

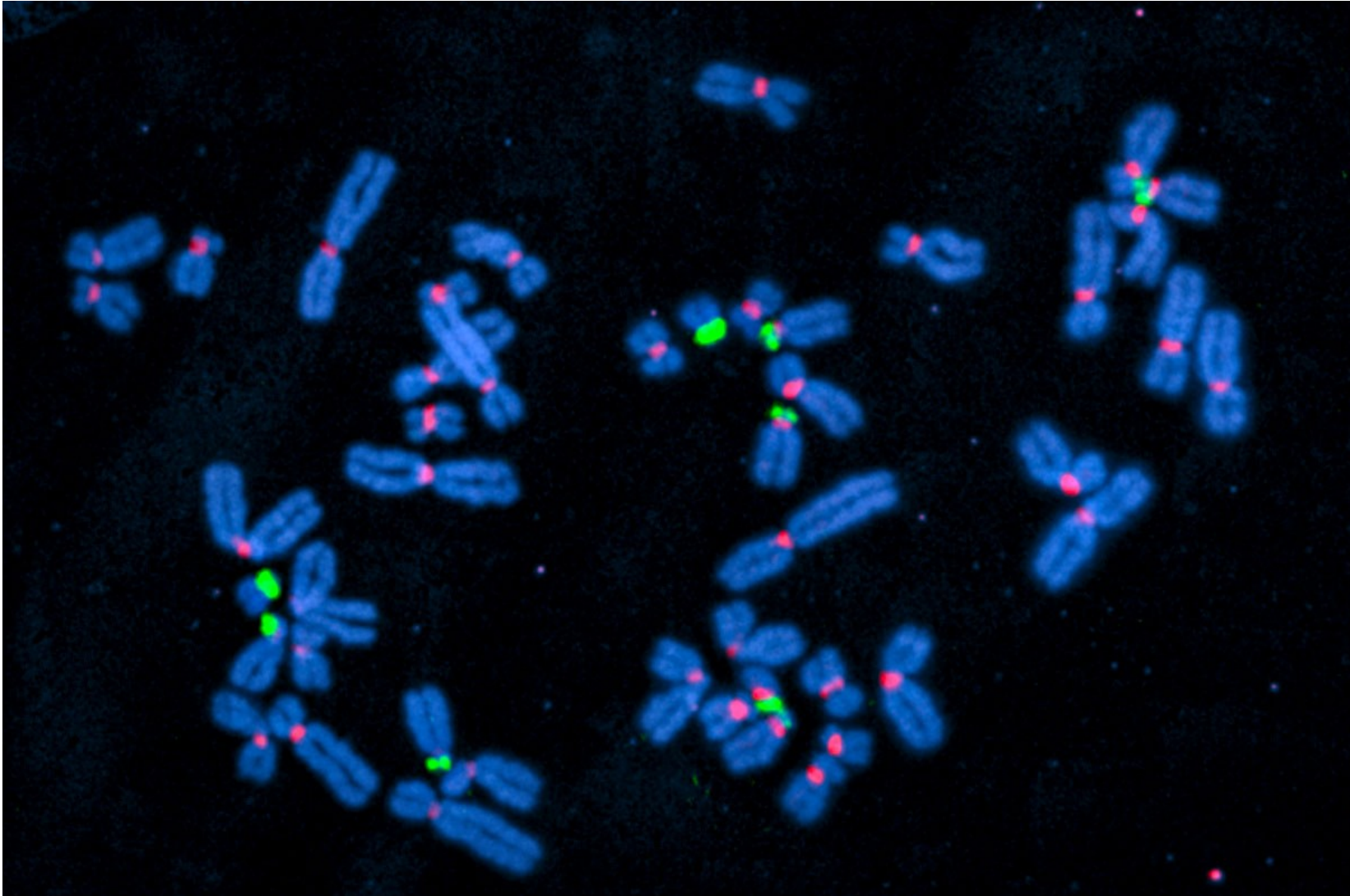
# Quantitative trait prediction using GCAE

Richel Bilderbeek  
2022-03-05

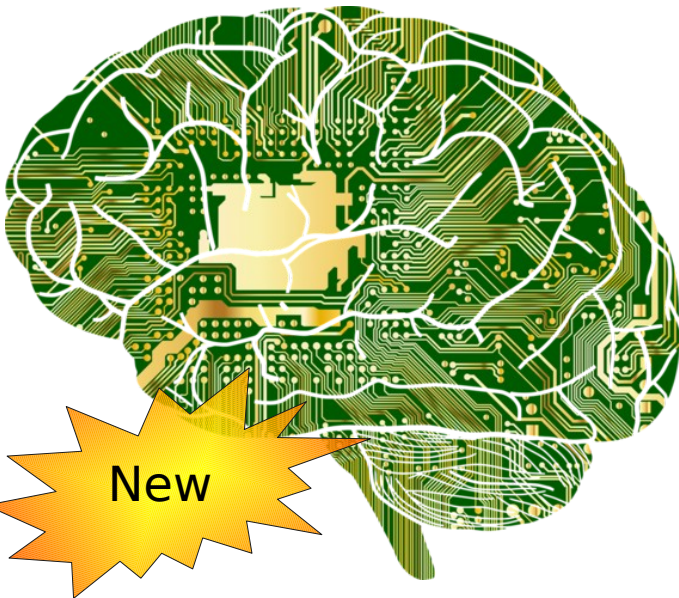
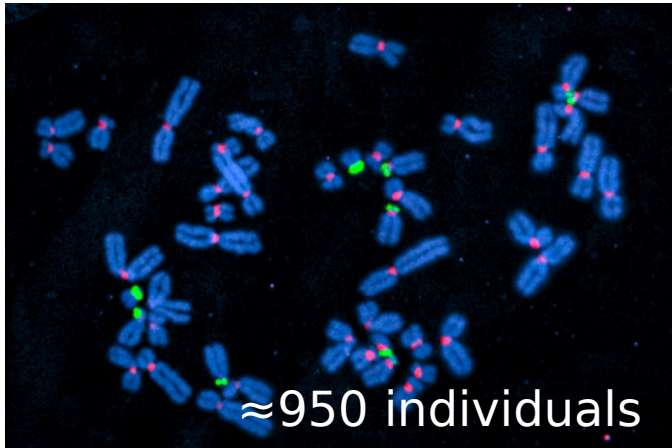


[https://github.com/richelbilderbeek/science\\_presentation\\_20220305](https://github.com/richelbilderbeek/science_presentation_20220305)

# Introduction

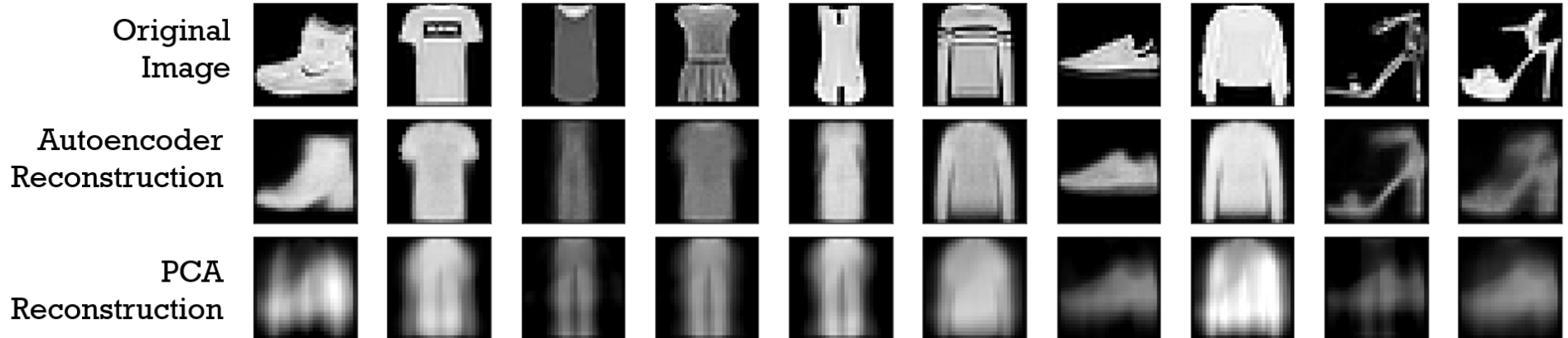
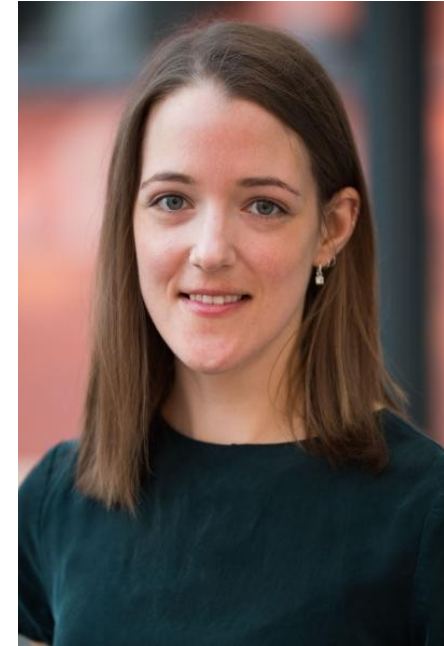


# Goal



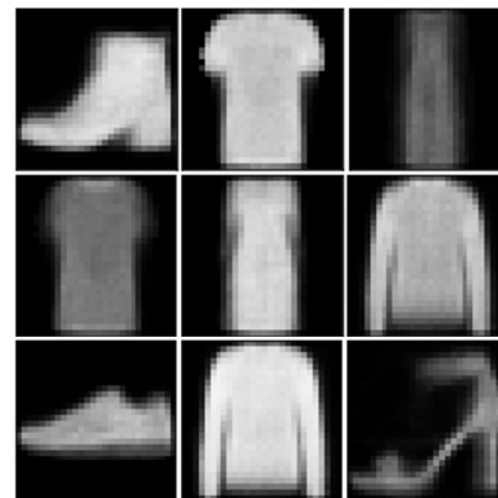
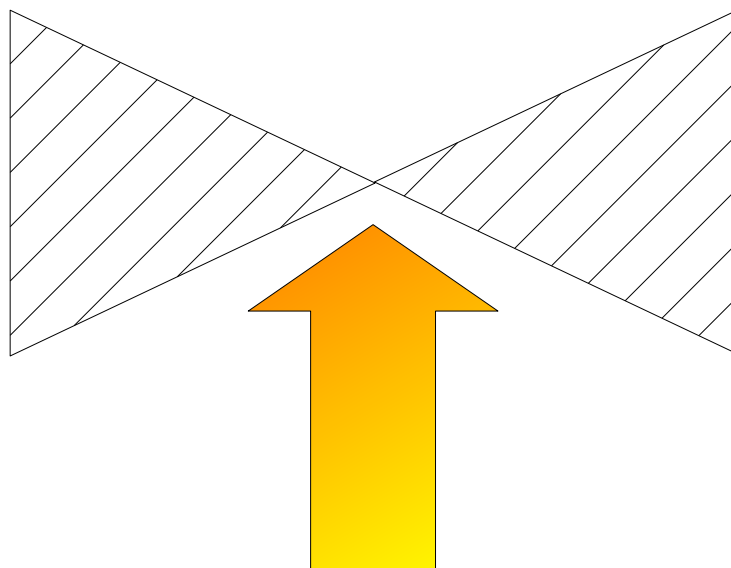
# Method: GCAE

