

llmNER: (Zero|Few)-Shot Named Entity Recognition, Exploiting the Power of Large Language Models

Fabián Villena^{a,b}, Luis Miranda^{b,c}, Claudio Aracena^{b,d}

^a*Department of Computer Science, University of Chile, Chile*

^b*Millennium Institute for Foundational Research on Data, Chile*

^c*Department of Computer Science, Catholic University of Chile, Chile*

^d*Faculty of Physical and Mathematical Sciences, University of Chile, Chile*

Abstract

Large language models (LLMs) allow us to generate high-quality human-like text. One interesting task in natural language processing (NLP) is named entity recognition (NER), which seeks to detect mentions of relevant information in documents. This paper presents `llmNER`, a Python library for implementing zero-shot and few-shot NER with LLMs; by providing an easy-to-use interface, `llmNER` can compose prompts, query the model, and parse the completion returned by the LLM. Also, the library enables the user to perform prompt engineering efficiently by providing a simple interface to test multiple variables. We validated our software on two NER tasks to show the library's flexibility. `llmNER` aims to push the boundaries of in-context learning research by removing the barrier of the prompting and parsing steps.

Keywords: Natural language processing, Named entity recognition, Large language models, In-context learning, Few-shot learning

Metadata

1. Motivation and significance

NLP has gained tremendous importance in recent years with the advent of Transformer-based pre-trained language models (PLM) [1]. These language models (LM) have become the new paradigm for NLP-based machine learning modeling because of their modularity and ease of transferring learning. Nowadays, one can fine-tune a off-the-shelf PLM to solve any NLP task using off-the-shelf [2].

Email addresses: `fvillena@imfd.cl` (Fabián Villena), `lmirandn@uc.cl` (Luis Miranda), `claudio.aracena@uchile.cl` (Claudio Aracena)

Nr.	Code metadata description	
C1	Current code version	0.1.0
C2	Permanent link to code/repository used for this code version	https://github.com/plncmm/llmner
C3	Permanent link to Reproducible Capsule	https://github.com/plncmm/llmner/blob/main/notebooks/1-example.ipynb
C4	Legal Code License	Apache License, Version 2.0
C5	Code versioning system used	git
C6	Software code languages, tools, and services used	python, LangChain
C7	Compilation requirements, operating environments & dependencies	Python \geq 3.10, git
C8	If available Link to developer documentation/manual	https://github.com/plncmm/llmner/blob/main/README.md/
C9	Support email for questions	fvillena@proton.me

Table 1: Code metadata (mandatory)

PLMs trained on web-scale unannotated text, such as BERT [3] and BERT-like models, such as RoBERTa [4] and DeBERTa [4] have been the *de facto* standard for solving NLP tasks. These pre-trained Transformer-based context-aware models have shown outstanding performance, but their size limits capacity. LLMs are PLMs with a significantly larger model size scale [5]; for example, BERT has a model size of 0.3×10^9 parameters, and the LLM GPT-3 [6], has 175×10^9 parameters. Additionally, it has been found that scaling PLMs improves the performance of the models on downstream tasks [7].

Besides superior performance in downstream tasks, LLMs show surprising and more important behaviors in solving complex tasks, called emergent abilities. Emergent abilities are aptitudes not present in small models but arise in LLMs [8] and include in-context learning, where a model can generate expected outputs to natural language instructions without additional training; instruction following, where a model fine-tuned using natural language instructions performs well on unseen tasks that are also described in the form of instructions; and *step-by-step reasoning*, where a model can solve complex problems by instructing the model involving intermediate reasoning steps for deriving the final answer.

In-context learning (ICL) is a paradigm that allows language models to learn tasks given only a few examples in the form of demonstration. Although

LLMs themselves show good ICL capabilities, some techniques, such as supervised instruction-tuning, are used to improve ICL performance. Instruction tuning enhances ICL model capabilities by updating model parameters by fine-tuning the model on task instructions [9].

Few-shot learning is a machine learning problem where the model only has limited examples with supervised information to solve the task [10]. This branch of task-agnostic pre-training and task-agnostic model architectures has recently proliferated in the NLP area. By exploiting ICL, an LM can be contextualized to solve a task by only giving zero, one, or few examples without fine-tuning. This paradigm of solving a task with few to no examples solves the issue of the need for large datasets for training models for new tasks but with a significant penalty in performance compared to state-of-the-art fine-tuned models [6].

NER is the task of finding spans of text that constitute named entities (anything that can be referred to with a proper name) and tagging the entity type. The four most common entity types are person, location, organization, and geopolitical entity. Sometimes, the term named entity is extended to include things that are not entities, including dates, times, and even numerical expressions [11].

