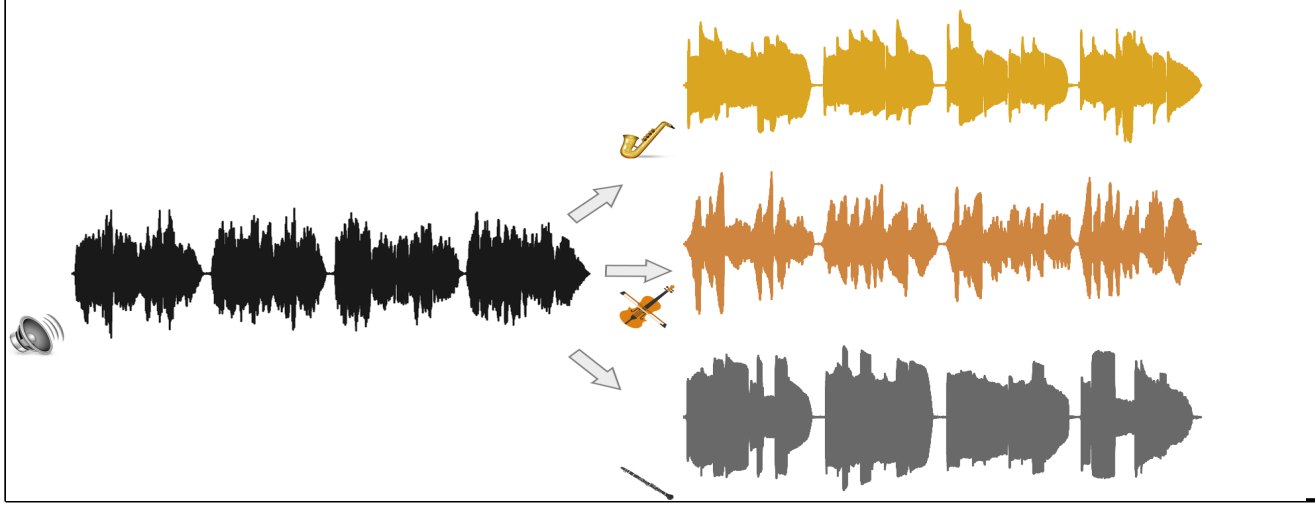


## Introduction

Single channel source separation (SCSS) has important applications (e.g. speech recognition). Approaches taken have been NMF, denoising autoencoders, and CNNs. Ability of CNNs to interpret local connectivity is appealing. Recently, [1] proposed a variant of an autoencoder that uses convolutional layers to process spatiotemporal information that achieved comparable results to previous architectures. We further investigated the performance of the convolutional autoencoder on the Bach10 dataset.

