# Parallelizable StackLSTM

Shuoyang Ding    Philipp Koehn

NAACL 2019 Structured Prediction Workshop
Minneapolis, MN, United States
June 7th, 2019
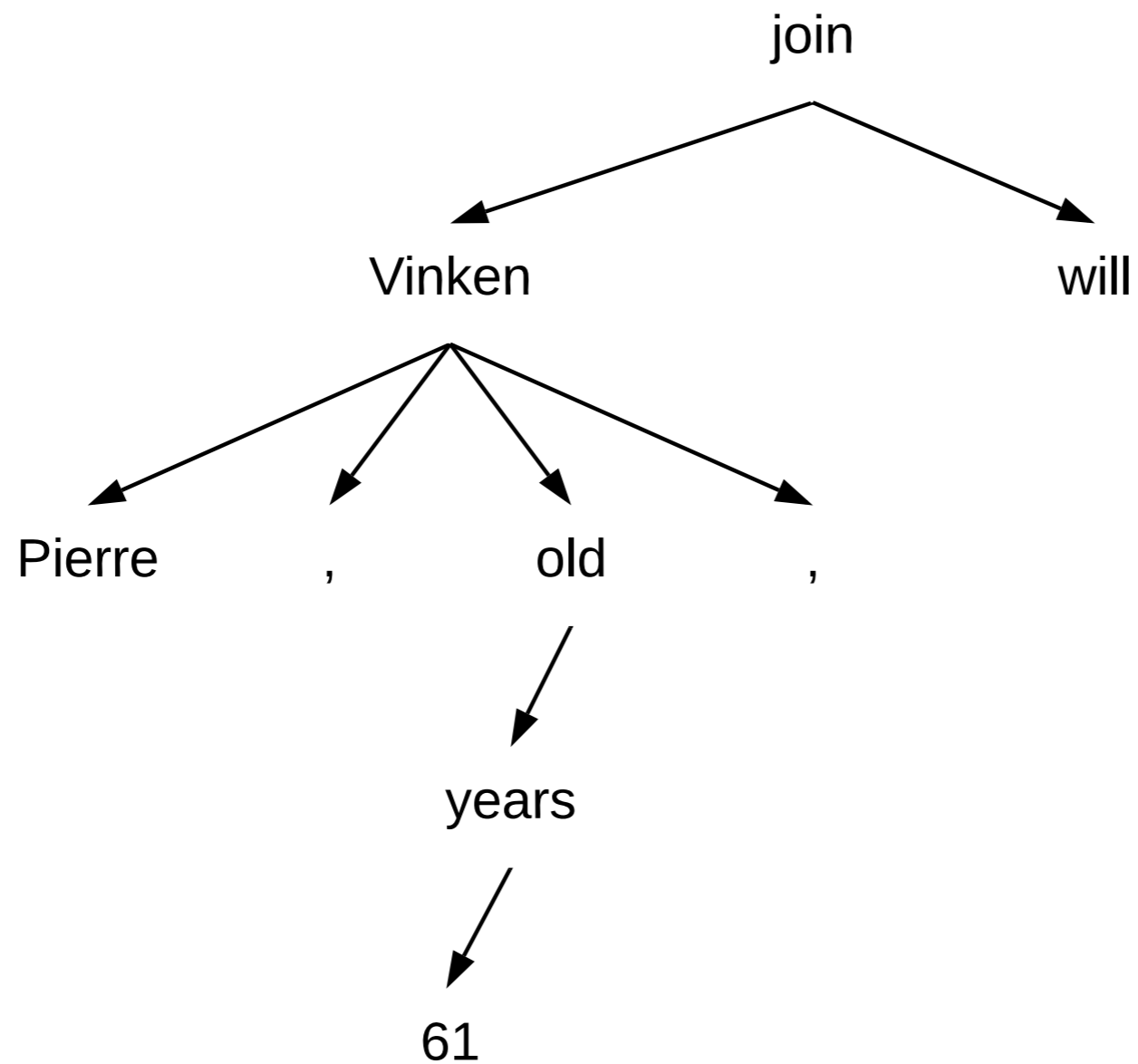
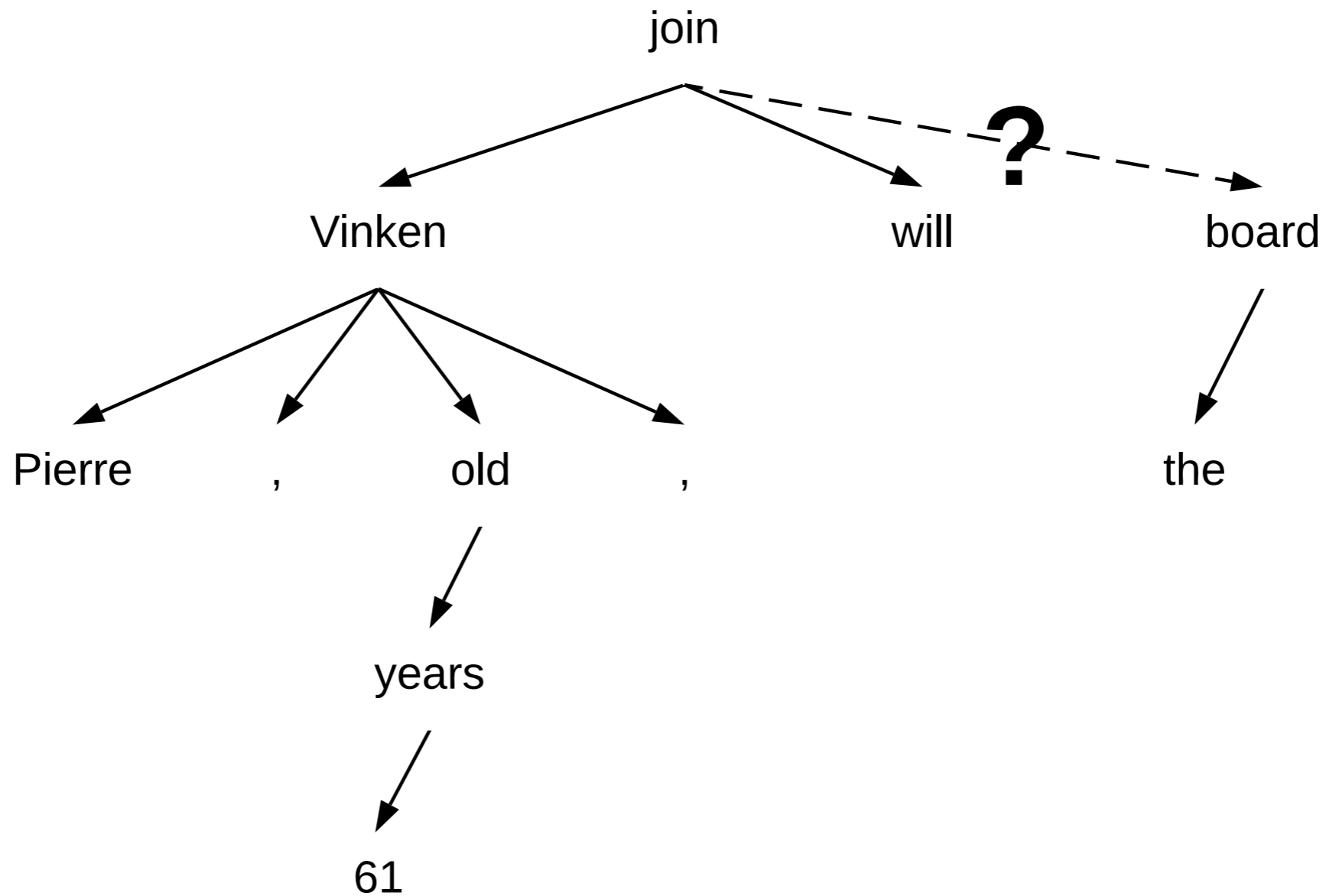JOHNS HOPKINS
UNIVERSITY

