# Deep Learning for Broad Coverage Semantics: SRL, Coreference, and Beyond

## Luke Zettlemoyer[†*]

Joint work with **Luheng He**[†], **Kenton Lee**[†], **Matthew Peters***, Christopher Clark[†], Matthew Gardner*, Mohit Iyyer*, Mandar Joshi[†], Mike Lewis[‡], Julian Michael[†], Mark Neumann*

[†] Paul G. Allen School of Computer Science & Engineering, University of Washington,
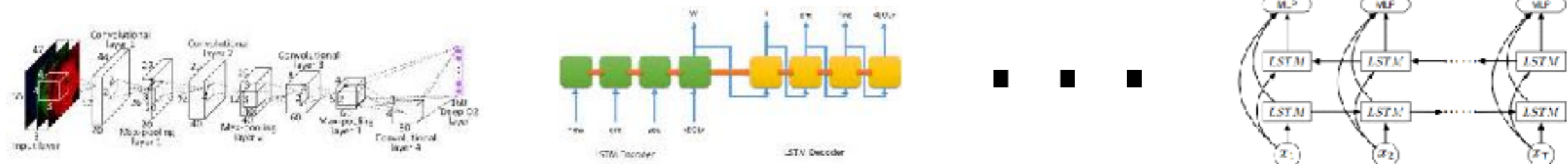[‡] Facebook AI Research
[*] Allen Institute for Artificial Intelligence

# Three Simple Steps that will Revolutionize Your ML Research

*Step 1: Gather lots of training data!*



*Step 2: Apply Deep Learning!!*



*Step 3: Observe Impressive Gains!!!*

# Broad Coverage Semantics

*Example Tasks:*

Coreference: clustering NPs

A fire in a Bangladeshi garment factory has left at least 37 people dead and 100 hospitalized. Most of the deceased were killed in the crush as workers tried to flee the blaze in the four-story building.

Semantic Role Labeling: who did what, etc.

| | |
|---|---|
| ARG0 | NASA |
| PRED | *observe* |
| ARG1 | an X-ray flare 400 times brighter than usual |
| TMP | On January 5, 2015 |

*Many applications:*

Question Answering

Information Extraction

Machine Translation

# Does the Recipe Work for Broad Coverage Semantics?

*Step 1: Gather lots of training data!*

**Challenge 1: Data is costly and limited
(e.g. linguists required to label
PennTreebank / OntoNotes)**

*Step 2: Apply Deep Learning!!*

**Challenge 2: Pipeline of structured
prediction problems with cascading errors
(e.g. POS->Parsing->SRL->Coref)**

*Step 3: Observe Impressive Gains!!!*

# New Learning Approaches

*New state-of-the-art results for two tasks:*

**Coreference:**

A fire in a Bangladeshi garment factory has left at least 37 people dead and 100 hospitalized. Most of the deceased were killed in the crush as workers tried to flee the blaze in the four-story building.

**Semantic Role Labeling:**

| | |
|---|---|
| ARG0 | NASA |
| PRED | *observe* |
| ARG1 | an X-ray flare 400 times brighter than usual |
| TMP | On January 5, 2015 |

*Common themes:*
- End-to-end training of deep neural networks
- No preprocessing (e.g., no POS, no parser, etc.)
- Large gains in accuracy with simpler models and no extra training data

# Semantic Role Labeling (SRL)

role label

who what when where why ...

*predicate*  →  argument

**subj**  **v**  **obj**  **prep**

The robot *broke* my favorite mug with a wrench.

breaker
ARG0

thing broken
ARG1

instrument
ARG2

**subj**  **v**  **prep**  **adv**

My mug *broke* into pieces immediately.

thing broken
ARG1

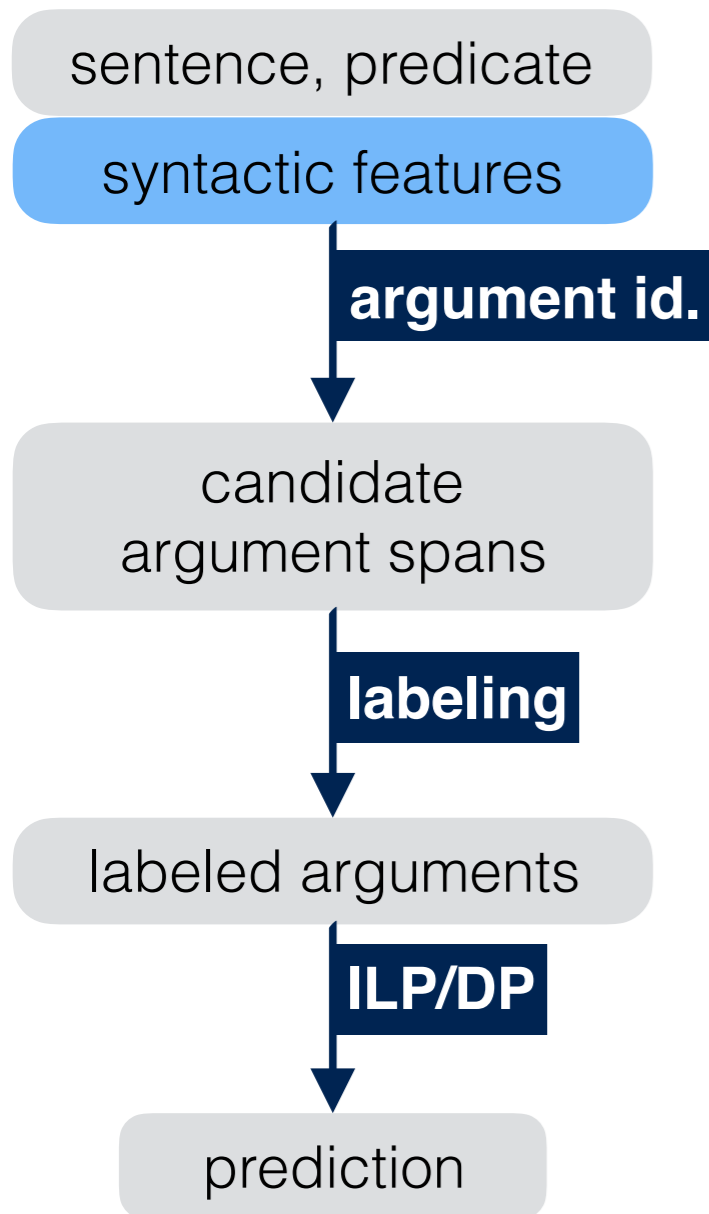pieces (final state)
ARG3

temporal
ARGM-TMP

Frame: *break.01*

| role | description |
|------|-------------|
| ARG0 | breaker |
| ARG1 | thing broken |
| ARG2 | instrument |
| ARG3 | pieces |
| ARG4 | broken away from what? |

# SRL is a hard problem …

- Over 10 years, F1 on PropBank:
  **80.3** (Toutanova et al, 2005) — **80.3** (FitzGerald et al, 2015)

