# Generating Synthetic Data
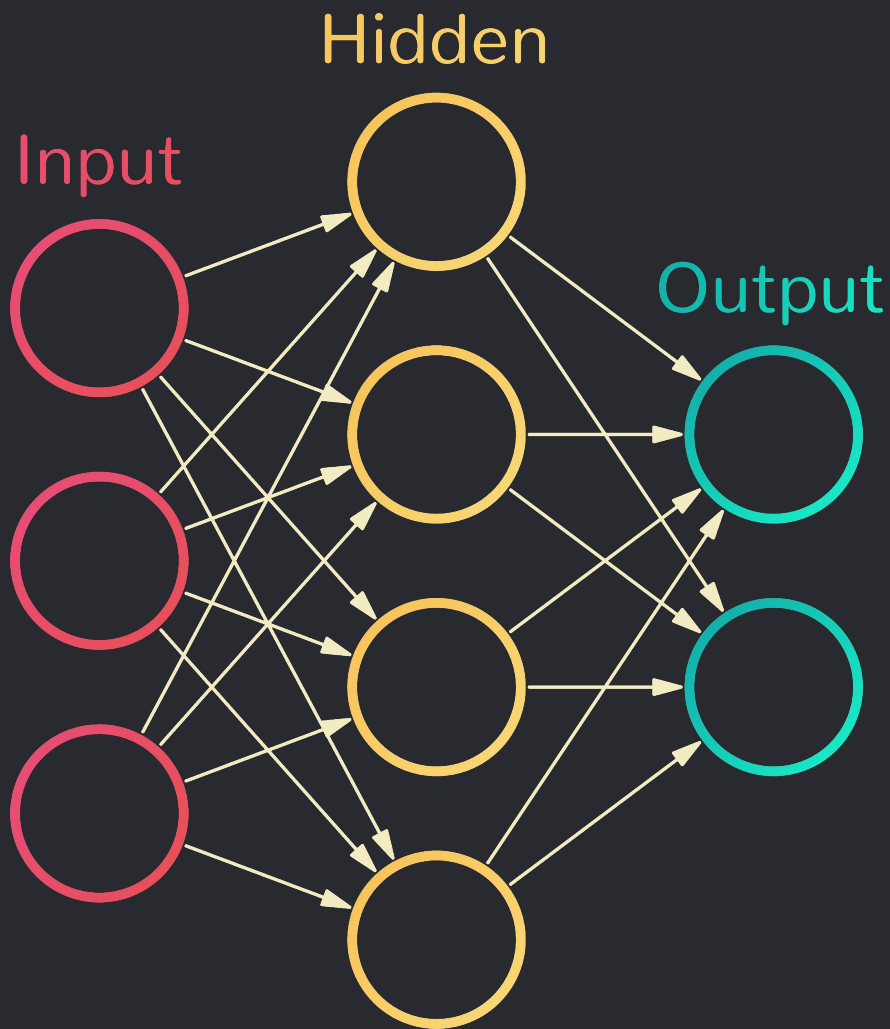
Rauf Verdiev, Roman Kotovich, Vera Terenteva, Elizaveta Terekhova, Laura Avinyó

SMTB

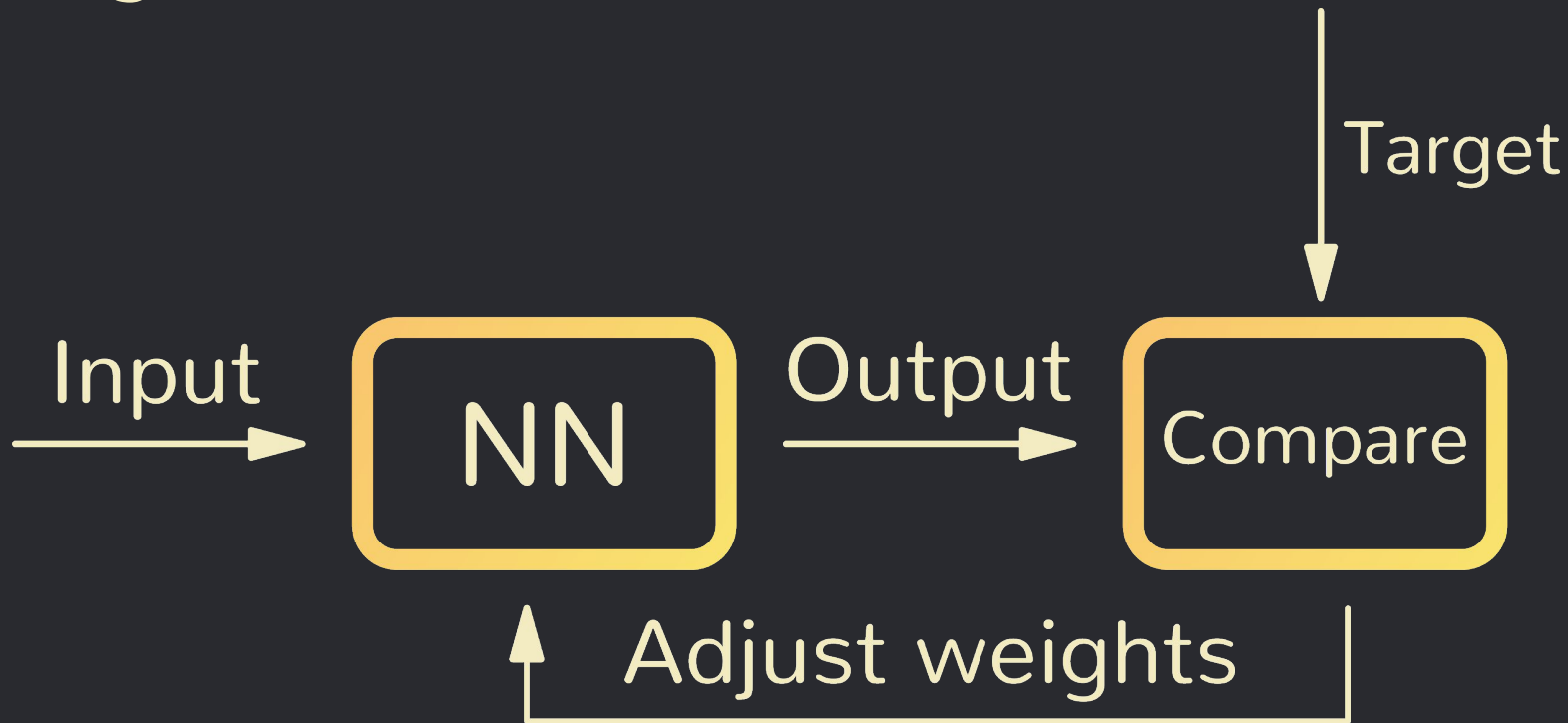# Project goal

- Generate synthetic data
    - HMM
    - GANs
- Analyze generated sequences
    - Zm
    - Decision tree
    - Shannon Entropy

