# A shortcut to the genetic basis of Neurodegenerative Diseases

Lior Sivan Schapiro

Supervisors: Prof. Michal Linial and Roni Rasnic

# The Timeline

**01**

**02**

**03**

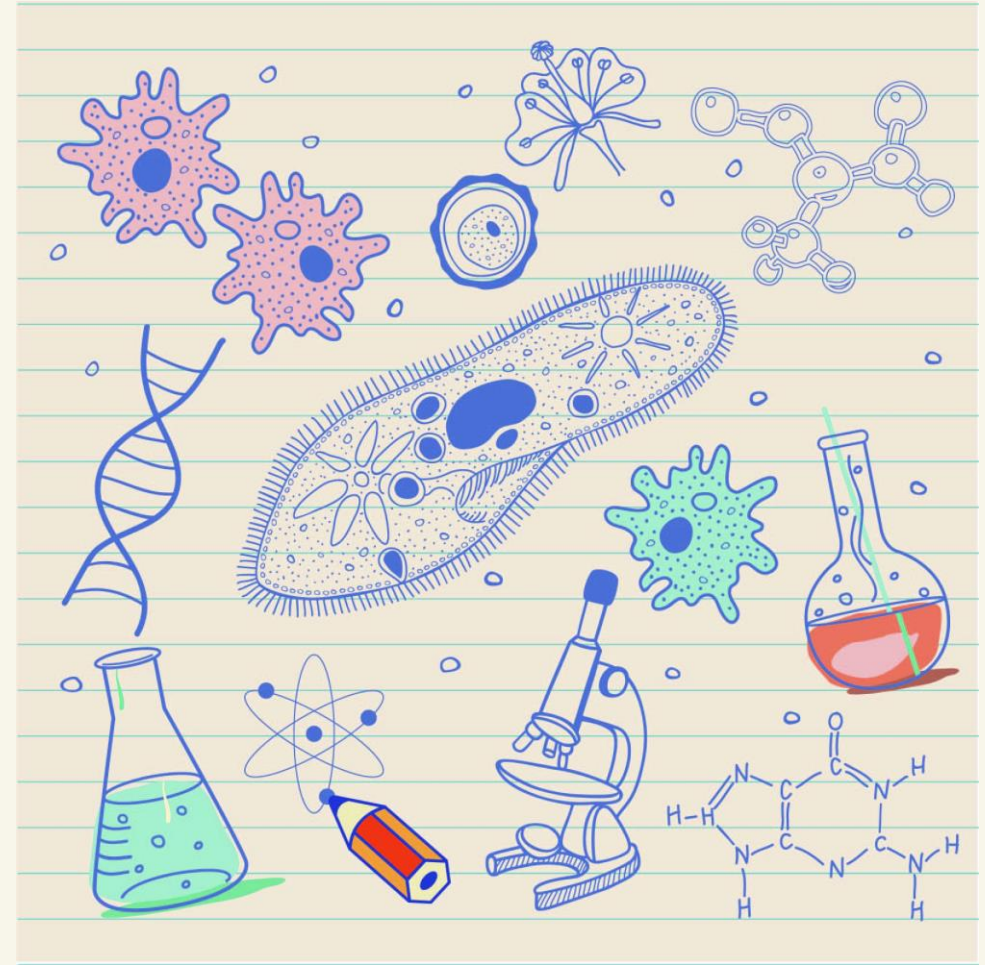**04**

**05**

**Biological Background**
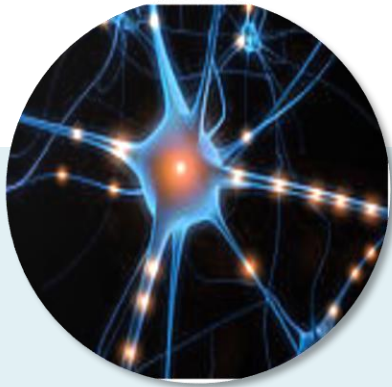
