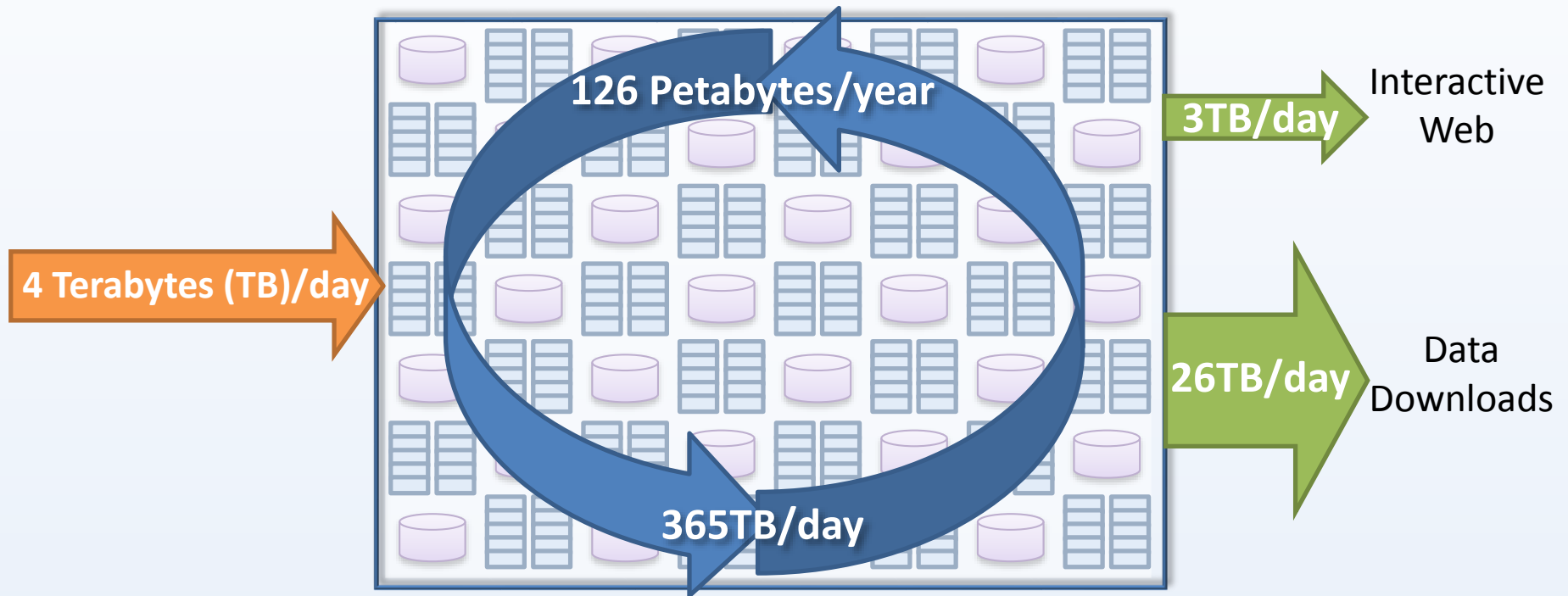# What is Big Data?

