

# Harvesting More Answer Spans from Paragraph beyond Annotation

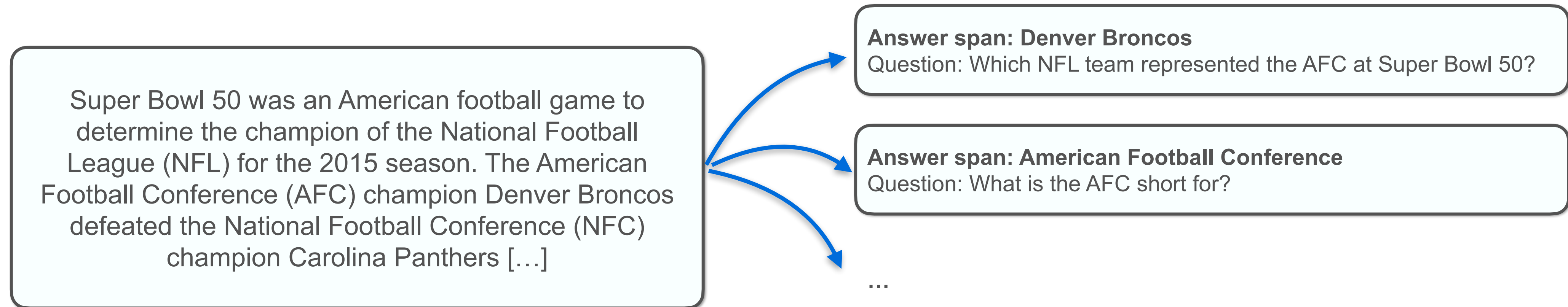
Qiaoben Bao, Jiangjie Chen, Linfang Liu, Jingping Liu, Jiaqing Liang, Yanghua Xiao

Fudan University



# Background

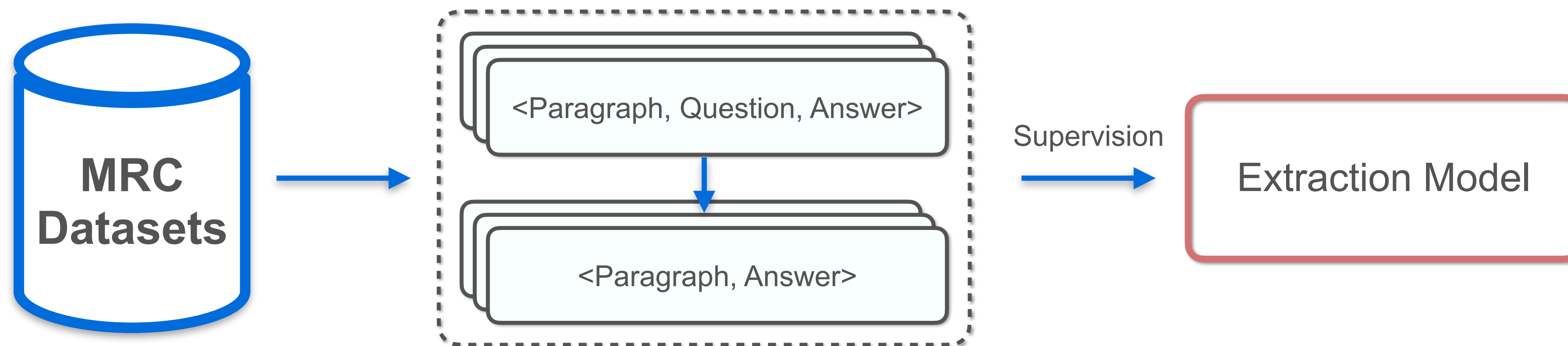
- Answer span extraction (AE) focuses on **identifying answer spans** from paragraphs



- AE has a wide range of both research and real-life applications
  - Facilitating information extraction
  - Data augmentation for MRC or QG
  - Building FAQs against documents
  - ...

