

Machine Translation for Humanitarian Response

Alp Öktem - Translators without Borders

May 26th, 2021

GourMeT Open NLP Workshop 2021



The next twenty minutes

- Translators without Borders (TWB)
- Gamayun - Language equality initiative
 - Language data collection
 - Machine translation
 - Use-case oriented evaluation
- Questions, feedback, discussion

Information is power!

in the right language

Information is power!



Helping people to get vital information, and be heard, in the languages they speak and understand



**TRANSLATORS
WITHOUT BORDERS**

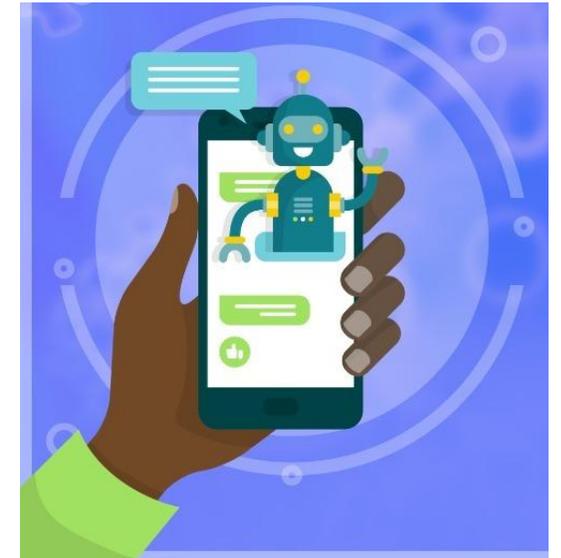
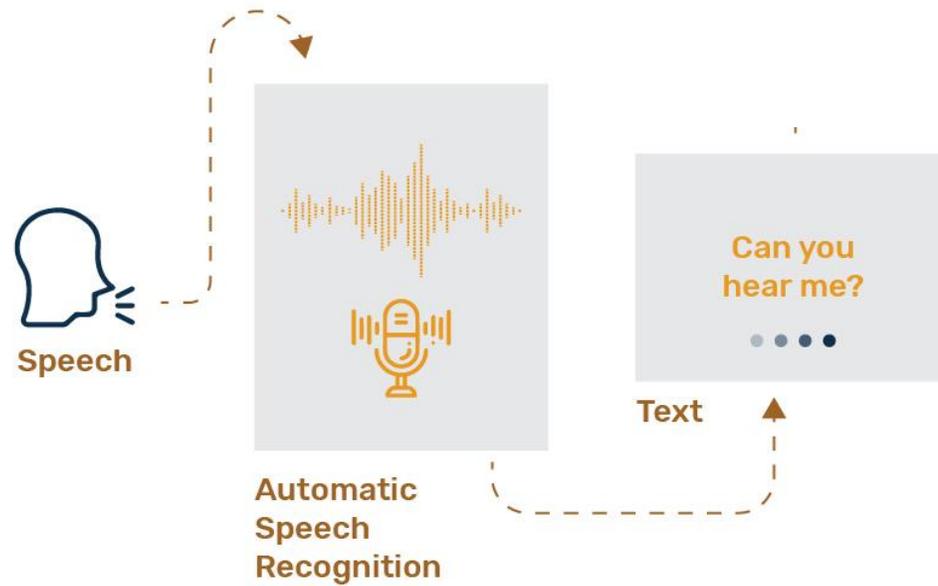