- Information Retrieval *not* Big Data

- Computing across TB's of data (Higgs boson) *is* Big Data

- Computing across credit card data to detect possible fraud is Big Data
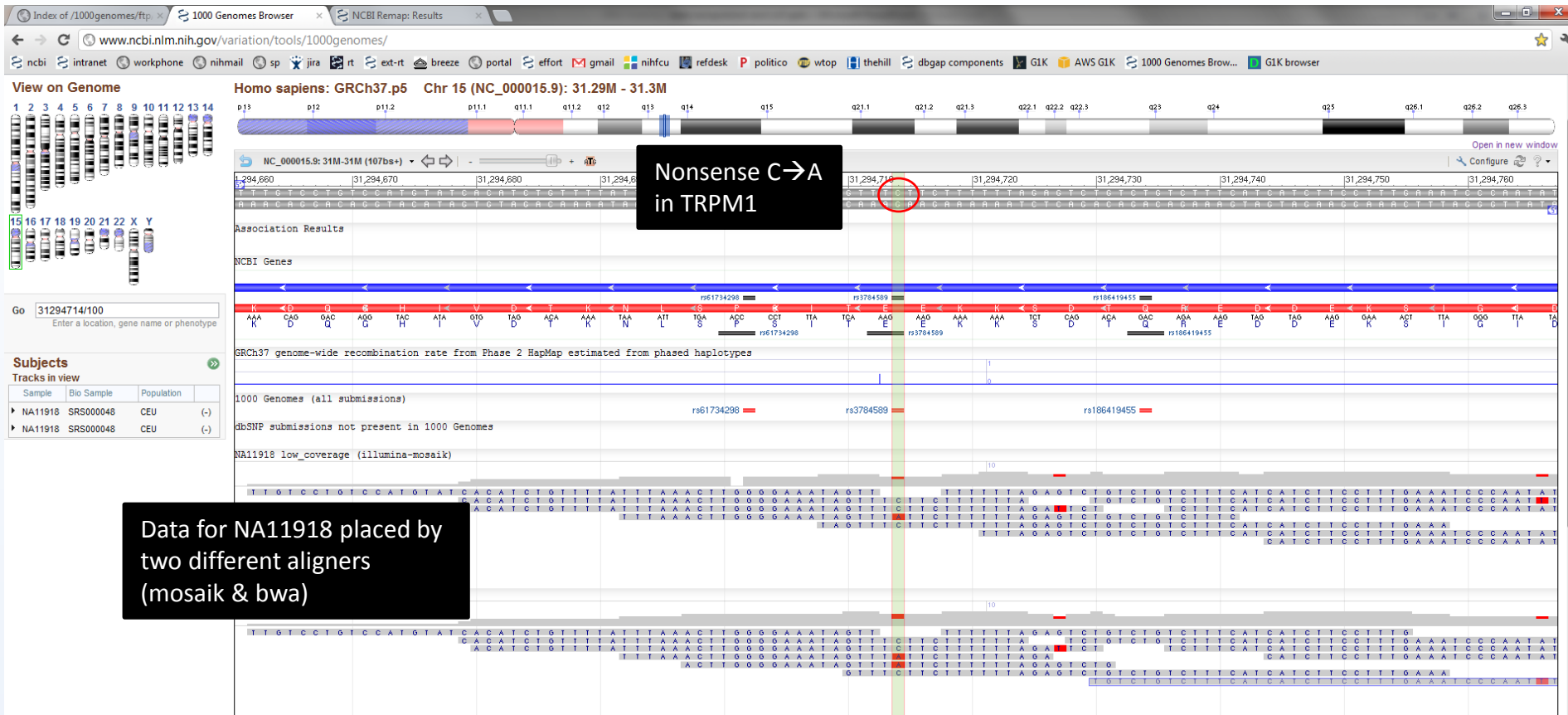
& Diapers == *Predictive Analytics*
http://www.nytimes.com/2012/02/19/magazine/shopping-habits.html?smid=pl-share

# Daily Data Processing at NCBI



**4 Terabytes (TB)/day** → 

**126 Petabytes/year**

**365TB/day**

**3TB/day** → Interactive Web

**26TB/day** → Data Downloads

2012 data…

# Stationary night blindness due to premature termination in TRPM1



Nonsense C→A in TRPM1

Data for NA11918 placed by two different aligners (mosaik & bwa)

All individual genotypes For rs3784589

# What is the Big Data Problem in Biology?
# Example: Reducing the 1000 Genome Dataset



**Submitted BAM**
Read IDs as strings
Original quality & recalibrated quality scores
Additional analysis tags

**250TB**

**cSRA (lossless)**
Read IDs as integers
40-level read qualities using recalibrated quality scores

**85TB**

**cSRA (lossy)**
8 level qualities for all sites
Uniform binning of recalibrated quality scores

**30TB**

**Variant Call Format (VCF)**
Genotype likelihoods for all variants **0.1TB**

Size (Terabytes)

300
250
200
150
100
50
0

Total Project Size     Lossless cSRA     Lossy cSRA     Analysis Genotypes

• Computing on reads is a short term technical challenge

• Computing on growing amounts of derived data is a long term challenge & opportunity

# NextGen Churn



**Basic Data Processing**

Mapping reads, calling SNPs, splices, peaks, etc.

**Derived Data**

**Genotypes, Genes, expression levels, motifs, etc.**

Etc….

