

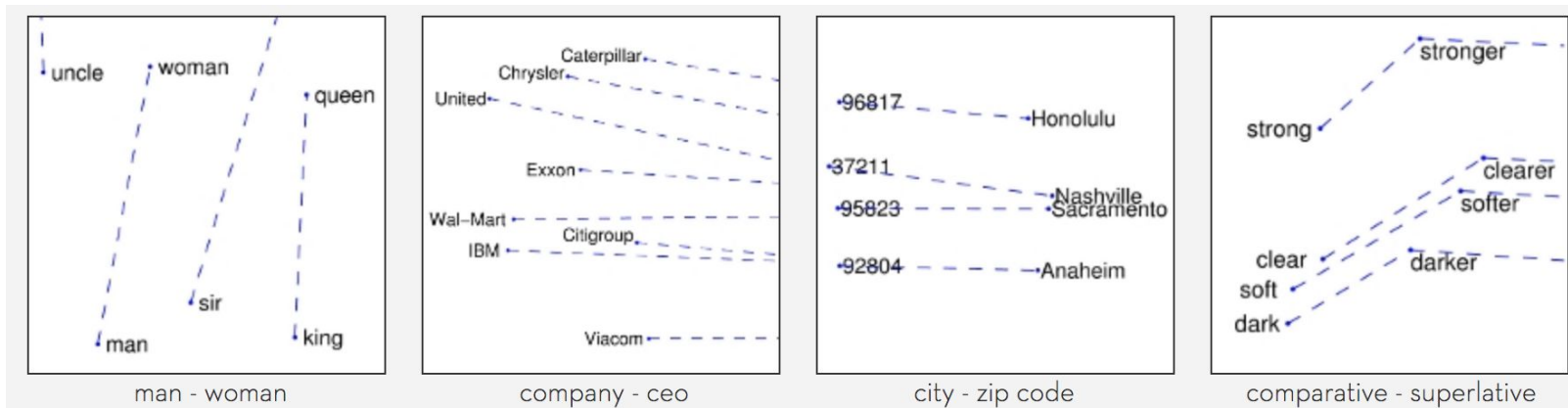
# **ELMo**

# **Embeddings from Language Models**

Mikhail Arkhipov - [arkhipov@yahoo.com](mailto:arkhipov@yahoo.com)  
Deep Learning Lab MIPT

# Embeddings

- tl;dr - a vector representation of words
- Every word gets a trainable vector; surprisingly, embeddings create a “geometry” for words.
- First step of neural language modeling;
- Typically, 100-300 dimensional



# Language modeling

By accurately assigning probability to a natural sequence (words or characters), you can improve:

- **Machine Translation:**  $p(\text{strong tea}) > p(\text{powerful tea})$
- **Speech Recognition:**  $p(\text{speech recognition}) > p(\text{speech wreck ignition})$
- **Question Answering / Summarization:**  
 $p(\text{President X attended ...})$  is higher for  $X=\text{Obama}$

# Traditional approaches

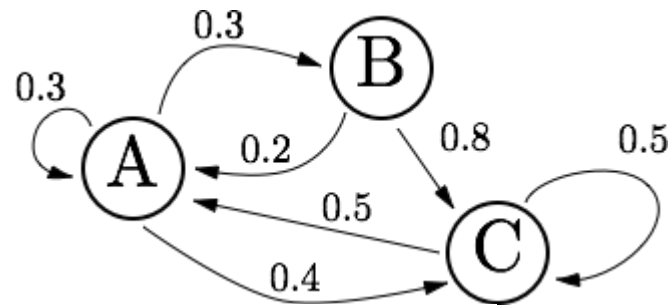
n-gram models and their adaptations:

$P(\text{I, saw, the, red, house})$

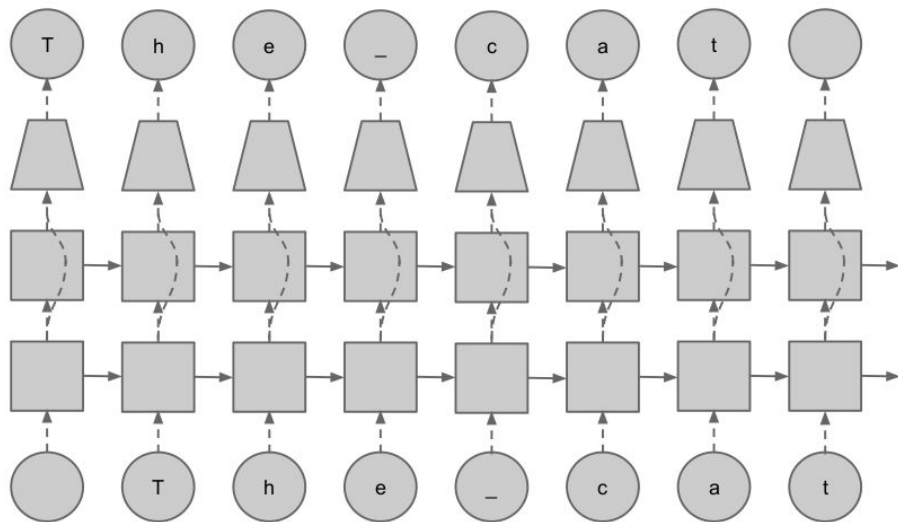
$\approx P(\text{I} \mid \langle s \rangle, \langle s \rangle)P(\text{saw} \mid \langle s \rangle, \text{I})P(\text{the} \mid \text{I, saw})P(\text{red} \mid \text{saw, the})P(\text{house} \mid \text{the, red})P(\langle /s \rangle \mid \text{red, house})$

Issues:

- What to do if n-gram has never been seen?
- How do you choose n for the n-grams?
- “Deep Learning is **amazing** because” v/s “Deep Learning is **awesome** because”

