
CEDILLE: A LARGE AUTOREGRESSIVE LANGUAGE MODEL IN FRENCH

Martin Müller*

Florian Laurent*

Cedille AI¹
hello@cedille.ai

ABSTRACT

Scaling up the size and training of autoregressive language models has enabled novel ways of solving Natural Language Processing tasks using zero-shot and few-shot learning. While extreme-scale language models such as GPT-3 offer multilingual capabilities, zero-shot learning for languages other than English remain largely unexplored. Here, we introduce Cedille, a large open source auto-regressive language model, specifically trained for the French language. Our results show that Cedille outperforms existing French language models and is competitive with GPT-3 on a range of French zero-shot benchmarks. Furthermore, we provide an in-depth comparison of the toxicity exhibited by these models, showing that Cedille marks an improvement in language model safety thanks to dataset filtering.

1 Introduction

Large autoregressive language models have drawn wide attention due to their zero-shot and few-shot capabilities, allowing them to be used for a wide variety of Natural Language Processing tasks without the need for task-specific finetuning or annotation data [1, 2]. Additionally, previous work highlights the improved sample and compute efficiency of larger models, generally justifying the move towards larger models [3].

Although large language models, such as GPT-3 [2], have been trained on multilingual corpuses, the performance on NLP tasks may vary significantly between languages. Assessing zero-shot performance in non-English languages is challenging due to the limited number of human-curated benchmarks available. However, with the exception of recent work in machine translation [4], multilingual models generally perform worse than mono- or bilingual language models [5].

Monolingual autoregressive language models in French have previously been proposed. GPT-fr [6] and PAGnol [7] have been trained on filtered versions of Common Crawl² and CCNet [8], respectively. Both works highlight the importance of deduplicating and filtering of pre-training data and use decoder-only transformer architectures, closely following the GPT models with model sizes reaching 1B and 1.5B parameters, respectively. It’s worth noting that these works do not directly compare performance against extreme-scale large multilingual models, such as GPT-3, in particular with regard to zero-shot tasks.

Previous work on the various encoding biases in large language models highlights the importance of dataset curation and documentation [9, 10]. Experiments conducted on GPT-3 (which has been trained on 570GB of text data from Common Crawl) show that the model may generate toxic sentences even when prompted with non-toxic text [11]. Although applying filtering of training data using automated toxicity scores may introduce classifier-specific biases [12], this technique remains more effective than

* Authors contributed equally, order is random

¹Coteries SA, EPFL Innovation Park, Lausanne, Switzerland

²<https://commoncrawl.org/>

decoder-based detoxification using methods such as swear word filters, PPLM [13], soft prompt tuning [14] or toxicity control tokens [15].

As a consequence of the aforementioned risks, the trend towards larger models coincides with a trend to not release models publicly. Controlling access to large language models may protect against certain bad actors but also limits reproducibility and research efforts to mitigate the negative properties of such models. In a push for building models in the open, EleutherAI, a grassroot collective of researchers, released GPT-J [16], a 6B parameter English language model. This model was trained on the Pile [20], a 825GB text corpus by the same collective.

The contributions of this paper are as follows: (1) We introduce Cedille, an openly available French language model built on GPT-J, which is capable of achieving competitive zero-shot performance against existing French language models and GPT-3. (2) We release the toxicity scores of the complete French C4 dataset, and (3) we provide a comparison of Cedille’s toxicity to other language models (including GPT-3).

2 Methods

2.1 Model architecture

Our model architecture is identical to GPT-J [16]. GPT-J uses a similar transformer architecture to the one used in 6.7B GPT-3 with three main differences: (1) No sparse attention patterns were used; (2) the dimension of the attention head was increased from 128 to 256; and (3) Rotary positional embeddings [17] were used instead of sinusoidal embeddings. See Table 1 for more details.