# Outline

- What is StackLSTM?

- Parallelization Problem

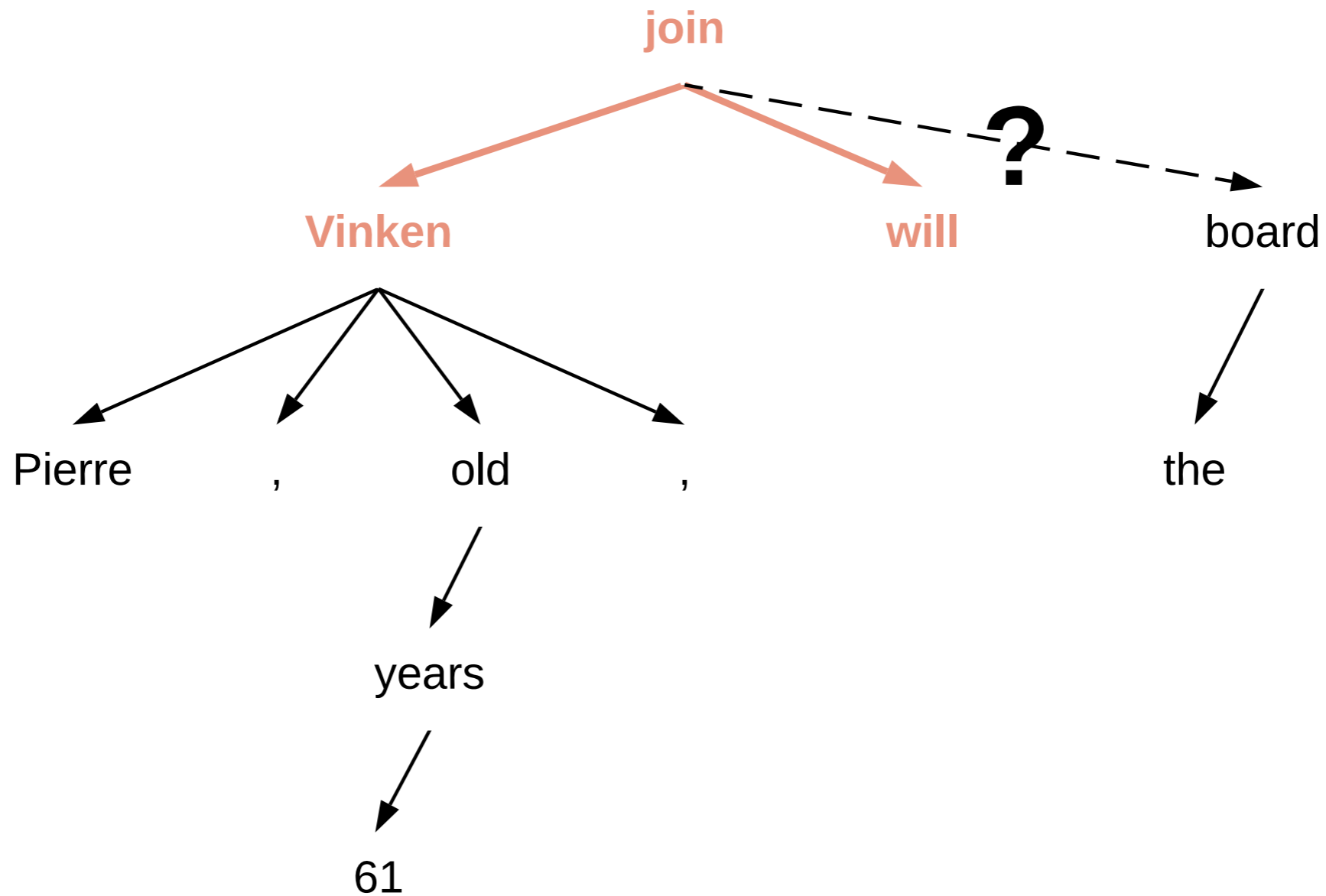- Homogenizing Computation

- Experiments
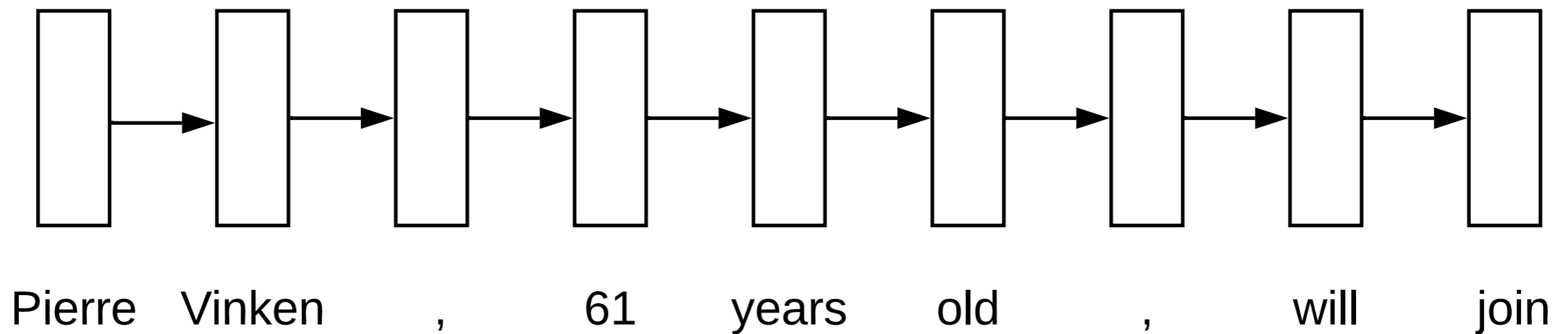
# What is StackLSTM?

# A Partial Tree

# Good Edge?

# Good Edge?

# LSTM?



Pierre   Vinken       ,       61     years     old        ,       will      join
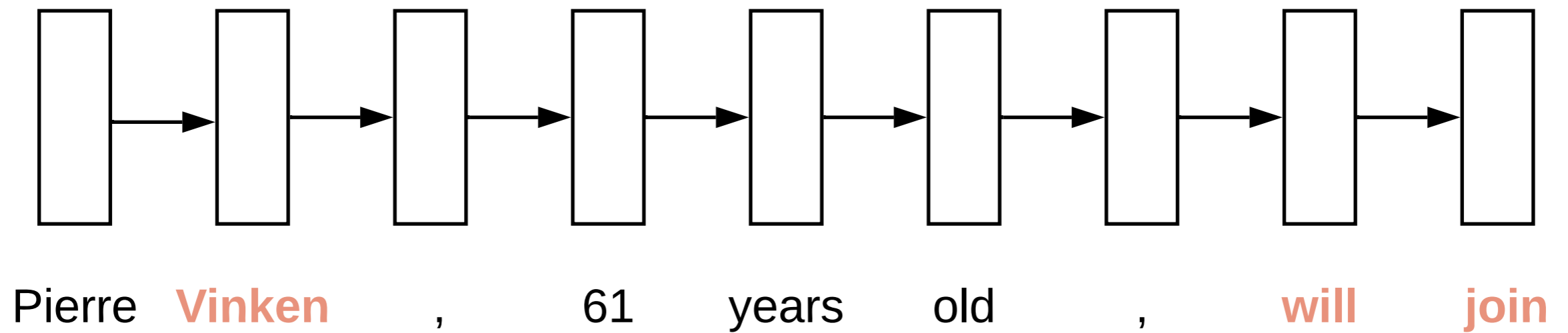
:(



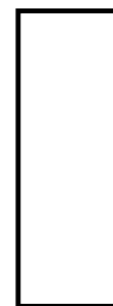Pierre **Vinken** , 61 years old , **will** **join**

