

The Role of Language Imbalance in Cross-lingual Generalisation: Insights from Cloned Language Experiments

Anonymous ACL submission

Abstract

Multilinguality is crucial for extending recent advancements in language modelling to diverse linguistic communities. To maintain high performance while representing multiple languages, multilingual models ideally align representations, allowing what is learned in one language to generalise to others. Prior research has emphasised the importance of parallel data and shared vocabulary elements as key factors for such alignment. In this study, we investigate an unintuitive novel driver of cross-lingual generalisation: language *imbalance*. In controlled experiments on perfectly equivalent cloned languages, we observe that the existence of a predominant language during training boosts the performance of less frequent languages and leads to stronger alignment of model representations across languages. Furthermore, we find that this trend is amplified with scale: with large enough models or long enough training, we observe that bilingual training data with a $\frac{1}{10}$ language split yields better performance on both languages than a balanced $\frac{5}{10}$ split. Building on these insights, we design training schemes that can improve performance in all cloned languages, even without altering the training data. As we extend our analysis to real languages, we find that infrequent languages still benefit from frequent ones, yet whether language imbalance causes cross-lingual generalisation there is not conclusive.

1 Introduction

In recent years, autoregressive language models (LMs) pretrained on massive text corpora have advanced the state of the art in NLP tasks across the board (Brown et al., 2020; Touvron et al., 2023a,b; Köpf et al., 2023). While most of the leading models are trained on English texts, multilingual capabilities are crucial to make these advances accessible to a broader user base with diverse linguistic backgrounds. Ideally, data in one language should improve these multilingual

models’ performance in others. Such multilingual models should thus display **cross-lingual generalisation**: by reusing circuits (Cammarata et al., 2020; Elhage et al., 2021) and aligning their internal representations across languages, they may generalise concepts learned in a language to another.¹

How can such cross-lingual generalisation be achieved? This has been a focus of much prior work. One previously identified driver of cross-lingual generalisation is **parallel training data**; empirical evidence shows that training the model on either parallel sentence pairs (Lample and Conneau, 2019) or on corpora which are comparable across languages (Dufter and Schütze, 2020) improves generalisation. Another driver of cross-lingual generalisation is the availability of **anchor points**, i.e., vocabulary elements that are shared between languages; these can be naturally occurring subwords (e.g., *computer* in English and *computador* in Portuguese may share the subword *comp*; Pires et al., 2019; Wu and Dredze, 2019), shared special tokens (e.g., mask or bos symbols; Dufter and Schütze, 2020), or even artificially inserted “code-switching” augmentations (Conneau et al., 2020b; Reid and Artetxe, 2022; K et al., 2020; Feng et al., 2022). Beyond these two drivers, **limited model capacity** has been found to improve generalisation by Dufter and Schütze (2020), but to constrain multilingual capabilities by Chang et al. (2023).

In this work, we identify a surprising new factor that can boost cross-lingual generalisation abilities: **language imbalance**. We first conduct

¹A circuit is typically defined as a subgraph of a neural network which performs a specific computation. E.g., a circuit could be responsible for computing “greater than” comparisons between numbers in English sentences (Hanna et al., 2024). If representations are aligned across languages (in terms of how they encode, e.g., numbers) and circuits are reused, a model should be able to apply what it learns in one language (e.g., “greater than” comparisons in English) to perform similar computations in another language (e.g., French).

076 experiments in a synthetic setting with perfectly
077 equivalent cloned languages; this allows us to in-
078 vestigate LMs’ generalisation abilities in isolation
079 from the effects of languages’ dissimilarities, giv-
080 ing us a rough upper bound on the generalisation
081 we should expect to see between real language
082 pairs. In this cloned language setting, we find that
083 having a dominant main language improves gen-
084 eralisation, significantly boosting the performance
085 of less frequent languages. Furthermore, we find
086 that this effect becomes stronger when we either
087 increase our model’s size or when we train it for
088 longer. Based on these insights, we design training
089 curricula that improve performance in all cloned
090 languages without any modifications to the training
091 data. In the second part of our paper, we investigate
092 to what extent our insights transfer to real language
093 pairs. While we find that lower resource languages
094 typically do benefit from higher resource ones,
095 the impact of language imbalance on cross-lingual
096 generalisation is much less clear in this more
097 realistic setting. Overall, our results suggest an
098 interesting attribute of model training dynamics:
099 in some settings, having a main language can lead
100 model components to be shared across languages.

101 2 Cross-lingual Generalisation

102 While natural languages differ widely in their typol-
103 ogical properties, any pair of languages will share
104 at least a few grammatical and syntactic patterns.
105 Further, as their semantics reflect the underlying
106 processes of our world, language pairs should also
107 have similarities in the types of messages their
108 users typically convey. Intuitively, this suggests
109 that what is learned about a language L_A should
110 be useful to model another language L_B , and vice
111 versa. The extent of such generalisation depends
112 not only on how similar the two languages are, but
113 also on the employed learning algorithm. We anal-
114 yse such generalisation here, with a focus on how
115 language imbalance influences multilingual LMs.

116 Intuitively, if a model generalises well across
117 languages, it should achieve better performance
118 in each language (in terms of, e.g., perplexity)
119 than a monolingual model trained on the same
120 data. Concretely, a model trained on a multi-
121 lingual dataset $\mathcal{D}_{\text{multi}} = \mathcal{D}_A \cup \mathcal{D}_B$ containing
122 languages L_A and L_B should perform better than
123 monolingual models trained only on \mathcal{D}_A or \mathcal{D}_B .
124 This becomes clear when using definitions from
125 information theory: $\mathcal{D}_{\text{multi}}$ contains at least as
126 much information about L_A as \mathcal{D}_A . However, such

127 a multilingual model could also perform worse.
128 This could happen, for instance, if the data from
129 different languages interfere with each other during
130 optimisation through conflicting gradient update
131 directions (Wang et al., 2020). It could also happen
132 if the model has limited capacity: the multilingual
133 model has to represent many languages, which
134 intuitively requires more capacity than a single
135 one, even if some parameters are shared across
136 them (Conneau et al., 2020a; Pfeiffer et al., 2022).

137 In an attempt to make models better across
138 many languages, many multilingual models these
139 days are trained on somewhat balanced data (Scao
140 et al., 2023; Faysse et al., 2024). In some of these
141 cases, low-resource languages are upsampled
142 to improve their performance under the model.
143 As mentioned above, however, while balancing
144 languages’ appearance in a model’s training set
145 should intuitively improve performance, this
146 is not necessarily true. In fact, (and perhaps
147 surprisingly) some recent large language models
148 trained in mostly English-focused settings perform
149 reasonably well in a large sample of languages
150 (Ahia et al., 2023; Blevins and Zettlemoyer, 2022;
151 Briakou et al., 2023). These models’ training
152 data is typically highly imbalanced, with only a
153 small fraction being composed of “non-English”
154 languages. It is thus unclear whether multilingual
155 models indeed benefit from training on datasets
156 with balanced languages (Ye et al., 2023).

157 In smaller training scales, the benefits of multi-
158 lingual training are better understood. In general, it
159 has been found that low-resource languages tend to
160 benefit from data in higher-resource languages, but
161 high-resource languages benefit much less from
162 each other (Conneau et al., 2020a; Chang et al.,
163 2023). It is, however, unclear what causes cross-
164 lingual generalisation in this case. Is the model
165 in fact able to generalise better in the imbalanced
166 setting? Or does the model generalise equally
167 well in the balanced case, but its capacity bottle-
168 necks performance in higher-resource languages,
169 stopping us from observing performance gains?

170 We investigate the role of language imbalance in
171 cross-lingual generalisation here. Notably, Wendler
172 et al. (2024) recently showed that LMs seem
173 to perform internal computations in an abstract
174 “concept space” which is closest to their main
175 language (English in this case); representations are
176 then mapped back into the input language only in
177 the models’ final layers. Alabi et al. (2024) observe
178 a similar trend when using language adapters.

3 Experimental Setup

In this section, we provide a brief overview over models, data, and metrics used; for more details, see App. A. Our code will be made available on GitHub. All of our experiments use GPT-2-style decoder-only transformers (Radford et al., 2019). We base our implementation on the Languni Kitchen codebase (Stanić et al., 2023), and unless otherwise noted, we use the gpt-small configuration with 85M non-embedding parameters, training on 1.2B tokens of English or French books. We use separate tokenisers for English and French. For some of our experiments, we treat their vocabularies as **disjoint** and do not merge them. If we merge subwords that occur in both vocabularies, we make this clear with the label **anchored**.

As our main evaluation metric, we report our models’ perplexity (PPL) on the test set. Further, we define three metrics that allow for easy comparison of monolingual and multilingual models. Let t_A and t_B be the number of tokens a multilingual model is trained on in languages L_A and L_B , respectively. We define monolingual token equivalence (MLTE) as the number of tokens that would be required by a monolingual model, trained only in either language L_A or L_B , to achieve the same perplexity as the multilingual model does in that language. To determine MLTE, we fit a simple scaling law to predict perplexity from the number of training tokens (e.g., t_A) using the results from our trained monolingual models (see App. B for details). Analogously, we define monolingual PPL equivalence (MLPE) as the perplexity a monolingual model would reach when trained on the same number of L_A tokens (i.e., t_A) as a given multilingual model. Finally, we define token efficiency (TEff) as the fraction between MLTE and the number of tokens used for multilingual training, e.g., $\text{TEff}_A = \frac{\text{MLTE}_A}{t_A}$. Intuitively, if $\text{TEff} > 1$, performance improves due to multilinguality, while if $\text{TEff} < 1$, multilinguality hurts performance.

4 Cloned Languages

In this section, we examine the model’s capability to generalise across perfectly equivalent **cloned languages**. We create a cloned language by duplicating the language model’s vocabulary; this allows us to encode each sequence in either the original language (using the original vocabulary) or in the cloned language (using the cloned vocabulary). This experimental paradigm was originally proposed by K et al. (2020) and Dufter and

Schütze (2020).² Formally, let L_{orig} be an “original” language with a vocabulary of subword units Σ ; we denote each subword $w \in \Sigma$. This language can be described by a probability distribution $p(\mathbf{w}_{\text{orig}})$, where $\mathbf{w}_{\text{orig}} \in \Sigma^*$. We clone language L_{orig} by creating multiple instantiations of it: L_1, L_2, \dots, L_N . These languages have vocabularies Σ_n , each of which has symbols that are equivalent to the original ones.³ Furthermore, these languages define probability distributions which are isometric to the original language. If we denote equivalence as $\mathbf{w}_n \stackrel{\circ}{=} \mathbf{w}_{\text{orig}}$ for $\mathbf{w}_n \in \Sigma_n^*$ and $\mathbf{w}_{\text{orig}} \in \Sigma^*$, we have $\mathbf{w}_n \stackrel{\circ}{=} \mathbf{w}_{\text{orig}} \implies p(\mathbf{w}_n) = p(\mathbf{w}_{\text{orig}})$.

Given dataset $\mathcal{D}_{\text{orig}} = \{\mathbf{w}_{\text{orig}}^{(k)}\}_{k=1}^K$ with $\mathbf{w}_{\text{orig}}^{(k)} \sim p(\mathbf{w}_{\text{orig}})$, we can now create a multilingual dataset $\mathcal{D}_{\text{multi}}$ by independently mapping each sequence to one of the cloned languages: For each $\mathbf{w}^{(k)}$, we first sample a language $L^{(k)} \sim p(L)$ from a categorical distribution over languages, then we map the sequence to $L^{(k)}$ by encoding it using the corresponding vocabulary. We can write $\mathcal{D}_{\text{multi}} = \bigcup_{n=1}^N \mathcal{D}_n$ where

$$\mathcal{D}_n = \left\{ \mathbf{w}_n^{(k)} \mid \mathbf{w}_n^{(k)} \stackrel{\circ}{=} \mathbf{w}_{\text{orig}}^{(k)} \text{ and } L^{(k)} = L_n \right\}$$

denotes the subset in language L_n .

