

The African Languages Lab: A Global Low-Resource Language Collaborative Approach to Advancing NLP for African Languages

Anonymous ACL submission

Abstract

The digital revolution has left behind hundreds of millions of speakers of low-resource languages (LRLs), particularly in Africa, creating a critical gap in global information access and technological representation. We introduce the African Languages Lab (All Lab), a systematic approach to advancing NLP capabilities for African LRLs through a coordinated research framework. Our work introduces (1) a quality-controlled data collection pipeline, yielding the largest validated multi-modal speech and text dataset for African LRLs spanning 40 languages, encompassing 500 GB of combined parallel text and 4,000 hours of aligned speech data and (2) experiments demonstrating how our custom dataset, combined with QLoRA, achieves improvements across multiple metrics (up to +49.8 chrF++, +60.2 BLEU, and +0.28 COMET points) compared to the base model. Our work establishes a sustainable framework for expanding NLP capabilities to historically underserved languages while fostering local research capacity through structured mentorship and collaboration. We will release our data for research.

1 Introduction

Most of the artificial intelligence (AI) and natural language processing (NLP) research today focuses on about 20 of the 7000 languages of the world, leaving the vast majority of languages understudied (Magueresse et al., 2020). Without a clearly established definition, LRLs are languages that exist at the periphery of the digital transformation, characterized by three critical deficits: (1) a scarcity of machine-readable corpora (2) limited personalized computational technologies and trained language models (3) insufficient representation in global research communities (Nigatu et al., 2024; Issaka et al., 2024; Magueresse et al., 2020). While often serving substantial speaker populations, these languages face significant challenges in participating

fully in the AI-driven information economy.

For Africa, the scale of this crisis is staggering: over 2,000 languages are spoken across Africa (nearly one-third of all languages worldwide). Yet, a stunning 88% of African languages are "severely underrepresented" or "completely ignored" in computational linguistics (Joshi et al., 2020). As shown in Figure 1, about 814 African languages are in danger of extinction. Countries like Nigeria, Cameroon, and Ivory Coast have 171, 75, and 65 languages facing extinction, respectively. This exclusion has far-reaching consequences, from poor educational and healthcare outcomes to preventing full participation in the digital economy (Laitin et al., 2019; Gessler and von der Wense, 2024)

This problem is compounded by a severe underrepresentation in the global NLP research community. Analyzing mentions of the top 10 global languages versus the top 10 African languages across major academic databases reveals that, on average, for every paper discussing African languages in multilingual LLM contexts, there are 20 papers on global languages in Google Scholar, 23 in CORE, 34 in arXiv, and 70 in IEEE (Table 1) and (Table 4 in the Appendix). This systematic underrepresentation reinforces a cycle of technological marginalization.

Contributing to broader efforts to bridge this systemic technological gap, we present the African Languages Lab (All Lab), an initiative to democratize NLP technology for African languages. Started in 2020, the All Lab operates through a tightly coordinated team of dedicated researchers who combine three innovative elements:

1. a systematic, quality-controlled data collection framework powered by our "All Voice" platform, which has enabled the creation of 500GB of validated multi-modal speech and text data,
2. state-of-the-art multilingual modeling tech-

niques optimized for low-resource scenarios, achieving an average +26.7 BLEU, +0.22 COMET, and chrF++ +17.31 improvement over existing baselines, and

3. a sustainable research-first approach that cultivates and empowers young researchers through formal mentorship and collaboration structures.

Our framework represents a step toward digital inclusion for millions of African language speakers, demonstrating that systematic, research-driven, and community-involved approaches can effectively bridge the technological divide while preserving linguistic diversity.

2 Related Work

Our work builds on three pillars of African NLP research: community-driven initiatives, advances in multilingual LLMs, and benchmarks and evaluation frameworks. We examine how these connected areas have shaped the landscape of African NLP.

2.1 Community-Driven Research Initiatives

The development of African NLP has been shaped by several complementary community and institutional efforts. Masakhane, comprising 2,947 Slack members as of January 2025, represents the largest, sustained, community-driven NLP initiatives for African languages (Orife et al., 2020). Complementing this work, the "Breaking the Unwritten Language Barrier" project addresses challenges specific to unwritten and under-documented languages (Adda et al., 2016). Their work on development for languages like Basaa, Myene, and Embosi has established methodological approaches for speech recognition in LRLs.

These community efforts have been supported by institutional initiatives providing essential infrastructure. The Lacuna Fund has enabled dataset development (Rathi et al., 2023), while Meta's No Language Left Behind project has contributed to advancing multilingual modeling capabilities (Team et al., 2022). Additional infrastructure support has come from Mozilla's Common Voice project (Ardila et al., 2020) for speech resources and the AI4D African Language Program (Siminyu et al., 2021) for benchmark development. The Deep Learning Indaba¹ has contributed to research capacity building through its convenings, while plat-

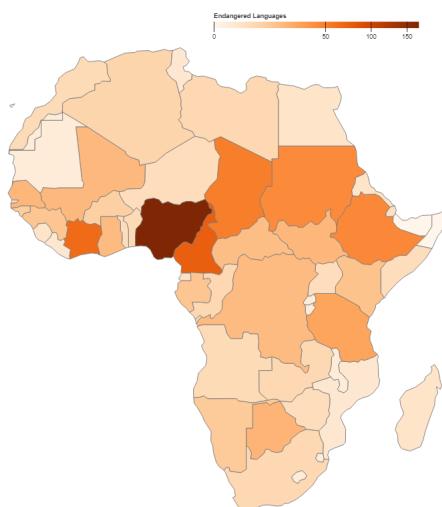


Figure 1: Illustrates the number of endangered languages in each African country using data from Ethnologue

forms like Lanfrica have improved resource discoverability and research sharing across the continent (Emezue and Dossou, 2020).

2.2 Advances in African Multilingual NLP

The evolution of multilingual LLMs has shown steady progress in language coverage and capabilities. Early approaches like mBERT (Muller et al., 2021) and XLM-R (Conneau et al., 2020) established initial benchmarks, supporting approximately 100 languages each. Subsequent developments included more focused models like mBART (Liu et al., 2020), mT5 (Xue et al., 2021), and XGLM (Ersoy et al., 2023), which traded broader language coverage for improved performance on specific language sets. The advent of massive LLMs further expanded these capabilities, with models like GPT-3, mGPT (Shliazhko et al., 2024), and BLOOM (Workshop et al., 2023) supporting varying numbers of African languages. Also, Glot500-m (Imani et al., 2023) extends support to 511 languages and the SERENGETI and Cheetah models supports about 517 African languages (Adebara et al., 2023, 2024). Additional progress has come from the Aya model, which demonstrates instruction-following capabilities across 101 languages (Üstün et al., 2024), and specialized models like AfroLM, which focuses on 23 African languages (Dossou et al., 2022).

While not specifically trained in African languages, English-centric LLMs such as GPT-4 (Ope-

¹<https://deeplearningindaba.com/>

High-Resource Languages					African Languages				
Language	GS	arXiv	IEEE	CORE	Language	GS	arXiv	IEEE	CORE
English	14,700	323	256	3,095	Swahili	617	10	3	114
Chinese	7,710	60	85	1,694	Hausa	261	1	0	49
Hindi	1,980	20	41	336	Yoruba	276	1	0	59
Spanish	4,240	29	24	908	Igbo	203	0	0	38
Arabic	3,150	25	24	616	Amharic	338	2	2	49
French	4,490	38	17	1,037	Oromo	104	1	1	21
Bengali	943	9	8	183	Berber	55	0	0	11
Portuguese	1,980	13	7	400	Zulu	175	1	1	38
Russian	2,950	19	16	611	Fula	20	0	0	7
Urdu	728	3	9	131	Malagasy	72	0	0	15

Table 1: Research visibility analysis comparing publication volumes for top 10 global languages versus top 10 African languages across major academic databases (Google Scholar (GS), arXiv, IEEE, and CORE) from 2020 to 2024. The stark contrast in publication volumes highlights the digital divide in academic research visibility.

nAI et al., 2024), Gemini (Bao et al., 2023), and Llama (Wendler et al., 2024) have shown capability in handling some African languages, (Robinson et al., 2023; Ojo et al., 2024; Zhu et al., 2024; Dong et al., 2024), though their performance generally does not match that of specialized models.

