

Findings from the Bambara - French Machine Translation Competition

Anonymous EACL submission

Abstract

Anonymous Organization (AO) hosted a low-resource machine translation (MT) competition with monetary prizes. The goals of the project were to raise awareness of the challenges in the low-resource MT domain, improve MT algorithms and data strategies, and support MT expertise development in the regions where people speak Bambara and other low-resource languages. The participants built Bambara to French and French to Bambara machine translation systems using data provided by AO as well as additional data resources that were shared amongst the competitors. This paper details each team's different approaches as well as motivation for ongoing work in both Bambara and the broader low-resource machine translation domain.

1 Competition Introduction

Anonymous Organization, hosted a low-resource machine translation competition that ended on February 15, 2023. The competition launched December 15, 2022 with interested participants being given access to a Github repository with a training set of parallel French-Bambara aligned sentences. The participants were also invited into a Slack community to share their approaches and data. Later, an additional development set was provided to the teams and fewer than 48 hours before the submission deadline, a test set was released for generating text output to be sent to the competition organizers to evaluate translation accuracy using BLEU scores (Post, 2018).

The goals of the competition were to improve not only French to Bambara and Bambara to French automated translation systems but support a transparent and collaborative community to work on these and other language pairs, especially those low-resource languages spoken by West Africans. 50 people joined the online community and fourteen people competed in 6 teams. The teams were

made of participants from Mali, Senegal, Namibia, Nigeria, Ireland, Germany, Russia, Spain, France, the US, and the UK. Many of the participants speak or have working knowledge of a "low-resource language" or a language that does not have the digital resources that support highly accurate Natural Language Processing tool development.

2 Related Work

Current state-of-the-art low-resource MT is surveyed in (Haddow et al., 2022). Google Translate has implemented more low-resource languages into their language library sharing innovations in blog-posts: ¹ and ².

MT for the Bambara - French language pair has been explored in recent years by (Tapo et al., 2020) and (Leventhal et al., 2020). This work is in part motivated by an increased financial and cultural focus on bringing machine learning to the Sahel region (Diarra and Leventhal, 2020).

2.1 Evaluation

MT can be evaluated by automated and manual methods. In this competition, we used automated tools to evaluate the closeness of translations to a gold standard. We used BLEU scores with sacreBLEU (Post, 2018) and (Papineni et al., 2002).

2.2 Datasets

The organizers provided a training dataset of aligned parallel Bambara - French sentences from the medical and dictionary domains. Each line is one sentence.

The dataset composition is as follows; Training data = 3150 lines, Dev data = 460 lines, and Test data = 460 lines. In addition to the competition data, all participants were encouraged to gather, utilize, and share additional resources with other

¹<https://www.teachyoubackwards.com/>

²<https://translate.googleblog.com/2010/05/five-more-languages-on.html>

members of the competition community. These resources are shown in Table 1.

Dataset	Teams
MAFAND (Adelani et al., 2022)	All Teams
NLLB-SEED (Team et al., 2022)	All Teams
FLORES (Goyal et al., 2022)	All Teams
BAYELEMABAGA (Vydrin et al., 2022)	All Teams
XP3 (Muennighoff et al., 2022a)	D
Wikipedia (Wikimedia, 2023)	A

Table 1: Additional Bambara datasets used by the different teams

2.3 Machine Learning

Technique	Reference
BART	(Lewis et al., 2019)
BLOOM-z 560M, mt0-small	(Muennighoff et al., 2022b)
byt5	(Xue et al., 2021a)
DeltaLM	(Ma et al., 2021)
HuggingFace	(Wolf et al., 2020)
LION optimizer	(Chen et al., 2023)
LoRA	(Hu et al., 2021)
M2M100 model	(Fan et al., 2020)
MarianNMT/Opus-MT	(Junczys-Dowmunt et al., 2018)
mt5	(Xue et al., 2021b)
NLLB model	Team et al., 2022
PEFT library	(Mangrulkar et al., 2022)
Sockeye	(Hieber et al., 2020)

Table 2: Techniques and models used by the different teams

Table 2 shows the different techniques and models used by the teams with transformer (Vaswani et al., 2017) and BERT models (Mishra et al., 2022; Sheshadri et al., 2023) inspiring much of the development.

3 The Machine Translation Systems

Six teams submitted system output that could be evaluated using sacreBLEU. An additional team built an MT system but did not submit output. They were not scored and did not place but their insights from training and error analysis are included in this paper. In the following sections, each team first describes their methodology and model, and then discusses their error analysis. See Table 1 for the datasets used by each team.

3.1 Team A

3.1.1 Models

We used an additional dataset from Wikipedia (Wikimedia, 2023) which provided us with an extra 892 lines of data. Next, we made a list of

MT models that contained Bambara and French in their dataset during pre-training. As a result, we started with the NLLB-200 (Team et al., 2022) pre-trained model. We fine-tuned both the 600M and the 1.3B (in order to test the impact of scaling on model capacity) parameter versions, from the Huggingface Hub. We found the NLLB model to be under-performing. Next, we switched to an M2M-100 (Fan et al., 2020) model after we discovered it had fine-tuned multilingual MT models separately for each language direction, which outperformed NLLB-200 (Adelani et al., 2022)

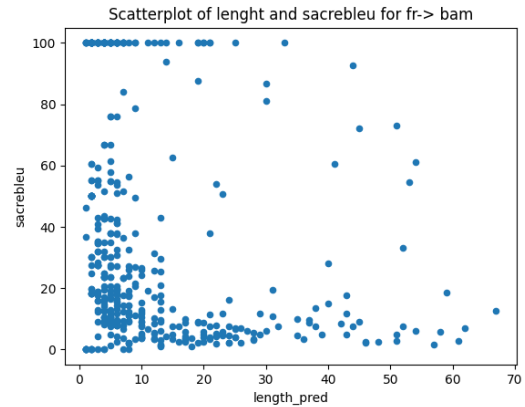


Figure 1: Scatterplot showing length of predicted sentences against sentence BLEU scores.