Gamayun

Language equality initiative



Machine Translation



Chatbots



Language data collection

Machine translation

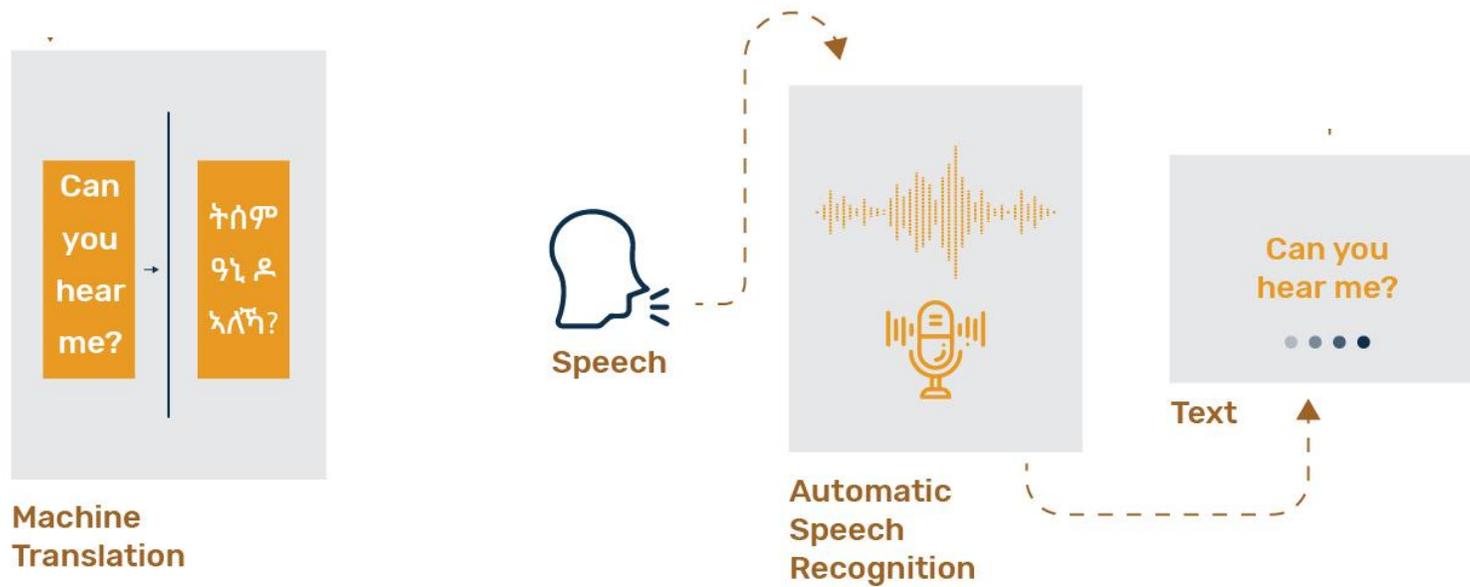
Use-case oriented evaluation



Language data collection



Language data collection



Machine Translation

Automatic Speech Recognition

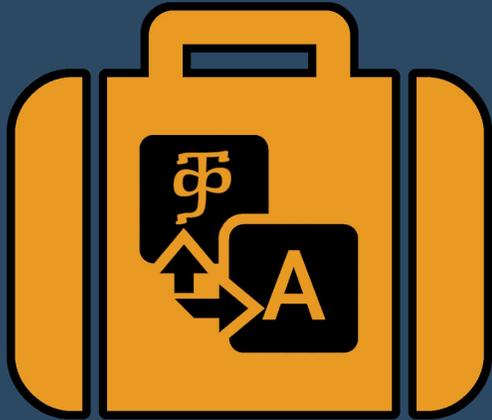
Text

Parallel data

Audio data

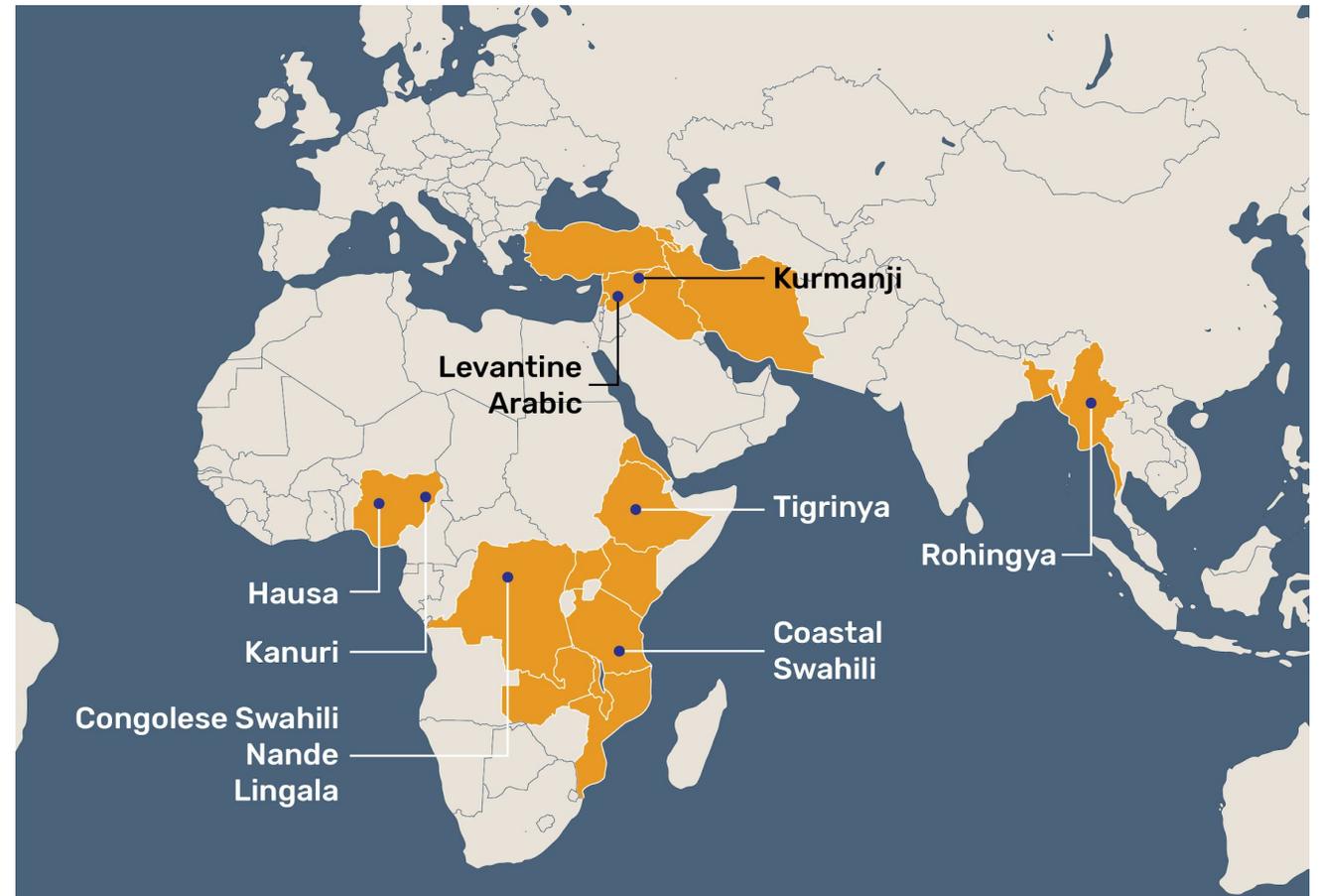
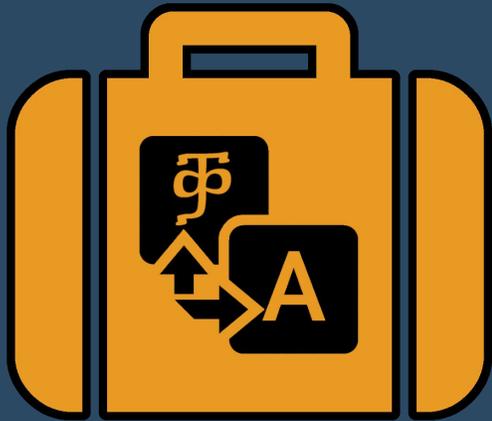


Gamayun Kits



- ❖ Language data for zero/low-resource languages
- ❖ General domain sentences
 - from Tatoeba corpus
 - in English, Spanish, French
- ❖ Translated for parallel data → MT
- ❖ Recorded for audio data → ASR, TTS

Gamayun Kits



❖ Currently: Hausa, Kanuri, Rohingya, Coastal/Congolesse Swahili, Nande

❖ <http://gamayun.translatorswb.org/data/>



TICO-19



SUPPORT TICO-19



Volunteer



Translate



Share

- ❖ **Translation Initiative for COVID-19**
- ❖ *Partners: George Mason University, Carnegie Mellon U., Johns Hopkins U., Google, Microsoft, Facebook, Amazon, Translated and Appen*
- ❖ **Covid-19 domain data for:**
 - Translation memories
 - Building specialized MT
 - Terminology



❖ COVID-19 domain source content from:

- Medical conversations
- PubMed articles
- News
- Travel restriction announcements
- Organization announcements
- Informative Articles



	Amharic	Arabic	Bengali	Kurdish Sorani	Dinka	English	Spanish LA	Farsi	French	Nigerian Fulfulde	Hausa	Hindi	I
Amharic (am)													1
Arabic (ar)							zipped .tmx		zipped .tmx			zipped .tmx	2
English (en)		zipped .tmx	zipped .tmx	zipped .tmx	zipped .tmx		zipped .tmx	zipped .tmx	zipped .tmx	zipped .tmx	zipped .tmx	zipped .tmx	2
Spanish LA (es-LA)		zipped .tmx							zipped .tmx			zipped .tmx	2
Farsi (fa)													
French (fr)		zipped .tmx					zipped .tmx					zipped .tmx	2
Hindi (hi)		zipped .tmx	zipped .tmx				zipped .tmx		zipped .tmx				2
Indonesian (id)		zipped .tmx					zipped .tmx		zipped .tmx			zipped .tmx	
Kurdish Kurmanji (ku)				zipped .tmx									
Portuguese Brazilian (pt-BR)		zipped .tmx					zipped .tmx		zipped .tmx			zipped .tmx	2

- 105 Translation Memories,
- 38 languages,
- Based on 3000 translated sentences
- <https://tico-19.github.io/>

Antonios Anastasopoulos, Alessandro Cattelan, Zi-Yi Dou, Marcello Federico, Christian Federmann, Dmitriy Genzel, Francisco Guzmán, Junjie Hu, Macduff Hughes, Philipp Koehn, Rosie Lazar, Will Lewis, Graham Neubig, Mengmeng Niu, Alp Öktem, Eric Paquin, Grace Tang, Sylwia Tur. **TICO-19: The Translation Initiative for COvid-19**. In: NLP COVID-19 Workshop (Part 2) @ EMNLP. 2020 November 19-20; Online.

