# Collective Entity Disambiguation with Structured Gradient Tree Boosting
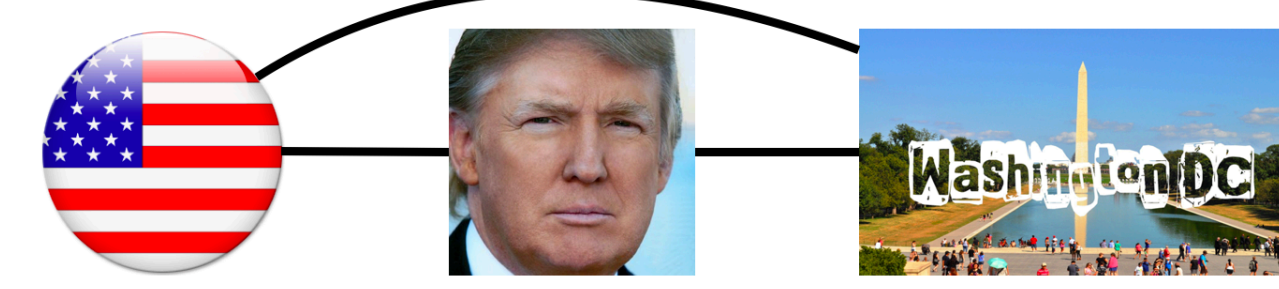
## Yi Yang, Ozan Irsoy, and Kazi Shefaet Rahman
### Bloomberg

**Bloomberg**
Engineering

---

## COLLECTIVE ENTITY DISAMBIGUATION

### Entity mentions are ambiguous

$\mathbf{x}$   US president Trump left Washington

$\mathbf{y}$

- Local and global context for disambiguation

### Structured prediction

- Entity dependencies:

- Inference: $\hat{\mathbf{y}} = \arg\max S(\mathbf{x}, \mathbf{y})$
- Learning:

$$\max S(\mathbf{x}, \quad + \quad + \quad )$$

---

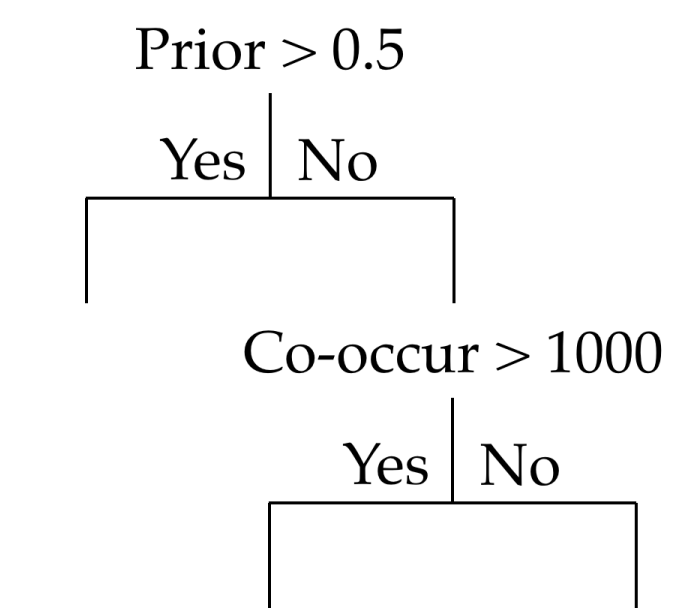## WHY GRADIENT TREE BOOSTING

### Heterogeneous features: $\phi(\mathbf{x}, \mathbf{y})$

Prior (   | "Trump" ) $\in [0, 1]$

Co-occur (  ,   ) $\in \{0, 1, ..., 10000, ...\}$

- Ideal models can handle:
  - Categorical features and count data
  - Nonlinear relationships between features

### Regression-tree-based models

Prior > 0.5
Yes | No
Co-occur > 1000
Yes | No

### Challenges

- Long-term dependencies between entities
- Approximate inference algorithms

---

## MODEL

### Structured Gradient Tree Boosting (SGTB)

$$S(\mathbf{x}, \mathbf{y}) = \sum_{t=1}^{T} F(\mathbf{x}, y_t, \mathbf{y}_{1:t-1})$$

$$p(\mathbf{y}|\mathbf{x}) = \frac{\exp\{\sum_{t=1}^{T} F(\mathbf{x}, y_t, \mathbf{y}_{1:t-1})\}}{Z(\mathbf{x})}$$

