# Leveraging discourse information effectively for authorship attribution

Elisa Ferracane, Su Wang, Raymond J. Mooney

University of Texas at Austin

# Task

- **Authorship Attribution:** identify the author of a text, given a set of author-labeled training texts.

# Authorship Attribution

- **Neural networks** (e.g., character-level CNNs) have proven very powerful…

  - capture stylometric cues at the surface level

| | |
|---|---|
| "My very photogenic mother died in a freak accident **(**picnic, lightning**)** when I was three…" | *Lolita*, Nabokov |
| "But what principally attracted attention of Nicholas, was the old gentleman's eye… Grafted upon the quaintness and oddity of his appearance, was something…" | *Nichola Nickleby*, Dickens |

# Authorship Attribution

- Authors also have particular **rhetorical** styles…

- But how do you incorporate discourse into a neural net?

# Our Contributions

1) How can you *featurize* discourse information?

2) How can you *integrate* discourse information into the network?

3) Can discourse help in SOTA model (bigram character CNN)?

# Q1: How can you *featurize* discourse information?

- Use an entity grid model (Barzilay & Lapata, 2008) with either:

  - grammatical relations, or

  - RST discourse relations

# Q1: How can you *featurize* discourse information?

(1) My father was a clergyman of the north of England, who was deservedly respected by all who knew him; and, in his younger days, lived pretty comfortably on the joint income of a small incumbency and a snug little property of his own.

(2) My mother, who married him against the wishes of her friends, was a squire's daughter, and a woman of spirit.

(3) In vain it was represented to her, that if she became the poor parson's wife, she must relinquish her carriage and her lady's-maid, and all the luxuries and elegancies of affluence; which to her were little less than the necessaries of life.

# Q1: How can you *featurize* discourse information?

(1) **My father** was a clergyman of the north of England, **who** was deservedly respected by all who knew **him**; and, in his younger days, lived pretty comfortably on the joint income of a small incumbency and a snug little property of his own.

(2) My mother, who married **him** against the wishes of her friends, was a squire's daughter, and a woman of spirit.

(3) In vain it was represented to her, that if she became the **poor parson**'s wife, she must relinquish her carriage and her lady's-maid, and all the luxuries and elegancies of affluence; which to her were little less than the necessaries of life.

# Q1: How can you *featurize* discourse information?

(1) **My father** was a clergyman of the north of England, **who** was deservedly respected by all who knew **him**; and, in his younger days, lived pretty comfortably on the joint income of a small incumbency and a snug little property of his own.

(2) **My mother**, who married **him** against the wishes of her friends, was a squire's daughter, and a woman of spirit.

(3) In vain it was represented to **her**, that if **she** became the **poor parson**'s wife, **she** must relinquish her carriage and her lady's-maid, and all the luxuries and elegancies of affluence; which to **her** were little less than the necessaries of life.

# Q1: How can you *featurize* discourse information?

father    mother

| | |
|---|---|
| (1) | |
| (2) | |
| (3) | |

row: sentence
column: salient entity

Barzilay and Lapata (2008)

# Q1: How can you *featurize* discourse information?

(1) **[My father]SUBJECT** was a clergyman of the north of England, who was deservedly respected by all who knew him; and, in his younger days, lived pretty comfortably on the joint income of a small incumbency and a snug little property of his own.

(2) **[My mother]SUBJECT**, who married **[him]OBJECT** against the wishes of her friends, was a squire's daughter, and a woman of spirit.

(3) In vain it was represented to her, that if **[she]SUBJECT** became the **[poor parson]OTHER**'s wife, she must relinquish her carriage and her lady's-maid, and all the luxuries and elegancies of affluence; which to her were little less than the necessaries of life.

# Q1: How can you *featurize* discourse information?

(1) **[My father]SUBJECT** was a clergyman of the north of England, who was deservedly respected by all who knew him; and, in his younger days, lived pretty comfortably on the joint income of a small incumbency and a snug little property of his own.

(2) **[My mother]SUBJECT**, who married **[him]OBJECT** against the wishes of her friends, was a squire's daughter, and a woman of spirit.

(3) In vain it was represented to her, that if **[she]SUBJECT** became the **[poor parson]OTHER**'s wife, she must relinquish her carriage and her lady's-maid, and all the luxuries and elegancies of affluence; which to her were little less than the necessaries of life.

# Q1: How can you *featurize* discourse information?

|  | father | mother |
|---|---|---|
| (1) | S | - |
| (2) | O | S |
| (3) | X | S |

Grammatical relations

13

# Q1: How can you *featurize* discourse information?

- Discourse relations:

  - Rhetorical Structure Theory (RST)

    - Divide a document into elementary discourse units (EDUs), usually clauses

    - Organize EDUs into a **tree** structure:

      - edges are discourse relation types

      - node in a relation can be either the nucleus (more "important") or satellite

14

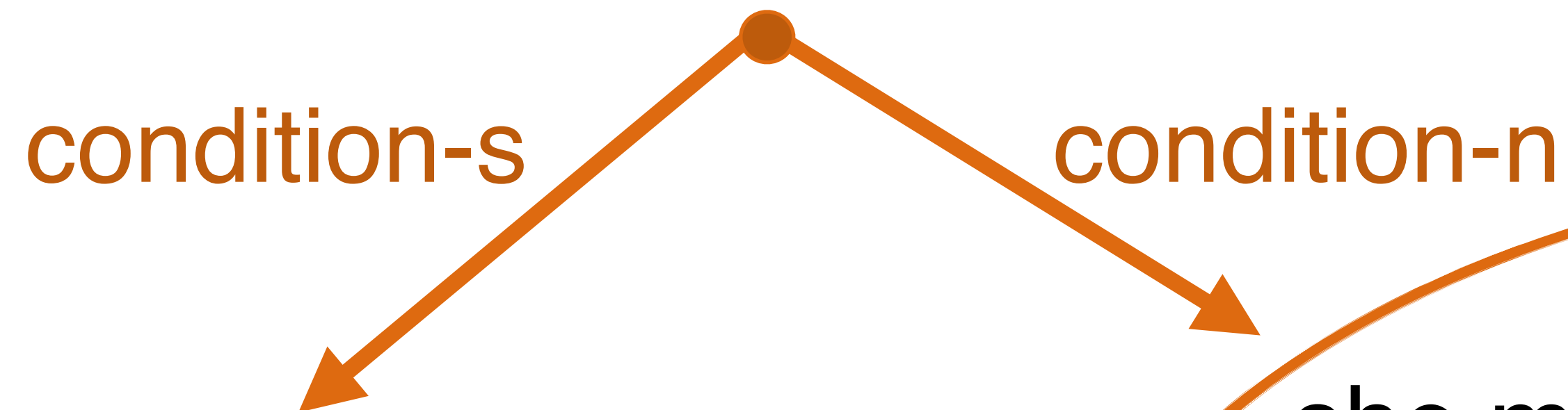# Q1: How can you *featurize* discourse information?