# Background

- The current work of AE relies on the annotation from MRC datasets
  - There is currently no dataset dedicated to AE tasks
  - MRC dataset contains  $\langle \text{Paragraph}, \text{Question}, \text{Answer} \rangle$  triples and can be easily converted to  $\langle \text{Paragraph}, \text{Answer} \rangle$  pairs for AE



# Challenge

---

- Is the annotation from MRC sufficient for the AE task?
  - MRC datasets are only required to extract limited answer spans (usually 5) for each paragraph
  - The unannotated candidate spans may also be valid answer spans
- We analyze two well-known MRC datasets
  - SQuAD (Rajpurkar et al.) & DROP\* (Dua et al.)

\* DROP contains three types of answers and we only consider the extractive examples where the answer is a span from the original paragraph.

# Challenge

---

- Re-annotate 50 paragraphs from each dataset
  - We define the missing rate  $\gamma$  as

$$\gamma = \frac{|\mathcal{M}|}{|\mathcal{S}_p \cup \mathcal{M}|}$$

- $\mathcal{S}_p$  is the positive labeled samples in paragraph
- $\mathcal{M}$  is the unlabeled samples in paragraph

# Challenge

- Re-annotate 50 paragraphs from each dataset
  - We define the missing rate  $\gamma$  as

$$\gamma = \frac{|\mathcal{M}|}{|S_p \cup \mathcal{M}|}$$

- $S_p$  is the positive labeled samples in paragraph
  - $\mathcal{M}$  is the unlabeled samples in paragraph
- Both datasets contain a comparable number of positive answer spans not annotated among unlabeled candidate spans

Dataset	#Sentences	$ S_p $	$ \mathcal{M} $	$\gamma$
SQuAD	237	219	207	48.59%
DROP	445	296	492	62.44%



# Challenge

- This MRC annotation procedure ignores other detailed key information that would also be helpful for readers to understand the context.

## Paragraph #1:

[...] In 1815, the British government selected Saint Helena as the place of detention of Napoleon Bonaparte. He was taken to the island [...]

### Golden answer spans:

Answer span: 1815

Corresponding question: What year was Napoleon Bonaparte taken to the island?

Answer span: Napoleon Bonaparte

Corresponding question: The British government detained who in Saint Helena?

### Unlabeled answer spans:

Answer span: British

Corresponding question: Which government sent Napoleon Bonaparte to Saint Helena?

Answer span: Saint Helena

Corresponding question: Where was Napoleon Bonaparte imprisoned?

## Paragraph #2:

[...] In 1882, Albert Zahm built an [...]. Around 1899, Professor Jerome Green became the first American to [...]. In 1931, Father Julius Nieuwland performed early work [...]

### Golden answer spans:

Answer span: 1882

Corresponding question: In what year did Albert Zahm begin comparing aeronautical models at Notre Dame?

Answer span: Professor Jerome Green

Corresponding question: Which professor sent the first wireless message in the USA?

...

