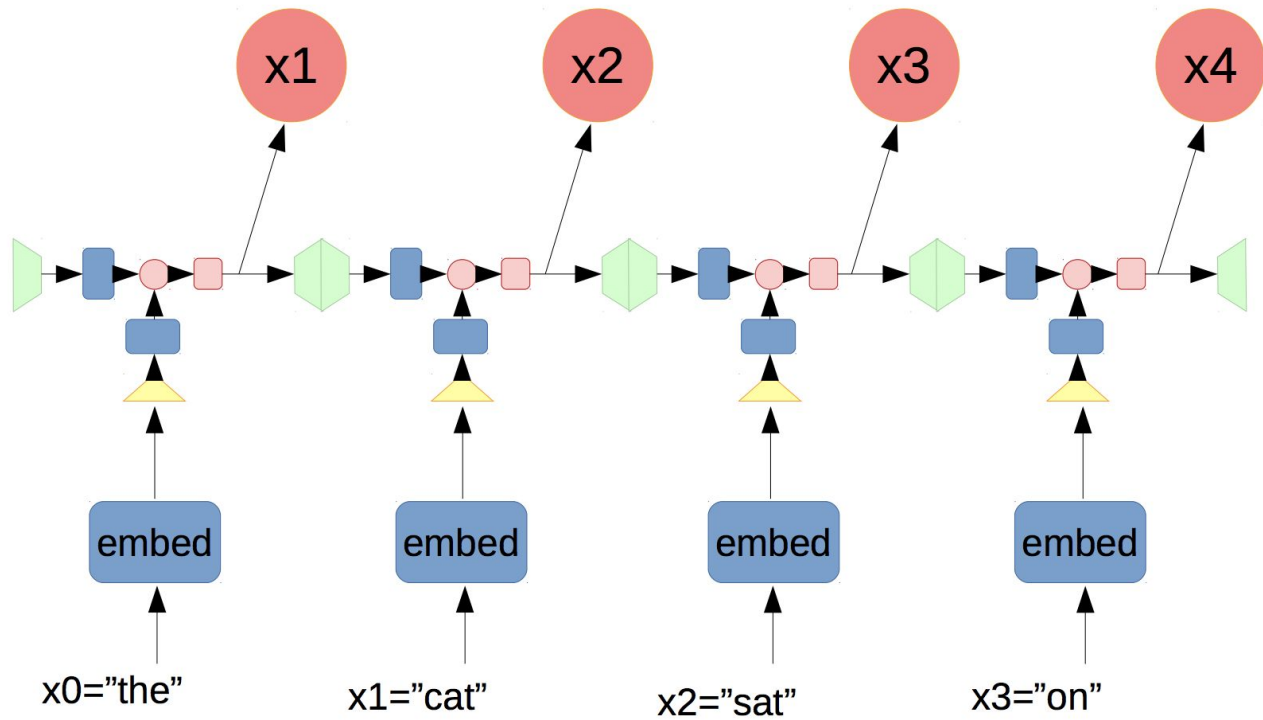# Lecture 2: CNN and vanishing gradient

**Radoslav Neychev**
**Ivan Provilkov**
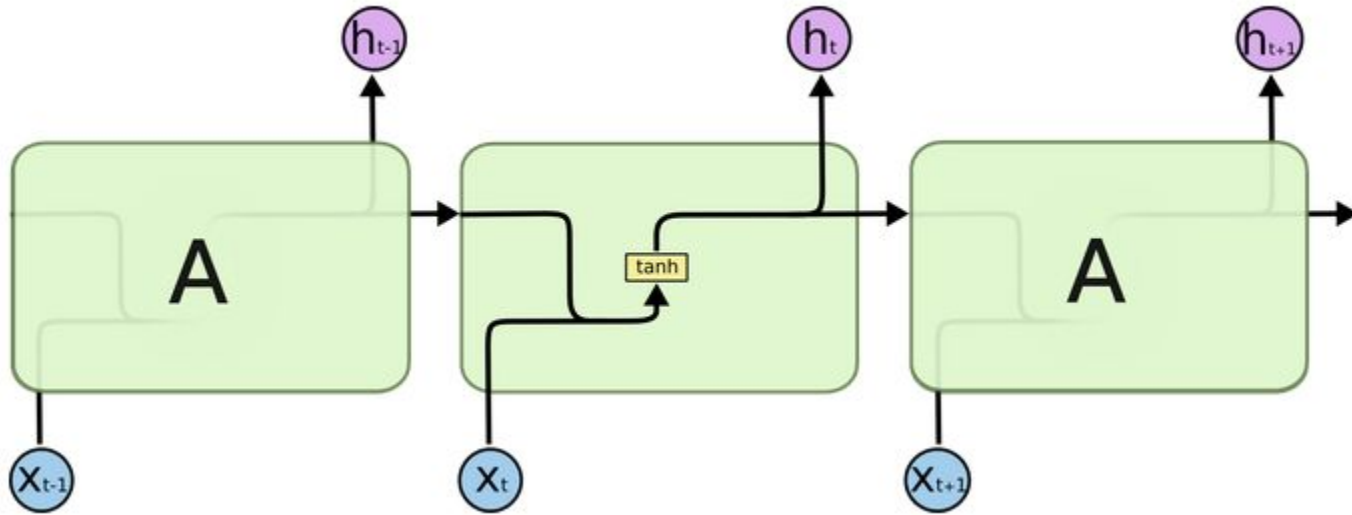
# Outline

- Simple RNN recap
- Complex RNN:
  - Vanishing gradient
  - Exploding gradient
  - LSTM/GRU
  - Gradient clipping
  - Skip connections
  - Residual networks as ensembles
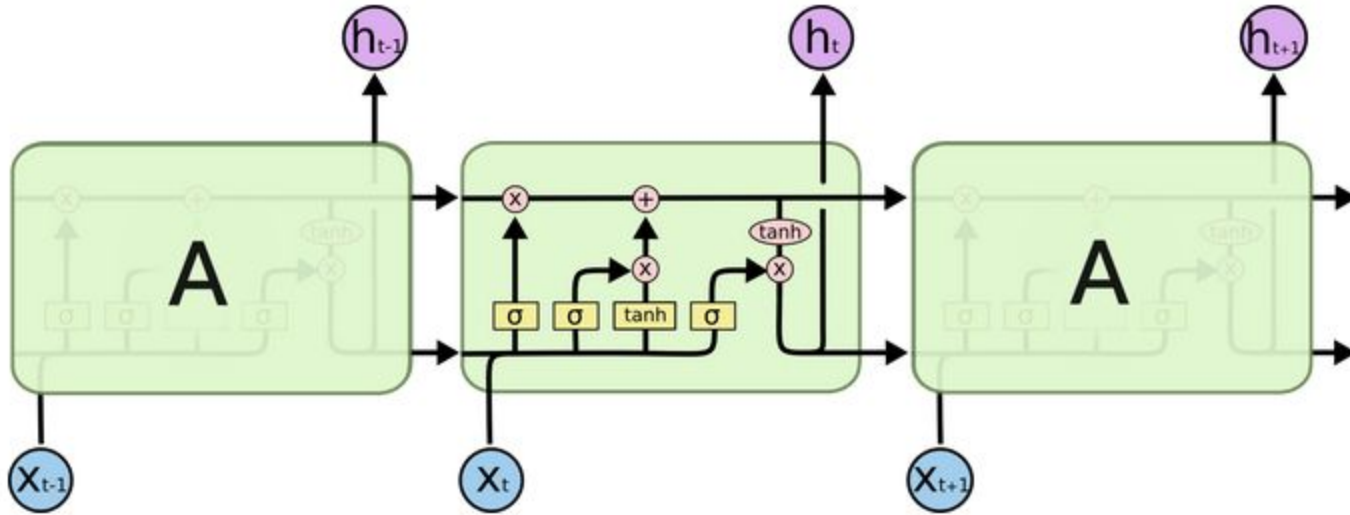- CNNs for text
- Text segmentation

x1  x2  x3  x4

embed  embed  embed  embed

x0="the"  x1="cat"  x2="sat"  x3="on"

Source: https://colah.github.io/posts/2015-08-Understanding-LSTMs/

Source: https://colah.github.io/posts/2015-08-Understanding-LSTMs/

# Vanishing gradient problem

$$J^{(4)}(\theta)$$

$h^{(1)}$ $\quad W \quad$ $h^{(2)}$ $\quad W \quad$ $h^{(3)}$ $\quad W \quad$ $h^{(4)}$

Based on: Lecture by Abigail See, CS224n Lecture 7

# Vanishing gradient problem



$$\frac{\partial J^{(4)}}{\partial \boldsymbol{h}^{(1)}} = \;\; ?$$

Based on: Lecture by Abigail See, CS224n Lecture 7

# Vanishing gradient problem



$$\frac{\partial J^{(4)}}{\partial \boldsymbol{h}^{(1)}} = \frac{\partial \boldsymbol{h}^{(2)}}{\partial \boldsymbol{h}^{(1)}} \times \frac{\partial J^{(4)}}{\partial \boldsymbol{h}^{(2)}}$$

chain rule!

Based on: Lecture by Abigail See, CS224n Lecture 7

# Vanishing gradient problem



$$\frac{\partial J^{(4)}}{\partial \boldsymbol{h}^{(1)}} = \frac{\partial \boldsymbol{h}^{(2)}}{\partial \boldsymbol{h}^{(1)}} \times \qquad \frac{\partial \boldsymbol{h}^{(3)}}{\partial \boldsymbol{h}^{(2)}} \times \frac{\partial J^{(4)}}{\partial \boldsymbol{h}^{(3)}}$$

chain rule!

