

Can Python Save Next Generation Sequencing?

Chris Mueller
Life Technologies
June 30, 2010
SciPy 2010 Austin, TX





Human Genome Project



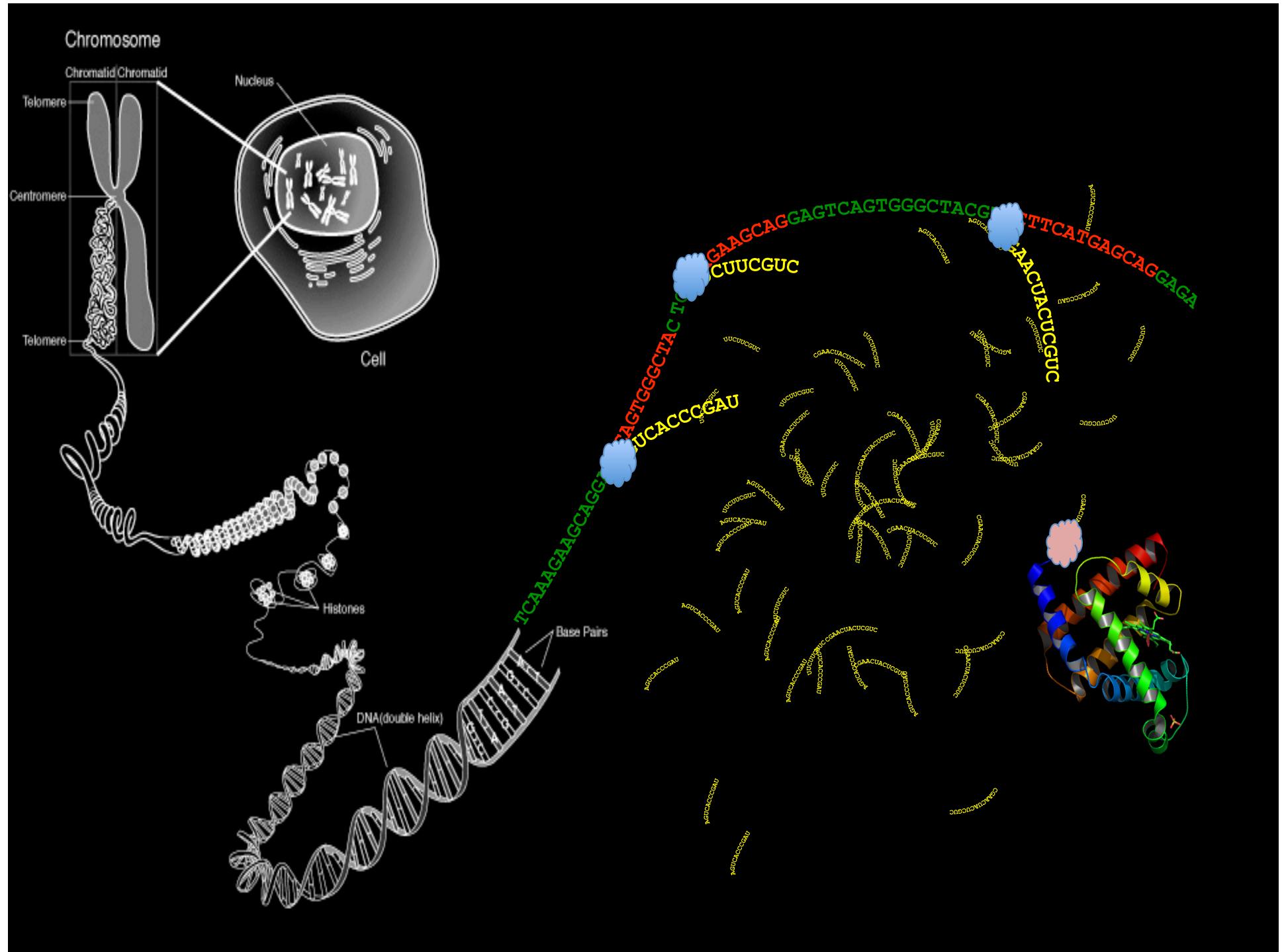
10 Years
Thousands of Sequencers
\$3,000,000,000
21 Billion Base Pairs (Gbp)