## Genomic Convolutional Auto Encoder Dimensionality reduction



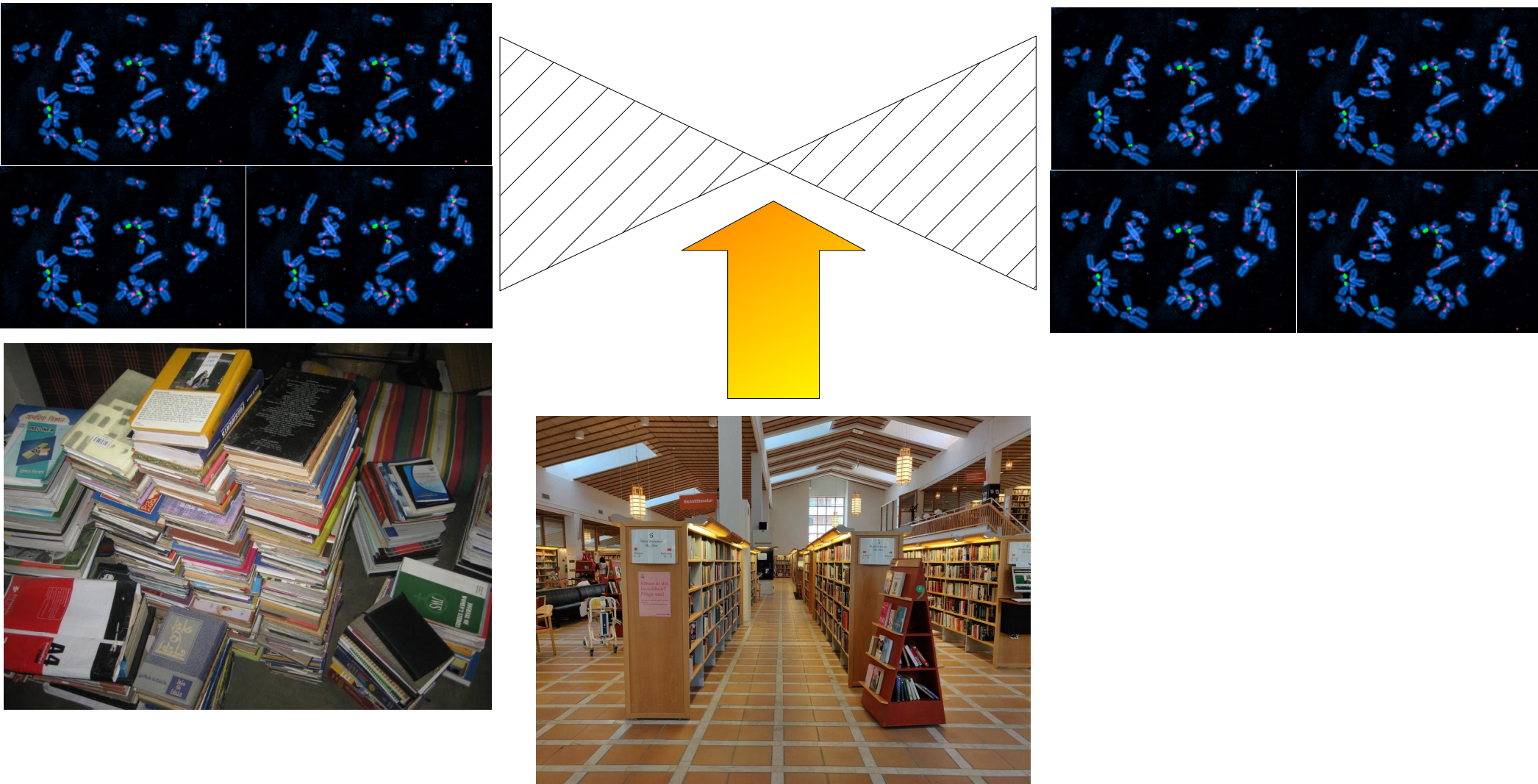


# What an auto-encoder does

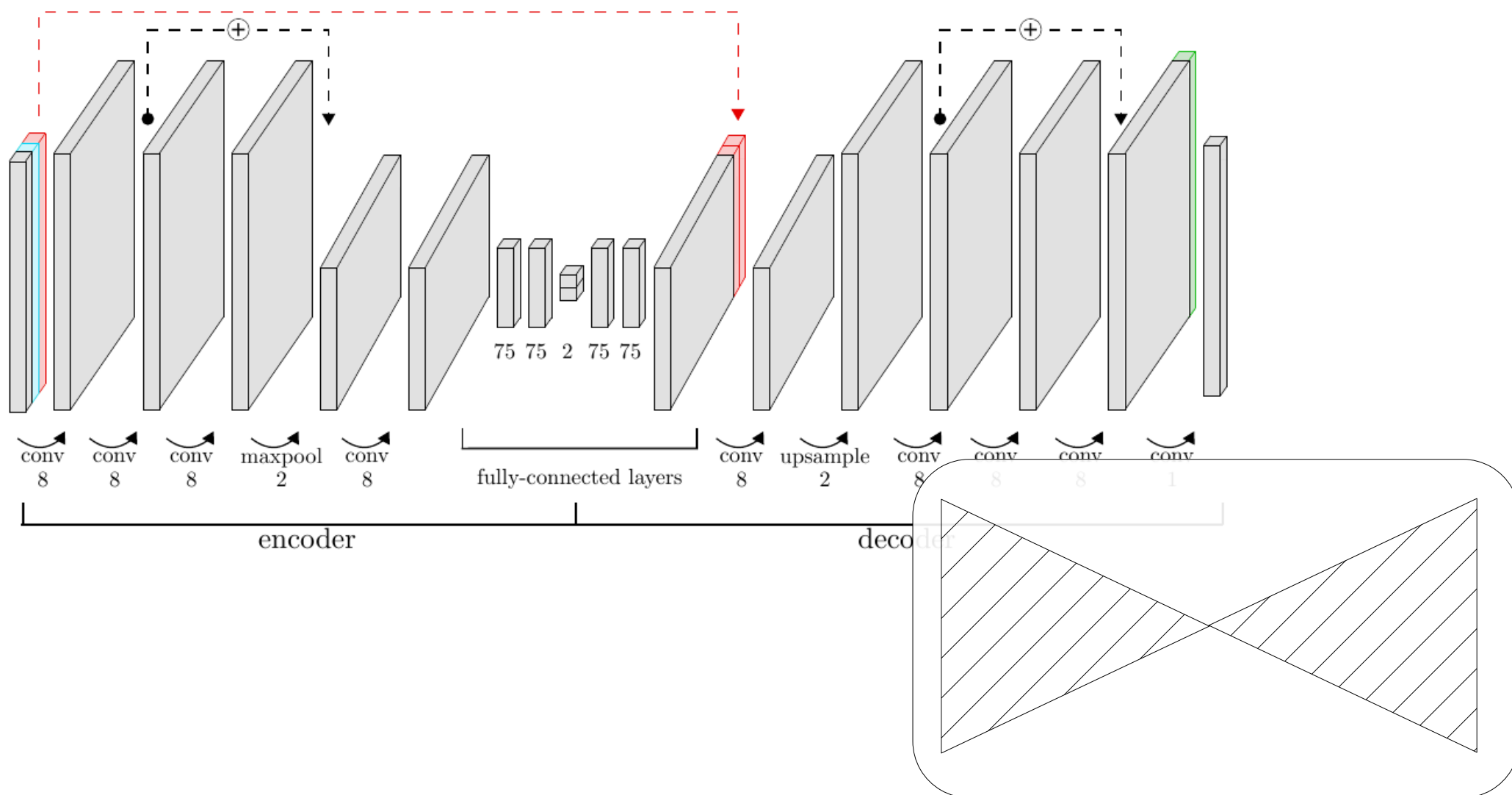


Latent layer with  
dimension  
reduced data

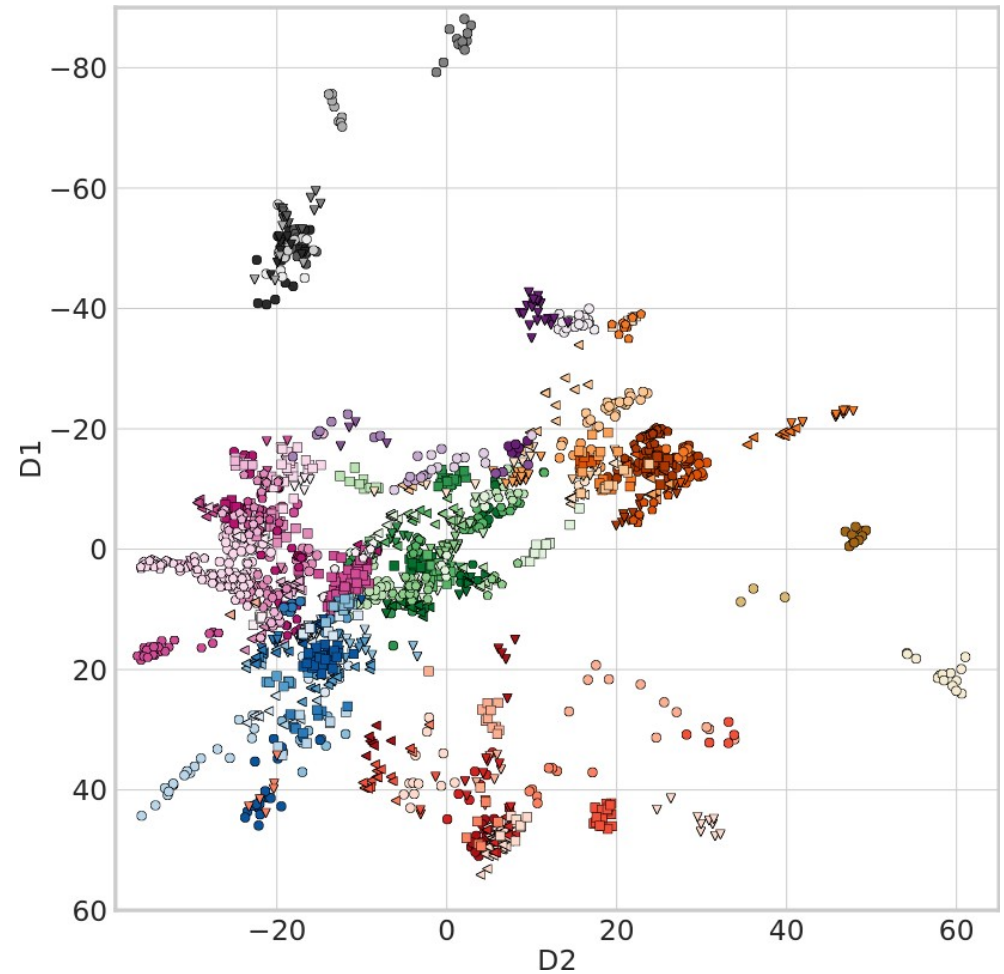
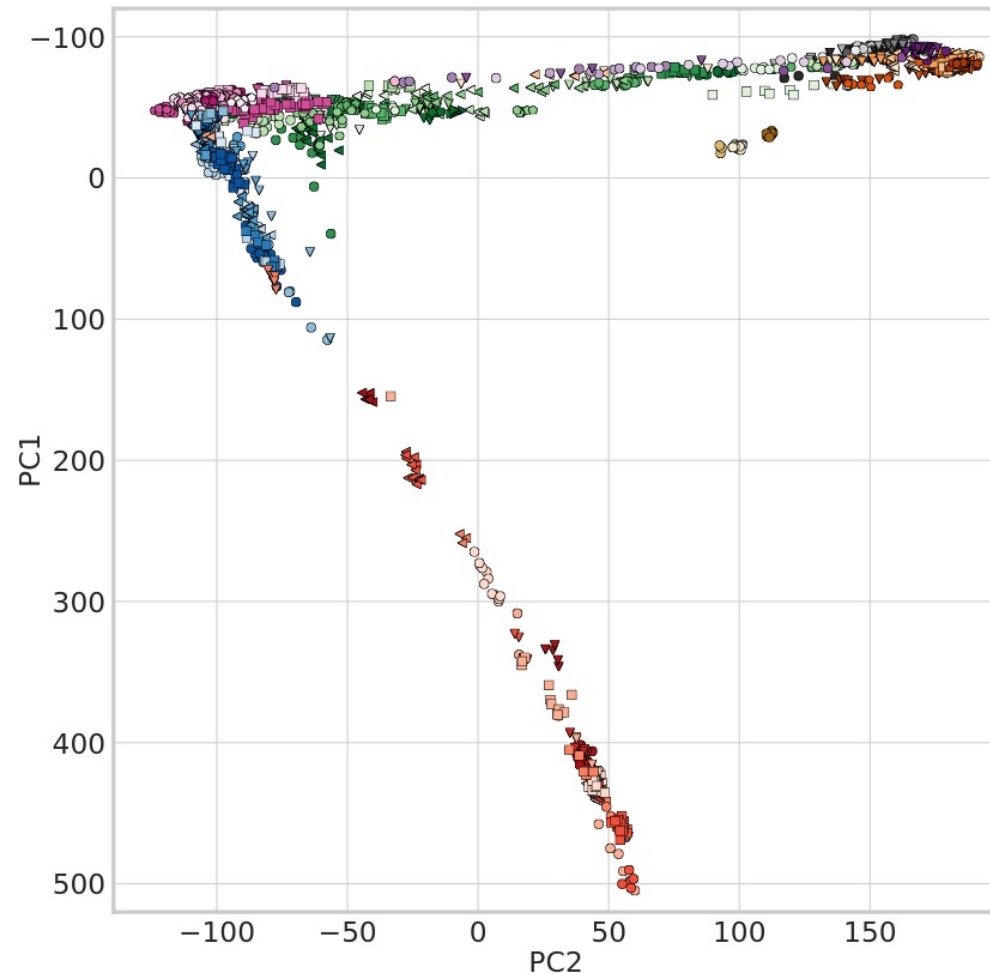
# GCAE is an auto-encoder



