

# CS 224S / LINGUIST 285

## Spoken Language Processing

Andrew Maas  
Stanford University  
Spring 2022

### **Lecture 5: Deep Learning Preliminaries**

# Outline for today

---

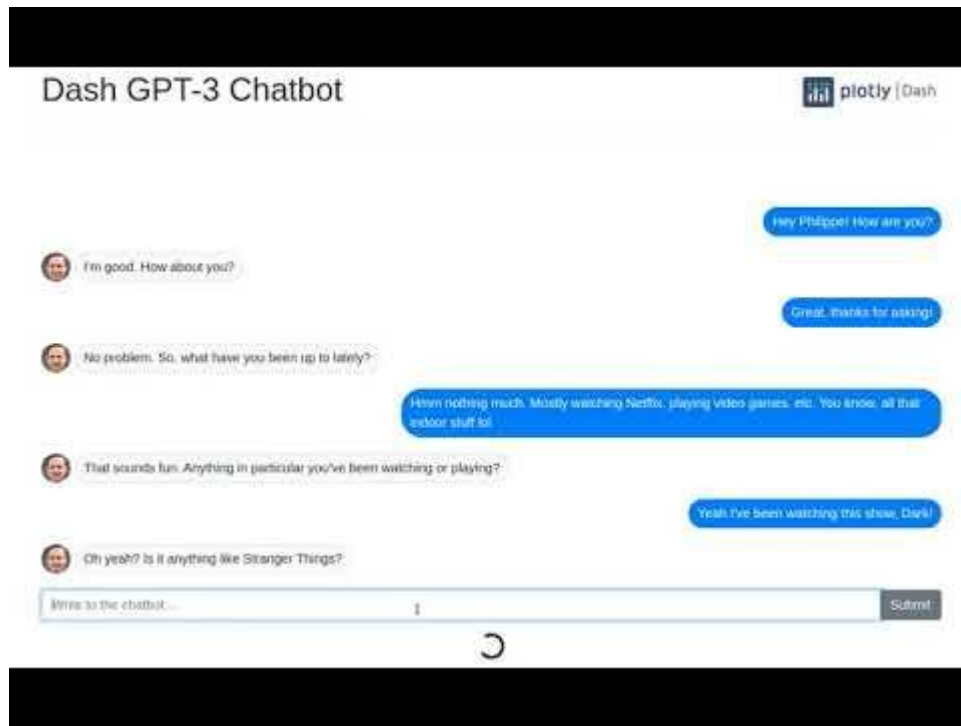
HW1 due now. HW2 posted to course website tonight

- Neural network basics: Hidden layers and gradient computation
- Recurrent NNs
  - Introduction
  - Encoder-decoder
- Attention
  - Introduction
  - Multi-headed attention
- Transformers: stacking attention

# Chatbots in the wild

---

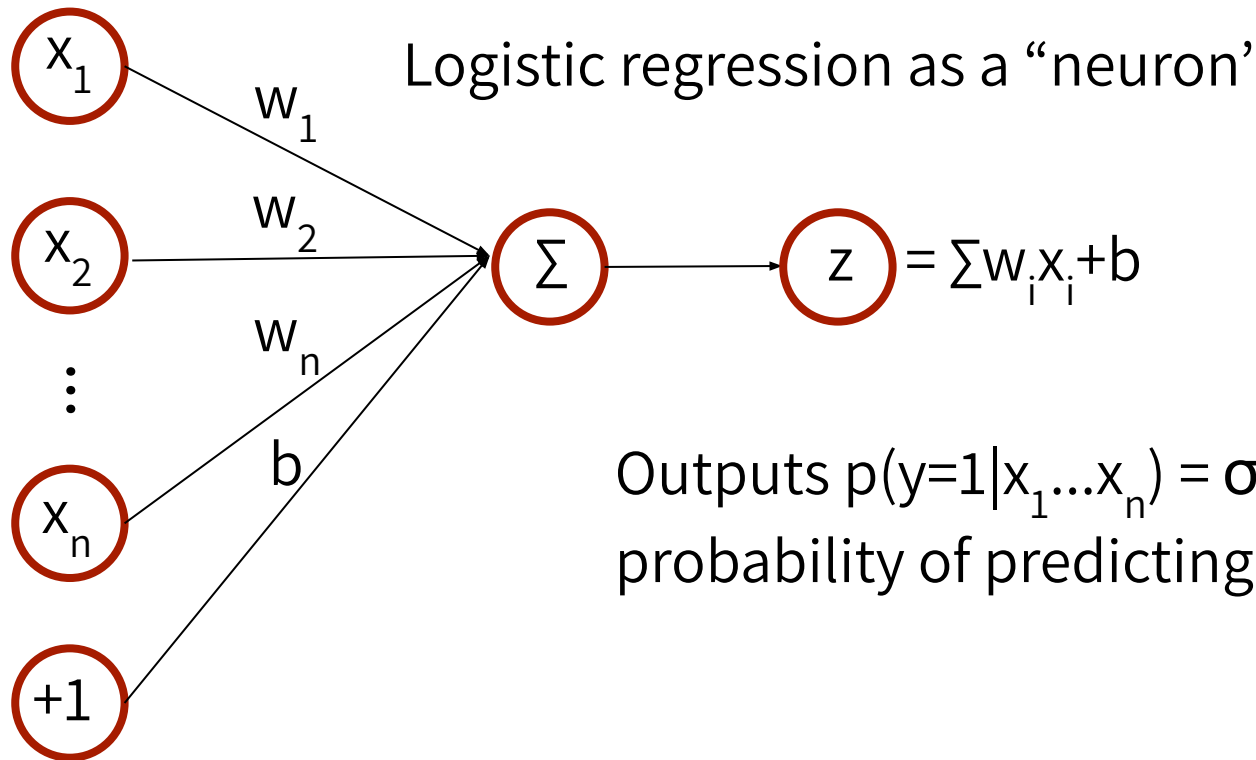
Large encoder decoder transformer models are doing wonders for realistic chatbots.



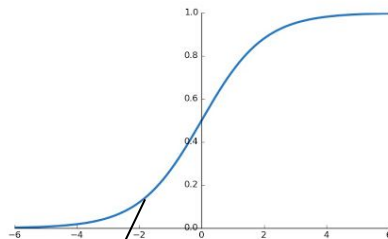
# Neural Network Review

---

# Logistic Regression



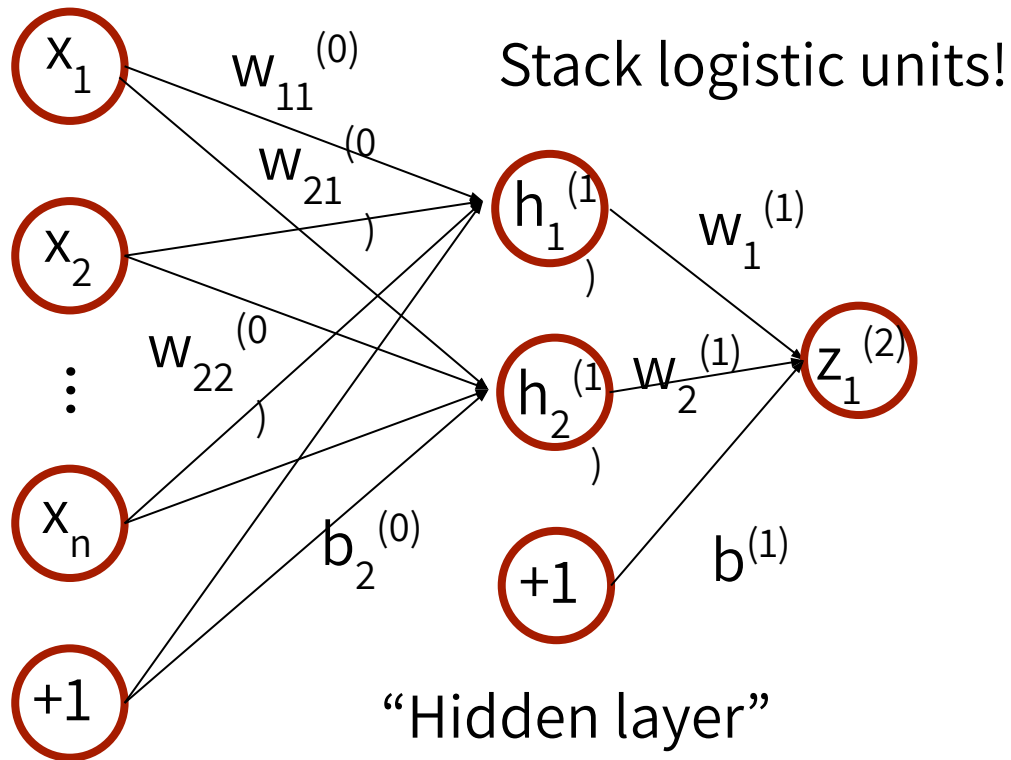
nonlinearity



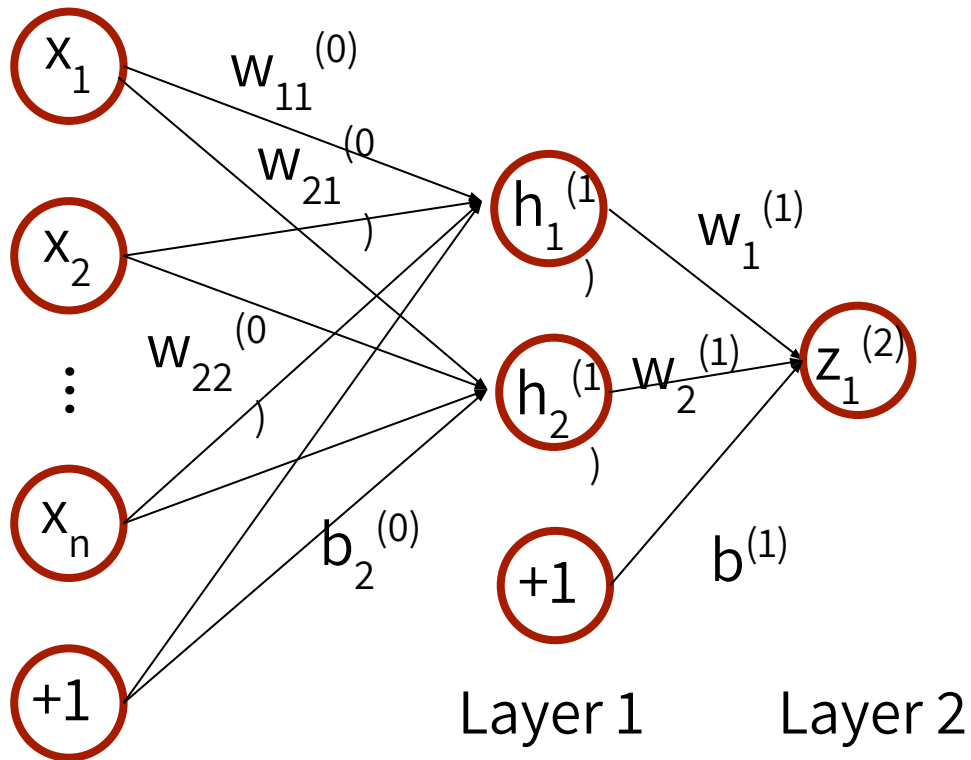
Outputs  $p(y=1|x_1...x_n) = \sigma(z)$ , the probability of predicting class 1!

# Multi-layer Perceptrons

---



# Multi-layer Perceptrons



## Forward Propagation

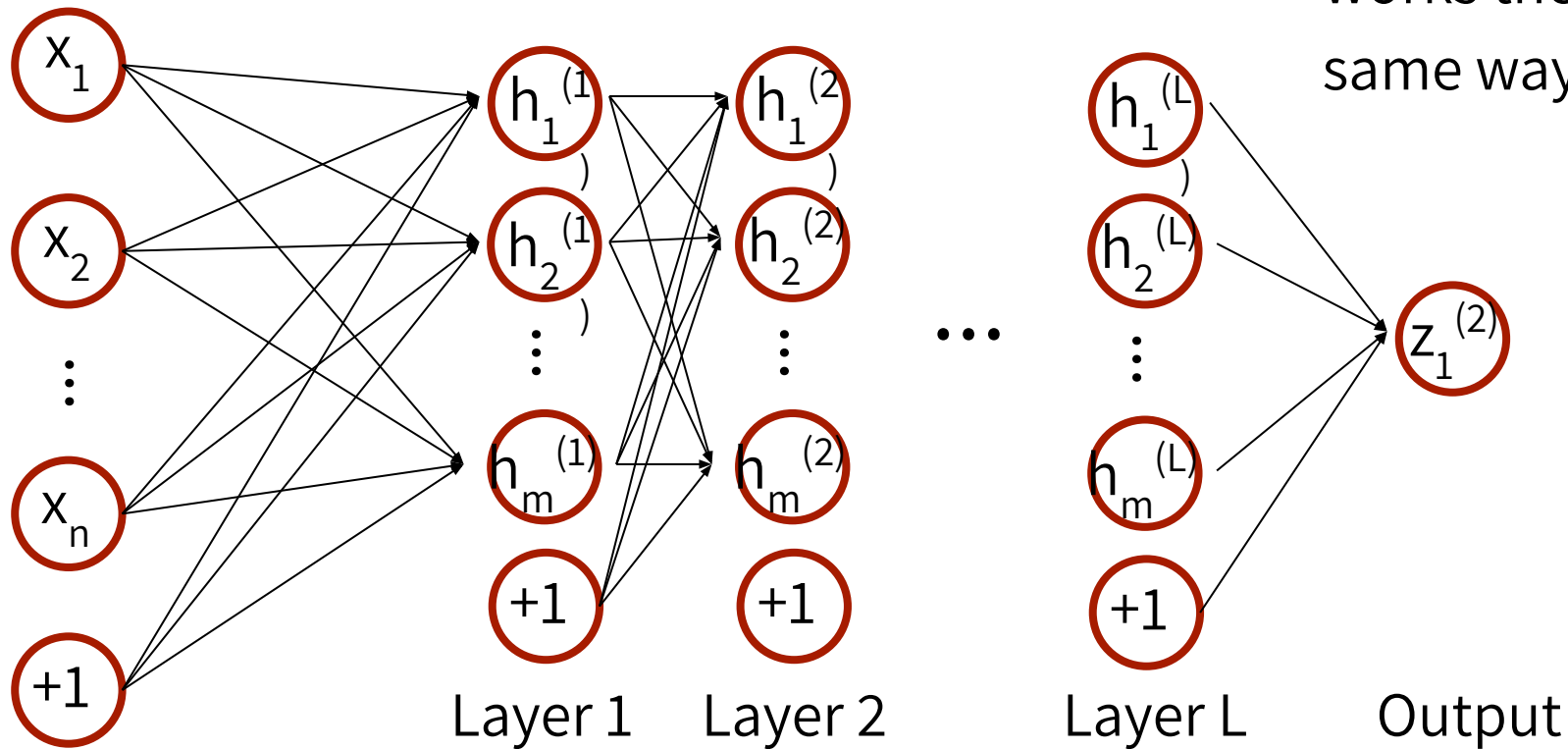
$$z_j^{(l+1)} = \sum w_{ij}^{(l)} a_i^{(l)} + b_j^{(l)}$$

$$h_j^{(l)} = \sigma(z_j^{(l)}) \text{ “activation”}$$

$$\theta = \{ w_{ij}^{(l)}, b_j^{(l)} \text{ for all } i, j, l \}$$

^ “parameters”

