

Research on Named Entity Recognition from Patent Texts with Local Large Language Model

Yu Chi¹ **Liang Chen(Speaker)**^{1*} **Haiyun Xu**²

1. Institute of Scientific and Technical Information of China, Fuxing Road 15, 100038 Beijing, P.R.China

2. Business school, Shandong University of Technology, Xinchunxi 266, 255000 Zibo, P.R.China

Apr. 24th, 2024

EEKE-AII workshop at iConference2024, Changchun, China and Online

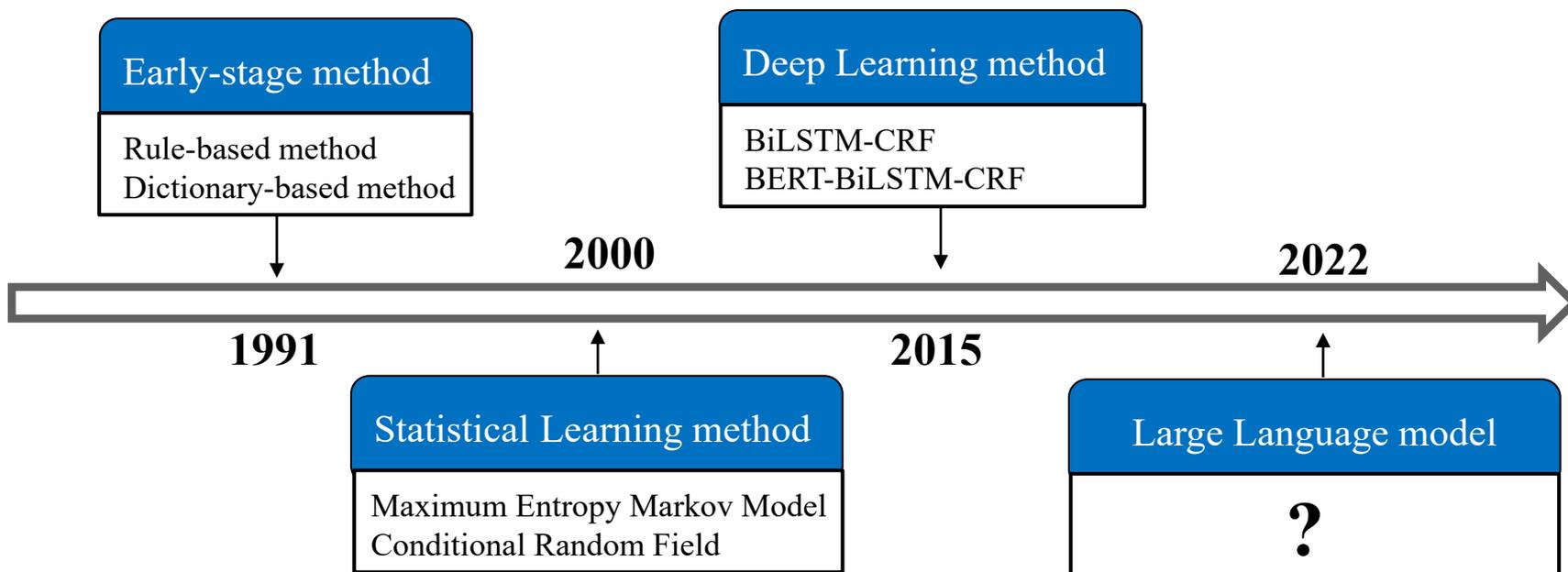


Content

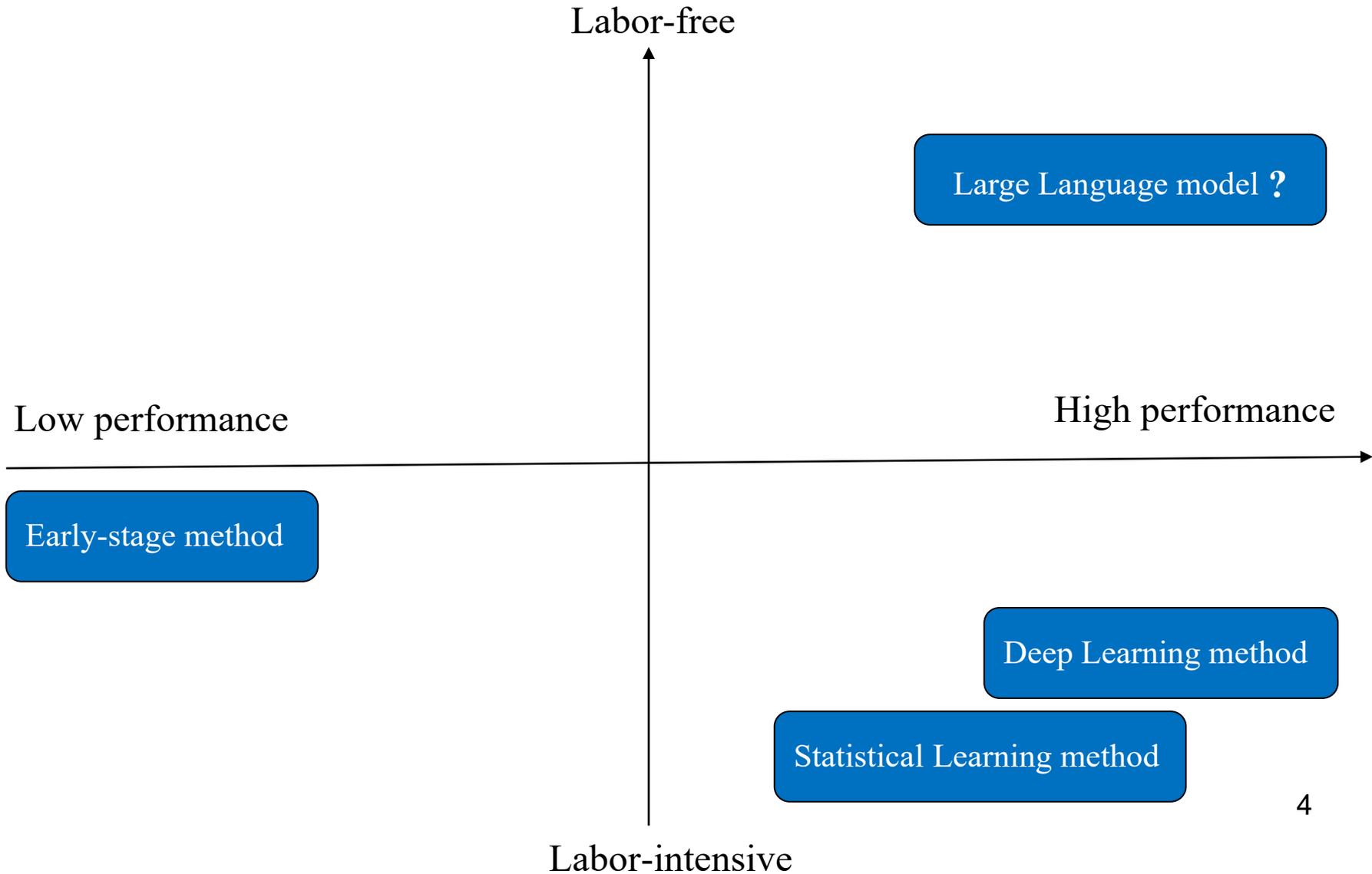
1. Background
2. Methodology & Experimental Results
3. Conclusion

