

5,000 Words is the Optimal Training Size for Low-Resource Morphology

Anonymous ACL submission

Abstract

Utilizing NLP to assist data annotation remains a challenge for low-resource languages. This study shows that the optimal minimal training size for morphology is about 5,000 labeled tokens. After that, annotation results in diminishing returns in terms of model performance. We assess performance improvement in relation to annotated dataset size for a neural Transformer model, a pre-trained llama3 model, and a non-neural CRF model, all of which are utilized in an active learning loop. We analyze corpus diversity (via type-token ratio) to better understand how sample diversity impacts active learning model improvement and find raw TTR scores indicate peak performance.

1 Introduction

Creating initial digital resources in a language requires primarily manual methods that cannot scale to the pace of language endangerment threatening roughly half the world’s 7,100+ languages (Eberhard et al., 2025; Seifart et al., 2018) nor to the volume of annotated data demanded for high-performing natural language processing (NLP) models. Even though optimal performance does not wholly depend on data set size (Liu and Prud’hommeaux, 2022; Palmer, 2009), NLP must begin with some amount of annotated data. Unfortunately for low resource languages, “low resource” does not refer just to limited digital language data (Liu et al., 2022). Underfunded academic and community linguists must budget their resources and carefully determine when they have the minimal amount of annotated data that will make optimal contribution towards their goals. The question then is: how much training data should we start with?

This explores how to determine the minimal optimal training data size if no annotated data yet exists. To this end, we use four computational morphology tasks (see 4.1) as the test cases, leveraging *documentary collections from linguistic fieldwork* in six

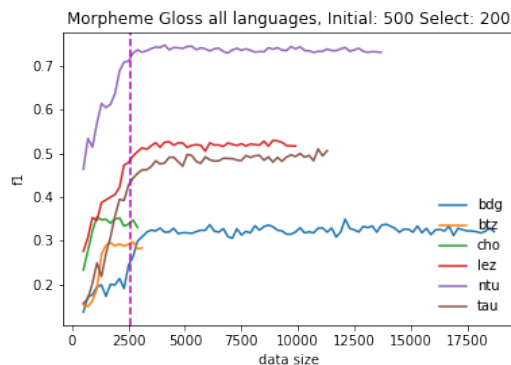


Figure 1: Annotating training data for morphology quickly reaches a point of diminishing returns in terms of model performance. For morpheme glossing with the Transformer the minimal optimal training size is about 2,000 tokens. On average across models, languages, and tasks, the rate of model improvement begins to flatten as total training size grows to about 5,000 tokens.

diverse endangered languages. We assume that at some point model performance becomes only marginally better than a model trained on less data. At this point, model performance begins to flatten and annotation efforts yield diminishing returns (e.g., the plateauing trends in Figure 1). Adopting an active learning framework (Settles, 2009), we ask:

- Q1: How quickly do manual annotation efforts yield diminishing returns in terms of model performance? Moreover, in an active learning scenario, what impact does initial size and selection size have on model improvement?
- Q2: What is the minimal optimal size of training data for morphology?
- Q3: How might the minimal optimal size be identified by the lexical diversity of training data?

Since large scale annotation is cost prohibitive for minority language communities, identifying the

earliest point of diminishing returns—the minimal optimal training size—guides field linguists who wish to give back to their host language community by making optimal contribution to technological development in the language.

2 Related Work

Investigations into the impact of data size on NLP offer little concrete guidance for budgeting manual annotation or the optimal initial size when developing early digital resources in a new language. If novices search the internet, they will find the most popular answers are “The more the better!” or “It depends”.¹² Model performance benefits from lots of relevant examples but the ideal number cannot be calculated *a priori* because so many factors can determine how many examples are enough. Popular guidelines to estimate optimal dataset size are premised on the needs of the model, rather than on available resources (Dorfman, 2022; gaumutrapremi, 2025). The “10x Rule” considers computing resources, not available data resources, and recommends having 10 times as many training examples as parameters. Though correct, these answers are unhelpful for low-resource languages.

In low-resource contexts, NLP models and techniques should adjust to available data, not the other way around. Successful adjusted techniques include synthetic data augmentation (Silfverberg et al., 2017; Bergmanis and Goldwater, 2019; Anatasopoulos and Neubig, 2019), or cross-lingual transfer learning (Baumann and Pierrehumbert, 2014; Gu et al., 2018), and active learning (Palmer, 2009). Some research (Vieira et al., 2024; Le Scao and Rush, 2021; Debnath et al., 2021) looks not at minimal training size, but at the most effective number of examples for few-shot prompting, finetuning, or transfer learning from multilingual models. For example, Vieira et al. (2024) found a minimum of 5k examples is optimal for fine-tuning large language models (LLMs) for machine translation. Unfortunately, transfer learning or fine-tuning can yield unreliable results for low-resource languages because even the largest multilingual models fall far short of representing the world’s language diversity (Debnath et al., 2021; Costa-Jussà et al., 2022). Additionally, such techniques are best applied only after some minimal training data exists in the lan-

¹www.quora.com/How-much-data-is-required-for-machine-learning

²www.geeksforgeeks.org/how-much-data-are-sufficient-to-train-my-machine-learning-model

guage or a closely related language, which is not the case for many communities.

Some work does not assume an initial amount of data exists provide a few insights. Pérez-Mayos et al. (2021) conducted a cost-benefit analysis of unsupervised dataset size on downstream syntactic tasks. Their conclusion simply states “more data is better”. In machine translation Maillard et al. (2023) found that across 39 languages if the data is high quality, then diminishing returns being after 6,000 sentence pairs. High quality data is not always obtainable in documentary settings, where field collections tends to be noisy. This number is close to our findings. Looking at unsupervised learning, Eiselen and Gaustad (2023) found that when applying embedding models to part of speech tagging, morphologically complex languages need substantially more data than languages with simpler word structures. This pattern holds fairly well in our experiments, with the polysynthetic languages generally requiring more annotated data.