**Phase I**

**Phase II**

**Results**

**Future Work**

# 01

# Biological Background

# Neurodegenerative Diseases

Cells in the brain lose function over time and ultimately die

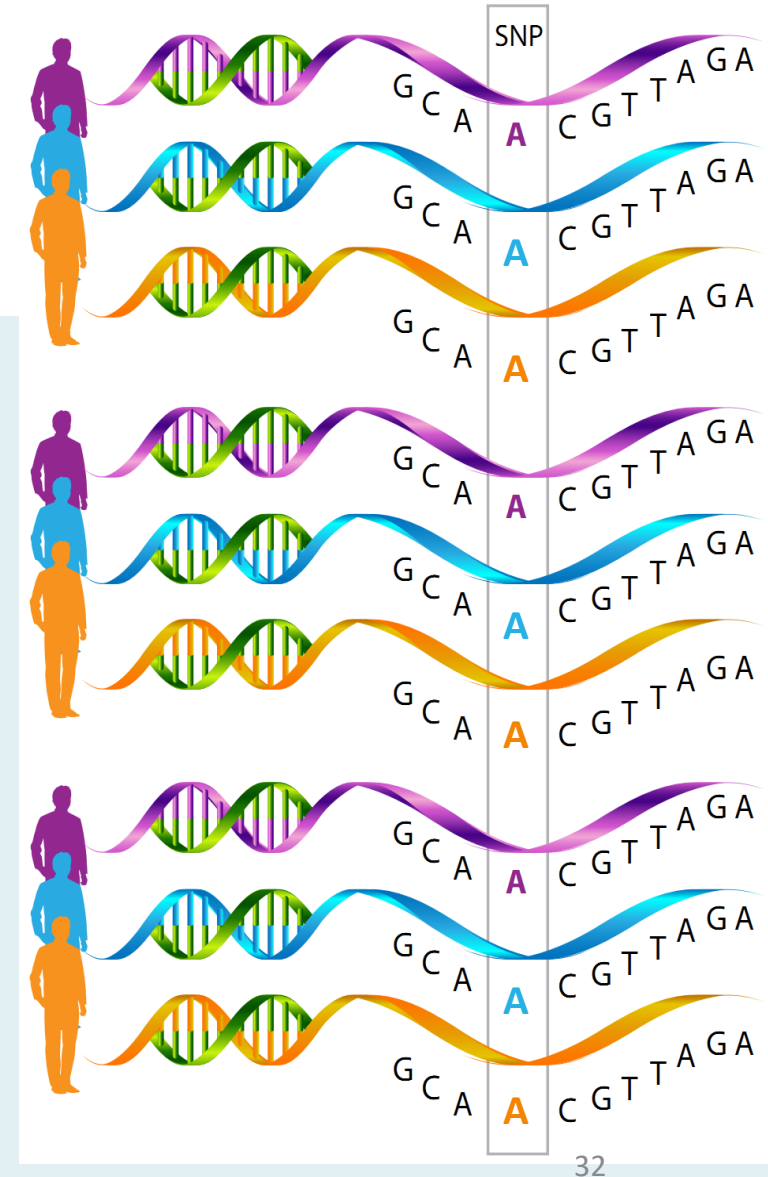**Examples-** Alzheimer's disease, Parkinson's disease (PD), ALS

