# Gene Diffusion

Philip Kenneweg[1]     Rangaram Dandinasivara Raghuram[1]     Alexander Schönhuth[1]     Barbara Hammer [1]
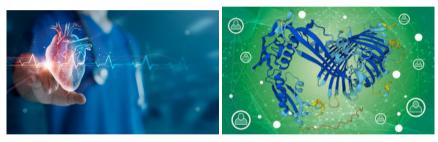
[1]University Bielefeld

24-02-2023, Bielefeld

Introduction

# Motivation

There is great promise in **disease prediction**, **functional understanding**, **drug synthesis** and other applications of artificial intelligence in the field of genetics.



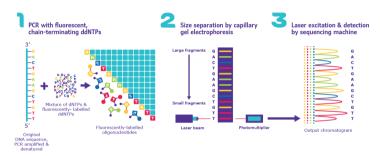(a) disease prediction



(b) protein folding

# Challenges

Very sensitive personal data, which should not be distributed.

# Challenges

**1** PCR with fluorescent, chain-terminating ddNTPs

**2** Size separation by capillary gel electrophoresis

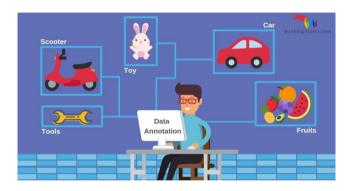**3** Laser excitation & detection by sequencing machine

Genetic data is expensive to obtain, even though sequencing it is getting cheaper each year.

# Challenges

Annotated data is costly to obtain and labels are unreliable, Especially true for medical data.

# Idea

We want to adapt learning paradigms from the recently, successful area of natural language
and image processing. I.e. processing large scales of unlabeled data which enable models
to be trained with few labeled examples.

# Idea

We are taking a closer look at **diffusion models**. These kind of models can help us with:
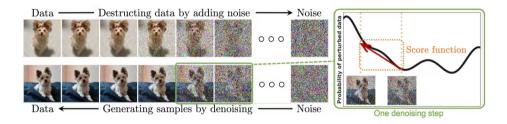
▶ **Semi-supervised learning**: Through the diffusion process we can learn on data without labels.

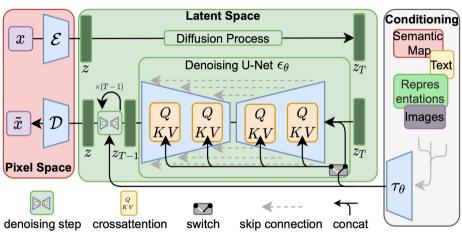▶ **Privacy**: Generalizing from sensitive personal data to synthetic data.
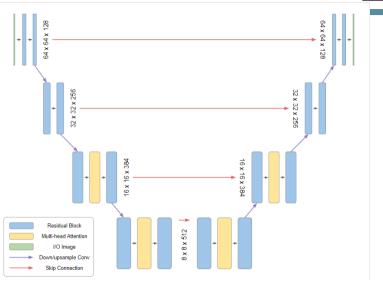
# Background

# Diffusion Models

Data —— Destructing data by adding noise —→ Noise

Data ←— Generating samples by denoising —— Noise

Probability of perturbed data

Score function

One denoising step

# Diffusion Models

# Architecture Unet 1D with Attention

Legend:
- Residual Block
- Multi-head Attention
- I/O Image
- Down/upsample Conv
- Skip Connection

Block dimensions: 64 x 64 x 128, 32 x 32 x 256, 16 x 16 x 384, 8 x 8 x 512

# Genetic Data SNPs
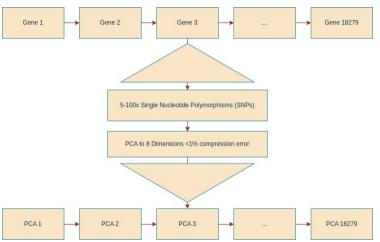
Individual 1
Maternal ...CGATATTCC**T**ATCGAATGTC...
Paternal ...CGATATTCC**C**ATCGAATGTC...

Individual 2
Maternal ...CGATATTCC**C**ATCGAATGTC...
Paternal ...CGATATTCC**C**ATCGAATGTC...

Individual 3
Maternal ...CGATATTCC**T**ATCGAATGTC...
Paternal ...CGATATTCC**T**ATCGAATGTC...

Individual 4
Maternal ...CGATATTCC**C**ATCGAATGTC...
Paternal ...CGATATTCC**T**ATCGAATGTC...

# Genetic Data preprocessing

# Latent Space Analysis

To evaluate the Diffusion Model it is paramount to look at a variety of metrics. Since high dimensional gene pca data is not understandable for humans we use a variety of other metrics.

- ▶ loss curves
- ▶ UMAP
- ▶ classification task performance (ALS disease detection)
- ▶ diversity measures (closest sample etc.)
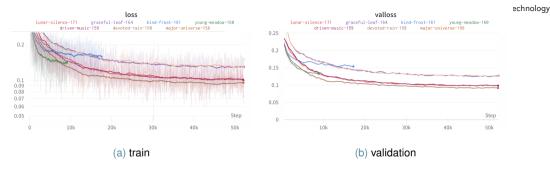
# Loss Curves

(a) train



(b) validation

Figure: Performance differences for different network architectures

In general:

▶ Bigger = Better (no overfitting observed)

▶ some additional preprocessing beneficial (a custom dense layer for each pca)
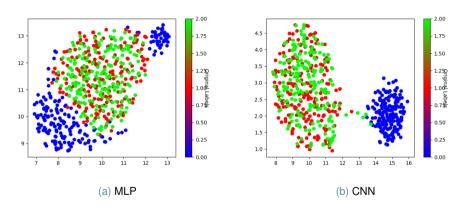
# UMAPs

(a) MLP

(b) CNN

Figure: UMAPs for different network backbones using cosine similarity as a distance metric.
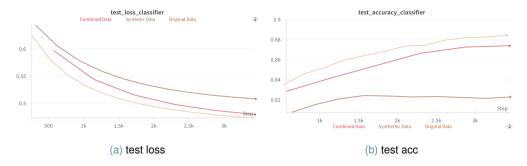Euclidean distance UMAPs are not informative.

# ALS classifier performance

(a) test loss



(b) test acc

Figure: Comparison between synthetic and real training data. Performance is measured on a hold out test set.

UNIVERSITÄT
BIELEFELD
Faculty of Technology
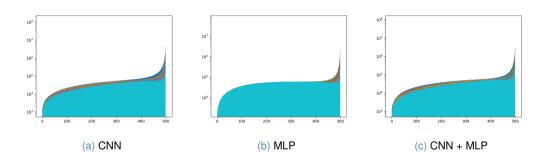


(a) CNN

(b) MLP

(c) CNN + MLP

Figure: Loss curves during training for different Model Architectures.
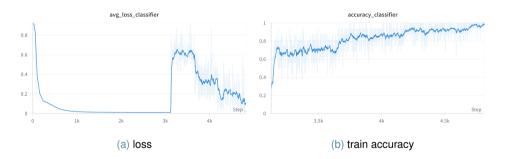
(a) loss

(b) train accuracy

Figure: Loss and Accuracy for Transformer during pre-training and fine-tuning.

# Discussion

# Improvements?

- ▶ More data
- ▶ Different validation task (ALS seems problematic)
- ▶ ... your ideas?

Thank you!