Importantly, cloned languages are perfectly equivalent, having the same syntax, semantics, and distribution. They differ only in the symbols used to encode their vocabularies. Any generalisation we observe in this setting should thus serve as an upper bound on the potential to generalise across non-identical natural languages.⁴ In other words, if our model cannot generalise across cloned languages, we would have strong reason to believe it shouldn’t generalise across distinct languages. If we observe that a model can generalise across cloned languages, however, we may or may not observe the same to happen across non-cloned languages. We’ll investigate the latter in Section 5.

²K et al. (2020) perform duplication on the character IDs, i.e., before tokenisation, while Dufter and Schütze (2020) adopt an approach equivalent to ours. Both of these works term L_2 a “fake” language. Since there is no distinction between L_1 and L_2 , however, we call them cloned languages instead. Other related studies have investigated the effect of infinitely many cloned languages on LMs’ performance (Huang et al., 2023; Chen et al., 2023), or employed duplicated vocabularies at the token level to study their impact on LMs’ memorisation or performance (Kharitonov et al., 2021; Schäfer et al., 2024).

³Unless otherwise noted, these vocabularies are defined as disjoint sets in our experiments, meaning that no anchor points exist across languages.

⁴As for most of our experiments we consider cloned languages’ alphabets to be disjoint, in practice our results only upper bound the cross-lingual generalisation of models with no anchor points (i.e., with disjoint vocabularies).

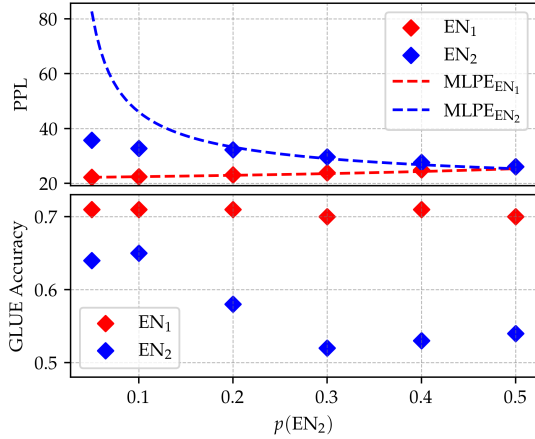


Figure 1: LM performance by imbalance ratio. (top) LM perplexity. (bottom) LM accuracy on GLUE; models were fine-tuned in EN_1 and evaluated on either EN_1 and EN_2 .

4.1 Generalisation

Due to the equivalence of cloned languages, one may expect language models to easily generalise across them. In that case, training a multilingual model on datasets \mathcal{D}_1 and \mathcal{D}_2 would lead to similar performance to training a monolingual model on the original dataset $\mathcal{D}_{\text{orig}}$ (note that $|\mathcal{D}_{\text{orig}}| = |\mathcal{D}_1| + |\mathcal{D}_2|$). We perform this experiment here, training either monolingual models on English (EN), or multilingual models on cloned English (EN_1 and EN_2), setting $p(\text{EN}_1) = p(\text{EN}_2) = 0.5$. Perhaps surprisingly, when training in this balanced multilingual setting, language modelling performance is significantly worse than in the monolingual setting (see Table 1, rows 2 & 4). In fact, one would obtain better performance training two monolingual models for half as many steps than training on this combined data. Training data in one language seems to hurt performance in the other language instead of boosting it. This indicates that the model is not able to generalise well across languages in this setting.

Takeaway 1. *The model does not generalise well across cloned languages given a 50/50 data split.*

4.2 Language Imbalance

How does the balance of the languages’ data affect generalisation performance? Will the multilingual model still underperform its monolingual equivalents when trained on an uneven language distribution? When varying the ratio of EN_1 to EN_2 data shown during training (while keeping the total number of training steps constant), we observe that the rarer “lower resource” language, here always EN_2 , benefits from the presence of a

dominant “main language”. Fig. 1 (left) shows that, under higher imbalance, the model’s performance on EN_2 becomes much better than that of a monolingual model trained on the same amount of EN_2 data. For example, when training in the 90/10 regime, we obtain a $\text{TEff}_{\text{EN}_2}$ of over 2 (see Table 1, row 5). Do these improvements translate to better cross-lingual generalisation on downstream tasks as well? We test this by fine-tuning models on the GLUE benchmark (Wang et al., 2019) in EN_1 only, and evaluating them on EN_1 and EN_2 . We observe that models trained under higher language imbalance indeed have significantly better EN_2 zero-shot performance (see Fig. 1 right). Together, these results suggest that cross-lingual generalisation is occurring.

Is this generalisation attained due to the model’s internal computations being shared across languages? To answer this question, we analyse how language imbalance affects the cross-lingual alignment of our models’ representations. Looking at the cosine similarity of equivalent subwords $w_1 \stackrel{\circ}{=} w_2$ in EN_1 and EN_2 , we find that similarity steadily increases with higher imbalance: in the 50/50 setting, embeddings are not aligned (exhibiting an average cosine similarity of 0.02), while, e.g., in the 90/10 setting, equivalent subwords are much more aligned, showing a similarity of 0.28 (details in App. C). Comparing the cosine similarity of hidden states when the LM is given equivalent sequences $w_1 \stackrel{\circ}{=} w_2$, we also observe stronger alignment for a model trained in the imbalanced 90/10 regime, compared to the 50/50 counterpart (see App. F). Interestingly, the cosine similarity between gradients is also higher in the imbalanced setting: when processing equivalent sequences, the gradients with respect to w_1 or w_2 have an average cosine similarity of 0.53 for the model trained in the 90/10 setting, compared to 0.07 in the 50/50 setting (see full plots of similarities per model component in App. G). This suggests that the gradient updates with respect to one language may benefit the optimisation process of that language’s cloned counterpart more when training under higher imbalance.

Takeaway 2. *Language imbalance improves generalisation and leads to representations which are more aligned across cloned languages.*

4.3 Many Languages

How does this trend transfer to settings with more than two languages? In such cases, sharing circuits across languages might be even more crucial due to the model’s limited capacity. Instead of cloning

Run Type	Row	Training Data				PPL		TEff	
		# Tokens	$p(\text{EN}_1)$	$p(\text{EN}_2)$	$p(\text{EN}_3), \dots, p(\text{EN}_{10})$	EN ₁	EN ₂	EN ₁	EN ₂
Monolingual	1	1.2B	100%	0%	0%	21.9	-	1	-
	2	$0.5 \times 1.2\text{B}$	100%	0%	0%	25.3	-	1	-
	3	$0.1 \times 1.2\text{B}$	100%	0%	0%	48.4	-	1	-
2 languages	4	1.2B	50%	50%	0%	26.1	26.1	0.89	0.89
	5	1.2B	90%	10%	0%	22.5	32.8	1.00	2.08
10 languages	6	1.2B	10%	10%	$10\%, \dots, 10\%$	35.5	35.7	1.69	1.67
	7	1.2B	50%	$\frac{1}{18}$	$\frac{1}{18}, \dots, \frac{1}{18}$	24.6	33.4	1.15	3.56
Schedule	8	1.2B	100% \downarrow 0%	0% \uparrow 100%	0%	>1B	31.4	-	0.47
	9	1.2B	90% \downarrow 10%	10% \uparrow 90%	0%	26.5	24.4	0.83	1.18
2x data	10	$2 \times 1.2\text{B}$	50%	50%	0%	23.3	23.3	0.73	0.73
	11	$2 \times 1.2\text{B}$	90% \downarrow 10%	10% \uparrow 90%	0%	22.8	20.4	0.83	1.60
3x data	12	$3 \times 1.2\text{B}$	50%	50%	0%	22.2	22.2	0.64	0.64
	13	$3 \times 1.2\text{B}$	90% \downarrow 10%	10% \uparrow 90%	0%	21.5	19.3	0.77	1.63

Table 1: Performance of language models trained on different compositions of EN₁ and EN₂. $a\% \downarrow b\%$ indicates an immediate decrease from $a\%$ down to $b\%$ halfway during training. Analogously, $a\% \uparrow b\%$ indicates an immediate increase.

the language only once, we now clone it nine times, obtaining in total 10 languages. In Table 1 (rows 6 & 7), we report the performance when sampling the languages in a balanced way and when having a much stronger main language.

Interestingly, when sampling uniformly, we obtain $\text{TEff} \approx 1.7$; performance is thus better than with a monolingual model trained on an equivalent amount of monolingual data (compare rows 6 & 3). This differs from our observations for the bilingual setting, where uniform language sampling performed worse than the equivalent monolingual models. Presumably, modelling this many languages effectively with limited model capacity may lead the model to share its circuits, improving cross-lingual generalisation (Dufter and Schütze, 2020). The limit of infinite languages (in which a model never observes the same language more than once) was analysed by Huang et al. (2023); interestingly, LMs still seem to learn the language, to some extent, even in that setting.

In the imbalanced setting where we sample a stronger “main language” 50% of the time, we observe even stronger performance on all languages. Despite the model seeing only roughly 67M tokens in each of the rarer languages (1/18 of all steps), it achieves **better** performance in these languages than in the uniform setting with 120M tokens (1/10 of all steps) per language. In fact, on the rarer languages, the model achieves $\text{TEff} \approx 3.6$, matching the performance of a monolingual model trained on 240M tokens.

Takeaway 3. *When training on many cloned languages, sampling a main language disproportionately improves generalisation.*

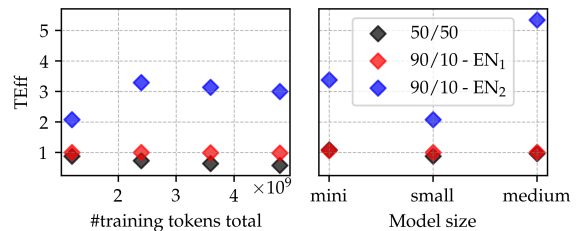


Figure 2: TEff as we train LMs with (left) more data, or (right) larger architectures. mini, small and medium denote GPT sizes in Languini (Stanić et al., 2023), with 11M, 85M, and 303M non-embedding parameters.

4.4 Effect of scaling

Model and data size are crucial factors for the performance of LMs. Here, we investigate how the previously identified trends are affected by scaling the model architecture or training data. Fig. 2 (left) shows that the effect of imbalance on cross-lingual generalisation appears to increase when we train on twice as much data (2.4B tokens instead of 1.2B), reaching $\text{TEff} > 3$; this corresponds to a “chinchilla-optimal” setup for our GPT small model (Hoffmann et al., 2022). At the same time, the TEff of the 50/50 setting seems to be decreasing under prolonged training. This might be caused by the heightened importance of model capacity under longer training, which may have a stronger impact on performance when representations are less aligned across languages. Overall, the disparity in effectiveness between the imbalanced and balanced settings grows with longer training. Remarkably, when training for 4.8B tokens, the 90/10 setting yields better performance in both languages, compared to the 50/50 setting.

When decreasing the model size, we also

observe higher performance benefits in the imbalanced setting (see Fig. 2 right), potentially due to the capacity argument described above. Interestingly, however, the effect of imbalance appears to be significantly stronger for larger models as well. When training a larger model with around 300M parameters (GPT medium in Languini; Stanić et al., 2023), in the 90% setting, we achieve better performance on both languages than under the 50% split. This might be because larger models generally exhibit better generalisation ability than smaller ones (Brown et al., 2020).

Takeaway 4. *Longer training and larger models lead to stronger performance benefits due to language imbalance.*

4.5 Language Sampling Schedule

Knowing that language imbalance boosts generalisation, how can we use this insight to train better models? Is there a way to leverage our insights in order to improve performance on two languages, even with the same training data? In Table 1 (rows 8, 9, 11, and 13), we report results when training with a language sampling schedule that ensures a language imbalance throughout all of training, but which still leads to an overall 50% split between EN₁ and EN₂ data seen by the model. We sample EN₁ with a higher probability during the first half of training. Then, we sample EN₂ more often to achieve a marginal split of 50%.

