

Evaluation of Transfer Learning for Polish with a text-to-text model

Anonymous ACL submission

Abstract

We present polT - a general purpose text-to-text model for Polish that can be fine-tuned on a variety of Natural Language Processing (NLP) tasks with a single training objective. Unsupervised denoising pre-training is performed efficiently by initializing the model weights with multi-lingual T5 (mT5) counterpart. We evaluate performance of polT, mT5, Polish BART (plBART) and Polish GPT-2 (papuGaPT2) on diverse downstream tasks such as: text-to-text KLEJ benchmark, en-pl machine translation, question answering and summarization. The polT scores top on all of these tasks except summarization where plBART is best. In general (except summarization), the larger the model the better the results. The encoder-decoder architectures prove to be better than decoder-only equivalent. Additionally, since summarization and question answering lack benchmark datasets for Polish language we describe in detail their construction and will make them publicly available.

1 Introduction

Recent years have brought significant progress in both natural language understanding (NLU) and natural language generation (NLG). Transformer architecture enabled efficient training of large-scale language models (Radford et al., 2019) and language understanding models (Devlin et al., 2019; Liu et al., 2019). On the other hand, transfer learning, which has been used in representation learning for years (Mikolov et al., 2013; Pennington et al., 2014; Devlin et al., 2019), has finally been successfully applied to text-to-text problem as well (Raffel et al., 2020; Lewis et al., 2020). The text-to-text framework takes text as input and produces new text as output. This unified view enables the use of same architecture, training procedure and decoding process for many NLP tasks such as classification, machine translation, summarization and question answering, to name a few. Raffel et al. (2020)

demonstrated that simplicity of this approach combined with scale can achieve state-of-the-art results on many benchmark datasets. Their T5 model was available only for English language, but more recently pre-trained multi-lingual architectures (Liu et al., 2020; Xue et al., 2021; Nagoudi et al., 2021a) and non-English counterparts (Carmo et al., 2020; Husein, 2018) were released. It is now known that models targeted for specific language perform better than multi-lingual one (Martin et al., 2020; Le et al., 2020; Chan et al., 2020; Mroczkowski et al., 2021; Virtanen et al., 2019; Nagoudi et al., 2021b). Moreover, specialized architectures are typically smaller due to significantly reduced vocabulary size and they can be trained efficiently via transfer from multi-lingual checkpoints (Arhipov et al., 2019; Mroczkowski et al., 2021). There were some attempts to pre-train Transformer-based models for generating Polish namely plBART¹ and papuGaPT2², but they lack detailed description and evaluation on benchmark datasets.

Our contributions are:

1. demonstration of the efficiency of pre-training procedure for transferring knowledge from multi-lingual to monolingual text-to-text models based on work by Arhipov et al. (2019); Mroczkowski et al. (2021),
2. comprehensive evaluation of text-to-text models on diverse tasks in Polish, such as text-to-text KLEJ benchmark (Rybak et al., 2020), machine translation, question answering and summarization
3. release of polT³ – a T5-based model for the Polish language, which achieves the best⁴ re-

¹<https://github.com/sdadas/polish-nlp-resources#bart>

²<https://huggingface.co/flax-community/papuGaPT2>

³<https://proper.download.site/model>

⁴among text-to-text models

076	sults on KLEJ benchmark, machine translation and question answering and second best results in summarization.	
077		
078		
079	4. construction of benchmark datasets in Polish domain for question answering and summarization.	
080		
081		
082	2 Polish T5 model	
083	There was no large-scale text-to-text model for Polish to benchmark transfer learning methods, so we trained 3 versions of polT model (small, base, and large) with T5 architecture. We prepared a Polish language tokenizer, initialized the model weights based on mT5, and pre-trained it on Polish corpora.	
084		
085		
086		
087		
088		
089	2.1 Architecture	
090	We followed the original mT5 (Xue et al., 2021) architectures with encoder and decoder stacks for all trained model sizes. We trained 3 variants of the T5 model: small (8 layers, 6 attention heads, hidden dimension 1024), base (12 layers, 12 attention heads, hidden dimension 768) and large (24 layers, 16 attention heads, hidden dimension 1024).	
091		
092		
093		
094		
095		
096		
097	2.2 Initialization	
098	Parameters for all layers except word embeddings were copied directly from public mT5 models (Xue et al., 2021). The difference in tokenizers vocabulary was addressed as described in (Arkipov et al., 2019; Mroczkowski et al., 2021). Embedding weights for tokens appearing in both dictionaries were copied directly. Otherwise, if target vocabulary token was missing in source vocabulary it was split into sub-tokens and the embedding was the average of sub-tokens embeddings.	
099		
100		
101		
102		
103		
104		
105		
106		
107		
108	2.3 Pre-training datasets	
109	Denosing pre-training was performed on a weighted mixture ⁵ of the following corpora: Polish Wikipedia, National Corpus of Polish (NKJP, Przepiórkowski et al. 2011), Wolne Lektury ⁶ , Polish Open Subtitles ⁷ and Common Crawl of Polish sites. Each dataset is described in Mroczkowski et al. (2021).	
110		
111		
112		
113		
114		
115		
		2.4 Tokenizer
		116
		The training dataset was tokenized into subword units using a sentencepiece unigram model (Kudo, 2018) with vocabulary size of 50k tokens. Unigram model was trained ⁸ on the most representative parts of our corpus, i.e subset (quarter of the whole corpus) of the NKJP, and Polish Wikipedia.
		117
		118
		119
		120
		121
		122
		2.5 Pre-training procedure
		123
		For unsupervised denosing pre-training we used the original T5 training scripts ⁹ , which are based on the Mesh TensorFlow framework (Shazeer et al., 2018). We used default options for span-corruption objective with a mean span length of 3 and corruption rate of 15% (Raffel et al., 2020).
		124
		125
		126
		127
		128
		129
		polT models were trained using AdaFactor optimizer (Shazeer and Stern, 2018) with constant warmup and inverse square root learning rate schedule (Raffel et al., 2020, section 3.1.2). We use the peak learning rate 5e-3, batch of 1048576 tokens, 0% dropout, and trained for 50k steps with 1024 warmup steps on a single preemptible TPU v3.
		130
		131
		132
		133
		134
		135
		136
		The learning rate was found in simple hypertuning procedure. We evaluated polT-base model snapshots saved after 10k training steps for pre-training learning rates in {1e-3, 2e-3, 5e-3}. For performance evaluation we used KLEJ benchmark tasks equipped with validation sets, namely: AR, NKJP-NER, PolEmo2.0-OUT, CDSC-E. Best model training was continued up to 50k steps.
		137
		138
		139
		140
		141
		142
		143
		144
		3 Downstream tasks
		145
		polT is a generic text-to-text model for Polish. It can be fine-tuned on many NLP tasks. We evaluated it by fine-tuning on text-to-text reformulation of KLEJ benchmark, machine translation, question answering and summarization. The same fine-tuning scheme was used for mT5 and other text-to-text or NLG models for Polish: plBART ¹⁰ and papuGaPT2 ¹¹ , in order to compare their performance. Trainable number of parameters and vocabulary size for each architecture is show in Table 1.
		146
		147
		148
		149
		150
		151
		152
		153
		154
		155
		⁸ https://github.com/google/sentencepiece
		⁹ https://github.com/google-research/text-to-text-transfer-transformer
		¹⁰ https://github.com/sdadas/Polish-nlp-resources#bart
		¹¹ https://huggingface.co/flax-community/papuGaPT2
	⁵ Probability of selecting given corpus in the mixture was proportional to the size of the dataset.	
	⁶ https://wolnelektury.pl/	
	⁷ https://www.opensubtitles.org/pl	

Model		# Params	Vocab. size
Small models	mT5	300M	250k
	polT	95M	50k
Base models	mT5	582M	250k
	polT	275M	50k
	plBART	139M	50k
	papuGaPT2	124M	50k
Large models	mT5	1.23B	250k
	polT	820M	50k

Table 1: Summary of total number of trainable parameters for pre-trained models for Polish.

3.1 General language understanding

We verified the quality of the assessed models in terms of general language understanding with KLEJ benchmark (Rybak et al., 2020). To compare the generative models, we used all of the given tasks, except for the CDSC-R regression. We formed the tasks into a text-to-text format (see Appendix A) proposed in (Raffel et al., 2020). During fine-tuning phase, the model was provided descriptive input containing specific prefix and the textual features with descriptive labels¹². The training objective was to generate greedily a specific label consisting of up to several tokens. Targets were generated over the whole vocabulary and only an exact match was treated as the correct answer.

3.1.1 Experiments

To ensure a fair comparison between models, we followed a simple experimental setup. The learning rate was the only hyperparameter we tuned (see Appendix A.3). By default, we trained the models for 1600 steps with a batch size of 64. Reported results are the median of 7 runs for small and base architectures and the median of 3 runs for the large one.

3.1.2 Results

The results for KLEJ benchmark (text-to-text format) are in Table 2. In general, encoder-decoder models dedicated to the Polish language perform significantly better at tasks in the KLEJ benchmark than the multi-lingual counterparts. The most striking performance difference occurs for the base version of the model on the PolEmo2.0-OUT task.

¹²We didn't explore multi-task learning in this setup with task-specific prompts

The difference between polT-base and mT5-base is over 9 pp. The difference on PolEmo2.0-OUT of more than 2pp is also visible for the other sizes. We hypothesize that domain adaptation necessary to solve this task reveals better Polish language understanding for monolingual T5. The superiority of the dedicated model is also visible in the PSC paraphrase identification task. Especially for the T5-small and T5-base models, the difference is 4.6pp and 3.9pp, respectively. The gap closes to 1pp for the large version of a T5 model. On the other hand, the multi-lingual mT5-large model obtained the best result on the NKJP-NER task, which shows that the task formulated in a classification form is more straightforward than its original version.

We observe a consistent increase in results with the size of the polT model in each task. The best encoder-decoder model is polT-large. It performs the best in 6 out of the 8 evaluated tasks from the KLEJ benchmark. However, it is not as good as the best Polish encoder, namely HerBERT-large. The difference between these architectures is, on average, 1pp. We observe the most significant performance gap of 6.8pp when moving from polT-small to polT-base. We found the biggest difference of around 20pp on *PSC* and *Czy wiesz?* tasks.

Regarding the smaller architectures, the best text-to-text model is polT-base. It is worth emphasizing that almost 2 times smaller plBART is very competitive. It is the best on 3 out of 8 tasks among the smaller models. The decoder only papuGaPT2 achieves on average results on par with polT-small but is the worst at dealing with the highly unbalanced tasks, *CBD* and *Czy wiesz?*.

We emphasize the variability of the results for the two *CBD* and *Czy wiesz?* tasks. The difference between the results for different architectures is as high as 30pp.

To sum up, the models dedicated to the Polish language are the best from the perspective of the KLEJ benchmark. Among them, the T5 models were the best.

3.2 Summarization

Summarization is a task of writing a shorter text that has all the key information of the longer text. Summaries can vary in length and be extractive (a selection of passages from the original text) or abstractive (written from scratch, though passages from the original text are not forbidden).

Model		AVG	NKJP-NER	CBD	Czy wiesz?	PolEmo2.0-IN	PolEmo2.0-OUT	AR	PSC	CDSC-E
Small models	mT5	74.3 ± 2.6	88.7	56.6	42.7	86.0	73.2	84.6	70.8	91.9
	polT	<u>76.8 ± 1.8</u>	<u>92.4</u>	<u>60.4</u>	<u>44.7</u>	<u>88.0</u>	<u>76.6</u>	<u>84.7</u>	<u>75.4</u>	<u>92.5</u>
Base models	mT5	82.4 ± 0.0	92.8	<u>65.6</u>	67.8	88.4	70.2	87.6	93.3	93.1
	papuGaPT2	76.5 ± 0.9	90.7	33.3	49.8	89.2	76.2	86.2	95.3	91.3
	plBART	81.9 ± 0.5	93.1	47.6	<u>68.5</u>	89.5	77.9	<u>88.0</u>	<u>97.9</u>	93.2
	polT	<u>83.4 ± 0.3</u>	<u>93.6</u>	62.3	63.8	<u>90.0</u>	<u>79.3</u>	87.8	97.2	<u>93.4</u>
	HerBERT	84.7 ± 0.4	94.5	66.4	64.3	90.9	80.4	87.7	98.9	94.5
Large models	mT5	84.9 ± 6.7	94.8	62.3	69.9	91.4	80.3	88.8	97.9	93.5
	polT	86.4 ± 0.3	94.5	70.0	69.4	92.9	82.6	88.9	98.9	94.0
	HerBERT	87.5 ± 0.2	96.4	72.0	75.8	92.2	81.8	89.1	98.9	94.1

Table 2: Evaluation results on KLEJ benchmark. AVG is the average score across all tasks. Scores are reported for test set and correspond to median values across 7 runs for small and base models and 3 runs for large version of models. The best scores for text-to-text models within each group are underlined, the best overall are in bold. Scores for the HerBERT model were taken from (Mroczkowski et al., 2021) and evaluated in a standard setting.

