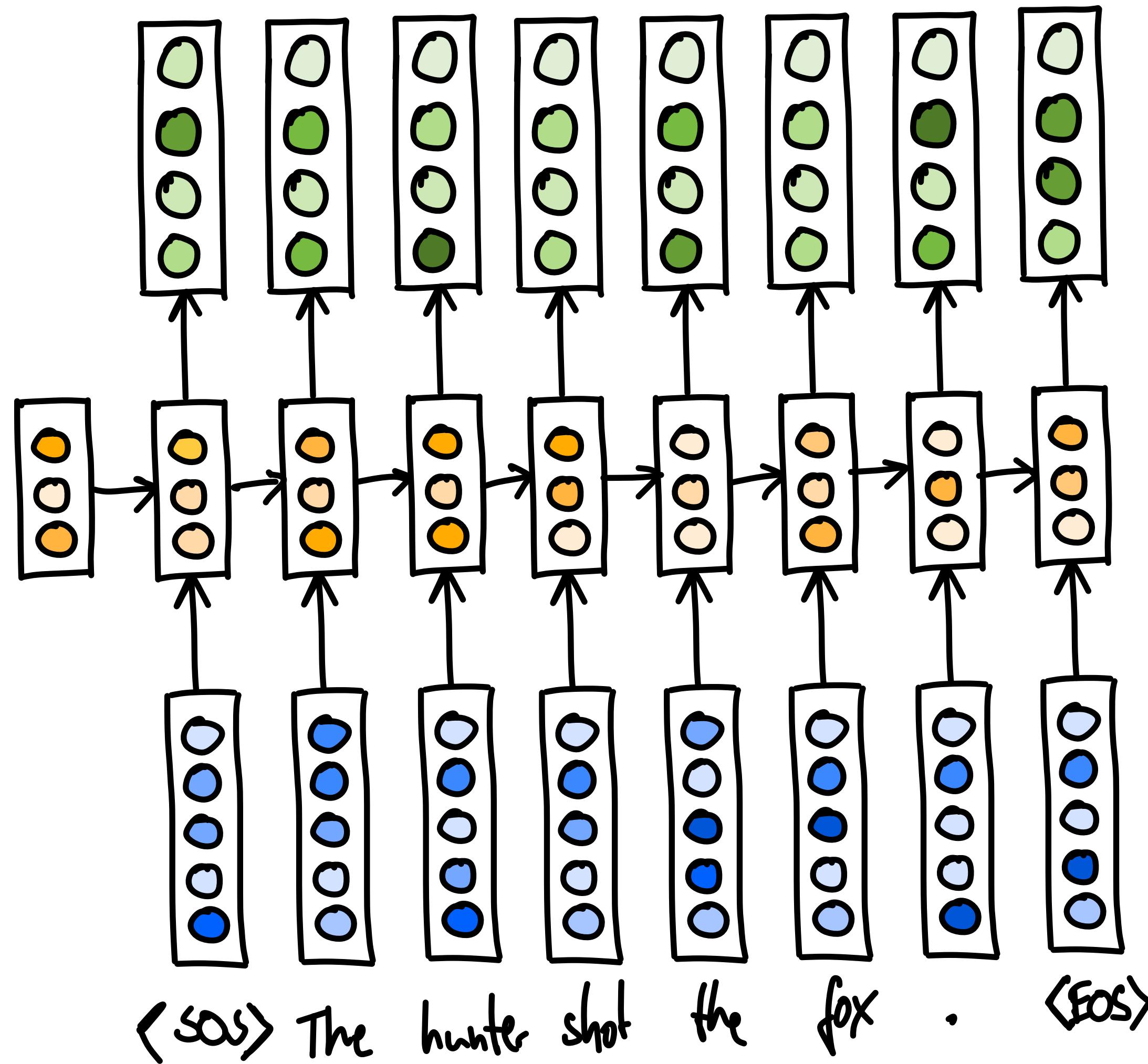


Advanced RNNs

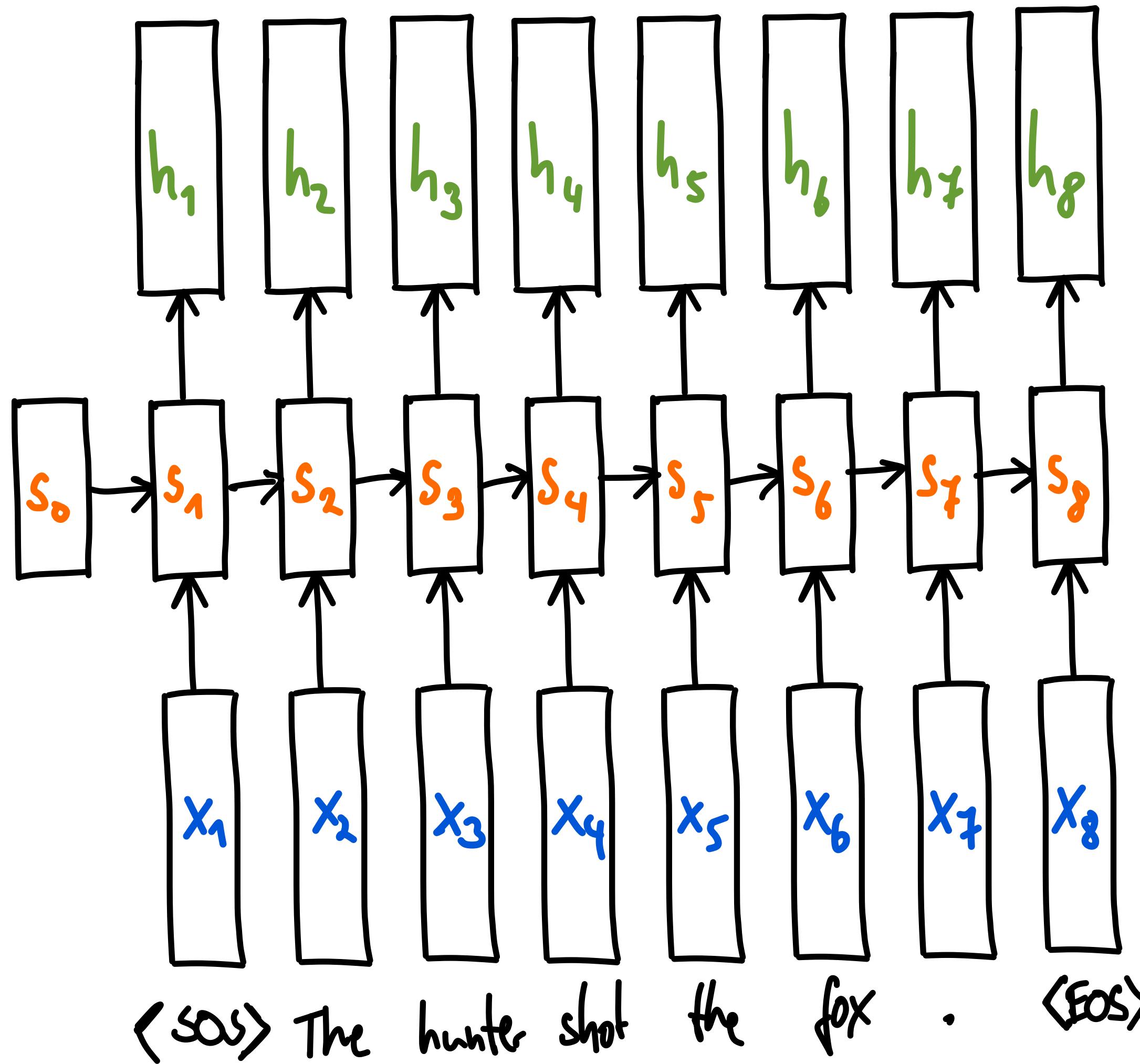
Tim Rocktäschel & Sebastian Riedel
COMP0087 Natural Language Processing



Backpropagation Through Time (BPTT)



Backpropagation Through Time (BPTT)



$$s_t = \tanh(\mathbf{W}[\mathbf{x}_t; \mathbf{s}_{t-1}] + \mathbf{b})$$

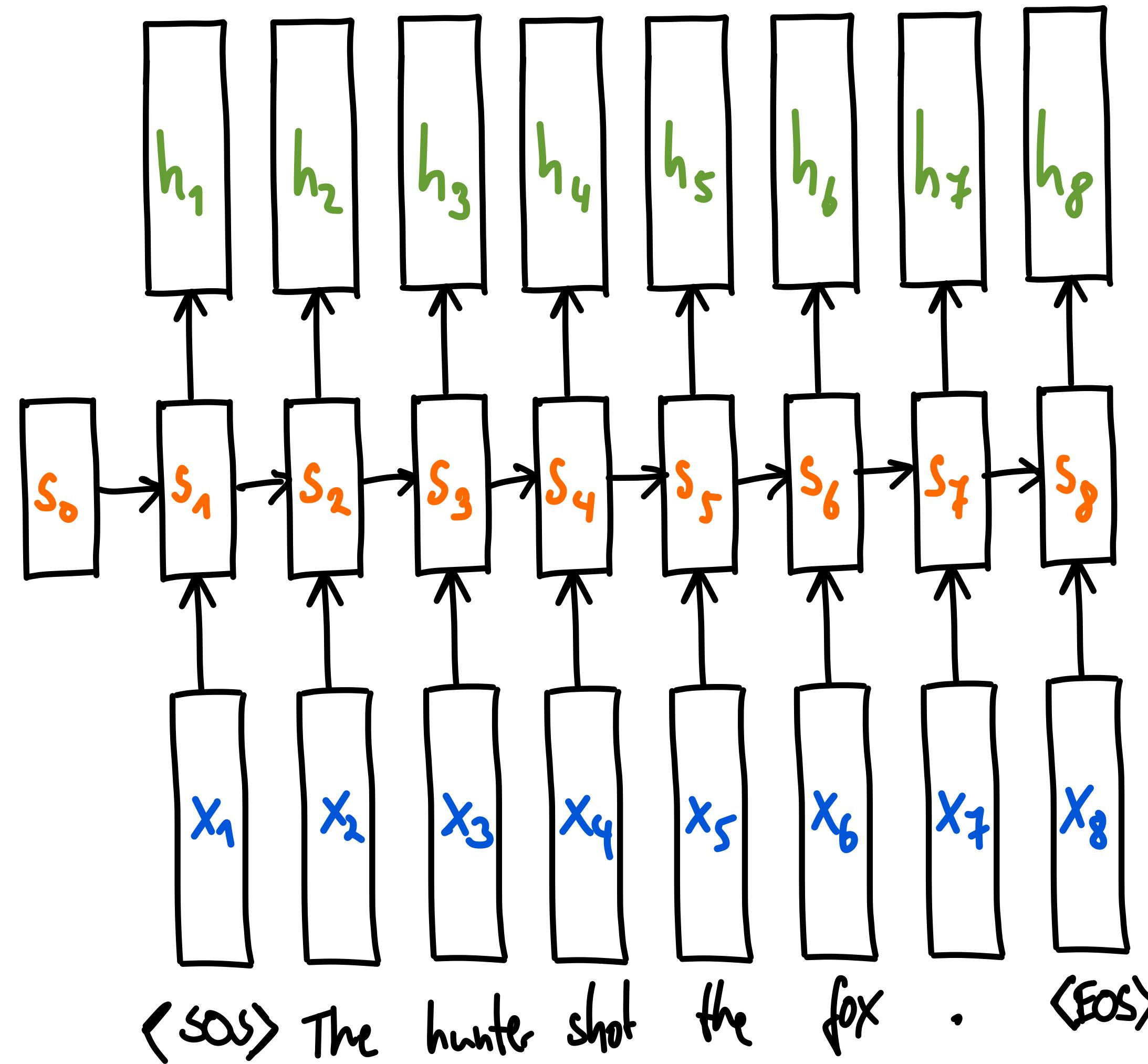
$$\mathbf{h}_t = \mathbf{s}_t$$

$$\frac{\partial L(w_N, \hat{w}_N)}{\partial s_t} = \frac{\partial L(w_N, \hat{w}_N)}{\partial s_N} \frac{\partial s_N}{\partial s_{N-1}} \dots \frac{\partial s_{t+1}}{\partial s_t}$$

$$= \frac{\partial L(w_N, \hat{w}_N)}{\partial s_N} \prod_{i=t+1}^N \frac{\partial s_i}{\partial s_{i-1}}$$

$$\frac{\partial L(w, \hat{w})}{\partial s_t} = \sum_{j=t}^N \frac{\partial L(w_j, \hat{w}_j)}{\partial s_j} \prod_{i=t+1}^j \frac{\partial s_i}{\partial s_{i-1}}$$

Vanishing and Exploding Gradients



$$\mathbf{u}_t = \mathbf{W}^x \mathbf{x}_t + \mathbf{W}^s \mathbf{s}_{t-1} + \mathbf{b}$$

$$\mathbf{s}_t = \tanh(\mathbf{u}_t)$$

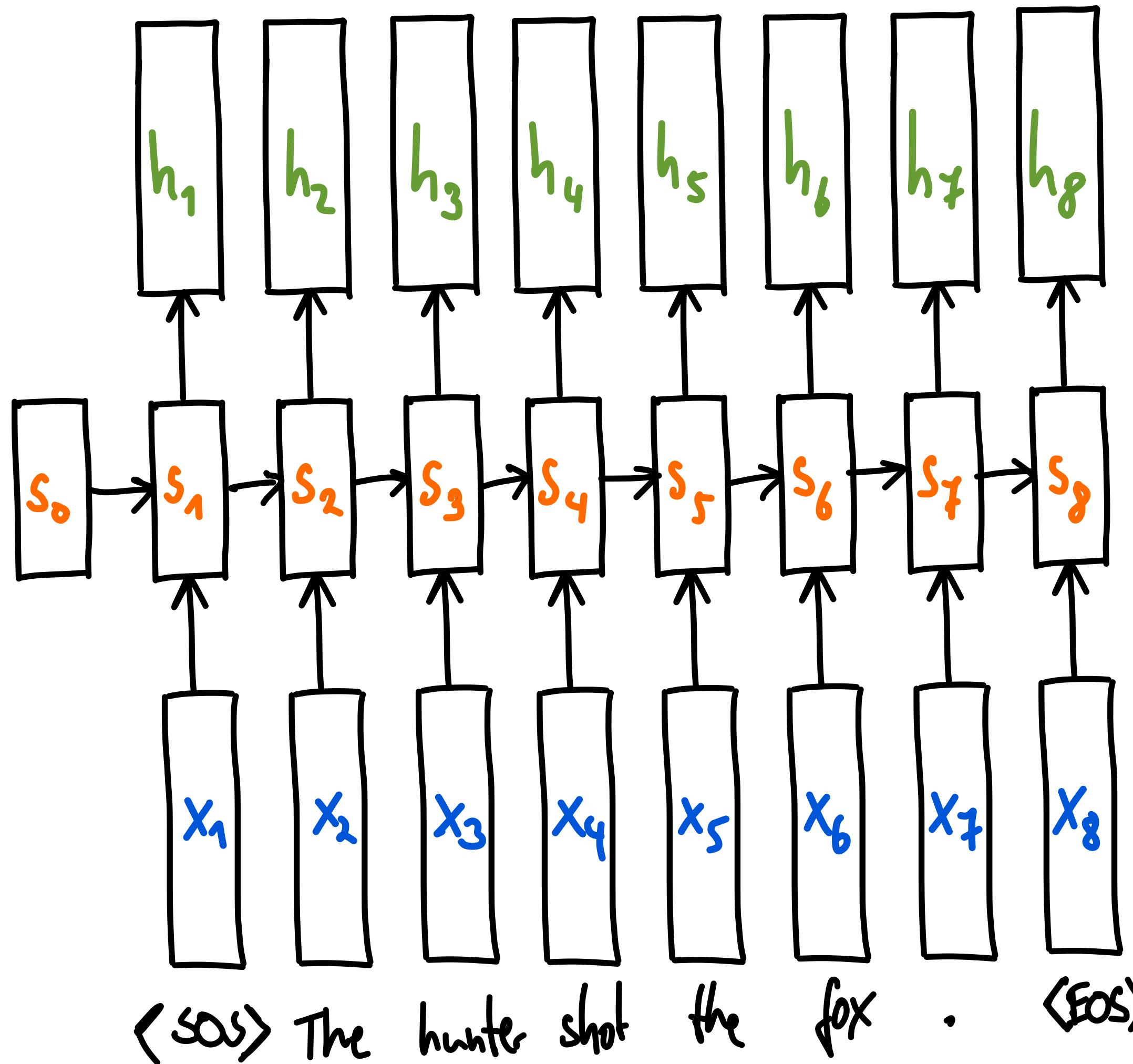
$$\mathbf{h}_t = \mathbf{s}_t$$

$$\frac{\partial L(w_N, \hat{w}_N)}{\partial \mathbf{s}_t} = \frac{\partial L(w_N, \hat{w}_N)}{\partial \mathbf{s}_N} \prod_{i=t+1}^N \frac{\partial \mathbf{s}_i}{\partial \mathbf{u}_i} \frac{\partial \mathbf{u}_i}{\partial \mathbf{s}_{i-1}}$$

$$\frac{\partial \mathbf{s}_i}{\partial \mathbf{u}_i} = \tanh'(\mathbf{u}_{i-1})^T, \quad \frac{\partial \mathbf{u}_i}{\partial \mathbf{s}_{i-1}} = \mathbf{W}^s$$

$$\frac{\partial L(w_N, \hat{w}_N)}{\partial \mathbf{s}_t} = \frac{\partial L(w_N, \hat{w}_N)}{\partial \mathbf{s}_N} \prod_{i=t+1}^N \tanh'(\mathbf{u}_{i-1})^T \mathbf{W}^s$$

