



# Hands on Embedded ML (Vision and Audio)



*Brian Plancher*  
*Harvard John A. Paulson School of Engineering and Applied Sciences*  
*brianplancher.com*



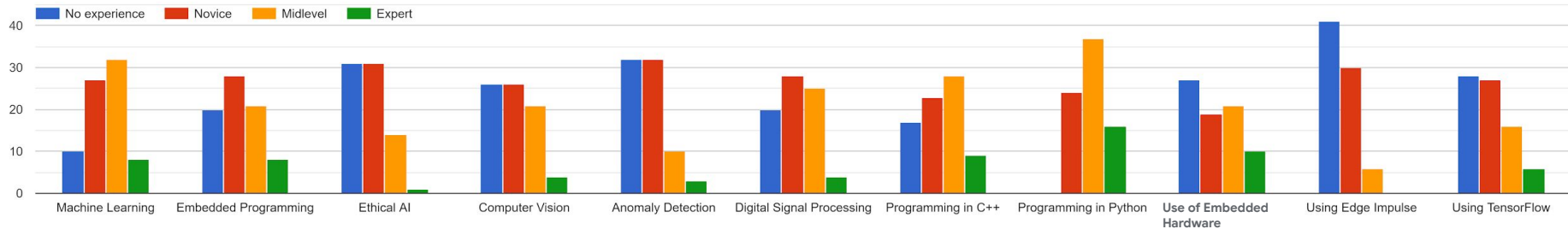
# Quick Disclaimer:

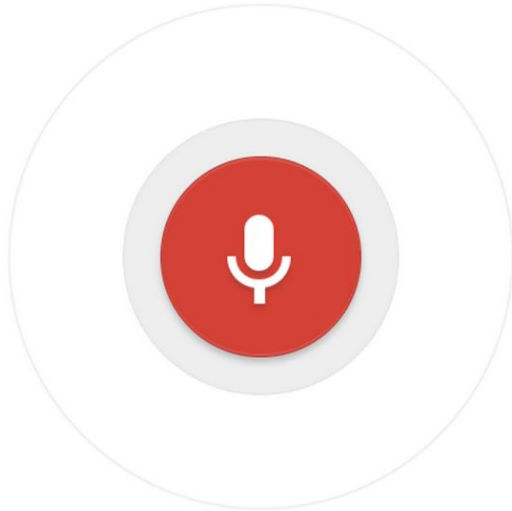
Today will be **both too fast**  
and **too slow!**

# Quick Disclaimer:

Today will be **both too fast**  
and **too slow!**

Do you have experience in?





# By the end of today: Hands-on Keyword Spotting

We will explore the **science** behind KWS and **collect data** and **train** our own custom model to recognize “yes” vs. “no” using **Edge Impulse**

# Today's Agenda

- Deep ML Background
- Hands-on Computer Vision: Thing Translator
- The Tiny Machine Learning Workflow
- Keyword Spotting (KWS) Data Collection
- KWS Preprocessing and Training
- Deployment Challenges and Opportunities for Embedded ML
- Summary

# Today's Agenda

- Deep ML Background

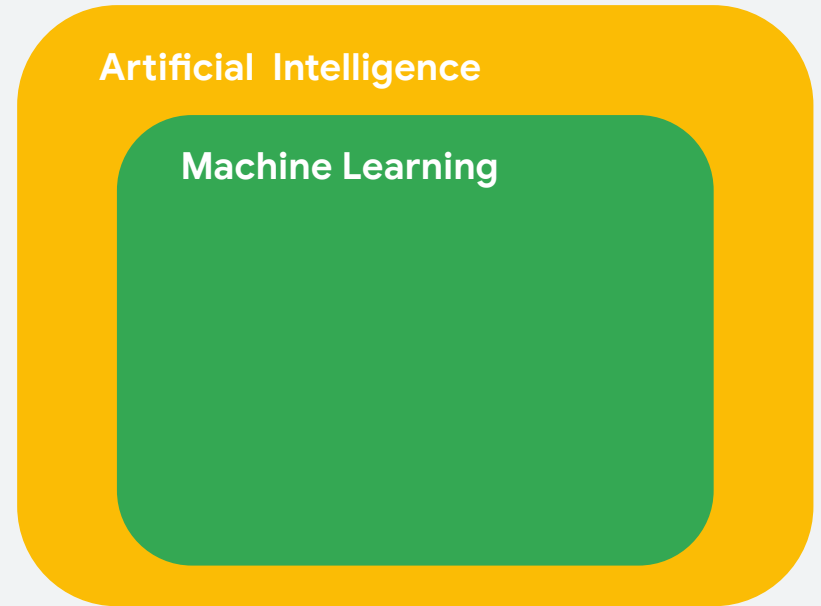
## How does (Deep) Machine Learning Work?

Exploring Deep ML through Computer Vision

- Hands-on Computer Vision: Thing Translator
- The Tiny Machine Learning Workflow
- Keyword Spotting (KWS) Data Collection
- KWS Preprocessing and Training
- Deployment Challenges and Opportunities for Embedded ML
- Summary

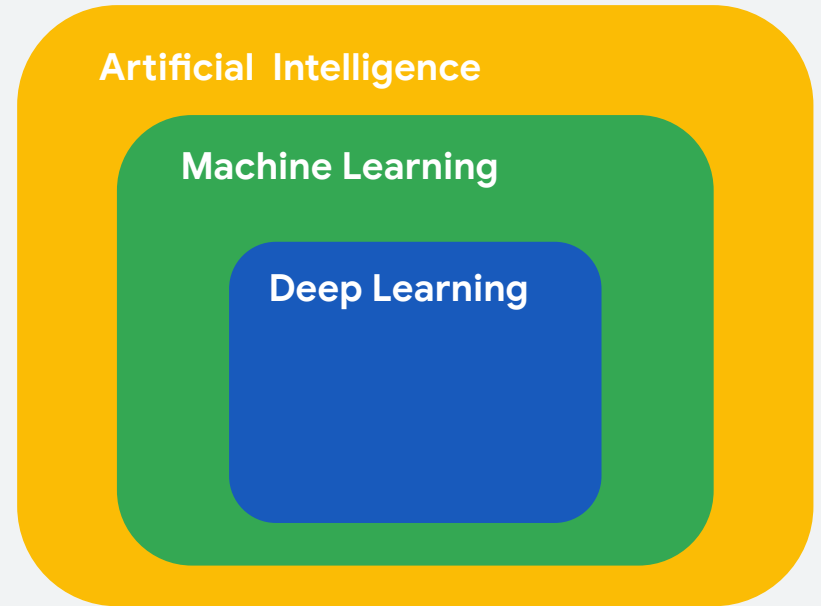
# What is Machine Learning?

1. **Machine Learning** is a subfield of **Artificial Intelligence** focused on developing algorithms that learn to **solve problems by analyzing data for patterns**



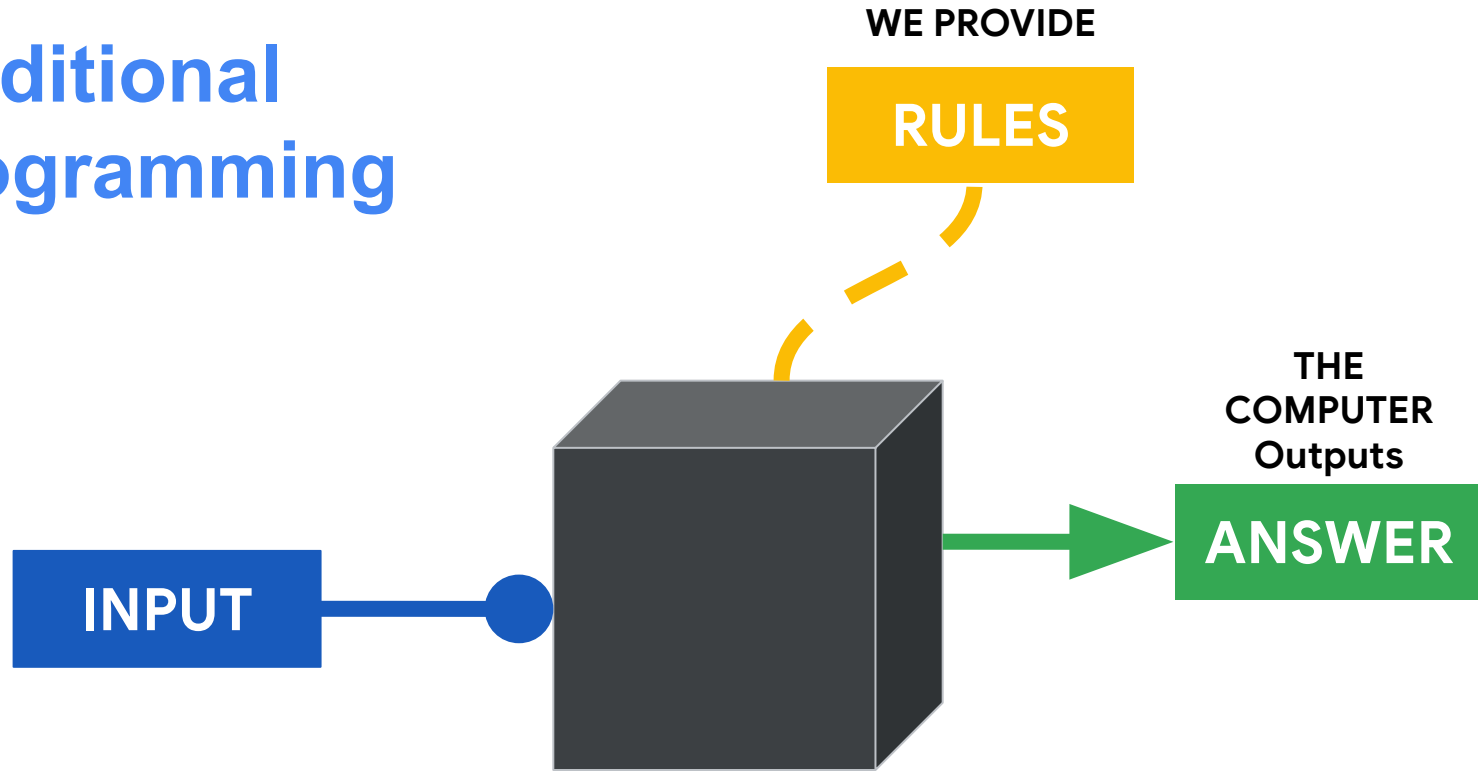
# What is (**Deep**) Machine Learning?

1. Machine Learning is a subfield of Artificial Intelligence focused on developing algorithms that learn to solve problems by analyzing data for patterns
2. **Deep Learning** is a type of **Machine Learning** that leverages **Neural Networks** and **Big Data**





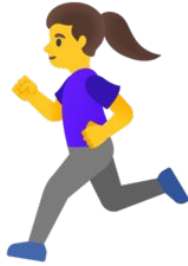
# Traditional Programming



# Let's try to figure out **what** she's doing?



```
if (speed < 4):  
    then walking
```



```
if (speed < 4):  
    then walking  
  
else:  
    running
```

**data** we can gather

input: speed

Write a rule

extend the rule

# Let's try to figure out **what** she's doing?



```
if (speed < 4):  
    then walking
```



```
if (speed < 4):  
    then walking
```

```
else:  
    running
```



```
if (speed < 4):  
    then walking
```

```
else if (speed < 12):  
    then running
```

```
else:  
    biking
```

# Let's try to figure out **what** she's doing?



```
if (speed < 4):  
    then walking
```



```
if (speed < 4):  
    then walking
```

```
else:  
    running
```



```
if (speed < 4):  
    then walking
```

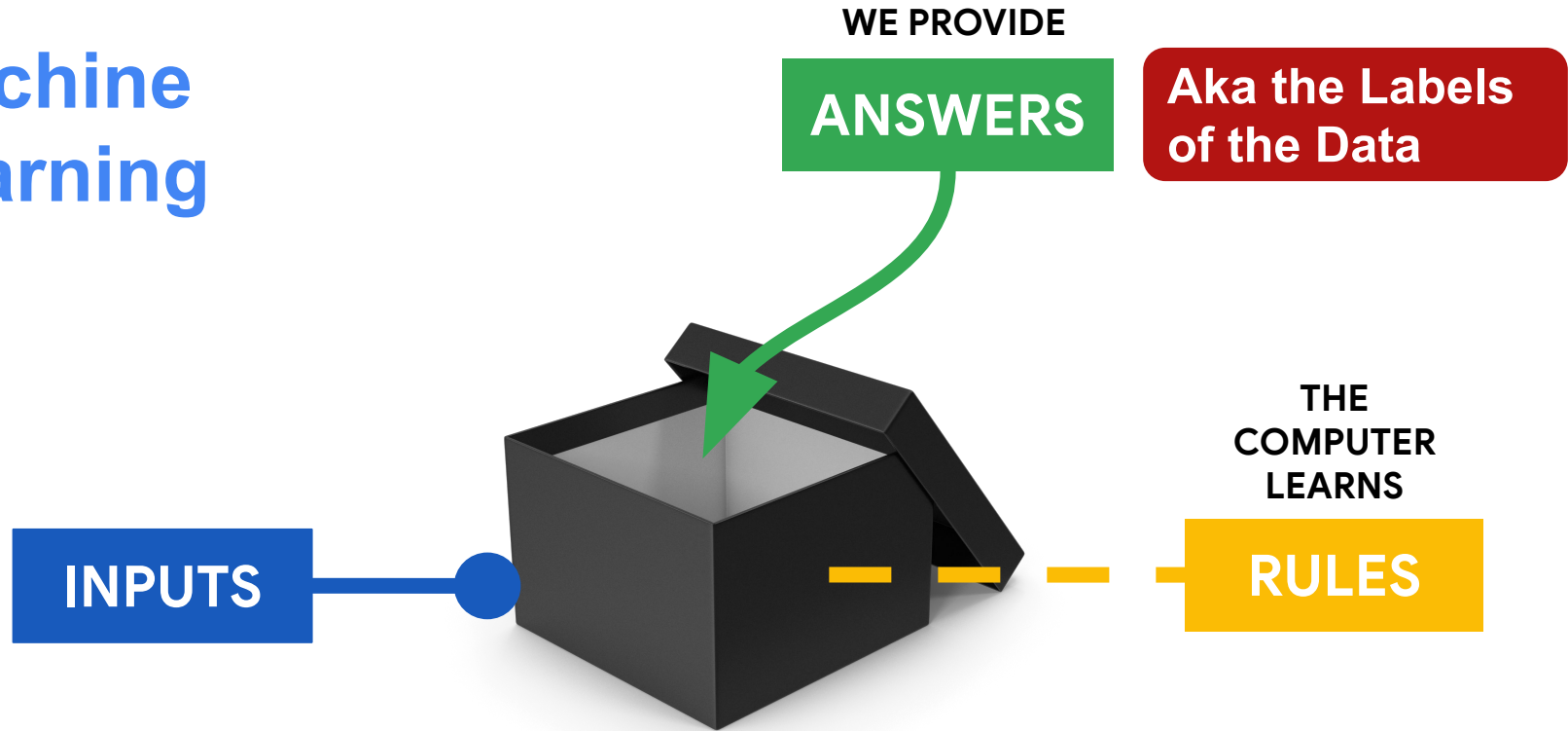
```
else if (speed < 12):  
    then running
```

```
else:  
    biking
```



?? **WHAT IS THIS** ??

# Machine Learning



Let's try to figure out **what** she's doing?



```
010101010010001110101
01010100101001001010
101010111010100101001
```



```
11110101001001010101
01010010100101010100
11010110010101001111
```



```
00001110101110101101
01010111101011010101
11010111111001001011
```



```
01111110101110101010
10101110101011010101
11111111100100001110
```

walking

running

biking

golfing

Let's try to figure out **what** she's doing?



```
010101010010001110101  
01010100101001001010  
10101011010100101001
```



```
11110101001001010101  
01010010100101010100  
11010110010101001111
```



```
00001110101110101101  
01010111101011010101  
11010111111001001011
```



```
01111110101110101010  
10101110101011010101  
11111111100100001110
```

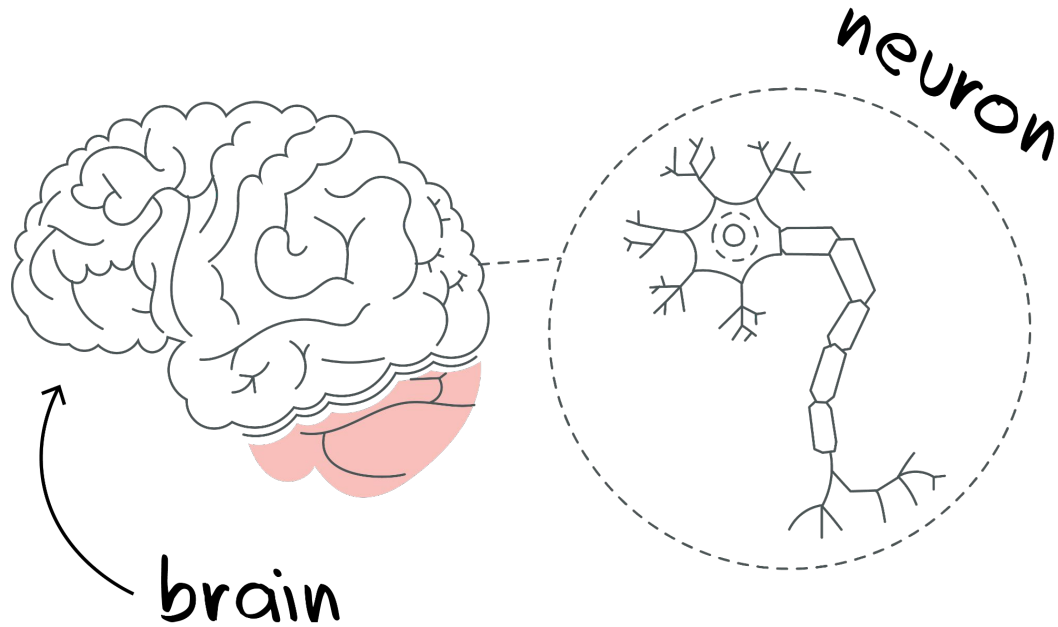
walking

running

biking

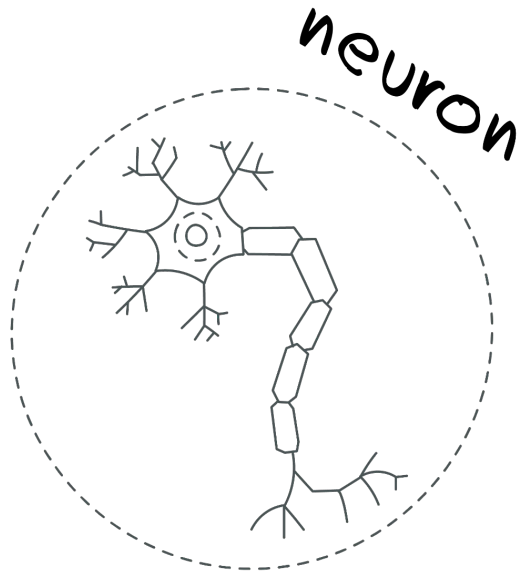
golfing

# What is a **neural network**?

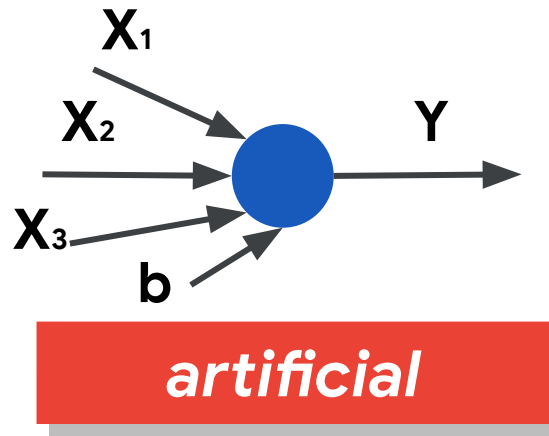
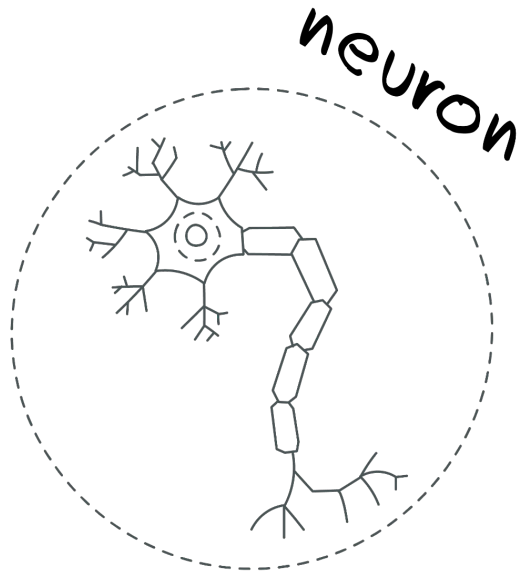




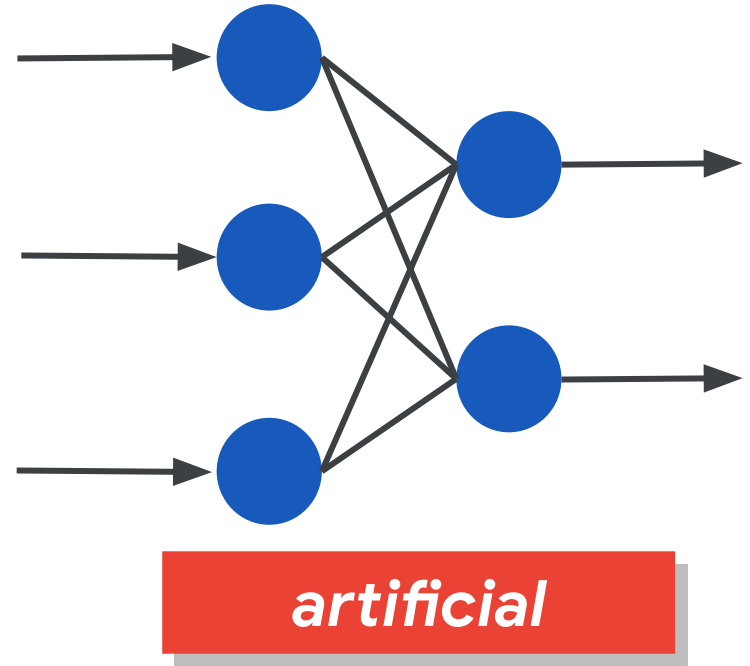
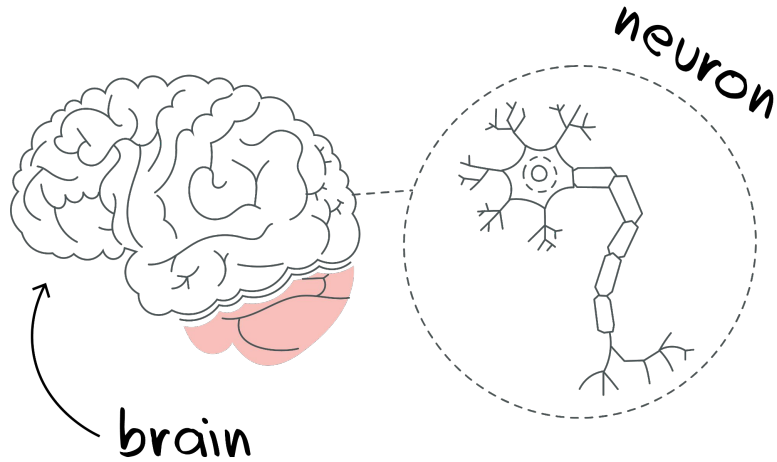
What is a **neural network**?



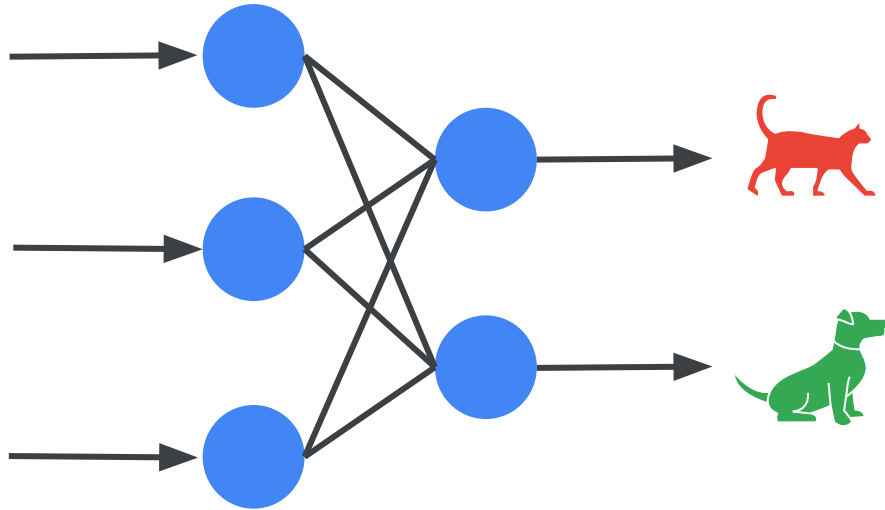
# What is a **neural network**?



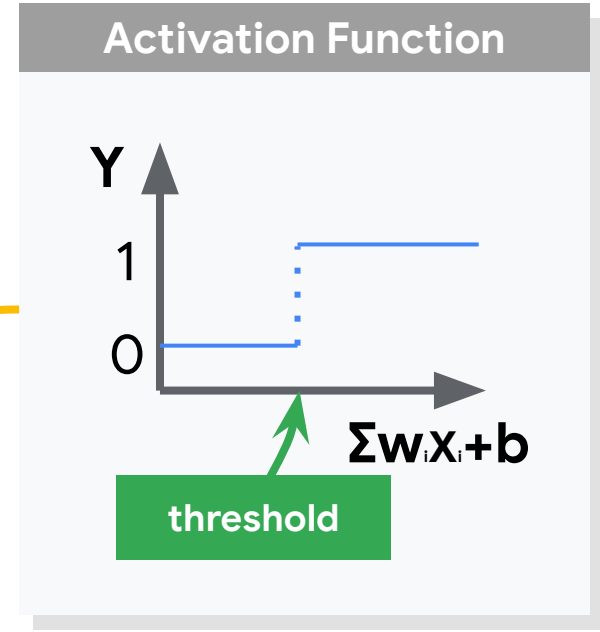
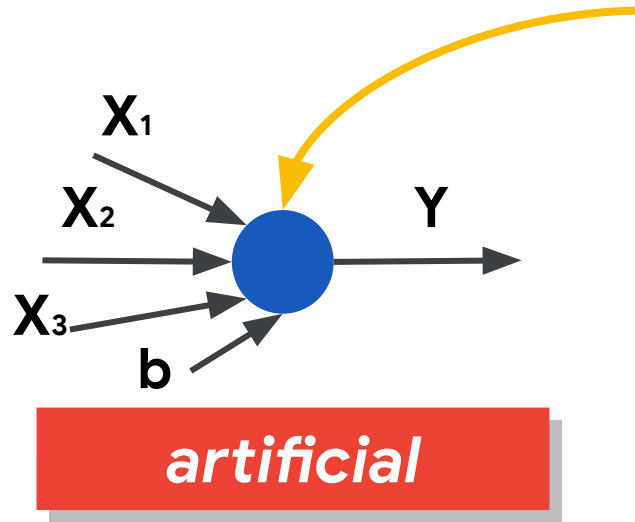
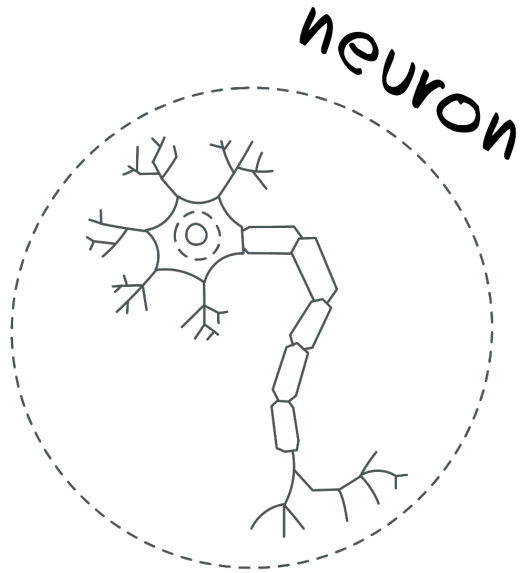
# What is a **neural network**?