3.2.1 Datasets

Allegro Articles (AA) Collection of over 33k articles from Allegro.pl¹³ - a popular e-commerce marketplace. They contain among others product reviews and shopping guides. Every article contains a title, lead (opening paragraph) and a body of text. We prepared 2 different summarization tasks: (1) *body2lead* - generate lead from a body of the article (2) *body+lead2title* - generate a title from a full article (lead and body). The tasks are entirely abstractive and differ in source and target length. The details of dataset construction and statistics are in Appendix B.1.

Polish Summaries Corpus (PSC) Collection of summaries for 569 news articles created by human annotators (Ogrodniczuk and Kopeć, 2014). The annotators created 5 extractive summaries (each by different author) for each article by choosing approximately 5%, 10%, or 20% of the original text. The subset of 154 articles was also supplemented with additional 5 abstractive summaries each, i.e., not created from the fragments of the original article. We constructed 3 summarization subsets: (1) *whole* - all summaries and articles (2)

extract - only extractive summaries and (3) *abstract* - only abstractive summaries. The details of dataset construction and statistics are in Appendix B.2.

3.2.2 Experiments

For a fair comparison between models, we designed a common experimental setup, where the learning rate is the only hyperparameter (details in Appendix B.4). By default, each model was trained for 10 epochs with a batch size of 32 and Adam optimizer (Kingma and Ba, 2014). The only exceptions were large versions of T5 models (all tasks) and papuGaPT2 (only on PSC task), which were trained for 4 epochs. Each experiment, including hyperparameter search, was repeated with 3 seed.

Trimming inputs Summarization models by its nature deal with long source texts and shorter target texts. Unfortunately, we had to substantially trim source texts due to model and memory limitations. Usually, the limit was 1024 tokens, but sometimes it was even lower. The PSC tasks (news articles domain) were affected the most. None of the models ever saw 100% of source text during training, the percentage varied between 26-52% depending on the model. Details can be found in Appendix B.3.

¹³<https://allegro.pl/artykuly>

3.2.3 Results

Results are shown in Table 3 which contains arithmetic mean of (f-measure) ROUGE-1, ROUGE-2 and ROUGE-L (Lin, 2004)¹⁴ for each model and task. More detailed metrics for each tasks can be found in Appendix D (Tables 10-14). In all cases, the best performing model was plBART, with one exception of AA *body+lead2title* task in which polT-large outperformed the others. Usually, some version of polT (small, base, or large) was the second after plBART. Considering the mean ROUGE averaged over all summarization tasks, plBART is the winner with an average ROUGE 0.293 ± 0.003 and the second is polT-small with the result 0.255 ± 0.062 . polT models are always better than mT5 and papuGaPT2. We hypothesize that the pre-training procedure that makes BART biased towards copying the source is a good heuristic in the news summarization task. Furthermore, we observed that large architectures are unstable, and their performance degrades.

Baseline To further evaluate the results, we computed a baseline: copying 3 leading sentences from source as candidate target summarization, following (Gliwa et al., 2019; See et al., 2017). Then we computed a second baseline *adaptive n*: copying n leading sentences from source that best match the average target length determined on the train set. On average $n = 11$ for PSC tasks, $n = 3$ for AA *body2lead* and $n = 1$ for AA *body+lead2title*. Adaptive baseline was much better than fixed $n = 3$.

The best model was usually 2-3pp above the adaptive baseline except for the PSC abstract task (0.1pp below the baseline) and AA *body+lead2title* (32pp above the baseline). However, not all the models outperformed the baseline. PSC tasks have strong baselines. In the news domain, the essential information is at the beginning of the article.

Upper bound estimates PSC dataset contains multiple summaries by different authors per source text. We evaluate human performance as ROUGE metrics of different summaries for the same source text. This upper bound includes both variability in the content and the length of the summary. The best

¹⁴We used native Python implementation of ROUGE score <https://github.com/google-research/google-research/tree/master/ROUGE> and replaced default tokenizer with nltk <http://www.nltk.org/> punkt word tokenizer for Polish language.

models are halfway between the baseline and upper bound for the PSC whole and PSC extract task. Interestingly, for PSC abstract, all: the baseline, the best model and the upper bound are roughly the same. This shows that ROUGE is not a perfect metric, in particular for abstractive summaries.

3.2.4 Discussion

The PSC dataset contains summaries of different lengths that may degrade performance since the model may be uncertain about the length of the summary. One idea to lift that ambiguity is to add a prefix that would specify the target summary length. Furthermore, different source contexts seen by each model made interpretation of the results more challenging. It was not clear whether the drop in metrics was due to model performance or less information. If the source was trimmed down significantly and the target was not affected, then the whole task turns into generation rather than summarization. On the other hand, copy baselines performed very well on PSC, and indeed, every model saw the leading sentences from the source during training. Almost none of them, except plBART, were able to learn valuable summaries.

After exploration, we found out that, on average, summaries generated by encoder-decoder architectures (polT, mT5, plBART) are shorter than targets on PSC and AA tasks. The papuGaPT2 model generated summaries comparable in length to targets. We also spotted some anomalies. The mT5-large on PSC whole task generated summaries of maximal length, much longer than the target. The issue was not present on PSC abstract task. Furthermore, the PSC abstract was particularly difficult for polT-base which generated extremely short summaries (samples can be found in Appendix G, H, I).

3.3 Question answering

Question answering (QA) systems enable users to automatically obtain accurate answers to questions asked in natural language. We distinguish reading comprehension QA, in which the system, apart from the question, also receives a passage, which may contain the correct answer, and open-domain QA, in which the system receives only the question itself and answers only using the previously collected knowledge.

3.3.1 Datasets

There is no standard dataset for training question answering system for the Polish language, so we

Model		ROUGE AVG	AA body2lead	AA body+Head2title	PSC whole	PSC extract	PSC abstract
Baselines	lead n=3 source sentences	17.0	12.4	6.8	22.0	23.4	20.3
	lead n (adaptive) source sentences	21.6	12.4	7.9	30.3	31.7	25.6
Upper bounds	human performance	-	-	-	<u>34.3</u>	<u>39.0</u>	25.4
Small models	mT5	<u>23.3 ± 0.4</u>	<u>13.0</u>	<u>34.2</u>	23.3	25.0	21.2
	polT	<u>25.5 ± 6.2</u>	<u>14.3</u>	<u>36.3</u>	<u>30.5</u>	23.2	23.2
Base models	papuGaPT2	15.0 ± 0.2	12.1	14.0	16.6	17.5	14.8
	mT5	20.2 ± 1.3	14.1	36.1	20.6	21.2	8.8
	plBART	29.3 ± 0.3	15.6	38.3	32.6	34.3	25.5
	polT	<u>23.1 ± 3.3</u>	<u>14.9</u>	<u>38.9</u>	25.2	24.7	11.7
Large models	mT5	15.3 ± 1.1	10.9	<u>33.5</u>	12.0	10.9	9.1
	polT	18.9 ± 1.8	12.1	39.4	17.5	15.1	10.6

Table 3: Evaluation results on all summarization tasks (test split, mean values across 3 runs). For simplicity only arithmetic mean ROUGE (f-measure) is reported, ROUGE-1, ROUGE-2 and ROUGE-L are in Appendix D (Tables 10-14). The best scores are in bold, results above the baseline are underlined. ROUGE AVG is arithmetic mean of mean ROUGE of all tasks per model.

combined several resources to create a QA task specifically for this evaluation.

MKQA Multi-lingual Knowledge Questions & Answers (MKQA) (Longpre et al., 2020) contains 10,000 queries sampled from the Google Natural Questions dataset (Kwiatkowski et al., 2019), manually translated from English into 26 typologically diverse languages, including Polish.

After manually inspecting a sample of the data, we removed many of the questions due to specific domain (TV series, movies, etc.), open-endedness ("why?" and "how?" questions), lacking the answer or the answer depending on when question is asked (e.g. "Who won World Cup this year?"). After filtration, we left 1875 useful questions.

Jeden z Dziesięciu is a Polish game show broadcast on Polish Television. The participants answer the host's questions from various domains. We gathered 1004 question-answer pairs¹⁵.

Poleval is an evaluation campaign for NLP tools for Polish. The 2021 edition contains the question

answering task¹⁶ with a dataset of 6000 questions and answers. We used the validation set (1000 questions) for training and the test-A set (2500 questions) for test. The test-B set was not released at the time of conducting the experiments.

Final dataset dataset combines the aforementioned datasets, which resulted in the training set of 3879 questions and the test set of 2500 questions.

3.3.2 Tasks

In the task of open-domain question answering we aim to assess the factual knowledge stored in pre-trained models. We expect that the evaluated models already have the required knowledge to answer the question and we use fine-tuning procedure only to guide the models on how to generate the answer (Prager, 2006; Roberts et al., 2020).

The next question answering task is similar to the well-known reading comprehension task: besides the question, the model additionally takes a passage of text which may contain the answer (Zeng et al., 2020). Since none of the used datasets contains

¹⁵We parsed <http://tvturnieje.blogspot.com/p/jeden-z-dziesieciu.html>

¹⁶<https://github.com/poleval/2021-question-answering>

such passages, we retrieved them on our own in the following way.

First, we manually annotated 10k question-passage pairs with an information if the passage contains the correct answer. The source of candidates for positives pairs were severalfold. We started with a simple Bag-of-Words retriever as well as a more sophisticated model such as Universal Sentence Encoder (Yang et al., 2019). Next, we train a neural retriever based on HerBERT-base model (Mroczkowski et al., 2021) and annotated the retrieved passages.

Overall out of 10k annotated pairs, there were 2215 matching passages for 1347 unique questions¹⁷. We combined them with "Czy wiesz?" dataset (Marcinczuk et al., 2013; Rybak et al., 2020) to train a neural retriever based on HerBERT-base model.

Finally, to prepare a dataset for the task, we used this retriever to find the top 10 best Wikipedia passages for all questions in our dataset. The goal of the task was to generate the answer based on a question and retrieved passages.

3.3.3 Experiments

Using the aforementioned training set, we finetuned the following models: polT-base, polT-wiki¹⁸, polT-large, mT5-base, plBART, papuGaPT2 and T5 with random weights (T5-random) as a reference point. We trained the models for 50 epochs with Adam (Kingma and Ba, 2014) optimizer except for polT-large, which we trained for 30 epochs. We performed the best learning rate search for each model in the set of {1e-2, 5e-3, 1e-3, 5e-4} with 1 seed value. We used a single NVIDIA A100 GPU for all the models. The best results for each model are presented in Table 4. In both tasks the best model was polT-large. In base size polT-wiki performs the best in both tasks.

Models trained on open-ended task learned mainly to answer yes/no questions. In general, they were able to generate a reasonable and fluent response but mostly incorrect. Models also tend to overfit to the training examples from the similar domain (e.g. it answers "Montmartre" for the question of "the highest peak of the Beskids range"). Finetuning the model on Wikipedia as the knowledge database improved the results on both tasks by 1.0pp and 0.3pp, respectively. However, directly providing retrieved passages improves the F1 score

by a magnitude of 20pp. Thus, it is evident that there is still plenty of room to improve text-to-text models from a knowledge database perspective.

