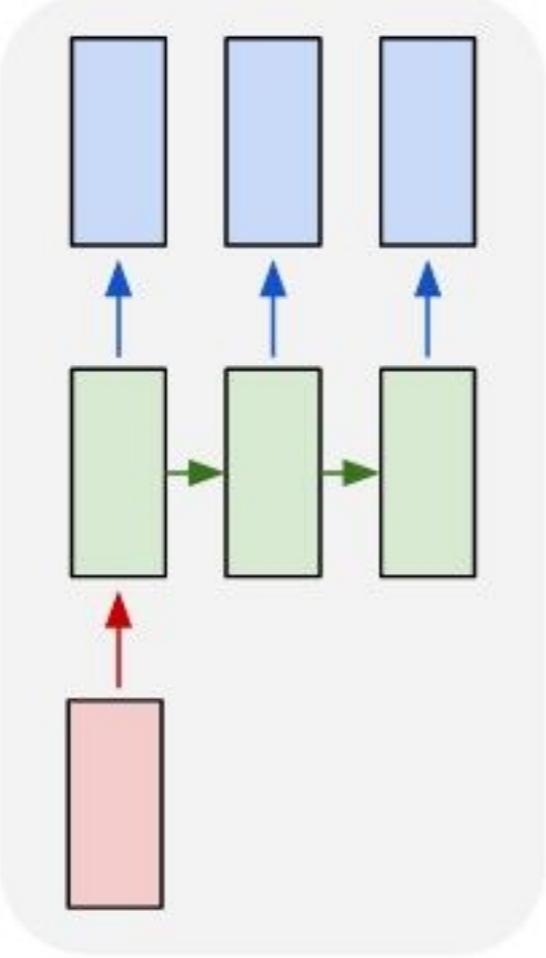
MLP

Variedad de estructuras utilizando RNNs

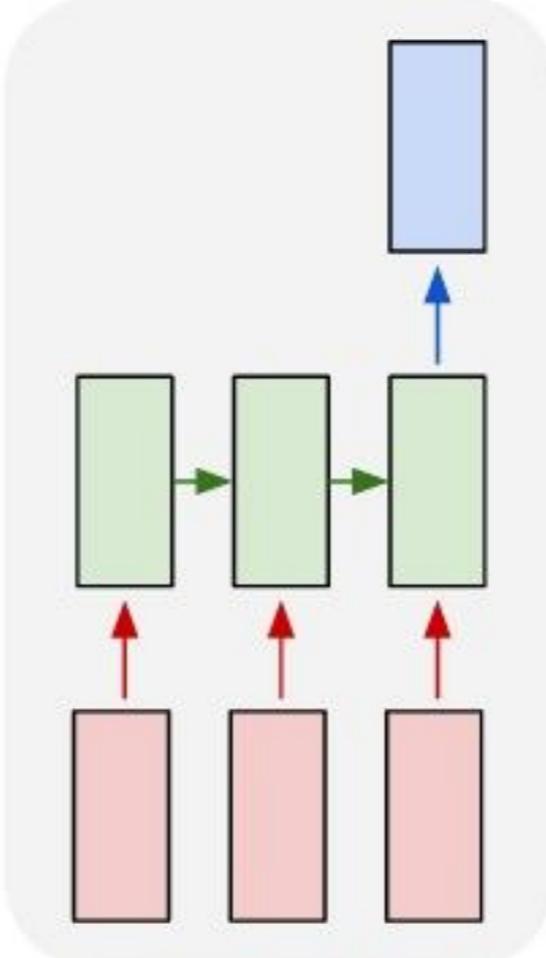
one to one



one to many



many to one



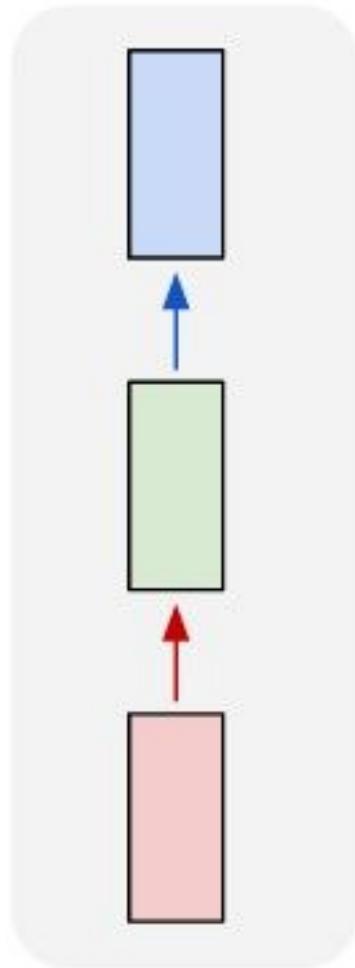
MLP

e.g. **Image Captioning**
image -> sequence of words

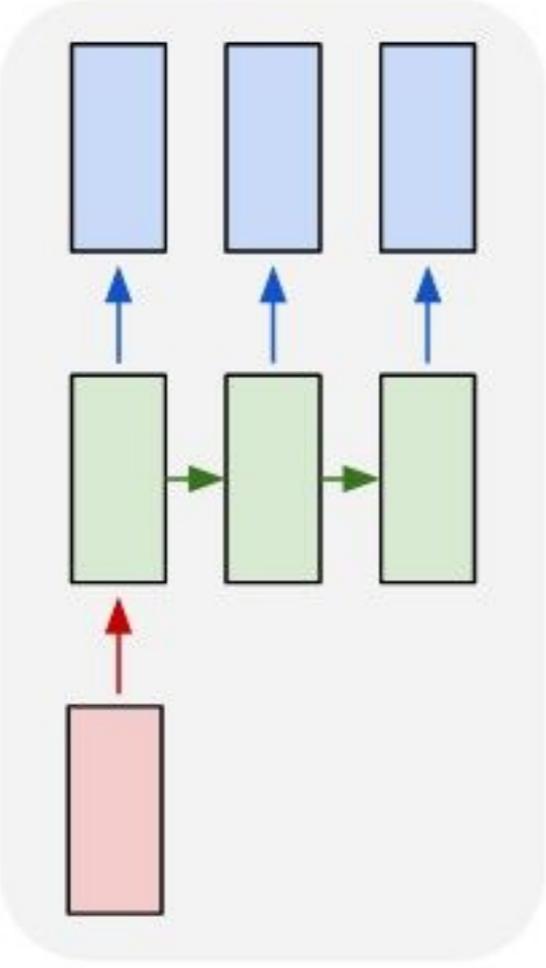
e.g. **Sentiment Classification**
sequence of words -> sentiment

Variedad de estructuras utilizando RNNs

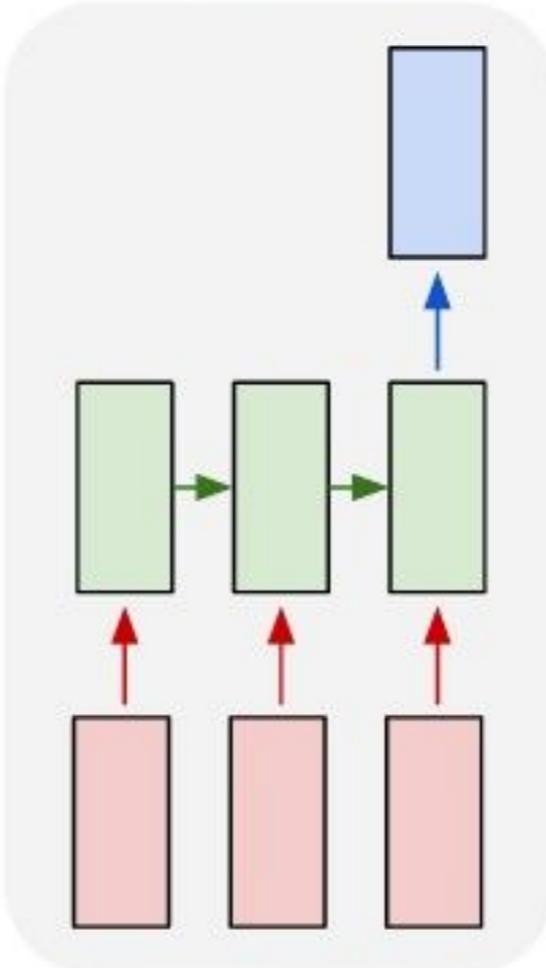
one to one



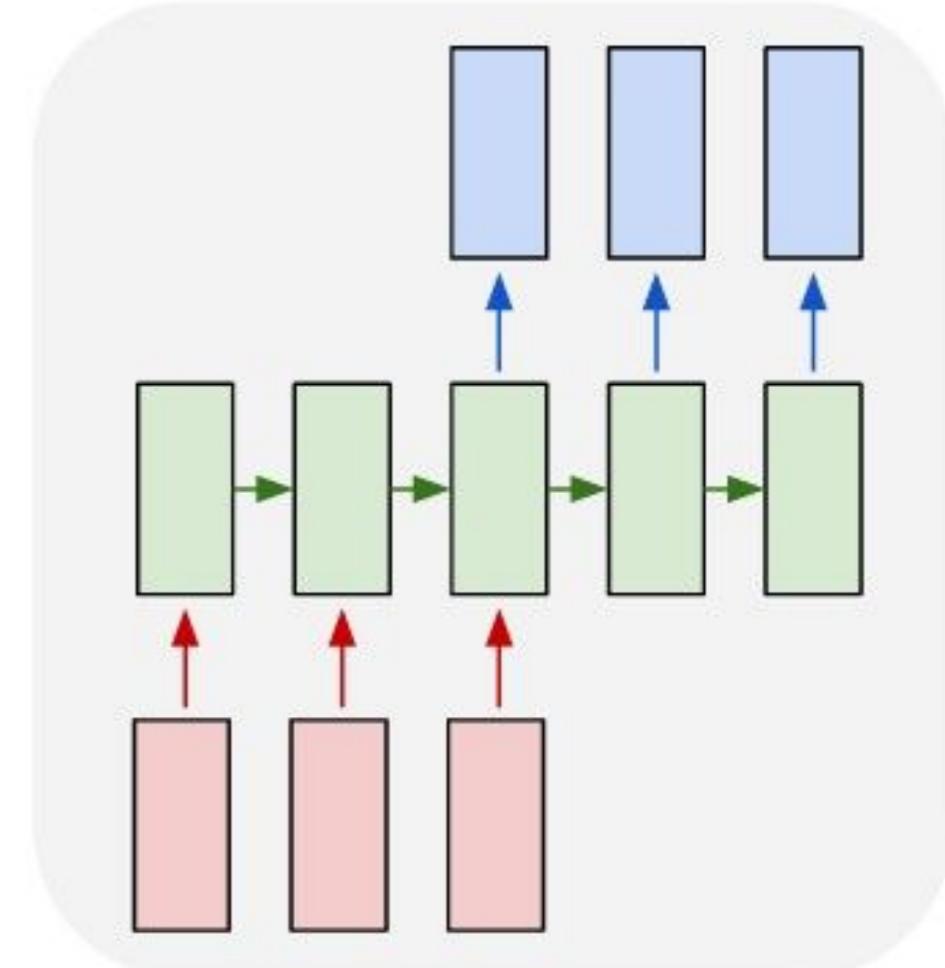
one to many



many to one



many to many



MLP

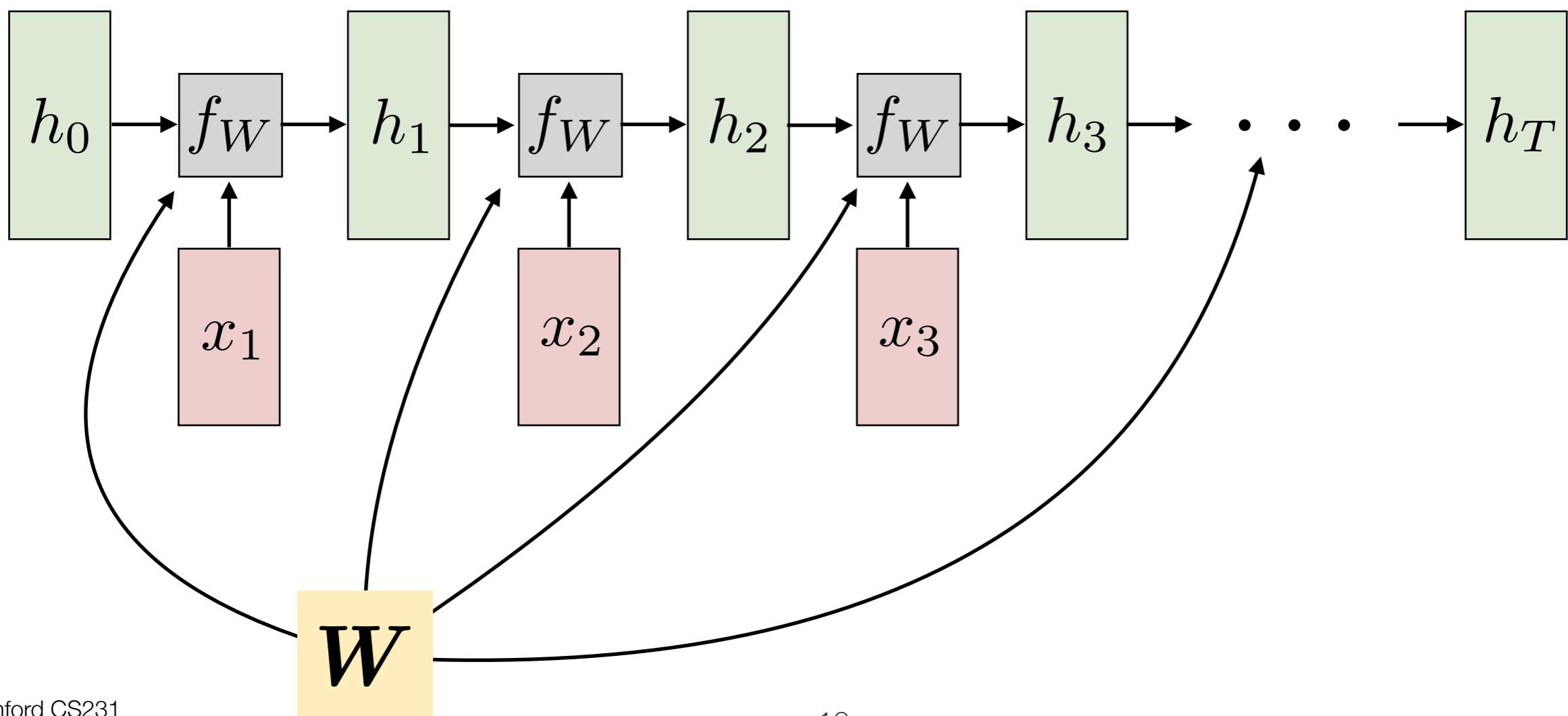
e.g. **Image Captioning**
image -> sequence of words

e.g. **Sentiment Classification**
sequence of words -> sentiment

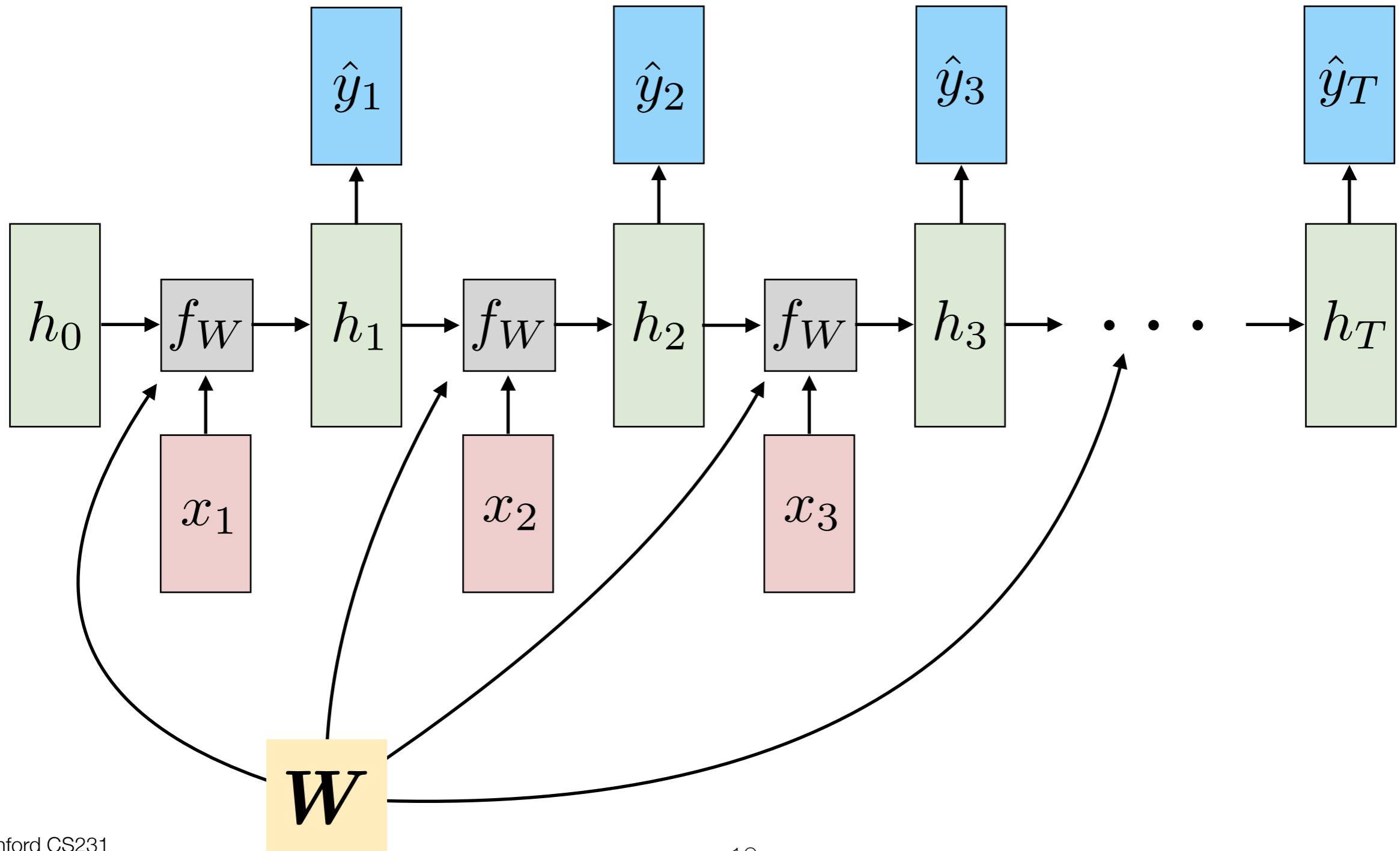
e.g. **Machine Translation**
seq. of words -> seq. of words

Múltiples predicciones a lo largo de la secuencia

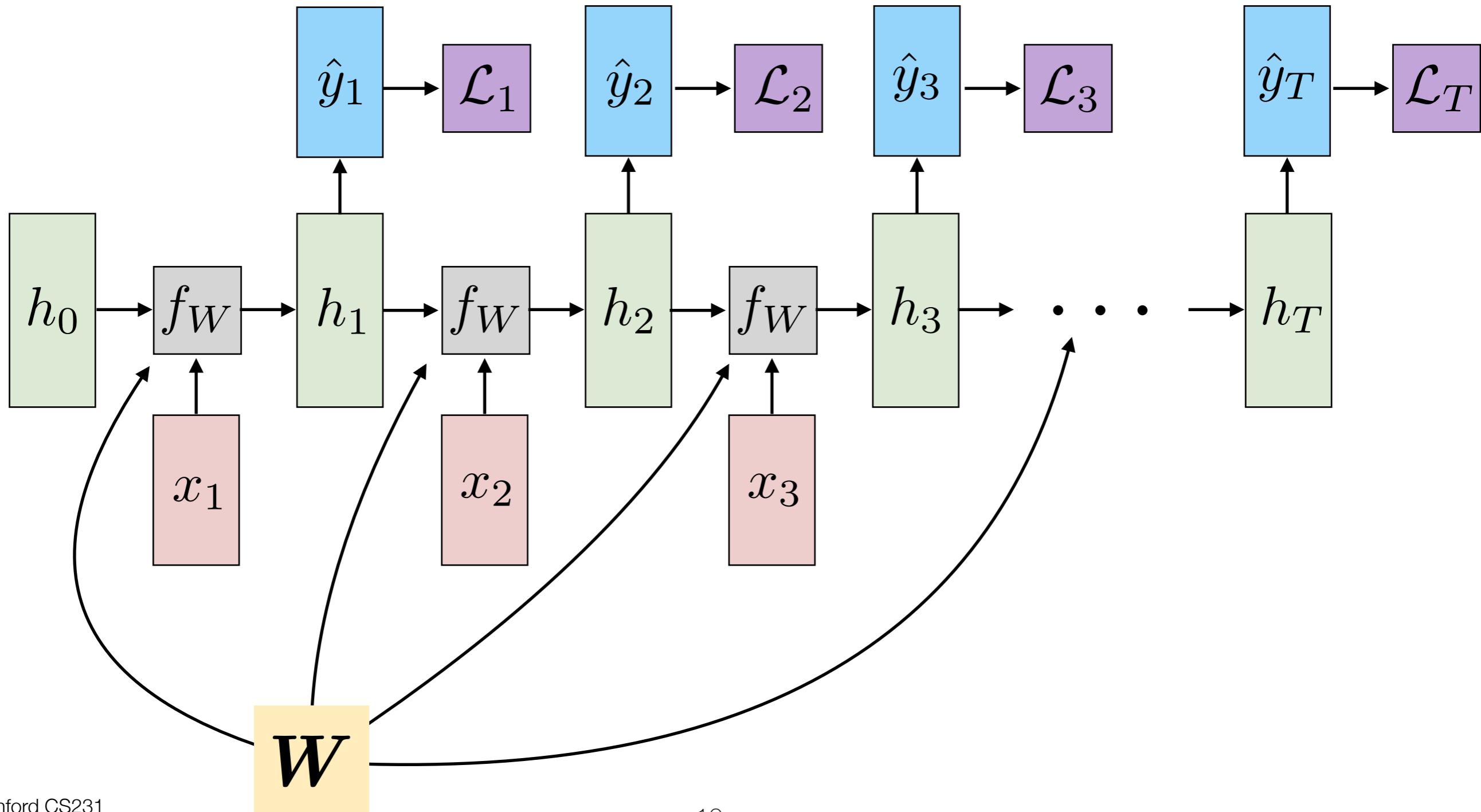
Múltiples predicciones a lo largo de la secuencia