*Biological Questions*

# As cost of computers decreased, overall investments increased...

# Economics of NextGen

*Grantees want to maximize impact of budget*

*Sequencing will continue to increase if* **total cost** *continues to drop:*

*sample prep + raw sequencing + IT*

# Big Data:

(In

**Ro**

Ema

htt

**What mess?**

When

# BIG DATA meets Small Signal

the signal has to be **teased out** of the data

- There are many ways to do this: a **small change** in the analysis details can cause a **large change** in the results.

- It is too easy to distort your findings, either by **fooling yourself** or **on purpose**.

# Systematic errors are routinely observed and have been reported in many papers

- Identification and correction of systematic error in high-throughput sequence data by
  Frazer Meacham, Dario Boffelli, Joseph Dhahbi, David IK Martin, Meromit Singer and Lior Pachter
      BMC Bioinformatics. 2011 Nov 21;12:451

  This paper shows the existence of systematic errors, even at high coverage, often strand dependent but not always

# Mismatch profiles are dominated by 'spikes' occurring at particular sequencing cycles (sequencing batch effect). Furthermore, each problematic cycle has a specific limited profile of mismatch types, adding to the bias

**mismatches per base position NB1 to NB70**



## This effect on Illumina HiSeq is sequencing lane dependent



100 + 100 bp paired end reads, Illumina, Read 1 ------→          ---------------Read 2 ----------------------→

*SEQC neuroblastoma study, L Shi, Fischer et al*

# Systematic errors generate noise in the low to intermediate allele fraction (1 to 30%), making identification of true SNPs hard in that area, even at high coverage

**Histogram of candidate variant allele fraction before filtering: massive presence of SNP candidates with low to intermediate variant allele fraction values**

RNA-Seq often differs from 'true' concentration by a factor 2 or more.
This affects all platforms, in particular Illumina and PGM/Proton/Solid



**Measured expression index (log2-based) and nominal concentration for ERCC control mRNA molecules**

The two red and green runs should be superimposed on the blue nominal concentration.
Yet some specific ERCC molecules are vastly different, e.g. ERCC116 is measured 16 fold below nominal.

*Data by Setterquist, LifeTech, QC for PGM*

HEALTH INDUSTRY | Updated May 4, 2012, 10:24

## Analytical Trend Trou

### Studying Up

While increasingly popular, observational studies often yield spurious results.

| 1990 – 2000 | 2001 – 2011 |
|---|---|

40 thousand studies

John Ioannides,
Stanford

Total: 79,619

Total: 263,557

**Number of studies coded as corrections in 1990 – 2000**

**185**
(.23%)

**Number of studies coded as corrections in 2001 – 2011**

**881**
(.33%)

**Top five fields that use observational studies**

Public, environmental, and occupational health — 10,585 studies

Medicine, general and internal — 8,958

Oncology — 5,356

Psychology — 5,315

Surgery — 4,785

**Top five fields that use observational studies**

Public, environmental, and occupational health — 27,240 studies

Medicine, general and internal — 23,545

Neurosciences and neurology — 19,613

Oncology — 17,522

Cardiovascular and cardiology — 16,141

*Abstracts of indexed studies were not started until 1991, therefore 1990 is an undercount

Source: Thomson Reuters Web of Science, an index of peer-reviewed journals

The Wall Street Journal

"That partly explains why observational studies in general can be replicated only 20% of the time, versus 80% for large, well-designed randomly controlled trials, says Dr. Ioannidis. Dr. Young, meanwhile, pegs the replication rate for observational data at an even lower 5% to 10%."

# Correlation between impact factor and retraction index.



Fang F C , Casadevall A Infect. Immun. 2011;79:3855-3859

Infection and Immunity

# Repeatability of published microarray gene expression analyses
Ioannidis et al. *Nature Genetics* **41**, 149 - 155 (2009)

# "Preclinical research generates many secondary publications, even when results cannot be reproduced"

◄ back to article

**Table 1: Reproducibility of research findings**
Preclinical research generates many secondary publications, even when results cannot be reproduced.

| Journal impact factor | Number of articles | Mean number of citations of non-reproduced articles[*] | Mean number of citations of reproduced articles |
|---|---|---|---|
| >20 | 21 | 248 (range 3–800) | 231 (range 82–519) |
| 5–19 | 32 | 169 (range 6–1,909) | 13 (range 3–24) |

Results from ten-year retrospective analysis of experiments performed prospectively. The term 'non-reproduced' was assigned on the basis of findings not being sufficiently robust to drive a drug-development programme.

