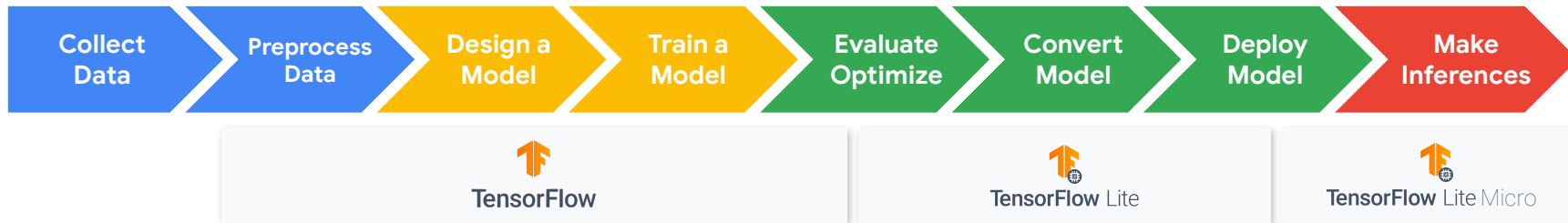
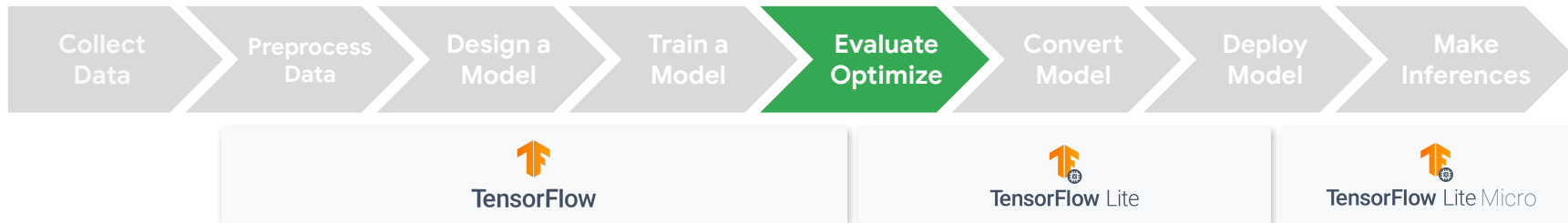


