

# Sequences and Recurrent Neural Nets

Oriol Vinyals  
Staff Research Scientist  
DeepMind  
UCL, Feb 2017

# Roadmap

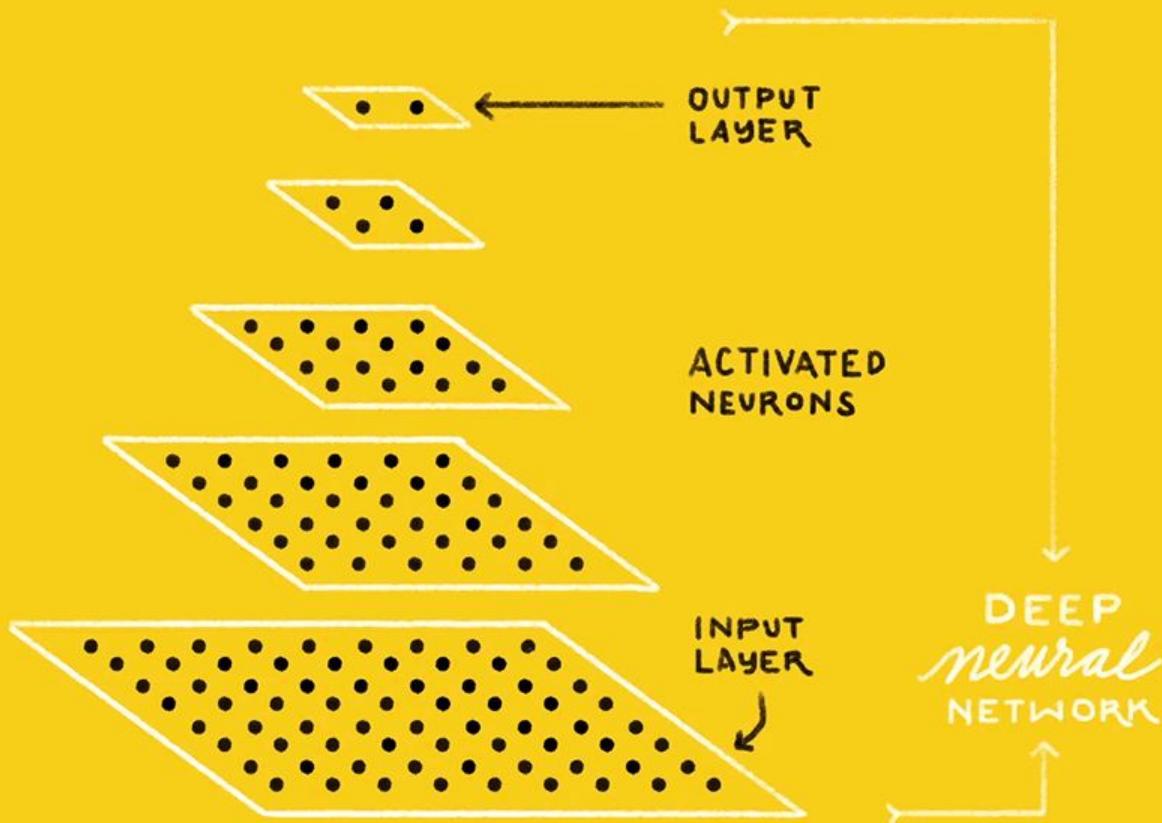
- I. Deep Learning in (almost) one slide
- II. Why should I care about Sequences?
- III. Recurrent Neural Networks Fundamentals
  - A. The chain rule
  - B. Vanishing gradients
  - C. LSTMs
- IV. Recurrent Networks as Sequence Decoders
  - A. Sequence-To-Sequence Learning & Applications
- V. Beyond RNNs

# Part I: Deep Learning Overview

IS THIS A  
**CAT or DOG?**



CAT   DOG



# A simple recipe

1. Lots of data

# A simple recipe

1. Lots of data
2. Deep, BIG DNNs/CNNs/RNNs

# A simple recipe

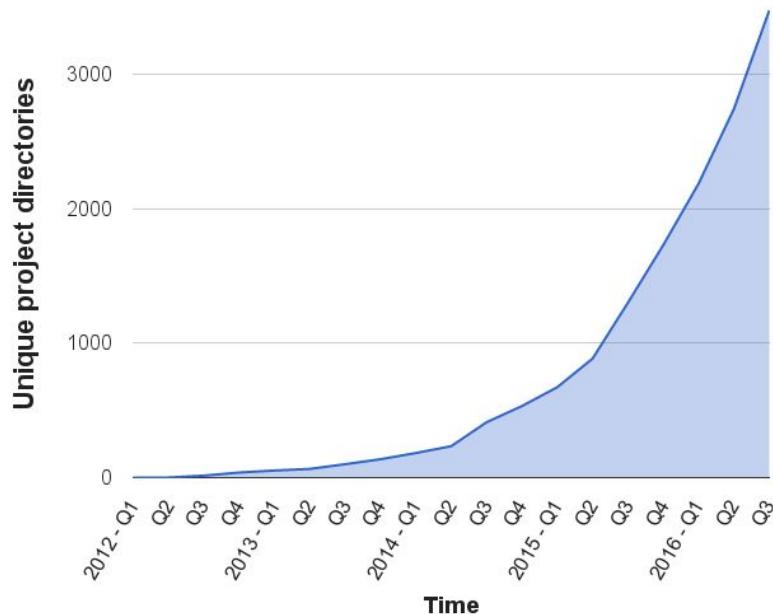
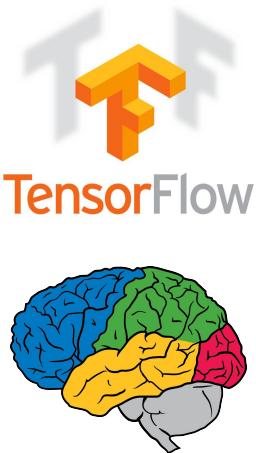
1. Lots of data
2. Deep, BIG DNNs/CNNs/RNNs
3. Lots of GPUs

# A simple recipe

1. Lots of data
2. Deep, BIG DNNs/CNNs/RNNs
3. Lots of GPUs
4. PROFIT!

# Growing Use of Deep Learning at Google

*"Google will soon be a big LSTM"*, Jürgen Schmidhuber

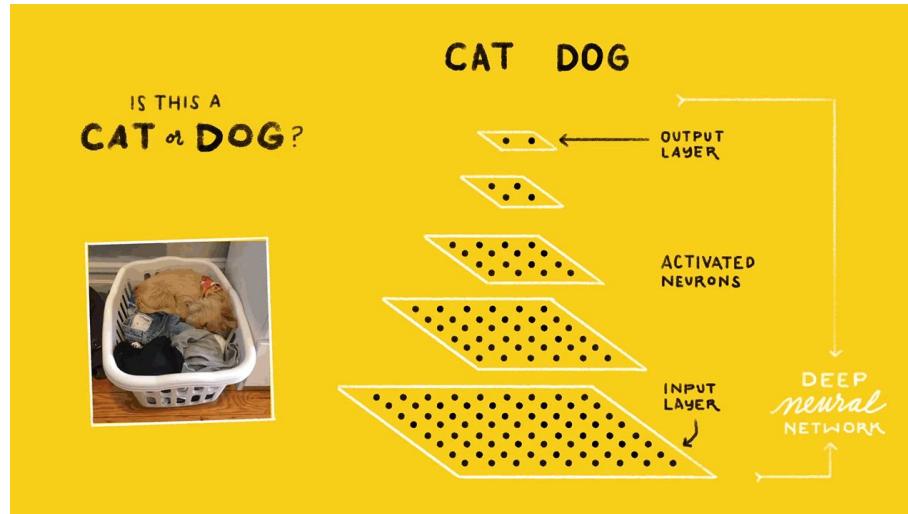


**Across many products/areas:**

- Android
- Apps
- drug discovery
- Gmail
- Image understanding
- Maps
- Photos
- Robotics research
- Search
- Speech
- Translate
- YouTube
- ... many others ...

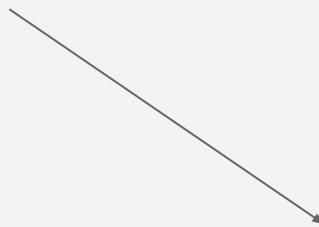
# Are we done yet?

- Neural Nets struggle with sequences
- Inputs and outputs must be fixed-sized vectors



# Part II: Why Sequences?

[Part, II:, Why, Sequences?]



## Part II: Why Sequences?

[Part, II:, Why, Sequences?]

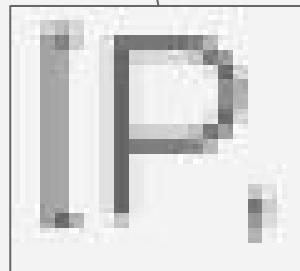
[P, a, r, t]

Part II: Why Sequences?

[Part, II:, Why, Sequences?]

[P, a, r, t]

Part II: Why Sequences?



Part II: Why Sequences?

[Part, II:, Why, Sequences?]

Part II: Why Sequences?

[Part, II:, Why, Sequences?]

[P, a, r, t]

Part II: Why Sequences?

[Part, II:, Why, Sequences?]

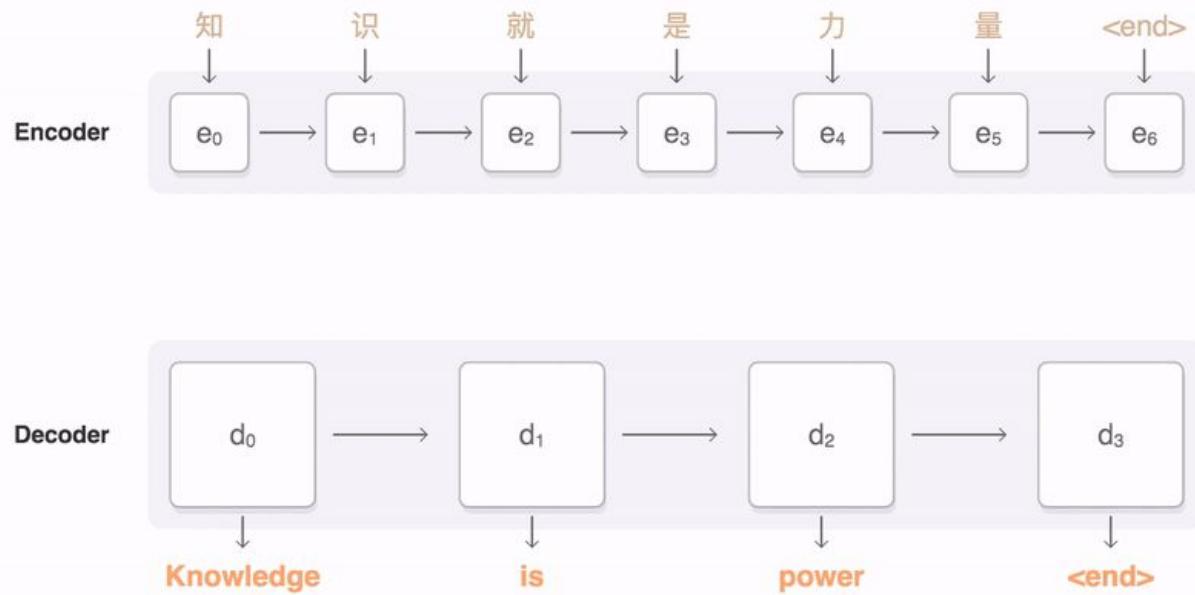
[P, a, r, t]

Part II: Why Sequences?

[P,

# Mix and Match: Google Neural MT

Wu et al, 2016 (Kalchbrenner et al, 2013; Sutskever et al, 2014;  
Cho et al, 2014; Bahdanau et al, 2014; ...)



# Math / Code [Karpathy, 2015]

*Proof.* Omitted.  $\square$

**Lemma 0.1.** Let  $\mathcal{C}$  be a set of the construction.

Let  $\mathcal{C}$  be a gerber covering. Let  $\mathcal{F}$  be a quasi-coherent sheaves of  $\mathcal{O}$ -modules. We have to show that

$$\mathcal{O}_{\mathcal{O}_X} = \mathcal{O}_X(\mathcal{L})$$

*Proof.* This is an algebraic space with the composition of sheaves  $\mathcal{F}$  on  $X_{\text{étale}}$  we have

$$\mathcal{O}_X(\mathcal{F}) = \{\text{morph}_1 \times_{\mathcal{O}_X} (\mathcal{G}, \mathcal{F})\}$$

where  $\mathcal{G}$  defines an isomorphism  $\mathcal{F} \rightarrow \mathcal{F}$  of  $\mathcal{O}$ -modules.  $\square$

**Lemma 0.2.** This is an integer  $\mathcal{Z}$  is injective.

*Proof.* See Spaces, Lemma ??.

**Lemma 0.3.** Let  $S$  be a scheme. Let  $X$  be a scheme and  $X$  is an affine open covering. Let  $\mathcal{U} \subset \mathcal{X}$  be a canonical and locally of finite type. Let  $X$  be a scheme. Let  $X$  be a scheme which is equal to the formal complex.

The following to the construction of the lemma follows.

Let  $X$  be a scheme. Let  $X$  be a scheme covering. Let

$$b : X \rightarrow Y' \rightarrow Y \rightarrow Y \rightarrow Y' \times_X Y \rightarrow X.$$

be a morphism of algebraic spaces over  $S$  and  $Y$ .

*Proof.* Let  $X$  be a nonzero scheme of  $X$ . Let  $X$  be an algebraic space. Let  $\mathcal{F}$  be a quasi-coherent sheaf of  $\mathcal{O}_X$ -modules. The following are equivalent

- (1)  $\mathcal{F}$  is an algebraic space over  $S$ .
- (2) If  $X$  is an affine open covering.

Consider a common structure on  $X$  and  $X$  the functor  $\mathcal{O}_X(U)$  which is locally of finite type.  $\square$

```
static void num_serial_settings(struct tty_struct *tty)
{
    if (tty == tty)
        disable_single_st_p(dev);
    pci_disable_spool(port);
    return 0;
}

static void do_command(struct seq_file *m, void *v)
{
    int column = 32 << (cmd[2] & 0x80);
    if (state)
        cmd = (int)(int_state ^ (in_8(&ch->ch_flags) & Cmd) ? 2 : 1);
    else
        seq = 1;
    for (i = 0; i < 16; i++) {
        if (k & (1 << 1))
            pipe = (in_use & UMXTHREAD_UNCCA) +
                ((count & 0x00000000fffff8) & 0x0000000f) << 8;
        if (count == 0)
            sub(pid, ppc_md.kexec_handle, 0x20000000);
        pipe_set_bytes(i, 0);
    }
    /* Free our user pages pointer to place camera if all dash */
    subsystem_info = &of_changes[PAGE_SIZE];
    rek_controls(offset, idx, &offset);
    /* Now we want to deliberately put it to device */
    control_check_polarity(&context, val, 0);
    for (i = 0; i < COUNTER; i++)
        seq_puts(s, "policy ");
}
```

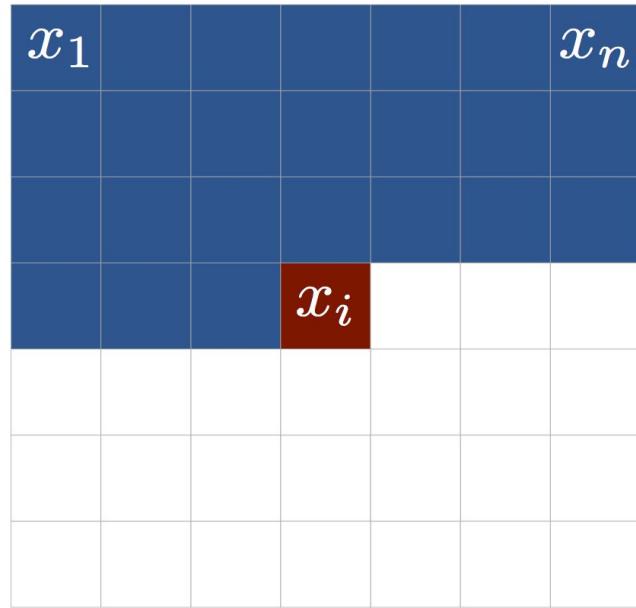


*Human:* A young girl asleep on the sofa cuddling a stuffed bear.

*NIC:* A close up of a child holding a stuffed animal.

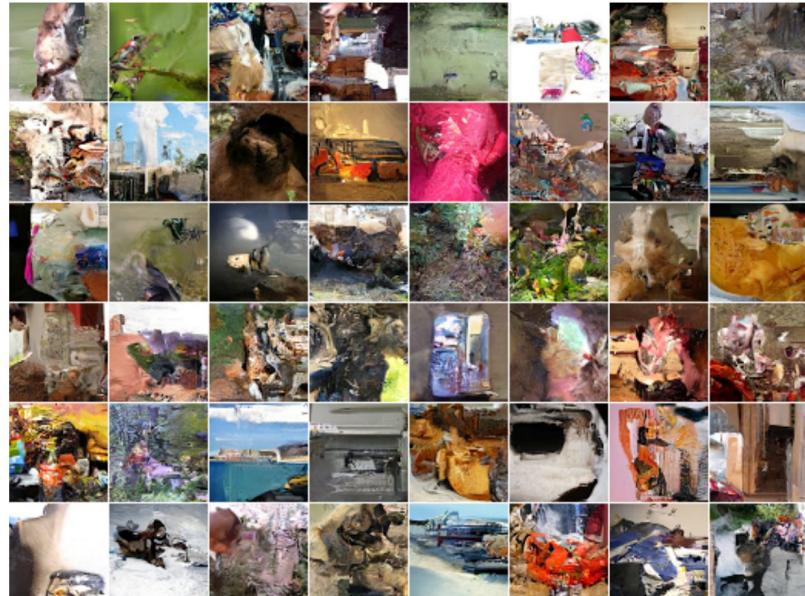
*NIC:* A baby is asleep next to a teddy bear.

# Pixels as a sequence



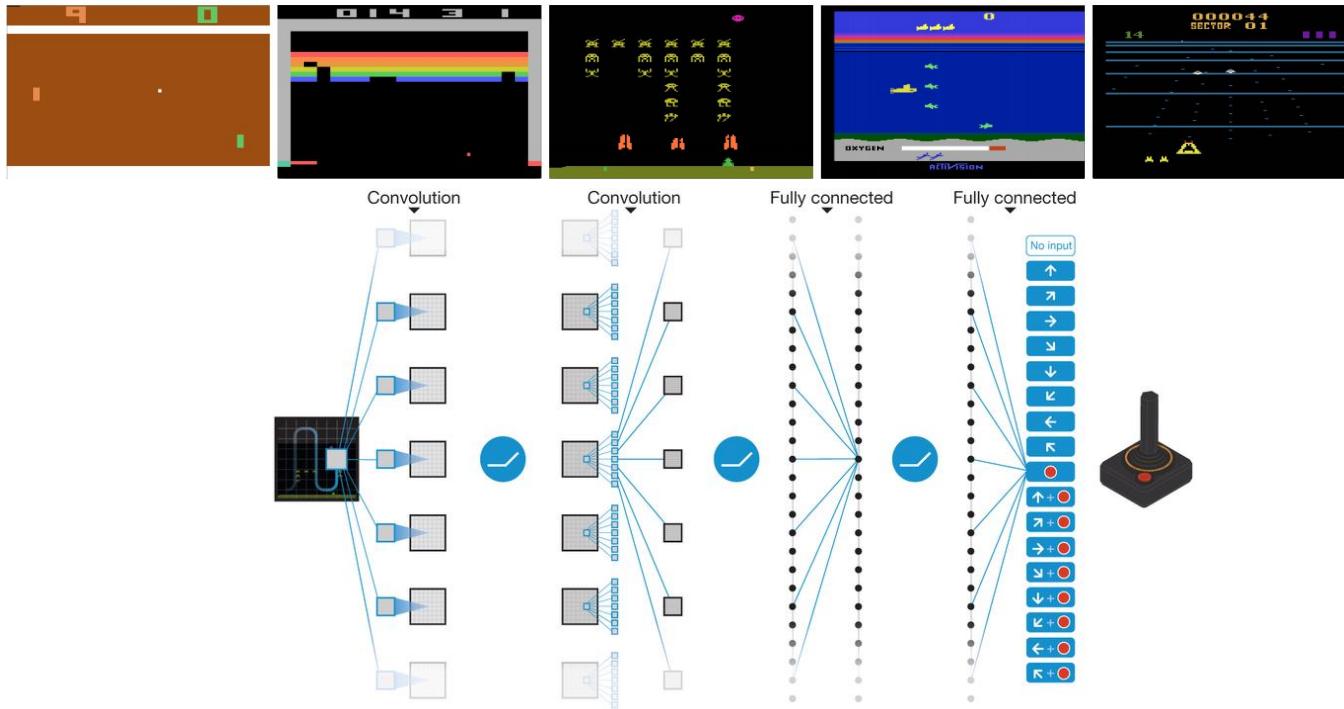
$$p(\mathbf{x}) = \prod_{i=1}^{n^2} p(x_i | x_1, \dots, x_{i-1})$$

The Chain Rule (more on it later)



# Deep Reinforcement Learning

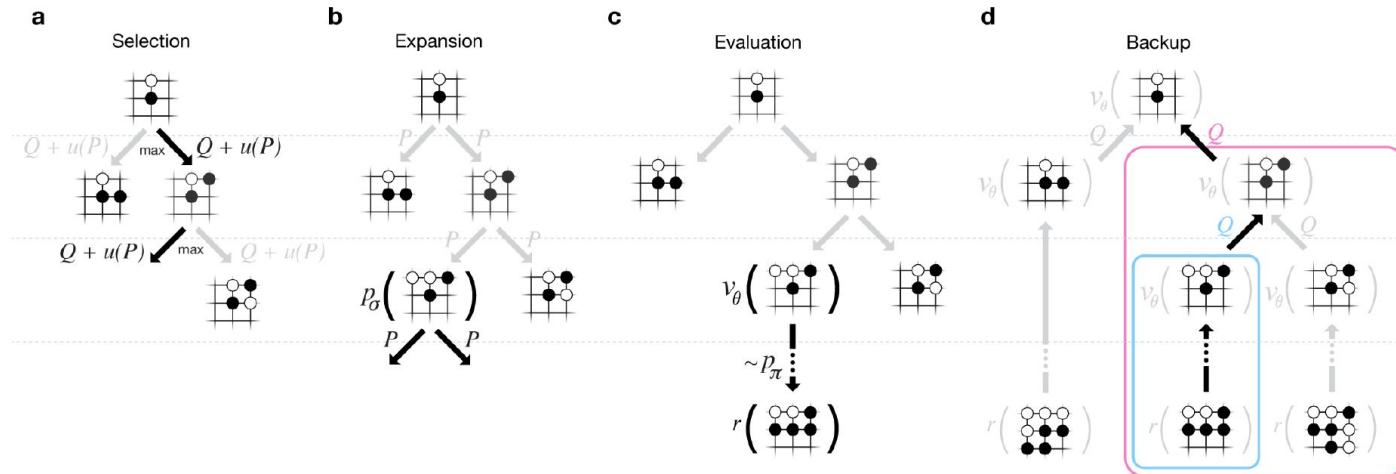
**Human-level control through deep reinforcement learning,**  
Mnih et al, Nature 2015



# Deep Reinforcement Learning

**Mastering the game of Go with deep neural networks and tree search.**

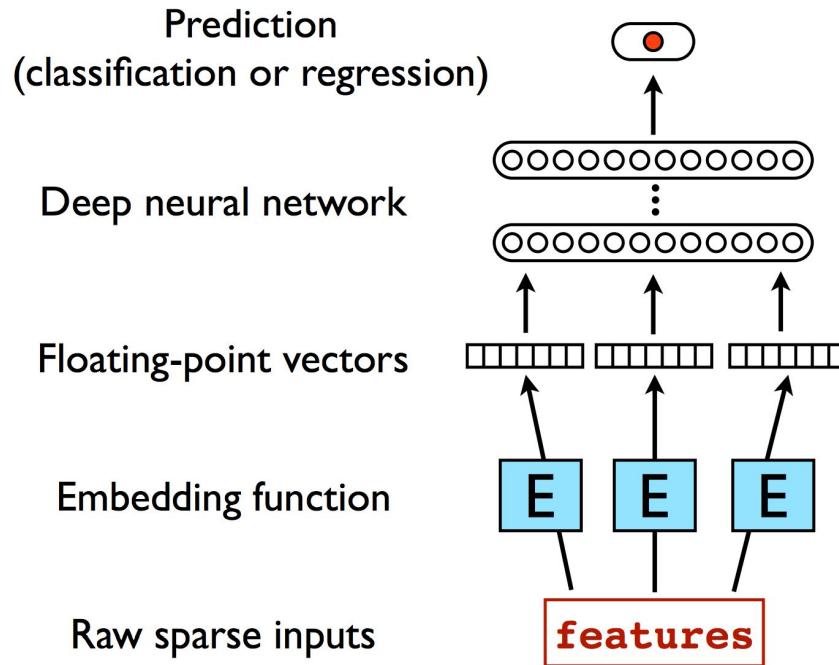
Silver et al, Nature 2016



# Part III: Recurrent Neural Networks

## Fundamentals

# Word Vectors Aside



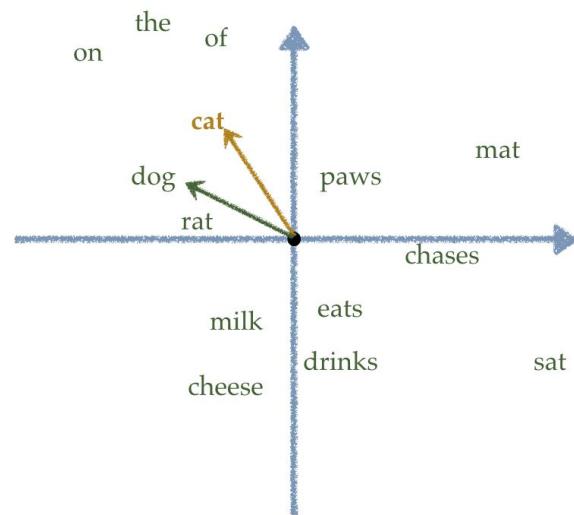
# Embedding == Sparse matrix multiplication

Let  $v = [0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 1 \ 0 \ 0 \ 0 \ 0 \dots \ 0 \ 0 \ 0 \ 0]$

each position is a word ( $\text{length}(v) == \text{vocab size}$ )

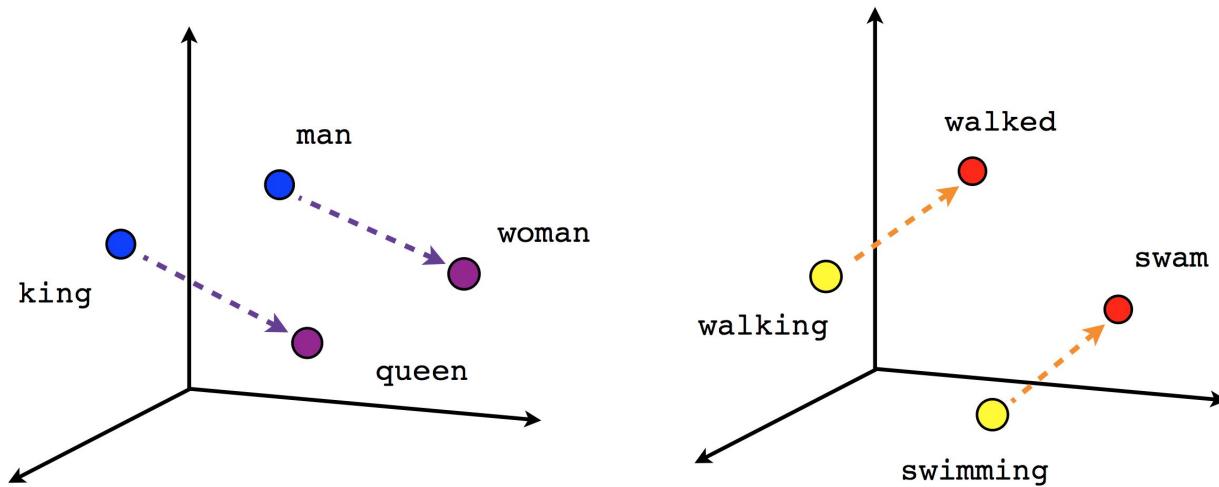
To embed a word we simply do:

$$e = \mathbf{W}v$$



where  $\mathbf{W}$  is a  $D \times \text{vocab size}$  matrix (one  $D$  dim. vector for each word)

# Vector arithmetic



Solve analogies with vector arithmetic!

$$V(\text{queen}) - V(\text{king}) \approx V(\text{woman}) - V(\text{man})$$

$$V(\text{queen}) \approx V(\text{king}) + (V(\text{woman}) - V(\text{man}))$$

# Part III.A The Chain Rule

# Basics

Supervised Learning: many input / output examples ( $\mathbf{x}$ ,  $\mathbf{y}$ )

Find mapping  $f$  s.t.

$$\mathbf{y} \approx f(\mathbf{x})$$

Define

$$\mathcal{L}(\boldsymbol{\theta}) = \sum_{i=1}^N \ell(f_{\boldsymbol{\theta}}, \{\mathbf{x}_i, \mathbf{y}_i\})$$

and learn by optimizing

$$\boldsymbol{\theta} = \arg \min_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta})$$

# Motivating Example: Language

KING LEAR:

O, if you were a feeble sight, the courtesy of your law,  
Your sight and several breath, will wear the gods  
With his heads, and my hands are wonder'd at the deeds,  
So drop upon your lordship's head, and your opinion  
Shall be against your honour.

