# Transformers

Marta Szuwarska, Mateusz Nizwantowski

musimy jakis plan zrobic

- data loader (introduction, silence, unknown, mfcc)
- Transformer basic architecture
- Transformer results
- other models listed
- for ever model introduction + results
- for best model plots

musimy jakis plan zrobic

- data loader (introduction, silence, unknown, mfcc)
- struktura plikow, how we handle silence
- mfcc
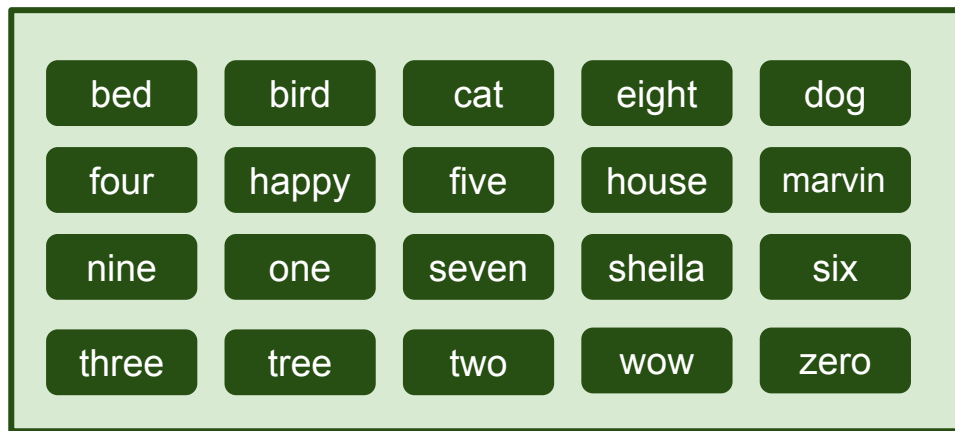- dlugo nam sie liczy slabe wyniki

# Original classes

| | | | | |
|---|---|---|---|---|
| bed | bird | cat | eight | dog |
| four | happy | five | house | marvin |
| nine | one | seven | sheila | six |
| three | tree | two | wow | zero |

| | | | | |
|---|---|---|---|---|
| up | left | yes | go | on |
| down | right | no | stop | off |

silence

# Original classes

bed | bird | cat | eight | dog
four | happy | five | house | marvin
nine | one | seven | sheila | six
three | tree | two | wow | zero

→ unknown

up | left | yes | go | on
down | right | no | stop | off

silence

# Modified classes

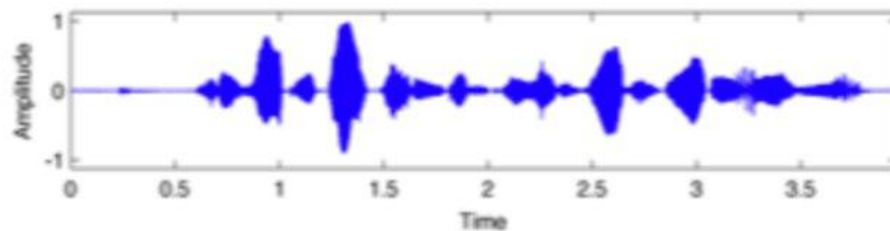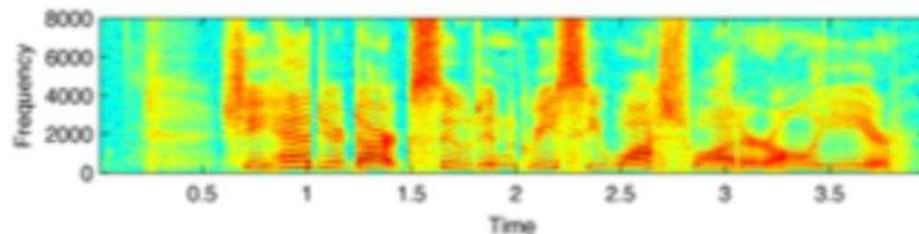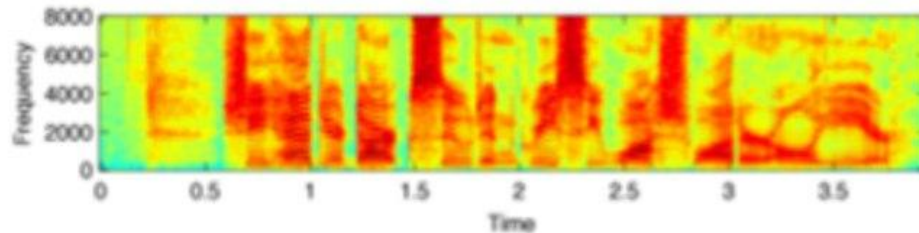| | | | | |
|---|---|---|---|---|
| up | left | yes | go | on |
| down | right | no | stop | off |

