# Machine Learning for early disease diagnostics
## Himani Agrawal, PhD

Galvanize

Simpatica Medicine

Tremors

Stiffness

Changes in speech, voice and swallowing

Balance problems

Trouble with handwriting

Slowness of movement

Parkinson's disease:
- disease of the nervous system
- affects more than 10 million people
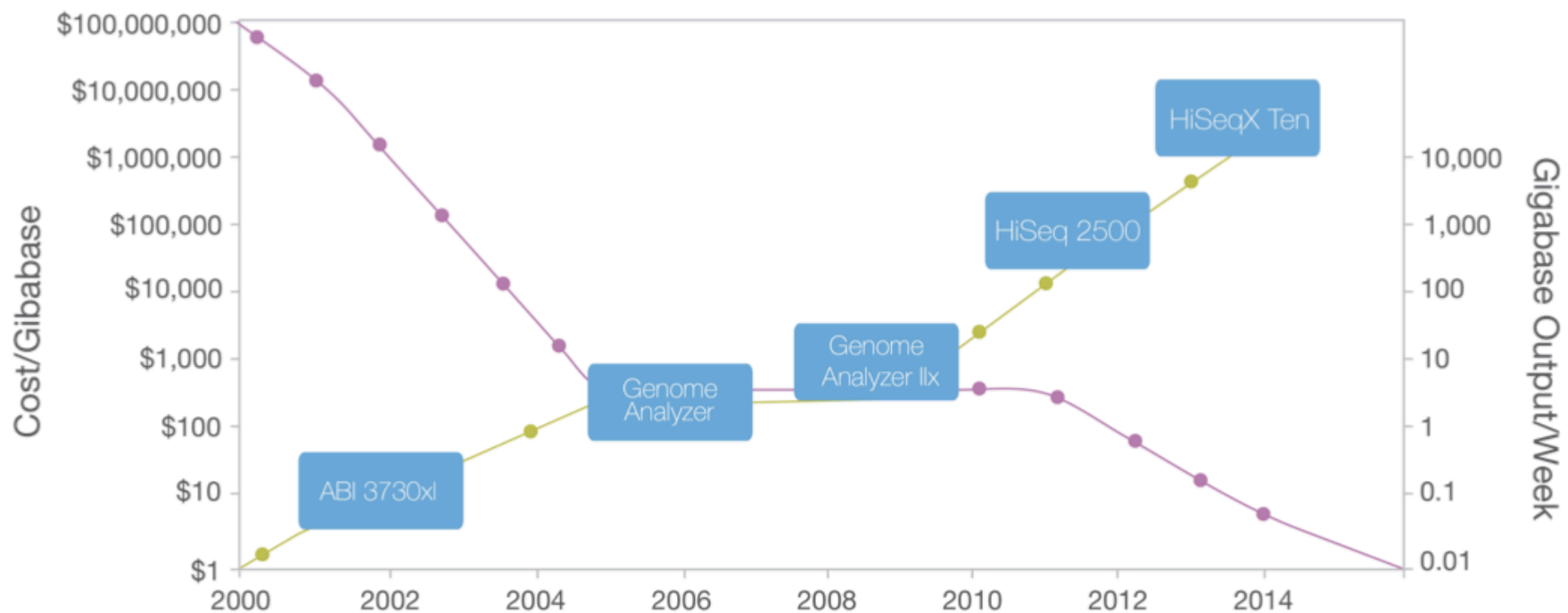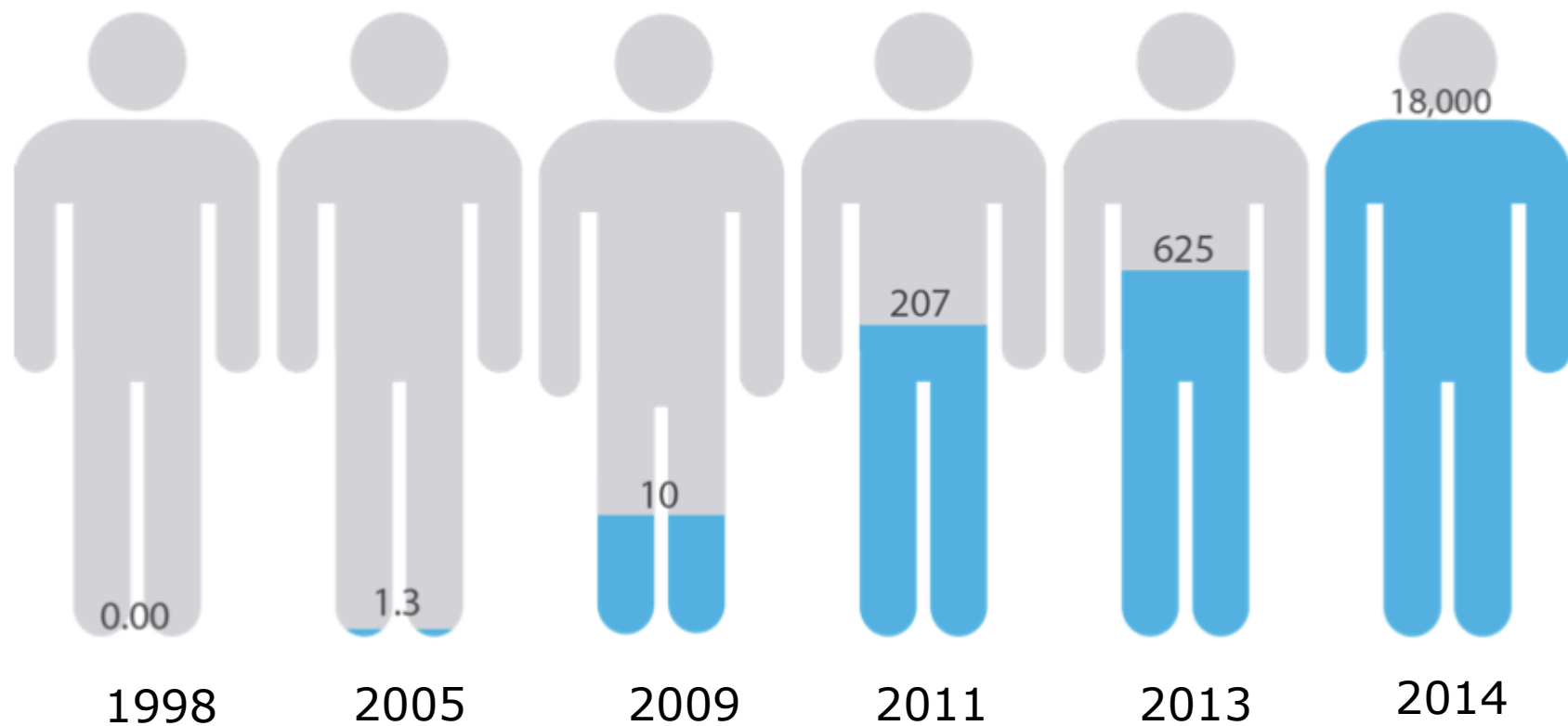- only 5% people diagnosed before age 50
- reduced dopamine

