

Deep Learning Course

Lecture 7: Recurrent Neural Networks

Valentin Malykh

Based on Stanford CS224d

Overview

- Traditional language models
- RNNs
- RNN language models
- Important training problems and tricks
 - Vanishing and exploding gradient problems
- RNNs for other sequence tasks
- Bidirectional and deep RNNs

Language Models

A language model computes a probability for a sequence of words: $P(w_1, \dots, w_T)$

- Useful for machine translation
 - Word ordering:
 $p(\text{the cat is small}) > p(\text{small the is cat})$
 - Word choice:
 $p(\text{walking home after school}) > p(\text{walking house after school})$

Traditional Language Models

- Probability is usually conditioned on window of n previous words
- An incorrect but necessary Markov assumption!

$$P(w_1, \dots, w_m) = \prod_{i=1}^m P(w_i | w_1, \dots, w_{i-1}) \approx \prod_{i=1}^m P(w_i | w_{i-(n-1)}, \dots, w_{i-1})$$

- To estimate probabilities, compute for unigrams and bigrams (conditioning on one/two previous word(s):

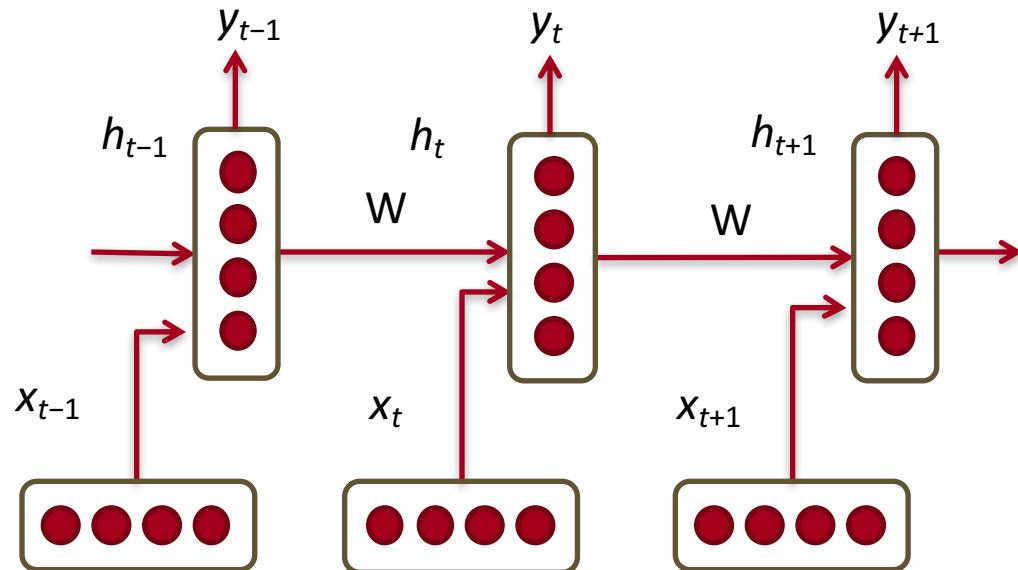
$$p(w_2|w_1) = \frac{\text{count}(w_1, w_2)}{\text{count}(w_1)} \quad p(w_3|w_1, w_2) = \frac{\text{count}(w_1, w_2, w_3)}{\text{count}(w_1, w_2)}$$

Traditional Language Models

- Performance improves with keeping around higher n-grams counts and doing smoothing and so-called backoff (e.g. if 4-gram not found, try 3-gram, etc)
- There are A LOT of n-grams!
→ Gigantic RAM requirements!
- Recent state of the art: *Scalable Modified Kneser-Ney Language Model Estimation* by Heafield et al.:
“Using one machine with 140 GB RAM for 2.8 days, we built an unpruned model on 126 billion tokens”

Recurrent Neural Networks!

- RNNs tie the weights at each time step
- Condition the neural network on all previous words
- RAM requirement only scales with number of words



Recurrent Neural Network language model

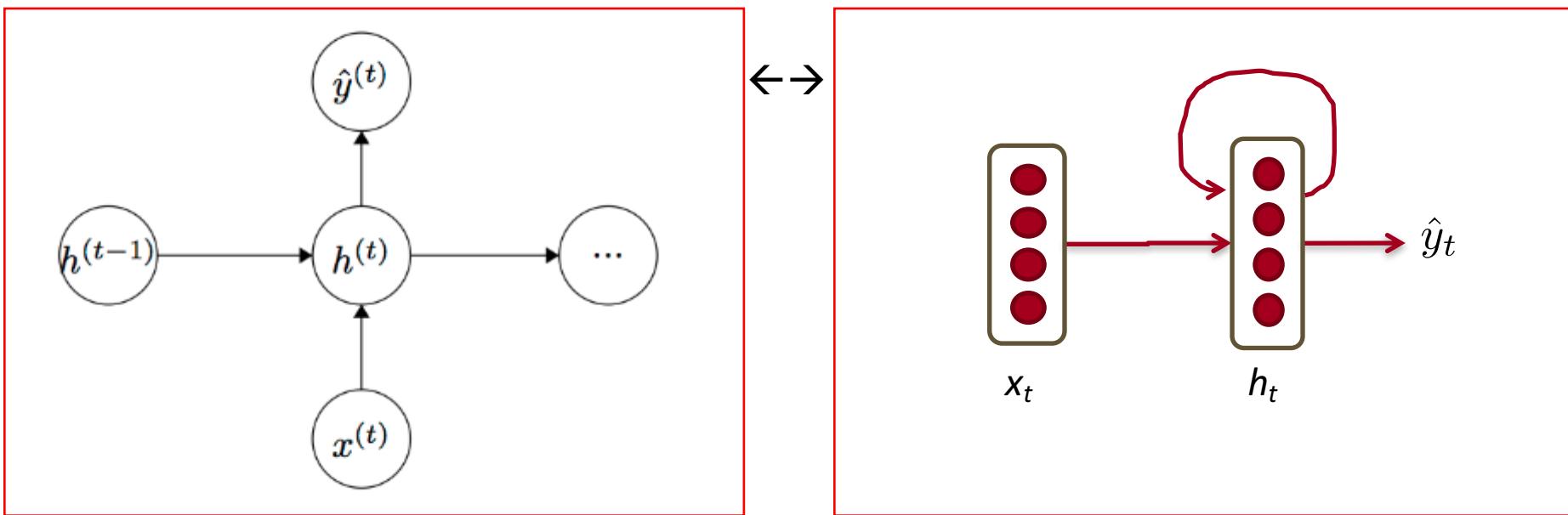
Given list of word **vectors**: $x_1, \dots, x_{t-1}, x_t, x_{t+1}, \dots, x_T$

At a single time step:

$$h_t = \sigma \left(W^{(hh)} h_{t-1} + W^{(hx)} x_{[t]} \right)$$

$$\hat{y}_t = \text{softmax} \left(W^{(S)} h_t \right)$$

$$\hat{P}(x_{t+1} = v_j \mid x_t, \dots, x_1) = \hat{y}_{t,j}$$



Recurrent Neural Network language model

Main idea: we use the same set of W weights at all time steps!

Everything else is the same:

$$\begin{aligned} h_t &= \sigma \left(W^{(hh)} h_{t-1} + W^{(hx)} x_{[t]} \right) \\ \hat{y}_t &= \text{softmax} \left(W^{(S)} h_t \right) \\ \hat{P}(x_{t+1} = v_j \mid x_t, \dots, x_1) &= \hat{y}_{t,j} \end{aligned}$$

$h_0 \in \mathbb{R}^{D_h}$ is some initialization vector for the hidden layer at time step 0

$x_{[t]}$ is the column vector of L at index $[t]$ at time step t

$$W^{(hh)} \in \mathbb{R}^{D_h \times D_h} \quad W^{(hx)} \in \mathbb{R}^{D_h \times d} \quad W^{(S)} \in \mathbb{R}^{|V| \times D_h}$$

Recurrent Neural Network language model

$\hat{y} \in \mathbb{R}^{|V|}$ is a probability distribution over the vocabulary

Same cross entropy loss function but predicting words instead of classes

$$J^{(t)}(\theta) = - \sum_{j=1}^{|V|} y_{t,j} \log \hat{y}_{t,j}$$

Recurrent Neural Network language model

Evaluation could just be negative of average log probability over dataset of size (number of words) T:

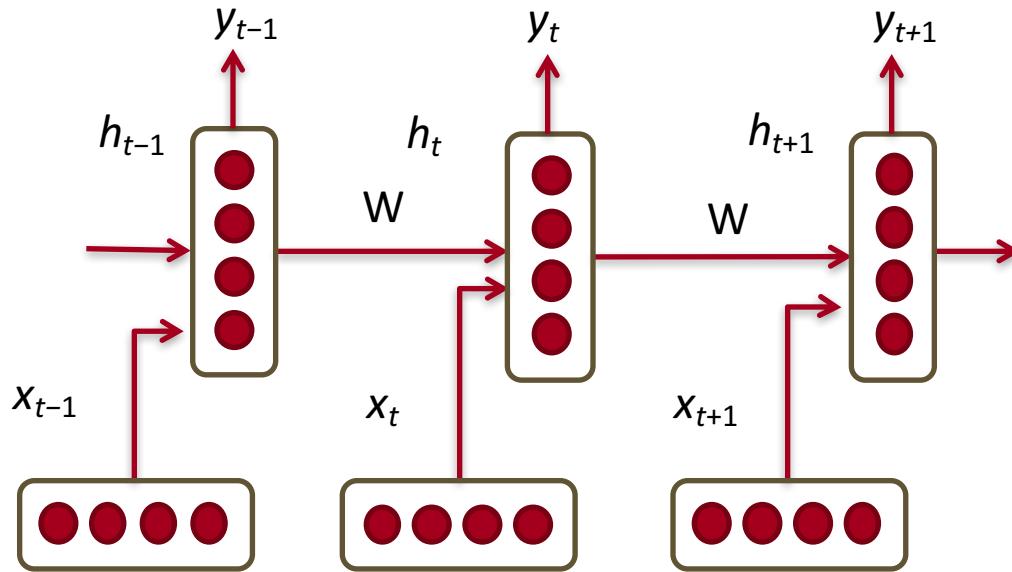
$$J = -\frac{1}{T} \sum_{t=1}^T \sum_{j=1}^{|V|} y_{t,j} \log \hat{y}_{t,j}$$

But more common: Perplexity: 2^J

Lower is better!

Training RNNs is hard

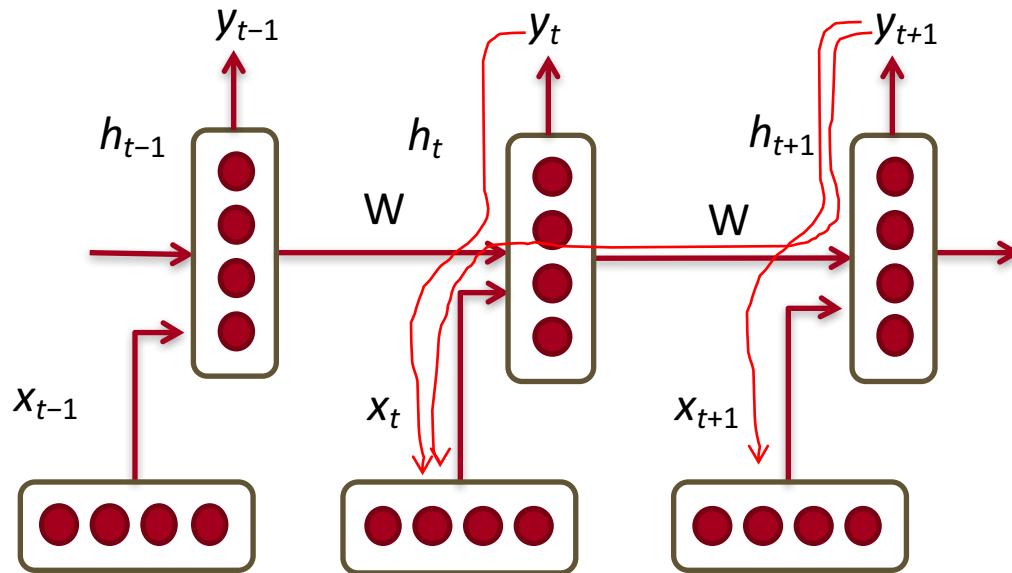
- Multiply the same matrix at each time step during forward prop



- Ideally inputs from many time steps ago can modify output y
- Take $\frac{\partial E_2}{\partial W}$ for an example RNN with 2 time steps! Insightful!