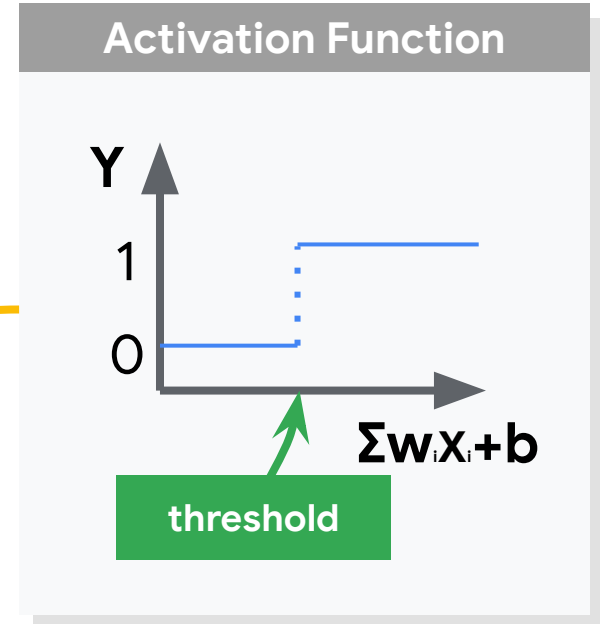
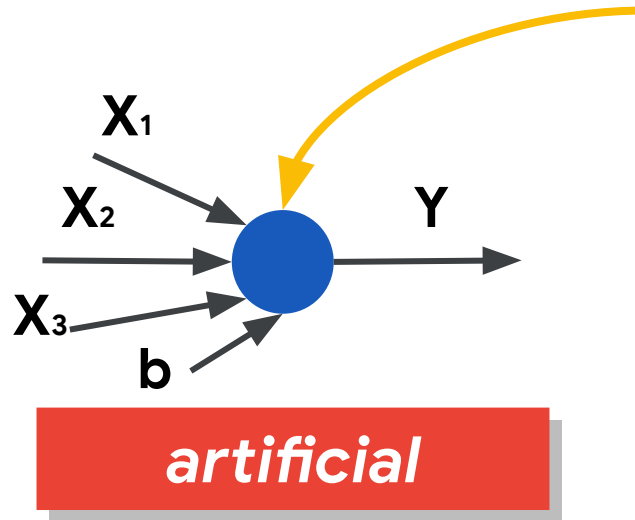
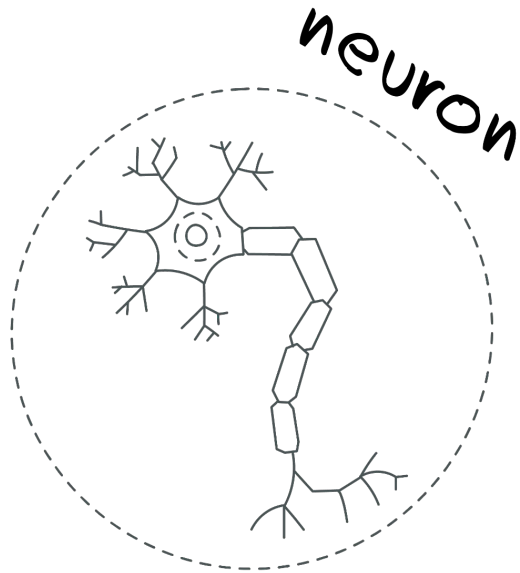
# What is a **neural network**?



# What is a **neural network**?



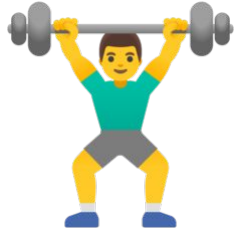
# What is a **neural network**?



**$Y = \sum w_i x_i + b$**

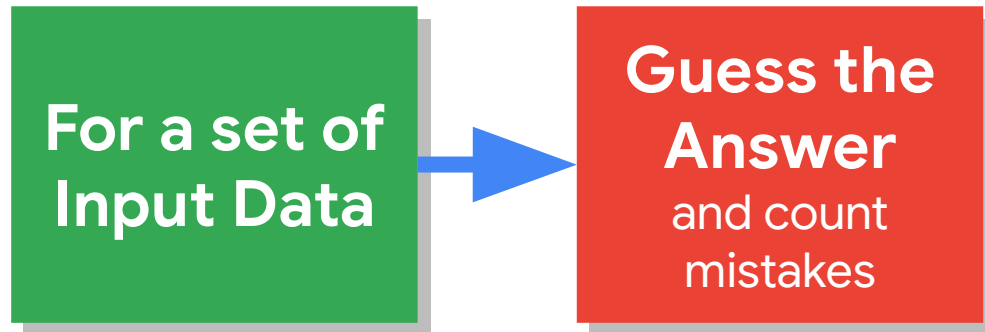
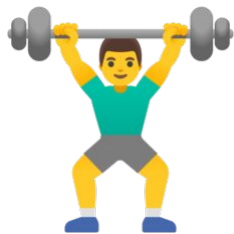
So training the model is finding the right values for  $w$  and  $b$

# Training the machine



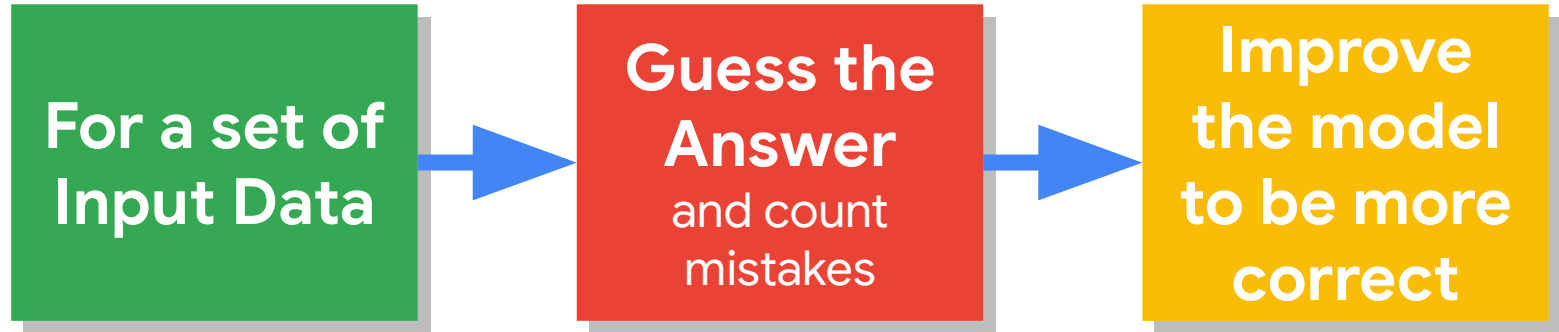
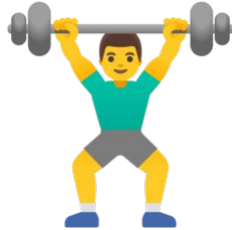
For a set of  
Input Data

# Training the machine

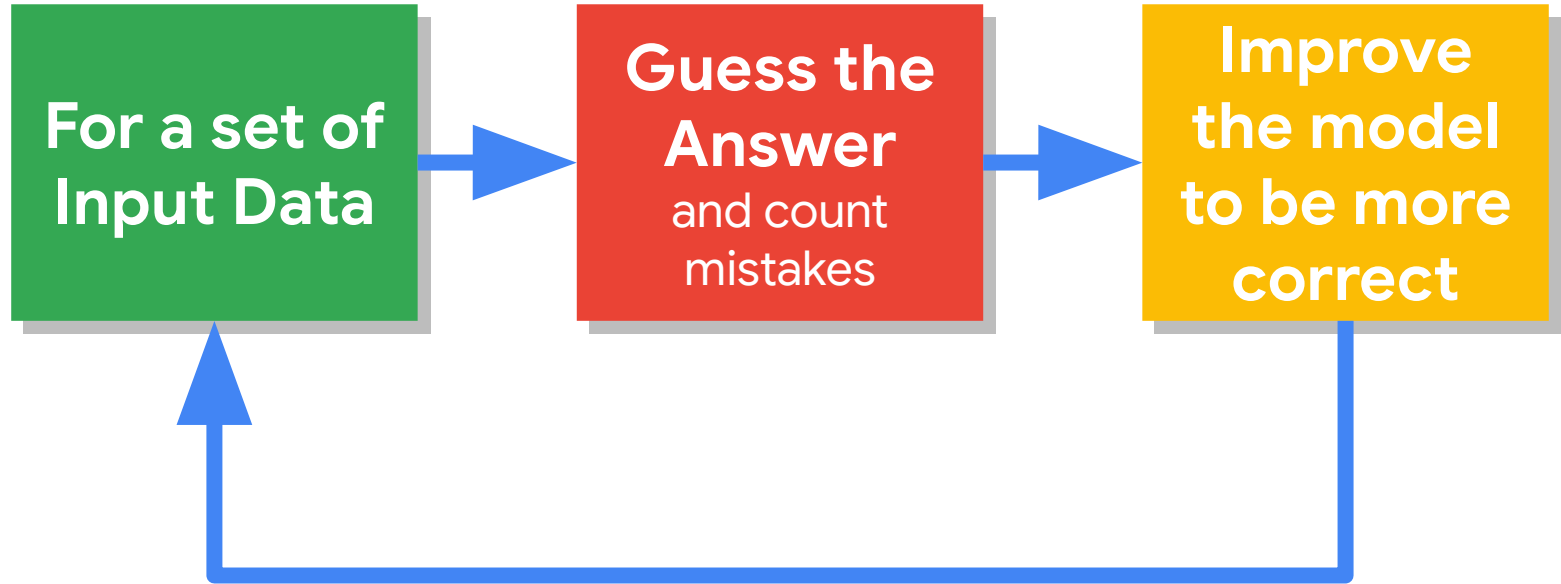
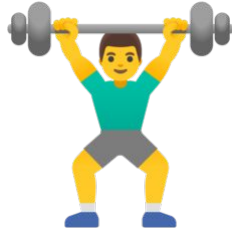




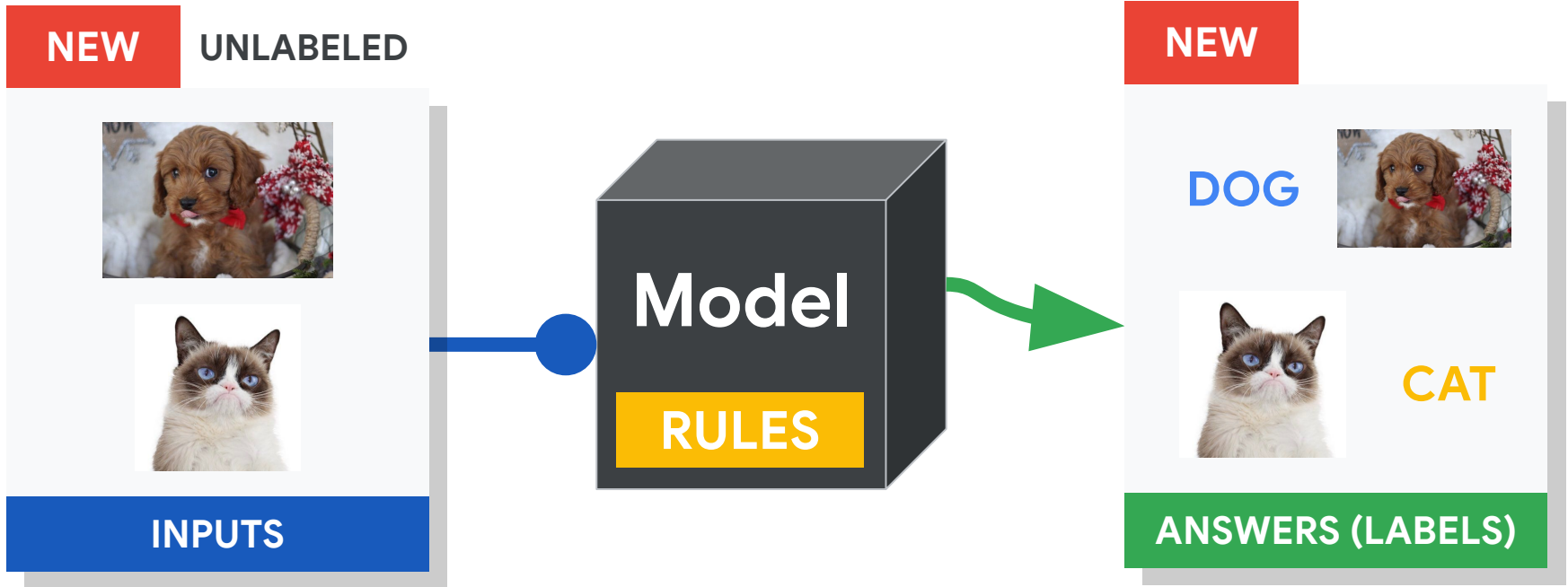
# Training the machine



# Training the machine



After it's **learned** use it for **inference**:



# What is a **neural network**?



To learn more about the **math behind neural network training** there is a nice series of videos here:  
[https://www.youtube.com/playlist?list=PLZHQObOWTQDNU6R1\\_67000Dx\\_ZCJB-3pi](https://www.youtube.com/playlist?list=PLZHQObOWTQDNU6R1_67000Dx_ZCJB-3pi)

*artificial*

# Today's Agenda

- Deep ML Background

How does (Deep) Machine Learning Work?

## Exploring Deep ML through Computer Vision

- Hands-on Computer Vision: Thing Translator
- The Tiny Machine Learning Workflow
- Keyword Spotting (KWS) Data Collection
- KWS Preprocessing and Training
- Deployment Challenges and Opportunities for Embedded ML
- Summary

# Computer Vision is Hard

**What color are the pants and the shirt?**



Slide Credit: Hamilton Chong

# Computer Vision is Hard



Slide Credit: Hamilton Chong

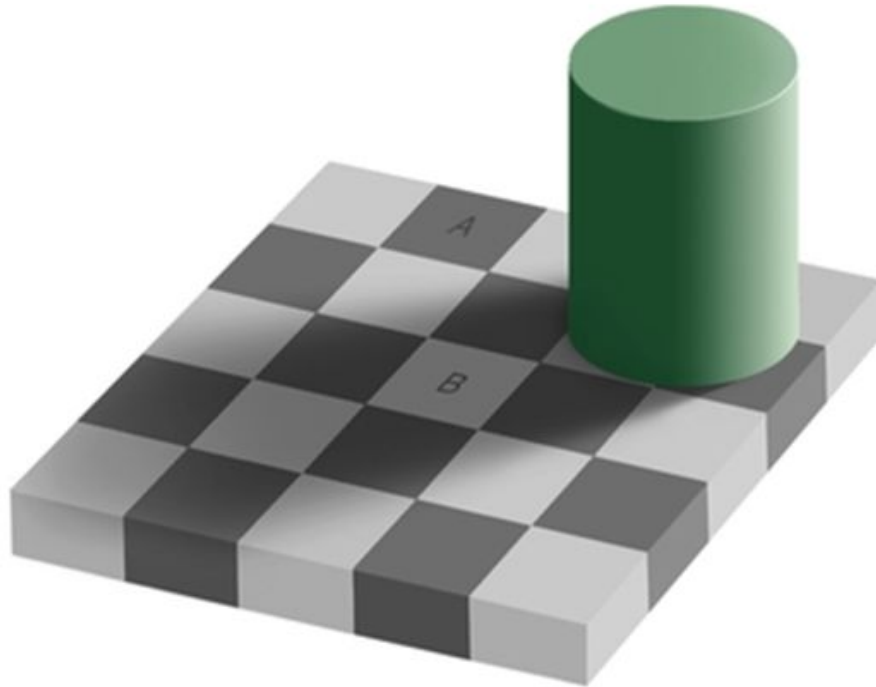
# Computer Vision is Hard



Slide Credit: Hamilton Chong

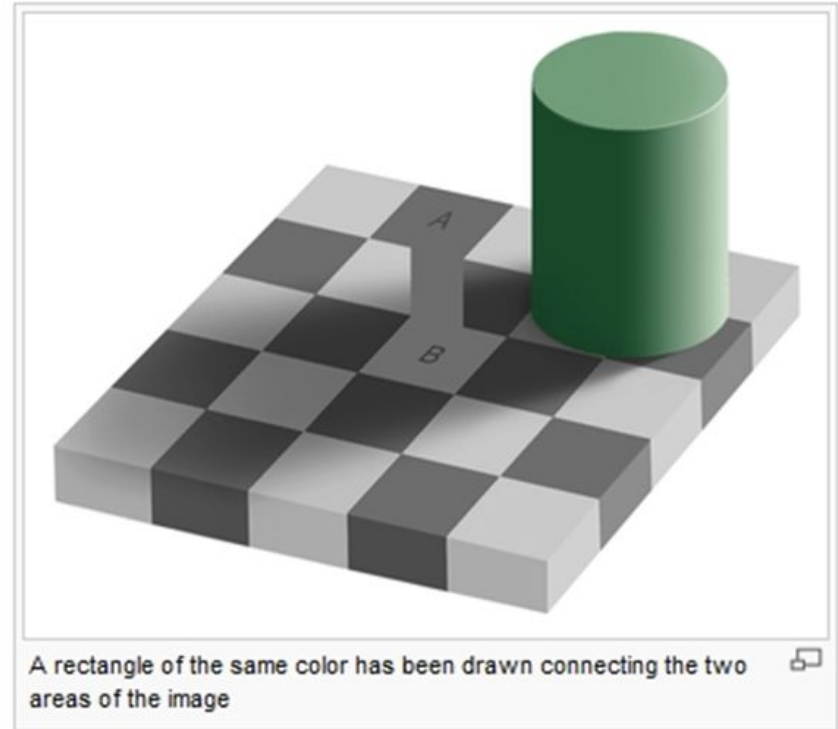
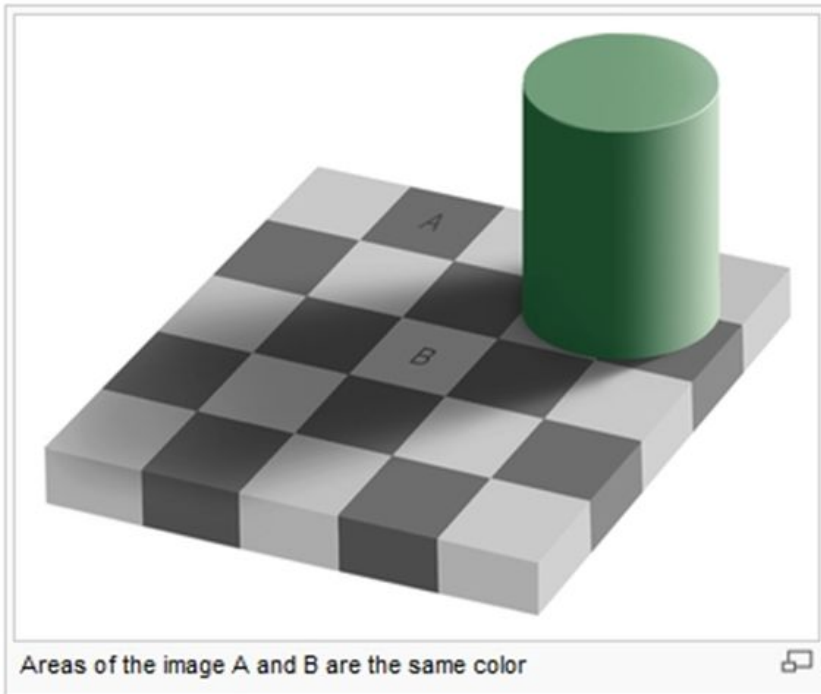


# Computer Vision is Hard



**Is square  
A or B  
darker in  
color?**

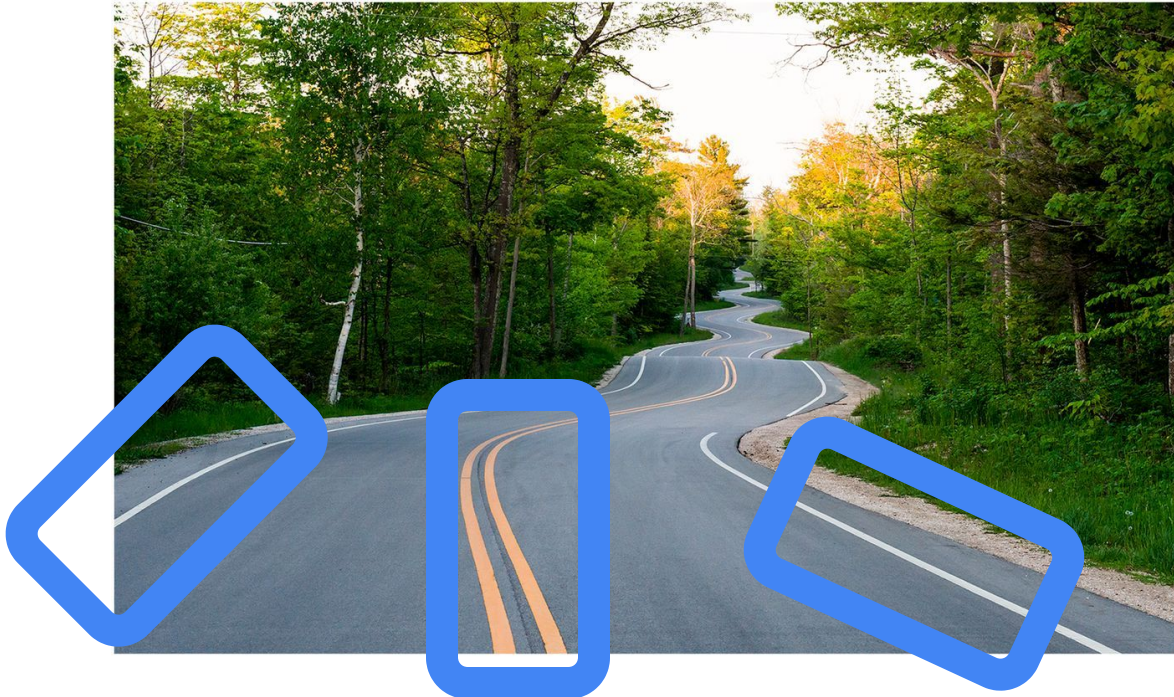
# Computer Vision is Hard



What **Features** of the image might be important for self driving cars?



What **Features** of the image might be important for self driving cars?

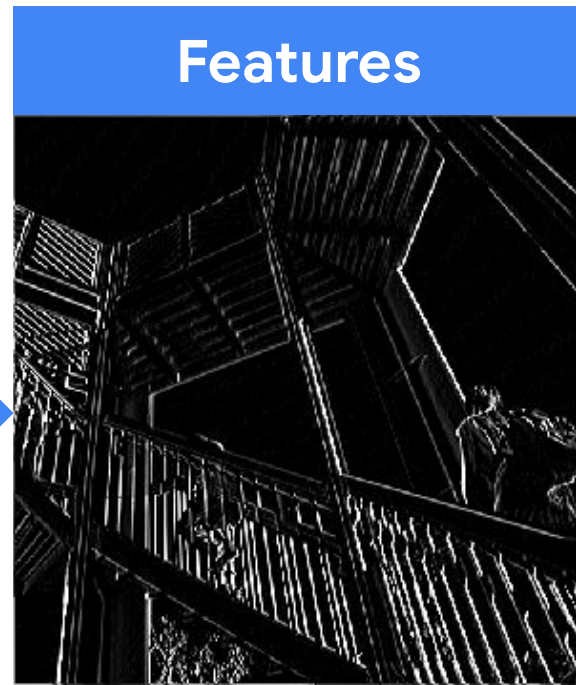


**Maybe  
straight  
lines to  
see the  
lanes  
of the  
road?**

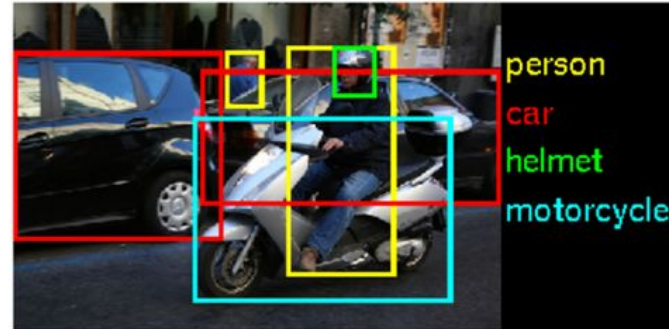
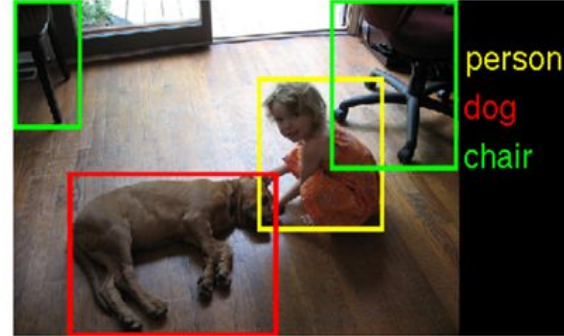
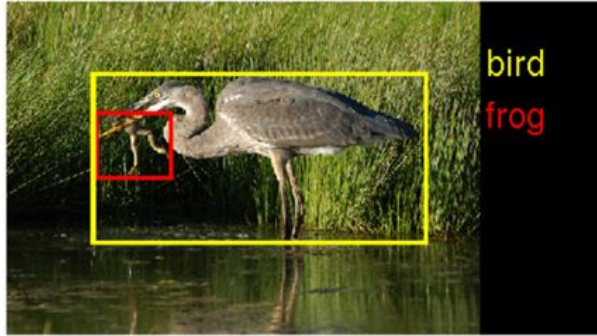
# Features can be found with **Convolutions**



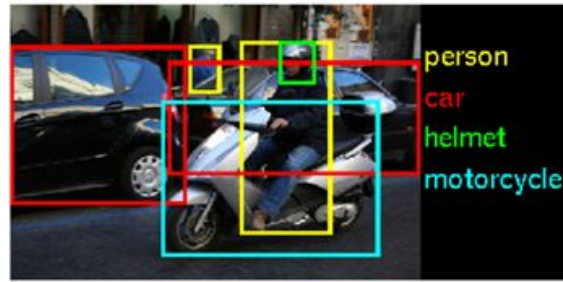
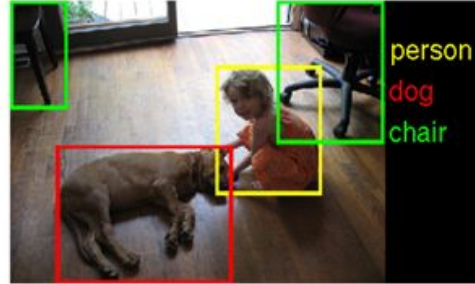
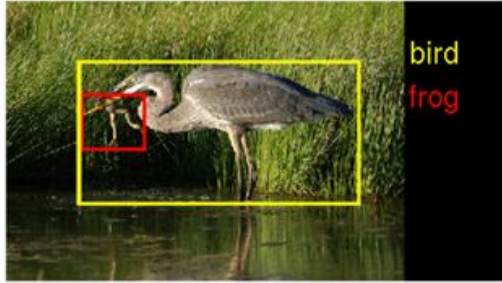
-1	0	1
-2	0	2
-1	0	1



# What features are needed for Object Detection?



# What features are needed for Object Detection?



The ImageNet Challenge provided 1.2 million examples of 1,000 **labeled** items and challenged algorithms to learn from the data and then was tested on another 100,000 images

# What features are needed for Object Detection?

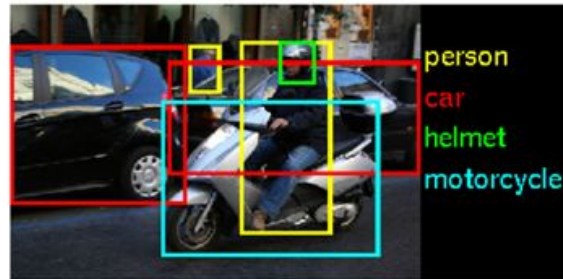
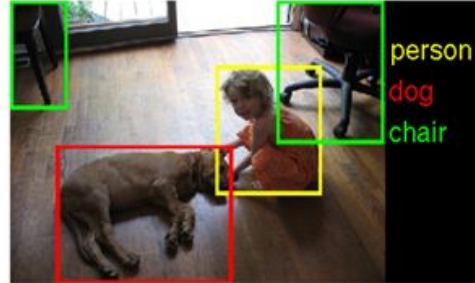
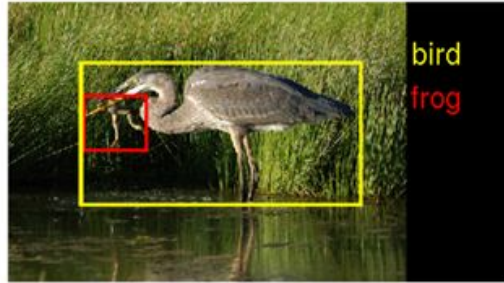


**Regression, Clustering, etc.**

**Vertical Lines, Horizontal Lines,  
Changes in Color, Changes in  
Focus, etc.**



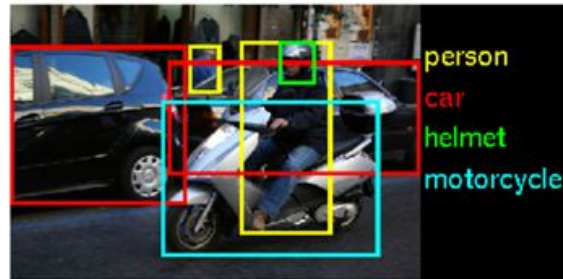
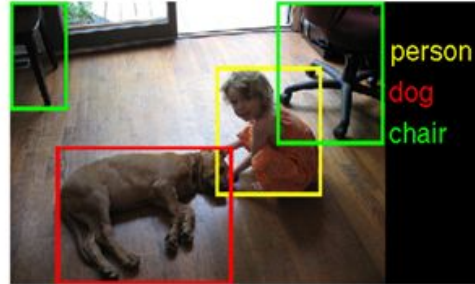
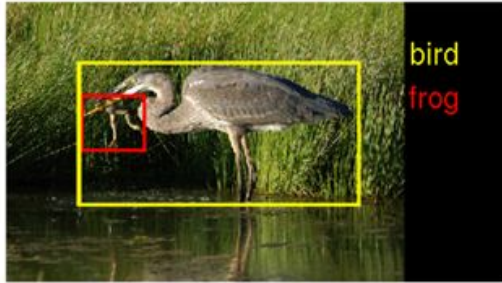
# What features are needed for Object Detection?



