# Neural Constituency Parsing
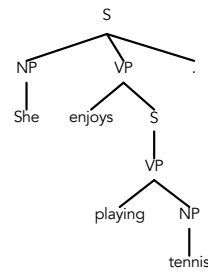
Berkeley
N L P

Dan Klein
CS 288

---

## Syntactic Parsing

*She enjoys playing tennis.*

3

---

## Syntactic Parsing



4

---

## Historical Trends



Single Parser          Multi-Modal / Additional Data

5

## Output Correlations

```
              S
         /    |    \
       NP     VP     .
        |    /  \
      She enjoys  S
                   |
                   VP
                  /  \
            playing   NP
                       |
                     tennis
```

6

## Grammars

S → NP VP

```
              S
         /    |    \
       NP     VP     .
        |    /  \
      She enjoys  S
                   |
                   VP
                  /  \
            playing   NP
                       |
                     tennis
```

VP[*enjoys*] : S[*playing*]

NP^S → *she*

7

## Input-Output Correlations

*She enjoys playing tennis.*

8

## Span-Based Parsing

```
                 S
        /     /   \      \
      NP     VP          .
       |    /  \
     She enjoys  S
                  |
                  VP
                 /  \
           playing   NP
                      |
                    tennis
```

| S |
|---|
| NP | VP |
| S–VP |
| NP |

She   enjoys   playing   tennis   .

9

## Parsing as Span Classification

VP

She enjoys playing tennis .

10

## Routing with LSTMs

She enjoys playing tennis .