The risk of being affected increases dramatically with age

# Constraint

Rank of genic functional intolerance to mutation variants using evolutionary conservation of protein sequences within species.
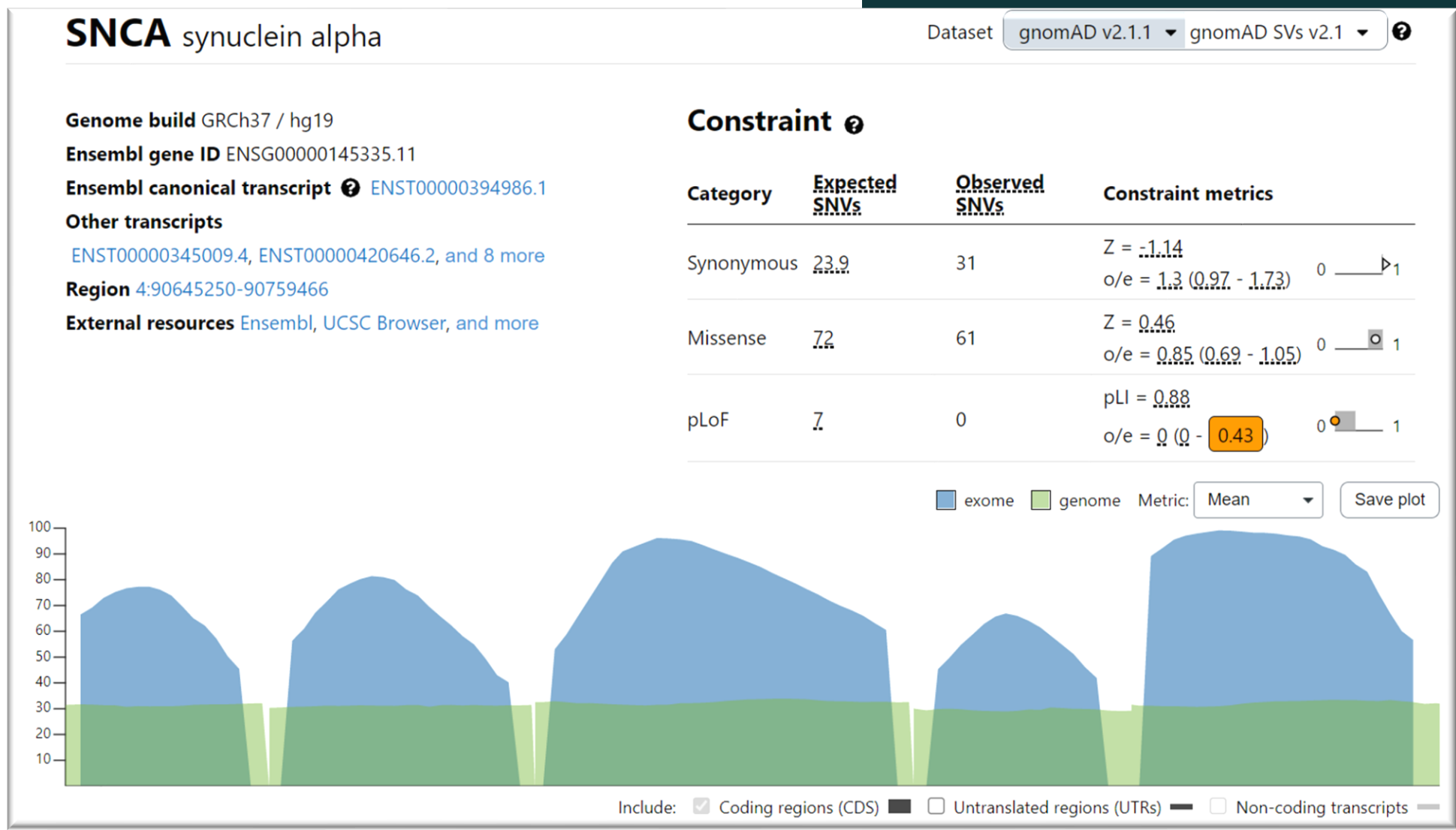
More than 150,000 human genomes (exomes) where tested

# 02

## Phase I

Can we learn from a healthy population about the genetic basis of Neurodegenerative Diseases?

# GnomAD data

# GnomAD data

change that doesn't change the amino acid

# GnomAD data



change that doesn't change the amino acid

change of a single nucleotide

# GnomAD data

change that doesn't
change the amino acid

change of a single nucleotide

# GnomAD data



change that doesn't change the amino acid

change of a single nucleotide

changes that interrupt the protein function

# Re-calculation of the pLoF score – WHY?



| Category | Expected SNVs | Observed SNVs | Constraint metrics |
|---|---|---|---|
| pLoF | 7 | 0 | pLI = 0.88<br>o/e = 0 (0 - 0.43) |

# Re-calculation of the pLoF score – WHY?

0 : 10370

0 : 19952

1 : 32299

**Population Frequencies** ❓

| Population | Allele Count | Allele Number |
|---|---|---|
| ▸ Ashkenazi Jewish | 0 | 10370 |
| ▸ East Asian | 0 | 19952 |
| ▸ European (Finnish) | 0 | 25122 |
| ▸ South Asian | 0 | 30616 |
| ▸ European (non-Finnish) | 4 | 129196 |
| ▸ Other | 4 | 7228 |
| ▸ Latino/Admixed American | 22 | 35440 |
| ▸ African/African-American | 384 | 24972 |
| XX | 223 | 129500 |
| XY | 191 | 153396 |
| **Total** | **414** | **282896** |

