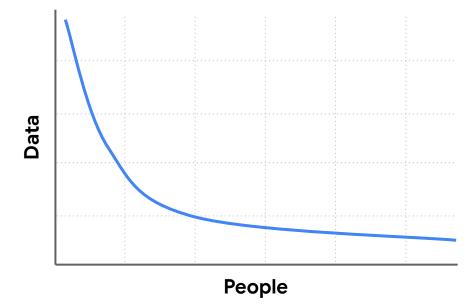
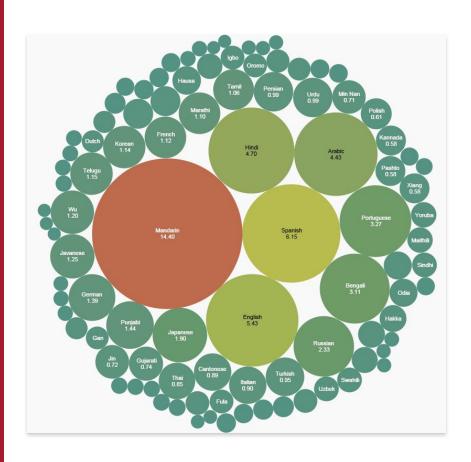
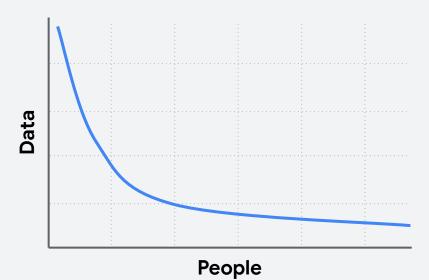
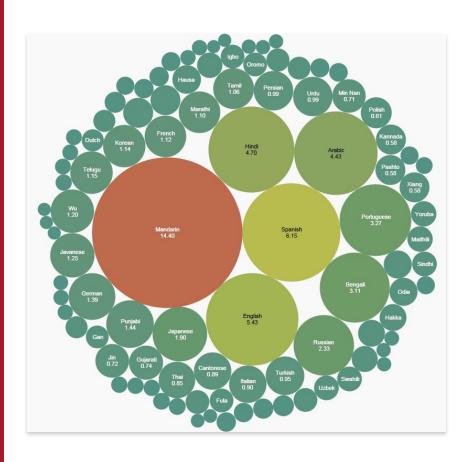
Crowdsourcing Data for the Long Tail









- Speech commands for the whole planet?
- For more than just voice assistants

Cost Model v. Community Model?



