# Bioinformatics

## from **genotype** to **phenotype**

April 3rd. **Data Rave. @damiankao**

# What is it

**bioinformatics**

**computational biology**
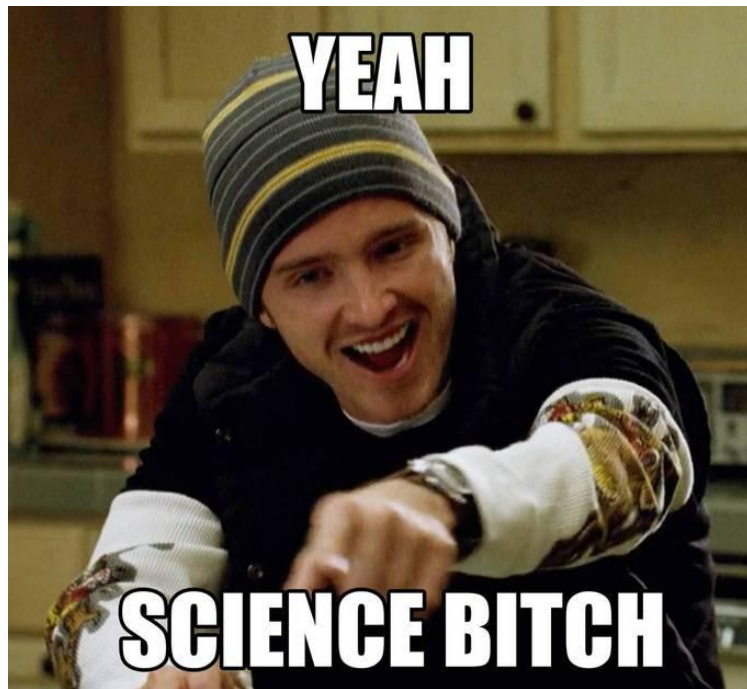
**quantitative biology**

**biostatistics**

# What is it

bioinformatics

computational biology

Science

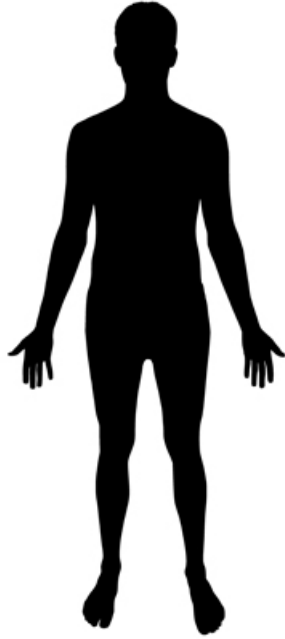biostatistics


YEAH

SCIENCE BITCH

# Genotype vs phenotype

**every cell in an organism contains a full set of genome**

```
ATGAGATAGAGATAGCCCTA
TACTCTATCTCTATCGGGAT

GACTAGATAAAGACAGA..
CTGATCTATTTCTGTCT..
```

- 4 bases (A, G, T, C)
- human contains ~3 billion bases
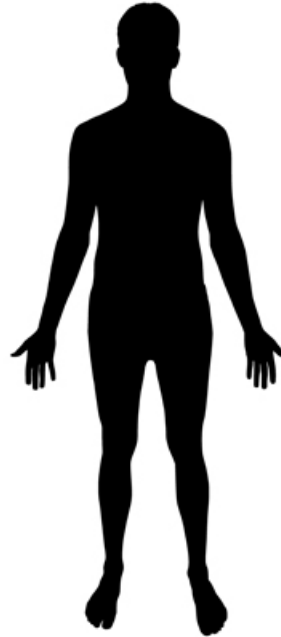- genotype is information

# Genotype vs **phenotype**

- Phenotype is a vague term
- It is any feature or trait that is NOT genotype
- Hair color, height, blood type, rate of production of a certain protein, fertility….etc
- Can be hard to quantify
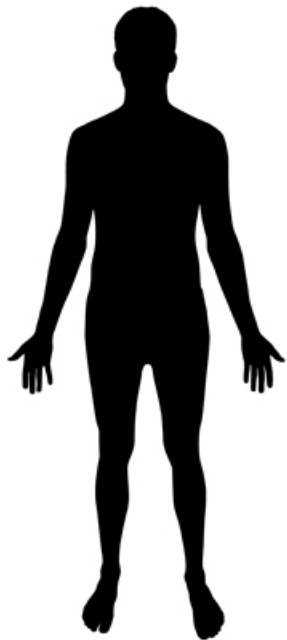
# **Genotype to phenotype**

ATGAGATAGAGATAGCCCTA
TACTCTATCTCTATCGGGAT

GACTAGATAAAGACAGA..
CTGATCTATTTCTGTCT..

- ~ 3 billion bases for humans
- ~ 1 gigabyte of data so far

- How do you quantify all these dimensions?
- Probably a lot of data?

# Genotype to phenotype



**phenotype** data is algorithmically compressed into **genotype** data

greatest compression algorithm ever??

ATGAGATAGAGATAGCCCTA
TACTCTATCTCTATCGGGAT

GACTAGATAAAGACAGA..
CTGATCTATTTCTGTCT..

