

# Entering the era of mega-genomics

Michael Schatz

Simons Center for Quantitative Biology

Feb 7, 2010

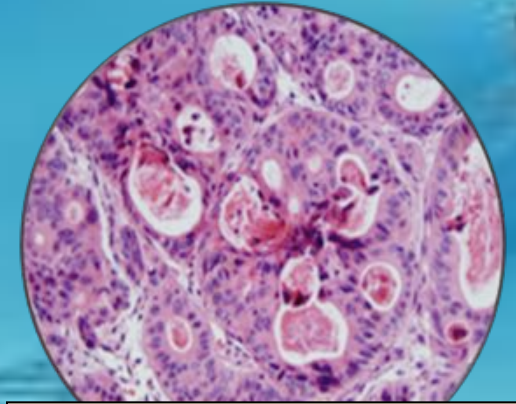
Pioneer/Dupont Des Moines



# Schatz Lab Overview



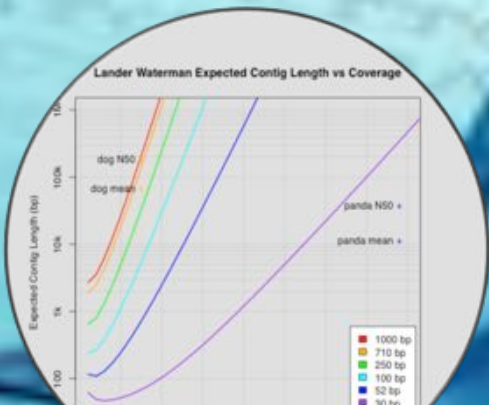
Computation



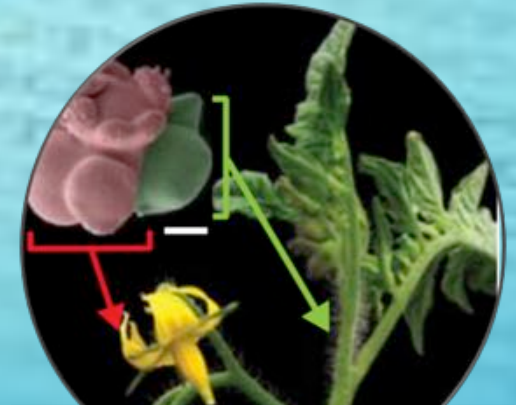
Human Genetics



Sequencing



Modeling



Plant Genomics

# Outline



1. Milestones in genomics
2. 21<sup>st</sup> Century Mega-Genomics
  1. Quantitative Traits and Measurements
  2. Parallel & Cloud Computing
3. Hadoop Applications for Genomics
  1. Kmer counting
  2. Mapping & Jnomics
  3. Assembly & Contrail

# Milestones in Genomics



Observations of 29,000 pea plants and 7 traits

Generation				in Verhältniss gestellt:		
	<i>A</i>	<i>Aa</i>	<i>a</i>	<i>A</i>	<i>Aa</i>	<i>a</i>
1	1	2	1	1	2	1
2	6	4	6	3	2	3
3	28	8	28	7	2	7
4	120	16	120	15	2	15
5	496	32	496	31	2	31
<i>n</i>				$2^n - 1$	2	$2^n - 1$

Seed		Flower	Pod		Stem	
Form	Cotyledons	Color	Form	Color	Place	Size
Grey & Round	Yellow	White	Full	Yellow	Axial pods, Flowers along	Long (6-7ft)
White & Wrinkled	Green	Violet	Constricted	Green	Terminal pods, Flowers top	Short (1-1ft)
1	2	3	4	5	6	7