For training NER models, one needs a large *corpus* of task-specific annotated text, but constructing an annotated *corpus* is both time-consuming and expensive. Also, after the *corpus* is annotated, one needs to fine-tune a model (only if a foundation model is available for the domain), a process that is also expensive. We propose exploiting LLMs' ICL ability to solve NER tasks using zero and few-shot learning without *corpora* annotation and model fine-tuning. Zero and few-shot learning NER could be helpful when no annotated *corpora* are available; no base models are available to fine-tune; when one needs to prototype a NER-based system before spending time and money annotating *corpora* and training models; or as a pre-annotation step to aid human annotation processes.

To leverage the ICL ability of LLMs, one must prompt the model using instructions in natural language to demonstrate the context of the task to the model. In the same way, the model will return a text completion also in natural language [12]. The problem with this pipeline is that different prompts can lead to different results. Also, one must ensure a consistent output to parse the response into a machine-processable output. Ensuring a consistent output and the parsing of the completion can be wrapped into a high-level programming interface.

We propose `llmNER`, a Python library that wraps the prompting and completion parsing process into an easy-to-use interface. Our library implements multiple prompting methods, multiple answer shape parsers, a flexi-

ble prompt templating for adapting to other languages or domains, and the possibility to perform zero-shot and few-shot learning. This tool facilitates the process of prompt engineering for solving NER tasks.

There are many ways to prompt LLMs, and each method has a different impact on their performance and ability to solve specific tasks. For example, to perform machine translation, we must design the prompt to instruct the model to complete the input text with the explicit translation. In this case, we do not need to post-process the completion because the completion is the prediction itself. On the other hand, to perform text classification, we must design the prompt to make the model complete the text with a single word corresponding to an instance of the label space [12]. Even though prompting methods for some NLP tasks, such as text classification [13, 14, 15, 16] or generation [6, 17, 18], prompting methods for NER still need to be explored. Xie et al. [19] explored multiple prompting methods and only a JSON formatted answer shape for solving zero-shot NER tasks, where the most relevant are the following: a vanilla strategy, which we will call the **single-turn prompting method**, where they directly ask the model to annotate the entity mentions; and a decomposed-question-answering strategy, which we will call **multi-turn prompting method** where they broke down the task into simpler subproblems by asking the model to annotate the text one entity at a time and syntactic augmentation when they first ask the model to analyze the syntactic structure of the text and then perform NER. Cui et al. [20] used a different multi-turn prompting method, where the model scores candidate text spans based on how well they fit alongside the entity classes. Finally, Wang et al. [21] used a single-turn prompting method, prompting the model to echo the input document with in-line annotations where special symbols enclosed the entity mentioned inside the input text.

NER remains a tough challenge for ICL, falling far behind [22, 19, 20, 21] state-of-the-art fine-tuned models [23], but shows capabilities for universal information extraction without training data [24]; therefore, a necessity arises to develop research to improve the results of this paradigm.

Multiple domain-specific applications of NER through ICL have been investigated; for example, in the biomedical and clinical domains [25], this paradigm has been used to extract information for cancer research [26], rare diseases [27], and adverse drug events [28]; in the legal domain to extract contract [29] and legal violation [30] information; in the social media domain to detect offensive text [31] and in human resources domain to extract skill information [32].

