# Managing Big Data Reaching Back to the 11ᵗʰ Century

*Scott Sorensen*

ancestry.com

**ancestry**.com™

Our mission is to help everyone
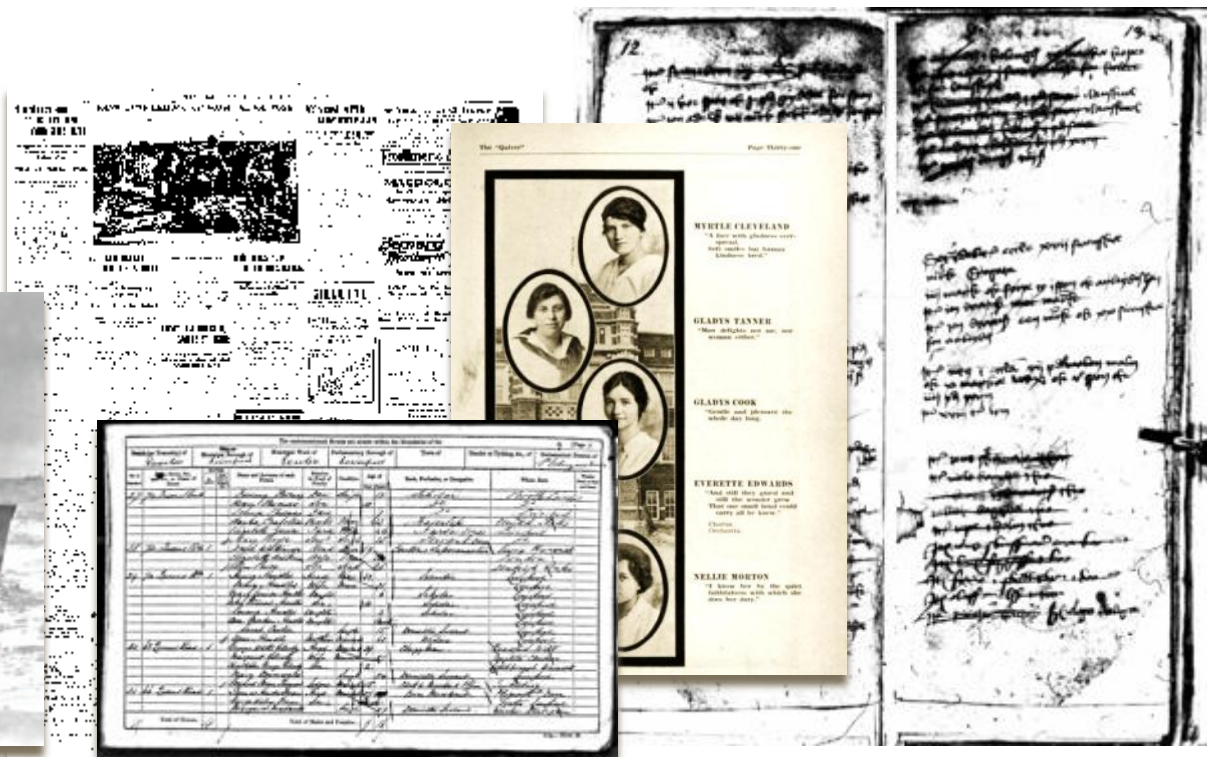discover, preserve and share
their family history.

# Data is our product

## It's the "aha" moment of a discovery that drives our business!

ancestry.com

# World's largest online family history resource

- Over 30,000 historical content collections

- Records dating back to 11$^{th}$ century

- 12 billion records and images

- 10 petabytes

# User contributed content and structure

- 50 million family trees

- More than 5 billion profiles

- 200 million stories and photos

ancestry.com

# Behavioral data



## Next Best Discovery Algorithm

- 40 million searches/day

- 10 million people added to trees/day



Alamy

ancestry.com

# The math behind our big data equation

2,020 nodes
572 marriage edges
2,910 family edges

# We've barely scratched the surface

- Making the site more social through sharing

- Mobile extends the core users experience and attracts a new demographic

- New experiences like AncestryDNA

ancestry.com

# Our transition to Hadoop



## Machine Learning for Information Retrieval

MapReduce

| Random forest (R) | Random forest (R) | Random forest (R) | Random forest (R) |

