# Takeaways from Course 2

# Tiny Machine Learning Apps

# Tiny Machine Learning Apps

Collect Data → Preprocess Data → Design a Model → Train a Model → Evaluate Optimize → Convert Model → Deploy Model → Make Inferences

# Tiny Machine Learning Apps

| Collect Data | Preprocess Data | Design a Model | Train a Model | Evaluate Optimize | Convert Model | Deploy Model | Make Inferences |
|---|---|---|---|---|---|---|---|

**TensorFlow**

**TensorFlow** Lite

**TensorFlow** Lite Micro

**ML** Code

# Life cycle of **ML**

DATA FIXES

DATA NEEDS

**Data Collection**
*Continuous input stream*

Raw data

**Data Ingestion**
*Prep data for downstream ML apps*

Indexed data

**Data Analysis, Curation**
*Inspect/select the right data*

Selected data

**Data Labelling**
*Annotate data*

Labeled data

**Data Validation**
*Verify data is usable through pipeline*

Validated data

**Data Preparation**
*Prep data for ML uses (split, versioning)*

Online Performance

ML ready Datasets

**ML System Deployment**
*Deploy ML system to production*

ML Certificate

**ML System Validation**
*Validate ML system for deployment*

KPIs

**Model Evaluation**
*Compute model KPIs*

Models

**Model Training**
*Use ML algos to create models*

Online ML System

Validated ML System

**Acoustic Sensors**
Ultrasonic, **Microphones**, Geophones, Vibrometers

**Image Sensors**
Thermal, **Image**

**Motion Sensors**
Gyroscope, Radar, **Accelerometer**

Minimum

**Course 2:** End-to-end **TinyML** application design

**AI Infrastructure**

**Data Engineering**

Model Engineering

Model Deployment

Product Analytics

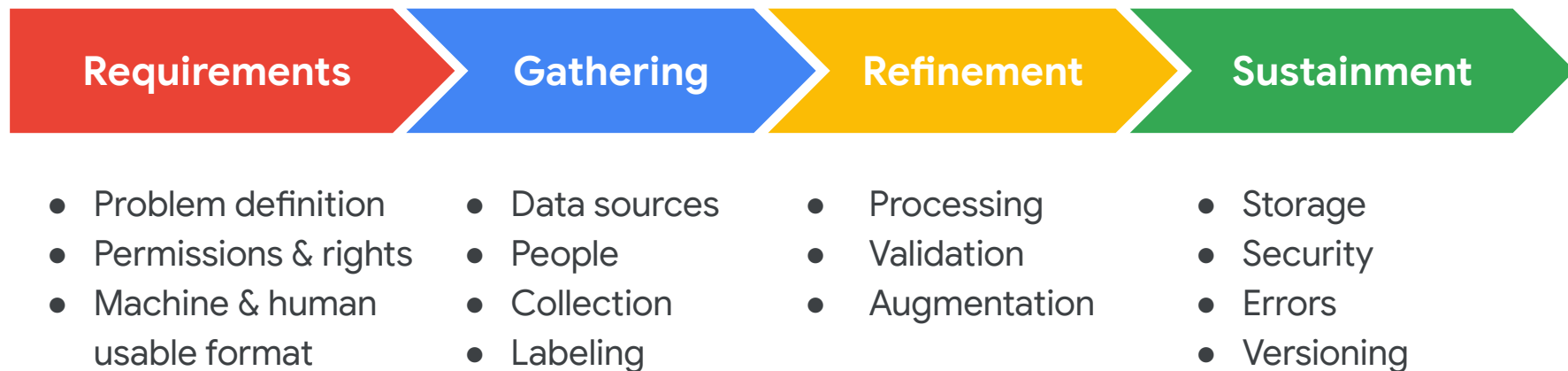**Collect Data** → Preprocess Data → Design a Model → Train a Model → Evaluate Optimize → Convert Model → Deploy Model → Make Inferences