# StackLSTM

- An LSTM whose states are stored in a stack

- Computation is conditioned on the stack operation

*Dyer et al. (2015)*
*Ballesteros et al. (2017)*
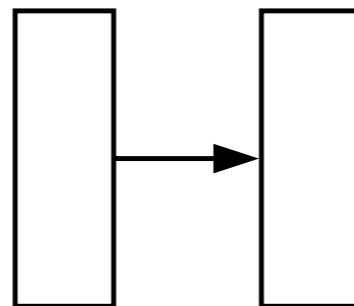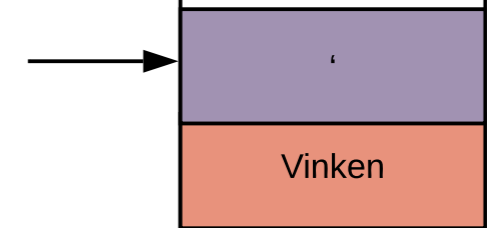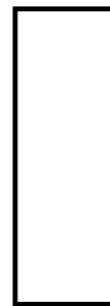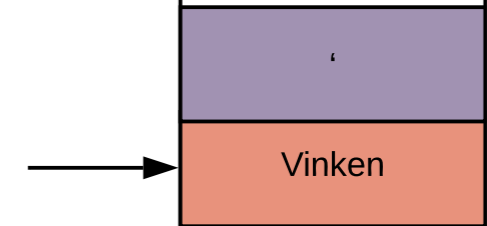
# StackLSTM

Vinken

Vinken

# Push ,

Vinken

Vinken        ,

JOHNS HOPKINS
UNIVERSITY

# Pop

Vinken

Vinken        ,

# Push 61



Vinken    ,    61

# Push years



Vinken , 61 years

# Push old

# Pop

old

years

61

Vinken

Vinken   ,   61   years   old

# Pop

old

years

61

Vinken

Vinken , 61 years old

JOHNS HOPKINS
UNIVERSITY

# Pop

old

years

61

Vinken

Vinken      ,      61     years     old

# Push ,



old

years

,

Vinken

Vinken    ,    61    years    old    ,

# Pop

old

years

,

Vinken

Vinken    ,    61    years    old    ,

# Push will



Vinken    ,    61    years    old    ,    will

# Push join

old

join

will

Vinken

Vinken , 61 years old , will join

JOHNS HOPKINS
UNIVERSITY

# Parallelization Problem

# LSTM

# LSTM

# Batched LSTM

# Batched… StackLSTM?

# Wouldn't it be nice if...

# Homogenizing Computation

# Push

- read the stack top
  hidden state h_{p(t)};

- perform LSTM forward
  computation with x(t)
  and h_{p(t)};

- write new hidden state
  to h_{p(t) + 1};

- update stack top pointer
  p(t+1) = p(t) + 1;

# Push

- **read the stack top hidden state h_{p(t)};**

- perform LSTM forward computation with x(t) and h_{p(t)};

- write new hidden state to h_{p(t) + 1};

- update stack top pointer p(t+1) = p(t) + 1;

Vinken

# Push

- read the stack top hidden state $h_{p(t)}$;

- **perform LSTM forward computation with x(t) and $h_{p(t)}$;**

- write new hidden state to $h_{p(t) + 1}$;

- update stack top pointer $p(t+1) = p(t) + 1$;

LSTM Op

Vinken

# Push

- read the stack top hidden state h_{p(t)};

- perform LSTM forward computation with x(t) and h_{p(t)};

- **write new hidden state to h_{p(t) + 1};**

- update stack top pointer p(t+1) = p(t) + 1;

Vinken

# Push

- read the stack top hidden state h_{p(t)};

- perform LSTM forward computation with x(t) and h_{p(t)};

- write new hidden state to h_{p(t) + 1};

- **update stack top pointer p(t+1) = p(t) + 1;**

# Pop

- update stack top pointer
  `p(t+1) = p(t) - 1;`

# Pop

- **update stack top pointer**
  **p(t+1) = p(t) - 1;**

Vinken

# Observation 1

- read the stack top
  hidden state h_{p(t)};

- perform LSTM forward
  computation with x(t)
  and h_{p(t)};

- write new hidden state
  to h_{p(t) + 1};

- update stack top pointer
  p(t+1) = p(t) + 1;

- update stack top pointer
  p(t+1) = p(t) - 1;

# Observation 1

- read the stack top hidden state h_{p(t)};

- perform LSTM forward computation with x(t) and h_{p(t)};

- write new hidden state to h_{p(t) + 1};

- update stack top pointer p(t+1) = p(t) **+ op**;

**Use op = +1 for push and op = -1 for pop**

- update stack top pointer p(t+1) = p(t) **+ op**;

# Observation 1

The computation performed for Pop operation is a subset of Push operation.

# Observation 2

Is it safe to do the other computations for push for pop as well?

JOHNS HOPKINS
UNIVERSITY

# Observation 2

- read the stack top
  hidden state h_{p(t)};

- perform LSTM forward
  computation with x(t)
  and h_{p(t)};

- write new hidden state
  to h_{p(t) + 1};

- update stack top pointer
  p(t+1) = p(t) + op;