### Unlabeled answer spans:

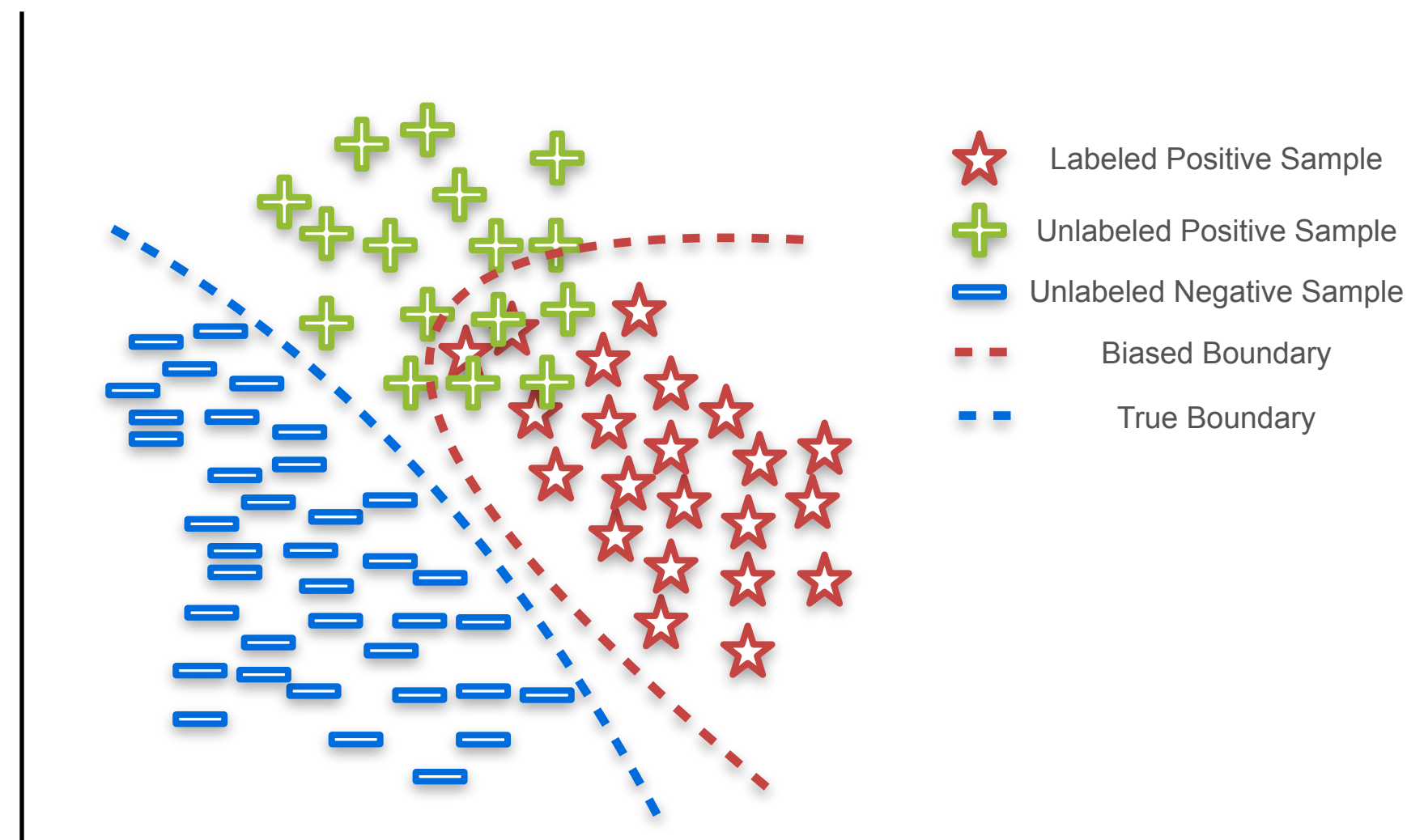
Answer span: 1931

Answer span: Albert Zahm

...

# Challenge

- This MRC annotation procedure ignores other detailed key information that would also be helpful for readers to understand the context.
  - Previous work
    - Directly treating unlabeled data as negative one may lead to the wrong decision boundary





# Challenge

---

- This MRC annotation procedure ignores other detailed key information that would also be helpful for readers to understand the context.
  - Previous work
    - Directly treating unlabeled data as negative one may lead to the wrong decision boundary
  - Our goal
    - Narrows the discrepancy between AE and MRC to extract more answer spans from paragraph
    - For answer span that is not labeled, we can automatically evaluate its quality

# Method

---

- Basic idea

- Formulate AE task as a positive-unlabeled learning problem
  - Split the risk estimator into the positive part and the negative part

$$R_\ell = \pi \mathbb{E}_{\mathbf{x}, y=1} \ell(f(\mathbf{x}), 1) + (1 - \pi) \mathbb{E}_{\mathbf{x}, y=0} \ell(f(\mathbf{x}), 0)$$

