# Cloud Computing and the DNA Data Race

## Michael Schatz

October 22, 2010
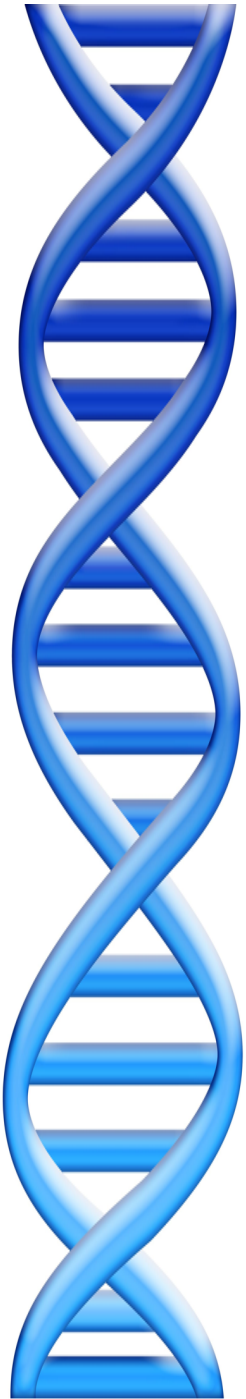CSHL Advanced Sequencing Course

CSH

# Outline

## Part 1: Theory

1. Genome Assembly by Analogy
2. DNA Sequencing and Genomics
3. Sequence Analysis in the Clouds
   1. Sequence Alignment
   2. Mapping & Genotyping
   3. Genome Assembly

## Part 2: Practice

1. AWS Mini-Tutorial
2. Hadoop Mini-Tutorial

# Shredded Book Reconstruction

- Dickens accidentally shreds the first printing of <u>A Tale of Two Cities</u>
  - Text printed on 5 long spools



- How can he reconstruct the text?
  - 5 copies x 138, 656 words / 5 words per fragment = 138k fragments
  - The short fragments from every copy are mixed together
  - Some fragments are identical

# Greedy Reconstruction

It was the best of

age of wisdom, it was

best of times, it was

it was the age of

it was the age of

it was the worst of

of times, it was the

of times, it was the

of wisdom, it was the

the age of wisdom, it

the best of times, it

the worst of times, it

times, it was the age

times, it was the worst

was the age of wisdom,

was the age of foolishness,

was the best of times,

was the worst of times,

wisdom, it was the age

worst of times, it was

---

It was the best of

was the best of times,

the best of times, it

best of times, it was

of times, it was the

of times, it was the

times, it was the worst

times, it was the age

The repeated sequence make the correct reconstruction ambiguous

- It was the best of times, it was the [worst/age]

Model sequence reconstruction as a graph problem.

# de Bruijn Graph Construction

- $D_k = (V, E)$
  - V = All length-k subfragments (k < l)
  - E = Directed edges between consecutive subfragments
    - Nodes overlap by k-1 words

Original Fragment

| It was the best of |
| --- |

Directed Edge

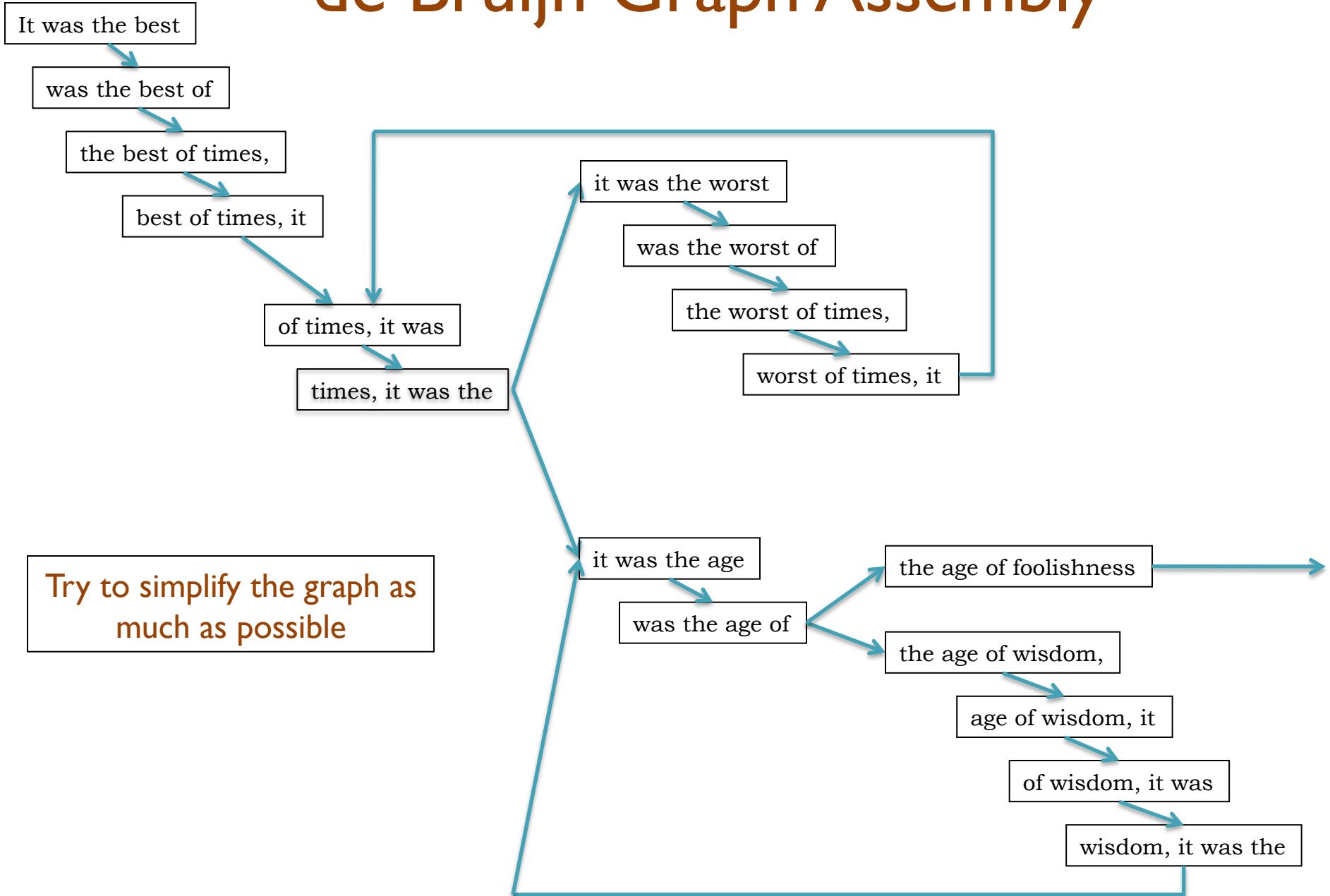| It was the best | → | was the best of |
| --- | --- | --- |