# Re-calculation of the pLoF score

# Re-calculation of the pLoF score

$$observedSNV_{gene} = \sum_{altarnate \in var} \sum_{lof\_var \in geneX} variant$$

## Assumption

- One of the PLOF mutation:

    transcript ablation

    Splice acceptor variant

    Splice donor variant

    Stop gained

    frameshift variant

- Alternate allele > 0 → total ++

- $\frac{Total\ allele\ number}{2}$ s.t every person has 1

➡ **2 lists of 18,000 genes with their sum of alternate allele count and total number of allele count**

# Fisher exact and Bonferroni Correction

$$p = \frac{(a+b)!\,(c+d)!\,(a+c)!\,(b+d)!}{a!\,b!\,c!\,d!\,n!}$$

$$\alpha^* = 1 - \left(1 - \frac{\alpha}{n}\right)^n$$

$p = p - value$
$a, b, c, d = values\ of\ table$
$n = set\ size$
$\alpha = given\ alpha$

43

# Using Fisher p-value as pLoF score

| | altearante allele = num of Lof var | total count = total vars | |
|---|---|---|---|
| $non-$ $neuro$ | 20 | 20,000 | 20020 |
| $candidates$ | 30 | 30,000 | 30030 |
| | 50 | 50000 | 50,050 |

18k genes $\times$

**P-value**

# Gene's constraint distribution of European



Legend:
- P-value = 1
- P-value > 0.05
- P-value < 0.05
- P-value << 0.05
- P-value = 0

Pie chart values: 93%, 1%, 4%, 1%, 1%

# Gene's constraint distribution of European



Legend:
- P-value = 1
- P-value > 0.05
- P-value < 0.05
- P-value << 0.05
- P-value = 0