Number of parameters	6,053,381,344
Number of layers N	28
Model dimensions d_{model}	4096
Feed-forward dimension d_{ff}	16,384
Number of attention heads n_{heads}	16
Head dimension d_{head}	256
Context size	2048
Vocab size	50,257

Table 1: Cedille model details.

2.2 Training data

Cedille is trained on a filtered version of the French part of the multilingual C4 (mC4) dataset [18], which contains 332M documents or 1.1TB of uncompressed text. mC4 is

extracted from 71 Common Crawl snapshots (years 2013 to 2020) and uses CLD3³, a small feed-forward neural network, for language identification. mC4 filtered out pages of less than three lines of at least 200 characters.

We apply two different forms of filtering to the dataset 1) toxicity filtering using the Detoxify model [19] and 2) loss filtering using the FlauBERT model [20]. For both filtering steps we compute the metric on a per document level of the entire base dataset. In some cases chunking the documents into splits of 1200 characters was necessary due to the fixed context size of the used models. Chunks smaller than 600 characters were not evaluated. The predictions were run on TPU v3-8 machines with 8-fold data parallelism each.

Each percentile as well as the tails of both the loss and the toxicity distribution were sampled and manually inspected to find suitable cut-off values for filtering. The inspection of these samples revealed that both toxicity and loss values were appropriate⁴. We removed documents corresponding to a toxicity score higher than 0.5, corresponding to 0.25% of the content (0.8M documents). For the loss filtering we considered the loss distribution of each of the 2048 files and removed documents below a 0.2 percentile loss (corresponding to a loss value of roughly 4.5) and above an absolute loss value of 10. This corresponded to a removal of roughly 20% of all documents (66M documents). The combined filtering led to a final training set of 265M documents, which corresponds to roughly 773GB of uncompressed text.

The text was then run through the `fix_text` method of the Python library `ftfy` [21] using NFKC normalization and encoded using the unmodified GPT-2 tokenizer. Documents were simply concatenated and split into samples of 2049 tokens. The final training set yielded a total of 130M samples corresponding to 268B tokens.

2.3 Training process

Cedille was trained starting from the official GPT-J model checkpoint using the mesh-transformer-jax codebase [22]. Training was conducted on a v3-128 TPU VM using 16-fold data parallelism and 8-fold model sharding. For all our experiments we used an effective batch size of 256. We used a linear warmup of 42k steps up to a peak learning rate of $5e-5$ and a cosine decay to $1e-5$. Weight decay was set to 0.1. Cedille was trained for 150k steps, which corresponds to 0.3 epochs on the training set or 78.7B tokens. The starting and final training perplexities were 6.13 and 3.89, respectively. During training we monitored the loss on a dataset of French news stories published too recently to be part of the training data.

³<https://github.com/google/clD3>

⁴Despite the positive visual inspection a bug in the loss computation was discovered much later in the analysis. Further investigation revealed that roughly 10% of samples were wrongly included in the final dataset as a result. Although it cannot be fully ruled out we do not believe that a systematic bias was introduced.

2.4 Evaluation

Zero-shot performance was evaluated using a forked version of the lm-evaluation-harness codebase [23]. In particular, we added a different way of evaluating perplexity using strides (see section 3.1), implemented the various benchmarks discussed in this work, and integrated the mesh-transformer-jax library (for evaluating checkpoints on TPUs) and the Pagnol model families. Benchmarking was conducted on v3-8 TPU VMs and on A100 GPUs.

Toxicity evaluation was conducted using a modified version of the real-toxicity-prompts codebase⁵. The main difference is the use of the Detoxify model in order to predict toxicity (see section 4). Our adapted codebase is available at <https://github.com/coterie/real-toxicity-prompts>.