Based on: Lecture by Abigail See, CS224n Lecture 7

# Vanishing gradient problem

$$J^{(4)}(\theta)$$

$$\boldsymbol{h}^{(1)} \quad W \quad \boldsymbol{h}^{(2)} \quad W \quad \boldsymbol{h}^{(3)} \quad W \quad \boldsymbol{h}^{(4)}$$

$$\frac{\partial J^{(4)}}{\partial \boldsymbol{h}^{(1)}} = \frac{\partial \boldsymbol{h}^{(2)}}{\partial \boldsymbol{h}^{(1)}} \times \qquad \frac{\partial \boldsymbol{h}^{(3)}}{\partial \boldsymbol{h}^{(2)}} \times \qquad \frac{\partial \boldsymbol{h}^{(4)}}{\partial \boldsymbol{h}^{(3)}} \times \frac{\partial J^{(4)}}{\partial \boldsymbol{h}^{(4)}}$$

chain rule!

Based on: Lecture by Abigail See, CS224n Lecture 7

# Vanishing gradient problem

$$J^{(4)}(\theta)$$

Vanishing gradient problem:

*When the derivatives are small, the gradient signal gets smaller and smaller as it backpropagates further*



$h^{(1)}$ $\quad$ $W$ $\quad$ $h^{(2)}$ $\quad$ $W$ $\quad$ $h^{(3)}$ $\quad$ $W$ $\quad$ $h^{(4)}$

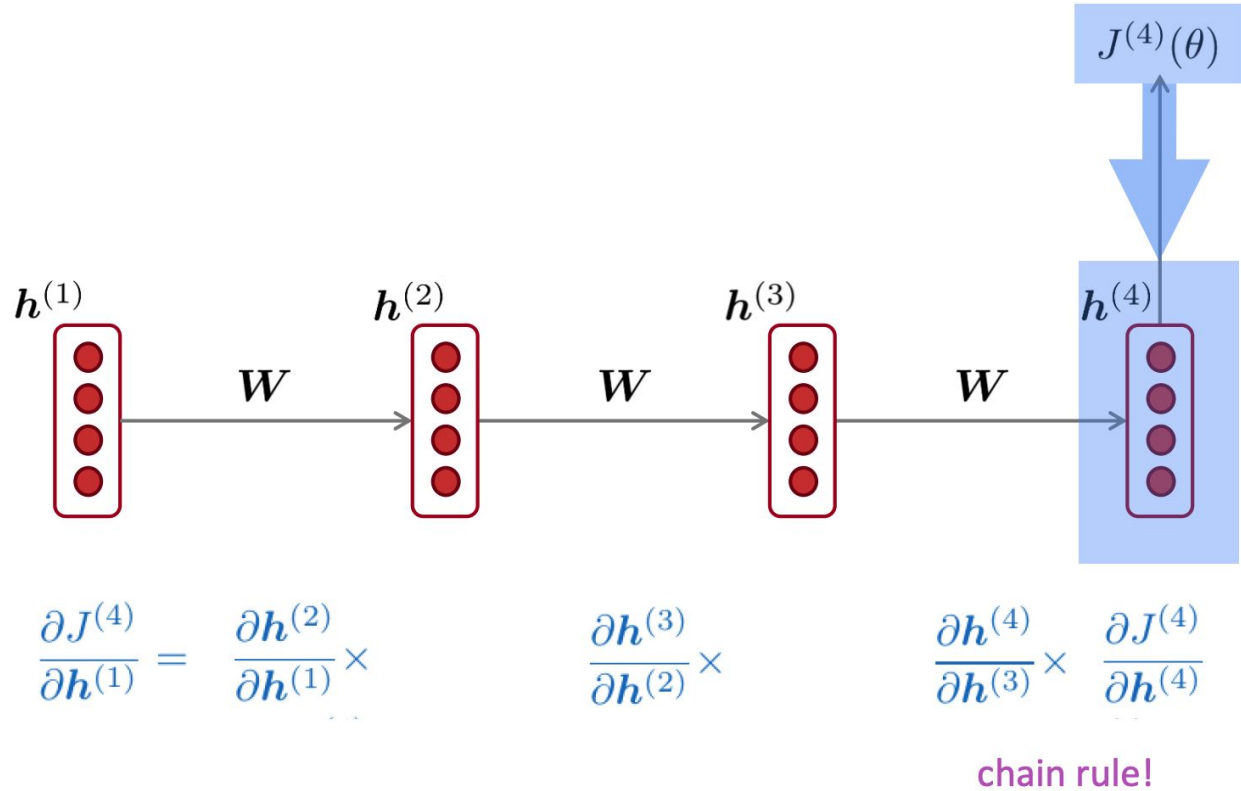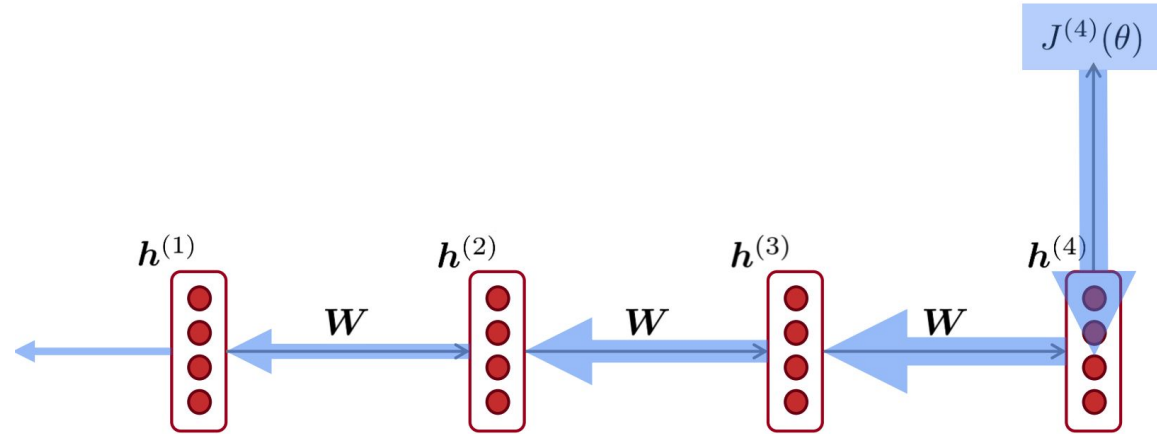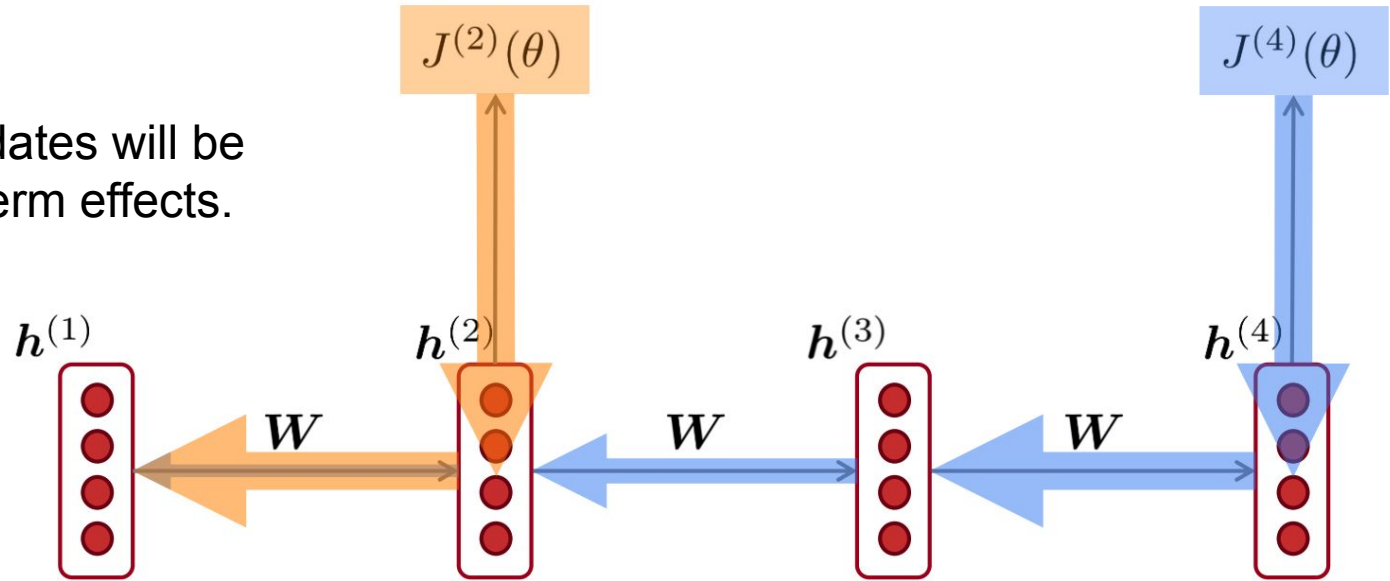$$\frac{\partial J^{(4)}}{\partial h^{(1)}} = \boxed{\frac{\partial h^{(2)}}{\partial h^{(1)}}} \times \boxed{\frac{\partial h^{(3)}}{\partial h^{(2)}}} \times \boxed{\frac{\partial h^{(4)}}{\partial h^{(3)}}} \times \frac{\partial J^{(4)}}{\partial h^{(4)}}$$

What happens if these are small?

More info:  "On the difficulty of training recurrent neural networks", Pascanu et al, 2013
http://proceedings.mlr.press/v28/pascanu13.pdf

11

Gradient signal from far away is lost because it's much smaller than from close-by.

So model weights updates will be based only on short-term effects.