# Metrics for KWS







# What **metrics** matter?



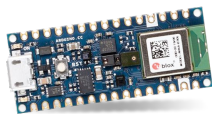
Accuracy



Efficiency



Beyond Model Metrics



# What **metrics** matter?



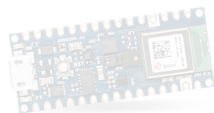
Accuracy



Efficiency

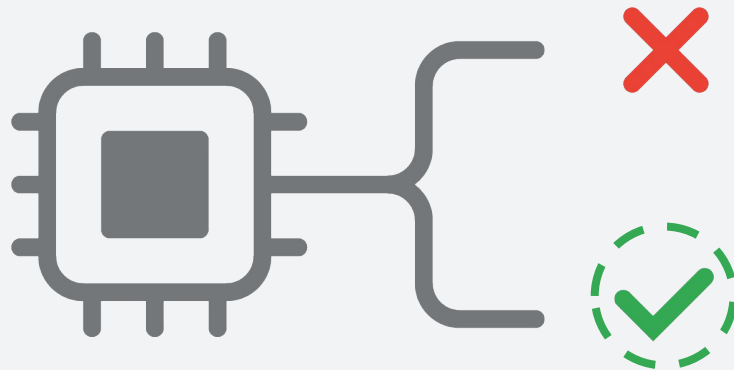


Beyond Model Metrics



# False Positive

Did **NOT** say keyword  
but device **DOES**  
trigger



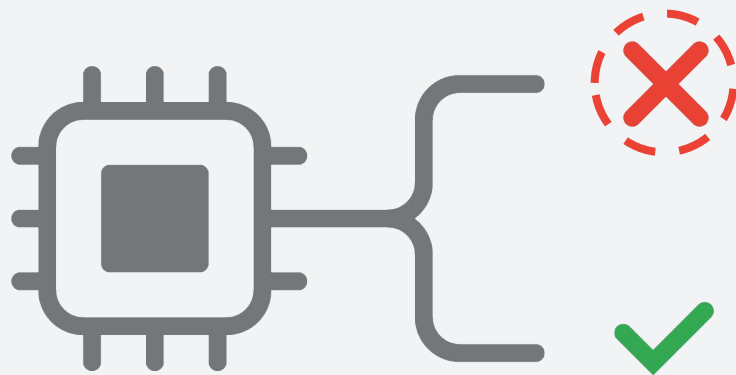
REAL



**ERROR**

# False Negative

**DID** say keyword  
but device **DOES**  
**NOT** trigger



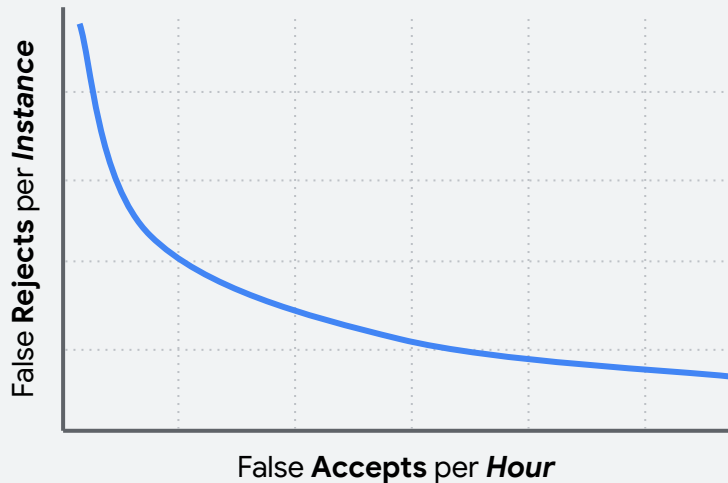
REAL



**ERROR**

# False Positive and False Negative

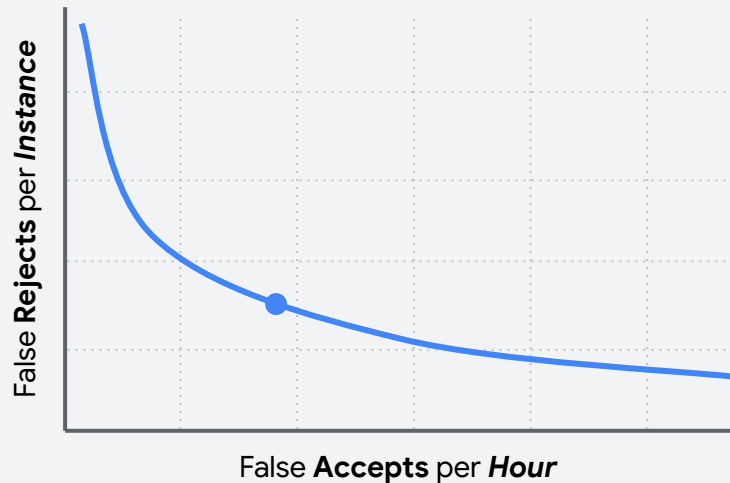
- Accuracy is measured as a tradeoff between **false accept rate** (FAR) and **false reject rate** (FRR)





