# Machine Learning Overview

CS Clinic Auburn-USDA 2025

# Task Description

Given a file containing time-series EPG data, predict a waveform type for each timestep of the series.

Input: ~70 EPG recordings, some with multiple probes

- **EPG Signal [volts] vs Time**
- Measurement settings [current type, resistance, voltage]

Output

- Waveform Type vs Time (ex: NP, NP, NP, …, J, J, … K, …, L, …, Z, … NP)

# Approaches

We are viewing this as a sequence-to-sequence prediction task.

$$[7.4, 8.9, 9.0, 9.1, \dots ] \rightarrow [NP, NP, J, J, \dots ]$$

To make the task easier, we have a variety of techniques we can utilize, including

- Preprocessing
- Windowing
- Feature extraction
- Data augmentation

Finally, we can use post-processing to make the output sequences easier for a human to edit.
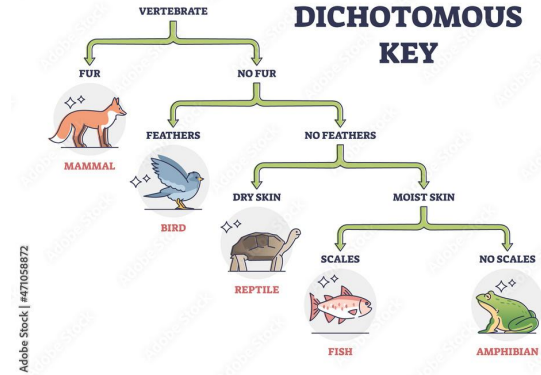
# Models

# Random Forests

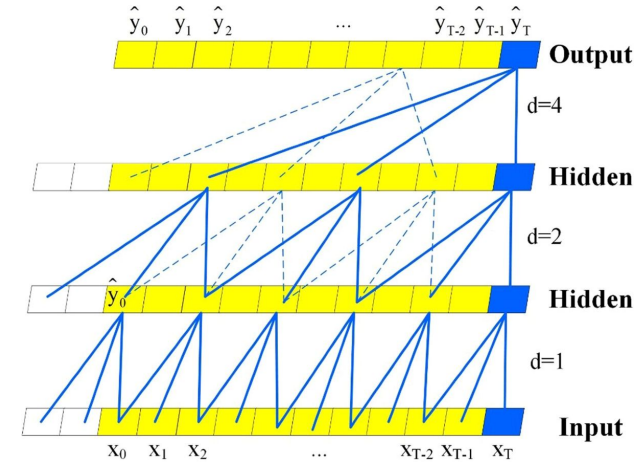This method largely replicates that of Willett et al.

Decision Tree: like a dichotomous key, but automatically generated for a dataset to maximize the likelihood it makes a correct classification.

Random Forests Model: a large collection of decision trees all trained slightly differently from each other (this gives them some diversity in the way they "think"). We poll all of them and we go with the majority to make our decision.

# Temporal Convolutional Network (TCN)

While there are many flavors of TCN, our TCN learns to extracts features from the input data in short segments and then feeds them into a neural network. It only processes data at a granular level.
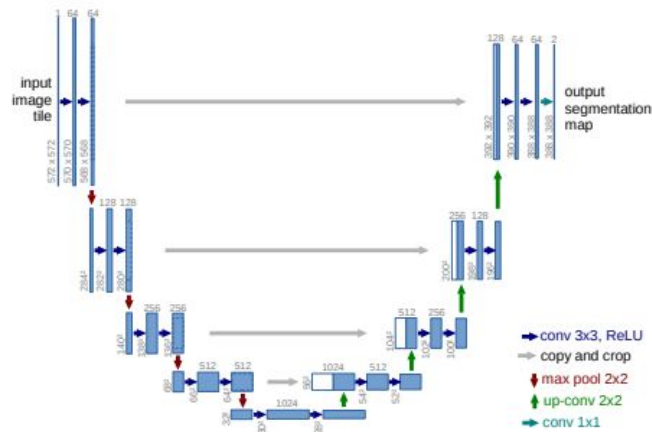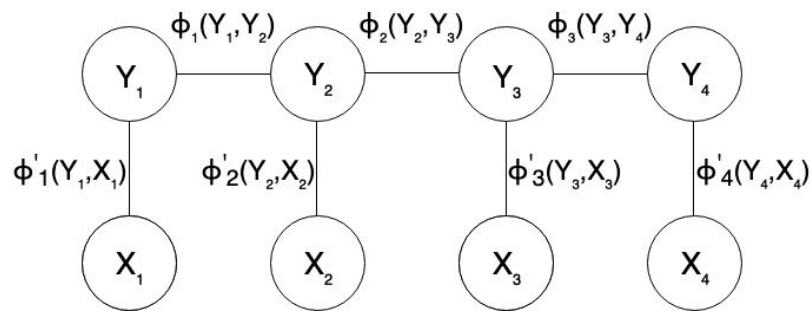
# UNet