In 2010 teams had  
**75-50%** error

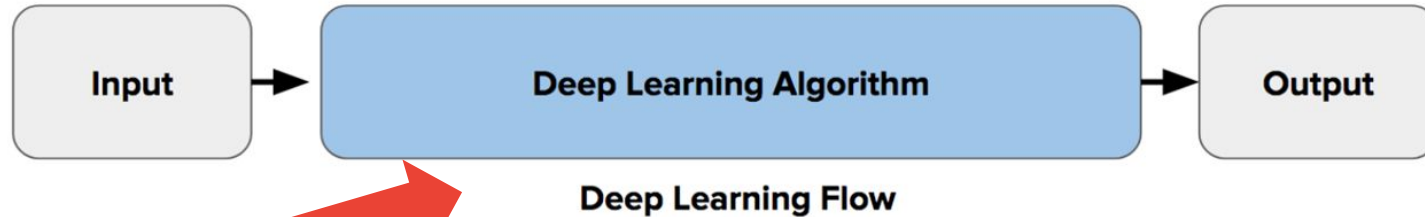
In 2011 teams had  
**75-25%** error

# What features are needed for Object Detection?



In 2012 still no team  
had less than 25%  
error barrier except  
**AlexNet at 15%**

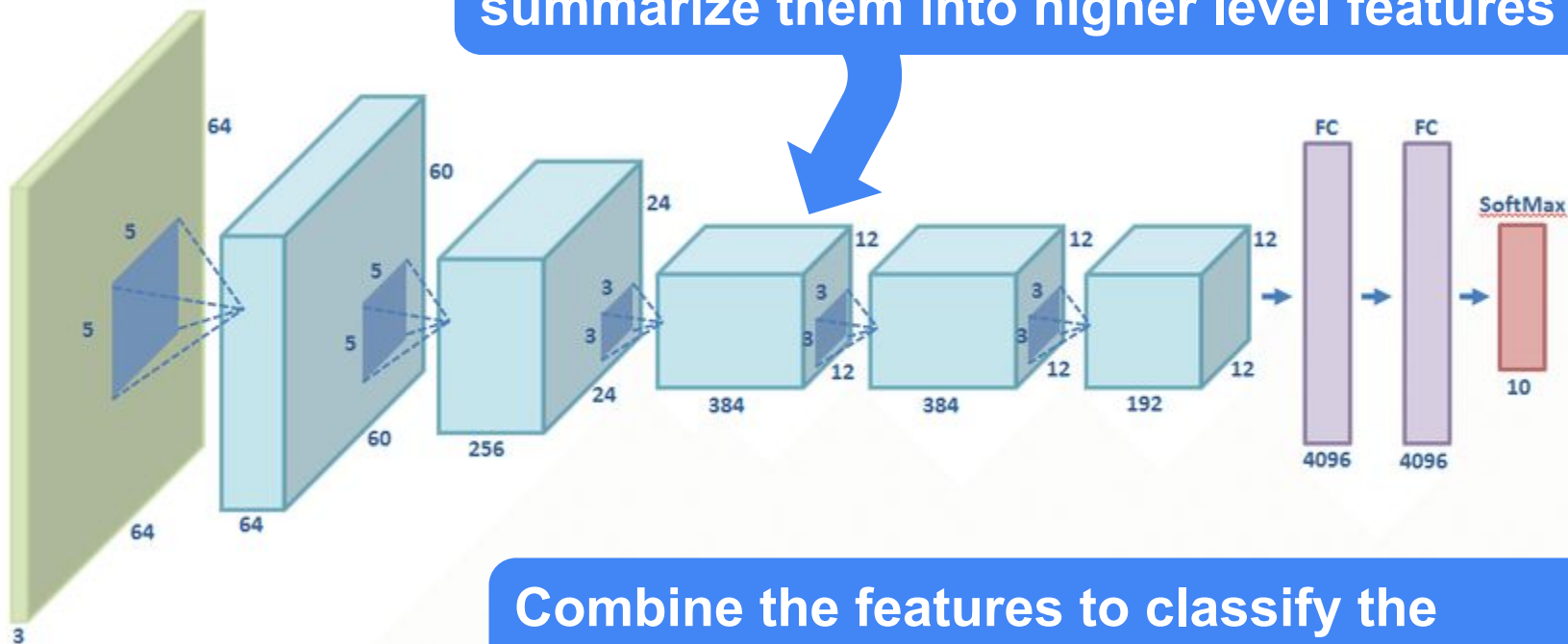
# What features are needed for Object Detection?



**Let the computer figure out its own features and how to combine them!**

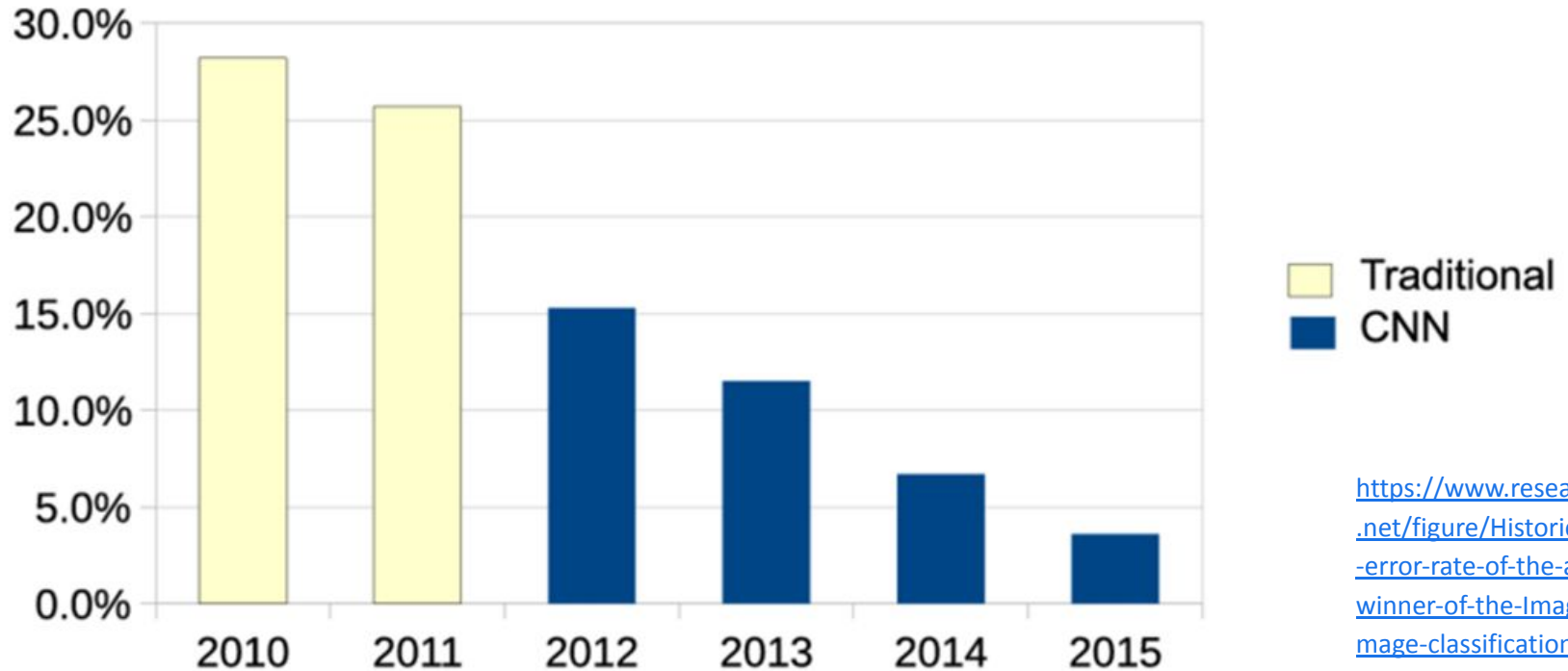
# AlexNet

Use convolutions to find features and the summarize them into higher level features



Combine the features to classify the various objects in the dataset

# What features are needed for Object Detection?



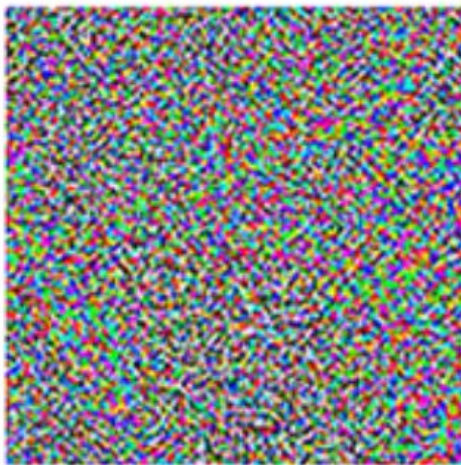
[https://www.researchgate.net/figure/Historical-top5-error-rate-of-the-annual-winner-of-the-ImageNet-image-classification\\_fig7\\_303992986](https://www.researchgate.net/figure/Historical-top5-error-rate-of-the-annual-winner-of-the-ImageNet-image-classification_fig7_303992986)

# A word of caution...

Ackerman "Hacking the Brain With Adversarial Images"



+  $\epsilon$



=



"panda"

57.7% confidence

There is **no model** of  
the world semantically  
just mathematically

"gibbon"

99.3% confidence

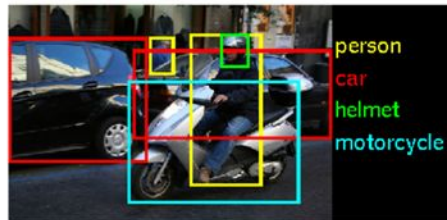
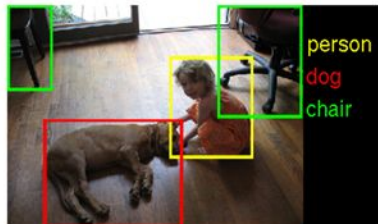
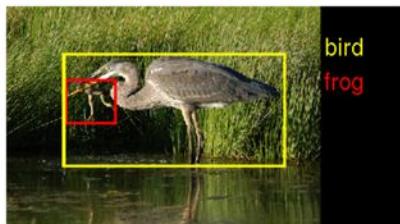
# Today's Agenda

- Deep ML Background
- **Hands-on Computer Vision: Thing Translator**
- The Tiny Machine Learning Workflow
- Keyword Spotting (KWS) Data Collection
- KWS Preprocessing and Training
- Deployment Challenges and Opportunities for Embedded ML
- Summary

# The **Thing** Translator

Open On Your Phone

<https://thing-translator.appspot.com/>

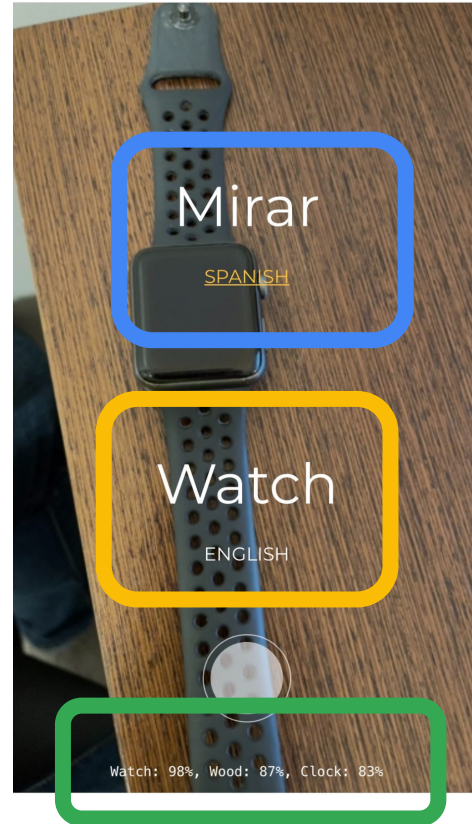




# The **Thing** Translator

[https://thing-translator.  
appspot.com/](https://thing-translator.appspot.com/)

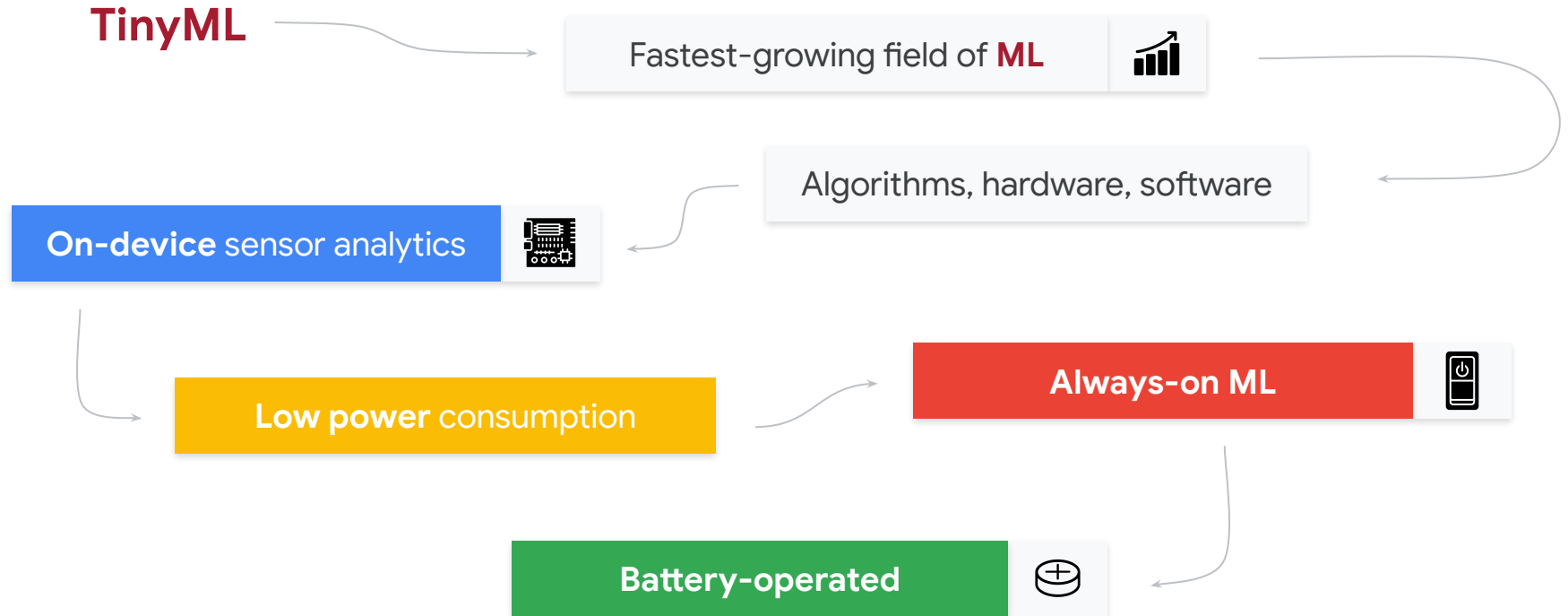
**Open On Your Phone**

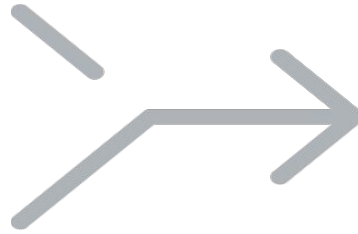
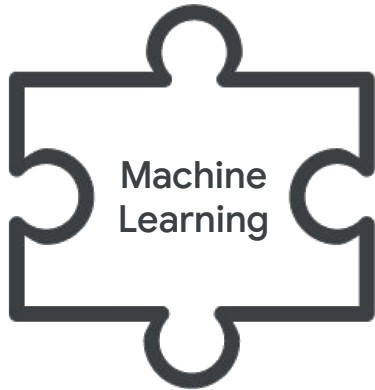
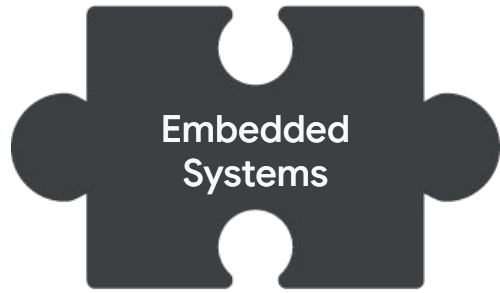


# Today's Agenda

- Deep ML Background
- Hands-on Computer Vision: Thing Translator
- **The Tiny Machine Learning Workflow**
- Keyword Spotting (KWS) Data Collection
- KWS Preprocessing and Training
- Deployment Challenges and Opportunities for Embedded ML
- Summary

# What is Embedded Machine Learning (**TinyML**)?





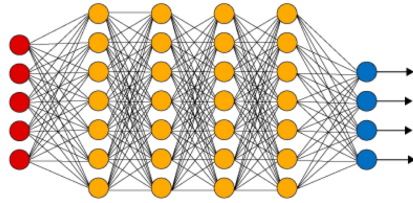
# TinyML



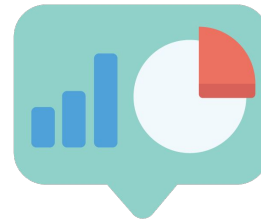
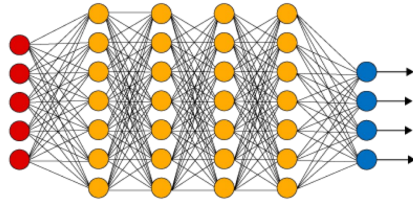
# The **TinyML** Workflow



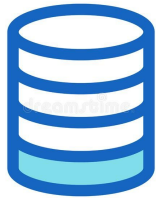
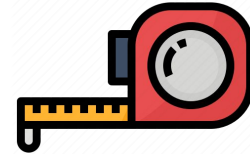
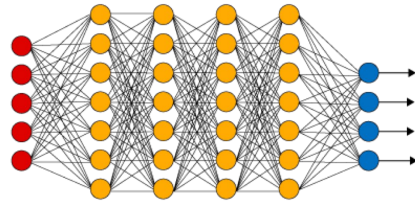
# The **TinyML** Workflow



# The **TinyML** Workflow

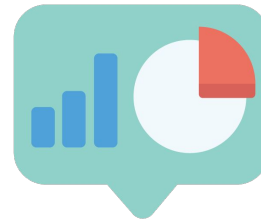
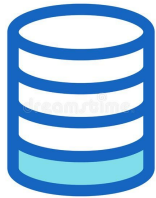
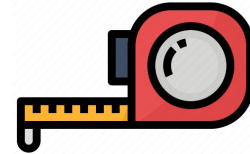
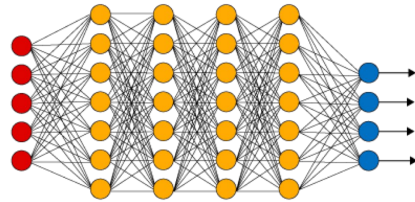


# The **TinyML** Workflow

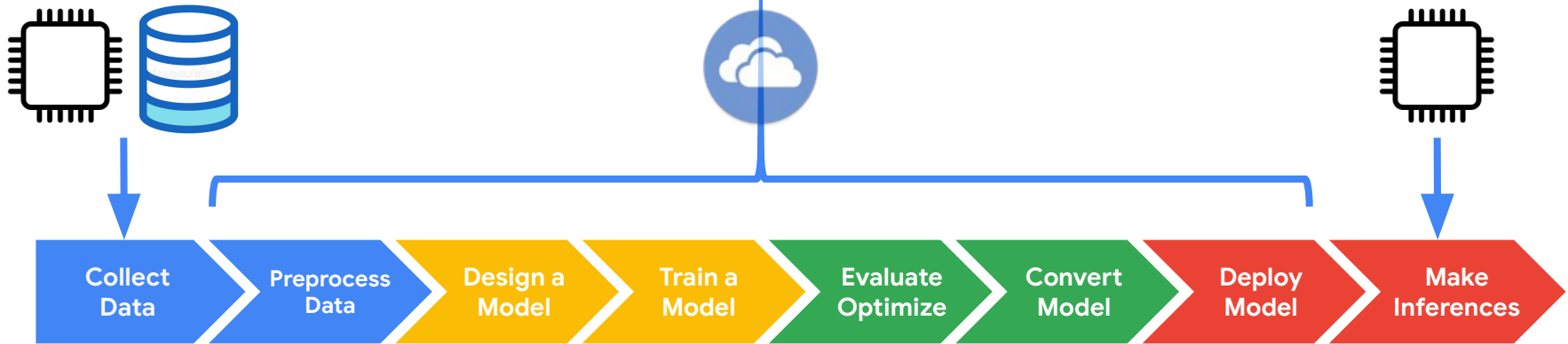




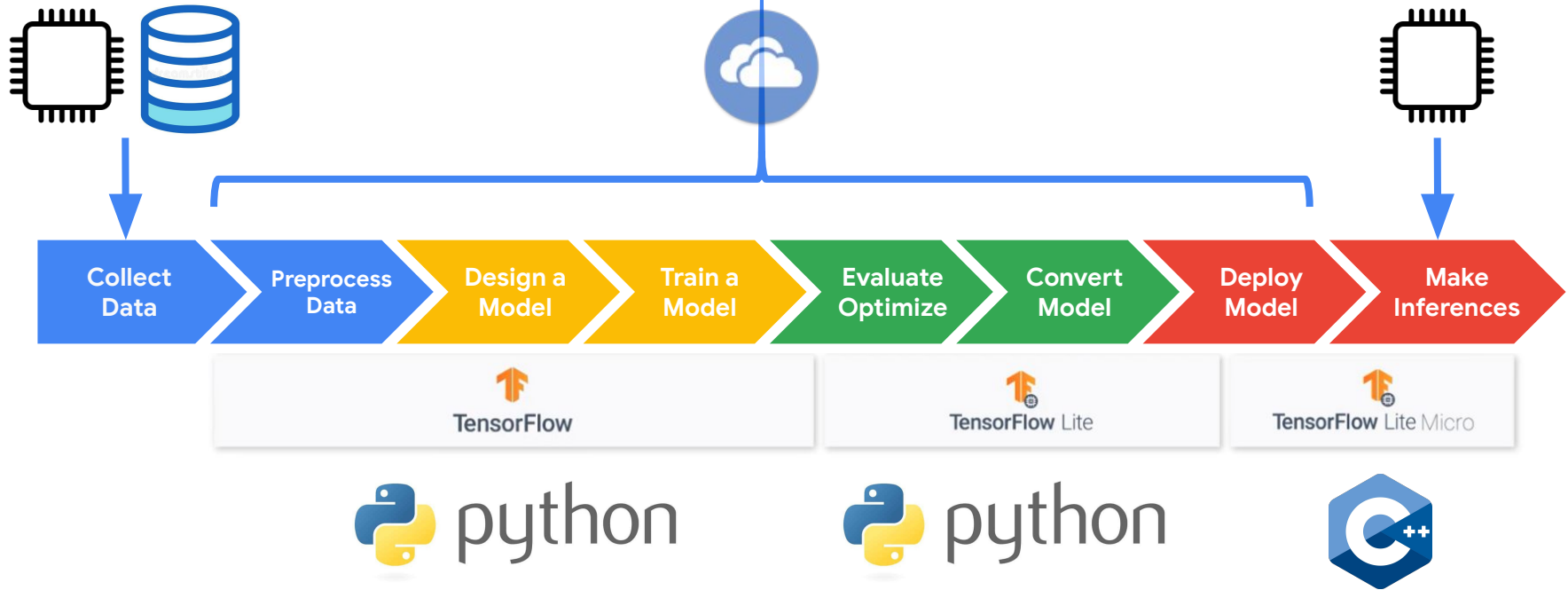
# The **TinyML** Workflow (“What”)



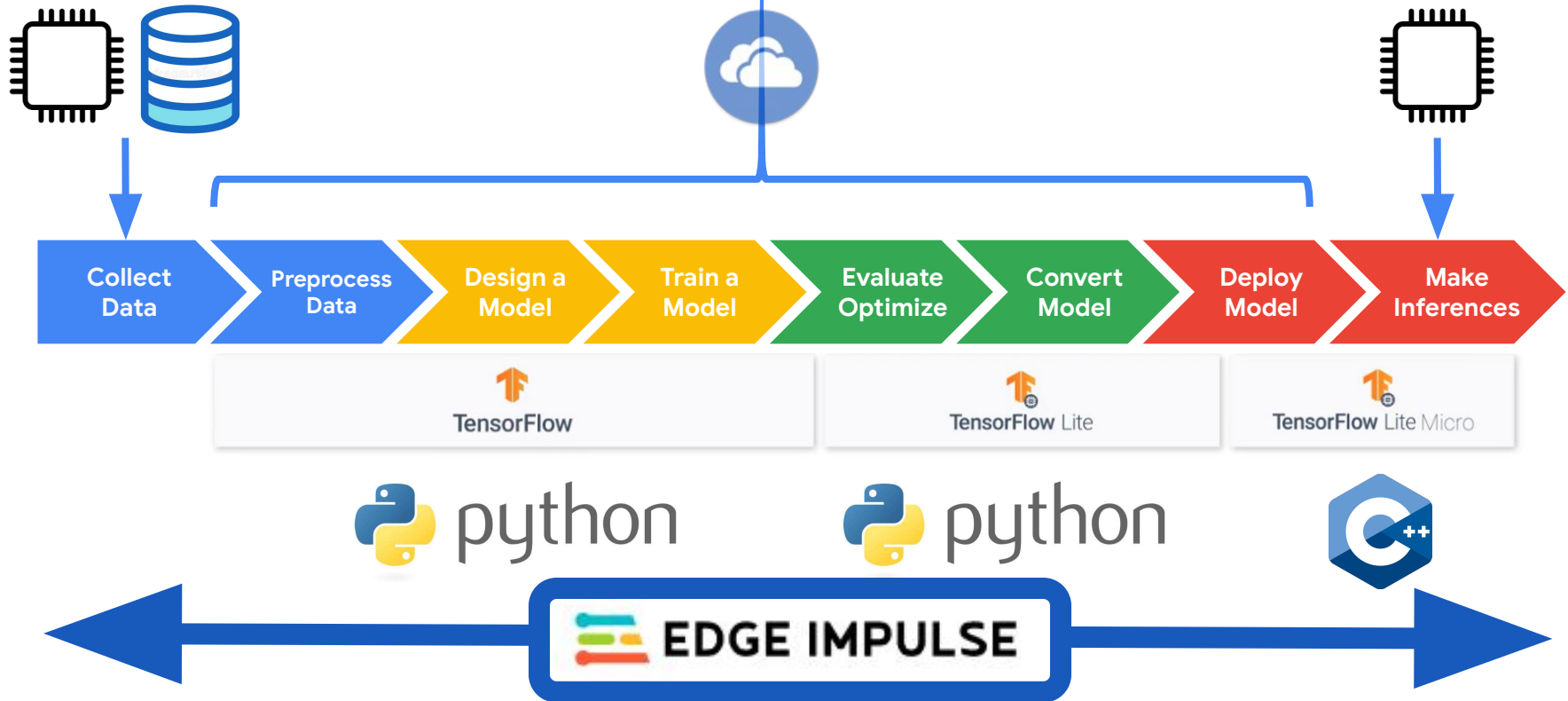
# The **TinyML** Workflow (“Where”)



# The **TinyML** Workflow (“How”)



# The **TinyML** Workflow (“How”)

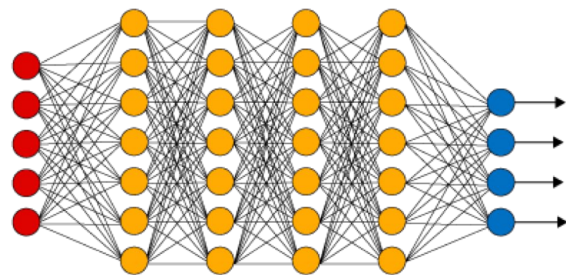


# Today's Agenda

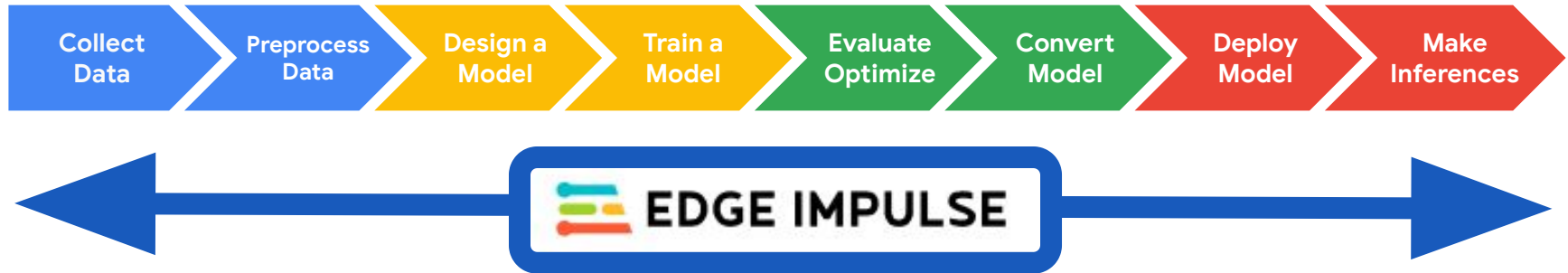
- Deep ML Background
- Hands-on Computer Vision: Thing Translator
- The Tiny Machine Learning Workflow
- **Keyword Spotting (KWS) Data Collection**
  - A Quick Primer on Data Engineering
  - Hands-on KWS Data Collection with Edge Impulse
- KWS Preprocessing and Training
- Deployment Challenges and Opportunities for Embedded ML
- Summary

# Keyword Spotting in One Slide

If we **pick a simple task** to only identifying a **few key words** we can then use a **small model** and train it with **little data** and fit it onto an **embedded device**



# The **TinyML** Workflow



# Today's Agenda

- Deep ML Background
- Hands-on Computer Vision: Thing Translator
- The Tiny Machine Learning Workflow
- Keyword Spotting (KWS) Data Collection

## **A Quick Primer on Data Engineering**

### Hands-on KWS Data Collection with Edge Impulse

- KWS Preprocessing and Training
- Deployment Challenges and Opportunities for Embedded ML
- Summary



# Data Engineering for KWS (Part 1)

## (How to collect good data)

# Data Engineering for KWS (Part 1)

(How to collect good data)

**Who** will use your  
ML model?

- What **languages** will they speak?
- What **accents** will they have?
- Will they use **slang** or formal diction?