1% 4% 1% 1%

93%

Most of the genes have p-value = 1 → can change and have number of variants

# Constraint in different populations



European = 1249

Ashkenazi Jewish = 242

Asian = 567

# Constraint in different populations

European

Ashkenazi jewish

1249

172

242

94

110

328

567

European = 1249

Ashkenazi Jewish = 242

Asian = 567

Asian

## 94 overlap genes

# Ontology results of the most suspected genes

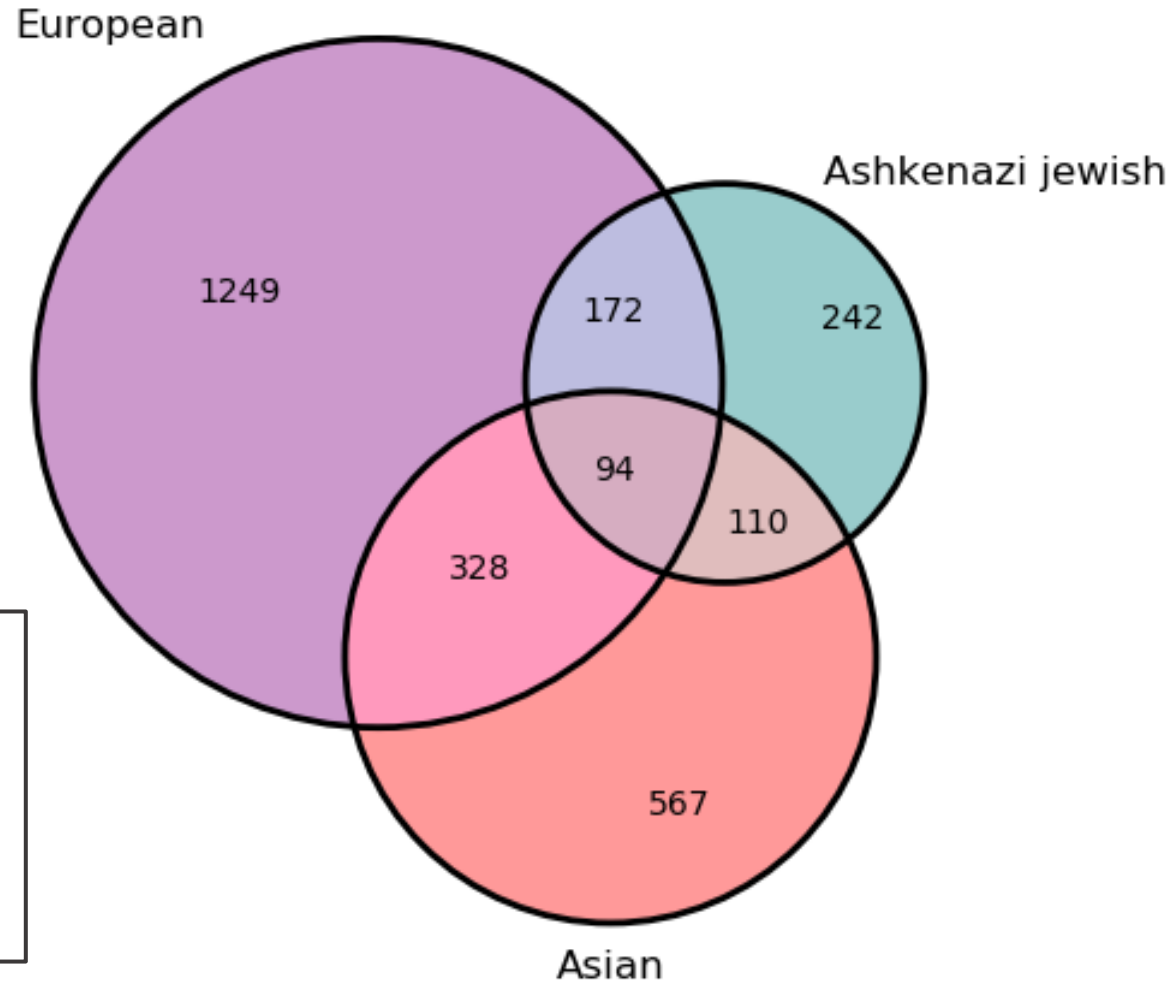| GO biological process complete | Homo sapiens (REF) # | upload_1 (▼ Hierarchy NEW! ?) # | expected | Fold Enrichment | +/- | raw P value | FDR |
|---|---|---|---|---|---|---|---|
| detection of chemical stimulus involved in sensory perception of smell | 441 | 28 | 4.28 | 6.54 | + | 1.19E-14 | 1.87E-10 |
| ⤷detection of chemical stimulus involved in sensory perception | 486 | 28 | 4.72 | 5.93 | + | 1.15E-13 | 4.53E-10 |
| ⤷detection of stimulus involved in sensory perception | 554 | 30 | 5.38 | 5.58 | + | 6.49E-14 | 3.41E-10 |
| ⤷sensory perception | 973 | 31 | 9.45 | 3.28 | + | 8.87E-09 | 1.55E-05 |
| ⤷nervous system process | 1380 | 36 | 13.40 | 2.69 | + | 7.30E-08 | 1.15E-04 |
| ⤷system process | 2040 | 39 | 19.81 | 1.97 | + | 4.30E-05 | 4.85E-02 |
| ⤷detection of stimulus | 718 | 33 | 6.97 | 4.73 | + | 2.55E-13 | 8.05E-10 |
| ⤷detection of chemical stimulus | 522 | 28 | 5.07 | 5.52 | + | 5.99E-13 | 1.58E-09 |
| ⤷sensory perception of chemical stimulus | 542 | 28 | 5.26 | 5.32 | + | 1.42E-12 | 3.20E-09 |
| ⤷sensory perception of smell | 468 | 28 | 4.54 | 6.16 | + | 4.77E-14 | 3.77E-10 |
| G protein-coupled receptor signaling pathway | 1329 | 34 | 12.91 | 2.63 | + | 2.79E-07 | 4.00E-04 |
| cellular component organization | 5775 | 29 | 56.08 | .52 | - | 7.90E-06 | 9.58E-03 |
| ⤷cellular component organization or biogenesis | 5999 | 29 | 58.26 | .50 | - | 2.06E-06 | 2.70E-03 |
| cellular nitrogen compound metabolic process | 3407 | 13 | 33.09 | .39 | - | 4.63E-05 | 4.87E-02 |
| ⤷cellular metabolic process | 7570 | 35 | 73.51 | .48 | - | 3.53E-09 | 6.97E-06 |

# 03

# Phase II

Can Early Onset of Complex Diseases

be a hint for Etiology?