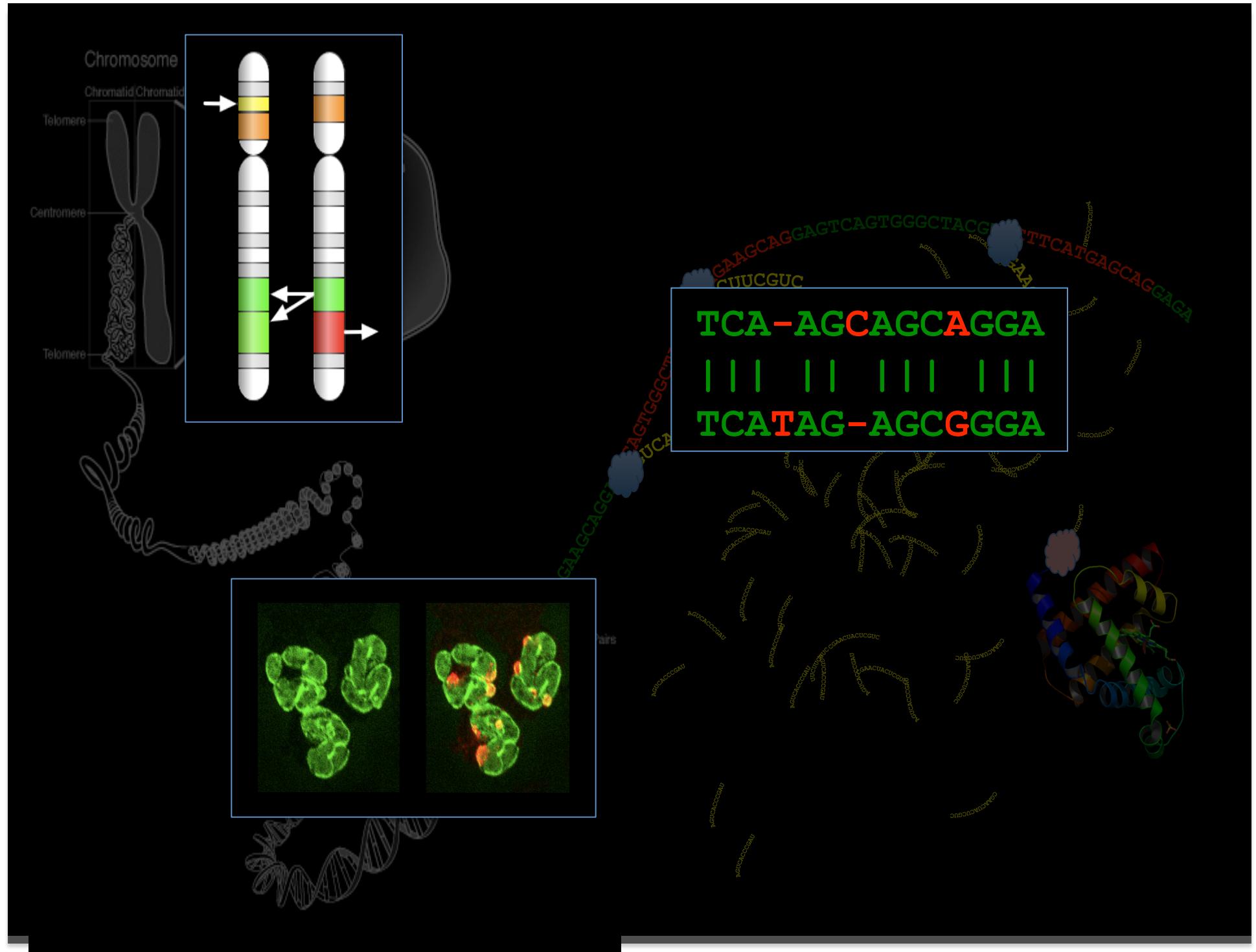
Modern Sequencing



2 weeks
One Sequencer
\$6,000
100-200 Gbp

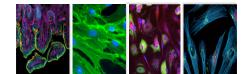
life
technologies™



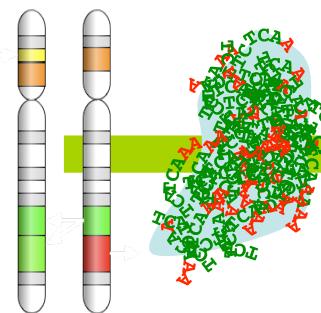




Next Generation Sequencing (NGS)



Next Generation Sequencing (NGS) instruments such as Life Technologies SOLiD™ sequence genomic material in millions of small pieces, enabling a high-level of throughput and sequencing depth.