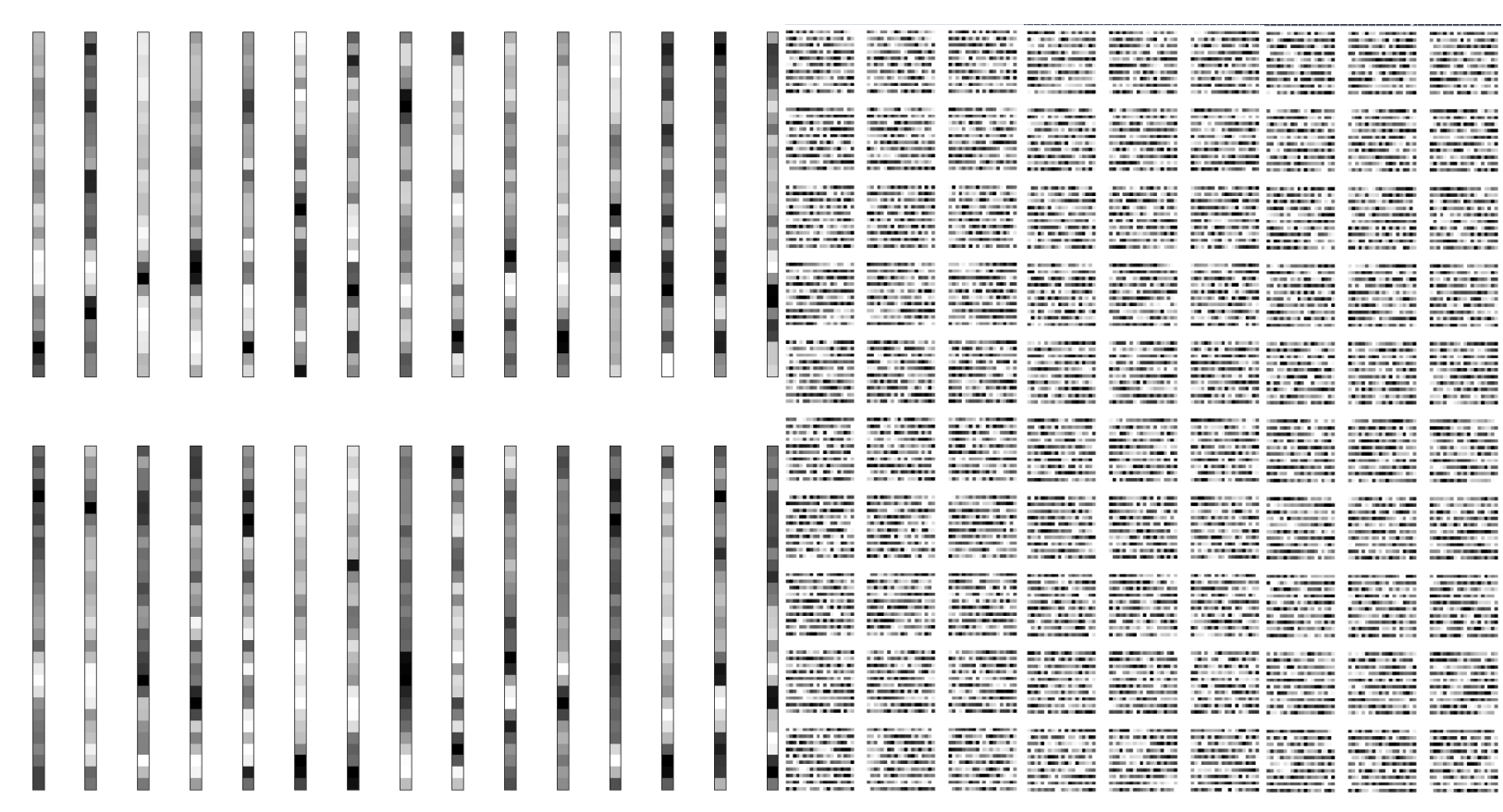


## Dataset

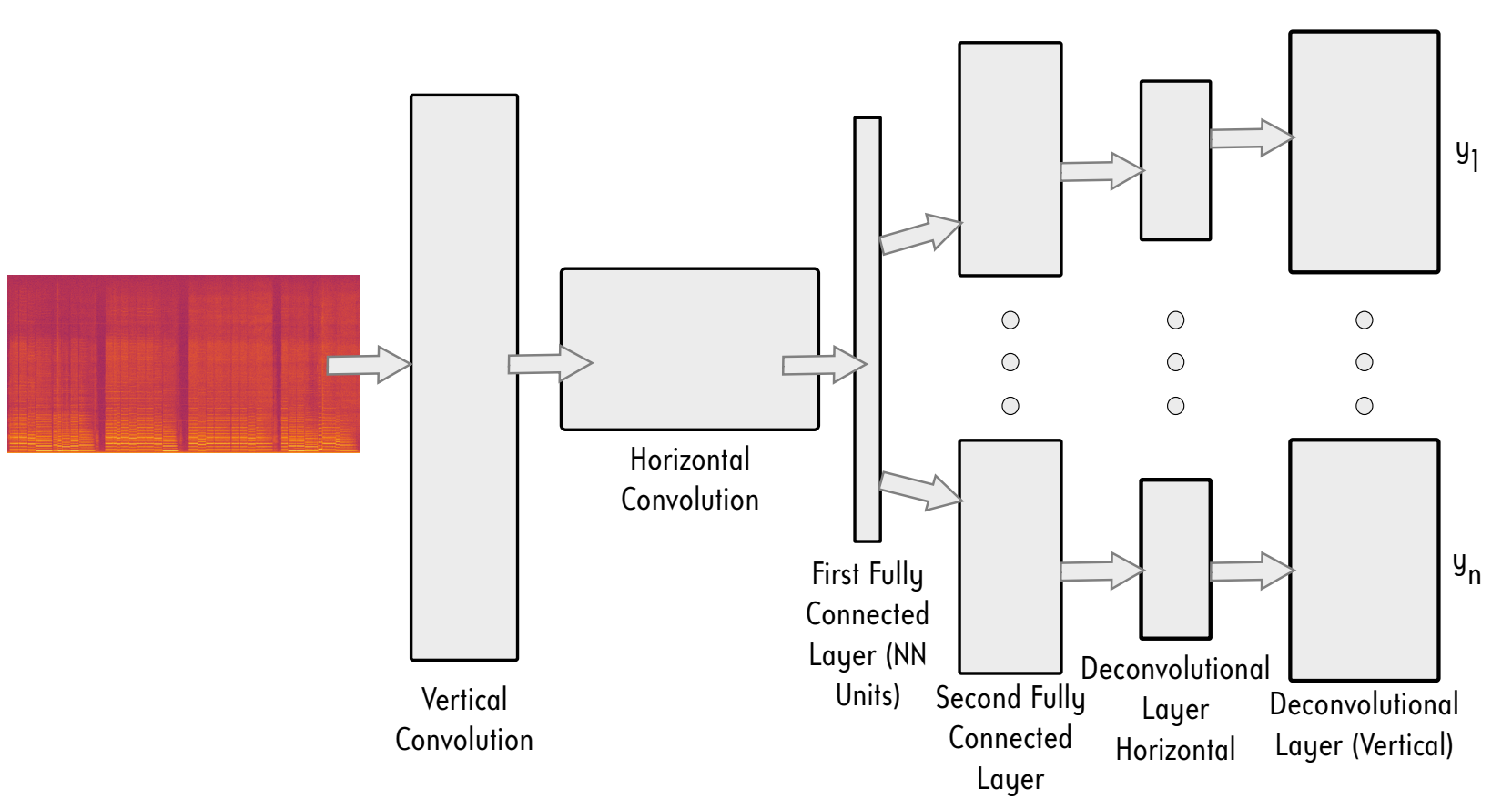
- Bach10 Dataset
- 10 songs
  - Contains full song and individual sources
  - Individual sources include:
    - (1) Bassoon
    - (2) Clarinet
    - (3) Saxophone
    - (4) Violin