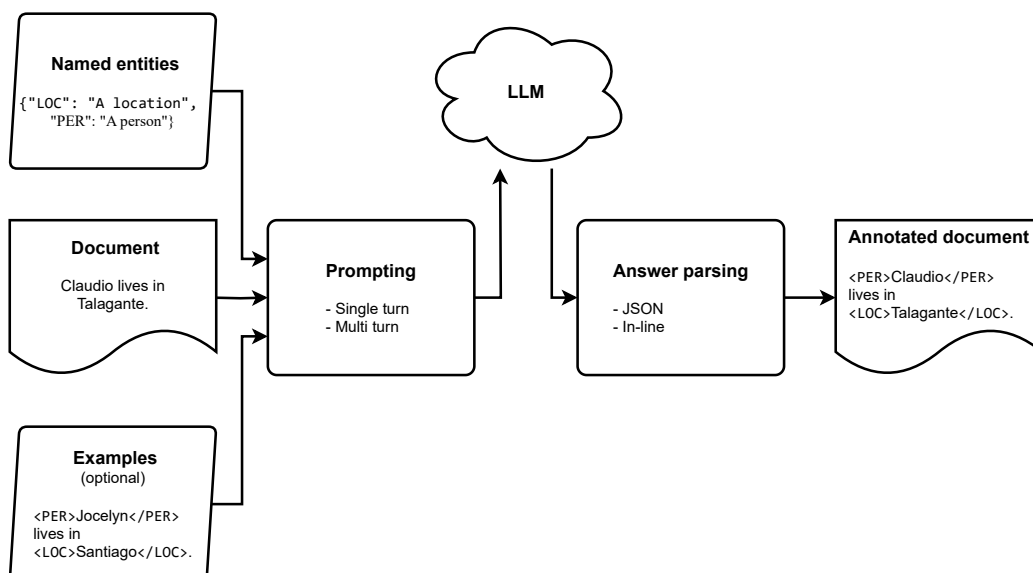
2. Software description

11mNER streamlines the prompting and answers parsing steps for NER through zero-shot and few-shot learning; in its most basic usage, the user only defines in natural language the entities to recognize, and then the library takes the model completion and parses the answer into a machine-readable abstract object. This tool also implements multiple prompting methods, answer shape parsers, and POS augmentation to streamline the prompt engineering process. The library's interface is based on the well-known machine learning library `scikit-learn` [33] to ease the learning process of using 11mNER. As the backend LLM, the library can use any model platform that exposes an OpenAI-compatible API endpoint.

2.1. Software architecture

The library components are divided into the prompting methods and the answer parsers. The named entity definitions and few-shot examples provided by the user are compiled as a text prompt using a prompting method. Then the completion performed by the LLM is parsed using an answer parser, returning an object containing the named entity mentions. A diagram of this architecture is in Figure 1.

Figure 1: Diagram of the architecture of the library and the pipeline for Named Entity Recognition.



To assist the user in the prompt engineering phase of a zero-shot or few-shot learning NER pipeline, the library implements the following prompting methods:

Single-turn In this prompting method, the model is asked to annotate the entity mentions in a single response as in this summarized prompt: *Annotate the mentions of the entities location and person: {response}*.

Multi-turn In this prompting method, the model is asked to annotate one entity at a time [19], exploiting the chain-of-thought following ability of LLMs, as in this summarized prompt: *Annotate the mentions of the entity location: {response}. And now annotate the mentions of the entity person: {response}*.

Step-by-step annotation In this option, the entity mentions are stored in a set at each turn, and at the end, this set is used as the final annotations. This option is useful when a specific text span can be associated with multiple entities. The user can also define the delimiters [21] that enclose the entity mentions at each turn.

Final-step annotation In this option, the model is prompted to annotate all the entities after the last turn, as in the single-turn prompting method.

Also, the library implements the following answer parsers:

In-line In this answer shape, the model is prompted to echo the exact input text, but with the entity mentions enclosed with in-line tags named after each entity class as in `<person>Fei-Fei Li</person> is a female scientist born in <country>China</country>`, where *Fei-Fei Li* is an instance of the Person entity class, and *China* is an instance of the Country entity class.

Custom delimiters The in-line tags can be defined by the user. For example, the user can define `@@` as an opening tag and `##` as a closing tag for an entity. This mode can be used only in a multi-turn setting.

JSON In this answer shape, the model is prompted to compose the annotations as a JSON object, where each key is the entity name and each value is a list of text mentions, as in `{"person":["Fei-Fei Li"], "country":["China"]}`.

Finally, the user can augment the prompt with part-of-speech (POS) data. This syntactic information, by default, is added through ICL by prompting the model to add the data.

2.2. Software functionalities

11mNER specialized in NER using the power of LLMs. Its main functionalities are:

- Zero-shot learning for NER: the user can input a document (or multiple documents) and a list of entities to detect with their descriptions. 11mNER will generate a prompt to query a selected LLM and parse the response to extract recognized entities.
- Few-shot learning for NER: Similar to the previous functionality, but the user can also input examples of NER to give a better context to the LLM when querying it.
- There are multiple options for querying and parsing, such as single or multi-turn prompting and in-line or JSON parsers, as explained in the previous section.
- POS tags can augment the context for prompting a selected LLM. These tags can be generated by a function (from an ad-hoc library such as NLTK or Spacy) or by a prompt to a LLM.
- When the user passes multiple documents, the library can parallelize the querying process to get results faster.

2.3. Sample code snippets analysis

To display the ease of use of 11mNER, Figure 2 and 3 show examples for zero-shot and few-shot cases, respectively.

3. Illustrative examples

To demonstrate the library’s functionality and usefulness, we assessed its performance on two public datasets in English and Spanish: CoNLL 2003 and CoNLL 2002. We measured the performance in all four entity classes: LOC (Location), MISC (Miscellaneous), ORG (organization), and PER (Person). Currently, the state-of-the-art (SOTA) for both benchmarks is a fine-tuned method proposed by Wang et al. [23].