3 Languages and Data

Our data sets are extracted from realistic, resource-constrained collections of interlinearized glossed texts (IGTs) compiled by documentary linguists in six low-resource languages and shared for research purposes. Table 1 summarizes the original data. These IGTs represent some of the only digital resources available in these languages. We do not publish the data on our Github to reduce possible risks to the minority language communities who provided the data. Most the data is publicly available elsewhere or by contacting the authors. We use only locally run models to reduce the risk that the data may be harvested and used for commercial purposes.

To create our datasets, we first constructed the test set by randomly extracting 200 morphologically annotated types (unique word forms), after removing punctuation and digits (unless the digits represent phonemes in the writing system). The training sets contain all remaining morphologically annotated tokens (non-unique words) after all copies of the test words were removed. There is no overlap between the training and test sets. This leads to an overestimation of the data because in realistic linguistic fieldwork, naturalistic language will contain many repeated words. We believe it is better to annotate too much data when chances of organizing additional annotation efforts may be

158	slim. Relevant information and our code is avail-	Alaska and has been taught at a school in the one	201
159	able on Github. ³	village in the Yukon Territory of Canada. Despite	202
160	Batak Alas-Kluet [btz] is an vigorous Austrone-	its status, it is critically endangered, with about 100	203
161	sian language spoken by 100,000 people on the In-	speakers (Lovick, 2020). Its complex morphology	204
162	donesian island of Sumatra (Eberhard et al., 2025).	features non-continuous lexical, derivational, and	205
163	The selected data is an unpublished corpus from	inflectional prefixes on verbs. The selected corpus	206
164	the Alas-Kluet dialect and features reduplication,	was collected in 2006-2019 in Alaska mostly	207
165	infixation, and circumfixation.	representing speakers of the Tetlin and Northway	208
166	Bonggi [bdg] is a Western Austronesian language	dialects. The primary data is preserved at the	209
167	spoken by some 1,500 people in on Banggi and	Alaska Native Language Archive. ⁷	210
168	Balambangan islands of Malaysia (Eberhard et al.,	4 Experiments	211
169	2025). It morphologically marks nominal case and	4.1 Four morphological tasks	212
170	verbal agreement (Boutin, 1996).	Four morphological tasks add crucial information	213
171	Choctaw [cho] is spoken by less than 10,000 peo-	for early scientific documentation of a language	214
172	ple primarily in Mississippi and Oklahoma. It has	and automating the tasks supports technological de-	215
173	complex morphology but the orthographic repre-	velopment and the creation language maintenance	216
174	sentation in the selected data breaks long words	materials. The tasks are:	217
175	into shorter orthographic units (Broadwell, 2006).	(1) Canonical Morpheme Segmentation (Mott	218
176	The original documents are in the United States	et al., 2020): decomposing a word into its	219
177	National Archive and are available for research pur-	component underlying morphemes	220
178	poses per the usage policy at Oklahoma Historical	input: churches → output: church+s	221
179	Society (OHS). ⁴ The annotated data can be made		222
180	available upon request.	(2) Morpheme Glossing (Okabe and Yvon,	223
181	Lezgi [lez] (Lezgian) is a Nakh-Daghestanian	2023): converting a sequence of compo-	224
182	(Northeast Caucasian) language spoken by over	nent morphemes to technical translation	225
183	500,000 speakers in Russia and Azerbaijan (Eber-	(glosses) ⁸	226
184	hard et al., 2025). The corpus used is from the	input: church+s → output: church+PL	227
185	Qusar dialect in Azerbaijan (Donet, 2014). The		228
186	public parts of the corpus are archived in ELAR	(3) Joint segmentation and glossing (Moeller	229
187	(The Endangered Languages Archive). ⁵ Lezgi is a	and Hulden, 2021) treats (1) and (2) as a	230
188	highly agglutinative language with overwhelmingly	single task, imitating common practice by	231
189	suffixing morphology (Haspelmath, 1993).	linguists.	232
190	Natügu [ntu] belongs to the Reefs-Santa Cruz	input: churches →	233
191	group in the Austronesian family. It is spoken by	output: church#religious.body + s#PL	234
192	about 5,000 people in the Temotu Province of the		235
193	Solomon Islands. It has primarily agglutinative	(4) Inflection (McCarthy et al., 2019): provid-	236
194	morphology with complex verb structures (Næss	ing only glosses for the inflectional affixes	237
195	and Boerger, 2008). The selected corpus contains	input: churches → output: PL	238
196	transcribed narratives and a large written text. ⁶		239
197	Upper Tanana [tau] (Nee'aaneegn', Nabesna,	4.2 Models	240
198	Tabesna) belongs to the Northern Dene	We compare a neural, non-neural, and large lan-	241
199	(Athabaskan) family. It is an official lan-	guage model (LLM). We focus on the Transformer	242
200	guage of Alaska, spoken in the eastern interior of	model because it is expected to give lowest per-	243

³withheld in review version because it is not anonymized

⁴gateway.okhistory.org/help/faq/rights

⁵The direct link is <http://hdl.handle.net/2196/439f3247-66f3-4b76-b5e1-8f07f0d7b5a1>

⁶A Natqgu grammar and much of the text data are available at <https://www.langxmelanesia.com/tilp>

⁷<https://www.uaf.edu/anla>

⁸3SG ← third person singular, PRES ← present tense

Language	ISO 639-3	Family	Typology	Tokens
Alas-Kluet	btz	Austronesian	Agglutinative	3,839
Bonggi	bdg	Austronesian	Non-concatenative	21,616
Choctaw	cho	Muskogean	Agglutinative	7,793
Lezgi	lez	Nakh-Daghestanian	Agglutinative	13,953
Natügu	ntu	Austronesian	Agglutinative	16,544
Upper Tanana	tau	Athabaskan	Polysynthetic	14,099

Table 1: Languages in our study; tokens are non-unique words with morphological annotation.

overestimates the minimal optimal amount of training data.

