

# Language Technology for All

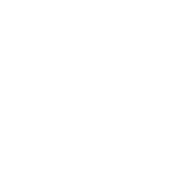
Industry Initiatives to Serve Low Resource Languages

Blaise Hylak

## About Myself



Blaise Hylak is a localization industry professional with six years of experience. A graduate of Villanova University, he rose through the ranks from intern to Program Manager. He holds a master's degree in Technical Communication and Localization from the University of Strasbourg. He speaks at language industry events nationwide, manages teams for local/national clients as a Program Manager at Core Alive Communications, Inc., and participates as a consultant/researcher to local, national, and international organizations requiring guidance on their current tech stacks, DEI, and how to enhance processes. He has been a member of ATA's Language Technology Division Leadership Council for the past three years.



## Paradigm Shift in the LOC Industry

"Equal access representation that is [...] pushed by governmental regulations [...] (make clients) more interested in ensuring [...] (that) access to all kinds of [...] baseline services or information is guaranteed" (Beccalotto, 2023).

Simona Beccalotto, Head of TAUS' Human Language Project Operations

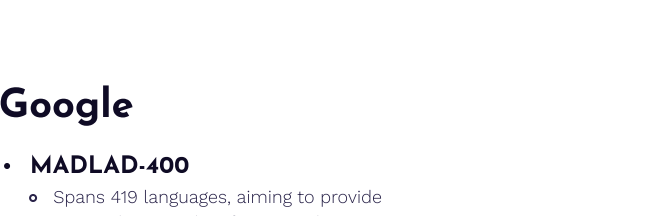
## Legislation Mandating Language Access

- United States of America (USA)
  - Americans with Disabilities Act (ADA)
  - Affordable Care Act (ACA)
  - Civil Rights Act of 1964, Title VI
  - Executive Order 13166 (2000)
- European Union (EU)
  - Open-source tools and models for Regional or Minority Languages
  - European Accessibility Act (EAA)

## Tech Bias Toward Low Resource Languages

- In a December 2023 exclusive interview for my thesis, Don DePalma of CSA Research provided an excellent graph produced by CSA Research that illustrates all training data by language group in the Common Crawl as of May 2023.
- 85.4% of the data is for European Languages
- English alone astonishingly accounts for nearly half of all training data

## Research Overview



## Principle Findings

Company	Key Initiatives
Meta	Meta is developing high-quality MT and NLP tools. <b>No Language Left Behind (NLLB)</b> , Seamless
Google	Google is developing high-quality datasets and developing a quality S2S tool. <b>MADLAD-400</b> , <b>Translator3</b>
Microsoft	Microsoft is developing high-quality data. <b>Project ELLORA</b>
IIITAUS	TAUS is developing high-quality data. <b>Human Language Project (HLP)</b>
Translators without Borders (TWB)	TWB is developing high-quality data. <b>Language Data Initiative</b>

## Meta

- No Language Left Behind (NLLB)**
  - Uses Human-Translated Datasets and Tools to Create Large Bilingual Datasets
  - Achieved a 44% improvement in BLEU scores, advancing the goal of a universal translation system
  - Open-sourced tools and models for wider community use

## Google

- MADLAD-400**
  - Spans 419 languages, aiming to provide comprehensive data for MT and NLP research
  - Includes 3 trillion clean tokens and 100 billion words, with a focus on LRLs
  - Models trained on MADLAD-400 have shown competitive performance with larger models

## Microsoft

- Project ELLORA**
  - Prioritizes Indian Languages:** This initiative focuses on Indian languages with limited digital presence, such as Gondi, Mundari, and Idu Mishmi
  - Partnerships in Data Collection:** Microsoft collaborates with local communities to gather and preserve language data, ensuring cultural sensitivity and accuracy
  - Tailored Digital Resources:** They develop digital dictionaries, translation services, and educational tools specifically for these languages, fostering digital inclusion and language preservation

## TAUS

- Human Language Project (HLP)**
  - Focuses on creating data for machine translation (MT) and speech-to-speech translation (S2S)
  - Involves crowd-sourced data collection from diverse global communities
  - Has expanded to cover over 30 languages across 20 countries

## Translators without Borders (CLEAR Global)

- Language Data Initiative**
  - 56 datasets covering almost 60 countries
    - Curated, cleaned, and reformatted the data to be more accessible for humanitarian and developmental purposes
  - Learn which languages are spoken where through the use of state-of-the-art language data globally with the Global Language Data Review
  - Pre-formatted and translated questions for language data collection

## Language Data for AI (LD4AI)

- Since there are no reliable estimates available for the LD4AI market, LD4AI is treated as an emerging sub-sector of the overall AI training data market.
- According to a market report by Grand View Research (2023), the AI training data market, including LD4AI, is expected to reach \$8.83 billion by 2030 and expand at a CAGR of 22.7% from 2023 to 2030\*

## DATAFORCE BY TRANSPERFECT

TransPerfect, the world's largest LSP, redirect their efforts towards AI data solutions.

