

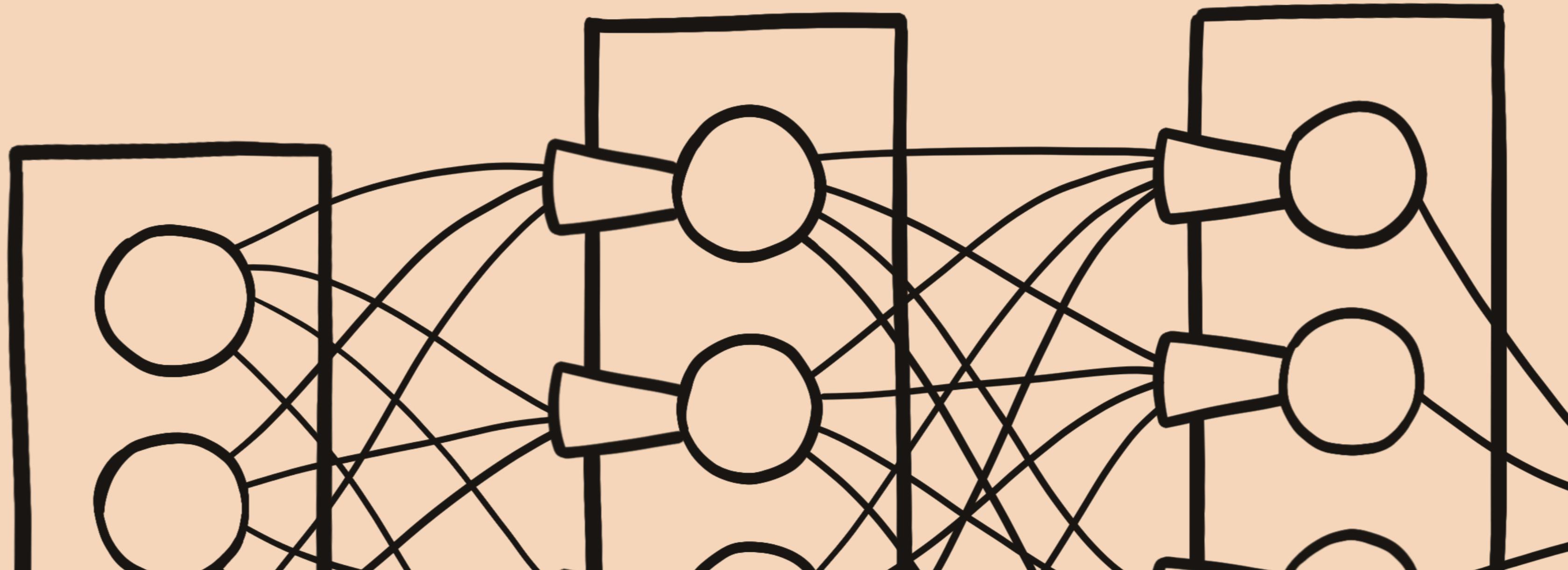
Can Neural Networks Make Me a Better Parent?