We solved both benchmarks with zero-shot and few-shot paradigms using multiple LLMs with all of the hyperparameter combinations available in our library and compared the performance of each combination. We measured the relaxed F_1 score by each entity class in the benchmark.

Tables 2 and 3 present the results of the CoNLL 2003 and CoNLL 2002 NER tasks, utilizing the models and some hyperparameter combinations exposed by the library.

```
# INPUT:

import os

os.environ["OPENAI_API_KEY"] = "<your OpenAI api key>"
from llmner import ZeroShotNer

entities = {
    "person": "A person name, it can include first and last names
    , for example: John, Fabian or Mark Miranda",
    "organization": "An organization name, it can be a
    company, a government agency, etc.",
    "location": "A location name, it can be a city, a country
    , etc.",}

model = ZeroShotNer()
model.contextualize(entities=entities)

model.predict(["Fei-Fei Li is a female scientist born in
    China."])

# OUTPUT:

AnnotatedDocument(
    text="Fei-Fei Li is a female scientist born in China.
    ",
    annotations={
        Annotation(start=0, end=10, label="person"),
        Annotation(start=41, end=46, label="location"),
    },
)
```

Figure 2: Example usage of llmNER for zero-shot NER.

```

# INPUT:

examples = [
    AnnotatedDocument(
        text="Elon Musk is the owner of the US company Tesla"
        ,
        annotations={
            Annotation(start=30, end=32, label="location"),
            Annotation(start=0, end=9, label="person"),
            Annotation(start=41, end=46, label="organization"
            ),
        },
    ),
    AnnotatedDocument(
        text="Bill Gates is the owner of Microsoft",
        annotations={
            Annotation(start=0, end=10, label="person"),
            Annotation(start=27, end=36, label="organization"
            ),
        },
    ),
]

model = FewShotNer()
model.contextualize(entities=entities, examples=examples)
model.predict(["Pedro Pereira is the president of Peru and
the owner of Walmart."])

# OUTPUT:

    AnnotatedDocument(
        text="Pedro Pereira is the president of Peru and the
owner of Walmart.",
        annotations={
            Annotation(start=0, end=13, label="person"),
            Annotation(start=34, end=39, label="location"),
            Annotation(start=56, end=63, label="organization"
            ),
        },
    ),

```

Figure 3: Example usage of 11mNER for few-shot NER.

Both tables show the superior performance of GPT 3.5 across both datasets, as this model is the most tested across literature and achieves state-of-the-art in multiple tasks [34]. As expected, incorporating few-shot examples into the prompts led to a general performance enhancement because it is known that LLM’s performance scales with the number of examples given in the prompt [35, 6].

Regarding the prompting method we used. For the English dataset, the best-performing method was the multi-turn for zero-shot learning, a result that Xie et al. [19] also reported. The best prompting method for few-shot learning was the single-turn, but only by a small margin. This situation can be attributed to the fact that the most decisive contribution to the model’s decision is the few-shot examples, not the prompting method. For the Spanish dataset, the best performing prompting strategy for zero-shot learning was single-turn by a small margin and multi-turn for few-shot learning. This result shows that the performance may vary across datasets, and the user should perform hyperparameter optimization for their experiments.

Considering the answer shape, our proposed in-line answer shape generally outperforms JSON; this is important because almost all the literature shapes the answers as JSON and does not explore multiple answer shapes in the prompt engineering phase.

In general, the models performed better for the English NER benchmark (CoNLL 2003); most of the text on the LLM’s training *corpora* are in English; therefore, it is easier for those models to solve tasks in the English language. Results show that the SOTA performance on the respective benchmarks, represented by the ACE model [23], outperforms all our results by a significant margin. Despite our comparative standing being considerably lower than the SOTA, we exploit a different paradigm. The SOTA model exploits the model fine-tuning paradigm, which uses a task-specific training set to adapt the model parameters for the specific downstream task, and we are using a general-knowledge model, only fine-tuned for instruction following to solve a downstream task by only describing the task in natural language; hence both results cannot be easily compared. The main contribution of our present work and our main goal is not the performance of the strategies but the release of a library to ease the research of NER through in-context learning.