Vanishing and Exploding Gradients



$$\mathbf{u}_t = \mathbf{W}^x \mathbf{x}_t + \mathbf{W}^s \mathbf{s}_{t-1} + \mathbf{b}$$

$$\mathbf{s}_t = \tanh(\mathbf{u}_t)$$

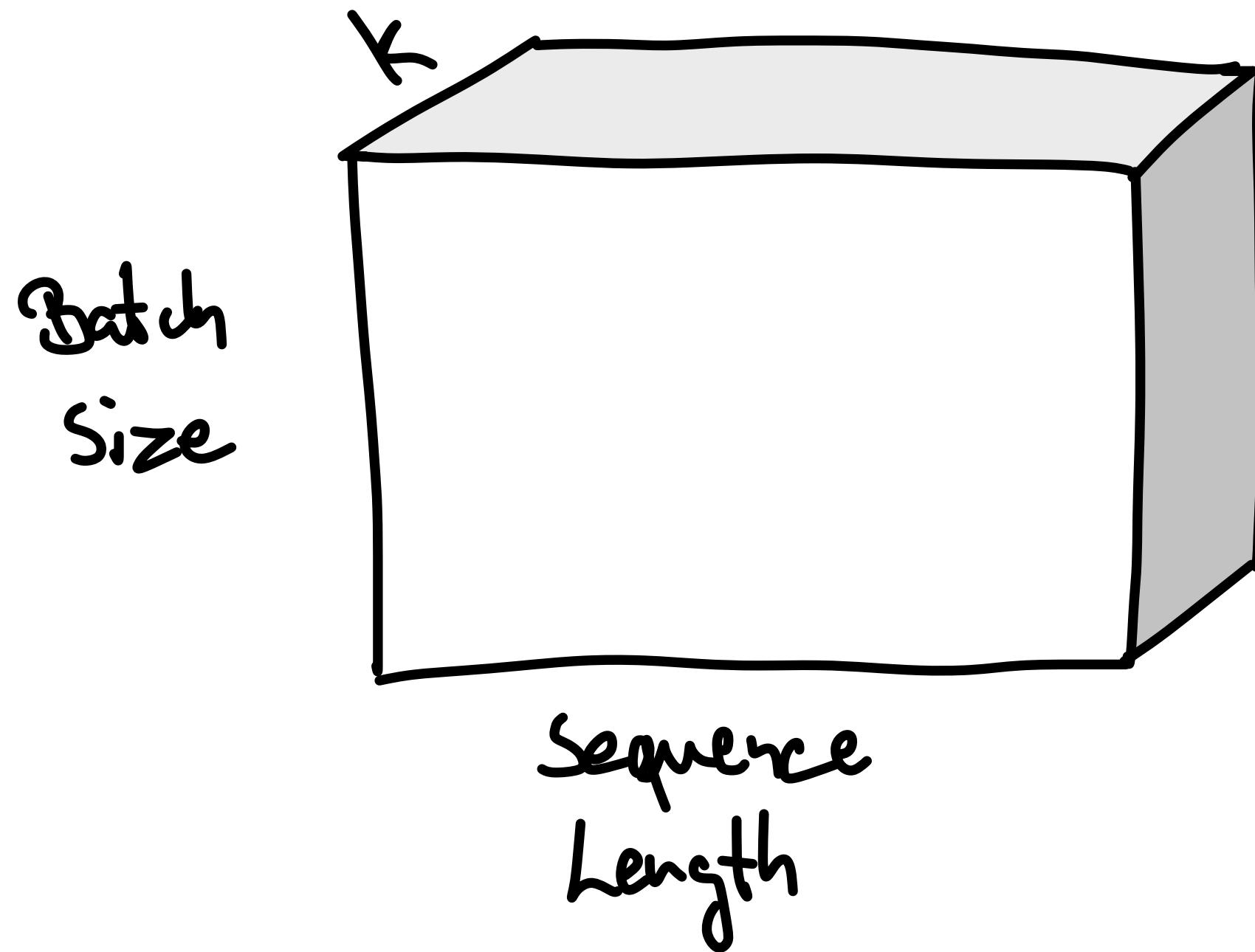
$$\mathbf{h}_t = \mathbf{s}_t$$

$$\frac{\partial L(w_N, \hat{w}_N)}{\partial s_t} = \frac{\partial L(w_N, \hat{w}_N)}{\partial s_N} \prod_{i=t+1}^N \tanh'(\mathbf{u}_{i-1})^T \mathbf{W}^s$$

Note the repeated multiplication of \mathbf{W}^s
 Let λ be the largest eigenvalue of \mathbf{W}^s

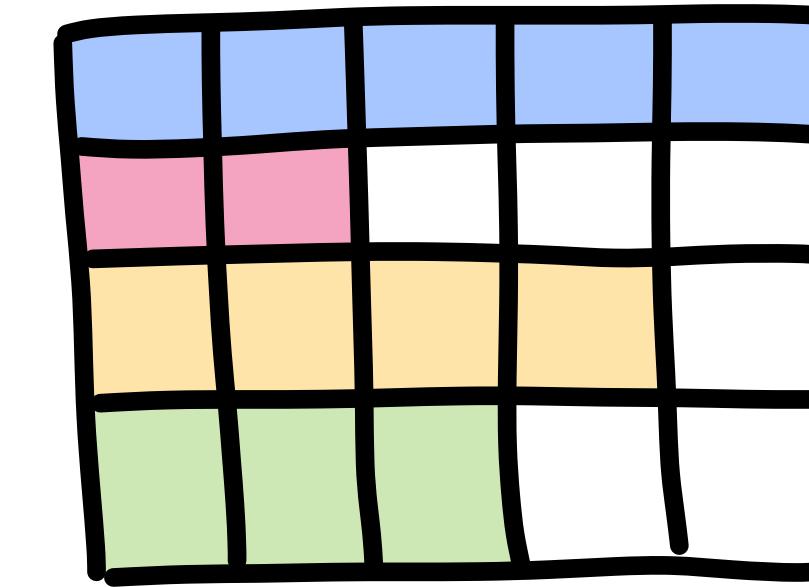
- Gradient explodes if $\lambda > 1$
- Gradient vanishes if $\lambda < 1$
- Gradient is stable if $\lambda = 1$

Batch Processing



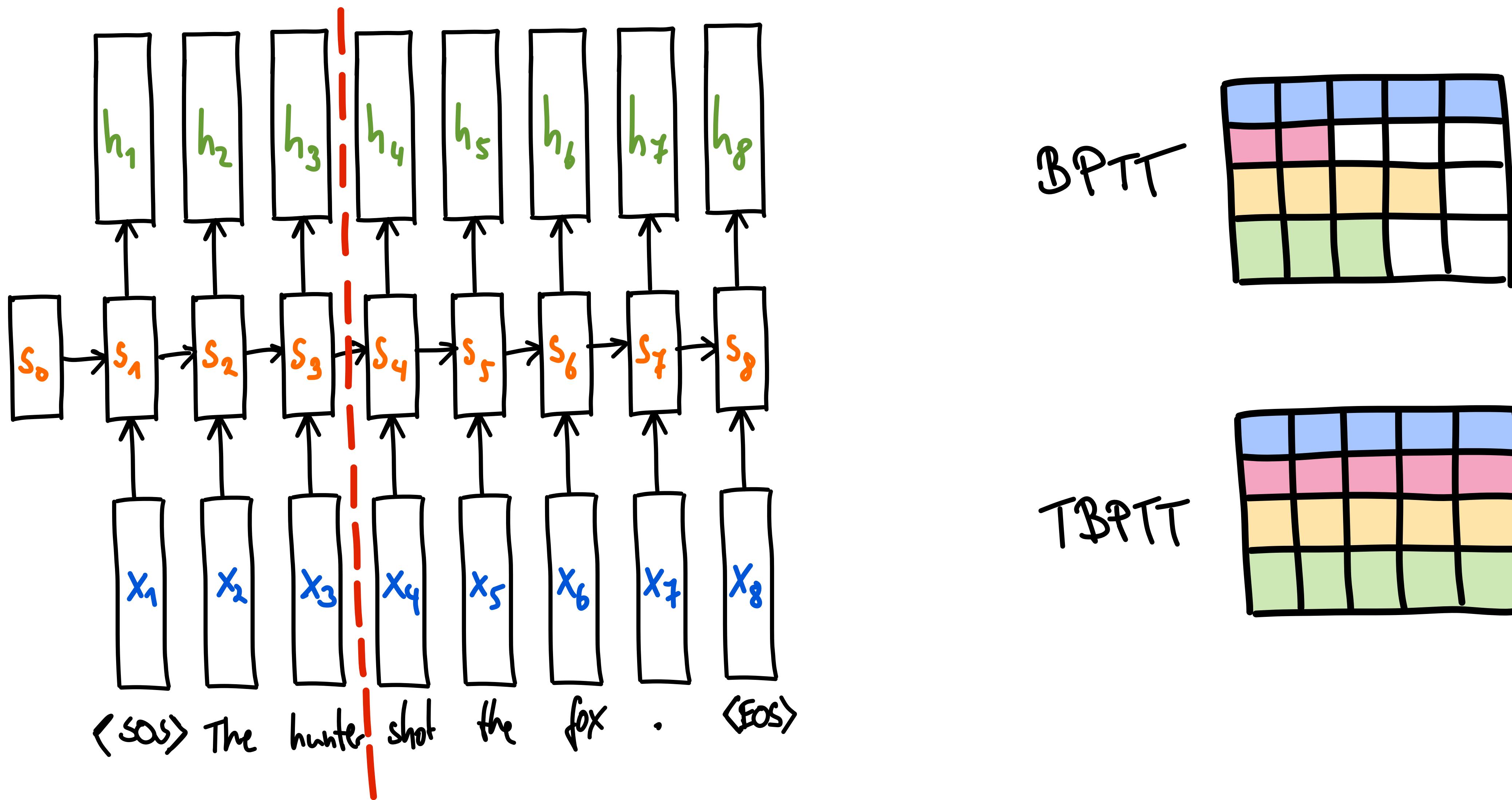
$$\mathbf{S}_t = \tanh(\mathbf{X}_t \mathbf{W}_x^T + \mathbf{S}_{t-1} \mathbf{W}_s^T + \mathbf{b})$$

BPTT

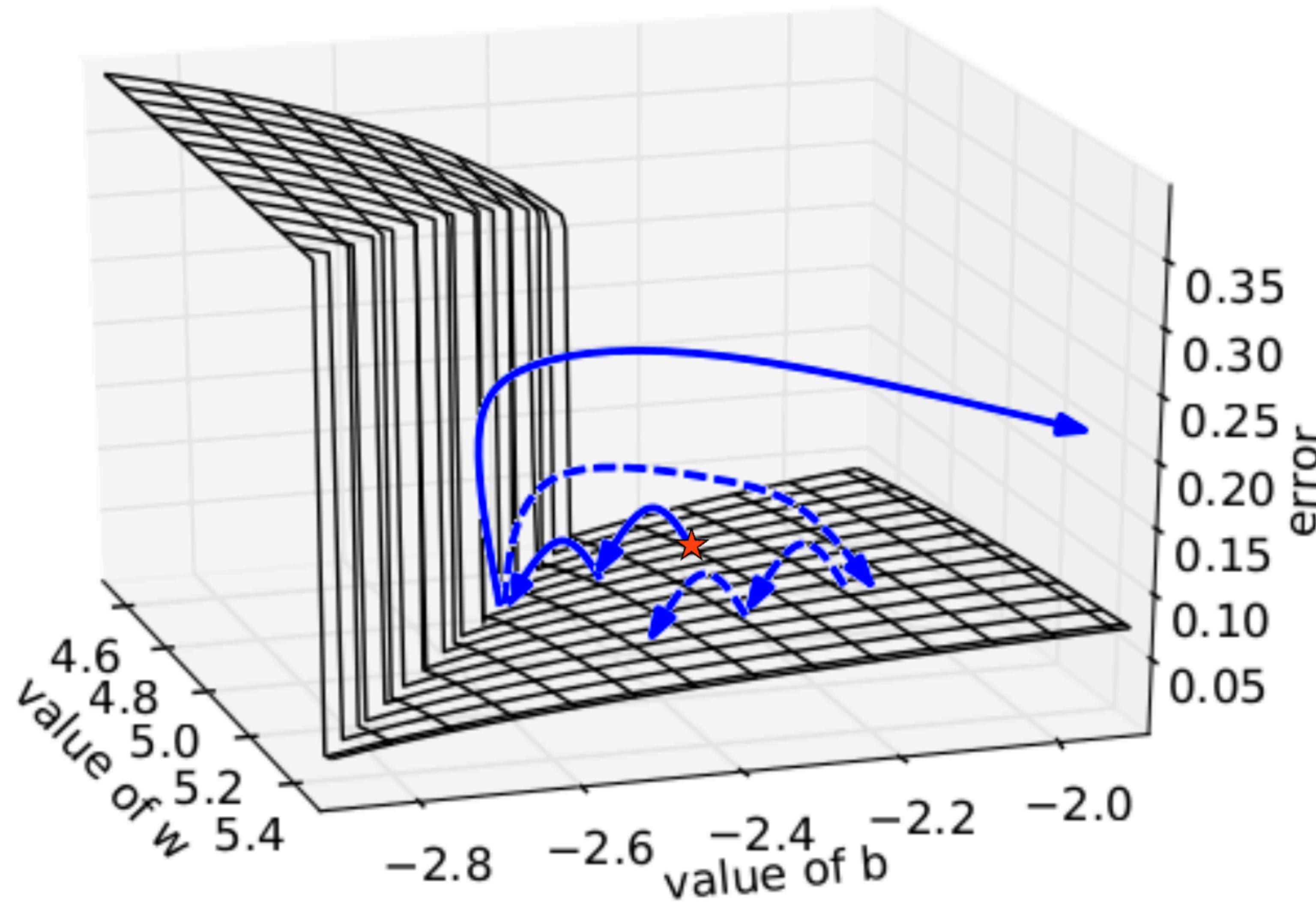


```
def f(self, X, S):
    X = torch.cat([X, S]) # -- [batch_size x 2*k]
    S = torch.tanh(torch.einsum("ij,kj->ki", [self.W, X]) + self.b) # -- [batch_size x k]
    H = S
    return H, S
```

