

Evaluation of few-shot learning capabilities in polish language models

Tsimur Hadeliya Dariusz Kajtoch

ML Research at Allegro, Poznań, Poland

Introduction

Recent research reports that pre-trained language models can effectively solve natural language problems using only few examples. This approach called few-shot learning (FSL) and gained much popularity in recent years. However, vast majority of research in few-shot learning conducted exclusively for English, while other languages remains unexplored. To address this gap for polish language, we conducted experiments using two main approaches in FSL: In-context learning (ICL) and Parameter Efficient Fine-Tuning (PEFT). To get relevant and reliable results we stick to classification tasks during the experiments and construct few-shot classification benchmark based on publicly available datasets.

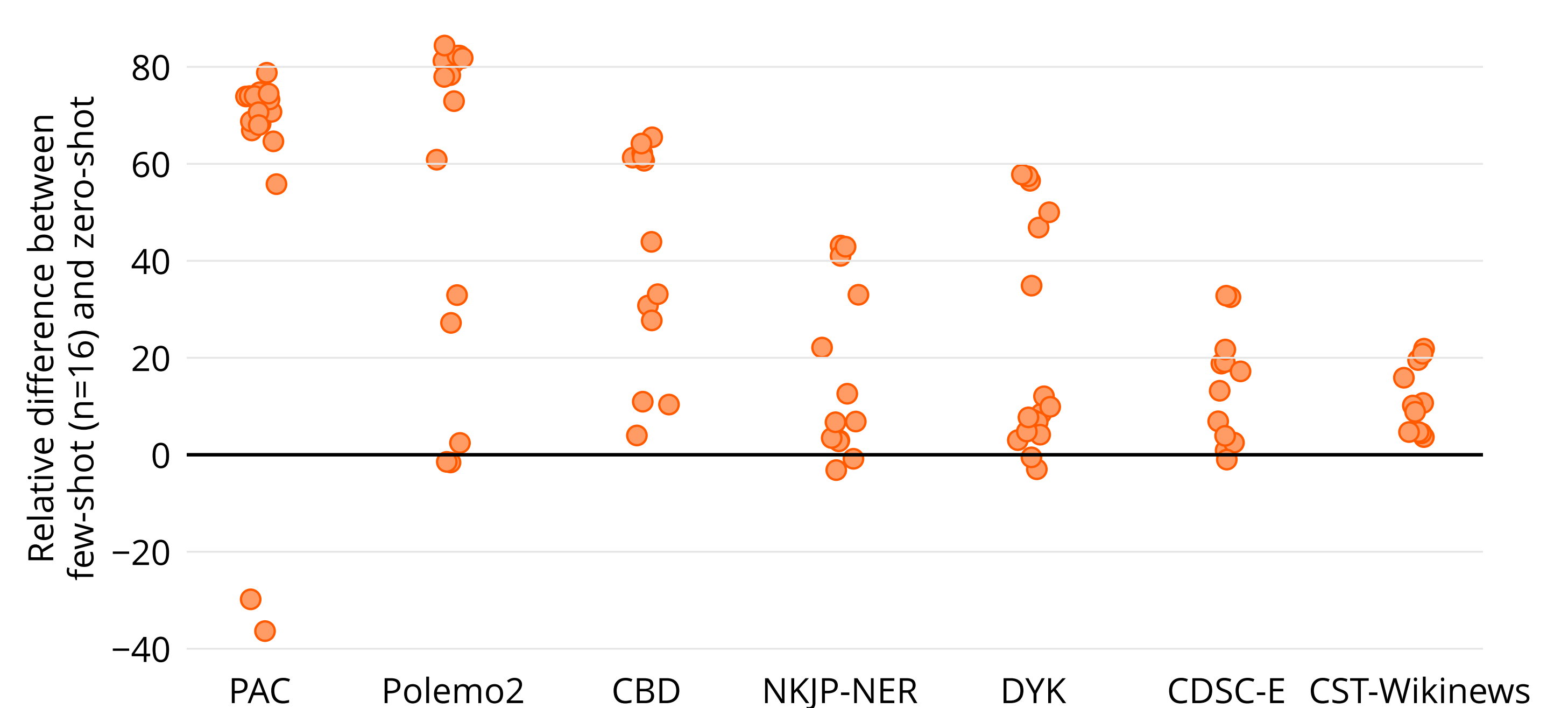
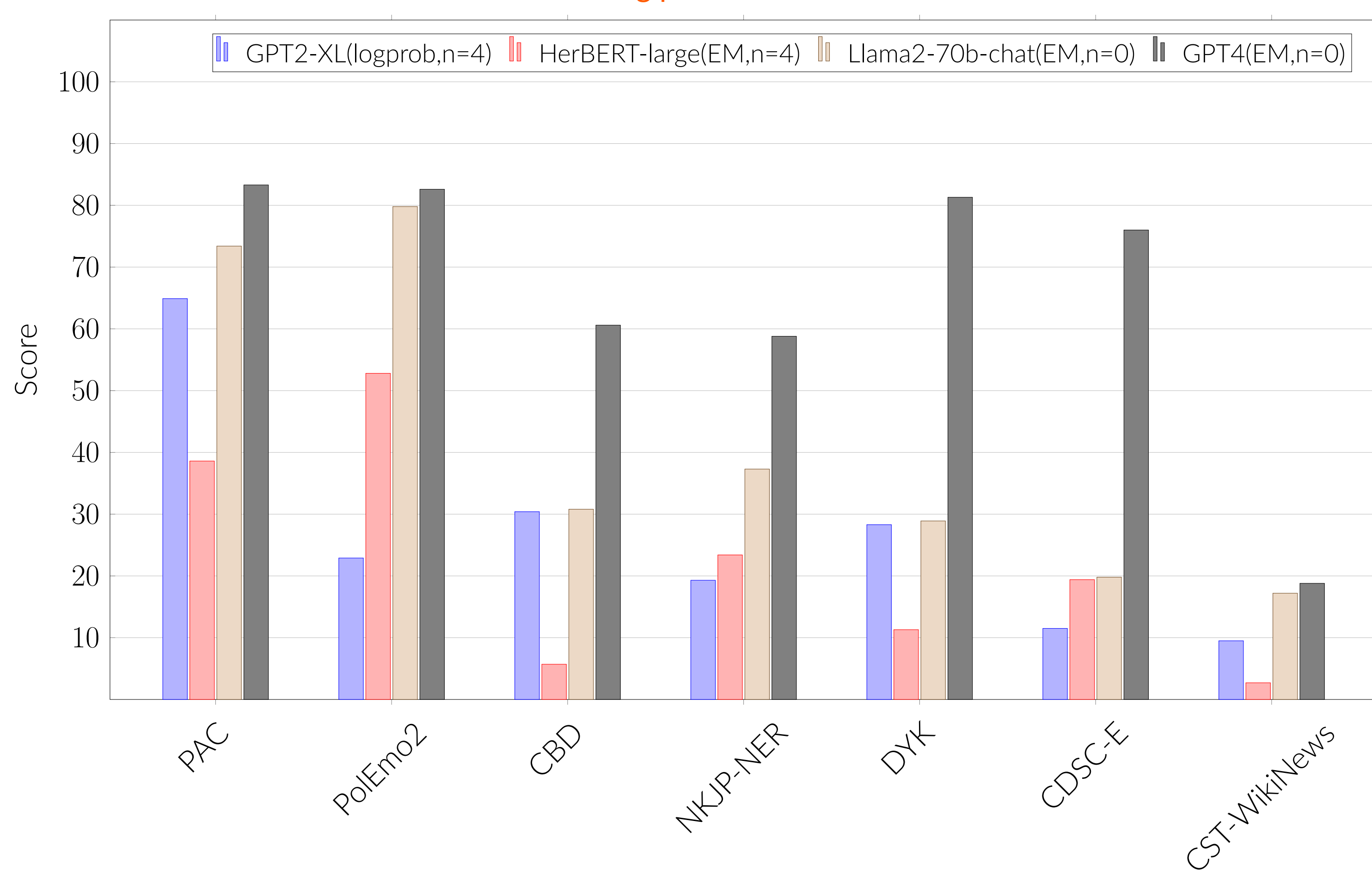
Benchmark

Name	Input	#Classes	Metrics	Avg. len	Domain
Lepiszczce					
PAC	text	2	F1-binary	185	legal texts
DYK	text pair	2	F1-binary	288	Wikipedia
CDSC-E	text pair	3	Accuracy	144	image captions
Polemo2	text	4	Accuracy	758	online reviews
KLEJ					
CBD	text	2	F1-binary	93	social media
NKJP-NER	text	6	Accuracy	85	national corpus
Other					
CST-Wikinews	text pair	12	Accuracy	232	Wikinews

Main findings

- **Few-shot approaches shows weaker performance compared to baseline models**
Except GPT-4, all models and methods shows much weaker performance compared to full fine-tuning and even the simplest baselines.
- **PEFT outperforms ICL almost for all models**
ICL methods is useful with LLM, but can't guarantee stable performance/
- **Polish LM from the box is not suitable to in-context learning methods**
The main problems are: short context window and high sensitivity for input format.
- **Instruction tuning of LLM transfers to Polish language**
Significant difference is observed between base and chat version of Llama-2-70b models.

