# 291K
# Deep Learning for Machine Translation
# Recurrent Neural Network
# Sequence to sequence learning

Lei Li

UCSB

10/4/2021

# Outline

- Language Modeling
- Recurrent Neural Network
- Long-short term memory network (LSTM)
- Gated Recurrent Unit (GRU)
- Attention
- Encoder-decoder framework
- LSTM Seq2seq

# **Language Modeling**

- Given a sentence y, estimate the probability

$$P(y) = \prod_t P(y_{t+1} \mid y_1 \ldots y_t)$$

$$- P(y_{t+1} \mid y_1 \ldots y_t) = f_\theta(y_1, \ldots, y_t)$$

$$p(y_6 \mid y_1, \ldots, y_5)$$

The cat sits on a ___
$y_1$  $y_2$  $y_3$  $y_4$  $y_5$  $y_6$

| | |
|---|---|
| mat | 0.15 |
| rug | 0.13 |
| chair | 0.08 |
| hat | 0.05 |
| dog | 0.01 |

# **Predict Next Token Probability**

There are many methods to predict the next token:

- N-gram: assuming $p\left(x_t \mid x_1, \ldots, x_{t-1}\right) = p\left(x_t \mid x_{t-k}, \ldots, x_{t-1}\right)$, and estimate it directly

- Context MLP: use DNN to estimate $p\left(x_t \mid x_{t-k}, \ldots, x_{t-1}\right)$

- CNN-LM (previous lecture)
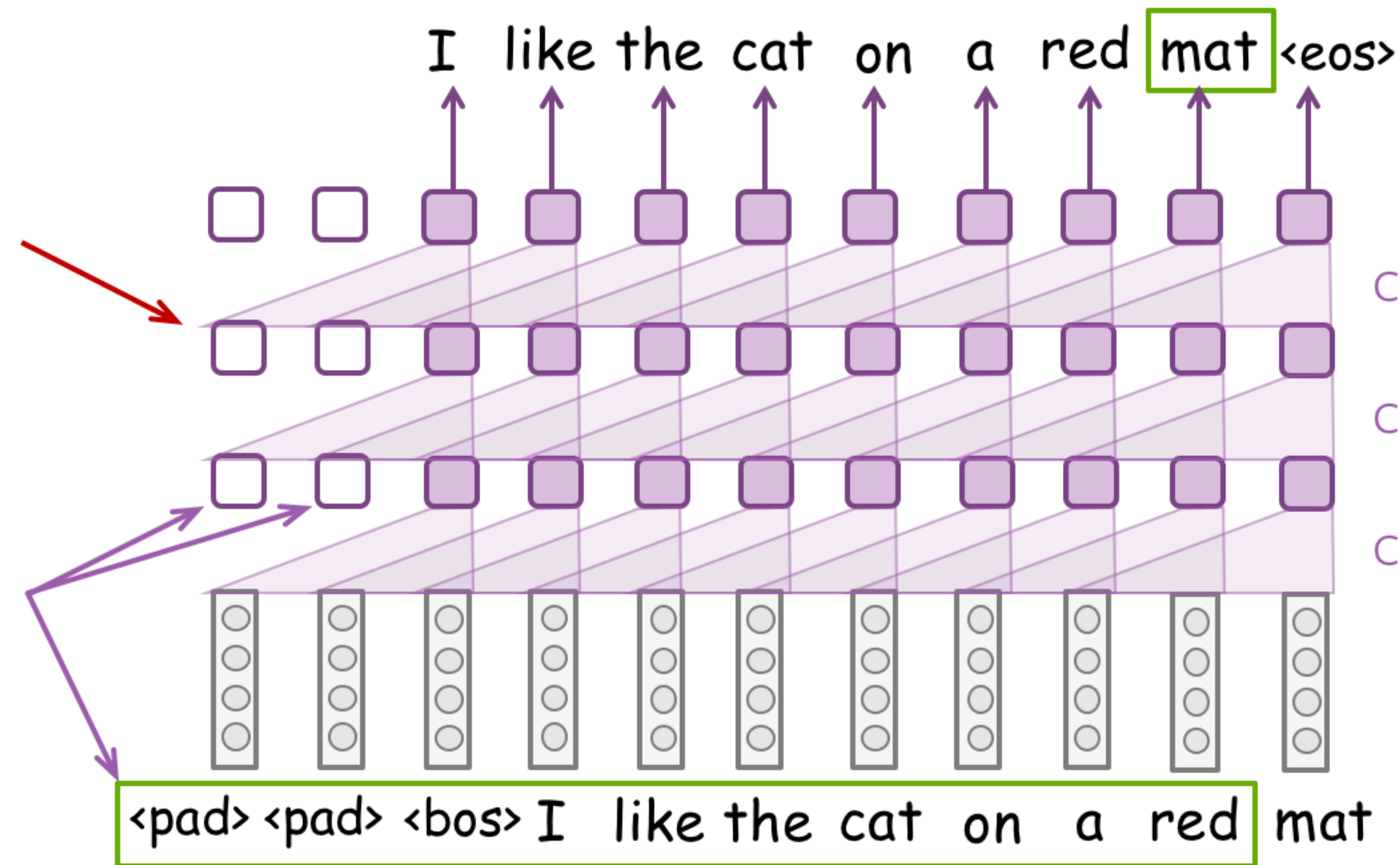- RNN-LM, LSTM, GRU
- GPT

# CNN Language Model (recap)

$$P(y_{t+1} \mid y_1, \ldots, y_t) \approx \mathrm{CNN}_\theta(y_{t-k}, \ldots, y_t)$$

predict the next token

But, limited context

No pooling between convolutions: do not want to lose positional information

Padding to shift tokens: we need to prevent information flow from future tokens

I like the cat on a red mat <eos>

convolution

convolution

convolution

Stack several convolutions

<pad> <pad> <bos> I like the cat on a red mat

condition on the previous tokens

# Limitation of CNN-LM

- CNN-LM only has a fixed-length receptive field
  – probability of next token only dependent on a fixed-size context
- But sentences are of variable length
- How to handle sentences with variable length?
- Idea:
  – adding memory to network
  – adaptive updating memory

# Recurrent Memory

- Introduce memory representation
- RNN-LM: use RNN to estimate

$$p\left(x_t \mid x_1, \ldots, x_{t-1}\right) = \text{softmax}(W \cdot h_t)$$
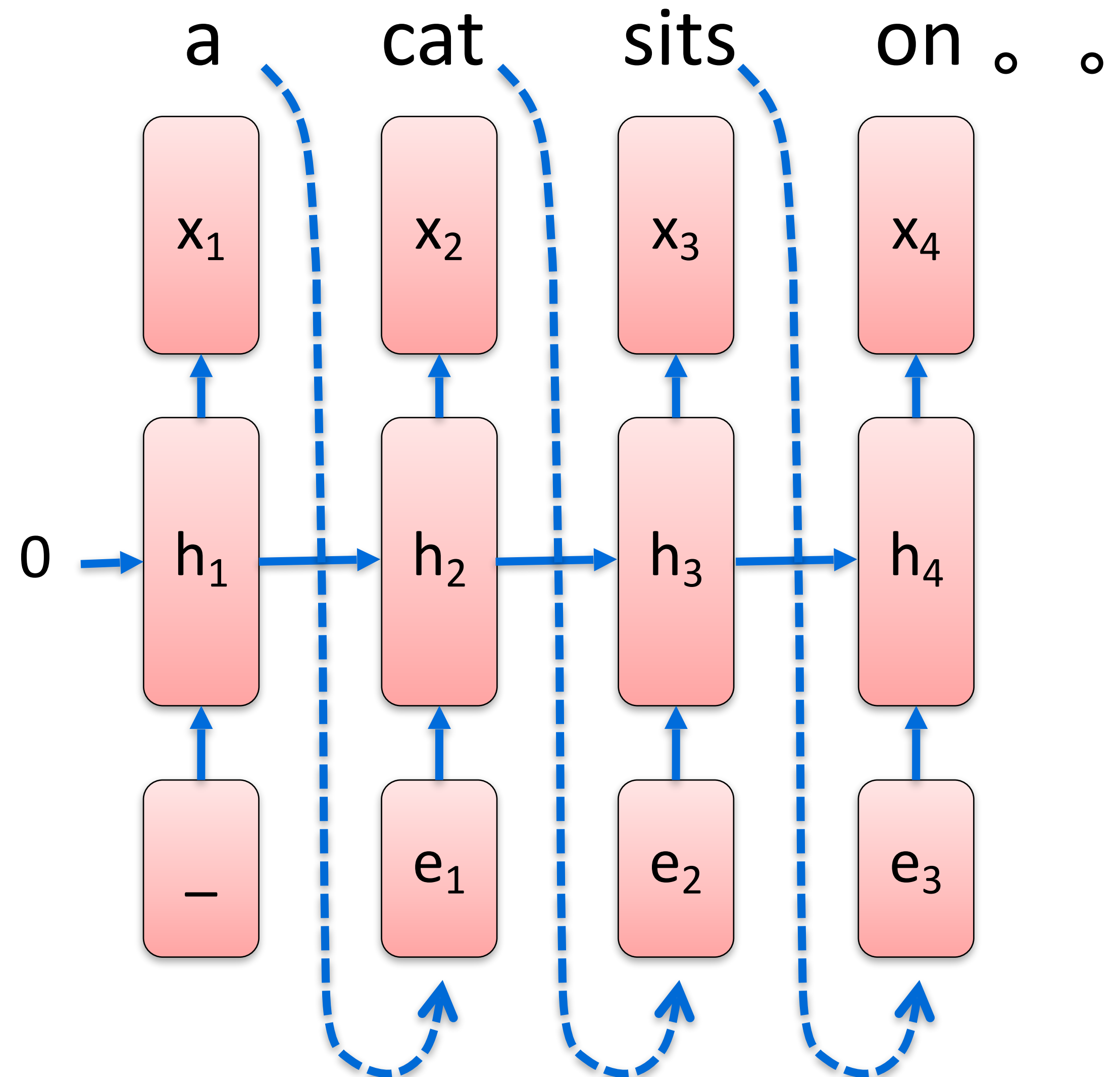
$$h_t = RNN(h_{t-1}, \ Emb(x_{t-1}))$$