3 Tasks

3.1 Perplexity

Model	#params	Byte-PPL	Token-PPL
GPT-3 (ada)	1.3B ^a	1.930	7.952
GPT-3 (babbage)	6.7B	1.973	6.447
GPT-3 (curie)	13B	1.809	5.082
GPT-3 (davinci)	175B	1.656	3.993
GPT-J	6.05B	1.746	5.797
Cedille	6.05B	1.646	3.932
Pagnol (small)	124M	1.852	17.802
Pagnol (medium)	335M	1.775	14.623
Pagnol (large)	773M	1.725	12.791
GPT-fr (base)	1B	2.090	11.882

Table 2: Byte-level and token-level perplexity scores on the WikiText-fr benchmark (lower is better).

^aOpenAI hasn’t officially disclosed the size of the models provided by their API, however recent experiments suggest the mapping presented in the table [24].

Zero-shot perplexity was evaluated on the test subset of the WikiText-fr⁶ dataset [6], containing articles from the French Wikipedia which are part of the “quality articles” or “good articles” categories, similar to the English WikiText-103 dataset [25]. The test set contains 589k words or 3.7M characters of cleaned French text from 60 articles. We evaluated perplexity by concatenating the text without further preprocessing and using a sliding window approach [26] with a stride of 512 tokens. Therefore models with a context window of 1024 tokens (GPT-fr, Pagnol) had 512 tokens of context, whereas models with a context window of 2048 tokens had 1536 tokens of context. Table 2 shows

the summed log likelihoods both normalized by number of characters and by number of tokens. Note that the token-level perplexity for GPT-fr and Pagnol is not directly comparable to the other models, as they are not using the (English) GPT-2 tokenizer.

Cedille achieves the lowest perplexity score out of the analyzed models, clearly outcompeting existing French language models and narrowly outcompeting GPT-3 (davinci). Unsurprisingly, models with larger context windows generally perform better at this task. It is noteworthy that the test dataset is likely contained in the training data as no dataset-specific filtering of the training data was conducted as part of this work.

3.2 Summarization

We evaluated the summarization capabilities on the OrangeSum benchmark, as introduced in the BARThez work [27] as a French equivalent of XSum [28]. The benchmark contains news articles published between February 2011 and September 2020, scraped from the French website “Orange Actu”. The models were given the news article in the test subset using the following prompt:

{article text}\nPour résumer :

The models were tasked to generate 100 tokens using top- k of 2 and a temperature of 1, following the methodology in [1]. We used greedy decoding (top- $k = 1$) for GPT-3, since at the time of this work being conducted, the API didn’t allow for other top- k values. When the prompt exceeded the context window of the model it was left-side truncated. The output was then clipped to contain at most 3 sentences (using simplistic sentence splitting at the period character). Table 3 shows the ROUGE score [29] of the output compared to the title of the corresponding articles.

Model	R_1	R_2	R_L
GPT-3 (ada)	13.95	4.75	11.59
GPT-3 (babbage)	4.62	1.76	3.86
GPT-3 (curie)	5.28	2.21	4.42
GPT-3 (davinci)	15.49	5.82	13.05
GPT-J	14.46	4.72	11.68
Cedille	14.74	4.83	11.86
Pagnol (small)	8.52	1.61	7.24
Pagnol (medium)	8.98	1.86	7.55
Pagnol (large)	9.19	1.85	7.71
GPT-fr (base)	10.15	2.60	8.27

Table 3: Performance of summarization in French. Shown are the ROUGE scores on the OrangeSum dataset (higher is better).

Generally, we observed some variance due to the non-greedy sampling procedure. However, computational limi-

⁵<https://github.com/allenai/real-toxicity-prompts>

⁶https://huggingface.co/datasets/asi/wikitext_fr

tations and cost made it difficult to estimate this variance. We also observed that the choice of the prefix (“Pour résumer :”) strongly influences the scores. Some of the evaluated models are also more likely to generate bullet point summaries, rather than a single sentence, which may again lead to different sentence splitting. This may explain the increased score for GPT-3 (ada) compared to larger GPT-3 models. Nevertheless, the scores provided in Table 3 give some rough indication of summarization performance.