To gain further insight into the challenges posed by certain sentence characteristics in our MT model, we conducted an analysis of the per-sentence BLEU scores plotted against the length of the predicted sentences. Initially, we postulated that our MT model would perform better with shorter sentences and perform worse with longer sentences. However, as illustrated in Figure 1, which presents a scatterplot of the lengths of the predicted sentence against their sentence BLEU scores, our model struggled even with shorter sentences. This led us to reconsider our hypothesis and explore the possibility that our model was underfitting. Next, we decided to investigate the potential benefits of implementing backtranslation.

3.1.2 Team A’s Backtranslation Approach

Several papers have highlighted the positive effect of backtranslation (Sennrich et al., 2016a; Ponceles et al., 2018; Zhang et al., 2020; Dossou and Emezue, 2020; Fan et al., 2020; Emezue and Dossou, 2021; Adelani et al., 2022; Team et al., 2022). Inspired by random online backtranslation (Zhang

Algorithm 1 Team A’s Backtranslation Approach

$n_epochs \leftarrow$ number of fine-tuning epochs
 $D_{train} \leftarrow$ training dataset of French- Bambara parallel sentences
 $D_{bam}^{wiki} \leftarrow$ 892 monolingual cleaned sentences from Wikipedia.
 $D_{fr} \leftarrow$ dataset of French sentences only. For our case it was gathered by taking the French instances of D_{train}
 $D_{bam} \leftarrow$ dataset of Bambara sentences only. For our case it was gathered by taking the Bambara instances of D_{train} and additional monolingual sentences from D_{bam}^{wiki}
 $M_{fr \rightarrow bam}^0 \leftarrow$ fine-tuned MT model of (Adelani et al., 2022) for French \rightarrow Bambara
 $M_{bam \rightarrow fr}^0 \leftarrow$ fine-tuned MT model of (Adelani et al., 2022) for Bambara \rightarrow French.
 $D_{train}^0 \leftarrow D_{train}$.
for $k \leftarrow [0, 1, 2 \dots n]$ **do**
 $M_{fr \rightarrow bam}^{k+1} \leftarrow$ fine-tune $M_{fr \rightarrow bam}^k$ on D_{train}^k for n_epochs epochs.
 $M_{bam \rightarrow fr}^{k+1} \leftarrow$ fine-tune $M_{bam \rightarrow fr}^k$ on D_{train}^k for n_epochs epochs.
 $D_{bam}^k \leftarrow$ generated synthetic translations to Bambara from D_{fr} using $M_{fr \rightarrow bam}^{k+1}$.
 $D_{fr}^k \leftarrow$ generated synthetic translations to French from D_{bam} using $M_{bam \rightarrow fr}^{k+1}$.
 $D_{train}^{k+1} \leftarrow$ concatenated training dataset gotten from $D_{train}^0 \cup \{D_{bam}^k \leftrightarrow D_{fr}\} \cup \{D_{fr}^k \leftrightarrow D_{bam}\}$
end for

et al., 2020), we created our version, explained in Algorithm 0, to help our model better utilize the training dataset, and the 892 monolingual Bambara sentences from Wikipedia. Our approach, dubbed "cyclic translation" (Lam et al., 2021), would theoretically enable the model to leverage the available training and monolingual dataset by compelling the MT model for each direction, at each step k , to learn from a concatenation of the original training dataset, its synthetically generated sentences, and those generated by the MT model of the opposite

direction in the previous step.

Despite its potential benefits, implementing backtranslation presented several challenges. First, it was a difficult process to set up, particularly in achieving a high degree of automation and reducing the need for human intervention. Secondly, it was computationally expensive and time-consuming, as each iteration of the backtranslation process involved working with three times more data than the previous iteration. Consequently, we were only able to complete one backtranslation successfully.

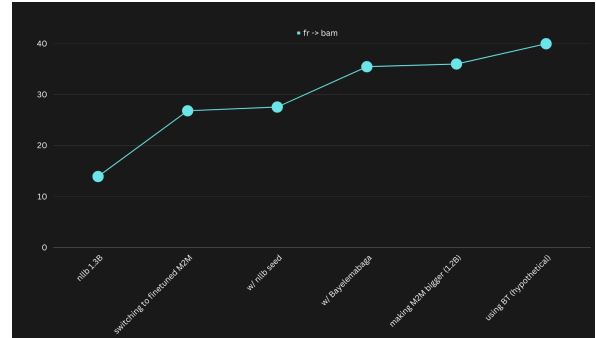


Figure 2: Timeline of Team A efforts and BLEU score on dev set. The chart begins with our use of NLLB, switches to fine-tuned M2M, incorporates NLLB Seed dataset, then includes the BAYELEMABAGA dataset, and ends in our hypothetical performance using our cyclic backtranslation approach.

We included a potential impact in Figure 2 which shows the timeline of our activities and their corresponding evaluation results on the French \rightarrow Bambara direction.

One of the major challenges facing machine translation for African languages is the limited availability of high-quality datasets (V et al., 2020; Caswell et al., 2021; Adelani et al., 2022). This became apparent in our study, where the use of the BAYELEMABAGA dataset resulted in a significant increase in the performance of our MT model. The scarcity of such resources highlights the need for continued efforts to develop and curate datasets for African languages, which could significantly improve the performance of machine translation models for African languages.

3.2 Team B

3.2.1 Model

We used a pre-trained MarianMT transformer model (Junczys-Dowmunt et al., 2018) which was pretrained for Romance languages to English due to the non-existence of Bam-Fr pre-trained weights

for the MarianMT model. The model was then trained using a set of hyperparameters which were inspired by findings from (Araabi and Monz, 2020) and (Van Biljon et al., 2020) where the authors found the hyperparameters that would achieve the highest BLEU scores when dealing with low-resource languages. Although our implementation was limited due to insufficient computing power (i.e., we were unable to increase attention heads without the GPU crashing during training).

We use the following set of hyperparameters; optimizer: adam, learning rate: $2e^{-5}$, beta 1:0.9, beta 2: 0.999, epsilon: $1e^7$, batch size: 64, and attention heads: 8.