A tale in three acts

Andrew Hao

@andrewhao

ahao@lyft.com

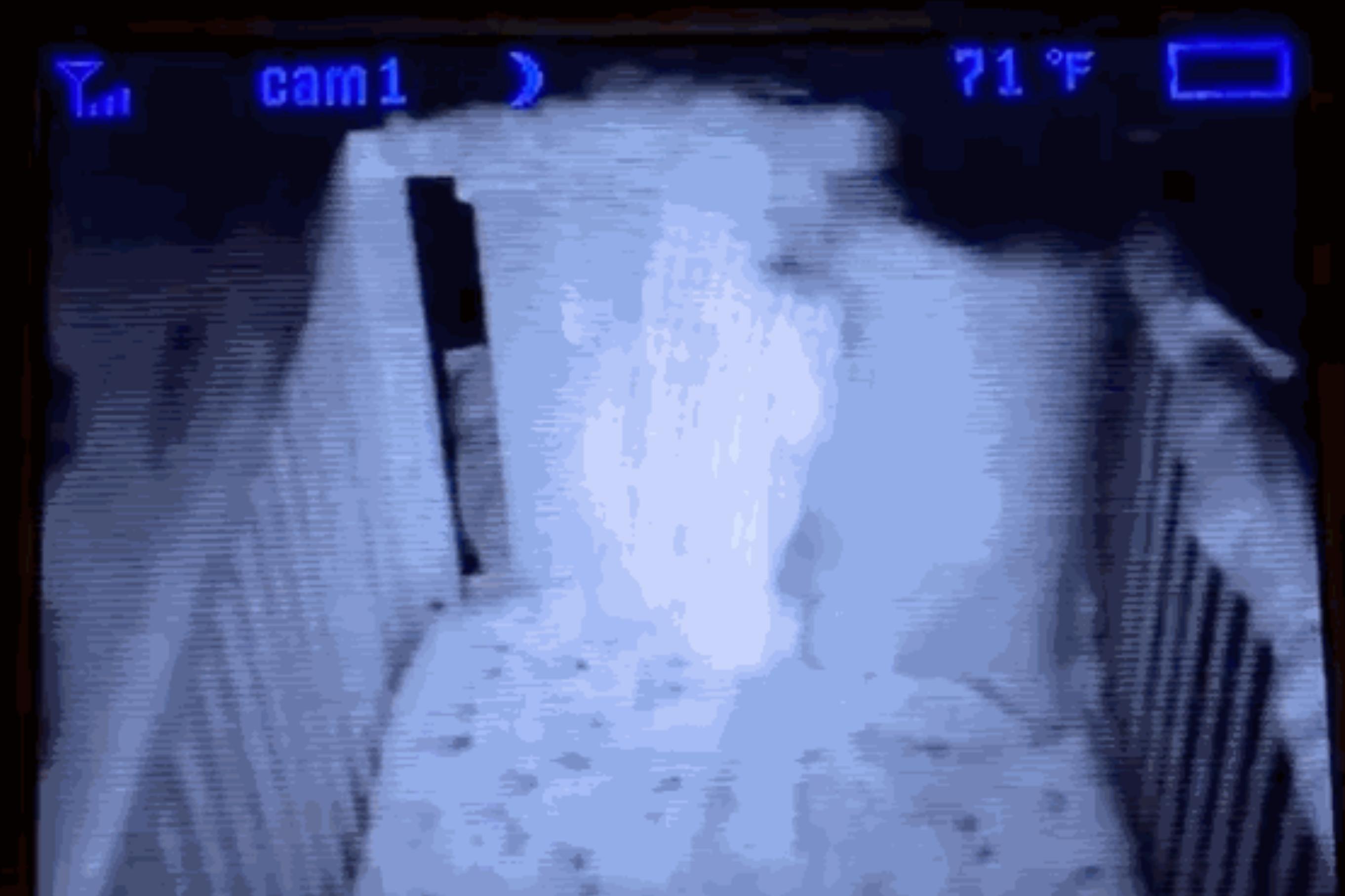


Act I

Night & Chaos

put the little one down





helplessness

frustration

anger

despair

I need datp

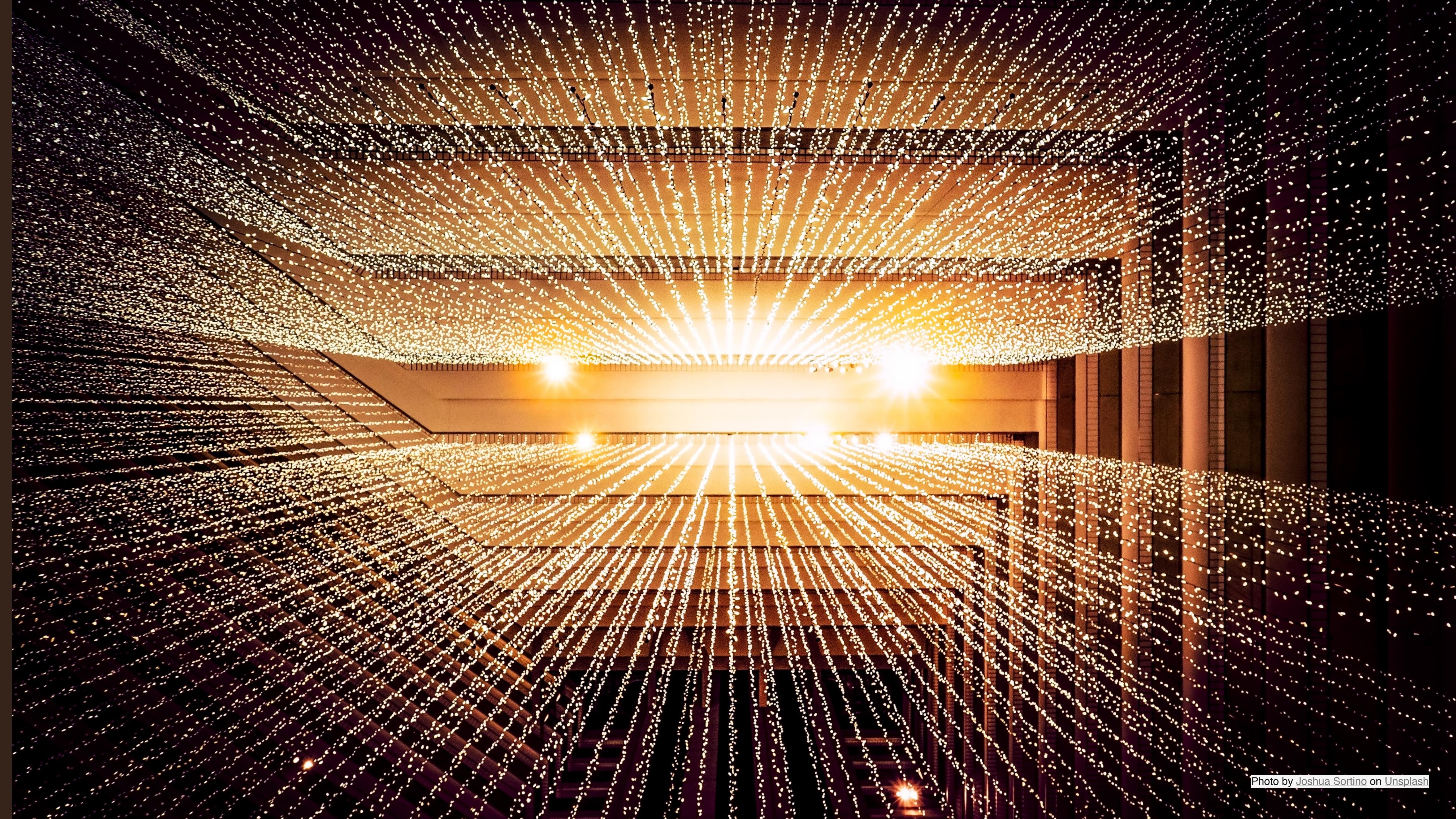


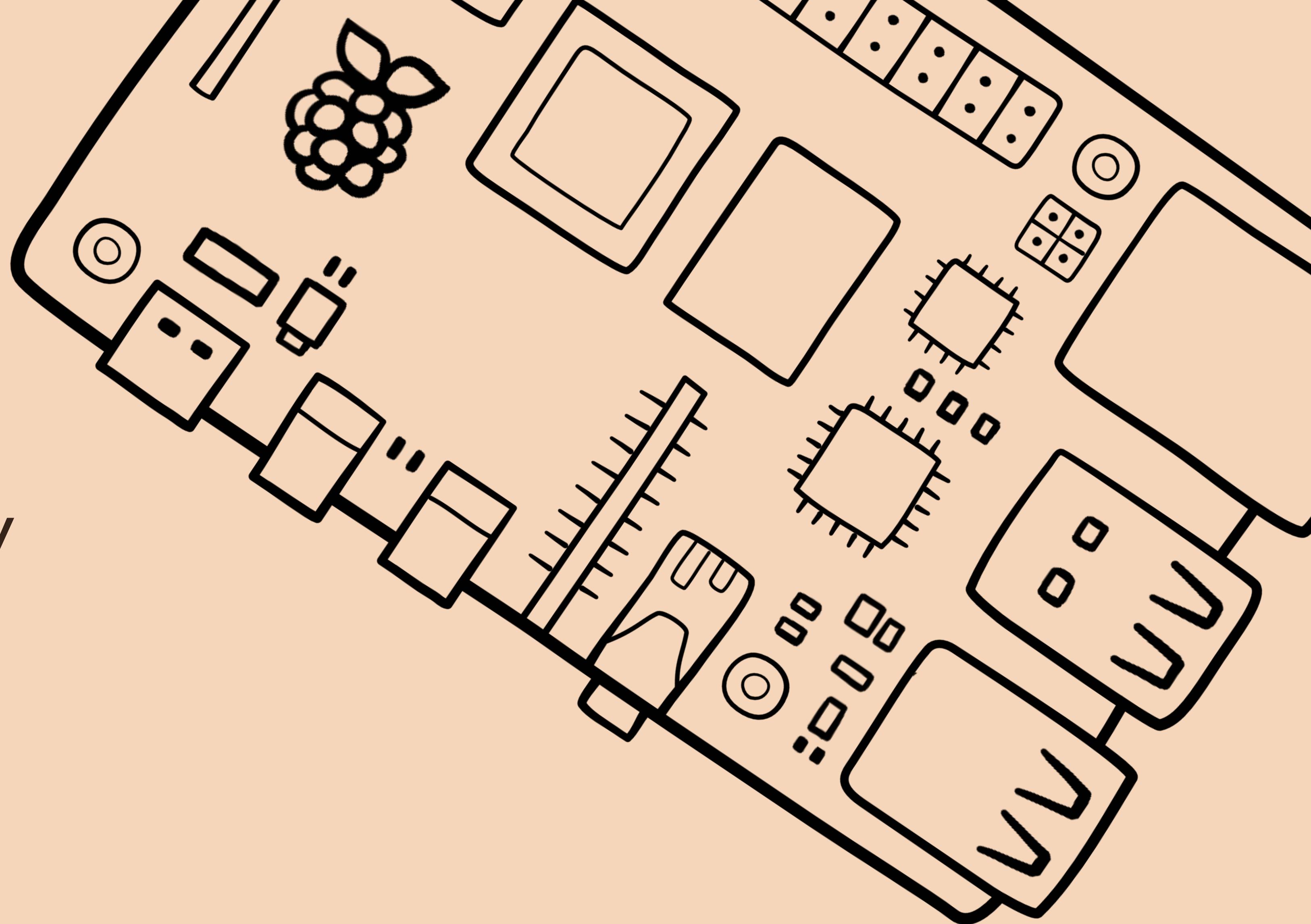
Photo by Joshua Sortino on Unsplash

Act II

Enter the Machine

Parts List

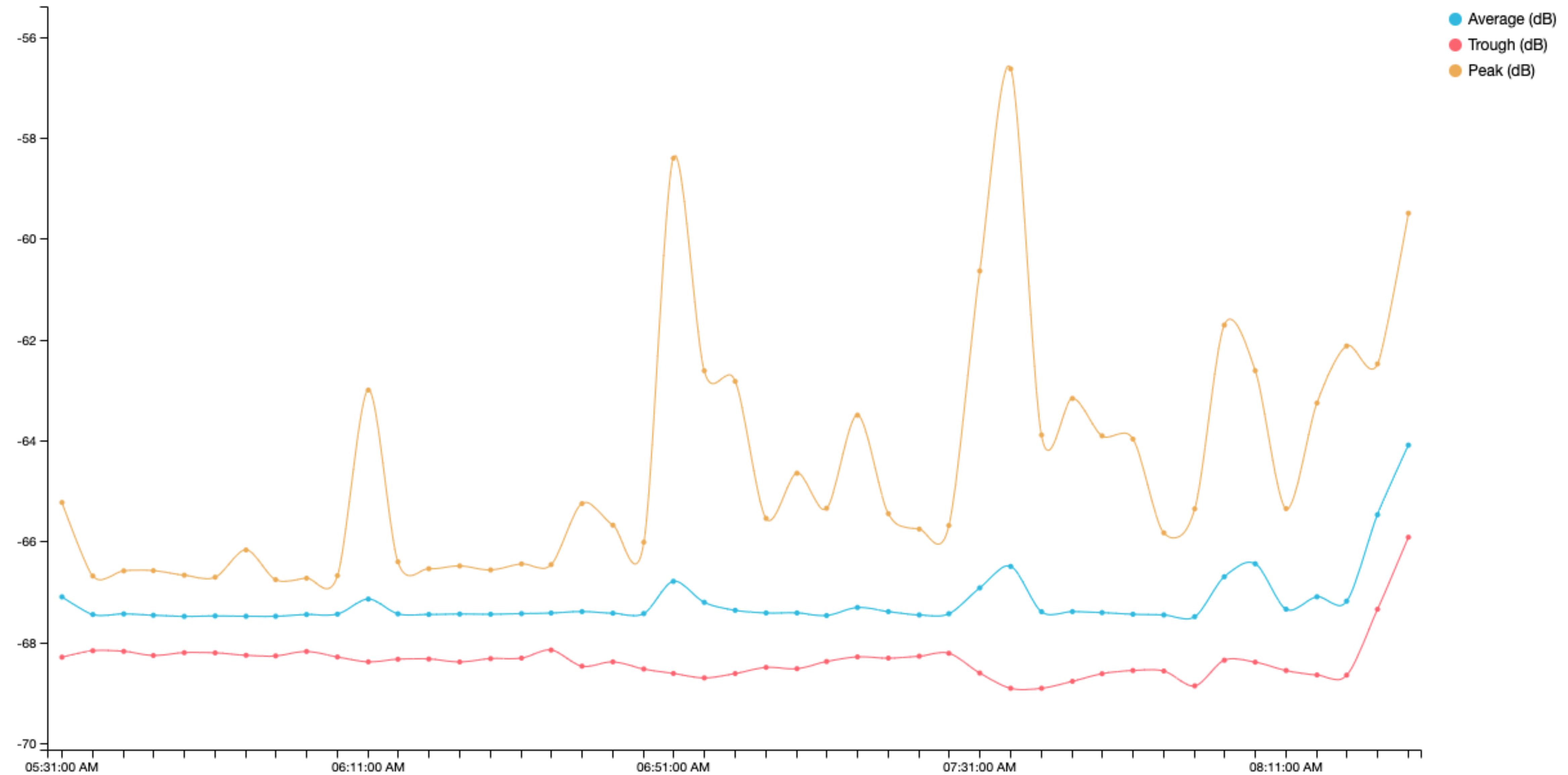
- Raspberry Pi
- USB microphone
- DHT22 temp/humidity sensor

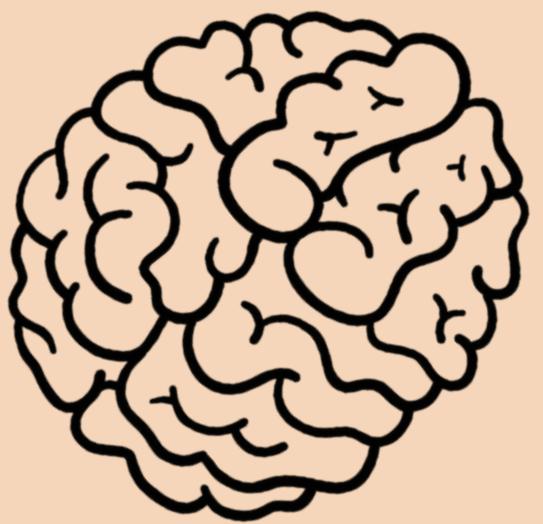


```
$ arecord --device=hw:1,0 --format S16_LE --rate 22050 -c1 -d 10 "${RECORDING_FILE}"  
$ sox -V3 ${RECORDING_FILE} -n stats 2>&1 | grep dB
```

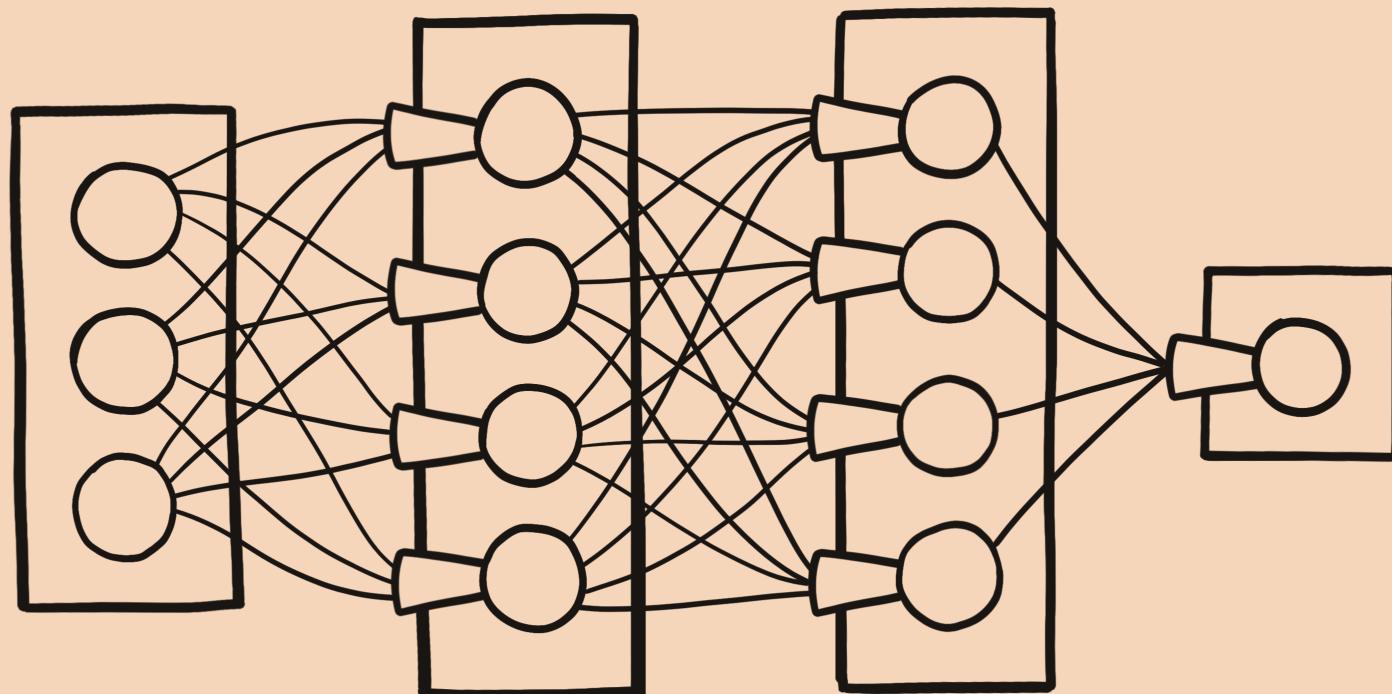
Pk lev dB	-57.05
RMS lev dB	-67.10
RMS Pk dB	-64.94
RMS Tr dB	-68.26

Last Day Noise





Deep learning & neural nets



Simple Audio Recognition | Te X

webcache.googleusercontent.com/search?q=cache:IAxOQJXHducJ:https://www.tensorflow.org/tutorials/sequences/audio...

TensorFlow [Install](#) [Learn](#) [API](#) [Resources](#) [Community](#) [More](#) [Search](#) [Language](#) [GitHub](#)

Overview [Tutorials](#) [Guide](#) [TF 2.0 RC](#)

Get started with TensorFlow

Learn and use ML

Research and experimentation

ML at production scale

Generative models

Distributed training

Images

Sequences

Load data

Data representation

Non-ML

TensorFlow 2.0 RC is available [Learn more](#)

TensorFlow > Learn > TensorFlow Core > Tutorials ★★★★★

Simple Audio Recognition

This tutorial will show you how to build a basic speech recognition network that recognizes ten different words. It's important to know that real speech and audio recognition systems are much more complex, but like MNIST for images, it should give you a basic understanding of the techniques involved. Once you've completed this tutorial, you'll have a model that tries to classify a one second audio clip as either silence, an unknown word, "yes", "no", "up", "down", "left", "right", "on", "off", "stop", or "go". You'll also be able to take this model and run it in an Android application.