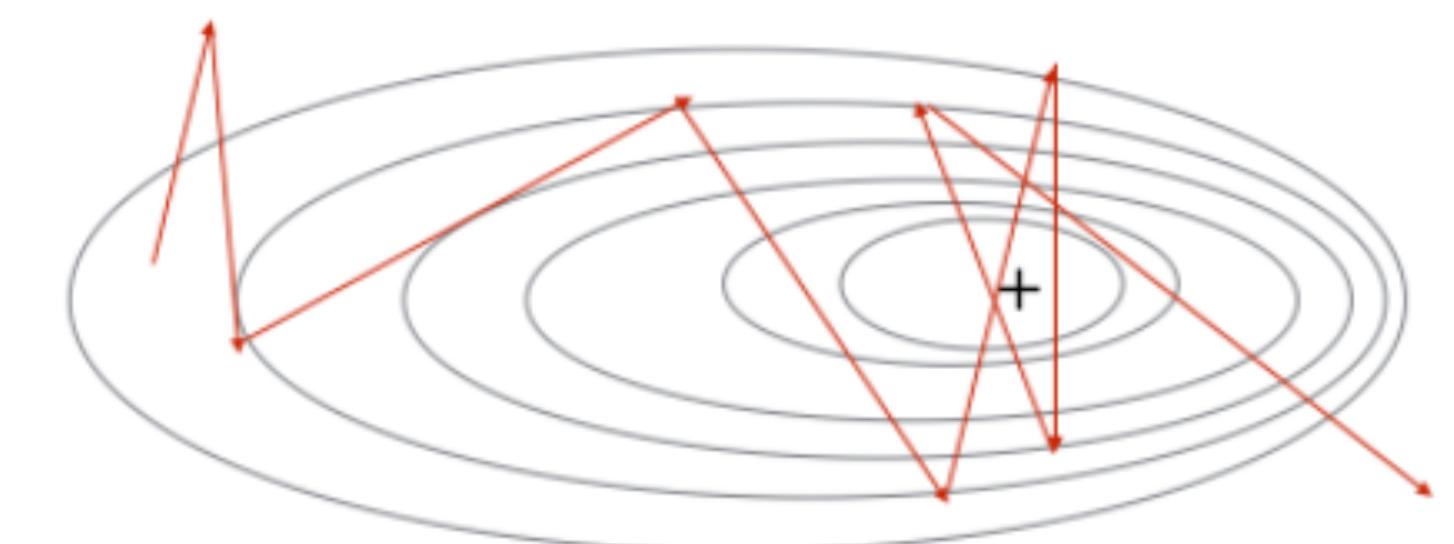
Truncated BPTT



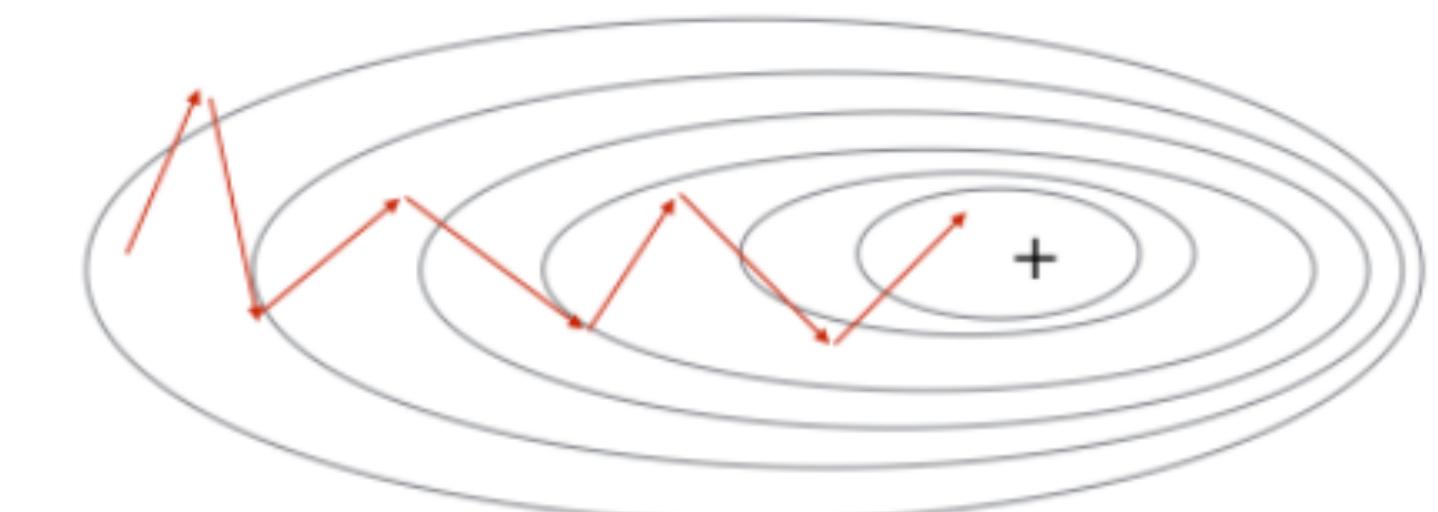
Gradient Clipping



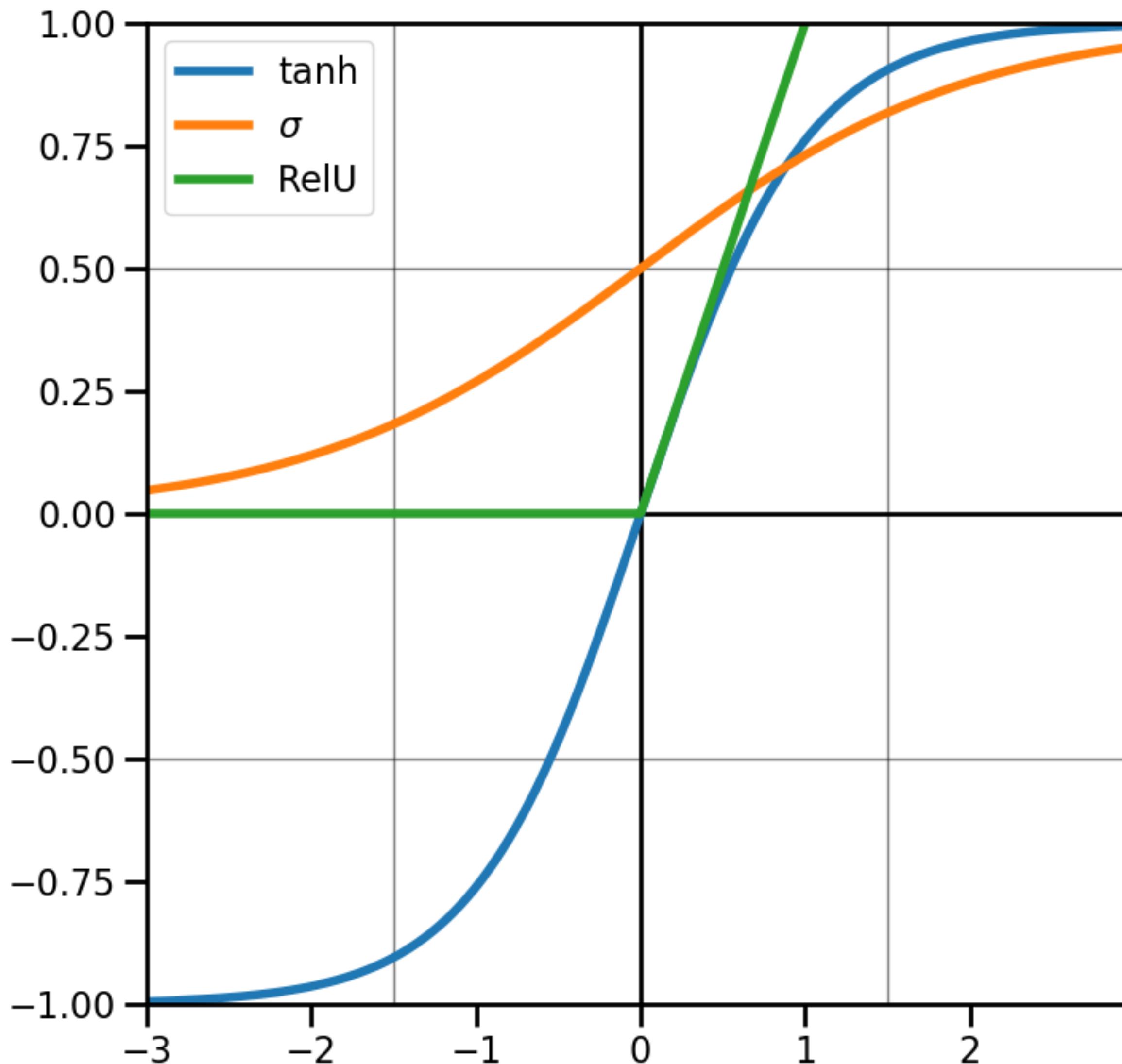
Without gradient clipping



With gradient clipping

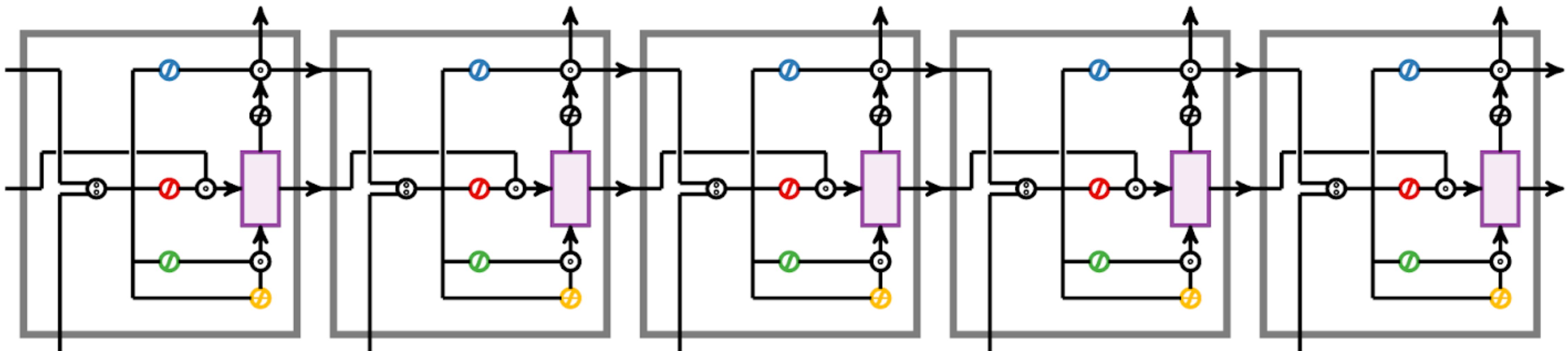


Activation Functions

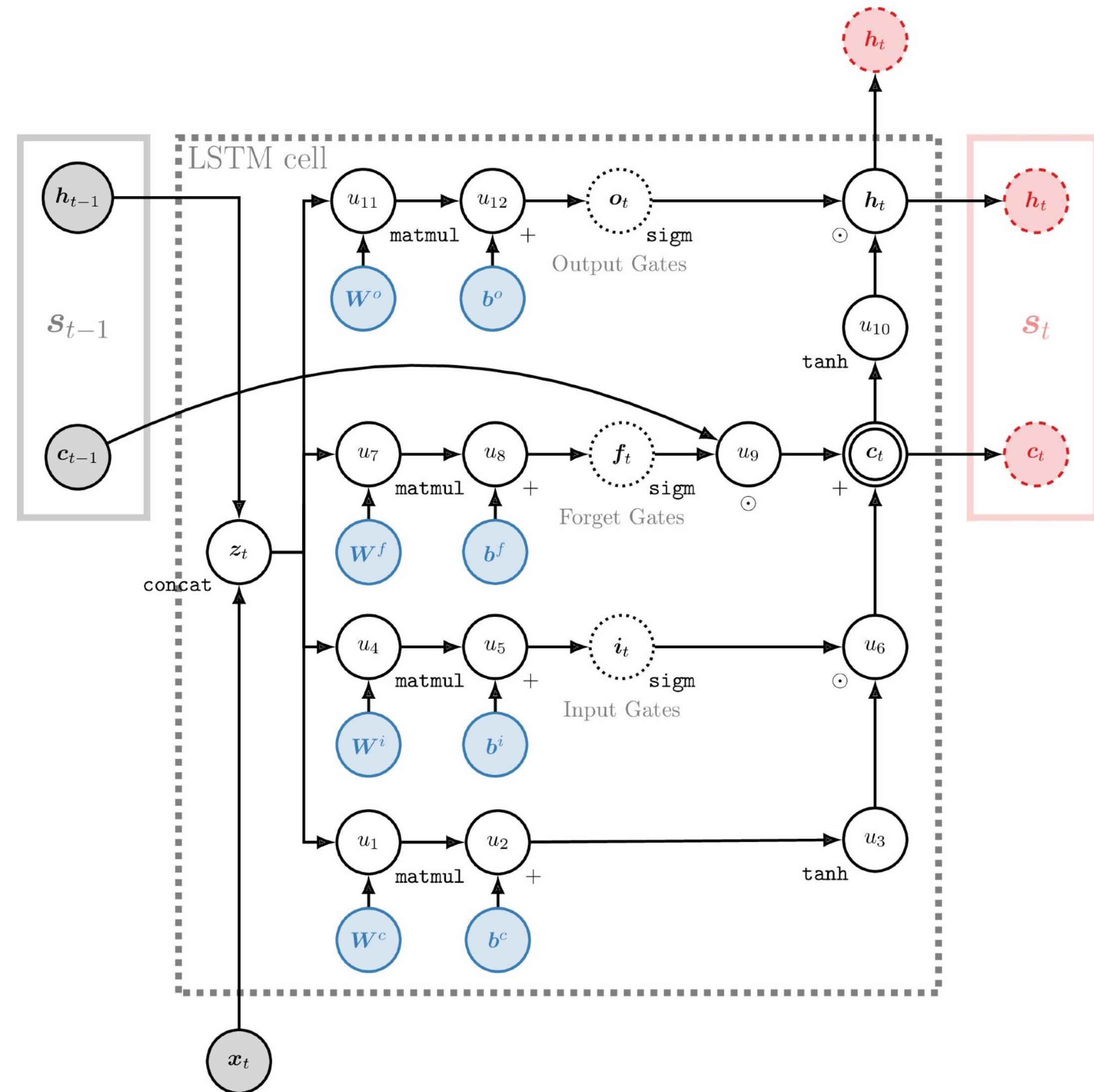


Long Short-term Memory

- Dedicated **memory cell**
- **Input gate**, **forget gate** and **output gate** control flow **into** and out of memory