- RNN cell can be
  - Simple feedforward neural network
  - Long-short term memory
  - Gated recurrent units

# Recurrent Neural Network

$$p(x_t \mid x_1, \ldots, x_{t-1}) = \text{softmax}(U \cdot h_t)$$

$$h_t = \sigma \left( W \cdot \begin{bmatrix} h_{t-1} \\ x_t \end{bmatrix} + b \right)$$



Elman, Finding Structure in Time. Cog. Sci. 1990.     Mikolov et al, Recurrent neural network based language model. Interspeech 2010.     8

# Training RNN-LM

- Risk:
  - Loss: cross-entropy for every next-token given prefix context
  - CE(x_t+1, f(x_1, …, x_t))
- SGD
  - Calculate gradient: Back-propogation through time (BPTT)
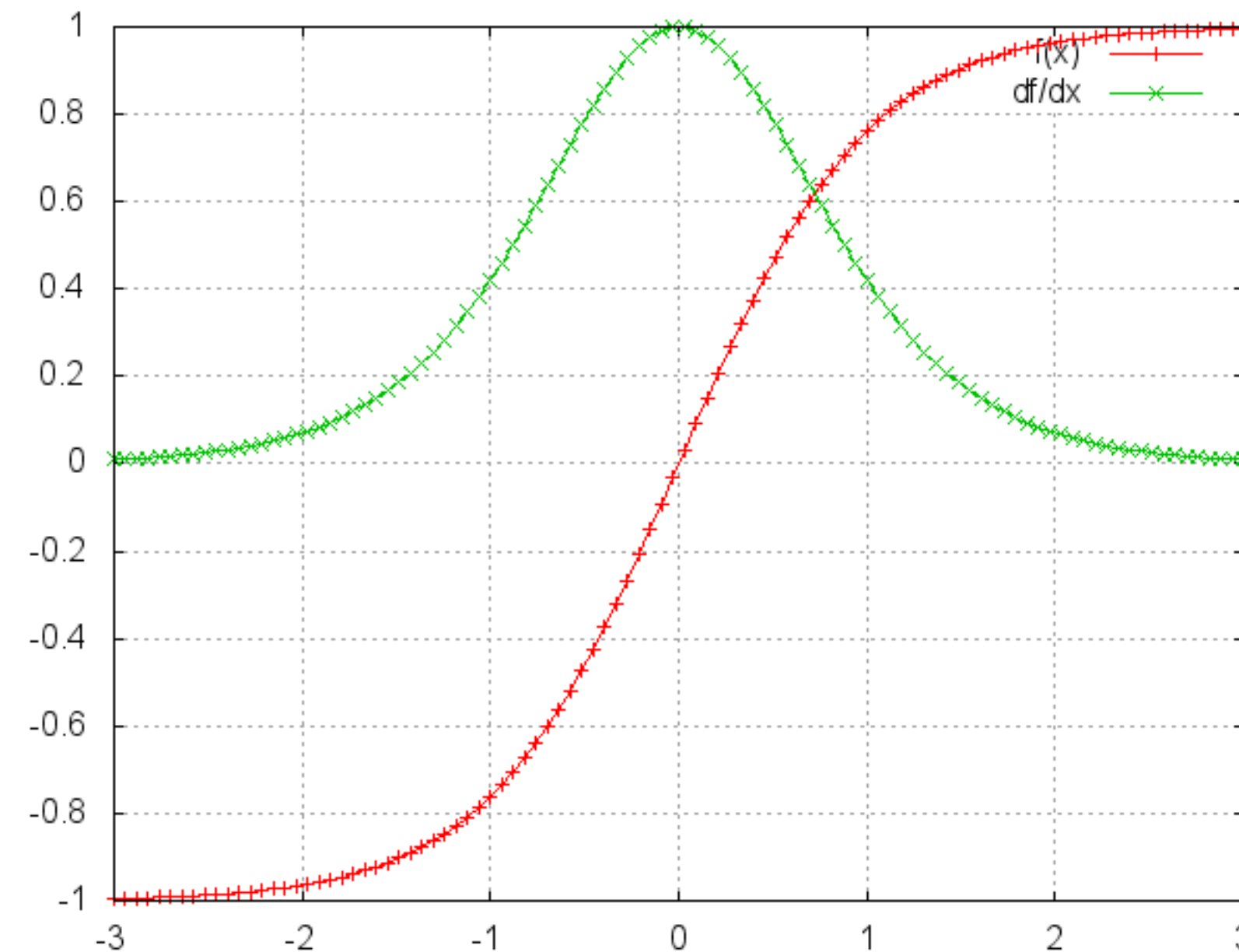  - $\nabla E_t$