4. Impact

11mNER’s main characteristic is its ease of use for performing zero-shot and few-shot NER alongside its flexible interface for performing prompt engineering. Regarding those strengths, our library can help push the boundaries of

Model		LOC	MISC	ORG	PER	μF_1		
Zero-shot	Single-turn	GPT 3.5	67.0	29.0	51.9	79.9	62.9	
		JSON Mixtral 8x7B	44.0	7.3	29.0	55.2	36.8	
		Llama 2 70B	38.6	4.5	37.1	58.4	33.5	
	In-line	GPT 3.5	62.1	14.6	59.1	89.3	59.8	
		Mixtral 8x7B	52.3	4.8	34.6	73.6	39.2	
		Llama 2 70B	28.3	2.6	18.6	28.5	23.5	
	Multi-turn	JSON	GPT 3.5	74.3	3.1	41.1	90.1	68.8
			Mixtral 8x7B	47.0	5.8	37.1	56.4	41.0
			Llama 2 70B	41.6	4.8	27.6	77.1	46.0
		In-line	GPT 3.5	76.1	28.3	52.4	81.2	72.9
			Mixtral 8x7B	43.6	2.7	17.8	57.2	36.0
			Llama 2 70B	35.2	1.8	16.3	41.7	31.2
Few-shot	Single-turn	GPT 3.5	65.7	34.2	55.6	80.6	62.6	
		JSON Mixtral 8x7B	52.7	15.2	36.0	77.8	46.9	
		Llama 2 70B	49.6	13.7	31.5	79.6	48.3	
	In-line	GPT 3.5	75.5	48.9	62.3	88.8	73.5	
		Mixtral 8x7B	59.4	20.1	51.2	78.6	54.6	
		Llama 2 70B	54.5	19.0	38.8	62.8	48.1	
	Multi-turn	JSON	GPT 3.5	60.3	17.4	59.4	75.3	58.0
			Mixtral 8x7B	47.2	15.8	31.3	79.7	54.6
			Llama 2 70B	59.3	7.8	48.1	86.9	57.0
		In-line	GPT 3.5	67.8	17.6	51.3	88.2	72.1
			Mixtral 8x7B	40.2	2.0	12.6	36.8	30.4
			Llama 2 70B	25.8	3.5	16.9	30.6	20.9

Table 2: 11mNER F_1 performance results (in percentage) over CoNLL 2003 (English NER task). μF_1 column represents micro-averaged F_1 over all entity classes.

Model		LOC	MISC	ORG	PER	μF_1		
Zero-shot	Single-turn	GPT 3.5	73.4	8.3	56.2	85.1	64.0	
		JSON Mixtral 8x7B	50.9	7.8	49.2	59.2	44.9	
		Llama 2 70B	56.3	6.4	38.1	74.8	39.6	
	In-line	GPT 3.5	68.9	10.8	56.8	90.3	57.2	
		Mixtral 8x7B	47.7	4.3	36.3	68.2	35.4	
		Llama 2 70B	37.4	1.7	10.3	15.7	20.0	
	Multi-turn	JSON	GPT 3.5	67.6	6.6	59.1	75.3	55.4
			Mixtral 8x7B	47.4	0.0	39.5	56.7	42.0
			Llama 2 70B	52.6	6.3	54.0	64.2	39.9
		In-line	GPT 3.5	60.6	20.7	66.7	78.5	63.7
			Mixtral 8x7B	50.2	4.4	52.8	58.8	39.6
			Llama 2 70B	27.5	3.7	22.9	33.1	23.0
Few-shot	Single-turn	GPT 3.5	64.6	19.8	52.1	85.1	59.0	
		JSON Mixtral 8x7B	60.3	14.8	38.6	78.3	47.3	
		Llama 2 70B	66.9	12.7	23.5	77.1	41.1	
	In-line	GPT 3.5	64.5	15.4	54.4	87.2	61.4	
		Mixtral 8x7B	75.5	15.6	63.0	80.2	56.5	
		Llama 2 70B	52.0	17.5	50.5	62.9	49.1	
	Multi-turn	JSON	GPT 3.5	71.6	18.2	61.3	88.7	58.5
			Mixtral 8x7B	54.7	13.1	41.7	79.8	45.1
			Llama 2 70B	54.3	9.2	33.2	60.0	34.5
		In-line	GPT 3.5	67.8	15.2	72.2	92.6	68.8
			Mixtral 8x7B	38.9	4.0	25.3	48.0	29.3
			Llama 2 70B	43.2	7.5	0.0	53.7	25.6

Table 3: 11mNER F_1 performance results (in percentage) over CoNLL 2002 (Spanish NER task). μF_1 column represents micro-averaged F_1 over all entity classes.

ICL research by removing the barrier of the prompting and parsing steps. Without this barrier, ICL can be easily used in applied machine learning, such as LLM research on clinical NLP, improving the research in areas adjacent to machine learning and NLP.

Our library has been used on clinical NLP as a pre-annotation step for streamlining human annotations of corpora for privacy-preserving clinical NLP [36, 37].