# False Positive and False Negative

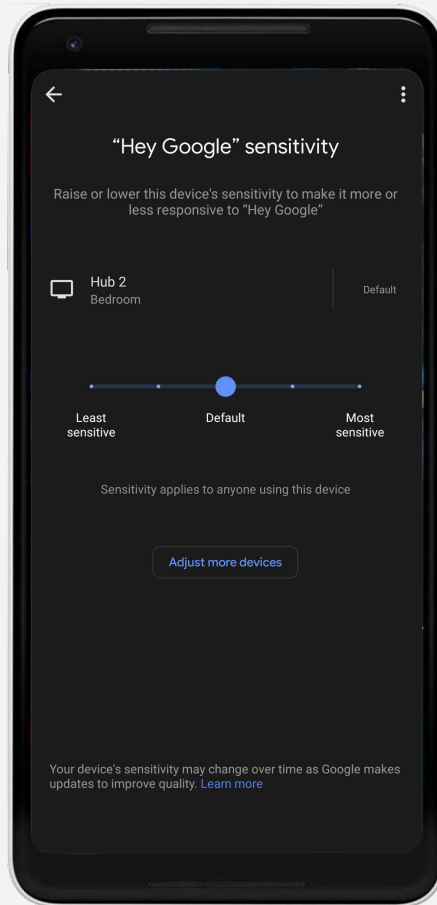
- Accuracy is measured as a tradeoff between **false accept rate** (FAR) and **false reject rate** (FRR)
- Choose an **operating point**

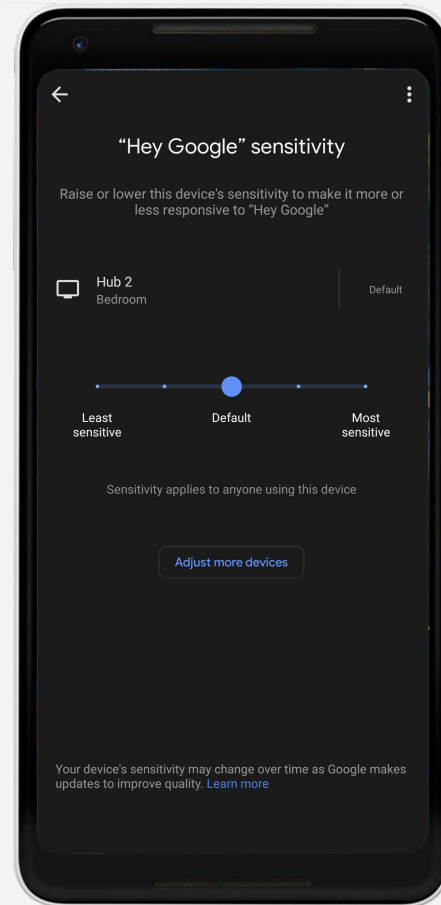
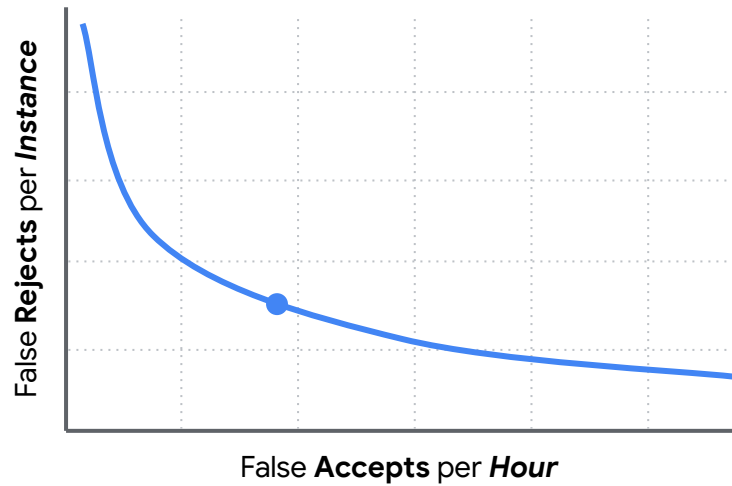


# Operating Point

***False Accept Rate*** and ***False Reject Rate*** are measured on **particular audio**—however

- Your phone might be in your pocket, purse, or backpack most of the day
- Your smart speaker might be next to a TV, or next to where your family eats, or in a relatively quiet bedroom
- You might not be a native speaker of English





# Latency

- Model must be **fast enough** to keep up with the speech input
- The model must run fast enough to be **responsive** to the end user