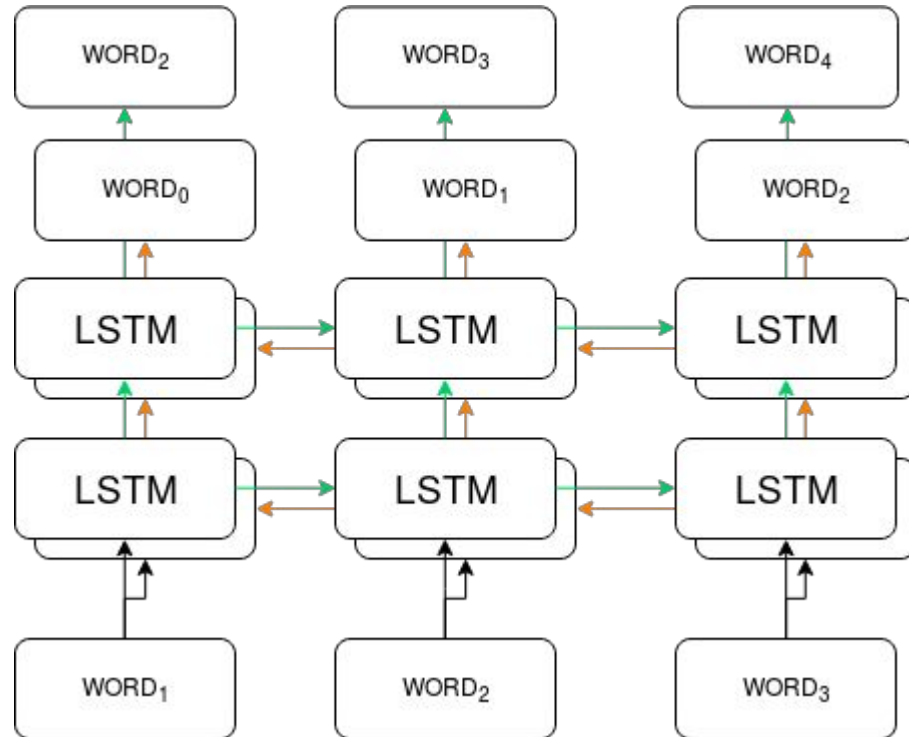


# Language modeling with Recurrent Networks

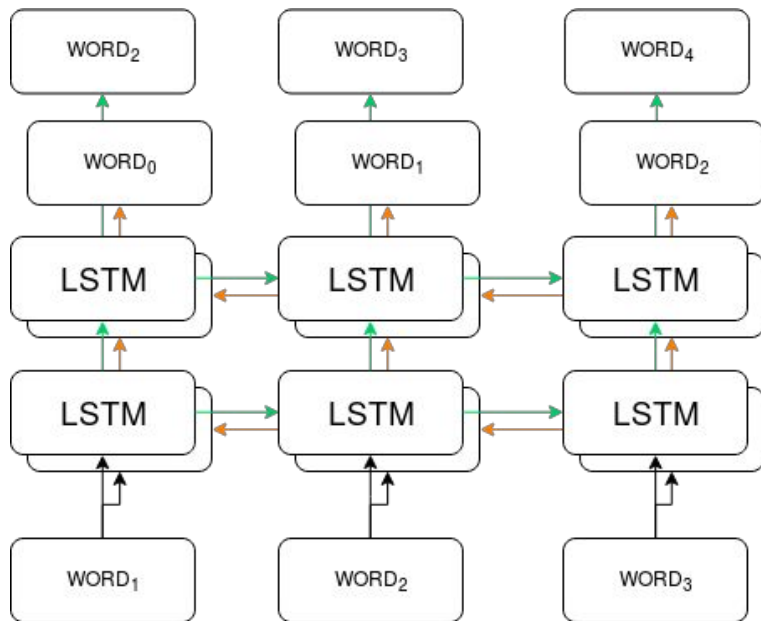


$$p(t_1, t_2, \dots, t_N) = \prod_{k=1}^N p(t_k \mid t_1, t_2, \dots, t_{k-1}).$$

# Bi-directional Language Model

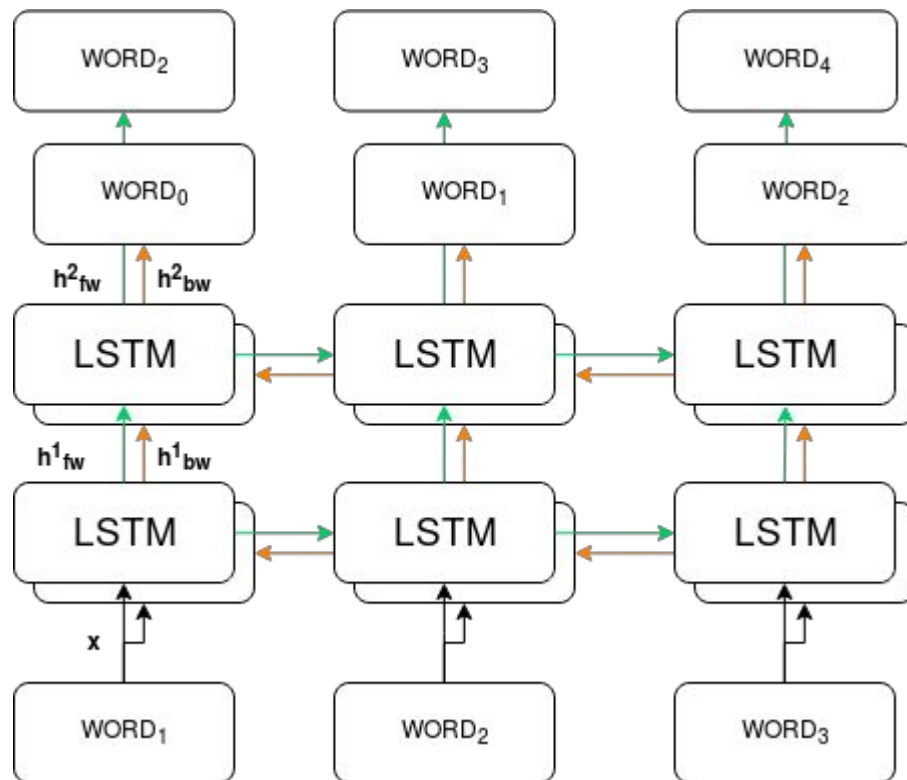


# Bi-directional Language Model



$$\sum_{k=1}^N ( \log p(t_k \mid t_1, \dots, t_{k-1}; \Theta_x, \vec{\Theta}_{LSTM}, \Theta_s) \\ + \log p(t_k \mid t_{k+1}, \dots, t_N; \Theta_x, \overleftarrow{\Theta}_{LSTM}, \Theta_s) ).$$