Model

ancestry.com

# How we're using Hadoop



1. Machine Learning

2. Predictive analytics

3. Natural Language Processing and Entity Extraction

4. DNA Processing

ancestry.com

## Spit in a tube, pay $99, learn your past

Autosomal DNA tests

Samples from over 200,000 people

700,000 SNPs for each sample

10,000,000 4th cousin matches

**Discover Your Ethnicity**

Find out if you're part Irish, Native American, or maybe Cameroonian.

NEWLY UPDATED

**Connect with new relatives**

Imagine meeting a 3rd cousin for the 1st time.

**Family history is in our DNA**

Even more powerful when combined with Ancestry.com.

ancestry.com

# Estimating IBD (matching)

- We identify "long" DNA segments shared by two individuals.

- These segments are said to be Identical-by-Descent (IBD) and identify recent shared genetic ancestry.



Sample #1

Sample #2

Sites essentially either match or they don't

IBD estimation is based on long sequences of consecutive matches

We have a statistical model, based on real pedigree data, that maps match lengths to relationship

# Network effect & cousin matches

# Algorithms in the pipeline

GERMLINE



**Hadoop Cluster** (20 x 4 slots x 96g)



Server  Server  Server  Server  Server  Server  Server  Server  Server  Server

**1) Map Reduce**

| Admixture | Admixture | Admixture |
| Admixture | Admixture | Admixture |
| Admixture | Admixture | Admixture |

ancestry.com

# GERMLINE run times (in hours)

ancestry.com

# Projected GERMLINE run times (in hours)

# The Input



Starbuck : ACTGACCTAGTTGAC
Adama    : TTAAGCCTAGTTGAC



**Kara Thrace, aka Starbuck**

- Ace viper pilot
- Has a special destiny
- Not to be trifled with

**Admiral Adama**

- Admiral of the Colonial Fleet
- Routinely saves humanity from destruction

# Separate into words



                               0      1      2

Starbuck : ACTGA CCTAG TTGAC

Adama    : TTAAG CCTAG TTGAC

ancestry.com

# Build the hash table



|  | 0 | 1 | 2 |
|---|---|---|---|
| Starbuck : | ACTGA | CCTAG | TTGAC |
| Adama : | TTAAG | CCTAG | TTGAC |

ACTGA_0 : Starbuck
TTAAG_0 : Adama
CCTAG_1 : Starbuck, Adama
TTGAC_2 : Starbuck, Adama

ancestry.com

## Iterate through genome and find matches



|       | 0     | 1     | 2     |
|-------|-------|-------|-------|
| Starbuck : ACTGA | CCTAG | TTGAC |
| Adama    : TTAAG | CCTAG | TTGAC |

ACTGA_0 : Starbuck

TTAAG_0 : Adama

CCTAG_1 : Starbuck, Adama

TTGAC_2 : Starbuck, Adama

**Starbuck and Adama match from position 1 to position 2**

# Does that mean they're related?



## ...maybe

ancestry.com

# But wait… what about Baltar?

Baltar : TTAAGCCTAGGGGCG



**Gaius Baltar**

- Handsome
- Genius
- Kinda evil

# The jermline way

## Step one : Update the hash table.

| | Starbuck | Adama |
|---|---|---|
| 2_ACTGA_0 | 1 | |
| 2_TTAAG_0 | | 1 |
| 2_CCTAG_1 | 1 | 1 |
| 2_TTGAC_2 | 1 | 1 |

**← Already stored in HBase**

Baltar : TTAAG CCTAG GGGCG

**← New sample to add**

Key : [CHROMOSOME]_[WORD]_[POSITION]
Cell value : A byte set to 1, denoting that the user has that word at that position on that chromosome

# The *germline* way

## Step two : Find matches, update the results table

|  | 2_Starbuck | 2_Adama |
|---|---|---|
| 2_Starbuck |  | { (1, 2), ...} |
| 2_Adama | { (1, 2), ...} |  |

