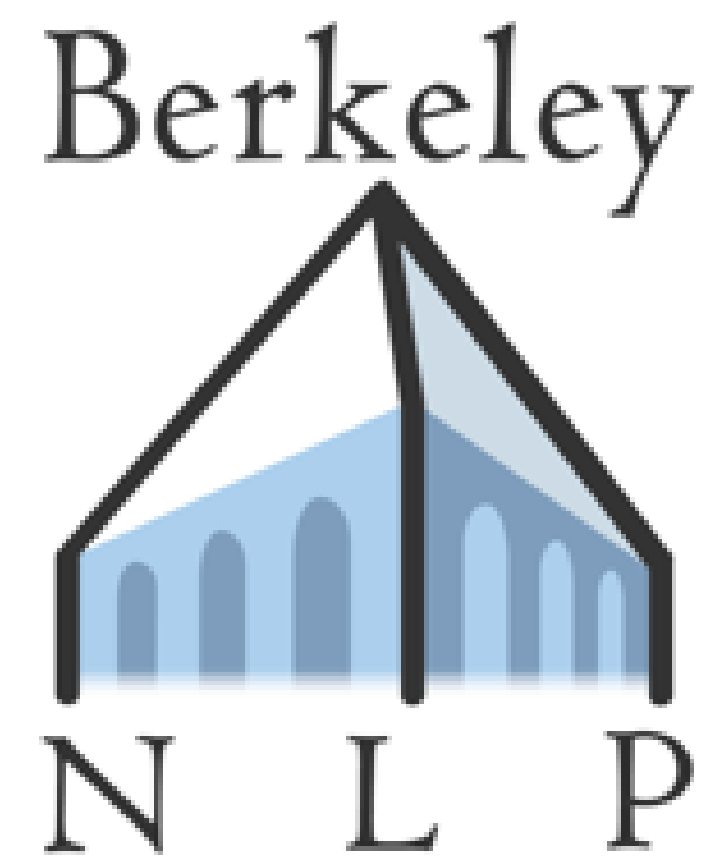


Neural Constituency Parsing



Dan Klein
CS 288

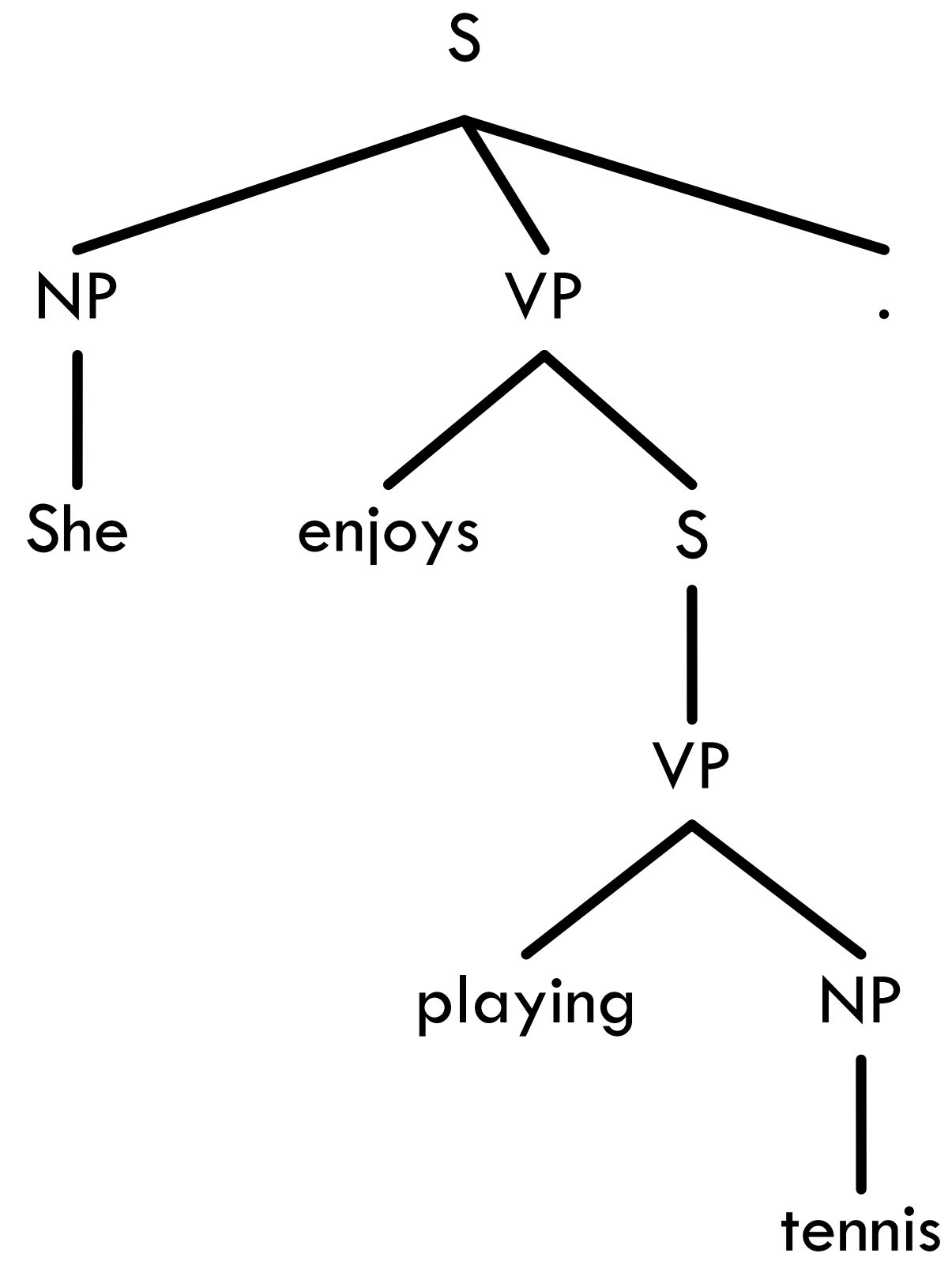


Syntactic Parsing

She enjoys playing tennis.