In the passages subtask we encountered a similar trimming issue as in summarization task 3.2.2. Because of papuGaPT2 input size limitation we couldn't pass the entire passage to the model. That's why we repeated polT-wiki and plBART training using 1024 input size setting. The results for polT and plBART were not significantly different than those with not size-limited input, what confirms their superiority over decoder-only architecture in this task.

Task		open	passages
Base models	polT	20.9	42.8
	polT-wiki	21.9	43.1
	papuGaPT2	18.5	24.0
	mT5	17.2	39.8
	plBART	17.4	37.1
	T5-random	10.9	8.0
Large models	polT	26.5	51.3

Table 4: Accuracy scores for open questions and questions followed by the generated passages evaluated using Poleval *testA* set.

3.4 Machine translation

Machine translation is one of the most popular applications of text-to-text models. Therefore, we finetuned polT on parallel corpora consisting of pairs of English and Polish sentences and evaluate performance in both en→pl and pl→en directions.

3.4.1 Datasets

We used datasets collected for the news translation competition organized as part of the 5th Conference on Machine Translation¹⁹ (WMT20). Specifically, we used 4 out of 5 parallel corpora available to the competition's participants: EuroParl v10, TildeRapid, WikiTitles v2, and ParaCrawl v5.1. The corpora are described by Barrault et al. (2020). We did not use the 5th dataset, WikiMatrix, which is known to introduce more noise than useful knowledge (Caswell et al., 2021). To all 4 corpora, we applied similar filtering as Jónsson et al. (2020).

¹⁷We release the dataset at anonymized

¹⁸polT-base additionally fine-tuned using Wikipedia corpus

¹⁹<https://www.statmt.org/wmt20/translation-task.html>

3.4.2 Experiments

We evaluated the models on the WMT20 development set using the BLEU score (Papineni et al., 2002). For generation, we used beam search with 5 beams and a maximal sequence length of 100. Input sentences were also limited to 100 tokens. We compared polT, mT5, plBART, and papuGaPT2. The models were trained in both directions at the same time, with examples fed to the models alternately. Since polT tokenizer was trained on Polish-only data, we trained a separate tokenizer for WMT20 data. It was a unigram language model (Kudo, 2018) with a vocabulary of 32k tokens based on our training corpora with the same filtering. In order to initialize embeddings of tokens from the WMT20 tokenizer that are not present in the polT tokenizer, we applied a technique based on Arkhipov et al. (2019, section 3). To have a fair comparison with mT5, plBART and papuGaPT2, we applied the same conversion to those models. To distinguish between en→pl and pl→en directions, we prepended source sentences with either <2en> or <2pl> special tokens.

For training we used Adam optimizer (Kingma and Ba, 2014) with learning rate 1e-1 for small models and 1e-2 for base and large ones, gradient accumulation of 8 batches, 10k steps of a linear warm-up, and inverse square root learning rate schedule (Raffel et al., 2020, section 3.1.2). We used different batch sizes for different model sizes to accommodate as many sentences in a batch as possible. Specific batch sizes are listed in Appendix C. We used a single NVIDIA A100 GPU for all the models. Small and base models were trained for 5 epochs and large for 2.

3.4.3 Results

The results are presented in Table 5. polT generates translations superior to mT5 in direction en→pl but not in pl→en. This may be due to the fact that polT was trained on Polish-only data and, therefore, the ability to generate English text deteriorates. polT also outperforms plBART and papuGaPT2. We repeated polT-base and plBART experiments for 3 different seeds. Variance of the polT results was up to 0.1 BLEU, while the variance of the plBART results was slightly higher, up to 0.5 BLEU.

It should be stated that the results reported for polT are not up to the level of the WMT20 winning team’s results (Krubinski et al., 2020). However, polT-large performance in the en→pl direction is slightly better than their bi-directional baseline. In

		Vocab	en→pl	pl→en
Small models	mT5	mT5	20.5	25.0
	polT	polT	20.3	24.7
	mT5	wmt20	<u>20.8</u>	24.8
	polT		<u>20.8</u>	<u>25.4</u>
Base models	mT5	mT5	22.4	27.0
	polT	polT	<u>23.2</u>	26.7
	mT5	wmt20	22.5	26.7
	polT		22.7.	26.8
	plBART	plBART	21.2	25.4
		wmt20	21.6	26.3
Large models	papu-GaPT2	papuGaPT2	21.2	25.5
		wmt20	22.0	26.2
	mT5	wmt20	24.8	29.0
	polT		25.5	28.9

Table 5: BLEU scores for the WMT20 devset. Best results in each category (small, base and large) are underlined. Best overall results are shown in bold.

this work, we did not attempt to achieve state-of-the-art on those datasets, but to compare polT with similar architectures in similar settings.

4 Conclusion

In this work, we introduced a novel text-to-text model for Polish, polT. It outperforms mT5 on the KLEJ benchmark, summarization, en-pl machine translation, and question answering. It is better than plBART and papuGaPT2, except for summarization, where plBART is the best. We should state that the overall performance of plBART is impressive, considering that it has almost twice as few parameters as polT-base. We observed that the larger the model, the better the results (except the summarization), and encoder-decoder architectures are better than decoder only. We efficiently pre-trained polT by initializing the weights from mT5 checkpoints with no exhaustive training. Lastly, since summarization and question answering lack Polish language benchmark datasets, we described their construction and released them publicly.

References

Mikhail Arkhipov, Maria Trofimova, Yuri Kuratov, and Alexey Sorokin. 2019. [Tuning multilingual trans-](#)

692	Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Man-	<i>the 40th Annual Meeting of the Association for Com-</i>	748
693	dar Joshi, Danqi Chen, Omer Levy, Mike Lewis,	<i>putational Linguistics</i> , pages 311–318, Philadelphia,	749
694	Luke Zettlemoyer, and Veselin Stoyanov. 2019.	Pennsylvania, USA. Association for Computational	750
695	Roberta: A robustly optimized bert pretraining ap-	Linguistics.	751
696	proach. <i>arXiv preprint arXiv:1907.11692</i> .		
697	Shayne Longpre, Yi Lu, and Joachim Daiber. 2020.	Jeffrey Pennington, Richard Socher, and Christopher	752
698	Mkqa: A linguistically diverse benchmark for multi-	Manning. 2014. GloVe: Global vectors for word	753
699	lingual open domain question answering.	representation . In <i>Proceedings of the 2014 Confer-</i>	754
700		<i>ence on Empirical Methods in Natural Language</i>	755
701	Michał Marcinczuk, Marcin Ptak, Adam Radziszewski,	<i>Processing (EMNLP)</i> , pages 1532–1543, Doha,	756
702	and Maciej Piasecki. 2013. Open dataset for de-	Qatar. Association for Computational Linguistics.	757
703	velopment of polish question answering systems.		
704	In <i>Proceedings of the 6th Language & Technology</i>	John Prager. 2006. Open-domain question: Answering.	758
705	<i>Conference: Human Language Technologies as a</i>	<i>Found. Trends Inf. Retr.</i> , 1(2):91–231.	759
706	<i>Challenge for Computer Science and Linguistics</i> ,		
707	Wydawnictwo Poznańskie, Fundacja Uniwersytetu	Adam Przepiórkowski, Mirosław Bańko, Rafał L	760
	im. Adama Mickiewicza.	Górski, Barbara Lewandowska-Tomaszczyk, Marek	761
		Łaziński, and Piotr Pęzik. 2011. National corpus of	762
		polish. In <i>Proceedings of the 5th language & tech-</i>	763
708	Louis Martin, Benjamin Muller, Pedro Javier Or-	<i>nology conference: Human language technologies</i>	764
709	tiz Suárez, Yoann Dupont, Laurent Romary, Éric	<i>as a challenge for computer science and linguistics</i> ,	765
710	de la Clergerie, Djamé Seddah, and Benoît Sagot.	pages 259–263.	766
711	2020. CamemBERT: a tasty French language model .		
712	In <i>Proceedings of the 58th Annual Meeting of the</i>	Alec Radford, Jeffrey Wu, Rewon Child, David Luan,	767
713	<i>Association for Computational Linguistics</i> , pages	Dario Amodei, Ilya Sutskever, et al. 2019. Lan-	768
714	7203–7219, Online. Association for Computational	guage models are unsupervised multitask learners.	769
715	Linguistics.	<i>OpenAI blog</i> .	770
716	Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Cor-	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine	771
717	rado, and Jeff Dean. 2013. Distributed representa-	Lee, Sharan Narang, Michael Matena, Yanqi Zhou,	772
718	tions of words and phrases and their compositionality.	Wei Li, and Peter J Liu. 2020. Exploring the lim-	773
719	In <i>Advances in neural information processing</i>	its of transfer learning with a unified text-to-text	774
720	<i>systems</i> , pages 3111–3119.	transformer. <i>Journal of Machine Learning Research</i> ,	775
		21:1–67.	776
721	Robert Mroczkowski, Piotr Rybak, Alina Wróblewska,	Adam Roberts, Colin Raffel, and Noam Shazeer. 2020.	777
722	and Ireneusz Gawlik. 2021. HerBERT: Efficiently	How much knowledge can you pack into the param-	778
723	pretrained transformer-based language model for	eters of a language model? In <i>Proceedings of the</i>	779
724	Polish . In <i>Proceedings of the 8th Workshop on Balto-</i>	<i>2020 Conference on Empirical Methods in Natural</i>	780
725	<i>Slavic Natural Language Processing</i> , pages 1–10,	<i>Language Processing (EMNLP)</i> , pages 5418–5426,	781
726	Kiyv, Ukraine. Association for Computational Lin-	Online. Association for Computational Linguistics.	782
727	guistics.		
728	El Moatez Billah Nagoudi, Wei-Rui Chen, Muham-	Piotr Rybak, Robert Mroczkowski, Janusz Tracz, and	783
729	mad Abdul-Mageed, and Hasan Cavusoglu. 2021a.	Ireneusz Gawlik. 2020. KLEJ: Comprehensive	784
730	IndT5: A text-to-text transformer for 10 indigenous	benchmark for Polish language understanding . In	785
731	languages . In <i>Proceedings of the First Workshop on</i>	<i>Proceedings of the 58th Annual Meeting of the Asso-</i>	786
732	<i>Natural Language Processing for Indigenous Lan-</i>	<i>ciation for Computational Linguistics</i> , pages 1191–	787
733	<i>guages of the Americas</i> , pages 265–271, Online. As-	1201, Online. Association for Computational Lin-	788
734	sociation for Computational Linguistics.	guistics.	789
735	El Moatez Billah Nagoudi, AbdelRahim Elmadany,	Abigail See, Peter J. Liu, and Christopher D. Manning.	790
736	and Muhammad Abdul-Mageed. 2021b. Arat5:	2017. Get to the point: Summarization with pointer-	791
737	Text-to-text transformers for arabic language under-	generator networks . In <i>Proceedings of the 55th An-</i>	792
738	standing and generation .	<i>nuual Meeting of the Association for Computational</i>	793
		<i>Linguistics (Volume 1: Long Papers)</i> , pages 1073–	794
		1083, Vancouver, Canada. Association for Computa-	795
739	Maciej Ogrodniczuk and Mateusz Kopeć. 2014. The	tional Linguistics.	796
740	Polish summaries corpus . In <i>Proceedings of the</i>		
741	<i>Ninth International Conference on Language Re-</i>	Noam Shazeer, Youlong Cheng, Niki Parmar, Dustin	797
742	<i>sources and Evaluation (LREC’14)</i> , pages 3712–	Tran, Ashish Vaswani, Penporn Koanantakool, Peter	798
743	3715, Reykjavik, Iceland. European Language Re-	Hawkins, HyoukJoong Lee, Mingsheng Hong, Cliff	799
744	sources Association (ELRA).	Young, et al. 2018. Mesh-tensorflow: deep learning	800
		for supercomputers. In <i>Proceedings of the 32nd In-</i>	801
745	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-	<i>ternational Conference on Neural Information Pro-</i>	802
746	Jing Zhu. 2002. Bleu: a method for automatic eval-	<i>cessing Systems</i> , pages 10435–10444.	803
747	uation of machine translation . In <i>Proceedings of</i>		

