

Looking forward by looking back: Applying lessons from 20 years of African language technology

Martin Benjamin

École Polytechnique Fédérale de Lausanne
Lausanne, Switzerland
martin@kamusi.org

Mohomodou Houssouba

Songhay.org
Basel, Switzerland
mohomodou.houssouba@unibas.ch

Abstract

This paper takes a frank look at what has and has not been achieved in African language technology during the past two decades. Several questions are addressed: What was the status of technology for African languages 20 years ago? What were the major initiatives during that time? What were their successes and failures? What can we learn from these experiences? How does this inform the work that we are planning going forward? Examining in particular the history of Swahili, it is argued that technology projects have often achieved their expressed aims, but have collectively not significantly advanced the normalization of African languages as operable within the technical sphere, even while Africa has become blanketed with mobile technology. It is argued that future projects will succeed only by asserting the goal that technology of 2035 must be fully operational in users' primary languages, and gearing policy, funding, and individual project efforts toward gathering and deploying linguistic data for a large number of African languages to meet cutting-edge technologies as they emerge.

Keywords: HLT, localization, African languages, digital content, lexicography, Wikipedia, Swahili, equity, policy

1. African languages and technology in 1995

In 1995, when graphical browsers were first opening the Internet to dial-up home-computer users in some Northern countries, African languages had virtually no representation in the technical sphere. The computer support department of a major American university questioned when a project proposed for a large African language would ever be used by its native speakers. At a time when Tanzania had just one telephone per each 300-odd inhabitants (and only 2000 mobile phones),¹ and Mali had 17,164 phones in the entire country,² the notion that Africa would have more than two-thirds of a billion mobile devices twenty years later³ was unthinkable. Arguably, this failure to imagine the future led to chronic unpreparedness; while the number of devices and reach of networks available to Africans have grown exponentially, the ability to use those resources in and for even the largest of the continent's 2000 languages has barely gotten off the ground. In this paper, we review some of the major technological developments in support of African languages during that time,

discuss the obstacles and opportunities those experiences reveal, and make suggestions for building those lessons into work going forward.

ICT in 1995 was dominated by English, a consequence of the locations of the predominant early developers and mass markets. European and Asian languages improved their resources over the next two decades in response to market forces, government initiatives, and academic interest. It is difficult to remember that in 1995, rendering standard European diacritical characters such as “á” demanded contortions such as “á”. African characters such as Yoruba “ò” were literally impossible to type. With essentially no electronic wordlists for African languages, basic tools such as spellcheckers or OCR could not exist. While the full text search of known Internet pages that became possible with the release of WebCrawler in 1994 was in principle language neutral, African language content was almost entirely absent online and thus search was irrelevant. Before the release of the Netscape graphical web browser at the end of 1994, few organizations had even contemplated having a “homepage”, to say nothing of translated versions of their interface or content. Windows 95 was available to the brave⁴ for French, German, Italian, Russian, Arabic, Hebrew, Japanese, Korean, and Traditional Chinese, and a Swedish version of Word was spotted in Dar es Salaam by 1996, but neither Microsoft nor most other corporations attempted localization for Swahili or other large African languages until a decade ago at the

¹ Tanzania Communications Regulatory Authority, “Trends of Telephone Subscription”, <https://www.tcra.go.tz/index.php/25-statistics>, accessed 9/9/15

² The World Bank Databank, “Fixed Telephone Subscriptions”, <http://data.worldbank.org/indicator/IT.MLT.MAIN?page=3>, accessed 9/9/15

³ The World Bank Databank, “World Development Indicators: Power and communications, Table 5.11”, <http://wdi.worldbank.org/table/5.11>, accessed 9/9/15

⁴ Setting up Windows 95 for Multiple Languages, <http://www.ficorp.com/multi/>, accessed 9/9/15