# “Deep” Neural Networks





# Objective Function

---

Depends on the task! Examples:

Binary classification

Label:  $y \in \{0,1\}$

Objective:

$$p = \text{sigmoid}(z)$$

$$J_{\theta} = y \log p + (1-y) \log(1-p)$$

Multiclass classification

Label:  $y \in \{1, \dots, K\}$

Objective:

$$p_{1:K} = [p_1, \dots, p_K] = \text{softmax}(z_{1:K})$$

$$J_{\theta} = -\sum y_c^{\text{onehot}} \log p_c$$

Regression

Label:  $y \in \mathbb{R}^d$

Objective:

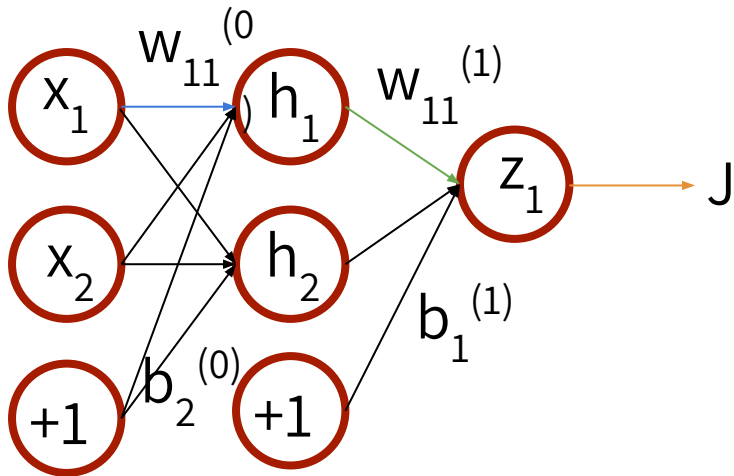
$$J_{\theta} = \text{sum}(\text{sqrt}(z - y)) / d$$

Now we have to optimize!

# Backpropagation

Let's do  $dJ/dw_{11}^{(0)}$  as an example:

$$\frac{dJ}{dw_{11}^{(0)}} = \frac{dJ}{dz_1} \frac{dz_1}{dw_{11}^{(0)}} = \frac{dJ}{dz_1} \left( \frac{dz_1}{dh_1} \frac{dh_1}{dw_{11}^{(0)}} \right)$$



Use chain rule!

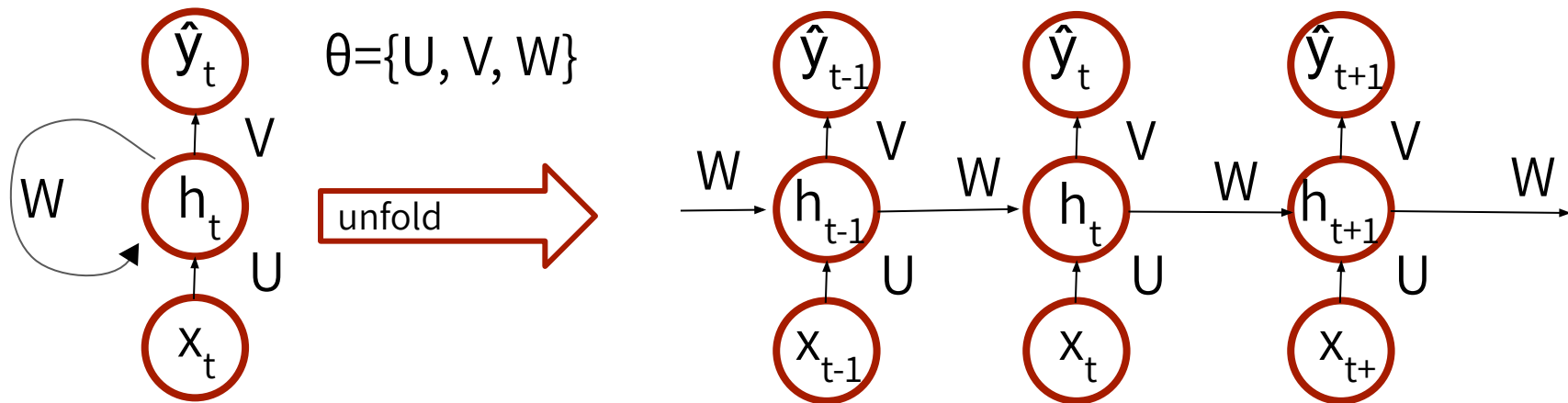
For a fixed objective  $J$  **and** a fixed architecture, you can compute a closed form for every derivative.

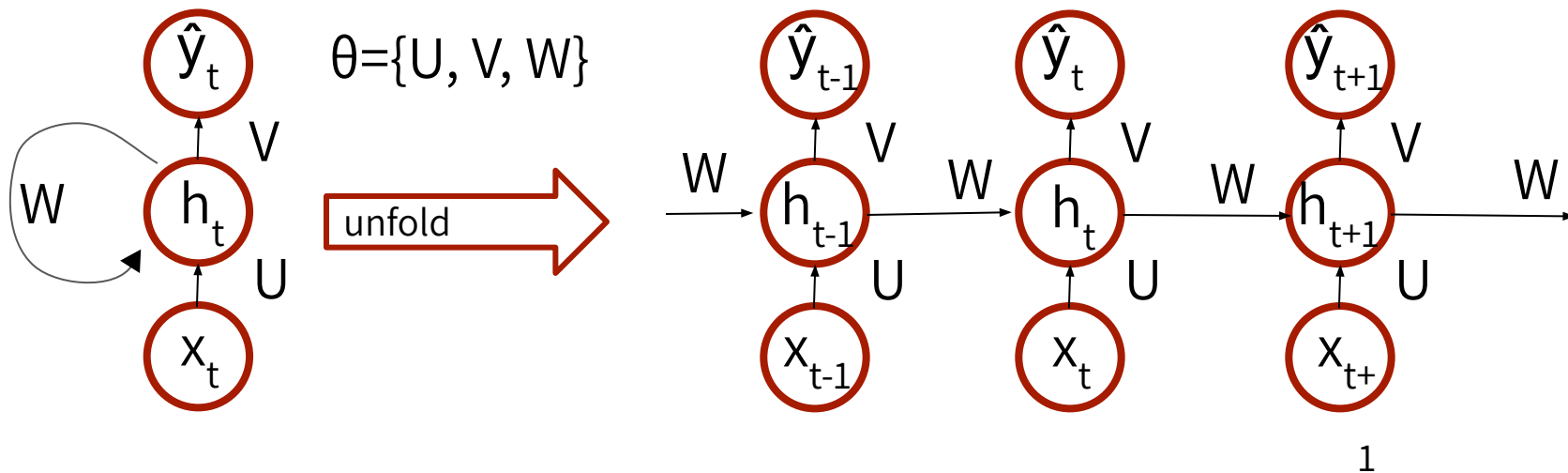
# Recurrent NNs

---

# Recurrent NNs

- Input is a sequence of tokens  $(x_1, x_2, \dots, x_T)$
- Output is a sequence of tokens  $(y_1, y_2, \dots, y_T)$
- Goal is map  $x_t$  to a “hidden state”  $h_t$  (a real-valued vector)
- Think of  $h_t$  is a nonlinear summary of  $(x_1, \dots, x_t)$
- Use  $h_{t-1}$  and  $x_t$  to predict  $y_t$





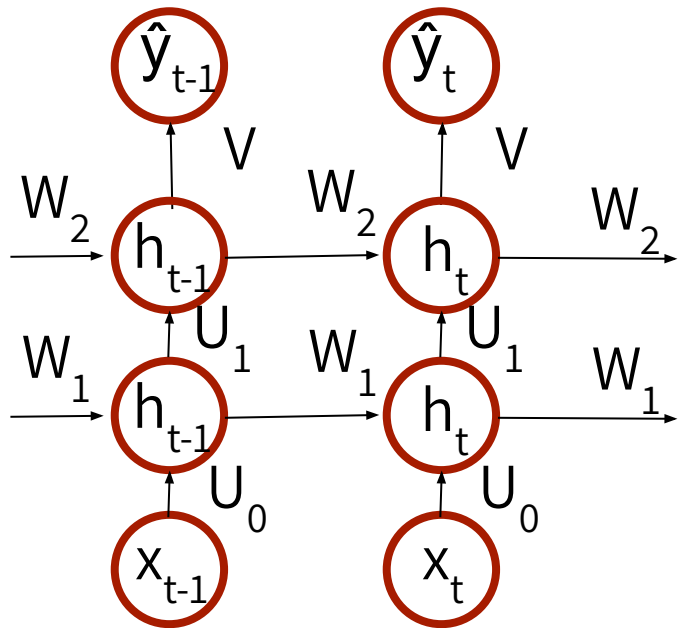
- U, V, W are shared over all timesteps (same ones!)
- The model does a forward pass for every single timestep:  $f_{\theta}(x_t, h_{t-1})$
- Objective (assume x and y are discrete tokens)

$$J(x, y, \theta) = - \sum_t \log p(y_t | x_1, \dots, x_t) = - \sum_t \text{CrossEntropy}(y_t, f_{\theta}(x_t, h_{t-1}))$$

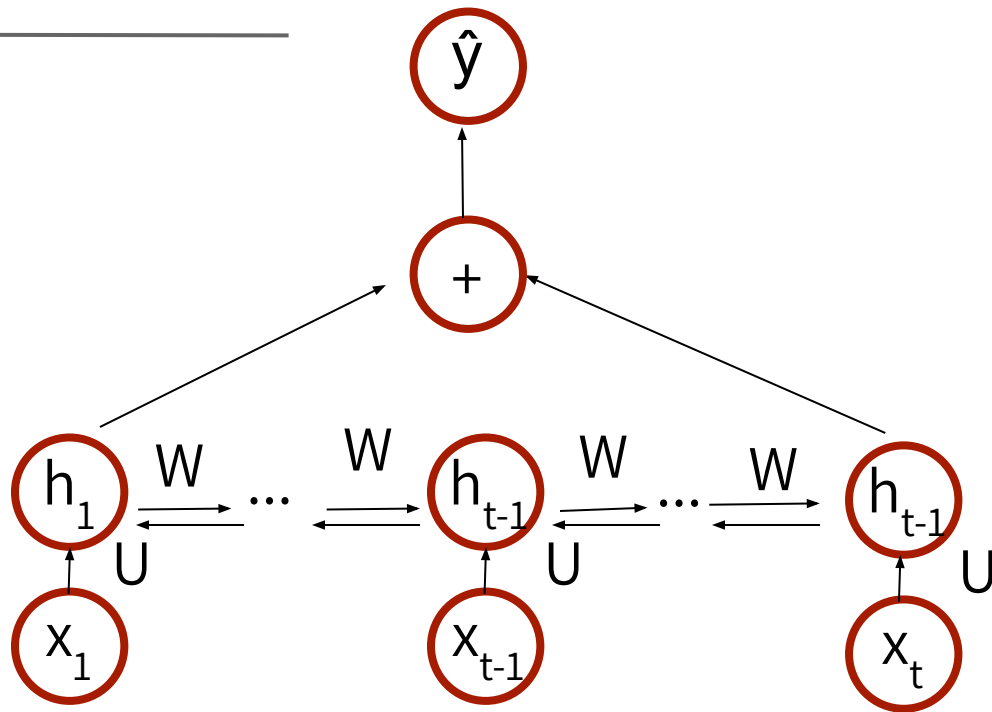
- Backpropagation through time (BPTT) [Rumelhart et. al. 1986, Werbos 1990]

