

The PLLuM Instruction Corpus

Piotr Pezik¹, Filip Żarnecki¹, Konrad Kaczyński¹, Anna Cichosz¹, Zuzanna Deckert¹,
 Monika Garnys¹, Izabela Grabarczyk¹, Wojciech Janowski¹, Sylwia Karasińska¹,
 Aleksandra Kujawiak¹, Piotr Misztela¹, Maria Szymańska¹, Karolina Walkusz¹, Igor Siek¹,
 Maciej Chrabąszcz², Anna Kołos², Agnieszka Karlińska², Karolina Seweryn²,
 Aleksandra Krasnodębska², Paula Betscher², Zofia Cieślińska², Katarzyna Kowol²,
 Artur Wilczek², Maciej Trzciński², Katarzyna Dziewulska², Roman Roszko³,
 Tomasz Bernaś³, Jurgita Vaičenonienė³, Danuta Roszko³, Paweł Levchuk³, Paweł Kowalski³,
 Irena Prawdzic-Jankowska³, Marek Kozłowski⁴, Sławomir Dadas⁴, Rafał Poświata⁴,
 Alina Wróblewska⁵, Katarzyna Krasnowska-Kieraś⁵, Maciej Ogorodniczuk⁵, Michał Rudolf⁵,
 Piotr Rybak⁵, Karolina Saputa⁵, Joanna Wołoszyn⁵, Marcin Oleksy⁶, Bartłomiej Koptyra⁶,
 Teddy Ferdinan⁶, Stanisław Woźniak⁶, Maciej Piasecki⁶, Paweł Walkowiak⁶,
 Konrad Wojtasik⁶, Arkadiusz Janz⁶, Przemysław Kazienko⁶, Julia Moska⁶, Jan Kocon⁶

¹ University of Łódź

² NASK National Research Institute

³ Institute of Slavic Studies PAS

⁴ National Information Processing Institute

⁵ Institute of Computer Science PAS

⁶ Wrocław Tech

Correspondence: piotr.pezik@uni.lodz.pl

Abstract

This paper describes the instruction dataset used to fine-tune a set of transformer-based large language models (LLMs) developed in the PLLuM (Polish Large Language Model) project. We present a functional typology of the organic, converted, and synthetic instructions used in PLLuM and share some observations about the implications of using human-authored versus synthetic instruction datasets in the linguistic adaptation of base LLMs. Additionally, we release the first representative subset of the PLLuM instruction corpus (PLLuMIC), which we believe to be useful in guiding and planning the development of similar datasets for other LLMs.

1 Introduction

The Polish Large Language Model (PLLuM) was a project funded by the Polish Ministry of Digital Affairs in 2024¹. Its main delivery was a ‘family’ of language-adapted, fine-tuned, and aligned large language models (LLMs) ranging in size from 8 to 70 billion parameters as summarized in Table 1. One of the central tasks of the project was the design of an original corpus of instructions that could be

used to develop the basic interactive capabilities of the target models. Several challenges that became apparent in creating such a dataset formed the core motivation for this study.

First, the datasets used to fine-tune both proprietary and open-weight LLMs are usually withheld by their developers. This, in turn, makes the replication of LLM capabilities difficult. Second, the composition of stand-alone instruction datasets (i.e., datasets released independently of any specific LLM) is usually poorly documented. Such resources are typically constructed opportunistically, often as a conflation of other datasets, and even when they follow a predefined typology, the accompanying documentation is often too sparse to serve as a reliable foundation for designing similar datasets for LLM development.

Furthermore, while the role of instruction datasets in LLM development is generally acknowledged, there is relatively little published research on how different types of instructions impact the capabilities of original, published models. Another clear research gap is related to the growing trend of large-scale instruction distillation from so-called *strong LLMs*. While this approach offers a convenient shortcut to creating instruction datasets, it

¹See <http://pllm.org.pl>.

comes with unobvious limitations, especially in the context of linguistic and cultural adaptation of LLMs, which was central to the PLLuM project. Paradoxically, there is also the converse issue of human authors of instructions having to learn the style of LLM responses, which has recently emerged as a new register of language² as a result of human interactions with popular LLMs. Finally, deriving instructions from annotated corpora and structured knowledge sources (i.e., *automatic instructions*), while promising in many respects, introduces the risk of distorting the overall balance between organic, automatic, and synthetic instructions (see our definitions below) in the dataset and consequently also the behaviour of the resulting model.³

This paper presents the instruction datasets used to fine-tune the PLLuM models. We describe our functional typology of the organic, automatic, and synthetic instructions in the context of the LLM research issues signalled above. We also release a representative subset of the PLLuM instruction corpus (PLLuMIC) to provide potential guidance and inspiration for developing similar datasets for other LLMs.

2 LLM Fine-tuning

In the context of Large Language Model (LLM) development, the term *instructions* refers to single- or multi-turn question-and-answer pairs, i.e. (Q_i, A_i) , which exemplify the format, style, and functional content of interactions between the model and its users:

$$I = \{(Q_1, A_1), (Q_2, A_2), \dots, (Q_n, A_n)\} \quad (1)$$

Instructions can be categorized with respect to their origin as:

- *Organic*, i.e. authored by humans, including experts and trained annotators. Manual instructions can also be crowd-sourced or collected from human prompts in interactions with existing LLMs.
- *Converted*, i.e., derived from annotated corpora, knowledge sources, etc.
- *Synthetic*, i.e. distilled more or less directly from existing LLMs through manual or automated prompting techniques.

²We use the term register in the sense of a functional genre or variety of language (Conrad, 2023).

³Automatically converted instructions tend to be very repetitive. In large quantities, they may affect the conversational fluency of a fine-tuned LLM.

Hybrid scenarios for acquiring instructions are also possible in that synthetic instructions can be verified and corrected manually, organic and automatically converted instructions can be enhanced by LLMs, etc. As we explain below, each of these three basic sources of acquiring instructions has its advantages and limitations.

The fine-tuning of pre-trained models on instruction datasets remains a crucial step in the development of LLMs based on the transformer architecture. Although base models can perform certain tasks and interpolate between knowledge items attested in pre-training, it is clear that balanced, high-quality datasets of instructions are indispensable resources in text-to-text LLM development workflows (Longpre et al., 2023).

3 Availability and Transparency of Instruction Datasets

Although the basic steps of developing LLMs, such as fine-tuning on instructions, are widely researched, there is relatively little practical information about the composition of datasets used in real-world model-building projects. For various legal and business-related reasons, high-quality text corpora, instructions, and preferences are often withheld or inadequately documented by vendors and publishers of closed and open LLMs. We provide an overview of the transparency of instruction datasets used to develop a number open-weight models in Appendix A⁴. In short, the vast majority of such models are provided without the instructions used to fine-tune them and with very little if any documentation about such resources.

On the other hand, open instruction datasets (often developed independently of any particular LLM), tend to be largely opportunistic⁵. The definition and compilation of balanced instruction corpora remains a major methodological challenge for any team developing an original instruction fine-tuned LLM. Even in projects which utilize large-scale distillation of skills and knowledge from existing models, a general functional typology of human-LLM interactions is required to design the corpus of instructions.

⁴A useful distinction is made between open-source and open-weight models, where the latter are usually provided without key data resources.

⁵Appendix A contains a summary of open instruction datasets availability and documentation.

Model name	Base model
PLLuM-12B-nc-chat	Mistral-Nemo-Base-2407
PLLuM-8x7B-nc-chat	Mixtral-8x7B-v0.1
Llama-PLLuM-8B-chat	Llama-3.1-8B
Llama-PLLuM-70B-chat	Llama-3.1-70B

Table 1: A subset of the models adapted, fine-tuned, and aligned in PLLuM.

Category	Proportion
Knowledge (QA)	43%
Generation	25%
Extraction	6%
Programming	6%
Conversational	4%
NLP	3%
Adversarial	3%
Visualization	3%
Data manipulation	3%
Chain of Thought	2%
Translation	1%
Identity	1%

Table 2: High-level PLLuMIC composition with their respective approximate representation in the full organic component of the corpus.

4 The Composition of PLLuMIC

In the following section, we introduce the structure of the PLLuM Instruction Corpus (henceforth PLLuMIC). Although by design, the bulk of the corpus consists of (1) hand-crafted, high-quality organic instructions curated by a team of trained annotators, we also explored the value of (2) instructions distilled from existing LLMs and (3) converted from annotated corpora database and text repositories.

4.1 Organic Instructions

Instructions annotated by professional human annotators hired for the project formed the primary component of PLLuMIC. We refer to such instructions as *organic* to distinguish them from synthetic and automatic or ‘converted’ instructions. They were either written from scratch by single or multiple annotators to fill the above-mentioned categories or adapted from open datasets. Adaptation of open-source datasets (e.g. CREAK (Onoe et al., 2021) (3591 samples), ECQA (Aggarwal et al., 2021) (1033 samples), QED (Lamm et al., 2020) (1855 samples)) was an effective initial strategy, but this

approach showed significant limitations with time. Many of the adapted samples contained low-quality, simplistic, or erroneous instructions, but with some effort invested in their corrections, they proved to be useful for both fine-tuning and evaluation purposes.

The annotation process was subject to rigorous quality control, described in more detail in Appendix C.1.

Prompt-response Instructions We started the core manual annotation phase with a small ad-hoc typology covering mostly simple prompt-response interactions such as factual knowledge and commonsense reasoning question-answering and several generative subtypes, i.e. short text composition prompts with relatively long expected output. Some extractive tasks, such as summarization and keyphrase identification, were also considered in this initial phase. The resulting high-level composition of PLLuMIC is outlined in Table 2. A more detailed account of the PLLuMIC typology is given in Appendix D.

Dialogue instructions One of the stages in the development of PLLuMIC was the shift from simple prompt-response instructions to multi-turn dialogues. Although prompt-response turns are prototypical instructions, they fail to capture more sophisticated conversational scenarios such as role-playing, context-sensitivity, and multi-turn prompting, whereby several stages of interaction are required to specify and solve a task at hand. As the dataset size became sufficient to fine-tune early versions of our LLMs, instructions and multi-turn dialogues were also gathered through human-model interactions. The responses generated by intermediate fine-tuned models were carefully validated and refined before being included in the instruction dataset. We created a subset of over 3,500 dialogues with an average of approximately 12 turns per conversation.

A subset of our typology also features instructions in other languages, mostly Ukrainian, Lithuanian, Russian and Belarussian.

4.2 Synthetic Instructions

To extend the range of tasks and topical domains covered by human annotators, we generated an experimental subset of high-quality instructions using selected LLMs with limited human supervision. To this end, we gradually devised a map of topical domains (Appendix D.2.1 & D.2.2), and depending on the type of skill or knowledge, we used different multi-step generation-pipelines involving minimal human supervision and several locally-hosted LLMs. The main two types of synthetic instructions included in PLLuMIC were focused on *Knowledge distillation*, *RAG* and *Context-injected NLP* tasks.

4.2.1 Knowledge Distillation

Starting with a manually constructed list of topics and subtopics; for each topic, human annotators compiled a series of hand-written subject prompts that were subsequently injected into a meta-prompt generating a question; then LLM-generated questions were fed into a meta-prompt to generate the answer. Meta-prompts at each pipeline step contained detailed specifications of the desired content, style, and format. All prompt-answer pairs in this phase were generated and validated with the permissively licensed Mixtral8x22b-instruct model.

4.2.2 RAG Instructions

To optimize PLLuM models for Retrieval Augmented Generation (RAG), especially in the domain of public administration, we compiled a comprehensive set of instructions and preferences from documents available on Polish government websites in the *gov.pl* domain. These included mainly administrative guides, and structured informational pamphlets covering a range of issues such as applying for identity documents, business activity, taxes, residence registration, and others. We prepared three sets of questions: (1) *regular questions* that are likely to be answered by information contained in the indexed documents, (2) *adversarial questions* intended to trick the model into providing unacceptable answers and (3) *unrelated questions*, which were completely unrelated to the topic of the documents and therefore should be ignored. The regular and adversarial questions were generated by a strong LLM for each fragment of the document while unrelated questions were sampled from various QA datasets and reviewed by annotators. Afterwards, for each fragment, the top 5 documents were retrieved via a pipeline composed of the *bge-*

m3 retriever⁶ and the *bge-reranker-v2-m3* reranker⁶ (Chen et al., 2023a). For each set of questions and retrieved documents, we generated an answer using Llama-3.3-70B (serving as the strong LLM) and treated it as a reference answer. We also generated preferred answers to be used in the model alignment phase, using a weaker model Llama-3.1-8B. To avoid overfitting our generic models, we limited the set of RAG instructions to 5,000 in the SFT phase and 5,000 preferences in the alignment phase. The final training set contained 80% regular questions, 14% adversarial questions, and 6% unrelated questions.

4.2.3 Context-injected NLP

Text samples extracted from open-source collections were injected into system prompts containing detailed specifications of NLP tasks, such as named entity recognition, classification, semantic similarity, translation, etc. Special effort was invested in defining the desired structured output formats such as JSON, CSV, XML, etc. Pairs of system prompts and LLM-generated answers were subsequently validated for compliance with the constraints defined inside the system prompt.

4.2.4 Limitations of Mass-distillation

Despite the current trend to use large-scale distillation techniques in both LLM training and inference, we attempted to control and mitigate certain synthetic data limitations throughout the development of PLLuMIC. First, many LLMs are governed by licenses restricting data generation for derivative model development. These legal constraints and a lack of expertise in constructing original instruction and alignment datasets may lead to long-term over-dependence on existing LLMs. Second, a distillation of aligned models can propagate biases and pre-existing preferences, potentially compromising our model balancing and neutrality definitions. Furthermore, poorly controlled recursive distillation may lead to model degradation or even collapse (Shumailov et al., 2024). Finally, as discussed in Section 5.1, we observe significant negative transfer effects in language-adapted models while transfer learning and task interpolation are fundamental properties of generative language models.