Language data collection

Language data collection

Machine translation



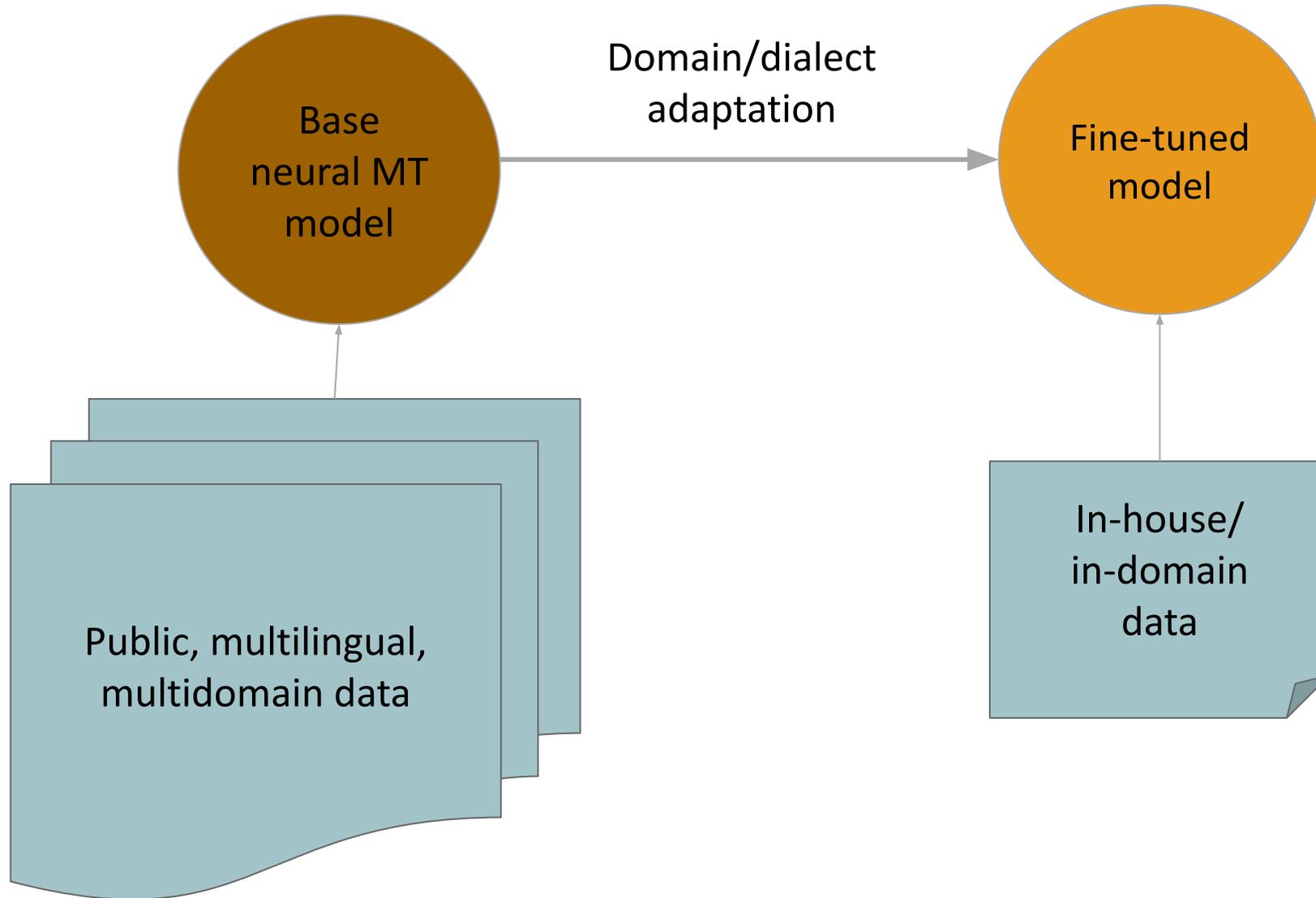
Machine translation

Objectives

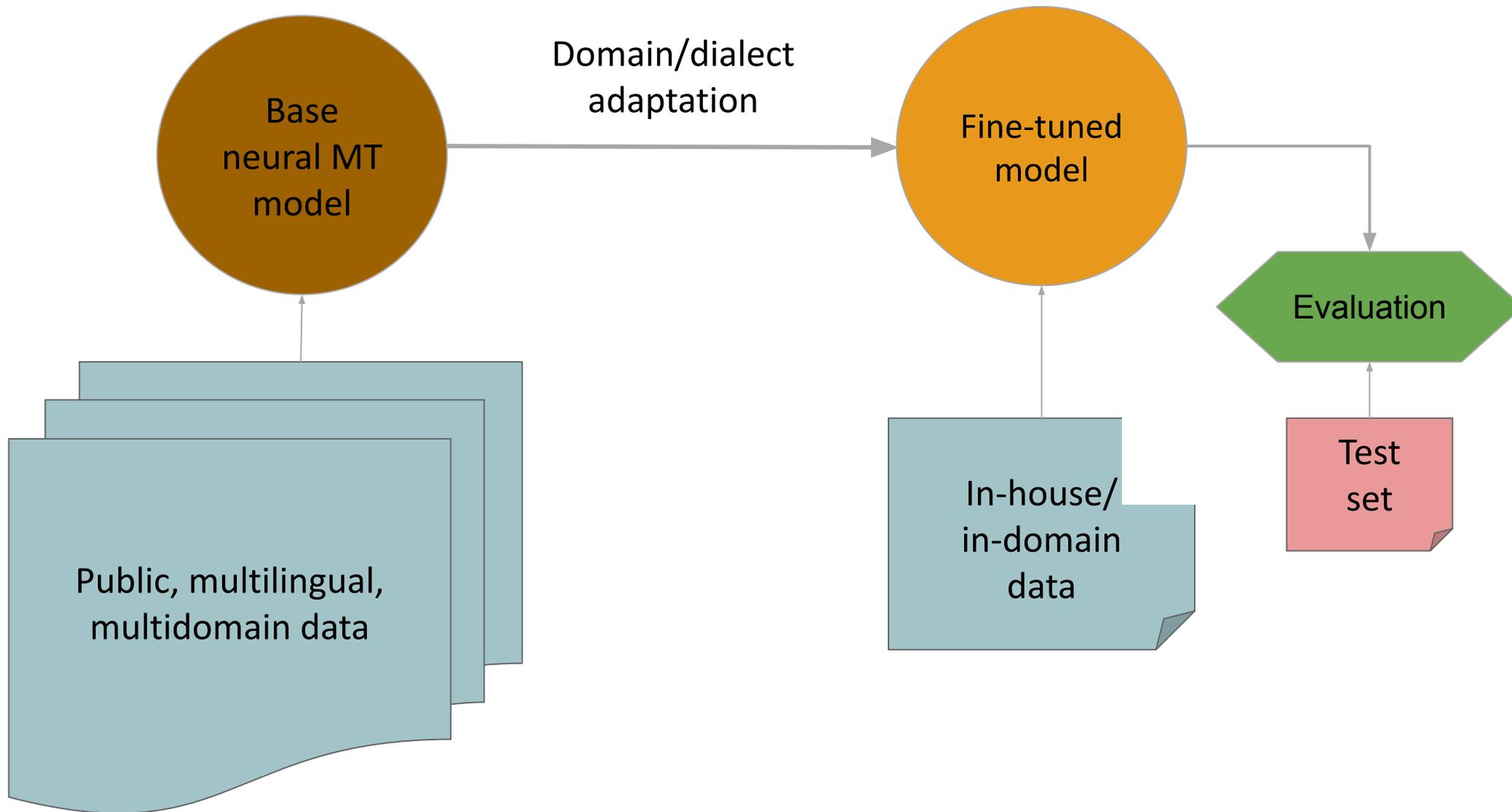
- Closing the technology gap for marginalized languages
- *Empowering* professional/non-professional translators
- Accelerate the flow of vital information in crises
- Automatic monitoring of humanitarian data