- Locally constructed graph reveals the global sequence structure
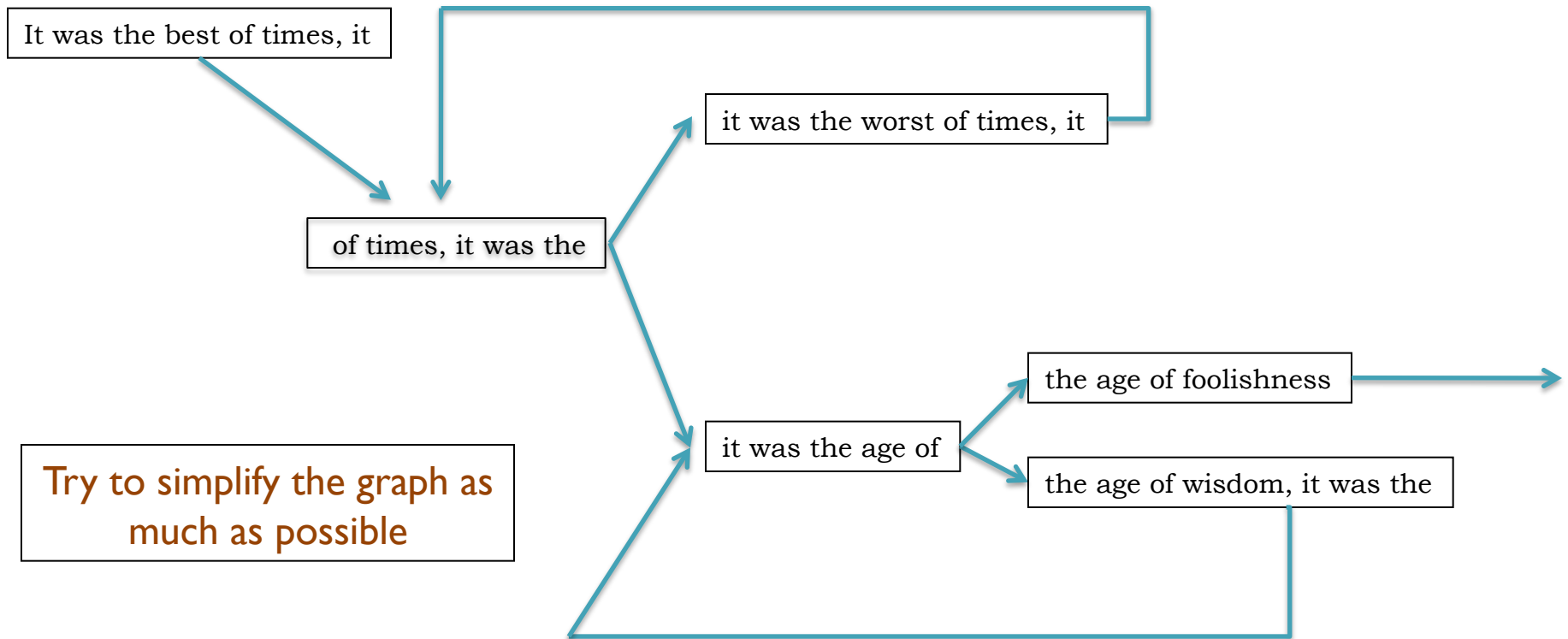  - Overlaps between sequences implicitly computed

de Bruijn, 1946
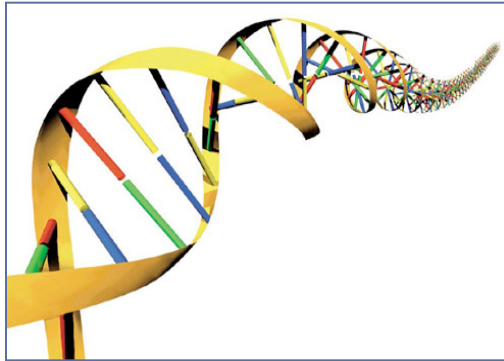Idury and Waterman, 1995
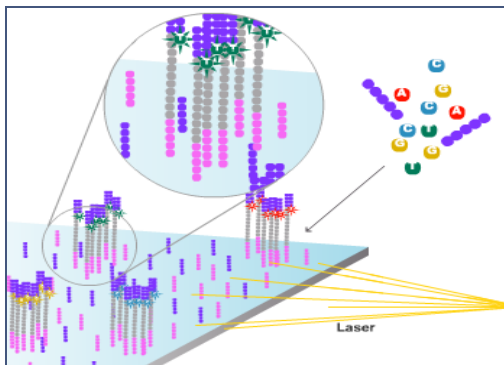Pevzner, Tang, Waterman, 2001

# de Bruijn Graph Assembly

It was the best

was the best of

the best of times,

best of times, it

of times, it was

times, it was the

it was the worst

was the worst of

the worst of times,

worst of times, it

it was the age

was the age of

the age of foolishness

the age of wisdom,

age of wisdom, it

of wisdom, it was

wisdom, it was the

Try to simplify the graph as much as possible

# de Bruijn Graph Assembly

It was the best of times, it

of times, it was the

it was the worst of times, it

Try to simplify the graph as much as possible

it was the age of

the age of foolishness

the age of wisdom, it was the

# Molecular Biology & DNA Sequencing



Genome of an organism encodes the genetic information in long sequence of 4 DNA nucleotides: ACGT

- Bacteria: ~3 million bp
- Humans: ~3 billion bp



Current DNA sequencing machines can sequence millions of short (25-500bp) reads from random positions of the genome

- Per-base error rate estimated at 1-2% (Simpson *et al*, 2009)

ATCTGATAAGTCCCAGGACTTCAGT

GCAAGGCAAACCCGAGCCCAGTTT

TCCAGTTCTAGAGTTTCACATGATC

GGAGTTAGTAAAAGTCCACATTGAG

Like Dickens, we can only sequence small fragments of the genome at once.