- Many interesting challenges:
  Syntactic alternation
  Prepositional phrase attachment
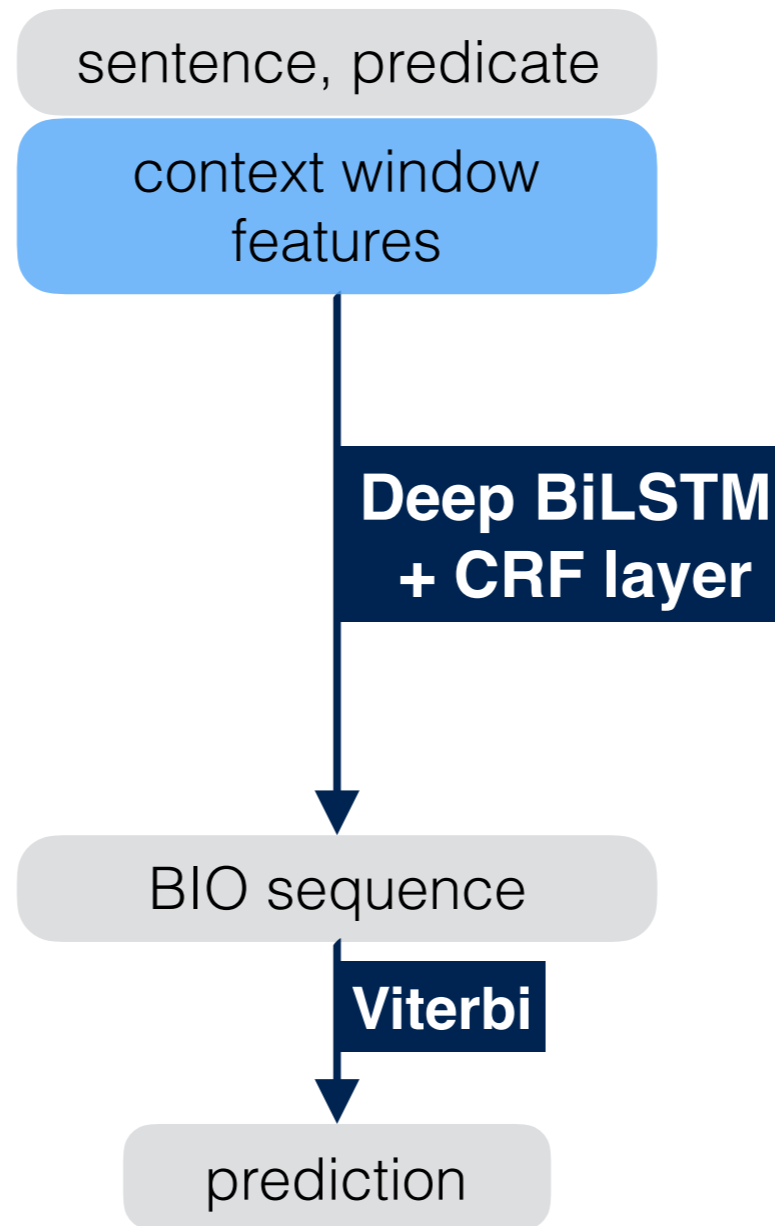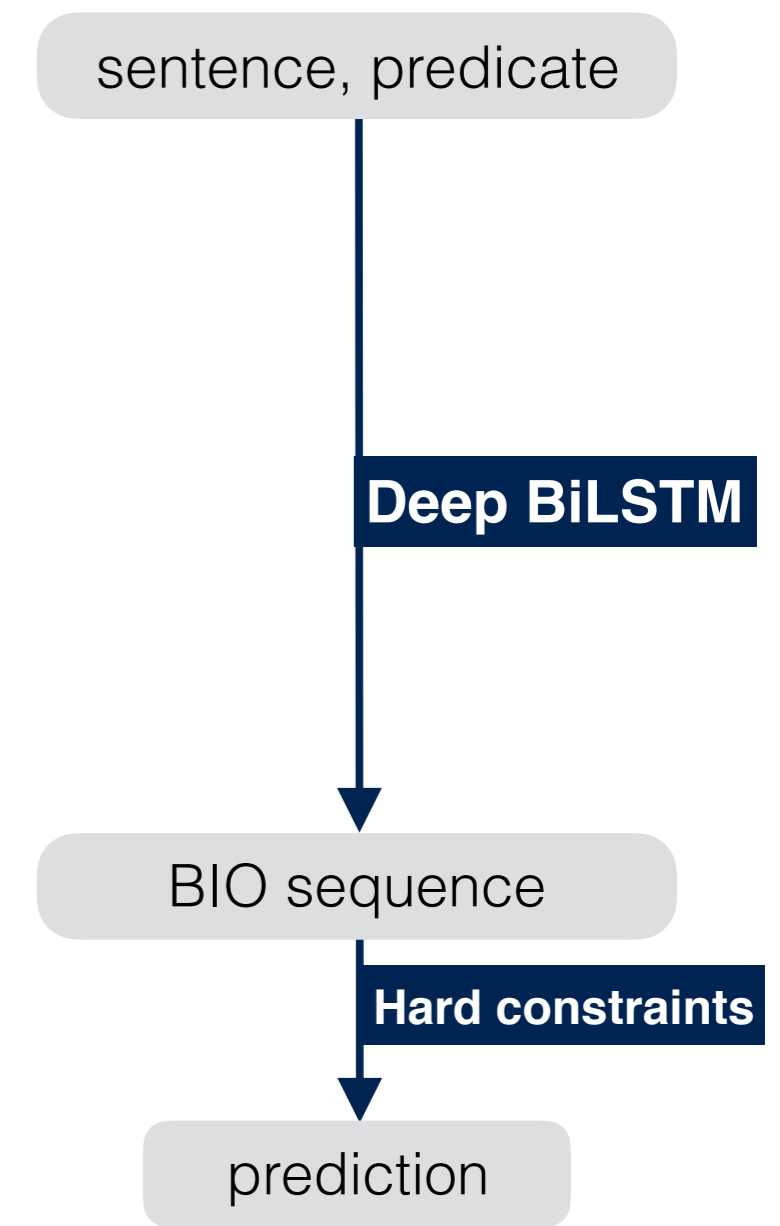  Long-range dependencies and common sense

# SRL Systems

| Pipeline Systems | End-to-end Systems | **\*This work** |
|---|---|---|

**Pipeline Systems**

- sentence, predicate
- syntactic features
  - → **argument id.**
- candidate argument spans
  - → **labeling**
- labeled arguments
  - → **ILP/DP**
- prediction

Punyakanok et al., 2008
Täckström et al., 2015
FitzGerald et al., 2015

**End-to-end Systems**

- sentence, predicate
- context window features
  - → **Deep BiLSTM + CRF layer**
- BIO sequence
  - → **Viterbi**
- prediction

Collobert et al., 2011
Zhou and Xu, 2015
Wang et. al, 2015

**\*This work**

- sentence, predicate
  - → **Deep BiLSTM**
- BIO sequence
  - → **Hard constraints**
- prediction

He et al., 2017

# SRL as BIO Tagging Problem

Input (sentence and predicate):

The | cats | *love* | hats | .