1. Background

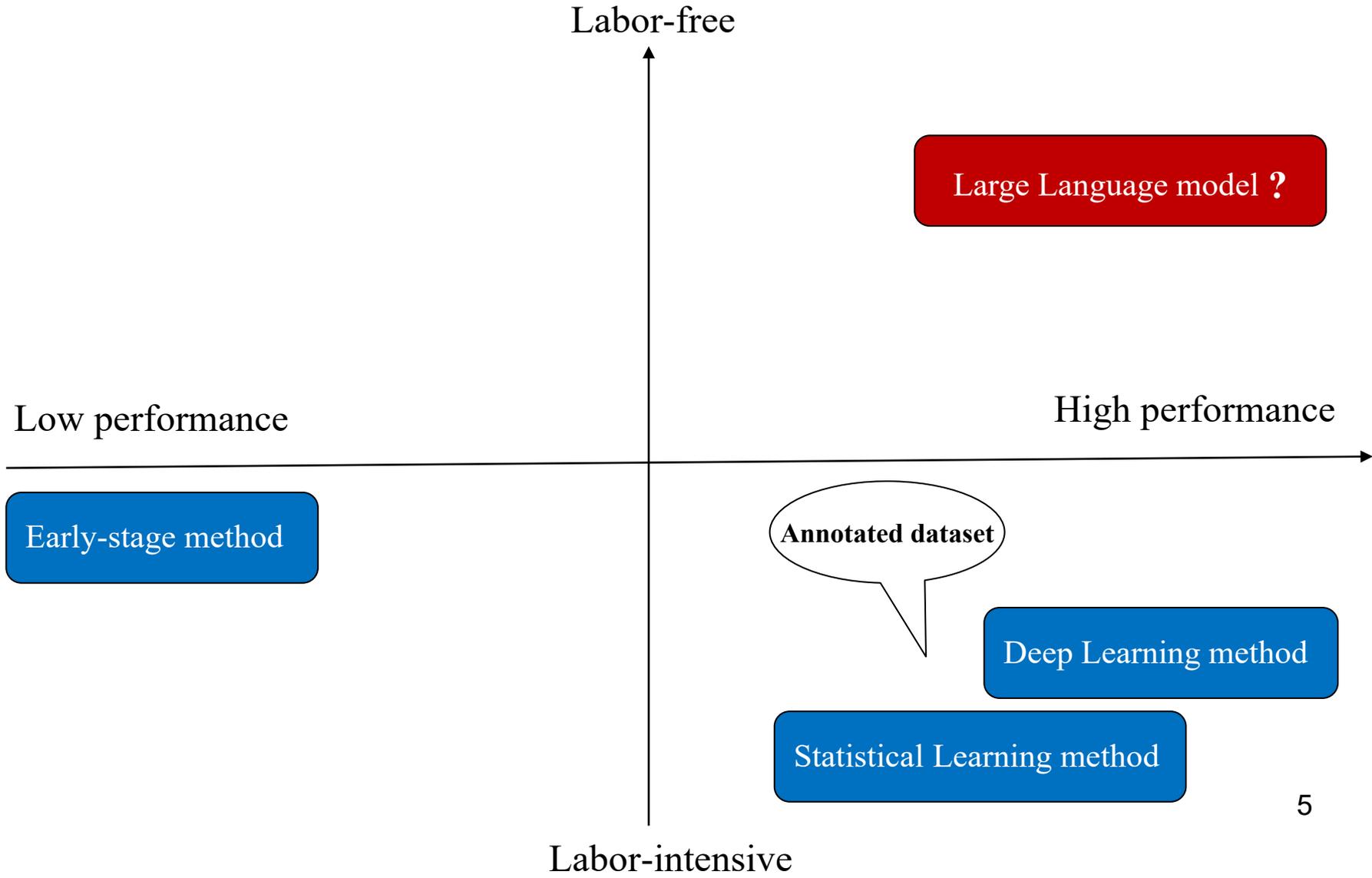
The technological evolution of NER methods



