EACL 2026
MOROCCO
Palais Des Congres, Rabat
March 24 - 29, 2026

UBC a place of mind
THE UNIVERSITY OF BRITISH COLUMBIA

# BeDiscovER: The Benchmark of Discourse Understanding in the Era of Reasoning Language Models

Chuyuan Li, Giuseppe Carenini

Department of Computer Science
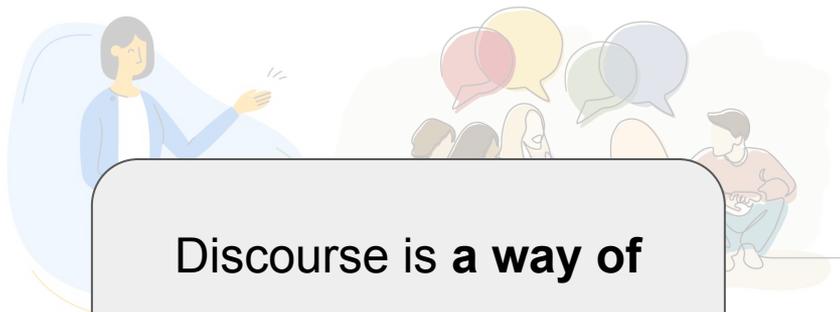The University of British Columbia

# Why discourse understanding?

Speech

Dialogue

Scientific papers

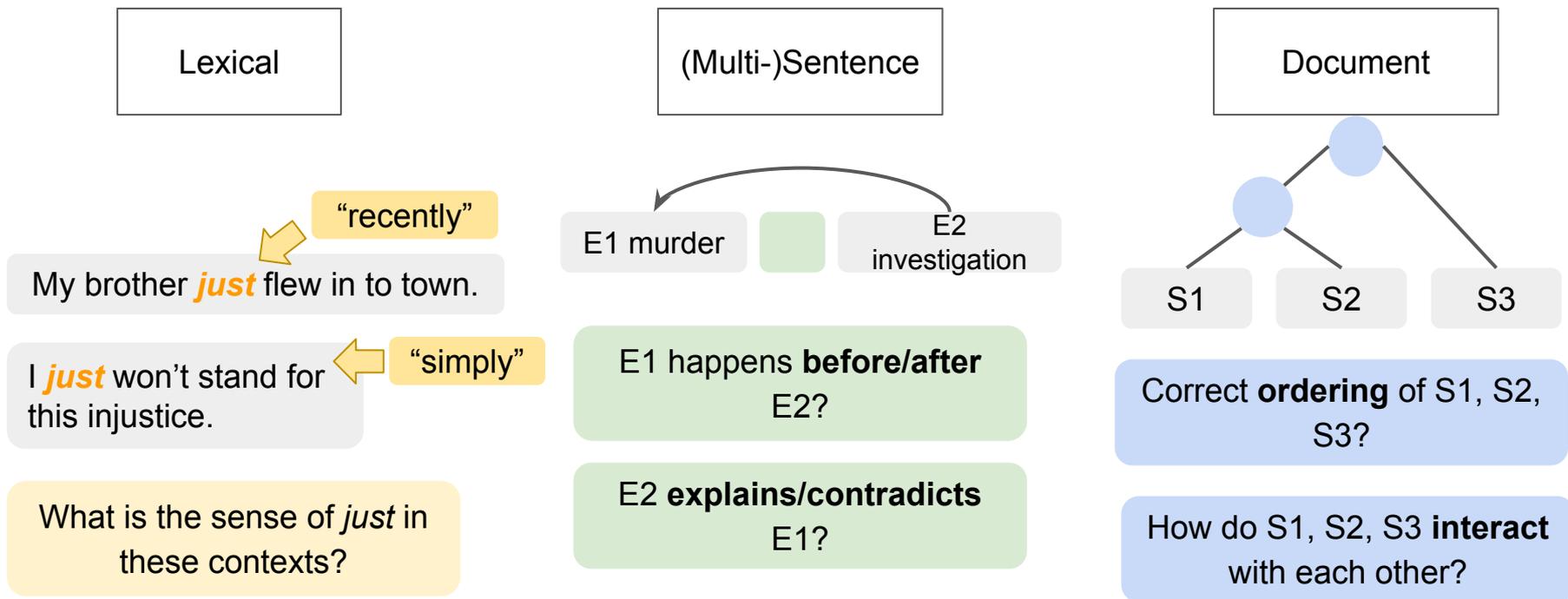News articles

Rhetoric / Persuasion

Discourse is **a way of speaking.**

**But also:**

- Text structure
- Communicative / behavior intentions
- Temporal reasoning
- Causal reasoning
- …

# Why discourse understanding? – a few examples

## Lexical

"recently"

My brother *just* flew in to town.

