# A 1000-language Collaborative Universal Dictionary and Universal Translator
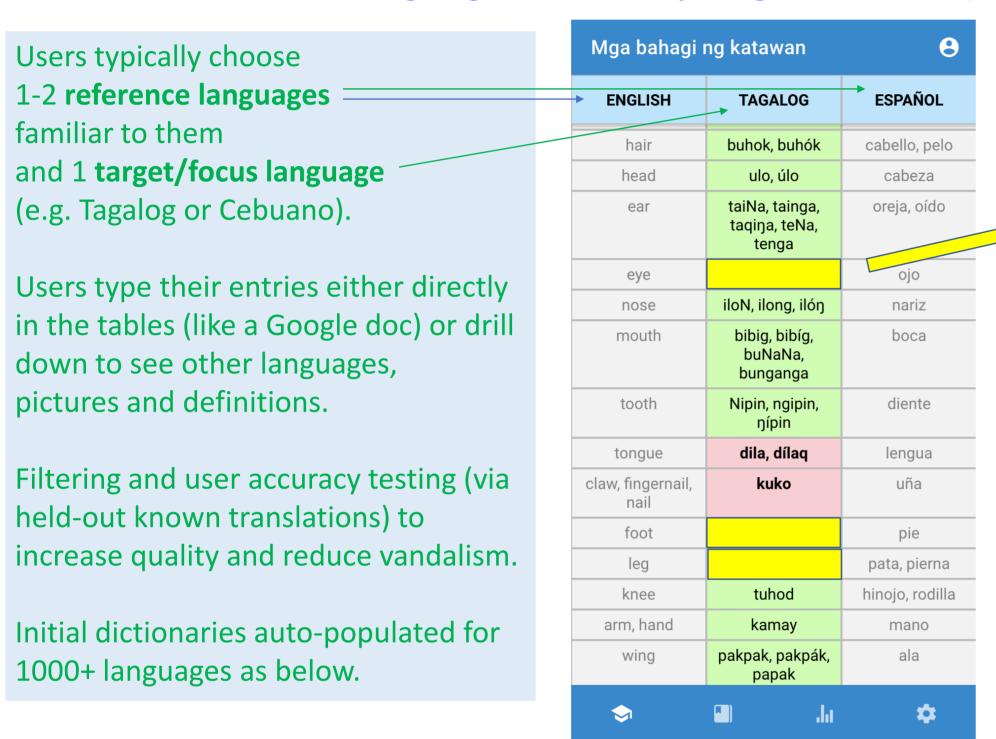
**David Yarowsky**, Arya D. McCarthy, Garrett Nicolai, Winston Wu, Aaron Mueller, Dylan Lewis, Yingqi Ding, Abhinav Nigam, Emre Ozgu,   Debanik Purkayastha, James Scharf and Kenneth Zheng

**Johns Hopkins University**

`yarowsky@jhu.edu`

We present JHU's Universal Dictionary and Universal Translator, covering 1000+ world languages in a broadly-accessible Android/IOS mobile phone and web-browser app, with 1,000,000+ planet-wide language pairs and 100's of under-resourced languages which have never had access to a substantial dictionary or machine translation capability. In addition to providing immediate access to a base vocabulary of 1500-20000 core vocabulary lemmas in all 1000+ languages, this novel app actively engages its users to contribute collaboratively to the universal dictionary in an easy-to-use and efficient way, with automatic suggestion of possible translations based on sound-shift transductions from related languages and pan-linguistic compositional constructions.