2.3 Benchmarks and Evaluation Frameworks

The development of evaluation frameworks has enabled systematic measurement of progress in African NLP. They span multiple task domains: MasakhaNER provides Named Entity Recognition datasets for 10 languages (Adelani et al., 2021), AfriSenti offers sentiment analysis benchmarks across 14 languages (Muhammad et al., 2023), and AFROMT establishes standardized translation benchmarks for 8 languages (Reid et al., 2021).

More targeted evaluation resources include NaijaSenti for Nigerian languages (Muhammad et al., 2022) and Kencorpus for Kenyan languages (Wanjawa et al., 2023). IrokoBench provides human-translated datasets across 17 typologically-diverse African languages (Adelani et al., 2025). These Africa-focused frameworks complement broader initiatives like FLORES200 (Team et al., 2022), the Aya Dataset (Singh et al., 2024b), and Global-MMLU (Singh et al., 2024a).

Despite these developments, significant challenges remain in African NLP research (Adebara and Abdul-Mageed, 2022; Issaka et al., 2024). Our work builds upon these foundations while addressing several key limitations in existing approaches.

3 Methodology

3.1 Datasets

All Voices Platform. To address the fundamental challenge of data scarcity in African languages, we developed All Voices, a mobile-first platform that stands as the only solution specifically designed for data collection in any LRL. The platform’s innovative approach enables direct translation between LRLs without requiring English as an intermediary, addressing a critical gap in the existing data collection infrastructure. Also, All Voices distinguishes itself through its multimodal capabilities, supporting text and audio data collection and validation. The platform features an intuitive, user-friendly interface that encourages broad participation, complemented by gamification elements, including a global leaderboard system that promotes user engagement. Importantly, All Voices ² is open and free to everyone, aligning with our mission to democratize language technology development.

The platform’s architecture, built using React-Native ³ and Firebase ⁴, integrates user authentication and analytics, translation corpus management, and quality control components. Our authentication system provides comprehensive user profiling, tracking contributor demographics and expertise through quantifiable metrics, including successful translations and community validation scores. This system implements OAuth 2.0 authentication and role-based access control to ensure data integrity and user privacy. The translation corpus manage-

²[inserted_after_accept](#)

³<https://reactnative.dev/>

⁴<https://firebase.google.com/>

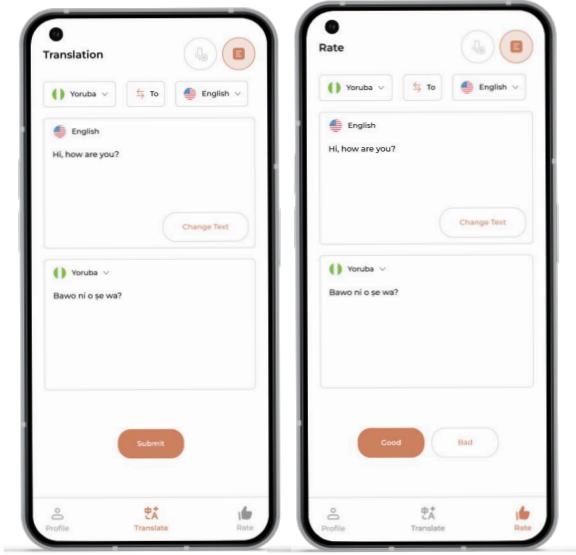


Figure 2: The All Voices platform interface demonstrating its dual functionality: direct text translation from English to Yoruba (left panel) and community-driven translation validation system (right panel).

ment provides centralized storage for textual and audio translations with associated metadata, securing all data through AES-256 encryption at rest and TLS 1.3 in transit. Translations undergo peer review requiring both a minimum threshold of positive validation (>5 upvotes) and an acceptable error margin (<3 downvotes) to achieve verified status. A key innovation is our recursive translation pipeline: verified translations become eligible source material for subsequent translations, creating a multiplicative effect in data collection.

Data Collection and Processing. Our dataset development methodology combines crowd-sourced translations through All Voices (Figure 2) with carefully curated open-source corpora. We integrate validated translations from our platform with established datasets, including NLLB (Team et al., 2022), CCMATRIX (Wenzek et al., 2019), OpenSubtitles (Tiedemann, 2016), MultiCCAligned (El-Kishky et al., 2020), ParaCrawl (Bañón et al., 2020), XLEnt (El-Kishky et al., 2021), MultiParaCrawl(Bañón et al., 2020), LinguaTools-WikiTitles (Tiedemann, 2012), and CCAigned (El-Kishky et al., 2020). Additionally, we collect new datasets through community partners.

Our data processing implements a robust two-tier approach combining general normalization with language-specific processing. The general

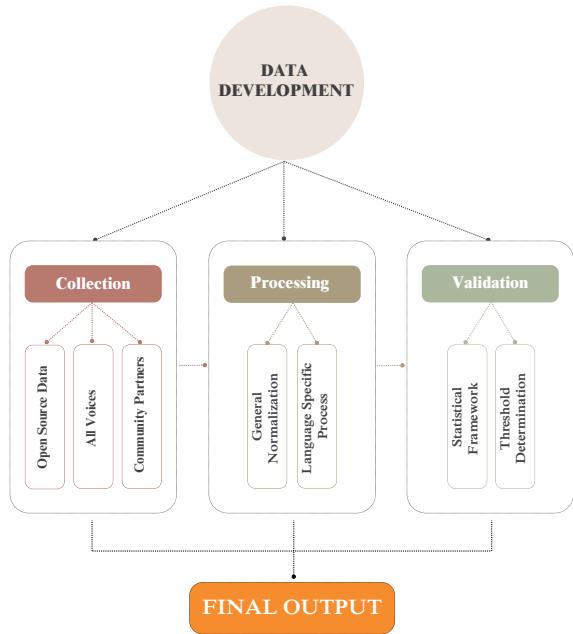


Figure 3: End-to-end data processing pipeline for African language corpus development, illustrating: (a) multi-source data collection, (b) language-specific preprocessing, and (c) validation frameworks for ensuring dataset integrity.

normalization phase addresses universal text artifacts through Unicode normalization, character encoding standardization, and structural cleaning, including HTML removal and symbol standardization. The language-specific processing phase implements specialized handling for African language features, including morphological analysis, script variant normalization, and tone mark standardization, with custom rule sets developed for specific language families.

Next, our translation validation methodology implements a robust statistical framework for assessing translation quality through quantitative analysis of character-level distributions. The validation metric employs character ratio analysis between source and target texts, computed as the ratio of target text length to source text length. We analyze these ratios using z-score normalization within language-specific distribution, enabling the detection of statistical outliers while accounting for natural variations in text length across different language pairs. This approach is augmented with character overlap detection to identify potential artifacts or inappropriate text preservation, particularly crucial for languages sharing similar orthographic features.

Also, the threshold determination process implements an adaptive sampling methodology. For each language pair, we establish baseline distributions through initial sampling of 10,000 translation pairs, employing Kernel Density Estimation for robust distribution modeling. This approach effectively captures the non-Gaussian characteristics frequently observed in cross-lingual character distributions. Thresholds are dynamically computed using a modified Tukey method with an adaptive multiplier. This adaptive threshold mechanism automatically calibrates to language-specific characteristics, implementing more stringent filtering for language pairs exhibiting consistent ratios while allowing appropriate flexibility for pairs with inherently higher variability. The resulting validation framework effectively identifies and filters anomalous translations while maintaining sensitivity to legitimate linguistic variations across diverse African language families. The processed datasets are structured following HuggingFace⁵ Dataset specifications, enabling seamless API integration.