804	Noam Shazeer and Mitchell Stern. 2018. Adafactor: Adaptive learning rates with sublinear memory cost.	A.3 Learning rate tuning	852
805		We performed the best learning rate search for each model in the set of {2e-5, 4e-5, 6e-5, 8e-5, 1e-4, 2e-4, 4e-4, 6e-4, 8e-4, 1e-3}. For learning rate tuning, we used tasks with the validation set provided. The designated learning rate value was used for all test runs on all KLEJ tasks.	853
806	Antti Virtanen, Jenna Kanerva, Rami Ilo, Jouni Luoma, Juhani Luotolahti, Tapio Salakoski, Filip Ginter, and Sampo Pyysalo. 2019. Multilingual is not enough: Bert for finnish.	B Summarization task	854
807		B.1 Allegro Reviews (AA)	855
808		The raw dataset is a collection of articles scrapped from https://allegro.pl/artykuly website. Each example contains title, lead, body, information about category tree, and other metadata. In a pre-processing stage, we removed markdown formatting and normalized white spaces. This dataset serves as a good benchmark for highly abstractive summarization/generation tasks since every article contains a title or lead that can be viewed as a loose abstract of the article. Moreover, every article is written by a professional editor, and they are relatively short hence fit almost entirely into 512 token context.	856
809		We constructed 2 versions of the abstractive summarization task, which affect the summary ratio. In <i>body2lead</i> task, we use a body of the article as a source and lead as a generation target. In <i>body+lead2title</i> we use concatenated lead and body as source (full article) and generate title as a target.	857
810	Linting Xue, Noah Constant, Adam Roberts, Mihir Kale, Rami Al-Rfou, Aditya Siddhant, Aditya Barua, and Colin Raffel. 2021. mT5: A massively multilingual pre-trained text-to-text transformer. In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 483–498, Online. Association for Computational Linguistics.	For evaluation purposes, we divided the whole data source into 80% train set and 20% test set in a stratified way using top-level category in Allegro category tree as a label. Then, we saved 10% of the train set for validation (also stratified split) and hyperparameter tuning. Table 7 contains summary of average source/target length and size of each split.	858
811		B.2 Polish Summaries Corpus (PSC)	859
812		Original dataset (Ogrodniczuk and Kopeć, 2014) contains extractive or abstractive summaries and other metadata for news articles prepared by annotators ²⁰ . Each article has a few corresponding summaries that depend on the annotator, summary ratio (5, 10, or 20% of the original text), and type (extractive or abstractive). Due to the small size of the entire dataset, we included all summaries in the final dataset. In this way, the source contains duplicates and targets near-duplicates. Moreover,	860
813		²⁰ We used publicly available version https://huggingface.co/datasets/polsum	861
814			862
815			863
816			864
817			865
818			866
819	Yinfei Yang, Daniel Cer, Amin Ahmad, Mandy Guo, Jax Law, Noah Constant, Gustavo Hernandez Abrego, Steve Yuan, Chris Tar, Yun-Hsuan Sung, et al. 2019. Multilingual universal sentence encoder for semantic retrieval. <i>arXiv preprint arXiv:1907.04307</i> .		867
820			868
821			869
822			870
823			871
824			872
825	Changchang Zeng, Shaobo Li, Qin Li, Jie Hu, and Jianjun Hu. 2020. A survey on machine reading comprehension: Tasks, evaluation metrics and benchmark datasets.		873
826			874
827			875
828			876
829	A Converting KLEJ benchmark to text-to-text format		877
830			878
831	A.1 Task formulation		879
832	We transformed KLEJ tasks in a text-to-text format where each input sentence (one or two in KLEJ tasks) is given a prompt describing its role in the task semantics. Additionally, a task identification token could prefix whole input, but in all experiments in the paper, we did not use such prompt because we were training models for one task at a time. Finally, labels were converted into semantically significant text tokens. Input features and task labels are given in the Table 6.		880
833			881
834			882
835			883
836			884
837			885
838			886
839			887
840			888
841			889
842	A.2 Evaluation		890
843	At test time model is given input formatted as described above, and the goal is to generate a label. A sequence of target tokens is generated until EOS token generation or when it reaches max target length. We calculate the max target length per task, and it corresponds to the longest label after the tokenization process. The generated sequence of tokens is converted into text. We count for a positive when it matches precisely the corresponding label.		891
844			892
845			893
846			894
847			895
848			896
849			897
850			898
851			899

Task	Prefix 1	Prefix 2	Labels
NKJP-NER	zdanie	-	geograficzna, brak, organizacja, osoba, miejsce, czas
CBD	zdanie	-	neutralna, przemoc
Czy wiesz?	pytanie	odpowiedź	fałsz, prawda
PolEmo2.0	zdanie	-	niejednoznaczny, negatywny, pozytywny, neutralny
AR	zdanie	-	1.0, 2.0, 3.0, 4.0, 5.0
PSC	streszczenie 1	streszczenie 2	nie_parafraza, parafraza
CDSC-E	zdanie 1	zdanie 2	neutralny, wynikanie, sprzeczność

Table 6: Prompts and labels used in text-to-text formulation of KLEJ benchmark tasks. In case of encoder-decoder architectures, source was prepared as <Prefix 1>: text 1 <Prefix 2>: text 2. For decoder architecture we additionally append [SEP] token at the end.

our train set contains target summaries of different ratios that may affect learning and prediction.

We constructed 3 versions of the summarization task. The *extract* subset contains only extractive summaries, the *abstract* subset contains only abstractive summaries and *whole* contains both of them.

For evaluation purposes, we divided the whole data source into 80% train set and 20% test set in a stratified way using summary ratio, summary type, and article category as a label. Then, we saved 10% of the train set for validation (also stratified split) and hyperparameter tuning. Table 7 contains summary of average source/target length and size of each split.

Task	AA body2lead	AA body+lead2title	PSC whole	PSC extract	PSC abstract
Domain	E-commerce		News articles		
Average length (characters)					
Source	3.9k	4.1k	10.6k		
Target	0.3k	0.05k	1.3k		
Number of examples					
Train	24.4k		0.8k	6.1k	1.7k
Dev	2.7k		0.9k	0.7k	0.2k
Test	6.8k		2.2k	1.7k	0.5k

Table 7: Overview of summarization tasks with total number of examples for each split and average length (in characters) of source and target

B.3 Trimming inputs

The BART and GPT-2 architectures use fixed positional encoding, which is limited to 1024 tokens by default. On the other hand, the T5 model uses relative positional encoding and can process sequences of any length in principle. In our case, we are bounded by memory requirements. Thus, we fixed 1024 tokens as the maximal sequence length for all models. Even then, large architectures required further trimming of input tensors due to memory errors. Consequently, we were not able to fit the whole target and source texts into the model inputs.

Since each model uses a different tokenizer, the context used during training and prediction was different for individual models. The effect is more pronounced for source texts because they are much longer than targets. We were always able to input 98-100% of target tokens into the model and assumed that target trimming did not substantially affect the results. Usually, only a small number of outliers (extremely long targets) were trimmed.

On the other hand, source trimming could have affected the models. The final number of tokens to which the source was trimmed and the percentage of the source tokens seen by the model are shown in Table 8. The PSC task was affected the most. For example, mT5-large saw on average only 25% of the source. Similarly, the papuGaPT2 saw on average only 26% of the source. Less affected were plBART and polT models, which saw above 50% of the source. Most importantly, none of the models ever saw 100% of source texts during training for PSC tasks.

B.4 Learning rate tuning

On AA task, we performed best learning rate search in the range {1e-5, 2e-5, 4e-5, 6e-5, 8e-5, 1e-4, 2e-4, 4e-4, 6e-4, 8e-4, 1e-3}. On PSC task, we per-

Task		PSC	AA body2lead	AA body+lead2title
Small models	mT5	33%	85%	80%
	polT	51%	96%	95%
Base models	papuGaPT2	26% (512)	92% (767)	96% (959)
	mT5	33%	85%	80%
	plBART	52%	96%	95%
	polT	51%	96%	95%
Large models	mT5	25% (768)	85%	80%
	polT	51%	96%	95%

Table 8: Average percentage of tokens from source text that fit into model. Differences are due to a tokenizer, model capacity, and memory limitations. By default, source text was trimmed to 1024 tokens, and the exceptions are indicated (in parentheses) in the table.

formed learning rate search for each model on PSC whole subset and used the same learning rates for PSC extract and PSC abstract subsets as well. The search was geometrical, with the start point at $4e-4$, i.e. at first subset $\{2e-4, 4e-4, 8e-4\}$ was examined, each learning rate with 3 seed runs. If the middle point was better than the two other, the search completed, otherwise either $\{1e-4, 2e-4, 4e-4\}$ or $\{4e-4, 8e-4, 1.6e-3\}$ learning rates were examined during the next step depending on the trend. As the result, learning rates from $5e-5$ to $2.56e-2$ were selected for different models. The only exception was mT5-large model which was unstable and required finer granularity $\{2.5e-5, 5e-5, 7.5e-5\}$ was used at the end of the search. Our procedure of learning rate tuning was motivated by budget constraints.

C Machine translation batch sizes

All translation experiments were run on a single NVIDIA A100 GPU with 40 GB of memory. To utilize all the available memory we used different batch sizes for different model architectures. Specific batch sizes are shown in Table 9.

	Model	Vocab	Batch size
Small models	mT5	mT5	25
	polT	polT	50
Base models	mT5	wmt20	50
	polT	mT5	10
		polT	25
	mT5	wmt20	25
Large models	plBART	plBART wmt20	40
	papu-GaPT2	papuGaPT2 wmt20	25
Large models	mT5	wmt20	10
	polT		

Table 9: Machine translation batch sizes

D Summarization ROUGE scores

	Model	ROUGE-1	ROUGE-2	ROUGE-L	ROUGE
Baselines	lead n=3 source sentences	20.3	2.8	14.1	12.4
	lead n (adaptive, avg n=3) source sentences	20.2	2.8	14.1	12.4
Small models	mT5	19.7	<u>4.1</u>	<u>15.2</u>	<u>13.0 ± 0.1</u>
	polT	<u>21.5</u>	<u>4.9</u>	<u>16.5</u>	<u>14.3 ± 0.1</u>
Base models	papuGaPT2	20.2	1.6	12.6	12.1 ± 0.1
	mT5	<u>21.4</u>	4.6	<u>16.3</u>	<u>14.1 ± 0.2</u>
	polT	<u>22.3</u>	<u>5.2</u>	<u>17.1</u>	<u>14.9 ± 0.1</u>
	plBART	24.0	5.1	17.7	15.6 ± 0.1
Large models	mT5	16.7	2.7	13.4	10.9 ± 2.4
	polT	19.2	2.7	<u>14.4</u>	12.1 ± 2.6

Table 10: Evaluation results on Allegro Articles body2lead task (test split, mean values across 3 runs). The best scores are in bold, results above the baseline are underlined. The ROUGE score reported in the last column is arithmetic mean of ROUGE-1, ROUGE-2 and ROUGE-L scores. All reported ROUGE scores are f-measures.

	Model	ROUGE-1	ROUGE-2	ROUGE-L	ROUGE
Baselines	lead n=3 source sentences	9.3	2.9	8.1	6.8
	lead n (adaptive, avg n=1) source sentences	10.4	3.6	9.7	7.9
Small models	mT5	<u>39.3</u>	<u>25.1</u>	<u>38.2</u>	<u>34.2 ± 0.3</u>
	polT	<u>41.4</u>	<u>27.2</u>	<u>40.1</u>	<u>36.3 ± 0.1</u>
Base models	papuGaPT2	<u>18.8</u>	<u>5.2</u>	<u>18.0</u>	<u>14.0 ± 0.1</u>
	mT5	<u>41.3</u>	<u>27.0</u>	<u>40.1</u>	<u>36.1 ± 1.2</u>
	plBART	<u>44.0</u>	<u>29.6</u>	<u>41.3</u>	<u>38.3 ± 0.1</u>
	plBART [†]	<u>44.0</u>	<u>29.7</u>	<u>42.4</u>	<u>38.7 ± 0.5</u>
	polT	<u>44.1</u>	<u>29.7</u>	<u>42.8</u>	<u>38.9 ± 0.3</u>
	polT [†]	<u>43.9</u>	<u>29.4</u>	<u>42.6</u>	<u>38.7 ± 0.2</u>
Large models	mT5	<u>38.4</u>	<u>24.7</u>	<u>37.3</u>	<u>33.5 ± 1.8</u>
	polT	44.6	30.2	43.3	39.4 ± 0.4

Table 11: Evaluation results on Allegro Articles body+lead2title task (test split, mean values across 3 runs). The best scores are in bold, results above the baseline are underlined. The ROUGE score reported in the last column is arithmetic mean of ROUGE-1, ROUGE-2 and ROUGE-L scores. All reported ROUGE scores are f-measures. (†) models were trained with the same input and target truncation as papuGaPT2.