# Exercise: Gradient for RNN

# Back-propagation for RNN (python)

```python
def bptt(self, x, y):
    T = len(y)
    # Perform forward propagation
    o, s = self.forward_propagation(x)
    # We accumulate the gradients in these variables
    dLdU = np.zeros(self.U.shape)
    dLdV = np.zeros(self.V.shape)
    dLdW = np.zeros(self.W.shape)
    delta_o = o
    delta_o[np.arange(len(y)), y] -= 1.
    # For each output backwards...
    for t in np.arange(T)[::-1]:
        dLdV += np.outer(delta_o[t], s[t].T)
        # Initial delta calculation: dL/dz
        delta_t = self.V.T.dot(delta_o[t]) * (1 - (s[t] ** 2))
        # Backpropagation through time (for at most self.bptt_truncate steps)
        for bptt_step in np.arange(max(0, t-self.bptt_truncate), t+1)[::-1]:
            # Add to gradients at each previous step
            dLdW += np.outer(delta_t, s[bptt_step-1])
            dLdU[:,x[bptt_step]] += delta_t
            # Update delta for next step dL/dz at t-1
            delta_t = self.W.T.dot(delta_t) * (1 - s[bptt_step-1] ** 2)
    return [dLdU, dLdV, dLdW]
```
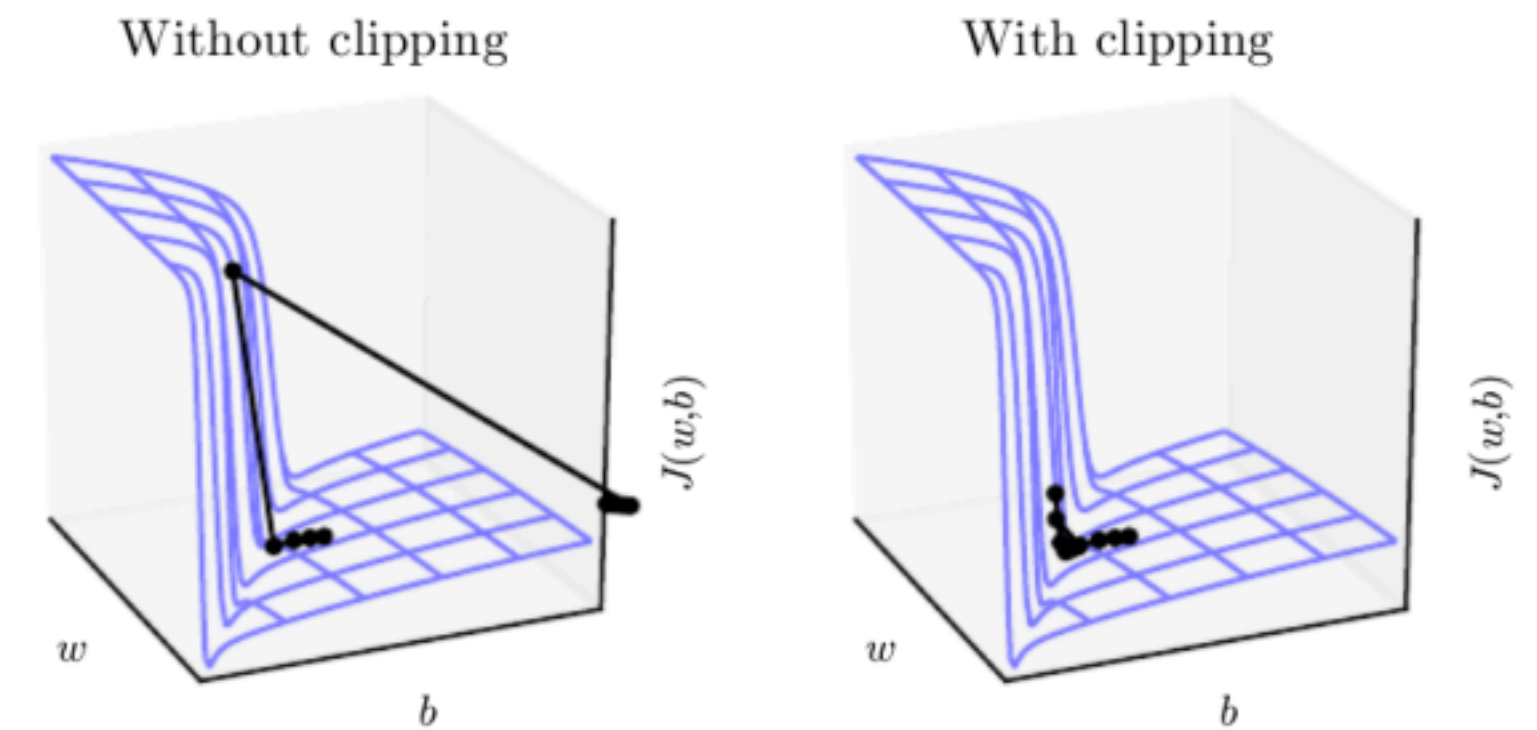
# Computational Issue: Gradient Vanishing

- tanh has derivative close to zero at both ends



Pascanu et al. On the difficulty of training recurrent neural networks. ICML 2013
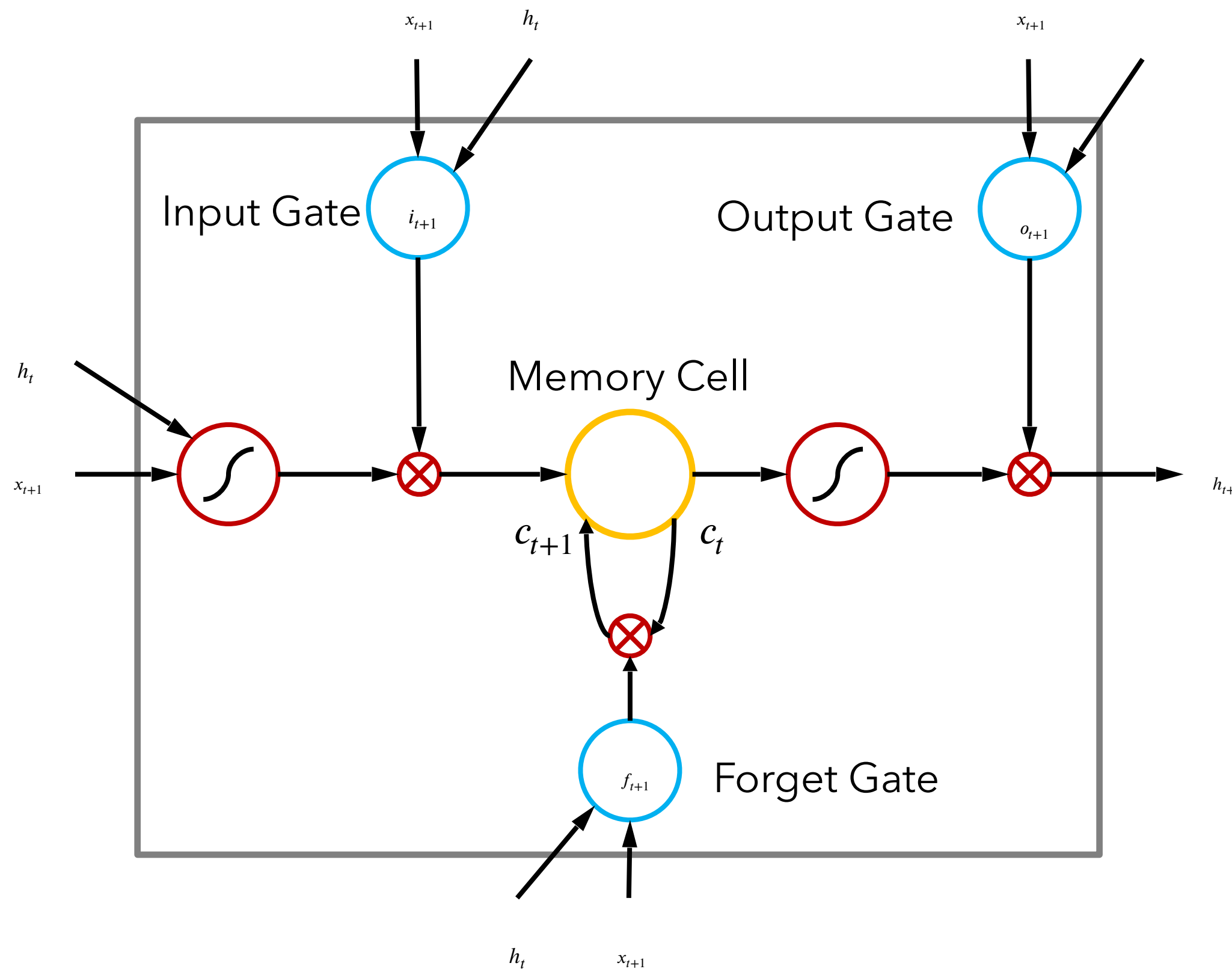
# **Gradient Exploding**

- Use gradient clipping
- Two options: clip by absolute value or rescale norm
- if $|g| > \eta, \hat{g} \leftarrow \eta$
- if $|g| > \eta, \hat{g} \leftarrow \dfrac{\eta}{|g|} g$

# Long-Short Term Memory (LSTM)

- Replace cell with more advanced one
- Adaptively memorize short and long term information