# GCAE is highly tunable

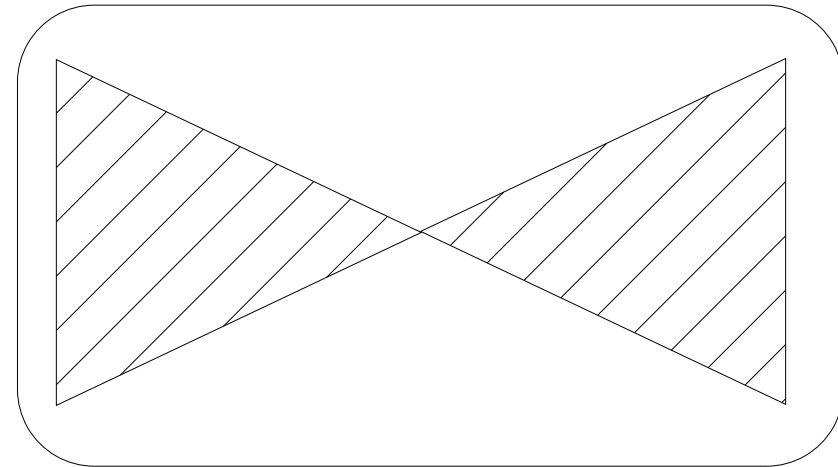
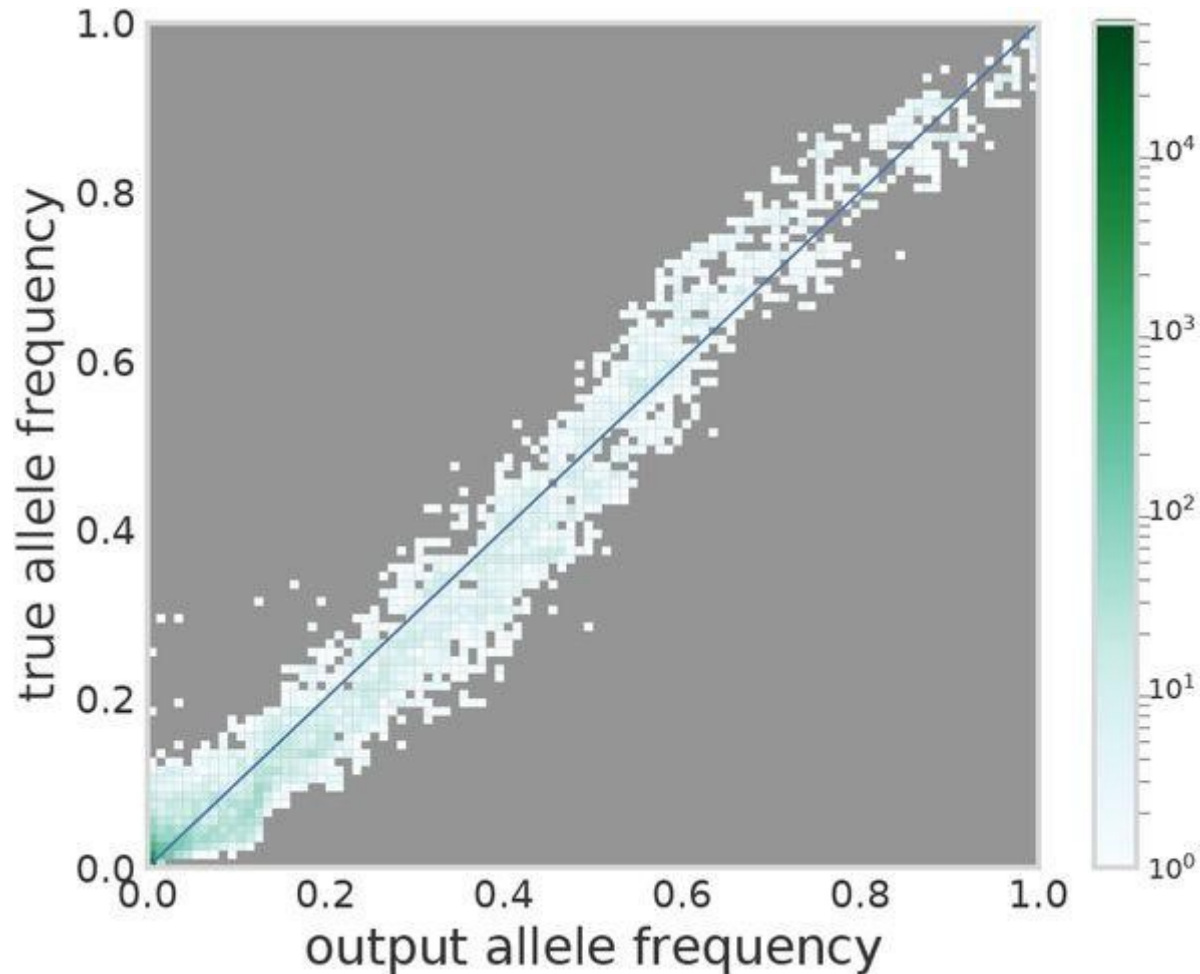


# GCAE can do non-linear dimensionality reduction



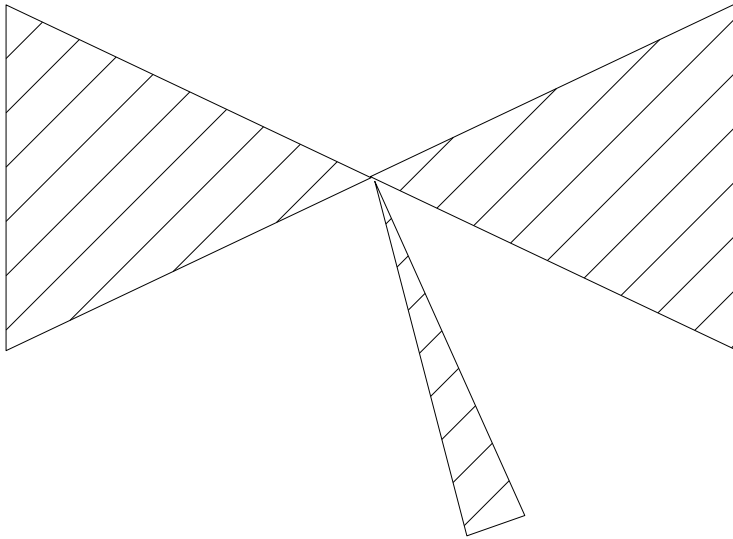


# GCAE cannot do rare alleles

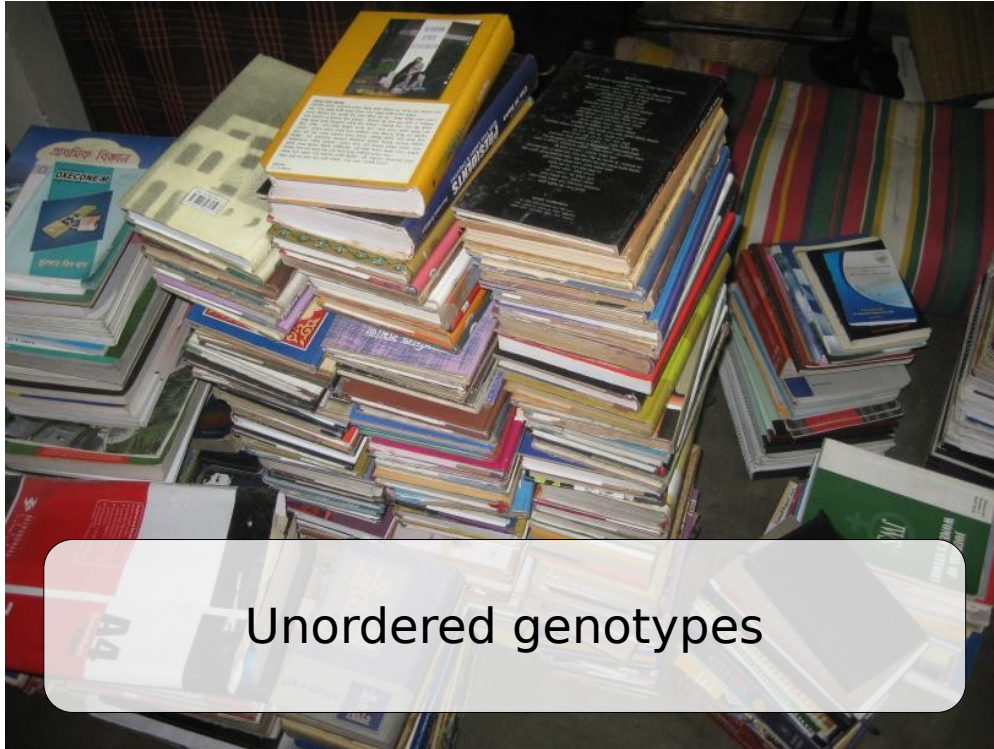


# GCAE may do trait prediction

## Untested quantitative trait prediction extension



# Analogy



Unordered genotypes



Ordered genotypes

Predict the number of horses, spaceships, etc.

Protein concentrations

# What has been done

## GCAE

Can actually run with trait prediction

Can run on Bianca

## **gcaer** works with **GCAE**

Tested to be correct

Proper error messages

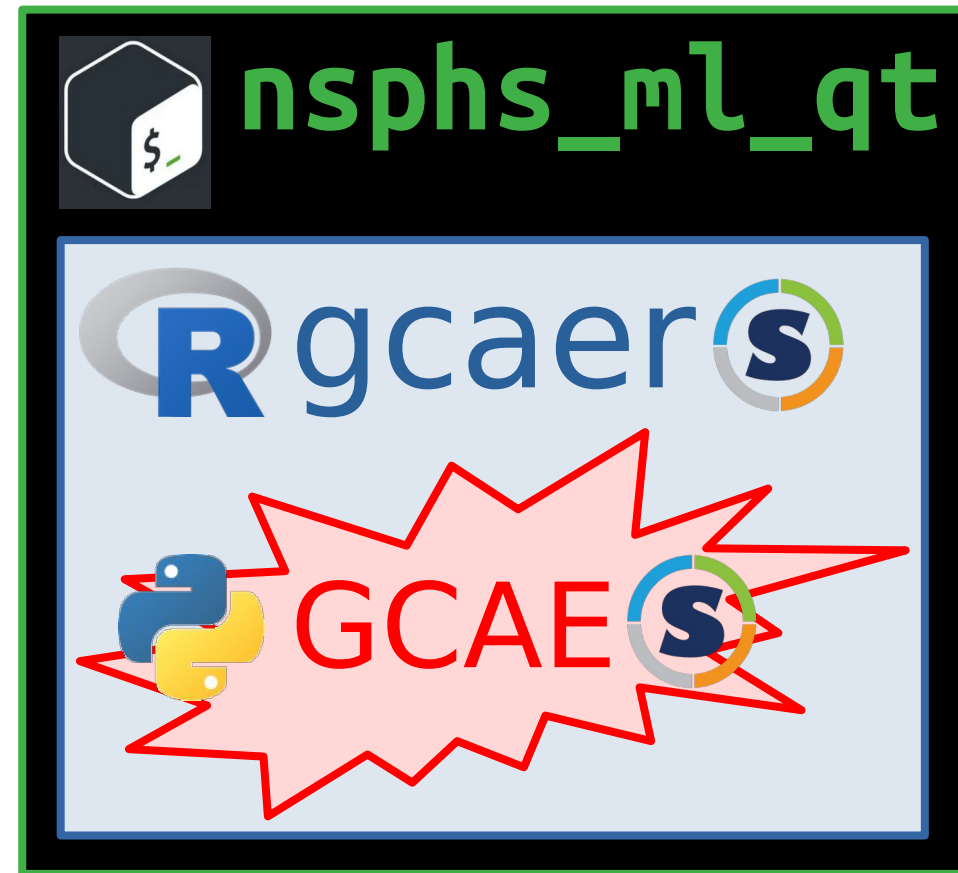
Extend functionality

## **nsphs\_ml\_qt** does experiments

Scripts re-used in many contexts

Works locally, on Rackham, on Bianca

You can run these today





# First experiment

**Simulate simplest dataset**

**Run GCAE as-is**

**Evaluate performance**

# First experiment: simulated dataset

**1 monogenic trait with  $H^2 = h^2 = 1$**

**1k individuals**

**MAF = 0.499**

For PLINK users:  
why not just use  
0, 1 and 2  
for trait values?