# Improvements to RNN

---

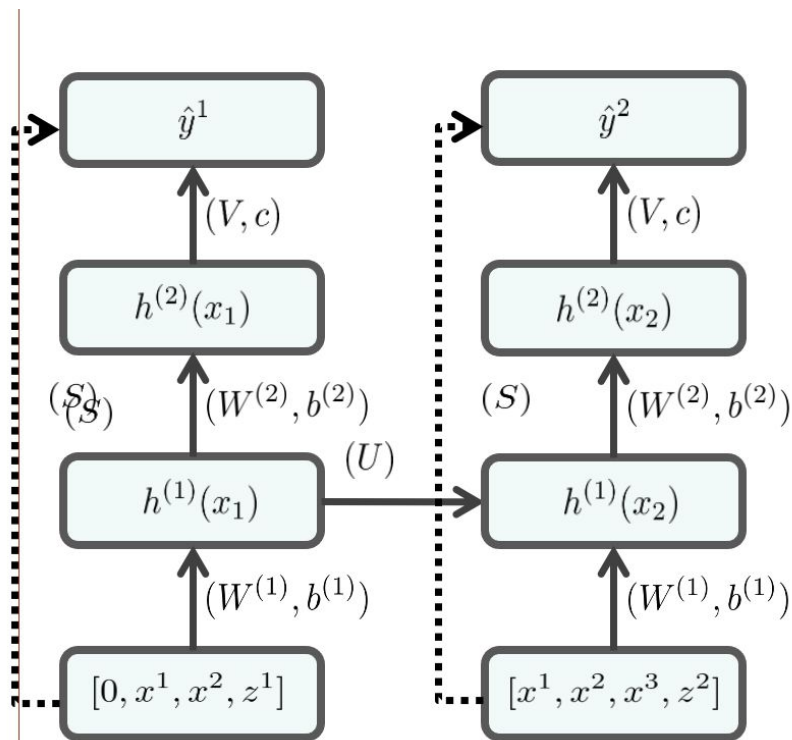


Stacked RNNs



Bi-directional RNNs

# Recurrent networks for speech denoising



**Output Layer**

$$\hat{y}^2 = Vh^{(2)}(x_2) + Sx_2 + c$$

**Hidden Layer**

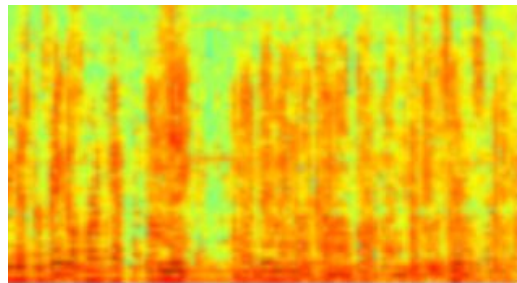
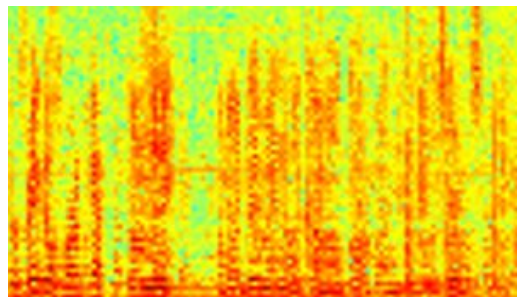
$$h^{(2)}(x_2) = \sigma(W^{(2)}h^{(1)}(x_2) + b^{(2)})$$

**Hidden Layer**

$$h^{(1)}(x_2) = \sigma(W^{(1)}x_2 + b^{(1)} + Uh^{(1)}(x_1))$$

Noisy Input  $X$

Noise Estimate  $Z$



# Encoder-Decoder Models

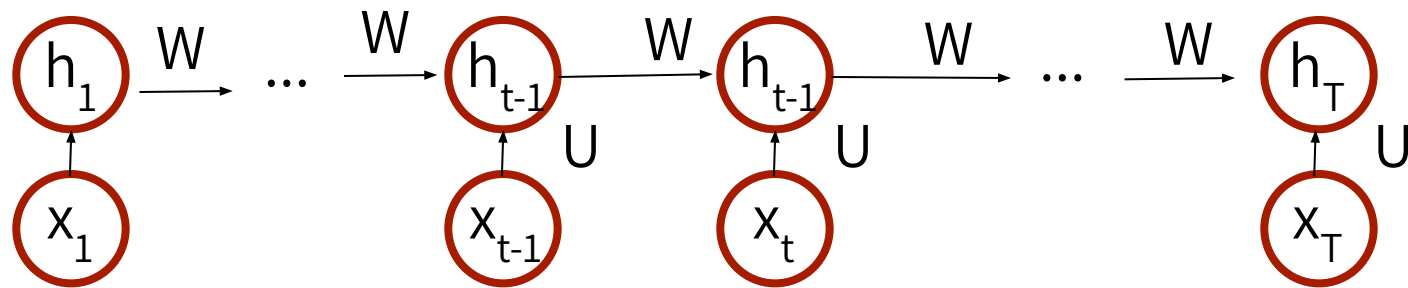
---



# RNN Encoder

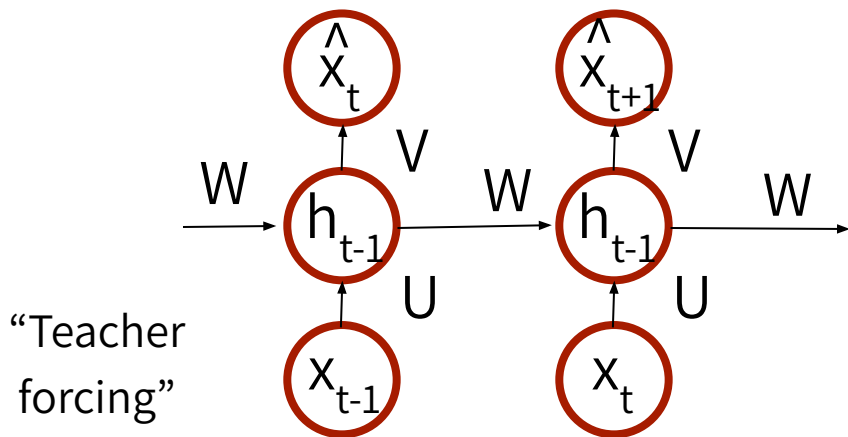
---

Input is a sequence of tokens  $(x_1, x_2, \dots, x_T)$ . Goal is to summarize the sequence into a single vector.

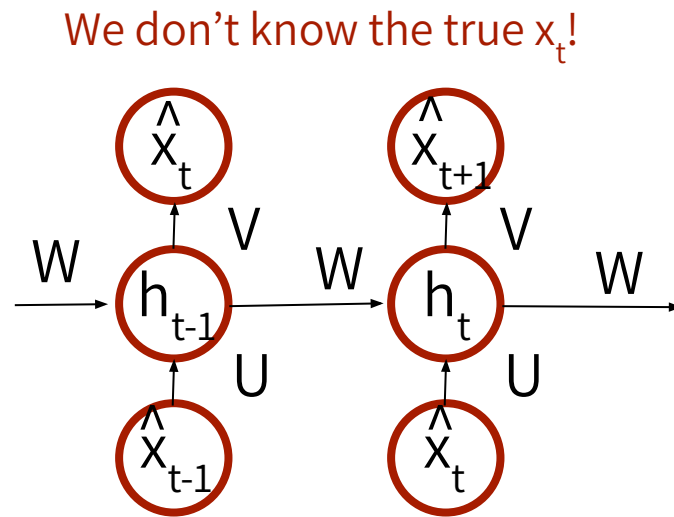


# RNN Decoder

Input is a sequence of tokens  $(x_1, x_2, \dots, x_T)$ , no output sequence. The model should learn  $p(x_t | x_1, \dots, x_{t-1})$ .



Training

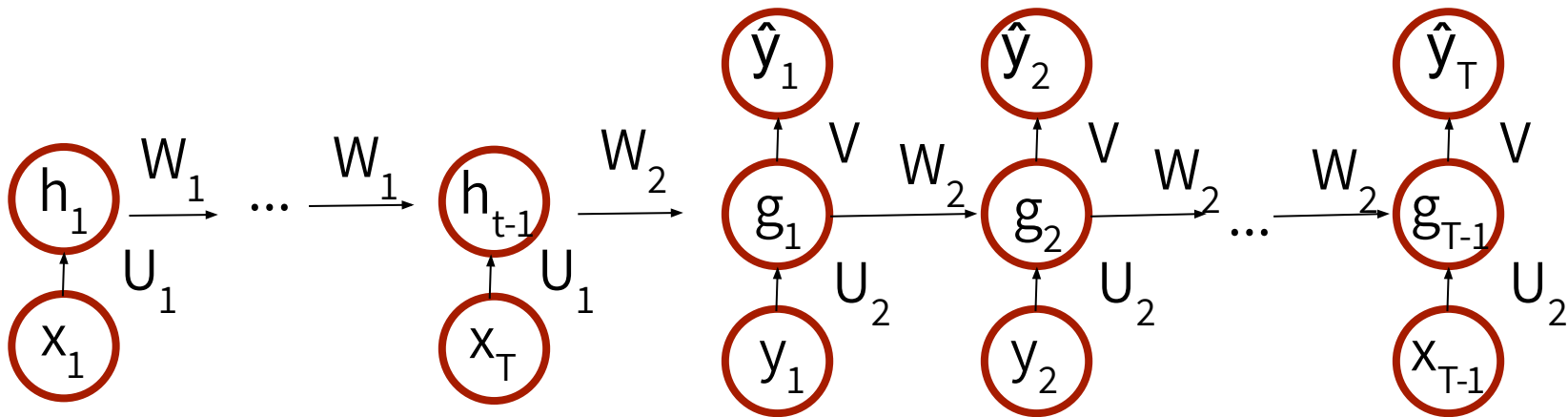


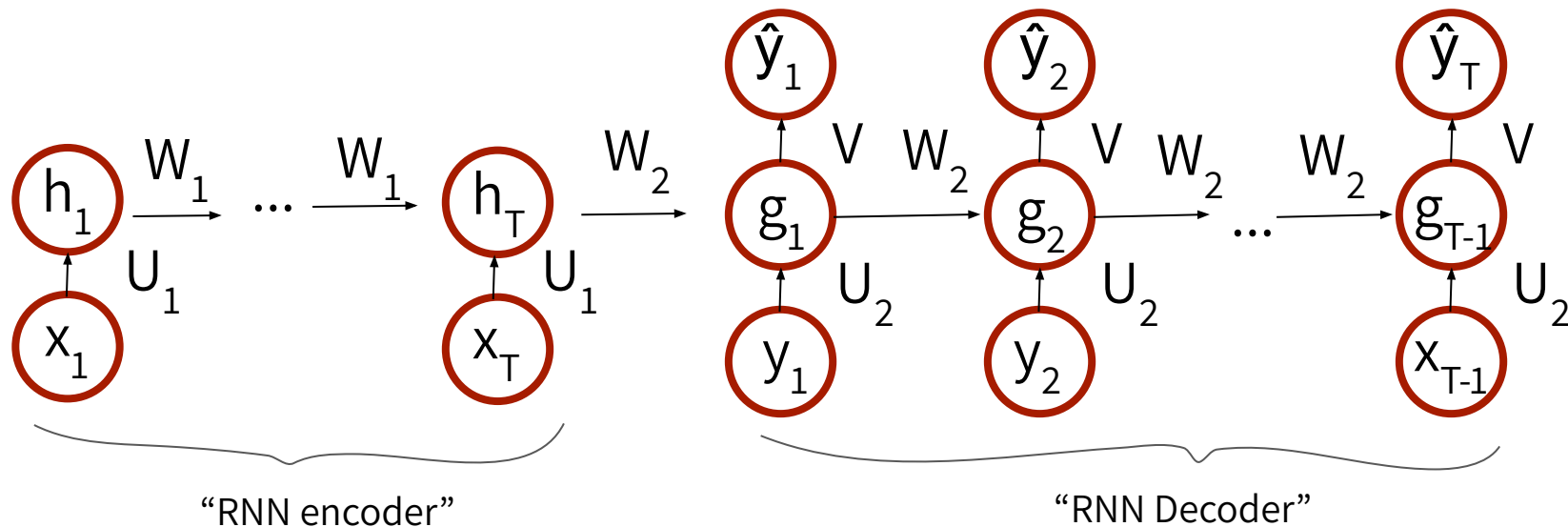
Generation

# Seq2Seq [Cho et . al. 2014, Sustekever et. al. 2014]

---

- All that work and we have our first encoder-decoder model!
- Given an input and output sequence  $(x_1, x_2, \dots, x_t), (y_1, y_2, \dots, y_t)$ , the model should capture  $p(y_t | y_1, \dots, y_{t-1}, x_1, \dots, x_T)$ . It gets to see all of  $x$ !





- Two different RNNs glue'd together (separate parameters)
- One of them encodes  $(x_1, \dots, x_T)$  into a summary vector,  $h_T$
- The other one uses  $h_T$  to initialize a language model
- Train this just like an RNN language model ( $x$  = speaker 1,  $y$  = speaker 2)

# Limitations of RNNs

---

- If you have a really long sequence (e.g.  $T=1000$ ), hard to believe  $h_{854}$  will remember  $x_2$
- $h_t$  captures “local” info since it has to predict  $y_t$  (little incentive to remember  $h_{t-100}$ ).

But long range dependencies are important.

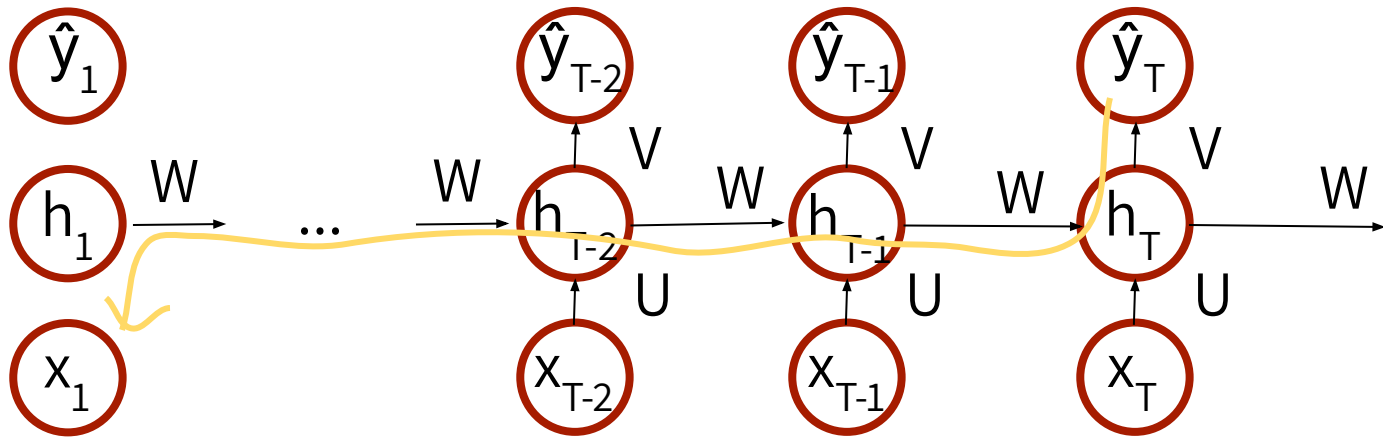
**Example:** ... the flights the airline was cancelling were full.