# Vanishing gradient problem

Based on: Lecture by Abigail See, CS224n Lecture 7

# Exploding gradient problem

- If the gradient becomes too big, then the SGD update step becomes too big:

$$\theta^{new} = \theta^{old} - \overbrace{\alpha}^{\text{learning rate}} \underbrace{\nabla_\theta J(\theta)}_{\text{gradient}}$$

- This can cause bad updates: we take too large a step and reach a bad parameter configuration (with large loss)

- In the worst case, this will result in Inf or NaN in your network (then you have to restart training from an earlier checkpoint)

Based on: Lecture by Abigail See, CS224n Lecture 7

# Exploding gradient solution

- Gradient clipping: if the norm of the gradient is greater than some threshold, scale it down before applying SGD update
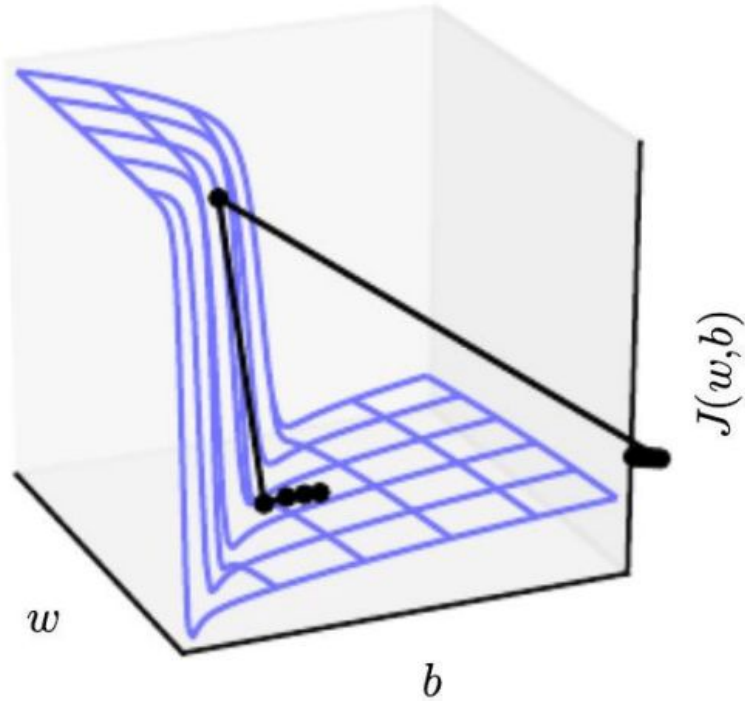
---

**Algorithm 1** Pseudo-code for norm clipping

$\hat{\mathbf{g}} \leftarrow \frac{\partial \mathcal{E}}{\partial \theta}$

**if** $\|\hat{\mathbf{g}}\| \geq threshold$ **then**

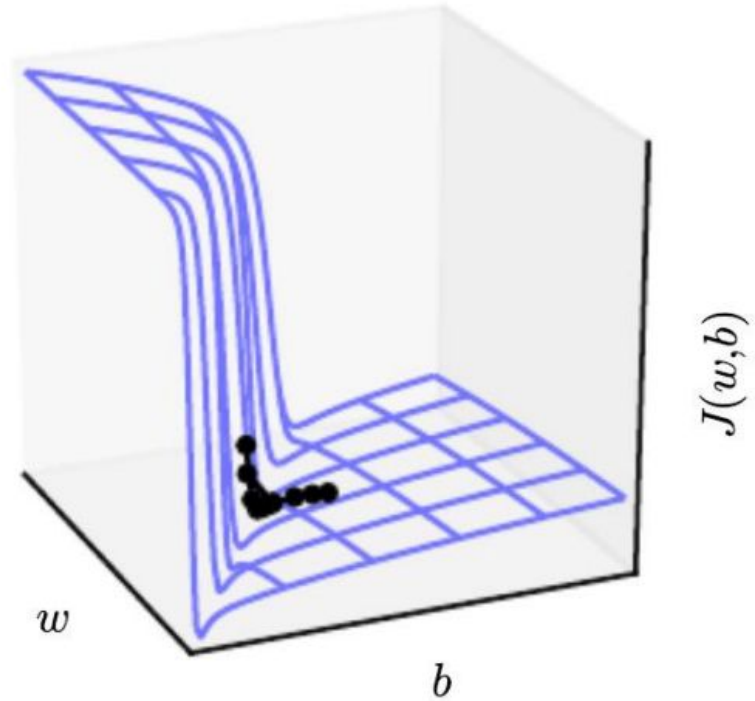$\quad \hat{\mathbf{g}} \leftarrow \frac{threshold}{\|\hat{\mathbf{g}}\|} \hat{\mathbf{g}}$

**end if**

---

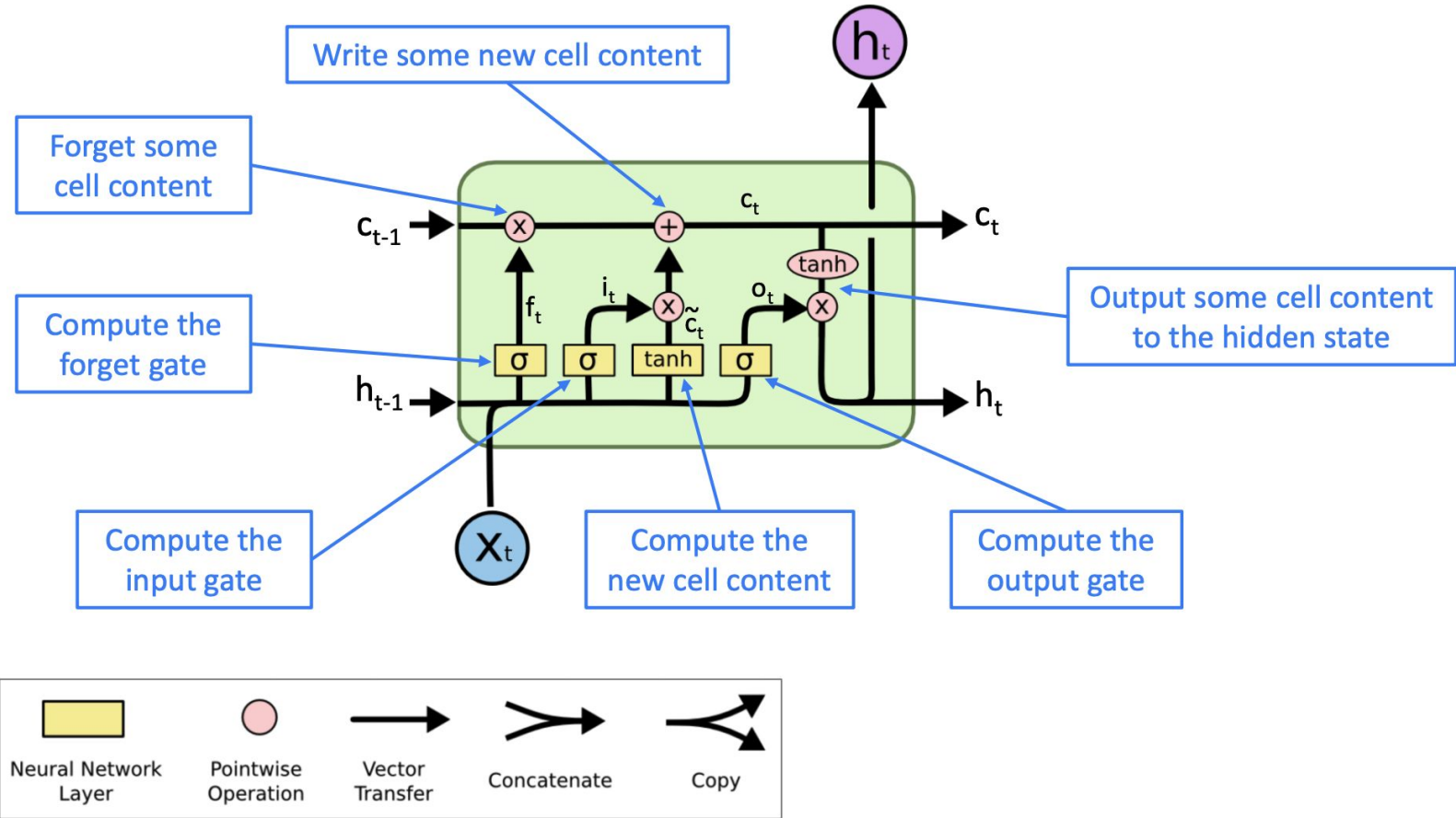- Intuition: take a step in the same direction, but a smaller step

Based on: Lecture by Abigail See, CS224n Lecture 7

# Exploding gradient solution

## Without clipping



## With clipping

Based on: Lecture by Abigail See, CS224n Lecture 7

# Vanishing gradient: LSTM

Write some new cell content

Forget some cell content

Compute the forget gate

Compute the input gate

Compute the new cell content

Compute the output gate

Output some cell content to the hidden state

$h_t$

$c_{t-1}$ $\times$ $+$ $c_t$ $c_t$

$f_t$ $i_t$ $\times$ $\tilde{c}_t$ $o_t$ tanh $\times$

$\sigma$ $\sigma$ tanh $\sigma$

$h_{t-1}$ $h_t$

$x_t$

| | | | | |
|---|---|---|---|---|
| Neural Network Layer | Pointwise Operation | Vector Transfer | Concatenate | Copy |