if she became the poor parson's wife, she must relinquish her carriage and her lady's-maid, and all the luxuries and elegancies of affluence; which to her were little less than the necessaries of life.

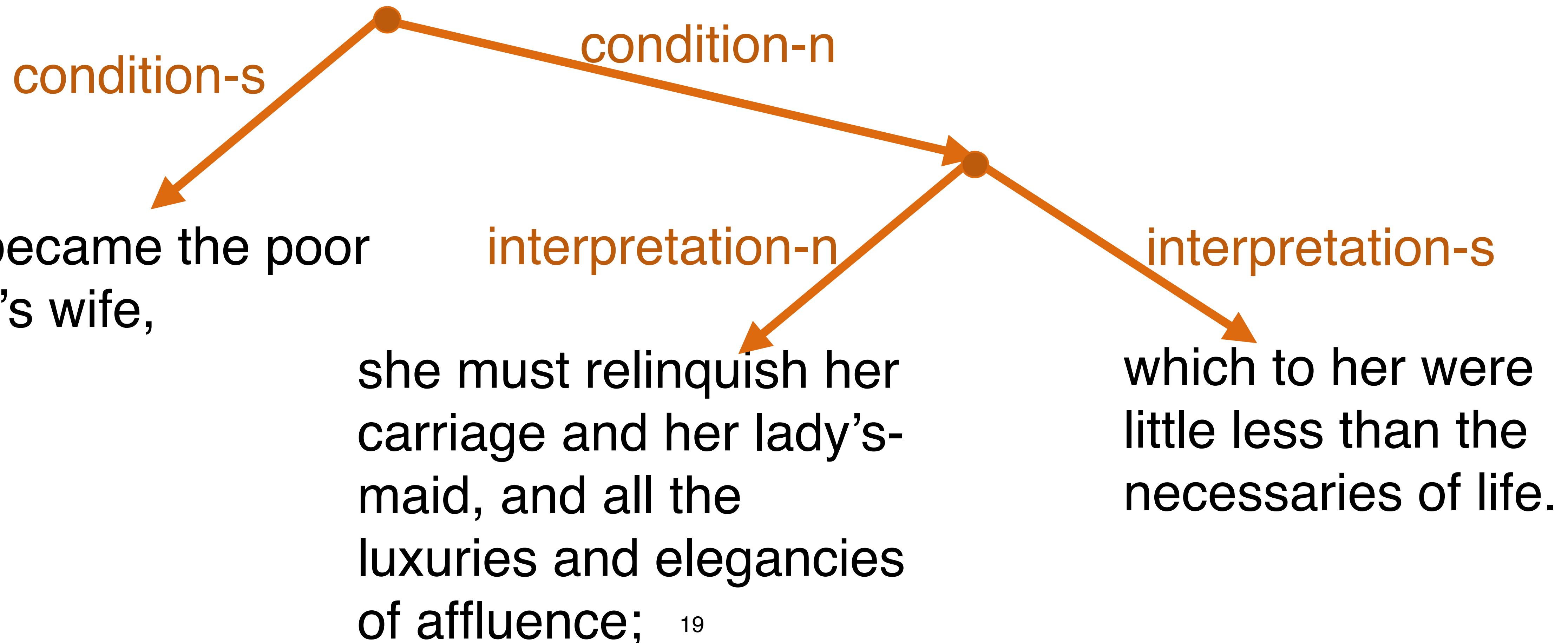# Q1: How can you *featurize* discourse information?

if she became the poor parson's wife, she must relinquish her carriage and her lady's-maid, and all the luxuries and elegancies of affluence; which to her were little less than the necessaries of life.

# Q1: How can you *featurize* discourse information?

if she became the poor parson's wife,

she must relinquish her carriage and her lady's-maid, and all the luxuries and elegancies of affluence;