BIO output:

(**B**egin, **I**nside, **O**utside)

B-ARG0 | I-ARG0 | B-V | I-ARG1 | O

Final SRL output:

ARG0 | V | ARG1

(4) Viterbi decoding with hard constraints

| B-ARG0 | 0.4 |
|--------|-----|
| I-ARG0 | 0.05 |
| B-ARG1 | 0.5 |
| I-ARG1 | 0.03 |

| B-ARG0 | 0.1 |
|--------|-----|
| I-ARG0 | 0.5 |
| B-ARG1 | 0.1 |
| I-ARG1 | 0.2 |

| B-ARG0 | 0.001 |
|--------|-------|
| I-ARG0 | 0.001 |
| B-ARG1 | 0.001 |
| ... | ... |
| B-V | 0.95 |

| B-ARG0 | 0.1 |
|--------|-----|
| I-ARG0 | 0.1 |
| B-ARG1 | 0.7 |
| I-ARG1 | 0.2 |

(3) Variational dropout

(2) Highway connections

(1) Deep BiLSTM tagger

the [ ]     cats [ ]     love [V]     hats [ ]

[He et al, 2017]

# Other Implementation Details …



- 8 layer BiLSTMs with 300D hidden layers.

- 100D GloVe embeddings, updated during training.

- **Orthonormal initialization** for LSTM weight matrices (Saxe et al., 2013)

- 5 model ensemble with **product-of-experts** (Hinton 2002)

- Trained for 500 epochs.

# CoNLL 2005 Results

■ WSJ Test ■ Brown (out-domain) Test *:Ensemble models

| | Ours* 2017 | Ours 2017 | Zhou 2015 | FitzGerald* 2015 | Täckström 2015 | Toutanova* 2008 | Punyakanok* 2008 |
|---|---|---|---|---|---|---|---|
| WSJ Test | 84.6 | 83.1 | 82.8 | 80.3 | 79.9 | 80.3 | 79.4 |
| Brown (out-domain) Test | 73.6 | 72.1 | 69.4 | 72.2 | 71.3 | 68.8 | 67.8 |

F1

BiLSTM models    Pipeline models

# Ablations

## (single model, on CoNLL05 Dev)

Full model    No highway    No orthonormal init.    No dropout

Without dropout, model overfits at ~300 epochs.

Without orthonormal initialization, the deep model learns very slowly

F1 on Dev. Set

Num. Epochs

# Error Breakdown

## Oracle Transformations

Fix Label:

**[We]** _fly_ to NYC tomorrow.

~~ARG0~~

ARG1

Labeling error 29%

Split/Merge span:

ARG1

I _eat_ **[pasta with delight]**.

ARG1    ARGM-MNR

**[pasta] [with delight]**

ARG1    ARGM-MNR

I _eat_ **[pasta] [with broccoli]**.

ARG1

**[pasta with broccoli]**

Attachment error 25%

# Labeling Errors

Confusion matrix for labeling errors (column normalized)

| pred. \ gold | A0 | A1 | A2 | A3 | ADV | DIR | LOC | MNR | PNC | TMP |
|---|---|---|---|---|---|---|---|---|---|---|
| A0 | - | 55 | 11 | 13 | 4 | 0 | 0 | 0 | 0 | 0 |
| A1 | 78 | - | 46 | 0 | 0 | 22 | 11 | 10 | 25 | 14 |
| A2 | 11 | 23 | - | 48 | 15 | 56 | 33 | 41 | 25 | 0 |
| A3 | 3 | 2 | 2 | - | 4 | 0 | 0 | 0 | 25 | 14 |
| ADV | 0 | 0 | 0 | 4 | - | 0 | 15 | 29 | 25 | 36 |
| DIR | 0 | 0 | 5 | 4 | 0 | - | 11 | 2 | 0 | 0 |
| LOC | 5 | 9 | 12 | 0 | 4 | 0 | - | 10 | 0 | 14 |
| MNR | 3 | 0 | 12 | 26 | 33 | 0 | 0 | - | 0 | 21 |
| PNC | 0 | 3 | 5 | 4 | 0 | 11 | 4 | 2 | - | 0 |
| TMP | 0 | 8 | 5 | 0 | 41 | 11 | 26 | 6 | 0 | - |

- ARG2 is often confused with certain adjuncts (DIR, LOC, MNR), why?

---

**Predicate: *move***

**Arg0-PAG**: *mover*
**Arg1-PPT**: *moved*
**Arg2-GOL**: *destination*
**Arg3-VSP**: *aspect, domain in which arg1 moving*

---

**Predicate: *cut***

**Arg0-PAG**: *intentional cutter*
**Arg1-PPT**: *thing cut*
**Arg2-DIR**: *medium, source*
**Arg3-MNR**: *instrument, unintentional cutter*
**Arg4-GOL**: *beneficiary*

---

**Predicate: *strike***

**Arg0-PAG**: *Agent*
**Arg1-PPT**: *Theme(-Creation)*
**Arg2-MNR**: *Instrument*

---

- **Argument-adjunct distinctions** are difficult even for expert annotators!

# PP Attachment

Wrong PP attachment
(attach high)

Arg1 (NP)     Arg2 (PP)

Wrong SRL spans

Sumimoto ***financed*** the acquisition from Sears

merge

Arg1 (NP)

Correct PP attachment
(attach low)

Correct SRL spans

## Takeaway

— Traditionally hard tasks, such as **argument-adjunct** distinction and **PP attachment decisions** are still challenging!

# New Learning Approaches

*New state-of-the-art results for two tasks:*

## Coreference:

A fire in a Bangladeshi garment factory has left at least 37 people dead and 100 hospitalized. Most of the deceased were killed in the crush as workers tried to flee the blaze in the four-story building.

## Semantic Role Labeling:

| | |
|---|---|
| ARG0 | NASA |
| PRED | *observe* |
| ARG1 | an X-ray flare 400 times brighter than usual |
| TMP | On January 5, 2015 |

*Common themes:*
- End-to-end training of deep neural networks
- No preprocessing (e.g., no POS, no parser, etc.)
- Large gains in accuracy with simpler models and no extra training data

# Coreference Resolution

| Input document |
|---|
| A fire in a Bangladeshi garment factory has left at least 37 people dead and 100 hospitalized. Most of the deceased were killed in the crush as workers tried to flee the blaze in the four-story building. |

# Coreference Resolution

| Input document |
|---|
| A fire in a Bangladeshi garment factory has left at least 37 people dead and 100 hospitalized. Most of the deceased were killed in the crush as workers tried to flee the blaze in the four-story building. |

| Cluster #1 | A fire in a Bangladeshi garment factory | the blaze in the four-story building |
|---|---|---|

# Coreference Resolution

| Input document |
|---|
| A fire in a Bangladeshi garment factory has left at least 37 people dead and 100 hospitalized. Most of the deceased were killed in the crush as workers tried to flee the blaze in the four-story building. |

| Cluster #1 | A fire in a Bangladeshi garment factory | the blaze in the four-story building |
|---|---|---|
| Cluster #2 | a Bangladeshi garment factory | the four-story building |