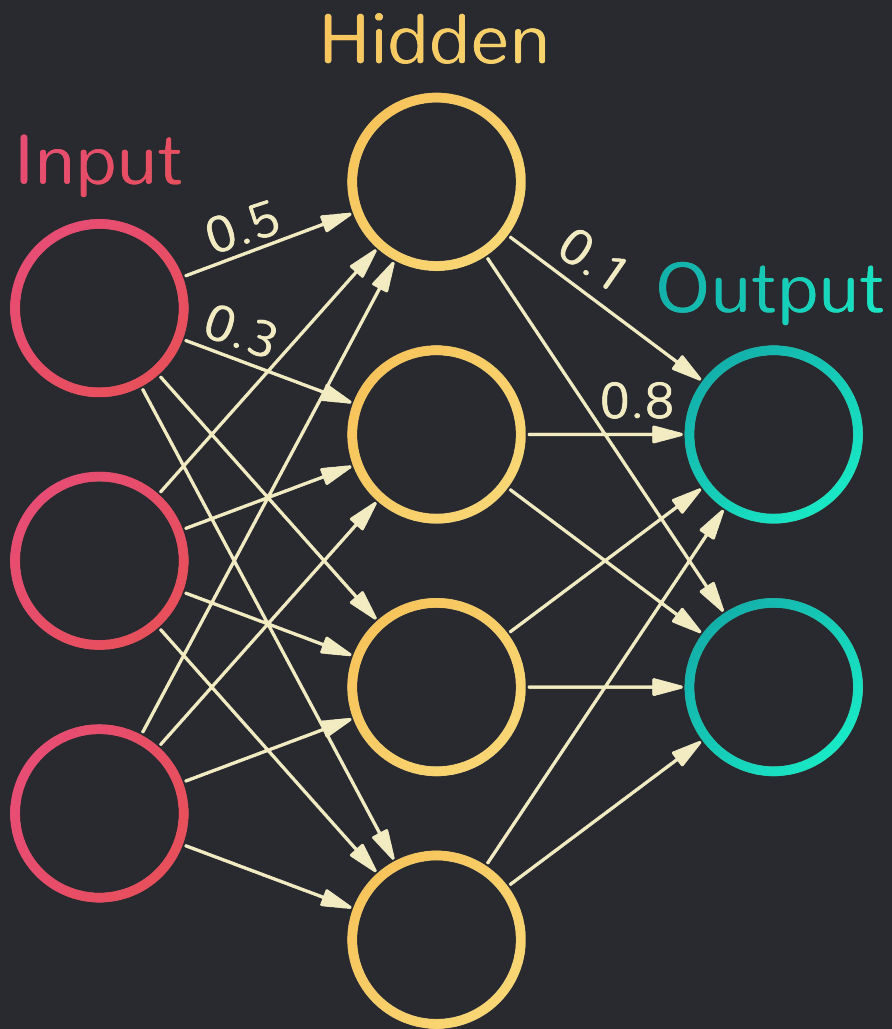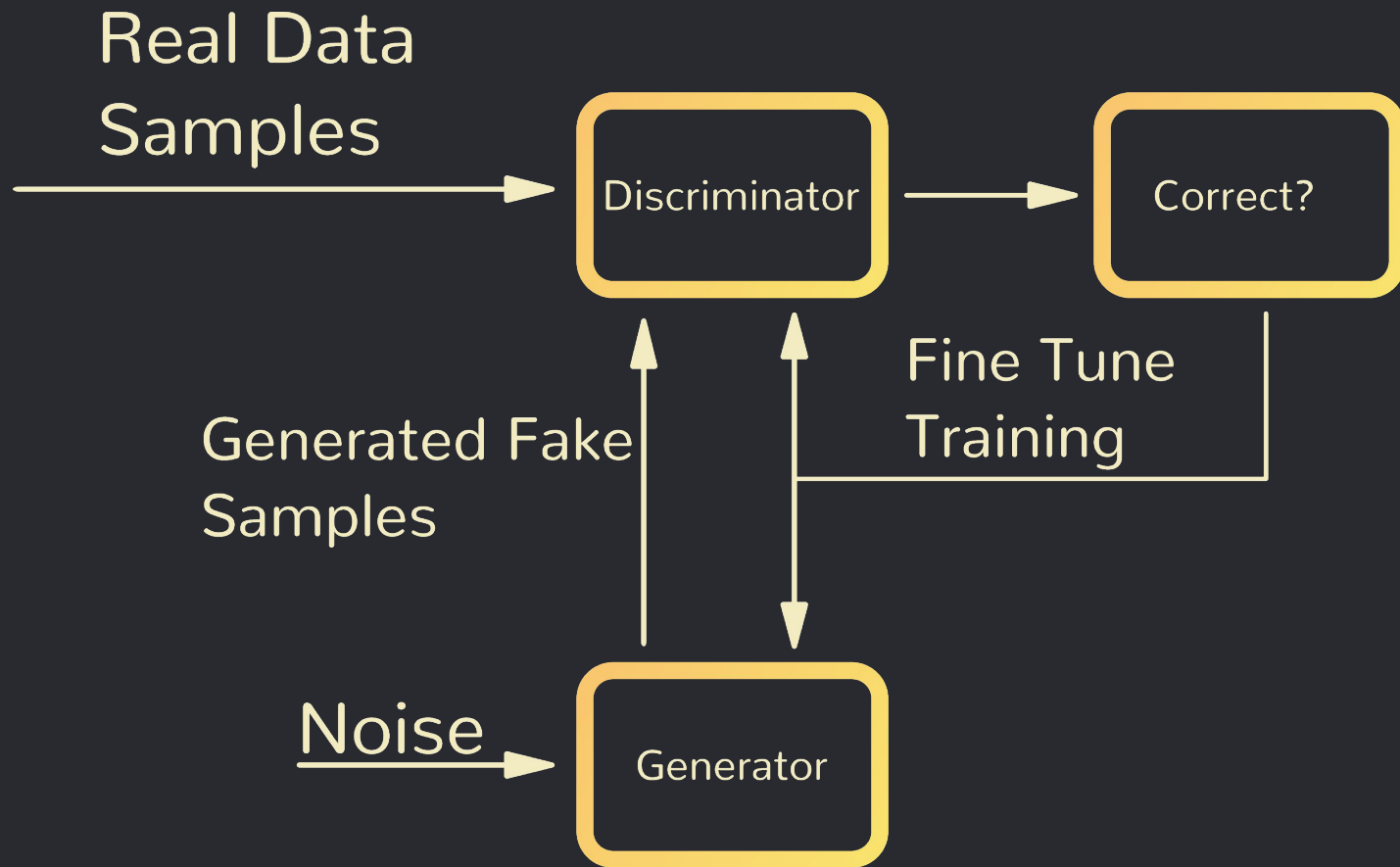Why did we do this?

# Training

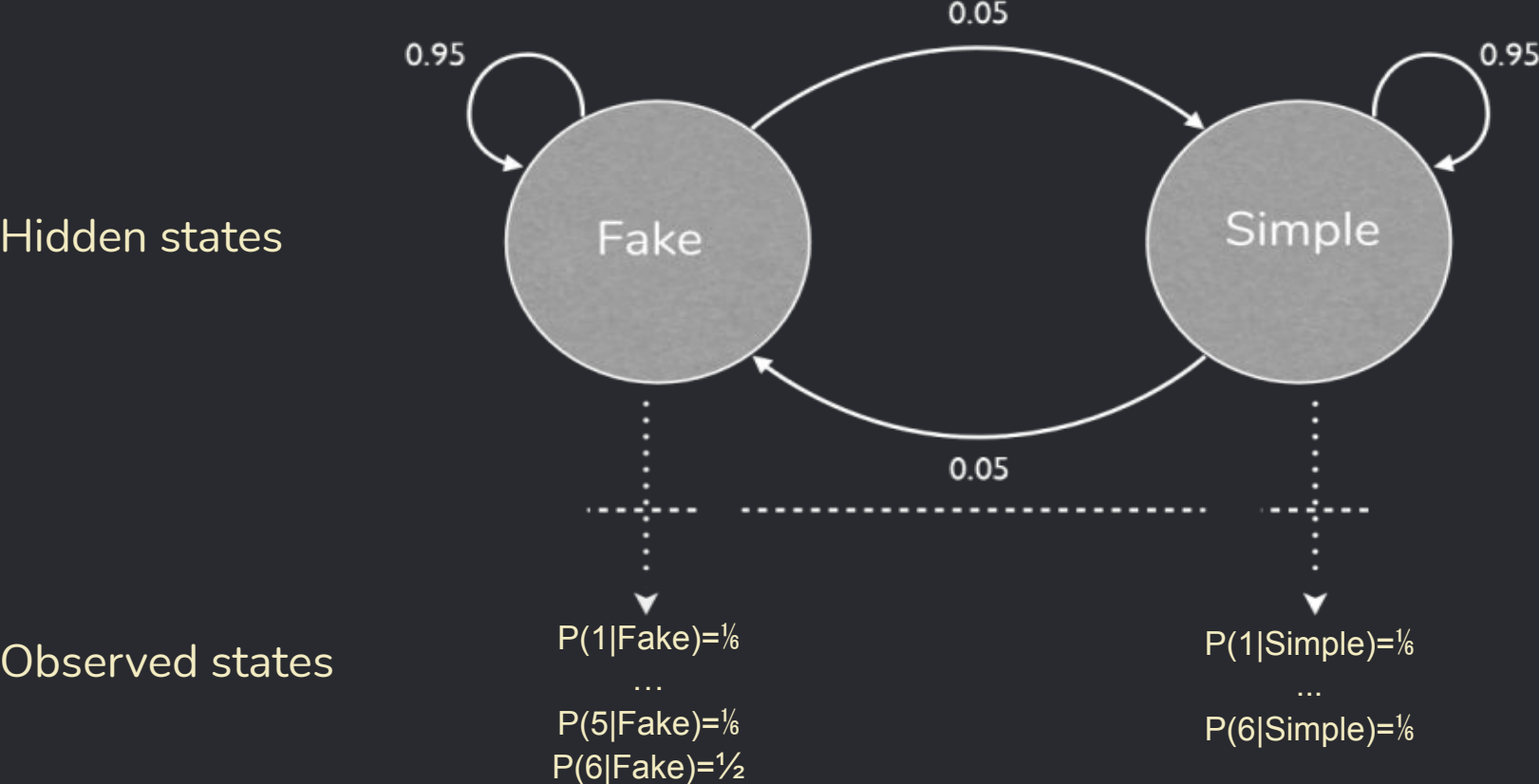Input → **NN** → Output → **Compare** ← Target

Adjust weights

# GANs

**G**enerative
**A**dversarial
**N**eural Networks

Real Data
Samples

Discriminator

Correct?