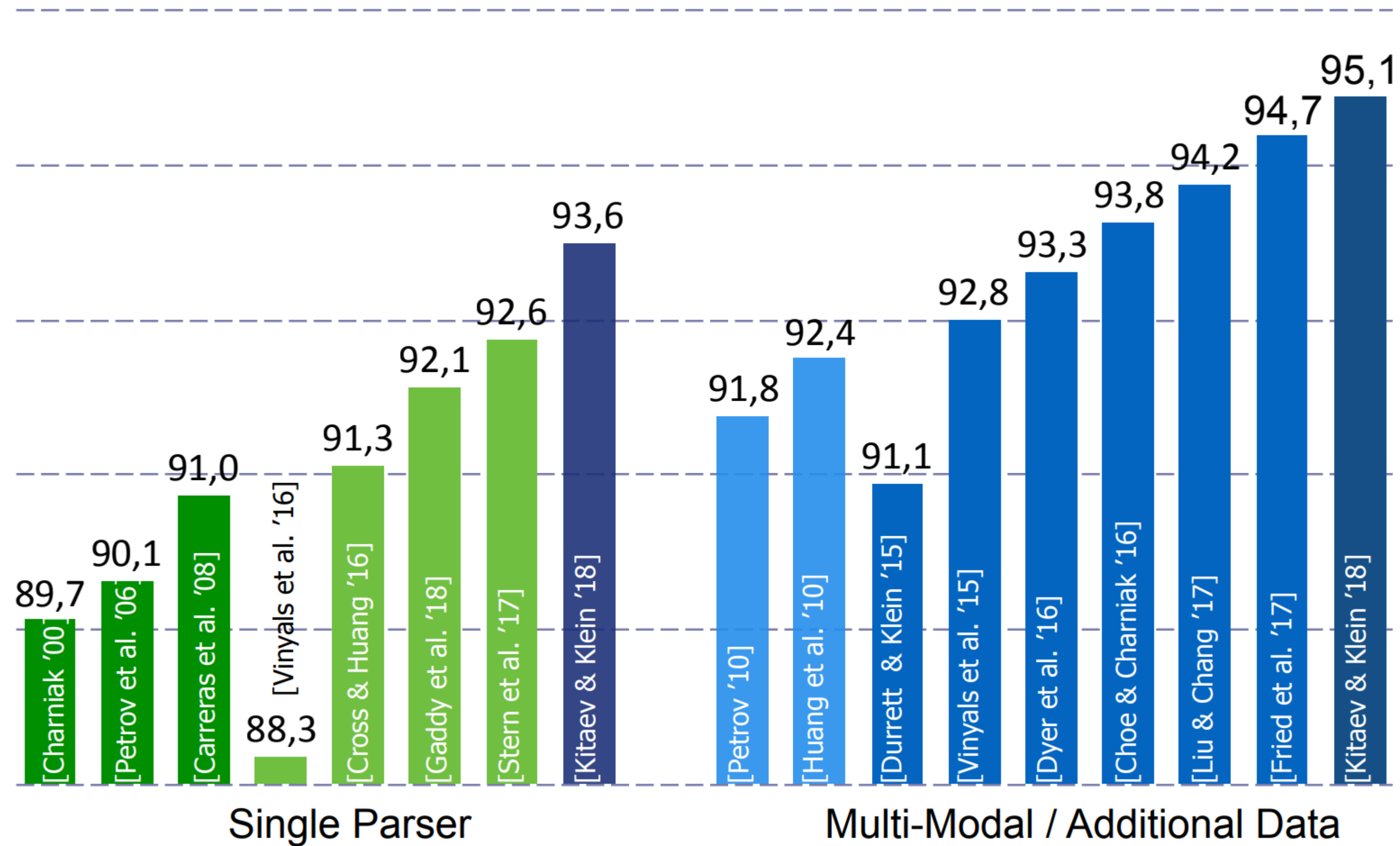


Syntactic Parsing



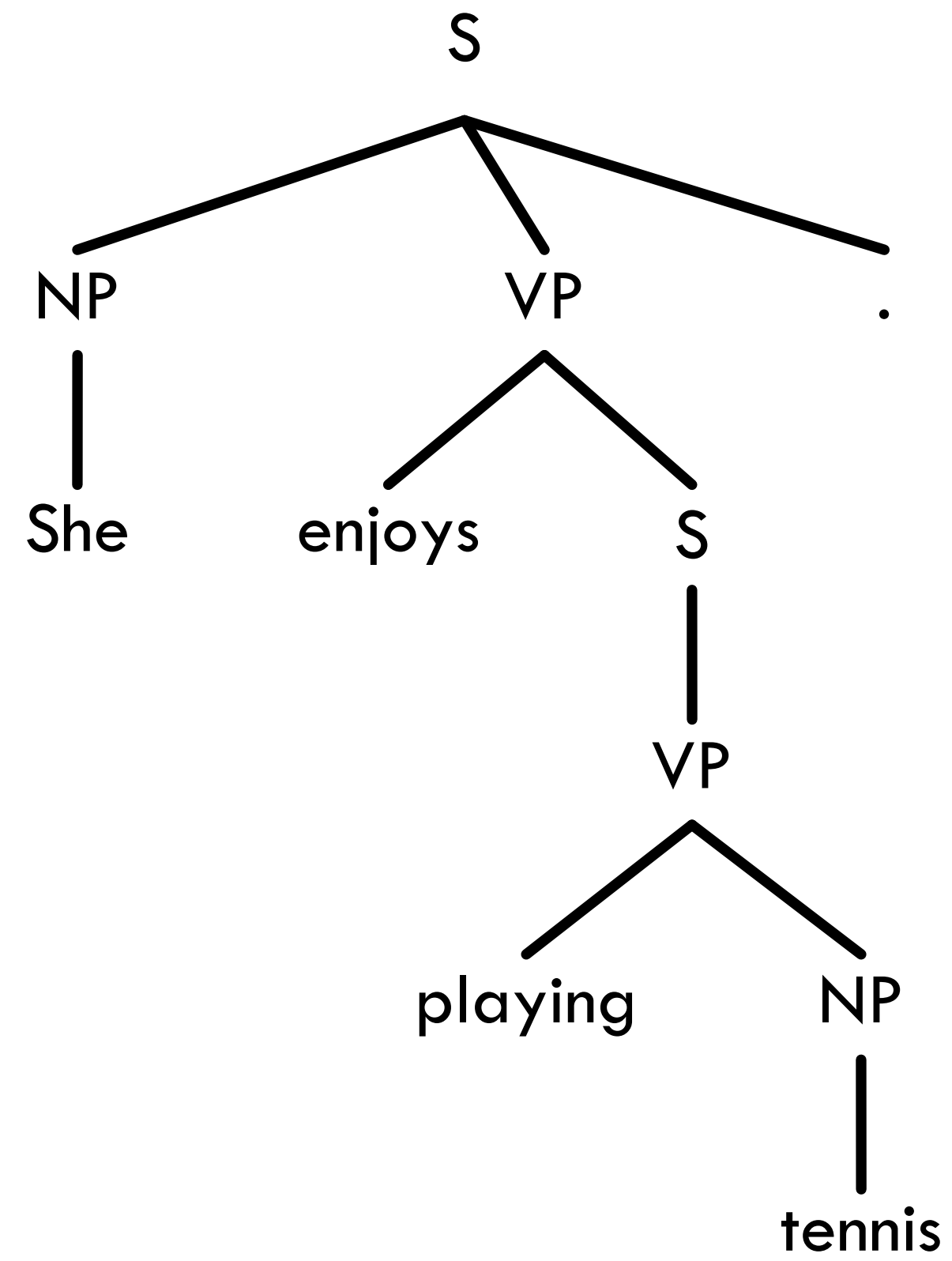


Historical Trends





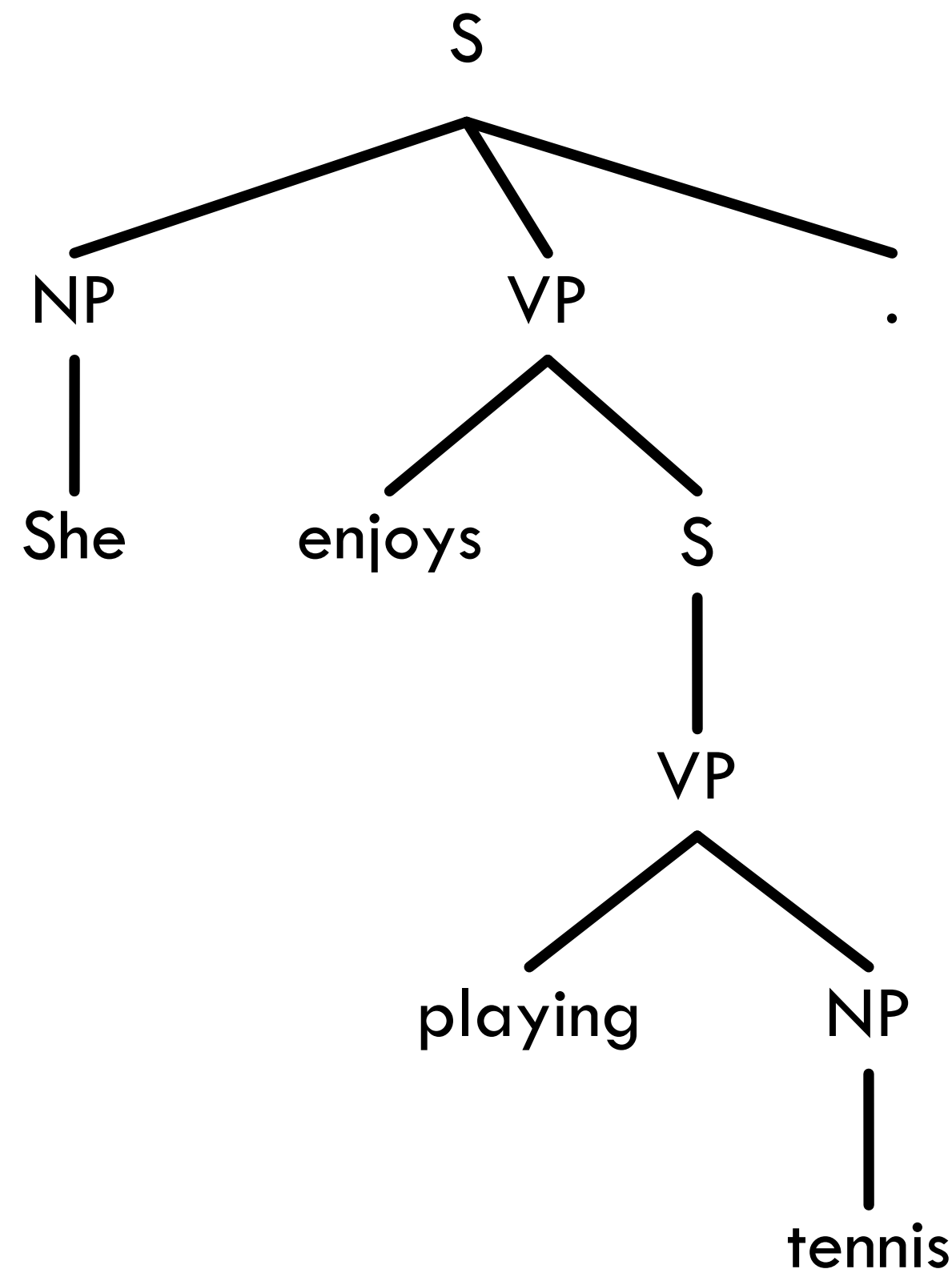
Output Correlations





Grammars

$S \rightarrow NP VP$



$VP[enjoys] : S[playing]$

$NP^S \rightarrow she$

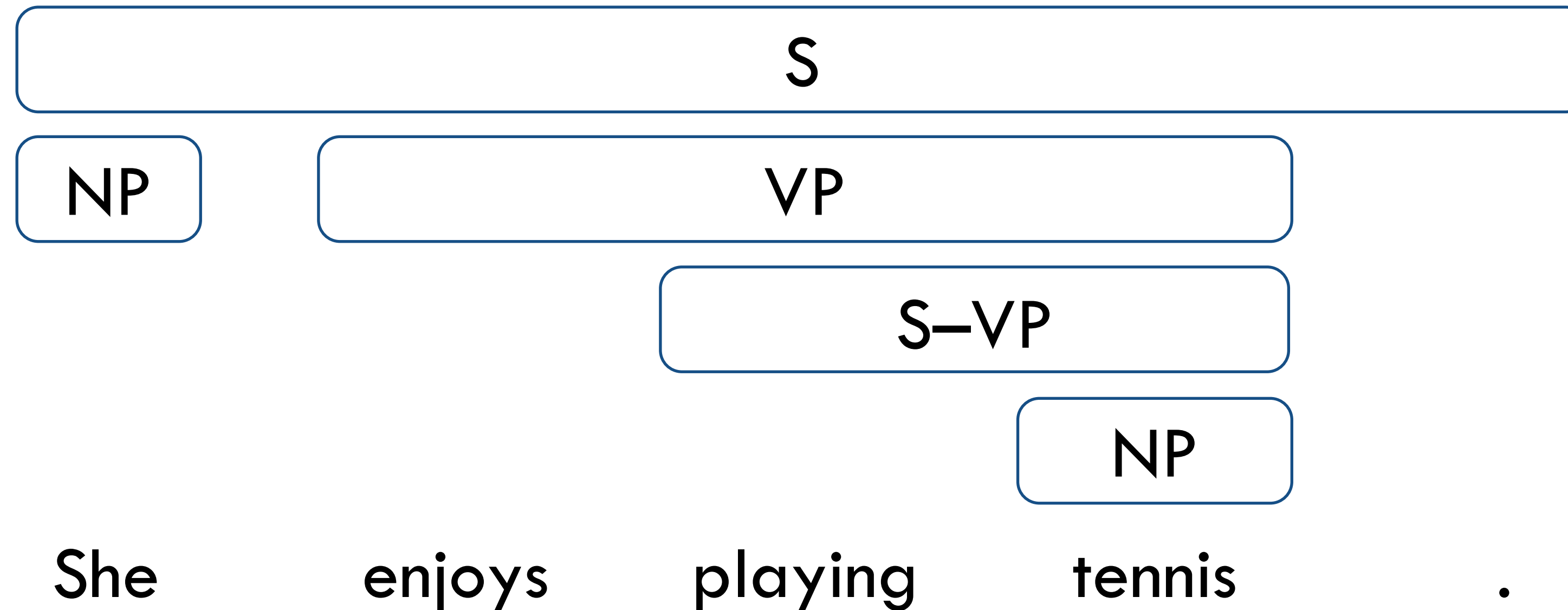
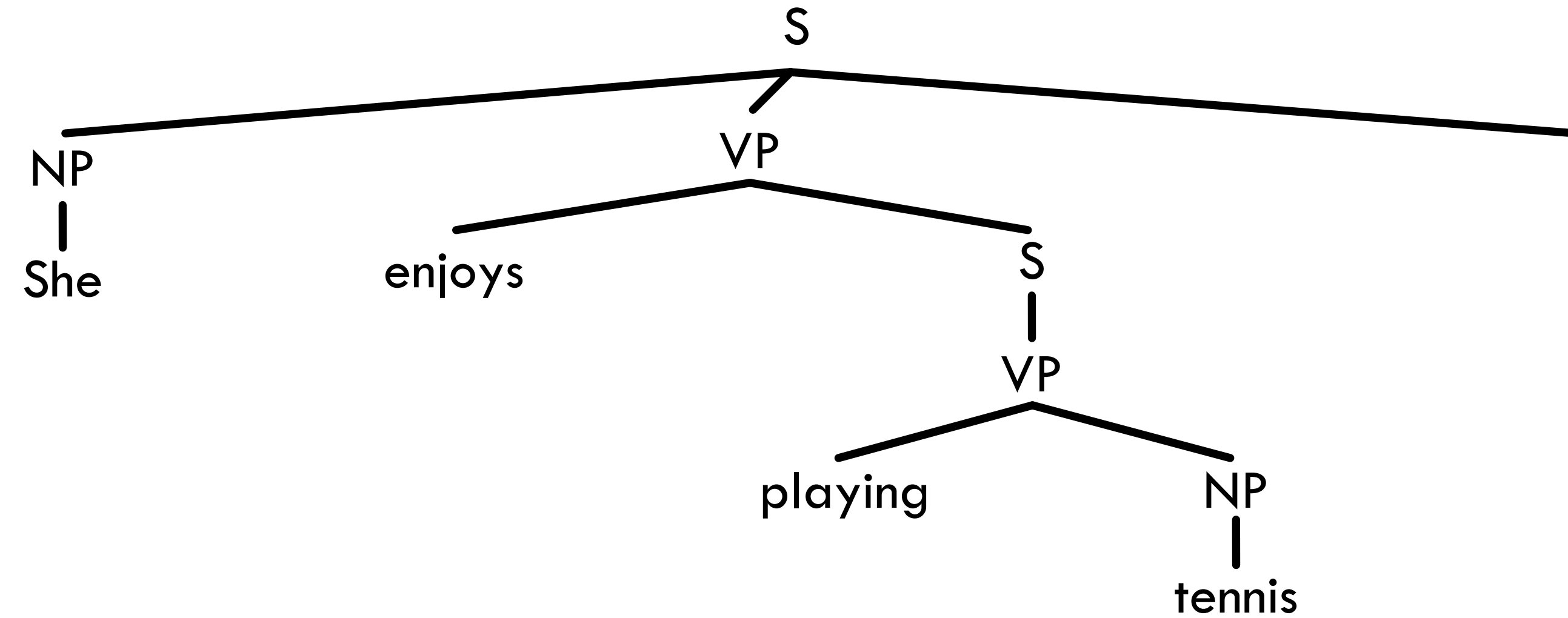


Input-Output Correlations

She enjoys playing tennis.

