

Recurrent Neural Networks (RNNs)

Naeemullah Khan

naeemullah.khan@kaust.edu.sa



جامعة الملك عبد الله
للعلوم والتقنية
King Abdullah University of
Science and Technology

KAUST Academy
King Abdullah University of Science and Technology

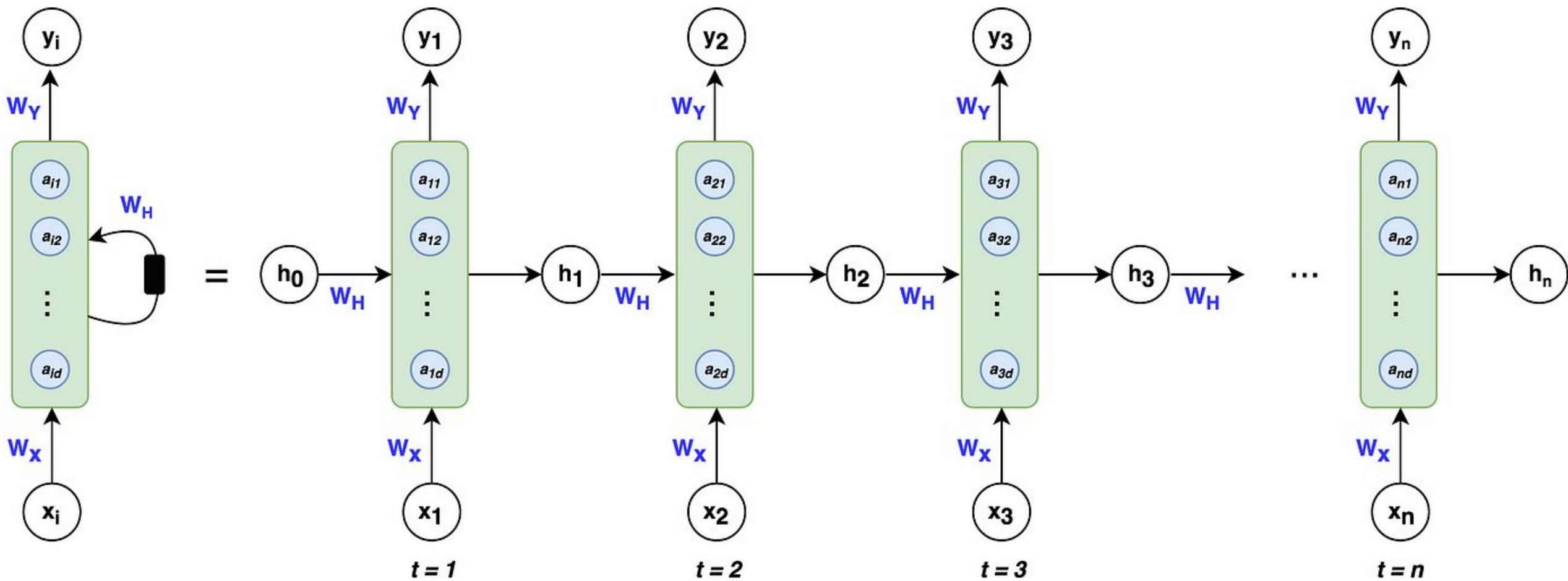


Table of Contents

1. Motivation
2. Learning Outcomes
3. Introduction
4. RNN Architectures
5. RNN Training Challenges
6. Applications of RNNs
7. Limitations of RNNs
8. Future Directions
9. Summary
10. References

Why Do We Need RNNs?

- ▶ Traditional feedforward networks don't handle **sequential data** effectively.
- ▶ Many applications (**language modeling, time-series prediction, speech recognition**) require memory of past inputs.
- ▶ RNNs enable **temporal dynamic behavior** by maintaining hidden states.