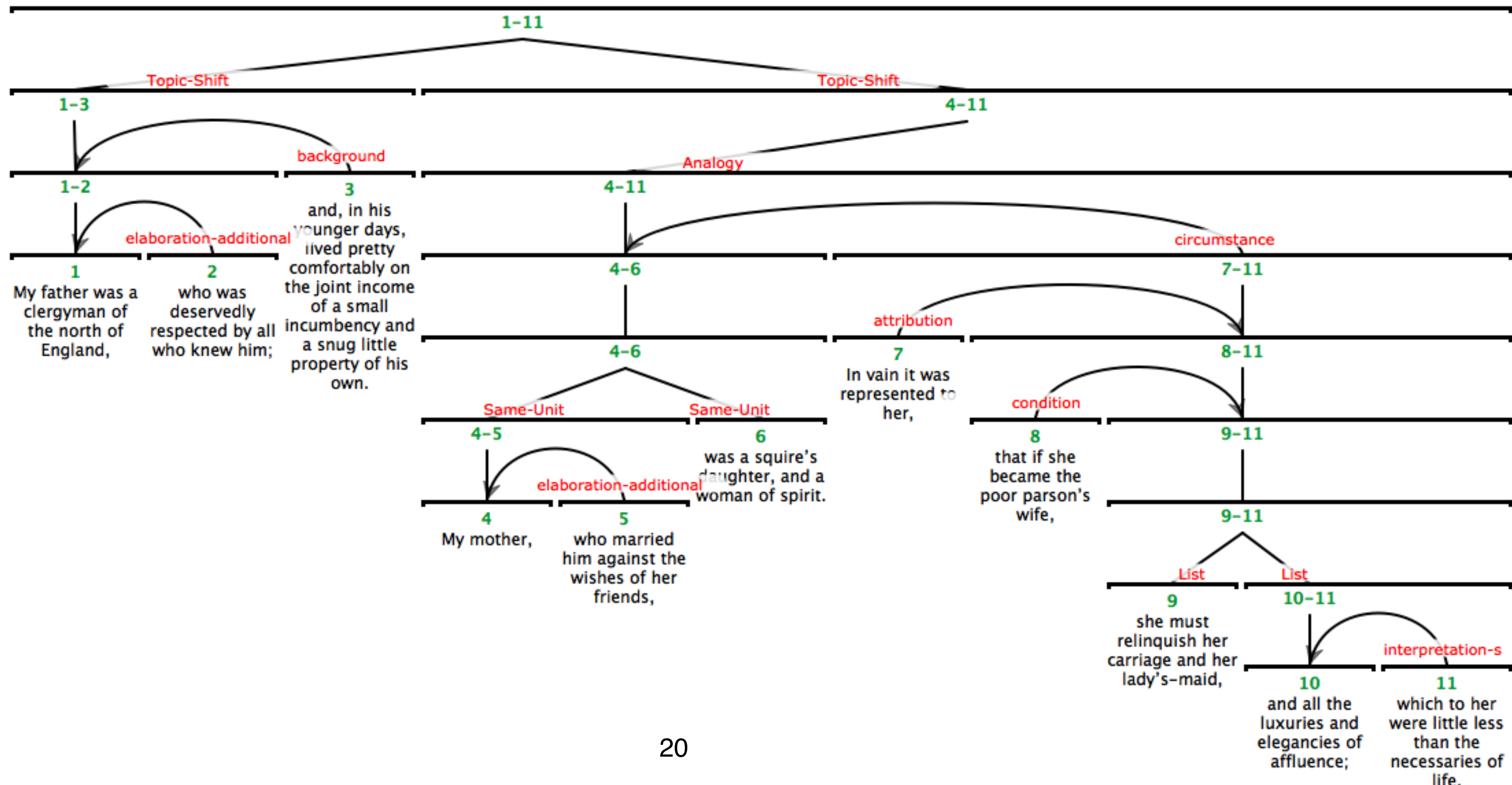which to her were little less than the necessaries of life.

# Q1: How can you *featurize* discourse information?

condition-s

condition-n

if she became the poor parson's wife,

she must relinquish her carriage and her lady's-maid, and all the luxuries and elegancies of affluence;

which to her were little less than the necessaries of life.

18

# Q1: How can you *featurize* discourse information?

condition-s

condition-n

if she became the poor parson's wife,

interpretation-n

interpretation-s

she must relinquish her carriage and her lady's-maid, and all the luxuries and elegancies of affluence;

which to her were little less than the necessaries of life.

# Q1: How can you *featurize* discourse information?

# Q1: How can you *featurize* discourse information?

| | father | mother |
|---|---|---|
| (1) | background.N, TopicShift, elaboration.S, background.S | - |
| (2) | elaboration.S | elaboration.N, circumstance.N, TopicShift |
| (3) | condition.N | attribution.S, condition.N, interpretation.S |

RST discourse relations

21

Feng and Hirst (2014)

# Q2: How can you *integrate* discourse information into the network?

- Use probability vector

- Use embeddings!

# Q2: How can you *integrate* discourse information into the network?

CNN without discourse



Ruder et al., 2016; Shrestha et al., 2017, Sari et al., 2017
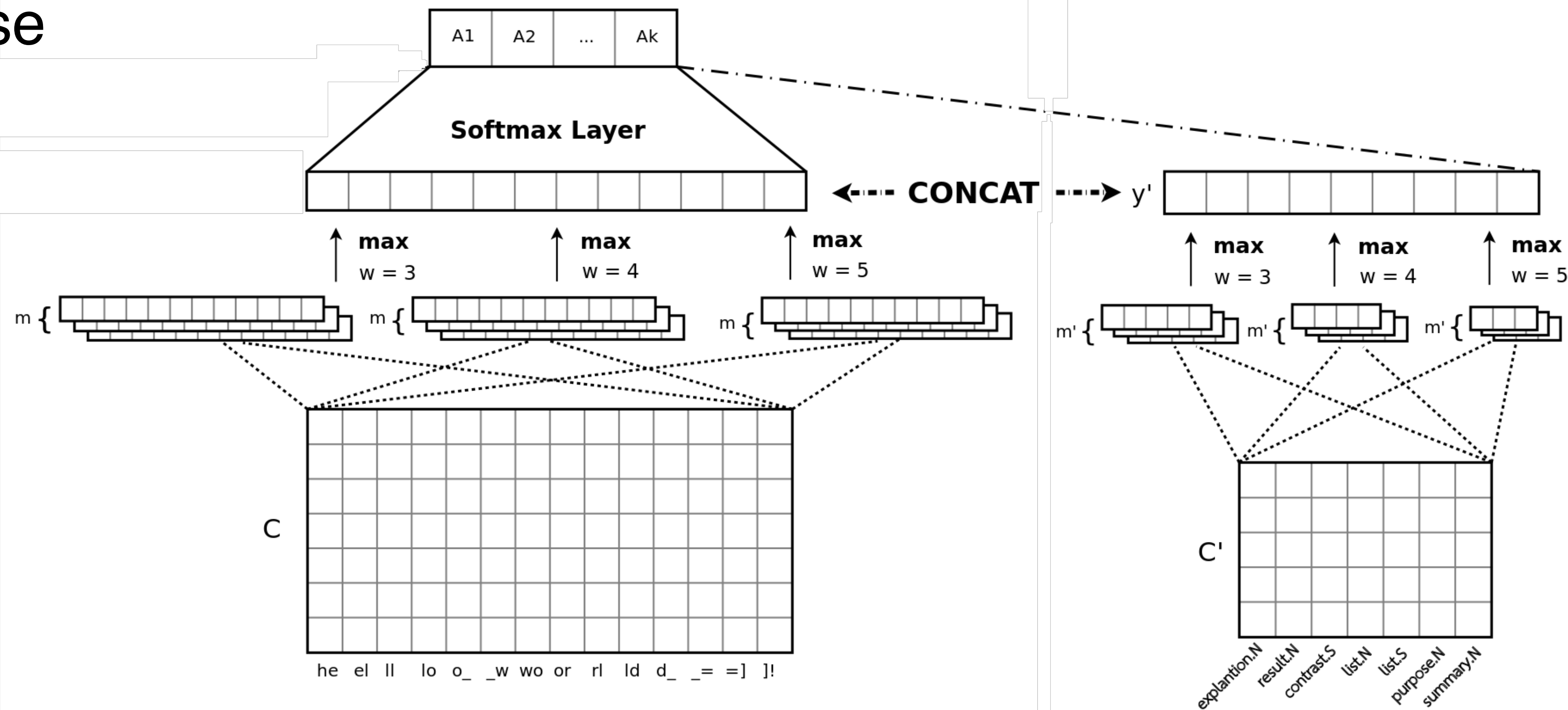
# Q2: How can you *integrate* discourse information into the network?