# Coreference Resolution

| | Input document |
|---|---|
| | A fire in a Bangladeshi garment factory has left at least 37 people dead and 100 hospitalized. Most of the deceased were killed in the crush as workers tried to flee the blaze in the four-story building. |

| **Cluster #1** | A fire in a Bangladeshi garment factory | the blaze in the four-story building |
|---|---|---|
| **Cluster #2** | a Bangladeshi garment factory | the four-story building |
| **Cluster #3** | at least 37 people | the deceased |

# Two Subproblems

| **Input document** |
|---|
| A fire in a Bangladeshi garment factory has left at least 37 people dead and 100 hospitalized. Most of the deceased were killed in the crush as workers tried to flee the blaze in the four-story building. |

**Mention detection** →

| |
|---|
| A fire in a Bangladeshi garment factory |
| at least 37 people |
| … |
| the four-story building |

**Mention clustering**

| **Cluster #1** | A fire in a Bangladeshi garment factory | the blaze in the four-story building |
|---|---|---|
| **Cluster #2** | a Bangladeshi garment factory | the four-story building |
| **Cluster #3** | at least 37 people | the deceased |

# Previous Approach: Rule-based pipeline

**Input document**

A fire in a Bangladeshi garment factory has left at least 37 people dead and 100 hospitalized.

**Syntactic parser**

**Hand-engineered rules**

| Candidate mentions |
| --- |
| A fire in a Bangladeshi garment factory |
| garment |
| factory |
| at least 37 people dead and 100 hospitalized |
| … |

| Mention #1 | Mention #2 | Coreferent? |
| --- | --- | --- |
| A fire in a Bangladeshi garment factory | garment | ✓/✗ |
| garment | factory | ✓/✗ |
| factory | at least 37 people dead and 100 hospitalized | ✓/✗ |
| … | … | ✓/✗ |

# Previous Approach: Rule-based pipeline



**Input document**

A fire in a Bangladeshi garment factory has left at least 37 people dead and 100 hospitalized.

**Syntactic parser**

**Hand-engineered rules**

Mention clustering: main source of improvement for many years!

- Haghighi and Klein (2010)
- Raghunathan et al. (2010)
- …
- Clark & Manning (2016)

| | Mention #2 | Coreferent? |
|---|---|---|
| | garment | ✓/✗ |
| | factory | ✓/✗ |
| | least 37 people dead and 100 hospitalized | ✓/✗ |
| at | … | ✓/✗ |

# Previous Approach: Rule-based pipeline

**Input document**

A fire in a Bangladeshi garment factory has left at least 37 people dead and 100 hospitalized.

**Syntactic parser**

**Hand-engineered rules**

Relies on parser for:
- mention detection
- syntactic features for clustering (e.g. head words)

| | Coreferent? |
|---|---|
| A fire in a Bangladeshi garment factory | |
| garment | |
| factory | |
| at least 37 people dead and 100 hospitalized | |
| … | |

| | | Coreferent? |
|---|---|---|
| A fire in a Bangladeshi garment factory | garment | ✓/✗ |
| garment | factory | ✓/✗ |
| factory | at least 37 people dead and 100 hospitalized | ✓/✗ |
| … | … | ✓/✗ |

# End-to-end Approach

- Consider all possible spans

- Learn to rank antecedent spans

- Factored model to prune search space

# Key Idea: Span Representations

# Key Idea: Span Representations

Span representation

the Postal Service

Bidirectional LSTM

Word & character embeddings

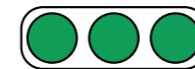General   Electric   said   the   Postal   Service   contacted   the   company

# Key Idea: Span Representations



Boundary representations

Span representation

the Postal Service

Bidirectional LSTM

Word & character embeddings

General    Electric    said    the    Postal    Service    contacted    the    company

# Key Idea: Span Representations

Attention mechanism
to learn headedness

the Postal Service
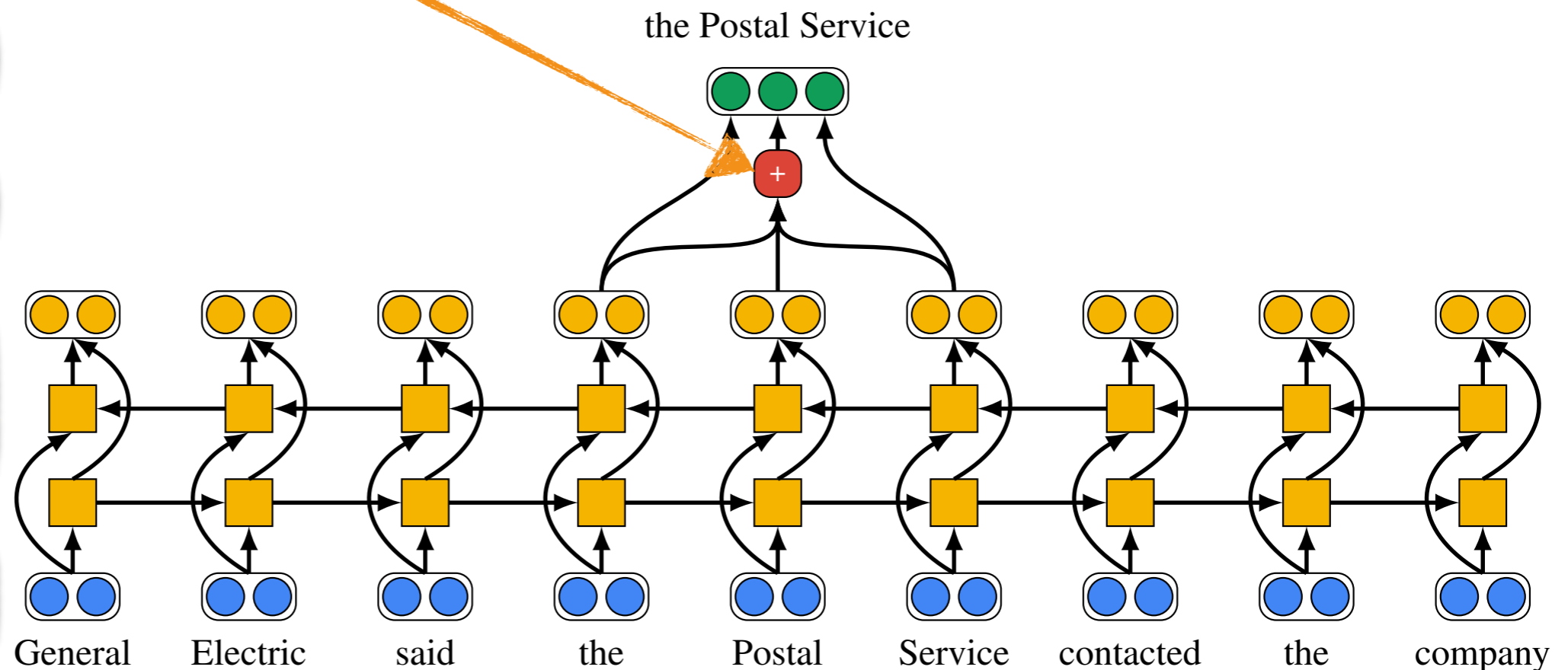
Span representation

Head-finding attention

Bidirectional LSTM

Word & character
embeddings

General   Electric   said   the   Postal   Service   contacted   the   company

# Key Idea: Span Representations

# Mention Ranking

Every span independently chooses an antecedent

| Input document |
| --- |
| A fire in a Bangladeshi garment factory has left at least 37 people dead and 100 hospitalized. Most of the deceased were killed in the crush as workers tried to flee the blaze in the four-story building. Witnesses say the only exit door was on the ground floor, and that it was locked when the fire broke out. |