- Use software to analyze the sequences
- Modern Biology requires Computational Biology

# The DNA Data Race

| Year | Genome | Technology | Cost |
|------|--------|------------|------|
| 2001 | Venter *et al.* | Sanger (ABI) | $300,000,000 |
| 2007 | Levy *et al.* | Sanger (ABI) | $10,000,000 |
| 2008 | Wheeler *et al.* | Roche (454) | $2,000,000 |
| 2008 | Ley *et al.* | Illumina | $1,000,000 |
| 2008 | Bentley *et al.* | Illumina | $250,000 |
| 2009 | Pushkarev *et al.* | Helicos | $48,000 |
| 2009 | Drmanac *et al.* | Complete Genomics | $4,400 |

(Pushkarev *et al.*, 2009)

Sequencing a single human genome uses ~100 GB of compressed sequence data in billions of short reads.

~20 DVDs / genome

# The DNA Data Tsunami



Use massive amounts of sequencing to explore the genetic origins of life

Our best (only) hope is to use many computers:

- Parallel Computing aka Cloud Computing

- Now your programs will crash on 1000 computers instead of just 1 ☺

# Amazon Web Services

http://aws.amazon.com

- All you need is a credit card, and you can immediately start using one of the largest datacenters in the world

- Elastic Compute Cloud (EC2)
  - On demand computing power
    - Support for Windows, Linux, & OpenSolaris
    - Starting at 8.5¢ / core / hour

- Simple Storage Service (S3)
  - Scalable data storage
    - 10¢ / GB upload fee, 15¢ / GB monthly fee

# Cloud Computing Spectrum

Embarrassingly
Parallel

Loosely
Coupled

Tightly
Coupled



Batch Computing

MapReduce/DryadLINQ

MPI/MapReduce/Pregel

Alignment
HMM Scoring

Genotyping
K-mer Counting

Graph Analysis
Genome Assembly

Scheduling +
Load Balance

Embarrassingly Parallel +
Parallel Communication

Loosely Coupled +
Parallel Algorithm Design

# Embarrassingly Parallel

- Batch computing
  - Each item is independent
  - Split input into many chunks
  - Process each chunk separately on a different computer

- Challenges
  - Distributing work, load balancing, monitoring & restart

- Technologies
  - Condor, Sun Grid Engine
  - Amazon Simple Queue

# Elementary School Dance

# Loosely Coupled

- Divide and conquer
  - Independently process many items
  - Group partial results
  - Scan partial results into final answer

- Challenges
  - Batch computing challenges
  - + Shuffling of huge datasets

- Technologies
  - Hadoop, Elastic MapReduce, Dryad
  - Parallel Databases

# Junior High Dance

# Hadoop MapReduce

http://hadoop.apache.org

- MapReduce is the parallel distributed framework invented by Google for large data computations.
  - Data and computations are spread over thousands of computers, processing petabytes of data each day (Dean and Ghemawat, 2004)
  - Indexing the Internet, PageRank, Machine Learning, etc…
  - Hadoop is the leading open source implementation

- Benefits
  - Scalable, Efficient, Reliable
  - Easy to Program
  - Runs on commodity computers

- Challenges
  - Redesigning / Retooling applications
    - Not Condor, Not MPI
    - Everything in MapReduce

# K-mer Counting

- Application developers focus on 2 (+1 internal) functions
  - Map: input ➜ key:value pairs
  - Shuffle: Group together pairs with same key
  - Reduce: key, value-lists ➜ output