Generated Fake
Samples

Fine Tune
Training

Noise

Generator

# Hidden Markov Models (HMM)

# Hidden Markov Models (HMM)

1 2 4 5 1 3 2 4 5 6 3 2 4     P = E(1|T)*T(T|T)*E(2|T)*T(T|T)*...*E(3|T)*T(T|F)*E(2|F)*T(F|F)*...

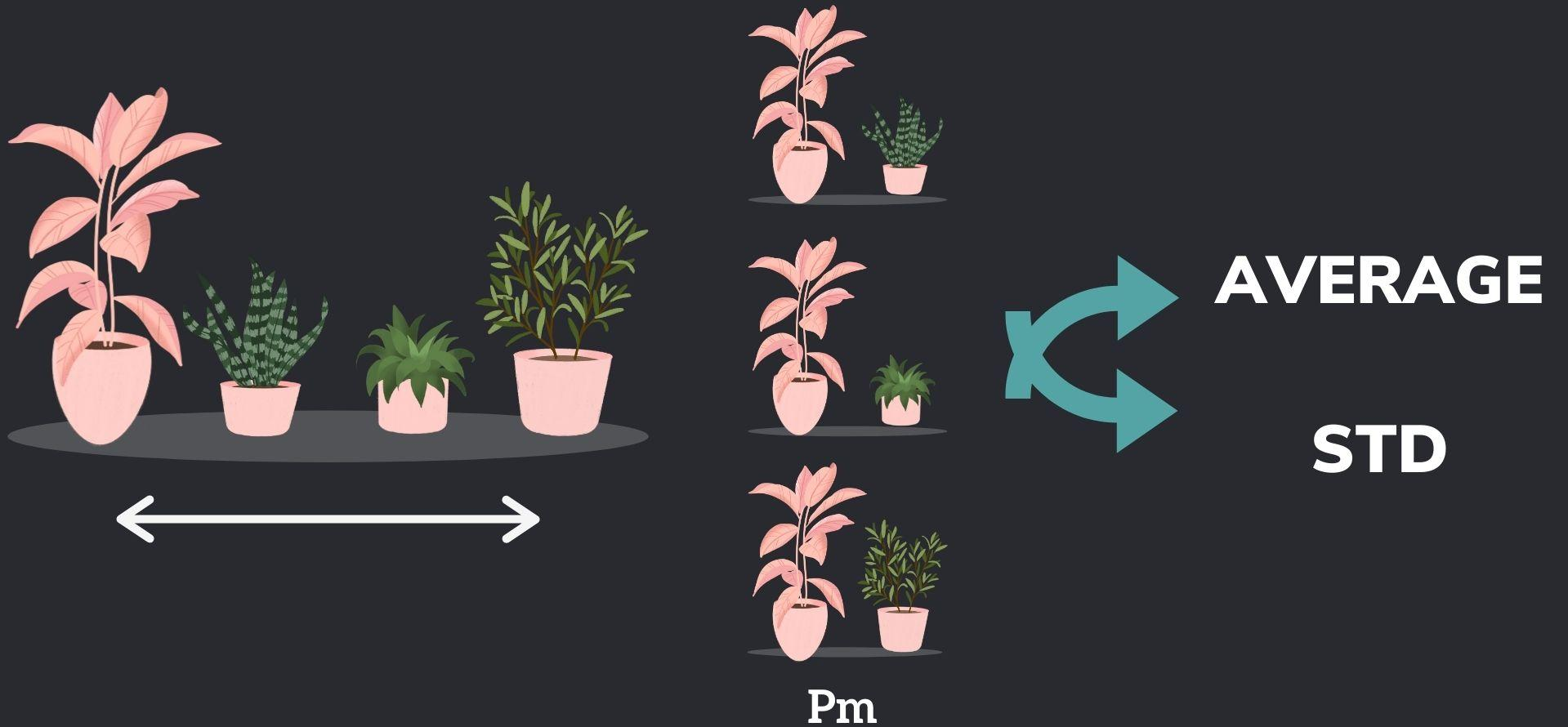3 1 5 2 5 6 3 2 4 1 5 6 6     P = E(3|T)*T(T|F)*E(1|F)*T(F|T)*E(5|T)*T(T|F)*E(2|F)*T(F|T)*...

We can generate some sequences!

Emission match = [ 1: {A=0.01; B=0.06; ...}, 2:{A=0.07; B=0.05; ...}, ..., 39:{A=0.01; B=0.09; ...} ]

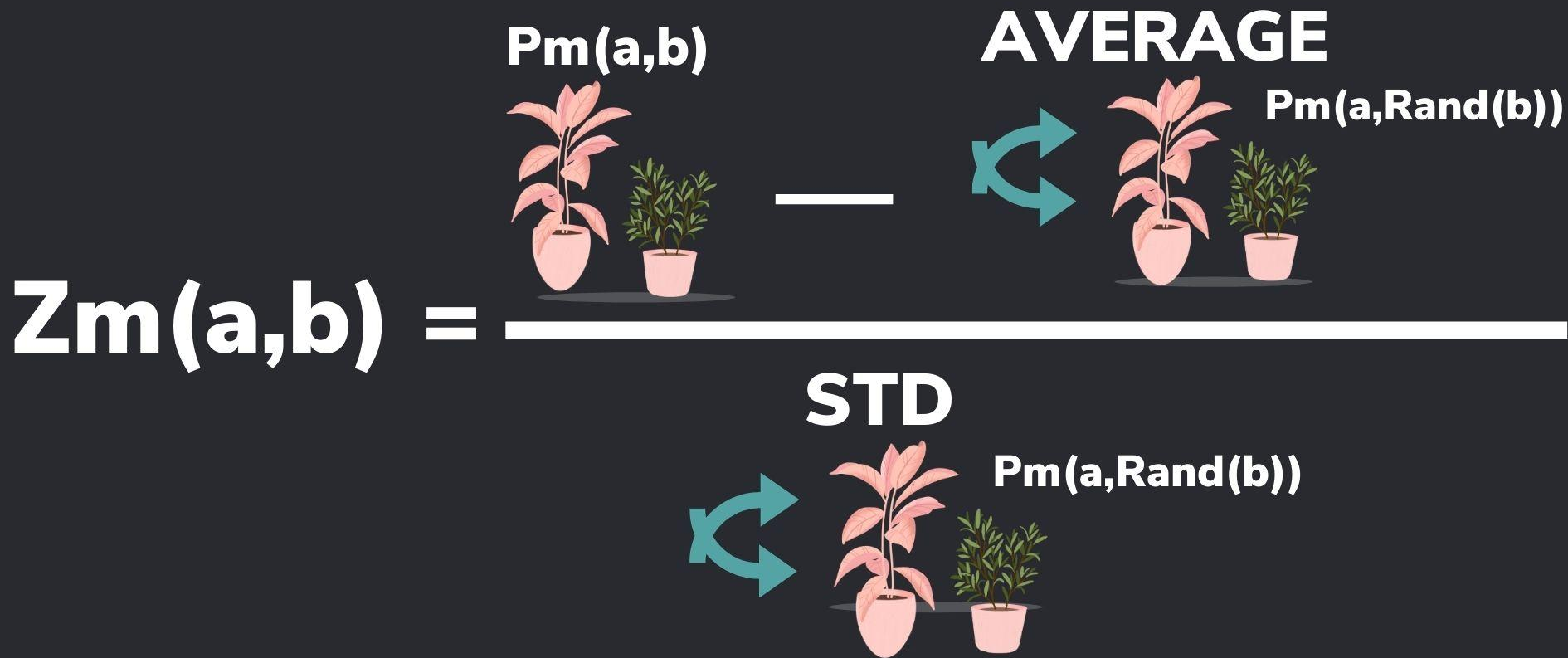Insertion match = [ 1: {A=0.03; B=0.05; ...}, 2:{A=0.04; B=0.06; ...}, ..., 39:{A=0.07; B=0.05; ...} ]