$$i_{t+1} = \sigma(M_{ix}x_{t+1} + M_{ih}h_t + b_i)$$
$$f_{t+1} = \sigma(M_{fx}x_{t+1} + M_{fh}h_t + b_f)$$
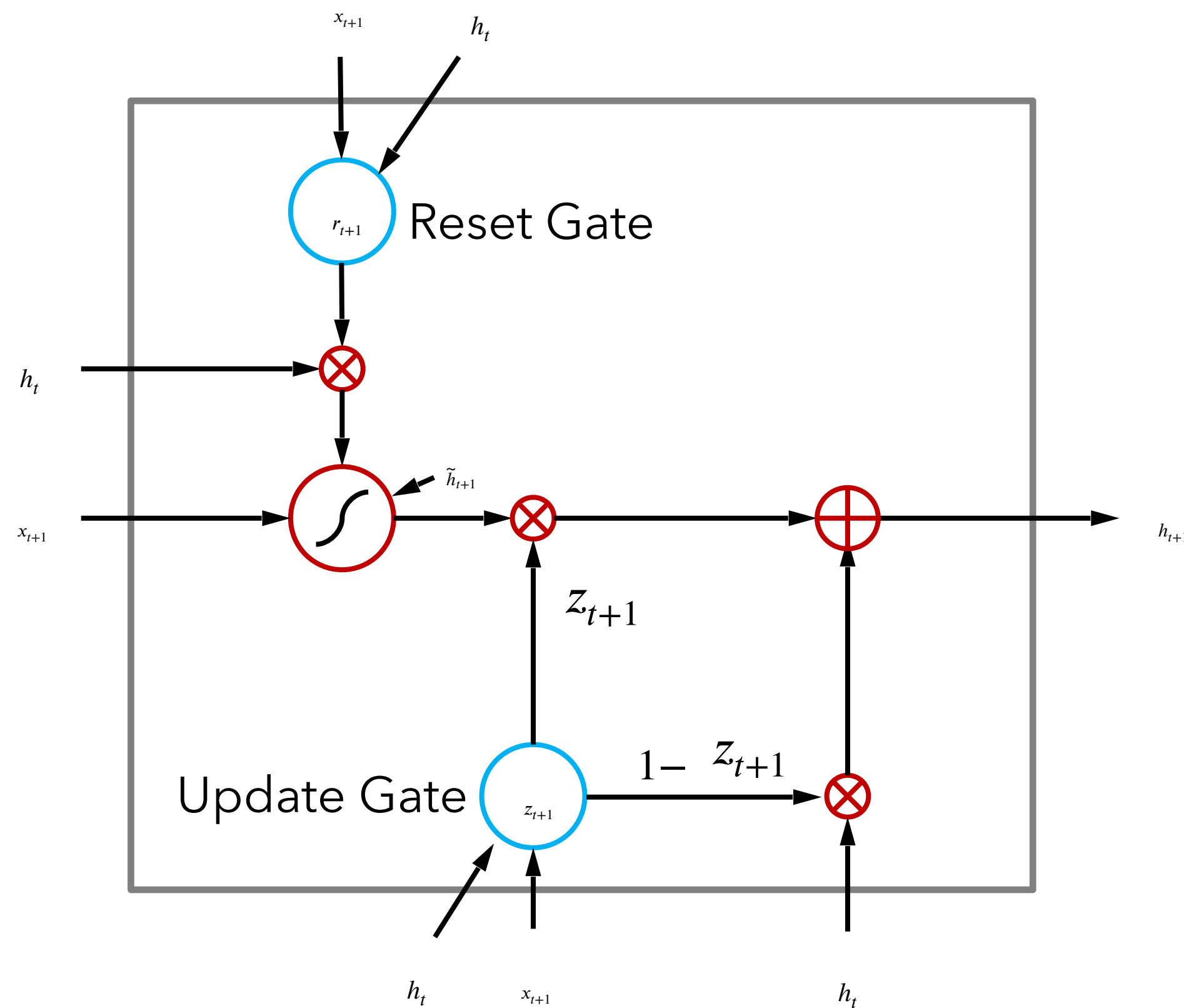$$o_{t+1} = \sigma(M_{ox}x_{t+1} + M_{oh}h_t + b_o)$$

$$a_{t+1} = \tanh(M_{cx}x_{t+1} + M_{ch}h_t + b_a)$$

$$c_{t+1} = \mathbf{\textcolor{green}{f_{t+1} \otimes c_t}} + i_{t+1} \otimes a_{t+1}$$
$$h_{t+1} = o_{t+1} \otimes \tanh(c_{t+1})$$

Hochreiter & Schmidhuber. Long Short-Term Memory, 1997
Gers et al. Learning to Forget: Continual Prediction with LSTM. 2000

# Gated Recurrent Unit (GRU)

- Adaptively memorize short and long term information
- like LSTM, but fewer parameters



Input: $x_t$
Memory: $h_t$

$$r_{t+1} = \sigma(M_{rx}x_{t+1} + M_{rh}h_t + b_r)$$
$$z_{t+1} = \sigma(M_{zx}x_{t+1} + M_{zh}h_t + b_z)$$

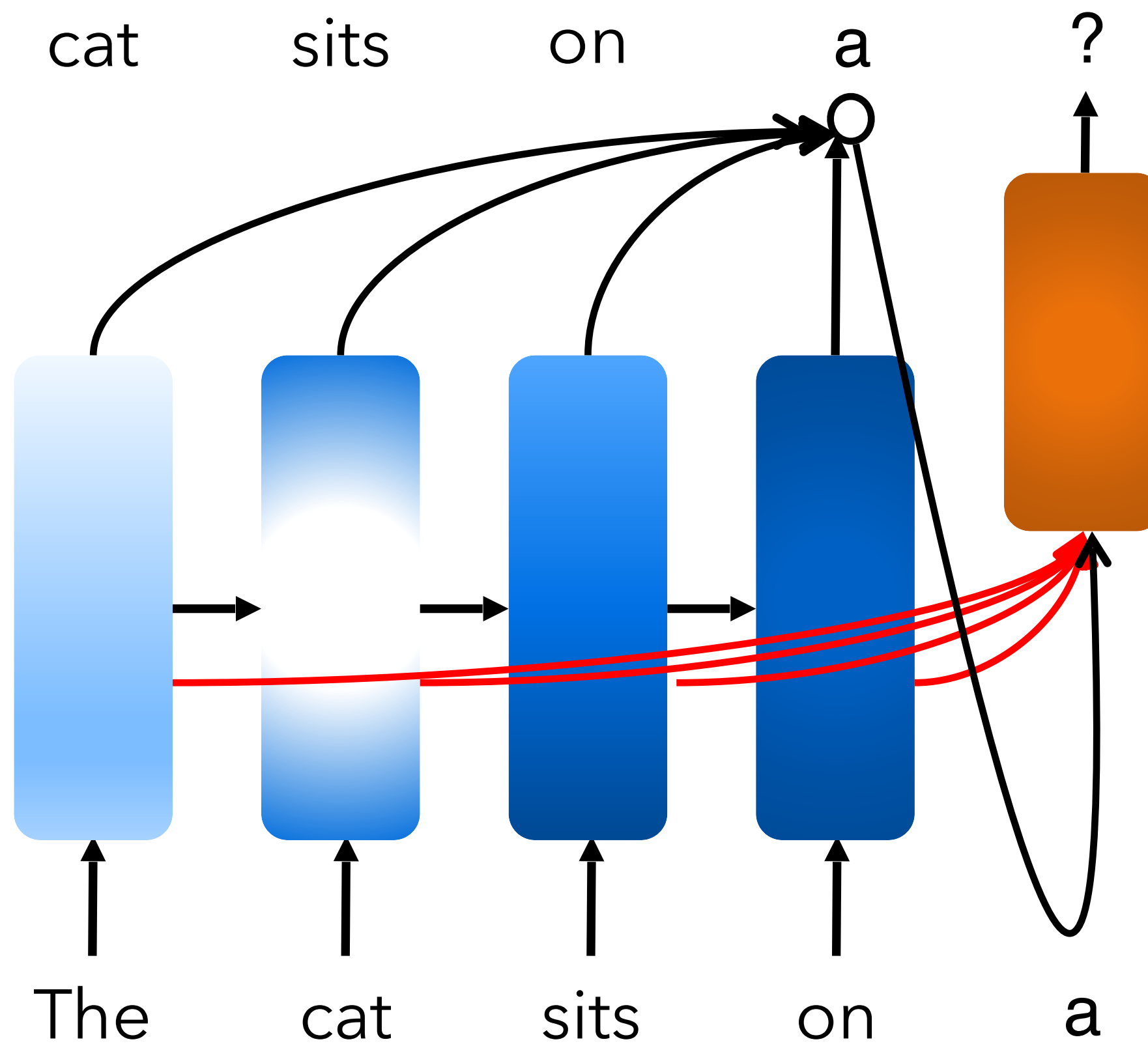$$\tilde{h}_{t+1} = \tanh(M_{hx}x_{t+1} + M_{hh}(r_{t+1} \otimes h_t) + b_h)$$

$$h_{t+1} = z_{t+1} \otimes \tilde{h}_{t+1} + (1 - z_{t+1}) \otimes h_t$$

Cho et al. Learning Phrase Representations using RNN Encoder–Decoder for Statistical Machine Translation. 2014

# Expressive power of RNN-LM

Perplexity: $PPL = P(x_1, \ldots, x_N)^{-\frac{1}{N}} = \exp(-\frac{1}{N}\sum_{n=1}^{N} \log P(x_n | x_1 \ldots x_{n-1}))$

