# Quantitative Biology Bootcamp

Michael Schatz, Justin Kinney, Mickey Atwal

Aug 27, 2014
WSBS

# Unsolved Questions in Biology

- What is your genome sequence?
- How does your genome compare to my genome?

- Where are the genes and how active are they?
- How does gene activity change during development?
- How does splicing change during development?

- How does methylation change during development?
- How does chromatin change during development?
- How does is your genome folded in the cell?
- How do proteins bind and regulate genes?

- What virus and microbes are living inside you?
- How do your mutations relate to disease?
- What drugs and treatments should we give you?

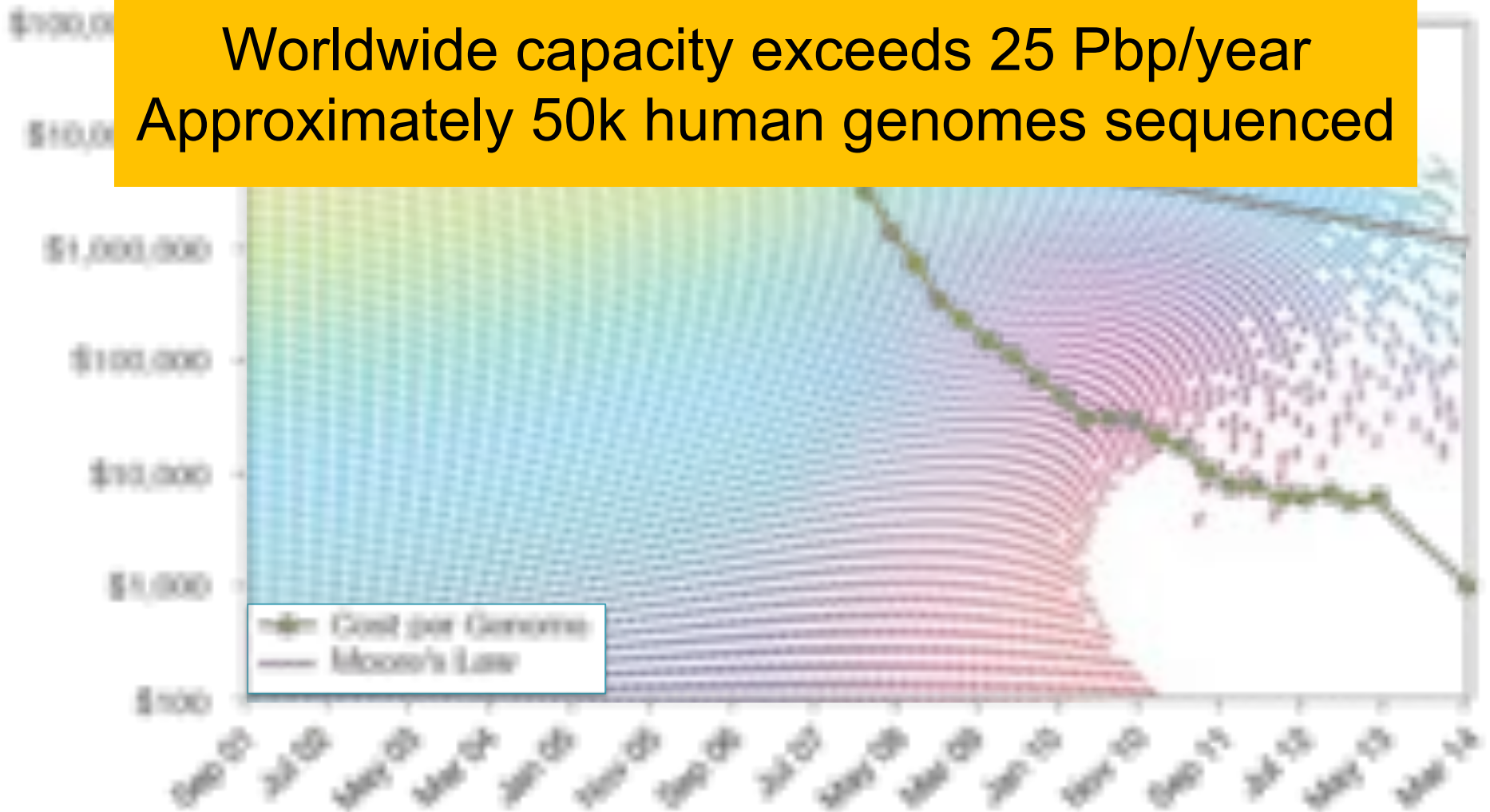- ***Plus thousands and thousands more***

