

RECURRENT NEURAL NETWORKS

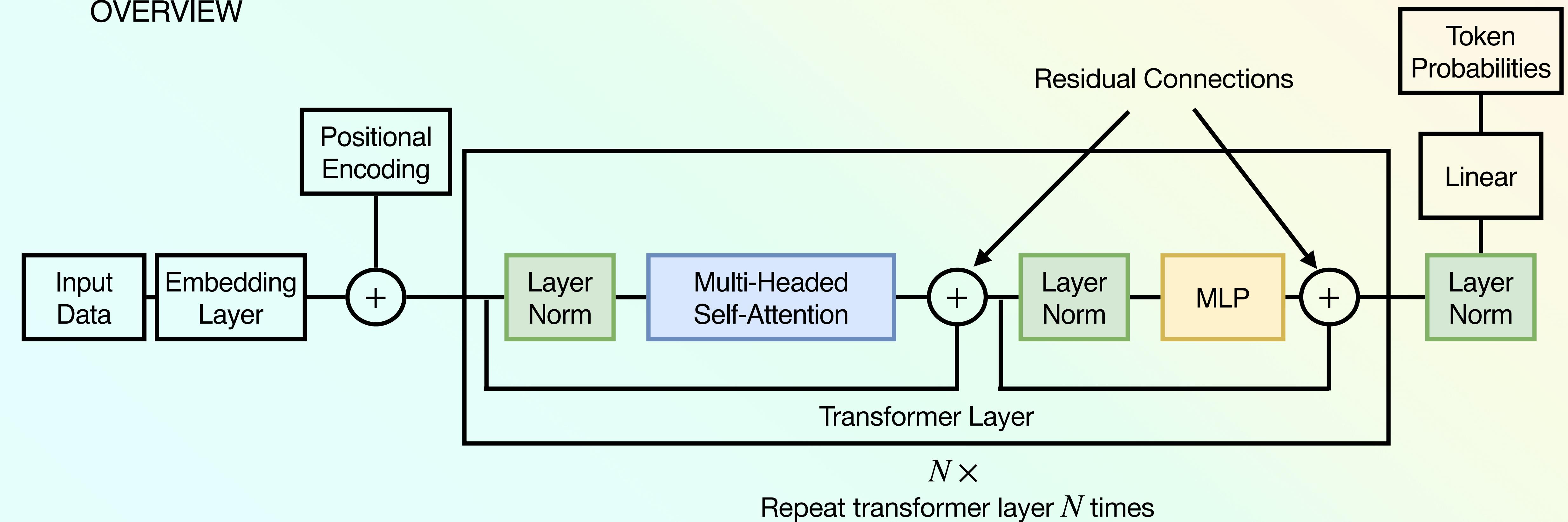
TODAY'S CLASS

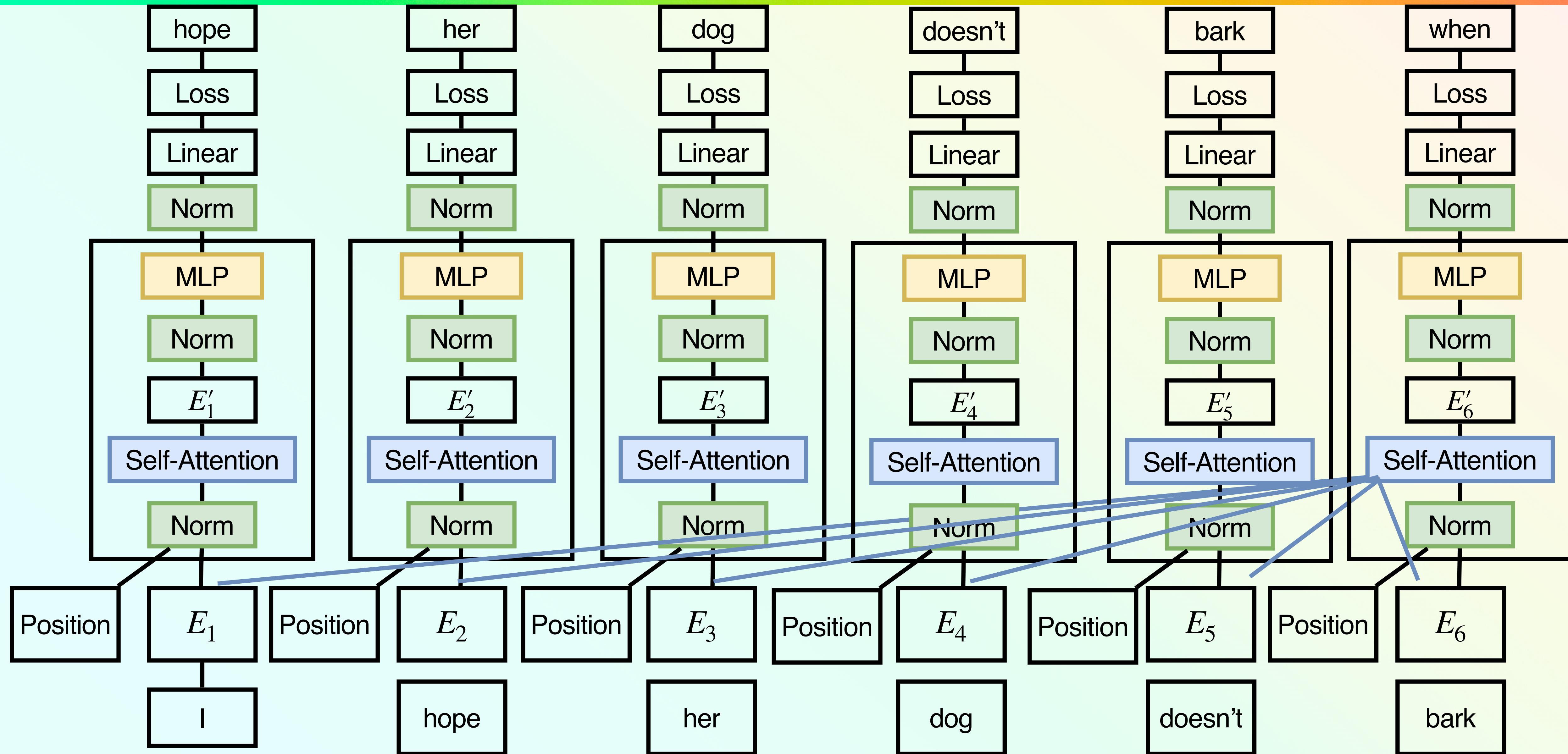
GOALS

1. Review transformers and discuss a few properties
2. Motivate RNNs, Define them, give the basics, and their problems
3. Cover LSTMs and what they do better than RNNs

TRANSFORMERS

OVERVIEW





GPT

DEFINITION

Generative Pretrained Transformer

Train on a MASSIVE dataset

Model predictions are “generally” “good”

Fine-tune on new tasks

Pretraining is a smart weight initialization for the new task

Chatbots have extra steps

QUIZ

TRANFORMERS

فَلَمَّا
أَتَاهُمْ
الْكِتَابَ،
لَمْ يَنْظُرُوا

HOW THEY PROCESS INFORMATION

- ⌚ Self-attention – brings in information from other points in the sequence
- ⌚ MLP – performs nonlinear operations to combine and transform this information
- ⌚ More attention heads → more information from the sequence can be brought into each layer
- ⌚ Bigger MLP → more processing that can happen at each layer
- ⌚ Multiple transformer layers – sequentially process more information about the sequence
- ⌚ The last layer must contain all relevant information from the sequence to predict the next token

TRANFORMERS

HOW THEY PROCESS INFORMATION

A transformer with T layers can perform the same computations as a finite Autonoma on a T -length sequence.

Short-cut solutions exist and are easily learned with standard training paradigms

$O(\log T)$ depth solution always exists

$O(1)$ solutions are common

These shortcut solutions are brittle to out-of-distribution sequences:

Training on math operations for one and two-digit numbers, then trying for three-digit numbers

QUIZ

TRANFORMERS

HOW THEY PROCESS INFORMATION

Some operations are hard to perform “ $3+4\times 10+6=$ ”

Predicting 49 as the next two tokens is hard.

The network has to compute the answer and determine that the right thing to output is just the number.

TRANFORMERS

HOW THEY PROCESS INFORMATION

Some operations are hard to perform “ $3+4\times 10+6=$ ”

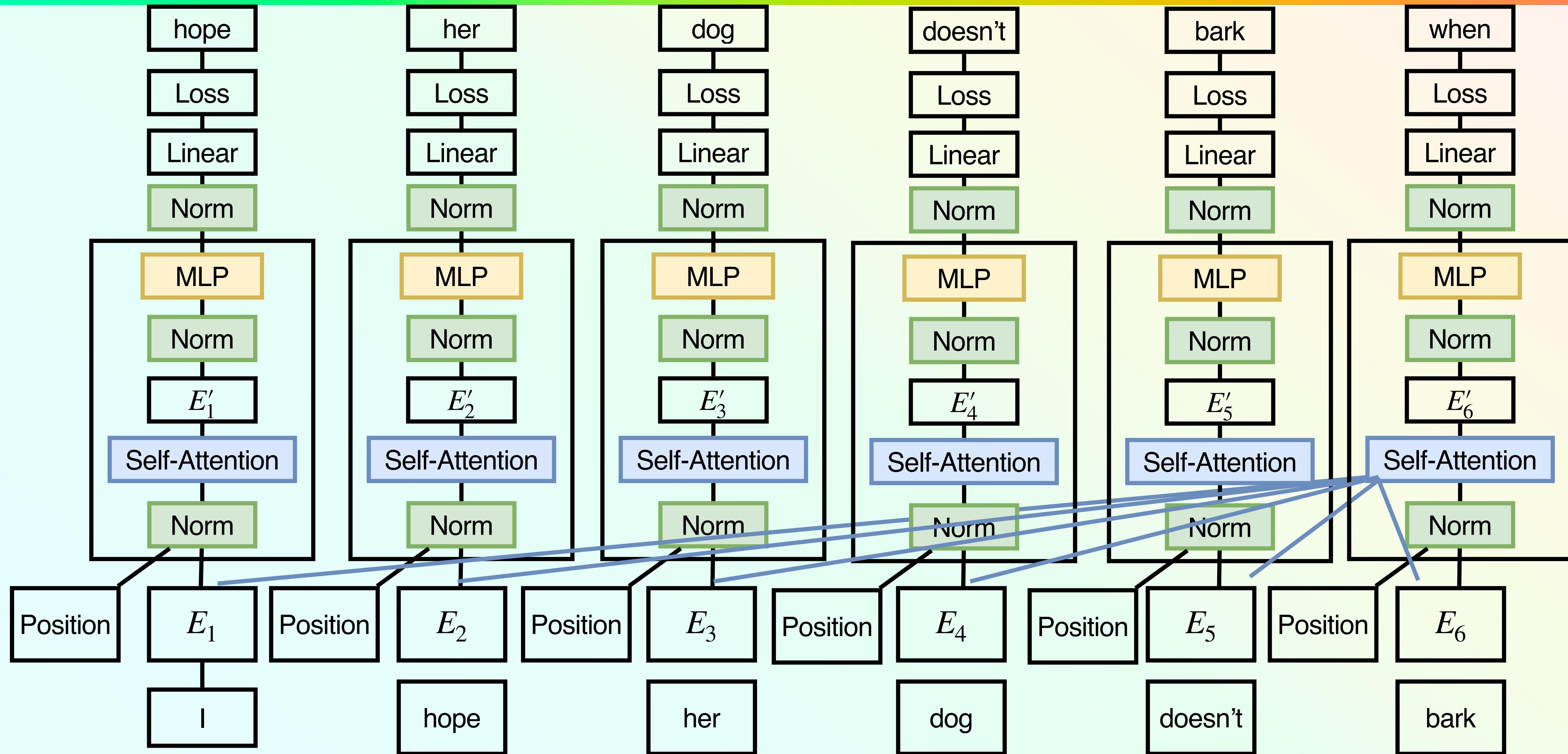
Predicting 49 as the next two tokens is hard.