4.3 Converted Instructions

Prompt-response pairs can also be automatically created from annotated corpora (e.g., treebanks,

⁶<https://github.com/FlagOpen/FlagEmbedding/tree/master>

named-entity datasets, etc.) and other resources, including machine-readable dictionaries and ontologies. Members of the PLLuM consortium used their experience in developing various types of NLP datasets to convert instances of these datasets into instructions. Responses were often extracted from annotation layers using handwritten question-and-answer templates. For some datasets,⁷ several prompt formats were prepared. When mapping an example into a single-turn instruction, a prompt format was randomly selected.

Notably, we only took train splits of these datasets for training, while validation splits were optionally used for internal evaluation.

Table 10 provides examples of specific subsets of converted instructions. Although this approach allows for the efficient large-scale production of instructions, the resulting data is often highly repetitive, reducing the fine-tuned model’s conversational versatility. To mitigate this, we imposed strict limits on the number of converted instructions obtained from each resource.⁸

5 Language Adaptation Experiments

Developing a carefully curated instruction corpus makes it possible to analyse how different instruction types’ diversity, quality, and quantity impact the performance of fine-tuned models. In this section, we demonstrate that fine-tuning organically sourced instructions enhances the model’s capabilities, particularly in areas where continued pre-training on textual data and mass distillation from *strong* seemingly multilingual LLMs achieve sub-optimal proficiency such as language and culture-specific forms of written communication.

5.1 Base Model Adaptation and Fine-tuning

One of the primary goals of the PLLuM project was to adapt existing base models (see Table 1), to better support understanding and generation of native Polish texts. This was partly achieved through (a) continued pre-training of the base models on a corpus of approx. 150 billion tokens, compiled from diverse textual sources, and (b) using a subset of this data to *anneal* the resulting model. At the same time, we observed that certain functional

⁷The text classification tasks from Polish Summaries Corpus, DYK, PolEmo2, Polish CBD, Polish Paraphrase Corpus, CDSC-E, 8tags, and NKJP-NER. More details can be found in Table 10.

⁸By default, only a maximum of 1,000 instructions were converted from a single resource.

types of texts, which may be particularly important for the intended use of the model, are either underrepresented in the raw pre-training data or, when included, are substandard in terms of style, grammar, and formatting.

For example, while the pre-training data included Polish e-mails and high-quality style guidelines for e-mail writing, the actual e-mail messages found in the raw corpus data often exhibited non-standard spelling and inconsistent punctuation. To illustrate, normative Polish e-mail style dictates that the second line of a message should begin with a lowercase letter if the first line contains an addressative form followed by a comma, as in:

Szanowny Panie,

*chciałbym uzyskać informację w sprawie
wymiany licznika energii...*

Although this differs from the English convention, where the second line is always capitalized, both capitalization patterns appear in naturally occurring Polish e-mails. Another subtle and frequently ignored prescriptive rule in Polish e-mail writing is the avoidance of a comma between complementary closings and newline signatures, as in:

Pozdrawiam

Jan Kowalski

Again, this differs from English-language e-mail conventions, where complementary closings are usually separated with a comma from signatures. Although early versions of our SFT models alternated between both conventions when prompted to produce e-mails, we found that a subset of less than 100 high-quality e-mail writing instructions was sufficient to imprint the above-mentioned (and several other) guidelines in the fine-tuned model. Based on our experience, the need for idiomatic handwritten instructions becomes particularly evident in adapting multi-lingual base models, which were pre-trained mostly on languages other than Polish. We found that special care is required when using fine-tuned LLMs for distilling language-specific instructions. The following is an example of a GPT-4 style Polish email message that illustrates the latter point:

Szanowny Panie Profesorze,
Mam nadzieję, że ten email zastanie
Pana w dobrym zdrowiu i nastroju.

Type	Quantity Train	Proportion Train	Quantity Total
Organic	38,106	49.12%	47,295
Converted	33,789	43.56%	33,789
Synthetic	5,679	7.32%	5,679

Table 3: Sources of instructions in PLLuM – Structure of training dataset and total quantity.

Chciałabym/chciałbym uprzejmie poprosić o możliwość umówienia się na krótką konsultację w najbliższą środę o godzinie 11:00. (...) Z góry dziękuję za poświęcony czas i rozwazenie mojej prośby. Z poważaniem, Twoje imię i nazwisko

While the email successfully fulfills the communicative goal of scheduling an appointment (as specified in the prompt), it also contains several instances of negative linguistic transfer from English. Beyond the formatting and punctuation violations mentioned earlier, the opening sentence directly translates a formulaic English email introduction (*I hope this message finds you in good health*), which sounds unidiomatic in Polish. This exemplifies a broader issue in LLM transfer learning: while models are designed to generalize across languages, their ability to transfer knowledge and skills can sometimes manifest as unintended stylistic interference. Transformer-based models tend to transfer stylistic conventions from languages best represented in the pre-training phase, leading to non-idiomatic outputs in the target language. This is similar to the human-like transfer of syntactic and pragmatic constructions from native or otherwise predominant language (Selinker, 1969).

5.2 Alignment

Model alignment on preference-based datasets, where chosen and rejected response pairs are annotated according to human preferences, aims to teach the model appropriate behaviours, particularly in responding to controversial and potentially harmful prompts. For the alignment of the PLLuM models, we used a dataset of over 40,000 manually annotated instructions, derived from three distinct annotation methodologies:

- Rating-based annotation – each response was assessed according to a predefined metric, with the higher-rated response designated as the preferred (chosen) response.

- Ranking-based annotation – four responses were ranked according to response quality.
- Dialog-based annotation – annotators engaged in multi-turn interactive conversations with the model, selecting the most appropriate responses.

The prompts were primarily created manually and did not overlap directly with PLLuMIC, although they were based on a similar typology, with a strong emphasis on safety-related prompts. Responses were generated by various open models, including the PLLuM ones. In cases where no response met the established criteria, annotators (over 50 different persons in total) provided their own responses.

Our experimental results indicate that, within the scope of this study, the most effective alignment method was the Odds Ratio Preference Optimization (ORPO) algorithm (Hong et al., 2024), which integrates alignment with instruction tuning through a specifically designed loss function. While previous research suggests that employing such a loss function obviates the need for SFT, our findings demonstrate that applying ORPO after SFT still yielded superior performance compared to alternative approaches such as KTO (Ethayarajh et al., 2024), DPO (Rafailov et al., 2024), and PPO (Schulman et al., 2017).

Through alignment training, model safety behaviours improved significantly, enabling our models to proactively address adversarial inputs and provide well-reasoned, defensible explanations, confirmed by red-teaming evaluation results (see 5.3). At the same time, we observed the well-documented trade-off between safety and helpfulness—a tendency for models to overly refuse to respond (Bai et al., 2022), even in the case of non-adversarial prompts (those that do not contain harmful content or encourage unsafe behaviour). While factuality and linguistic correctness remained mostly consistent with those achieved through SFT, verbosity increased, with models demonstrating a tendency to generate more elaborate responses, even in cases

where a more concise answer would have been sufficient.

At the same time, we also observed negative language transfer at this stage, stemming from the preference for responses generated by models trained predominantly on English-language data. Without additional linguistic adjustments, alignment on the preference-based dataset occasionally resulted in grammatical, lexical, and stylistic inconsistencies reflecting English-language rules, many of which had already been addressed during the SFT phase (e.g., punctuation in emails). This highlights the fact that for both SFT and alignment in low- and mid-resourced languages, manual human annotation, evaluation, and quality assurance are indispensable for maintaining proper language standards.

5.3 Evaluation

The linguistic adaptation of the base models listed in Table 1 was evaluated on the Polish Linguistic and Cultural Competency Benchmark (PLCC) (Dadas et al., 2025), which consists of 600 questions covering topics such as Polish history, geography, culture, tradition, art, entertainment, grammar, and vocabulary. The answers to the questions are assessed using an IFEval evaluation scheme (Zhou et al., 2023b). It is important to note that the PLLuM instruction corpus was developed independently from this benchmark. The base models were fine-tuned for 3 epochs using the AdamW optimizer (weight decay: 0.1) with a learning rate of 1e-5, a cosine scheduler (1% warmup), and a cumulative batch size of 128 on PLLuM instructions with a maximum sequence length of 16 384 tokens. The loss values were calculated only on the response turn of the instructions. The training was performed in a multi-node configuration⁹ using DeepSpeed ZeRO Stage 3 optimization.

Table 4 shows the PLCC scores obtained for the different stages of linguistic adaptation. Each group of evaluated models consists of:

1. The original reference instruction-following model, e.g. Mistral-Nemo-Instruct-2407,
2. The reference base model fine-tuned on PLLuMIC, e.g. Mistral-Nemo-2407+PLLuMIC,
3. A model continually pre-trained on Polish texts and fine-tuned on PLLuMIC, e.g.

PLLuM-12B-nc-instruct,

4. The continually pre-trained model, fine-tuned on PLLuMIC, aligned on PLLuM human preferences, e.g. PLLuM-12B-nc-chat.

Several conclusions emerge from this evaluation. Firstly, the continually pre-trained models consistently outperform their base counterparts across all four architectures. Secondly, fine-tuning on PLLuMIC is only effective for models that have undergone continual pretraining; otherwise, fine-tuning even degrades the model performance. In other words, fine-tuning on Polish instructions requires a model sufficiently primed on Polish data in the pre-training phase. To further examine the relationship between these two training phases, we have conducted additional ablation experiments, described in the Appendix ???. Finally, models aligned with human preferences achieve slightly higher benchmark scores than their instruction-fine-tuned predecessors. This might be partly because aligned models usually generate longer responses, which increases their chances of meeting the inclusion criteria of IFEval-style benchmarks. For example, if a model frequently paraphrases or summarizes parts of its responses, it is more likely to include words that match the benchmark criteria, thus improving its overall score.

The general knowledge capabilities (in contrast to the more cultural or linguistic competences) of the models fine-tuned on the PLLuM instruction corpus was also evaluated on the LLMzSzŁ benchmark (see Table 6), which is “a collection of Polish national exams, including both academic and professional tests extracted from the archives of the Polish Central Examination Board” (Jassem et al., 2025). Interestingly, the performance of our LLama-PLLuM-70B-chat model, which was fine-tuned on our instructions is only 2.71 points lower than the performance of LLama-3.3-70B-Instruct, which is reported to have been fine-tuned on millions of manually crafted instructions (AI@Meta, 2024).

Finally, the red-teaming evaluation results of our models are summarized in Table 5. The evaluation was conducted on 18,656 harmful prompts for the attack success rate (ASR) metric and 9,724 non-harmful samples for the false-refusal rate (FRR) metric (Krasnodębska et al., 2025). Both datasets cover 14 hazard categories defined by the Llama-Guard taxonomy (Inan et al., 2023). Additionally, they were generated using 10 different attack styles

⁹We used NVIDIA H100 nodes maintained by the Wrocław Centre for Networking and Supercomputing.

Model	PLCC ↑
Mistral-Nemo-Instruct-2407	23.00
Mistral-Nemo-2407+PLLMIC	22.33
PLLM-12B-nc-instruct	56.33
PLLM-12B-nc-chat	59.50
Mixtral-8x7B-Instruct-v0.1	35.33
Mixtral-8x7B-v0.1+PLLMIC	32.17
PLLM-8x7B-nc-instruct	67.17
PLLM-8x7B-nc-chat	68.17
Llama-3.1-8B-Instruct	22.67
Llama-3.1-8B+PLLMIC	24.67
Llama-PLLM-8B-instruct	58.00
Llama-PLLM-8B-chat	60.67
Llama-3.1-70B-Instruct	47.83
Llama-3.1-70B+PLLMIC	38.67
Llama-PLLM-70B-instruct	65.17
Llama-PLLM-70B-chat	66.33
Qwen-Max	50.83
GPT-4	59.50
Grok-2-1212	66.00
DeepSeek-v3	69.17
DeepSeek-R1	76.00
O1-2024-12-17	89.17

Table 4: Linguistic adaptation rate as evaluated on the **PLCC**: Polish Linguistic and Cultural Competency Benchmark

Model	ASR ↓	FRR ↓
Mistral-Nemo-Instruct-2407	21.85	0.62
PLLM-12B-nc-base	72.80	10.90
PLLM-12B-nc-instruct	77.61	0.62
PLLM-12B-nc-chat	1.03	3.31
Mixtral-8x7B-Instruct-v0.1	31.86	0.59
PLLM-8x7B-nc-base	74.35	6.95
PLLM-8x7B-nc-instruct	70.63	0.56
PLLM-8x7B-nc-chat	0.78	8.69
Llama-3.1-8B-Instruct	19.66	0.86
Llama-PLLM-8B-base	80.02	3.86
Llama-PLLM-8B-instruct	78.60	1.2
Llama-PLLM-8B-chat	0.76	5.27
Llama-3.1-70B-Instruct	22.27	0.36
Llama-PLLM-70B-base	76.35	2.01
Llama-PLLM-70B-instruct	70.69	0.36
Llama-PLLM-70B-chat	0.79	5.22

Table 5: Red-teaming evaluation results.

Model	LLMzSzŁ ↑
Llama-PLLM-8B-chat	47.68
PLLM-12B-nc-chat	53.40
PLLM-8x7B-nc-chat	60.52
Llama-PLLM-70B-chat	64.42
Meta-Llama-3.1-8B-Instruct	47.41
Mixtral-8x7B-Instruct-v0.1	49.46
Bielik-11B-v2.1-Instruct	57.52
Llama-3.3-70B-Instruct	67.13

Table 6: Academic performance as evaluated on **LLMzSzŁ**: a comprehensive LLM benchmark for Polish

inspired by the "Rainbow Teaming framework" ([Samvelyan et al., 2024](#)). For the ASR, the Llama-Guard model was utilized to assess the percentage of unsafe responses, whereas for the FRR, we prompted one of our trained models to obtain the proportion of refusals to benign queries. In general, the PLLuM models fine-tuned on instructions are characterized by a relatively higher ASR and a lower FRR than their derivatives aligned on human preferences.

6 PLLuMIC Public Sample

Apart from evaluating the impact of the different phases of model training on its linguistic adaptation, we release a representative subset of the organic PLLuM instruction corpus. Overall the first release of the dataset contains a total of 1278 human-authored instructions, spanning across 12 types, 126 subtypes and 34 topics. Substantial effort has been made to ensure a wide diversity and high quality, both reflected in each data sample. Appendix [D.1](#) details the subset's typology.

