

# CLASSIFICATION OF ENVIRONMENTAL SOUND USING IOT SENSORS

Jon Nordby [jon@soundsensing.no](mailto:jon@soundsensing.no)

November 19, 2019

# INTRODUCTION

# JON NORDBY

Internet of Things specialist

- B.Eng in **Electronics**
- 9 years as **Software** developer. **Embedded + Web**
- M. Sc in **Data Science**

Now:

- CTO at Soundsensing
- Machine Learning Consultant



sound sensing

## What we do

---

### Noise Monitoring for outdoor and indoor environments

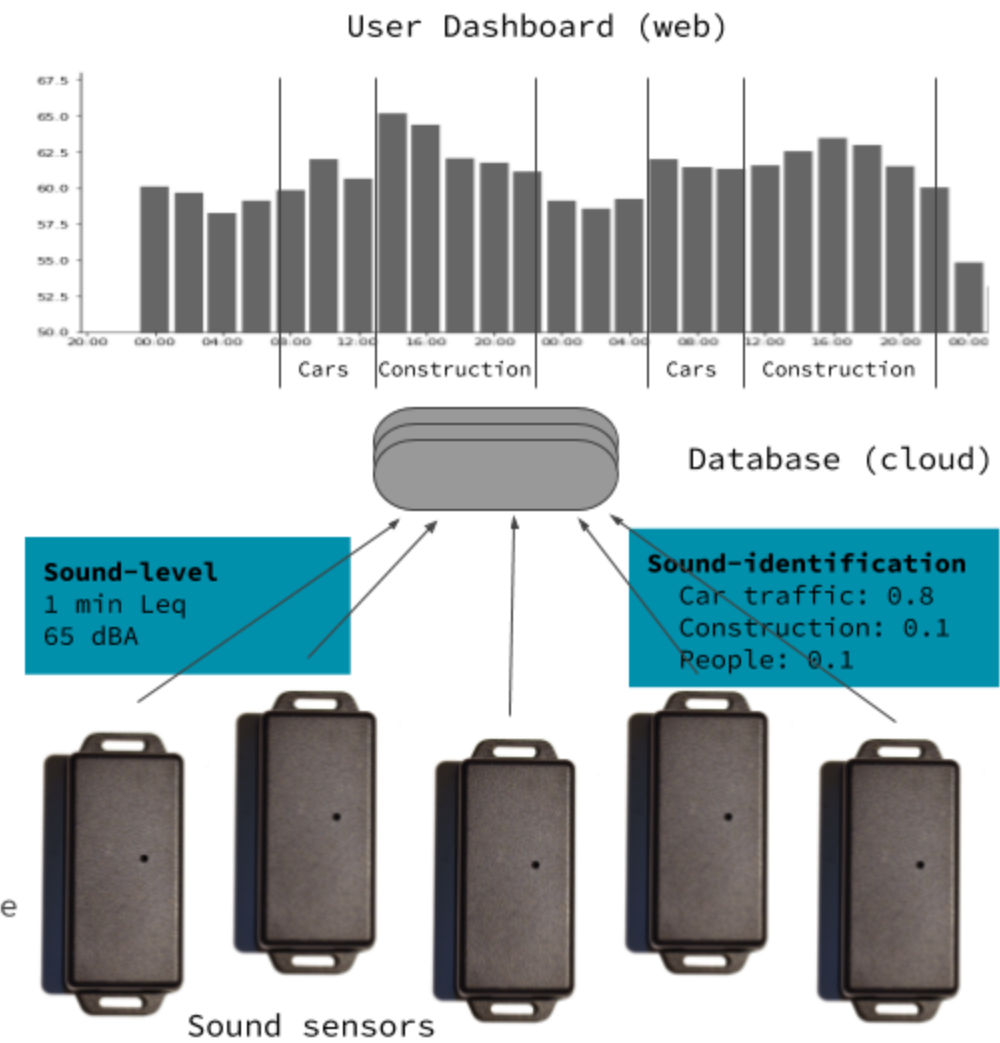
- Smart cities
- Workspaces/offices
- Hotels/AirBnB
- Music venues

#### Designed for Privacy:

Sound is not transmitted or stored.  
- only the noise information

#### Our R&D focus

How to provide the best information possible about noise using IoT sensors and machine learning



Pilot projects with customers Now - 2020

# THESIS

*Environmental Sound Classification on Microcontrollers using Convolutional Neural Networks*



Master's Thesis 2019 30 ECTS  
Faculty of Science and Technology

**Environmental Sound Classification  
on Microcontrollers  
using Convolutional Neural Networks**

Jon Nordby  
Master of Science in Data Science

Report & Code: <https://github.com/jonnor/ESC-CNN-microcontroller>

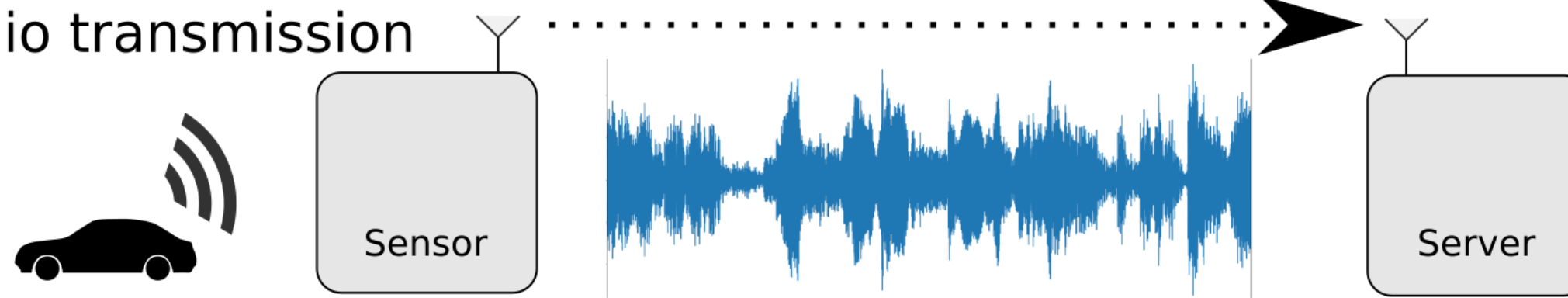
# WIRELESS SENSOR NETWORKS

- Want: Wide and dense coverage
- Need: Sensors need to be low-cost
- **Opportunity:** Wireless reduces costs
- **Challenge:** Power consumption



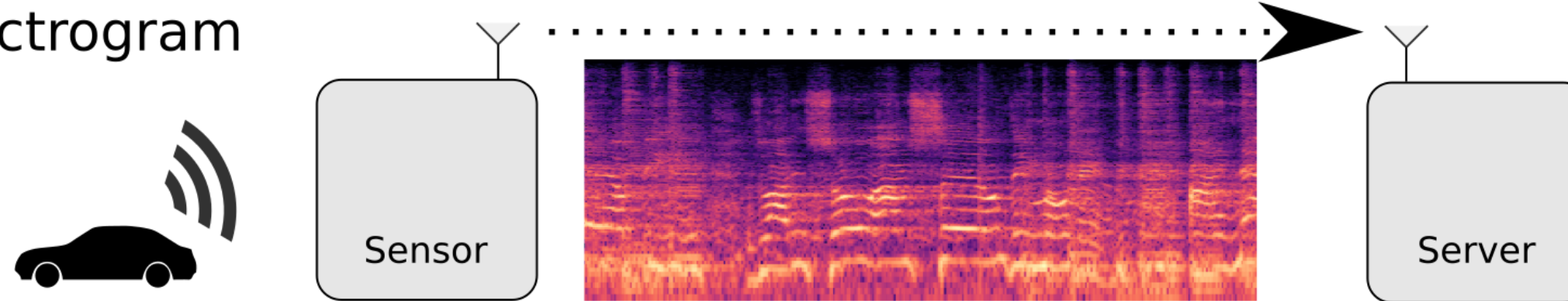
# SENSOR NETWORK ARCHITECTURES

A) Audio transmission



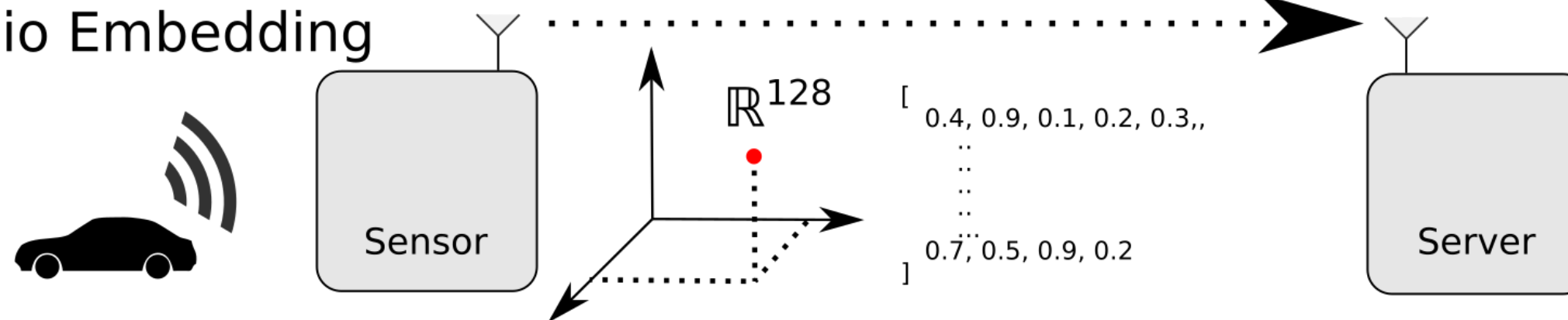
Airconditioner	0.11
Engine idling	0.05
<b>Car horn</b>	<b>0.88</b>
Children playing	0.22
Dog barking	0.12
Siren	0.09
Street Music	0.30
Drilling	0.07
Jackhammer	0.04

B) Spectrogram



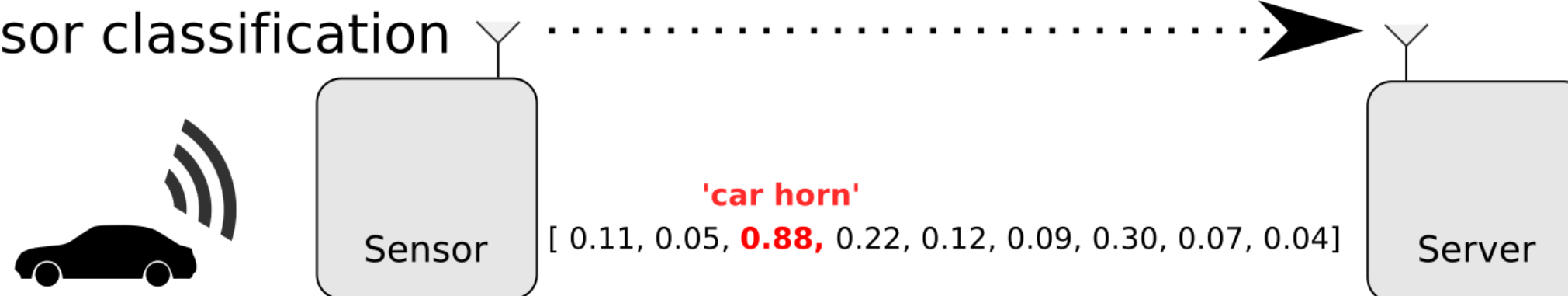
Airconditioner	0.11
Engine idling	0.05
<b>Car horn</b>	<b>0.88</b>
Children playing	0.22
Dog barking	0.12
Siren	0.09
Street Music	0.30
Drilling	0.07
Jackhammer	0.04