chr	id	posg	pos	ref	alt
:---	:-----	-----:	---:	:---	:---
1	snp_1	0	1	A	C

FID	IID	additive
:---	:---	-----:
A	1	9.424778
A	2	6.283185
A	3	3.141593
A	4	9.424778
A	5	9.424778

	snp_1
:--	-----:
1	0
2	1
3	2
4	0
5	0

fam	id	pat	mat	sex	pheno
:---	:--	:---	:---	-----:	-----:
A	1	0	0	1	-9
A	2	0	0	1	-9
A	3	0	0	1	-9
A	4	0	0	1	-9
A	5	0	0	1	-9

# First experiment: PLINK results

## PLINK measures association perfectly

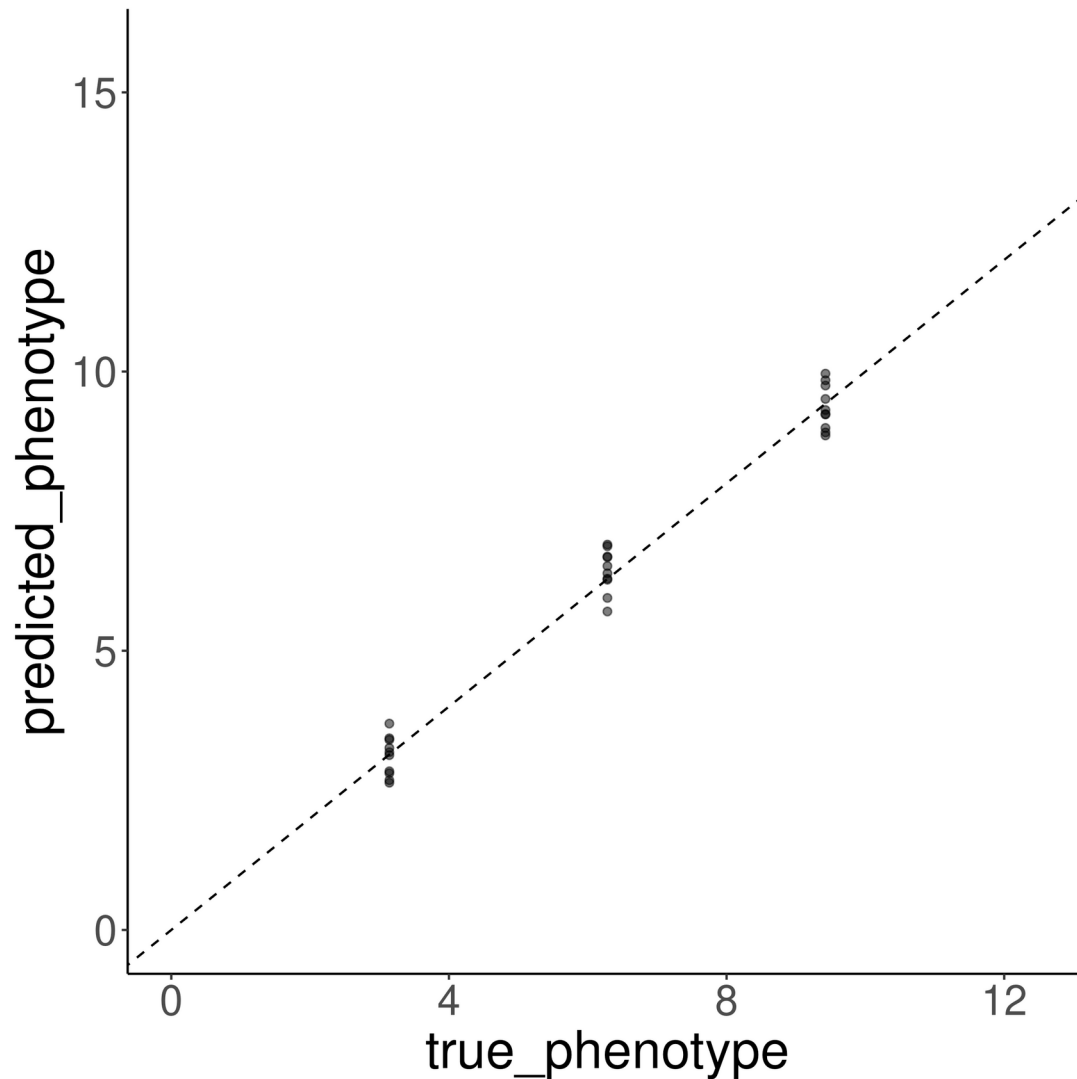
CHR	SNP	BP	NMISS	BETA	SE	R2	T	P
---	:---	--:	-----:	-----:	--:	--:	-----:	--:
1	snp_1	1	1000	-3.142	0	1	-116700000	0

FID	IID	additive
:---	:---	-----:
A	1	9.424778
A	2	6.283185
A	3	3.141593
A	4	9.424778
A	5	9.424778

	snp_1
:--	-----:
1	0
2	1
3	2
4	0
5	0

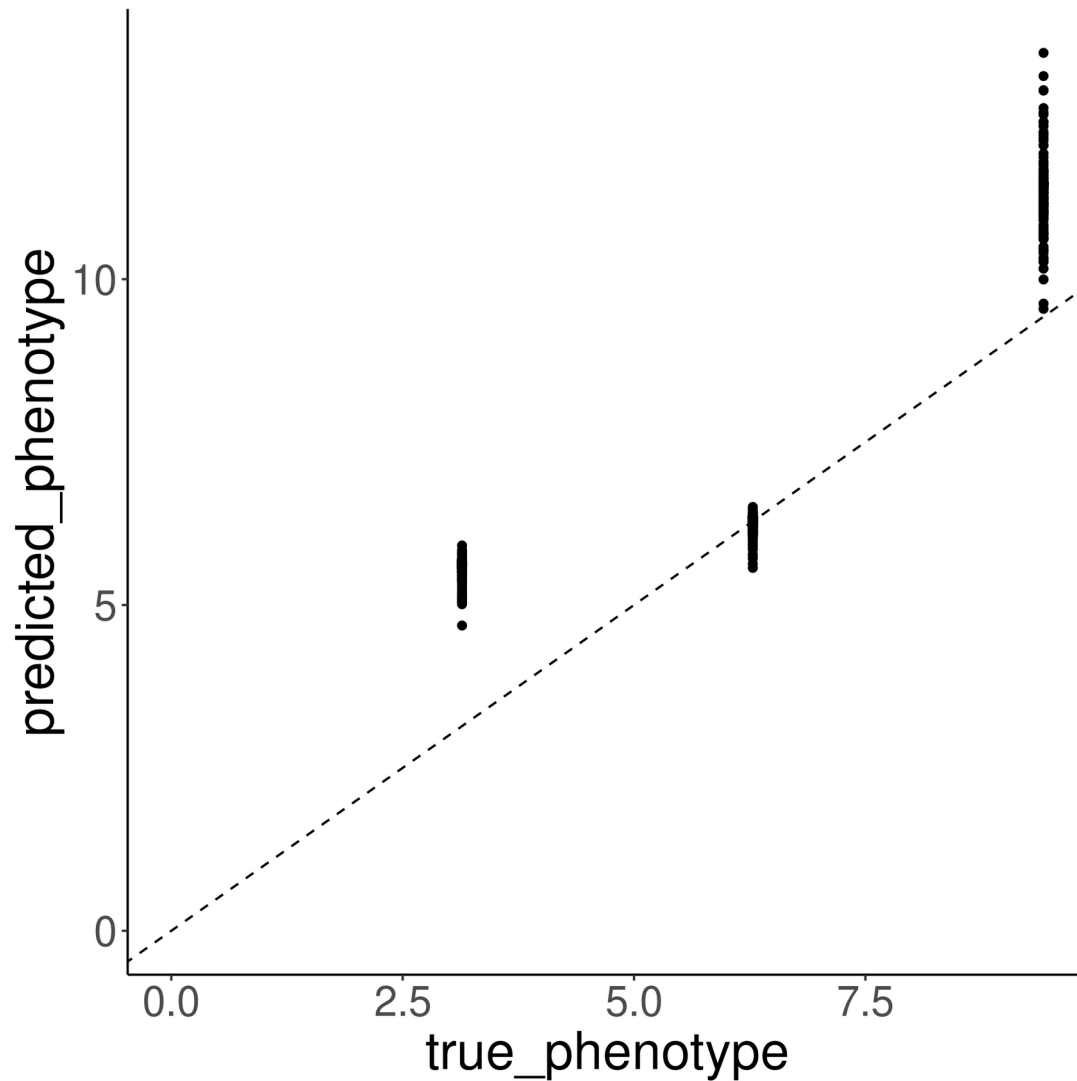
For PLINK users:  
why not just use  
0, 1 and 2  
for trait values?

# First experiment: expected GCAE predictions





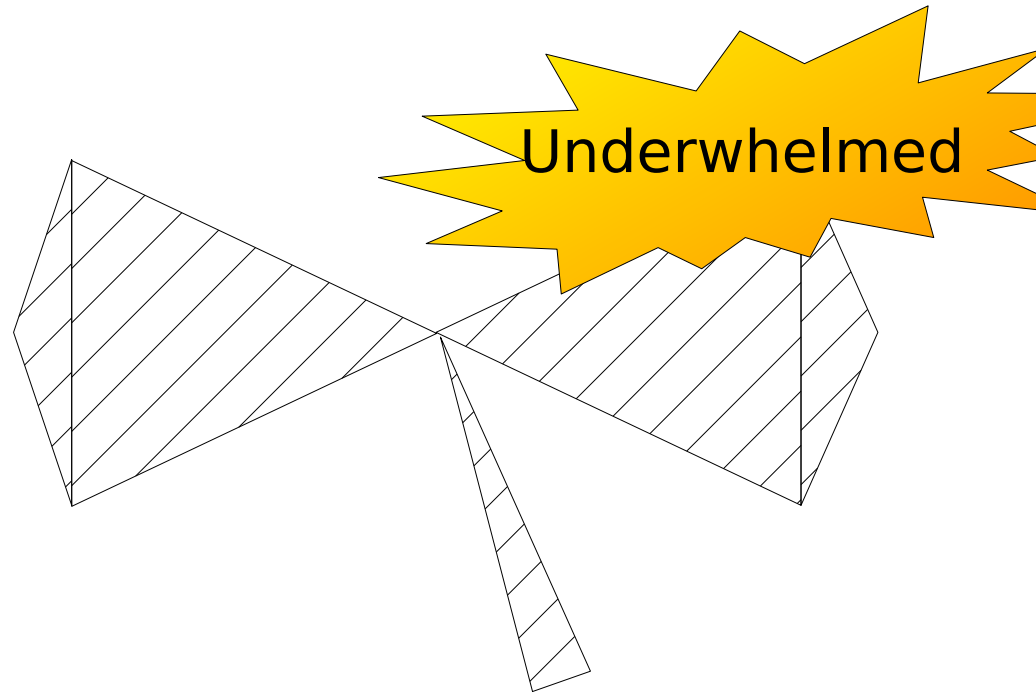
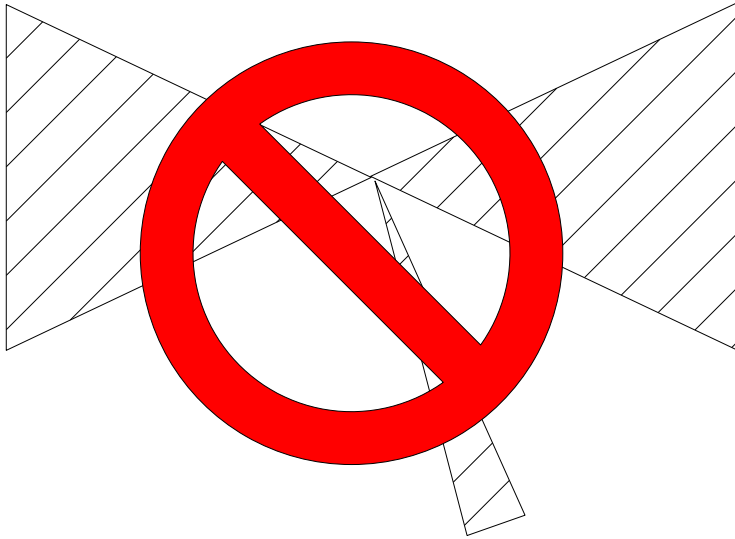
# First experiment: GCAE predictions