11

## Routing with LSTMs
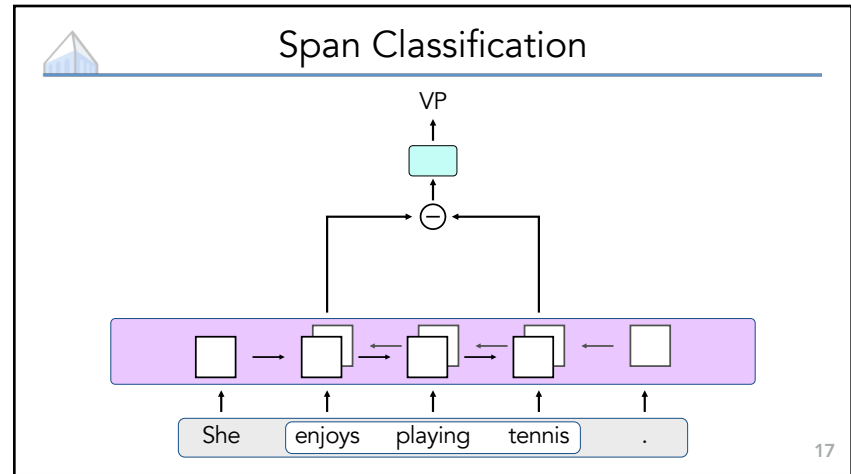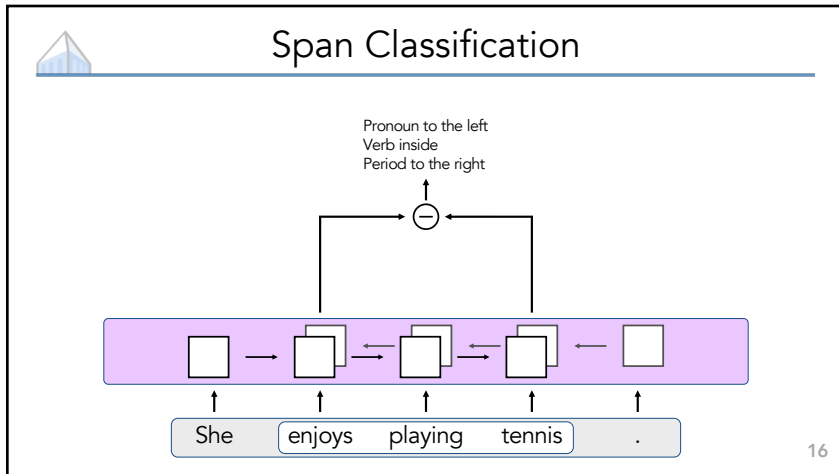
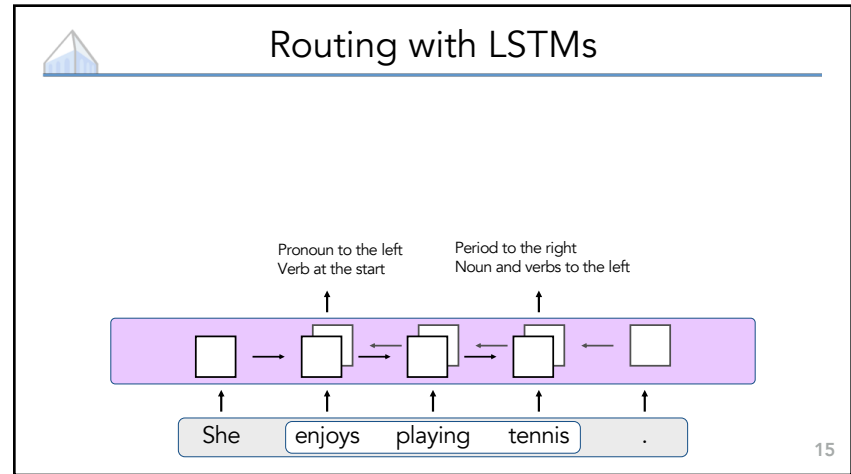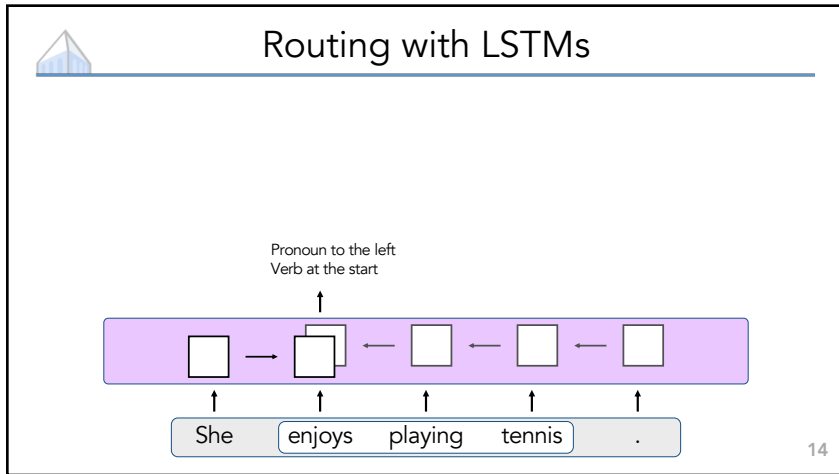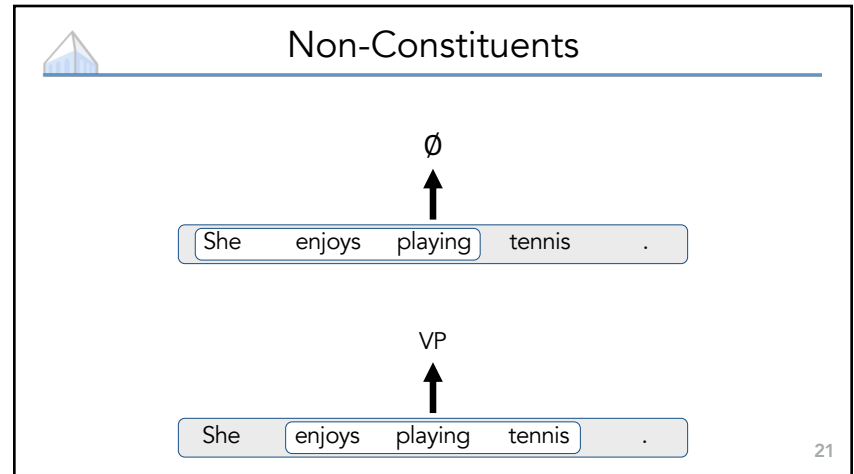Pronoun to the left

She enjoys playing tennis .

12

## Routing with LSTMs

Verb at the start

She enjoys playing tennis .

13

## Routing with LSTMs

Pronoun to the left
Verb at the start

She  enjoys  playing  tennis  .

