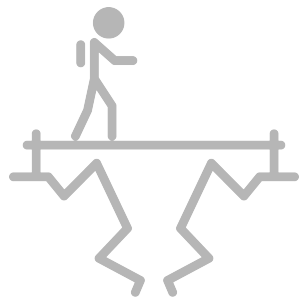


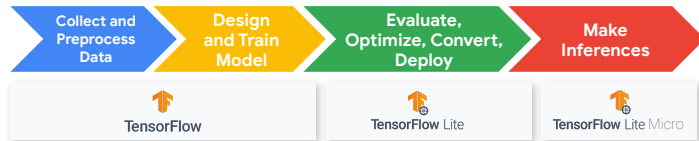
Keyword Spotting Datasets



What are we going to learn?



Challenges with
Keyword
Spotting



The Keyword
Spotting ML
Pipeline



Hands-on training
of a Keyword
Spotting Model

What does it mean to
have a **good** dataset?



So how do we **build** a **good** dataset?

- Who are the **users**?

So how do we **build** a **good** dataset?

- Who are the **users**?
- What do they **need**?

So how do we **build** a **good** dataset?

- Who are the **users**?
- What do they **need**?
- What **task** are they trying to solve?

So how do we **build** a **good** dataset?

- Who are the **users**?
- What do they **need**?
- What **task** are they trying to solve?
- How do they **interact** with the system?

So how do we **build** a **good** dataset?

- Who are the **users**?
- What do they **need**?
- What **task** are they trying to solve?
- How do they **interact** with the system?
- How does the **real world** make this hard?

Speech Commands: A Dataset for Limited-Vocabulary Speech Recognition

Pete Warden
Google Brain
Mountain View, California
`petewarden@google.com`

April 2018

The Speech Commands Dataset

- Recorded as individual **words** not sentences

The Speech Commands Dataset

- Recorded as individual **words** not sentences
- 1000-4000 examples of each word

The Speech Commands Dataset

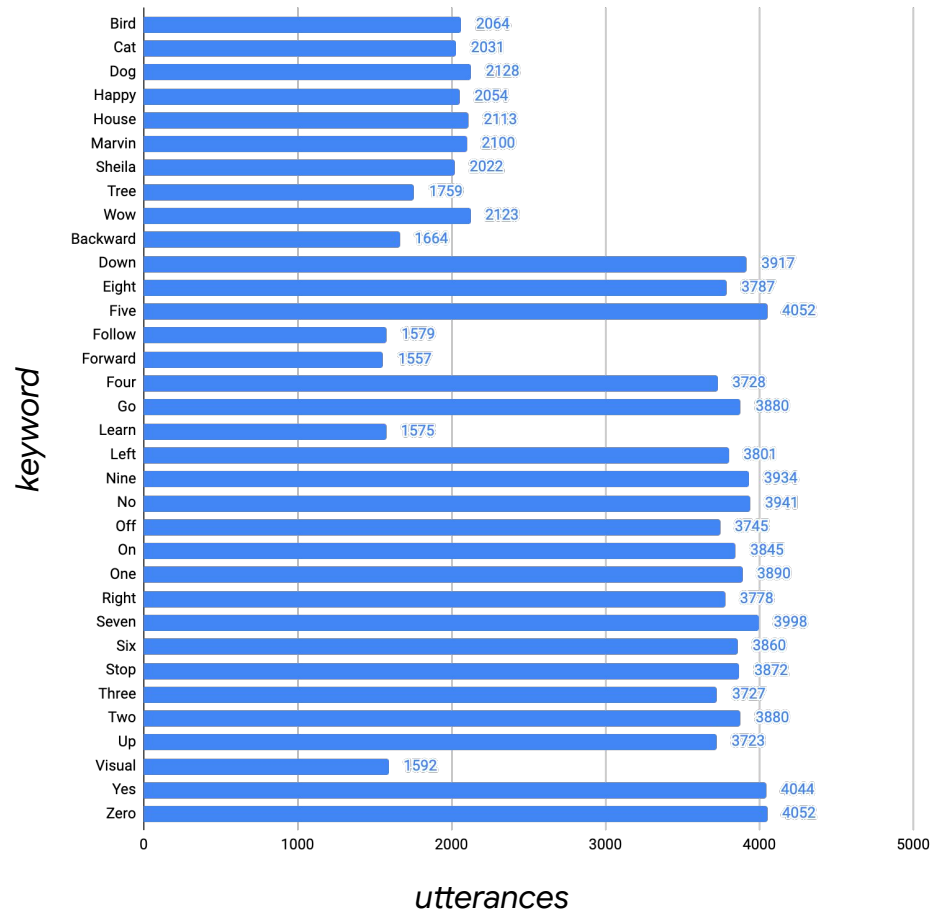
- Recorded as individual **words** not sentences
- 1000-4000 examples of each word
- >2,500 volunteers