# Mention Ranking

- Reason over all possible spans

- Assign an antecedent to every span

$$y_3 \in \{\epsilon, 1, 2\}$$

| | Span | Antecedent |
|---|---|---|
| 1 | A | $y_1$ |
| 2 | A fire | $y_2$ |
| 3 | A fire in | $y_3$ |
| … | … | … |
| M | out | $y_M$ |

# Mention Ranking

- Reason over all possible spans

- Assign an antecedent to every span

$$y_3 \in \{\epsilon, 1, 2\}$$

| | Span | Antecedent |
|---|---|---|
| 1 | A | $y_1$ |
| 2 | A fire | $y_2$ |
| 3 | A fire in | $y_3$ |
| … | … | … |
| M | out | $y_M$ |

$\epsilon$ : no coreference link

# Mention Ranking

- Reason over all possible spans

- Assign an antecedent to every span

| | Span | Antecedent |
|---|---|---|
| 1 | A | $y_1$ |
| 2 | A fire | $y_2$ |
| 3 | A fire in | $y_3$ |
| ... | ... | ... |
| M | out | $y_M$ |

$$y_3 \in \{\epsilon, 1, 2\}$$

Coreference link from span 1 to span 3

# Mention Ranking

- Reason over all possible spans

- Assign an antecedent to every span

| | Span | Antecedent |
|---|---|---|
| 1 | A | $y_1$ |
| 2 | A fire | $y_2$ |
| 3 | A fire in | $y_3$ |
| … | … | … |
| M | out | $y_M$ |

$$y_3 \in \{\epsilon, 1, 2\}$$

Coreference link from span 2 to span 3

# Example Clustering

| Input document |
|---|
| A fire in a Bangladeshi garment factory has left at least 37 people dead and 100 hospitalized. Most of the deceased were killed in the crush as workers tried to flee the blaze in the four-story building. Witnesses say the only exit door was on the ground floor, and that it was locked when the fire broke out. |

| Span | Antecedent ($y_i$) |
|---|---|
| A | $\epsilon$ |
| A fire | $\epsilon$ |
| ... | ... |
| a Bangladeshi garment factory | $\epsilon$ |
| ... | ... |
| the four-story building | a Bangladeshi garment factory |
| ... | ... |
| out | $\epsilon$ |

# Example Clustering

| Input document |
|---|
| A fire in a Bangladeshi garment factory has left at least 37 people dead and 100 hospitalized. Most of the deceased were killed in the crush as workers tried to flee the blaze in the four-story building. Witnesses ...d floor, and that it was locked when the fire broke out. |

**Not a mention**

| Span | Antecedent ( $y_i$ ) |
|---|---|
| A | $\epsilon$ |
| A fire | $\epsilon$ |
| ... | ... |
| a Bangladeshi garment factory | $\epsilon$ |
| ... | ... |
| the four-story building | a Bangladeshi garment factory |
| ... | ... |
| out | $\epsilon$ |

# Example Clustering

| Input document |
|---|
| A fire in a Bangladeshi garment factory has left at least 37 people dead and 100 hospitalized. Most of the deceased were killed in the crush as workers tried to flee the blaze in the four-story building. Witnesses say the only exit door was on the ground floor, and that it was locked when the fire broke out. |

| Span | Antecedent ( $y_i$ ) |
|---|---|
| ... | ... |
| a Bangladeshi garment factory | $\epsilon$ |
| ... | ... |
| the four-story building | a Bangladeshi garment factory |
| ... | ... |
| out | $\epsilon$ |

No link with previously occurring span

# Example Clustering

| Input document |
|---|
| A fire in a Bangladeshi garment factory has left at least 37 people dead and 100 hospitalized. Most of the deceased were killed in the crush as workers tried to flee the blaze in the four-story building. Witnesses say the only exit door was on the ground floor, and that it was locked when the fire broke out. |

| Span | Antecedent ($y_i$) |
|---|---|
| A | $\epsilon$ |
| A fire | $\epsilon$ |
| ... | ... |
| | $\epsilon$ |
| ... | ... |
| the four-story building | a Bangladeshi garment factory |
| ... | ... |
| out | $\epsilon$ |

Predicted coreference link

# Span Ranking Model

$$P(y_1, \ldots, y_M \mid D) = \prod_{i=1}^{M} P(y_i \mid D)$$

$$= \prod_{i=1}^{M} \frac{e^{s(i,y_i)}}{\sum_{y' \in \mathcal{Y}(i)} e^{s(i,y')}}$$

Factor coreference score $s(i,j)$ to enable span pruning:

$$s(i,j) = \begin{cases} s_\mathrm{m}(i) + s_\mathrm{m}(j) + s_\mathrm{a}(i,j) & j \neq \epsilon \\ 0 & j = \epsilon \end{cases}$$

# Span Ranking Model

$$P(y_1, \ldots, y_M \mid D) = \prod_{i=1}^{M} P(y_i \mid D)$$

$$\frac{e^{s(i,y_i)}}{\sum_{y' \in \mathcal{Y}(i)} e^{s(i,y')}}$$

Is this span a mention?

Factor coreference score $s(i,j)$ to enable span pruning:

$$s(i,j) = \begin{cases} s_{\mathrm{m}}(i) + s_{\mathrm{m}}(j) + s_{\mathrm{a}}(i,j) & j \neq \epsilon \\ 0 & j = \epsilon \end{cases}$$

# Span Ranking Model

$$P(y_1, \ldots, y_M \mid D) = \prod_{i=1}^{M} P(y_i \mid D)$$

$$= \prod_{i=1}^{M} \frac{e^{s(i, y_i)}}{\ldots}$$

Is span j an antecedent of span i?

Factor coreference score $s(i, j)$ to enable span pruning:

$$s(i, j) = \begin{cases} s_{\mathrm{m}}(i) + s_{\mathrm{m}}(j) + s_{\mathrm{a}}(i, j) & j \neq \epsilon \\ 0 & j = \epsilon \end{cases}$$

# Span Ranking Model

$$P(y_1, \ldots, y_M \mid D) = \prod_{i=1}^{M} P(y_i \mid D)$$

$$= \prod_{i=1}^{M} \frac{e^{s(i,y_i)}}{\sum_{y' \in \mathcal{Y}(i)} e^{s(i,y')}}$$

Factor coreference score $s(i,j)$ to enable span pruning:

$$s(i,j) = \begin{cases} s_{\mathrm{m}}(i) + s_{\mathrm{m}}(j) + s_{\mathrm{a}}(i,j) & j \neq \epsilon \\ 0 & j = \epsilon \end{cases}$$

Dummy antecedent
has a fixed zero score

# Experimental Setup

**Dataset**: English OntoNotes (CoNLL-2012)

**Genres**: Telephone conversations, newswire, newsgroups, broadcast conversation, broadcast news, weblogs

**Documents**: 2802 training, 343 development, 348 test

Longest document has 4009 words!