unknown

silence

**Time Domain Waveform**

**Spectrogram**

**MFCC Spectrogram**

# Early transformer architecture

```
Layer (type:depth-idx)                   Output Shape          Param #
================================================================================
SpeechCommandTransformer                 [16, 12]              256
├─MelSpectrogram: 1-1                    [16, 64, 101]         --
│    └─Spectrogram: 2-1                  [16, 201, 101]        --
│    └─MelScale: 2-2                     [16, 64, 101]         --
├─AmplitudeToDB: 1-2                     [16, 64, 101]         --
├─Sequential: 1-3                        [16, 64, 64, 101]     --
│    └─Conv2d: 2-3                       [16, 32, 64, 101]     320
│    └─BatchNorm2d: 2-4                  [16, 32, 64, 101]     64
│    └─ReLU: 2-5                         [16, 32, 64, 101]     --
│    └─Conv2d: 2-6                       [16, 64, 64, 101]     18,496
│    └─BatchNorm2d: 2-7                  [16, 64, 64, 101]     128
│    └─ReLU: 2-8                         [16, 64, 64, 101]     --
├─Linear: 1-4                            [16, 6464, 256]       16,640
├─TransformerEncoder: 1-5                [16, 6465, 256]       --
│    └─ModuleList: 2-9                   --                    --
│    │    └─TransformerEncoderLayer: 3-1 [16, 6465, 256]       527,104
│    │    └─TransformerEncoderLayer: 3-2 [16, 6465, 256]       527,104
│    │    └─TransformerEncoderLayer: 3-3 [16, 6465, 256]       527,104
│    │    └─TransformerEncoderLayer: 3-4 [16, 6465, 256]       527,104
├─Linear: 1-6                            [16, 12]              3,084
================================================================================
Total params: 2,147,404
Trainable params: 2,147,404
Non-trainable params: 0
Total mult-adds (Units.GIGABYTES): 1.96
================================================================================
Input size (MB): 1.02
Forward/backward pass size (MB): 4607.58
Params size (MB): 4.38
Estimated Total Size (MB): 4612.98
================================================================================
```
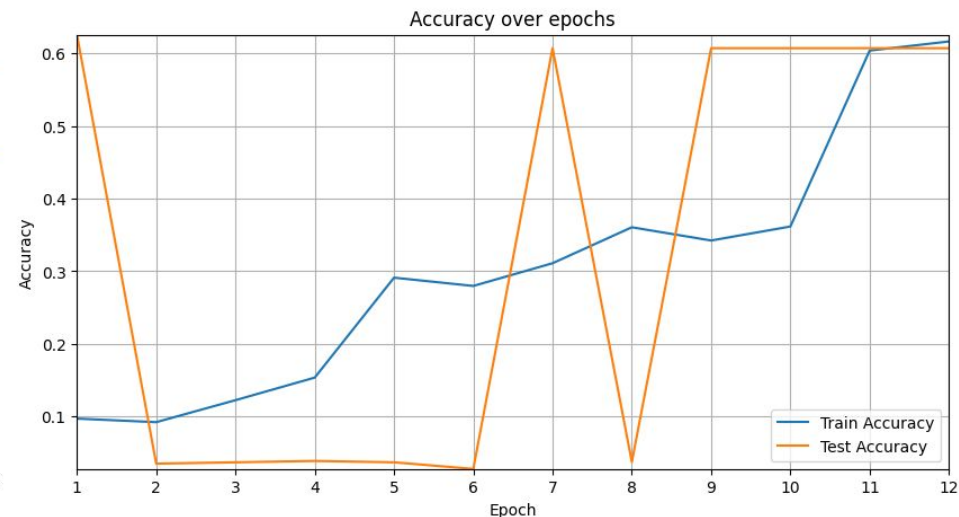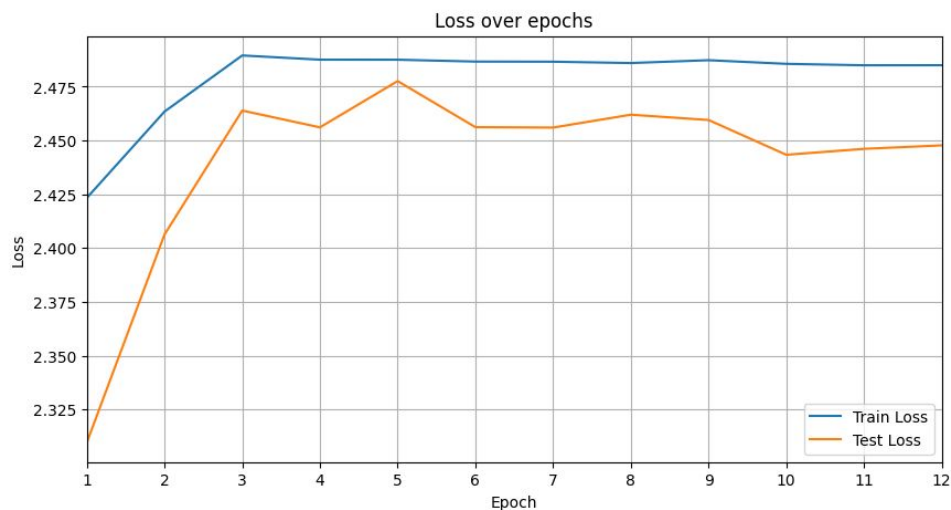
Epoch 1/20:    0%|                | 1/3292 [00:17<15:37:26, 17.09s/it,

OutOfMemoryError: CUDA out of memory.
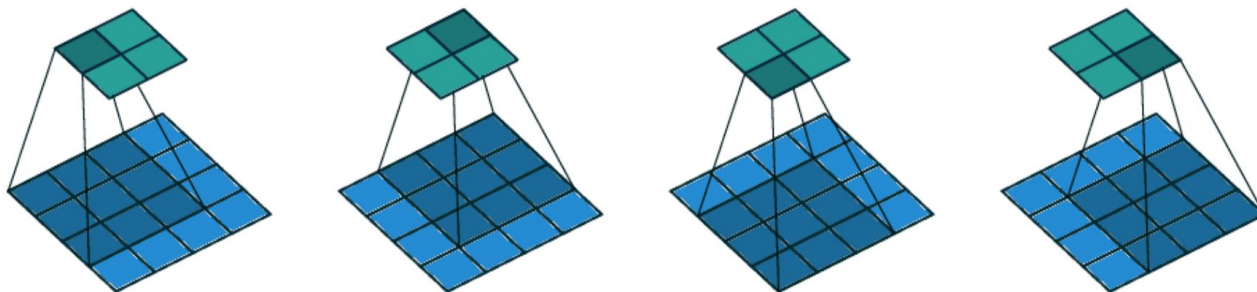
# Early transformer results



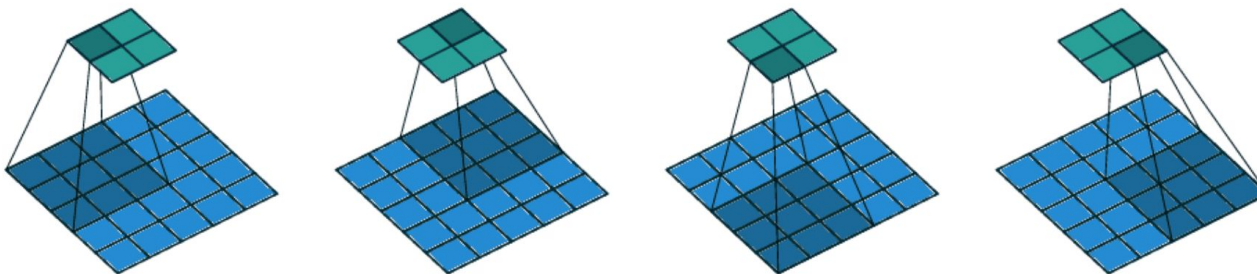Best test accuracy: 60.71%, corresponding train accuracy: 34.23%