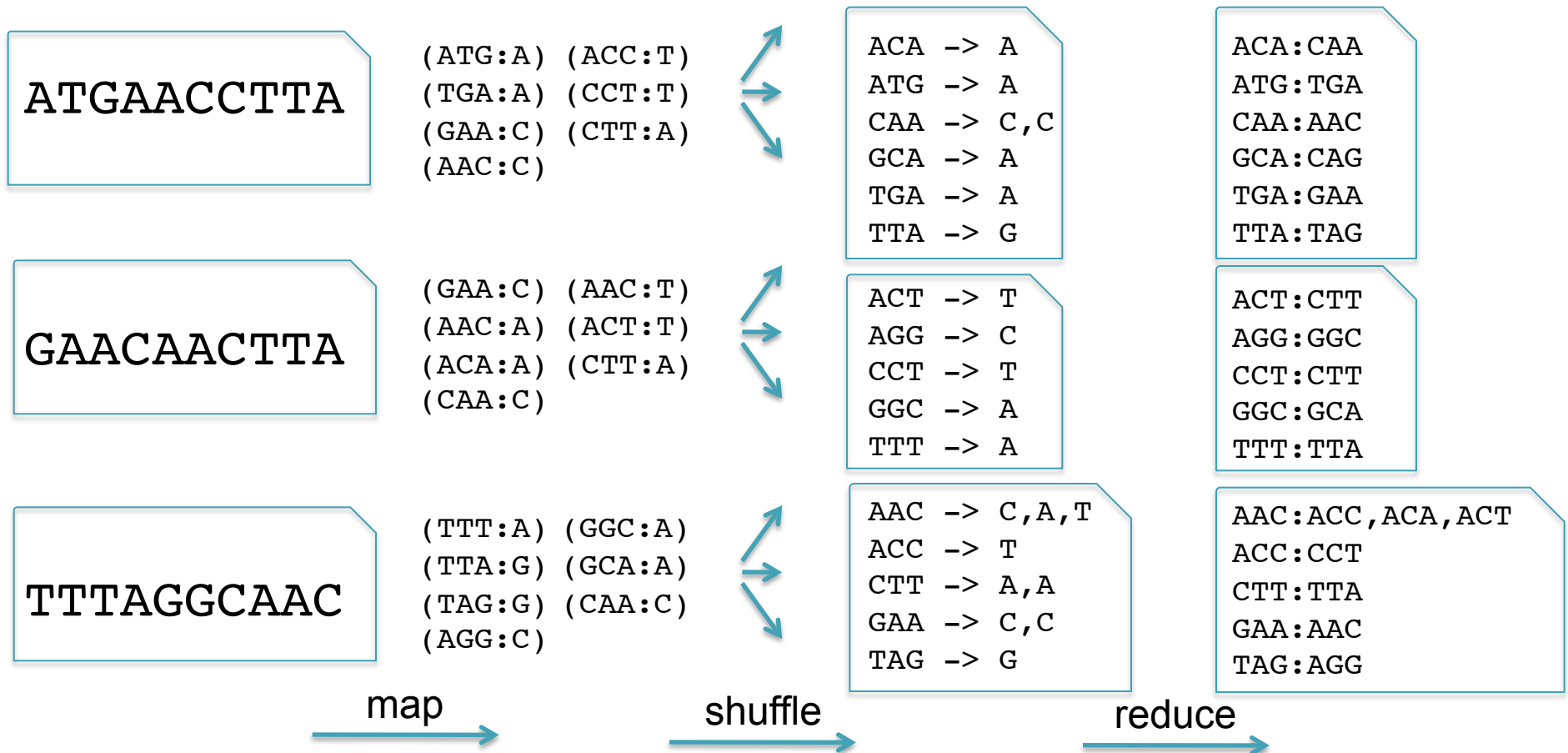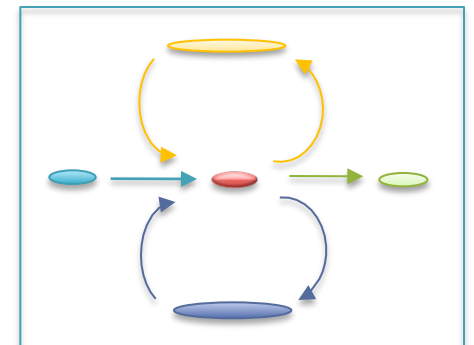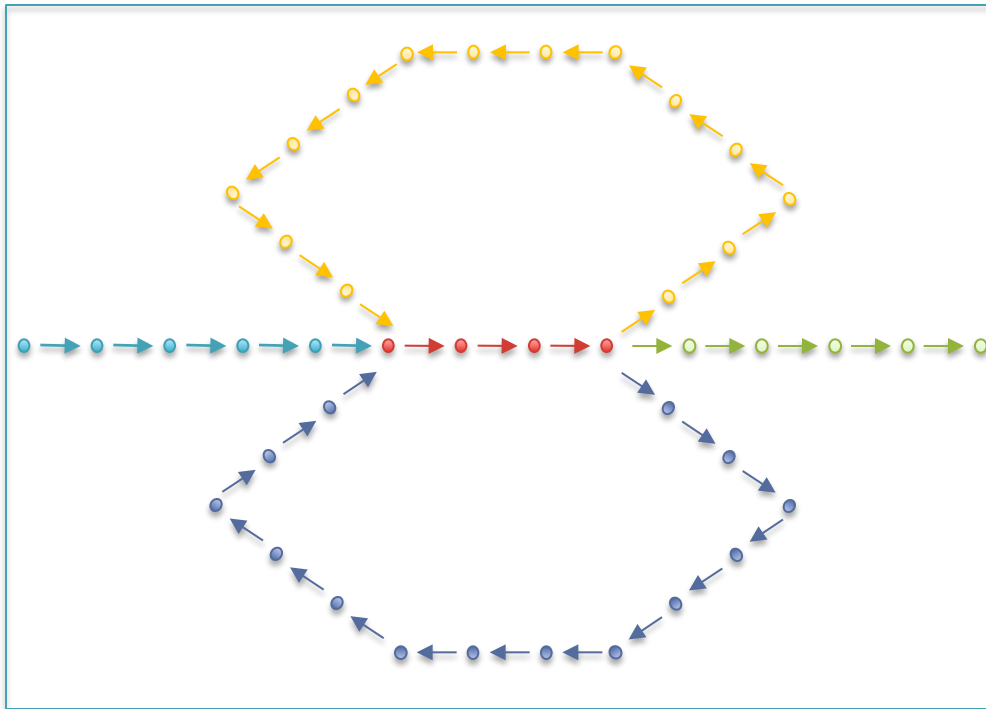Map, Shuffle & Reduce
All Run in Parallel

ATGAACCTTA

```
(ATG:1) (ACC:1)
(TGA:1) (CCT:1)
(GAA:1) (CTT:1)
(AAC:1) (TTA:1)
```

```
ACA -> 1
ATG -> 1
CAA -> 1,1
GCA -> 1
TGA -> 1
TTA -> 1,1,1
```

```
ACA:1
ATG:1
CAA:2
GCA:1
TGA:1
TTA:3
```

GAACAACTTA

```
(GAA:1) (AAC:1)
(AAC:1) (ACT:1)
(ACA:1) (CTT:1)
(CAA:1) (TTA:1)
```

```
ACT -> 1
AGG -> 1
CCT -> 1
GGC -> 1
TTT -> 1
```

```
ACT:1
AGG:1
CCT:1
GGC:1
TTT:1
```

TTTAGGCAAC

```
(TTT:1) (GGC:1)
(TTA:1) (GCA:1)
(TAG:1) (CAA:1)
(AGG:1) (AAC:1)
```

```
AAC -> 1,1,1,1
ACC -> 1
CTT -> 1,1
GAA -> 1,1
TAG -> 1
```

```
AAC:4
ACC:1
CTT:1
GAA:1
TAG:1
```

map          shuffle          reduce

# Hadoop Architecture



- ## Hadoop Distributed File System (HDFS)
  - Data files partitioned into large chunks (64MB), replicated on multiple nodes
  - Computation moves to the data, rack-aware scheduling

- ## Hadoop MapReduce system won the 2009 GreySort Challenge
  - Sorted 100 TB in 173 min (578 GB/min) using 3452 nodes and 4x3452 disks

# Short Read Mapping

Identify variants

```
                                                                    GGTATAC…
          …CCATAG       TATGCGCCC      CGG A AATTT  CGGTATAC
          …CCAT       CTATATGCG          TCGG A AATT    CGGTATAC
          …CCAT  GGCTATATG         CTATCGG A A A    GCGGTATA
Subject   …CCA  AGGCTATAT         CCTATCGG A      TTGCGGTA   C…
          …CCA  AGGCTATAT     GCCCTATCG       TTTGCGGT     C…
          …CC    AGGCTATAT     GCCCTATCG   A AATTTGC     ATAC…
          …CC  TAGGCTATA  GCGCCCTA        A AATTTGC  GTATAC…

Reference  …CCATAGGCTATATGCGCCCTATCGG C AATTTGCGGTATAC…
```

