



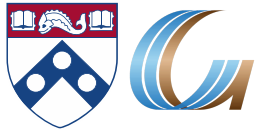
---

# Zero-shot Event Extraction via Transfer Learning: Challenges and Insights

Qing Lyu, Hongming Zhang, Elior Sulem, Dan Roth

Department of Computer & Information Science

University of Pennsylvania



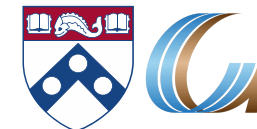
- An **event** is represented as a **trigger** + several **arguments**.
- Example from [ACE-2005](#):

Event type: TRANSFER-OWNERSHIP

China has purchased two nuclear submarines from Russia last month.

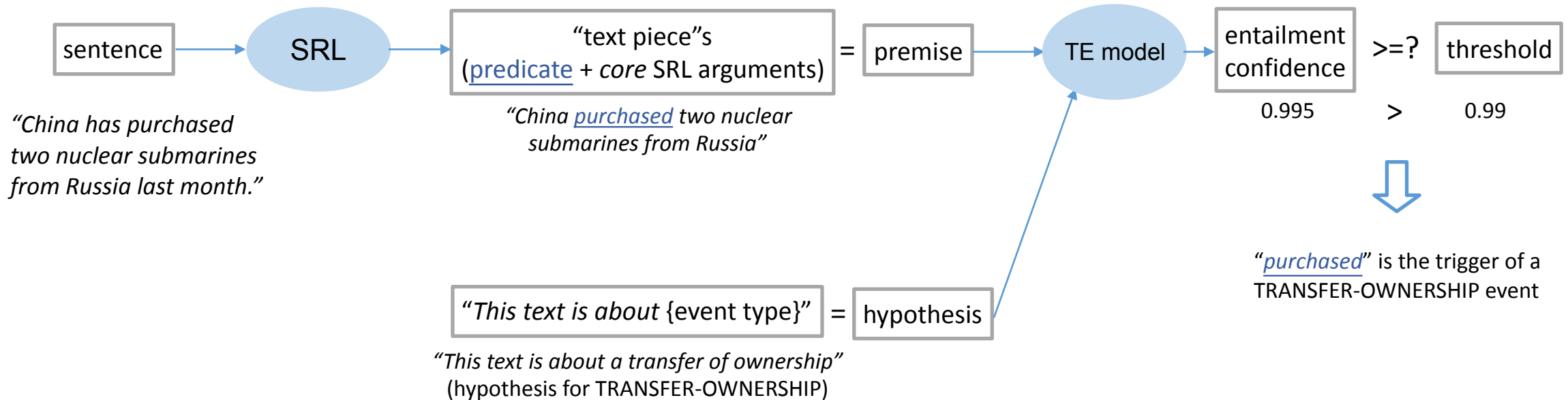
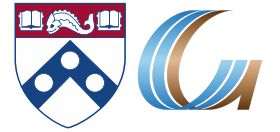
Buyer-Arg      Trigger                      Artifact-Arg                      Seller-Arg      Time-Arg

- **Event Extraction (EE)** = Trigger Identification (TI) + Trigger Classification (TC)  
+ Argument Identification (AI) + Argument Classification (AC)



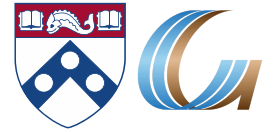
- Predominant approaches: **supervised**, both expensive & inflexible.
- Recent efforts explored zero-shot event extraction, usually requiring some event types to be **seen** ([Huang et al., 2018](#)) / only dealing with triggers or arguments **alone** ([Peng et al., 2016](#); [Liu et al., 2020](#)).
- Their performance is still far from supervised methods, but little is known about **why**.
- Our work:
  - Proposes a zero-shot event extraction system that tackles **both** triggers and arguments without any **event training data**, via transfer learning from Question Answering (**QA**) / Textual Entailment (**TE**).
  - Provides insights into the remaining challenges behind the performance gap.