Sample prep breaks the DNA or RNA into short segments that are attached 500 million to 1 billion beads.



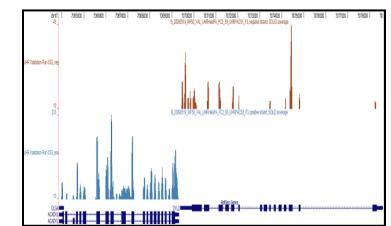
10 GB – 2 TB of raw data is transferred from the instrument to the analysis cluster.

```
>187_29_706_F3  
T2330201030313112312302220311112320021  
0100122001102  
>187_29_800_F3  
T3112001221322200222213012112112211203  
2220323121202  
>187_29_824_F3  
T2221113002302013323132330231030313112  
3123022201211  
>187_29_829_F3  
T2330201000313012312302220311112012212  
3202132301212  
>187_29_858_F3  
T2330201030313112312302220311112322212  
3122122321212  
>187_29_885_F3  
T2330201030313112312302220311112122001  
3212122021222
```

The sequence from each bead is reported in a data file. Files can be over 100 GB.



The bead sequences are either assembled into a new genome or “mapped” to a reference genome.



The mapped reads from RNA samples can be further analyzed to determine which genes are active in the sample.

The Goal: Medical Genomics

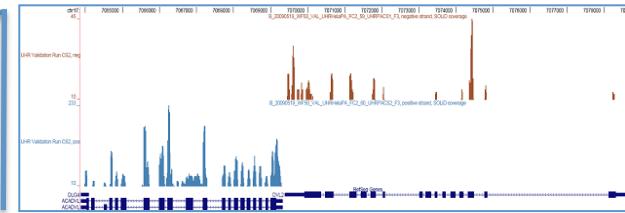
Acquire a Reference



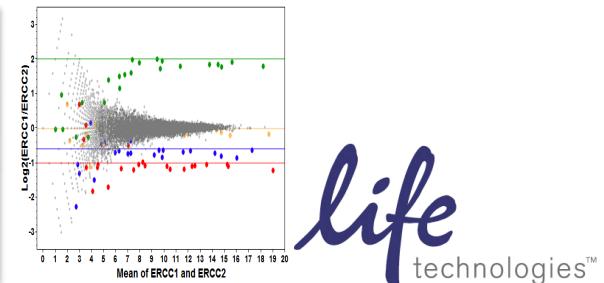
Sequence the diseased genome

GATCAACTTG AGGCCAGCCT GACCAACGTG GCAAAACTCC ATCTCTACTA
ATTAGCTGA GCGTGGTGC ACACACCTGT CATCCCAGT ACTCAGGAGG
AGATGCCCT GACCCCTGGGA GACAGAGGTT GCAGTGAGCT GAGATTCCAT
TAGCCTGGGT GACAGAGTGA ATGGGAGGGA GGAAAAAAA AAAAAGGAA
AGGGAGGAG CCTAGGATGG GGAAGGCTCA CCAGAAGTGG ATGCAAAGAG
GAGCTATTC TATTTGCCCTA GGAAAGAAAA ACCTCCAGAA ACCTGGCCTT
GCCGAGGCC TCCAGGAAG CCAGGCAGAC CCTGCTCCCTG CTCTGACCCC