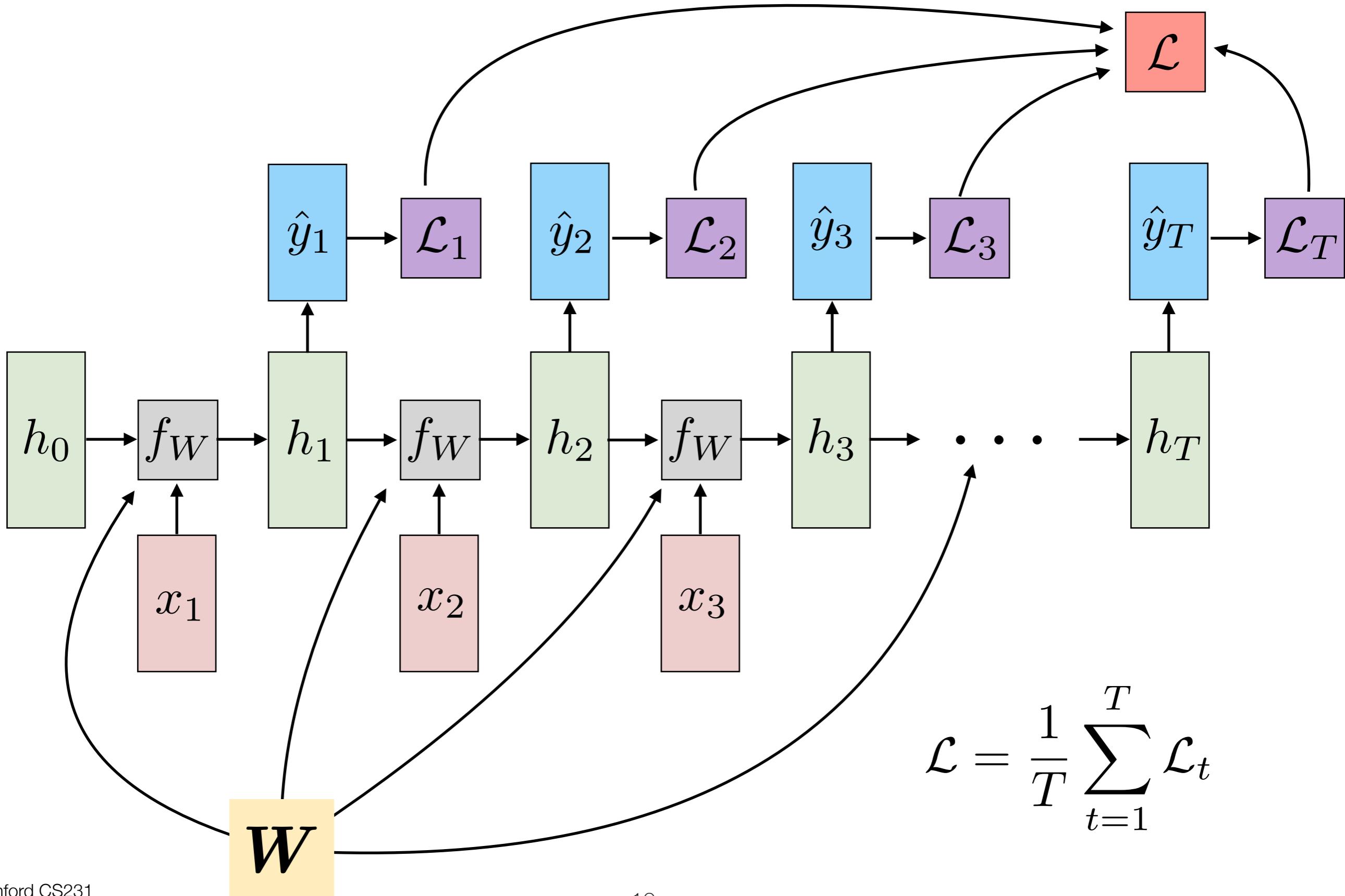
Múltiples predicciones a lo largo de la secuencia



Múltiples predicciones a lo largo de la secuencia

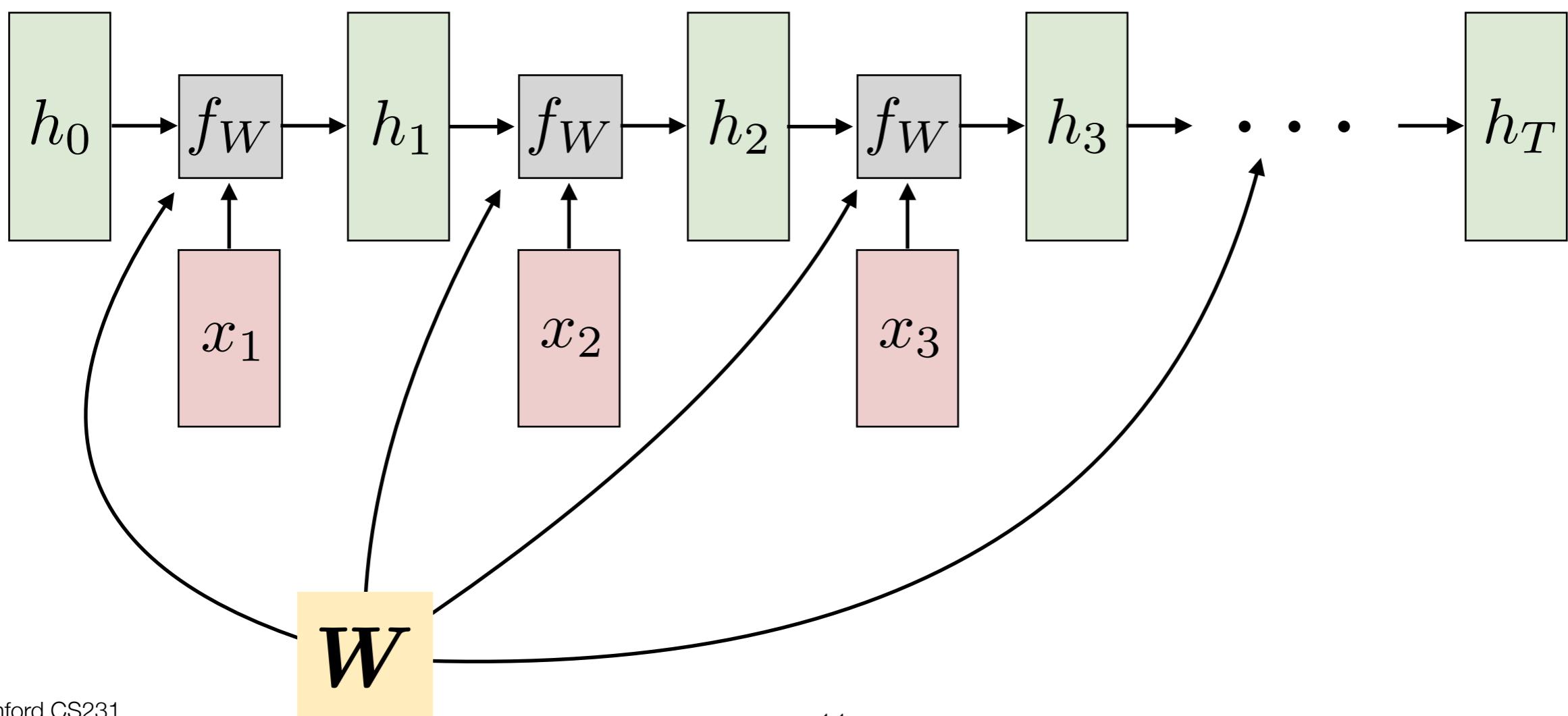


Múltiples predicciones a lo largo de la secuencia

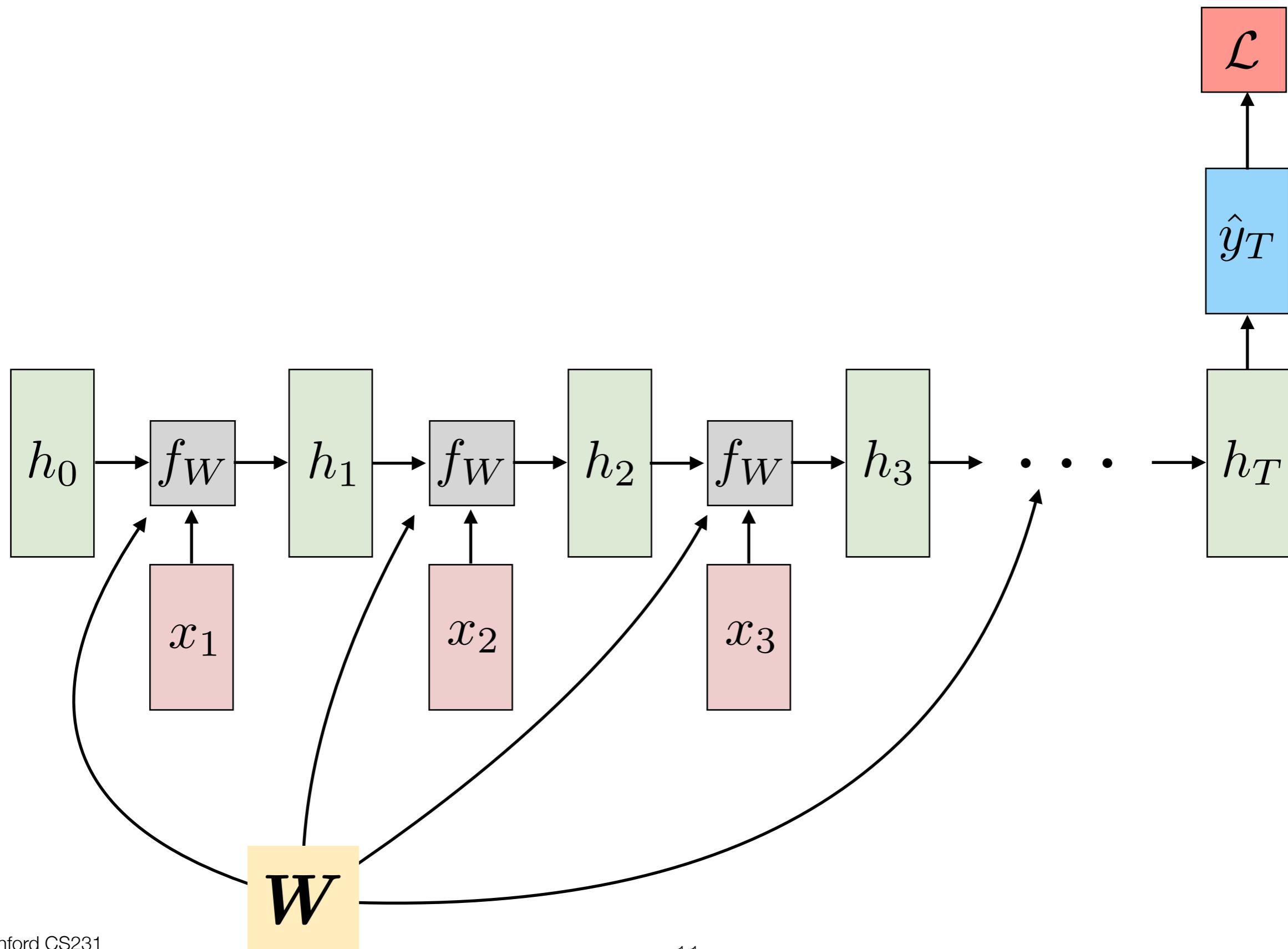


Única predicción al final de la secuencia

Única predicción al final de la secuencia

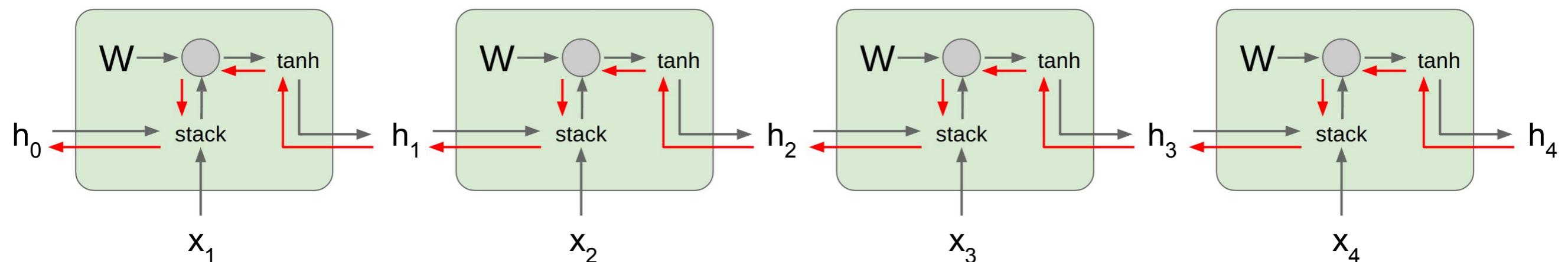


Única predicción al final de la secuencia



Limitaciones RNN básica

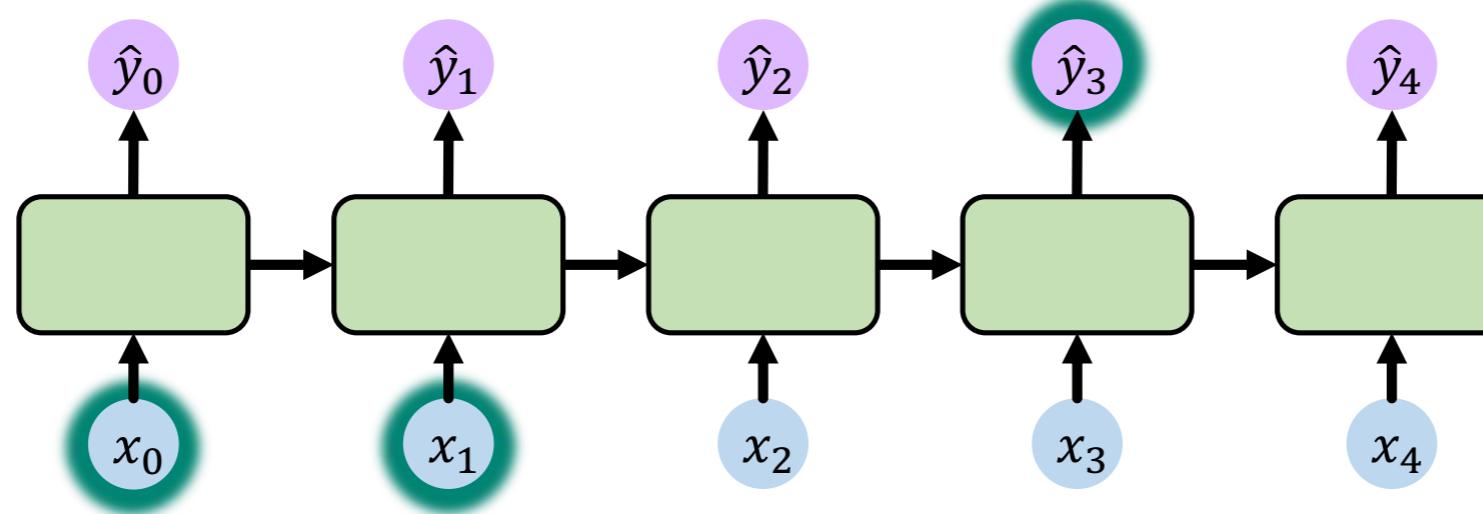
Un único modo de operación en el que cada nueva entrada modifica el estado



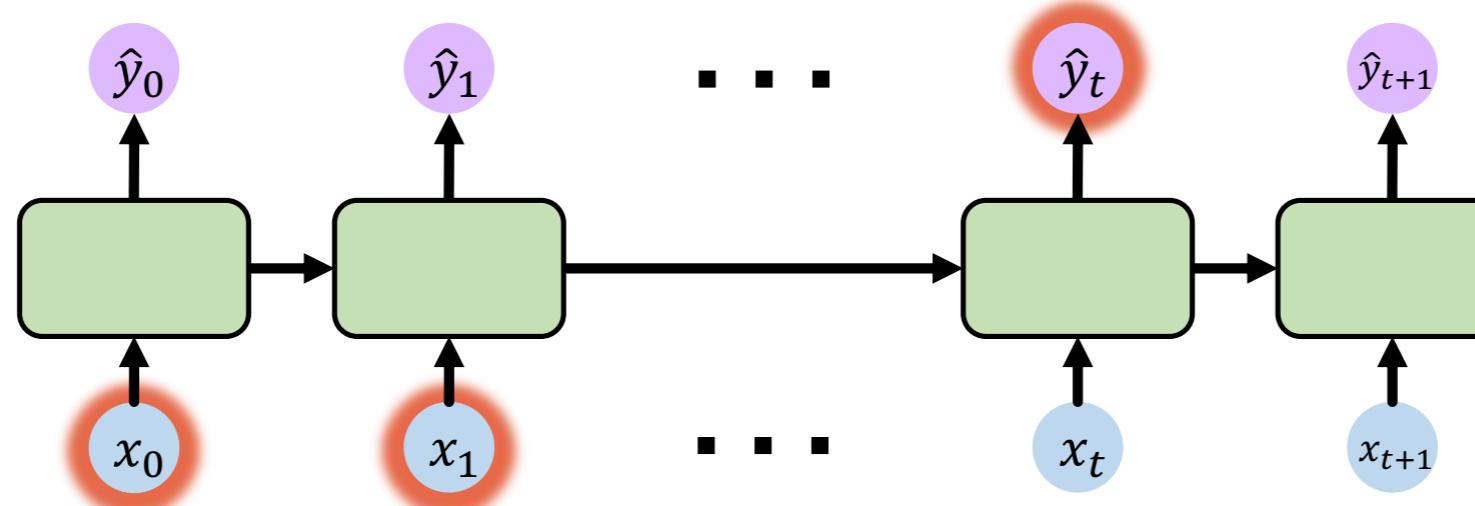
La red carece de mecanismos para recordar eventos pasados o ignorar observaciones presentes

Limitaciones RNN básica

“The clouds are in the ___”



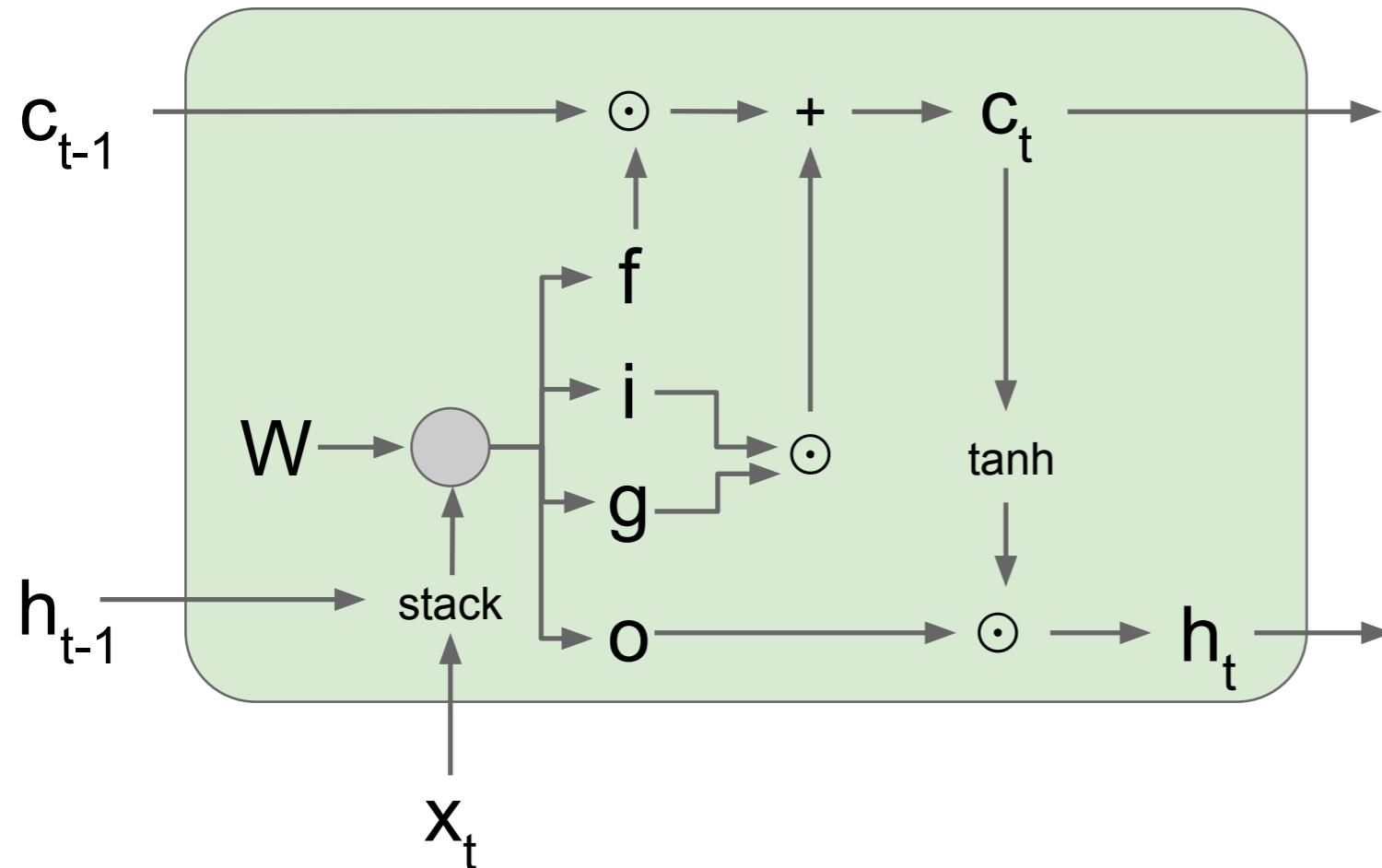
“I grew up in France, ... and I speak fluent ___”