Transformer. For all settings, we employ the Transformer architecture, TRANSFORMER_TINY⁹, from the Python package fairseq. This model architecture has two encoder-decoder layers and two self-attention heads. The embedding dimension size and that of the feed-forward layer are both 64. This is motivated by Liu and Dorr (2024) which show TRANSFORMER_TINY to be the most effective sequence-to-sequence (seq2seq) model for morphological segmentation in low-resource scenarios across 19 typologically distinct languages. All models were trained using the default parameters, maintaining consistency across experimental settings. All training was performed with an NVIDIA A100 GPU with a single CPU core and 8GB of RAM.

CRF. As a comparison to the encoder-decoder architecture, we adopt conditional random field (CRF), a non-neural log-linear model (Lafferty et al., 2001) that has been applied to morphological segmentation (Cotterell et al., 2016b), inflection (Cotterell et al., 2016a), and glossing (Okabe and Yvon, 2023). Prior work has shown that CRF can perform well for morphological segmentation (Liu and Dorr, 2024), often the task of interest was *surface* segmentation, where the linear concatenation of morphemes refers only to the orthography. Surface segmentation misses important meanings (e.g. that the prefix of the following words is underlyingly identical: *il+legal*, *im+possible*, *in+complete*), as opposed to *canonical* segmentation, which ignores surface orthographic changes during word formation that may be due to phonological processes (e.g. the rendering of the plural suffix without the epenthetic *e* in example (1)). Since our morphological data only consists of canonical segmentation and since seq2seq models have gradually become the default architec-

tures for morphological inflection (Kodner et al., 2022), we apply the CRF only to morphological glossing. We employ first-order CRF, which treats glossing as a sequence labeling task: given a sequence of morphemes in a word, the model predicts each morpheme’s English translation. Our hand-curated feature set includes its character-level and morpheme-level unigrams and bigrams, the linear position of the morpheme in the word, and whether the word itself is monomorphemic. All CRF models are implemented with python-crfsuite using a single CPU core and 8GB of RAM.¹⁰

LLM. We additionally assess the tasks on llama3 (Dubey et al., 2024), an open-source LLM. We used the ollama package to run llama3-8b locally with an NVIDIA GeForce RTX 2080 Ti, thereby avoiding any potential data contamination that could result from relying on a closed-source model. The LLM system prompt is stated in Appendix B.

4.3 Machine-in-the-Loop

Machine-in-the-loop methods have been shown to be effective for digital language resources in under-documented languages by leveraging a less accurate but still helpful NLP model to assist manual annotation (Palmer, 2009; Bird and Yibarbuk, 2024; Moeller and Arpe, 2024). We simulate a machine-in-the-loop approach to morphological analysis by applying active learning (AL), which derives better model performance with less training data by selecting additional batches for manual annotation in informed ways (Settles, 2009; Tharwat and Schenck, 2023). First, a model is trained on available annotated data and used to pre-annotate the remaining unannotated data. Information from the model is assigned to each pre-annotated example and that estimate is used to select a sampling from the data to present to humans experts for manual work. Our selection strategy is least confidence uncertainly sampling (Culotta and McCallum, 2005) which se-

⁹<https://fairseq.readthedocs.io/en/latest/models.html>

¹⁰<https://python-crfsuite.readthedocs.io/en/latest/>

lects the N least confident examples based on the model’s entropy or marginal probability estimate. In our case, the new annotated instances are retrieved and added to training set from the IGT to simulate human annotation after each AL cycle.

To determine the relationship between combined training size and model performance, the Transformer and CRF models are trained with different initial training sizes ($\{500, 1,000, 1,500, 2,000, 2,500\}$) which represent all annotation efforts before the AL process begins. To determine whether the model improvement is affected by the number of new tokens selected by the least confidence strategy, different selection sizes ($\{50, 100, 200\}$) are added to the training set at each AL cycle.¹¹ The total training size is the sum of the initial training size and selected data from all previous AL iterations. If the initial training size is 500 and the selection size is 200, the initial model is trained on 500 tokens, the second on 700, and so on.

LLM. For the pre-trained LLM strategy approach, we use a slightly adjusted active learning strategy that leverages the MiniLM-L8 embeddings with furthest-point sampling (FPS), an explore-only batch acquisition function from the botorch package (Balandat et al., 2020). This differs from the explore-only AL approach with the CRF and transformer models, which leverage the statistical uncertainty metrics inherent to their training and evaluation. We use in-context learning (ICL) to present llama3-8b with the desired task (eg. morpheme glossing) and relevant training examples, and iteratively assess F1 accuracy for each sample in the test set. For training dataset sizes consisting of greater than 100 examples, we find that the LLM context window can easily be exceeded. We therefore use retrieval augmented generation (RAG) (Lewis et al., 2020) with a MiniLM-L6 embedding model (Wang et al., 2020) to select the 100 most relevant training examples (according to embedding cosine distance) to use for ICL on each test instance. We assess our LLM+RAG approach on the morpheme glossing task across available languages, with an initial randomly selected set of 500 samples. We use our FPS

¹¹The sizes were chosen based on a pilot study that tested initial training sizes of 50, 100, and 500 tokens and a selection size of 25 tokens and found that an initial training size less than 500 is counterproductive. As the total combined initial+selection size grows past 500, the model performance can actually be lower compared to a model trained on the same total tokens but starting with 500 tokens. We chose the other selection sizes based on what seemed reasonable for an average person to annotate in an hour or less.

acquisition function to acquire subsequent samples in batches of 200, assessing to a maximum of 2500 tokens. Note that the size of the ICL samples remains fixed at 100, but the samples provided to the LLM change and improve as more tokens are iteratively acquired.