month,¹³ with normal traffic in the vicinity of 700 visits per day, including both new and returning users. Many people certainly land elsewhere on the site via search engines or other links, but the landing page is the entry point for intentional searching for Swahili encyclopedia entries. The engagement between the Swahili world and its Wikipedia, orders of magnitude lower than the circulation of major Swahili newspapers, cannot be seen as notable. Reasons for the low uptake rates might include: (1) content still being too thin to be useful, as a substantial percentage of the articles are about locally irrelevant topics; (2) a lack of knowledge that the resource exists, as Wikimedia has done very little publicity, and no education ministry has worked to promote its use or production as part of the curriculum; (3) low access - students using one of the few computers in a school (if there are any at all) are not at leisure to peruse an online knowledge base, and few homes in East Africa have permanent network subscriptions; (4) people have not adopted the habit of going online to seek information in any language, which will undoubtedly change for some in the next few years; (5) the information that is available on the English Wikipedia is much more extensive, and the Kenyans and Tanzanians who are educated enough to be using a technology resource are likely to have the ability to read English at some level. Reasons for low editorship could include all of the above, added to the lack of awareness that Wikipedia can be edited, poor roadmaps to editing features, or frustrating experiences with the system.¹⁴ To that, one must add motivation: it is the rare user who donates content creation in any language, such as contributing a user review on TripAdvisor, because most people feel their time is better spent doing things other than providing bits of information to others.

A common thread among the efforts mentioned above is that they have been shipped without much attention (except to the extent that Google offers users the option to switch to Swahili when they log in from East Africa) to engaging the end user on their own terms. For example, Firefox could have been distributed on CDs that vendors would be encouraged to “pirate”, rather than quietly announced as available for download. Instead, the projects generally focus on the technical work, and leave consumer uptake to the wind. This is a matter of both the focus of the teams involved, and that the limited budgets of the projects dry up by the point of product delivery. Given that the potential market is largely unaware of the existence of the products, much less how to acquire and make use

!!

¹³ “Wikipedia Article Traffic Statistics”, <http://stats.grok.se/sw/201505/Mwanzo>, accessed 16/9/2015.

¹⁴ The new WYSIWYG editing tool is much more user-friendly than learning wiki markup, but still involves a substantial amount of learning, without evident tutorials, that would intimidate most novices.!

of them, it should not be too surprising that they have not become widely popular.

4.1 Future work - caught in the past or future shaping?

With twenty years of experience to guide us, we propose a new direction for African language technology. The first element should be, in fact, to have a twenty-year view toward what technical facilities will be available in Africa twenty years from now, and what is required starting now to make those functional for African languages. This view should have an assertive focus toward language equity; average consumers should expect to use their technology in their language, whether in Seattle or Soweto. Voice commands in Bambara for smart home technology, for example, should be considered a requirement, not a fantasy - if a product works in English and Chinese but not Akan or Malagasy, it is broken and not fit for sale in Ghana or Madagascar.

Such a declaration is, of course, much easier to state than to accomplish, involving both political commitment, financial investment, and concerted technological development. By and large, African states have rarely led the charge. The intergovernmental institution, ACALAN, has a small budget and a limited voice. Thus, ACALAN articulates goals for technological milestones¹⁵ that it may not be in a position to implement unless it can garner influential support from politics and civil society. In this regard, corporations, philanthropies, and international agencies should be brought into the conversation. To date, most organizations that should recognize language equity as a key to reaching economic and social goals hold, at best, to the notion that some tools can be localized to a few major African languages. Thus, Facebook’s founder can go to India and assert that language is important in reaching billions of potential subscribers in developing countries,¹⁶ but leave that work to unpaid volunteers while investing billions in such technologies as virtual reality gaming.¹⁷ International organizations, meanwhile, do not have language on their agenda, because that falls out of the scope usually considered important for successful activities in fields such as health or environmental conservation. Involving bilateral, multilateral and non-governmental organizations in the production of

!!

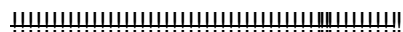
¹⁵ Project on African Languages and Cyberspace, <http://www.acalan.org/eng/projects/cyberspace.php>

¹⁶ Mark Zuckerberg Facebook post, 9 October 2014, <https://www.facebook.com/zuck/posts/10101687201975511>

¹⁷ African language game interfaces, much less games tailored to African consumers, are not within the visual horizon of the industry. “What languages to localize your game into”, http://gamasutra.com/blogs/JacobStempniewicz/20150619/244998/What_languages_to_localize_your_game_into.php, 19 June 2015.!

language tools that may be of use to their communications or their beneficiaries can raise their awareness of the aid proper language technology can provide for their missions, as well as raising the needed funds.