The network has to compute the answer and determine that the right thing to output is just the number.

Let the network show its work so the computation is easier:

“What is $3+4\times 10+6$? Show your work?” —> ok so we first multiply 4 with 10, which equals 40...

Allows the network to use its output as working memory. – Scratch-pad training



TRANSFORMERS

COMPUTATIONAL COST

Self-attention requires $O(T^2)$ computation for sequence length $T \leftarrow$ *Expensive*

Long context is important for considering huge amounts of information or staying on topic

$O(T^2)$ requires too much computation and memory

TRANSFORMERS

COMPUTATIONAL COST

Self-attention requires $O(T^2)$ computation for sequence length T

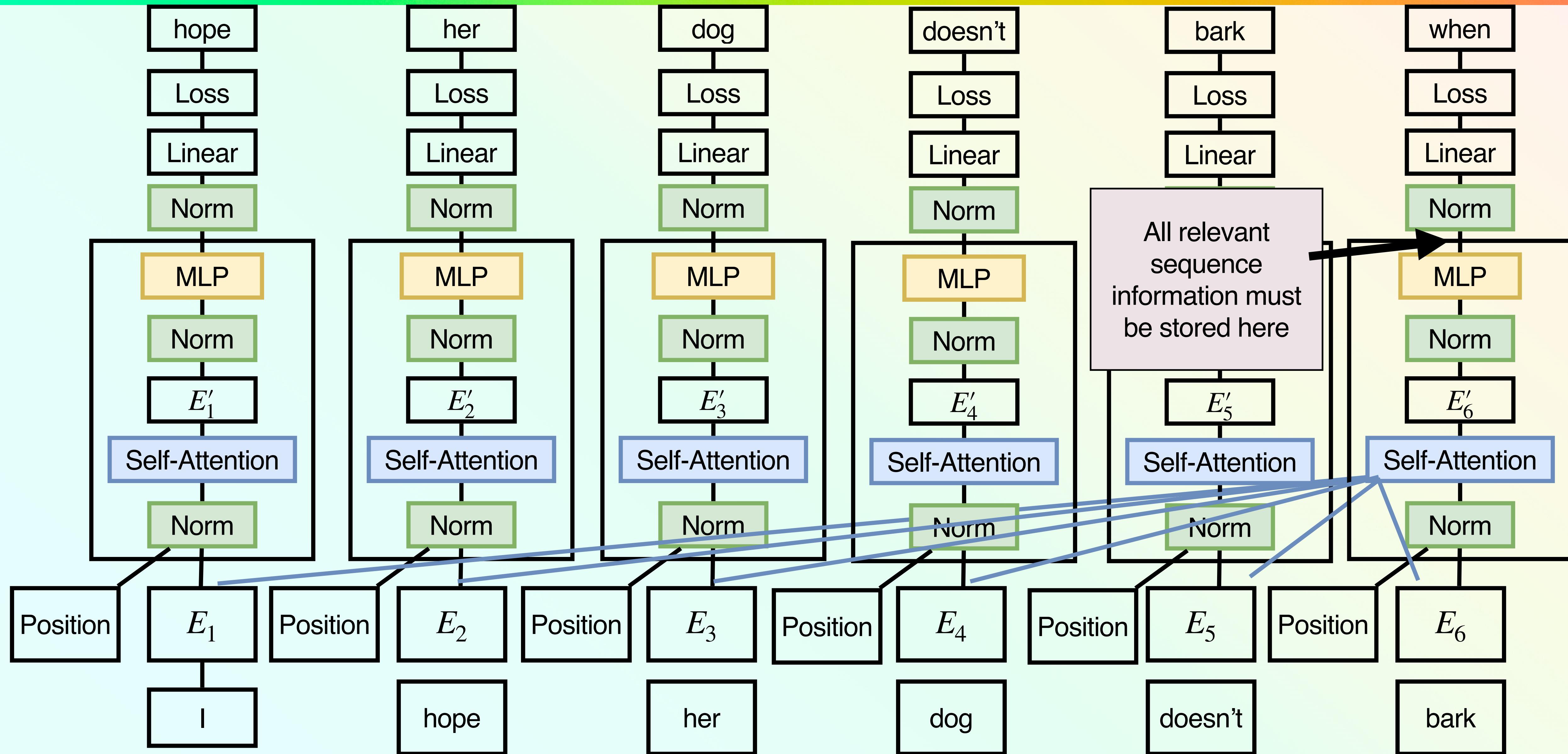
Long context is important for considering huge amounts of information or staying on topic

$O(T^2)$ requires too much computation and memory

Solution:

Attention mechanisms that scale better with T , e.g., $T \log(T)$

Different kinds of neural networks that have constant $O(1)$ memory and computation



RECURRENT NEURAL NETWORKS

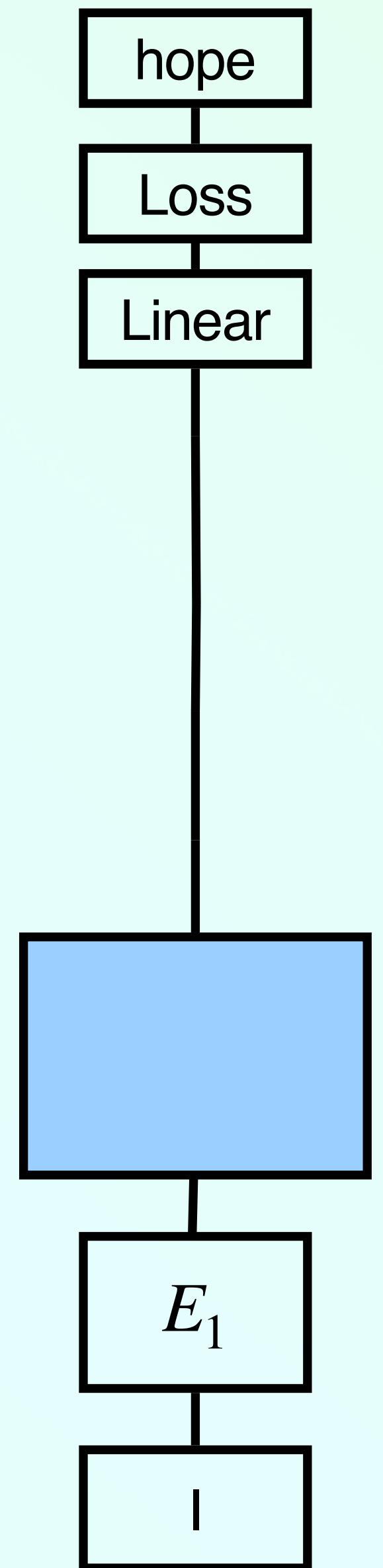
IDEA

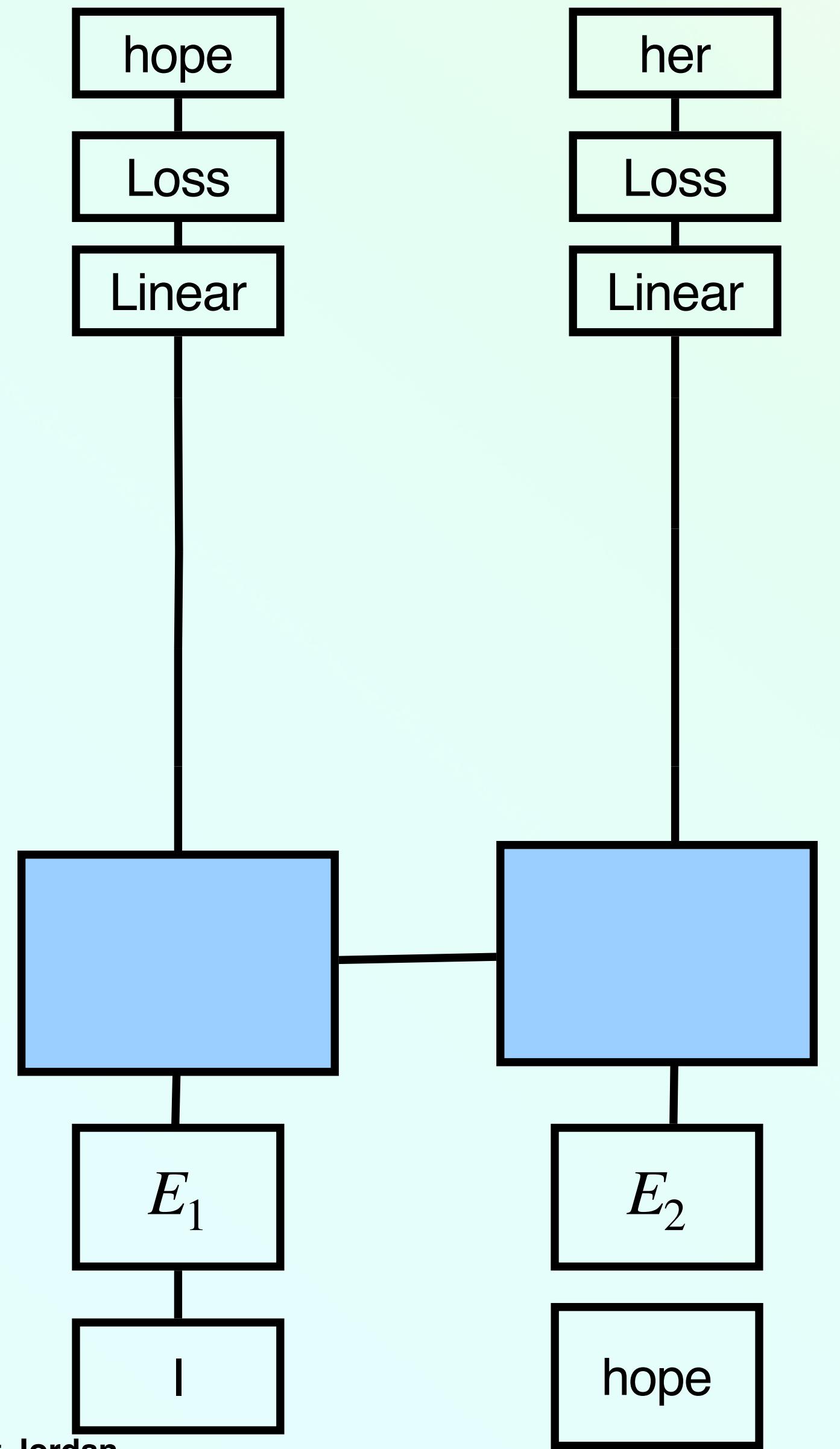
Transformer:

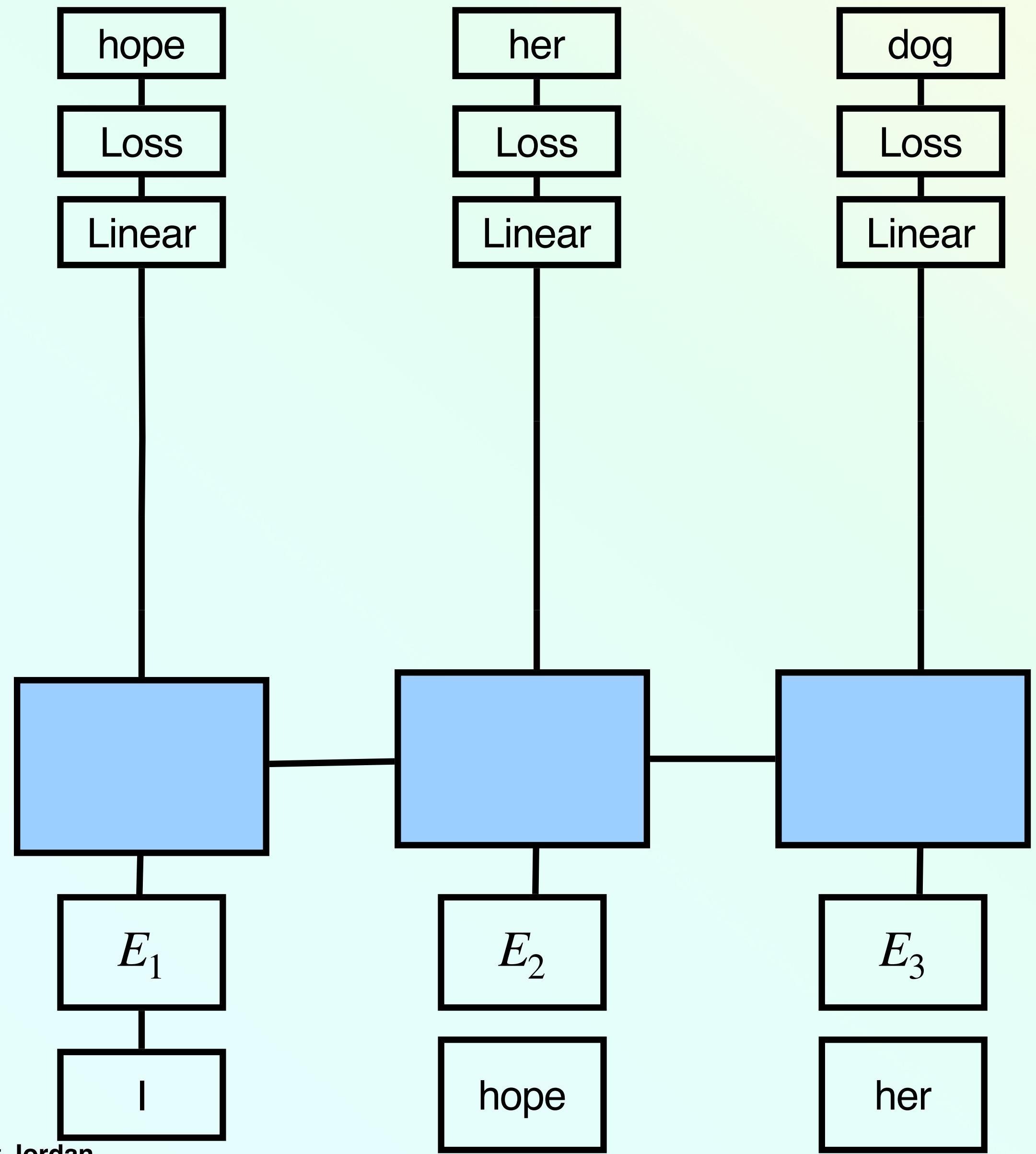
- sequentially include more information with each layer
- last layer includes all relevant information for **just** the next token prediction

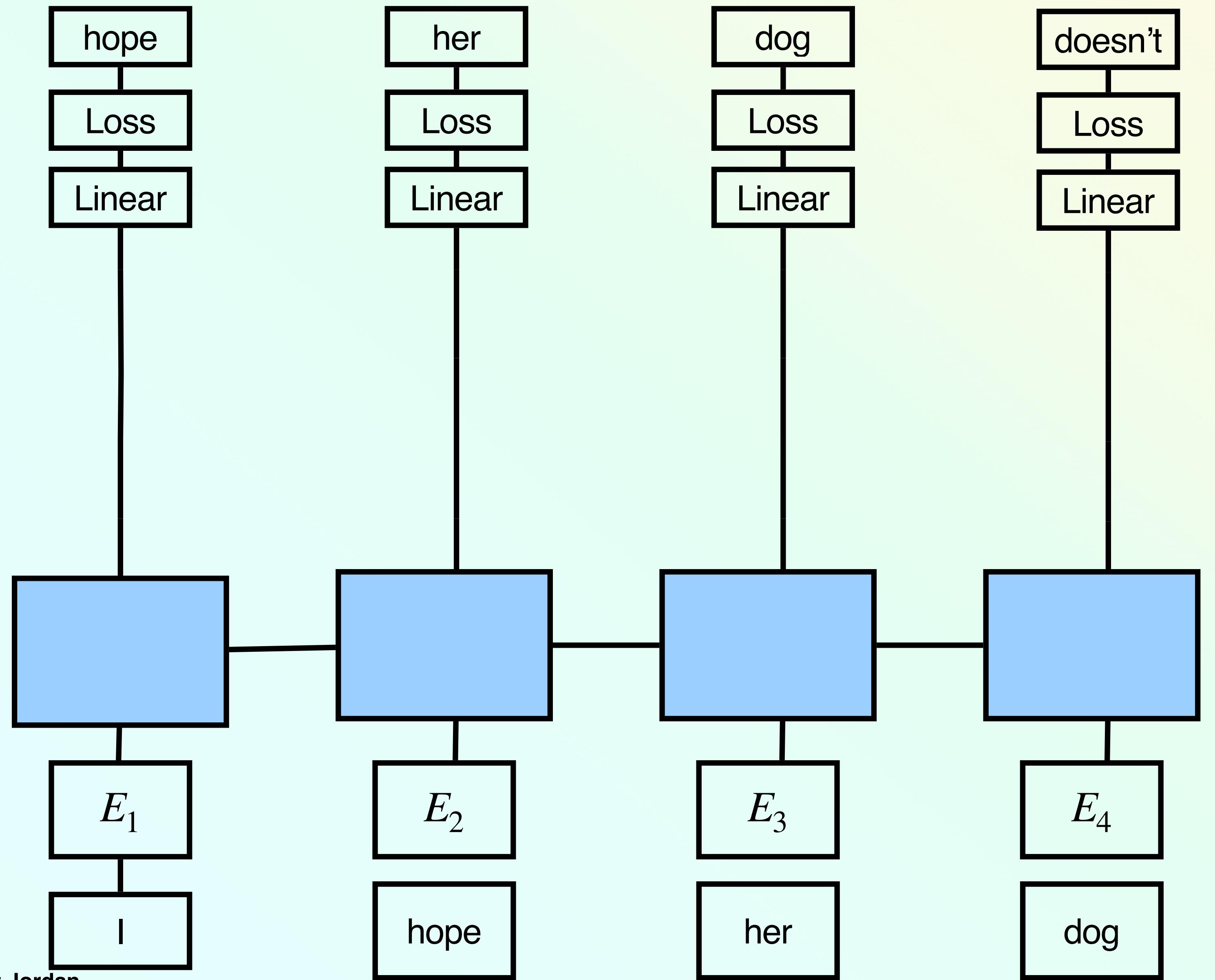
Idea:

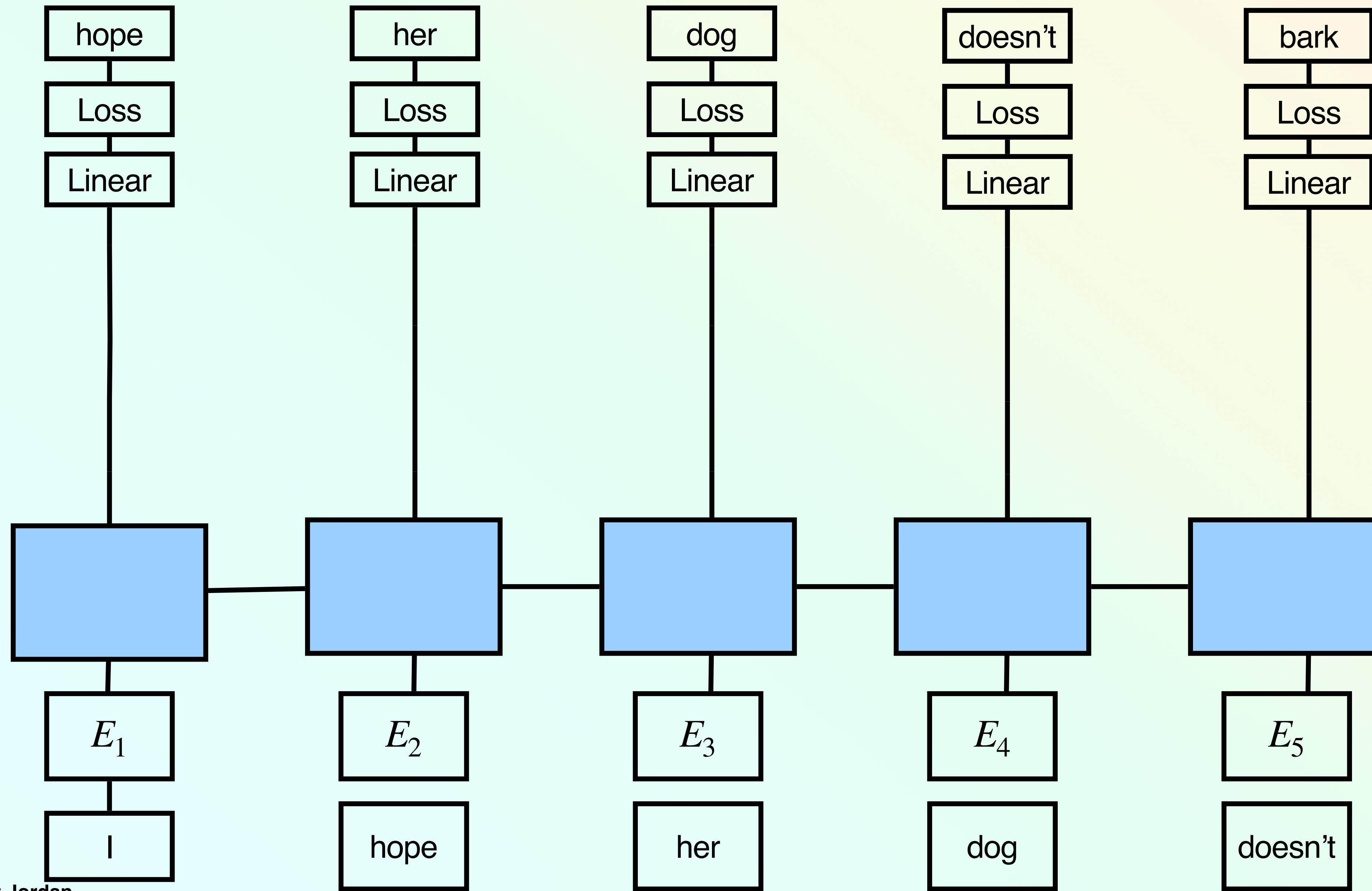
Have a layer that remembers all relevant information for **all** the next tokens