Preparation

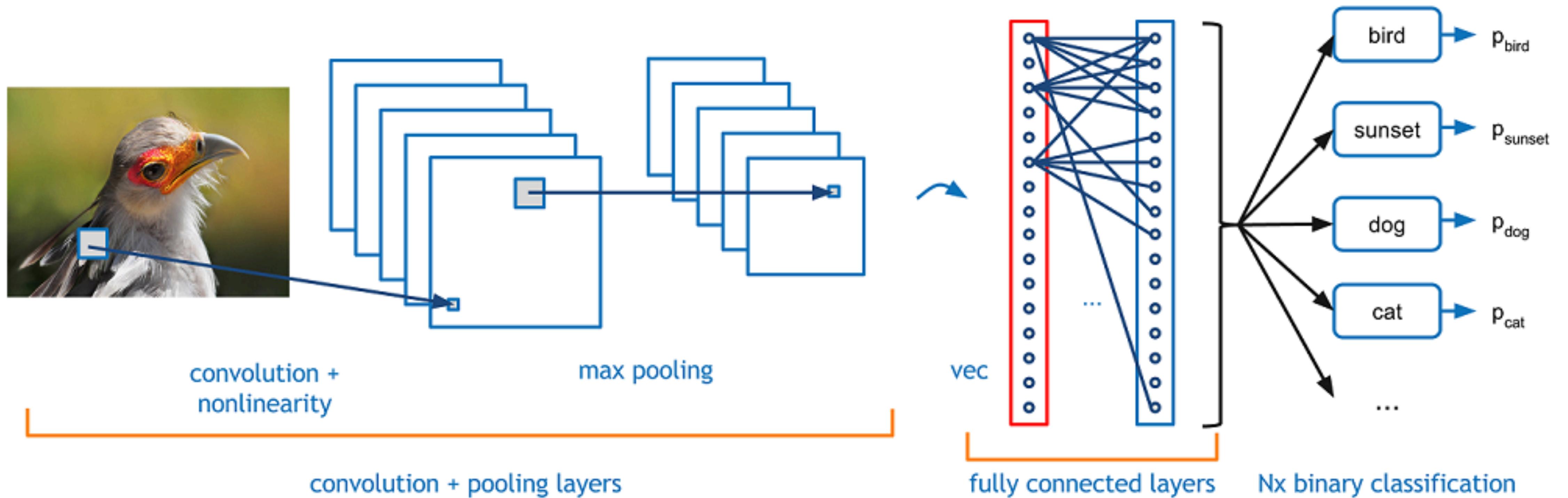
You should make sure you have TensorFlow installed, and since the script downloads over 1GB of training data, you'll need a good internet connection and enough free space on your machine. The training process itself can take several hours, so make sure you have a machine available for that long.

Training

Contents

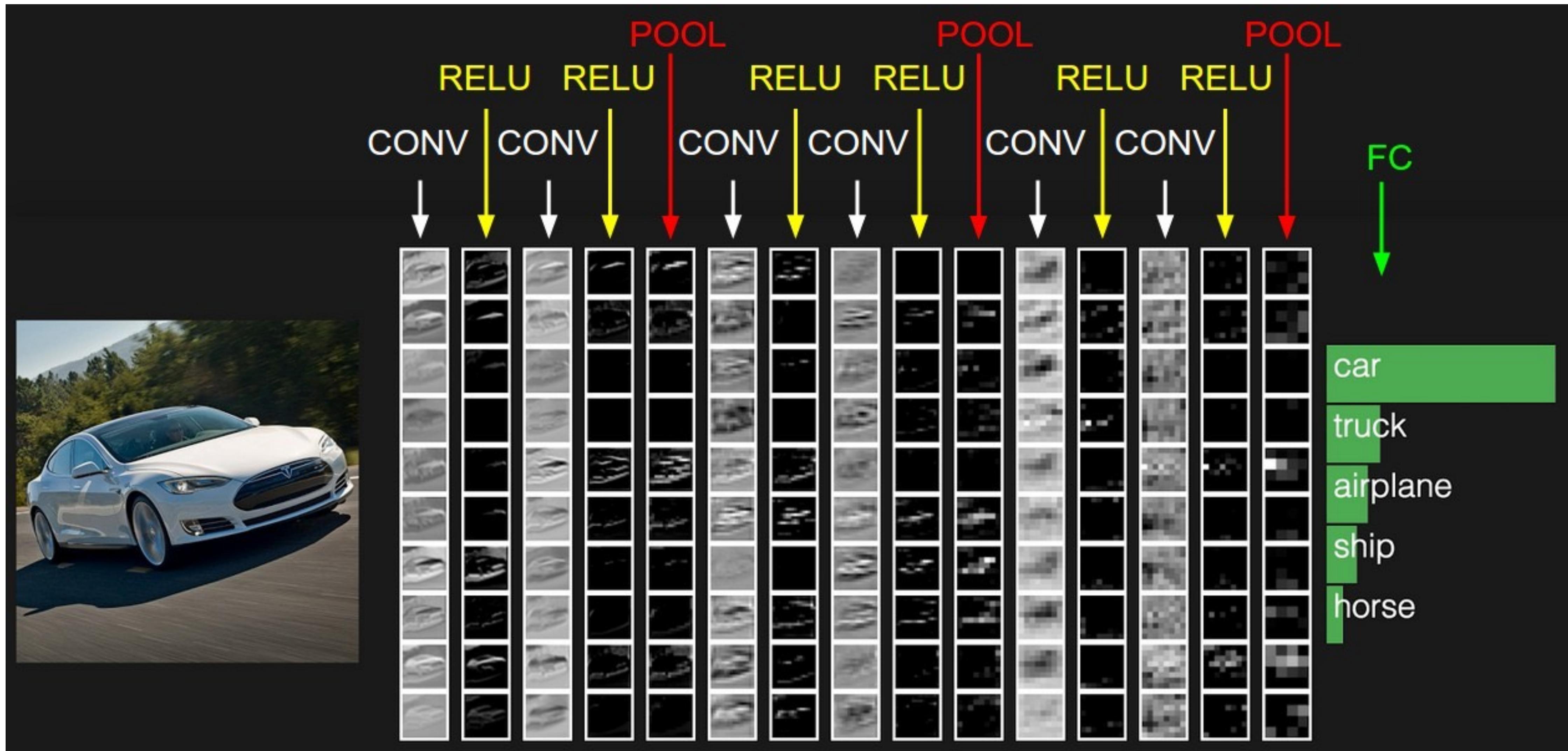
- Preparation
- Training
- Confusion Matrix
- Validation
- Tensorboard
- Training Finished
- Running the Model in an Android App
- How does this Model Work?
- Streaming Accuracy
- RecognizeComm...
- Advanced Training
- Custom Training Data
- Unknown Class
- Background Noise
- Silence
- Time Shifting

convolutional neural network

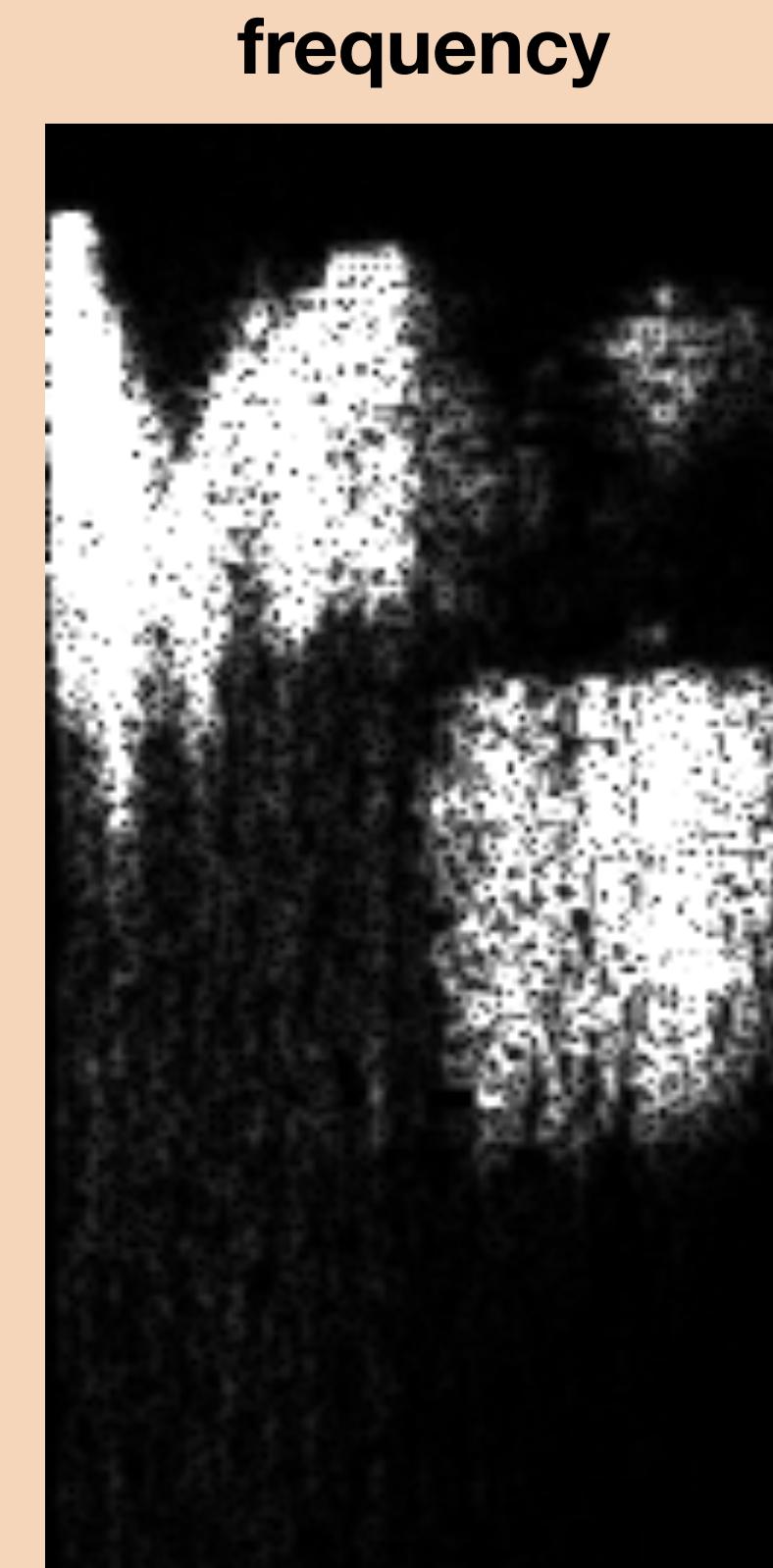


<https://www.fastcompany.com/3037882/how-flickrs-deep-learning-algorithms-see-whats-in-your-photos>

