# Comparison of Neural Networks for Segmentation of Vocalizations

David Nicholson, Emory University. nickledave.github.io  🐦 nicholdav    ⬛ nickledave
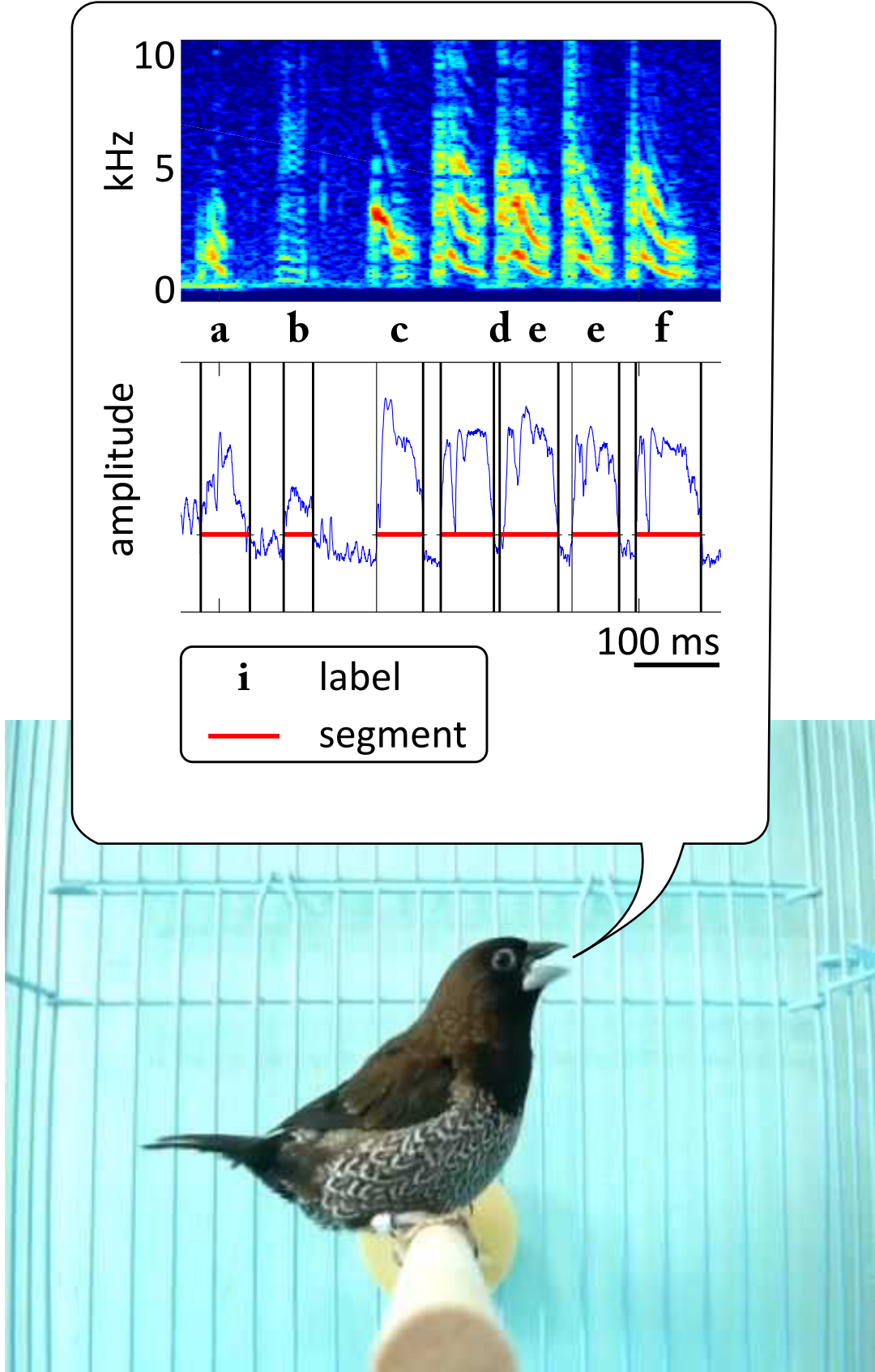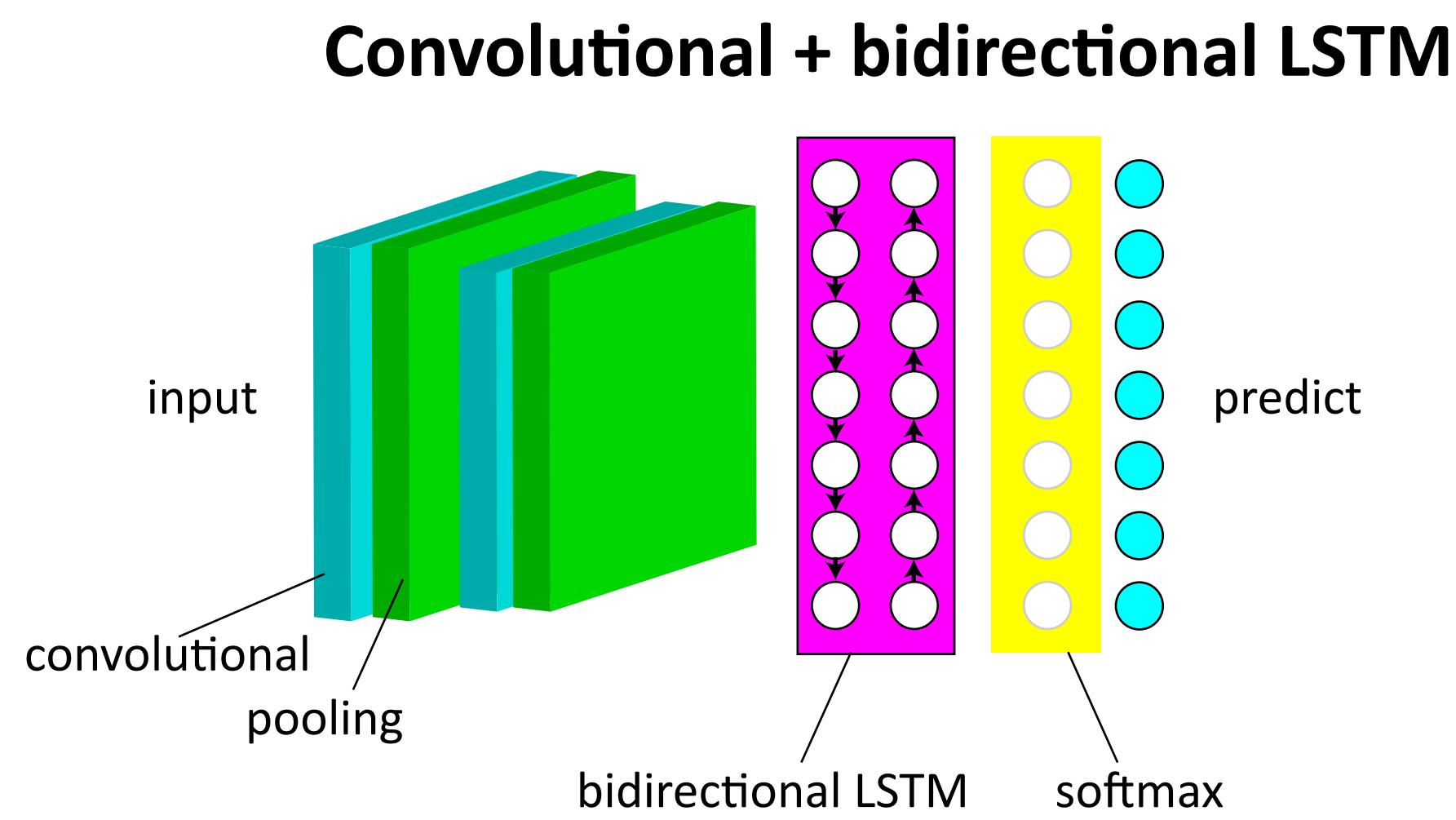