When showing exclusively EN₁ at first, and then showing only EN₂ (100% 0/0 or 100%; row 8), we observe bad overall performance. By the end of training, perplexity on EN₁ is very high, presumably due to catastrophic forgetting (McCloskey and Cohen, 1989; French, 1999). Further, EN₂ does not seem to benefit from the EN₁ data, achieving very low performance, which might be due to the lower learning rate in the second half of training.⁵

On the other hand, if we avoid catastrophic forgetting, making sure that the model encounters at least some samples of both languages at every point in training, via a 90% 10% or 90 split (first sampling languages with ratio 90%, and then switching to 10% after half of training), we can mitigate these issues. On our standard training set (1.2B tokens, row 9), we observe almost equivalent performance to uni-

⁵Chen et al. (2023) find that an equivalent setting can still be beneficial when using many more languages: they periodically reinitialise the learned embeddings (which is equivalent to switching to a new cloned language) and obtain models that are better adaptable to new languages.

form language sampling on EN₁, but significantly improved performance on EN₂. Under longer training, these benefits become more pronounced: this language schedule improves performance on both languages compared to the simple 50% setting (compare row 10 vs 11 and row 12 vs 13).

Takeaway 5. *Compared to uniform language sampling, an imbalanced ratio throughout training can lead to better results on all languages, even if the overall language split remains balanced.*

5 Real Languages

To verify whether the insights from our cloned-language experiments hold in a more natural setting, we now run experiments with multilingual models on English (EN) and French (FR).

5.1 Generalisation

In the cloned setting, we observed no significant generalisation when training on a balanced language mix (i.e., TEff < 1, representations were unaligned, and zero-shot GLUE accuracy on EN₂ was bad). Similarly, when sampling EN and FR data uniformly, we also obtain TEff < 1. A multilingual model’s performance is thus worse than a monolingual model trained only in the same EN or FR data (see Table 2, row 7). Notably, prior work has identified anchors (shared vocabulary items across languages) help generalisation (Dufter and Schütze, 2020; Pires et al., 2019; Wu and Dredze, 2019). We thus experiment with similarly merging vocabulary items shared between EN and FR, and confirm this helps performance (compare Table 2, row 7 vs 11). We run more experiments analysing the impact of anchor points in both cloned and real languages, see App. D. Note that, with an anchored vocabulary, generalisation across EN and FR is not necessarily upper bounded by our results on disjoint cloned languages. In fact, in the 50% setting, we observe a marginally higher TEff for EN–FR models with an anchored vocabulary than for EN₁–EN₂ models where we used disjoint vocabularies (compare Table 1 row 4 and Table 2 row 11).

5.2 Language Imbalance

Analogous to the cloned setting, we observe that an imbalanced EN/FR ratio leads to improved performance (TEff > 1), on the rarer language (see Table 2, rows 7-9 & 11-13). This is the case for both, a 90% and a 10% EN/FR ratio. Fig. 3 shows PPL and TEff in EN and FR as a function of the language imbalance. We observe that large

Run Type	Row	Training Data			PPL		TEff	
		# Tokens	$p(\text{EN})$	$p(\text{FR})$	EN	FR	EN	FR
Monolingual	1	1.2B	100%	0%	21.9	-	1	-
	2	$0.5 \times 1.2\text{B}$	100%	0%	25.3	-	1	-
	3	$0.1 \times 1.2\text{B}$	100%	0%	48.4	-	1	-
	4	1.2B	0%	100%	-	16.0	-	1
	5	$0.5 \times 1.2\text{B}$	0%	100%	-	18.4	-	1
	6	$0.1 \times 1.2\text{B}$	0%	100%	-	34.1	-	1
Multilingual disjoint vocabs	7	1.2B	50%	50%	26.4	19.4	0.85	0.82
	8	1.2B	90%	10%	22.5	31.9	1.00	1.05
	9	1.2B	10%	90%	43.5	16.4	1.10	0.97
	10	1.2B	90% \downarrow 10%	10% \uparrow 90%	29.1	20.5	0.60	0.66
Multilingual anchored vocabs	11	1.2B	50%	50%	26.0	19.0	0.91	0.88
	12	1.2B	90%	10%	22.5	29.0	1.00	1.27
	13	1.2B	10%	90%	39.5	16.5	1.33	0.96
	14	1.2B	90% \downarrow 10%	10% \uparrow 90%	28.9	19.3	0.61	0.83
	15	1.2B	90% \downarrow 10% \uparrow 50% \rightarrow 50%	10% \uparrow 90% \downarrow 50% \rightarrow 50%	26.4	18.5	0.85	1.00
	16	1.2B	95% \downarrow 35% \rightarrow 35% \rightarrow 35%	5% \uparrow 65% \rightarrow 65% \rightarrow 65%	26.3	18.7	0.86	0.95
2x data	17	$2 \times 1.2\text{B}$	50%	50%	23.0	16.9	0.79	0.76
	18	$2 \times 1.2\text{B}$	90% \downarrow 10%	10% \uparrow 90%	26.1	17.1	0.44	0.70
3x data	19	$3 \times 1.2\text{B}$	50%	50%	21.8	16.0	0.70	0.67
	20	$3 \times 1.2\text{B}$	90% \downarrow 10%	10% \uparrow 90%	25.1	16.2	0.35	0.63

Table 2: Performance of language models trained on different compositions of EN and FR. $a\% \rightarrow b\% \rightarrow c\% \rightarrow d\%$ indicates a four stage language schedule, switching immediately between, e.g., $c\%$ and $d\%$ after 75% of training.

imbalances generally seem to yield $\text{TEff} > 1$; the worst TEff is reached with a balanced EN/FR ratio. These trends are in line with our findings in the cloned setting. However, especially with disjoint vocabularies, the observed performance benefits due to generalisation are less significant. Presumably, this is due to EN and FR not being equivalent and thus generally allowing less generalisation.

Does imbalance again improve generalisation due to a better alignment of the model’s representations in the two languages? As in the cloned language setting, we investigate the cosine similarity between the models’ hidden states when processing “equivalent” sequences in the two languages. For real languages, however, we do not have access to perfectly equivalent sequences. Instead, we mimic this scenario using parallel translated sequences in the two languages, which should contain roughly similar properties. Differently from the cloned language setting, we do not observe higher hidden state similarities for models trained on imbalanced data (see App. F). Further, we find that gradient similarities barely differ across balanced and imbalanced settings when using disjoint vocabularies. For the anchored vocabulary they are even marginally higher in the balanced setting (see App. G). We thus do not find evidence that the improved TEff in the imbalanced setting is caused by a stronger alignment of model updates across languages in this setting. A possible reason for this discrepancy could be that, at the scales of our experiments, LMs tend to rely on

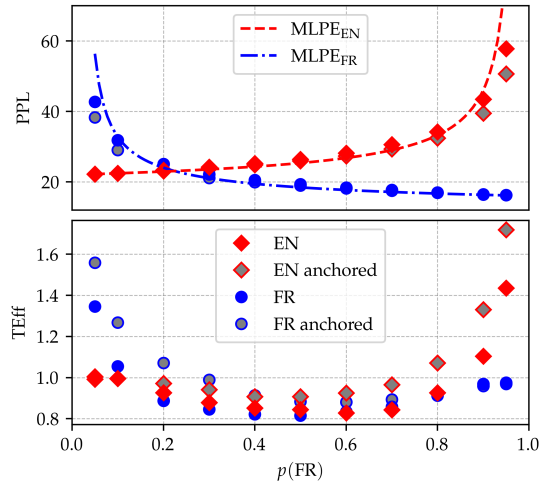


Figure 3: LM performance on EN and FR by imbalance ratio.

language specific surface-level features (which are shared by cloned languages, but not by distinct real languages) and show less understanding of complex semantics which might be more generalisable. Future research might thus consider investigating these trends at larger scales.

Takeaway 6. *Imbalanced multilinguality boosts the performance of real low-resource languages. However, this effect is weaker here than for cloned languages. Further, for real languages, we do not find evidence of language imbalance leading to representations which are more cross-lingually aligned.*

5.3 Effect of Scaling

In the cloned setting, we observed that prolonging training significantly decreased TEff in the 50/50 setting. We hypothesised that this might be caused by a stronger influence of the limited model capacity with longer training, and poor sharing of representations between languages. As EN and FR are distinct languages that require at least some language specific representations, we might expect this trend to be even more pronounced for these languages. However, compared to the cloned setting, prolonging training leads to a smaller decline in TEff in the 50/50 setting here. Presumably, the anchored vocabulary allows for better generalisation compared to the cloned setting, despite the languages being distinct.

Further, unlike in the cloned setting, longer training significantly decreases the TEff of the lower-resource language in the imbalanced setting here (see Fig. 4). In fact, the 90/10 TEff even falls below 1, approaching the TEff of the 50/50 setting. This suggests that language imbalance might not improve generalisation across distinct real languages. Still, when scaling up the model, we observe an increase of almost 2x in the TEff of the lower-resource language (see Fig. 4). This is in line with our cloned languages observations, although the effect is weaker.

Takeaway 7. *Performance benefits for real low-resource languages tend to decrease or vanish with longer training. Larger models, however, appear to yield higher performance benefits in both the balanced and imbalanced setting.*

5.4 Language Sampling Schedule

For equivalent cloned languages, we found that an imbalanced language sampling schedule can lead to improvements upon simple uniform sampling. If this held for real languages as well, it could have important practical implications for future multilingual LM training. However, whereas a 90/10/10/90 schedule yielded strong performance on cloned languages, matching or outperforming the 50/50 setting, this is not the case for EN and FR (see Table 2, row 10 vs 7 and row 14 vs 11). Furthermore, in line with the observations above, longer training does not make this schedule more effective, but instead increases its gap to the performance of the 50/50 setting (see rows 17-20).

The discrepancy between these results and the ones in cloned languages might be explained by

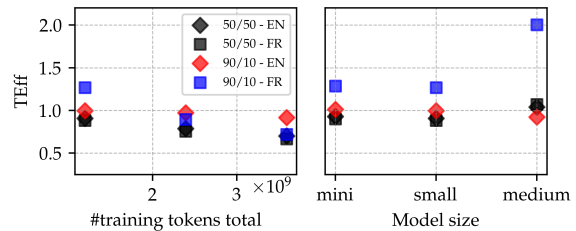


Figure 4: TEff of models on EN and FR with anchored vocab as we train them with (left) more data, or (right) larger architectures.

the reduced effect of imbalance on the generalisation and representation alignment in real languages. The schedules may be enough to force LMs to share circuits across cloned languages, but not across real ones. To investigate if this negative result was a particular property of our chosen schedule, we explore other more complex scheduling options.⁶ In general, none of the tested schedules appears to outperform the 50/50 setting (see rows 15, 16) on both languages. However, more complex 4-stage schedules, can obtain better performance on one language while incurring a slight performance drop in the other. Intriguingly, this allows trading off the performance of different languages without altering the training data.

Takeaway 8. *For real languages, we do not find improvements on all languages due to the tested language schedules. However, they allow for trading off performance in different languages.*

6 Conclusion

We ran experiments to measure cross-lingual generalisation in both a controlled setting with cloned English languages, as well as with English and French. In both settings, we find that, without vocabulary overlap, our models do not show strong cross-lingual generalisation when trained on a balanced language set. However, when training on an imbalanced mix of languages, we observe increased performance compared to monolingual settings. For cloned languages, we find that this can be explained by a higher alignment of the model’s representations across languages, which indicates circuit reuse and improved cross-lingual generalisation. Yet, at the scales of our experiments, such a correlation is less evident in real languages. While our findings allow us to design an imbalanced language schedule that yields improved performance in the cloned setting, further research is required to extend these improvements to real-world settings.

