# Using Existing Datasets for TinyML

# **Don't *collect*** from scratch

Data collection is **difficult**!
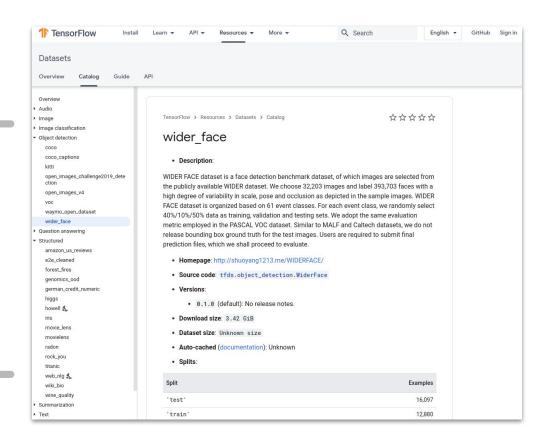
- Can we *reuse* existing data?

**What's available?** **What's missing?**

# TensorFlow
## Datasets Catalog

*Audio*
*Image*
*Image Classification*
*Object Detection*
*Question Answering*
*Structured*
*Summarization*
*Text*
*Translate*
*Video*

# TinyML
# Person Detection

- **Visual Wake Words**: a new dataset built from Common Objects in Context (**COCO**)
  - *people* v. *no people*

**Repurposing** existing datasets for **TinyML** tasks is a powerful concept

# **Don't** *learn* from scratch

- **Transfer** learning
- **Pretrained** models: your "AI Data Labeling Assistant"
- **Generate** your own data
  - Simulations
  - ML models