3.2 Model Development

To evaluate our dataset’s utility and establish baselines for 15 randomly selected African language translations, we experimented with Llama-3.2-3B-Instruct (Grattafiori et al., 2024). We chose this model as our base architecture due to its demonstrated multilingual capabilities and efficient parameter scaling, making it particularly suitable for low-resource scenarios. Given the computational constraints and the need to train on multiple languages efficiently, we implemented parameter-efficient fine-tuning using Quantization-aware Low-Rank Adaptation (QLoRA) (Üstün et al., 2024). Specifically, we employed 4-bit quantization with a LoRA rank of 4, which our preliminary experiments showed to provide an optimal balance between memory efficiency and performance.

The training pipeline was implemented using the Transformers library, with the model fine-tuned using a supervised approach with a consistent instruction format: ”Translate the following English text to X:”, where X represents the target African language. We configured the model with a maximum input sequence length of 512 tokens and limited the maximum generation length to 600 tokens. To manage computational resources effectively while ensuring comprehensive coverage, we capped the

training data at 1 million random examples per language. The training process utilized an NVIDIA H100 GPU with a batch size 16, and gradient accumulation steps of 16. We trained each model for 2 epochs using a learning rate 5.0×10^{-5} with cosine decay, warm-up ratio of 0.15, and BF16 mixed precision. For generation, we used a beam size of 5, temperature of 0.7, and top-p of 0.9.

3.3 Evaluation Metrics

Model performance was evaluated using a complementary set of metrics: BiLingual Evaluation Understudy (BLEU) (Wieting et al., 2019), which measures n-gram precision, METEOR (Banerjee and Lavie, 2005) which accounts for word stems and synonyms, COMET (Rei et al., 2020) which leverages multilingual embeddings to assess semantic similarity, and chrF++ (Wang et al., 2025), which operates on character-level n-grams to better capture morphological variations common in African languages. Together, these metrics comprehensively assess translation quality across different linguistic aspects. Also, the test set was capped at 5,000 random samples per language for consistent evaluation across all languages while maintaining diversity.

4 Results

4.1 Datasets

Through extensive data collection efforts spanning multiple sources, we have compiled a comprehensive dataset encompassing 40 African languages. Our analysis reveals a nuanced stratification with four distinct tiers of resource availability:

1. **Primary resource languages (>40 GB):** This tier includes Arabic, Swahili, and Hausa, representing languages that have benefited from sustained digitization efforts and strong institutional support.
2. **Established digital languages (20-40 GB):** This category includes Afrikaans, Bemba, Amharic, and Xhosa. These languages demonstrate robust digital presence, likely due to consistent documentation and preservation initiatives.
3. **Emerging digital languages (5-20 GB):** A substantial group including Luganda, Lingala, Yoruba, and Malagasy. These languages show emerging digital footprints but still face resource constraints.

⁵<https://huggingface.co/>

Language	Size	Language	Size	Language	Size	Language	Size
Arabic	85.6 GB	Amharic	25.87 GB	Zulu	17.96 GB	Rundi	372.43 MB
Swahili	44.73 GB	Bemba	26.69 GB	Malagasy	16.41 GB	Kikongo	418.38 MB
Hausa	40.73 GB	Xhosa	24.65 GB	Tswana	16.54 GB	Mossi	245.45 MB
Afrikaans	35.70 GB	Ewe	18.52 GB	Luganda	15.96 GB	Tshiluba	282.56 MB
Twi	31.00 GB	Sesotho	18.55 GB	Lingala	13.39 GB	Bambara	200.60 MB
Yoruba	12.05 GB	Kinyarwanda	12.54 GB	Igbo	6.28 GB	Umbundu	147.64 MB
Wolof	8.77 GB	Kikuyu	5.15 GB	Oromo	5.22 GB	Berber	118.63 MB
Chewa	4.85 GB	Shona	4.19 GB	Somali	4.19 GB	Krio	40.51 MB
Fon	1.59 GB	Tigrinya	1.24 GB	Mandinka	33.13 MB	Ngambay	47.61 MB
Fula	404.58 MB	Kanuri	3.92 MB	Fang	36.89 KB	Kiluba	75.20 KB

Table 2: Dataset composition across our 40 African languages, showing the distribution of combined speech and text corpus in GB.

Language	chrF++		COMET		BLEU	
	Llama3B	Finetuned	Llama3B	Finetuned	Llama3B	Finetuned
Amharic	9.2	23.2	0.44	0.72	4.2	21.0
Fula	15.7	42.3	0.38	0.64	2.2	21.3
Yoruba	17.0	19.8	0.35	0.62	1.44	3.2
Igbo	11.5	61.3	0.43	0.68	11.0	71.2
Oromo	18.7	27.7	0.38	0.60	2.70	12.7
Swahili	29.4	43.0	0.59	0.70	19.0	21.0
Hausa	19.7	43.2	0.40	0.67	2.21	48.6
Twi	18.1	52.8	0.41	0.69	3.53	61.1
Shona	17.3	29.2	0.39	0.57	2.20	38.0
Somali	18.8	37.0	0.38	0.64	1.32	40.5
Kinyarwanda	19.1	27.7	0.41	0.57	2.82	8.23
Ewe	16.4	21.4	0.35	0.61	1.84	16.2
Bambara	14.8	34.0	0.36	0.56	1.48	28.3
Wolof	14.2	28.2	0.37	0.53	1.39	32.5
Luganda	24.2	32.9	0.39	0.57	7.44	41.5

Table 3: Performance comparison between base Llama3B and Finetuned Llama3B models across different metrics (chrF++, COMET, and BLEU) for various African languages. Higher scores indicate better performance.

4. Resource-constrained languages (<5 GB):

The majority of languages in our dataset fall into this category, including widely spoken languages like Bambara, Kikongo, and Fang. This tier reveals critical gaps in digital infrastructure, even for languages with significant speaker populations.

Our analysis highlights a stark digital divide, with the top 3 languages accounting for about 35% of our data. This disparity underscores the urgent need for targeted resource development efforts, particularly for languages with substantial speaker populations but limited digital presence.

4.2 Models

Our evaluation of Llama-3.2-3B-Instruct and its fine-tuned variant across 15 African languages revealed nuanced patterns in translation performance. The base model exhibited varying degrees of effectiveness across languages, with chrF++ scores ranging from 9.2 to 29.4, COMET scores from 0.35 to 0.59, and BLEU scores from 1.32 to 19.0. Notably, Swahili demonstrated good performance across all metrics (chrF++ 29.4, COMET 0.59, BLEU 19.0), likely attributable to its substantial presence in the pre-training corpus and rich digital resources. Fine-tuning yielded substantial but heterogeneous improvements across the language set. The most remarkable gains were observed in Igbo, where performance increased dramatically across

372
373
374
375
376
377
378

379
380
381
382
383
384

385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400

401 all metrics (chrF++ from 11.5 to 61.3, COMET
402 from 0.43 to 0.68, BLEU from 11.0 to 71.2), and
403 Twi (chrF++ from 18.1 to 52.8, COMET from
404 0.41 to 0.69, BLEU from 3.53 to 61.1). These
405 improvements suggest that targeted fine-tuning can
406 effectively bridge the performance gap for African
407 LRLs. However, the impact of fine-tuning varied
408 significantly across languages. While languages
409 like Hausa and Fula showed substantial improve-
410 ments (chrF++ increases of 23.5 and 26.6 points
411 respectively), others like Yoruba and Ewe exhib-
412 ited more modest gains (chrF++ increases of 2.8
413 and 5.0 points). This variability in improvement
414 patterns suggests complex interactions between
415 model architecture, pre-training data distribution,
416 and language-specific characteristics that warrant
417 further investigation. COMET scores, which corre-
418 late better with human judgments, showed consist-
419 ent improvements across all languages after fine-
420 tuning, with gains ranging from 0.11 to 0.28 points.
421 This uniform improvement in COMET scores, de-
422 spite varying changes in chrF++ and BLEU, indi-
423 cates that fine-tuning enhances translation quality
424 in ways that may not be fully captured by surface-
425 level metrics.