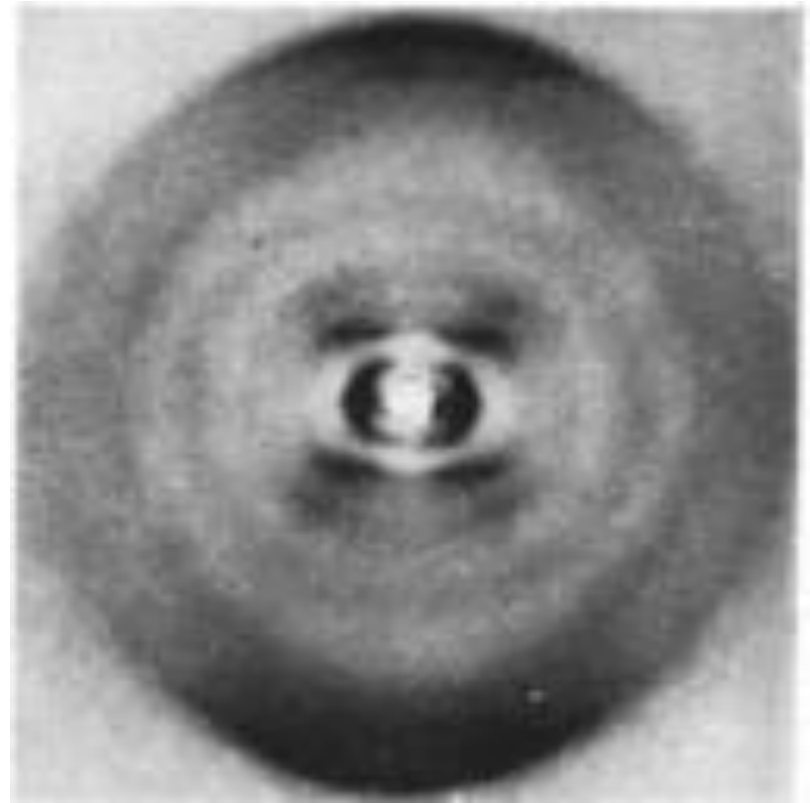
[http://en.wikipedia.org/wiki/Experiments\\_on\\_Plant\\_Hybridization](http://en.wikipedia.org/wiki/Experiments_on_Plant_Hybridization)

**Versuche über Pflanzen-Hybriden. Verh. Naturforsch (Experiments in Plant Hybridization)**  
Mendel, G. (1866). Ver. Brünn 4: 3–47 (in English in 1901, J. R. Hortic. Soc. 26: 1–32).

# Milestones in Genomics

## ***The origin and behavior of mutable loci in maize***

McClintock, B (1950) *Proceedings of the National Academy of Sciences*. 36:344–55.



## ***Molecular structure of nucleic acids; a structure for deoxyribose nucleic acid***

Watson JD, Crick FH (1953). *Nature* 171: 737–738.

# Milestones in Genomics: Zeroth Generation Sequencing

Nature Vol. 265 February 24 1977 687

---

## articles

---

### Nucleotide sequence of bacteriophage $\Phi$ X174 DNA

F. Sanger, G. M. Air\*, B. G. Barrell, N. L. Brown\*, A. R. Coulson, J. C. Fiddes, C. A. Hutchison III\*, P. M. Slocombe\* & M. Smith\*

MRC Laboratory of Molecular Biology, Hills Road, Cambridge CB2 2QH, UK

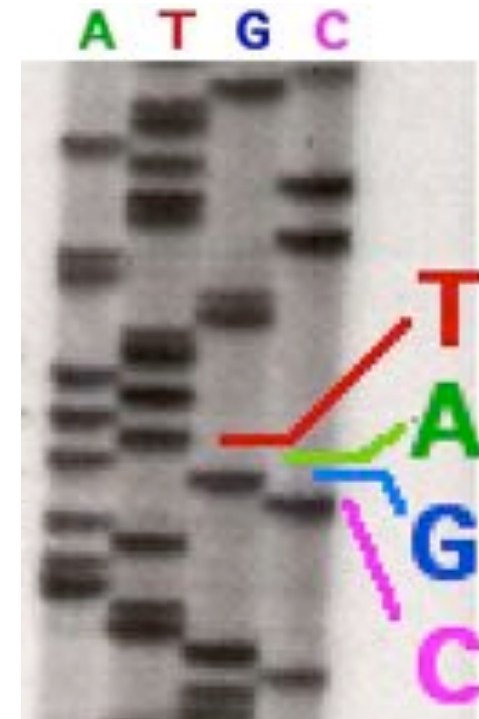
*A DNA sequence for the genome of bacteriophage  $\Phi$ X174 of approximately 5,375 nucleotides has been determined using the rapid and simple 'plus and minus' method. The sequence identifies many of the features responsible for the production of the proteins of the nine known genes of the organism, including initiation and termination sites for the proteins and RNAs. Two pairs of genes are coded by the same region of DNA using different reading frames.*