# Data Engineering for KWS (Part 1)

## (How to collect good data)

**Who** will use your ML model?

- What **languages** will they speak?
- What **accents** will they have?
- Will they use **slang** or formal diction?

**Where** will your ML model be used?

- Will there be **background noise**?
- **How far** will users be from the microphone?
- Will there be **echos**?

# Data Engineering for KWS (Part 1)

## (How to collect good data)

**Who** will use your ML model?

- What **languages** will they speak?
- What **accents** will they have?
- Will they use **slang** or formal diction?

**Where** will your ML model be used?

- Will there be **background noise**?
- **How far** will users be from the microphone?
- Will there be **echos**?

**Why** will your ML model be used?  
**Why** those Keywords?

- What **tone of voice** will be used?
- Are your **keywords commonly** used? (aka will you get a lot of false positives)
- What about false negatives?

# Data Engineering for KWS (Part 1)

(How to collect good data)

There are a lot more things to consider to **eliminate bias** and **protect privacy** when collecting data that we will talk about in future sessions!

ML model be used?

**Why** those Keywords?

the **you get a lot of false positives** (aka will

- What about false negatives?

# Tips and Tricks for Custom KWS

- Pick **uncommon words** for Keywords
- Record lots of “**other words**”
- Record in the **location** you are going to be **deploying**
- Get **your end users** to help you build a dataset
- Record with the same **hardware** you will **deploy**
- Always **test** and then **improve** your dataset and model

# Tips and Tricks for Custom KWS

Today we are just working on a demo so to give our demo the the best chance of working we will:

1. **Stay in one spot** (we're cheating)
2. **Only record ourselves**
3. **Use common words (yes, no)**
4. **Only test ourselves**

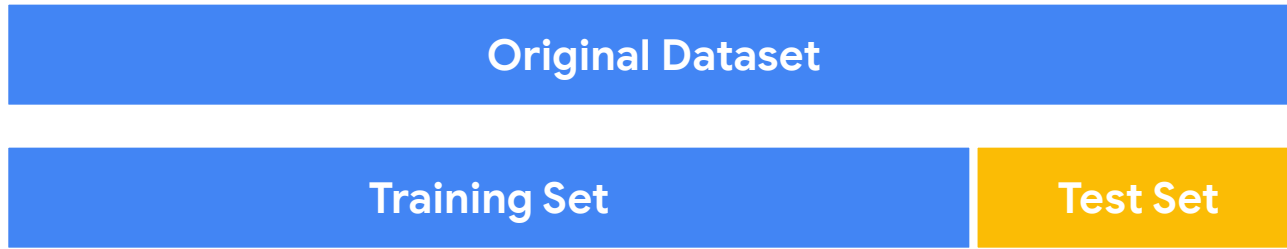
# Data Engineering for KWS (Part 2)

(how to test with our data)



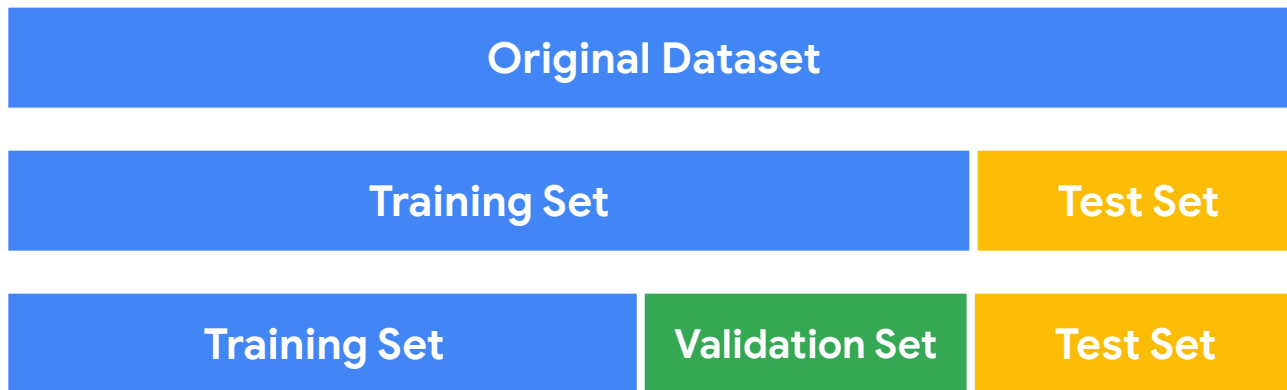
# Data Engineering for KWS (Part 2)

(how to test with our data)



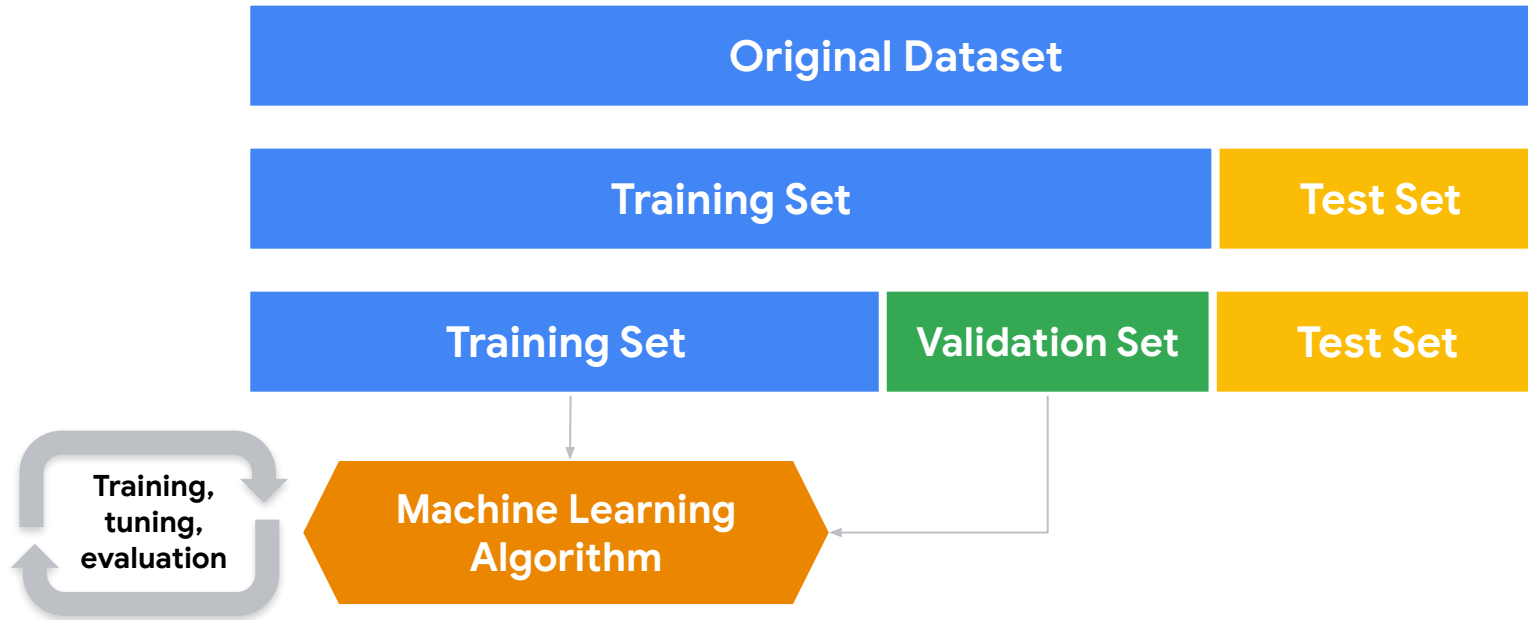
# Data Engineering for KWS (Part 2)

(how to test with our data)



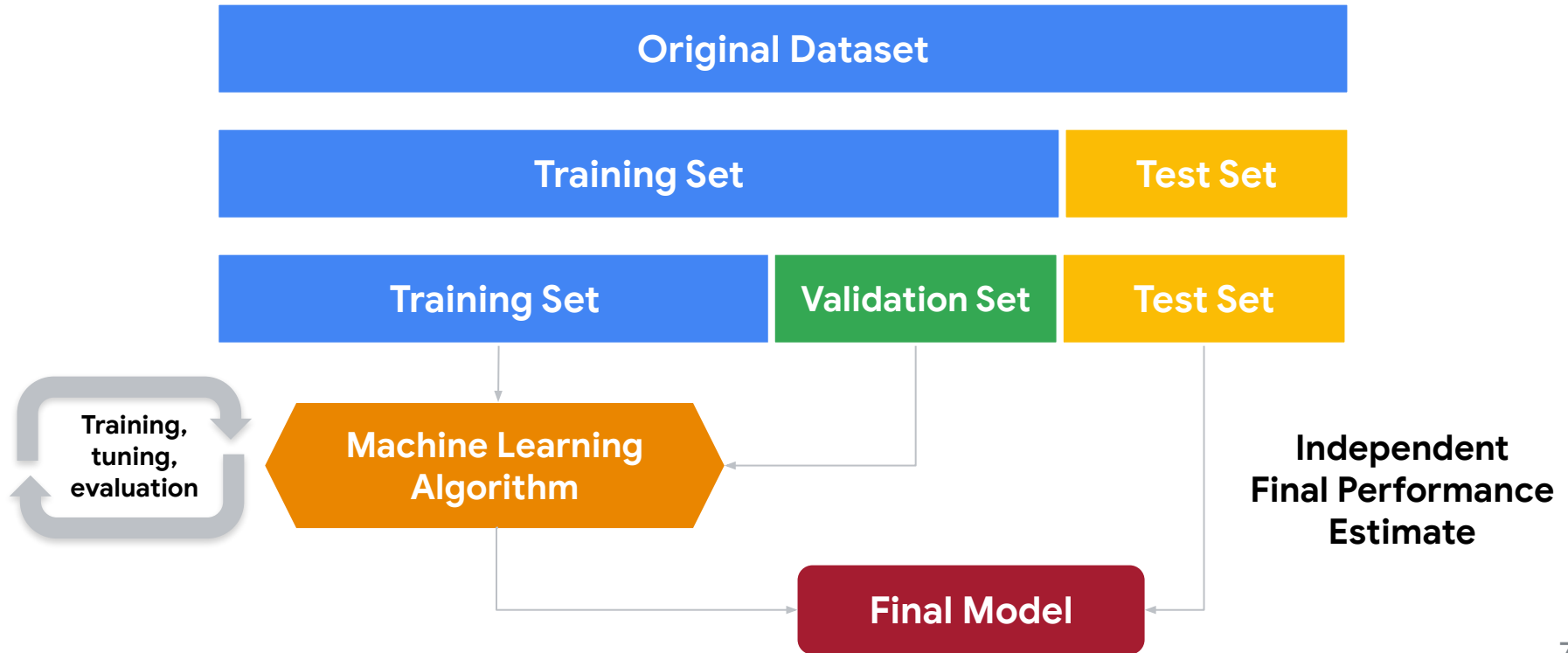
# Data Engineering for KWS (Part 2)

## (how to test with our data)



# Data Engineering for KWS (Part 2)

## (how to test with our data)



# Today's Agenda

- Deep ML Background
- Hands-on Computer Vision: Thing Translator
- The Tiny Machine Learning Workflow
- Keyword Spotting (KWS) Data Collection

A Quick Primer on Data Engineering

## **Hands-on KWS Data Collection with Edge Impulse**

- KWS Preprocessing and Training
- Deployment Challenges and Opportunities for Embedded ML
- Summary

# The **TinyML** Workflow using **Edge Impulse**

Today we'll also collect all of our data using Edge Impulse...

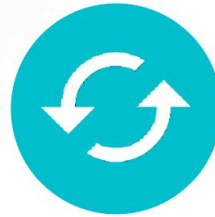
...and deploy to your cell phone as well



**Dataset**



**Impulse**



**Test**



**Deploy**

# Create an Edge Impulse Account

1. Create an Edge Impulse account:  
<https://studio.edgeimpulse.com/signup>
2. Validate your email by clicking the link in the email sent to your account's email address

 EDGE IMPULSE

Log in

[Forgot your password?](#)

Log in

Don't have an account? [Sign up](#)



Start building embedded  
machine learning  
models today.

© 2021 EdgeImpulse Inc. All rights reserved


# Select project

Select your Edge Impulse project, or create a new one.

NAME

COLLABORATORS

WELCOME!

 Your profile

ORGANIZATIONS

 Harvard University

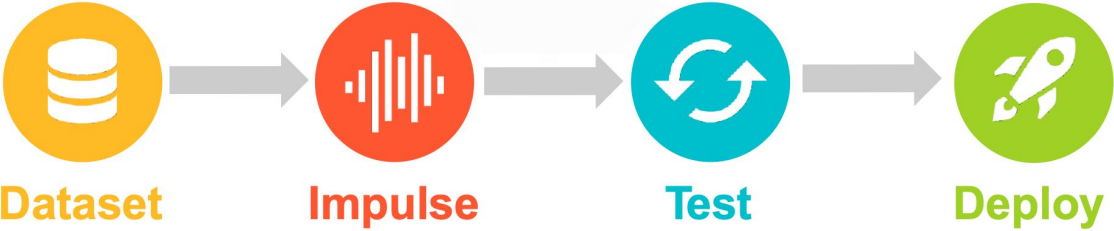
PROJECTS



 Create new project



# Edge Impulse Project Dashboard



- Dashboard
- Devices
- Data acquisition
- Impulse design
- Create impulse
- EON Tuner
- Retrain model
- Live classification
- Model testing
- Versioning
- Deployment



## **Activity:** Create a Keyword Spotting Dataset

Collect **~30 samples each** of the following classes of data:

- **Keyword #1 “yes”** (label: yes) (length: 2 seconds)
- **Keyword #2 “no”** (label: no) (length: 2 seconds)
- **“Unknown” words** that are not the keyword **and background noise** (label: unknown) (length: 2 seconds)

- Dashboard
- Devices
- Data acquisition
- Impulse design
  - Create impulse
- EON Tuner
- Retrain model
- Live classification
- Model testing

This is your Edge Impulse project. From here, you acquire new training data, design impulses and train models.

## Creating your first impulse (0% complete)



### Acquire data

Every Machine Learning project starts with data. You can capture data from [your device](#) or import data you already collected.

 [LET'S COLLECT SOME DATA](#)



### Design an impulse

Teach the model to interpret previously unseen data, based on historical data. Use this to categorize new data, or to find anomalies in sensor readings.

## Collect data

You can collect data from development boards, from your own devices, or by uploading an existing dataset.



### Connect a fully supported development board

Get started with real hardware from a wide range of silicon vendors - fully supported by Edge Impulse.

[Browse dev boards](#)



### Use your mobile phone

Use your mobile phone to capture movement, audio or images, and even run your trained model locally. No app required.

[Show QR code](#)



### Use your computer

Capture audio or images from your webcam or microphone, or from an external audio device.

[Collect data](#)



### Data from any device with the data forwarder

Capture data from any device or development board over a serial connection, in 10 lines of code.

[Show docs](#)



### Upload data

Already have data? You can upload your existing datasets directly in WAV, JPG, PNG, CBOR, CSV or JSON format.

[Go to the uploader](#)

## Collect data

You can collect data from development boards, from your own devices, or by uploading an existing dataset.



### Connect a fully supported development board

Get started with real hardware from a wide range of silicon vendors - fully supported by Edge Impulse.

[Browse dev boards](#)



**Point your phone camera at the QR code and open the link!**

images, and even

[Show QR code](#)



phone, or from an

[Collect data](#)



Capture data from any device or development board over a serial connection, in 10 lines of code.

[Show docs](#)



### Upload data

Already have data? You can upload your existing datasets directly in WAV, JPG, PNG, CBOR, CSV or JSON format.

[Go to the uploader](#)



### Connected as phone\_kunh8zjd

You can collect data from this device from the **Data acquisition** page in the Edge Impulse studio.

📷 Collecting images?

🎤 Collecting audio?

~ Collecting motion?



### Connected as phone\_kunh8zjd

You can collect data from this device from the **Data acquisition** page in the Edge Impulse studio.

📷 Collecting images?

🎤 Collecting audio?

📊 Collecting measurements?



### Connected as phone\_kunh8zjd

You can collect data from this device from the **Data acquisition** page in the Edge Impulse studio.

 Collecting images?

 Collecting audio?

 Collecting motion?



### Data collection

Label: goodbye

Length: 3s.

Category: split

 Start recording

Audio captured with current settings: 0s





### Connected as phone\_kunh8zjd

You can collect data from this device from the **Data acquisition** page in the Edge Impulse studio.

 Collecting images?

 Collecting audio?


 Collecting motion?



### Data collection


Label: goodbye

length: 3s.

 Start recording

Audio captured with current settings: 0s

smartphone.edgeimpulse.com



Connected as phone\_kunh8zjd

You can collect data from this device from the **Data acquisition** page in the Edge Impulse studio.

Collecting images?

Collecting audio?

Collecting motion?

smartphone.edgeimpulse.com

Data collection

Label: goodbye Length: 3s.  
Category: split

Start recording

Audio captured with current settings: 0s

smartphone.edgeimpulse.com

Data collection

Label: goodbye Length: 3s.  
Category: split

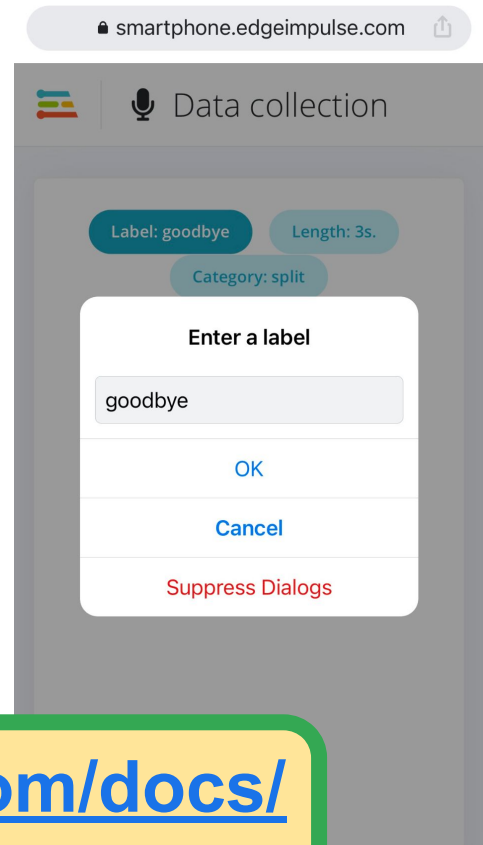
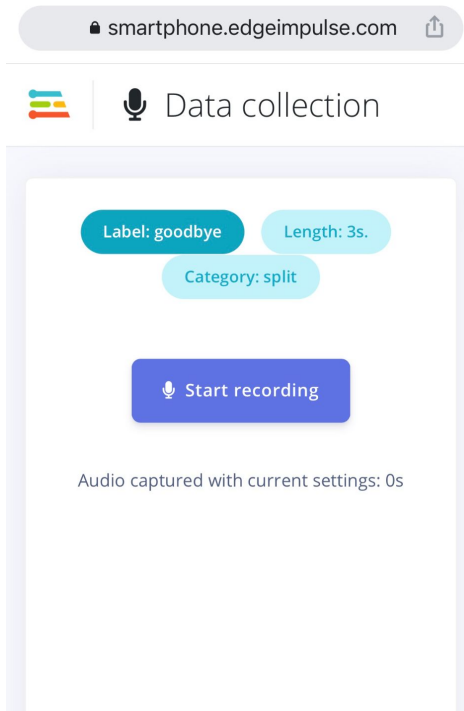
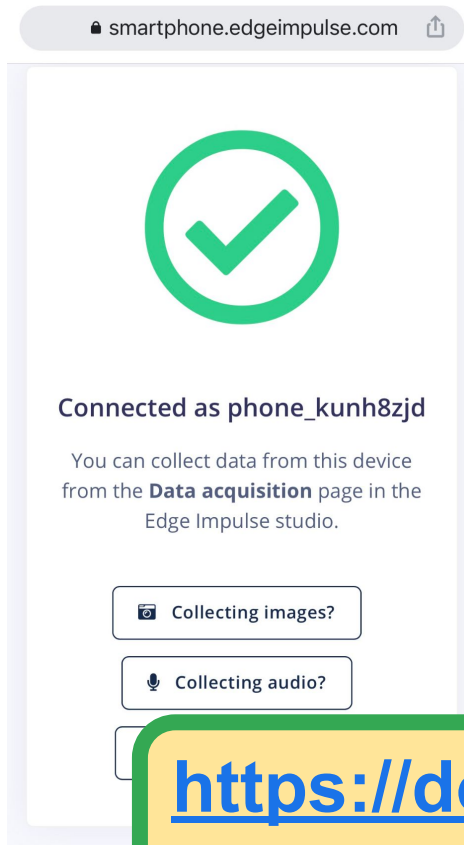
Enter a label

goodbye

OK

Cancel

Suppress Dialogs



<https://docs.edgeimpulse.com/docs/using-your-mobile-phone>

Training data

Test data

**Did you know?** You can capture data from any device or [upload your existing datasets](#) - [Show options](#)

DATA COLLECTED

1m 27s



TRAIN / TEST

86% / 14%

Collected data

SAMPLE NAME	LABEL	ADDED	Duration	Actions
unknown.2hvfrrhdt	unknown	Today, 16:45:06	2s	⋮
unknown.2hvfrrd4u	unknown	Today, 16:45:02	2s	⋮
unknown.2hvfrr8a4	unknown	Today, 16:44:57	2s	⋮
unknown.2hvfrrq4	unknown	Today, 16:44:47	2s	⋮
unknown.2hvfrrq15	unknown	Today, 16:44:43	2s	⋮
unknown.2hvfrrqmr3	unknown	Today, 16:44:39	2s	⋮
unknown.2hvfrrqj1g	unknown	Today, 16:44:35	2s	⋮
unknown.2hvfrrq9bn	unknown	Today, 16:44:25	2s	⋮

- Rename
- Edit label
- Move to test set
- Disable
- Crop sample
- Split sample
- Download
- Download (.WAV)
- Delete

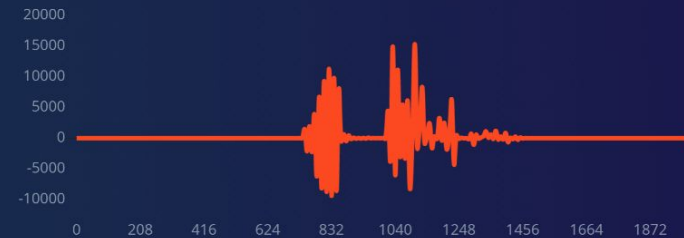
Record new data

Connect using WebUSB

No devices connected to the remote management API.

RAW DATA

unknown.2hvfrrhdt



audio

0:00 / 0:00 [Progress bar] [Volume icon] [More options icon]

- Dashboard
- Devices
- Data acquisition**
- Impulse design
  - Create impulse
- EON Tuner
- Retrain model
- Live classification
- Model testing
- Versioning
- Deployment

GETTING STARTED

- Documentation
- Forums

<https://docs.edgeimpulse.com/docs/using-your-mobile-phone>

## **Activity:** Create a Keyword Spotting Dataset

Collect **~30 samples each** of the following classes of data:

- **Keyword #1 “yes”** (label: yes) (length: 2 seconds)
- **Keyword #2 “no”** (label: no) (length: 2 seconds)
- **“Unknown” words** that are not the keyword **and background noise** (label: unknown) (length: 2 seconds)

**Also take a quick break! We'll  
resume in 10 minutes!**

# Today's Agenda

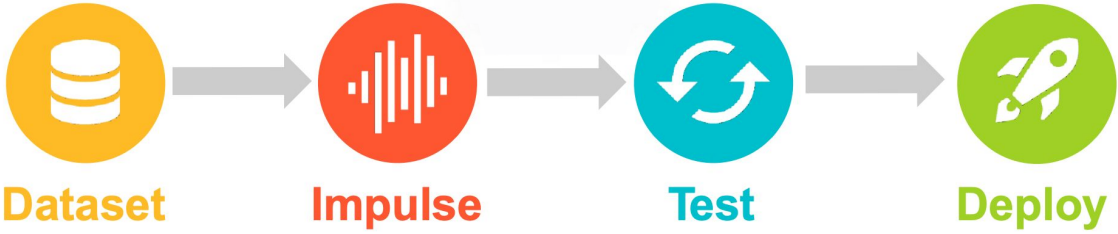
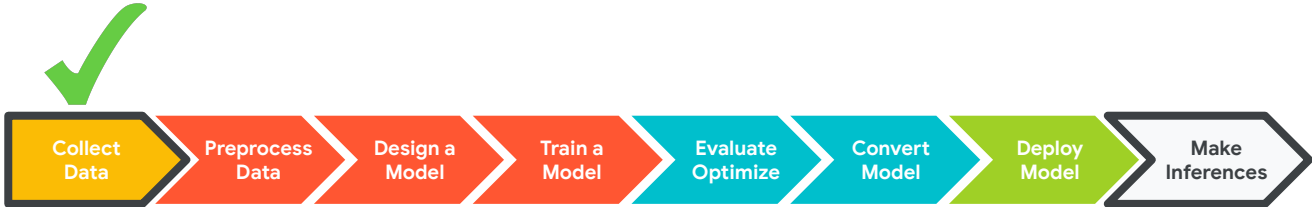
- Deep ML Background
- Hands-on Computer Vision: Thing Translator
- The Tiny Machine Learning Workflow
- Keyword Spotting (KWS) Data Collection
- **KWS Preprocessing and Training**

## **Preprocessing (for KWS)**

Hands-on Preprocessing and Training with Edge Impulse

- Deployment Challenges and Opportunities for Embedded ML
- Summary

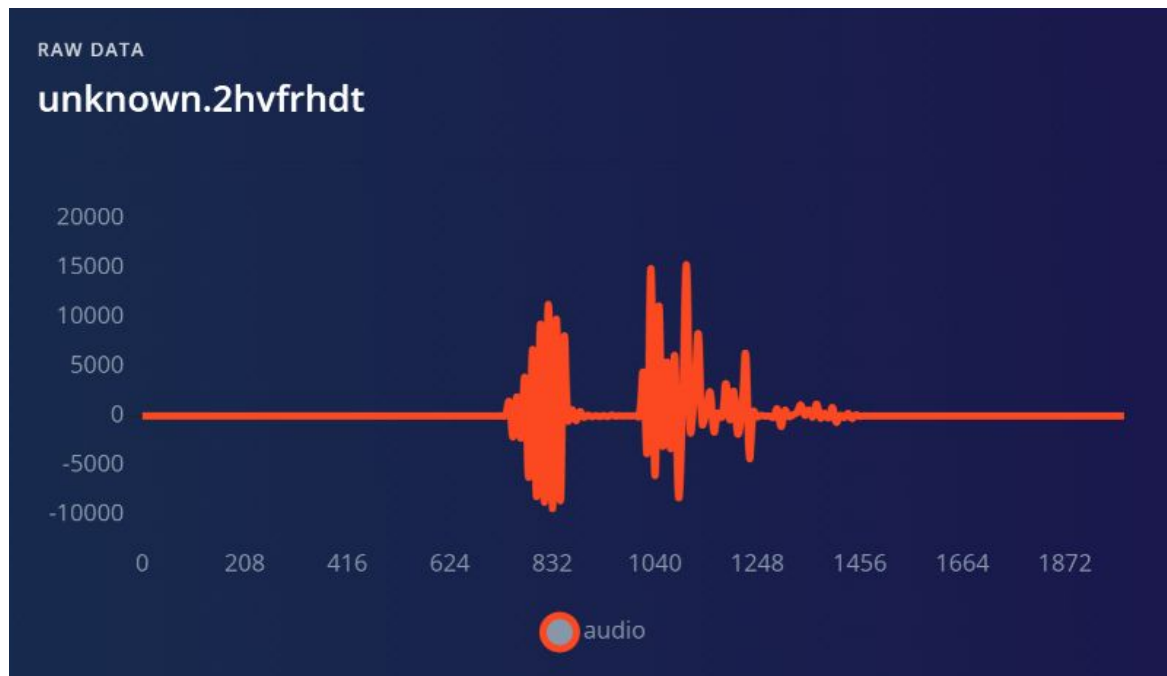
# Edge Impulse Project Dashboard



- Dashboard
- Devices
- Data acquisition
- Impulse design
- Create impulse
- EON Tuner
- Retrain model
- Live classification
- Model testing
- Versioning
- Deployment



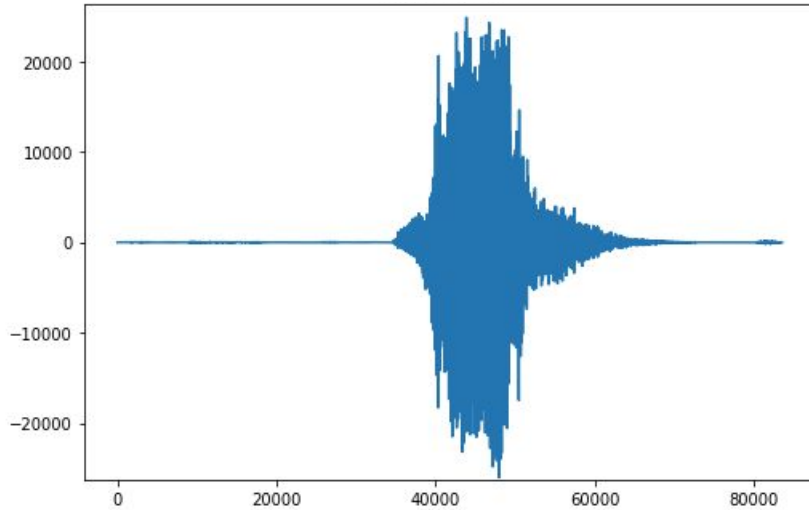
Why might we want to **preprocess** data and not send the raw data to the neural network?