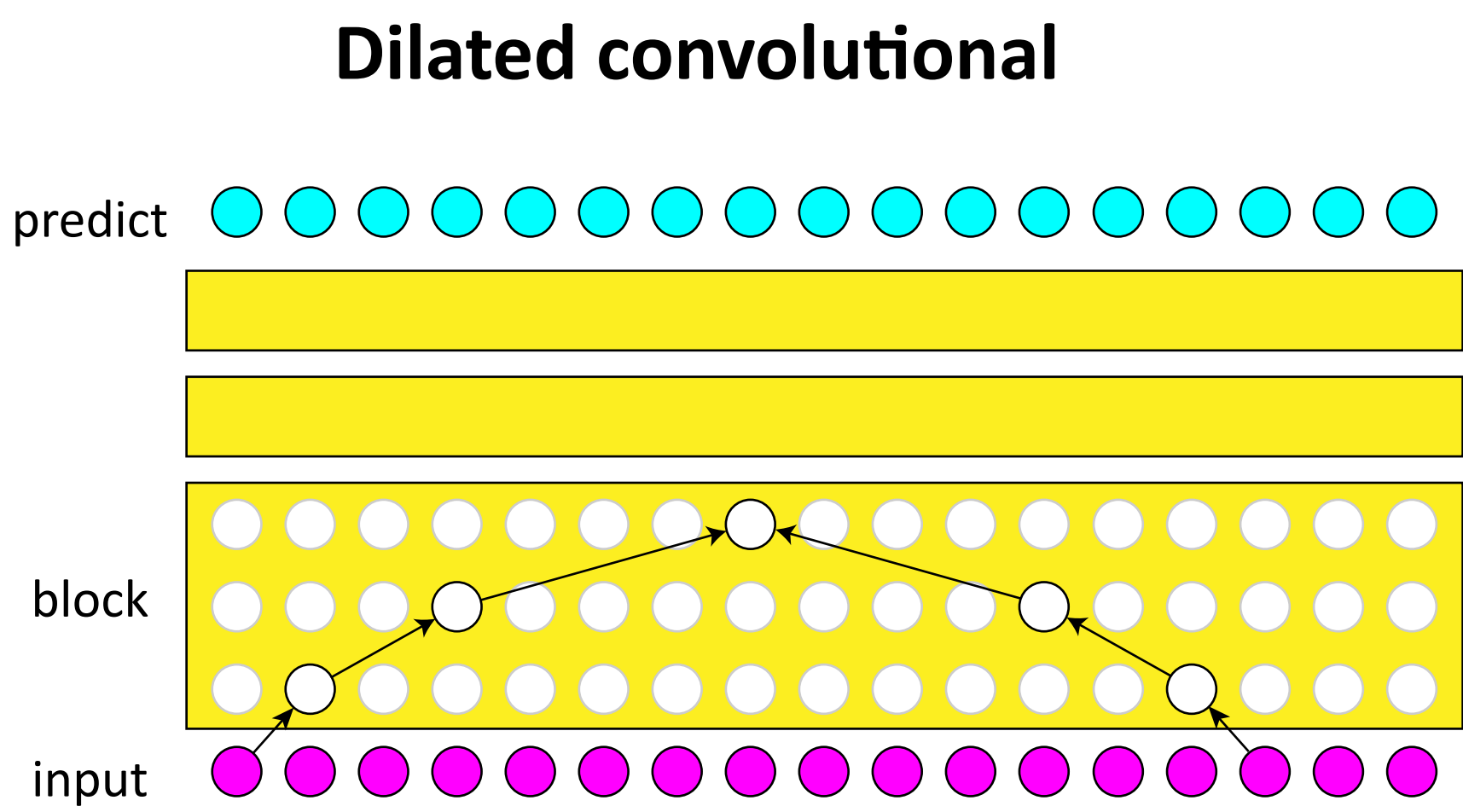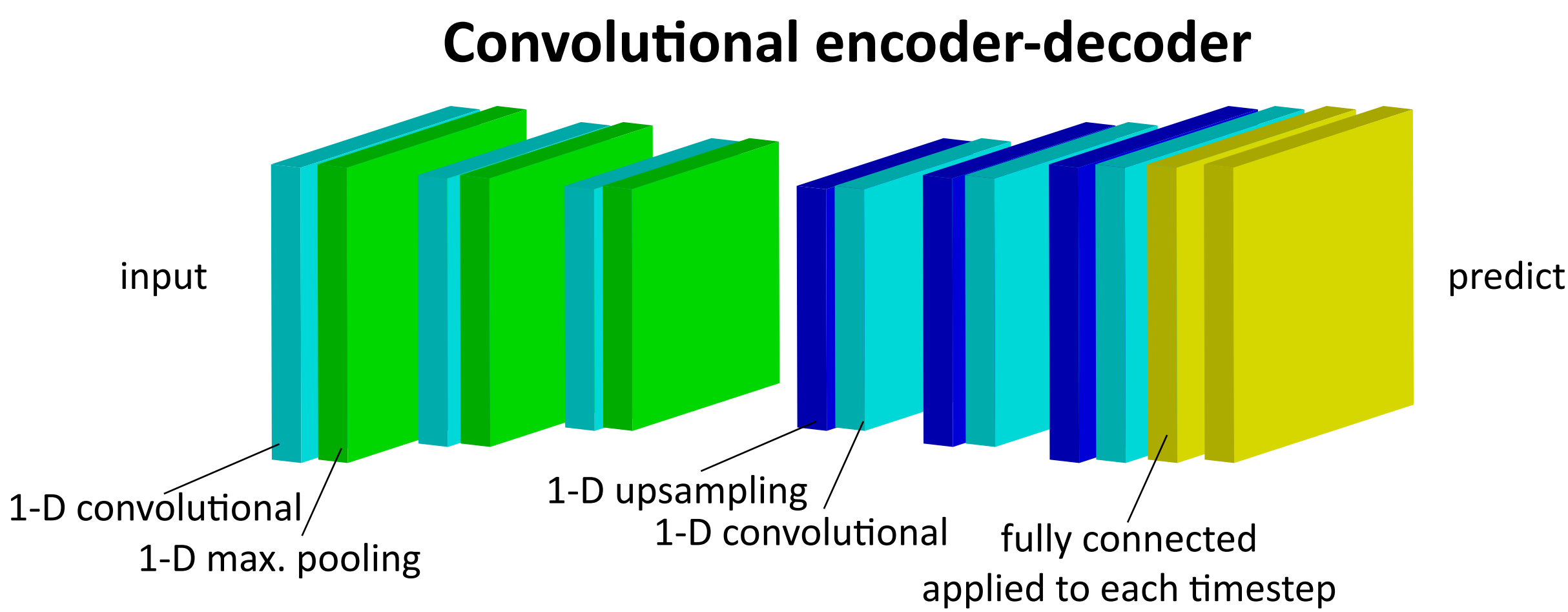
## Introduction