Long Short-term Memory



$$\mathbf{h}_{t-1}, \mathbf{c}_{t-1} = \mathbf{s}_{t-1}$$

$$\mathbf{z}_t = [\mathbf{x}_t; \mathbf{h}_{t-1}]$$

$$\mathbf{i}_t = \sigma(\mathbf{W}^i \mathbf{z}_t + \mathbf{b}^i)$$

$$\mathbf{f}_t = \sigma(\mathbf{W}^f \mathbf{z}_t + \mathbf{b}^f)$$

$$\mathbf{o}_t = \sigma(\mathbf{W}^o \mathbf{z}_t + \mathbf{b}^o)$$

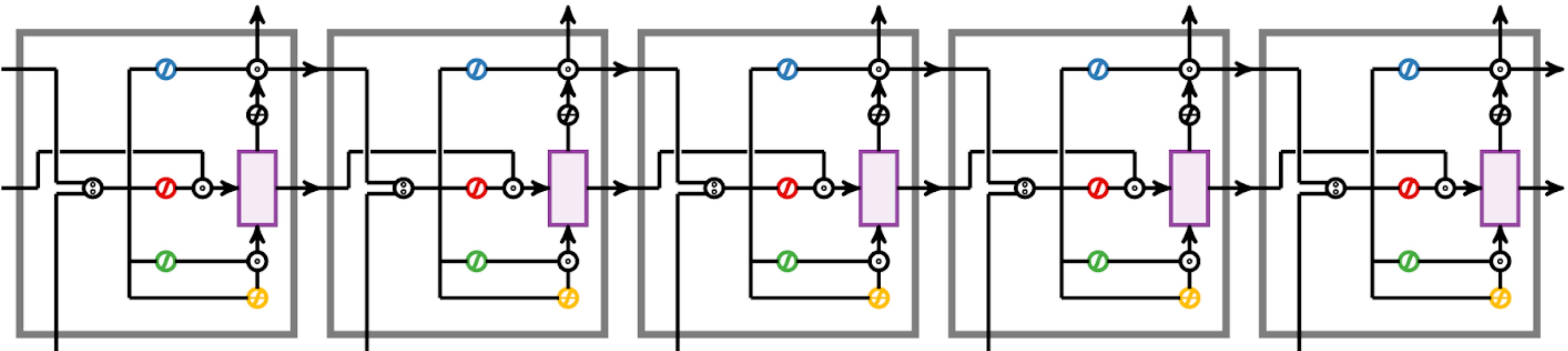
$$\mathbf{c}_t = \mathbf{f}_t \odot \mathbf{c}_{t-1} + \mathbf{i}_t \odot \tanh(\mathbf{W}^c \mathbf{z}_t + \mathbf{b}^c)$$

$$\mathbf{h}_t = \mathbf{o}_t \odot \tanh(\mathbf{c}_t)$$

$$\mathbf{s}_t = [\mathbf{h}_t; \mathbf{c}_t]$$

Long Short-term Memory

- Dedicated **memory cell**
- **Input gate**, **forget gate** and **output gate** control flow **into** and out of memory

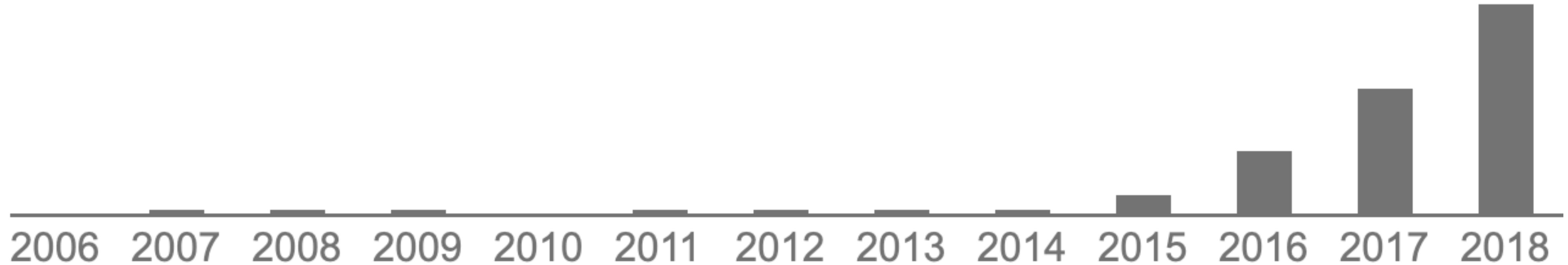


$$\mathbf{c}_t = \mathbf{f}_t \odot \mathbf{c}_{t-1} + \mathbf{i}_t \odot \tanh(\mathbf{W}^c \mathbf{z}_t + \mathbf{b}^c)$$

$$\frac{\partial \mathbf{c}_t}{\partial \mathbf{c}_{t-1}} = \frac{\partial}{\partial \mathbf{c}_{t-1}} [\mathbf{f}_t \odot \mathbf{c}_{t-1}] + \frac{\partial}{\partial \mathbf{c}_{t-1}} [\mathbf{i}_t \odot \tanh(\mathbf{W}^c \mathbf{z}_t + \mathbf{b}^c)]$$

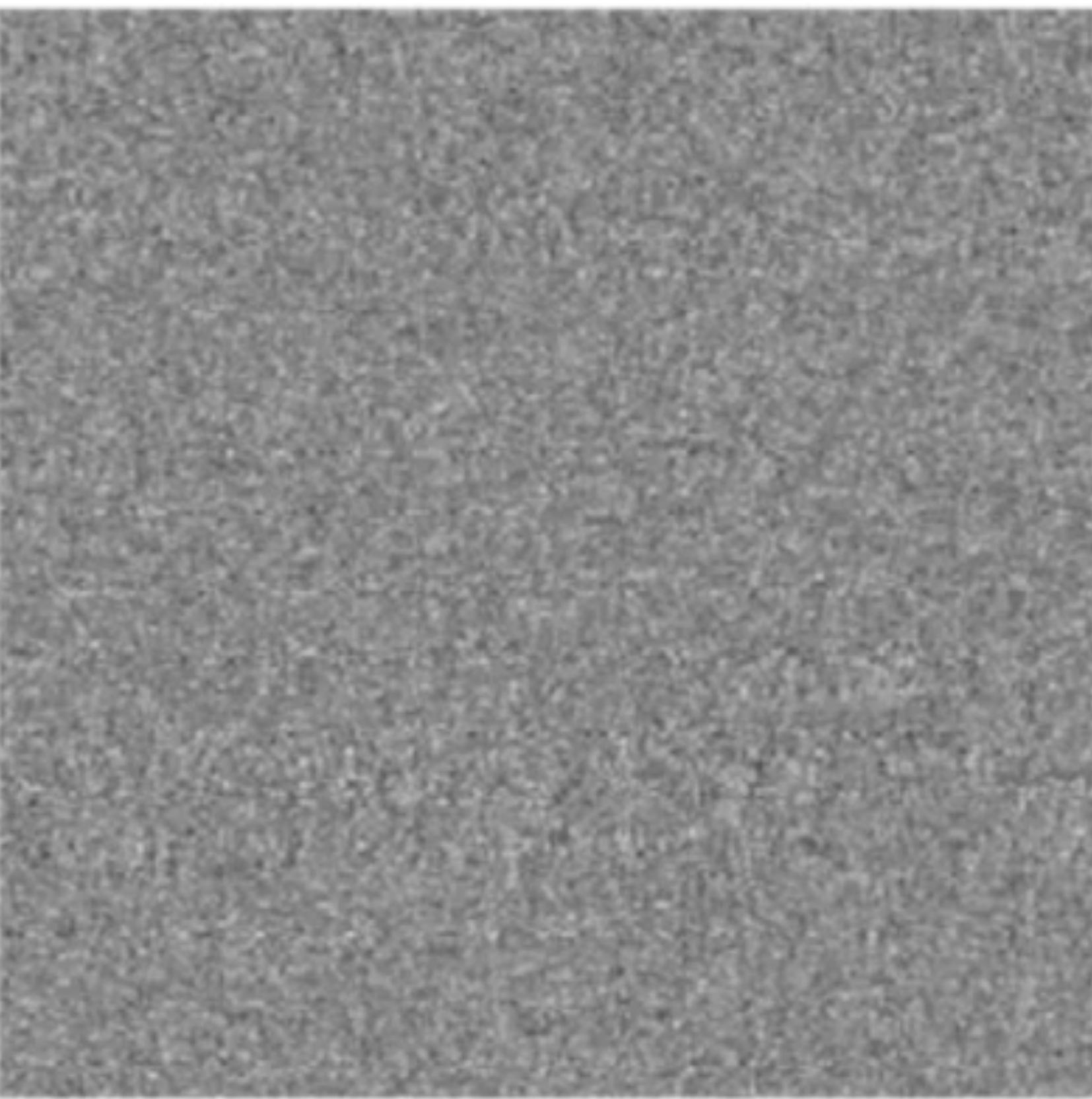
Long Short-term Memory

Cited by 15513

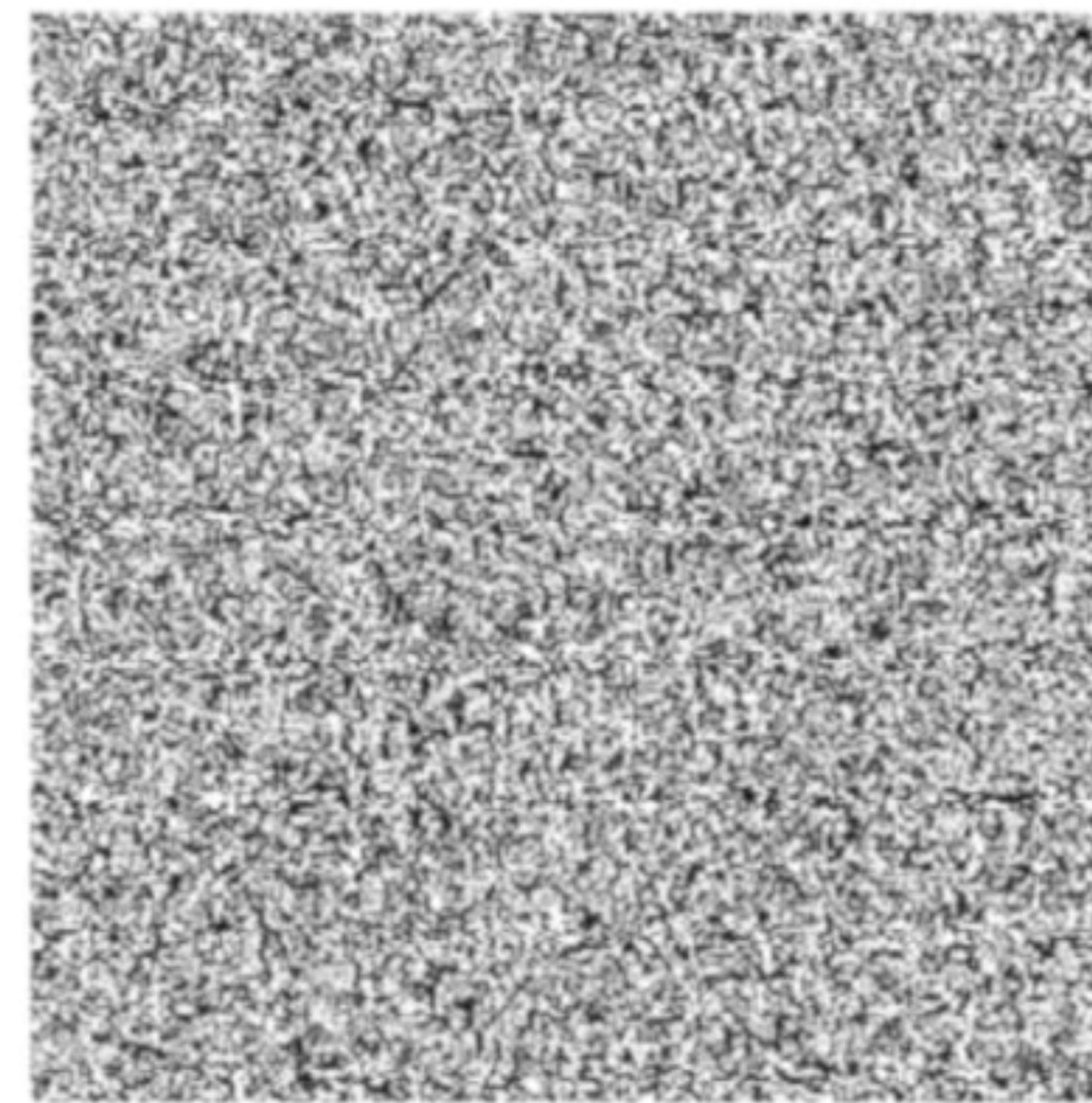


Mitigates Vanishing and Exploding Gradients

127



127



LSTM Efficiency

$$\mathbf{h}_{t-1}, \mathbf{c}_{t-1} = \mathbf{s}_{t-1}$$

$$\mathbf{z}_t = [\mathbf{x}_t; \mathbf{h}_{t-1}]$$

$$\mathbf{i}_t = \sigma(\mathbf{W}^i \mathbf{z}_t + \mathbf{b}^i)$$

$$\mathbf{f}_t = \sigma(\mathbf{W}^f \mathbf{z}_t + \mathbf{b}^f)$$

$$\mathbf{o}_t = \sigma(\mathbf{W}^o \mathbf{z}_t + \mathbf{b}^o)$$

$$\mathbf{c}_t = \mathbf{f}_t \odot \mathbf{c}_{t-1} + \mathbf{i}_t \odot \tanh(\mathbf{W}^c \mathbf{z}_t + \mathbf{b}^c)$$

$$\mathbf{h}_t = \mathbf{o}_t \odot \tanh(\mathbf{c}_t)$$

$$\mathbf{s}_t = [\mathbf{h}_t; \mathbf{c}_t]$$

$$\mathbf{h}_{t-1}, \mathbf{c}_{t-1} = \mathbf{s}_{t-1}$$

$$\mathbf{z}_t = [\mathbf{x}_t; \mathbf{h}_{t-1}]$$

$$\mathbf{u}_t = \mathbf{W} \mathbf{z}_t + \mathbf{b}$$

$$\mathbf{i}_t = \sigma(\mathbf{u}_t[1..k])$$

$$\mathbf{f}_t = \sigma(\mathbf{u}_t[k+1..2k])$$

$$\mathbf{o}_t = \sigma(\mathbf{u}_t[2k+1..3k])$$

$$\mathbf{c}_t = \mathbf{f}_t \odot \mathbf{c}_{t-1} + \mathbf{i}_t \odot \tanh(\mathbf{u}_t[3k+1..4k])$$