# Strided convolution

stride=1



https://www.baeldung.com/wp-content/uploads/sites/4/2023/10/Screenshot-2023-10-10-at-1.11.45-PM.png

stride=2



https://www.baeldung.com/wp-content/uploads/sites/4/2023/10/Screenshot-2023-10-10-at-1.11.23-PM.png

# Strided transformer architecture - original vs modified

```
================================================================
Layer (type:depth-idx)              Output Shape        Param #
================================================================
SpeechCommandTransformer            [16, 31]            256
├─MelSpectrogram: 1-1               [16, 64, 101]       --
│    └─Spectrogram: 2-1             [16, 201, 101]      --
│    └─MelScale: 2-2                [16, 64, 101]       --
├─AmplitudeToDB: 1-2                [16, 64, 101]       --
├─Sequential: 1-3                   [16, 64, 16, 26]    --
│    └─Conv2d: 2-3                  [16, 32, 32, 51]    320
│    └─BatchNorm2d: 2-4             [16, 32, 32, 51]    64
│    └─ReLU: 2-5                    [16, 32, 32, 51]    --
│    └─Conv2d: 2-6                  [16, 64, 16, 26]    18,496
│    └─BatchNorm2d: 2-7             [16, 64, 16, 26]    128
│    └─ReLU: 2-8                    [16, 64, 16, 26]    --
├─Linear: 1-4                       [16, 416, 256]      16,640
├─TransformerEncoder: 1-5           [16, 417, 256]      --
│    └─ModuleList: 2-9              --                  --
│    │    └─TransformerEncoderLayer: 3-1  [16, 417, 256]  527,104
│    │    └─TransformerEncoderLayer: 3-2  [16, 417, 256]  527,104
│    │    └─TransformerEncoderLayer: 3-3  [16, 417, 256]  527,104
│    │    └─TransformerEncoderLayer: 3-4  [16, 417, 256]  527,104
├─Linear: 1-6                       [16, 31]            7,967
================================================================
Total params: 2,152,287
Trainable params: 2,152,287
Non-trainable params: 0
Total mult-adds (M): 148.75
================================================================
Input size (MB): 1.02
Forward/backward pass size (MB): 307.11
Params size (MB): 4.40
Estimated Total Size (MB): 312.53
================================================================
```

```
================================================================
Layer (type:depth-idx)              Output Shape        Param #
================================================================
SpeechCommandTransformer            [16, 12]            256
├─MelSpectrogram: 1-1               [16, 64, 101]       --
│    └─Spectrogram: 2-1             [16, 201, 101]      --
│    └─MelScale: 2-2                [16, 64, 101]       --
├─AmplitudeToDB: 1-2                [16, 64, 101]       --
├─Sequential: 1-3                   [16, 64, 16, 26]    --
│    └─Conv2d: 2-3                  [16, 32, 32, 51]    320
│    └─BatchNorm2d: 2-4             [16, 32, 32, 51]    64
│    └─ReLU: 2-5                    [16, 32, 32, 51]    --
│    └─Conv2d: 2-6                  [16, 64, 16, 26]    18,496
│    └─BatchNorm2d: 2-7             [16, 64, 16, 26]    128
│    └─ReLU: 2-8                    [16, 64, 16, 26]    --
├─Linear: 1-4                       [16, 416, 256]      16,640
├─TransformerEncoder: 1-5           [16, 417, 256]      --
│    └─ModuleList: 2-9              --                  --
│    │    └─TransformerEncoderLayer: 3-1  [16, 417, 256]  527,104
│    │    └─TransformerEncoderLayer: 3-2  [16, 417, 256]  527,104
│    │    └─TransformerEncoderLayer: 3-3  [16, 417, 256]  527,104
│    │    └─TransformerEncoderLayer: 3-4  [16, 417, 256]  527,104
├─Linear: 1-6                       [16, 12]            3,084
================================================================
Total params: 2,147,404
Trainable params: 2,147,404
Non-trainable params: 0
Total mult-adds (M): 148.68
================================================================
Input size (MB): 1.02
Forward/backward pass size (MB): 307.10
Params size (MB): 4.38
Estimated Total Size (MB): 312.51
================================================================
```
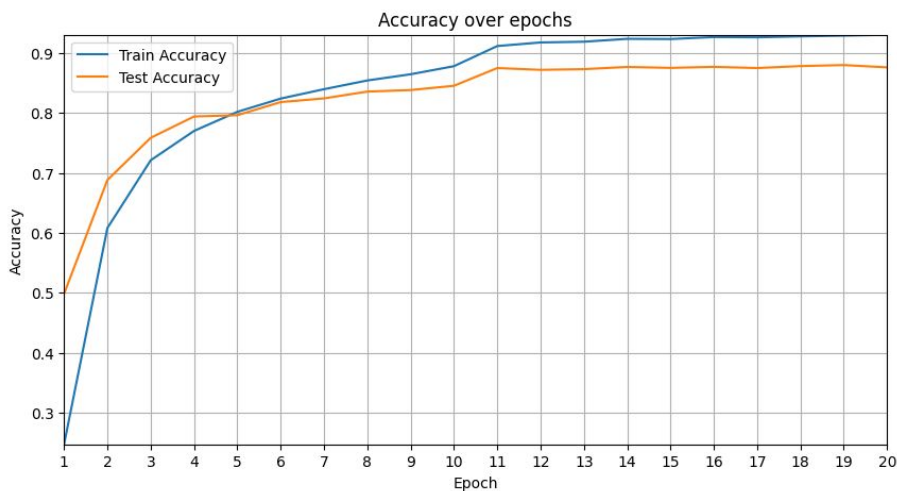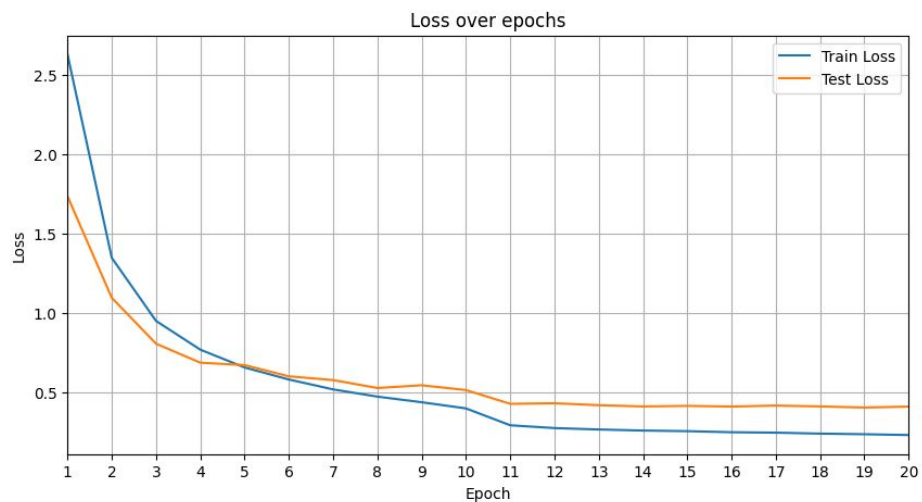
# Transformer - original classes