The most effective way to institute this general call will be through specific actions. To this end, kamusi.org and songhay.org are embarking on new directions in response to the experiences of the past two decades. As with the other projects discussed above, these two digital lexicography initiatives have not yet lived up to their potential, either in terms of usage of published resources or contributions to their growth. In part, this stems from chronic financial shortfalls that make it difficult to maintain existing systems, much less invest in new development. Kamusi faces the particular technical challenge of Big Data, with now tens of millions of internal links that must be computed and available to search engine indexing without bringing down its underpowered server.¹⁸ The projects have also neglected publicity, on the dubious assumption that speakers will find their way to web resources for their languages; with Africa engaging the digital sphere with handheld devices, projects designed around the wired web do not find the audience where they reside. They have also been weak in encouraging their users to contribute. While Kamusi especially has always made it clear that data is open to editing, the project has not been aggressive about recruiting the public to contribute. Moreover, the editing interface is a web form that asks for a lot of different types of lexical information¹⁹; while no more difficult than booking a hotel room online, experience training college students in Burundi shows that the platform is intimidating, and too bulky given the slow connection speeds and frequent power outages many Africans confront. Consequently, the two dictionary projects have come together and are completely revising their approach to public interactions. Because most Africans engage through Facebook and mobile devices (Rivron 2012), software development is now focusing on those platforms. Rather than ask for masses of data about each word, participants are asked highly targeted questions about particular linguistic elements. A lot of attention will be given to fostering language communities, including diasporic populations who have good connectivity and the desire to give something back to the people in their homeland, with kasahorow.org organizing online user groups beginning with two dozen languages. Through intensive and repetitive public outreach, the project seeks to make the idea of contributing to and accessing linguistic data normative for currently less-resourced languages.

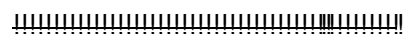


¹⁸ The Kamusi Big Data Beta, https://kamusi.org/big_data_beta

¹⁹ How Kamusi Works, <https://kamusi.org/how-kamusi-works>

On the other side, partnerships are needed between African language specialists and other technologists who are developing tools for the deployment of linguistic data. For example, linguistic data from one group can be used for morphological analysis by another, which can be the basis for computer assisted translation from a third that supplies live chat services from a fourth. Much NLP technology is rooted in European languages for three reasons: (1) that is where the funding is, (2) that is where the researchers are comfortable, and (3) that is where data is available. By resolving the third issue, African languages become attractive to researchers as an under-explored opportunity to make impactful advances. By resolving the second issue, showing that high-quality African language data can be produced at low cost and implemented within cutting technology by leading research groups, the idea of African languages as full participants in the technology sphere shifts from esoteric fantasy to evident fundability. The challenge now is to implement these optimistic plans in a financial and policy environment where African language equity does not factor as a realistic goal, despite a technically viable roadmap to success.

A review of academic presentations at international conferences regarding language technology shows that a small number of groups are hard at work on a few important topics for a few African languages, principally Arabic, Swahili, Yoruba, Amharic, Gikuyu, and several languages in South Africa. Using conference participation as a proxy for research support, the number of attendees focused on African languages is notably low, whether the conference topic is Computational Linguistics, Language Resources, or Language Documentation.²⁰ Frequently, African languages are best represented at such events by a few attendees in workshops for under-resourced or endangered languages. The major biennial meeting that focuses on African language technology is AfLaT (aflat.org), with numbers that showcase how small the research community is. Without a large cadre of researchers and institutions pushing forward on a number of topics and languages, it is hard to see a significant change in the status quo.



²⁰ Examples of proceedings from relevant events include: 52nd Annual Meeting of the Association for Computational Linguistics (2014), <http://acl2014.org/acl2014/mainconferenceprogram.html>; 9th International Conference on Language Resources and Evaluation (2014), <http://www.lrec-conf.org/proceedings/lrec2014/index.html>; 4th International Conference on Language Documentation and Conservation (2015), <http://scholarspace.manoa.hawaii.edu/handle/10125/35354>

Can African language technology emerge from the shallows? The past twenty years have largely missed preparing for the present - most Africans have no expectation of ever seeing the sorts of services in their languages that are routinely available today to most Europeans. The task is akin to designing urban light rail; plans must be in motion today to have an effective system two decades hence. If technologists continue with modest projects that address limited current objectives, we will find ourselves twenty years from now with African language technology that conquers some of the pressing issues of 2015. To address the needs of 2035 requires an insistent vision of a future where technology helps erase the language barriers it currently embeds. By stating the broad goal of technology for African languages reaching parity with where better-resourced languages in Europe and Asia will be twenty years hence, we can see the target. If this vision can be converted to policy and funding support for the development of linguistic data and tools throughout the continent, then perhaps we can not just see the target, but actually reach it.