4.4 Optimal minimal training size

We hypothesize that for probabilistic supervised models an certain amount of training data exists after which the costs to annotate data are no longer justified because annotation efforts yield diminishing in terms of model performance. Correspondingly, we define the *optimal minimal amount of training data* (hereafter, optimal minimal size) as the point where the rate of model improvement starts to decrease as the amount of annotated data increases monotonically. We estimate the optimal minimal size in three steps. First, we plot the training size x for all combinations of initial and selection sizes and the corresponding model performance y . Second, we find the point of maximum curvature in the plotted line, i.e. the “knee”, using the kneefinder¹² algorithm, which identifies the point that has the maximum distance from a line passing between the first and last $[x, y]$ coordinates. In reality, the knee is not a single point but an area where the curve begins to flatten. So as the third step, we average the knees across the three selection sizes of each initial size; this average provides a better approximation of the area of maximum curvature since the selection size has minimal impact on the $[x, y]$ curvature (see sec. 5.1). We do this for the Transformer and the CRF results. For LLM assessment, we identify knees on the single selection size and acquisition rate.

4.5 Data diversity

The method of calculating the minimal optimal size described above requires already having more than the minimal optimal size. A more generalizable method would identify the optimal minimal size for any given dataset as annotation efforts are ongoing. In search of such a method, we examine training set diversity with the Transformer and CRF. As a measure of diversity, we adapt a simple measure of lexical diversity, the type to token ratio (TTR). In corpus linguistics, “type” refers to a distinct word form and “token” refers to all word forms in a data set (including duplicates). TTR is calculated by

¹²<https://github.com/vlavorini/kneefinder>

417 dividing the total number of types by the total number of tokens and multiplying by 100 to achieve a percentage. For example, the TTR of the sentence “They saw eye to eye.” is 80% because eye is repeated twice, giving ,000 types in a sentence with 5 tokens. We calculate the TTR of each combined training set after every AL cycle and compare the curve trajectories to the model performance trajectory. We measured TTR on the input and output for each task in examples (1)-(4), but the TTR of outputs did not provide useful information, so we consider only the TTR of the inputs in 5.4.

429 5 Results

430 Having compared all combinations of initial and selection data sizes in an AL scenario for all four tasks and all languages, we present the results and discuss each research question.

434 5.1 Q1: Annotation yields diminishing returns

435 As Figure 1 indicates, annotation efforts rather quickly yield diminishing returns in model performance across every language and nearly every task. This tendency towards a logarithmic relationship (with an upper bound) between data size and model performance holds less strongly for inflection and for the two smaller datasets. Alas-Kluet [btz] and Choctaw [cho] do not seem to contain their optimal minimal size, particularly for inflection. This is more clearly seen with the CRF which tends to take longer to show leveling of performance (and yields slightly higher peak F1 scores).

447 **Initial training size** With AL, the amount of initial training data only seems to impact how quickly optimal model performance is reached. For example, starting Alas-Kluet glossing at 500 tokens (and a selection size of 100), the model tops out after ten AL cycles (1,500 total training tokens). Starting with 2000 tokens requires only three iterations to reach the same performance, but more total training tokens.¹³

456 **Selection size** Overall, smaller selection size have a slight positive impact on the rate of model improvement. With the Transformer a selection sizes of 50 reaches the highest overall model performance about 30% of the time and occasionally levels out slightly earlier. With the CRF, however, the smaller selection sizes almost always achieve

¹³Figures showing all settings and training combinations can be found on our GitHub repository.

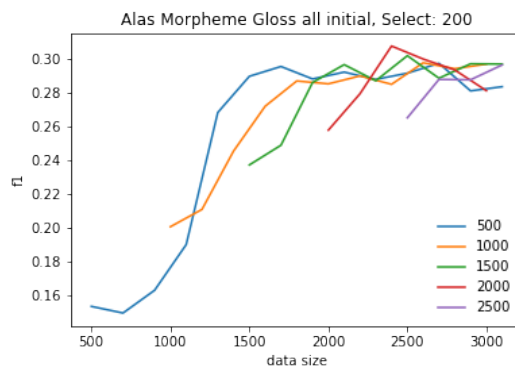


Figure 2: All initial training sizes achieve similar model performance.

463 best performance and level out earlier (about 65% and 45% of of the experiments, respectively).

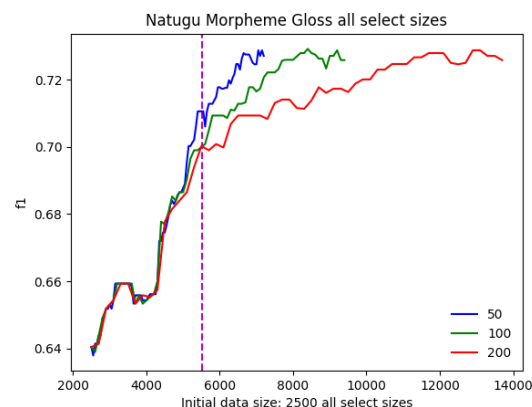


Figure 3: The rate of model improvement diverges slightly with the CRF depending on selection size during AL but does not impact peak performance.

465 5.2 Q2: Optimal minimal size is 5k

466 Considering all combined sizes with the Transformer and CRF, the optimal minimal training size ranges 1,000-12,000 (see Appendix A). The higher amounts are due to larger initial training sizes. The average optimal minimal size is just under 4,000 training tokens with the Transformer and around 5,000 tokens with the CRF. The non-neural model also yields slightly higher F1 scores which is typical of low-resource settings, unless/until strategies such as data augmentation are applied. Alas-Kluet and Choctaw have datasets smaller than these averages and may not contain the optimal minimal size. However, even without those two languages, the average point of diminishing returns from annotation increases slightly to about 4,500 tokens (Transformer) and 6,000 (CRF) tokens. Overall,

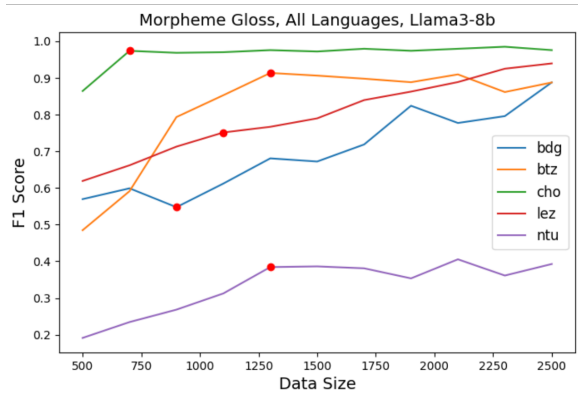


Figure 4: Llama3-8b performance on morpheme glossing with active learning selection and retrieval of relevant examples for ICL. Red points show the approximate region where a knee may be located.