Sequence the expressed RNA

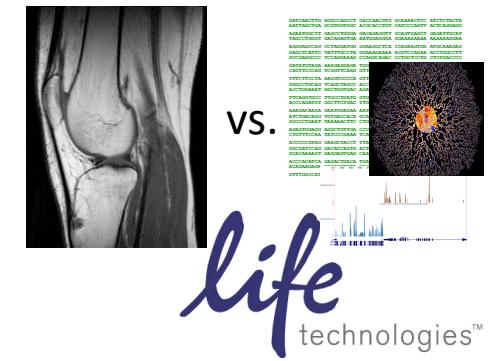
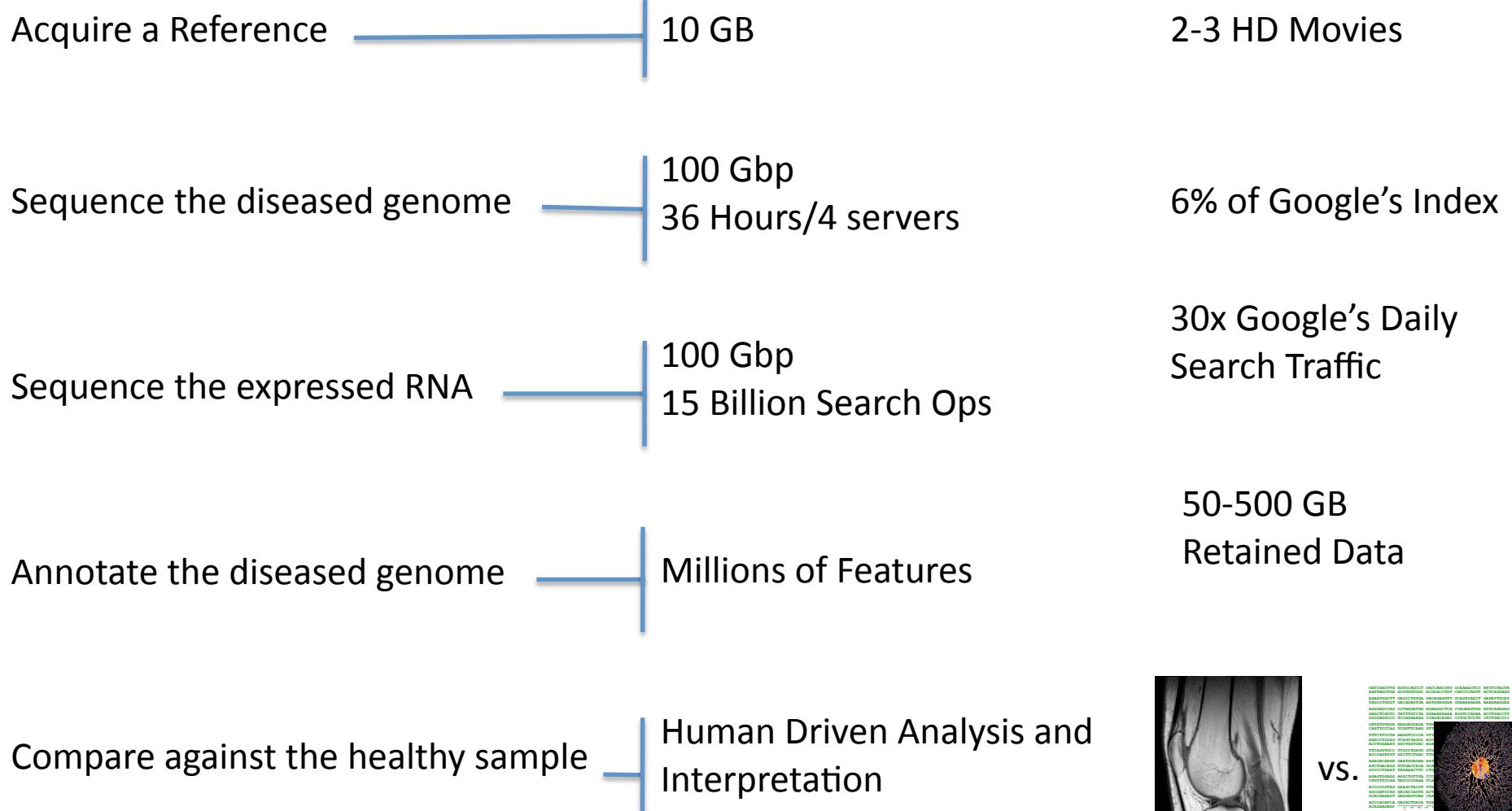