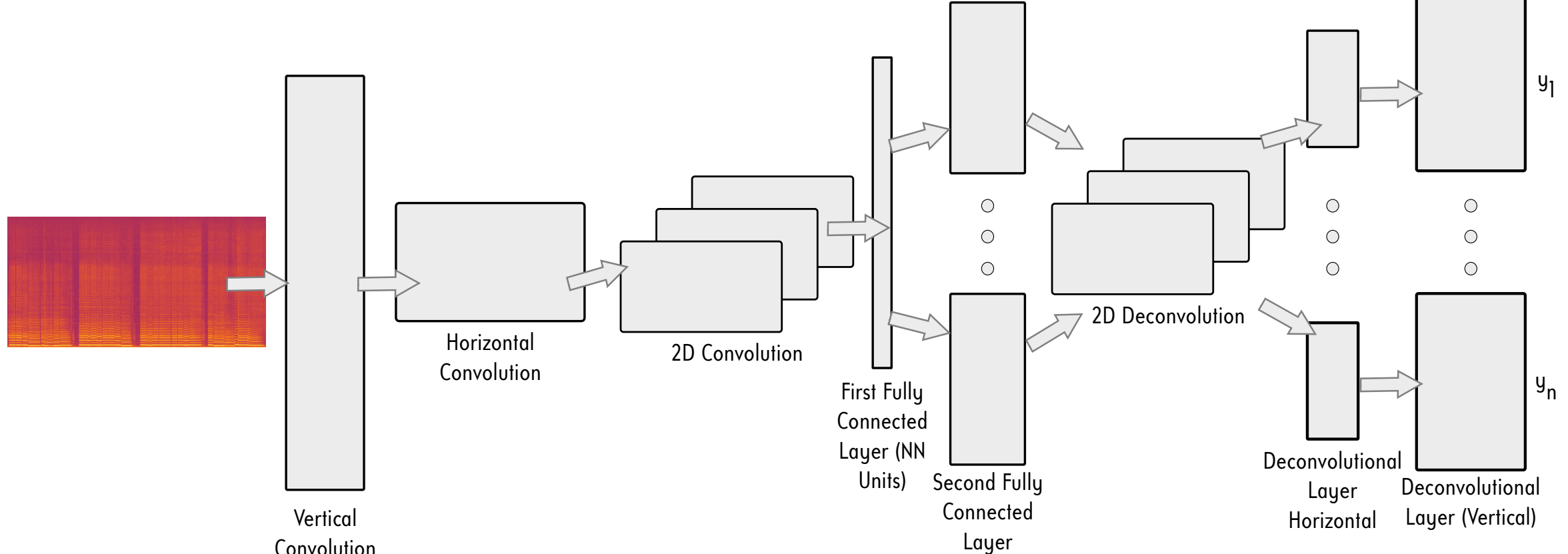
## Filters



## General Network Architecture



## Novel Network Architecture