14

## Routing with LSTMs

Pronoun to the left        Period to the right
Verb at the start          Noun and verbs to the left

She  enjoys  playing  tennis  .

15

## Span Classification

Pronoun to the left
Verb inside
Period to the right

She  enjoys  playing  tennis  .

16

## Span Classification

VP

She  enjoys  playing  tennis  .

17

4

## Span Classification

VP

She    enjoys    playing    tennis    .

18

## Span Classification

VP

She    enjoys    playing    tennis    .

19

## Span Classification

S

⋮

∅

VP

She    enjoys    playing    tennis    .

20

## Non-Constituents

∅

She    enjoys    playing    tennis    .

VP

She    enjoys    playing    tennis    .

21

5

## … But Will We Get a Tree Out?



## Reconciliation



## Does It Work?



F1 (English, dev)

## 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 predicted from the neural representation of the span



26

## 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

27

## 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

28

## 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 |

29

7

## 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

30

## 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**

31

## Routing with Transformers

**Query:**
**verb**

She    enjoys    playing    tennis    .
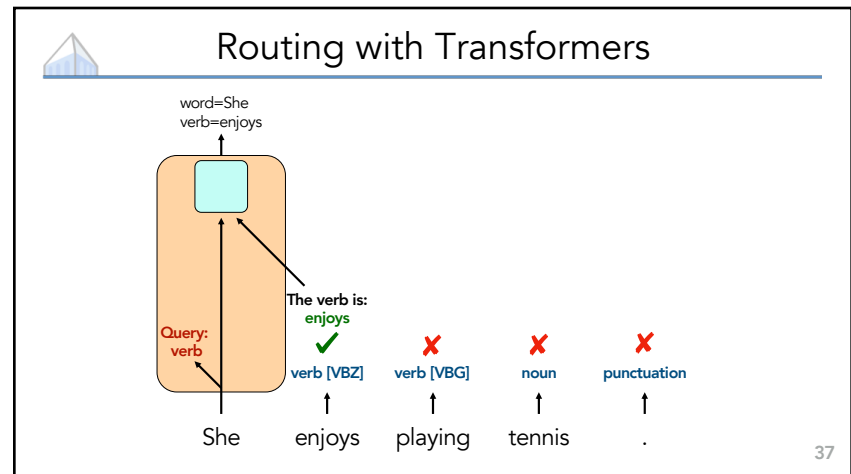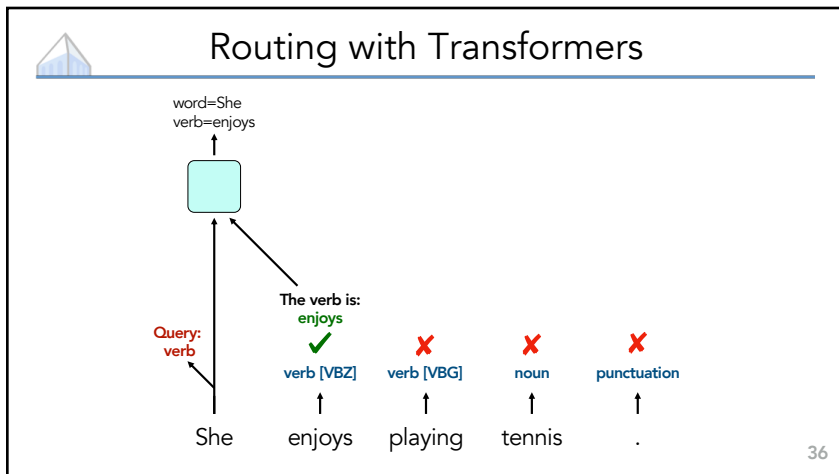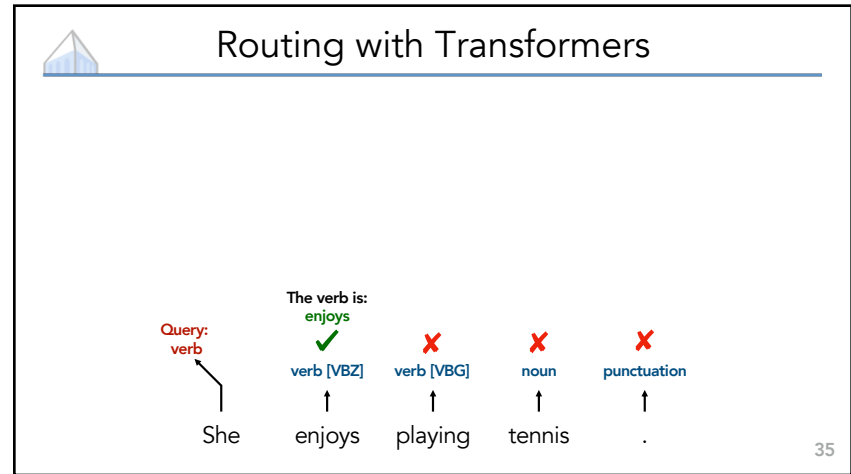
32