From Ronneberger, et. al 2015

Originally used for labeling medical images, UNet is a neural network that applies learned filters to the data. Along the way, the data gets compressed and then uncompressed.

This allows the data to be processed at both a high and low level before it is labelled.

# Conditional Random Fields (CRF)

Typically used for labeling text data and images, CRFs allow us to take outputs from any of the previous models and make sure they are coherent given the transitions.
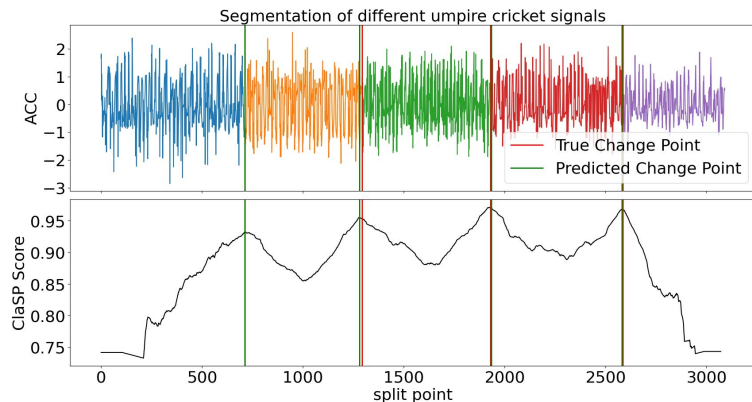
# ClaSP + HMMs

Instead of labeling each timestep individually (or in windows), we could first determine which intervals all contain the same waveform type and then label those intervals.

Based on what we have gathered, this is more akin to how humans label data. However, for these methods to work well we need to know how many segments there should be, which is often not possible.

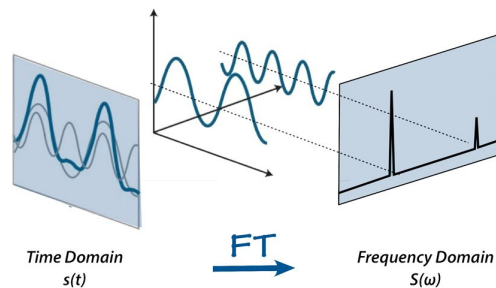# Pre- and post-processing

# Feature Engineering

Fourier Transform

- Extract the frequencies that make up the EPG signal

ClaSP

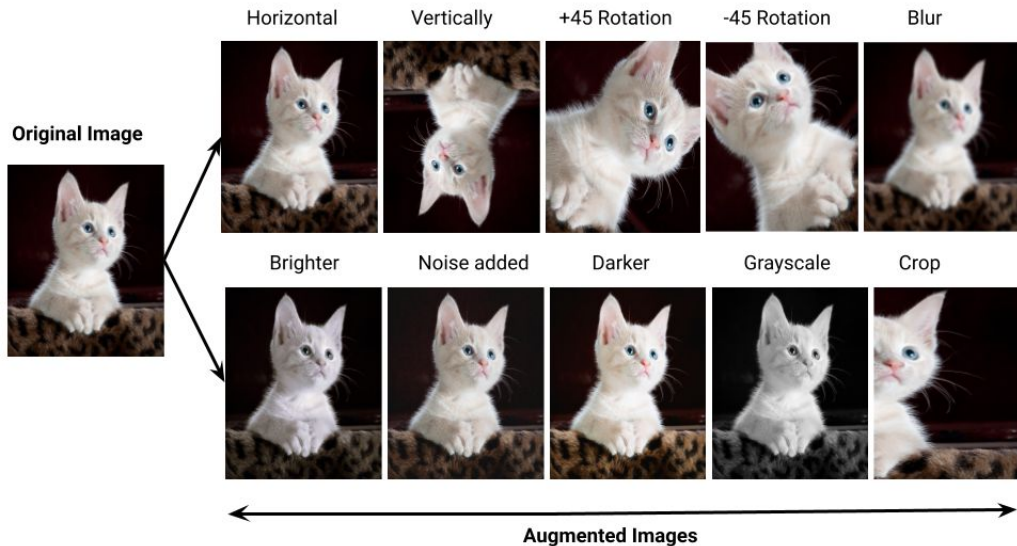The filters in the CNNs can be thought of as learning to extract meaningful features from the EPG signal