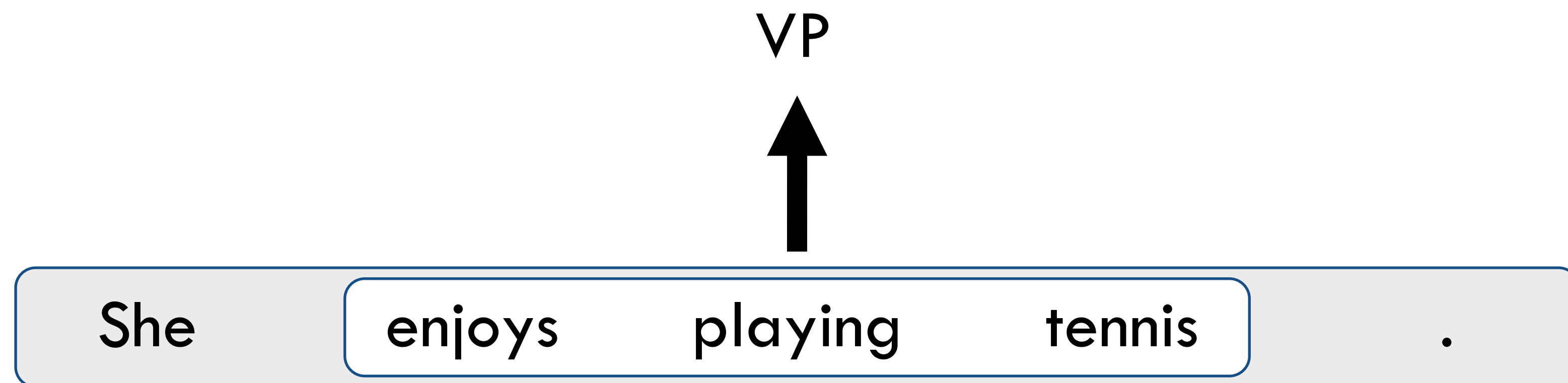


Span-Based Parsing



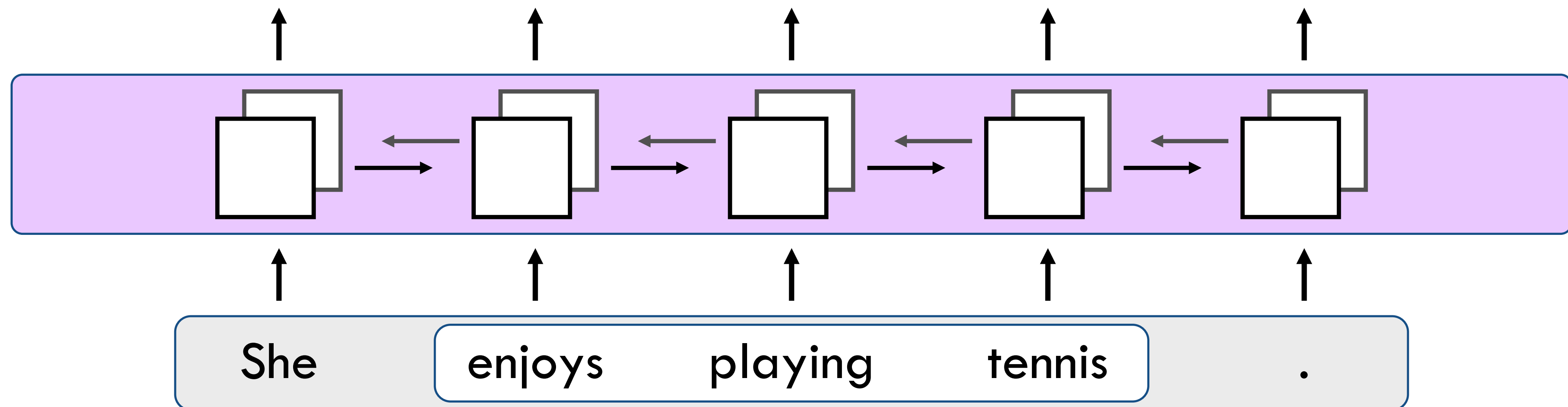


Parsing as Span Classification



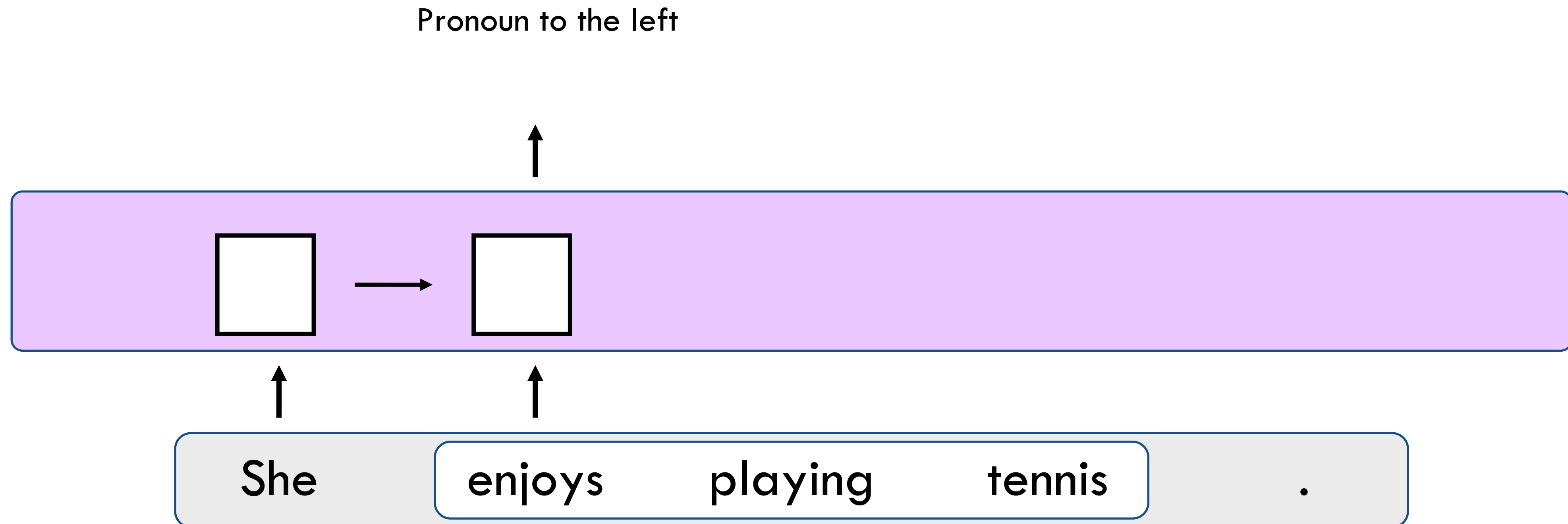


Routing with LSTMs



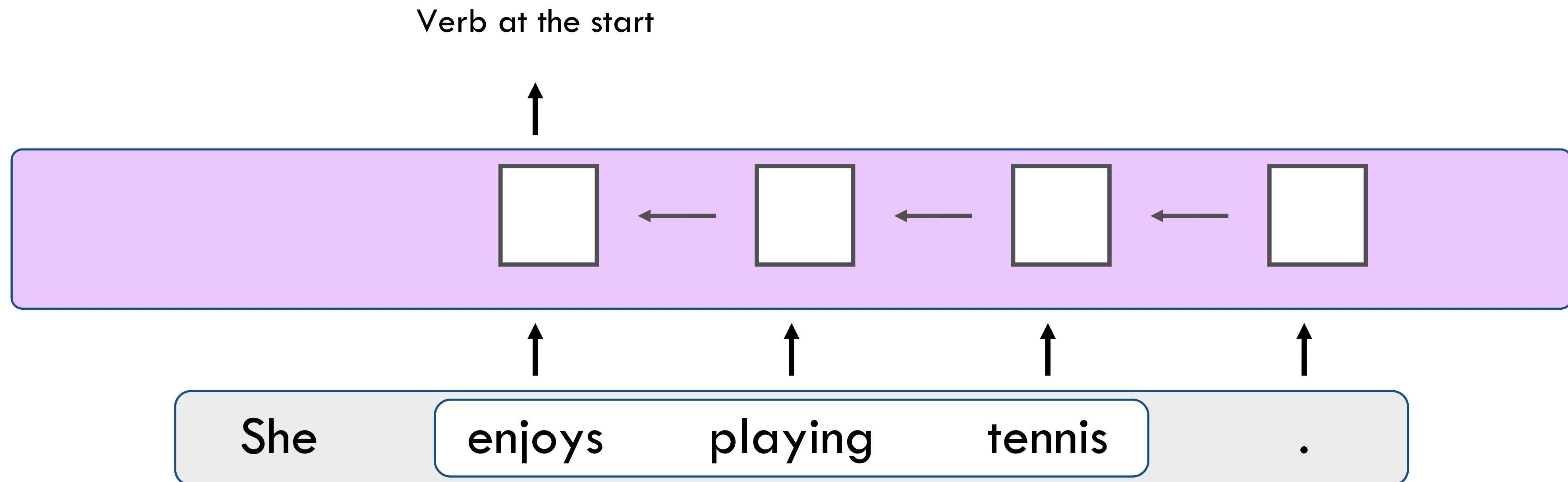


Routing with LSTMs



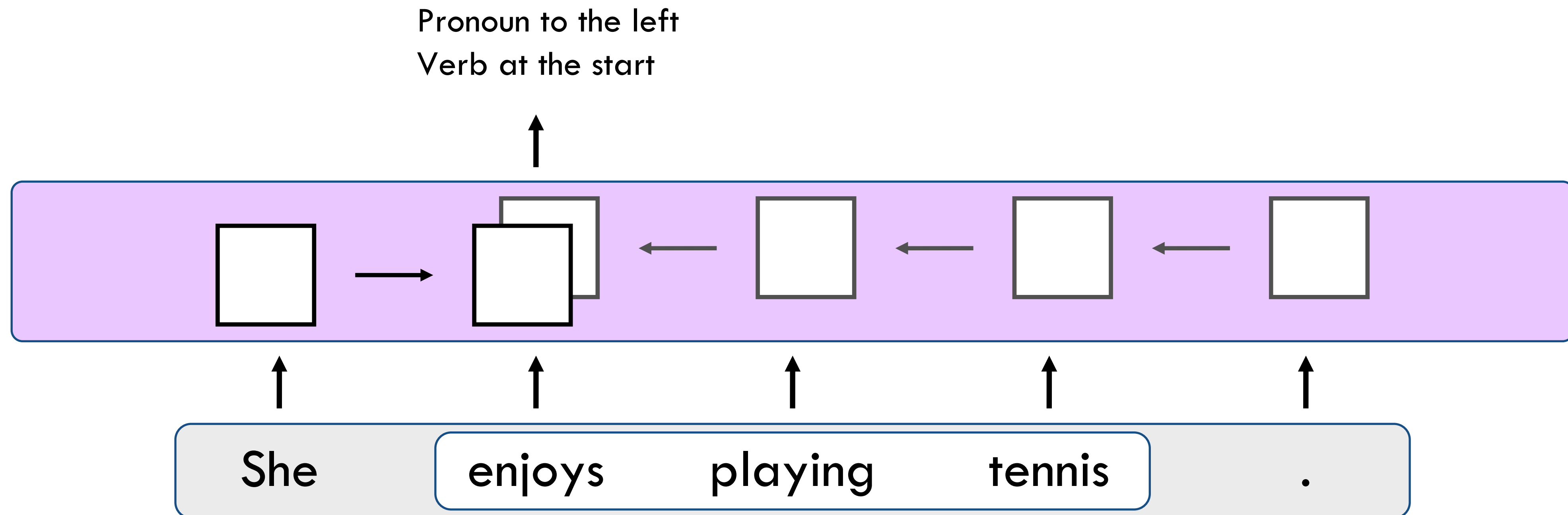


Routing with LSTMs



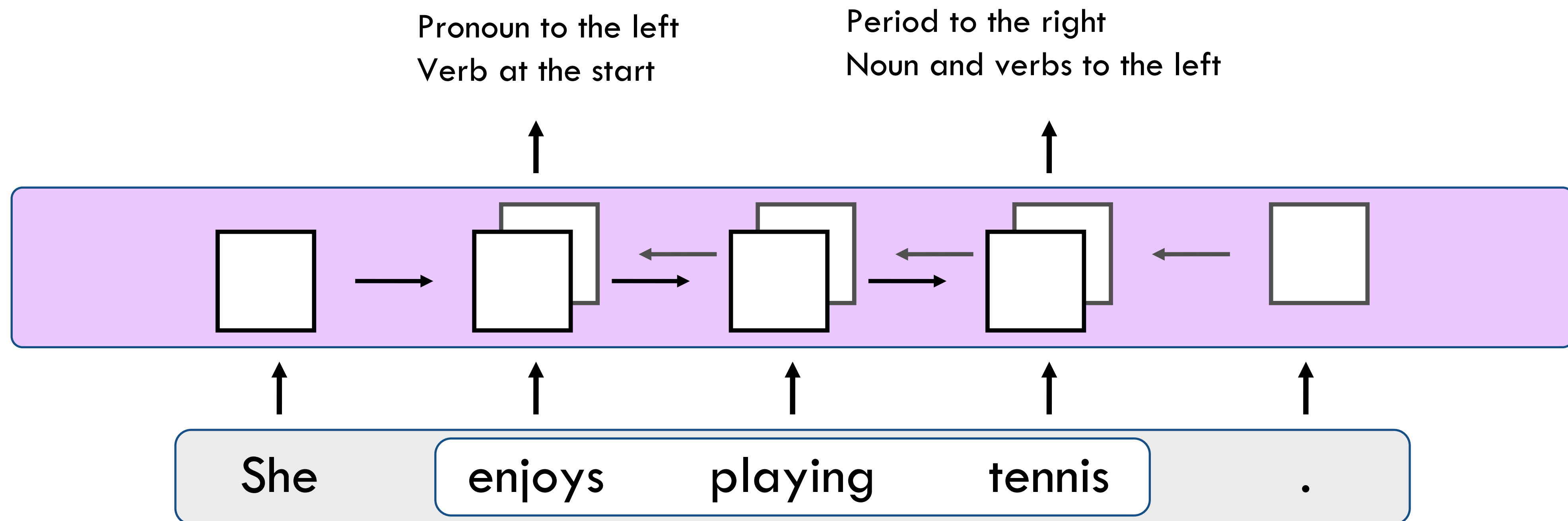


Routing with LSTMs



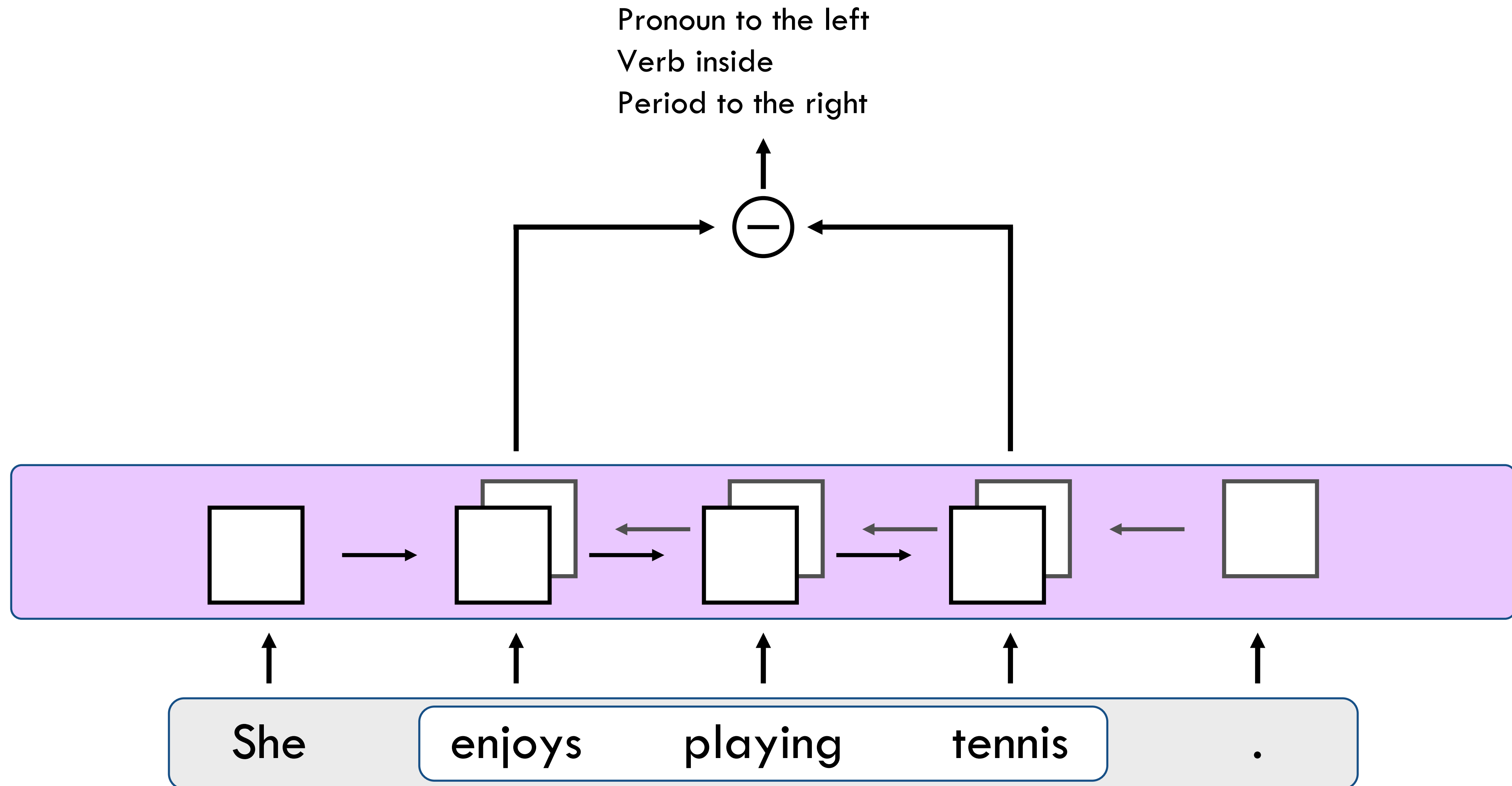


Routing with LSTMs