Transition = [ 1: { m->m: 0.92; m->i: 0.10; m->d: 0.11; i->m: 0.70; i->i: 0.20; d->m: 0.90; d->d: 0.01}, ... ]
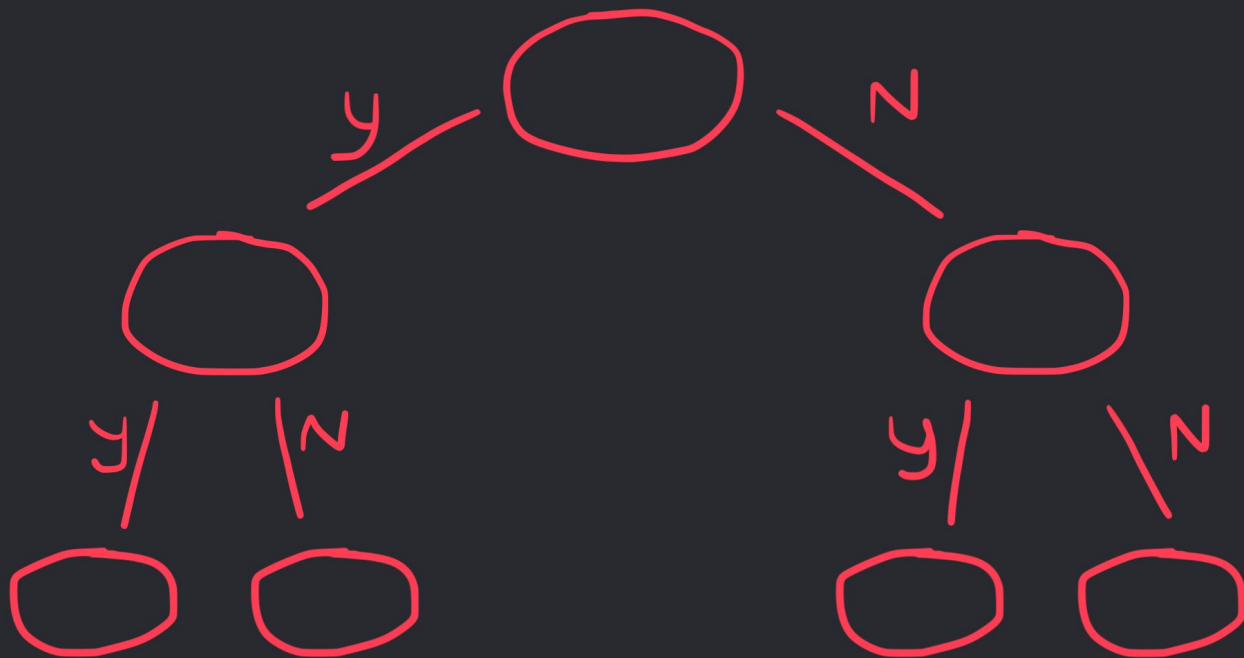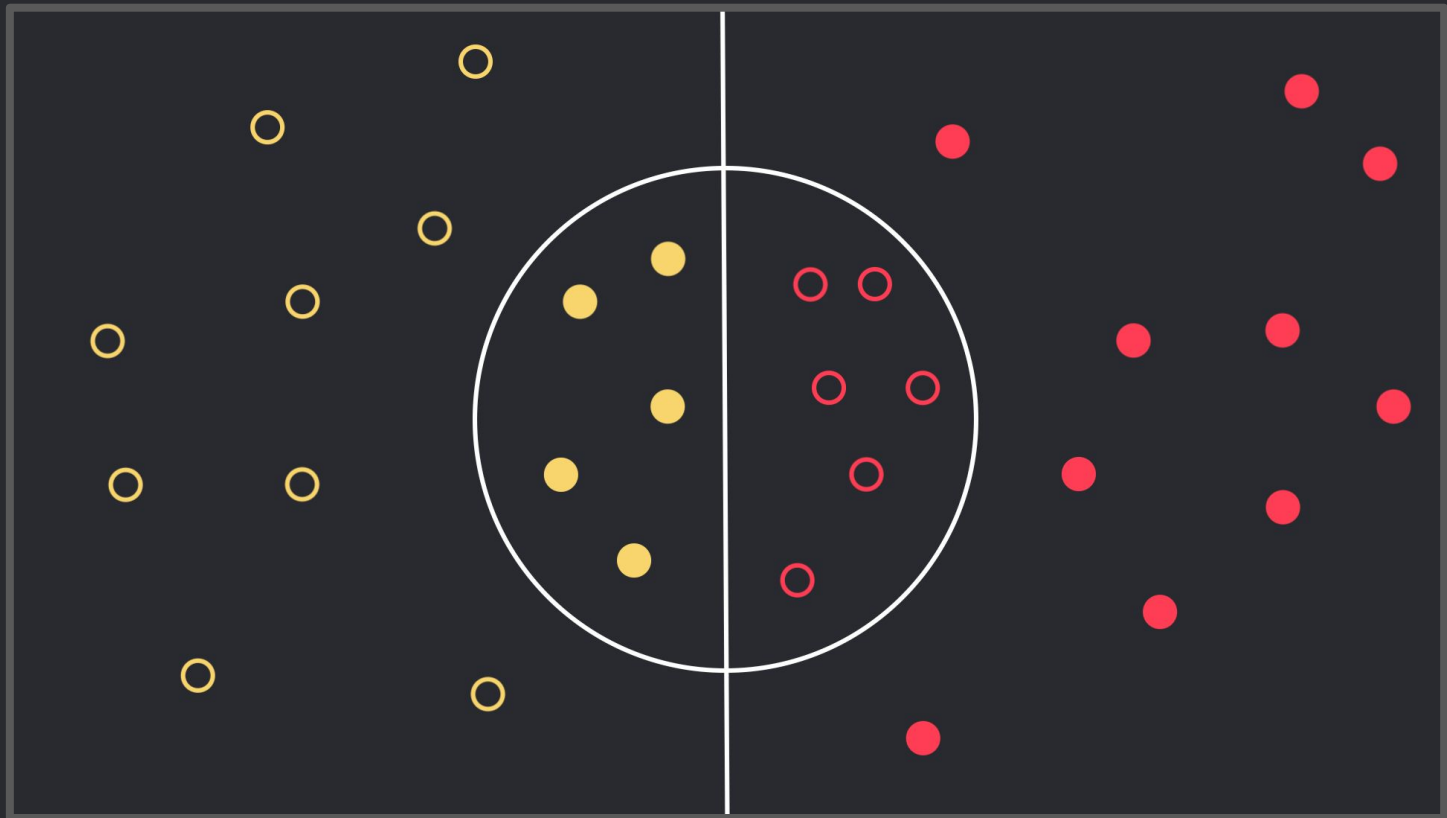
HOW WE CHECK DATA?