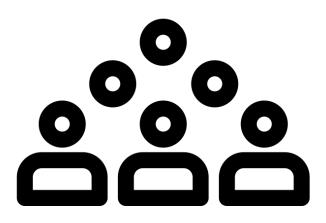
Limited Scale

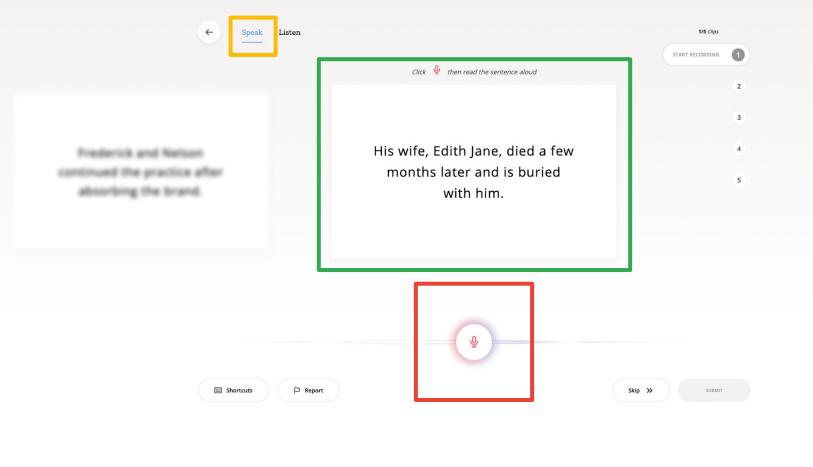


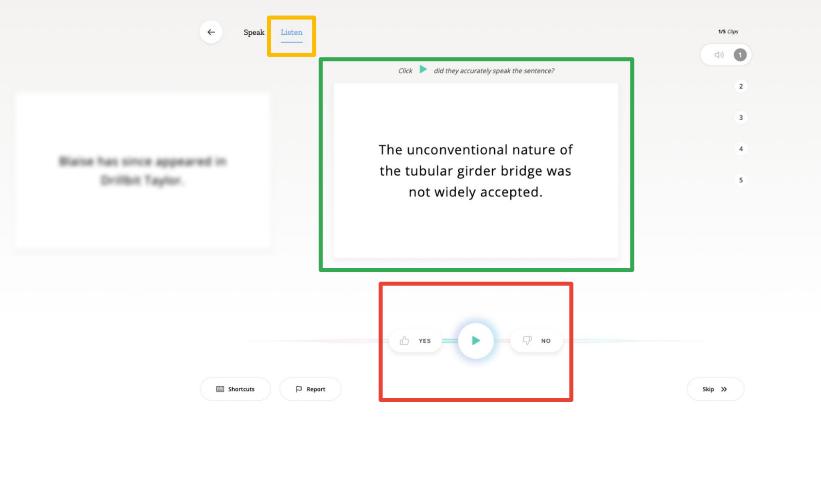
Social Good

https://commonvoice.mozilla.org

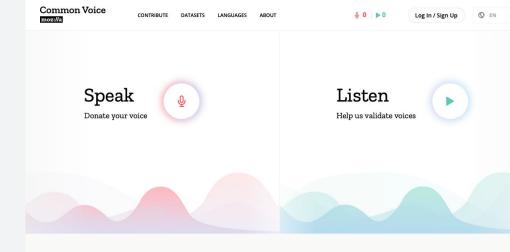
• Crowdsourcing platform







- Crowdsourcing platform
- Over 50,000 volunteers



Common Voice is Mozilla's initiative to help teach machines how real people speak.

Voice is natural, voice is human. That's why we're excited about creating usable voice technology for our machines. But to create voice systems, developers need an extremely lage amount of voice data.

Most of the data used by large companies isn't available to the majority of people. We think that stifles innovation. So we've launched Common Voice, a project to help make voice recognition open and accessible to everyone.

READ MORE





- Crowdsourcing platform
- Over 50,000 volunteers
- 54 different languages

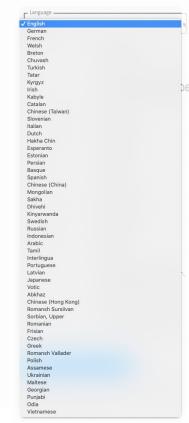
We're building

an open source, multi-language dataset of voices that anyone can use to train speech-enabled applications.

We believe that large, publicly available voice datasets will foster innovation and healthy commercial competition in machine-learning based speech technology.

Common Voice's multi-language dataset is already the largest publicly available voice dataset of its kind, but it's not the only one

Look to this page as a reference hub for other open source voice datasets and, as Common Voice continues to grow, a home for our release updates.



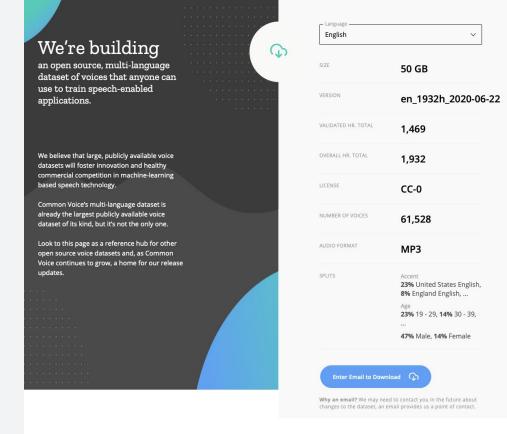


Recorded Hours



What's inside the Common Voice dataset?

- Crowdsourcing platform
- Over 50,000 volunteers
- 54 different languages
- Goal: speech recognition for all languages on the planet





File **Structure**

Valid

 At least 2 people listen to them, and the majority of those listeners say the audio matches the text

Invalid

 At least 2 listeners, and the majority say the audio does not match the clip

Other

 All other clips, i.e., fewer than 2 votes, or those that have equal valid and invalid votes, are labelled "other"



Interesting Attributes

- Permissive license
- Many contributors
- Comes with metadata

We're building

an open source, multi-language dataset of voices that anyone can use to train speech-enabled applications.

We believe that large, publicly available voice datasets will foster innovation and healthy commercial competition in machine-learning based speech technology.

Common Voice's multi-language dataset is already the largest publicly available voice dataset of its kind, but it's not the only one

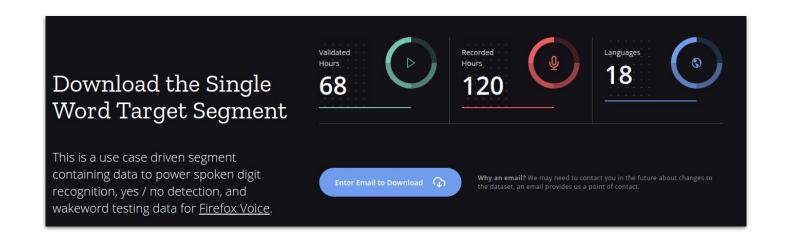
Look to this page as a reference hub for other open source voice datasets and, as Common Voice continues to grow, a home for our release undates.





What's inside the Common Voice





Single Word Target Segment

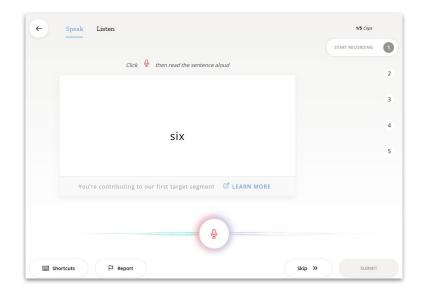
A speech commands-style dataset for 18 languages

- "Yes" // "no"
- "hey" & "Firefox"
- **digits** 0-9



ASR **Diversity** and **Reach**

- Common Voice
 - Permissive license
 - Minority languages
- Ease-of-use, wide reach
 - Browser-based
 - Community can add new languages
- You can contribute!



Common Voice Data Structure