3.2.2 Error analysis

Due to limited computing power, we were not able to fully train our MT model until convergence. In is plausible our model could have achieved higher accuracy or lower bias with more iterations of gradient descent. We also were not able to fine-tune our hyper-parameters as much as we would have liked.

In the seq2seq translation output, one word would get repeated multiple times back-to-back. This hallucination could be reduced by using a model that was pre-trained in French, so it would know from experience that French sentences do not normally include back-to-back repeated words.

There were words that appear infrequently in the training set and were mistranslated very frequently. With more time in this competition, this could have been alleviated with Byte Pair Encoding (BPE).

3.2.3 Discussion

While the existing literature suggests that Transformer models typically need a large training corpus to do well, our model suggests otherwise. With minor modifications made to the transformer out-of-the-box, the Transformer seq2seq model was still able to achieve a BLEU of 14.81 despite a limited training corpus, lack of pre-trained Bambara experience, computing power, and hyper-parameter tuning. In hindsight, we should have used a model that was pre-trained for Bambara to any Romance language, because it would be easier to learn Bambara to French if it had been pre-trained in Bambara to English, for example. We believe that Bambara is so different from anything that the model had been pre-trained on that it had a hard time learning a very different language with such a small dataset.

3.3 Team C

3.3.1 Model

Our team has previously worked on MT tasks on languages such as French, Reunionese Creole, Portuguese, Umbundu, and Kimbundu, where we observed sub-optimal outcomes when training an autoregressive generative Transformer model, either encoder-only or decoder-only, starting from scratch. Hence, for the given task, we wanted to use a Sequence to Sequence model with prior training on the Bambara language. We evaluated the following models on the development datasets; mt0-small (Muennighoff et al., 2022b), BLOOM-z 560M (Muennighoff et al., 2022b), NLLB 600M distilled Team et al., 2022, NLLB 1.3B Team et al., 2022, NLLB 1.3B distilled Team et al., 2022, and NLLB 3.3B Team et al., 2022.

Upon evaluating the dev dataset, NLLB 600M distilled and NLLB 1.3B distilled exhibited superior performance. However, due to computational limitations even with our optimizations, training the NLLB 3B version would have been impossible. For an auto-regressive/instruction model, BLOOM-z exhibited more potential than mt0-small, and after two epochs, it produced acceptable scores. Nevertheless, it appears that general-purpose models of such small sizes do not rival specialized Sequence to Sequence models of similar dimensions, especially in a "low-resource" scenario.

We chose to concentrate our scarce GPU hours to the most two promising models, NLLB 600M and NLLB 1.3B, both distilled models, and fine-tune them until the competition deadline. This provided an avenue to utilize and fine-tune distilled models. 1.3B distilled was better than not distilled models. Without fine-tuning, by using the default HuggingFace *generate* method, the 600M distilled model had a BLEU score of 19.8157 and 17.9217 for BAM to FR and FR to BAM, respectively. And the non fine-tuned distilled 1.3B model had 24.5496 and 25.5610 for BAM to FR and FR to BAM, respectively. Both were tested on the dev corpus provided by the challenge organizers. Please see Table 3 for the BLEU scores using different models and steps.

The hyperparameters used for fine-tuning the NLLB Models are Optimizer: Adafactor, Learning rate: $1e^{-04}$, Batch size (1.3B model): 4, Batch size (600M model): 10, Gradient acc. (1.3B model): 16, and Gradient acc. (600M model): 10.

Model/Steps	BAM → FR	FR → BAM
600M/3000 steps	21.7641	18.8674
600M/6000 steps	21.5270	21.3773
600M/9000 steps	21.3773	17.8374
1.3B/1500 steps	20.3349	17.8032
1.3B/3000 steps	18.6542	17.6243
1.3B/4500 steps	24.2556	19.3324
1.3B/6000 steps	25.3816	18.7743
1.3B/7500 steps	26.0991	18.1205

Table 3: BLEU Scores on development set (Team C)

3.3.2 Error Analysis

We made a challenging discovery during this competition. In the NLLB paper, the src and tgt sequences are fed to the model with this scheme: (src_sequence, src_lang) for the source sequence and (tgt_lang, tgt_sequence) for the target sequence. On the other hand, the NLLB tokenizer in the HuggingFace transformer tokenizes the pair of sequences as (src_sequence, src_lang) and (tgt_sequence, tgt_lang). Once we fixed this issue, the sacreBLEU scores of our finetuned NLLB models started to improve, consistently with the decrease of the loss, and with the quality differences that we could observe. However, we discovered and fixed this issue less than 24 hours before the deadline, and we had lost quite a bit of time by trying other fixes. As French is our native tongue, and a member of our group has some understanding of Bambara, we were able to compare the outputs of the model, to the targets of the development set, and before this discovery, the BLEU scores of our finetuned models were underwhelming and inconsistent with the steadily decreasing loss on the dev set, and our observations of the outputs. After this fix, the BLEU scores started to make more sense, even if we did not resolve the difference in behaviour between the two translation directions. Indeed, the French to Bambara was getting marginally better in terms of BLEU scores, while Bambara to French was dramatically worse than the base performance.

3.3.3 Discussion

For our next MT project, we are seeking access to a large language model-based system. We believe it would be a good idea to try the performance of few-shot prompting on these LLMs, because we have seen that the most promising model is still very limited for languages like Bambara.

Since Bambara, like many languages, is primary spoken, we will try speech-based approaches in future work. These approaches will potentially have more impact and be more useful to these communi-

ties, most of whom do not write their languages.

3.4 Team D

3.4.1 Models

For pre-trained models, we explored several models available on the HuggingFace Hub, including M2M-100 (Fan et al., 2020), NLLB (Team et al., 2022), mT5 (Xue et al., 2021b) and byt5 (Xue et al., 2021a) models each pre-trained by the Masakhane Organization (V et al., 2020). Each model was evaluated on the dev-set provided by the organizers with respect to the BLEU score. The M2M-100 (Fan et al., 2020) was chosen as a starting point since it scored the highest. It is a 483 million parameters distilled version of the original 1.2 billion parameters encoder-decoder transformer.

