

Fast(er) **Exact Decoding** and **Global Training** for Transition-Based Dependency Parsing via a Minimal Feature Set

Tianze Shi*

Liang Huang†

Lillian Lee*



* Cornell
University

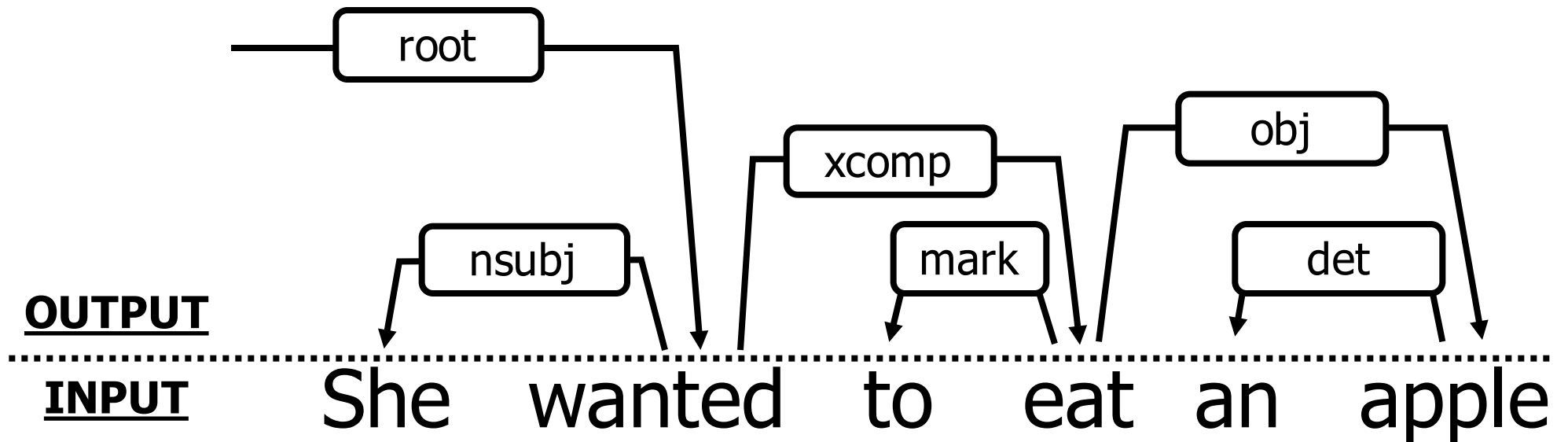


† Oregon State
University

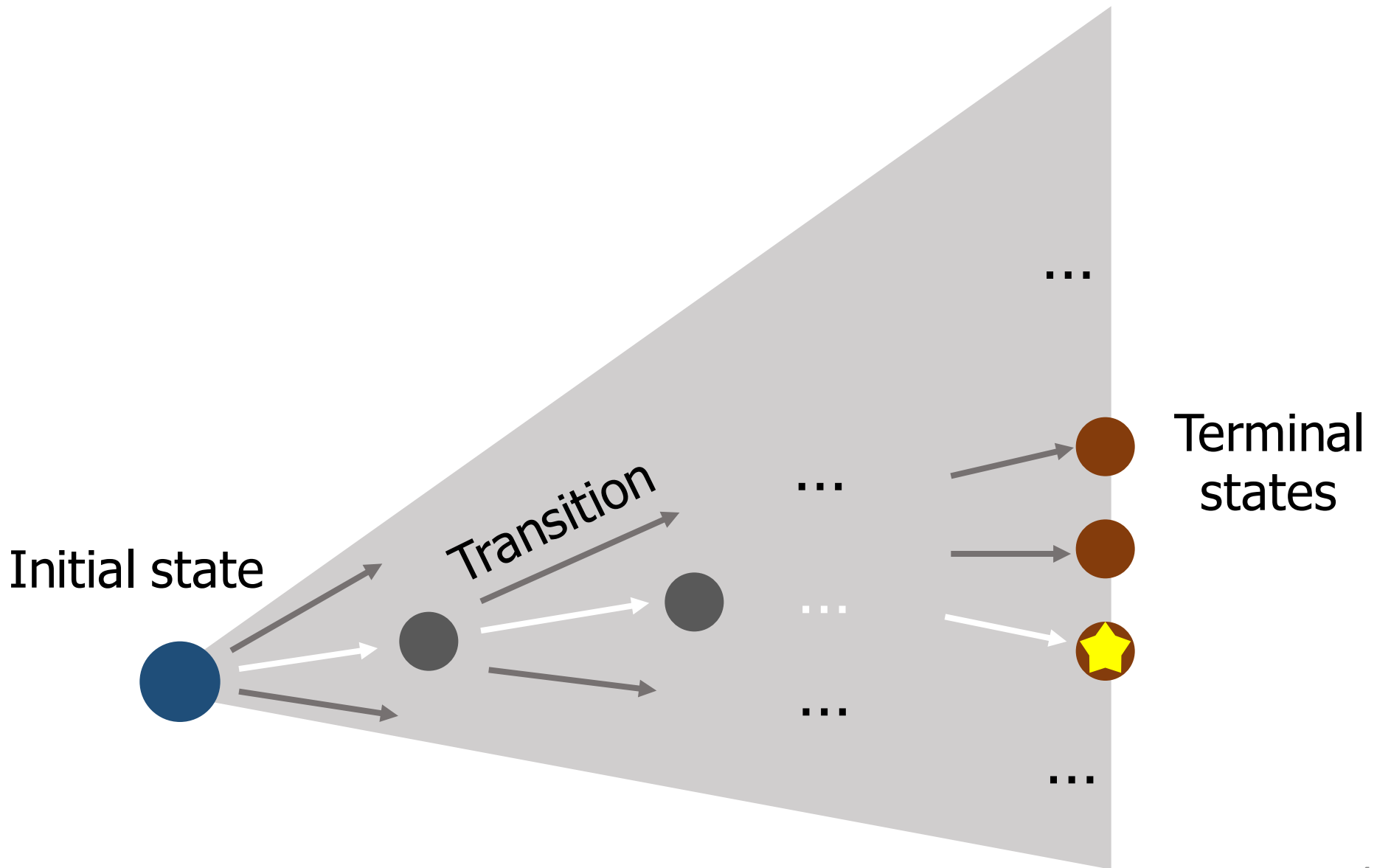
Short Version

- Transition-based dependency parsing has an exponentially-large search space
- $O(n^3)$ exact solutions exist
- In practice, however, we needed rich features $\implies O(n^6)$
- (This work) with bi-LSTMs, now we can do $O(n^3)$!
- And we get state-of-the-art results

Dependency Parsing



Transition-based Dependency Parsing



Transition-based Dependency Parsing

Goal:

$\max \text{score}(\bullet \rightarrow \bullet \rightarrow \bullet \rightarrow \dots \rightarrow \bullet)$

$= \max \sum \text{score}(\bullet \rightarrow \bullet)$

Initial state



Transition



...

...

...



Terminal
states

Exact Decoding

- Dynamic programming (Huang and Sagae, 2010; Kuhlmann, Gómez-Rodríguez and Satta, 2011)

... since transition (sub-)sequences can be reused

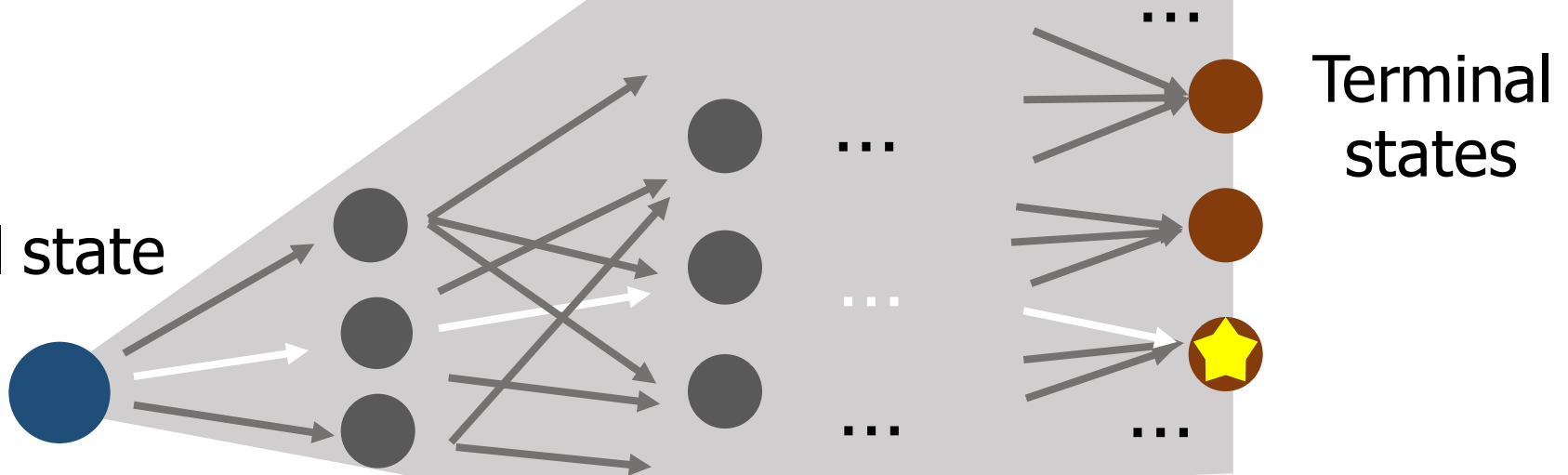
Exact Decoding

Goal:

$\max \text{score}(\bullet \rightarrow \bullet \rightarrow \bullet \rightarrow \dots \rightarrow \bullet)$

$= \max \sum \text{score}(\bullet \rightarrow \bullet)$

Initial state



Exponential to polynomial

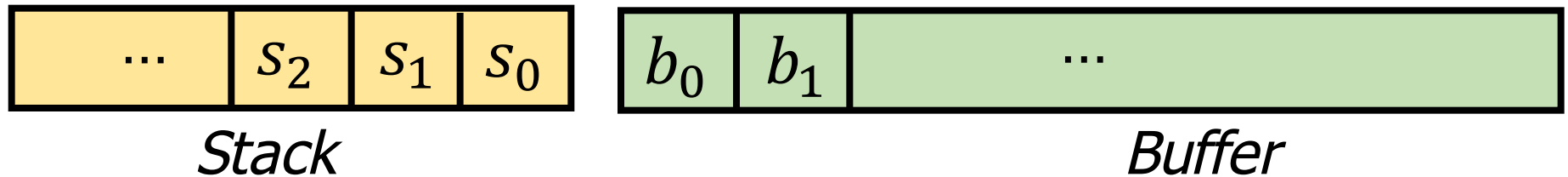
Transition Systems

	DP Complexity	# Transitions	
Arc-standard	$O(n^4)$	3	} In our paper
Arc-eager	$O(n^3)$	4	
<u>Arc-hybrid</u>	$O(n^3)$	3	

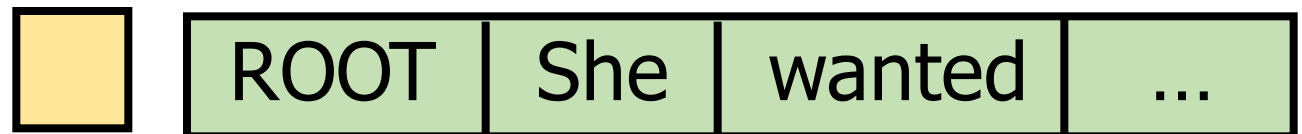
Presentational convenience

Arc-hybrid Transition System

Search State



Initial State



Terminal State

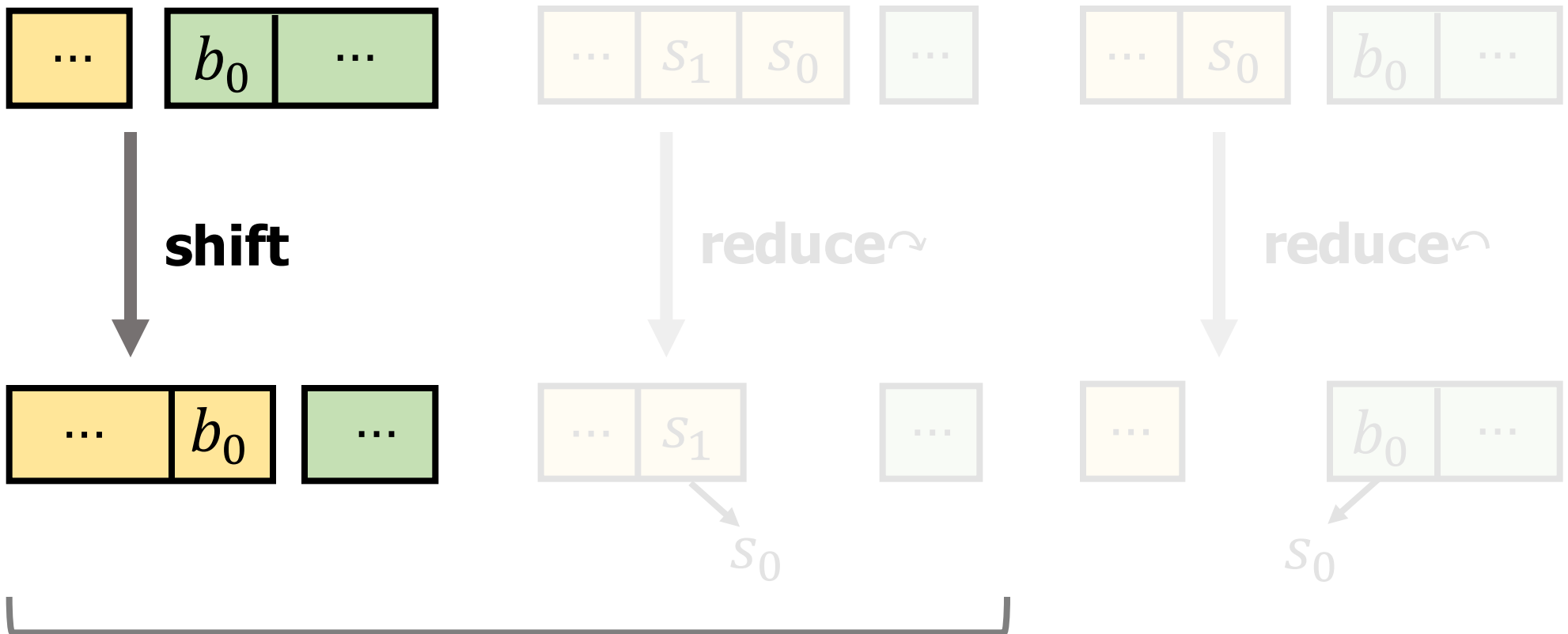


(Yamada and Matsumoto, 2003)

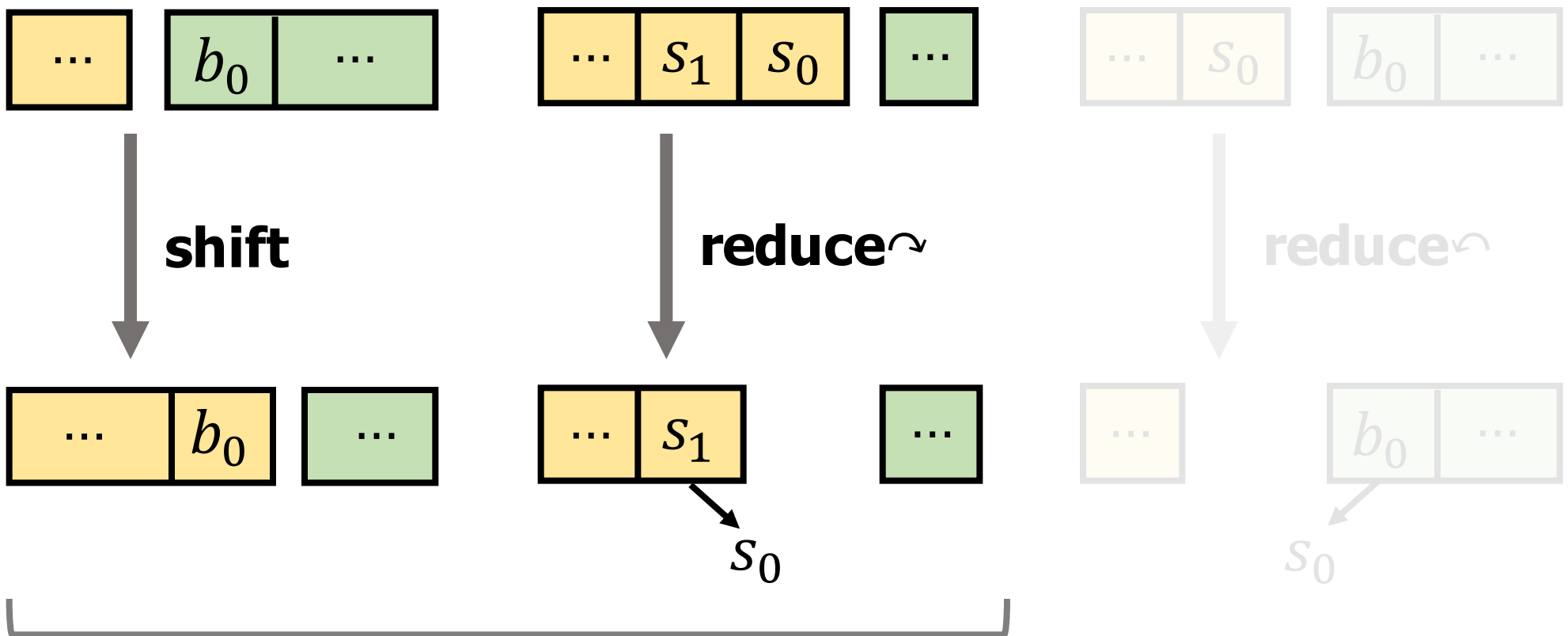
(Gómez-Rodríguez et al., 2008)