# The plummeting cost of genome sequencing, rise of genome data

# Human genomes sequenced annually



0.00 — 1998
1.3 — 2005
10 — 2009
207 — 2011
625 — 2013
18,000 — 2014

# RNA Sequencing

# Machine learning for disease detection



Machine learning-based brain disease diagnosis

Brain Health

ML for early detection

Pre-Clinical
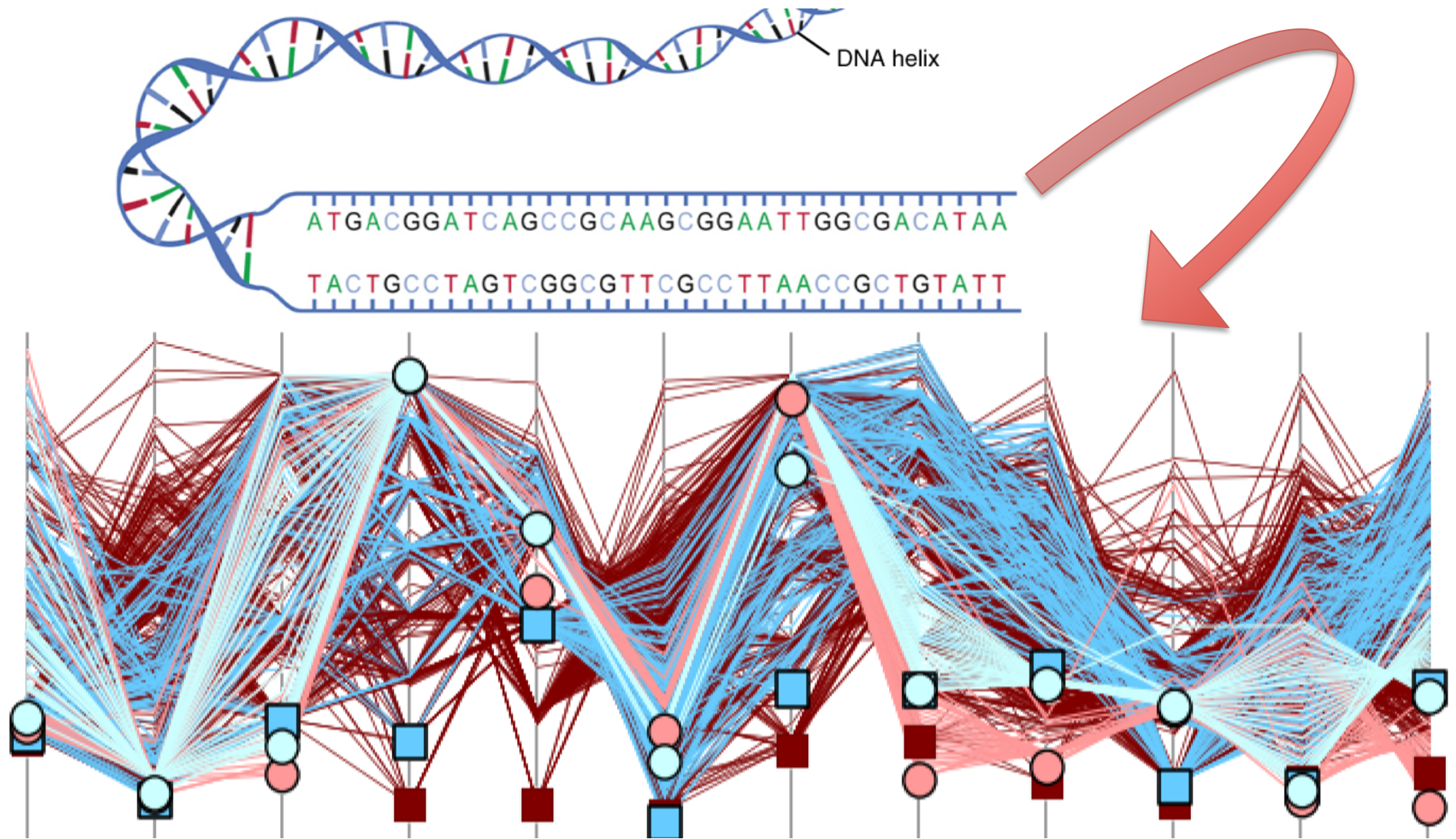
Major Symptoms

Parkinson's disease

BIG DATA

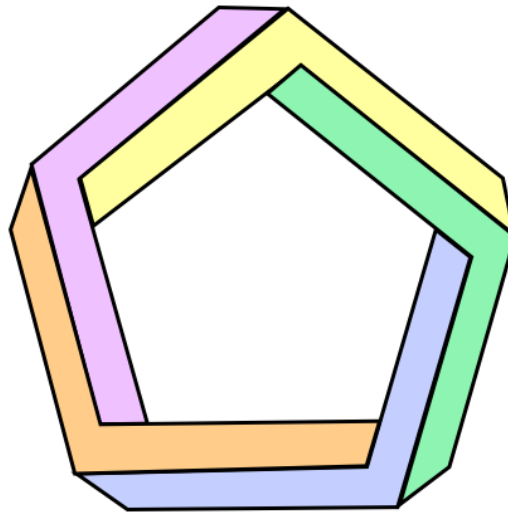http://mlcenter.postech.ac.kr/healthcare

# High dimensional genomics data

# High dimensional genomics data
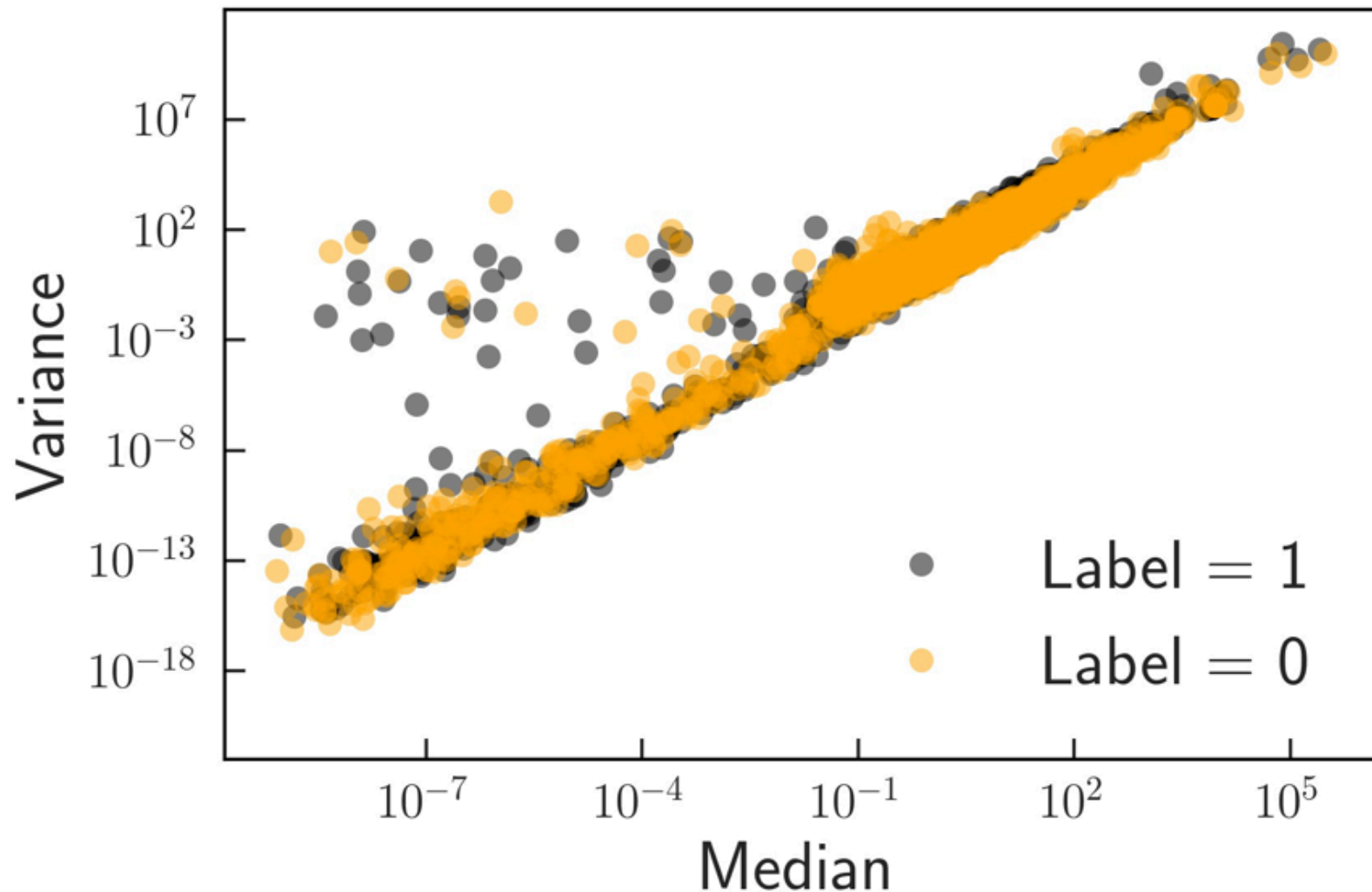