Who wrote these lines?

- A. William Shakespeare
- B. William Shakespeare's ghostwriter
- C. Ben Johnson
- D. Molière (translation)
- E. Andrej Karpathy's recurrent neural network

# n-grams

$$P(w_t | w_{t-1}, \dots, w_{t-n+1})$$

what to cook with broccoli and \_

what to cook with broccoli and **beef**

what to cook with broccoli and **butter**

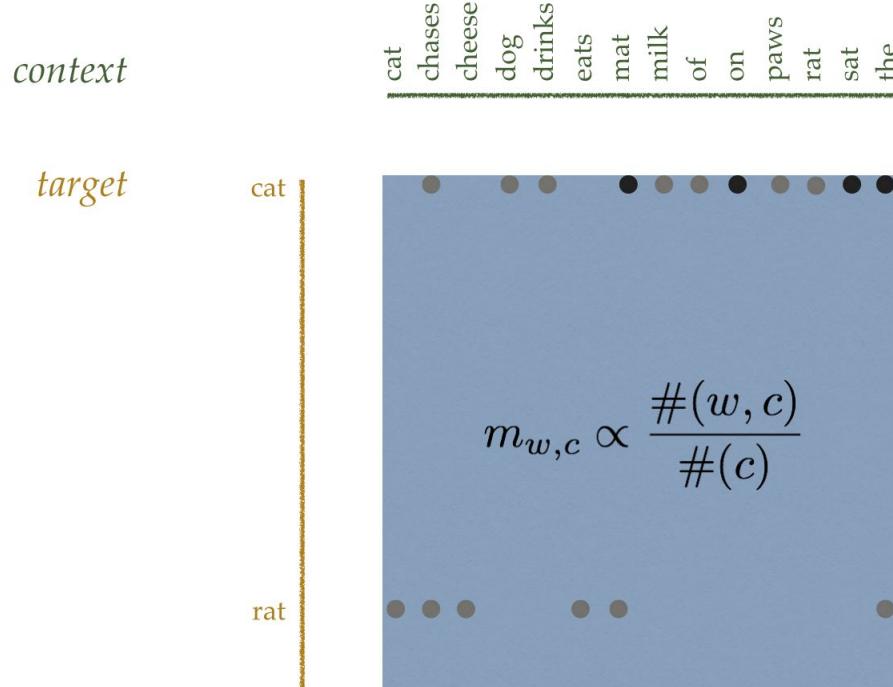
what to cook with broccoli and **blenders**

what to cook with broccoli and **boomboxes**

# Modelling “P”

<i>context</i>						<i>target</i>	$P(w_t   w_{t-1}, w_{t-2}, \dots w_{t-5})$
the	cat	sat	on	the	mat		0.15
$w_{t-5}$	$w_{t-4}$	$w_{t-3}$	$w_{t-2}$	$w_{t-1}$	$w_t$		
the	cat	sat	on	the	rug		0.12
the	cat	sat	on	the	hat		0.09
the	cat	sat	on	the	dog		0.01
the	cat	sat	on	the	the		0
the	cat	sat	on	the	sat		0
the	cat	sat	on	the	robot		?
the	cat	sat	on	the	printer		?

# 2-grams



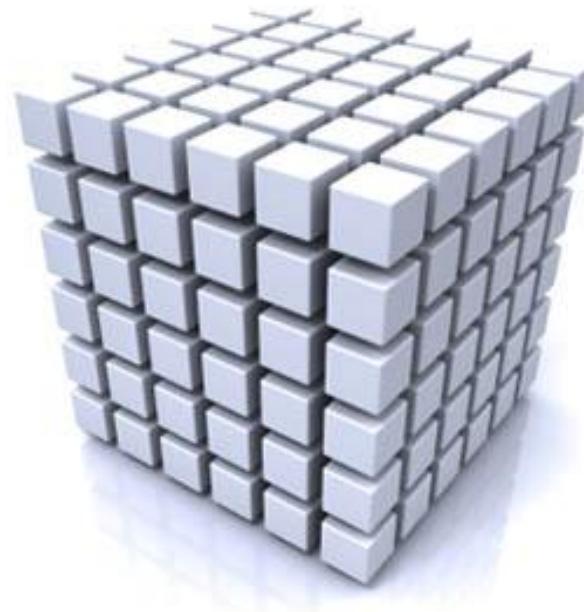
the **cat** sat on the mat  
the **cat** drinks milk  
the dog chases the **cat**  
the paws of the **cat**

the cat chases the **rat**  
the **rat** eats cheese  
the **rat** eats the mat

# n-grams

Vocabulary: 10K

Table size:  $10000^n$



# n-grams

$$P(w_1, w_2, \dots, w_{T-1}, w_T) \approx \prod_{t=1}^T P(w_t | w_{t-1}, \dots, w_{t-n+1})$$

the	cat	sat	on	the	mat	$P(w_1)$
the	<b>cat</b>	sat	on	the	mat	$P(w_2   w_1)$
the	cat	<b>sat</b>	on	the	mat	$P(w_3   w_2, w_1)$
the	cat	sat	<b>on</b>	the	mat	$P(w_4   w_3, w_2)$
the	cat	sat	on	<b>the</b>	mat	$P(w_5   w_4, w_3)$
the	cat	sat	on	the	<b>mat</b>	$P(w_6   w_5, w_4)$

# The Chain Rule

$$P(w_1, w_2, \dots, w_{T-1}, w_T) = \prod_{t=1}^T P(w_t | w_{t-1}, w_{t-2}, \dots, w_1)$$

the	cat	sat	on	the	mat	$P(w_1)$
the	<b>cat</b>	sat	on	the	mat	$P(w_2   w_1)$
the	cat	<b>sat</b>	on	the	mat	$P(w_3   w_2, w_1)$
the	cat	sat	<b>on</b>	the	mat	$P(w_4   w_3, w_2, w_1)$
the	cat	sat	on	<b>the</b>	mat	$P(w_5   w_4, w_3, w_2, w_1)$
the	cat	sat	on	the	<b>mat</b>	$P(w_6   w_5, w_4, w_3, w_2, w_1)$

# Proof

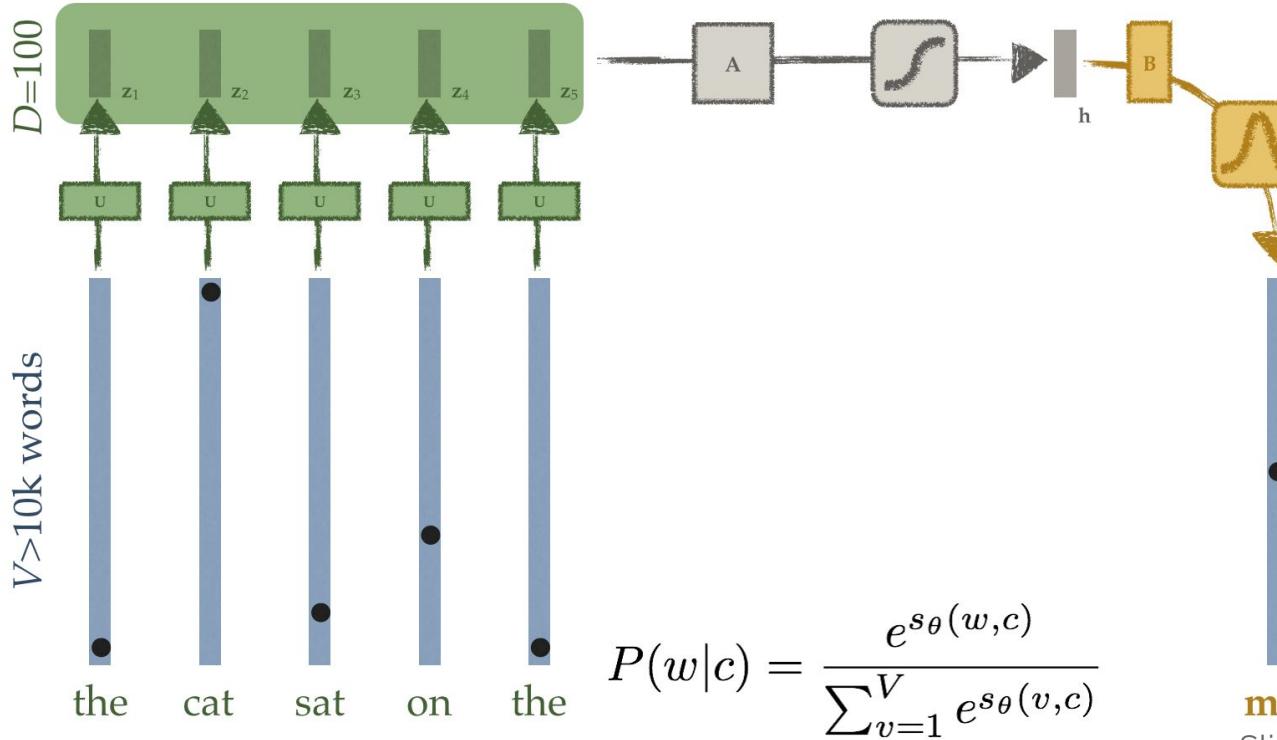
$$\sum_{w_1, \dots, w_T} p(w_1, \dots, w_T) = 1$$

$$\begin{aligned} & \sum_{w_1, \dots, w_T} \prod_{t=1}^T p(w_t | w_1, \dots, w_{t-1}) = \\ & \sum_{w_1, \dots, w_T} p(w_T | w_1, \dots, w_{T-1}) \prod_{t=1}^{T-1} p(w_t | w_1, \dots, w_{t-1}) = \\ & \sum_{w_1, \dots, w_{T-1}} p(w_{T-1} | w_1, \dots, w_{T-2}) \prod_{t=1}^{T-2} p(w_t | w_1, \dots, w_{t-1}) = \\ & \quad \dots \\ & \sum_{w_1, \dots, w_T} \prod_{t=1}^T p(w_t | w_1, \dots, w_{t-1}) = 1 \end{aligned}$$

# A Key Insight: vectorizing context

[Yoshua Bengio et al. (2001, 2003), "A Neural Probabilistic Language Model", *JMLR*;  
Andriy Mnih and Geoff Hinton, "Three new graphical models for statistical language modeling", *ICML*]

$$p(w_t | w_1, \dots, w_{t-1}) = p_\theta(w_t | f_\theta(w_1, \dots, w_{t-1}))$$

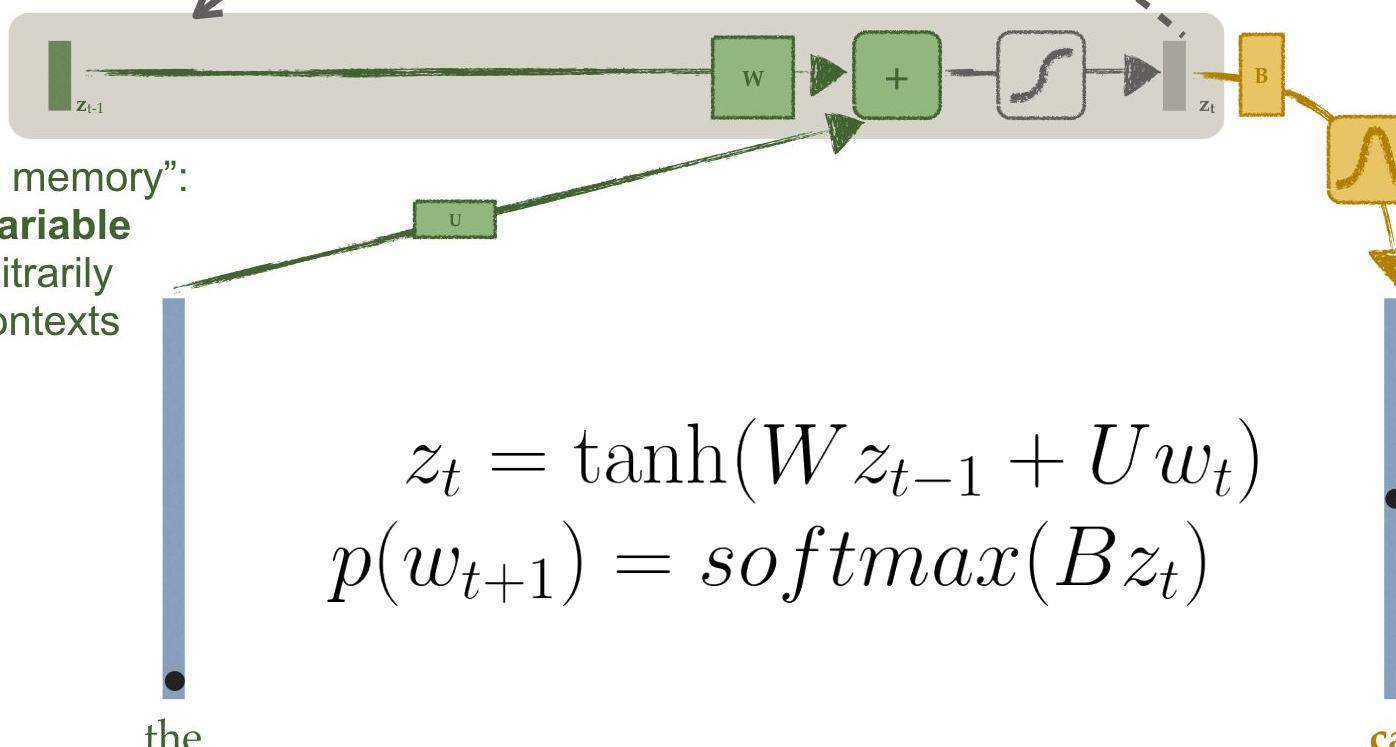


mat

Slide Credit: Piotr Mirowski

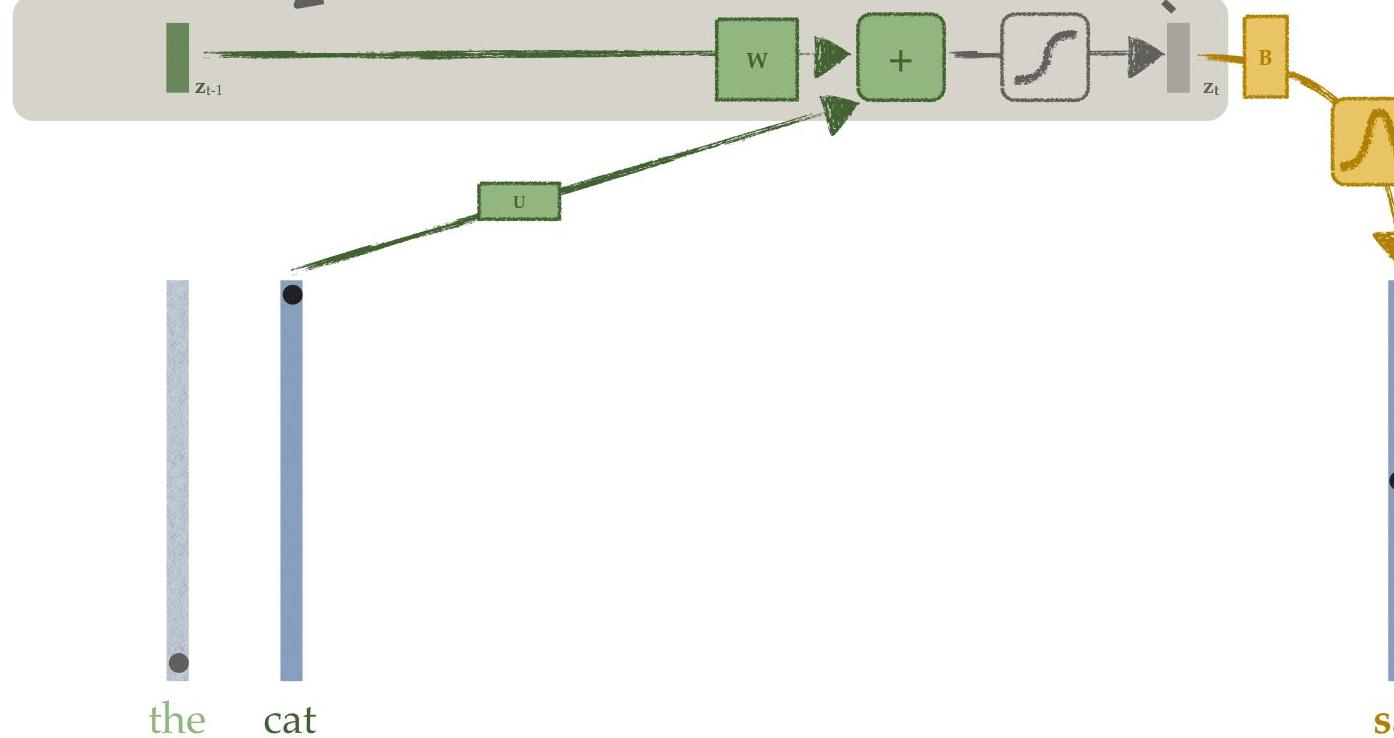
# Recurrent Neural Network Language Models

[Jeffrey L Elman (1991) "Distributed representations, simple recurrent networks and grammatical structure", *Machine Learning*; Tomas Mikolov et al. (2010) "Recurrent neural network based language model", *INTERSPEECH*]



$$\begin{aligned} z_t &= \tanh(Wz_{t-1} + Uw_t) \\ p(w_{t+1}) &= softmax(Bz_t) \end{aligned}$$

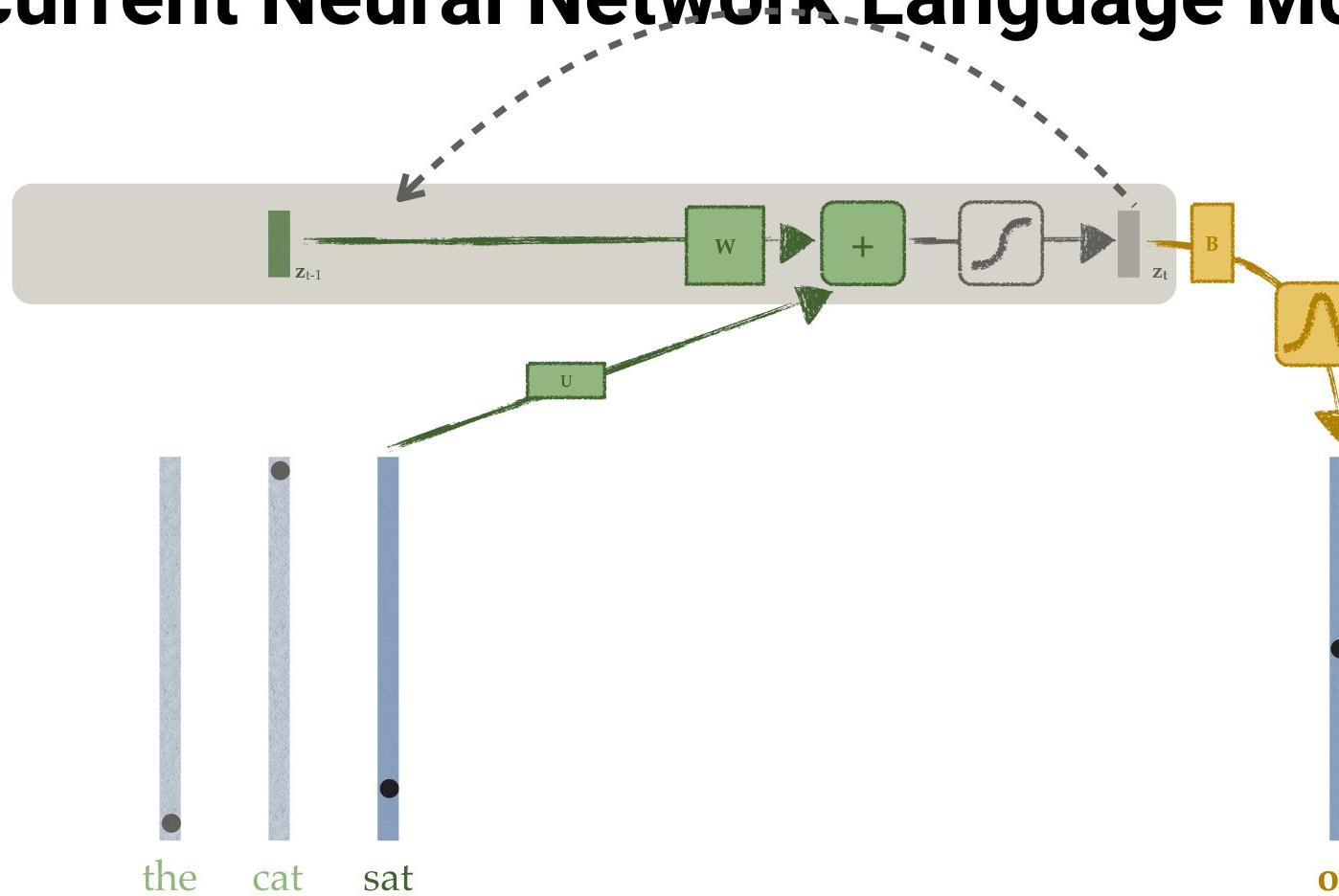
# Recurrent Neural Network Language Models



**sat**

Slide Credit: Piotr Mirowski

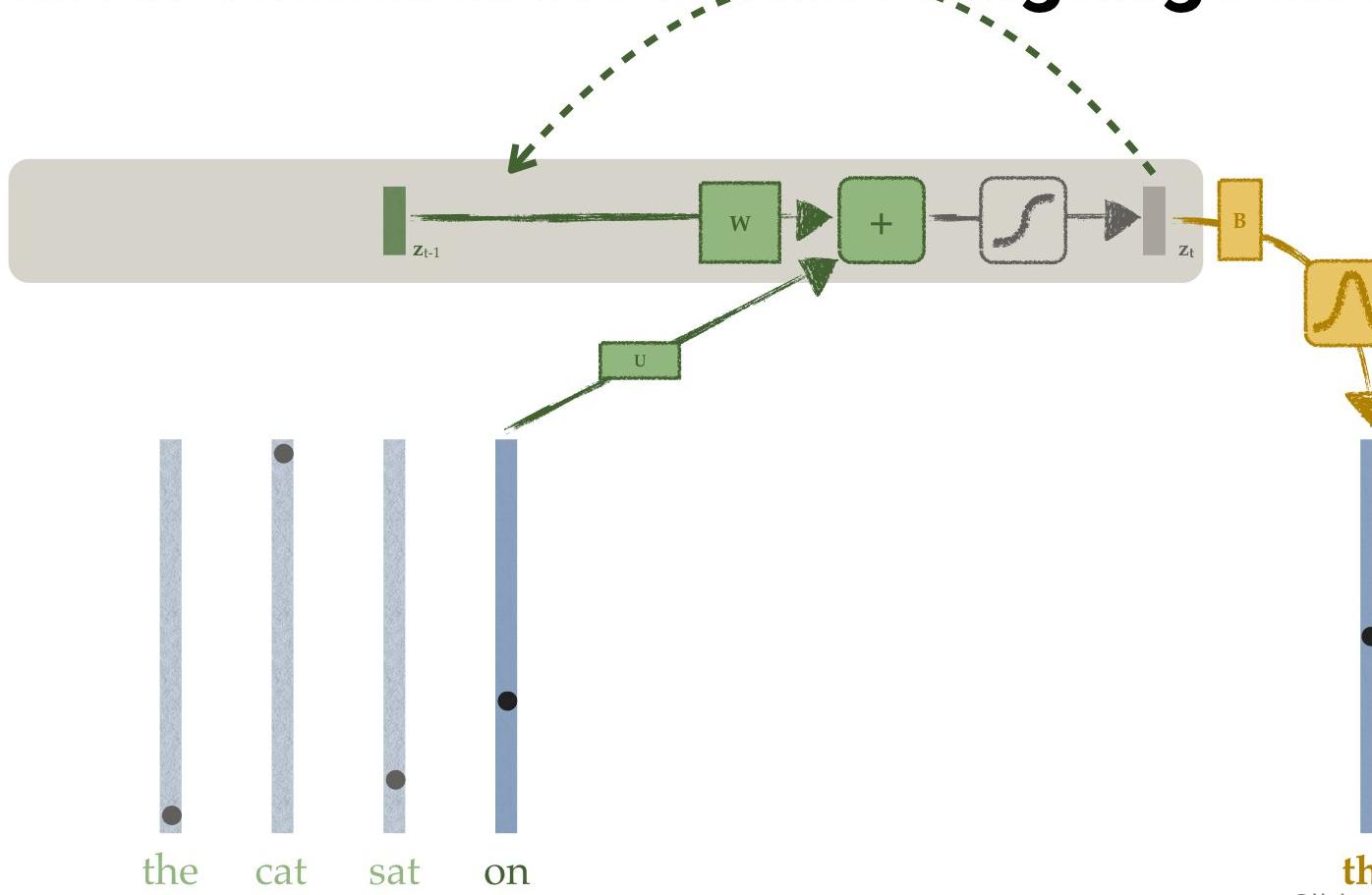
# Recurrent Neural Network Language Models