Annotate the diseased genome

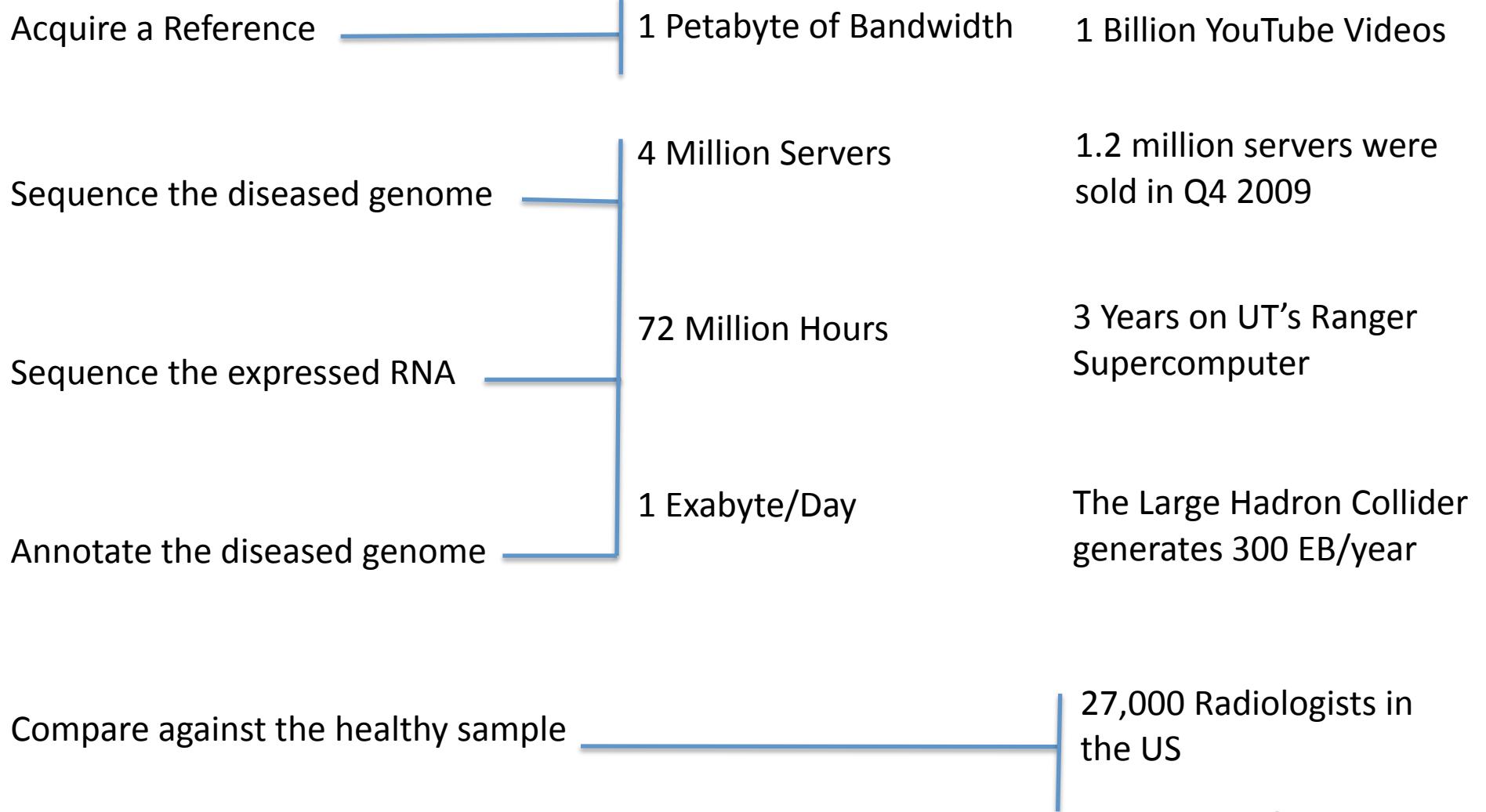


Compare against the healthy sample

Personal Genomics



Population Genomics





So...

Can Python save next generation sequencing?

Probably not on its own...

Better questions:

How can Python help next generation sequencing?