⁶Future research might design these more carefully, also analysing the interplay of language- and learning rate schedule

639 Limitations

640 There are several limitations of our work, many of
641 which present opportunities for future research.

642 **Data and model size.** While we conduct exper-
643 iments with varying data (up to 4.8B tokens) and
644 model size (up to 336M parameters), it is uncertain
645 whether the identified trends also apply at the scale
646 of modern large language models. Additionally, for
647 more capable models, cross-lingual generalisation
648 might be relevant in different aspects, with, e.g.,
649 semantics playing a larger role. As the semantic
650 content communicated in different languages might
651 be easily transferable, this might impact generali-
652 sation dynamics.

653 **Languages.** We only run experiments on English
654 and French, two Indo-European languages. Further
655 work could consider more languages and investi-
656 gate the impact of language similarity in results
657 more broadly.

658 **Model architecture.** We run most of our exper-
659 iments on a Transformer decoder (we also mea-
660 sure embedding alignment for simpler Word2Vec
661 models). Future research could analyse the effects
662 of architecture in more depth to better understand
663 the drivers of representation alignment. [Conneau
664 et al. \(2020b\)](#), e.g., find that shared parameters in
665 the top layers lead to better cross-lingual transfer.
666 In our Word2Vec experiments, we do not observe
667 improvements in representation alignment due to
668 language imbalance (see Fig. 6), presumably due
669 to no parameters being shared between the two lan-
670 guages. Would this change when adding a shared
671 layer to the Word2Vec model?

672 **Downstream performance.** In our evaluation we
673 mainly rely on perplexity as a metric, with a sin-
674 gle experiment on GLUE accuracy. It might be
675 insightful to analyze effects on downstream task
676 performance more broadly.

677 **Quantifying generalisation.** In this work, we
678 mainly measure cross-lingual generalisation by
679 comparing the performance of multilingual mod-
680 els with that of monolingual models trained on
681 the same amount of data in the given language.
682 If a multilingual model on languages L_A and L_B
683 requires fewer L_A tokens to reach a given perplex-
684 ity on L_A than a monolingual model, we speak
685 of cross-lingual generalisation, knowing that per-
686 formance on L_A must have been boosted by data

in language L_B . Future work could formalise this
measure and aim to model/quantify the relationship
between the number of training tokens seen in a lan-
guage L_B and performance in another language L_A ,
depending on model size, language imbalance, lan-
guage similarity, anchor points, and other factors.
An accurate model of these relationships could be
of substantial practical value.

References

- Orevaoghene Ahia, Sachin Kumar, Hila Gonen, Jungo Kasai, David Mortensen, Noah Smith, and Yulia Tsvetkov. 2023. [Do all languages cost the same? Tokenization in the era of commercial language models](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 9904–9923, Singapore. Association for Computational Linguistics.
- Jesujoba O. Alabi, Marius Mosbach, Matan Eyal, Dietrich Klakow, and Mor Geva. 2024. [The hidden space of transformer language adapters](#). *arXiv preprint arXiv:2402.13137*.
- Terra Blevins and Luke Zettlemoyer. 2022. [Language contamination helps explain the cross-lingual capabilities of english pretrained models](#). In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 3563–3574.
- Eleftheria Briakou, Colin Cherry, and George Foster. 2023. [Searching for needles in a haystack: On the role of incidental bilingualism in PaLM’s translation capability](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 9432–9452.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. [Language models are few-shot learners](#). In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc.
- Nick Cammarata, Shan Carter, Gabriel Goh, Chris Olah, Michael Petrov, Ludwig Schubert, Chelsea Voss, Ben Egan, and Swee Kiat Lim. 2020. [Thread: Circuits](#). *Distill*.
- Tyler A. Chang, Catherine Arnett, Zhuowen Tu, and Benjamin K Bergen. 2023. [When is multilinguality a curse? language modeling for 250 high-and low-resource languages](#). *arXiv preprint arXiv:2311.09205*.

742	Yihong Chen, Kelly Marchisio, Roberta Raileanu, David Adelani, Pontus Lars Erik Saito Stenetorp, Sebastian Riedel, and Mikel Artetxe. 2023. Improving language plasticity via pretraining with active forgetting . <i>Advances in Neural Information Processing Systems</i> , 36:31543–31557.	799
743		800
744		801
745		
746		
747		
748	Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2020a. Unsupervised cross-lingual representation learning at scale . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 8440–8451. Association for Computational Linguistics.	802
749		803
750		804
751		805
752		
753		
754		
755		
756	Alexis Conneau, Shijie Wu, Haoran Li, Luke Zettlemoyer, and Veselin Stoyanov. 2020b. Emerging cross-lingual structure in pretrained language models . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 6022–6034. Association for Computational Linguistics.	806
757		807
758		808
759		809
760		810
761		
762		
763	Philipp Dufter and Hinrich Schütze. 2020. Identifying elements essential for BERT’s multilinguality . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 4423–4437. Association for Computational Linguistics.	811
764		812
765		813
766		814
767		815
768		816
769		817
770		818
771		819
772		820
773		821
774		
775		
776		
777		
778		
779	Nelson Elhage, Neel Nanda, Catherine Olsson, Tom Henighan, Nicholas Joseph, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, Nova DasSarma, Dawn Drain, Deep Ganguli, Zac Hatfield-Dodds, Danny Hernandez, Andy Jones, Jackson Kernion, Liane Lovitt, Kamal Ndousse, Dario Amodei, Tom Brown, Jack Clark, Jared Kaplan, Sam McCandlish, and Chris Olah. 2021. A mathematical framework for transformer circuits . <i>Transformer Circuits Thread</i> .	822
780		823
781		824
782		825
783		
784		
785		
786		
787		
788		
789		
790		
791	Manuel Faysse, Patrick Fernandes, Nuno Guerreiro, António Loison, Duarte Alves, Caio Corro, Nicolas Boizard, João Alves, Ricardo Rei, Pedro Martins, et al. 2024. CroissantLLM: A truly bilingual french-english language model . <i>arXiv preprint arXiv:2402.00786</i> .	826
792		827
793		828
794		829
795		
796		
797		
798		
799		
800		
801		
802		
803		
804		
805		
806		
807		
808		
809		
810		
811		
812		
813		
814		
815		
816		
817		
818		
819		
820		
821		
822		
823		
824		
825		
826		
827		
828		
829		
830		
831		
832		
833		
834		
835		
836		
837		
838		
839		
840		
841		
842		
843		
844		
845		
846		
847		
848		
849		
850		
851		
852		
853		
854		
855		

856	Taku Kudo and John Richardson. 2018. SentencePiece: A simple and language independent subword tokenizer and detokenizer for neural text processing . In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing: System Demonstrations</i> , pages 66–71.	913
857		914
858		915
859		916
860		917
861		918
862	Guillaume Lample and Alexis Conneau. 2019. Cross-lingual language model pretraining . <i>arXiv preprint arXiv:1901.07291</i> .	919
863		920
864		921
865	Michael McCloskey and Neal J. Cohen. 1989. Catastrophic interference in connectionist networks: The sequential learning problem . In <i>Psychology of Learning and Motivation</i> , volume 24, pages 109–165. Academic Press.	922
866		923
867		924
868		925
869		926
870	Tomás Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013. Efficient estimation of word representations in vector space . In <i>1st International Conference on Learning Representations, Workshop Track Proceedings</i> , Scottsdale, Arizona, USA.	927
871		928
872		929
873		930
874		931
875	Jonas Pfeiffer, Naman Goyal, Xi Lin, Xian Li, James Cross, Sebastian Riedel, and Mikel Artetxe. 2022. Lifting the curse of multilinguality by pre-training modular transformers . In <i>Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 3479–3495, Seattle, United States. Association for Computational Linguistics.	932
876		933
877		934
878		935
879		936
880		937
881		938
882		939
883		940
884	Telmo Pires, Eva Schlinger, and Dan Garrette. 2019. How multilingual is multilingual BERT? In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 4996–5001.	941
885		942
886		943
887		944
888	PleIAs. 2024. French-PD-Books dataset. https://huggingface.co/datasets/PleIAs/French-PD-Books . Accessed in 01/2024, Hugging Face Datasets library.	945
889		946
890		947
891		948
892	Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners . <i>OpenAI Blog</i> .	949
893		950
894		951
895		952
896	Machel Reid and Mikel Artetxe. 2022. PARADISE: Exploiting parallel data for multilingual sequence-to-sequence pretraining . In <i>Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 800–810, Seattle, United States. Association for Computational Linguistics.	953
897		954
898		955
899		956
900		957
901		958
902		959
903		960
904	Teven Le Scao, Angela Fan, Christopher Akiki, Elie Pavlick, Suzana Ilić, Daniel Hesslow, Roman Castagné, Alexandra Sasha Luccioni, François Yvon, Matthias Gallé, Jonathan Tow, Alexander M. Rush, Stella Biderman, Albert Webson, Pawan Sasanka Amanamanchi, Thomas Wang, Benoît Sagot, Niklas Muennighoff, Albert Villanova del Moral, Olatunji Ruwase, Rachel Bawden, Stas Bekman, Angelina McMillan-Major, Iz Beltagy, Huu Nguyen, Lucile	961
905		962
906		963
907		964
908		965
909		966
910		967
911		968
912		969
	Saulnier, Samson Tan, Pedro Ortiz Suarez, Victor Sanh, Hugo Laurençon, Yacine Jernite, Julien Launay, Margaret Mitchell, Colin Raffel, Aaron Gokaslan, Adi Simhi, Aitor Soroa, Alham Fikri Aji, Amit Alfassy, Anna Rogers, Ariel Kreisberg Nitzav, Canwen Xu, Chenghao Mou, Chris Emezue, Christopher Klamm, Colin Leong, Daniel van Strien, David Ifeoluwa Adelani, Dragomir Radev, Eduardo González Ponferrada, Efrat Levkovizh, Ethan Kim, Eyal Bar Natan, Francesco De Toni, Gérard Dupont, Germán Kruszewski, Giada Pistilli, Hady Elsahar, Hamza Benyamina, Hieu Tran, Ian Yu, Idris Abdulmumin, Isaac Johnson, Itziar Gonzalez-Dios, Javier de la Rosa, Jenny Chim, Jesse Dodge, Jian Zhu, Jonathan Chang, Jörg Frohberg, Joseph Tobing, Joydeep Bhattacharjee, Khalid Almubarak, Kimbo Chen, Kyle Lo, Leandro Von Werra, Leon Weber, Long Phan, Loubna Ben allal, Ludovic Tanguy, Manan Dey, Manuel Romero Muñoz, Maraim Masoud, María Grandury, Mario Šaško, Max Huang, Maximin Coavoux, Mayank Singh, Mike Tian-Jian Jiang, Minh Chien Vu, Mohammad A. Jauhar, Mustafa Ghaleb, Nishant Subramani, Nora Kassner, Nuru-laqilla Khamis, Olivier Nguyen, Omar Espejel, Ona de Gibert, Paulo Villegas, Peter Henderson, Pierre Colombo, Priscilla Amuok, Quentin Lhoest, Rheza Harliman, Rishi Bommasani, Roberto Luis López, Rui Ribeiro, Salomey Osei, Sampo Pyysalo, Sebastian Nagel, Shamik Bose, Shamsuddeen Hassan Muhammad, Shanya Sharma, Shayne Longpre, So-maieh Nikpoor, Stanislav Silberberg, Suhas Pai, Sydney Zink, Tiago Timponi Torrent, Timo Schick, Tristan Thrush, Valentin Danchev, Vassilina Nikoulina, Veronika Laippala, Violette Lepercq, Vrinda Prabhu, Zaid Alyafeai, Zeerak Talat, Arun Raja, Benjamin Heinzerling, Chenglei Si, Davut Emre Taşar, Elizabeth Salesky, Sabrina J. Mielke, Wilson Y. Lee, Abheesht Sharma, Andrea Santilli, Antoine Chaffin, Arnaud Stiegler, Debajyoti Datta, Eliza Szczechla, Gunjan Chhablani, Han Wang, Harshit Pandey, Hendrik Strobelt, Jason Alan Fries, Jos Rozen, Leo Gao, Lintang Sutawika, M Saiful Bari, Maged S. Al-shaibani, Matteo Manica, Nihal Nayak, Ryan Teehan, Samuel Albanie, Sheng Shen, Srulik Ben-David, Stephen H. Bach, Taewoon Kim, Tali Bers, Thibault Fevry, Trishala Neeraj, Urmish Thakker, Vikas Raunak, Xiangru Tang, Zheng-Xin Yong, Zhiqing Sun, Shaked Brody, Yallow Uri, Hadar Tojarieh, Adam Roberts, Hyung Won Chung, Jaesung Tae, Jason Phang, Ofir Press, Conglong Li, Deepak Narayanan, Hatim Bourfoune, Jared Casper, Jeff Rasley, Max Ryabinin, Mayank Mishra, Minjia Zhang, Mohammad Shoeybi, Myriam Peyrounette, Nicolas Patry, Nouamane Tazi, Omar Sanseviero, Patrick von Platen, Pierre Cornette, Pierre François Lavallée, Rémi Lacroix, Samyam Rajbhandari, Sanchit Gandhi, Shaden Smith, Stéphane Requena, Suraj Patil, Tim Dettmers, Ahmed Baruwā, Amanpreet Singh, Anastasia Cheveleva, Anne-Laure Ligozat, Arjun Subramonian, Aurélie Névól, Charles Lovering, Dan Garrette, Deepak Tunuguntla, Ehud Reiter, Ekaterina Taktasheva, Ekaterina Voloshina, Eli Bogdanov, Genta Indra Winata, Hailey Schoelkopf, Jan	970
		971
		972
		973
		974
		975