426 5 Conclusion

427 We presented the African Languages Lab (All Lab),
428 introducing a comprehensive initiative for advanc-
429 ing NLP capabilities in African languages through
430 three key innovations:

- 431 1. the All Voices platform, which enabled the
432 creation of the largest validated multi-modal
433 dataset spanning 40 African languages,
- 434 2. experimental validation demonstrating im-
435 provements through QLoRA fine-tuning (av-
436 erage gains of +26.7 BLEU, +0.22 COMET,
437 and +17.31 chrF++), and
- 438 3. a sustainable research development program
439 that has successfully mentored fifteen early-
440 career researchers.

441 Our results demonstrate that targeted data collec-
442 tion and parameter-efficient fine-tuning can sig-
443 nificantly improve performance for low-resource
444 African languages, while our open-source approach
445 establishes a replicable framework for expanding
446 NLP capabilities to other underserved languages.

447 6 Limitations

448 6.1 Computational and Model Constraints

449 Our experimental evaluation was conducted us-
450 ing a single base model (Llama-3.2-3B-Instruct)
451 with QLoRA fine-tuning, primarily due to compu-
452 tational constraints. While this approach allowed
453 us to demonstrate the potential of our dataset, it
454 may not fully represent the optimal architecture
455 for African language processing. Additionally, we
456 were only able to evaluate on a subset of our col-
457 lected languages and data, which may not fully
458 represent the diversity of African languages.

459 6.2 Data Quality and Validation

460 Our data cleaning and standardization procedures
461 rely heavily on automated approaches due to lim-
462 ited access to native speakers across all 40 lan-
463 guages. While we implement grounded statistical
464 validation methods, this automation may miss sub-
465 tle linguistic nuances, dialectal variations, and cul-
466 tural contexts that human validators would catch.
467 Also, the quality assessment of our dataset, may
468 not fully capture potential biases or quality issues.

469 6.3 Platform and Infrastructure

470 The All Voices platform, while innovative, cur-
471 rently operates primarily through mobile interfaces,
472 which may limit participation from communities
473 with different technology preferences or access pat-
474 terns. The platform’s quality control mechanisms,
475 while systematic, may inadvertently favor certain
476 linguistic varieties over others.

477 These limitations inform our ongoing work and
478 highlight important areas for future research in
479 African NLP. They also underscore the need for
480 continued investment in computational resources,
481 human expertise, and infrastructure development
482 to support comprehensive technology development
483 for African languages.

484 7 Ethics Statement and Broader Impacts

485 7.1 Research Capacity Building

486 The All Lab maintains a structured research de-
487 velopment program which has successfully men-
488 tored fifteen early-career researchers across four
489 institutions through a comprehensive twelve-week
490 curriculum. The program implements a four-phase
491 model: foundation building, guided research, in-
492 dependent project development, and project com-
493 pletion. Each researcher receives one-on-one men-

494 torship from experienced researchers, with high-
495 performing participants transitioning into extended
496 research roles. This approach has shown measur-
497 able success, evidenced by achievements such as
498 the Ozy Genius Award recognition, while establish-
499 ing a sustainable pipeline for African NLP talent
500 development.

501 7.2 Sustainable Societal Impact

502 The All Lab’s work advances several United Na-
503 tions Sustainable Development Goals, particularly
504 in education and inequality reduction. Our plat-
505 form democratizes access to digital resources for
506 millions of speakers of African LRLs, enabling
507 communities to preserve their linguistic heritage
508 while participating in the digital economy. This
509 work addresses historical technological disparities
510 through three key mechanisms: increasing digital
511 representation of marginalized languages, enabling
512 community-led content creation, and facilitating
513 open knowledge transfer.

514 7.3 Challenges and Mitigation Strategies

515 We acknowledge several critical challenges in our
516 work. The commercial viability of LRL technolo-
517 gies remains limited, affecting sustainable devel-
518 opment. Data collection and evaluation face sig-
519 nificant hurdles due to limited digital presence and
520 infrastructure constraints. To address these chal-
521 lenges, we implement: (1) partnerships with aca-
522 demic and industry stakeholders to ensure resource
523 sustainability, (2) a quality-controlled data collec-
524 tion framework that balances automation with hu-
525 man validation, and (3) structured community en-
526 gagement programs to ensure cultural and linguis-
527 tic authenticity.

528 7.4 Future Directions and Long-term Impact

529 Our work establishes a framework for sustain-
530 able LRL technology development guided by the
531 African philosophy of Ubuntu (a philosophy that
532 affirms the positive values of inclusivity, com-
533 munity, difference, anti-racism, hospitality, and
534 openness to others). We outline a clear roadmap
535 for future development, including expanding lan-
536 guage coverage, optimizing model architectures for
537 African LRLs, and strengthening research partner-
538 ships. This approach not only advances technical
539 capabilities but also contributes to cultural preser-
540 vation, educational advancement, and economic
541 inclusion for hundreds of millions of speakers of

542 low-resource African languages. Through these ini-
543 tiatives, we aim to contribute to addressing the sys-
544 tematic underrepresentation of African languages
545 in NLP while establishing replicable methodolo-
546 gies for LRL technology development globally.
547 Our work demonstrates that technical innovation
548 in NLP can directly contribute to broader societal
549 goals while maintaining rigorous research.

550 References

Gilles Adda, Sebastian Stürker, Martine Adda-Decker,
551 Odette Ambouroue, Laurent Besacier, David Bla-
552 chon, Hélène Bonneau-Maynard, Pierre Godard, Fa-
553 timia Hamlaoui, Dmitry Idiatov, Guy-Noël Kouarata,
554 Lori Lamel, Emmanuel-Moselly Makasso, Annie
555 Rialland, Mark Van de Velde, François Yvon, and
556 Sabine Zerbian. 2016. *Breaking the unwritten lan-*
557 *guage barrier: The bulb project.* *Procedia Computer*
558 *Science*, 81:8–14. SLTU-2016 5th Workshop on Spo-
559 *ken Language Technologies for Under-resourced lan-*
560 *guages* 09-12 May 2016 Yogyakarta, Indonesia.

Ife Adebara and Muhammad Abdul-Mageed. 2022.
561 [Towards afrocentric nlp for african languages:](#)
562 [Where we are and where we can go.](#) *Preprint*,
563 arXiv:2203.08351.

Ife Adebara, AbdelRahim Elmadany, and Muhammad
564 Abdul-Mageed. 2024. [Cheetah: Natural language](#)
565 [generation for 517 african languages.](#) *Preprint*,
566 arXiv:2401.01053.