The vanishing/exploding gradient problem

- Multiply the same matrix at each time step during backprop



The vanishing gradient problem - Details

- Similar but simpler RNN formulation:

$$\begin{aligned} h_t &= Wf(h_{t-1}) + W^{(hx)}x_{[t]} \\ \hat{y}_t &= W^{(S)}f(h_t) \end{aligned}$$

- Total error is the sum of each error at time steps t

$$\frac{\partial E}{\partial W} = \sum_{t=1}^T \frac{\partial E_t}{\partial W}$$

- Hardcore chain rule application:

$$\frac{\partial E_t}{\partial W} = \sum_{k=1}^t \frac{\partial E_t}{\partial y_t} \frac{\partial y_t}{\partial h_t} \frac{\partial h_t}{\partial h_k} \frac{\partial h_k}{\partial W}$$

The vanishing gradient problem - Details

- Similar to backprop but less efficient formulation
- Useful for analysis we'll look at:

$$\frac{\partial E_t}{\partial W} = \sum_{k=1}^t \frac{\partial E_t}{\partial y_t} \frac{\partial y_t}{\partial h_t} \boxed{\frac{\partial h_t}{\partial h_k}} \frac{\partial h_k}{\partial W}$$

- Remember: $h_t = Wf(h_{t-1}) + W^{(hx)}x_{[t]}$
- More chain rule, remember:

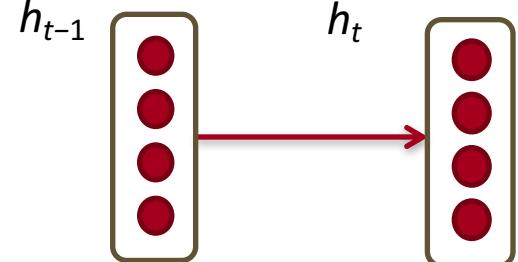
$$\frac{\partial h_t}{\partial h_k} = \prod_{j=k+1}^t \frac{\partial h_j}{\partial h_{j-1}}$$

- Each partial is a Jacobian:

$$\frac{d\mathbf{f}}{d\mathbf{x}} = \begin{bmatrix} \frac{\partial \mathbf{f}}{\partial x_1} & \dots & \frac{\partial \mathbf{f}}{\partial x_n} \end{bmatrix} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \dots & \frac{\partial f_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \dots & \frac{\partial f_m}{\partial x_n} \end{bmatrix}$$

The vanishing gradient problem - Details

- From previous slide: $\frac{\partial h_t}{\partial h_k} = \prod_{j=k+1}^t \frac{\partial h_j}{\partial h_{j-1}}$
- Remember: $h_t = Wf(h_{t-1}) + W^{(hx)}x_{[t]}$



- To compute Jacobian, derive each element of matrix: $\frac{\partial h_{j,m}}{\partial h_{j-1,n}}$

$$\frac{\partial h_t}{\partial h_k} = \prod_{j=k+1}^t \frac{\partial h_j}{\partial h_{j-1}} = \prod_{j=k+1}^t W^T \text{diag}[f'(h_{j-1})]$$

- Where: $\text{diag}(z) = \begin{pmatrix} z_1 & & & 0 \\ & z_2 & & \\ & & \ddots & \\ 0 & & & z_{n-1} & z_n \end{pmatrix}$

Check at home
that you understand
the diag matrix
formulation

The vanishing gradient problem - Details

- Analyzing the norms of the Jacobians, yields:

$$\left\| \frac{\partial h_j}{\partial h_{j-1}} \right\| \leq \|W^T\| \|\text{diag}[f'(h_{j-1})]\| \leq \beta_W \beta_h$$

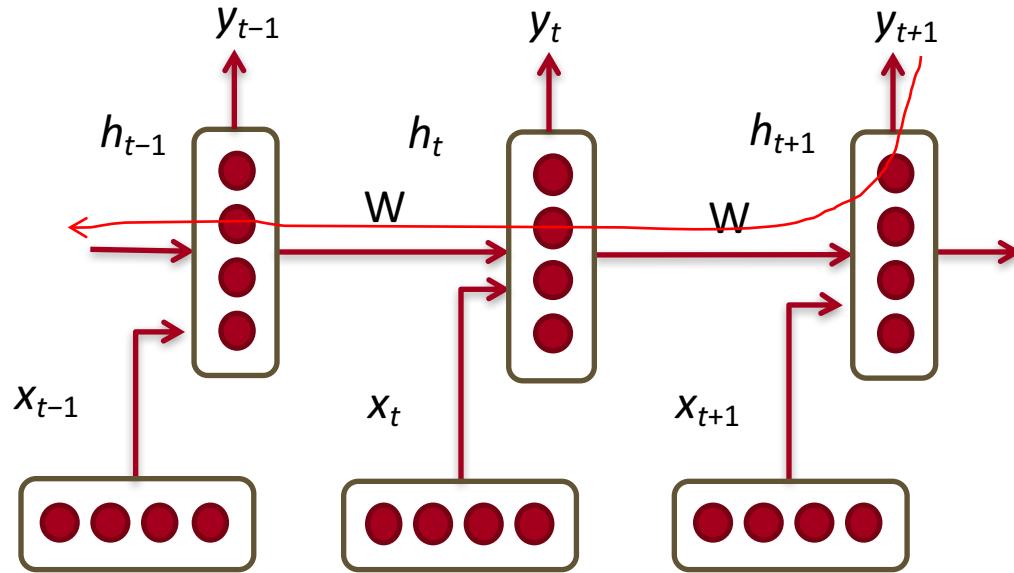
- Where we defined $\bar{\cdot}$'s as upper bounds of the norms
- The gradient is a product of Jacobian matrices, each associated with a step in the forward computation.

$$\left\| \frac{\partial h_t}{\partial h_k} \right\| = \left\| \prod_{j=k+1}^t \frac{\partial h_j}{\partial h_{j-1}} \right\| \leq (\beta_W \beta_h)^{t-k}$$

- This can become very small or very large quickly [Bengio et al 1994], and the locality assumption of gradient descent breaks down. → **Vanishing or exploding gradient**

Why is the vanishing gradient a problem?

- The error at a time step ideally can tell a previous time step from many steps away to change during backprop



The vanishing gradient problem for language models

- In the case of language modeling or question answering words from time steps far away are not taken into consideration when training to predict the next word
- Example:

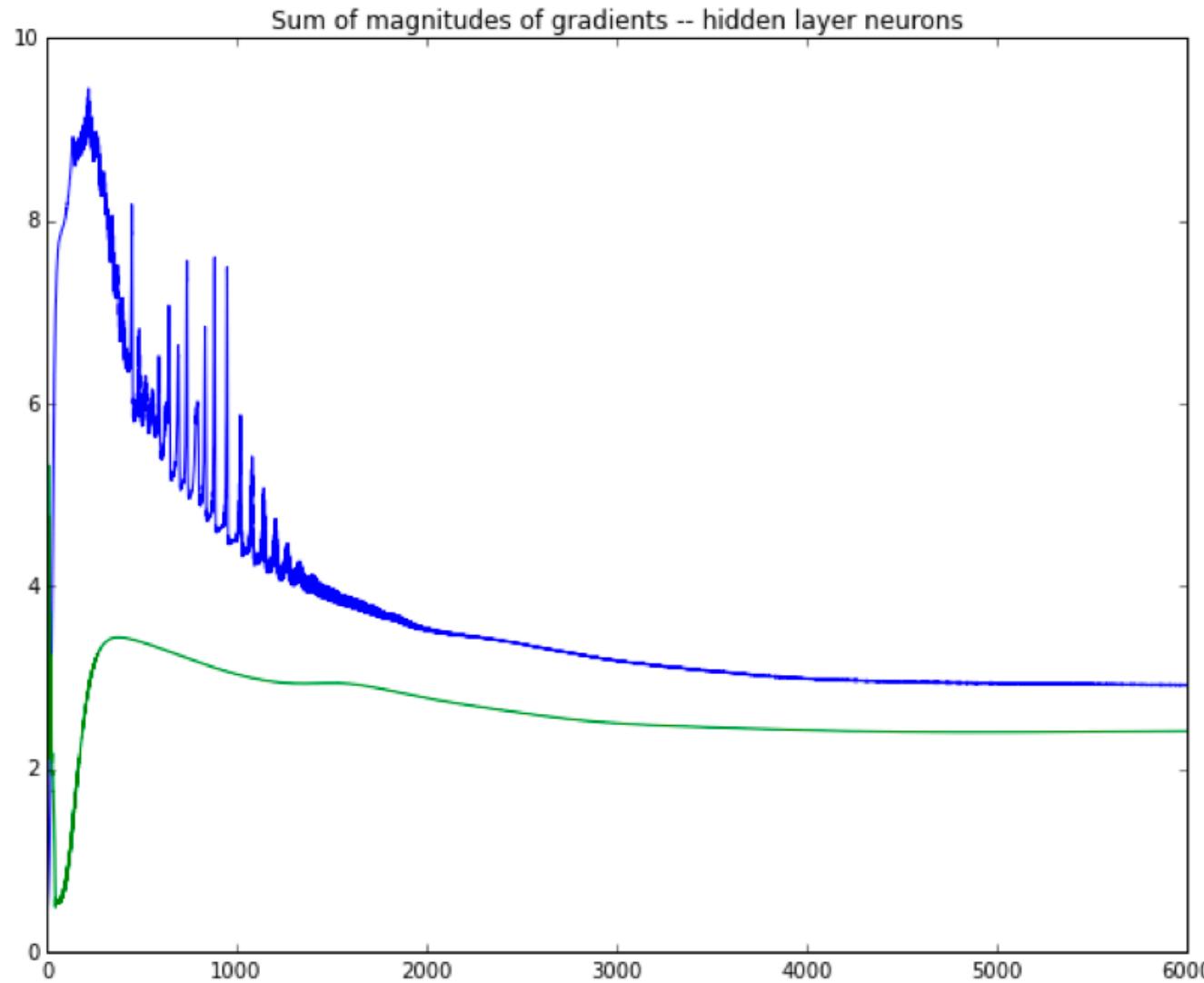
Jane walked into the room. John walked in too. It was late in the day. Jane said hi to _____

IIPython Notebook with vanishing gradient example

- Example of simple and clean NNet implementation
- Comparison of sigmoid and ReLu units
- A little bit of vanishing gradient

```
In [21]: plt.plot(np.array(rect_array[:6000]),color='blue')
plt.plot(np.array(sigm_array[:6000]),color='green')
plt.title('Sum of magnitudes of gradients -- hidden layer neurons')
```

```
Out[21]: <matplotlib.text.Text at 0x10a331310>
```



Trick for exploding gradient: clipping trick

- The solution first introduced by Mikolov is to clip gradients to a maximum value.

Algorithm 1 Pseudo-code for norm clipping the gradients whenever they explode

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

- Makes a big difference in RNNs.

Gradient clipping intuition

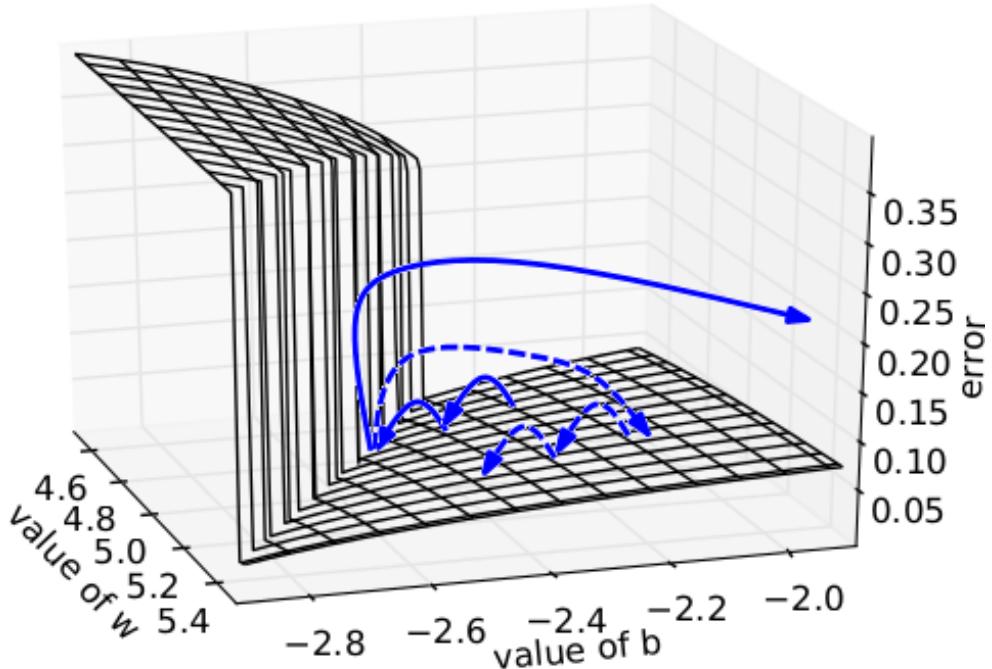
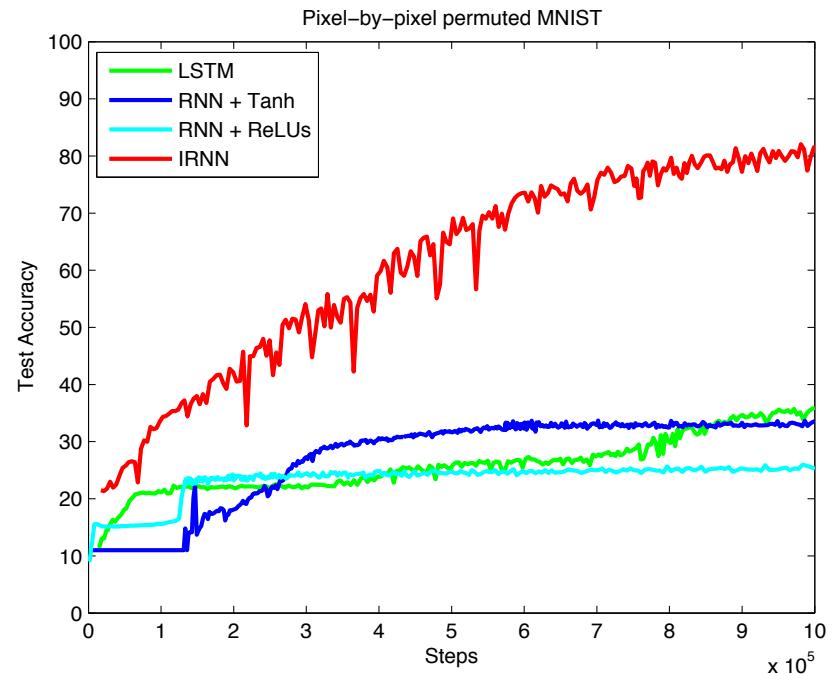


Figure from paper:
On the difficulty of
training Recurrent Neural
Networks, Pascanu et al.
2013

- Error surface of a single hidden unit RNN,
- High curvature walls
- Solid lines: standard gradient descent trajectories
- Dashed lines gradients rescaled to fixed size

For vanishing gradients: Initialization + ReLUs!