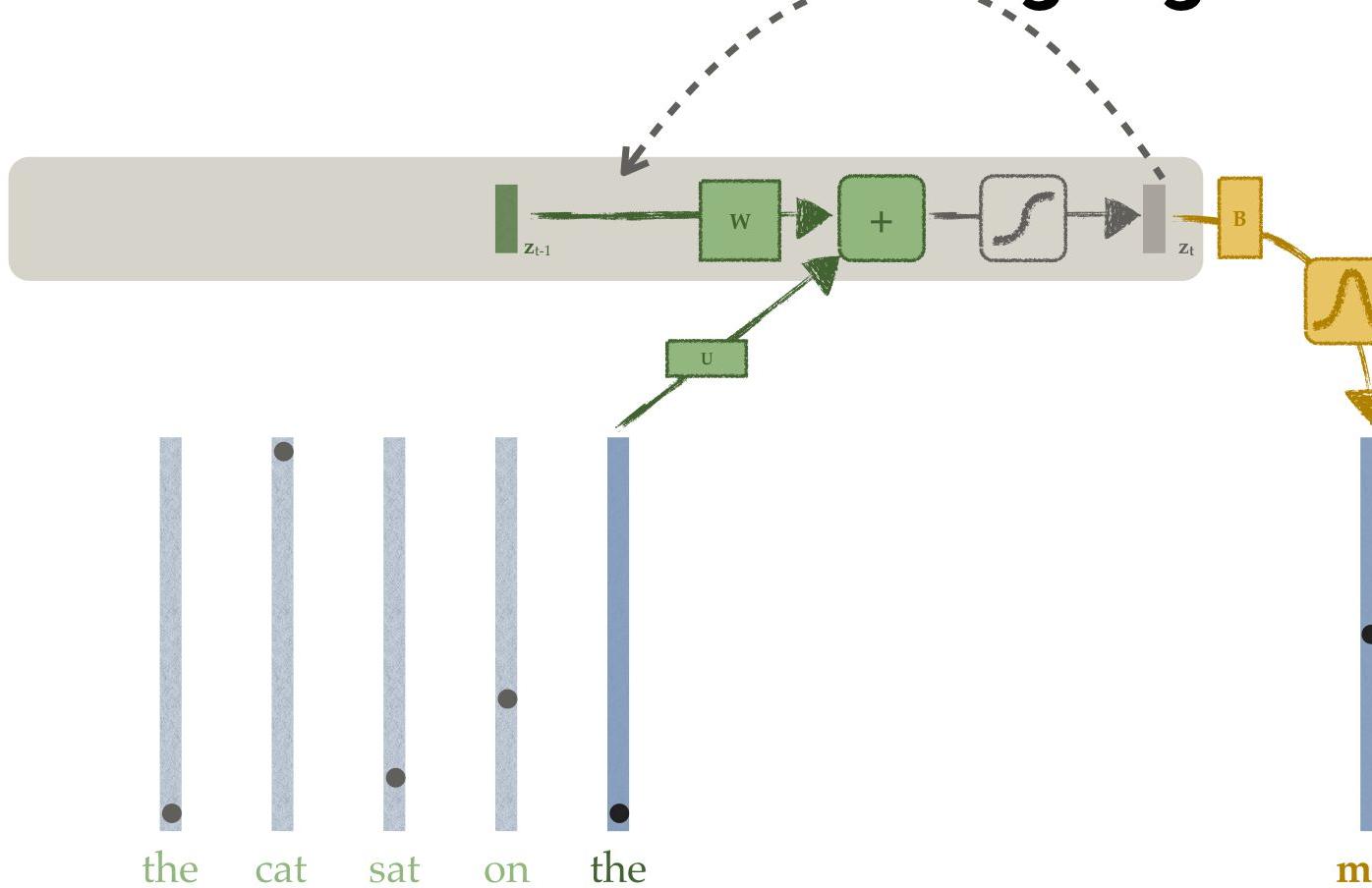
on

Slide Credit: Piotr Mirowski

# Recurrent Neural Network Language Models



# Recurrent Neural Network Language Models



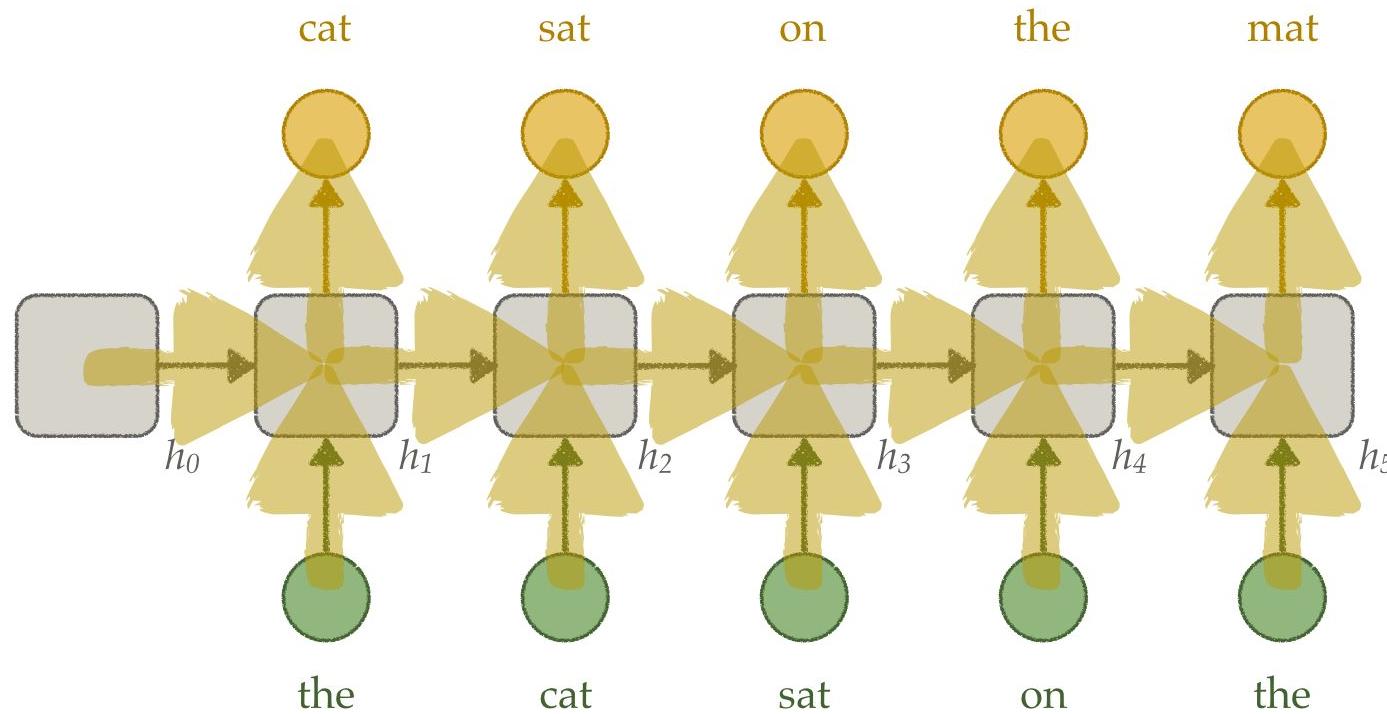
mat

Slide Credit: Piotr Mirowski

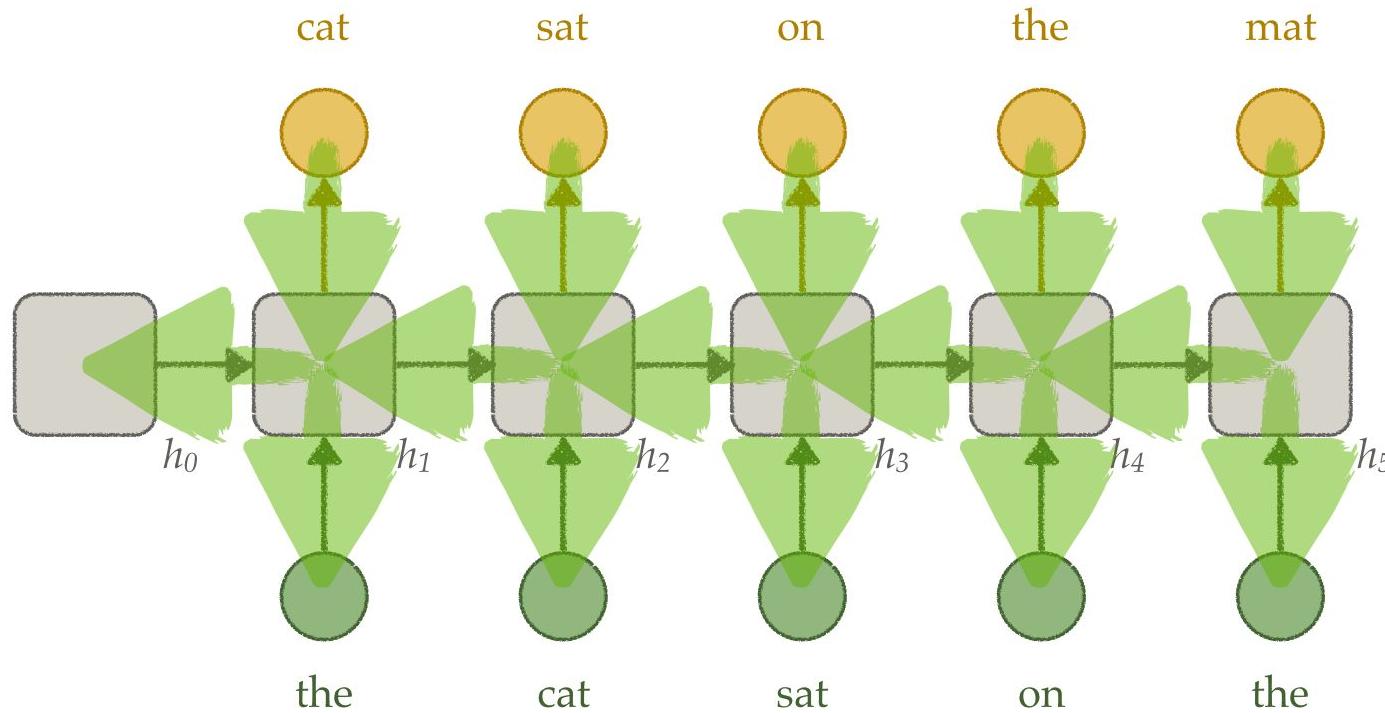
# What do we Optimize?

$$\theta^* = \arg \max_{\theta} E_{w \sim data} \log P_{\theta}(w_1, \dots, w_T)$$

# Recurrent Neural Network Language Models



# Recurrent Neural Network Language Models



# Deriving BPTT

$$z_t = \tanh(W z_{t-1} + U w_t)$$

$$p(w_{t+1}) = \text{softmax}(B z_t)$$

$$L(W, U, B) = - \sum_{t=1}^T \log p_{W,U,B}(w_t | w_1, \dots, w_{t-1})$$

$$L_t = - \log p(w_t | w_1, \dots, w_{t-1})$$

## Deriving BPTT - B (softmax)

$$L_t = -\log p(w_t | w_1, \dots, w_{t-1})$$

$$\frac{\delta L_t}{\delta B} = \frac{\delta L_t}{\delta p_t} \frac{\delta p_t}{\delta l_t} \frac{\delta l_t}{\delta B}$$

$$\frac{\delta L_t}{\delta B} = (p_t - \text{label}_t) z_t^T$$

$$\frac{\delta L}{\delta B} = \sum_{t=1}^T (p_t - \text{label}_t) z_t^T$$

# Deriving BPTT - W (recurrent)

$$z_t = \tanh(W z_{t-1} + U w_t)$$
$$p(w_{t+1}) = \text{softmax}(B z_t)$$

$$\frac{\delta L_t}{\delta W} = \frac{\delta L_t}{\delta p_t} \frac{\delta p_t}{\delta z_t} \frac{\delta z_t}{\delta W}$$

$$\frac{\delta L_t}{\delta W} = \sum_{k=1}^t \frac{\delta L_t}{\delta p_t} \frac{\delta p_t}{\delta z_t} \frac{\delta z_t}{\delta z_k} \frac{\delta z_k}{\delta W}$$

# Part III.B Vanishing Gradients

# Vanishing (& exploding) Gradients

[Sepp Hochreiter (1991) "Untersuchungen zu dynamischen neuronalen Netzen", *Diploma TUM*;

Yoshua Bengio et al. (1994) "Learning Long-Term Dependencies with Gradient Descent is Difficult", *IEEE Transactions on Neural Networks*]

- n-grams have VERY strong vanishing gradients : )
- RNNs *can* encode long term, but *learning* is hard
- Does long-term ALWAYS matter?
  - Language has strong correlations short term (but also long!)
  - Seq2seq models *must* remember (later)
  - Learning programs

# A simple example

$$h_t = w * h_{t-1}$$

$$h_t = w^t * h_o$$

$h_t \rightarrow \infty$  if  $|w| > 1$

$h_t \rightarrow 0$  if  $|w| < 1$

# A simple example

$$h_t = w * h_{t-1}$$

$$h_t = w^t * h_o$$

$$h_t \rightarrow \infty \text{ if } |w| > 1$$

$$h_t \rightarrow 0 \text{ if } |w| < 1$$

$$h_t = \tanh(w * h_{t-1})$$

$$h_t \rightarrow \text{bounded for all } w : )$$

$$\delta h_t / \delta h_{t-T} = (1 - \tanh^2(w * h_{t-1})) * w * \delta h_{t-1} / \delta h_{t-T}$$

# A simple example

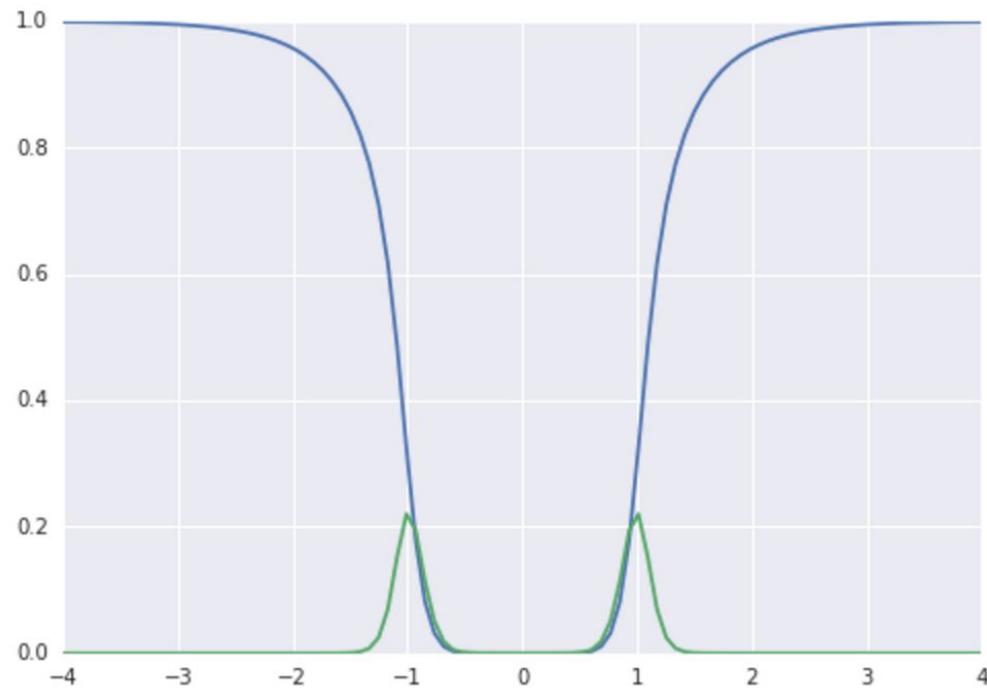
```
import numpy as np
from matplotlib import pyplot

def forward_backward_prop(w, T):
    hs = [0.5]
    for _ in range(T):
        hs.append(np.tanh(w*hs[-1]))
    dh = 1
    for t in range(T):
        dh = (1-hs[-1-t]*hs[-1-t])*w*dh
    return hs[-1], dh

T = 10
wlim = 4

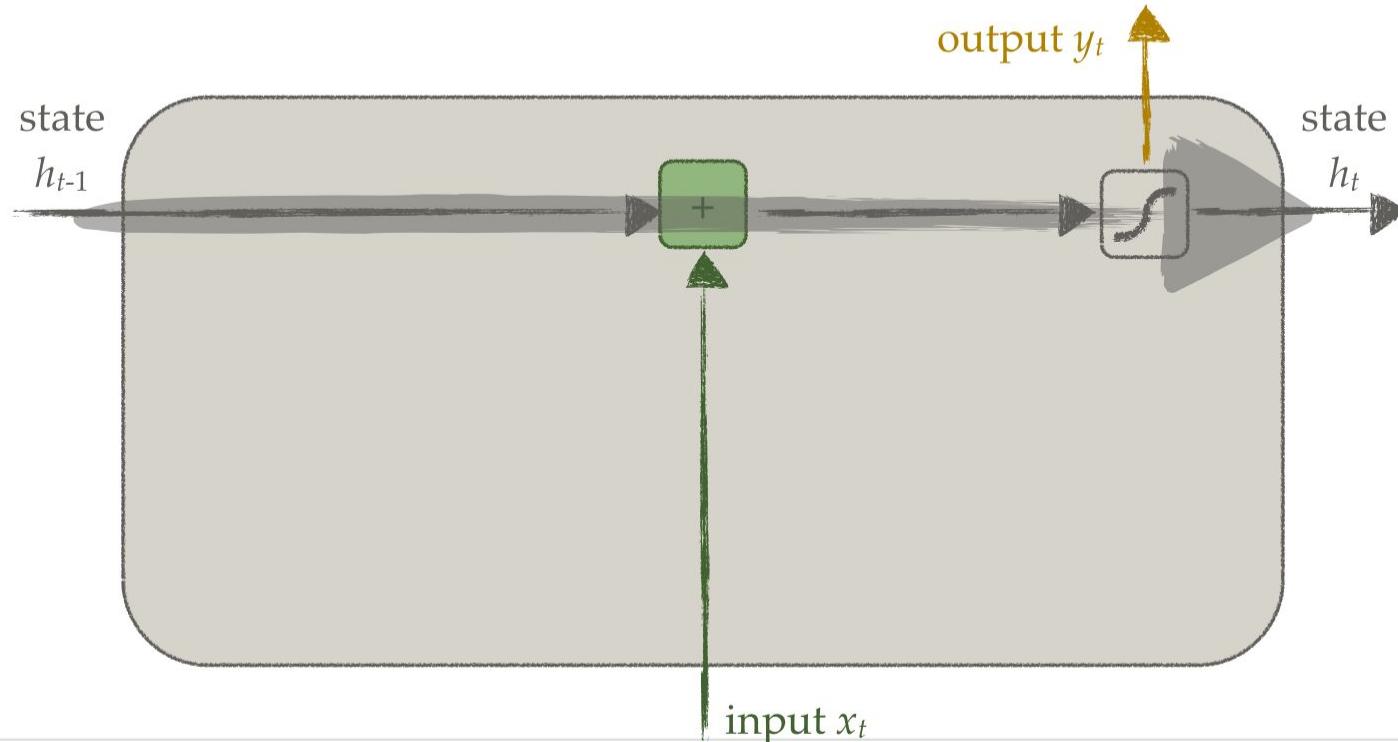
ws = np.linspace(-wlim, wlim, 100)
res = []
for w in ws:
    res.append(forward_backward_prop(w,T))

pyplot.plot(ws, [r[0] for r in res])
pyplot.plot(ws, [r[1] for r in res])
```



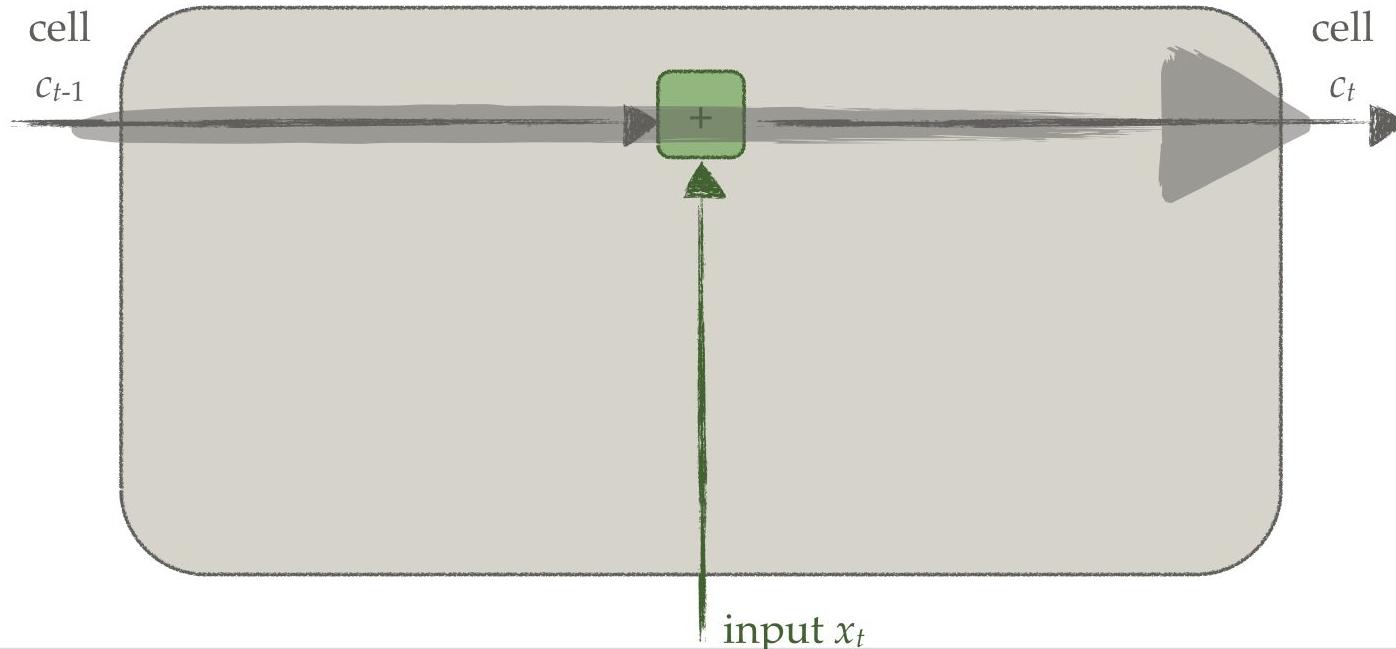
## Part III.C LSTMs

# LSTMs



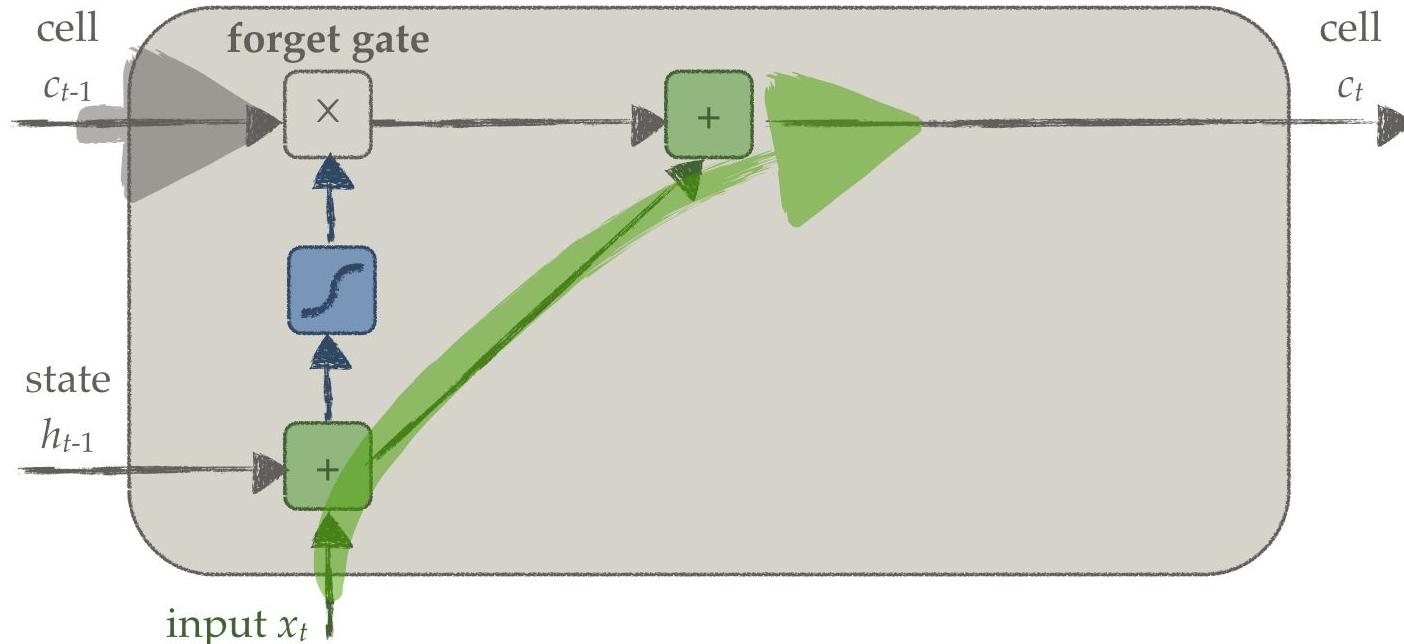
# Requirement #1: linear cell

[Sepp Hochreiter and Jürgen Schmidhuber (1997) “Long Short-Term Memory”, *Neural Computation*;  
Alex Graves (2013a) “Generating sequences with recurrent neural networks”, arXiv 1308.0850]

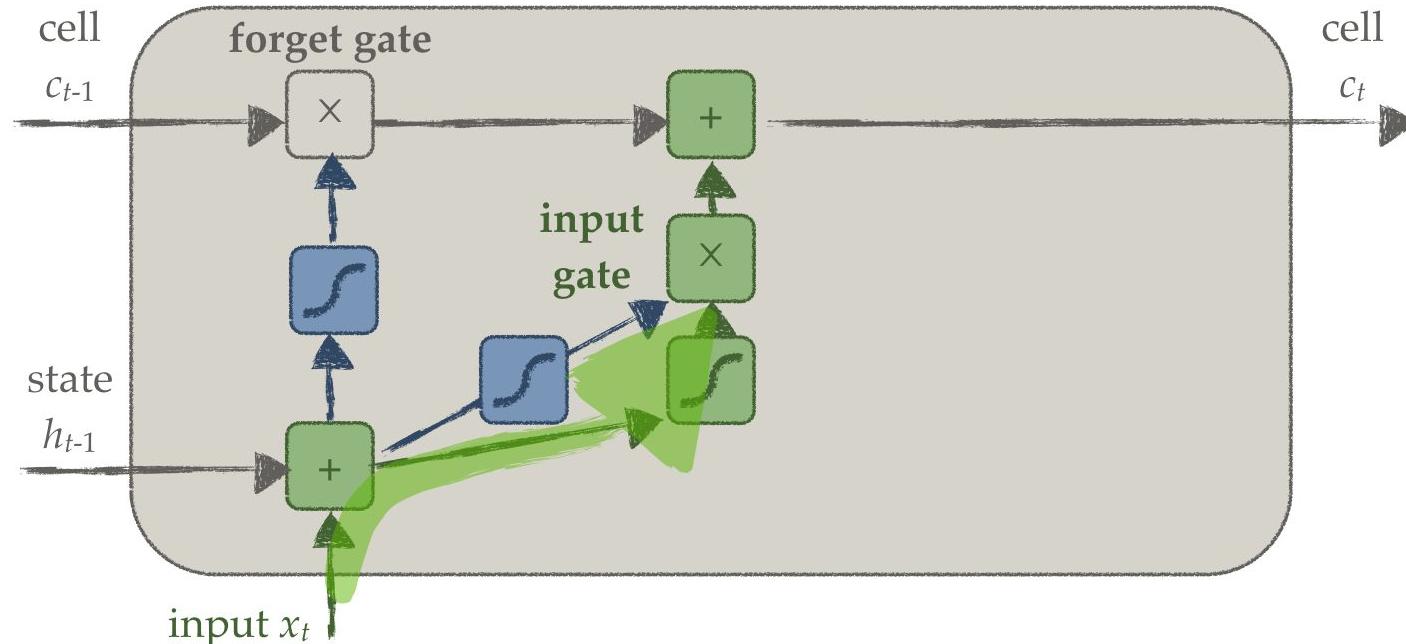


# Requirement #2: forget information

[Sepp Hochreiter and Jürgen Schmidhuber (1997) “Long Short-Term Memory”, *Neural Computation*; Alex Graves (2013a) “Generating sequences with recurrent neural networks”, arXiv 1308.0850]