Needs Visualization

Needs veracity

High variability

High volume

High variety

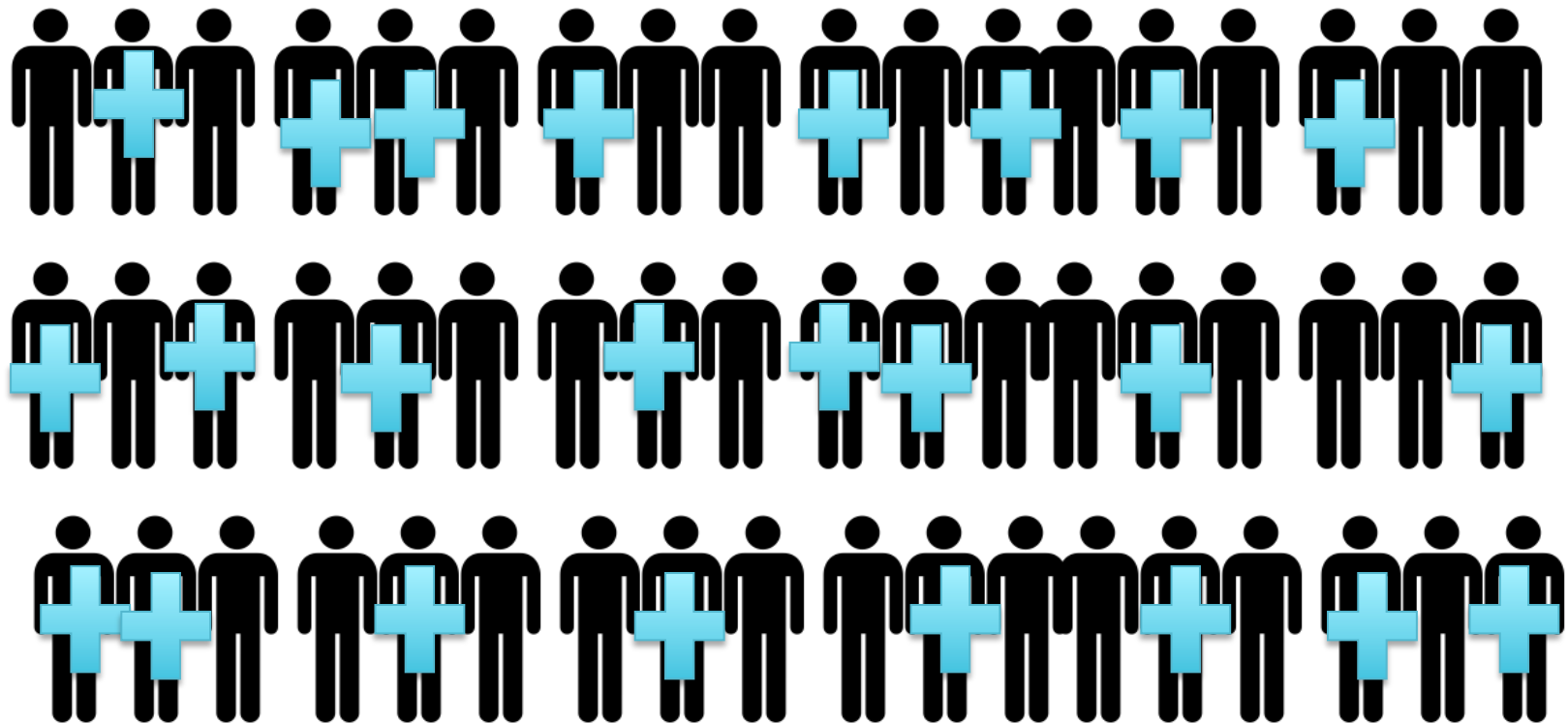# High variability: genomics data

# High variety, high volume: medical data

- Different formats: text, numeric, paper, digital, pictures, videos, multimedia