CNN with discourse probability vector

# Q2: How can you *integrate* discourse information into the network?
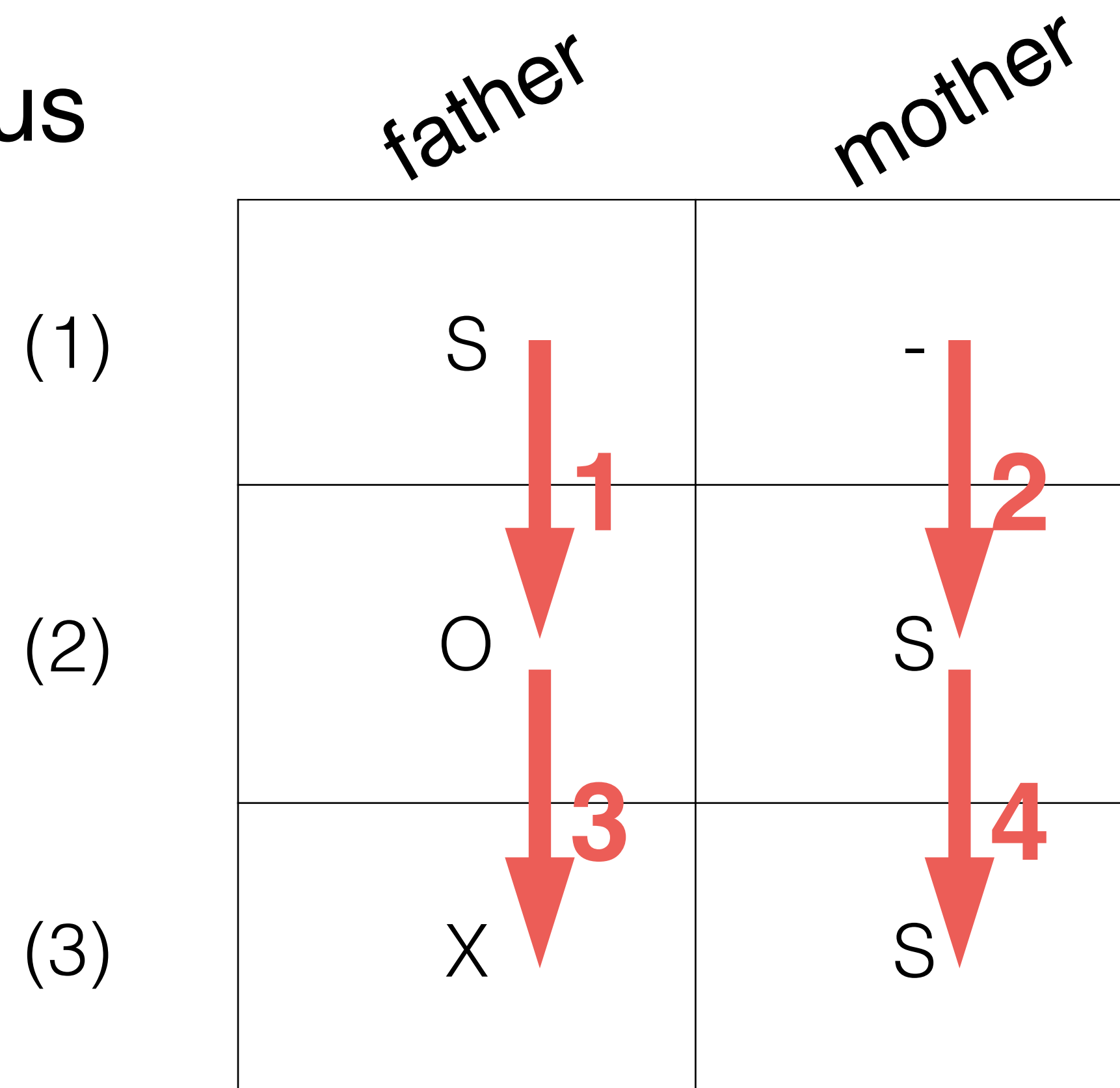
CNN with discourse embeddings

# Q2: How can you *integrate* discourse information into the network?

- Use embeddings

  - Local vs. Global

  - Local: how are entities changing across **contiguous** sentences?

  - **Global**: how is each entity changing across a **document**?