Long Short Term Memory (LSTMs) networks

Hochreiter and Schmidhuber, Neural Computation 1997

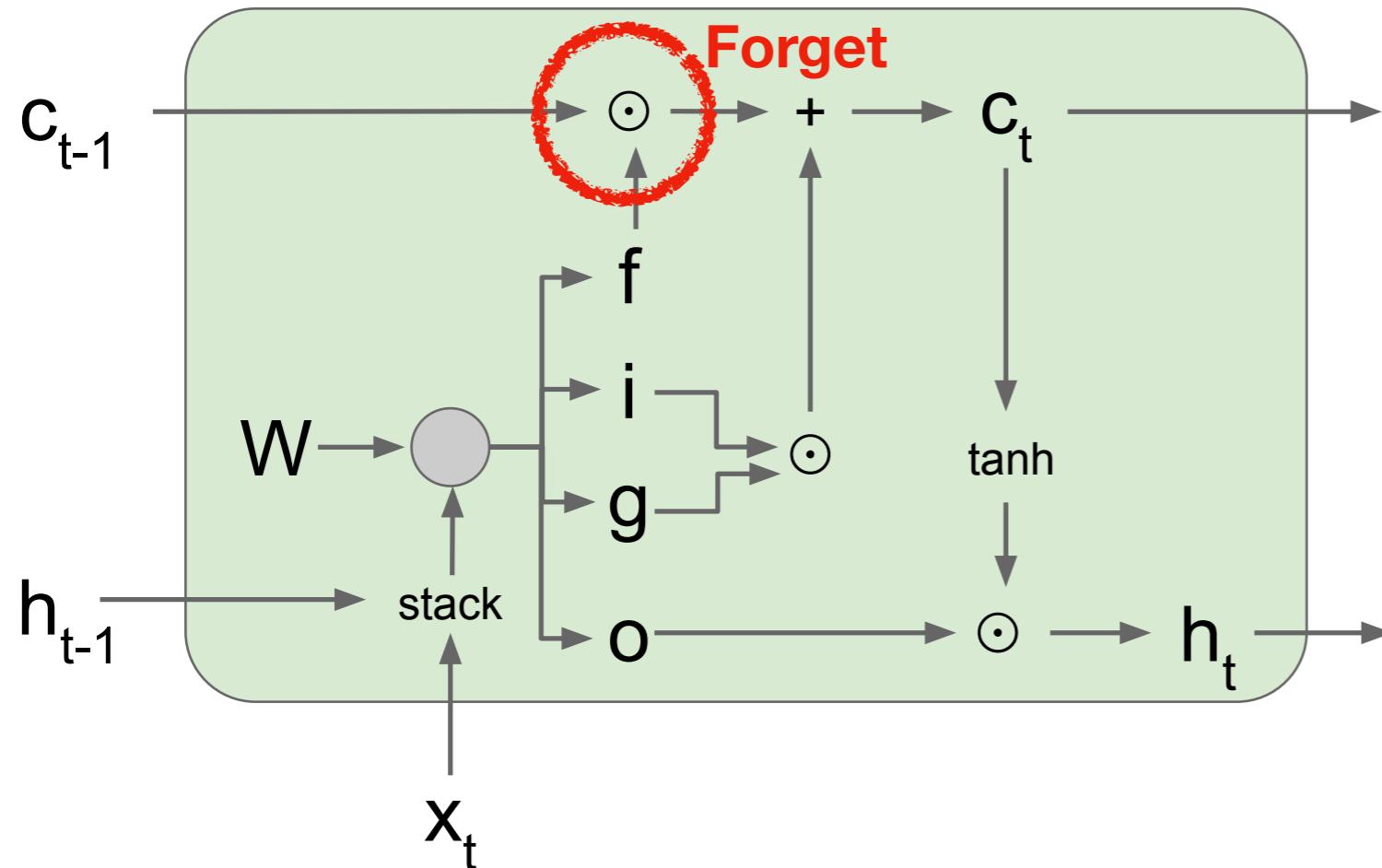
LSTMs: mecanismo de puertas



$$\begin{aligned} f &= \sigma(W_{fx}x_t + W_{fh}h_{t-1} + b_f) \\ i &= \sigma(W_{xi}x_t + W_{hi}h_{t-1} + b_i) \\ o &= \sigma(W_{ox}x_t + W_{oh}h_{t-1} + b_o) \\ g &= \tanh(W_{cx}x_t + W_{ch}h_{t-1} + b_c) \\ c_t &= f \circ c_{t-1} + i \circ \tilde{c}_t \\ h_t &= o \circ \tanh(c_t) \end{aligned}$$

Las puertas controlan la propagación de la información a lo largo del tiempo

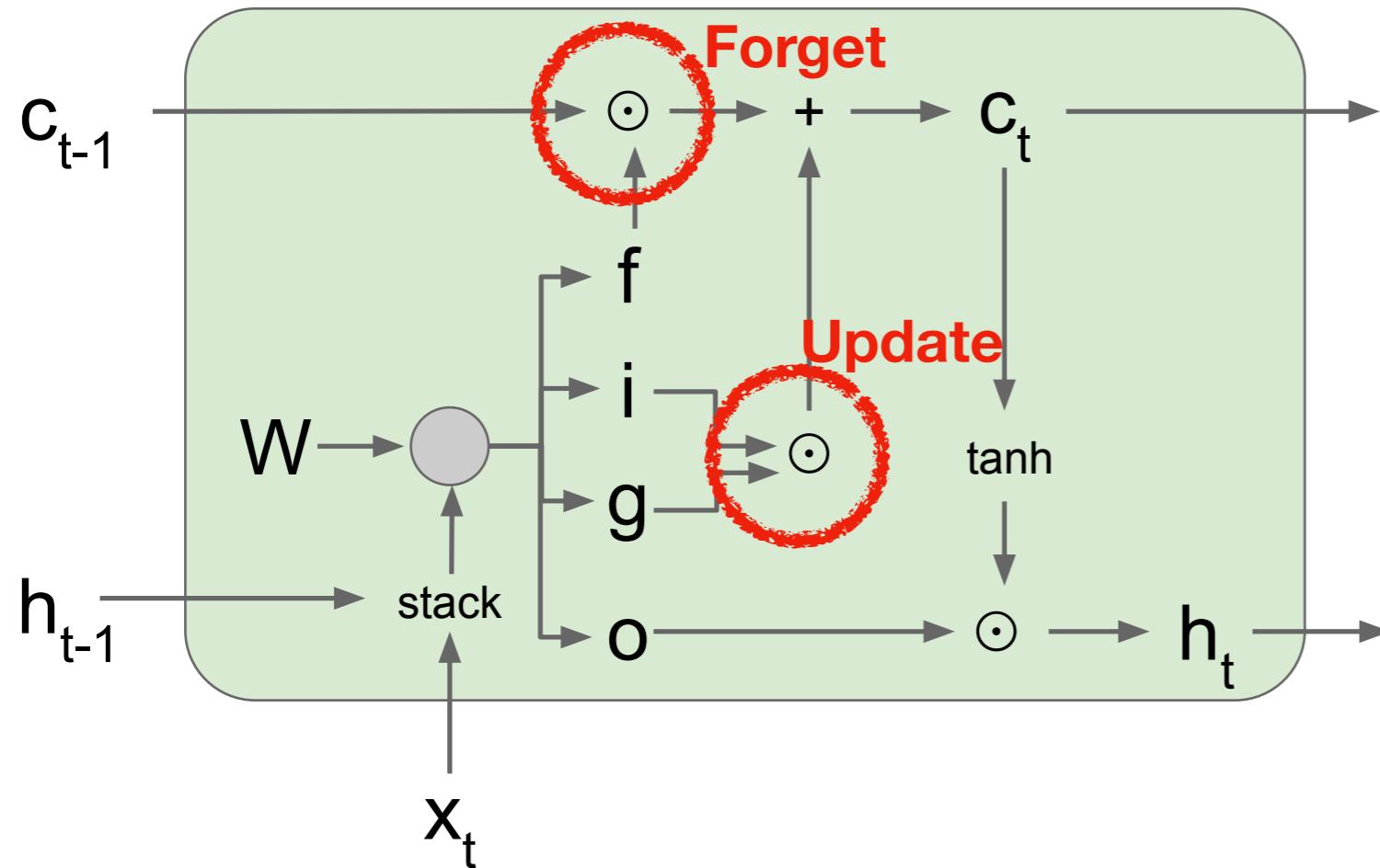
LSTMs: mecanismo de puertas



$$f = \sigma(W_{fx}x_t + W_{fh}h_{t-1} + b_f)$$
$$i = \sigma(W_{xi}x_t + W_{hi}h_{t-1} + b_i)$$
$$o = \sigma(W_{ox}x_t + W_{oh}h_{t-1} + b_o)$$
$$g = \tanh(W_{cx}x_t + W_{ch}h_{t-1} + b_c)$$
$$c_t = f \circ c_{t-1} + i \circ \tilde{c}_t$$
$$h_t = o \circ \tanh(c_t)$$

Las puertas controlan la propagación de la información a lo largo del tiempo

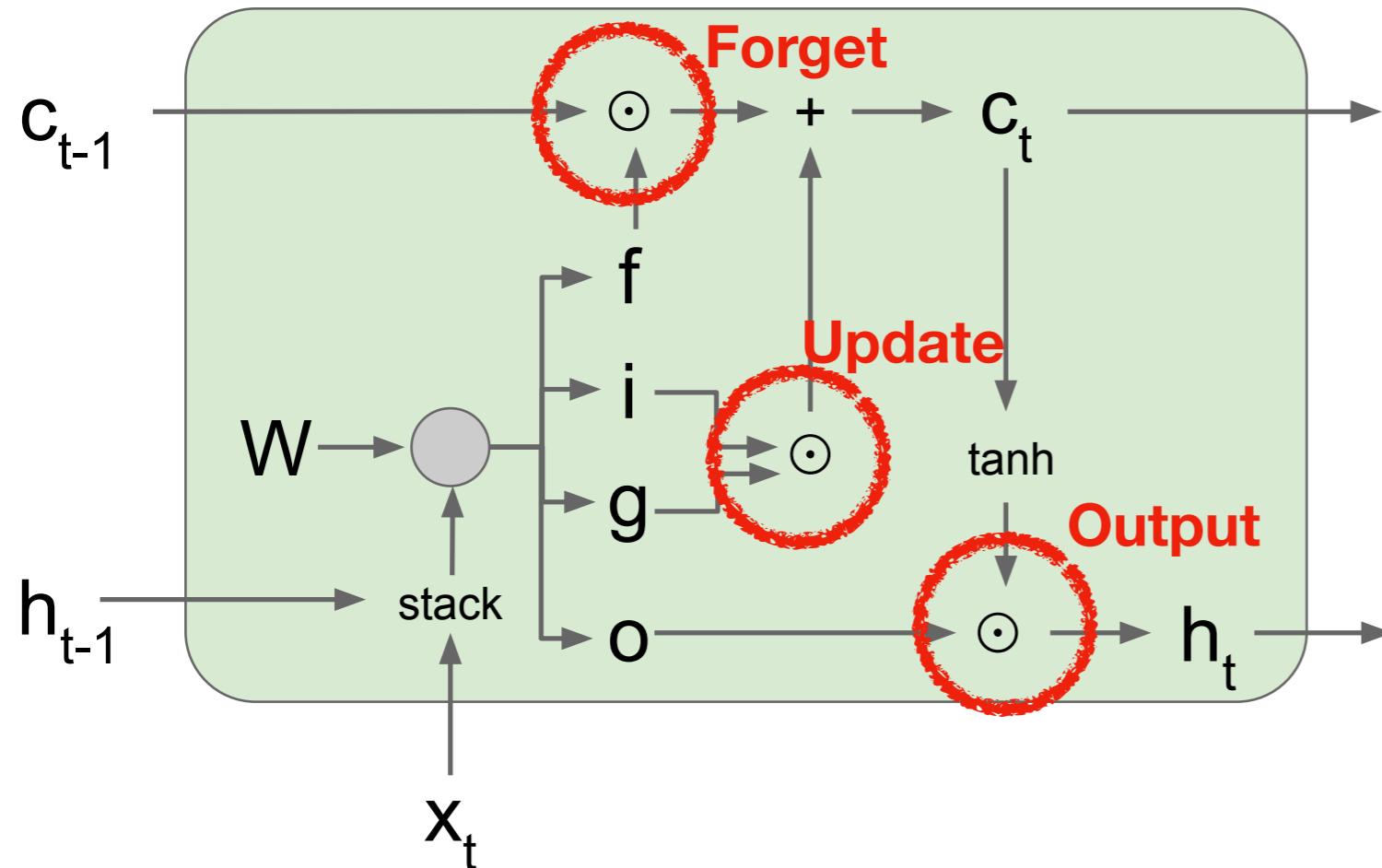
LSTMs: mecanismo de puertas



$$f = \sigma(W_{fx}x_t + W_{fh}h_{t-1} + b_f)$$
$$i = \sigma(W_{xi}x_t + W_{hi}h_{t-1} + b_i)$$
$$o = \sigma(W_{ox}x_t + W_{oh}h_{t-1} + b_o)$$
$$g = \tanh(W_{cx}x_t + W_{ch}h_{t-1} + b_c)$$
$$c_t = f \circ c_{t-1} + i \circ \tilde{c}_t$$
$$h_t = o_t \circ \tanh(c_t)$$

Las puertas controlan la propagación de la información a lo largo del tiempo

LSTMs: mecanismo de puertas

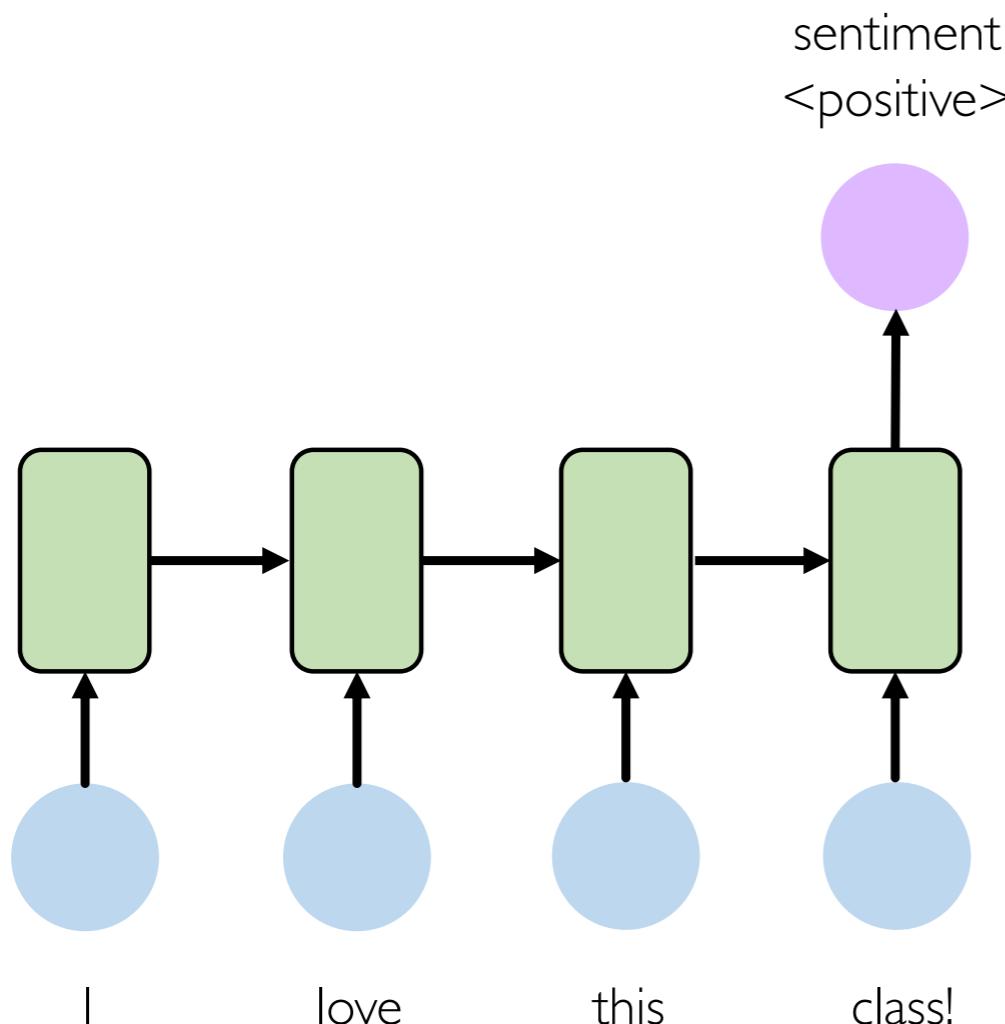


$$f = \sigma(W_{fx}x_t + W_{fh}h_{t-1} + b_f)$$
$$i = \sigma(W_{xi}x_t + W_{hi}h_{t-1} + b_i)$$
$$o = \sigma(W_{ox}x_t + W_{oh}h_{t-1} + b_o)$$
$$g = \tanh(W_{cx}x_t + W_{ch}h_{t-1} + b_c)$$
$$c_t = f \circ c_{t-1} + i \circ \tilde{c}_t$$
$$h_t = o \circ \tanh(c_t)$$

Las puertas controlan la propagación de la información a lo largo del tiempo

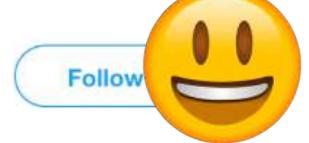
RNNs en NLP

Sentiment Analysis



Tweet sentiment classification

 **Ivar Hagendoorn**
@IvarHagendoorn



The [@MIT](#) Introduction to [#DeepLearning](#) is definitely one of the best courses of its kind currently available online introtodeeplearning.com

12:45 PM - 12 Feb 2018

 **Angels-Cave**
@AngelsCave

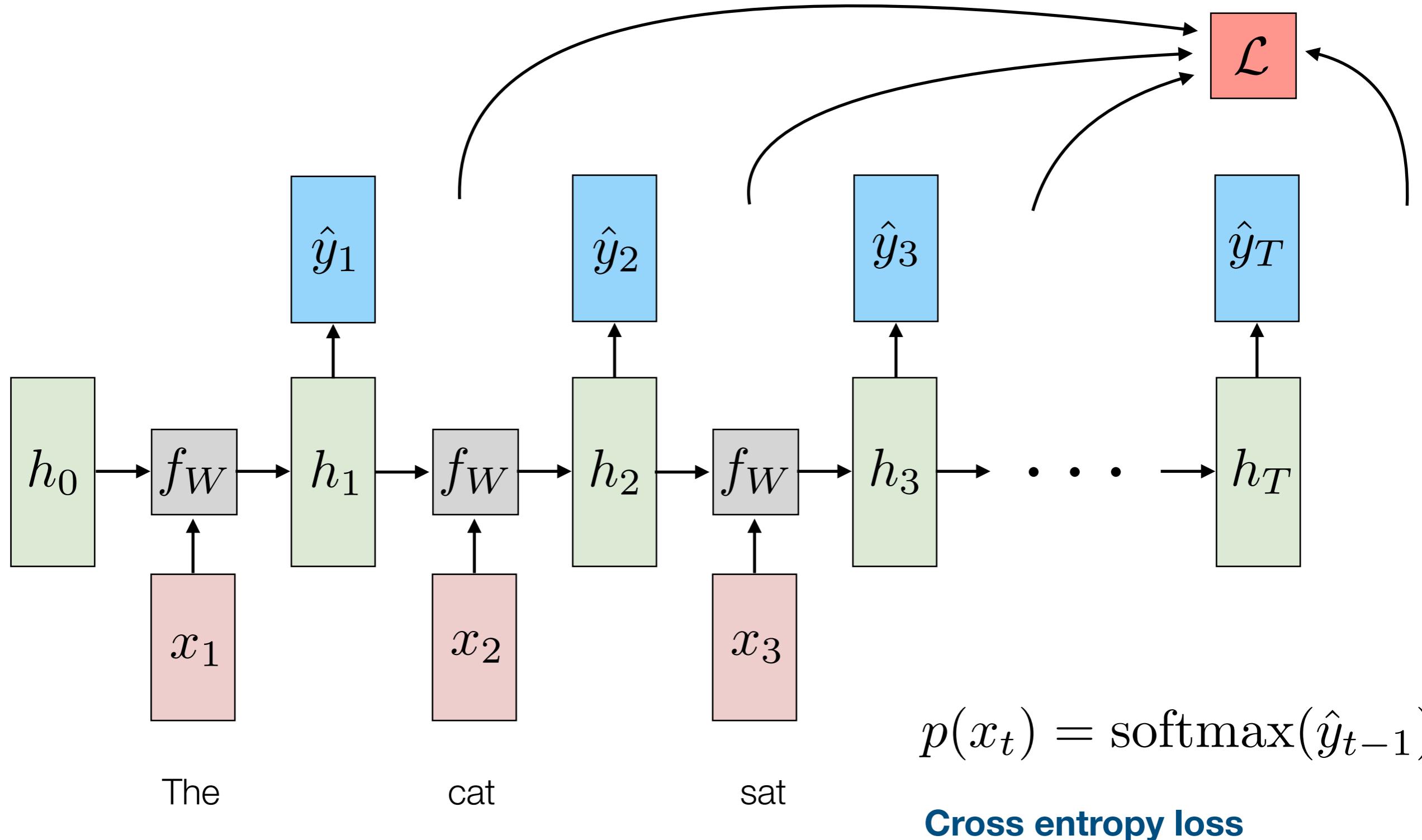


Replying to [@Kazuki2048](#)

I wouldn't mind a bit of snow right now. We haven't had any in my bit of the Midlands this winter! :(