Pilot languages

- Levantine Arabic, Tigrinya, Congolese Swahili



General methodology



General methodology



Levantine Arabic

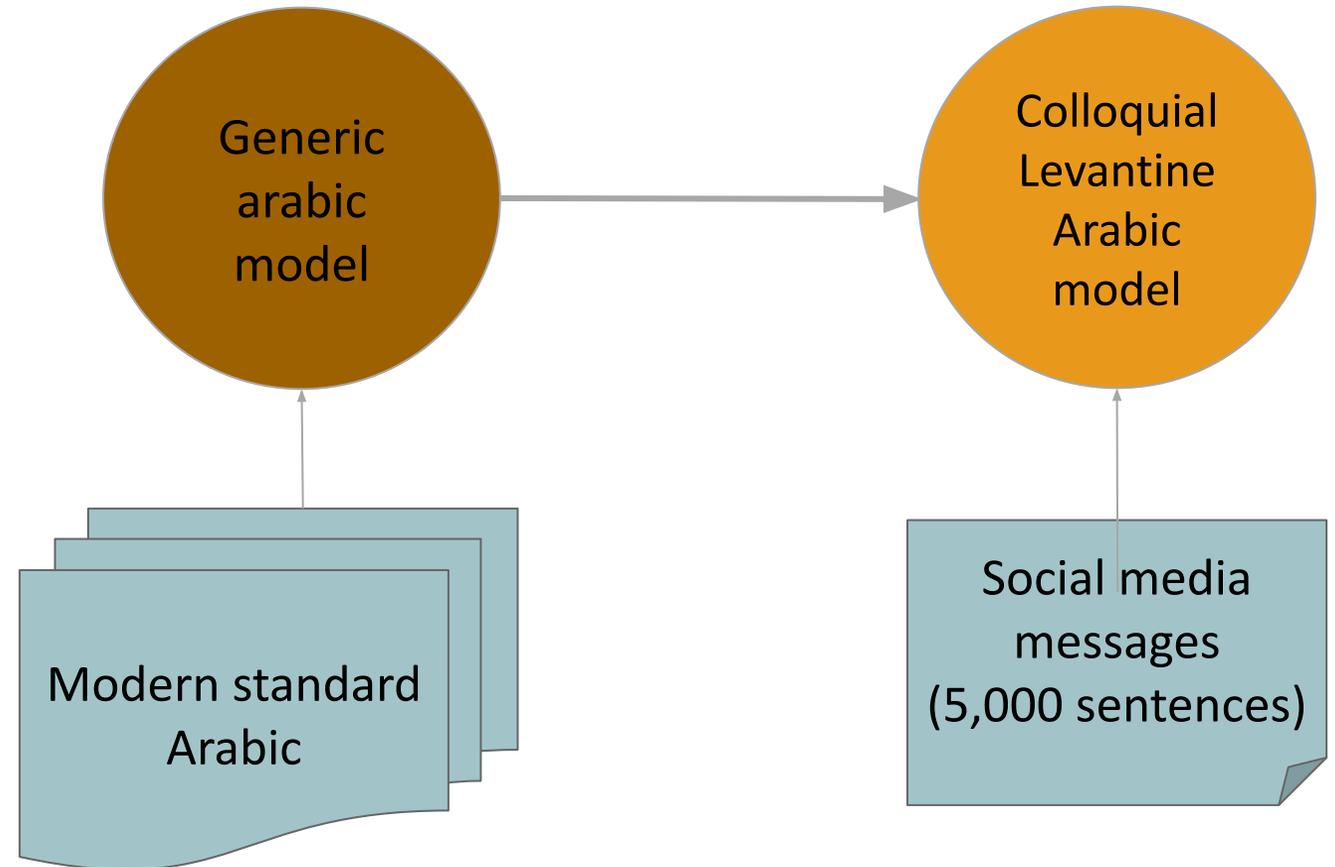
Social media monitoring



Alp Öktem, Muhannad Albayk Jaam, Eric DeLuca, Grace Tang. *Gamayun – Language Technology for Humanitarian Response*. In: 2020 IEEE Global Humanitarian Technology Conference (GHTC.) 2020 November 1: Virtual.

Levantine Arabic

Social media monitoring





Evaluation

- ❖ BLEU scoring with respect to reference translation
 - Test set: 200 messages
 - Domain adaptation: 19.5% -> 24.8%
 - Google MT: 21.2%