Fine-tuning on the challenge dataset was promising but the model validation loss curves showed overfitting despite fine-tuning for weight decay, small learning rate with decreasing linear schedule, warmup, and dropout. In addition, the BLEU score would not exceed 15 on the dev set, but upon manual investigation, the produced translations were shallow and sometimes semantically unrelated to the ground truth.

3.4.2 Error analysis

We examined the generated translations for common issues such as mis-translations, omissions, and word order errors. The resulting training process consisted of two steps: fine-tuning on the extended dataset and a step involving the challenge data. Yielding a BLEU score of **27** on the dev set, this approach produced a better result than fine-tuning on a mix of both extended and challenge data. The challenge data would then be under-represented which would allow for a low BLEU score since the model is evaluated on a dev-set from the challenge data distribution and not the extended one.

That score was further improved by changing the generation algorithm type and number of beams, resulting in the final dev BLEU score of **28.93** seen in Figure 3. This improved the score by 2 points.

Error analysis showed the gap in BLEU score between the dev set medical data and dictionary data. An average of 10 points of difference was reported from one distribution to the other, which could be explained by two main differences: that in sequence length (the dictionary data was notably shorter) and in vocabulary distribution (the medical data was more domain-specific).

3.4.3 Discussion

In addition to the data in Table 1, we extended our training data by processing a many-to-Bambara dataset from BigScience: the Bambara split of XP3-all (Muennighoff et al., 2022a). XP3-all contains 265,180 many-to-Bambara lines, but we only included the French-to-Bambara subset, and enriched it with the English-to-Bambara subset that was translated with the opus-mt-en-fr model from Helsinki-NLP (Tiedemann and Thottingal, 2020) resulting in 8,377 additional lines of training data.

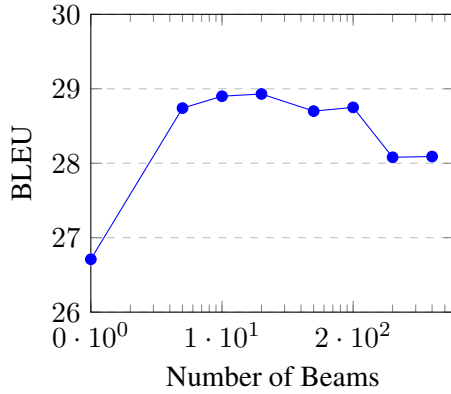


Figure 3: BLEU as a function of the number of beams. Note that a value of one necessarily implies greedy decoding while bigger values correspond to the beam-search algorithm. Not surprisingly, the score dramatically improves before plateauing around 10 and reaching diminishing returns. Notably, the optimum is reached at 15 and increasing the number of beams further has a negative impact on the score.

In the future, we would spend more time automating tasks, including hyper-parameter tuning, to improve the efficiency of the system. Notably, the cross-entropy loss function is only a differentiable proxy for the metric we are trying to optimize: the BLEU score (which is not differentiable). With the recent success of Reinforcement Learning Techniques in neural language generation tasks (Stiennon et al., 2020), we plan to further fine-tune the model using the BLEU metric as a task reward, similar to Pinto et al., 2023.

In the future we will explore techniques, such as the recently introduced PEFT (Mangrulkar et al., 2022) library, which allows for fine-tuning of LLMs on very small datasets using parameter efficient fine-tuning methods. IA3 (Liu et al., 2022), Prompt-Tuning (Lester et al., 2021), Prefix-Tuning (Liu et al., 2021), and Low Rank Adaptation (Hu et al., 2021) methods are currently leveraged to train large models efficiently on as few as 10 ex-

amples. In comparison to classic fine-tuning that involves training all the weights of the model, these methods have the added advantage of achieving similar (sometimes even better) results by training only a small subset of the weights (by freezing the pre-trained weights and adding trainable adapter weights in the case of LoRA and IA3). We therefore expect these methods to be increasingly used for any low-resource task in the near future.

Moreover, it seems that the Adam optimizer has finally found a worthy, artificially evolved (Chen et al., 2023) opponent. We look forward to test it using the parameters of this task.

Finally, we would suggest the use of learned metrics for the evaluation of the translations instead of the BLEU metric - that ignores synonyms and idioms - building on the works of Zhang et al., 2019. Of course such models are not yet trained on Bambara but Eddine et al., 2021 seems to offer at least part of the solution, and an alternative would simply be computing the cosine-distance between the embedding representation of the produced translation and that of the reference thanks to Reimers and Gurevych, 2020.

3.5 Team E

3.5.1 Models

All of our models were trained using Sockeye (Hieber et al., 2020). In this task we decided to concentrate on models "from scratch" and utilized 4 checkpoints averaging model parameters in our system. We averaged the parameters of the best 4 checkpoints, which helped to improve results. In addition we used BPE for word segmentation (Sennrich et al., 2016b).

3.5.2 Error analysis

We performed error analysis based on the BLEU metric, and used it as an optimized metric while training. We also used the sacreBLEU tool (Post, 2018).

3.5.3 Discussion

There are other, extended techniques such as back translation and pre-trained models that we anticipate using for future research. In addition, we also plan to add additional training datasets that were provided and used by the other teams.

3.6 Team F

3.6.1 Models

We looked at several approaches concurrently, first pre-training a bilingual Bambara to French Seq-2-Seq-based foundational model with a lower quality dataset, inspired by the BART (Lewis et al., 2019) technique, then fine-tuning it with a higher-quality dataset. This approach yielded non-optimal translations and performance, with all the scores being sub-8 BLEU (it was also resource-heavy, and time-consuming). We worked on a model fine-tuned with the DeltaLM (Ma et al., 2021) pre-trained model, the base, and the large versions, although the training was never completed, crashing with various adjustments primarily linked to computing resource problems.

We were able to double our performance when we retrained with the NLLB200 (Team et al., 2022) 600M parameters pre-trained model, with a learning rate of 2, batch size of 512, and training steps of 20k with the lower-quality dataset, with and without DABA-assisted pre-processing³.