Based on: Lecture by Abigail See, CS224n Lecture 7

**Forget gate:** controls what is kept vs forgotten, from previous cell state

**Input gate:** controls what parts of the new cell content are written to cell

**Output gate:** controls what parts of cell are output to hidden state

**New cell content:** this is the new content to be written to the cell

**Cell state:** erase ("forget") some content from last cell state, and write ("input") some new cell content

**Hidden state:** read ("output") some content from the cell

**Sigmoid function:** all gate values are between 0 and 1

$$f^{(t)} = \sigma\left(W_f h^{(t-1)} + U_f x^{(t)} + b_f\right)$$

$$i^{(t)} = \sigma\left(W_i h^{(t-1)} + U_i x^{(t)} + b_i\right)$$

$$o^{(t)} = \sigma\left(W_o h^{(t-1)} + U_o x^{(t)} + b_o\right)$$

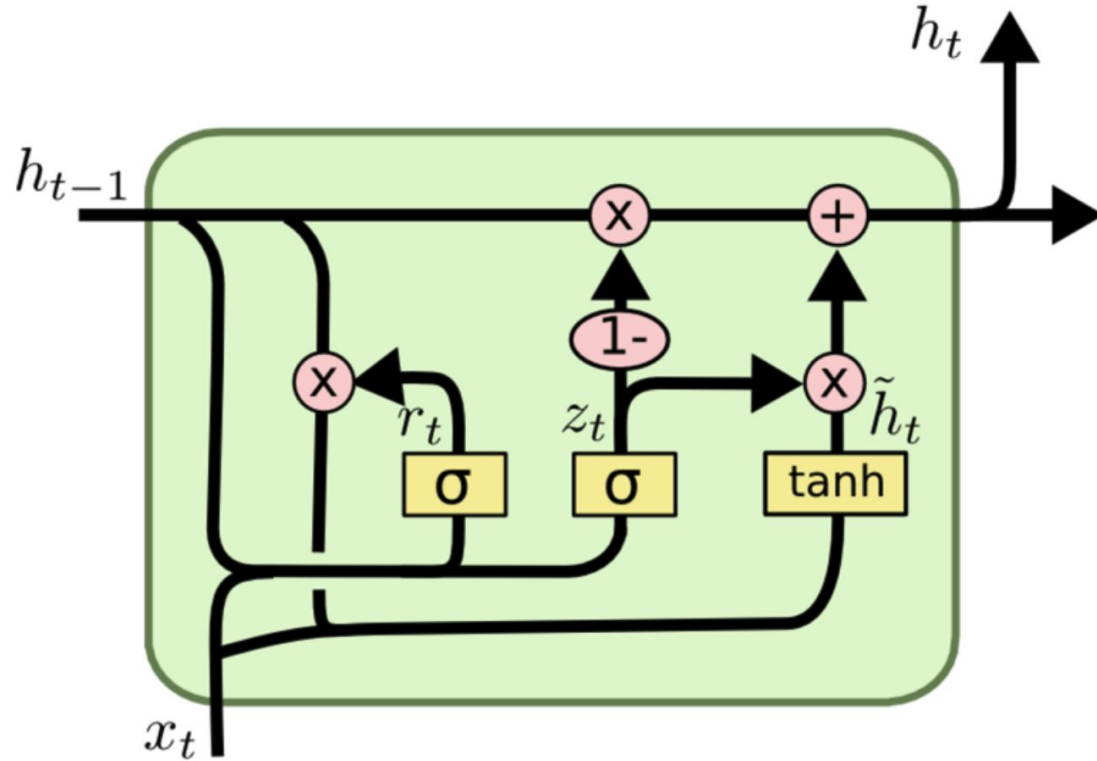$$\tilde{c}^{(t)} = \tanh\left(W_c h^{(t-1)} + U_c x^{(t)} + b_c\right)$$

$$c^{(t)} = f^{(t)} \circ c^{(t-1)} + i^{(t)} \circ \tilde{c}^{(t)}$$

$$h^{(t)} = o^{(t)} \circ \tanh c^{(t)}$$

All these are vectors of same length $n$

Gates are applied using element-wise product

17

**Update gate:** controls what parts of hidden state are updated vs preserved

$$u^{(t)} = \sigma\left(W_u h^{(t-1)} + U_u x^{(t)} + b_u\right)$$

**Reset gate:** controls what parts of previous hidden state are used to compute new content

$$r^{(t)} = \sigma\left(W_r h^{(t-1)} + U_r x^{(t)} + b_r\right)$$

**New hidden state content:** reset gate selects useful parts of prev hidden state. Use this and current input to compute new hidden content.

$$\tilde{h}^{(t)} = \tanh\left(W_h(r^{(t)} \circ h^{(t-1)}) + U_h x^{(t)} + b_h\right)$$

$$h^{(t)} = (1 - u^{(t)}) \circ h^{(t-1)} + u^{(t)} \circ \tilde{h}^{(t)}$$

**Hidden state:** update gate simultaneously controls what is kept from previous hidden state, and what is updated to new hidden state content

**How does this solve vanishing gradient?**
Like LSTM, GRU makes it easier to retain info long-term (e.g. by setting update gate to 0)

19

# Vanishing gradient: LSTM vs GRU

- LSTM and GRU are both great
    - GRU is quicker to compute and has fewer parameters than LSTM
    - There is no conclusive evidence that one consistently performs better than the other
    - LSTM is a good default choice (especially if your data has particularly long dependencies, or you have lots of training data)

**Rule of thumb**: start with LSTM, but switch to GRU if you want something more efficient

# Vanishing gradient in non-RNN

Vanishing gradient is present in **all** deep neural network architectures.

- Due to chain rule / choice of nonlinearity function, gradient can become vanishingly small during backpropagation
- Lower levels are hard to train and are trained slower
- **Potential solution**: direct (or skip-) connections (just like in ResNet)
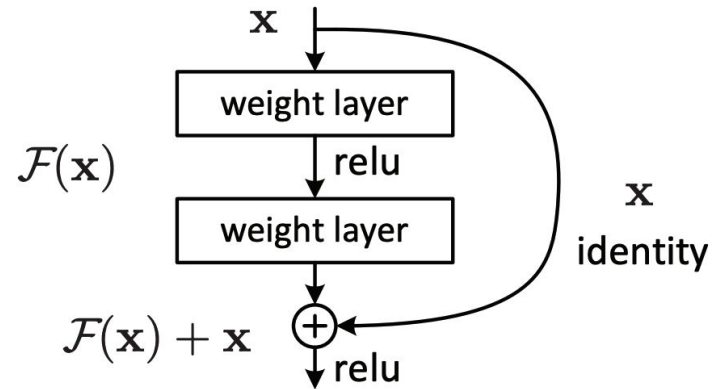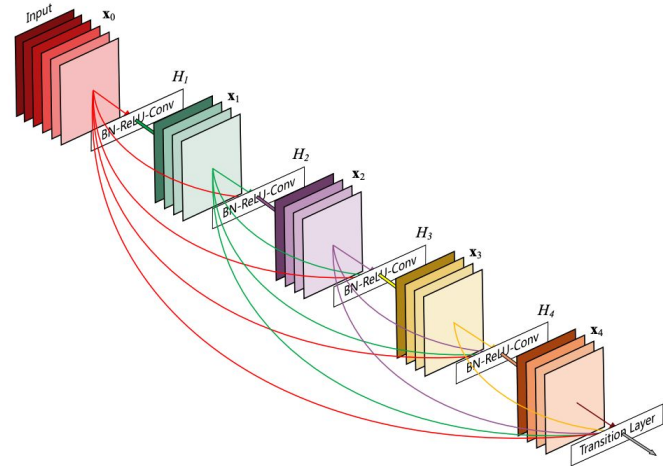
Figure 2. Residual learning: a building block.

Source: "Deep Residual Learning for Image Recognition", He et al, 2015. https://arxiv.org/pdf/1512.03385.pdf
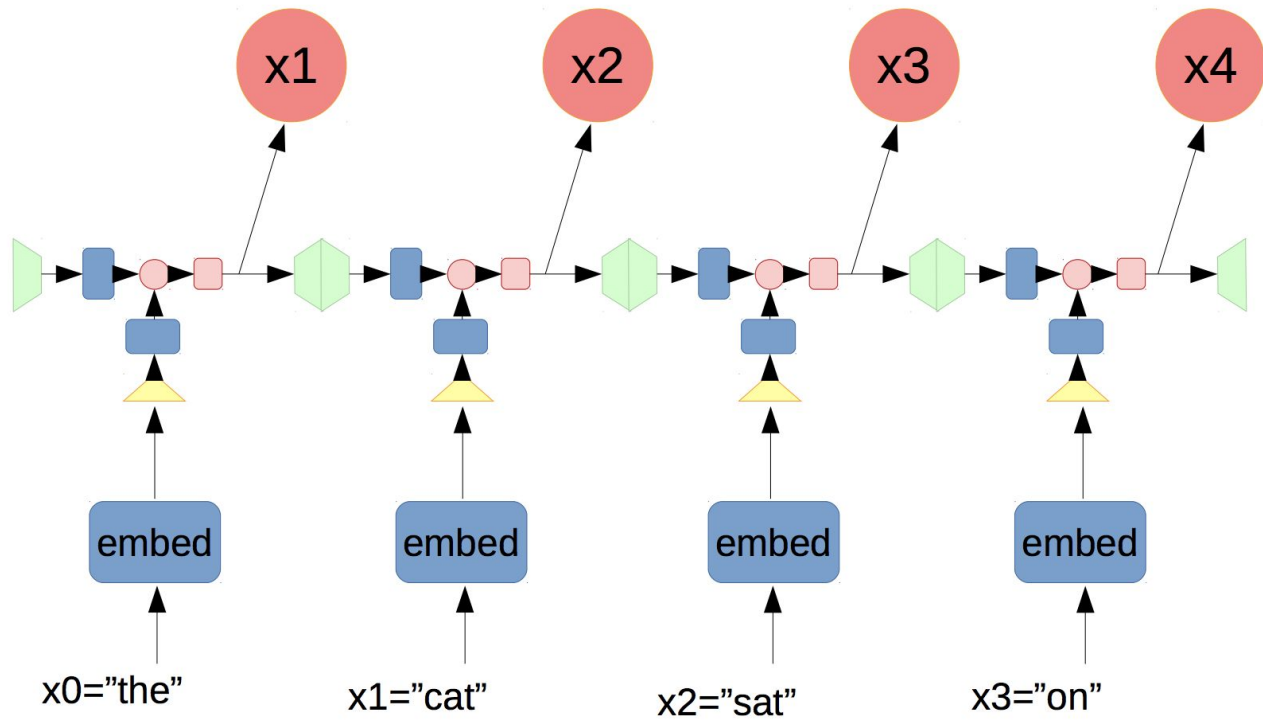
# Vanishing gradient in non-RNN

Vanishing gradient is present in **all** deep neural network architectures.

- Due to chain rule / choice of nonlinearity function, gradient can become vanishingly small during backpropagation
- Lower levels are hard to train and are trained slower
- **Potential solution**: dense connections (just like in DenseNet)

Source: "Densely Connected Convolutional Networks", Huang et al, 2017 https://arxiv.org/pdf/1608.06993.pdf

# Vanishing gradient in non-RNN

Vanishing gradient is present in **all** deep neural network architectures.