- Model $F$ using Gradient Tree Boosting.
- Optimize with Functional gradient descent:

$$F_m(\mathbf{x}, y_t, \mathbf{y}_{1:t-1}) = F_{m-1}(\mathbf{x}, y_t, \mathbf{y}_{1:t-1}) - \eta_m g_m(\mathbf{x}, y_t, \mathbf{y}_{1:t-1})$$
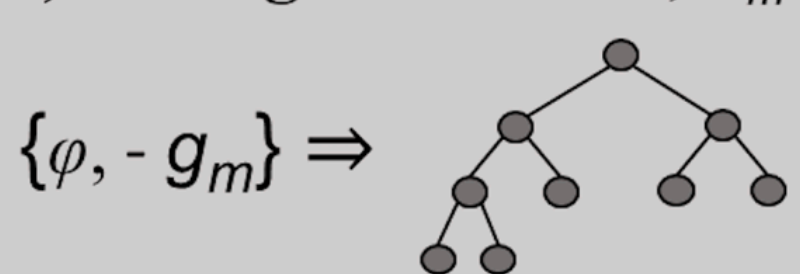
**(i) Beam search using $F_{m-1}$**

["United States", "Donald Trump", "Washington, D.C."]
["United States", "Donald Trump", "Washington (state)"]
["Us Weekly", "Trump, CO", "Washington, D.C."]
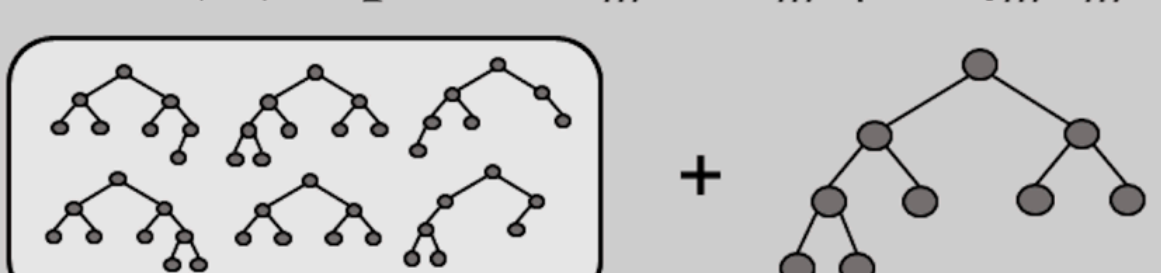["Us (Novel)", "Trump University", "Washington (state)"]

↓

**(ii) Compute functional gradient points, $g_m$**

$g_m(\mathbf{x}$, "Washington, D.C.", ["United States", "Donald Trump"])
$g_m(\mathbf{x}$, "Washington (state)", ["United States", "Donald Trump"])
$g_m(\mathbf{x}$, "Washington, D.C.", ["Us Weekly", "Trump, CO"])
$g_m(\mathbf{x}$, "Washington (state)", ["Us (novel)", "Trump University"])