• Scientists study how songbirds learn and produce their song to understand how our brains learn and control behaviors like speech

• To get results from experiments with songbirds, scientists **segment** song into elements called *syllables*, and then **label** these syllables to extract acoustic parameters such as pitch

• Supervised machine learning can automate labeling *but fails when noise or experimental conditions impair segmentation*

• Here I compare neural networks that both segment song and classify syllables
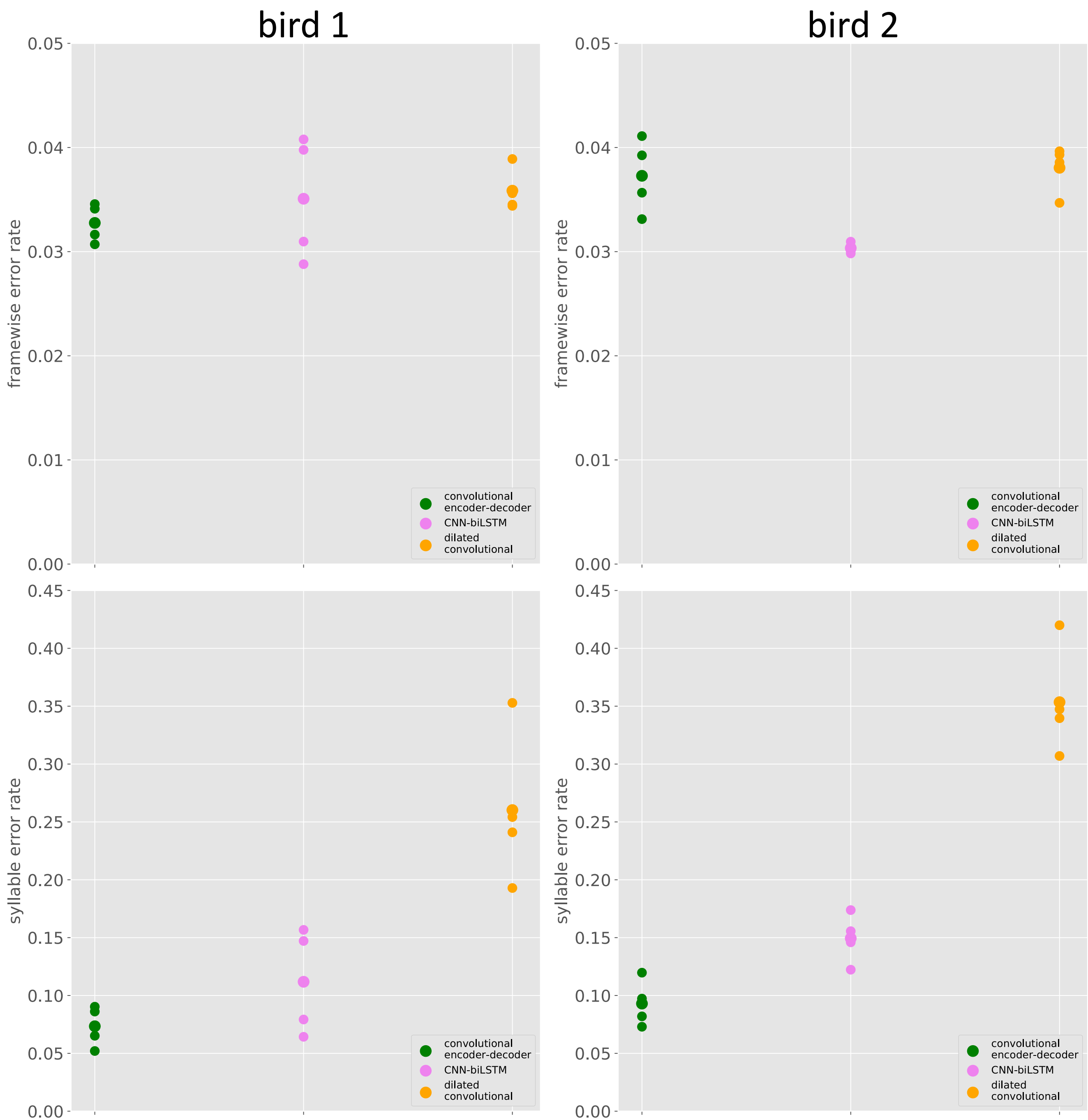


## Methods

• Compare three different neural network architectures:

### Convolutional encoder-decoder



input — predict
1-D convolutional
1-D max. pooling
1-D upsampling
1-D convolutional
fully connected applied to each timestep

### Dilated convolutional



predict
block
input

### Convolutional + bidirectional LSTM



input — predict
convolutional
pooling
bidirectional LSTM
softmax

• Benchmark on song from public repositories:
https://figshare.com/articles/Bengalese_Finch_song_repository/4805749
https://figshare.com/articles/BirdsongRecognition/3470165

## Results



• Framewise error not significantly different between the three network types

• Syllable error:
encoder-decoder < CNN-bidirectional LSTM < dilated convolutional

## Conclusion

• Fully convolutional networks are competitive with recurrent networks, giving roughly the same framewise error rate, and in the case of the encoder-decoder network, lower syllable error rate than the CNN-bidirectional LSTM

• *and* fully convolutional networks require much less time to train (1 hour for encoder-decoder and dilated convolutional versus 3-4 hours for CNN-bidirectional LSTM)

## References

René, Colin Lea Michael D. Flynn, and Vidal Austin Reiter Gregory D. Hager. "Temporal convolutional networks for action segmentation and detection." (2017).

https://github.com/colincsl/TemporalConvolutionalNetworks

https://github.com/yardencsGitHub/tf_syllable_segmentation_annotation

https://github.com/NickleDave/tf_syllable_segmentation_annotation/network