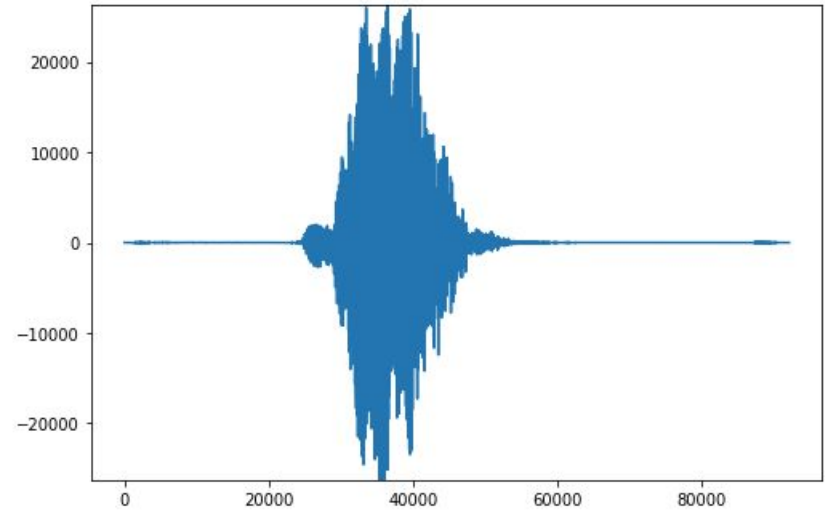


# Can you tell these two signals apart?

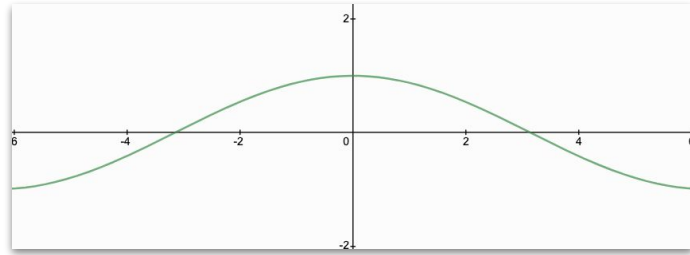
**“Yes”** (*spoken loudly*)



**“No”** (*spoken loudly*)

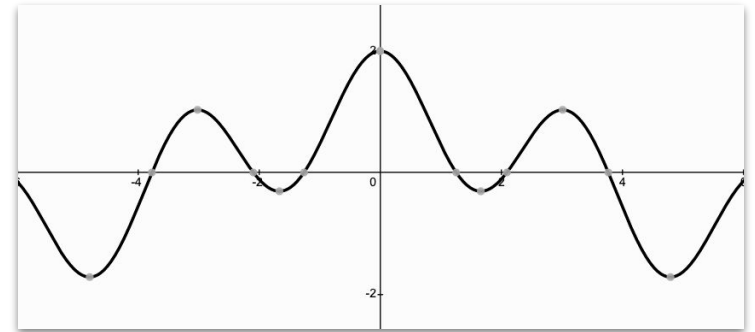
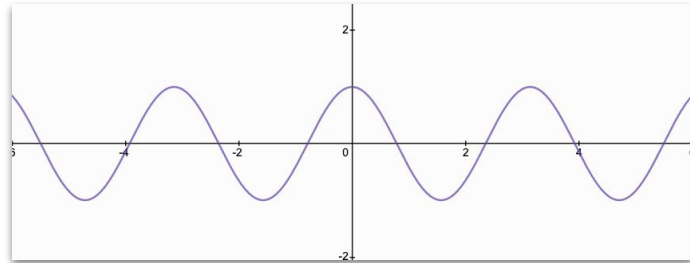


# Signal Components?



+

=



# Signal Components?

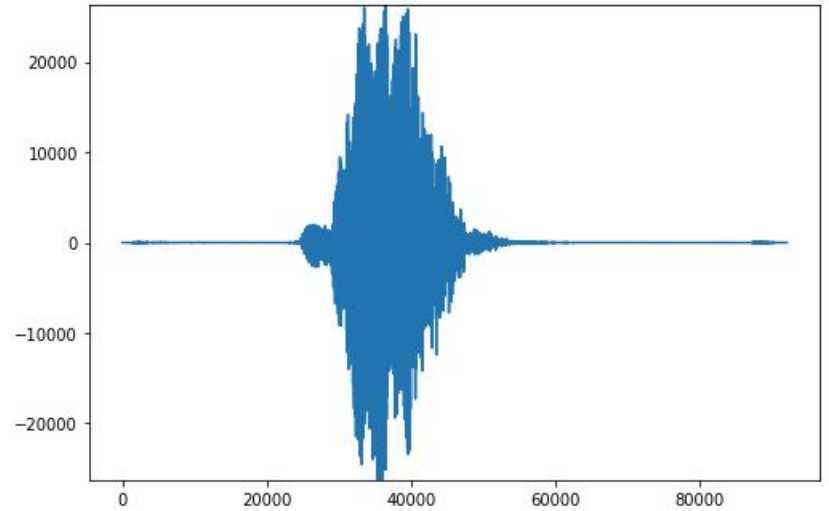
?

+

?

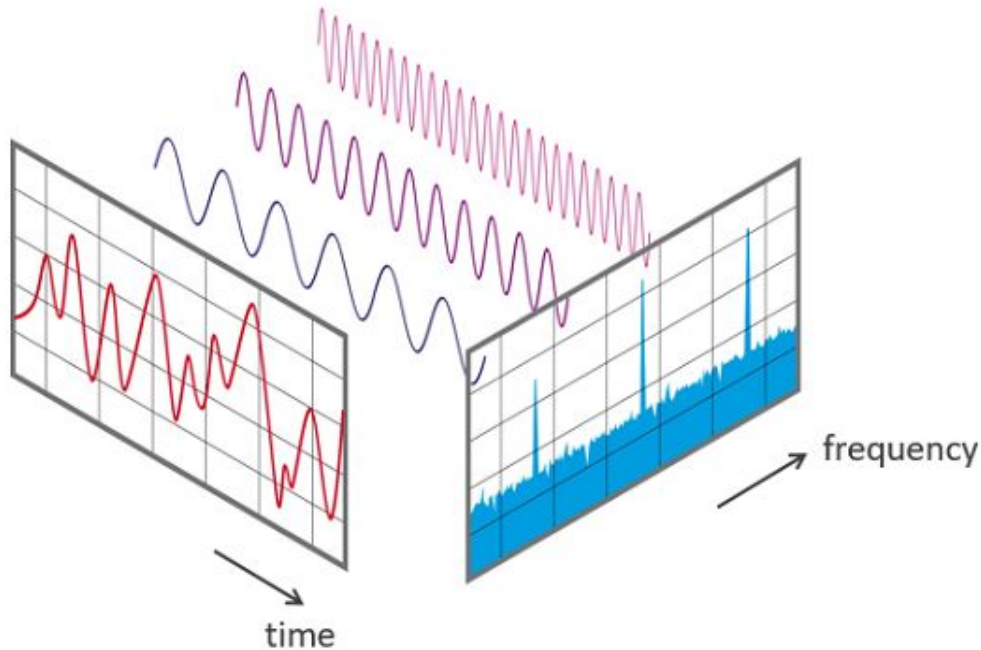
=

“No” (*spoken loudly*)



# Fast Fourier Transform:

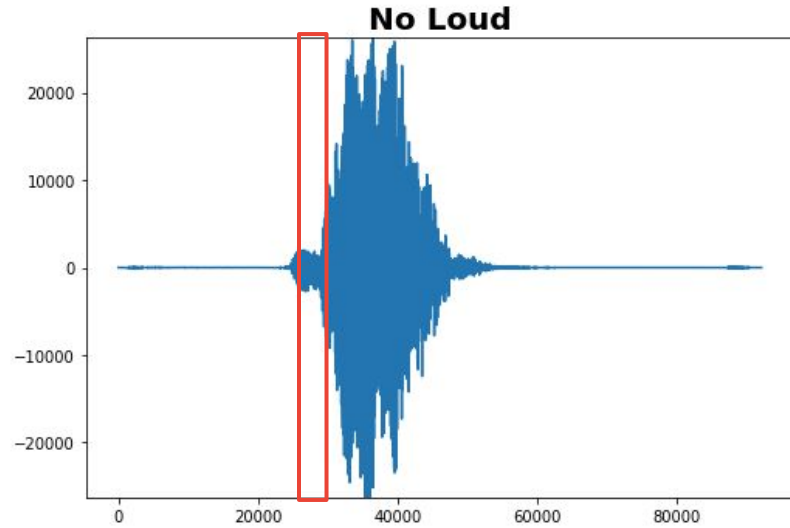
extract the frequencies from a signal



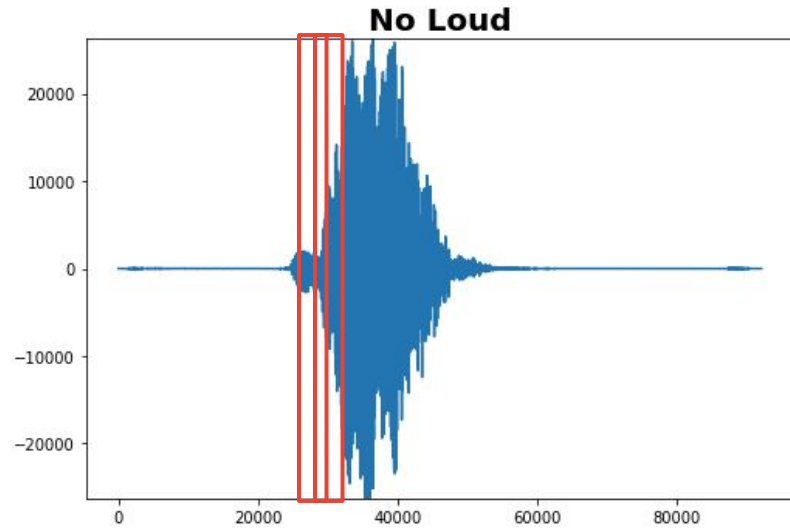
# Fast Fourier Transform



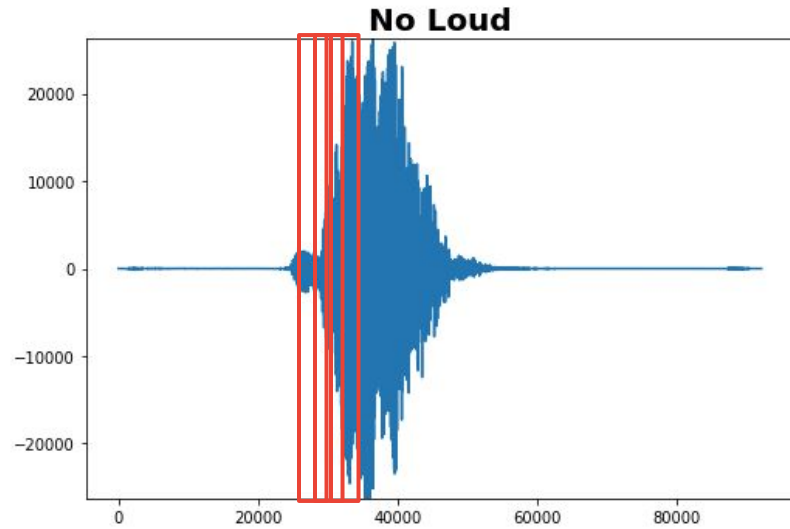
# Building a Spectrogram using FFTs



# Building a Spectrogram using FFTs

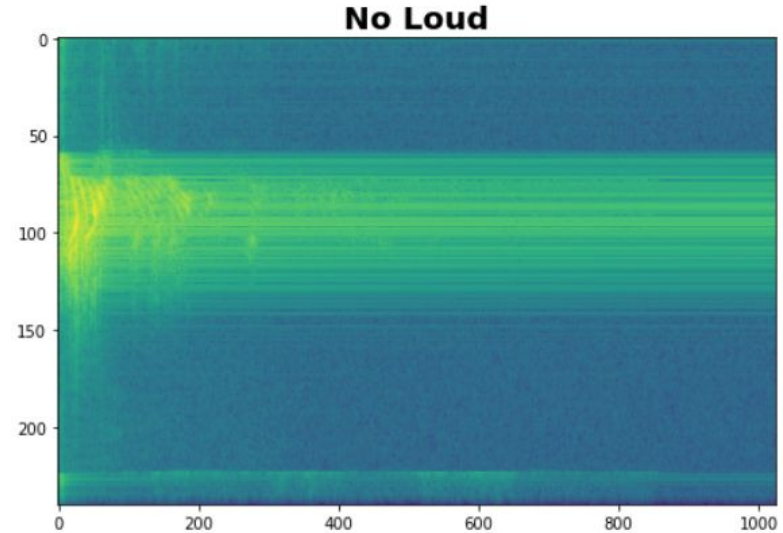
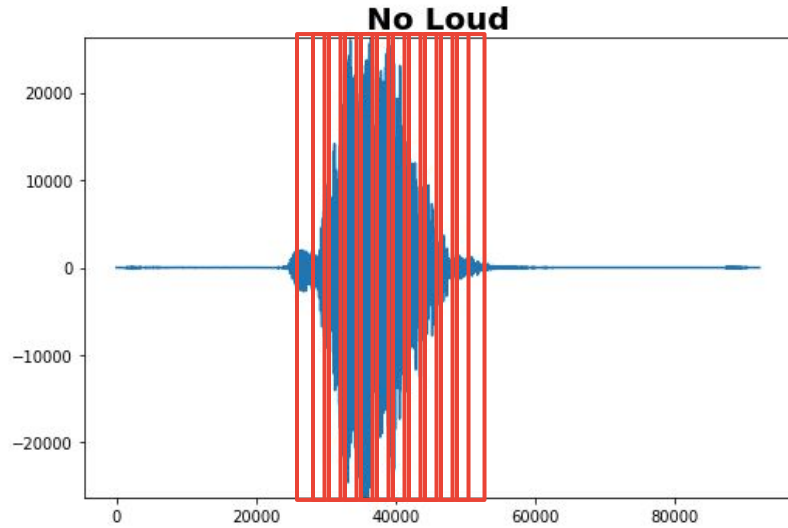


# Building a Spectrogram using FFTs



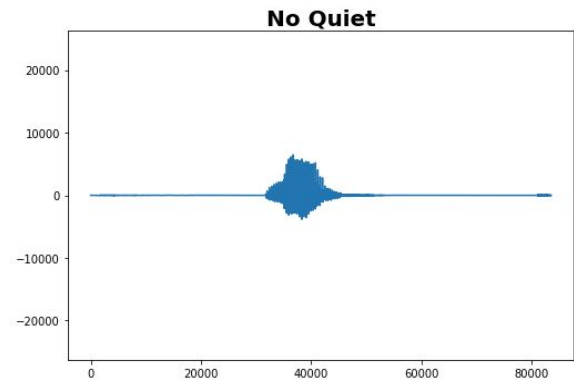
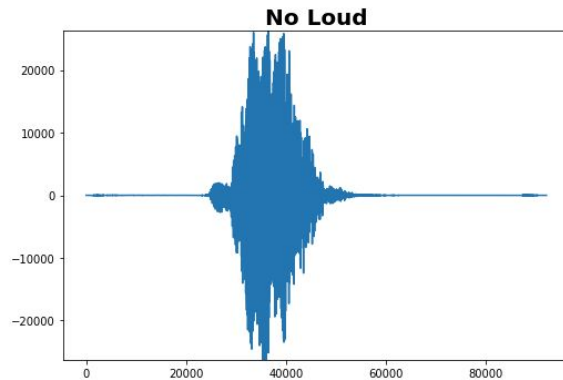
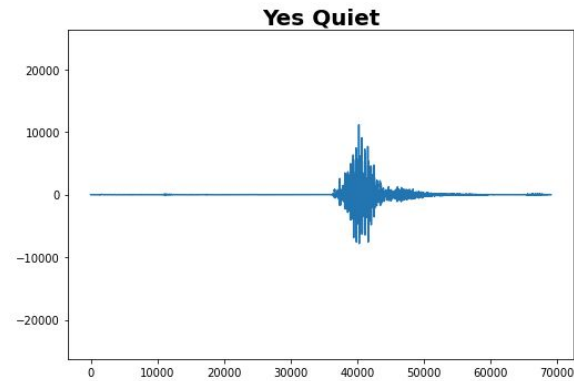
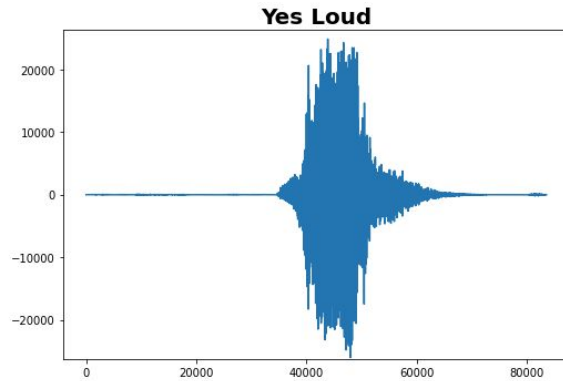


# Building a Spectrogram using FFTs

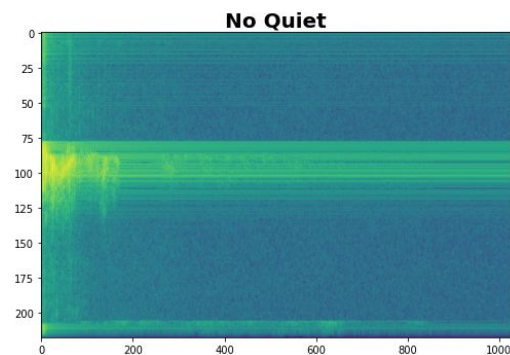
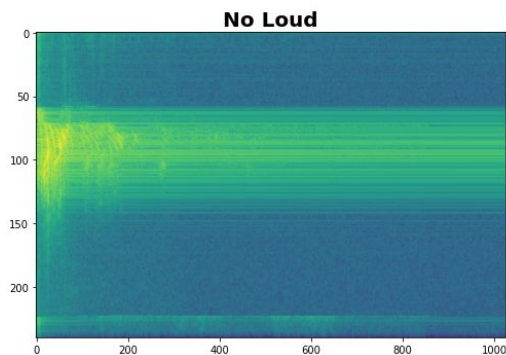
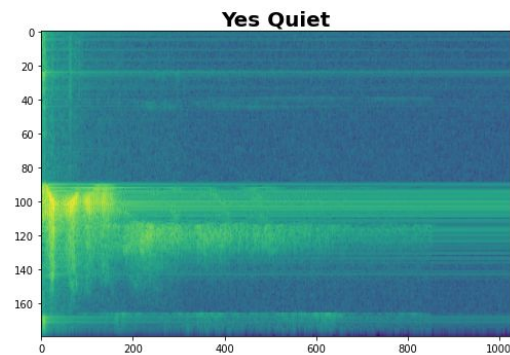
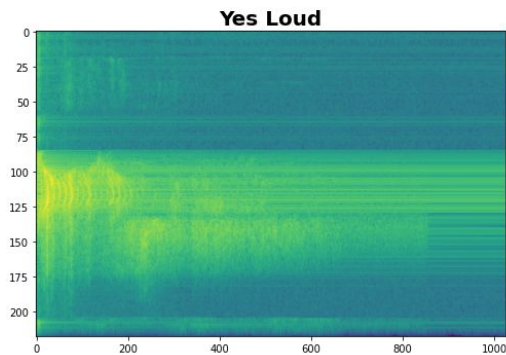


Essentially if you **stack up all the FFTs in a row** then you get the **Spectrogram** (time vs. frequency with color indicating intensity)

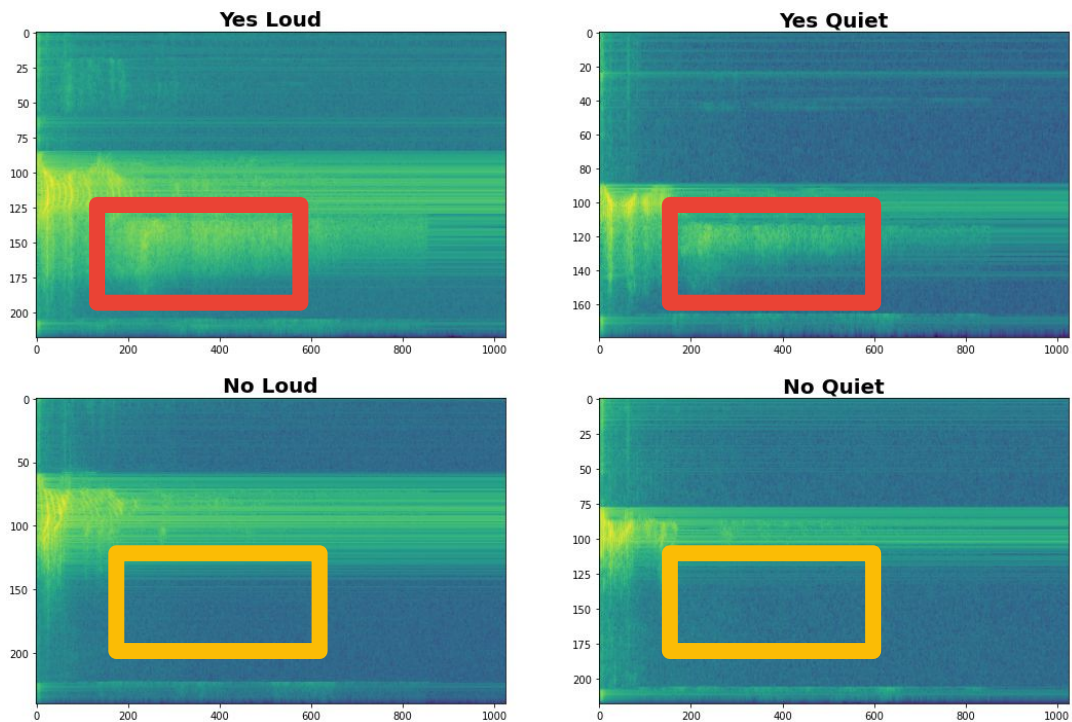
# Spectrograms help differentiate the data



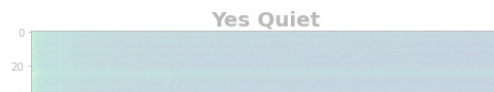
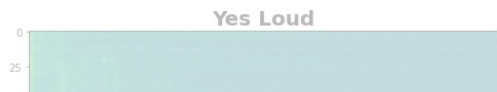
# Spectrograms help differentiate the data



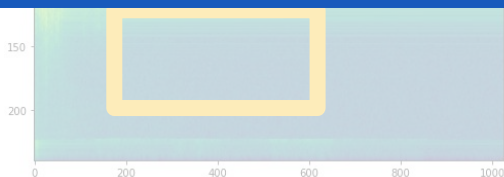
# Spectrograms help differentiate the data



# Data Preprocessing: Spectrograms

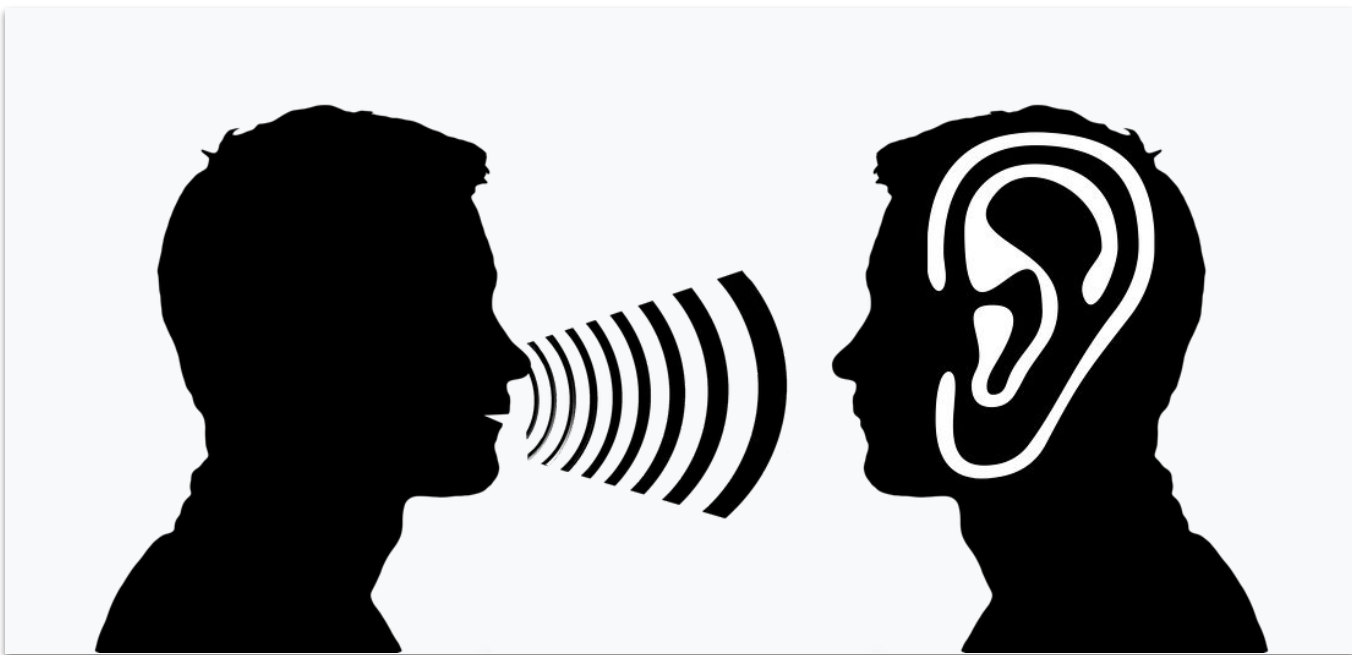


A spectrogram is also effectively an **image** that we can use as an input to a neural network!

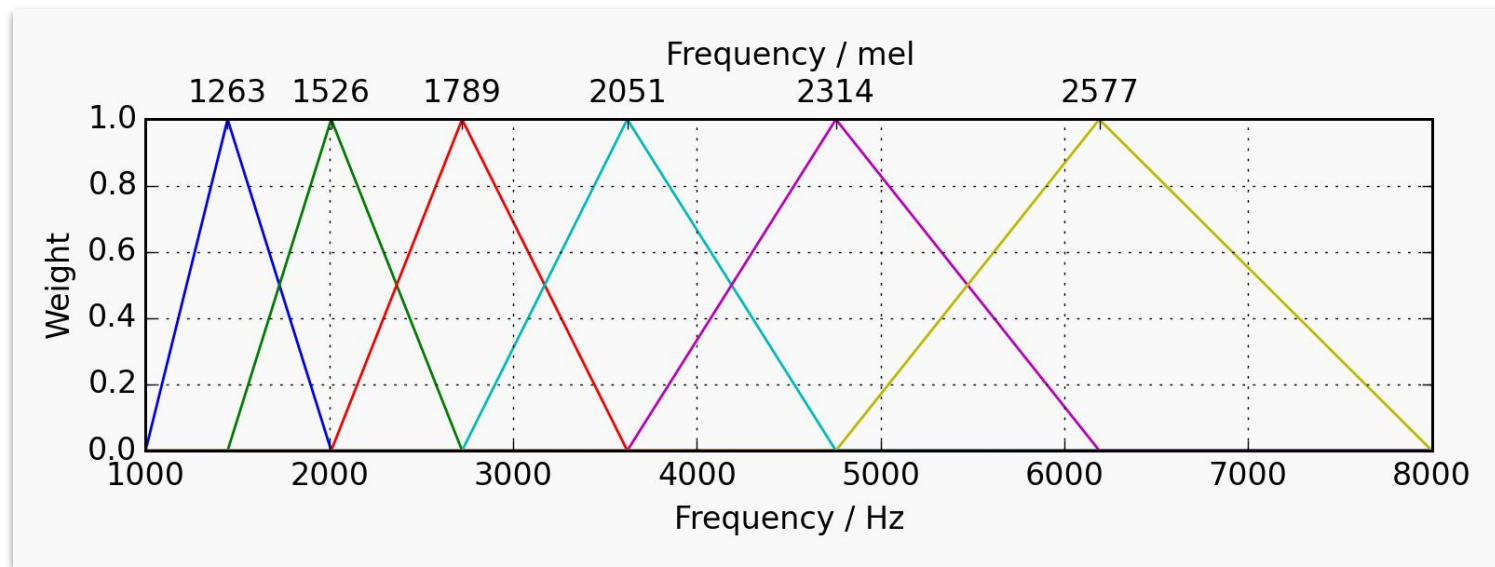


Can we do **better** than a spectrogram?

Can we take **domain knowledge** into account?