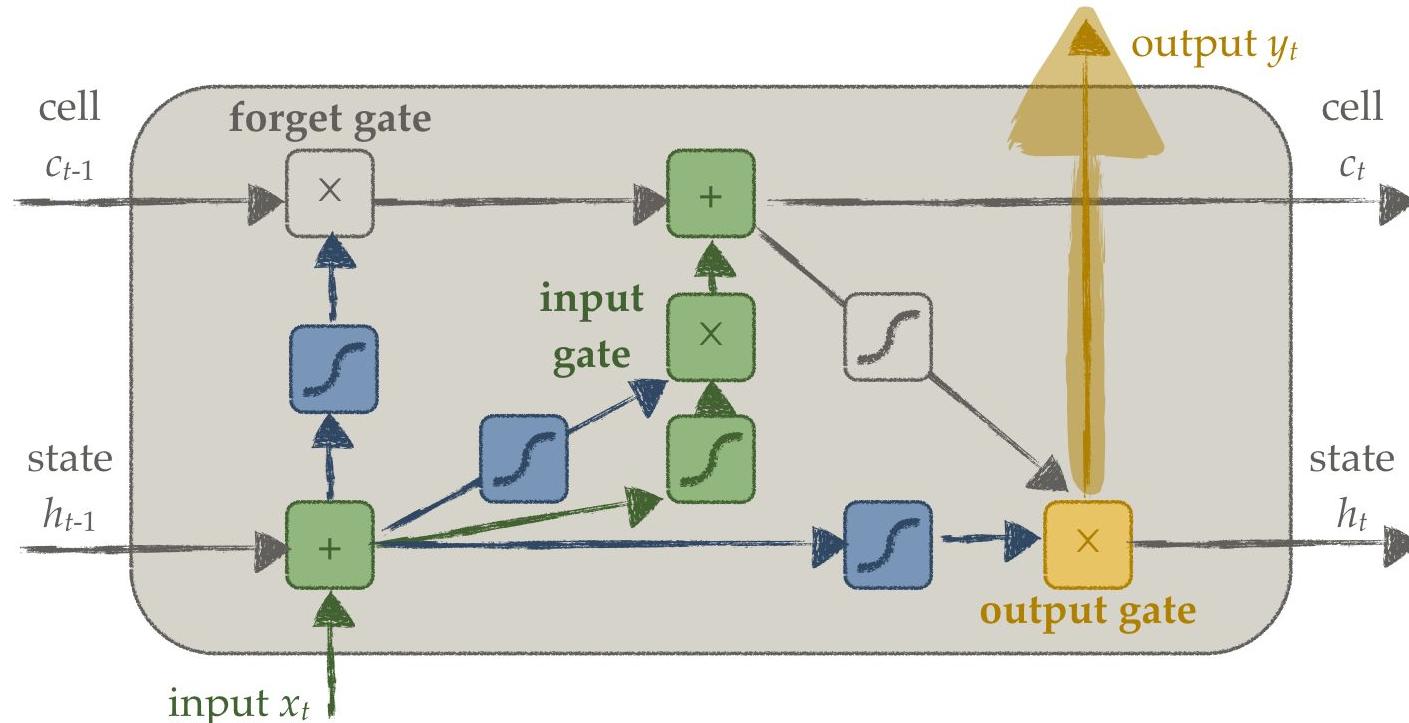
# Requirement #3: ignore inputs



[Hochreiter and Schmidhuber (1997); Graves (2013a)]

Learning Sequences – Piotr Mirowski

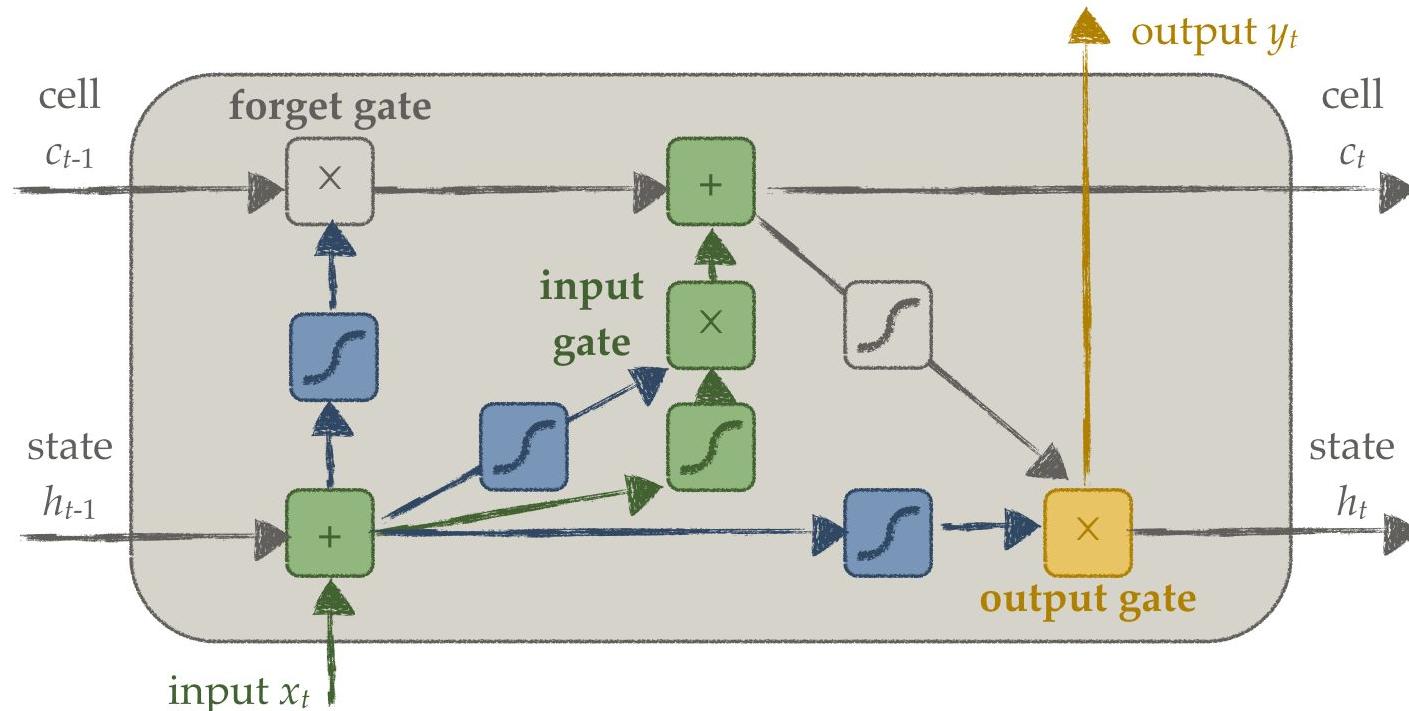
# Requirement #4: control outputs



[Hochreiter and Schmidhuber (1997); Graves (2013a)]

Learning Sequences – Piotr Mirowski

# Long Short-Term Memory (LSTM)



[Hochreiter and Schmidhuber (1997); Graves (2013a)]

Learning Sequences – Piotr Mirowski

# LSTM

[Hochreiter et al, 1997][Gers et al, 1999]

$$\begin{aligned}
 i_t &= W_{ix}x_t + W_{ih}h_{t-1} + b_i \\
 j_t &= W_{jx}x_t + W_{jh}h_{t-1} + b_j \\
 f_t &= W_{fx}x_t + W_{fh}h_{t-1} + b_f \\
 o_t &= W_{ox}x_t + W_{oh}h_{t-1} + b_o \\
 c_t &= \sigma(f_t) \odot c_{t-1} + \sigma(i_t) \odot \tanh(j_t) \\
 h_t &= \sigma(o_t) \odot \tanh(c_t)
 \end{aligned}$$

Enables  
long term  
dependencies  
to flow



```

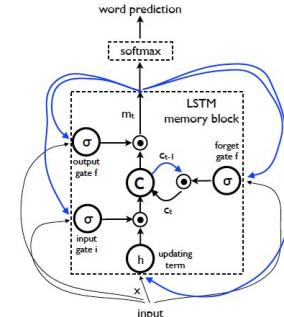
def __call__(self, inputs, state, scope=None):
    """Long short-term memory cell (LSTM)."""
    with vs.variable_scope(scope or type(self).__name__): # "BasicLSTMCell"
        # Parameters of gates are concatenated into one multiply for efficiency.
        c, h = array_ops.split(1, 2, state)
        concat = linear([inputs, h], 4 * self._num_units, True)

        # i = input_gate, j = new_input, f = forget_gate, o = output_gate
        i, j, f, o = array_ops.split(1, 4, concat)

        new_c = c * sigmoid(f + self._forget_bias) + sigmoid(i) * tanh(j)
        new_h = tanh(new_c) * sigmoid(o)

    return new_h, array_ops.concat(1, [new_c, new_h])

```



# How about exploding gradients? [Pascanu et al, 2013]

---

## Algorithm 1 Pseudo-code for norm clipping

---

```
 $\hat{\mathbf{g}} \leftarrow \frac{\partial \mathcal{E}}{\partial \theta}$ 
if  $\|\hat{\mathbf{g}}\| \geq threshold$  then
     $\hat{\mathbf{g}} \leftarrow \frac{threshold}{\|\hat{\mathbf{g}}\|} \hat{\mathbf{g}}$ 
end if
```

---

# Summary

- Sequences are important, everything is a sequence
- Recurrent Nets are a way forward, but they pose some challenges
  - Chain Rule makes us not assume independence
  - Vectorizing context is key for feasibility / model size
  - Vanishing/Exploding gradients
    - LSTMs is better at this, but by no means “solves” all problems
  - Expensive to compute

# Part IV: Recurrent Networks as Sequence Decoders

# Mini-Turing Test

Some of the obese people lived five to eight years longer than others.

Abu Dhabi is going ahead to build solar city and no pollution city.

Or someone who exposes exactly the truth while lying.

VIERA , FLA . -- Sometimes, Rick Eckstein dreams about baseball swings.

For decades, the quintessentially New York city has elevated its streets to the status of an icon.

The lawsuit was captioned as United States ex rel.

# Mini-Turing Test

Some of the obese people lived five to eight years longer than others.

Abu Dhabi is going ahead to build solar city and no pollution city.

Or someone who exposes exactly the truth while lying.

VIERA , FLA . -- Sometimes, Rick Eckstein dreams about baseball swings.

For decades, the quintessentially New York city has elevated its streets to the status of an icon.

The lawsuit was captioned as United States ex rel.

# Mini-Turing Test

Some of the obese people lived five to eight years longer than others.

Abu Dhabi is going ahead to build solar city and no pollution city.

Or someone who exposes exactly the truth while lying

VIERA , FLA . -- Sometimes, Rick Eckstein dreams about baseball swings.

For decades, the quintessentially New York city has elevated its streets to the status of an icon.

The lawsuit was captioned as United States ex rel.

# Results

Measure: perplexity (~branching factor (low == good))

Model	Perplexity	#params (Billions)
5-gram (smoothed) (2013)	67.6	1.8
RNN + Maxent (2013)	51.3	20 (!)
Big Ensemble (2013)	<b>43.8</b>	LOTS
RNN (ICLR, 2015)	68.3	4.1
RNN + 5-gram (ICLR, 2015)	<b>42.0</b>	LOTS
Sparse NNMF (IS, 2015)	52.9	33 (!)

# Results

Measure: perplexity (~branching factor (low == good))

Model	Perplexity	#params (Billions)
5-gram (smoothed) (2013)	67.6	1.8
RNN + Maxent (2013)	51.3	20 (!)
Big Ensemble (2013)	43.8	LOTS
RNN (ICLR, 2015)	68.3	4.1
RNN + 5-gram (ICLR, 2015)	42.0	LOTS
Sparse NNMF (IS, 2015)	52.9	33 (!)
<b>Ours (LSTM) (2016)</b>	<b>30.0</b>	<b>1.0</b>
<b>Ours (Ensemble) (2016)</b>	<b>23.7</b>	LOTS

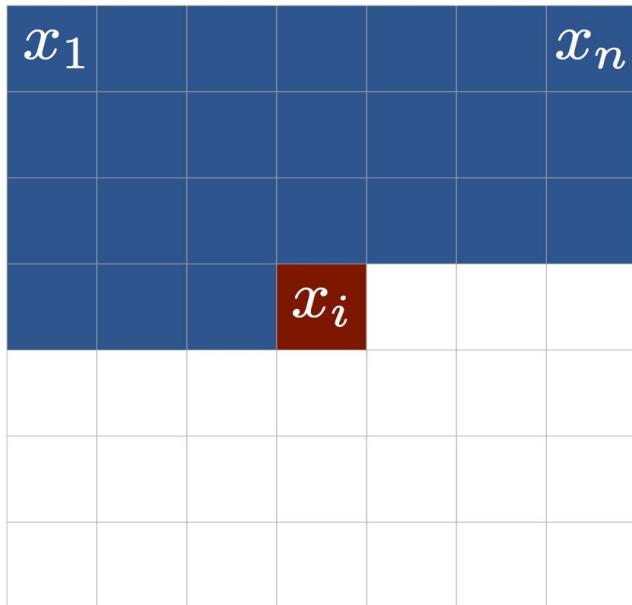
# Results

Measure: perplexity (~branching factor (low == good))

Model	Perplexity	#params (Billions)
5-gram (smoothed) (2013)	67.6	1.8
RNN + Maxent (2013)	51.3	20
Big Ensemble (2013)	45.3	LOTS
RNN (ICLR, 2015)	68.3	4.1
RNN + 5-gram (ICLR, 2015)	42.0	LOTS
Sparse NNMF (IS, 2015)	52.9	33 (!)
<b>Ours (LSTM) (2016)</b>	<b>30.0</b>	<b>1.0</b>
<b>Ours (Ensemble) (2016)</b>	<b>23.7</b>	LOTS

**HUMANS: ???**

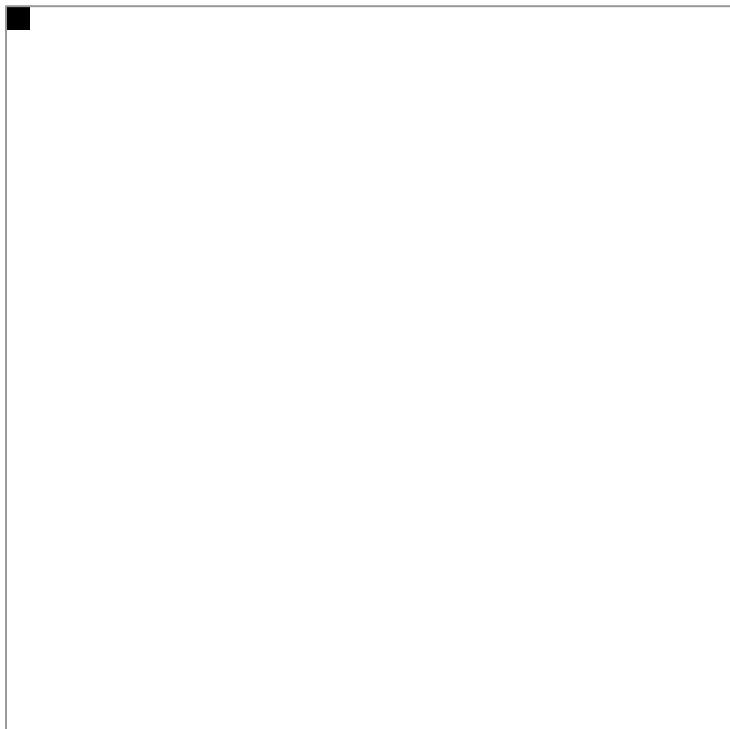
# PixelRNN - Model



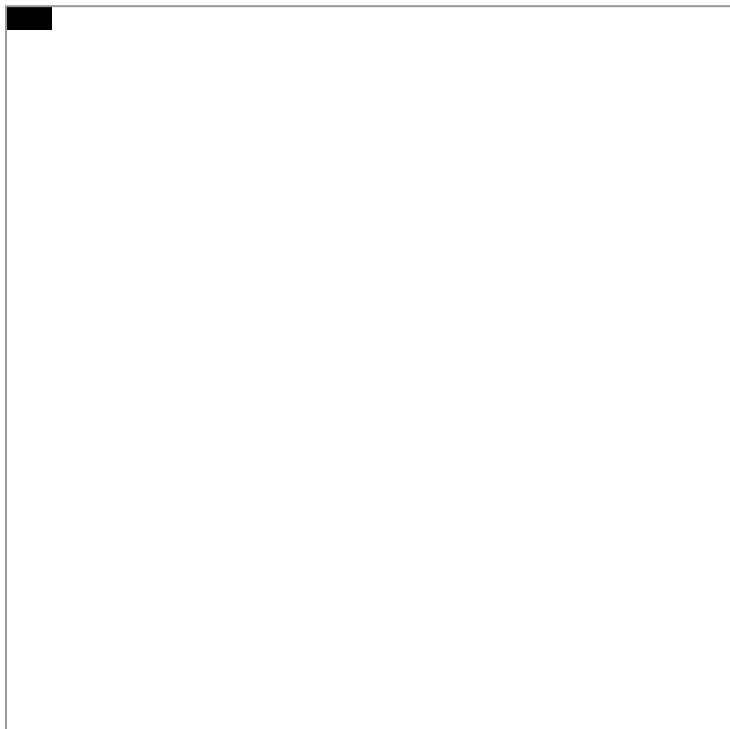
$$p(\mathbf{x}) = \prod_{i=1}^{n^2} p(x_i | x_1, \dots, x_{i-1})$$

- Fully visible
- Similar to language models with RNNs
- Model pixels with Softmax