What flights??? What airline???

If you have building a chatbot, you might need to remember things from long ago.

# Vanishing Gradients

---



- Suppose  $T \rightarrow \infty$ . Say we want to calculate  $dy_t / dx_1$
- What happens if  $dh_t / dh_{t-1} < 1$  for all  $t$ ?
- Impact of  $x_{t-s}$  on  $y_t$  decreases as  $s$  increases.

# Attention

---

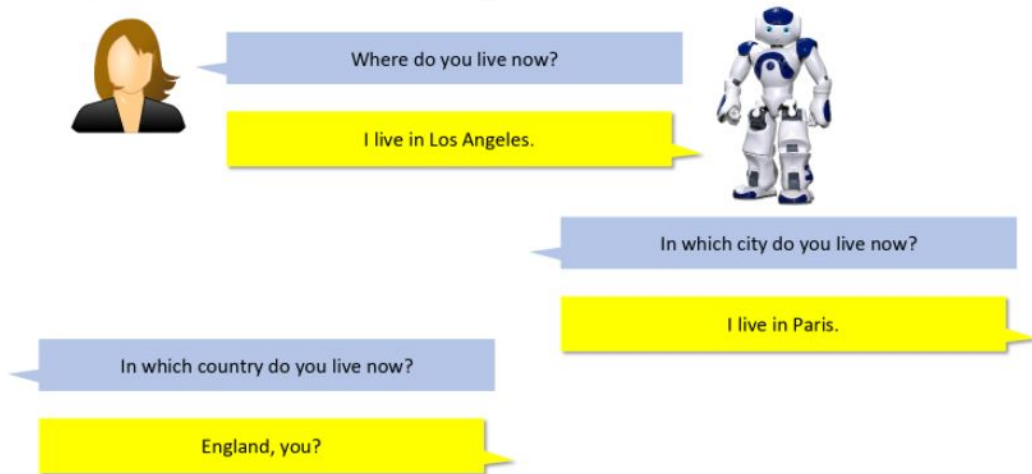
# Attention [Bahdanau et. al. 2014]

---

**Holy grail:** capturing long term dependencies.

- Vanishing gradient problem & local dependency of RNNs.
- Attention gets at this more directly (and simply).

## Speaker Consistency



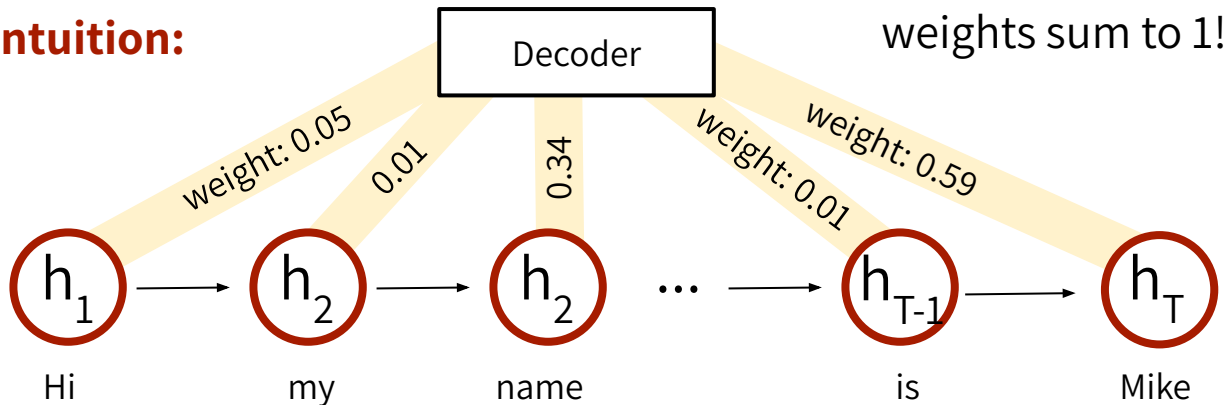


# Attention [Bahdanau et. al. 2014]

**Holy grail:** capturing long term dependencies.

- Vanishing gradient problem & local dependency of RNNs.
- Attention gets at this more directly (and simply).

**Intuition:**



## Attention mechanism

$q \in \mathbb{R}^d$ : query ;  $k_1, \dots, k_T \in \mathbb{R}^d$ : keys ;  $v_1, \dots, v_T \in \mathbb{R}^d$ : values

## Attention mechanism

$q \in \mathbb{R}^d$ : query ;  $k_1, \dots, k_T \in \mathbb{R}^d$ : keys ;  $v_1, \dots, v_T \in \mathbb{R}^d$ : values

**Step 1:** define a similarity function  $\text{sim}(q, k_t)$ .

$$\text{sim}(q, k_t) = w_2 \text{relu}(w_1[q; k_t] + b_1) + b_2 \quad [\text{Bahdanau et. al. 2014}] \quad (\text{MLP})$$

$$\text{sim}(q, k_t) = q^T W k_t \quad [\text{Luong et. al. 2015}]$$

(Bilinear)

$$\text{sim}(q, k_t) = q^T k_t \quad [\text{Luong et. al. 2015}]$$

(dot-pdt)

$$\text{sim}(q, k_t) = q^T k_t / \text{sqrt}(d) \quad [\text{Luong et. al. 2015}] \quad (\text{scaled dot-pdt})$$



## Attention mechanism

$\mathbf{q} \in \mathbb{R}^d$ : query ;  $\mathbf{k}_1, \dots, \mathbf{k}_T \in \mathbb{R}^d$ : keys ;  $\mathbf{v}_1, \dots, \mathbf{v}_T \in \mathbb{R}^d$ : values

**Step 1:** define a similarity function  $\text{sim}(\mathbf{q}, \mathbf{k}_t)$ .

**Step 2:** compute attention weights  $a_t$ .

$$a_t = \frac{\exp\{\text{sim}(\mathbf{q}, \mathbf{k}_t)\}}{\sum_{s=1}^T \exp\{\text{sim}(\mathbf{q}, \mathbf{k}_s)\}}$$

Note  $a_t \in [0, 1]$  for all  $t$   
Also  $\sum_t a_t = 1$

## Attention mechanism

$q \in \mathbb{R}^d$ : query ;  $k_1, \dots, k_T \in \mathbb{R}^d$ : keys ;  $v_1, \dots, v_T \in \mathbb{R}^d$ : values

**Step 1:** define a similarity function  $\text{sim}(q, k_t)$ .

**Step 2:** compute attention weights  $a_t$ .

**Step 3:** attend to values vectors.

$$c = \sum_{t=1}^T a_t v_t \quad \text{weighted linear combo of values!}$$

# Multi-headed Attention [Vaswani et. al. 2017]

---

What if I want to pay attention to different things at the same time!?

This is my big red dog, Clifford.

Content-based

This is my big red dog, Clifford.

Description-based

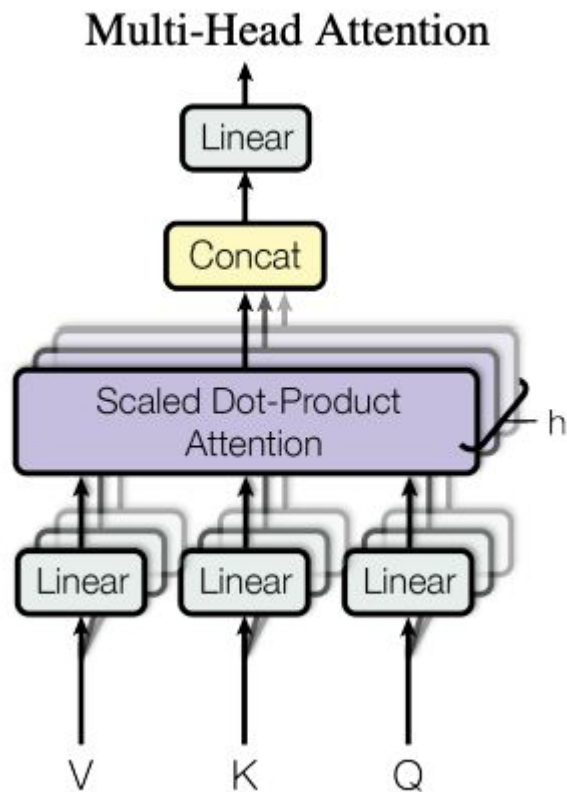
This is my big red dog, Clifford.

Reference-based

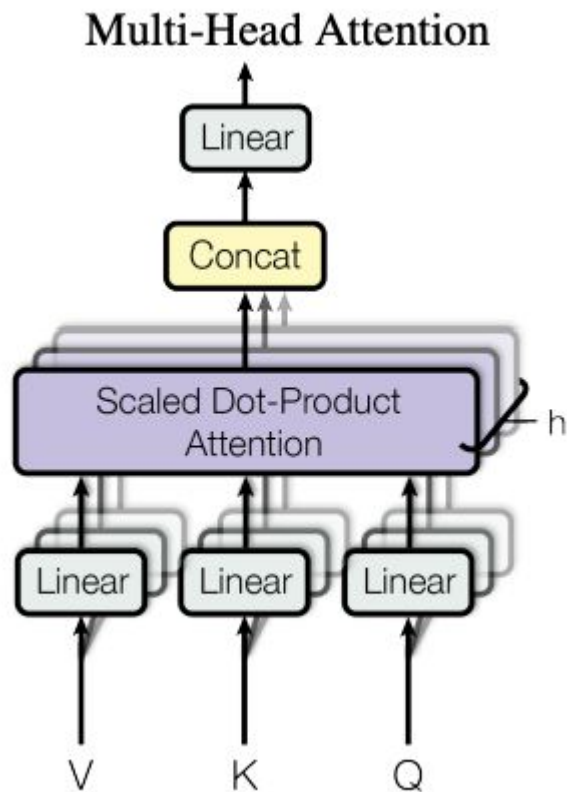
This is my big red dog, Clifford.

What's useful depends on the task. How do I pick what to do?

Idea [Vaswani et. al. 2017]: Don't pick. Pay attention as if you had “multiple heads”.



Idea [Vaswani et. al. 2017]: Don't pick. Pay attention as if you had “multiple heads”.



Pick  $H$  heads.

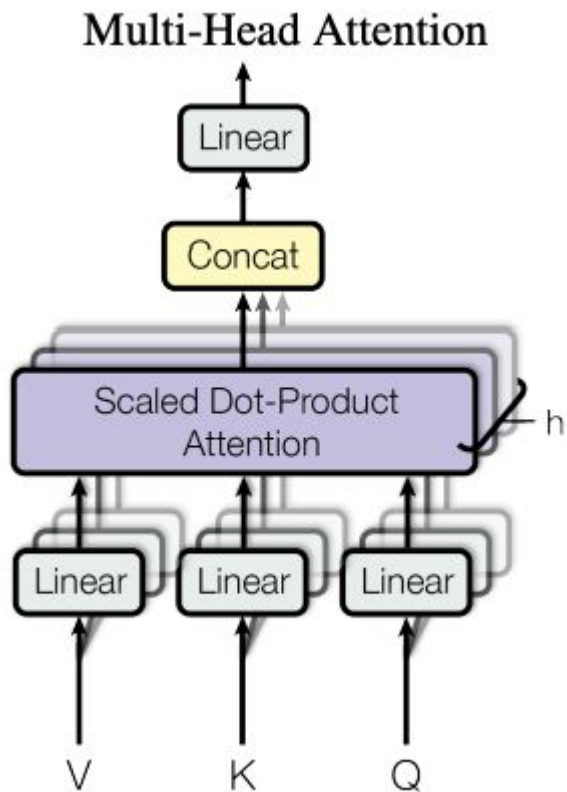
For  $h$  in range( $H$ ):

$$\begin{aligned} q^{(h)} &= \text{MLP}_h(q) \\ k_t^{(h)} &= \text{MLP}_h(k_t) \text{ for all } t \\ v_t^{(h)} &= \text{MLP}_h(v_t) \text{ for all } t \end{aligned}$$





Idea [Vaswani et. al. 2017]: Don't pick. Pay attention as if you had “multiple heads”.



Pick H heads.