# First experiment: conclusions

**Proof of concept works, prediction underperforms**

**Hypothesis:**



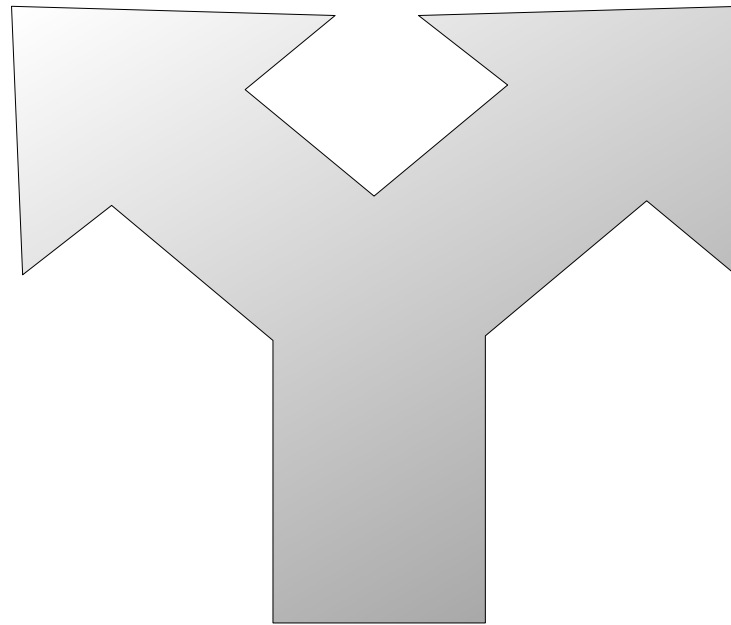
# Decision

## **Tune the network to the problem**

try a simpler neural network  
GCAE focused

## **Tune the problem to the network**

try harder problems  
application focused



# First NSPHS experiment

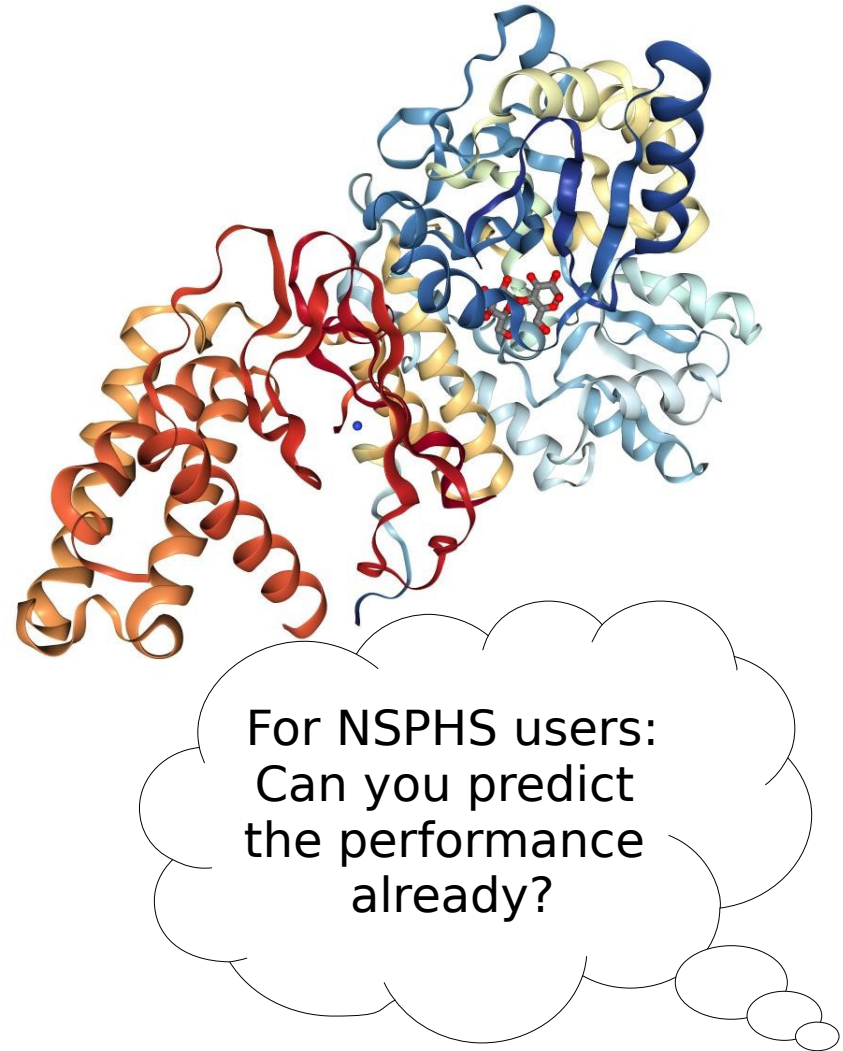
## Use NSPHS and first protein concentration