The genome of bacteriophage  $\Phi$ X174 is a single-stranded, circular DNA of approximately 5,400 nucleotides coding for nine known proteins. The order of these genes, as determined by genetic techniques<sup>1-3</sup>, is A-B-C-D-E-F-F'-G-H. Genes F, G and H code for structural proteins of the virus capsid, and gene J (as defined by sequence work) codes for a small basic protein

strand DNA of  $\Phi$ X has the same sequence as the mRNA and, in certain conditions, will bind ribosomes so that a protected fragment can be isolated and sequenced. Only one major site was found. By comparison with the amino acid sequence data it was found that this ribosome binding site sequence coded for the initiation of the gene G protein<sup>4</sup> (positions 2,362-2,413).

At this stage sequencing techniques using primed synthesis with DNA polymerase were being developed<sup>5,6</sup> and Schott<sup>7</sup> synthesised a decanucleotide with a sequence complementary to part of the ribosome binding site. This was used to prime into the intercistronic region between the F and G genes, using DNA polymerase and <sup>32</sup>P-labelled triphosphates<sup>8</sup>. The ribo-substitution technique<sup>9</sup> facilitated the sequence determination of the labelled DNA produced. This decanucleotide-primed system was also used to develop the plus and minus method<sup>10</sup>. Suitable synthetic primers are, however, difficult to prepare and an

**1977**  
**1<sup>st</sup> Complete Organism**  
**Bacteriophage  $\phi$ X174**  
**5375 bp**



Radioactive Chain Termination  
5000bp / week / person

<http://en.wikipedia.org/wiki/File:Sequencing.jpg>  
<http://www.answers.com/topic/automated-sequencer>

**Nucleotide sequence of bacteriophage  $\phi$ X174 DNA**  
Sanger, F. et al. (1977) *Nature*. 265: 687 - 695

# Milestones in Genomics: First Generation Sequencing



**1995**

Fleischmann *et al.*  
1<sup>st</sup> Free Living Organism  
TIGR Assembler. 1.8Mbp



**2000**

Myers *et al.*  
1<sup>st</sup> Large WGS Assembly.  
Celera Assembler. 116 Mbp



**2001**

Venter *et al.* / IHGSC  
Human Genome  
Celera Assembler. 2.9 Gbp

ABI 3700: 500 bp reads x 768 samples / day = 384,000 bp / day.

"The machine was so revolutionary that it could decode in a single day the same amount of genetic material that most DNA labs could produce in a year." J. Craig Venter

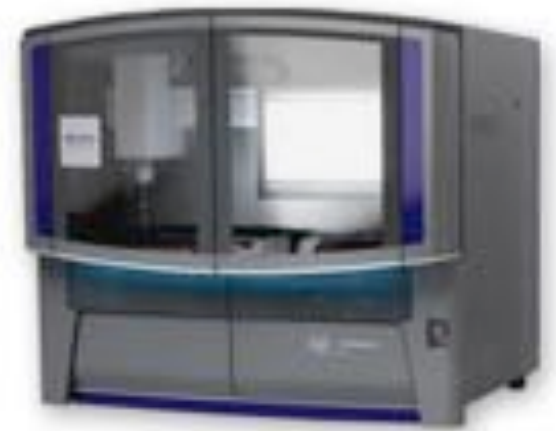
# Milestones in Genomics: Second Generation Sequencing



**2004**  
454/Roche  
*Pyrosequencing*  
Current Specs (Titanium):  
1M 400bp reads / run =  
1 Gbp / day



**2007**  
Illumina  
*Sequencing by Synthesis*  
Current Specs (HiSeq 2000):  
2.5B 100bp reads / run =  
60Gbp / day



**2008**  
ABI / Life Technologies  
*SOLiD Sequencing*  
Current Specs (5500xl):  
5B 75bp reads / run =  
30Gbp / day



# Milestones in Genomics: Third Generation Sequencing



**2010**

Ion Torrent

*Postlight Sequencing*

Current Specs (Ion 318):