	Model	ROUGE-1	ROUGE-2	ROUGE-L	ROUGE
Baselines	lead n=3 source sentences	28.2	16.5	21.3	22.0
	lead n (adaptive, avg n=11) source sentences	39.0	23.8	28.3	30.3
Upper bounds	human performance	<u>41.8</u>	<u>27.4</u>	<u>33.8</u>	<u>34.3</u>
Small models	mT5	29.4	17.0	23.5	23.3 ± 0.5
	polT	36.7	<u>24.3</u>	<u>30.4</u>	<u>30.5 ± 0.1</u>
Base models	papuGaPT2	26.9	7.6	15.5	16.7 ± 0.1
	mT5	27.4	14.1	20.3	20.6 ± 0.2
	plBART	39.0	26.3	32.4	32.6 ± 0.1
	polT	32.9	18.1	24.5	25.2 ± 0.5
Large models	mT5	20.6	3.7	11.5	12.0 ± 0.3
	polT	23.1	12.0	17.4	17.5 ± 0.6

Table 12: Evaluation results on Polish Summary Corpus (whole) task (test split, mean values across 3 runs). The best scores are in bold, results above the baseline are underlined. The ROUGE score reported in the last column is arithmetic mean of ROUGE-1, ROUGE-2 and ROUGE-L scores. All reported ROUGE scores are f-measures.

	Model	ROUGE-1	ROUGE-2	ROUGE-L	ROUGE
Baselines	lead n=3 source sentences	29.5	18.3	22.6	23.4
	lead n (adaptive, avg n=11) source sentences	40.0	25.5	29.7	31.7
Upper bounds	human performance	<u>45.1</u>	<u>33.3</u>	<u>38.8</u>	<u>39.0</u>
Small models	mT5	30.4	19.3	25.3	25.0 ± 0.2
	polT	27.7	18.3	23.6	23.2 ± 13.8
Base models	papuGaPT2	26.9	8.5	15.9	17.1 ± 0.2
	mT5	27.6	15.0	20.9	21.2 ± 0.5
	plBART	39.9	28.6	34.4	34.3 ± 0.1
	polT	31.9	18.0	24.1	24.7 ± 0.5
Large models	mT5	19.3	2.6	10.9	10.9 ± 1.4
	polT	22.4	7.3	15.6	15.1 ± 2.6

Table 13: Evaluation results on Polish Summary Corpus (extract) task (test split, mean values across 3 runs). The best scores are in bold, results above the baseline are underlined. The ROUGE score reported in the last column is arithmetic mean of ROUGE-1, ROUGE-2 and ROUGE-L scores. All reported ROUGE scores are f-measures.

	Model	ROUGE-1	ROUGE-2	ROUGE-L	ROUGE
Baselines	lead n=3 source sentences	0.271	0.142	0.197	0.203
	lead n (adaptive, avg n=11) source sentences	0.354	0.178	0.236	0.256
Upper bounds	human performance	0.346	0.167	<u>0.250</u>	0.254
Small models	mT5	0.280	0.149	0.207	0.212 ± 0.005
	polT	0.304	0.165	0.227	0.232 ± 0.019
Base models	papuGaPT2	0.260	0.057	0.139	0.152 ± 0.002
	mT5	0.120	0.050	0.095	0.088 ± 0.025
	plBART	0.335	0.177	0.254	0.255 ± 0.007
	polT	0.156	0.077	0.120	0.117 ± 0.073
Large models	mT5	0.161	0.019	0.094	0.091 ± 0.008
	polT	0.156	0.051	0.112	0.106 ± 0.011

Table 14: Evaluation results on Polish Summary Corpus (abstract) task (test split, mean values across 3 runs). The best scores are in bold, results above the baseline are underlined. The ROUGE score reported in the last column is arithmetic mean of ROUGE-1, ROUGE-2 and ROUGE-L scores. All reported ROUGE scores are f-measures.

E	Example predictions on Allegro Reviews <i>body+lead2title</i>	974
1.	Gold Summary: Speedminton – badminton na wietrzne dni	975
	poIT-base: Jak zacząć grać w speedminton?	976
	plBART: Speedminton – sport dla każdego	977
	papuGaPT2: Squa – jak zacząć?	978
	mT5-base: Jak zacząć grać w speedmintona?	979
	poIT-large: Speedminton – sport dla aktywnych	980
2.	Gold Summary: Urządzamy wnętrze po skandynawsku	981
	poIT-base: Jak urządzić mieszkanie w stylu skandynawskim?	982
	plBART: Jak urządzić mieszkanie w stylu skandynawskim?	983
	papuGaPT2: Skandynawski styl w wersji budżetowej – jak go urządzić?	984
	mT5-base: Jak urządzić wnętrze w stylu skandynawskim?	985
	poIT-large: Jak urządzić mieszkanie w stylu skandynawskim?	986
3.	Gold Summary: Czy warto robić domowe peelingi – poradnik	987
	poIT-base: Najlepsze peelingi do domu	988
	plBART: Najlepsze peelingi do twarzy	989
	papuGaPT2: Peeslingi – najlepszy sposób na przesuszoną skórę	990
	mT5-base: Jak zrobić peelingi w domu?	991
	poIT-large: Jak zrobić peeling w domu?	992
4.	Gold Summary: Szerokie spodnie – must have sezonu wiosna/lato 2015	993
	poIT-base: Wiosenne trendy – szerokie spodnie	994
	plBART: Modne szerokie spodnie – jak je nosić?	995
	papuGaPT2: Wiosenne stylizacje ze schullami w roli głównej	996
	mT5-base: Wiosenny trend – szerokie spodnie	997
	poIT-large: Wiosenne spodnie w stylu lat 70.	998
5.	Gold Summary: Przygotuj się na lodowisko	999
	poIT-base: Jak rozpocząć sezon na łyżwy?	1000
	plBART: Jak zacząć przygodę z łyżwami?	1001
	papuGaPT2: Jakie ubrania na sezon kąpielowy?	1002
	mT5-base: Jak zacząć sezon na łyżwy?	1003
	poIT-large: Jak przygotować się do sezonu na łyżwy?	1004
6.	Gold Summary: Wózek widłowy niczym limuzyna sterowana iPadem	1005
	poIT-base: Wózki widłowe zwiększające możliwości swojego działania	1006
	plBART: Wózki widłowe – zwiększamy wydajność i oszczędzamy paliwo	1007
	papuGaPT2: Nowoczesny wózek widłowy – dlaczego jest tak ważny?	1008
	mT5-base: Czym są nowoczesne technologie w wózkach widłowych?	1009
	poIT-large: Automatyzacja wózków widłowych	1010

F Example predictions on Allegro Reviews *body2lead*

- 1011
1012
1013
1014
1. **Gold Summary:** Każdy czuje respekt przed królową nauk. Matematyka spędza sen z powiek niejednemu uczniowi. Są jednak sposoby zachęcenia dziecka do nauki matematyki i sprawienia, że lepiej ją zrozumie i polubi.

1015
1016
1017

poIT-base: Matematyka to jedna z najprzyjemniejszych dziedzin nauki. Uczymy się ją w szkole i w domu. Nauka przez zabawę może być świetną rozrywką dla całej rodziny. Jakie gadzety mogą pomóc dziecku w nauce liczenia?

1018
1019

plBART: Nauka matematyki nie musi być trudna. Wystarczy odrobina chęci i kilka akcesoriów, które pomogą dziecku w odkrywaniu tajemnic matematyki.

1020
1021
1022

papuGaPT2: Nauka poprzez zabawę to jeden z najlepszych sposobów na rozwijanie zdolności poznawczych i umiejętności logicznego myślenia. Podpowiadamy, co warto zabrać ze sobą do szkoły i na co zwrócić uwagę przy wyborze odpowiednich akcesoriów.

1023
1024
1025

mT5-base: Nauka matematyki to doskonały sposób na spędzenie wolnego czasu w gronie najbliższych. Dzieci chętnie sięgają po pomoce naukowe, które zachęcą do nauki przez zabawę i doświadczenie.

1026
1027

poIT-large: Matematyka nie musi być nudna. Nauka liczenia może być wspaniałą zabawą, która wspomogę rozwój dziecka. Podpowiadamy, jak zachęcić dziecko do nauki matematyki.

- 1028
1029
1030
1031
2. **Gold Summary:** Ciepłe, letnie dni aż się proszą o aktywną zabawę na dworze. Dlatego jeśli planujesz przyjęcie dla dziecka i jego przyjaciół, koniecznie postaw na garden party. Tego typu imprezy są nie tylko modne, ale także łatwe w przygotowaniu i posprzątaniu po nich. Jeśli nie masz pomysłu, jak je przygotować, podpowiadamy.

1032
1033

poIT-base: Ozdabianie ogrodu to nie tylko świetna zabawa, ale także świetna rozrywka. Podpowiadamy, jak przygotować ogrodowe przyjęcie dla dziecka.

1034
1035
1036
1037

plBART: Dzieci uwielbiają spędzać czas na świeżym powietrzu. To doskonała okazja do zorganizowania im urodzinowego przyjęcia, które zachwyci ich nie tylko wyglądem, ale także kreatywnością i zapewni wiele pozytywnych emocji. Podpowiadamy, jak przygotować party dla maluchów w cenie do 50 zł.

1038
1039
1040

papuGaPT2: Organizujesz przyjęcie urodzinowe dla dziecka? Zastanawiasz się, co mu kupić, aby było mu jak najlepiej. Podpowiadamy, na co zwrócić uwagę, wybierając dekoracje do ogrodu. Sprawdź nasze propozycje.

1041
1042
1043

mT5-base: Oгородowe przyjęcie to wyjątkowe wydarzenie w życiu każdego dziecka. To świetna okazja do spędzenia czasu na świeżym powietrzu. Warto więc zadbać o to, aby było przyjemne i kolorowe. Podpowiadamy, jakie dekoracje będą najlepsze dla twojego malucha.

1044
1045

poIT-large Ozdabianie ogrodu to nie tylko świetna zabawa, ale także świetna rozrywka. Podpowiadamy, jak przygotować ogrodowe przyjęcie dla dziecka.

G Example predictions on PSC (whole)

1046

1. **Gold Summary:** Ostatnia znaczna podwyżka stóp procentowych NBP oraz gwałtowne przyspieszenie inflacji wywołały krytykę poczyniła Rady Polityki Pieniężnej. Rada chciałaby podjąć dyskusję z głosami krytycznymi. nikt nie przewidywał, że w 1999 r. inflacja przewyższy cel inflacyjny NBP. Świadczy to dobitnie o tym, że głównymi czynnikami przyspieszającymi inflację były szoki podażowe na rynku żywności i paliw. Inny komentator stawia z kolei zarzut, że Rada niepotrzebnie zwlekała z decyzją o podwyżce stóp do listopada. Rada listopadową decyzję nie dlatego podjęła w listopadzie, i w takiej skali, że "zaspała" w październiku, ale dlatego, by podjąć ją w takiej właśnie skali w listopadzie. Dla każdego, kto wie, co to jest efektywność mechanizmu transmisji impulsów polityki pieniężnej do gospodarki, decyzja taka jest zrozumiała. Każdy ma prawo twierdzić, że dokonane podniesienie stóp NBP jest za duże. Pamiętać jednak musi, że jego twierdzenie sprowadza się do tezy: skala dokonanej podwyżki stóp procentowych doprowadzi do przestrzelenia celu inflacyjnego w dół. Czy rację mają ci, którzy krytykują teraz Radę, dowiemy się za rok. Przejdźmy do zarzutu zbyt dużej redukcji stóp procentowych 20 stycznia 1999 r. Otóż, redukcja rzeczywiście była zbyt duża bądź niepotrzebna. Problem polega jednak na tym, że o tym wiemy dopiero teraz. skala obniżki stóp procentowych dostosowana była do przewidywanego w 1999 r. przebiegu zjawisk makroekonomicznych. Było to działanie zgodne z zasadą forward looking. Czy Rada miała przesłanki do wcześniejszej zmiany swoich przewidywań? Bardzo mało i późno. o problemach ZUS dowiedzieliśmy się dopiero w lipcu. Gwałtowny wzrost cen żywności rozpoczął się od sierpnia. wzrosły ceny paliw. Z żadnych analiz nie wynikała taka skala wzrostu cen na rynku żywności i paliw. W tym kontekście stawianie Radzie zarzutu, że będąc odpowiedzialną za walkę z inflacją, stara się zrzucić winę na rząd i wydarzenia na rynku międzynarodowym, jest nierzetelne. Nie można wina strażnicy pożarnej za lekkomyślność dziecka z zapalkami, jeśli okaże się sprawna w gaszeniu pożaru. nie można ignorować krótkookresowego wpływu na inflację czynników pozapieniężnych. nie da się na trwałe zmniejszyć inflacji bez budowy sprawnych i konkurencyjnych rynków paliw czy energii.