- Initialize $W^{(*)}$'s to identity matrix I and $f(z) = \text{rect}(z) = \max(z, 0)$
- → Huge difference!



- Initialization idea first introduced in *Parsing with Compositional Vector Grammars*, Socher et al. 2013
- New experiments with recurrent neural nets 2 weeks ago (!) in *A Simple Way to Initialize Recurrent Networks of Rectified Linear Units*, Le et al. 2015

Perplexity Results

KN5 = Count-based language model with Kneser-Ney smoothing & 5-grams

Table 2. Comparison of different neural network architectures on Penn Corpus (1M words) and Switchboard (4M words).

Model	Penn Corpus		Switchboard	
	NN	NN+KN	NN	NN+KN
KN5 (baseline)	-	141	-	92.9
feedforward NN	141	118	85.1	77.5
RNN trained by BP	137	113	81.3	75.4
RNN trained by BPTT	123	106	77.5	72.5

Table from paper *Extensions of recurrent neural network language model* by Mikolov et al 2011

Problem: Softmax is huge and slow

Trick: Class-based word prediction

$$\begin{aligned} p(w_t | \text{history}) &= p(c_t | \text{history})p(w_t | c_t) \\ &= p(c_t | h_t)p(w_t | c_t) \end{aligned}$$

Table 3. Perplexities on Penn corpus with factorization of the output layer by the class model. All models have the same basic configuration (200 hidden units and BPTT=5). The Full model is a baseline and does not use classes, but the whole 10K vocabulary.

Classes	RNN	RNN+KN5	Min/epoch	Sec/test
30	134	112	12.8	8.8
50	136	114	9.8	6.7
100	136	114	9.1	5.6
200	136	113	9.5	6.0
400	134	112	10.9	8.1
1000	131	111	16.1	15.7
2000	128	109	25.3	28.7
4000	127	108	44.4	57.8
6000	127	109	70	96.5
8000	124	107	107	148
Full	123	106	154	212

One last implementation trick

- You only need to pass backwards through your sequence once and accumulate all the deltas from each E_t

Sequence modeling for other tasks

- Classify each word into:
 - NER
 - Entity level sentiment in context
 - opinionated expressions
- Example application and slides from paper *Opinion Mining with Deep Recurrent Nets* by Irsoy and Cardie 2014

Opinion Mining with Deep Recurrent Nets

Goal: Classify each word as

direct subjective expressions (DSEs) and
expressive subjective expressions (ESEs).

DSE: Explicit mentions of private states or speech events expressing private states

ESE: Expressions that indicate sentiment, emotion, etc. without explicitly conveying them.

Example Annotation

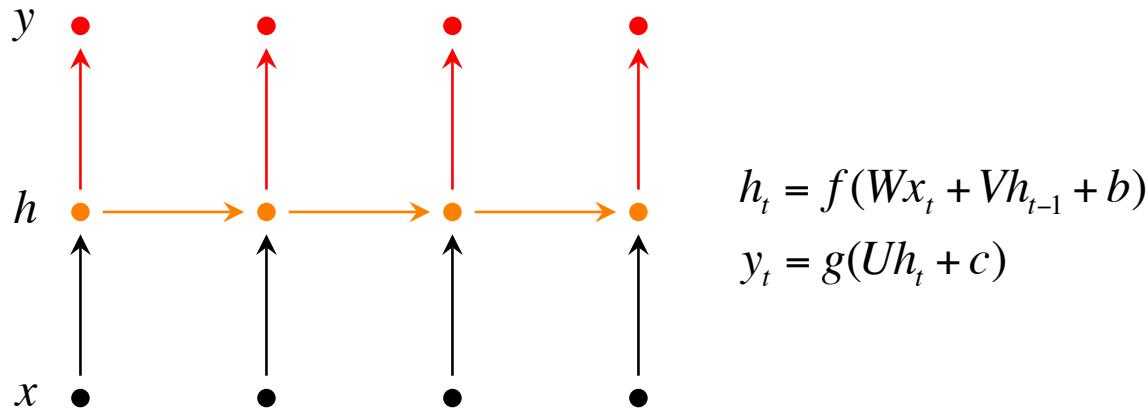
In BIO notation (tags either begin-of-entity (`B_X`) or continuation-of-entity (`I_X`)):

The committee, [as usual]_{ESE}, [has refused to make any statements]_{DSE}.

The	committee	,	as	usual	,	has
O	O	O	<code>B_ESE</code>	<code>I_ESE</code>	O	<code>B_DSE</code>
refused	to	make	any	statements	.	
<code>I_DSE</code>	<code>I_DSE</code>	<code>I_DSE</code>	<code>I_DSE</code>	<code>I_DSE</code>	O	

Approach: Recurrent Neural Network

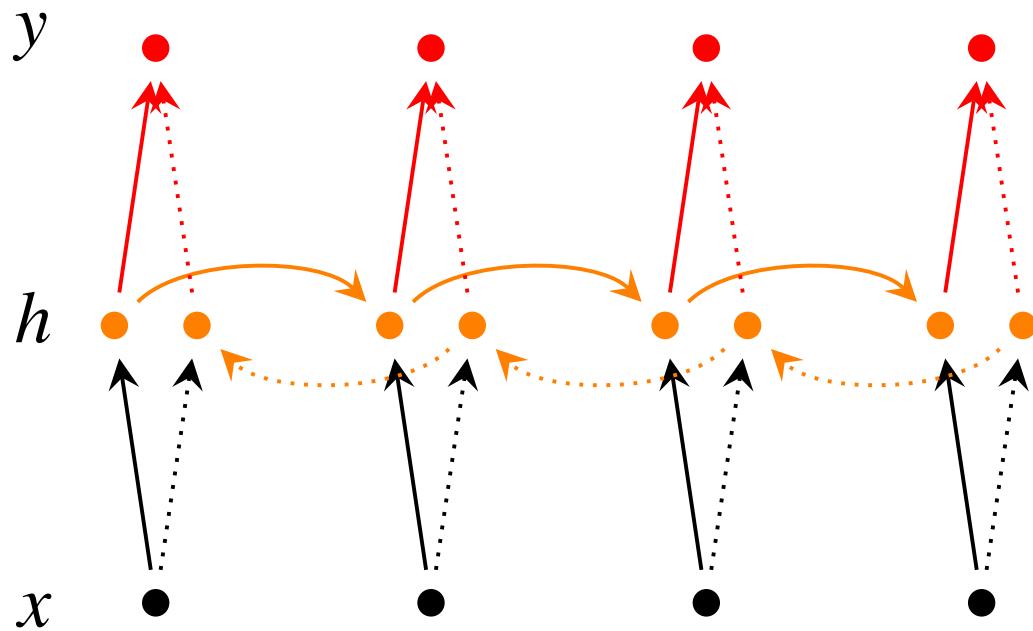
- Notation from paper (so you get used to different ones)



- x represents a token (word) as a vector.
- y represents the output label (B, I or O) – g = softmax !
- h is the memory, computed from the past memory and current word. It summarizes the sentence up to that time.

Bidirectional RNNs

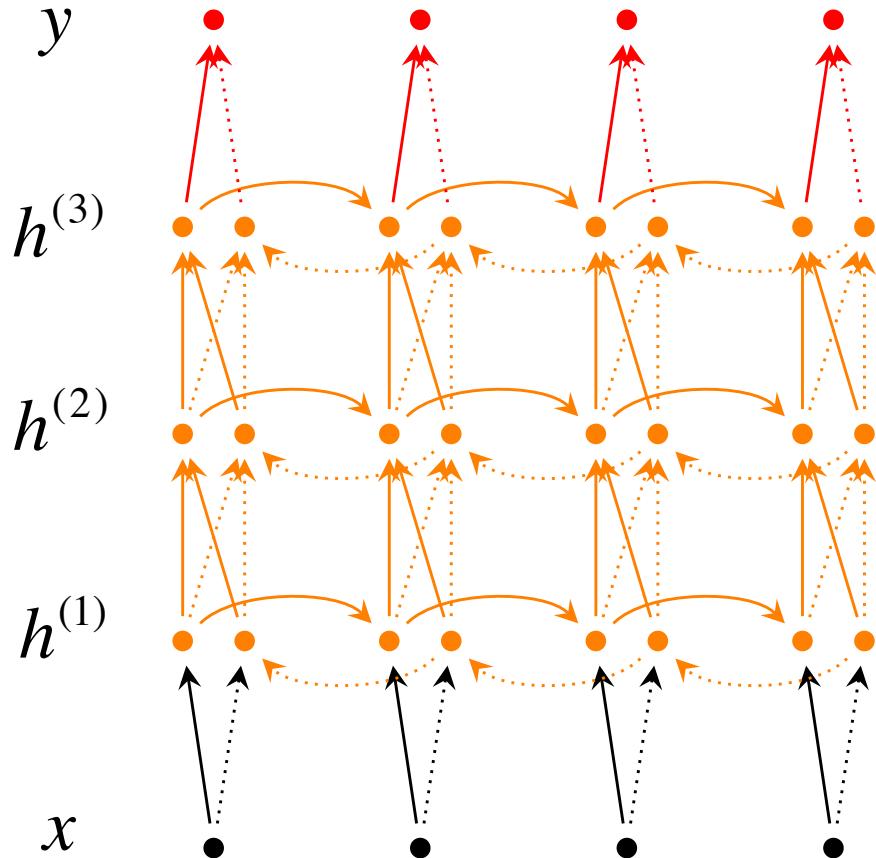
Problem: For classification you want to incorporate information from words both preceding and following



$$\begin{aligned}\vec{h}_t &= f(\vec{W}x_t + \vec{V}\vec{h}_{t-1} + \vec{b}) \\ \overleftarrow{h}_t &= f(\overleftarrow{W}x_t + \overleftarrow{V}\overleftarrow{h}_{t+1} + \overleftarrow{b}) \\ y_t &= g(U[\vec{h}_t; \overleftarrow{h}_t] + c)\end{aligned}$$

$h = [\vec{h}; \overleftarrow{h}]$ now represents (summarizes) the past and future around a single token.

Deep Bidirectional RNNs



$$\begin{aligned}\overrightarrow{h}_t^{(i)} &= f(\overrightarrow{W}^{(i)} \overrightarrow{h}_{t-1}^{(i-1)} + \overrightarrow{V}^{(i)} \overrightarrow{h}_{t-1}^{(i)} + \overrightarrow{b}^{(i)}) \\ \overleftarrow{h}_t^{(i)} &= f(\overleftarrow{W}^{(i)} \overleftarrow{h}_{t+1}^{(i-1)} + \overleftarrow{V}^{(i)} \overleftarrow{h}_{t+1}^{(i)} + \overleftarrow{b}^{(i)}) \\ y_t &= g(U[\overrightarrow{h}_t^{(L)}; \overleftarrow{h}_t^{(L)}] + c)\end{aligned}$$

Each memory layer passes an intermediate sequential representation to the next.