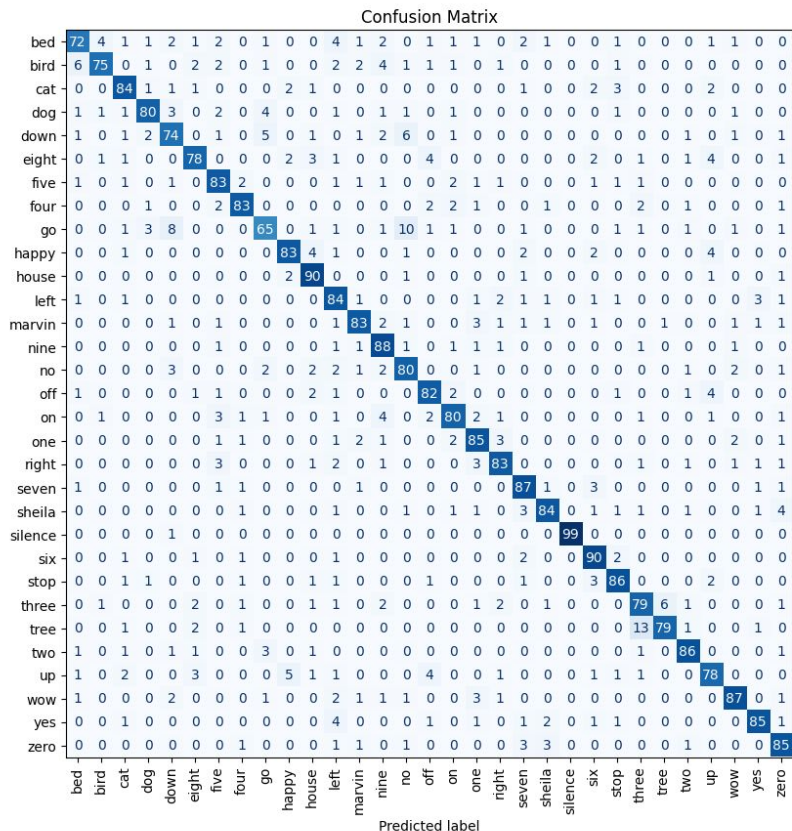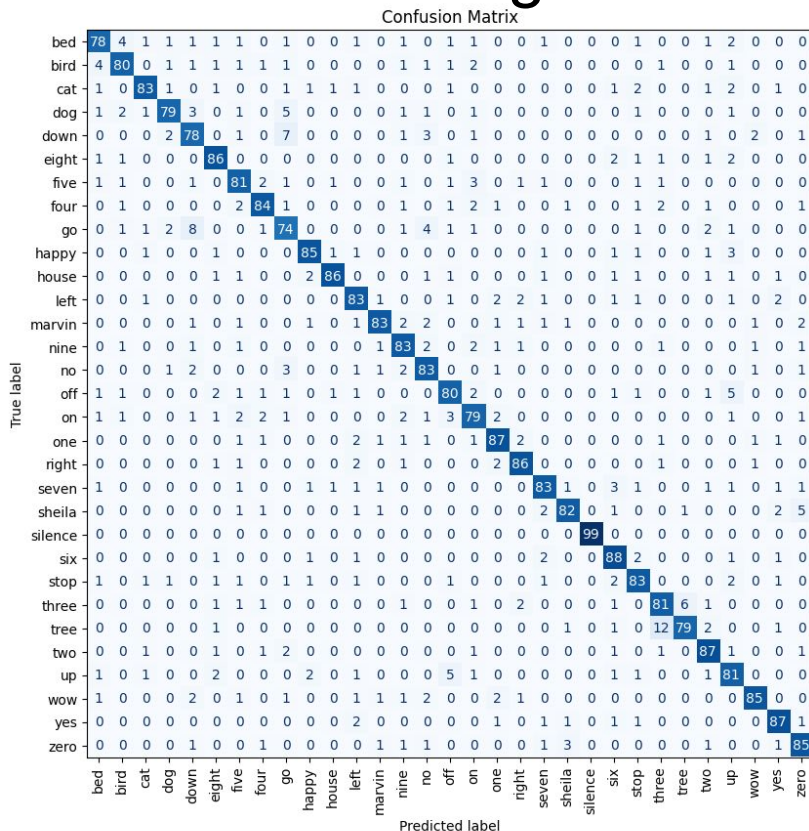| no. | optimizer | learning rate | epochs | train acc | test acc | train time |
|-----|-----------|---------------|--------|-----------|----------|------------|
| 1 | AdamW | 0.000005 | 10 | 87.62 | 83.88 | 4h 25m |
| 2 | Adam | 0.001 | 9 | 3.44 | 3.53 | 4h 50m |
| 3 | AdamW | 0.000005 | 20 | 93.04 | 87.60 | 14h 17m |

All experiments were conducted with stride = 2, batch size = 16, embedding dimension = 256, StepLR scheduler (step size = 10, gamma = 0.1), no position embedding and weight decay = 0.001 for AdamW optimizer.