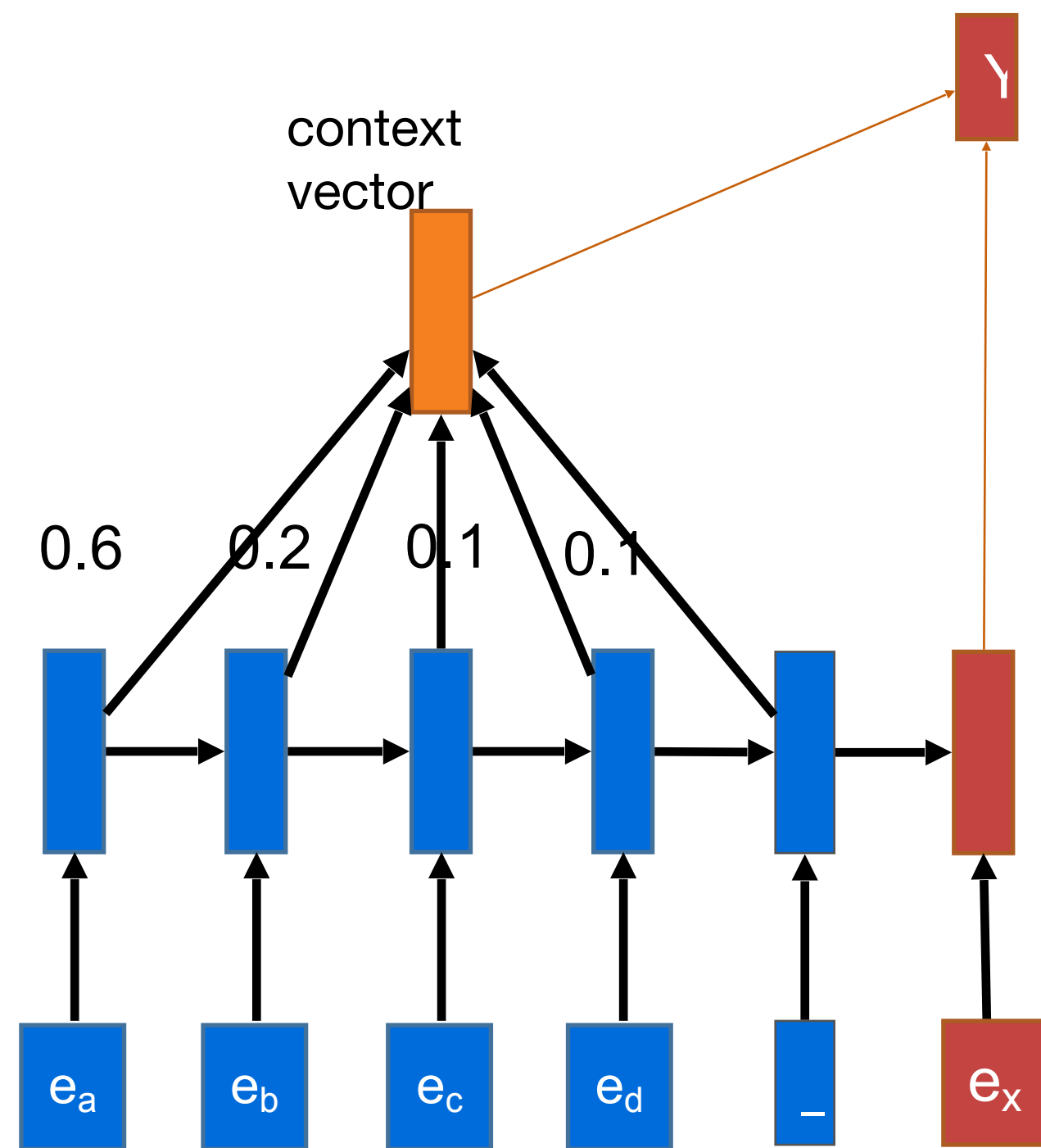| MODEL | TEST PERPLEXITY | NUMBER OF PARAMS [BILLIONS] |
|---|---|---|
| SIGMOID-RNN-2048 (JI ET AL., 2015A) | 68.3 | 4.1 |
| INTERPOLATED KN 5-GRAM, 1.1B N-GRAMS (CHELBA ET AL., 2013) | 67.6 | 1.76 |
| SPARSE NON-NEGATIVE MATRIX LM (SHAZEER ET AL., 2015) | 52.9 | 33 |
| RNN-1024 + MAXENT 9-GRAM FEATURES (CHELBA ET AL., 2013) | 51.3 | 20 |
| LSTM-512-512 | 54.1 | 0.82 |
| LSTM-1024-512 | 48.2 | 0.82 |
| LSTM-2048-512 | 43.7 | 0.83 |
| LSTM-8192-2048 (NO DROPOUT) | 37.9 | 3.3 |
| LSTM-8192-2048 (50% DROPOUT) | 32.2 | 3.3 |
| 2-LAYER LSTM-8192-1024 (BIG LSTM) | 30.6 | 1.8 |
| BIG LSTM+CNN INPUTS | **30.0** | **1.04** |
| BIG LSTM+CNN INPUTS + CNN SOFTMAX | 39.8 | **0.29** |
| BIG LSTM+CNN INPUTS + CNN SOFTMAX + 128-DIM CORRECTION | 35.8 | **0.39** |
| BIG LSTM+CNN INPUTS + CHAR LSTM PREDICTIONS | 47.9 | **0.23** |

Jozefowicz et al. Exploring the limits of language modelling, 2016

# Attention

A context vector c will be predicted before, which represents the related source context for current predicted word.
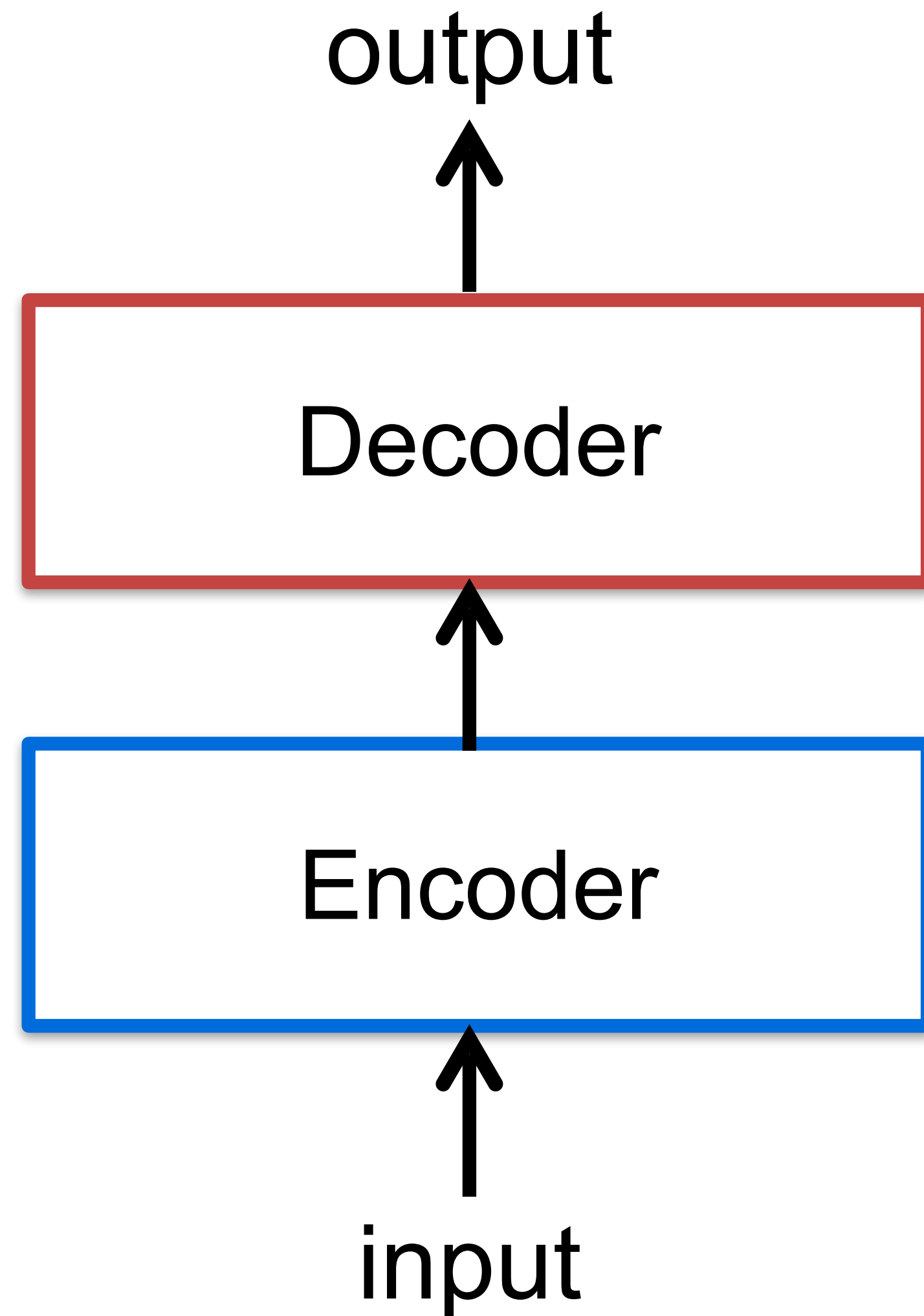
$$\alpha_{nj} = \text{Softmax}(D(s_n, h_{1\dots n-1})) = \frac{\exp(D(s_n, h_j))}{\sum_k \exp(D(s_n, h_k))}$$

$$c_n = \sum_j \alpha_{nj} h_j$$

The probability of word y_i is computed as:

$$p(y_i) \propto \exp(Wh_i) \quad \Rightarrow \quad p(y_i) \propto \exp(Wh_i + Vc_i)$$

Mnih et al. Recurrent Models of Visual Attention. 2014.