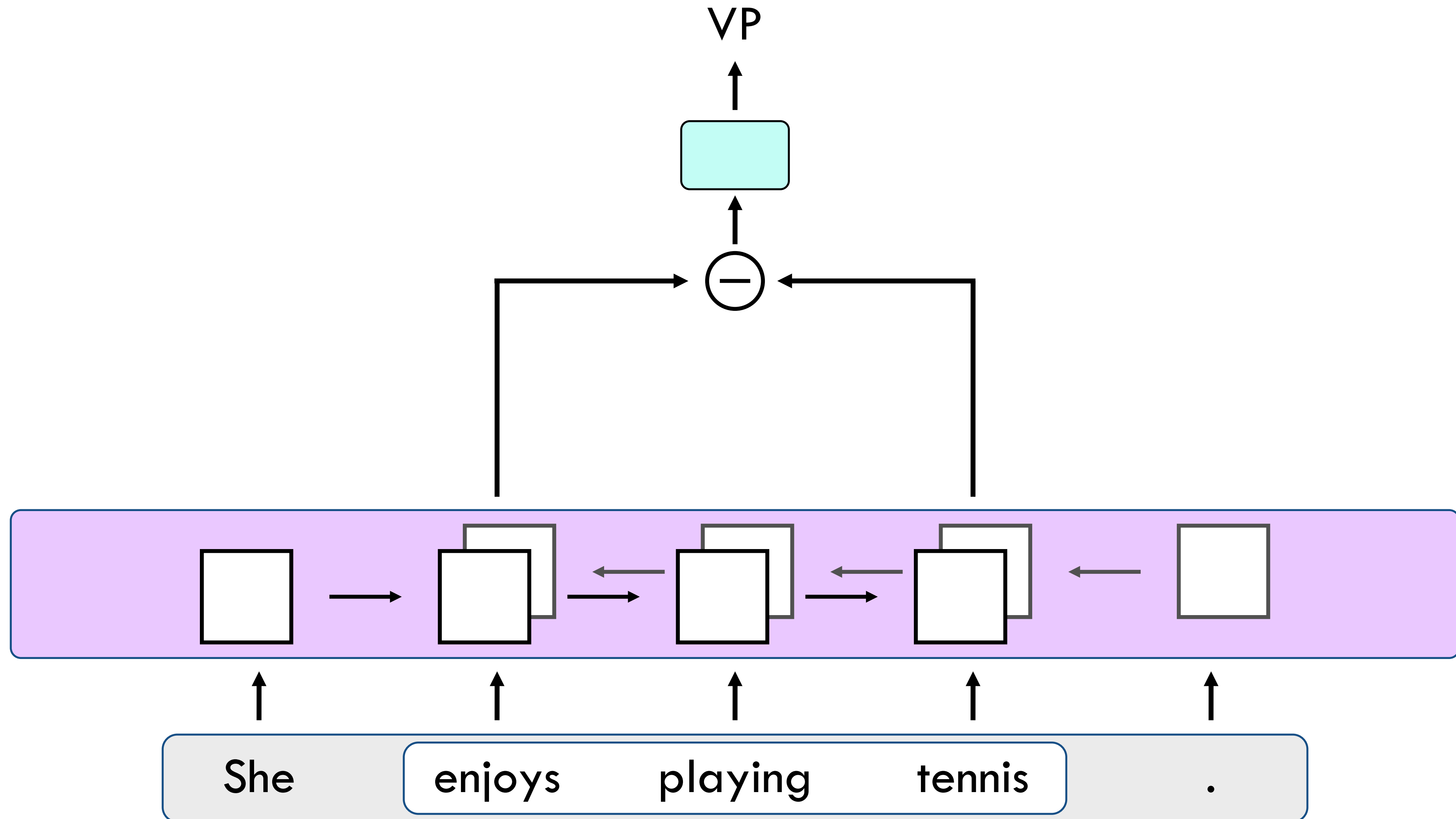


Span Classification



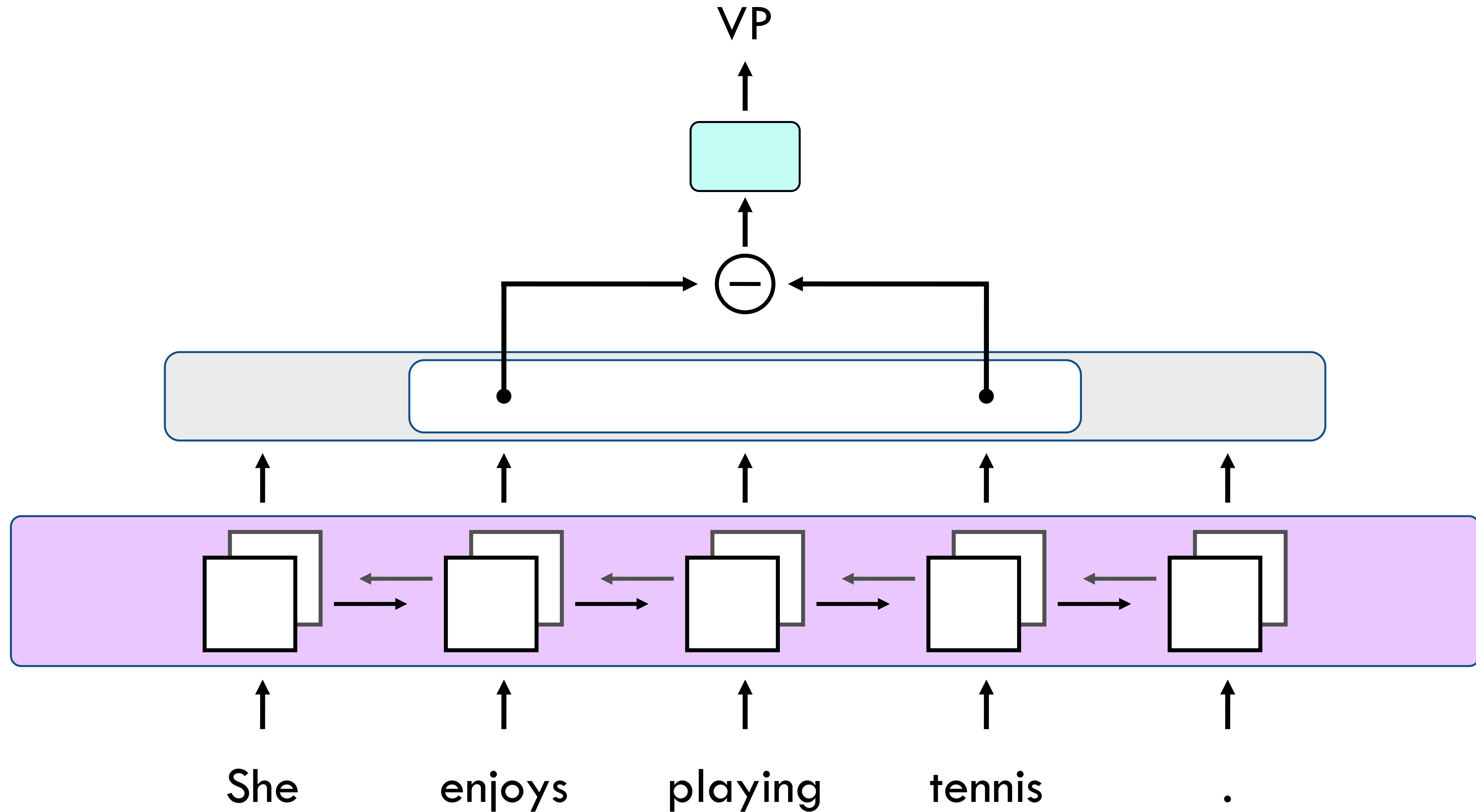


Span Classification



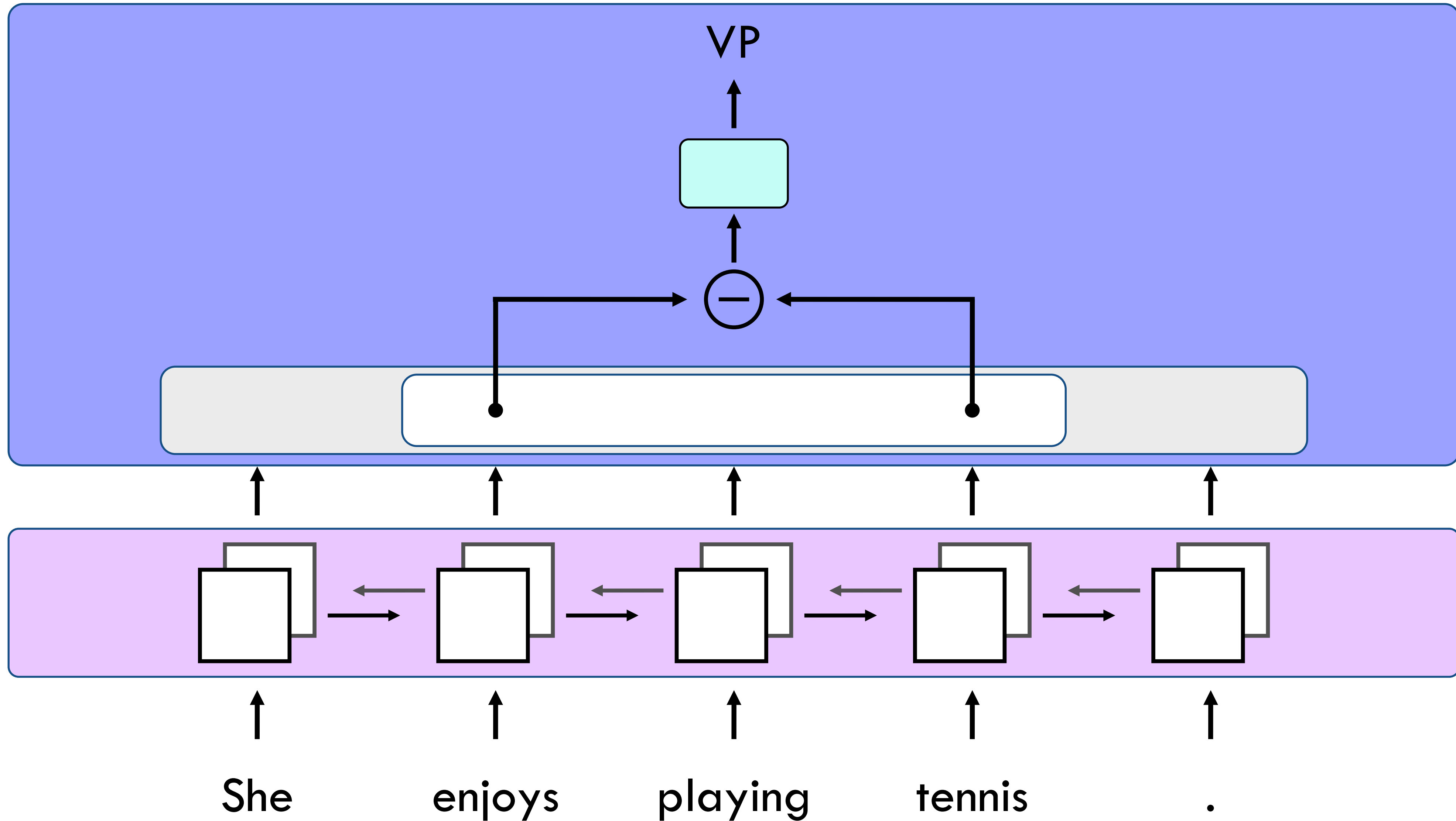


Span Classification



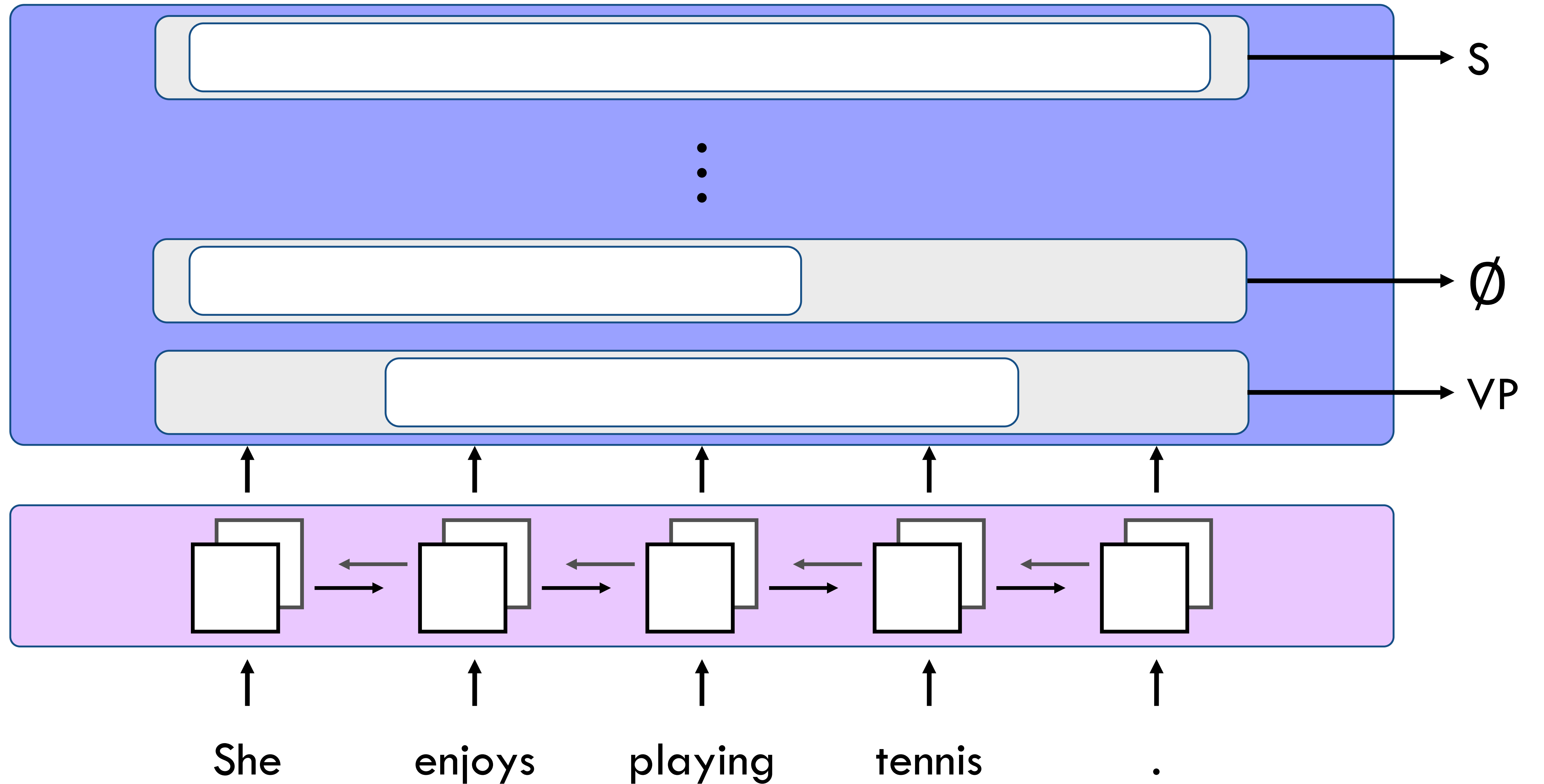


Span Classification



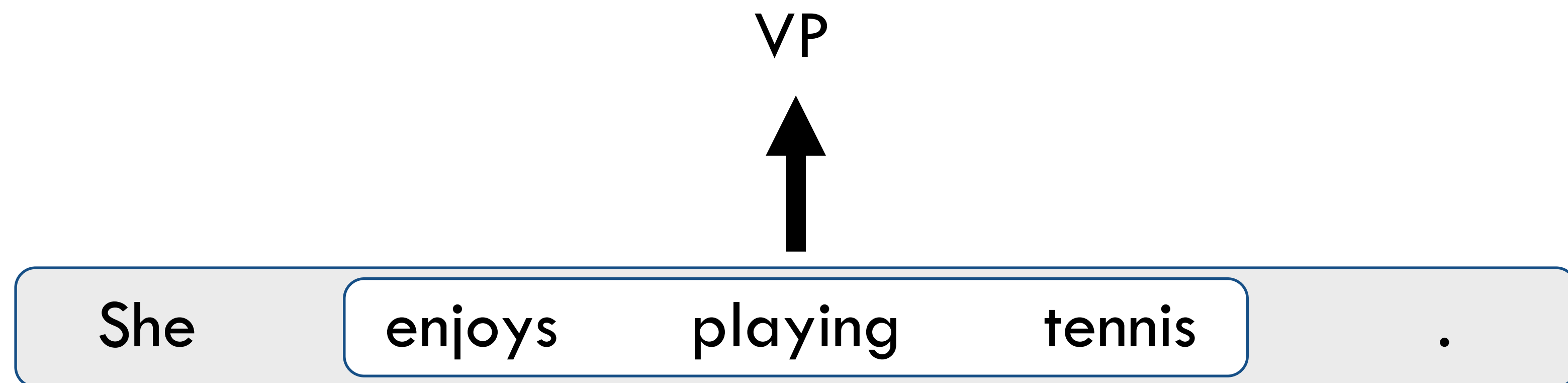
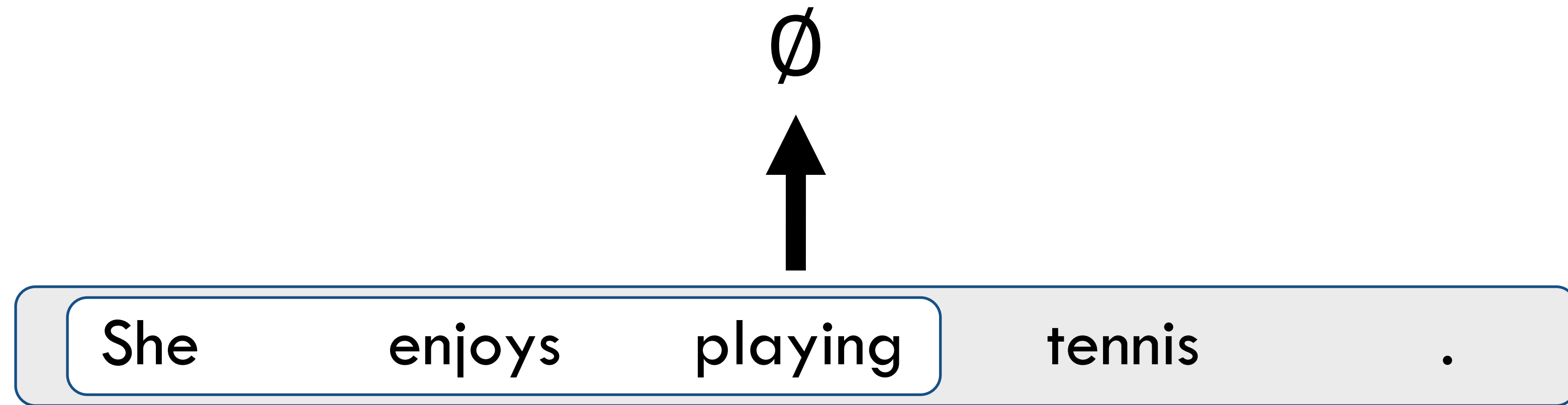


Span Classification



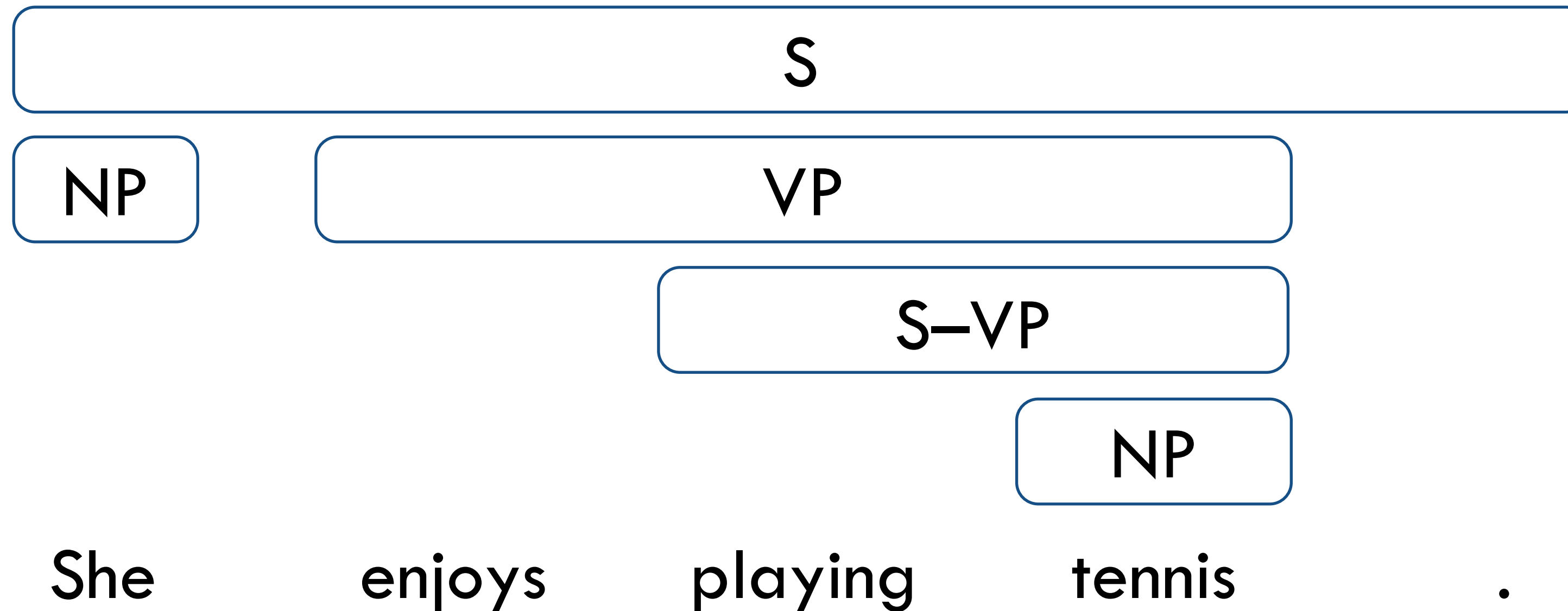
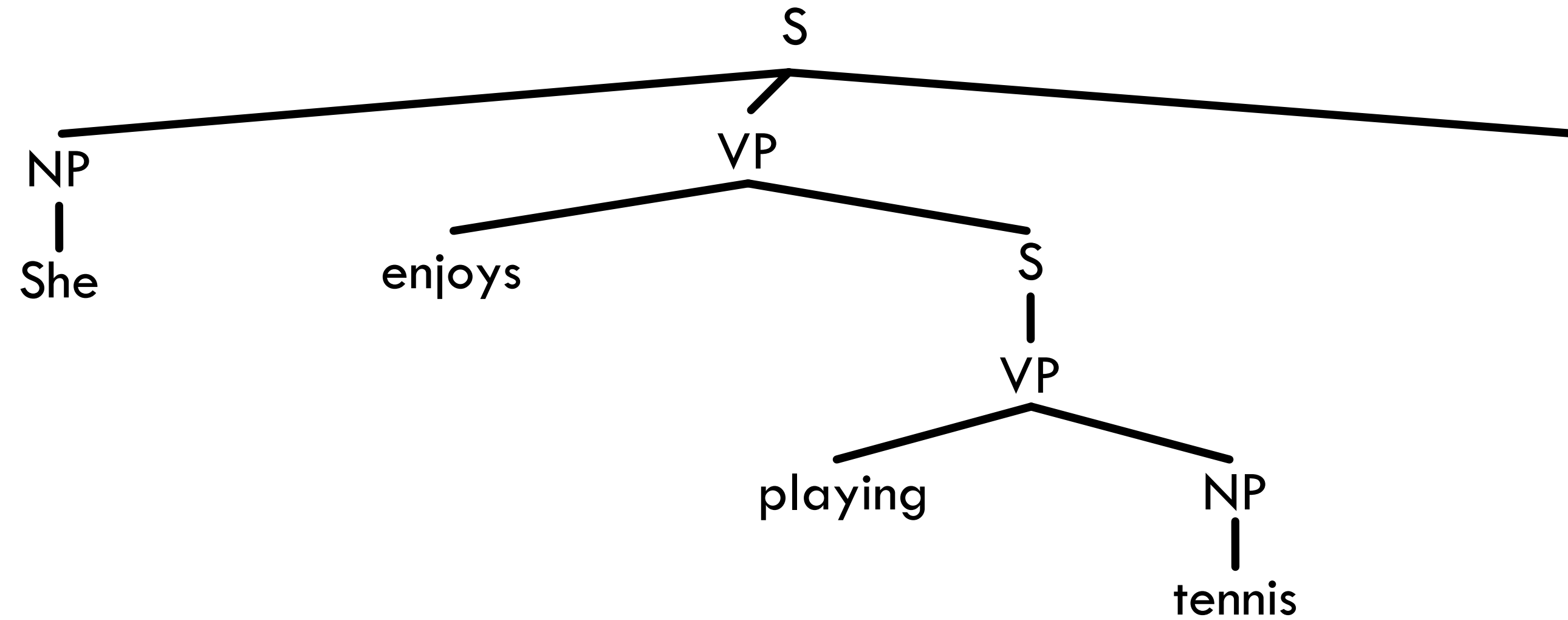