Data Collection	Data Annotation	Data Relevance	Localize Chatbots
"Gather high-quality data for your unique model training/evaluation needs. Configuration options" (TransPerfect, 2024)	"DataForce accelerates your range of labeling processes with our range of human annotator services at scale" (TransPerfect, 2024)	"DataForce supports and improves your relevance models, accuracy, and recall to ensure the content you showcase to your customers is relevant, culturally acceptable, and geographically precise" (TransPerfect, 2024)	"Create chatbots that sound human and are culturally appropriate" (TransPerfect, 2024)

Data Moderation	Transcription	User Studies	Gen AI Training
"A multicultural and multilingual solution for your moderation needs" (TransPerfect, 2024)	"Scale speech and audio recognition capabilities with DataForce" (TransPerfect, 2024)	"DataForce utilizes its global footprint in over 46 countries to build your personal-level experience and collect the data you need through piloted, situational, and custom user studies" (TransPerfect, 2024)	"Whether you are developing new foundational models, such as LLMs, or customizing an existing model for a new use case, DataForce tailors solutions that address the unique data challenges your organization faces" (TransPerfect, 2024)

## LLMs for LRLs?

- A study found that LLMs, including ChatGPT, struggle significantly with LRLs compared to high-resource languages
- New approaches like Linguistically-Diverse Prompting (LDP) have been developed to help LLMs better handle LRLs by leveraging their strengths in high-resource languages, particularly English
- Despite advancements, LLMs continue to underperform in low-resource settings, often failing to access traditional machine learning models in tasks like machine translation and named-entity recognition
- While there have been technical improvements, the performance gap between LRLs and high-resource languages for LLMs remains significant, highlighting the need for continued research and development

## LRLs Can Jailbreak GPT-4..

- Researchers in a study jailbroke GPT-4 by translating unclear English input into a low resource language and consequently inputting that output into GPT-4. In essence, requesting a "back translation" (English to a low resource language and then back to English)
- GPT-4 engages with the unsafe translated inputs and provides actionable items that can get the users towards their harmful goals 79% of the time\*

## Conclusion

Generating/gathering high quality data is essential for language technology initiatives, in regards to low resource languages, this is a particular challenge that must be addressed and budgeted.

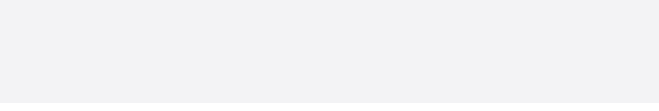
However, the matter of supporting low resource languages represents an interesting industry crossroads. Since governments worldwide are increasingly mandating language access via legislation, this also implies FUNDING. The financial incentives alone are a compelling reason to continue supporting initiatives for low resource languages besides mere compliance of the law.