Ife Adebara, AbdelRahim Elmadany, Muhammad
567 Abdul-Mageed, and Alcides Alcoba Inciarte. 2023.
568 [SERENGETI: Massively multilingual language](#)
569 [models for Africa.](#) In *Findings of the Association for*
570 *Computational Linguistics: ACL 2023*, pages 1498–
571 1537, Toronto, Canada. Association for Compu-
572 *tational Linguistics.*

David Ifeoluwa Adelani, Jade Abbott, Graham Neu-
573 big, Daniel D’souza, Julia Kreutzer, Constantine Lig-
574 nos, Chester Palen-Michel, Happy Buzaaba, Shruti
575 Rijhwani, Sebastian Ruder, Stephen Mayhew, Is-
576 rael Abebe Azime, Shamsuddeen H. Muhammad,
577 Chris Chinene Emezue, Joyce Nakatumba-Nabende,
578 Perez Ogayo, Aremu Anuoluwapo, Catherine Gitau,
579 Derguene Mbaye, Jesujoba Alabi, Seid Muhie Yi-
580 mam, Tajuddeen Rabiu Gwadabe, Ignatius Ezeani,
581 Rubungo Andre Niyongabo, Jonathan Mukiibi, Ver-
582 rah Otiende, Iroro Orife, Davis David, Samba Ngom,
583 Tosin Adewumi, Paul Rayson, Mofetoluwa Adeyemi,
584 Gerald Muriuki, Emmanuel Anebi, Chiamaka Chuk-
585 wunike, Nkiruka Odu, Eric Peter Wairagala, Samuel
586 Oyerinde, Clemencia Siro, Tobius Saul Bateesa,
587 Temilola Oloyede, Yvonne Wambui, Victor Akin-
588 ode, Deborah Nabagereka, Maurice Katusiime, Ayo-
589 dele Awokoya, Mouhamadane MBOUP, Dibora Ge-
590 breyohannes, Henok Tilaye, Kelechi Nwaike, De-
591 gaga Wolde, Abdoulaye Faye, Blessing Sibanda, Ore-
592 vaoghene Ahia, Bonaventure F. P. Dossou, Kelechi

598	Ogueji, Thieno Ibrahima DIOP, Abdoulaye Diallo, Adewale Akinfaderin, Tendai Marengereke, and Salomey Osei. 2021. Masakhaner: Named entity recognition for african languages . <i>Transactions of the Association for Computational Linguistics</i> , 9:1116–1131.	600	Andrew Caines. 2015. The geographic diversity of nlp conferences . <i>MAREK REI</i> , arXiv:1503.06733. Version 2.	601	658
602		603		659	
604	David Ifeoluwa Adelani, Jessica Ojo, Israel Abebe Azime, Jian Yun Zhuang, Jesujoba O. Alabi, Xuanli He, Millicent Ochieng, Sara Hooker, Andiswa Bukula, En-Shiun Annie Lee, Chiamaka Chukwuneke, Happy Buzaaba, Blessing Sibanda, Godson Kalipe, Jonathan Mukiibi, Salomon Kabongo, Foutse Yuehgooh, Mmabidi Setaka, Lolwethu Ndolela, Nkiruka Odu, Rooweither Mabuya, Shamsuddeen Hassan Muhammad, Salomey Osei, Sokhar Samb, Tadesse Kebede Guge, Tombekai Vangoni Sherman, and Pontus Stenetorp. 2025. Irokobench: A new benchmark for african languages in the age of large language models . <i>Preprint</i> , arXiv:2406.03368.	605	Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2020. Unsupervised cross-lingual representation learning at scale . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 8440–8451, Online. Association for Computational Linguistics.	606	661
607		608		662	
609		610		663	
611		612		664	
613		614		665	
615		616		666	
617	Cynthia Jayne Amol, Evelyn Asiko Chimoto, Rose Delilah Gesicho, Antony M. Gitau, Naome A. Etori, Caringtone Kinyanjui, Steven Ndung'u, Lawrence Moruye, Samson Otieno Ooko, Kavengi Kitonga, Brian Muhiya, Catherine Gitau, Antony Ndolo, Lilian D. A. Wanzare, Albert Njoroge Kahira, and Ronald Tombe. 2024. State of nlp in kenya: A survey . <i>Preprint</i> , arXiv:2410.09948.	618	Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Jingyuan Ma, Rui Li, Heming Xia, Jingjing Xu, Zhiyong Wu, Tianyu Liu, Baobao Chang, Xu Sun, Lei Li, and Zhifang Sui. 2024. A survey on in-context learning . <i>Preprint</i> , arXiv:2301.00234.	619	670
620		621		671	
622		623		672	
624		625		673	
626	Rosana Ardila, Megan Branson, Kelly Davis, Michael Kohler, Josh Meyer, Michael Henretty, Reuben Morais, Lindsay Saunders, Francis Tyers, and Gregor Weber. 2020. Common voice: A massively-multilingual speech corpus . In <i>Proceedings of the Twelfth Language Resources and Evaluation Conference</i> , pages 4218–4222, Marseille, France. European Language Resources Association.	627	Bonaventure F. P. Dossou, Atnafu Lambebo Tonja, Oreen Yousuf, Salomey Osei, Abigail Oppong, Iyanoluwa Shode, Oluwabusayo Olufunke Awoyomi, and Chris Emezue. 2022. AfroLM: A self-active learning-based multilingual pretrained language model for 23 African languages . In <i>Proceedings of The Third Workshop on Simple and Efficient Natural Language Processing (SustaiNLP)</i> , pages 52–64, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.	628	675
629		630		676	
631		632		677	
633	Satanjeev Banerjee and Alon Lavie. 2005. METEOR: An automatic metric for MT evaluation with improved correlation with human judgments . In <i>Proceedings of the ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization</i> , pages 65–72, Ann Arbor, Michigan. Association for Computational Linguistics.	634	Ahmed El-Kishky, Vishrav Chaudhary, Francisco Guzmán, and Philipp Koehn. 2020. CCAligned: A massive collection of cross-lingual web-document pairs . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP 2020)</i> , pages 5960–5969, Online. Association for Computational Linguistics.	635	685
636		637		686	
638		639		687	
640		641		688	
642	Marta Bañón, Pinzhen Chen, Barry Haddow, Kenneth Heafield, Hieu Hoang, Miquel Esplà-Gomis, Mikel L. Forcada, Amir Kamran, Faheem Kirefu, Philipp Koehn, Sergio Ortiz Rojas, Leopoldo Pla Sempere, Gema Ramírez-Sánchez, Elsa Sarriàs, Marek Strelec, Brian Thompson, William Waites, Dion Wiggins, and Jaume Zaragoza. 2020. ParaCrawl: Web-scale acquisition of parallel corpora . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 4555–4567, Online. Association for Computational Linguistics.	643	Ahmed El-Kishky, Adithya Renduchintala, James Cross, Francisco Guzmán, and Philipp Koehn. 2021. XLEnt: Mining a large cross-lingual entity dataset with lexical-semantic-phonetic word alignment . In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 10424–10430, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.	644	692
645		646		693	
647		648		694	
649		650		695	
651		652		696	
653	Guangsheng Bao, Zebin Ou, and Yue Zhang. 2023. GEMINI: Controlling the sentence-level summary style in abstractive text summarization . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 831–842, Singapore. Association for Computational Linguistics.	654	Chris Chinene Emezue, Sanchit Gandhi, Lewis Tunstall, Abubakar Abid, Josh Meyer, Quentin Lhoest, Pete Allen, Patrick Von Platen, Douwe Kiela, Yacine Jernite, Julien Chaumont, Merve Noyan, and Omar Sanseviero. 2023. AfroDigits: A community-driven spoken digit dataset for african languages . <i>Preprint</i> , arXiv:2303.12582.	655	704
656		657		705	
658		659		706	
660		661		707	
662		663		708	
664		665		709	
666		667		710	
667		668		711	
668		669		712	
669		670		713	
670		671		714	