- Due to chain rule / choice of nonlinearity function, gradient can become vanishingly small during backpropagation
- Lower levels are hard to train and are trained slower
- **Potential solution(but not actually for that problem)**: dense connections (just like in DenseNet)

**Conclusion:**
   *Though vanishing/exploding gradients are a general problem, RNNs are particularly unstable due to the repeated multiplication by the same weight matrix* [Bengio et al, 1994]. Gradients magnitude drops exponentially with connection length.

Source:  "Learning Long-Term Dependencies with Gradient Descent is Difficult", Bengio et al. 1994, http://ai.dinfo.unifi.it/paolo//ps/tnn-94-gradient.pdf

x1  x2  x3  x4

embed  embed  embed  embed
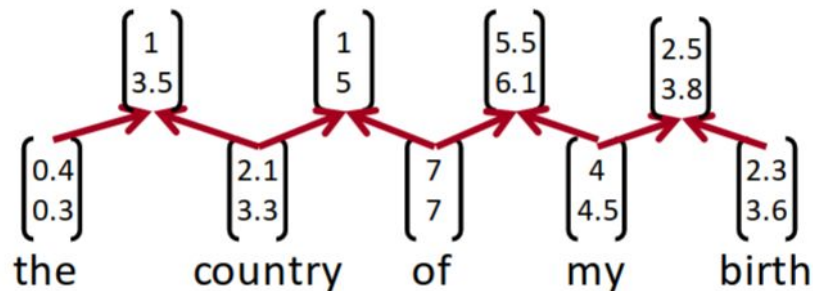
x0="the"  x1="cat"  x2="sat"  x3="on"

# From RNN to CNN

- RNN: Get compositional vectors for grammatical phrases only

- CNN: What if we compute vectors for every possible phrase?
  - Example: *"the country of my birth"* computes vectors for:
    - *the country, country of, of my, my birth, the country of, country of my, of my birth, the country of my, country of my birth*



- Regardless of whether it is grammatical
- Wouldn't need parser
- Not very linguistically or cognitively plausible

Based on: Lecture by Richard Socher 5/12/16, http://cs224d.stanford.edu

# From RNN to CNN

- Imagine using only bigrams



$$\begin{bmatrix} 1 \\ 3.5 \end{bmatrix} \quad \begin{bmatrix} 1 \\ 5 \end{bmatrix} \quad \begin{bmatrix} 5.5 \\ 6.1 \end{bmatrix} \quad \begin{bmatrix} 2.5 \\ 3.8 \end{bmatrix}$$

$$\begin{bmatrix} 0.4 \\ 0.3 \end{bmatrix} \quad \begin{bmatrix} 2.1 \\ 3.3 \end{bmatrix} \quad \begin{bmatrix} 7 \\ 7 \end{bmatrix} \quad \begin{bmatrix} 4 \\ 4.5 \end{bmatrix} \quad \begin{bmatrix} 2.3 \\ 3.6 \end{bmatrix}$$

the    country   of    my    birth

- Same operation as in RNN, but for every pair

$$p = \tanh\left(W \begin{bmatrix} c_1 \\ c_2 \end{bmatrix} + b\right)$$

- Can be interpreted as convolution over the word vectors
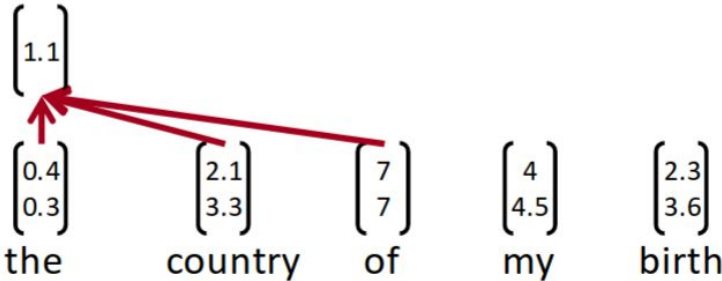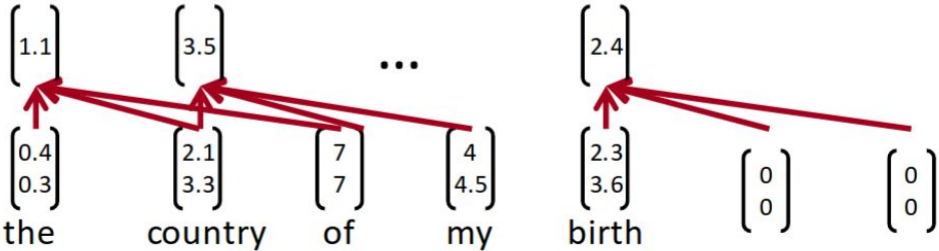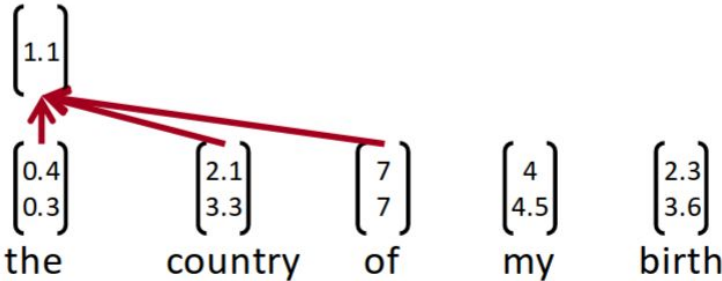
Based on: Lecture by Richard Socher 5/12/16, http://cs224d.stanford.edu

# From RNN to CNN

- Simple convolution + pooling
- Window size may be different (2 or more)
- The feature map based on bigrams:

$$c_i = f(\mathbf{w}^T \mathbf{x}_{i:i+h-1} + b)$$

$$\mathbf{c} = [c_1, c_2, \ldots, c_{n-h+1}] \in \mathbb{R}^{n-h+1}$$

Based on: Lecture by Richard Socher 5/12/16, http://cs224d.stanford.edu
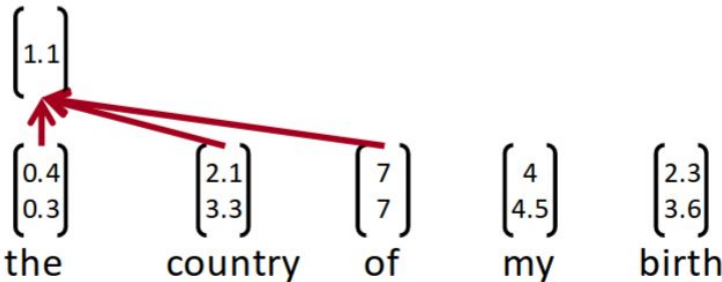
# One layer CNN

- Simple convolution + pooling
- Window size may be different (2 or more)
- The feature map based on bigrams:

$$c_i = f\left(\mathbf{w}^T \mathbf{x}_{i:i+h-1} + b\right)$$

$$\mathbf{c} = [c_1, c_2, \ldots, c_{n-h+1}] \in \mathbb{R}^{n-h+1}$$

Based on: Lecture by Richard Socher 5/12/16, http://cs224d.stanford.edu

# One layer CNN

- Simple convolution + pooling
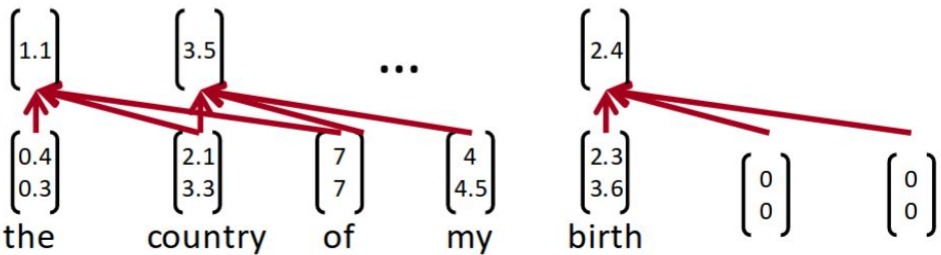- Window size may be different (2 or more)
- The feature map based on bigrams:

$$c_i = f(\mathbf{w}^T \mathbf{x}_{i:i+h-1} + b)$$

$$\mathbf{c} = [c_1, c_2, \ldots, c_{n-h+1}] \in \mathbb{R}^{n-h+1}$$

What's next?

Based on: Lecture by Richard Socher 5/12/16, http://cs224d.stanford.edu

# One layer CNN

- Simple convolution + pooling
- Window size may be different (2 or more)
- The feature map based on bigrams:

$$c_i = f(\mathbf{w}^T \mathbf{x}_{i:i+h-1} + b)$$

$$\mathbf{c} = [c_1, c_2, \ldots, c_{n-h+1}] \in \mathbb{R}^{n-h+1}$$

What's next?

We need more features!

Based on: Lecture by Richard Socher 5/12/16, http://cs224d.stanford.edu

- Feature representation is based on some applied filter:

$$\mathbf{c} = [c_1, c_2, \ldots, c_{n-h+1}] \in \mathbb{R}^{n-h+1}$$

- Let's use pooling over the time axis: $\quad \hat{c} = \max\{\mathbf{c}\}$

- Now the length of **c** is irrelevant!

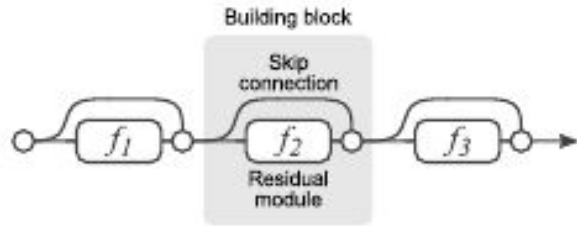- So we can use filters based on unigrams, bigrams, tri-grams, 4-grams, etc.