"simply"

I *just* won't stand for this injustice.

What is the sense of *just* in these contexts?

## (Multi-)Sentence

E1 murder          E2 investigation

E1 happens **before/after** E2?

E2 **explains/contradicts** E1?

## Document

S1      S2      S3

Correct **ordering** of S1, S2, S3?

How do S1, S2, S3 **interact** with each other?

Sheffield et al., Is it JUST semantics? a case study of discourse particle understanding in LLMs. Findings ACL 2025.

# Why discourse understanding? – a few examples

| Lexical | (Multi-)Sentence | Document |
|---------|------------------|----------|

*Discourse understanding requires **lexical & semantic, temporal, rhetorical, commonsense**… knowledge.*

***How well do modern LLMs understand discourse?***

My brother **just** flew

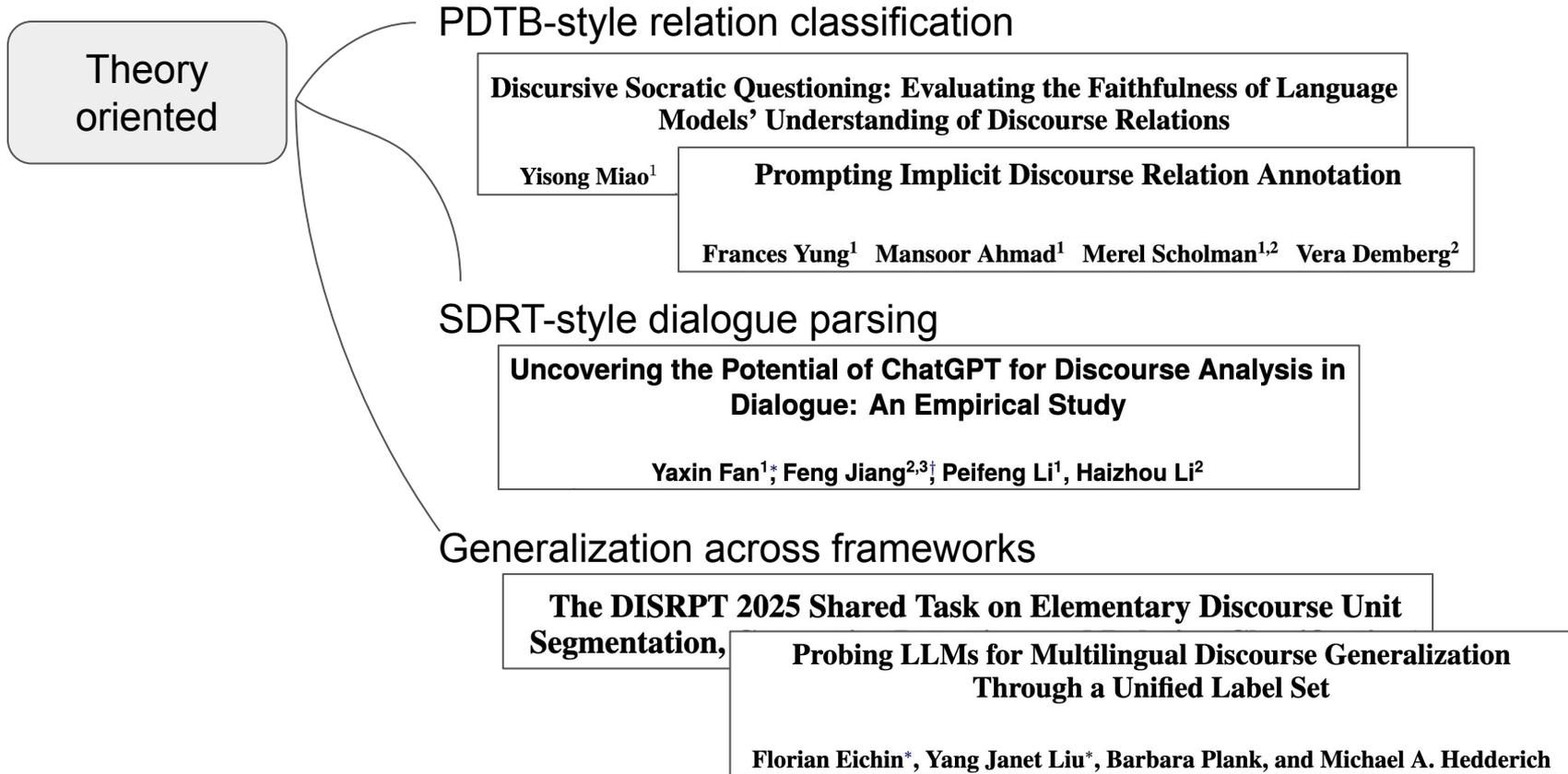"simply"

