# Indirect Supervision from Natural Language Inference
## Indirectly Supervised Natural Language Processing (Part I)

Wenpeng Yin

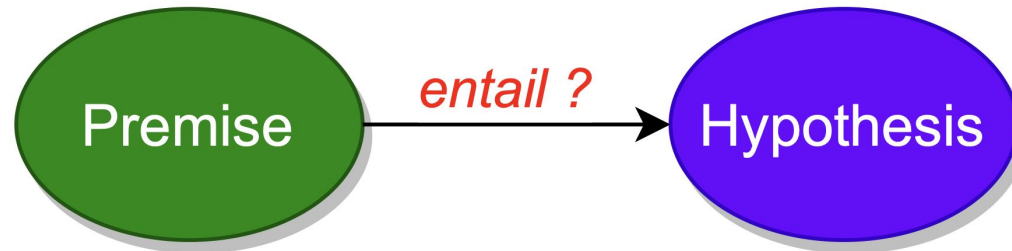Department of Computer Science and Engineering

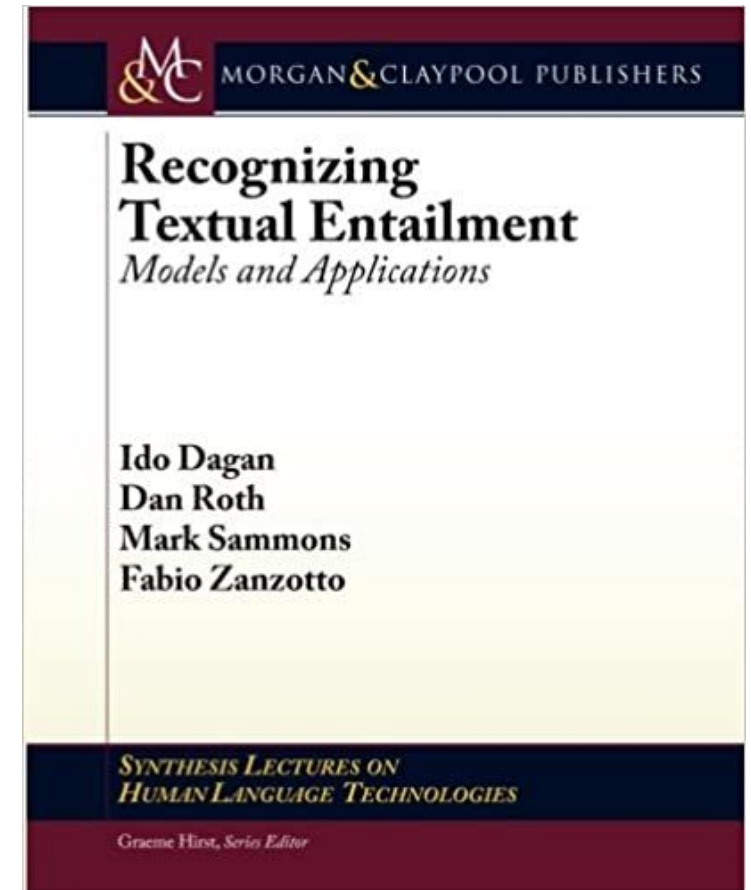Penn State University

July 2023

ACL Tutorials

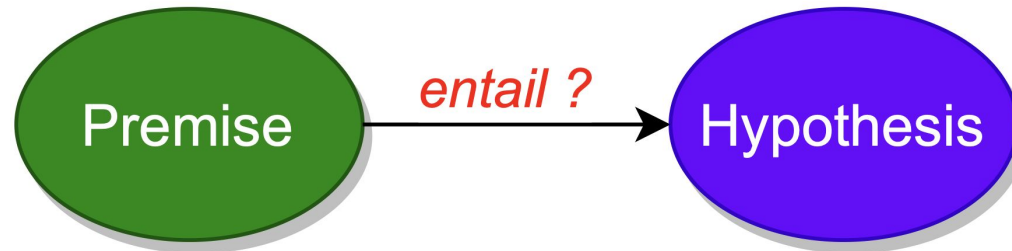Indirectly Supervised Natural Language Processing

# Natural Language Inference for NLP

❏ Textual Entailment (Dagan et al., 2006)



❏ Textual entailment, a unified inference framework for NLP
  ❏ zero-shot text classification (Yin et al., 2019)
  ❏ summarization (Falke et al., 2019)
  ❏ QA & Coreference (Yin et al., 2020)
  ❏ relation extraction (Xia et al., 2021)
  ❏ entity typing (Li et al., 2022)
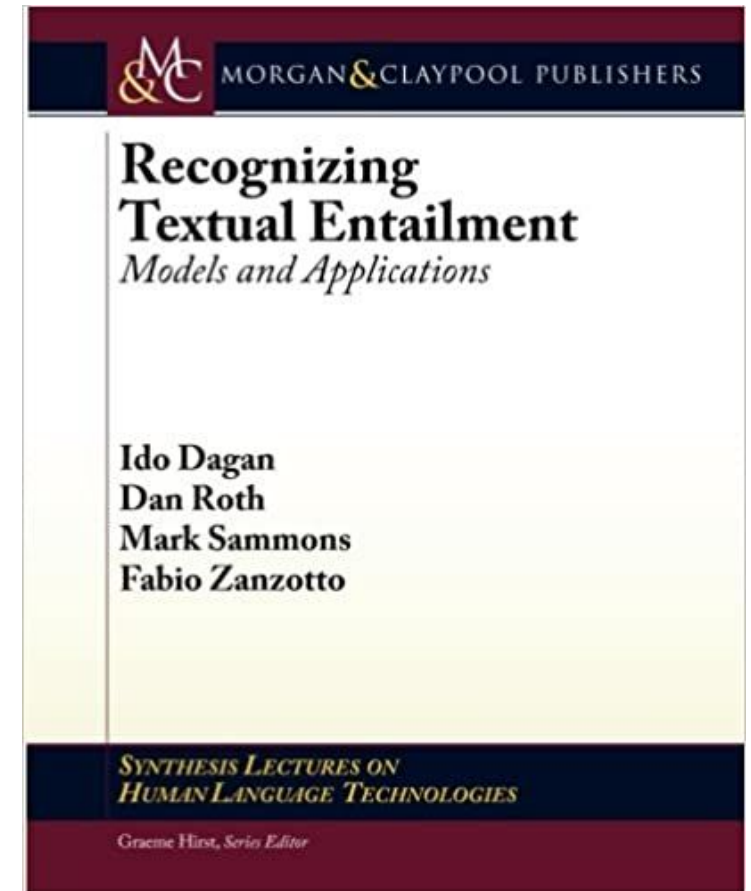  ❏ …

# Natural Language Inference for NLP
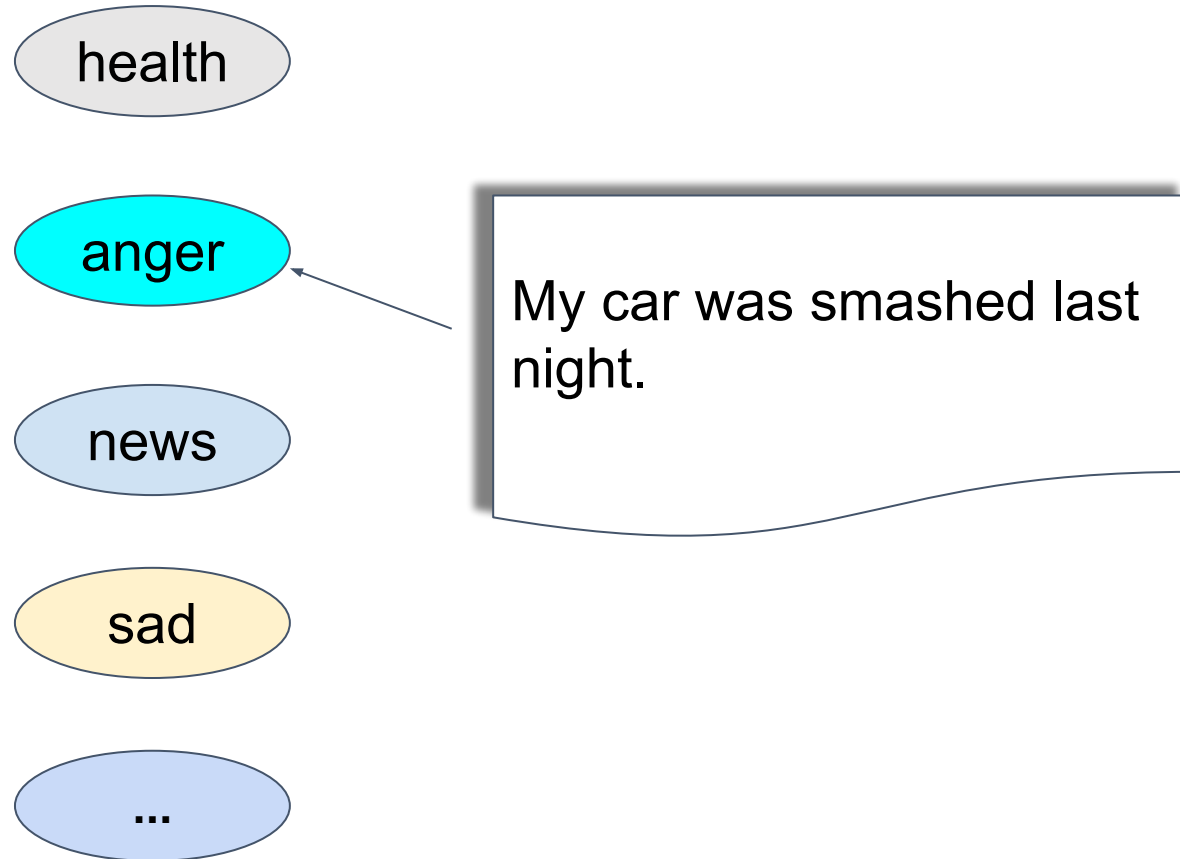
❏ Textual Entailment (Dagan et al., 2006)



❏ Textual entailment, a unified inference framework for NLP
   ❏ zero-shot text classification (Yin et al., 2019)
   ❏ summarization (Falke et al., 2019)
   ❏ QA & Coreference (Yin et al., 2020)
   ❏ relation extraction (Xia et al., 2021)
   ❏ entity typing (Li et al., 2022)
   ❏ …

❏ "Textual Entailment" was referred to as "Natural Language Inference (NLI)" (Bowman et al., 2015)

## Example task: Text Classification

health

anger

news

sad

...