Data

- MPQA 1.2 corpus (Wiebe et al., 2005)
- consists of 535 news articles (11,111 sentences)
- manually labeled at the phrase level
- Evaluation: F1

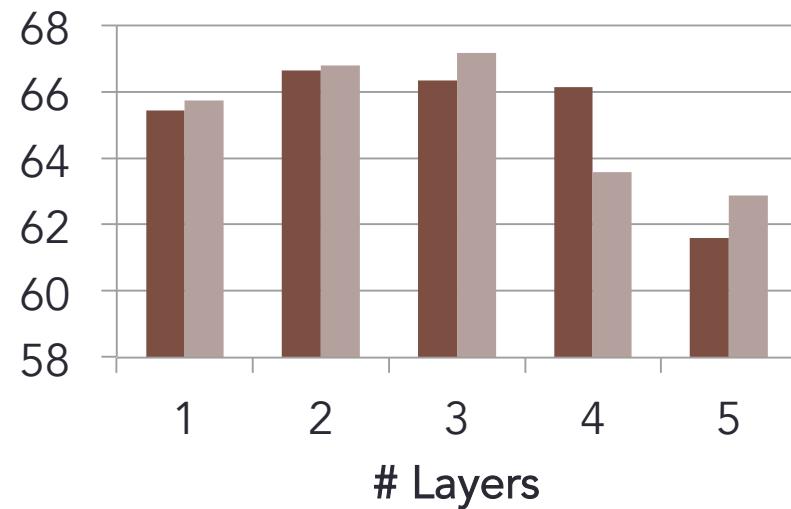
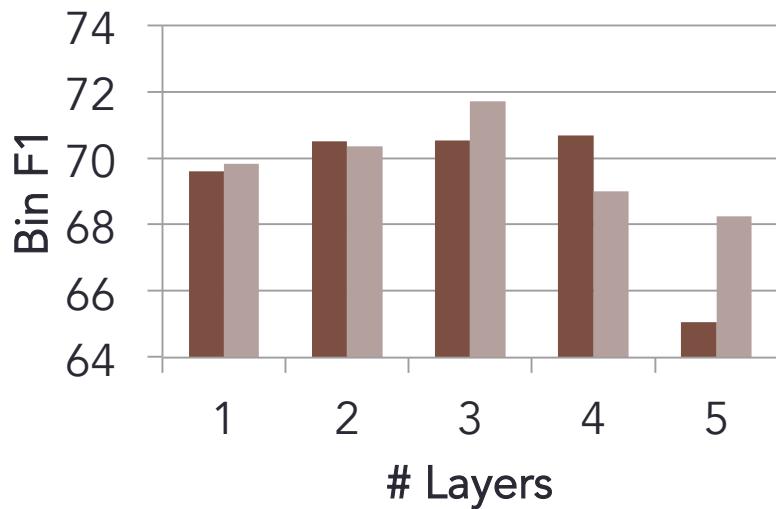
Total population	Condition positive	Condition negative
Test outcome positive	True positive	False positive (Type I error)
Test outcome negative	False negative (Type II error)	True negative

$$\text{precision} = \frac{tp}{tp + fp}$$

$$\text{recall} = \frac{tp}{tp + fn}$$

$$F1 = 2 \cdot \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}}$$

Evaluation



■ 24k
■ 200k

Recap

- Recurrent Neural Network is one of the best deepNLP model families
- Training them is hard because of vanishing and exploding gradient problems
- They can be extended in many ways and their training improved with many tricks (more to come)
- Next: Most important and powerful RNN extensions with LSTMs and GRUs

Overview

- Recap of most important concepts & equations
- Wrap up: Deep RNNs and F1 evaluation
- Machine translation
- Fancy RNN Models tackling MT:
 - Gated Recurrent Units by Cho et al. (2014)
 - Long-Short-Term-Memories
by Hochreiter and Schmidhuber (1997)

Advanced, cutting edge, blast from the past

Recap of most important concepts

Word2Vec
$$J_t(\theta) = \log \sigma(u_o^T v_c) + \sum_{j \sim P(w)} [\log \sigma(-u_j^T v_c)]$$

Glove
$$J(\theta) = \frac{1}{2} \sum_{i,j=1}^W f(P_{ij})(u_i^T v_j - \log P_{ij})^2$$

Nnet & Max-margin
$$s = U^T f(Wx + b)$$

$$J = \max(0, 1 - s + s_c)$$

Recap of most important concepts

Recurrent Neural Networks

$$\begin{aligned} h_t &= \sigma \left(W^{(hh)} h_{t-1} + W^{(hx)} x_{[t]} \right) \\ \hat{y}_t &= \text{softmax} \left(W^{(S)} h_t \right) \end{aligned}$$

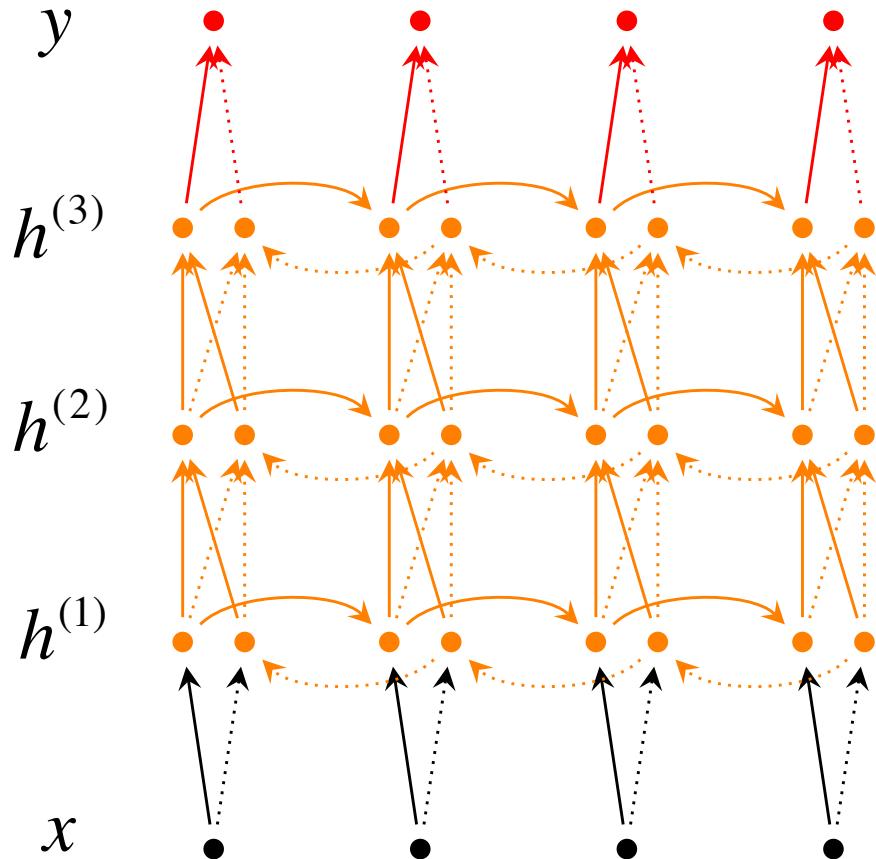
Cross Entropy Error

$$J^{(t)}(\theta) = - \sum_{j=1}^{|V|} y_{t,j} \log \hat{y}_{t,j}$$

Mini-batched SGD

$$\theta^{new} = \theta^{old} - \alpha \nabla_{\theta} J_{t:t+B}(\theta)$$

Deep Bidirectional RNNs by Irsoy and Cardie

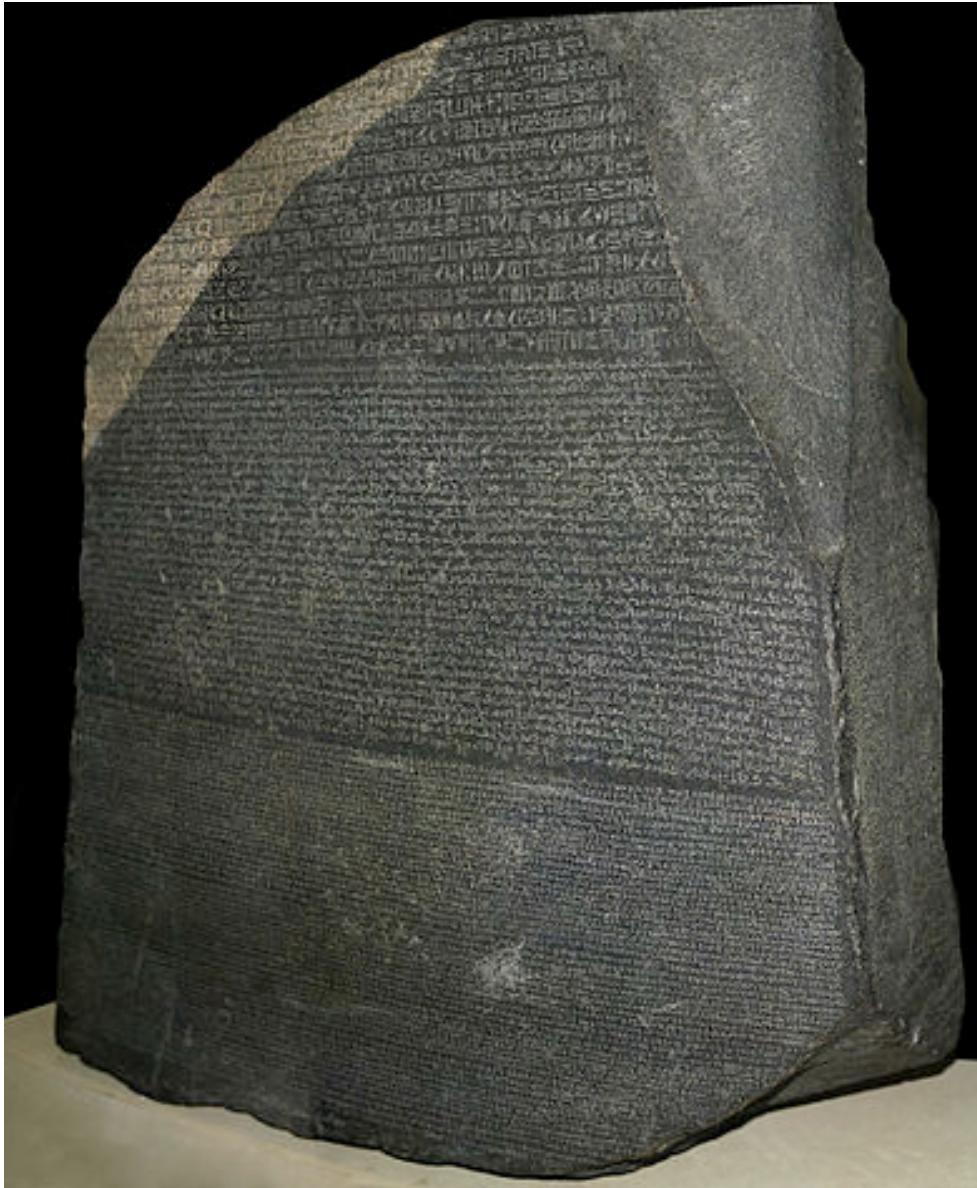


$$\begin{aligned}\overrightarrow{h}_t^{(i)} &= f(\overrightarrow{W}^{(i)} \overrightarrow{h}_{t-1}^{(i-1)} + \overrightarrow{V}^{(i)} \overrightarrow{h}_{t-1}^{(i)} + \overrightarrow{b}^{(i)}) \\ \overleftarrow{h}_t^{(i)} &= f(\overleftarrow{W}^{(i)} \overleftarrow{h}_t^{(i-1)} + \overleftarrow{V}^{(i)} \overleftarrow{h}_{t+1}^{(i)} + \overleftarrow{b}^{(i)}) \\ y_t &= g(U[\overrightarrow{h}_t^{(L)}; \overleftarrow{h}_t^{(L)}] + c)\end{aligned}$$

Each memory layer passes an intermediate sequential representation to the next.

Machine Translation

- Methods are statistical
- Use parallel corpora
 - European Parliament
- First parallel corpus:
 - Rosetta Stone →
- Traditional systems are very complex

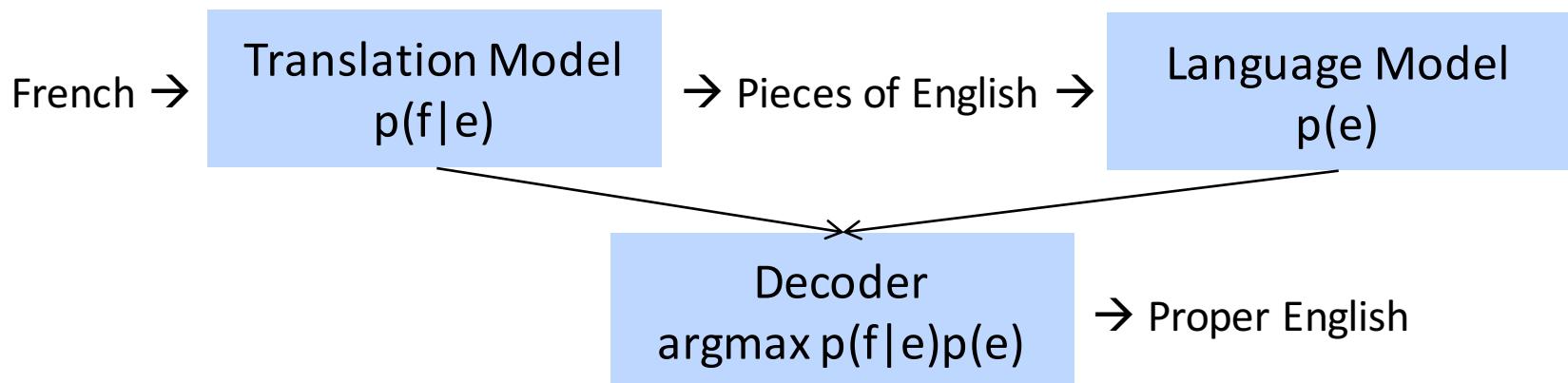


Current statistical machine translation systems

- Source language f , e.g. French
- Target language e , e.g. English
- Probabilistic formulation (using Bayes rule)