Pm

AVERAGE

STD

# How to grow a decision tree?

Accuracy, Precision and Recall

# Shannon Entropy

The Shannon entropy can measure the uncertainty of a random process

Its lower values imply less uncertainty

For each column of multiple-sequence alignment (MSA) Shannon is calculated as:

$$s.e. = - \sum_{i=1}^{20} p(x_i) log_{20} p(x_i)$$

p(x$_i$) - is the probability to find the amino acid *i* at the column of MSA

In terms of columns in MSA its low values represent highly conservative positions, whereas high entropy shows diversity at certain position
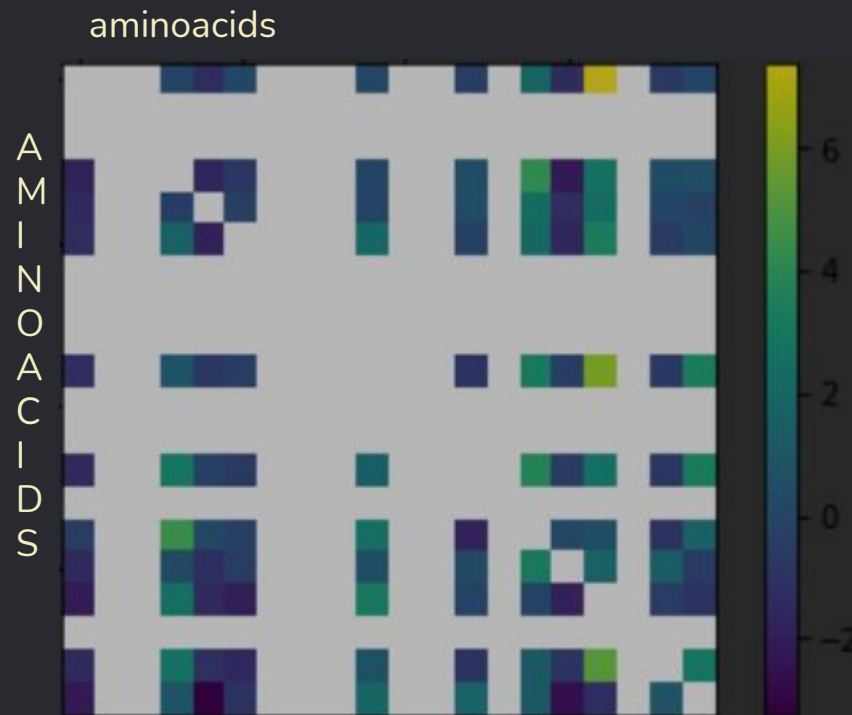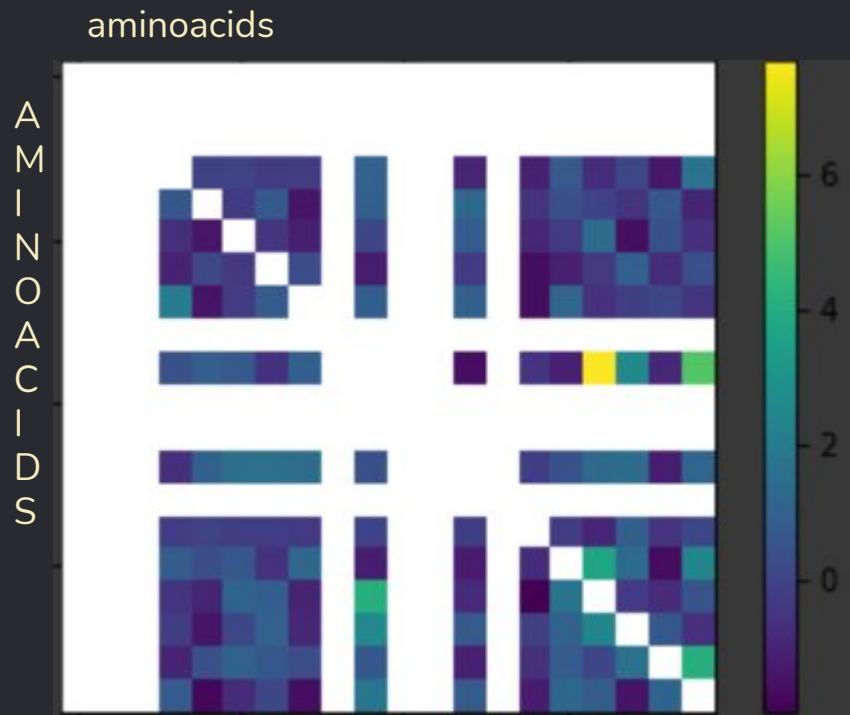
# Project goal

1.  Generate Synthetic data with the different methods
    a.  GANs
    b.  HMM
    c.  Random Baseline
2.  Test how good or bad the data is:
    a.  ML discriminator -> Overall quality
    b.  Shannon's Entropy -> Functional
    c.  Zm -> 3D structure
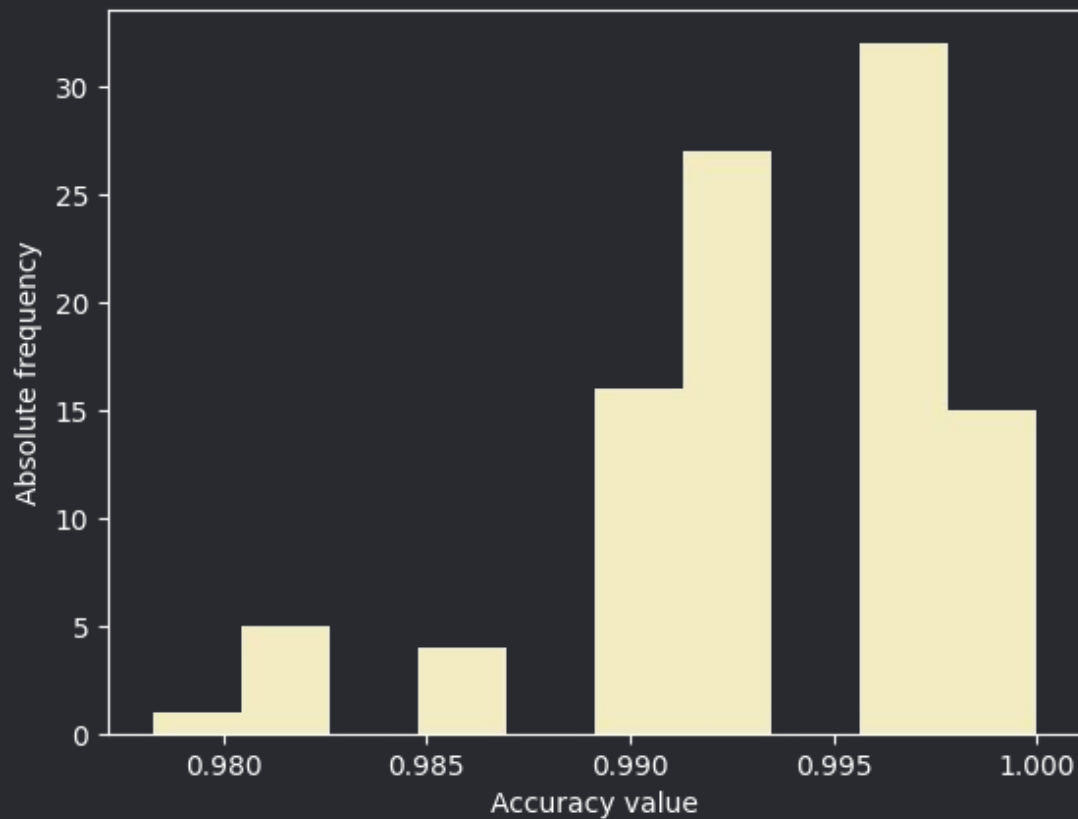
Why did we do this?
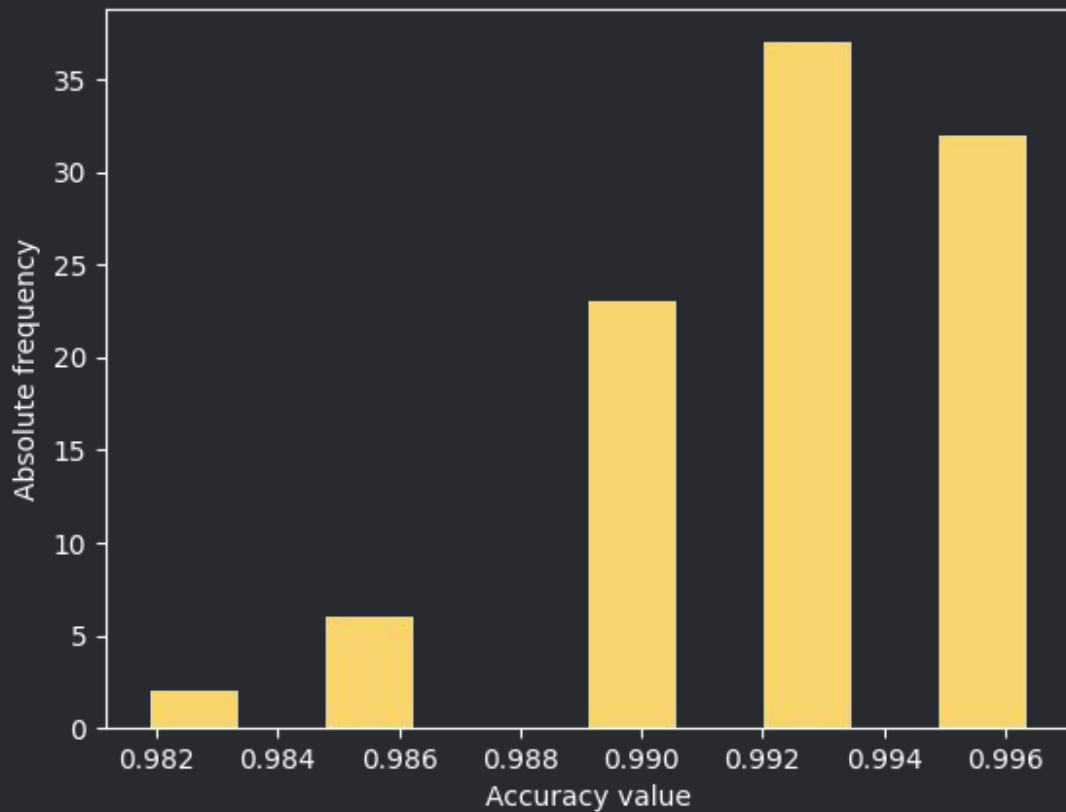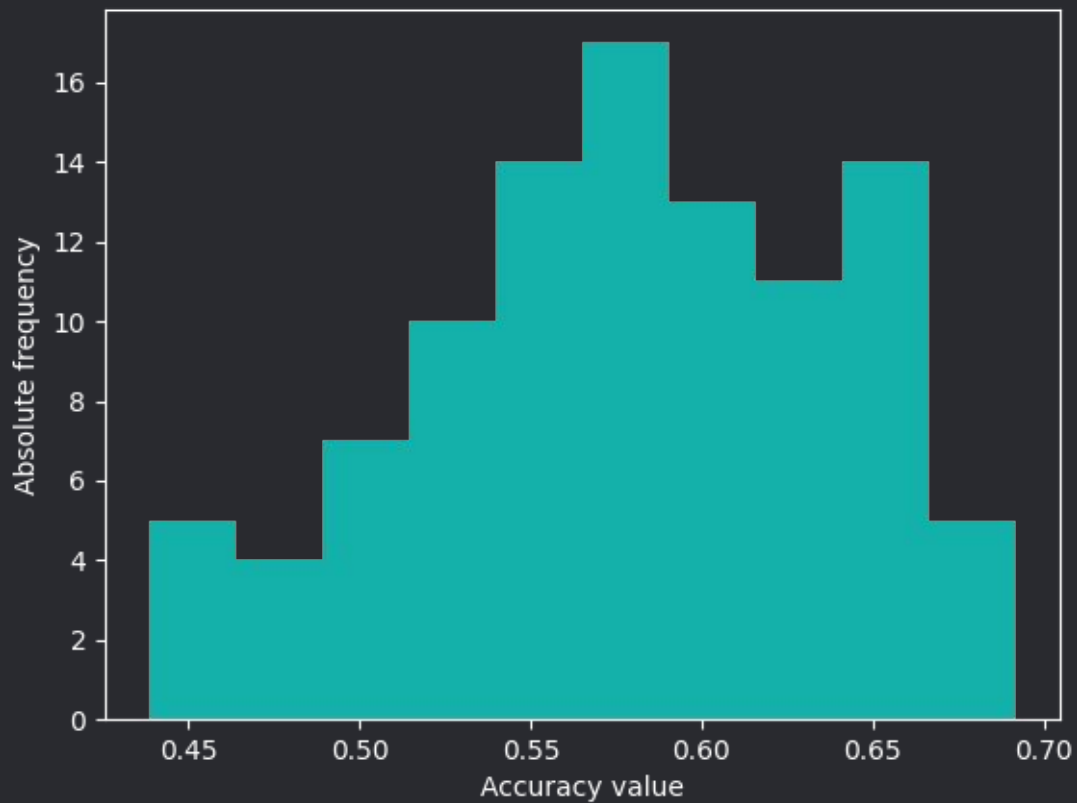
# Results