## Routing with Transformers

**Query:**
**verb**

verb [VBZ]    verb [VBG]    noun    punctuation

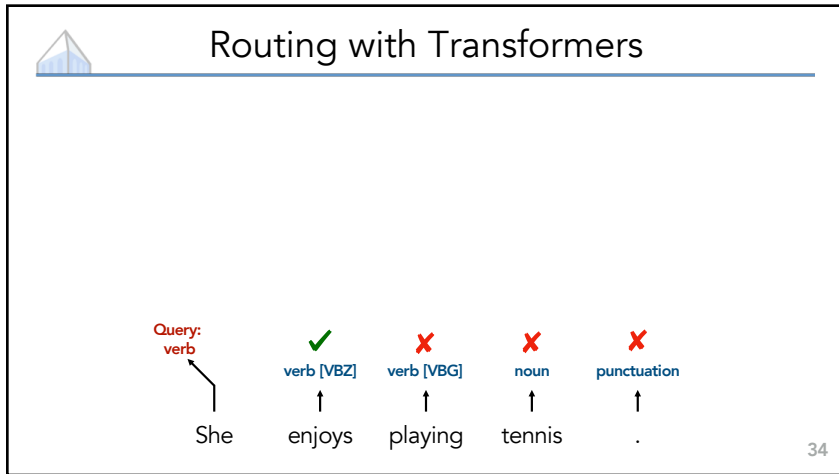She    enjoys    playing    tennis    .

33

Routing with Transformers

Query: verb

✔ verb [VBZ]    ✘ verb [VBG]    ✘ noun    ✘ punctuation

She    enjoys    playing    tennis    .

34



Routing with Transformers

The verb is: enjoys

Query: verb

✔ verb [VBZ]    ✘ verb [VBG]    ✘ noun    ✘ punctuation

She    enjoys    playing    tennis    .

35



Routing with Transformers

word=She
verb=enjoys

The verb is: enjoys

Query: verb

✔ verb [VBZ]    ✘ verb [VBG]    ✘ noun    ✘ punctuation

She    enjoys    playing    tennis    .

36



Routing with Transformers

word=She
verb=enjoys

The verb is: enjoys

Query: verb

✔ verb [VBZ]    ✘ verb [VBG]    ✘ noun    ✘ punctuation

She    enjoys    playing    tennis    .

37

## What Helps?



LSTM — 92
Self-Attentive — 93
+Factored — 94

F1 (English, dev)

38

## Results: Multilingual



Björkelund et al. (2014)   Coavoux and Crabbé (2017)   Cross and Huang (2016)   Ours