We obtained another peak in performance when we unfroze the model and then tuned it with the competition dataset with the same configuration for a downstream understanding of the type of text used for the competition (although we are suspecting over-fitting here). We have seen similar results from both directions Bambara to French and French to Bambara.

3.6.2 Error analysis

We knew that Bambara is a complex and morphologically sophisticated language. Bambara and French have a one sentence to many translation scheme, where one sentence can have multiple interpretations in the other language. Additionally, with Bambara being predominantly a spoken language there are many fluidities that only native speakers can pick up from translations, compared to a more structured language. We chose to weigh human evaluation higher than automated metrics. Both evaluation techniques coupled gave an insight into the overall performance of our models. We came up with our own defined method for manual evaluation, described as follows: For every model trained, we sampled 50 lines from our test set and classified each line into three classes manually BAD, ACCEPTABLE, and GOOD. Where BAD was given a value of 0; it is chosen when the

³<https://github.com/maslinych/daba>

Team Member 1		
Model	BAG	BLEU
FR → BAM	35	8.86
BAM → FR	41	11.26
Team Member 2		
Model	BAG	BLEU
FR → BAM	63	10.05
BAM → FR	84	16.12
Team Member 3		
Model	BAG	BLEU
FR → BAM	75	13.74
BAM → FR	29	9.74
Team Member 4		
Model	BAG	BLEU
FR → BAM	57	15.86
BAM → FR	62	13.10

Table 4: Results for each team member

hypothesis does not relay any information from the source or is a bad translation. ACCEPTABLE was given a value of 1; it is chosen when the hypothesis is a literal translation of the source without context. GOOD was given a value of 2; it is chosen when the hypothesis is an accurate translation of the source with context.

Each member of the team evaluated a batch of 50 lines per model, given the source text, a reference translation, and the hypothesis generated by the model. They were tasked to evaluate the manual score and to compute the BLEU score of the batch, for a comparative analysis of the two results.

We acknowledge the subjective nature of human evaluation, therefore should state that this was used to guide our analysis of the performance of our models for the competition, further investigations are needed to validate its viability.

3.6.3 Discussion

Bambara’s complexity made it challenging to find the best possible approach, as each aspect of the training required analysis. From pre-processing to evaluation, we found that fine-tuning with the NLLB200 600M model to be more performant. The most significant aspect in our method was the human-in-the loop approach, where coupling human annotation and automated metrics was the primary indicator that informed our decisions during the competition.

3.7 Team G

3.7.1 Model

We experimented with straightforward transformer-based models and utilized the attention mechanism,

which enables one component of the model to concentrate on another part of the model, it was chosen. Due to the issue of vanishing gradient and the weakness of limited levels of parallelization, respectively, both recurrent neural networks (RNNs) and Long Short Term Memory (LSTM) were not considered (Vaswani et al., 2017). The selected transformer model was Facebook/nllb-200-distilled-600M (Team et al., 2022), which was fine-tuned on the training dataset, which allowed for the design of the encoder, latent representation, and decoder. By using semi-supervised learning, the decoder fed features to the model, and to ensure convergence, the team employed 100 epochs.

3.7.2 Error Analysis

By using Google Translation, the team was able to avoid having a native speaker as a teammate. In the future, a native speaker will be a part of the team.

3.7.3 Discussion

Beyond needing additional compute and a powerful internet connection, we would like to consider other alternative models for cross-validation.

Team Name	BLEU Score (BAM to FR)	BLEU Score (FR to BAM)
Team A	16.31	17.45
Team C	13.12	11.1
Team D	19.05	N/A
Team E	7.54	8.06
Team B	14.81	N/A
Team F	5.82	N/A

Table 5: BLEU score results by team for Bambara - French and French - Bambara, with placement ordering.

4 Results and Discussion

Table 5 shows the BLEU scores for both Bambara to French and French to Bambara translations. Not all of the teams attempted both translation directions and the scores were averaged across both language pairs to determine the winners.

This MT competition aimed to increase research in low-resource language machine translation by providing training and evaluation data and supporting community-building around scientific transparency. Community-building included teams being constructed from individuals with complementary skills and all relevant training data discovered by one team being shared amongst the teams.

Nonetheless, there were key themes to the submissions. All of the teams used the same core data sets with two teams bootstrapping alternatives as

shown in Table 1. Additional data provided a significant advantage in this low-resource situation. From a machine learning perspective, many of the teams shared similar approaches with effectively utilizing the M2M-100 model (Fan et al., 2020) as the differentiator between the top performing teams. Interestingly, the NLLB-200 (Adelani et al., 2022) model comparatively under-performed. We believe this is because the M2M-100 model had fine-tuned MT models separately for each language direction.

Subsequent insights were that the winning team used a backtranslation approach called "cyclic translation" and another successful team used a beam search optimization. Also, we learned that smaller distilled models could beat larger models with limited amounts of data (i.e., fine-tuning distilled models yields more accurate results).

Only two of the teams had members that spoke Bambara but many participants are speakers of other low-resource languages and hope to extend their experience with MT system development to languages that their families and friends speak.

5 Conclusion and Future Work

Because of this competition, researchers have successfully implemented innovative low-resource machine translation systems. These implementations are extensible to other language pairs, which is helpful since low-resource languages continue to face numerous challenges in terms of research focus and funding. We believe this competition has not only supported increased visibility of the Bambara language, but it has also showcased the talent that is working on using creative techniques to address these technical challenges globally.

We would like to extend this work by holding another competition. Ideally, the next competition will utilize low-resource language automatic speech recognition data. Finally, we would like to provide greater financial support to the participating teams by sponsoring free access to computational resources. This even playing field could better illuminate which machine learning models are the highest performers.

Limitations

We aim to address these limitations in future work.

1. There are known limitations of BLEU for meaningful evaluation including with how

well BLEU corresponds to human evaluation of language correctness and naturalness. In the future we would like to conduct human evaluation of the MT competition output.