- Given a reference and many subject reads, report one or more "good" end-to-end alignments per alignable read

  - Find where the read most likely originated

  - Fundamental computation for many assays

    - Genotyping              RNA-Seq              Methyl-Seq
    - Structural Variations   Chip-Seq             Hi-C-Seq

- Desperate need for scalable solutions

  - Single human requires >1,000 CPU hours / genome

# Crossbow

http://bowtie-bio.sourceforge.net/crossbow

- ## Align billions of reads and find SNPs
  - Reuse software components: Hadoop Streaming

- ## Map: Bowtie (Langmead *et al.*, 2009)
  - Find best alignment for each read
  - Emit (chromosome region, alignment)

- ## Shuffle: Hadoop
  - Group and sort alignments by region

- ## Reduce: SOAPsnp (Li *et al.*, 2009)
  - Scan alignments for divergent columns
  - Accounts for sequencing error, known SNPs

# Performance in Amazon EC2

http://bowtie-bio.sourceforge.net/crossbow

| | Asian Individual Genome | | |
|---|---|---|---|
| **Data Loading** | 3.3 B reads | 106.5 GB | $10.65 |
| **Data Transfer** | 1h :15m | 40 cores | $3.40 |
| | | | |
| **Setup** | 0h : 15m | 320 cores | $13.94 |
| **Alignment** | 1h : 30m | 320 cores | $41.82 |
| **Variant Calling** | 1h : 00m | 320 cores | $27.88 |
| | | | |
| **End-to-end** | 4h : 00m | | $97.69 |

Analyze an entire human genome for ~$100 in an afternoon.
Accuracy validated at >99%

**Searching for SNPs with Cloud Computing.**
Langmead B, Schatz MC, Lin J, Pop M, Salzberg SL (2009) *Genome Biology.* **10**:R134

# Map-Shuffle-Scan for Genomics



**Cloud Computing and the DNA Data Race.**
Schatz, MC, Langmead B, Salzberg SL (2010) *Nature Biotechnology* 28:691-693

# Tightly Coupled

- Computation that cannot be partitioned
  - Graph Analysis
  - Molecular Dynamics
  - Population simulations

- Challenges
  - Loosely coupled challenges
  - + Parallel algorithms design

- Technologies
  - MPI
  - MapReduce, Dryad, Pregel

# Short Read Assembly

**Reads**

AAGA
ACTT
ACTC
ACTG
AGAG
CCGA
CGAC
CTCC
CTGG
CTTT
...



de Bruijn Graph

Potential Genomes

AAGACTCCGACTGGGACTTT

AAGACTGGGACTCCGACTTT

- Genome assembly as finding an Eulerian tour of the de Bruijn graph
  - Human genome: >3B nodes, >10B edges

- The new short read assemblers require tremendous computation
  - Velvet (Zerbino & Birney, 2008) serial: > 2TB of RAM
  - ABySS (Simpson *et al.*, 2009) MPI: 168 cores x ~96 hours
  - SOAPdenovo (Li *et al.*, 2010) pthreads: 40 cores x 40 hours, >140 GB RAM

# K-mer Counting

- Application developers focus on 2 (+1 internal) functions
  - Map: input ➜ key:value pairs
  - Shuffle: Group together pairs with same key
  - Reduce: key, value-lists ➜ output

Map, Shuffle & Reduce
All Run in Parallel

ATGAACCTTA

```
(ATG:1) (ACC:1)
(TGA:1) (CCT:1)
(GAA:1) (CTT:1)
(AAC:1) (TTA:1)
```

```
ACA -> 1
ATG -> 1
CAA -> 1,1
GCA -> 1
TGA -> 1
TTA -> 1,1,1
```

```
ACA:1
ATG:1
CAA:2
GCA:1
TGA:1
TTA:3
```

GAACAACTTA

```
(GAA:1) (AAC:1)
(AAC:1) (ACT:1)
(ACA:1) (CTT:1)
(CAA:1) (TTA:1)
```

```
ACT -> 1
AGG -> 1
CCT -> 1
GGC -> 1
TTT -> 1
```

```
ACT:1
AGG:1
CCT:1
GGC:1
TTT:1
```

TTTAGGCAAC

```
(TTT:1) (GGC:1)
(TTA:1) (GCA:1)
(TAG:1) (CAA:1)
(AGG:1) (AAC:1)
```

```
AAC -> 1,1,1,1
ACC -> 1
CTT -> 1,1
GAA -> 1,1
TAG -> 1
```

```
AAC:4
ACC:1
CTT:1
GAA:1
TAG:1
```

map          shuffle          reduce

# Graph Construction

- Application developers focus on 2 (+1 internal) functions
  - Map: input ➜ key:value pairs
  - Shuffle: Group together pairs with same key
  - Reduce: key, value-lists ➜ output

Map, Shuffle & Reduce
All Run in Parallel

```
ATGAACCTTA
```
```
(ATG:A) (ACC:T)
(TGA:A) (CCT:T)
(GAA:C) (CTT:A)
(AAC:C)
```
```
ACA -> A
ATG -> A
CAA -> C,C
GCA -> A
TGA -> A
TTA -> G
```
```
ACA:CAA
ATG:TGA
CAA:AAC
GCA:CAG
TGA:GAA
TTA:TAG
```

```
GAACAACTTA
```
```
(GAA:C) (AAC:T)
(AAC:A) (ACT:T)
(ACA:A) (CTT:A)
(CAA:C)
```
```
ACT -> T
AGG -> C
CCT -> T
GGC -> A
TTT -> A
```
```
ACT:CTT
AGG:GGC
CCT:CTT
GGC:GCA
TTT:TTA
```

```
TTTAGGCAAC
```
```
(TTT:A) (GGC:A)
(TTA:G) (GCA:A)
(TAG:G) (CAA:C)
(AGG:C)
```
```
AAC -> C,A,T
ACC -> T
CTT -> A,A
GAA -> C,C
TAG -> G
```
```
AAC:ACC,ACA,ACT
ACC:CCT
CTT:TTA
GAA:AAC
TAG:AGG
```

map                    shuffle                    reduce

# Graph Compression

- ## After construction, many edges are unambiguous
  - Merge together compressible nodes
  - Graph physically distributed over hundreds of computers