# Transformer - original classes - best model

# Transformer - original classes - best model

# Transformer - modified classes

| no. | optimizer | learning rate | epochs | scheduler | pos embed | train acc | test acc | train time |
|-----|-----------|---------------|--------|-----------|-----------|-----------|----------|------------|
| 1 | AdamW | 0.000005 | 10 | StepLR | no | 69.88 | 69.67 | 4h 6m |
| 2 | Adam | 0.001 | 20 | StepLR | no | 37.33 | 60.71 | 1h 30m |
| 3 | AdamW | 0.000005 | 20 | StepLR | no | 85.23 | 83.04 | 4h 50m |
| 4 | AdamW | 0.000005 | 20 | - | no | 84.75 | 77.84 | 5h 3m |
| 5 | AdamW | 0.000005 | 20 | StepLR | yes | 82.24 | 80.91 | 6h 54m |
| 6 | AdamW | 0.000005 | 20 | CosineAnnealingLR | no | 85.83 | 82.11 | 14h 16m |

All experiments were conducted with stride = 2, batch size = 16, embedding dimension = 256, 20 epochs and weight decay = 0.001 for AdamW optimizer.

# Transformer - modified classes - best model

# Transformer - modified classes - best model
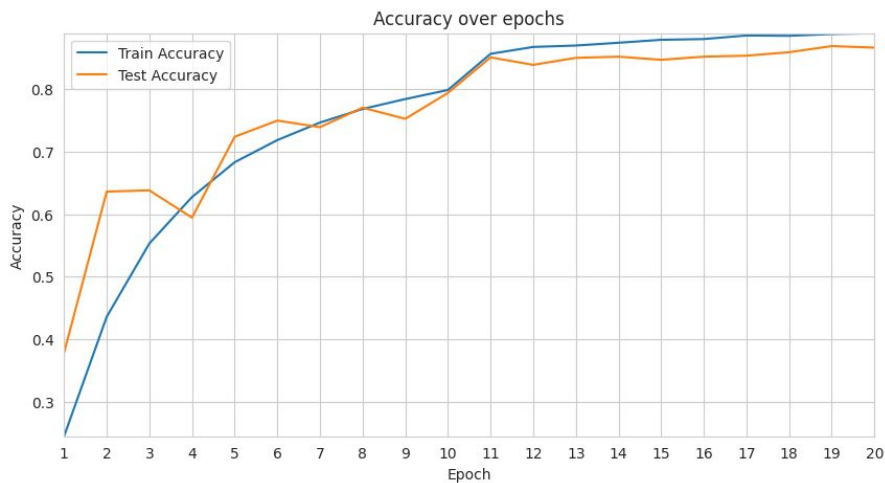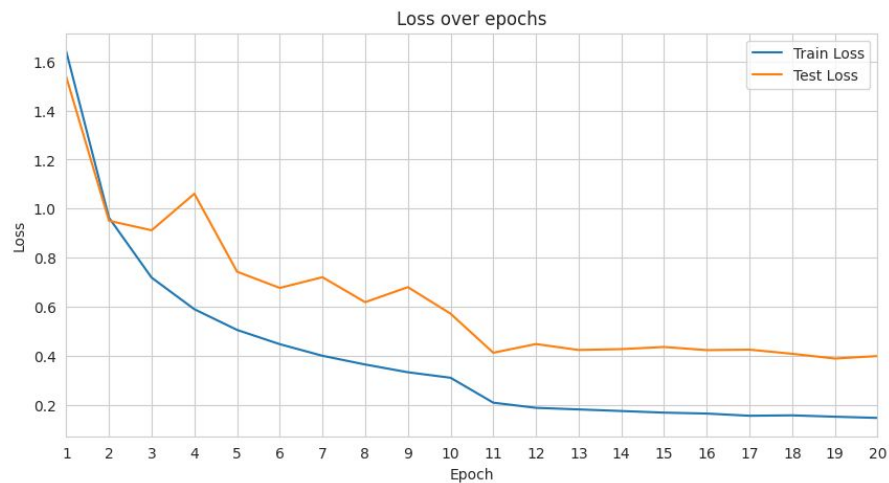
# Transformer - different architecture

```
======================================================================
Layer (type:depth-idx)                  Output Shape          Param #
======================================================================
SpeechCommandTransformer                [16, 12]              256
├─MelSpectrogram: 1-1                   [16, 64, 101]         --
│    └─Spectrogram: 2-1                 [16, 201, 101]        --
│    └─MelScale: 2-2                    [16, 64, 101]         --
├─AmplitudeToDB: 1-2                    [16, 64, 101]         --