Already stored in HBase

Baltar and Adama match from position 0 to position 1
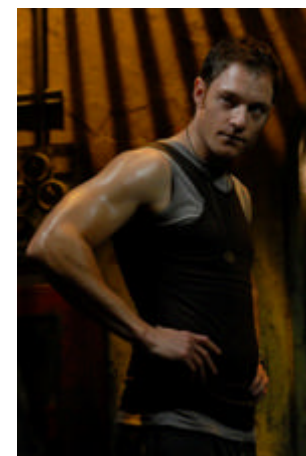Baltar and Starbuck match at position 1

New matches to add

Key : [CHROMOSOME]_[USER ID]
Cell value : A list of ranges where the two users match on a chromosome

ancestry.com

# The jermline way

| Hash Table | | | |
|---|---|---|---|
| | Starbuck | Adama | Baltar |
| 2_ACTGA_0 | 1 | | |
| 2_TTAAG_0 | | 1 | 1 |
| 2_CCTAG_1 | 1 | 1 | 1 |
| 2_TTGAC_2 | 1 | 1 | |
| 2_GGGCG_2 | | | 1 |

| Results Table | | | |
|---|---|---|---|
| | 2_Starbuck | 2_Adama | 2_Baltar |
| 2_Starbuck | | { (1, 2), …} | { (1), …} |
| 2_Adama | { (1, 2), …} | | { (0,1), …} |
| 2_Baltar | { (1), …} | { (0,1), …} | |