[*]Source of citations: Google Scholar, May 2011.

**Tables index**

- Biomarkers – most highly cited studies overestimated effect sizes (*JAMA* 2011; 305(21): 2200-2210)

- Faculty & Trainee survey at MD Anderson – 50% had experienced at least one case of irreproducibility (PLOS One 2013 May 15; 8(5))

- More first-in-class small molecule drugs approved between 1999-2008 identified by "classical" methods than genomics approaches (J. Biomol Screen 2013 Dec; 18(10): 1143-55)

- "Koch's Postulates" for assigning causality between genetic variants & disease phenotypes (Cell 2013 Sep 26; 155(1) : 21-6)

# Lab Mistakes Hobble Cancer Studies
# But Scientists Slow to Take Remedies

"Cancer experts seeking to solve the problem have found that a <u>fifth to a third or more</u> of cancer cell lines tested were mistakenly identified—with researchers unwittingly studying the wrong cancers, slowing progress toward new treatments and wasting precious time and money."

"…Dr. Masters, in a study of scientific papers published between 2000 and 2004, found nearly a 1,000 citations of the same contaminated cancer lines revealed in Dr. Gartler's 1966 findings, which have since been replicated many times using more advanced techniques."

Citations of T24 bladder cancer cells referred to as normal endothelial cells.
*Nature Reviews Cancer* **10**, (June 2010) | doi:10.1038/nrc2852



Nature Reviews | Cancer

# Dr. John Snow - location of pumps and cholera deaths, London, England, 1854

**LETTUCE**
Canada, Chile, Dominican Republic, Mexico, Peru, USA

**CROUTONS**
Argentina, Australia, Brazil, Canada, China, France, India, Mexico, Netherlands, Poland, Russia, Switzerland, Uruguay, USA, Vietnam

## The Well-Traveled Salad.
### Do You Know Where Your Food Has Been?

As consumers, many of us fail to recognize that even our domestic and local food supplies are part of a global network. The daily activity of consuming food directly links our health as humans to the health of crops and produce, food animals, and the environments in which they are produced.

**CUCUMBERS**
Canada, Honduras, India, Mexico, Spain, USA

**TOMATOES**
Canada, Dominican Republic, Holland, Israel, Italy, Mexico, USA

**FETA CHEESE**
Canada, Denmark, Egypt, Germany, Greece, Israel, Italy, Turkey, UK, USA

**ONIONS**
Canada, China, Germany, India, USA

**VINAIGRETTE**
Argentina, Brazil, Canada, Chile, China, France, Germany, Greece, India, Indonesia, Italy, Mexico, Morocco, Peru, Portugal, Spain, Thailand, Tunisia, Turkey, USA, Vietnam

**OLIVES**
Greece, Israel, Mexico, Spain, USA

**SPROUTS**
Argentina, Australia, Bangladesh, Canada, China, Egypt, France, India, Morocco, Nepal, Pakistan, South Africa, Spain, Turkey, USA

**MANDARIN ORANGES**
Israel, Mexico, Morocco, South Africa, Spain

A "One Health" approach to food safety—bringing together expertise and resources from the clinical, veterinary, wildlife health, and ecology communities—has the potential to reveal the sources, pathways, and factors driving the outbreaks of foodborne illness and possibly prevent them from occurring in the first place.

**NOTE**: Countries are listed in alphabetical order and not by volume of export.

**INSTITUTE OF MEDICINE**
OF THE NATIONAL ACADEMIES

www.iom.edu

# Foodborne Illness

- ~48 million US cases annually (CDC)
  - 128,000 hospitalized
  - 3,000 deaths
- Trends show little evidence of progress

# DNA Forensics…



Functional prediction can be developed and refined more slowly from this base.

