

Do Language Models Learn about Legal Entity Types during Pretraining?

Claire Barale¹, Michael Rovatsos¹, Nehal Bhuta²
School of Informatics¹, School of Law², The University of Edinburgh

Problem definition

- During pretraining, LMs learn signals from the corpus: **what kind of signals and cues do they capture?**
Factual and abstract? Single token and multi-token entities?
Generic and legal vocabulary? Semantic and syntactic signals?
- Legal entity typing is a foundation: (i) for incidental supervision and text labeling, (ii) for downstream tasks, (iii) for text comprehension
- Task:** Entity Typing (predict the entity type) zero-shot
- Dataset:** AsyLex, a dataset of refugee decisions
- Models:** Encoder-only vs. Decoder-only

Objective: Evaluating the quality of entity knowledge learned during pretraining in LMs, a surrogate for legal knowledge

Key takeaways

- Llama2 performs well on certain generic entity types
- Optimized prompt templates are necessary
- Law-oriented LMs show inconsistent performance
- LMs can also type multi-token entities accurately
- Entities that belong to sub-domains of the law are difficult
- Llama2 frequently overlooks syntactic cues, contrary to BERT-based architectures

Vocabulary Overlap in the Pretraining Corpora

Vocabulary overlap (%) between the pretraining corpora. Gen stands for Generic and is sampled from sources similar to RoBERTa's pretraining corpus. Vocabularies are created with the top 10K most frequent tokens in a sample of 50K documents per model

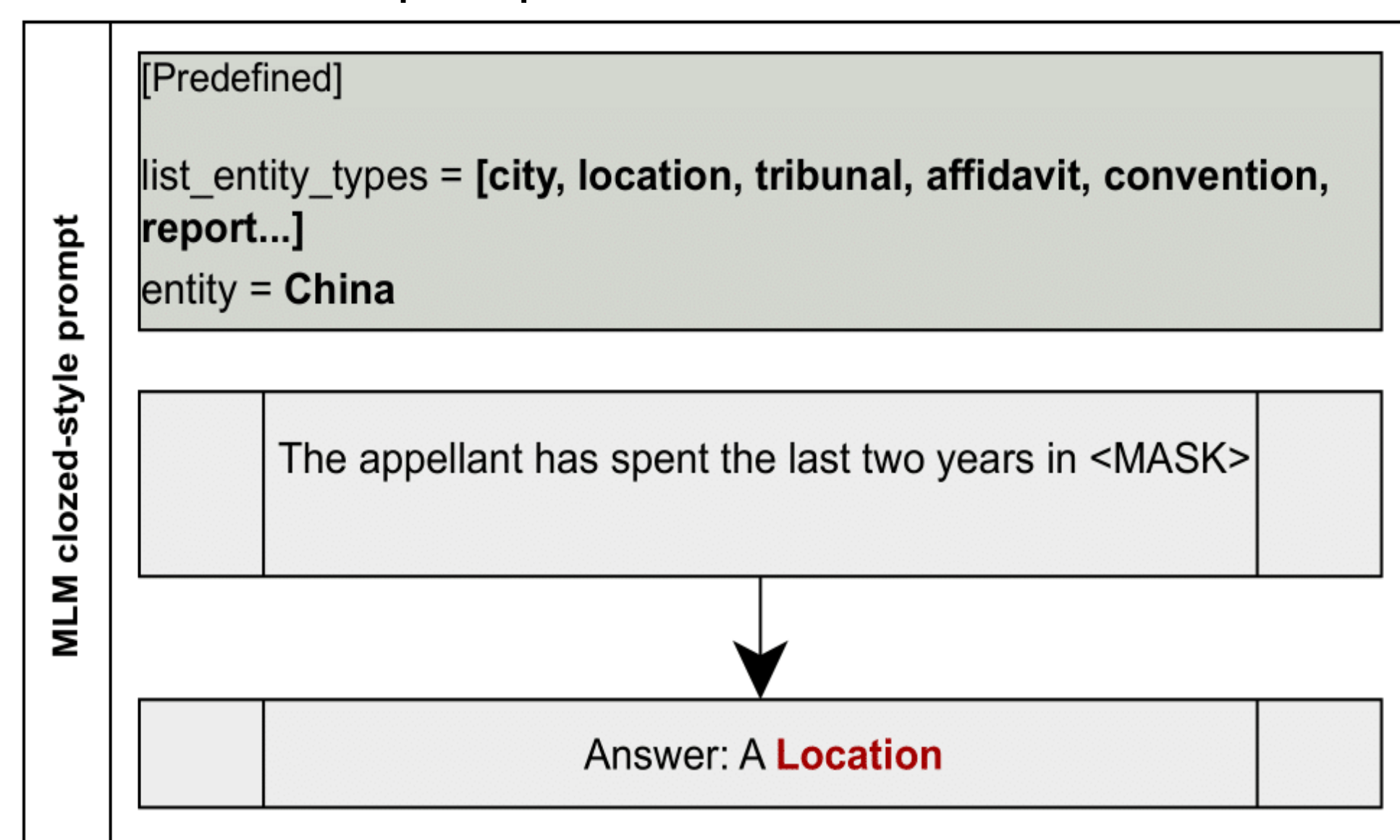
Gen	100.0	34.8	46.2	41.2
CH	34.8	100.0	55.7	44.5
PoL	46.2	55.7	100.0	55.1
LexLM	41.2	44.5	55.1	100.0
	Gen	CH	PoL	LexLM