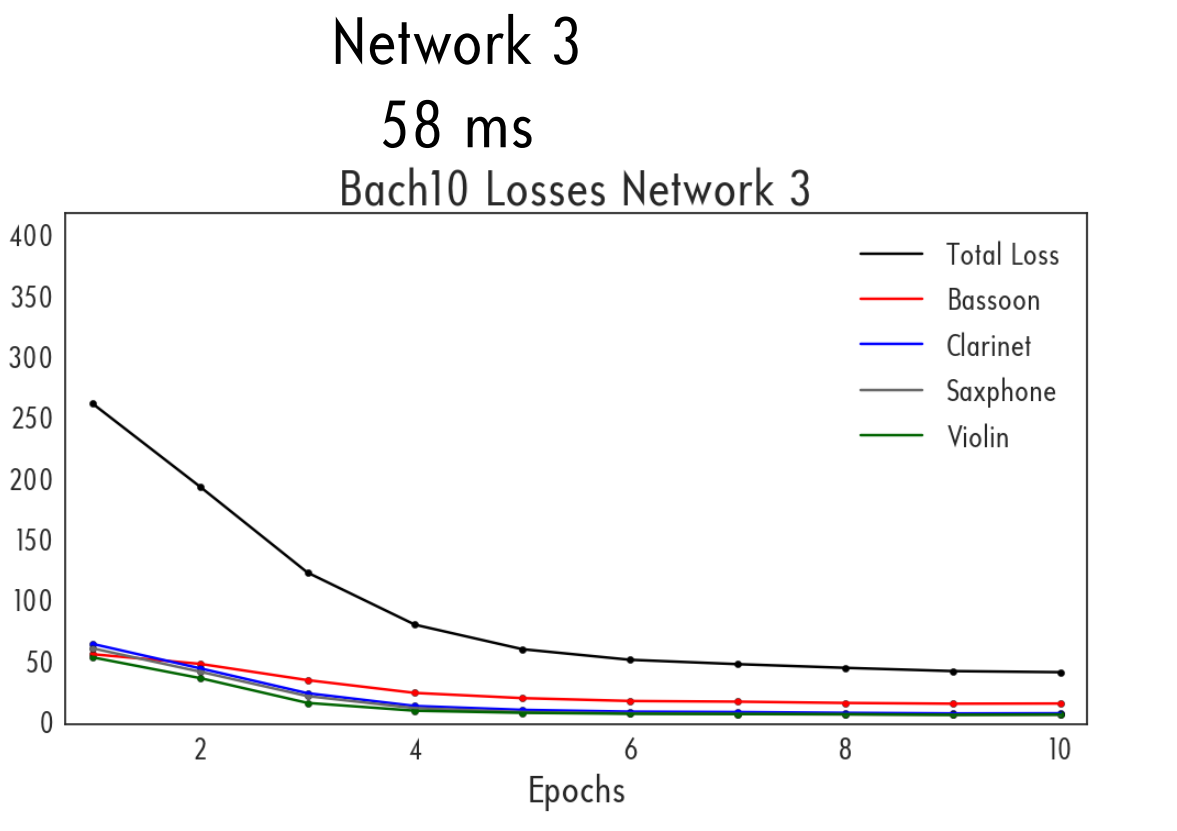
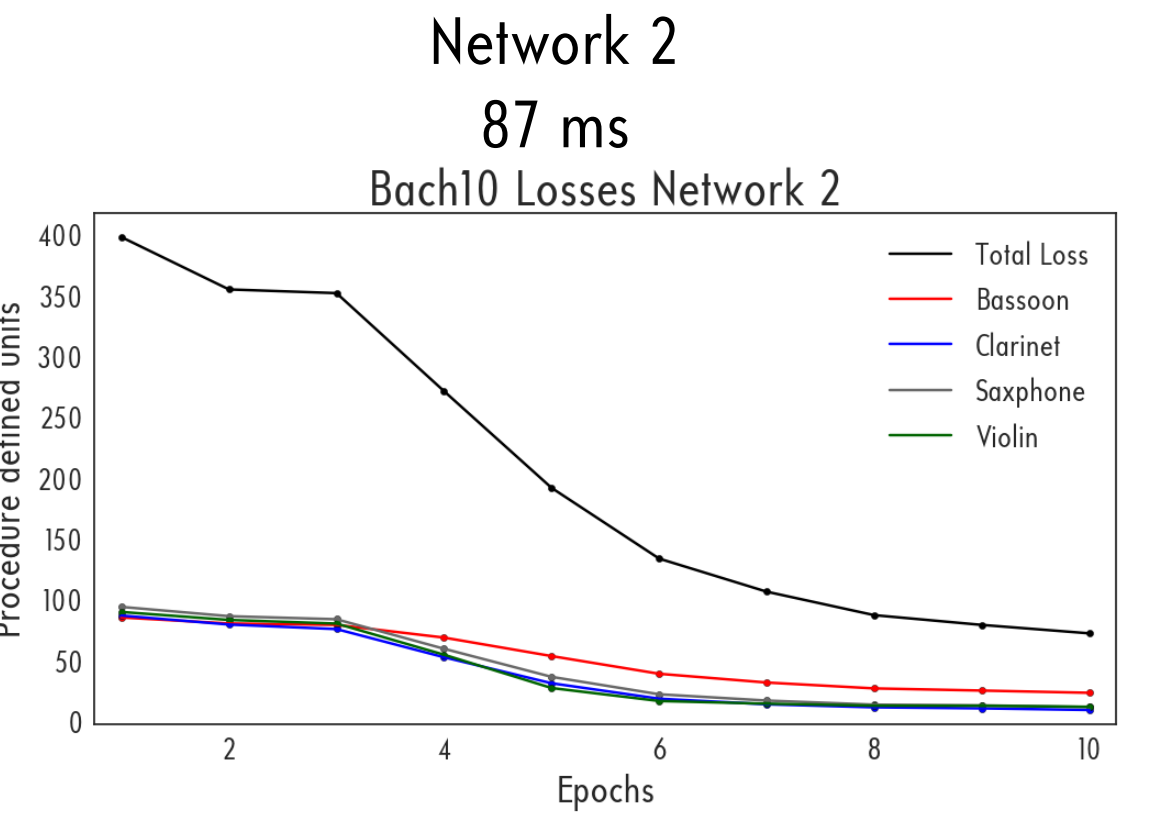
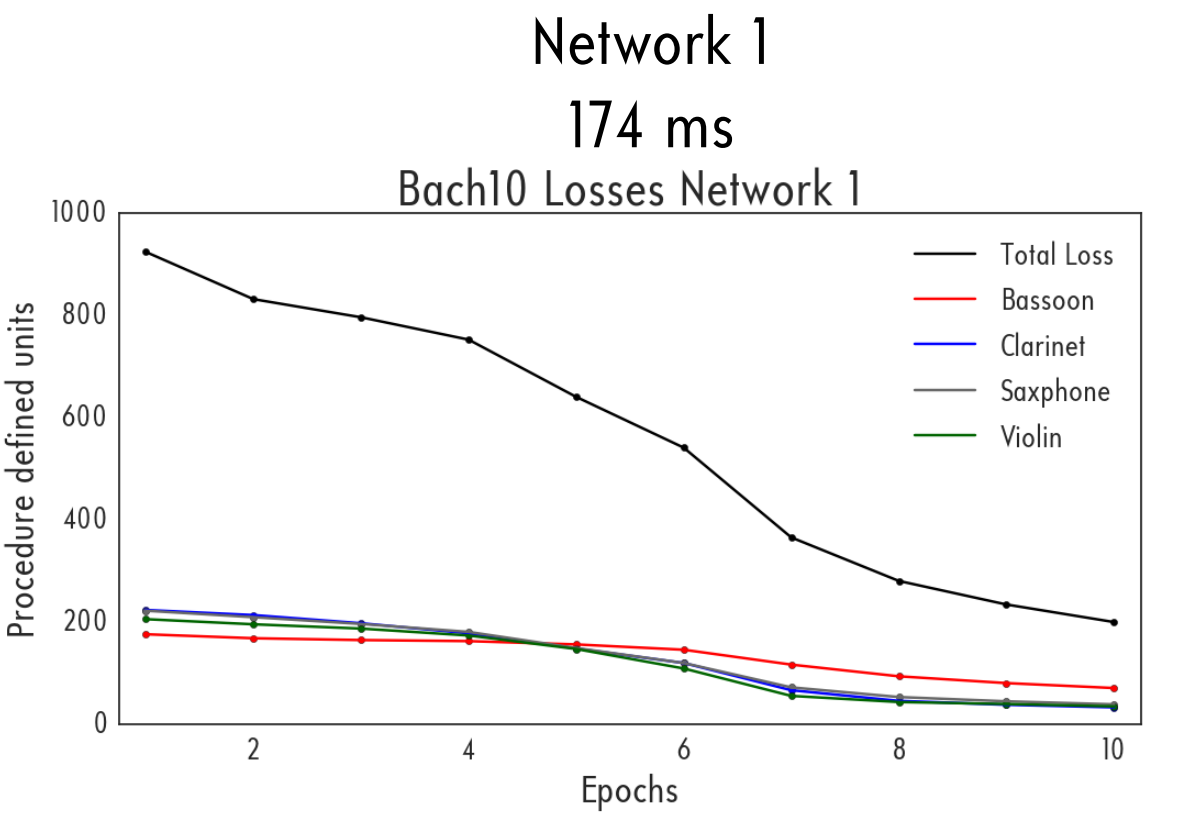