↓

**(iii) Fit regression tree, $h_m$**

$\{\varphi, -g_m\} \Rightarrow$

↓

**(iv) Update $F_m = F_{m-1} + \eta_m h_m$**

+

---

## INFERENCE

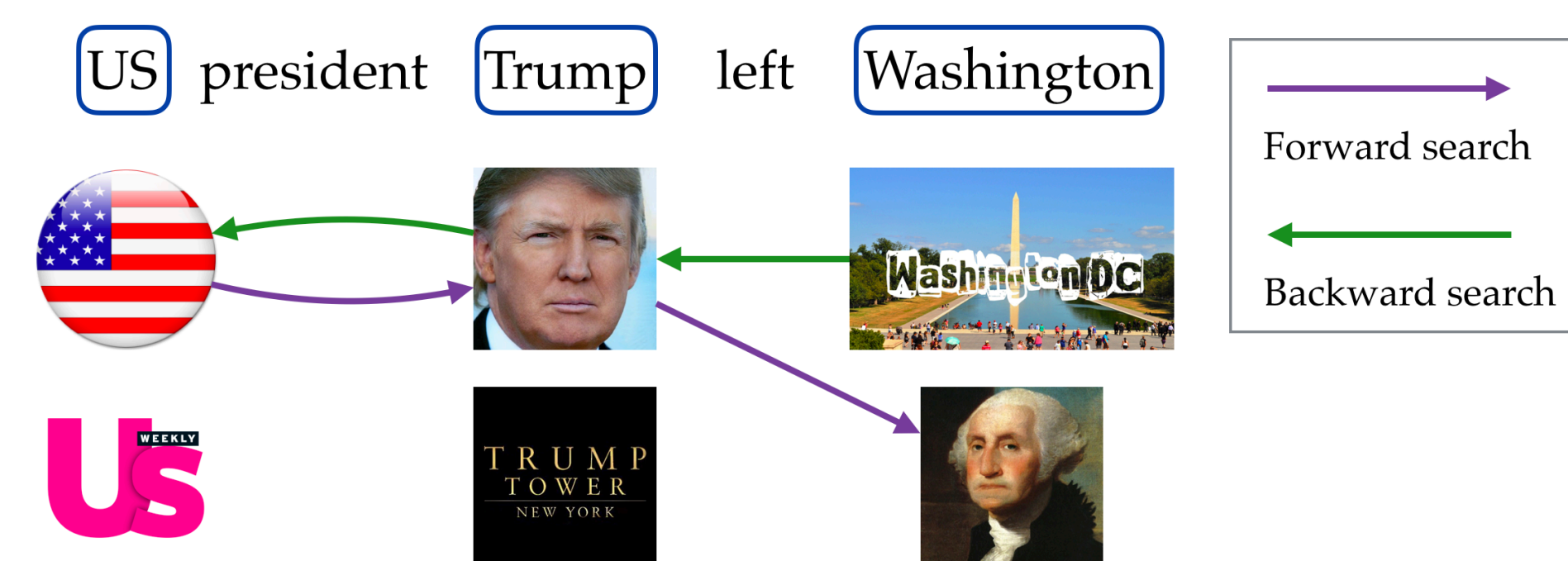### Computing $g_m$ requires inference

$$g_m(\mathbf{x}, y_t, \mathbf{y}_{1:t-1}) = \frac{\partial L(\mathbf{y}^*, S(\mathbf{x}, \mathbf{y}))}{\partial F(\mathbf{x}, y_t, \mathbf{y}_{1:t-1}))}$$

$$= p(\mathbf{y}_{1:t}|\mathbf{x}) - \mathbf{1}[\mathbf{y}_{1:t} = \mathbf{y}_{1:t}^*]$$

$$L(\mathbf{y}^*, S(\mathbf{x}, \mathbf{y})) = \log Z(\mathbf{x}) - S(\mathbf{x}, \mathbf{y}^*)$$

- Exact inference is intractable.
- Beam search is used to approximate $p(\mathbf{y}_{1:t}|\mathbf{x})$.

### Bi-directional beam search

US president Trump left Washington

→ Forward search
← Backward search

- New scoring function: $S(\mathbf{x}, \mathbf{y}) = (\overrightarrow{S} + \overleftarrow{S})/2$