- Electronic Health Records

- Medical Images

- Fitness tracking from mobile apps

- ….etc.
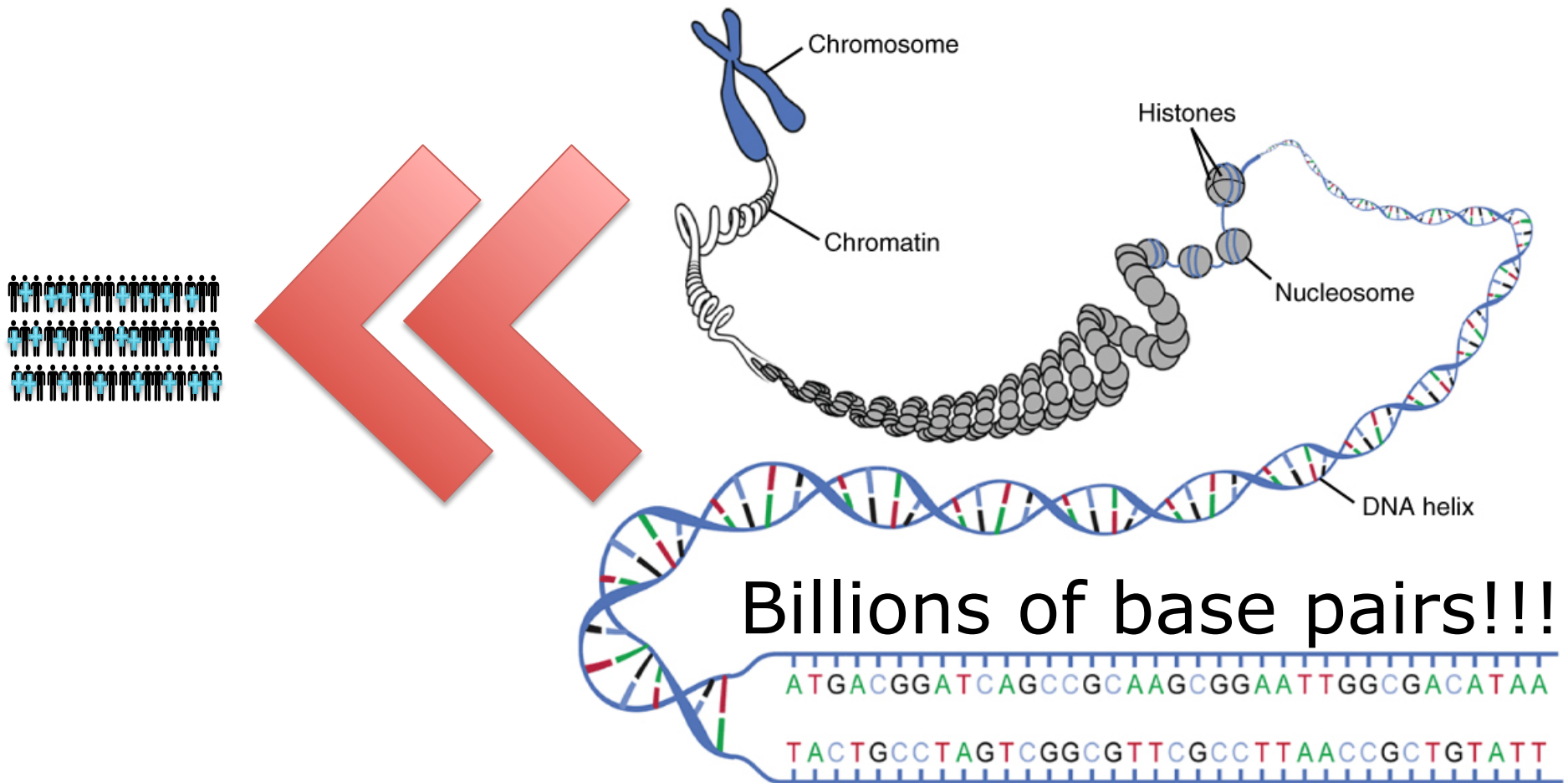
# Dataset characteristics

- 53 patients, 60000 features (gene expression level),
- ✚ implies presence of Parkinson