7 Conclusions

We believe that our description of the PLLuM Instruction Corpus along with its public subset can be used to design and complement manual and automated annotation work in other LLM projects. Thanks to iterative instruction corpus development and continual evaluation we established that effective fine-tuning on language-specific instructions requires models to first undergo sufficient continual pre-training on the target language. We also identified several cases of negative linguistic transfer, where conventions from dominant languages (particularly English) can interfere with idiomatic text generation in the target language. Such interference

may occur both in the process during fine-tuning and alignment on synthetic instructions. This highlights the need for high-quality, language-specific organic instructions in linguistically adapted LLMs.

8 Availability

The manually annotated sample of PLLuMIC can be accessed at <https://huggingface.co/datasets/pelcra/PLLuMIC>. We are planning to release its synthetic extension (PLLuMIC-syn-ext) separately.

9 Acknowledgments

The work reported in this paper was funded by several grants:

- The continued pre-training of the 8B, 12B, and 70B models reported in Table 4, most of their fine-tuning and alignment were performed on the WCSS HPC infrastructure as part of an earmarked grant (1/WI/DBiL/2023) from the Polish Ministry of Digital Affairs.
- The continued pre-training of variants of the 8x7B and 12B models reported above were performed on the ACC Cyfronet AGH infrastructure under a grant no. PLG/2024/017788.
- The first edition of the PLLuMIC subset released with this paper was developed after the completion of the PLLuM project and supported by the grant CLARIN-BIZ-bis (FENG.02.04-IP.04-0004/24).

References

- Marah Abdin, Jyoti Aneja, Harkirat Behl, Sébastien Bubeck, Ronen Eldan, Suriya Gunasekar, Michael Harrison, Russell J. Hewett, Moján Javaheripan, Piero Kauffmann, James R. Lee, Yin Tat Lee, Yuanzhi Li, Weishung Liu, Caio C. T. Mendes, Anh Nguyen, Eric Price, Gustavo de Rosa, Olli Saarikivi, and 8 others. 2024. *Phi-4 Technical Report*. *arXiv preprint*. ArXiv:2412.08905 [cs].
- Shourya Aggarwal, Divyanshu Mandowara, Vishwajeet Agrawal, Dinesh Khandelwal, Parag Singla, and Dinesh Garg. 2021. Explanations for CommonsenseQA: New Dataset and Models. In *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)*, Online. Association for Computational Linguistics.
- AI@Meta. 2024. *Llama 3 model card*.
- Ebtessam Almazrouei, Hamza Alobeidli, Abdulaziz Alshamsi, Alessandro Cappelli, Ruxandra Cojocaru, Mérouane Debbah, Étienne Goffinet, Daniel Hesslow, Julien Launay, Quentin Malartic, Daniele Mazzotta, Badreddine Noune, Baptiste Pannier, and Guilherme Penedo. 2023. *The Falcon Series of Open Language Models*. *arXiv preprint*. ArXiv:2311.16867 [cs].
- Yuvanesh Anand, Zach Nussbaum, Brandon Duderstadt, Benjamin Schmidt, and Andriy Mulyar. 2023. Gpt4all: Training an assistant-style chatbot with large scale data distillation from gpt-3.5-turbo. <https://github.com/nomic-ai/gpt4all>.
- Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, and 1 others. 2022. Training a helpful and harmless assistant with reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*.
- Faeze Brahman, Sachin Kumar, Vidhisha Balachandran, Pradeep Dasigi, Valentina Pyatkin, Abhilasha Ravichander, Sarah Wiegreffe, Nouha Dziri, Khyathi Chandu, Jack Hessel, Yulia Tsvetkov, Noah A. Smith, Yejin Choi, and Hannaneh Hajishirzi. 2024. *The art of saying no: Contextual noncompliance in language models*. *Preprint*, arXiv:2407.12043.
- Bartosz Broda, Michał Marcińczuk, Marek Maziarz, Adam Radziszewski, and Adam Wardyński. 2012. *KPWr: Towards a free corpus of Polish*. In *Proceedings of the Eighth International Conference on Language Resources and Evaluation (LREC’12)*, pages 3218–3222, Istanbul, Turkey. European Language Resources Association (ELRA).
- Linzheng Chai, Shukai Liu, Jian Yang, Yuwei Yin, Ke Jin, Jiaheng Liu, Tao Sun, Ge Zhang, Changyu Ren, Hongcheng Guo, Zekun Wang, Boyang Wang, Xianjie Wu, Bing Wang, Tongliang Li, Liqun Yang, Sufeng Duan, and Zhoujun Li. 2024. *Mceval: Massively multilingual code evaluation*. *Preprint*, arXiv:2406.07436.
- Jianlv Chen, Shitao Xiao, Peitian Zhang, Kun Luo, Defu Lian, and Zheng Liu. 2023a. *Bge m3-embedding: Multi-lingual, multi-functionality, multi-granularity text embeddings through self-knowledge distillation*. *Preprint*, arXiv:2309.07597.
- Wenhu Chen, Ming Yin, Max Ku, Pan Lu, Yixin Wan, Xueguang Ma, Jianyu Xu, Xinyi Wang, and Tony Xia. 2023b. Theoremqa: A theorem-driven question answering dataset. In *The 2023 Conference on Empirical Methods in Natural Language Processing*.
- Aleksandra Chrabrowa, Łukasz Dragan, Karol Grzegorczyk, Dariusz Kajtoch, Mikołaj Koszowski, Robert Mroczkowski, and Piotr Rybak. 2022. *Evaluation of transfer learning for Polish with a text-to-text model*. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 4374–4394, Marseille, France. European Language Resources Association.

- Mike Conover, Matt Hayes, Ankit Mathur, Jianwei Xie, Jun Wan, Sam Shah, Ali Ghodsi, Patrick Wendell, Matei Zaharia, and Reynold Xin. 2023. [Free dolly: Introducing the world’s first truly open instruction-tuned llm](#).
- Susan Conrad. 2023. [Register in corpus linguistics: the role and legacy of Douglas Biber](#). *Corpus Linguistics and Linguistic Theory*, 19(1):7–21.
- Ślawomir Dadas. 2022. [Training effective neural sentence encoders from automatically mined paraphrases](#). In *2022 IEEE International Conference on Systems, Man, and Cybernetics (SMC)*, pages 371–378.
- Ślawomir Dadas, Michał Perelkiewicz, and Rafał Poświata. 2020. [Evaluation of sentence representations in Polish](#). In *Proceedings of the 12th Language Resources and Evaluation Conference*, pages 1674–1680, Marseille, France. European Language Resources Association.
- Ślawomir Dadas, Małgorzata Grębowiec, Michał Perelkiewicz, and Rafał Poświata. 2025. [Evaluating polish linguistic and cultural competency in large language models](#). *Preprint*, arXiv:2503.00995.
- DeepSeek-AI. 2024. [Deepseek-v2: A strong, economical, and efficient mixture-of-experts language model](#). *Preprint*, arXiv:2405.04434.
- DeepSeek-AI, Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Daya Guo, Dejian Yang, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, and 181 others. 2024. [DeepSeek-V3 Technical Report](#). *arXiv preprint*. ArXiv:2412.19437 [cs].
- Ning Ding, Yulin Chen, Bokai Xu, Yujia Qin, Zhi Zheng, Shengding Hu, Zhiyuan Liu, Maosong Sun, and Bowen Zhou. 2023. [Enhancing chat language models by scaling high-quality instructional conversations](#). *Preprint*, arXiv:2305.14233.
- Guanting Dong, Keming Lu, Chengpeng Li, Tingyu Xia, Bowen Yu, Chang Zhou, and Jingren Zhou. 2024. [Self-play with execution feedback: Improving instruction-following capabilities of large language models](#). *Preprint*, arXiv:2406.13542.
- Kawin Ethayarajh, Winnie Xu, Niklas Muennighoff, Dan Jurafsky, and Douwe Kiela. 2024. [Kto: Model alignment as prospect theoretic optimization](#). *Preprint*, arXiv:2402.01306.
- Team Falcon-LLM. 2024. [The falcon 3 family of open models](#).
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, and 542 others. 2024. [The llama 3 herd of models](#). *Preprint*, arXiv:2407.21783.
- Dirk Groeneveld, Iz Beltagy, Pete Walsh, Akshita Bhagia, Rodney Kinney, Oyvind Tafjord, Ananya Harsh Jha, Hamish Ivison, Ian Magnusson, Yizhong Wang, Shane Arora, David Atkinson, Russell Author, Khyathi Raghavi Chandu, Arman Cohan, Jennifer Dumas, Yanai Elazar, Yuling Gu, Jack Hessel, and 24 others. 2024. [Olmo: Accelerating the science of language models](#). *Preprint*, arXiv:2402.00838.
- Włodzimierz Gruszczyński, Dorota Adamiec, Renata Bronikowska, Witold Kieraś, Emanuel Modrzejewski, Aleksandra Wieczorek, and Marcin Woliński. 2022. [The electronic corpus of 17th- and 18th-century Polish texts](#). *Language Resources and Evaluation*, 56(1):309–332.
- Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. 2021. [Measuring mathematical problem solving with the math dataset](#). *NeurIPS*.
- Jiwoo Hong, Noah Lee, and James Thorne. 2024. [Orpo: Monolithic preference optimization without reference model](#). *Preprint*, arXiv:2403.07691.
- Hakan Inan, Kartikeya Upasani, Jianfeng Chi, Rashi Rungta, Krithika Iyer, Yuning Mao, Michael Tontchev, Qing Hu, Brian Fuller, Davide Testuggine, and Madiam Khabsa. 2023. [Llama guard: Llm-based input-output safeguard for human-ai conversations](#). *Preprint*, arXiv:2312.06674.
- Hamish Ivison, Yizhong Wang, Valentina Pyatkin, Nathan Lambert, Matthew Peters, Pradeep Dasigi, Joel Jang, David Wadden, Noah A. Smith, Iz Beltagy, and Hannaneh Hajishirzi. 2023. [Camels in a changing climate: Enhancing lm adaptation with tulu 2](#). *Preprint*, arXiv:2311.10702.
- Arkadiusz Janz, Agnieszka Dziob, Marcin Oleksy, and Joanna Baran. 2022. [A unified sense inventory for word sense disambiguation in polish](#). In *Computational Science – ICCS 2022*, pages 682–689, Cham. Springer International Publishing.
- Arkadiusz Janz, Grzegorz Kostkowski, and Marek Maziarz. 2021. [Constructing vesnet: Mapping lod thesauri onto princeton wordnet and polish wordnet](#). In *Advances in Computational Collective Intelligence*, pages 608–620, Cham. Springer International Publishing.
- Arkadiusz Janz, Dominik Kurowski, Joanna Baran, Julia Moska, Tomasz Bernaś, and Marcin Oleksy. 2024. [Refining natural language inferences using cross-document structure theory](#). In *Computational Collective Intelligence*, pages 263–276, Cham. Springer Nature Switzerland.
- Krzysztof Jassem, Michał Ciesiółka, Filip Graliński, Piotr Jabłoński, Jakub Pokrywka, Marek Kubis, Monika Jabłońska, and Ryszard Staruch. 2025. [Llmzszl: a comprehensive llm benchmark for polish](#). *Preprint*, arXiv:2501.02266.

- Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Lélio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2023. *Mistral 7b*. Preprint, arXiv:2310.06825.
- Albert Q. Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, Gianna Lengyel, Guillaume Bour, Guillaume Lample, Lélio Renard Lavaud, Lucile Saulnier, Marie-Anne Lachaux, Pierre Stock, Sandeep Subramanian, Sophia Yang, and 7 others. 2024. *Mixtral of Experts*. arXiv preprint, ArXiv:2401.04088 [cs].
- Agnieszka Karlińska, Piotr Miłkowski, Paulina Czwordon-Lis, Bartłomiej Kopytyna, and Jan Kocoń. 2024. *Comprehensive sentiment analysis of polish book reviews using large and small language models*.
- W. Kieraś, M. Marciak, M. Łaziński, M. Woliński, K. Bojałkowska, W. Eźlakowski, Ł. Kobyliński, D. Komosińska, K. Krasnowska-Kieraś, M. Rudolf, A. Tomaszewska, J. Wołoszyn, and N. Zawadzka-Paluktau. 2024. *Korpus Współczesnego Języka Polskiego. Dekada 2011–2020. Język Polski*.
- Witold Kieraś and Marcin Woliński. 2018. *Manually annotated corpus of Polish texts published between 1830 and 1918*. In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*, pages 3854–3859, Paris, France. European Language Resources Association (ELRA).
- Łukasz Kobyliński, Witold Kieraś, and Szymon Rynkun. 2021. *PolEval 2021 Task 3: Post-correction of OCR Results*. In *Proceedings of the PolEval 2021 Workshop*, pages 85–91, Warszawa. Institute of Computer Science, Polish Academy of Sciences.
- Jan Kocoń, Piotr Miłkowski, and Monika Zaśko-Zielńska. 2019. *Multi-level sentiment analysis of PolEmo 2.0: Extended corpus of multi-domain consumer reviews*. In *Proceedings of the 23rd Conference on Computational Natural Language Learning (CoNLL)*, pages 980–991, Hong Kong, China. Association for Computational Linguistics.
- Anna Kolos, Inez Okulska, Kinga Głabińska, Agnieszka Karlińska, Emilia Wiśnios, Paweł Ellerik, and Andrzej Prałat. 2024. Ban-pl: A polish dataset of banned harmful and offensive content from wykop.pl web service. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 2107–2118.
- Aleksandra Krasnodębska, Maciej Chrabaszcz, and Wojciech Kusa. 2025. Rainbow-teaming for the polish language: A reproducibility study. In *Proceedings of the TrustNLP: Fifth Workshop on Trustworthy Natural Language Processing at NAACL*. Accepted.
- Bespoke Labs. 2025. Bespoke-stratos: The unreasonable effectiveness of reasoning distillation. <https://www.bespokelabs.ai/blog/bespoke-stratos-the-unreasonable-effectiveness-of-reasoning-distillation>. Accessed: 2025-01-22.
- Nathan Lambert, Jacob Morrison, Valentina Pyatkin, Shengyi Huang, Hamish Ivison, Faeze Brahman, Lester James V. Miranda, Alisa Liu, Nouha Dziri, Shane Lyu, Yuling Gu, Saumya Malik, Victoria Graf, Jena D. Hwang, Jiangjiang Yang, Ronan Le Bras, Oyvind Tafjord, Chris Wilhelm, Luca Soldaini, and 4 others. 2024a. Tülu 3: Pushing frontiers in open language model post-training.
- Nathan Lambert, Jacob Morrison, Valentina Pyatkin, Shengyi Huang, Hamish Ivison, Faeze Brahman, Lester James V. Miranda, Alisa Liu, Nouha Dziri, Shane Lyu, Yuling Gu, Saumya Malik, Victoria Graf, Jena D. Hwang, Jiangjiang Yang, Ronan Le Bras, Oyvind Tafjord, Chris Wilhelm, Luca Soldaini, and 4 others. 2024b. Tülu 3: Pushing frontiers in open language model post-training.
- Matthew Lamm, Jennimaria Palomaki, Chris Alberti, Daniel Andor, Eunsol Choi, Livio Baldini Soares, and Michael Collins. 2020. *Qed: A framework and dataset for explanations in question answering*. Preprint, arXiv:2009.06354.
- Ariel N. Lee, Cole J. Hunter, and Nataniel Ruiz. 2024. *Platypus: Quick, cheap, and powerful refinement of llms*. Preprint, arXiv:2308.07317.
- Peng Li, Yeye He, Dror Yashar, Weiwei Cui, Song Ge, Haidong Zhang, Danielle Rifinski Fainman, Dongmei Zhang, and Surajit Chaudhuri. 2023. *Table-gpt: Table-tuned gpt for diverse table tasks*. Preprint, arXiv:2310.09263.
- Wing Lian, Bleys Goodson, Eugene Pentland, Austin Cook, Chanvichet Vong, and "Teknium". 2023a. *Openorca: An open dataset of gpt augmented flan reasoning traces*.
- Wing Lian, Guan Wang, Bleys Goodson, Eugene Pentland, Austin Cook, Chanvichet Vong, and "Teknium". 2023b. *Slimorca: An open dataset of gpt-4 augmented flan reasoning traces, with verification*.
- Hunter Lightman, Vineet Kosaraju, Yura Burda, Harry Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. 2023. Let's verify step by step. arXiv preprint arXiv:2305.20050.
- Wei Liu, Weihao Zeng, Keqing He, Yong Jiang, and Junxian He. 2024. *What makes good data for alignment? a comprehensive study of automatic data selection in instruction tuning*. In *The Twelfth International Conference on Learning Representations*.