(Kuhlmann et al., 2011)

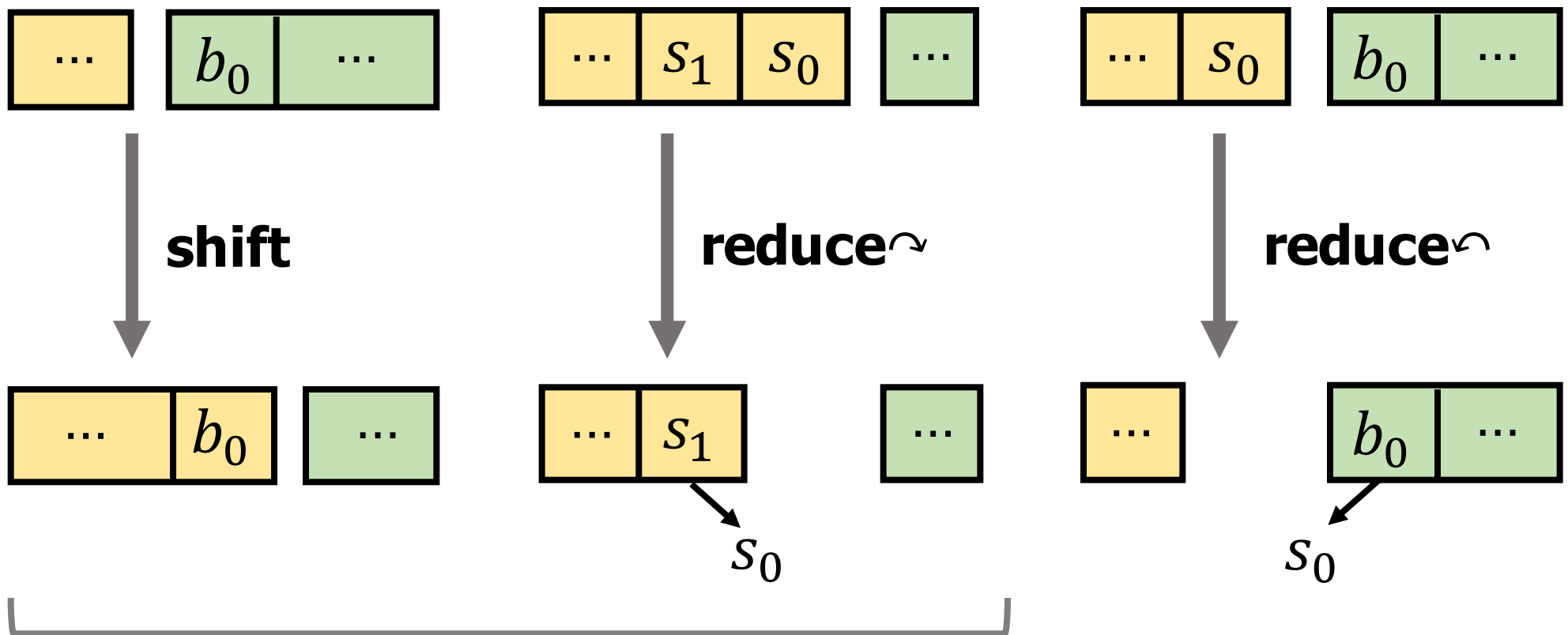
Arc-hybrid Transition System



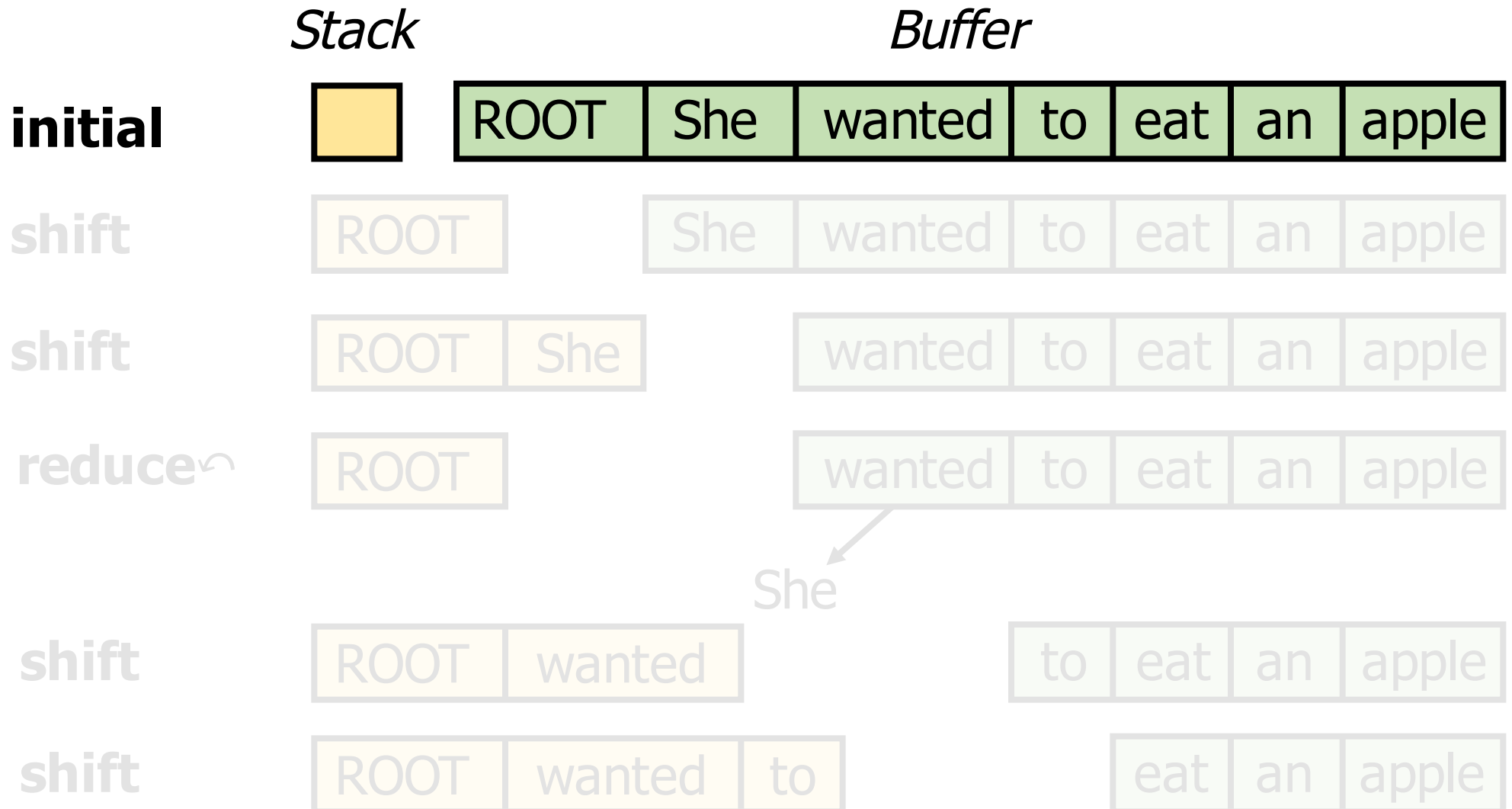
Arc-hybrid Transition System



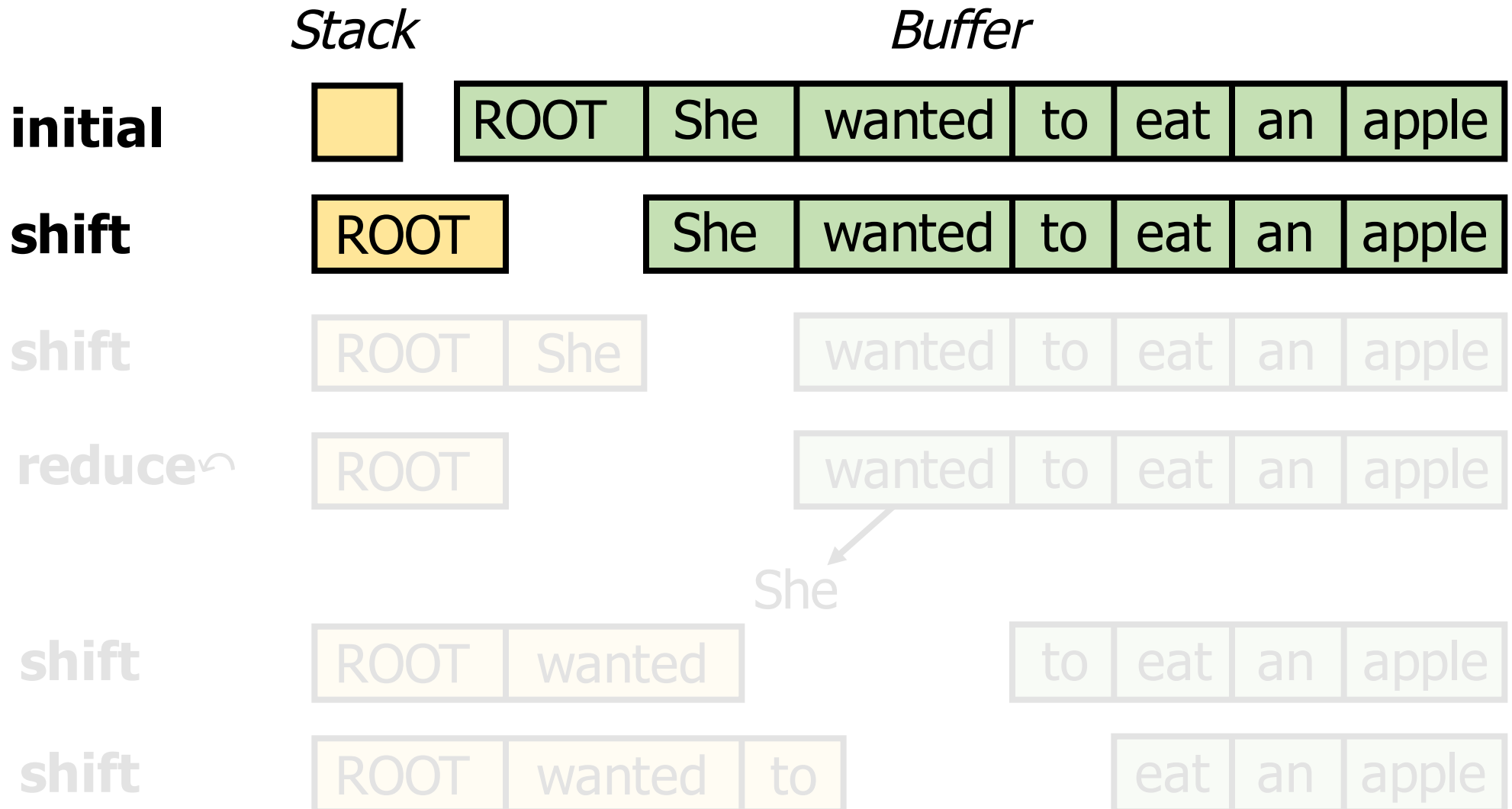
Arc-hybrid Transition System



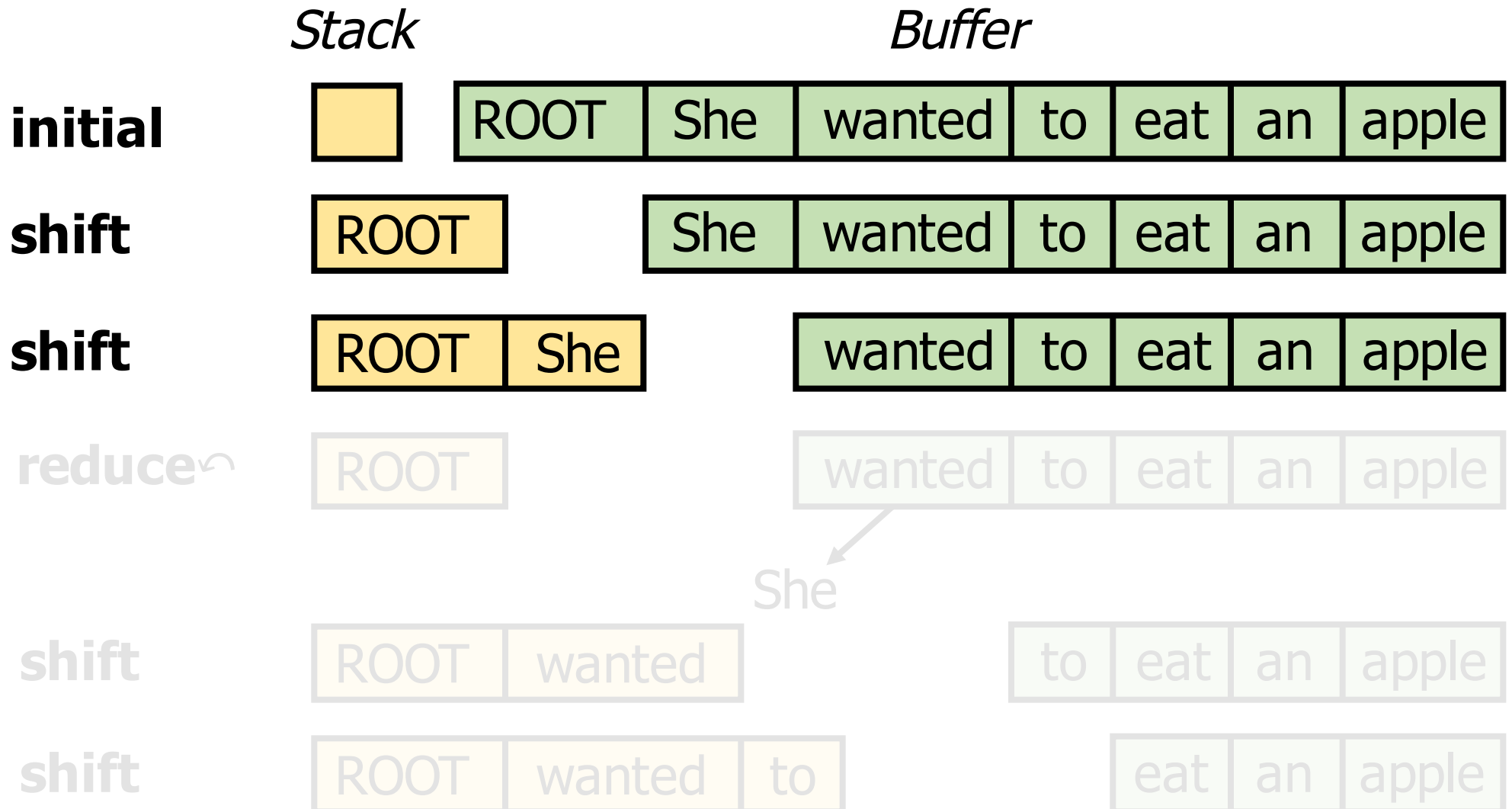
Arc-hybrid Transition System



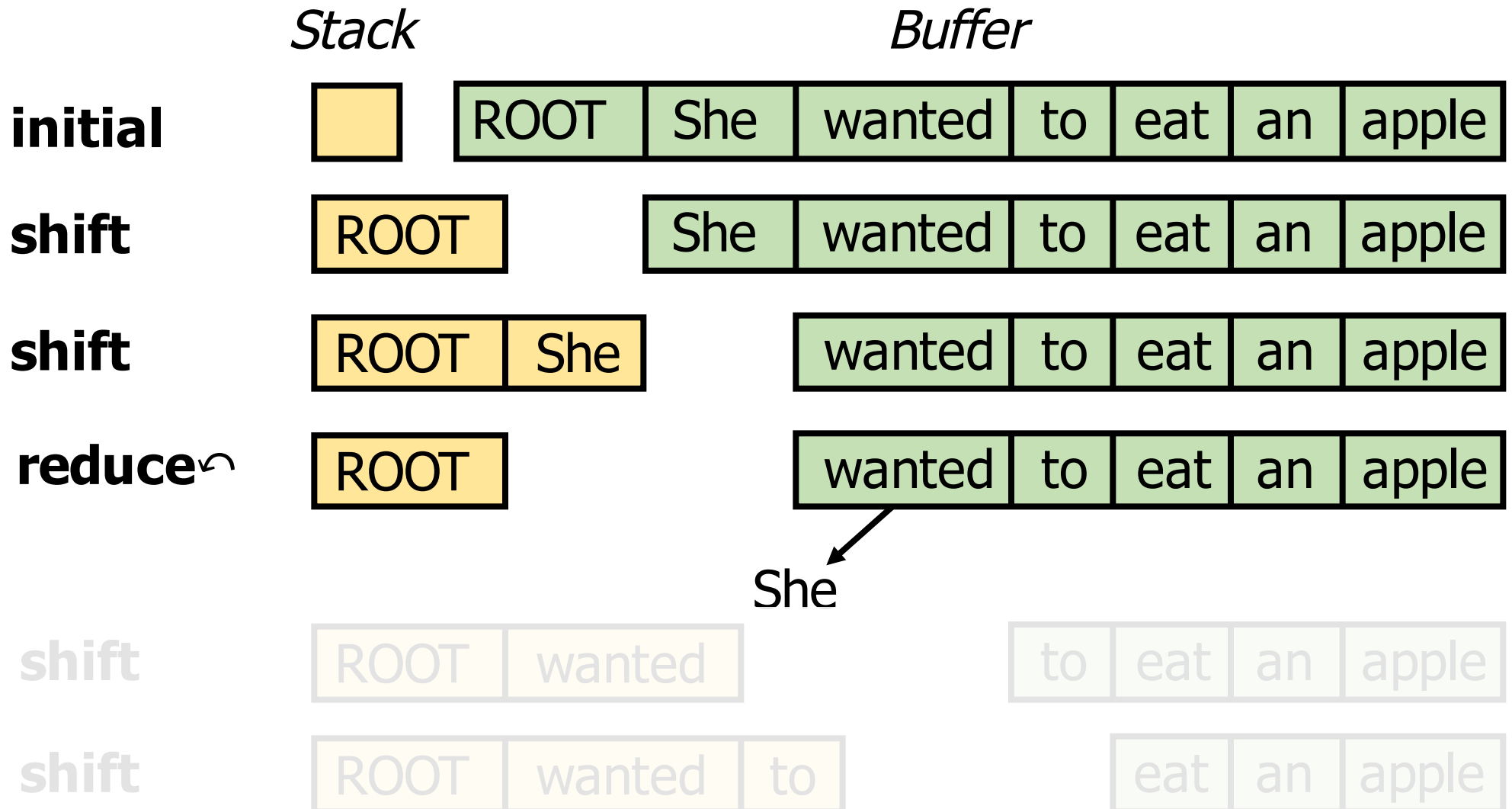
Arc-hybrid Transition System



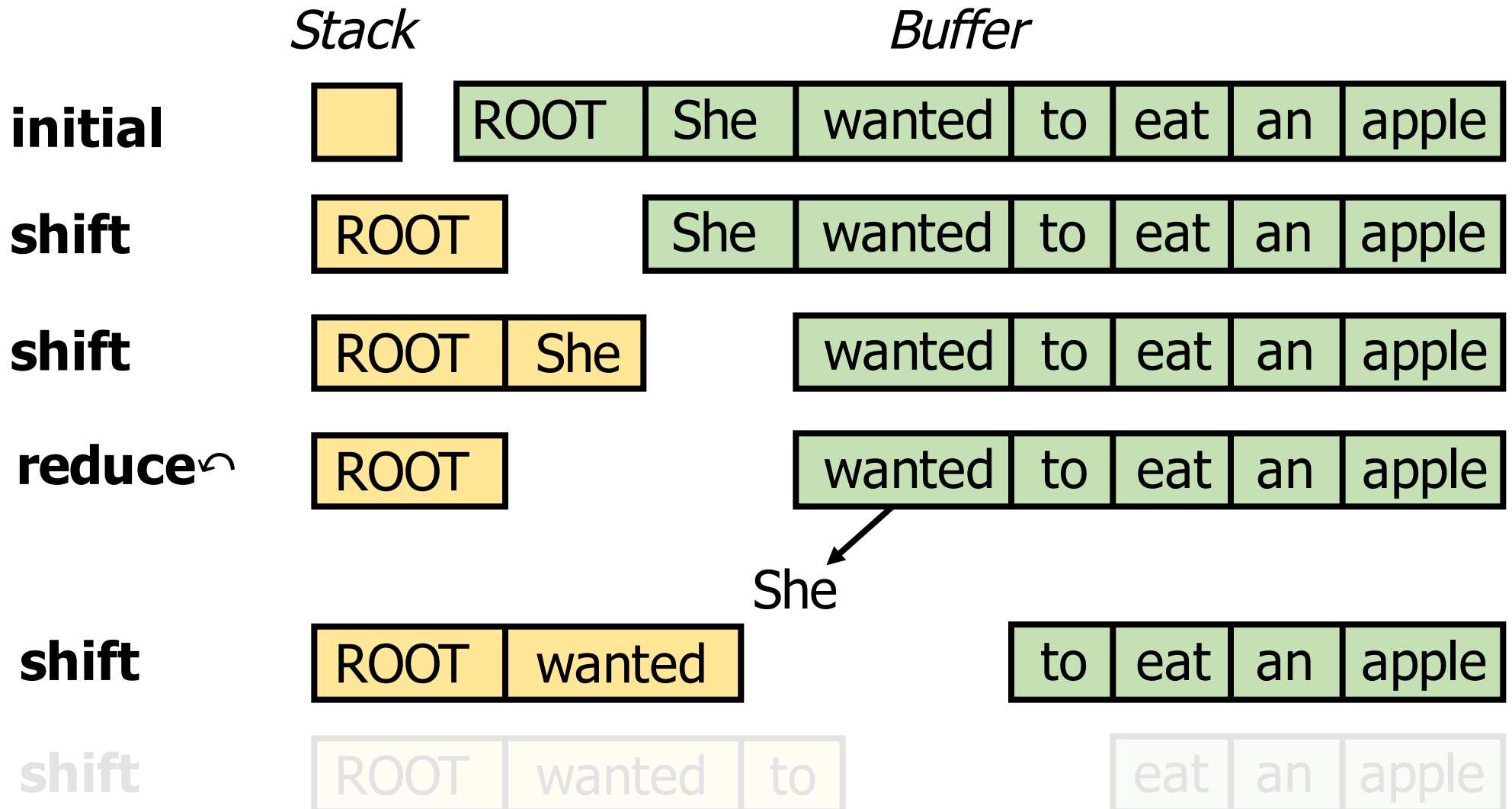
Arc-hybrid Transition System



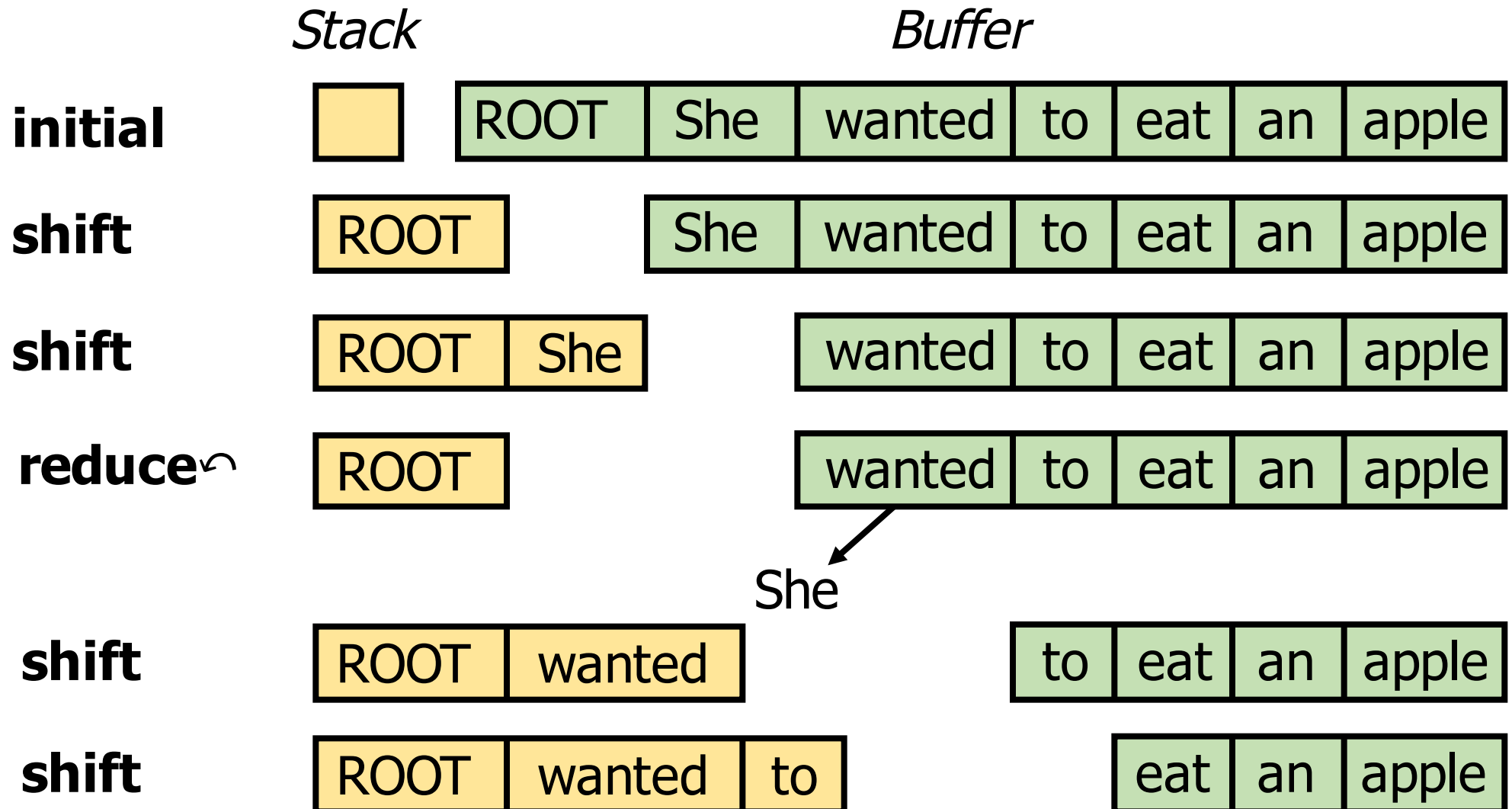
Arc-hybrid Transition System



Arc-hybrid Transition System



Arc-hybrid Transition System

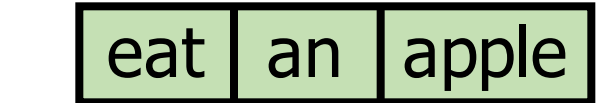
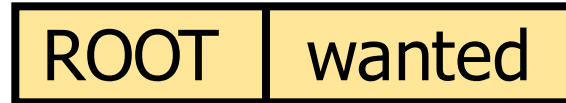


Arc-hybrid Transition System

Stack

Buffer

reduce ↷



to

shift



shift



reduce ↷



an

shift

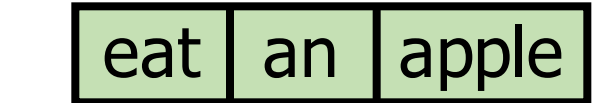
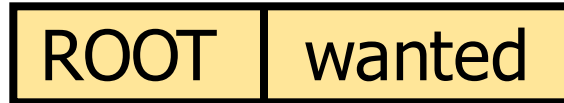


Arc-hybrid Transition System

Stack

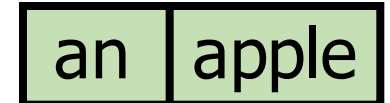
Buffer

reduce ↷



to

shift



shift



reduce ↷



an

shift

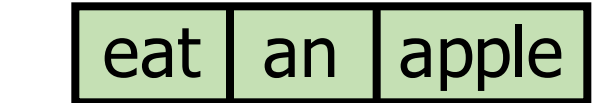
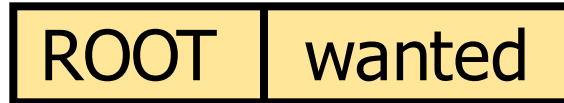


Arc-hybrid Transition System

Stack

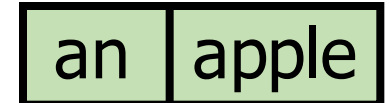
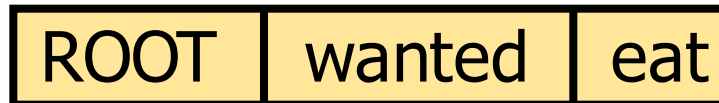
Buffer

reduce ↷

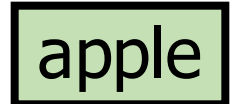


to

shift



shift



reduce ↷



an

shift

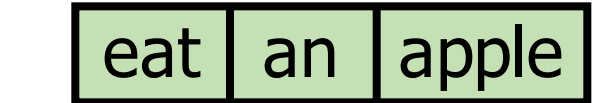
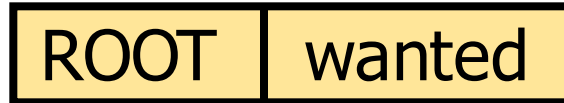


Arc-hybrid Transition System

Stack

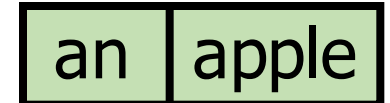
Buffer

reduce ↷



to

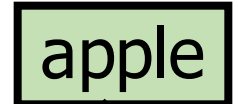
shift



shift



reduce ↷



an

shift

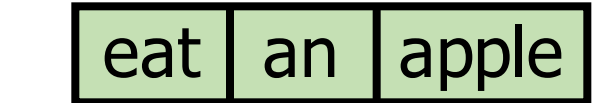
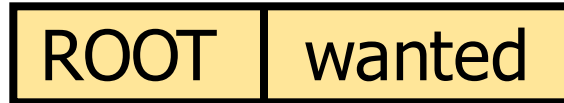