$$\mathbf{h}_t = \mathbf{o}_t \odot \tanh(\mathbf{c}_t)$$

$$\mathbf{s}_t = [\mathbf{h}_t; \mathbf{c}_t]$$

Gated Recurrent Unit

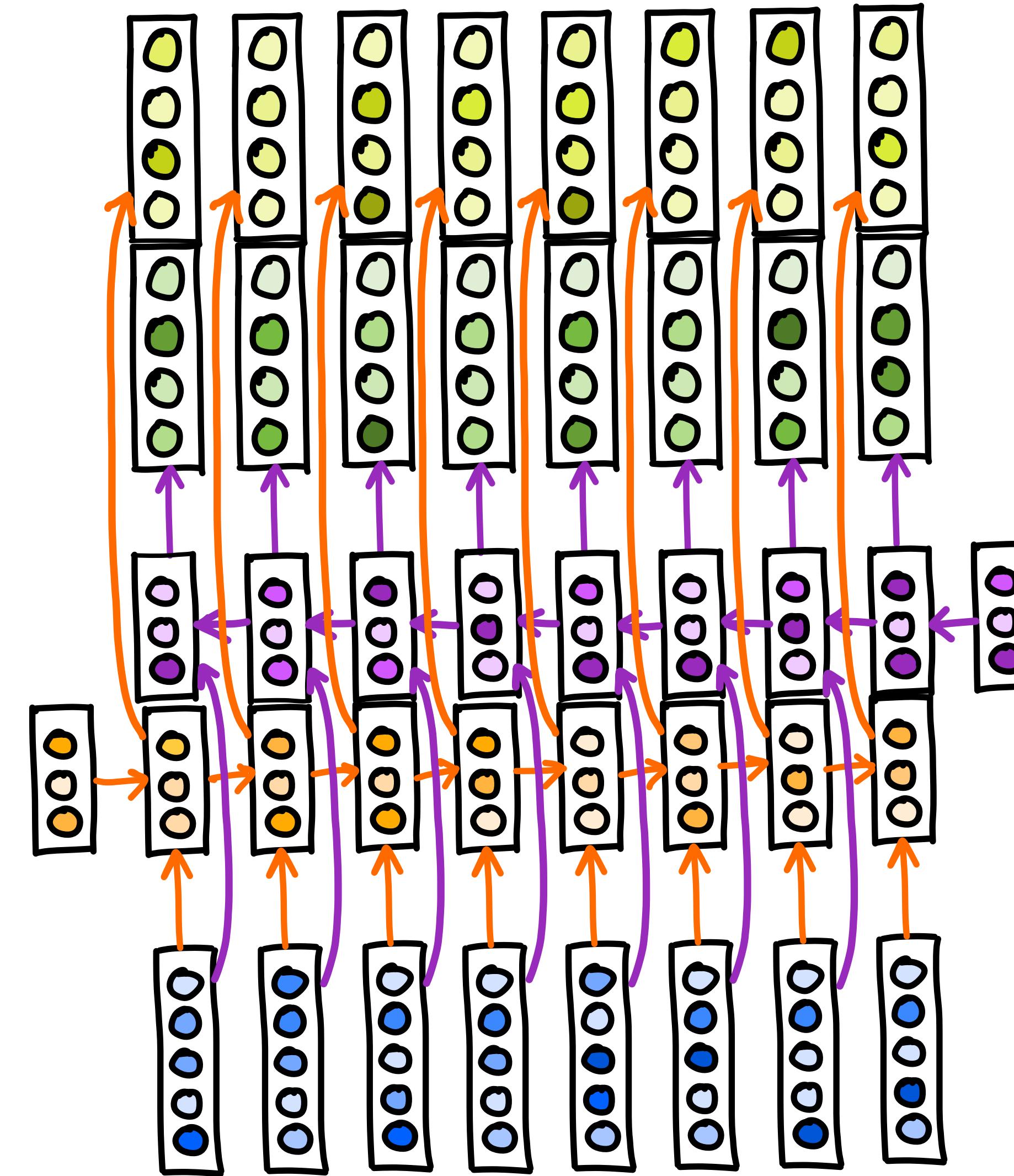
$$\mathbf{z}_t = \sigma(\mathbf{W}^z[\mathbf{x}_t; \mathbf{s}_{t-1}] + \mathbf{b}^z)$$

$$\mathbf{r}_t = \sigma(\mathbf{W}^r[\mathbf{x}_t; \mathbf{s}_{t-1}] + \mathbf{b}^r)$$

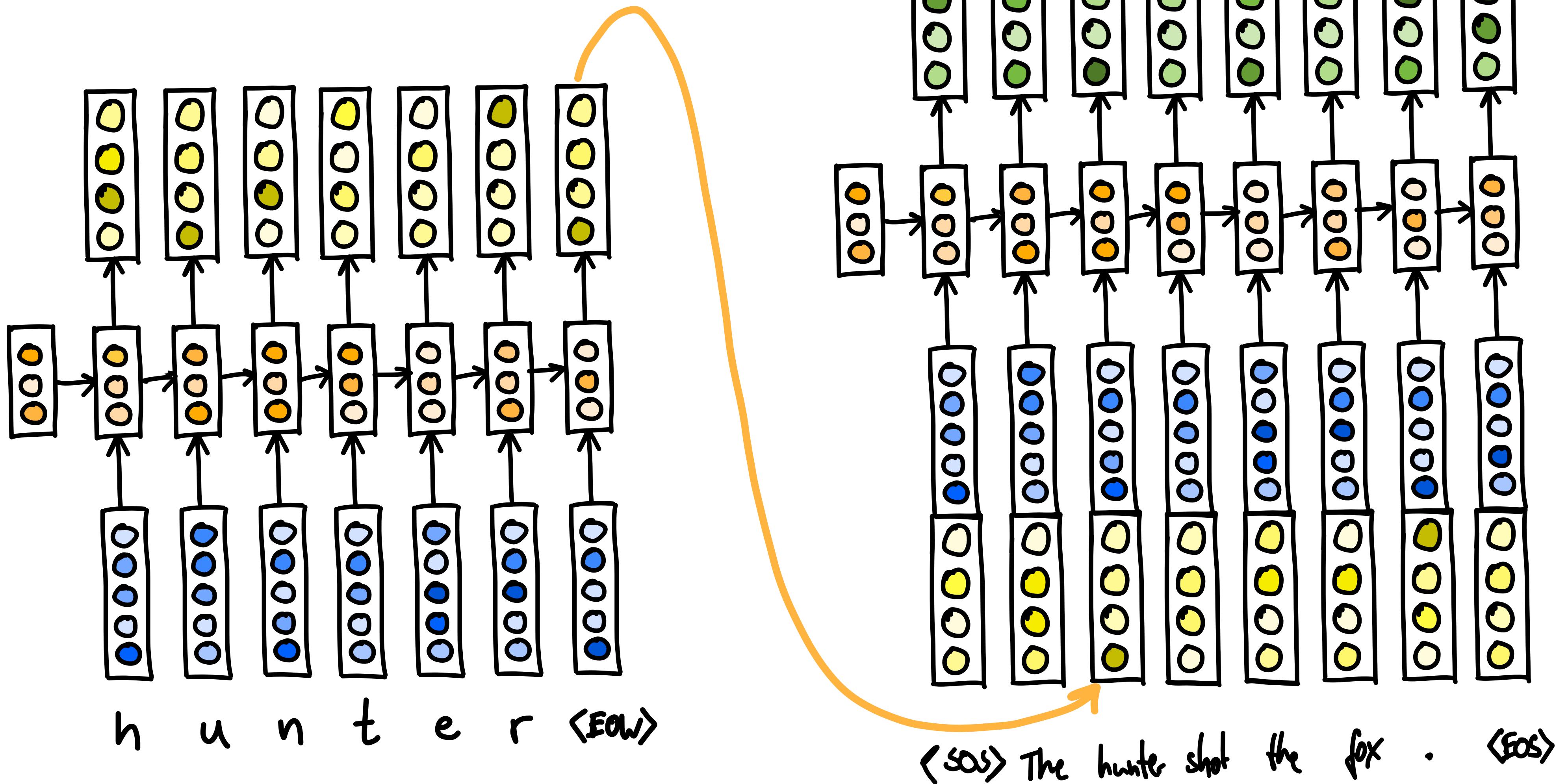
$$\mathbf{s}_t = (1 - \mathbf{z}_t) \odot \mathbf{s}_{t-1} + \mathbf{z}_t \odot \sigma(\mathbf{W}^x[\mathbf{x}_t; \mathbf{r}_t \odot \mathbf{s}_{t-1}] + \mathbf{b}^x)$$