- Shayne Longpre, Le Hou, Tu Vu, Albert Webson, Hyung Won Chung, Yi Tay, Denny Zhou, Quoc V. Le, Barret Zoph, Jason Wei, and Adam Roberts. 2023. *The flan collection: Designing data and methods for effective instruction tuning*. *Preprint*, arXiv:2301.13688.
- Pan Lu, Swaroop Mishra, Tony Xia, Liang Qiu, Kai-Wei Chang, Song-Chun Zhu, Oyvind Tafjord, Peter Clark, and Ashwin Kalyan. 2022. Learn to explain: Multimodal reasoning via thought chains for science question answering. In *The 36th Conference on Neural Information Processing Systems (NeurIPS)*.
- Michał Marcińczuk, Marcin Ptak, Adam Radziszewski, and Maciej Piasecki. 2013. Open dataset for development of polish question answering systems. In *Proceedings of the 6th Language & Technology Conference: Human Language Technologies as a Challenge for Computer Science and Linguistics*. Wydawnictwo Poznańskie, Fundacja Uniwersytetu im. Adama Mickiewicza.
- Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, Pouya Tafti, Léonard Hussenot, Pier Giuseppe Sessa, Aakanksha Chowdhery, Adam Roberts, Aditya Barua, Alex Botev, Alex Castro-Ros, Ambrose Sloane, Amélie Héliou, and 88 others. 2024. *Gemma: Open Models Based on Gemini Research and Technology*. *arXiv preprint*. ArXiv:2403.08295 [cs].
- Jinjie Ni, Fuzhao Xue, Kabir Jain, Mahir Hitesh Shah, Zangwei Zheng, and Yang You. 2023. Instruction in the wild: A user-based instruction dataset. <https://github.com/XueFuzhao/InstructionWild>.
- Nvidia, Bo Adler, Niket Agarwal, Ashwath Aithal, Dong H. Anh, Pallab Bhattacharya, Annika Brundyn, Jared Casper, Bryan Catanzaro, Sharon Clay, Jonathan Cohen, Sirshak Das, Ayush Dattagupta, Olivier Delalleau, Leon Derczynski, Yi Dong, Daniel Egert, Ellie Evans, Aleksander Ficek, and 63 others. 2024. *Nemotron-4 340B Technical Report*. *arXiv preprint*. ArXiv:2406.11704 [cs].
- Maciej Ogrodniczuk, Katarzyna Głowińska, Mateusz Kopeć, Agata Savary, and Magdalena Zawiśawska. 2016. Polish coreference corpus. In *Human Language Technology. Challenges for Computer Science and Linguistics*, pages 215–226, Cham. Springer International Publishing.
- Maciej Ogrodniczuk and Mateusz Kopeć. 2014. *The Polish Summaries Corpus*. In *Proceedings of the Ninth International Conference on Language Resources and Evaluation, LREC 2014*, pages 3712–3715, Reykjavík, Iceland. European Language Resources Association (ELRA).
- Yasumasa Onoe, Michael J.Q. Zhang, Eunsol Choi, and Greg Durrett. 2021. Creak: A dataset for common-sense reasoning over entity knowledge. *OpenReview*.
- OpenAI. 2022. *Chatgpt: Optimizing language models for dialogue*. Accessed: 2025-03-18.
- OpenAI. 2023. *Gpt-3.5: Generative pre-trained transformer 3.5*. Accessed: 2025-03-18.
- Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. 2022. *Training language models to follow instructions with human feedback*. *Preprint*, arXiv:2203.02155.
- Baolin Peng, Chunyuan Li, Pengcheng He, Michel Galley, and Jianfeng Gao. 2023. *Instruction tuning with gpt-4*. *Preprint*, arXiv:2304.03277.
- Piotr Pęzik, Gosia Krawentek, Sylwia Karasińska, Paweł Wilk, Paulina Rybińska, Anna Cichosz, Angelika Peljak-Łapińska, Mikołaj Deckert, and Michał Adamczyk. 2022. *DiaBiz – an annotated corpus of Polish call center dialogs*. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 723–726, Marseille, France. European Language Resources Association.
- Adam Przepiórkowski, Mirosław Bańko, Rafał L. Górski, and Barbara Lewandowska-Tomaszczyk, editors. 2012. *Narodowy Korpus Języka Polskiego*. Wydawnictwo Naukowe PWN, Warsaw.
- Michał Ptaszynski, Agata Pieciukiewicz, Paweł Dybala, Paweł Skrzek, Kamil Soliwoda, Marcin Fortuna, Gniewosz Leliwa, and Michał Wroczynski. 2023. Expert-annotated dataset to study cyberbullying in polish language. *Data*, 9(1):1.
- Piotr Pęzik. 2016. *Exploring phraseological equivalence with Paralela*. In Ewa Gruszczyńska and Agnieszka Leńko-Szymańska, editors, *Polish-language Parallel Corpora*, pages 67–81. Instytut Lingwistyki Stosowanej UW, Warsaw.
- Piotr Pęzik, Gosia Krawentek, Sylwia Karasińska, Paweł Wilk, Paulina Rybińska, Anna Cichosz, Angelika Peljak-Łapińska, Mikołaj Deckert, and Michał Adamczyk. 2022. *DiaBiz*. CLARIN-PL digital repository.
- Zheng Lin Qingyi Si. 2023. *Alpaca-cot: An instruction fine-tuning platform with instruction data collection and unified large language models interface*. Institute of Information Engineering, Chinese Academy of Sciences, Beijing, China.
- Shanghaoran Quan, Tianyi Tang, Bowen Yu, An Yang, Dayiheng Liu, Bofei Gao, Jianhong Tu, Yichang Zhang, Jingren Zhou, and Junyang Lin. 2024. *Language models can self-lengthen to generate long texts*. *Preprint*, arXiv:2410.23933.
- Qwen, An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin

- Yang, Jiaxi Yang, Jingren Zhou, Junyang Lin, and 24 others. 2025. *Qwen2.5 technical report*. Preprint, arXiv:2412.15115.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D. Manning, and Chelsea Finn. 2024. *Direct preference optimization: Your language model is secretly a reward model*. Preprint, arXiv:2305.18290.
- Nazneen Rajani, Lewis Tunstall, Edward Beeching, Nathan Lambert, Alexander M. Rush, and Thomas Wolf. 2023. No robots.
- Abhinav Rastogi, Xiaoxue Zang, Srinivas Sunkara, Raghav Gupta, and Pranav Khaitan. 2020. Towards scalable multi-domain conversational agents: The schema-guided dialogue dataset. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 8689–8696.
- Piotr Rybak, Piotr Przybyła, and Maciej Ogrodniczuk. 2024. *PolQA: Polish question answering dataset*. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 12846–12855, Torino, Italia. ELRA and ICCL.
- Zygmunt Saloni, Marcin Woliński, Robert Wołosz, Włodzimierz Gruszczyński, and Danuta Skowrońska. 2015. *Słownik gramatyczny języka polskiego*, 3rd edition. Warsaw.
- Mikayel Samvelyan, Sharath Chandra Raparthy, Andrei Lupu, Eric Hambro, Aram H. Markosyan, Manish Bhatt, Yuning Mao, Minqi Jiang, Jack Parker-Holder, Jakob Foerster, Tim Rocktäschel, and Roberta Raileanu. 2024. *Rainbow teaming: Open-ended generation of diverse adversarial prompts*. Preprint, arXiv:2402.16822.
- Victor Sanh, Albert Webson, Colin Raffel, Stephen H. Bach, Lintang Sutawika, Zaid Alyafeai, Antoine Chaffin, Arnaud Stiegler, Teven Le Scao, Arun Raja, Manan Dey, M Saiful Bari, Canwen Xu, Urmish Thakker, Shanya Sharma Sharma, Eliza Szczechla, Taewoon Kim, Gunjan Chhablani, Nihal Nayak, and 22 others. 2022. *Multitask prompted training enables zero-shot task generalization*. Preprint, arXiv:2110.08207.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. *Proximal policy optimization algorithms*. Preprint, arXiv:1707.06347.
- Larry Selinker. 1969. Language transfer. *General linguistics*, 9(2):67.
- Ilia Shumailov, Zakhar Shumaylov, Yiren Zhao, Nicolas Papernot, Ross Anderson, and Yarin Gal. 2024. Ai models collapse when trained on recursively generated data. *Nature*, 631(8022):755–759.
- Snowflake AI Research. 2024. *Snowflake Arctic Cookbook Series: Arctic’s Approach to Data*.
- Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023. Stanford alpaca: An instruction-following llama model.
- OpenThoughts Team. 2025. Open Thoughts. <https://open-thoughts.ai>.
- Teknium. 2023. *Openhermes 2.5: An open dataset of synthetic data for generalist llm assistants*.
- Ryszard Tuora, Aleksandra Zwierzchowska, Natalia Zawadzka-Paluekta, Cezary Klamra, and Łukasz Kobylinski. 2023. *Poquad — the polish question answering dataset — description and analysis*.
- Tamás Váradi, Bence Nyéki, Svetla Koeva, Marko Tadić, Vanja Štefanec, Maciej Ogrodniczuk, Bartłomiej Nitój, Piotr Pęzik, Verginica Barbu Mititelu, Elena Irimia, Maria Mitrofan, Dan Tufiș, Radovan Garabík, Simon Krek, and Andraž Repar. 2022. *Introducing the CURLICAT corpora: Seven-language domain specific annotated corpora from curated sources*. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 100–108, Marseille, France. European Language Resources Association.
- David Wadden, Kejian Shi, Jacob Morrison, Aakanksha Naik, Shruti Singh, Nitzan Barzilay, Kyle Lo, Tom Hope, Luca Soldaini, Shannon Zejiang Shen, Doug Downey, Hannaneh Hajishirzi, and Arman Cohan. 2024. *Sciriff: A resource to enhance language model instruction-following over scientific literature*. Preprint, arXiv:2406.07835.
- Guan Wang, Sijie Cheng, Xianyuan Zhan, Xiangang Li, Sen Song, and Yang Liu. 2024a. *Openchat: Advancing open-source language models with mixed-quality data*. Preprint, arXiv:2309.11235.
- Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A. Smith, Daniel Khashabi, and Hannaneh Hajishirzi. 2022a. *Self-instruct: Aligning language model with self generated instructions*.
- Yizhong Wang, Swaroop Mishra, Pegah Alipoormolabashi, Yeganeh Kordi, Amirreza Mirzaei, Anjana Arunkumar, Arjun Ashok, Arut Selvan Dhanasekaran, Atharva Naik, David Stap, Eshaan Pathak, Giannis Karamanolakis, Haizhi Gary Lai, Ishan Purohit, Ishani Mondal, Jacob Anderson, Kirby Kuznia, Krima Doshi, Maitreya Patel, and 21 others. 2022b. *Super-naturalinstructions: Generalization via declarative instructions on 1600+ nlp tasks*. Preprint, arXiv:2204.07705.
- Zhilin Wang, Yi Dong, Olivier Delalleau, Jiaqi Zeng, Gerald Shen, Daniel Egert, Jimmy J. Zhang, Makesh Narsimhan Sreedhar, and Oleksii Kuchaiev. 2024b. *Helpsteer2: Open-source dataset for training top-performing reward models*. Preprint, arXiv:2406.08673.