3.3 Question Answering (QA)

Question answering (QA) was evaluated on FQuAD (French Question Answering Dataset) [30], a dataset inspired by the English SQuAD equivalent [31]. The models were evaluated on the validation subset, which contains 3188 human-curated question-answer pairs, based on 768 high-quality French Wikipedia articles.

Model	F1	Exact match (%)
GPT-3 (ada)	19.09	4.48
GPT-3 (babbage)	26.16	8.81
GPT-3 (curie)	39.49	17.84
GPT-3 (davinci)	-	-
GPT-J	26.14	6.96
Cedille	34.59	12.23
Pagnol (small)	10.66	0.43
Pagnol (medium)	13.80	0.84
Pagnol (large)	17.67	2.72
GPT-fr (base)	15.15	2.03

Table 4: Question-answering F1 and exact match scores in French on the FQuAD benchmark (higher is better).

The models were evaluated using the SQuAD v2 metric [31], which also takes into consideration “no answer” probabilities, i.e. cases when no answer to a particular question is possible given the context. The models were tasked to generate 100 tokens and at most 1 sentence using greedy sampling and the following prompt:

Titre: {title}\nContexte: {context}\n\nQuestion: {question}\n\nRéponse:

The “no answer” probabilities were calculated against the string:

{prompt} Sans réponse.

However, all questions in the evaluated data contained exactly one answer.

The results in Table 4 show that GPT-3 is very competitive on this task, with GPT-3 (curie) outperforming Cedille and all other evaluated models. GPT-3 (davinci) was not evaluated on this task for cost reasons, as OpenAI did not support our request for funding at the time of writing. The

results may be contrasted to a finetuned version of Camembert [32] which yields F1 of 88% and best match of 78% on this dataset [30].

3.4 Translation

Zero-shot translation was evaluated for the language pair English and French on the WMT14 dataset [33]. Traditionally, such benchmarks are evaluated using the BLEU score [34]. The datasets contains 3003 samples each and are provided by the sacrebleu library [35]. The zero-shot task is formulated using the following pattern:

{source_lang} phrase: {text}\n{target_lang} phrase:

Where source_lang and target_lang are French and English, respectively, depending on the direction. Greedy sampling is used to generate 256 tokens. The output was clipped to at most 1 sentence.

Cedille outperforms other models for the direction English to French, highlighting the strong French writing capabilities (see Table 5). Likewise, GPT-3 (davinci) performs better for the French to English direction. Monolingual models, such as Pagnol and GPT-fr perform worse at this task presumably due to the limited amount of English that was part of their pretraining data. Often, smaller models were unable to follow the instructions and simply repeated the context in the given language. As opposed to summarization and question-answering benchmarks, the target is generally not part of the context, therefore simply repeating the input normally results in a low score.

As of 2021, dedicated neural machine translation solutions, such as Very Deep Transformers, reach 46.4 BLEU for English to French translation [36].

Model	BLEU (en→fr)	BLEU (fr→en)
GPT-3 (ada)	2.71	16.64
GPT-3 (babbage)	3.20	24.56
GPT-3 (curie)	13.45	27.15
GPT-3 (davinci)	20.40	27.70
GPT-J	14.71	26.06
Cedille	24.89	20.59
Pagnol (small)	0.76	1.20
Pagnol (medium)	1.07	1.48
Pagnol (large)	1.06	3.47
GPT-fr (base)	1.47	1.57

Table 5: BLEU scores for translation on WMT14 for the English-French language pair (higher is better).

4 Toxicity analysis

In order to evaluate the toxicity of the model we closely followed the work conducted in [11]. We studied the case

of unprompted (i.e. conditioned only on a start-of-sentence token) and prompted generation.

The original work in [11] used the Perspective API, a service that uses machine learning classifiers to estimate the perceived toxicity of text. In this work, we employ the Detoxify tool [19] instead. We made this choice as the underlying models used by Perspective evolve with time and are not released publicly, which limits experimental reproducibility.