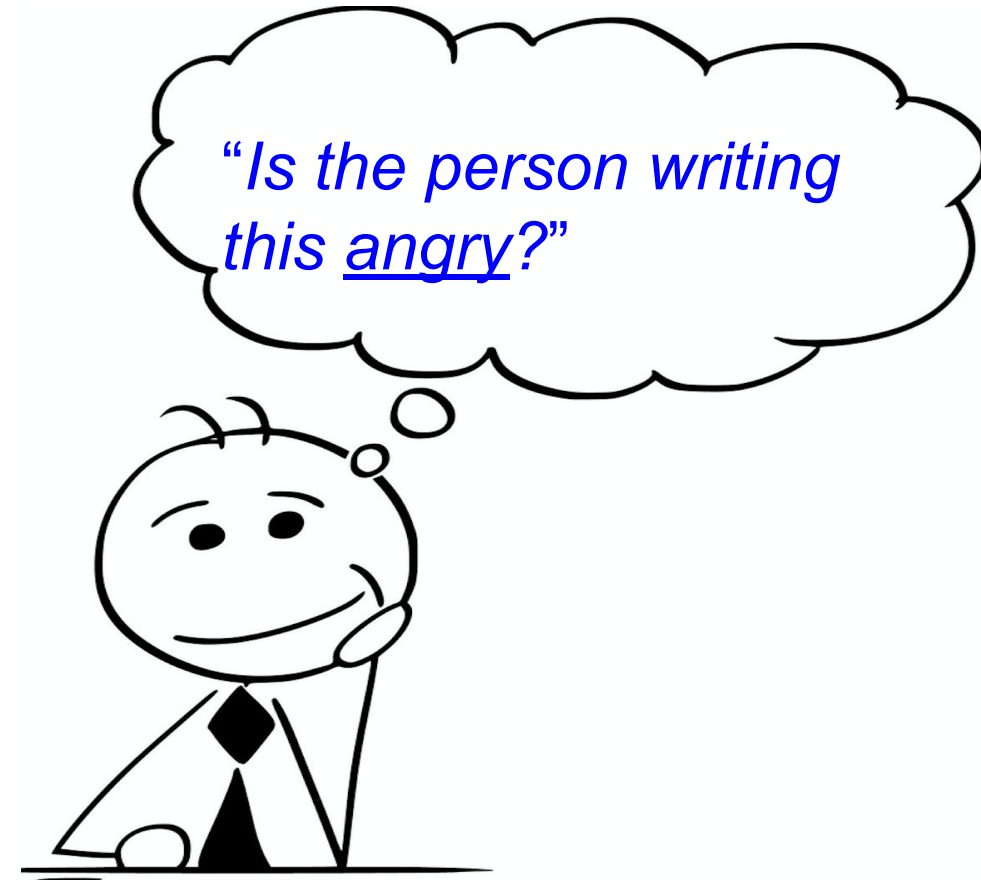My car was smashed last night.

**Example task: Text Classification**

health

anger ← My car was smashed last night.

news

sad

...

"*Is the person writing this angry?*"

Humans implicitly build a hypothesis then infer its truth value.

**Example task: Text Classification**

My car was smashed last night.

The person is angry.

The person is not angry.

The classification problem is **inference (text → hypothesis)**

**Example task: Question Answering**

The road to Grandpa's house was long and winding. [...]. Finally, Jimmy arrived at Grandpa's house and knocked. Grandpa answered the door with a smile and welcomed Jimmy inside. They sat leisurely by the fire and talked about the insects. They watched the lightning bugs light up as night came.

Where do Jimmy and his Grandpa sit?
A) On insects
B) Outside
C) By the fire
D) On the path

## Example task: Question Answering

The road to Grandpa's house was long and winding. [...]. Finally, Jimmy arrived at Grandpa's house and knocked. Grandpa answered the door with a smile and welcomed Jimmy inside. They sat leisurely by the fire and talked about the insects. They watched the lightning bugs light up as night came.

Where do Jimmy and his Grandpa sit?

A) On insects

B) Outside

C) By the fire

D) On the path

Because "by the fire" and "sat" co-occur in the same sentence? NO!

## Example task: Question Answering

The road to Grandpa's house was long and winding. [...]. Finally, Jimmy arrived at Grandpa's house and knocked. <u>Grandpa answered the door with a smile and welcomed Jimmy inside. They sat leisurely by the fire and talked about the insects</u>. They watched the lightning bugs light up as night came.

Where do Jimmy and his Grandpa sit?
A) On insects
B) Outside
C) <u>By the fire</u>
D) On the path

<span style="color:blue">can infer the meaning</span>

<u>"Jimmy and his Grandpa sit by the fire"</u>

**Example task: Question Answering (QA)**

The road to Grandpa's house was long and winding. [...]. Finally, Jimmy arrived at Grandpa's house and knocked. Grandpa answered the door with a smile and welcomed Jimmy inside. They sat leisurely by the fire and talked about the insects. They watched the lightning bugs light up as night came.

Where do Jimmy and his Grandpa sit?

A) On insects
B) Outside
C) By the fire
D) On the path

QA is **inference (text → hypothesis)**

**Example task: Coreference Resolution**

?

The **trophy** would not fit in the brown **suitcase** because <u>it</u> was too big.

?

**Example task: Coreference Resolution**

?

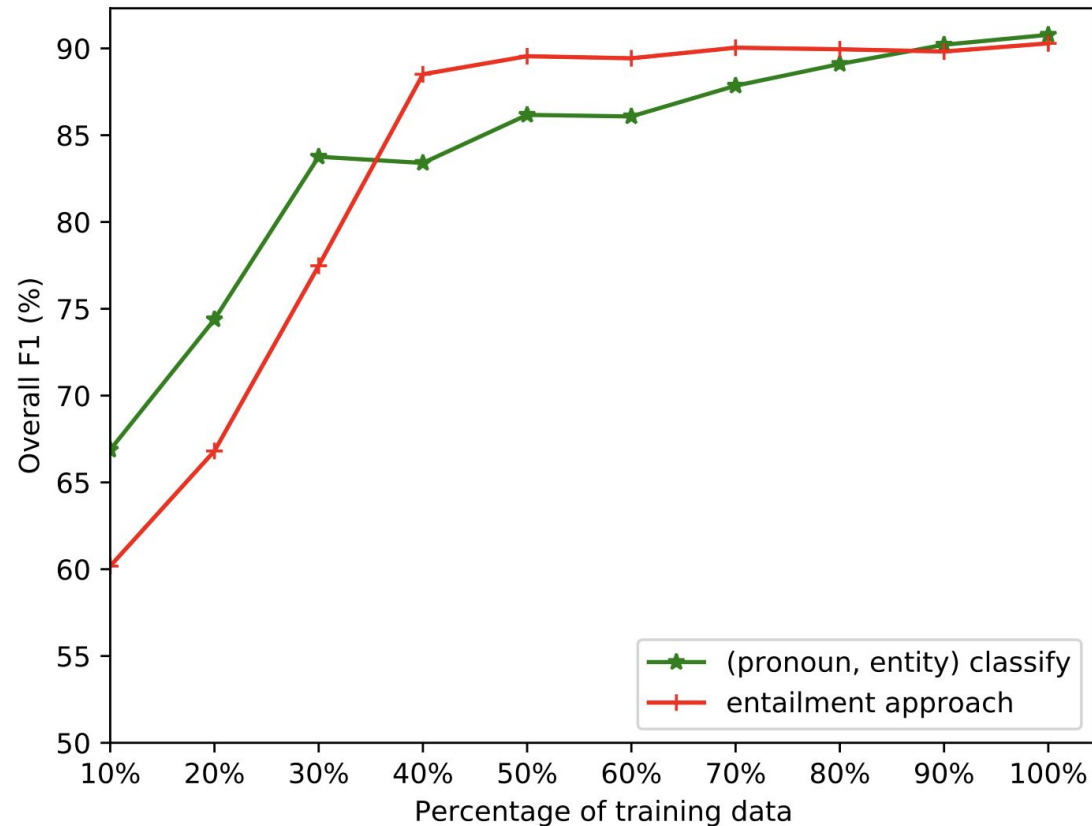The **trophy** would not fit in the brown **suitcase** because <u>it</u> was too big.