- Zhilin Wang, Yi Dong, Jiaqi Zeng, Virginia Adams, Makesh Narsimhan Sreedhar, Daniel Egert, Olivier Delalleau, Jane Polak Scowcroft, Neel Kant, Aidan Swope, and Oleksii Kuchaiev. 2023. **Helpsteer: Multi-attribute helpfulness dataset for steerlm**. *Preprint*, arXiv:2311.09528.
- Jason Wei, Maarten Bosma, Vincent Y. Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M. Dai, and Quoc V. Le. 2022. **Finetuned language models are zero-shot learners**. *Preprint*, arXiv:2109.01652.
- Marcin Woliński and Elżbieta Hajnicz. 2021. **Składnica: a constituency treebank of Polish harmonised with the Walenty valency dictionary**. *Language Resources and Evaluation*, 55:209–239.
- Alina Wróblewska and Katarzyna Krasnowska-Kieraś. 2017. Polish evaluation dataset for compositional distributional semantics models. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 784–792.
- Canwen Xu, Daya Guo, Nan Duan, and Julian McAuley. 2023. **Baize: An open-source chat model with parameter-efficient tuning on self-chat data**. *Preprint*, arXiv:2304.01196.
- Zhangchen Xu, Fengqing Jiang, Luyao Niu, Yun-tian Deng, Radha Poovendran, Yejin Choi, and Bill Yuchen Lin. 2024. **Magpie: Alignment data synthesis from scratch by prompting aligned llms with nothing**. *Preprint*, arXiv:2406.08464.
- An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jialong Tang, Jialin Wang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Ma, and 43 others. 2024. **Qwen2 technical report**. *Preprint*, arXiv:2407.10671.
- Longhui Yu, Weisen Jiang, Han Shi, Jincheng Yu, Zhengying Liu, Yu Zhang, James T. Kwok, Zhenguo Li, Adrian Weller, and Weiyang Liu. 2024. **Metamath: Bootstrap your own mathematical questions for large language models**. *Preprint*, arXiv:2309.12284.
- Xiang Yue, Tuney Zheng, Ge Zhang, and Wenhua Chen. 2024. **Mammoth2: Scaling instructions from the web**. *Advances in Neural Information Processing Systems*.
- Bo-Wen Zhang, Yan Yan, Lin Li, and Guang Liu. 2024. **Infinitymath: A scalable instruction tuning dataset in programmatic mathematical reasoning**. *Preprint*, arXiv:2408.07089.
- Jianguo Zhang, Kun Qian, Zhiwei Liu, Shelby Heinecke, Rui Meng, Ye Liu, Zhou Yu, Silvio Savarese, and Caiming Xiong. 2023. **Dialogstudio: Towards richest and most diverse unified dataset collection for conversational ai**. *arXiv preprint arXiv:2307.10172*.
- Hanyu Zhao, Li Du, Yiming Ju, Chengwei Wu, and Tengfei Pan. 2024a. **Beyond iid: Optimizing instruction learning from the perspective of instruction interaction and dependency**. *Preprint*, arXiv:2409.07045.
- Wenting Zhao, Xiang Ren, Jack Hessel, Claire Cardie, Yejin Choi, and Yuntian Deng. 2024b. **Wildchat: 1m chatgpt interaction logs in the wild**. *Preprint*, arXiv:2405.01470.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Tianle Li, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zhuohan Li, Zi Lin, Eric. P Xing, Joseph E. Gonzalez, Ion Stoica, and Hao Zhang. 2023. **Lmsys-chat-1m: A large-scale real-world llm conversation dataset**. *Preprint*, arXiv:2309.11998.
- Tianyu Zheng, Ge Zhang, Tianhao Shen, Xueling Liu, Bill Yuchen Lin, Jie Fu, Wenhua Chen, and Xiang Yue. 2025. **Opencodeinterpreter: Integrating code generation with execution and refinement**. *Preprint*, arXiv:2402.14658.
- Yaowei Zheng, Richong Zhang, Junhao Zhang, Yanhan Ye, Zheyuan Luo, Zhangchi Feng, and Yongqiang Ma. 2024. **Llamafactory: Unified efficient fine-tuning of 100+ language models**. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 3: System Demonstrations)*, Bangkok, Thailand. Association for Computational Linguistics.
- Chunting Zhou, Pengfei Liu, Puxin Xu, Srini Iyer, Jiao Sun, Yuning Mao, Xuezhe Ma, Avia Efrat, Ping Yu, Lili Yu, Susan Zhang, Gargi Ghosh, Mike Lewis, Luke Zettlemoyer, and Omer Levy. 2023a. **LIMA: Less Is More for Alignment**. *arXiv preprint*, ArXiv:2305.11206 [cs].
- Jeffrey Zhou, Tianjian Lu, Swaroop Mishra, Siddhartha Brahma, Sujoy Basu, Yi Luan, Denny Zhou, and Le Hou. 2023b. **Instruction-Following Evaluation for Large Language Models**. *arXiv preprint*. Version Number: 1.

A Availability of Instruction Datasets

A.1 Selected LLMs and their source datasets

Table 7 summarizes the availability status of instruction datasets for a number of open-weight models. More specifically, Llama 3.1 (Grattafiori et al., 2024) offers a general description of the data preparation process, including the sampling of synthetic and human-annotated instructions, but does not share the data itself.

OpenChat utilizes the acclaimed ShareGPT¹⁰ dataset of prompts and responses generated by OpenAIs GPT-3.5 and GPT-4. The accompanying paper (Wang et al., 2024a) offers a condensed analysis of the data distribution and quality.

Qwen 2.5 (Qwen et al., 2025) uses various datasets synthesized according to their guidelines during the post pre-training phase (Quan et al., 2024; Dong et al., 2024), either sampling existing datasets (e.g. (Chai et al., 2024)) or using web-scraped data as input.

Qwen2 (Yang et al., 2024) gives out more details concerning the human-annotation process in the data generation process. However, no ready-to-use data was published alongside these models.

OLMO’s (Groeneveld et al., 2024) fine-tuning is based on the Tulu2 dataset (Ivison et al., 2023), recently developed into Tulu3 (Lambert et al., 2024a). Tulu2 consists of publicly available datasets such as FLAN (Wei et al., 2022), No Robots (Rajani et al., 2023) and WildChat (Zhao et al., 2024b).

The authors of Mixtral 8x7B (Jiang et al., 2024) offer no description of the employed instruction data, while the publication of Mistral 7B (Jiang et al., 2023) points to loosely defined “instruction datasets publicly available on the Hugging Face repository. No analysis can be found in either of the papers.

Falcon (Falcon-LLM, 2024; Almazrouei et al., 2023) is predominantly focused on pre-training but its various *-instruct* versions utilize mainly Baize (Xu et al., 2023) dataset, sourcing also from other GPT4-based online data repositories, such as GPT4All (Anand et al., 2023) or GPTeacher.

Gemma (Mesnard et al., 2024) uses an undisclosed teacher-model to generate answers for synthetic and human-made prompts, joining it with a “mixture of internal and external public data” but offers no insight as for the composition of these datasets.

Microsoft’s Phi model (Abdin et al., 2024) uses

undisclosed publicly available datasets to generate synthetic responses for supervised fine-tuning (SFT) and gives no detailed description of their contents.

In contrast, Dolly (Conover et al., 2023) use their open-source dolly-bricks-15k dataset comprising 15k human-generated prompt-response pairs inspired by InstructGPT (Ouyang et al., 2022), accompanied by extensive documentation, including annotation guidelines.

For their V3 model (DeepSeek-AI et al., 2024), Deepseek uses other iterations of models such as V2.5 (DeepSeek-AI, 2024) or R1 as “expert models” to generate instruction data from scratch.

Nvidia’s Nemotron (Nvidia et al., 2024) relies on their own Helpsteer2 dataset (Wang et al., 2024b) and Mixtral-8x7b’s abilities to generate synthetic instruction data. The data generation process is documented, but only the seed dataset is available.

¹⁰No longer available online.

Table 7: Availability of instruction datasets for selected open LLMs. We characterize LLMs especially based on the availability of information concerning the annotation process and synthetic data generation (SDG). Ideally, we would expect the final instruction mix used in SFT to be fully documented (e.g. exact proportions of each instruction type).

Model	Instructions	Documentation
LLaMA 3.1 (Grattafiori et al., 2024)	Unavailable	SDG & annotation process explained
Mixtral (Jiang et al., 2024)	Unavailable	(Jiang et al., 2023) points to “instruction datasets publicly available on the Hugging Face repository”
Falcon3 (Falcon-LLM, 2024)	Unavailable	None
Gemma (Mesnard et al., 2024)	Unavailable	High-level desc. of SDG process
Nemotron (Nvidia et al., 2024)	As input for SDG	SDG process explained
Snowflake-arctic (Snowflake AI Research, 2024)	Unavailable	None
Phi-3.5/4 (Abdin et al., 2024)	Unavailable	High-level desc. of SDG process
Dolly (Conover et al., 2023)	Available	Annotation process explained
OpenChat (Wang et al., 2024a)	Unavailable	Low
Qwen 2.5 (Qwen et al., 2025)	Unavailable	Low
DeepSeek (DeepSeek-AI et al., 2024)	Unavailable	SDG process explained
OLMO (Groeneveld et al., 2024)	Available	Inventory of component datasets

B Stand-alone Instruction Datasets

B.1 Transparency and Representativeness

Several instruction datasets available as open stand-alone collections have also been used in more experimental LLM research projects. One of the early large-scale resources of instructions is the original FLAN (Wei et al., 2022) dataset, followed by the FLAN Collection (Longpre et al., 2023) dataset published by Google Research. The OpenOrca¹¹ (Lian et al., 2023a) dataset complemented FLAN with explanation traces and step-by-step thought processes from GPT-3.5 (OpenAI, 2023) and GPT-4 (Peng et al., 2023). Other frequently utilized datasets include Databrick’s Dolly (Conover et al., 2023) 15k and LIMA (Zhou et al., 2023a) 1k dataset both consisting of curated and hand-written examples, LMSYS-Chat-1M (Zheng et al., 2023) with 1 million human-LLM conversations, WebInstruct (Yue et al., 2024) with 10 million instruction pairs

harvested and refined from the web, UltraChat (Ding et al., 2023) comprising 1.5 million multi-turn dialogues generated by ChatGPT (OpenAI, 2022) from C4 data, SelfInstruct (Wang et al., 2022a) containing 52k synthetic instructions bootstrapped from 175 hand-written examples by GPT-3 (Ouyang et al., 2022). Similarly, the Stanford Alpaca dataset (Taori et al., 2023) contains 52k instructions created with OpenAI’s text-davinci-003 model (OpenAI, 2023) and it was subsequently reconstructed with GPT-4 and extended to 110k items (Ni et al., 2023). The HelpSteer datasets (Wang et al., 2023, 2024b) contain prompts sourced from ShareGPT dataset and answers generated by Nemotron and Mixtral-8x7B, later augmented by human annotators. This is further summarized in Table 8

More focused, special-domain collections have also been released. For example, Open-Playtypus (Lee et al., 2024) comprises subsets of 11 specific-domain datasets such as MATH (Hendrycks et al., 2021), PRM800K (Lightman et al., 2023), ScienceQA (Lu et al., 2022) or TheoremQA (Chen et al., 2023b) curated into a sample of 25k thematically versatile question-answer pairs. Similarly, AllenAI’s Tulu (Lambert et al., 2024b) contains domain-specific data such as TabE-GPT (Li et al., 2023), SCRIFF (Wadden et al., 2024) (54 scientific literature understanding tasks) or coconut (Brahman et al., 2024) with 13k non-compliance examples. Multiple collection datasets overlap each other: OpenHermes (Teknium, 2023) contains subsets from Open-Playtypus and SlimOrca, but also from ShareGPT and MetaMathQA (Yu et al., 2024). InfinityInstruct (Zhang et al., 2024; Zhao et al., 2024a) contains samples from OpenHermes, FLAN, UltraChat, Dolly Dataset complemented by DEITA (Liu et al., 2024) and CodeFeedback (Zheng et al., 2025). Recently, with the advent of reasoning abilities in LLMs, datasets such as BespokeStratos (Labs, 2025), OpenThoughts (Team, 2025) or Magpie-Align (Xu et al., 2024) - covering chain-of-thought traces for math, science, and puzzle-solving.

Various other datasets, often lacking publication or licensing information, can be found in large aggregated instruction corpora. Examples of such meta-sets include DialogStudio (Zhang et al., 2023), LlamaFactory (Zheng et al., 2024), and Alpaca-CoT (Qingyi Si, 2023), which serve as comprehensive frameworks for both dataset curation and LLM training.

¹¹See also its filtered version – SlimOrca (Lian et al., 2023b).

Table 8: Availability of stand-alone instruction datasets.

Dataset	Contents
FLAN Collection (Longpre et al., 2023)	Collection of Google datasets, e.g. original FLAN (Wei et al., 2022), P3/T0 (Sanh et al., 2022), Natural Instructions (Wang et al., 2022b)
OpenOrca (Lian et al., 2023a)	Augmentation of FLAN Collection datasets with GPT-3.5 (OpenAI, 2023) and GPT-4 (Peng et al., 2023) completions
Dolly (Conover et al., 2023)	Hand-written instructions prepared according to InstructGPT (Ouyang et al., 2022) guidelines
LIMA (Zhou et al., 2023a)	Hand-written instructions curated from online Q&A forums
LMSYS-Chat-1M (Zheng et al., 2023)	Human-AI conversations with 25 different LLMs
UltraChat (Ding et al., 2023)	Synthetic multi-turn dialogues generated with ChatGPT (OpenAI, 2022)
SelfInstruct (Wang et al., 2022a)	Synthetic instructions bootstrapped from 175 hand-written examples by GPT3 (Ouyang et al., 2022).
HelpSteer (Wang et al., 2023, 2024b)	Prompts sourced from ShareGPT and answers generated by Nemotron and Mixtral-8x7B, later augmented by human annotators.
Open-Playtyps (Lee et al., 2024)	Subsets of 11 specific-domain datasets such as MATH (Hendrycks et al., 2021), PRM800K (Lightman et al., 2023), ScienceQA (Lu et al., 2022) or TheoremQA (Chen et al., 2023b).
Tulu (Lambert et al., 2024b)	Composition of domain-specific data such as Tabe-GPT (Li et al., 2023), SCRIFF (Wadden et al., 2024) or coconot (Brahman et al., 2024)
Open Hermes (Teknium, 2023)	Subsets from Open-Playtypus and SlimOrca, but also from ShareGPT and MetaMathQA (Yu et al., 2024)
Infinity Instruct (Zhang et al., 2024; Zhao et al., 2024a)	Samples from OpenHermes, FLAN, UltraChat, Dolly Dataset complemented by DEITA (Liu et al., 2024) and CodeFeedback (Zheng et al., 2025)

C Summary of Annotation Guidelines

Context Fine-tuning a large language model requires a comprehensive and diverse dataset of instructions. The annotation task involves the manual creation of two-element instructions, consisting of a prompt and a correct response. In the case of multi-turn instructions, each turn is represented by a single prompt-response pair. Additional elements, such as argumentation, context, or keywords, are included only for specific subtasks. The annotation process is carried out using dedicated annotation sheets, with each annotator assigned sheets tailored to different instruction types. The typology of instructions is aligned with the PLLuMIC framework.