References

- Badenhorst, J., Van Heerden C., Davel M. H., and Barnard E. (2011). "Collecting and evaluating speech recognition corpora for 11 South African languages", In: *Language Resources and Evaluation*, (45) 3, 289-309.
- Barnard, E., Davel, M., Van Huyssteen, G. (2010). "Speech technology for information access: a South African case study". AAAI Symposium on Artificial Intelligence, AAAI Spring Symposium Series, Stanford University, USA.
- Bearth. Th., Bonato, J., Geitlinger, K., Coray-Dapretto, L., Möhlig, W.J.G. and Olver, Th. (Eds.) (2009). *African Languages in Global Society*. Cologne: Rüdiger Köppe.
- Besha R., (2009). "Regional and local languages as resources of human development in the age of globalization," in Bearth et al. (1-13).
- Bosch, S. (2014). "Towards an Integrated E-Dictionary Application – The Case of an English to Zulu Dictionary of Possessives", In: *Proceedings of the XVI EURALEX International Congress: The User in Focus*, Bolzano (Italy), 739-748.
- Chiarcos, C., Fiedler, I., Grubic, M., Hartmann, K., Ritz, J., Schwarz, A., Zeldes, A., and Zimmermann, M. (2011). "Information structure in African languages: corpora and tools", In: *Language Resources and Evaluation*, (45) 3, 361-374.
- De Pauw, G., de Schryver G., and van de Loo, J., (2012). "Resource-Light Bantu Part-of-Speech Tagging, Proceedings of the workshop on Language technology for normalisation of less-resourced languages", In: *Proceedings of the Workshop on Language Technology for Normalisation of Less-Resourced Languages*, Istanbul, Turkey, 85-92.
- Diki-Kidiri, M. (2007). *Comment assurer la présence d'une langue dans le cyberspace*. Paris: UNESCO.
- Gasser, M. (2012). "Toward a Rule-Based System for English-Amharic Translation", In "*Proceedings of the Workshop on Language Technology for Normalisation of Less-Resourced Languages*, Istanbul, Turkey, 41-46.
- Gelas H., Abate S.T., Besacier L. and Pellegrino F. (2011). Evaluation of crowdsourcing transcriptions for African languages. In: *HLTD 2011*.
- Grover, A., van Huyssteen, G.B., and Pretorius, M. (2011). "The South African Human Language Technology Audit", in *Language Resources and Evaluation*, (45) 3, 271-288.
- HLTD (2011). *Proceedings of the Conference of Human Language Technology for Development (2011)*. Alexandria: Bibliotheca Alexandrina.
- Katushemerwe, F. (2013). *Computational Morphology and Bantu Language Learning: An Implementation for Runyakitara*, Ph.D. Thesis, Groningen University, Netherlands.
- Littell, P., Price, K., and Levin, L. (2014), "Morphological Parsing of Swahili using Crowdsourced Lexical Resources", In: *Proceedings of the Ninth International Conference on Language Resources and Evaluation*, Reykjavik (Iceland).
- Mali Ministry of Digital Economy (2014). "Plan Mali numérique 2020: Stratégie nationale de l'économie numérique". Bamako: Ministère de l'Économie Numérique, de l'Information et de la Communication. (Approved by Government 21 May, 2015).
- Mali Ministry of Education (2014). "Document de politique linguistique du Mali". Bamako: Ministère de l'Éducation Nationale. (Approved by Government 3 December, 2014).
- Mohomodou Houssouba (2007), "La concomitance du français et des langues nationales au Mali" in Françoise Argot-Dutard, ed, *Le français: des mots de chacun, une langue pour tous*. (3e Lyriades). Rennes: Presses universitaires de Rennes.
- Osborn, O. (2010). *African Languages in a Digital Age: Challenges and opportunities for indigenous language computing*. Cape Town: HSRC Press.
- Rivron, V. (2012). "The Use of Facebook by the Eton of Cameroon," In: *NET.LANG: Toward a Multilingual Cyberspace*, 160-165.
- Sinha, C. and Hyma, R. (2013). ICTs and Social Inclusion. In: Elder, L. et al (Eds), *Connecting ICTs to Development: The IDRC Experience*. London: Anthem.
- Vannini L. and Le Crosnier, H. (2012). *NET.LANG: Toward a Multilingual Cyberspace*. Caen: C&F Éditions and MAAYA Network.