C) Audio Embedding



Airconditioner	0.11
Engine idling	0.05
<b>Car horn</b>	<b>0.88</b>
Children playing	0.22
Dog barking	0.12
Siren	0.09
Street Music	0.30
Drilling	0.07
Jackhammer	0.04

D) Sensor classification



Airconditioner	0.11
Engine idling	0.05
<b>Car horn</b>	<b>0.88</b>
Children playing	0.22
Dog barking	0.12
Siren	0.09
Street Music	0.30
Drilling	0.07
Jackhammer	0.04



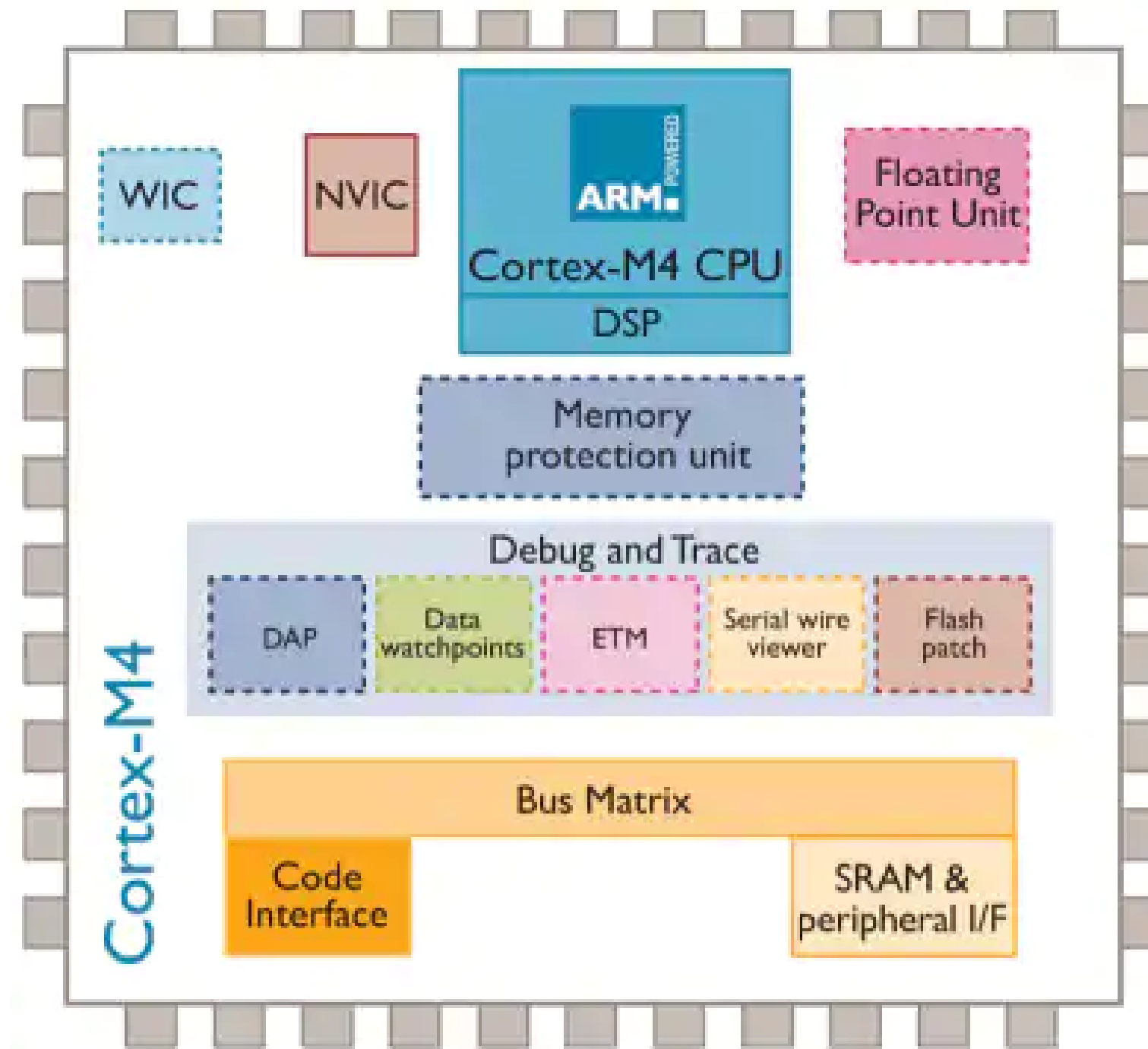
# AUDIO MACHINE LEARNING ON LOW- POWER SENSORS

# WHAT DO YOU MEAN BY LOW-POWER?

Want: 1 year lifetime for palm-sized battery

Need:  $<1\text{mW}$  system power

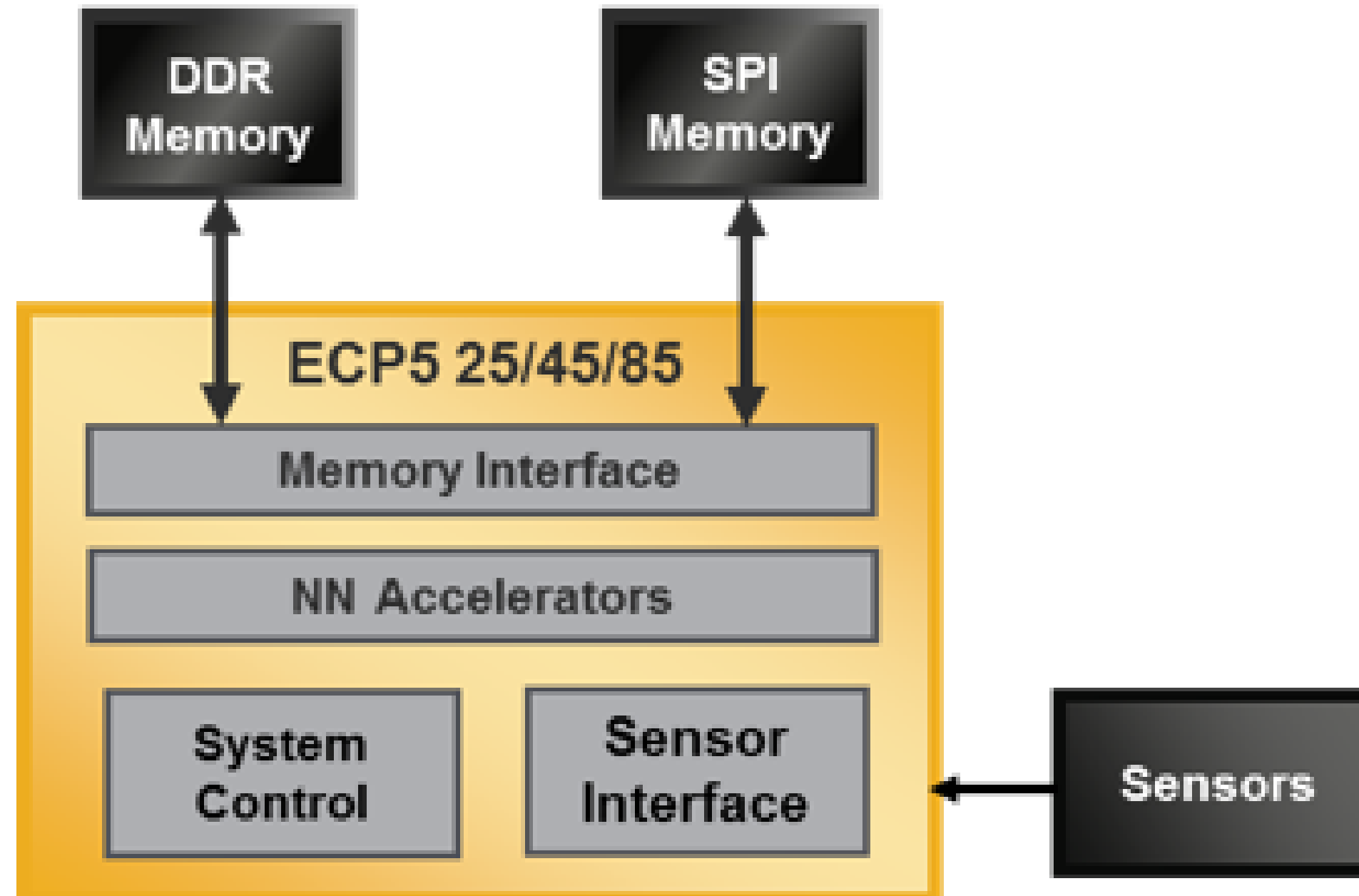
# GENERAL PURPOSE MICROCONTROLLER



STM32L4 @ 80 MHz. Approx **10 mW**.

- TensorFlow Lite for Microcontrollers (Google)
- ST X-CUBE-AI (ST Microelectronics)

# FPGA

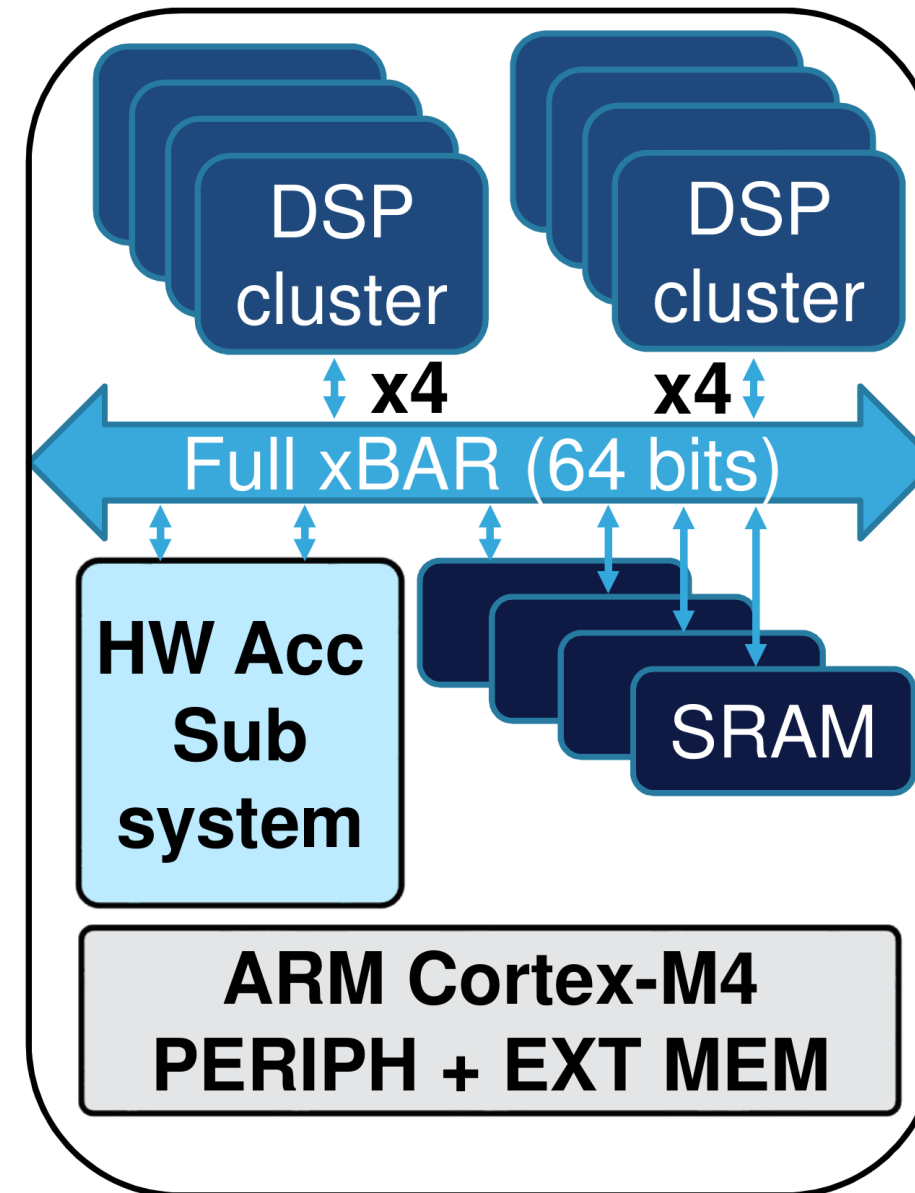


Lattice ICE40 UltraPlus with Lattice sensAI

Human presence detection. VGG8 on 64x64 RGB image, 5 FPS: 7 mW.