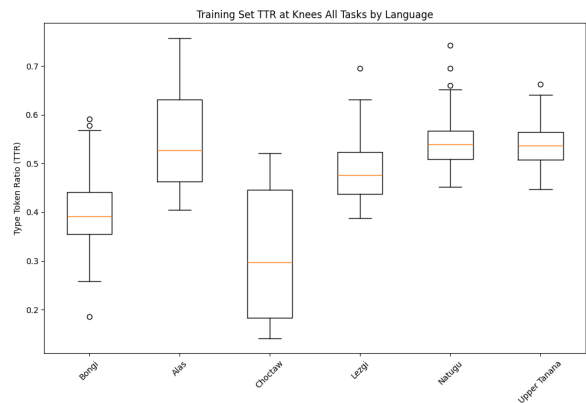


Figure 5: Average TTR scores at point of performance plateauing suggest a raw TTR score of 50% or below indicates the minimal optimal training size is reached.

the average optimal minimal training size is about 5,000 tokens.

5.3 Q1 & Q2: LLM results

Based on preliminary assessments with our pre-trained LLM, we find that for most languages the peak F1 score is significantly higher than the statistical or fine-tuned transformer strategies. Figure 4 shows the F1 scores across languages for the morpheme glossing task, ranging in dataset size between 500 and 2500 labeled tokens. Natügu is a notable exception, with F1 score peaking at around 0.78 for the Transformer model but at only around 0.40 for the LLM. We also find that, according to the kneefinder algorithm, the data size at the point of diminishing returns may be far less than 5,000 tokens for the morpheme glossing task. This is especially the case for Choctaw, where the pre-trained LLM appears saturated at an F1 score of around 0.95 after being provided only 700 labeled tokens. While the Transformer model also peaks at around 700 labelled tokens for the morpheme glossing task, the F1 score remains at only 0.35. In future efforts, we will assess the pre-trained LLM across additional tasks (segmentation and inflection) with larger data sizes to determine whether, for languages like Bonggi and Lezgi, the pre-trained LLM reaches a point of diminishing returns at a later point.

5.4 Q3: TTR indicates optimal training size

After examining changes in data diversity, we find that raw measures of TTR is a useful proxy for optimal minimal size with the Transformer and CRF. Each task follows a unique pattern of TTR change

before showing a final steady decline and the patterns hold across languages and models. Except for segmentation, however, these patterns are not reliable indicators of optimal minimal size. However, raw TTR score is a reasonably reliable indicator that one approaching the optimal minimal size of any given dataset because this size has a very strong tendency to when the TTR score declines to 50% or less especially in the larger datasets, as can be seen in Figure 5. These tendency is more obvious after eliminating the two smallest datasets.

6 Conclusion

This work ascertains the minimum amount of labeled low-resource language data that ought to be acquired manually before expecting statistical, neural, and large language models to provide optimal linguistic assistance. Applying active learning to four morphology tasks in six diverse low-resource languages and using available annotated documentary data, we compare model performance of CRF, Transformer, and llama3. We find that the average minimal optimal training size is about 5,000 words. Our work suggests that raw TTR scores of 50% or less can indicate when the optimal minimal size has been reached during the AL process.

A note on Choctaw and Alas-Kluet: The two data sets under 10k available tokens reached performance plateau more quickly than the larger data sets with the Transformer or not all with the CRF. We conclude that 2-3,000 training tokens is not enough even though we see a leveling of model performance in those datasets. They showed less clear correlation between TTR scores and optimal minimal size. The larger datasets, on the other

549 hand, peaked fairly consistently shortly after 4,000
550 (Transformer) or 6,000 (CRF) training tokens and
551 did not show marked increase even after another
552 10 thousand tokens.

553 Interestingly, 5,000 tokens corresponds nicely
554 with the size of two IGT collections (btz and lez)
555 created by graduate student language documenta-
556 tion projects. (Lezgi’s annotation has since ex-
557 panded.) Although too small to be conclusive we
558 tentatively conclude that early documentary efforts
559 may be producing the minimal optimal training
560 data needed to support initial NLP development.
561 We also recommend starting with an initial size of
562 500 and using an active learning approach with a
563 smaller selection size of 50 or 100. We acknowl-
564 edge smaller selection batches demands more fre-
565 quent computing time and other logistical matters
566 which could slow the journey to optimal perfor-
567 mance. In realistic scenarios, the practicality of
568 frequent retraining should be considered.

569 This paper does not answer a related question,
570 namely “When with AI be useful for my language?”
571 An optimal trade-off between data size and model
572 performance does mean the model is yet useful.
573 At this point we suggest academic and community
574 linguist should cease annotation efforts and seek
575 collaboration with experts in low-resource NLP.
576 Our ongoing work is evaluating when a models
577 becomes useful as an annotation assistant.

578 Limitations

579 Due to time constraints we were not able to run
580 k-fold or average across random seeds. This means
581 the numbers and scores may not hold well as an
582 average performance. We were also not able to run
583 the pre-trained LLM for larger dataset sizes and
584 across segmentation and inflection tasks. Although
585 the set of languages is diverse in morphological,
586 geographical, and genealogical categories, six is
587 a small number of languages. The conclusions
588 we draw should be tested further on other *realistic*
589 *documentary field data*.