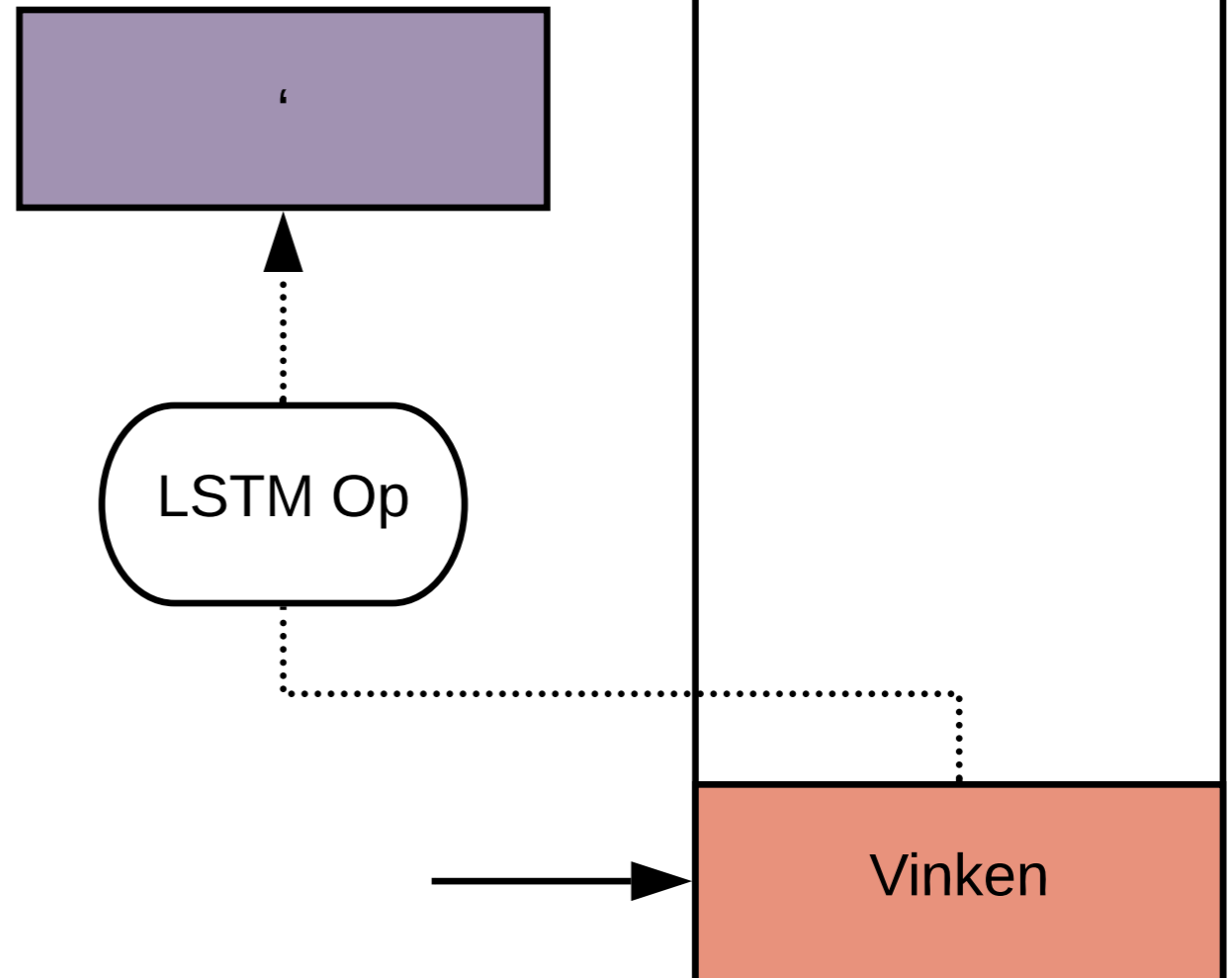- update stack top pointer
  p(t+1) = p(t) + op;

# Observation 2

- read the stack top
  hidden state h_{p(t)};

- perform LSTM forward
  computation with x(t)
  and h_{p(t)};

- **write new hidden state
  to h_{p(t) + 1};**
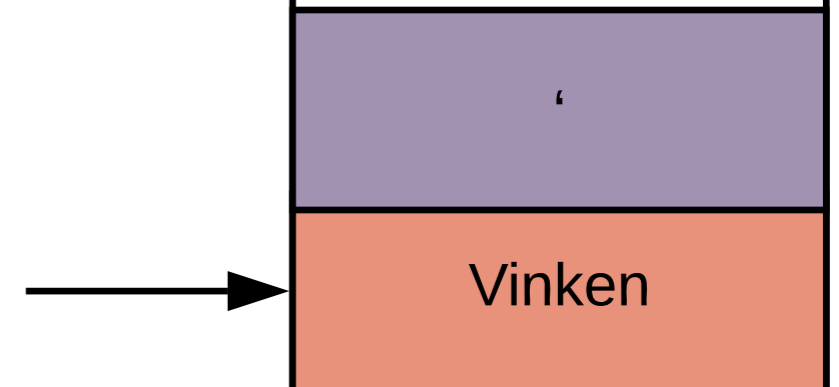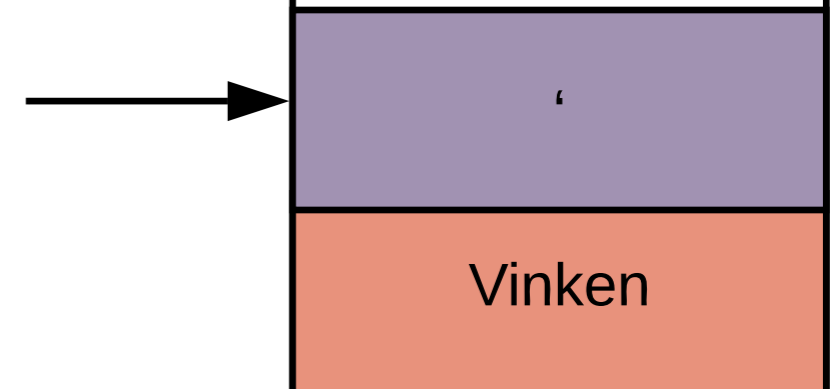
- update stack top pointer
  p(t+1) = p(t) + op;