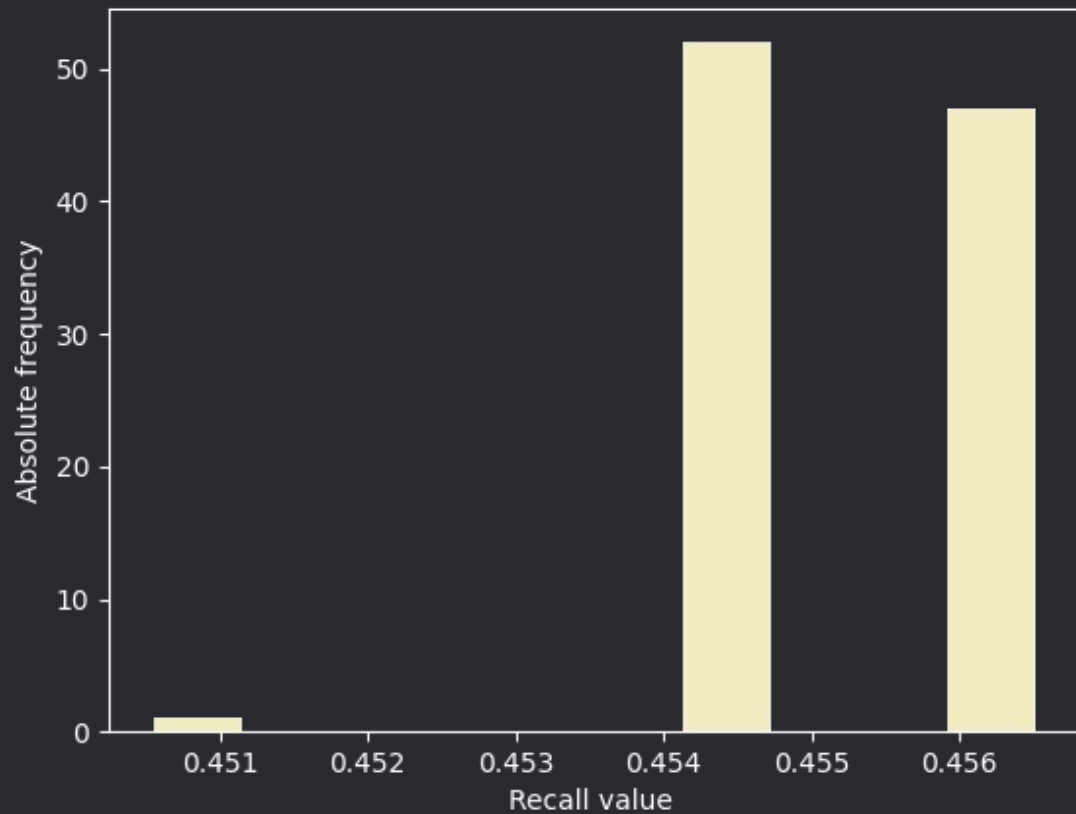# Zm results on generated and real data
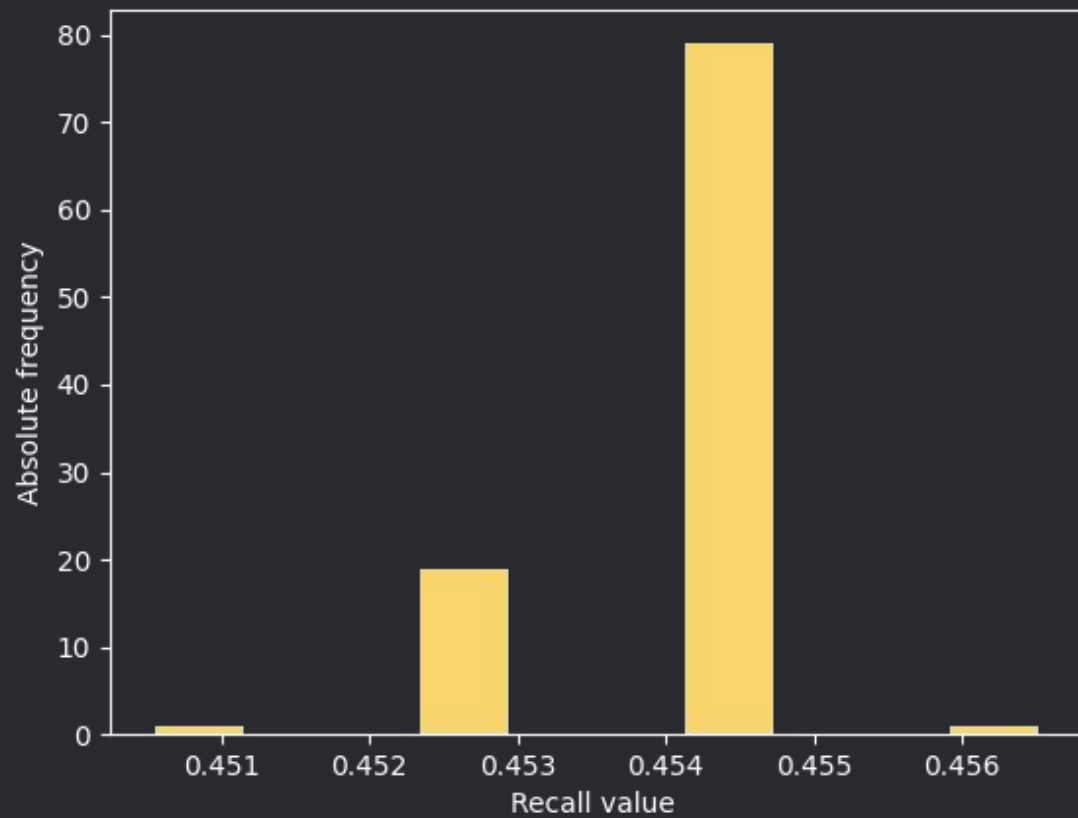
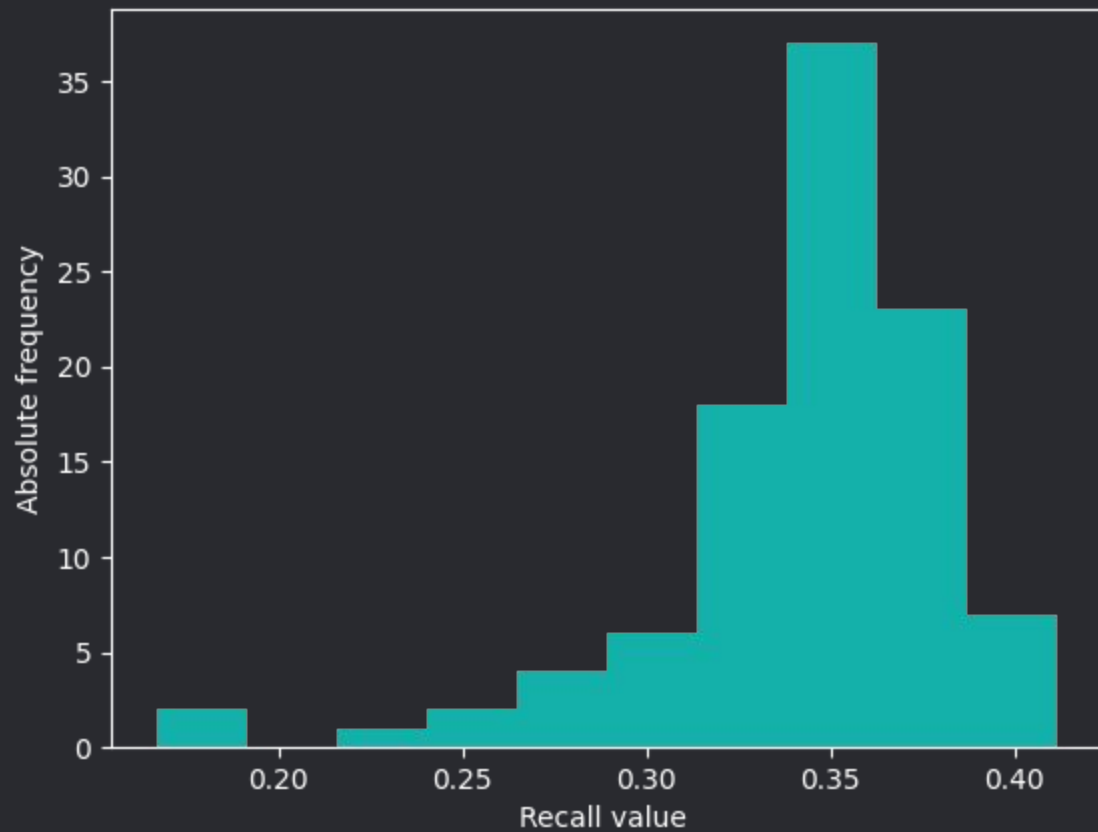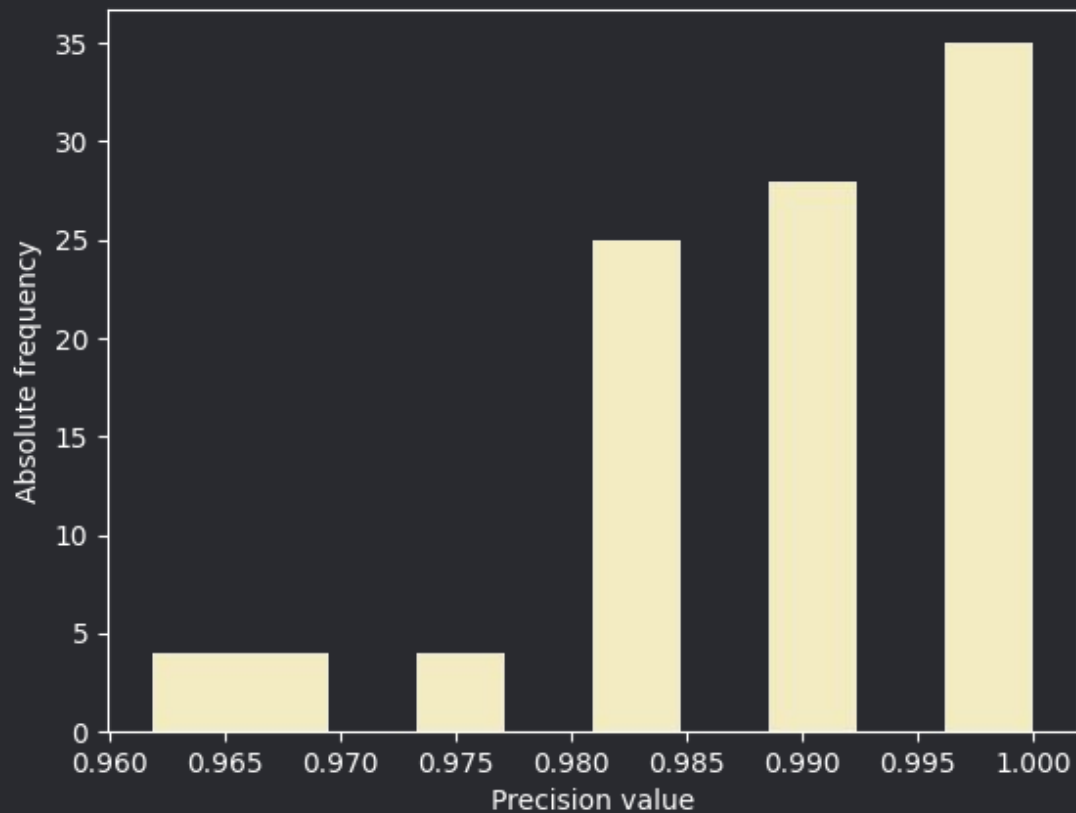# Accuracy - baseline