# Parkinson in numbers



**60%**  **40%**

### Age distribution



**10M**
people worldwide

**60k**
new cases every year

**52B$**
every year

# GWAS - genome wide association study



**Genotype**

rs12349

AACGCTGCTAAT

**Phenotype**

Cohort

# GWAS - genome wide association study

# What data did we have from GWAS?

**Late: 1236** variants
→ 43 genes

**Early: 2254** variants
→ 161 genes

**4** overlap genes (all are known as related)

# FABRIC

(**F**unctional **A**lteration **B**ias **R**ecovery **I**n **C**oding-regions)

- Assess the **impact** of mutations on gene/protein **function**

- Find genes more damaged than expected
  - Given this number of *random* mutations

# FABRIC

# FABRIC

# FABRIC



C

All possible single-nucleotide
variants in the gene

CAAAACAAAAAA...
GCCGCGCGCCGC...
TGTTTTGTGTTG...
ATGCGATCTGCT...

Expected distribution

# FABRIC



C

All possible single-nucleotide variants in the gene

CAAAACAAAAAA...
GCCGCGCGCCGC...
TGTTTTGTGTTG...
ATGCGATCTGCT...

Expected distribution

D

Observed-to-expected analysis

p = 2E-13

Z = -1.62

# FABRIC output

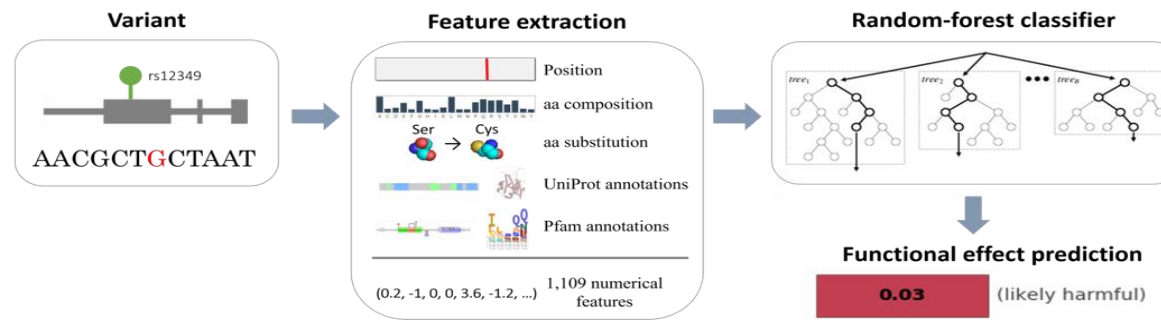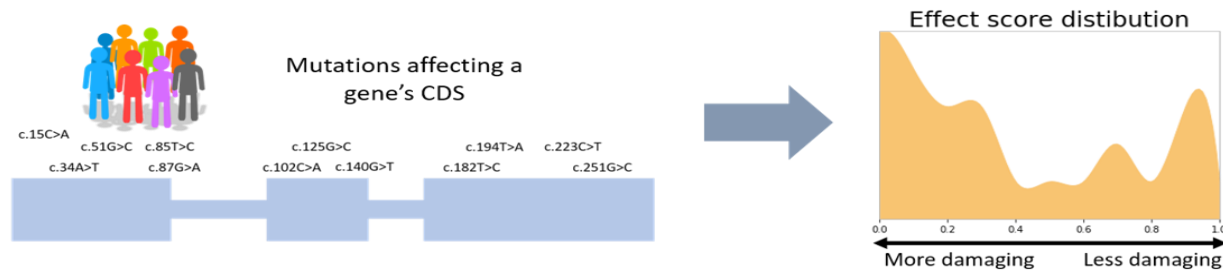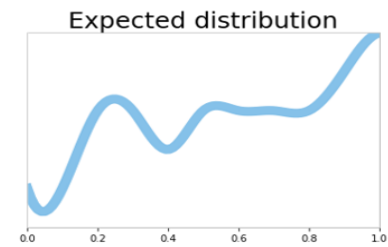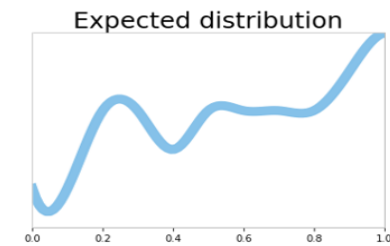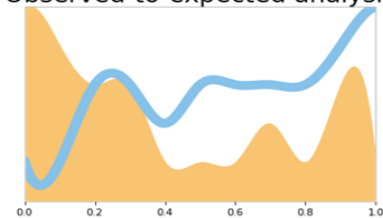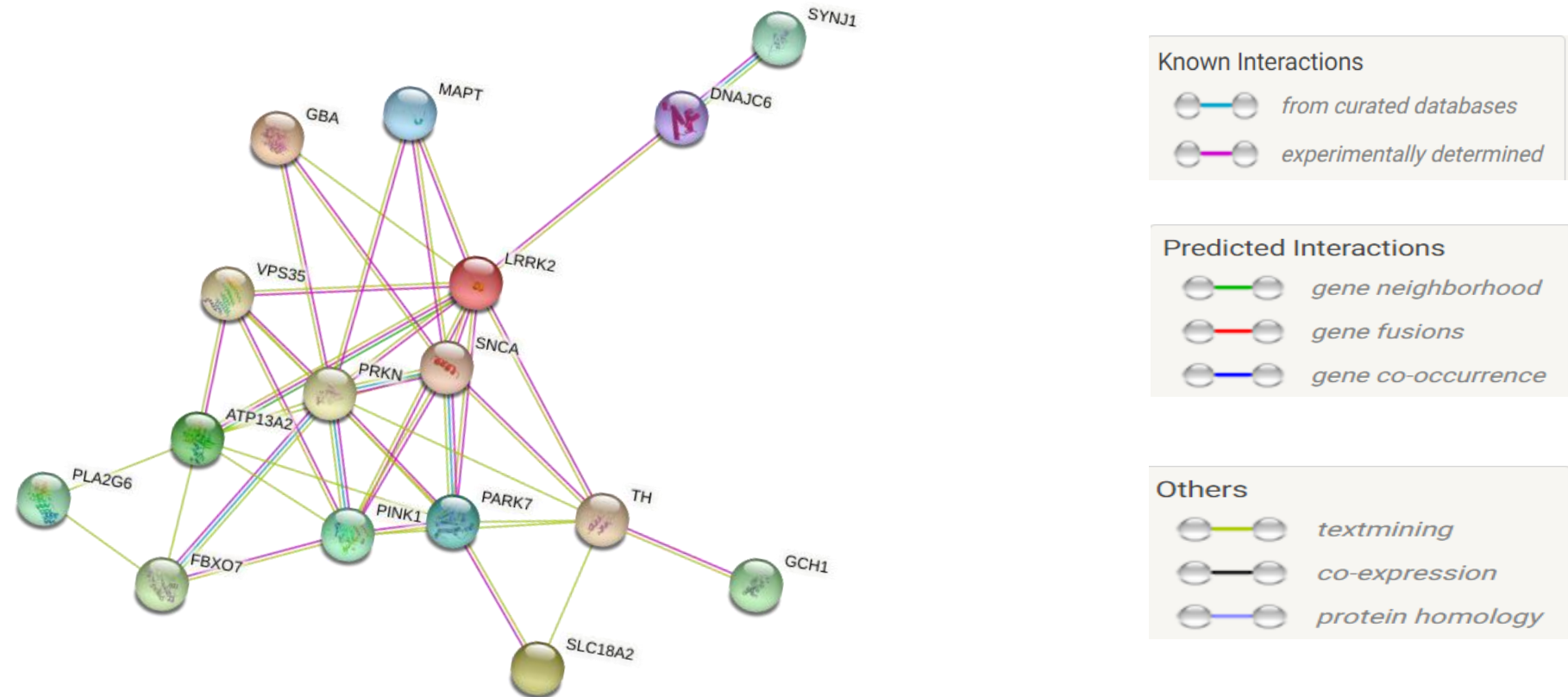| | uniprot_id | symbol | chr | overall_z_value | overall_pval | overall_fdr_significance | overall_fdr_qval |