*A Pathogen Genome Is The Fingerprint*

Search for phylogenetic signal at the level of SNPs

# Why do we need WGS? To Shift the Paradigm from a "low resolution" Public Health Approach to A Real-Time "high resolution" Approach



**Clinical ID and fingerprint**

CDC

**Identify Food and confirm PFGE Fingerprint**

FDA/FSIS

**Product enters commerce**

Source of contaminati
identified too late

Number of cases

Days

**Current Genometrakr network
7 state Laboratories + 11 FDA-ORA**

Network of Sequencers

# Next-Generation Lab Response vs. Conventional Lab Response

Isolate arrives at lab

Library Prep
Whole Genome
Sequencing

Illumina
Miseq

Upload to NCBI
Alignment and Clade diagram produced

AGCTAGA
CTACGAA

High Resolution tracking
information avaiable
to public health

**Outbreak** → 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 **Day 15**

Isolate arrives at lab

Serotyping

PFGE

Pattern    PFGE-Xbal
022
041
042
043
044

Bionumeric Analysis

**Conventional Lab Response**

Low Resolution tracking information avaiable to public

FDA

## WGS In its First Official Regulatory Action



**First use of Genometrakr network and WGS approach to support regulatory action and positive public health outcome in real-time** *Listeria monocytogenes project with CDC, FDA, NCBI, USDA March 2014*

**> 6 SNPs** →

Isolates from cheese facility, distributed product, and patients who consumed product

OUTGROUP_PNUSAL000140
CFSAN008989_Clinical_CA