I **just** won't stand for this injustice.

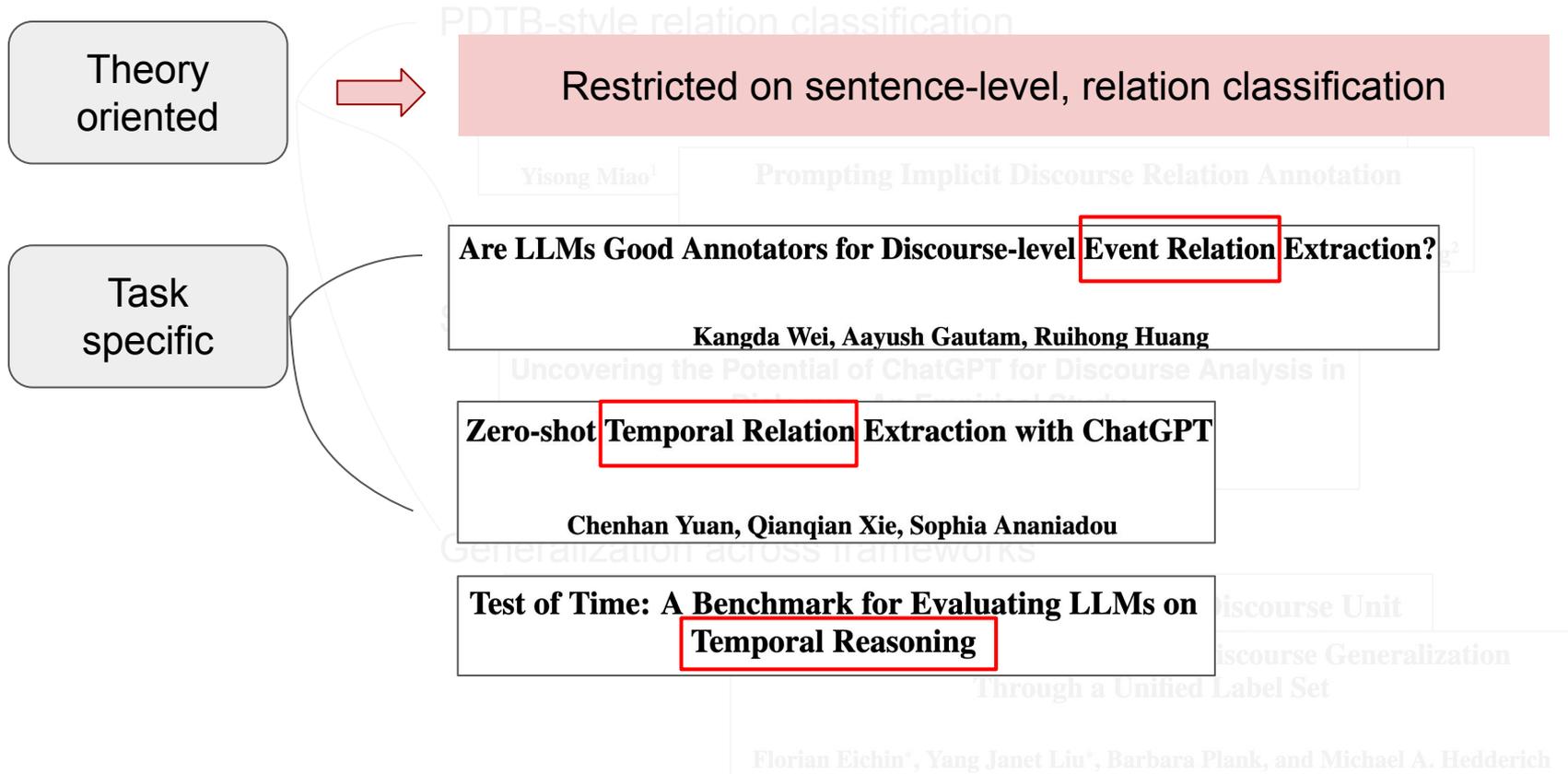S2          S3

Correct **ordering** of S1, S2, S3?

with each other?

## Theory oriented

### PDTB-style relation classification

**Discursive Socratic Questioning: Evaluating the Faithfulness of Language Models' Understanding of Discourse Relations**

Yisong Miao[1]

**Prompting Implicit Discourse Relation Annotation**

Frances Yung[1]    Mansoor Ahmad[1]    Merel Scholman[1,2]    Vera Demberg[2]

### SDRT-style dialogue parsing

**Uncovering the Potential of ChatGPT for Discourse Analysis in Dialogue: An Empirical Study**

Yaxin Fan[1]*, Feng Jiang[2,3]†, Peifeng Li[1], Haizhou Li[2]