# ELMo

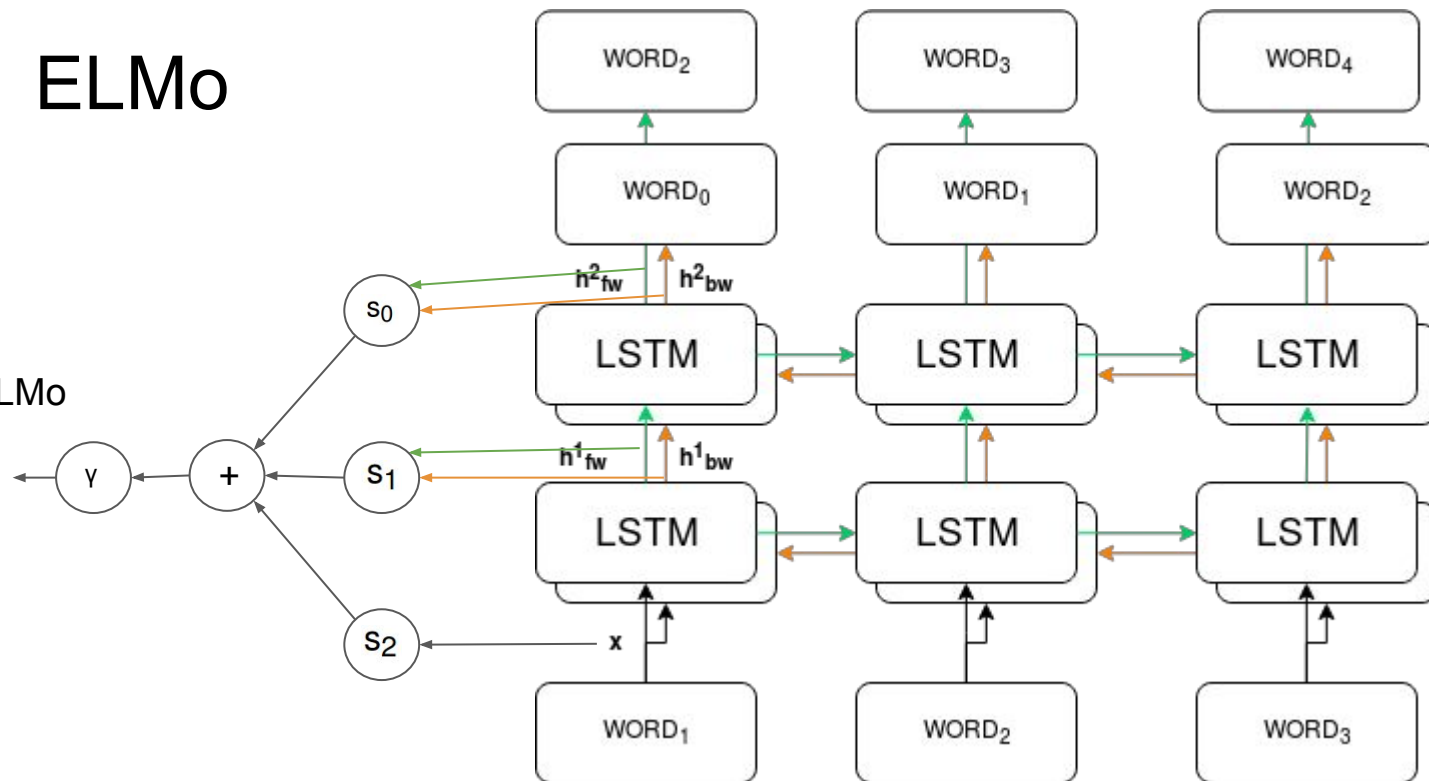


$$\begin{aligned}
 R_k &= \{\mathbf{x}_k^{LM}, \vec{\mathbf{h}}_{k,j}^{LM}, \overleftarrow{\mathbf{h}}_{k,j}^{LM} \mid j = 1, \dots, L\} \\
 &= \{\mathbf{h}_{k,j}^{LM} \mid j = 0, \dots, L\},
 \end{aligned}
 \quad \text{ELMo}_k^{\text{task}} = E(R_k; \Theta^{\text{task}}) = \gamma^{\text{task}} \sum_{j=0}^L s_j^{\text{task}} \mathbf{h}_{k,j}^{LM}.$$