?

➔ The trophy would not fit in the brown suitcase because <u>**trophy**</u> was too big. (True)

➔ The trophy would not fit in the brown suitcase because <u>**suitcase**</u> was too big. (False)

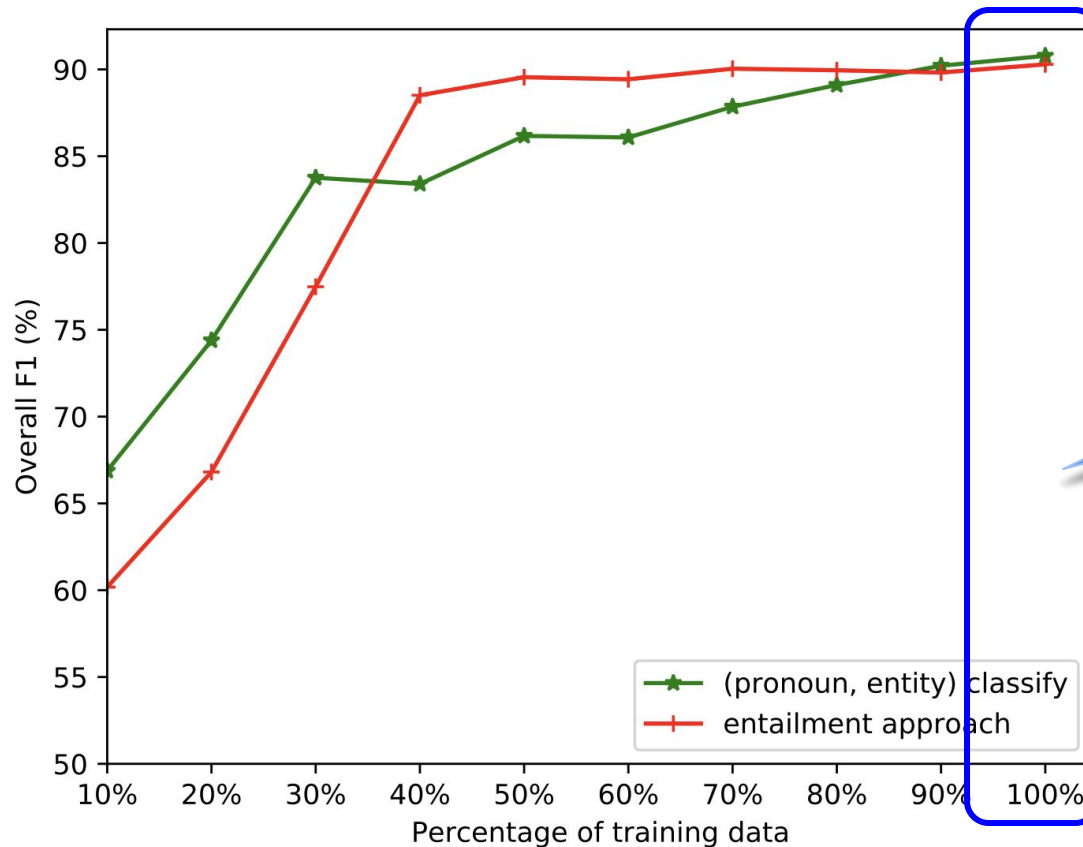Coreference resolution is **inference (text → hypothesis)**

# Why and when to convert NLP to NLI?



**Task**: conference resolution
**Dataset**: GAP (Webster et al., 2018)
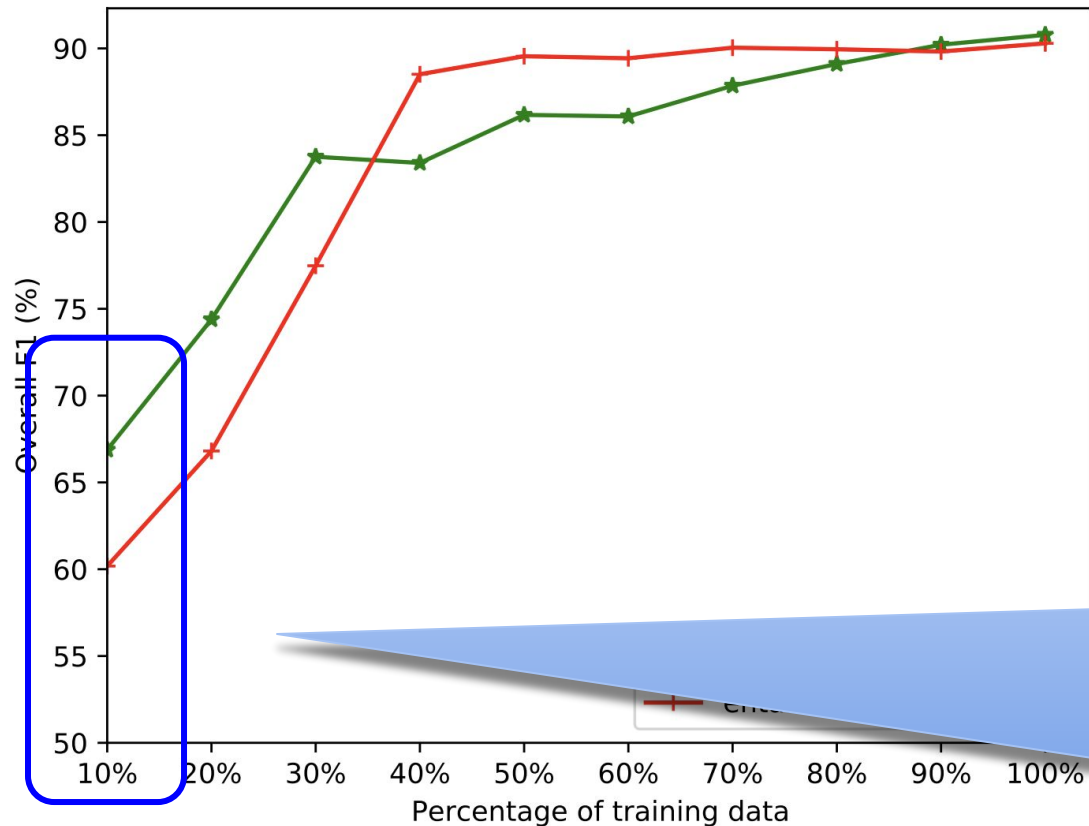
# Why and when to convert NLP to NLI?



When an NLP task has rich annotations:
classical classifiers ≈ NLI

**Task**: conference resolution
**Dataset**: GAP (Webster et al., 2018)

# Why and when to convert NLP to NLI?



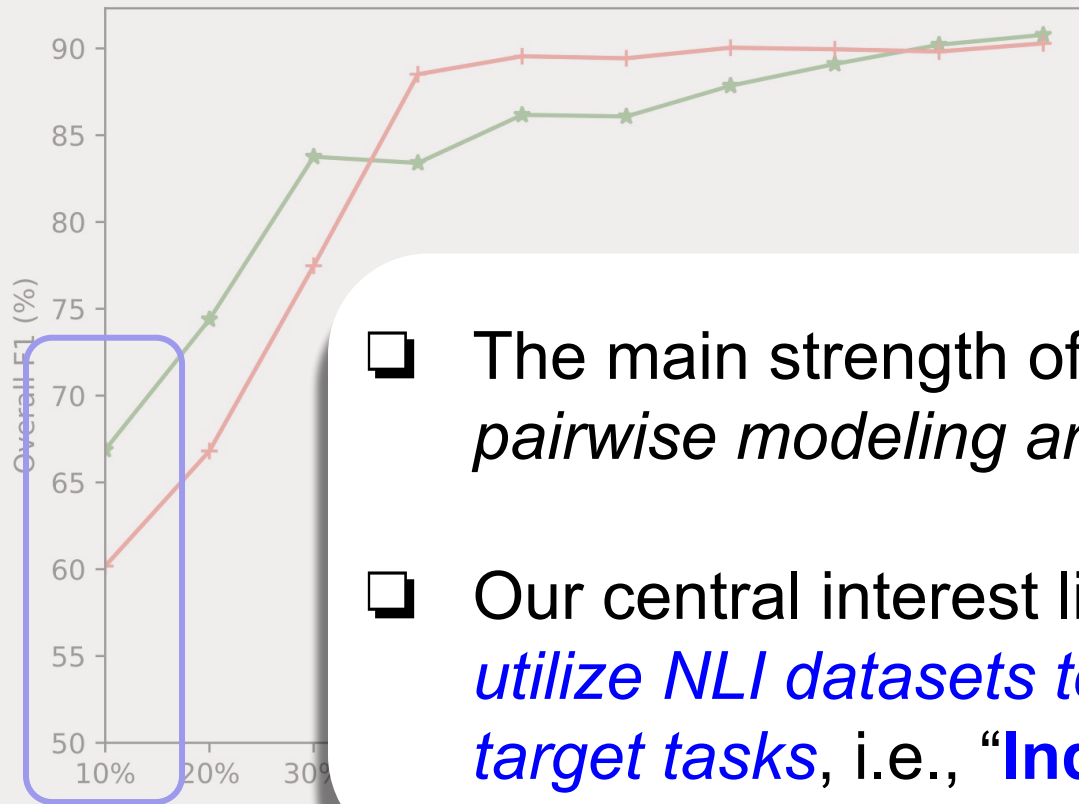**Task**: conference resolution
**Dataset**: GAP (Webster et al., 2018)

In reality, what we truly care about is "limited annotation"

❑ Both standard classifiers and NLI exhibit poor performance.