Examples Where Order Matters:




- ▶ Translating “I am happy” vs “Happy I am”
- ▶ Predicting next stock price based on past trends
- ▶ Understanding a sentence word-by-word

Key Idea

Add a feedback loop to remember previous computations — introducing memory into neural networks.

By the end of this session, you should be able to:

- ▶ Understand the structure and working of Recurrent Neural Networks (RNNs)
- ▶ Recognize various RNN architectures (One-to-Many, Many-to-One, Many-to-Many)
- ▶ Explain the concept of shared parameters in RNNs
- ▶ Evaluate RNN limitations and how advanced models improve them
- ▶ Explore potential future directions of sequential models

	Sequence	$P(\text{Sequence})$
  J'ai vu le match de foot 	I saw the game of soccer	4.5 e-5
	I saw the soccer game	<u>6.0 e-5</u>
	I saw the soccer match	4.6 e-5
	Saw I the game of soccer	2.6 e-9

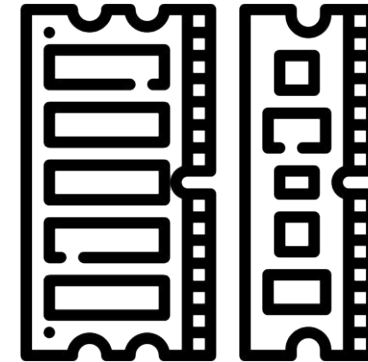
$$P(w_2|w_1) = \frac{\text{count}(w_1, w_2)}{\text{count}(w_1)} \longrightarrow \text{Bigrams}$$

$$P(w_3|w_1, w_2) = \frac{\text{count}(w_1, w_2, w_3)}{\text{count}(w_1, w_2)} \quad \text{Trigrams}$$

$$P(w_1, w_2, w_3) = P(w_1) \times P(w_2|w_1) \times P(w_3|w_2)$$

- Large N-grams to capture dependencies between distant words
- Need a lot of space and RAM

- N-grams consume a lot of memory
- Different types of RNNs are the preferred alternative



What is an RNN?

A neural network with loops — allowing information to persist.

Core Elements:

- ▶ **Hidden State** h_t : captures memory of previous inputs
- ▶ **Input** x_t , **Output** y_t
- ▶ Same weights used across time steps (parameter sharing)

Mathematical Formulation:

$$h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t + b_h)$$

$$y_t = W_{hy}h_t + b_y$$

This recurrence allows information to propagate through time.