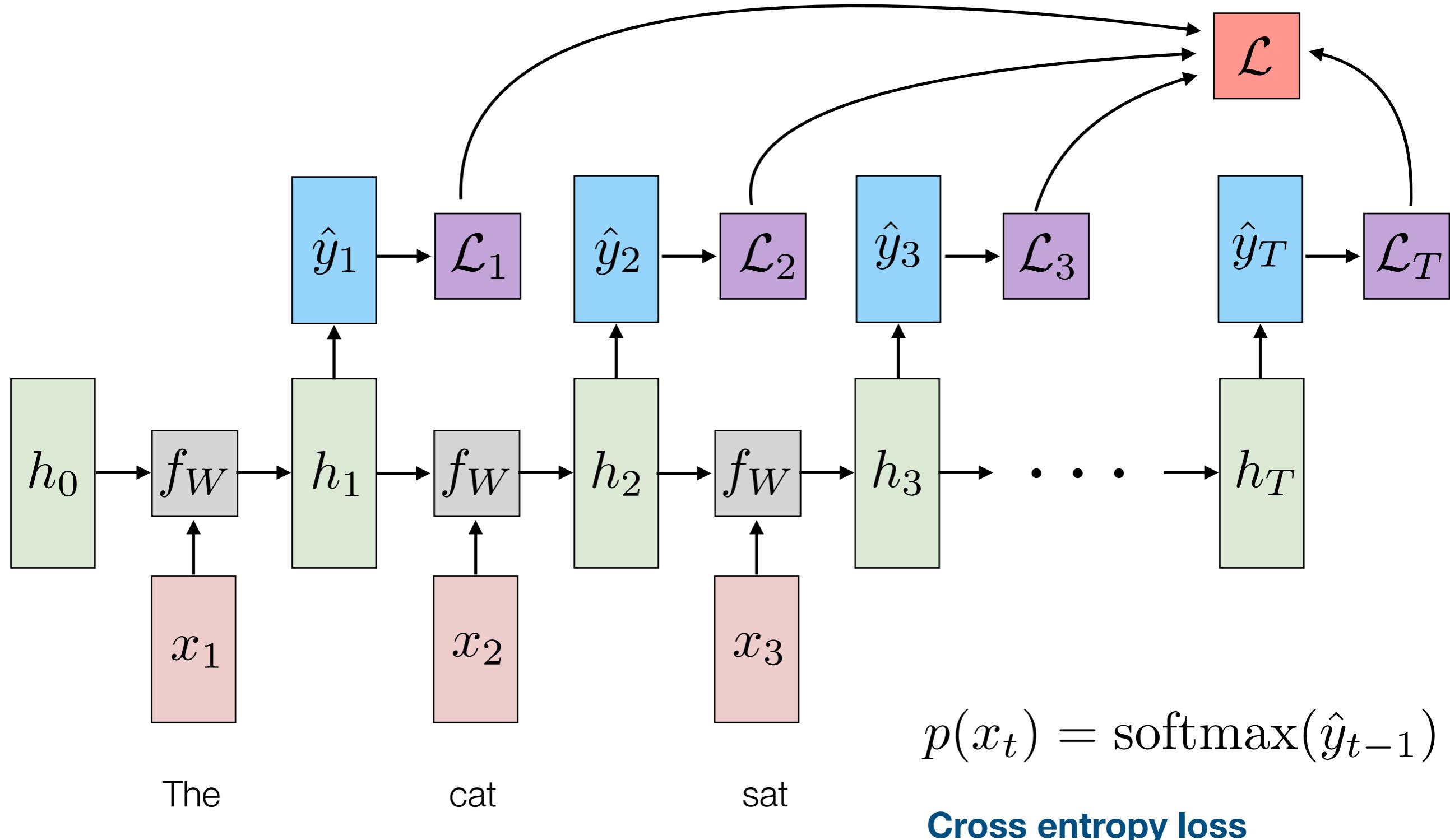
2:19 AM - 25 Jan 2019

Generación de texto con RNNs (Entrenamiento)



Tomas Mikolov, Martin Karafiat, Lukas Burget, Jan Cernocky, and Sanjeev Khudanpur. 2010.
Recurrent Neural Network Based Language Model. InProceedings of INTERSPEECH

Generación de texto con RNNs (Entrenamiento)



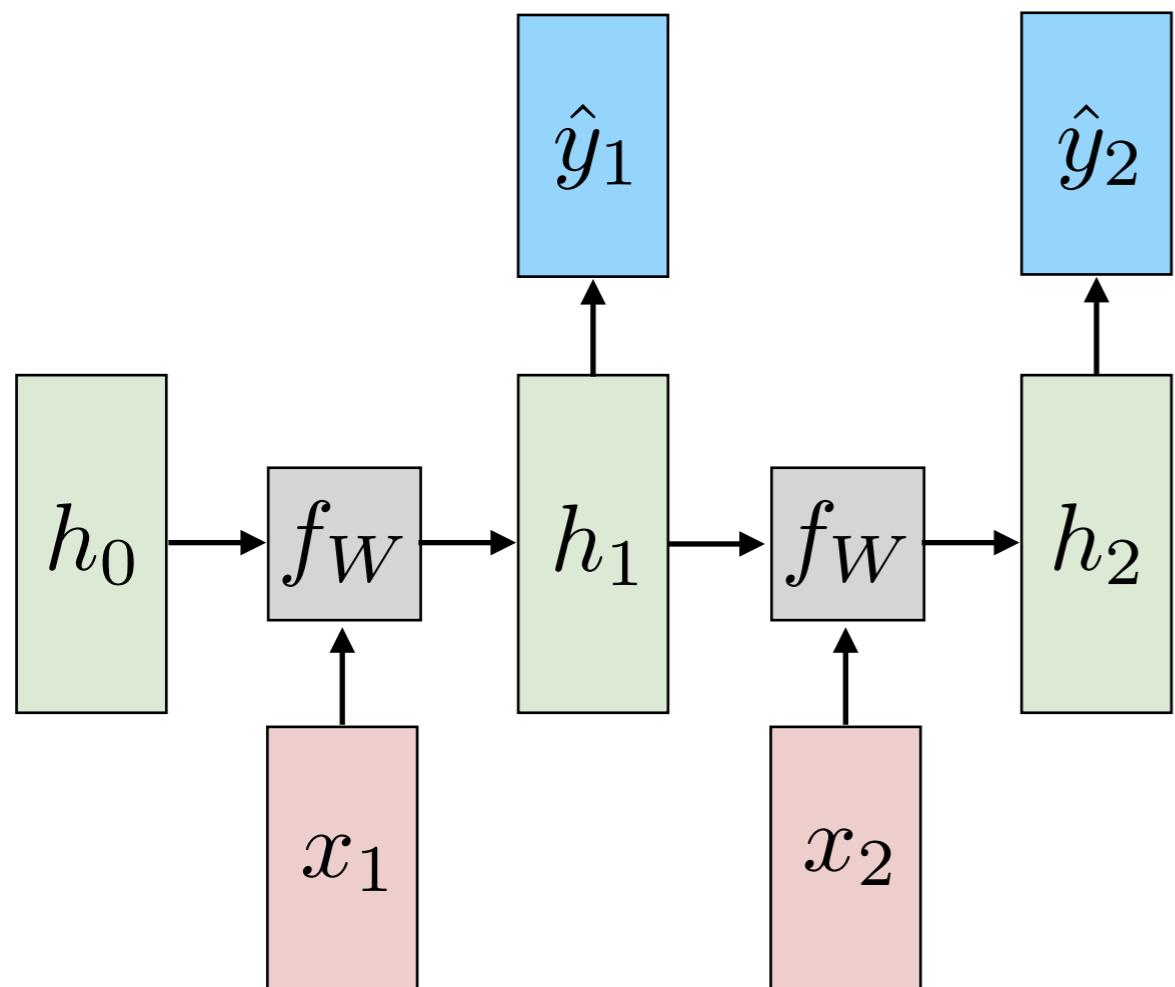
Tomas Mikolov, Martin Karafiat, Lukas Burget, Jan Cernocky, and Sanjeev Khudanpur. 2010.
Recurrent Neural Network Based Language Model. InProceedings of INTERSPEECH

Generación de texto con RNNs (Test)

Muestreo secuencial ...

Generación de texto con RNNs (Test)

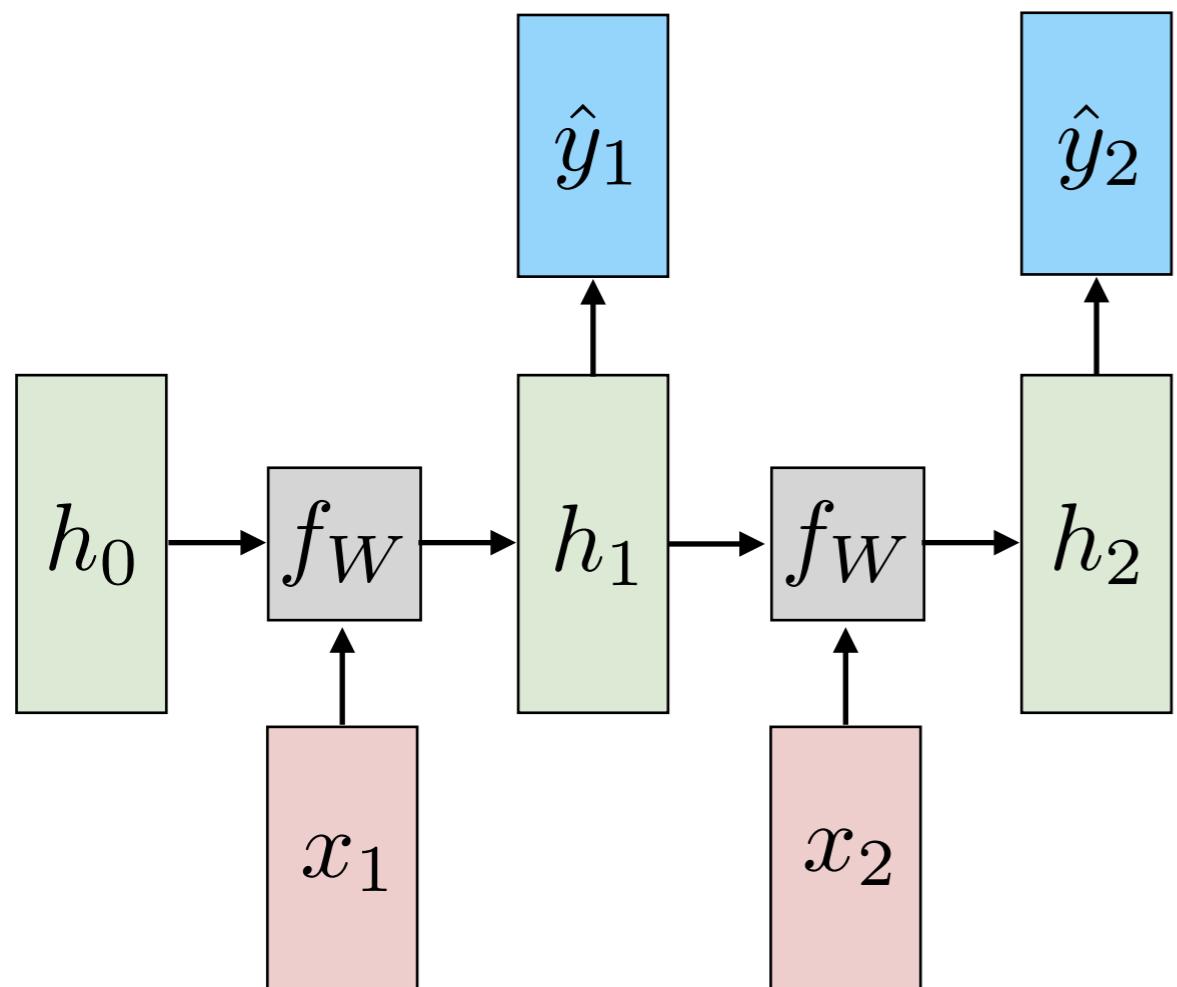
Muestreo secuencial ...



Tomas Mikolov, Martin Karafiat, Lukas Burget, Jan Cernocky, and Sanjeev Khudanpur. 2010.
Recurrent Neural Network Based Language Model. InProceedings of INTERSPEECH

Generación de texto con RNNs (Test)

Muestreo secuencial ...

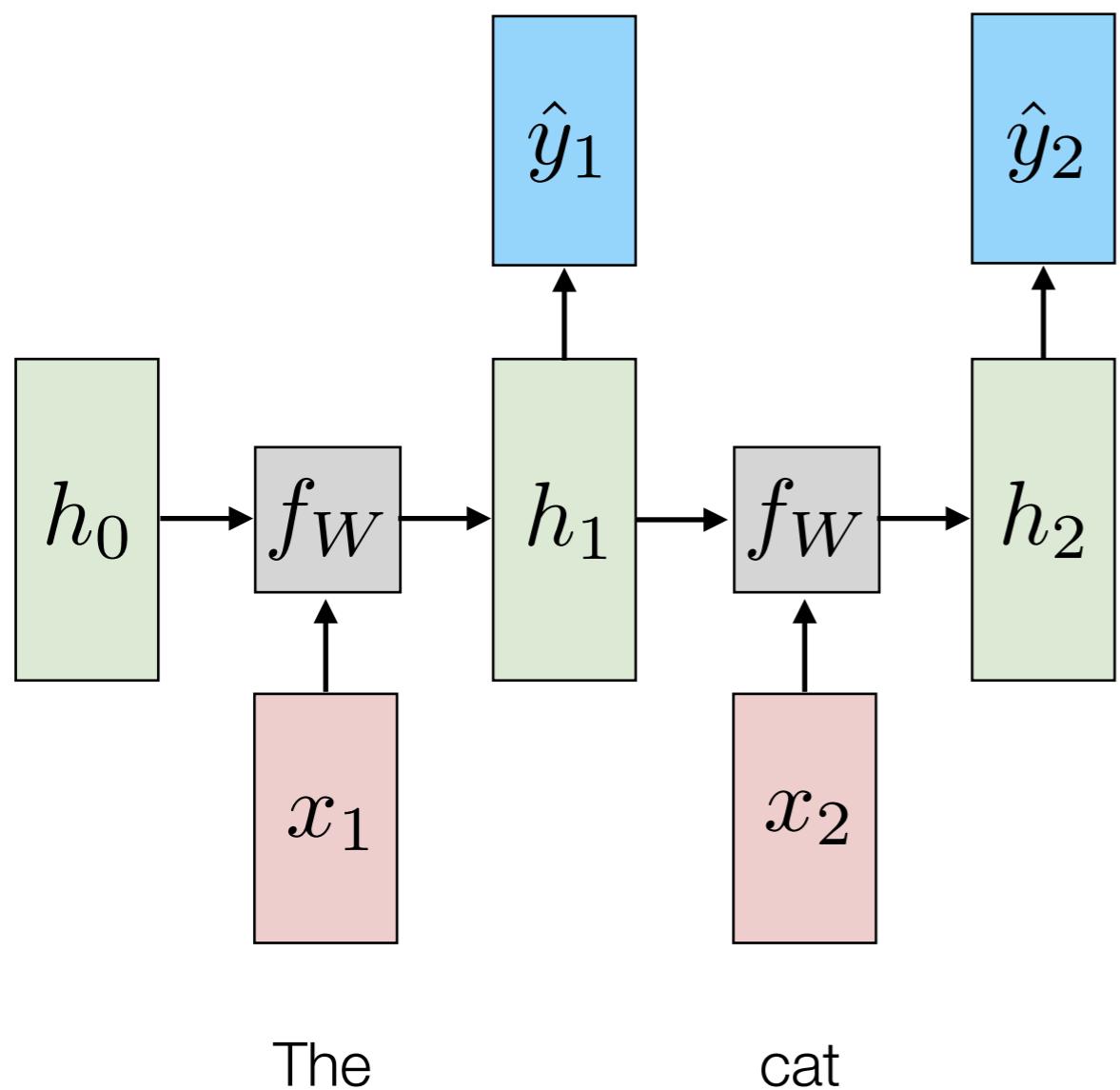


The

Tomas Mikolov, Martin Karafiat, Lukas Burget, Jan Cernocky, and Sanjeev Khudanpur. 2010.
Recurrent Neural Network Based Language Model. InProceedings of INTERSPEECH

Generación de texto con RNNs (Test)

Muestreo secuencial ...

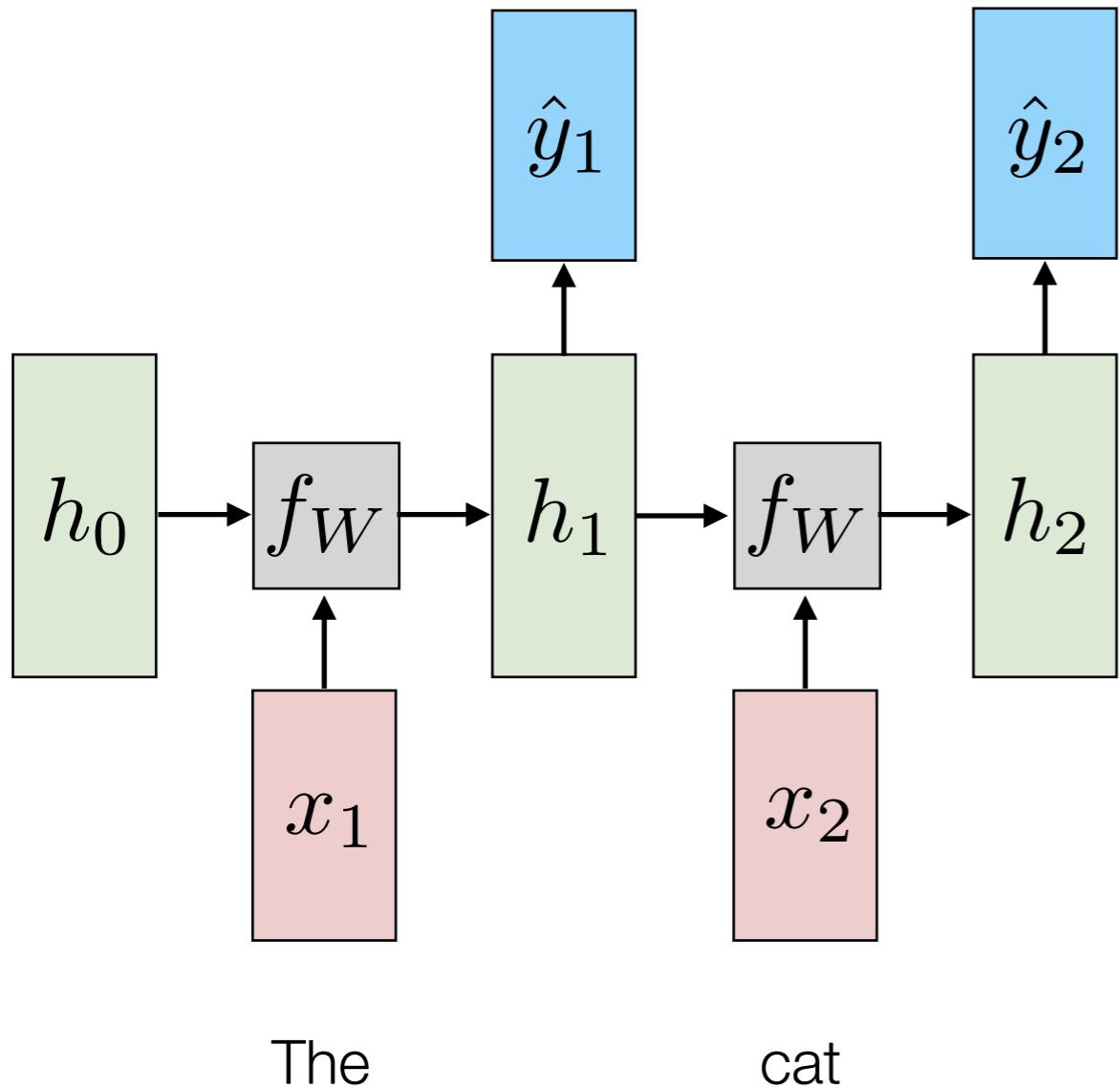


Tomas Mikolov, Martin Karafiat, Lukas Burget, Jan Cernocky, and Sanjeev Khudanpur. 2010.
Recurrent Neural Network Based Language Model. InProceedings of INTERSPEECH

Generación de texto con RNNs (Test)

Muestreo secuencial ...