# But wait … what about Zarek, Roslin, Hera, and Helo?

ancestry.com
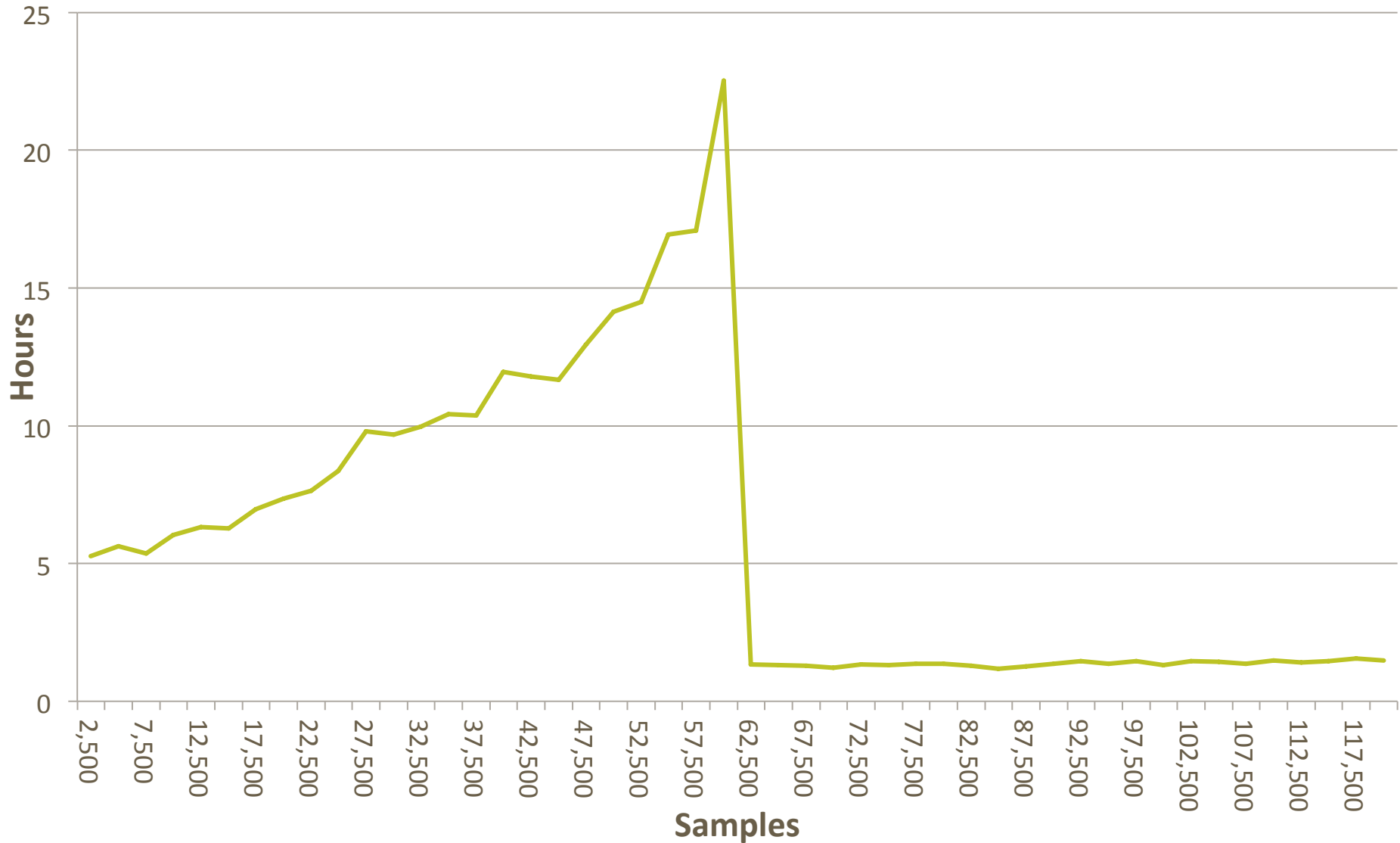
# Run them in parallel with Hadoop!
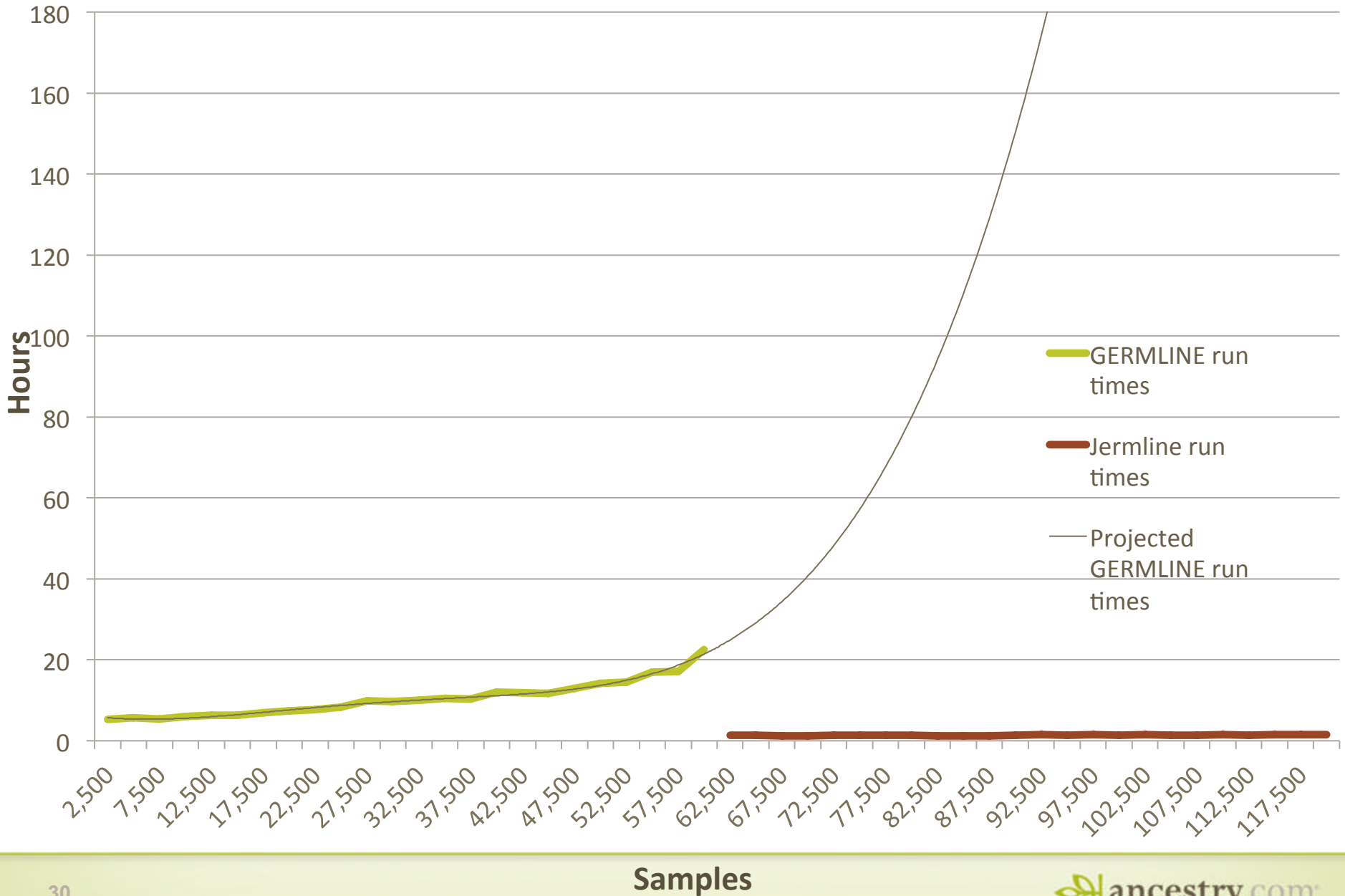