1047

1048

1049

1050

1051

1052

1053

1054

1055

1056

1057

1058

1059

1060

1061

1062

1063

1064

1065

1066

1067

1068

1069

1070

mT5-small: Ostatnia znaczna podwyżka stóp procentowych NBP oraz gwałtowne przyspieszenie inflacji wywołały ze strony niektórych ekonomistów, analityków i komentatorów krytykę poczyniła Rady Polityki Pieniężnej. niestety, trudno podjąć merytoryczną dyskusję z głosami, szczególnie w sytuacji niepowodzeń. Rada listopadową decyzję podejmowała w całkowitym spokoju, bez żadnej nerwowości. Piędą zwracać uwagę na to, że polityka pieniężna była zbyt pożądanym, zważywszy rolę, jaką odgrywają. W tym kontekście stawianie Radzie zarzutów można było doprowadzić do końca września.

1071

1072

1073

1074

1075

1076

1077

polT-small: Ostatnia podwyżka stóp procentowych NBP oraz gwałtowne przyspieszenie inflacji wywołały krytykę poczyniła Rady Polityki Pieniężnej. Rada listopadową decyzję podejmowała w całkowitym spokoju. Dla każdego, kto wie, co to jest efektywność mechanizmu transmisji impulsów polityki pieniężnej do gospodarki i jakie jest znaczenie oczekiwań w tym procesie, decyzja taka jest zrozumiała. redukcja stóp w styczniu rzeczywiście była zbyt duża bądź w ogóle niepotrzebna.

1078

1079

1080

1081

1082

papuGaPT2: redukcja stóp w styczniu była zbyt duża bądź w ogóle niepotrzebna. skala krytycznych wniosków była dość trudna, stawiano wiele pytań, do których nie podejmowano ostatecznych decyzji. Rada miała przesłanki do wcześniejszej zmiany swoich przewidywań bardzo późno. informacje o wszystkich wydarzeniach dostępne były kilka tygodni po ich wystąpieniu. niską przewidywalność polityki gospodarczej rządu widać było gołym okiem. decyzja taka jest zrozumiała. redukcja stóp była koniecznym posunięciem, aby zbliżyć się do rynku. jednak Rada przyjmowała wówczas za wiarygodną zapowiedź zaostrzenia polityki fiskalnej oraz przewidywała powolny wzrost eksportu. starano się reanimować kryzys. pojawiła się propozycja, żeby Rada zaczęła mówić jednym głosem. Dziś grono to znacznie się powiększyło. z kolejnych wypowiedzi zaczęli wycofać wnioski o rewizję nadzwyczajną. Rada listopadową decyzję podejmowała w całkowitym spokoju, bez żadnej nerwowości. Dla każdego, kto wie, co to jest efektywność mechanizmu transmisji impulsów polityki do gospodarki i jakie jest znaczenie oczekiwań w tym procesie, decyzja o podwyżce stóp nie jest żadnym donosem. To oznacza duże prawdopodobieństwo żywiołowego narastania deficytu sektora pub-

1083

1084

1085

1086

1087

1088

1089

1090

1091

1092

1093

1094

1095

1096 liczne. Ryzyko destabilizacji rynku jest wysokie. nie wolno zrezygnować z powrotu na przyzwoitą
1097 ścieżkę wzrostu ani z poprawienia równowagi zewnętrznej. Trzeba bardziej powściągliwie określić
1098 możliwości wycofania się inflacji z polityki monetarnej. jeśli spojrzymy na deprecjację parytetu,
1099 okazuje się, że zestrzelenie nastąpiło nie z inicjatywy samego ROR-u, lecz z chęci ponownego oży-
1100 wienia poprzez odpartyjnienie poprzez podjęcie działań, których w tej sytuacji jeszcze nie jesteśmy
1101 w stanie podjąć.

1102 **mT5-base:** Nie da się na trwałe zmniejszyć inflację bez budowy sprawnych i konkuren-
1103 cyjnych rynków paliw czy energii Szukanie winnego RYS. MARCIN CHUDZIK BoGUSŁAW
1104 GRABOWSKI Ostatnia znaczna podwyżka stóp procentowych NBP oraz gwałtowne przyspieszenie
1105 inflacji wywołały ze strony niektórych ekonomistów krytykę poczynąń Rady Polityki Pieniężnej.
1106 Niestety, trudno podjąć merytoryczną dyskusję z głosami, szczególnie w sytuacji niepowodzeń w
1107 ograniczaniu inflacją.

1108 **pIBART:** Ostatnia znaczna podwyżka stóp procentowych NBP oraz gwałtowne przyspieszenie
1109 inflacji wywołały ze strony niektórych komentatorów krytykę poczynąń Rady Polityki Pieniężnej.
1110 Rada listopadową decyzję podejmowała w całkowitym spokoju. Dla każdego, kto wie, co to
1111 jest efektywność mechanizmu transmisji impulsów polityki pieniężnej do gospodarki i jakie jest
1112 znaczenie oczekiwań w tym procesie, decyzja taka jest w pełni zrozumiała. redukcja stóp w styczniu
1113 rzeczywiście była zbyt duża bądź niepotrzebna. o tym wiemy dopiero teraz. stawianie Radzie
1114 zarzutu, że będąc odpowiedzialną za walkę z inflacją, stara się zrzucić winę na rząd, jest nierzetelne.

1115 **polT-base:** Rząd Polityki Pieniężnej w poczuciu odpowiedzialności przywiązuje wielkie znaczenie
1116 do przejrzystości swoich działań od początku istnienia. Dlatego chciałaby podjąć dyskusję z głosami
1117 krytycznymi, szczególnie w sytuacji niepowodzeń w ograniczaniu inflacji. Znaczna część ostatnich
1118 krytycznych wypowiedzi wobec Rady sprowadzała się właściwie do trzech zarzutów: o to, że ostatnia
1119 podwyżka stóp procentowych była "spóźnioną, zbyt nerwową i przesadną" reakcją Rady, której
1120 " optymizm co do inflacji trwał zbyt długo", i która "jakby próbowała nadrobić zaległości w walce z
1121 inflacją".

- 1122 2. **Gold Summary:** Dwaj słynni bokserzy zawodowi Amerykanin Evander Holyfield i Brytyjczyk
1123 Lennox Lewis szykują się do rewanżowego pojedynku. Ich pierwszy pojedynek zakończył się,
1124 haniebnym według fachowców, remisem. Wygrany będzie królem boksu zawodowego. Historia
1125 boksu jest pełna ekscytujących rewanży, mają je również na swoim koncie Holyfield i Lewis. Obecny
1126 faworyt Lewis, jest aktualnym mistrzem świata WBC. Początkowy faworyt Holyfield, ma dwa pasy
1127 mistrzowskie WBA i IBF. Jeśli przegra Holyfield przegra tym samym Don King, najslawniejszego
1128 promotora boksu zawodowego. Walkę transmituje Canal Plus.

1129 **mT5-small:** W najbliższą sobotę w Las Vegas dojdzie do kolejnego "rewanżu stulecia". Zmierzą
1130 się w nim dwaj najwybitniejsi obecnie bokserzy wagi ciężkiej, Amerykanin Evander Holyfield i
1131 Brytyjczyk Lennox Lewis. Osiem miesięcy temu, 13 marca, w nowojorskiej Madison Square Garden
1132 ich pierwsza walka zakończyła się kontrowersyjnym remisem. Stawką rewanżowego pojedynku
1133 będą pasy mistrzowskie trzech najważniejszych organizacji boksu zawodowego: WBA, WBC, IBF.

1134 **polT-small:** W najbliższą sobotę w Las Vegas dojdzie do kolejnego "rewanżu stulecia". Zmierzą
1135 się w nim dwaj najwybitniejsi obecnie bokserzy wagi ciężkiej, Amerykanin Evander Holyfield i
1136 Brytyjczyk Lennox Lewis. Osiem miesięcy temu, 13 marca, w nowojorskiej Madison Square Garden
1137 ich pierwsza walka zakończyła się kontrowersyjnym remisem. Stawką rewanżowego pojedynku
1138 będą pasy mistrzowskie trzech najważniejszych organizacji boksu zawodowego: WBA, WBC, IBF.
1139 Ostatnim posiadaczem trzech pasów w wadze ciężkiej był Riddick Bowe. Od siedmiu lat żaden z
1140 pięściarzy nie może się uważać za króla tej najbardziej cenionej kategorii.

1141 **papuGaPT2:** W nocy z 13 na 14 listopada w Thomas and Mack Center w "Ready World Press"
1142 zmierzą się dwaj znakomici zawodnicy. Ich walkę będzie można oglądać na żywo w Canal Plus od
1143 godz. 3:25. Retransmisja w niedzielę będzie miała miejsce już raz trzeci w historii walki, która miała
1144 być największym wydarzeniem ostatnich lat. Transmitowany przez Jamesa "Bustera" pojedynek

Sonny'ego Listona z Bobem Fitzsimmonsem w 1899 r. zakończył się dyskwalifikacją jednego z zawodników. Walka okazała się skandalem. Trzy niezależne postępowania wyjaśniające, które	1145
prowadzono po tej walce, niewiele dały. Nikomu nie udowodniono winy. \nBoks jest tak stary jak	1146
świat. Już Homer opiewał walki na gołe pięści. Wielkim zwolennikiem walk bokserów był rzymski	1147
cesarz Kaligula. W 1350 roku, gdy naprzeciw siebie stanęli Jeffries i O'Donnel, dwóch zawodników	1148
zadało przeciwnikowi dwa razy więcej ciosów, niż spalonych. Dwudziestotrzyletni wówczas Buster	1149
trafił prawym prostym. Teraz, na dziesięć metrów, już w pierwszej rundzie Floyda Pattersona spotkał	1150
się z Joem Frazierem, Muhammadem Alim, Tunneyem i Michaelem Moorerem. O 3.00 rano, kiedy	1151
stoczyli już swoje ostatnie rugie, przyszedł czas na jeszcze jedną walkę z "Milionerem", "Drakula" z	1152
Jackiem Dempseyem.	1153
	1154
mT5-base: W najbliższą sobotę w Las Vegas dojdzie do kolejnego "rewanżu stulecia". Zmierzą	1155
się w nim dwaj najwybitniejsi obecnie bokserzy wagi ciężkiej, Amerykanin Evander Holyfield i	1156
Brytyjczyk Lennox Lewis. Osiem miesięcy temu, 13 marca, w nowojorskiej Madison Square Garden	1157
ich pierwsza walka zakończyła się kontrowersyjnym remisem. Stawką rewanżowego pojedynku	1158
będą pasy mistrzowskie trzech najważniejszych organizacji boksu zawodowego: WBA, WBC, IBF.	1159
pIBART: W najbliższą sobotę w Las Vegas dojdzie do kolejnego "rewanżu stulecia". Zmierzą	1160
się w nim dwaj najwybitniejsi obecnie bokserzy wagi ciężkiej, Amerykanin Evander Holyfield i	1161
Brytyjczyk Lennox Lewis. Osiem miesięcy temu, 13 marca, w nowojorskiej Madison Square Garden	1162
ich pierwsza walka zakończyła się kontrowersyjnym remisem. Stawką rewanżowego pojedynku	1163
będą pasy mistrzowskie trzech najważniejszych organizacji boksu zawodowego: WBA, WBC, IBF.	1164
Ostatnim posiadaczem trzech pasów w wadze ciężkiej był Riddick Bowe. Od siedmiu lat żaden z	1165
pięściarzy nie może się uważać za króla tej najbardziej cenionej kategorii.	1166
polT-base: W najbliższą sobotę w Las Vegas dojdzie do kolejnego "rewanżu stulecia". Zmierzą	1167
się w nim dwaj najwybitniejsi obecnie bokserzy wagi ciężkiej, Amerykanin Evander Holyfield i	1168
Brytyjczyk Lennox Lewis. Osiem miesięcy temu, 13 marca, w nowojorskiej Madison Square Garden	1169
ich pierwsza walka zakończyła się kontrowersyjnym remisem. Stawką rewanżowego pojedynku	1170
będą pasy mistrzowskie trzech najważniejszych organizacji boksu zawodowego: WBA, WBC, IBF.	1171
Ostatnim posiadaczem trzech pasów w wadze ciężkiej był Riddick Bowe. Od siedmiu lat żaden z	1172
pięściarzy nie może się uważać za króla tej najbardziej cenionej kategorii.	1173