# High School Dance

# Warmup Exercise

- Who here was born closest to October 22?
  - You can only compare to 1 other person at a time



Find winner among 64 teams in just 6 rounds

# Fast Path Compression

## Challenges

- Nodes stored on different computers
- Nodes can only access direct neighbors

## Randomized List Ranking

- Randomly assign (H)/[T] to each compressible node
- Compress (H)➔[T] links



Initial Graph: 42 nodes

**Randomized Speed-ups in Parallel Computation.**
Vishkin U. (1984) *ACM Symposium on Theory of Computation. 230-239.*

# Fast Path Compression

## Challenges

- Nodes stored on different computers
- Nodes can only access direct neighbors

## Randomized List Ranking

- Randomly assign ⒣/☐T to each compressible node
- Compress ⒣➔☐T links



Round 1: 26 nodes (38% savings)

**Randomized Speed-ups in Parallel Computation.**
Vishkin U. (1984) *ACM Symposium on Theory of Computation. 230-239.*

# Fast Path Compression

## Challenges

– Nodes stored on different computers

– Nodes can only access direct neighbors

## Randomized List Ranking

– Randomly assign (H)/ [T] to each compressible node

– Compress (H)→[T] links



Round 2: 15 nodes (64% savings)

**Randomized Speed-ups in Parallel Computation.**
Vishkin U. (1984) *ACM Symposium on Theory of Computation. 230-239.*

# Fast Path Compression

## Challenges

– Nodes stored on different computers

– Nodes can only access direct neighbors

## Randomized List Ranking

– Randomly assign (H)/[T] to each compressible node

– Compress (H)→[T] links



Round 2: 8 nodes (81% savings)

**Randomized Speed-ups in Parallel Computation.**
Vishkin U. (1984) *ACM Symposium on Theory of Computation. 230-239.*

# Fast Path Compression

## Challenges

- Nodes stored on different computers
- Nodes can only access direct neighbors

## Randomized List Ranking

- Randomly assign (H)/ T to each compressible node
- Compress (H)→T links



Round 3: 6 nodes (86% savings)

**Randomized Speed-ups in Parallel Computation.**
Vishkin U. (1984) *ACM Symposium on Theory of Computation. 230-239.*

# Fast Path Compression

## Challenges

– Nodes stored on different computers

– Nodes can only access direct neighbors

## Randomized List Ranking

– Randomly assign (H)/ [T] to each compressible node

– Compress (H)→[T] links

## Performance

– Compress all chains in $\log(S)$ rounds

Round 4: 5 nodes (88% savings)

**Randomized Speed-ups in Parallel Computation.**
Vishkin U. (1984) *ACM Symposium on Theory of Computation. 230-239.*

# Node Types



Isolated nodes (10%)

Tips (46%)

Bubbles/Non-branch (9%)

Dead Ends (.2%)

Half Branch (25%)

Full Branch (10%)

(Chaisson, 2009)

# Contrail

http://contrail-bio.sourceforge.net

## Scalable Genome Assembly with MapReduce

- *Genome: E. coli* K12 MG1655, 4.6Mbp

- *Input:* 20.8M 36bp reads, 200bp insert (~150x coverage)

- *Preprocessor*: Quality-Aware Error Correction

| | Initial | Compressed | Error Correction | Resolve Repeats | Cloud Surfing |
|---|---|---|---|---|---|
| N | 5.1 M | 245,131 | 2,769 | 1,909 | 300 |
| Max | 27 bp | 1,079 bp | 70,725 bp | 90,088 bp | 149,006 bp |
| N50 | 27 bp | 156 bp | 15,023 bp | 20,062 bp | 54,807 bp |

**Assembly of Large Genomes with Cloud Computing.**
Schatz MC, Sommer D, Kelley D, Pop M, *et al. In Preparation.*

# E. coli Assembly Quality

Incorrect contigs: Align at < 95% identity or < 95% of their length

| Assembler | Contigs ≥ 100bp | N50 (bp) | Incorrect contigs |
|---|---|---|---|
| Contrail PE | 300 | 54,807 | 4 |
| Contrail SE | 529 | 20,062 | 0 |
| SOAPdenovo PE | 182 | 89,000 | 5 |
| ABySS PE | 233 | 45,362 | 13 |
| Velvet PE | 286 | 54,459 | 9 |
| EULER-SR PE | 216 | 57,497 | 26 |
| SSAKE SE | 931 | 11,450 | 38 |
| Edena SE | 680 | 16,430 | 6 |

# Contrail

http://contrail-bio.sourceforge.net

De novo Assembly of the Human Genome

- *Genome:* African male NA18507 (SRA000271, Bentley *et al.*, 2008)
- *Input: 3.5*B 36bp reads, 210bp insert (~40x coverage)

| | Initial | Compressed | Error Correction | Resolve Repeats | Cloud Surfing |
|---|---|---|---|---|---|
| N | >7 B | >1 B | 4.2 M | 4.1 M | 3.3 M |
| Max | 27 bp | 303 bp | 20,594 bp | 20,594 bp | 20,594 bp |
| N50 | 27 bp | < 100 bp | 995 bp | 1,050 bp | 1,427 bp* |

**Assembly of Large Genomes with Cloud Computing.**
Schatz MC, Sommer D, Kelley D, Pop M, *et al. In Preparation.*

# Hadoop for NGS Analysis

## Quake



Quality-aware error correction of short reads

*Correct 97.9% of errors with 99.9% accuracy*

http://www.cbcb.umd.edu/software/quake/

(Kelley, Schatz, Salzberg, 2010*)

## CloudBurst



Highly Sensitive Short Read Mapping with MapReduce

*100x speedup mapping on 96 cores @ Amazon*

http://cloudburst-bio.sf.net

(Schatz, 2009)