ings of the Association for Computational Linguistics: EMNLP 2023	2650–2666	Singapore.	777
Association for Computational Linguistics.			778
Luke Gessler and Katharina von der Wense. 2024.	NLP for language documentation: Two reasons for the gap between theory and practice.	In Proceedings of the 4th Workshop on Natural Language Processing for Indigenous Languages of the Americas (AmericasNLP 2024), pages 1–6, Mexico City, Mexico.	779
Association for Computational Linguistics.			780
Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mittra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Al-lonsius, Daniel Song, Danielle Pintz, Danny Livshits, Danny Wyatt, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Francisco Guzmán, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Govind Thattai, Graeme Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Is-han Misra, Ivan Evtimov, Jack Zhang, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Jun-teng Jia, Kalyan Vasuden Alwala, Karthik Prasad, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini, Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Kushal Lakhotia, Lauren Rantala-Yeary, Laurens van der Maaten, Lawrence Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo, Lukas Blecher, Lukas Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Mannat Singh, Manohar Paluri, Marcin Kardas, Maria Tsimpoukelli, Mathew Oldham, Mathieu Rita, Maya Pavlova, Melanie Kam-badur, Mike Lewis, Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal, Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Ning Zhang, Olivier Duchenne, Onur Çelebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Va-sic, Peter Weng, Prajjwal Bhargava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong, Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnic, Roberta Raileanu, Rohan Maheswari,	Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ron-nie Polidor, Roshan Sumbaly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sa-hana Chennabasappa, Sanjay Singh, Sean Bell, Seo-hyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Shar-an Narang, Sharath Raparth, Sheng Shen, Shengye Wan, Shruti Bhosale, Shun Zhang, Simon Van-denhende, Soumya Batra, Spencer Whitman, Sten Sootla, Stephane Collot, Suchin Gururangan, Syd-ney Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas Scialom, Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh Ramanathan, Viktor Kerkez, Vincent Gonguet, Vir-ginie Do, Vish Vogeti, Vítor Albiero, Vladan Petro-vic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whi-ney Meers, Xavier Martinet, Xiaodong Wang, Xi-aofang Wang, Xiaoqing Ellen Tan, Xide Xia, Xin-feng Xie, Xuchao Jia, Xuewei Wang, Yaelle Gold-schlag, Yashesh Gaur, Yasmine Babaei, Yi Wen, Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao, Zacharie Delpierre Coudert, Zheng Yan, Zhengxing Chen, Zoe Papakipos, Aaditya Singh, Aayushi Sri-vastava, Abha Jain, Adam Kelsey, Adam Shajnfeld, Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand, Ajay Menon, Ajay Sharma, Alex Boesenberg, Alexei Baevski, Allie Feinstein, Amanda Kallet, Amit Sangani, Amos Teo, Anam Yunus, Andrei Lupa, Andres Alvarado, Andrew Caples, Andrew Gu, Andrew Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchan-dani, Annie Dong, Annie Franco, Anuj Goyal, Apara-jita Saraf, Arkabandhu Chowdhury, Ashley Gabriel, Ashwin Bharambe, Assaf Eisenman, Azadeh Yaz-dan, Beau James, Ben Maurer, Benjamin Leonhardi, Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Han-cock, Bram Wasti, Brandon Spence, Brani Stojkovic, Brian Gamido, Britt Montalvo, Carl Parker, Carly Burton, Catalina Mejia, Ce Liu, Changhan Wang, Changkyu Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai, Chris Tindal, Christoph Fe-ichtenhofer, Cynthia Gao, Damon Civin, Dana Beaty, Daniel Kreymer, Daniel Li, David Adkins, David Xu, Davide Testuggine, Delia David, Devi Parikh, Diana Liskovich, Didem Foss, Dingkang Wang, Duc Le, Dustin Holland, Edward Dowling, Eissa Jamil, Elaine Montgomery, Eleonora Presani, Emily Hahn, Emily Wood, Eric-Tuan Le, Erik Brinkman, Este-ban Arcaute, Evan Dunbar, Evan Smothers, Fei Sun, Felix Kreuk, Feng Tian, Filippos Kokkinos, Firat Ozgenel, Francesco Caggioni, Frank Kanayet, Frank Seide, Gabriela Medina Florez, Gabriella Schwarz, Gada Badeer, Georgia Swee, Gil Halpern, Grant Herman, Grigory Sizov, Guangyi, Zhang, Guna Lakshminarayanan, Hakan Inan, Hamid Shojanazeri, Han Zou, Hannah Wang, Hanwen Zha, Haroun Habeeb, Harrison Rudolph, Helen Suk, Henry As-pegren, Hunter Goldman, Hongyuan Zhan, Ibrahim Damlaj, Igor Molybog, Igor Tufanov, Ilias Leontiadis, Irina-Elena Veliche, Itai Gat, Jake Weissman, James Geboski, James Kohli, Janice Lam, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jen-nifer Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe	781	
			782
			783
			784
			785
			786
			787
			788
			789
			790
			791
			792
			793
			794
			795
			796
			797
			798
			799
			800
			801
			802
			803
			804
			805
			806
			807
			808
			809
			810
			811
			812
			813
			814
			815
			816
			817
			818
			819
			820
			821
			822
			823
			824
			825
			826
			827
			828
			829
			830
			831
			832
			833
			834
			835
			836
			837
			838
			839
			840