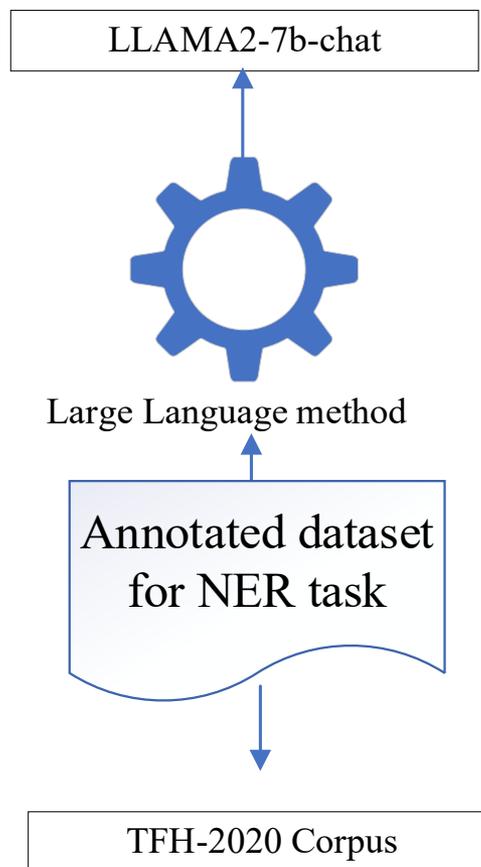
1. Background



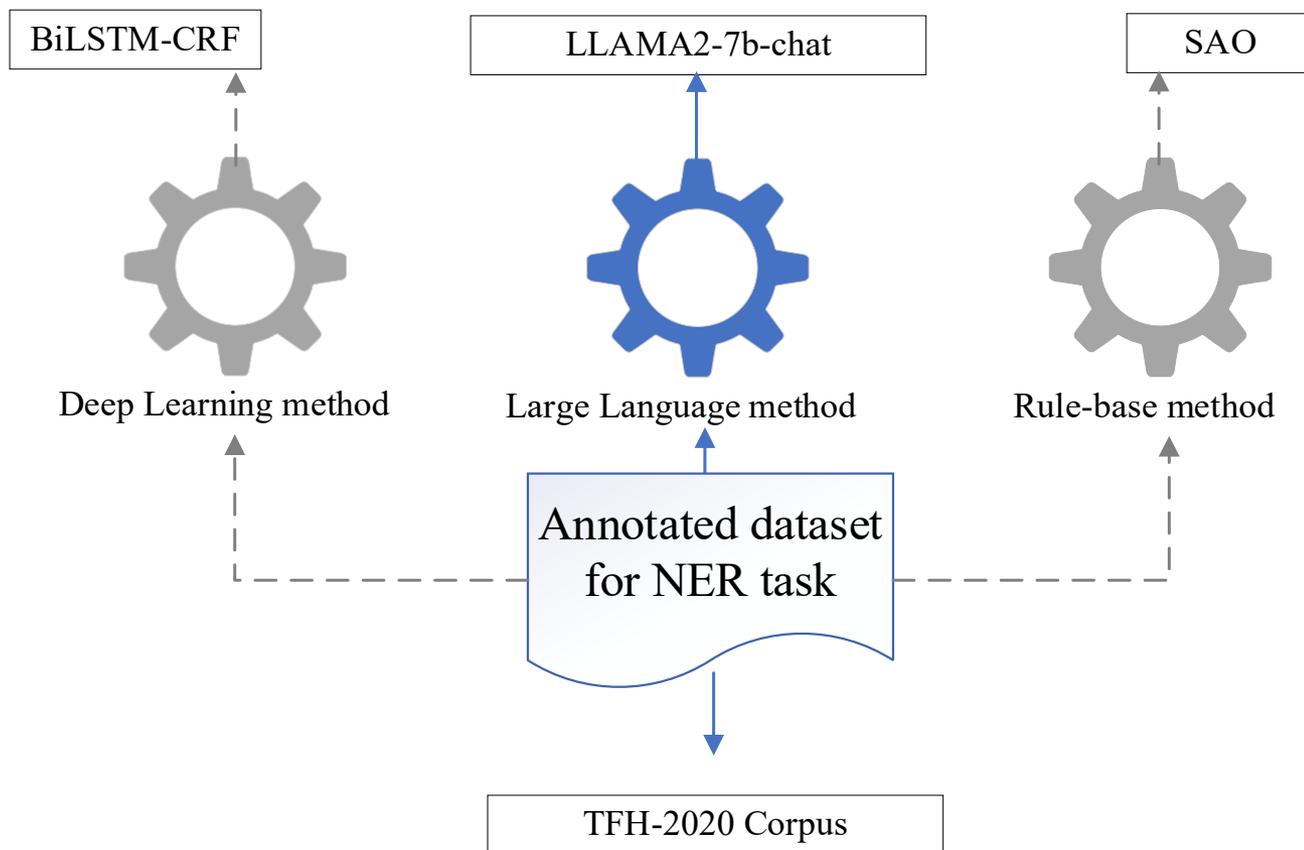
1. Background



2. Methodology & Experiment



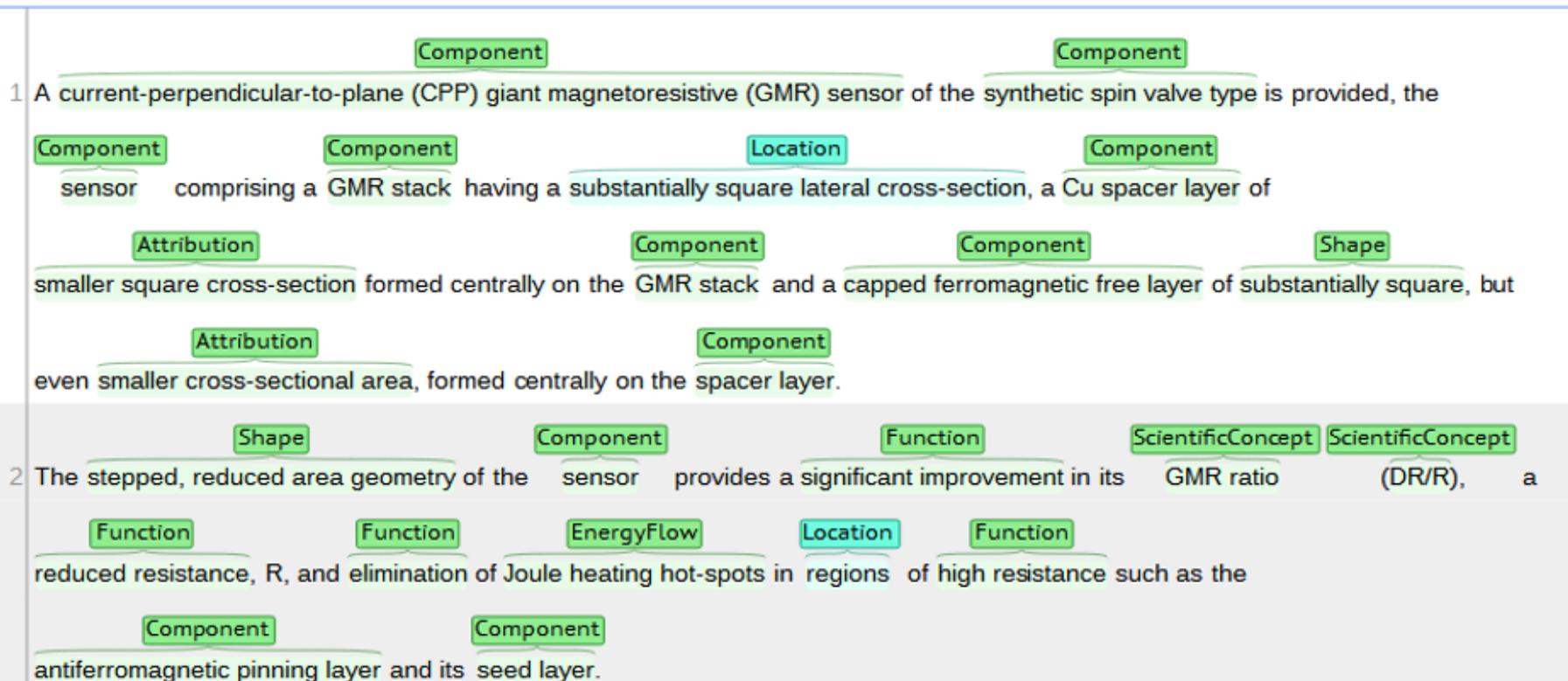
2. Methodology & Experiment



2. Methodology & Experiment

Description of TFH-2020 patent dataset:

- **Domain:** Thin film head technology in hard disk drive
- **Data source:** USPTO
- **Number of Entity mentions:** 22833
- **Number of Sentence:** 3996