Based on: Lecture by Richard Socher 5/12/16, http://cs224d.stanford.edu

# Another example from Kim (2014) paper



wait
for
the
video
and
do
n't
rent
it

n x k representation of sentence with static and non-static channels

Convolutional layer with multiple filter widths and feature maps

Max-over-time pooling

Fully connected layer with dropout and softmax output
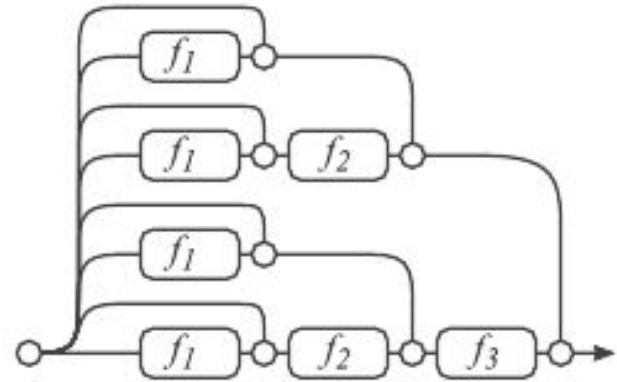
Based on: Lecture by Richard Socher 5/12/16, http://cs224d.stanford.edu

# Another view on ResNets and vanishing gradient

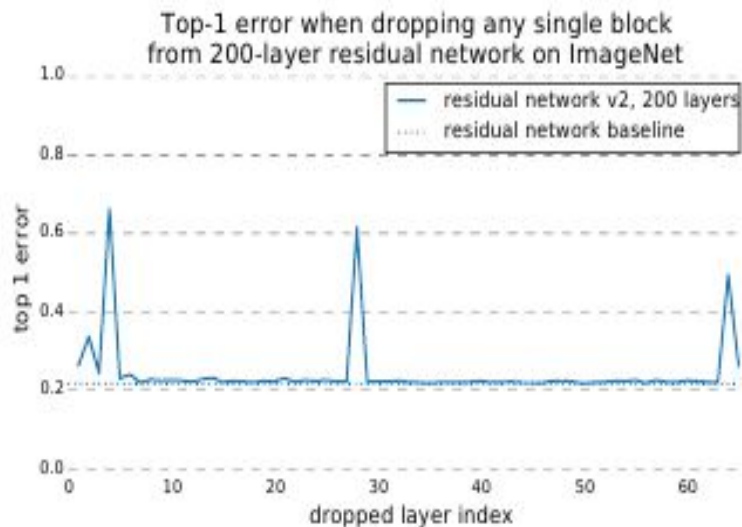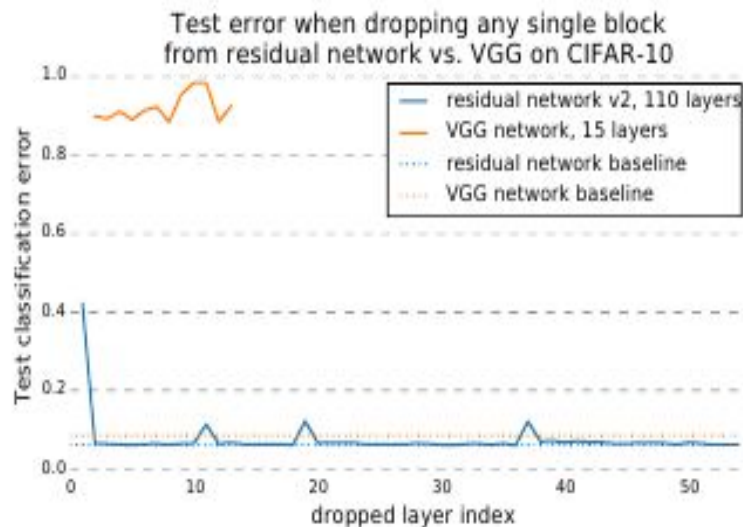"**Residual Networks Behave Like Ensembles of Relatively Shallow Networks**"



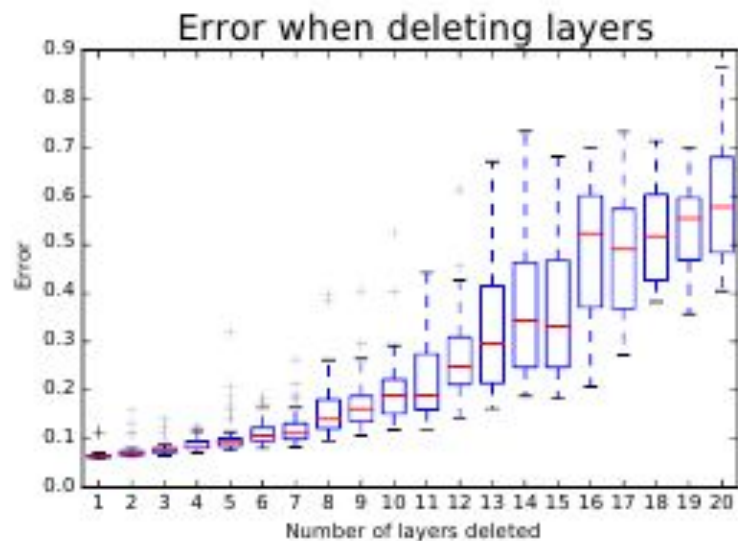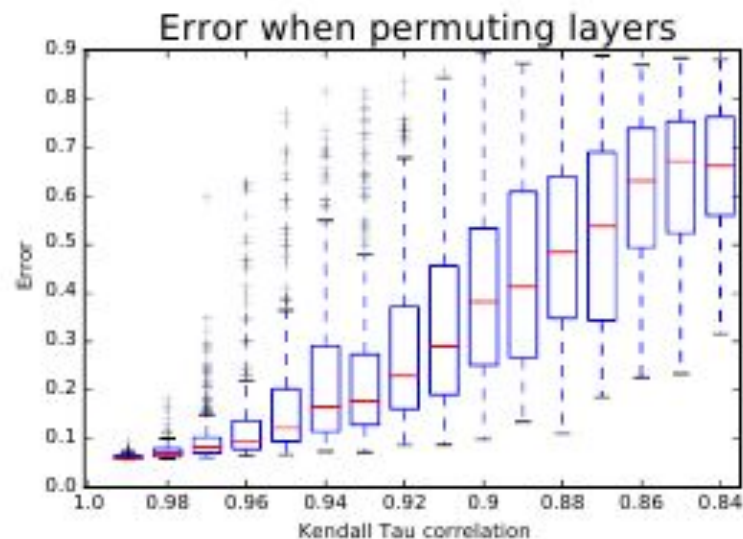(a) Conventional 3-block residual network

(b) Unraveled view of (a)

Source: https://arxiv.org/pdf/1605.06431.pdf

# "**Residual Networks Behave Like Ensembles of Relatively Shallow Networks**"



Test error when dropping any single block from residual network vs. VGG on CIFAR-10

- residual network v2, 110 layers
- VGG network, 15 layers
- residual network baseline
- VGG network baseline

Top-1 error when dropping any single block from 200-layer residual network on ImageNet

- residual network v2, 200 layers
- residual network baseline

https://arxiv.org/pdf/1605.06431.pdf

34

## "**Residual Networks Behave Like Ensembles of Relatively Shallow Networks**"



(a) Error when deleting layers

(b) Error when permuting layers

# "**Residual Networks Behave Like Ensembles of Relatively Shallow Networks**"

Fixed-size vector of features

Class label /
Probability distribution

From YSDA nlp course 2018

# Recurrent neural networks for texts

From YSDA nlp course 2018

# Convolutional neural networks for texts

$(100,9)$

$x_0$ The  $x_1$ quick  $x_2$ brown  $x_3$ fox  $x_4$ jumps  $x_5$ over  $x_6$ the  $x_7$ dog  $x_8$ .
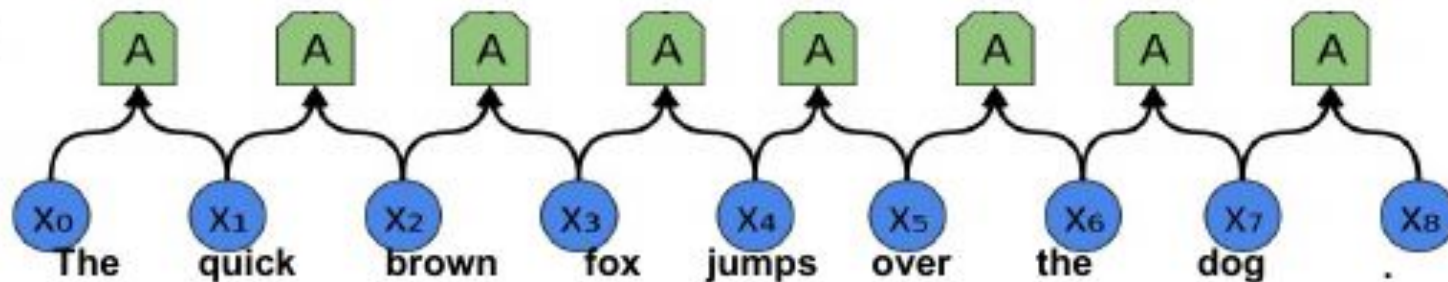
# CNN for texts

A convolution kernel is a tensor of size
[output dim, embedding dim, kernel size]

1d-convolution
32x(100x2)

(100,9)