|---|---|---|---|---|---|---|---|
| 1 | uniprot_id | symbol | chr | overall_z_value | overall_pval | overall_fdr_significance | overall_fdr_qval |
| 2 | S4R3Y5 | MTRNR2L1 | 1 | -2.220891562 | 0.111111111 | FALSE | 0.228958549 |
| 3 | Q4KMX7 | FAM106A | 17 | -1.28556089 | 0.061311263 | FALSE | 0.142174949 |
| 4 | Q5JQF8 | | X | -1.090634791 | 0.068965517 | FALSE | 0.156388068 |
| 5 | A0A075B6P5 | IGKV2-28 | 2 | -1.076343589 | 0.432989691 | FALSE | 0.613729799 |
| 6 | P01593 | IGKV1D-3 | 2 | -0.850532119 | 0.121932632 | FALSE | 0.245920771 |
| 7 | Q86YR6 | POTED | 21 | -0.845584393 | 0.086815871 | FALSE | 0.188505458 |
| 8 | S4R3P1 | MTRNR2L4 | 4 | -0.837559411 | 0.234137303 | FALSE | 0.404099638 |
| 9 | Q8NHZ8 | CDC26 | 9 | -0.834219789 | 0.00085341 | TRUE | 0.003540289 |
| 10 | O95013 | OR4F21 | 8 | -0.772694291 | 0.247840382 | FALSE | 0.420351121 |
| 11 | Q9BTY7 | HGH1 | 8 | -0.759830842 | 0.488584475 | FALSE | 0.661896051 |
| 12 | P01624 | IGKV3-15 | 2 | -0.743705832 | 0.641509434 | FALSE | 0.78273921 |
| 13 | Q5EBN2 | TRIM61 | 4 | -0.674288599 | 0.540540541 | FALSE | 0.705391255 |
| 14 | P01597 | IGKV1-39 | 2 | -0.628736726 | 0.626506024 | FALSE | 0.772964802 |
| 15 | Q8NH02 | OR2T29 | 1 | -0.62486385 | 0.219727393 | FALSE | 0.385156066 |
| 16 | Q6NT46 | GAGE2A | X | -0.622272709 | 0.291996584 | FALSE | 0.472044923 |