$$\tilde{x}_3 \sim p(x_3) = \text{softmax}(\hat{y}_2)$$

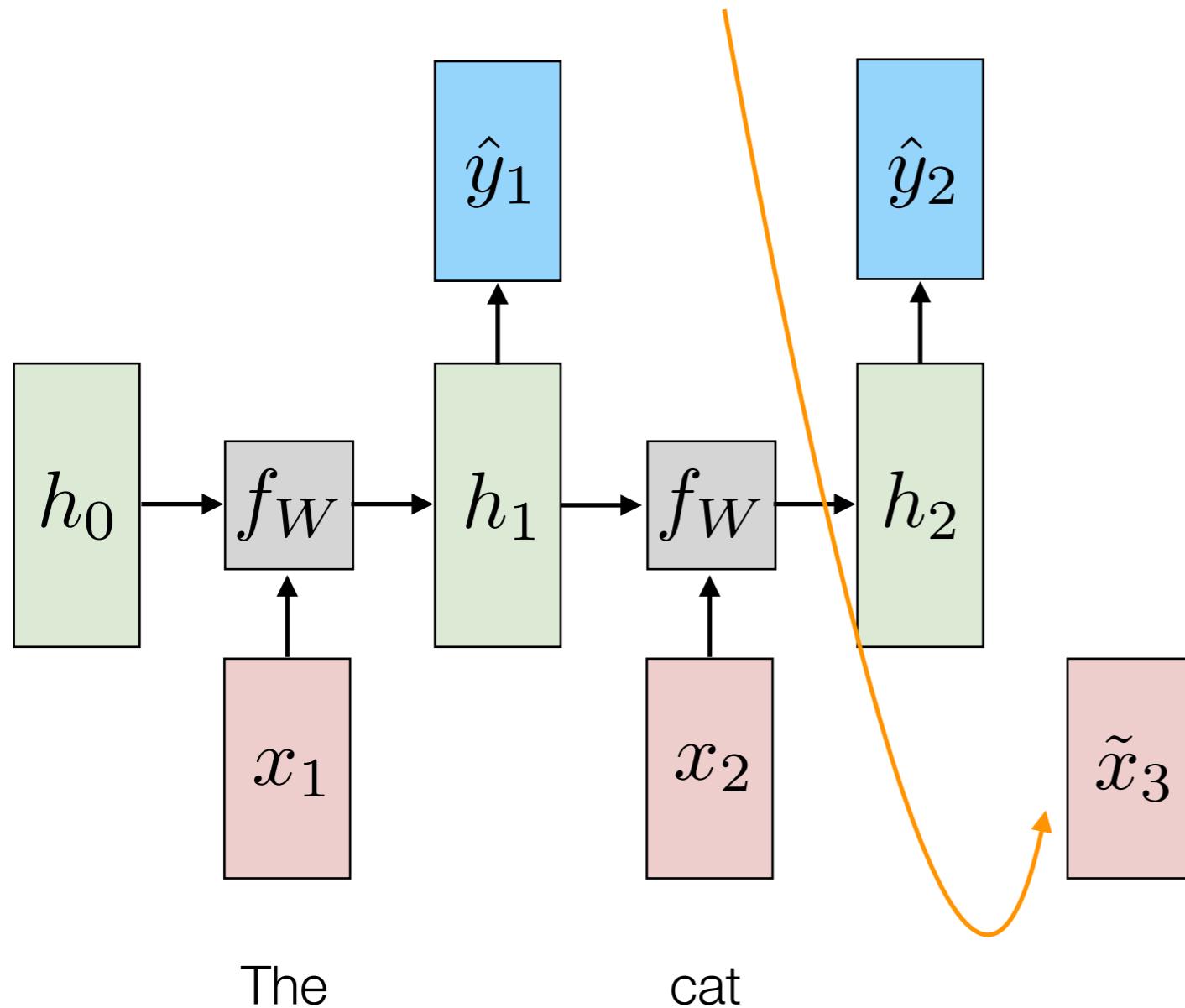


Tomas Mikolov, Martin Karafiat, Lukas Burget, Jan Cernocky, and Sanjeev Khudanpur. 2010.
Recurrent Neural Network Based Language Model. InProceedings of INTERSPEECH

Generación de texto con RNNs (Test)

Muestreo secuencial ...

$$\tilde{x}_3 \sim p(x_3) = \text{softmax}(\hat{y}_2)$$

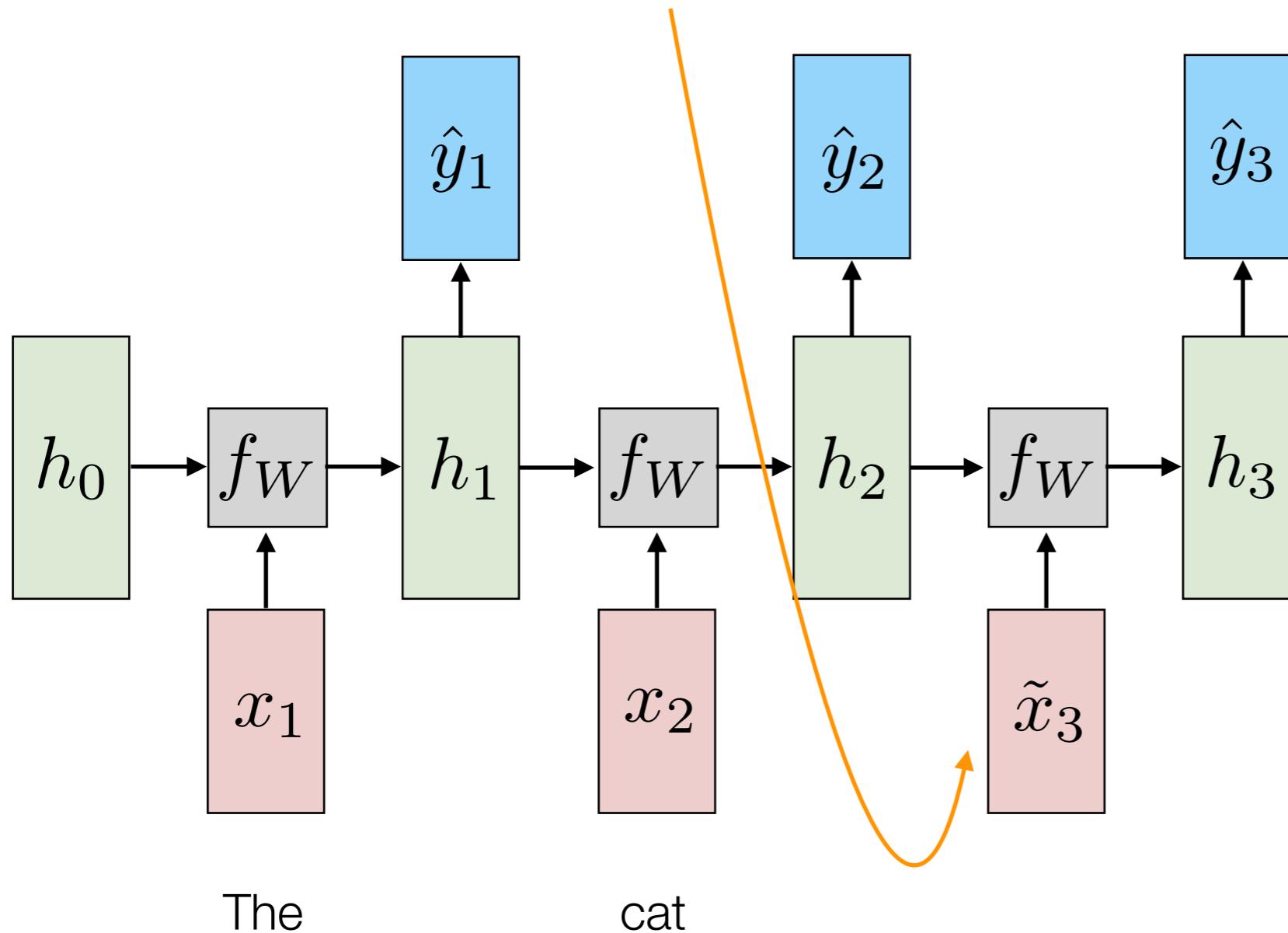


Tomas Mikolov, Martin Karafiat, Lukas Burget, Jan Cernocky, and Sanjeev Khudanpur. 2010.
Recurrent Neural Network Based Language Model. InProceedings of INTERSPEECH

Generación de texto con RNNs (Test)

Muestreo secuencial ...

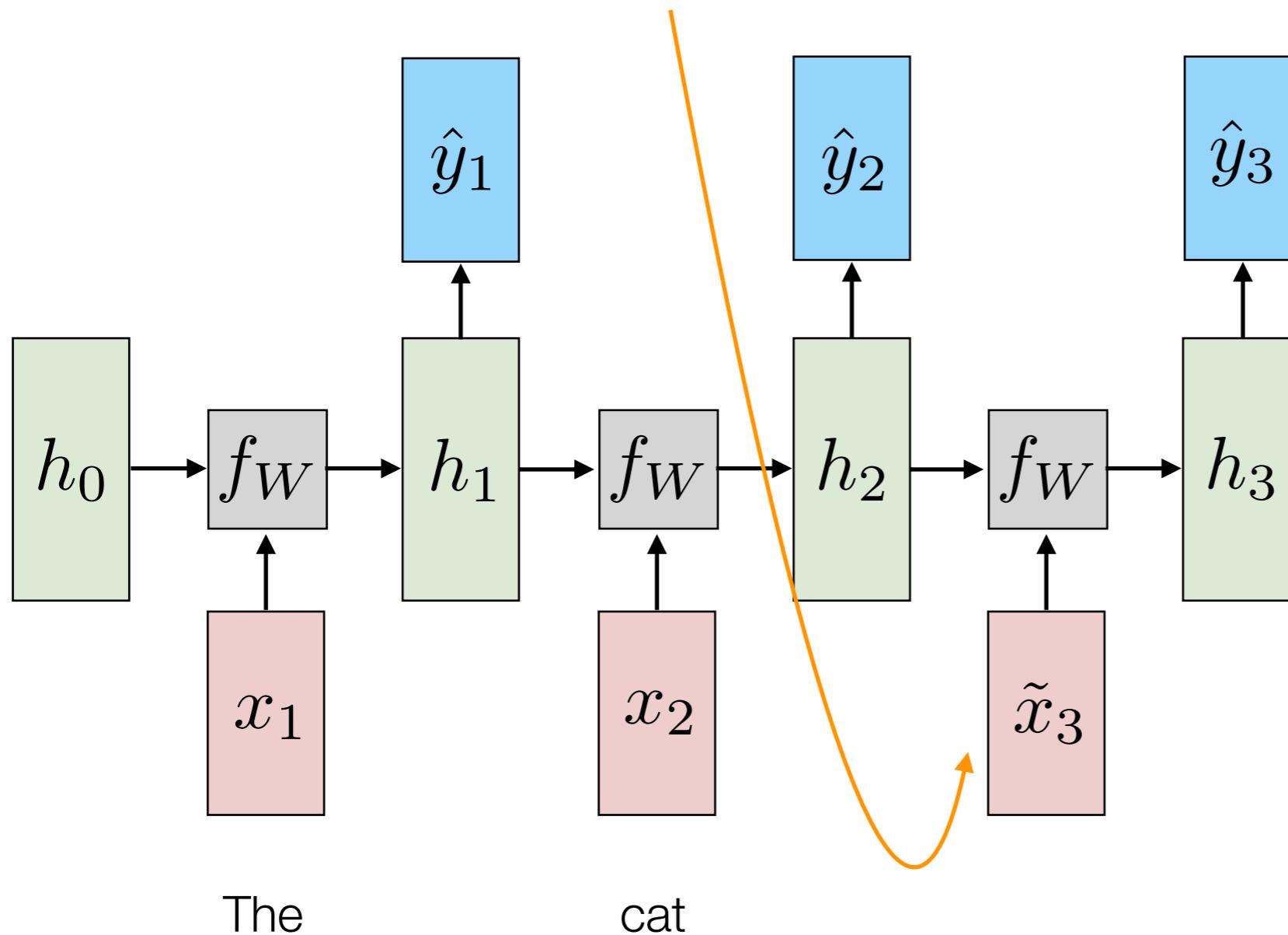
$$\tilde{x}_3 \sim p(x_3) = \text{softmax}(\hat{y}_2)$$



Tomas Mikolov, Martin Karafiat, Lukas Burget, Jan Cernocky, and Sanjeev Khudanpur. 2010.
Recurrent Neural Network Based Language Model. InProceedings of INTERSPEECH

Generación de texto con RNNs (Test)

Muestreo secuencial ...

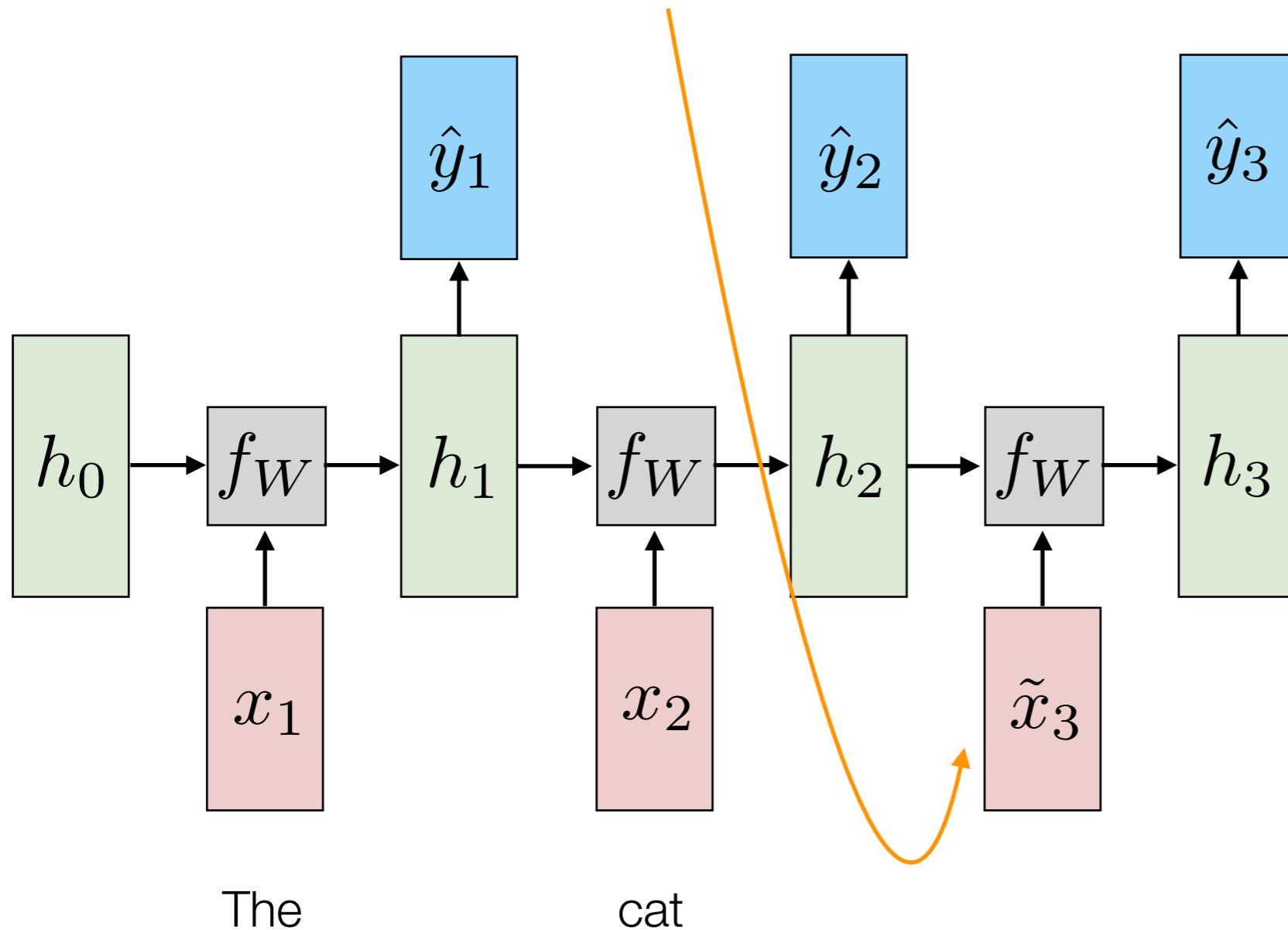


Tomas Mikolov, Martin Karafiat, Lukas Burget, Jan Cernocky, and Sanjeev Khudanpur. 2010.
Recurrent Neural Network Based Language Model. InProceedings of INTERSPEECH

Generación de texto con RNNs (Test)

Muestreo secuencial ...

$$\tilde{x}_4 \sim p(x_4) = \text{softmax}(\hat{y}_3)$$

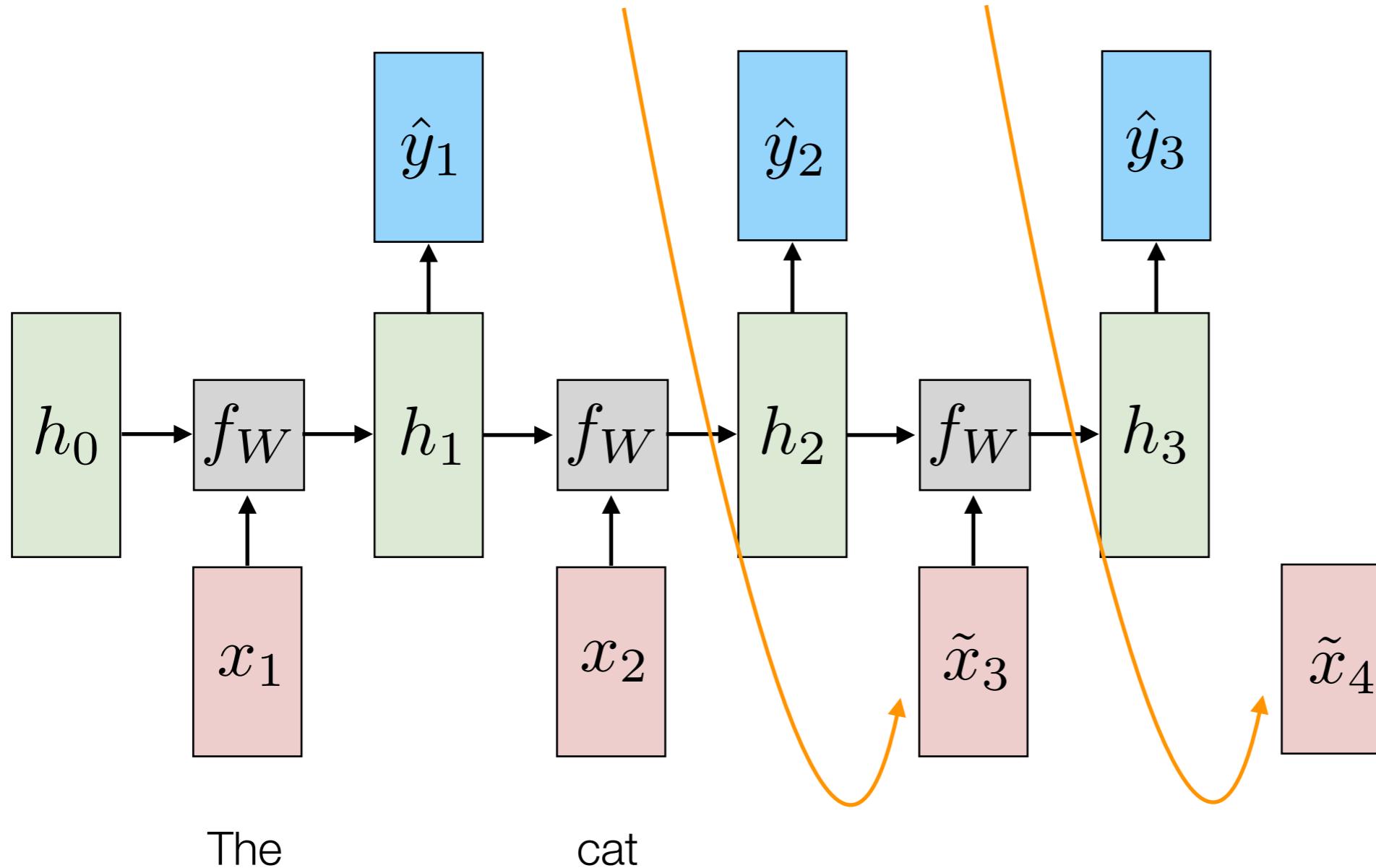


Tomas Mikolov, Martin Karafiat, Lukas Burget, Jan Cernocky, and Sanjeev Khudanpur. 2010.
Recurrent Neural Network Based Language Model. InProceedings of INTERSPEECH

Generación de texto con RNNs (Test)

Muestreo secuencial ...