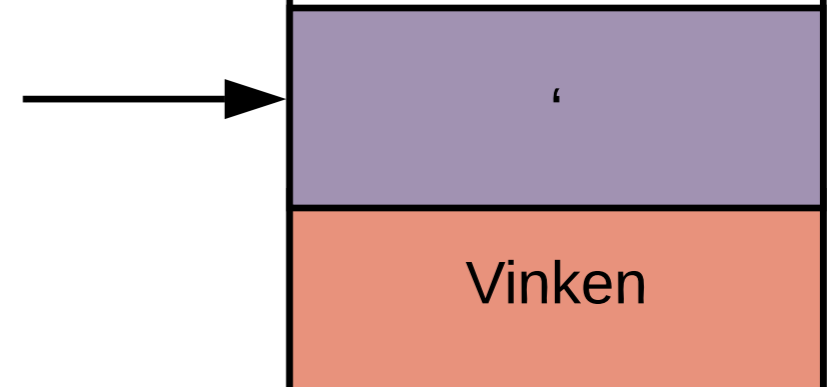- update stack top pointer
  p(t+1) = p(t) + op;

# Observation 2

A write will always happen before the stack top pointer advances.

# Observation 2

If one wants to write anything in the higher position than the current stack top pointer...

JOHNS HOPKINS
UNIVERSITY

# Observation 2

If one wants to write anything in the higher position than the current stack top pointer...

**Just do it!**

# Observation 2

- read the stack top
  hidden state h_{p(t)};

- perform LSTM forward
  computation with x(t)
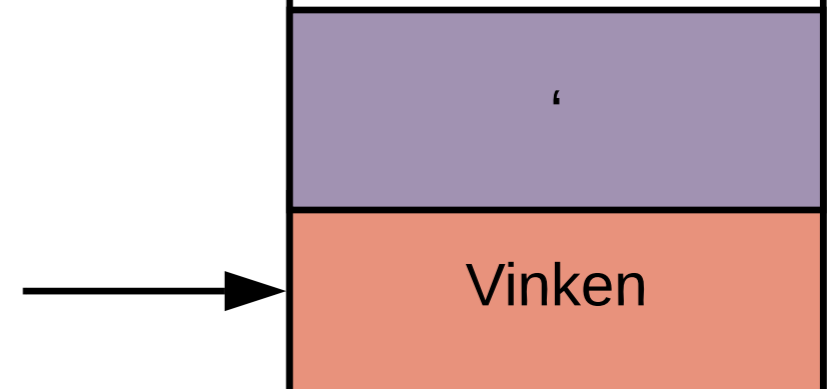  and h_{p(t)};

- write new hidden state
  to h_{p(t) + 1};

- update stack top pointer
  p(t+1) = p(t) + op;

- update stack top pointer
  p(t+1) = p(t) + op;

# Observation 2

- read the stack top hidden state h_{p(t)};

- perform LSTM forward computation with x(t) and h_{p(t)};

- write new hidden state to h_{p(t) + 1};

- update stack top pointer p(t+1) = p(t) + op;

- read the stack top hidden state h_{p(t)};

- perform LSTM forward computation with x(t) and h_{p(t)};

- write new hidden state to h_{p(t) + 1};

- update stack top pointer p(t+1) = p(t) + op;

# Done!

- read the stack top hidden state h_{p(t)};

- perform LSTM forward computation with x(t) and h_{p(t)};

- write new hidden state to h_{p(t) + 1};