- $f$  is the classifier
- $\ell$  is loss function
- $\pi$  is the prior distribution of positive samples
- We do not have labeled negative samples for calculating  $(1 - \pi) \mathbb{E}_{\mathbf{x}, y=0} \ell(f(\mathbf{x}), 0)$

# Method

---

- Basic idea

- Formulate AE task as a positive-unlabeled learning problem

- Split the risk estimator into the positive part and the negative part

$$R_\ell = \pi \mathbb{E}_{\mathbf{x}, y=1} \ell(f(\mathbf{x}), 1) + (1 - \pi) \mathbb{E}_{\mathbf{x}, y=0} \ell(f(\mathbf{x}), 0)$$

- Re-estimate the negative part with positive samples and unlabeled samples

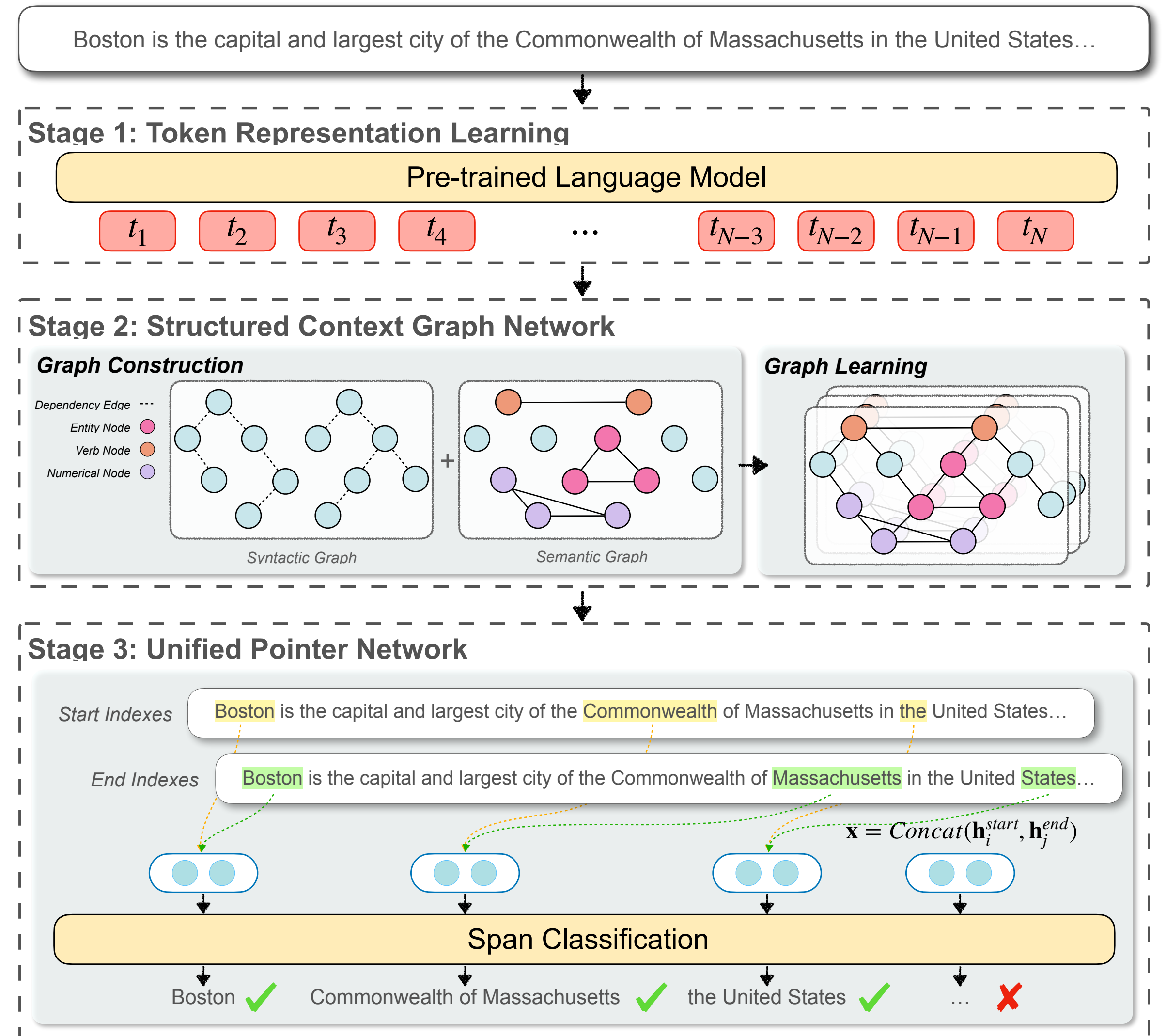
$$(1 - \pi) \mathbb{E}_{\mathbf{x}, y=0} \ell(f(\mathbf{x}), 0) = \mathbb{E}_{\mathbf{x}} \ell(f(\mathbf{x}), 0) - \pi \mathbb{E}_{\mathbf{x}, y=1} \ell(f(\mathbf{x}), 0)$$

