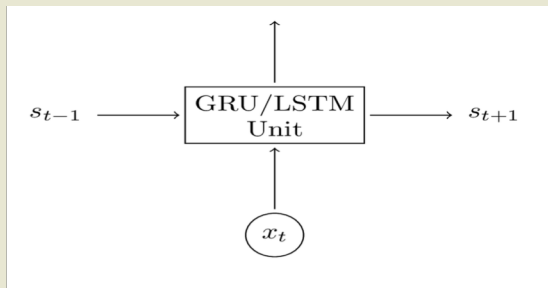
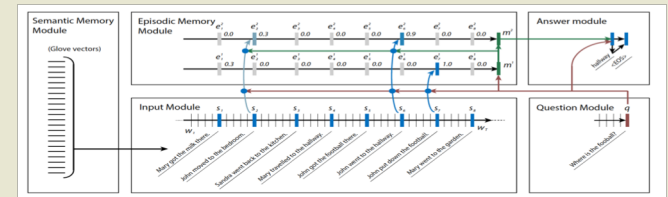




Recurrent Neural Networks



Time Series



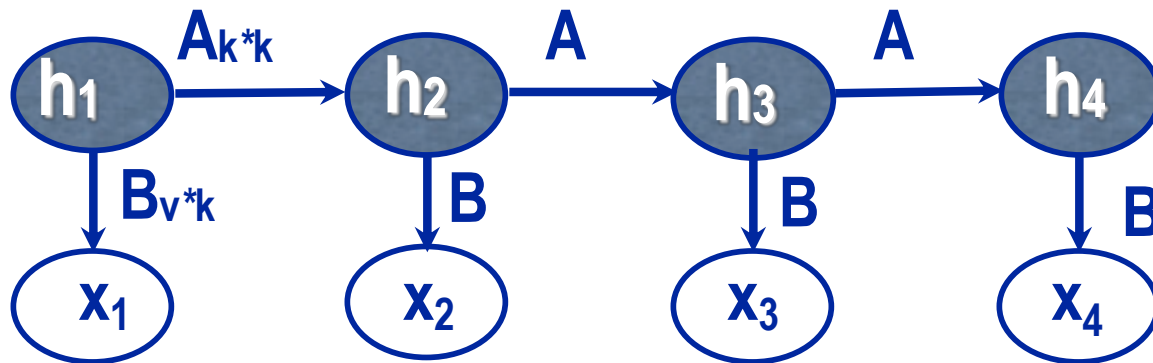
Recurrent Neural Nets

- ◆ **Generalize HMMs or Linear Dynamical Systems**
 - Hidden state dynamical models, but *nonlinear*
- ◆ **Needed if you have inputs of varying length**
 - E.g. sequence of observations
 - speech
 - text
 - robots
 - power plants, chemical plants, data centers



Standard HMM

- ◆ HMM learning problem: Estimate A and B



A = Markov transition matrix
 B = emission matrix

- ◆ Estimation done via EM
 - Or spectral methods
- ◆ History is forgotten with an exponential decay



Simple Recurrent Neural Net

$$s_t = \tanh(Ux_t + Ws_{t-1})$$

$$o_t = \text{softmax}(Vs_t)$$

x_t = input (e.g. a word)

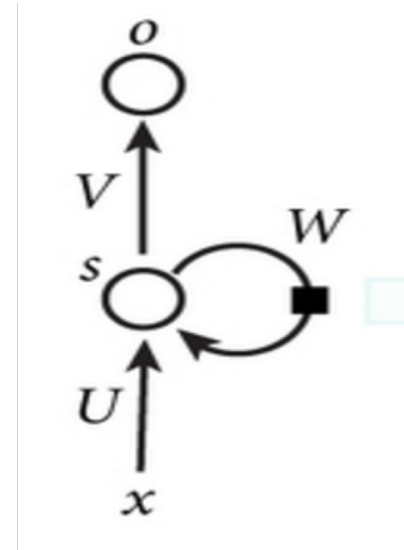
s_t = hidden state

o_t = output (e.g. probability of the next word)

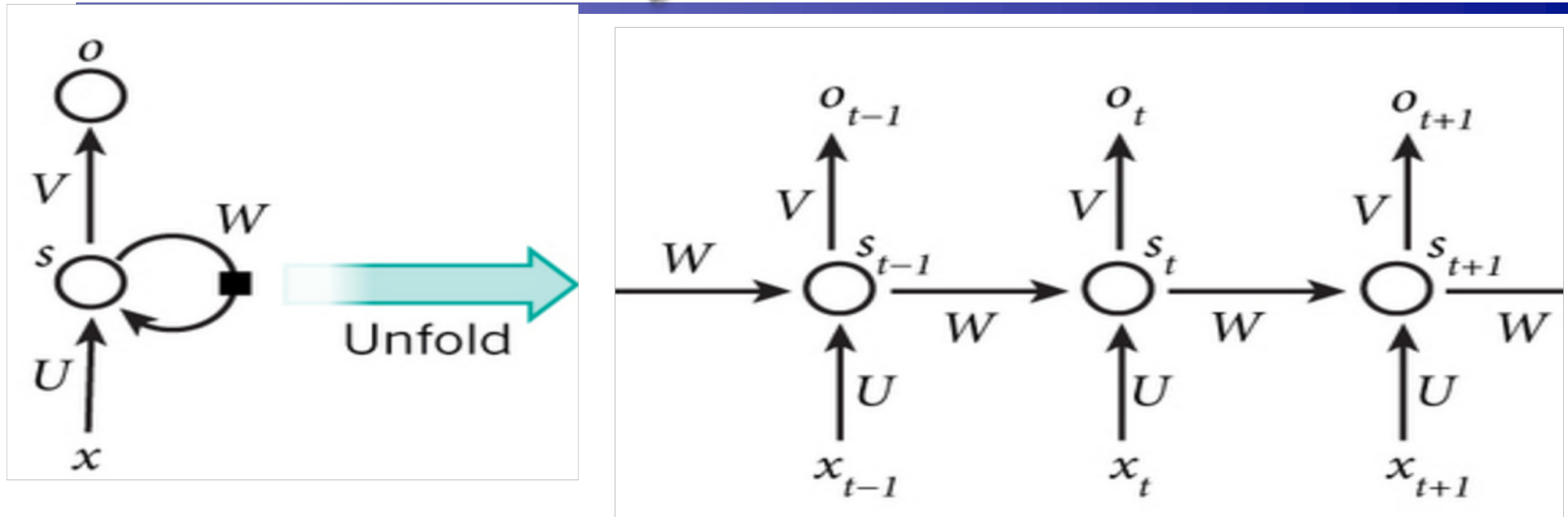
y_t = true value (e.g. x_{t+1})