❑ NLI performs even worse because NLI is generally a more challenging task when the availability of labeled examples is severely limited.

*Universal Natural Language Processing with Limited Annotations: Try Few-shot Textual Entailment as a Start* (Yin et al., EMNLP'2020)

# Why and when to convert NLP to NLI?



❏ The main strength of NLI is **NOT** its *pairwise modeling architecture*.

❏ Our central interest lies in its potential to *utilize NLI datasets to supervise different target tasks*, i.e., "**Indirect Supervision**"
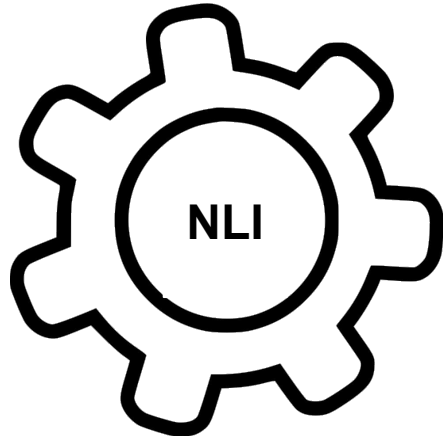
**Task**: conference resolution
**Dataset**: GAP (Webster et al., 2018)

*Universal Natural Language Processing with Limited Annotations: Try Few-shot Textual Entailment as a Start* (Yin et al., EMNLP'2020)
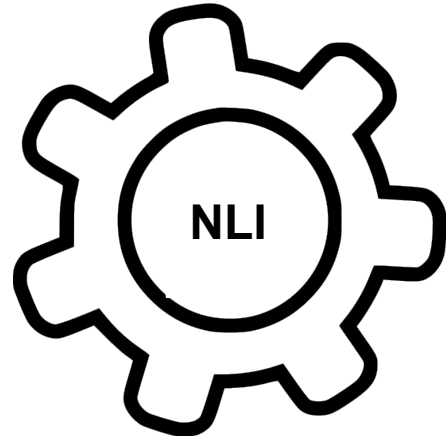
❏ Implementation
& Applications

❏ Benefits

❏ Challenges &
Solutions

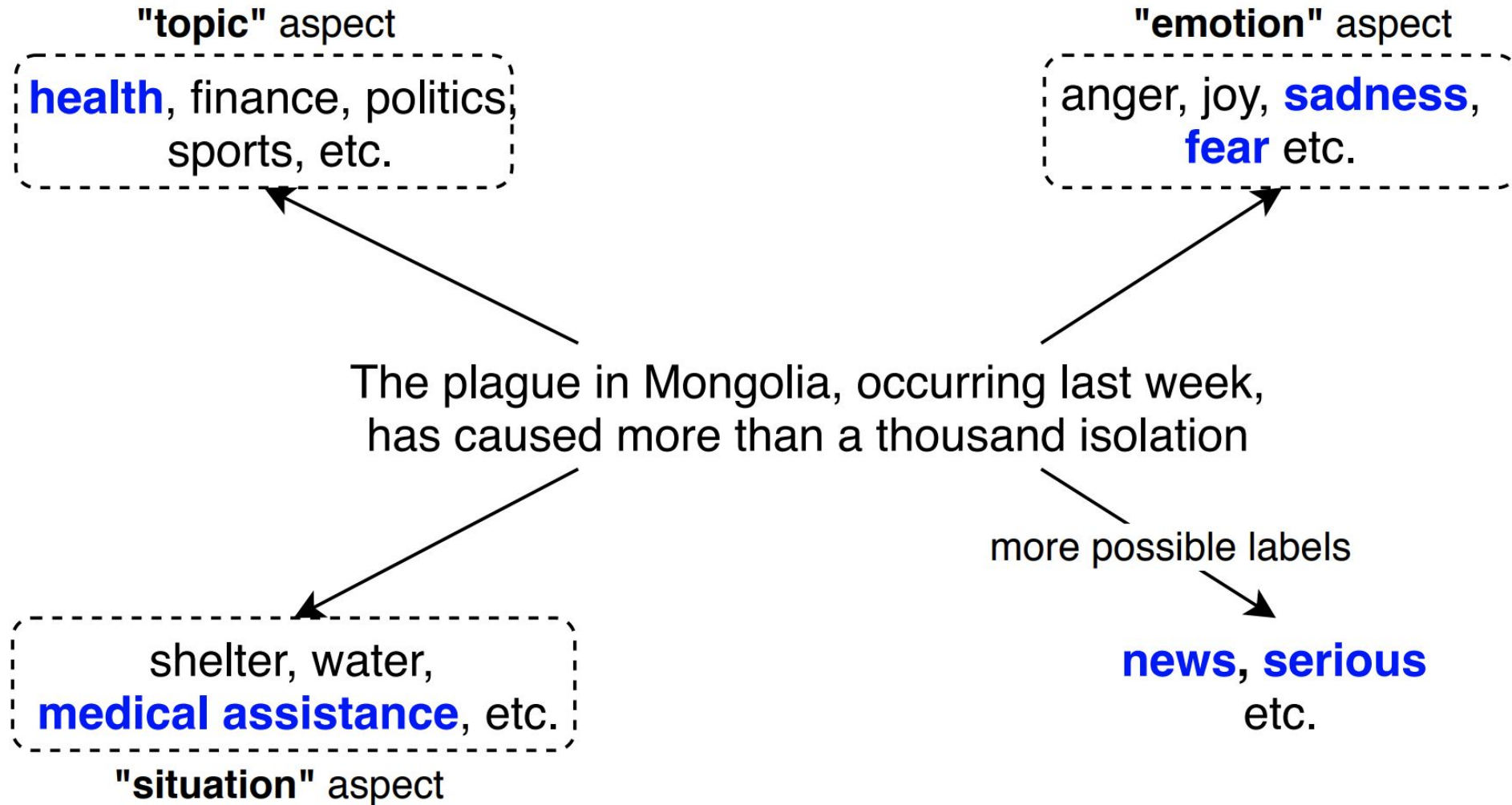# NLI-based indirect supervision: outline

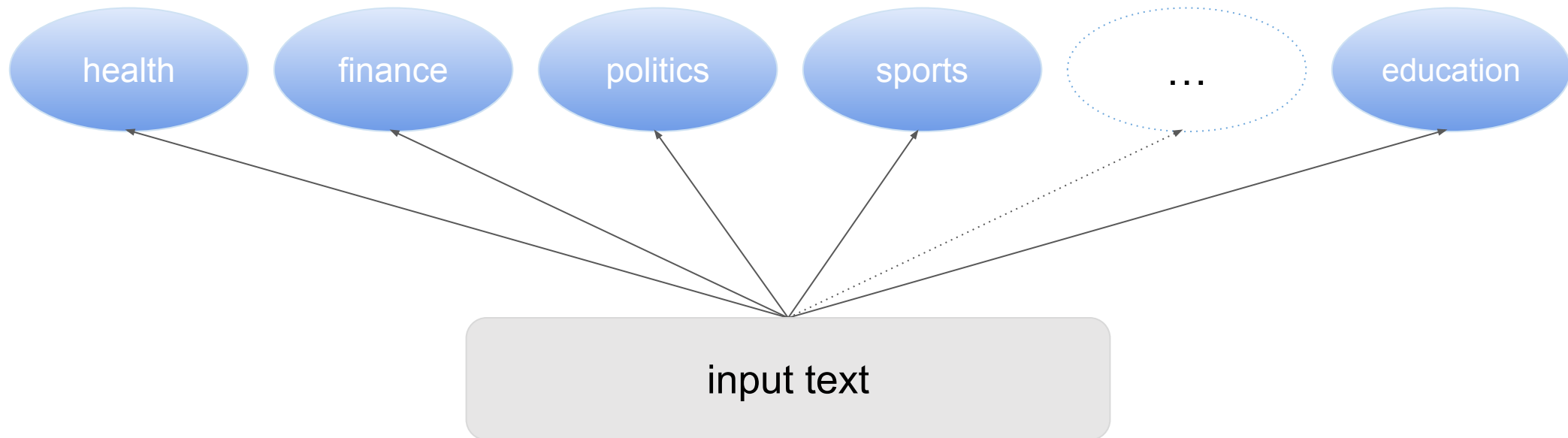❏ Implementation & Applications

❏ Benefits

❏ Challenges & Solutions

# Text classification task

**"topic" aspect**

**health**, finance, politics sports, etc.

**"emotion" aspect**

anger, joy, **sadness**, **fear** etc.

The plague in Mongolia, occurring last week, has caused more than a thousand isolation

shelter, water, **medical assistance**, etc.

**"situation" aspect**

more possible labels

**news, serious** etc.

*Benchmarking Zero-shot Text Classification: Datasets, Evaluation and Entailment Approach* (Yin et al., 2019)

# How do conventional classifiers work

real labels are converted into indices; label semantics are missing

# Issues of conventional classifiers

labeled data

labeled data

...

labeled data

0

1

...