976	Christoph Kalo, Jekaterina Novikova, Jessica Zosa	cate subwords in language modelling. <i>arXiv preprint</i>	1039
977	Forde, Jordan Clive, Jungo Kasai, Ken Kawamura,	<i>arXiv:2404.06508</i> .	1040
978	Liam Hazan, Marine Carpuat, Miruna Clinciu, Na-		
979	joung Kim, Newton Cheng, Oleg Serikov, Omer	Rico Sennrich, Barry Haddow, and Alexandra Birch.	1041
980	Antverg, Oskar van der Wal, Rui Zhang, Ruochen	2016. Neural machine translation of rare words with	1042
981	Zhang, Sebastian Gehrmann, Shachar Mirkin, Shani	subword units . In <i>Proceedings of the 54th Annual</i>	1043
982	Pais, Tatiana Shavrina, Thomas Scialom, Tian Yun,	<i>Meeting of the Association for Computational Lin-</i>	1044
983	Tomasz Limisiewicz, Verena Rieser, Vitaly Protasov,	<i>guistics (Volume 1: Long Papers)</i> , pages 1715–1725,	1045
984	Vladislav Mikhailov, Yada Pruksachatkun, Yonatan	Berlin, Germany. Association for Computational Lin-	1046
985	Belinkov, Zachary Bamberger, Zdeněk Kasner, Al-	guistics.	1047
986	lice Rueda, Amanda Pestana, Amir Feizpour, Ammar		
987	Khan, Amy Faranak, Ana Santos, Anthony Hevia,	Aleksandar Stanić, Dylan Ashley, Oleg Serikov, Louis	1048
988	Antigona Unldreaj, Arash Aghagol, Arezoo Abdol-	Kirsch, Francesco Faccio, Jürgen Schmidhuber,	1049
989	lahi, Aycha Tammour, Azadeh HajiHosseini, Bahareh	Thomas Hofmann, and Imanol Schlag. 2023. The	1050
990	Behroozi, Benjamin Ajibade, Bharat Saxena, Car-	languini kitchen: Enabling language modelling re-	1051
991	los Muñoz Ferrandis, Daniel McDuff, Danish Con-	search at different scales of compute . <i>arXiv preprint</i>	1052
992	tractor, David Lansky, Davis David, Douwe Kiela,	<i>arXiv:2309.11197</i> .	1053
993	Duong A. Nguyen, Edward Tan, Emi Baylor, Ez-		
994	inwanne Ozoani, Fatima Mirza, Frankline Onon-	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier	1054
995	iwu, Habib Rezanejad, Hessie Jones, Indrani Bhat-	Martinet, Marie-Anne Lachaux, Timothée Lacroix,	1055
996	tacharya, Irene Solaiman, Irina Sedenko, Isar Ne-	Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal	1056
997	jadgholi, Jesse Passmore, Josh Seltzer, Julio Bonis	Azhar, Aurelien Rodriguez, Armand Joulin, Edouard	1057
998	Sanz, Livia Dutra, Mairon Samagaio, Maraim El-	Grave, and Guillaume Lample. 2023a. LLaMA:	1058
999	badri, Margot Mieskes, Marissa Gerchick, Martha	Open and efficient foundation language models .	1059
1000	Akinlolu, Michael McKenna, Mike Qiu, Muhammed	<i>arXiv preprint arXiv:2302.13971</i> .	1060
1001	Ghuri, Mykola Burynok, Nafis Abrar, Nazneen Ra-		
1002	jani, Nour Elkott, Nour Fahmy, Olanrewaju Samuel,	Hugo Touvron, Louis Martin, Kevin Stone, Peter Al-	1061
1003	Ran An, Rasmus Kromann, Ryan Hao, Samira Al-	bert, Amjad Almahairi, Yasmine Babaei, Nikolay	1062
1004	izadeh, Sarmad Shubber, Silas Wang, Sourav Roy,	Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti	1063
1005	Sylvain Viguier, Thanh Le, Tobi Oyebade, Trieu Le,	Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton	1064
1006	Yoyo Yang, Zach Nguyen, Abhinav Ramesh Kashyap,	Ferrer, Moya Chen, Guillem Cucurull, David Esiobu,	1065
1007	Alfredo Palasciano, Alison Callahan, Anima Shukla,	Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller,	1066
1008	Antonio Miranda-Escalada, Ayush Singh, Benjamin	Cynthia Gao, Vedanuj Goswami, Naman Goyal, An-	1067
1009	Beilharz, Bo Wang, Caio Brito, Chenxi Zhou, Chirag	thony Hartshorn, Saghar Hosseini, Rui Hou, Hakan	1068
1010	Jain, Chuxin Xu, Clémentine Fourrier, Daniel León	Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa,	1069
1011	Periñán, Daniel Molano, Dian Yu, Enrique Manjava-	Isabel Kloumann, Artem Korenev, Punit Singh Koura,	1070
1012	cas, Fabio Barth, Florian Fuhrmann, Gabriel Altay,	Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Di-	1071
1013	Giyaseddin Bayrak, Gully Burns, Helena U. Vrabec,	ana Liskovich, Yinghai Lu, Yuning Mao, Xavier Mar-	1072
1014	Imane Bello, Ishani Dash, Jihyun Kang, John Giorgi,	tinnet, Todor Mihaylov, Pushkar Mishra, Igor Moly-	1073
1015	Jonas Golde, Jose David Posada, Karthik Ranga-	bog, Yixin Nie, Andrew Poulton, Jeremy Reizen-	1074
1016	sai Sivaraman, Lokesh Bulchandani, Lu Liu, Luisa	stein, Rashi Rungta, Kalyan Saladi, Alan Schelten,	1075
1017	Shinzato, Madeleine Hahn de Bykhovetz, Maiko	Ruan Silva, Eric Michael Smith, Ranjan Subrama-	1076
1018	Takeuchi, Marc Pàmies, Maria A Castillo, Mari-	nian, Xiaoqing Ellen Tan, Binh Tang, Ross Tay-	1077
1019	anna Nezhurina, Mario Sängner, Matthias Samwald,	lor, Adina Williams, Jian Xiang Kuan, Puxin Xu,	1078
1020	Michael Cullan, Michael Weinberg, Michiel De	Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan,	1079
1021	Wolf, Mina Mihaljcic, Minna Liu, Moritz Freidank,	Melanie Kambadur, Sharan Narang, Aurelien Ro-	1080
1022	Myungsun Kang, Natasha Seelam, Nathan Dahlberg,	driguez, Robert Stojnic, Sergey Edunov, and Thomas	1081
1023	Nicholas Michio Broad, Nikolaus Muellner, Pascale	Scialom. 2023b. Llama 2: Open foundation and fine-	1082
1024	Fung, Patrick Haller, Ramya Chandrasekhar, Renata	tuned chat models . <i>arXiv preprint arXiv:2307.09288</i> .	1083
1025	Eisenberg, Robert Martin, Rodrigo Canalli, Rosaline		
1026	Su, Ruisi Su, Samuel Cahyawijaya, Samuele Garda,	Alex Wang, Amanpreet Singh, Julian Michael, Felix	1084
1027	Shlok S Deshmukh, Shubhanshu Mishra, Sid Ki-	Hill, Omer Levy, and Samuel R. Bowman. 2019.	1085
1028	blawi, Simon Ott, Sinee Sang-aroonsiri, Srishti Ku-	GLUE: A multi-task benchmark and analysis plat-	1086
1029	mar, Stefan Schweter, Sushil Bharati, Tanmay Laud,	form for natural language understanding . In <i>Internat-</i>	1087
1030	Théo Gigant, Tomoya Kainuma, Wojciech Kusa, Ya-	<i>ional Conference on Learning Representations</i> .	1088
1031	nis Labrak, Yash Shailesh Bajaj, Yash Venkatraman,		
1032	Yifan Xu, Yingxin Xu, Yu Xu, Zhe Tan, Zhongli	Zirui Wang, Zachary C Lipton, and Yulia Tsvetkov.	1089
1033	Xie, Zifan Ye, Mathilde Bras, Younes Belkada, and	2020. On negative interference in multilingual mod-	1090
1034	Thomas Wolf. 2023. BLOOM: A 176B-parameter	els: Findings and a meta-learning treatment . In	1091
1035	open-access multilingual language model . <i>arXiv</i>	<i>Proceedings of the 2020 Conference on Empirical</i>	1092
1036	<i>preprint arXiv:2211.05100</i> .	<i>Methods in Natural Language Processing (EMNLP)</i> ,	1093
		pages 4438–4450.	1094
1037	Anton Schäfer, Thomas Hofmann, Imanol Schlag, and	Chris Wendler, Veniamin Veselovsky, Giovanni Monea,	1095
1038	Tiago Pimentel. 2024. On the effect of (near) dupli-	and Robert West. 2024. Do Llamas work in English?	1096

- 1097 [On the latent language of multilingual transformers.](#)
1098 *arXiv preprint arXiv:2402.10588.*
- 1099 Shijie Wu and Mark Dredze. 2019. [Beto, bentz, be-](#)
1100 [cas: The surprising cross-lingual effectiveness of](#)
1101 [BERT](#). In *2019 Conference on Empirical Methods in*
1102 *Natural Language Processing and 9th International*
1103 *Joint Conference on Natural Language Processing,*
1104 *EMNLP-IJCNLP 2019*, pages 833–844.
- 1105 Jiacheng Ye, Xijia Tao, and Lingpeng Kong. 2023. [Lan-](#)
1106 [guage versatilists vs. specialists: An empirical revis-](#)
1107 [iting on multilingual transfer ability.](#) *arXiv preprint*
1108 *arXiv:2306.06688.*

A Experimental Setup

Model. We use a GPT-2-style decoder-only transformer architecture in our experiments (Radford et al., 2019). Unless otherwise noted, we instantiate our model with 12 layers and a hidden size of 768, which results in 85M non-embedding parameters; this corresponds to Languini’s gpt-small configuration. We follow previous work and train our models with sequence length 512, batch size 128, the Adam optimiser (Kingma and Ba, 2015), and a cosine learning rate schedule from $6e-4$ to $6e-6$ with 500 warmup steps.