Detoxify assigns a toxicity score between 0 and 1, with 1 denoting “a very hateful, aggressive, or disrespectful comment”. We refer to content with a score > 0.5 as “toxic”. We use the “multilingual” Detoxify model from release v0.4.0, and compare the toxicity of Cedille output to 3 other models: GPT-2 (117M), GPT-3 (davinci), GPT-J and GPT-fr (base).

4.1 Unprompted toxicity

For the unprompted toxicity we analyze the expected maximum toxicity, i.e. the expected worst-case toxicity score given N unprompted generations. Figure 1 shows bootstrap estimates (1000 iterations) of the expected maximum toxicity for N generations with variance bounds as shades.

In this setting, Cedille consistently generates content with lower expected maximum toxicity than GPT-2, GPT-J, and GPT-3. After 100 generations, this value is under 0.5 for GPT-fr and Cedille (0.41 and 0.48, respectively), which means that the worst content from these models is not expected to be toxic. This is in contrast with the other models, for which maximum expected toxicity values are 0.64, 0.54 and 0.56.

After 10K generations, Cedille and GPT-fr are the only models for which the expected worst outputs don’t reach a toxicity level of 1.0. We expect all other models to have at least one output that is maximally toxic as detected by Detoxify. Generally the two models that perform best are GPT-fr and Cedille, which were both trained on carefully filtered datasets, pointing to the importance of dataset curation when considering the safety of language models.

Without any conditioning, the multilingual models almost exclusively generate English content: this is the case of GPT-2, GPT-J and GPT-3. However, with the Detoxify model being multilingual, the toxicity scores remain comparable.

4.2 Prompted toxicity

For prompted toxicity we used a set of 50 French prompts with values of toxicity spanning the full range, with a mean of 0.34. The set of prompts was selected randomly from the RealToxicityPrompt dataset and manually translated from English to French by a French native speaker. We used a smaller number of prompts than in [11] due to limited computing resources. The French prompts cause the multilingual models (GPT-2, GPT-J and GPT-3) to gener-

ate French content. For each prompt, each model generates 50 completions. We used nucleus sampling with $p = 0.9$ to generate up to 20 tokens per continuation, following the protocol from [11].

Table 6 shows two properties: 1) the expected maximum toxicity over 25 generations (with standard deviations in parentheses) and 2) the empirical probability of generating toxic text at least once among 25 generations.

Model	Exp. max tox.	Prob. toxicity
GPT-2 ^a	0.63 (0.23)	0.66
GPT-3 (davinci)	0.68 (0.27)	0.74
GPT-J	0.73 (0.26)	0.78
Cedille	0.66 (0.27)	0.72
GPT-fr (base)	0.73 (0.27)	0.78

Table 6: Toxicity of prompted generations.

^aUpon manual inspection, it appeared that GPT-2 is unable to generate sensible French content, and as such the resulting toxicity values can’t be compared to other models.

For both properties, Cedille outperforms the other models. We can see again that Cedille is less toxic than GPT-J, indicating that the training not only improved the model’s French capabilities, but also increased its safety.

5 Conclusions

In this work we introduced Cedille, a large auto-regressive French language model. Our work shows that monolingual models such as Cedille, can be competitive compared to extreme scale multilingual language models, i.e. GPT-3. Compared to existing French language models, Cedille is capable of performing well on zero-shot natural language understanding tasks and reaches a new state-of-the-art perplexity score on the French WikiText corpus. Lastly, our approach of toxicity filtering of the training data led to a decrease in both maximum toxicity as well as the likelihood of toxic output.

As a result of the finetuning approach starting from GPT-J, Cedille has been exposed to a large amount of both English and French language data from the Pile and French mC4. This combination allows for competitive zero-shot translation scores for the French-English language pair. Early experiments indicate that finetuning an existing English language model and adapting it to French is more efficient even with considerable compute and data investments (see appendix).