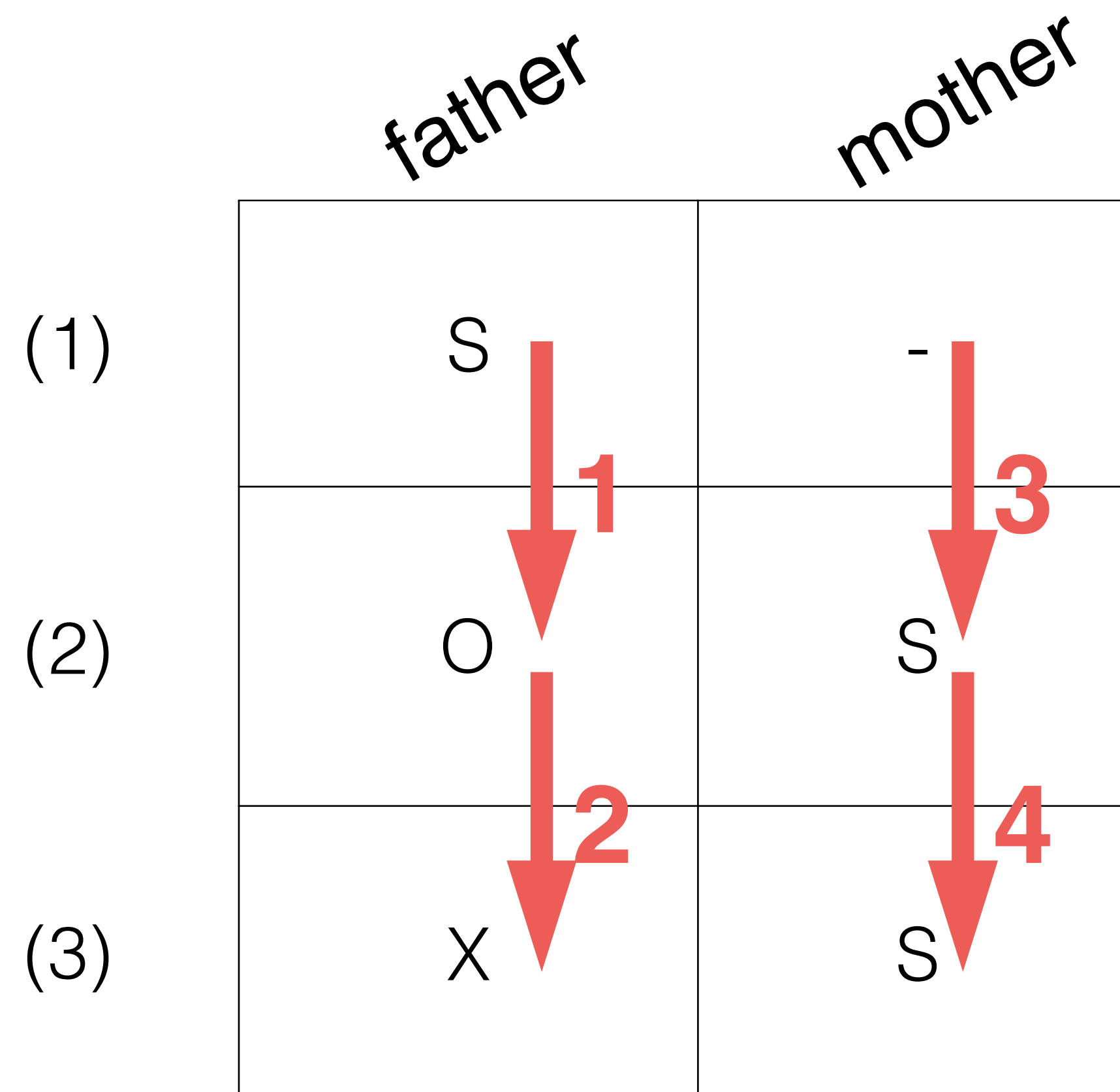# Q2: How can you *integrate* discourse information into the network?

**Local:** by contiguous sentences



Sequence: so, -s, ox, ss

|  | father | mother |
|---|---|---|
| (1) | S | - |
|  | ↓ **1** | ↓ **2** |
| (2) | O | S |
|  | ↓ **3** | ↓ **4** |
| (3) | X | S |

# Q2: How can you *integrate* discourse information into the network?

**Global:** by entity



|  | father | mother |
|---|---|---|
| (1) | S | - |
|  | ↓ **1** | ↓ **3** |
| (2) | O | S |
|  | ↓ **2** | ↓ **4** |
| (3) | X | S |

Sequence: so, ox, -s, ss

# Datasets

| Dataset | # authors | mean words/ auth | mean words/ text |
|---------|-----------|------------------|------------------|
| IMDB62 | 62 | 349,004 | 349 |
| Novel-50 | 50 | 709,880 | 2,000 |

# Results



1) How to *featurize*?
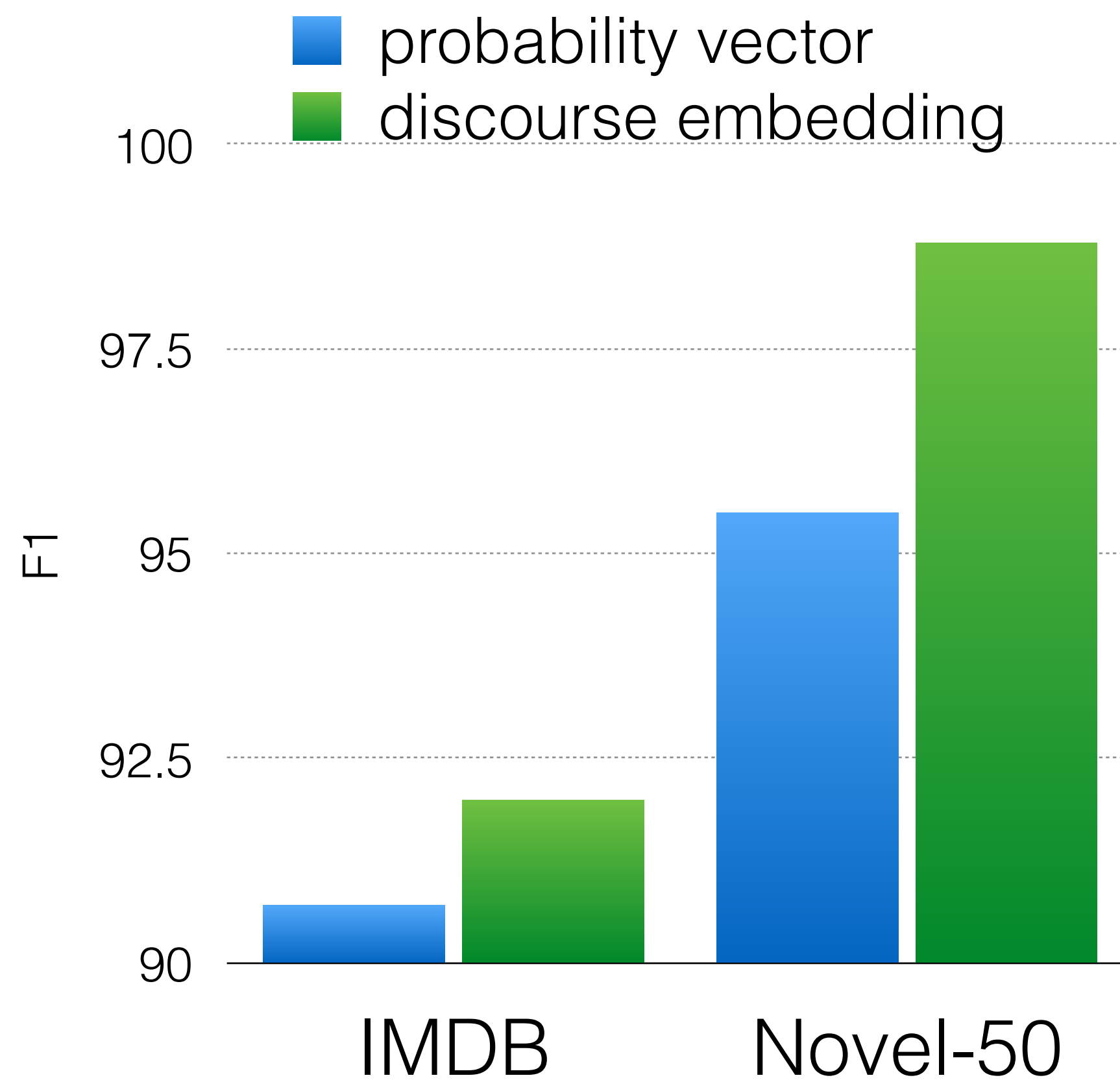grammatical relations
vs.
RST discourse relations