5. Conclusions

The developed Python library offers an accessible framework for conducting zero and few-shot NER, validated in Spanish and English using CoNLL 2003 and CoNLL 2002 benchmarks. This comprehensive tool encompasses various prompting methods, answer shape parsers, and POS augmentation techniques, allowing researchers to explore diverse prompt engineering configurations.

Our results are not comparable with the SOTA because our approach diverges significantly from conventional NER methodologies. Rather than relying on task-specific fine-tuning, we leverage a general-knowledge model, fine-tuned solely for instruction following. Consequently, direct comparison with SOTA results is challenging due to the inherent paradigm differences.

The release of the library to the research community aims to foster reproducibility and innovation in NLP research, enabling researchers to explore and validate various settings, thus advancing the collective understanding of NER through in-context learning.

Acknowledgements

This work was funded by ANID Chile National Doctoral Scholarships 21211659 (CA) and 21220200 (FV).

References

- [1] X. Qiu, T. Sun, Y. Xu, Y. Shao, N. Dai, X. Huang, Pre-trained models for natural language processing: A survey, *Science China Technological Sciences* 63 (2020) 1872–1897. URL: <https://doi.org/10.1007/s11431-020-1647-3>. doi:10.1007/s11431-020-1647-3.
- [2] J. Dodge, G. Ilharco, R. Schwartz, A. Farhadi, H. Hajishirzi, N. Smith, Fine-tuning pretrained language models: Weight initializations, data orders, and early stopping, 2020. [arXiv:2002.06305](https://arxiv.org/abs/2002.06305).

- [3] J. Devlin, M.-W. Chang, K. Lee, K. Toutanova, BERT: Pre-training of deep bidirectional transformers for language understanding, in: Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers), Association for Computational Linguistics, Minneapolis, Minnesota, 2019, pp. 4171–4186. URL: <https://aclanthology.org/N19-1423>. doi:10.18653/v1/N19-1423.
- [4] P. He, X. Liu, J. Gao, W. Chen, Deberta: Decoding-enhanced bert with disentangled attention, in: International Conference on Learning Representations, 2021. URL: <https://openreview.net/forum?id=XPZiaotutsD>.
- [5] W. X. Zhao, K. Zhou, J. Li, T. Tang, X. Wang, Y. Hou, Y. Min, B. Zhang, J. Zhang, Z. Dong, Y. Du, C. Yang, Y. Chen, Z. Chen, J. Jiang, R. Ren, Y. Li, X. Tang, Z. Liu, P. Liu, J.-Y. Nie, J.-R. Wen, A survey of large language models, arXiv preprint arXiv:2303.18223 (2023). URL: <http://arxiv.org/abs/2303.18223>.
- [6] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, S. Agarwal, A. Herbert-Voss, G. Krueger, T. Henighan, R. Child, A. Ramesh, D. Ziegler, J. Wu, C. Winter, C. Hesse, M. Chen, E. Sigler, M. Litwin, S. Gray, B. Chess, J. Clark, C. Berner, S. McCandlish, A. Radford, I. Sutskever, D. Amodei, Language models are few-shot learners, in: H. Larochelle, M. Ranzato, R. Hadsell, M. Balcan, H. Lin (Eds.), Advances in Neural Information Processing Systems, volume 33, Curran Associates, Inc., 2020, pp. 1877–1901. URL: https://proceedings.neurips.cc/paper_files/paper/2020/file/1457c0d6bfc4967418bfb8ac142f64a-Paper.pdf.
- [7] J. Kaplan, S. McCandlish, T. Henighan, T. B. Brown, B. Chess, R. Child, S. Gray, A. Radford, J. Wu, D. Amodei, Scaling laws for neural language models, CoRR abs/2001.08361 (2020). URL: <https://arxiv.org/abs/2001.08361>. arXiv:2001.08361.
- [8] J. Wei, Y. Tay, R. Bommasani, C. Raffel, B. Zoph, S. Borgeaud, D. Yogatama, M. Bosma, D. Zhou, D. Metzler, E. H. Chi, T. Hashimoto, O. Vinyals, P. Liang, J. Dean, W. Fedus, Emergent abilities of large language models, Transactions on Machine Learning Research (2022). URL: <https://openreview.net/forum?id=yzkSU5zdWd>, survey Certification.