Photo by Benh Lieu Song

# Parallelism with Hadoop

- Batches are usually about a thousand people.

- Each mapper takes a single chromosome for a single person.

- MapReduce Jobs :

    Job #1 : Match Words

    - o Updates the hash table

    Job #2 : Match Segments

    - o Identifies areas where the samples match

# Run times for matching (in hours)

# Run times for matching (in hours)



**Hours**

**Samples**

Legend:
- GERMLINE run times
- Jermline run times
- Projected GERMLINE run times

ancestry.com

# AncestryDNA – Cast of characters

## Scientists

Think they can code:

- Linux

- MySQL

- PERL and/or Python

## Software Engineers

Think they are Scientists:

- Math

- Statistics

- Read science papers

# Pressures of a startup business

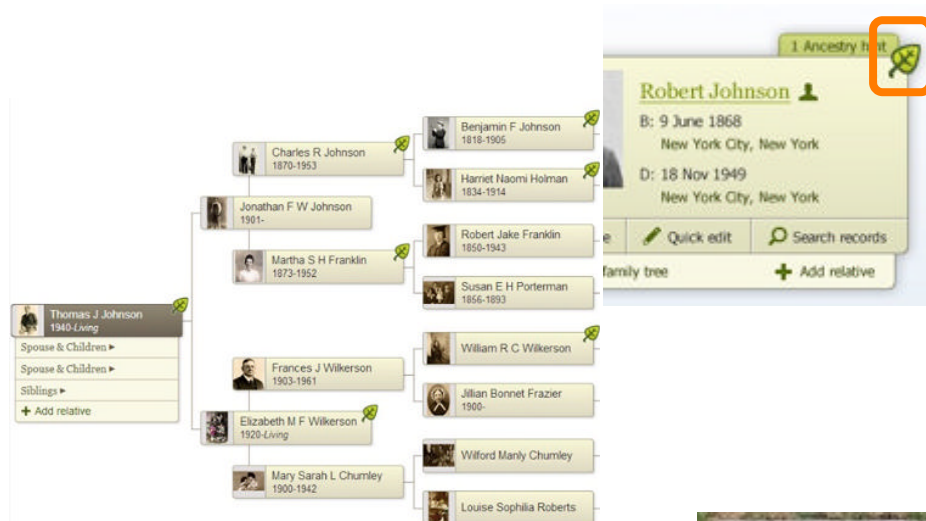– Release a product, learn, and then scale

# Other lessons learned

- Prototyping is key to overcoming resistance to change

- Technical architecture is heavily influenced by people organization

- Developing a team of experienced Hadoop users can often be done using internal employees

- A culture of experimentation and innovation yields the best results

ancestry.com

# Using Hadoop to drive scalable results

- Machine learning and predictive analytics

- Entity extraction and product development

- DNA pipeline processing

# Questions?

Tech Roots blog - http://blogs.ancestry.com/techroots/