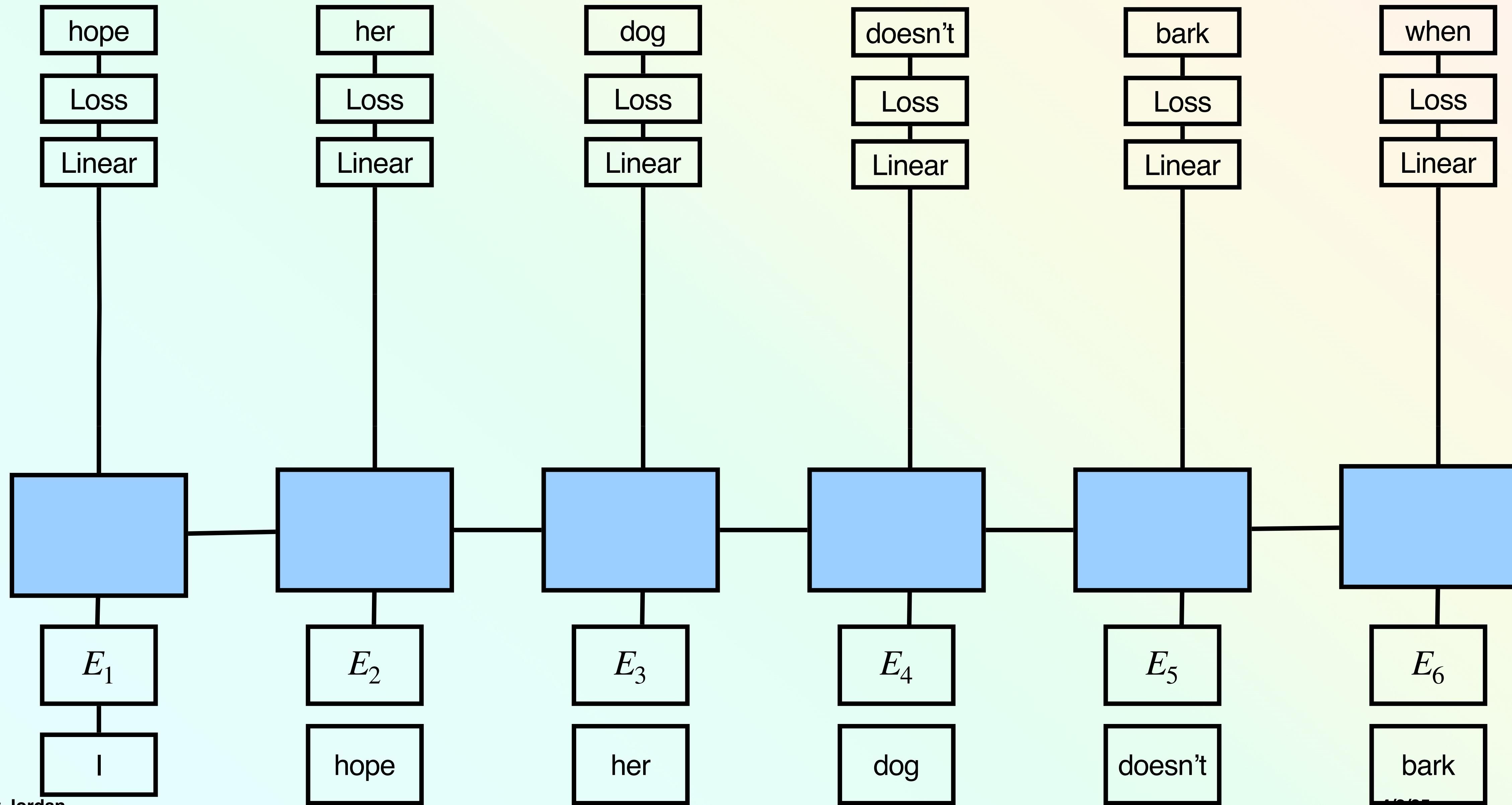


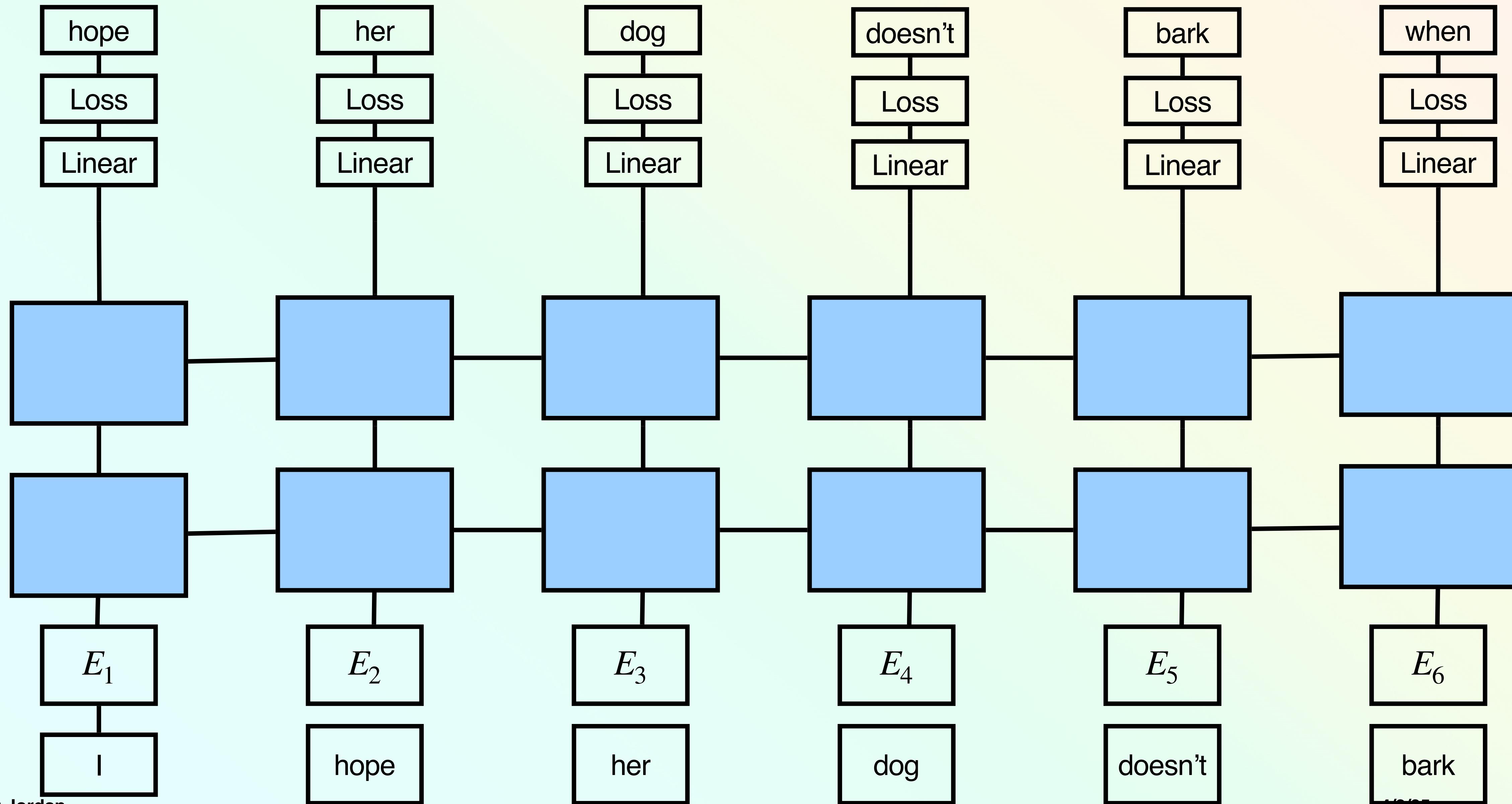


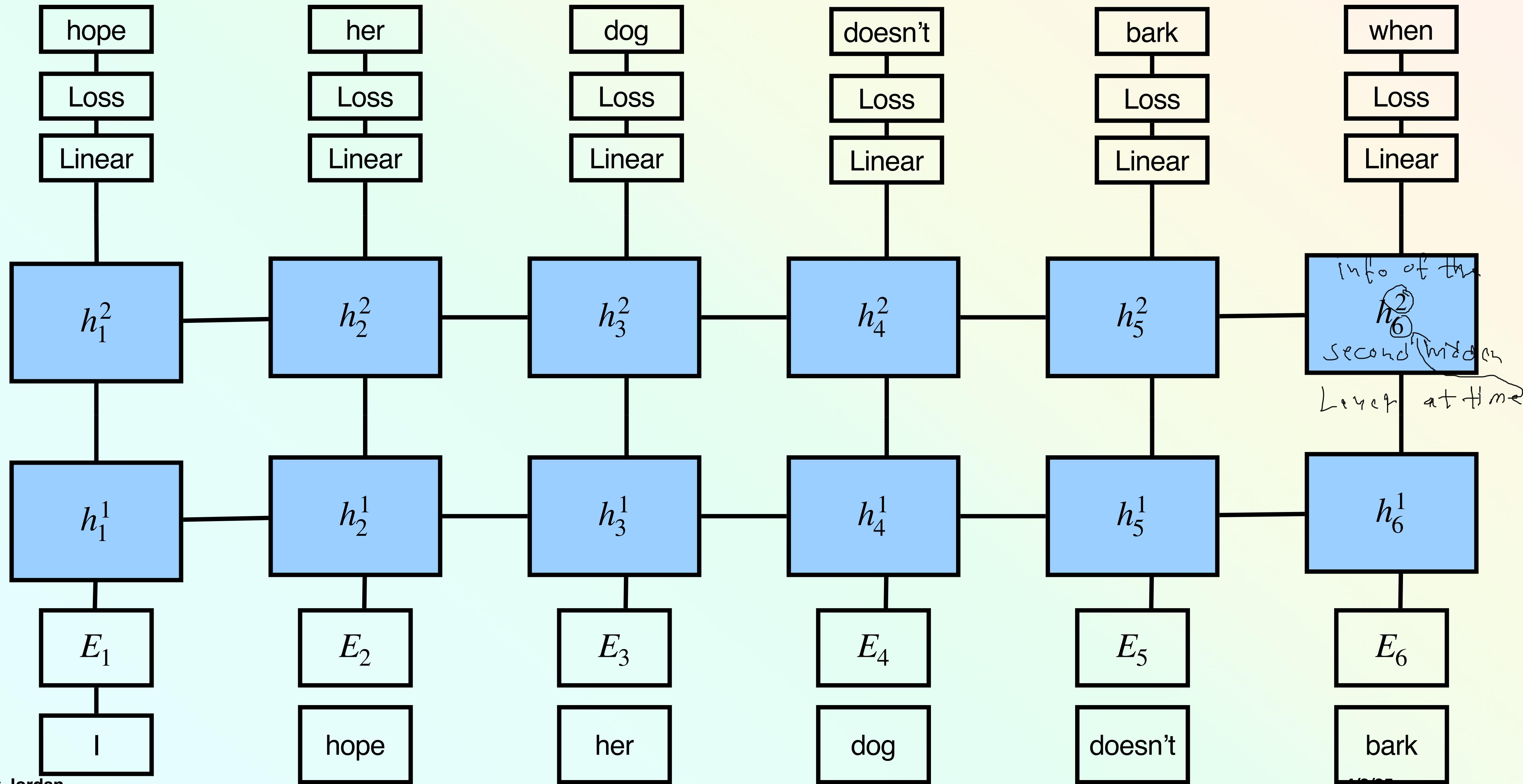






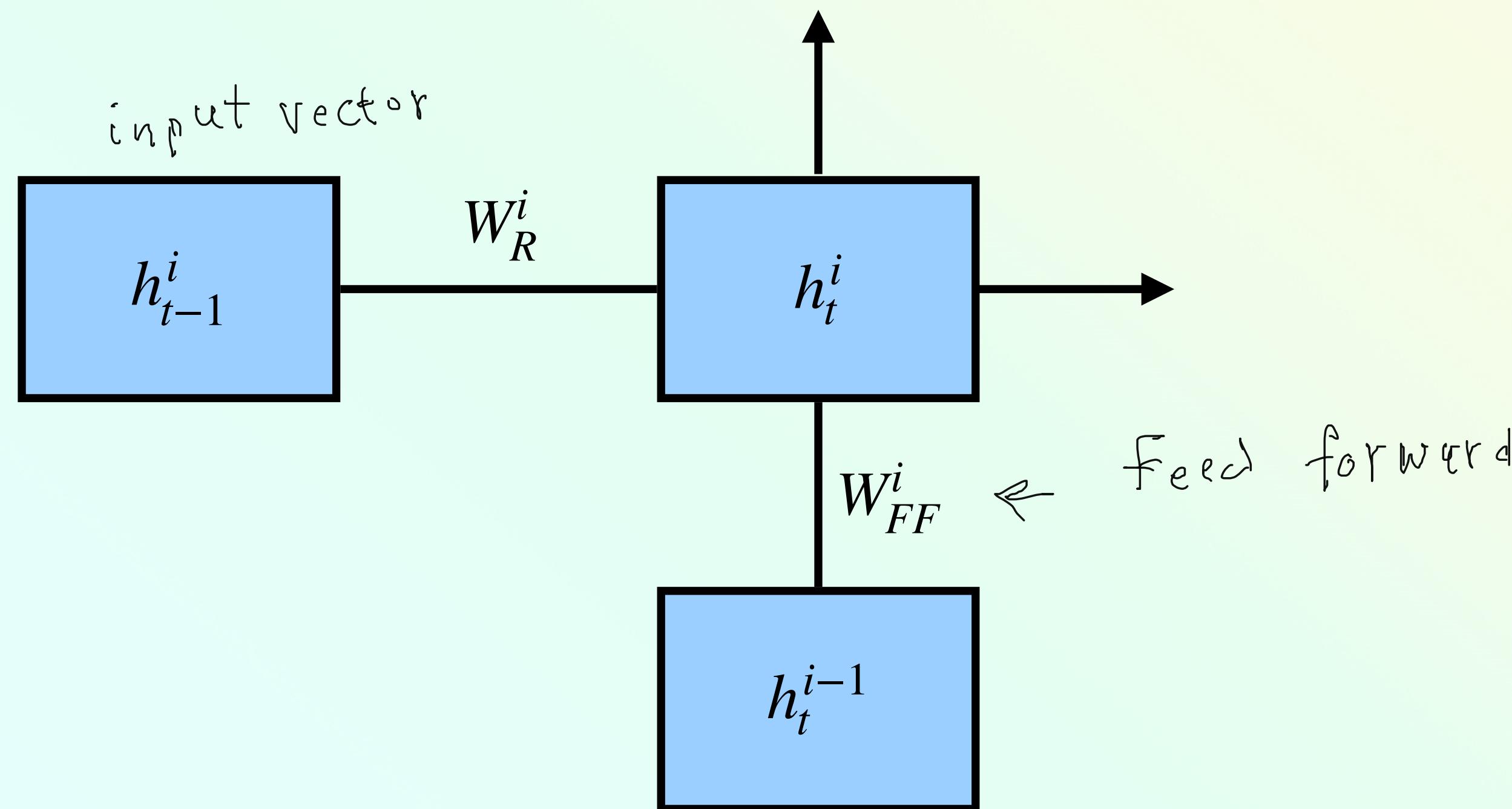






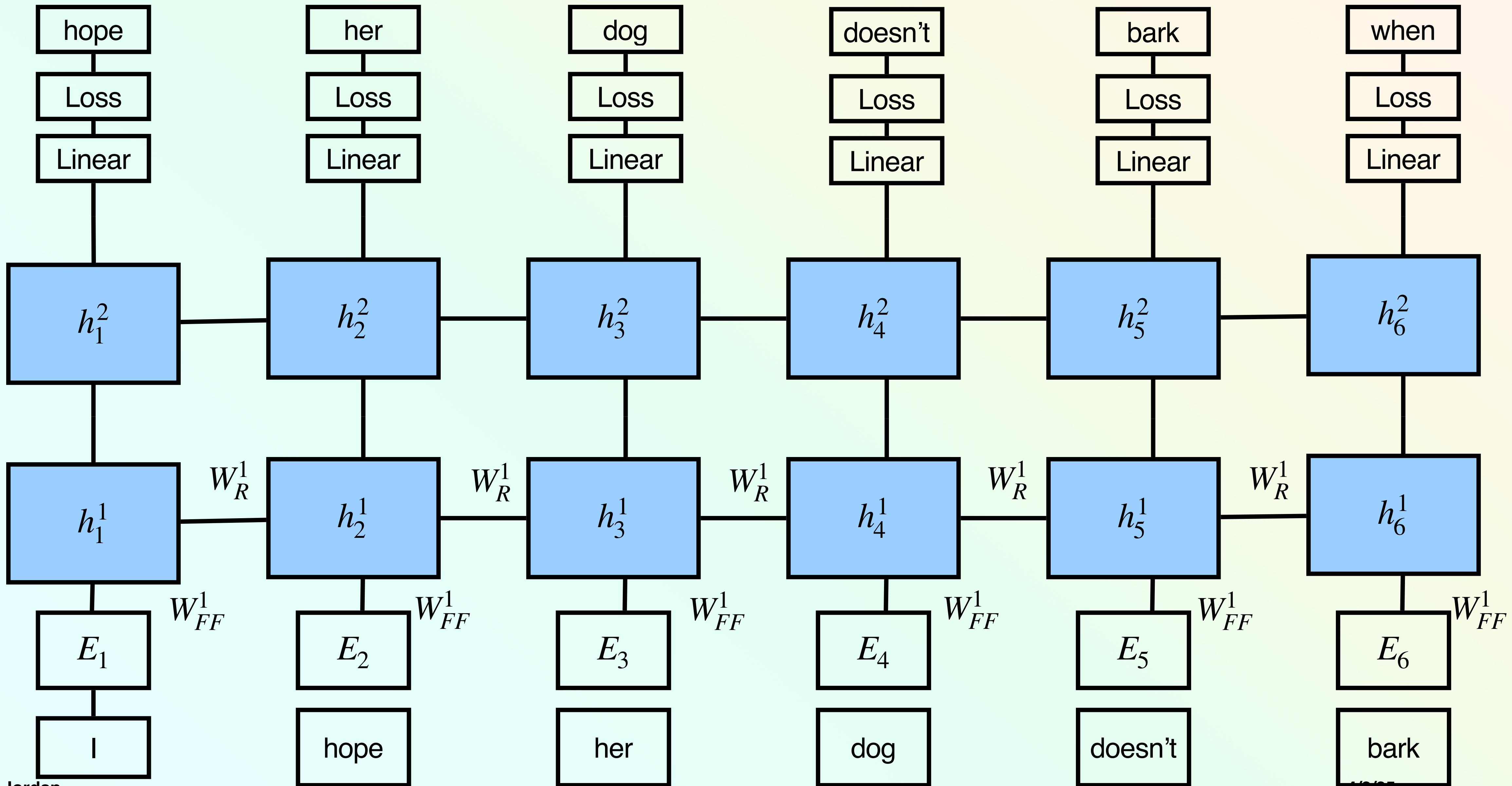
RNNs

LAYER

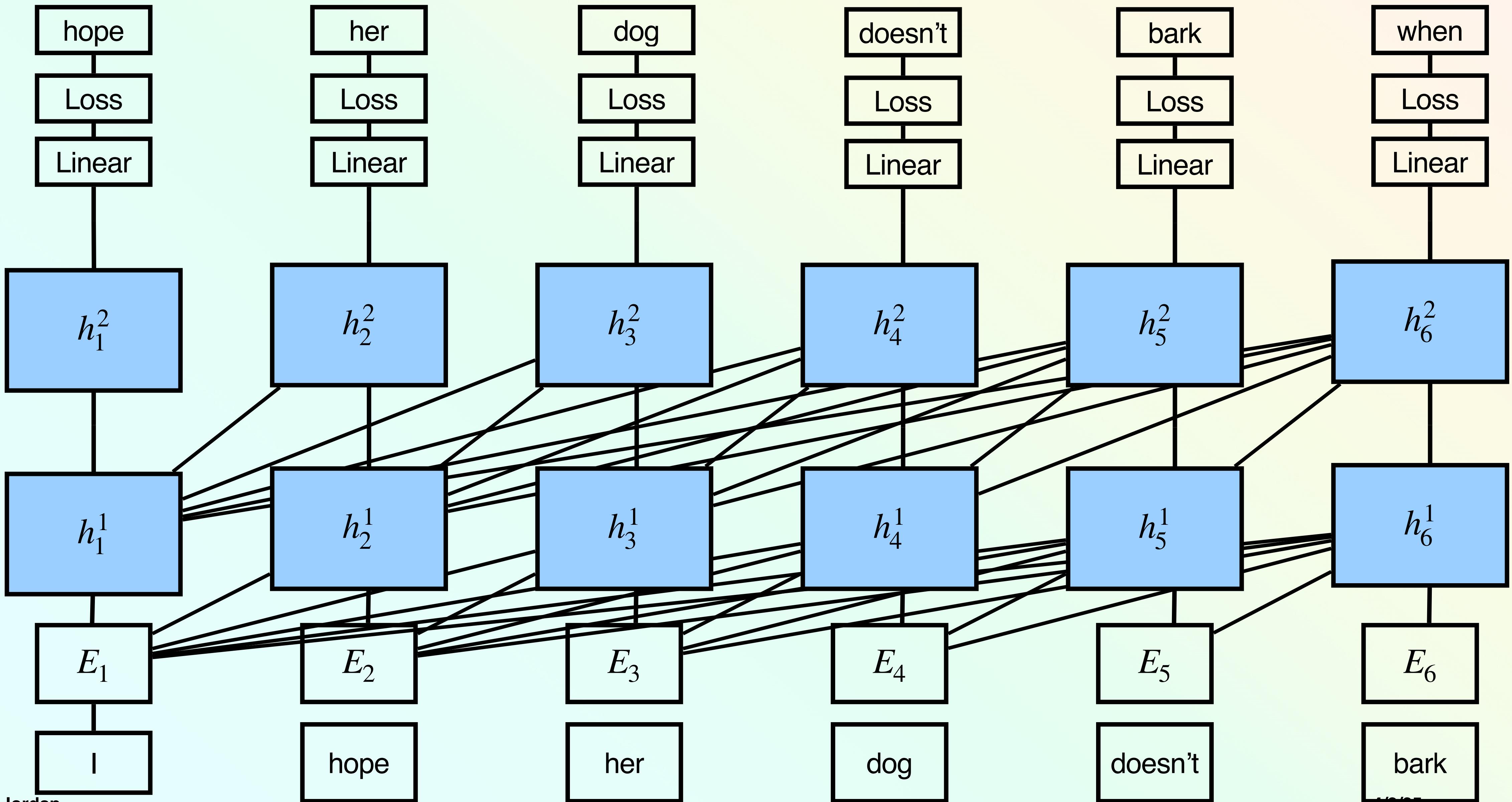