# Results



1) How to *featurize*?
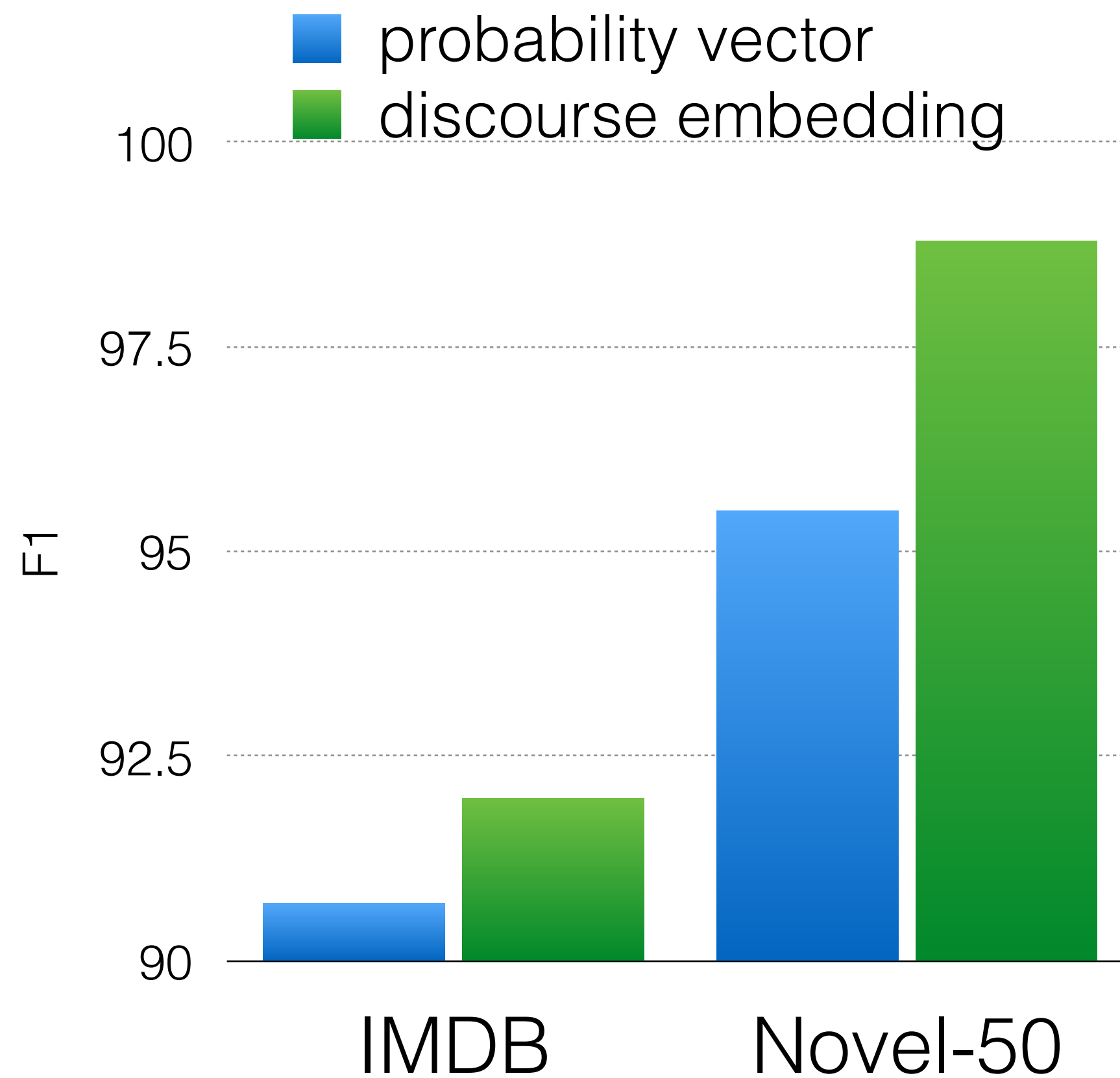grammatical relations
vs.
**RST discourse relations**