# Speech Commands: A Dataset for Limited-Vocabulary Speech Recognition

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

April 2018

# Data Engineering

| Requirements | Gathering | Refinement | Sustainment |
|---|---|---|---|

- Problem definition
- Permissions & rights
- Machine & human usable format

- Data sources
- People
- Collection
- Labeling

- Processing
- Validation
- Augmentation

- Storage
- Security
- Errors
- Versioning

# Data **Collection** and **Processing**

## Someone scraped 40,000 Tinder selfies to make a facial dataset for AI experiments

**Natasha Lomas**   @riptari   /   7:21 PM EDT • April 28, 2017

**Update:** A Tinder spokesperson has now provided the following statement:

We take the security and privacy of our users seriously and have tools and systems in place to uphold the integrity of our platform. It's important to note that Tinder is free and used in more than 190 countries, and the images that we serve are profile images, which are available to anyone swiping on the app. We are always working to improve the Tinder experience and continue to implement measures against the automated use of our API, which includes steps to deter and prevent scraping.

This person has violated our terms of service (Sec. 11) and we are taking appropriate action and investigating further.

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

Data collection is **difficult**!

- Can we *reuse* existing data?

**What's available?**

**What's missing?**

# Visual Wake Words **Dataset**

## Visual Wake Words Dataset

Aakanksha Chowdhery, Pete Warden, Jonathon Shlens,
Andrew Howard, Rocky Rhodes
Google Research
{chowdhery, petewarden, shlens, howarda, rocky}@google.com

# Visual Wake Words **Dataset**



*Label:* "person"

*Label:* "person"

# Common Voice

- **Crowdsourcing** platform
- Over **50,000 volunteers**

# Bias and **Market Forces**

**AI Infrastructure**

**Data Engineering**

Model Engineering

Model Deployment

Product Analytics

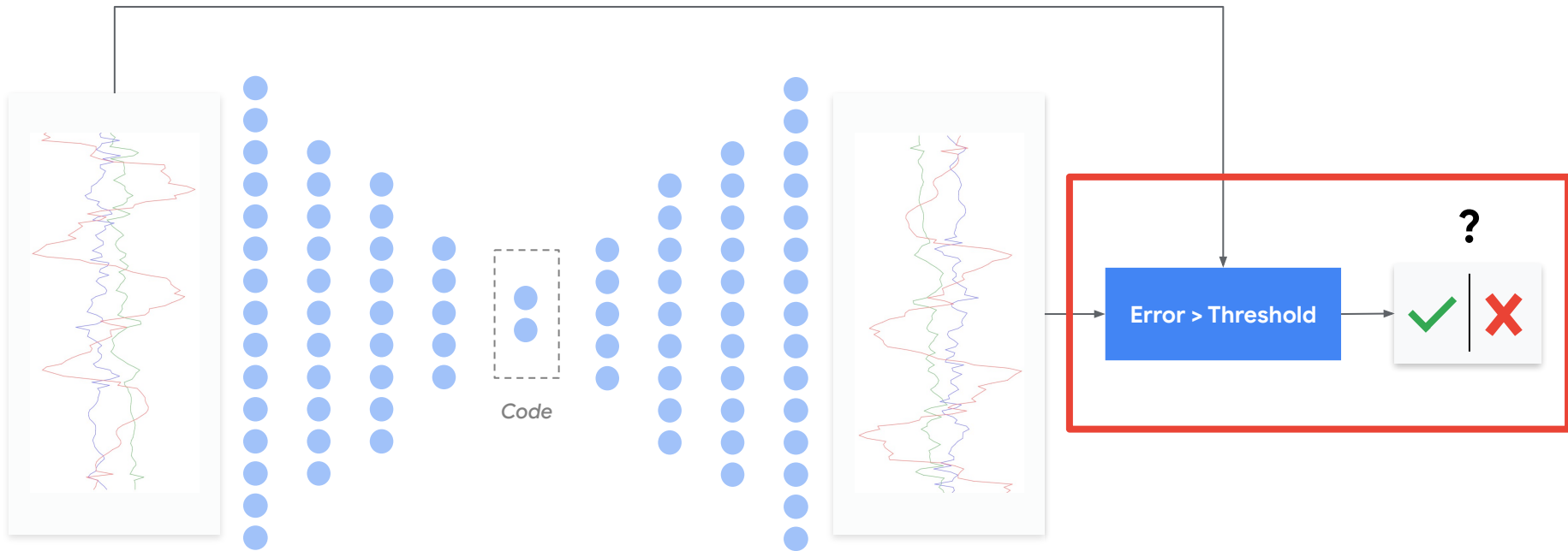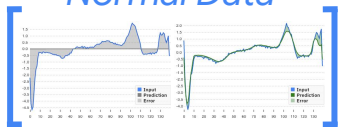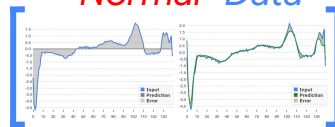| Collect Data | **Preprocess Data** | Design a Model | Train a Model | Evaluate Optimize | Convert Model | Deploy Model | Make Inferences |