- Finally, we can calculate  $R_\ell$  by estimating the prior distribution  $\pi$

# Method

- Framework: SCOPE

- PLM-based token representation
  - Encode sub-token with PLM
  - Token representation is averaged by sub-tokens
- GNN-based information prorogation
  - Syntactic edges
  - Semantic edges
- PU classifier
  - Unified pointer network
  - PU loss



# Evaluation

- Question-worthy score: Questionability + Worthiness
  - Questionability

## Questionability

If a span is askable, there exists at least one question that can be answered by this span with a high probability.

- Evaluated by a QG-QA model
  - The question generation model first generates questions based on the given paragraph and extracted answer span
  - The question answering model then scores the answer spans against the generated question

– Worthiness



# Evaluation

---

- Question-worthy score: Questionability + Worthiness
  - Questionability
  - Worthiness

## Worthiness

If a span is worthy of asking, it contains more information for people to ask a question.

- Evaluated by an extractive summarization model
  - Extractive summarization model score each sentence an informative score
  - For each candidate span, we define its worthiness as the informative score of the sentence it locates

# Experiments

- Conventional metrics
  - Our proposed framework **extracts more high-quality answer spans** on both datasets

Model	SQuAD				DROP			
	Precision	Recall	F1	Avg. spans	Precision	Recall	F1	Avg. spans
<i>ENT</i>	<i>13.63</i>	<i>40.41</i>	<i>20.39</i>	<i>12.82</i>	<i>6.31</i>	<i>52.58</i>	<i>11.27</i>	<i>43.90</i>
ENT Classifier (BERT <sub>base</sub> )	<b>48.55</b>	20.37	28.70	1.81	36.90	19.10	25.17	2.73
└ (SpanBERT <sub>base</sub> )	47.90	21.09	29.29	1.90	<b>38.62</b>	20.52	26.80	2.80
└ (RoBERTa <sub>base</sub> )	47.57	20.61	28.76	1.87	35.54	20.90	26.32	3.10
Sequence Tagger (BERT <sub>base</sub> )	44.39	25.96	32.76	2.53	30.07	20.60	24.45	3.61
└ (SpanBERT <sub>base</sub> )	46.96	25.98	33.45	2.39	35.24	20.52	25.94	3.07
└ (RoBERTa <sub>base</sub> )	48.43	25.19	33.14	2.25	34.91	19.33	24.88	2.92
Boundary-aware NER (Zheng et al., 2019)	32.80	18.67	23.79	2.46	34.40	7.23	11.95	1.11
BiFlaG (Luo and Zhao, 2020)	36.50	25.97	30.35	3.08	38.13	20.03	26.27	2.77
MRC NER (BERT <sub>base</sub> ) (Li et al., 2020)	37.71	25.84	30.67	2.96	29.53	21.20	24.68	3.78
└ (SpanBERT <sub>base</sub> )	37.04	27.94	31.85	3.26	31.79	20.52	24.94	3.40
└ (RoBERTa <sub>base</sub> )	39.59	27.39	32.38	2.99	32.38	22.43	26.50	3.65
SCOPE (BERT <sub>base</sub> )	36.96	39.99	38.41	4.68	30.74	32.32	31.51	5.54
└ (SpanBERT <sub>base</sub> )	41.15	39.70	<b>40.41</b>	4.17	33.95	35.84	34.87	5.56
└ (RoBERTa <sub>base</sub> )	36.10	<b>45.19</b>	40.14	<b>5.41</b>	33.51	<b>37.08</b>	<b>35.21</b>	<b>5.83</b>