Arc-hybrid Transition System

Stack

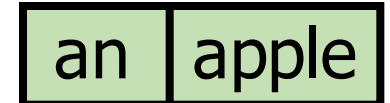
Buffer

reduce ↷

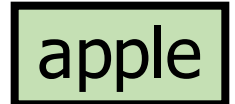


to

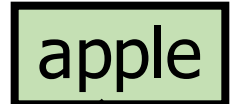
shift



shift



reduce ↷



an

shift



Arc-hybrid Transition System

Stack

Buffer

reduce \rightsquigarrow



apple



reduce \rightsquigarrow



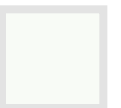
eat



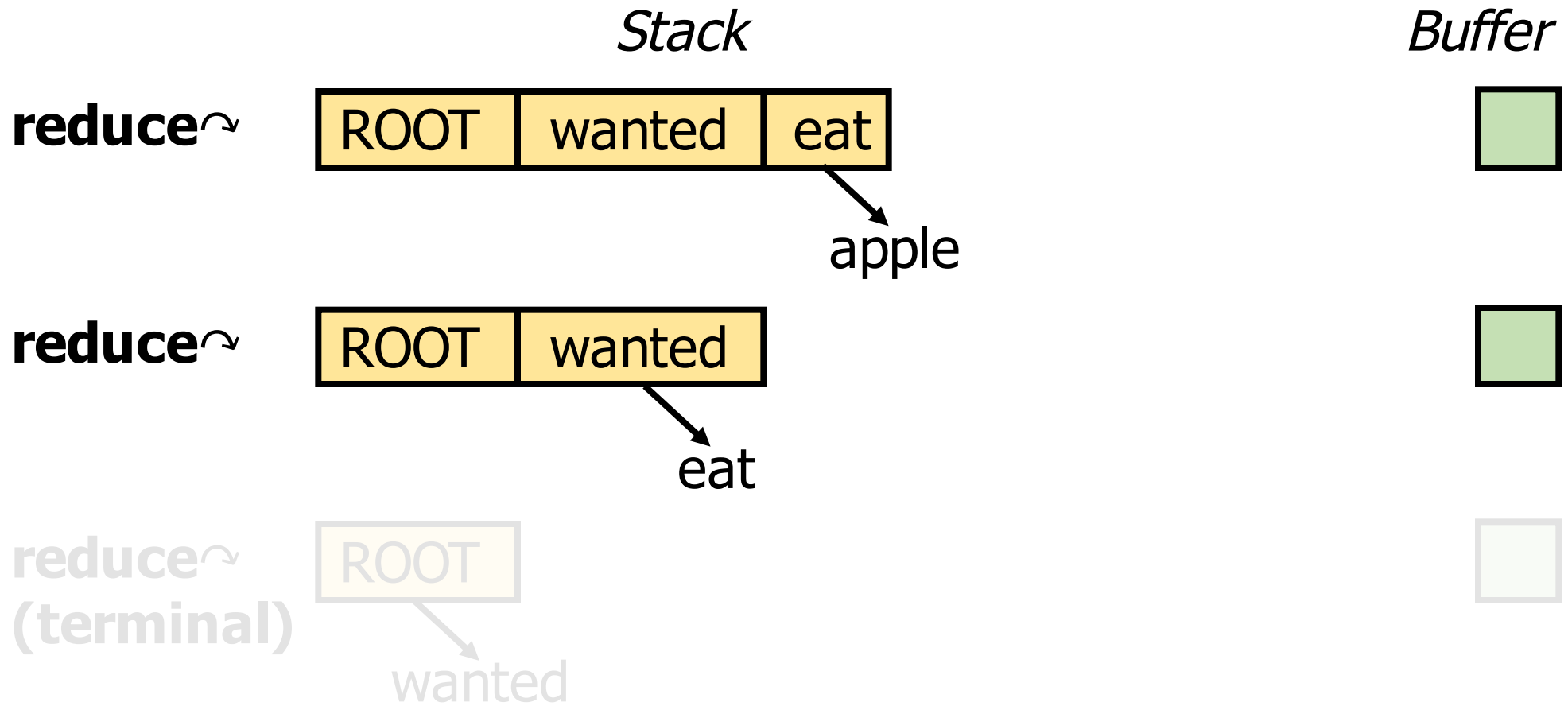
reduce \rightsquigarrow
(terminal)