# Results



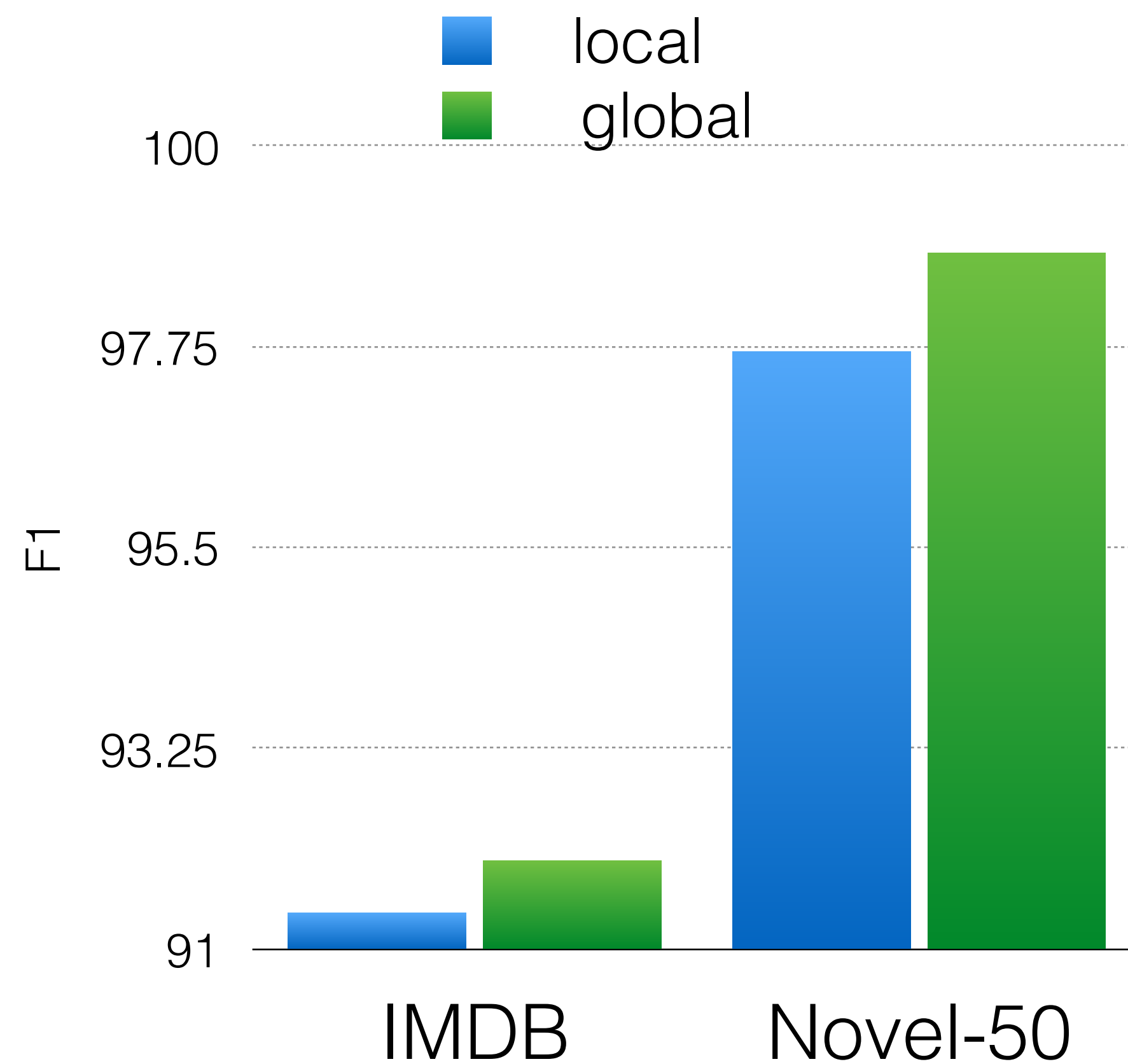2) How to *integrate*?
probability vector
vs.
discourse embedding

# Results



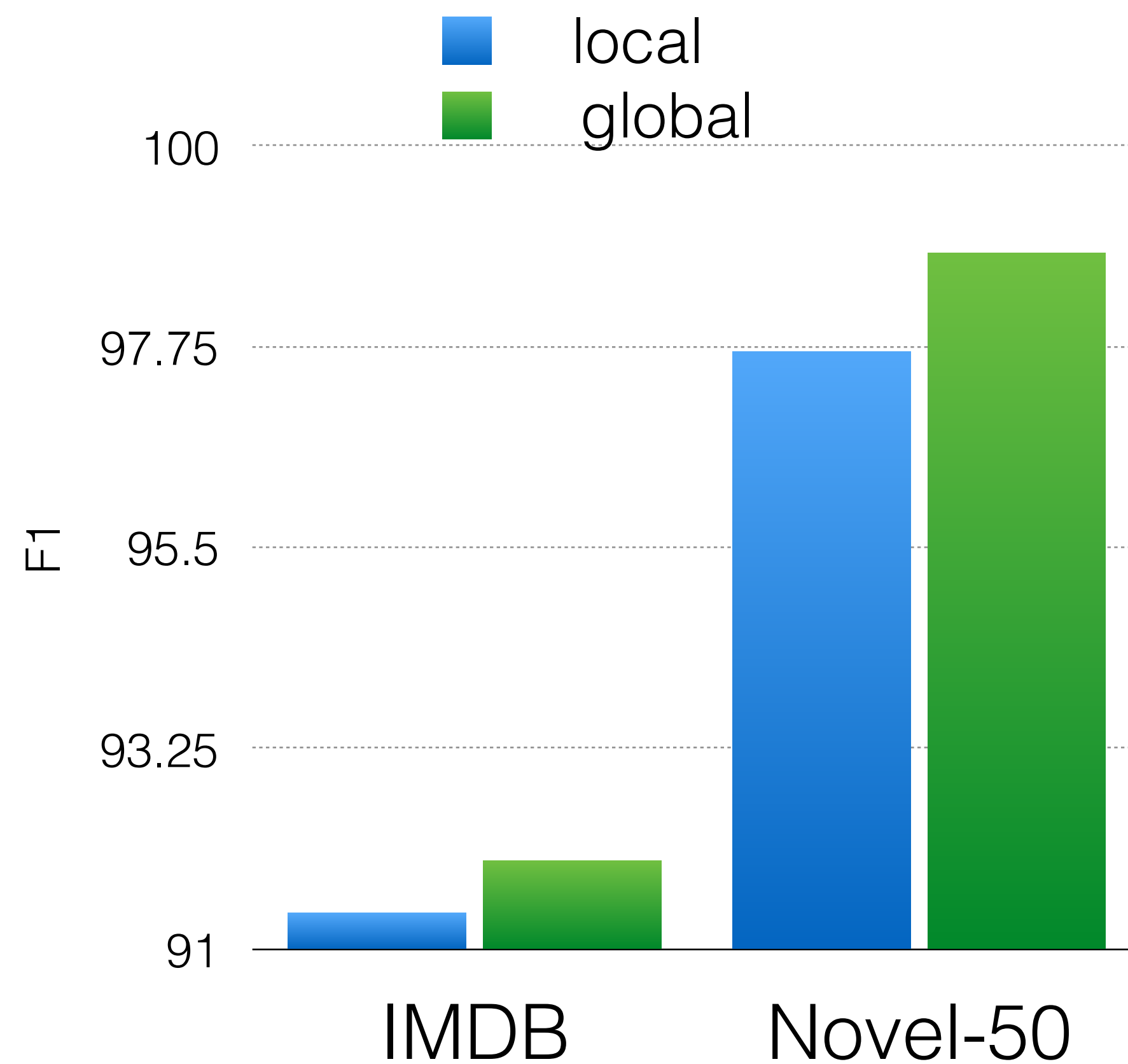F1 chart — legend: probability vector (blue), discourse embedding (green). Y-axis: 90, 92.5, 95, 97.5, 100. X-axis categories: IMDB, Novel-50.