$$\mathbf{h}_t = \mathbf{s}_t$$

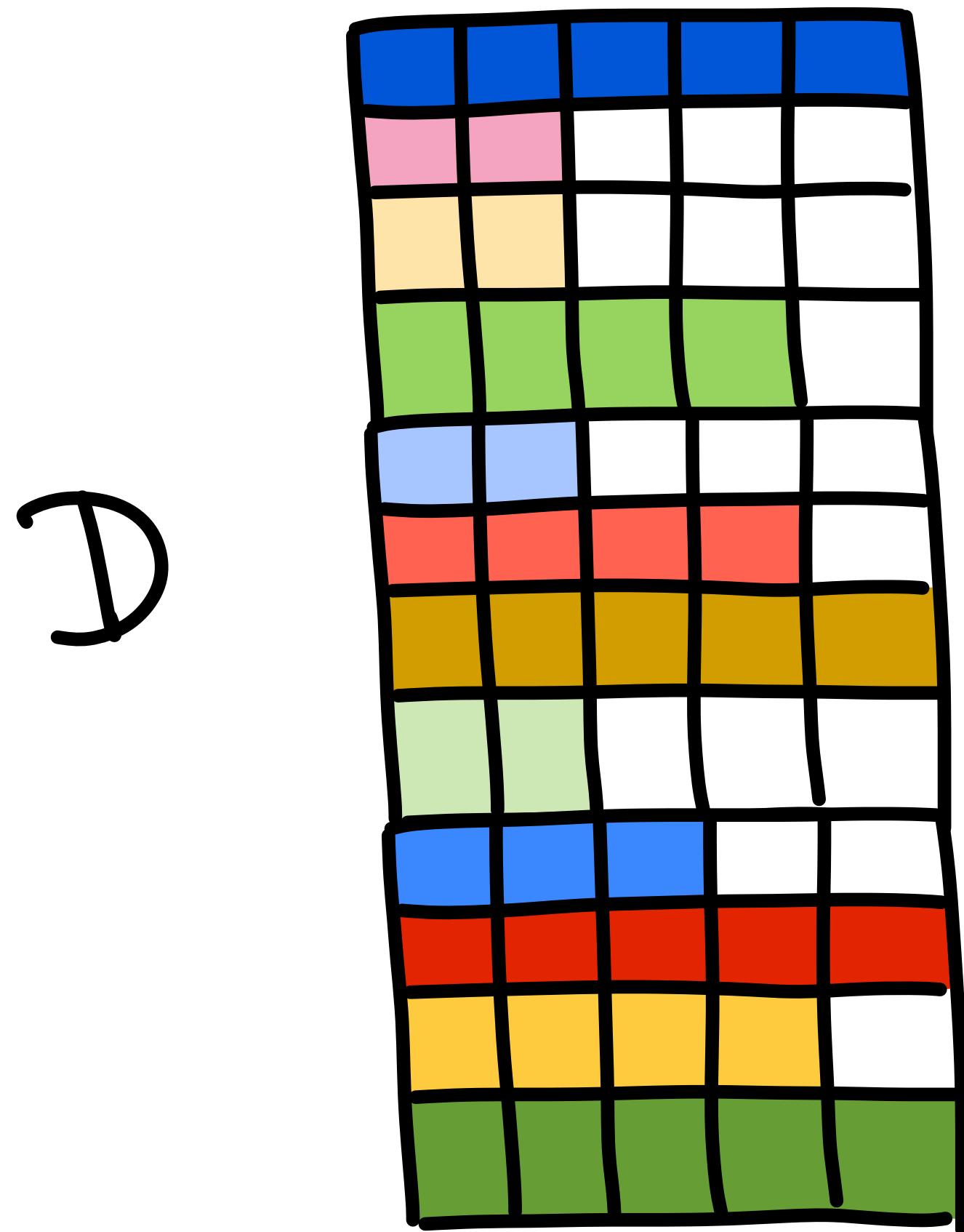
Bidirectional RNNs



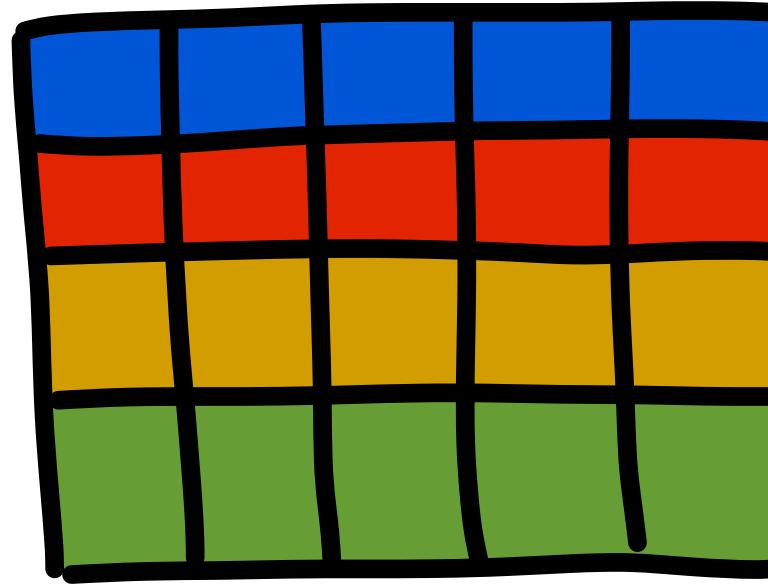
Character-level Models



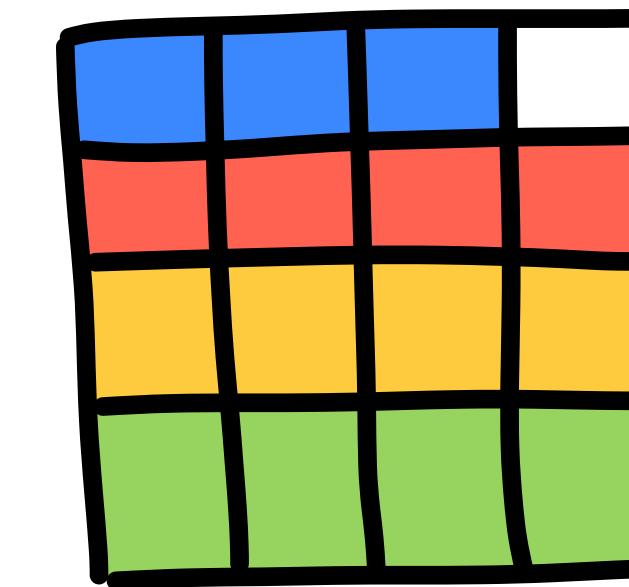
Bucketing



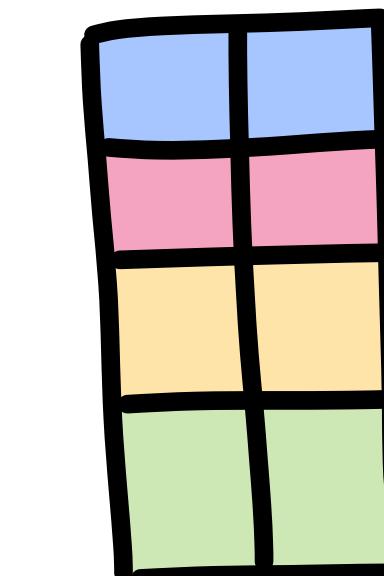
\mathcal{B}_5



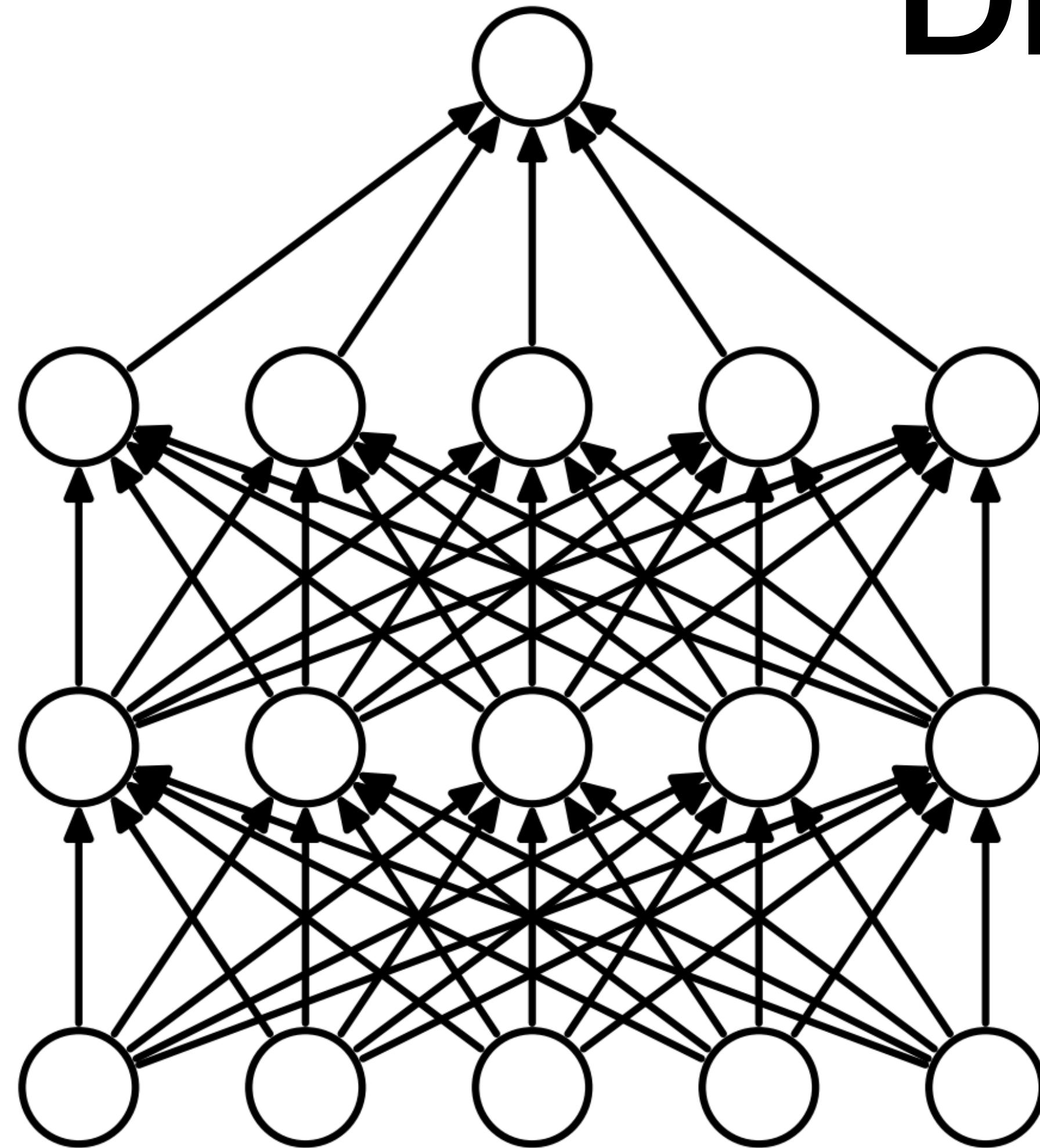
$\mathcal{B}_{3..4}$



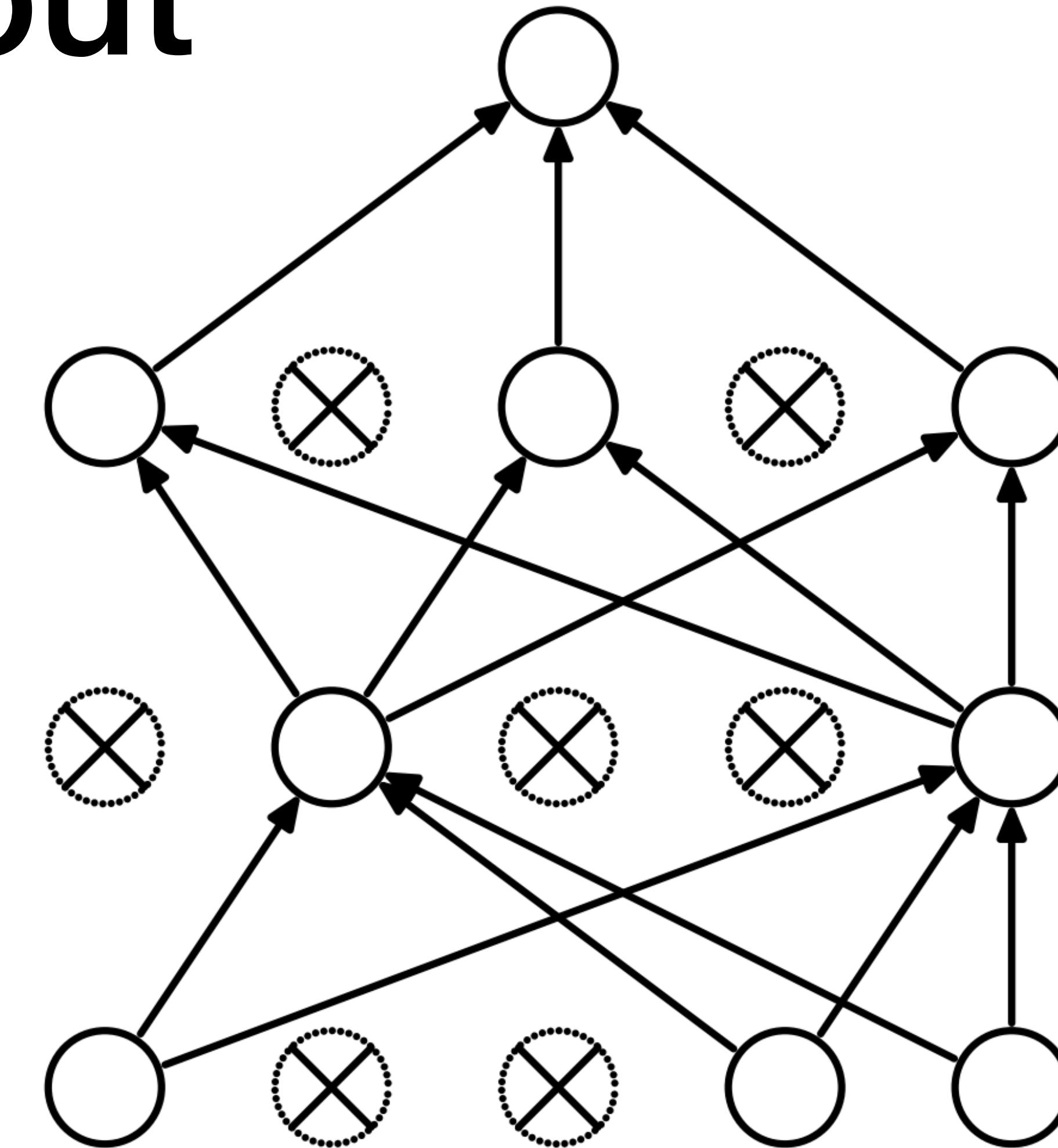
\mathcal{B}_2



Dropout

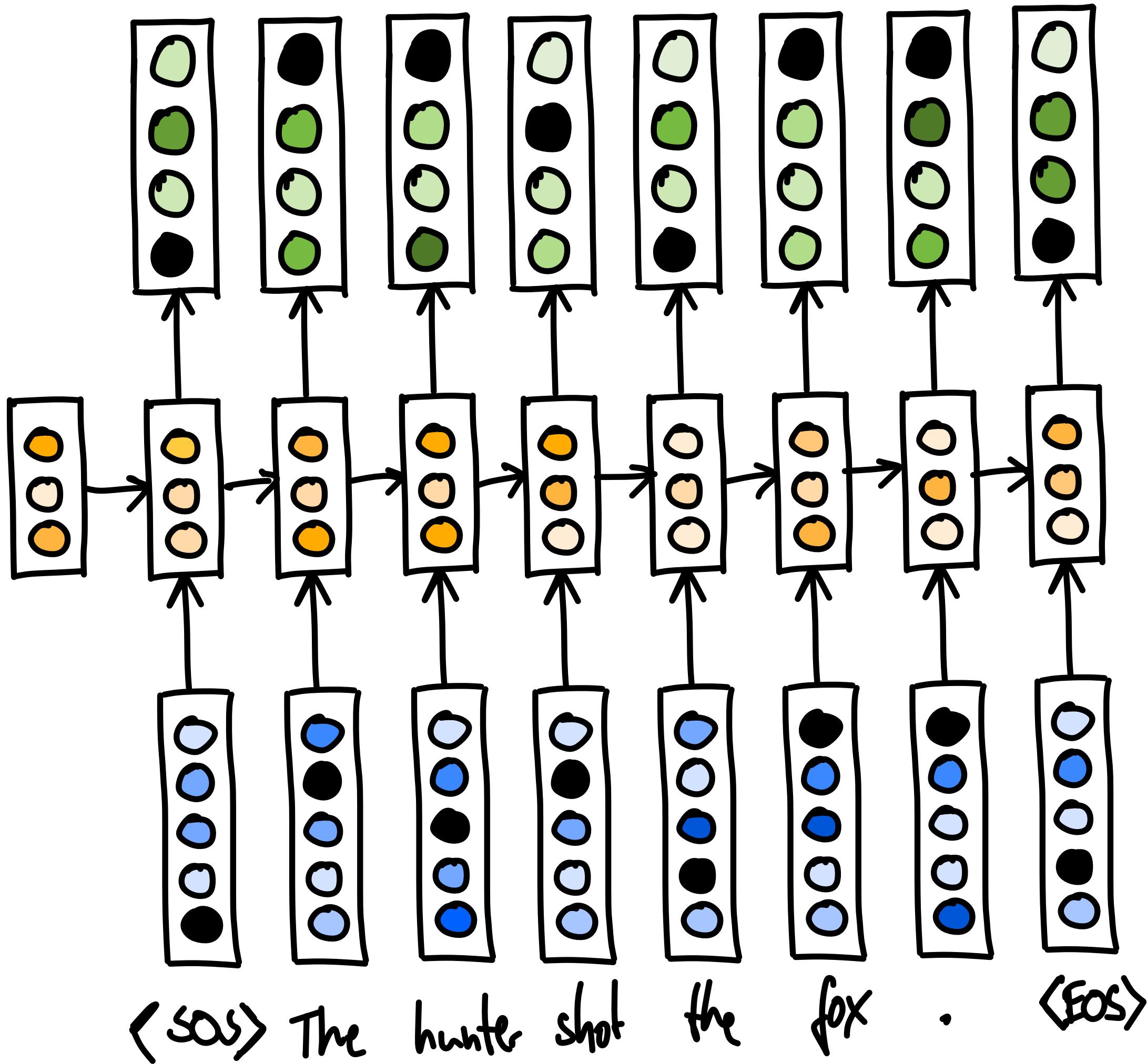


(a) Standard Neural Net

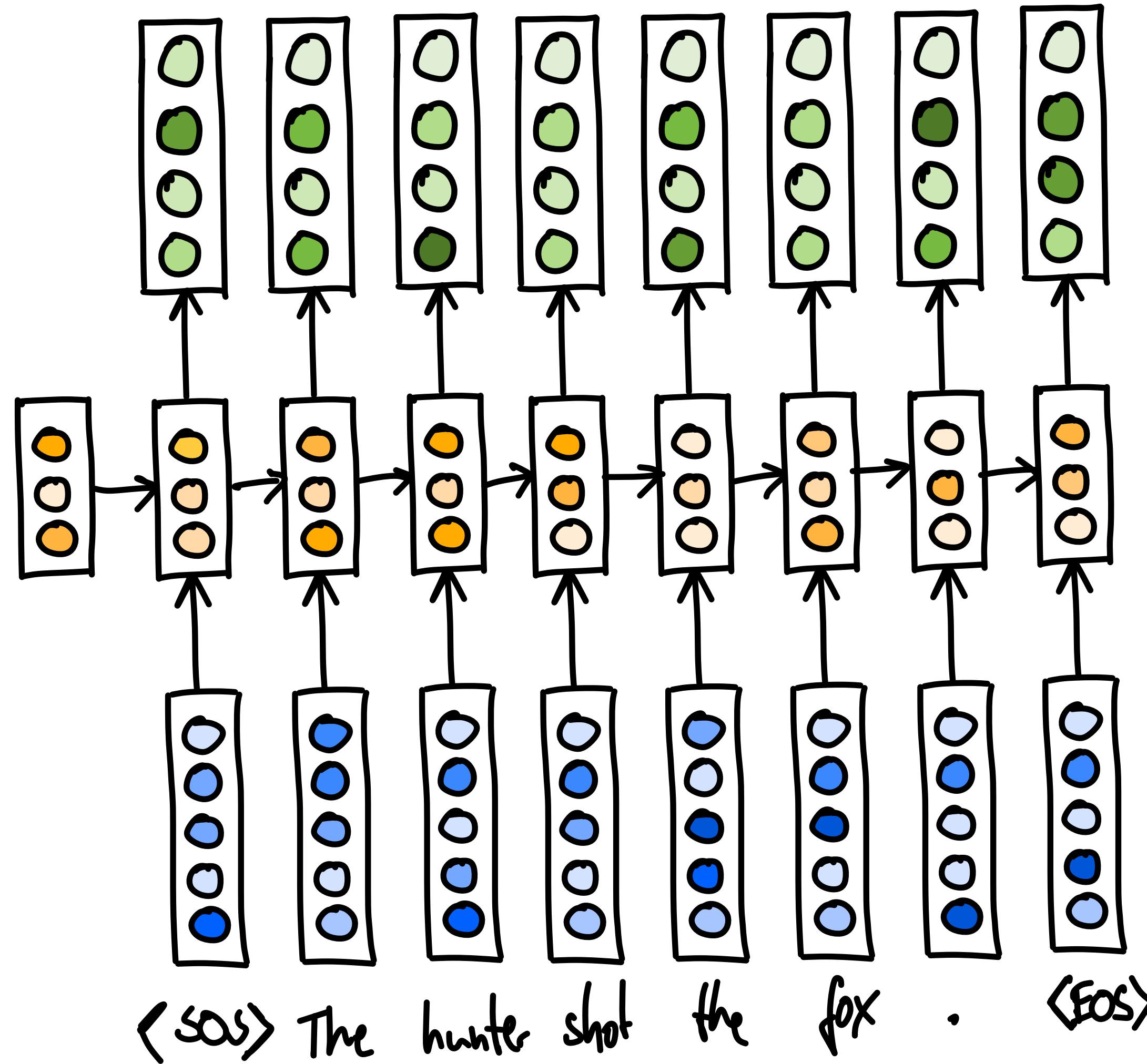


(b) After applying dropout.

Dropout



Recurrent Dropout?



Applications

- Language Modeling
- Recognizing Textual Entailment
- Named Entity Recognition
- Machine Translation
- Question Answering
- Dialog Modeling
- Language Generation
- Sentence Summarization
- Paraphrasing
- Sentiment Analysis
- ...