Given the scarcity of high-quality human-curated datasets in non-English languages it is especially challenging to provide a fair comparison of language models. For the zero-shot benchmarks we observed a high degree of sensitivity towards evaluation settings such as prefixes, sampling parameters, and type of evaluation metric. The scores

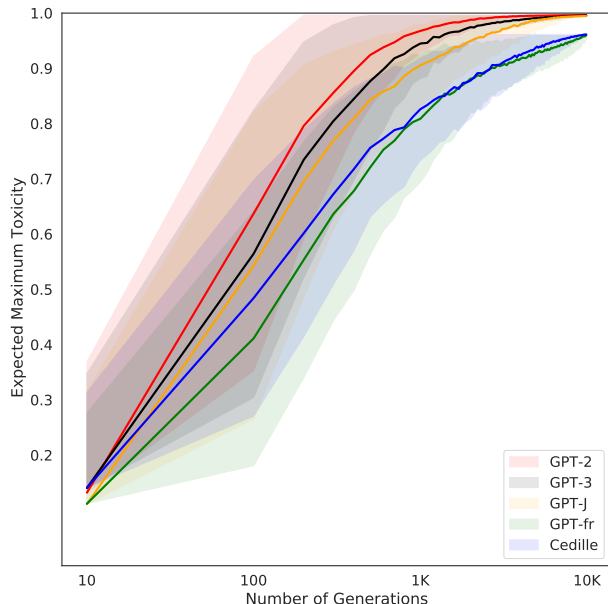


Figure 1: Unprompted expected maximum toxicity against increasing numbers of generations.

should therefore only be considered as a rough guidance and model performance may be highly task specific. In this work we haven’t provided performance metrics for other NLP tasks such as text classification or word sense disambiguation. Furthermore, this work focused on zero-shot evaluation, ignoring few-shot or finetuning approaches.

Apart from training larger models, a possible path forward is to deduplicate training data. This method has been shown to improve end-task performance significantly [8, 37] but was not conducted as part of this work. In order to further reduce language model toxicity, a possible direction is the integration of human feedback in the training process in order to reduce toxic output generation [38].

Data availability. Cedille is available under the MIT License on the Hugging Face model hub: <https://huggingface.co/Cedille/fr-boris>, and on our GitHub repository: <https://github.com/coterie/cedille-ai>. Regarding the French mC4 toxicity scores and toxicity analysis code, please refer to: <https://github.com/coterie/real-toxicity-prompts>.

Funding. This work was funded by, and conducted at, Coterie SA⁷. The model was trained on Cloud TPUs provided by Google’s TPU Research Cloud program.

Acknowledgments. We thank Sébastien Flury and François Bochatay for their guidance and feedback. Tiago Castanheiro, Flavien Bonvin and Livio Gamassia implemented the web-based Playground used to evaluate the

model. Tiago Castanheiro, Flavien Bonvin, Sacha Toufani, Livio Gamassia, and Kasper Andkjaer tested out multiple versions of the model. Sébastien Von Roth designed the Cedille logo as well as the visual design of the Playground and Cedille website⁸. Sonja Dossenbach assembled the dataset of recent French news. We are grateful to EleutherAI for publicly releasing the GPT-J model and offering us support on their Discord server⁹. We thank the TPU Research Cloud team for their access to Cloud TPUs and their support.

References

- [1] Alec Radford et al. “Language models are unsupervised multitask learners”. In: *OpenAI blog* 1.8 (2019), p. 9.
- [2] Tom B Brown et al. “Language models are few-shot learners”. In: *arXiv preprint arXiv:2005.14165* (2020).
- [3] Jared Kaplan et al. “Scaling laws for neural language models”. In: *arXiv preprint arXiv:2001.08361* (2020).
- [4] Chau Tran et al. “Facebook AI WMT21 news translation task submission”. In: *arXiv preprint arXiv:2108.03265* (2021).
- [5] Naveen Arivazhagan et al. “Massively multilingual neural machine translation in the wild: Findings and challenges”. In: *arXiv preprint arXiv:1907.05019* (2019).