Data. For the English settings, we use Languini’s default datasets to train and evaluate our models. These are English books from a filtered version of the books3 subset from the Pile (Gao et al., 2020). The train set consists of a total of 23.9B tokens, while the test set contains i.i.d. books, with a total of roughly 11M tokens. This data is tokenised into a vocabulary of size 16k, obtained using a BPE tokeniser trained with SentencePiece (Gage, 1994; Sennrich et al., 2016; Kudo and Richardson, 2018). For our experiments in French, we use the French-PD-Books dataset (PleIAS, 2024), to which we apply the preprocessing pipeline of the Languini Kitchen, but for French. We train a separate BPE tokeniser on this French dataset, using a 16k-sized vocabulary. Depending on the experiment, the French and English vocabularies are either kept separate (disjoint) or merged (anchored). Unless otherwise noted, we train our models for 18,265 steps—i.e., the first 1.2B tokens in our dataset; this corresponds to a GPT small model trained for 6h on an RTX 3090 GPU, the Languini GPT small 6h setting (Stanić et al., 2023). For experiments where we compare hidden representations or gradients on parallel French–English or cloned English sequences, we use data from the Europarl parallel corpus (Koehn, 2005).

Evaluation. When evaluating PPL (from which we also compute MLPE, MLTE and TEff) on the held-out test set, we want to ensure sufficient context for all predictions. To this end, we use a sliding window with steps of 128: we fill in a 512 tokens context, ignore the model’s outputs on the initial 384, and evaluate it only using the last 128 tokens.

B Fitted Scaling Laws

1132

To predict the performance of monolingual models depending on the amount of tokens they are trained on, we fit a power law curve to predict the relationship between number of training tokens and perplexity for models of all three sizes and for both languages (see Fig. 5).

1133

1134

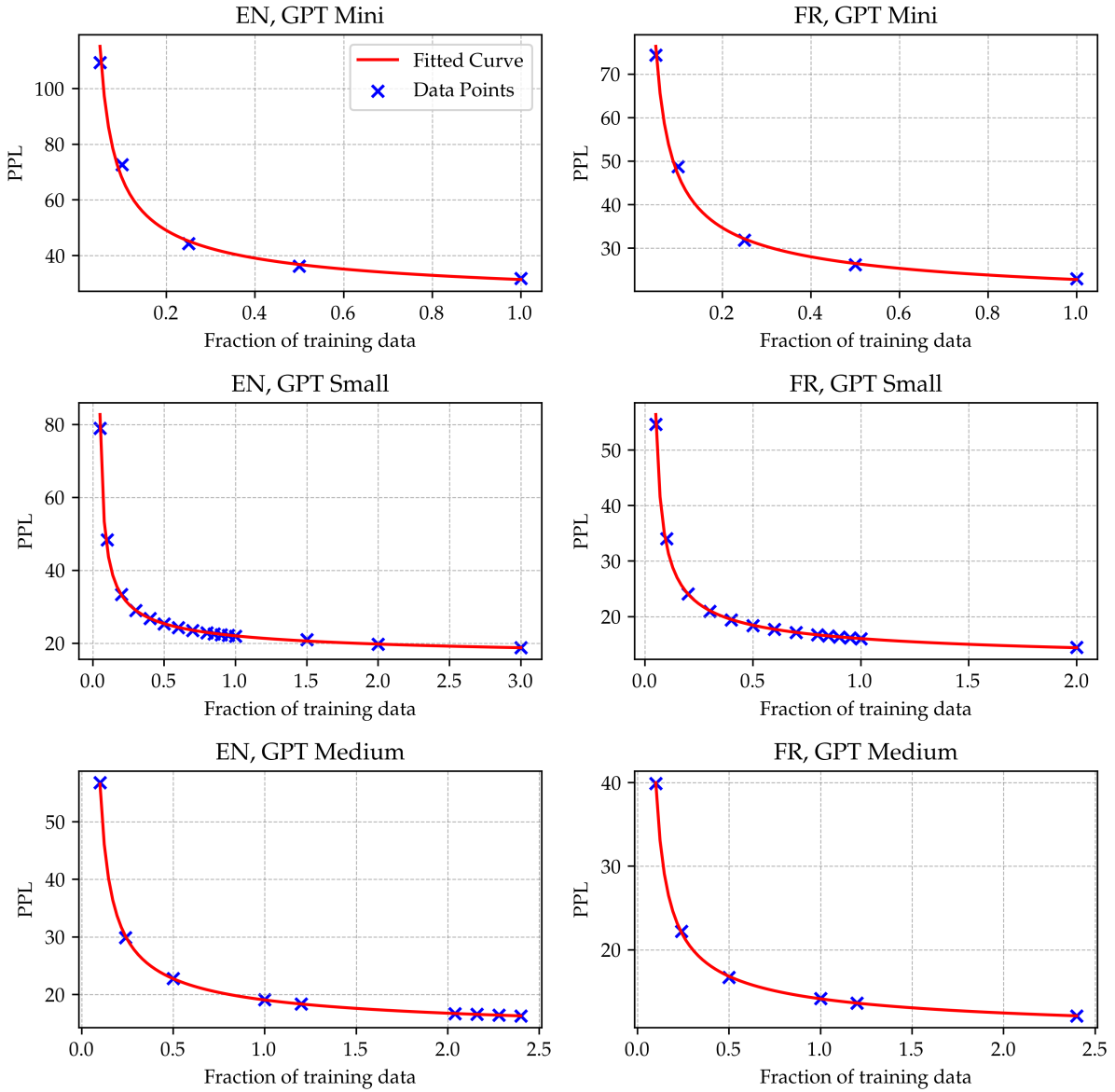


Figure 5: Fitted power laws curves predicting perplexity depending on the fraction of training tokens (compared to our standard 1.2B tokens) for different languages and model sizes.

1135

C Alignment of EN_1 and EN_2 Representations

1137

1138

1139

1140

1141

1142

1143

1144

While, under balanced language sampling, embeddings of corresponding subwords are not much more aligned than embeddings of random pairs, we observe an increase in cosine similarity with increasing language imbalance: from 0.02 for $50/50$ to 0.28 for $95/5$ (see Fig. 6). Fig. 7 shows that this alignment is higher for frequent subwords. This seems natural: at initialisation, subword embeddings are random and not aligned. Then, they become more and more aligned over the course of training.

Interestingly, the embeddings of a simple word2vec (Mikolov et al., 2013) model do not show stronger alignment under higher imbalance. This might be due to a lack of shared parameters between the languages (Conneau et al., 2020b).

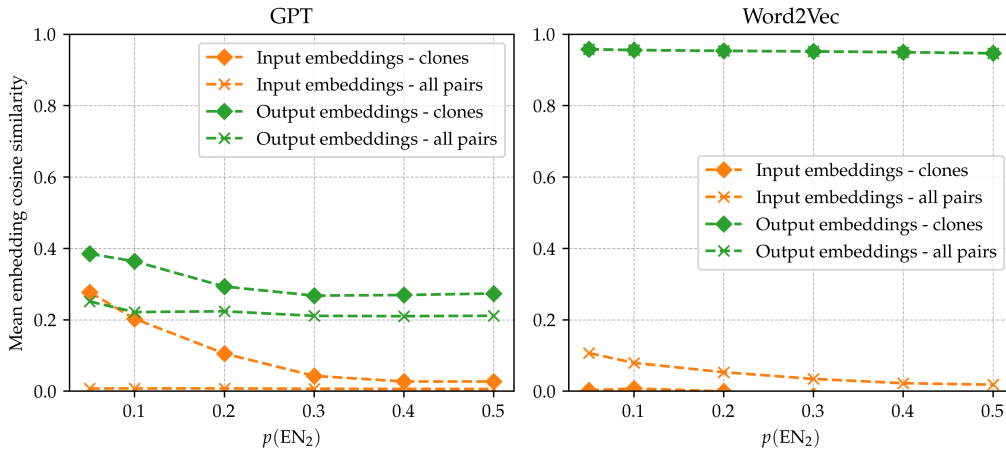


Figure 6: Embedding cosine similarity of corresponding duplicate subwords from EN_1 and EN_2 and random pairs to control for anisotropy. Left: our GPT model. Right: Word2vec embeddings trained on the same data (computed with Gensim).

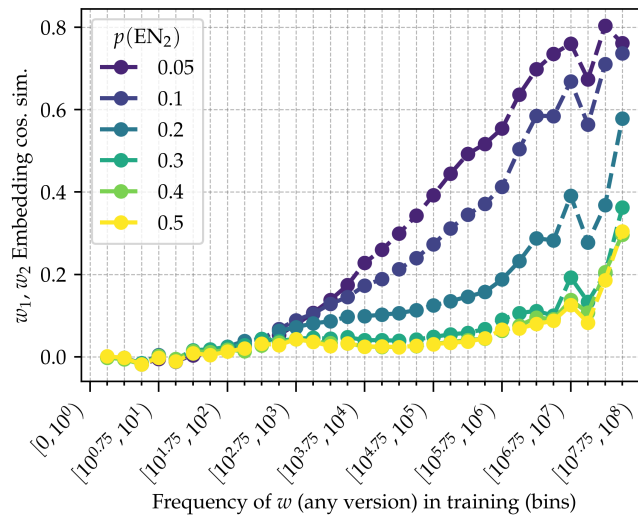


Figure 7: Embedding cosine similarity of corresponding cloned subwords $w_1 \doteq w_2$ from EN_1 and EN_2 , by frequency.

D Anchor Points

1145

D.1 Anchors on Cloned Languages

1146

As described earlier, previous works found that anchor points—i.e., lexical items which overlap between languages—can lead to better generalisation and alignment of representations (Dufter and Schütze, 2020; Pires et al., 2019; Wu and Dredze, 2019). In our cloned setting, we can investigate this in a controlled manner by varying the number of vocabulary elements we duplicate. While in the experiments described above we created EN_2 by duplicating the entire vocabulary, we now duplicate only a fraction. The remaining vocabulary is shared between EN_1 and EN_2 . In this experiment, we observe that a small number of anchor points already significantly boosts model performance (see Fig. 8), which indicates improved generalisation.

1147

1148

1149

1150

1151

1152

1153

1154

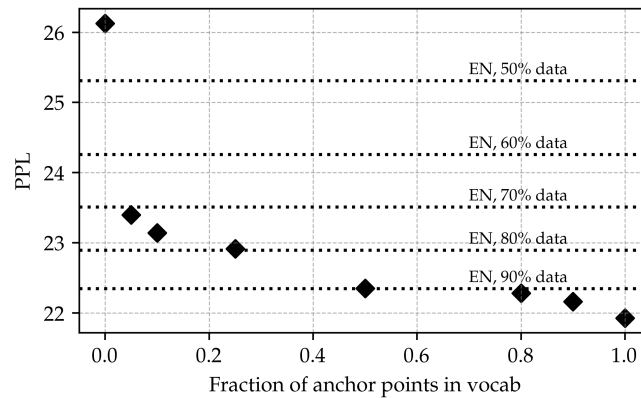


Figure 8: Perplexity by percentage of anchor points, i.e., overlap between EN_1 and EN_2 vocabularies. Models trained on balanced EN_1/EN_2 split.

D.2 Anchors on Real Languages

1155

English and French vocabularies naturally overlap, having common subwords. These shared elements potentially act as anchors, facilitating better cross-lingual generalisation. However, the effectiveness of such anchor points may be moderated by semantic differences; for instance, a shared subword might carry a different meaning or connotations in English and French, affecting its utility as an anchor. Despite these nuances, anchor points appear to boost generalisation between real languages: when we merge the EN and FR vocabularies, we obtain better performance on both languages (compare Table 2, row 7 vs 11) as well as higher alignment of gradients (see App. G). This aligns with our findings from the cloned language setting (see App. D.1). Given these benefits, it is natural to use an anchored (i.e., merged) vocabulary when possible.⁷

1156

1157

1158

1159

1160

1161

1162

1163

1164

⁷In practice, this is usually achieved by training a tokeniser on multilingual data, instead of merging monolingually trained vocabularies.