$$h_t^i = f^i(h_t^{i-1}, h_{t-1}^i, W_{FF}^i, W_R^i) = \sigma\left(h_t^{i-1}W_{FF}^{i\top} + h_{t-1}^iW_R^{i\top}\right)$$

RNN Block Diagram



Transformer Block Diagram



RNNs

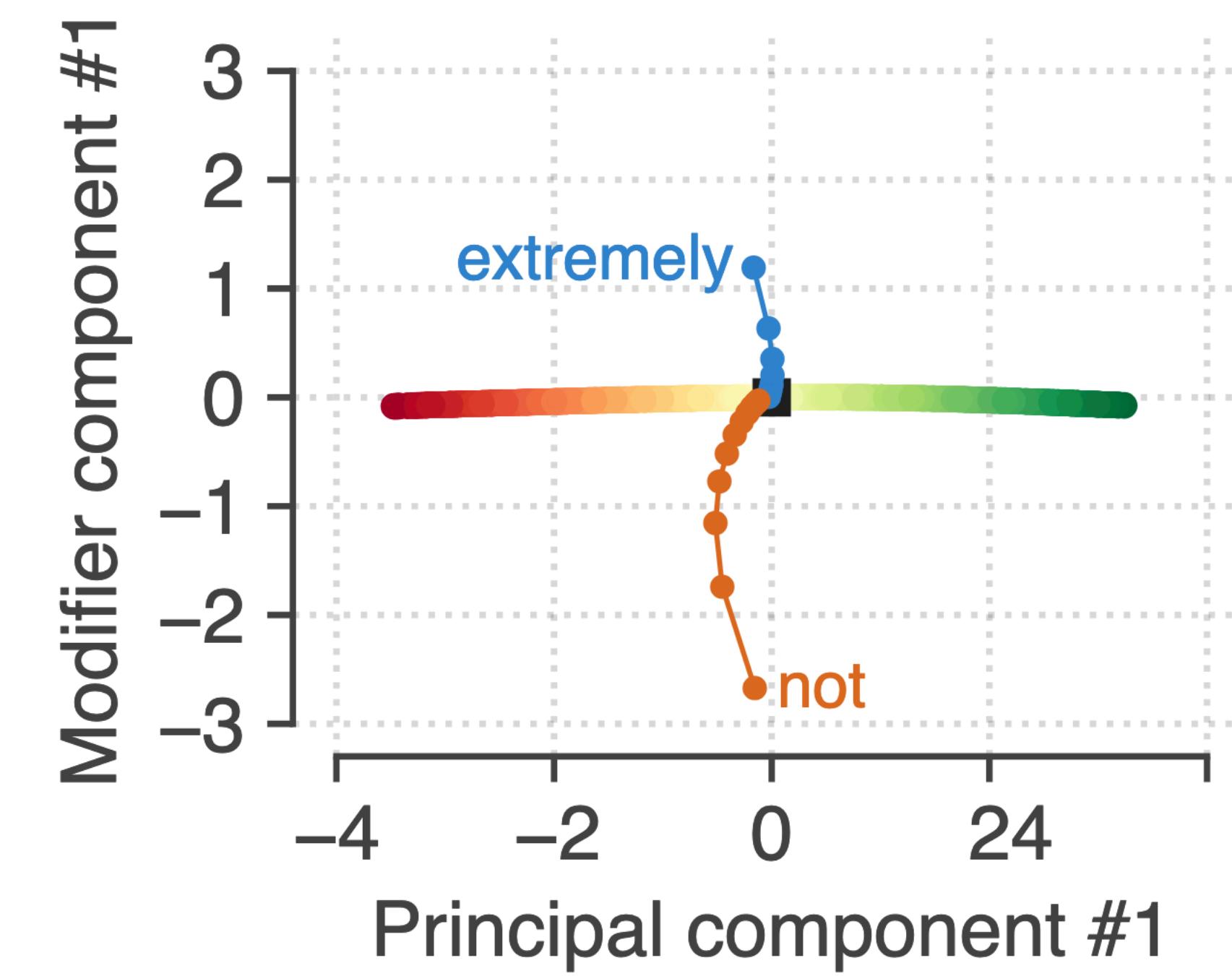
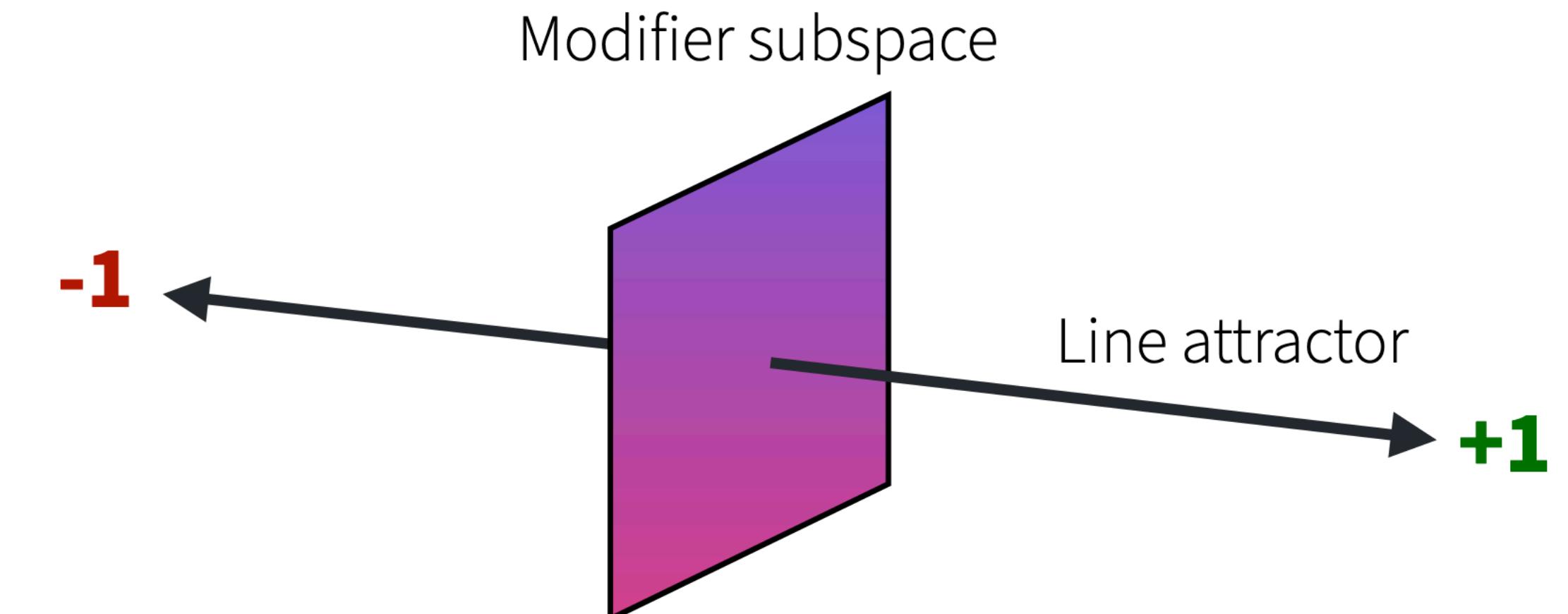
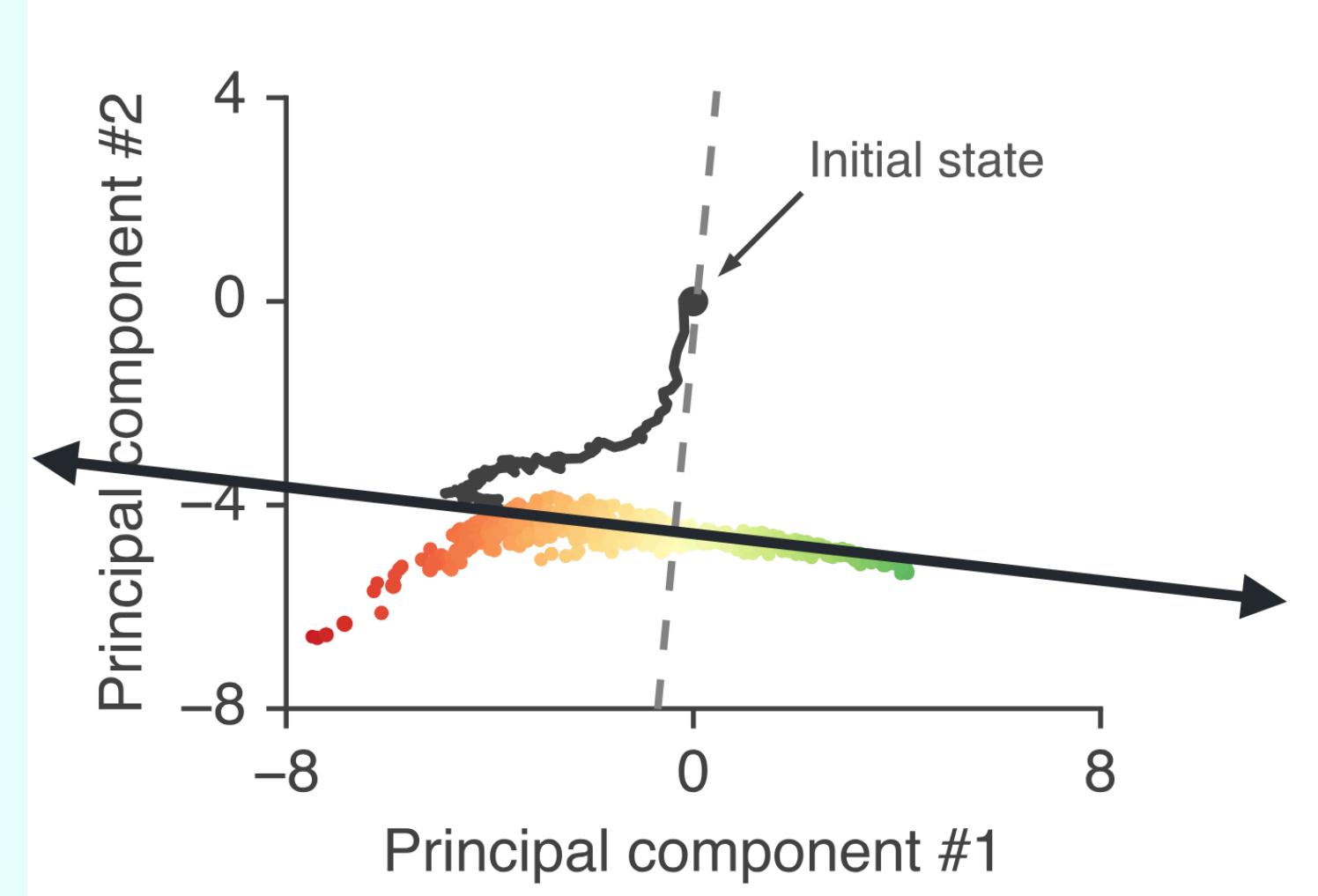
PROPERTIES