Network Parameters (unless noted otherwise):

- Max # epochs: 10-20
- Initial learning rate: 0.001
- Training songs: First 5
- Input feature size: 1025

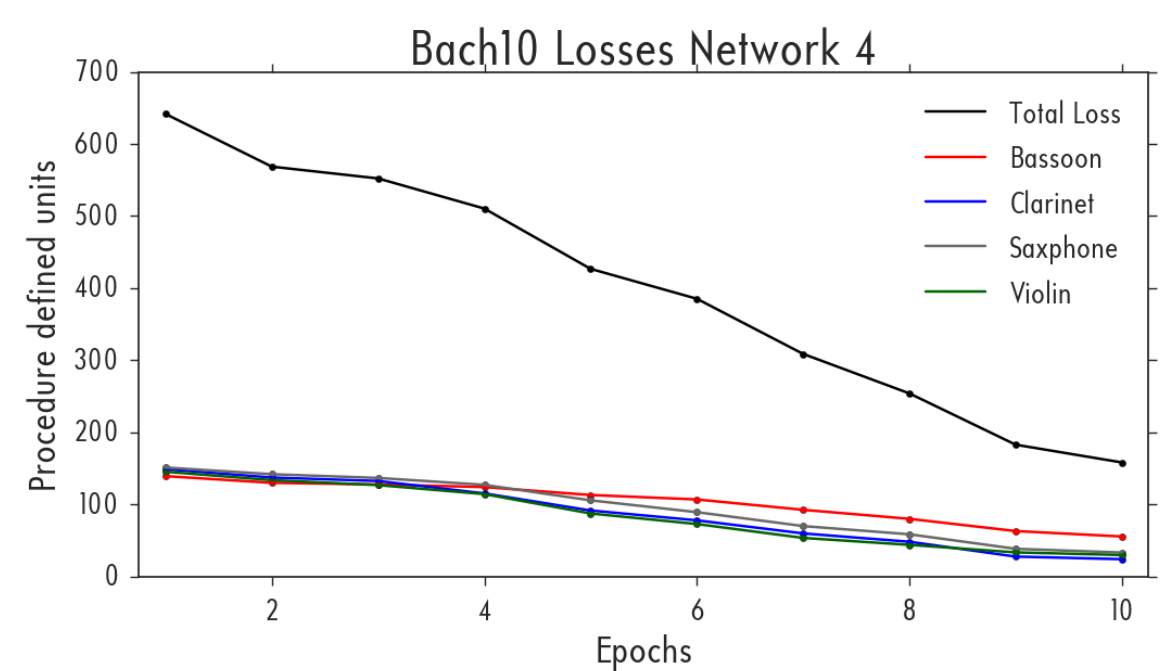
$$m_n(f) = \frac{|\hat{y}_n(f)|}{\sum_{n=1}^N |\hat{y}_n(f)|} \quad \tilde{y}_n(f) = m_n(f)x(f) \quad L_{sq} = \sum_{i=1}^N \|\tilde{y}_n - y_n\|^2$$