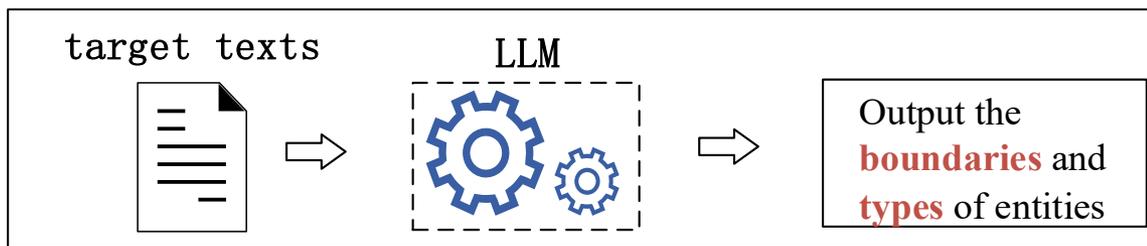
2.Methodology& Experiment

URL: https://github.com/awesome-patent-mining/TFH_Annotated_Dataset

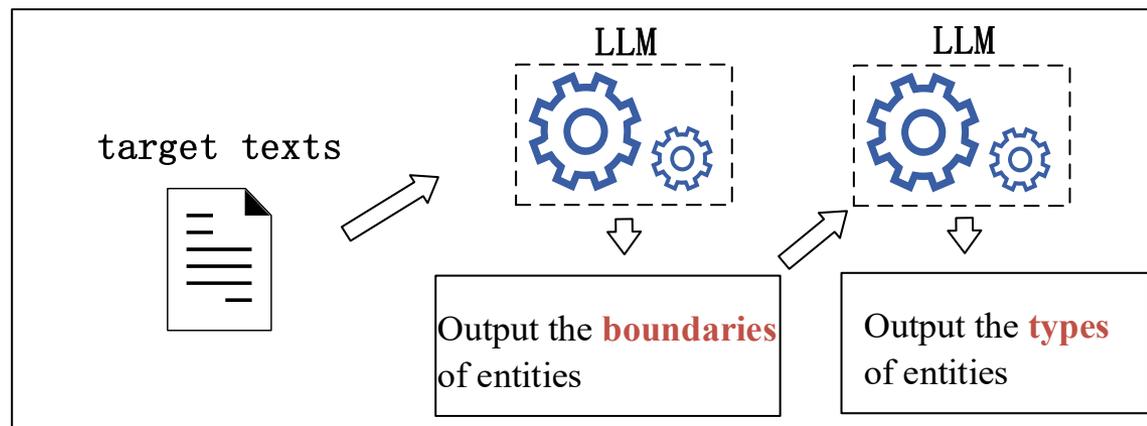
Type	Comment	Example
Physical flow	Substance that flows freely	The etchant solution has a suitable solvent additive such as glycerol or methyl cellulose
Information flow	Information data	A camera using a film having a magnetic surface for recording magnetic data thereon
Energy flow	Entity relevant to energy	Conductor is utilized for producing writing flux in magnetic yoke
Measurement	Method of measuring something	The curing step takes place at the substrate temperature less than 200.degree
Value	Numerical amount	The curing step takes place at the substrate temperature less than 200.degree
Location	Place or position	The legs are thinner near the pole tip than in the back gap region
State	Particular condition at a specific time	The MR elements are biased to operate in a magnetically unsaturated mode
Effect	Change caused an innovation	Magnetic disk system permits accurate alignment of magnetic head with spaced tracks
Function	Manufacturing technique or activity	A magnetic head having highly efficient write and read functions is thereby obtained
Shape	The external form or outline of something	Recess is filled with non-magnetic material such as glass
Component	A part or element of a machine	A pole face of yoke is adjacent edge of element remote from surface
Attribution	A quality or feature of something	A pole face of yoke is adjacent edge of element remote from surface
Consequence	The result caused by something or activity	This prevents the slider substrate from electrostatic damage
System	A set of things working together as a whole	A digital recording system utilizing a magnetoresistive transducer in a magnetic recording head
Material	The matter from which a thing is made	Interlayer may comprise material such as Ta
Scientific concept	Terminology used in scientific theory	Peak intensity ratio represents an amount hydrophilic radical
Other	Not belongs to the above entity types	Pressure distribution across air bearing surface is substantially symmetrical side

2. Methodology & Experiment

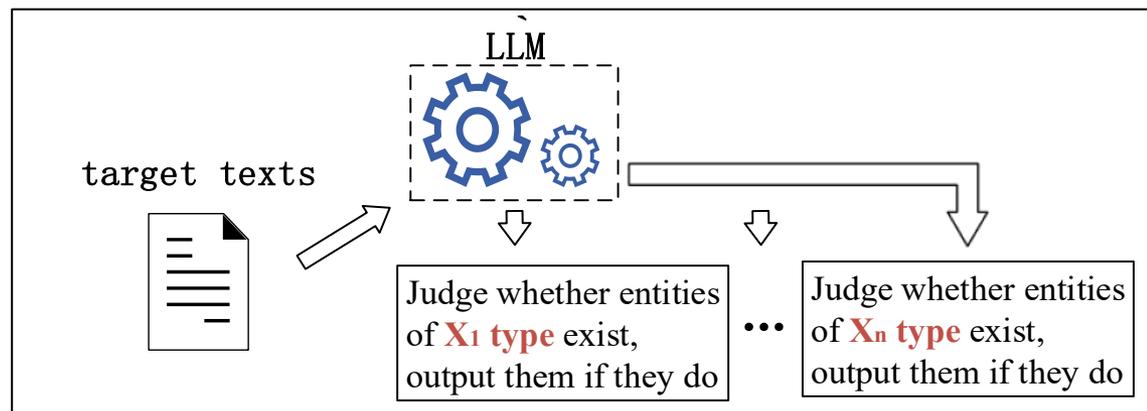
(1) Baseline Prompt



(2) Two-step Prompt



(3) multi-entity-type Prompt



2.Methodology& Experiment

Baseline Prompt Example: Follow Alpaca Instruction Template

Below is an instruction that describes a task, paired with an input that provides further context. Write a response that appropriately completes the request.

Instruction:

The following is a content of the patent for a hard disk head, the task is to identify named entities, including 12 types of entities suitable for describing process methods and processes: system, part, function, effect, consequence, attribute, measurement, numerical value, orientation, material, shape, scientific concept, and 4 types of entities that highlight the state of existence and physical characteristics: Physical flows, information flows, energy flows, states, as well as entities that still exist that are not covered, are classified as 'other'. Extract entities from the following text and return it in json format.

To help you understand the task better, some examples are provided. Please use these examples as a reference to format your entity list correctly.

input: A thin film structure suitable for use as a shield for a read element of a transducing head has a first ferromagnetic layer, a second ferromagnetic layer, a spacer layer and a bias layer.