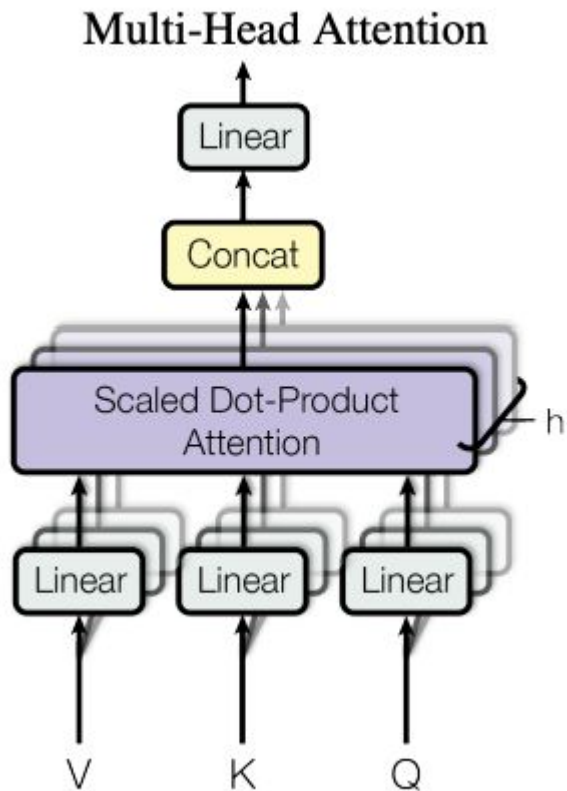
For  $h$  in range(H):

$$q^{(h)} = \text{MLP}_h(q)$$
$$k_t^{(h)} = \text{MLP}_h(k_t) \text{ for all } t$$
$$v_t^{(h)} = \text{MLP}_h(v_t) \text{ for all } t$$

$$a_t^{(h)} = \frac{\exp\{q^{(h)\top} k_t^{(h)} / \sqrt{d}\}}{\sum_{s=1}^T \exp\{q^{(h)\top} k_s^{(h)} / \sqrt{d}\}}$$

$$c^{(h)} = \sum_{t=1}^T a_t^{(h)} v_t^{(h)}$$

Idea [Vaswani et. al. 2017]: Don't pick. Pay attention as if you had “multiple heads”.



Pick H heads.

For  $h$  in range(H):

$$q^{(h)} = \text{MLP}_h(q)$$
$$k_t^{(h)} = \text{MLP}_h(k_t) \text{ for all } t$$
$$v_t^{(h)} = \text{MLP}_h(v_t) \text{ for all } t$$

$$a_t^{(h)} = \frac{\exp\{q^{(h)T} k_t^{(h)} / \sqrt{d}\}}{\sum_{s=1}^T \exp\{q^{(h)T} k_s^{(h)} / \sqrt{d}\}}$$

$$c^{(h)} = \sum_{t=1}^T a_t^{(h)} v_t^{(h)}$$

$$c_{\text{all}} = \text{concat}(c^{(1)}, \dots, c^{(H)})$$

$c = \text{linear}(c_{\text{all}})$  # project to smaller dimension

# Self Attention [Vaswani et. al. 2017]

---

Simple idea: query, keys and values in attention are the same.

Take some sequence  $(x_1, x_2, \dots, x_T)$ . For every  $t$ , build:

$$q_t = \text{MLP}(x_t) \text{ , } k_t = \text{MLP}(x_t) \text{ , } v_t = \text{MLP}(x_t)$$

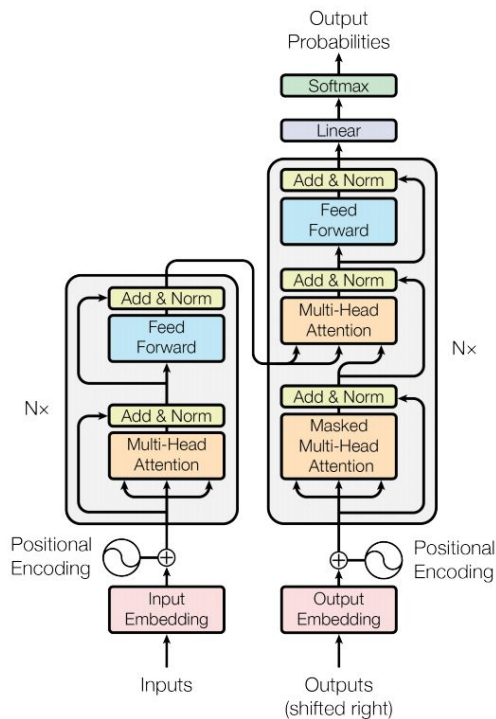
Then we can extract a sequence:

$$c_t = \sum_{t=1}^T a_t v_t \quad \longrightarrow \quad (c_1, c_2, \dots, c_T)$$

# Transformers (SOTA)

---

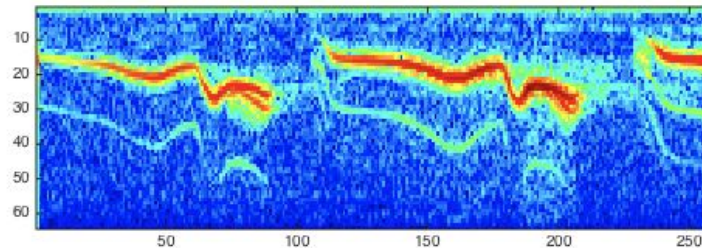
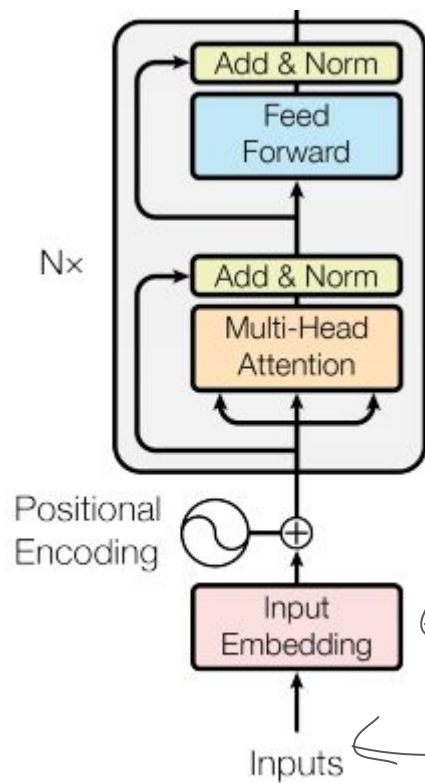
# Transformer: Self-attention Nets



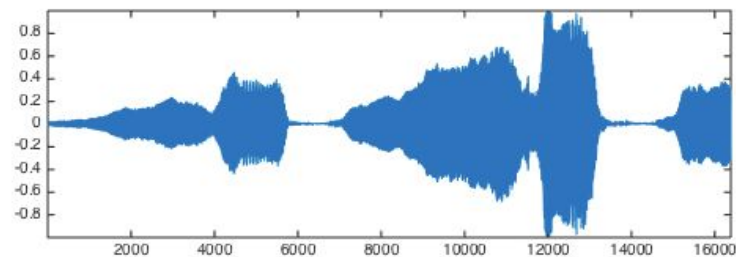
- A transformer layer is composed of an encoder and a decoder.
- Both use the same building blocks.

Similar to RNNs, here...

- Encoder sees  $(x_1, x_2, \dots, x_T)$  and outputs a hidden sequence  $(h_1, h_2, \dots, h_T)$ .
- Decoder sees  $(x_1, x_2, \dots, x_{T-1})$  and outputs predictions for  $(x_2, x_3, \dots, x_T)$ .

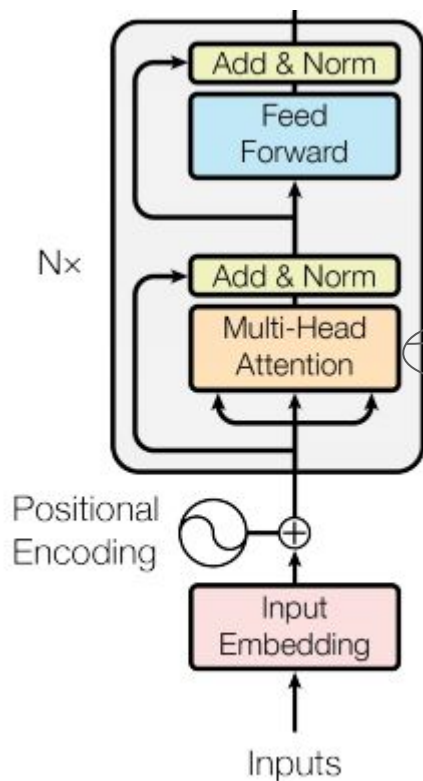


$(e_1, e_2, \dots, e_T)$  e.g. mel spectrogram



$(x_1, x_2, \dots, x_T)$  e.g. a waveform

Transformer encoder

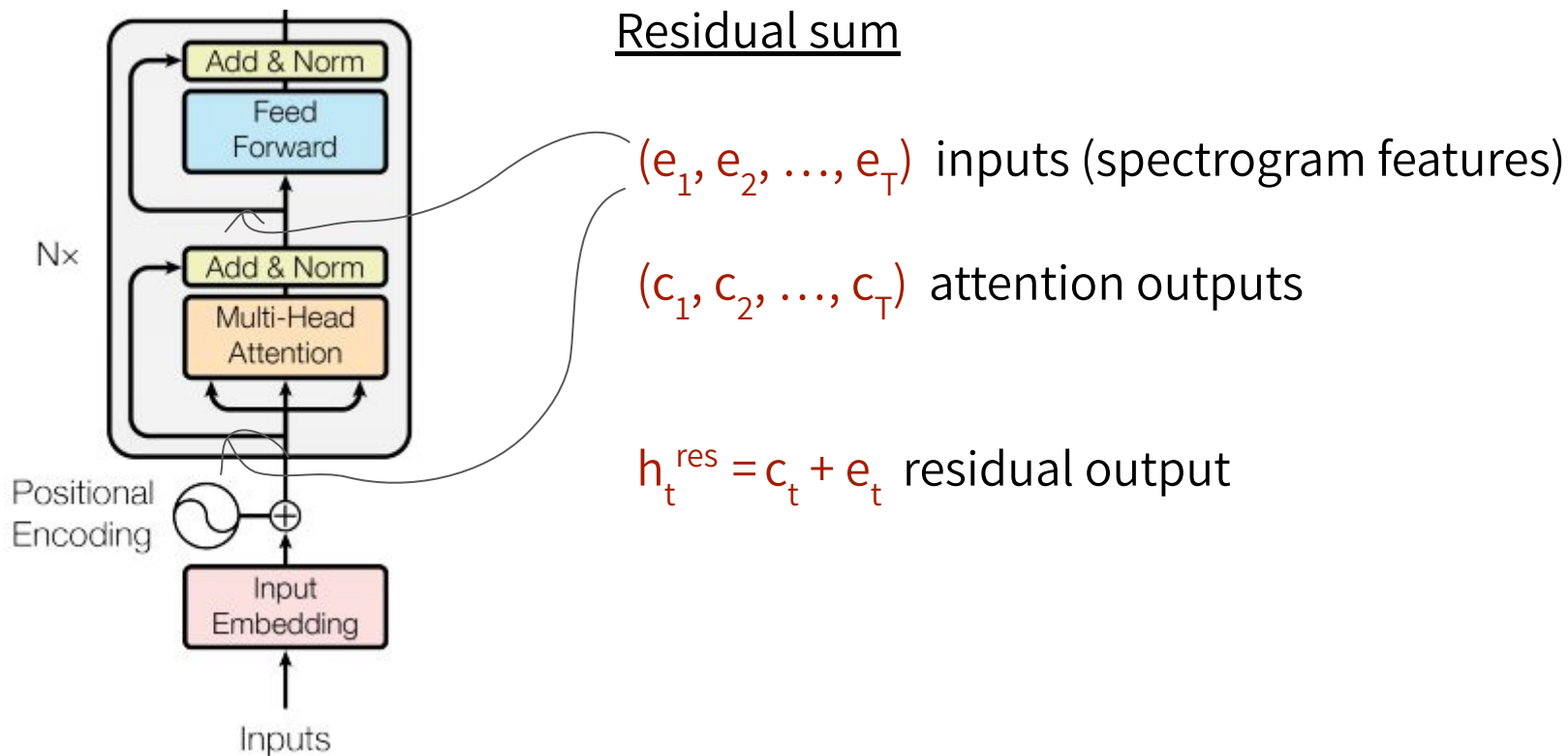


We've seen this! We do self-attention with  $H$  heads on the lookup embeddings.

Inputs  $(e_1, e_2, \dots, e_T)$ , each  $e_t$  is now a vector!

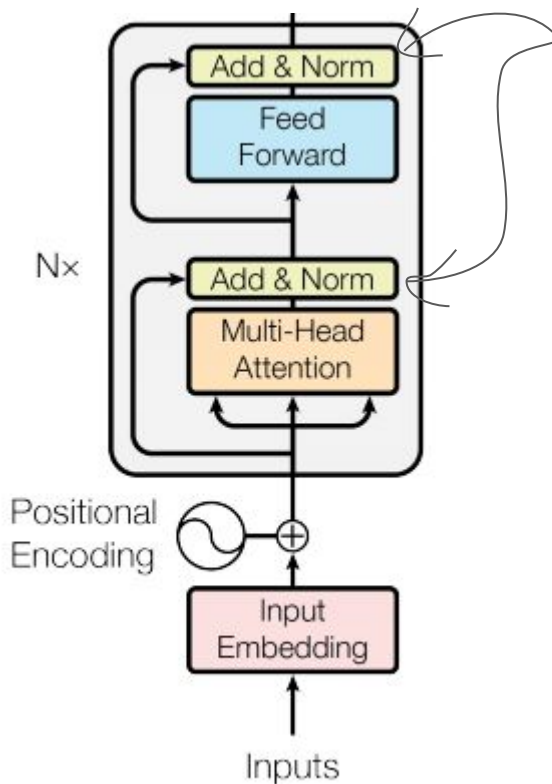
Outputs  $(c_1, c_2, \dots, c_T)$  each  $c_t \in \mathbb{R}^d$

Transformer encoder



Transformer encoder





## Layer normalization

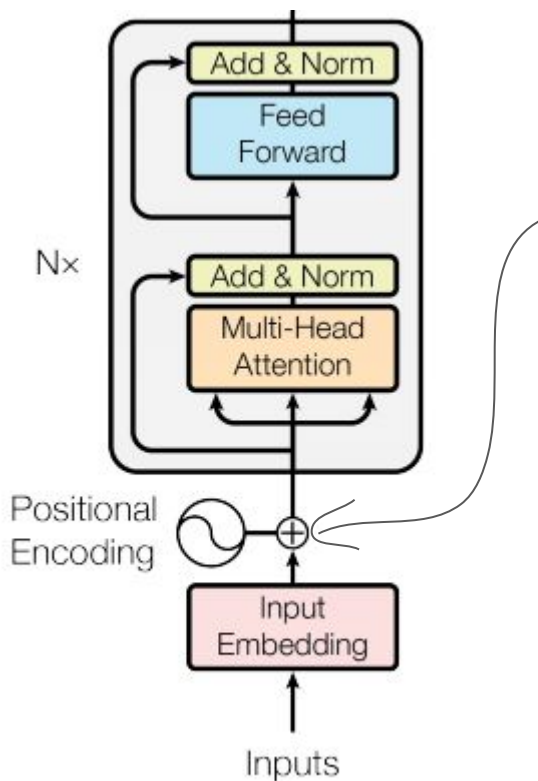
$$h^{\text{norm}} = \frac{h^{\text{res}} - E[h^{\text{res}}]}{\sqrt{\text{Var}[h^{\text{res}}] + \epsilon}} * \gamma + \square$$

Note  $\{\gamma, \square\} \subseteq \theta$  e.g. learnable parameters.

$$h^{\text{res}} = (h_1^{\text{res}}, h_2^{\text{res}}, \dots, h_T^{\text{res}}).$$

- The mean and variance are over the sequence of size T.
- Not like batch norm (which is over a batch of examples). This is only on 1 example.

