# What is data?

Jeffrey Leek, Assistant Professor of Biostatistics
Johns Hopkins Bloomberg School of Public Health

# Definition of data

" Data are values of qualitative or quantitative variables, belonging to a set of items. "

http://en.wikipedia.org/wiki/Data

# Definition of data

"Data are values of qualitative or quantitative variables, belonging to a set of items."

http://en.wikipedia.org/wiki/Data

**Set of items**: Sometimes called the population; the set of objects you are interested in

# Definition of data

"Data are values of qualitative or quantitative variables, belonging to a set of items."

http://en.wikipedia.org/wiki/Data

**Variables**: A measurement or characteristic of an item.

# Definition of data

**"** Data are values of <span style="color:red">qualitative</span> or <span style="color:red">quantitative</span> variables, belonging to a set of items. **"**

http://en.wikipedia.org/wiki/Data

**Qualitative**: Country of origin, sex, treatment

**Quantitative**: Height, weight, blood pressure

# Raw versus processed data

**Raw data**

- The original source of the data

- Often hard to use for data analyses

- Data analysis *includes* processing

- Raw data may only need to be processed once

http://en.wikipedia.org/wiki/Raw_data

**Processed data**

- Data that is ready for analysis

- Processing can include merging, subsetting, transforming, etc.

- There may be standards for processing
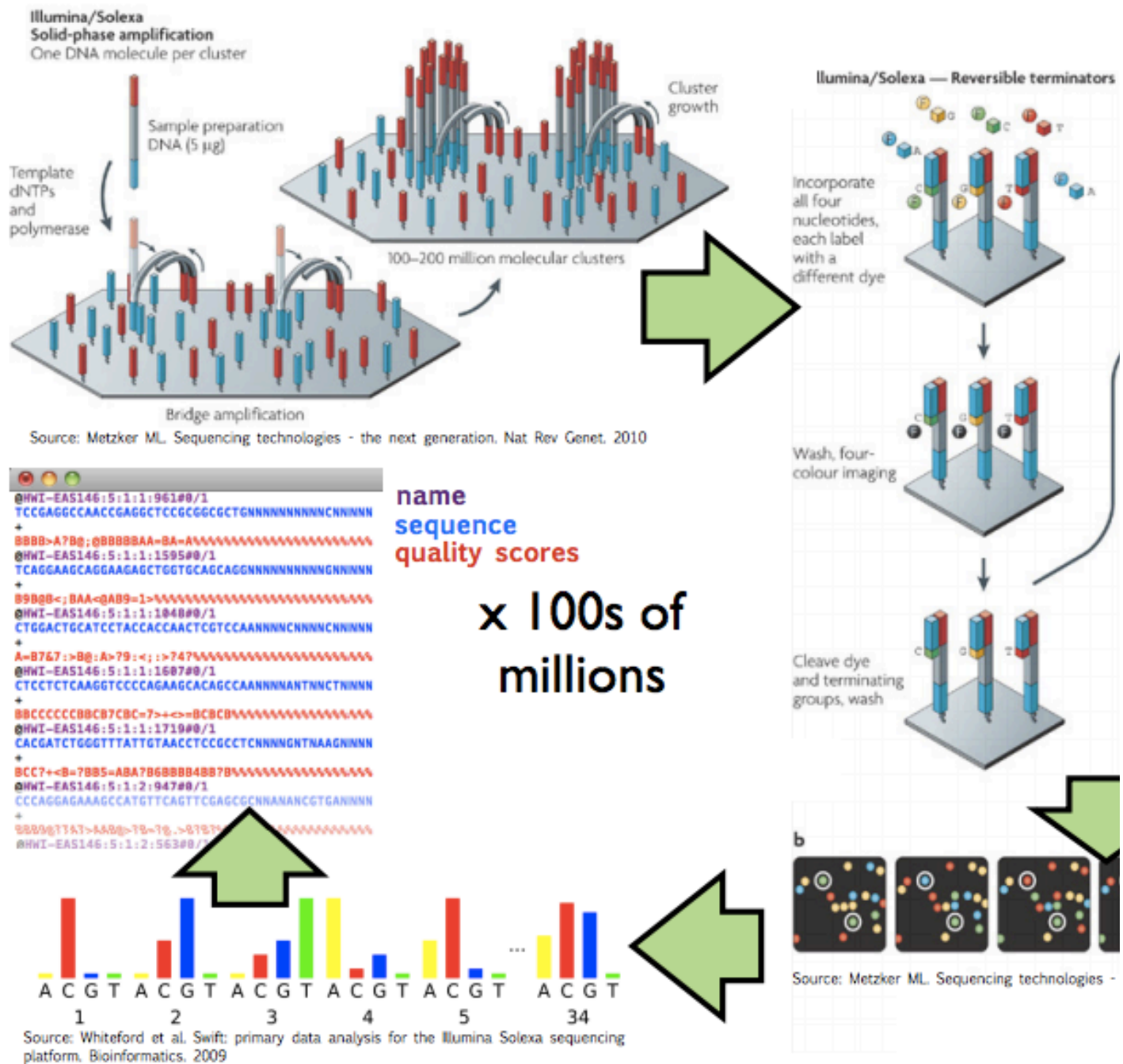
- All steps should be recorded

http://en.wikipedia.org/wiki/Computer_data_processing

# An example of a processing pipeline



http://www.illumina.com.cn/support/sequencing/sequencing_instruments/hiseq_1000.asp