N-1

input text

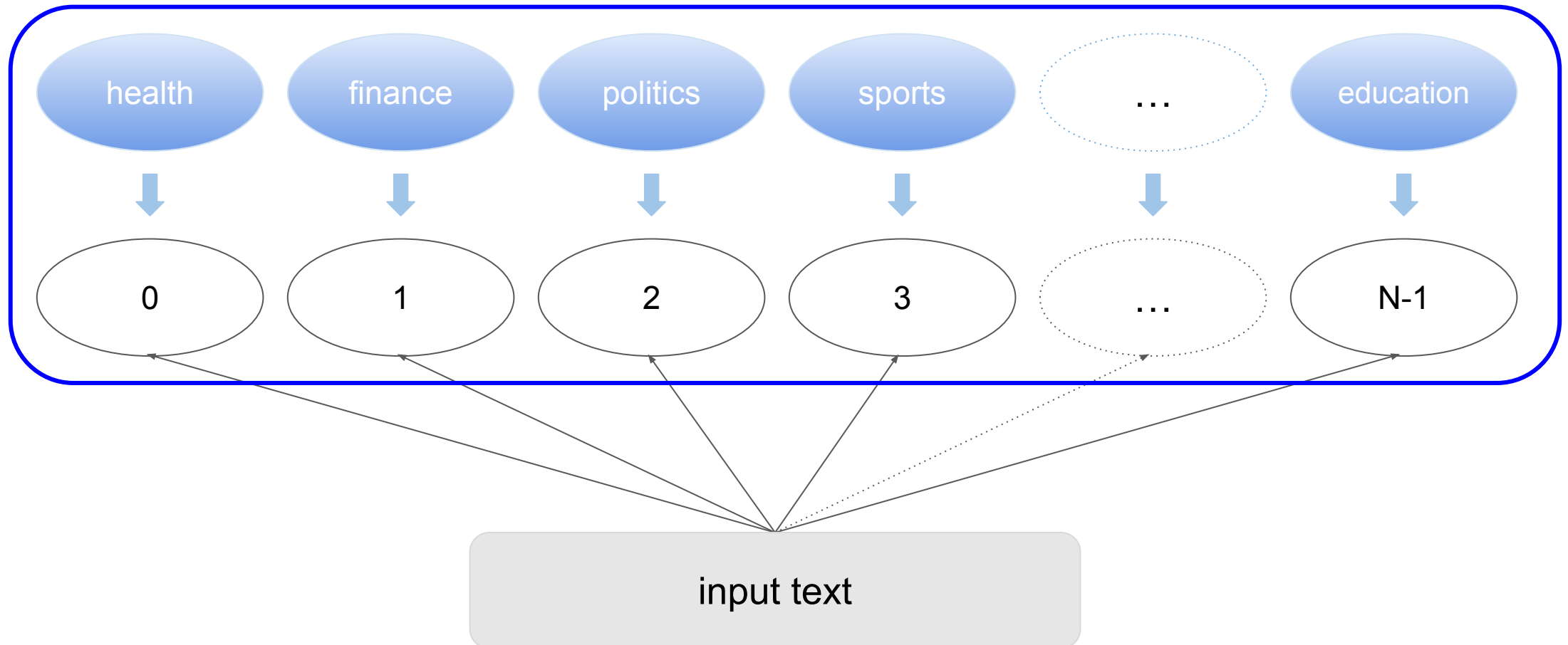1) needs a large number of labeled examples to learn existing types

# Issues of conventional classifiers

1) needs a large number of labeled examples to learn existing types

2) **cannot generalize to new types**

0   1   …   N-1   action

input text

# Issues of conventional classifiers



The conventional approach is inapplicable to low-annotation text classification—we need indirect supervision

**Zero-shot text classification**

**Natural language inference**

0
(health)

1
(anger)

2
(accident)

3
(crime)

...

My car was smashed last night.

(Yin et al., 2019; Xia et al., 2021, Xu et al., 2022, etc)

**Zero-shot text classification**

**Natural language inference**

0
(health)

1
(anger)

2
(accident)

3
(crime)

...

Premise

My car was smashed last night.

Hypothesis

This text expresses anger.

(Yin et al., 2019; Xia et al., 2021, Xu et al., 2022, etc)

28

**Zero-shot text classification**

**Natural language inference**

0
(health)

1
(anger)

2
(accident)

3
(crime)

...

Premise

My car was smashed last night.

Hypothesis

This text expresses anger.

This text is about accident.

(Yin et al., 2019; Xia et al., 2021, Xu et al., 2022, etc)

**Zero-shot text classification**

**Natural language inference**

0
(health)

1
(anger)

2
(accident)

3
(crime)

...

Premise

My car was smashed last night.

Hypothesis

This text expresses anger.

This text is about accident.

We can keep the original label strings, or use other related words
— **label verbalizer** (Schick and Schütze, 2020)

**NLI** datasets:
- ❏ MNLI ([Williams et al., 2018](#))
- ❏ ANLI ([Nie et al., 2020](#))
- ❏ SNLI ([Bowman et al., 2015](#))
- ❏ DocNLI ([Yin et al., 2021](#))
- ❏ SciTail ([Khot et al., 2018](#))
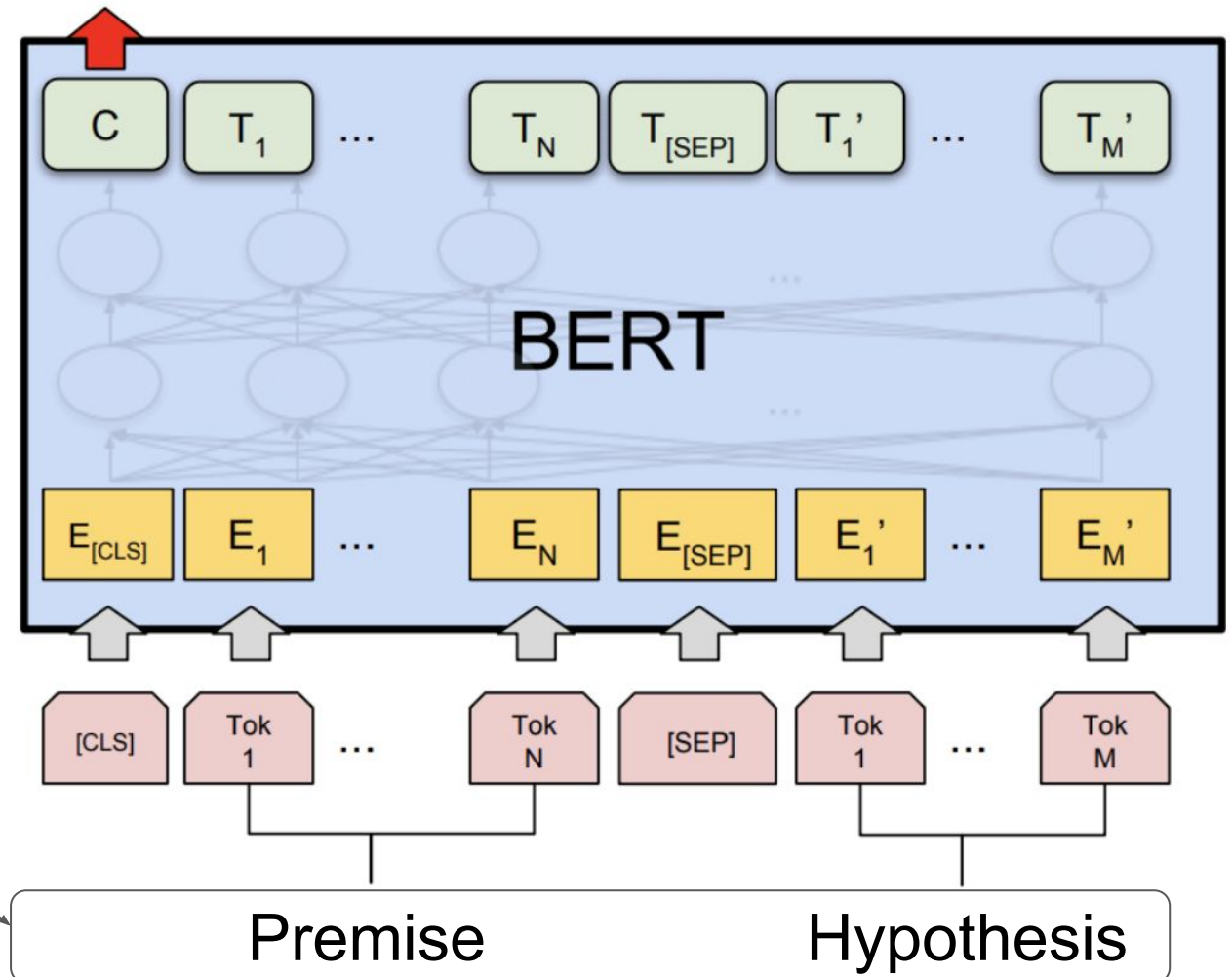- ❏ RTE ([Dagan et al. 2005](#))
- ❏ ...