Softmax $\sigma(\mathbf{z})$ transforms the K-dimensional real valued output \mathbf{z} to a distribution – like logistic regression

$$\sigma(\mathbf{z})_j = \frac{e^{z_j}}{\sum_{k=1}^K e^{z_k}} \quad \text{for } j = 1, \dots, K.$$



Like HMMs, unroll RNNs in time



x_t = input (e.g. a word)

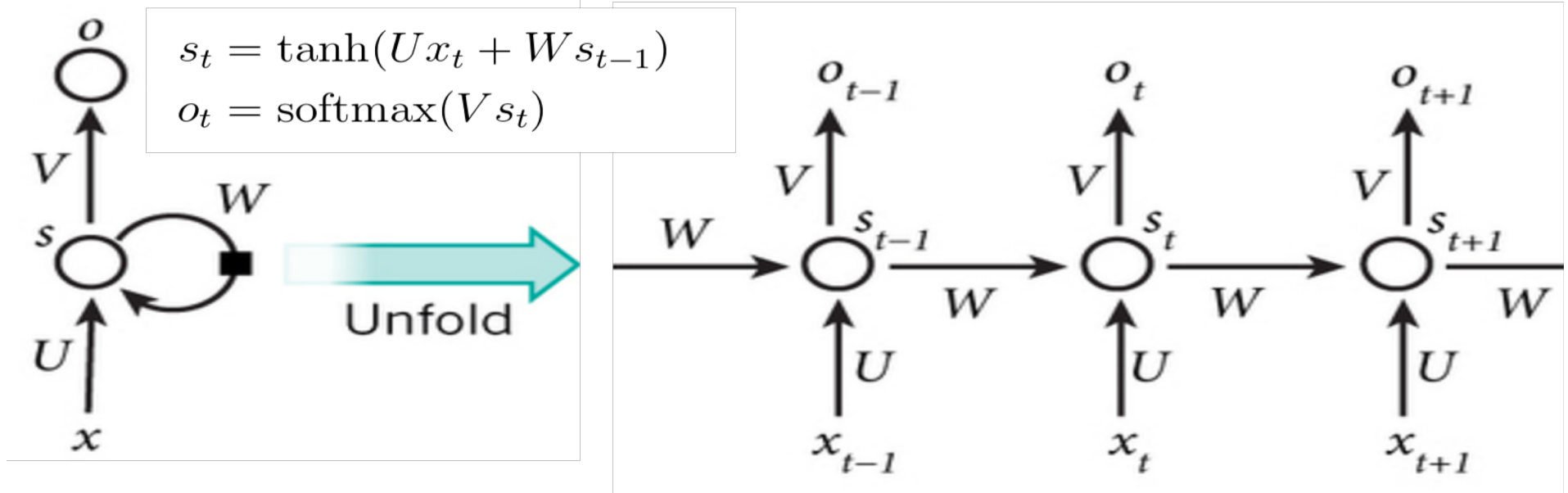
s_t = hidden state

o_t = output (e.g. probability of the next word)