841	Cummings, Jon Carvill, Jon Shepard, Jonathan Mc-	904
842	Phie, Jonathan Torres, Josh Ginsburg, Junjie Wang,	905
843	Kai Wu, Kam Hou U, Karan Saxena, Kartikay Khan-	906
844	delwal, Katayoun Zand, Kathy Matosich, Kaushik	907
845	Veeraraghavan, Kelly Michelena, Keqian Li, Kir-	908
846	ran Jagadeesh, Kun Huang, Kunal Chawla, Kyle	
847	Huang, Lailin Chen, Lakshya Garg, Lavender A,	
848	Leandro Silva, Lee Bell, Lei Zhang, Liangpeng	
849	Guo, Licheng Yu, Liron Moshkovich, Luca Wehrst-	
850	edt, Madian Khabsa, Manav Avalani, Manish Bhatt,	
851	Martynas Mankus, Matan Hasson, Matthew Lennie,	
852	Matthias Reso, Maxim Groshev, Maxim Naumov,	
853	Maya Lathi, Meghan Keneally, Miao Liu, Michael L.	
854	Seltzer, Michal Valko, Michelle Restrepo, Mihir Pa-	
855	tel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark,	
856	Mike Macey, Mike Wang, Miquel Jubert Hermoso,	
857	Mo Metanat, Mohammad Rastegari, Munish Bansal,	
858	Nandhini Santhanam, Natascha Parks, Natasha	
859	White, Navyata Bawa, Nayan Singhal, Nick Egebo,	
860	Nicolas Usunier, Nikhil Mehta, Nikolay Pavlovich	
861	Laptev, Ning Dong, Norman Cheng, Oleg Chernoguz,	
862	Olivia Hart, Omkar Salpekar, Ozlem Kalinli, Parkin	
863	Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pe-	
864	dro Rittner, Philip Bontrager, Pierre Roux, Piotr	
865	Dollar, Polina Zvyagina, Prashant Ratanchandani,	
866	Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel	
867	Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu	
868	Nayani, Rahul Mitra, Rangaprabhu Parthasarathy,	
869	Raymond Li, Rebekkah Hogan, Robin Battey, Rocky	
870	Wang, Russ Howes, Ruty Rinott, Sachin Mehta,	
871	Sachin Siby, Sai Jayesh Bondu, Samyak Datta, Sara	
872	Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov,	
873	Satadru Pan, Saurabh Mahajan, Saurabh Verma,	
874	Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lind-	
875	say, Shaun Lindsay, Sheng Feng, Shenghao Lin,	
876	Shengxin Cindy Zha, Shishir Patil, Shiva Shankar,	
877	Shuqiang Zhang, Shuqiang Zhang, Sinong Wang,	
878	Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala,	
879	Stephanie Max, Stephen Chen, Steve Kehoe, Steve	
880	Satterfield, Sudarshan Govindaprasad, Sumit Gupta,	
881	Summer Deng, Sungmin Cho, Sunny Virk, Suraj	
882	Subramanian, Sy Choudhury, Sydney Goldman, Tal	
883	Remez, Tamar Glaser, Tamara Best, Thilo Koehler,	
884	Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim	
885	Matthews, Timothy Chou, Tzook Shaked, Varun	
886	Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai	
887	Mohan, Vinay Satish Kumar, Vishal Mangla, Vlad	
888	Ionescu, Vlad Poenaru, Vlad Tiberiu Mihailescu,	
889	Vladimir Ivanov, Wei Li, Wenchen Wang, Wen-	
890	wen Jiang, Wes Bouaziz, Will Constable, Xiaocheng	
891	Tang, Xiaojian Wu, Xiaolan Wang, Xilun Wu, Xinbo	
892	Gao, Yaniv Kleinman, Yanjun Chen, Ye Hu, Ye Jia,	
893	Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi,	
894	Youngjin Nam, Yu, Wang, Yu Zhao, Yuchen Hao,	
895	Yundi Qian, Yunlu Li, Yuzi He, Zach Rait, Zachary	
896	DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang,	
897	Zhiwei Zhao, and Zhiyu Ma. 2024. <i>The llama 3 herd</i>	
898	<i>of models</i> . Preprint, arXiv:2407.21783.	
899	Ayyoob Imani, Peiqin Lin, Amir Hossein Kargaran,	949
900	Silvia Severini, Masoud Jalili Sabet, Nora Kass-	950
901	ner, Chunlan Ma, Helmut Schmid, André Martins,	951
902	François Yvon, and Hinrich Schütze. 2023. <i>Glot500:</i>	952
903	<i>Scaling multilingual corpora and language models to</i>	953
904	500 languages	954
905	. In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 1082–1117, Toronto, Canada. Association for Computational Linguistics.	955
906		956
907		
908		
909	Sheriff Issaka, Zhaoyi Zhang, Mihir Heda, Keyi Wang,	909
910	Yinka Ajibola, Ryan DeMar, and Xuefeng Du. 2024.	910
911	<i>The ghanaian nlp landscape: A first look</i> . Preprint,	911
912	arXiv:2405.06818.	912
913	Pratik Joshi, Sebastin Santy, Amar Budhiraja, Kalika	913
914	Bali, and Monojit Choudhury. 2020. <i>The state and</i>	914
915	915	
916	916	
917	6282–6293, Online. Association for Computational	917
918	Linguistics.	918
919		919
920	David D. Laitin, Rajesh Ramachandran, and Stephen L.	920
921	Walter. 2019. <i>The legacy of colonial language poli-</i>	921
922	922	
923	923	
924	924	
925	272.	925
926	Yinhan Liu, Jiatao Gu, Naman Goyal, Xian Li, Sergey	926
927	Edunov, Marjan Ghazvininejad, Mike Lewis, and	927
928	Luke Zettlemoyer. 2020. <i>Multilingual denoising pre-</i>	928
929	929	
930	930	
931	931	
932	Alexandre Magueresse, Vincent Carles, and Evan	932
933	Heetderks. 2020. <i>Low-resource languages: A re-</i>	933
934	934	
935	arXiv:2006.07264.	935
936	Shamsuddeen Hassan Muhammad, Idris Abdulkumin,	936
937	Abinew Ali Ayele, Nedjma Ousidhoum, David Ife-	937
938	oluwa Adelani, Seid Muhie Yimam, Ibrahim Sa’id	938
939	Ahmad, Meriem Beloucif, Saif M. Mohammad, Se-	939
940	bastian Ruder, Oumaima Hourrane, Pavel Brazdil,	940
941	Felermino Dário Mário António Ali, Davis David,	941
942	Salomey Osei, Bello Shehu Bello, Falalu Ibrahim,	942
943	Tajuddeen Gwadabe, Samuel Rutunda, Tadesse Be-	943
944	lay, Wendimu Baye Messelle, Hailu Beshada Balcha,	944
945	Sisay Adugna Chala, Hagos Tesfahun Gebremichael,	945
946	Bernard Opoku, and Steven Arthur. 2023. <i>Afrisenti:</i>	946
947	947	
948	948	
949	Shamsuddeen Hassan Muhammad, David Ifeoluwa	949
950	Adelani, Ibrahim Said Ahmad, Idris Abdulkumin,	950
951	Bello Shehu Bello, Monojit Choudhury, Chris C.	951
952	Emezeue, Anuoluwapo Aremu, Saheed Abdul, and	952
953	Pavel Brazdil. 2022. <i>Naijasenti: A nigerian twitter</i>	953
954	954	
955	955	
956	956	
957	Benjamin Muller, Antonios Anastasopoulos, Benoît	957
958	Sagot, and Djamé Seddah. 2021. <i>When being un-</i>	958
959	959	
960	960	

961	In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 448–462, Online. Association for Computational Linguistics.	1021
962		1022
963		1023
964		1024
965		1025
966		1026
967		1027
968		1028
969		1029
970		1030
971		1031
972		1032
973		1033
974		1034
975		1035
976		1036
977		1037
978		1038
979		1039
980		1040
981		1041
982		1042
983		1043
984		1044
985		1045
986		1046
987		1047
988		1048
989		1049
990		1050
991		1051
992		1052
993		1053
994		1054
995		1055
996		1056
997		1057
998		1058
999		1059
1000		1060
1001		1061
1002		1062
1003		1063
1004		1064
1005		1065
1006		1066
1007		1067
1008		1068
1009		1069
1010		1070
1011		1071
1012		1072
1013		1073
1014		1074
1015		1075
1016		1076
1017		1077
1018		1078
1019		1079
1020		1080
	Simón Posada Fishman, Juston Forte, Isabella Fulford, Leo Gao, Elie Georges, Christian Gibson, Vik Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo-Lopes, Jonathan Gordon, Morgan Grafstein, Scott Gray, Ryan Greene, Joshua Gross, Shixiang Shane Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris, Yuchen He, Mike Heaton, Johannes Heidecke, Chris Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele, Brandon Houghton, Kenny Hsu, Shengli Hu, Xin Hu, Joost Huizinga, Shantanu Jain, Shawn Jain, Joanne Jang, Angela Jiang, Roger Jiang, Haozhun Jin, Denny Jin, Shino Jomoto, Billie Jonn, Hee woo Jun, Tomer Kaftan, Łukasz Kaiser, Ali Kamali, Ingmar Kanitscheider, Nitish Shirish Keskar, Tabarak Khan, Logan Kilpatrick, Jong Wook Kim, Christina Kim, Yongjik Kim, Jan Hendrik Kirchner, Jamie Kirov, Matt Knight, Daniel Kokotajlo, Łukasz Kondraciuk, Andrew Kondrich, Aris Konstantinidis, Kyle Kosic, Gretchen Krueger, Vishal Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan Leike, Jade Leung, Daniel Levy, Chak Ming Li, Rachel Lim, Molly Lin, Stephanie Lin, Mateusz Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue, Anna Makanju, Kim Malfacini, Sam Manning, Todor Markov, Yaniv Markovski, Bianca Martin, Katie Mayer, Andrew Mayne, Bob McGrew, Scott Mayer McKinney, Christine McLeavey, Paul McMillan, Jake McNeil, David Medina, Aalok Mehta, Jacob Menick, Luke Metz, Andrey Mishchenko, Pamela Mishkin, Vinnie Monaco, Evan Morikawa, Daniel Mossing, Tong Mu, Mira Murati, Oleg Murk, David Mély, Ashvin Nair, Reiichiro Nakano, Rajeev Nayak, Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh, Long Ouyang, Cullen O’Keefe, Jakub Pachocki, Alex Paino, Joe Palermo, Ashley Pantuliano, Giambattista Parascandolo, Joel Parish, Emy Parparita, Alex Passos, Mikhail Pavlov, Andrew Peng, Adam Perelman, Filipe de Avila Belbute Peres, Michael Petrov, Henrique Ponde de Oliveira Pinto, Michael Pokorny, Michelle Pokrass, Vitchyr H. Pong, Tolly Powell, Alethea Power, Boris Power, Elizabeth Proehl, Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh, Cameron Raymond, Francis Real, Kendra Rimbach, Carl Ross, Bob Rotsted, Henri Roussez, Nick Ryder, Mario Saltarelli, Ted Sanders, Shibani Santurkar, Girish Sastry, Heather Schmidt, David Schnurr, John Schulman, Daniel Selsam, Kyla Sheppard, Toki Sherbakov, Jessica Shieh, Sarah Shoker, Pranav Shyam, Szymon Sidor, Eric Sigler, Maddie Simens, Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin Sokolowsky, Yang Song, Natalie Staudacher, Felipe Petroski Such, Natalie Summers, Ilya Sutskever, Jie Tang, Nikolas Tezak, Madeleine B. Thompson, Phil Tillet, Amin Tootoonchian, Elizabeth Tseng, Preston Tuggle, Nick Turley, Jerry Tworek, Juan Felipe Cerón Uribe, Andrea Vallone, Arun Vijayvergiya, Chelsea Voss, Carroll Wainwright, Justin Jay Wang, Alvin Wang, Ben Wang, Jonathan Ward, Jason Wei, CJ Weinmann, Akila Welihinda, Peter Welinder, Jiayi Weng, Lilian Weng, Matt Wiethoff, Dave Willner, Clemens Winter, Samuel Wolrich, Hannah Wong, Lauren Workman, Sherwin Wu, Jeff Wu, Michael Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qiming Yuan, Wojciech Zaremba, Rowan Zellers, Chong	1081
		1082
		1083
		1084