How can next generation sequencing help Python?

NGS Data and Workflow Components

EDA

Novel Apps

Interactions

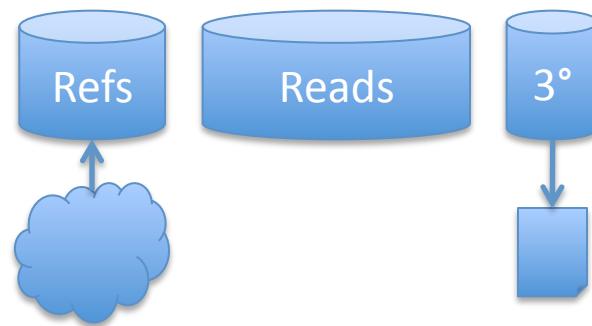
Expression

Variation

Annotation

Assembly

Mapping



Exploratory Analysis

Undirected analysis and unforeseen sequencing applications

Standard Scientific Workflows

Understand the structure and function of genomic elements

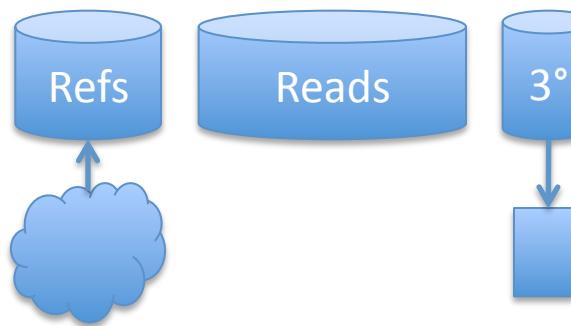
Fundamental Algorithms

Transform the raw data into scientifically relevant forms

References and Data

Reference genomes, domain-specific data sets, raw data, analysis results

Algorithms and Scale

<u>Data and Workflows</u>	<u>Algorithms</u>	<u>Considerations</u>
EDA Novel Apps	Information visualization, data-mining	Interactivity is essential.
Interactions Expression	Clustering, statistical models, network analysis	Data is smaller, but analysis may require round trips back to reads.
Variation Annotation		
Assembly Mapping	Graphs, Index schemes, dynamic programming	Generally I/O bound with opportunities for parallelism.
	Standard data formats are text-based	Can easily span 100s of TBs for a small lab. References are in constant flux.

Software and Hardware

Data and Workflows

EDA

Novel Apps

Interactions

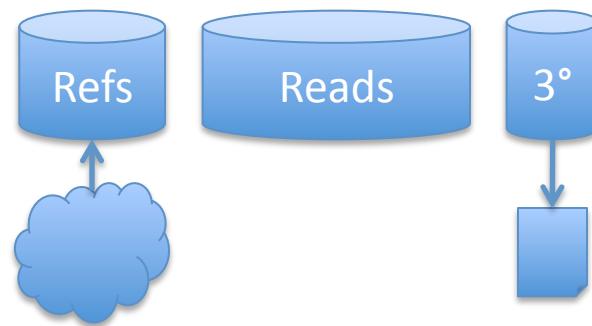
Expression

Variation

Annotation

Assembly

Mapping



Software

Genome browsers, Matlab,
Stats tools, R, etc

Scripting languages, data
analysis libraries

Graph libraries, pipeline
management software, job
schedulers

Databases, ORM tools,
Flat files

Hardware

Workstations, laptops with
fast access to storage

Single nodes with fast
access to storage

Clusters with high-memory
nodes, fast access to storage.

Distributed file systems,
Fast SANs

Python and NGS

Data and Workflows

Current

Potential

EDA

Novel Apps

Galaxy, user-developed tools

GUIs, Disco

Interactions

Expression

User tools, utilities shipped with assemblers and mappers, NGS libraries (HTSeq), NumPy, SciPy

More libraries!