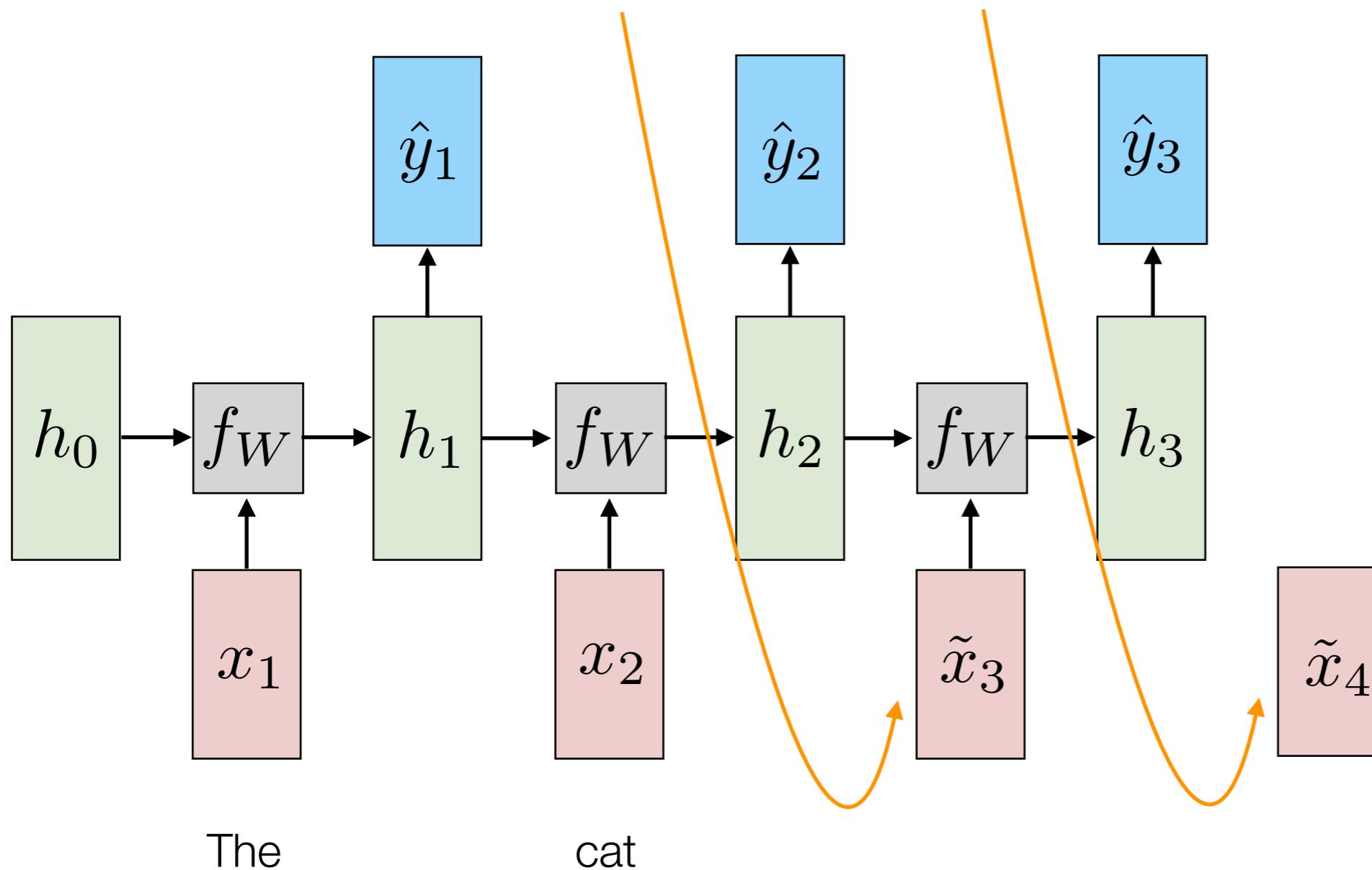
$$\tilde{x}_4 \sim p(x_4) = \text{softmax}(\hat{y}_3)$$



Tomas Mikolov, Martin Karafiat, Lukas Burget, Jan Cernocky, and Sanjeev Khudanpur. 2010.
Recurrent Neural Network Based Language Model. InProceedings of INTERSPEECH

Generación de texto con RNNs (Test)

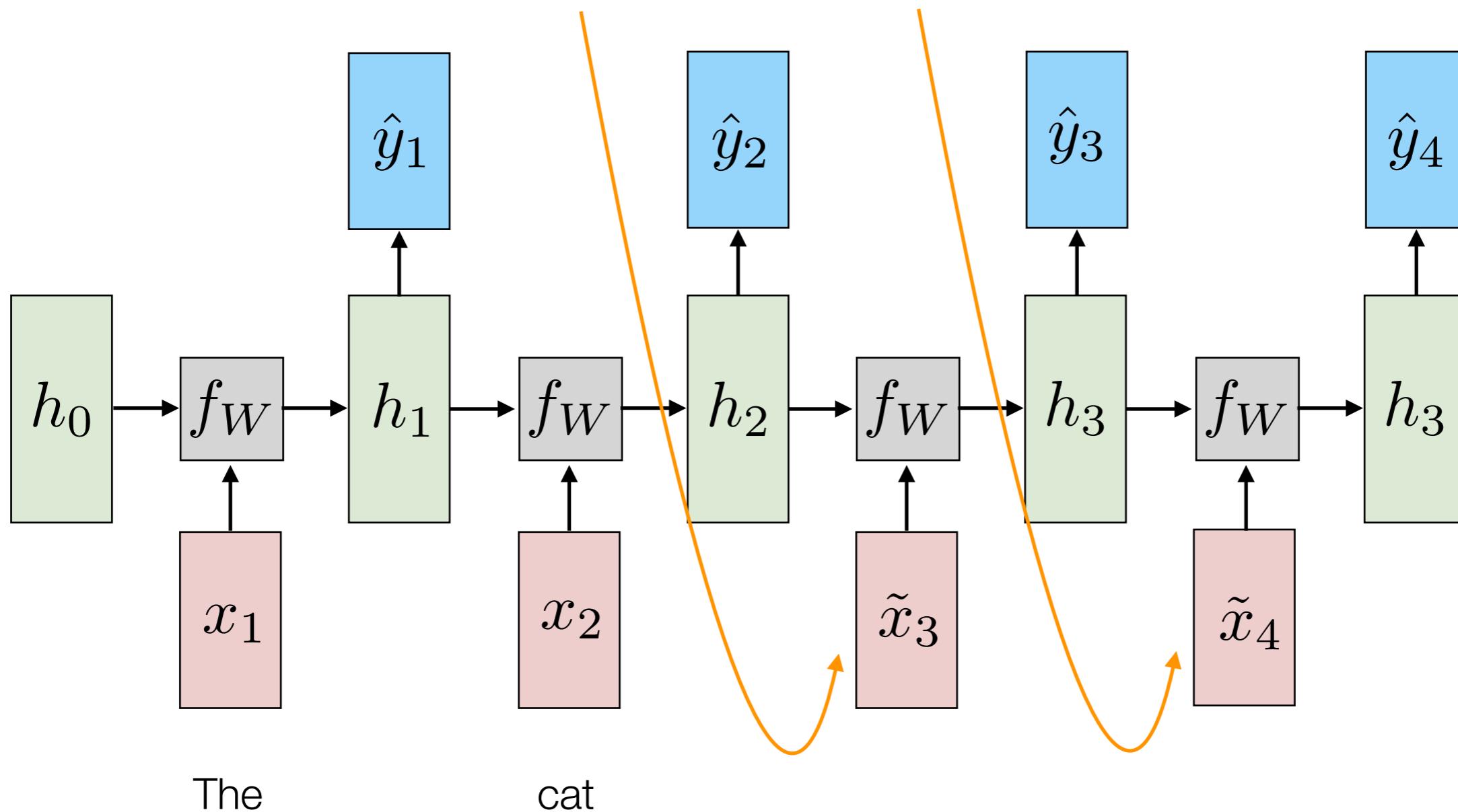
Muestreo secuencial ...



Tomas Mikolov, Martin Karafiat, Lukas Burget, Jan Cernocky, and Sanjeev Khudanpur. 2010.
Recurrent Neural Network Based Language Model. InProceedings of INTERSPEECH

Generación de texto con RNNs (Test)

Muestreo secuencial ...

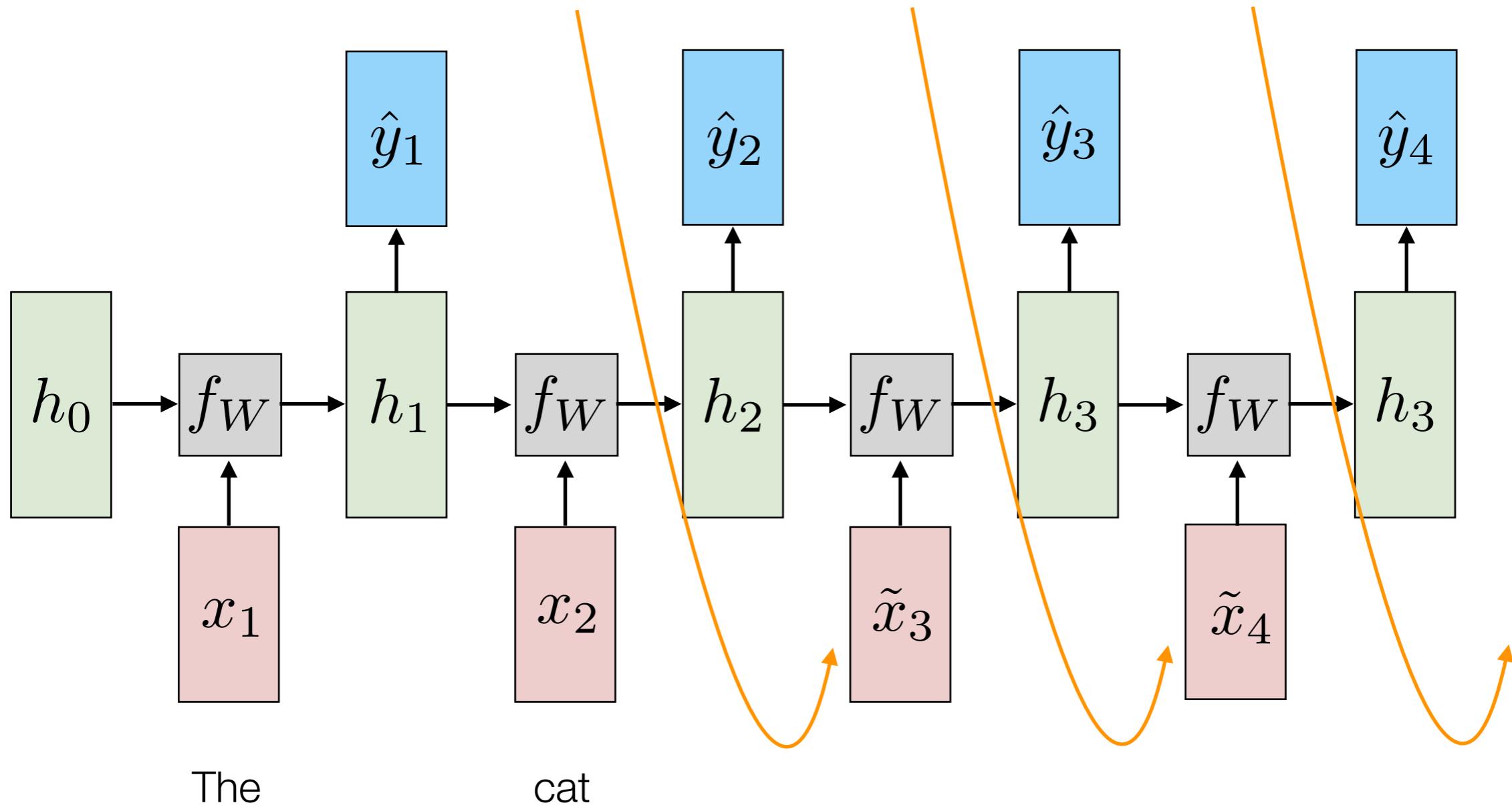


Tomas Mikolov, Martin Karafiat, Lukas Burget, Jan Cernocky, and Sanjeev Khudanpur. 2010.
Recurrent Neural Network Based Language Model. InProceedings of INTERSPEECH

Generación de texto con RNNs (Test)

Muestreo secuencial ...

$$\tilde{x}_5 \sim p(x_5) = \text{softmax}(\hat{y}_4)$$

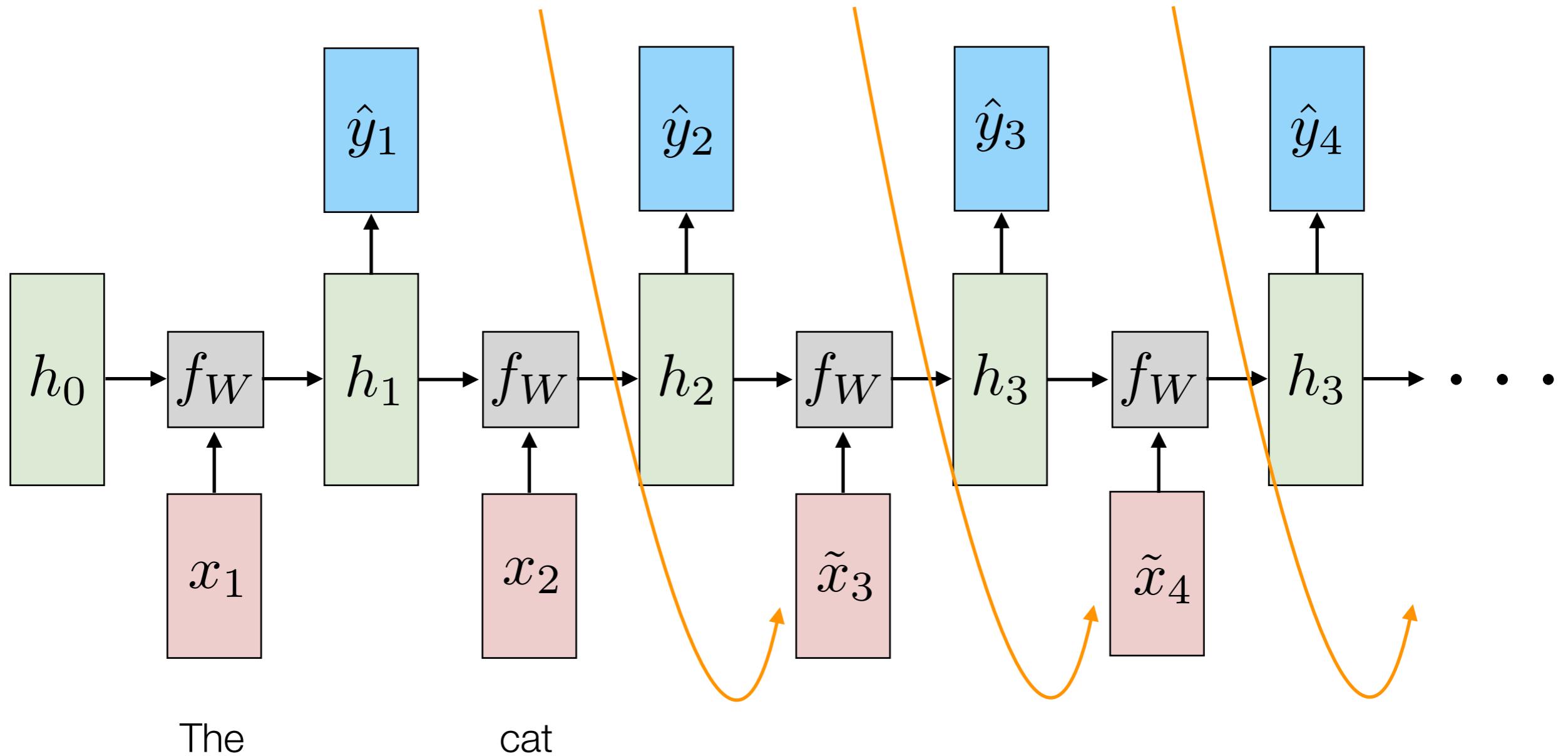


Tomas Mikolov, Martin Karafiat, Lukas Burget, Jan Cernocky, and Sanjeev Khudanpur. 2010.
Recurrent Neural Network Based Language Model. InProceedings of INTERSPEECH

Generación de texto con RNNs (Test)

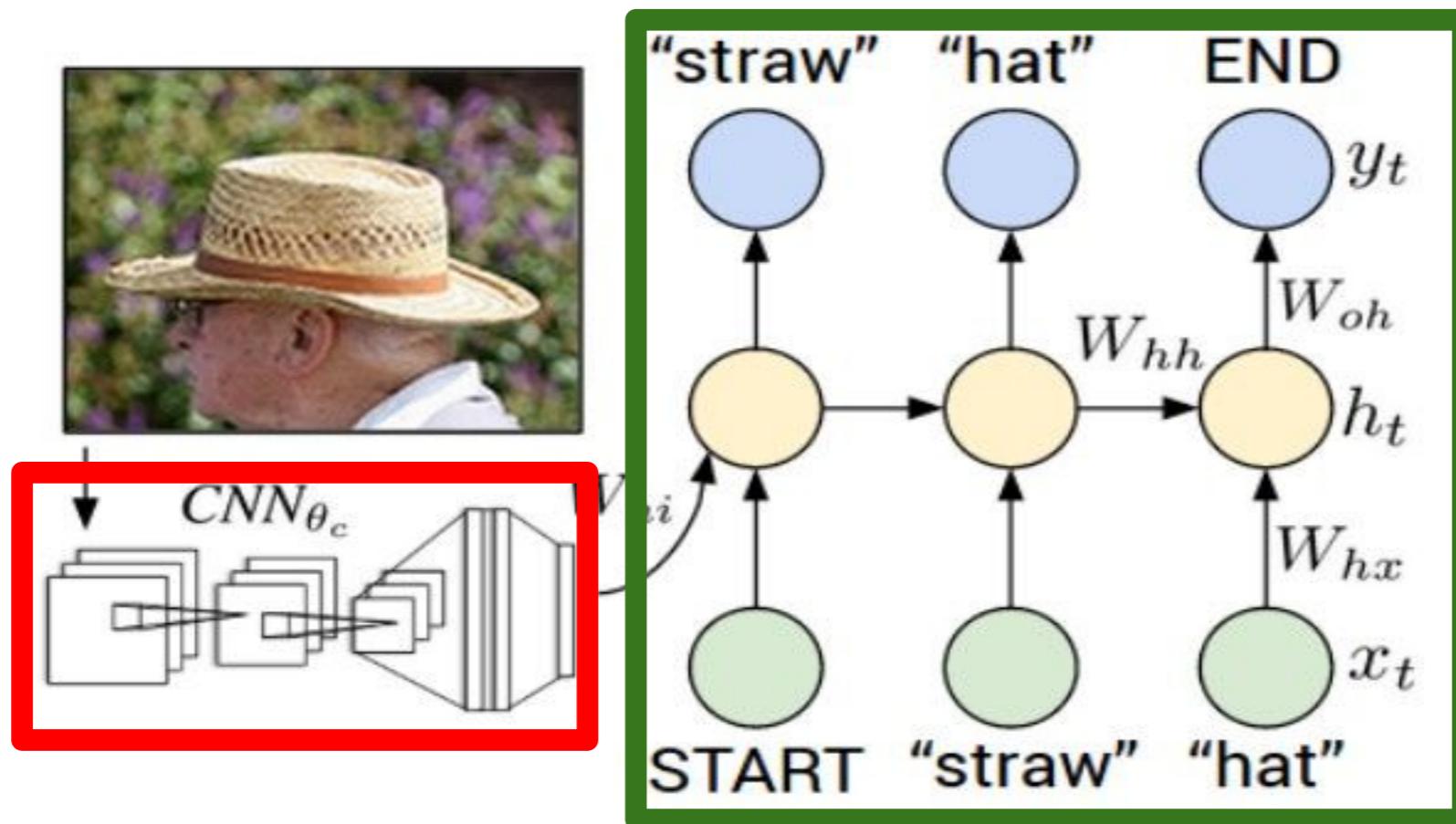
Muestreo secuencial ...

$$\tilde{x}_5 \sim p(x_5) = \text{softmax}(\hat{y}_4)$$



Tomas Mikolov, Martin Karafiat, Lukas Burget, Jan Cernocky, and Sanjeev Khudanpur. 2010.
Recurrent Neural Network Based Language Model. InProceedings of INTERSPEECH

Recurrent Neural Network



Convolutional Neural Network

Explain Images with Multimodal Recurrent Neural Networks, Mao et al.

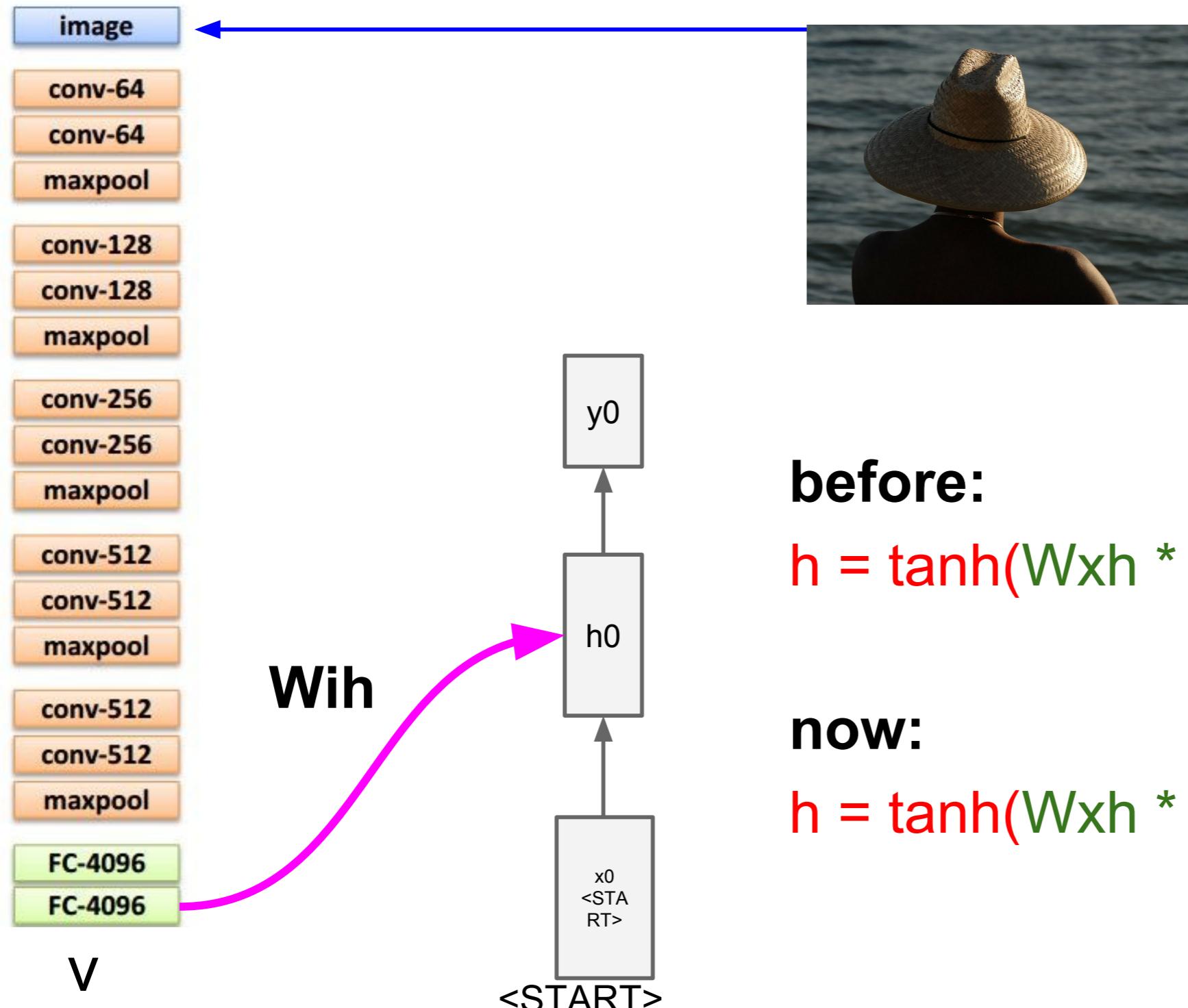
Deep Visual-Semantic Alignments for Generating Image Descriptions, Karpathy and Fei-Fei

Show and Tell: A Neural Image Caption Generator, Vinyals et al.

Long-term Recurrent Convolutional Networks for Visual Recognition and Description, Donahue et al.