## References

- Hylak, A. (2023, November 13). Microsoft Research project helps languages survive - and thrive. Microsoft. <https://www.microsoft.com/en-us/data/features/microsoft-research-project-helps-languages-survive-and-thrive/>
- Elasimoli, M. (2023, March 10). Google sheds light on 1000+ languages universal speech model. Slator. <https://www.slator.com/news/google-sheds-light-on-1000-languages-universal-speech-model/>
- Albarno, S. (2023, May 23). Meta challenges Whisper with massively multilingual speech launch. Slator. <https://www.slator.com/news/meta-challenges-whisper-with-massively-multilingual-speech-launch/>
- TransPerfect. (2024). DataForce by TransPerfect. <https://www.dataforce.ai/>
- Grand View Research. (2023). AI Training Dataset Market Size, Share and Trends Analysis Report. By Type (Text, Image/Video, Audio), By Vertical (IT, E-commerce, Government, Healthcare, BFSI), By Region, and Segment Forecasts, 2023 - 2030.
- Goyal, N., Li, Q., & Viriyals, O. (2021, June 23). Few-shot learning for machine translation. Google AI Blog. <https://ai.googleblog.com/2021/06/23/few-shot-learning-for-machine-translation.html>
- Johnson, M., Schuster, M., Le, Q. V., Krikun, M., Wu, Y., Chen, Z., & Dean, J. (2017). Google's multilingual neural machine translation system: Enabling zero-shot translation. Transactions of the Association for Computational Linguistics, 5, 339-351. <https://www.aclweb.org/anthology/P17-1024/>
- Zoph, B., & Le, Q. V. (2016). Neural architecture search with reinforcement learning. arXiv:1611.01027 [cs.LG].
- Sharoni, R., Johnson, M., & Firsirot, O. (2019). Machine search neural machine translation. In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Vol. 1, pp. 3874-3884).
- Sennrich, R., Haddow, B., & Birch, A. (2016). Improving neural machine translation models with monolingual data. In Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Vol. 1, pp. 86-96). <https://www.aclweb.org/anthology/P16-1009/>
- Fadawi, M., Bissazza, A., & Monz, C. (2017). Data augmentation for low-resource neural machine translation. arXiv:1705.00440
- Lampis, G., Cornaro, A., Denoyer, L., & Ferraro, M. A. (2018). Unsupervised machine translation using monolingual corpora only. In International Conference on Learning Representations. <https://arxiv.org/abs/1804.043>
- Settles, B. (2009). Active learning literature survey. University of Wisconsin, Madison. <https://cit.berkeley.edu/~settles/activelearning.pdf>
- Finn, C., Abbeel, P., & Levine, S. (2017). Model-agnostic meta-learning for fast adaptation of deep networks. In Proceedings of the 34th International Conference on Machine Learning (pp. 1126-1136).
- Heifler, J., Kamath, A., Ruckh, A., Chu, K., & Gurevych, I. (2020). AdapterFusion: Non-destructive, task-aware transfer learning. In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics (pp. 1152-1162). <https://www.aclweb.org/anthology/2020.acl-main/8/>
- Zaidan, O. F., Callison-Burch, C., & Poesio, M. (2011). Crowdsourcing annotation for machine translation: A tale of two approaches. In Proceedings of the Conference on Empirical Methods in Natural Language Processing (pp. 104-113). <https://www.aclweb.org/anthology/D11-1010/>
- CSA Research. "Percentage of Training Data by Language Group." Graph in Open Definition: "Investing with the Ethics of AI." World's megatrends 2023. <https://www.dataforce.ai/files/mega-trends-2023-information/investing-with-the-ethics-of-ai-2023/>
- Hong, Zheng-Yin, Cristina Menghini, and Stephen H. Bach. "Low-Resource Languages Jailbreak GPT-4." Submitted October 3, 2023. PDF. Accessed December 14, 2023. <https://arxiv.org/pdf/2310.0446.pdf>
- CLEAR Global. "Language Maps and Data." 2023. <https://www.clear-global.com/maps/>
- Lampis, G., Simona. Virtual interview by Blaise Hylak. December 20, 2023.
- DePalma, Donald. Virtual interview by Blaise Hylak. December 16, 2023.
- TAUS Human Language Project. Accessed January 14, 2024. <https://taus.com/>
- Kudurgatta, Sneha, Isaac Caswell, Bao Zhang, Xavier Garcia, Christopher A. Choquette-Choo, Katherine Lei, Derrick Xin, Aditya Kusupati, Romi Stella, Ankur Bapna, and Othan Firat. "MADLAD-400: A Multilingual And Document-Level Large Audited Dataset." Preprint. Accessed December 2023. <https://arxiv.org/abs/2310.14623>
- NLLB Team, Marta R. Costa-jussa, James Cross, Our Gilets, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elna Kabassi, Lam Dan, Daniel L. Kish, Jean-Marc Lavielle, Arina Sun, Shikhar Wang, Guillaume Wenzek, Al Youngblood, et al. "No Language Left Behind: Scaling Human-Centered Machine Translation." Accessed December 14, 2023. <https://arxiv.org/pdf/2207.04722.pdf>
- Johnson, M., Schuster, M., Le, Q. V., Krikun, M., Wu, Y., Chen, Z., & Manning, C. D. (2017). "Linguistically-Diverse Prompting (LDP) and Advancements in the Low-Resource Languages." In Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics: Tutorial Abstracts, June 2014. <https://arxiv.org/abs/1406.5592>
- Robison, Nathaniel R., Perez Ogayo, David R. Mortensen, and Graham Neubig. "ChatGPT MT: Competitive for High- (but not Low-) Resource Languages." September 14, 2023. PDF. Accessed December 14, 2023. <https://arxiv.org/pdf/2309.06231.pdf>
- Court, Sara, and Micha Elsner. "Shortcomings of LLMs for Low-Resource Translation: Retrieval and Understanding are Both the Problem." arXiv, June 2024. <https://arxiv.org/abs/2406.16269>
- Chhaywalya, Samrat, Holy Lovénia, and Pascale Fung. "LLMs Are Few-Shot In-Context Low-Resource Language Learners." arXiv, May 2024. <https://arxiv.org/abs/2403.08932>

Wanna continue the conversation? Connect with me on LinkedIn!

This PDF is a static version of a live Storydoc. For the most up-to-date version and a fully interactive experience, please visit the provided link.



Made with Storydoc