Audio ML approx **1 mW**

# NEURAL NETWORK CO-PROCESSORS



Project Orlando (ST Microelectronics), expected 2020  
2.9 TOPS/W. AlexNet, 1000 classes, 10 FPS. 41 mWatt  
Audio models probably < 1 mWatt.

# ON-EDGE CLASSIFICATION OF NOISE

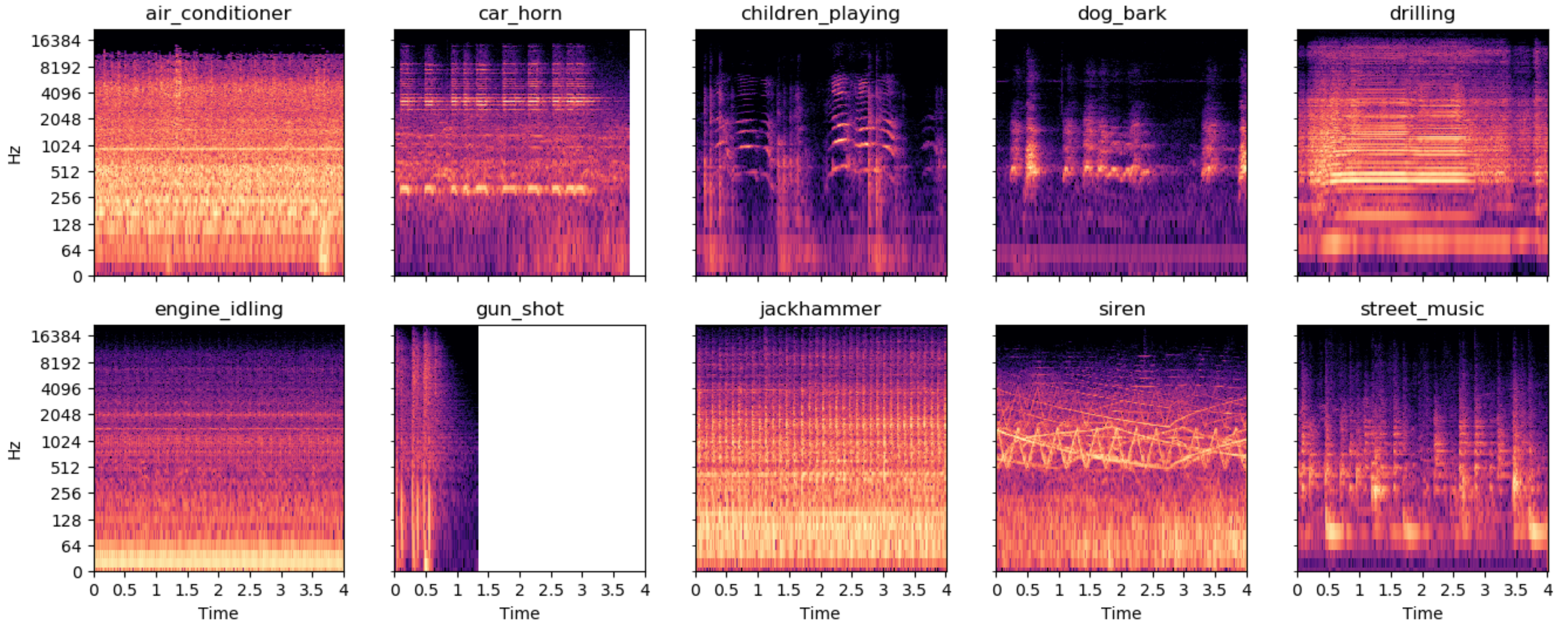


# ENVIRONMENTAL SOUND CLASSIFICATION

*Given an audio signal of environmental sounds,  
determine which class it belongs to*

- Widely researched. 1000 hits on Google Scholar
- Datasets. Urbansound8k (10 classes), ESC-50, AudioSet (632 classes)
- 2017: Human-level performance on ESC-50

# URBANSOUND8K



# EXISTING WORK

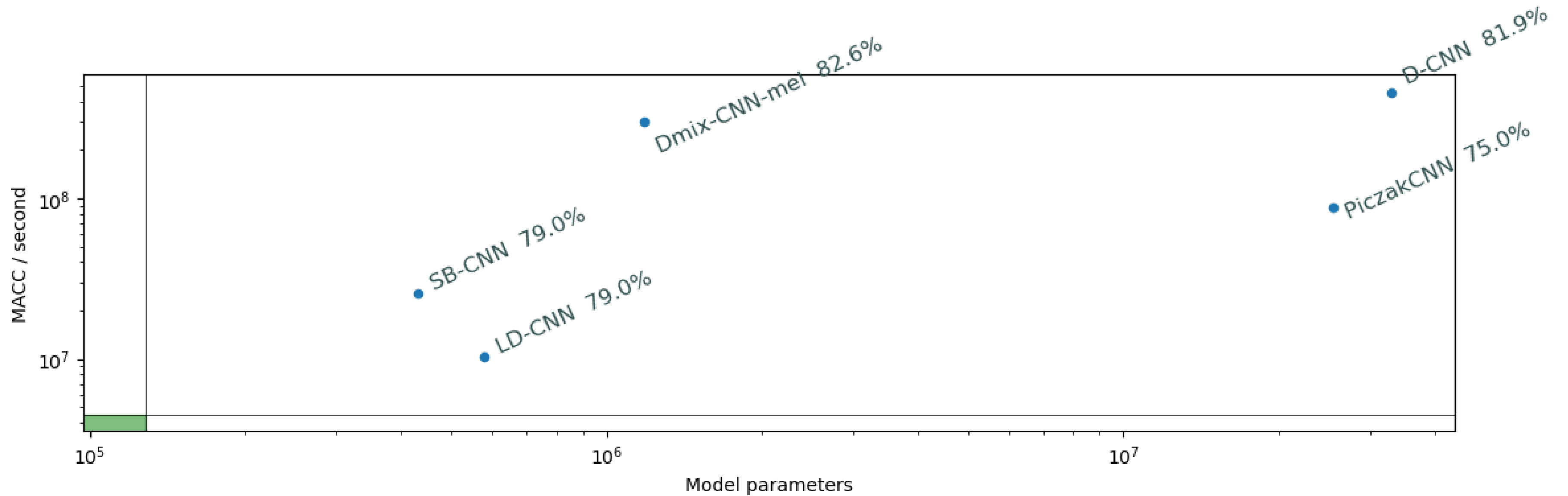
- Convolutional Neural Networks dominate
- Techniques come from image classification
- Mel-spectrogram input standard
- End2end models: getting close in accuracy
- “Edge ML” focused on mobile-phone class HW
- “Tiny ML” (sensors) just starting

# MODEL REQUIREMENTS

With 50% of STM32L476 capacity:

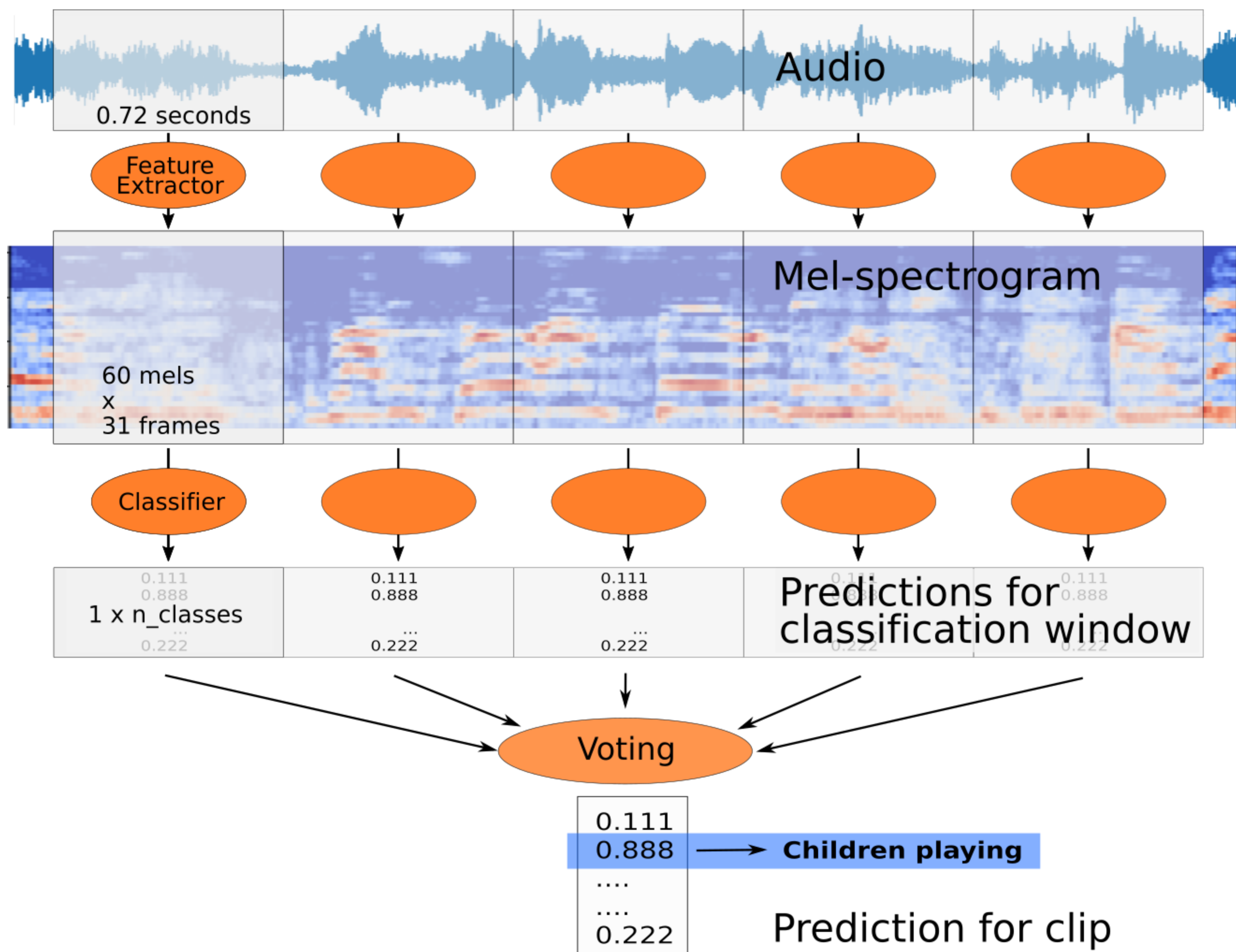
- 64 kB RAM
- 512 kB FLASH memory
- 4.5 M MACC/second