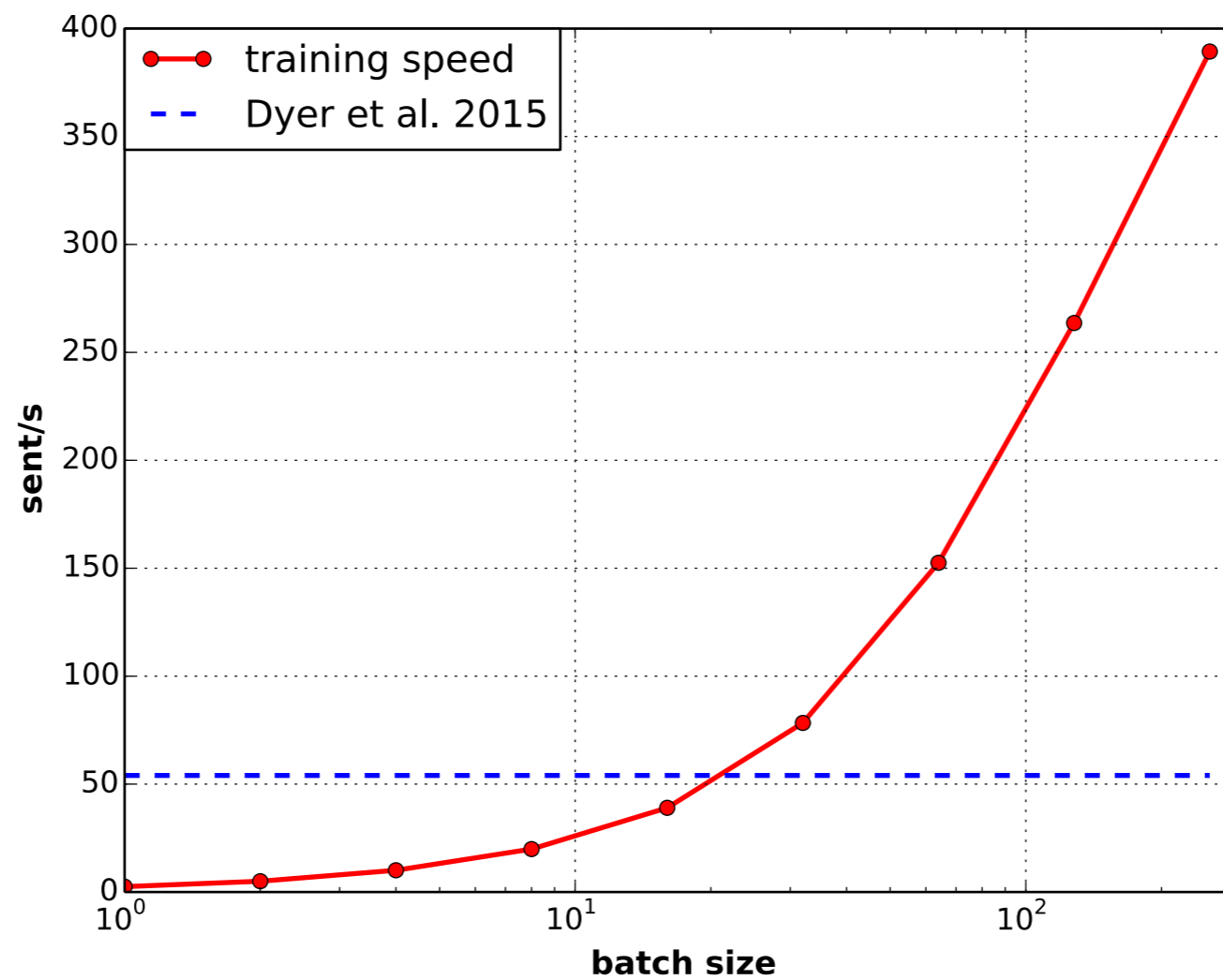- update stack top pointer p(t+1) = p(t) + op;

# Experiments

# Benchmark

Transition-based dependency parsing on Stanford Dependency Treebank