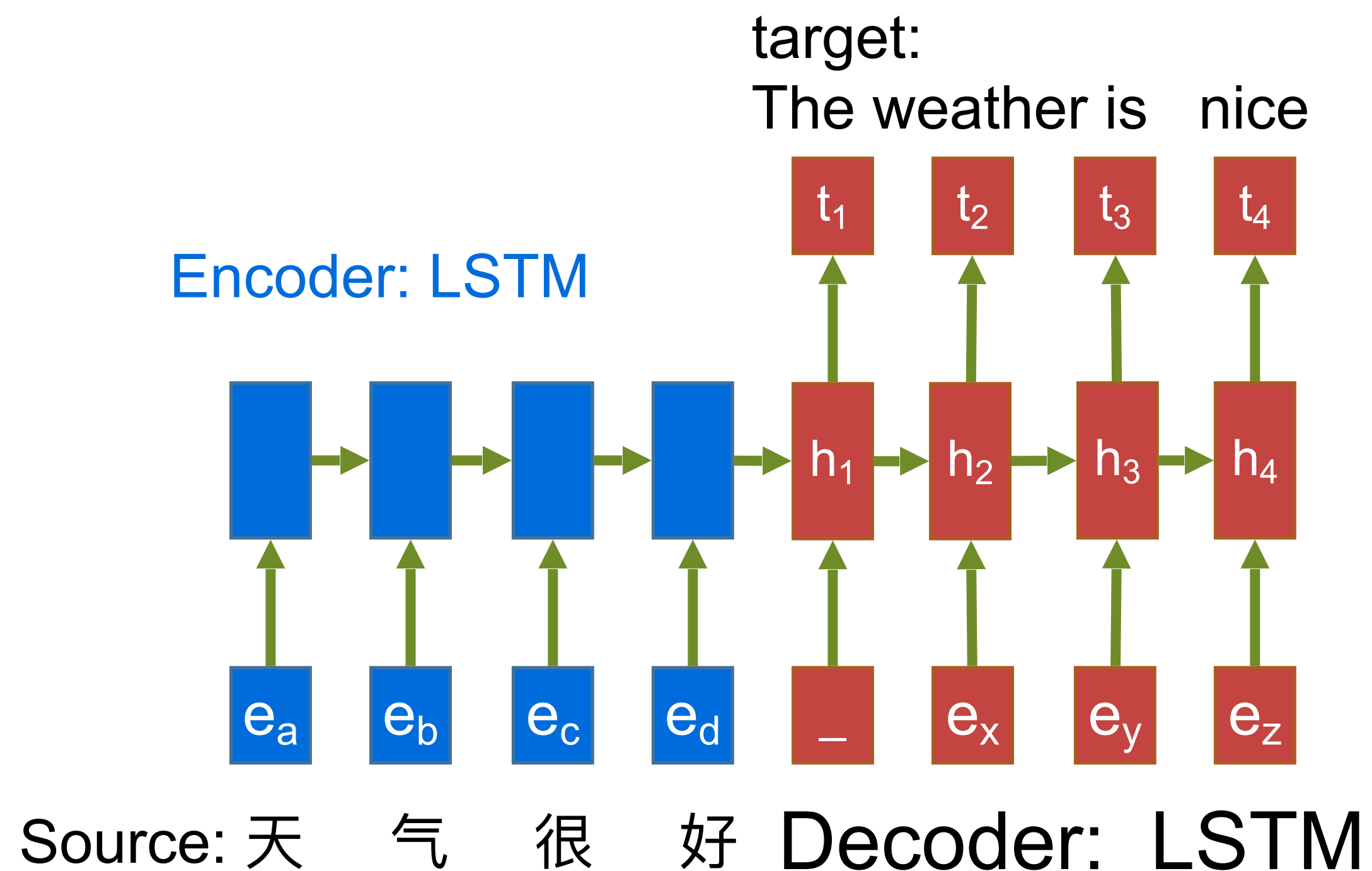# Encoder-decoder framework

output

Decoder

Encoder

input

A generic formulation

ImageCaption

Text-to-Image Generation

ASR (speech-to-text)

MT (text-to-text)

# Sequence To Sequence (Seq2seq)

- Machine translation as directly learning a function mapping from source sequence to target sequence
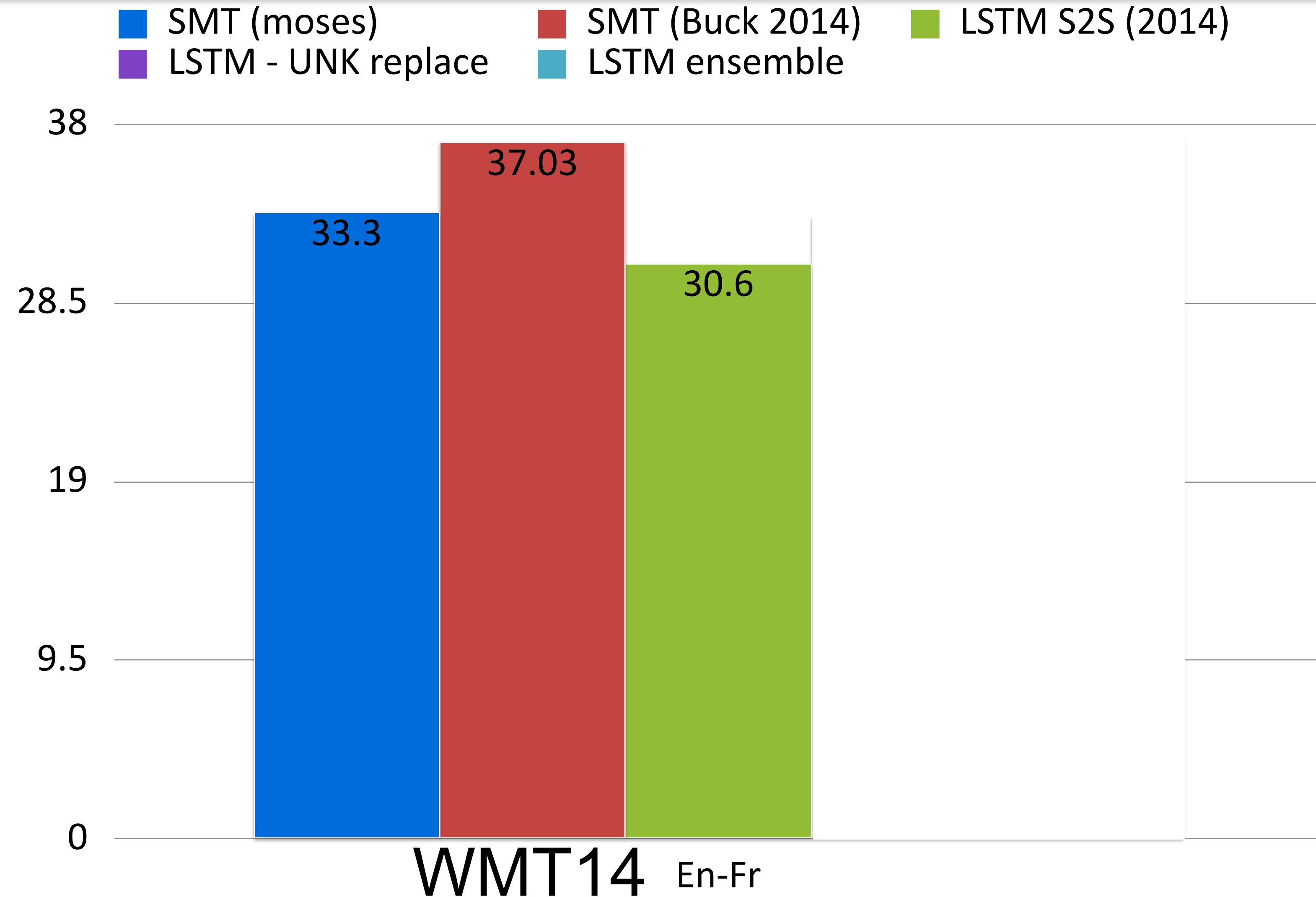


$$P(Y|X) = \prod P(y_t \mid y_{<t}, x)$$

Training loss: Cross-Entropy

$$l = -\sum_n \sum_t \log f_\theta(x_n, y_{n,1}, \ldots, y_{n,t-1})$$

Teacher-forcing during training.
(pretend to know groundtruth for prefix)