Pick 100k random variants

Adrenomedullin

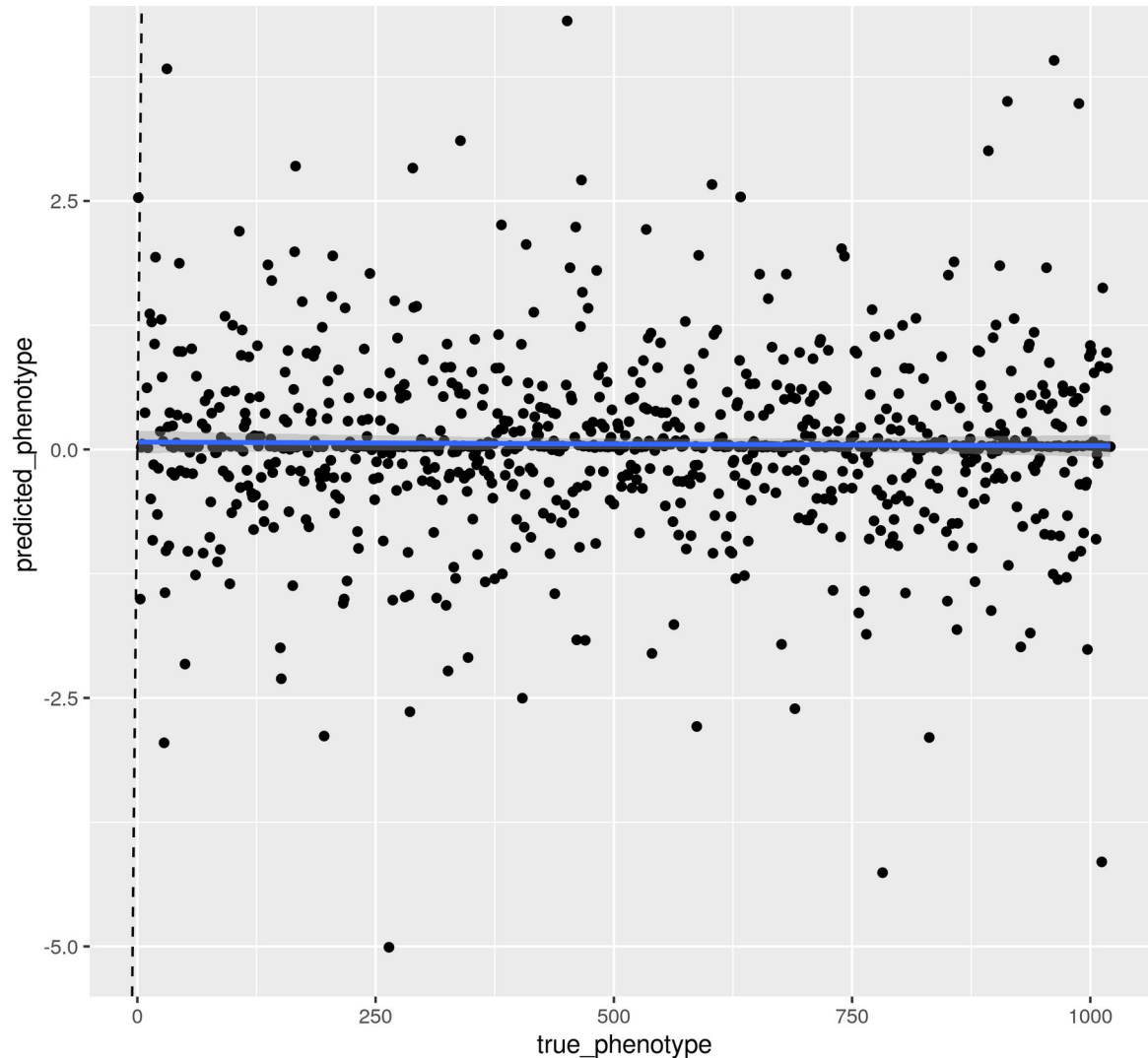
## Run GCAE as-is

## Evaluate performance





# First NSPHS experiment: results



**Runs!**

**60 hours**

**1k epochs**

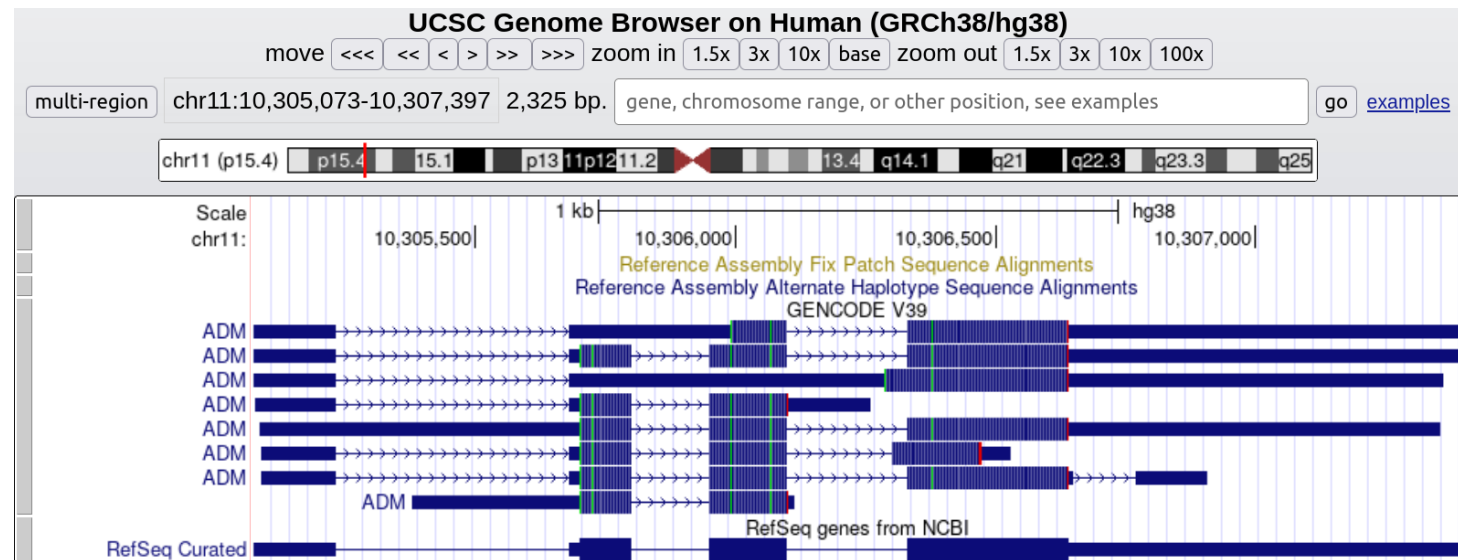
**Unsure about  
learning trajectory**

# First NSPHS: conclusion

## The result is just noise

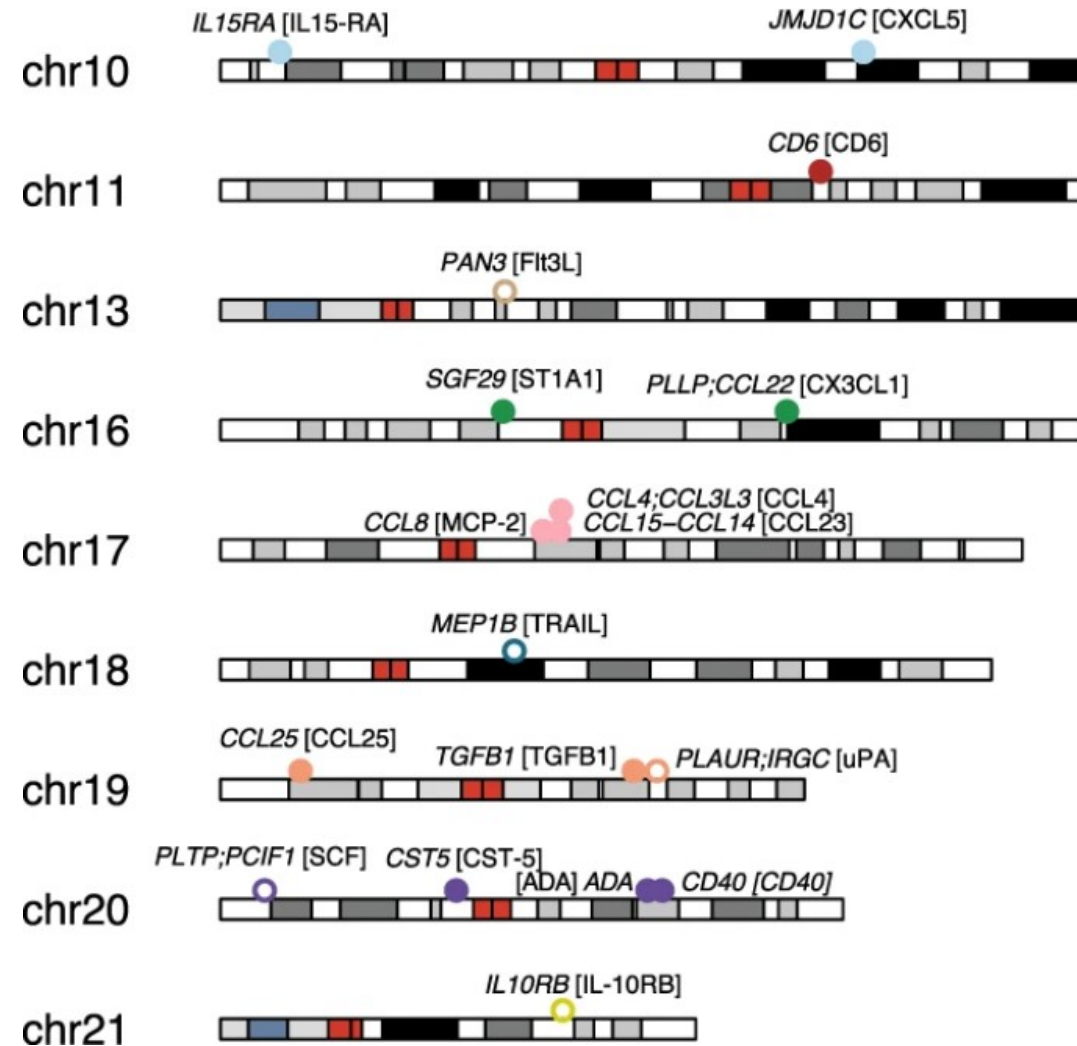
The gene for adrenomedullin is monogenic, 52 amino acids

## But hey, it works!



# First NSPHS: conclusion

**\* [ ] Use the SNPs known to have an association, e.g. 0.5 Mb downstream and upstream around known cis-regulatory elements, #5**



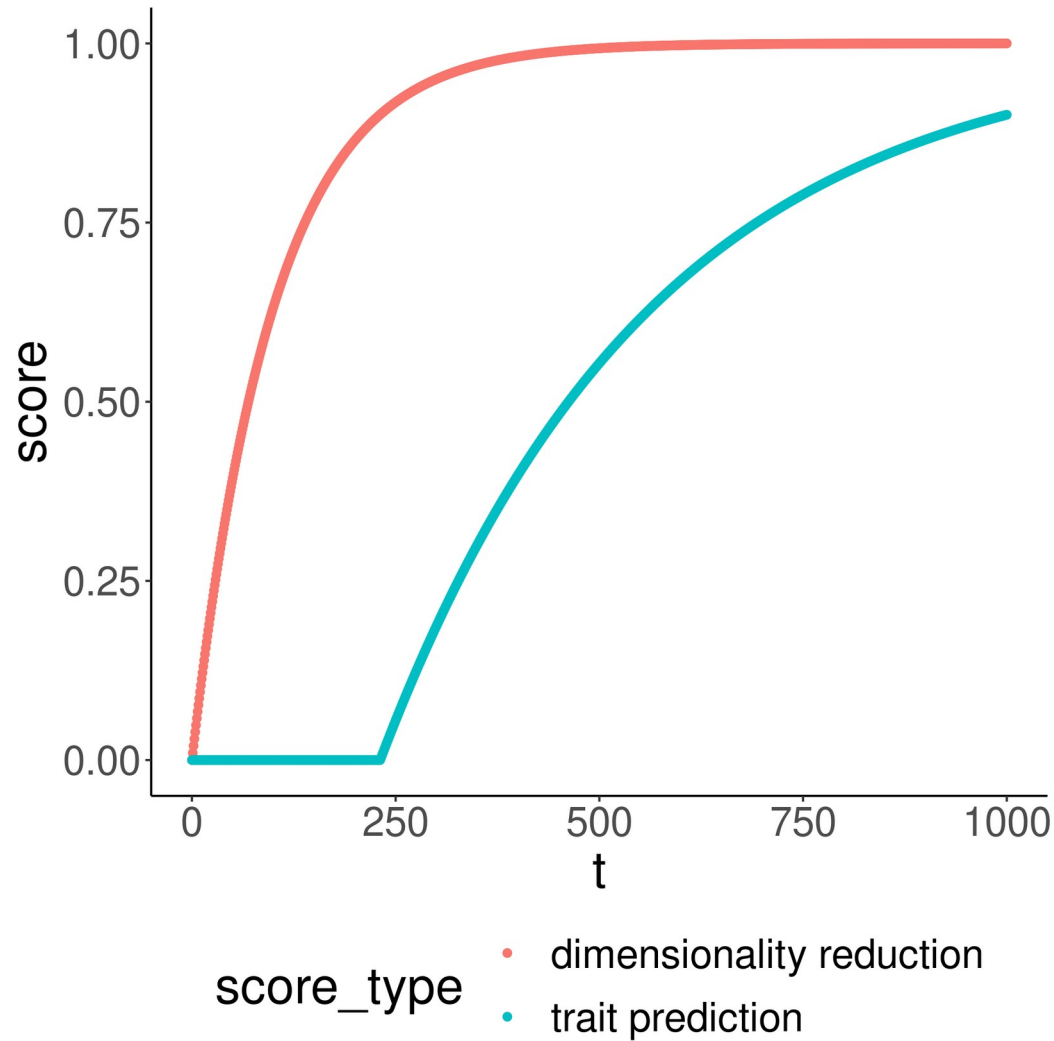
# Learning

Learning proceeds  
until death  
and only then  
does it stop

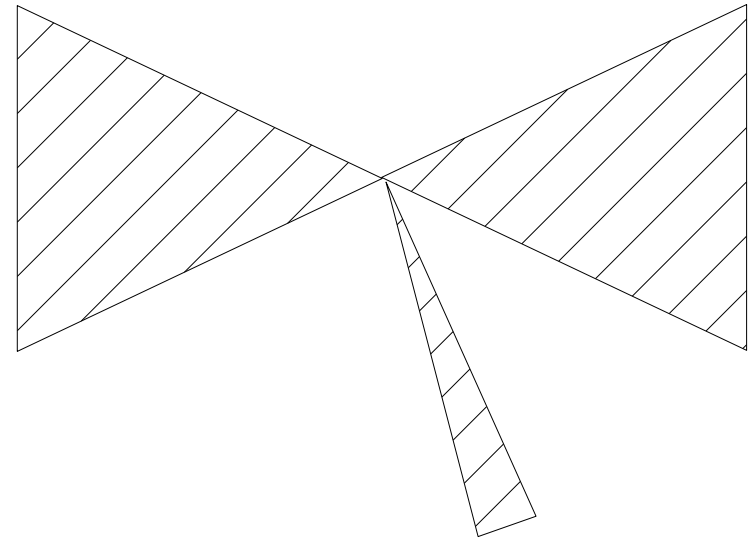




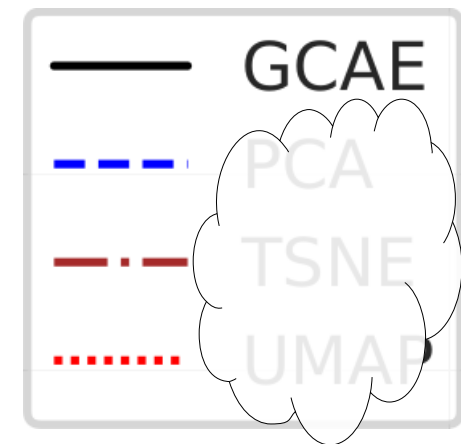
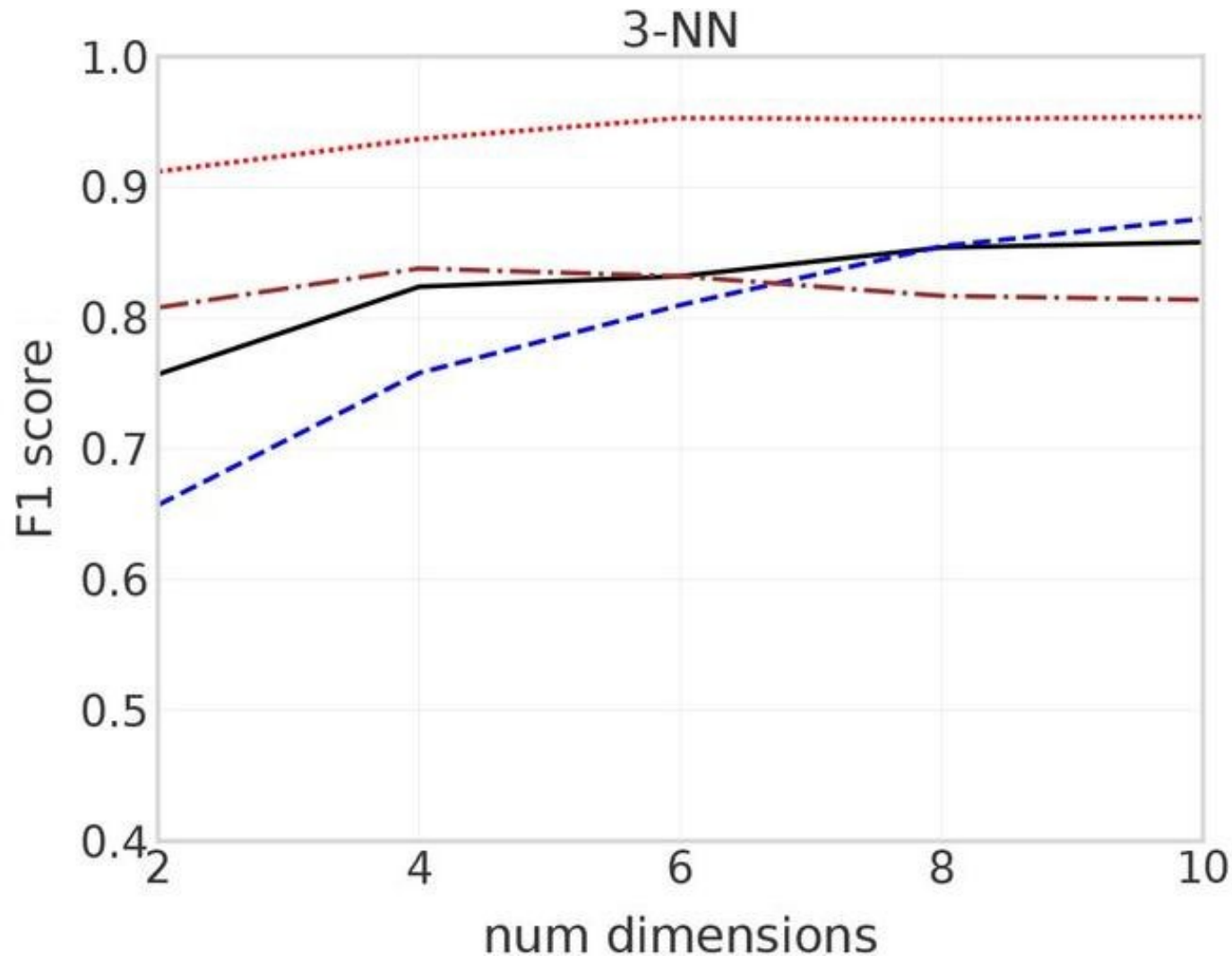
# Learning trajectory