Variation

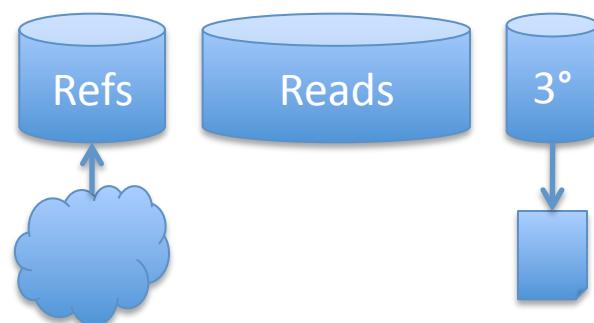
Annotation

Assembly

Mapping

Rapid prototyping, pipeline management, utilities

Multi-processing, queuing to build out pipeline manager



Python parsers for various formats, BioPython, SAMTools
Pygr for reference and annotation access

Disco or Hadoop for distributed read management

Example: Chromosome 20 Expression Analysis



A portion of the Chromosome 20 expression map created using reads from SOLiD™ at Life Technologies' Austin site. The full map spans over 260 feet when printed sequentially. Chromosome 20 represents 2.1% of the human genome.

Visualization

Custom Python tool rendering to PDF using ReportLab



14 Billion reads from three sample types. 1400+ files comprising 4TB of sequence data

Expression

LifeTech WT pipeline for mapping and expression analysis



Human Genome 36.3
Human RefSeq 39

Mapping

Custom Python scripts for result aggregation

life technologies™

Annotation



Cultural Challenges

As recently as five years ago, computing in the Life Sciences consisted of pencils, lab notebooks, and the occasional Excel spreadsheet. Needless to day, NGS caught the community completely off guard.

- Goals
 - Most users care about the science and just want the technology to work
- Skills
 - Few biologists and bioinformaticians are also experts in HPC and tera-scale data mining
- Expectations
 - Excel provides answers instantly. NGS analysis should, too.
- Cost
 - I just spent \$500k on an instrument and I need to spend how much on a computer and developer???



Final Thoughts

- NGS analysis is in its infancy
 - Fundamental methods are still being developed
 - The scientific community is still coming to grips with its potential and challenges
 - There won't be any silver bullets in the near future
- NGS analysis is complicated by the scale of data
 - Traditional supercomputing doesn't help much
- Development paradigms that simplify large data processing will succeed in this space
- Many Python projects show promise, but there's still work to be done!



Thank You!

Life Technologies Austin

Bioinformatics

Jeff Schageman
Joel Brockman
Penn Whitley
Dan Williams

Transcriptomics

Sheila Heater
Kelli Bramlett
Diane Ilsley

The Powers that Be

Bob Setterquist
Tim Sendera

Life Technologies Global

Software and HPC

Lee Jones
Patrick LeGresley
Somalee Datta
Asim Siddiqui
Aaron Kitzmiller

IT

Michael Moore
Antoine Uzzeni
David Morgan
Joanna Curlee

Sounding Boards

Peter Wang
Glen Otero
Travis Oliphant

© 2010 Life Technologies Corporation. All rights reserved. The trademarks mentioned herein are the property of Life Technologies Corporation or their respective owners. Void where prohibited, prohibited where void. For Research Use Only. Not intended for any animal or human therapeutic or diagnostic use.





What People are Using

- Compiled Languages:
 - Equally distributed between C, C++, and Java
- Interpreted Languages:
 - Almost everyone uses Perl, with about half the respondents using R or Python
 - But, for daily usage, about 31% used Python and 47% used Perl
- For Libraries:
 - 75% used BioPerl on a regular basis. SciPy is the next most used at 33%. BioPython was low at 16%
- For Statistics:
 - R (70%), Excel (53%), and Perl (35%) were the most common
- For Visualization:
 - Excel was the most common with R and GNUPlot also getting multiple responses



Compute Environment

- Peta-scale storage environment
- Multi-core processors and multiple nodes
 - though not as many as you think
- 64+ GB of RAM helps
- Heterogeneous processors
 - SIMD/GPUs
- Multiple languages
 - C, C++, and Java are common for ‘fast’ code
 - Many tools ship with Perl or Python utilities
- The cloud will matter at some point



Ideal Stack

- Must be useable by a wide range of developers and scientists
- Standard data formats
- Reference Management
- Optimized kernel operations
- Optimized record-based operations
- Cluster aware pipeline management
- Horizontally scalable record storage
- Database integration for LIMs
- Interactive visualization