Levantine Arabic

Social media monitoring





Levantine Arabic

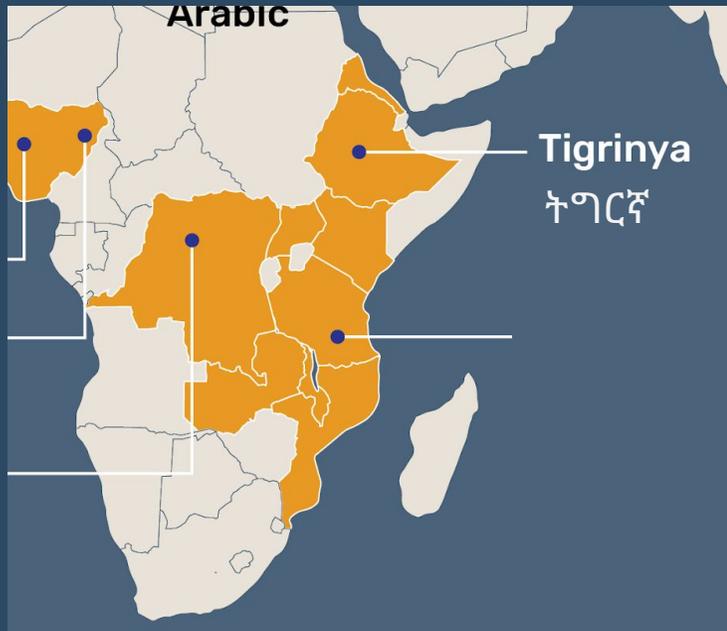
Social media monitoring



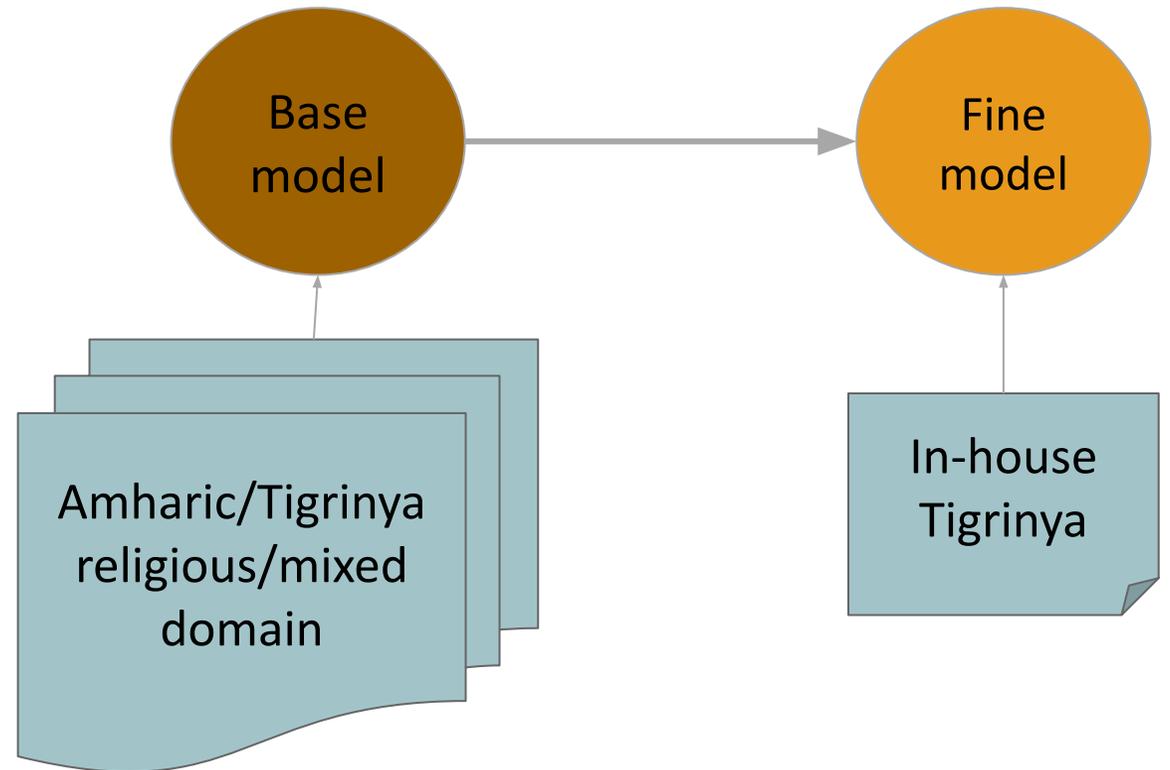
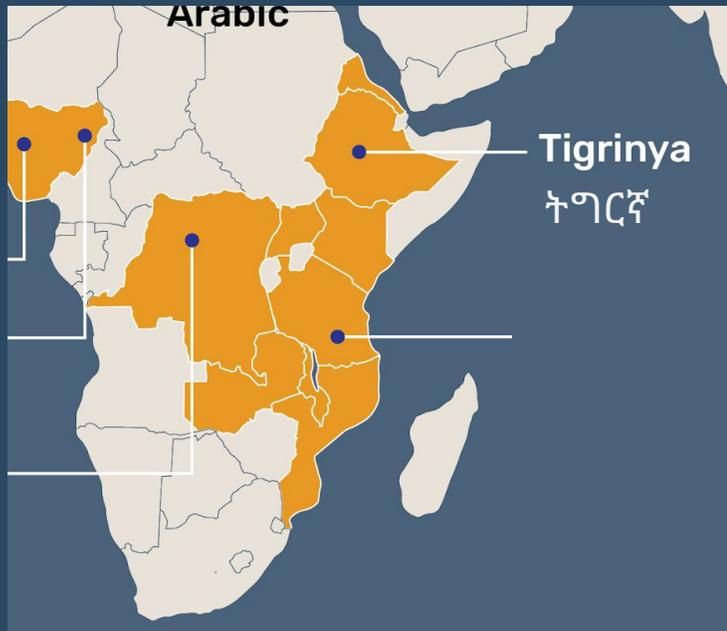
- ❖ “الغلاء انقطاع في المياه” - “Excessive water cuts and high prices”
 - TWB-MT: “high prices , water cuts”
 - Google MT: “Expensive interruption in the water”
- ❖ “الظروف المادية” - “Financial conditions”
 - TWB-MT: “Financial conditions”
 - Google MT: “Physical conditions”



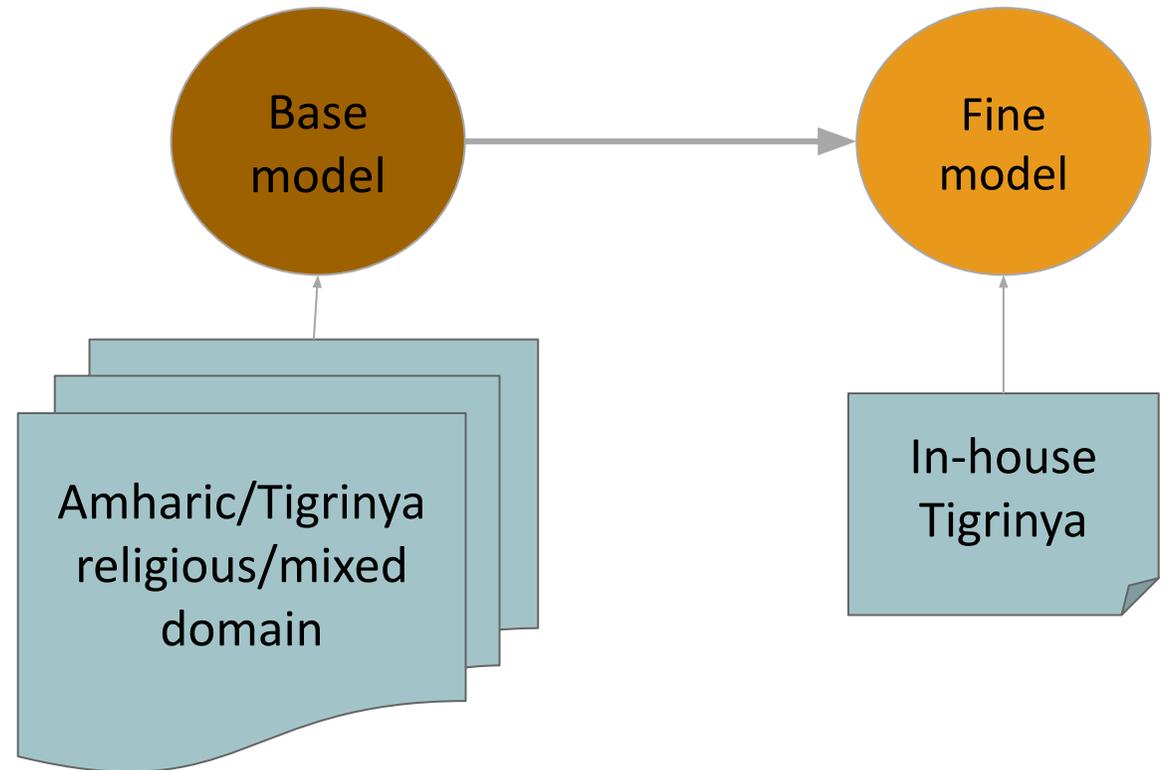
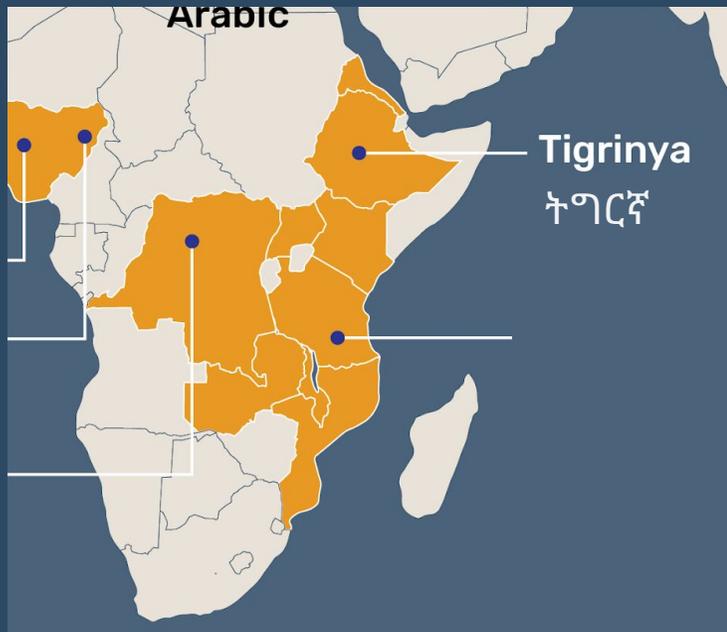
Tigrinya