### Generalization across frameworks

**The DISRPT 2025 Shared Task on Elementary Discourse Unit Segmentation,**

**Probing LLMs for Multilingual Discourse Generalization Through a Unified Label Set**

Florian Eichin*, Yang Janet Liu*, Barbara Plank, and Michael A. Hedderich

5

Theory oriented

Task specific

PDTB-style relation classification

Restricted on sentence-level, relation classification

Yisong Miao[1] **Prompting Implicit Discourse Relation Annotation**

**Are LLMs Good Annotators for Discourse-level Event Relation Extraction?**

Kangda Wei, Aayush Gautam, Ruihong Huang

Uncovering the Potential of ChatGPT for Discourse Analysis in

**Zero-shot Temporal Relation Extraction with ChatGPT**

Chenhan Yuan, Qianqian Xie, Sophia Ananiadou

Generalization across frameworks

**Test of Time: A Benchmark for Evaluating LLMs on Temporal Reasoning**

Discourse Unit

Discourse Generalization Through a Unified Label Set

Florian Eichin*, Yang Janet Liu*, Barbara Plank, and Michael A. Hedderich

6

# Why discourse understanding? – existing studies

**Theory oriented** → Restricted on sentence-level, relation classification

**Task specific** → Focus on certain aspects of discourse understanding

**Discourse benchmarks**

**Disco-Bench: A Discourse-Aware Evaluation Benchmark for Language Modelling**

Longyue Wang, Zefeng Du, Donghuai Liu, Cai Deng, Dian Yu, Haiyun Jiang, Yan Wang, Leyang Cui, Shuming Shi, Zhaopeng Tu[*]
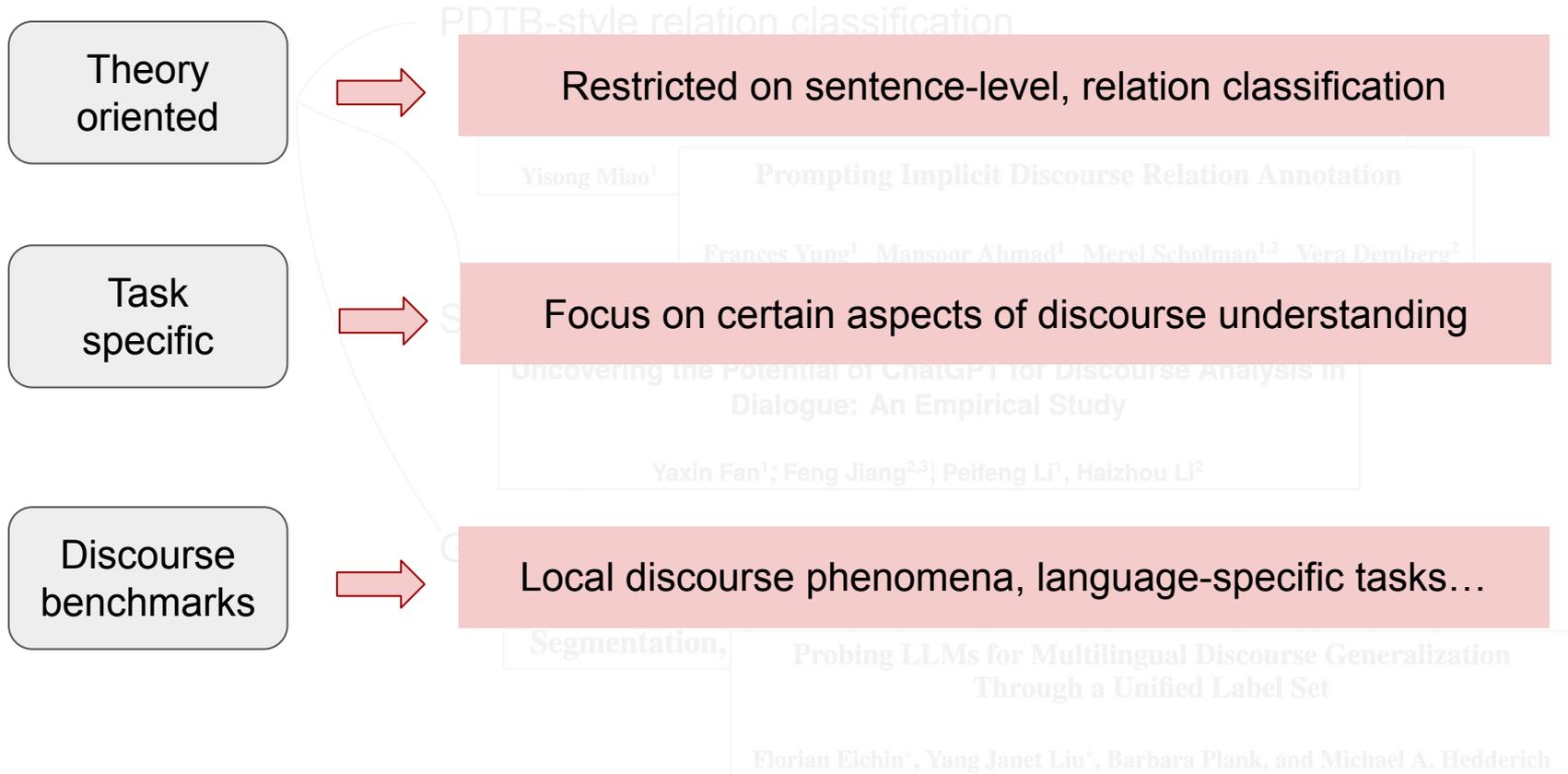
**Evaluation Benchmarks and Learning Criteria for Discourse-Aware Sentence Representations**

Mingda Chen[2*]  Zewei Chu[1*]  Kevin Gimpel[2]