# Approach: Trigger Extraction

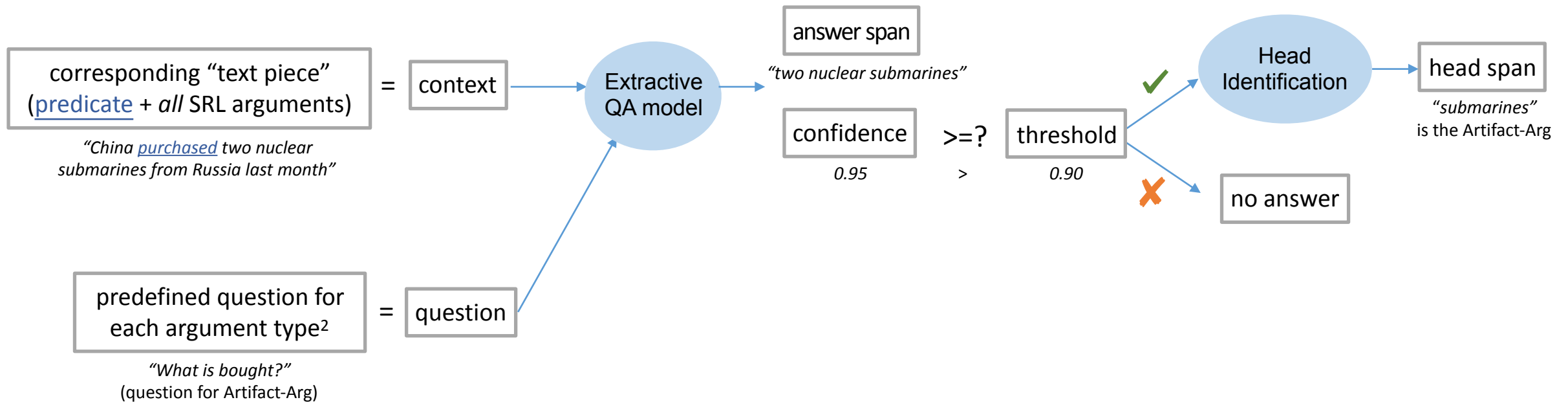


<sup>1</sup>An alternative uses Yes/No QA instead of TE, which is similar and thus not illustrated.

# Approach: Argument Extraction



For each extracted trigger (span + event type),  
e.g. “*purchased*” + TRANSFER-OWNERSHIP



<sup>2</sup>The questions are written based on the definition of each event type.

- **Dataset:** ACE-2005 (LDC2006T06), ERE (LDC2015E29)
- **Settings:**
  - *scratch*: the system performs all subtasks without any gold annotation
  - *gold TI*: gold trigger spans are given
  - *gold TI+TC*: gold trigger spans and types are given
- **Pretrained models<sup>3</sup>:**
  - Architecture: BERT/**RoBERTa**/BART - base/**large**
  - Pretraining data:

Target EE Subtask	Pretraining Dataset	Pretraining Task
Trigger Extraction	<b>MNLI</b> (Williams et al., 2018)	TE
	<b>BoolQ</b> (Clark et al., 2019)	Yes/No QA
Argument Extraction	<b>QAMR</b> (Michael et al., 2018)	Extractive QA
	<b>SQuAD2.0</b> (Rajpurkar et al., 2018)	Extractive QA

Table 1: Datasets used to pretrain the TE/QA models.

<sup>3</sup>Optimal configuration highlighted in **green**.

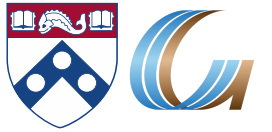
Setting	System	TI	TI+TC	AI	AI+AC
scratch (supervised)	Lin et al. 20	78.2	74.7	59.2	56.8
scratch (zero-shot)	Huang et al. 18	55.6	49.1	<b>27.8</b>	15.8
	Zhang et al. 20	<b>58.3</b>	<b>53.5</b>	16.3	6.3
	Ours	45.5	41.7	27.0	<b>16.8</b>
gold TI (zero-shot)	Huang et al. 18	-	33.5	-	14.7
	Zhang et al. 20	-	82.9	-	-
	Ours	-	<b>83.7</b>	<b>38.9</b>	<b>24.2</b>
gold TI+TC (zero-shot)	Liu et al. 20	-	-	-	25.8
	Ours	-	-	<b>44.3</b>	<b>27.4</b>

Table 2: The F1 score on ACE-2005. SOTA results among zero-shot methods are in boldface.

Setting	System	TI	TI+TC	AI	AI+AC
scratch <b>(supervised)</b>	Lin et al. 20	68.4	57.0	50.1	46.5
scratch		39.8	31.8	23.0	15.0
gold TI	<b>Ours</b>	-	58.4	30.8	18.8
gold TI+TC <b>(zero-shot)</b>		-	-	47.9	27.5

Table 3: The F1 score on the ERE. The optimal model is chosen on ACE dev and directly evaluated on ERE.