- Constant memory and computation requirements at inference time
 - It only needs to know h_{t-1} not all $h_{k < t}$
- Essentially, an MLP layer that is repeated at every time steps
- Need to learn some h_0^i (can be set to 0) \leftarrow could be an extra parameter to fine tune
- The hidden layer has to learn to remember all relevant information for the future
 - This is hard
 - “thinking” about all the information needed in the future all the time

RNNs

WHAT THEY LEARN

Trained for sentiment analysis



RNNs

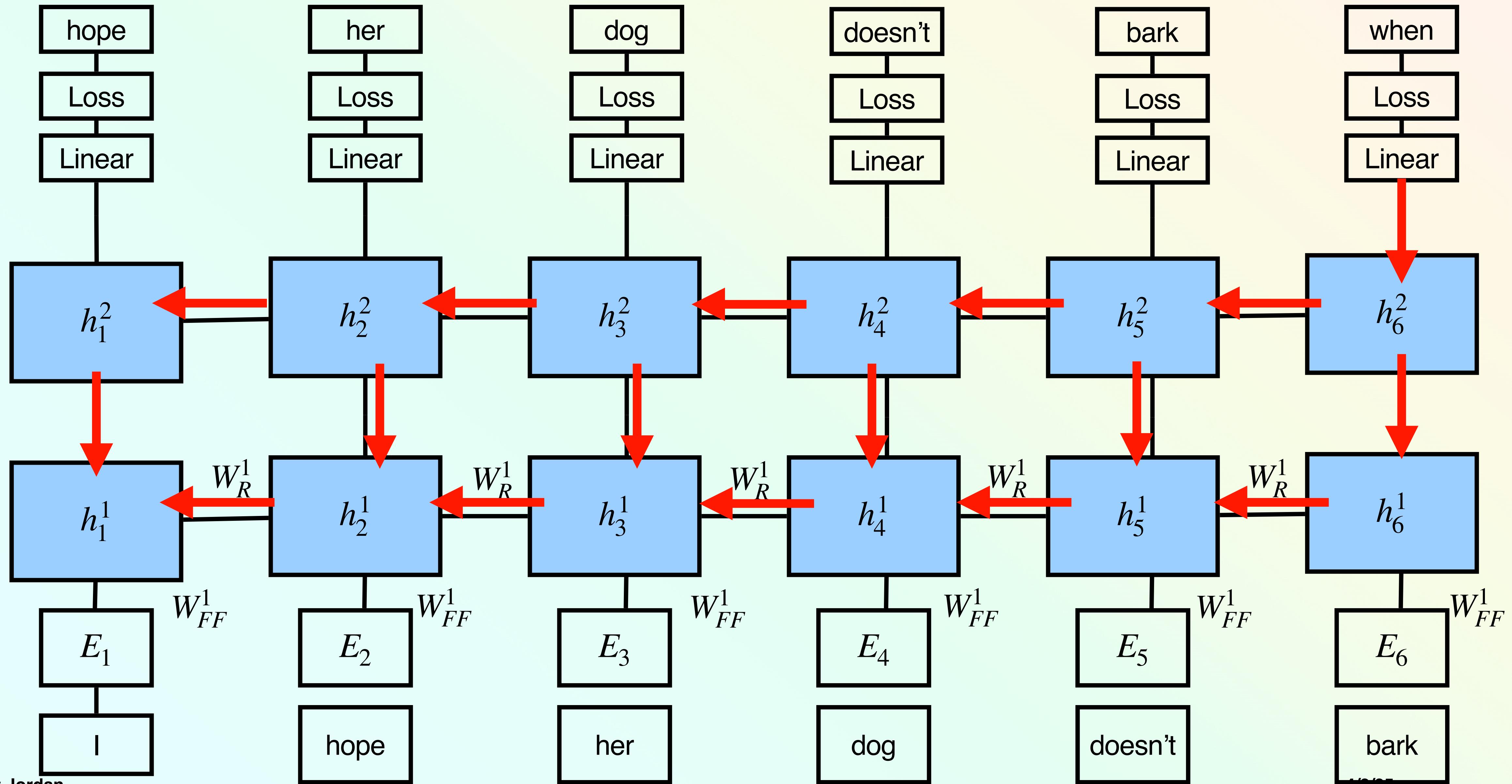
WHAT THEY LEARN

- Learns a context vector
- Context vectors get modified based on new inputs
- Context decays to some baseline representation
 - Slowly forgetting past inputs

RNNs

TRAINING

- Need to identify long-term relationships
- Each time step is like having another layer in an MLP
 - 500 length sequence is a 500 layer MLP
- **Vanishing and Exploding Gradients!**



LONG SHORT-TERM MEMORY

LSTM RNN

- Problem: the derivative for h_1^i can vanish when coming from h_T^i
- Idea: create path through the network that has less shrinking on derivatives over time