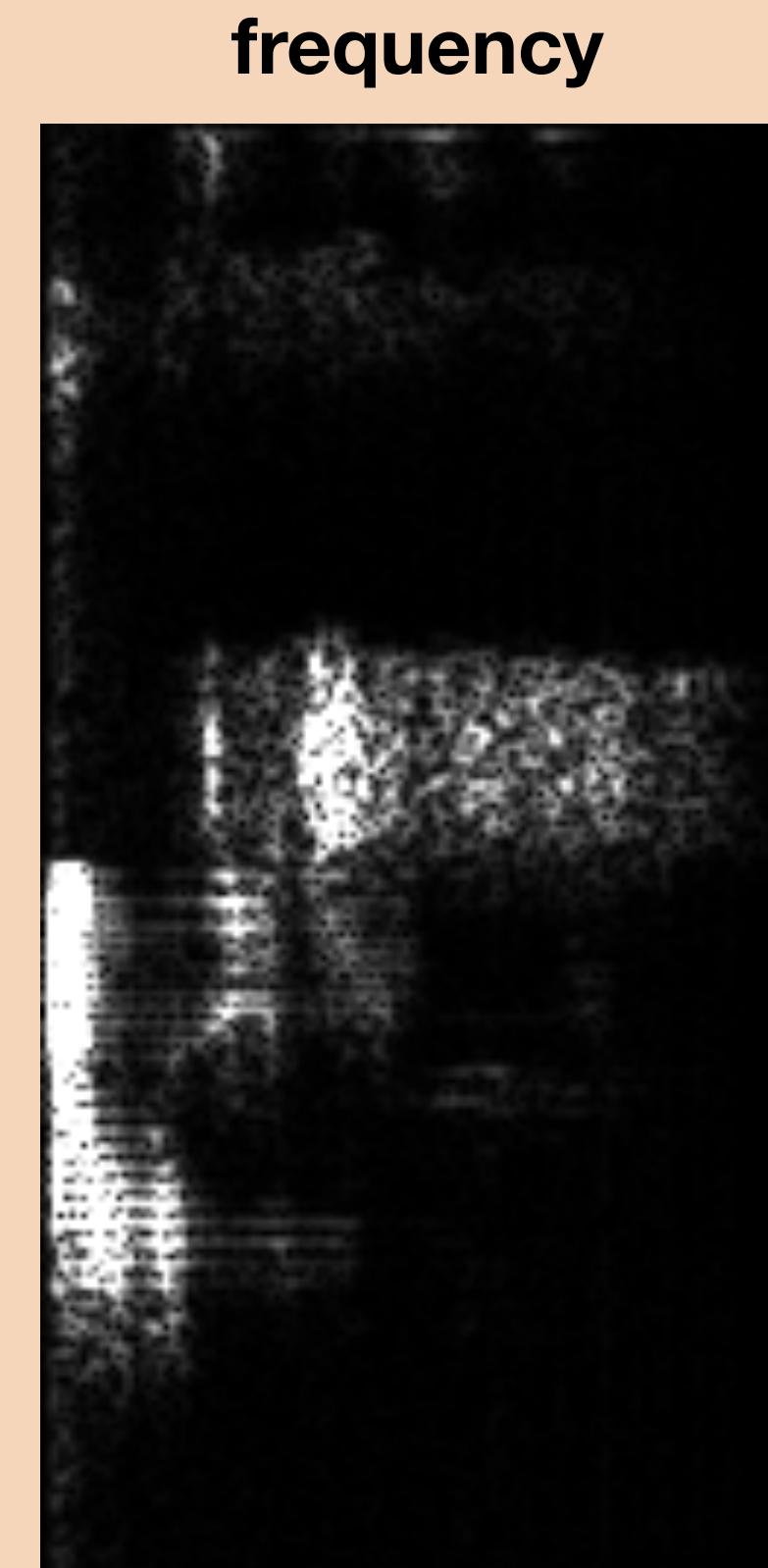
CNNs



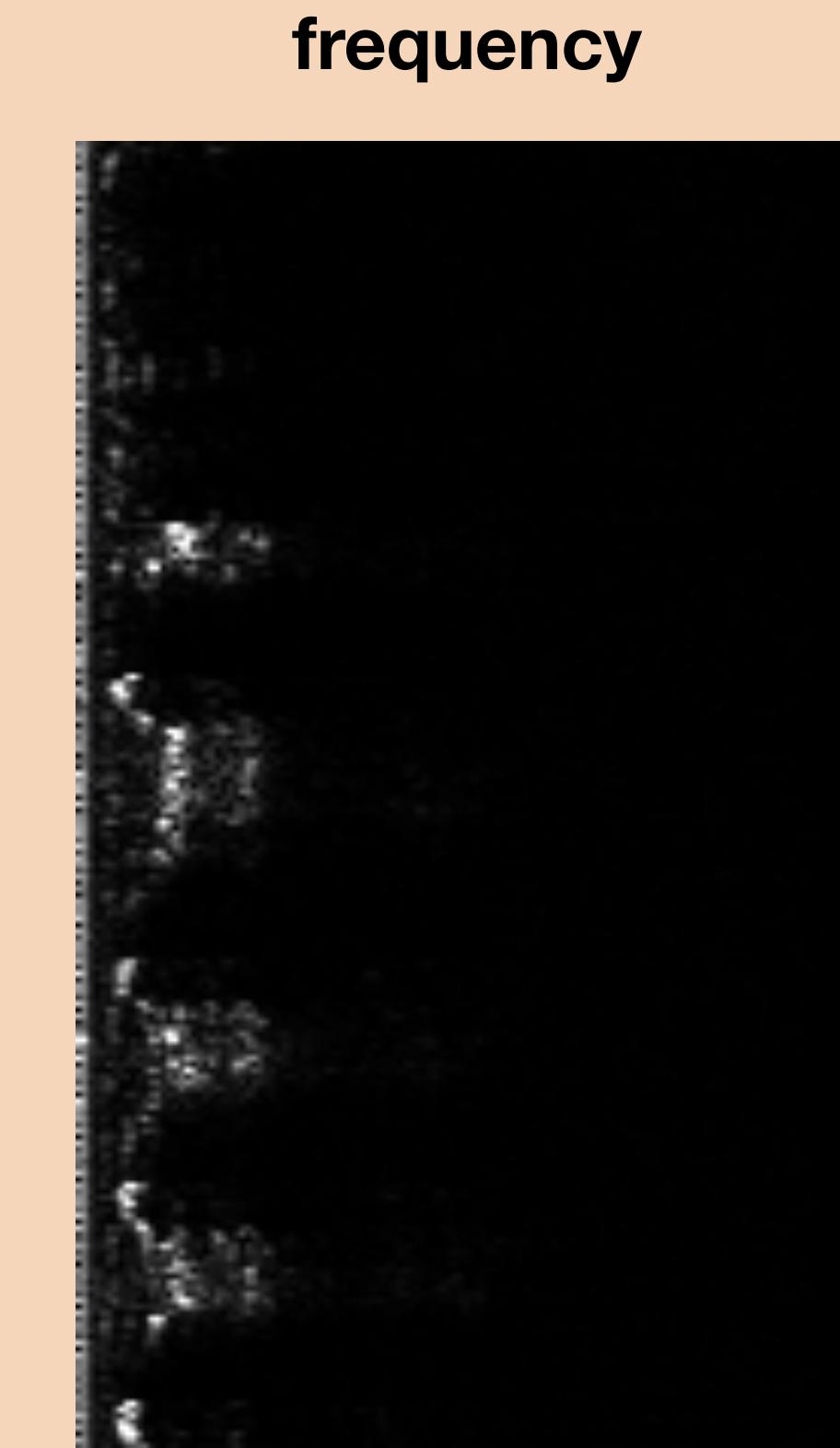
audio spectrogram = image



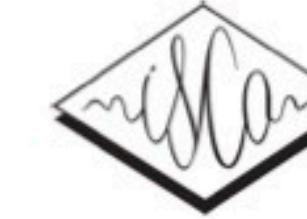
“yes”



“Sheila”



“wah, wah, wah”



Convolutional Neural Networks for Small-footprint Keyword Spotting

Tara N. Sainath, Carolina Parada

Google, Inc. New York, NY, U.S.A

{tsainath, carolinap}@google.com

Abstract

We explore using Convolutional Neural Networks (CNNs) for a small-footprint keyword spotting (KWS) task. CNNs are attractive for KWS since they have been shown to outperform DNNs with far fewer parameters. We consider two different applications in our work, one where we limit the number of multiplications of the KWS system, and another where we limit the number of parameters. We present new CNN architectures to address the constraints of each applications. We find that the CNN architectures offer between a 27-44% relative improvement in false reject rate compared to a DNN, while fitting into the constraints of each application.

1. Introduction

With the rapid development of mobile devices, speech-related technologies are becoming increasingly popular. For example, Google offers the ability to search by voice [1] on Android phones, while personal assistants such as Google Now, Apple's Siri, Microsoft's Cortana and Amazon's Alexa, all utilize speech recognition to interact with these systems. Google has enabled a fully hands-free speech recognition experience, known as "Ok Google" [2], which continuously listens for specific keywords to initiate voice input. This keyword spotting (KWS) system runs on mobile devices, and therefore must have

duce feature variation. While DNNs of sufficient size could indeed capture translational invariance, this requires large networks with lots of training examples. CNNs on the other hand capture translational invariance with far fewer parameters by averaging the outputs of hidden units in different local time and frequency regions.

We are motivated to look at CNNs for KWS given the benefits CNNs have shown over DNNs with respect to improved performance and reduced model size [4, 5, 6]. In this paper, we look at two applications of CNNs for KWS. First, we consider the problem where we must limit the overall computation of our KWS system, that is parameters and multiplies. With this constraint, typical architectures that work well for CNNs and pool in frequency only [8], cannot be used here. Thus, we introduce a novel CNN architecture which does not pool but rather strides the filter in frequency, to abide within the computational constraints issue. Second, we consider limiting the total number of parameters of our KWS system. For this problem, we show we can improve performance by pooling in time and frequency, the first time this has been shown to be effective for speech without using multiple convolutional blocks [5, 9].

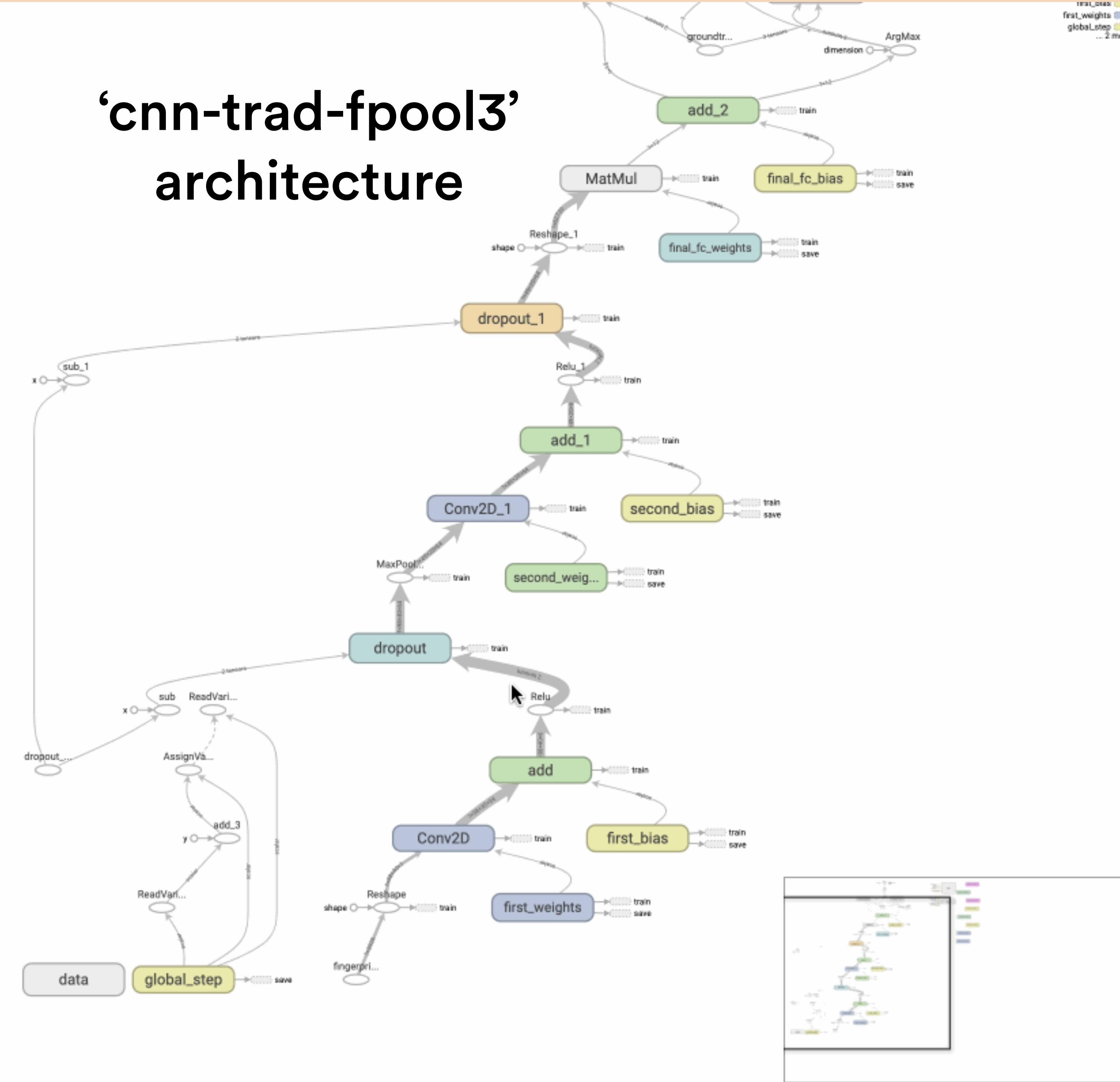
We evaluate our proposed CNN architectures on a KWS task consisting of 14 different phrases. Performance is measured by looking at the false reject (FR) rate at the operating threshold of 1 false alarm (FA) per hour. In the task where we

```

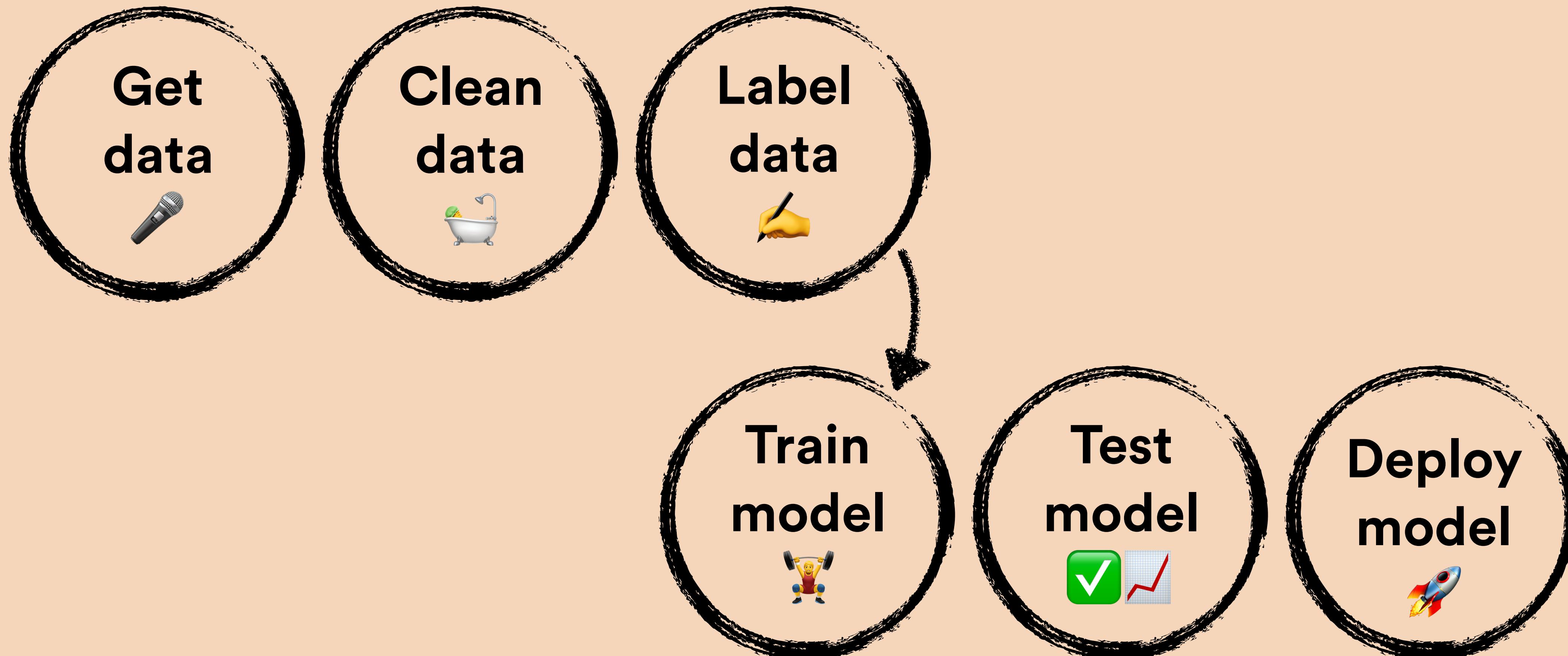
(fingerprint_input)
  v
[Conv2D]<- (weights)
  v
[BiasAdd]<- (bias)
  v
[Relu]
  v
[MaxPool]
  v
[Conv2D]<- (weights)
  v
[BiasAdd]<- (bias)
  v
[Relu]
  v
[MaxPool]
  v
[MatMul]<- (weights)
  v
[BiasAdd]<- (bias)
  v