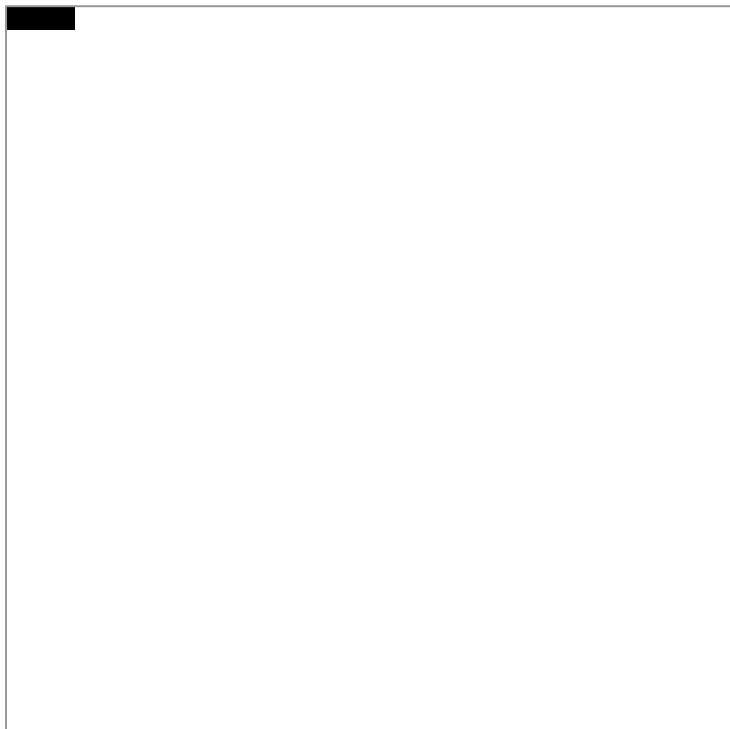
# Softmax Sampling



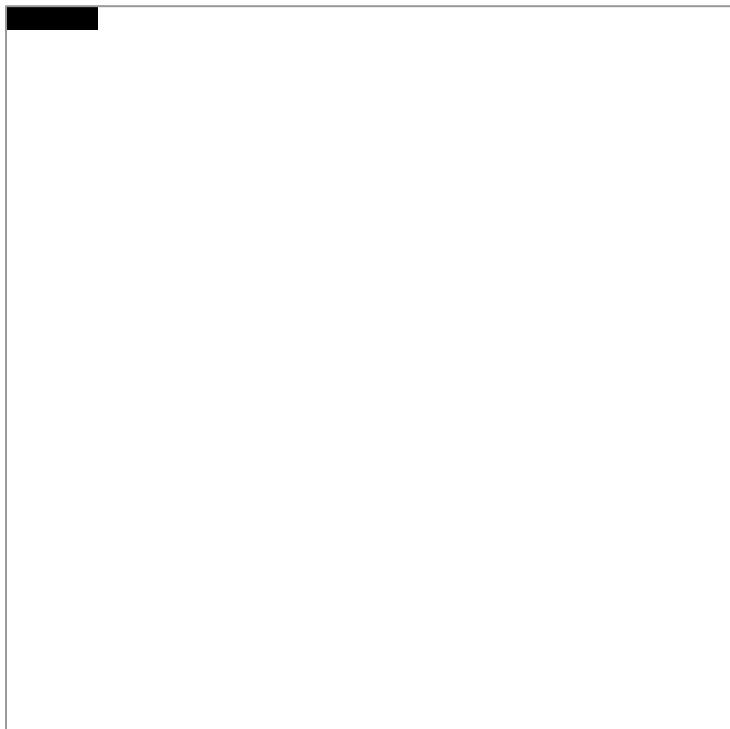
# Softmax Sampling



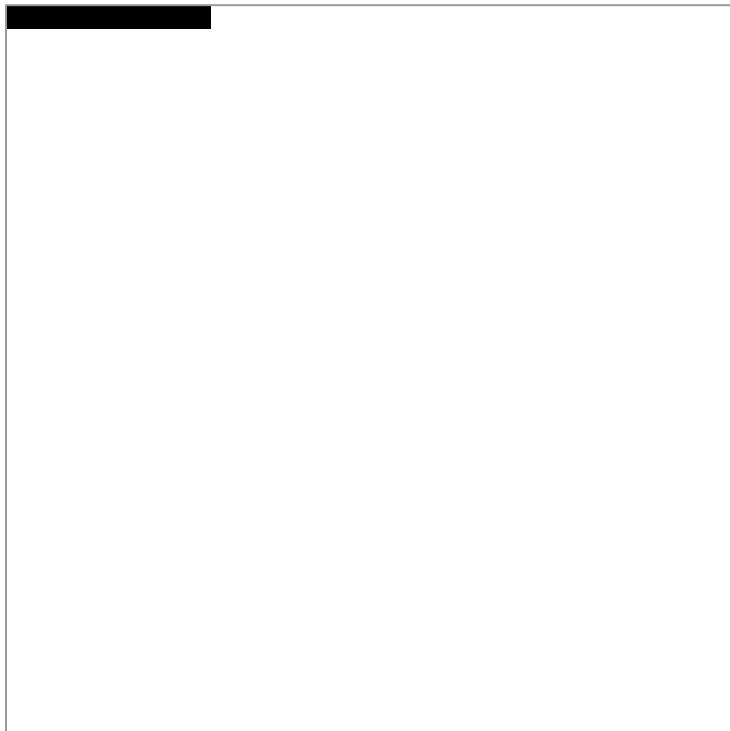
# Softmax Sampling



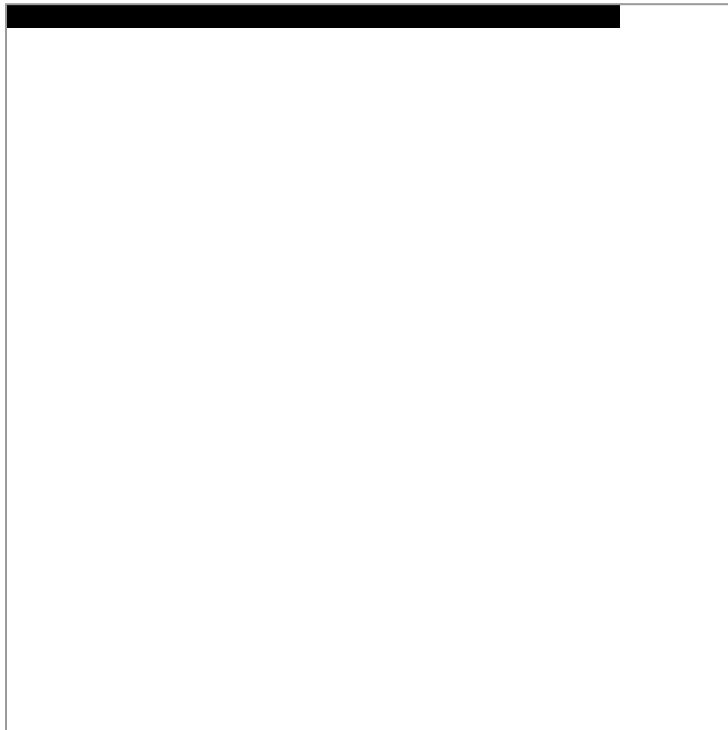
# Softmax Sampling



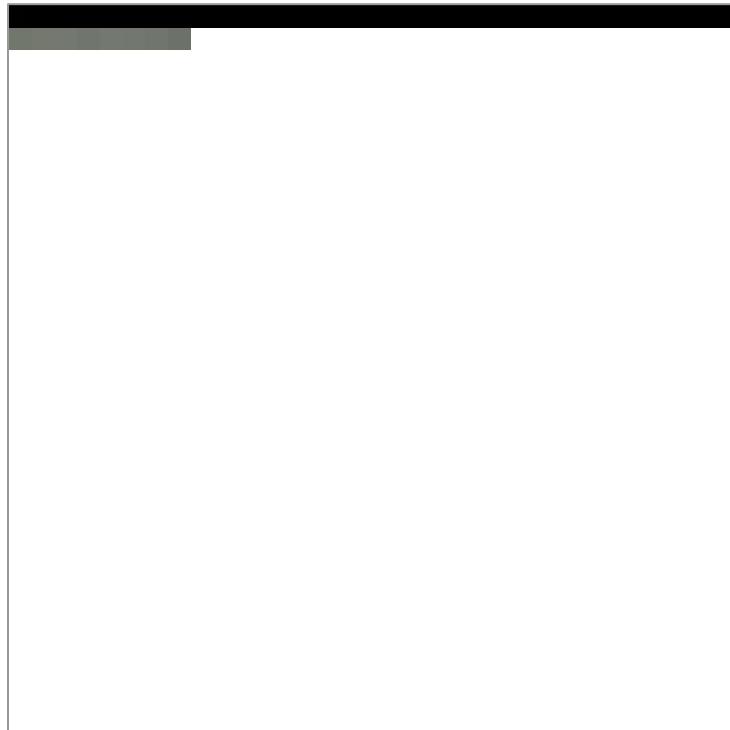
# Softmax Sampling



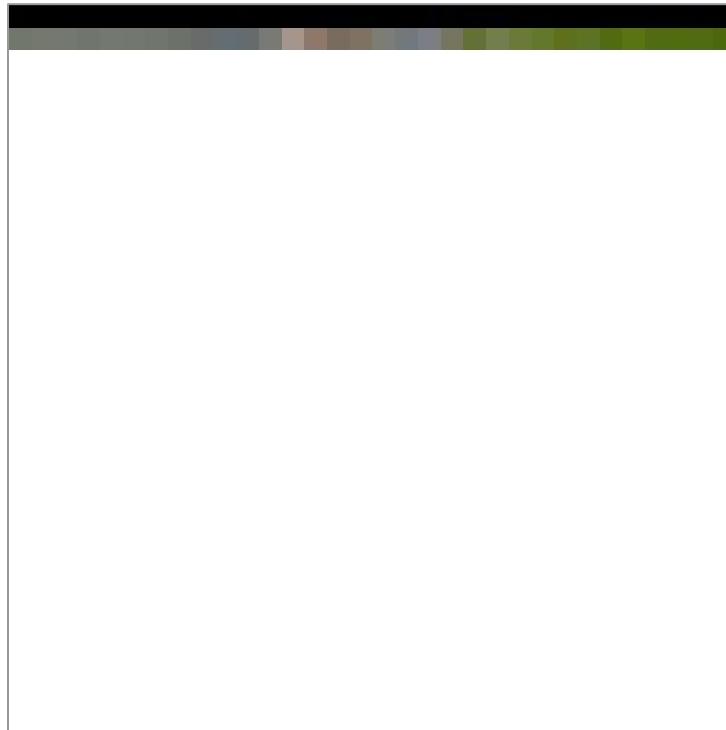
# Softmax Sampling



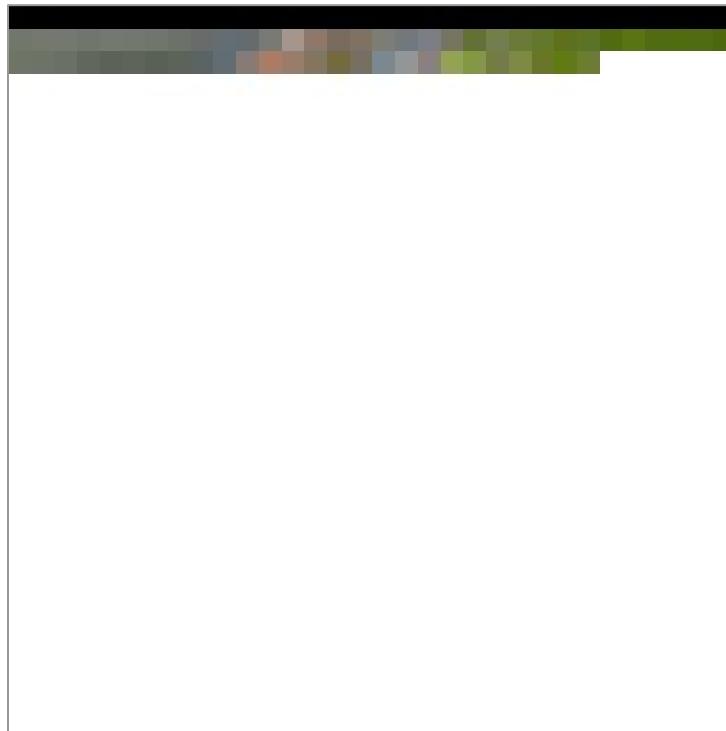
# Softmax Sampling



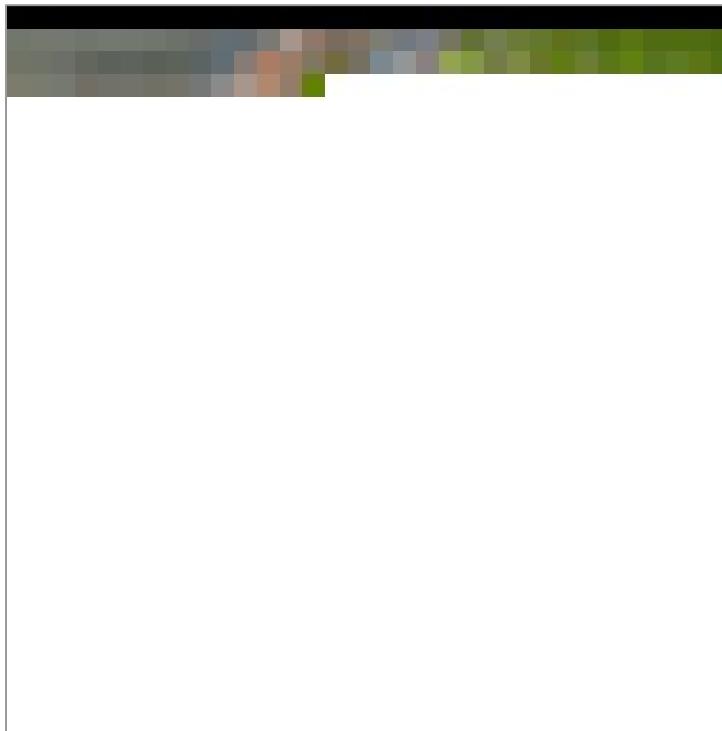
# Softmax Sampling



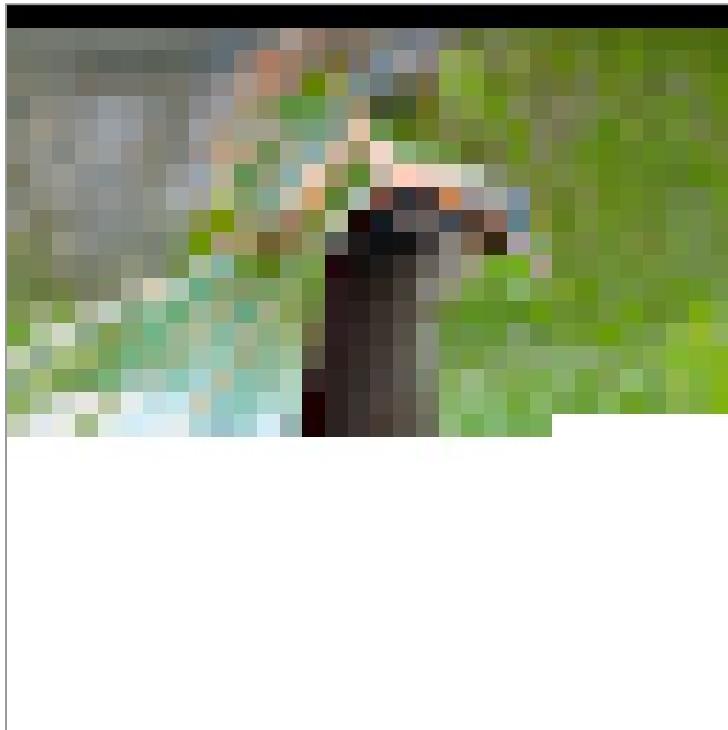
# Softmax Sampling



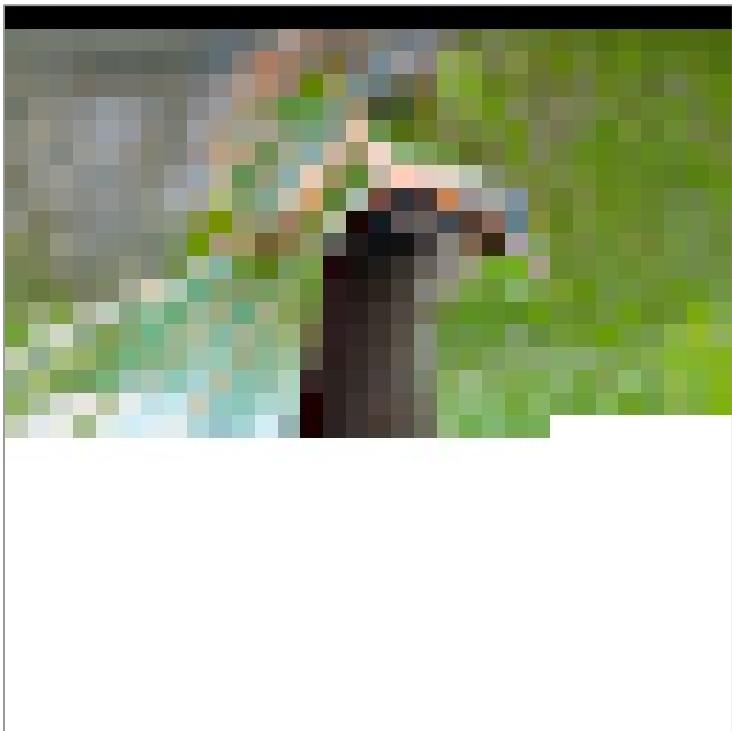
# Softmax Sampling



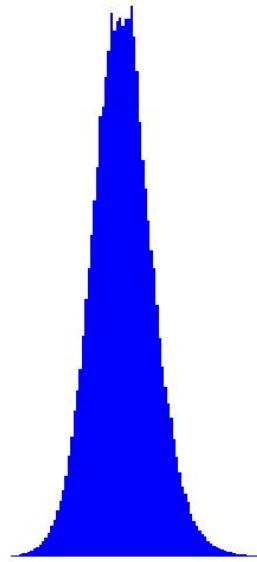
# Softmax Sampling



# Softmax Sampling

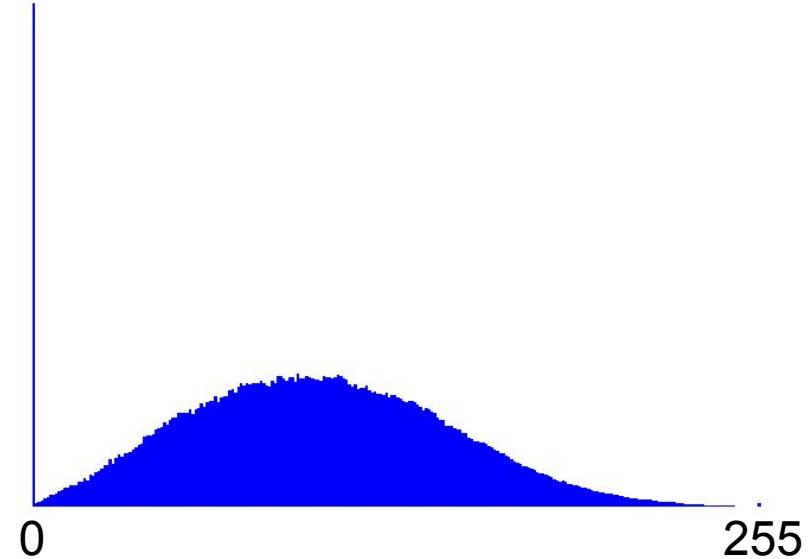
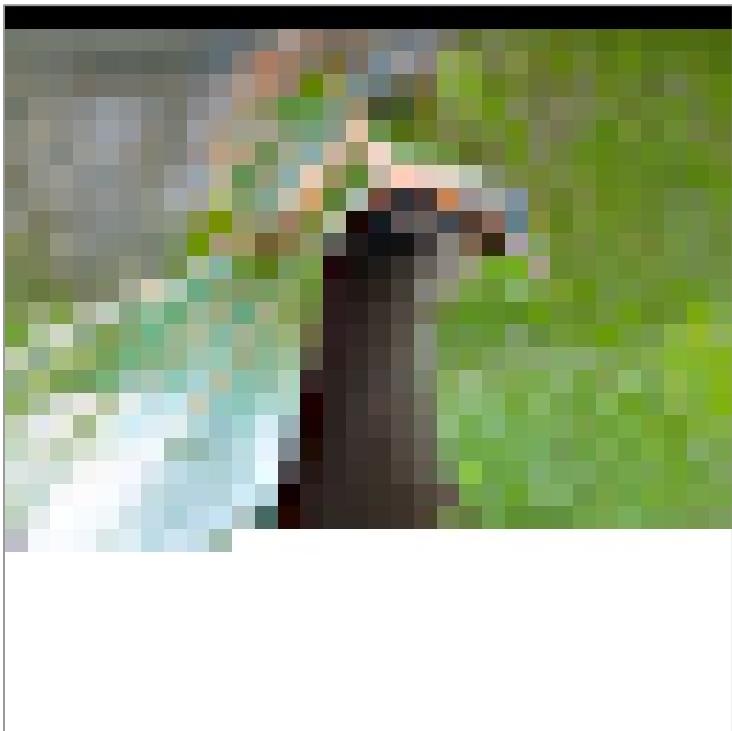


0

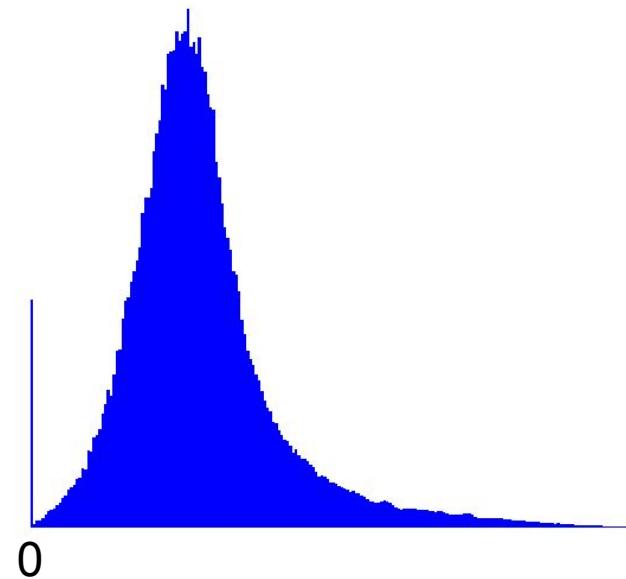
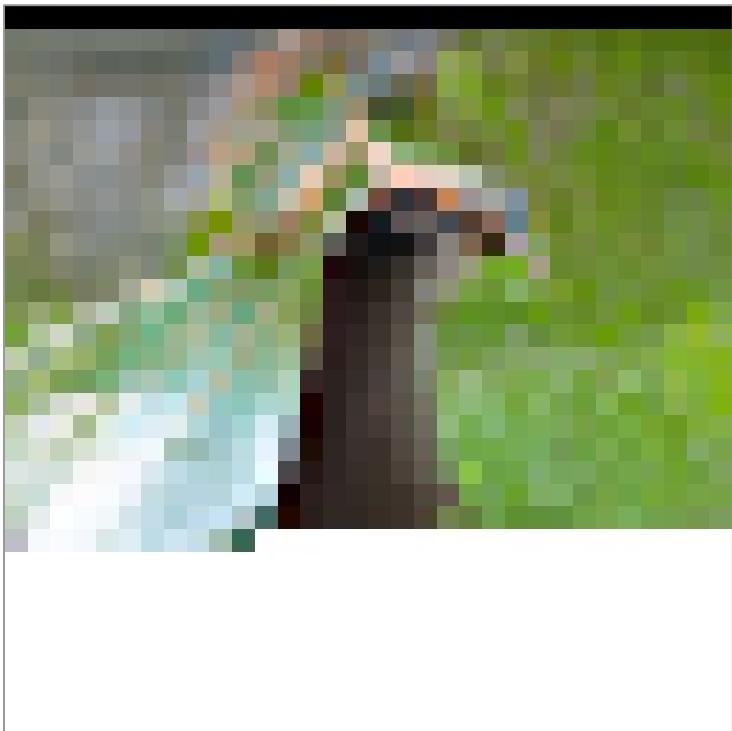


255

# Softmax Sampling

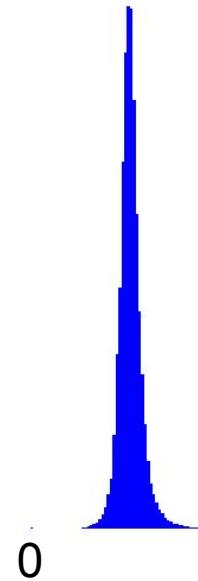
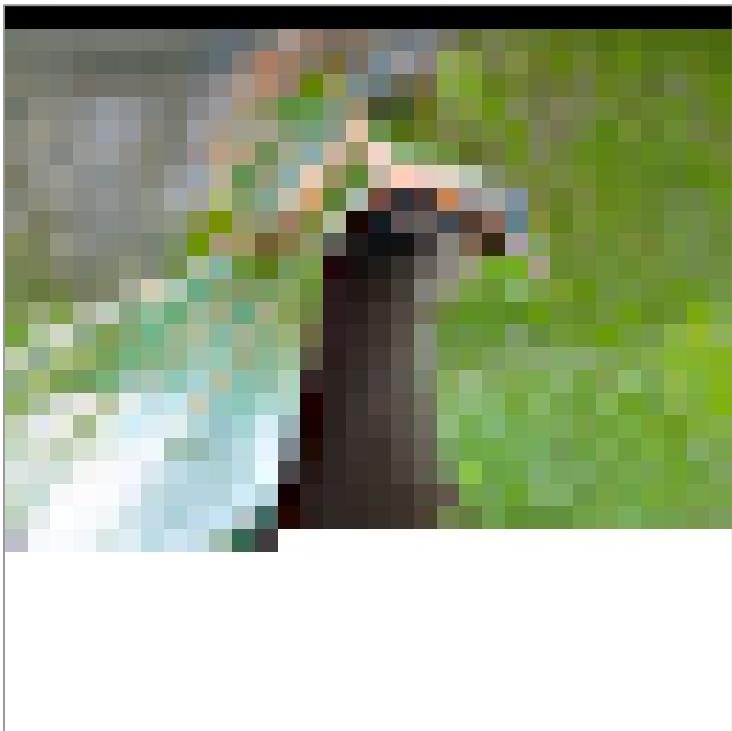


# Softmax Sampling



255

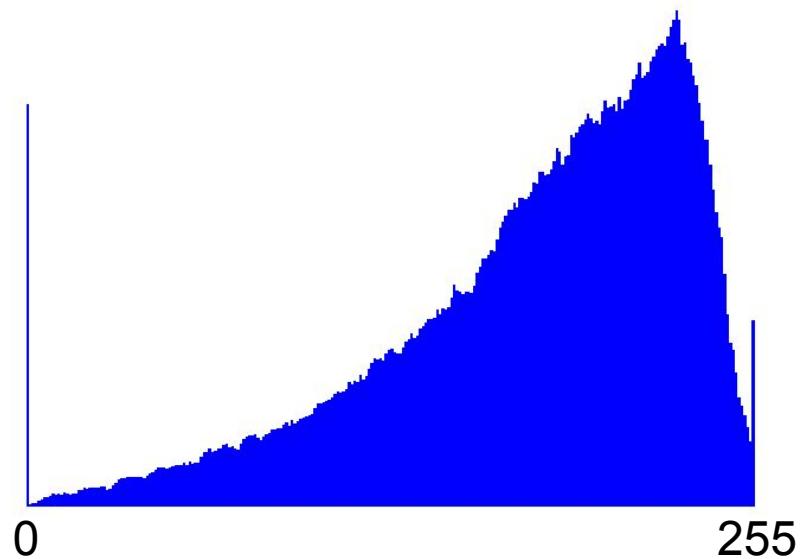
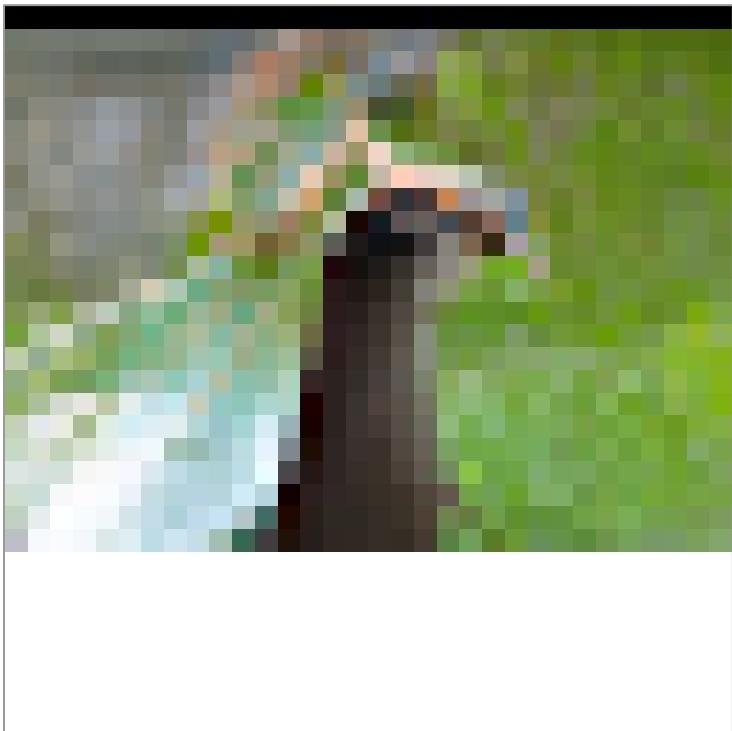
# Softmax Sampling



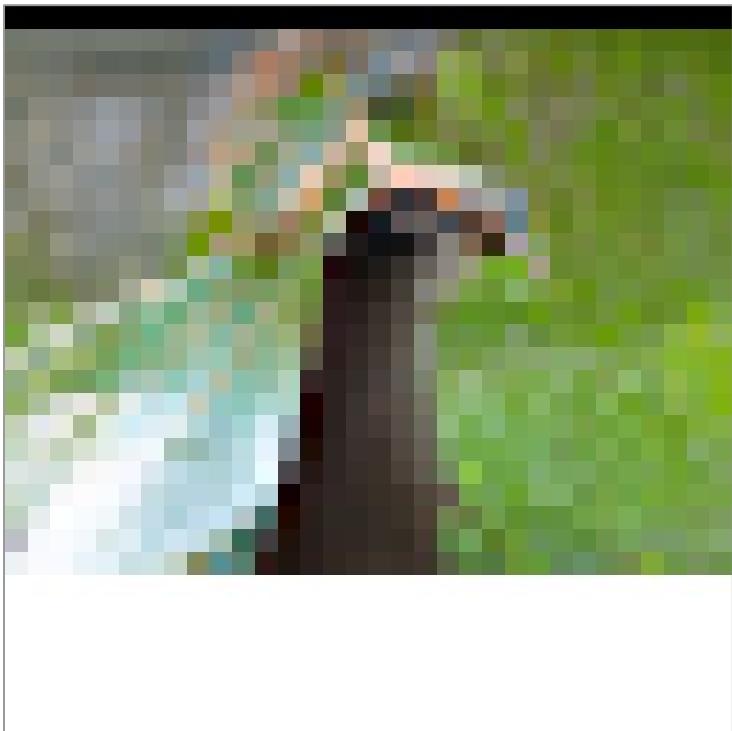
0

255

# Softmax Sampling

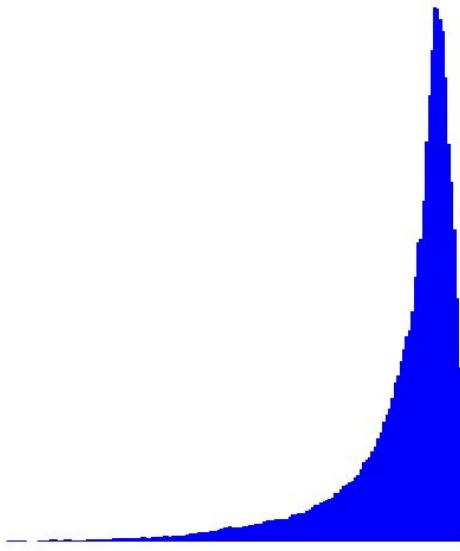


# Softmax Sampling

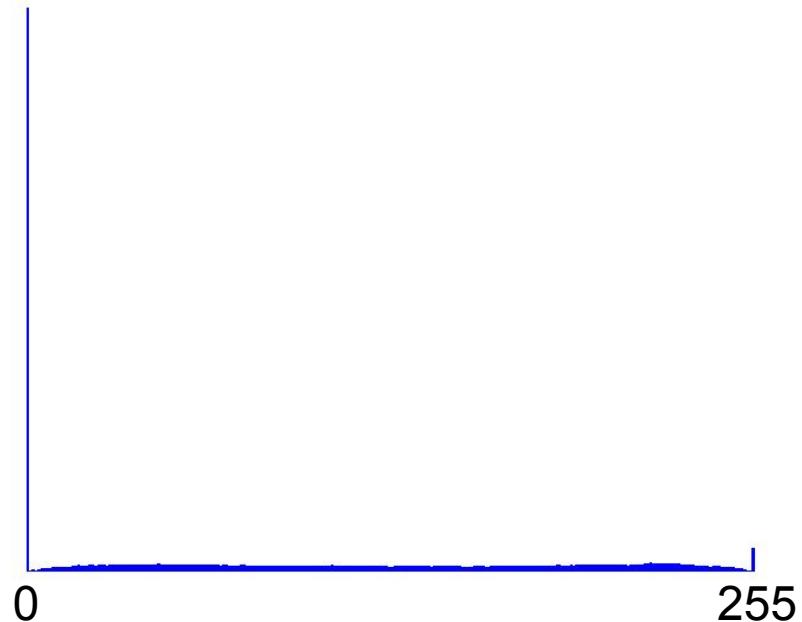
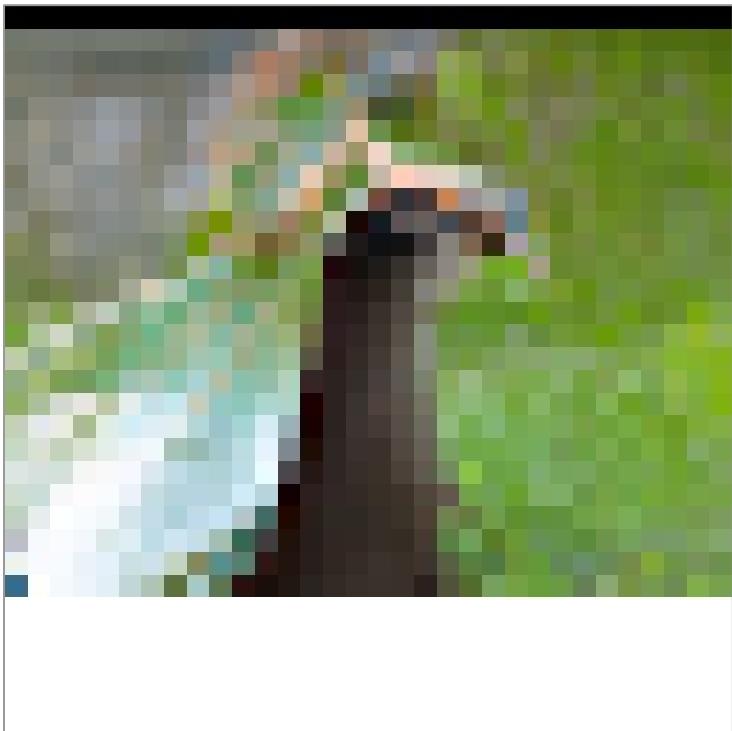


0

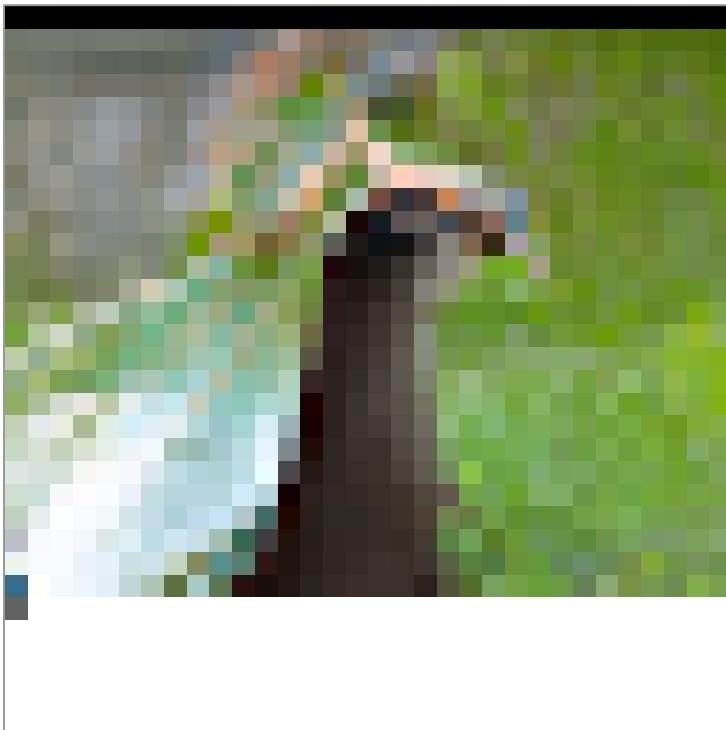
255



# Softmax Sampling

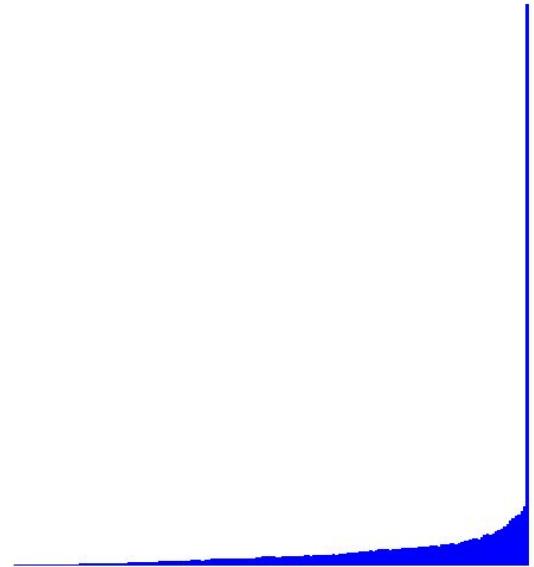


# Softmax Sampling

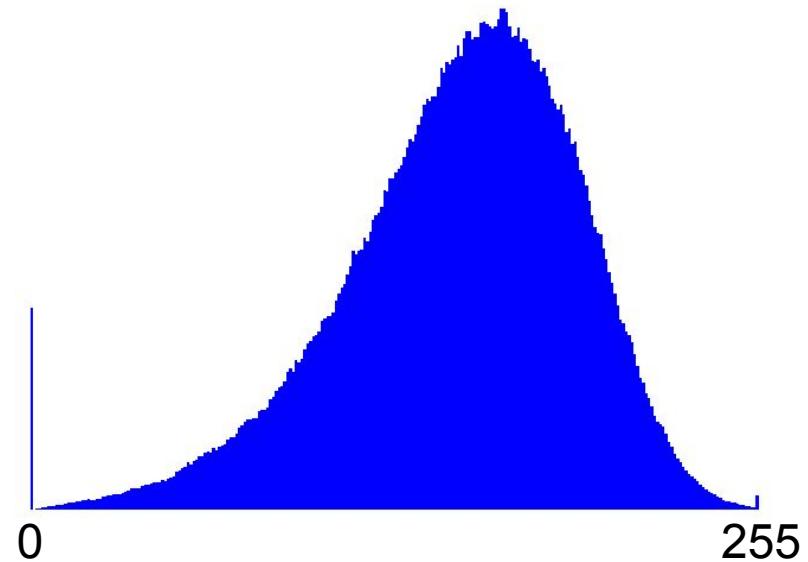
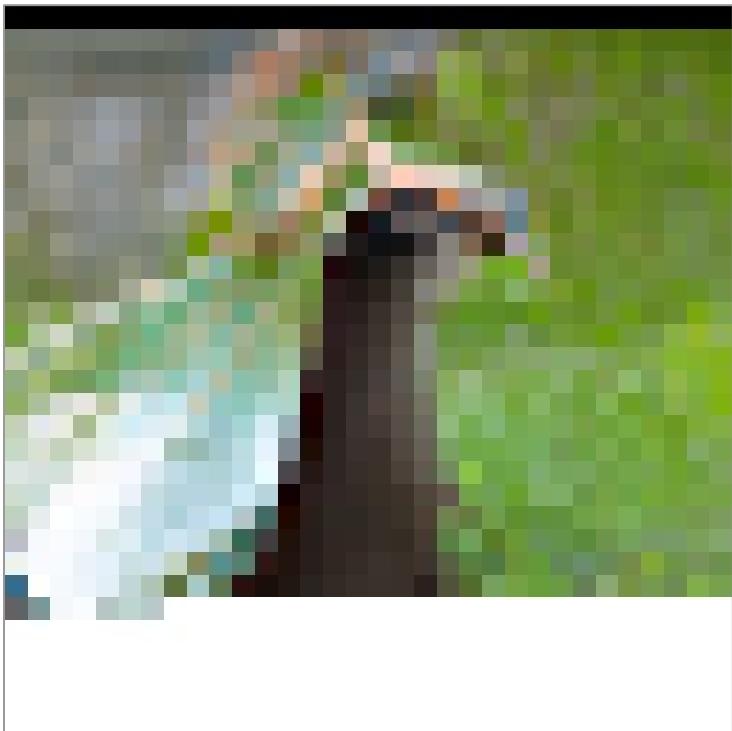