In-context learning performance of selected models



Difference in performance between the few-shot setting (n = 16) and the zero-shot setting for the GPT-3.5. Results were collected for all our templates.

Model	AVG	PAC	Polemo2	CBD	NKJP-NER	DYK	CDSC-E	CST-Wikinews
Metric used		F1	Acc	F1	Acc	F1	Acc	Acc
Baseline								
Random Guessing	26.6 ± 0.4	57.5 ± 0.7	24.8 ± 1.0	20.9 ± 1.6	16.8 ± 0.8	25.5 ± 1.7	33.6 ± 0.7	6.9 ± 0.9
Most frequent	42.6	80.6	41.3	23.6	34.3	28.9	74.4	15.4
/Full FT/ HerBERT-large	79.9 ± 0.6	91.1 ± 0.0	90.9 ± 0.0	53.2 ± 3.2	94.0 ± 0.0	68.8 ± 2.1	93.4 ± 0.0	67.9 ± 1.0
PEFT								
/SF/ SBERT-large	47.1	68.8 ± 6.5	69.9 ± 10.7	44.4 ± 5.4	30.7 ± 6.5	27.7 ± 2.9	72.3 ± 5.7	16.2 ± 3.2
/SF/ HerBERT-large	44.1	70.7 ± 11.9	46.0 ± 11.2	42.5 ± 10.9	25.8 ± 11.6	40.6	67.2 ± 14.5	15.7 ± 1.9
/SF/ RoBERTa-large	47.6	66.8 ± 16.1	83.6 ± 2.7	44.5 ± 10.6	35.5 ± 4.4	26.1	62.2 ± 10.5	14.3 ± 2.5
/LP/ SBERT-large	43.4	67.5 ± 7.0	60.3 ± 4.6	40.1 ± 4.3	30.4 ± 5.9	27.2 ± 1.8	62.8 ± 5.8	15.3 ± 3.0
/FT/ SBERT-large	31.2	33.6 ± 31.3	47.0 ± 5.8	32.0 ± 30.7	28.1 ± 7.7	6.4 ± 10.3	61.0 ± 17.6	10.1 ± 4.3
/LP/ Ada	40.9	72.9 ± 6.3	55.2 ± 5.0	30.7 ± 4.2	29.1 ± 2.5	25.4 ± 4.1	58.2 ± 11.7	14.8 ± 3.9
/LP/ DaVinci	42.7	67.6 ± 7.0	58.9 ± 1.3	36.9 ± 9.7	30.1 ± 6.6	29.9 ± 2.7	60.1 ± 11.2	15.5 ± 3.5
/LP/ Gecko	37.7	61.8 ± 4.5	42.9 ± 6.4	23.7 ± 1.3	24.8 ± 4.3	25.0 ± 3.8	68.8 ± 6.3	17.0 ± 2.8
In-context learning								
GPT-3.5 (EM) (n=0)	55.4	82.2 ₀₁₄	81.6 ₀₀₅	50.0 ₀₄₆	44.9 ₀₀₁	53.1 ₀	62.9 ₀	13.3 ₀₀₃
GPT-3.5 (EM) (n=16)	59.5	73.9 ± 3.6	81.9 ± 2.1	64.1 ± 1.9	46.1 ± 2.9	64.1 ± 1.8	66.7 ± 7.6	19.8 ± 2.7
GPT-4 (EM) (n=0)	<u>65.9</u>	83.3 ± 0	82.6 ± 0.009	60.6 ± 0.028	58.8 ± 0.004	81.3 ± 0.002	76.0 ± 0	18.8 ± 0.003
Llama-2-70b-chat (EM) (n=0)	41.0	73.4 ± 0.090	79.8 ± 0.002	30.8 ± 0.122	37.3 ± 0.087	28.9 ± 0	19.8 ± 0	17.2 ± 0.016
Llama-2-70b (EM) (n=0)	14.6	41.8 ± 0.576	11.3 ± 0.678	0.3 ± 0.662	21.4 ± 0.324	19.2 ± 0.055	6.6 ± 0	1.6 ± 0.859
Bison-text (EM) (n=0)	52.2	80.2 ± 0.006	80.7 ± 0.009	42.6 ± 0.077	47.5 ± 0.027	61.6 ± 0.016	35.0 ± 0.001	17.7 ± 0.003
Bison-text (EM) (n=16)	-	83.7 ± 1.4	81.8 ± 2.8	-	45.7 ± 3.6	76.6 ± 0.7	66.4 ± 7.3	19.5 ± 2.2
Krakowiak-7b (EM) (n=0)	20.5	38.0 ± 0.624	28.4 ± 0.056	0.5 ± 0.687	23.9 ± 0.002	24.4 ± 0.03	19.0 ± 0.002	9.1 ± 0.109
GPT-2-xl (EM) (n=0)	15.0	23.8 ± 0.334	20.0 ± 0.274	16.9 ± 0.293	21.3 ± 0.96	14.9 ± 0.332	8.1 ± 0.144	0.0 ± 0.0
GPT-2-xl (EM) (n=4)	28.9	65.1 ± 0.107	32.5 ± 0.110	29.3 ± 0.109	16.1 ± 0.33	35.6 ± 0.286	12.1 ± 0.76	11.6 ± 0.76
GPT-2-xl (EM) (n=16)	-	68.4 ± 0.49	-	19.9 ± 0.87	22.1 ± 0.66	23.1 ± 0.228	10.2 ± 0.50	11.5 ± 0.66
GPT-2-xl (logprob) (n=0)	20.2	45.6 ± 0.270	20.3 ± 0.58	3.6 ± 0.47	24.9 ± 0.16	29.0 ± 0.182	11.9 ± 0.51	6.2 ± 0.53
GPT-2-xl (logprob) (n=4)	26.7	64.9 ± 0.49	22.9 ± 0.43	30.4 ± 0.133	19.3 ± 0.100	28.3 ± 0.30	11.5 ± 0.68	9.5 ± 0.39
HerBERT-large (iter) (n=0)	11.2	20.0 _{44.7}	40.0 _{54.8}	0.0 ₀	5.4 _{12.1}	0.0 ₀	13.2 _{12.2}	0.0 ₀
HerBERT-large (iter) (n=4)	22.0	38.6 _{36.1}	52.8 _{7.5}	5.7 _{11.2}	23.4 _{2.6}	11.3 _{8.9}	19.4 _{1.0}	2.7 _{1.1}