**PLMs:**
- ❏ BERT, RoBERTa, T5, etc.

**Training pairs**
- ❏ are from the NLI datasets;
- ❏ are not specific to any target labels

(Yin et al., 2019; Xia et al., 2021, Xu et al., 2022, etc)

**PLMs**:
- ❏ pretrained NLI model

**Training pairs**
- ❏ **are from the target task** (if available);
  - topic classification (Yin et al., 2019)
  - multi-choice QA and corefernece (Yin et al., 2020)
  - intent identification (Xia et al., 2021)
  - stance detection (Xu et al., 2022)
  - ultra-fine entity typing (Li et al., 2022)
  - event argument extraction (Sainz et al., 2022)
  - (biomedical) relation extraction (Sainz et al., 2021; Xu et al., 2023)
  - …
- ❏ **are in the NLI format**

True/False

# Performance



Ultra-fine entity typing ([Li et al., TACL 2022](#))
**NLI can handle zero-shot & few-shot**

# Performance



Legend: MLMET, LITE

F-1 Score vs Label Frequency:
- 0 shot: MLMET 10.8, LITE 32.9
- 1~5 shot: MLMET 23.5, LITE 33.6
- 6~10 shot: MLMET 27.3, LITE 38.9

Ultra-fine entity typing (Li et al., TACL 2022)
**NLI can handle zero-shot & few-shot**



Legend: our pretrained entail, dataless learning (Chang et al., 2008)

Zero-shot text classification (Yin et al., EMNLP 2019). **One NLI system can handle various zero-shot tasks**

# Indirect supervision from NLI: benefits

❏ **Reduce task-specific annotation** requirements and address few-shot and zero-shot scenarios in a unified approach.



(Li et al., TACL 2022)

# Indirect supervision from NLI: benefits

❏ Reduce task-specific annotation requirements and address few-shot and zero-shot scenarios in a unified approach.

❏ **Facilitate cross-task transferability**, encompassing not only NLI to target tasks but also task A to task B.

| | original task | domain | premise length | hypothesis length |
|---|---|---|---|---|
| ANLI | NLI | various (wiki, news, etc.) | multi-sentence (20~94 words) | single sentence (4~18 words) |
| SQuAD | QA | wiki | paragraph (27~237 words) | single sentence (6~22 words) |
| DUC (2001) | summarization | news | doc. (124~879 words) | multi-sent (80~100 words) |
| CNN/Daily Mail | summarization | news | doc. (247~652 words) | 3~4 sent. (40~50 words) |
| Curation | summarization | news | doc. (229~842 words) | multi-sent (64~279 words) |

| | | FEVER binary | MCTest v160 | MCTest v500 |
|---|---|---|---|---|
| | random | 50.00 | 25.00 | 25.00 |
| pretrain | MNLI | 86.64 | 75.41 | 70.66 |
| | ANLI | 87.51 | 82.50 | 78.66 |
| | DocNLI | 88.84 | 90.00 | 85.83 |
| | +finetune | **89.44** | **90.83** | **90.66** |
| Prior state-of-the-art | | – | 80.00 | 75.50 |

DocNLI (Yin et al., 2021) converted QA, summarization tasks as NLI-style source tasks
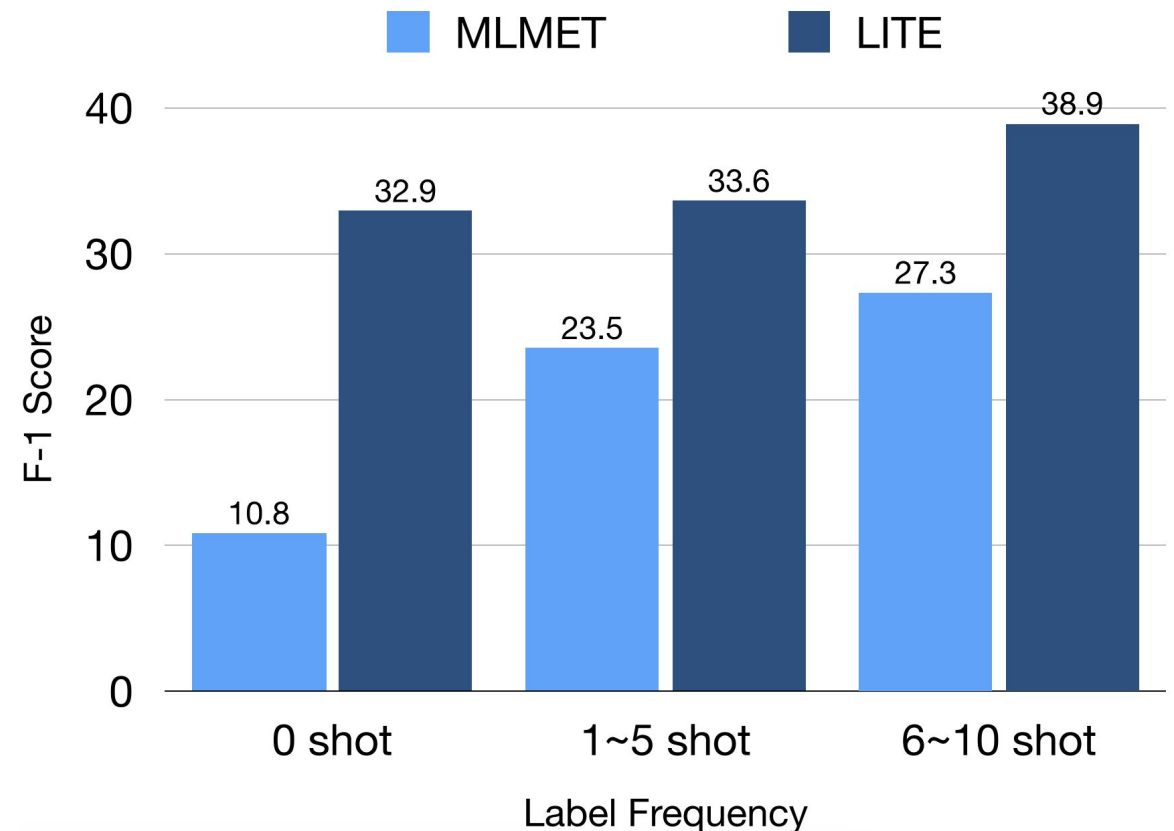
DocNLI generalizes to distinct target tasks

# Indirect supervision from NLI: benefits

❏ Reduce task-specific annotation requirements and address few-shot and zero-shot scenarios in a unified approach.

❏ Facilitate cross-task transferability, encompassing not only NLI to target tasks but also task A to task B.

❏ Enhance the feasibility and potential of employing smaller PLMs.

*"OpenStance: Real-world Zero-shot Stance Detection"* (Xu et al., CoNLL'22)

❏ Zero-shot

❏ RoBERTa (355M) + weak&NLI > GPT-3 (175B)

❏ Reduce task-specific annotation requirements and address few-shot and zero-shot scenarios in a unified approach.

❏ Facilitate cross-task transferability, encompassing not only NLI to target tasks but also task A to task B.

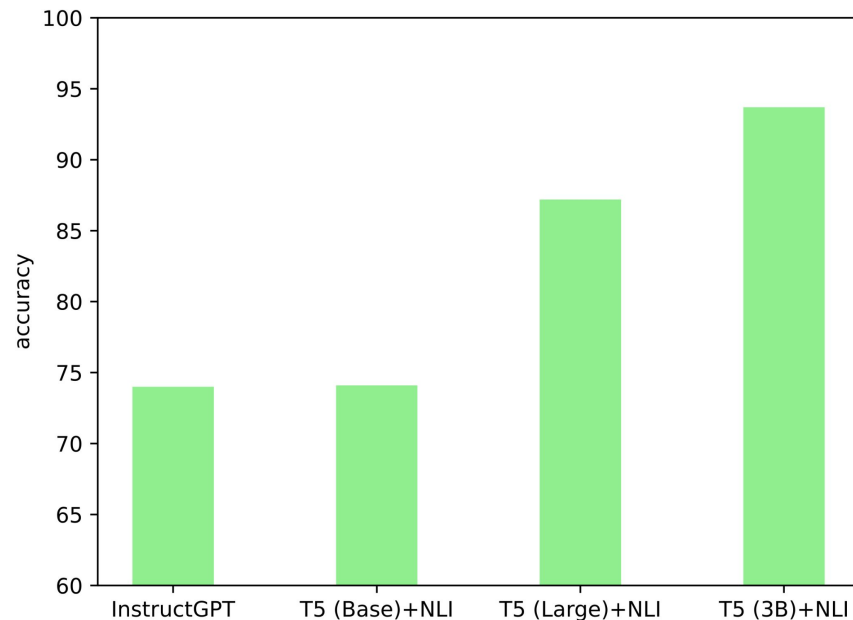❏ Enhance the feasibility and potential of employing smaller PLMs.