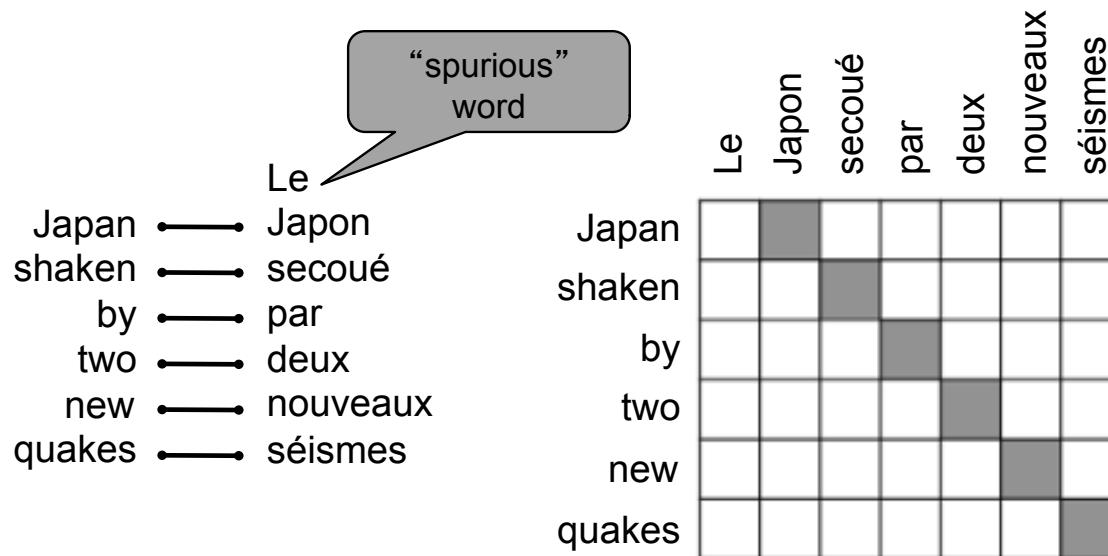
$$\hat{e} = \operatorname{argmax}_e p(e|f) = \operatorname{argmax}_e p(f|e)p(e)$$

- Translation model $p(f|e)$ trained on parallel corpus
- Language model $p(e)$ trained on English only corpus (lots, free!)



Step 1: Alignment

Goal: know which word or phrases in source language would translate to what words or phrases in target language? → Hard already!



Alignment examples from Chris Manning/CS224n

Step 1: Alignment

“zero fertility” word
not translated

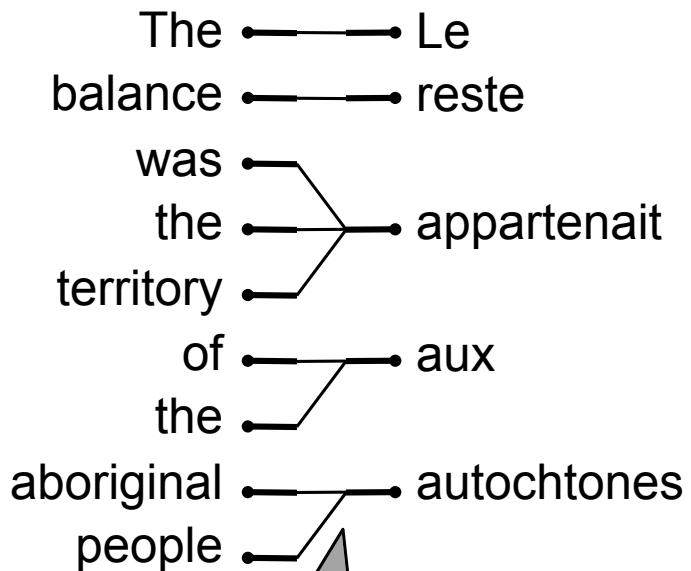
And Le
the programme
program a
has été
been mis
implemented en
 application

one-to-many
alignment

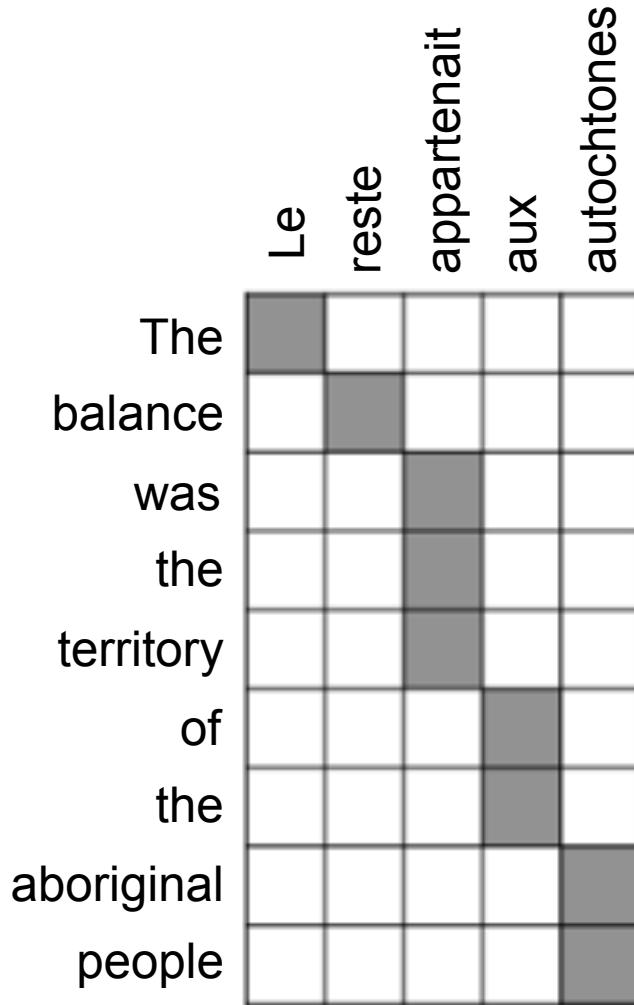
Le	programme	a	été	mis	en	application
And						
the						
program						
has						
been						
implemented						

Step 1: Alignment

Really hard :/



many-to-one
alignments



Step 1: Alignment

The Les
poor pauvres
don't sont
have démunis
any
money

many-to-many
alignment

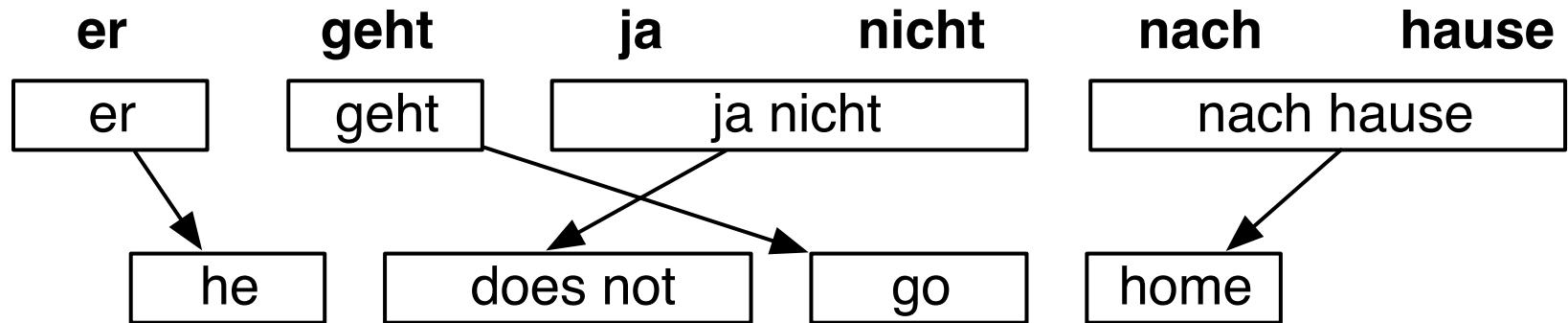
Les pauvres
sont démunis

The pauvres
poor démunis
don' t démunis
have pauvres
any pauvres
money pauvres

phrase
alignment

Step 1: Alignment

- We could spend an entire lecture on alignment models
- Not only single words but could use phrases, syntax
- Then consider reordering of translated phrases