Non-Constituents



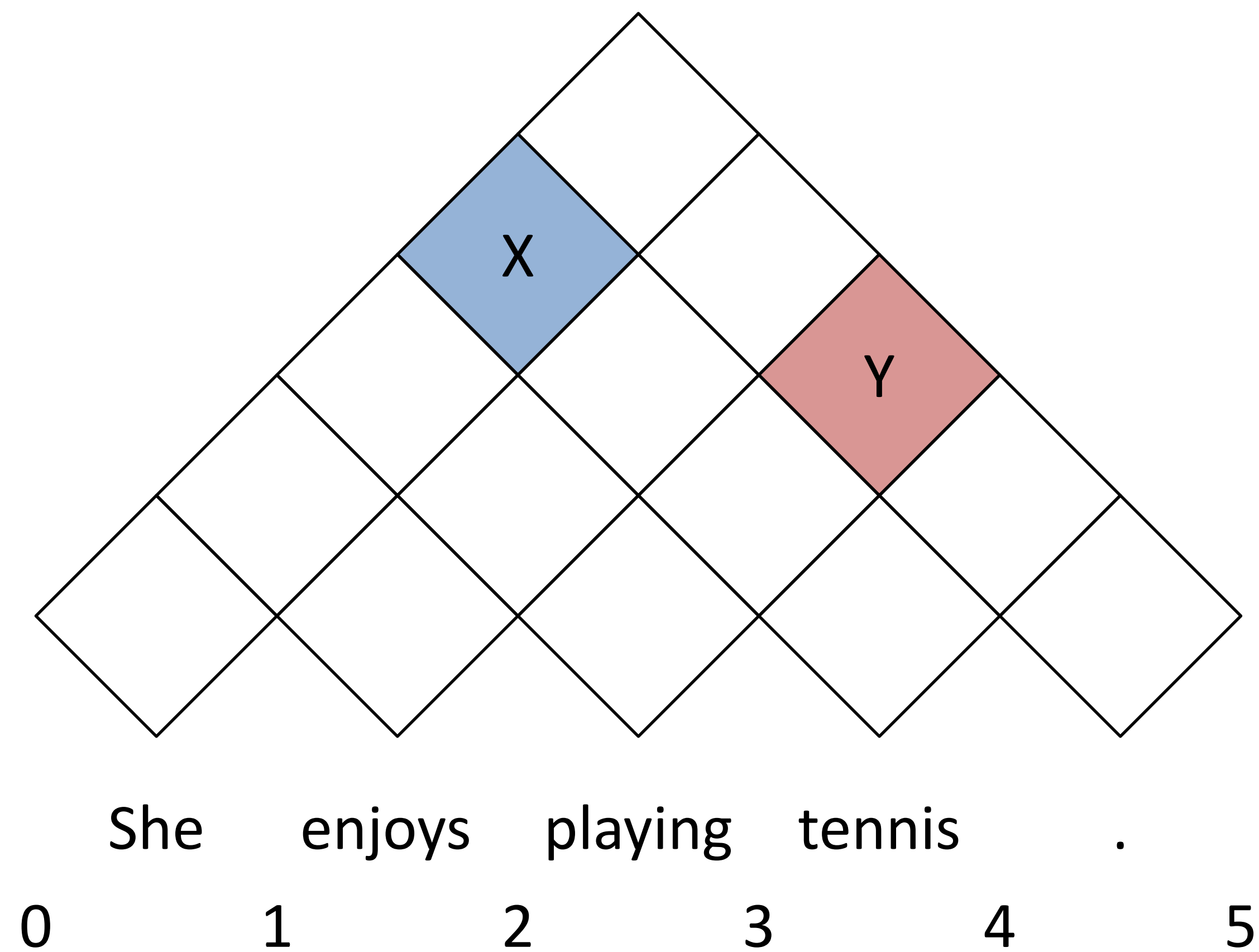


... But Will We Get a Tree Out?



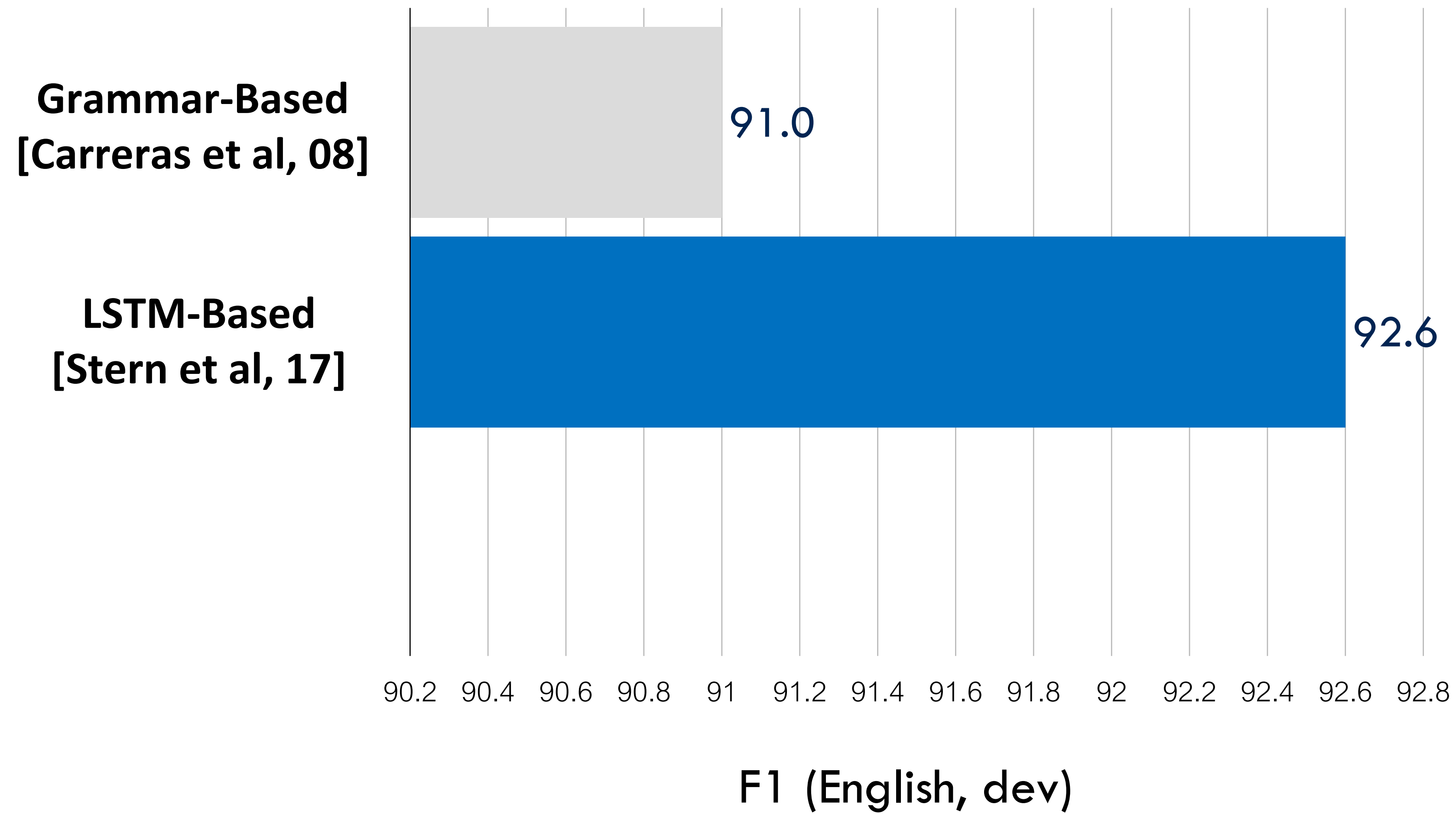


Reconciliation





Does It Work?





What's Going on in There?

Neural parsers no longer have much of the model structure provided to classical parsers.

How do they perform so well without it?



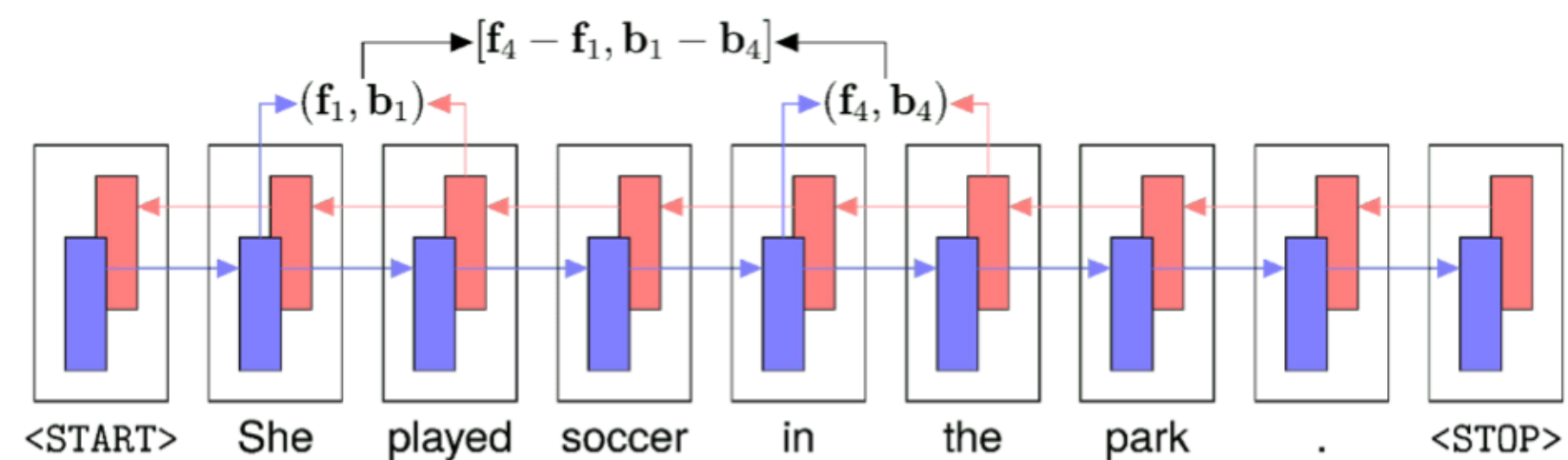
What's Going on in There?

Why don't we need a grammar?

Adjacent tree labels are redundant with LSTM features

If we can predict surrounding tree labels from our LSTM representation of the input, then this information doesn't need to be provided explicitly by grammar production rules

We find that for **92.3%** of spans, the label of the span's parent can be predicted from the neural representation of the span





What's Going on in There?

Do we need tree constraints?

Not for F1

Many neural parsers no longer model output correlations with grammar rules, but still use output correlations from tree constraints

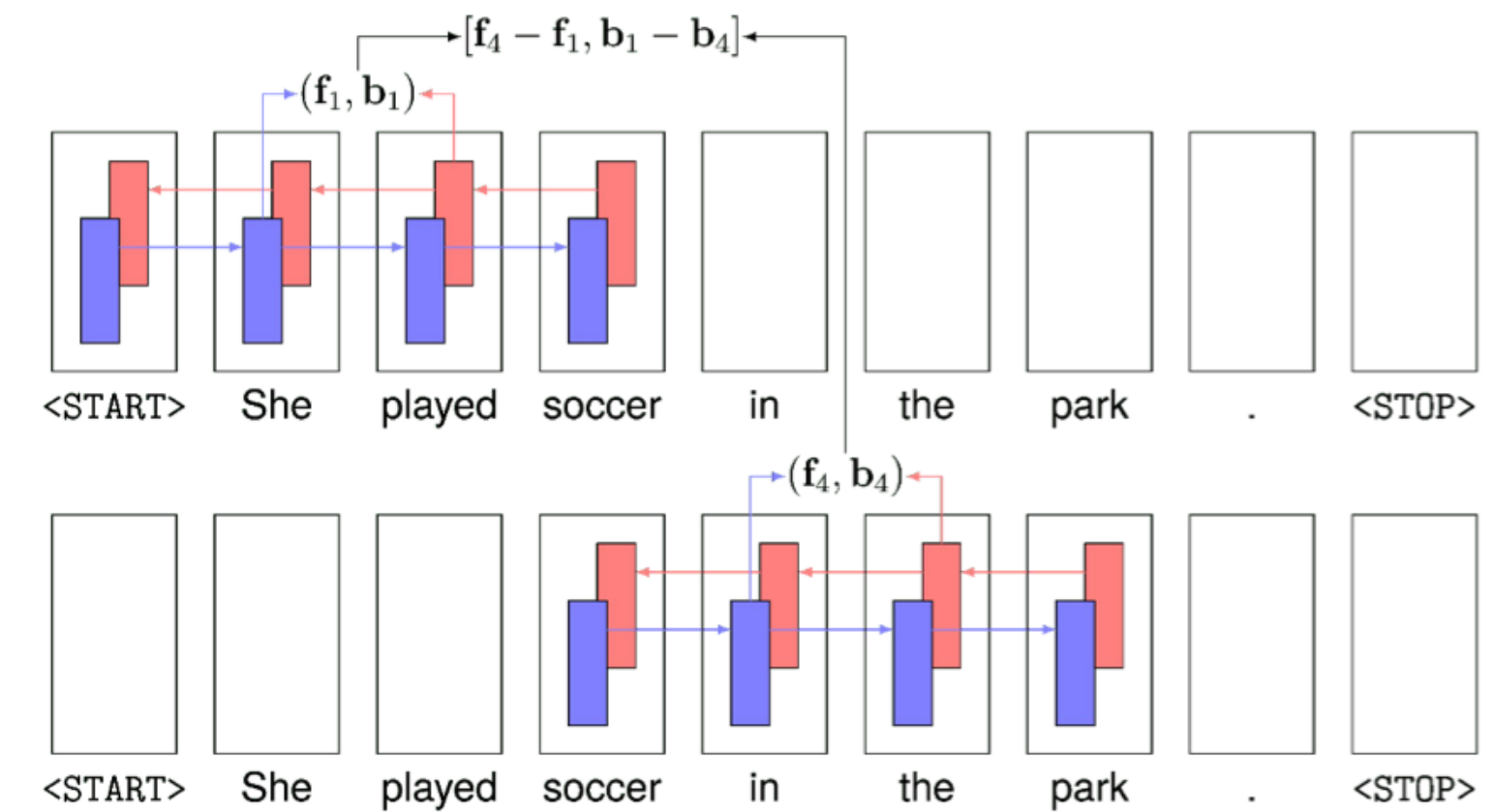
Predicting span brackets independently gives **nearly identical performance** on PTB development set F1 and produces valid trees for **94.5%** of sentences



What's Going on in There?

Is distant context important?

Yes!



Almost a full point of F1 is lost by truncating context 5 words away from span endpoints and half a point with 10 words



What's Going on in There?

What word representations do we need?

A character LSTM is sufficient

Word Only	91.44
Word and Tag	92.09
Character LSTM Only	92.24
Character LSTM and Word	92.22
Character LSTM, Word, and Tag	92.24



What's Going on in There?

What about lexicon features?

The character LSTM captures the same information

Heavily engineered lexicons used to be critical to good performance, but neural models typically don't use them

Word features from the Berkeley Parser (Petrov and Klein 2007) can be predicted with over **99.7%** accuracy from the character LSTM representation



What's Going on in There?

Do LSTMs introduce useful inductive bias compared to feedforward networks?

Yes!

We compare a truncated LSTM with feedforward architectures that are given the same inputs

The LSTM outperformed the best feedforward by **6.5 F1**



Routing with Transformers

Query:
verb

She

enjoys

playing

tennis

.



Routing with Transformers

Query:
verb

She

verb [VBZ]

enjoys

verb [VBG]

playing

noun

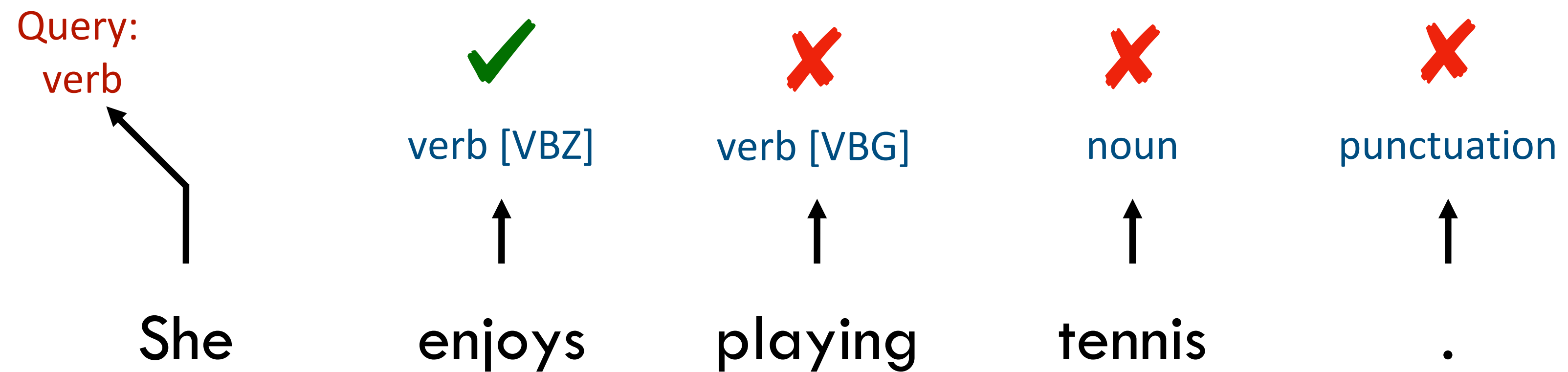
tennis

punctuation

.