# An example of a processing pipeline



http://www.cbcb.umd.edu/~hcorrada/CMSC858B/lectures/lect22_seqIntro/seqIntro.pdf

# What do raw data look like?

```
@HWI-EAS121:4:100:1783:550#0/1
CGTTACGAGATCGGAAGAGCGGTTCAGCAGGAATGCCGAGACGGATCT(
+HWI-EAS121:4:100:1783:550#0/1
aaaaa`b_aa`aa`YaX]aZ`aZM^Z]YRa]YSG[[ZREQLHESDHND
@HWI-EAS121:4:100:1783:1611#0/1
GGGTGGGCATTTCCACTCGCAGTATGGGTTGCCGCACGACAGGCAGCG(
+HWI-EAS121:4:100:1783:1611#0/1
a```^\__`_````^a``a`^a_^__]a_]\]`a_____`_^^`]X]_]:
@HWI-EAS121:4:100:1783:322#0/1
CGTTTATGTTTTTGAATATGTCTTATCTTAACGGTTATATTTTAGATG'
+HWI-EAS121:4:100:1783:322#0/1
abaa`^aaaaabbbaababbbbb`bbbb_bbbbbbb`bbbaV^_a`
@HWI-EAS121:4:100:1783:1394#0/1
GGGTCTTTATTGGTCTGGTGATCCCCCATATTCTCCGGTTGTGTGGTT
+HWI-EAS121:4:100:1783:1394#0/1
```` [aa\b^^[]aabbb][`a_abbb`a``bbbbbabaabaaaab_VZ;
@HWI-EAS121:4:100:1783:207#0/1
CCCTGGGAGATCGGAAGAGCGGTTCAGCAGGAATGCCGAGACCGATCT(
+HWI-EAS121:4:100:1783:207#0/1
abba`Xa\^\\`aa]ba__bba[a_O_a`aa`aa`a]^V]X_a^YS\R/
@HWI-EAS121:4:100:1783:455#0/1
GGGTAATTCAGGGACAATGTAATGGCTGCACAAAAAAATACATCTTTC;
+HWI-EAS121:4:100:1783:455#0/1
abb_babbabaabbbbbbbbbbbbbbbbbba\`b`\abbbabbbbabbbb[
```

http://brianknaus.com/software/srtoolbox/s_4_1_sequence80.txt

# What do raw data look like?



https://dev.twitter.com/docs/api/1/get/blocks/blocking

# What do raw data look like?



```
---------------------------  ALLERGIES  ---------------------------       -------

ast Updated: 01 Dec 2011 @ 0851                                        Last Up

                                                                       Medicat
llergy Name:          TRIMETHOPRIM                                     Instruc
ocation:              DAYT29                                           GRAPEFR
ate Entered:          09 Mar 2011                                      Status:
eaction:                                                               Refills
llergy Type:          DRUG                                             Last Fi
A Drug Class:         ANTI-INFECTIVES,OTHER                            Initial
bserved/Historical:  HISTORICAL                                        Quantit
omments:              The reaction to this allergy was MILD (NO SQUELAE)  Days S
                                                                       Pharmac
llergy Name:          TRAMADOL                                         Prescri
ocation:              DAYT29
ate Entered:          09 Mar 2011                                      Medicat
eaction:              URINARY RETENTION                                Instruc
llergy Type:          DRUG                                             Status:
A Drug Class:         NON-OPIOID ANALGESICS                            Refills
bserved/Historical:  HISTORICAL                                        Last Fi
omments:              gradually worsening difficulty emptying bladder  Initial
```

http://blue-button.github.com/challenge/

# What do processed data look like?



1. Each variable forms a column

2. Each observation forms a row

3. Each table/file stores data about one kind of observation (e.g. people/hospitals).

http://vita.had.co.nz/papers/tidy-data.pdf

Leek, Taub, and Pineda 2011 PLoS One

12/19

# How much is there?



[http://mashable.com/2011/06/28/data-infographic/](http://mashable.com/2011/06/28/data-infographic/)

# So what about big data?

# Depends on your perspective

# Why big data now?

An Experimental Study of the Small World Problem*

JEFFREY TRAVERS

Harvard University

AND

STANLEY MILGRAM

The City University of New York

Arbitrarily selected individuals $(N=296)$ in Nebraska and B... to generate acquaintance chains to a target person in M... ing "the small world method" (Milgram, 1967). Sixty-fou... the target person. Within this group the mean number... tween starters and targets is 5.2. Boston starting chains re...

[Travers and Milgram (1969) Sociometry](#)

# Why big data now?



arXiv.org > physics > arXiv:0803.0939

**Physics > Physics and Society**

## Planetary–Scale Views on an Ins

Jure Leskovec, Eric Horvitz

*(Submitted on 6 Mar 2008)*

We present a study of anonymized data capturing a mc Microsoft Messenger instant-messaging system. We ex dynamics of large numbers of people, rathe properties of 30 billion onversations amor 240 milli million nodes and 1.3 billion undirected ed on multiple aspects of the dataset and synthesized gra We investigate on a planetary-scale the oft-cited the average path length among Messenger users s 6.6 when they have similar age, language, and locati duration than conversations with the same gender.

[Leskovec and Horvitz WWW '08](#)

# Big or small - you need the right data

" The data may not contain the answer. The combination of some data and an aching desire for an answer does not ensure that a reasonable answer can be extracted from a given body of data... "

Tukey

# Big or small - you need the right data

" ...no matter how big the data are. "

[Leek](#)