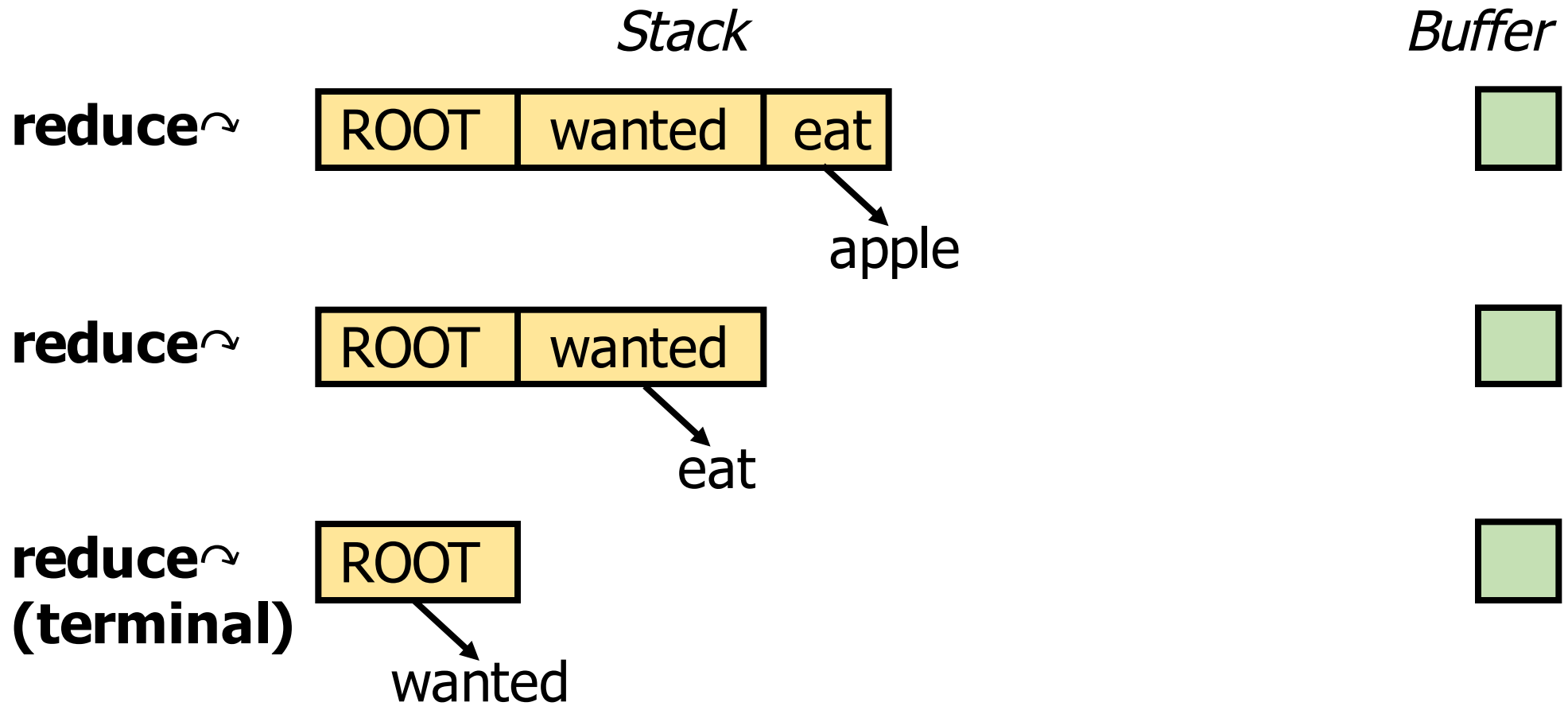
wanted



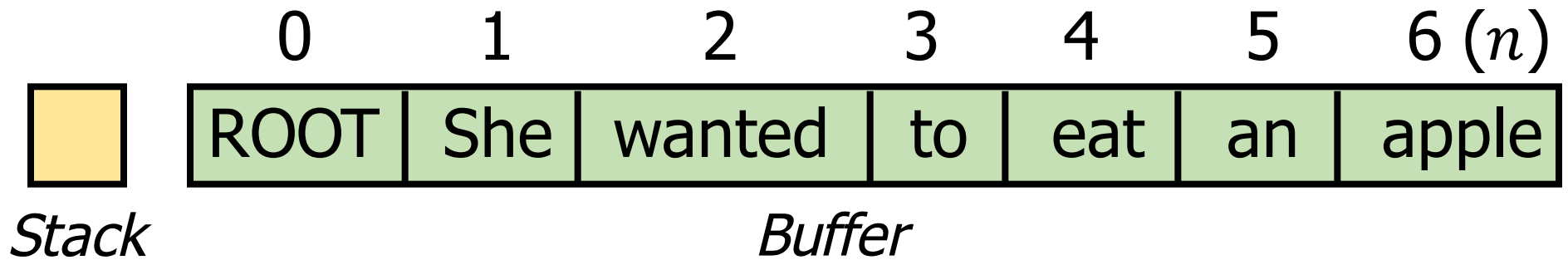
Arc-hybrid Transition System



Arc-hybrid Transition System



Deduction System for Arc-hybrid

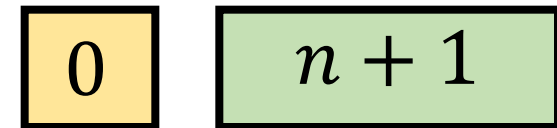


- Deduction Item



$[i, j]$

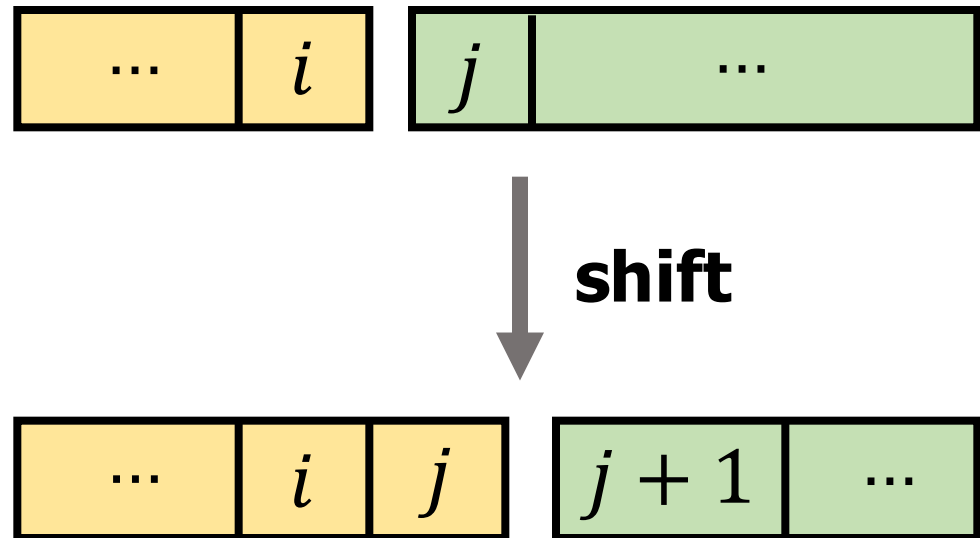
- Goal



$[0, n + 1]$

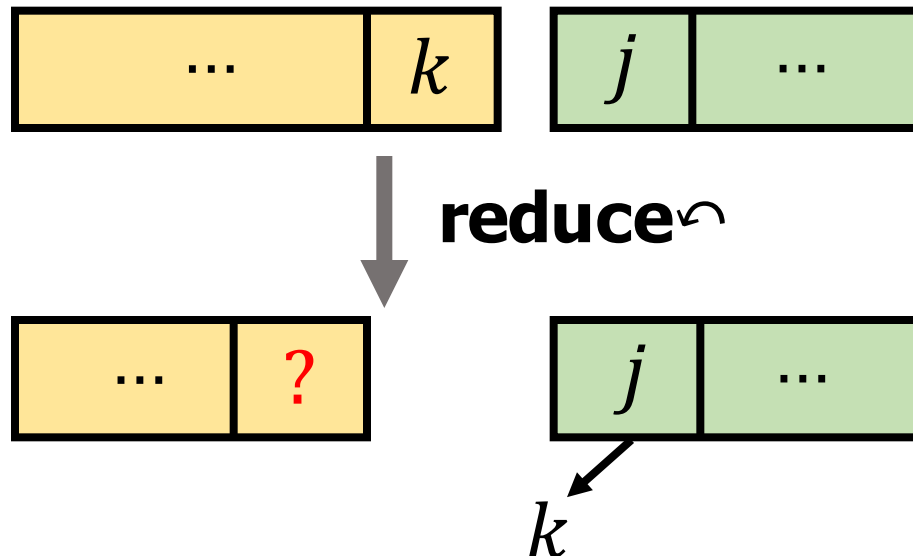
Deduction System for Arc-hybrid

$$\text{shift} \frac{[i, j]}{[j, j + 1]}$$



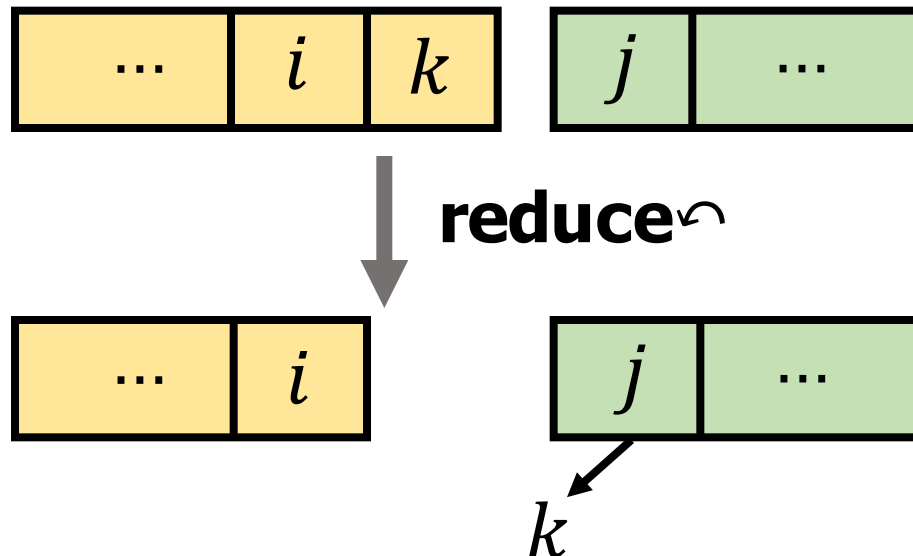
Deduction System for Arc-hybrid

reduce[↷] $\frac{[k, j]}{[?, j]}$

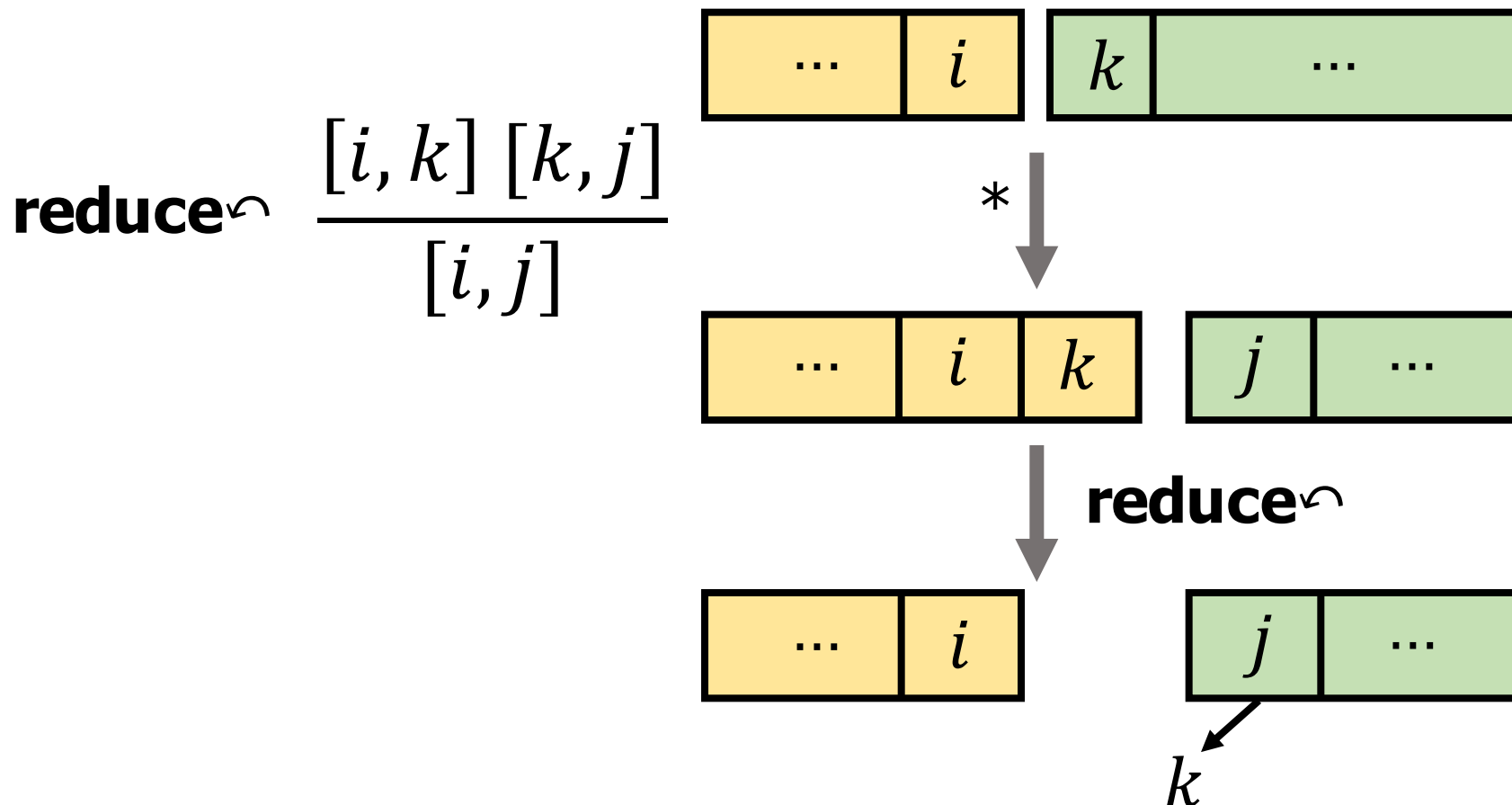


Deduction System for Arc-hybrid

reduce[↷] $\frac{[k, j]}{[i, j]}$



Deduction System for Arc-hybrid

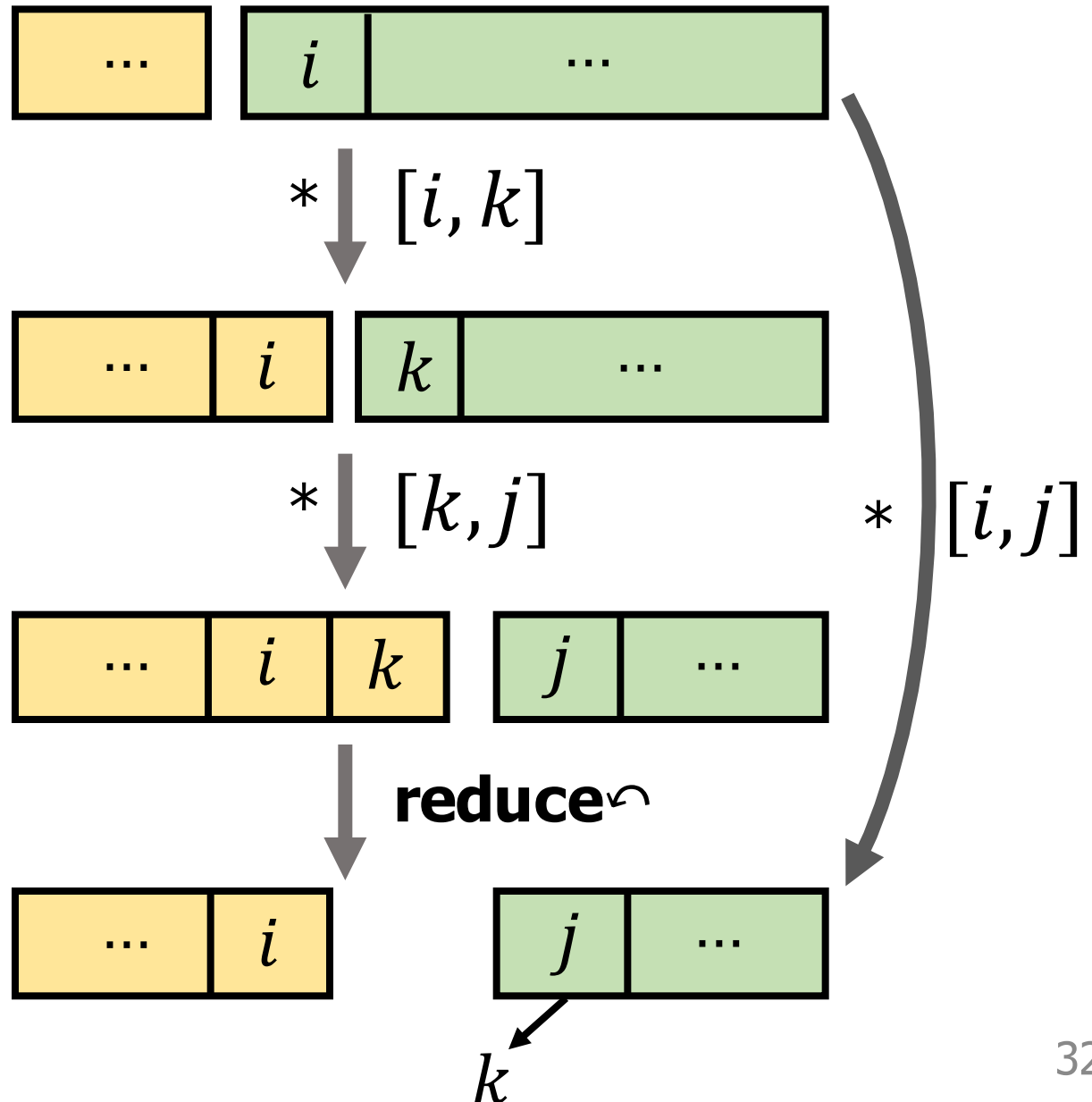


Deduction System for Arc-hybrid

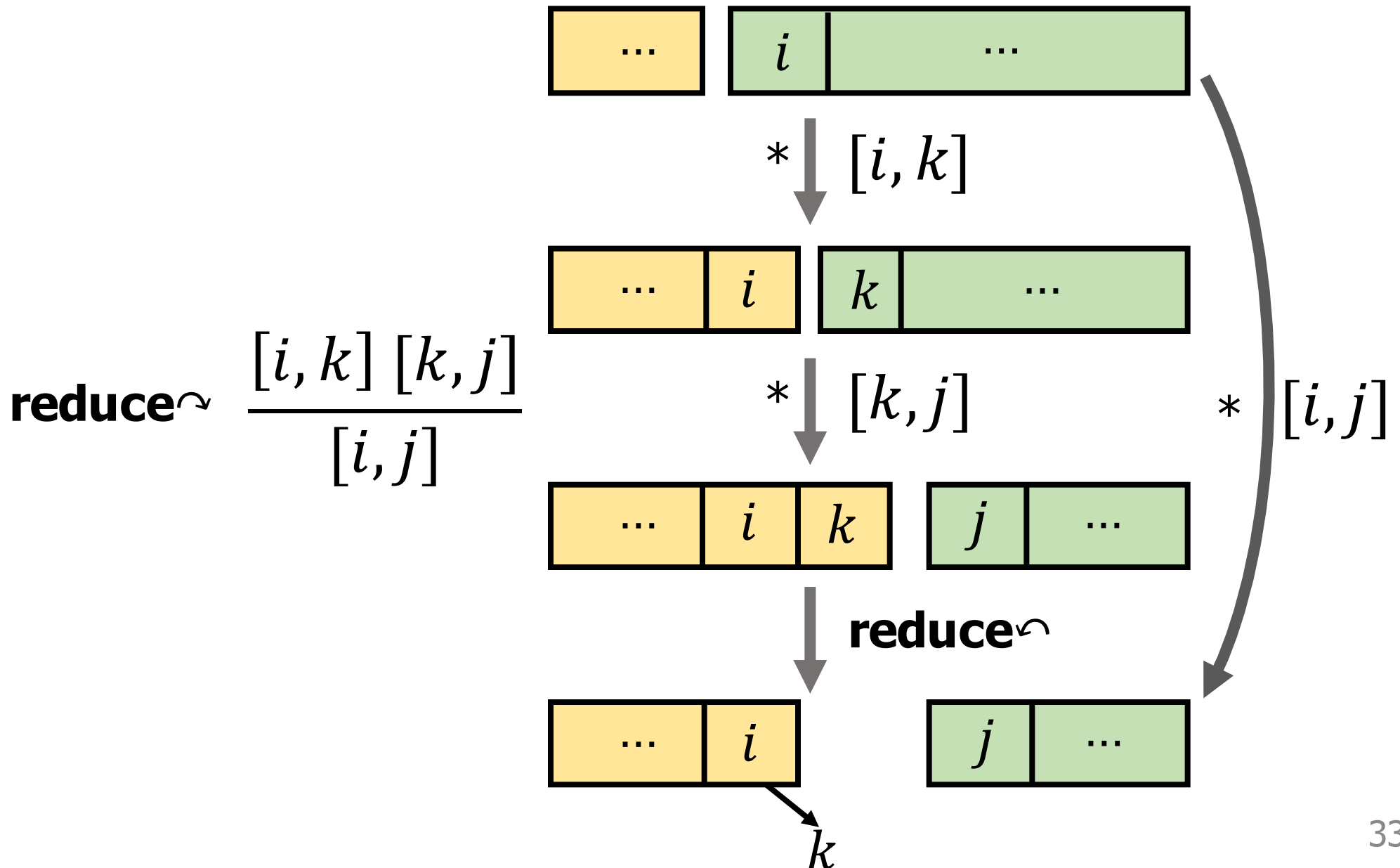
In Kuhlmann et al.
(2011)'s notation

reduce ↷

$$\frac{[i, k] \quad [k, j]}{[i, j]}$$



Deduction System for Arc-hybrid



Deduction System for Arc-hybrid

shift
$$\frac{[i, j]}{[j, j + 1]}$$

Goal: $[0, n + 1]$

reduce \curvearrowright
$$\frac{[i, k] [k, j]}{[i, j]} \quad k \curvearrowright j$$

reduce \curvearrowleft
$$\frac{[i, k] [k, j]}{[i, j]} \quad i \curvearrowleft k$$

$$O(n^3)$$

Time Complexity in Practice

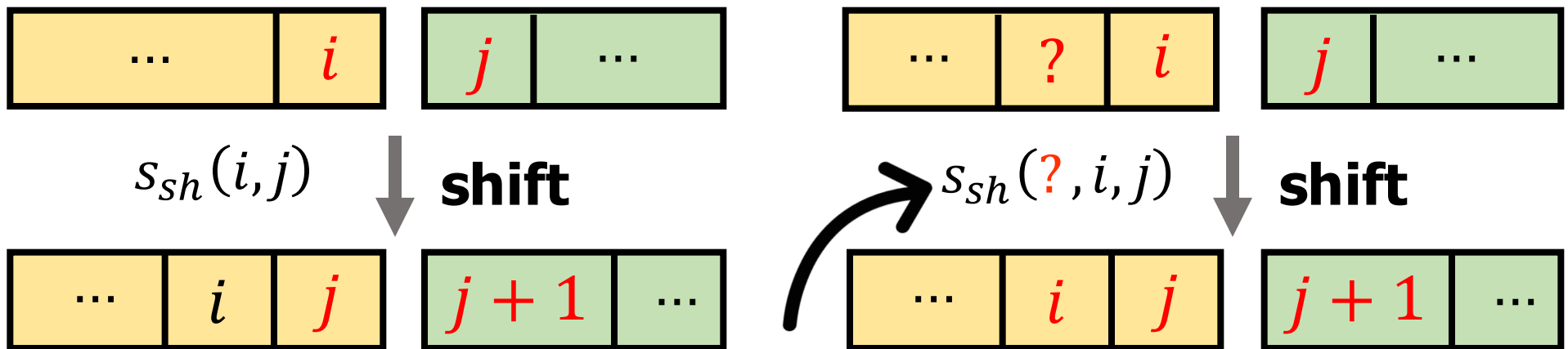
- Complexity depends on feature representation!
- Typical feature representation:
 - Feature templates look at specific *positions* in the stack and in the buffer

Time Complexity in Practice

- Compare the following features



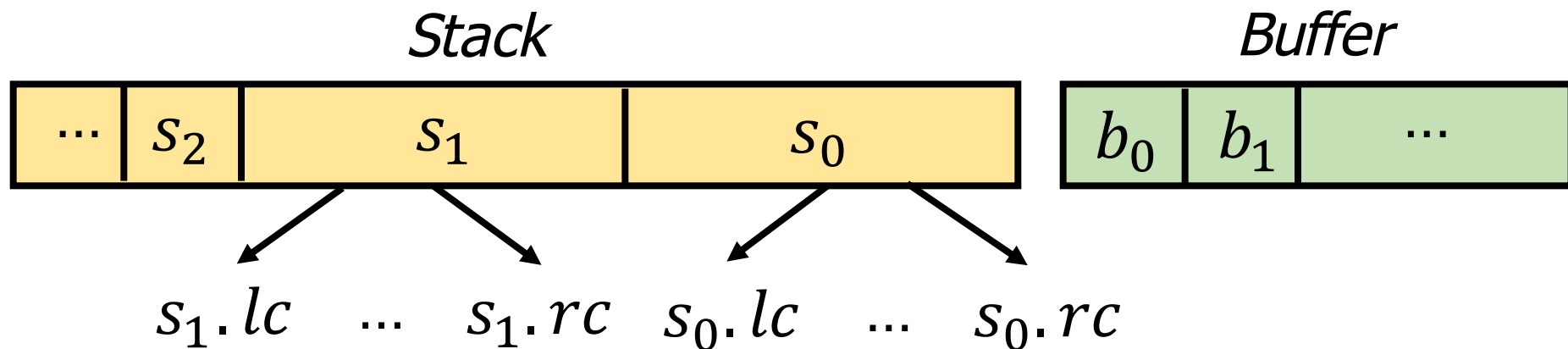
- Time complexities are different!!!



Information about s_1 is not available, needs extra bookkeeping

Time Complexity in Practice

- Complexity depends on feature representation!
- Typical feature representation:
 - Feature templates look at specific *positions* in the stack and in the buffer
- Best-known complexity in practice: $O(n^6)$
(Huang and Sagae, 2010)



Best-known Time Complexities (recap)

$O(n^3)$

Theoretical

Gap:
Feature
representation

$O(n^6)$

Practical

In Practice, Instead of Exact Decoding ...

- Greedy search (Nivre, 2003, 2004, 2008; Chen and Manning, 2014)
- Beam search (Zhang and Clark, 2011; Weiss et al., 2015)
- Best-first search (Sagae and Lavie, 2006; Sagae and Tsujii, 2007; Zhao et al., 2013)
- Dynamic oracles (Goldberg and Nivre, 2012, 2013)
- “Global” normalization on the beam (Zhou et al., 2015; Andor et al., 2016)
- Reinforcement learning (Lê and Fokkens, 2017)
- Learning to search (Daumé III and Marcu, 2005; Chang et al., 2016; Wiseman and Rush, 2016)
- ...

How Many Positional Features Do We Need?

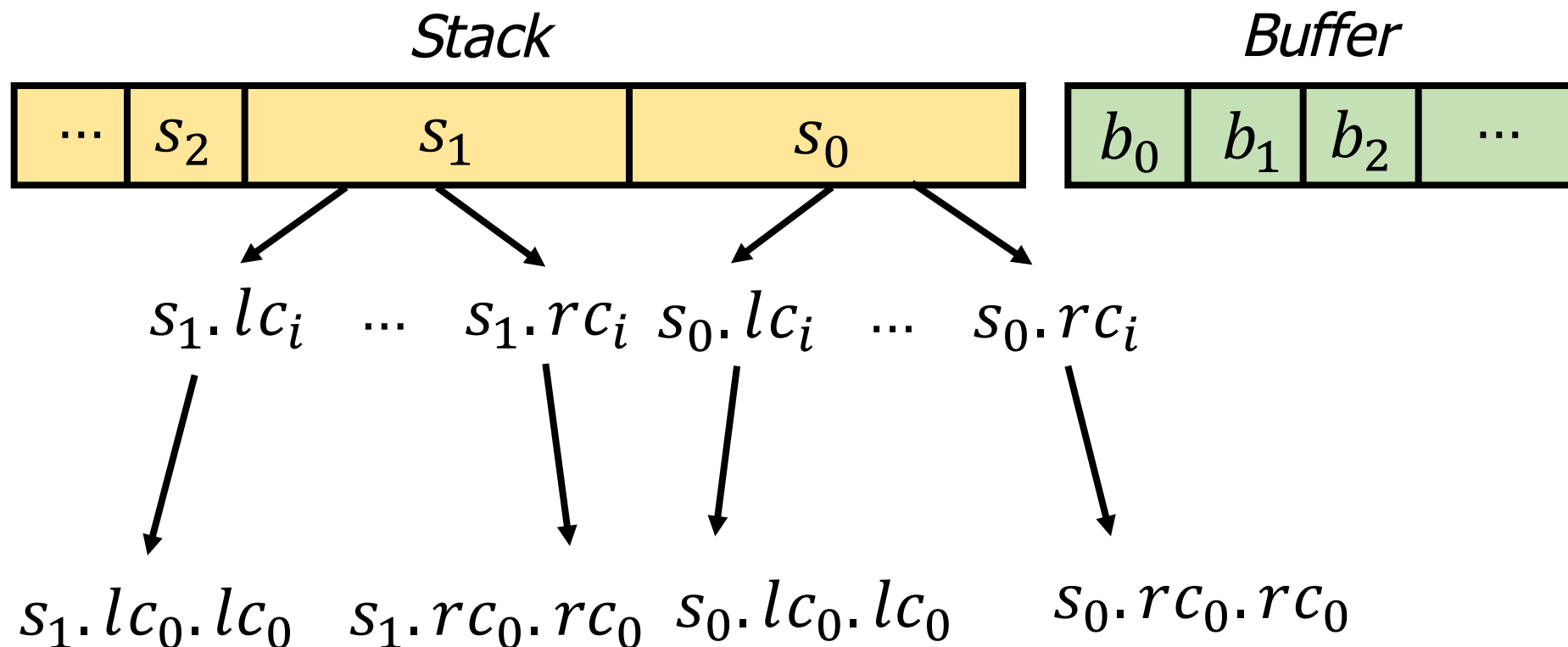
Non-neural (manual engineering)



Chen and Manning (2014)

How Many Positional Features Do We Need?