Tigrinya



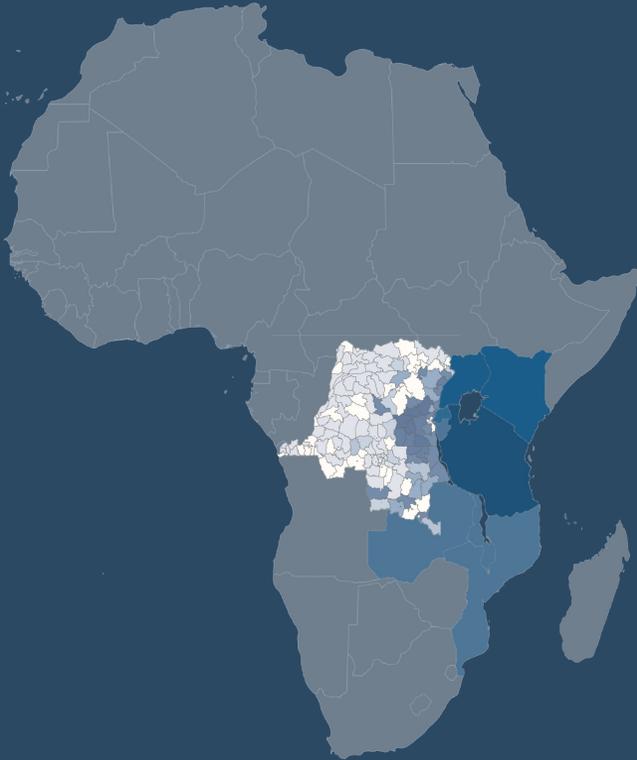
Tigrinya



- ❖ Tigrinya-to-English: 23.60% BLEU
 - *+1.32% with cross-lingual transfer*
- ❖ English-to-Tigrinya: 9.92% BLEU



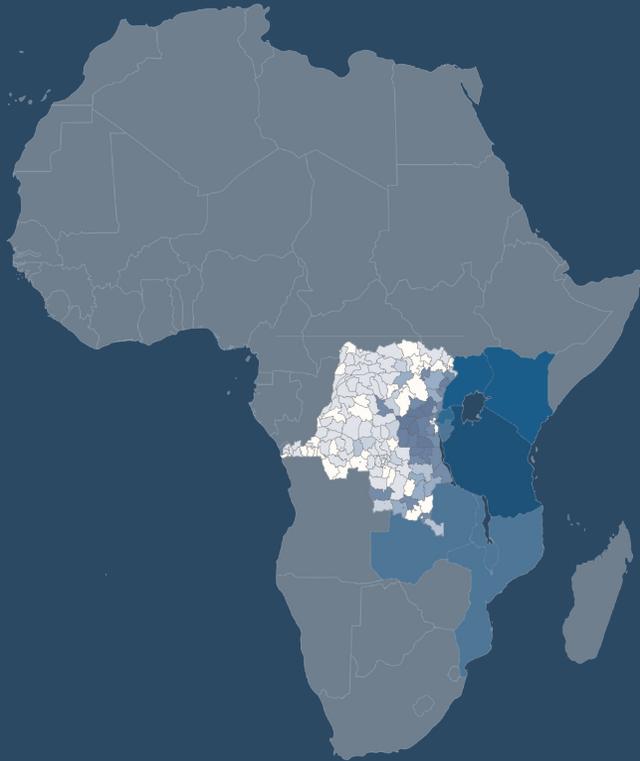
Congolese Swahili



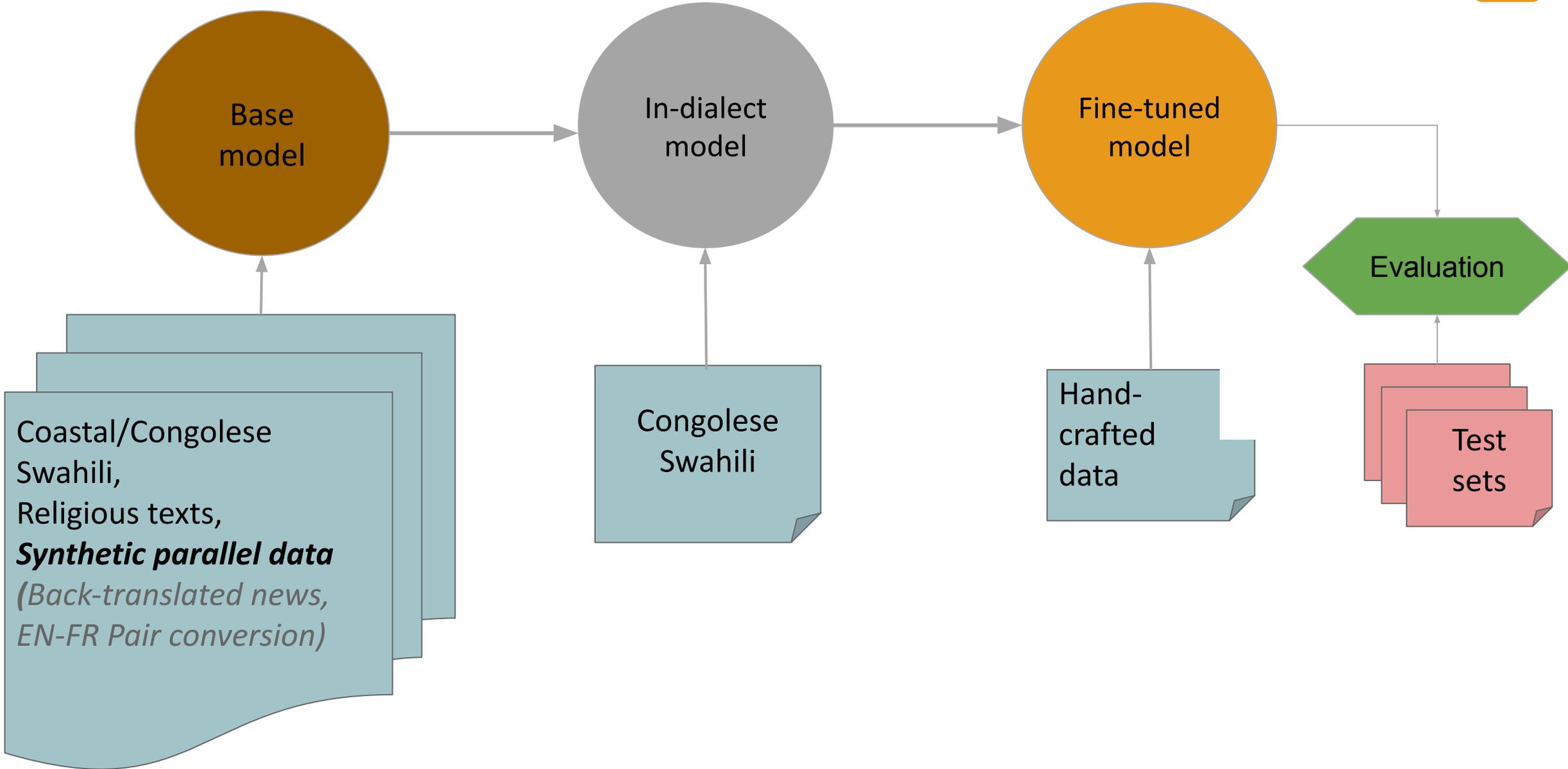
Alp Öktem, Eric DeLuca, Rodrigue Bashizi, Eric Paquin, Grace Tang. *Congolese Swahili Machine Translation for Humanitarian Response*. Accepted: Africa-NLP workshop at the 16th Conference of the European Chapter of the Association for Computational Linguistics (EACL2021.) 2021 April 19-23; Virtual.