**DiscoTrack: A Multilingual LLM Benchmark for Discourse Tracking**
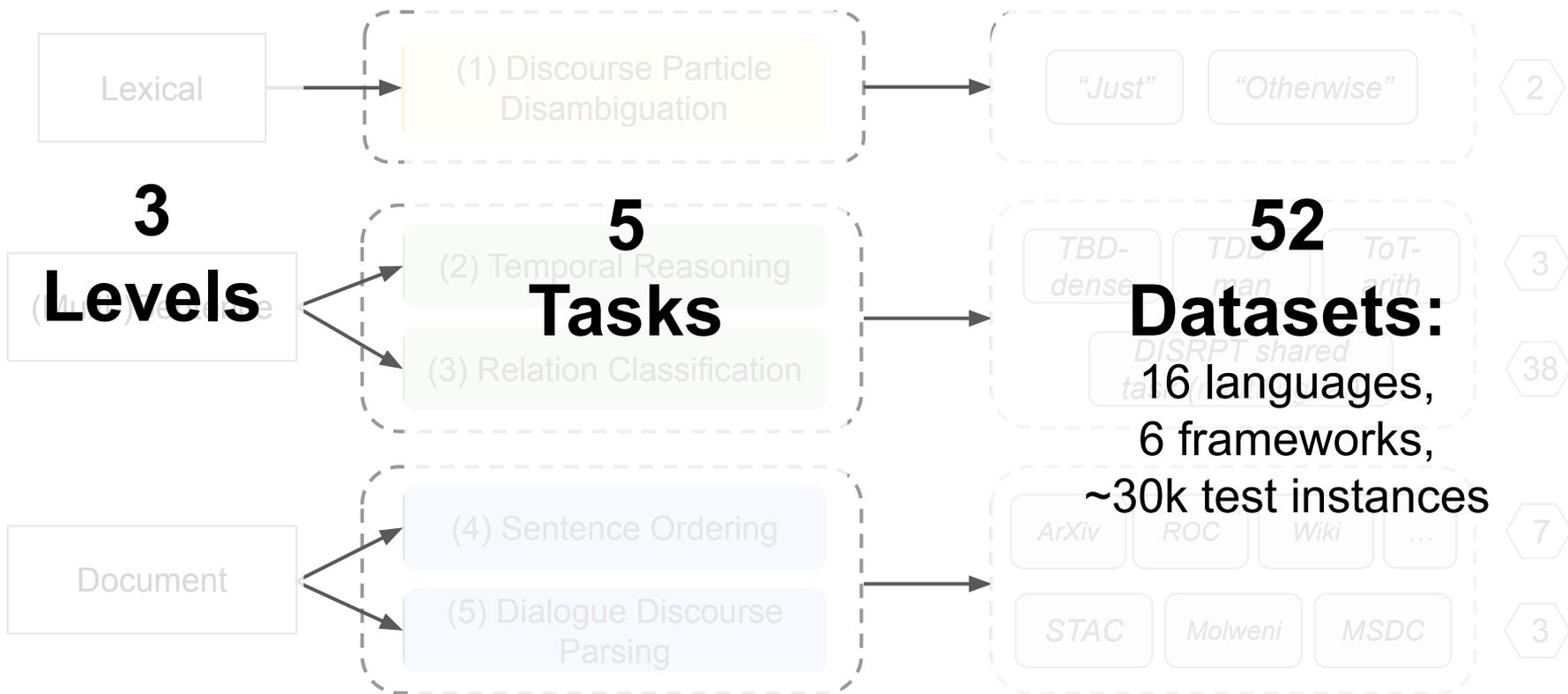
Lanni Bu    Lauren Levine    Amir Zeldes

# Why discourse understanding? – existing studies

**Theory oriented** → Restricted on sentence-level, relation classification

PDTB-style relation classification

Yisong Miao[1]    **Prompting Implicit Discourse Relation Annotation**

Frances Yung[1]  Mansoor Ahmad[1]  Merel Scholman[1,2]  Vera Demberg[2]

**Task specific** → Focus on certain aspects of discourse understanding

Uncovering the Potential of ChatGPT for Discourse Analysis in Dialogue: An Empirical Study

Yaxin Fan[1]; Feng Jiang[2,3]; Peifeng Li[1], Haizhou Li[2]

**Discourse benchmarks** → Local discourse phenomena, language-specific tasks…

Segmentation,    **Probing LLMs for Multilingual Discourse Generalization Through a Unified Label Set**

Florian Eichin*, Yang Janet Liu*, Barbara Plank, and Michael A. Hedderich

# Our proposal: BeDiscovER!

**Lexical** → **(1) Discourse Particle Disambiguation** → *"Just"* *"Otherwise"* ⟨2⟩

**(Multi-)Sentence** → **(2) Temporal Reasoning** **(3) Relation Classification** → *TBD-dense* *TDD-man* *ToT-arith* ⟨3⟩ *DISRPT shared task (multilingual)* ⟨38⟩

**Document** → **(4) Sentence Ordering** **(5) Dialogue Discourse Parsing** → *ArXiv* *ROC* *Wiki* *...* ⟨7⟩ *STAC* *Molweni* *MSDC* ⟨3⟩