Recognizing Textual Entailment

A man is crowd surfing at a concert

- The man is at a football game **Contradiction**
- The man is drunk **Neutral**
- The man is at a concert **Entailment**

A wedding party is taking pictures

- There is a funeral **Contradiction**
- They are outside **Neutral**
- Someone got married **Entailment**

State-of-the-Art Before 2015

[Lai and Hockenmaier, 2014, Jimenez et al., 2014, Zhao et al., 2014, Beltagy et al., 2015 etc.]

- Engineered natural language processing pipelines
- Various external resources
- Specialized subcomponents
- Extensive manual creation of features:
 - Negation detection, word overlap, part-of-speech tags, dependency parses, alignment, unaligned matching, chunk alignment, synonym, hypernym, antonym, denotation graph...

Data Matters!

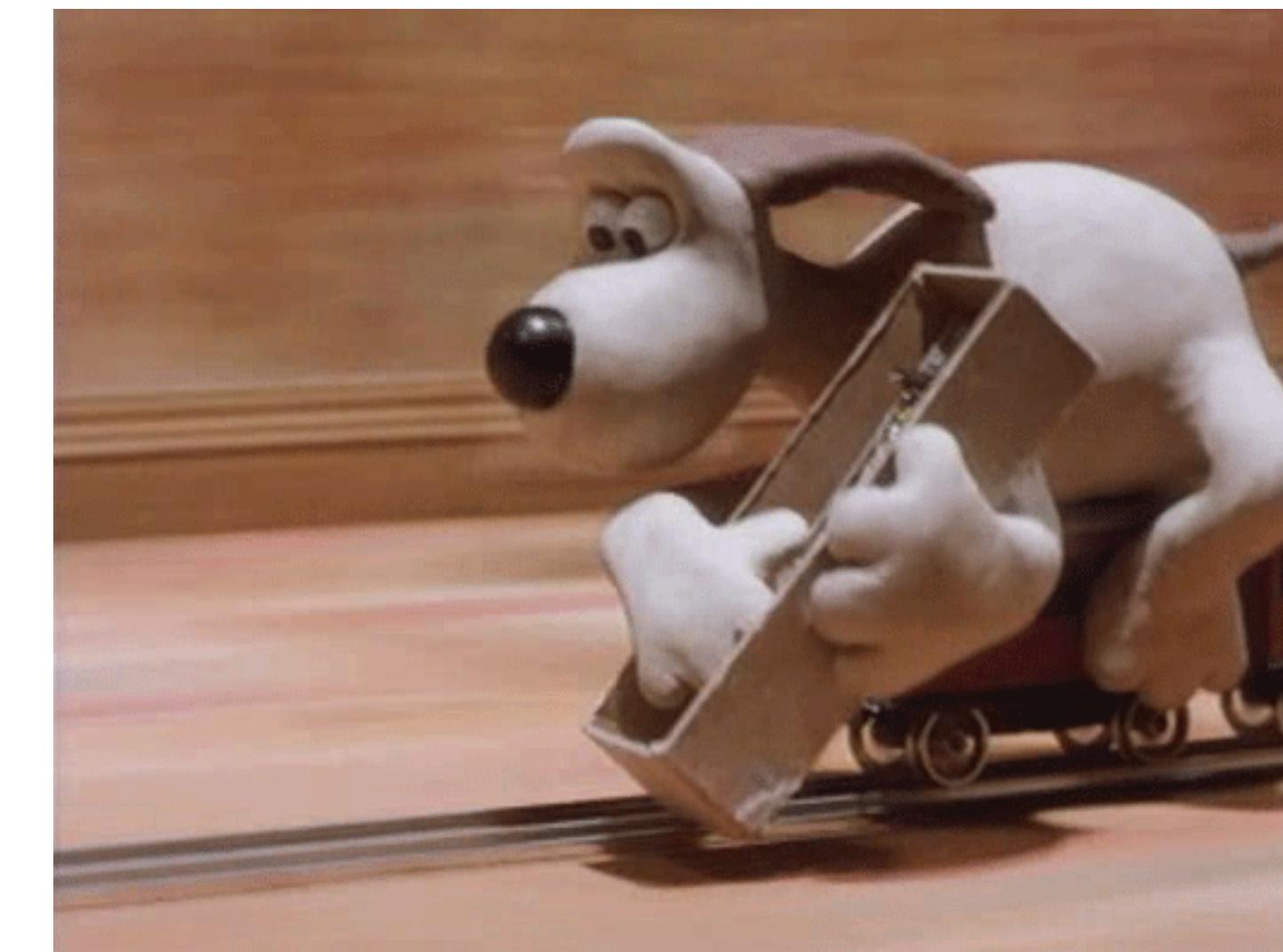
Neural Networks on SICK corpus
[Marelli et al., SemEval 2014]

10k sentence pairs, partly synthetic

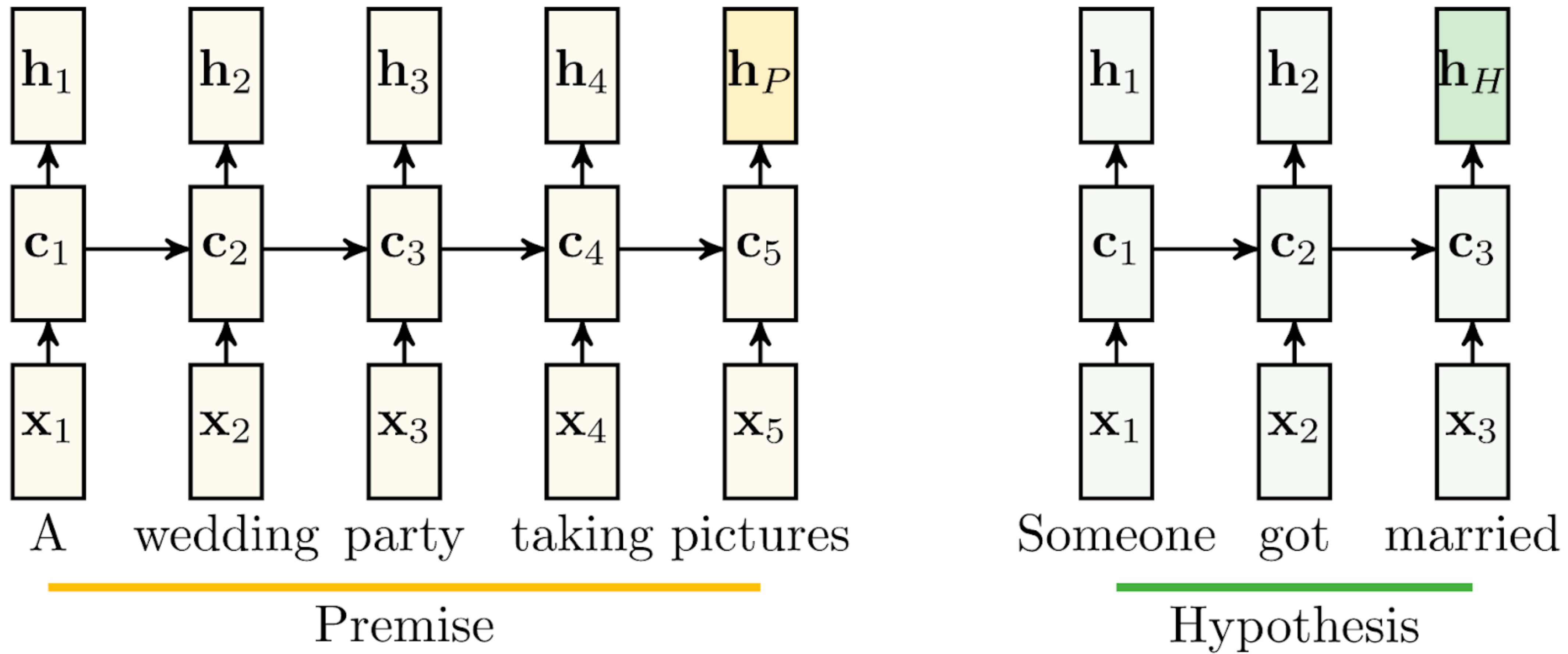


Neural Networks on SNLI corpus
[Bowman et al., EMNLP 2015]

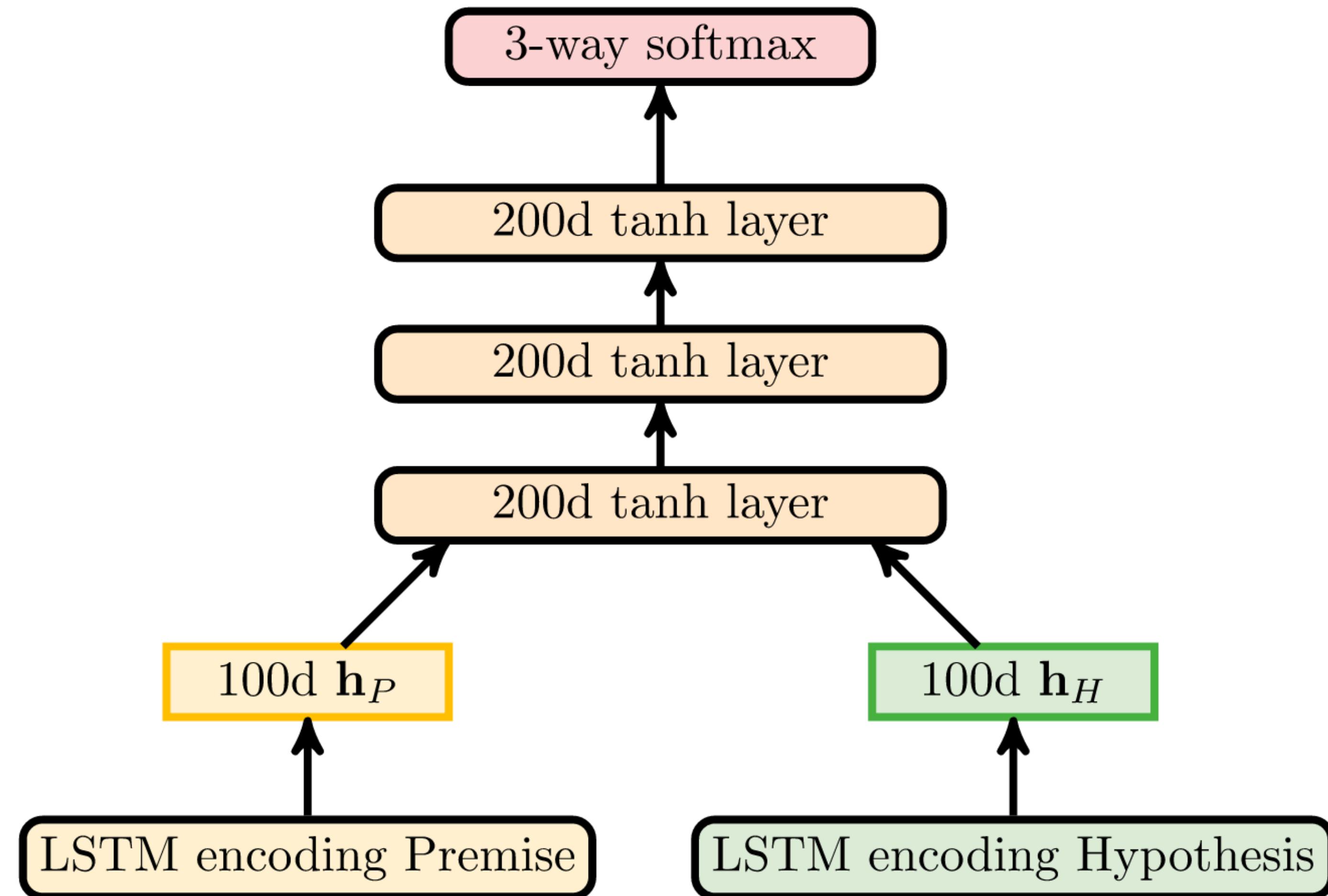
570k sentence pairs from Mechanical Turkers
EMNLP 2015 “best data set or resource” award!