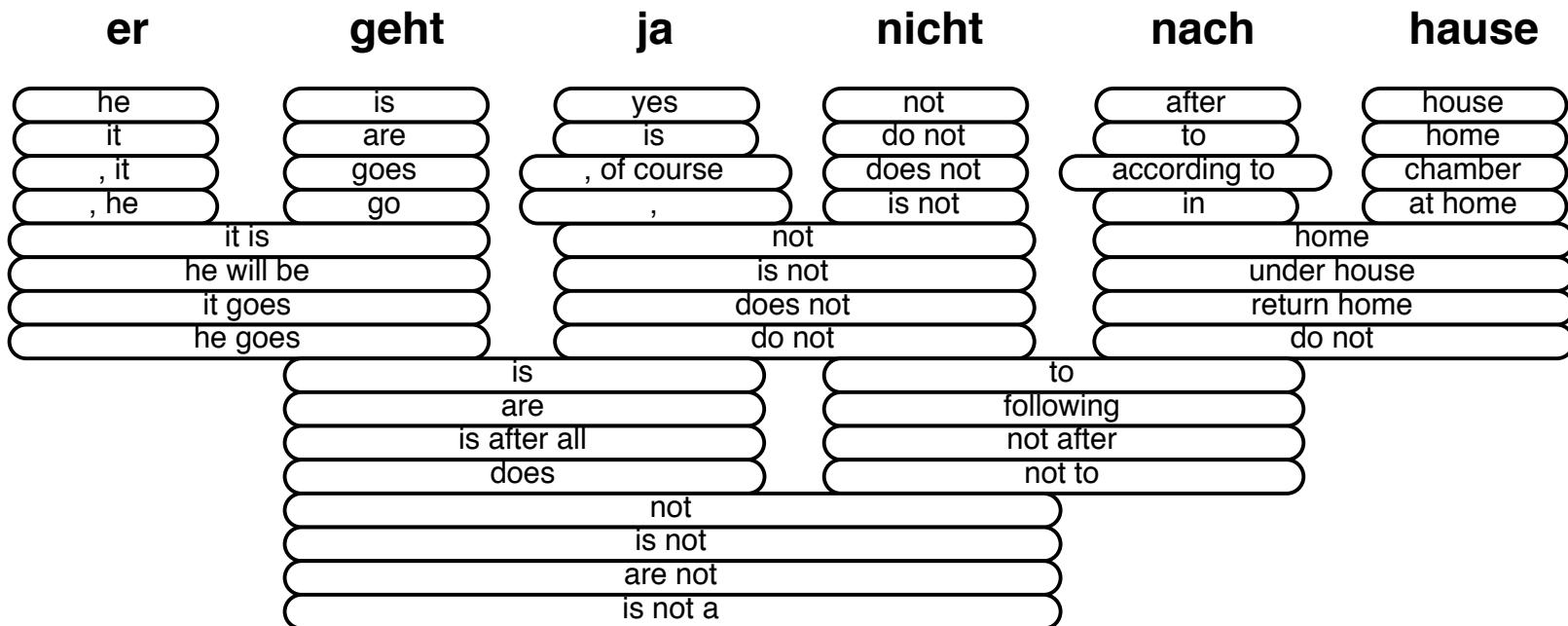


Example from Philipp Koehn

After many steps

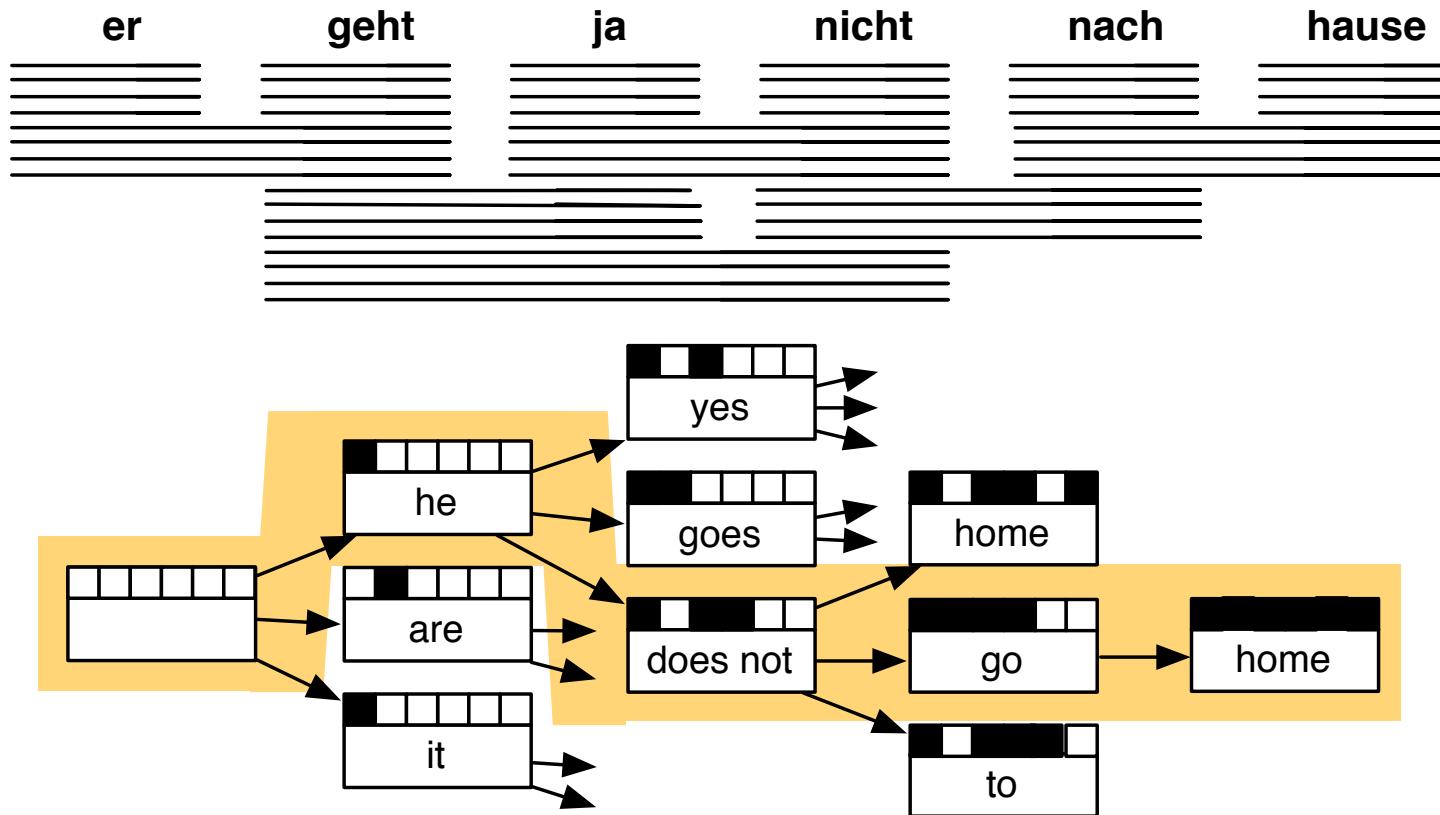
Each phrase in source language has many possible translations resulting in large search space:

Translation Options



Decode: Search for best of many hypotheses

Hard search problem that also includes language model

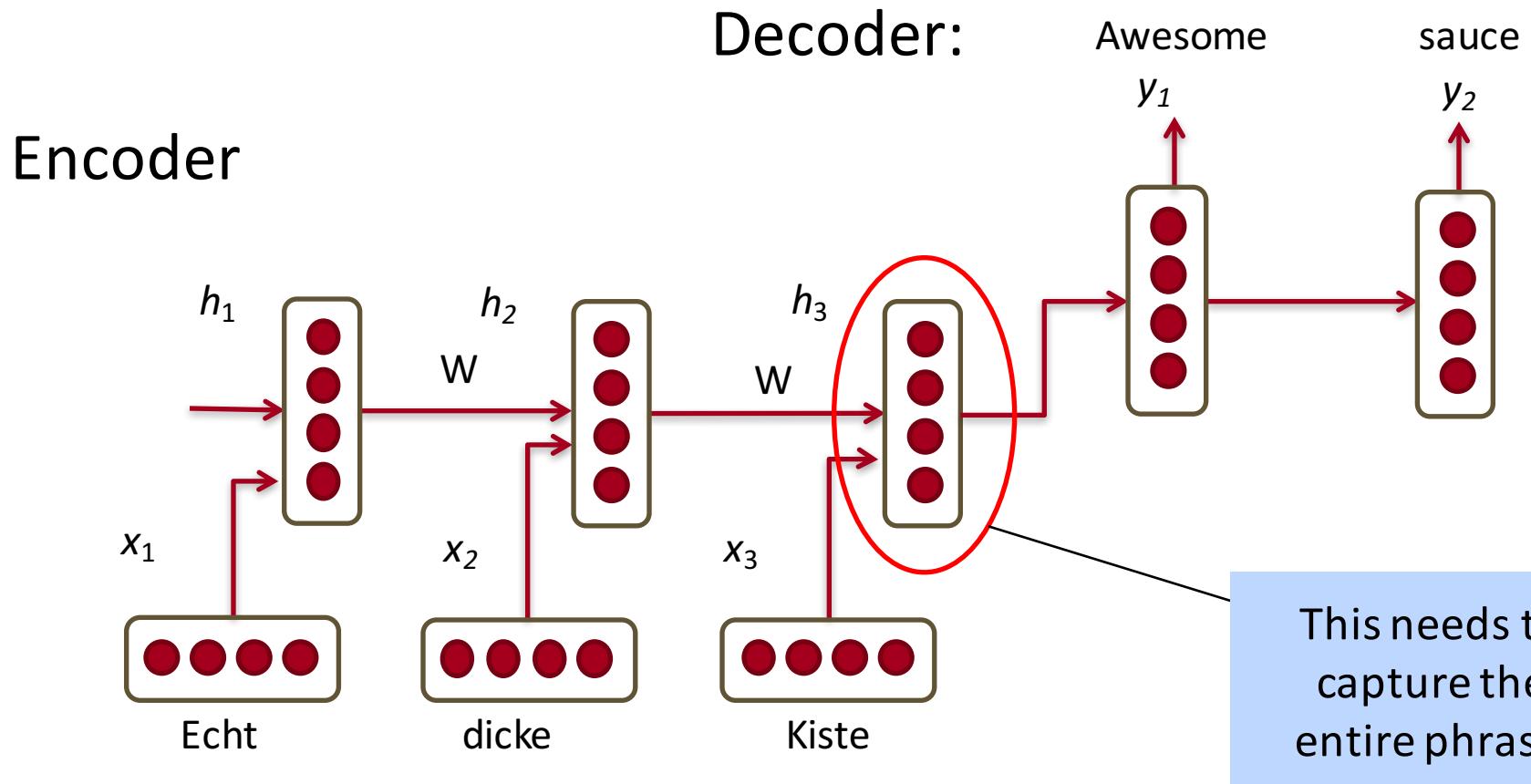


Traditional MT

- Skipped hundreds of important details
- A lot of human feature engineering
- Very complex systems
- Many different, independent machine learning problems

Deep learning to the rescue! ... ?

Maybe, we could translate directly with an RNN?



MT with RNNs – Simplest Model

Encoder: $h_t = \phi(h_{t-1}, x_t) = f\left(W^{(hh)}h_{t-1} + W^{(hx)}x_t\right)$

Decoder: $h_t = \phi(h_{t-1}) = f\left(W^{(hh)}h_{t-1}\right)$

$$y_t = \text{softmax}\left(W^{(S)}h_t\right)$$

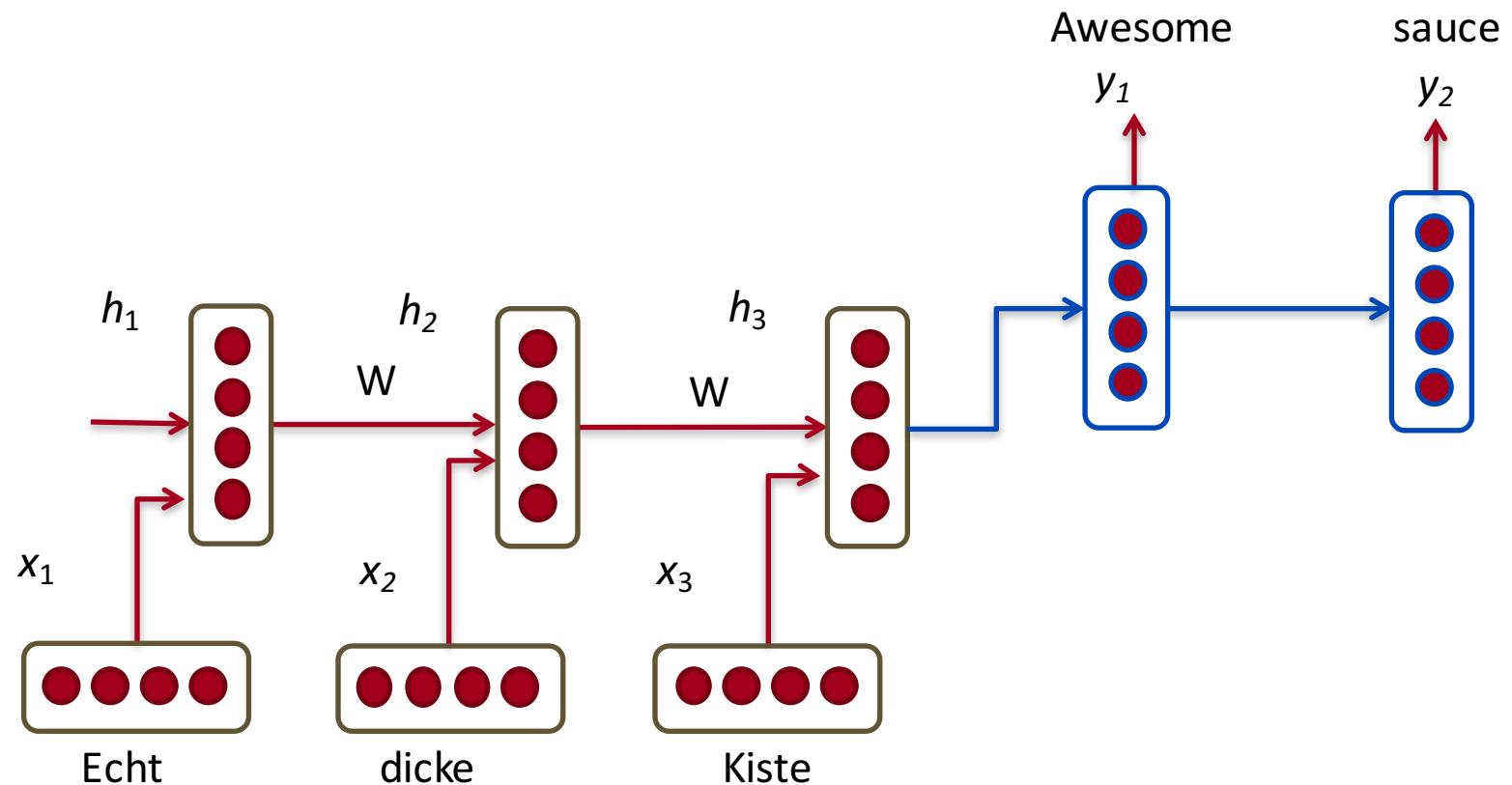
Minimize cross entropy error for all target words
conditioned on source words

$$\max_{\theta} \frac{1}{N} \sum_{n=1}^N \log p_{\theta}(y^{(n)} | x^{(n)})$$

It's not quite that simple

RNN Translation Model Extensions

1. Train different RNN weights for encoding and decoding



RNN Translation Model Extensions

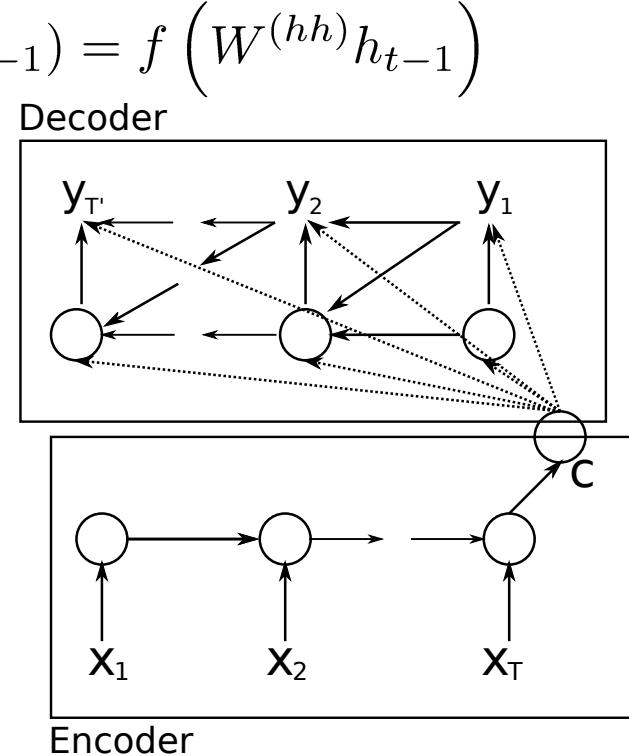
Notation: Each input of \mathcal{A} has its own linear transformation matrix. Simple:

$$h_t = \phi(h_{t-1}) = f\left(W^{(hh)} h_{t-1}\right)$$

2. Compute every hidden state in decoder from

- Previous hidden state (standard)
- Last hidden vector of encoder $c = h_T$
- Previous predicted output word y_{t-1}