## Myrna

Cloud-scale differential gene expression for RNA-seq

*Expression of 1.1 billion RNA-Seq reads in <2 hours for ~$66*



(Langmead, Hansen, Leek, 2010)

http://bowtie-bio.sf.net/myrna/

## AMOS

Searching for SNPs in the Turkey Genome

*Scan the de novo assembly to find 920k hetrozygous alleles*



(Dalloul et al, 2010)

http://amos.sf.net

# Summary

- Surviving the data deluge means computing in parallel
    - Cloud computing is an attractive platform for large scale sequence analysis and computation

- Use the right tool for the job
    - Embarassingly parallel = Condor/Hadoop
    - Loosely coupled = Hadoop/Dyrad
    - Tightly coupled = MPI/Hadoop

- Emerging technologies are a great start, but we need continued research
    - A word of caution: new technologies are new

# Acknowledgements

## Advisor

Steven Salzberg

## UMD Faculty

Mihai Pop, Art Delcher, Amitabh Varshney,
Carl Kingsford, Ben Shneiderman,
James Yorke, Jimmy Lin, Dan Sommer

## CBCB Students

Adam Phillippy, Cole Trapnell,
Saket Navlakha, Ben Langmead,
James White, David Kelley

# Break

# Outline

## Part 1: Theory
1. Genome Assembly by Analogy
2. DNA Sequencing and Genomics
3. Sequence Analysis in the Clouds
   1. Sequence Alignment
   2. Mapping & Genotyping
   3. Genome Assembly

## Part 2: Practice
1. AWS Mini-Tutorial
2. Hadoop Mini-Tutorial

# A Brief History of the Amazon Cloud

- Urban Legend
  - Additional capacity added every fall for the holiday shopping season, underutilized rest of the year…

- Official Story
  - Amazon is a technology company
    - Different divisions of Amazon share computation
  - Amazon Web Services is the 3$^{rd}$ Business Division
    - Retail & Seller Businesses

amazon.com

# Amazon Web Services

http://aws.amazon.com

- All you need is a credit card, and you can immediately start using one of the largest datacenters in the world

- Elastic Compute Cloud (EC2)
  - On demand computing power
    - Support for Windows, Linux, & OpenSolaris
    - Starting at 8.5¢ / core / hour

- Simple Storage Service (S3)
  - Scalable data storage
    - 10¢ / GB upload fee, 15¢ / GB monthly fee

# Other Services



Everything you need to run a large scale service & analysis suite in the clouds.

# Cloud Solutions for DNA Sequence Analysis

- Rapid & elastic deployment of vast computation
  - Accessible, Reproducible, Reliable, Collaborative

- Why not?
  - Transfer: 200 GBs takes 1 hr – 2 weeks
  - Privacy & security: Excellent… with care
  - Expertise: Computing on 1000 cores is complex
  - Price: The meter is always running

  - Who will be on the line for making it all work?
    - Psychological and Institutional paradigm shift

# EC2 Architecture

- Very large pool of machines
  - Effectively infinite resources
  - High-end servers with many cores and many GB RAM

- Machines run in a virtualized environment
  - Amazon can subdivide large nodes into smaller instances
  - You are 100% protected from other users on the machine
  - You get to pick the operating system, all installed software

# Instance Types

| Type | Price / hr | CPU | Resources |
|------|-----------|-----|-----------|
| **Micro** <br> *Web service* | 2¢ | 1 core @ 1 ECU | .6 GB RAM <br> 10 GB Disk |
| **Standard** <br> *Light Tasks* | 8.5¢ | 1 core @ 1 ECU | 1.7 GB RAM <br> 160 GB Disk |
| **Extra Large** <br> *Mapping w/BWA* | 68¢ | 4 core @ 2 ECU | 15 GB RAM <br> 1.7 TB Disk |
| **HighCPU XL** <br> *Mapping w/Crossbow* | 68¢ | 8 cores @ 2.5 ECU | 7 GB RAM <br> 1.7 TB Disk |
| **HighMem Quad XL** <br> *Assembly* | $2.00 | 8 cores @ 3.25 ECU | 68.4 GB RAM <br> 1.7 TB Disk |

ECU = EC2 Compute Unit. Approximately 1.0 – 1.2 GHz Intel Xeon from 2007

Reserved Instances make it cheaper for consistent use.
Pay for what you use, rounded UP to the next full hour

# Amazon Machine Images



- A few Amazon sponsored images
  - Suse Linux, Windows

- Many Community Images & Appliances
  - Crossbow: Hadoop, Bowtie, SOAPsnp
  - CloudBioLinux.com: Appliance for Genomics

- Build you own
  - Completely customize your environment
  - You results could be totally reproducible

# Amazon S3

- S3 provides persistent storage for large volumes of data
  - Very high speed connection from S3 to EC2 compute nodes
  - Public data sets include s3://1000genomes

- Tiered pricing by volume
  - Pricing starts at 15¢ / GB / month
  - 5.5¢ / GB / month for over 5 PB
  - Pay for transfer in and out of Amazon

- Import/Export service for large volumes
  - FedEx your drives to Amazon

# Getting Started

http://docs.amazonwebservices.com/AWSEC2/latest/GettingStartedGuide/

# Signing Up

# AWS Management Console

# Running your First Cloud Analysis

1. Pick your AMI
   – Machine Image: Operating System & Tools
2. Pick your instance type & quantity
   – Micro - High-Memory Quadruple Extra Large
3. Pick your credentials
   – SSH Keys
4. Configure your Firewall
   – Protect your servers
5. Launch!