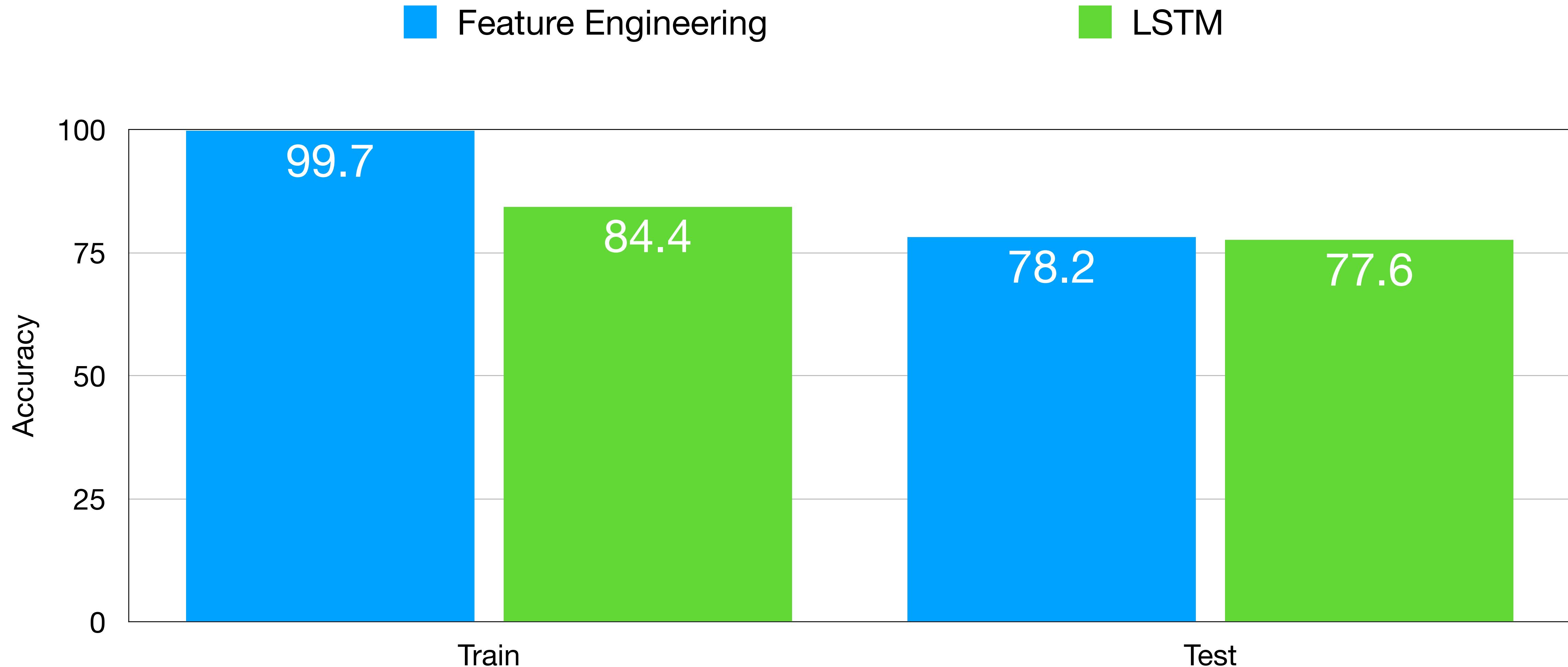
Independent Sentence Encoding



Independent Sentence Encoding



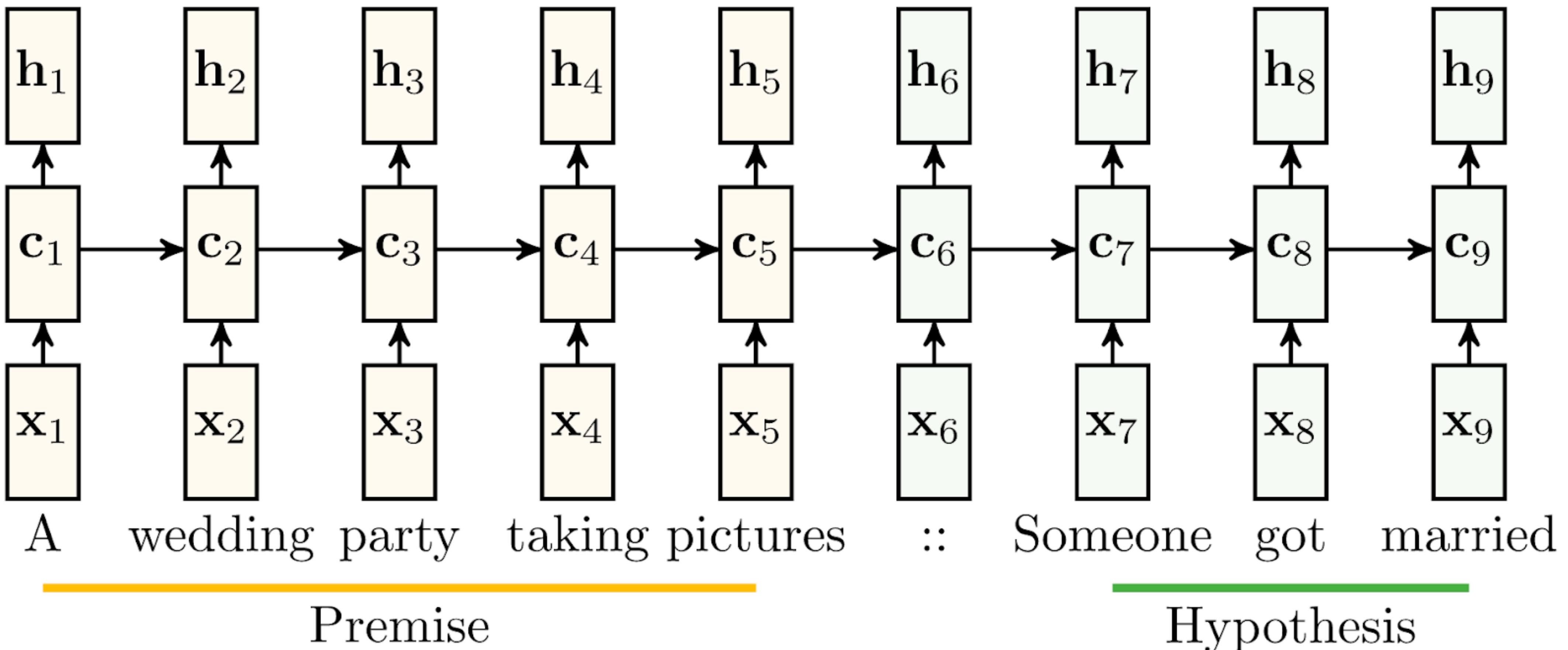
Results



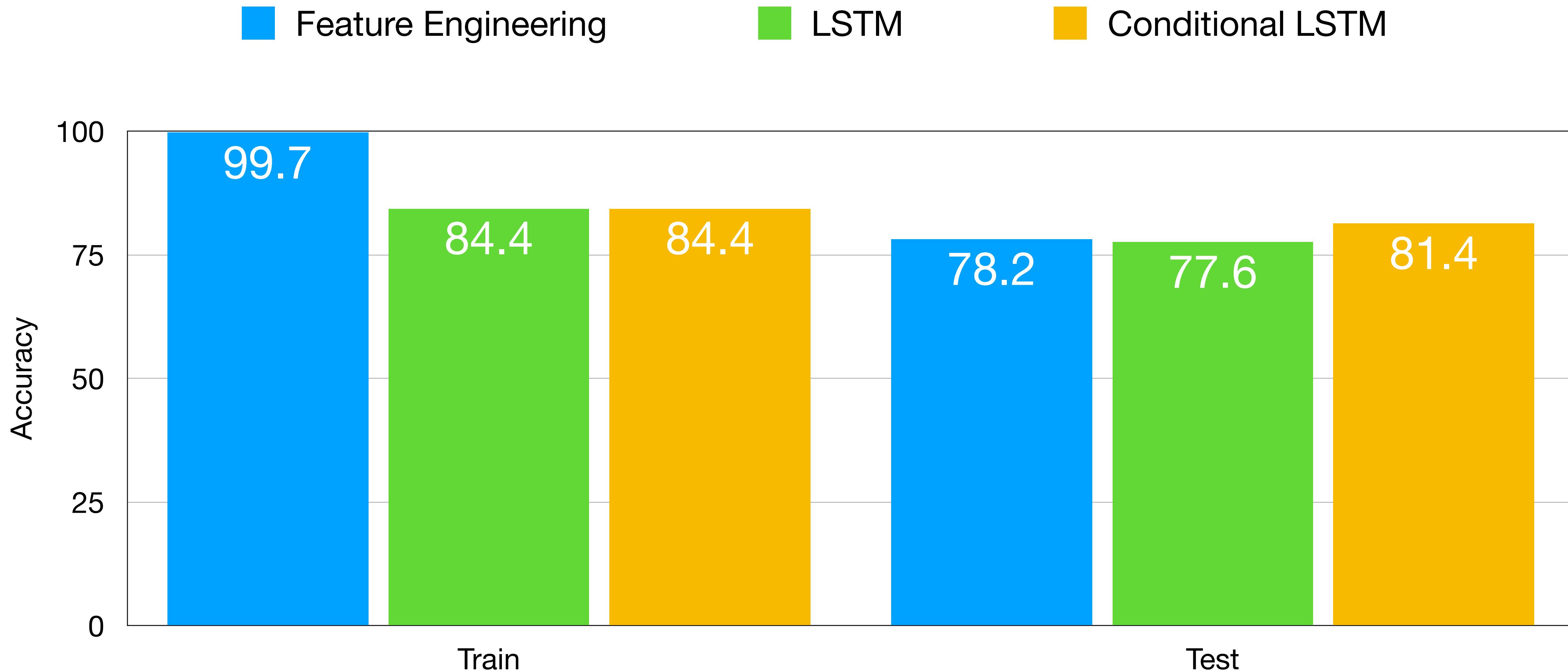
“You can’t cram the meaning of a whole %&!\$# sentence into a single \$&!#* vector!”

– Raymond J. Mooney

Conditional Encoding



Results



Named Entity Recognition

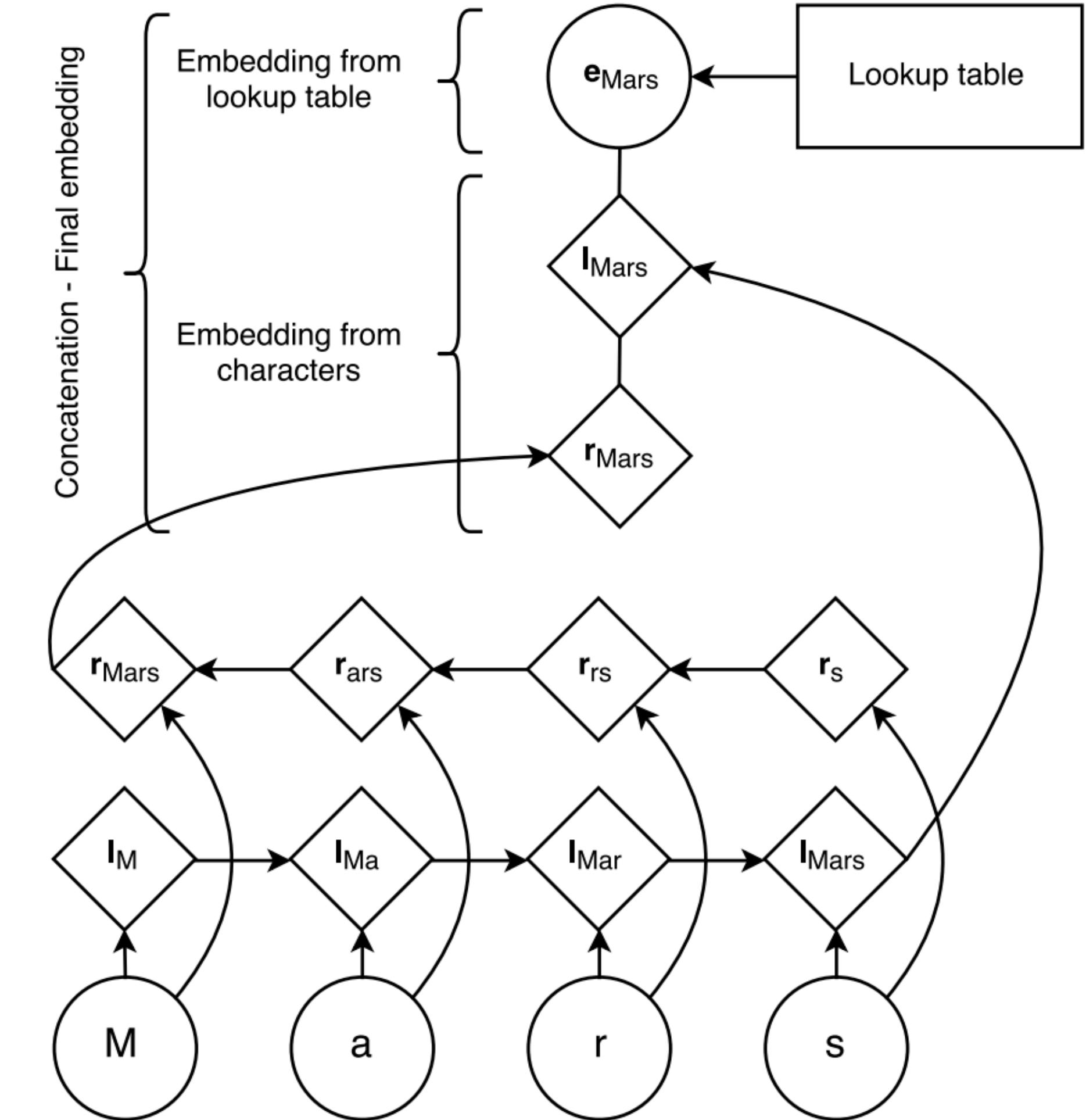
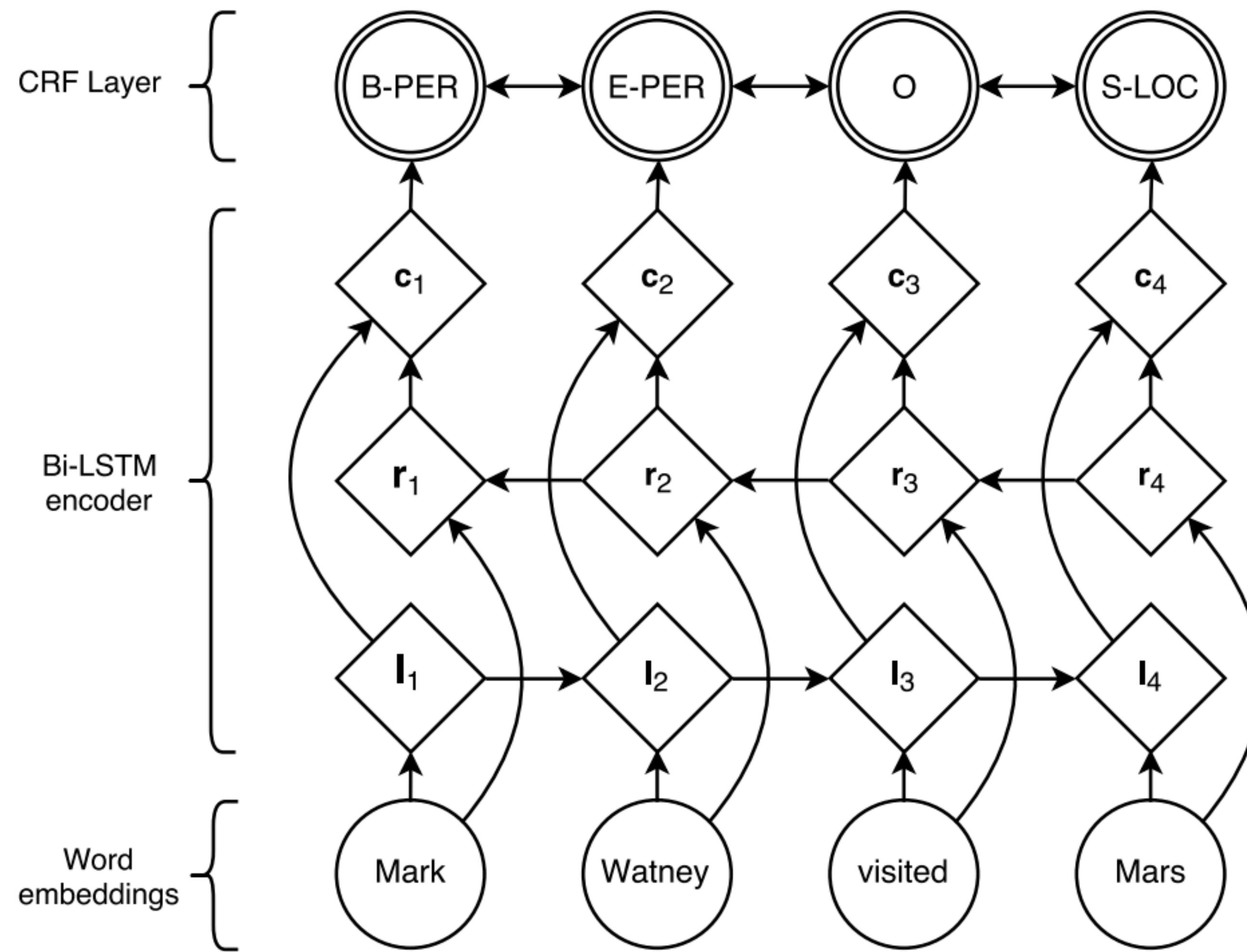
Chase Manhattan and its merger partner J.P.Morgan and Citibank, which was involved in moving about \$100 million for Raul Salinas de Gortari, brother of a former Mexican president, to banks in Switzerland, are also expected to sign on.

The diagram illustrates Named Entity Recognition (NER) annotations for the given sentence. It uses colored boxes to highlight entities and labels above them with their corresponding entity types. The annotations are as follows:

- "Chase Manhattan" is highlighted with a blue box and labeled "Organization".
- "J.P.Morgan" is highlighted with a blue box and labeled "Organization".
- "Citibank" is highlighted with a blue box and labeled "Org".
- "\$100 million" is highlighted with a green box and labeled "Money".
- "Raul Salinas de Gortari" is highlighted with an orange box and labeled "Person".
- "former Mexican president" is highlighted with an orange box and labeled "Person".
- "Switzerland" is highlighted with a green box and labeled "GPE".

A red arrow labeled "Family" points from the "Person" entity under "Raul Salinas de Gortari" to the "Person" entity under "former Mexican president".

Bidirectional Character-aware LSTM



Results

Collobert et al. (2011) Lin and Wu (2009) Passos et al. (2014)
Chiu and Nichols (2015) Lample et al. (2016)