├─Sequential: 1-3                       [16, 64, 16, 25]      --
│    └─Conv2d: 2-3                      [16, 32, 64, 101]     320
│    └─BatchNorm2d: 2-4                 [16, 32, 64, 101]     64
│    └─ReLU: 2-5                        [16, 32, 64, 101]     --
│    └─MaxPool2d: 2-6                   [16, 32, 32, 50]      --
│    └─Conv2d: 2-7                      [16, 64, 16, 25]      18,496
│    └─BatchNorm2d: 2-8                 [16, 64, 16, 25]      128
│    └─ReLU: 2-9                        [16, 64, 16, 25]      --
├─Linear: 1-4                           [16, 400, 256]        16,640
├─TransformerEncoder: 1-5               [16, 401, 256]        --
│    └─ModuleList: 2-10                 --                    --
│    │    └─TransformerEncoderLayer: 3-1 [16, 401, 256]       527,104
│    │    └─TransformerEncoderLayer: 3-2 [16, 401, 256]       527,104
│    │    └─TransformerEncoderLayer: 3-3 [16, 401, 256]       527,104
│    │    └─TransformerEncoderLayer: 3-4 [16, 401, 256]       527,104
├─Linear: 1-6                           [16, 12]              3,084
======================================================================
Total params: 2,147,404
Trainable params: 2,147,404
Non-trainable params: 0
Total mult-adds (M): 168.68
======================================================================
Input size (MB): 1.02
Forward/backward pass size (MB): 335.41
Params size (MB): 4.38
Estimated Total Size (MB): 340.82
======================================================================
```
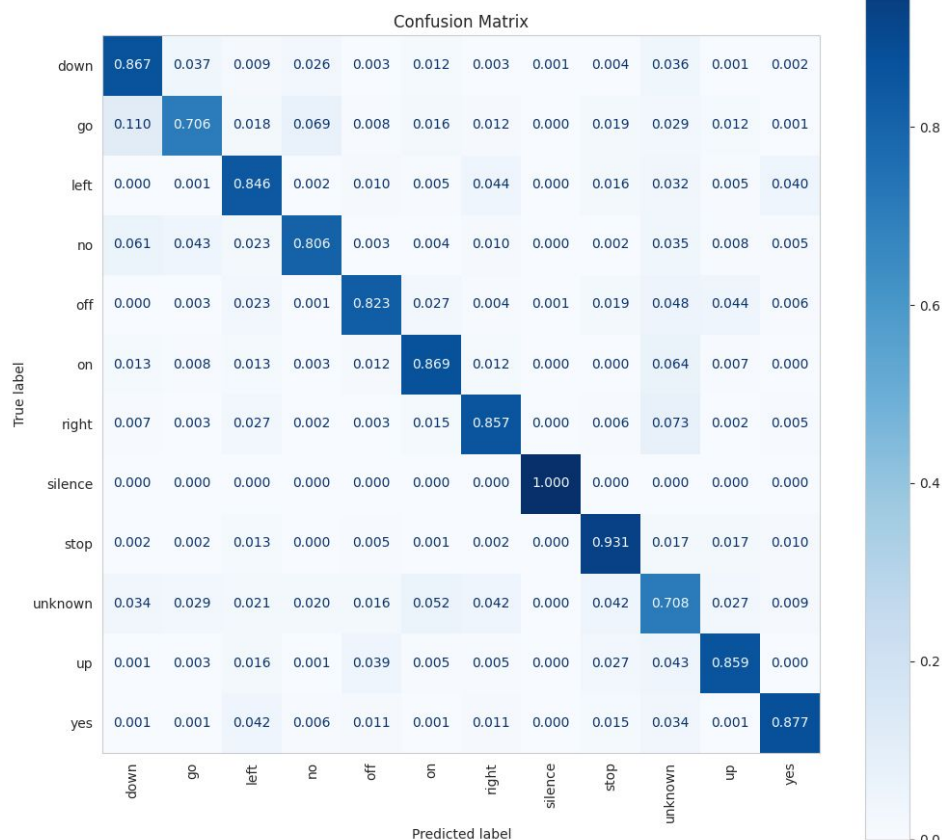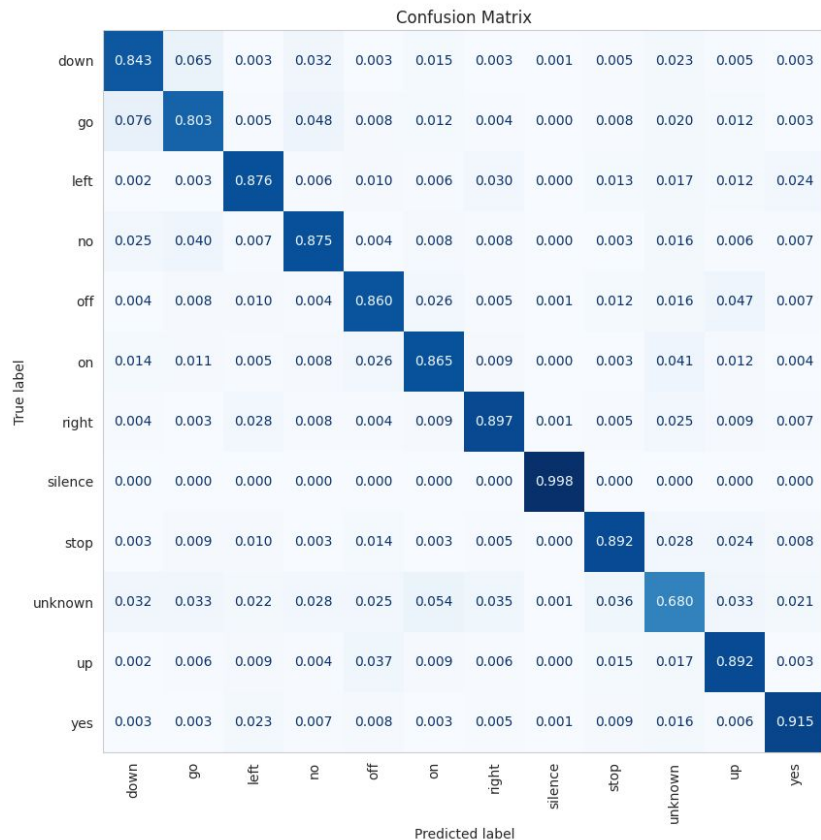
Best test accuracy:
86.67%
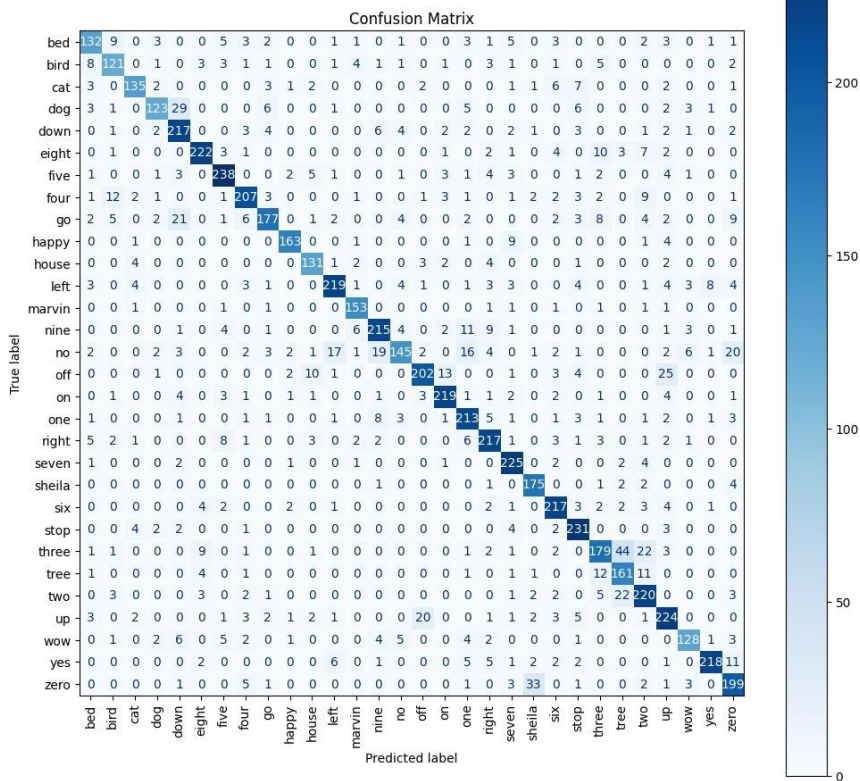
Corresponding train accuracy:
88.98%