<http://www.nature.com/nature/journal/v521/n7553/full/nature14539.html>



Like HMMs, unroll RNNs in time

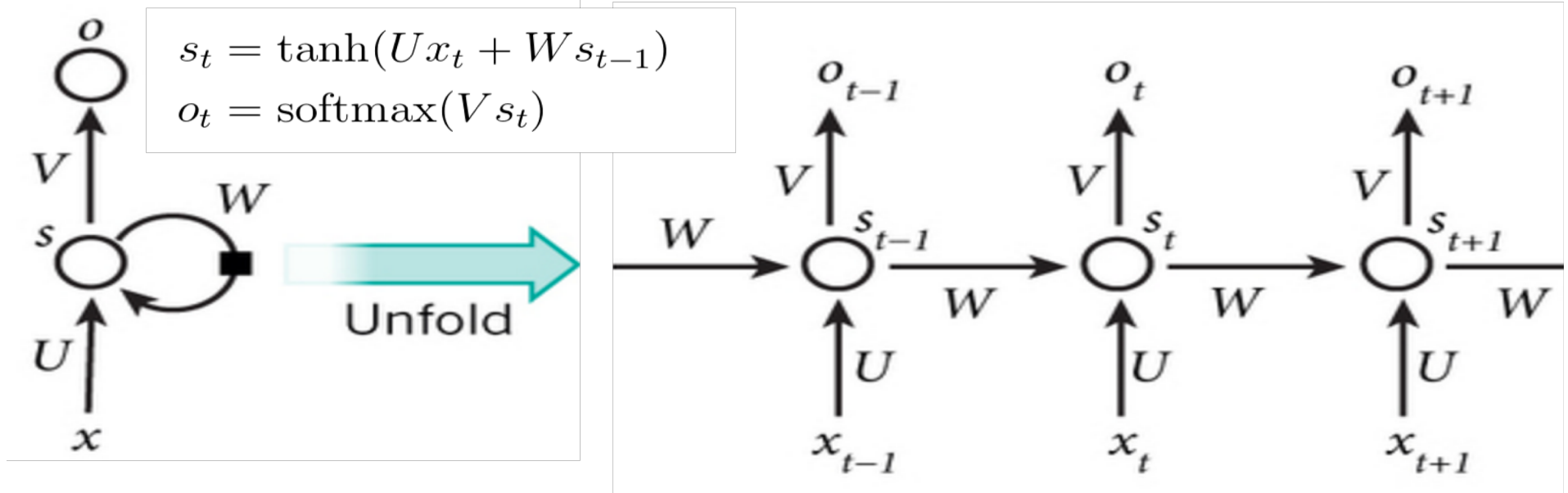


x_t = input (e.g. a word) - v
 s_t = hidden state - k
 o_t = output - v

What are the dimensions of U , W , V ?
 U $k \times v$ W $k \times k$ V $v \times k$



Like HMMs, unroll RNNs in time

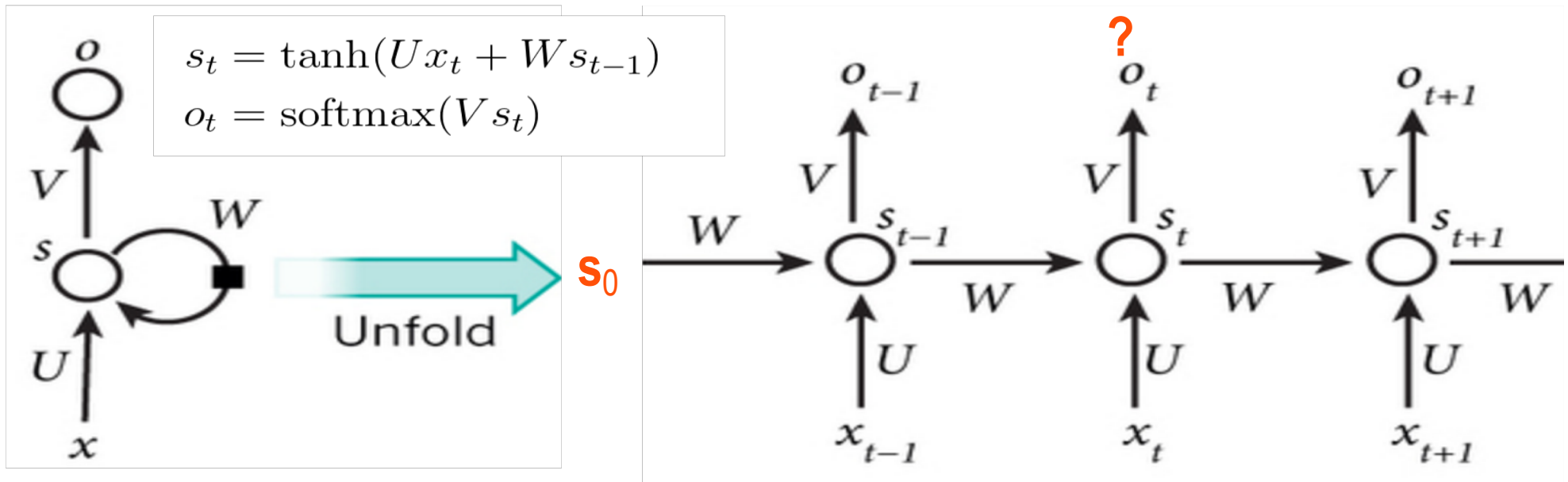


x_t = input (e.g. a word) - v
 s_t = hidden state - k
 o_t = output - v

What is the usual loss function?
 $-\sum_t \log(o_t[y_t])$ - *est. prob. of truth*
where $y_t=i$ gives the true label



Like HMMs, unroll RNNs in time



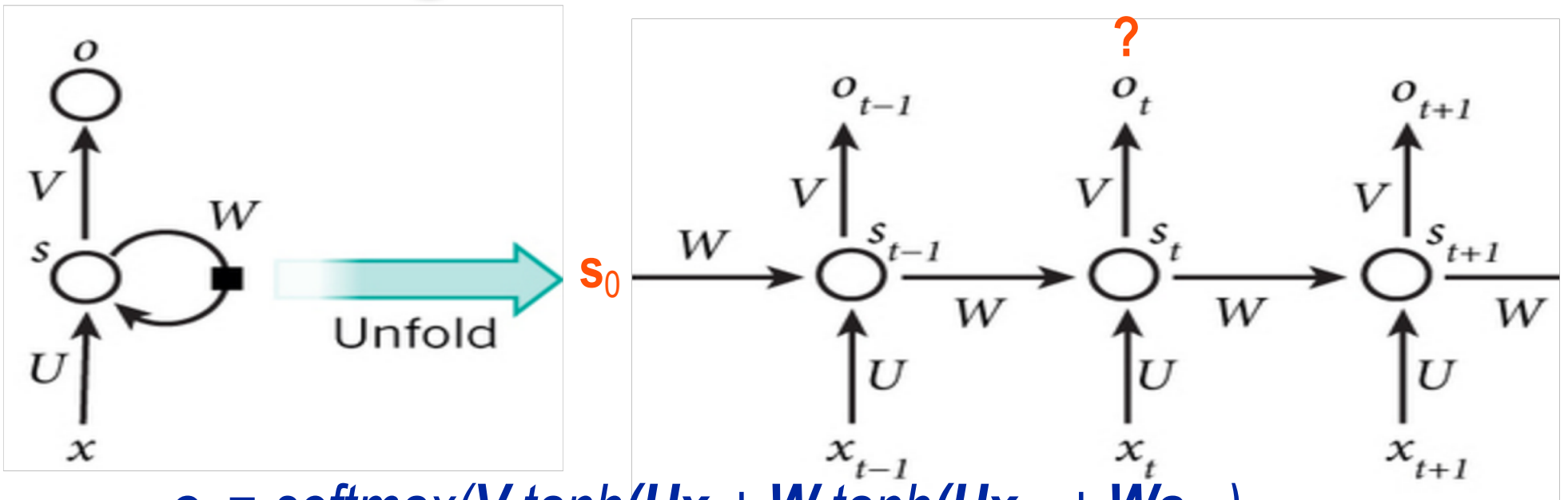
x_t = input - v
 s_t = hidden state - k
 o_t = output - v