Advantages of RNNs

Nour was supposed to study with me. I called her but she **did not** answer

want

respond

choose

want

have

ask

attempt

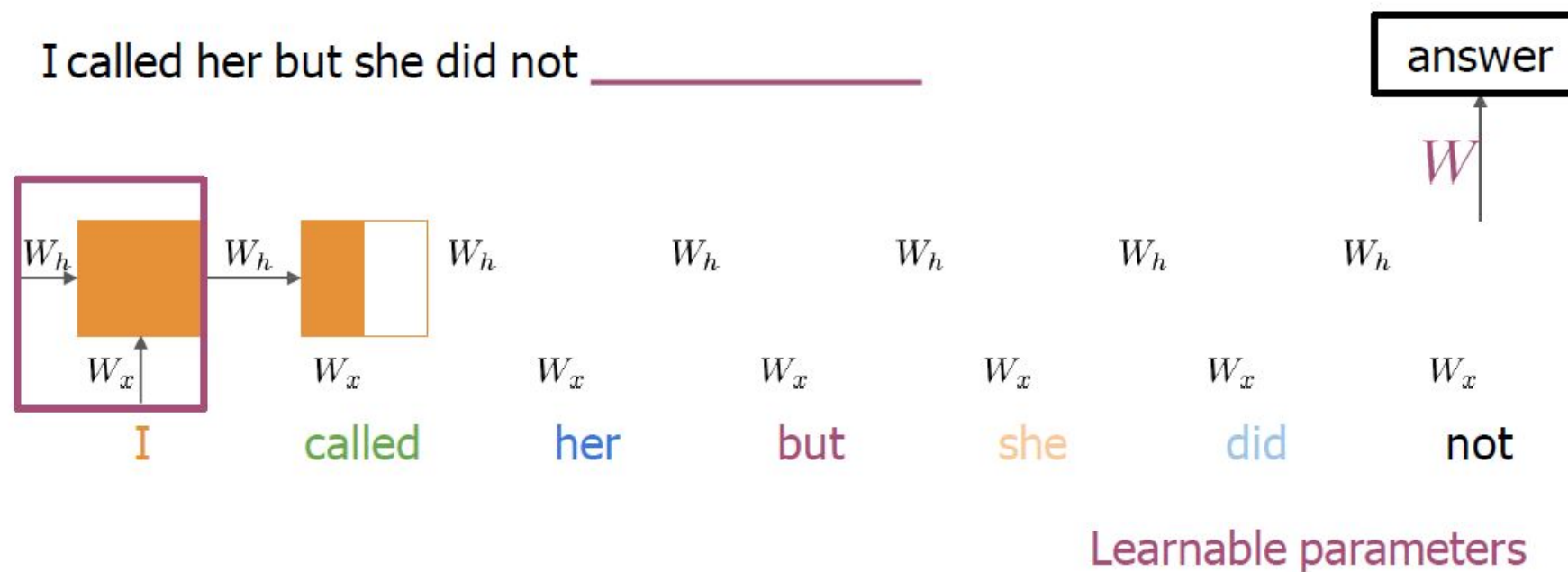
answer

know

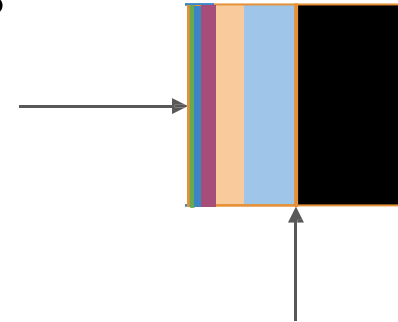
RNNs look at every previous word

Similar probabilities with
trigram

RNNs Basic Structure



- RNNs model relationships among distant words
- In RNNs a lot of computations share parameters



RNN Architecture: Unrolled View

Unrolled RNN:

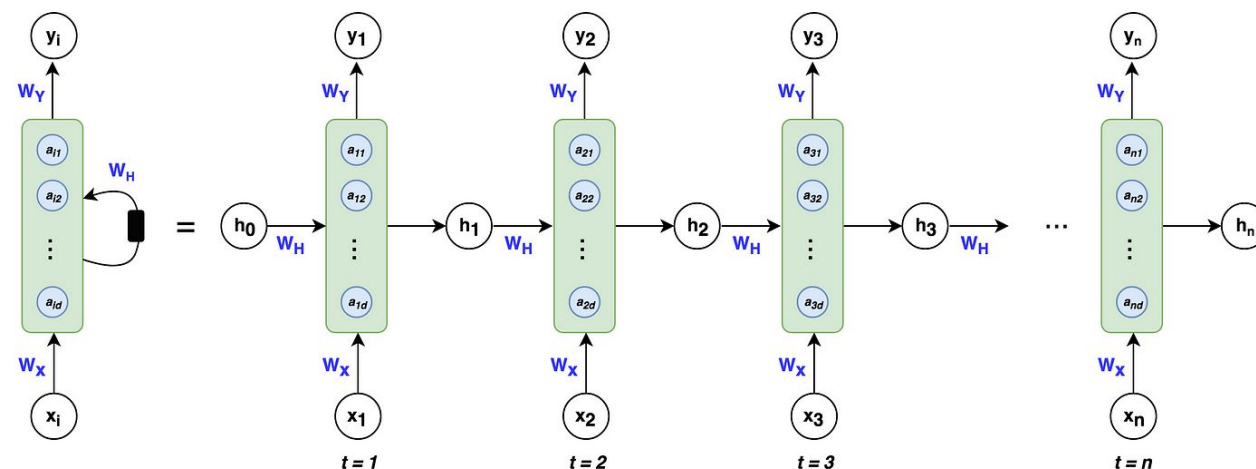
An RNN is essentially a chain of repeating neural network modules, one for each time step.