- **Remaining challenges:** Manually annotated in 100 wrong predictions
- **Error attribution:**
  - **Model-Error:** the intrinsic fragility of pretrained TE/QA models
  - **Usage-Error:** our usage of the models
  - **Task-Error:** the task itself
- **Ablation study:**

To isolate their individual impact, we alter *certain conditions* that have caused the target error<sup>4</sup>, and see *how many errors are corrected after when predicting again*.

<sup>4</sup>See Section 5 of our paper for details.

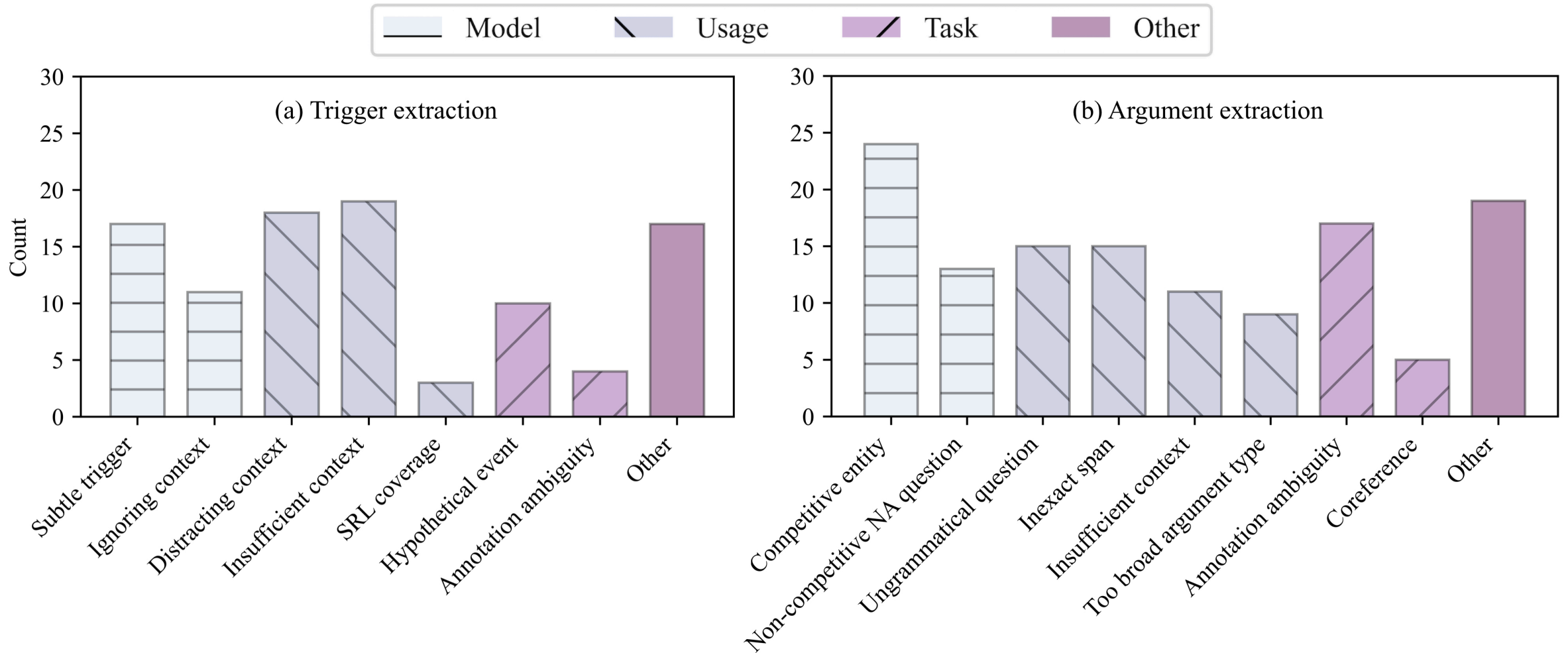


Figure 1: Error types in trigger and argument extraction in 100 wrong predictions. The count sum exceeds 100 since a prediction can contain multiple types of error.