Output: {"System": ["thin film structure"], "Component": ["shield", "read element", "transducing head", "first ferromagnetic layer", "second ferromagnetic layer", "spacer layer", "bias layer"]}

examples

You need to extract the entity from the following text. Ensure that the data is returned accurately in the format of the provided example.

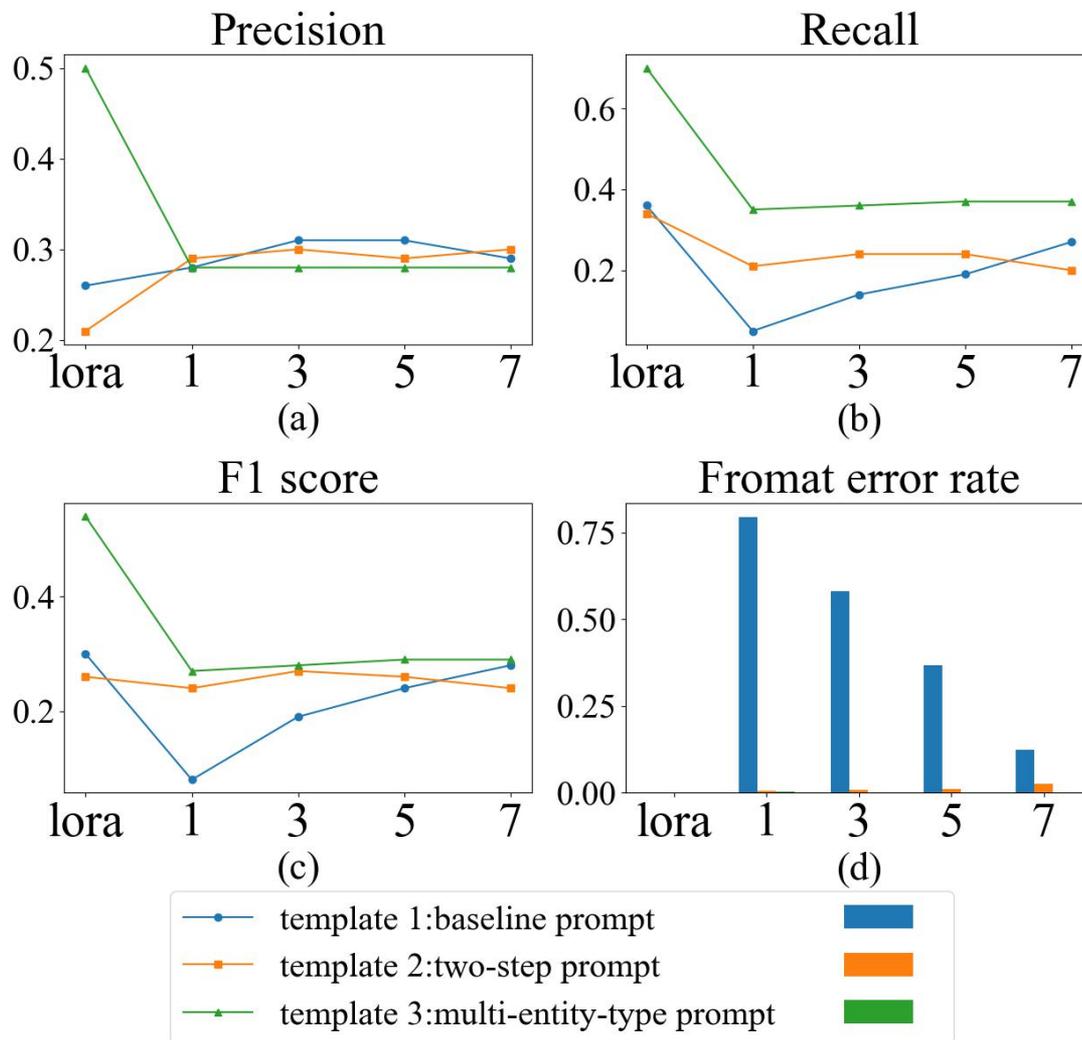
Input:

This invention relates to a multi-layer lithographically fabricated device used to produce improved thin-film recording heads. It further relates to a focused particle beam system for milling a recording head pole-tip assembly without irradiating a sensitive structure, e.g. a read head, of the recording head.