# EXISTING MODELS



Green: Feasible region

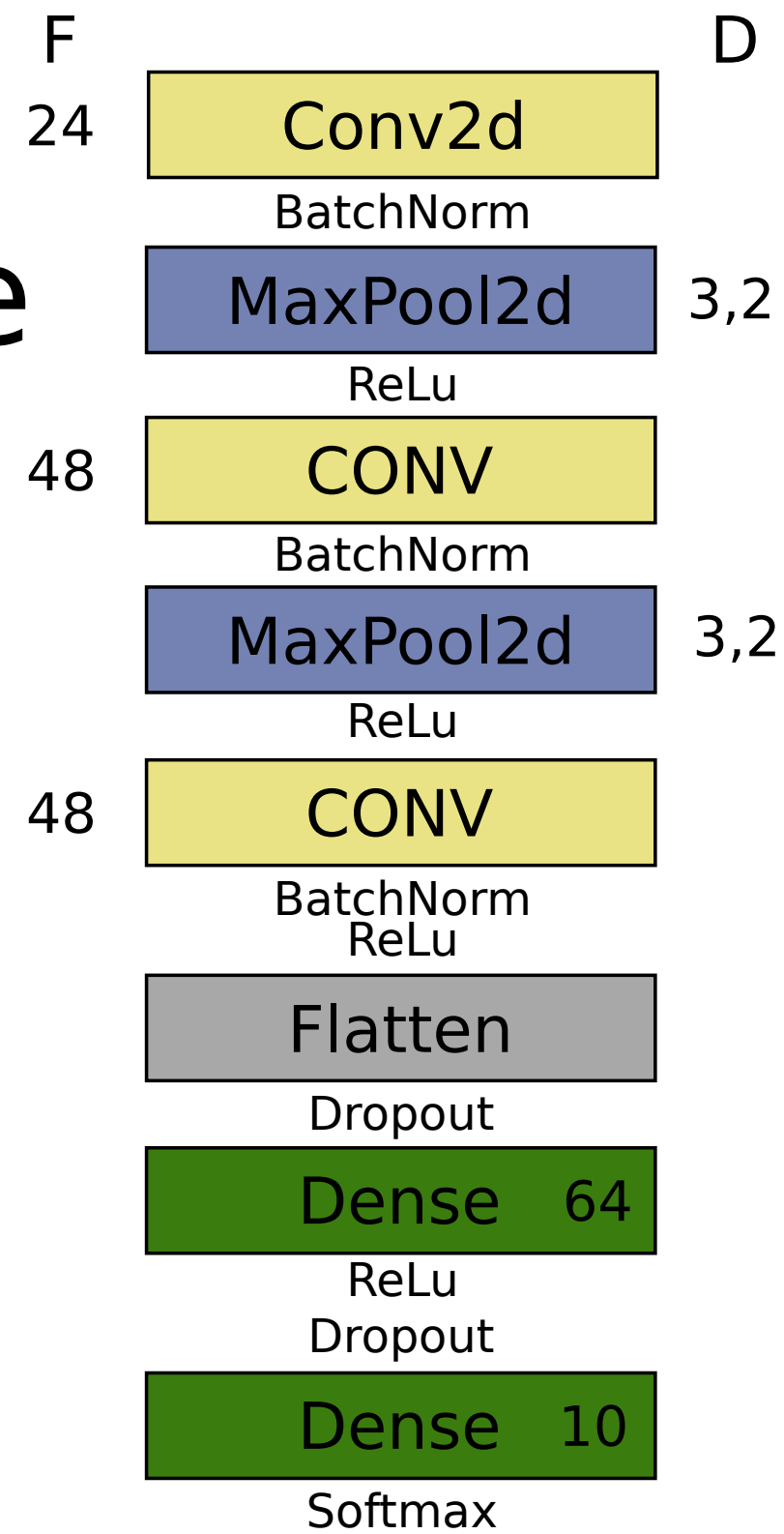
eGRU: running on ARM Cortex-M0 microcontroller, accuracy 61% with **non-standard** evaluation



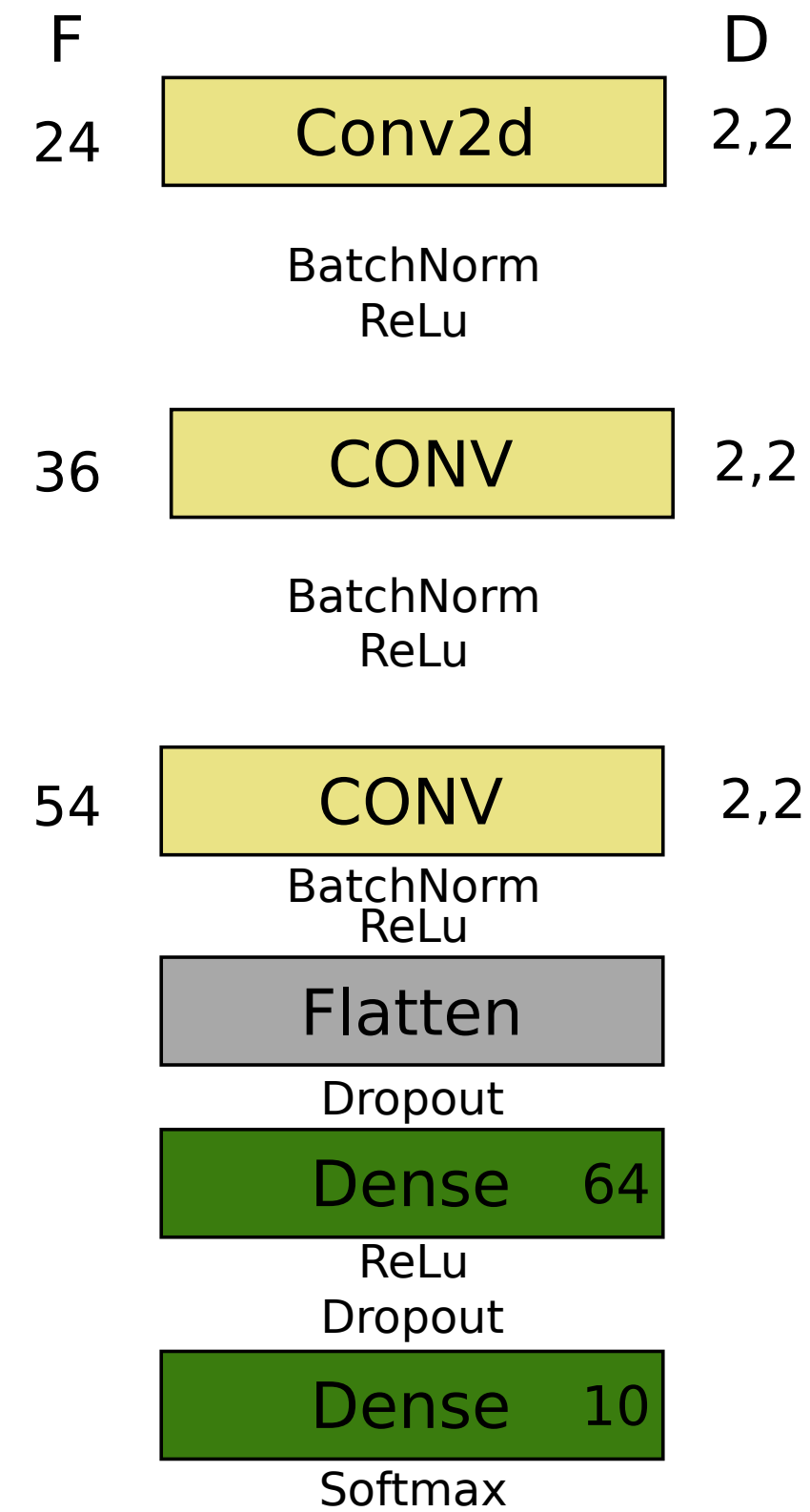


# MODELS

## Baseline



## Stride

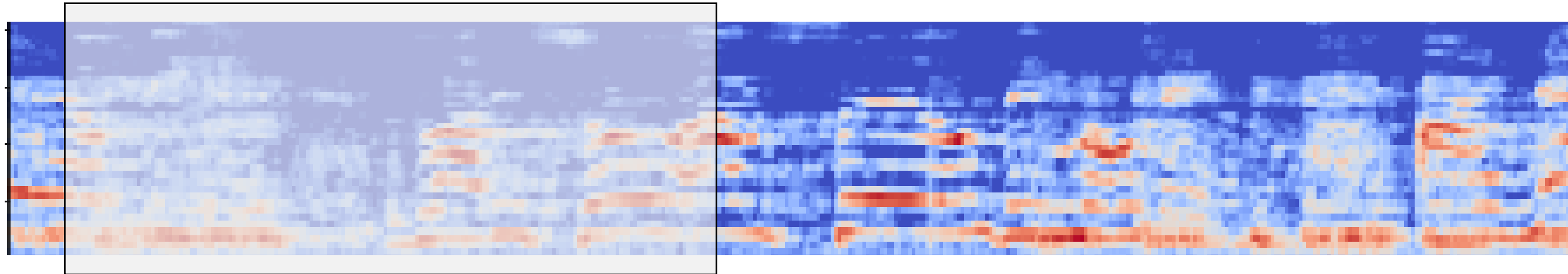




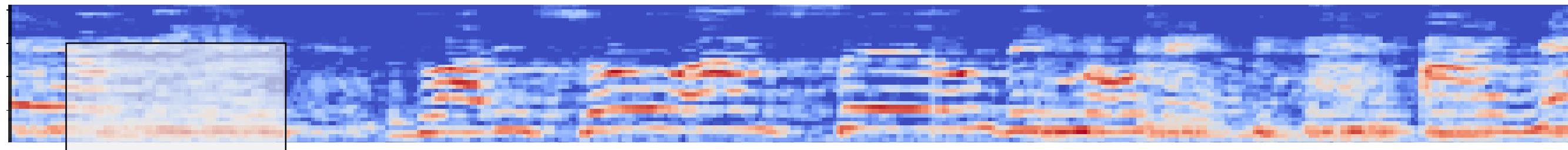
# STRATEGIES FOR SHRINKING CONVOLUTIONAL NEURAL NETWORK

# REDUCE INPUT DIMENSIONALITY

44.1kHz, 2 seconds, 128x128

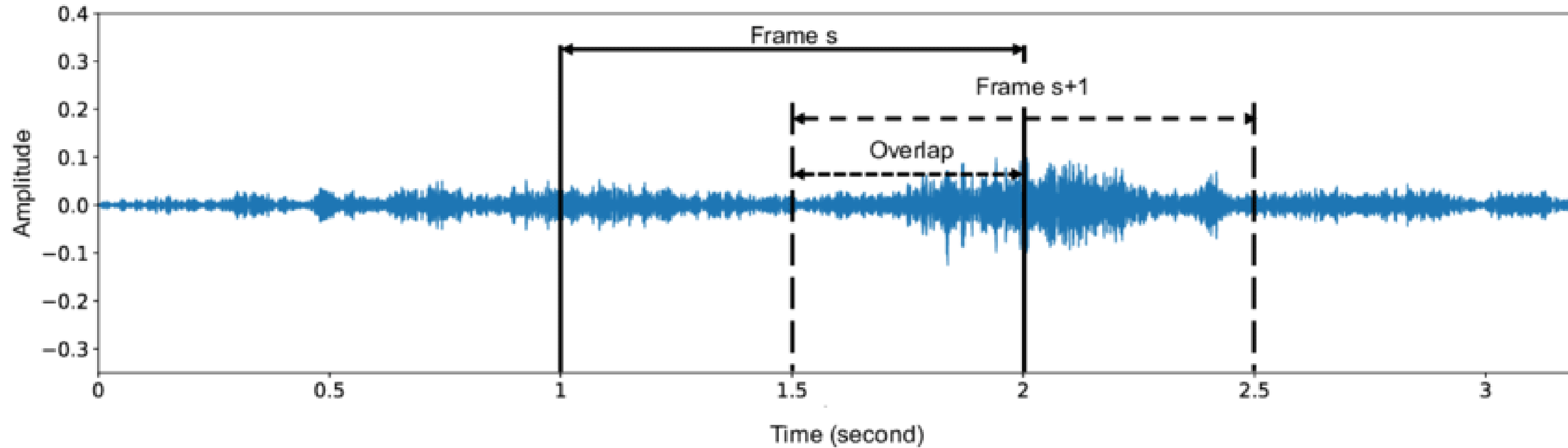


16kHz, 0.75 seconds, 32x32



- Lower frequency range
- Lower frequency resolution
- Lower time duration in window
- Lower time resolution