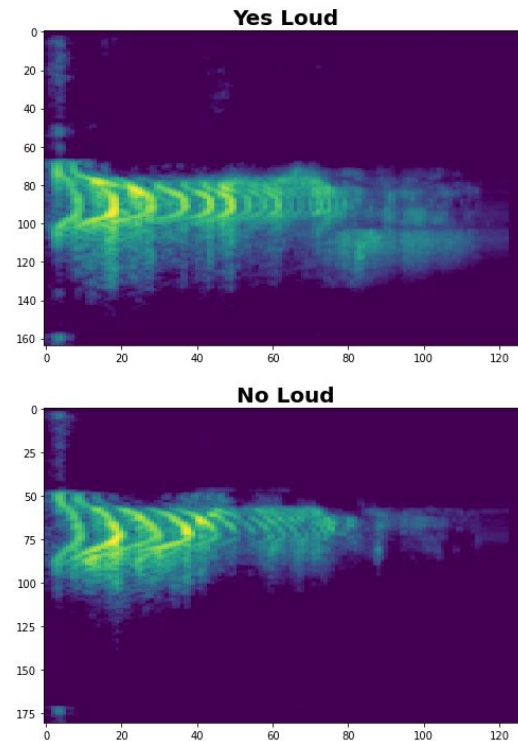
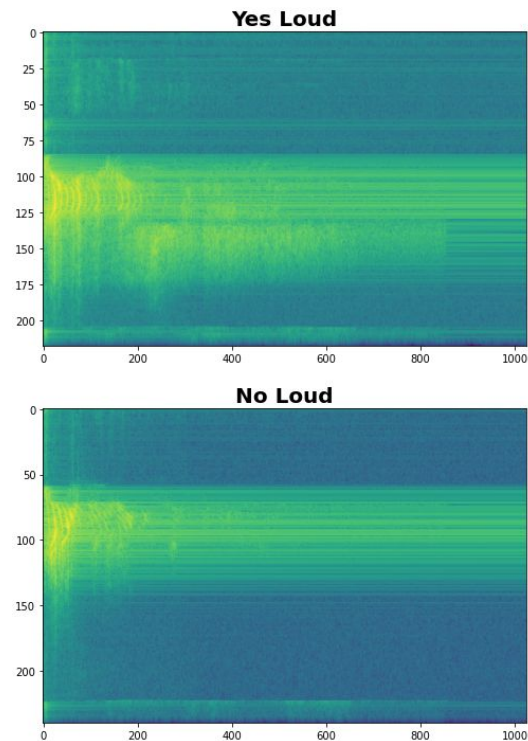


# Mel Filterbanks

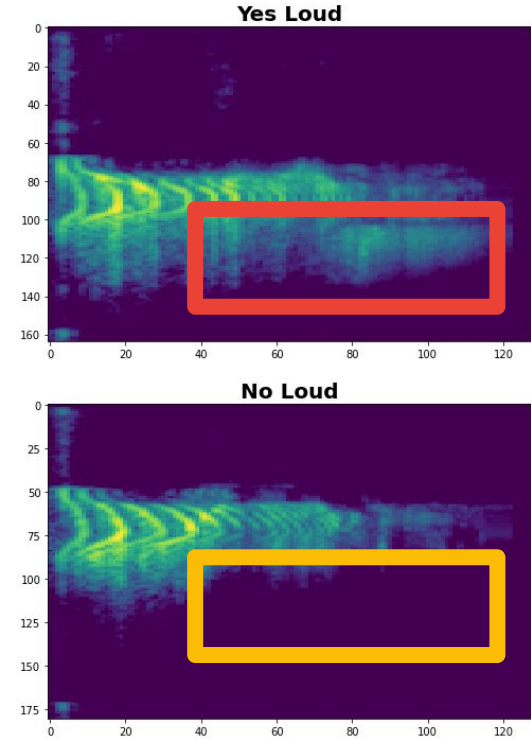
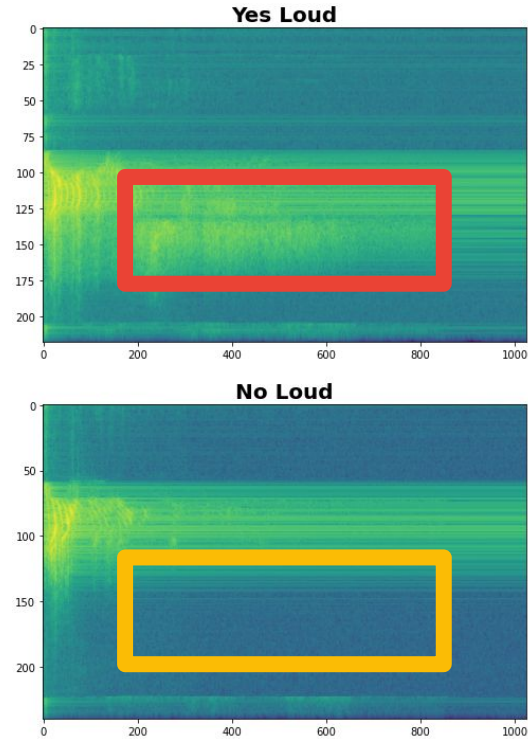




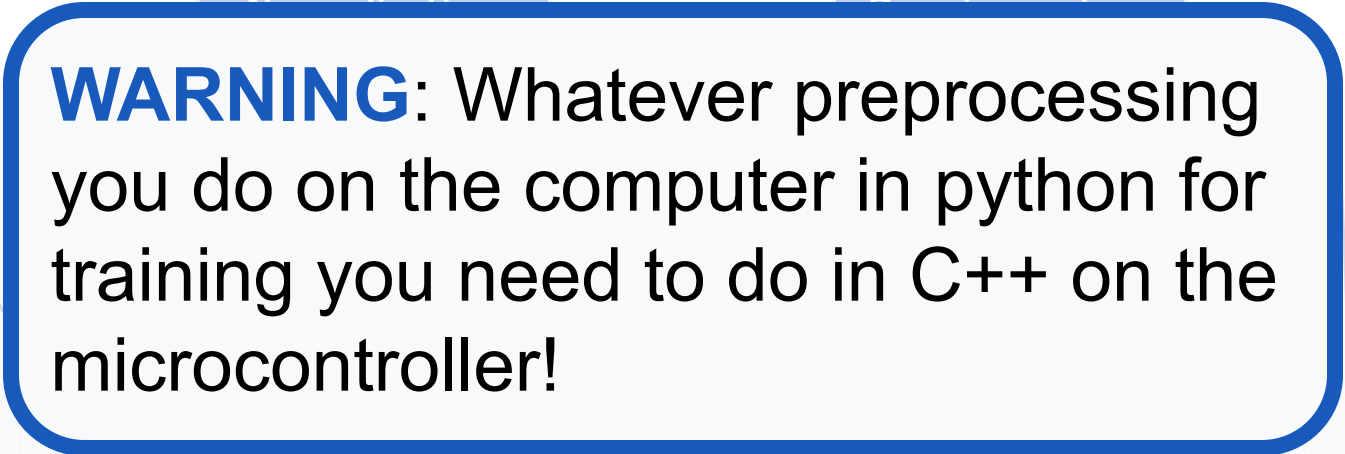
# Spectrograms v. MFCCs



# Spectrograms v. MFCCs



# Additional Feature Engineering



**WARNING:** Whatever preprocessing you do on the computer in python for training you need to do in C++ on the microcontroller!

# Today's Agenda

- Deep ML Background
- Hands-on Computer Vision: Thing Translator
- The Tiny Machine Learning Workflow
- Keyword Spotting (KWS) Data Collection
- **KWS Preprocessing and Training**

Preprocessing (for KWS)

## **Hands-on Preprocessing and Training with Edge Impulse**

- Deployment Challenges and Opportunities for Embedded ML
- Summary

 An impulse takes raw data, uses signal processing to extract features, and then uses a learning block to classify new data.

 Dashboard

 Devices

 Data acquisition

 Create impulse

 Retrain model

 Live classification

 Model testing


 Versioning

 Deployment



GETTING STARTED



 Documentation



 Forums


**Time series data** 


Axes  
audio


Window size   
  
1000 ms.


Window increase   
  
500 ms.


Frequency (Hz)   
16000 

Zero-pad data 



  
Add a processing block

  
Add a learning block

**Output features** 


Save Impulse

 An impulse takes raw data, uses signal processing to extract features, and then uses a learning block to classify new data.



- Dashboard
- Devices
- Data acquisition
- Impulse design
  - Create impulse
- EON Tuner
- Retrain model
- Live classification
- Model testing
- Versioning
- Deployment



---



- GETTING STARTED
  - Documentation
  - Forums


**Time series data** 


Axes  
audio


Window size   
 1000 ms.


Window increase   
 500 ms.


Frequency (Hz)   
 

Zero-pad data 



  
Add a processing block

  
Add a learning block

**Output features** 

Save Impulse

**⚡ Add a processing block**

Recommended based on your inputs

DESCRIPTION	AUTHOR	RECOMMENDED
<b>Audio (MFCC)</b> Extracts features from audio signals using Mel Frequency Cepstral Coefficients, great for human voice.	Edgelpulse Inc. ★	<b>Add</b>
<b>Audio (MFE)</b> Extracts a spectrogram from audio signals using Mel-filterbank energy features, great for non-voice audio.	Edgelpulse Inc. ★	Add
<b>Flatten</b> Flatten an axis into a single value, useful for slow-moving averages like temperature data, in combination with other blocks.	Edgelpulse Inc.	Add
<b>Image</b> Preprocess and normalize image data, and optionally reduce the color depth.	Edgelpulse Inc.	Add
<b>Spectral Analysis</b> Great for analyzing repetitive motion, such as data from accelerometers. Extracts the frequency and power characteristics of a signal over time.	Edgelpulse Inc.	Add
<b>Spectrogram</b> Extracts a spectrogram from audio or sensor data, great for non-voice audio or data with continuous frequencies.	Edgelpulse Inc.	Add

We'll keep things simple today and just add an MFCC but/and in future projects you can:

- **create your own blocks**
- **use multiple blocks**

<https://docs.edgeimpulse.com/docs/custom-blocks>

An impulse takes raw data, uses signal processing to extract features, and then uses a learning block to classify new data.

- Dashboard
- Devices
- Data acquisition
- Impulse design
  - Create impulse

- EON Tuner
- Retrain model
- Live classification
- Model testing
- Versioning
- Deployment

GETTING STARTED

- Documentation
- Forums

### Time series data

Axes  
audio

Window size ?  
1000 ms.

Window increase ?  
500 ms.

Frequency (Hz) ?  
16000

Zero-pad data ?

### Audio (MFCC)

Name  
MFCC

Input axes  
 audio

Add a learning block

### Output features

Save Impulse

Add a processing block



## Add a learning block ✕

Some learning blocks have been hidden based on the data in your project.

DESCRIPTION

AUTHOR

RECOMMENDED

### Classification (Keras)

Learns patterns from data, and can apply these to new data. Great for categorizing movement or recognizing audio.

EdgeImpulse Inc.



Add

### Regression (Keras)

Learns patterns from data, and can apply these to new data. Great for predicting numeric continuous values.

EdgeImpulse Inc.

Add

Cancel

Add a processing block

### Time series data



Axes

audio

Window size



1000 ms.

Window increase



500 ms.

Frequency (Hz)

16000



Zero-pad data



### Audio (MFCC)



Name

MFCC

Input axes



### Classification (Keras)



Name

NN Classifier

Input features



Output features

3 (no, unknown, yes)

### Output features



3 (no, unknown, yes)

Save Impulse



Add a processing block



Add a learning block

Successfully stored impulse. Configure the signal processing and learning blocks in the navigation bar.

Dashboard

Devices

Impulse design

- Create impulse
- MFCC
- NN Classifier

Retrain model

Live classification

Model testing

Versioning

Deployment

GETTING STARTED

### Time series data

Axes

audio

Window size

1000 ms.

Window increase

500 ms.

Frequency (Hz)

16000

Zero-pad data

### Audio (MFCC)

Name

MFCC

Input axes

audio

### Classification (Keras)

Name

NN Classifier

Input features

MFCC

Output features

3 (no, unknown, yes)

### Output features

3 (no, unknown, yes)

Save Impulse



#1 ▼ Click to set a description for this version

Parameters

Generate features

### Training set

Data in training set	1m 24s
Classes	3 (no, unknown, yes)
Window length	1000 ms.
Window increase	500 ms.
Training windows	126

Generate features

### Feature explorer ?

No features generated yet.

#1 Click to set a description for this version

Parameters

Generate features

### Training set

Data in training set	1m 24s
Classes	3 (no, unknown, yes)
Window length	1000 ms.
Window increase	500 ms.
Training windows	126

Generate features

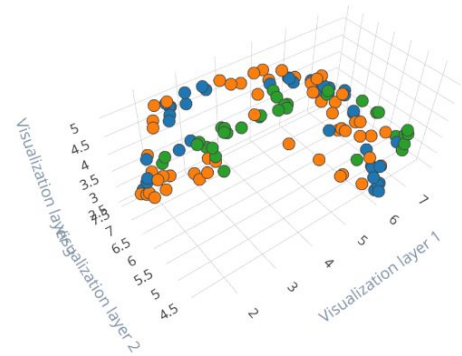
### Feature generation output

Sat Oct 16 17:25:45 2021 Construct embedding  
Still running...  
completed 0 / 500 epochs

### Feature explorer (126 samples)

X Axis: Visualization layer 1  
Y Axis: Visualization layer 2  
Z Axis: Visualization layer 3

- no
- unknown
- yes



- Dashboard
- Devices
- Data acquisition
- Impulse design
  - Create impulse
  - MFCC
  - NN Classifier
- EON Tuner
- Retrain model
- Live classification
- Model testing
- Versioning
- Deployment

GETTING STARTED

#1 [Click to set a description for this version](#)

Parameters

Generate features

Training set

Data in training set	1m 24s
Classes	3 (no, unknown, yes)
Window length	1000 ms.
Window increase	500 ms.
Training windows	126

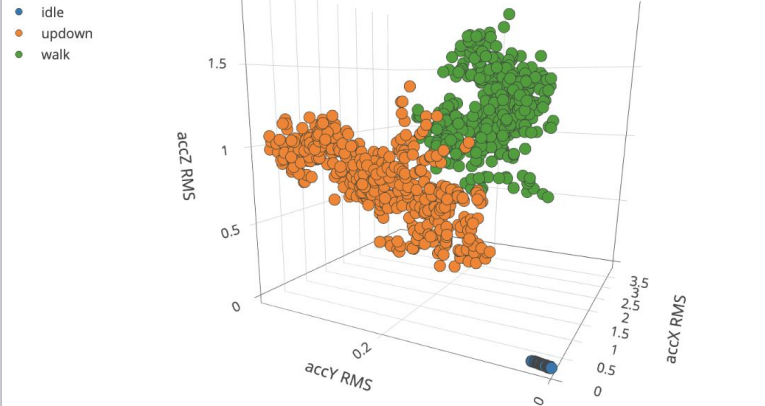
Generate features

Feature generation output

Sat Oct 16 17:25:45 2021 Construct embedding  
Still running...  
completed 0 / 500 epochs

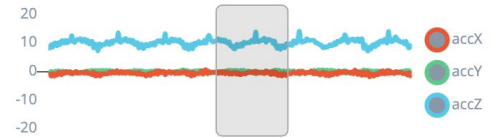
Feature explorer (1,506 samples)

X Axis: accX RMS  
Y Axis: accY RMS  
Z Axis: accZ RMS



updown.9.1cjh52qu  
Window: 4608 - 6608 ms.  
Label: updown

[View features](#)



- Dashboard
- Devices
- Data acquisition
- Impulse design
  - Create impulse
  - MFCC
  - NN Classifier
- EON Tuner
- Retrain model
- Live classification
- Model testing
- Versioning
- Deployment

GETTING STARTED

#1 Click to set a description for this version

Parameters

Generate features

### Training set

Data in training set	1m 24s
Classes	3 (no, unknown, yes)
Window length	1000 ms.
Window increase	500 ms.
Training windows	126

Generate features

### Feature generation output

Sat Oct 16 17:25:45 2021 Construct embedding  
Still running...  
completed 0 / 500 epochs

### Feature explorer (126 samples)

X Axis

Y Axis

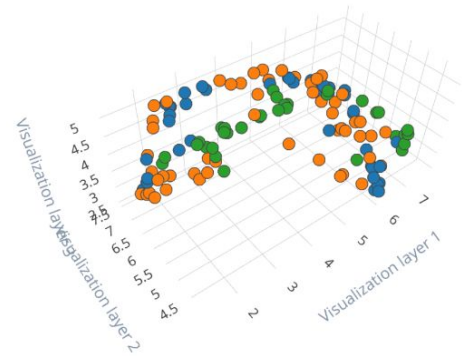
Z Axis

Visualization layer 1

Visualization layer 2

Visualization layer 3

- no
- unknown
- yes



Dashboard

Devices

Data acquisition

Impulse design

Create impulse

MFCC

NN Classifier

Retrain model

Live classification


Model testing

Versioning

Deployment

GETTING STARTED



 Dashboard Devices Data acquisition Impulse design Create impulse MFCC NN Classifier EON Tuner

## Neural Network settings

### Training settings

Number of training cycles Learning rate 

### Audio training options

Data augmentation  Switch to Keras (expert) mode Edit as iPython notebook

# Model Design with Edge Impulse

Pre-made neural network  
“blocks” that you can add!

### Neural Network settings

Training settings

Number of training cycles ⓘ

Learning rate ⓘ

Minimum confidence rating ⓘ

Neural network architecture

Input layer (637 features)

Reshape layer (13 columns)

1D conv / pool layer (30 neurons, 5 kernel size)

1D conv / pool layer (10 neurons, 5 kernel size)

Flatten layer

Add an extra layer

Output layer (5 features)

# Model Design with Edge Impulse

“Expert” mode to write your own TensorFlow code

## Neural network architecture

```
1 import tensorflow as tf
2 from tensorflow.keras.models import Sequential
3 from tensorflow.keras.layers import Dense, InputLayer,
  Dropout, Conv1D, Conv2D, Flatten, Reshape, MaxPooling1D,
  MaxPooling2D, BatchNormalization
4 from tensorflow.keras.optimizers import Adam
5 sys.path.append('./resources/libraries')
6 import ei_tensorflow.training
7
8 # model architecture
9 model = Sequential()
10 channels = 1
11 columns = 13
12 rows = int(input_length / (columns * channels))
13 model.add(Reshape((rows, columns, channels), input_shape
  =(input_length, )))
14 model.add(Conv2D(8, kernel_size=3, activation='relu',
  kernel_constraint=tf.keras.constraints.MaxNorm(1),
  padding='same'))
15 model.add(MaxPooling2D(pool_size=2, strides=2, padding
  ='same'))
16 model.add(Dropout(0.25))
17 model.add(Conv2D(16, kernel_size=3, activation='relu',
  kernel_constraint=tf.keras.constraints.MaxNorm(1),
  padding='same'))
18 model.add(MaxPooling2D(pool_size=2, strides=2, padding
  ='same'))
19 model.add(Dropout(0.25))
20 model.add(Flatten())
21 model.add(Dense(classes, activation='softmax', name='y_pred'
  ))
```

Start training

## Neural network architecture

Architecture presets ⓘ 1D Convolutional (Default) 2D Convolutional

Input layer (650 features)

Reshape layer (13 columns)

1D conv / pool layer (8 neurons, 3 kernel size, 1 layer)

Dropout (rate 0.25)

1D conv / pool layer (16 neurons, 3 kernel size, 1 layer)

Dropout (rate 0.25)

Flatten layer

Add an extra layer

Output layer (3 features)

Start training

## Neural network architecture

```
1 import tensorflow as tf
2 from tensorflow.keras.models import Sequential
3 from tensorflow.keras.layers import Dense, InputLayer, Dropout, Conv1D, Conv2D,
  Flatten, Reshape, MaxPooling1D, MaxPooling2D, BatchNormalization,
  TimeDistributed
4 from tensorflow.keras.optimizers import Adam
5
6 # model architecture
7
8 model.add(Reshape((int(input_length / 13), 13), input_shape=(input_length, )))
9 model.add(Conv1D(8, kernel_size=3, activation='relu', padding='same'))
10 model.add(MaxPooling1D(pool_size=2, strides=2, padding='same'))
11
12 model.add(Conv1D(16, kernel_size=3, activation='relu', padding='same'))
13 model.add(MaxPooling1D(pool_size=2, strides=2, padding='same'))
14 model.add(Dropout(0.25))
15 model.add(Flatten())
16 model.add(Dense(classes, activation='softmax', name='y_pred'))
17
18 # this controls the learning rate
19 opt = Adam(lr=0.005, beta_1=0.9, beta_2=0.999)
20 # this controls the batch size, or you can manipulate the tf.data.Dataset objects
  yourself
21 BATCH_SIZE = 32
22 train_dataset = train_dataset.batch(BATCH_SIZE, drop_remainder=False)
23 validation_dataset = validation_dataset.batch(BATCH_SIZE, drop_remainder=False)
24 callbacks.append(BatchLoggerCallback(BATCH_SIZE, train_sample_count))
25
26 # train the neural network
27 model.compile(loss='categorical_crossentropy', optimizer=opt, metrics=['accuracy'])
28 model.fit(train_dataset, epochs=100, validation_data=validation_dataset, verbose=2,
  callbacks=callbacks)
```

## Neural network architecture

Architecture presets ⓘ 1D Convolutional (Default) 2D Convolutional

Input layer (650 features)

Reshape layer (13 columns)

1D conv / pool layer (8 neurons, 3 kernel size, 1 layer)

Dropout (rate 0.25)

1D conv / pool layer (16 neurons, 3 kernel size, 1 layer)

Dropout (rate 0.25)

Flatten layer

Add an extra layer

Output layer (3 features)

Start training

## Neural network architecture

```
1 import tensorflow as tf
2 from tensorflow.keras.models import Sequential
3 from tensorflow.keras.layers import Dense, InputLayer, Dropout, Conv1D, Conv2D,
  Flatten, Reshape, MaxPooling1D, MaxPooling2D, BatchNormalization,
  TimeDistributed
4 from tensorflow.keras.optimizers import Adam
5
6 # model architecture
7 model = Sequential()
8 model.add(Reshape((int(input_length / 13), 13), input_shape=(input_length, )))
9 model.add(Conv1D(8, kernel_size=3, activation='relu', padding='same'))
10 model.add(MaxPooling1D(pool_size=2, strides=2, padding='same'))
11 model.add(Dropout(0.25))
12 model.add(Conv1D(16, kernel_size=3, activation='relu', padding='same'))
13 model.add(MaxPooling1D(pool_size=2, strides=2, padding='same'))
14 model.add(Dropout(0.25))
15 model.add(Flatten())
16 model.add(Dense(classes, activation='softmax', name='y_pred'))
17
18 # this controls the learning rate
19 opt = Adam(lr=0.005, beta_1=0.9, beta_2=0.999)
20 # this controls the batch size, or you can manipulate the tf.data.Dataset objects
  yourself
21 BATCH_SIZE = 32
22 train_dataset = train_dataset.batch(BATCH_SIZE, drop_remainder=False)
23 validation_dataset = validation_dataset.batch(BATCH_SIZE, drop_remainder=False)
```

For now just stick with the defaults but/and you can easily design **any model** you want and use **any optimizer** you want using **TensorFlow!**

Input layer (650 features)

1D co

1D co

**WARNING:** if you want to deploy to a microcontroller make sure you only use Ops supported by TensorFlow Lite Micro!  
[https://github.com/tensorflow/tflite-micro/blob/main/tensorflow/lite/micro/all\\_ops\\_resolver.cc#L22](https://github.com/tensorflow/tflite-micro/blob/main/tensorflow/lite/micro/all_ops_resolver.cc#L22)

Output layer (3 features)

Start training

```
1 import tensorflow as tf
2 from tensorflow.keras.models import Sequential
3 from tensorflow.keras.layers import Dense, InputLayer, Dropout, Conv1D, Conv2D,
  Flatten, Reshape, MaxPooling1D, MaxPooling2D, BatchNormalization,
  TimeDistributed
4 from tensorflow.keras.optimizers import Adam
5
6 # model architecture
```

```
...shape=(input_length, ))
padding='same'))
'same'))
padding='same'))
'same'))
pred'))
the tf.data.Dataset objects
remainder=False)
E._drop_remainder=False)
...y')
e=2,
```

easily design **any model** you want and use **any optimizer** you want using **TensorFlow!**

## Neural network architecture

Architecture presets ⓘ 1D Convolutional (Default) 2D Convolutional

Input layer (650 features)

Reshape layer (13 columns)

1D conv / pool layer (8 neurons, 3 kernel size, 1 layer)

Dropout (rate 0.25)

1D conv / pool layer (16 neurons, 3 kernel size, 1 layer)

Dropout (rate 0.25)

Flatten layer

Add an extra layer

Output layer (3 features)

Start training

## Neural network architecture

```
1 import tensorflow as tf
2 from tensorflow.keras.models import Sequential
3 from tensorflow.keras.layers import Dense, InputLayer, Dropout, Conv1D, Conv2D,
  Flatten, Reshape, MaxPooling1D, MaxPooling2D, BatchNormalization,
  TimeDistributed
4 from tensorflow.keras.optimizers import Adam
5
6 # model architecture
7 model = Sequential()
8 model.add(Reshape((int(input_length / 13), 13), input_shape=(input_length, )))
9 model.add(Conv1D(8, kernel_size=3, activation='relu', padding='same'))
10 model.add(MaxPooling1D(pool_size=2, strides=2, padding='same'))
11 model.add(Dropout(0.25))
12 model.add(Conv1D(16, kernel_size=3, activation='relu', padding='same'))
13 model.add(MaxPooling1D(pool_size=2, strides=2, padding='same'))
14 model.add(Dropout(0.25))
15 model.add(Flatten())
16 model.add(Dense(classes, activation='softmax', name='y_pred'))
17
18 # this controls the learning rate
19 opt = Adam(lr=0.005, beta_1=0.9, beta_2=0.999)
20 # this controls the batch size, or you can manipulate the tf.data.Dataset objects
  yourself
21 BATCH_SIZE = 32
22 train_dataset = train_dataset.batch(BATCH_SIZE, drop_remainder=False)
23 validation_dataset = validation_dataset.batch(BATCH_SIZE, drop_remainder=False)
```

For now just stick with the defaults but/and you can easily design **any model** you want and use **any optimizer** you want using **TensorFlow!**

## Training output

```
Epoch 95/100  
4/4 - 0s - loss: 0.1044 - accuracy: 0.9500 - val_loss: 0.2934 - val_accuracy: 0.9231  
Epoch 96/100  
4/4 - 0s - loss: 0.0256 - accuracy: 1.0000 - val_loss: 0.3830 - val_accuracy: 0.8846  
Epoch 97/100  
4/4 - 0s - loss: 0.0523 - accuracy: 0.9800 - val_loss: 0.4366 - val_accuracy: 0.8462  
Epoch 98/100  
4/4 - 0s - loss: 0.0451 - accuracy: 0.9800 - val_loss: 0.4265 - val_accuracy: 0.8846  
Epoch 99/100  
4/4 - 0s - loss: 0.0514 - accuracy: 0.9900 - val_loss: 0.3926 - val_accuracy: 0.8846  
Epoch 100/100  
4/4 - 0s - loss: 0.0348 - accuracy: 0.9900 - val_loss: 0.3571 - val_accuracy: 0.9231  
Finished training
```



Training Set



Validation Set



# Final Accuracy



Model

Model version: ?

Quantized (int8) ▾

Last training performance (validation set)

**%** ACCURACY  
92.3%

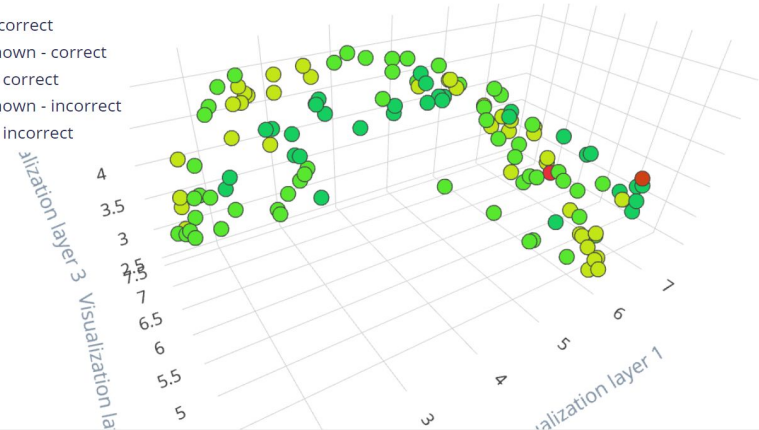
**LOSS**  
0.27

Confusion matrix (validation set)

	NO	UNKNOWN	YES
NO	100%	0%	0%
UNKNOWN	9.1%	90.9%	0%
YES	0%	11.1%	88.9%
F1 SCORE	0.92	0.91	0.94

Feature explorer (full training set) ?

- no - correct
- unknown - correct
- yes - correct
- unknown - incorrect
- yes - incorrect



Final Accuracy

Accuracy Breakdown

Model

Model version: ?

Quantized (int8)

Last training performance (validation set)

ACCURACY  
92.3%

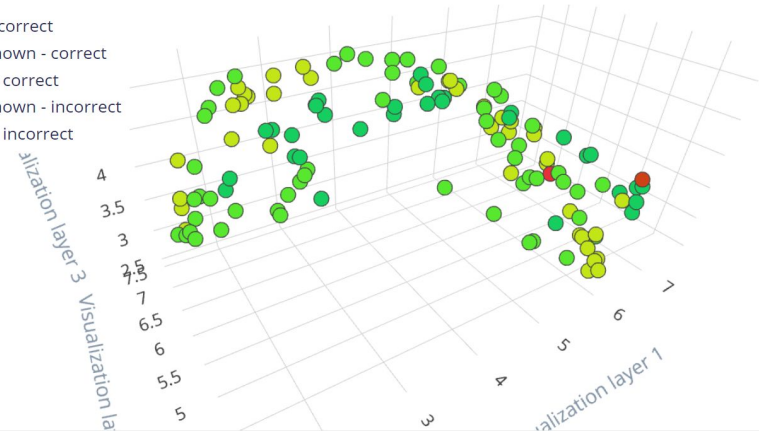
LOSS  
0.27

Confusion matrix (validation set)

	NO	UNKNOWN	YES
NO	100%	0%	0%
UNKNOWN	9.1%	90.9%	0%
YES	0%	11.1%	88.9%
F1 SCORE	0.92	0.91	0.94

Feature explorer (full training set) ?