Fractal Dimension, Hurst Exponent, Etc…

- These all measure the "complexity" of data. Other papers have used them but we have not found them to improve performance



Time Domain
s(t)

FT

Frequency Domain
S(ω)

# Data Augmentation
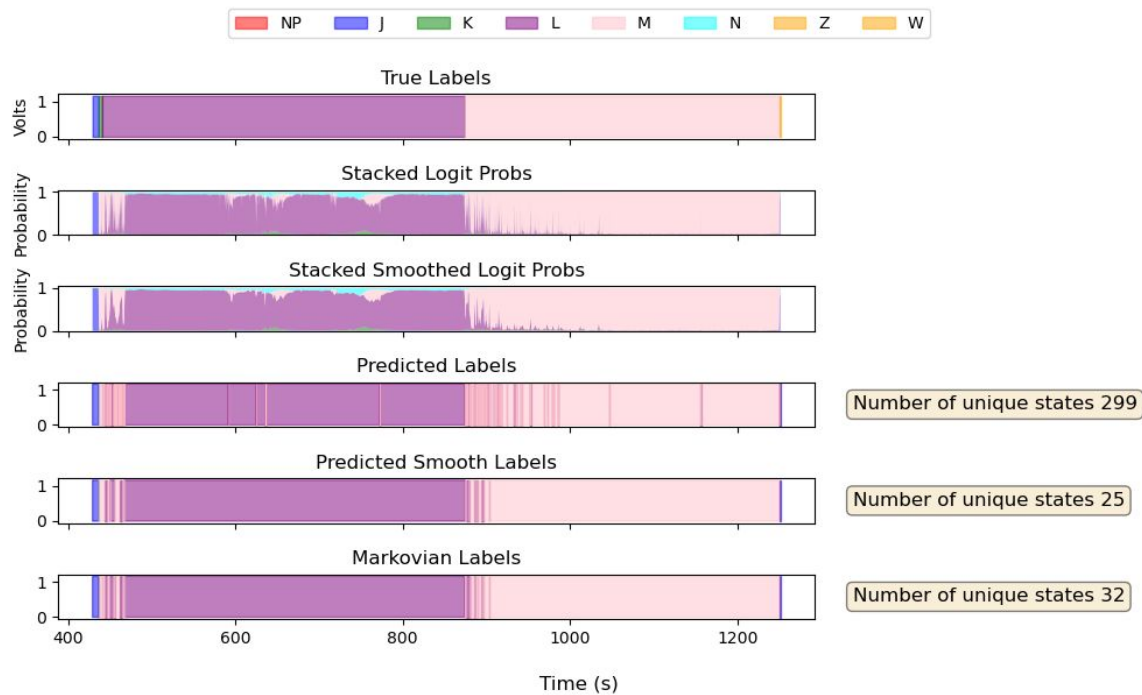
- Augment_concat_self
- Augment_concat_other
- Augment_warp_overall
- Augment_warp_by_state
- Augment_change_amplitude_overall
- Augment_change_amplitude_by_state
- Augment_noise_voltage_overall
- Augment_noise_voltage_by_state
- Augment_franken

Combine these together!

# Post-processing

- Hidden Markov Model (later, add semi)
- Smoothing

# Questions?