E Larger Models and More Data

Fig. 9 and Fig. 10 contain results for the full array of model- and dataset size combinations we ran for cloned languages and for English and French, respectively.

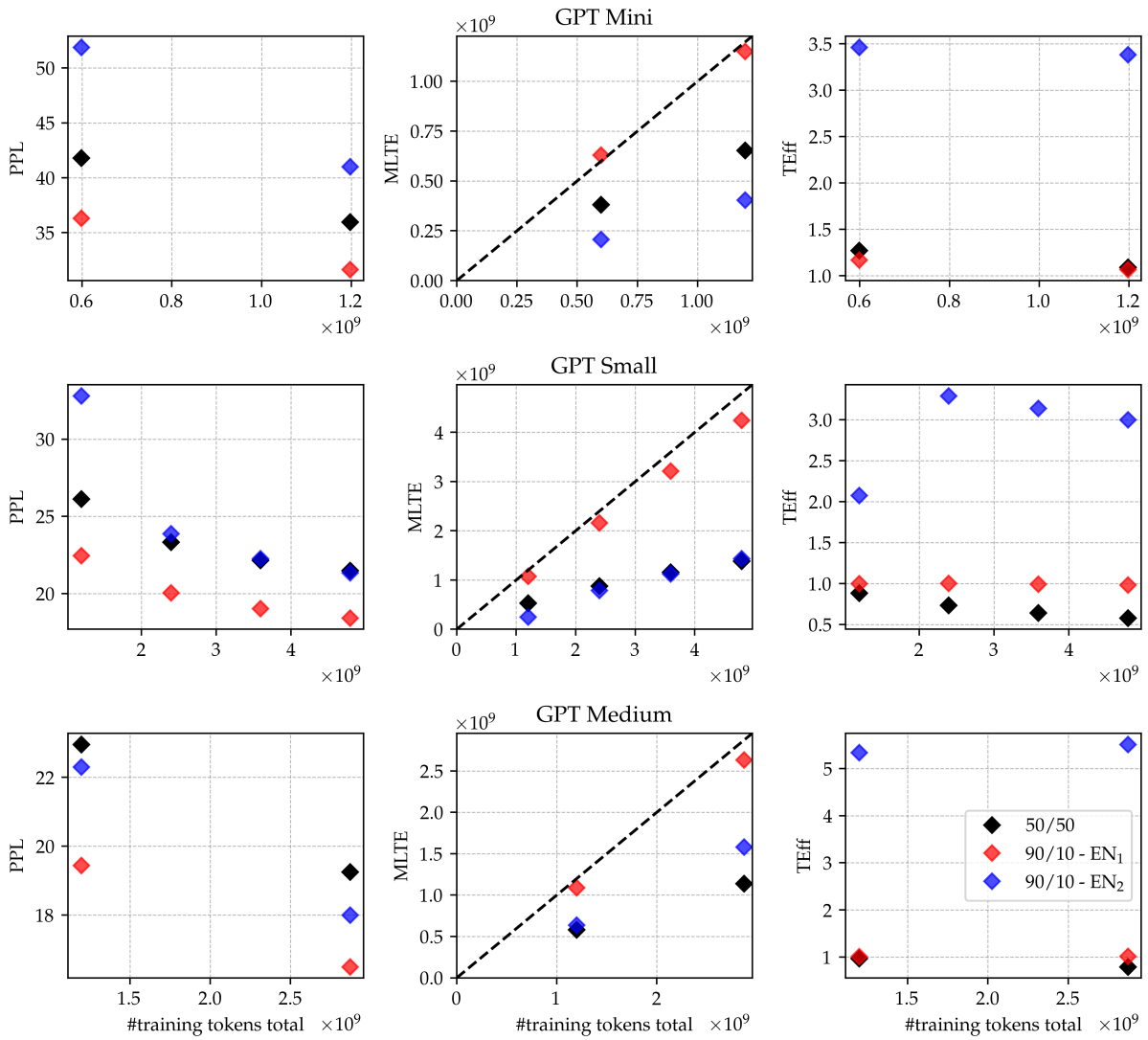


Figure 9: Performance with balanced and imbalanced EN₁ and EN₂ data for different configurations of model- and dataset size

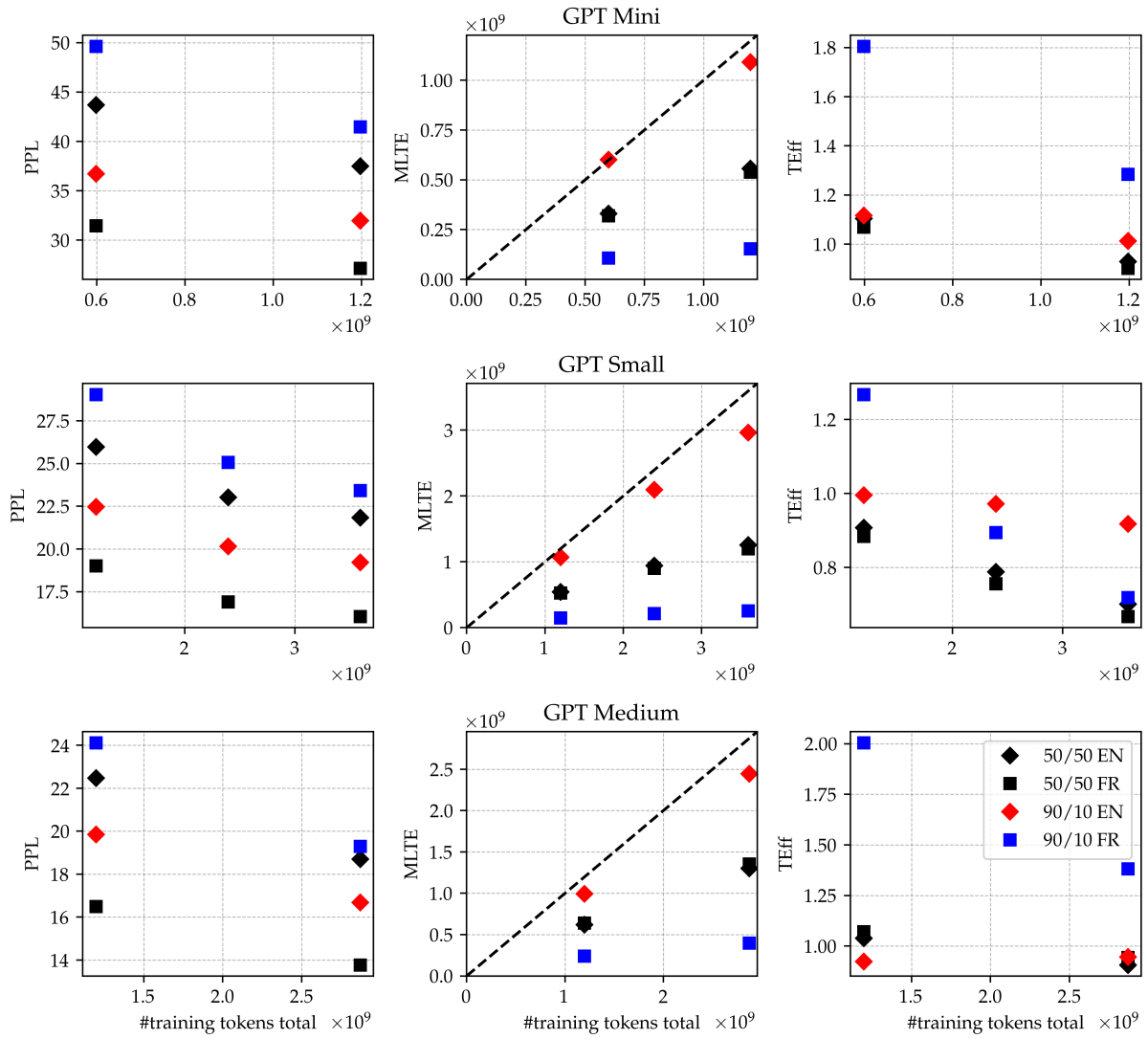


Figure 10: Performance with balanced and imbalanced EN and FR data for different configurations of model- and dataset size. Using anchored vocabulary.

F Hidden State Similarity

1169
1170
1171
1172
1173
1174
1175
1176
1177

Here, we compare the hidden states of our model when processing parallel sequences, both in cloned languages (see Table 3) and in English and French (see Table 4). I.e., for a given trained model and parallel sequences w_a and w_b , we first feed w_a through the model, then w_b , and finally compute the cosine similarities for the hidden states of pairs of corresponding tokens from w_a and w_b (see App. H for details on how these pairs are determined). We use 500 parallel sequences obtained from the Europarl parallel corpus (Koehn, 2005). For cloned languages, we observe that hidden states of the model trained under higher language imbalance generally have higher cosine similarity than the those of the model trained in a balanced setting. For English and French such a trend is less clear. Interestingly, however, an anchored vocabulary seems to lead to slightly higher similarities of the hidden states.

Training Data		Layer											
$p(\text{EN}_1)$	$p(\text{EN}_2)$	1	2	3	4	5	6	7	8	9	10	11	12
50%	50%	0.55	0.79	0.83	0.88	0.85	0.83	0.78	0.66	0.56	0.46	0.25	-0.21
90%	10%	0.86	0.93	0.96	0.96	0.96	0.96	0.96	0.95	0.94	0.90	0.67	0.11
Δ		0.31	0.14	0.13	0.09	0.11	0.14	0.18	0.28	0.38	0.44	0.42	0.32

Table 3: Hidden states’ cosine similarity when LM is fed equivalent inputs in cloned languages. Similarity is computed per token (i.e., comparing pairs of equivalent tokens).

Training Data		Layer												
$p(\text{EN})$	$p(\text{FR})$	1	2	3	4	5	6	7	8	9	10	11	12	
Disjoint	50%	50%	0.68	0.80	0.84	0.88	0.86	0.84	0.80	0.75	0.62	0.53	0.34	-0.15
	90%	10%	0.71	0.83	0.88	0.87	0.86	0.84	0.81	0.74	0.69	0.57	0.40	-0.17
	Δ		0.03	0.03	0.04	-0.01	0.00	0.00	0.01	0.00	0.07	0.04	0.06	-0.03
Anchored	50%	50%	0.73	0.84	0.88	0.91	0.89	0.88	0.85	0.78	0.71	0.61	0.36	0.10
	90%	10%	0.78	0.87	0.89	0.91	0.89	0.87	0.84	0.77	0.72	0.63	0.28	0.06
	Δ		0.05	0.03	0.01	0.00	0.00	-0.01	0.00	-0.01	0.00	0.02	-0.08	-0.04

Table 4: Hidden states’ cosine similarity for parallel inputs in EN and FR for anchored and disjoint vocabularies. We first match which tokens correspond to each other in the two languages, and then compare their representations (see App. H).

G Gradient Similarity

1178
1179
1180
1181
1182
1183
1184

Here, we compare the cosine similarity of trained models' gradients with respect to parallel sequences in two different (possibly cloned) languages. For cloned languages, the alignment between gradients is significantly higher for the model trained in the imbalanced $90/10$ setting (see Fig. 11). For EN and FR data, this does not seem to be the case, whether the vocabulary is anchored (see Fig. 12) or disjoint (see Fig. 13). However, under the anchored vocabulary, the gradient similarities appear to be generally higher, suggesting better cross-lingual representation alignment.