**Input** : input document $\mathbf{x}$, candidate sequences $\{\mathbf{y}\}$,
        joint scoring function $S(\mathbf{x}, \mathbf{y}_{t_1:t_2})$
**Output**: beam sequence set $C$
$C \leftarrow \emptyset$
**while** *not converged* **do**
   // forward beam search
   **for** $t = 1, \cdots, T$ **do**
     $C^{(F)} \leftarrow \text{top-B}_{\mathbf{y}_{1:t}}[S(\mathbf{x}, \mathbf{y}_{1:t}) + S(\mathbf{x}, \mathbf{y}_{T:t})]$
     // add gold subsequence
     $C^{(F)} \leftarrow C^{(F)} \cup \{\mathbf{y}_{1:t}^*\}$
     $C \leftarrow C \cup C^{(F)}$
   **end**
   // backward beam search
   **for** $t = T, \cdots, 1$ **do**
     $C^{(B)} \leftarrow \text{top-B}_{\mathbf{y}_{T:t}}[S(\mathbf{x}, \mathbf{y}_{T:t}) + S(\mathbf{x}, \mathbf{y}_{1:t})]$
     // add gold subsequence
     $C^{(B)} \leftarrow C^{(B)} \cup \{\mathbf{y}_{T:t}^*\}$
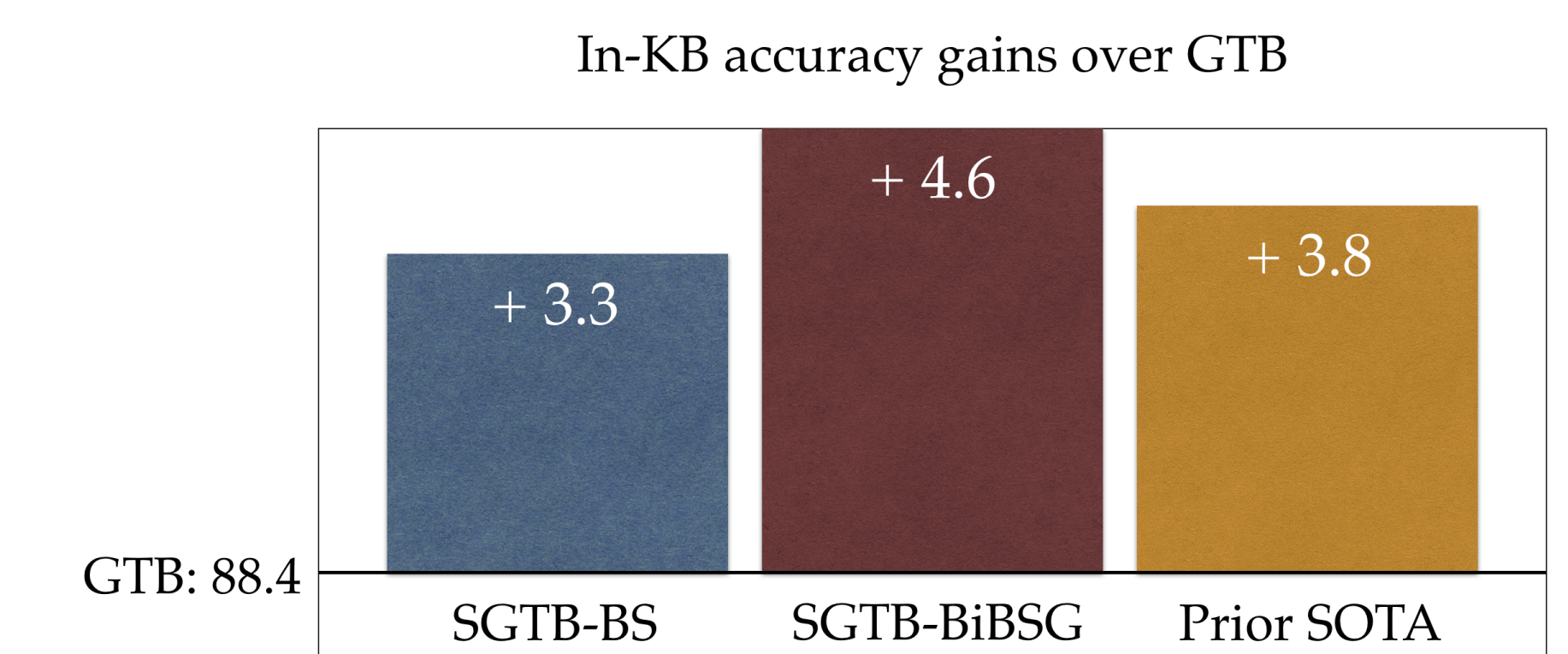     $C \leftarrow C \cup C^{(B)}$
   **end**
**end**