- [9] Q. Dong, L. Li, D. Dai, C. Zheng, Z. Wu, B. Chang, X. Sun, J. Xu, L. Li, Z. Sui, A survey on in-context learning, 2023. [arXiv:2301.00234](https://arxiv.org/abs/2301.00234).
- [10] Y. Wang, Q. Yao, J. T. Kwok, L. M. Ni, Generalizing from a few examples: A survey on few-shot learning, *ACM Comput. Surv.* 53 (2020). URL: <https://doi.org/10.1145/3386252>. doi:10.1145/3386252.
- [11] D. Jurafsky, J. H. Martin, *Speech and Language Processing: An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition*, 3rd ed. draft ed., 2023. URL: <https://web.stanford.edu/~jurafsky/slp3/>.
- [12] P. Liu, W. Yuan, J. Fu, Z. Jiang, H. Hayashi, G. Neubig, Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing, *ACM Computing Surveys* 55 (2023) 1–35. URL: <https://doi.org/10.1145/3560815>. doi:10.1145/3560815.
- [13] W. Yin, J. Hay, D. Roth, Benchmarking zero-shot text classification: Datasets, evaluation and entailment approach, in: K. Inui, J. Jiang, V. Ng, X. Wan (Eds.), *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, Association for Computational Linguistics, Hong Kong, China, 2019, pp. 3914–3923. URL: <https://aclanthology.org/D19-1404>. doi:10.18653/v1/D19-1404.
- [14] T. Gao, A. Fisch, D. Chen, Making pre-trained language models better few-shot learners, in: C. Zong, F. Xia, W. Li, R. Navigli (Eds.), *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, Association for Computational Linguistics, Online, 2021, pp. 3816–3830. URL: <https://aclanthology.org/2021.acl-long.295>. doi:10.18653/v1/2021.acl-long.295.
- [15] K. Hambardzumyan, H. Khachatrian, J. May, WARP: Word-level Adversarial ReProgramming, in: C. Zong, F. Xia, W. Li, R. Navigli (Eds.), *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, Association for Computational Linguistics, Online, 2021, pp. 4921–4933. URL: <https://aclanthology.org/2021.acl-long.381>. doi:10.18653/v1/2021.acl-long.381.

- [16] B. Lester, R. Al-Rfou, N. Constant, The power of scale for parameter-efficient prompt tuning, in: M.-F. Moens, X. Huang, L. Specia, S. W.-t. Yih (Eds.), Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, Online and Punta Cana, Dominican Republic, 2021, pp. 3045–3059. URL: <https://aclanthology.org/2021.emnlp-main.243>. doi:10.18653/v1/2021.emnlp-main.243.
- [17] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, I. Sutskever, Language models are unsupervised multitask learners, 2019. URL: <https://api.semanticscholar.org/CorpusID:160025533>.
- [18] T. Schick, H. Schütze, Few-shot text generation with natural language instructions, in: M.-F. Moens, X. Huang, L. Specia, S. W.-t. Yih (Eds.), Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, Online and Punta Cana, Dominican Republic, 2021, pp. 390–402. URL: <https://aclanthology.org/2021.emnlp-main.32>. doi:10.18653/v1/2021.emnlp-main.32.
- [19] T. Xie, Q. Li, J. Zhang, Y. Zhang, Z. Liu, H. Wang, Empirical study of zero-shot ner with chatgpt, 2023. arXiv:2310.10035.
- [20] L. Cui, Y. Wu, J. Liu, S. Yang, Y. Zhang, Template-based named entity recognition using BART, in: C. Zong, F. Xia, W. Li, R. Navigli (Eds.), Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021, Association for Computational Linguistics, Online, 2021, pp. 1835–1845. URL: <https://aclanthology.org/2021.findings-acl.161>. doi:10.18653/v1/2021.findings-acl.161.
- [21] S. Wang, X. Sun, X. Li, R. Ouyang, F. Wu, T. Zhang, J. Li, G. Wang, Gpt-ner: Named entity recognition via large language models, 2023. arXiv:2304.10428.
- [22] C. Pang, Y. Cao, Q. Ding, P. Luo, Guideline learning for in-context information extraction, in: H. Bouamor, J. Pino, K. Bali (Eds.), Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, Singapore, 2023, pp. 15372–15389. URL: <https://aclanthology.org/2023.emnlp-main.950>. doi:10.18653/v1/2023.emnlp-main.950.
- [23] X. Wang, Y. Jiang, N. Bach, T. Wang, Z. Huang, F. Huang, K. Tu, Automated concatenation of embeddings for structured prediction, in:

- C. Zong, F. Xia, W. Li, R. Navigli (Eds.), Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), Association for Computational Linguistics, Online, 2021, pp. 2643–2660. URL: <https://aclanthology.org/2021.acl-long.206>. doi:10.18653/v1/2021.acl-long.206.
- [24] D. Dukić, J. Šnajder, Do not (always) look right: Investigating the capabilities of decoder-based large language models for sequence labeling, 2024. arXiv:2401.14556.
- [25] Y. Hu, Q. Chen, J. Du, X. Peng, V. K. Keloth, X. Zuo, Y. Zhou, Z. Li, X. Jiang, Z. Lu, K. Roberts, H. Xu, Improving large language models for clinical named entity recognition via prompt engineering, Journal of the American Medical Informatics Association (2024). URL: <http://dx.doi.org/10.1093/jamia/ocad259>. doi:10.1093/jamia/ocad259.
- [26] X. Jiang, K. Zhang, T. Huang, B. Malin, T. Osterman, Q. Long, mcodegpt: Enhancing cancer research through zero-shot information extraction from clinical free text data (2024). URL: <http://dx.doi.org/10.21203/rs.3.rs-3882757/v1>. doi:10.21203/rs.3.rs-3882757/v1.
- [27] C. Shyr, Y. Hu, L. Bastarache, A. Cheng, R. Hamid, P. Harris, H. Xu, Identifying and extracting rare diseases and their phenotypes with large language models, Journal of Healthcare Informatics Research (2024). URL: <http://dx.doi.org/10.1007/s41666-023-00155-0>. doi:10.1007/s41666-023-00155-0.
- [28] S. Modi, K. Azhar Kasmiran, N. Mohd Sharef, M. Yunus Sharum, Extracting adverse drug events from clinical notes: A systematic review of approaches used, Journal of Biomedical Informatics (2024) 104603. URL: <http://dx.doi.org/10.1016/j.jbi.2024.104603>. doi:10.1016/j.jbi.2024.104603.
- [29] H. S. Adibhatla, P. Baswani, M. Shrivastava, Fine-grained contract ner using instruction based model, 2024. URL: <https://arxiv.org/abs/2401.13545>. doi:10.48550/ARXIV.2401.13545.
- [30] D. Bernsohn, G. Semo, Y. Vazana, G. Hayat, B. Hagag, J. Niklaus, R. Saha, K. Truskovskiy, Legallens: Leveraging llms for legal violation identification in unstructured text, 2024. arXiv:2402.04335.
- [31] A. Paraschiv, T. A. Ion, M. Dascalu, Offensive text span detection in romanian comments using large language models, Information 15

- (2023) 8. URL: <http://dx.doi.org/10.3390/info15010008>. doi:10.3390/info15010008.
- [32] K. C. Nguyen, M. Zhang, S. Montariol, A. Bosselut, Rethinking skill extraction in the job market domain using large language models, 2024. [arXiv:2402.03832](https://arxiv.org/abs/2402.03832).
- [33] L. Buitinck, G. Louppe, M. Blondel, F. Pedregosa, A. Mueller, O. Grisel, V. Niculae, P. Prettenhofer, A. Gramfort, J. Grobler, R. Layton, J. VanderPlas, A. Joly, B. Holt, G. Varoquaux, API design for machine learning software: experiences from the scikit-learn project, in: ECML PKDD Workshop: Languages for Data Mining and Machine Learning, 2013, pp. 108–122. URL: <https://arxiv.org/abs/1309.0238>.
- [34] J. Ye, X. Chen, N. Xu, C. Zu, Z. Shao, S. Liu, Y. Cui, Z. Zhou, C. Gong, Y. Shen, J. Zhou, S. Chen, T. Gui, Q. Zhang, X. Huang, A comprehensive capability analysis of gpt-3 and gpt-3.5 series models, 2023. [arXiv:2303.10420](https://arxiv.org/abs/2303.10420).
- [35] Y. Song, T. Wang, P. Cai, S. K. Mondal, J. P. Sahoo, A comprehensive survey of few-shot learning: Evolution, applications, challenges, and opportunities, *ACM Computing Surveys* 55 (2023) 1–40. URL: <http://dx.doi.org/10.1145/3582688>. doi:10.1145/3582688.
- [36] C. Aracena, L. Miranda, T. Vakili, F. Villena, T. Quiroga, F. Núñez-Torres, V. Rocco, J. Dunstan, A privacy-preserving corpus for occupational health in spanish: Evaluation for NER and classification tasks, in: The 6th Clinical Natural Language Processing Workshop, 2024. URL: <https://openreview.net/forum?id=GLsmr53Hwi>.
- [37] J. Dunstan, T. Vakili, L. Miranda, F. Villena, C. Aracena, P. Vera, S. V. Valenzuela, V. Rocco, A pseudonymized corpus of occupational health narratives for clinical entity recognition in spanish (2024). URL: <http://dx.doi.org/10.21203/rs.3.rs-3826527/v1>. doi:10.21203/rs.3.rs-3826527/v1.