Transformer encoder



- Unlike RNNs, transformers have no order!
- But speech is left-to-right so it might be useful to tell the model that.

## Positional Encodings

Input:  $(x_1, x_2, x_3, \dots, x_T)$

Position:  $(1, 2, 3, \dots, T)$

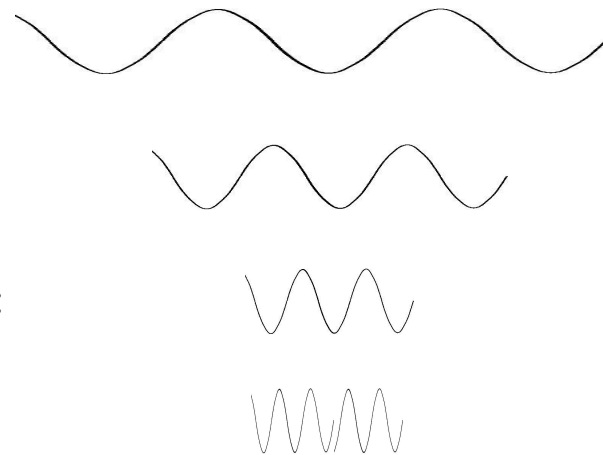
But we can be a bit more clever:

$$PE(t, 2i) = \sin(t/10000^{2i/d})$$

$$PE(t, 2i+1) = \cos(t/10000^{2i/d})$$

$$PE(t) = [PE(t, 0), PE(t, 1), \dots, PE(t, d)]$$

Add embedding of t-th token  $e_t = e_t + PE(t)$ .



Transformer encoder

## Positional Encodings

Input:  $(x_1, x_2, x_3, \dots, x_T)$

Position:  $(1, 2, 3, \dots, T)$

But we can be a bit more clever:

$$PE(t, 2i) = \sin(t/10000^{2i/d})$$

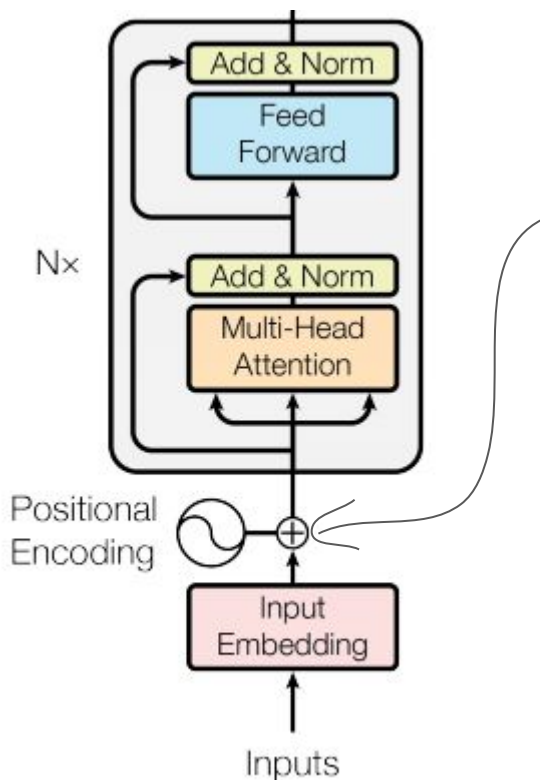
$$PE(t, 2i+1) = \cos(t/10000^{2i/d})$$

$$PE(t) = [PE(t, 0), PE(t, 1), \dots, PE(t, d)]$$

Add embedding of  $t$ -th token  $e_t = e_t + PE(t)$ .

- Assigns every timestep a unique waveform
- No need to specify maximum length

Transformer encoder



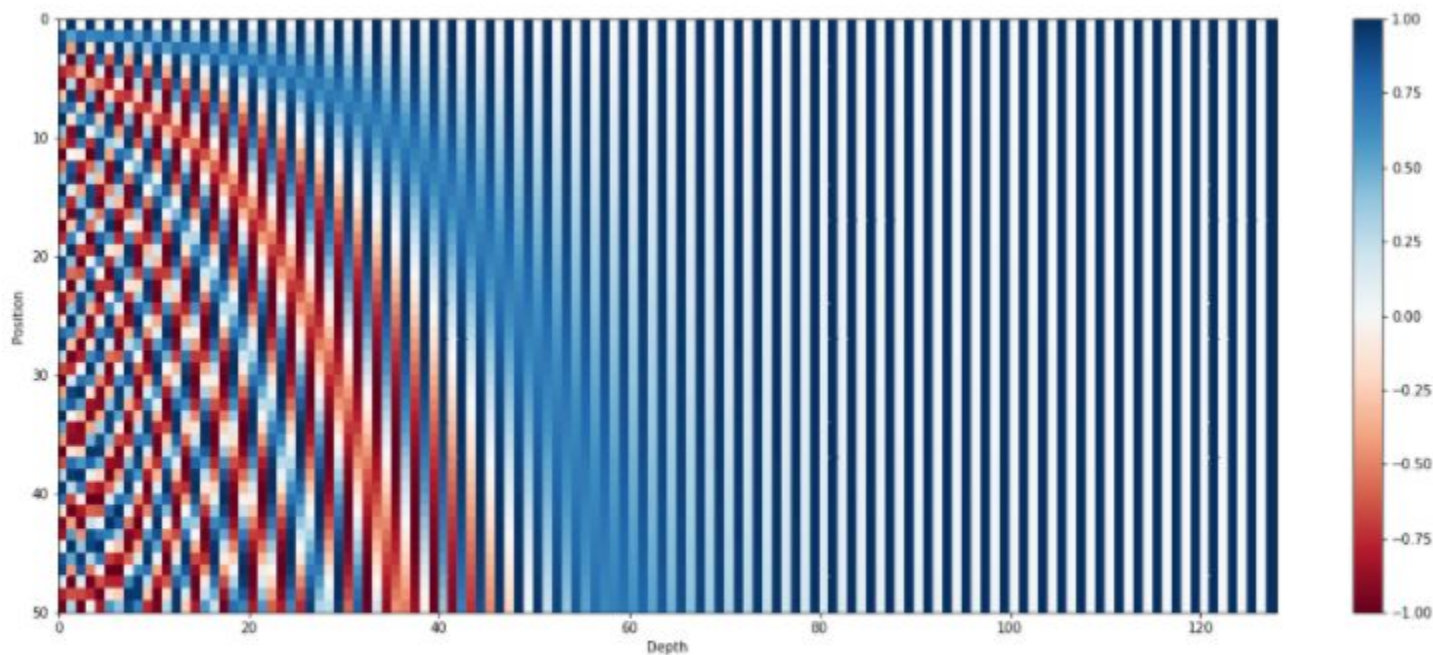
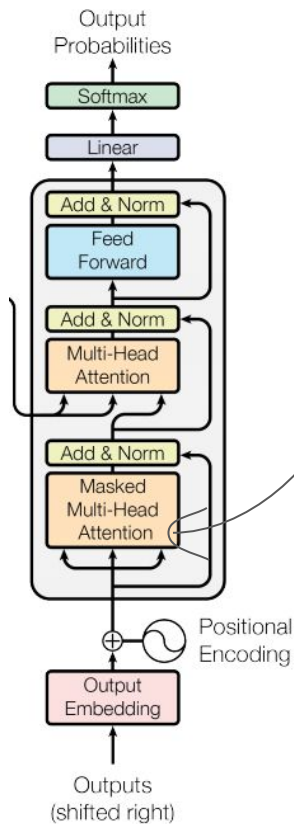


Figure 2 - The 128-dimensional positional encoding for a sentence with the maximum length of 50. Each row represents the embedding vector  $\vec{p}_t$



## Masked Multi-head Attention:

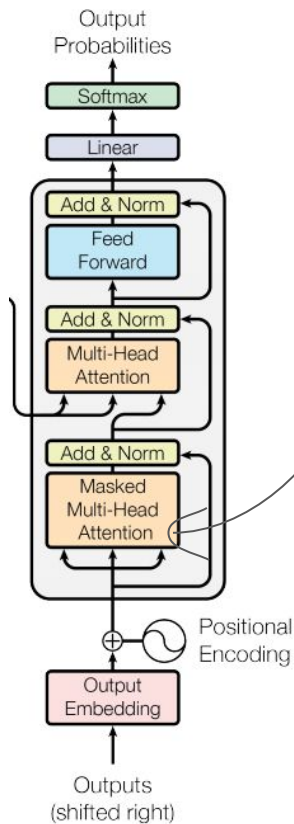
We can't do exactly what we do in the encoder b/c we don't want to bleed future info.

$$(x_1, x_2, \dots, x_{t-1}, \textcolor{red}{x_t}, x_{t+1}, \dots, x_{T-1}, x_T)$$

current

Cheating if we see this b/c in test time, we don't have access to  $> t+1$

Transformer decoder

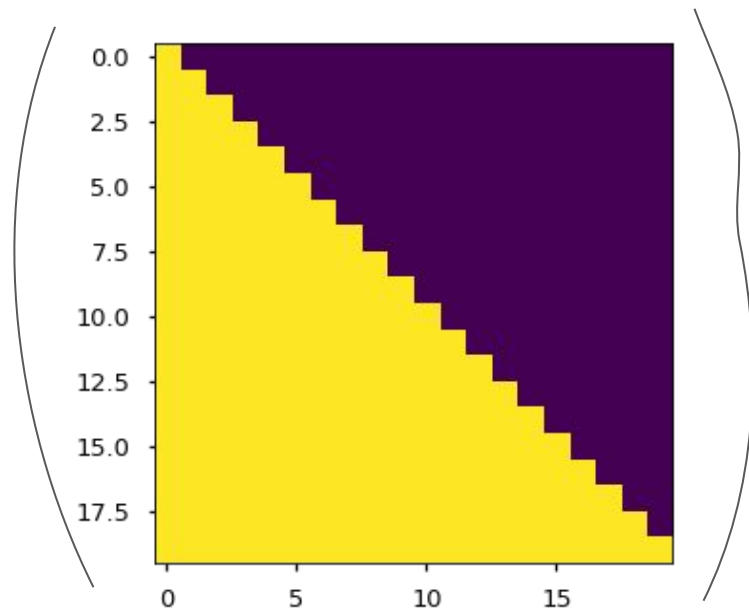


## Masked Multi-head Attention:

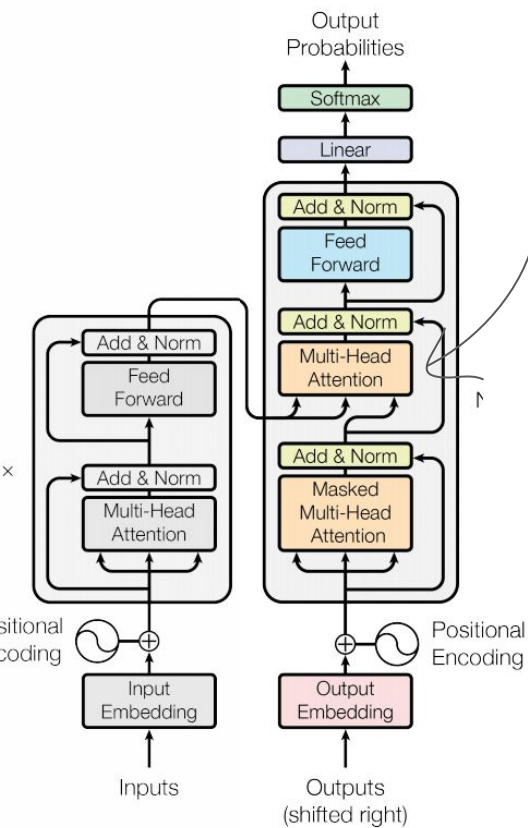
Recall,  $q$ : query ;  $k_1, \dots, k_T$ : keys ;  $v_1, \dots, v_T$ : values

$$\text{maskedsim}(q, k_t, m) = m^T (q^T k_t) / \text{sqrt}(d)$$

$m =$



Transformer decoder



## Encoder Multi-head Attention:

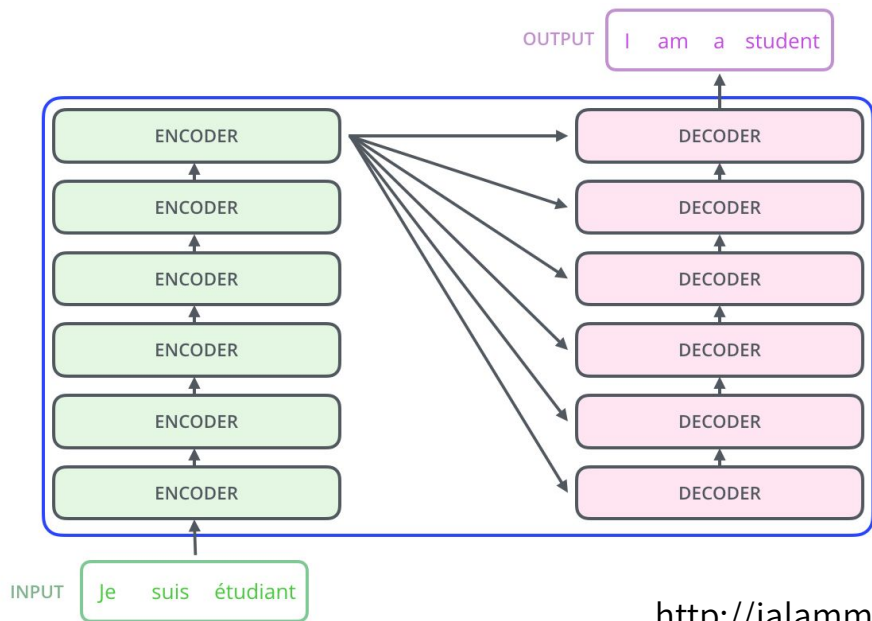
- Output of encoder:  $(h_1^{enc}, h_2^{enc}, \dots, h_T^{enc})$ .
- Use this for keys and values in attention.
- Query vectors come from decoder.
- This blends information from encoder into the decoder.
- Note: no bleeding problem here!
- SOTA speech recognition network is a variant of this idea for online+offline training with weight sharing

Transformer decoder

# Stacked Transformers

---

The trend is make things deep. A single transformer encoder or decoder returns a sequence of the same signature as the input.