# REDUCE OVERLAP

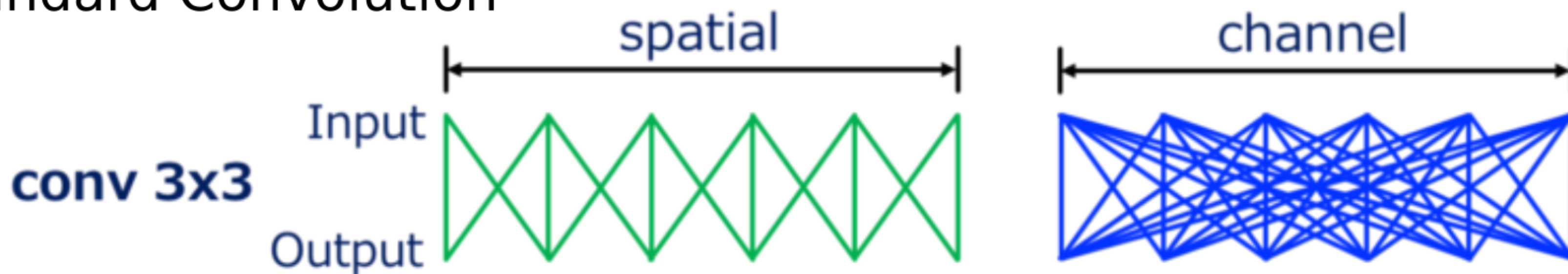


Models in literature use 95% overlap or more. 20x penalty in inference time!

Often low performance benefit. Use 0% (1x) or 50% (2x).

# DEPTHWISE-SEPARABLE CONVOLUTION

## Standard Convolution

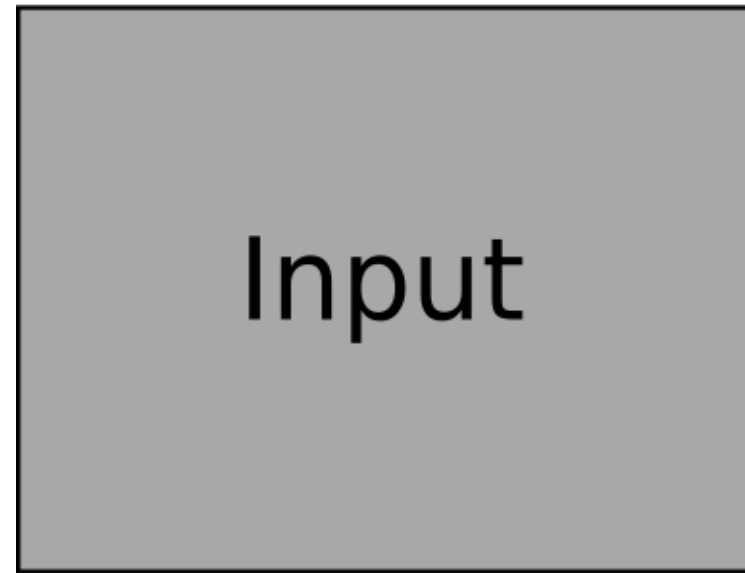


## Depthwise Separable Convolution



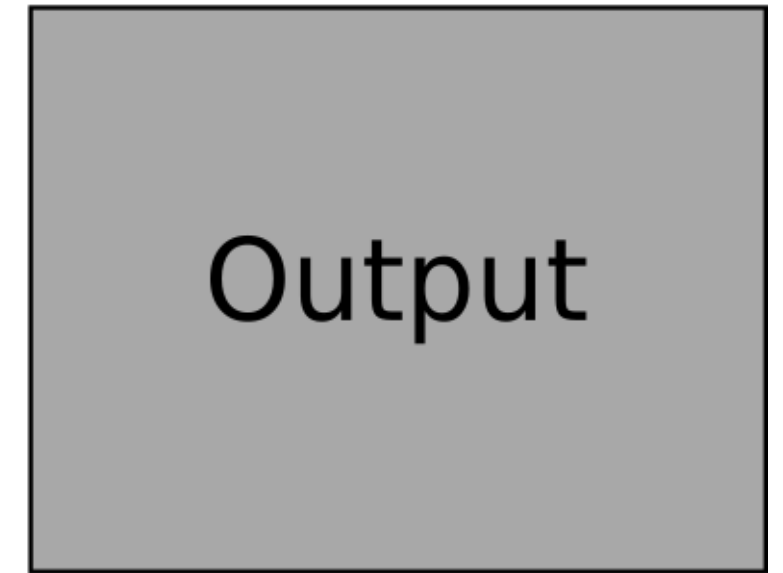
# SPATIALLY-SEPARABLE CONVOLUTION

## Standard Convolution

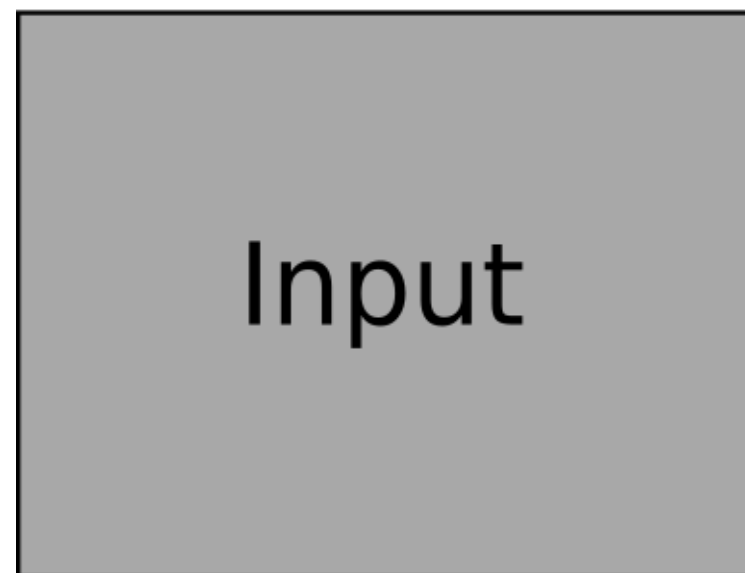


3x3  
convolution

$$\begin{bmatrix} 1 & 0 & -1 \\ 2 & 0 & -2 \\ 1 & 0 & -1 \end{bmatrix}$$

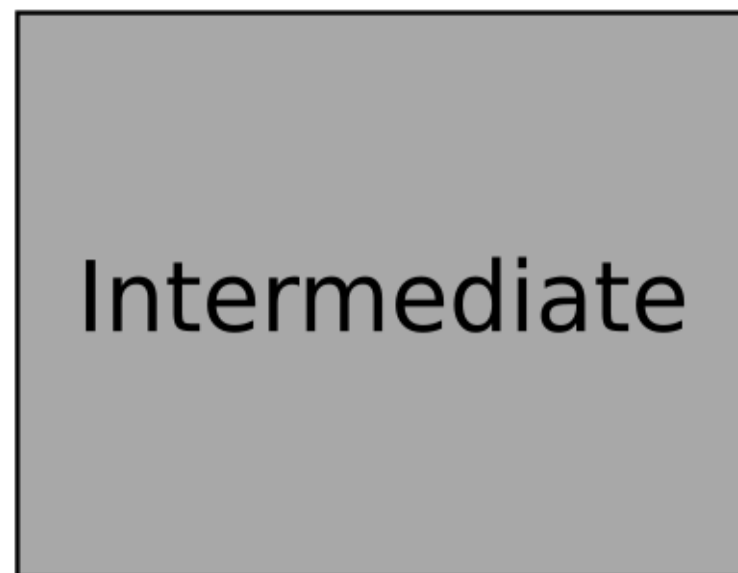


## Spatially Separable Convolution



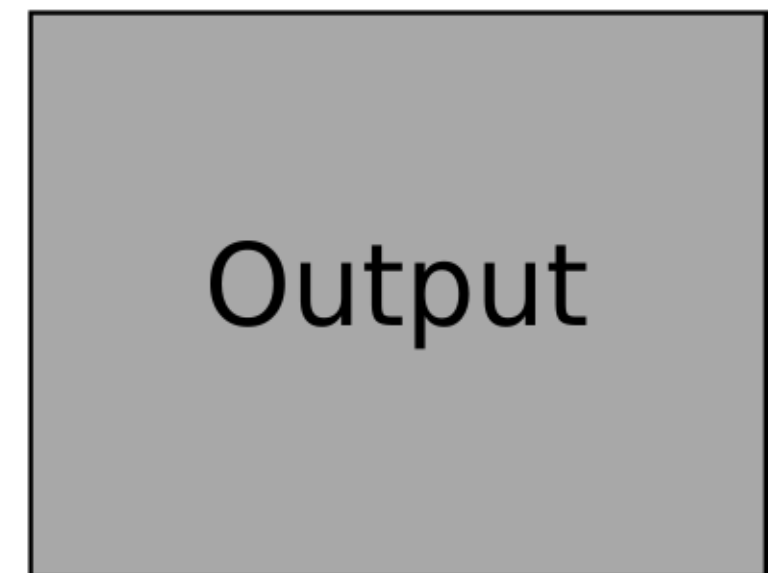
3x1  
convolution

$$\begin{bmatrix} 1 \\ 2 \\ 1 \end{bmatrix}$$



1x3  
convolution

$$[+1 \quad 0 \quad -1]$$



# Downsampling Using Max-Pooling

7	2	5	1
3	6	5	8
4	4	4	1
6	3	7	2



Maxpool  
2x2 filter  
2x2 stride

7	8
6	7

Wasteful? Computing convolutions, then throwing away 3/4 of results!