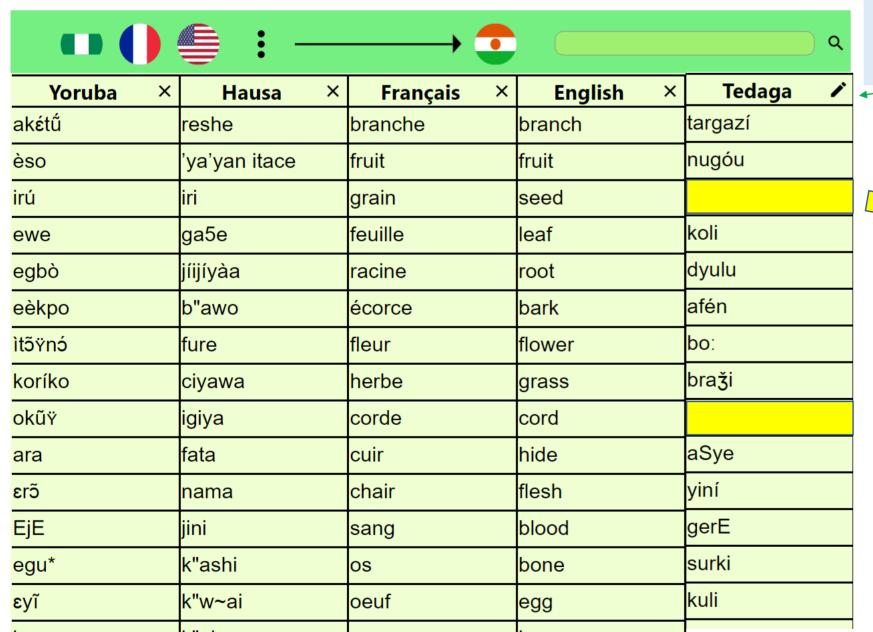## Collaborative 1000+ Language Dictionary Augmentation (Android App):

Users typically choose 1-2 **reference languages** familiar to them and 1 **target/focus language** (e.g. Tagalog or Cebuano).

Users type their entries either directly in the tables (like a Google doc) or drill down to see other languages, pictures and definitions.

Filtering and user accuracy testing (via held-out known translations) to increase quality and reduce vandalism.

Initial dictionaries auto-populated for 1000+ languages as below.



## Collaborative 1000+ Language Dictionary Augmentation (Web App):

**4 reference languages** familiar in Nigeria
**1 Target/focus language (Tedaga)**

**Drill-down window** (with photo etc.)



## Initial Dictionary Population (for 1000+ languages)

- Swadesh list
- Panlex
- Wiktionary
- Unimorph
- Universal Declaration of Human Rights
- Bible bitext (~3000+ core dictionary entries)

## Multiple models of word formation:

(1) cognate and sound shift models
(2) borrowing models
(3) derivational models (i.e. compositional translation of digger from dig +V:N(AGT) [and +V:N(INST)]
(4) universalized compositional models:

hospital = sick|sickness|disease  + house|place|institution

hospital = doctor|medicine        + house|place|institution

| da | syge\|hus | => hospital |
| sv | sjuk\|hus | => hospital |
| hu | kór\|házi | => hospital |
| hu | kór\|ház | => hospital |
| no | syke\|hus | => hospital |
| de | kranken\|haus | => hospital |
| nl | gast\|huis | => hospital |
| zh | 病院 | => hospital |

## Automatic compositional prediction:

| English | Gold | Lang | Hypotheses |
|---|---|---|---|
| skirt | yubka | azj | yubka, yubqa, jubka, jubqa, yubkɔ, Yubka, yubxa, übka, yübka, yubqə |
| fluorine | fluor, flüor | azj | ftor, ftar, flüor, ftor, faor, fdor, fluor, fdar, vlor, lor |
| food | gıda, qida | azj | gı, qı, qıda, gıda, gida, qısa, ğıda, ǧı, kida, qada |
| Greece | Юнанистан | kaz | Жунанастан, жунанастан, жананастан, Жананастан, Жунанастан, жунанастан, Жунанастан, жананастан, жунанастан, Жананастан |
| where | қайда | kaz | қайда, кайда, кажда, қажда, қайта, қаза, кайта, кőйда, кайта, шайда |
| cheese | сыр | kaz | сыр, сыл, сырт, сур, сілыр, сым, сір, сур, сыры, сір |
| wall | диvar | tat | дуал, двал, дуэл, дуғал, дуэл, дуаль, дгал, дваль, дуар, дуал |
| letter | harf, xäref, xat | tat | hät, tarf, xärf, xarf, Qät, hirf, harp, härp, harş, kärf |
| dove | yeni, yaña | tat | yaña, yaNa, yaNi, yañı, yañi, yene, yañge, yaNe, yeni, yange |
| weaving | dokuma, dokma | tuk | dokuma, dokamak, dokumak, dokuşmak, dokama, dokumaklyk, dokume, doku- mamak, dokulamak, dokuşma |
| cop | polis, polisiýa | tuk | polisiýa, politsiýa, polis, milisiýa, militsiýa, pilisiýa, polits, poliz, polisi, polys |
| shaman | şaman, şaman | tuk | şaman, shaman, saman, sheman, naman, kaman, şeman, şamen, sharman, haman |

(d) Turkic, 1-best  (e) Turkic, 10-best  (f) Turkic, MRR