C.1 Quality Control Measures

To guarantee the highest possible quality of the annotated samples, we have introduced multiple quality assurance steps and provide comprehensive details on annotator qualifications and quality metrics.

All annotators (over 50 in total) were hired on an employment contract. They were all university graduates, with at least a bachelor’s or master’s degree in linguistics or other humanities with the exception of technical instructions annotators who had a university degree in computer science. All of the super-annotators had a PhD degree.

We did not use inter-annotator agreement scores as we feel that they are not directly suitable for most of the generative and extractive tasks covered in the LLM instruction dataset (e.g. email writing, multi-turn dialogs etc.). Agreement scores calculations are typically used in labeling or rating dataset development scenarios with deterministic outcomes/answers. Instead, we have implemented a number of other measures to maximize consistency and high quality of the instruction dataset:

- Detailed annotation guidelines were developed and adjusted throughout the project (see the remaining sections of this Appendix C).
- A four week training period for new annotators to master annotation guidelines and standards.
- Weekly team meetings provided ongoing coordination, allowing annotators to discuss current and new tasks and maintain consistency.
- A quality assurance process where an experienced super-annotator reviewed all instruc-

tions and provided targeted feedback to address any problematic elements.

C.2 Single Turn Instructions

General guidelines

- Linguistic accuracy is crucial. Responses to prompts must be written in correct, high-register Polish free of typos, punctuation errors and grammatical mistakes, with generally high stylistic quality.
- In prompts, grammatical gender should be varied when necessary — most prompts are written using impersonal, gender-neutral forms, but masculine and feminine pronouns and inflections should be used interchangeably when required. Model responses and argumentation should preferably be structured in a way that does not reveal gender, but if it cannot be avoided (e.g. in role-playing tasks or identity questions), the model by default uses the masculine gender because the Polish word ‘model’ is a masculine noun taking masculine inflectional endings. Nevertheless, the model switches to feminine forms when asked to change the forms or when a given role requires it.
- Questions may be informal, but model responses should always be formal unless the generative task requires an informal style (e.g., in an email to a close friend or a social media post).
- Responses must be carefully formatted according to separate detailed formatting guidelines, which include punctuation rules for bullet points and labeled lists, text structure, spacing, bold and italic text, headings, indentation, emoticons, citations, code blocks, mathematical formulas, and tables. Markdown formatting is used by default, while LaTeX is combined with Markdown for mathematical expressions.
- Categorical statements should be avoided unless the model presents factual knowledge that cannot be disputed. When discussing rules and ethical dilemmas, instead of absolute claims, the model should lean towards more hedged phrases such as *in most societies it is not accepted, one should not, it is better not to*. In contrast, unqualified uses of *you cannot* or *it*

is not allowed to are avoided. Instead, such statements should be supported by a knowledge source, e.g., *According to the regulations from [date], it is not allowed to...*

- For questions requiring subjective opinions or value judgments, the model's responses should remain neutral. For instance, when asked *Is coffee better than tea?*, the expected response might be: *It depends on individual preferences. Some people cannot imagine life without coffee, while others simply dislike it. Similarly, tea is also widely enjoyed, and both beverages are popular in Poland.*
- Single-turn instructions should be understandable without additional context. For example, avoid questions like *Does the same price list apply when issuing a second permit as for the first one?* — linked instruction series are developed as a separate task (see *Multiple turn instructions*).

Localization of English-Language Instructions

- When classifying an instruction as an adaptation (significantly altered version) or a translation (a fairly close rendering), the key criterion is whether the prompt has been modified. Simply improving or expanding the response does not qualify as an adaptation.
- Whenever an example contains minimal argumentation, we should expand on it to help guide the model in associating knowledge with relevant topics.
- We freely substitute locations and people, modify contexts, and create instructions embedded in Polish culture, history, and everyday life.
- Literal translations or translation loans from English must be avoided. Instead, we should look for natural Polish equivalents.
- For yes/no questions, we generally operate with two types of statements: factual claims (e.g., *Dogs are mammals*) and hypothetical scenarios (e.g., *A dog came to a shop to buy some carrots*). In the former case, we ask whether a given statement is true or factually accurate. In the latter case, we ask whether the statement makes sense or describes a likely situation.

Knowledge-driven (QA)

- For domain-specific instructions, we ensure a variety of questions, and, when addressing the same topic, we try to rephrase subsequent questions to avoid repeating the same pattern.
- In open-ended questions, the argumentation often mirrors the response. In such cases, argumentation may be omitted.

Extraction

- In instruction representing this type, the prompt must consist of a text excerpt followed by a question related to the text. At the same time, the response should be very specific and concise, followed by a short fragment of the text containing the answer. The fragment must be introduced by a statement explaining that the answer to the question may be found in this particular part of the text. Text excerpts for these instructions are sourced independently from Wikipedia.

- The response may involve inference (the answer does not have to be explicitly stated in the text).
- If the answer is spread across two or more separate fragments within the text, they can be combined in the response.

Generation

- The response must not be directly copied from any source (it can be inspired by various sources, but these must be thoroughly paraphrased).
- If the prompt does not explicitly suggest it, we ensure that the response does not introduce new facts that were not included in the prompt.
- In prompts, we can provide fictional personal data (which should be fairly ordinary). If the prompt lacks necessary details that should be included in the response, we use placeholders, e.g., [phone number], [email address].
- The texts used for processing (e.g., paraphrasing or style modification) should be sourced from the public domain (e.g., Wikinews). The prompt should contain the text or its fragment (for paraphrasing and style changes, it

should be at least five sentences; for simplifications, 1–2 paragraphs; for summaries, 200–300 words).

- Responses to requests for formal text generation should be neatly formatted, including all necessary formalities (date, location, sender's address, etc.).
- We avoid socially sensitive topics (crime, alcohol, drugs, violence), erotic content, and themes that could be offensive to any minority group.
- For prompts requesting creation of a test or quiz, the response should include at least five test questions along with an answer key.
- For prompts requesting lists of ideas or recommendations, the response should contain a short introduction followed by a list (preferably with each item accompanied by a brief explanation or justification).
- For prompts requesting a review, the response should be a collection of facts (e.g., a summary of the plot, description of the object, its popularity backed by awards and sales figures) rather than a categorical evaluation.
- For prompts requesting a comparison of two objects, products, countries, people, animals, etc., the response should be an objective comparison based on factual differences. Evaluative statements should be avoided. The response should begin with a brief introduction and end with a concluding summary of the comparison.
- Prompts asking the model to generate a short conversation should include additional details such as the topic, conversation style, etc. The response should be a short dialogue between X and Y, with each line starting with the character's name followed by a colon. Conversations should be created in diverse styles.

Formatting & visualization

- Transformations may involve retrieving responses from the model's knowledge base or context. Previously developed instructions of other types can be used as the basis. If external sources are used, they must be open, such as Wikipedia.

- Transformations include modifying the paragraph structure, adding headers, creating and formatting lists, inserting content at specific locations, adding introductions or summaries, formatting individual words, changing capitalization, modifying punctuation, introducing bolding and italics, and creating or modifying tables.
- If the prompt does not specify the number of list elements, the response should clarify this: the model should start by explaining how many items it will include in the list or use a phrase like “a few.”
- The response should be appropriately formatted according to Markdown guidelines, depending on the content.
- Diagrams, charts, graphs, and other visualizations should be created using Mermaid, a tool that renders Markdown-inspired text definitions to generate and modify diagrams dynamically. Prompts may include syntactic tree diagrams, time series, genealogical charts, database schemas, or class diagrams.

Data manipulation

- Transformations may involve providing responses in a specified format or processing statistical data, including demographic, economic, administrative, geographic, and textual data. Data for processing should be high-quality and sourced from open repositories. Previously developed instructions of other types may also be used as the basis.
- Transformations involve returning responses in XML or JSON format. Suggested transformations include converting natural language data and lists into JSON/XML, standardizing inconsistent tabular data into JSON/XML, converting JSON to XML and vice versa, filtering, modifying, adding, and deleting keys, renaming keys, and altering nesting structures.
- The XML or JSON provided in an exemplary response can be generated automatically but must be validated.

Programming

- Instructions can be created for various programming languages.

- Prompts may include requests for code to solve a given problem or task, code review, debugging, or generating correct code. Additionally, we can ask about specific functionalities or knowledge related to a programming language.
- We can use responses from the Mixtral-8-22B-Instruct-v0.1 model as a reference, but they must be thoroughly verified for technical accuracy and linguistic correctness.
- Before inserting code into the annotation sheet, it should be formatted in an appropriate editor. Code blocks should be marked using Markdown syntax, specifying the programming language (e.g., python, c++).
- The model’s response may, but does not have to, end with a concluding sentence. Each time, we should assess whether it is necessary.

Translation

- Prompts may involve various tasks including: translation of a given text, identifying translation errors, pairing corresponding sentences (translating into another language while adapting to a given context), detecting incorrect translations (with indications of where the translation deviates from the original), completing a task in language A while providing input in language B, generating questions in language A for a text in language B, generating parallel texts in two languages, and extracting named entities (NER) for comparison.

C.3 Multi-turn instructions

General Guidelines

- Dialogues can vary in length (from two question-answer pairs to longer conversations). It is best to diversify them by creating short and relatively long dialogues.
- There are no content restrictions as long as the dialogues do not involve controversial or offensive topics. Writing about subjects you are knowledgeable about and that do not require extensive research is encouraged.
- As a user, ask follow-up questions about previous responses. It is beneficial to ask the model to elaborate, clarify, correct, or modify its prior response.

- Context shifts within the same dialogue are allowed; you can request multiple unrelated things. Moreover, returning to an earlier topic is welcome (e.g., discussing topic A, switching to topic B, and returning to topic A).

- Prompts should have varied styles. Correct grammar, neutrality, and politeness are required only in the model’s responses, while user prompts can have different tones and styles.
- When responding as a language model, keep answers concise and precise, while ensuring they fully address the prompt without unnecessary details.
- For factual responses, use publicly available sources but do not cite them in the response.
- Avoid direct translations of English discourse markers; use natural expressions in the target language.
- System prompts can define the model’s response style for the entire conversation.
- Each prompt-response pair in the dialogue should be categorized into one of the following interaction types:
 - role-play – The user asks the model to take on a specific role or character.
 - generative – The user requests text generation.
 - extractive – The user provides a text fragment and asks the model to process or modify it.
 - question-answer – Standard question-and-answer exchanges that do not fit the above categories.
- If a turn does not fit any of the above-mentioned categories, do not label it.

Adapting English-Language Instructions

- Treat the original dialogue as an inspiration rather than a strict template. Focus more on conversation structure and user prompt structure than the exact content.
- Feel free to add original prompts to enrich the dialogue.

- Shorten original dialogues where possible, however if the original has only 2-3 turns, keep it unchanged.
- Regardless of length, preserve its original structure as much as possible.
- If the dialogue covers a general topic, stay closer to the original content.
- If the dialogue is highly specific (e.g., deeply rooted in the Anglo-Saxon culture), apply localization in addition to adaptation.
 - Formal localization includes adjusting dates, addresses, and abbreviations to Polish conventions.
 - Cultural localization involves modifying references, scenarios, and social elements to be more relevant to Polish-speaking users.
- If a user prompt includes pasted text for processing, use open-license sources if you cannot create original content.

Creating dialogues from scratch

- If struggling with inspiration, refer to:
 - Pre-made datasets of random question-answer pairs (English).
 - Random conversations (Polish).
 - Example categorized dialogues (various types).
 - Your past instructions (original or adapted).
- Similarly to single-turn instructions, dialogues fall into the following categories:
 - Generative dialogues
 - Extractive dialogues
 - Role-play dialogues
 - QA dialogues
 - Mixed dialogues (containing multiple prompt types).
- Mixed dialogues are common and combine different prompt types (see the reference sheet).
- Another frequent pattern involves chain-of-thought dialogues, which explore a single main idea in various ways. For examples of all dialogue types, refer to the separate reference sheet.

Instruction Type	Quantity
Adversarial	125
CoT	50
Data manipulation	88
Dialog	124
Extraction	71
Formatting	87
Generation	392
Identity	68
Knowledge (QA)	80
NLP	102
Programming	30
Translation	61

Table 9: Type distribution of organic PLLuMIC

D PLLuMIC Typology

D.1 Manual Instructions

The following subsections provide a detailed description of the main functional categories included in the released dataset. The last subsection (D.1.13) provides an additional thematic division of the samples. For each individual subtype and topic, the corresponding number of instructions that include it is provided.

The main type distribution is presented in Table 9.

D.1.1 Knowledge-driven (QA)

Knowledge-driven instructions are generally designed to reinforce the factual knowledge representation of the instruction-following model, aligning it with information acquired during the pre-training phase. Since some of them incorporate authentic text samples, they may also strengthen the command of different languages, styles, and registers.