Each time step shares the same parameters:

- ▶ Input x_t
- ▶ Hidden state h_t
- ▶ Output y_t

Key Idea

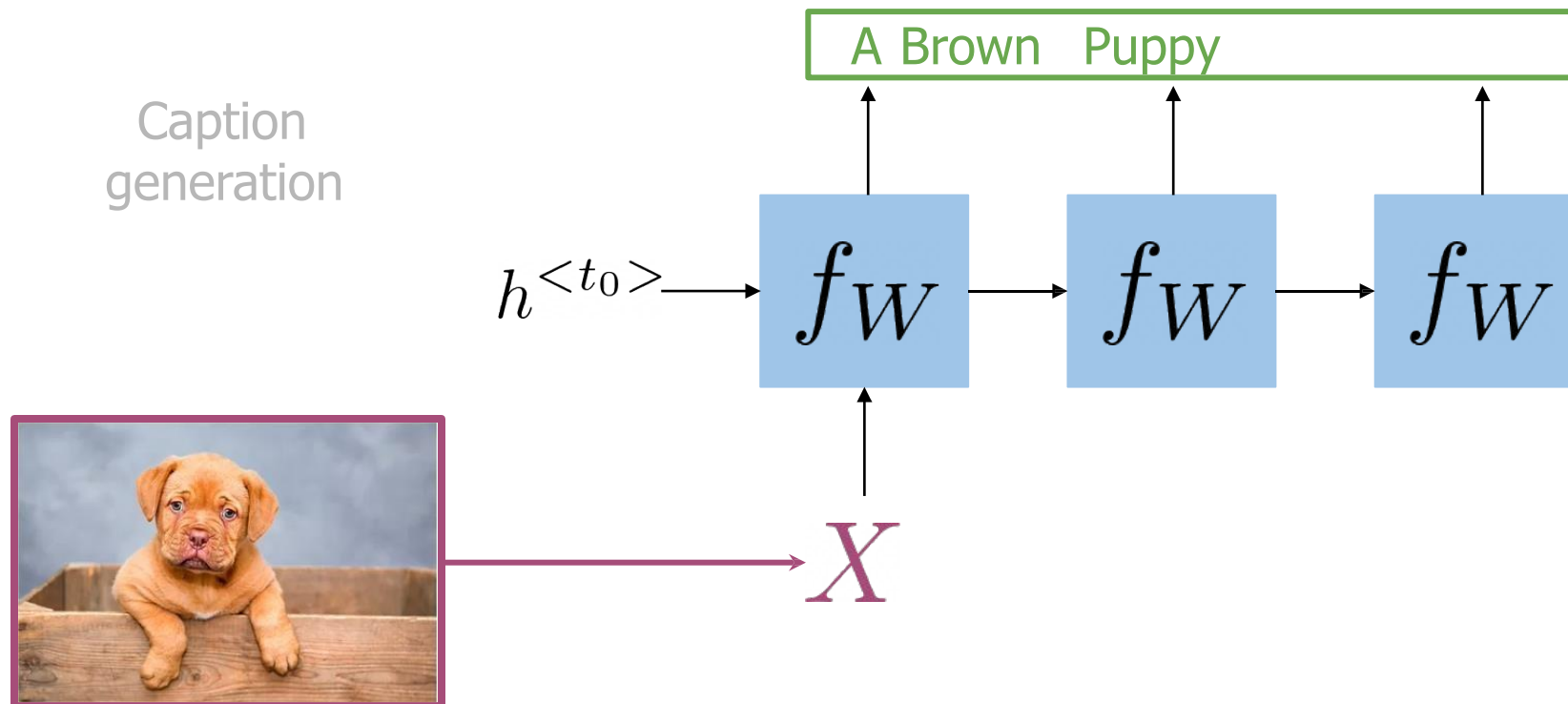
Temporal representation without increasing parameter count!