⁷<https://coterie.com>

⁸<https://cedille.ai>

⁹<https://discord.gg/zBGx3azzUn>

- [6] Antoine Simoulin and Benoit Crabbé. “Un modèle Transformer Génératif Pré-entraîné pour le français”. In: *Traitement Automatique des Langues Naturelles*. ATALA. 2021, pp. 245–254.
- [7] Julien Launay et al. “PAGNol: An Extra-Large French Generative Model”. In: *arXiv preprint arXiv:2110.08554* (2021).
- [8] Guillaume Wenzek et al. “Ccnnet: Extracting high quality monolingual datasets from web crawl data”. In: *arXiv preprint arXiv:1911.00359* (2019).
- [9] Emily M Bender et al. “On the Dangers of Stochastic Parrots: Can Language Models Be Too Big?” In: *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*. 2021, pp. 610–623.
- [10] Isaac Caswell et al. “Quality at a glance: An audit of web-crawled multilingual datasets”. In: *arXiv preprint arXiv:2103.12028* (2021).
- [11] Samuel Gehman et al. “RealToxicityPrompts: Evaluating neural toxic degeneration in language models”. In: *arXiv preprint arXiv:2009.11462* (2020).
- [12] Johannes Welbl et al. “Challenges in detoxifying language models”. In: *arXiv preprint arXiv:2109.07445* (2021).
- [13] Sumanth Dathathri et al. “Plug and play language models: A simple approach to controlled text generation”. In: *arXiv preprint arXiv:1912.02164* (2019).
- [14] Brian Lester, Rami Al-Rfou, and Noah Constant. “The power of scale for parameter-efficient prompt tuning”. In: *arXiv preprint arXiv:2104.08691* (2021).
- [15] Nitish Shirish Keskar et al. “Ctrl: A conditional transformer language model for controllable generation”. In: *arXiv preprint arXiv:1909.05858* (2019).
- [16] Ben Wang and Aran Komatsuzaki. *GPT-J-6B: A 6 Billion Parameter Autoregressive Language Model*. <https://github.com/kingoflolz/mesh-transformer-jax>. May 2021.
- [17] Jianlin Su et al. “Roformer: Enhanced transformer with rotary position embedding”. In: *arXiv preprint arXiv:2104.09864* (2021).
- [18] Linting Xue et al. “mT5: A massively multilingual pre-trained text-to-text transformer”. In: *arXiv preprint arXiv:2010.11934* (2020).
- [19] Laura Hanu and Unitary team. *Detoxify*. <https://github.com/unitaryai/detoxify>. 2020.
- [20] Hang Le et al. “Flaubert: Unsupervised language model pre-training for french”. In: *arXiv preprint arXiv:1912.05372* (2019).
- [21] Robyn Speer. *ftfy*. Zenodo. Version 5.5. 2019. DOI: 10.5281/zenodo.2591652. URL: <https://doi.org/10.5281/zenodo.2591652>.
- [22] Ben Wang. *Mesh-Transformer-JAX: Model-Parallel Implementation of Transformer Language Model with JAX*. <https://github.com/kingoflolz/mesh-transformer-jax>. May 2021.
- [23] Leo Gao et al. *A framework for few-shot language model evaluation*. Version v0.0.1. Sept. 2021. DOI: 10.5281/zenodo.5371628. URL: <https://doi.org/10.5281/zenodo.5371628>.
- [24] Leo Gao. *On the Sizes of OpenAI API Models*. <https://blog.eleuther.ai/gpt3-model-sizes/>. May 2021.
- [25] Stephen Merity et al. “Pointer sentinel mixture models”. In: *arXiv preprint arXiv:1609.07843* (2016).
- [26] *Perplexity of fixed-length models*. <https://huggingface.co/docs/transformers/perplexity>. Accessed: 2022-02-04.
- [27] Moussa Kamal Eddine, Antoine J-P Tixier, and Michalis Vazirgiannis. “BARThez: a skilled pre-trained french sequence-to-sequence model”. In: *arXiv preprint arXiv:2010.12321* (2020).
- [28] Shashi Narayan, Shay B Cohen, and Mirella Lapata. “Don’t give me the details, just the summary! topic-aware convolutional neural networks for extreme summarization”. In: *arXiv preprint arXiv:1808.08745* (2018).
- [29] Chin-Yew Lin. “Rouge: A package for automatic evaluation of summaries”. In: *Text summarization branches out*. 2004, pp. 74–81.
- [30] Martin d’Hoffschmidt et al. “FQuAD: French question answering dataset”. In: *arXiv preprint arXiv:2002.06071* (2020).
- [31] Pranav Rajpurkar et al. “SQuAD: 100,000+ questions for machine comprehension of text”. In: *arXiv preprint arXiv:1606.05250* (2016).
- [32] Louis Martin et al. “CamemBERT: a tasty french language model”. In: *arXiv preprint arXiv:1911.03894* (2019).
- [33] Ondřej Bojar et al. “Findings of the 2014 workshop on statistical machine translation”. In: *Proceedings of the ninth workshop on statistical machine translation*. 2014, pp. 12–58.
- [34] Kishore Papineni et al. “Bleu: a method for automatic evaluation of machine translation”. In: *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*. 2002, pp. 311–318.
- [35] Matt Post. “A Call for Clarity in Reporting BLEU Scores”. In: *Proceedings of the Third Conference on Machine Translation: Research Papers*. Belgium, Brussels: Association for Computational Linguistics, Oct. 2018, pp. 186–191. URL: <https://www.aclweb.org/anthology/W18-6319>.
- [36] Xiaodong Liu et al. “Very deep transformers for neural machine translation”. In: *arXiv preprint arXiv:2008.07772* (2020).
- [37] Katherine Lee et al. “Deduplicating training data makes language models better”. In: *arXiv preprint arXiv:2107.06499* (2021).
- [38] Long Ouyang et al. *Training language models to follow instructions with human feedback*. <https://openai.com/blog/instruction-following/>. Jan. 2022.