If $s_{t-2} = s_0$, what is o_t in terms of s_0 and x ?

$$\begin{aligned}
 o_t &= \text{softmax}(Vs_t) = \text{softmax}(V \tanh(Ux_t + Ws_{t-1})) \\
 &= \text{softmax}(V \tanh(Ux_t + W \tanh(Ux_{t-1} + Ws_{t-2})))
 \end{aligned}$$



RNN gradients



$$o_t = \text{softmax}(V \tanh(Ux_t + W \tanh(Ux_{t-1} + Ws_{t-2})))$$

Observe $y_t = i$ What is the stochastic gradient step?

$$Err = -\log(o_t[i])$$

Find $d Err/dV$, $d Err/dU$, $d Err/dW$



RNN Gradients

- ◆ $\mathbf{o}_t = \text{softmax}(\mathbf{V} \tanh(\mathbf{U}\mathbf{x}_t + \mathbf{W} \tanh(\mathbf{U}\mathbf{x}_{t-1} + \mathbf{W}\mathbf{s}_{t-2}))$
- ◆ **Observe** $y_t = i$ What is the stochastic gradient step?
- ◆ $\text{Err} = -\log(\mathbf{o}_t[i])$

$$\begin{aligned} d \text{Err}/d\mathbf{V} &= -(d \log(\mathbf{o}_t[i])/d\mathbf{o}_t[i]) \quad d\mathbf{o}_t[i]/d\mathbf{V} \\ &= -(1/\mathbf{o}_t[i]) \quad d \text{softmax}(\mathbf{z})/d\mathbf{z} \quad d\mathbf{z}/d\mathbf{V} \end{aligned}$$

$$\mathbf{z} = \mathbf{V} \tanh(\mathbf{U}\mathbf{x}_t + \mathbf{W} \tanh(\mathbf{U}\mathbf{x}_{t-1} + \mathbf{W}\mathbf{s}_{t-2}))$$

$$\begin{aligned} d \text{softmax}(\mathbf{z})/dz_j &= -1/(\sum_k e^{z_k})^2 e^{z_j} e^{z_k} \quad \text{for } k \text{ not equal to } j \\ &= -1/(\sum_k e^{z_k})^2 e^{2z_j} + e^{z_j}/(\sum_k e^{z_k}) \quad \text{for } k=j \end{aligned}$$

$$\sigma(\mathbf{z})_j = \frac{e^{z_j}}{\sum_{k=1}^K e^{z_k}} \quad \text{for } j = 1, \dots, K.$$

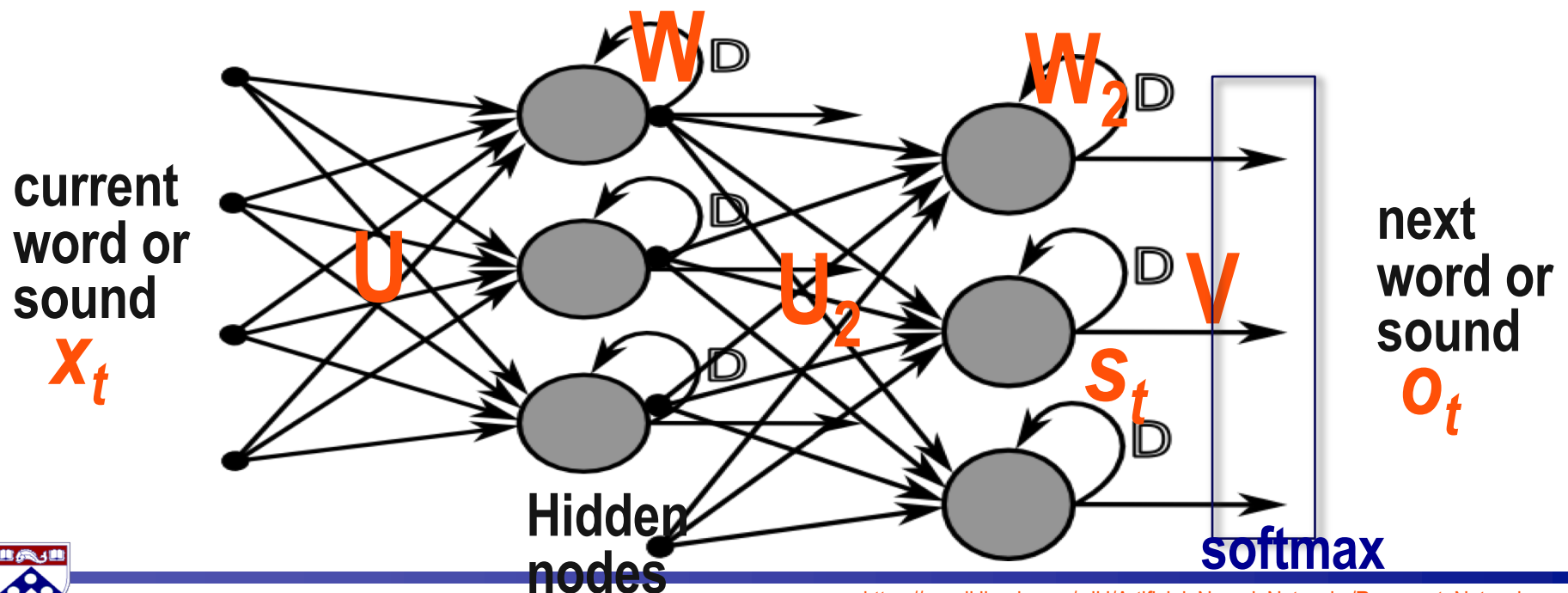


Recurrent Neural Nets (RNNs)