# What **metrics** matter?



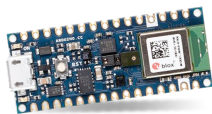
Accuracy



Efficiency



Beyond Model Metrics



# Latency

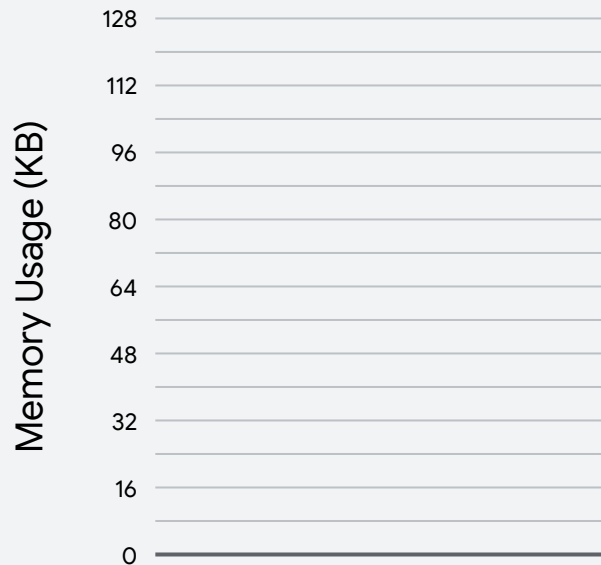
- Model must be **fast enough** to keep up with the speech input
- The model must run fast enough to be **responsive** to the end user
- But it must run efficiently on a *small* processor

**TinyML**



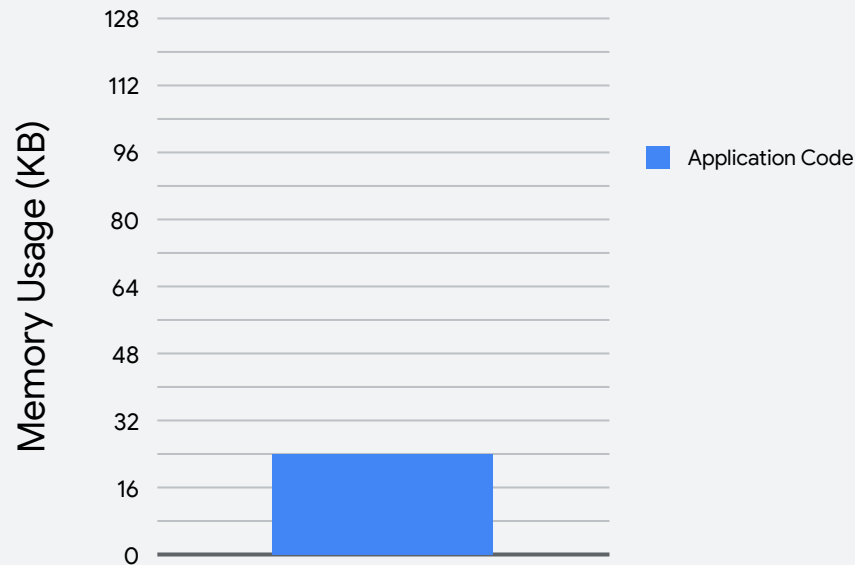
# Memory Usage

- Need to be **resource aware**
- **Less** compute
- **Less** memory
- Use **quantization**



# Memory Usage

- Need to be **resource aware**
- **Less** compute
- **Less** memory
- Use **quantization**