11M 300bp reads / run =

>1Gbp / day



**2011**

Pacific Biosciences

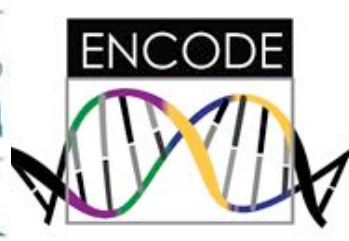
*SMRT Sequencing*

Current Specs (RS):

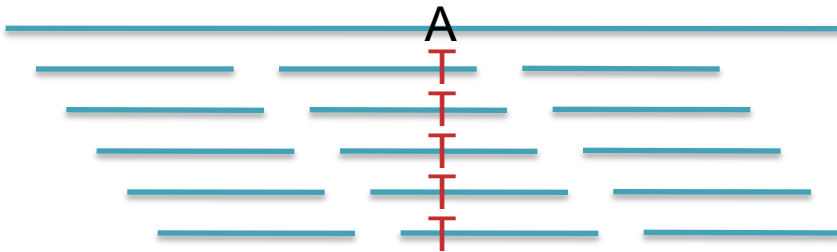
50k 2kbp reads / run =

>200Mbp / day

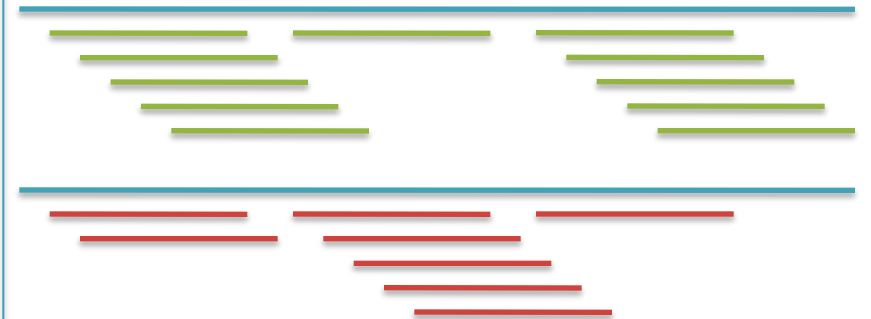
# Milestones in Genomics



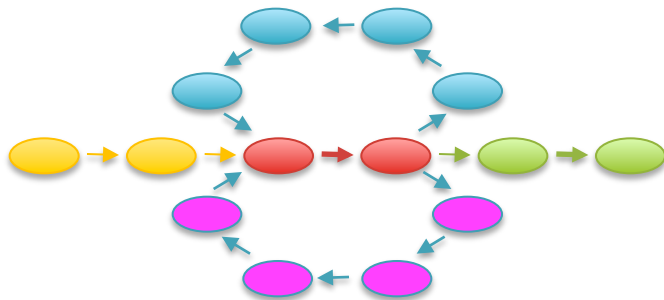
## Alignment & Variations



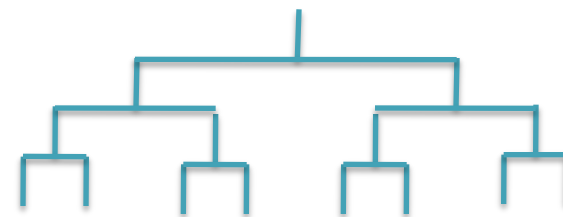
## Differential Analysis



## De novo Assembly



## Phylogeny & Modeling



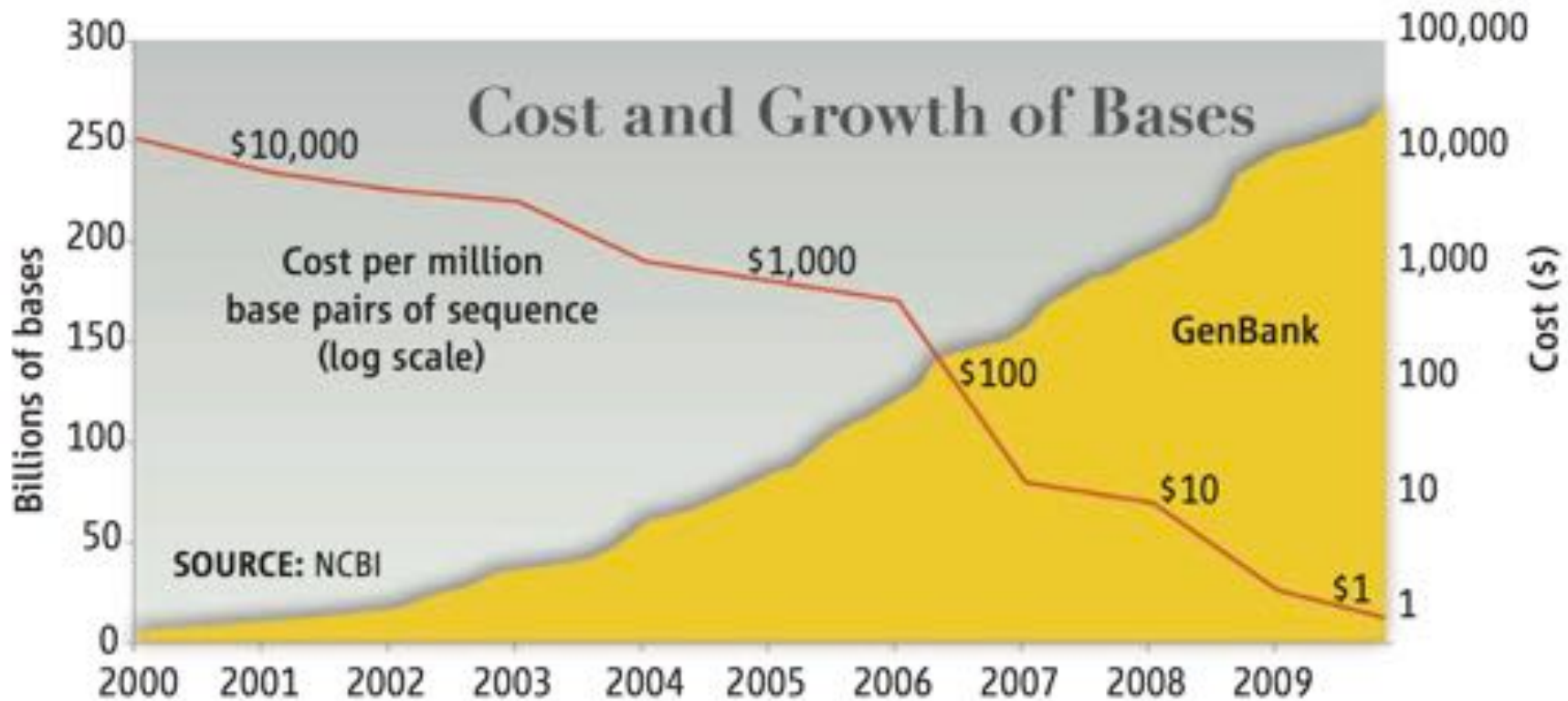
# Sequencing Centers



***Next Generation Genomics: World Map of High-throughput Sequencers***  
<http://pathogenomics.bham.ac.uk/hts/>

# DNA Data Tsunami

*Current world-wide sequencing capacity exceeds 13Pbp/year  
and is growing at 5x per year!*



**"Will Computers Crash Genomics?"**

Elizabeth Pennisi (2011) *Science*. 331(6018): 666-668.

# Mega-Genomics Challenges



The foundations of genomics will continue to be *observation, experimentation, and interpretation*

- Technology will continue to push the frontier
- Measurements will be made *digitally* over large populations, at extremely high resolution, and for diverse applications

## Rise in Quantitative Demands

1. *Experimental design*: selection, collection, tracking & metadata
  - Ontologies, LIMS, sample databases
2. *Observation*: measurement, storage, transfer, computation
  - Algorithms to overcome sensor errors & limitations, computing at scale
3. *Integration*: multiple samples, multiple assays, multiple analyses
  - Reproducible workflows, common formats, resource federation
4. *Discovery*: visualizing, interpreting, modeling
  - Clustering, data reduction, trend analysis