SUPPLEMENTARY MATERIAL

1 Experiments training from scratch

Given the amount of compute and data available, training from scratch rather than finetuning was considered. We experimented training Cedille from scratch using both the GPT-2 tokenizer (Cedille-fs-GPT2, vocab size 50,400) and the GPT-fr tokenizer (Cedille-fs-GPTfr, vocab size 50,000) for 60k steps using a peak learning rate of $1.2e-4$ and learning rate $1.2e-5$, and 7281 warm-up steps. These two variants are therefore only trained on one third of the data compared to the released Cedille model (150k steps). In order to have a fair comparison we show the result of Cedille after the same amount of steps (Cedille-60k). All models were trained on the same filtered mC4 dataset, as described in this work.

As shown in Table S1, Cedille-60k outperforms the from-scratch variants on the WikiText-fr benchmark. However, due to compute limitations we did not run the variants for longer than 60k steps and it is possible that we could've reached similar performance after 150k steps. Furthermore, both variants perform similarly, even though they are using a different tokenizer. Due to the variants performing very similarly, we conclude that even though a dedicated French tokenizer is a lot more efficient at encoding French text compared to the GPT-2 tokenizer, its benefit with regard to end-task performance was minimal in our experiments.

Model	PPL (byte)	PPL (token)
GPT-J	1.746	5.797
Cedille-60k	1.673	4.112
Cedille-fs-GPT2	1.794	4.972
Cedille-fs-GPTfr	1.775	6.856

Table S1: Byte-level and token-level perplexities for the WikiText-fr benchmark. Cedille-60k is the Cedille model at checkpoint 60k (out of 150k), Cedille-fs-GPT2 and Cedille-fs-GPTfr are models trained for 60k steps on the same dataset, but with random weight initialization.