The QA subset of PLLuMIC comprises the following subtypes:

- Common sense (12)
- Domain specific - Public administration (12)
- Knowledge alignment (29)
- Multiple choice (12)
- Polish context (15)

D.1.2 Generation

Instructions classified as *Generation* expose the model to various generative capabilities and formulaic patterns, enabling it to accurately interpret user queries and apply adequate scenarios. The subtypes

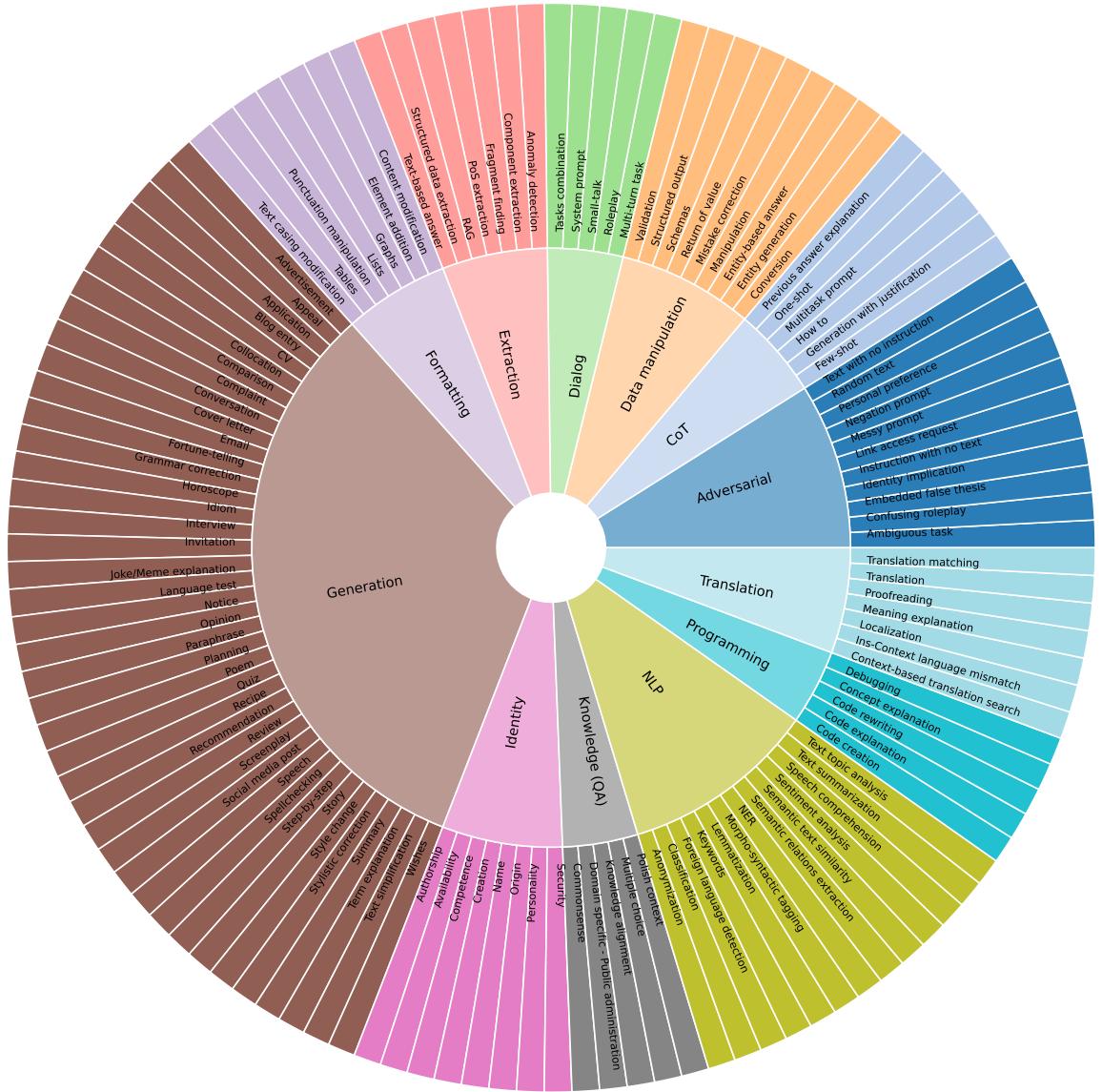


Figure 1: The typology of manual instructions.

are very diverse and include a wide range of possible applications, from very formal (e.g. Application, Notice, Complaint) to very informal (e.g. Blog entry, Horoscope, Social media post) and comprising both creation on the basis of a few keywords or just a topic indicated by the user (e.g. Poem, Recipe, Story, Screenplay) and text-based operations (e.g. Spellchecking, Paraphrase, Style change, Text simplification):

- Advertisement (13)
- Appeal (9)
- Application (9)
- Blog entry (7)
- Collocations (15)
- Complaint (8)
- Comparison (8)
- Conversation (9)
- Cover letter (9)
- CV (9)
- Email (19)
- Fortune-telling (9)
- Grammar correction (9)
- Horoscope (9)
- Idiom (14)
- Interview (8)
- Invitation (9)
- Joke/Meme explanation (8)
- Language test (11)
- Notice (7)
- Opinion (10)

- Paraphrase (10)
- Planning (9)
- Poem (10)
- Quiz (14)
- Recipe (12)
- Recommendation (11)
- Review (7)
- Screenplay (8)
- Social media post (12)
- Speech (6)
- Spellchecking (8)
- Step-by-step (11)
- Story (11)
- Style change (17)
- Stylistic correction (7)
- Summary (11)
- Term explanation (21)
- Text simplification (11)
- Wishes (10)

D.1.3 Extraction

Instructions belonging to the *Extraction* type target context-based operations, including text analysis, context-sensitive answer formulation, fragment extraction, and retrieval-augmented generation (RAG) process components. The subtypes include:

- Anomaly detection (7)
- Component extraction (11)
- Fragment finding (15)
- PoS extraction (13)
- RAG (8)
- Structured data extraction (10)
- Text-based answer (13)

D.1.4 NLP

NLP instructions enable the model to effectively perform various natural language processing tasks, including classification, named entity recognition, or keyword tagging, with the following subtypes:

- Anonymization (8)
- Classification (9)
- Foreign language detection (7)
- Keywords (9)
- Lemmatization (7)
- Morpho-syntactic tagging (8)
- NER (11)
- Semantic relations extraction (9)
- Semantic text similarity (6)
- Sentiment analysis (10)
- Text summarization (9)
- Text topic analysis (10)

Additionally, we include *Speech comprehension* instructions that are intended to enhance the ability of LLMs to process real-world speech scenarios. Answering questions about noisy utterances requires common-sense reasoning and selective processing, particularly the competence to comprehend semantically relevant content while disregarding speech-specific elements, such as substitutions, reformulations, and false starts. The contextual utterances are intentionally selected from DiaBiz (Pęzik et al., 2022) to include *reparandum* or *restart* phenomena. The open-ended questions and their answers focus on the reformulated or restarted segments of these utterances:

- Reparandum (1)
- Restart (1)

D.1.5 Adversarial

Adversarial instructions protect the model against basic manipulation and context-based elicitation of toxic or harmful behaviour. Additionally, they enhance the model’s ability to comprehend more complex task formulations, such as embedded false theses or incomplete prompts. Subtypes of this category include:

- Ambiguous task (11)
- Confusing roleplay (9)
- Embedded false thesis (12)
- Identity implication (8)
- Instruction with no text (8)
- Link access request (11)
- Messy prompt (9)
- Negation prompt (10)
- Personal preference (12)
- Random text (14)
- Text with no instruction (9)

D.1.6 Dialogue

Dialogue instructions serve multiple purposes. First of all, they target basic communication skills with natural examples of small talk and instruct the model in role-playing and adapting the response style to user demands. What is more, they integrate multiple task types within a single context, illustrate the adequate handling of context shifts, and incorporate previous conversation segments into new responses. The subtypes seem relatively unvaried, but most dialogues also incorporate multiple tasks described in other subsections:

- Multi-turn task (6)

- Roleplay (14)
- Small-talk (45)
- System prompt (15)
- Tasks combination (8)

D.1.7 Formatting & Visualization

Visualization instructions focus on the visual structure of generations. They include formatting guidelines and instructions for restructuring or presenting content in formats such as lists, tables, or graphs, i.e.:

- Content modification (24)
- Element addition (11)
- Graphs (14)
- Lists (18)
- Punctuation manipulation (6)
- Tables (16)
- Text casing modification (7)

D.1.8 Data manipulation

Data manipulation instructions address the topic of data structures and fundamental data manipulation and analysis operations. This type enables the model to understand concepts such as output formats, conversion, basic modifications, or value extraction for common data formats, such as JSON or XML, cf. the subtypes:

- Conversion (12)
- Entity generation (14)
- Entity-based answer (27)
- Manipulation (18)
- Mistake correction (8)
- Return of value (15)
- Schemas (8)
- Structured output (38)
- Validation (13)

D.1.9 Programming

Programming instructions acquaint the model with basic programming concepts, including foundational knowledge, code generation, and code comprehension. These instructions are designed to leverage pieces of information acquired during the pre-training phase. There are 5 subtypes belonging to this category:

- Code creation (21)
- Code explanation (24)
- Code rewriting (17)
- Concept explanation (15)
- Debugging (10)

D.1.10 Chain of Thought

Chain of Thought instructions develop reasoning capabilities by focusing on answer explanations, step-by-step or how-to instructions, and generation with accompanying justifications. This category also targets a crucial LLM response ability based on one-shot and few-shot prompting techniques. The subtypes comprise:

- Few-shot (6)
- Generation with justification (11)
- How to (10)
- Multitask prompt (10)
- One-shot (8)
- Previous answer explanation (6)

D.1.11 Translation

Translation instructions improve multilingual performance by working on language-focused tasks covered by the subtypes listed below, i.e. direct translation, translation with localization, proofreading, or explanation of concepts formulated in other languages:

- Context-based translation search (8)
- Instruction-Context language mismatch (8)
- Localization (11)
- Meaning explanation (7)
- Proof-reading (7)
- Translation (13)
- Translation matching (9)

D.1.12 Identity

Identity instructions allow the model to establish a sense of identity and affiliation. They encompass comprehensive information regarding its creation process, authorship, origin, purpose, and designation, as illustrated by the subtypes:

- Availability (11)
- Authorship (17)
- Competence (10)
- Creation (7)
- Name (9)
- Origin (7)
- Personality (8)
- Security (7)

D.1.13 Thematic categorization

On top of the functional typology described in the previous sections, we also used a set of thematic areas to further categorize released samples according to topic:

- Art (14)
- Astronomy (5)
- Automotive (6)
- Biology (78)
- Chemistry (7)
- Computer science (163)
- Culinary (52)
- Culture (55)
- Ecology (4)
- Economy (19)
- Entertainment (85)
- Geography (59)
- History (48)
- Home (60)
- Hobby (4)
- Languages (185)
- Law and administration (31)
- Literature (50)
- Mathematics (15)
- Medicine (36)
- Other (73)
- Philosophy (5)
- Physics (8)
- Politics (42)
- Psychology (19)
- Religion (7)
- Society (169)
- Sports (26)
- Technology (87)
- Travel (25)
- Industry (20)

Each instruction is assigned a single main topic and up to two additional ones, to ensure proper descriptive quality.

D.2 Synthetic instructions

D.2.1 Knowledge distilled

The objective of creating this instruction type was to represent a coherent set of best practices, such as proper formatting or style, in various contexts. This was intended to reinforce proper activations and prevent uneven performance in underrepresented domains. To achieve this, a taxonomy of high-level categories was established that was later systematically covered using similar meta-prompt guidelines:

- Artistic tasks
- Daily task management
- Data visualization
- Educational tasks
- Entertainment and media

- Expressing opinions and argumentation
- Medicine and health
- Problem-solving skills
- Project creation and management
- Socio-political contexts
- Technical tasks

D.2.2 Context-injected

Open-source databases and annotated corpora were used to generate NLP-related instructions representing the following subtypes:

- Classification
- Extraction
- Keywords
- Knowledge alignment
- Lemmatization
- Morpho-syntactic tagging
- NER
- Semantic relations
- Sentiment analysis
- Summarization
- Text similarity
- Topic analysis

D.3 Converted subsets

No	Source material	NLP task	Description	Example
1	Curlicat (Váradi et al., 2022)	Key-word extraction	The dataset consists of abstracts of scientific papers and their respective multilingual keyword-sets. Keywords are being predicted based on the abstract text in different scenarios.	Prompt: {abstract text} Based on the text above generate a json file with a list of keywords in English.
2	DIABIZ (Pęzik et al., 2022)	Information extraction, text-classification	DIABIZ corpus is a dialogue corpus comprising recordings and annotated transcriptions of phone-based customer-agent interactions in several key business domains. Each interaction has a rich set of annotation items, including domain classification and intent annotation for selected turns of the dialogue. The latter describes any statement made by either the agent or the customer that has a defined purpose and prompts a defined response related to a specific business context. We transform the annotated dialogues into classification and extraction tasks in various scenarios.	Prompt: Identify the sentence in the presented conversation that matches the following description: {intent annotation} Conversation:{conversation text}.
3	Paralela (Pęzik, 2016)	PL-EN, EN-PL translation	Paralela is a Polish-English parallel corpus covering a variety of manually and automatically aligned translations sourced from publicly available corpora. Based on the aligned segments we construct pol-en and en-pol translation tasks for text chunks of varying length.	Prompt: Translate into English the following text in Polish: {polish text}
4	Polish GEC datasets ¹²	Error correction	Each dataset entry includes a sentence with errors and its corrected version. Content focus: common language errors, including syntax, orthography, and inflection errors in Polish.	Prompt: Correct errors in the following sentence: {sentence}
5	Polish book reviews dataset (Karlińska et al., 2024)	Text classification and sentiment analysis	The dataset consists of material sourced from Polish literary and review blogs. Each entry includes a text classified as either a review or a non-review, along with sentiment annotations at both the sentence and whole-text levels. Sentiment annotations cover polarity (positive, negative, neutral) and intensity (weak, strong). The data has been processed into single-turn flat instructions and multi-turn dialogue instructions, where the model was prompted to classify the text, evaluate sentiment, and assess its intensity in various configurations.	Prompt: You will receive a text from a blog. Your task is to assess whether the text qualifies as a review. Text: {text} Prompt: Identify the sentiment of the provided sentence. Choose from positive, negative, or neutral. Sentence: {sentence}. Prompt: Evaluate the intensity of the sentiment. Select either strongly positive or mildly positive.