# DOWNSAMPLING USING STRIDED CONVOLUTION

stride: 2

7	7	2	5	1
7	7	2	5	1
3	3	6	5	8
4	4	4	4	1
6	6	3	7	2

5x5 input

\*

1	1	1
0	0	0
-1	-1	-1

3x3 filter

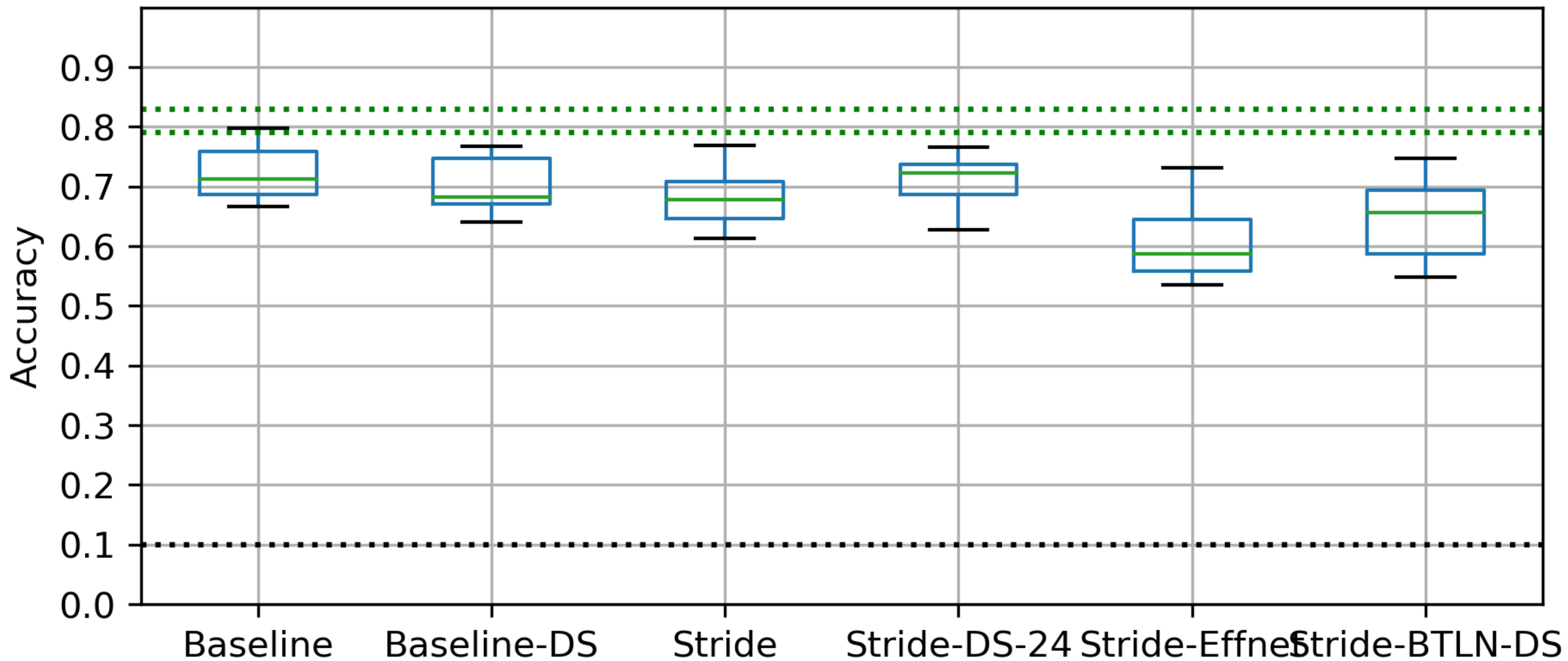
=


2x2 output

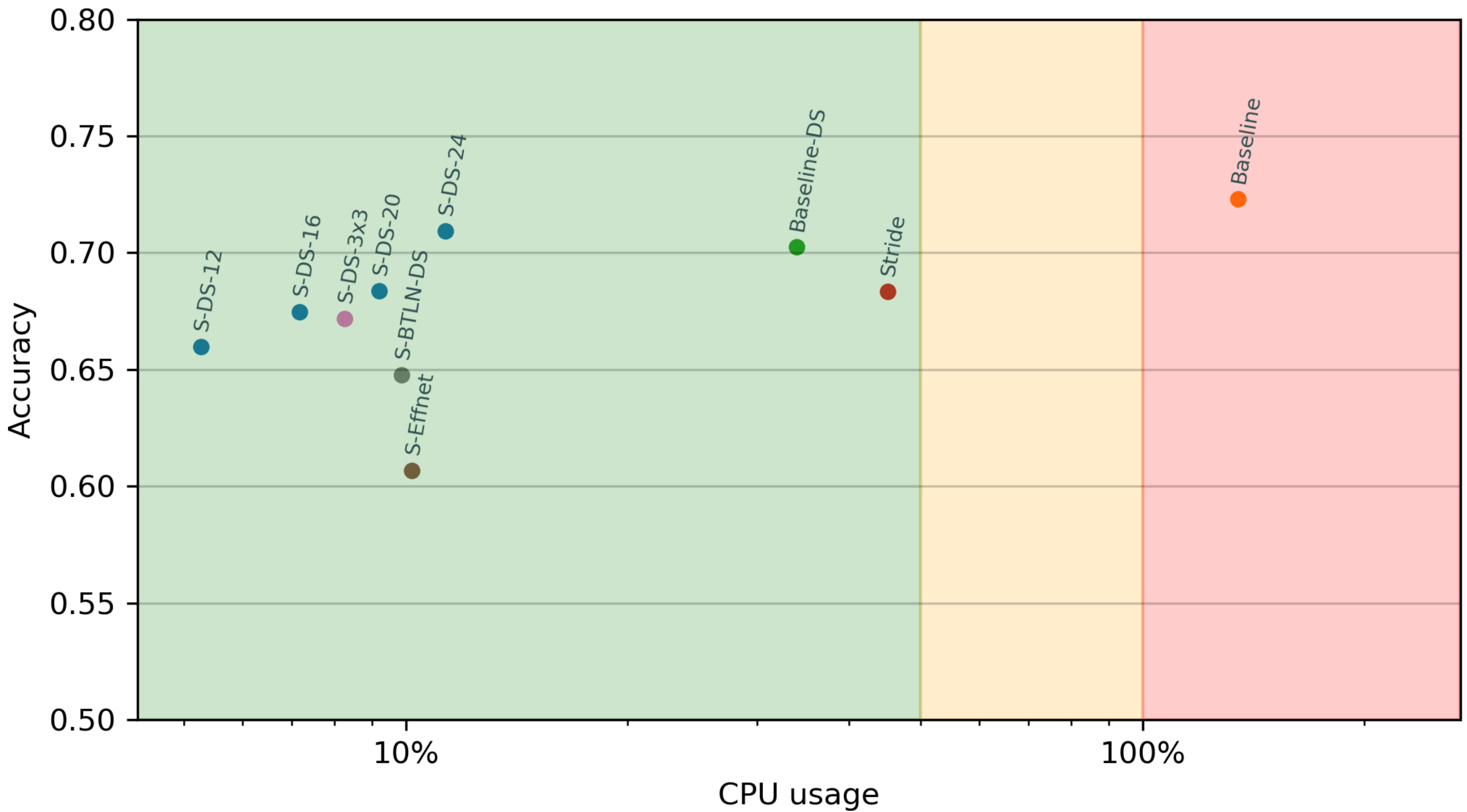




# MODEL COMPARISON



# PERFORMANCE VS COMPUTE



# QUANTIZATION

Inference can often use 8 bit integers instead of 32 bit floats

- 1/4 the size for weights (FLASH) and activations (RAM)
- 8bit **SIMD** on ARM Cortex M4F: 1/4 the inference time
- Supported in X-CUBE-AI 4.x (July 2019)

# CONCLUSIONS

- Able to perform Environmental Sound Classification at ~ 10mW power,
- Using *general purpose microcontroller*, ARM Cortex M4F
- Best performance: 70.9% mean accuracy, under 20% CPU load
- Highest reported Urbansound8k on microcontroller (over eGRU 62%)
- Best architecture: Depthwise-Separable convolutions with striding
- Quantization enables 4x bigger models (and higher perf)
- With dedicated Neural Network Hardware

# FURTHER RESEARCH

# WAVEFORM INPUT TO MODEL

- Preprocessing. Mel-spectrogram: **60** milliseconds
- CNN. Stride-DS-24: **81** milliseconds
- With quantization, spectrogram conversion is the bottleneck!
- Convolutions can be used to learn a Time-Frequency transformation.

Can this be faster than the standard FFT? And still perform well?

# ON-SENSOR INFERENCE CHALLENGES

- Reducing power consumption. Adaptive sampling
- Efficient training data collection in WSN. Active Learning?
- Real-life performance evaluations. Out-of-domain samples



# WRAPPING UP

# SUMMARY

- Noise pollution is a growing problem
- Wireless Sensor Networks can be used to quantify
- Noise Classification can provide more information
- Want high density of sensors. Need to be low cost
- On-sensor classification desirable for power/cost and privacy

# MORE RESOURCES

Machine Hearing. ML on Audio

- [github.com/jonnor/machinehearing](https://github.com/jonnor/machinehearing)

Machine Learning for Embedded / IoT

- [github.com/jonnor/embeddedml](https://github.com/jonnor/embeddedml)

Thesis Report & Code

- [github.com/jonnor/ESC-CNN-microcontroller](https://github.com/jonnor/ESC-CNN-microcontroller)

# QUESTIONS

?

Email: [jon@soundsensing.no](mailto:jon@soundsensing.no)

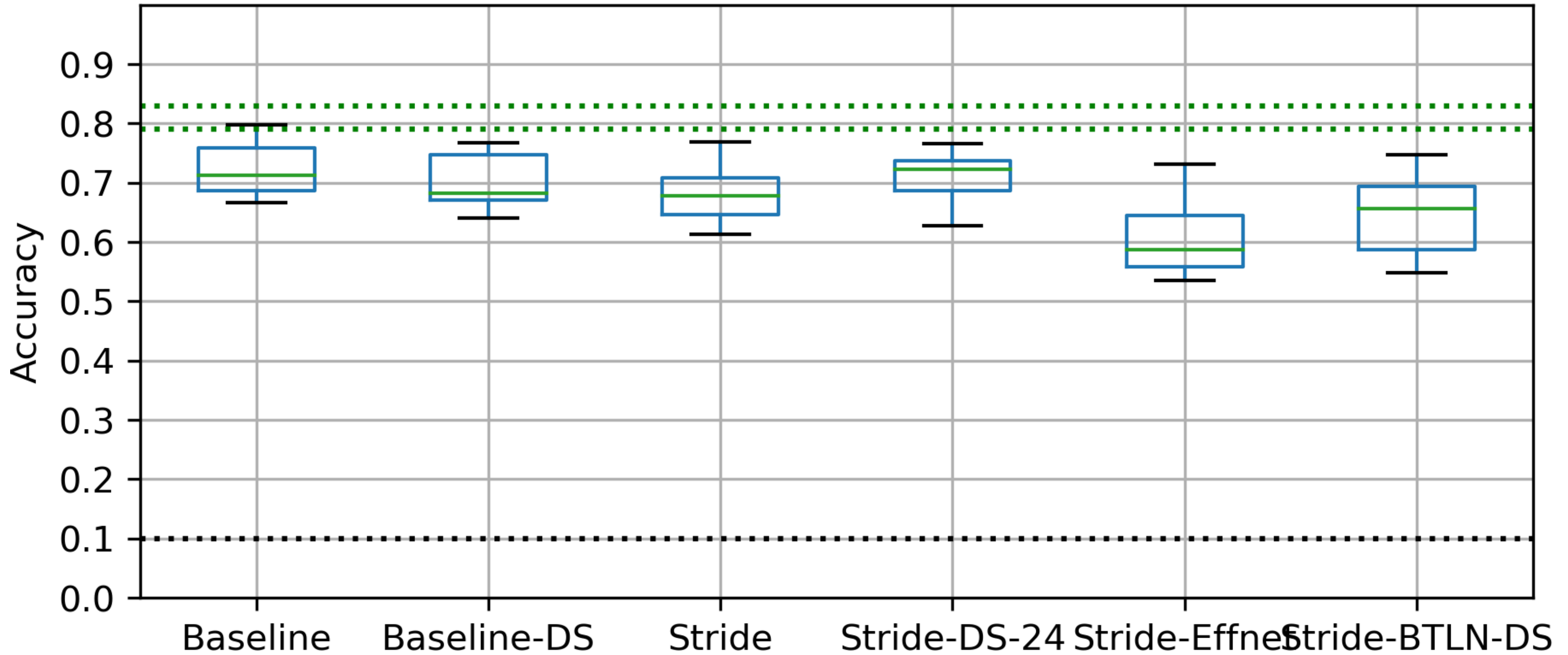
# COME TALK TO ME!