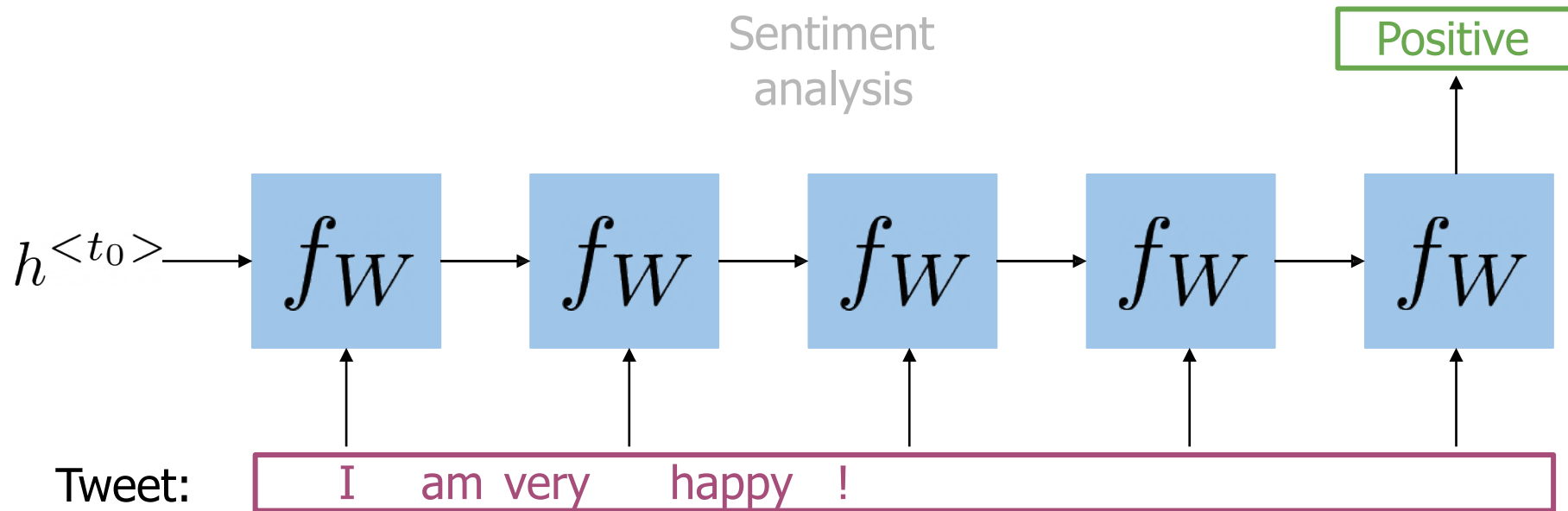
One to One



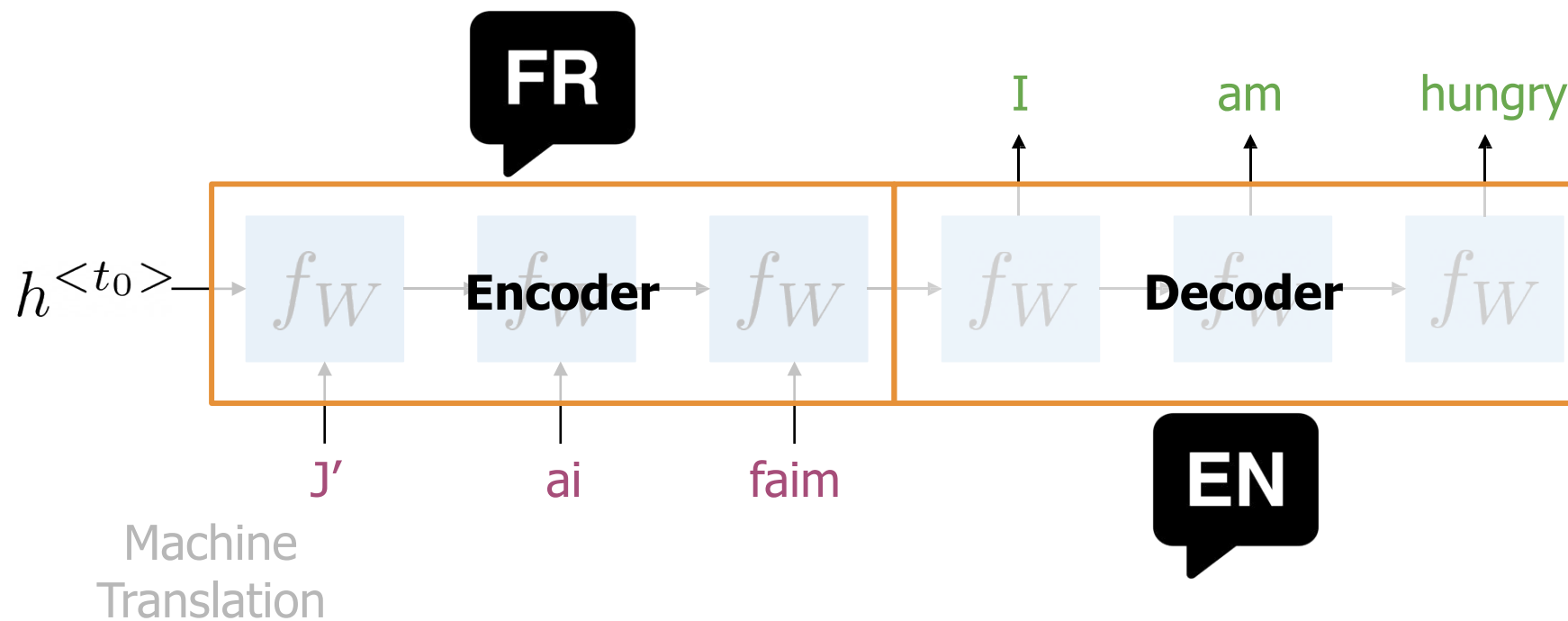
One to Many



Many to One



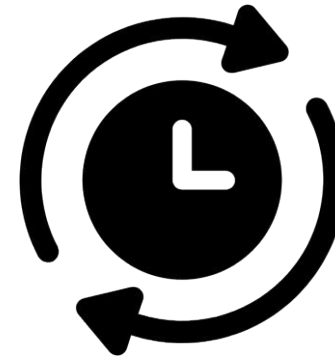