H Example predictions on PSC (extract)

- 1174
- 1175 1. **Gold Summary:** Politycy z pierwszych stron gazet szykują się do objęcia stanowisk prezydentów
1176 dużych miast. Dla nich gotowi są nawet zrezygnować z mandatów poselskich.
1177 W najbliższych tygodniach rozstrzygnie się, czy w jesiennych wyborach samorządowych prezydenci miast wybierani
1178 będą bezpośrednio przez mieszkańców. W Sejmie trwają prace nad projektem ustawy o bezpośrednim
1179 wyborze wójtów, burmistrzów i prezydentów. Ustawa ma poparcie SLD i PO. Jako kandydatów
1180 SLD na stanowisko w stołecznym ratuszu wymienia się Ryszarda Kalisza i Danutę Waniek. "Nie
1181 wyklucza" kandydowania w Warszawie Lech Kaczyński z Prawa i Sprawiedliwości. Trzecim
1182 liczącym się kandydatem może być Andrzej Olechowski, lider Platformy Obywatelskiej. Nie stanie
1183 w szranki natomiast były prezydent Paweł Piskorski.

1184 **mT5-small:** Dla nich gotowi są nawet zrezygnować z mandatów poselskich. W najbliższych
1185 tygodniach rozstrzygnie się, czy w jesiennych wyborach samorządowych prezydenci miast wybierani
1186 będą bezpośrednio przez mieszkańców. Jako kandydatów SLD na stanowisko w stołecznym ratuszu
1187 wymienia się Ryszarda Kalisza i Danutę Waniek. W partii krąży opinia, że Kalisz sam się zgłosił.

1188 **polT-small:** Politycy szykują się do objęcia stanowisk prezydentów dużych miast. Dla nich gotowi
1189 są nawet zrezygnować z mandatów poselskich. W najbliższych tygodniach rozstrzygnie się, czy
1190 w jesiennych wyborach samorządowych prezydenci miast wybierani będą bezpośrednio przez
1191 mieszkańców.

1192 **papuGaPT2:** Partia ma poważne problemy z przywództwem. Nie ma mocnego nazwiska.
1193 Na razie próbuje się rozwiązywać kolejne oświęcimskie kryzysy w drodze wykorzystywania
1194 kruczków prawnych. Niedawno wyszły na jaw kontakty radnego z półświatkiem. W Szczecinie
1195 najpoważniejszym z kandydatów jest Edmund Runowicz, znany wszystkim z Gliwic lider ZChN.
1196 Z racji wpływów jakie ma partyjniactwo w mieście, chcą go popierać. Próba stworzenia drugiego
1197 centrum politycznego nie powiodła się. SLD zrezygnował z pomysłu utworzenia szerokiej koalicji
1198 ugrupowań posierpniowych. Przy budżecie państwa 181,6 mld działalność biura wyborczego kosztuje
1199 0,041 proc. państwowych wydatków. Tymczasem w ramach walki z bezrobociem szuka ratunku
1200 w Urzędzie Miasta.

1201 **mT5-base:** W najbliższych tygodniach rozstrzygnie się, czy w jesiennych wyborach samorządowych
1202 prezydenci miast wybierani będą bezpośrednio przez mieszkańców. Tymczasem konkurenci sięgają
1203 po najlepsze kadry. Sojusz ma kłopot, bo nie ma drugiego tak dobrego kandydata. Jak nie Borowski,
1204 to kto? Jako kandydatów SLD na stanowisko w stołecznym ratuszu wymienia się Ryszarda Kalisza i
1205 Danutę Waniek. Te dwie osoby mają w stolicy duże poparcie.

1206 **plBART:** Politycy szykują się do objęcia stanowisk prezydentów dużych miast. Dla nich gotowi
1207 są nawet zrezygnować z mandatów poselskich. W najbliższych tygodniach rozstrzygnie się, czy
1208 w jesiennych wyborach samorządowych prezydenci miast wybierani będą bezpośrednio przez
1209 mieszkańców. Jako kandydatów SLD na stanowisko w stołecznym ratuszu wymienia się Ryszarda
1210 Kalisza i Danutę Waniek. "Nie wyklucza" kandydowania w Warszawie Lech Kaczyński z Prawa i
1211 Sprawiedliwości. Trzecim liczącym się kandydatem może być Andrzej Olechowski, lider Platformy
1212 Obywatelskiej.

1213 **polT-base:** Czy w jesiennych wyborach samorządowych prezydenci miast wybierani będą bezpośrednio
1214 przez mieszkańców. W Sejmie trwają prace nad projektem ustawy o bezpośrednim wyborze
1215 wójtów, burmistrzów i prezydentów. Ustawa ma poparcie głównych sił politycznych, jak SLD i
1216 PO, jest więc wielce prawdopodobne, iż zostanie uchwalona. Nic dziwnego, że już dzisiaj partie
1217 zastanawiają się, kim obsadzić najwyższe stanowiska w miastach. Dla niektórych polityków odejście
1218 do samorządu oznaczałoby konieczność złożenia mandatu posła, ale nie zraża ich to.

- 1219 2. **Gold Summary:** Ludzie całkowicie sparaliżowani będą mogli samą myślą sterować komputerem,
1220 by porozumieć się z otoczeniem. Niezwykle czułe detektory, odczytujące wybrane fale mózgowy,
1221 testowane są w laboratoriach USA, Japonii i Europy. Do łączności z urządzeniami elektronicznymi

można wykorzystać jakąkolwiek aktywność układu nerwowego - sygnały elektryczne wytwarzane ruchami gałek ocznych lub impulsy nerwów mięśni.	1222 1223
mT5-small: Uniwersytetu Emory w Atlancie wypróbowali u dwóch pacjentów system pozwalający sterować kursorem na ekranie monitora. Sterowanie komputerem samą myślą stało się możliwe, gdy uczeni lepiej opanowali odczytywanie i werbalizację sygnałów. Komunikowanie się przy użyciu takiego systemu jest dość żmudne, ale wyłącznie dla ludzi zdrowych. Dla osób całkowicie sparaliżowanych, nie mogących wykonywać żadnych ruchów, nawet mięśniami twarzy, jest jedynym sposobem porozumiewania.	1224 1225 1226 1227 1228 1229
poIT-small: Ludzie całkowicie sparaliżowani będą mogli samą myślą sterować komputerem, by porozumieć się z otoczeniem. Niezwykle czułe detektory, odczytujące wybrane fale mózgowe, testowane są w laboratoriach USA, Japonii i Europy. Do łączności z urządzeniami elektronicznymi można wykorzystać jakąkolwiek aktywność układu nerwowego - sygnały elektryczne wytwarzane ruchami gałek ocznych lub impulsy nerwów mięśni. Jeszcze niedawno wszelkie rozważania na ten temat były domeną wyłącznie autorów powieści fantastycznonaukowych. Obecnie naukowcy w kilkunastu ośrodkach prześcigają się w konstruowaniu urządzeń coraz sprawniej wykrywających komunikaty myślowe.	1230 1231 1232 1233 1234 1235 1236 1237
papuGaPT2: dzisiaj coraz więcej teoretyków na temat tego, co może się stać, jak będzie możliwe, za kilka lat. Trudno wyobrazić sobie wielką liczbę skomplikowaną, kręcących się po krzywej Kenii i Epsilon. Gdyby jednak wziąć pod uwagę normalny, a nie ten podkolorowany obraz powyborczej rzeczywistości, to rzeczywiście wyglądałby nędzny. Jeszcze niedawno takie pytanie byłoby niedorzeczne. Dziś naukowcy nie są już takimi sceptykami. \nNiektórzy jeszcze przed przystąpieniem Polski do UE zastanawiali się, czy w ogóle wchodzić w ten projekt. Nie wszyscy są entuzjastami nowego tworu, który stał się również przedmiotem marzeń i zazdrości milionów. A Polska? Jak wynika z sondaży, 77 proc. Amerykanów gotowych jest poprzeć likwidację lub gruntowną zmianę europejskich norm. Przeciwnicy wejścia do Unii Europejskiej obawiają się takiego myślenia, bo będzie on trwał tylko trzy i pół roku. Natomiast w rzeczywistości jest to proces bardzo długi.	1238 1239 1240 1241 1242 1243 1244 1245 1246 1247
mT5-base: Jeszcze niedawno wszelkie rozważania na ten temat były domeną wyłącznie autorów powieści fantastycznonaukowych. Obecnie uczeni w kilkunastu ośrodkach prześcigają się w konstruowaniu urządzeń coraz sprawniej wykrywających komunikaty myślowe. Niezwykle czułe detektory, odczytujące wybrane fale mózgowe, testowane są w laboratoriach USA, Japonii i Europy. Kilka takich rozwiązań jest testowanych w USA.	1248 1249 1250 1251 1252
pIBART: Ludzie całkowicie sparaliżowani będą mogli samą myślą sterować komputerem, by porozumieć się z otoczeniem. Niezwykle czułe detektory, odczytujące wybrane fale mózgowe, testowane są w laboratoriach USA, Japonii i Europy. Do łączności z urządzeniami elektronicznymi można wykorzystać jakąkolwiek aktywność układu nerwowego. Urządzenie ma jednak tę wadę, że wymaga użycia wszczepów, grożących powstaniem infekcji i uszkodzeniem mózgu.	1253 1254 1255 1256 1257
poIT-base: ZBIGNIEW WOJTASIŃSKI Ludzie całkowicie sparaliżowani będą mogli samą myślą sterować komputerem, by porozumieć się z otoczeniem. Niezwykle czułe detektory, odczytujące wybrane fale mózgowe, testowane są w laboratoriach USA, Japonii i Europy. Do łączności z urządzeniami elektronicznymi można wykorzystać jakąkolwiek aktywność układu nerwowego - sygnały elektryczne wytwarzane ruchami gałek ocznych lub impulsy nerwów mięśni.	1258 1259 1260 1261 1262

I Example predictions on PSC (abstract)