# Experiments

- Is the extracted span high-quality?
  - Automatic metrics
    - When SCOPE extracts more new spans, it also slightly outperforms baselines on question-worthy score

Model	Avg. $score_q$	Avg. $score_w$	Avg. $\epsilon$
Golden	80.07 (0.21)	38.09 (0.13)	59.08 (0.13)
ENT Classifier	<b>78.52</b> (0.21)	32.93 (0.13)	55.72 (0.12)
Sequence Tagger	74.34 (0.21)	35.66 (0.12)	55.00 (0.12)
Boundary-aware NER	70.02 (0.25)	35.30 (0.13)	52.66 (0.14)
BiFlaG	75.95 (0.22)	34.83 (0.13)	55.39 (0.13)
MRC NER	75.07 (0.21)	<b>36.09</b> (0.12)	55.58 (0.12)
SCOPE	76.94 (0.21)	35.60 (0.12)	<b>56.27</b> (0.12)

- Performance boost on down-stream QA tasks

Backbone Model	Exact Match	F1
BERT <sub>large</sub> (Devlin et al., 2019)	78.7	81.9
BERT <sub>large</sub> (Our implementation)	77.9	81.3
└ ENT Classifier	77.8 (−0.1)	81.1 (−0.2)
└ Sequence Tagger	79.0 (+1.1)	82.2 (+0.9)
└ Boundary-aware NER	77.3 (−0.6)	80.8 (−0.5)
└ BiFlaG	79.1 (+1.2)	82.1 (+0.8)
└ MRC NER	79.3 (+1.4)	82.6 (+1.3)
└ SCOPE	<b>79.9 (+2.0)</b>	<b>83.2 (+1.9)</b>



# Analysis

- Is the model sensitive to prior distribution?
  - SCOPE **has a consistent performance** gain with different backbone PLMs and prior distribution  $\pi$

	<b>BERT<sub>base</sub></b>				<b>SpanBERT<sub>base</sub></b>				<b>RoBERTa<sub>base</sub></b>			
	Precision	Recall	F1	Avg. spans	Precision	Recall	F1	Avg. spans	Precision	Recall	F1	Avg. spans
$\pi' \times 1.50$	<b>39.93</b>	35.02	37.31	3.79	<b>44.67</b>	35.33	39.46	3.42	<b>42.18</b>	36.83	39.32	3.78
$\pi' \times 1.75$	39.21	36.25	37.67	4.00	41.77	38.32	39.97	3.97	40.28	38.62	39.43	4.15
$\pi' \times 2.00^\dagger$	36.96	39.99	<b>38.41</b>	4.68	41.15	39.70	<b>40.41</b>	4.17	36.10	45.19	<b>40.14</b>	5.41
$\pi' \times 2.25$	29.32	47.65	36.30	7.03	37.99	43.04	40.36	4.90	33.31	<b>47.63</b>	39.20	6.18
$\pi' \times 2.50$	29.41	<b>49.17</b>	36.81	7.23	34.39	<b>46.64</b>	39.59	5.86	32.18	47.06	38.22	6.32

# Summarize

---

- Re-formulate current AE task as a PU learning problem
- Propose SCOPE, a Structured Context graph network with Positive-unlabeled learning, to extract more answer spans from paragraphs
- Propose question-worthy score for automatically evaluate the quality of answer spans



Thanks for listening