Learning a Recurrent Visual Representation for Image Caption Generation, Chen and Zitnick

Descripción automática de imágenes



before:

$$h = \tanh(W_{xh} * x + W_{hh} * h)$$

now:

$$h = \tanh(W_{xh} * x + W_{hh} * h + W_{ih} * v)$$

Descripción automática de imágenes

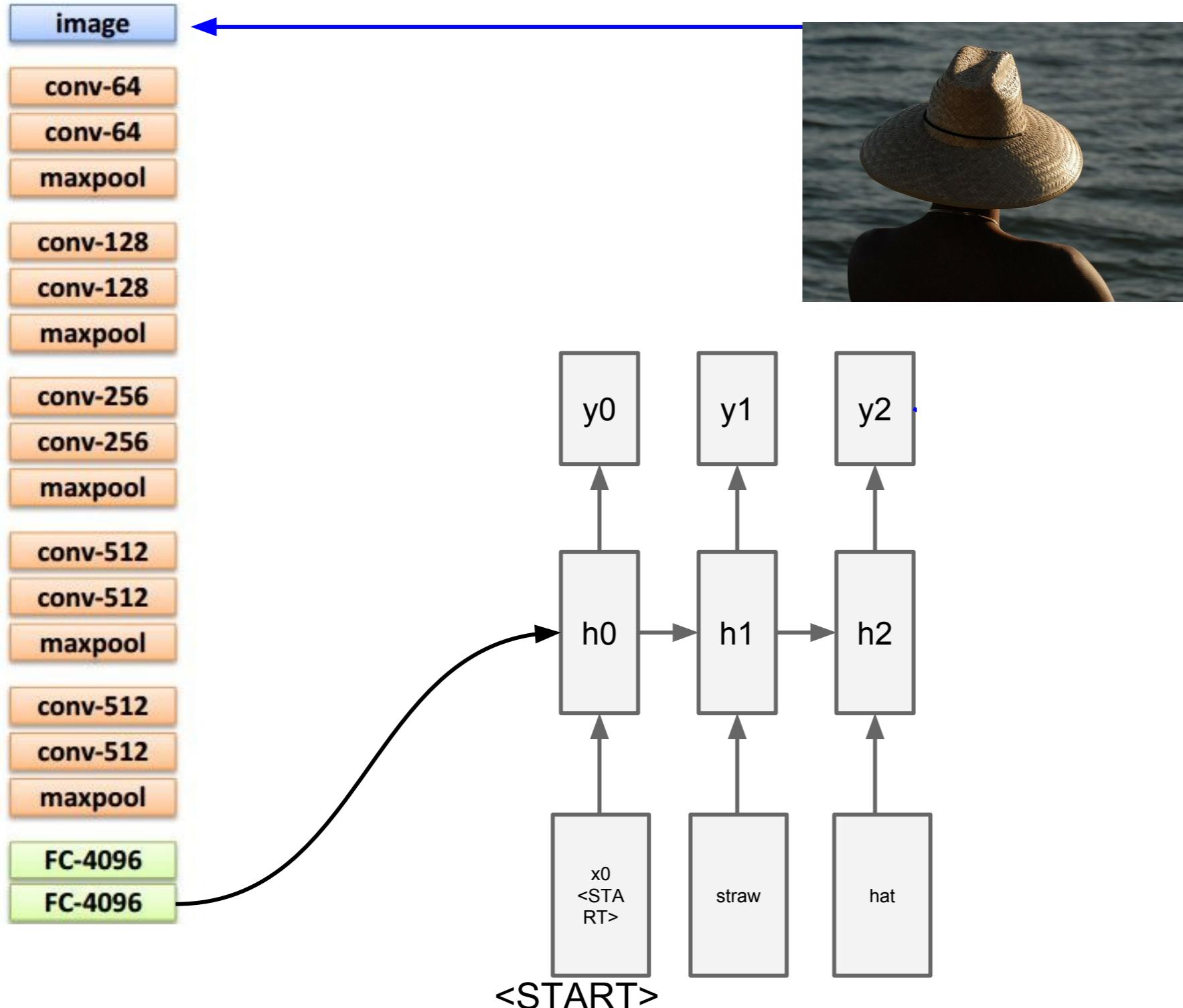


Image Captioning: Example Results

Captions generated using [neuraltalk2](#)
All images are CC0 Public domain:
[cat suitcase](#), [cat tree](#), [dog](#), [bear](#),
[surfers](#), [tennis](#), [giraffe](#), [motorcycle](#)



A cat sitting on a suitcase on the floor



A cat is sitting on a tree branch



A dog is running in the grass with a frisbee



A white teddy bear sitting in the grass



Two people walking on the beach with surfboards



A tennis player in action on the court



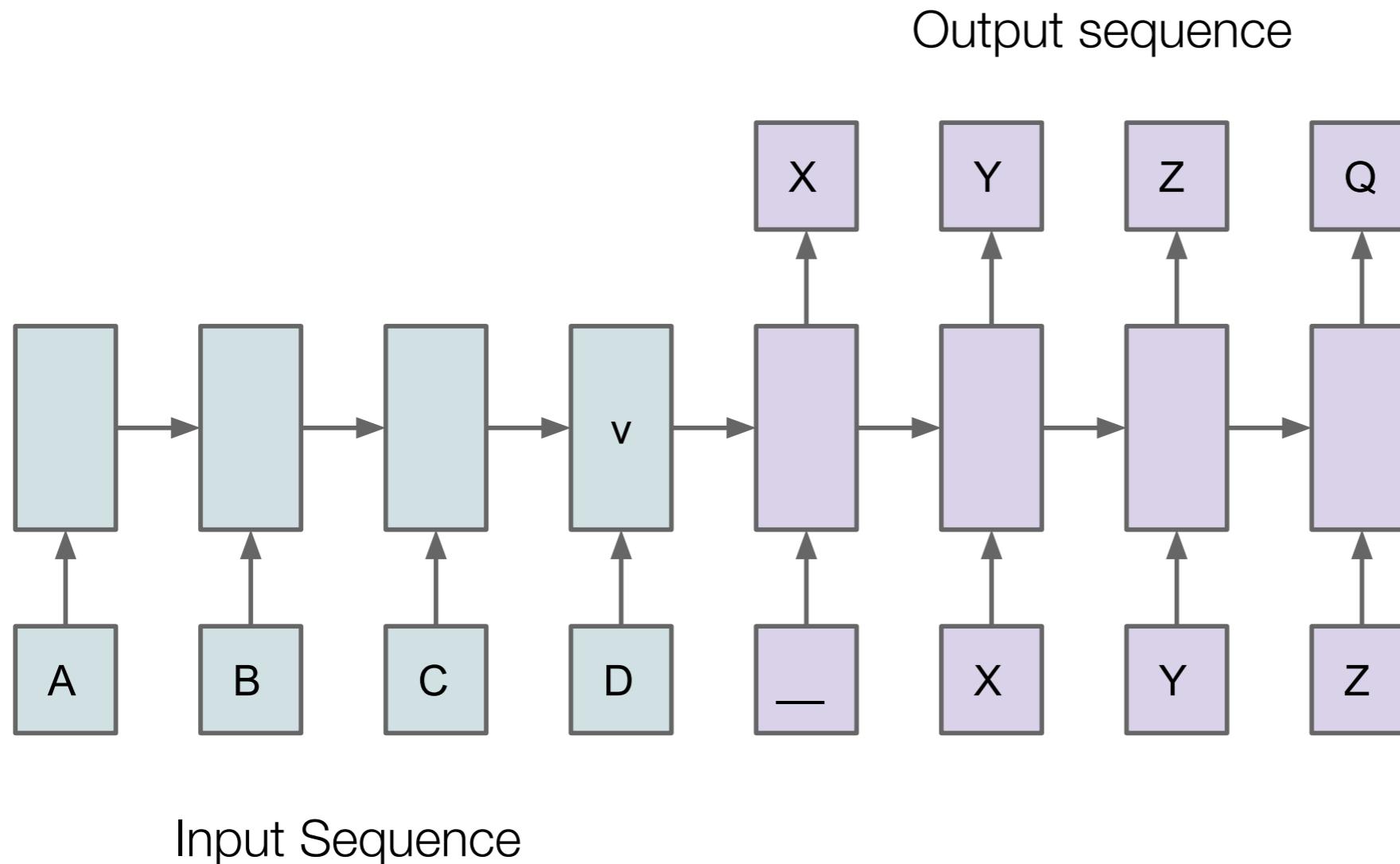
Two giraffes standing in a grassy field



A man riding a dirt bike on a dirt track

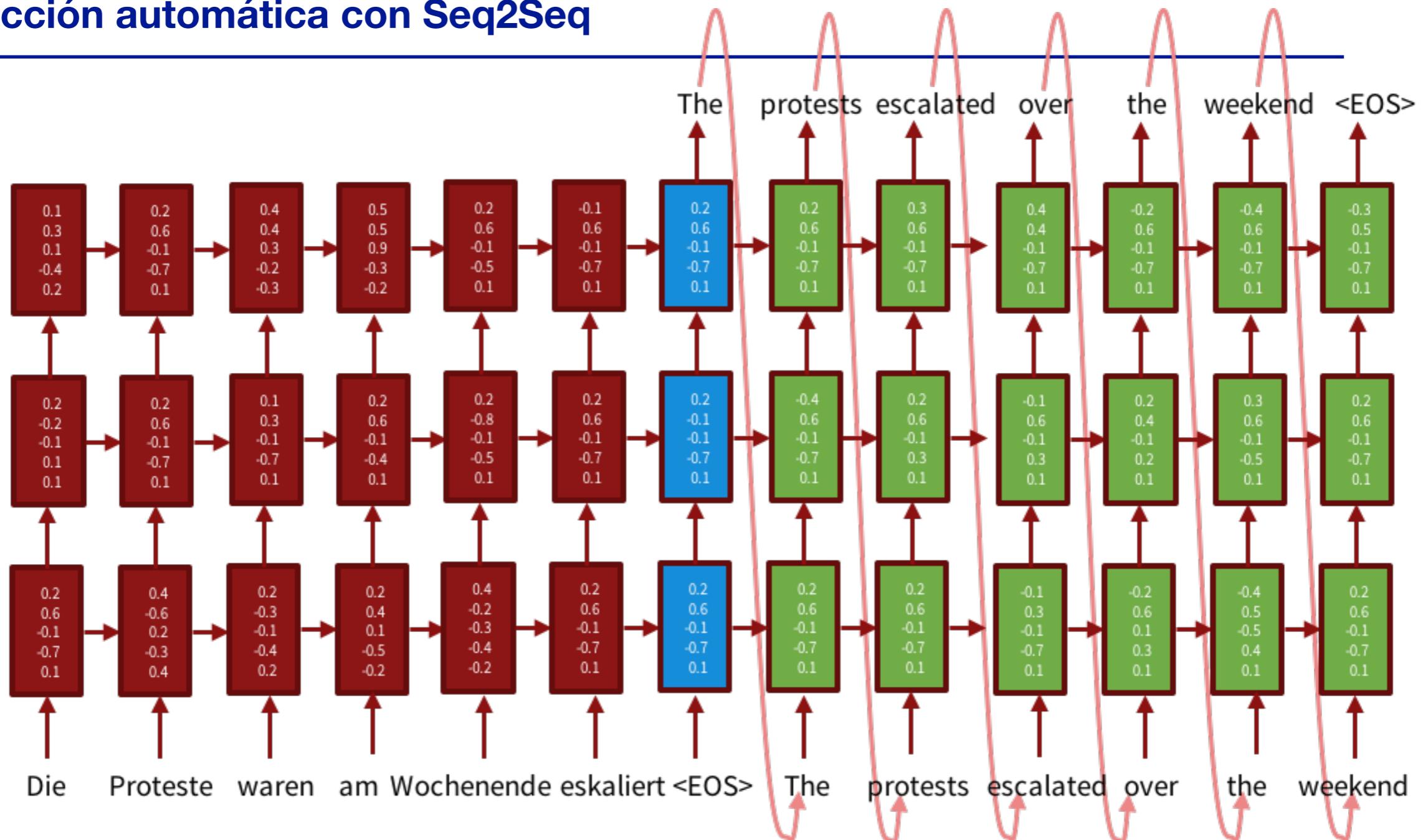
Seq2Seq y modelos de atención

Modelo Seq2Seq



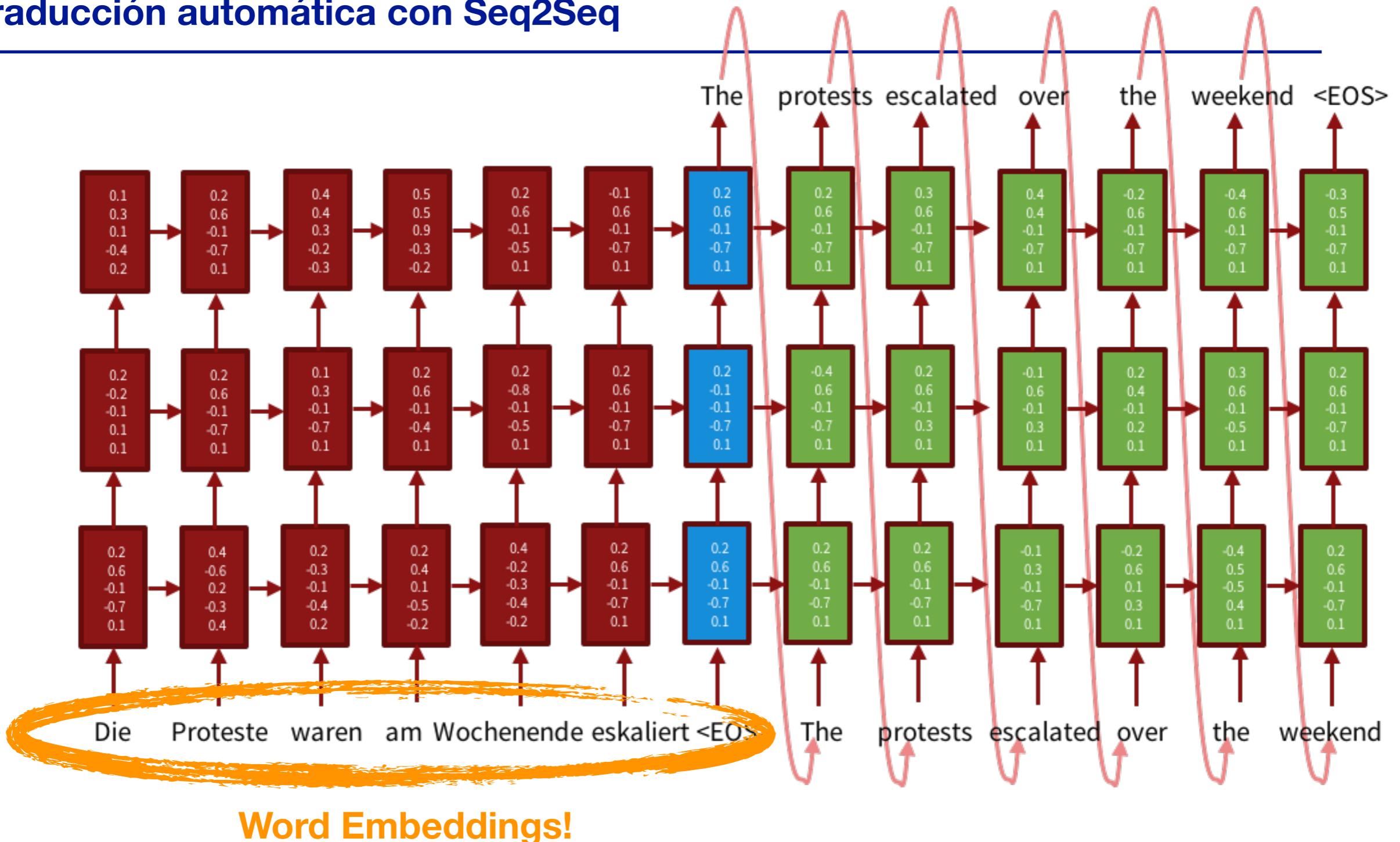
$$P(y_1, \dots, y_{T'} | x_1, \dots, x_T) = \prod_{t=1}^{T'} p(y_t | v, y_1, \dots, y_{t-1})$$

Traducción automática con Seq2Seq



1. Auli, M., et al. "Joint Language and Translation Modeling with Recurrent Neural Networks." *EMNLP* (2013)
2. Kalchbrenner, N., et al. "Recurrent Continuous Translation Models." *EMNLP* (2013)
3. Cho, K., et al. "Learning Phrase Representations using RNN Encoder-Decoder for Statistical MT." *EMNLP* (2014)
4. Sutskever, I., et al. "Sequence to Sequence Learning with Neural Networks." *NIPS* (2014)

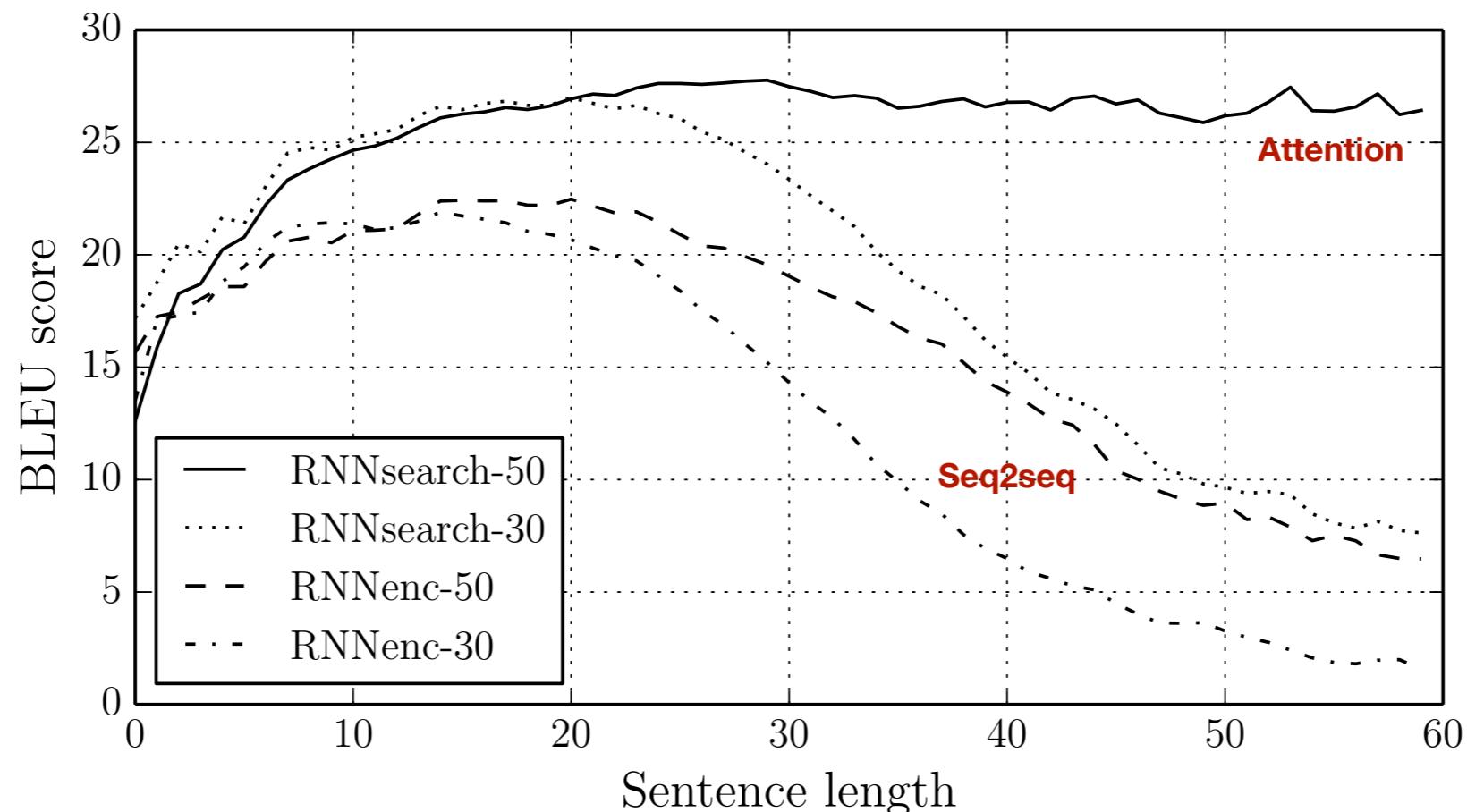
Traducción automática con Seq2Seq



1. Auli, M., et al. "Joint Language and Translation Modeling with Recurrent Neural Networks." *EMNLP* (2013)
2. Kalchbrenner, N., et al. "Recurrent Continuous Translation Models." *EMNLP* (2013)
3. Cho, K., et al. "Learning Phrase Representations using RNN Encoder-Decoder for Statistical MT." *EMNLP* (2014)
4. Sutskever, I., et al. "Sequence to Sequence Learning with Neural Networks." *NIPS* (2014)