1. **Gold Summary:** Anachroniczna jest koncepcja, u nas bynajmniej nierzadko wyznawana, że swoboda umów wyklucza ideę ochrony konsumenta. Przed kilkoma laty niemiecki Trybunał Konstytucyjny uznał, że nawet osoba pełnoletnia, samodzielna, nie poddana żadnemu przymusowi wymaga ochrony ręcząc za kredyt bankowy oraz że konieczna jest szczegółowa, wyczerpująca, niedwuznaczna, jasna informacja, ze wskazaniem na kwestie najbardziej niebezpieczne dla poręczyciela. Lepsza informacja dla konsumenta łagodzi bowiem nierówność pozycji rynkowej. W konsekwencji tych orzeczeń zmieniła się praktyka powszechnych sądów w Niemczech. Trzy kwestie zasługują tu na uwagę. Po pierwsze - umowy kredytowe i sytuacja poręczyciela doczekały się w Niemczech oceny TK, dokonywanej z konstytucyjnego punktu widzenia. Po drugie - TK za remedium na strukturalne zachwianie równowagi umownej uznał zwiększenie obowiązków informacyjnych kontrahenta konsumenta. Po trzecie wreszcie - opisywana sytuacja jest kolejnym przykładem tego, jak dalece anachroniczna jest koncepcja (u nas bynajmniej nierzadko wyznawana), że swoboda umów wyklucza ideę ochrony konsumenta. Z konstytucji nie można wyczytać, jakimi środkami i na jakim poziomie ma się chronić konsumenta, ale to nie jest jedyny możliwy sposób "użycia konstytucji" do takich celów. Ustawa zasadnicza da się użyć jako norma rozstrzygająca na wypadek kilku możliwych interpretacji jakiegoś przepisu: "prokonsumenckiej" (w zakresie wymienionych w konstytucji, szczególnie chronionych praw konsumenta) i "antykonsumenckiej" lub choćby "konsumencko neutralnej". Albo na wypadek interpretacji norm blankietowych, klauzul generalnych czy zwrotów niedookreślonych, które sąd musi odkodować, nadać im konkretną treść. Drugą ważną kwestią jest znaczenie informacji jako oręża konsumenta. W europejskim prawie wspólnotowym kamieniem węgielnym ochrony konsumenta jest informacja. Uważa się, że konsument wymaga ochrony, ponieważ jest źle poinformowany i na skutek tego nie może w prawdziwie wolny i nieskrępowany sposób decydować o swym "udziale na rynku". Stąd się biorą niezwykle rozbudowane w dyrektywach i ich implementacjach wewnątrz krajowych przepisy mówiące, o czym, kiedy i jak trzeba konsumenta informować. Nie tylko jednak sama obrona przez informację jest cechą charakterystyczną europejskiego prawa konsumenckiego. Występuje tu jeszcze wskazanie minimalnego poziomu treści umowy. Oznacza to, że poszczególne państwa, implementując dyrektywy konsumenckie we własnym porządku prawnym, nie mogą zejść niżej wskazanego tam poziomu ochrony. Ale tutaj tkwi pułapka: przyznanie w prawie krajowym zbyt wysokiego poziomu ochrony, może być uznane za rzeczywistą czy ukrytą praktykę dyskryminacyjną, ograniczającą dostęp towarów i usług na rynek wspólny. Nie tylko informacja i nie tylko określenie minimalnego poziomu uprawnień zagwarantowanych konsumentowi w dotyczącym danej umowy prawie wspólnotowym czy wewnętrznym są cechami charakterystycznymi europejskiego prawa konsumenckiego. Otóż cechą charakterystyczną współczesnego prawa umów jest to, że konsumentowi przysługuje prawo "wycofania się z transakcji". W praktyce i u nas ta nowa technika ochronna zaczyna się pojawiać: zna ją polskie prawo ubezpieczeniowe. Co charakterystyczne: tego rodzaju możliwość wycofania się z transakcji służy (w prawie europejskim) jedynie konsumentowi, nie jego kontrahentowi. I znów jest to jeden ze środków, który przyznaje się tylko stronie instytucjonalnie słabszej, używany, by zrównoważyć siłę rynkową konsumenta wobec profesjonalisty.

mT5-small: Pod koniec laty dwa orzeczenia niemieckiego Trybunału Konstytucyjnego zbulwersowały świat prawniczy. Trybunał uznał, że anachroniczna jest koncepcja, u nas bynajmniej nierzadko wyznawana, jako swoboda umów wyklucza ideę ochrony konsumenta. Po trzecie - umowy kredytowe i sytuacja poręczyciela doczekały się w Niemczech oceny TK, dokonywanej z konstytucyjnym punktu widzenia.

polT-small: Przed kilkoma laty dwa orzeczenia niemieckiego Trybunału Konstytucyjnego zbulwersowały świat prawniczy. Trybunał uznał, że nawet osoba pełnoletnia, samodzielna i niepoddana żadnemu przymusowi wymaga ochrony ręcząc za kredyt bankowy. Drugą istotną kwestią europejskiego prawa konsumenckiego jest wskazanie minimalnej wiedzy i doświadczenia między poręczycielem i bankiem. Innym przykładem środka ochrony konsumenta jest niemiecki Trybunał

Konstytucyjny. Moim zdaniem, konsumentowi przysługuje prawo wycofania się z transakcji bez negatywnych konsekwencji.	1313 1314
papuGaPT2: W porównaniu z prawem polskim ustawodawstwo i orzecznictwo krajów europejskich znacznie bardziej chroni prawa konsumentów, co w konsekwencji doprowadzić może do unikania przez banki i inne instytucje finansowe prokonsumenckie wszelkiej informacji o towarze i bezpiecznym sposobie korzystania z niego. Europejskie dyrektywy konsumenckie wskazują minimalny poziom treści umowy: są to przepisy zakazujące reklamy pewnych produktów czy usług (lekarstwa, alkohole, wolne zawody), używania zrozumiałego w określonych mediach (radio, tv), reklamy kierowanej do pewnych osób (dzieci) czy też zniósł dobrze ochronę przed wprowadzającymi w błąd informacjami na etykietach. Wreszcie cechą charakterystyczną współczesnego prawa umów jest to, iż konsument przez krótki czas (siedem do dziesięciu dni) ma prawo, bez żadnych negatywnych konsekwencji, zrezygnować z transakcji. Środek ten przysługuje tylko konsumentowi jako stronie instytucjonalnie słabszej; jego kontrahenta obowiązuje ogólna zasada mówiąca o "umowy powinny zostać dotrzymane". Zasada ta wyraża się w łacińskiej formule "in dubio contra proferentem" i jest ważnym instrumentem prawa cywilnego, który służy nie tylko ochronie konsumenta, ale także jako środek prawnokarnego (art. 76), który ma zrównoważyć siłę przyciągania i działać na rzecz zrekompensowania braków jego wiedzy. Stąd się biorą rozbudowane przepisy mówiące, o czym, kiedy i jak należy informować konsumenta o wszelkich aspektach umowy kredytowej, a także o możliwych niebezpieczeństwach i ryzykach, które mogą być nieświadomym wyrazem orientacji konsumenta w transakcjach rynkowych.	1315 1316 1317 1318 1319 1320 1321 1322 1323 1324 1325 1326 1327 1328 1329 1330 1331 1332
mT5-base: Rę swoboda umów wyklucza ideę ochrony konsumenta.	1333
plBART: Przed kilkoma laty dwa orzeczenia niemieckiego Trybunału Konstytucyjnego zbulwersowały świat prawniczy. Trybunał uznał, że nawet osoba pełnoletnia, samodzielna, nie poddana żadnemu przymusowi wymaga ochrony rękając za kredyt bankowy, bo nierównowaga wiedzy i doświadczenia między poręczycielem i bankiem jest zbyt wielka i tolerowanie jej zagraża jednej z konstytucyjnych zasad niemieckiego porządku prawnego: wyrażonej w art. 2 konstytucji autonomii woli. W konsekwencji orzeczeń TK zmieniła się praktyka niemieckich sądów, które umowy kredytowe doczekały się oceny TK, dokonywanej z konstytucyjnego punktu widzenia. Po drugie, TK uznał za remedium na strukturalne zachwianie równowagi umownej. Opisywana sytuacja jest kolejnym przykładem tego, jak dalece anachroniczna jest koncepcja (u nas bynajmniej nierzadko wyznawana), że swoboda umów wyklucza ideę ochrony konsumenta.	1334 1335 1336 1337 1338 1339 1340 1341 1342 1343
polT-base: W.	1344
2. Gold Summary: Praca Huberta Orłowskiego "Polnische Wirtschaft. Nowoczesny niemiecki dyskurs o Polsce" poświęcona jest stereotypom dotyczącym stosunków polsko-niemieckich, w tym jednemu z najsilniejszych - "polnische Wirtschaft". Autorstwo tego stereotypu przypisuje się Georgowi Fosterowi, uczonemu i podróżnikowi, który w swoich listach pisanych pod koniec XVIII w. terminem tym określał kombinację bałaganu, nieudolności, brudu i lenistwa. W stereotyp myślenia Niemców o Polsce wpisywały się wydarzenia historyczne. Zrywy niepodległościowe Polaków od powstania listopadowego do "Solidarności" traktowane były często jako przejaw polskiej skłonności do anarchii. We współczesnych Niemczech stereotyp "polnische Wirtschaft" jest wciąż żywy. Podtrzymują go zarówno media, jak i w literatura piękna.	1345 1346 1347 1348 1349 1350 1351 1352 1353
mT5-small: Nakładem wydawnictwa Huberta Orłowskiego poświęcona jest stereotypom dotyczącym stosunków polsko-niemieckich, ze szczególnym uwzględnieniem najsilniejszych z nich: „polnische Wirtschaft”. Autorstwo tego stereotypu przypisuje się niejakiemu Georgowi Fosterowi, uczonemu i podróżnikowi z Żuław Wiślanych. "Polnische Wirtschaft" oznaczającemu kombinację tych cech. Wśród warstw bogatszych widziano rozrzutność, równocześnie jednak także brak higieny, zacofanie, ignorancja.	1354 1355 1356 1357 1358 1359
polT-small: W książce profesora Huberta Orłowskiego poświęcona jest stereotypom dotyczącym stosunków polsko-niemieckich, ze szczególnym uwzględnieniem najsilniejszych z nich: "polnische	1360 1361

1362 Wirtschaft" oznaczającym kombinację bałaganu, nieudolności, brudu i złą gospodarkę.

1363 **papuGaPT2:** Od lat 40. literatura podtrzymywała mit szlachetnego Polaka, o którym zapomina
1364 młodsze pokolenie. Stereotyp ten nigdy by może nie zaistniał, gdyby nie głębokie różnice dzielące
1365 polskie cnoty szlacheckie od niemieckich cnót mieszczańskich, pogłębione różnicami między ka-
1366 tolicyzmem i protestantyzmem. Po klęsce powstania listopadowego, po upadku muru berlińskiego
1367 okazało się, że ustrój ten nie był w stanie zdobyć uznania i poparcia ludu, dlatego też w książkach,
1368 operetkach i filmach pełno jest dowcipów o polskiej gospodarce i polskim sejmikowaniu. Profesor
1369 Hubert Orłowski jest współautorem tzw. ustawy teczkowej, ustanawiającej Instytut Pamięci Naro-
1370 dowej, oraz ustawy o ochronie informacji niejawnych, określającej m.in. mianem "pacta servanda
1371 sunt" na której spoczywa cała odpowiedzialność.

1372 **mT5-base:** W niemieckich zbiorów przysłów o Polsce, ukazała się dwa lata temu w Wiesbaden
1373 nakładem wydawnictwa Otto Harrasowitza na polski.

1374 **plBART:** Praca profesora Huberta Orłowskiego "Polnische Wirtschaft. Nowoczesny niemiecki
1375 dyskurs o Polsce" poświęcona jest stereotypom dotyczącym stosunków polsko-niemieckich, ze
1376 szczególnym uwzględnieniem najsilniejszych z nich: „polnischer Wirthschaft”. Termin ten określa
1377 "nieporządek i złą gospodarkę, niedbalstwo i bałagan". Autorstwo tego, mającego długą historię
1378 i niesłuchaną żywotność, stereotypu przypisuje się niejakiemu Georgowi Fosterowi, uczonemu
1379 i podróżnikowi z Żuław Wiślanych. Naprawdę jednak formuła była już gotowa i niezłe znana
1380 wcześniej, choć w myśleniu Niemców (i nie tylko ich) zdomowała się dopiero ok. 1830 r. po edycji
1381 listów Fostera. W liście do swojego wierzyciela pisał z Wilna w 1784 roku: "O nieopisanym brudzie,
1382 lenistwie, opilstwie i nieudolności całej służby..., o niezdarności rzemieślników, ich niesłuchanie
1383 wiejskiej robocie, wreszcie o zadowoleniu Polaków [Polaken] z własnego bagienka, a także ich
1384 przywiązaniu do rodzinnych zwyczajów nie chcę pisać już nic więcej, aby nie przedłużać tego listu".

1385 **polT-base:** W.