$$s_t = \tanh(Ux_t + Ws_{t-1})$$

$$o_t = \text{softmax}(Vs_t)$$

Can use multiple layers



Gated RNNs

◆ Standard RNNs, like HMMs, tend to forget things exponentially fast

◆ Solution: Gated RNN

- Stores hidden state

$$z = \sigma(U^z x_t + W^z s_{t-1}) \quad z: \text{update gate}$$

$$r = \sigma(U^r x_t + W^r s_{t-1}) \quad r: \text{reset gate}$$

$$h = \tanh(U^h x_t + W^h (s_{t-1} \circ r))$$

$$s_t = (1-z) \circ h + z \circ s_{t-1} \quad s_t: \text{hidden state}$$

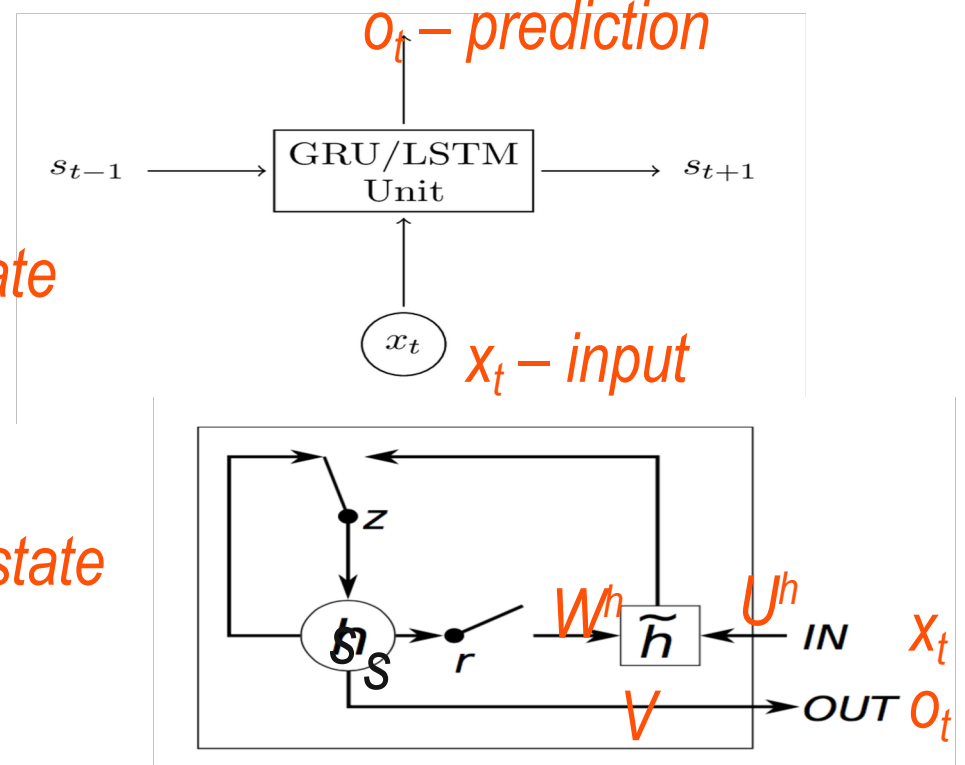
$r=0$ resets h

$z=1$ keeps state

$z=0$ updates it to h

$r=1$'s, $z=0$'s gives simple RNN

\circ is pointwise multiplication



Long Short Term Memory (LSTM)

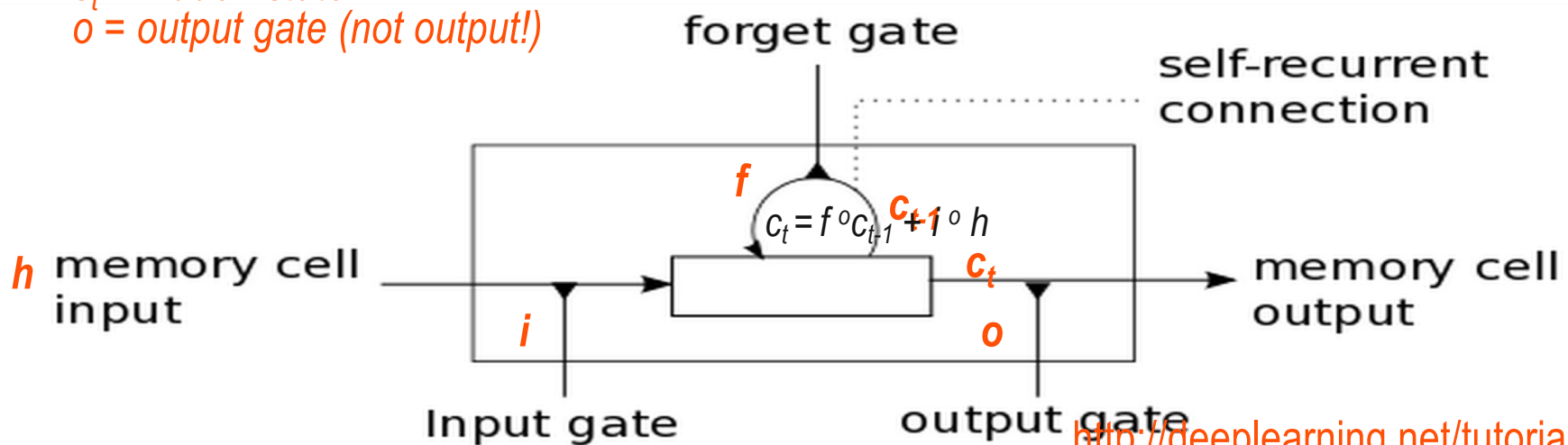
◆ LSTM is a kind of gated RNN

- Just with more, different gates
- Don't worry about what they are!!!