*"ZARA: Improving Few-Shot Self-Rationalization for Small Language Models"* (Chen et al., 2023)

❏ Few-shot

❏ NLI helps automatic data generation

❏ T5 (base)+NLI (2M) ≈ InstructGPT (175B)

❏ T5 (large)+NLI (7M) > InstructGPT (175B)

❏ T5 (3B)+NLI (2.7B) > InstructGPT (175B)

❏ Implementation & Applications



❏ Benefits



❏ Challenges & Solutions

# Challenge #1: domain discrepancy



Domains of the target problems (T) often differ from that of NLI datasets (S)

**Traditional solution**: pretrain on S +finetune on T (i.e., STILTS (Phang et al., 2018))

Domains of the target problems (T) often differ from that of NLI datasets (S)

**Traditional solution**: pretrain on S +finetune on T (i.e., STILTS (Phang et al., 2018))
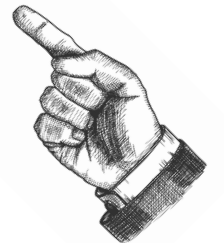
$$\text{loss } l = l_S + l_T$$

**Solution I**:

- ❏ T is few-shot → imitate few-shot learning on S (meta-learning)
- ❏ Novel strategy: predictions on T (or S) depend on the signals of both T and S
- ❏ A solution from the algorithm perspective

*Universal Natural Language Processing with Limited Annotations: Try Few-shot Textual Entailment as a Start* (Yin et al., EMNLP 2020)

loss $l_S$     loss $l_T$

mean     mean

batch of queries in S     batch of queries in T

$p_S^e$ | $p_S^n$ | $p_S^c$     $p_T^e$ | $p_T^n$ | $p_S^c$

class rep. in S     class rep. in T
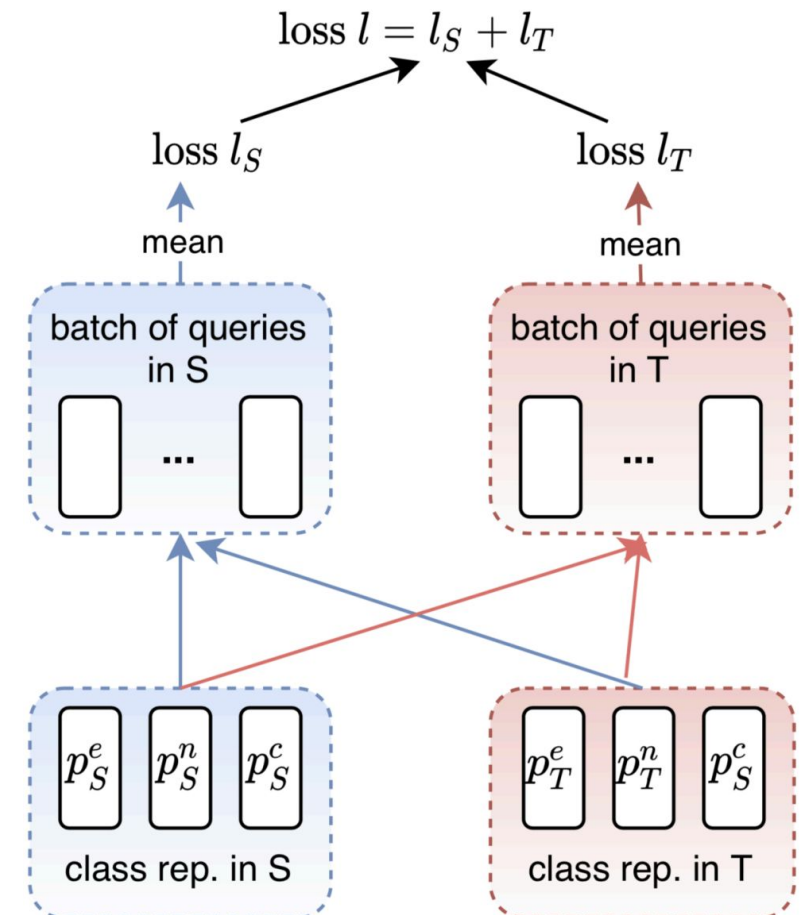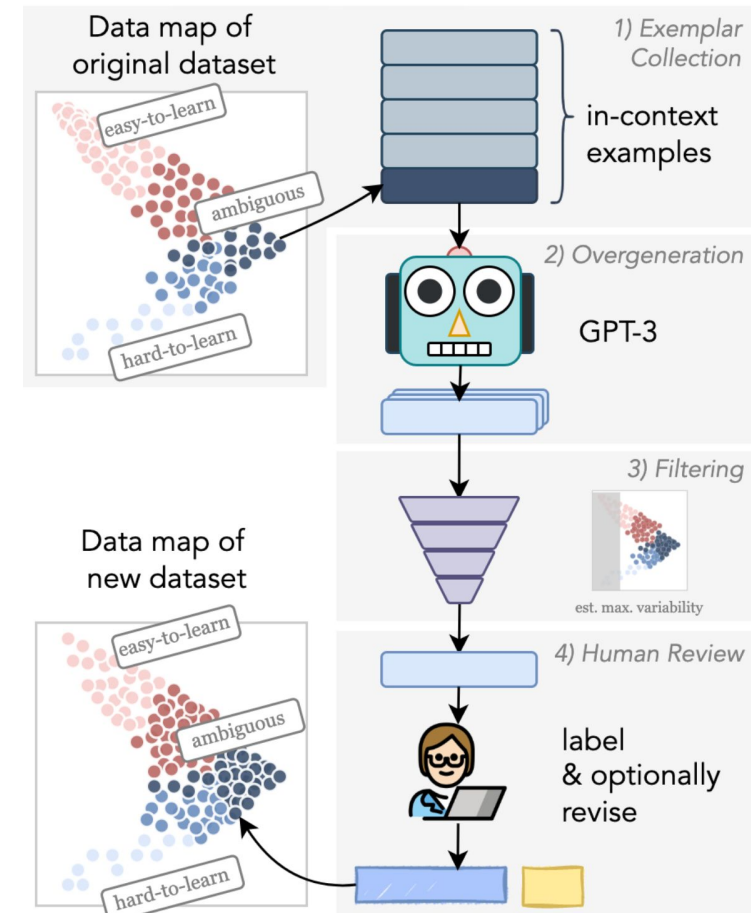
# Challenge #1: domain discrepancy

Domains of the target problems (T) often differ from that of NLI datasets (S)
    **Traditional solution**: pretrain on S +finetune on T (i.e., STILTS (Phang et al., 2018))

**Solution II**:

- ❏ Human annotations: correctness 👍, diversity 👎
- ❏ PLMs: creative writing
- ❏ Use GPT-3 to generate new examples for reasoning patterns that are challenging
- ❏ A solution from the data perspective

*WANLI: Worker and AI Collaboration for Natural Language Inference Dataset Creation*
(Liu et al., Findings of EMNLP 2022)



Data map of original dataset
- easy-to-learn
- ambiguous
- hard-to-learn

1) Exemplar Collection — in-context examples

2) Overgeneration — GPT-3

3) Filtering — est. max. variability

4) Human Review — label & optionally revise

Data map of new dataset
- easy-to-learn
- ambiguous
- hard-to-learn

Each input needs to infer a large number of output-specific hypotheses

❏ Ultra-fine entity typing (Choi et al. 2018): 10K labels take the NLI model (Li et al., 2022) 35 seconds for each test instance and about 19.4 hours to infer the entire test set.
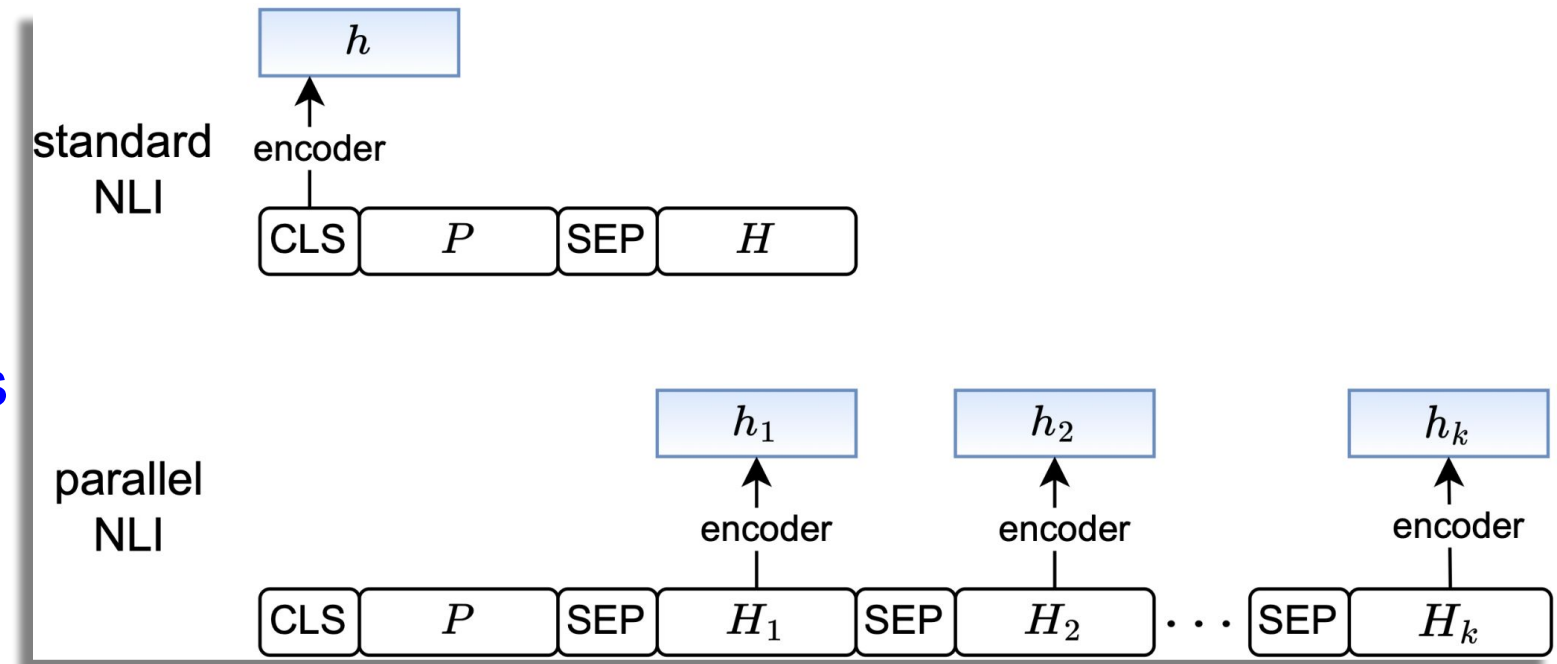
Each input needs to infer a large number of output-specific hypotheses
- ❏ Ultra-fine entity typing (Choi et al. 2018): 10K labels take the NLI model (Li et al., 2022) 35 seconds for each test instance and about 19.4 hours to infer the entire test set.