```

‘cnn-trad-fpool3’ architecture

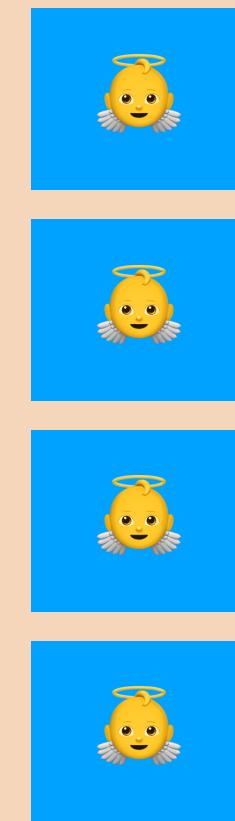
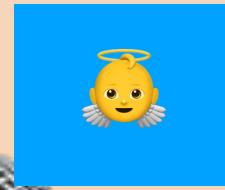


The Secret Life of ML Models





- Record data continuously in a cronjob
- Upload to the cloud

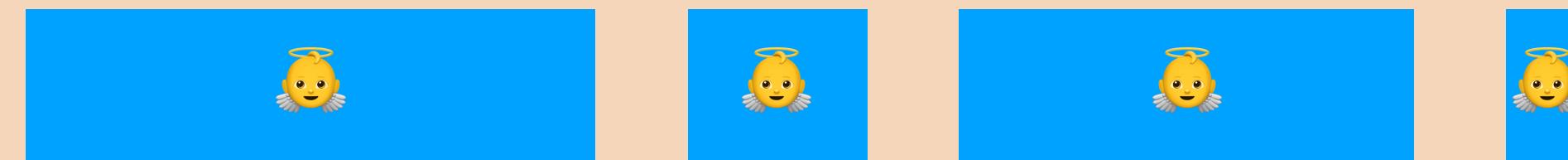


Clean data



```
# Concatenate  
sox crying/**/*.wav tmp/batch-crying.wav
```

- Align to a frame
- Normalize volume



Clean data



- Align to a frame
- Normalize volume

```
# Split
ffmpeg -i tmp/batch-crying.wav -
f segment -segment_time 10.1 -c
```



Clean data



- Align to a frame
- Normalize volume

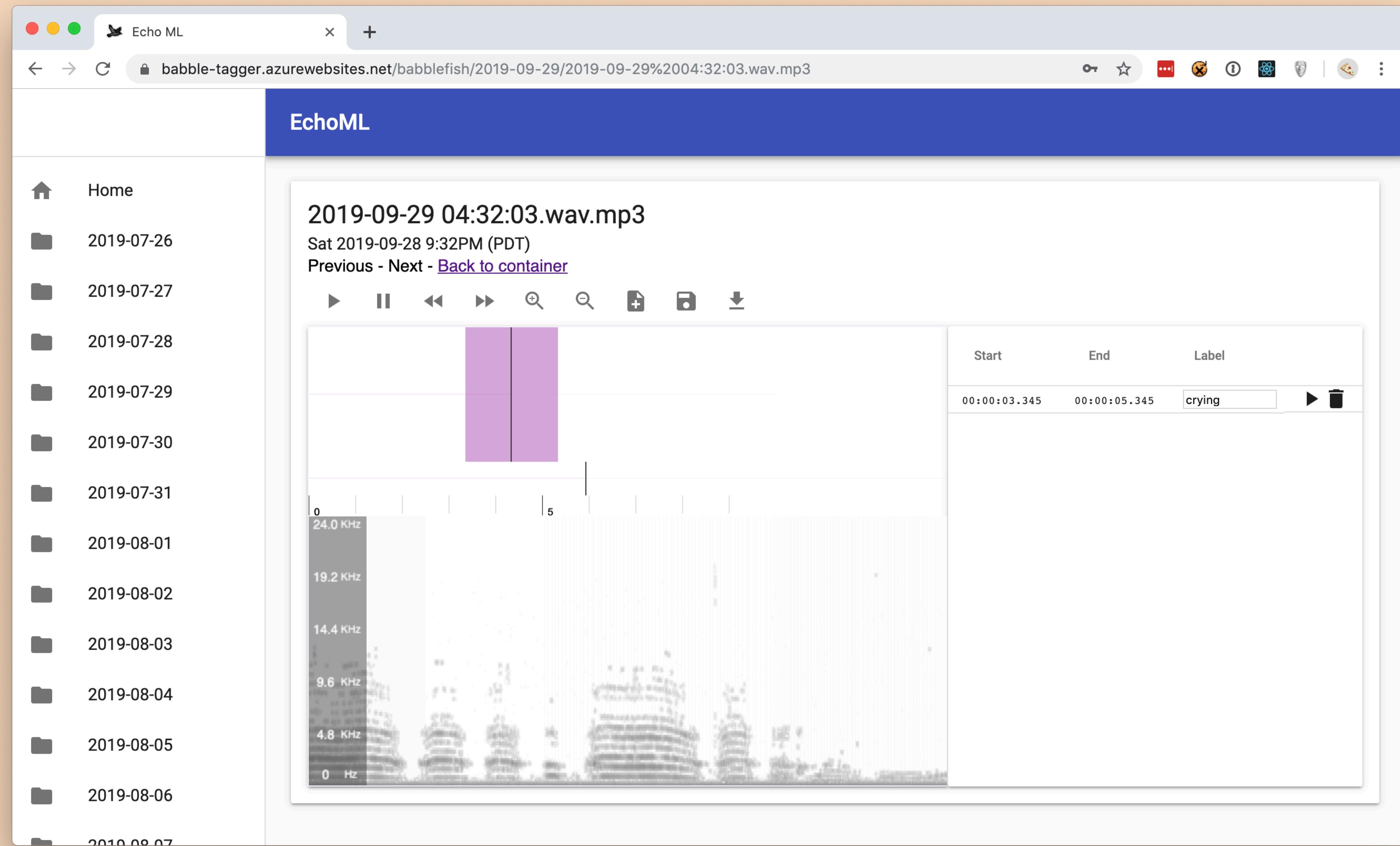
```
# Trim to exactly 10000ms, boost 45
# dB, resample @ 22050 hz
sox $in data/crying/$out vol 45 dB
trim 0 10 rate 22050
```



Label data



Using EchoML to classify data



Label data



- Bucket each audio sample in a folder by label

crying



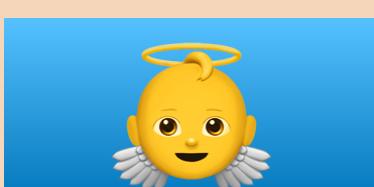
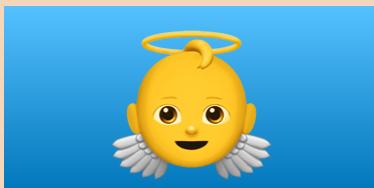
whining