Sutskever et al. Sequence to Sequence Learning with Neural Networks. 2014

# Performance (2014)



**Legend:** SMT (moses) | SMT (Buck 2014) | LSTM S2S (2014) | LSTM - UNK replace | LSTM ensemble

- SMT (moses): 33.3
- SMT (Buck 2014): 37.03
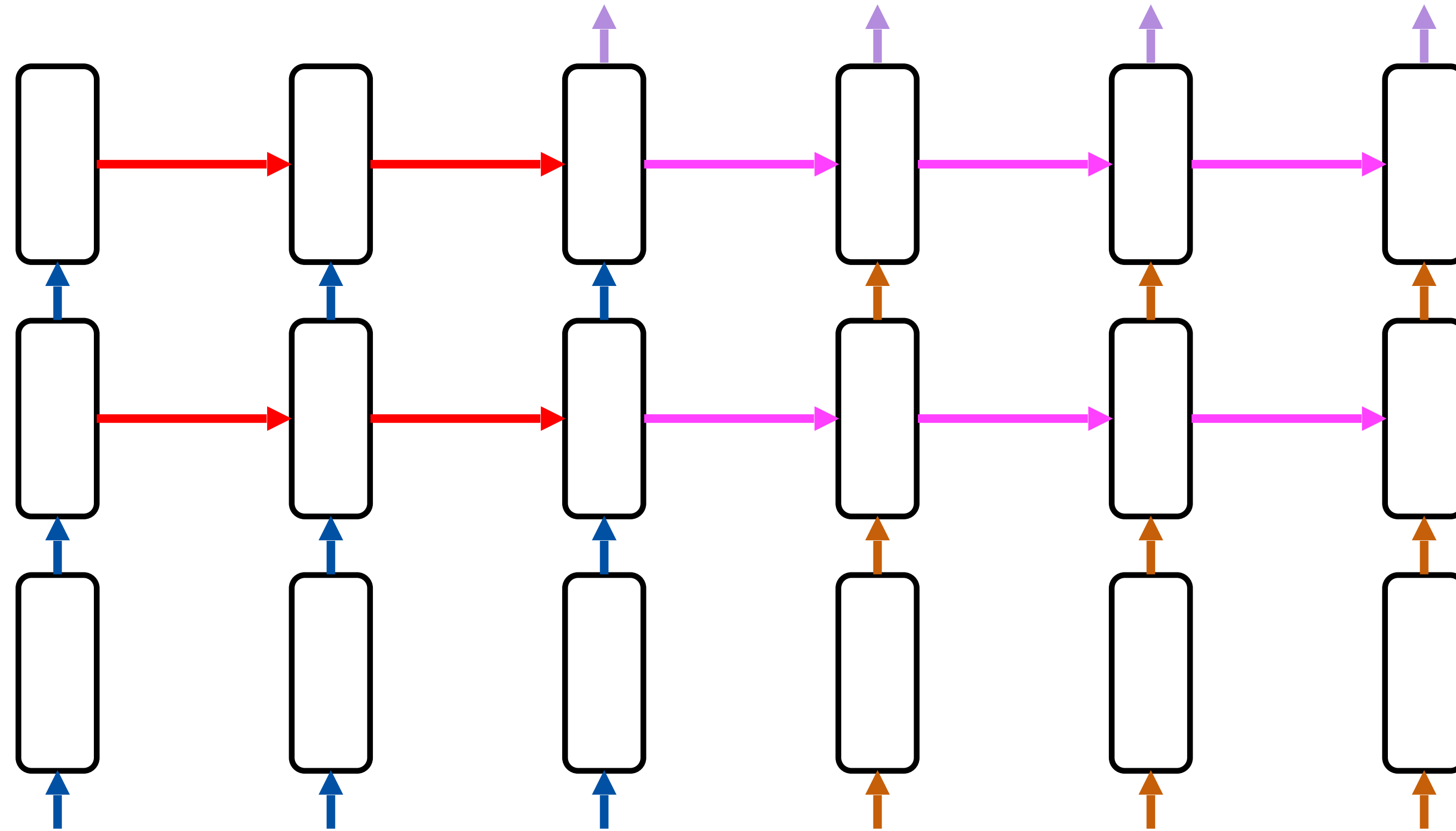- LSTM S2S (2014): 30.6