**Aggressive pruning**: Maximum span width, maximum sentence training, suppress spans with inconsistent bracketing, maximum number of antecedents

**Features**: distance between spans, span width

**Metadata**: speaker information, genre

# Coreference Results

# Coreference Results

# Coreference Results

Coreference Results

# Mention Recall



○ Raghunathan et al. (2010)   ○ Our model (actual threshold)   — Our model (various thresholds)
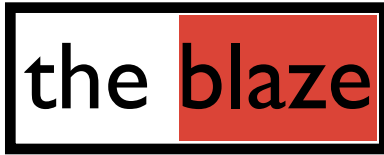
# Mention Recall

# Head-finding Agreement



% of constituent spans with predicted
heads that agree with syntactic heads

# Qualitative Analysis

☐ : Mention in a predicted cluster

🟥 : Head-finding attention weight

A fire in a Bangladeshi garment factory has left at least 37 people dead and 100 hospitalized. Most of the deceased were killed in the crush as workers tried to flee the blaze in the four-story building.

# Qualitative Analysis

□ : Mention in

■ : Head-findir

**Attention-based head finder facilitates soft similarity cues**

A fire in a Bangladeshi garment factory has left at least 37 people dead and 100 hospitalized. Most of the deceased were killed in the crush as workers tried to flee the blaze in the four-story building.

# Qualitative Analysis

☐ : Mention in a predicted cluster

🟥 : Hea...

**Good head-finding requires word-order information!**

A **fire** in a Bangladeshi garment factory has left at least 37 people dead and 100 hospitalized. Most of the deceased were killed in the crush as workers tried to flee the **blaze** in the four-story building.

# Common Error Case

The flight attendants have until 6:00 today to ratify labor concessions. The pilots' union and ground crew did so yesterday.

# Common Error Case



☐ : Mention in a predicted cluster

🟥 : Head-finding attention weight

The flight attendants have until 6:00 today to ratify labor concessions. The pilots' union and ground crew did so yesterday.

Conflating **relatedness** with **paraphrasing**

# Does the Recipe Work for Broad Coverage Semantics?

*Step 1: Gather lots of training data!*

**Challenge 1: Data is costly and limited
(e.g. linguists required to label
PennTreebank / OntoNotes)**

*Step 2: Apply Deep Learning!!*

**Challenge 2: Pipeline of structured
prediction problems with cascading errors
(e.g. POS->Parsing->SRL->Coref)**

*Step 3: Observe Impressive Gains!!!*

# Where Will the Data Come From???

**Option 1:** Semi-supervised learning

- E.g. word2vec and GloVe are in wide use
    [Mikolov et al., 2013; Pennington et al., 2014]

- Can we learn better word representations?

**Option 2:** Supervised learning

- Can we gather more direct forms of supervision?

# Learning Better Word Representations

**Goal:** Model contextualized syntax and semantics

$$R(w_i, w_1 \ldots w_n) \in \mathbb{R}^n$$

$R$(plays, "The robot plays piano.")

$$\neq$$

$R$(plays, "The robot starred in many plays.")

# Word Embeddings from a Language Model

**Step 1:** Train a large BiLM on unlabeled data



2 Layer Bidirectional LSTM

Character convolutions

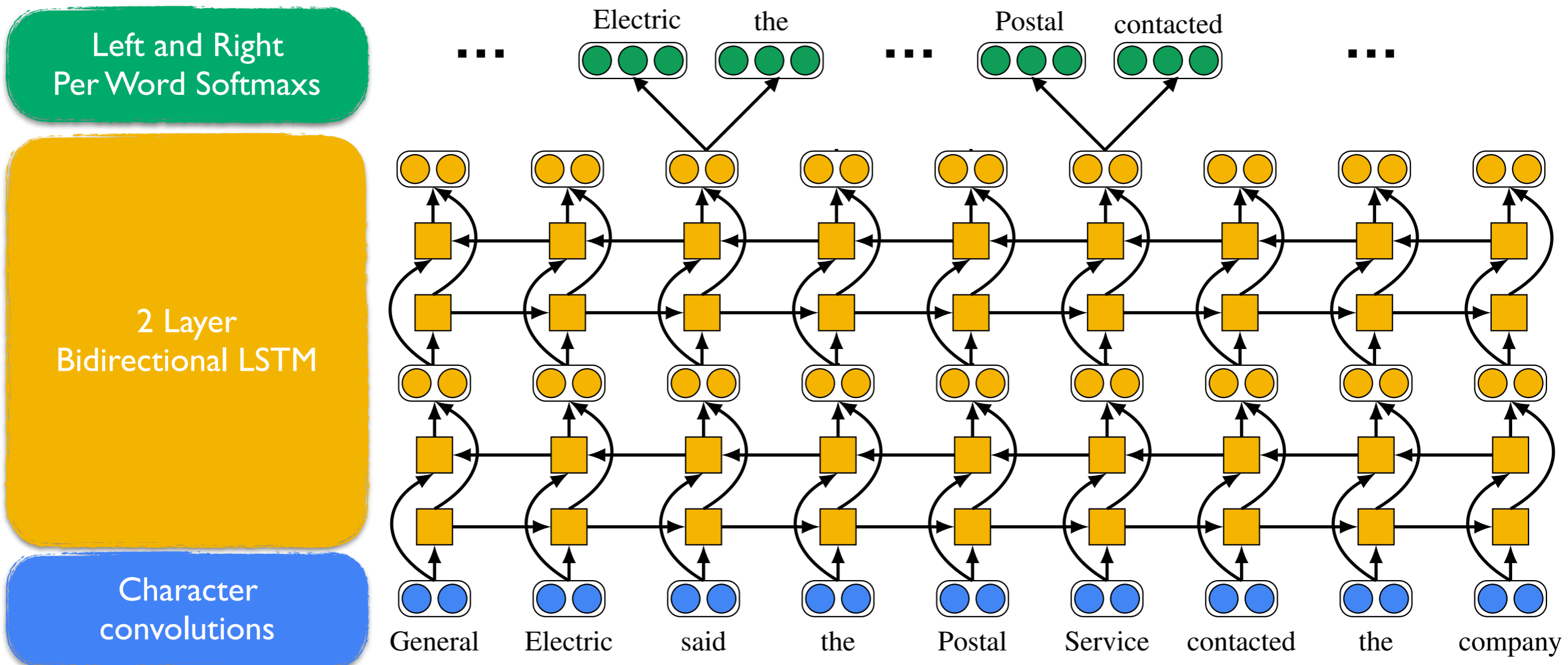General  Electric  said  the  Postal  Service  contacted  the  company
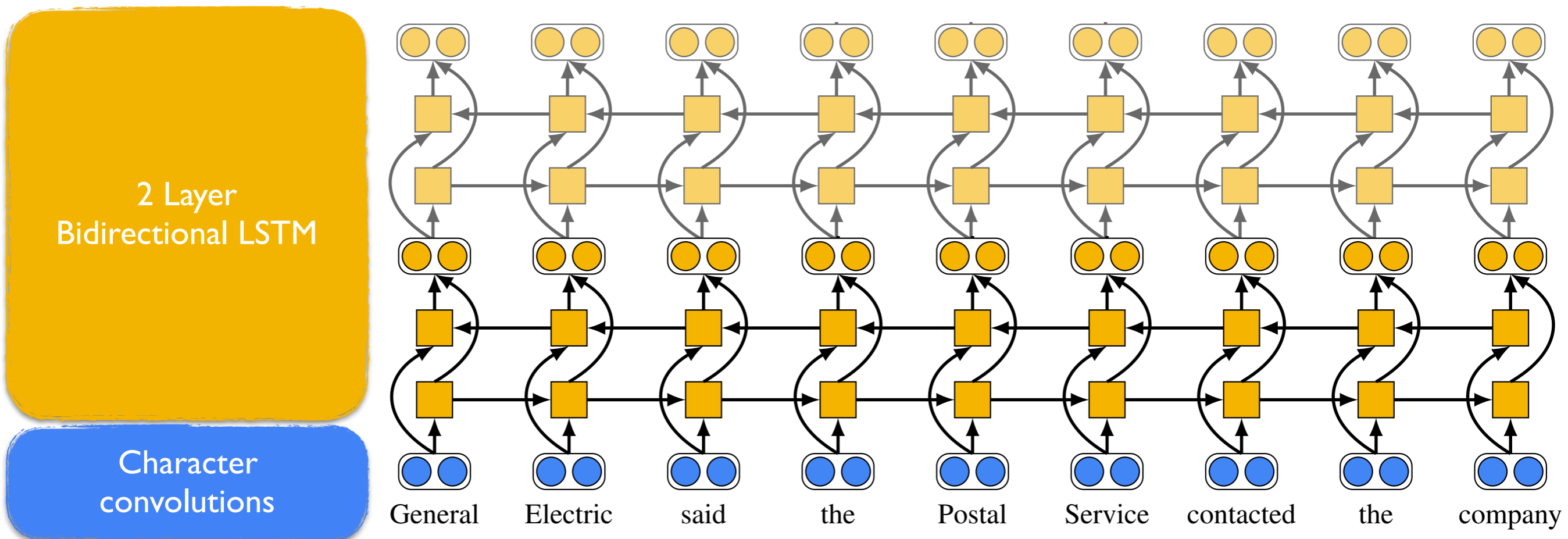
# Word Embeddings from a Language Model