Save some value for the hidden unit

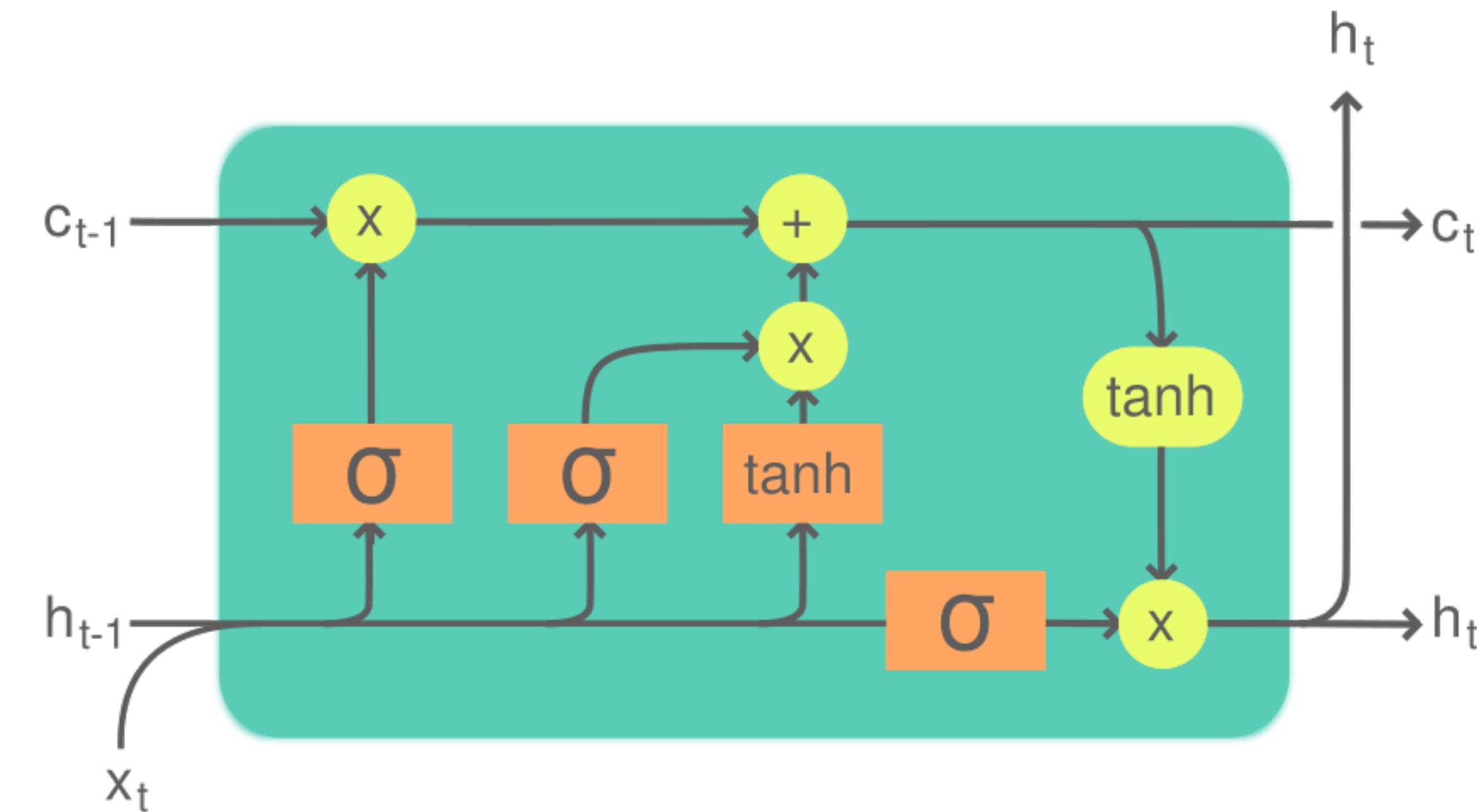
Add to this value with new information

When previous history becomes irrelevant, forget it (set it to saved value to 0)

Conditionally output relevant parts of the hidden information.

LONG SHORT-TERM MEMORY

LSTM RNN



Legend:

Layer



Componentwise



Copy



Concatenate



LONG SHORT-TERM MEMORY

LSTM RNN

c_t – contains the value that has less shrinkage over time

h_t – outputs the hidden state of the recurrent layer

$$c_t = c_{t-1} \odot F_t + I_t \odot \tilde{C}_t$$

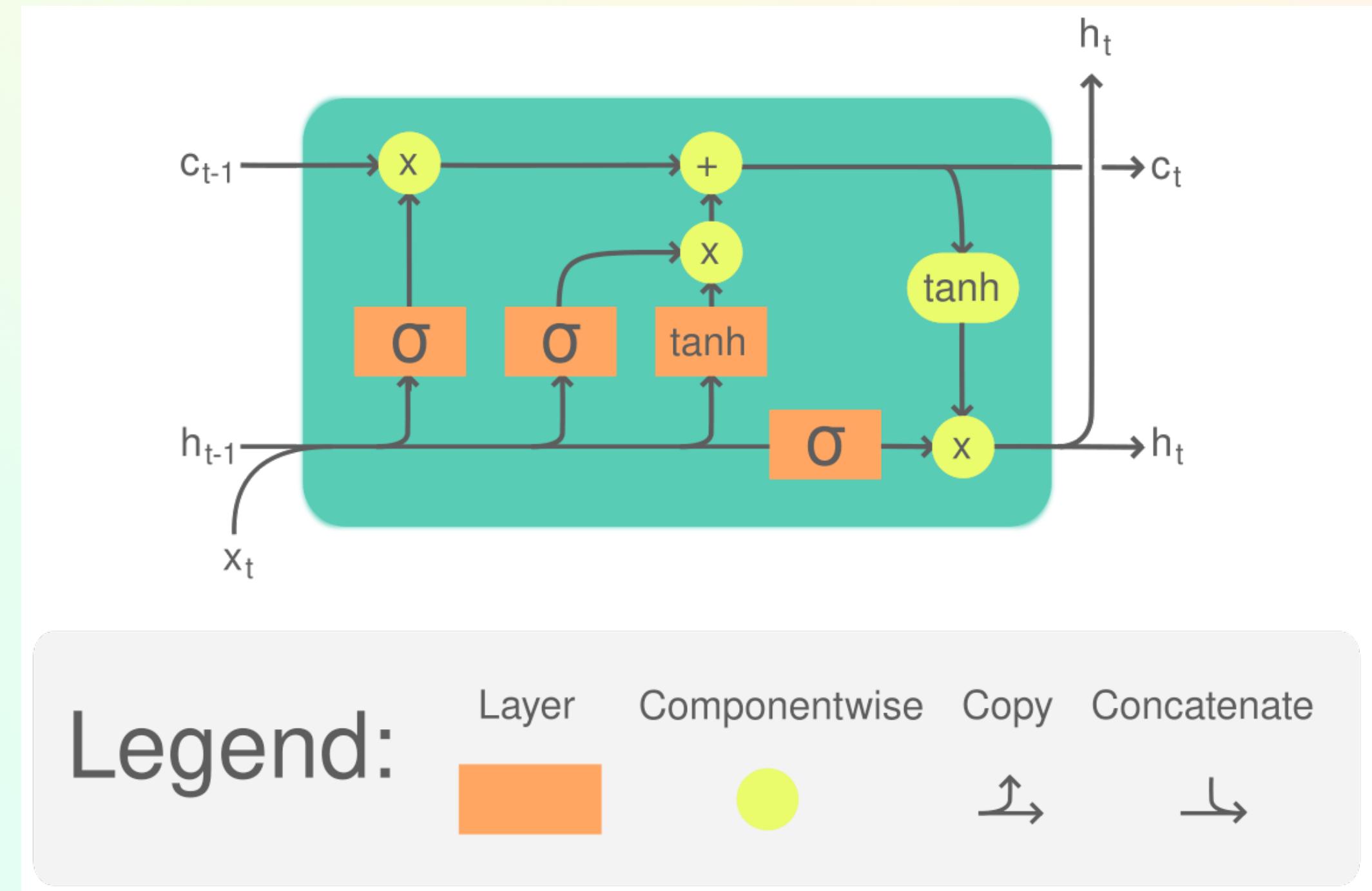
$F_t = \text{sigmoid}([h_{t-1}, x_t] W_F^\top)$ – forget gate

$I_t = \text{sigmoid}([h_{t-1}, x_t] W_I^\top)$ – increment gate

$\tilde{C}_t = \tanh([h_{t-1}, x_t] W_C^\top)$ – increment value

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

$O_t = \text{sigmoid}([h_{t-1}, x_t] W_O^\top)$ – output gate



LONG SHORT-TERM MEMORY

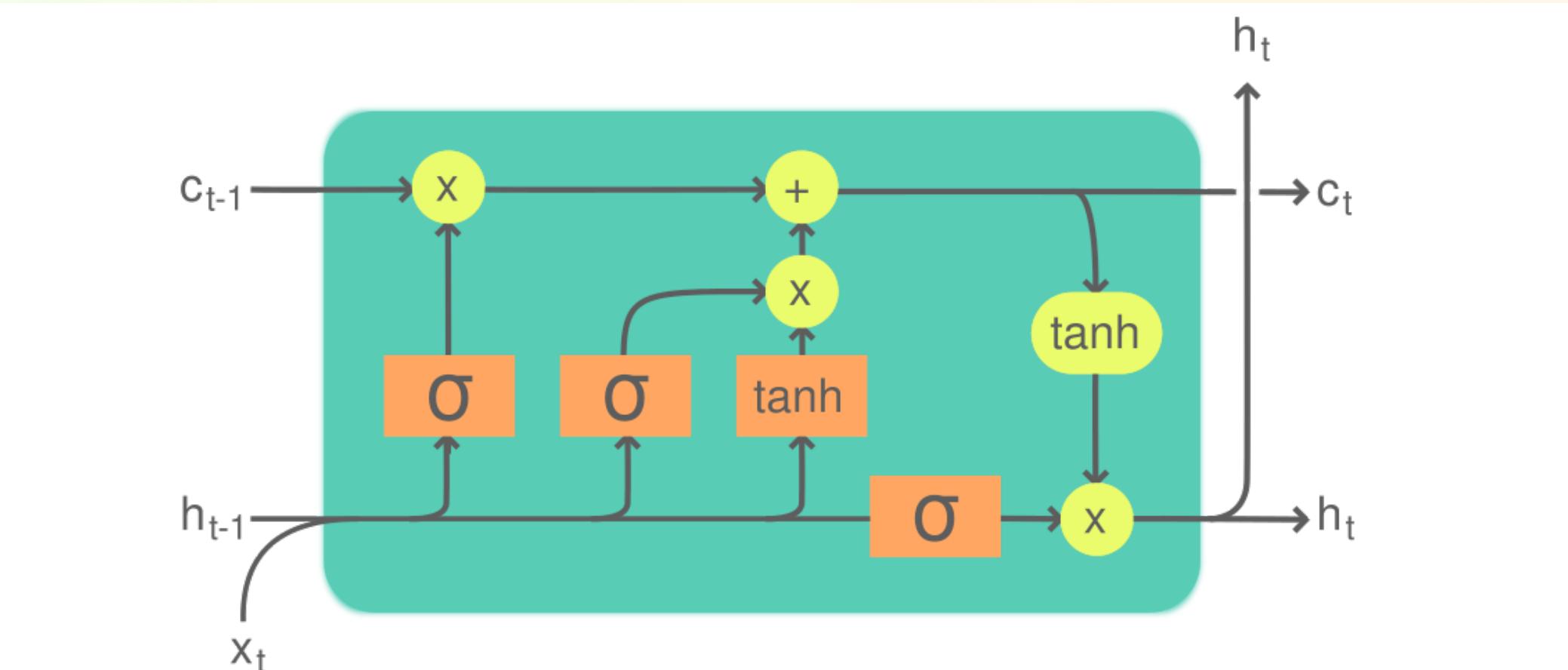
LSTM RNN

$$c_t = c_{t-1} \odot F_t + I_t \odot \tilde{C}_t$$

If $F \approx 1$ the derivative does not shrink much

$F \approx 0$ means prior information was not useful, and the derivative can be small.

LSTMs can learn relationships over longer sequences more easily than RNNs.



Legend:

Layer	Componentwise	Copy	Concatenate

RECURRENT NEURAL NETWORKS

SUMMARY

- Constant run time and memory at inference
- Stores all important information of the past in a single vector
- Optimizing over long sequences is difficult (vanishing/exploding gradients)
- LSTMs help with vanishing gradients

Transformer models can learn longer sequences.

Transformer models are parallelizable

RNN-type layers have to be computed sequentially over the sequence (slow to train)

NEXT CLASS

Reinforcement Learning