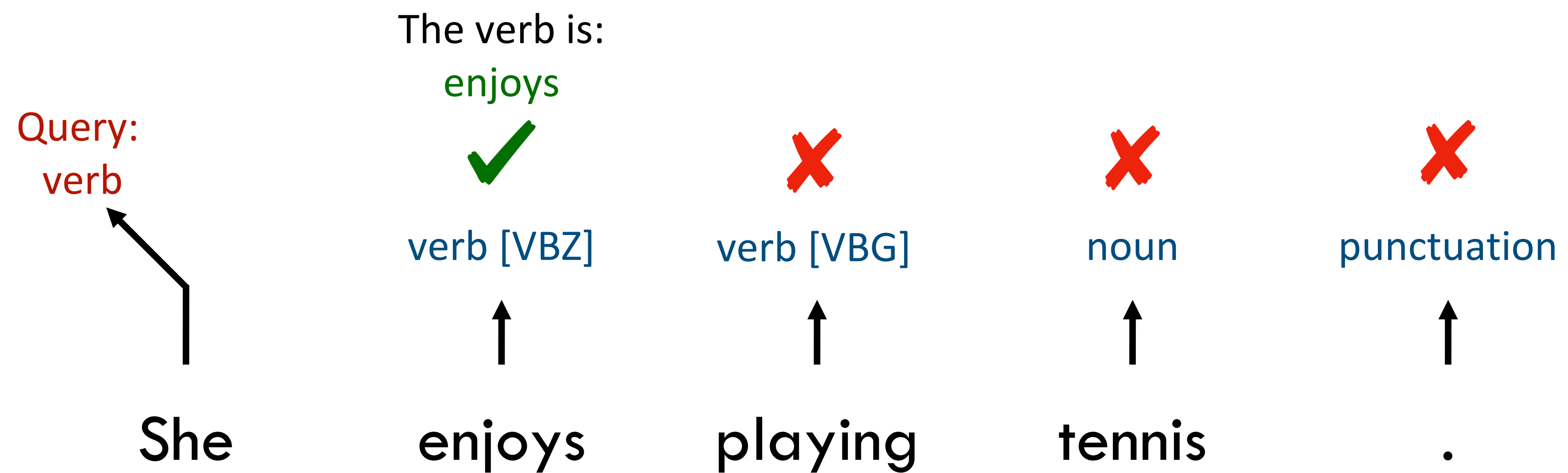


Routing with Transformers



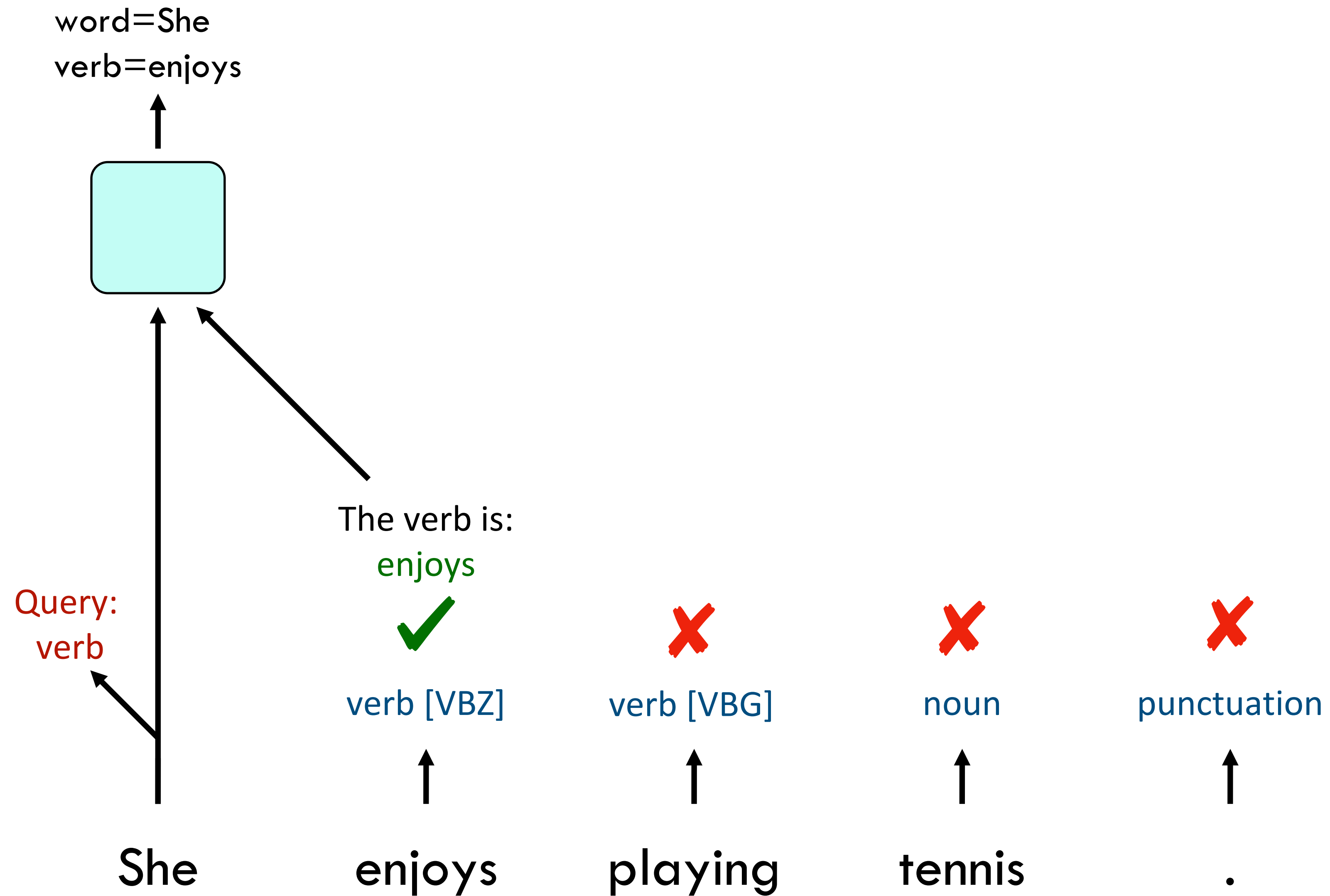


Routing with Transformers



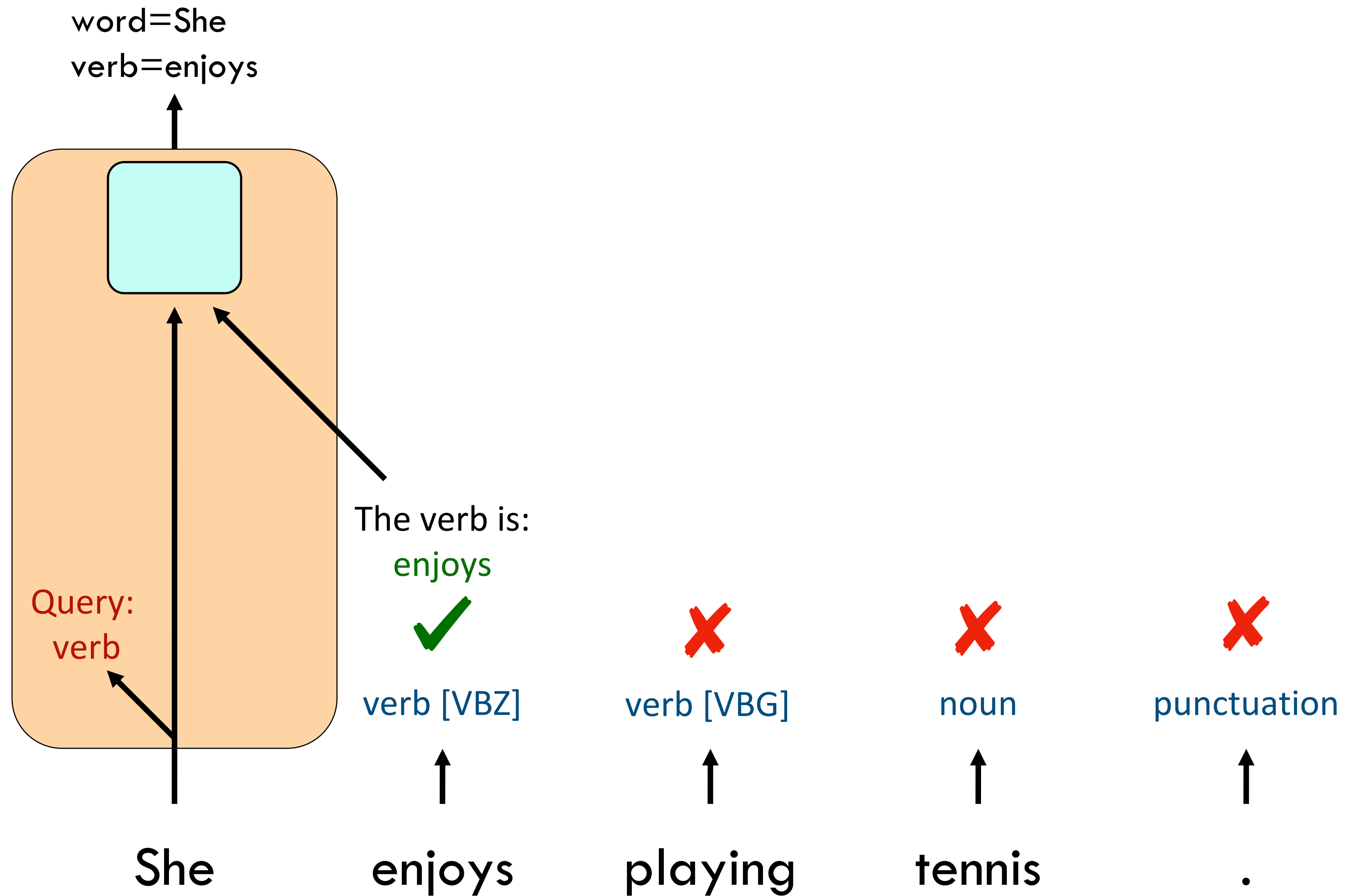


Routing with Transformers



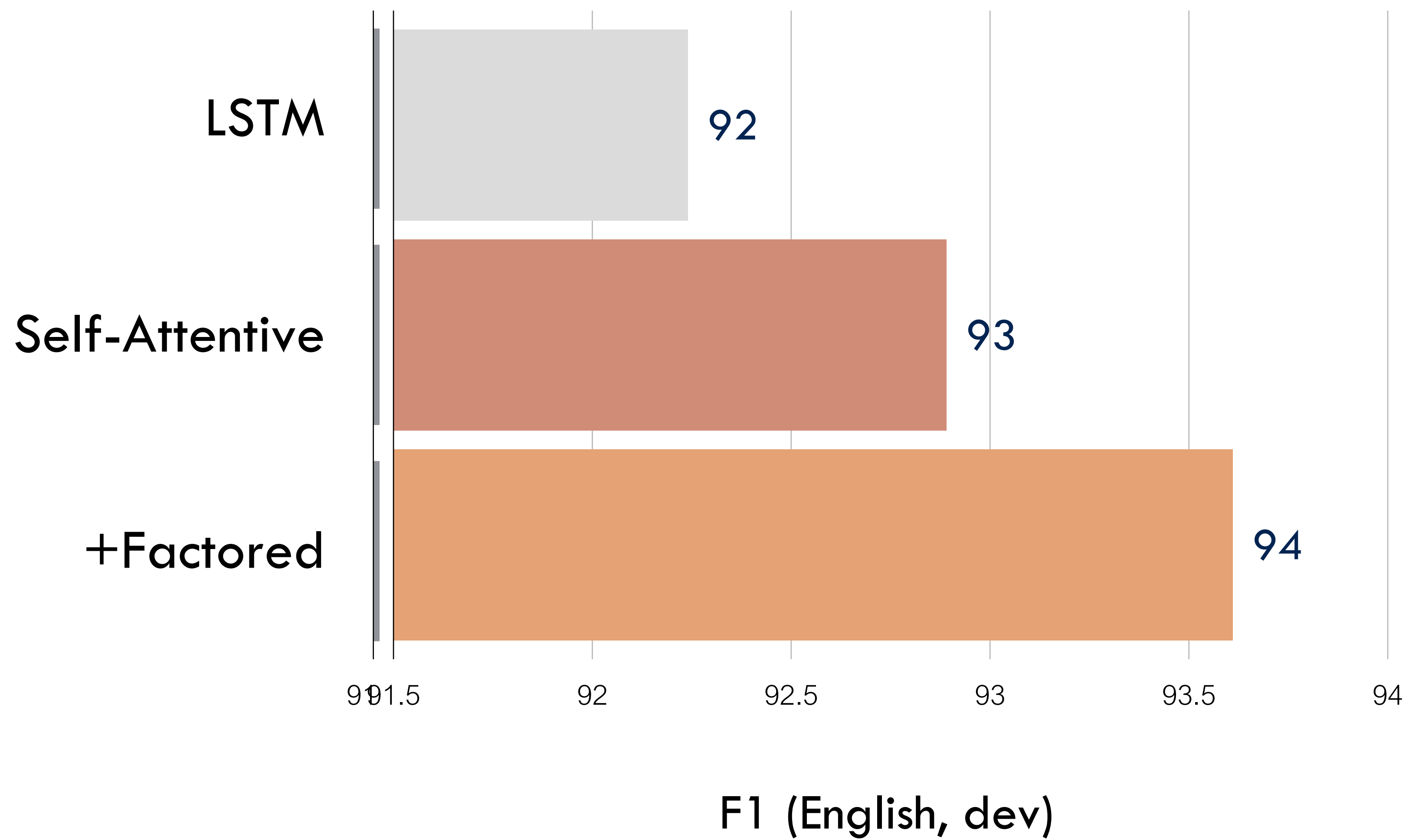


Routing with Transformers



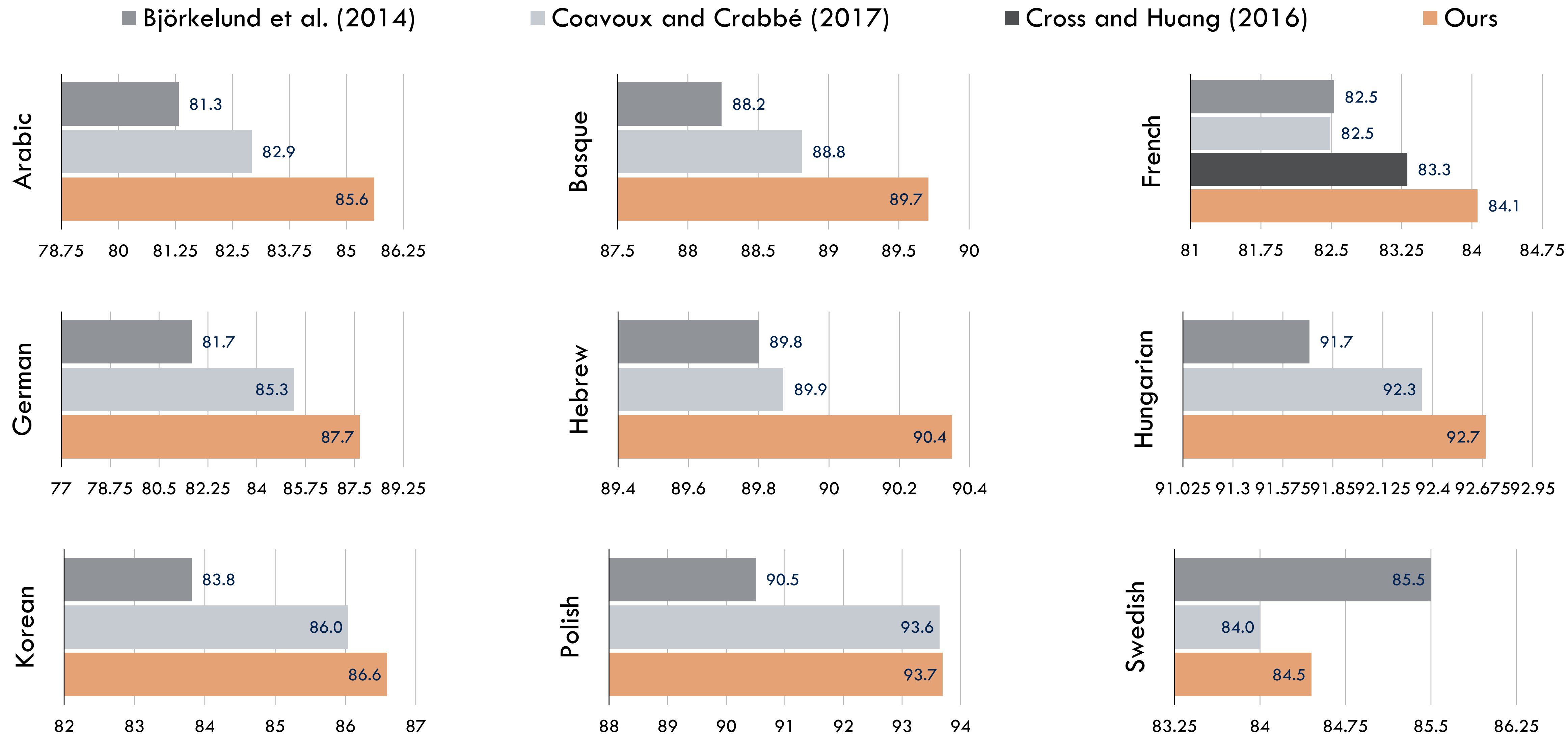


What Helps?