CFSAN009740 Environmental (spanish style cheese) NY
CFSAN010093 Environmental (swab) DE
CFSAN010098 Environmental (swab) DE
CFSAN010758_Fresh_Cheese_Curd_VA
CFSAN010088 Environmental (swab) DE
CFSAN010072_Cheese_MD
CFSAN009222_Clinical_MD
CFSAN010095 Environmental (swab) DE
CFSAN009226_Clinical_MD
CFSAN010075_Cheese_MD
CFSAN010097 Environmental (swab) DE
CFSAN010757_Fresh_Cheese_Curd_VA
CFSAN009229_Clinical_MD
CFSAN010972_Cheese
CFSAN010761_Fresh_Cheese_Curd_VA
CFSAN010762_Fresh_Cheese_Curd_VA
CFSAN010084_Fresh_Cheese_Curd_VA
CFSAN010078_Fresh_Cheese_Curd_VA
CFSAN010763_Fresh_Cheese_Curd_VA
CFSAN010756_Fresh_Cheese_Curd_VA
CFSAN010076_Cheese_MD
CFSAN010074_Cheese_MD
CFSAN010077_Cheese_MD
CFSAN010073_Cheese_MD
CFSAN010094 Environmental (swab) DE
CFSAN010089 Environmental (swab) DE
CFSAN010082_Fresh_Cheese_Curd_VA
CFSAN010759_Fresh_Cheese_Curd_VA
CFSAN010083_Fresh_Cheese_Curd_VA
CFSAN010079_Fresh_Cheese_Curd_VA
CFSAN010755_Fresh_Cheese_Curd_VA
CFSAN010090 Environmental (swab) DE
CFSAN010068_Cheese_MD
CFSAN010091 Environmental (swab) DE
CFSAN010973_Cheese
CFSAN010085_Fresh_Cheese_Curd_VA
CFSAN010096 Environmental (swab) DE
CFSAN010067__Fresh_Cheese_Curd_VA
CFSAN010081_Fresh_Cheese_Curd_VA
CFSAN010087_Fresh_Cheese_Curd_VA
CFSAN010760_Fresh_Cheese_Curd_VA
CFSAN010754_Fresh_Cheese_Curd_VA
CFSAN010092 Environmental (swab) DE
CFSAN010080_Fresh_Cheese_Curd_VA
CFSAN010069_Cheese_MD
CFSAN010070_Cheese_MD
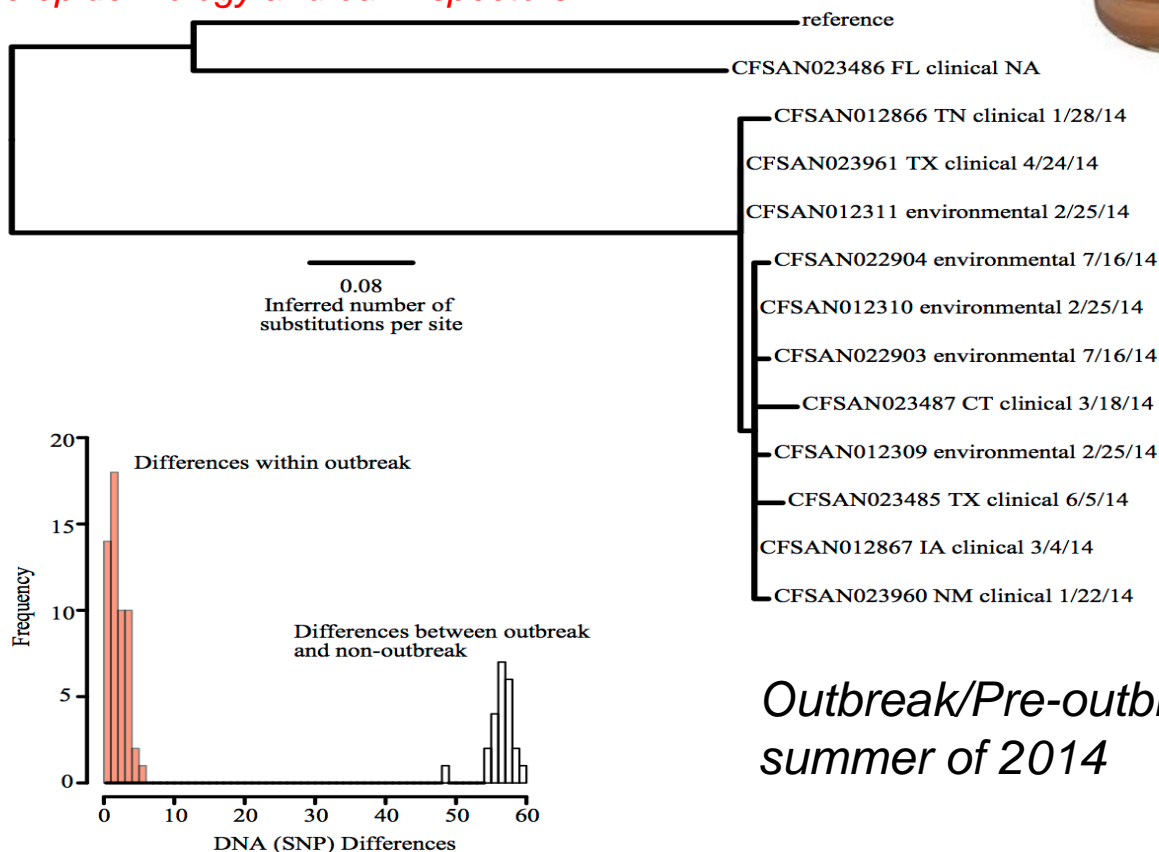CFSAN010086_Fresh_Cheese_Curd_VA
CFSAN010071_Cheese_MD