**Step 1:** Train a large BiLM on unlabeled data

# Word Embeddings from a Language Model

**Step 1:** Train a large BiLM on unlabeled data
**Step 2:** Compute linear function of pre-trained model

# Word Embeddings from a Language Model

**Step 1:** Train a large BiLM on unlabeled data
**Step 2:** Compute linear function of pre-trained model

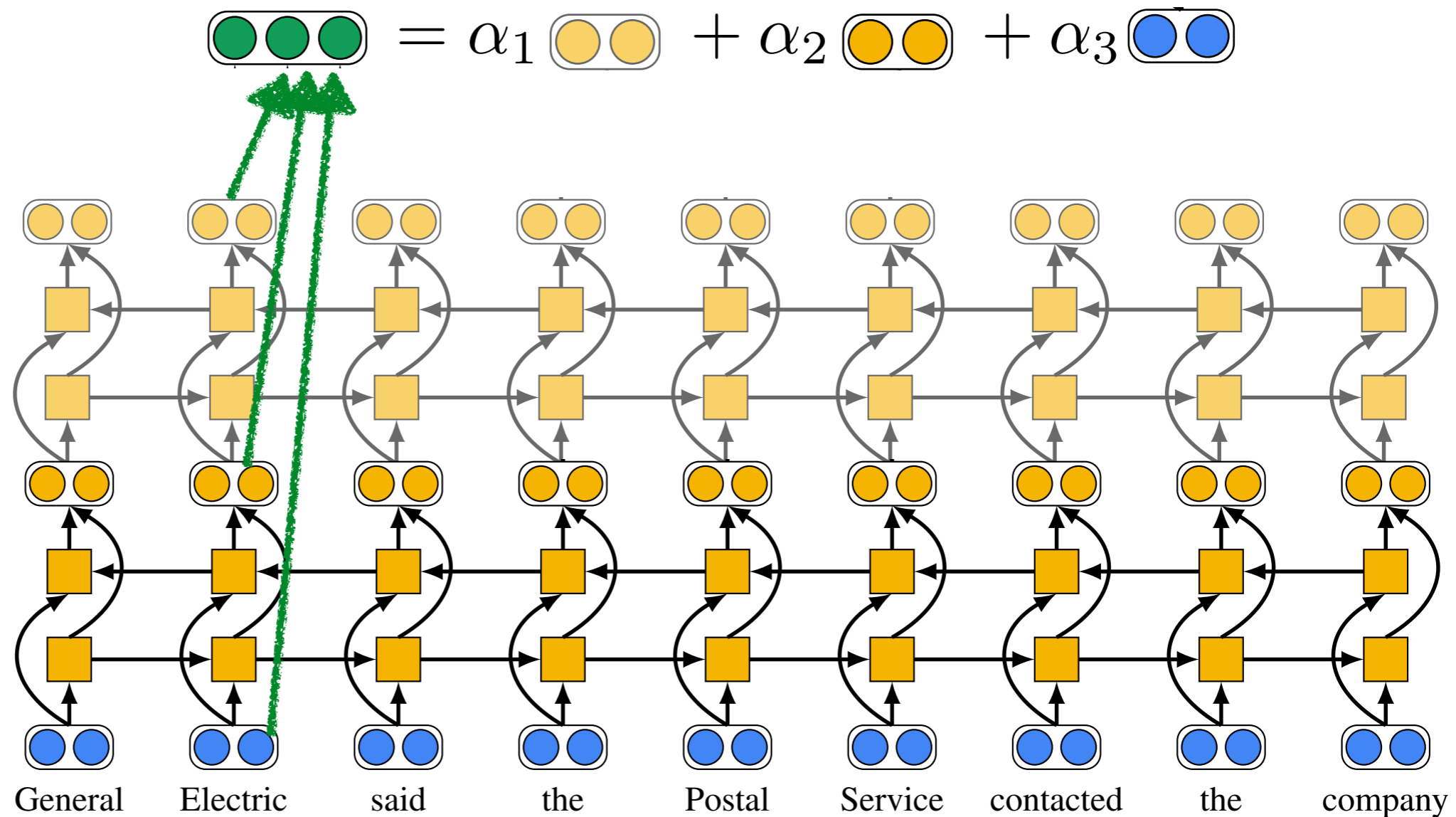# Word Embeddings from a Language Model

**Step 1:** Train a large BiLM on unlabeled data
**Step 2:** Compute linear function of pre-trained model
**Step 3:** Learn weights for each end task