Results: Multilingual





Pre-Training

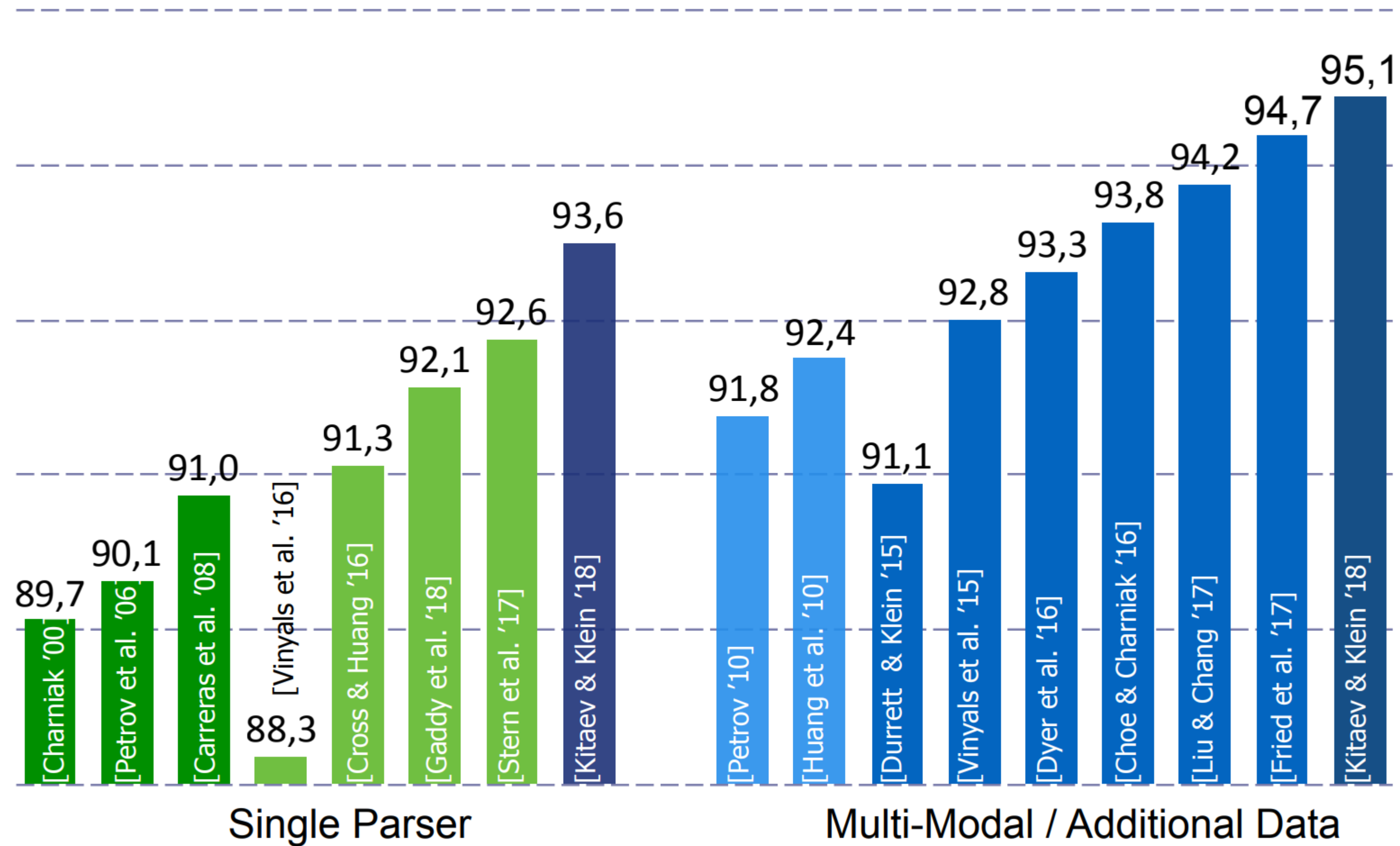
Problem: Input has more variation than output

Need to handle:

- Rare words not seen during training
- Word forms in morphologically rich languages
- Contextual paraphrase / lexical variation

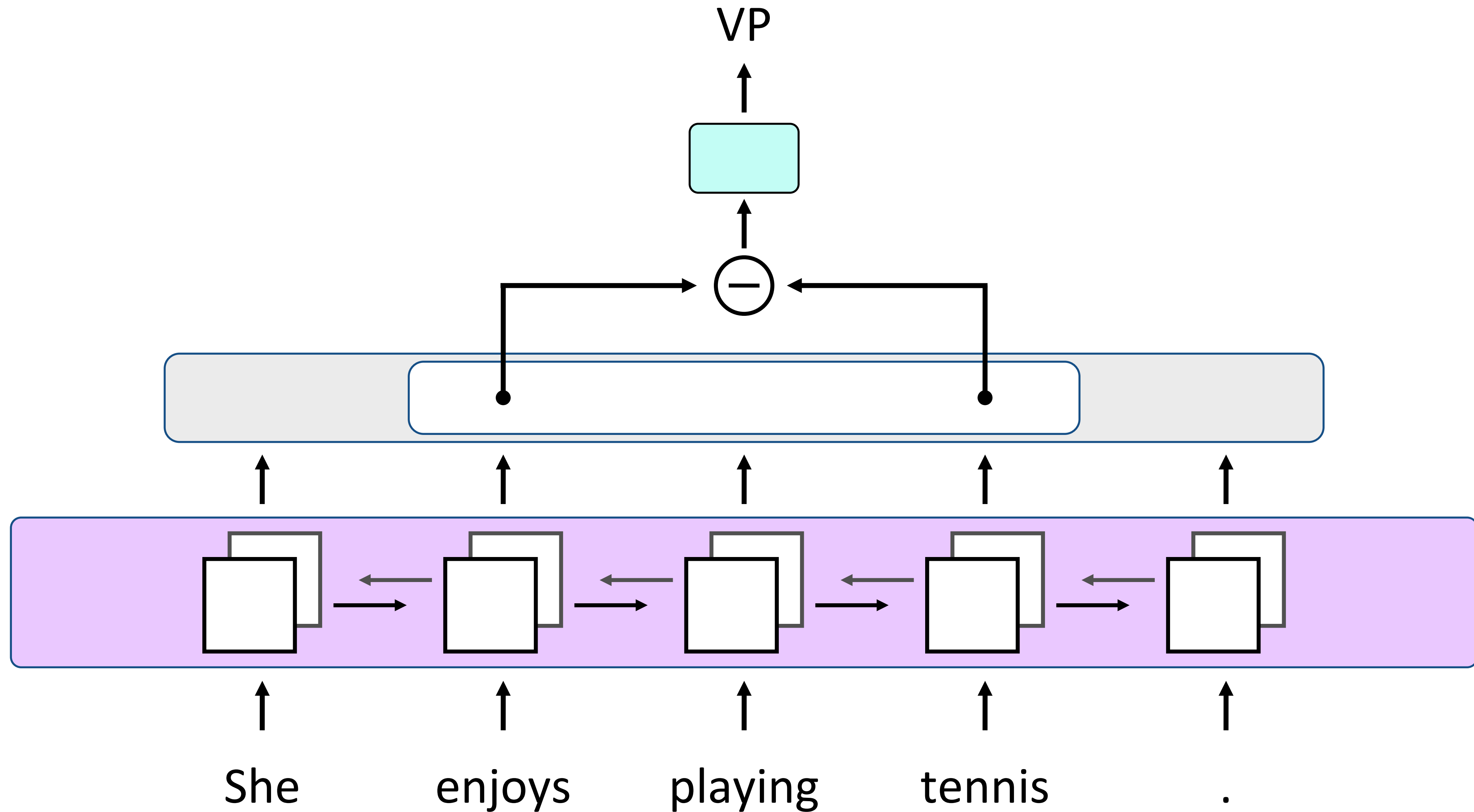


Historical Trends



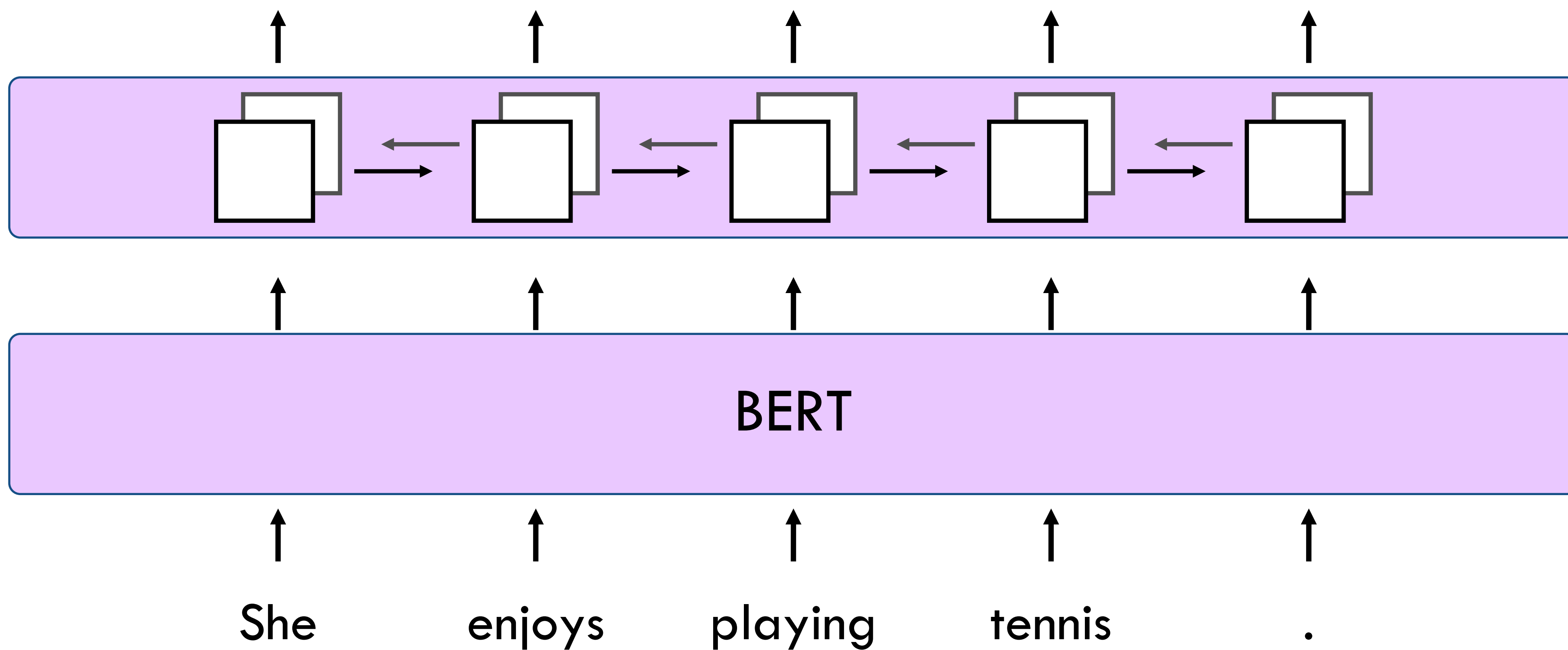


Parsing as Span Classification



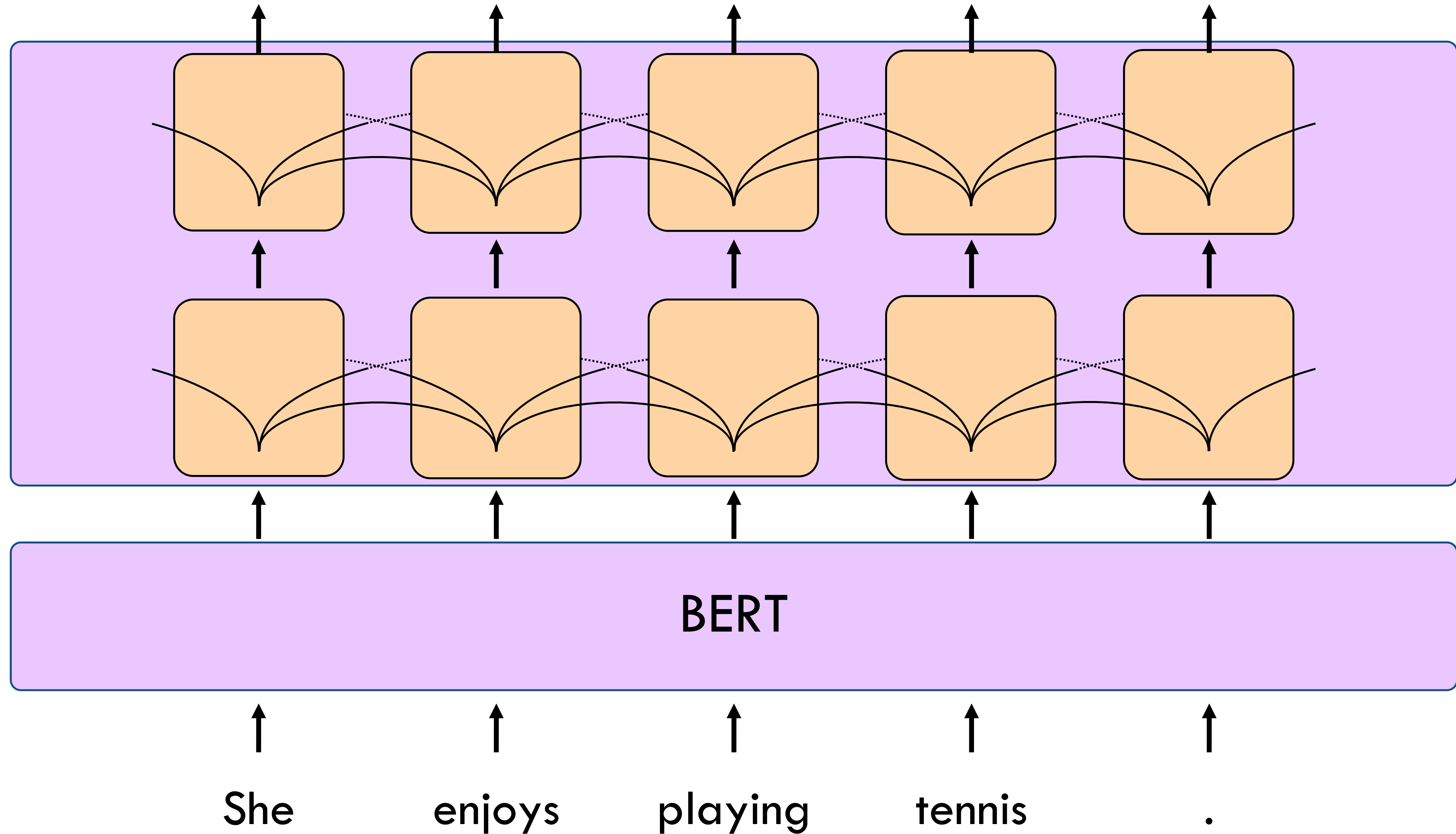


Pretraining





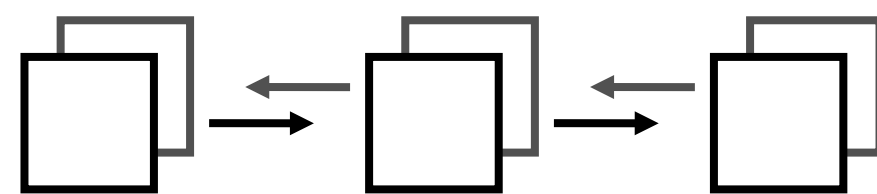
Architecture



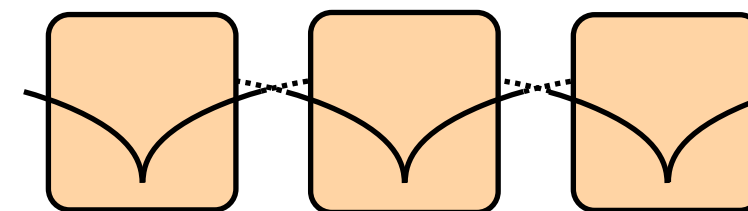


Encoder Architectures

LSTM



Self-Attention



No pre-training

92.08 F1

[Gaddy+ 2018]

93.55 F1

[Kitaev & Klein 2018]

Pre-training

95.13 F1
(with ELMo)

[Kitaev & Klein 2018]