# ELMo

ELMo

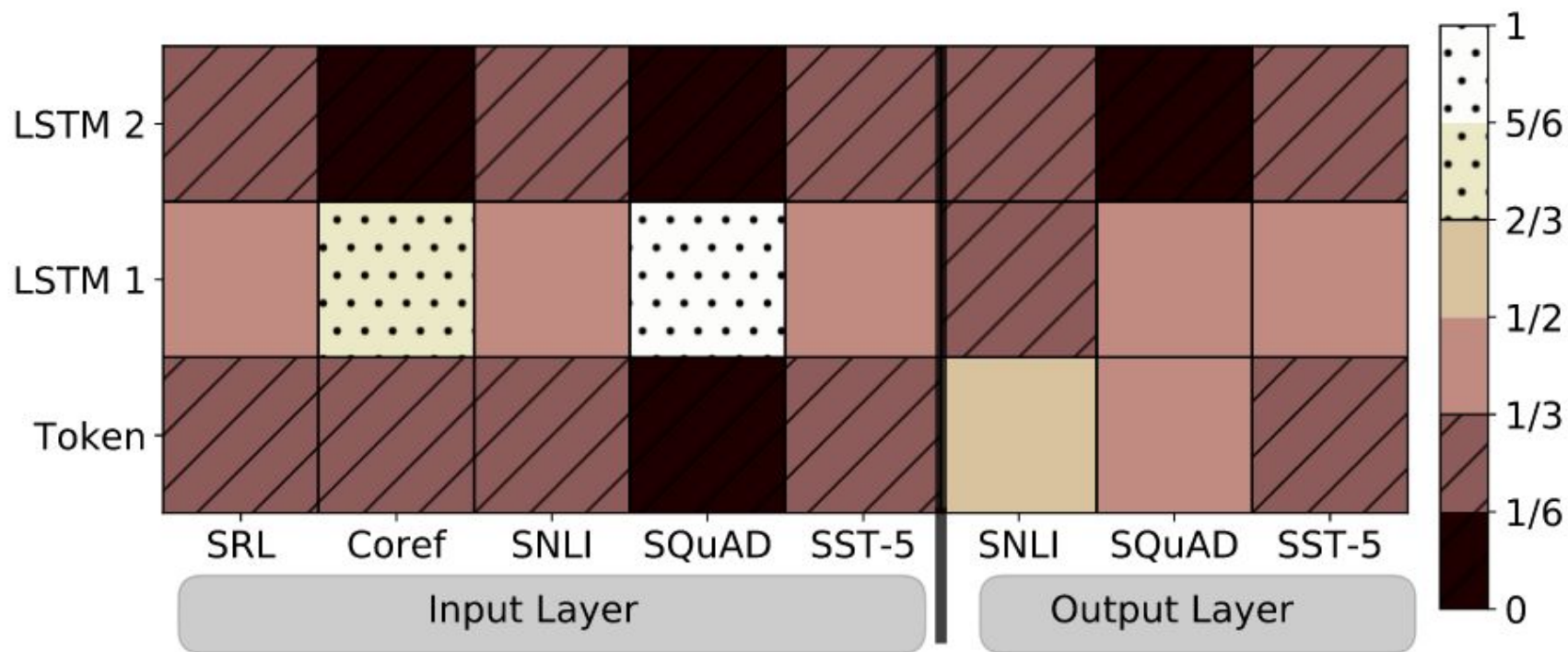


$$\begin{aligned}
 R_k &= \{\mathbf{x}_k^{LM}, \vec{\mathbf{h}}_{k,j}^{LM}, \overleftarrow{\mathbf{h}}_{k,j}^{LM} \mid j = 1, \dots, L\} \\
 &= \{\mathbf{h}_{k,j}^{LM} \mid j = 0, \dots, L\},
 \end{aligned}
 \quad
 \text{ELMo}_k^{task} = E(R_k; \Theta^{task}) = \gamma^{task} \sum_{j=0}^L s_j^{task} \mathbf{h}_{k,j}^{LM}.$$

# ELMo results

| TASK  | PREVIOUS SOTA        |                  | OUR<br>BASELINE | ELMo +<br>BASELINE | INCREASE<br>(ABSOLUTE/<br>RELATIVE) |
|-------|----------------------|------------------|-----------------|--------------------|-------------------------------------|
| SQuAD | Liu et al. (2017)    | 84.4             | 81.1            | 85.8               | 4.7 / 24.9%                         |
| SNLI  | Chen et al. (2017)   | 88.6             | 88.0            | $88.7 \pm 0.17$    | 0.7 / 5.8%                          |
| SRL   | He et al. (2017)     | 81.7             | 81.4            | 84.6               | 3.2 / 17.2%                         |
| Coref | Lee et al. (2017)    | 67.2             | 67.2            | 70.4               | 3.2 / 9.8%                          |
| NER   | Peters et al. (2017) | $91.93 \pm 0.19$ | 90.15           | $92.22 \pm 0.10$   | 2.06 / 21%                          |
| SST-5 | McCann et al. (2017) | 53.7             | 51.4            | $54.7 \pm 0.5$     | 3.3 / 6.8%                          |

# ELMo layer distribution



Spasibo