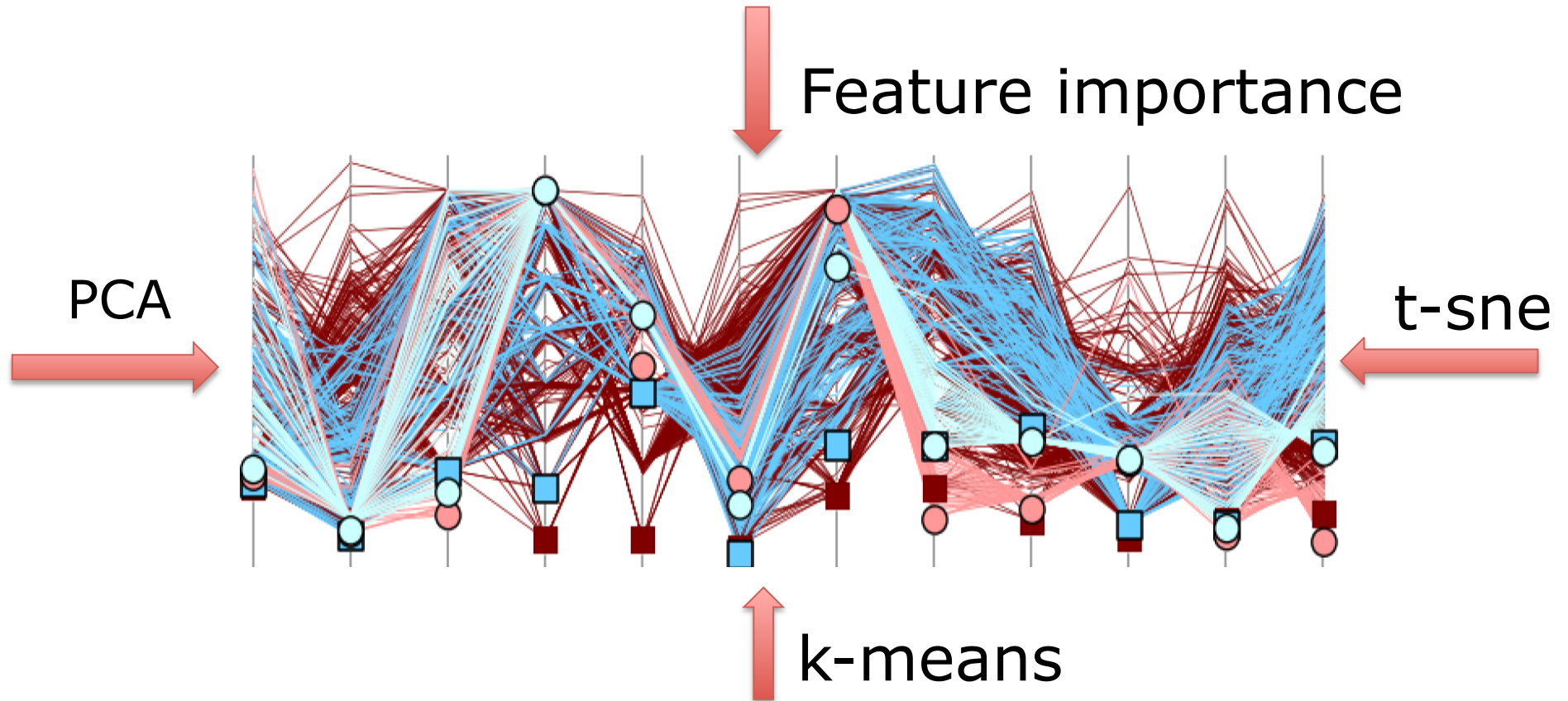# High dimensional analytics
# 53 << 60,000



Billions of base pairs!!!

… https://i.stack.imgur.com/ibE8u.jpg

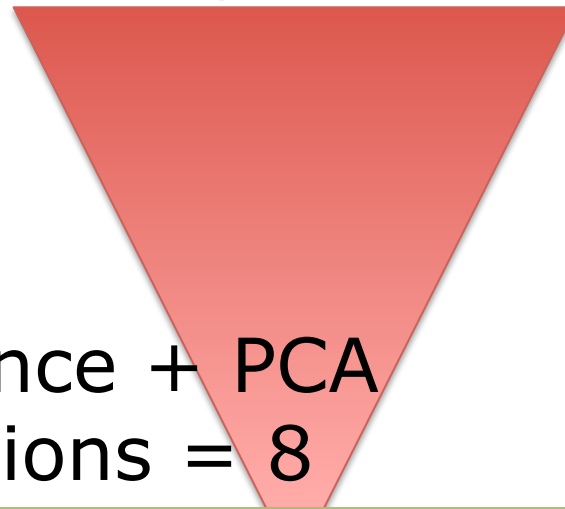# Dimensionality reduction: Smush the data !!!