# Data types across the NIH

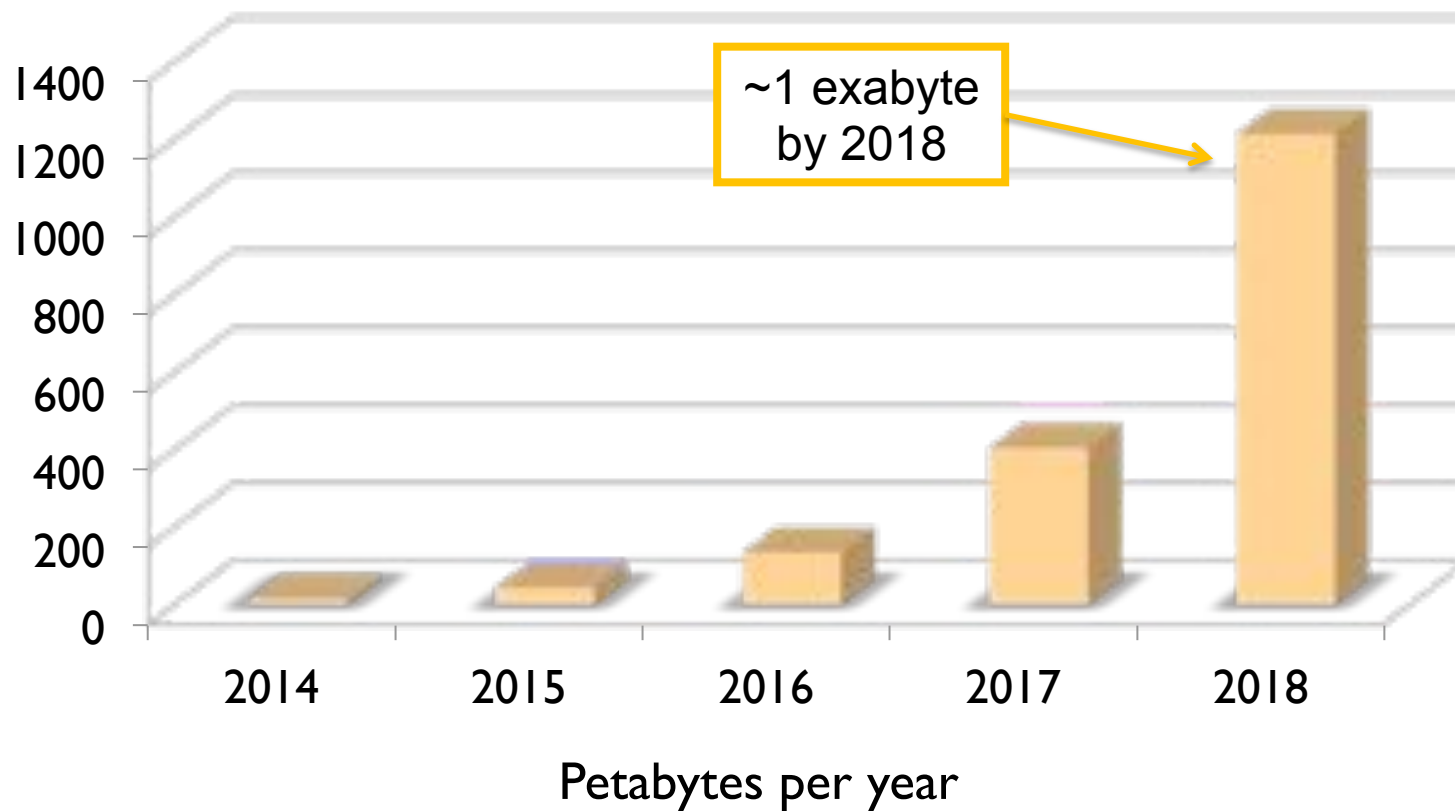

**Phil Bourne, Associate Director of Data Science for NIH**
http://www.slideshare.net/pebourne/wiki-mania080914

# Cost per Genome

Worldwide capacity exceeds 25 Pbp/year
Approximately 50k human genomes sequenced



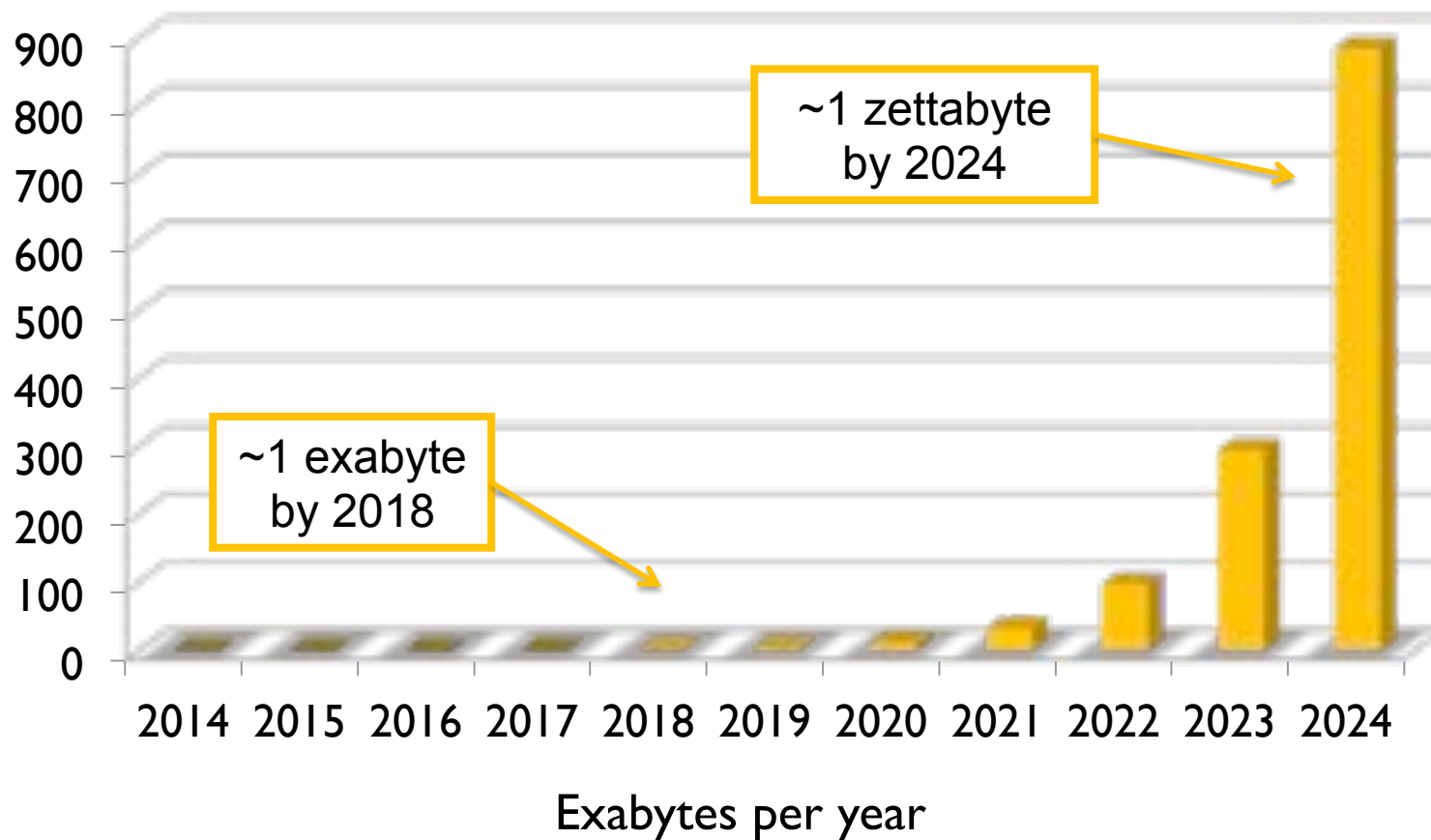http://www.genome.gov/sequencingcosts/

# DNA Data Tsunami

*Current world-wide sequencing capacity is growing at ~3x per year!*

# DNA Data Tsunami

*Current world-wide sequencing capacity is growing at ~3x per year!*