Arabic: 81.3 / 82.9 / 85.6
Basque: 88.2 / 88.8 / 89.7
French: 82.5 / 82.5 / 83.3 / 84.1
German: 81.7 / 85.3 / 87.7
Hebrew: 89.8 / 89.9 / 90.4
Hungarian: 91.7 / 92.3 / 92.7
Korean: 83.8 / 86.0 / 86.6
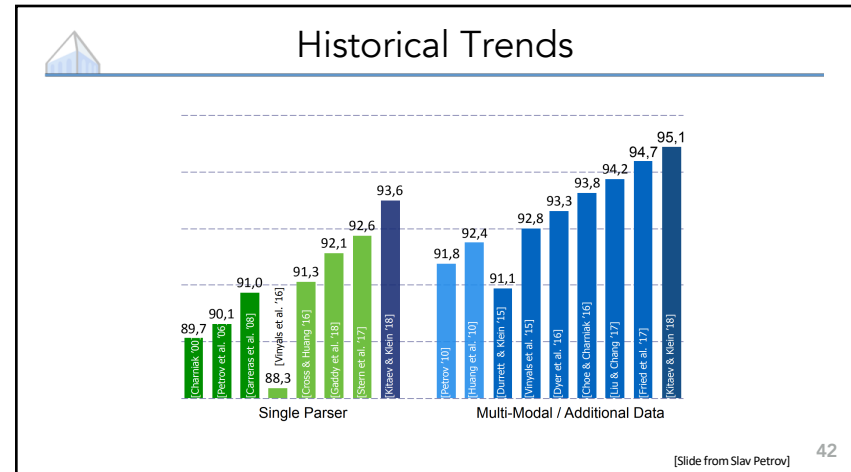Polish: 90.5 / 93.6 / 93.7
Swedish: 85.5 / 84.0 / 84.5

39

## Data Hunger

Problem: Input has more variation than output

Need to handle:
- Rare words not seen during training
- Word forms in morphologically rich languages

40

## Historical Trends



Single Parser: [Charniak '00] 89,7 / [Petrov et al. '06] 90,1 / [Carreras et al. '08] 88,3 / [Vinyals et al. '16] 91,0 / [Cross & Huang '16] 91,3 / [Gaddy et al. '18] 92,1 / [Stern et al. '17] 92,6 / [Kitaev & Klein '18] 93,6

Multi-Modal / Additional Data: [Petrov '10] 91,8 / [Huang et al. '10] 92,4 / [Durrett & Klein '15] 91,1 / [Vinyals et al. '15] 92,8 / [Dyer et al. '16] 93,3 / [Choe & Charniak '16] 93,8 / [Liu & Zhang '17] 94,2 / [Fried et al. '17] 94,7 / [Kitaev & Klein '18] 95,1

[Slide from Slav Petrov]
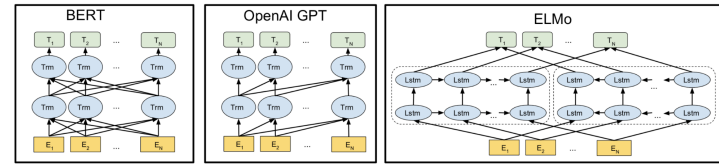
42

## Knowledge Modularity

- Knowledge modularity: Learn domain-general knowledge from one data source and use it solve specific problems elsewhere



## Context Embeddings and Pretraining

**Key Idea:** Embed contexts, not words. Use these embeddings for other tasks.

Example: BERT (Devlin et al., 2019) -- bidirectional Transformer trained on masked language modeling and next-sentence prediction



## Recent Explosion of Pretraining Work



## Parsing as Span Classification



She  enjoys  playing  tennis  .

47

**11**

## Pretraining



She  enjoys  playing  tennis  .

48

## Architecture



She  enjoys  playing  tennis  .

49

## 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) |

50

## Encoder Architectures

F1 Score (English)

| | |
|---|---|
| No pre-training | 93.6 |
| ELMo | 95.2 |
| BERT-base | 95.3 |
| BERT-large | 95.6 |
| XLNet-large | 96.0 |

92.25 93 93.75 94.5 95.25 96 96.75

Number of Parameters

| | |
|---|---|
| No pre-training | 26M |
| ELMo | 107M |
| BERT-base | 117M |
| BERT-large | 343M |
| XLNet-large | 361M |

M  100M  200M  300M  400M

51

# Results: Multilingual



Legend: Björkelund et al. (2014) · Coavoux and Crabbé (2017) · Cross and Huang (2016) · Kitaev and Klein (2018) · This work (one model per language) · This work (joint multilingual model)

52

---

# Does Structure Help?



(a) WSJ Test  (b) Brown All  (c) EWT All  (d) Genia All

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 | $\Delta$ Err. | F1 | $\Delta$ Err. | F1 | $\Delta$ Err. | F1 | $\Delta$ 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



- 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

Thank You!

Berkeley

N  L  P

nlp.cs.berkeley.edu