2. The importance of compute power was also evident in this competition but the MT systems were not compared in regards to computational resources. In future work we will support equal computational resources for all teams.
3. Bambara is a low-resource language and the amount of data needed to significantly improve MT is very large. Inconsistent Bambara orthographies might mitigate translation quality improvement even with additional data collection. There are very high rates of illiteracy for Malians (35%, the 5th highest in the world (Diarra and Leventhal, 2020)) and Bambara speakers. We would like to gather and translate spoken Bambara audio data to counter these challenges.

Ethics Statement

Any evaluation system that incorporates human workers motivates reflection on the ethical implications of their contribution. Two of the teams competing in the competition had members that were able to annotate their system’s output for translation quality due to their Bambara knowledge. This was part of their team’s evaluation efforts and all the team members had already consented to participate in the competition.

In addition to considering how participating in the competition affected the team members, this work also affects the many millions of Bambara speakers who have not historically had access to technology. A recent focus on Machine Learning by the Malian government aims to change that (Diarra and Leventhal, 2020). As a consequence, increasing awareness and access to MT data, tasks, and their applications has wide global impact.

Finally, due to the BLEU scores the competing teams produced, these current translation systems should not be used in critical situations where inaccurate translations could lead to harm.

Acknowledgements

The data used in this work was supported by funding from the Artificial Intelligence Journal promoting AI research and previously utilized for WMT

(Barrault et al., 2019). We would like to thank Anonymous Organization for their support of the machine translation competition including the monetary awards for the winners.

References

- David Adelani, Jesujoba Alabi, Angela Fan, Julia Kreutzer, Xiaoyu Shen, Machel Reid, Dana Ruiter, Dietrich Klakow, Peter Nabende, Ernie Chang, Tajudeen Gwadabe, Freshia Sackey, Bonaventure F. P. Dossou, Chris Emezue, Colin Leong, Michael Beukman, Shamsuddeen Muhammad, Guyo Jarso, Oreen Yousuf, Andre Niyongabo Rubungo, Gilles Hacheme, Eric Peter Wairagala, Muhammad Umair Nasir, Benjamin Ajibade, Tunde Ajayi, Yvonne Gitau, Jade Abbott, Mohamed Ahmed, Millicent Ochieng, Anuoluwapo Aremu, Perez Ogayo, Jonathan Mukiibi, Fatoumata Ouoba Kabore, Godson Kalipe, Derguene Mbaye, Allahsera Auguste Tapo, Victoire Memdjokam Koagne, Edwin Munkoh-Buabeng, Valencia Wagner, Idris Abdulmumin, Ayodele Awokoya, Happy Buzaaba, Blessing Sibanda, Andiswa Bukula, and Sam Manthalu. 2022. *A few thousand translations go a long way! leveraging pre-trained models for African news translation*. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 3053–3070, Seattle, United States. Association for Computational Linguistics.
- Ali Araabi and Christof Monz. 2020. Optimizing transformer for low-resource neural machine translation. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 3429–3435.
- Loïc Barrault, Ondřej Bojar, Marta R Costa-Jussa, Christian Federmann, Mark Fishel, and Yvette Graham. 2019. Findings of the 2019 conference on machine translation (wmt19). Association for Computational Linguistics (ACL).
- Isaac Caswell, Julia Kreutzer, Lisa Wang, Ahsan Wahaab, Daan van Esch, Nasanbayar Ulzii-Orshikh, Allahsera Auguste Tapo, Nishant Subramani, Artem Sokolov, Claytone Sikasote, Monang Setyawan, Supheakmongkol Sarin, Sokhar Samb, Benoît Sagot, Clara Rivera, Annette Rios Gonzales, Isabel Papadimitriou, Salomey Osei, Pedro Ortiz Suarez, Iroko Orife, Kelechi Ogueji, Andre Niyongabo Rubungo, Toan Q. Nguyen, Mathias Muller, Andre Matthias Muller, Shamsuddeen Hassan Muhammad, Nanda Firdausi Muhammad, Ayanda Mnyakeni, Jamshidbek Mirzakhlov, Tapiwanashe Matangira, Colin Leong, Nze Lawson, Sneha Kudugunta, Yacine Jernite, M. Jenny, Orhan Firat, Bonaventure F. P. Dossou, Sakhile Dlamini, Nisansa de Silva, Sakine cCabuk Balli, Stella Rose Biderman, Alessia Battisti, Ahmed Baruwa, Ankur Bapna, Pallavi N. Baljekar, Israel Abebe Azime, Ayodele Awokoya, Duygu Ataman, Orevaoghene Ahia, Oghenefego Ahia, Sweta Agrawal, and Mofetoluwa Adeyemi. 2021. Quality