590 References

591 Antonios Anastasopoulos and Graham Neubig. 2019.
592 [Pushing the Limits of Low-Resource Morphological](#)
593 [Inflection](#). In *Proceedings of the 2019 Conference on*
594 *Empirical Methods in Natural Language Processing*
595 *and the 9th International Joint Conference on Natu-*
596 *ral Language Processing (EMNLP-IJCNLP)*, pages
597 984–996, Hong Kong, China. Association for Com-
598 putational Linguistics.

Maximilian Balandat, Brian Karrer, Daniel Jiang,
Samuel Daulton, Ben Letham, Andrew G Wilson,
and Eytan Bakshy. 2020. Botorch: A framework
for efficient monte-carlo bayesian optimization. *Ad-*
vances in neural information processing systems,
33:21524–21538.

Peter Baumann and Janet Pierrehumbert. 2014. [Using](#)
[Resource-Rich Languages to Improve Morphologi-](#)
[cal Analysis of Under-Resourced Languages](#). *Pro-*
ceedings of the Ninth International Conference on
Language Resources and Evaluation (LREC-2014).

Toms Bergmanis and Sharon Goldwater. 2019. [Train-](#)
[ing Data Augmentation for Context-Sensitive Neural](#)
[Lemmatizer Using Inflection Tables and Raw Text](#).
In *Proceedings of the 2019 Conference of the North*
American Chapter of the Association for Computa-
tional Linguistics: Human Language Technologies,
Volume 1 (Long and Short Papers), pages 4119–4128,
Minneapolis, Minnesota. Association for Computa-
tional Linguistics.

Steven Bird and Dean Yibarbuk. 2024. Centering the
Mixtec speech community. In *Proceedings of the*
18th Conference of the European Chapter of the As-
sociation for Computational Linguistics (Volume 1:
Long Papers), pages 826–839.

Michael Boutin. 1996. Bonggi clitics. In *Pan-Asiatic*
Linguistics: Proceedings of the fourth International
Symposium on languages and linguistics, volume 1.

George Aaron Broadwell. 2006. *A Choctaw reference*
grammar. Studies in the anthropology of North
American Indians. University of Nebraska Press, Lin-
coln.

Marta R Costa-Jussà, James Cross, Onur Çelebi, Maha
Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe
Kalbassi, Janice Lam, Daniel Licht, Jean Maillard,
and 1 others. 2022. No language left behind: Scaling
human-centered machine translation. *arXiv preprint*
arXiv:2207.04672.

Ryan Cotterell, Christo Kirov, John Sylak-Glassman,
David Yarowsky, Jason Eisner, and Mans Hulden.
2016a. [The SIGMORPHON 2016 shared Task—](#)
[Morphological reinflection](#). In *Proceedings of the*
14th SIGMORPHON Workshop on Computational
Research in Phonetics, Phonology, and Morphology,
pages 10–22, Berlin, Germany. Association for Com-
putational Linguistics.

Ryan Cotterell, Arun Kumar, and Hinrich Schütze.
2016b. [Morphological segmentation inside-out](#). In
Proceedings of the 2016 Conference on Empirical
Methods in Natural Language Processing, pages
2325–2330, Austin, Texas. Association for Computa-
tional Linguistics.

Aron Culotta and Andrew McCallum. 2005. [Reduc-](#)
[ing Labeling Effort for Structured Prediction Tasks:.](#)
In *Proceedings, The Twentieth National Conference*
on Artificial Intelligence and the Seventeenth Innova-
tive Applications of Artificial Intelligence Conference,