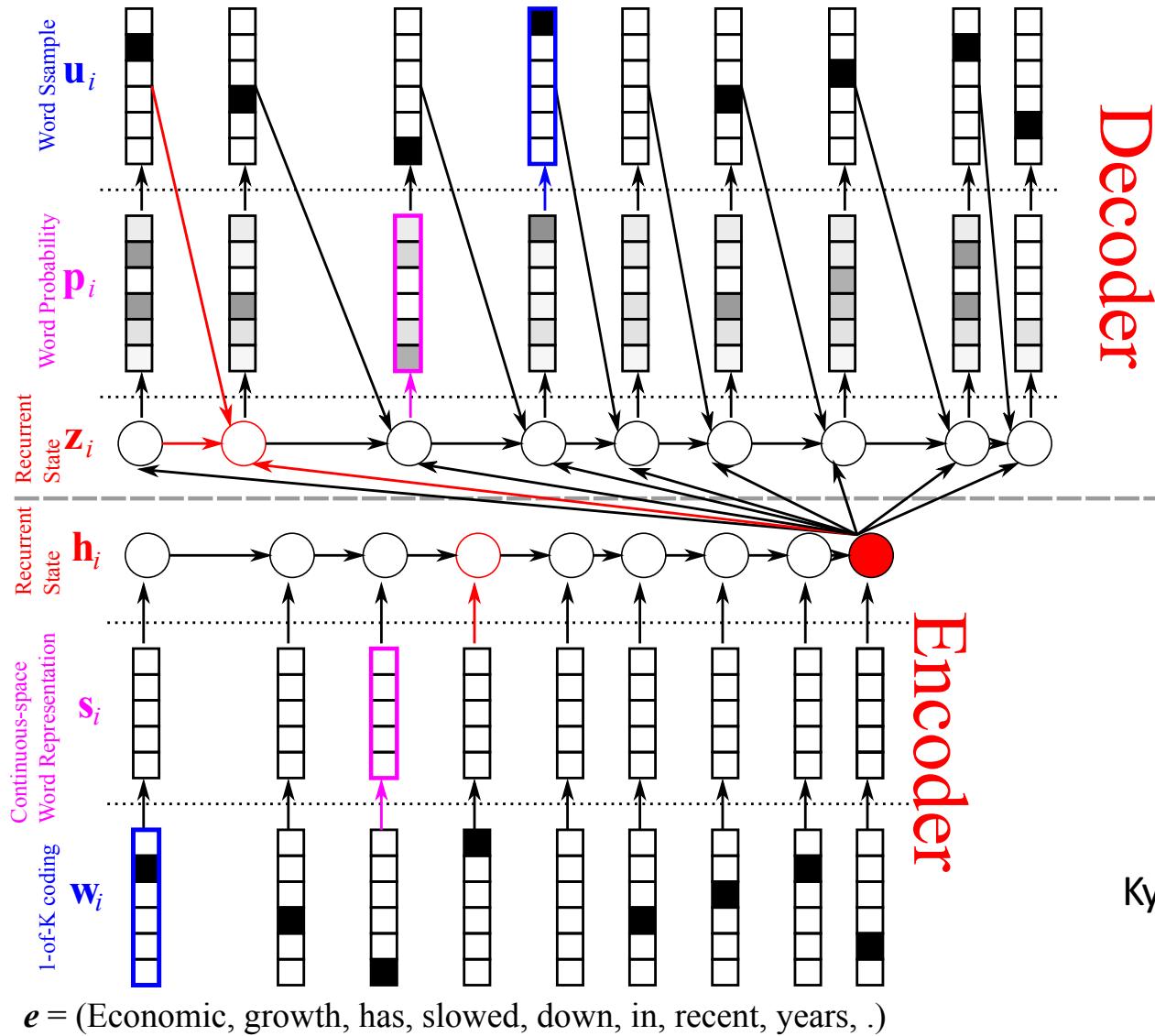
$$h_{D,t} = \phi_D(h_{t-1}, c, y_{t-1})$$



Cho et al. 2014

Different picture, same idea

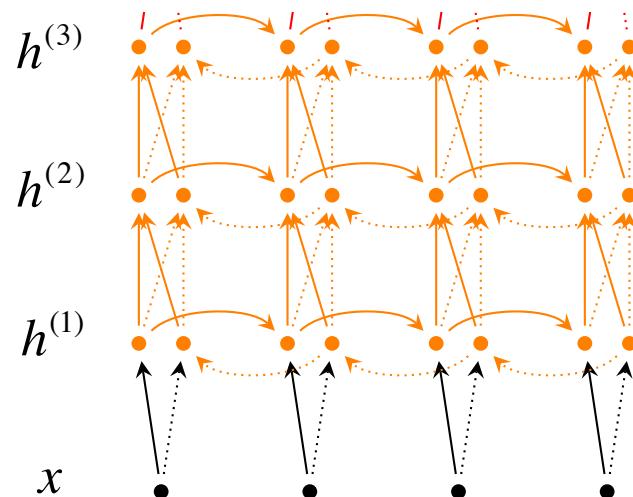
$f = (\text{La}, \text{croissance}, \text{économique}, \text{s'est}, \text{ralentie}, \text{ces}, \text{dernières}, \text{années}, .)$



Kyunghyun Cho et al. 2014

RNN Translation Model Extensions

3. Train stacked/deep RNNs with multiple layers
4. Potentially train bidirectional encoder
5. Train input sequence in reverse order for simple optimization problem: Instead of A B C → X Y, train with C B A → X Y



6. Main Improvement: Better Units

- More complex hidden unit computation in recurrence!
- Gated Recurrent Units (GRU) introduced by Cho et al. 2014
- Main ideas:
 - keep around memories to capture long distance dependencies
 - allow error messages to flow at different strengths depending on the inputs

GRUs

- Standard RNN computes hidden layer at next time step directly:
$$h_t = f \left(W^{(hh)} h_{t-1} + W^{(hx)} x_t \right)$$
- GRU first computes an update **gate** (another layer) based on current input word vector and hidden state

$$z_t = \sigma \left(W^{(z)} x_t + U^{(z)} h_{t-1} \right)$$

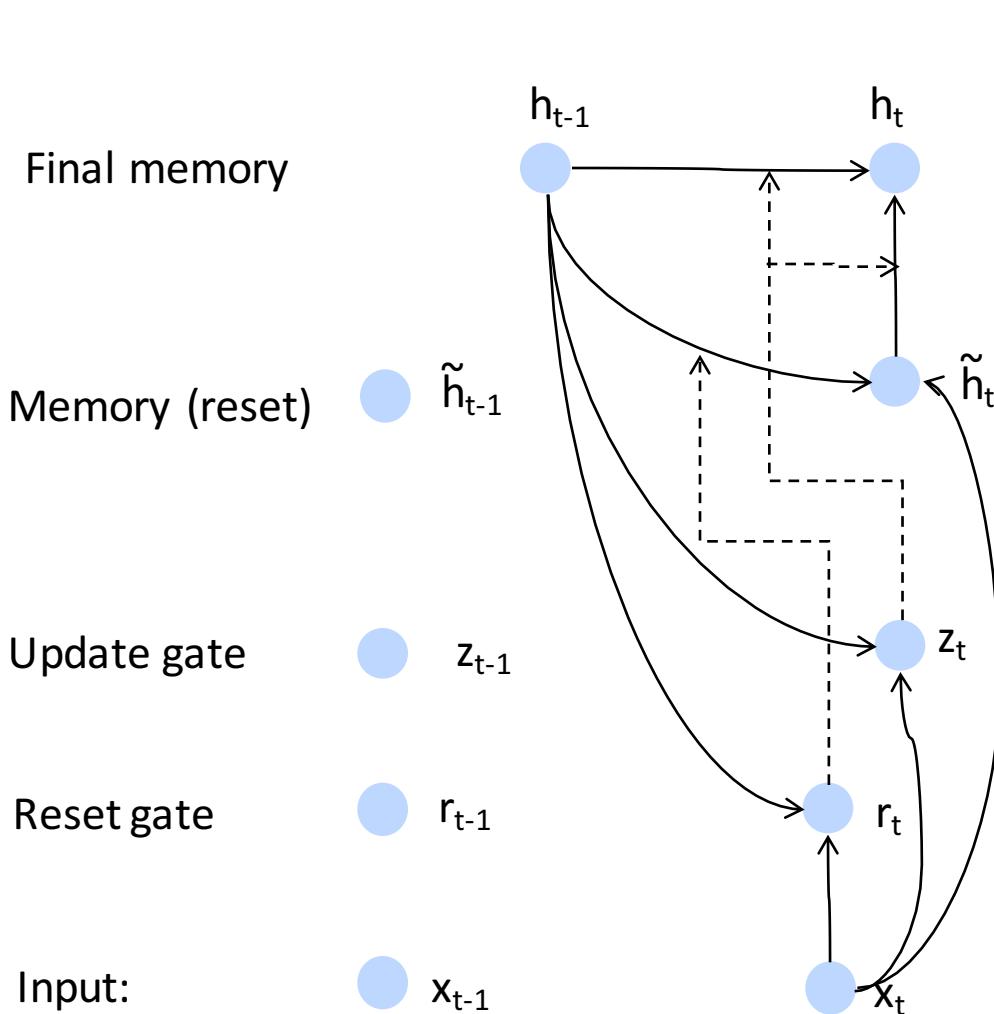
- Compute reset gate similarly but with different weights

$$r_t = \sigma \left(W^{(r)} x_t + U^{(r)} h_{t-1} \right)$$

GRUs

- Update gate
$$z_t = \sigma \left(W^{(z)} x_t + U^{(z)} h_{t-1} \right)$$
- Reset gate
$$r_t = \sigma \left(W^{(r)} x_t + U^{(r)} h_{t-1} \right)$$
- New memory content: $\tilde{h}_t = \tanh (W x_t + r_t \circ U h_{t-1})$
If reset gate unit is ~ 0 , then this ignores previous memory and only stores the new word information
- Final memory at time step combines current and previous time steps:
$$h_t = z_t \circ h_{t-1} + (1 - z_t) \circ \tilde{h}_t$$

Attempt at a clean illustration



$$z_t = \sigma \left(W^{(z)} x_t + U^{(z)} h_{t-1} \right)$$

$$r_t = \sigma \left(W^{(r)} x_t + U^{(r)} h_{t-1} \right)$$

$$\tilde{h}_t = \tanh (W x_t + r_t \circ U h_{t-1})$$

$$h_t = z_t \circ h_{t-1} + (1 - z_t) \circ \tilde{h}_t$$

GRU intuition

- If reset is close to 0, ignore previous hidden state
→ Allows model to drop information that is irrelevant in the future
- Update gate z controls how much of past state should matter now.
 - If z close to 1, then we can copy information in that unit through many time steps! **Less vanishing gradient!**
- Units with short-term dependencies often have reset gates very active

$$z_t = \sigma \left(W^{(z)} x_t + U^{(z)} h_{t-1} \right)$$

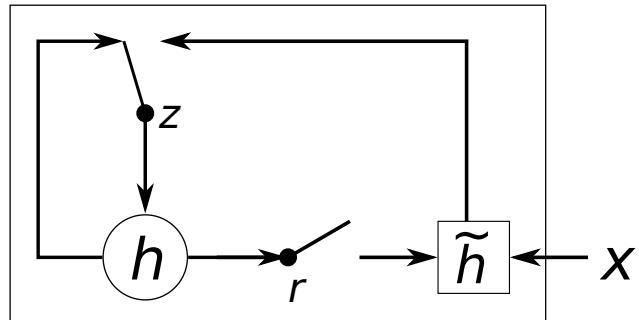
$$r_t = \sigma \left(W^{(r)} x_t + U^{(r)} h_{t-1} \right)$$

$$\tilde{h}_t = \tanh (W x_t + r_t \circ U h_{t-1})$$

$$h_t = z_t \circ h_{t-1} + (1 - z_t) \circ \tilde{h}_t$$

GRU intuition

- Units with long term dependencies have active update gates z
- Illustration:



- Derivative of $\frac{\partial}{\partial x_1} x_1 x_2$? → rest is same chain rule, but implement with **modularization** or automatic differentiation

$$z_t = \sigma(W^{(z)}x_t + U^{(z)}h_{t-1})$$

$$r_t = \sigma(W^{(r)}x_t + U^{(r)}h_{t-1})$$

$$\tilde{h}_t = \tanh(Wx_t + r_t \circ Uh_{t-1})$$

$$h_t = z_t \circ h_{t-1} + (1 - z_t) \circ \tilde{h}_t$$

Demo session

Long-short-term-memories (LSTMs)

