

'Your United States was normal': Has translation tech really made language learning redundant?

November 22 2023, by Ingrid Piller



Credit: Unsplash/CC0 Public Domain

Every day, millions of people start the day by posting a greeting on social media. None of them expect to be arrested for their friendly morning ritual.

But that's exactly what happened to a Palestinian construction worker in 2017, when the caption "يصبحهم" ("good morning") on his Facebook selfie was [auto-translated](#) as "attack them."

A human Arabic speaker would have immediately [recognized "يصبحهم" as an informal way to say "good morning"](#). Not so AI. Machines are notoriously bad at dealing with variation, a key characteristic of all human languages.

With recent advances in automated translation, the belief is taking hold that humans, particularly English speakers, no longer need to learn other languages. Why bother with the effort when Google Translate and a host of other apps can do it for us?

In fact, some Anglophone universities are making precisely this argument to [dismantle their language programs](#).

Unfortunately, [language](#) technologies are nowhere near being able to replace human language skills and will not be able to do so in the foreseeable future because machine language learning and [human language](#) learning differ in fundamental ways.

How machines learn languages

For [machine translation](#), algorithms are trained on large amounts of texts to find the probabilities of different patterns of words. These texts can be both monolingual and bilingual.

Bilingual [training](#) data comes in the form of human-translated parallel

texts. These are almost always based on the standard version of the training language, excluding dialects and slang phrases, as in the example above.

Diversity is a characteristic of all human languages, but diversity is a problem for machines. For instance, "deadly" means "causing death" in most varieties of English, and that is what appears in the training data.

The [Australian meaning](#) of "excellent" (from Aboriginal English) puts a spanner in the works. If you input "[Deadly Awards](#)" into any translation app, what you'll get in your target language is the equivalent of "death-causing awards."

How machines store languages

The internal linguistic diversity of English, as of any other language, is accompanied by great diversity across languages. Each language does things differently.

Tense, number or gender, for example, need to be grammatically encoded in some languages but not in others. Translating the simple English statement "I am a student" into German requires the inclusion of a grammatical gender marking and so will either end up as "I am a male student" or "I am a female student."

Furthermore, some languages are spoken by many people, have powerful nation states behind them, and are well resourced. Others are not.

"Well resourced" in the context of machine learning means that large digital corpora of training data are available.

The lists of language options [offered by automated translation tools](#)—like the list of 133 languages in which Google Translate is currently

available—erase all these differences and suggest that each option is the same.

AI speaks English

Nothing could be further from the truth. English is in a class of its own, with over 90% of the training data behind large language models [being in English](#).

The remainder comes from a few dozen languages, in which data of varying sizes are available. The majority of the world's 6,000+ languages are simply missing in action. Apps for some of these are now being created from [models "pre-trained" on English](#), which further serves to cement the dominance of English.

One consequence of inequalities in the training data is that translations into English usually sound quite good because the app can draw both on bilingual and monolingual training data. This doesn't mean they are accurate: one recent study found about half of all questions in Vietnamese were [incorrectly auto-translated as statements](#).

Machine-translated text into languages other than English is even more problematic and routinely riddled with mistakes. For instance, [COVID-19 testing information auto-translated into German](#) included invented words, grammatical errors, and inconsistencies.

What machine translation can and can't do

Machine translation is not as good as most people think, but it is useful to get the gist of web sites or be able to ask for directions in a tourist destination with the help of an app.

However, that is not where it ends. Translation apps are [increasingly used in high-stakes contexts, such as hospitals](#), where staff may attempt to bypass human interpreters for quick communication with patients who have limited proficiency in English.

This causes big problems when, for instance, a patient's discharge instructions state [the equivalent of "Your United States was normal"](#)—an error resulting from the abbreviation "US" being used for "ultrasound" in medical contexts.

Therefore, there is consensus that translation apps are suitable [only in risk-free or low-risk situations](#). Unfortunately, sometimes even a caption on a selfie can turn into a high-risk situation.

We need to cultivate human multilingual talent

Only humans can identify what constitutes a low- or high-risk situation and whether the use of machine [translation](#) may be appropriate. To make informed decisions, humans need to understand both how languages work and how machine learning works.

It could be argued that all the errors described here can be ironed out with more training data. There are two problems with this line of reasoning. First, AI already has more [training data](#) than any human will ever be able to ingest, yet makes mistakes no human with much lower levels of investment in their language learning would make.

Second, and more perniciously, training machines to do our language learning for us is incredibly costly. There are [the well-known environmental costs of AI](#), of course. But there is also the cost of dismantling language teaching programs.

If we let go of language programs because we can outsource simple

multilingual tasks to machines, we will never train humans to achieve advanced language proficiency. Even from the perspective of pure strategic national interest, the skills to communicate across language barriers in more risky contexts of economics, diplomacy or health care are essential.

Languages are diverse, fuzzy, variable, relational and deeply social. Algorithms are the opposite. By buying into the hype that [machines](#) can do our language work for us [we dehumanize what it means to use languages to communicate](#), to make meaning, to create relationships and to build communities.

This article is republished from [The Conversation](#) under a Creative Commons license. Read the [original article](#).

Provided by The Conversation

Citation: 'Your United States was normal': Has translation tech really made language learning redundant? (2023, November 22) retrieved 27 April 2024 from <https://phys.org/news/2023-11-states-tech-language-redundant.html>

This document is subject to copyright. Apart from any fair dealing for the purpose of private study or research, no part may be reproduced without the written permission. The content is provided for information purposes only.