- no - correct
- unknown - correct
- yes - correct
- unknown - incorrect
- yes - incorrect



# Confusion Matrix

	Actual Output = Yes	Actual Output = No
Predicted Output = Yes	<b># of True Positive</b>	<b># of False Positive <i>Type 1 Error</i></b>
Predicted Output = No	<b># of False Negative <i>Type 2 Error</i></b>	<b># of True Negative</b>

Final Accuracy

Accuracy Breakdown

Model

Model version: ?

Quantized (int8)

Last training performance (validation set)

ACCUACY  
92.3%

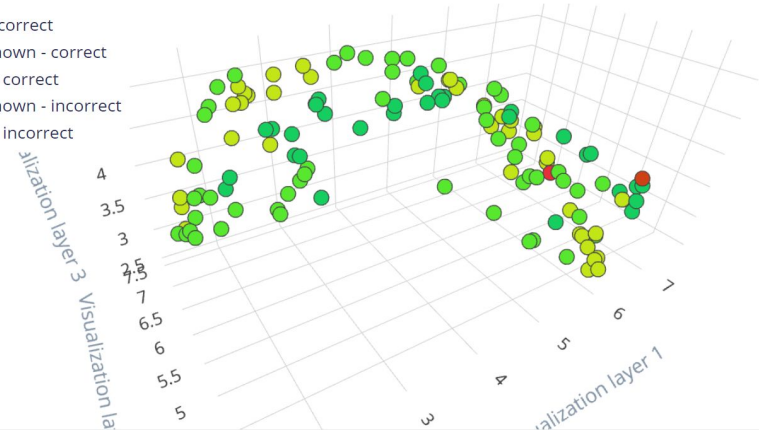
LOSS  
0.27

Confusion matrix (validation set)

	NO	UNKNOWN	YES
NO	100%	0%	0%
UNKNOWN	9.1%	90.9%	0%
YES	0%	11.1%	88.9%
F1 SCORE	0.92	0.91	0.94

Feature explorer (full training set) ?

- no - correct
- unknown - correct
- yes - correct
- unknown - incorrect
- yes - incorrect



# Final Accuracy



# Accuracy Breakdown



Model

Model version: ?

Quantized (int8) ▾

Last training performance (validation set)

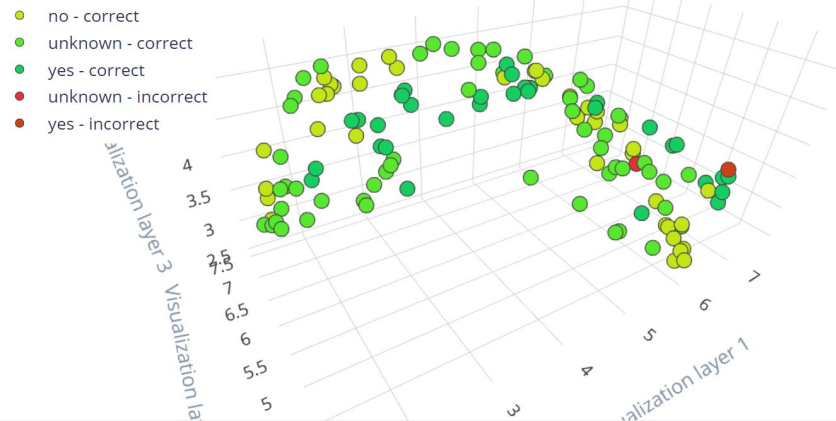
**%** ACCURACY  
92.3%

**LOSS**  
0.27

Confusion matrix (validation set)

	NO	UNKNOWN	YES
NO	100%	0%	0%
UNKNOWN	9.1%	90.9%	0%
YES	0%	11.1%	88.9%
F1 SCORE	0.92	0.91	0.94

Feature explorer (full training set) ?

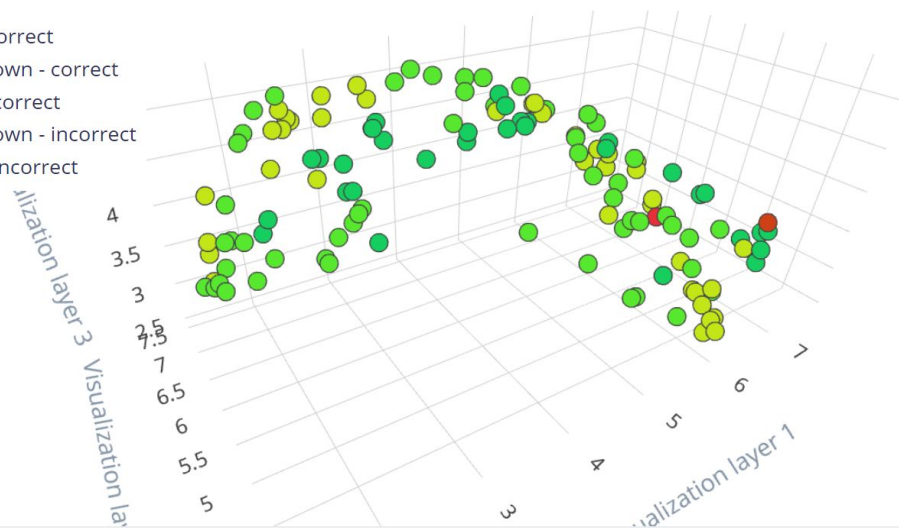


Final Accuracy

Accuracy Breakdown

### Feature explorer (full training set) ?

- no - correct
- unknown - correct
- yes - correct
- unknown - incorrect
- yes - incorrect



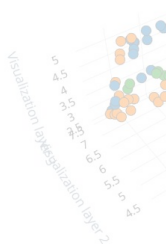
Quantized (int8)

	YES
	0%
	0%
	88.9%
	0.94

Feature explorer (126 samples)

X Axis: Visualization layer 1  
Y Axis: Visualization layer 2

- no
- unknown
- yes



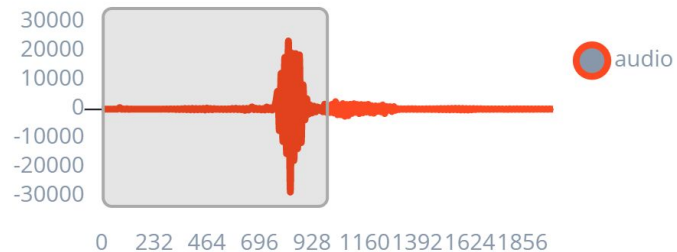
**yes.2hvfiruf**

Label: yes

Predicted: unknown

[View sample](#)

[View features](#)

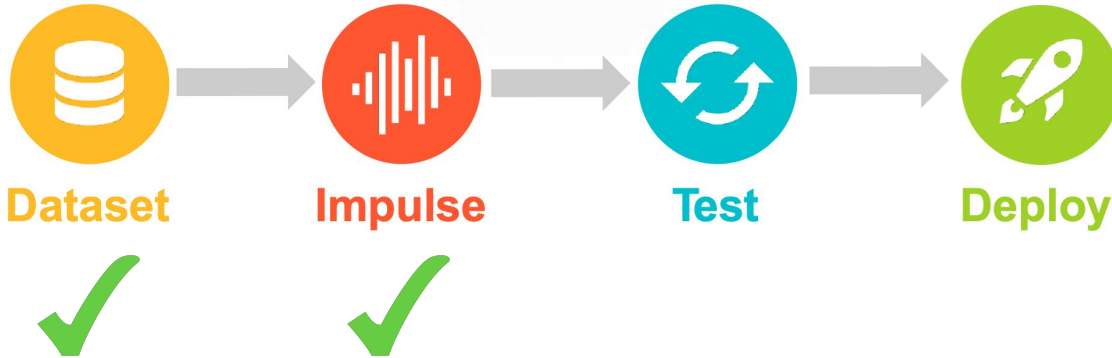
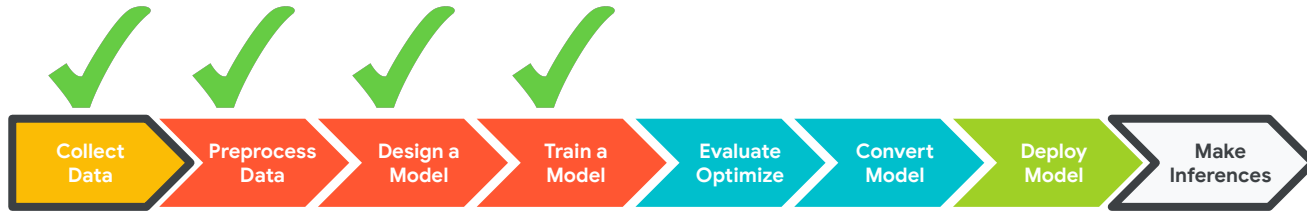


0 232 464 696 928 1160 1392 1624 1856

▶ 0:00 / 0:01



# Edge Impulse Project Dashboard



- ✓ Dashboard
- ✓ Devices
- ✓ Data acquisition
- ✓ Impulse design
- ✓ Create impulse
- ✓ MFCC
- ✓ NN Classifier
- EON Tuner
- Retrain model
- Live classification
- Model testing
- Versioning
- Deployment



# Today's Agenda

- Deep ML Background
- Hands-on Computer Vision: Thing Translator
- The Tiny Machine Learning Workflow
- Keyword Spotting (KWS) Data Collection
- KWS Preprocessing and Training
- **Deployment Challenges and Opportunities for Embedded ML**
- Summary



Even Lower power  
Even Lower bandwidth  
Even Lower cost



# Compute



# Memory





# Storage

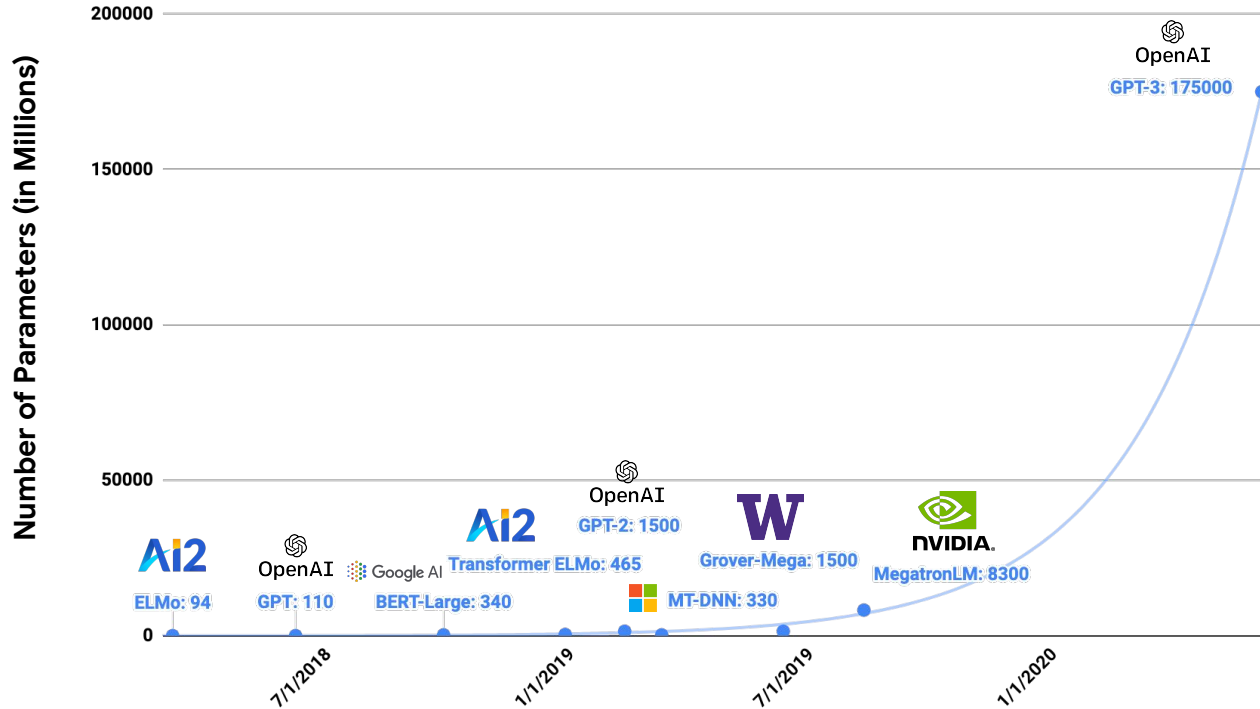


Microcontrollers have **slower** compute and **very little** memory and storage

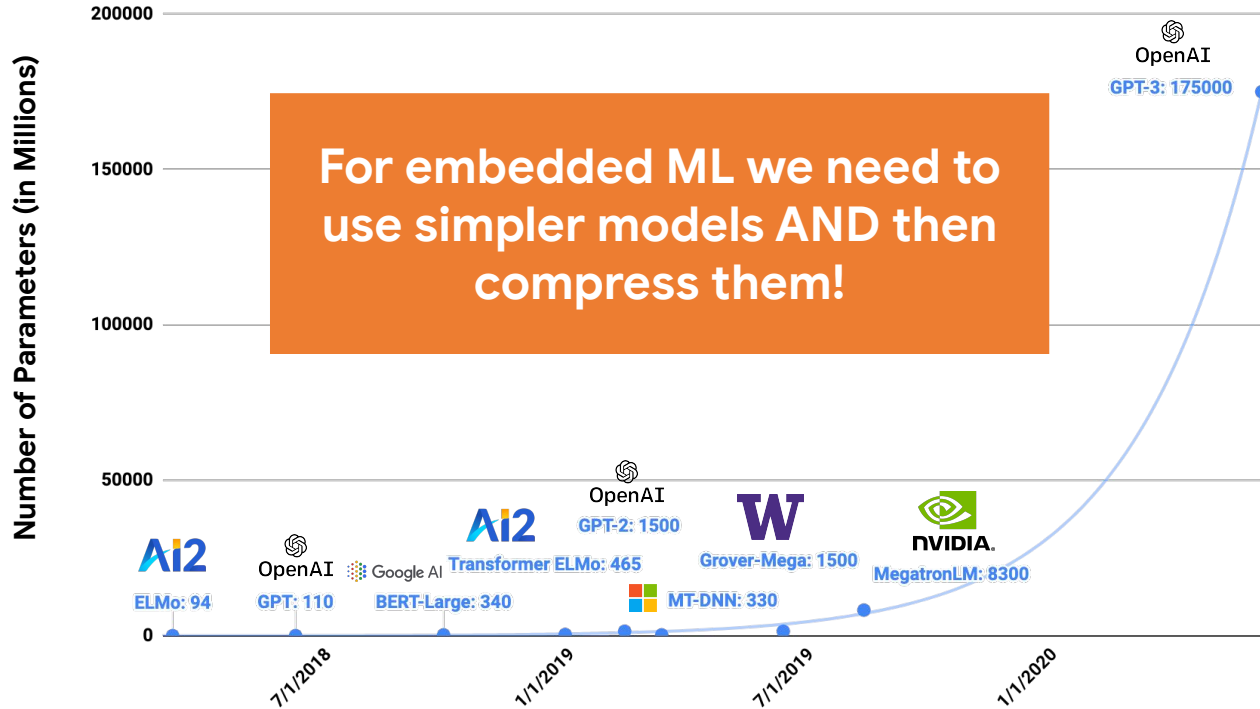
# Orders of Magnitude Difference

	Computer 		Microcontroller 
Compute	1GHz–4GHz	~10X	1MHz–400MHz
Memory	512MB–64GB	~10,000X	2KB–512KB
Storage	64GB–4TB	~100,000X	32KB–2MB

# ML Model Size Growth



# ML Model Size Growth



# Quantization

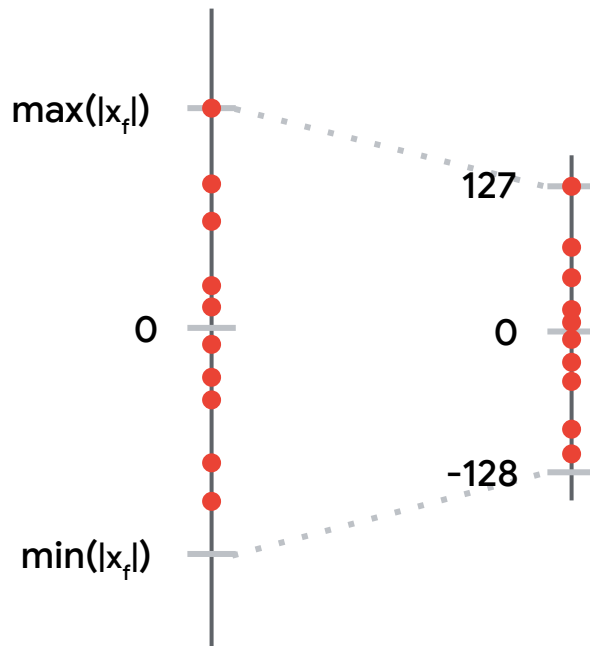
**Reduces the precision** of numbers used in a model which results in:

- **smaller model size**
- **faster computation**

# Reducing the Precision

*float32*

*int8*



**4 bytes per  
model  
parameter**

**1 byte per  
model  
parameter**

# Tradeoff

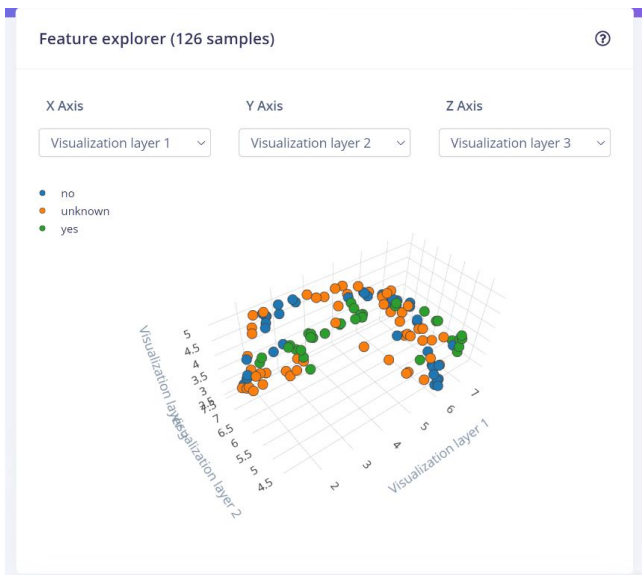
	Floating-point Baseline	After Quantization	Accuracy Drop
MobileNet v1 1.0 224	71.03%	69.57%	▼1.46%
MobileNet v2 1.0 224	70.77%	70.20%	▼0.57%
Resnet v1 50	76.30%	75.95%	▼0.35%



Final Accuracy



Accuracy Breakdown



Model

Model version: ?

Quantized (int8) ▾

Last training performance (validation set)

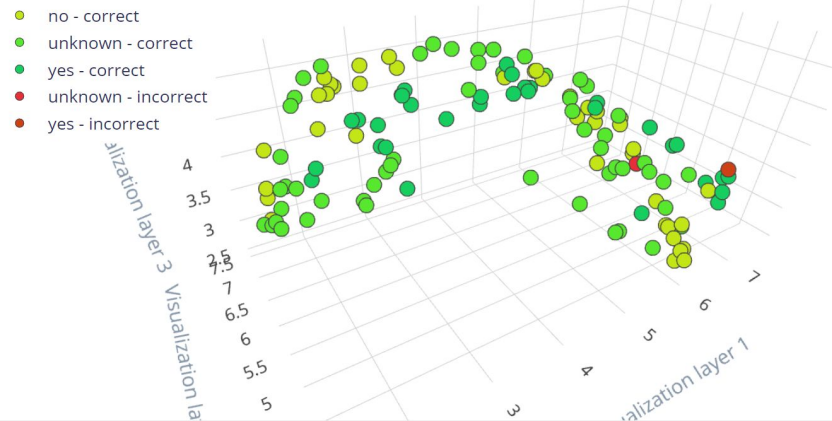
ACCURACY  
92.3%

LOSS  
0.27

Confusion matrix (validation set)

	NO	UNKNOWN	YES
NO	100%	0%	0%
UNKNOWN	9.1%	90.9%	0%
YES	0%	11.1%	88.9%
F1 SCORE	0.92	0.91	0.94

Feature explorer (full training set) ?



Final Accuracy



Model

Model version: ?

Quantized (int8)

Last training performance (validation set)



ACCURACY

92.3%



LOSS

0.27

Model

Model version: ?

Quantized (int8)

Last training performance (validation set)



ACCURACY

92.3%



LOSS

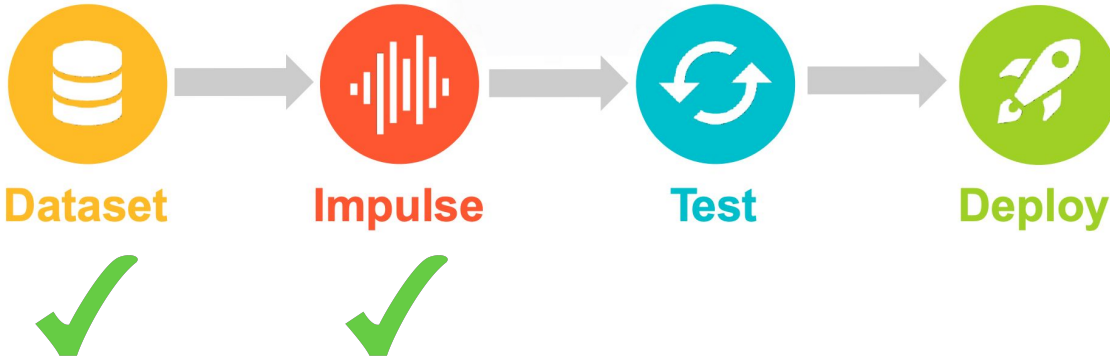
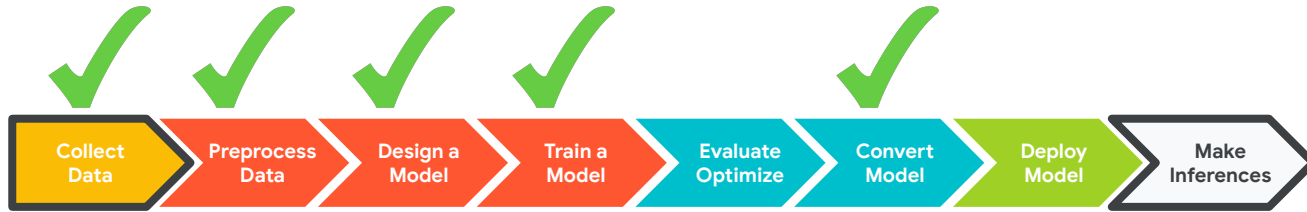
0.27

Quantized (int8)

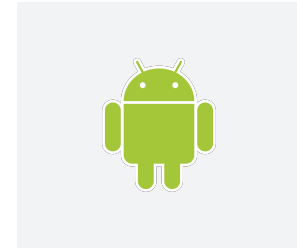
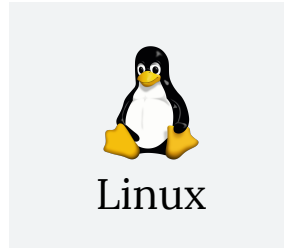
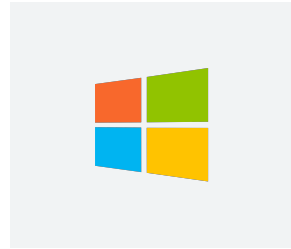
Unoptimized (float32)

Confusion matrix (validation set)

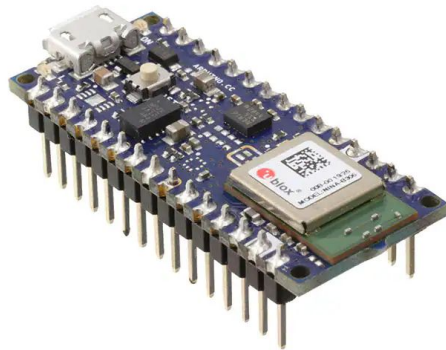
# Edge Impulse Project Dashboard



- Dashboard
- Devices
- Data acquisition
- Impulse design
- Create impulse
- MFCC
- NN Classifier
- EON Tuner
- Retrain model
- Live classification
- Model testing
- Versioning
- Deployment



Most **operating systems** come with many **libraries** and **applications** that make it **easy and portable to write code once** and then compile it in an optimized form for most computers (or smartphones)



Microcontrollers often require **custom code and compilation toolchains** to run optimally

# Edge Impulse simplifies deployment






Pick your destination / device  
and **deploy** the same model to  
any of them thanks to  
**collaboration with hardware  
vendors** and the use of  
**TensorFlow Lite Micro!**

## Deploy your impulse

You can deploy your impulse to any device. This makes the model run without an internet connection, minimizes latency, and runs with minimal power consumption. [Read more.](#)

## Create library

Turn your impulse into optimized source code that you can run on any device.

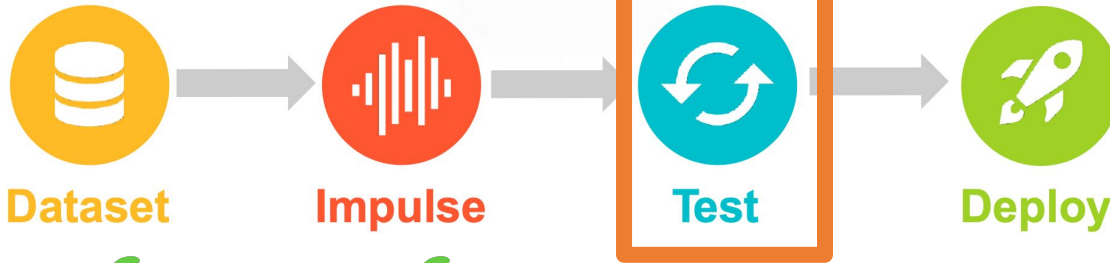
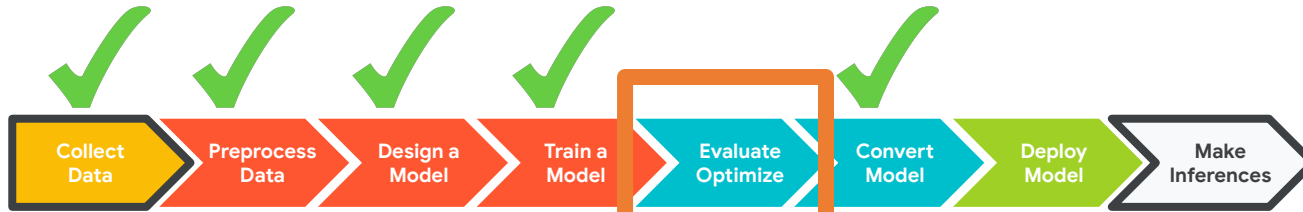
 C++ library	 Arduino library	 Cube.MX CMSIS-PAK
 WebAssembly	 TensorRT library	

## Build firmware

Or get a ready-to-go binary for your development board that includes your impulse.

 ST IoT Discovery Kit	 Arduino Nano 33 BLE Sense	 Eta Compute ECM3532 AI Sensor <span>END OF LIFE</span>
 <span>END OF LIFE</span>		

# Edge Impulse Project Dashboard

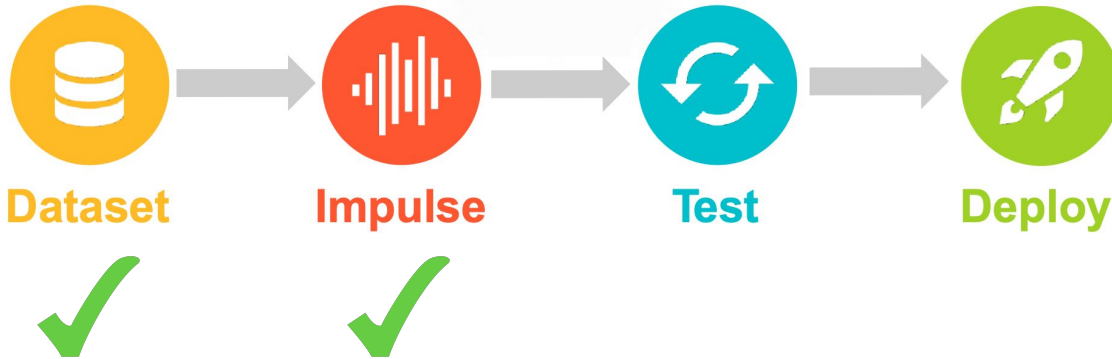
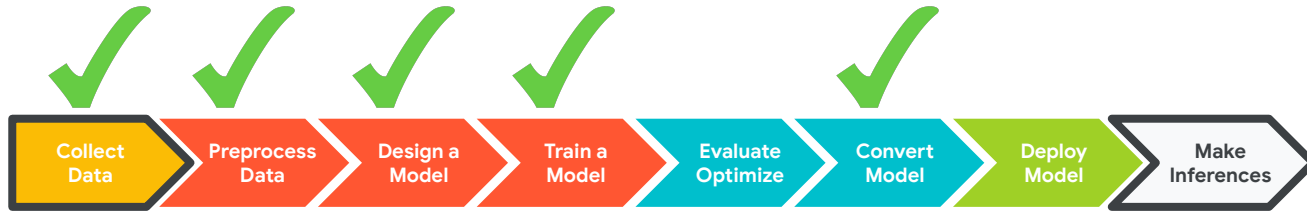


<https://www.edgeimpulse.com/blog/introducing-the-eon-tuner-edge-impulses-new-auto-ml-tool-for-embedded-machine-learning>

- ✓ Dashboard
- ✓ Devices
- ✓ Data acquisition
- ✓ Impulse design
- ✓ Create impulse
- ✓ MFCC
- ✓ NN Classifier

- EON Tuner
- Retrain model
- Live classification
- Model testing
- Versioning

# Edge Impulse Project Dashboard



- Dashboard
- Devices
- Data acquisition
- Impulse design
- Create impulse
- MFCC
- NN Classifier
- EON Tuner
- Retrain model
- Live classification
- Model testing
- Versioning

- Dashboard
- Devices
- Data acquisition
- Impulse design
  - Create impulse
  - MFCC
  - NN Classifier
- EON Tuner
- Retrain model
- Live classification
- Model testing
- Versioning
- Deployment**

## GETTING STARTED

- Documentation
- Forums

## Deploy your impulse

You can deploy your impulse to any device. This makes the model run without an internet connection, minimizes latency, and runs with minimal power consumption. [Read more.](#)

### Create library

Turn your impulse into optimized source code that you can run on any device.