**Solution**:

- ❏ Pairwise inference → group-wise inference

- ❏ Assumption: hypotheses exhibit a binary polarity irrespective of their competitors.



_Learning to Select from Multiple Options_ (Du et al., AAAI'23)
_Recall, Expand and Multi-Candidate Cross-Encode: Fast and Accurate Ultra-Fine Entity Typing_ (Jiang et al., Arxiv, 2022)
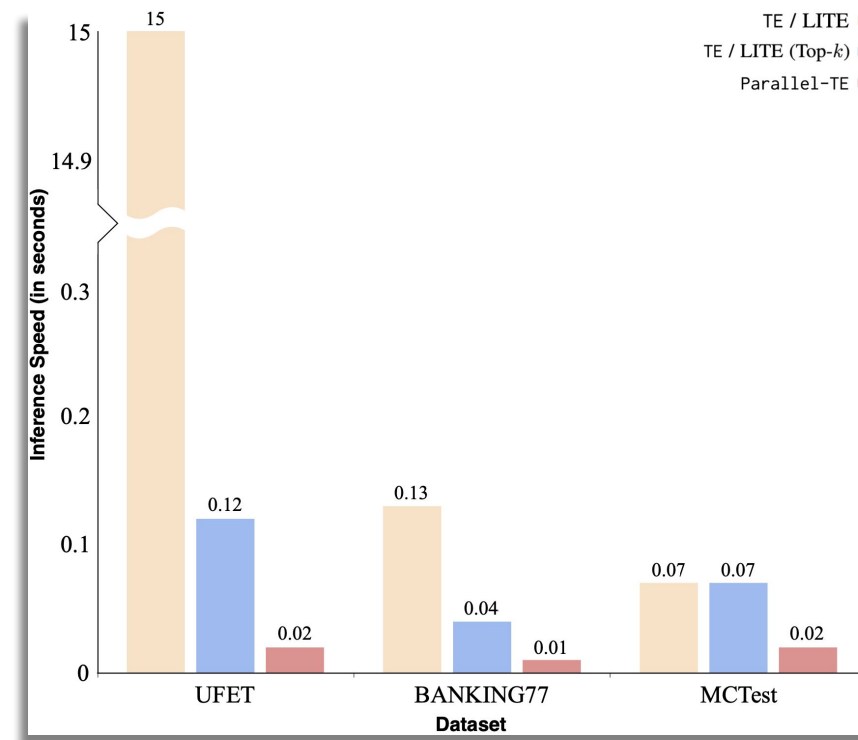
Each input needs to infer a large number of output-specific hypotheses

❑ Ultra-fine entity typing (Choi et al. 2018): 10K labels take the NLI model (Li et al., 2022) 35 seconds for each test instance and about 19.4 hours to infer the entire test set.

**Solution**:

❑ Pairwise inference → group-wise inference

❑ 15 seconds → 0.02 seconds (per example)



*Learning to Select from Multiple Options* (Du et al., AAAI'23)
*Recall, Expand and Multi-Candidate Cross-Encode: Fast and Accurate Ultra-Fine Entity Typing* (Jiang et al., Arxiv, 2022)

❏ We often pre-define the label set for classification tasks

❏ At times, we may want the model to generate some new labels for the input to "surprise" us

❏ We often pre-define the label set for classification tasks

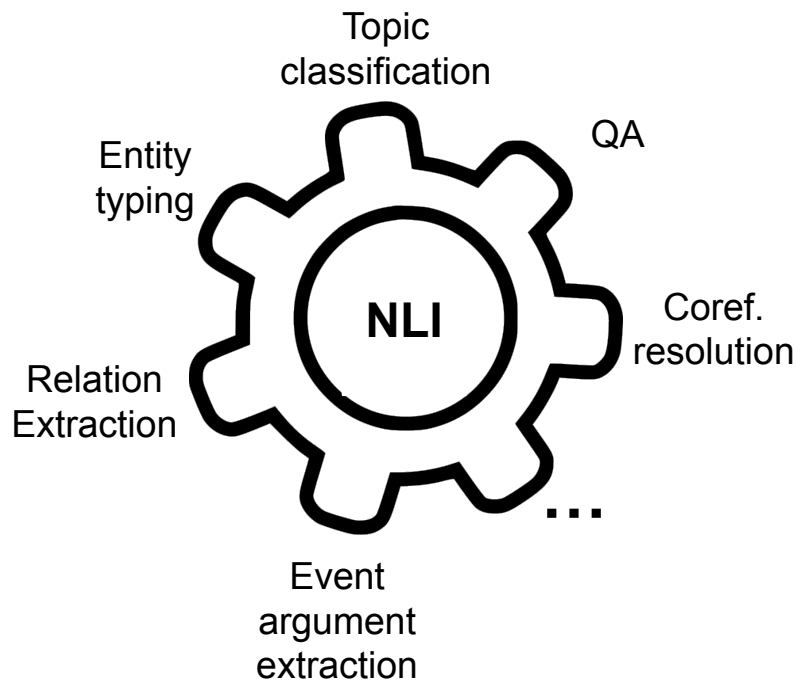❏ At times, we may want the model to generate some new labels for the input to "surprise" us



Will be addressed by the **next section** of our tutorial

# Recap of indirect supervision from NLI

## Implementation & Applications

Topic classification

Entity typing

QA

**NLI**

Coref. resolution

Relation Extraction

Event argument extraction

…

## Benefits

❏ Scarce-annotation NLP

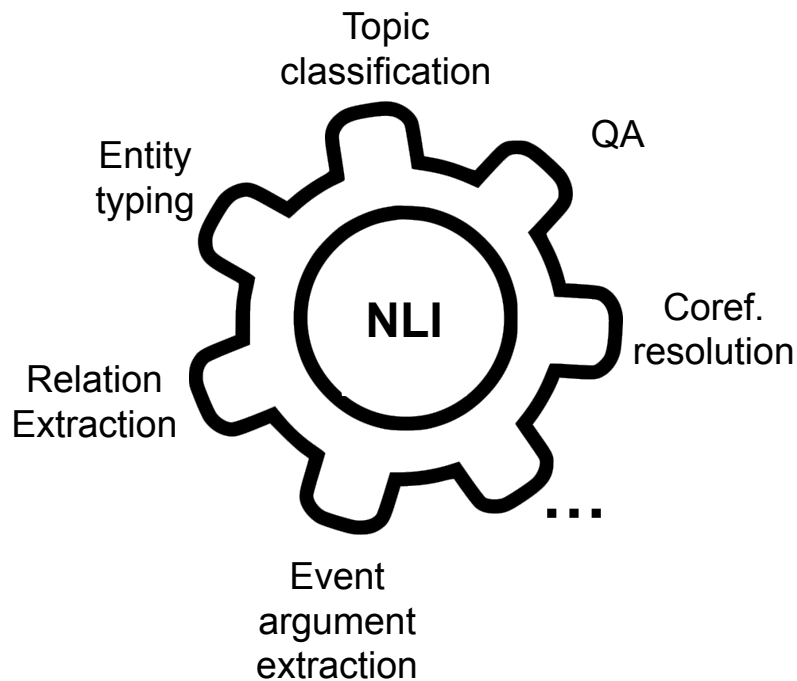❏ Cross-task transferability

❏ Maximize the potential of small PLMs

## Challenges & Solutions

❏ Domain discrepancy (solutions by algorithm and data threads)

❏ Inefficiency in testing (parallel-NLI)

❏ cannot discover new labels (next chapter…)

# Recap of indirect supervision from NLI

## Implementation & Applications

Topic classification

Entity typing

QA

**NLI**

Coref. resolution

Relation Extraction

...

Event argument extraction

## Benefits

❏ Scarce-annotation NLP

❏ Cross-task transferability

❏ Maximize the potential of small PLMs

## Challenges & Solutions

❏ Domain discrepancy (solutions by algorithm and data threads)

❏ Inefficiency in testing (parallel-NLI)

❏ cannot discover new labels (next chapter…)

*Thank You*