# Data Preprocessing: **Spectrograms**

# Depthwise Convolution *(3 Channel—e.g., RGB)*

includes **pointwise** conv

Reuse (freeze general feature extraction)

Train **only** last few layers

$W_{A1}$    $W_{A2}$    $W_{A3}$    $W_{A4}$    $W_{A5}$    $W_{A6}$    $W_{A7}$

Input A

Labels A

Learns *general features* irrespective of task

*Task-specific* features

Normal Data

"Normal" Data

Code

Error > Threshold

?

AI Infrastructure

Data Engineering

Model Engineering

**Model Deployment**

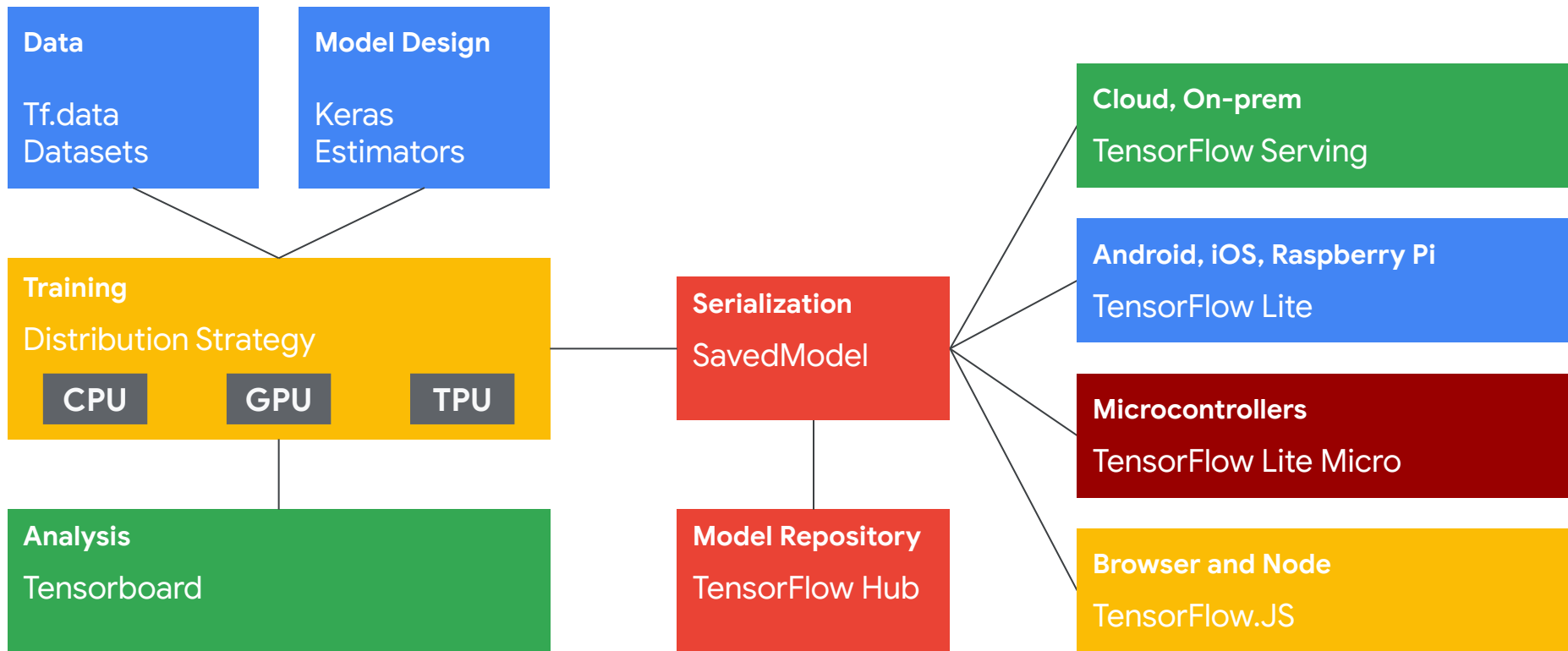Product Analytics

Collect Data → Preprocess Data → Design a Model → Train a Model → Evaluate Optimize → **Convert Model** → Deploy Model → Make Inferences

# **Common** Metrics



**Accuracy**

*Quantitative*

**Efficiency**

*Quantitative*

**User** Experience

*Qualitative*

# Latency

# Latency

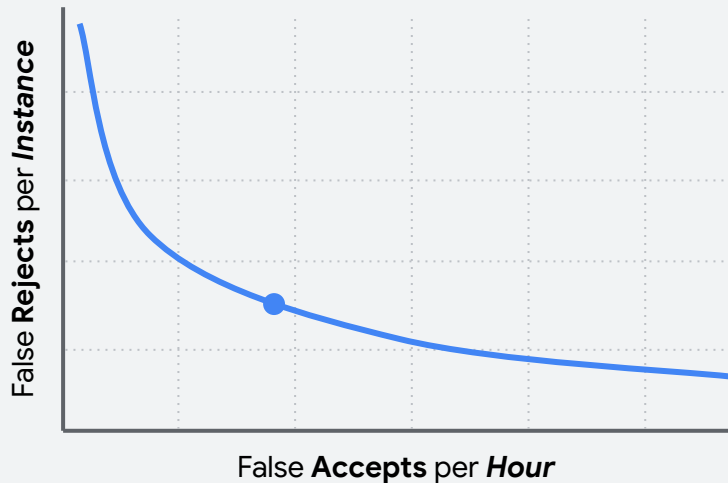Lower quality, but *faster* model?