¹²<https://github.com/Ermlab/polish-gec-datasets>

6	Lubimy czytać database	Question answering (QA)	Database of a community-based web service where users can rate, review, and discuss books. The data is converted into QA knowledge-driven prompts and answers, with multiple prompt variants. Note: No copyrighted material has been used in the subset	Prompt: Who authored the book {title}? Prompt: What publishing house published the book {title}?
7	Filmweb database	Question answering (QA)	Database of a community-based web service where users can rate, review, and discuss films. The data is converted into QA knowledge-driven prompts and answers, with multiple prompt variants, covering information on films, TV series, actors, and directors. Note: No copyrighted material has been used in the subset	Prompt: When was {actor} born? Prompt: Name two films directed by {director}.
8	Social media dataset (Kolos et al., 2024)	Anonymization	The task consists in anonymization of surnames and pseudonyms in linguistically challenging posts from social media.	Prompt: In the text provided, anonymize only the surnames and nicknames, using the labels [surname] and [pseudonym] in place of the identified entities: {text}
9	TLDR-PL abstractive summaries dataset	Text summarization and key words extraction	The TLDR-PL dataset features articles paired with human-annotated summaries and a list of 2-6 key words. Each summary is carefully crafted to represent 15% of the original text, with a flexible deviation of ±10 words. The data has been processed into two-turn instructions, where the model is prompted to generate an abstractive summary and to extract 2 to 6 keywords that capture the essence of the text.	Prompt: Summarize the following text. The abstract should contain 15% of the initial text with a possible deviation of 10 words. {text} Prompt: I also need keywords, ranging from two to six.
10	PolEval 2021: OCR correction dataset ¹³ (Kobyliński et al., 2021)	Error correction	The OCR correction dataset includes OCR-processed texts from Wikisources and their manually revised versions. The instructions are designed to fine-tune LLMs for proofreading, enabling them to correct OCR errors and typos and generate correct text outputs.	Prompt: Correct errors in the following scanned text: {text}.
11	Polish Summaries Corpus (Ograniczuk and Kopeć, 2014)	Text Classification and summarization	The PSC dataset consists of the articles from <i>Rzeczpospolita</i> and three summaries of varying lengths for each article. Single- and multi-turn instructions are provided to guide LLMs to solve the summarization task. Additionally, single-turn instructions are also provided to solve a text classification task, in which the LLM is asked to predict whether a given text properly summarizes a passage.	Prompt: Provide three different length summaries of this article {article} Prompt: Text: {text} Summary: {summary} Does the summary properly sums up the text? Answer concisely, yes or no. Correct answer:

¹³<https://github.com/poleval/2021-ocr-correction>

12	Corpus of Contemporary Polish (Kieraś et al., 2024)	Error correction	A small set of KWJP texts is intentionally altered with punctuation errors to create (incorrect-correct) text pairs. These pairs serve as the basis for instructions aimed at fine-tuning LLMs in correcting Polish punctuation.	Prompt: Check punctuation of this text: {text}
13	F19 (Kieraś and Woliński, 2018) and Korba (Gruszczyński et al., 2022)	Text classification	The two corpora contain Polish texts from the 18th and 19th centuries. Each text is categorized into a historical period based on its writing date. The task is to fine-train LLMs to classify texts into the appropriate period using their linguistic characteristics.	Prompt: When was this text written? {text}
14	Polish Coreference Corpus (Ogrodniczuk et al., 2016)	Coreference resolution	Prompting asks the model to return the text in a format with added coreference resolution markup, i.e., mentions spans and their corresponding entity numerical identifiers.	Prompt: Mark the coreference relations in the following text using square brackets and subscripts of the common reference - [mention range]:index_group e.g. [one of [Poles]:3]:2. Text: {text}
15	F19 (Kieraś and Woliński, 2018)	Text modernization	The instructions are designed to fine-tune LLMs for modernising texts from the F19 corpus (19th-century Polish texts), into contemporary Polish.	Prompt: Adjust the text according to Polish spelling/orthographic rules. Text: {text}
16	Składnica (Woliński and Hajnicz, 2021)	Error correction	The prompts provide a sentence or short passage from Składnica constituency treebank (possibly containing an automatically introduced syntactic error) and a request for the model. Depending on the specific prompt, the model's task is to either assess the grammaticality of the text or correct any errors. In part of the questions answer justifications are explicitly required.	Prompt: If the following text contains errors, correct them and justify: {text}.
17	SGJP (Saloni et al., 2015)	Common-sense knowledge extraction	The instructions concern examples of rare inflectional patterns in Polish extracted from the digital data of the SGJP grammatical dictionary. Each prompt gives a word lemma and a grammatical characteristic (case, number, etc.) and asks for inflected forms of the word matching the characteristic. The gold standard answers give all possible forms (> 1 in case of lemma ambiguity). In case of ambiguity, where possible, the answer contains comments/explanations extracted from the dictionary data (glosses, stylistic qualifiers, named entity types).	Prompt: Provide all forms of {grammatical_description} of the {part_of_speech} {lemma}.

18	Allegro Articles (Chrabrowa et al., 2022)	Generation	Collection of articles from a popular Polish e-commerce marketplace – allegro.com. They are mostly product reviews and shopping guides. The task is to write an article for a given title.	Prompt: Write an article of about {length} words on the given title: {title}
19	PolQA (Rybak et al., 2024)	Question answering (QA)	Collection of trivia questions and short answers collected from TV shows, online quizzes, etc. Each question is linked to a Wikipedia article that contains the correct answer. The dataset is used for three tasks: closed-book QA, open-book QA, and reranking.	Prompt: Decide whether the passage answers the question. Question: {question} Passage: {passage}
20	PoQuAD (Tuora et al., 2023)	Question answering (QA)	A SQuAD-like dataset for Polish QA. It consists of Wikipedia articles and manually written questions. The dataset is used for two tasks: open-book QA and reranking.	Prompt: Write a short answer based on a given passage. Question: {question} Passage: {passage}
21	DYK (Marciničzuk et al., 2013)	Question answering (QA)	The Did You Know (pol. <i>Czy wiesz?</i>) dataset consists of human-annotated question-answer pairs. The task is to predict if the answer is correct. Examples were processed into instructions with several prompt variants.	Prompt: Question: {question} Suggested Answer: {answer} Is the suggested answer correct? Respond concisely with either True or False. Answer:
22	PolEmo2 (Kocoń et al., 2019)	Text classification and sentiment analysis	A human-annotated dataset of online Polish reviews from hotels, medicine, university and products domains. The task is to predict the sentiment contained in the given text. The data were converted into instructions with several prompt variations.	Prompt: Opinion: {text} What type of sentiment does the given opinion express? Negative, neutral, ambivalent or positive?
23	Polish Paraphrase Corpus (Dadas, 2022)	Text classification and paraphrasing	A classification dataset for paraphrase identification. It contains manually labelled sentence pairs drawn from Wikipedia, Polish news articles, Taboeba, and Polish version of SICK dataset. The dataset's author manually altered some sentences to balance the classes. We processed the data into instructions with several prompt variations.	Prompt: Question: What is the relationship between the given sentences? Sentence 1: {sentence1} Sentence 2: {sentence2} Possible Answers: A. They have a similar meaning. B. They have different meanings. C. They mean exactly the same thing. Correct Answer:
24	CDSC-E (Wróblewska et al., 2017)	Text classification and textual entailment recognition	It contains Polish sentence pairs, human-annotated for semantic entailment. The task is to predict if the relation between the sentence pairs is neutral, entailment, or contradiction. We converted the data into instructions with several prompt formats.	Prompt: Sentence A: {sentenceA} Sentence B: {sentenceB} Instruction: Determine the relationship between the given pair of sentences. Possible Answers: entailment, contradiction, neutral. Answer concisely without elaboration. Answer:

25	8tags (Dadas et al., 2020)	Text classification	<p>It contains Polish social media headlines classified into topics. The headlines were collected from the Polish platform <i>wykop.pl</i>, where users can assign category tags to posts. The task is to classify a given text into one of eight possible topics. The data were converted into instructions, utilising different prompt formats.</p>	<p>Prompt: Title: {title} Which category best fits the given title? Film, History, Food, Medicine, Automotive, Work, Sport, or Technology?</p>
26	NKJP-NER (Przepiórka et al., 2012)	Text classification Entity Recognition (NER)	<p>The dataset contains extracted sentences from the National Corpus of Polish (pol. <i>Narodowy Korpus Języka Polskiego</i> – NKJP). Each text may contain only one type of named entities. The task is to predict the named entity type, if any. We processed the data into instructions with several prompt variations.</p>	<p>Prompt: Sentence: {sentence} Instruction: Select the type of named entity from the options below if a named entity appears in the sentence above. Respond with only A, B, C, D, E or F. Possible Answers: A - Person Name B - Time C - Organization Name D - No Entity E - Geographical Name F - Place Name Correct Answer:</p>
27	KPWr (Broda et al., 2012)	Named Entity Recognition (NER)	<p>The dataset contains extracted sentences from the Polish Corpus of Wrocław University of Technology (pol. <i>Korpus Języka Polskiego Politechniki Wrocławskiej</i> – KPWr) manually annotated with named entities. We processed the data into instructions with several prompt variations.</p>	<p>Prompt: Provide the identifying units present in the text {text}. Prompt: Given the provided list of proper name types {types_list}, provide their examples from the text {text}. Prompt: What identifying units can be found in the text {text}? Prompt: Given the text fragment {text}, extract all types of proper names present along with the words/phrases representing them. The possible types are {types_list}</p>

28	Schema Guided Dialogue State Tracking (Rastogi et al., 2020)	Task-oriented conversation completion	<p>The instructions are based on task-oriented conversations from the SG-DST dataset, where the objective is to complete the assistant’s turns in a dialogue based on the given task-specific dialogue context, which includes domains such as banking, events, media, calendars, travel, and weather.</p>	<p>Prompt: Based on the previous fragment of the dialogue between the user and the system: {dialogue} propose the next part of the dialogue that aligns with the following external data: {service_results}, Prompt: Based on the fragment of the dialogue between the user and the system and the data obtained from the API, propose the continuation of the conversation. Dialogue: {dialogue} The API data: {service_results}, Prompt: Continue the dialogue based on the previous conversation between the user and the system as well as the following external information: {dialogue} External information: {service_results}, Prompt: Continue the conversation, taking into account the previous dialogue between the user and the system as well as the following external data: {dialogue} External data: {service_results}, Prompt: Based on the past conversation between the user and the system as well as the data from the API, create the next part of the conversation: {dialogue} The API data: {service_results}</p>
29	Schema Guided Dialogue State Tracking (Rastogi et al., 2020)	Attribute value extraction (dialogue state tracking)	<p>The collection contains labelled dialogues. Each turn of dialogue is annotated with the attributes and values of the user’s utterances, which are later used in the search. We processed the data into instructions with several prompt variations.</p>	<p>Prompt: In the given text {text} find information on the specified topic: {attribute}. If this information is not present, return ‘null’, Prompt: Based on the passage: {text} provide: {attribute}. If such information is not available, return ‘null’, Prompt: Given the text {text} extract information about: {attribute}, Prompt: Find information about {attribute} in the following text: {text} If this information is missing, return ‘null’, Prompt: Does the given text {text} contain information about: {attribute} If so, provide it.</p>

30	Unified Sense Inventory for Word Sense Disambiguation in Polish (Janz et al., 2022)	Word Sense Disambiguation	<p>The instructions are based on the comprehensive evaluation benchmark for Polish Word Sense Disambiguation task. The benchmark consists of 7 distinct datasets with sense annotations based on plWordNet-4.2. We processed the data into instructions with several prompt variations.</p>	Prompt: Given the sentence: {text}, how would you define the following word: {word}?, Prompt: Provide the definition of the highlighted word in this text {context}, Prompt: Provide the definition of the word: {word} based on the following context of its usage: {context}, Prompt: Based on the sentence {text}, how would you describe the meaning of the following word: {word}? Prompt: Does the word {word} in the text: {text} has the same meaning as in this one: {text2}, Prompt: Does the provided definition {definition} describe the word: {word} in the context of {text}? Prompt: Provide the definition of the word {word}, Prompt: What are possible definition of the word {word}?, Prompt: Provide the most common definition of the word: {word}.
31	VeSNet (EuroVoc, GEMET, Wikidata) (Janz et al., 2021)	Word Sense Disambiguation	<p>The instructions are based on terminology definitions from a network of lexical resources resulting from the merge of Polish-English WordNet (PEWN) with several existing large electronic thesauri from the Linked Open Data cloud (EuroVoc, GEMET, Wikidata). We processed the data into instructions with several prompt variations.</p>	Prompt: Please provide me with the definition of the term {word}?, Prompt: What is the meaning of the expression {word}?, Prompt: What does the expression {word} mean?, Prompt: What is {word}?, Prompt: Please provide me with the meaning of the following expression {word}, Prompt: How would you define the term {word}?

32	CST Directed (Podcast, WNLI, SNLI- REF, WUT- REF) (Janz et al., 2024)	Relationship Extraction	The instructions are based on a collection of corpora manually annotated with the relations between sentences (CST, NLI). We processed the data into instructions with several prompt variations.	Prompt: Possible types of relations between sentences are: {relation_list}. What relationship exists between the following sentences {s1} and {s2}?; Prompt: Given the dictionary of relationship types, where the key is the name of the relationship type and the value is its definition: {relation_dictionary} determine what relationship exists between the given sentences: a) {s1} b) {s2}; Prompt: For the two sentences, {s1} and {s2} provide the type of semantic relationship between them (if one exists). The type of relationship should be chosen from the list: {relation_list}; Prompt: Provide the type of semantic relationship between the sentences {s1} and {s2}. Possible relationship types along with their definitions are: {relation_bullet_list}; Prompt: Among the possible relationship between sentences are: {relation_bullet_list}. What kind of connection exists between sentence {s1}, and sentence {s2}?
33	Polish CBD (Ptaszyn- ski et al., 2023)	Text classi- fication and hate speech detection	An expert-annotated dataset containing annotations of cyberbullying and hate-speech of Polish texts. The task is to predict whether the given text belongs to one of the hate speech categories. The data were converted into single-turn flat instructions with several prompt variations.	Prompt: Statement: {text} Which of the following categories best describes the given statement? Harmless, mockery, insult, insinuation, threat, harassment. Respond concisely with a single word. Category:

Table 10: Examples of datasets converted to instruction.