# Hadoop MapReduce

<http://hadoop.apache.org>

- MapReduce is Google's framework for large data computations
  - Data and computations are spread over thousands of computers
    - Indexing the Internet, PageRank, Machine Learning, etc... (Dean and Ghemawat, 2004)
    - 946PB processed in May 2010 (Jeff Dean at Stanford, 11.10.2010)
  - Hadoop is the leading open source implementation
    - Developed and used by Yahoo, Facebook, Twitter, Amazon, etc
    - GATK is an alternative implementation specifically for NGS
- Benefits
  - Scalable, Efficient, Reliable
  - Easy to Program
  - Runs on commodity computers
- Challenges
  - Redesigning / Retooling applications
    - Not Condor, Not MPI
    - Everything in MapReduce



# Hadoop for NGS Analysis



## CloudBurst

Highly Sensitive Short Read Mapping with MapReduce

*100x speedup mapping on 96 cores @ Amazon*

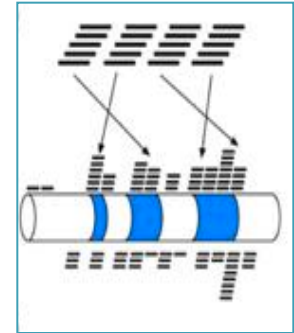
<http://cloudburst-bio.sf.net>

(Schatz, 2009)

## Myrna

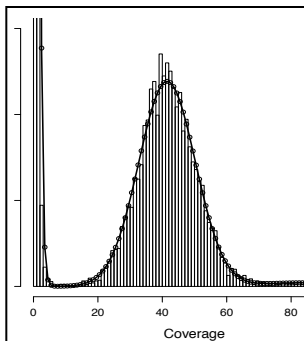
Cloud-scale differential gene expression for RNA-seq

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



(Langmead, Hansen, Leek, 2010)

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



## 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)

## Genome Indexing

Rapid Parallel Construction of Genome Index

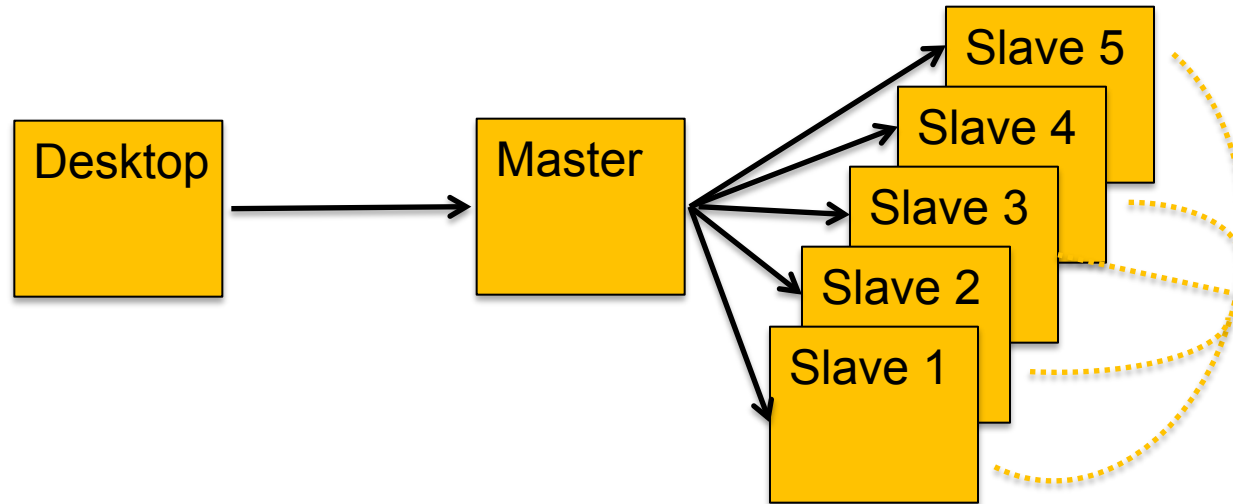
*Construct the BWT of the human genome in 9 minutes*

```
$GATTACA  
A$GATTAC  
ACA$GATT  
ATTACA$G  
CA$GATTA  
GATTACA£  
TACA$GAT  
TTACA$GA
```

(Menon, Bhat, Schatz, 2011\*)