<http://jalammar.github.io/illustrated-transformer>



# Resources

---

Neural network basics: <http://cs229.stanford.edu>, <https://www.deeplearningbook.org>

RNNs: <http://web.stanford.edu/class/cs224n/index.html#schedule> , Sequence to sequence learning with neural networks [Sutskever et. al. 2014]

Transformers: Attention is all you need [Vaswani et. al. 2017],  
<http://jalammar.github.io/illustrated-transformer>

Code: <https://huggingface.co/transformers>

CS224N deep learning tutorial videos (linked on course page)

Starting thinking about project ideas!

# Appendix

---

# Neural Chatbots

---

How do you train a neural network to chat?

# Neural Chatbots

---

How do you train a neural network to chat?

- Handwritten rules? (Elizabet)... interesting but wouldn't pass turing
- Finite state machines?... good for some tasks but too limited

# Neural Chatbots

---

How do you train a neural network to chat?

- Handwritten rules? (Elizabet)... interesting but wouldn't pass turing
- Finite state machines?... good for some tasks but too limited
- Treat it like a **translation** problem! End-to end neural networks.

Machine  
Translation

Hi how are you?



Hola! Cómo estás?

Chatbot

Hi how are you?



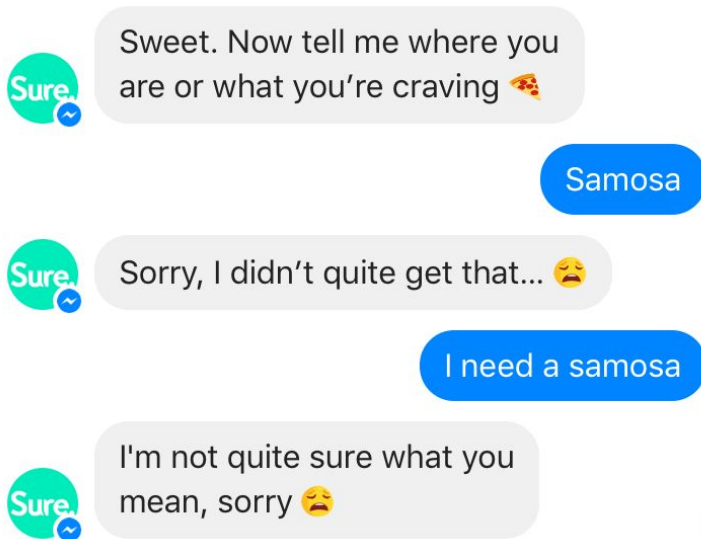
Not bad, you?

# Challenges of Chatbots

---

Even the best encoder decoder model doesn't "solve" chat bots.

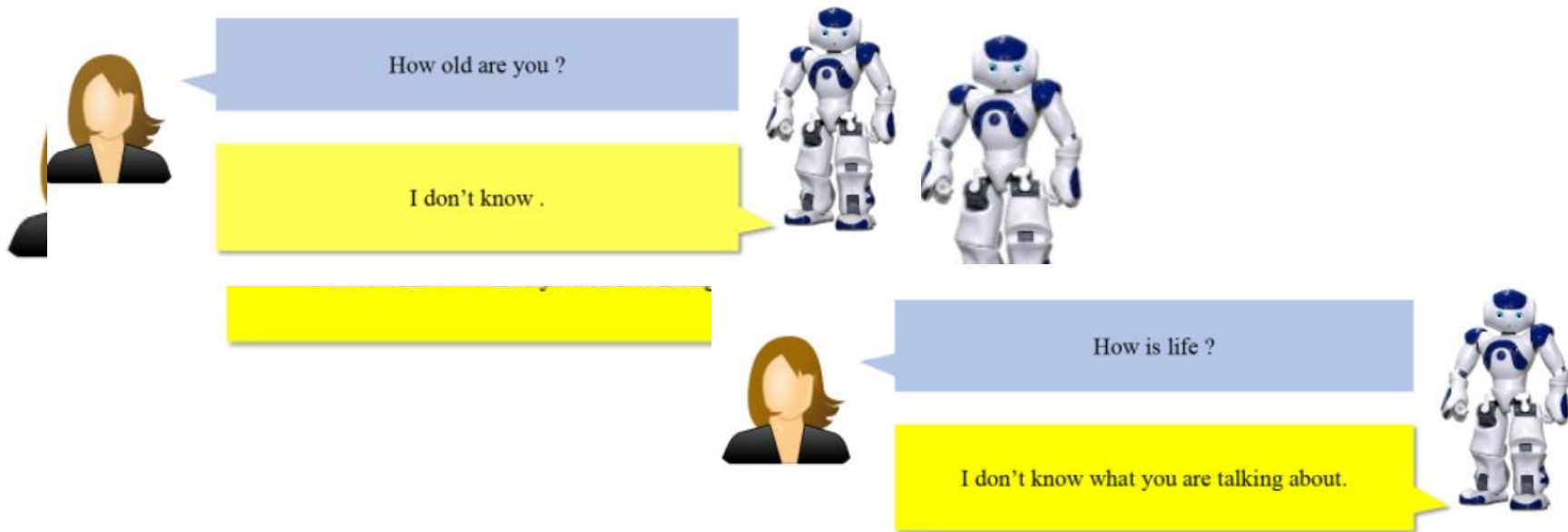
Here are 2 examples of failure modes and possible improvements.



# Generic Responses

---

Problem: Easy for the model to say something in domain but generic in response to everything e.g. “I don’t know” [Sordoni et. al. 2015]





# Generic Responses

---

Problem: Easy for the model to say something in domain but generic in response to everything e.g. “I don’t know” [Sordoni et. al. 2015]

A Solution: Auxiliary objectives!

Optimize for high mutual information between source and response!

$$J = \underbrace{-\log p(\text{“I don’t know”} \mid \text{“how’s life”})}_{\text{Regular objective}} + \underbrace{\text{MI}(\text{“i don’t know”}, \text{“how’s life”})}_{\text{Regularization}} \sim 0$$

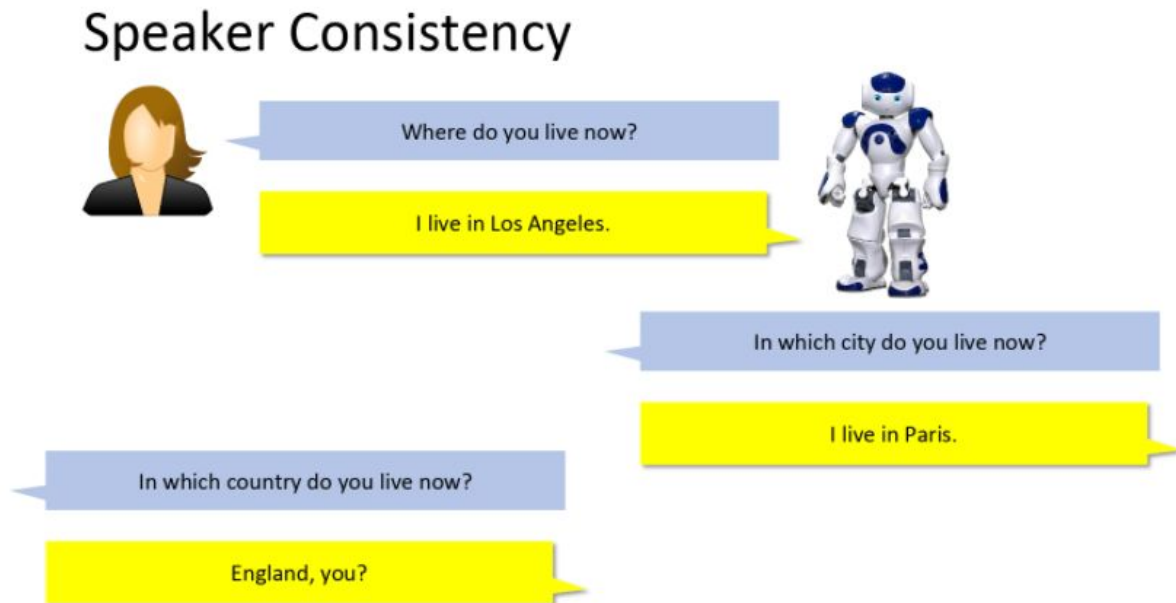
Regular objective

Regularization

# Speaker Consistency

---

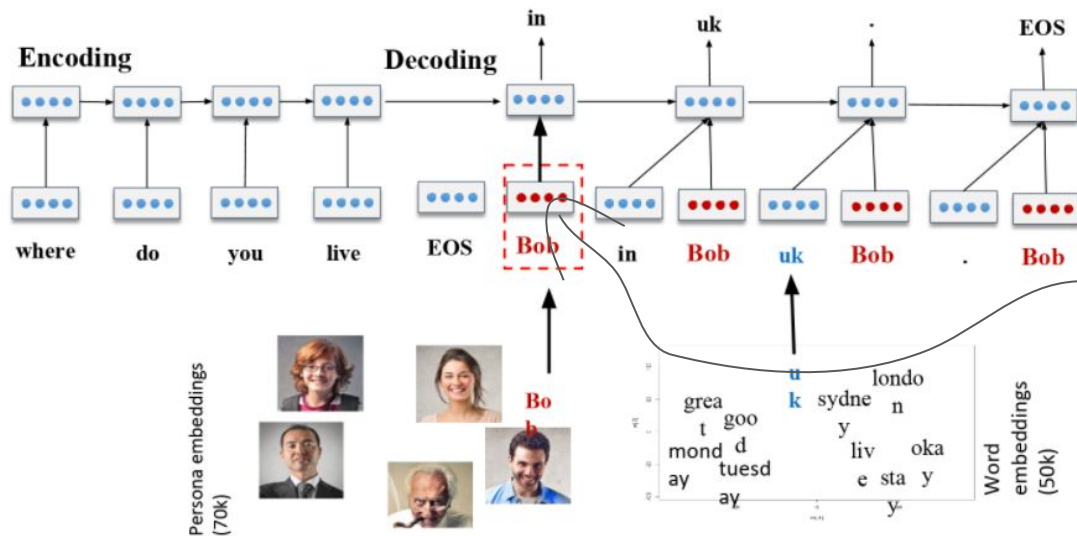
Problem: Just English sounding responses isn't enough. Chatbots should not contradict themselves factually.



# Speaker Consistency

Problem: Just English sounding responses isn't enough. Chatbots should not contradict themselves factually.

A Solution: Auxiliary information! Remember who you are talking to!

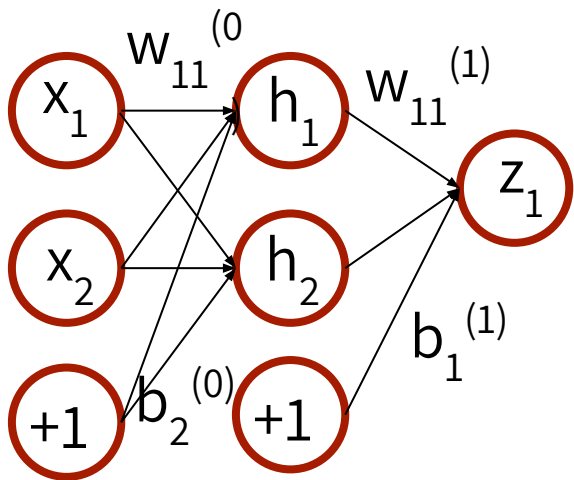


Can always add more info.