Response:

2.Methodology& Experiment

Experimental results of LLM with and without efficient fine-tuning



2.Methodology& Experiment

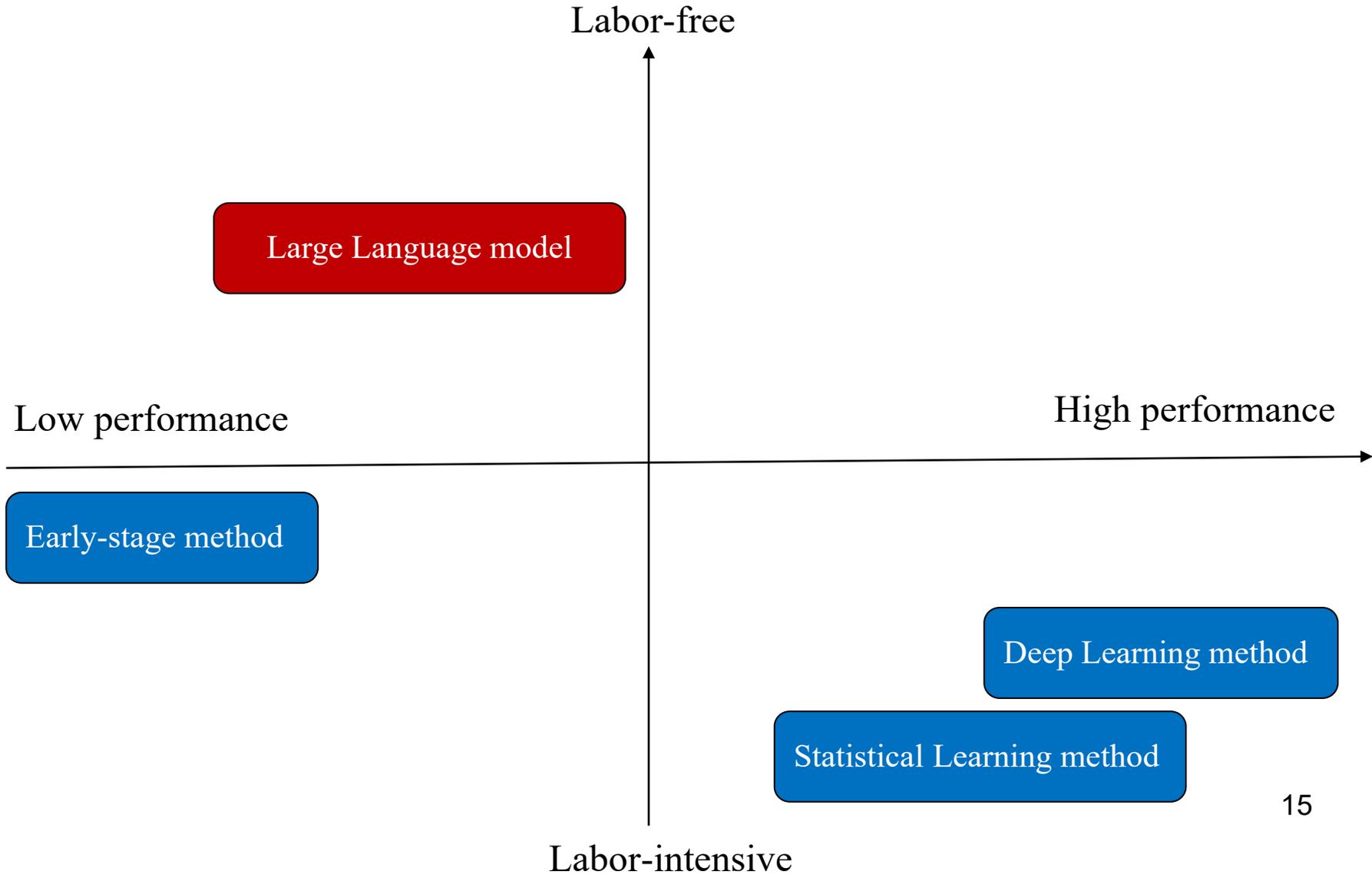
Experimental results of NER with different methods

Weighted-average	LLM with efficient fine-tuning		BiLSTM-CRF	PC-LDA	SAO
	YES	NO			
P (%)	50.0	28.0	78.5	20.5	3.0
R(%)	70.0	37.0	78.0	9.7	1.2
F1(%)	54.0	29.0	78.2	13.2	1.7

3. Conclusion

- ✓ Prompt template is crucial for NER task with LLM, it not only determines the quality of response yielded by LLM, but the performance of efficient fine-tuning method as well.
- ✓ For the task of NER from patent texts, there is a significant gap between LLM and the SOTA method, namely pretrain model with fine-tuning algorithm.
- ✓ Without annotated patent dataset, LLM will outperform previous methods, such as SAO, topic models by a large margin.

3. Conclusion





Thanks For Your Attention!
Q&A