# Genotype to phenotype



**Genetics = How does this compression work?**

# Layers of complexity

There are many layers of complexity above genotype that eventually leads to phenotype.

AGATCAGTTAGTCTAGTAGTGGCGCCCGCTAATATACGCGGC

**Transcription**

**Translation**

# Layers of complexity



translated protein → folded protein → cells → tissues → organs → organism → population → species → planet

# Sequencing the genome

- There are billions of bases and bunch of chromosomes


- In a perfect world, we would get a really long string of bases representing each chromosome.

# Sequencing the genome?

- Our sequencing technology isn't there yet (maybe soon)


- We can sequence a lot of short fragments

# **Sequencing the genome**



We take a sample from you (ie. cheek swab). This consists of millions of cells.

We extract the DNA from all of the cells.

We chop them up into random smaller sized fragments

Sequencing machine

# Sequencing the genome

Genome

AGATCAGTTAGTCTAGTAGTGGCGCCCGCTAATATACGCGGCGCGATTACTGTCTGTATAAGTATGTCGTGTGTAGTGCTGTCGTA

**As a result of sequencing, we get back small fragments of the genome**

# Result file (.fasta)

```
>Fragment01
AGGTTAGGTTTTAGCTTGATGCTTAGCTTGATGCAGTATTATGTATCGTATCGTATATGTCGA
>Fragment02
GGTTTTAGCTTGATGCTTAGCTTGATGCAAGTTAGTCGTTAGTCGTTAGTCGTAGTGATG
>Fragment03
TGCTTAGCTGCTTAGCTTGATGCAAGTTAGTCGTTGCTTAGCTTGCTTAGCTTCGTTAG
…
…
```

# Computational problems

## Assembly of fragments

**fragments**

**assembled**

## Alignment of fragments

**reference sequence**

**alignment of fragment to reference**

# Analysis problems

## Variant calling

**Bob** A G T C **C** G T T A T T **G** A C T T C G T A G T C

**Tom** G

**Mary** A G T C **C** G T T A T T **G** A C T T C G T A G T C

**Sam** G

A G T C **A** G T T A T T **T** A C T T C G T A G T C

## Correlation to phenotype

A G T C **A** G T T A T T **G** A C T T C G T A G T C

G

# Data engineering problems

- **Compression** of raw data.

- **Querying** the data.

- **Standardization** of formats.

- **Accessibility** of the data.

# DNA-seq

DNA sequencing (DNA-seq) is NOT the only type of sequencing that is being done.
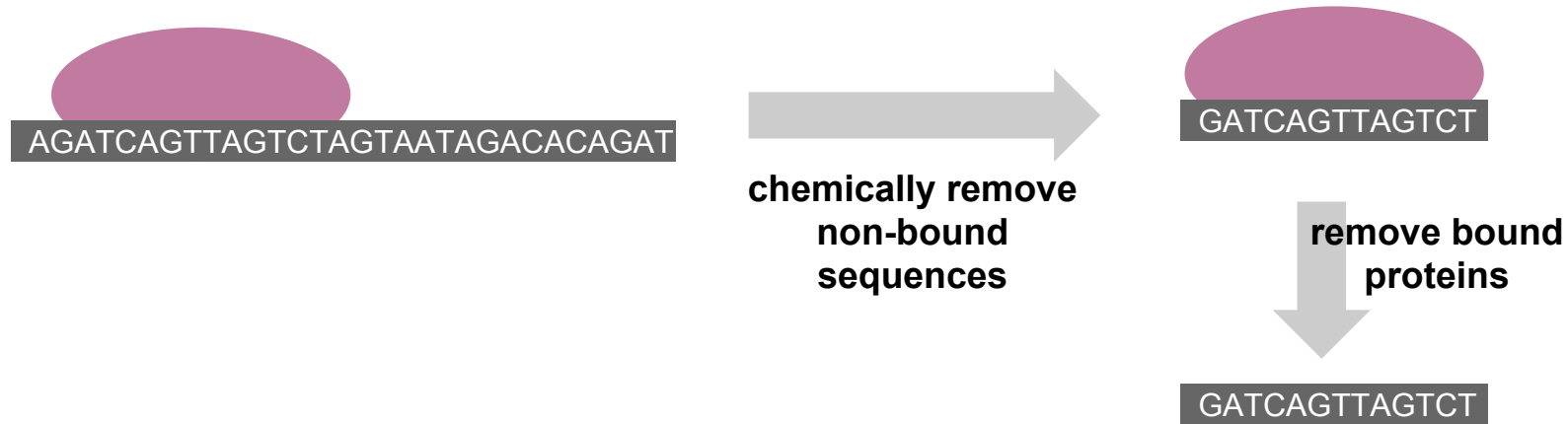
(RNA-seq, ChIP-seq, DNase-seq, RAD-seq, BisChIP-seq….. )

# RNA-seq

Sequencing transcripts can give us quantitative data on gene expression



AGATCAGTTAGTCTAGTAGTGGCGCCCGCTAATATACGCGGC

**Transcription**

# ChIP-seq

Sequencing bound DNA regions tells us where on the genome activity is happening

# Future

- Better sequencing technologies will make many problems non-issues

- Maybe we can co-opt quantify self data as phenotype data

- Open source data

- Questions??