**Metrics**

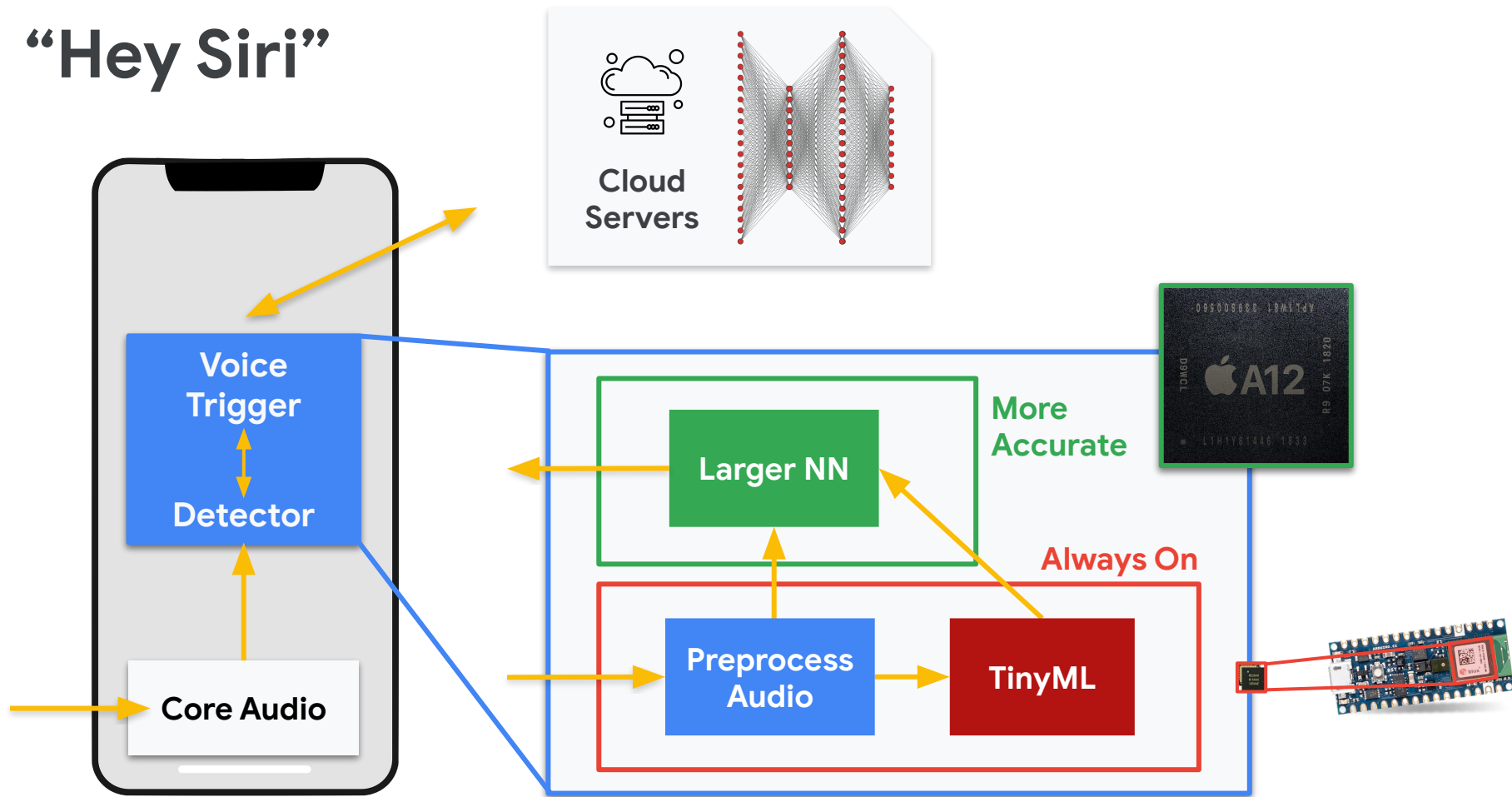| Accuracy | ~ |
| Efficiency (Latency) | ✓ |
| User Experience | ✓ |

# False Positive and False Negative

- Accuracy is measured as a tradeoff between **false accept rate** (FAR) and **false reject rate** (FRR)

False **Rejects** per *Instance*

False **Accepts** per *Hour*

"Hey Siri"

# Course Sequence

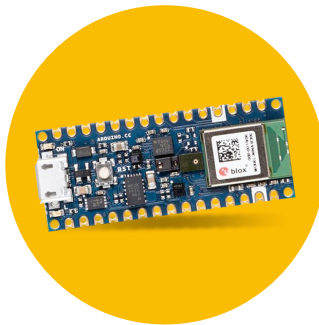**Course 1**

*Fundamentals of TinyML*

**Course 2**

*Applications of TinyML*

**Course 3**

*Deploying TinyML*



**✓ Learning** You will learn how to deploy models on a real microcontroller. Along the way you will explore the challenges unique to and amplified by **TinyML** (e.g., preprocessing, post-processing, dealing with resource constraints).