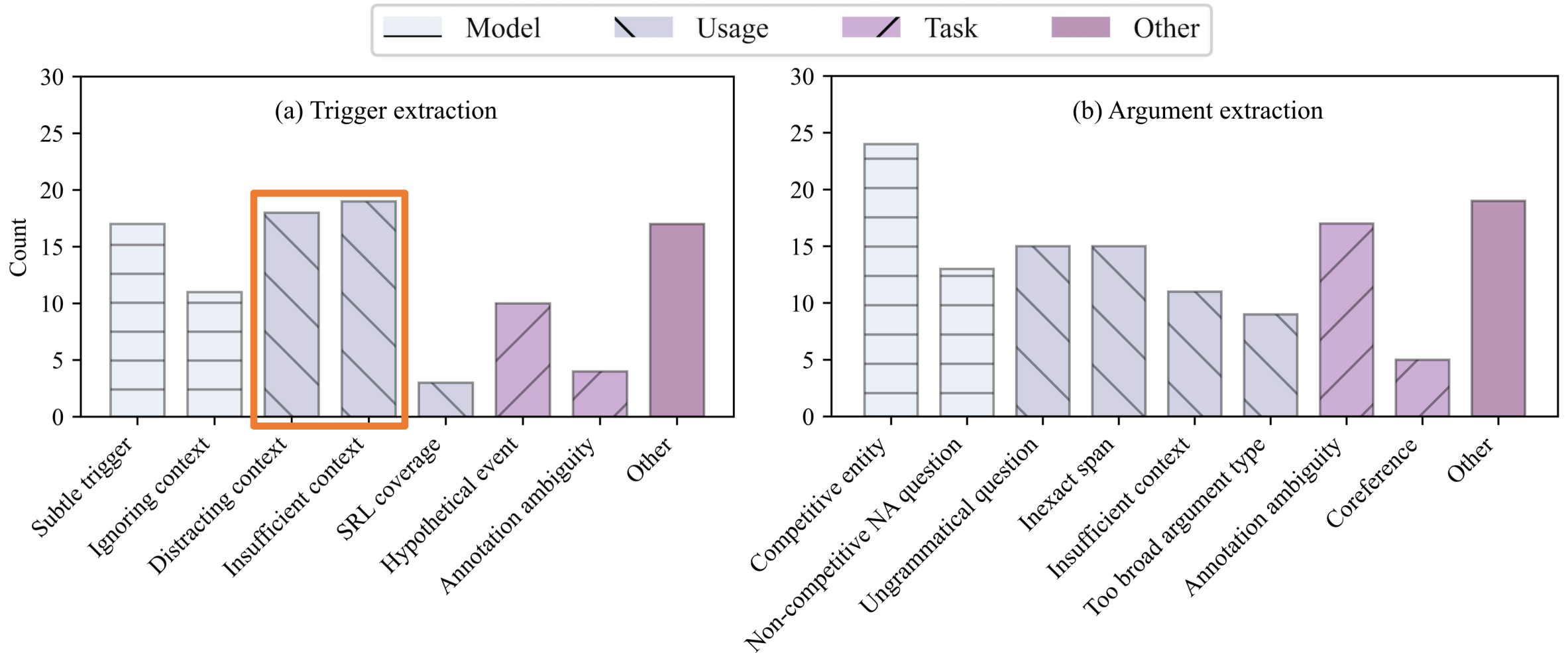
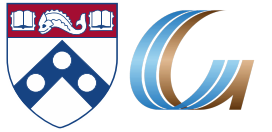


Figure 1: Error types in trigger and argument extraction in 100 wrong predictions. The count sum exceeds 100 since a prediction can contain multiple types of error.



- **Distracting Context (18%): Usage-Error**
  - e.g. “The woman’s parents ... found the decomposing body.”  
Gold type: Not a trigger      Predicted type: DIE
- **Insufficient Context (19%): Usage-Error**
  - e.g. “(Turkey sent 1,000 troops ... and said) it would send more”  
Gold type: TRANSPORT      Predicted type: TRANSFER-MONEY
- Ablation study: 18% **Distracting Context** errors and 59% **Insufficient Context** errors are corrected when predicting again.