dog_barking



city_noise



Train model



```
python app/train.py
--data_url=
--data_dir=./data
--wanted_words=room_empty,whining,crying
--sample_rate=22050
--clip_duration_ms=10000
--how_many_training_steps=400,50
--train_dir=./training
```



```
Default (tmux)

I1002 03:58:24.679805 4741027264 train.py:287] Step #430: rate 0.000100, accuracy 99.0%, cross entropy 0.04[95/3172]
INFO:tensorflow:Step #431: rate 0.000100, accuracy 99.0%, cross entropy 0.035736
I1002 03:59:02.294903 4741027264 train.py:287] Step #431: rate 0.000100, accuracy 99.0%, cross entropy 0.035736
INFO:tensorflow:Step #432: rate 0.000100, accuracy 98.0%, cross entropy 0.118002
I1002 03:59:40.747909 4741027264 train.py:287] Step #432: rate 0.000100, accuracy 98.0%, cross entropy 0.118002
INFO:tensorflow:Step #433: rate 0.000100, accuracy 98.0%, cross entropy 0.073030
I1002 04:00:19.569001 4741027264 train.py:287] Step #433: rate 0.000100, accuracy 98.0%, cross entropy 0.073030
INFO:tensorflow:Step #434: rate 0.000100, accuracy 97.0%, cross entropy 0.092732
I1002 04:00:55.971426 4741027264 train.py:287] Step #434: rate 0.000100, accuracy 97.0%, cross entropy 0.092732
INFO:tensorflow:Step #435: rate 0.000100, accuracy 100.0%, cross entropy 0.042072
I1002 04:01:35.398041 4741027264 train.py:287] Step #435: rate 0.000100, accuracy 100.0%, cross entropy 0.042072
INFO:tensorflow:Step #436: rate 0.000100, accuracy 100.0%, cross entropy 0.043324
I1002 04:02:13.348515 4741027264 train.py:287] Step #436: rate 0.000100, accuracy 100.0%, cross entropy 0.043324
INFO:tensorflow:Step #437: rate 0.000100, accuracy 100.0%, cross entropy 0.038082
I1002 04:02:56.560698 4741027264 train.py:287] Step #437: rate 0.000100, accuracy 100.0%, cross entropy 0.038082
INFO:tensorflow:Step #438: rate 0.000100, accuracy 98.0%, cross entropy 0.068574
I1002 04:03:35.357985 4741027264 train.py:287] Step #438: rate 0.000100, accuracy 98.0%, cross entropy 0.068574
INFO:tensorflow:Step #439: rate 0.000100, accuracy 99.0%, cross entropy 0.058624
I1002 04:04:13.681237 4741027264 train.py:287] Step #439: rate 0.000100, accuracy 99.0%, cross entropy 0.058624
INFO:tensorflow:Step #440: rate 0.000100, accuracy 99.0%, cross entropy 0.061955
I1002 04:04:51.182194 4741027264 train.py:287] Step #440: rate 0.000100, accuracy 99.0%, cross entropy 0.061955
INFO:tensorflow:Step #441: rate 0.000100, accuracy 97.0%, cross entropy 0.091335
I1002 04:05:31.424205 4741027264 train.py:287] Step #441: rate 0.000100, accuracy 97.0%, cross entropy 0.091335
INFO:tensorflow:Step #442: rate 0.000100, accuracy 100.0%, cross entropy 0.034299
I1002 04:06:07.957736 4741027264 train.py:287] Step #442: rate 0.000100, accuracy 100.0%, cross entropy 0.034299
INFO:tensorflow:Step #443: rate 0.000100, accuracy 99.0%, cross entropy 0.080818
I1002 04:06:49.358937 4741027264 train.py:287] Step #443: rate 0.000100, accuracy 99.0%, cross entropy 0.080818
INFO:tensorflow:Step #444: rate 0.000100, accuracy 99.0%, cross entropy 0.041448
I1002 04:07:28.075536 4741027264 train.py:287] Step #444: rate 0.000100, accuracy 99.0%, cross entropy 0.041448
INFO:tensorflow:Step #445: rate 0.000100, accuracy 100.0%, cross entropy 0.030128
I1002 04:08:06.771522 4741027264 train.py:287] Step #445: rate 0.000100, accuracy 100.0%, cross entropy 0.030128
INFO:tensorflow:Step #446: rate 0.000100, accuracy 99.0%, cross entropy 0.045735
I1002 04:08:47.032963 4741027264 train.py:287] Step #446: rate 0.000100, accuracy 99.0%, cross entropy 0.045735
INFO:tensorflow:Step #447: rate 0.000100, accuracy 99.0%, cross entropy 0.043301
I1002 04:09:25.154629 4741027264 train.py:287] Step #447: rate 0.000100, accuracy 99.0%, cross entropy 0.043301
INFO:tensorflow:Step #448: rate 0.000100, accuracy 97.0%, cross entropy 0.119577
I1002 04:10:01.199583 4741027264 train.py:287] Step #448: rate 0.000100, accuracy 97.0%, cross entropy 0.119577
INFO:tensorflow:Step #449: rate 0.000100, accuracy 97.0%, cross entropy 0.104553
```

Train mode



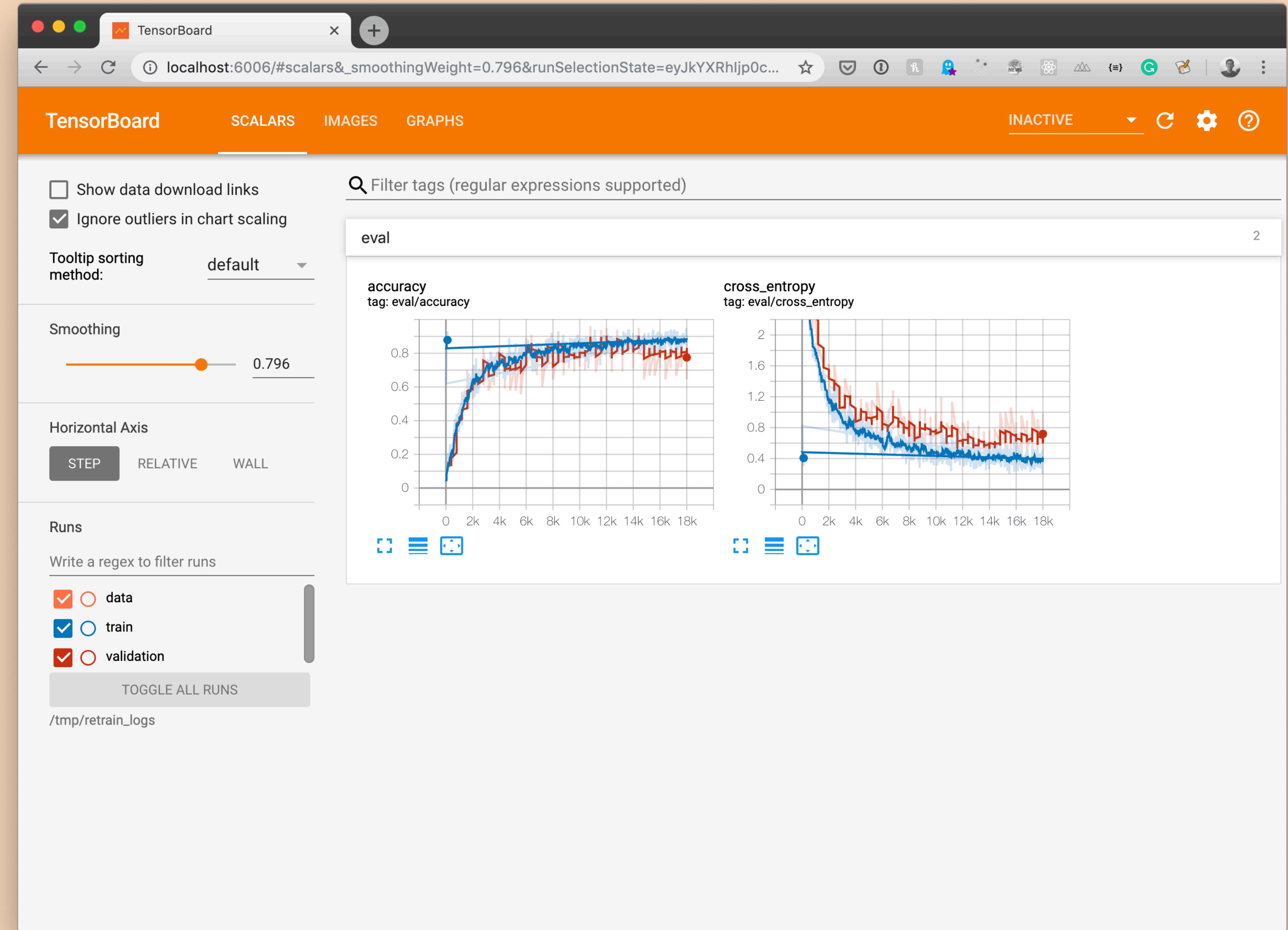
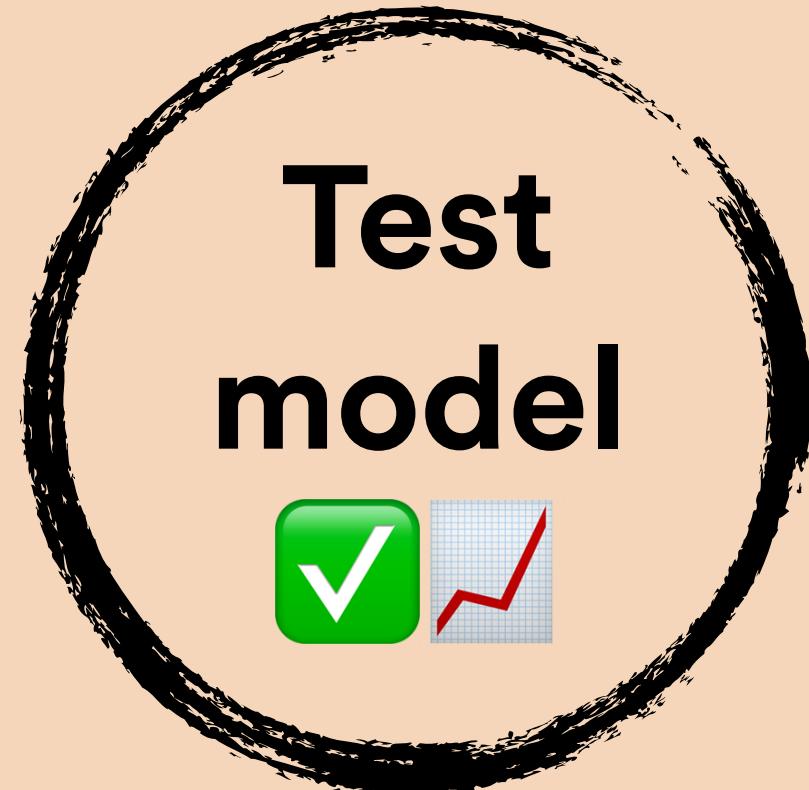
```
Default (tmux) [58/3172]

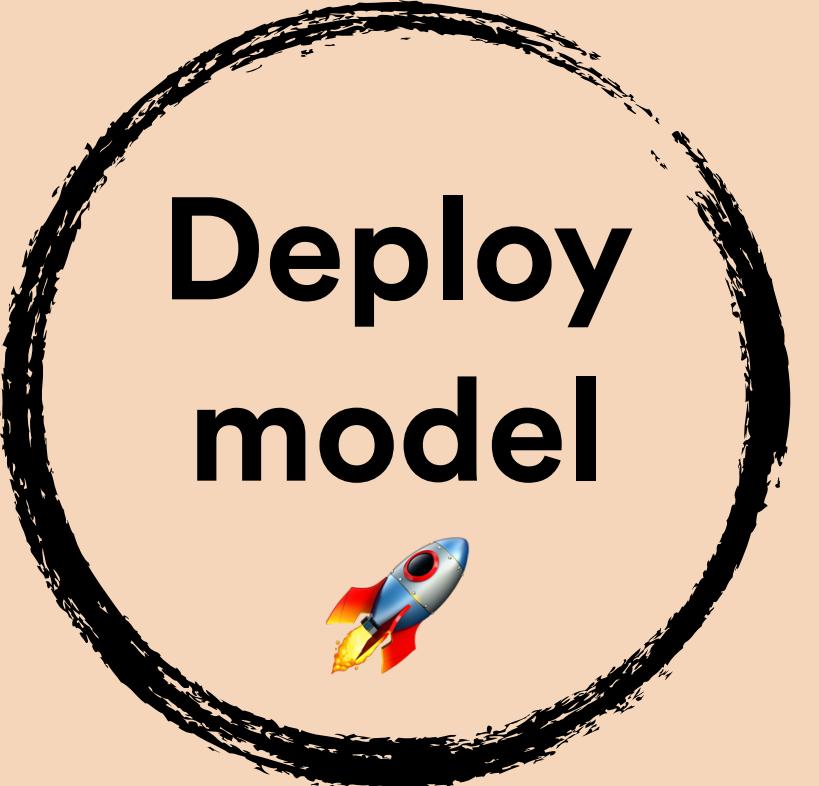
INFO:tensorflow:Step #449: rate 0.000100, accuracy 97.0%, cross entropy 0.104553
I1002 04:10:40.378780 4741027264 train.py:287] Step #449: rate 0.000100, accuracy 97.0%, cross entropy 0.104553
INFO:tensorflow:Step #450: rate 0.000100, accuracy 98.0%, cross entropy 0.079582
I1002 04:11:18.869298 4741027264 train.py:287] Step #450: rate 0.000100, accuracy 98.0%, cross entropy 0.079582
INFO:tensorflow:Confusion Matrix:
[[ 0  7  0  0  0]
 [ 0  7  0  0  0]
 [ 0  0 15  0  0]
 [ 0  0  0 33  0]
 [ 0  0  2  0 18]]
I1002 04:11:27.616060 4741027264 train.py:316] Confusion Matrix:
[[ 0  7  0  0  0]
 [ 0  7  0  0  0]
 [ 0  0 15  0  0]
 [ 0  0  0 33  0]
 [ 0  0  2  0 18]]
INFO:tensorflow:Step 450: Validation accuracy = 89.0% (N=82)
I1002 04:11:27.616290 4741027264 train.py:319] Step 450: Validation accuracy = 89.0% (N=82)
INFO:tensorflow:Saving to "./training/conv.ckpt-450"
I1002 04:11:27.616428 4741027264 train.py:331] Saving to "./training/conv.ckpt-450"
INFO:tensorflow:set_size=79
I1002 04:11:27.838648 4741027264 train.py:336] set_size=79
WARNING:tensorflow:Confusion Matrix:
[[ 0  7  0  0  0]
 [ 0  7  0  0  0]
 [ 0  0 10  0  0]
 [ 0  0  0 35  0]
 [ 0  0  0  1 19]]
W1002 04:11:36.138626 4741027264 train.py:357] Confusion Matrix:
[[ 0  7  0  0  0]
 [ 0  7  0  0  0]
 [ 0  0 10  0  0]
 [ 0  0  0 35  0]
 [ 0  0  0  1 19]]
WARNING:tensorflow:Final test accuracy = 89.9% (N=79)
W1002 04:11:36.139487 4741027264 train.py:359] Final test accuracy = 89.9% (N=79)
```