# Our proposal: BeDiscovER!

**3 Levels**

**5 Tasks**

**52 Datasets:**
16 languages,
6 frameworks,
~30k test instances

(1) Discourse Particle Disambiguation

"Just"   "Otherwise"   ⟨2⟩

Lexical

(2) Temporal Reasoning

(3) Relation Classification

TBD-dense   TDD-man   ToT-arith   ⟨3⟩

DISRPT shared task   ⟨38⟩

(4) Sentence Ordering

(5) Dialogue Discourse Parsing

ArXiv   ROC   Wiki   ...   ⟨7⟩

STAC   Molweni   MSDC   ⟨3⟩

Document

# BeDiscovER: Evaluation Setting

Open-ended Question-Answer Formatting

- **Unified evaluation pipeline**

- Classification tasks (1 2 3): fixed label space

- Parsing task (5): incremental generation task

Reasoning-oriented LLMs

GPT-5    Qwen3   DeepSeek-r1

Non reasoning-oriented LLMs

Llama-4

Qwen2.5
GPT4o

# BeDiscovER: Evaluation Setting

**System prompt:**
… Choose one of the following six labels: **[Exclusionary, Unelaboratory, Unexplanatory, Emphatic, Temporal, Adjective].**

**User prompt:**
My brother *just* flew in to town.

**Question:** What is the function of the discourse marker "just" in the sentence above?

Reasoning-oriented LLMs

Temporal

GPT-5    Qwen3  DeepSeek-r1

Qwen2.5
GPT4o

Llama-4

Non reasoning-oriented LLMs

# Performance: model scaling



Accuracy

How well do modern LLMs understand discourse
… in BeDiscovER tasks ?

Sentence
Ordering
(ArXiv abs)

70

60

50

40

30

20

4B

8B

Qwen3
14B

Qwen3
32B

Deep
Seek-r1

GPT-5 mini
(high)

LLMs

*Bigger the model, better the performance – expected!*

Sentence Ordering (ArXiv abs)

*Reasoning-oriented LLMs outperform non-reasoning optimized LLMs – expected!*

Sentence Ordering (ArXiv abs)

Accuracy

70
60
50
40
30
20

4B
8B
Qwen3 14B
Qwen3 32B
Deep Seek-r1
GPT-5 mini (high)

Reasoning (high)
Reasoning (low)

**Higher thinking effort does not yield better outcome**

*Models become verbose, yet, flat performance*

*… except for **arithmetic** temporal: long reasoning lead to big gains!*

Accuracy

80
70
60
50
40

Sentence

Dialogue

Temporal

Relation Classifi-cation

4B

8B

Qwen3 14B

Qwen3 32B

Deep Seek-r1

GPT-5 mini (high)

Supervised NAON-BART

LLMs

*Reasoning-oriented LLMs show markedly lower performance (~10–30%) compared to supervised models.*

# Analysis: Fine-grained sense disambiguation



Confusion Matrix for gpt-5-mini-2025-08-07 on TR task (TDD-Manual)

Temporal Reasoning

Relation Classifi-cation

Identifies *before/after* relations

but fails to capture *overlap* or **containment**.

# Analysis: Multilingual performance



Relation Classification

Supervised model show clear perform disparities (20%)

LLMs: nearly flat, lack of robust relation representation

low-resource languages

high-resource languages

# Analysis: Inter-task correlation

Performance similarity (Pearson r)

- For each task: aggregated results of all datasets
- Metric: accuracy
- Pairwise correlation: strong positive ($r \in [0.73, 0.98]$, $p < 0.05$)

- Lexicon & semantics
  - Rhetorical
  - Temporal
  - Logic
  - Arithmetic
  - Commonsense

*"Just"*

Relation classification

Overlap calculation:

$$Jaccard\ (A, B) = \frac{|A \cup B|}{|A \cap B|}$$

Knowledge overlap (Jaccard index)

# Analysis: Inter-task correlation



- ❖ DDP – dialogue discourse parsing
- ❖ SO – sentence ordering
- ❖ TR – temporal
- ❖ DR – relation
- ❖ DM – markers

# Summary: benchmark and evaluation baseline

Reasoning-oriented LLMs capture some discourse-level knowledge, especially **good in arithmetic** temporal reasoning.

**Struggle with subtle semantic phenomena** (e.g. rhetorical relations, for all languages)

Longer reasoning traces do not necessarily yield better outcomes → how to **improve the quality** of reasoning LLMs?

Data contamination? → **dynamic** benchmark, welcome contribution!