Congolese Swahili



- ❖ *“Once the test is negative, the family can take care of the funeral themselves”*
 - **(Coastal)** Mara tu upimaji ukiwa hasi, familia inaweza kushughulikia mazishi wao wenyewe.
 - **(Congo)** Ikiwa tu vipimovinaonesha kama ni mtu ambaye hakuhakikishwa ku kuwa na ugonjwa, familia inaweza kufanya mazishi yenyewe.



	SWC to French	French to SWC
Only Congolese Swahili	17.7	13.0
Coastal Swahili mix	19.4 (+1.7)	15.0 (+2)
Synthetic data augmentation	20.1 (+0.5)	16.5 (+1.5)

Congolese Swahili (SWC) machine translation evaluation (BLEU score percentages)
on TICO-19 test set

Language data collection

Machine translation

Language data collection

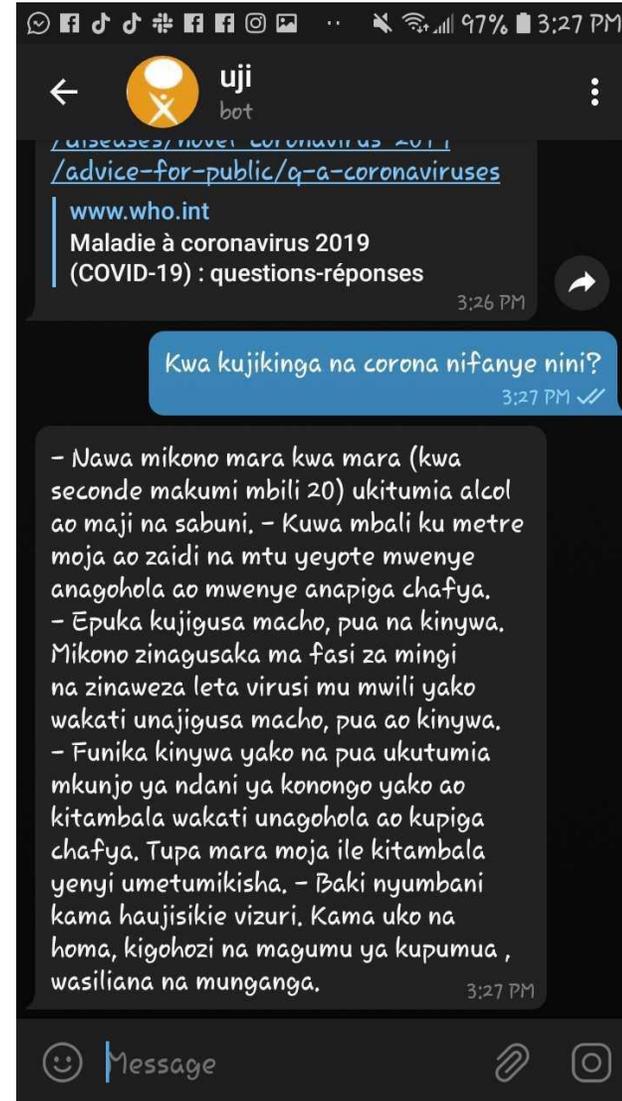
Machine translation

Use-case oriented evaluation

Use-case oriented human evaluation

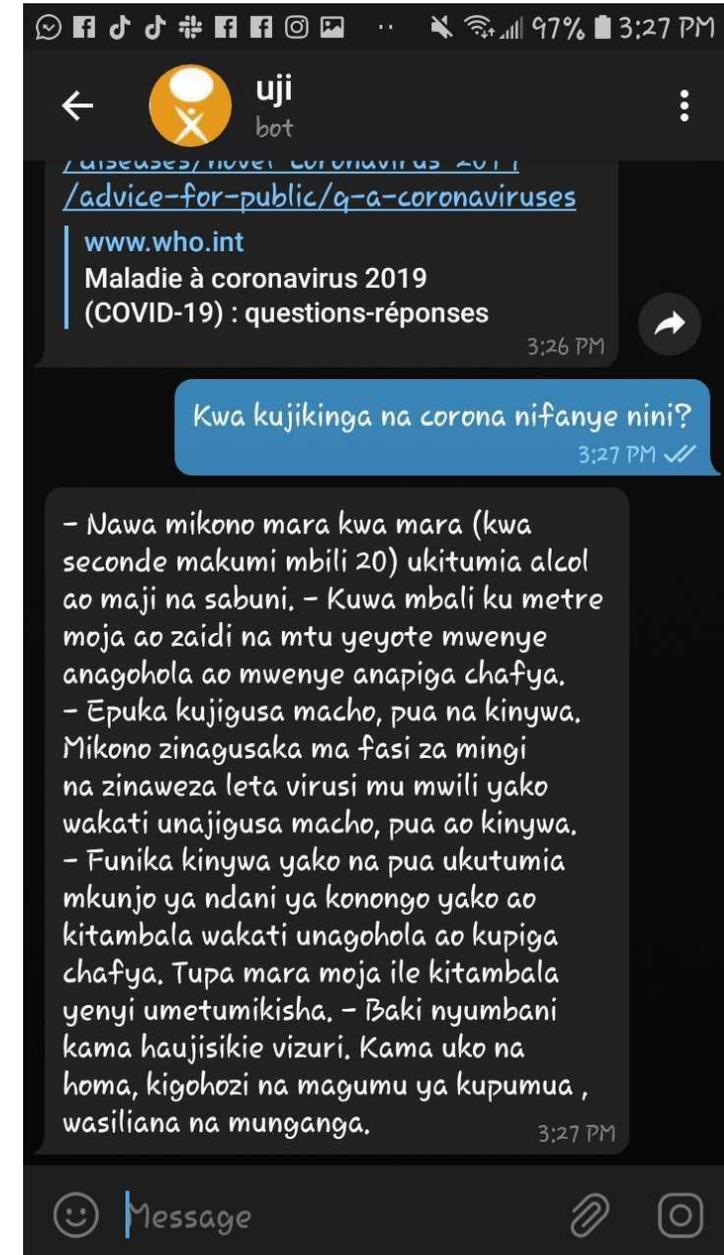


Talk to Uji - https://t.me/twb_uji_drc_bot



Use-case oriented human evaluation

- Post-editing (FRA->SWC)
- Direct assessment (SWC->FRA)



Use-case oriented human evaluation

- **Post-editing (FRA-→SWC)**
- Direct assessment (SWC-→FRA)

FR *Des cas de réinfection du COVID-19 ont été signalés mais sont rares. En général, la réinfection signifie qu'une personne a été infectée (est tombée malade) une fois, s'est rétablie, puis est redevenue infectée plus tard. Sur la base de ce que nous savons de virus similaires, certaines réinfections sont attendues.*

Raw-MT *Kesi za uambukizi wa COVID-19 zimeripotwa lakini hazipatikani kwa urahisi. Kwa ujumla, maambukizi inamaanisha kama mtu aliambukizwa (aliugua) mara moja, s' alirudishwa, na kisha akaambukizwa tena baadaye. Kutokana na kile tunachojua kuhusu virusi kama hivyo, maambukizo mengine yangali mbele.*

PE-MT *Kesi za **ma**ambukizi **ya** COVID-19 zimeripotwa lakini hazipatikani **mara nyingi**. Kwa ujumla, maambukizi inamaanisha kama mtu aliambukizwa (aliugua) mara moja, **akapona**, na kisha akaambukizwa tena baadaye. Kutokana na kile tunachojua kuhusu virusi kama hivyo, maambukizo mengine yangali mbele.*

