# Fast decoding in neural machine translation with Ray
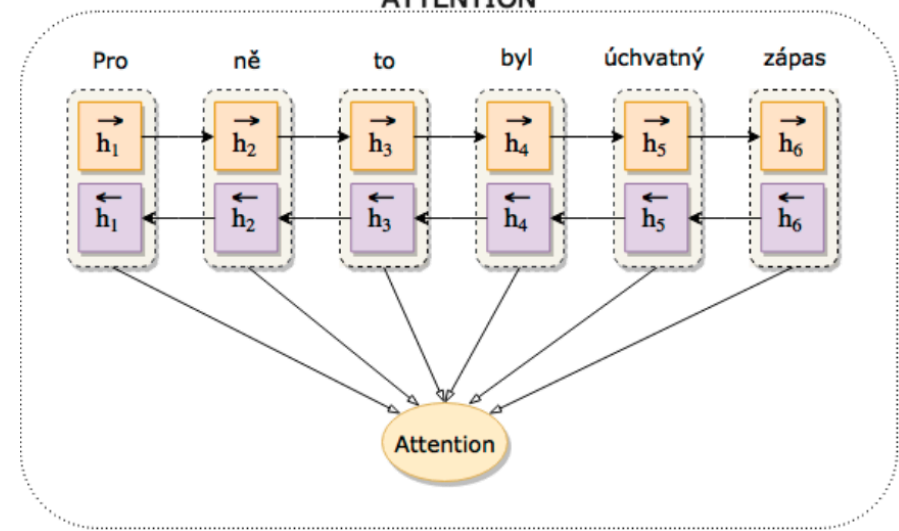
MAREK STRELEC
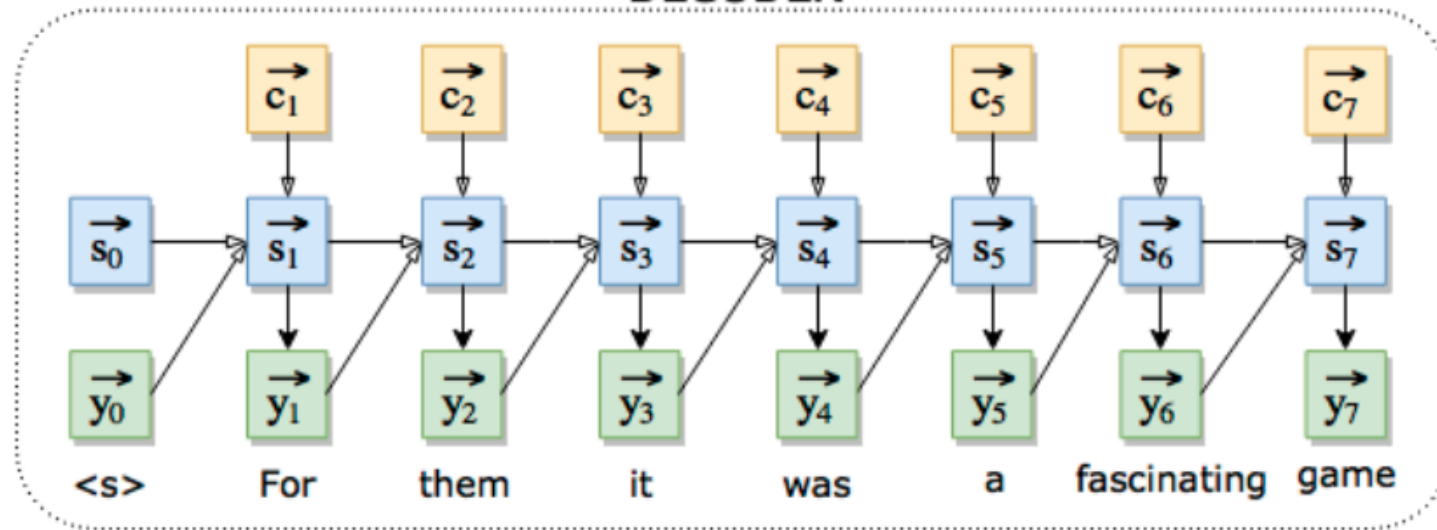
## ENCODER

Pro — $\overrightarrow{h_1}$, ně — $\overrightarrow{h_2}$, to — $\overrightarrow{h_3}$, byl — $\overrightarrow{h_4}$, úchvatný — $\overrightarrow{h_5}$, zápas — $\overrightarrow{h_6}$

$\overleftarrow{h_6}$, $\overleftarrow{h_5}$, $\overleftarrow{h_4}$, $\overleftarrow{h_3}$, $\overleftarrow{h_1}$, $\overleftarrow{h_1}$

C

## ATTENTION

Pro — $\overrightarrow{h_1}$, ně — $\overrightarrow{h_2}$, to — $\overrightarrow{h_3}$, byl — $\overrightarrow{h_4}$, úchvatný — $\overrightarrow{h_5}$, zápas — $\overrightarrow{h_6}$

$\overleftarrow{h_1}$, $\overleftarrow{h_2}$, $\overleftarrow{h_3}$, $\overleftarrow{h_4}$, $\overleftarrow{h_5}$, $\overleftarrow{h_6}$

Attention

## DECODER

$\overrightarrow{c_1}$, $\overrightarrow{c_2}$, $\overrightarrow{c_3}$, $\overrightarrow{c_4}$, $\overrightarrow{c_5}$, $\overrightarrow{c_6}$, $\overrightarrow{c_7}$

$\overrightarrow{s_0}$, $\overrightarrow{s_1}$, $\overrightarrow{s_2}$, $\overrightarrow{s_3}$, $\overrightarrow{s_4}$, $\overrightarrow{s_5}$, $\overrightarrow{s_6}$, $\overrightarrow{s_7}$

$\overrightarrow{y_0}$, $\overrightarrow{y_1}$, $\overrightarrow{y_2}$, $\overrightarrow{y_3}$, $\overrightarrow{y_4}$, $\overrightarrow{y_5}$, $\overrightarrow{y_6}$, $\overrightarrow{y_7}$

<s> — For — them — it — was — a — fascinating — game

# Time cost statistics for decoding

| Calculation Units | GPU | | CPU | |
|---|---|---|---|---|
| | Time(s) | Percentage | Time(s) | Percentage |
| Eq. (6): $s_j = f(e_{y^*_{j-1}}, s_{j-1}, c_j)$ | 551.07 | 75.73% | 1370.92 | 19.42% |
| Eq. (7): $t_j = g(e_{y^*_{j-1}}, c_j, s_j)$ | 88.25 | 12.13% | 277.76 | 3.93% |
| Eq. (8): $o_j = \mathbf{W}_o t_j$ | 25.33 | 3.48% | 2342.53 | 33.18% |
| Eq. (9): $\mathcal{D}_j = \text{softmax}(o_j)$ | 63.00 | 8.66% | 3069.25 | 43.47% |

Zhang, Wen, et al. (2018)

# Ray

- ❑ "A flexible, high-performance distributed execution framework"

- ❑ Implements a dynamic task graph computation model

- ❑ Global Control Store

- ❑ Bottom-up distributed scheduler

- ❑ Actor abstraction

# Steps

❑ Implement an NMT model in TensorFlow

❑ Train the model on a subset of parallel data (Europarl)

❑ Experiments
  ❑ Distributed batched translation
  ❑ Distributed Beam Search
  ❑ Dynamic Beam Search
  ❑ Heterogeneous environment

❑ Compare times and BLEU score

# Thank you!