719	at a glance: An audit of web-crawled multilingual datasets. <i>Transactions of the Association for Computational Linguistics</i> , 10:50–72.	Edward Hu, Yelong Shen, Phil Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Lu Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models .	774
720			775
721			776
722	Xiangning Chen, Chen Liang, Da Huang, Esteban Real, Kaiyuan Wang, Yao Liu, Hieu Pham, Xuanyi Dong, Thang Luong, Cho-Jui Hsieh, Yifeng Lu, and Quoc V. Le. 2023. Symbolic discovery of optimization algorithms .	Marcin Junczys-Dowmunt, Roman Grundkiewicz, Tomasz Dwojak, Hieu Hoang, Kenneth Heafield, Tom Neckermann, Frank Seide, Ulrich Germann, Alham Fikri Aji, Nikolay Bogoychev, André F. T. Martins, and Alexandra Birch. 2018. Marian: Fast neural machine translation in C++ . In <i>Proceedings of ACL 2018, System Demonstrations</i> , pages 116–121, Melbourne, Australia.	777
723			778
724			779
725			780
726			781
727	Haby Sanou Diarra and Michael Leventhal. 2020. Developing machine learning competence in africa in the francophone sahel region.		782
728			783
729			784
730	Bonaventure F. P. Dossou and Chris C. Emezue. 2020. Ffr v1.0: Fon-french neural machine translation. <i>arXiv preprint arXiv: Arxiv-2003.12111</i> .	Tsz Kin Lam, Shigehiko Schamoni, and Stefan Riezler. 2021. Cascaded models with cyclic feedback for direct speech translation. In <i>ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)</i> , pages 7508–7512. IEEE.	785
731			786
732			787
733	Moussa Kamal Eddine, Guokan Shang, Antoine J. P. Tixier, and Michalis Vazirgiannis. 2021. Frugalscore: Learning cheaper, lighter and faster evaluation metrics for automatic text generation .		788
734			789
735			790
736			791
737	Chris Chinenye Emezue and Bonaventure F. P. Dossou. 2021. MMTAfrica: Multilingual machine translation for African languages . In <i>Proceedings of the Sixth Conference on Machine Translation</i> , pages 398–411, Online. Association for Computational Linguistics.	Brian Lester, Rami Al-Rfou, and Noah Constant. 2021. The power of scale for parameter-efficient prompt tuning . <i>CoRR</i> , abs/2104.08691.	792
738			793
739			794
740			795
741			796
742	Angela Fan, Shruti Bhosale, Holger Schwenk, Zhiyi Ma, Ahmed El-Kishky, Siddharth Goyal, Man-deep Baines, Onur Celebi, Guillaume Wenzek, Vishrav Chaudhary, Naman Goyal, Tom Birch, Vitaliy Liptchinsky, Sergey Edunov, Edouard Grave, Michael Auli, and Armand Joulin. 2020. Beyond english-centric multilingual machine translation. <i>arXiv preprint</i> .	Michael Leventhal, Allahsera Tapo, Sarah Luger, Marcos Zampieri, and Christopher M Homan. 2020. Assessing human translations from french to bambara for machine learning: a pilot study. <i>arXiv preprint arXiv:2004.00068</i> .	797
743			798
744			799
745			800
746			801
747			802
748			803
749			804
750	W, Wilhelmina Nekoto, Vukosi Marivate, Tshinondiwa Matsila, Timi Fasubaa, Tajudeen Kolawole, Taiwo Fagbohungebe, Solomon Oluwole Akinola, Shamsuddee Hassan Muhammad, Salomon Kabongo, Salomey Osei, et al. 2020. Participatory research for low-resourced machine translation: A case study in african languages. <i>Findings of EMNLP</i> .	Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Ves Stoyanov, and Luke Zettlemoyer. 2019. Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. <i>arXiv preprint arXiv:1910.13461</i> .	805
751			806
752			807
753			808
754			809
755			810
756			811
757	Naman Goyal, Cynthia Gao, Vishrav Chaudhary, Peng-Jen Chen, Guillaume Wenzek, Da Ju, Sanjana Krishnan, Marc Aurelio Ranzato, Francisco Guzmán, and Angela Fan. 2022. The FLORES-101 evaluation benchmark for low-resource and multilingual machine translation. <i>Transactions of the Association for Computational Linguistics</i> , 10:522–538.	Xiao Liu, Kaixuan Ji, Yicheng Fu, Zhengxiao Du, Zhilin Yang, and Jie Tang. 2021. P-tuning v2: Prompt tuning can be comparable to fine-tuning universally across scales and tasks . <i>CoRR</i> , abs/2110.07602.	812
758			813
759			814
760			815
761			816
762			817
763			818
764	Barry Haddow, Rachel Bawden, Antonio Valerio Miceli Barone, Jindřich Helcl, and Alexandra Birch. 2022. Survey of low-resource machine translation . <i>Computational Linguistics</i> , 48(3):673–732.	Shuming Ma, Li Dong, Shaohan Huang, Dongdong Zhang, Alexandre Muzio, Saksham Singhal, Hany Hassan Awadalla, Xia Song, and Furu Wei. 2021. Deltalm: Encoder-decoder pre-training for language generation and translation by augmenting pretrained multilingual encoders. <i>arXiv preprint arXiv:2106.13736</i> .	819
765			820
766			821
767			822
768	Felix Hieber, Tobias Domhan, Michael Denkowski, and David Vilar. 2020. Sockeye 2: A toolkit for neural machine translation . In <i>Proceedings of the 22nd Annual Conference of the European Association for Machine Translation</i> , pages 457–458, Lisboa, Portugal. European Association for Machine Translation.	Sourab Mangrulkar, Sylvain Gugger, Lysandre Debut, Younes Belkada, and Sayak Paul. 2022. Peft: State-of-the-art parameter-efficient fine-tuning methods. urlhttps://github.com/huggingface/peft .	823
769			824
770			825
771			826
772			827
773			828