- We can make the units even more complex
- Allow each time step to modify
 - Input gate (current cell matters)
 - Forget (gate 0, forget past)
 - Output (how much cell is exposed)
 - New memory cell
- Final memory cell:
- Final hidden state:

$$i_t = \sigma \left(W^{(i)}x_t + U^{(i)}h_{t-1} \right)$$

$$f_t = \sigma \left(W^{(f)}x_t + U^{(f)}h_{t-1} \right)$$

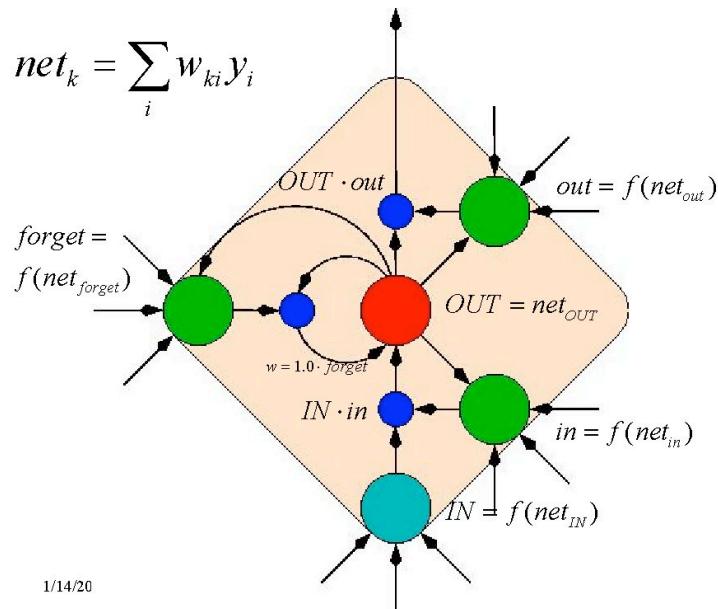
$$o_t = \sigma \left(W^{(o)}x_t + U^{(o)}h_{t-1} \right)$$

$$\tilde{c}_t = \tanh \left(W^{(c)}x_t + U^{(c)}h_{t-1} \right)$$

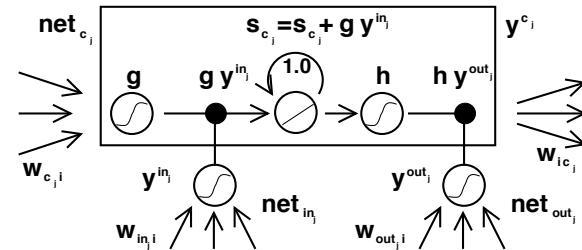
$$c_t = f_t \circ c_{t-1} + i_t \circ \tilde{c}_t$$

$$h_t = o_t \circ \tanh(c_t)$$

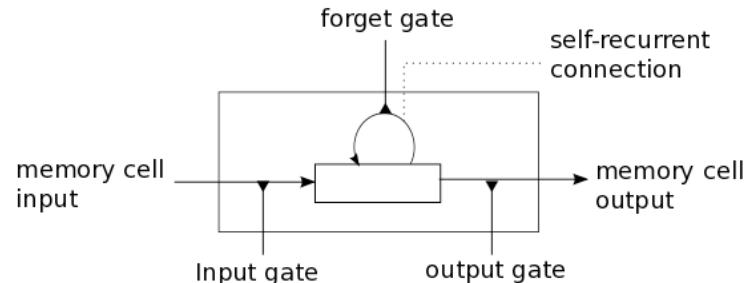
Illustrations all a bit overwhelming ;)



<http://people.idsia.ch/~juergen/lstm/sld017.htm>



Long Short-Term Memory by Hochreiter and Schmidhuber (1997)



<http://deeplearning.net/tutorial/lstm.html>

Intuition: memory cells can keep information intact, unless inputs makes them forget it or overwrite it with new input.
 Cell can decide to output this information or just store it

LSTMs are currently very hip!

- En vogue default model for most sequence labeling tasks
- Very powerful, especially when stacked and made even deeper (each hidden layer is already computed by a deep internal network)
- Most useful if you have lots and lots of data

Deep LSTMs compared to traditional systems 2015

Method	test BLEU score (ntst14)
Bahdanau et al. [2]	28.45
Baseline System [29]	33.30
Single forward LSTM, beam size 12	26.17
Single reversed LSTM, beam size 12	30.59
Ensemble of 5 reversed LSTMs, beam size 1	33.00
Ensemble of 2 reversed LSTMs, beam size 12	33.27
Ensemble of 5 reversed LSTMs, beam size 2	34.50
Ensemble of 5 reversed LSTMs, beam size 12	34.81

Table 1: The performance of the LSTM on WMT'14 English to French test set (ntst14). Note that an ensemble of 5 LSTMs with a beam of size 2 is cheaper than of a single LSTM with a beam of size 12.

Method	test BLEU score (ntst14)
Baseline System [29]	33.30
Cho et al. [5]	34.54
Best WMT'14 result [9]	37.0
Rescoring the baseline 1000-best with a single forward LSTM	35.61
Rescoring the baseline 1000-best with a single reversed LSTM	35.85
Rescoring the baseline 1000-best with an ensemble of 5 reversed LSTMs	36.5
Oracle Rescoring of the Baseline 1000-best lists	~45

Deep LSTMs (with a lot more tweaks) today

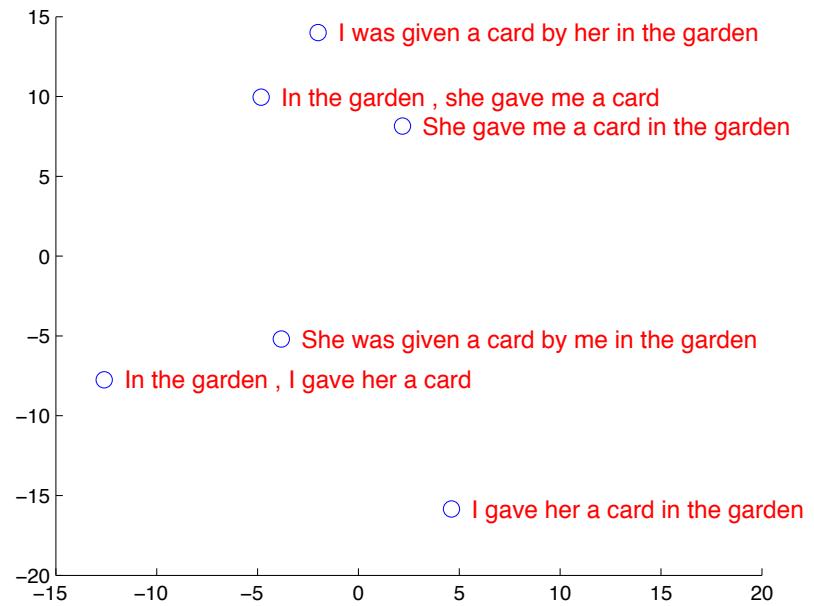
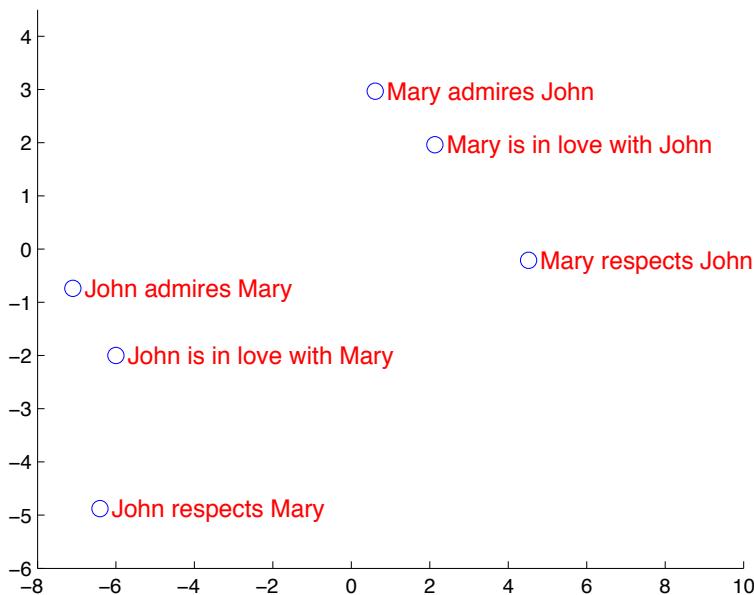
WMT 2016 competition results from yesterday

Scored Systems					
System	Submitter	System Notes	Constraint	Run Notes	BLEU
uedin-nmt-ensemble (Details)	rsennrich <i>University of Edinburgh</i>	BPE neural MT system with monolingual training data (back-translated), ensemble of 4, reranked with right-to-left model.	yes		34.8
metamind-ensemble (Details)	jekbradbury <i>Salesforce MetaMind</i>	Neural MT system based on Luong 2015 and Sennrich 2015, using Morfessor for subword splitting, with back-translated monolingual augmentation. Ensemble of 3 checkpoints from one run plus 1 Y-LSTM (see entry).	yes		32.8
uedin-nmt-single (Details)	rsennrich <i>University of Edinburgh</i>	BPE neural MT system with monolingual training data (back-translated), single model. (contrastive)	yes		32.2

KIT (Details)	niehues <i>KIT</i>	Phrase-based MT with NMT in rescoring	yes		29.7
uedin-pbt-wmt16-en-de (Details)	Matthias Huck <i>University of Edinburgh</i>	Phrase-based Moses	yes		29.1
Moses Phrase-Based (Details)	jhu-smt <i>Johns Hopkins University</i>	Phrase-based model, word clusters for all model components (LM, OSM, LR, sparse features), neural network joint model, large cc LM	yes	[26-7]	29.0
uedin-pbt-wmt16-en-de-contrastive (Details)	Matthias Huck <i>University of Edinburgh</i>	Phrase-based Moses (contrastive, 2015 system)	yes		29.0

Deep LSTM for Machine Translation

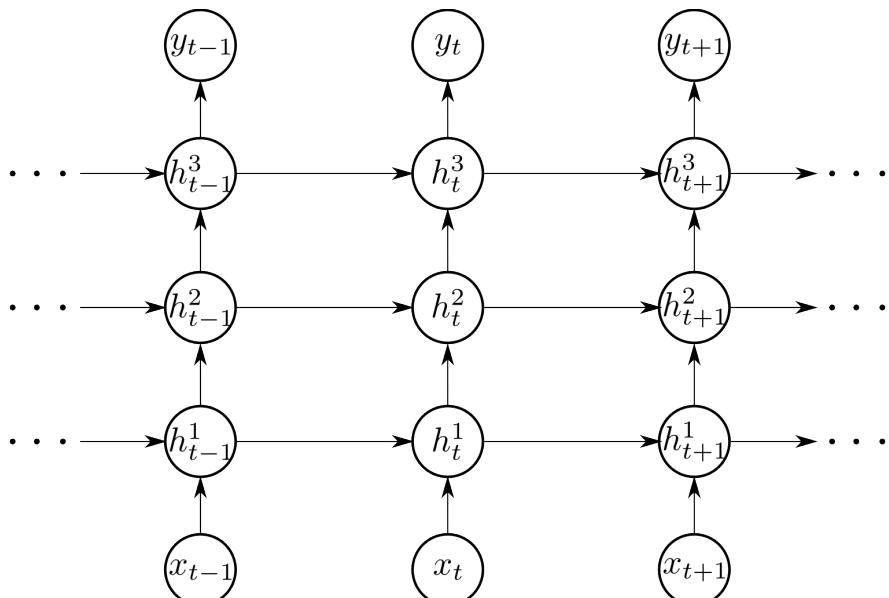
PCA of vectors from last time step hidden layer



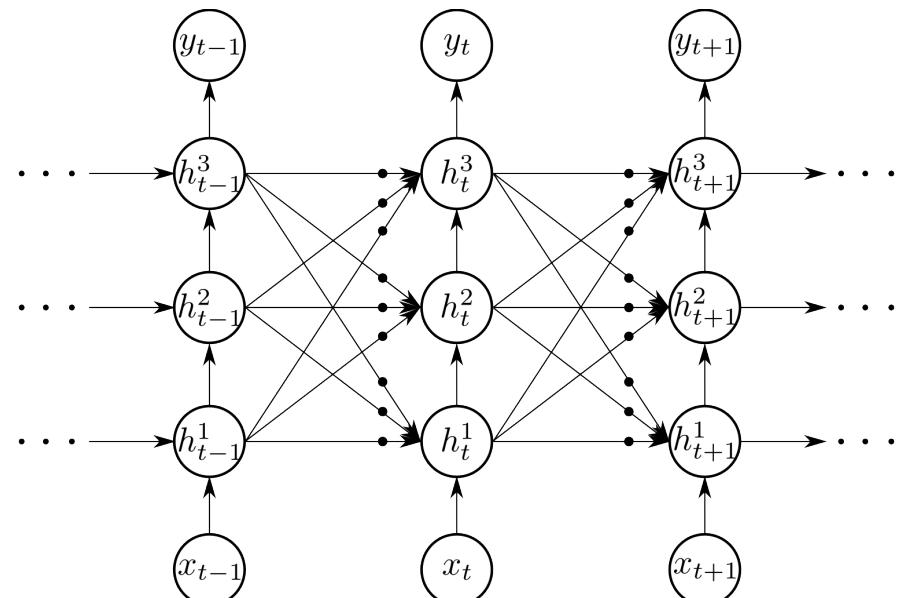
Sequence to Sequence Learning by Sutskever et al. 2014

Further Improvements: More Gates!

Gated Feedback Recurrent Neural Networks, Chung et al. 2015



(a) Conventional stacked RNN



(b) Gated Feedback RNN

Summary

- Recurrent Neural Networks are powerful
- A lot of ongoing work right now
- Gated Recurrent Units even better
- LSTMs maybe even better (jury still out)
- This was an advanced lecture → gain intuition, encourage exploration