Limitaciones Seq2Seq para traducción automática

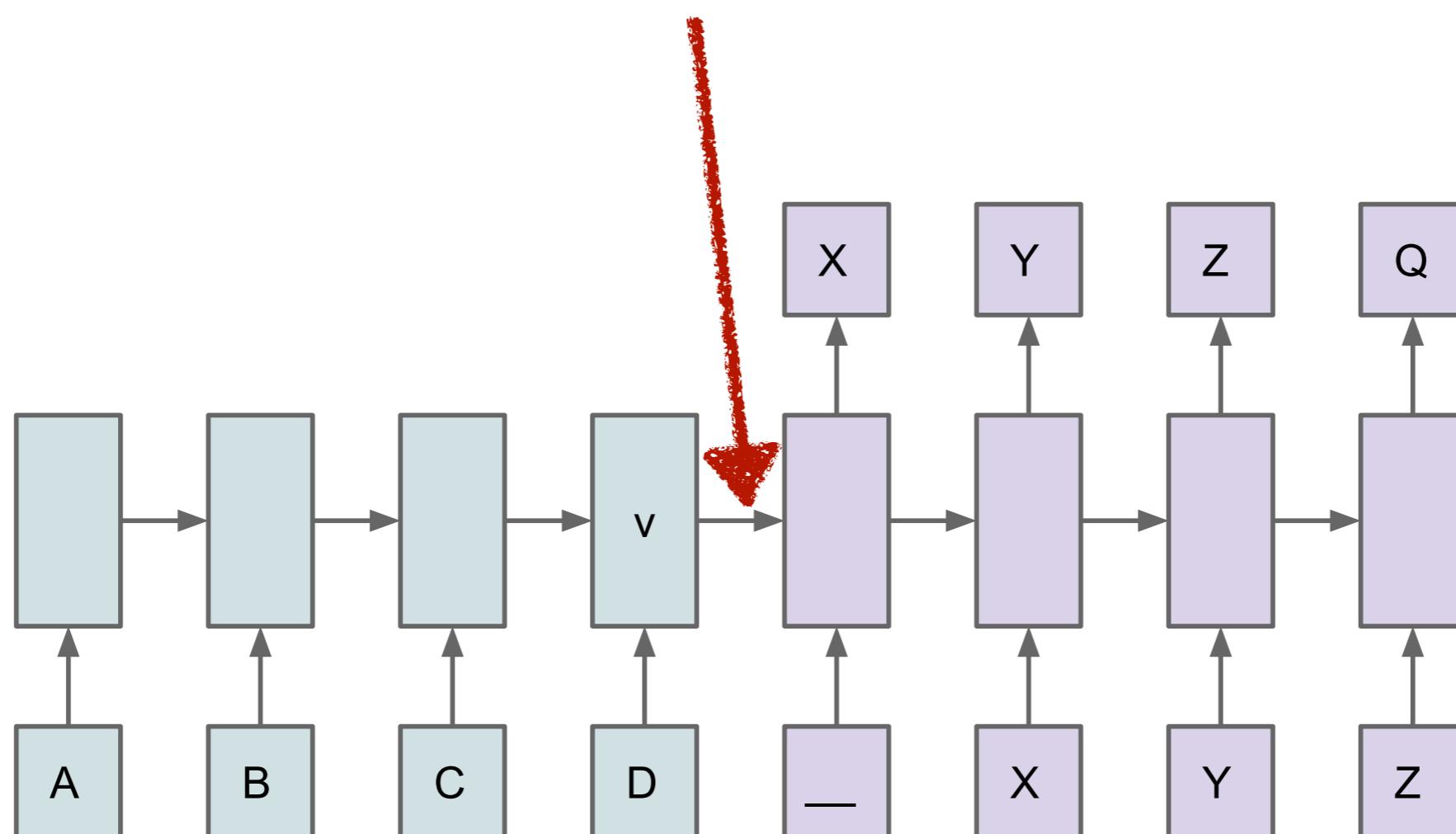
- Rendimiento cae rápidamente con la longitud de la secuencia
- BLEU (bilingual evaluation understudy) es uno métrica para evaluar la calidad de una traducción



1. Sutskever, I., et al. "Sequence to Sequence Learning with Neural Networks." *NIPS* (2014)
2. Bahdanau, D., et al. "Neural Machine Translation by Jointly Learning to Align and Translate." *ICLR* (2015)

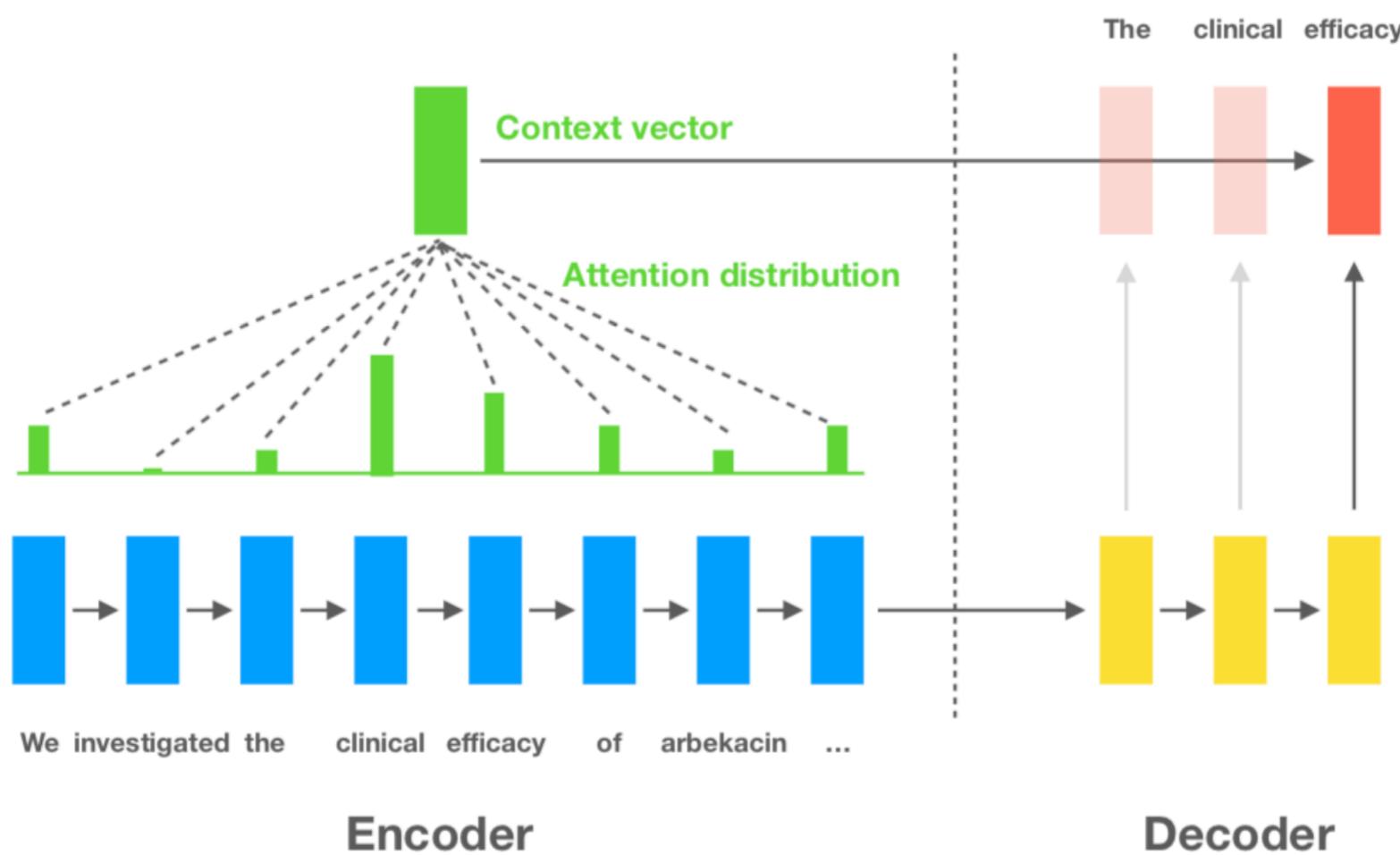
Limitaciones Seq2Seq para traducción automática

Un único embedding para representar toda la secuencia!



Seq2Seq con mecanismos de atención

- Permitimos al decodificador prestar **atención a diferentes partes de la secuencia de entrada** en cada paso de la generación de la secuencia de salida.
- Cada palabra generada por el decodificador estará condicionada a una **suma ponderada de los estados en el codificador**.

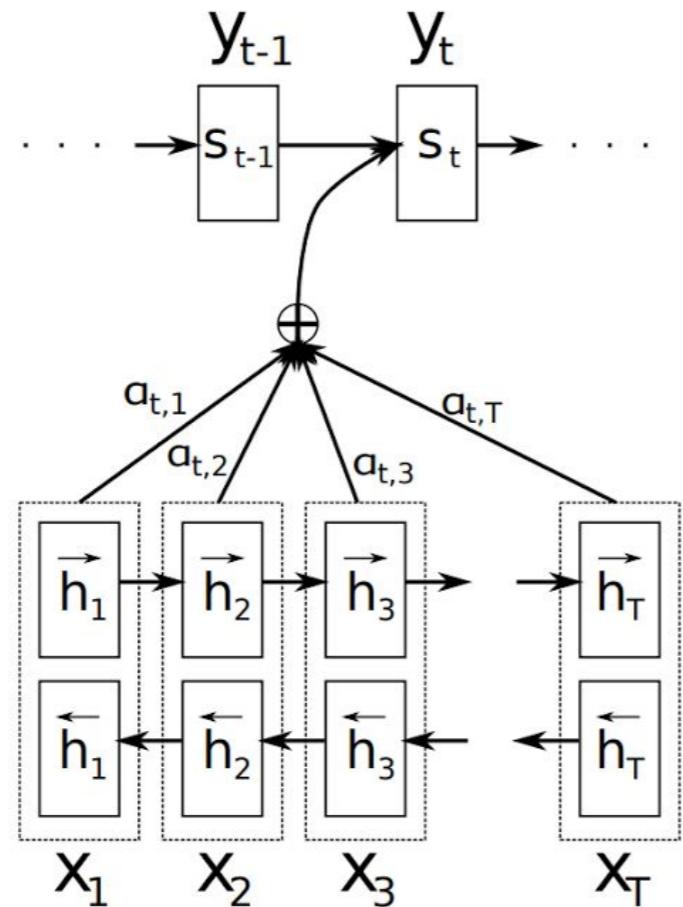


Seq2Seq con mecanismos de atención

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

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

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



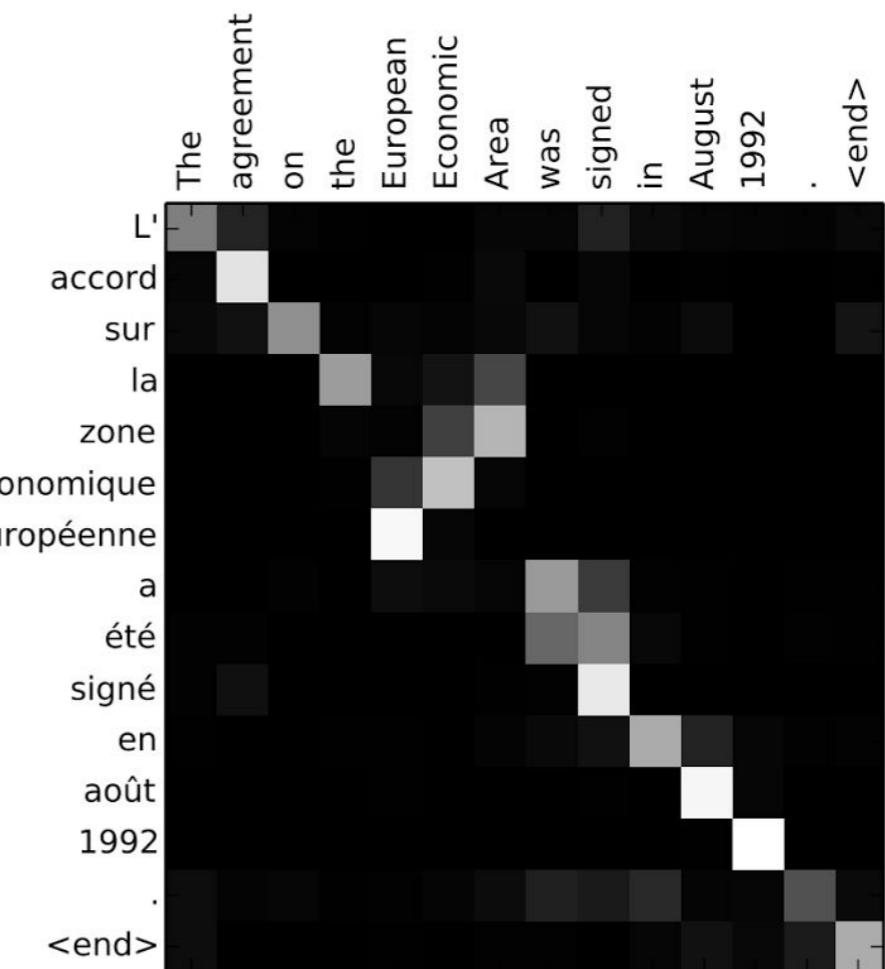
NEURAL MACHINE TRANSLATION BY JOINTLY LEARNING TO ALIGN AND TRANSLATE

Dzmitry Bahdanau

Jacobs University Bremen, Germany

KyungHyun Cho Yoshua Bengio*

Université de Montréal



Trasformer networks

Attention Is All You Need

Ashish Vaswani*
Google Brain
avaswani@google.com

Noam Shazeer*
Google Brain
noam@google.com

Niki Parmar*
Google Research
nikip@google.com

Jakob Uszkoreit*
Google Research
usz@google.com

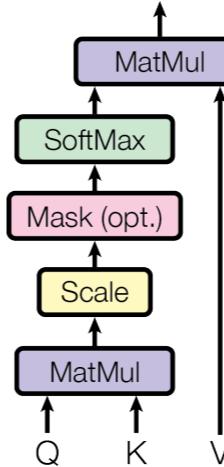
Llion Jones*
Google Research
llion@google.com

Aidan N. Gomez* †
University of Toronto
aidan@cs.toronto.edu

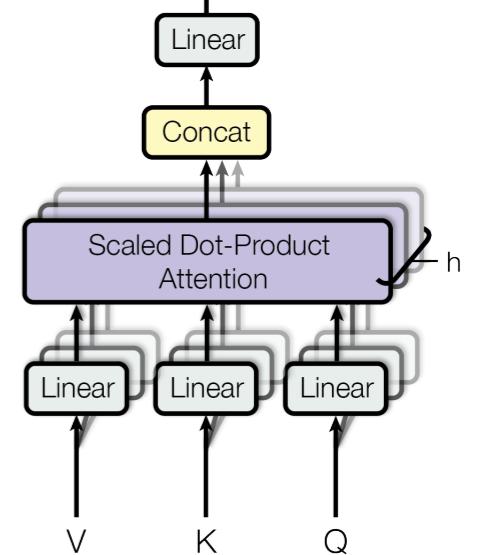
Lukasz Kaiser*
Google Brain
lukaszkaiser@google.com

Illia Polosukhin* ‡
illia.polosukhin@gmail.com

Scaled Dot-Product Attention



Multi-Head Attention

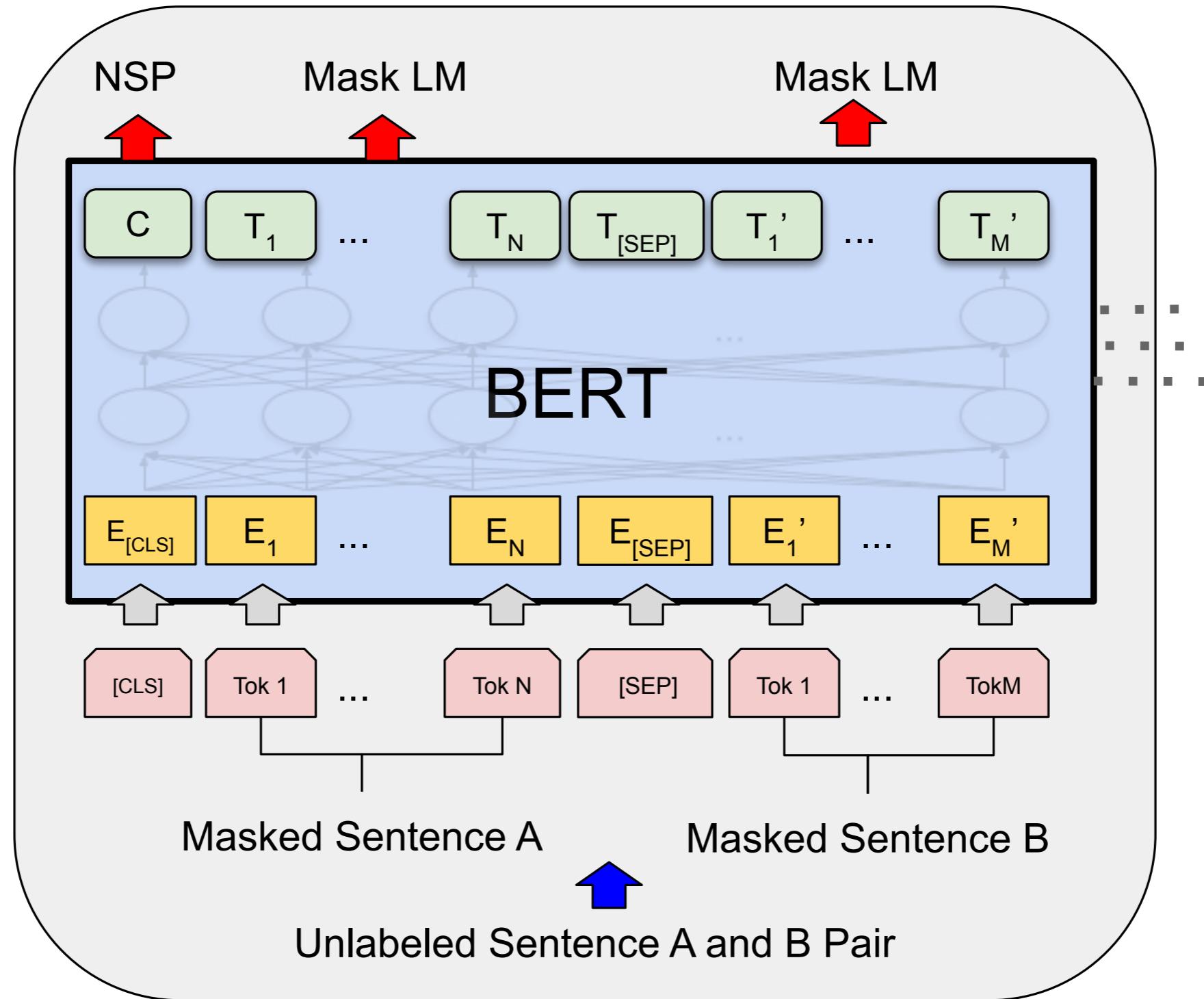


Dec 2017

The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention mechanism. We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely. Experiments on two machine translation tasks show these models to be superior in quality while being more parallelizable and requiring significantly less time to train. Our model achieves 28.4 BLEU on the WMT 2014 English-

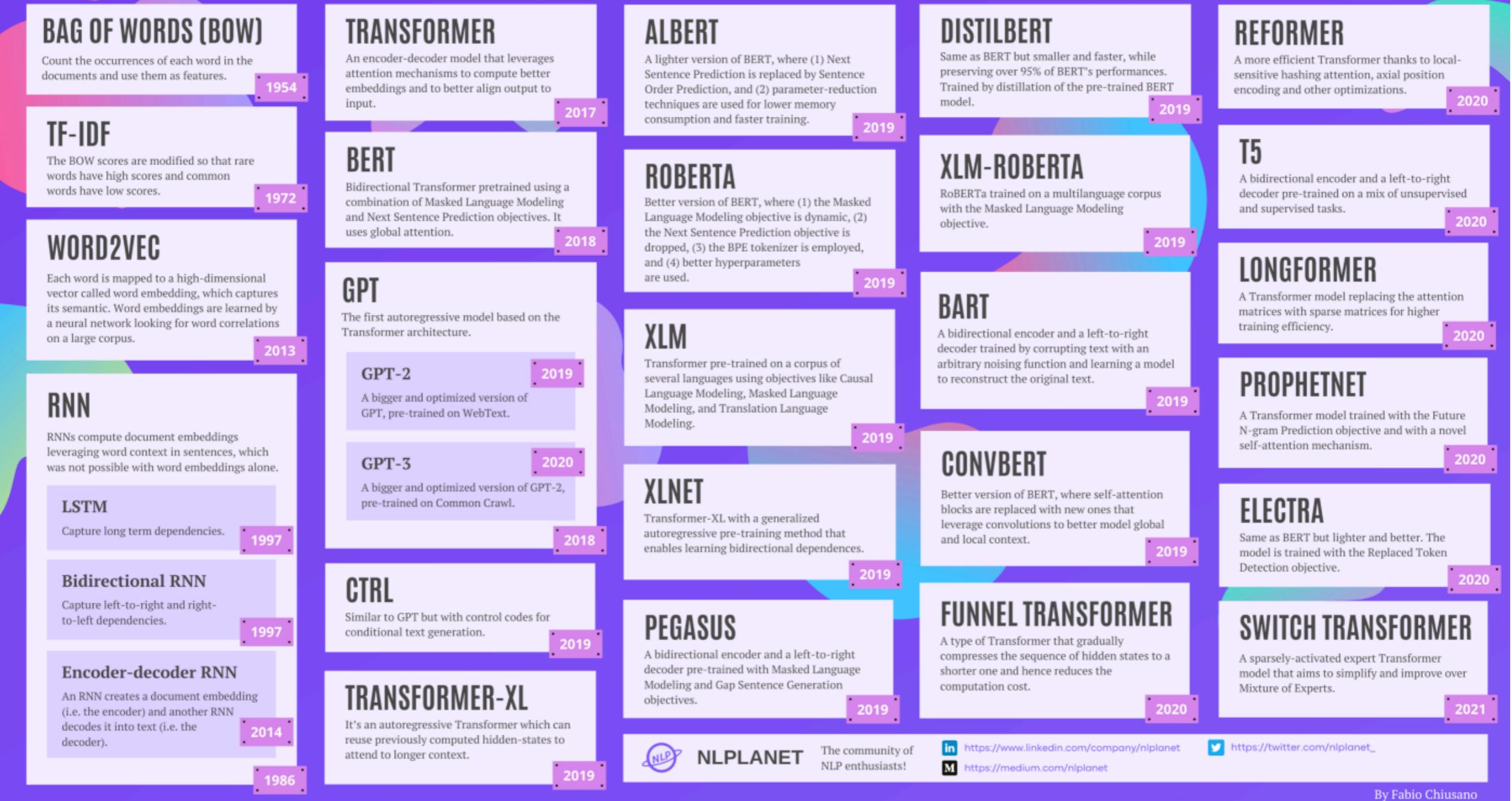
In this work we propose the Transformer, a model architecture eschewing recurrence and instead relying entirely on an attention mechanism to draw global dependencies between input and output.

BERT: Pre-training of Deep bidirectional Transformers for Language Understanding



Pre-training

AN NLP TIMELINE AND THE TRANSFORMER FAMILY



By Fabio Chiusano

Librería para NLP con transformer networks : Hugging Face library



The AI community building the future.

Build, train and deploy state of the art models powered by
the reference open source in natural language processing.

<https://huggingface.co>