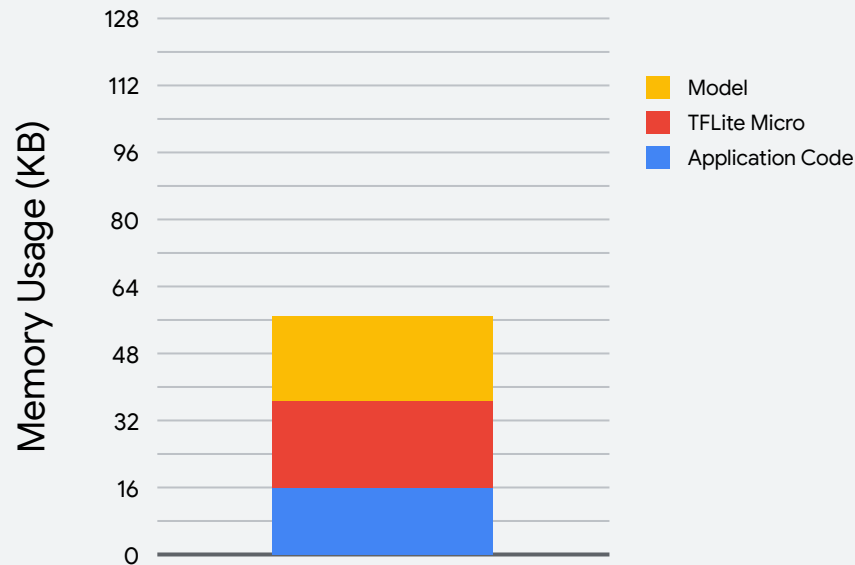
# Memory Usage

- Need to be **resource aware**
- **Less** compute
- **Less** memory
- Use **quantization**



# Memory Usage

- Need to be **resource aware**
- **Less** compute
- **Less** memory
- Use **quantization**



# Memory Usage

- Need to be **resource aware**
- **Less** compute
- **Less** memory
- Use **quantization**



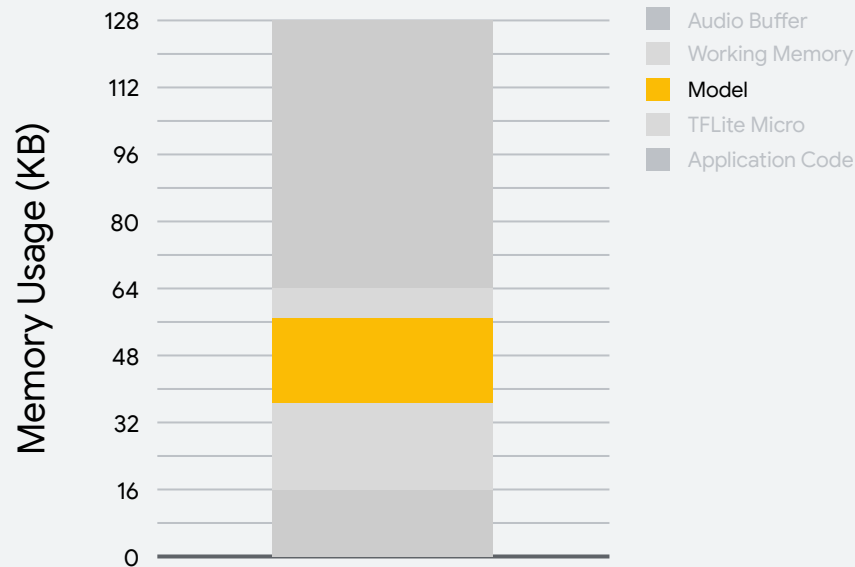
# Memory Usage

- Need to be **resource aware**
- **Less** compute
- **Less** memory
- Use **quantization**



# Memory Usage

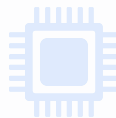
- Need to be **resource aware**
- **Less** compute
- **Less** memory
- Use **quantization**



# What **metrics** matter?



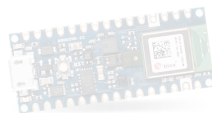
Accuracy



Efficiency



**Beyond Model Metrics**



# Beyond Model Metrics

Think about **quality of experience (QoE)**

What really defines that? It isn't just whether the model is performing well on a given dataset. It is about the **user experience**: how can we **assess** that?

# Beyond Model Metrics

Think about **quality of experience (QoE)**

What really defines that? It isn't just whether the model is performing well on a given dataset. It is about the **user experience**: how can we **assess** that?

- Have diverse **users** to test against?
- Test in different **backgrounds**?
- Add **noise** while training the models?



# So how can we **improve** our KWS Application?