babblefish on ↵ master [!?] via venv took 4h 59m 49s

0 ➤ 0 > Python - ➤ 1 > zsh ➤ 2 > [tmux] *

2019-10-02 < 09:41 Andrews-MacBook-P

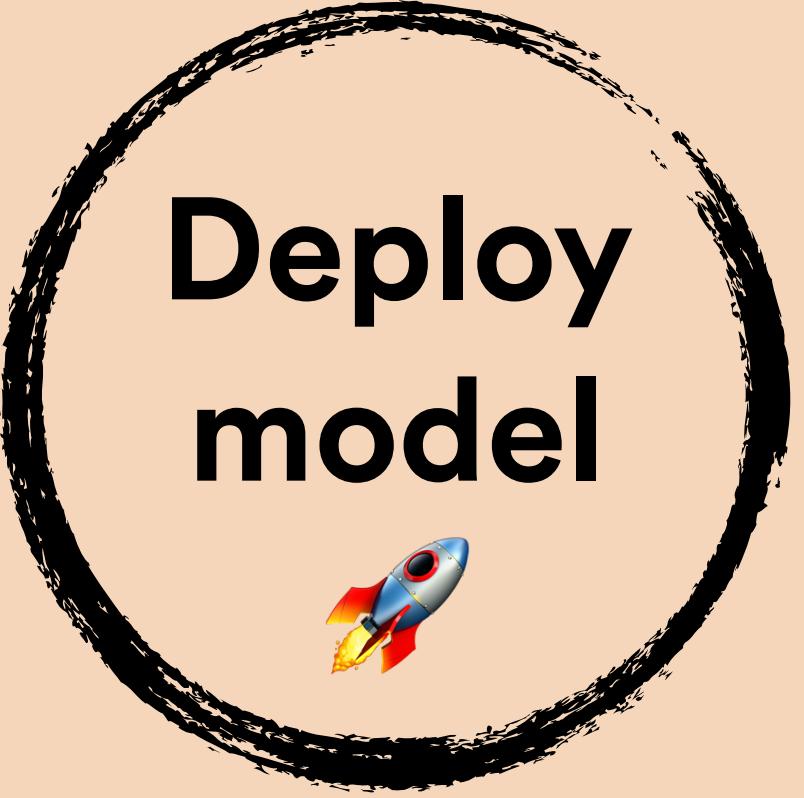




Deploy
model

Freeze the model

```
$ python app/freeze.py  
--start_checkpoint=./training/conv.ckpt-450  
--output_file=./graph.pb  
--clip_duration_ms=10000  
--sample_rate=22050  
--wanted_words=white_noise,room_empty,crying  
--data_dir=./data  
  
$ cp training/conv_labels.txt .
```



Deploy
model

Execute the model in prod

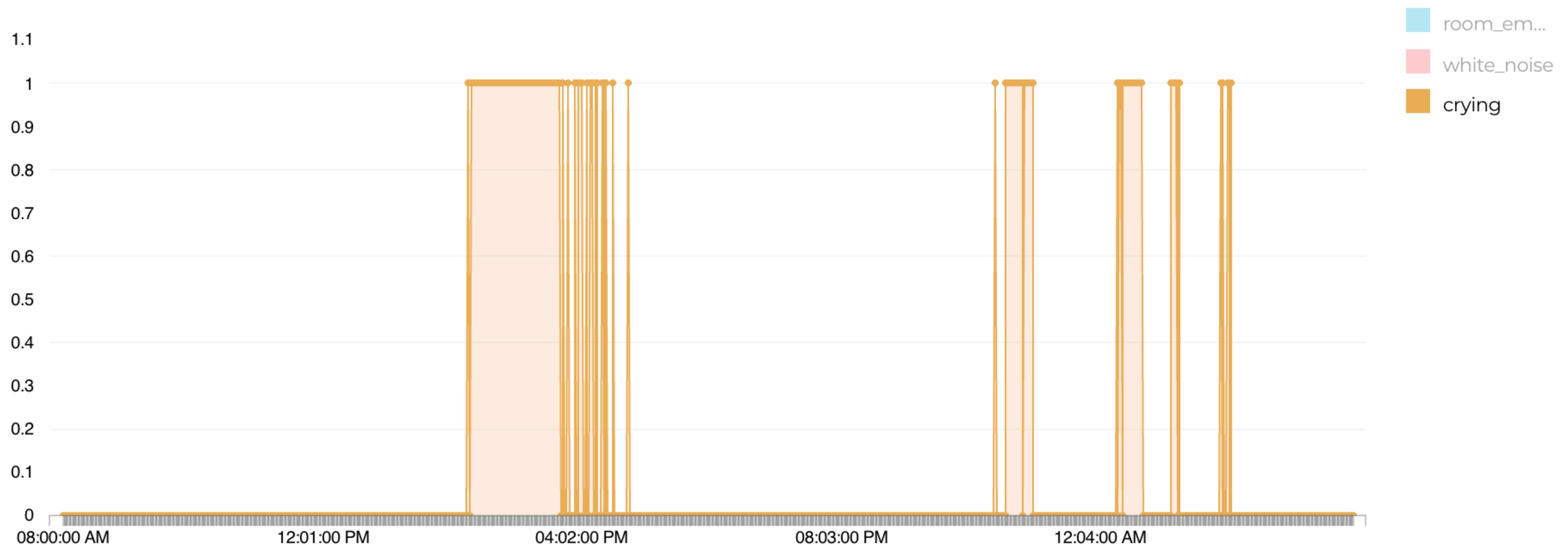
```
$ arecord --format S16_LE --rate 22050 -c1  
-d 10 $wav  
$ python app/label_wav.py  
  --graph=./graph.pb  
  --labels=./conv_labels.txt  
  --wav=$wav
```



```
> crying (score = 0.95350)  
room_empty (score = 0.01888)  
_silence_ (score = 0.01746)
```

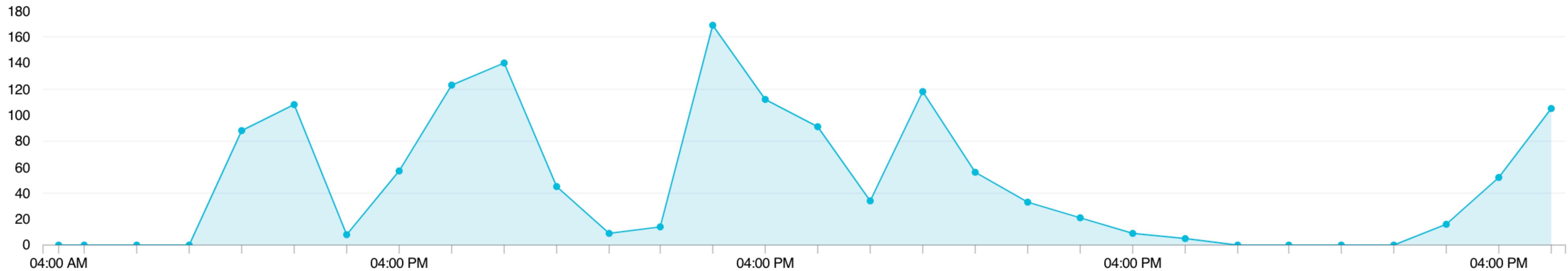
Will it work?

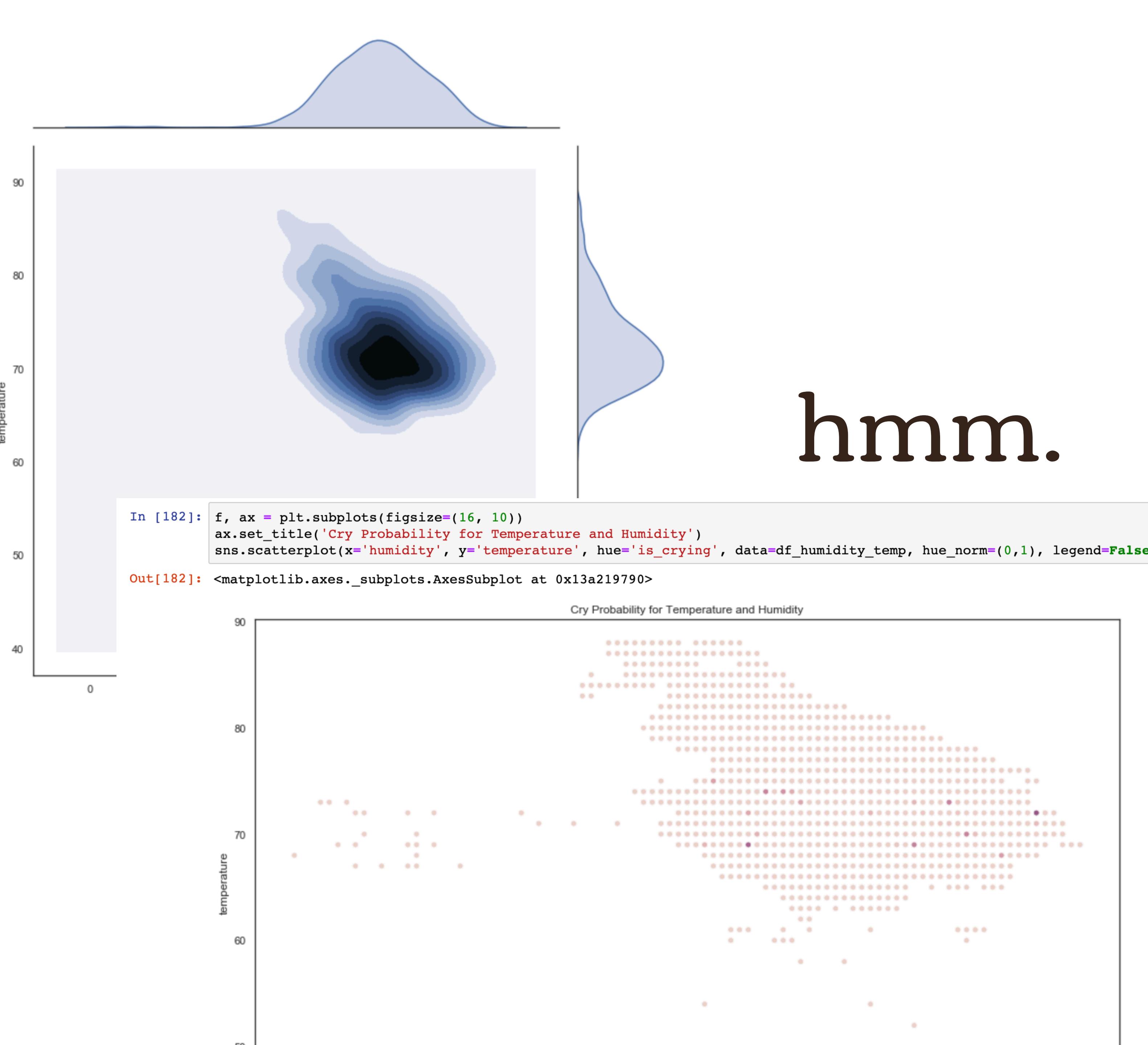
Yes!