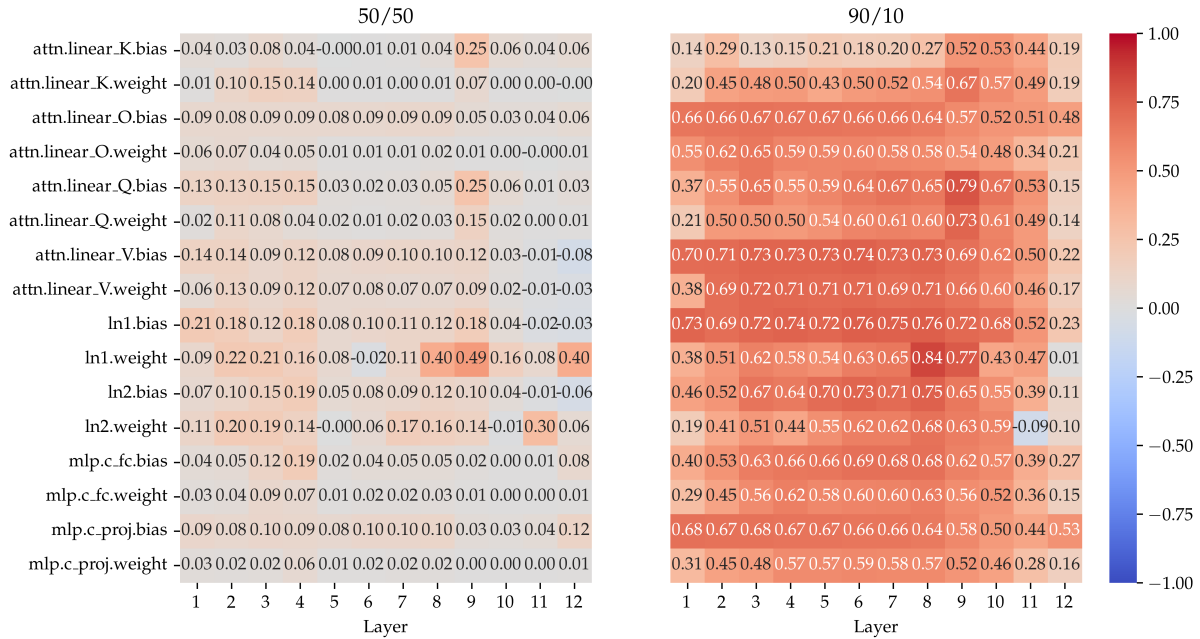


Figure 11: Similarity of gradients with respect to parallel sequences in EN_1 and EN_2 for models trained in balanced and imbalanced settings. Macro average for $50/50$: 0.07. Macro average for $90/10$: 0.53.

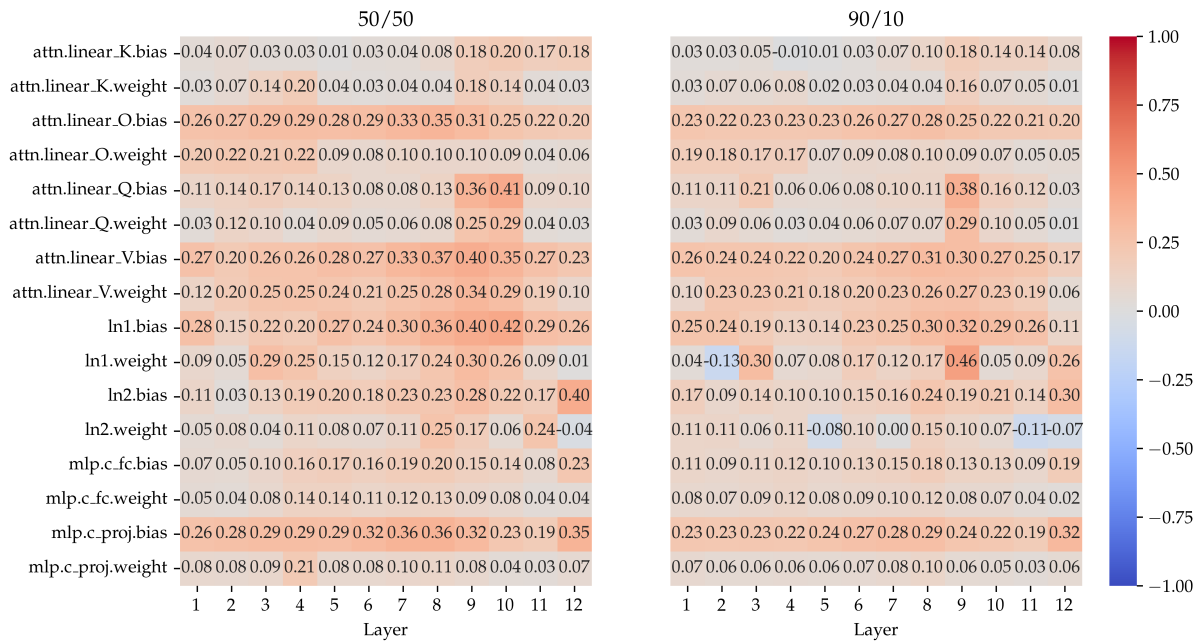


Figure 12: Similarity of gradients with respect to parallel sequences in EN and FR for models with anchored (i.e., merged) vocabulary, trained in balanced and imbalanced settings. Macro average for 50% : 0.17. Macro average for 90% : 0.14.

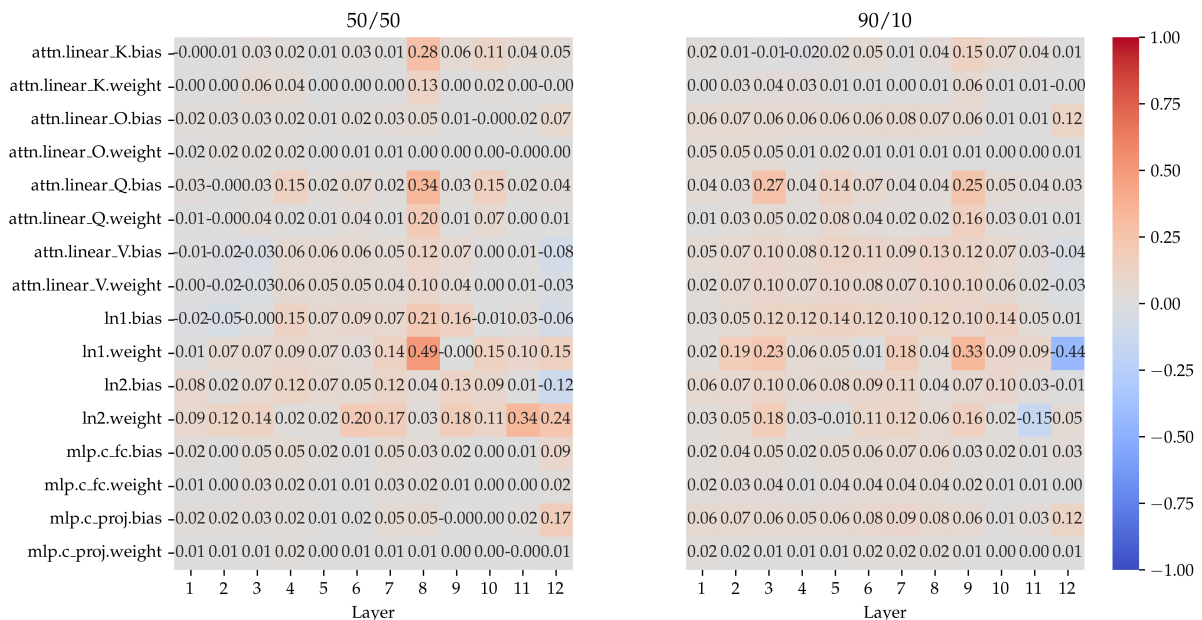


Figure 13: Similarity of gradients with respect to parallel sequences in EN and FR for models with disjoint vocabularies, trained in balanced and imbalanced settings. Macro average for 50% : 0.04. Macro average for 90% : 0.05.

H Matching Corresponding Tokens

1185

In our experiments in §5.2, we employ parallel sequences in different languages and compare both their hidden states’ and their gradients’ similarity. 1186 1187

When comparing gradients (see App. G), we adopt a setup that is analogous to the training process as we aim to understand how one language might affect optimisation of the other: we compute gradients with respect to a full sequence in each language, and then compare these sequence-level aggregated gradients. Analogously, during training, gradient updates are also aggregated for entire sequences. (In fact, during training, these updates are also aggregated for an entire batch, but we use a batch size of 1 for this evaluation.) 1188 1189 1190 1191 1192 1193

However, when comparing hidden states, we compare the individual representations of corresponding tokens in the two sequences. We first compute the cosine similarity of each equivalent token pair, and only then average over the sequence dimension; this provides us with a more informative signal. For parallel sequences $w_{EN_1} \stackrel{\circ}{=} w_{EN_2}$ in cloned languages, it is clear which token corresponds to which: At each given position t , we know that $w_{EN_1,t} \stackrel{\circ}{=} w_{EN_2,t}$ so we can simply compare the hidden states position by position (see Table 3). 1194 1195 1196 1197 1198 1199

Yet, this might not be the case for real languages EN and FR, e.g., due to differing word order or tokenisation. To ensure that we still compare the hidden states of tokens that approximately correspond to each other in the respective languages, we match them based on their cosine similarity scores. Concretely, we create a bipartite graph where the nodes consist of the tokens of the two sequences. For every pair of tokens $w_{EN,t}$ and $w_{FR,t'}$ we add an edge which is weighed by the mean cosine similarity of their hidden states across all layers. We then compute a maximum weight full matching in this graph.⁸ Such a matching maximises the average similarity across all token pairs. Indeed, the resulting token pairs appear to approximately correspond to each other (see Fig. 14). We can then compare the hidden states of these pairs (see Table 4). 1200 1201 1202 1203 1204 1205 1206 1207 1208

Notably, the cosine similarities of hidden states of corresponding EN and FR tokens $w_{EN,t}$ and $w_{FR,t'}$ computed in this way generally appear to be slightly higher than for corresponding tokens $w_{EN_1,t} \stackrel{\circ}{=} w_{EN_2,t}$ of cloned languages (compare Table 4 (disjoint) and Table 3). This might seem unexpected, given that $w_{EN_1,t}$ and $w_{EN_2,t}$ are perfectly equivalent but $w_{EN,t}$ and $w_{FR,t'}$ are generally not. Could this be an artifact of the employed matching strategy which always maximises the average similarity, potentially matching tokens that have very high similarity but are completely unrelated? If this is the case, we should also obtain higher similarity scores in the cloned setting when using the described matching strategy instead of comparing position by position. After running this experiment, we find that using the matching strategy the similarities under the 50/50 cloned language split are indeed marginally higher, although only in the last layers. Under the 90/10 split, however, we observe no notable changes. It thus seems that the proposed matching strategy does not artificially inflate similarity scores too strongly. 1209 1210 1211 1212 1213 1214 1215 1216 1217 1218 1219

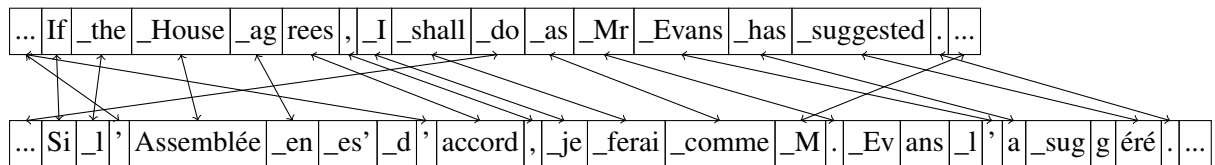


Figure 14: Computed matching for an example sentence using a model trained under 50/50 split with anchored vocabulary. Pointers to “...” denote a match with a token earlier or later in the sequence.

⁸We compute the matching using the NetworkX (Hagberg et al., 2008) implementation of the algorithm proposed by Karp (1978).