$X_0$ The  $X_1$ quick  $X_2$ brown  $X_3$ fox  $X_4$ jumps  $X_5$ over  $X_6$ the  $X_7$ dog  $X_8$ .
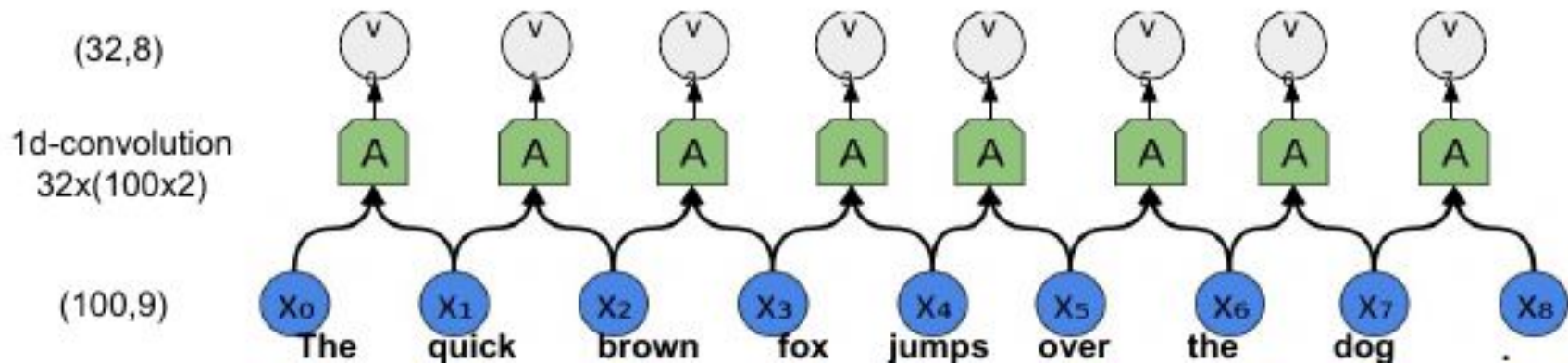
From YSDA nlp course 2018

$$\mathbf{v}_0 = \mathbf{A}(\mathbf{x}_0, \mathbf{x}_1)$$
$$v_{0i} = A_i(\mathbf{x}_0, \mathbf{x}_1) = \sum_j (K_{0ij} x_{0j} + K_{1ij} x_{1j})$$

(32,8)

1d-convolution
32x(100x2)

(100,9)

| The | quick | brown | fox | jumps | over | the | dog | . |

From YSDA nlp course 2018

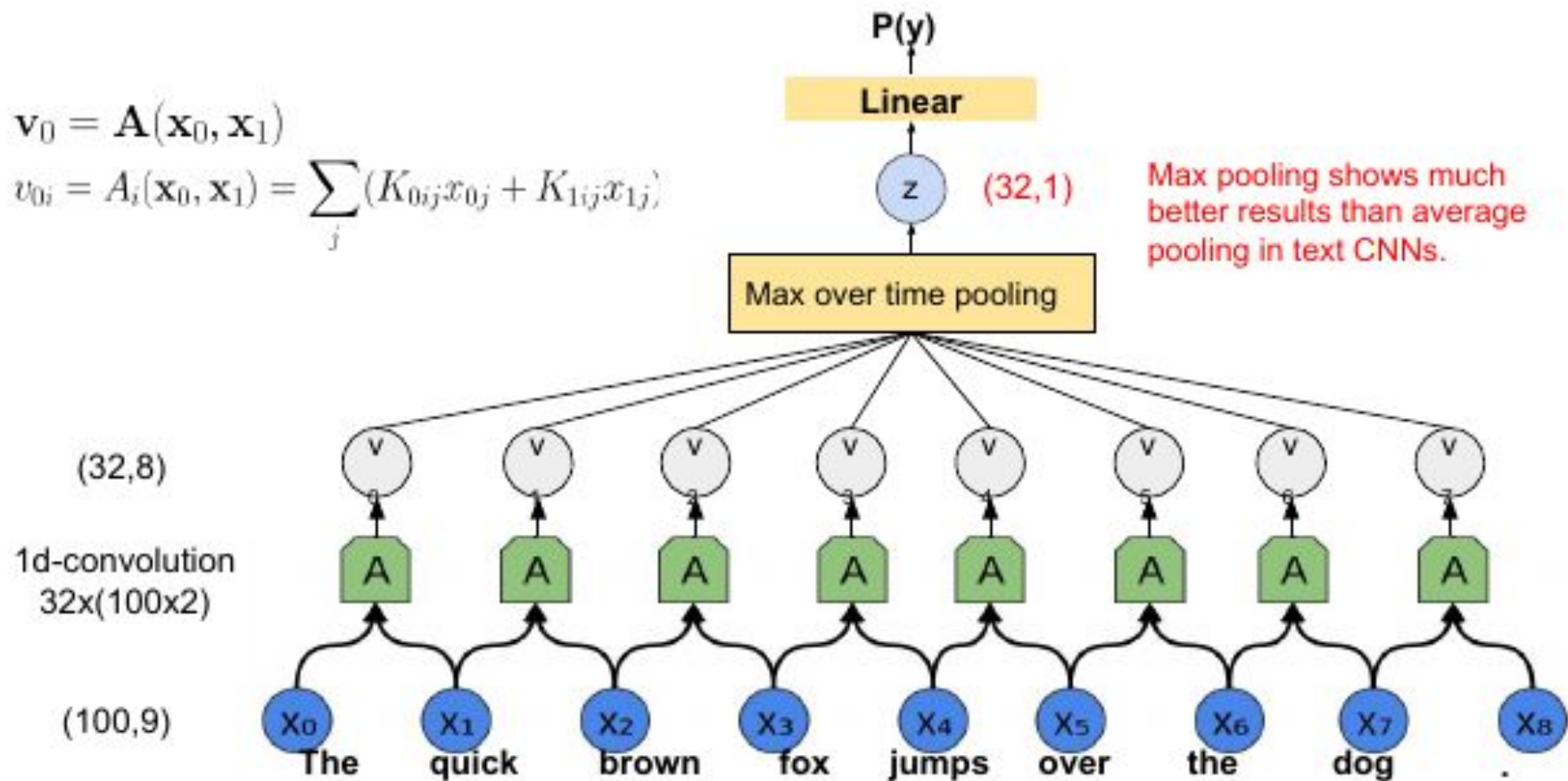$$\mathbf{v}_0 = \mathbf{A}(\mathbf{x}_0, \mathbf{x}_1)$$
$$v_{0i} = A_i(\mathbf{x}_0, \mathbf{x}_1) = \sum_j (K_{0ij}x_{0j} + K_{1ij}x_{1j})$$

From YSDA nlp course 2018

$$\mathbf{v}_0 = \mathbf{A}(\mathbf{x}_0, \mathbf{x}_1)$$

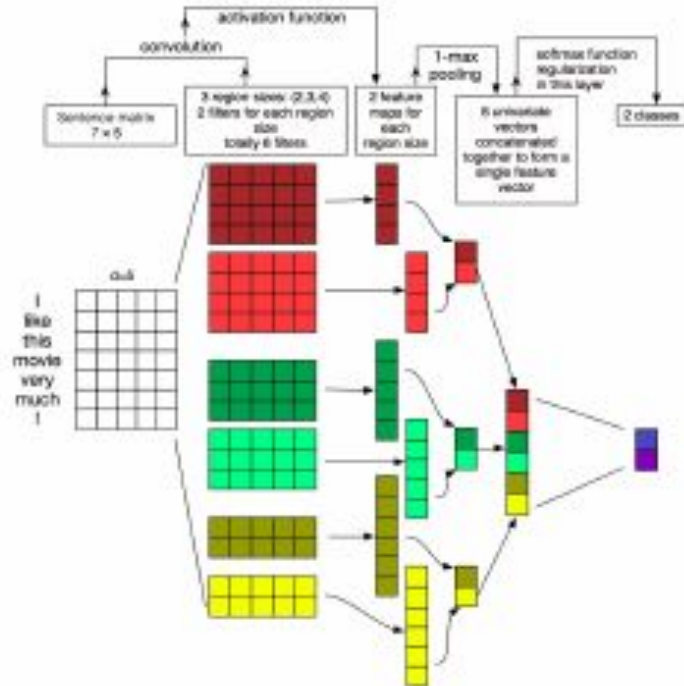$$v_{0i} = A_i(\mathbf{x}_0, \mathbf{x}_1) = \sum_j (K_{0ij} x_{0j} + K_{1ij} x_{1j})$$

P(y)

**Linear**

z   (32,1)

Max pooling shows much better results than average pooling in text CNNs.

Max over time pooling

(32,8)

1d-convolution
32x(100x2)

(100,9)

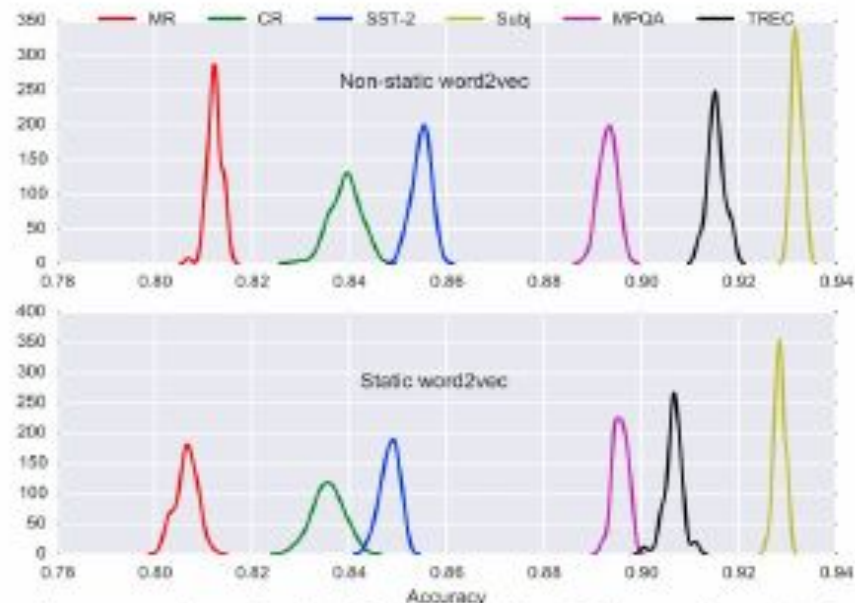| The | quick | brown | fox | jumps | over | the | dog | . |
|-----|-------|-------|-----|-------|------|-----|-----|---|
| X₀  | X₁    | X₂    | X₃  | X₄    | X₅   | X₆  | X₇  | X₈|