0.2

# And it gets even better….

*An increased degree of certainty that comes with matching strains of pathogens through whole genome sequencing allowed for detection of this Salmonella contamination event in nut butter across several states with low level contamination and a widely distributed product. In this case, WGS identifies the link and preempts an outbreak even w/o availability of food - it informs the epidemiology and our inspectors.*

Photo courtesy of the MaraNatha website

reference
CFSAN023486 FL clinical NA
CFSAN012866 TN clinical 1/28/14
CFSAN023961 TX clinical 4/24/14
CFSAN012311 environmental 2/25/14
CFSAN022904 environmental 7/16/14
CFSAN012310 environmental 2/25/14
CFSAN022903 environmental 7/16/14
CFSAN023487 CT clinical 3/18/14
CFSAN012309 environmental 2/25/14
CFSAN023485 TX clinical 6/5/14
CFSAN012867 IA clinical 3/4/14
CFSAN023960 NM clinical 1/22/14

0.08
Inferred number of
substitutions per site

RECALL

Differences within outbreak

Differences between outbreak
and non-outbreak

Frequency

DNA (SNP) Differences

*Outbreak/Pre-outbraek summer of 2014*

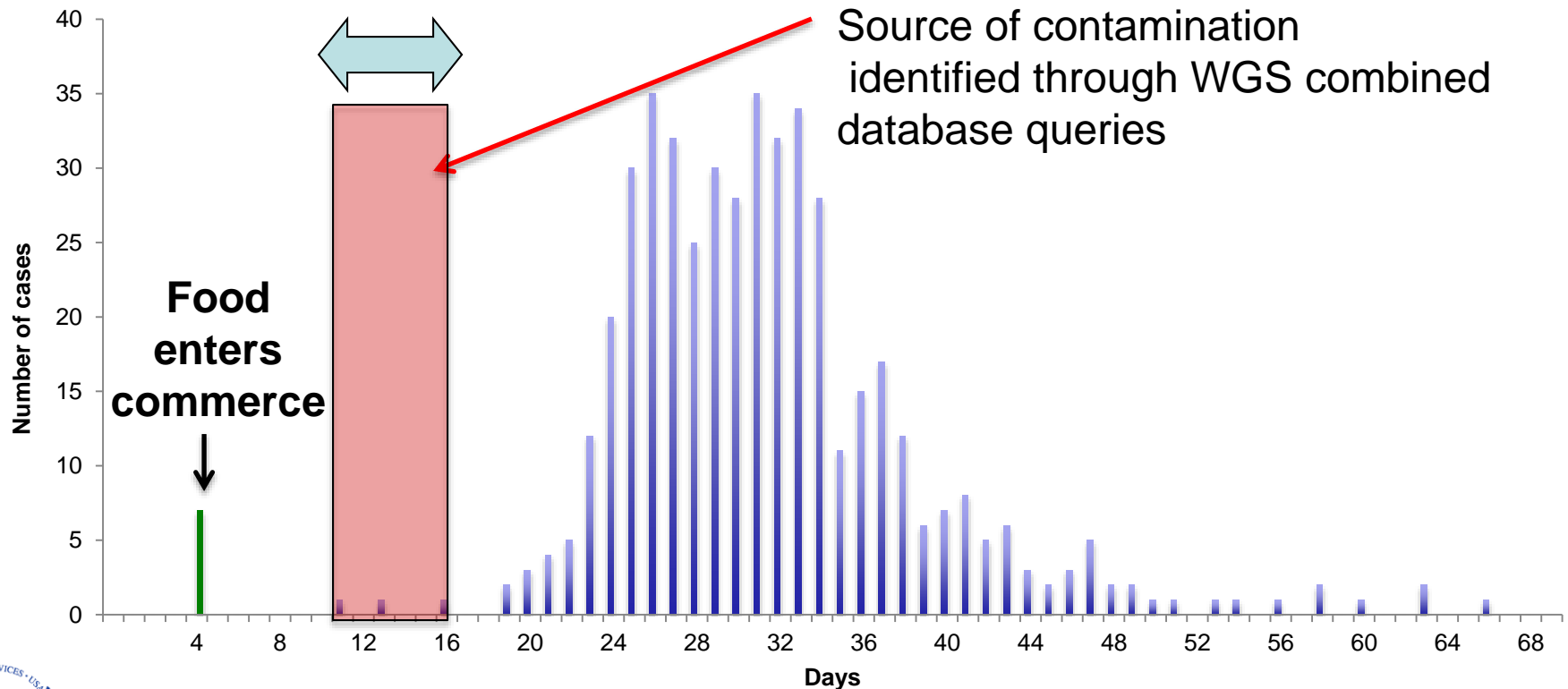# Health and Economic Impact of Active WGS-based Surveillance

**Comparison of Nspired and Sunland contamination events.**

- ❑ **Similar facilities – broad domestic distribution.**

- ❑ **Sunland 42 cases and 10 hospitalizations with as many as 1,260 illnesses unreported (Fall 2012)**

- ❑ **Nspired – 4 confirmed cases, 1 hospitalization (Summer 2014)**

- ❑ **WGS informed investigation prevented significant illness and hospitalizations**

    **– lower illness rate and treatment cost ($3000-$9000) + fraction of longterm and chronic onset complications associated with Salmonella infection (ie, Reiter's syndrome, GBS)**

# The New Microbiology Approach to Public Health

**Clinical ID WGS in real-time and in parallel
food and environmental WGS
FDA, CDC, FSIS, States**



Source of contamination
identified through WGS combined
database queries

# National Digital Immune system: Big Data?

- WGS + metadata

- PulseNet - ~50,000 isolates per year + Hospital based infections (AMR strains)  50-500K per year (??) +  Environmental samples

- Ecology & population genetics of pathogens...