# 1. Pick your AMIs

# CloudBioLinux

# 2. Pick your Instance Type

# 3. Pick your Credentials

# 4. Configure your Firewall

# 5. Launch!

# Monitoring your Server

# Connecting (1)

# Connecting (2)

# Calling SNPs in the Cloud ☺

```
chmod 400 mschatz.pem

scp -r -i mschatz.pem data.tgz ubuntu@ec2-174-129-123-73.compute-1.amazonaws.com:
ssh -i mschatz.pem ubuntu@ec2-174-129-123-73.compute-1.amazonaws.com

<remote>

ls

tar xzvf data.tgz
bowtie -S data/genomes/e_coli data/reads/e_coli_10000snp.fq ec_snp.sam
samtools view -bS -o ec_snp.bam ec_snp.sam
samtools sort ec_snp.bam ec_snp.sorted

samtools pileup -cv -f data/genomes/NC_008253.fna ec_snp.sorted.bam > snps

samtools index ec_snp.sorted.bam
samtools tview ec_snp.sorted.bam data/genomes/NC_008253.fna

exit

<local>

scp -i mschatz.pem ubuntu@ec2-174-129-123-73.compute-1.amazonaws.com:snps .
```

# 1000Genomes in the Cloud

```
s3cmd --configure

# cp data/.s3cfg .

s3cmd ls s3://1000genomes

s3cmd ls s3://1000genomes/Pilots_Bam/NA20828/

s3cmd get s3://1000genomes/Pilots_Bam/NA20828/*chr22* .

samtools view NA20828.SLX.maq.SRP000033.2009_09.chr22_1_49691432.bam
```

# Terminating



Total cost: 8.5¢

# Reflections

- Launching and managing virtual clusters with the AWS Console is quick and easy
  - Entirely scriptable using ec2 tools
  - iPhone App also available

- Things get really interesting on 168 cores
  - 1 week CPU = 1 hour wall

# Hadoop on AWS



Just 3 commands to bring up a 168 core (21 node) cluster & crunch terabytes:

$HADOOP/src/contrib/ec2/bin/hadoop-ec2 launch-cluster HADOOP 21

$HADOOP/src/contrib/ec2/bin/hadoop-ec2 <hadoop cmd> HADOOP

$HADOOP/src/contrib/ec2/bin/hadoop-ec2 terminate-cluster HADOOP

# Kmer Code

**kmer-map.pl**

```perl
#!/usr/bin/perl

my $K = 15;

while (<>)
{
 for (my $i = 0;
      $i < length($_)-$K;
      $i++)
 {
  my $kmer = substr($_,$i,$K);
  print "$kmer\t1\n";
 }
}
```

**kmer-reduce.pl**

```perl
#!/usr/bin/perl

my $mer = "";
my $cnt = 0;

while (<>)
{
  chomp;

  my ($curmer, $curcnt) = split /\t/, $_;

  if ($curmer ne $mer)
  {
    print "$mer\t$cnt\n" if ($cnt > 0);
    $mer = $curmer; $cnt = 0;
  }

 $cnt += $curcnt;
}

print "$mer\t$cnt\n" if ($cnt > 0);
```

# BashReduce

```
$ head -3 reads.txt
ATATTTTTTCTTGTTTTTTTATATCCACAAACTCTTT
CCACAAAATCAATACCTTGTGGAATAAAATTGTCCA
TATTTTTTTCTTGTTTTTTATATCCACAAACTCTTTT


$ cat reads.txt | ./kmer-map.pl | head -3
ATATTTTTTCTTGTT      1
TATTTTTTCTTGTTT      1
ATTTTTTTCTTGTTTT     1


$ cat reads.txt | ./kmer-map.pl | sort \
  | ./kmer-reduce.pl | head -3
AAAAAAAAGTAGCTA      44
AAAAAAAGTAGCTAT      44
AAAAAAGCAAATGTG      17
```

# Kmer Counting In Hadoop

```sh
#!/bin/sh

STREAMING=/usr/lib/hadoop-0.20/contrib/streaming/hadoop-
    streaming-0.20.2+320.jar

hadoop fs -mkdir /user/mschatz/kmertest/reads
hadoop fs -put reads.txt /user/mschatz/kmertest/reads

hadoop jar $STREAMING \
 -input  /user/mschatz/kmertest/reads \
 -output /user/mschatz/kmertest/kmers \
 -mapper ./kmer-map.pl \
 -reducer ./kmer-reduce.pl \
 -file ./kmer-map.pl \
 -file ./kmer-reduce.pl \
 -jobconf mapred.map.tasks=10 \
 -jobconf mapred.reduce.tasks=1

hadoop fs -cat /user/mschatz/kmertest/kmers/part-* | head -3
hadoop fs -rmr /user/mschatz/kmertest
```

# Hadoop Output

```
10/10/21 16:03:51 INFO mapred.FileInputFormat: Total input paths to process : 1
10/10/21 16:03:51 INFO streaming.StreamJob: getLocalDirs(): [/scratch0/hadoop/mapred/
    local]
10/10/21 16:03:51 INFO streaming.StreamJob: Running job: job_201009232028_2089
10/10/21 16:03:51 INFO streaming.StreamJob: To kill this job, run:
10/10/21 16:03:51 INFO streaming.StreamJob: /usr/lib/hadoop-0.20/bin/hadoop job  -
    Dmapred.job.tracker=szhdname01.umiacs.umd.edu:8021 -kill job_201009232028_2089
10/10/21 16:03:51 INFO streaming.StreamJob: Tracking URL: http://
    szhdname01.umiacs.umd.edu:50030/jobdetails.jsp?jobid=job_201009232028_2089
10/10/21 16:03:52 INFO streaming.StreamJob:  map 0%   reduce 0%
10/10/21 16:03:58 INFO streaming.StreamJob:  map 30%   reduce 0%
10/10/21 16:04:01 INFO streaming.StreamJob:  map 100%   reduce 0%
10/10/21 16:04:07 INFO streaming.StreamJob:  map 100%   reduce 20%
10/10/21 16:04:16 INFO streaming.StreamJob:  map 100%   reduce 100%
10/10/21 16:04:19 INFO streaming.StreamJob: Job complete: job_201009232028_2089
10/10/21 16:04:19 INFO streaming.StreamJob: Output: /user/mschatz/kmertest/kmers
```

**AAAAAAAAGTAGCTA          44**

**AAAAAAAGTAGCTAT          44**

**AAAAAAGCAAATGTG          17**

# Crossbow Webform

http://bowtie-bio.sf.net/crossbow/ui.html



- Enter your account info, manifest file, reference info, and pipeline settings
  - List of URLs to fastq files

- Crossbow
  - Parallel ftp
  - Parallel map
  - Parallel SNPs

# More Information

- Amazon Web Services
  - http://aws.amazom.com
  - http://aws.amazon.com/free

- Hadoop
  - http://hadoop.apache.org

- Crossbow & Bowtie
  - http://bowtie-bio.sf.net

# Thank You!

http://schatzlab.cshl.edu

@mike_schatz