## Varying Spectrogram Time Contexts

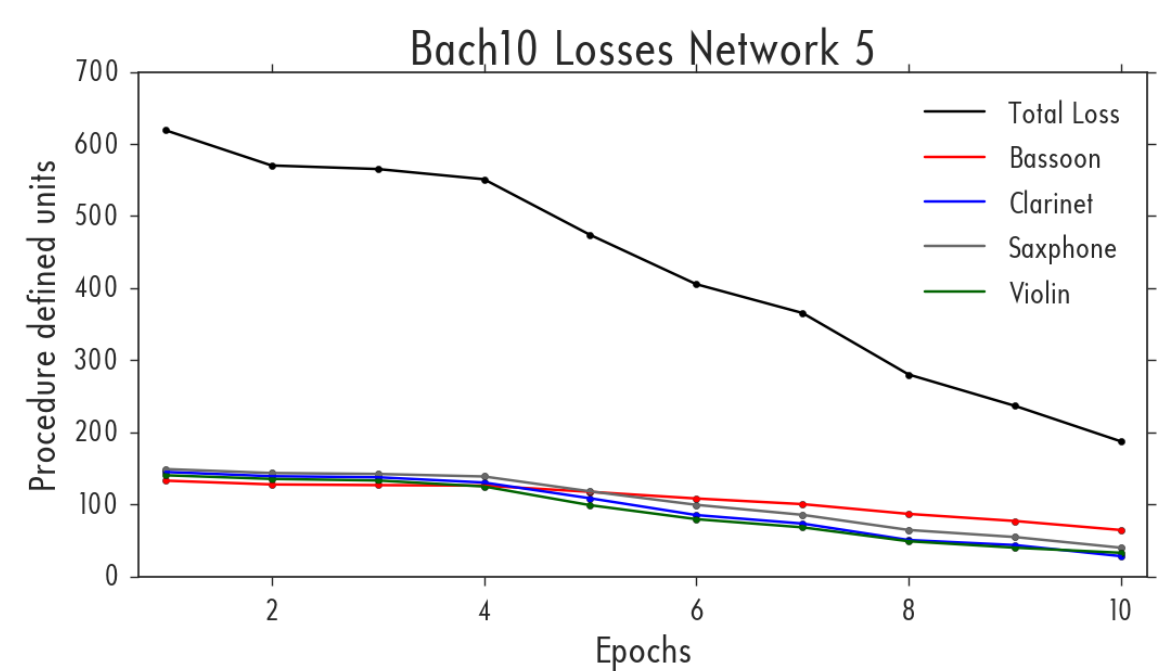


## Varying Fully Connected Layer Units

Network 4 Modification  
512 units fully connected layer

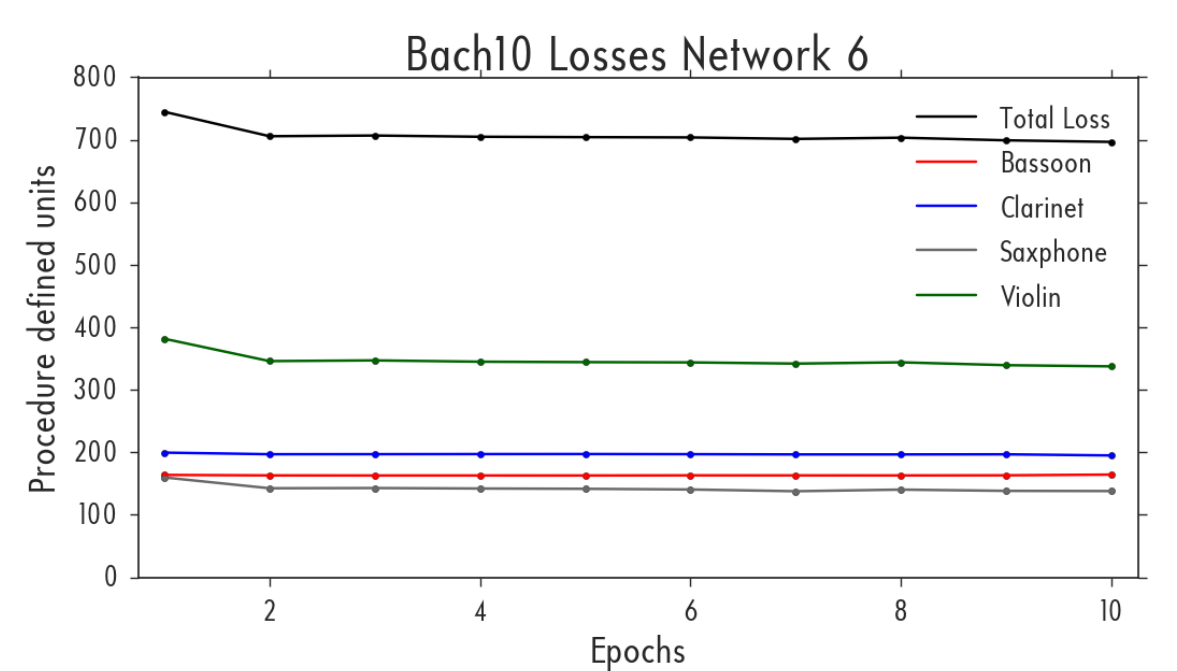


Network 5 Modification  
128 units fully connected layer



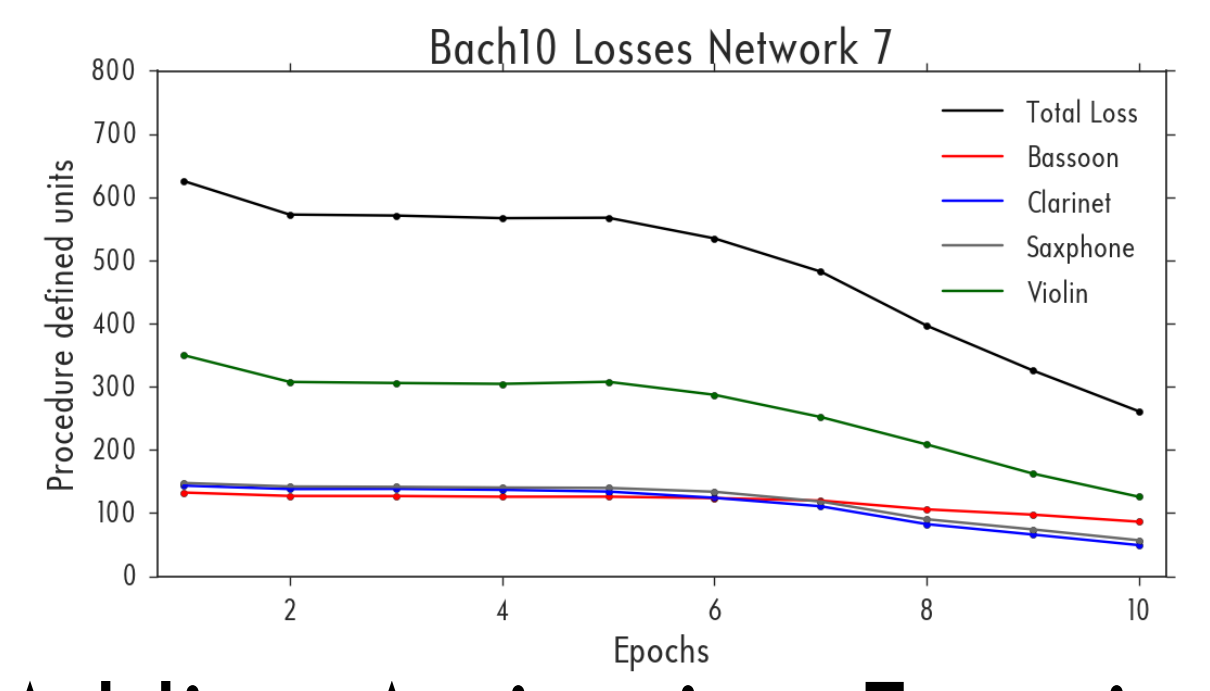
## Novel Architecture

Network 6 Modification  
Third convolutional/deconvolutional layer



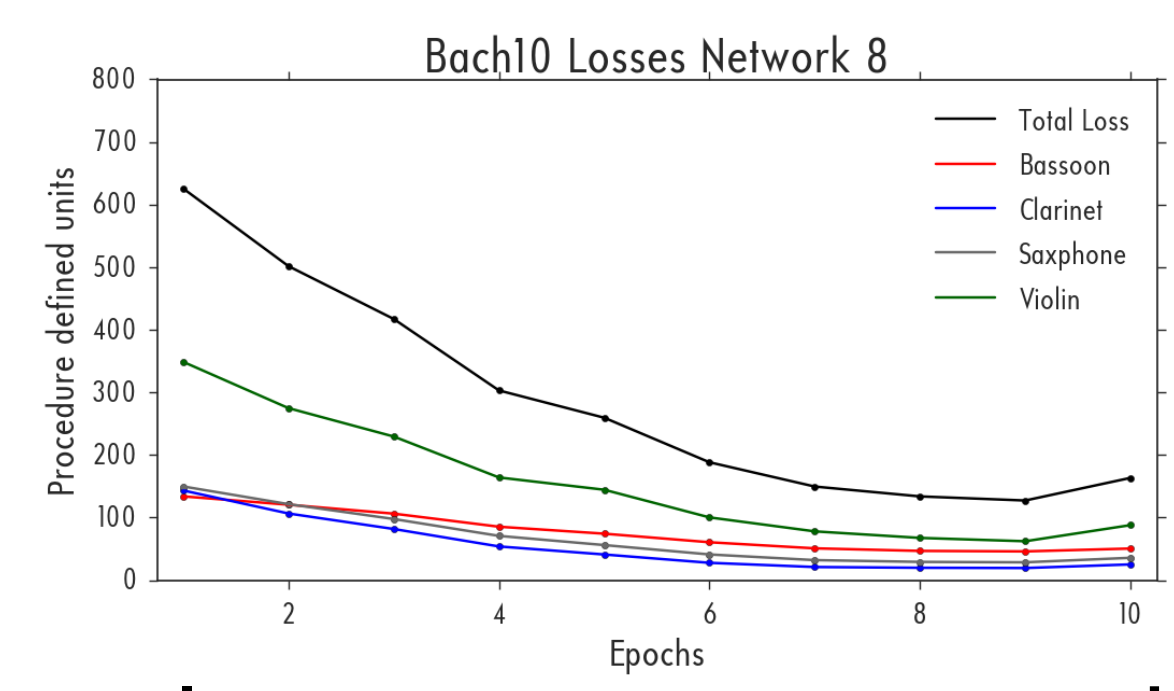
## Modifying Conv Layer Filters

Network 7 Modification  
64 convolution filters per layer



## Adding Activation Function

Network 8 Modification  
ReLU transfer function



## Conclusions & Future Works

All network variants learned the training set, as evidenced by decreasing, stabilizing learning curves. Surprisingly, network performance varied inversely with the time context of the spectrograms, suggesting time-dependent features only exist on a small scale in the dataset. Additionally, performance with respect to individual sources depended on the network architecture. Generally, the network learned the source signals of the violin and clarinet best, and of the bassoon the worst. However, when the convolutional processing was increased, the performance with respect to the bassoon and violin flipped.

Further investigation should extend these results to larger datasets, such as the DSD100, to test the scalability of our findings. The features of each source could be better understood by investigating gradients with respect to input.

## References

[1] Chandna, Prithish. "Monaural Audio Source Separation Using Deep Convolutional Neural Networks," 2017.  
[2] Grais, Emad M., and Mark D. Plumbley. "Single Channel Audio Source Separation Using Convolutional Denoising Autoencoders." arXiv:1703.08019 [Cs], March 23, 2017.  
\*Rodent Art Credit: Copyright (c) 2015 Etienne Ackermann \*\*Rodents were not used as a part of this study; however, they were petted intermittently for therapeutic purposes.