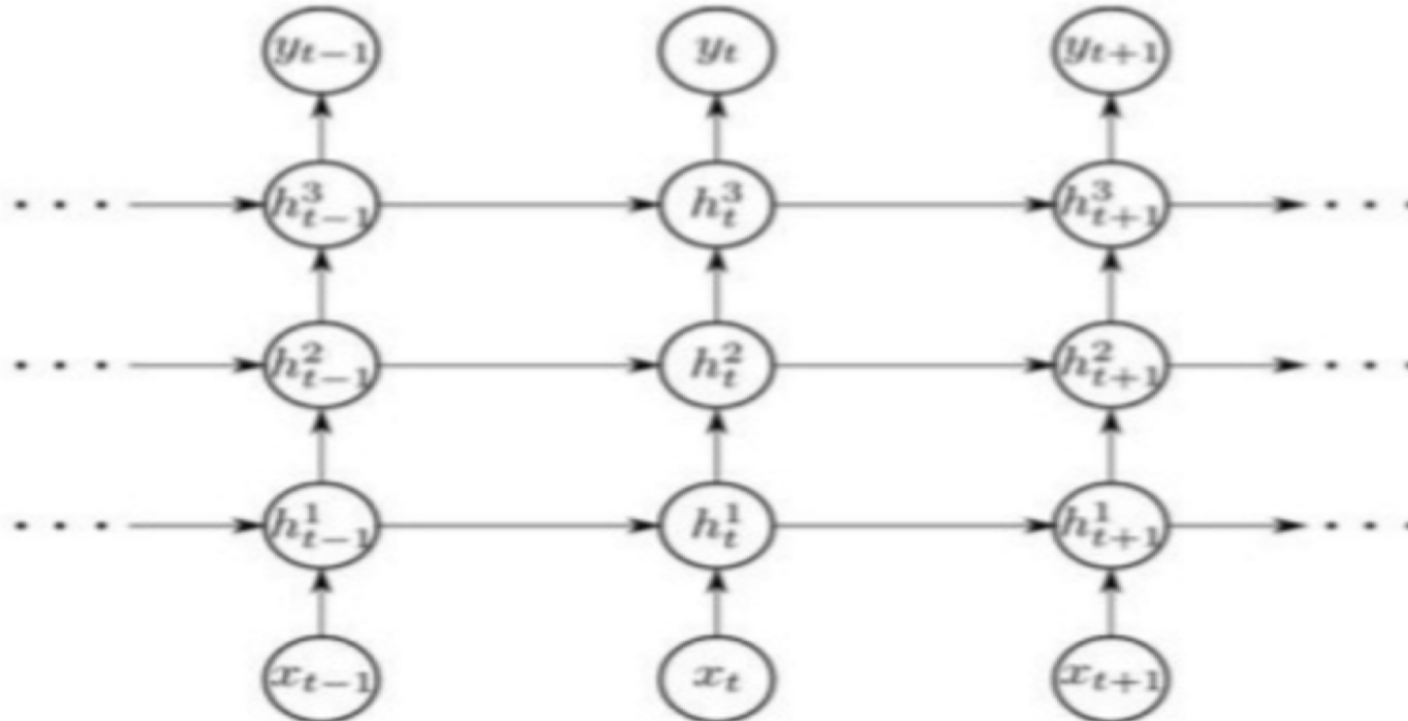
x_t – observation

s_t – hidden state

o = output gate (not output!)