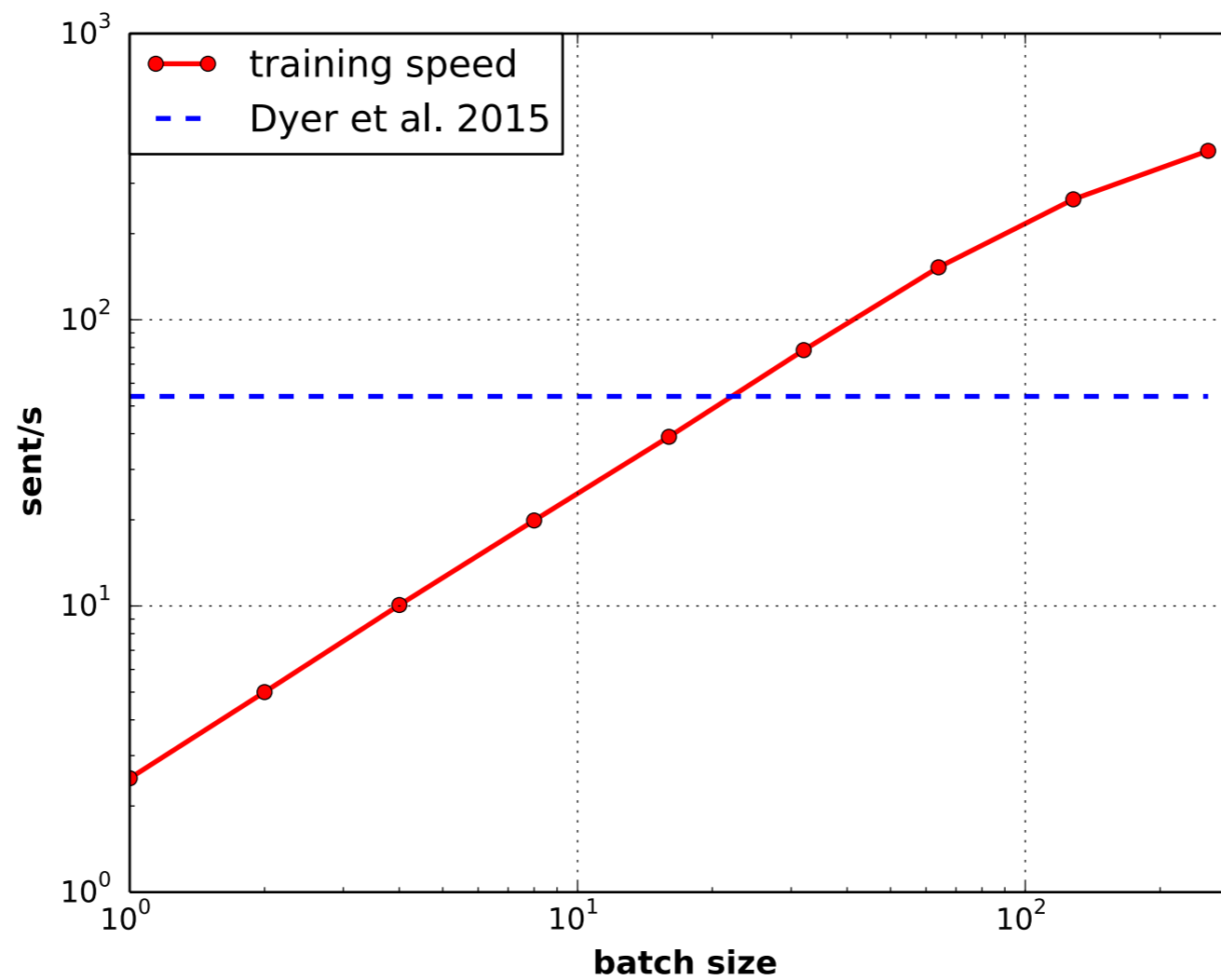PyTorch, Single K80 GPU

JOHNS HOPKINS
UNIVERSITY

# Hyperparameters

- Largely following Dyer et al. (2015); Ballesteros et al. (2017), except:

  - Adam w/ ReduceLROnPlateau and warmup

  - Arc-Hybrid w/o composition function

  - Slightly larger models (200 hidden, 200 state, 48 action embedding) perform better
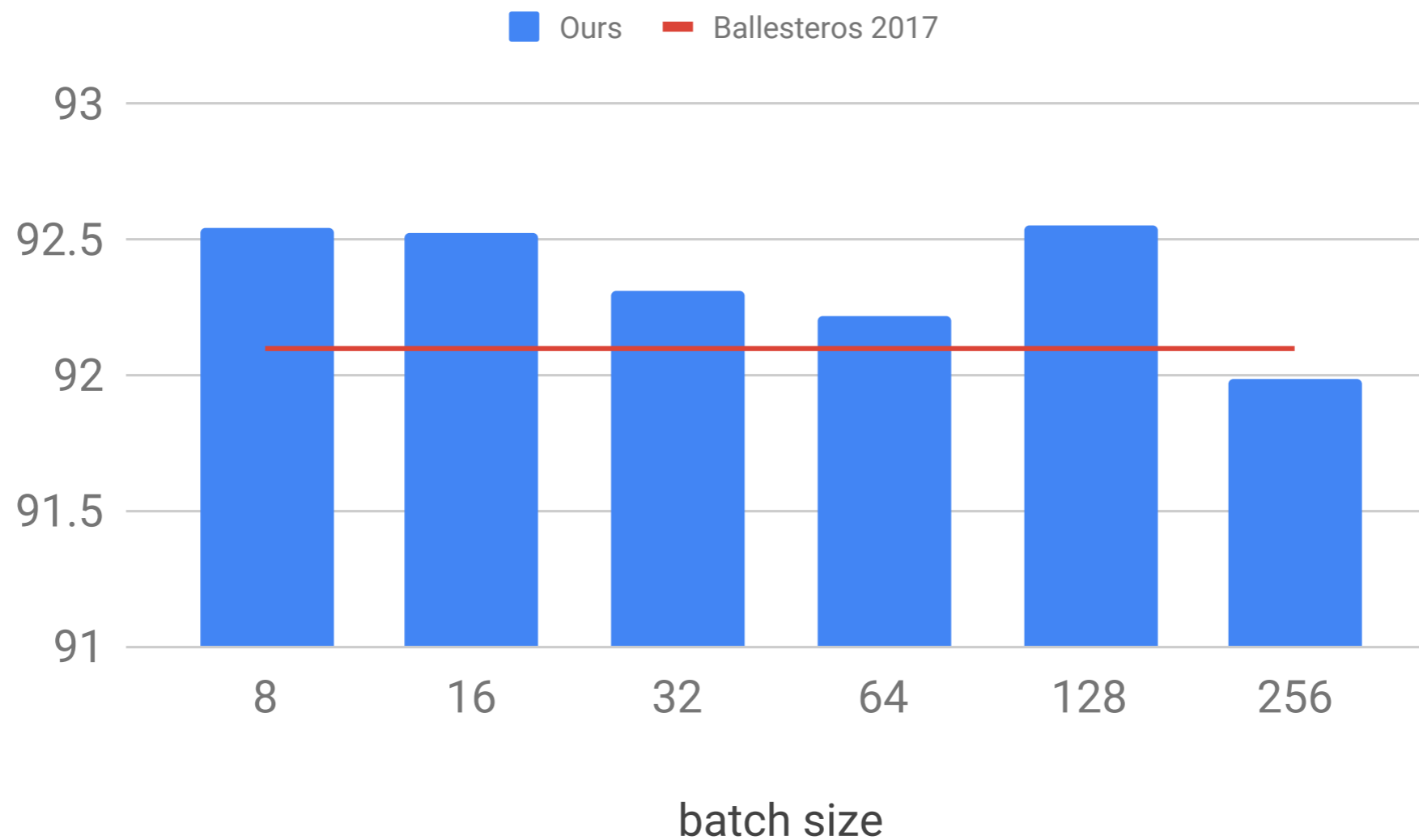
JOHNS HOPKINS
U N I V E R S I T Y

# Speed

# Speed

# Performance



Legend: Ours (blue), Ballesteros 2017 (red line)

X-axis (batch size): 8, 16, 32, 64, 128, 256
Y-axis: 91, 91.5, 92, 92.5, 93

JOHNS HOPKINS
U N I V E R S I T Y

# Conclusion

# Conclusion

- We propose a parallelization scheme for StackLSTM architecture.

- Together with a different optimizer, we are able to train parsers of comparable performance within 1 hour.

paper
code
slides

https://github.com/shuoyangd/hoolock

JOHNS HOPKINS
UNIVERSITY