BLEU	TER	ChrF
55.92	0.37	74.81



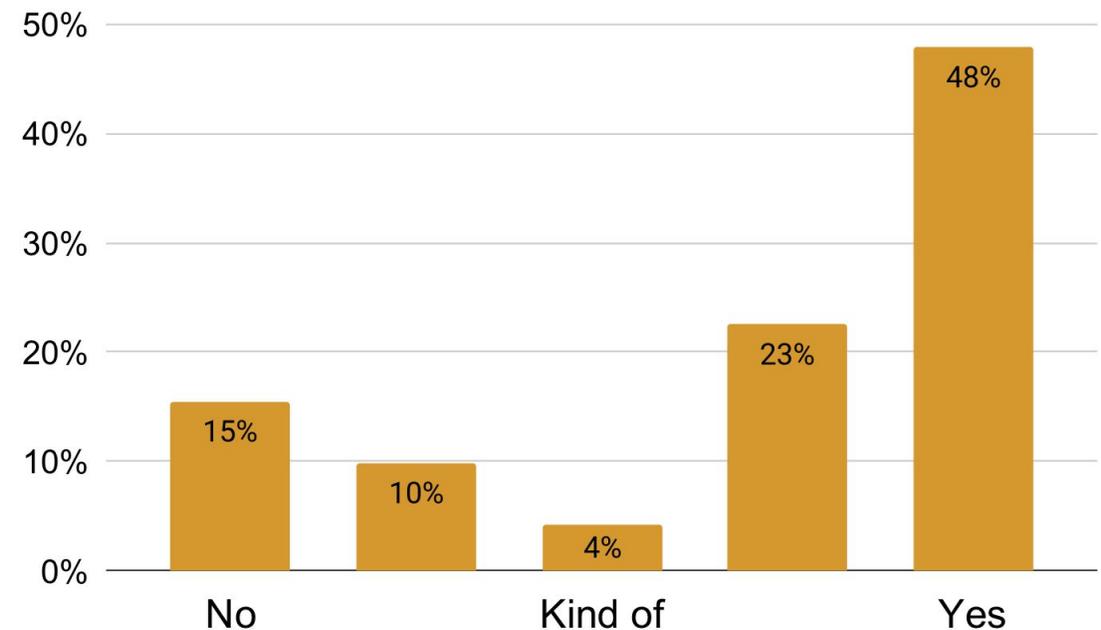
Use-case oriented human evaluation

- Post-editing (FRA-→SWC)
- **Direct assessment (SWC-→FRA)**

How good is the translation on a scale of 1 to 10?

➤ $6.3 \pm 3 / 10$

Does this translation convey the main meaning of the source text?





Language data collection

Machine translation

Use-case oriented evaluation



Resources

By admin · February 28, 2020

👁 463 🗨 0



Congolese Swahili speech corpus

5000 Audio samples recorded by 5 speakers. Sentences from Congolese Swahili mini kit.
Format: WAV
Size: 11 hours

 **Congolese Swahili audio mini-kit**
1 file(s) 144.47 MB [DOWNLOAD](#)

Tigrinya – English parallel text corpora

English sentences are sourced from Tatoeba repository and then translated into Tigrinya.
No. of sentences 5000

 **Gamayun Mini kit 5k Tigrinya – English**
1 file(s) 244.30 KB [DOWNLOAD](#)

Lingala – French parallel text corpora

French sentences are sourced from Tatoeba repository and then translated into Lingala.
No. of sentences 5000

 **Gamayun Mini kit 5k Lingala – French**
1 file(s) 292.46 KB [DOWNLOAD](#)

Congolese Swahili – French OpenNMT checkpoints

Bidirectional Congolese Swahili – French machine translation models

 **Congolese Swahili - French OpenNMT checkpoints**
1 file(s) 1,699.73 MB [DOWNLOAD](#)

Congolese Swahili – French parallel text corpora

French sentences are sourced from Tatoeba repository and then translated into Congolese Swahili.
No. of sentences 25305

 **Gamayun Mini kit 5k Congolese Swahili – French**
1 file(s) 296.98 KB [DOWNLOAD](#)

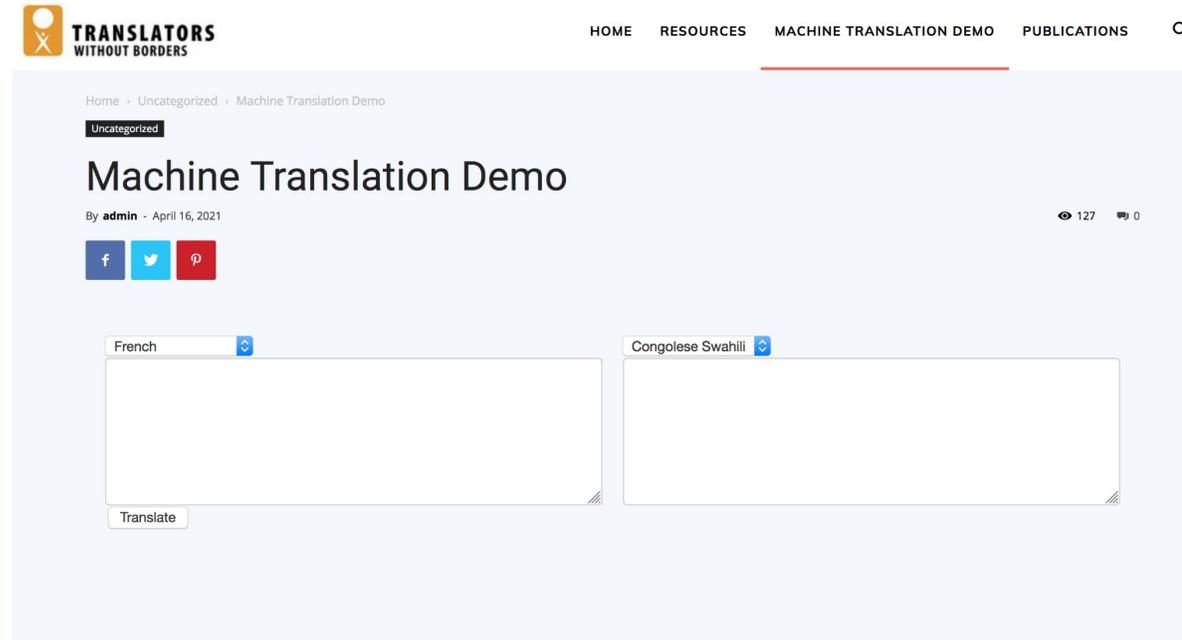
 **Gamayun Small kit 10k Congolese Swahili – French**
1 file(s) 445.24 KB [DOWNLOAD](#)

 **Gamayun Medium kit 15k (chunks 1-2) Congolese Swahili – French**
1 file(s) 401.84 KB [DOWNLOAD](#)

❖ Open data

- Hand-crafted training/testing data
- Model checkpoints

❖ Demo translator app



The screenshot shows the website header with navigation links: HOME, RESOURCES, MACHINE TRANSLATION DEMO, and PUBLICATIONS. The main content area is titled "Machine Translation Demo" and includes a "French" dropdown menu, a "Congolese Swahili" dropdown menu, and a "Translate" button. The page also displays the author "admin" and the date "April 16, 2021".

<http://gamayun.translatorswb.org>

Questions? Feedback?

<https://translatorswithoutborders.org>

alp@translatorswb.org

[@OktemAlp](https://twitter.com/OktemAlp)

[alpoktem.github.io](https://github.com/alpoktem)



**TRANSLATORS
WITHOUT BORDERS**