Many to Many



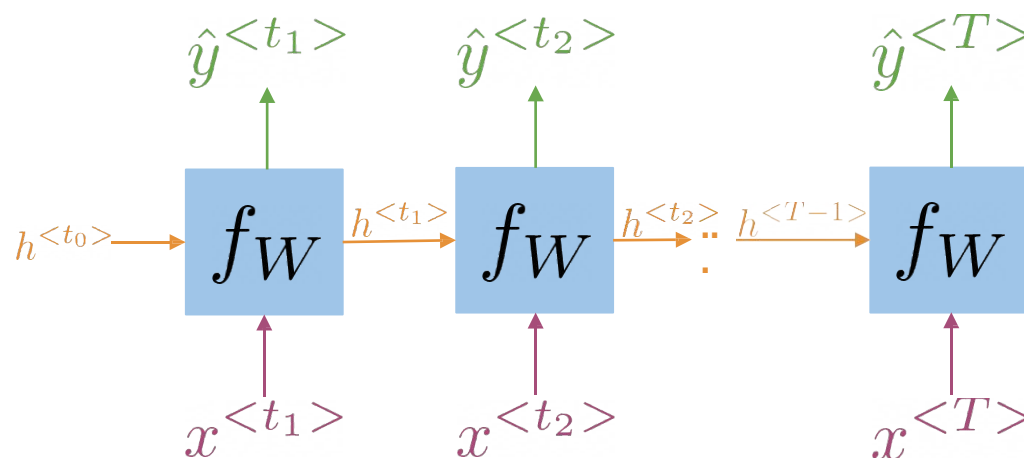
- RNNs can be implemented for a variety of NLP tasks
- Applications include Machine translation and caption generation