**Trait prediction can only be trained well after dimensionality reduction is good**



# Dimensionality reduction score



# Caveats in application

## **We cannot assign clusters:**

NSPHS is  $\approx 1$  cluster, i.e. Karesuando and Soppero

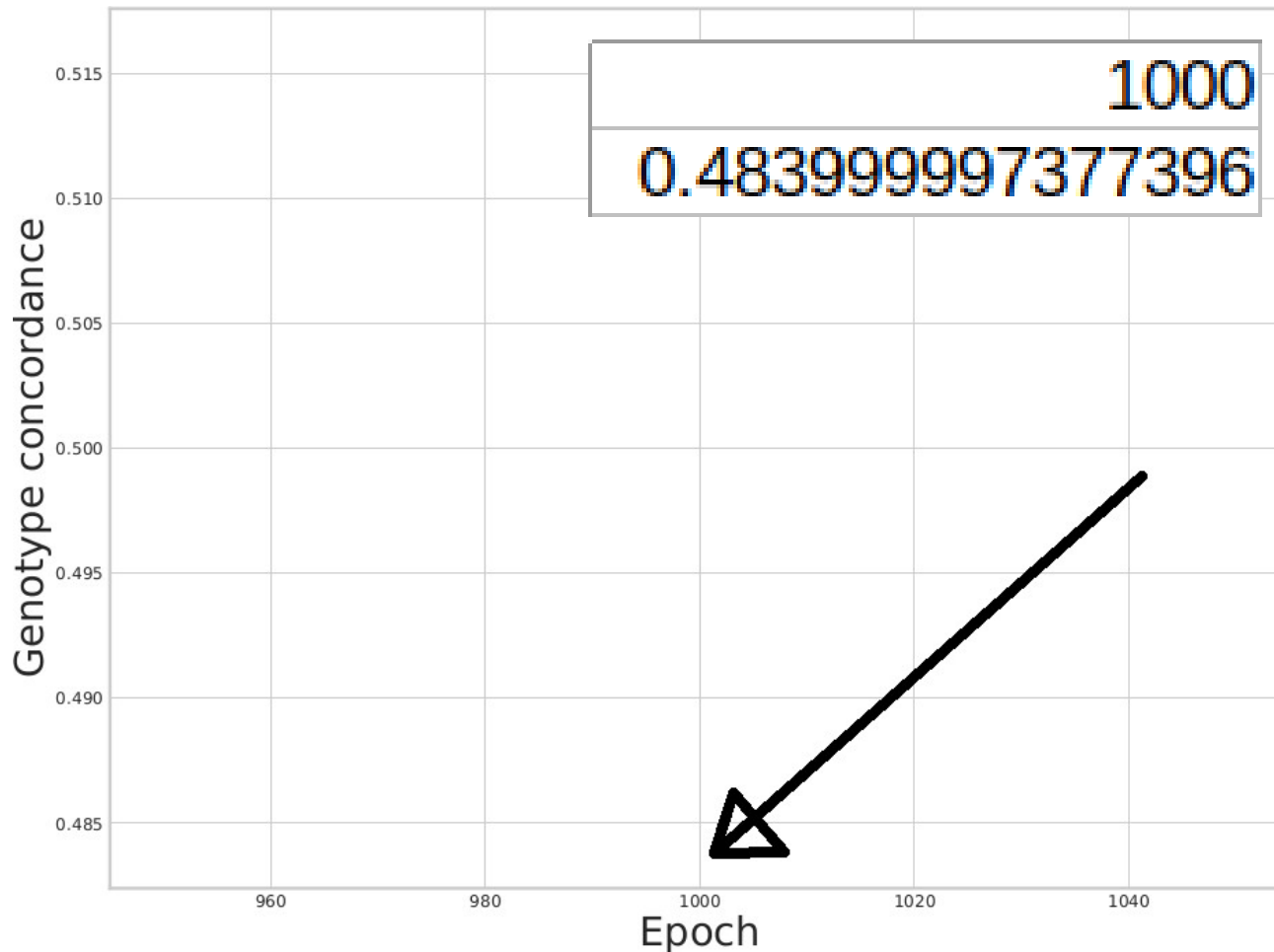
## **All we care about is that dimensionality reduction has reached equilibrium**

Using 2 clusters anyway may allow us to do so

## **Maybe we only care about the trait predictions**



# Dimension reduction quality: other way



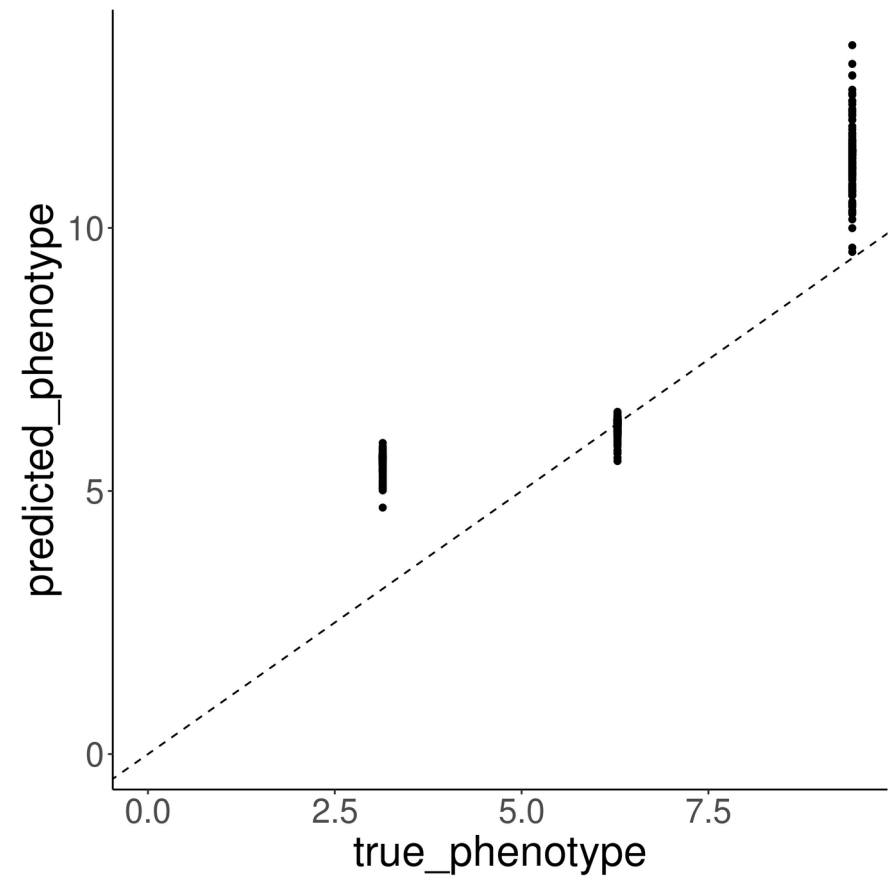
**[ ] Plot the genotype concordance through time, #17**

# Trait prediction score

## How well are phenotypes predicted?

... in time?

... comparable between protein concentrations?