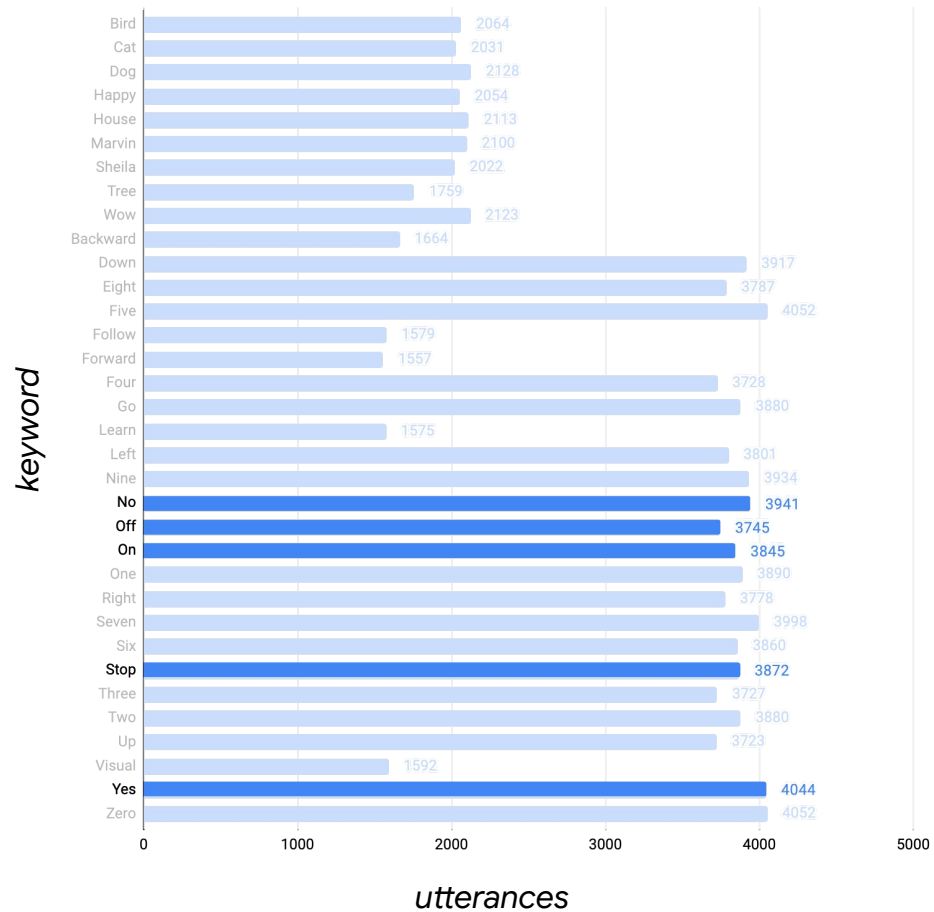
The Speech Commands Dataset

- Recorded as individual **words** not sentences
- 1000-4000 examples of each word
- >2,500 volunteers
- Representative of **real world audio** and includes background noise as well

The Speech Commands Dataset

- Recorded as individual **words** not sentences
- 1000-4000 examples of each word
- >2,500 volunteers
- Representative of **real world audio** and includes background noise as well
- **25 “IoT keywords”** + **10 “unknown words”** (with phonetic similarities: “three” vs “tree”)





Food for Thought

QC (Quality Control)

- Need to keep **only** what a human can hear
- Microphone issues
- **Noisy** backgrounds
(more on this soon)