**WMT14** En-Fr

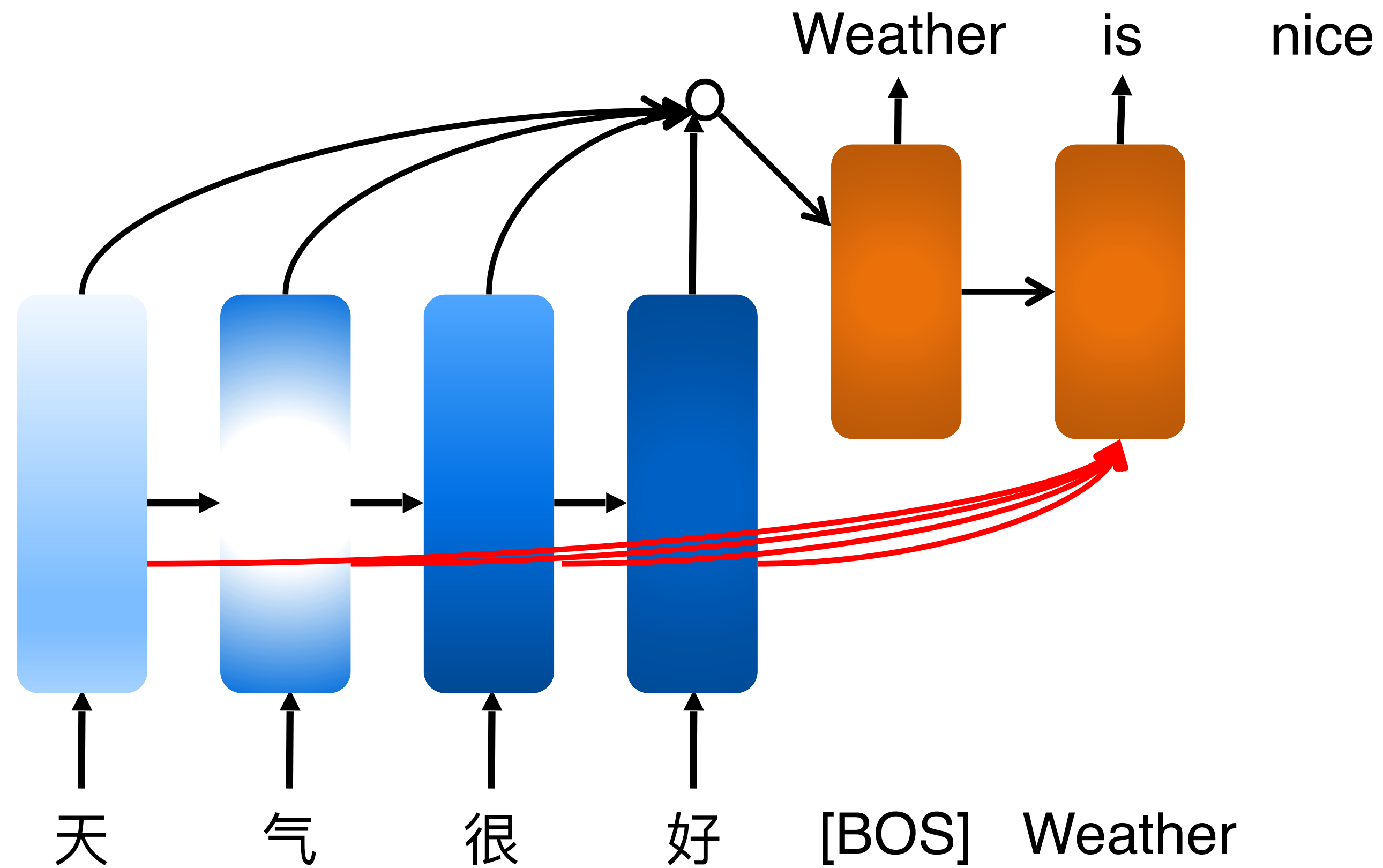Sutskever et al. Sequence to Sequence Learning with Neural Networks. 2014
Durrani et al. Edinburgh's Phrase-based Machine Translation Systems for WMT-14. 2014
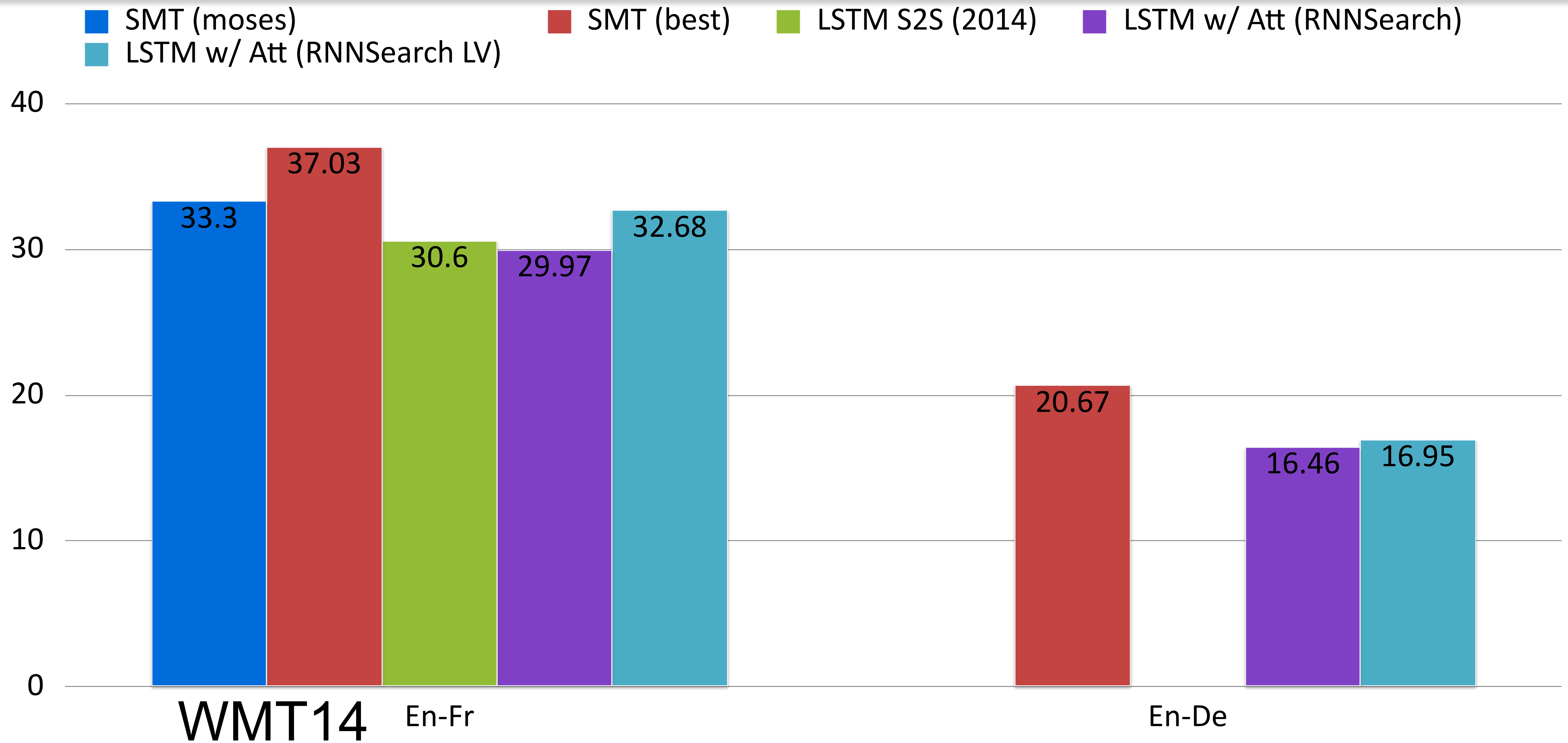
# Stacked LSTM for seq-2-seq
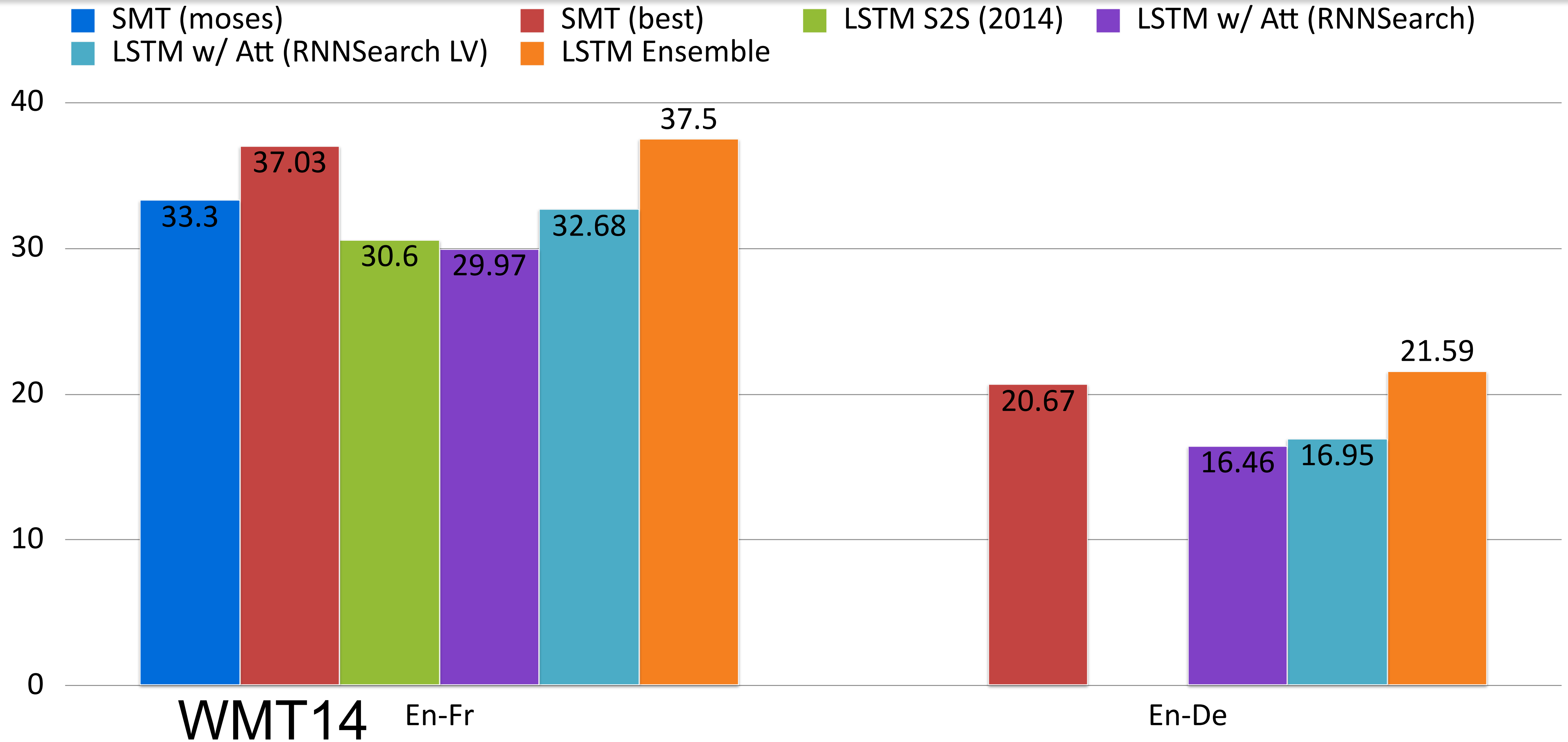
- More layers of LSTM

# LSTM Seq2seq with Attention

# LSTM Seq2Seq w/ Attention



Legend: SMT (moses), SMT (best), LSTM S2S (2014), LSTM w/ Att (RNNSearch), LSTM w/ Att (RNNSearch LV)

WMT14 En-Fr: SMT (moses) 33.3, SMT (best) 37.03, LSTM S2S (2014) 30.6, LSTM w/ Att (RNNSearch) 29.97, LSTM w/ Att (RNNSearch LV) 32.68

En-De: SMT (best) 20.67, LSTM w/ Att (RNNSearch) 16.46, LSTM w/ Att (RNNSearch LV) 16.95

Jean et al. On Using Very Large Target Vocabulary for Neural Machine Translation. 2015

# Performance with Model Ensemble



Luong et al. Effective Approaches to Attention-based Neural Machine Translation. 2015

# Reading

- Gers et al. Learning to Forget: Continual Prediction with LSTM. 2000

- Sutskever et al. Sequence to Sequence Learning with Neural Networks. 2014

- Bahdanau et al., Neural Machine Translation by Jointly Learning to Align and Translate. 2015

- Luong et al. Effective Approaches to Attention-based Neural Machine Translation. 2015