# Accuracy - GAN

# Accuracy - HMM

# Recall - baseline
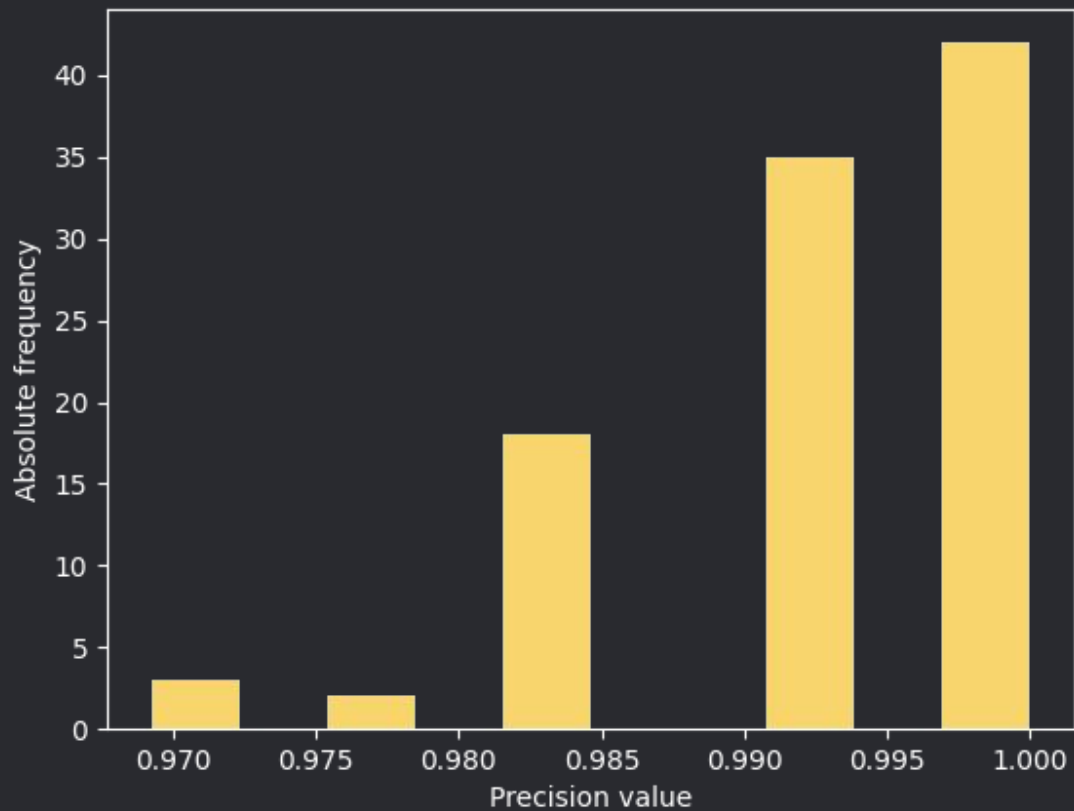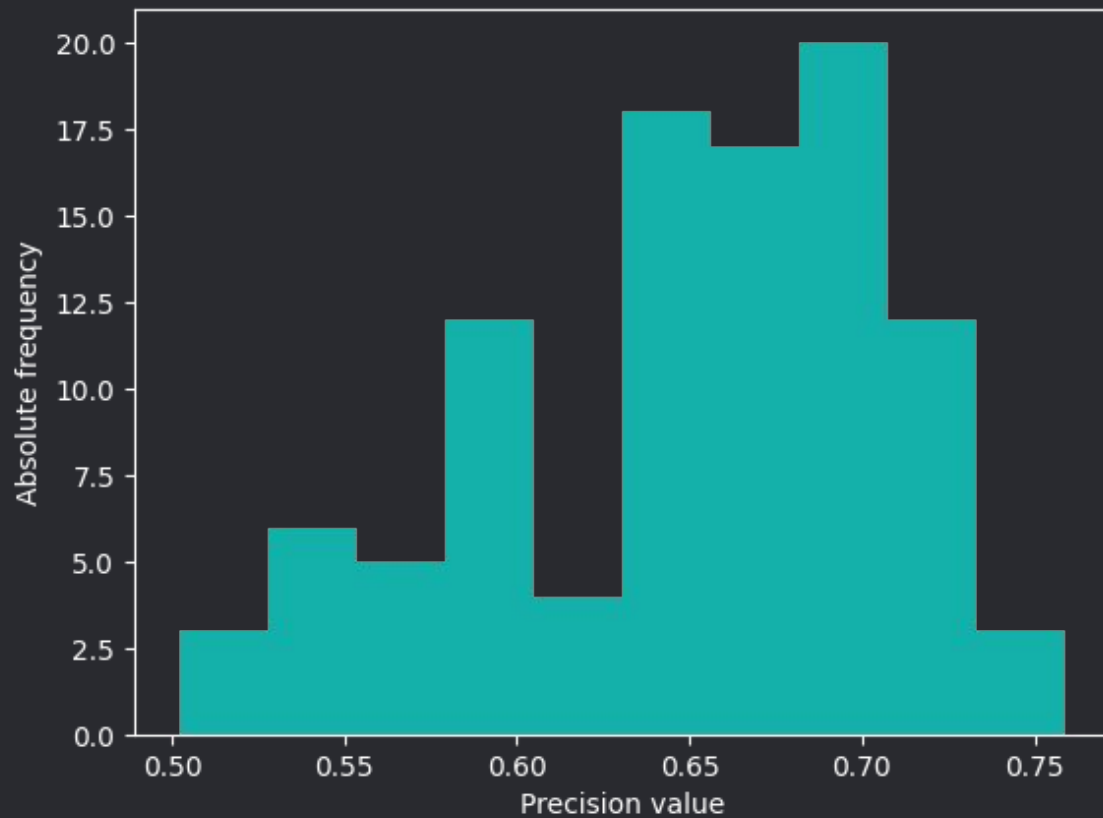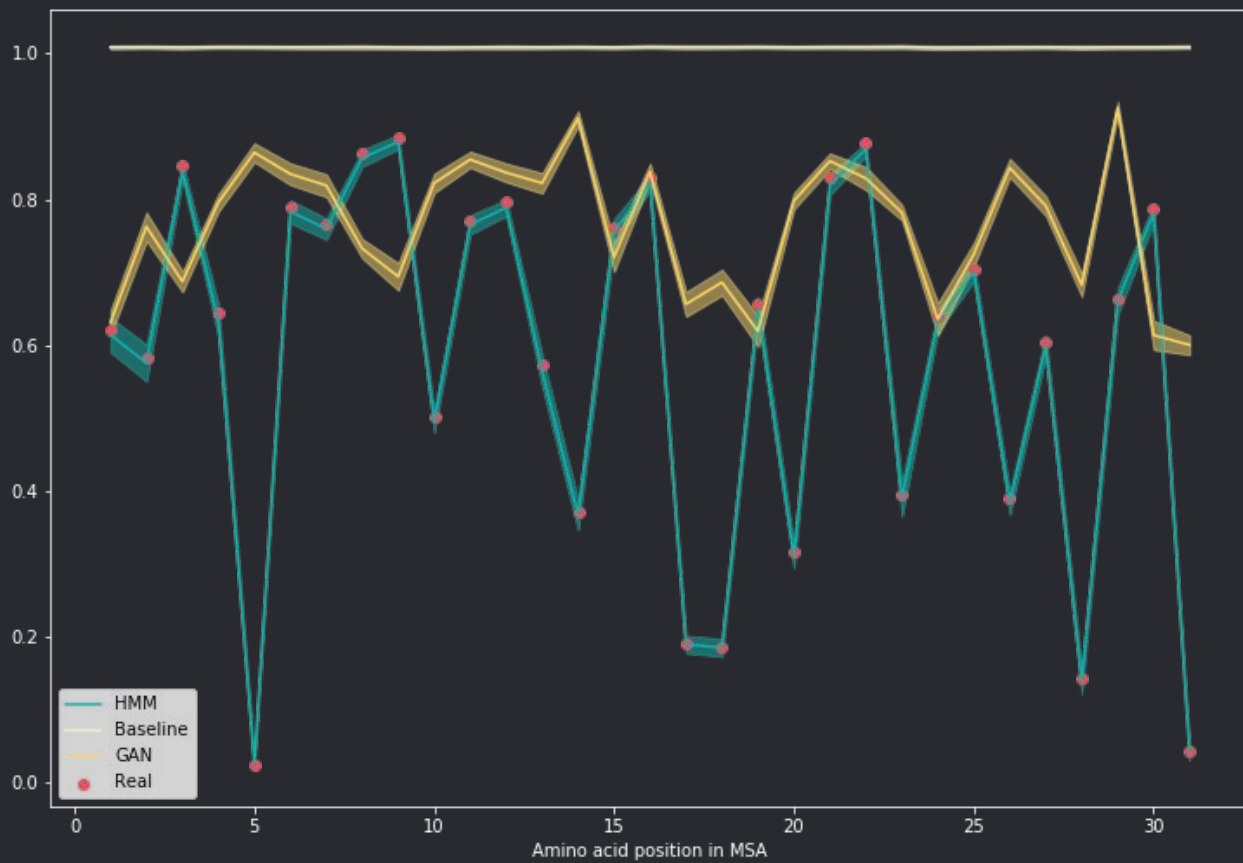
# Recall - GAN

# Recall - HMM

# Precision - baseline

# Precision - GAN

# Precision - HMM

# Shannon Entropy

# So did we do well?