1085	Zhang, Marvin Zhang, Shengjia Zhao, Tianhao Zheng, Juntang Zhuang, William Zhuk, and Barret Zoph. 2024. Gpt-4 technical report . <i>Preprint</i> , arXiv:2303.08774.	1142
1086		1143
1087		1144
1088		1145
1089		1146
1090		1147
1091		1148
1092		1149
1093		1150
1094		1151
1095		1152
1096		1153
1097		1154
1098		1155
1099		1156
1100		1157
1101		1158
1102		1159
1103		1160
1104		1161
1105		1162
1106		1163
1107		1164
1108		1165
1109		1166
1110		1167
1111		1168
1112		1169
1113		1170
1114		1171
1115		1172
1116		1173
1117		1174
1118		1175
1119		1176
1120		1177
1121		1178
1122		1179
1123		1180
1124		1181
1125		1182
1126		1183
1127		1184
1128		1185
1129		1186
1130		1187
1131		1188
1132		1189
1133		1190
1134		1191
1135		1192
1136		1193
1137		1194
1138		1195
1139		1196
1140		1197
1141		1198
1142		1199
1143		1200

1201	processing tasks. <i>Journal for Language Technology and Computational Linguistics</i> , 36(2):1–27.	1261
1202		1262
1203		1263
1204		1264
1205		1265
1206		1266
1207		1267
1208		1268
1209		1269
1210		1270
1211		1271
1212		1272
1213		1273
1214		1274
1215		1275
1216		1276
1217		1277
1218		1278
1219		1279
1220		1280
1221		1281
1222		1282
1223		1283
1224		1284
1225		1285
1226		1286
1227		1287
1228		1288
1229		1289
1230		1290
1231		1291
1232		1292
1233		1293
1234		1294
1235		1295
1236		1296
1237		1297
1238		1298
1239		1299
1240		1300
1241		1301
1242		1302
1243		1303
1244		1304
1245		1305
1246		1306
1247		1307
1248		1308
1249		1309
1250		1310
1251		1311
1252		1312
1253		1313
1254		1314
1255		1315
1256		1316
1257		1317
1258		1318
1259		1319
1260		1320
1261		1321
1262		1322
1263		1323
1264		1324

1325 Alfredo Palasciano, Alison Callahan, Anima Shukla,
 1326 Antonio Miranda-Escalada, Ayush Singh, Benjamin
 1327 Beilharz, Bo Wang, Caio Brito, Chenxi Zhou, Chirag
 1328 Jain, Chuxin Xu, Clémantine Fourrier, Daniel León
 1329 Periñán, Daniel Molano, Dian Yu, Enrique Manjavacas,
 1330 Fabio Barth, Florian Fuhrmann, Gabriel Altay,
 1331 Giyaseddin Bayrak, Gully Burns, Helena U. Vrabec,
 1332 Imane Bello, Ishani Dash, Jihyun Kang, John Giorgi,
 1333 Jonas Golde, Jose David Posada, Karthik Rangasai
 1334 Sivaraman, Lokesh Bulchandani, Lu Liu, Luisa
 1335 Shinzato, Madeleine Hahn de Bykhovetz, Maiko
 1336 Takeuchi, Marc Pàmies, Maria A Castillo, Marianna
 1337 Nezhurina, Mario Sänger, Matthias Samwald,
 1338 Michael Cullan, Michael Weinberg, Michiel De
 1339 Wolf, Mina Mihaljcic, Minna Liu, Moritz Freidank,
 1340 Myungsun Kang, Natasha Seelam, Nathan Dahlberg,
 1341 Nicholas Michio Broad, Nikolaus Muellner, Pascale
 1342 Fung, Patrick Haller, Ramya Chandrasekhar, Renata
 1343 Eisenberg, Robert Martin, Rodrigo Canalli, Rosaline
 1344 Su, Ruisi Su, Samuel Cahyawijaya, Samuele Garda,
 1345 Shlok S Deshmukh, Shubhanshu Mishra, Sid Kiblawi,
 1346 Simon Ott, Sinee Sang-aroon Siri, Srishti Kumar,
 1347 Stefan Schweter, Sushil Bharati, Tanmay Laud,
 1348 Théo Gigant, Tomoya Kainuma, Wojciech Kusa, Yannis
 1349 Labrak, Yash Shailesh Bajaj, Yash Venkatraman,
 1350 Yifan Xu, Yingxin Xu, Yu Xu, Zhe Tan, Zhongli
 1351 Xie, Zifan Ye, Mathilde Bras, Younes Belkada, and
 1352 Thomas Wolf. 2023. *Bloom: A 176b-parameter*
 1353 *open-access multilingual language model.* *Preprint*,
 1354 arXiv:2211.05100.

1355 Linting Xue, Noah Constant, Adam Roberts, Mihir Kale,
 1356 Rami Al-Rfou, Aditya Siddhant, Aditya Barua, and
 1357 Colin Raffel. 2021. *mT5: A massively multilingual*
 1358 *pre-trained text-to-text transformer.* In *Proceedings*
 1359 *of the 2021 Conference of the North American Chapter*
 1360 *of the Association for Computational Linguistics: Human*
 1361 *Language Technologies*, pages 483–498, Online. Association for Computational Linguistics.

1363 Wenhao Zhu, Hongyi Liu, Qingxiu Dong, Jingjing Xu,
 1364 Shujian Huang, Lingpeng Kong, Jiajun Chen, and
 1365 Lei Li. 2024. *Multilingual machine translation with*
 1366 *large language models: Empirical results and analysis.* *Preprint*, arXiv:2304.04675.

1368 Ahmet Üstün, Viraat Aryabumi, Zheng-Xin Yong, Wei-
 1369 Yin Ko, Daniel D’souza, Gbemileke Onilude, Neel
 1370 Bhandari, Shivalika Singh, Hui-Lee Ooi, Amr Kayid,
 1371 Freddie Vargus, Phil Blunsom, Shayne Longpre,
 1372 Niklas Müennighoff, Marzieh Fadaee, Julia Kreutzer,
 1373 and Sara Hooker. 2024. *Aya model: An instruction*
 1374 *finetuned open-access multilingual language model.* *Preprint*, arXiv:2402.07827.

Source	High-Resource	African	Ratio
GS (Link)	42,871	2,121	20.2
arXiv(Link)	539	16	33.7
IEEE (Link)	487	7	69.6
CORE(Link)	9,011	401	22.5

Table 4: Aggregate paper counts and ratios between high-resource and African languages (2020-2024). The ratio shows the disparity in research visibility, with higher numbers indicating greater inequality in representation. Search term: “multilingual” “X” “large language models”

A Appendices