0

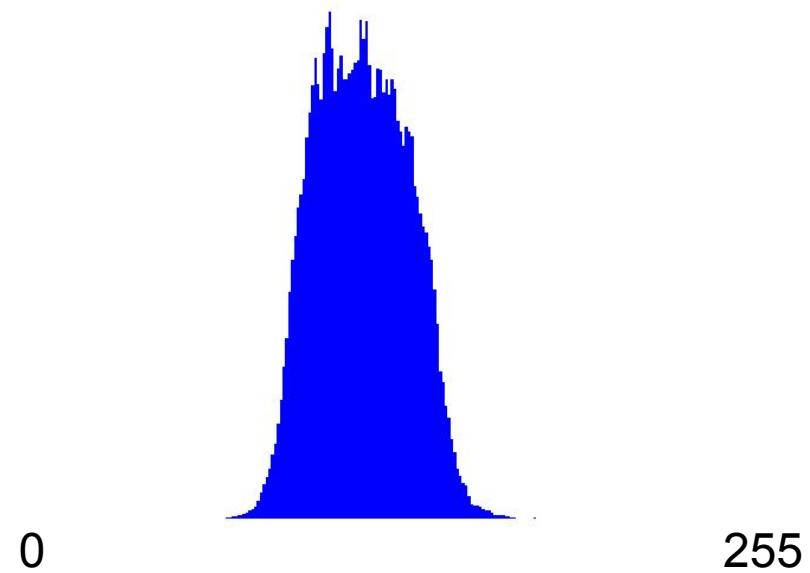
255



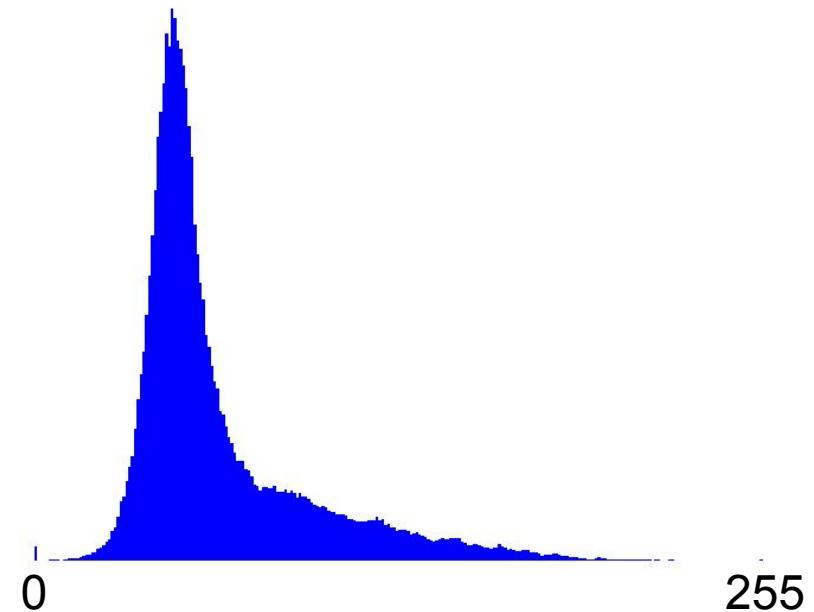
# Softmax Sampling



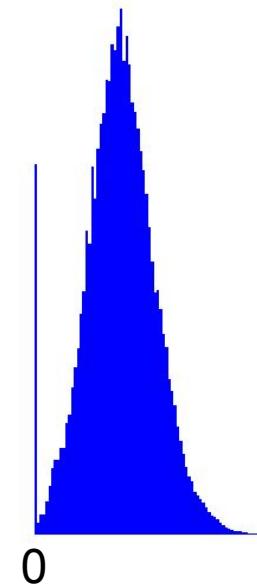
# Softmax Sampling



# Softmax Sampling



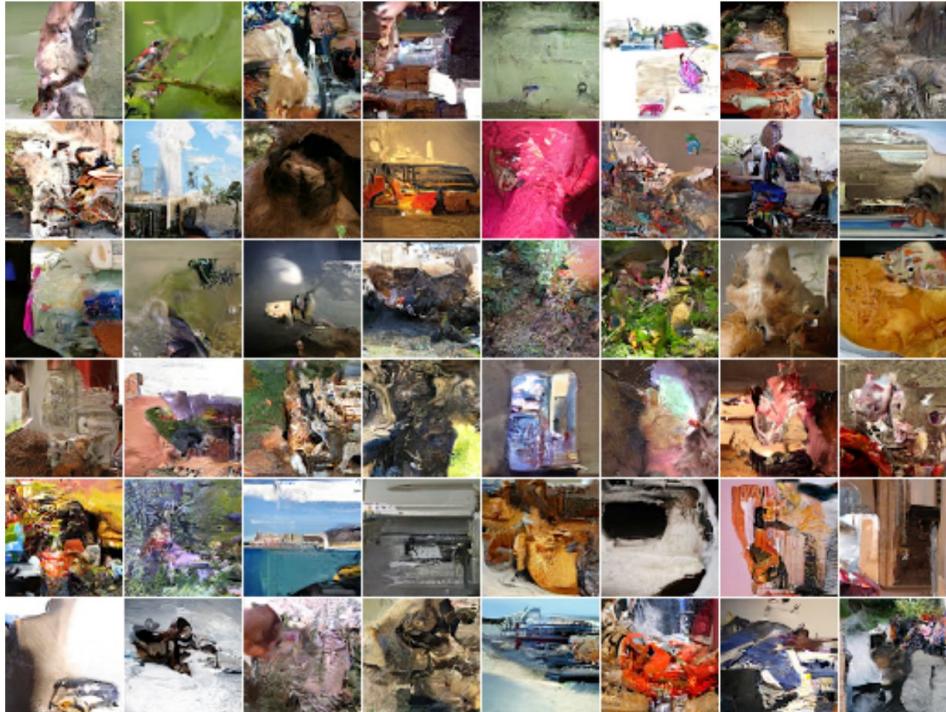
# Softmax Sampling



255

# Pixel RNN [van den Oord, Kalchbrenner, Kavukcuoglu, ICML16]

# Sequence of Words == Sequence of Pixels

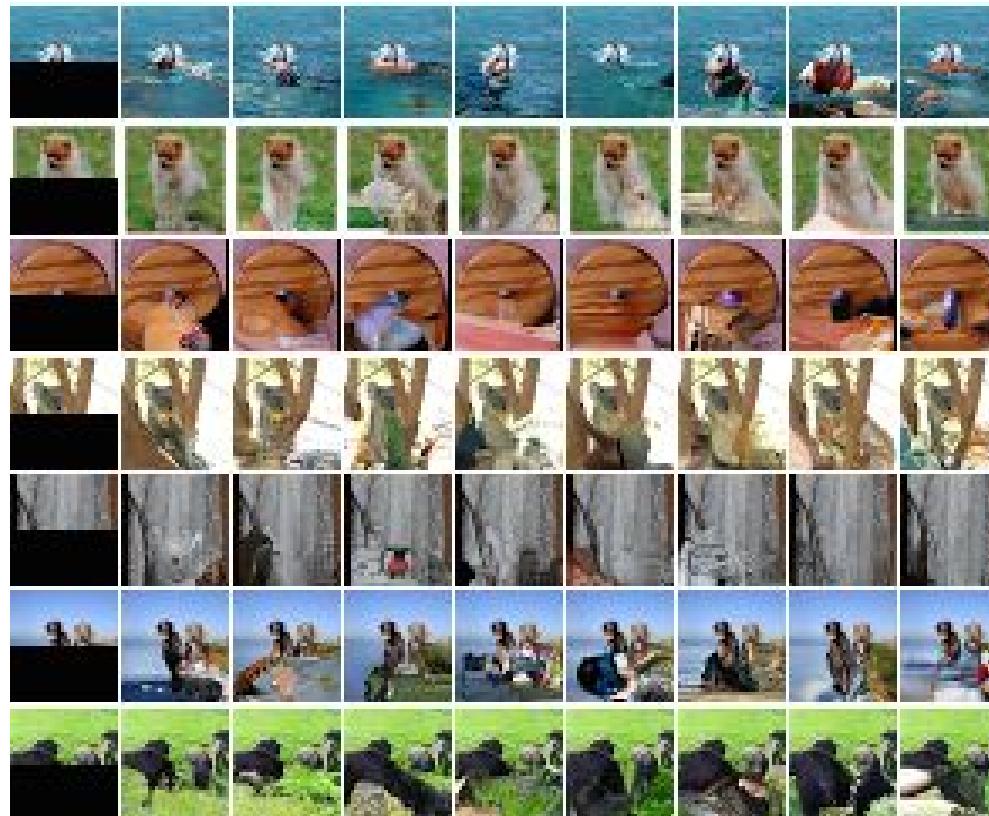


# occluded



**occluded**

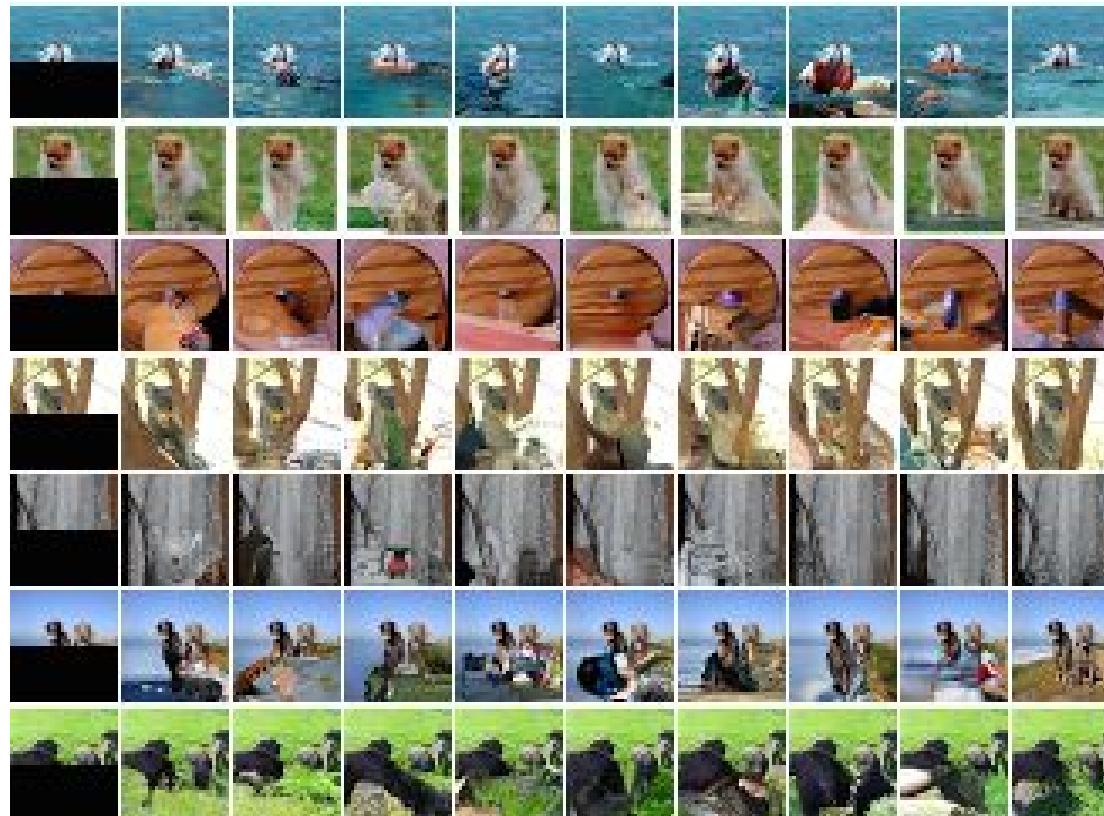
**completions**



**occluded**

**completions**

**original**



# Part IV.A: Sequence-to-Sequence & Applications

# Sequence-to-Sequence

1. MT [Kalchbrenner et al, EMNLP 2013][Cho et al, EMLP 2014]**[Sutskever & Vinyals & Le, NIPS 2014]**[Luong et al, ACL 2015][Bahdanau et al, ICLR 2015]
2. Image captions [Mao et al, ICLR 2015][Vinyals et al, CVPR 2015][Donahue et al, CVPR 2015][Xu et al, ICML 2015]
3. Speech [Chorowsky et al, NIPS DL 2014][Chan et al, arxiv 2015]
4. Parsing [Vinyals & Kaiser et al, NIPS 2015]
5. Dialogue [Shang et al, ACL 2015][Sordoni et al, NAACL 2015][Vinyals & Le, ICML DL 2015]
6. Video Generation [Srivastava et al, ICML 2015]
7. Geometry [Vinyals & Fortunato & Jaity, NIPS 2015]

# Sequence-to-Sequence



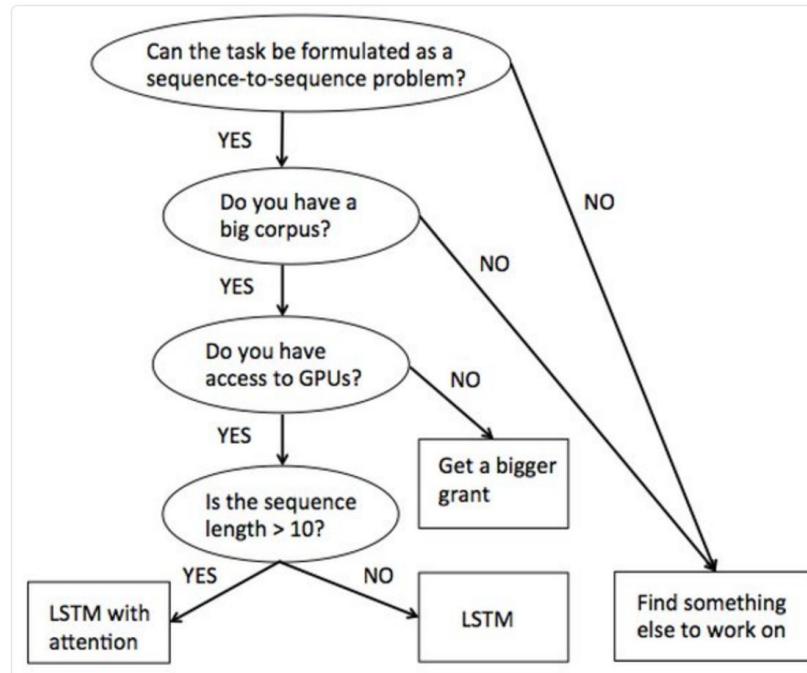
Isabelle Augenstein

@IAugenstei

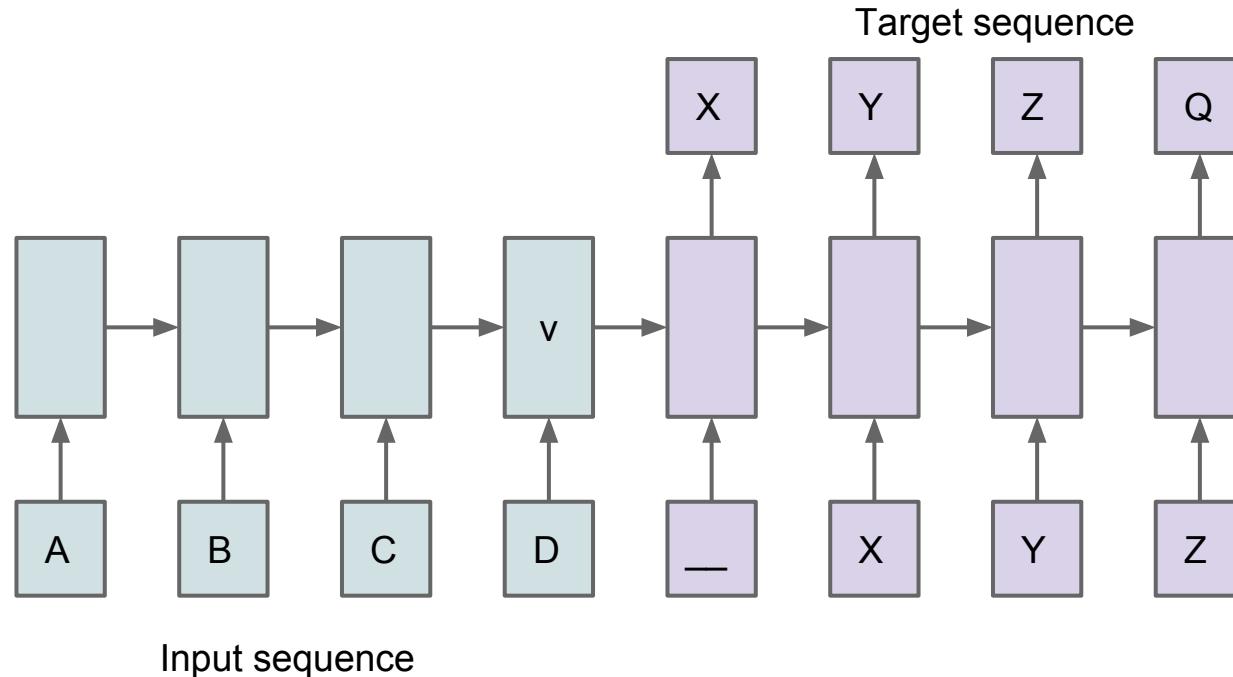


Follow

Brainstorming at [@uclmr](#) : how to recommend "what method works best for NLP task X"? Our baseline is ready [#NLProc](#)



# Main idea



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

# Training vs Inference

Training

$$\theta^* = \arg \max_{\theta} P(y_1, \dots, y_{T'} | x_1, \dots, x_T)$$

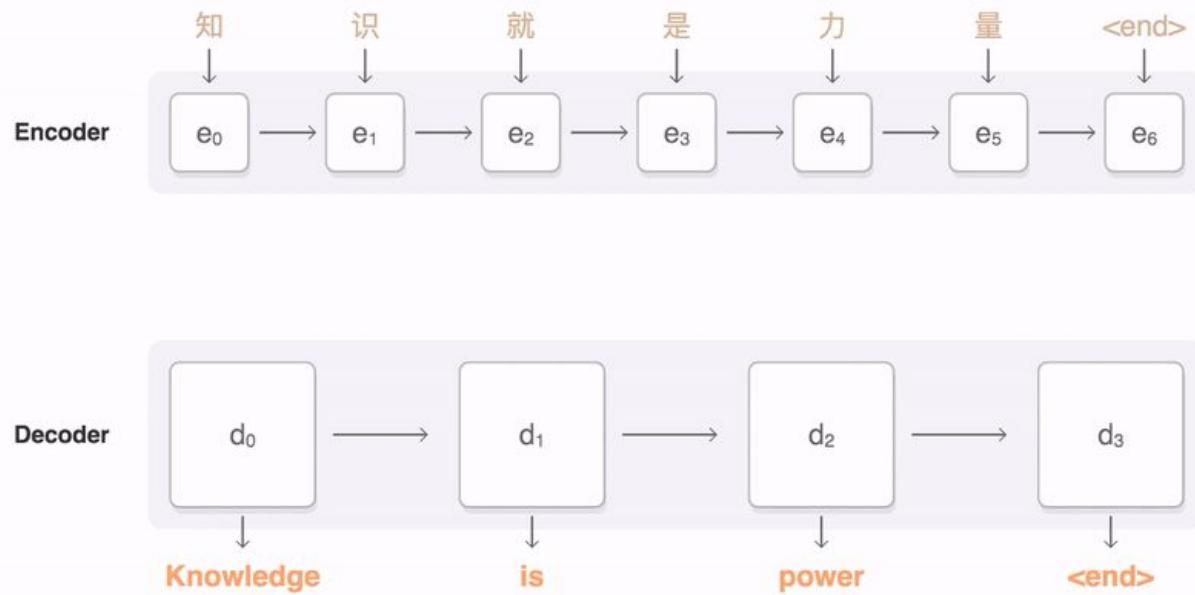
Inference

$$Y^* = \arg \max_Y P(Y | x_1, \dots, x_T)$$

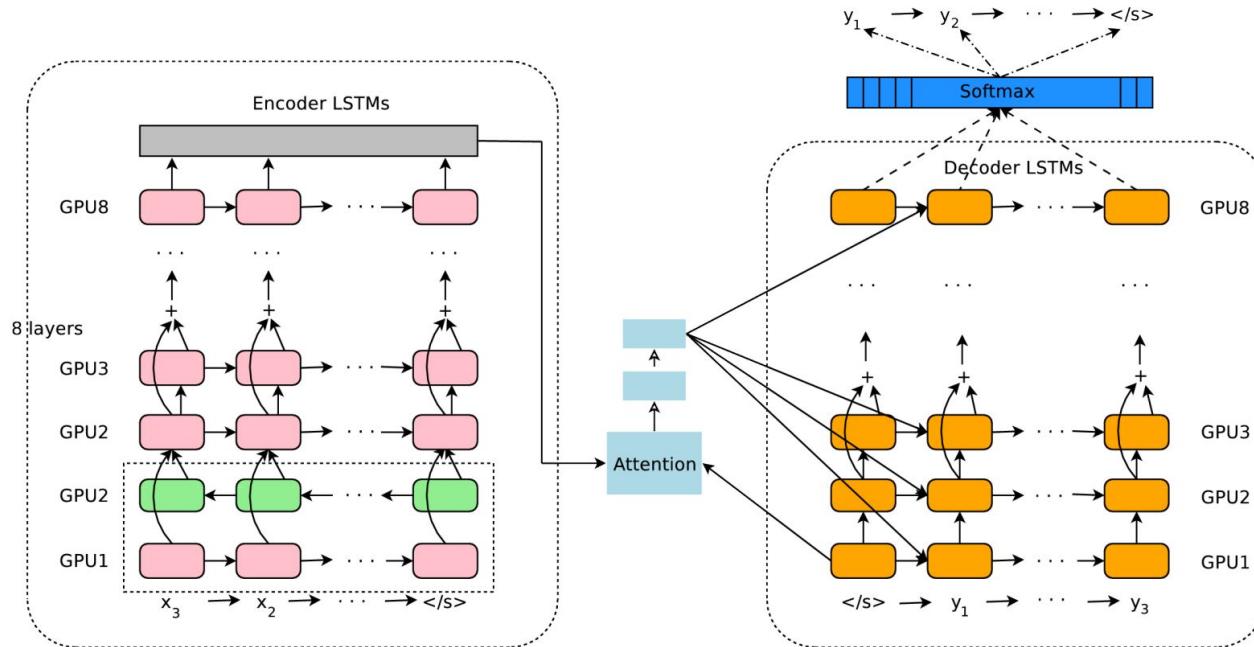
- Used beam search decoding over Y
- Small beam (~20) works fine

# Google Neural MT

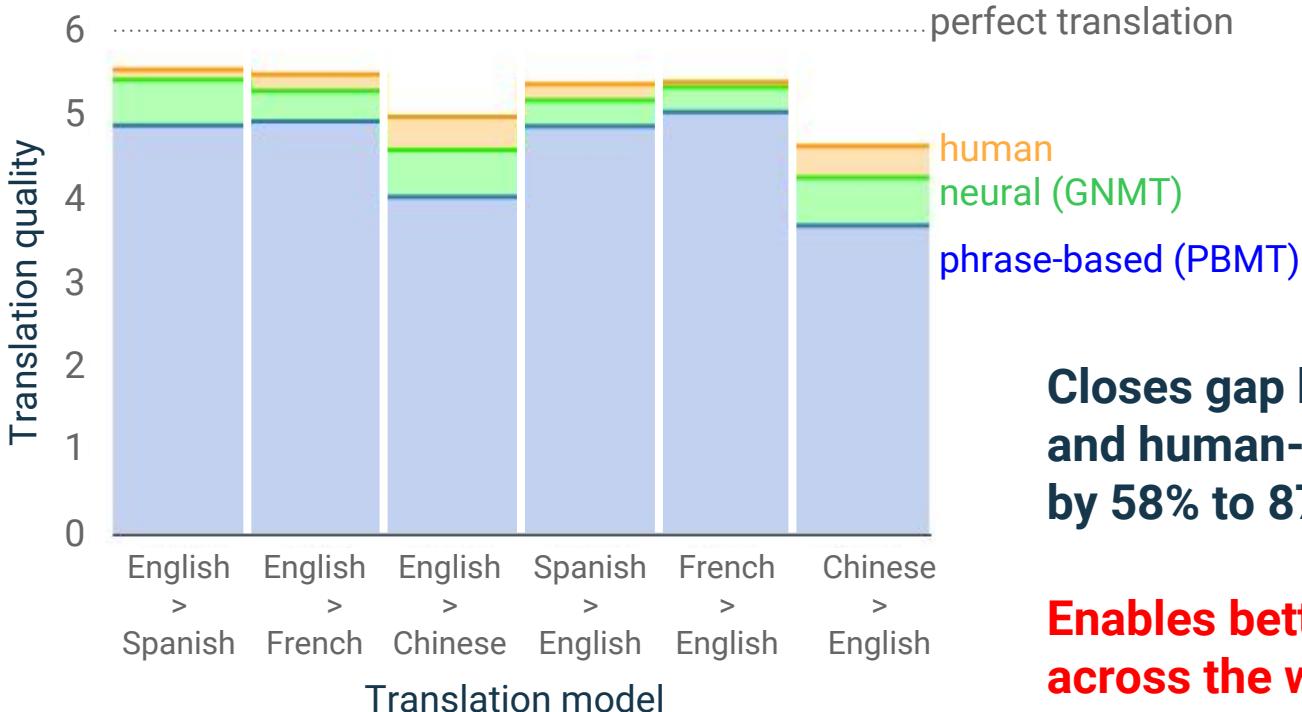
Wu et al, 2016 (Kalchbrenner et al, 2013; Sutskever et al, 2014;  
Cho et al, 2014; Bahdanau et al, 2014; ...)