---

## EXPERIMENTS

### Data

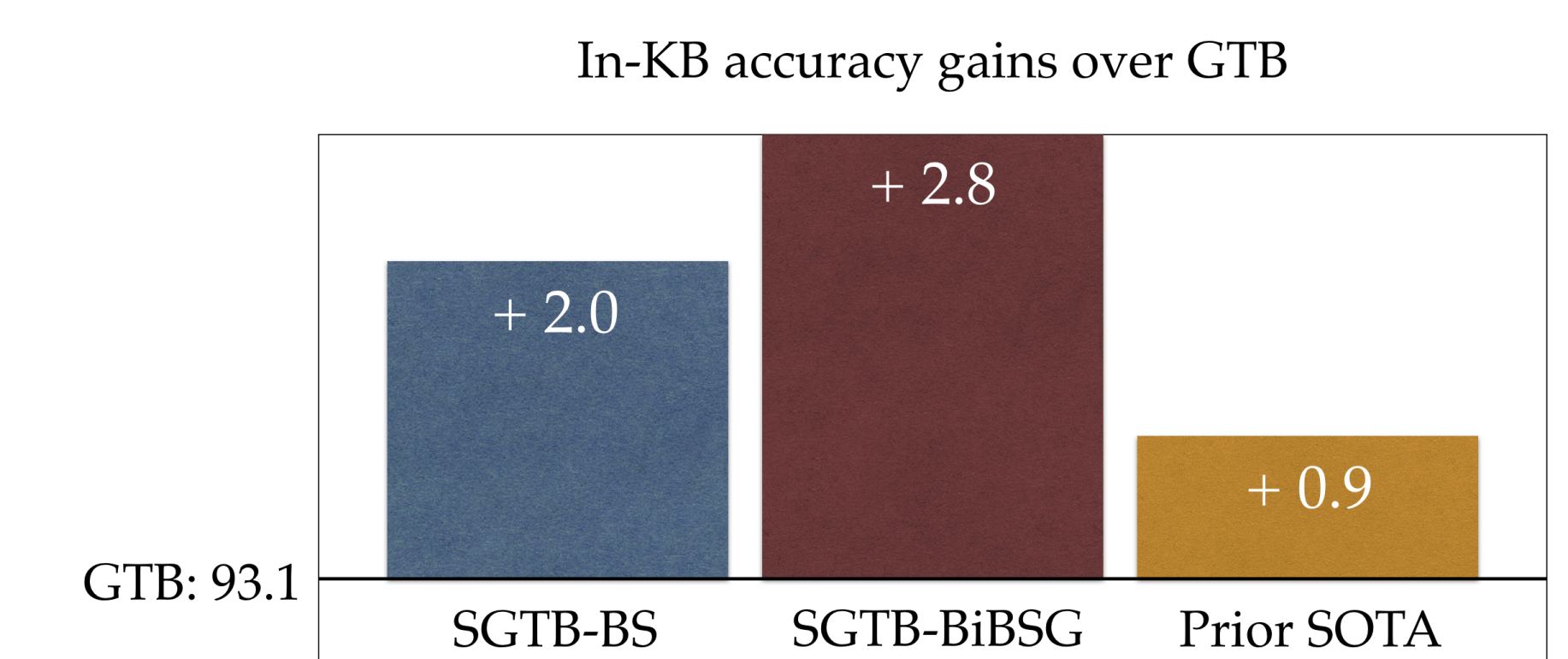| Class | Dataset | # mention | # document | # mention per document |
|---|---|---|---|---|
| Training | AIDA-train | 18,448 | 946 | 19.5 |
| Validation | AIDA-dev | 4,791 | 216 | 22.1 |
| In-domain testing | AIDA-test | 4,485 | 231 | 19.4 |
| Cross-domain testing | AQUAINT | 727 | 50 | 14.5 |
| | MSNBC | 656 | 20 | 32.8 |
| | ACE | 257 | 36 | 7.1 |
| | CWEB | 11,154 | 320 | 34.8 |
| | WIKI | 6,821 | 320 | 21.3 |