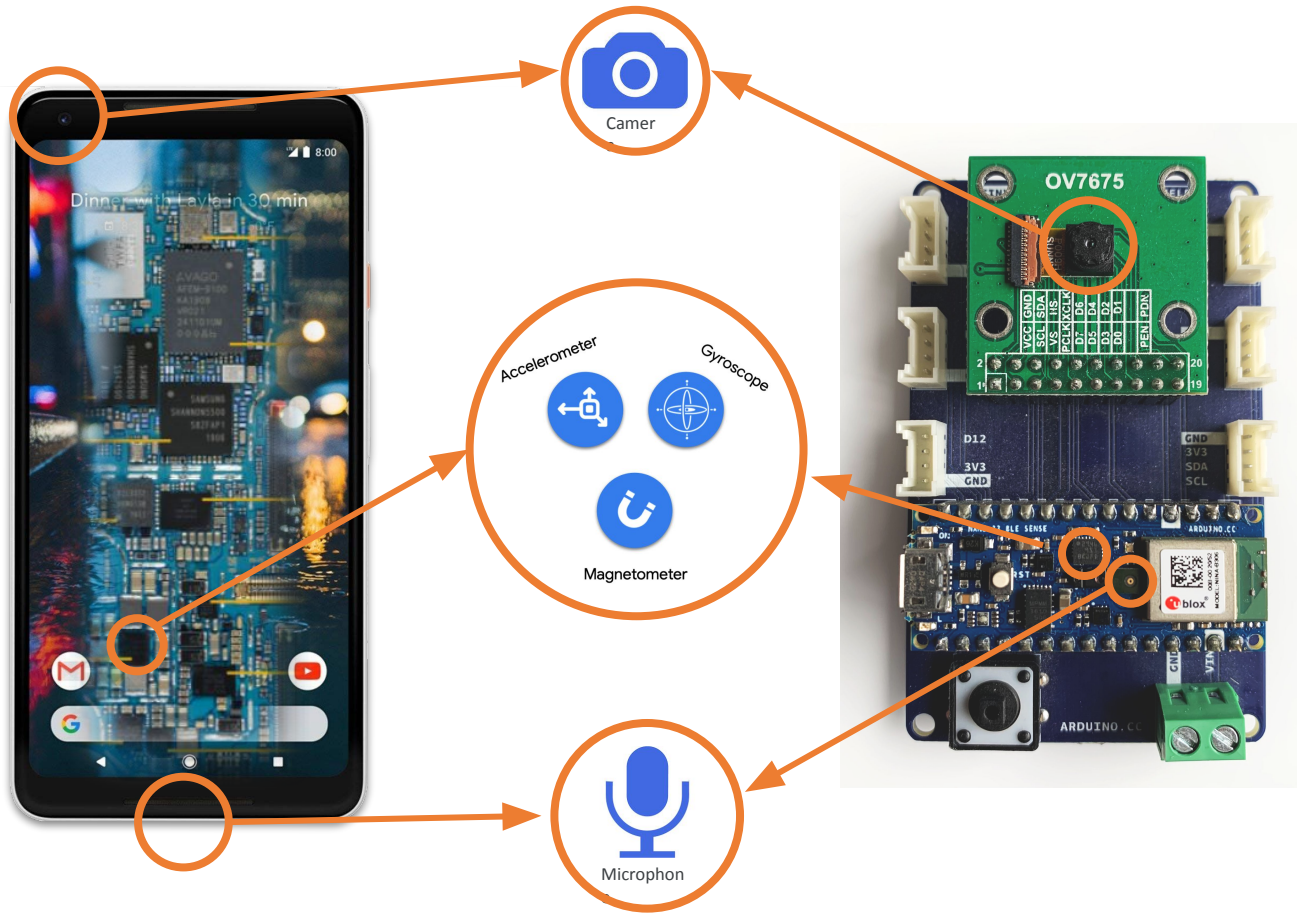
 C++ library	 Arduino library	 Cube.MX CMSIS-PACK
 WebAssembly	 TensorRT library	

### Build firmware

Or get a ready-to-go binary for your development board that includes your impulse.

 ST IoT Discovery Kit	 Arduino Nano 33 BLE Sense	<span style="background-color: #f08080; padding: 2px;">END OF LIFE</span> Eta Compute ECM3532 AI Sensor
<span style="background-color: #f08080; padding: 2px;">END OF LIFE</span>		





Devices

Impulse design

Create impulse

MFCC

NN Classifier

## Your devices

These are devices that are connected to the Edge Impulse

NAME

phone\_kunh8zjd

computer\_kq77e063

## Collect data

You can collect data from development boards, from your own devices, or by uploading an existing dataset.



## Connect a fully supported development board

Get started with real hardware from a wide range of silicon vendors - fully supported by Edge Impulse.

Browse dev boards



## Use your mobile phone

Use your mobile phone to capture movement, audio or images, and even run your trained model locally. No app required.

Show QR code

+ Connect a new device

REMOTE ...

LAST SEEN

camera, ...



Today, 16:24:48



camera



Jun 21 2021, 18:41:37




Devices

- Impulse design
  - Create impulse
  - MFCC
  - NN Classifier

## Your devices

These are devices that are connected to the Edge Impulse studio.

NAME

 phone_kunh8zjd
 computer_kq77e063

## Collect data

You can collect data from development boards, from your own devices, or by uploading an existing dataset.



## Connect a fully supported development board

Get started with real hardware from a wide range of silicon vendors - fully supported by Edge Impulse.

[Browse dev boards](#)

## Use your mobile phone

Use your mobile phone to capture movement, audio or images, and even run your trained model locally. No app required.

[Show QR code](#)[+ Connect a new device](#)

## Connected as phone\_kunh8zjd

You can collect data from this device from the **Data acquisition** page in the Edge Impulse studio.

[Collecting images?](#)[Collecting audio?](#)[Collecting motion?](#)[Switch to classification mode](#)[</> This client is open source.](#)

Devices

Impulse design

Create impulse

MFCC

NN Classifier

## Your devices

These are devices that are connected to the Edge Impulse studio.

NAME

phone\_kunh8zjd

computer\_kq77e063

## Collect data

You can collect data from development boards, from your own devices, or by uploading an existing dataset.



## Connect a fully supported development board

Get started with real hardware from a wide range of silicon vendors - fully supported by Edge Impulse.

Browse dev boards



## Use your mobile phone

Use your mobile phone to capture movement, audio or images, and even run your trained model locally. No app required.

Show QR code

+ Connect a new device



## Connected as phone\_kunh8zjd

You can collect data from this device from the **Data acquisition** page in the Edge Impulse studio.

Collecting images?

Collecting audio?

Collecting motion?

Switch to classification mode

</> This client is [open source](#).

## Classifier



## Building project...

Job started

Switch to data collection mode

</> This client is [open source](#).

Devices

+ Connect a new device

Collect data

You can collect data from development boards, from your own devices, or by uploading an existing dataset.



Connect a fully supported development board

Get started with real hardware from a wide range of silicon vendors - fully supported by Edge Impulse.

Browse dev boards



Use your mobile phone

Use your mobile phone to capture movement, audio or images, and even run your trained model locally. No app required.

Show QR code



Connected as phone\_kunh8zjd

You can collect data from this device from the **Data acquisition** page in the Edge Impulse studio.

Collecting images?

Collecting audio?

Collecting motion?

Switch to classification mode

This client is [open source](#).



Classifier



Building project...

Job started

Switch to data collection mode

This client is [open source](#).



yes

	NO	UNKNOWN	YES
178	0.07	0.06	0.87
177	0.00	0.05	0.94
176	0.00	0.06	0.94
175	0.00	0.10	0.90

# Deploy and Test your Model

Shows the **score** for (**confidence that the current sounds is**) each of the various keywords and unknown and bolds the highest score.



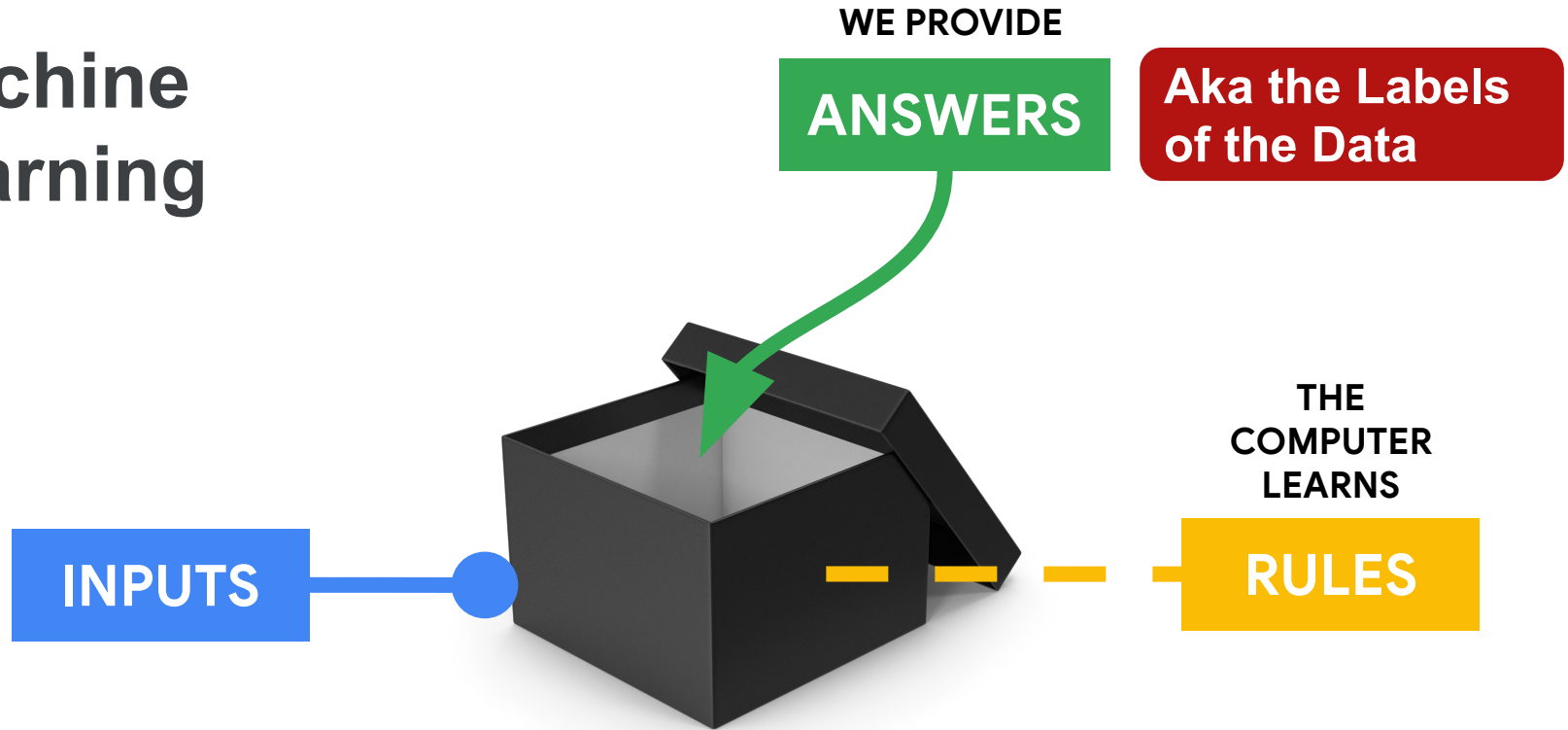
yes

	NO	UNKNOWN	YES
178	0.07	0.06	<b>0.87</b>
177	0.00	0.05	<b>0.94</b>
176	0.00	0.06	<b>0.94</b>
175	0.00	0.10	<b>0.90</b>

# Today's Agenda

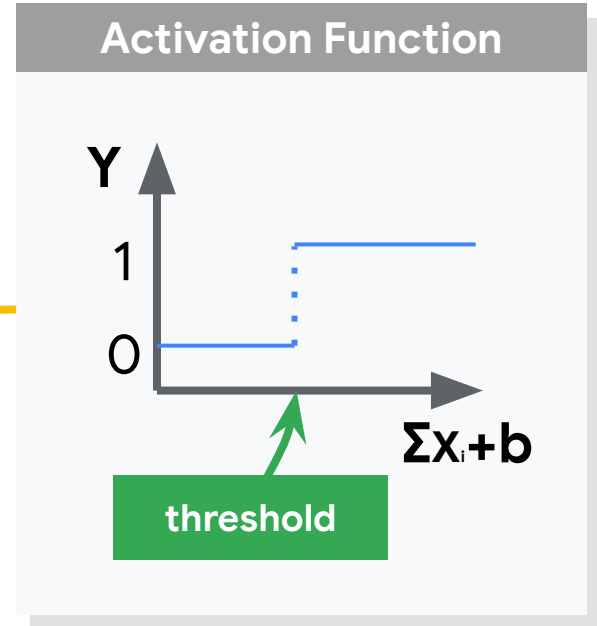
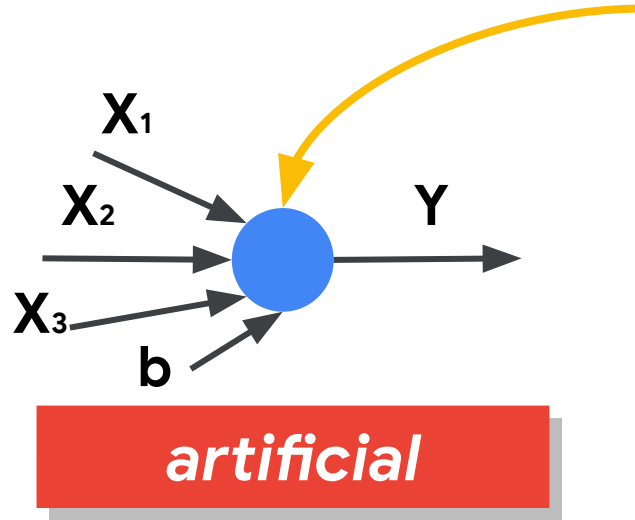
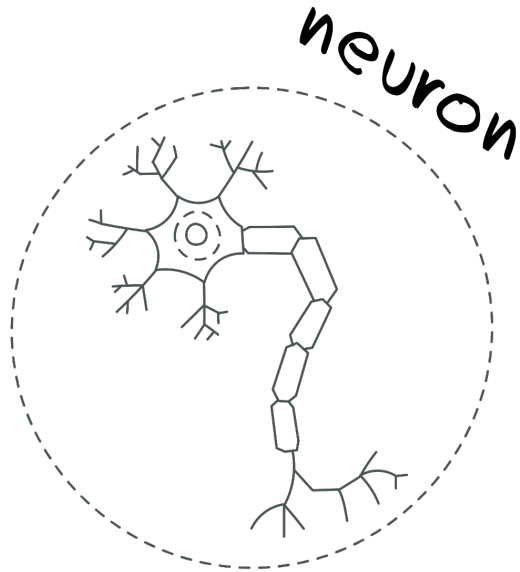
- Deep ML Background
- Hands-on Computer Vision: Thing Translator
- The Tiny Machine Learning Workflow
- Keyword Spotting (KWS) Data Collection
- KWS Preprocessing and Training
- Deployment Challenges and Opportunities for Embedded ML
- **Summary**

# Machine Learning



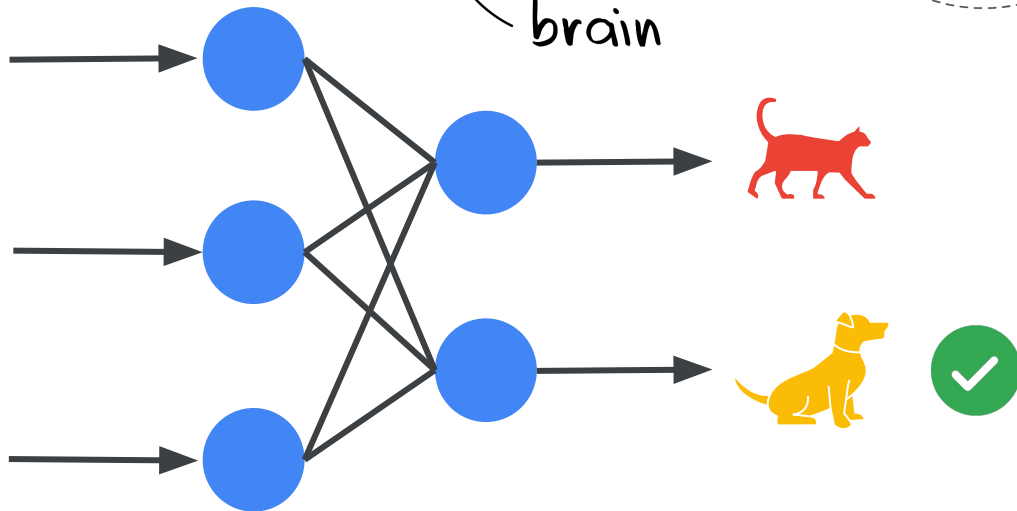
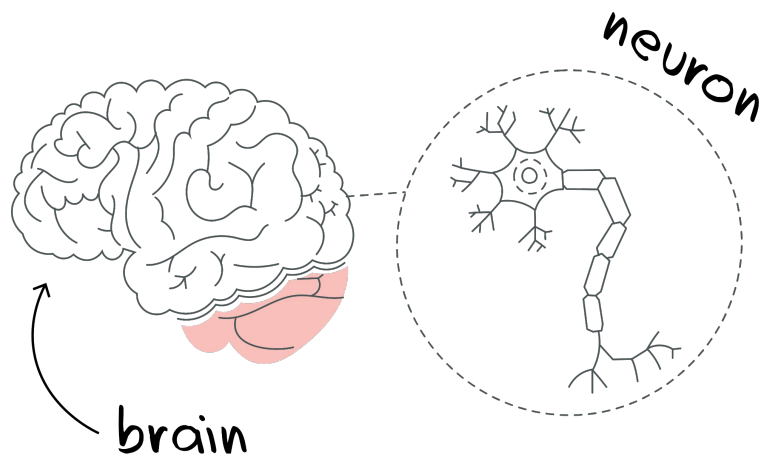


# What is a neural network?



$$Y = \Sigma w_i X_i + b$$
  
So training the model is finding the right values for  $w_i$  and  $b$

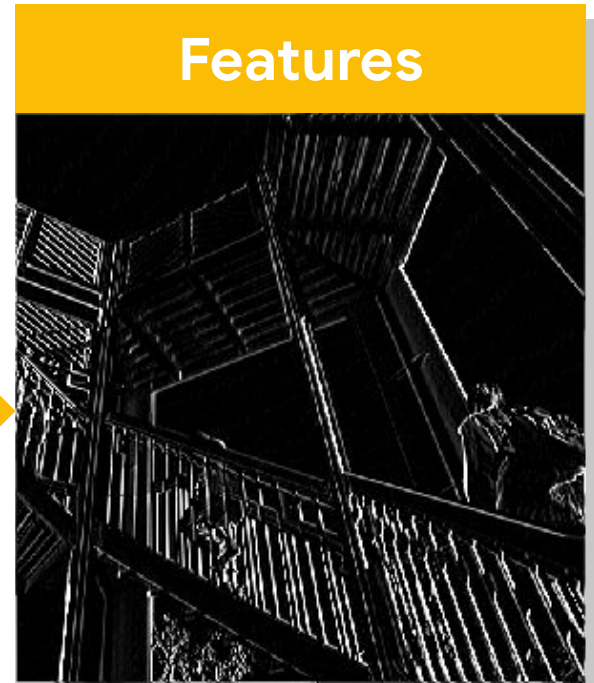
# Deep Learning with **Neural Networks**



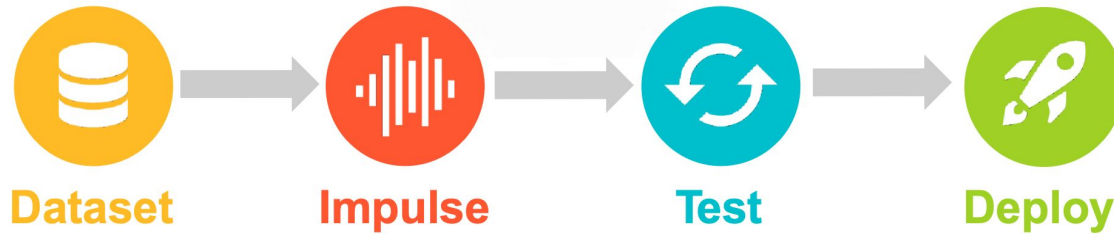
# Features can be found with **Convolutions**



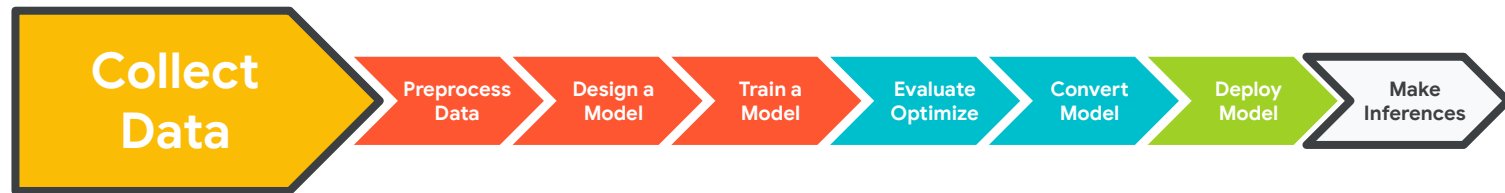
-1	0	1
-2	0	2
-1	0	1



# The **TinyML** Workflow



# The **TinyML** Workflow



**Who** will use your ML model?

**Where** will your ML model be used?

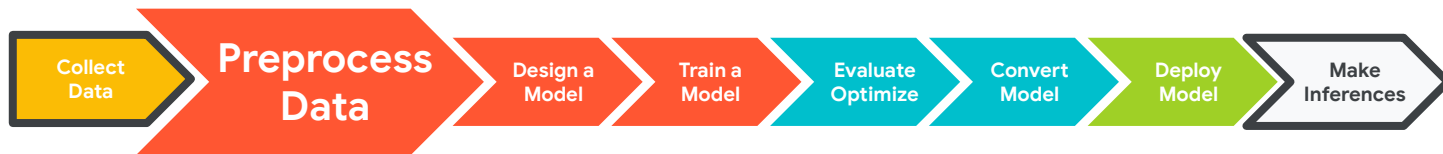
**Why** will your ML model be used?  
**Why** those Keywords?

Training Set

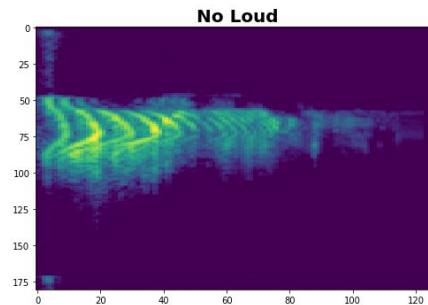
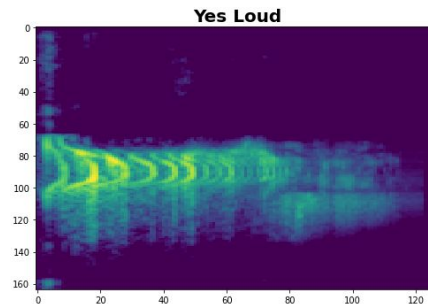
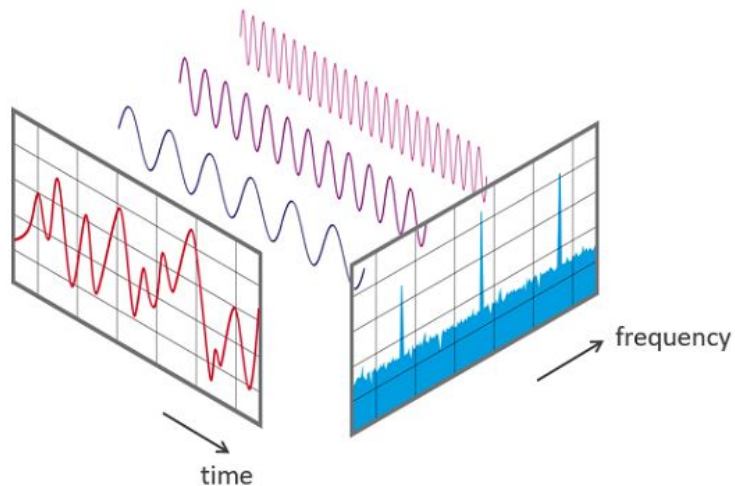
Validation Set

Test Set

# The **TinyML** Workflow



**FFT, Spectrogram, MFCC**



# The **TinyML** Workflow



<b>Confusion Matrix</b>	Actual Output = Yes	Actual Output = No
	Predicted Output = Yes	# of True Positive
Predicted Output = No	# of False Negative <i>Type 2 Error</i>	# of True Negative

# The **TinyML** Workflow



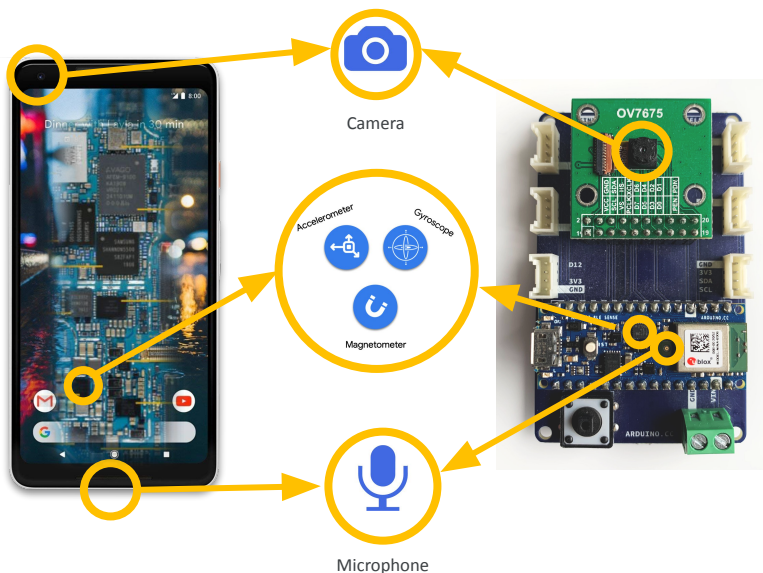
## Quantization

**Reduces the precision** of numbers used in a model which results in:

- **smaller model size**
- **faster computation**



# The **TinyML** Workflow



**Edge Impulse  
Simplifies  
Deployment**

# The **TinyML** Workflow



yes

Shows the **score** for (**confidence that the current sounds is**) each of the various keywords and unknown and **bolds the highest score**.

178	0.07	0.06	<b>0.87</b>
177	0.00	0.05	<b>0.94</b>
176	0.00	0.00	<b>0.94</b>



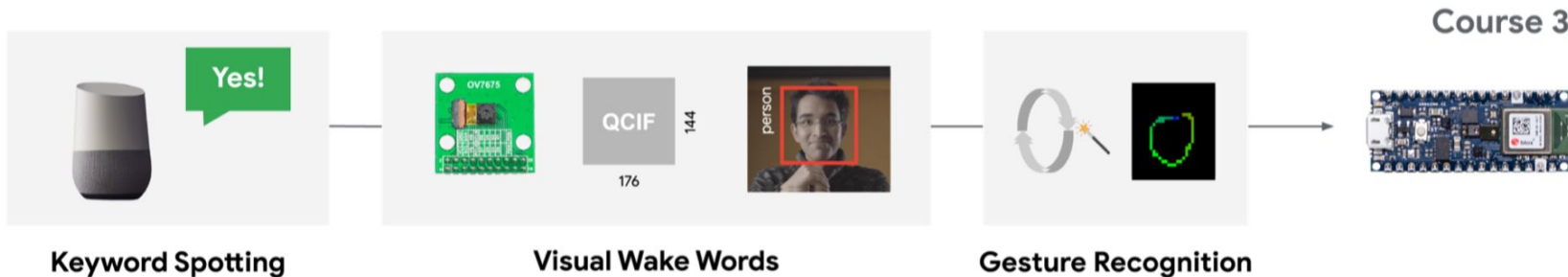
**Fundamentals of *TinyML***

- Course 1**
- Neural Network
  - Filters
  - Regression
  - Loss Function
  - Preprocessing
  - Data augmentation
  - Inference
  - Responsible AI
  - CNNs/ DNNs
  - Classification
  - Gradient Descent

**Applications of *TinyML***



**Deploying *TinyML***



**Managing *TinyML***

**Better Data = Better Models!**



# Hands on Embedded ML (Vision and Audio)



*Brian Plancher*  
*Harvard John A. Paulson School of Engineering and Applied Sciences*  
*brianplancher.com*