Research Questions

- How proficient are LMs at **acquiring knowledge** about **domain-specific** entities like legal entities during pretraining?
- Can this acquired knowledge be considered sufficiently reliable for tasks such as **annotating new datasets** or serving as an indirect source of supervision?
- How does the choice of **prompt type** impact the results obtained from knowledge queries?
- To what extent does the variation in acquired knowledge differ across entity types? What categories of factual knowledge can LMs retrieve, and in what instances do they make errors? **Which signals are captured**, and which are missed?
- Does **domain-specific pretraining** and jurisdiction-specific pretraining enhance the amount of factual knowledge compared to generic pretraining?

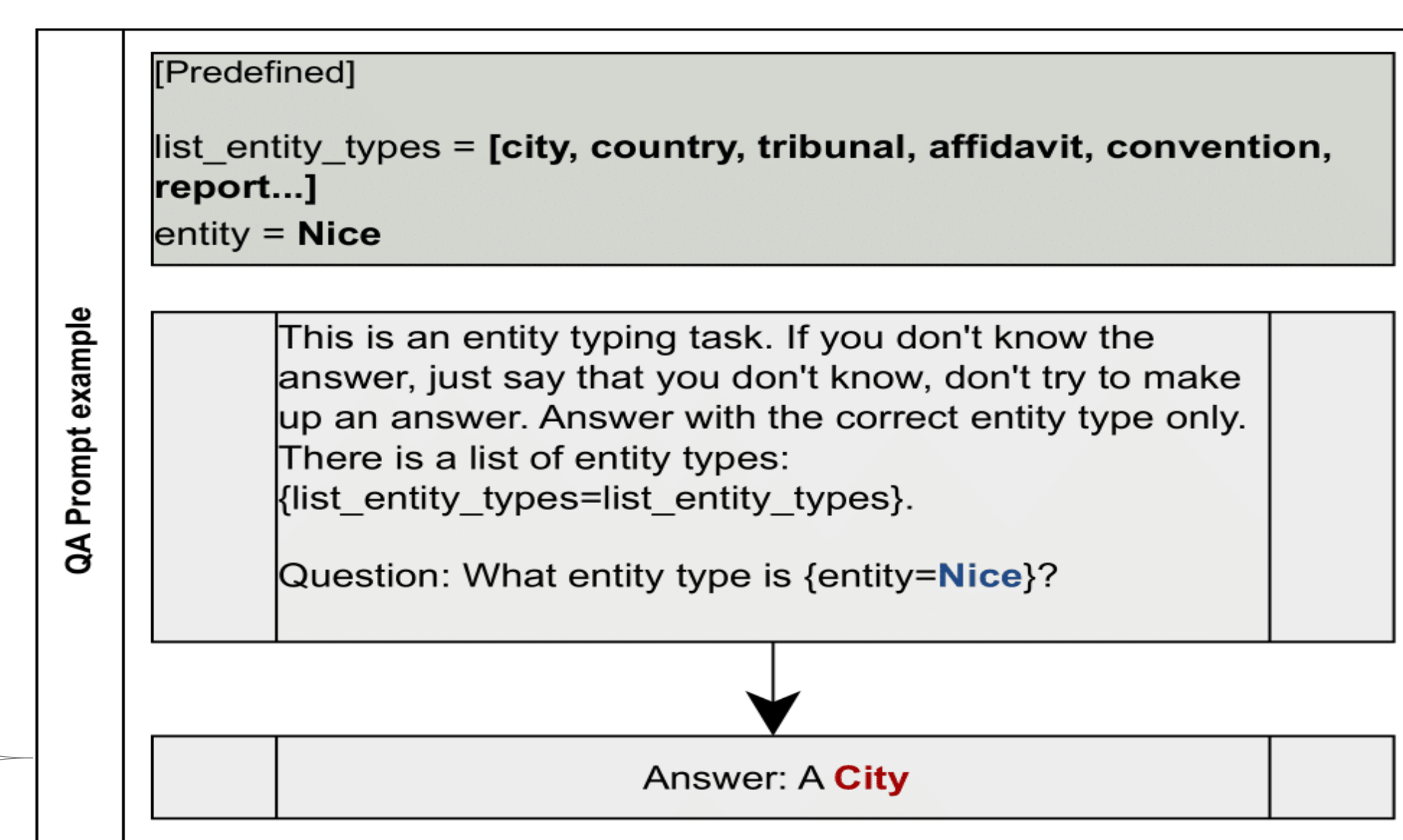
Prompt Templates and Models Used

Experiment 1: MLM Cloze prompts with BERT-based models



Model	Corpus	# Param.
ROBERTA	Generic	125M
DEBERTA-v3	Generic	86M
CASEHOLD	Harvard Case Law	110M
PILE OF LAW	US, Canada, ECtHR	340M
LEXLM	US, Canada, UK, India	125M