- Principal Component Analysis (PCA)
- Feature importance using Random Forest
- K-means clustering

# Machine learning model

Original Data

Feature importance + PCA
Reduced dimensions = 8
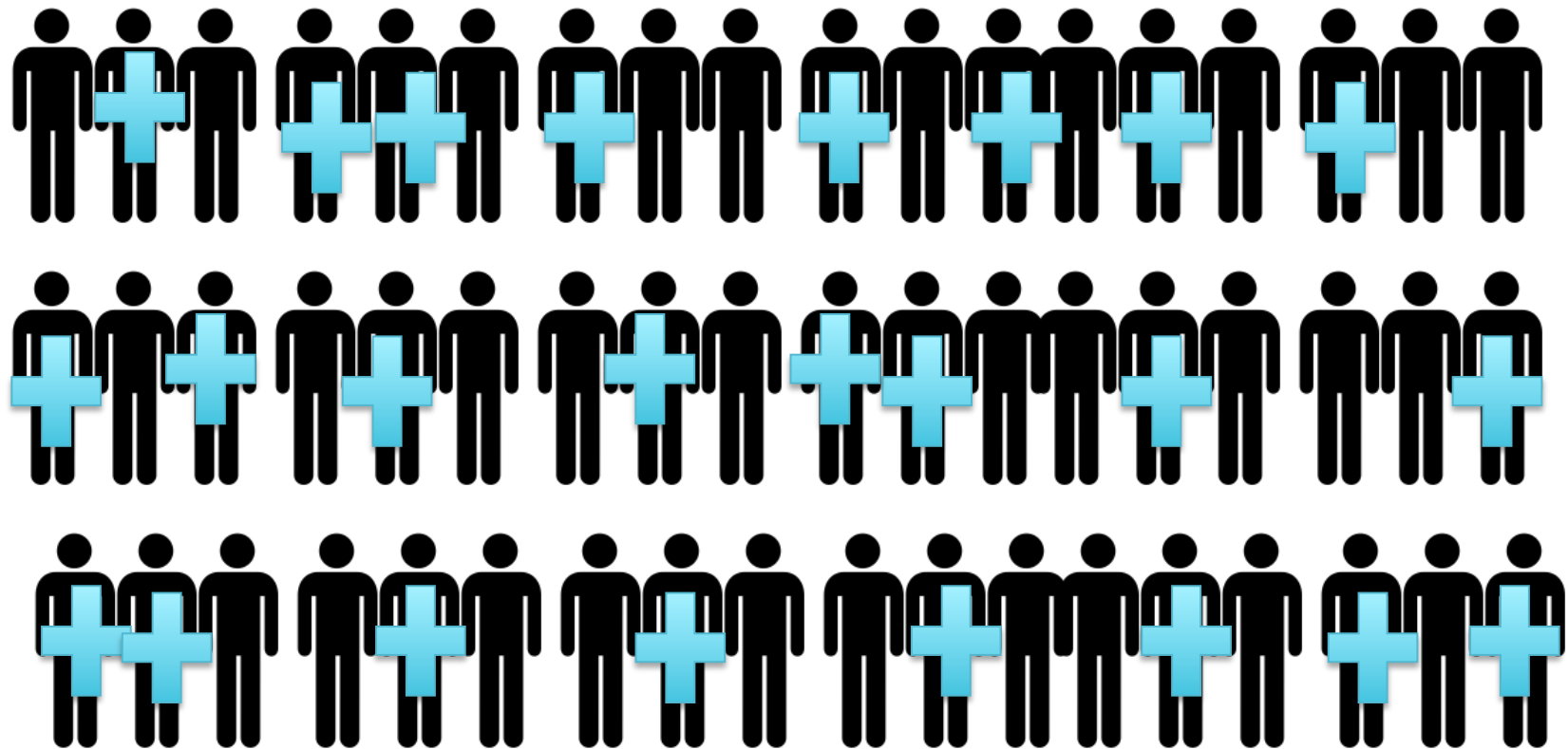
Gradient Boosting Classifier