769	Arya D. McCarthy, Ekaterina Vylomova, Shijie Wu,	Miikka Silfverberg, Adam Wiemerslage, Ling Liu, and	825
770	Chaitanya Malaviya, Lawrence Wolf-Sonkin, Garrett	Lingshuang Jack Mao. 2017. Data Augmentation for	826
771	Nicolai, Christo Kirov, Miikka Silfverberg, Sabrina J.	Morphological Reinflection . In <i>Proceedings of the</i>	827
772	Mielke, Jeffrey Heinz, Ryan Cotterell, and Mans	<i>CoNLL SIGMORPHON 2017 Shared Task: Univer-</i>	828
773	Hulden. 2019. The SIGMORPHON 2019 shared	<i>sals Morphological Reinflection</i> , pages 90–99.	829
774	task: Morphological analysis in context and cross-		
775	lingual transfer for inflection . In <i>Proceedings of the</i>	Alaa Tharwat and Wolfram Schenck. 2023. A Sur-	830
776	<i>16th Workshop on Computational Research in Pho-</i>	vey on Active Learning: State-of-the-Art, Practical	831
777	<i>netics, Phonology, and Morphology</i> , pages 229–244,	Challenges and Research Directions . <i>Mathematics</i> ,	832
778	Florence, Italy. Association for Computational Lin-	11(4):820.	833
779	guistics.		
780	Sarah Moeller and Antti Arppe. 2024. Machine-in-the-	Inacio Vieira, Will Allred, Séamus Lankford, Sheila	834
781	loop with documentary and descriptive linguists . In	Castilho, and Andy Way. 2024. How Much Data is	835
782	<i>Proceedings of the Seventh Workshop on the Use of</i>	Enough Data? Fine-Tuning Large Language Models	836
783	<i>Computational Methods in the Study of Endangered</i>	for In-House Translation: Performance Evaluation	837
784	<i>Languages</i> , pages 27–32, St. Julians, Malta. Associa-	Across Multiple Dataset Sizes . In <i>Proceedings of</i>	838
785	tion for Computational Linguistics.	<i>the 16th Conference of the Association for Machine</i>	839
786	Sarah Moeller and Mans Hulden. 2021. Integrating au-	<i>Translation in the Americas (Volume 1: Research</i>	840
787	tomated segmentation and glossing into documentary	<i>Track)</i> , pages 236–249, Chicago, USA. Association	841
788	and descriptive linguistics . In <i>Proceedings of the</i>	for Machine Translation in the Americas.	842
789	<i>4th Workshop on the Use of Computational Meth-</i>	Wenhui Wang, Furu Wei, Li Dong, Hangbo Bao, Nan	843
790	<i>ods in the Study of Endangered Languages Volume 1</i>	Yang, and Ming Zhou. 2020. Minilm: Deep self-	844
791	<i>(Papers)</i> , pages 86–95, Online. Association for Com-	attention distillation for task-agnostic compression	845
792	putational Linguistics.	of pre-trained transformers . <i>Advances in neural in-</i>	846
793	Justin Mott, Ann Bies, Stephanie Strassel, Jordan Kod-	<i>formation processing systems</i> , 33:5776–5788.	847
794	ner, Caitlin Richter, Hongzhi Xu, and Mitchell Mar-		
795	cus. 2020. Morphological segmentation for low re-		
796	source languages . In <i>Proceedings of the Twelfth Lan-</i>		
797	<i>guage Resources and Evaluation Conference</i> , pages		
798	3996–4002, Marseille, France. European Language		
799	Resources Association.		
800	Åshild Næss and Brenda H. Boerger. 2008. Reefs–Santa		
801	Cruz as Oceanic: Evidence from the Verb Complex .		
802	<i>Oceanic Linguistics</i> , 47(1):185–212.		
803	Shu Okabe and François Yvon. 2023. Towards multi-		
804	lingual interlinear morphological glossing . In <i>Find-</i>		
805	<i>ings of the Association for Computational Linguis-</i>		
806	<i>tics: EMNLP 2023</i> , pages 5958–5971, Singapore.		
807	Association for Computational Linguistics.		
808	Alexis Mary Palmer. 2009. <i>Semi-automated annota-</i>		
809	<i>tion and active learning for language documentation</i> .		
810	PhD Thesis, University of Texas at Austin.		
811	Laura Pérez-Mayos, Miguel Ballesteros, and Leo Wan-		
812	ner. 2021. How much pretraining data do language		
813	models need to learn syntax? In <i>Proceedings of the</i>		
814	<i>2021 Conference on Empirical Methods in Natural</i>		
815	<i>Language Processing</i> , pages 1571–1582, Online and		
816	Punta Cana, Dominican Republic. Association for		
817	Computational Linguistics.		
818	Frank Seifart, Nicholas Evans, Harald Hammarström,		
819	and Stephen C Levinson. 2018. Language document-		
820	ation twenty-five years on. <i>Language</i> , 94(4):e324–		
821	e345.		
822	Burr Settles. 2009. Active learning literature survey.		
823	<i>Tech. Rep. Computer Sciences Technical Report 1648,</i>		
824	<i>University of Wisconsin-Madison.</i>		

848

A Minimal Optimal Sizes

849

(see next page)

Table 2: The minimal optimal sizes for all four tasks with the Transformer; the minimal optimal size for each initial training size is the average of the knees (point of maximum curvature of plotted f1 scores) across all three selection sizes.

Language	ISO 639-3	Task	Initial Size	Optimal Size
Alas-Kluet	btz	morpheme segmentation	500	1666
Alas-Kluet	btz	morpheme segmentation	1000	1833
Alas-Kluet	btz	morpheme segmentation	1500	2333
Alas-Kluet	btz	morpheme segmentation	2000	2400
Alas-Kluet	btz	morpheme segmentation	2500	2700
Alas-Kluet	btz	morpheme glossing	500	1450
Alas-Kluet	btz	morpheme glossing	1000	1750
Alas-Kluet	btz	morpheme glossing	1500	2000
Alas-Kluet	btz	morpheme glossing	2000	2383
Alas-Kluet	btz	morpheme glossing	2500	2633
Alas-Kluet	btz	joint segmentation & glossing	500	1566
Alas-Kluet	btz	joint segmentation & glossing	1000	1866
Alas-Kluet	btz	joint segmentation & glossing	1500	2350
Alas-Kluet	btz	joint segmentation & glossing	2000	2583
Alas-Kluet	btz	joint segmentation & glossing	2500	2750
Alas-Kluet	btz	inflection	500	1500
Alas-Kluet	btz	inflection	1000	1666
Alas-Kluet	btz	inflection	1500	2100
Alas-Kluet	btz	inflection	2000	2500
Alas-Kluet	btz	inflection	2500	2750
Bonggi	bdg	morpheme segmentation	500	6583
Bonggi	bdg	morpheme segmentation	1000	6283
Bonggi	bdg	morpheme segmentation	1500	6300
Bonggi	bdg	morpheme segmentation	2000	6033
Bonggi	bdg	morpheme segmentation	2500	6200
Bonggi	bdg	morpheme glossing	500	3316
Bonggi	bdg	morpheme glossing	1000	3383
Bonggi	bdg	morpheme glossing	1500	3450
Bonggi	bdg	morpheme glossing	2000	3966
Bonggi	bdg	morpheme glossing	2500	4583
Bonggi	bdg	inflection	500	4116
Bonggi	bdg	inflection	1000	3200
Bonggi	bdg	inflection	1500	4083
Bonggi	bdg	inflection	2000	6000

Continued on next page

Table 2: The minimal optimal sizes for all four tasks with the Transformer; the minimal optimal size for each initial training size is the average of the knees (point of maximum curvature of plotted f1 scores) across all three selection sizes. (Continued)