~1 zettabyte
by 2024

~1 exabyte
by 2018

900
800
700
600
500
400
300
200
100
0

2014 2015 2016 2017 2018 2019 2020 2021 2022 2023 2024

Exabytes per year

# How much is a zettabyte?

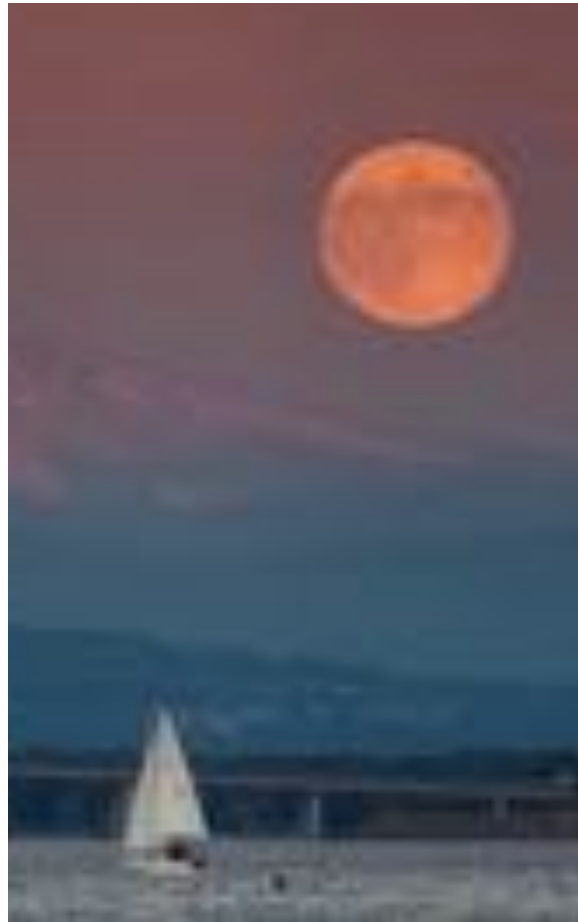| Unit | Size |
|---|---:|
| Byte | 1 |
| Kilobyte | 1,000 |
| Megabyte | 1,000,000 |
| Gigabyte | 1,000,000,000 |
| Terabyte | 1,000,000,000,000 |
| Petabyte | 1,000,000,000,000,000 |
| Exabyte | 1,000,000,000,000,000,000 |
| Zettabyte | 1,000,000,000,000,000,000,000 |