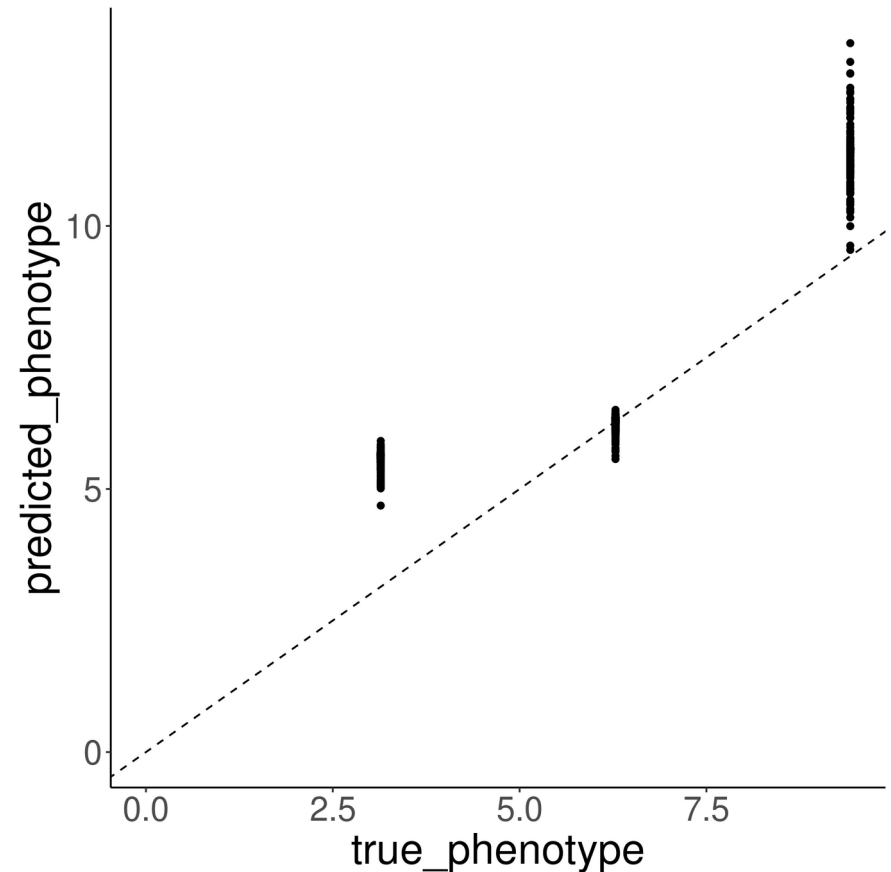


# Trait prediction score 1/2: MSE

**MSE: Mean squared error**

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$$

**Mean (vertical distance to identity line)<sup>2</sup>**



 **Scale matters**

# Trait prediction score 2/2: NMSE

## NMSE: Normalized Mean Squared Error

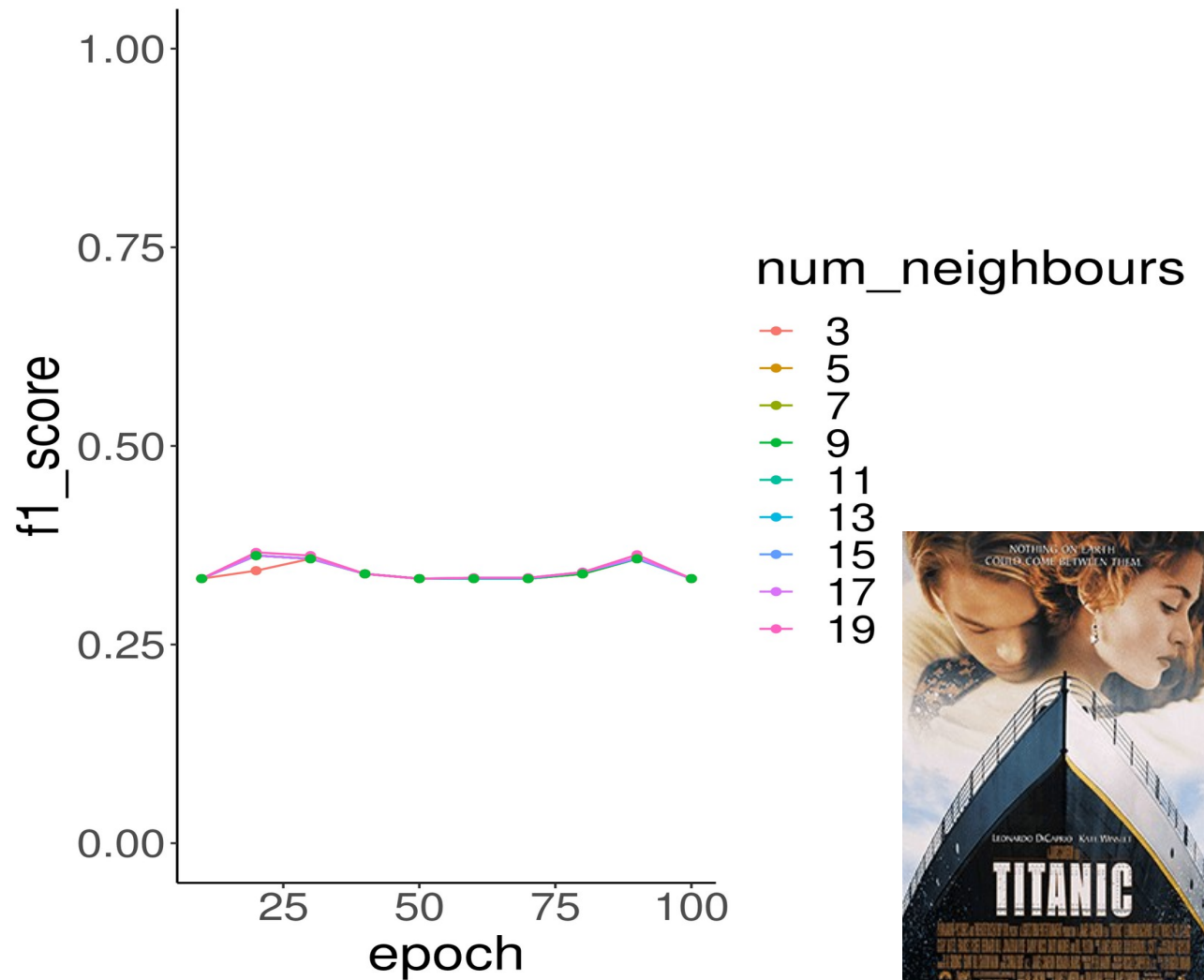
```
mean <- mean(true_values)
sd <- sd(true_values)
testthat::expect_true(sd > 0.0)
normalized_true_values <- (true_values - mean) / sd
normalized_estimated_values <- (estimated_values - mean) / sd
```





# Progress: get learning trajectory

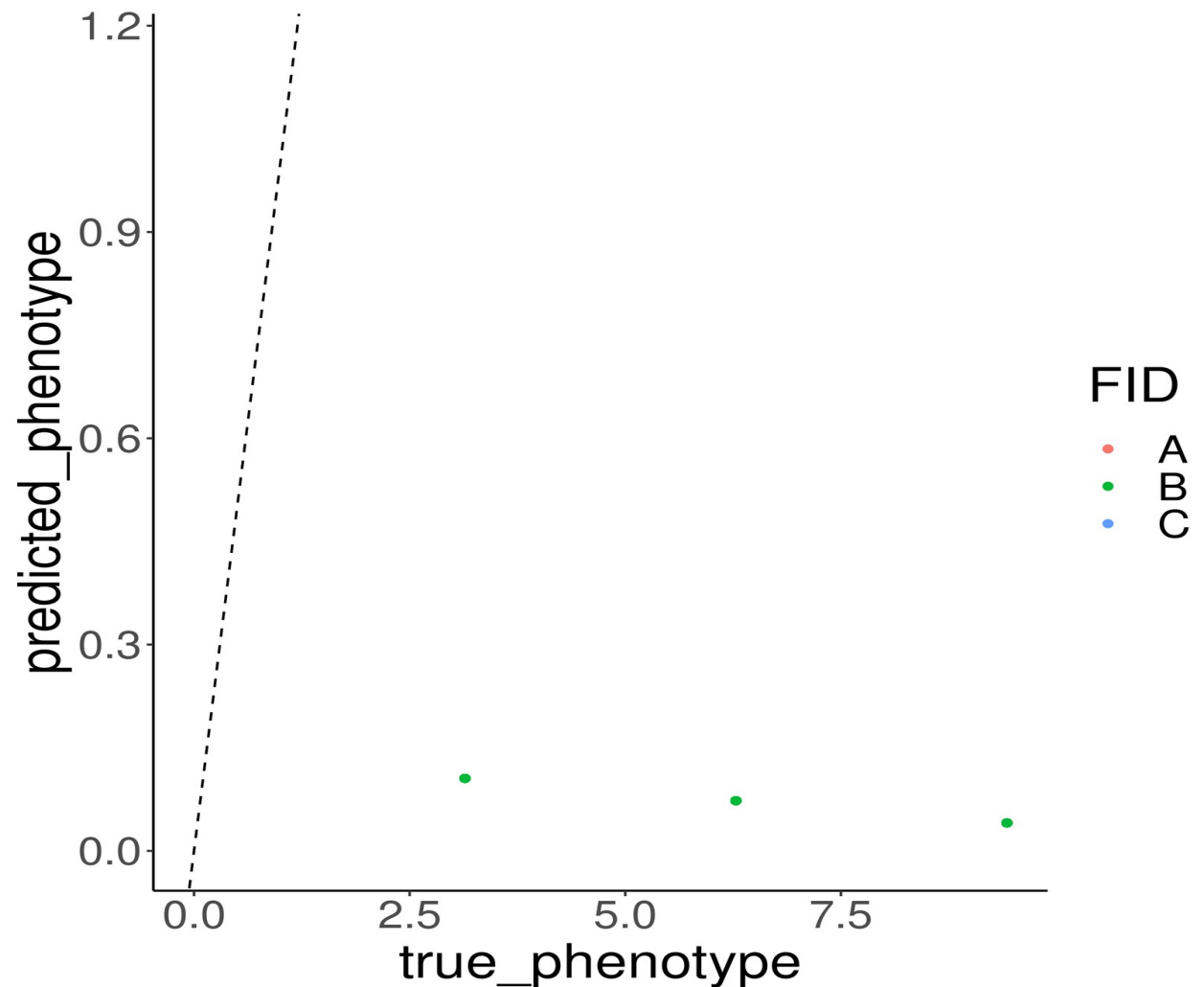
[x] Plot the F1 scores through time, #6



# Progress: get learning trajectory

**[x] Plot  
phenotype  
prediction  
through time,  
#21**

**[ ] Plot NMSE  
through time, #7**



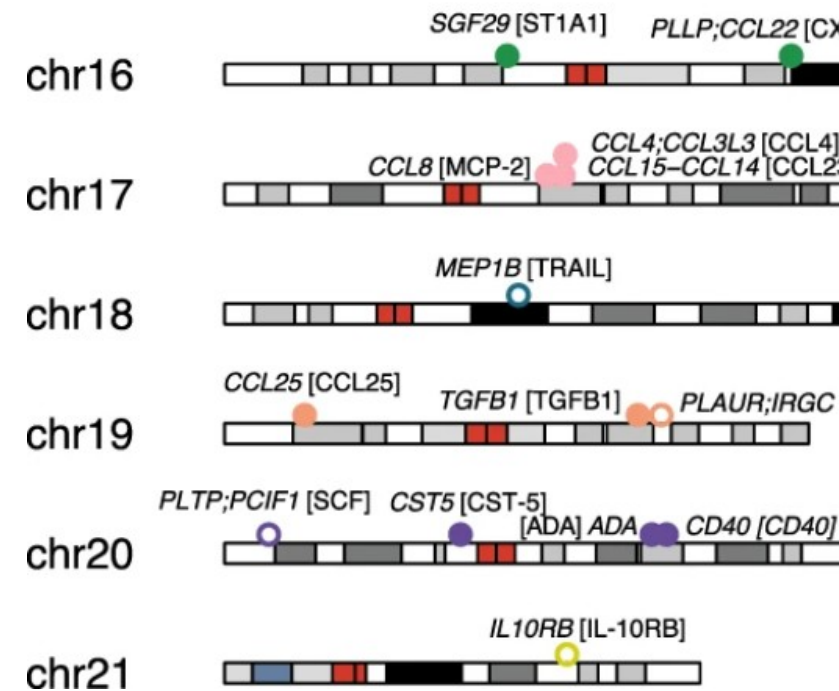
# Conclusion

## The autoencoder needs a hard problem

Hence, use NSPHS with useful genetic regions



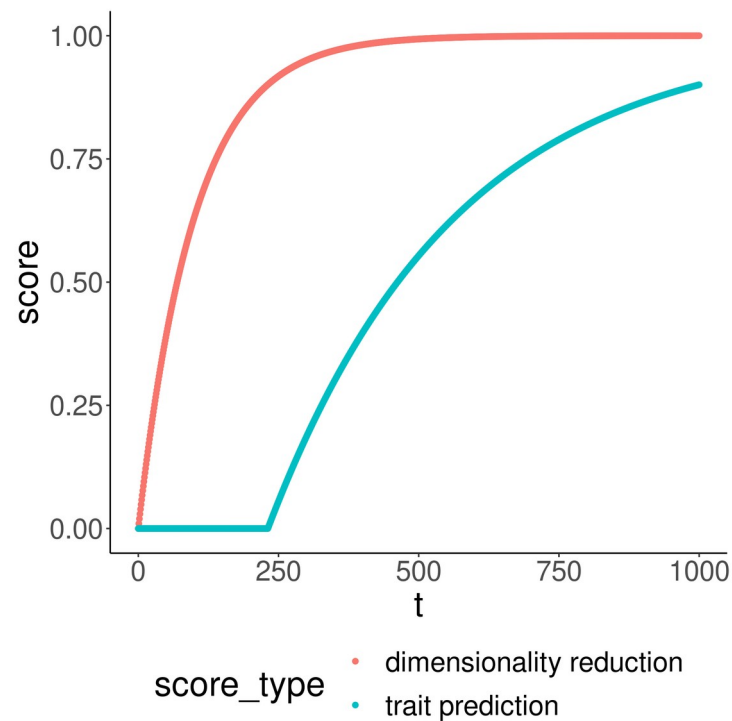
Underwhelmed



# Conclusion

## Need to determine when GCAE is done learning

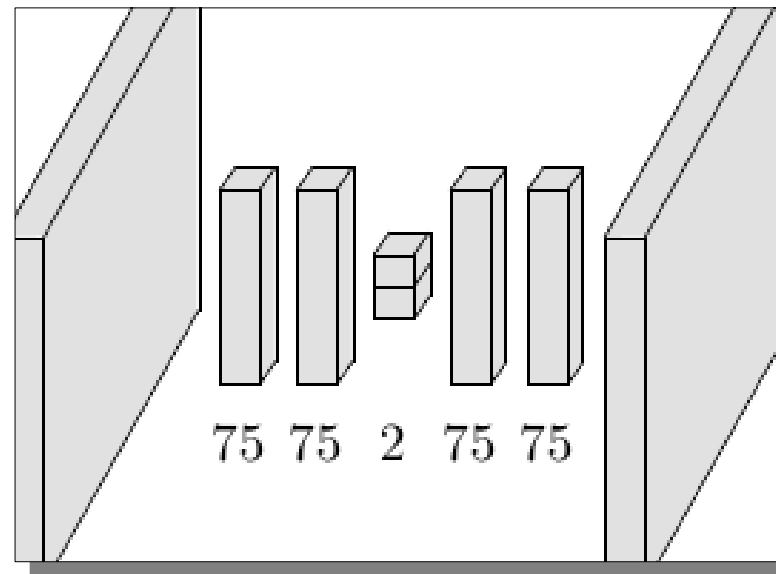
Hence, need to select/devise measures for that



# Conclusion

## The autoencoder then probably needs tuning

With a latent layer of 2 neurons (hence 2 dimensions), the autoencoder may underfit challenging data



# Discussion

**We cannot predict where/where the autoencoder outperforms existing methods ...**

**but it will be when:**

- Data of sufficient complexity

- Non-linear relations between the principal components

- Common alleles

- (Noisy data)



# Questions?



[https://github.com/richelbilderbeek/science\\_presentation\\_20220305](https://github.com/richelbilderbeek/science_presentation_20220305)

<https://youtu.be/ldwPcy263IU>