| 17 | B0FP48 | UPK3BL | 7 | -0.620947386 | 0.372608163 | FALSE | 0.554923327 |
| 18 | Q9UGB4 | C20orf187 | 20 | -0.618681218 | 0.19980723 | FALSE | 0.357822628 |
| 19 | A6NI03 | TRIM64B | 11 | -0.605519673 | 0.236615436 | FALSE | 0.40708471 |
| 20 | Q8NG35 | | 8 | -0.577461066 | 0.370034572 | FALSE | 0.552833868 |
| 21 | P0CV98 | TSPY3 | Y | -0.568874036 | 0.614886731 | FALSE | 0.765105904 |
| 22 | A6NE82 | MBD3L3 | 19 | -0.566757434 | 0.456831032 | FALSE | 0.635034059 |
| 23 | Q9UND3 | NPIPA1 | 16 | -0.559052222 | 0.470638017 | FALSE | 0.646784763 |
| 24 | O43261 | DLEU1 | 13 | -0.544280904 | 0.117836461 | FALSE | 0.239435363 |
| 25 | Q96P64 | AGAP4 | 10 | -0.544044917 | 0.36951088 | FALSE | 0.552235365 |

# Network of Parkinson genes identified by FABRIC



25 genes were marked as "significant" (positive selection) for early onset, none for late onset

# 04

# What's Next?

# What's Next?

# What's Next?

Run fabric on the 94 constraint genes in order to test how damaged, they are.

# What's Next?

Run fabric on the 94 constraint genes in order to test how damaged, they are.

Compare between other population & Run on more diseases like Alzheimer

# What's Next?

Run fabric on the 94 constraint genes in order to test how damaged, they are.

Compare between other population & Run on more diseases like Alzheimer

Analyze exomes from UK bio-bank & build predictor for individuals

# **Michal, Amir and Roni**

# Questions?