75 % cross-validated accuracy
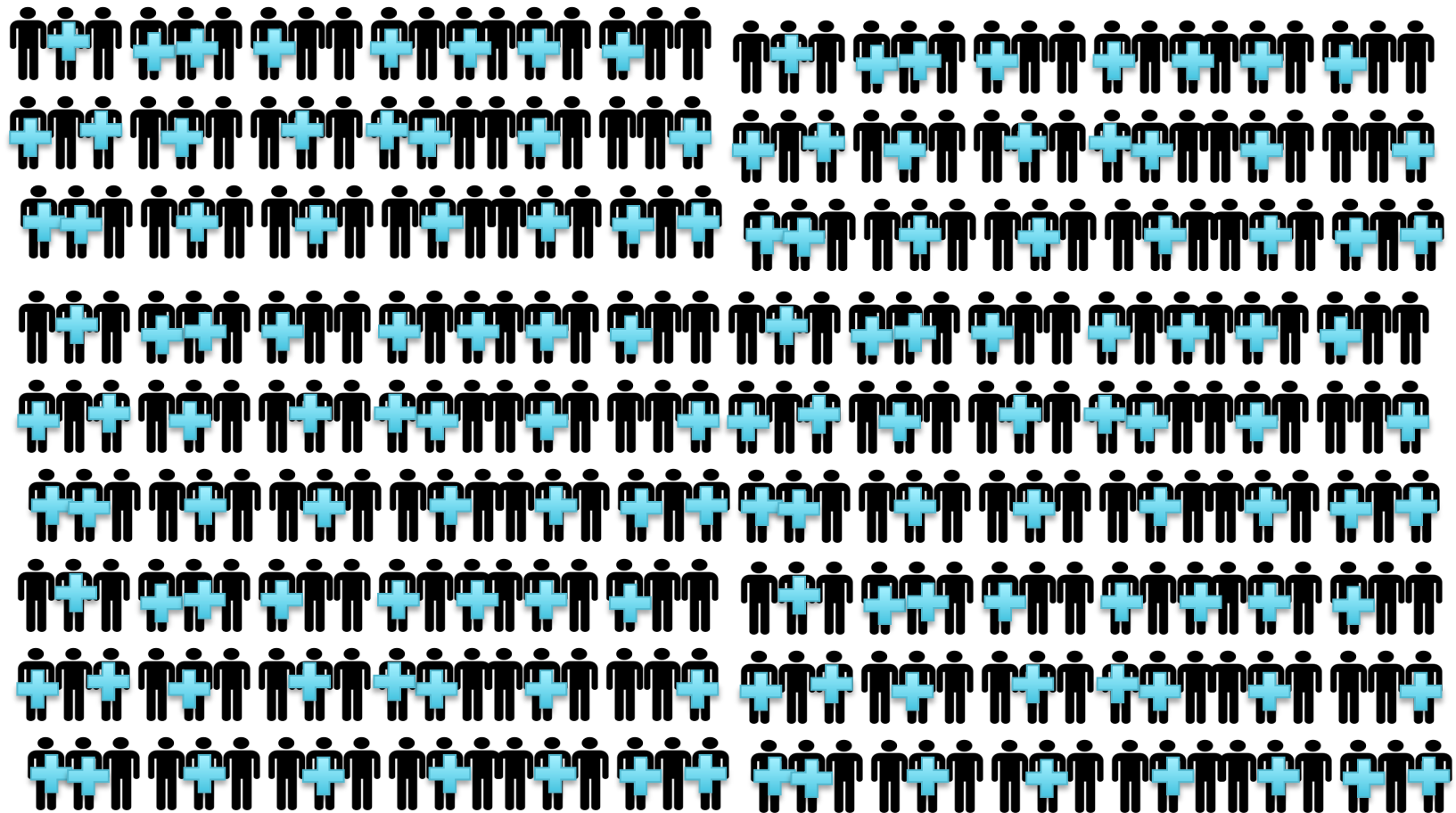
# Importance

- There is no definitive diagnostic of Parkinson's disease
- Doctor looks into medical history, neurological exam, major symptoms
- Response to Parkinson's medications for further evidence
- Imaging system for dopamine levels
- But no accurate diagnosis, particularly in early stages --- ML can help!

# Future work

# Small sample size

# Higher sample size, more robust model

# Incorporate published findings into ML model

**New gene shown to cause Parkinson's disease**

Third gene definitely linked to disease in patients from North America, Asia

| | |
|---|---|
| *Date:* | June 6, 2016 |
| *Source:* | Northwestern University |
| *Summary:* | Scientists have discovered a new cause of Parkinson's disease -- mutations in a gene called TMEM230. This appears to be the third gene definitively linked to confirmed cases of the common movement disorder. |

Mutant genes responsible for Parkinson's disease

Weidong Le, Stanley H Appel ✉

⊞ **Show more**

t rights and content

## Genetics of Parkinson disease

Nathan Pankratz[1] and Tatiana Foroud[1]

[1]Department of Medical and Molecular Genetics, Indiana University School of Medicine, Indianapolis, Indiana

# Thank you!!

Questions

Web site
hagrawal.me.uh.edu

e-MAIL
1991himani@gmail.com

LinkedIN
https://www.linkedin.com/in/
himaniagrawalgalvanize/

VISIT ME
AT