From YSDA nlp course 2018

43

# CNN for texts: Improvements



- Use convolutional layers with different kernel size, separate max-pooling over time and concatenation.
- K-max pooling: take not 1 but k highest activations in their original order.
  E.g. $(0,1,3,2,0,1,4,1) \rightarrow (3,2,4)$

Zhang et al.
https://arxiv.org/abs/1510.03820

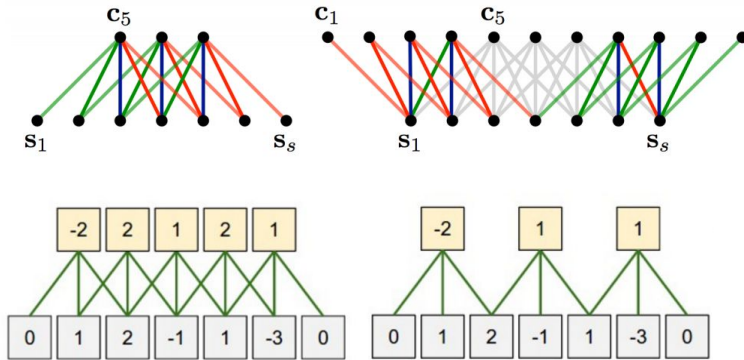From YSDA nlp course 2018,

# CNN for texts: Improvements



Accuracy density plots for non-static w2v (upper) and static w2v (lower) [for 10-fold CV over the 100 replications ]

- Use convolutional layers with different kernel size, separate max-pooling over time and concatenation.
- K-max pooling: take not 1 but k highest activations in their original order.
E.g. (0,1,3,2,0,1,4,1) → (3,2,4)
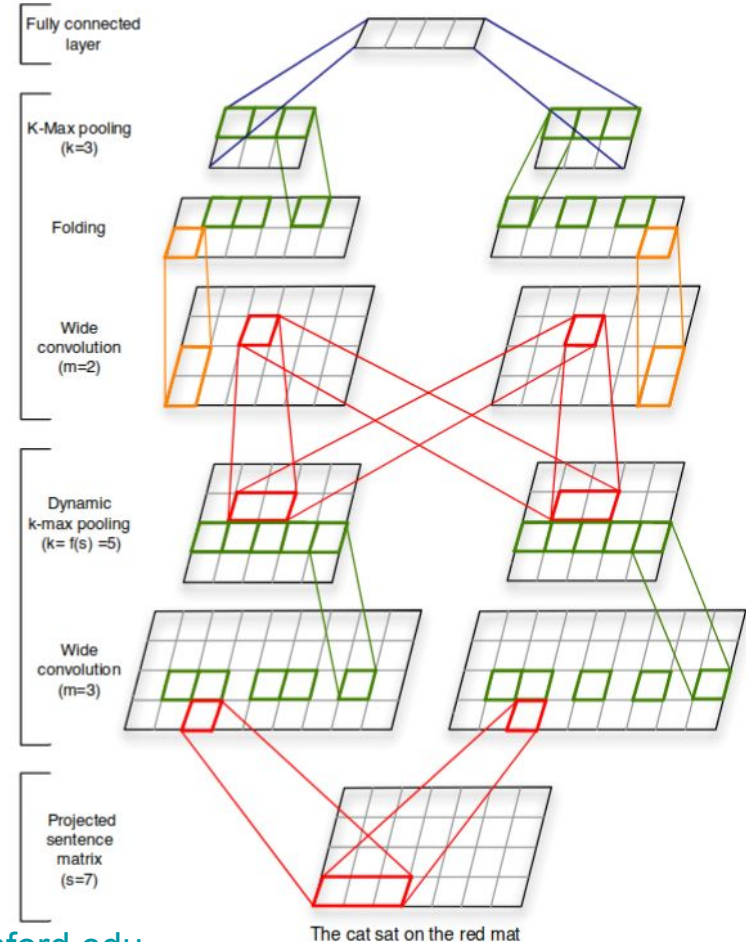- Use pre-trained word vectors only for embedding layer initialization, train it jointly with model

Zhang et al.
https://arxiv.org/abs/1510.03820

From YSDA nlp course 2018,

- Narrow vs wide convolution (stride and zero-padding)

- Complex pooling schemes over sequences
- Great readings (e.g. Kalchbrenner et. al. 2014)

Based on: Lecture by Richard Socher 5/12/16, http://cs224d.stanford.edu

46

# CNN applications

P( f l e )

- Neural machine translation: CNN as encoder, RNN as decoder
- Kalchbrenner and Blunsom (2013) "Recurrent Continuous Translation Models"
- One of the first neural machine translation efforts

S

e

csm

e

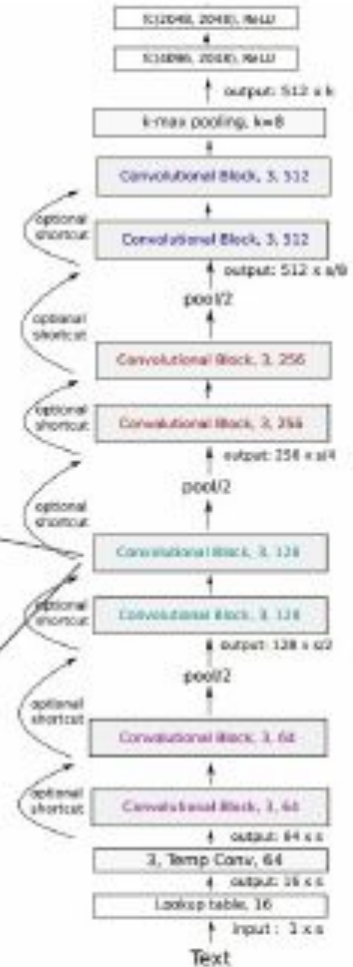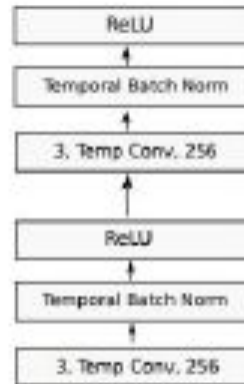Based on: Lecture by Richard Socher 5/12/16, http://cs224d.stanford.edu

# Deep Convolutional networks for texts

Q: Can we get some quality points just stacking much more layers?

A: It does make sense in case character-level convolutional architectures.

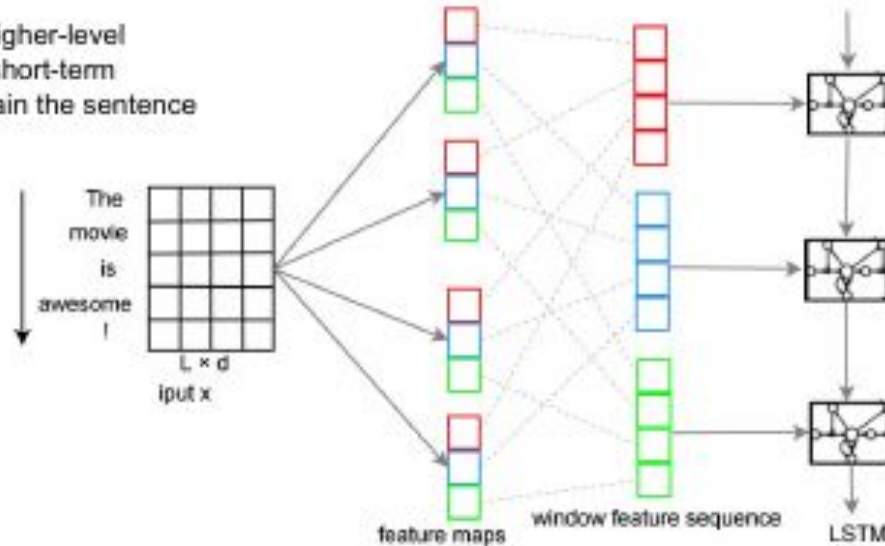VDCNN [Conneau et al. 2015] ~ ResNet-like network with 29 conv. layers

# We can combine CNN and RNN together

C-LSTM [Zhou et al. 2015]

[conv.]->[LSTM]

C-LSTM utilizes CNN to extract a sequence of higher-level phrase representations, and are fed into a long short-term memory recurrent neural network (LSTM) to obtain the sentence representation.

Based on: Lecture by Richard Socher 5/12/16, http://cs224d.stanford.edu

# Text augmentations

Like with images we can increase our training corpora size with augmentations.

Examples:

- Text deformations: mix some texts, change order…
- Reformulations
- Word dropout

# Text segmentation

Open vocabulary problem: In NLP language vocabulary is usually very big. To produce a good quality with particular word the algorithm should see a lot of examples with it. This is a big problem for rare words.

There are two extreme approaches for vocabulary modelling:

- Char level: small vocabulary, a lot of examples with each element, slow training, long sequences during encoding and decoding
- Each word is a new item in vocabulary: Big vocabulary, small number of examples for rare words, fast training, short sequences during encoding and decoding

# Text segmentation: Balance, BPE

We can balance vocabulary size with length of sequence

Bait Pair Encoding (BPE) (Sennrich et al.): Let's split rare words into subwords, while leave frequent sequences as a one token.

1) Compute merge table: Starting from characters let's one by one merge the most frequent symbols into one symbol until reaching desired vocabulary size.
2) During inference let's greedily (priority=number of step, when this pair was added in (1)) apply merge rules

https://arxiv.org/pdf/1508.07909.pdf

- We have not got out-of-vocabulary words, because we start from all characters.
- We can balance vocabulary size with decoding efficiency

  Example:

  "mother" -> (BPE) mother

  "sweetish" -> (BPE) sweet  ish

  "asft" -> (BPE) as  f  t

- Vanishing gradient is present not only in RNNs
    - Use some kind of memory or skip-connections
- LSTM and GRU are both great
    - GRU is quicker, LSTM catch more complex dependencies
- Rule of thumb: start with LSTM, but switch to GRU if you want something more efficient
- Clip your gradients
- Combining RNN and CNN worlds? Why not ;)

That's all. Feel free to ask any questions.