Recurrent Nets can be stacked

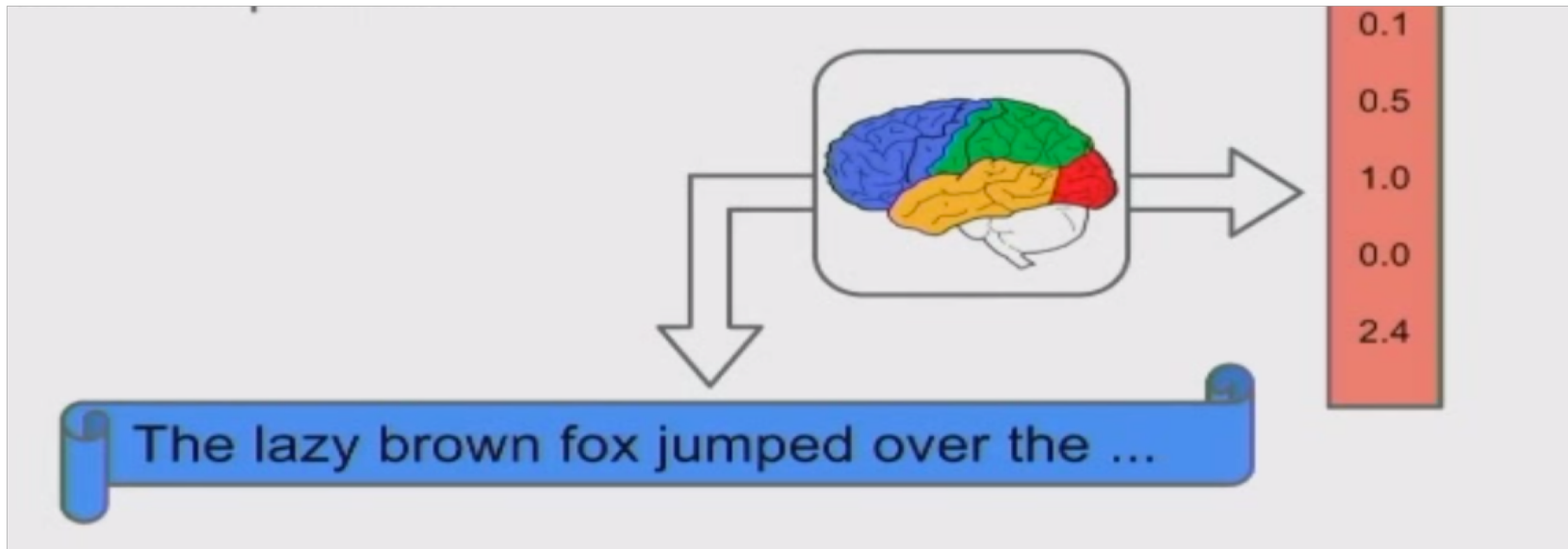


Recurrent Neural Nets

- ◆ Predict a label for each observation
 - $y_t = f(\mathbf{x}_t, \mathbf{s}_t)$
- ◆ Predict the next observation given past observations
 - $y_t = \mathbf{x}_{t+1} = f(\mathbf{x}_t, \mathbf{s}_t)$
- ◆ Or map one sequence to another sequence
 - An encoder
 - sentence (sequence of words) to vector
 - A decoder
 - vector to sentence (sequence of words)



LSTM encodes a sentence

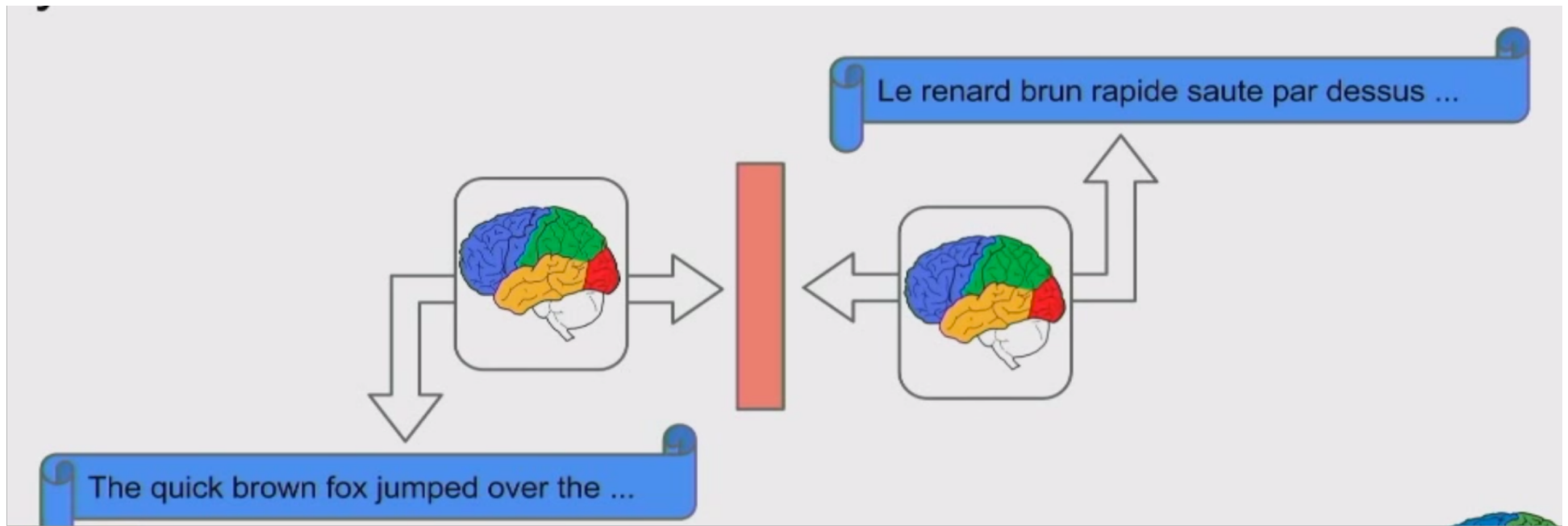


Jeff Dean, google

https://www.youtube.com/watch?v=90-S1M7Ny_o&spfreload=1



Sequence to sequence (Seq2seq)