- Chen and Manning (2014)



How Many Positional Features Do We Need?

Non-neural (manual engineering)



Chen and Manning (2014)



Stack LSTM

(Dyer et al., 2016)



Bi-LSTM

Kiperwasser and
Goldberg (2016)

Cross and
Huang (2016)

More tree-structure information



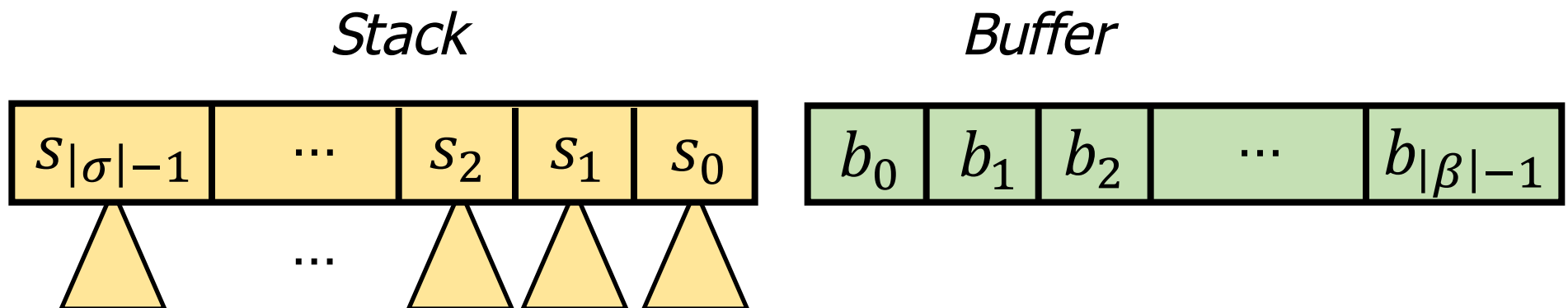
Exponential DP

Slow DP

Fast DP

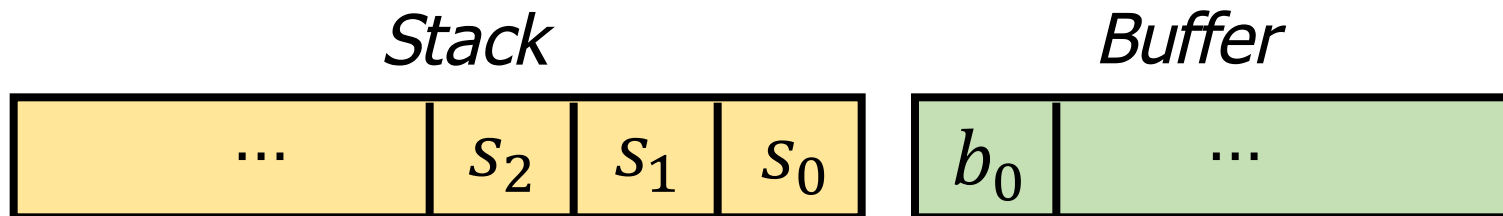
How Many Positional Features Do We Need?

- LSTMs can be used to encode the entire stack and buffer (Dyer et al., 2016)

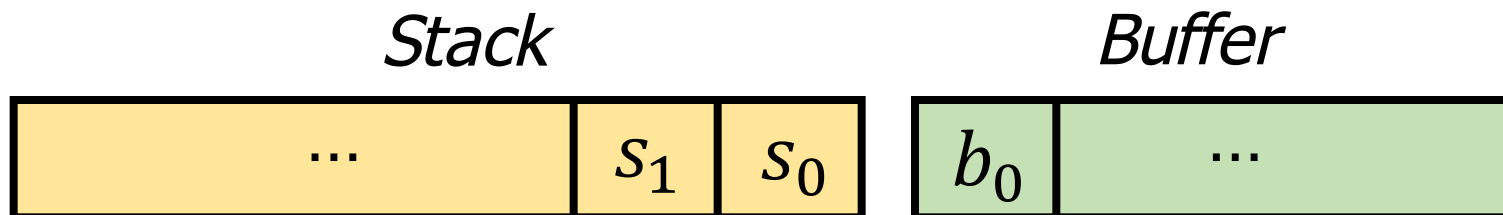


How Many Positional Features Do We Need?

- Bi-LSTMs give compact feature representations (Kiperwasser and Goldberg, 2016; Cross and Huang, 2016)
- Features used in Kiperwasser and Goldberg (2016)



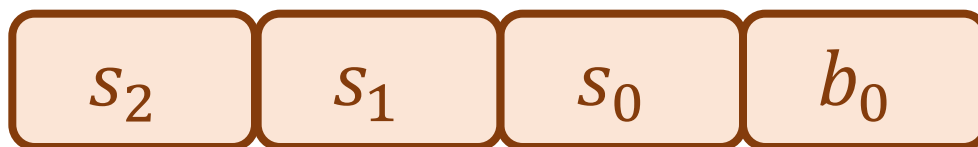
- Features used in Cross and Huang (2016)



Model Architecture

s_{sh}, s_{re}, s_{re}

Multi-layer perceptron

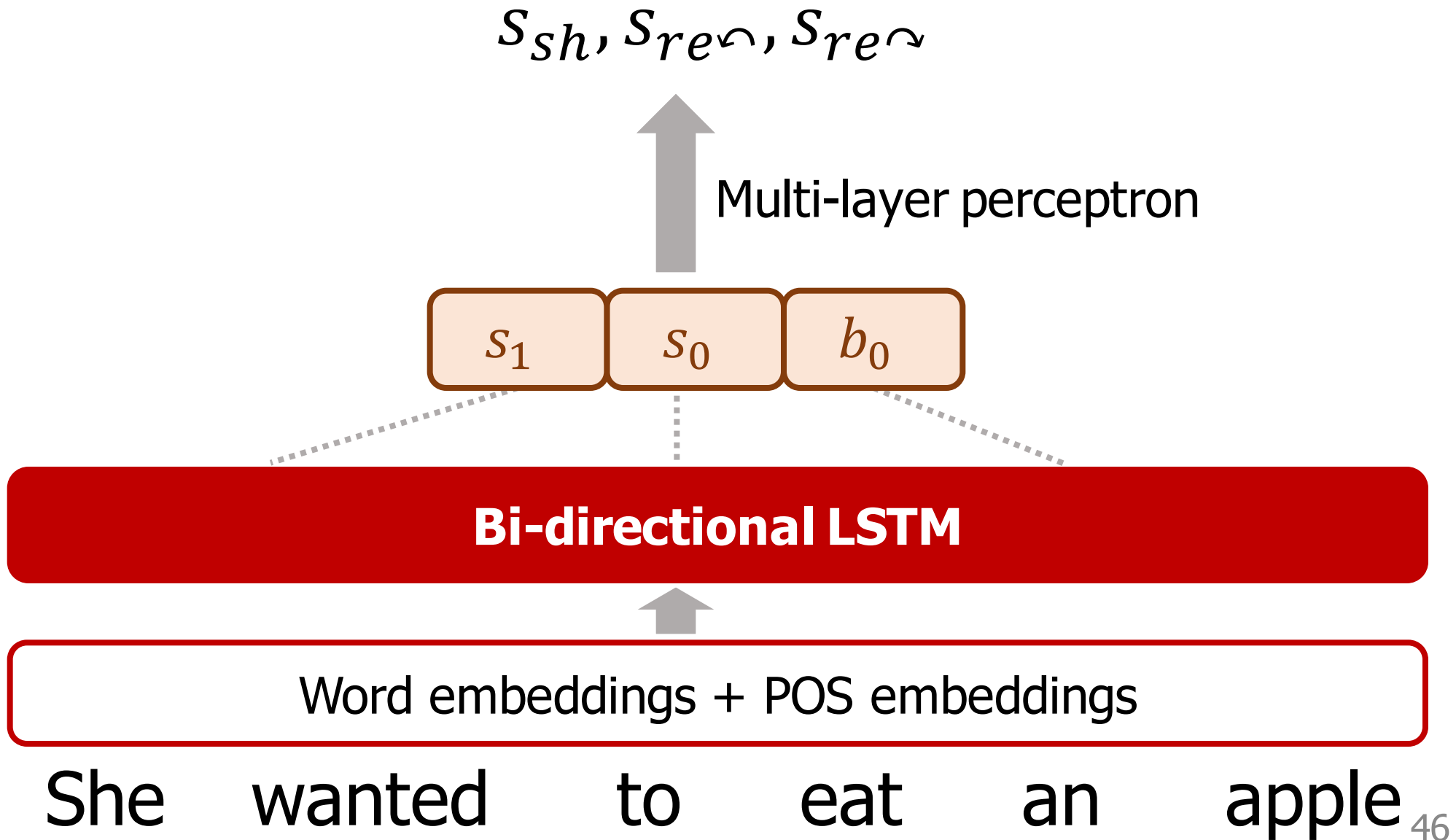


Bi-directional LSTM

Word embeddings + POS embeddings

She wanted to eat an apple₄₅

Model Architecture



Model Architecture

s_{sh}, s_{re}, s_{re}

Multi-layer perceptron



Bi-directional LSTM

Word embeddings + POS embeddings

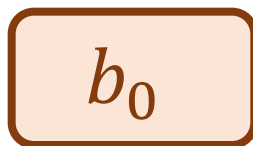
She wanted to eat an apple₄₇

Model Architecture

s_{sh}, s_{re}, s_{re}



Multi-layer perceptron



Bi-directional LSTM

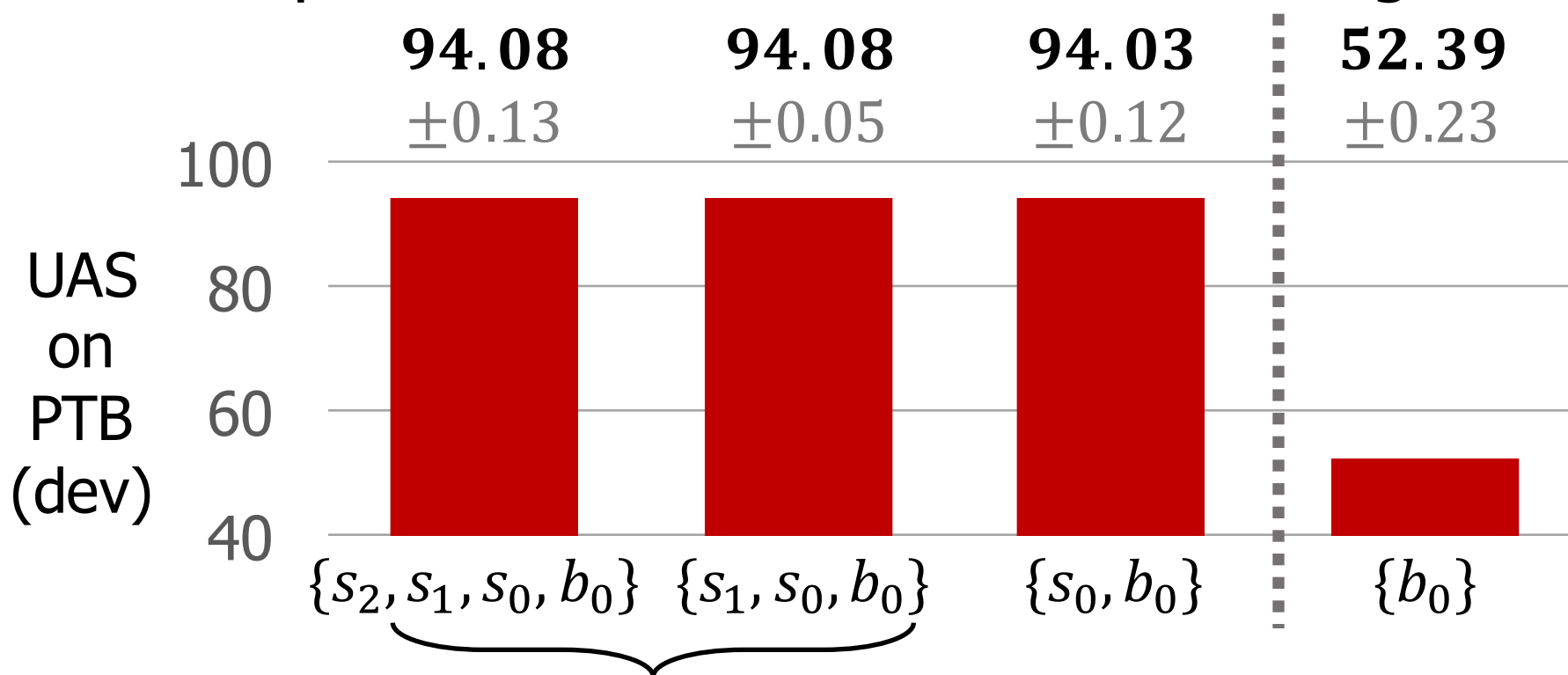


Word embeddings + POS embeddings

She wanted to eat an apple₄₈

How Many Positional Features Do We Need?

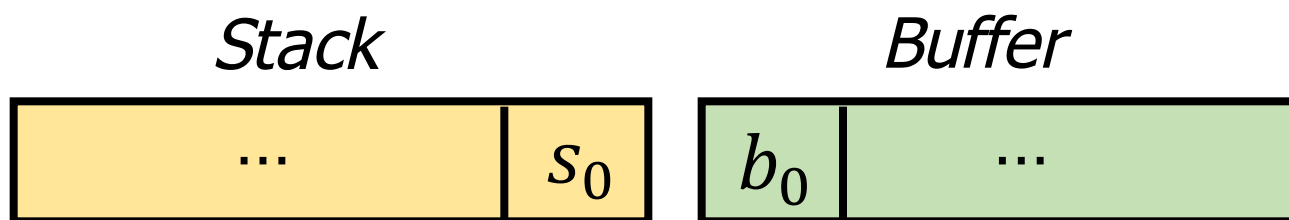
- We answer the question empirically
... experimented with greedy decoding
- Two positional feature vectors are enough!



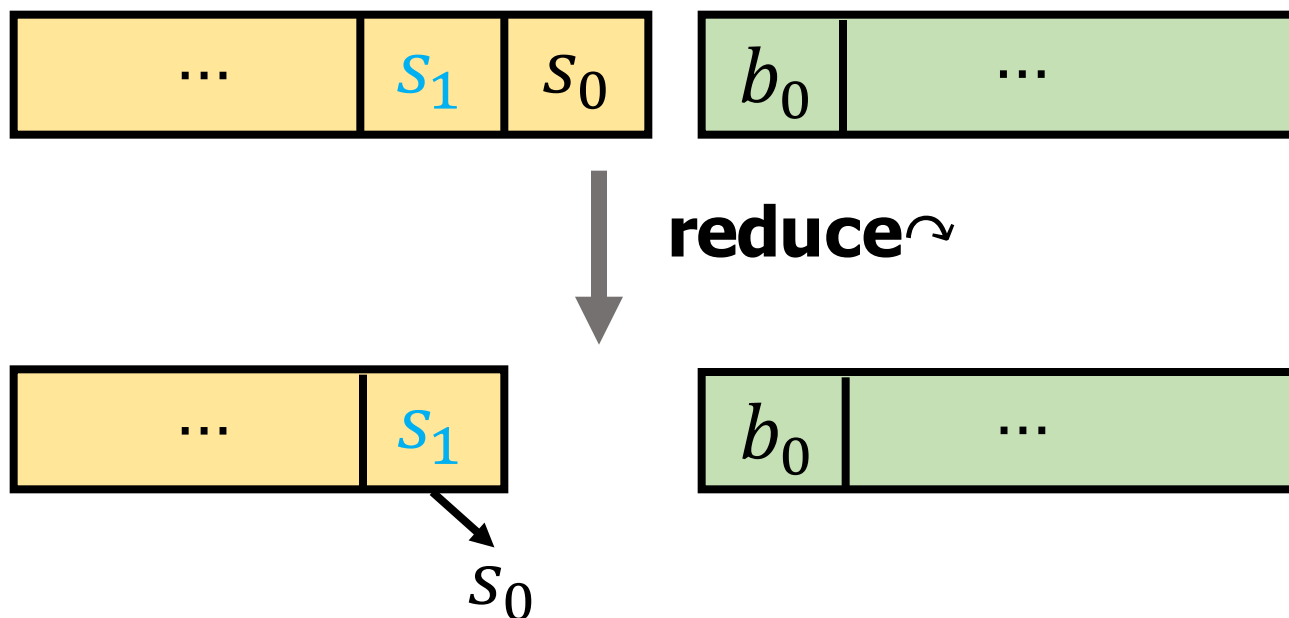
Considered in prior work

How Many Positional Features Do We Need?