- Noise Monitoring sensors. Pilot projects for 2020?
- Environmental Sound, Wireless Sensor Networks for Audio. Research partnering?
- “On-edge” / Embedded Device ML. Happy to advise!

Email: [jon@soundsensing.no](mailto:jon@soundsensing.no)

# THESIS RESULTS

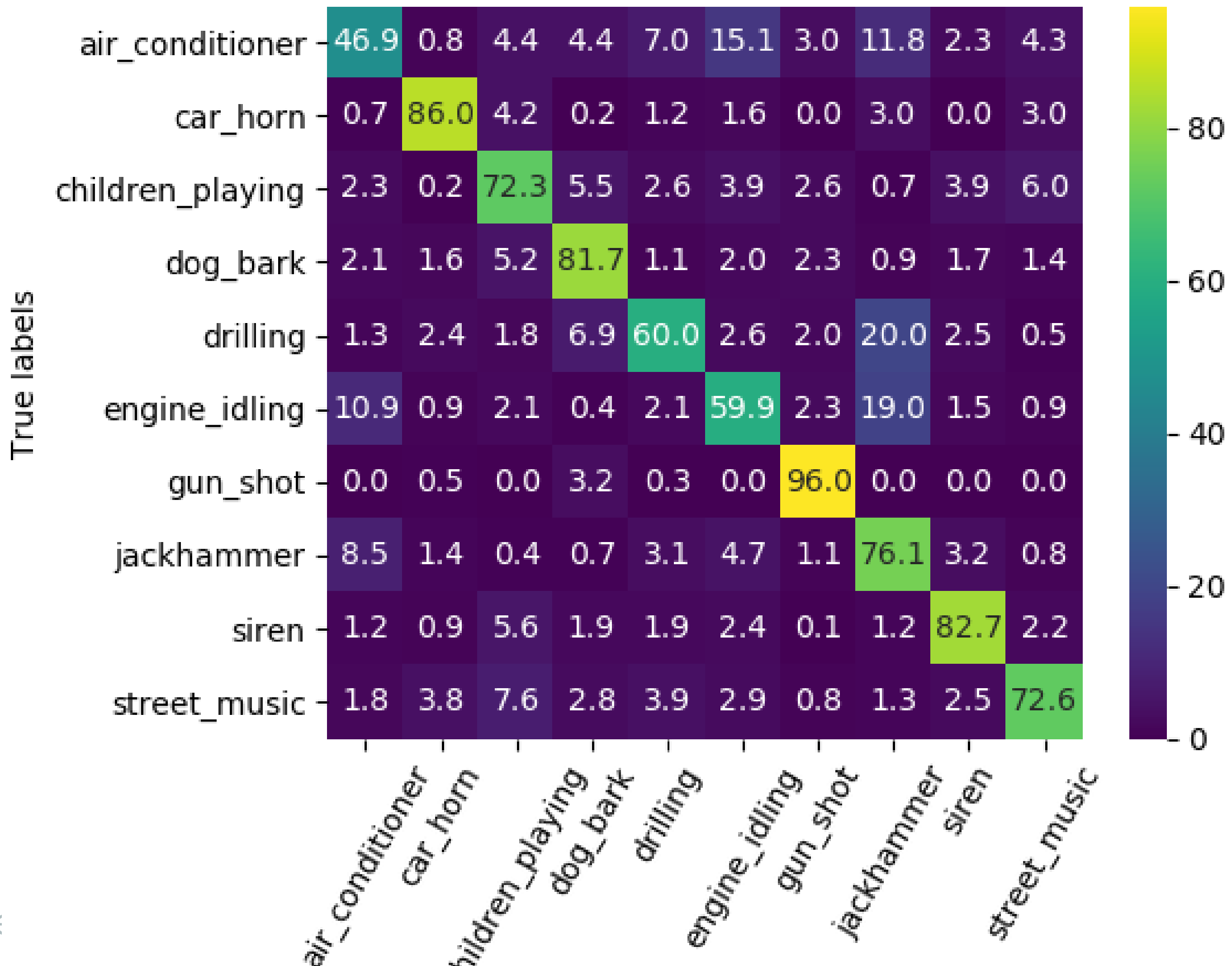
# MODEL COMPARISON



# LIST OF RESULTS

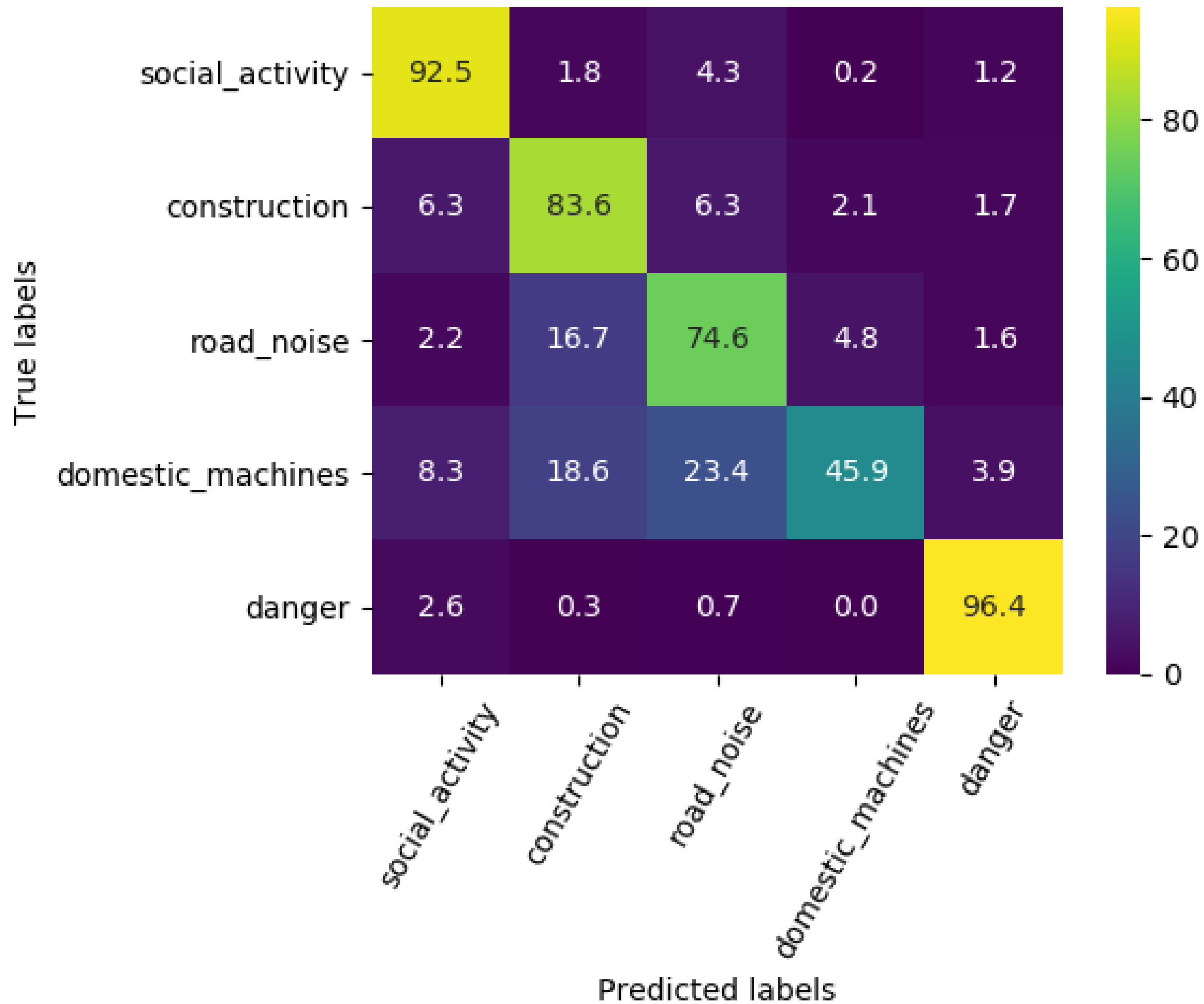
Model	CPU use	Accuracy	FG Accuracy	BG Accuracy
Baseline	971 ms	72.3% $\pm$ 4.6	78.3% $\pm$ 7.1	60.5% $\pm$ 7.7
Baseline-DS	244 ms	70.2% $\pm$ 4.7	76.1% $\pm$ 7.5	58.6% $\pm$ 8.2
Stride	325 ms	68.3% $\pm$ 5.2	74.1% $\pm$ 6.6	56.6% $\pm$ 8.0
Stride-BTLN-DS	71 ms	64.8% $\pm$ 7.1	69.5% $\pm$ 8.2	55.3% $\pm$ 8.9
Stride-DS-12	38 ms	66.0% $\pm$ 6.0	72.6% $\pm$ 6.5	53.3% $\pm$ 9.1
Stride-DS-16	51 ms	67.5% $\pm$ 5.6	73.3% $\pm$ 7.7	56.2% $\pm$ 8.3
Stride-DS-20	66 ms	68.4% $\pm$ 5.2	75.0% $\pm$ 7.4	55.2% $\pm$ 10.0
Stride-DS-24	81 ms	70.9% $\pm$ 4.3	75.8% $\pm$ 6.3	61.8% $\pm$ 6.8
Stride-DS-3x3	59 ms	67.2% $\pm$ 6.5	73.0% $\pm$ 7.4	55.8% $\pm$ 9.1
Stride-Effnet	73 ms	60.7% $\pm$ 6.6	66.9% $\pm$ 7.9	48.7% $\pm$ 8.3