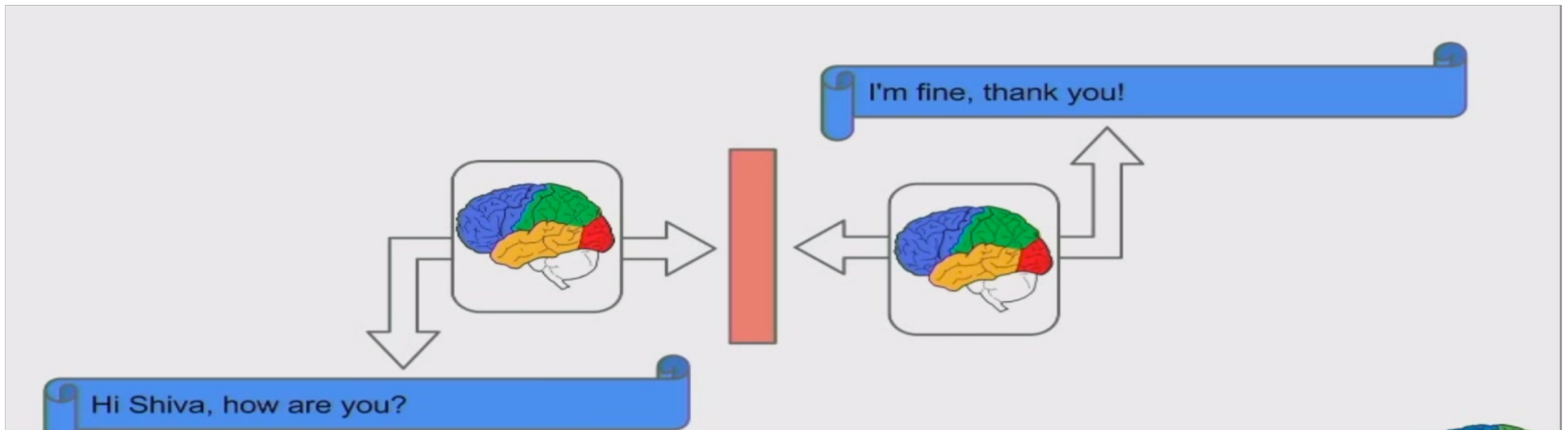


Encode and Decode = translate

Jeff Dean, google



Seq2seq chatbot



Encode and Decode = chatbot

Jeff Dean, google



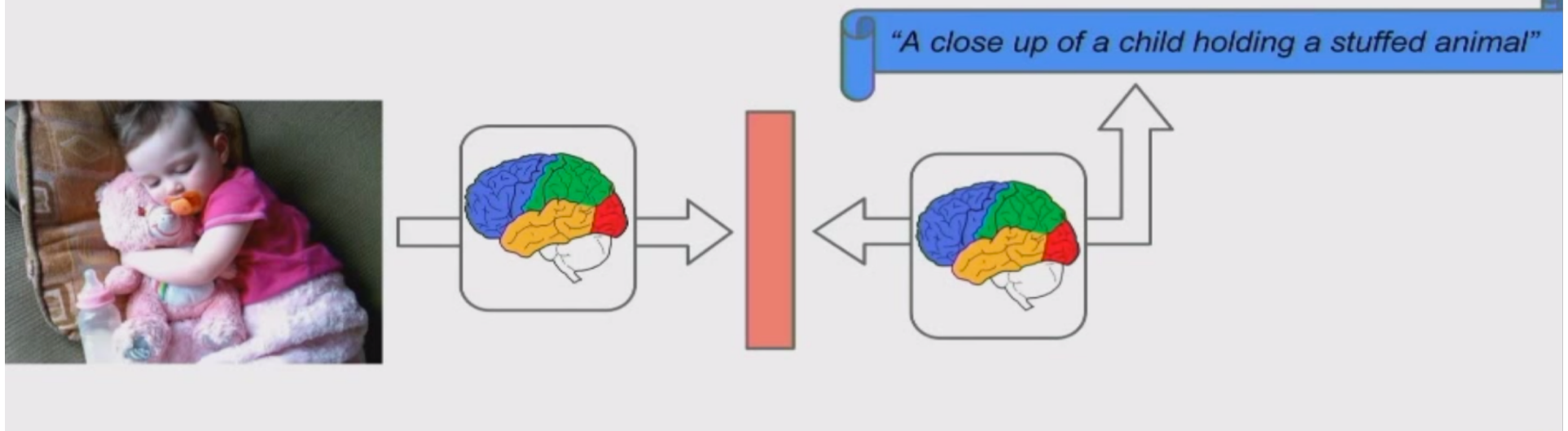
It works (sort of)

```
User> hello , i have a problem with my machine
BrainStop> hi
BrainStop> this is shiva .
BrainStop> how are you doing today ?
User> hi shiva , how are you ?
BrainStop> i ' m fine , thank you !
BrainStop> how may i assist you today ?
User> i want to access using vpn
BrainStop> currently the computer is connected to corp network ?
User> no
BrainStop> check this solution :
BrainStop> https : // INTERNAL_HOSTNAME_REDACTED / articles / connecting - to -
disconnecting - from - ssl - vpn # heading - 4
User> thanks , bye !
BrainStop> thank you .
```



Can map images to words

Connect sequence and image models, you get a captioning system



It works (sort of)



A man holding a tennis racquet on a tennis court.



Two pizzas sitting on top of a stove top oven



A group of young people playing a game of Frisbee



A man flying through the air while riding a snowboard



Language inputs to RNN

- ◆ Words (“one-hot”)
- ◆ Characters (“one-hot”)
- ◆ Word embeddings
 - Typically 300 dimensional



Attention-based Q&A

by *ent423* , *ent261* correspondent updated 9:49 pm et , thu march 19 , 2015 (*ent261*) a *ent114* was killed in a parachute accident in *ent45* , *ent85* , near *ent312* , a *ent119* official told *ent261* on wednesday . he was identified thursday as special warfare operator 3rd class *ent23* , 29 , of *ent187* , *ent265* . `` *ent23* distinguished himself consistently throughout his career . he was the epitome of the quiet professional in all facets of his life , and he leaves an inspiring legacy of natural tenacity and focused

...

ent119 identifies deceased sailor as **X** , who leaves behind a wife

by *ent270* , *ent223* updated 9:35 am et , mon march 2 , 2015 (*ent223*) *ent63* went familial for fall at its fashion show in *ent231* on sunday , dedicating its collection to `` mamma '' with nary a pair of `` mom jeans '' in sight . *ent164* and *ent21* , who are behind the *ent196* brand , sent models down the runway in decidedly feminine dresses and skirts adorned with roses , lace and even embroidered doodles by the designers ' own nieces and nephews . many of the looks featured saccharine needlework phrases like `` i love you ,

...

X dedicated their fall fashion show to moms



Dynamic Network Summary

- ◆ **Gated Neural Nets generalize HMMs, Kalman filters**
 - But are far more powerful!
- ◆ **They have replaced HMMs for speech to text and machine translation**
- ◆ **Lots of black magic “engineering”**
 - Unclear what matters about the network structure
 - Number and size of layers, regularization
 - Forms of gating (LSTM ...), attention ...
 - Gradient descent is tricky
- ◆ **Good software: tensorflow, pytorch ...**