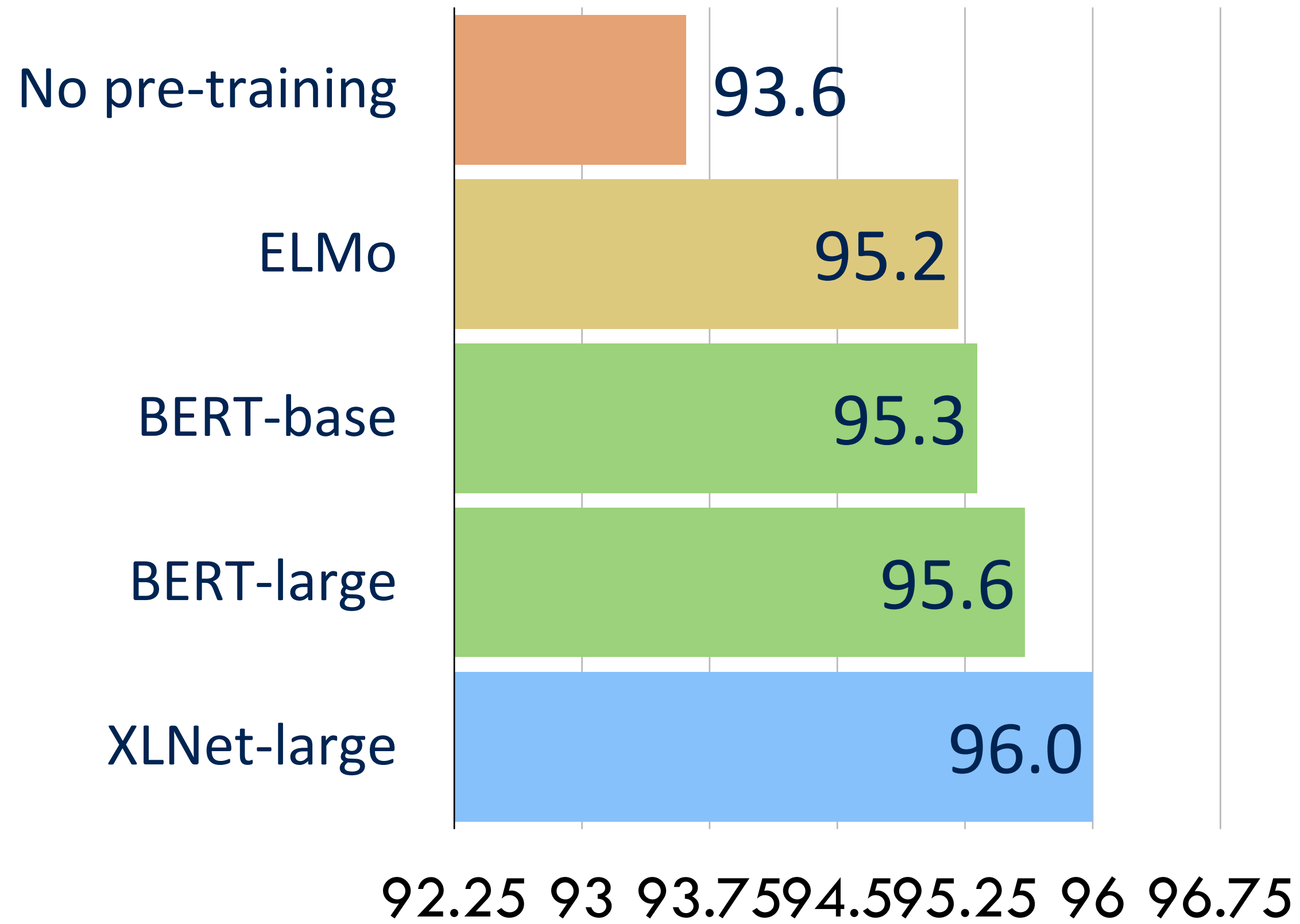
95.60 F1
(with BERT)

[Kitaev et al 2019]

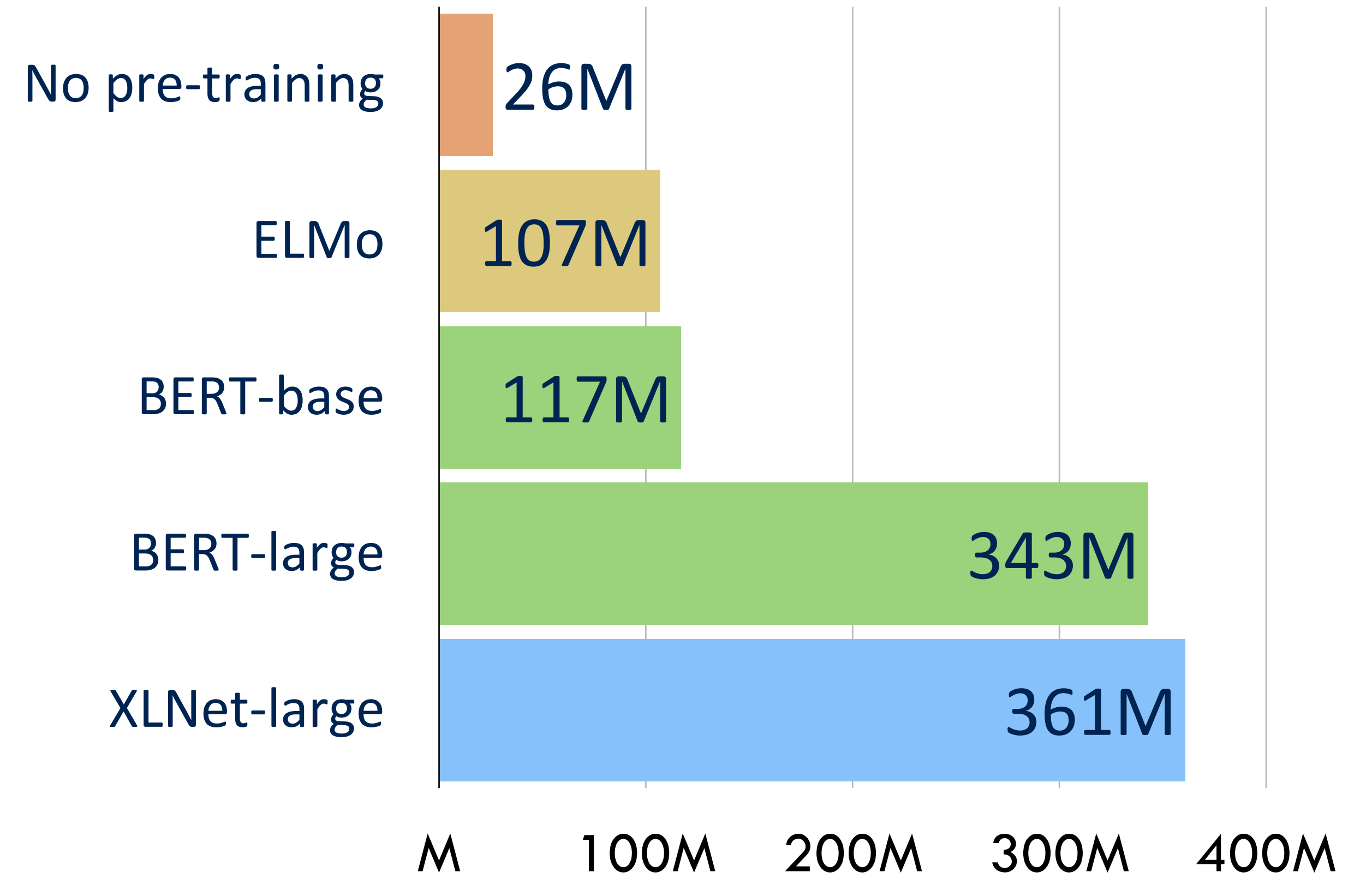


Encoder Architectures

F1 Score (English)

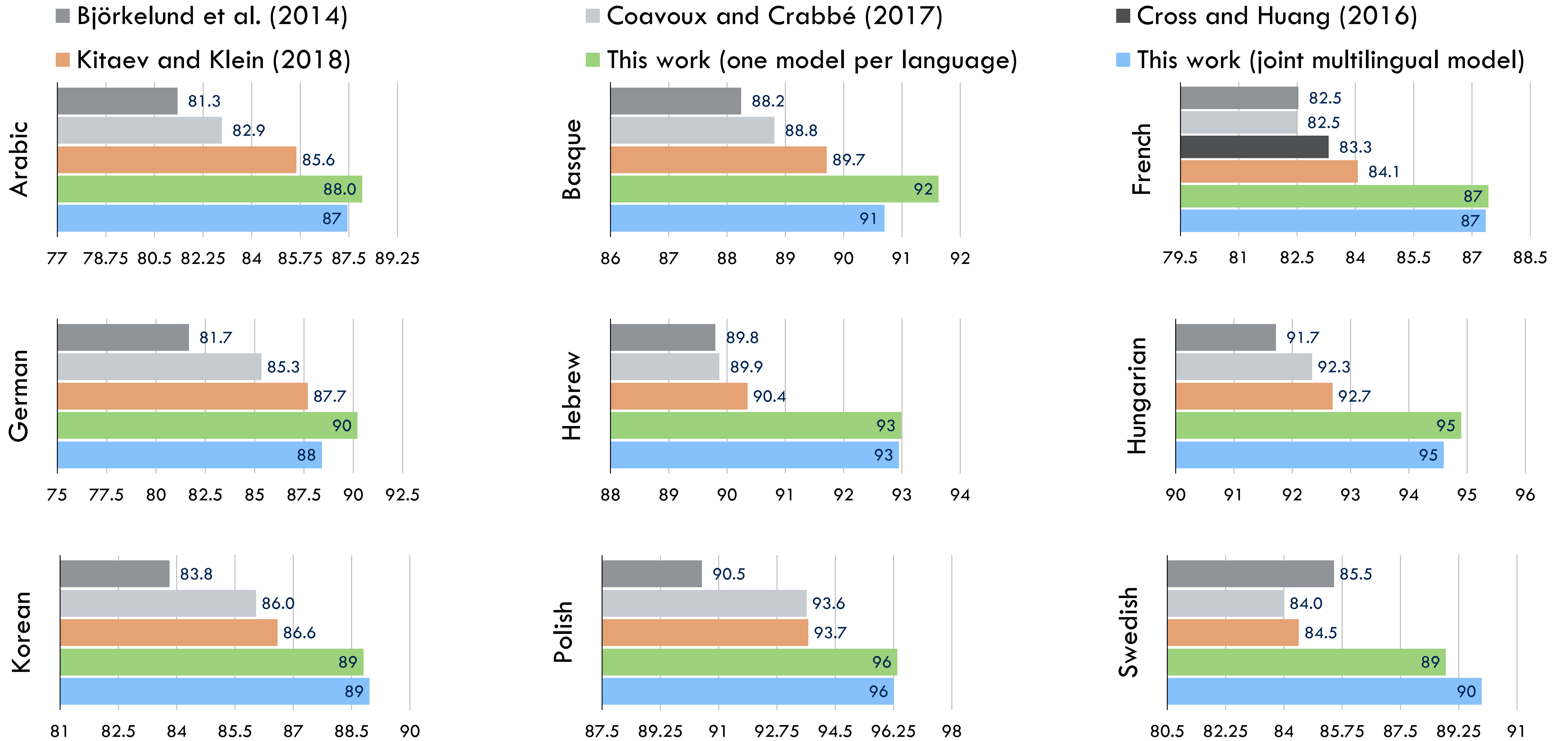


Number of Parameters





Results: Multilingual





Does Structure Help?

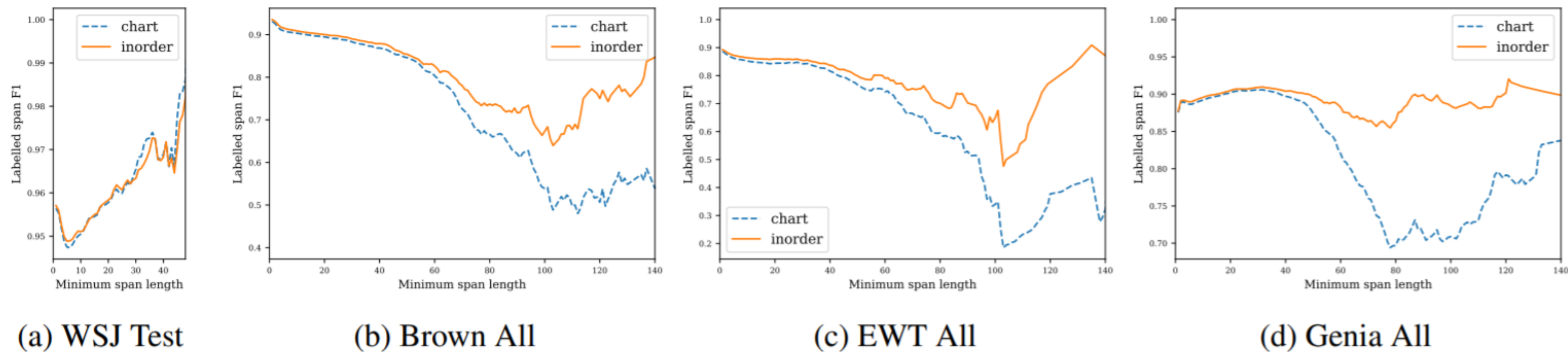


Figure 1: Labelled bracketing F1 versus minimum span length for the English corpora. F1 scores for the In-Order parser with BERT (orange) and the Chart parser with BERT (cyan) start to diverge for longer spans.



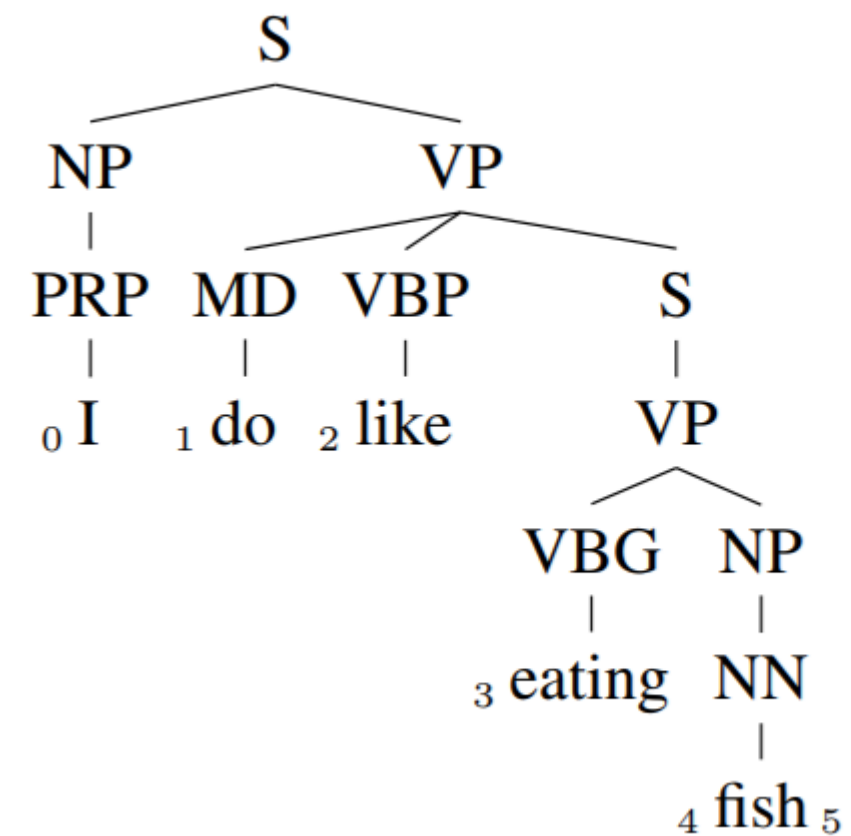
Out of Domain Parsing

	Berkeley		BLLIP		In-Order		Chart	
	F1	Δ Err.	F1	Δ Err.	F1	Δ Err.	F1	Δ Err.
WSJ Test	90.06	+0.0%	91.48	+0.0%	91.47	+0.0%	93.27	+0.0%
Brown All	84.64	+54.5%	85.89	+65.6%	85.60	+68.9%	88.04	+77.7%
Genia All	79.11	+110.2%	79.63	+139.1%	80.31	+130.9%	82.68	+157.4%
EWT All	77.38	+127.6%	79.91	+135.8%	79.07	+145.4%	82.22	+164.2%

Neural parsers improve out-of-domain numbers, but not more than in-domain numbers



Other Neural Constituency Parsers



steps	structural action	label action	stack after	bracket
1-2	sh(I/PRP)	label-NP	0△1	0NP ₁
3-4	sh(do/MD)	nolabel	0△1△2	
5-6	sh(like/VBP)	nolabel	0△1△2△3	
7-8	comb	nolabel	0△1△3	
9-10	sh(eating/VBG)	nolabel	0△1△3△4	
11-12	sh(fish/NN)	label-NP	0△1△3△4△5	4NP ₅
13-14	comb	label-S-VP	0△1△3△5	3S ₅ , 3VP ₅
15-16	comb	label-VP	0△1△5	1VP ₅
17-18	comb	label-S	0△5	0S ₅

- Back to at least Henderson 1998!
- Recent directions:
 - Shift-Reduce, eg Cross and Huang 2016
 - SR/Generative, eg Dyer et al 2016 (RNNG)
 - In-Order Generative, eg Liu and Zhang 2017