### Setup

- Metrics
  - In-KB accuracy
  - Bag-of-Title (BoT) F1 score
- Competing systems
  - Gradient Tree Boosting (GTB)
  - SGTB with Beam search (SGTB-BS)
  - Bidir. BS using Gold path (SGTB-BiBSG)
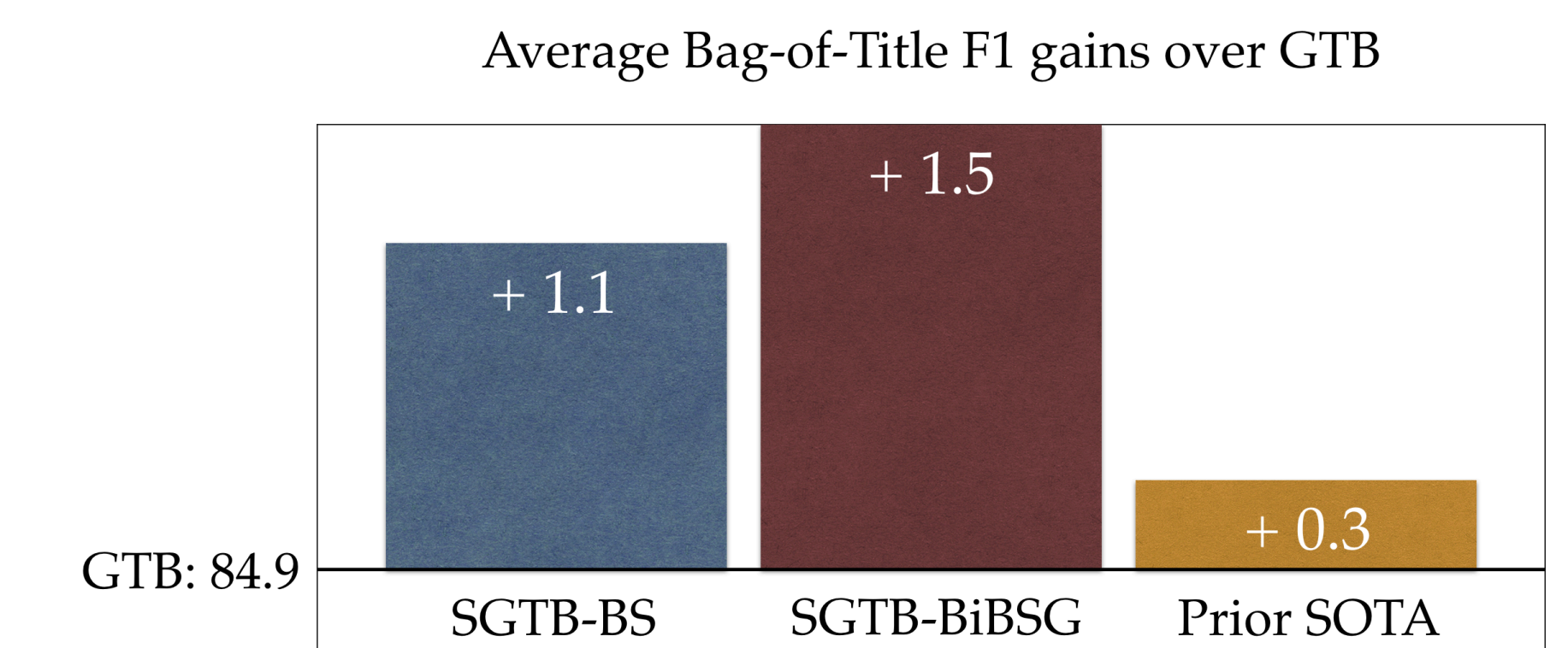  - Previous state-of-the-art (SOTA) systems

### In-domain results

In-KB accuracy gains over GTB

SGTB-BS: +3.3   SGTB-BiBSG: +4.6   Prior SOTA: +3.8
GTB: 88.4

- Generate candidates with PPRforNED

In-KB accuracy gains over GTB

SGTB-BS: +2.0   SGTB-BiBSG: +2.8   Prior SOTA: +0.9
GTB: 93.1

### Cross-domain results

Average Bag-of-Title F1 gains over GTB

SGTB-BS: +1.1   SGTB-BiBSG: +1.5   Prior SOTA: +0.3
GTB: 84.9

---

## SUMMARY

- We present a novel Structured Gradient Tree Boosting (SGTB) model for collectively disambiguating entities in a document.
- SGTB combines structured learning with Gradient Tree Boosting to produce globally optimal entity assignments for all the mentions.
- We present Bidirectional Beam Search with Gold path (BiBSG), an efficient approximate inference algorithm tailored for SGTB.
- SGTB achieves state-of-the-art (SOTA) results on popular entity disambiguation datasets of different domains.