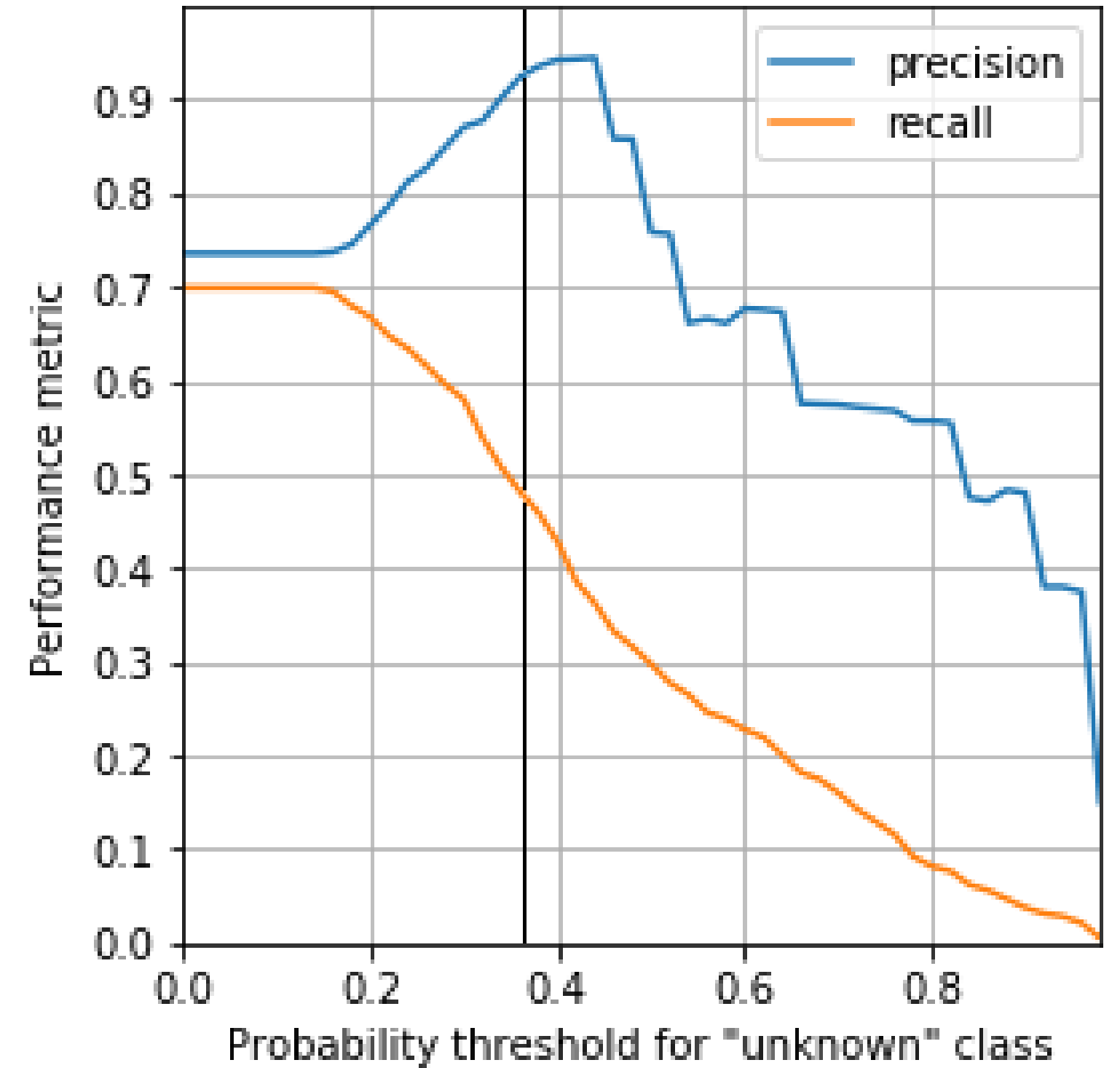
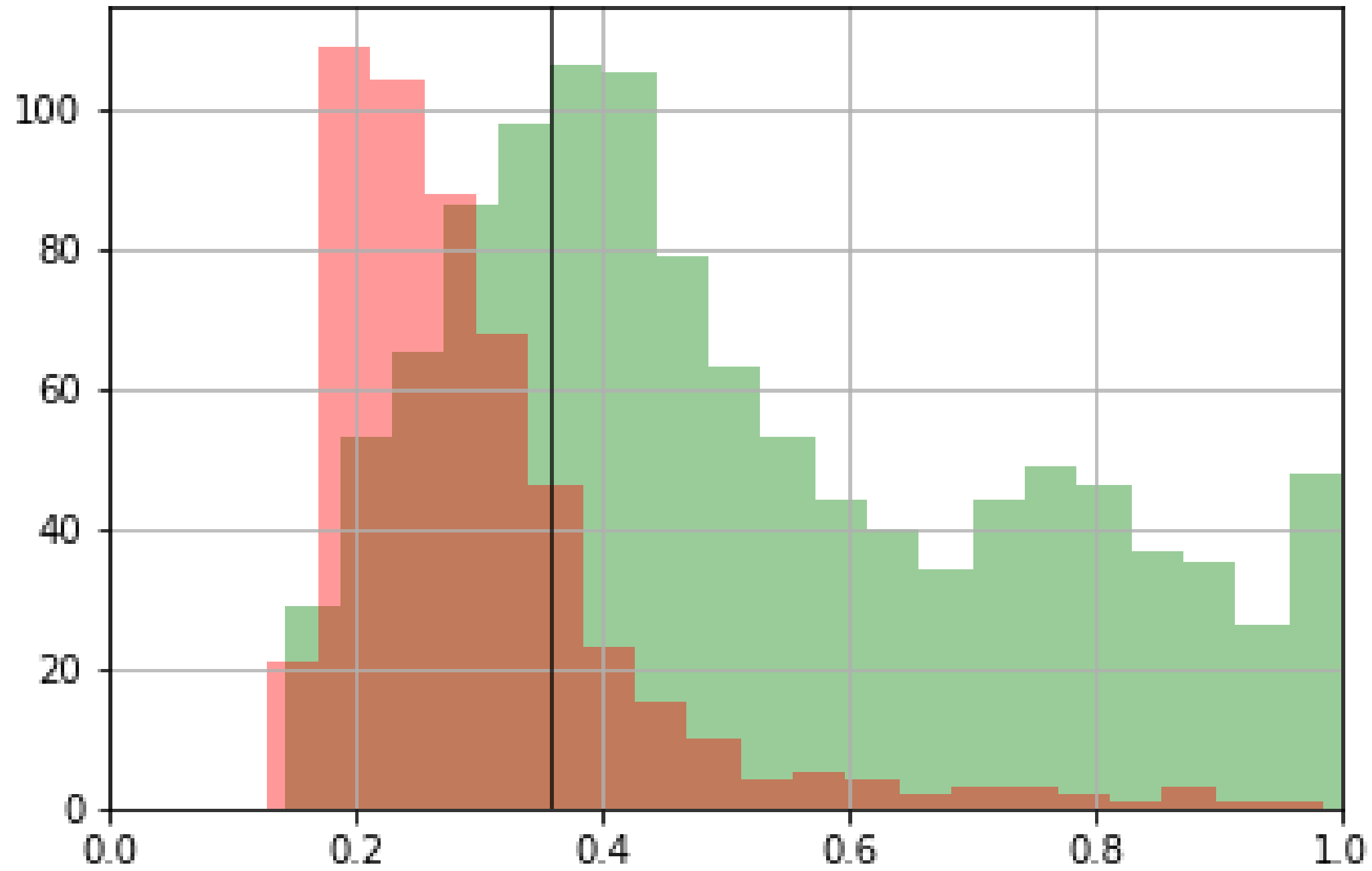


Predicted labels

# GROUPED CLASSIFICATION



# UNKNOWN CLASS



# EXPERIMENTAL DETAILS

# ALL MODELS

Model	Downsample	Convolution	L	F	MACC	RAM	FLASH
Baseline	maxpool 3x2	standard	3	24	10185 K	35 kB	405 kB
Baseline-DS	maxpool 3x2	DS	3	24	1567 K	55 kB	96 kB
Stride	stride 2x2	standard	3	22	2980 K	55 kB	372 kB
Stride-BTLN-DS	stride 2x2	BTLN-DS	3	22	445 K	47 kB	80 kB
Stride-DS-12	stride 2x2	DS	3	12	208 K	27 kB	88 kB
Stride-DS-16	stride 2x2	DS	3	16	291 K	36 kB	118 kB
Stride-DS-20	stride 2x2	DS	3	20	380 K	45 kB	149 kB
Stride-DS-24	stride 2x2	DS	3	24	477 K	54 kB	180 kB
Stride-DS-3x3	stride 2x2	DS	4	24	318 K	54 kB	95 kB
Stride-Effnet	stride 2x2	Effnet	3	22	468 K	47 kB	125 kB

# METHODS

Standard procedure for Urbansound8k

- Classification problem
- 4 second sound clips
- 10 classes
- 10-fold cross-validation, predefined
- Metric: Accuracy



# TRAINING SETTINGS

---

Samplerate (Hz)	22050
Melfilter bands	60
FFT length (samples)	1024
FFT hop (samples)	512
Classification window	31
Minibatch size	400
Epochs	100
Training samples/epoch	30000
Validation samples/epoch	5000
Learning rate	0.005
Nesterov momentum	NaN

---

# TRAINING

- NVidia RTX2060 GPU 6 GB
- 10 models x 10 folds = 100 training jobs
- 100 epochs
- 3 jobs in parallel
- 36 hours total

# EVALUATION

For each fold of each model

1. Select best model based on validation accuracy
2. Calculate accuracy on test set

For each model

- Measure CPU time on device

# YOUR MODEL WILL TRICK YOU

And the bugs can be hard to spot

# FAIL: INTEGER TRUNCATION

features: Fix integer truncation

Especially combined with meanstd normalization  
which makes range (-3,3),  
this accidentally removed a lot of details

```
----- microesc/features.py -----  
index 9d16d2a..6d9ca59 100644  
@@ -152,7 +154,7 @@ def load_sample(sample, settings, feature_dir, window_frames,  
    if window_frames is None:  
        padded = mels  
    else:  
-        padded = numpy.full((n_mels, window_frames), 0)  
+        padded = numpy.full((n_mels, window_frames), 0.0, dtype=float)  
    inp = mels[:, 0:min(window_frames, mels.shape[1])]  
    padded[:, 0:inp.shape[1]] = inp
```

# FAIL. DROPOUT LOCATION

models: Fix Dropout location

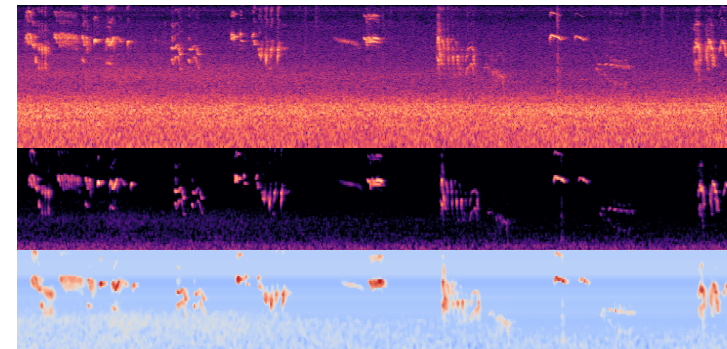
Suprised this was able to train at all before,  
as whole classes must have been dropped

Training and validation loss now follow eachother much more closely

```
----- microesc/models/sbcnn.py -----  
index 2e95242..3986bcf 100644  
@@ -98,12 +98,12 @@ def backend_dense1(x, n_classes, fc=64, regularization=0.001, dropout=0.5):  
    """  
    x = Flatten()(x)  
+   x = Dropout(dropout)(x)  
    x = Dense(fc, kernel_regularizer=l2(regularization))(x)  
    x = Activation('relu')(x)  
-   x = Dropout(dropout)(x)  
  
-   x = Dense(n_classes, kernel_regularizer=l2(regularization))(x)  
    x = Dropout(dropout)(x)  
+   x = Dense(n_classes, kernel_regularizer=l2(regularization))(x)  
    x = Activation('softmax')(x)  
    return x
```

# BACKGROUND

# MEL-SPECTROGRAM





# NOISE POLLUTION

Reduces health due to stress and loss of sleep

In Norway

- 1.9 million affected by road noise (2014, SSB)
- 10'000 healthy years lost per year (Folkehelseinstituttet)

In Europe

- 13 million suffering from sleep disturbance (EEA)
- 900'000 DALY lost (WHO)

# NOISE MAPPING

Simulation only, no direct measurements