2) How to *integrate*?
probability vector
vs.
**discourse embedding**

# Results



2) How to *integrate*?
local
vs.
global

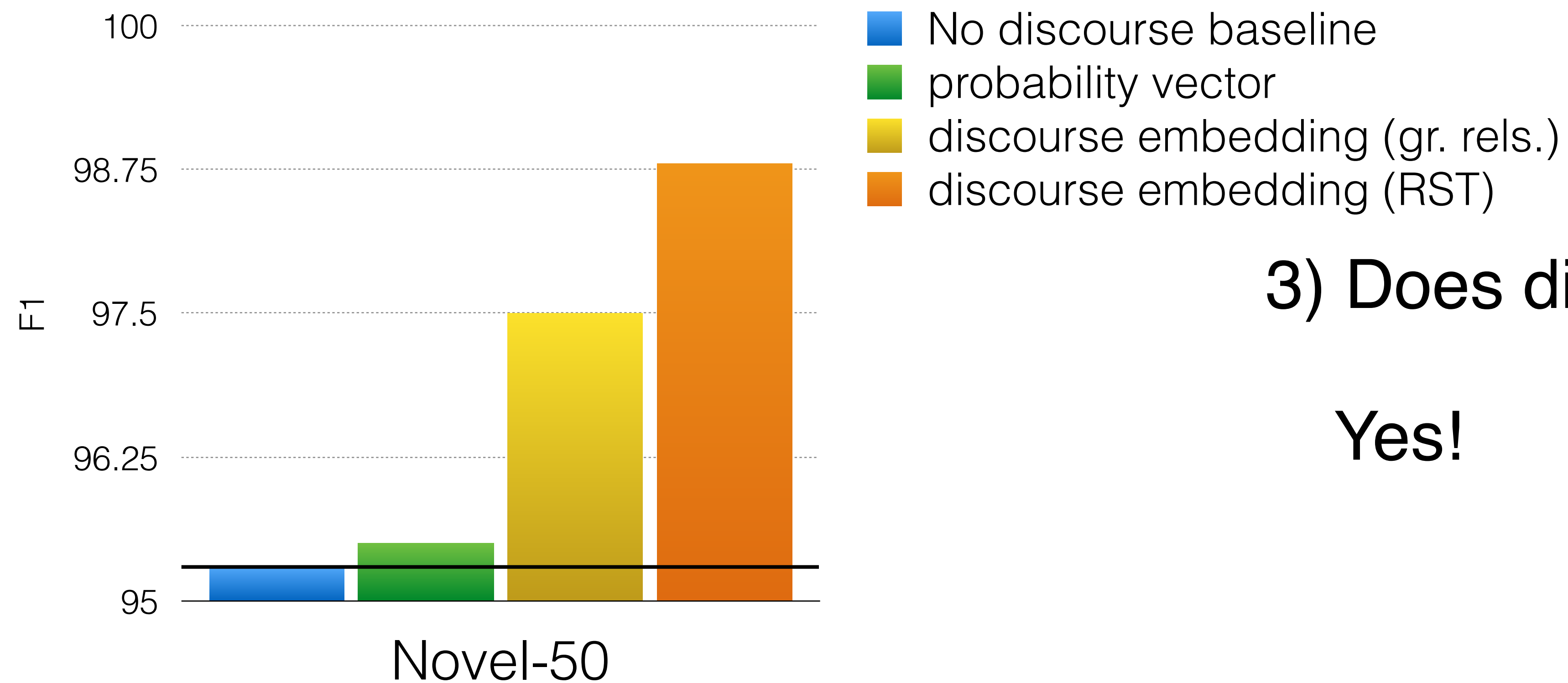# Results



2) How to *integrate*?
local
vs.
**global**

35

# Results



- ■ No discourse baseline
- ■ probability vector
- ■ discourse embedding (gr. rels.)
- ■ discourse embedding (RST)

## 3) Does discourse help?

## It depends…

# Results



No discourse baseline
probability vector
discourse embedding (gr. rels.)
discourse embedding (RST)

3) Does discourse help?

Yes!

# Error Analysis

- The least-represented author (Ambrose Bierce) obtains the biggest improvement from discourse:

  —Discourse feature is more **robust** with smaller, fewer samples compared to character bigrams

- Two authors who gained large improvements from discourse wrote a variety of genres (e.g., both supernatural horror and love stories)

  —Character bigrams can't generalize well to the different vocabularies, but discourse captures the similar rhetorical style

# Conclusion

- Discourse **improves** authorship attribution over a strong baseline of character-level CNN

- Embeddings of RST discourse relations at the global level perform the best

- Works better on longer documents

# Thank you!

Leveraging discourse information effectively for authorship attribution

elisa@ferracane.com