- How RNNs propagate information (Through time!)
- How RNNs make predictions



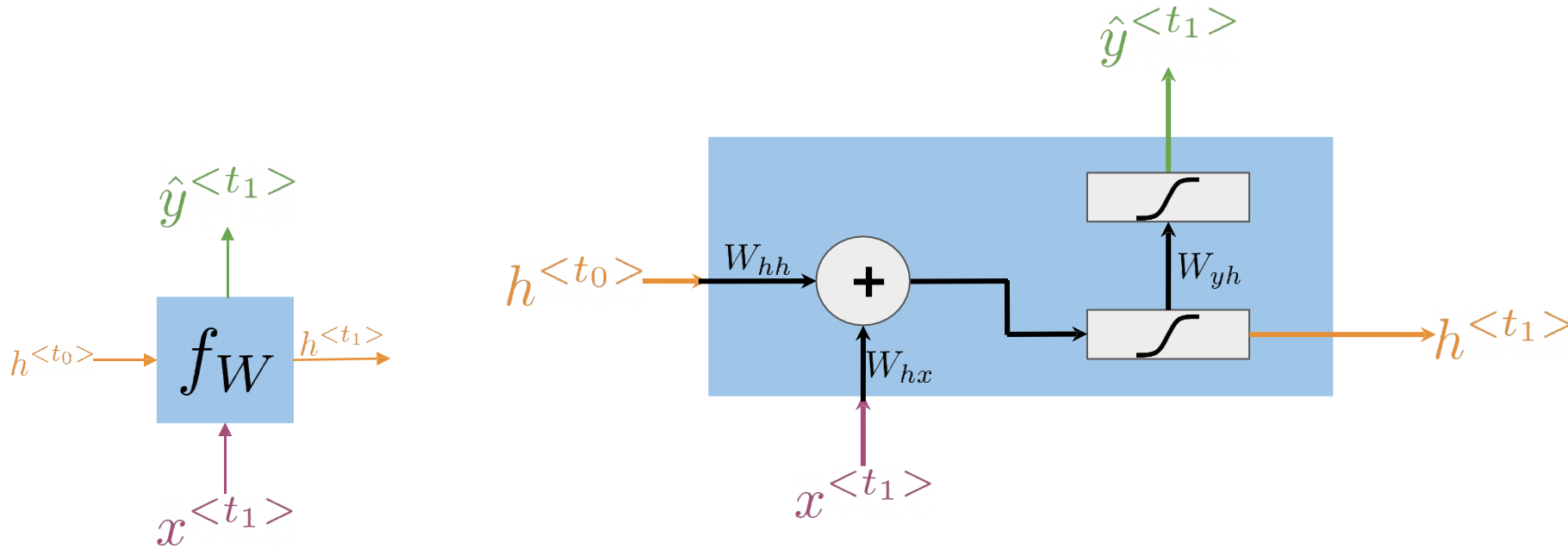
A Vanilla RNN



$$h^{<t>} = g(W_h[h^{<t-1>}, x^{<t>}] + b_h)$$
$$\hat{y}^{<t>} = g(W_{yh}h^{<t>} + b_y)$$

$$h^{<t>} = g(W_{hh}h^{<t-1>} + W_{hx}x^{<t>} + b_h)$$

A Vanilla RNN



$$h^{<t>} = g(W_{hh}h^{<t-1>} + W_{hx}x^{<t>} + b_h)$$
$$\hat{y}^{<t>} = g(W_{yh}h^{<t>} + b_y)$$

Backpropagation Through Time (BPTT):

- ▶ Training RNNs involves **unfolding** the network across time steps.
- ▶ Standard backpropagation is applied through this unrolled structure.