LLAMA2-7B	Generic	7B

Experiment 2: Llama2 QA prompts



Results

Gen groups the results of RoBERTa and DeBERTa-v3, CH refers to CaseHOLD

Model	Gen	CH	PoL	LexLM	Llama2
GENERIC	11.59	8.51	20.42	11.97	63.26
GEN LEGAL	17.98	20.89	15.40	12.48	29.52
REFUGEE LAW	10.01	5.93	4.68	4.92	13.03

Entity types prediction F1 scores averaged on 3 aggregated groups:

- Generic: location, date, norp (adjectives)
- Gen Legal - Legal entities applicable to most legal domains: org, law (citations), claimant_info, procedure, doc_evidence, law_case (precedents)
- Refugee Law - Legal entities specific to refugee law: credibility, determination, explanation, legal_ground, law_report (NGO reports)

Failure Cases Analysis

	Experiment 1	Experiment 2
RANDOM PREDICTION	70.71	22.22
CONTEXTUALLY ACCURATE	12.43	-
CLOSELY RELATED	16.86	18.52
FALSE NEGATIVE	-	33.33
PROMPT ERROR	-	25.93

Examples

under <mask> of the Republic of China, they cannot take on a second citizenship.

Predicted: lawsuit. Gold: law

What is female claimant? Predicted: female claimant. Gold: gender (claimant information)

Scan here for paper