# Transformer - different architecture - even better results

# Transformer - different architecture - even better results

# GRU - baseline model - 1D - original classes



Confusion Matrix

| Layer (type:depth-idx) | Output Shape | Param # |
|---|---|---|
| SpeechCommandGRU | [32, 30] | -- |
| ├─Sequential: 1-1 | [32, 40, 39] | -- |
| │    └─Conv1d: 2-1 | [32, 32, 39] | 2,592 |
| │    └─ReLU: 2-2 | [32, 32, 39] | -- |
| │    └─Conv1d: 2-3 | [32, 40, 39] | 3,880 |
| │    └─ReLU: 2-4 | [32, 40, 39] | -- |
| ├─GRU: 1-2 | [32, 39, 256] | 427,008 |
| ├─Linear: 1-3 | [32, 30] | 7,710 |

Total params: 441,190
Trainable params: 441,190
Non-trainable params: 0
Total mult-adds (Units.MEGABYTES): 541.23

Input size (MB): 0.20
Forward/backward pass size (MB): 3.28
Params size (MB): 1.76
Estimated Total Size (MB): 5.25
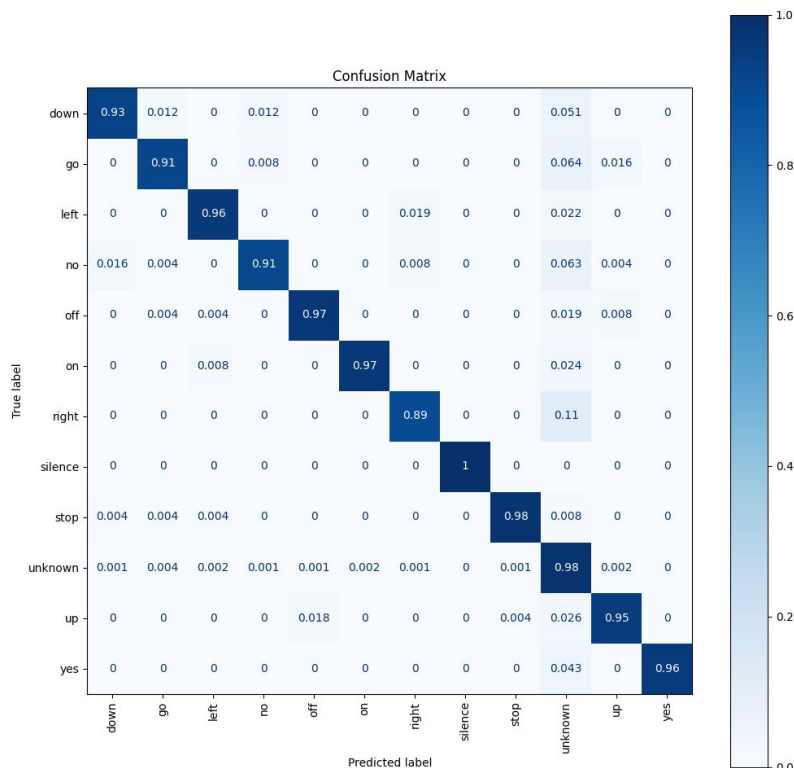
Test Accuracy: 82.31%