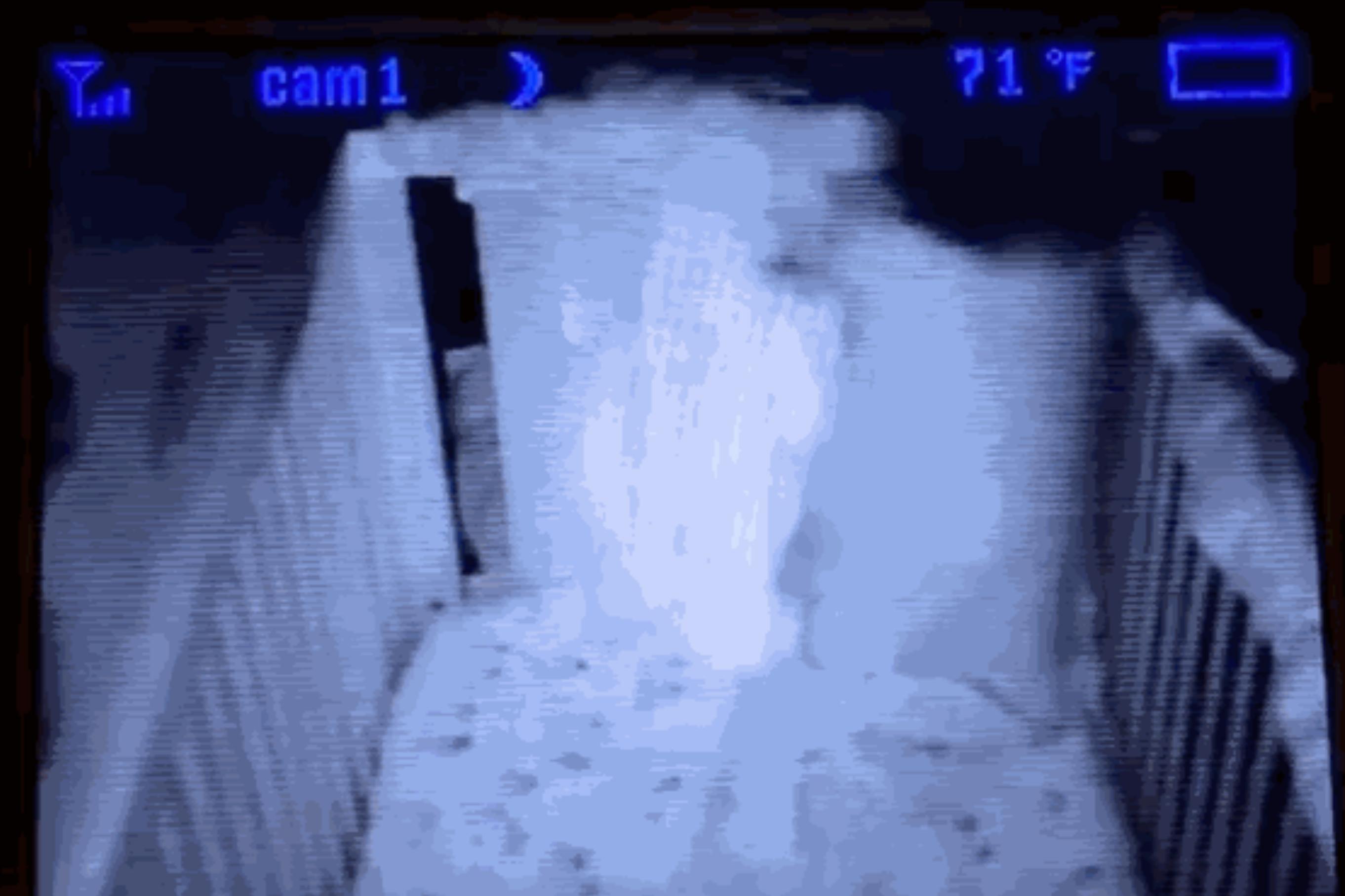


..yes?

Daily Cry Minutes (30-Day)



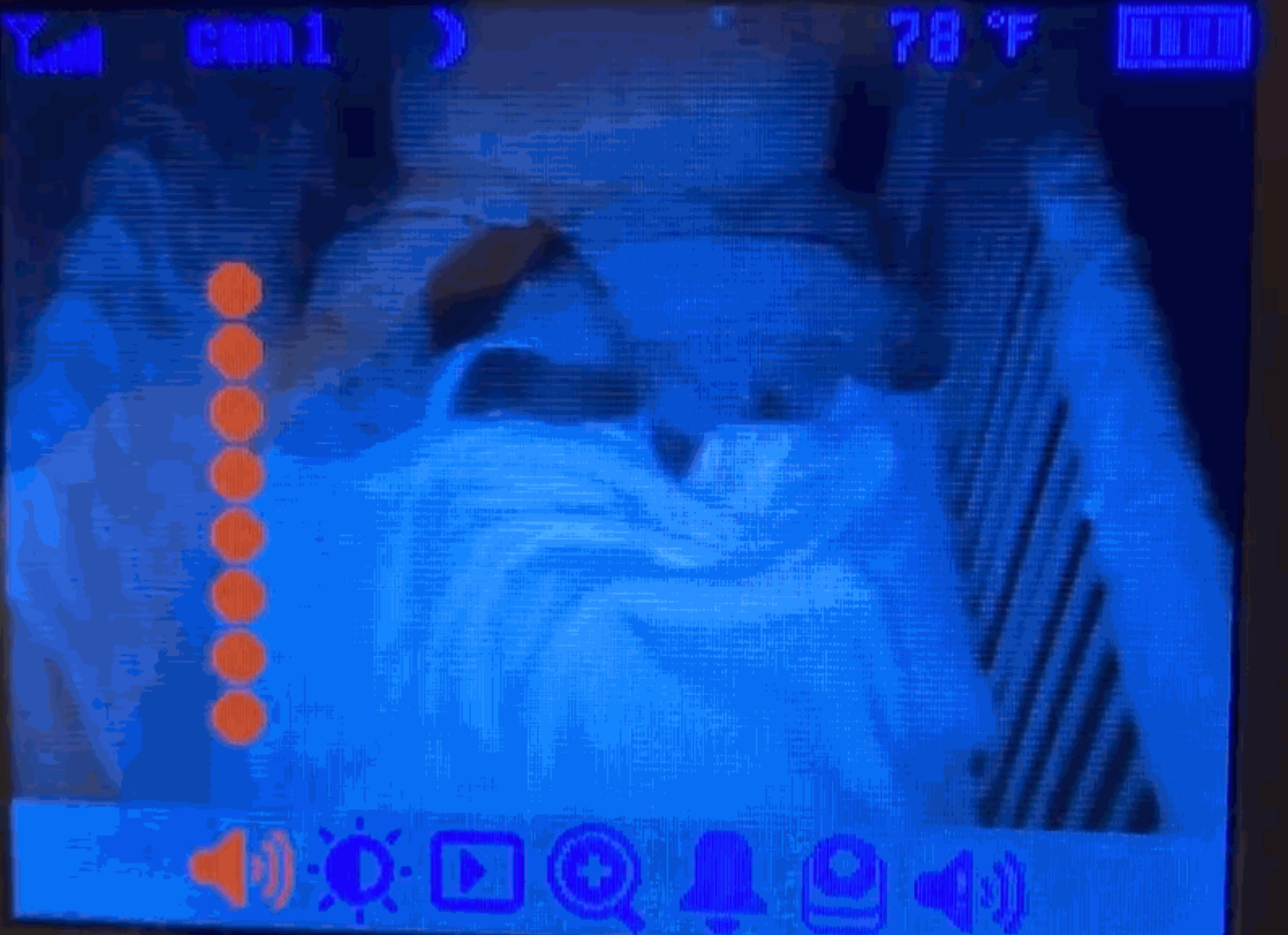




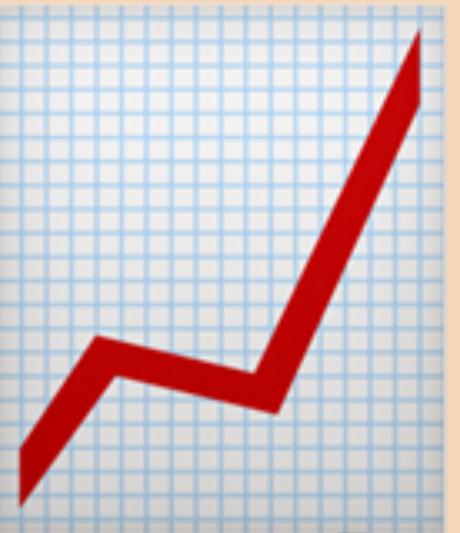
Act III

The Illusion of Insight

We just looked at each other



**Systems,
Models,
Metrics,
Funnels**

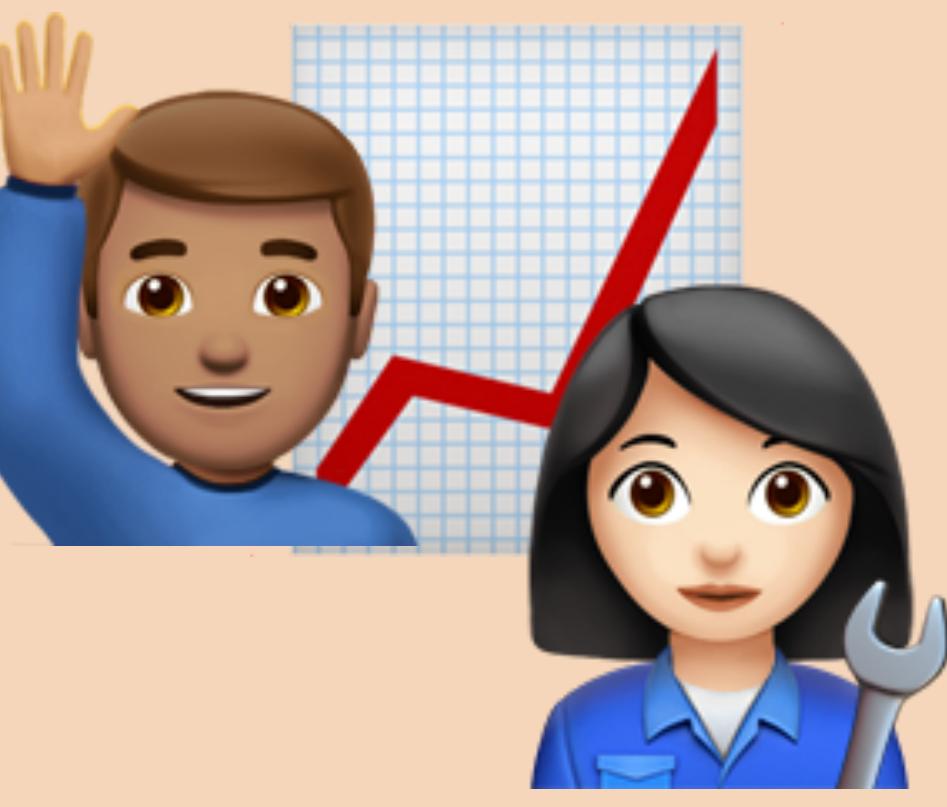


**When we lose insights about the humans we serve,
It's ultimately bad for business**

In a Big Data world,
Let's re-emphasize the human
experience



**Systems,
Models,
Metrics,
Funnels**



**Are our assumptions right?
How do we find out?**

Human-centered
UX research
Anthropology
Lean Startup: Customer interviews

Human-centered ML?



Jenna

- Customer
- Recent college graduate
- Aspiring indie electronica producer
- Wants to pay off college loans
- Challenges: ...

Photo by Philip Martin on [Unsplash](#)



Reeve

- Internal stakeholder
- Director, Digital Marketing
- Concerned about ad budget spend

Personas

Photo by Jhon David on [Unsplash](#)



Lydia

- Customer
- Civil engineer
- Mother of one & caretaker of her aging parents
- Wants to save for her daughter's college fund
- Challenges: ...

We need human-centered
design at all stages of the ML
lifecycle

Let's start by connecting with
people first





Thank you

github.com/andrewhao/babblefish

Andrew Hao
@andrewhao
ahtao@lyft.com