# Google Neural MT



# Google Neural MT



**Closes gap between old system  
and human-quality translation  
by 58% to 87%**

**Enables better communication  
across the world**



*Human:* A young girl asleep on the sofa cuddling a stuffed bear.

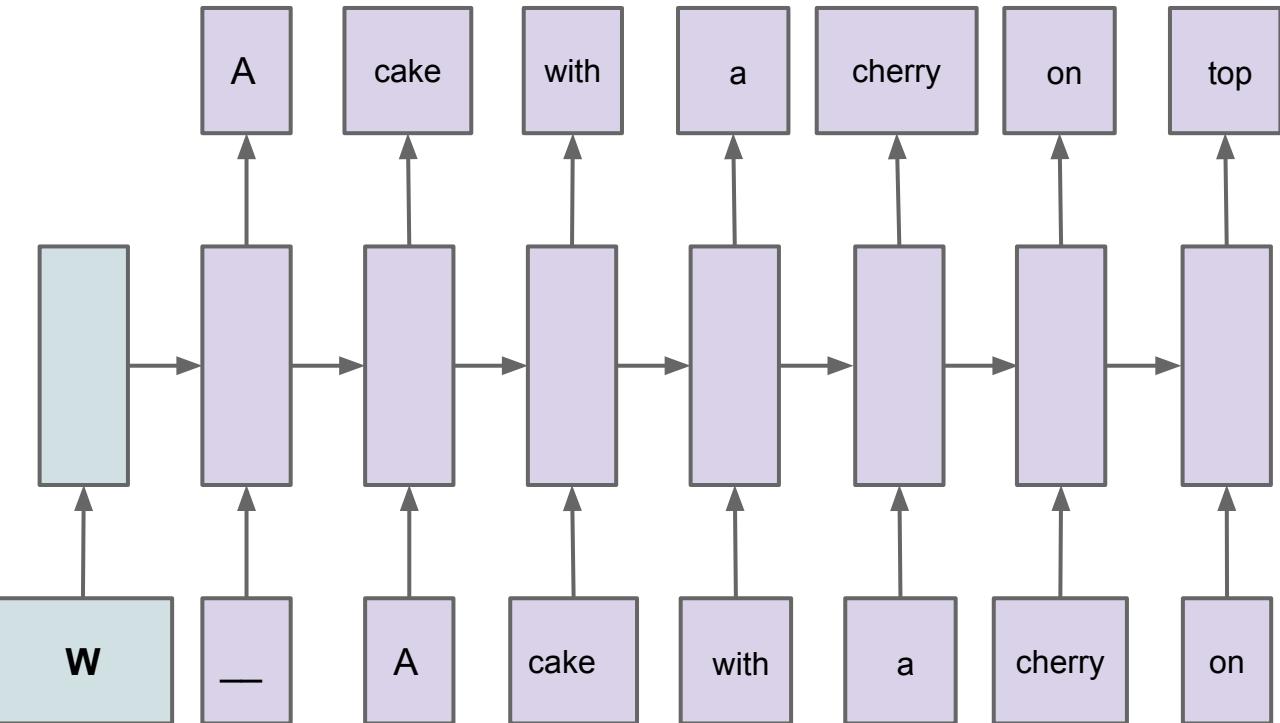
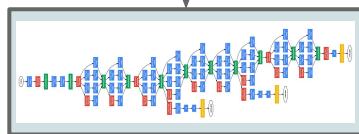
*NIC:* A close up of a child holding a stuffed animal.

*NIC:* A baby is asleep next to a teddy bear.

# How to do Image Captions?

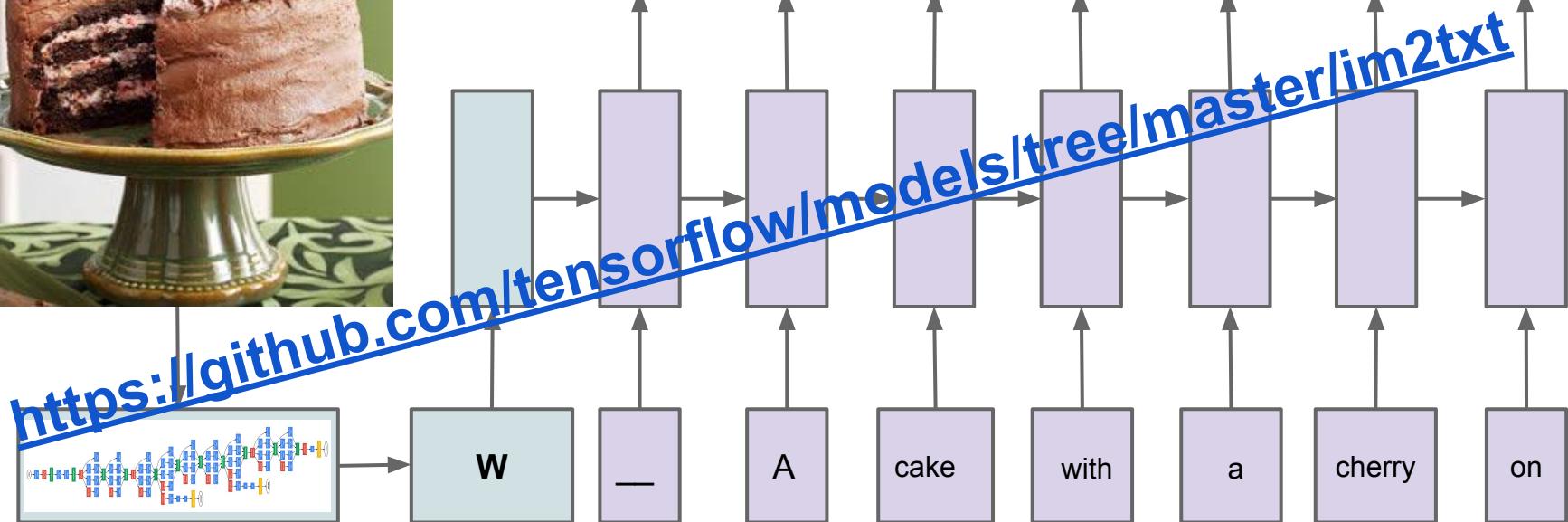
P(English | French)

# Captioning



$$\theta^* = \arg \max_{\theta} p(S|I)$$

# Captioning



$$\theta^* = \arg \max_{\theta} p(S|I)$$



*Human: A close up of two bananas with bottles in the background.*

*BestModel: A bunch of bananas and a bottle of wine.*

*InitialModel: A close up of a plate of food on a table.*



*Human: A view of inside of a car where a cat is laying down.*

*BestModel: A cat sitting on top of a black car.*

*InitialModel: A dog sitting in the passenger seat of a car.*



*Human: A brown dog laying in a red wicker bed.*

*BestModel: A small dog is sitting on a chair.*

*InitialModel: A large brown dog laying on top of a couch.*



*Human: A man outside cooking with a sub in his hand.*

*BestModel: A man is holding a sandwich in his hand.*

*InitialModel: A man cutting a cake with a knife.*



*Human: Someone is using a small grill to melt his sandwich.*

*BestModel: A person is cooking some food on a grill.*

*InitialModel: A pizza sitting on top of a white plate.*



*Human: A woman holding up a yellow banana to her face.*

*BestModel: A woman holding a banana up to her face.*

*InitialModel: A close up of a person eating a hot dog.*



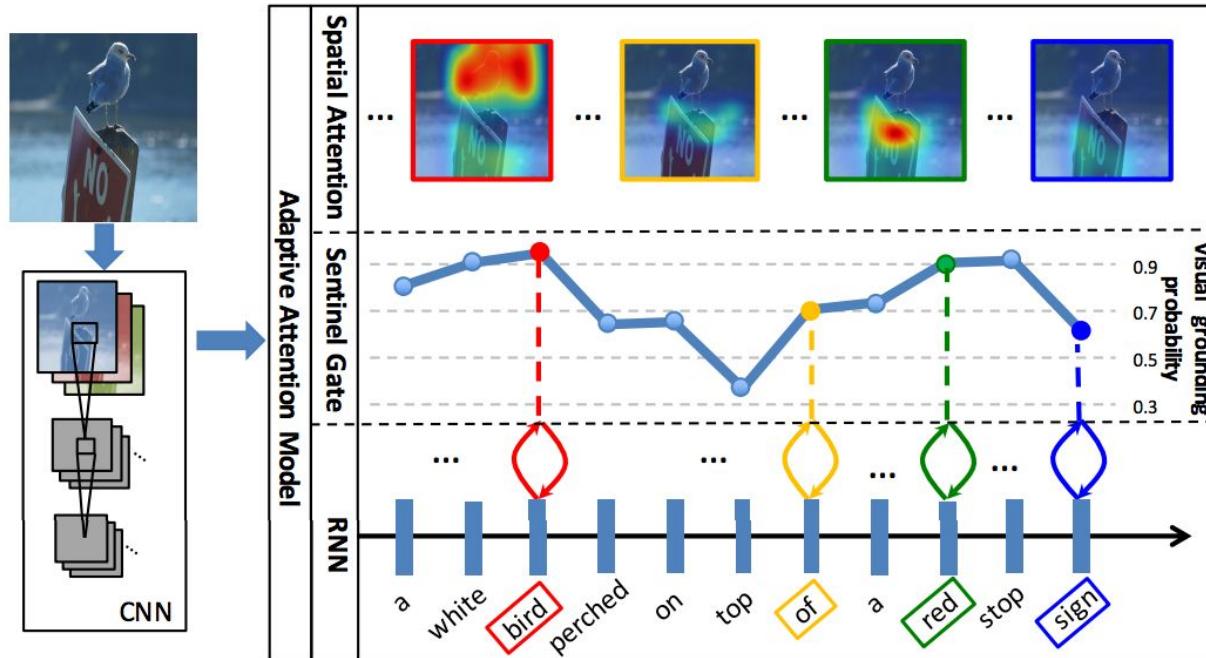
*Human: A blue , yellow and red train travels across the tracks near a depot.*

*BestModel: A blue and yellow train traveling down train tracks.*

*InitialModel: A train that is sitting on the tracks.*

# Advanced Networks: Captioning

Figure from [Lu et al, 2016]



# Conditional Pixel CNN

[van den Oord, Kalchbrenner, Vinyals, Espeholt, Graves, Kavukcuoglu, NIPS16]



Geyser



Hartebeest



Grey whale



Tiger

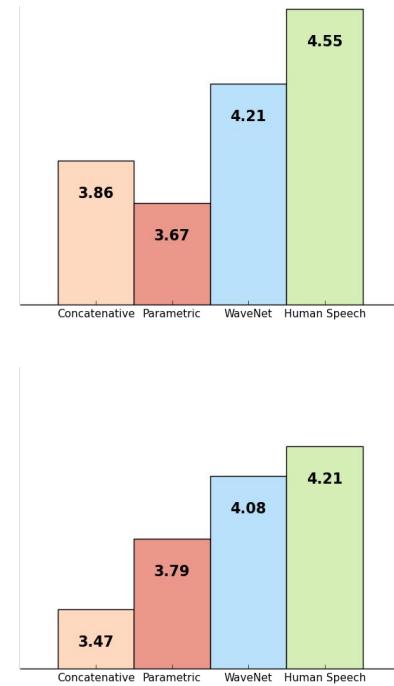
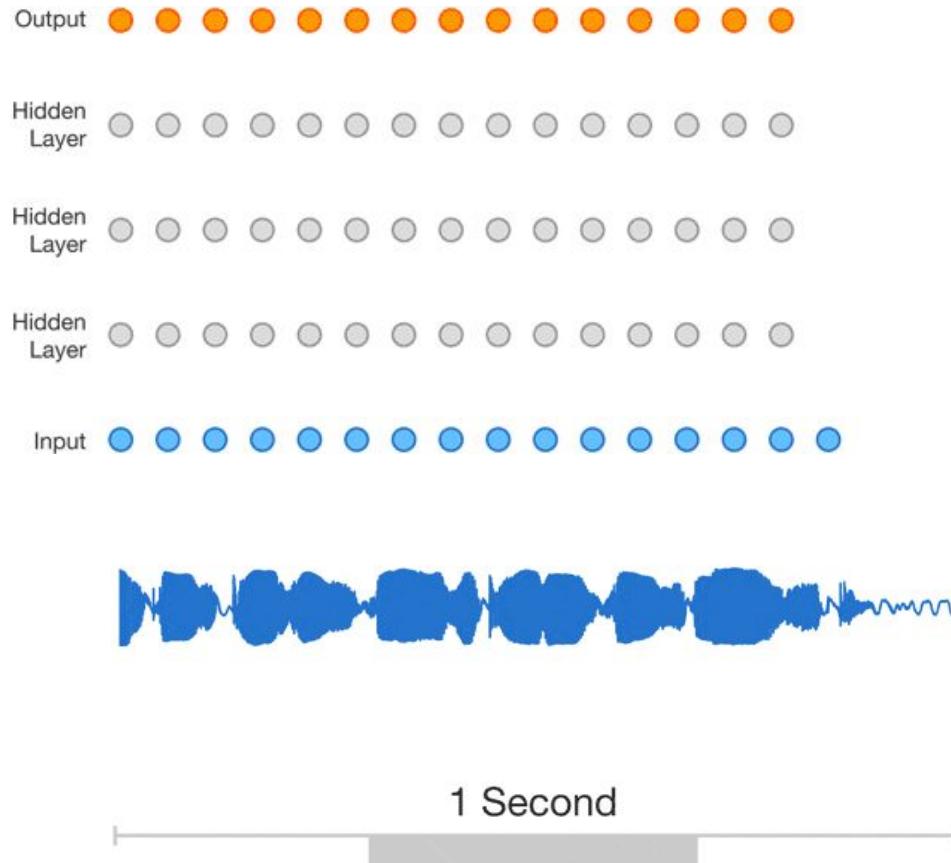


EntleBucher (dog)



Yellow lady's slipper (flower)

# WaveNets [van den Oord et al, 2016]



# Neural Conversational Models

- Take movie subtitles (~900M words) or IT HelpDesk chats
- Predict the next dialog from history

i got to go .

no .

i get too emotional when i drink .

have another beer . i 've got to get up early .

no , you don 't . sit down .

i get too emotional when i drink .

will you have another beer ?

i 've got to go !

why ?

i got to get up early in the morning .

you 're drunk .

and emotional !

you got to go .

# Conditional Language [Vinyals, Le, ICML DL WS 2015]

Human: what do you think about tesla ? Machine: he 's a good conductor .

Human: What do you think about messi ? Machine: he 's a great player .

Human: what do you think about cleopatra ? Machine: oh , she 's very regal .

Human: who is skywalker ? Machine: he is a hero .

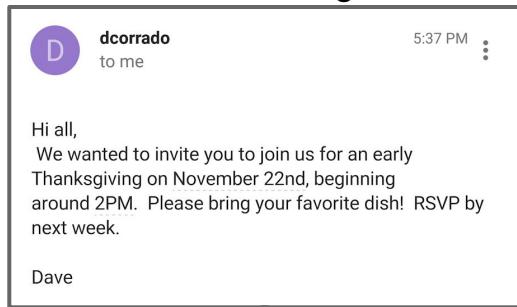
Human: is sky blue or black ? Machine: blue .

Human: does a cat have a tail ? Machine: yes .

Human: does a cat have a wing ? Machine: no

# Smart Reply

Incoming Email



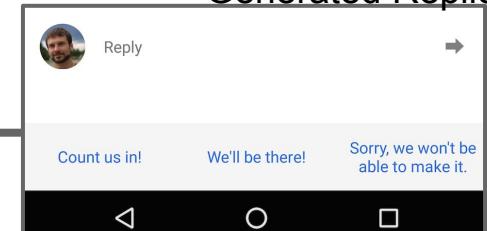
Google Research Blog  
- Nov 2015

Small  
Feed-Forward  
Neural Network

Activate  
Smart Reply?  
**yes/no**

Generated Replies

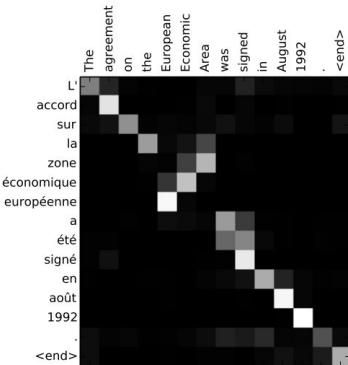
Deep Recurrent  
Neural Network



# Part V: Beyond RNNs

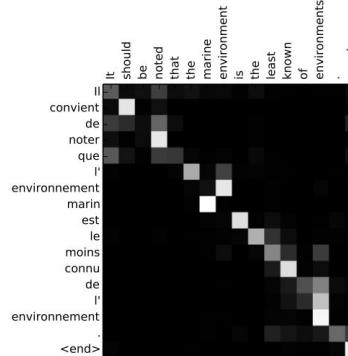
# Attention (Bahdanau et al)

L'accord sur la zone économique européenne a été signé en août 1992.  
<end>

An attention matrix visualization for a French sentence. The matrix is a 10x10 grid where darker shades indicate higher attention weights between words. The words are: L'accord, sur, la, zone, économique, européenne, a, été, signé, en, août, 1992, <end>. The diagonal shows high attention along it, with significant off-diagonal attention from 'zone' to 'économique' and 'européenne', and from 'signé' to 'en'.

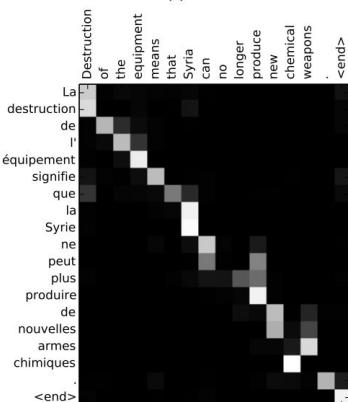
(a)

The agreement on the European Economic Area was signed in August 1992.  
<end>

An attention matrix visualization for an English sentence. The matrix is a 10x10 grid showing attention weights between words. The words are: The, agreement, on, the, European, Economic, Area, was, signed, in, August, 1992, <end>. The attention is primarily on the diagonal, with some cross-attention between 'European' and 'Economic'.

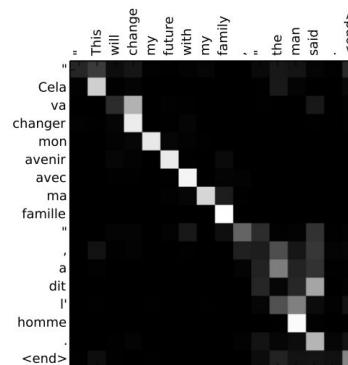
(b)

La destruction de l'équipement signifie que la Syrie ne peut plus produire de nouvelles armes chimiques.  
<end>

An attention matrix visualization for a French sentence. The matrix is a 10x10 grid showing attention weights. The words are: La, destruction, de, l'équipement, signifie, que, la, Syrie, ne, peut, plus, produire, de, nouvelles, armes, chimiques, <end>. The attention is concentrated on the first few words, with 'destruction' and 'équipement' having high attention to each other.

(c)

This will change my future with my family.  
' " Cela va changer mon avenir avec ma famille ' , a dit l'homme  
<end>

An attention matrix visualization for an English sentence. The matrix is a 10x10 grid showing attention weights. The words are: This, will, change, my, future, with, my, family, ' " Cela, va, changer, mon, avenir, avec, ma, famille, ' , a, dit, l'homme, <end>. The attention is on the first few words, with 'change' and 'my' showing significant cross-attention.

(d)

# Visual Attention (Xu et al)

Figure 3. Examples of attending to the correct object (*white* indicates the attended regions, *underlines* indicate the corresponding word)



A woman is throwing a frisbee in a park.



A dog is standing on a hardwood floor.



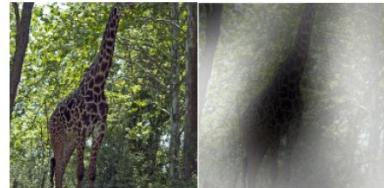
A stop sign is on a road with a mountain in the background.



A little girl sitting on a bed with a teddy bear.

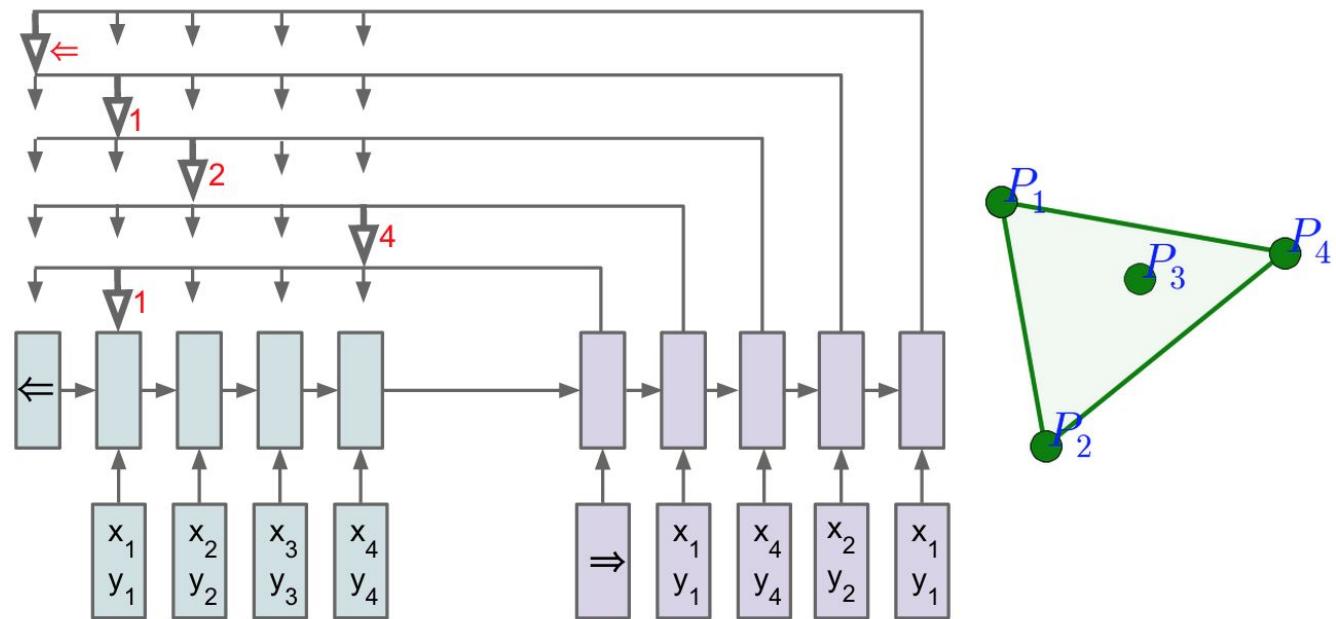


A group of people sitting on a boat in the water.



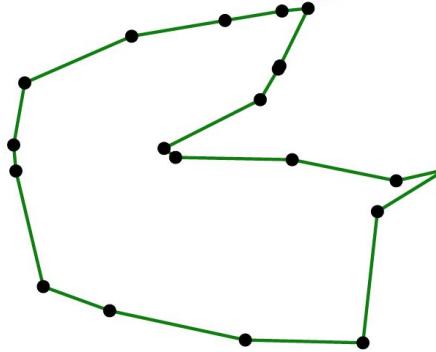
A giraffe standing in a forest with trees in the background.

# Pointer Networks

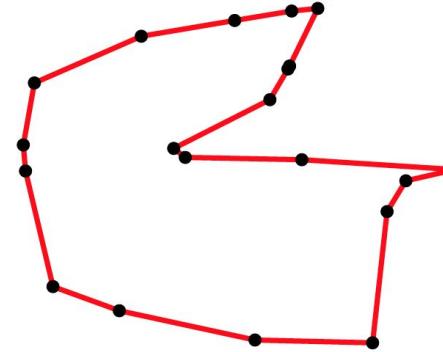


# TSP: NP hard!

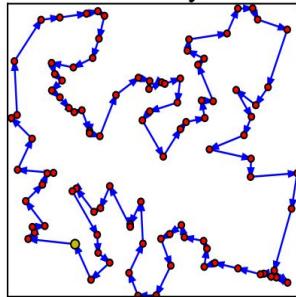
Ground Truth: tour length is 3.518



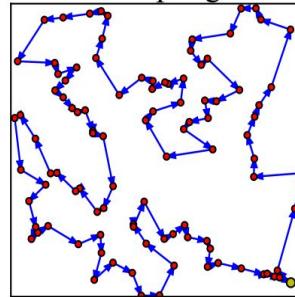
Predictions: tour length is 3.523



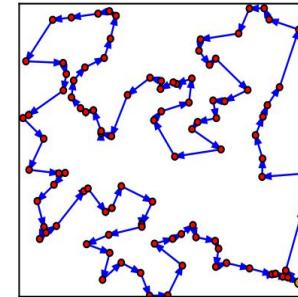
RL pretraining  
-Greedy



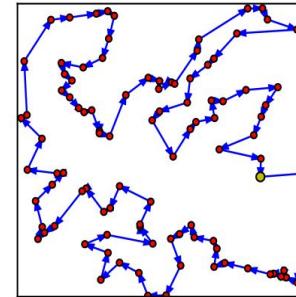
RL pretraining  
-Sampling



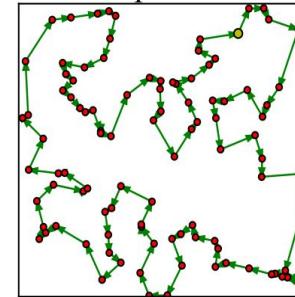
RL pretraining  
-Active Search



Active Search

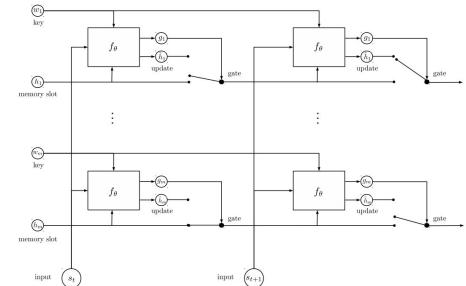
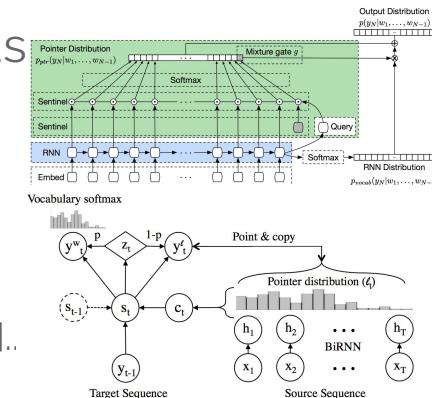


Optimal



# Frontiers: Natural Language Understanding

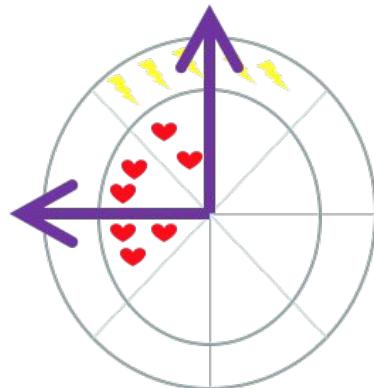
- Good progress on dealing with rarely seen / unknown words  
[Luong et al, 2015][Gulcehre et al, 2016][Merity et al, 2016][Wu et al, 2016]...
- Limited domain (e.g. BABI) tasks increasingly solved  
[Kadlec et al, 2016][Dhingra et al, 2016][Trischler et al, 2016][Henaff et al, 2016]..
- We can't generate long / coherent text
- We can't summarize (or write : )) a book