# How much is a zettabyte?



100 GB / Genome
4.7GB / DVD
~20 DVDs / Genome

X

10,000,000,000 Genomes

=

1ZB Data
200,000,000,000 DVDs

150,000 miles of DVDs
~ ½ distance to moon

Both currently ~100Pb
And growing exponentially

# Unsolved Questions in Biology

- What is your genome sequence?
- 
- 
- 
- 
- 
- 
- 
- 
- 
- 
- *Plus thousands and thousands more*

The instruments provide the data, but none of the answers to any of these questions.

**What software and systems will?**

**And who will create them?**

# Who is a Data Scientist?



http://en.wikipedia.org/wiki/Data_science
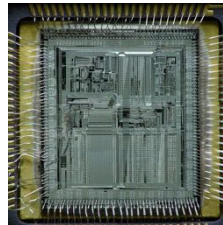
# What is a computer?

[hardware]



**Hard Drive**
Permanent Storage – 1TB
(big, slow, cheap)

**RAM**
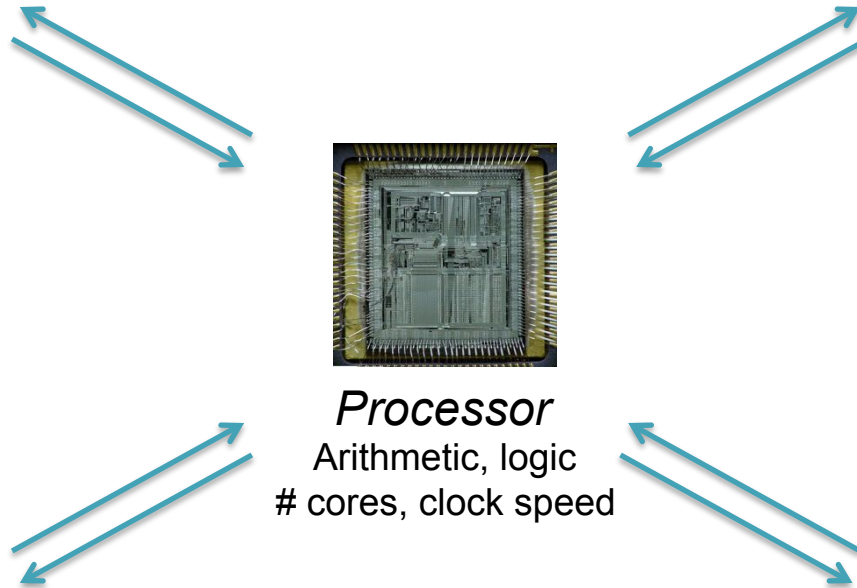Working Storage – 8 GB
(small, fast, expensive)

**Processor**
Arithmetic, logic
# cores, clock speed

**Display**
Human Interface

**Network**
Computer Interface
Home: 10Mb/s, CSHL: 1Gb/s
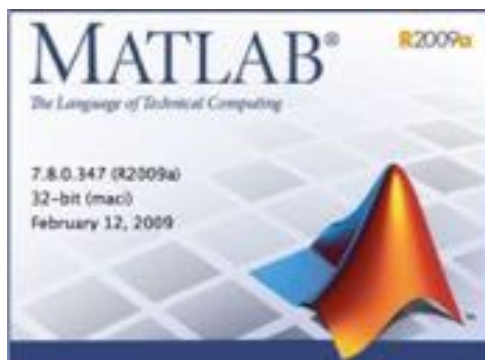
# What is a computer?
## [software]



**Office Applications**
Presentations, Documents
Simple statistics and plots

**Files / Data**
Papers, sequences,
measurements

**Operating System**
Mission Control
Windows, Mac, Unix, iOS

**Scientific Applications**
Specialized Analysis
Commercial

**Code / Scripts**
Research Applications
Academic

# Programming 101



A software program is like sheet music for the orchestra inside your computer
Static, written representations of an active process

# Programming with Python