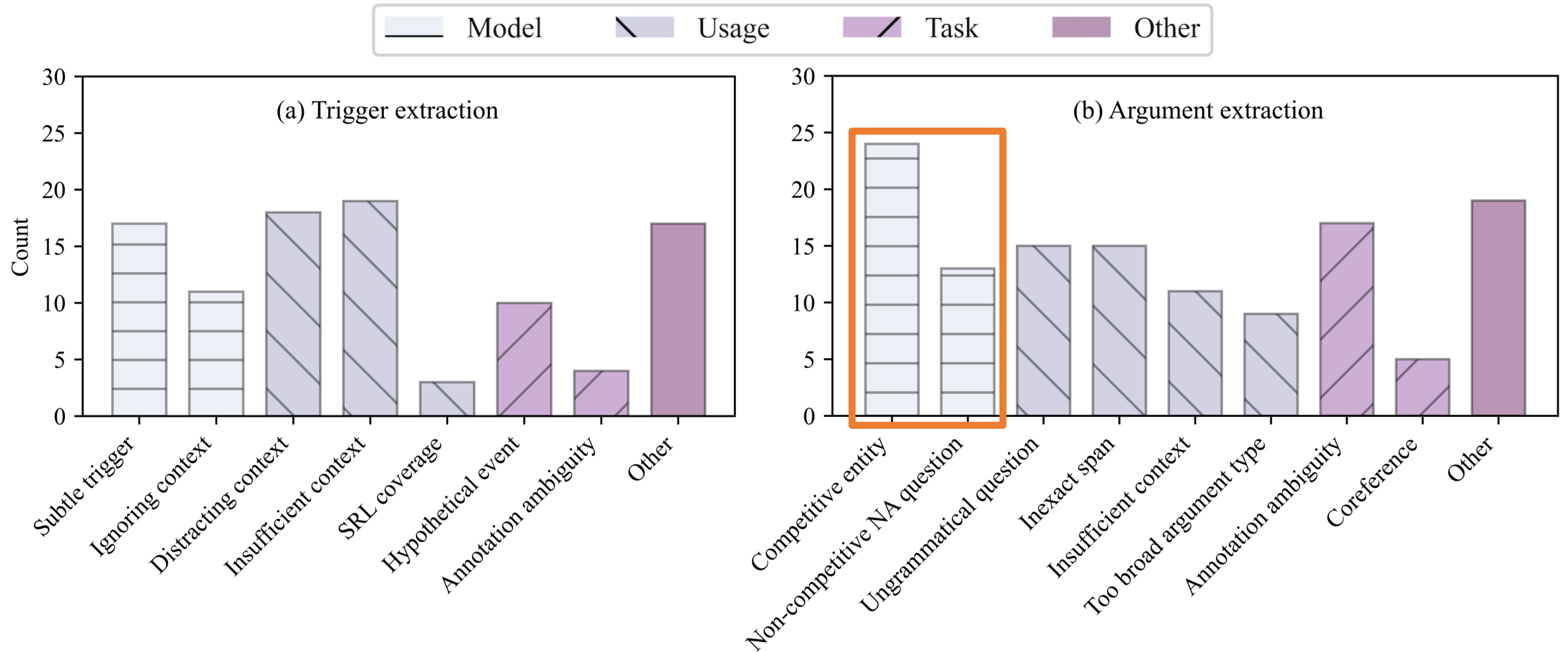
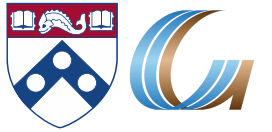


Figure 1: Error types in trigger and argument extraction in 100 wrong predictions. The count sum exceeds 100 since a prediction can contain multiple types of error.



- **“Competitive” Entity (24%): Model-Error**

- e.g. “A unit meets in confidential sessions to review terrorist activities in Europe.”

Question for *Place-Arg*: “Where is the meeting?”

Gold answer: No Answer      Predicted answer: “Europe”

- **Non-competitive NA Questions (19%): Model-Error**

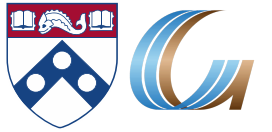
- e.g. “Iraqi forces responded with artillery fire.”

Question for *Time-Arg*: “When is the fire?”

Gold answer: No Answer      Predicted answer: “artillery”

- Ablation study: Adding training data on NA questions (SQuAD2.0) even hurts the performance<sup>5</sup>.

<sup>5</sup>We propose and test three hypotheses behind this. See Section 5.2.2 of our paper for details.



- We propose the first complete zero-shot event extraction system via transfer learning from TE and QA.
- While QA/TE models perform exceptionally well on standard benchmarks (SQuAD, QAMR, MNLI), they do not generalize as expected when being used on event extraction datasets.
- We analyze the limited success and several main challenges of this promising approach, and point out future research directions.

Thank you for listening!