829	Niklas Muennighoff, Thomas Wang, Lintang Sutawika,	Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel	883
830	Adam Roberts, Stella Biderman, Teven Le Scao,	Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford,	884
831	M Saiful Bari, Sheng Shen, Zheng-Xin Yong, Hai-	Dario Amodei, and Paul F Christiano. 2020. Learn-	885
832	ley Schoelkopf, Xiangru Tang, Dragomir Radev, Al-	ing to summarize with human feedback . In <i>Ad-</i>	886
833	ham Fikri Aji, Khalid Almubarak, Samuel Albanie,	<i>Advances in Neural Information Processing Systems</i> ,	887
834	Zaid Alyafeai, Albert Webson, Edward Raff, and	volume 33, pages 3008–3021. Curran Associates,	888
835	Colin Raffel. 2022a. Crosslingual generalization	Inc.	889
836	through multitask finetuning .		
837	Niklas Muennighoff, Thomas Wang, Lintang Sutawika,	Allahsera Auguste Tapo, Bakary Coulibaly, Sébastien	890
838	Adam Roberts, Stella Biderman, Teven Le Scao,	Diarra, Christopher Homan, Julia Kreutzer, Sarah	891
839	M Saiful Bari, Sheng Shen, Zheng-Xin Yong, Hai-	Luger, Arthur Nagashima, Marcos Zampieri, and	892
840	ley Schoelkopf, Xiangru Tang, Dragomir Radev, Al-	Michael Leventhal. 2020. Neural machine translation	893
841	ham Fikri Aji, Khalid Almubarak, Samuel Albanie,	for extremely low-resource african languages: A case	894
842	Zaid Alyafeai, Albert Webson, Edward Raff, and	study on bambara. <i>arXiv preprint arXiv:2011.05284</i> .	895
843	Colin Raffel. 2022b. Crosslingual generalization		
844	through multitask finetuning .	NLLB Team, Marta R. Costa-jussà, James Cross, Onur	896
845	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-	Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Hef-	897
846	Jing Zhu. 2002. Bleu: a method for automatic evalu-	ernan, Elahe Kalbassi, Janice Lam, Daniel Licht,	898
847	ation of machine translation. In <i>Proceedings of the</i>	Anna Sun Jean Maillard, Skyler Wang, Guillaume	899
848	<i>40th annual meeting of the Association for Computa-</i>	Wenzek, Al Youngblood, Bapi Akula, Loic Bar-	900
849	<i>tional Linguistics</i> , pages 311–318.	rault, Gabriel Mejia Gonzalez, Prangthip Hansanti,	901
850	André Susano Pinto, Alexander Kolesnikov, Yuge Shi,	John Hoffman, Semarley Jarrett, Kaushik Ram	902
851	Lucas Beyer, and Xiaohua Zhai. 2023. Tuning com-	Sadagopan, Dirk Rowe, Shannon Spruit, Chau	903
852	puter vision models with task rewards .	Tran, Pierre Andrews, Necip Fazil Ayan, Shruti	904
853	Alberto Poncelas, Dimitar Shterionov, Andy Way,	Bhosale, Sergey Edunov, Angela Fan, Cynthia	905
854	Gideon Maillette de Buy Wenniger, and Peyman Pass-	Gao, Vedanuj Goswami, Francisco Guzmán, Philipp	906
855	ban. 2018. Investigating backtranslation in neural	Koehn, Alexandre Mourachko, Christophe Ropers,	907
856	machine translation. <i>arXiv preprint arXiv: Arxiv-</i>	Safiyyah Saleem, Holger Schwenk, and Jeff Wang.	908
857	<i>1804.06189</i> .	2022. No language left behind: Scaling human-	909
858	Matt Post. 2018. A call for clarity in reporting BLEU	centered machine translation.	910
859	scores . In <i>Proceedings of the Third Conference on</i>	Jörg Tiedemann and Santhosh Thottingal. 2020. OPUS-	911
860	<i>Machine Translation: Research Papers</i> , pages 186–	MT – building open translation services for the world .	912
861	191, Brussels, Belgium. Association for Computa-	In <i>Proceedings of the 22nd Annual Conference of</i>	913
862	tional Linguistics.	<i>the European Association for Machine Translation</i> ,	914
863	Nils Reimers and Iryna Gurevych. 2020. Making	pages 479–480, Lisboa, Portugal. European Associa-	915
864	monolingual sentence embeddings multilingual us-	tion for Machine Translation.	916
865	ing knowledge distillation . In <i>Proceedings of the</i>	Elan Van Biljon, Arnu Pretorius, and Julia Kreutzer.	917
866	<i>2020 Conference on Empirical Methods in Natural</i>	2020. On optimal transformer depth for low-	918
867	<i>Language Processing</i> . Association for Computational	resource language translation. <i>arXiv preprint</i>	919
868	Linguistics.	<i>arXiv:2004.04418</i> .	920
869	Rico Sennrich, Barry Haddow, and Alexandra Birch.	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob	921
870	2016a. Improving neural machine translation models	Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz	922
871	with monolingual data. <i>ACL</i> .	Kaiser, and Illia Polosukhin. 2017. Attention is all	923
872	Rico Sennrich, Barry Haddow, and Alexandra Birch.	you need. <i>Advances in neural information processing</i>	924
873	2016b. Neural machine translation of rare words	<i>systems</i> , 30.	925
874	with subword units . In <i>Proceedings of the 54th An-</i>	Valentin Vydrin, Jean-Jacques Meric, Kirill Maslin-	926
875	<i>Annual Meeting of the Association for Computational</i>	sky, Andriy Rovenchak, Allahsera Auguste Tapo, Se-	927
876	<i>Linguistics (Volume 1: Long Papers)</i> , pages 1715–	bastien Diarra, Christopher Homan, Marco Zampieri,	928
877	1725, Berlin, Germany. Association for Computa-	and Michael Leventhal. 2022. Machine learn-	929
878	tional Linguistics.	ing dataset development for manding languages.	930
879	Shailashree K Sheshadri, Deepa Gupta, and Marta R	urlhttps://github.com/robotsmali-ai/datasets .	931
880	Costa-Jussà. 2023. A voyage on neural machine	Foundation Wikimedia. 2023. Wikimedia downloads .	932
881	translation for indic languages. <i>Procedia Computer</i>	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien	933
882	<i>Science</i> , 218:2694–2712.	Chaumond, Clement Delangue, Anthony Moi, Pier-	934
		ric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz,	935
		Joe Davison, Sam Shleifer, Patrick von Platen, Clara	936
		Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le	937
		Scao, Sylvain Gugger, Mariama Drame, Quentin	938
		Lhoest, and Alexander M. Rush. 2020. Transform-	939
		ers: State-of-the-art natural language processing . In	940

Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations, pages 38–45, Online. Association for Computational Linguistics.

Linting Xue, Aditya Barua, Noah Constant, Rami Al-Rfou, Sharan Narang, Mihir Kale, Adam Roberts, and Colin Raffel. 2021a. [Byt5: Towards a token-free future with pre-trained byte-to-byte models](#).

Linting Xue, Noah Constant, Adam Roberts, Mihir Kale, Rami Al-Rfou, Aditya Siddhant, Aditya Barua, and Colin Raffel. 2021b. [mT5: A massively multilingual pre-trained text-to-text transformer](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 483–498, Online. Association for Computational Linguistics.

Biao Zhang, Philip Williams, Ivan Titov, and Rico Senrich. 2020. [Improving massively multilingual neural machine translation and zero-shot translation](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 1628–1639, Online. Association for Computational Linguistics.

Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. 2019. [Bertscore: Evaluating text generation with bert](#).