# Frontiers: Complex Decision Making in Games

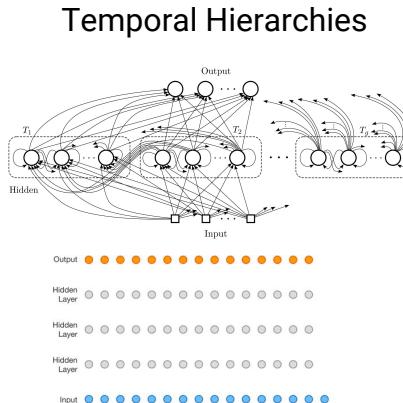
Credits: Berkeley Overmind, AIIDE competition, BWAPI

Moore and Kummerfeld, Berkeley CS294 Class Report, 2011

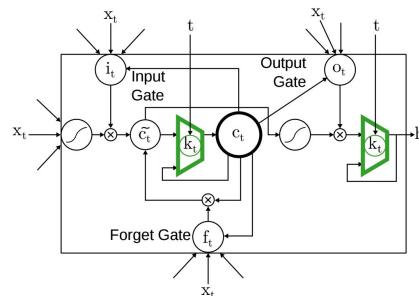


# Common Challenge: Long Term Dependencies

- Thousands / Millions of steps
  - Recursive programs
  - Learning to Learn
  - Complex Decision Making
  - NLU



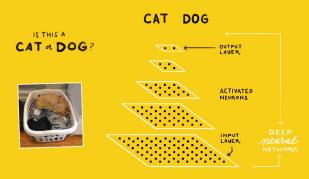
Jaderberg et al, 2016



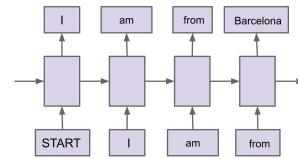
Neil et al, 2016

# The Deep Learning Toolbox

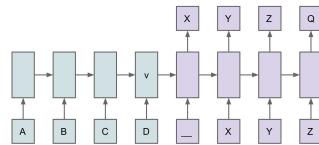
## Feed forward models



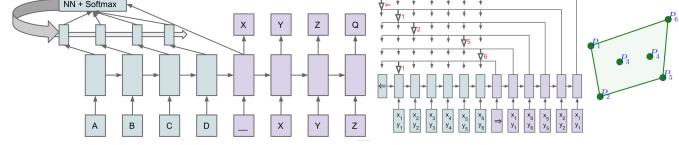
## Sequence Prediction



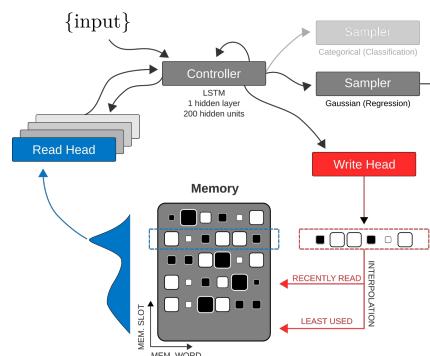
## Seq2Seq



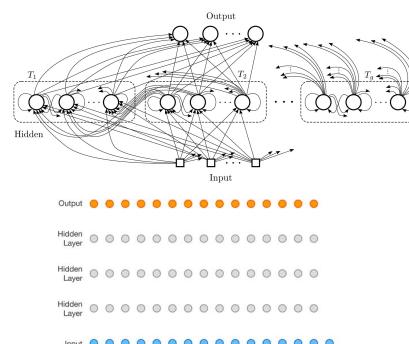
## Attention/Pointers



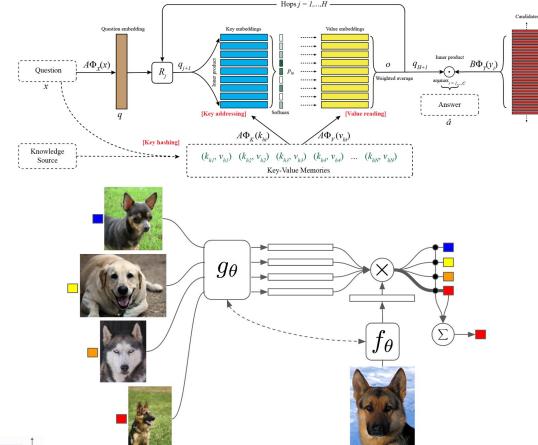
## Read/Write memories



## Temporal Hierarchies



## Key,Value memories



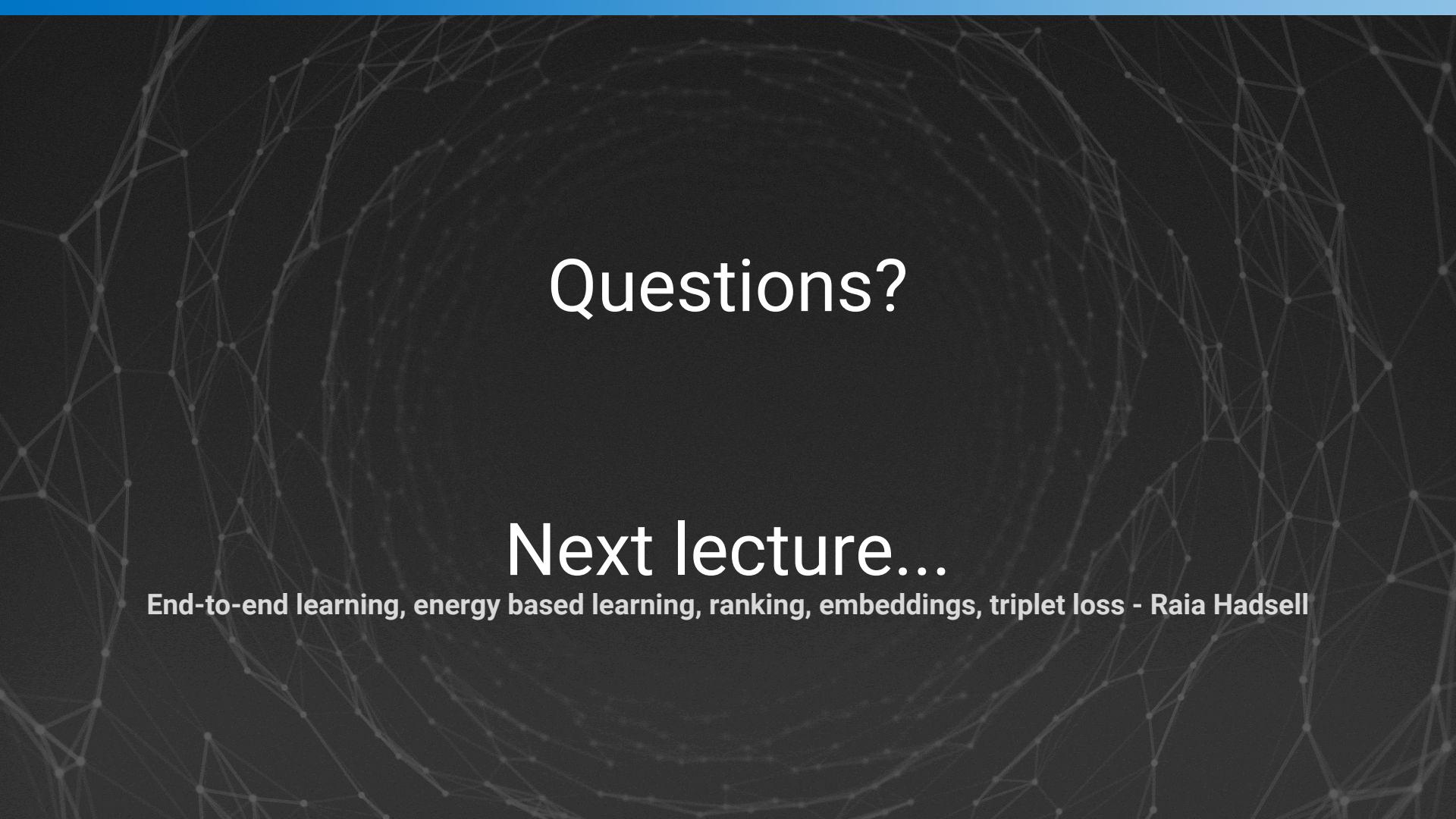
## Recurrent Architectures



Figure credits: Jeff Dean, Chris Olah, Santoro et al 2016, Koutnik et al 2014, van den Oord et al 2016, Miller et al 2016, Vinyals et al 2016

# Summary

- Sequences are important, everything is a sequence
- Recurrent Nets are a way forward, but they pose some challenges
  - Vanishing/Exploding gradients
  - Expensive to compute
- Conditional models are very powerful, and many applications “just worked”  
when seq2seq and similar were trained on large amounts of data
- Attention/memory, longer term memory... many challenges remain!

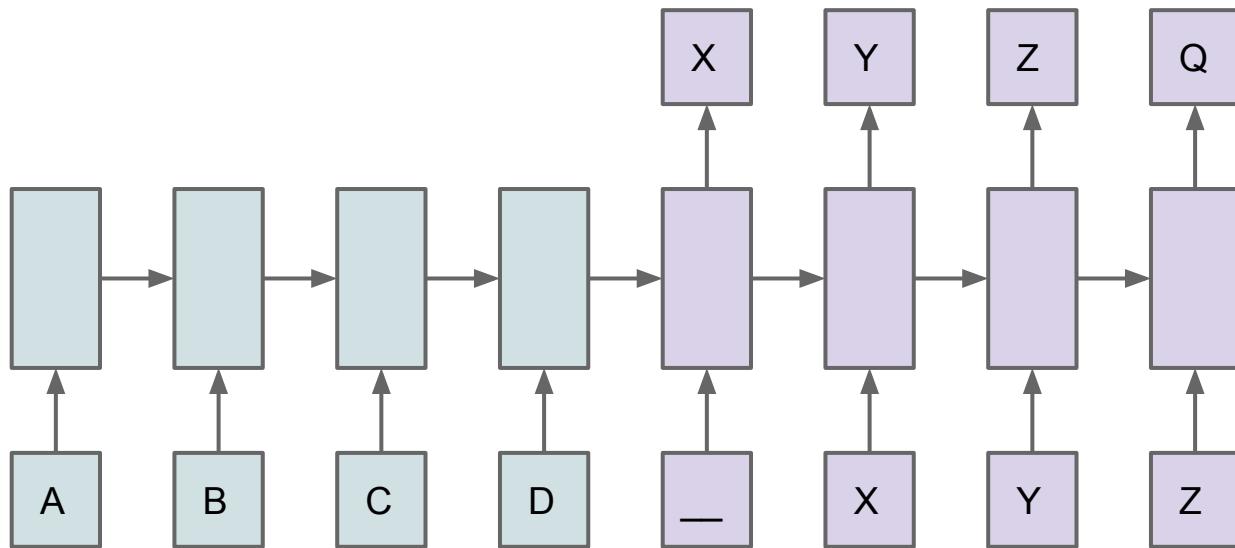


# Questions?

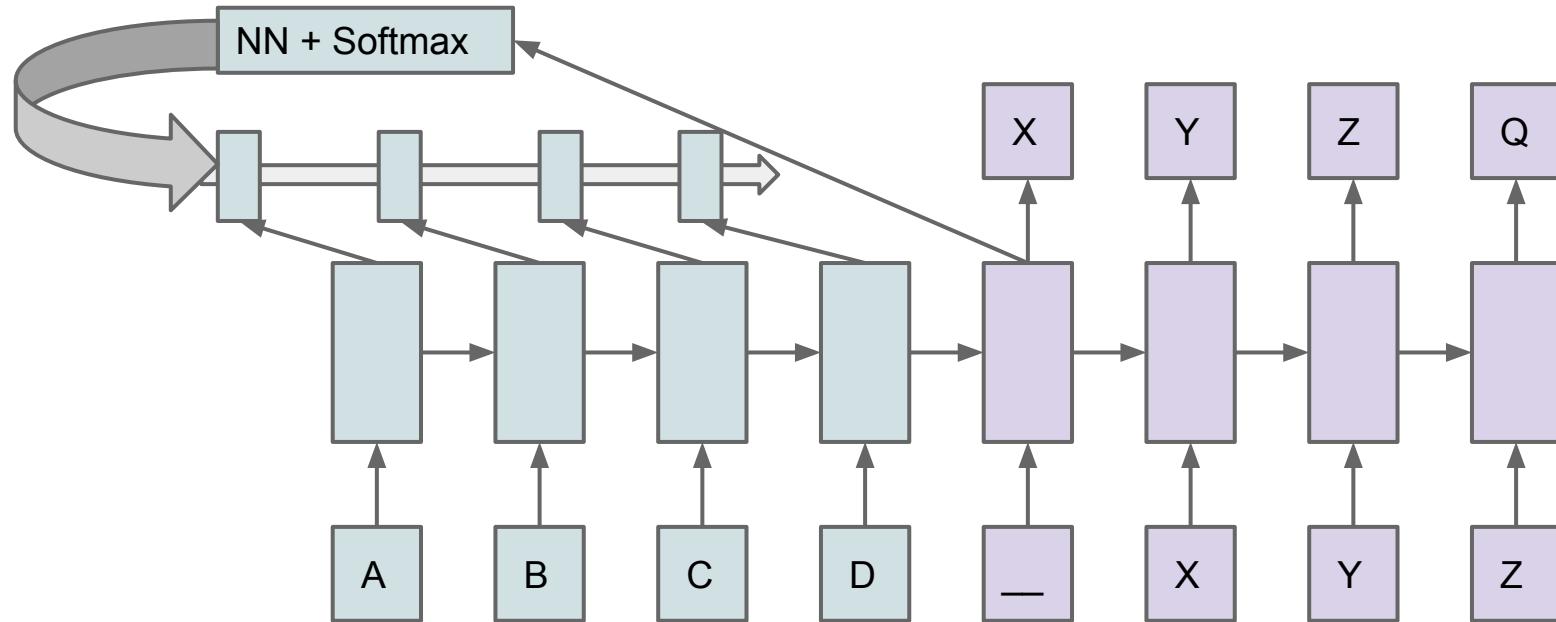
## Next lecture...

End-to-end learning, energy based learning, ranking, embeddings, triplet loss - Raia Hadsell

# Attention



# Attention



# Position Based Attention

Inputs: "I am a cat."

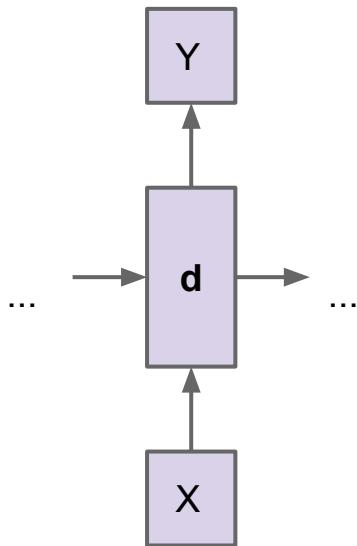
Input RNN states:  $\mathbf{e}_1 \mathbf{e}_2 \mathbf{e}_3 \mathbf{e}_4$

Decoder RNN state:  $\mathbf{d}$

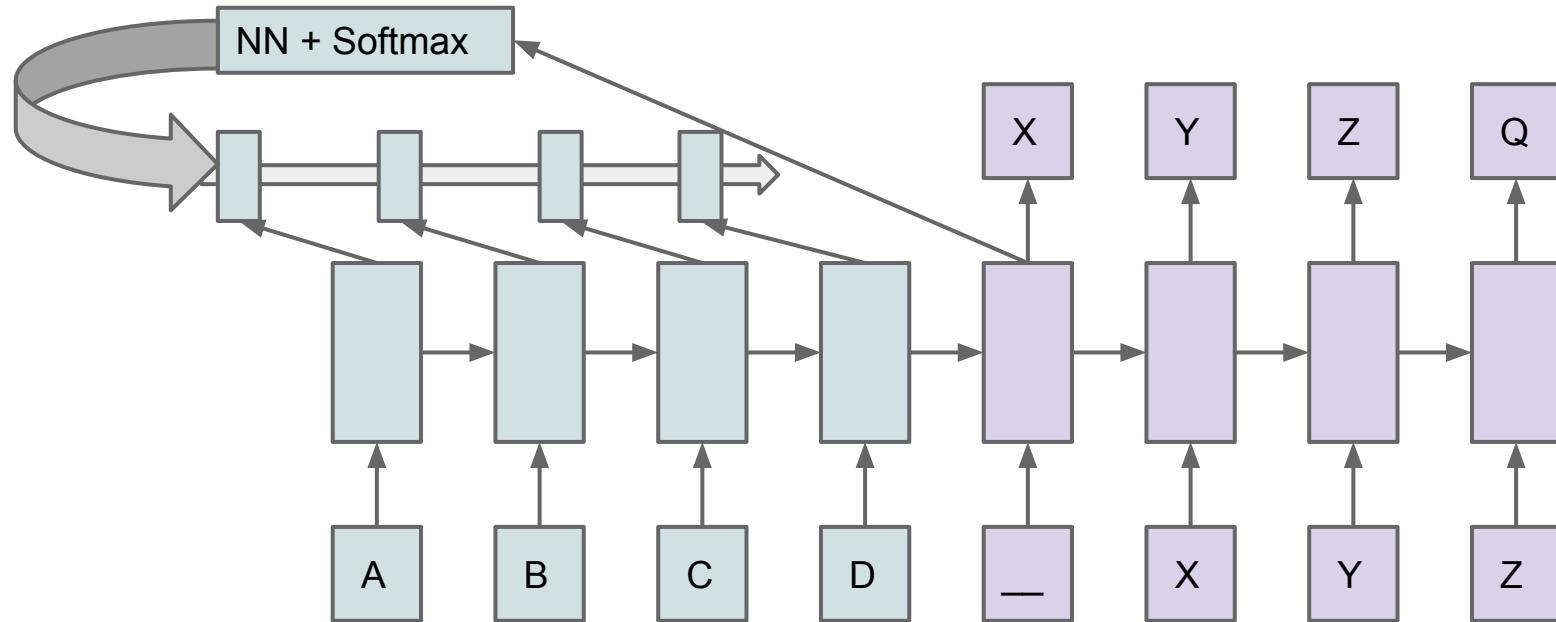
Map  $\mathbf{d}$  to  $\mathbf{a} = [0.1 \ 0.9 \ 0.0 \ 0.0]$  E.g.  $\mathbf{a} = s(\mathbf{Wd})$

Compute:  $\mathbf{d}' = 0.1^*\mathbf{e}_1 + 0.9^*\mathbf{e}_2 + 0^*\mathbf{e}_3 + 0^*\mathbf{e}_4$

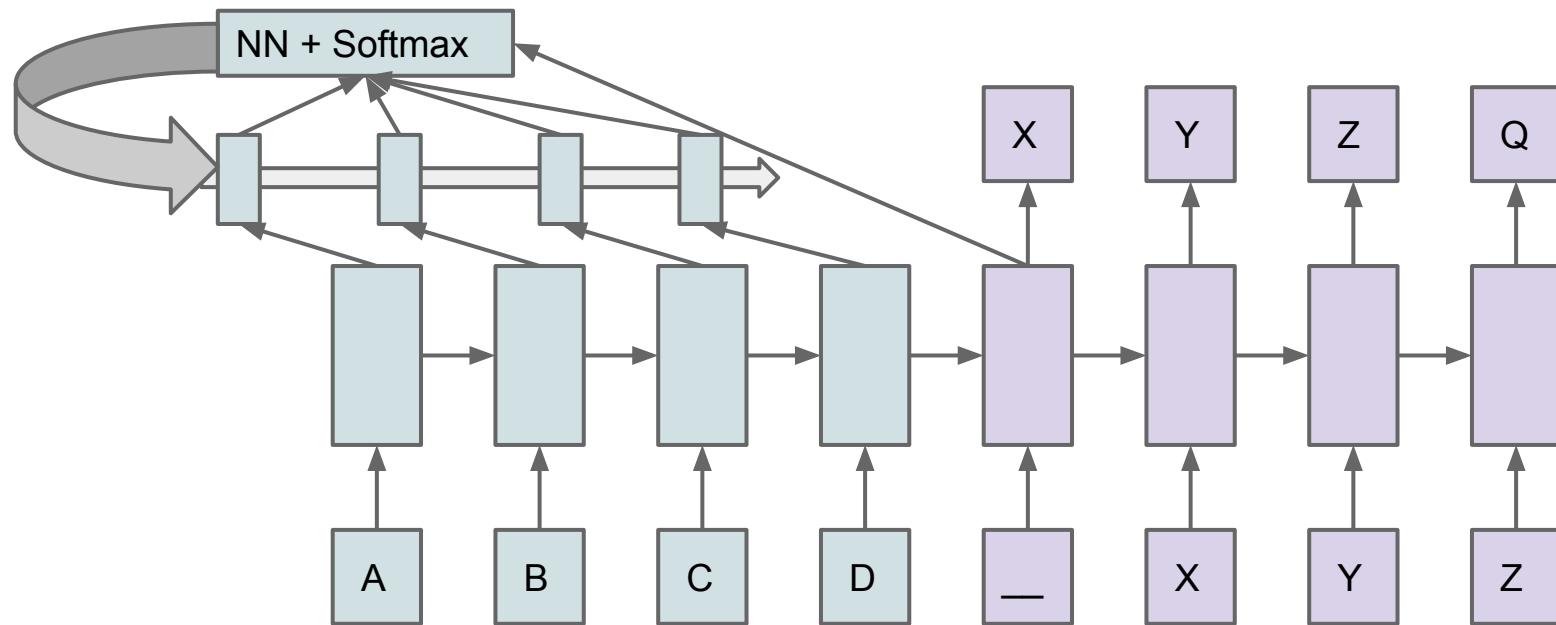
Set  $\mathbf{d} = [\mathbf{d} \ \mathbf{d}']$



# Position Based Attention



# Content Based Attention



# Content Based Attention

Inputs: "I am a cat."

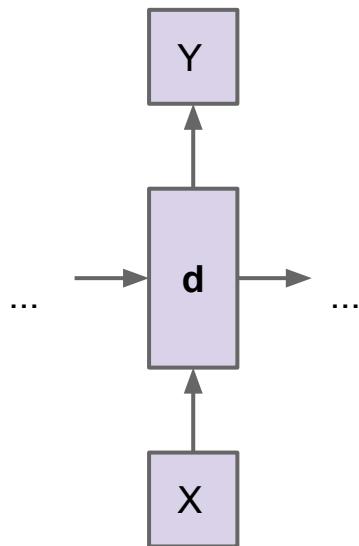
Input RNN states:  $\mathbf{e}_1 \mathbf{e}_2 \mathbf{e}_3 \mathbf{e}_4$

Decoder RNN state:  $\mathbf{d}$

Compute  $\mathbf{d}^T \mathbf{e}_i$ : [0.0 0.05 0.9 0.05]

Compute:  $\mathbf{d}' = 0^* \mathbf{e}_1 + 0.05^* \mathbf{e}_2 + 0.9^* \mathbf{e}_3 + 0.05^* \mathbf{e}_4$

Set  $\mathbf{d} = [\mathbf{d} \mathbf{d}']$



# Content Based Attention

Attention [Bahdanau, Cho and Bengio, 2014]

$$u_j = v^T \tanh(W_1 e_j + W_2 d) \quad j \in (1, \dots, n)$$

$$a_j = \text{softmax}(u_j) \quad j \in (1, \dots, n)$$

$$d' = \sum_{j=1}^n a_j e_j$$

## Neural Turing Machine & Memory Networks:

Given a **key** vector  $\mathbf{k}$  and a **strength** scalar  $s$  emitted by the controller, get a weighting  $(w_1, \dots, w_N)$  over memory locations where

$$w_i = \frac{e^{\hat{w}_i}}{\sum_j e^{\hat{w}_j}} \quad \hat{w}_i = \log(1 + e^s) C(\mathbf{k}, \mathbf{m}_i)$$
$$C(\mathbf{k}, \mathbf{m}_i) = \frac{\mathbf{k} \cdot \mathbf{m}_i}{\|\mathbf{k}\| \|\mathbf{m}_i\|}$$

and  $C(\mathbf{k}, \mathbf{m}_i)$  is the **cosine similarity** between the key and the **word**  $\mathbf{m}_i$  at memory location  $i$

# Frontiers: Learning Programs

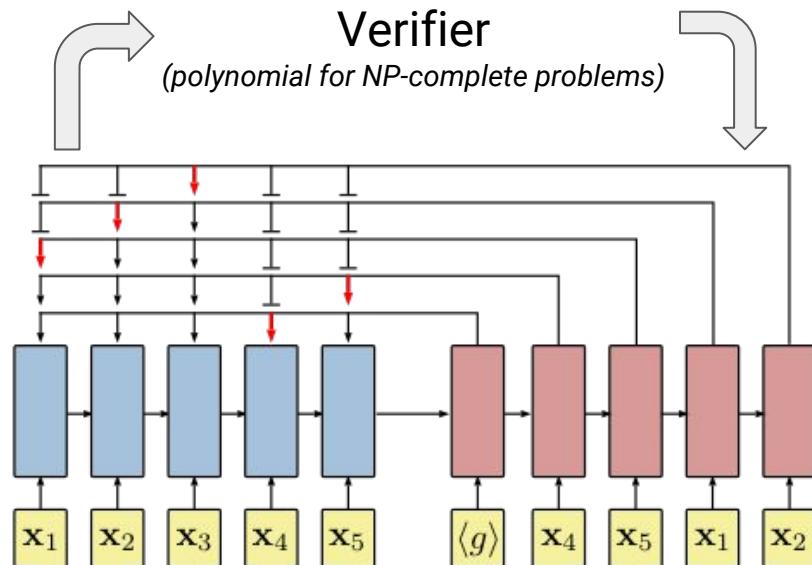
LOTS of related work: [Zaremba et al, 2014][Reed et al, 2015][Neelakantan et al, 2015][Vinyals et al, 2015]...

RL-PtrNets [Bello, Pham, et al 2016]

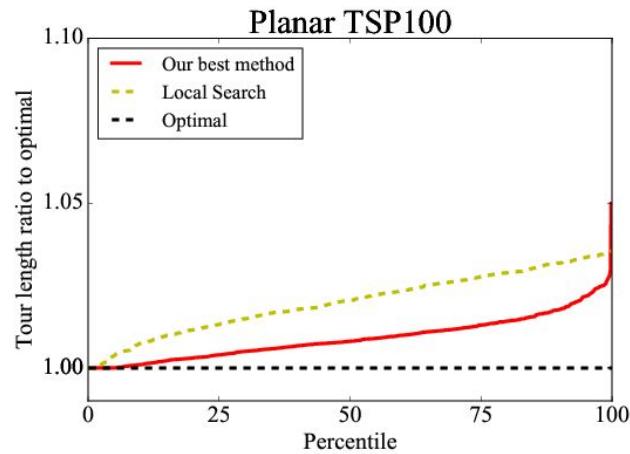
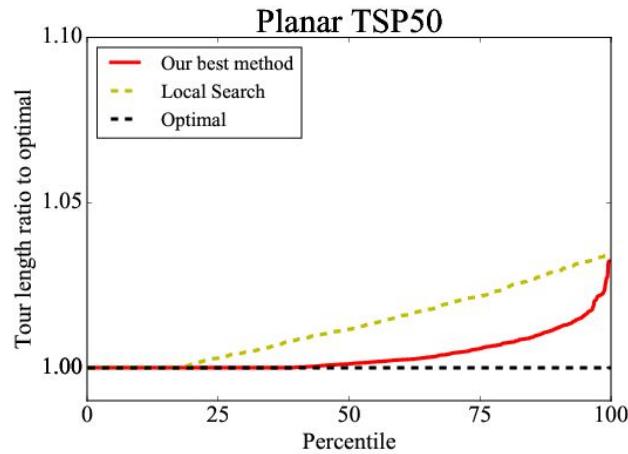
Actor-critic training using the expected combinatorial optimization objective

Different configurations:

- *Train* on training set
- *Search* at inference time
- *Refine* parameters at inference time



# Frontiers: Learning Programs



Average tour lengths on TSP100 ( $\sim 10^{157}$  solutions)

Christofides: 9.18 | Local search: 7.99 | RL-PtrNets: 7.83 | Optimal: 7.77

# Frontiers: Complex Decision Making in Games

Credits: Berkeley Overmind, AIIDE competition, BWAPI

Predicting Game Futures with RNNs [Vinyals, 2010, unpublished]