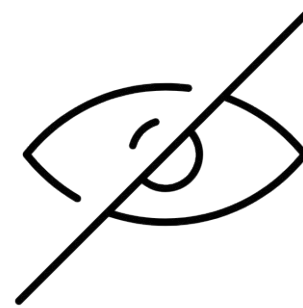
Problems:

- ▶ **Vanishing Gradients:** Gradients shrink as they are propagated back, making it hard to learn long-term dependencies.
- ▶ **Exploding Gradients:** Gradients grow exponentially, leading to unstable updates.

Solutions Preview (covered in future modules):

- ▶ Use of LSTM and GRU architectures
- ▶ Gradient Clipping

- Hidden states propagate information through time
- Basic recurrent units have two inputs at each time: $h^{<t-1>}$ $x^{<t>}$



Natural Language Processing

- ▶ Language modeling
- ▶ Named Entity Recognition
- ▶ Machine Translation

Time-Series Forecasting

- ▶ Stock prediction
- ▶ Weather forecasting

Audio Processing

- ▶ Speech recognition
- ▶ Music generation

Cognitive Modeling

- ▶ Simulating memory in brain-like systems

- ▶ Sequential computation — hard to parallelize
- ▶ Forget long-term dependencies
- ▶ Slow training due to sequential nature
- ▶ Struggle with varying-length sequences

Key Developments

- ▶ **LSTM and GRU:** Designed to address memory and gradient issues
- ▶ **Transformers:** Non-recurrent, highly parallelizable models

Future Directions: What's Beyond Vanilla RNNs?

- ▶ **LSTM (Long Short-Term Memory):** Overcomes vanishing gradients
- ▶ **GRU (Gated Recurrent Unit):** Simpler than LSTM, efficient gating
- ▶ **Attention Mechanisms:** Focus on relevant parts of the input sequence
- ▶ **Transformers & Self-Attention:** Replace recurrence with parallelizable attention
- ▶ **Neural ODEs:** Model continuously evolving hidden states

Hybrid Models:

- ▶ Combine RNNs, CNNs, and Attention for complex tasks (e.g., video, multimodal text)

- ▶ **RNNs introduce memory into neural nets for sequence modeling**
- ▶ **Use shared weights across time steps**
- ▶ **Architectures like One-to-Many, Many-to-One fit various tasks**
- ▶ **RNNs face training and memory limitations**
- ▶ **Advances like LSTM, GRU, and Transformers push beyond RNNs**

These slides have been adapted from

- Younes Mourri & Lukasz Kaiser, [Natural Language Processing Specialization, DeepLearning.AI](#)

Core Papers:

- ▶ Elman, J. L. (1990). Finding structure in time. *Cognitive Science*.
- ▶ Hochreiter, S., & Schmidhuber, J. (1997). Long Short-Term Memory. *Neural Computation*.
- ▶ Mikolov, T., Karafiát, M., Burget, L., Cernocký, J., & Khudanpur, S. (2010). Recurrent neural network based language model. *Interspeech*.
- ▶ Bengio, Y., Simard, P., & Frasconi, P. (1994). Learning long-term dependencies with gradient descent is difficult. *IEEE Transactions on Neural Networks*.
- ▶ Sutskever, I., Vinyals, O., & Le, Q. V. (2014). Sequence to Sequence Learning with Neural Networks. *NeurIPS*.