- Our minimal feature set works



- Counter-intuitive, but works for *greedy decoding*



Implication of Minimal Feature Set

- The bare deduction items already contain enough information to extract features
- We don't need extra book keeping
- Leads to the first $O(n^3)$ implementation of global decoders!

How Many Positional Features Do We Need?

Non-neural (manual engineering)



Chen and Manning (2014)



Stack LSTM

(Dyer et al., 2016)



Bi-LSTM

Kiperwasser and
Goldberg (2016)

Cross and
Huang (2016)

Our work

More tree-structure information



Exponential DP

Slow DP

Fast DP

Fast(er) DP

Best-known Time Complexities (recap)

$O(n^3)$

Theoretical

Gap:
Feature
representation

$O(n^6)$

Practical

Our contribution

$O(n^3)$

Theoretical

Bi-directional
LSTMs



$O(n^3)$

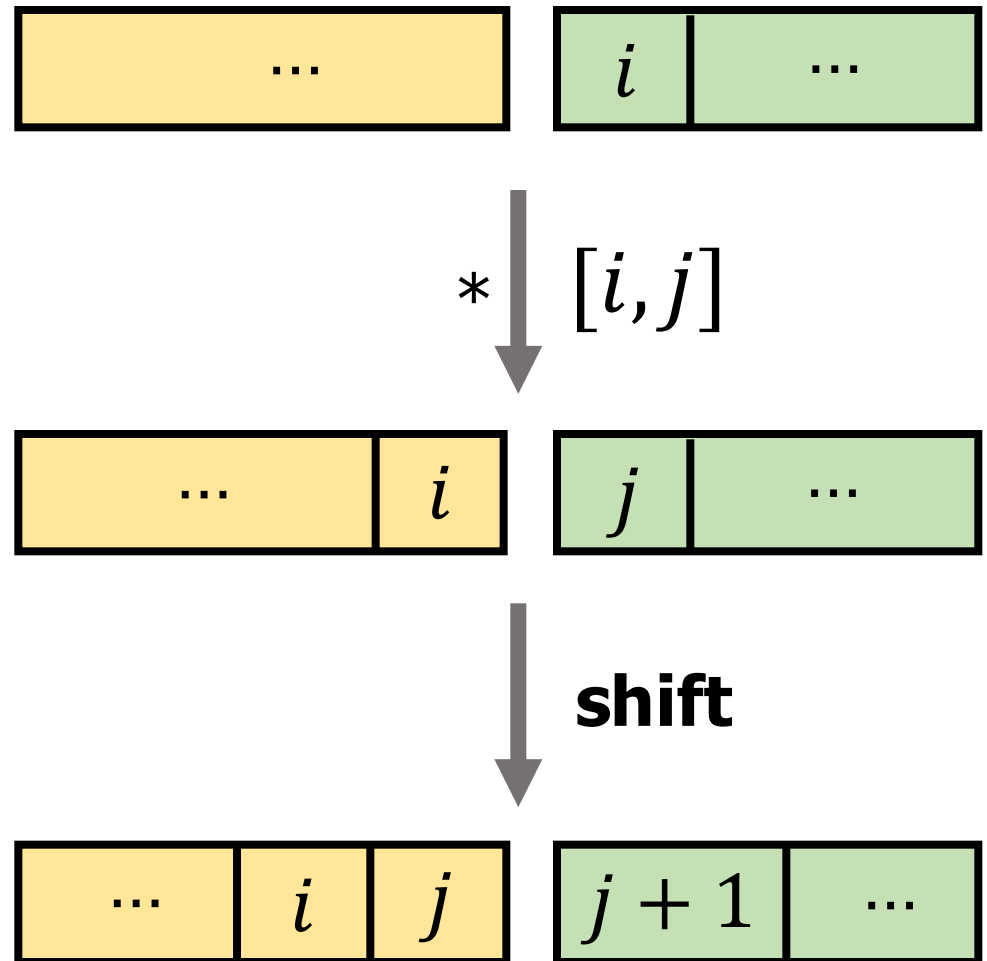
~~$O(n^6)$~~

Practical

Decoding

Score of the sub-sequence

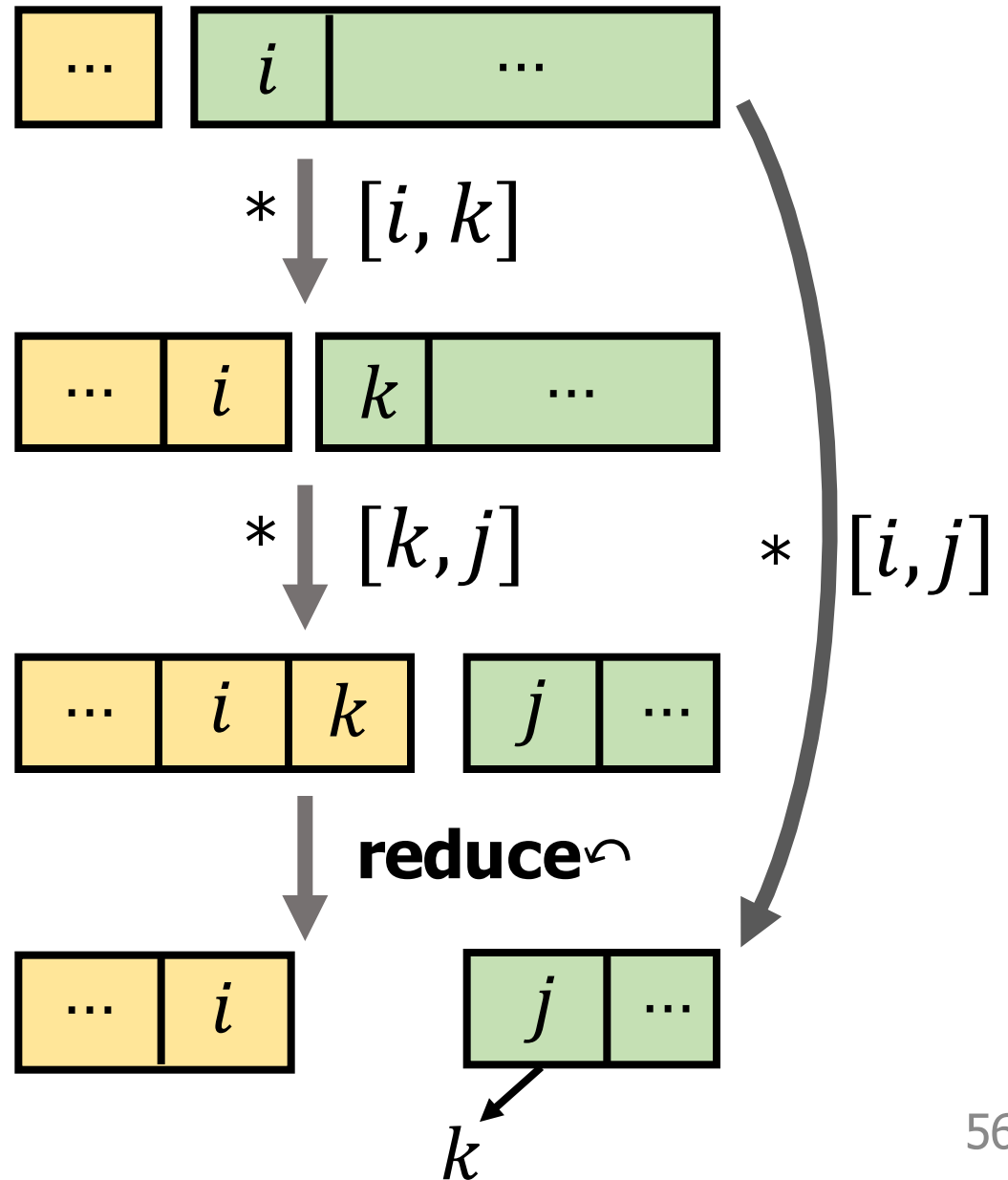
shift $\frac{[i, j]: v}{[j, j + 1]: 0}$



Decoding

$$\text{reduce} \curvearrowright \frac{[i, k]: v_1 \quad [k, j]: v_2}{[i, j]: v_1 + v_2 + \Delta}$$

$$\Delta = s_{sh}(i, k) + s_{re} \curvearrowright (k, j)$$

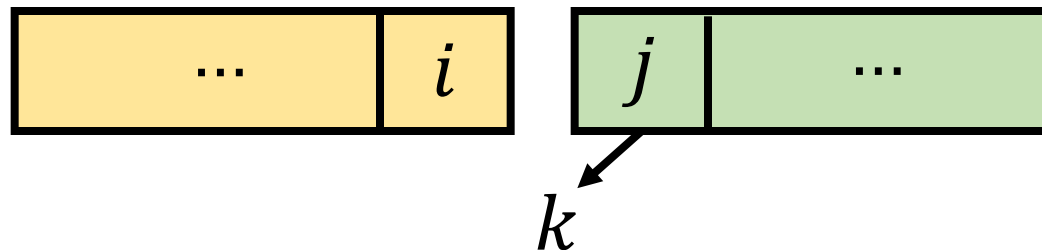


Training

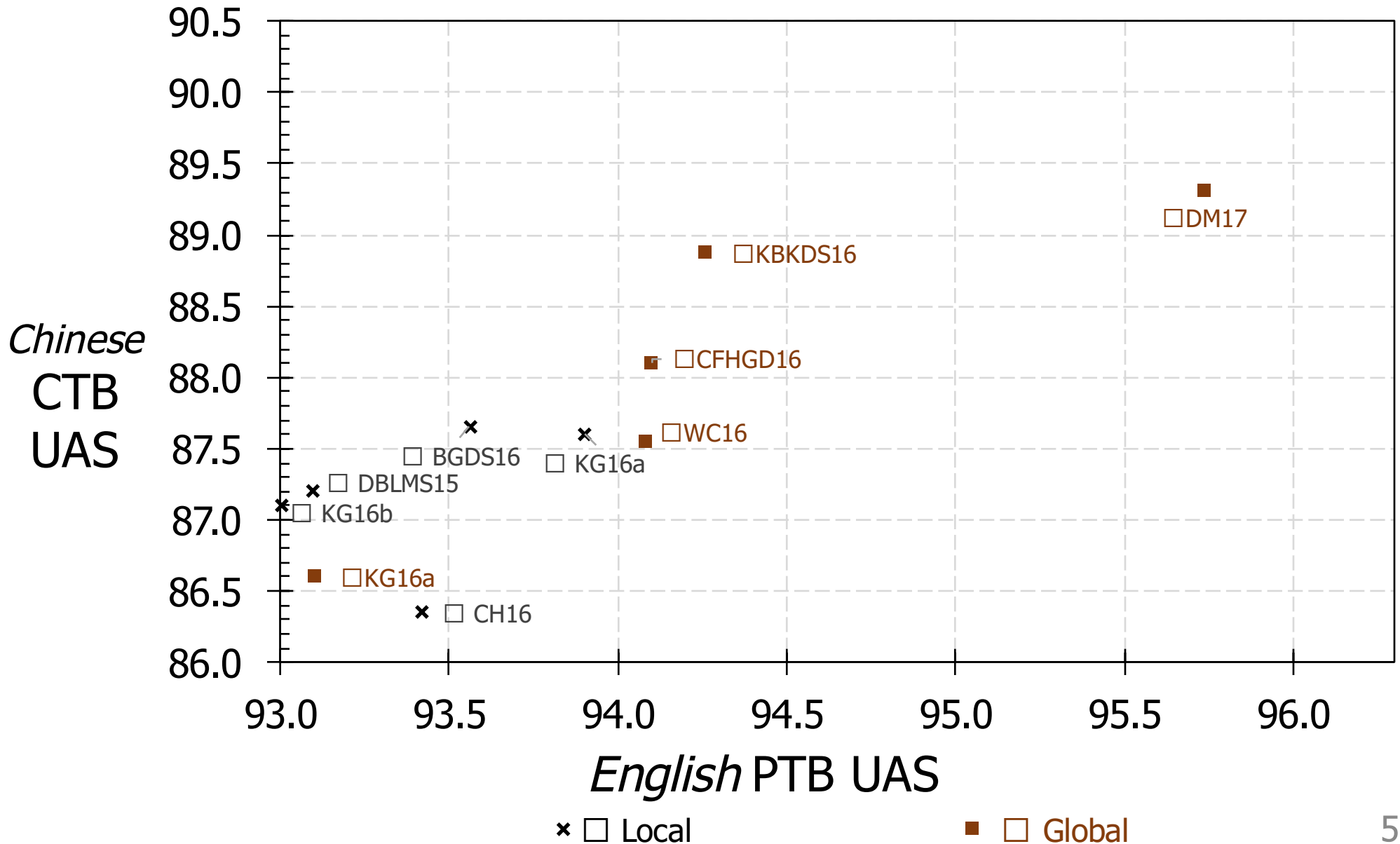
- Cost-augmented decoding (Taskar et al., 2005)

$$\max \text{score}(\bullet \rightarrow \bullet \rightarrow \dots \rightarrow \bullet) + \text{cost}(\begin{matrix} \bullet \rightarrow \bullet \rightarrow \dots \rightarrow \bullet \\ \bullet \rightarrow \bullet \rightarrow \dots \rightarrow \star \end{matrix}) - \text{score}(\bullet \rightarrow \bullet \rightarrow \dots \rightarrow \star)$$

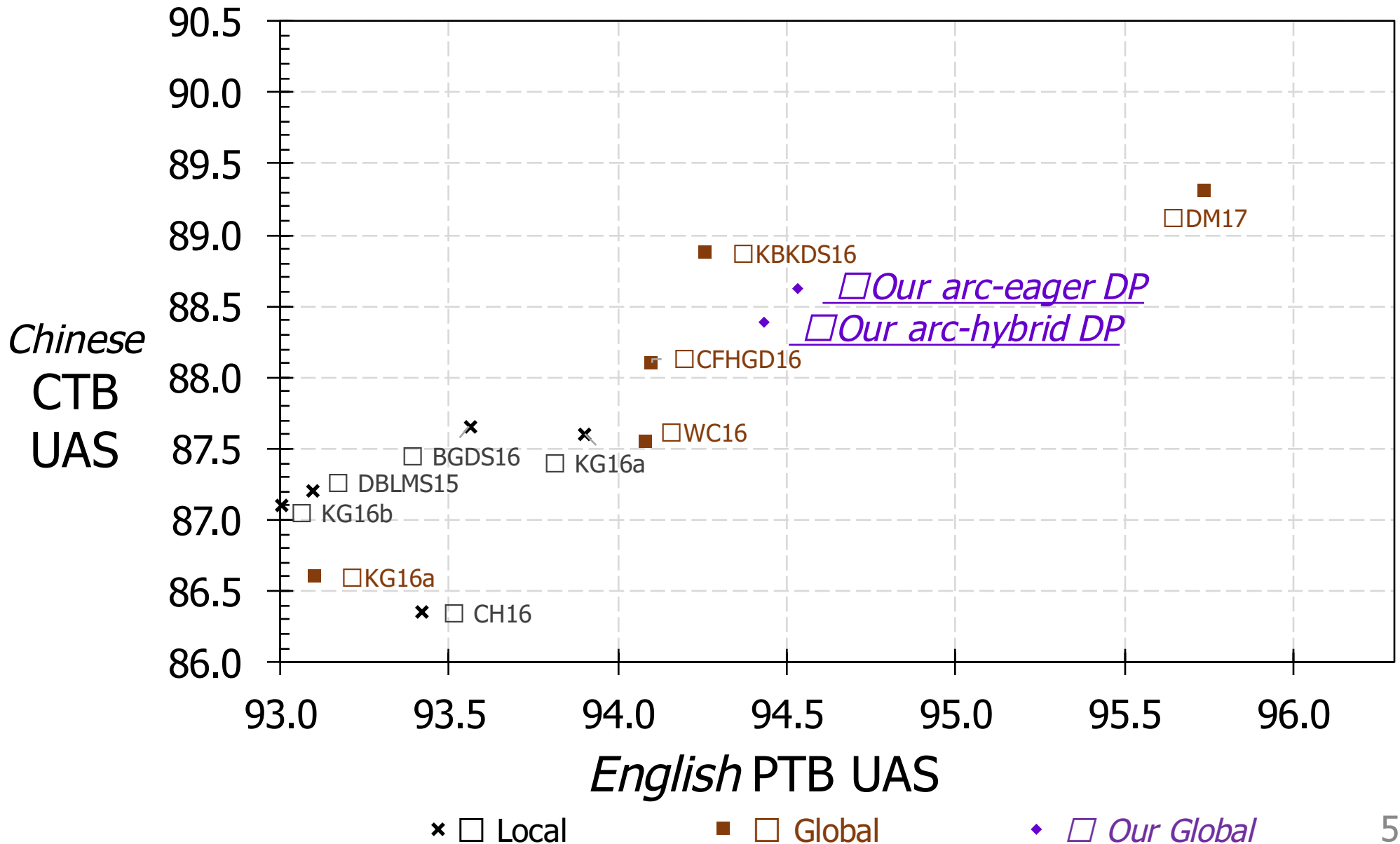
$$\text{reduce}^{\curvearrowright} \frac{[i, k]: v_1 \quad [k, j]: v_2}{[i, j]: v_1 + v_2 + s_{sh}(i, k) + s_{re}^{\curvearrowright}(k, j) + \mathbf{1}(\text{head}(k) \neq j)}$$