Language	ISO 639-3	Task	Initial Size	Optimal Size
Bonggi	bdg	inflection	2500	6050
Bonggi	bdg	joint segmentation & glossing	500	6583
Bonggi	bdg	joint segmentation & glossing	1000	6316
Bonggi	bdg	joint segmentation & glossing	1500	7016
Bonggi	bdg	joint segmentation & glossing	2000	6366
Bonggi	bdg	joint segmentation & glossing	2500	5750
Choctaw	cho	morpheme segmentation	500	2816
Choctaw	cho	morpheme segmentation	1000	3233
Choctaw	cho	morpheme segmentation	1500	3766
Choctaw	cho	morpheme segmentation	2000	3833
Choctaw	cho	morpheme segmentation	2500	4550
Choctaw	cho	morpheme glossing	500	1016
Choctaw	cho	morpheme glossing	1000	1333
Choctaw	cho	morpheme glossing	1500	1783
Choctaw	cho	morpheme glossing	2000	2116
Choctaw	cho	morpheme glossing	2500	2716
Choctaw	cho	joint segmentation & glossing	500	1383
Choctaw	cho	joint segmentation & glossing	1000	1450
Choctaw	cho	joint segmentation & glossing	1500	2116
Choctaw	cho	joint segmentation & glossing	2000	2333
Choctaw	cho	joint segmentation & glossing	2500	2716
Choctaw	cho	inflection	500	966
Choctaw	cho	inflection	1000	1383
Choctaw	cho	inflection	1500	2216
Choctaw	cho	inflection	2000	2316
Choctaw	cho	inflection	2500	2750
Lezgi	lez	morpheme segmentation	500	5483
Lezgi	lez	morpheme segmentation	1000	5583
Lezgi	lez	morpheme segmentation	1500	5250
Lezgi	lez	morpheme segmentation	2000	5583
Lezgi	lez	morpheme segmentation	2500	5950
Lezgi	lez	morpheme glossing	500	2700
Lezgi	lez	morpheme glossing	1000	2966
Lezgi	lez	morpheme glossing	1500	3000

Continued on next page

Table 2: The minimal optimal sizes for all four tasks with the Transformer; the minimal optimal size for each initial training size is the average of the knees (point of maximum curvature of plotted f1 scores) across all three selection sizes. (Continued)

Language	ISO 639-3	Task	Initial Size	Optimal Size
Lezgi	lez	morpheme glossing	2000	3316
Lezgi	lez	morpheme glossing	2500	3783
Lezgi	lez	joint segmentation & glossing	500	4416
Lezgi	lez	joint segmentation & glossing	1000	4450
Lezgi	lez	joint segmentation & glossing	1500	4416
Lezgi	lez	joint segmentation & glossing	2000	4600
Lezgi	lez	joint segmentation & glossing	2500	4483
Lezgi	lez	inflection	500	3300
Lezgi	lez	inflection	1000	2566
Lezgi	lez	inflection	1500	2850
Lezgi	lez	inflection	2000	4016
Lezgi	lez	inflection	2500	4033
Natügu	ntu	morpheme segmentation	500	5700
Natügu	ntu	morpheme segmentation	1000	5950
Natügu	ntu	morpheme segmentation	1500	5566
Natügu	ntu	morpheme segmentation	2000	6000
Natügu	ntu	morpheme segmentation	2500	6083
Natügu	ntu	morpheme glossing	500	2516
Natügu	ntu	morpheme glossing	1000	3350
Natügu	ntu	morpheme glossing	1500	3633
Natügu	ntu	morpheme glossing	2000	3933
Natügu	ntu	morpheme glossing	2500	4683
Natügu	ntu	joint segmentation & glossing	500	5233
Natügu	ntu	joint segmentation & glossing	1000	5300
Natügu	ntu	joint segmentation & glossing	1500	4683
Natügu	ntu	joint segmentation & glossing	2000	5050
Natügu	ntu	joint segmentation & glossing	2500	5383
Natügu	ntu	inflection	500	3483
Natügu	ntu	inflection	1000	4533
Natügu	ntu	inflection	1500	5200
Natügu	ntu	inflection	2000	5283
Natügu	ntu	inflection	2500	5433
Upper Tanana	tau	morpheme segmentation	500	5716
Upper Tanana	tau	morpheme segmentation	1000	5616

Continued on next page

Table 2: The minimal optimal sizes for all four tasks with the Transformer; the minimal optimal size for each initial training size is the average of the knees (point of maximum curvature of plotted f1 scores) across all three selection sizes. (Continued)

Language	ISO 639-3	Task	Initial Size	Optimal Size
Upper Tanana	tau	morpheme segmentation	1500	5933
Upper Tanana	tau	morpheme segmentation	2000	5883
Upper Tanana	tau	morpheme segmentation	2500	5400
Upper Tanana	tau	morpheme glossing	500	3116
Upper Tanana	tau	morpheme glossing	1000	3066
Upper Tanana	tau	morpheme glossing	1500	3400
Upper Tanana	tau	morpheme glossing	2000	3383
Upper Tanana	tau	morpheme glossing	2500	4050
Upper Tanana	tau	joint segmentation & glossing	500	4083
Upper Tanana	tau	joint segmentation & glossing	1000	4616
Upper Tanana	tau	joint segmentation & glossing	1500	4066
Upper Tanana	tau	joint segmentation & glossing	2000	4483
Upper Tanana	tau	joint segmentation & glossing	2500	4766
Upper Tanana	tau	inflection	500	3600
Upper Tanana	tau	inflection	1000	2883
Upper Tanana	tau	inflection	1500	3266
Upper Tanana	tau	inflection	2000	3983
Upper Tanana	tau	inflection	2500	4633

855 **B LLM Prompting**

856 The llama3 system prompt is provided as follows,
857 with curly braces representing where the language,
858 task, and input/output pairs are specified.

859 You are a linguist specializing in the language
860 {language}. You are given words and tasked with
861 {task}. Here is a training set of {N} word inputs
862 (left) and their expected {task} outputs (right):

863 {input A} {output A}

864 {input B} {output B}

865 ...

866 Now, I will give you a list of words. For each
867 word, provide only the {task} output on a separate
868 line, in the same order as the input words. Do not
869 include the input words in your response. Output
870 one result only. It is absolutely critical that you
871 only return the output, no other text or formatting.