# Backpropagation

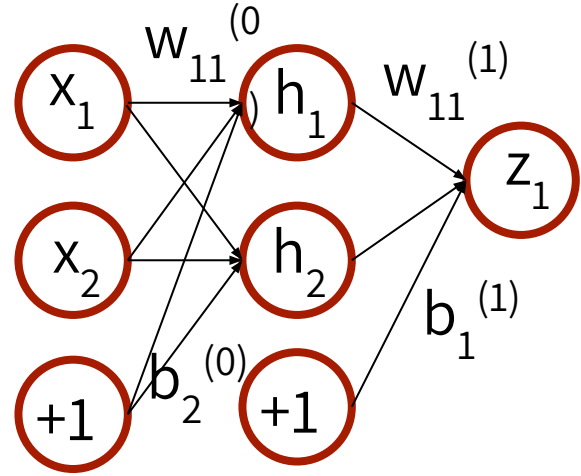
---

Do a simple example with a small MLP.



Compute  $\nabla_{\theta} L(x, y, \theta)$  of objective wrt parameters.

$$\nabla_{\theta} J(x, y, \theta) = \begin{bmatrix} dJ/dw_{11}^{(0)} & dJ/dw_{11}^{(1)} \\ dJ/dw_{12}^{(0)} & dJ/dw_{21}^{(1)} \\ dJ/db_1^{(0)} & dJ/db_1^{(1)} \\ dJ/dw_{21}^{(0)} & \\ dJ/dw_{22}^{(0)} & \\ dJ/db_2^{(0)} & \end{bmatrix}$$



# Autodifferentiation

---

- Deriving these by hand is annoying.
- If you have new objective functions, this could be really intractable.
- **Idea:** if you manually specify the derivative for a set of “basic” operations, you can calculate derivative of complicated functions using chain rule.



Magic :)



```
model.train()
for batch_idx, (data, target) in enumerate(train_loader):
    data, target = data.to(device), target.to(device)
    optimizer.zero_grad()
    output = model(data)
    loss = F.nll_loss(output, target)
    loss.backward()
    optimizer.step()
    if batch_idx % args.log_interval == 0:
        print('Train Epoch: {} [{}/{}] ({:.0f}%) \t Loss: {:.6f}'.format(
            epoch, batch_idx * len(data), len(train_loader.dataset),
            100. * batch_idx / len(train_loader), loss.item()))
    if args.dry_run:
        break
```



# Long Short Term Memory

---

LSTMs have the ability to “forget” information and “store” information that could be useful later [Schmidhuber 1997].

**Forget gate:**  $f_t = \sigma(U_f h_{t-1} + W_f x_t)$  Pick what to “forget”

$k_t = c_{t-1} * f_t$  Do the “forgetting”

**Add gate:**  $g_t = \sigma(U_g h_{t-1} + W_g x_t)$  Usual RNN function

$i_t = \sigma(U_i h_{t-1} + W_i x_t)$  Pick what to “add”

$j_t = g_t * i_t$  Do the “adding”

$c_t = j_t + k_t$  context is some of last context and some new stuff

Forget gate:  $f_t = \sigma(U_f h_{t-1} + W_f x_t)$  Pick what to “forget”

$k_t = c_{t-1} * f_t$  Do the “forgetting”

Add gate:  $g_t = \sigma(U_g h_{t-1} + W_g x_t)$  Usual RNN function

$i_t = \sigma(U_i h_{t-1} + W_i x_t)$  Pick what to “add”

$j_t = g_t * i_t$  Do the “adding”

$c_t = j_t + k_t$  context is some of last context and some new stuff

Output gate:  $o_t = \sigma(U_o h_{t-1} + W_o x_t)$  Pick what to use for current timestep

$H_t = o_t * \tanh(c_t)$  Do the partitioning

Forget gate:  $f_t = \sigma(U_f h_{t-1} + W_f x_t)$  Pick what to “forget”

$k_t = c_{t-1} * f_t$  Do the “forgetting”

Add gate:  $g_t = \sigma(U_g h_{t-1} + W_g x_t)$  Usual RNN function

$i_t = \sigma(U_i h_{t-1} + W_i x_t)$  Pick what to “add”

$j_t = g_t * i_t$  Do the “adding”

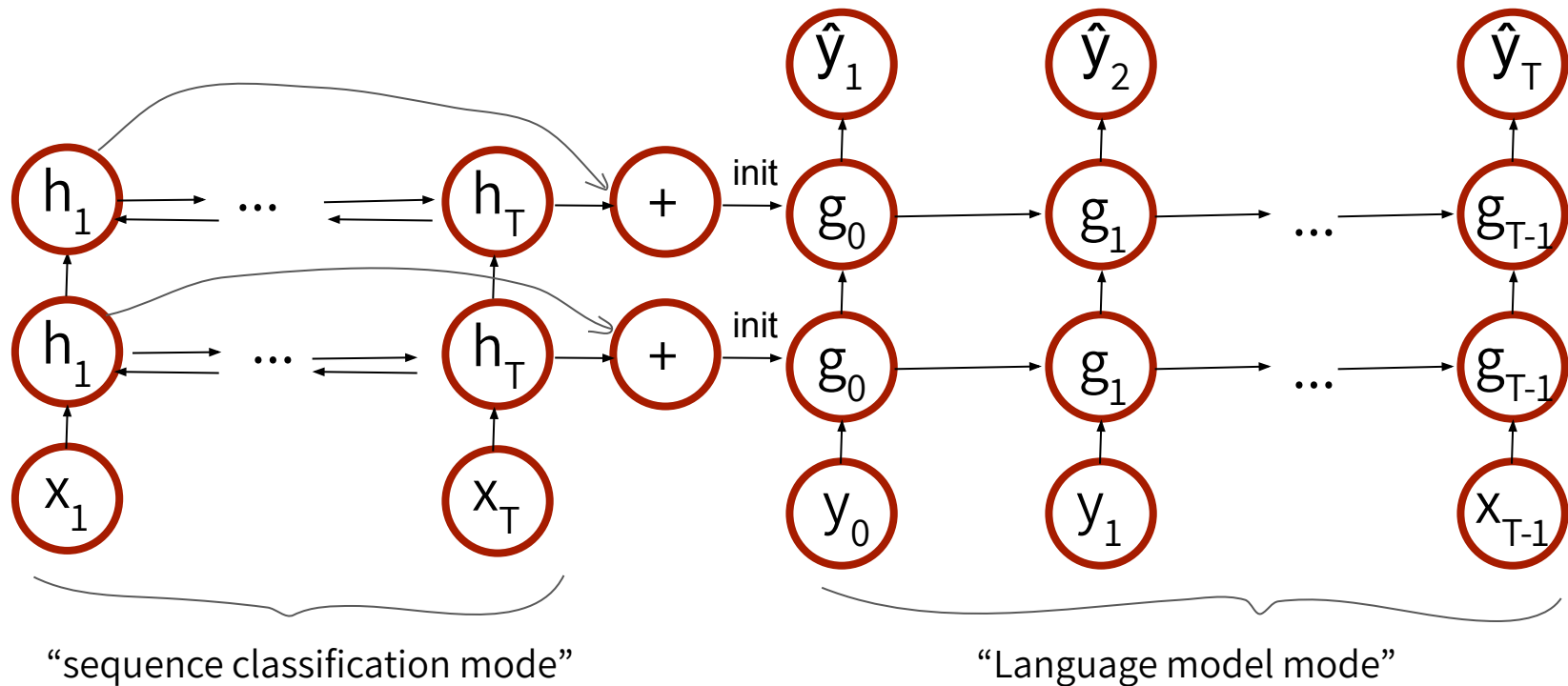
$c_t = j_t + k_t$  context is some of last context and some new stuff

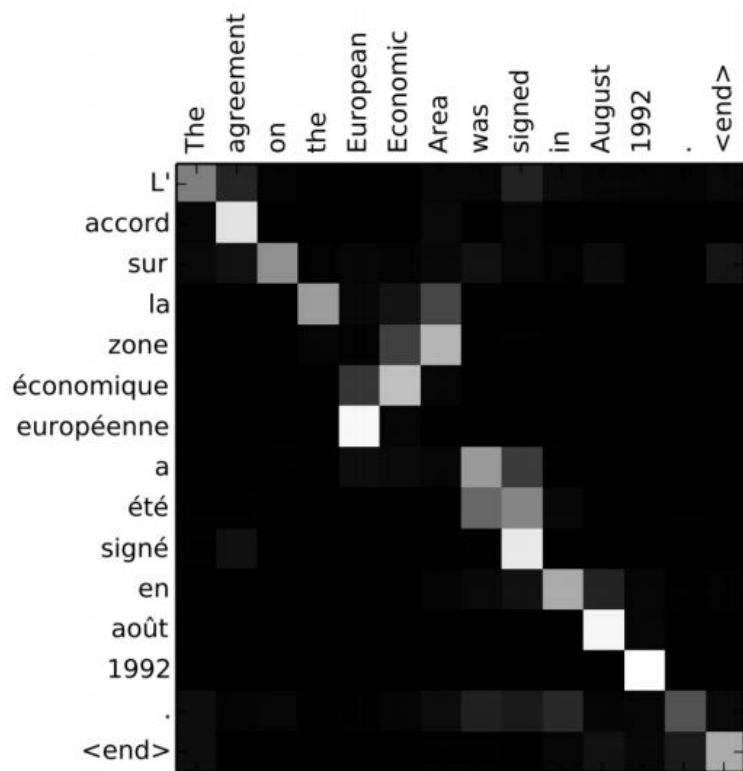
Output gate:  $o_t = \sigma(U_o h_{t-1} + W_o x_t)$  Pick what to use for current timestep

$H_t = o_t * \tanh(c_t)$  Do the partitioning

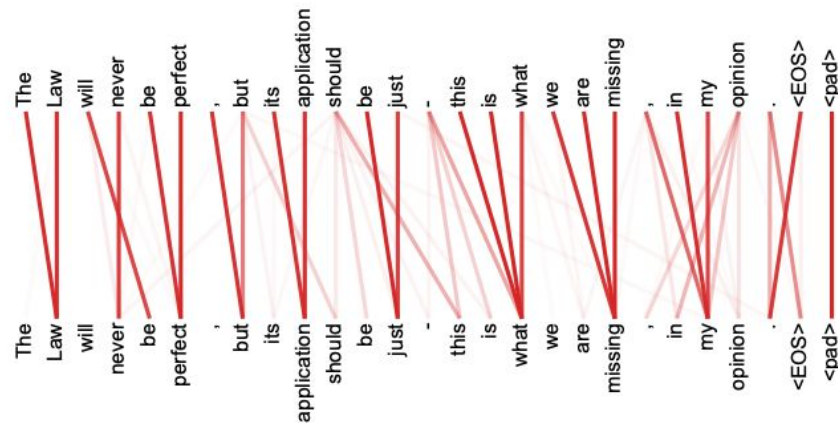
**This is horribly complicated but the intuition is good: separate what  $h_t$  is good for. Some of it is good for right now ( $y_t$ ); some of it is good for later!**

Classic deep learning: add more layers!





Regular Attention

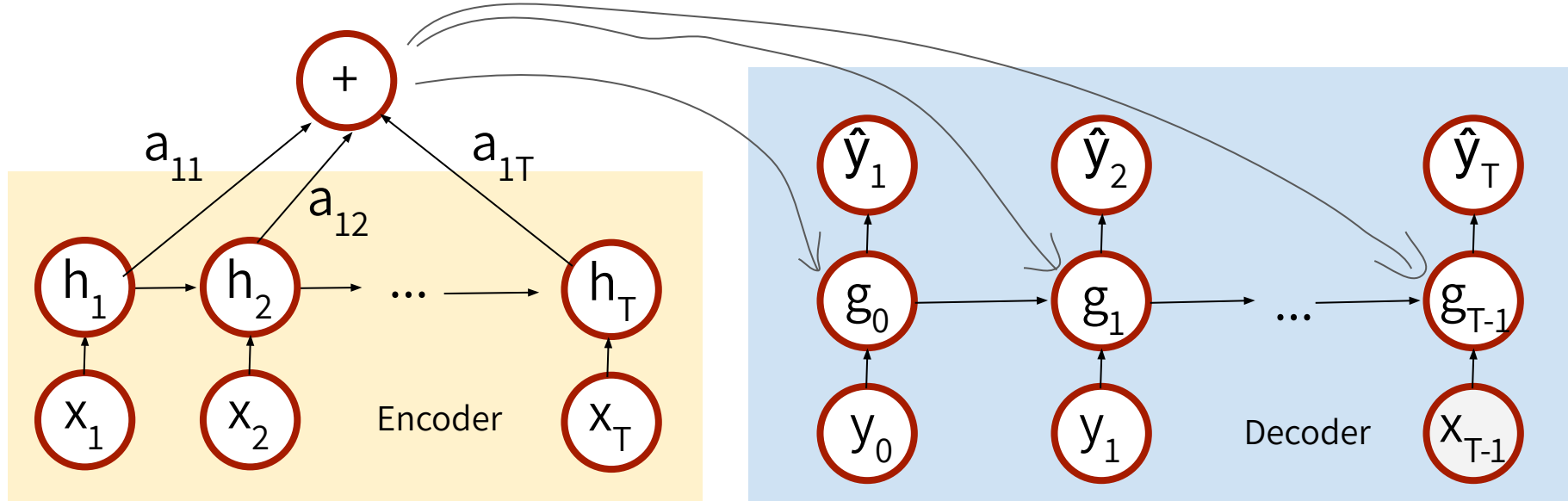


Self Attention

# Attention-Seq2seq

---

**Intuition:** Every timestep in decoder has its own attention to  $h_1, \dots, h_T$



- Suppose we are at timestep **\*\*i\*\*** in the decoder.
- Use  $g_{i-1}$  as query. Also  $h_1, \dots, h_T$  doubles as keys and values (self-attn)!

$$a_{it} = \frac{\exp\{\text{sim}(g_{i-1}, h_t)\}}{\sum_{s=1}^T \exp\{\text{sim}(g_{i-1}, h_s)\}}$$

$$c_i = \sum_{t=1}^T a_{it} h_t$$

$$g_i = \text{decoder\_rnn}(y_i, g_{i-1}, c_i)$$