# GRU - baseline model - 1D - modified classes



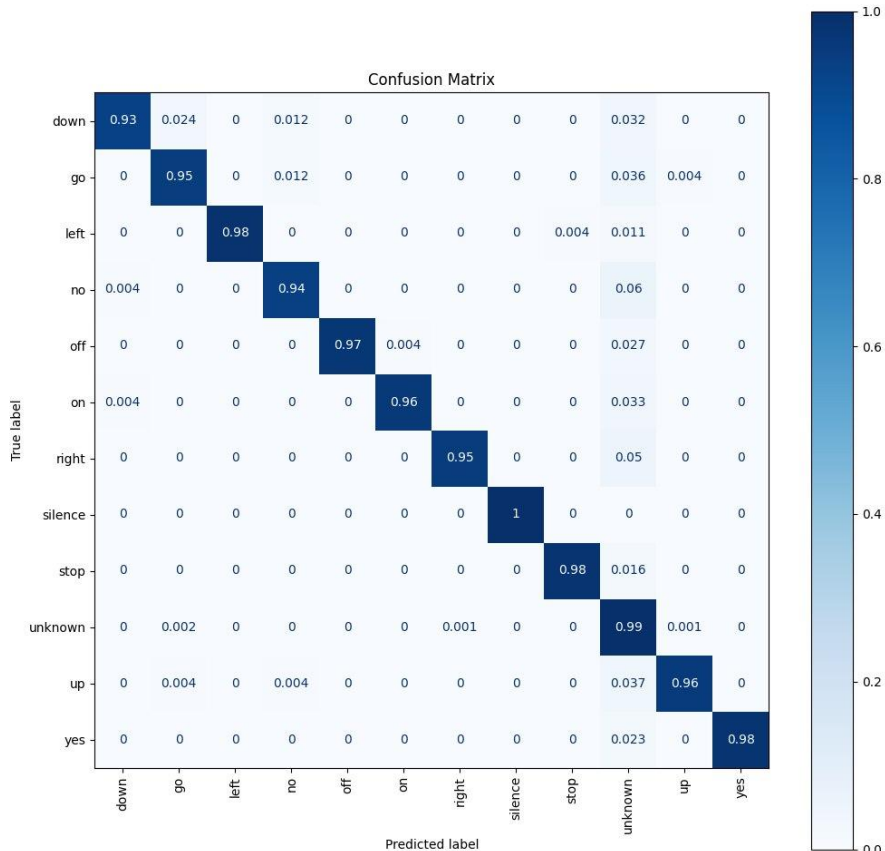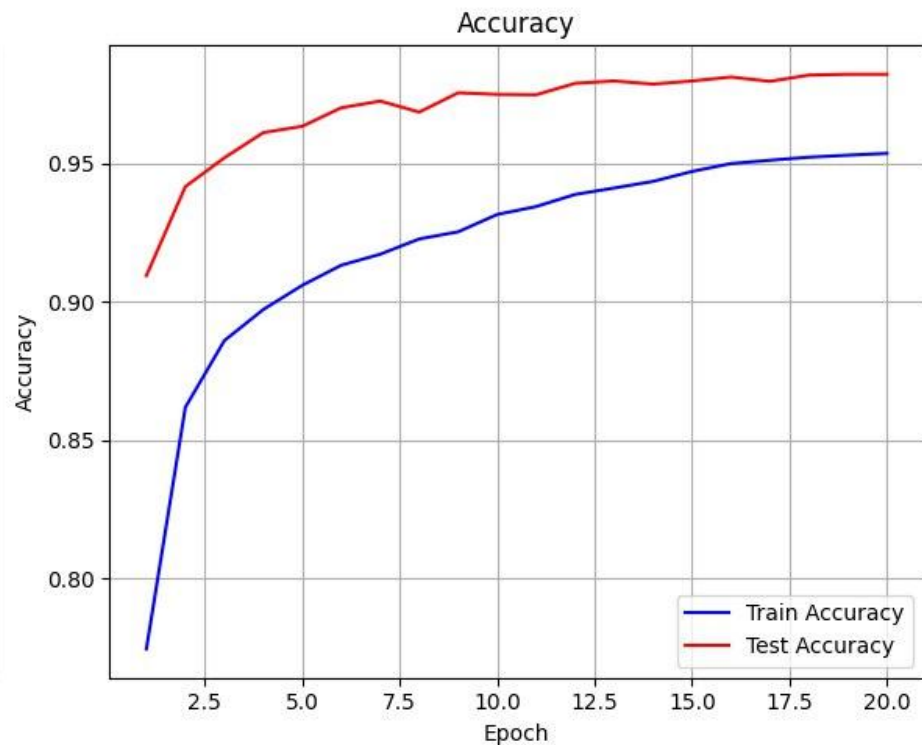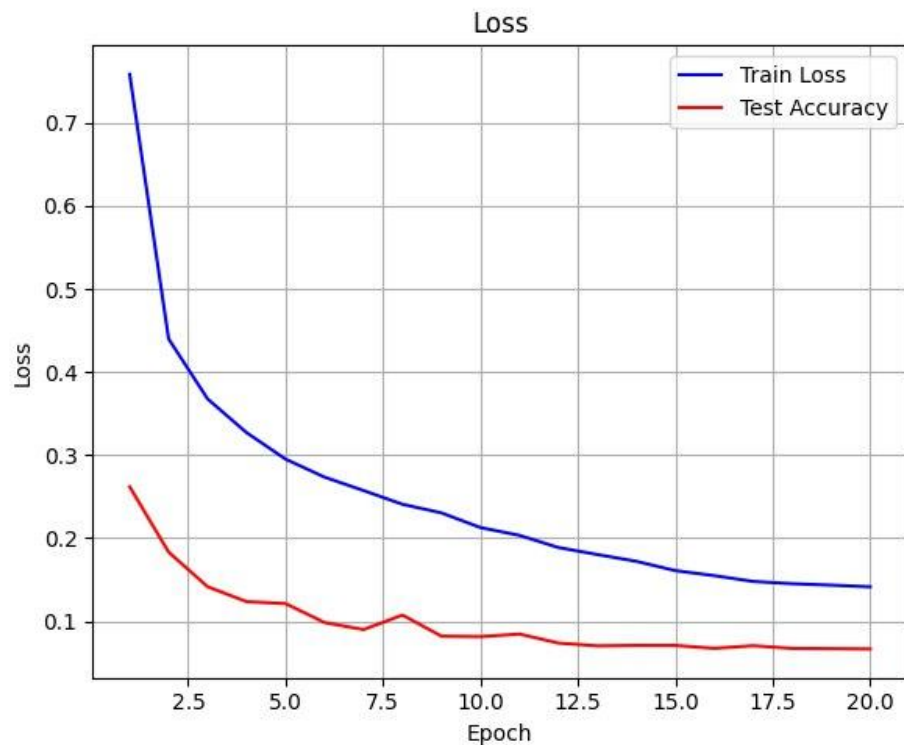Test Accuracy: 93.6%

# GRU - 2D - modified classes


Confusion Matrix

```
Layer (type:depth-idx)              Output Shape          Param #
===============================================================
SpeechCommandCRNN                   [2, 12]               --
├─Sequential: 1-1                   [2, 128, 40, 101]     --
│    └─Conv2d: 2-1                  [2, 32, 40, 101]      320
│    └─BatchNorm2d: 2-2             [2, 32, 40, 101]      64
│    └─ReLU: 2-3                    [2, 32, 40, 101]      --
│    └─Conv2d: 2-4                  [2, 64, 40, 101]      18,496
│    └─BatchNorm2d: 2-5             [2, 64, 40, 101]      128
│    └─ReLU: 2-6                    [2, 64, 40, 101]      --
│    └─Conv2d: 2-7                  [2, 128, 40, 101]     73,856
│    └─BatchNorm2d: 2-8             [2, 128, 40, 101]     256
│    └─ReLU: 2-9                    [2, 128, 40, 101]     --
├─GRU: 1-2                          [2, 101, 256]         4,328,448
├─Dropout: 1-3                      [2, 256]              --
├─Linear: 1-4                       [2, 12]               3,084
===============================================================
Total params: 4,424,652
Trainable params: 4,424,652
Non-trainable params: 0
Total mult-adds (Units.GIGABYTES): 1.62
===============================================================
Input size (MB): 0.03
Forward/backward pass size (MB): 29.37
Params size (MB): 17.70
Estimated Total Size (MB): 47.10
===============================================================
```

Test Accuracy: 96.93%

# Gru - modified classes - best model

# Gru - modified classes - best model

# Gru - modified classes - best model - over 10 runs

| set | Accuracy | | Loss | |
|---|---|---|---|---|
| | **mean** | **std** | **mean** | **std** |
| train | 95.42% | 0.09% | 0.1352 | 0.0026 |
| test | 98.02% | 0.08% | 0.0721 | 0.0021 |
| validation | **97.62%** | 0.07% | 0.0892 | 0.0046 |

Train time: under 4 minutes

# Thank you for your "attention"

Are there any questions?

# Sources

1. *[Attention is all you need](#)*

2. *[1b3b transformer series](#)*

3. [Andrej Karpathy transformers explained](#)

4. [Post on medium about speech recognition](#)

5. [Pytorch tutorial on speech commands recognition](#)