Best Single System Results

# SOTA For Many Others Tasks

■ Previous SOTA    ■ Baseline    ■ Baseline+LM

| | SNLI | SQuAD | Coref | SRL | NER | Sentiment (SST) |
|---|---|---|---|---|---|---|
| Previous SOTA | 88.1 | 84.3 | 67.2 | 81.7 | 91.9 | 53.7 |
| Baseline | 88 | 81.1 | 67.2 | 81.4 | 90.2 | 51.4 |
| Baseline+LM | 88.7 | 85.3 | 70.4 | 84.6 | 92.2 | 54.7 |

# What Does it Learn?

**Semantics:**

- Supervised WSD task [Miller et al.,1994]

- Use N-th layer in NN classifier

**Syntax:**

- Label POS corpus [Marcus et al., 1993]

- Learn classifier on N-th layer

# Where Will the Data Come From???

**Option 1:** Semi-supervised learning

- E.g. word2vec and GloVe are in wide use
  [Mikolov et al., 2013; Pennington et al., 2014]

- Can we learn better word representations?

**Option 2:** Supervised learning

- Can we gather more direct forms of supervision?

# A First Data Step: QA-SRL

- Introduce a **new SRL** formulation with **no frame or role inventory**

- Use **question-answer pairs** to model verbal predicate-argument relations

- Annotated **over 3,000 sentences in weeks** with **non-expert**, part-time annotators

- Showed that this data is **high-quality** and **learnable**

[He et al, 2015]

# Previous Method: Annotation with Frames

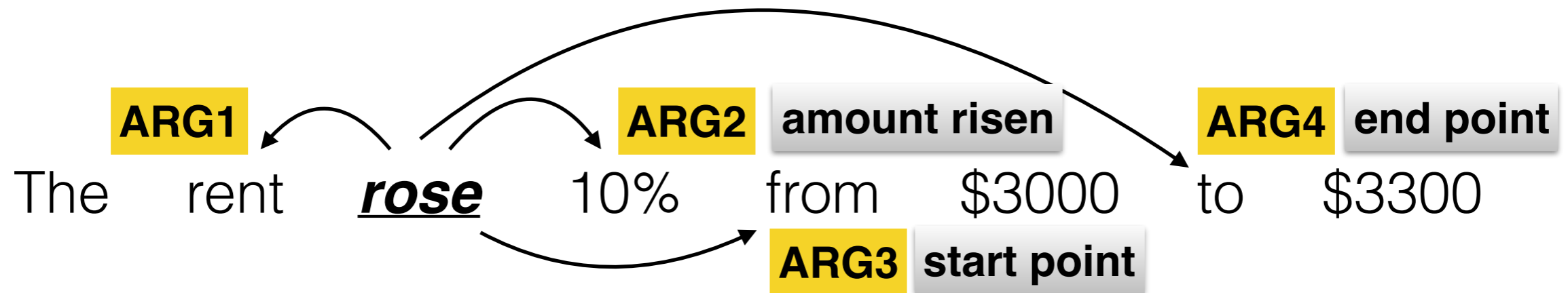ARG1 ← The rent **_rose_** 10% → ARG2 **amount risen** from $3000 to $3300 → ARG4 **end point**

ARG3 **start point**

Frameset: rise.01 , *go up*

   **Arg1-**: *Logical subject, patient, thing rising*
   **Arg2-EXT**: *EXT, amount risen*
   **Arg3-DIR**: *start point*
   **Arg4-LOC**: *end point*
   **Argm-LOC**: *medium*

- Depends on pre-defined frame inventory, requires syntactic parses
- Annotators need to:
   1) Identify the Frameset
   2) Find arguments in the parse
   3) Assign labels accordingly
- If frame doesn't exist, create new

The Proposition Bank: An Annotated Corpus of Semantic Roles, Palmer et al., 2005
http://verbs.colorado.edu/propbank/framesets-english/rise-v.html

# Our Annotation Scheme

They ***increased*** the rent this year .

Who increased something ?

They

What is increased ?

the rent

When is something increased ?

this year

# Cost and Speed



- Part-time freelancers from upwork.com (hourly rate: $10)
- ~2h screening process for native English proficiency

# Wh-words vs. PropBank Roles

| | Who | What | When | Where | Why | How | HowMuch |
|---|---|---|---|---|---|---|---|
| **ARG0** | 1575 | 414 | 3 | 5 | 17 | 28 | 2 |
| **ARG1** | 285 | 2481 | 4 | 25 | 20 | 23 | 95 |
| **ARG2** | 85 | 364 | 2 | 49 | 17 | 51 | 74 |
| **ARG3** | 11 | 62 | 7 | 8 | 4 | 16 | 31 |
| **ARG4** | 2 | 30 | 5 | 11 | 2 | 4 | 30 |
| **ARG5** | 0 | 0 | 0 | 1 | 0 | 2 | 0 |
| **AM-ADV** | 5 | 44 | 9 | 2 | 25 | 27 | 6 |
| **AM-CAU** | 0 | 3 | 1 | 0 | 23 | 1 | 0 |
| **AM-DIR** | 0 | 6 | 1 | 13 | 0 | 4 | 0 |
| **AM-EXT** | 0 | 4 | 0 | 0 | 0 | 5 | 5 |
| **AM-LOC** | 1 | 35 | 10 | 89 | 0 | 13 | 11 |
| **AM-MNR** | 5 | 47 | 2 | 8 | 4 | 108 | 14 |
| **AM-PNC** | 2 | 21 | 0 | 1 | 39 | 7 | 2 |
| **AM-PRD** | 1 | 1 | 0 | 0 | 0 | 1 | 0 |
| **AM-TMP** | 2 | 51 | 341 | 2 | 11 | 20 | 10 |

**Advantages**
- Easily explained
- No pre-defined roles, few syntactic assumption
- Can capture implicit arguments
- Generalizable across domains

**Limitations**
- Only modeling verbs (for now)
- Not annotating verb senses directly
- Can have multiple equivalent questions

**Challenges**
- What questions to ask?
- How much data do we need?
- Can we generalize to other tasks, such as coref?

# Does the Recipe Work for Broad Coverage Semantics?

*Step 1: Gather lots of training data!*

✓ **Challenge 1: Data is costly and limited (e.g. linguists required to label PennTreebank / OntoNotes)**

*Step 2: Apply Deep Learning!!*

✓ **Challenge 2: Pipeline of structured prediction problems with cascading errors (e.g. POS->Parsing->SRL->Coref)**

*Step 3: Observe Impressive Gains!!!*

# Contributions

**Models**

- End-to-end deep learning for SRL and coreference

- No preprocessing (e.g. no parser or POS tagger)

**Data**

- Contextualized word embeddings from a language model

- First steps towards scalable data annotation

# The End: Questions?

**Future Directions**

- Multi-task learning, given architectural similarities

- Multi-lingual should work, in theory…

- Need to scale up data annotation efforts, and focus on out of domain performance

**Recent Release**

- AllenNLP: Deep Learning Semantic NLP toolkit

- See demos and code at AllenNLP.org