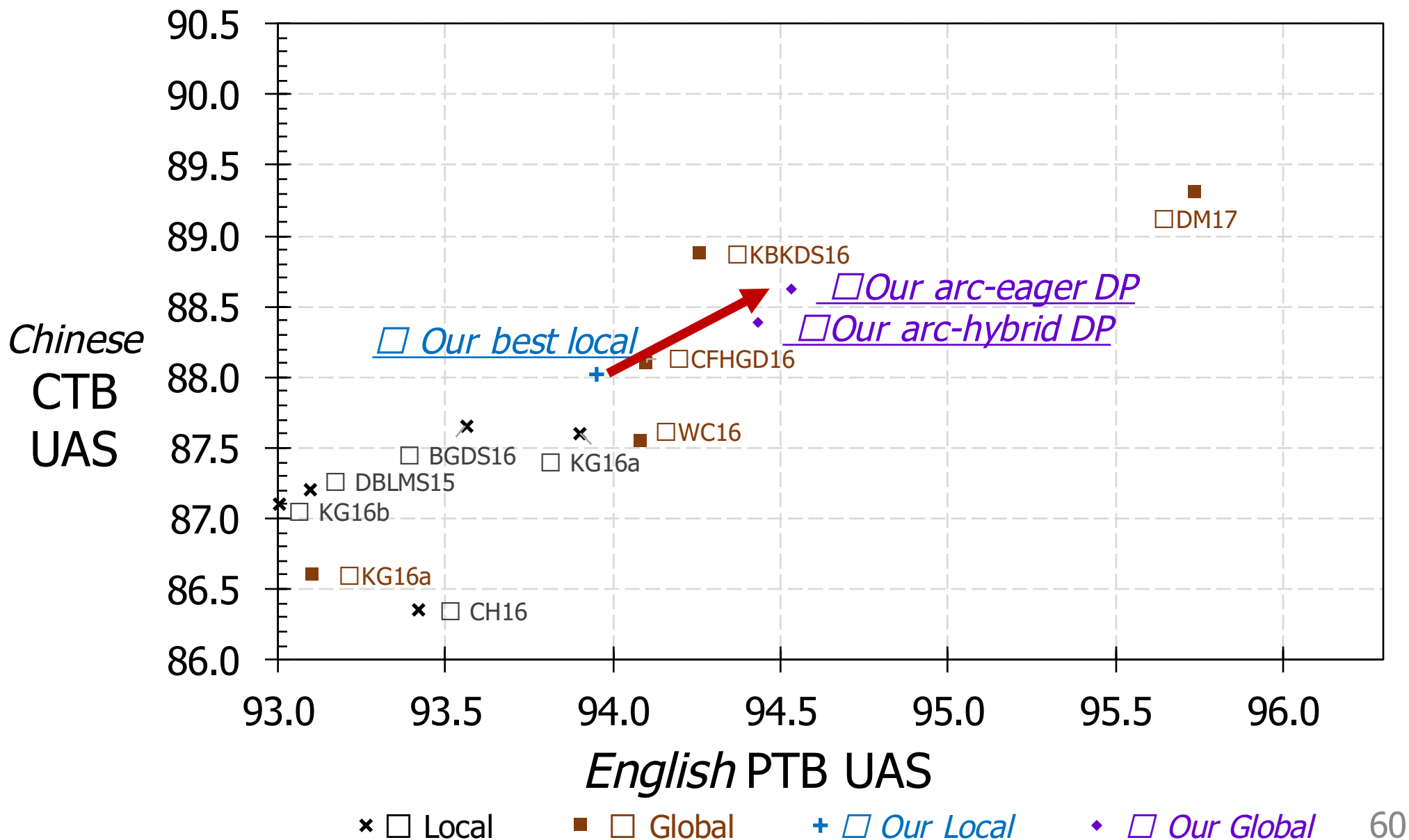
Comparing with State-of-the-art



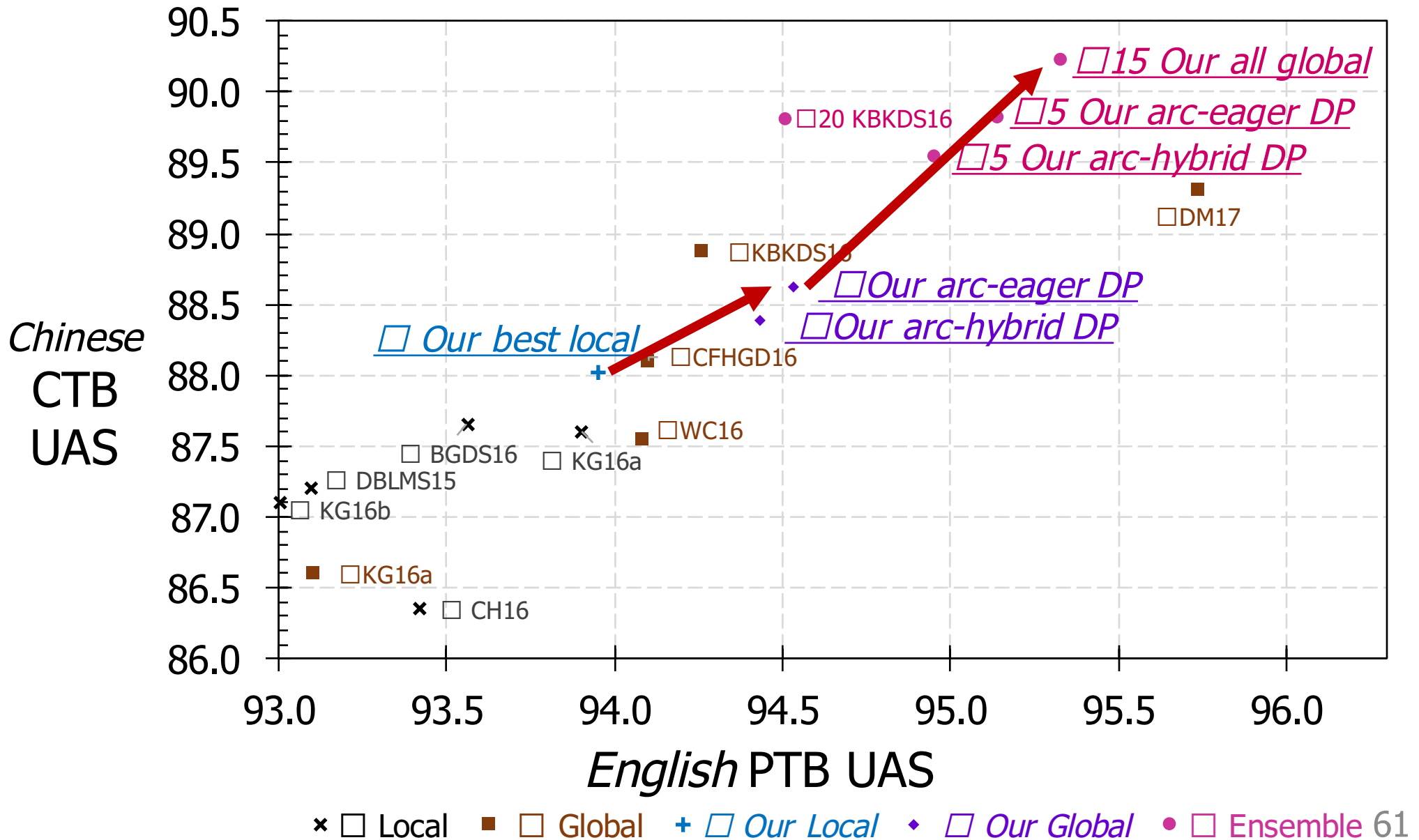
Comparing with State-of-the-art



Comparing with State-of-the-art

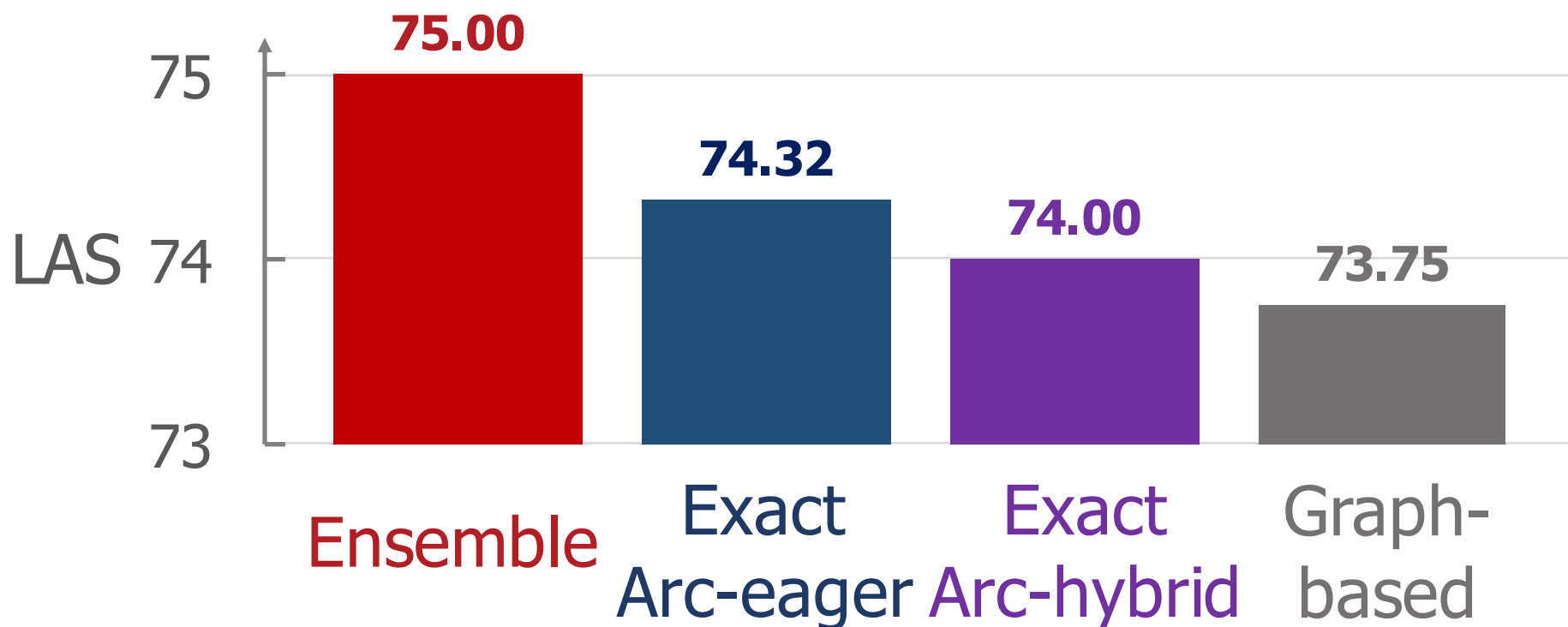


Comparing with State-of-the-art



Results – CoNLL'17 Shared Task

- Macro-average of 81 treebanks in 49 languages
- 2nd–highest overall performance




Conclusion

- Bi-LSTM feature set is minimal yet highly effective
- First $O(n^3)$ implementation of exact decoders
- Global training and decoding gave high performance

More in Our Paper

- Description and analysis of three transition systems (arc-standard, arc-hybrid, arc-eager)
- CKY-style representations of the deduction systems
- Theoretical analysis of the global methods
 - Arc-eager models can “simulate” arc-hybrid models
 - Arc-eager models can “simulate” edge-factored models

Fast(er) Exact Decoding and Global Training for Transition-Based Dependency Parsing via a Minimal Feature Set

 <https://github.com/tzshi/dp-parser-emnlp17>

Tianze Shi*

Liang Huang†

Lillian Lee*

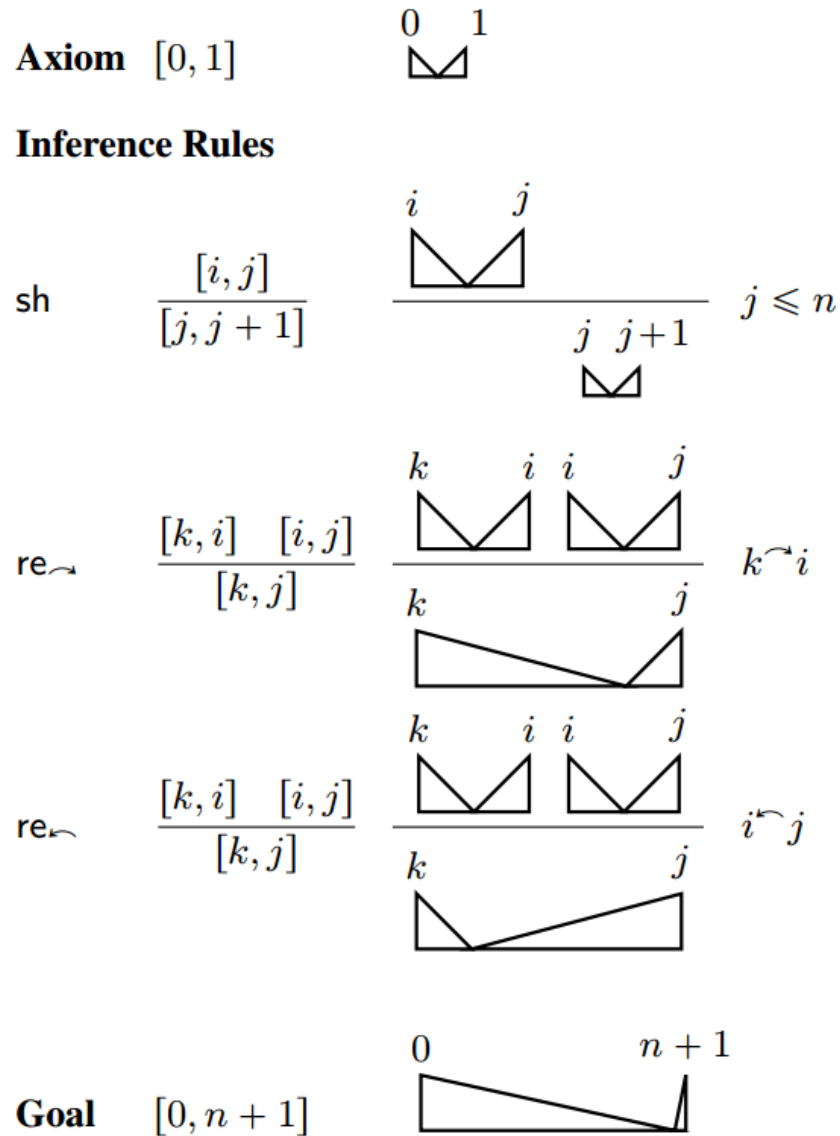


* Cornell
University



† Oregon State
University

CKY-style Visualization



(b) Arc-hybrid