Performance on test data
Underline numbers - best across method (PEFT or ICL); Bold numbers - best across two methods (PEFT and ICL)

Methodology

We use F1-binary and Accuracy metrics to follow approach used in KLEJ benchmark and have comparable results with models tested on this datasets.

To obtain reliable results, we conduct 5 experiments with different seeds and calculate mean with standard deviation (reported under-line).

Baseline

- **Random guessing** - Sample label from training data distribution.
- **Most frequent** - Use the most frequent label from train dataset as constant prediction. This method shows imbalance in datasets.

PEFT

- **/SF/** - fine-tuning with SetFit method.
- **/LP/** - Linear probing. Logistic regression on top of LM representations.
- **/FT/** - Head-based Fine-tuning.
- **/Full FT/** - Fine-tuning on all training data.

In-context learning

- **eval method:**
 - **EM** - Exact match. Check label in generated substring. If no label is matched, output special label and calculate as wrong prediction.
 - **logprob** - Calculate log probability of sequence with given label. Choose sequence with highest probability. Always choose one of the proposed labels.
 - **iter** - Iteratively add "[MASK]" token to generate label sequence. When generated EM approach is used.
- **(n=k)** - number of demonstrations (k) used in prompt.

References

- [1] Augustyniak et al., 2022, this is the way: designing and compiling lepszczce, a comprehensive nlp benchmark for polish.
- [2] Bach et al., 2022, PromptSource: An integrated development environment and repository for natural language prompts.
- [3] Radford et al., 2019, language models are unsupervised multitask learners.
- [4] Rybak et al., 2020, KLEJ: Comprehensive benchmark for Polish language understanding.
- [5] Touvron et al., 2023, llama 2: Open foundation and fine-tuned chat models.