<http://code.google.com/p/genome-indexing/>

# System 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



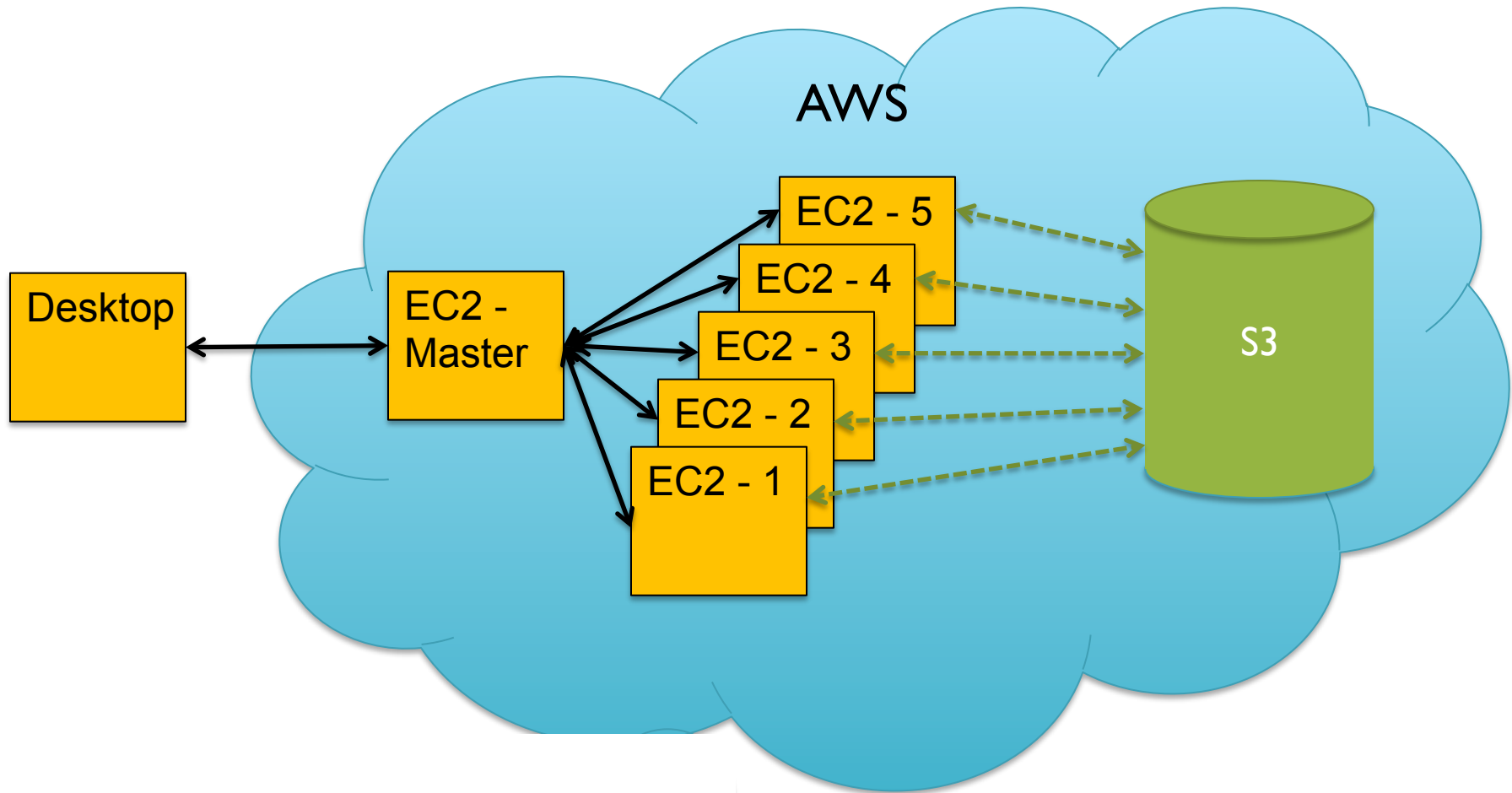
# Amazon Web Services

<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
- Simple Storage Service (S3)
  - Scalable data storage
- Plus many, many more



# Hadoop on AWS



- If you don't have 1000s of machines, rent them from Amazon
  - After machines pool up, ssh to master as if it was a local machine.
  - Use S3 for persistent data storage, with very fast interconnect to EC2.

# Parallel Algorithm Spectrum

## Embarrassingly Parallel



Map-only  
Each item is Independent

## Loosely Coupled



MapReduce  
Independent-Sync-Independent

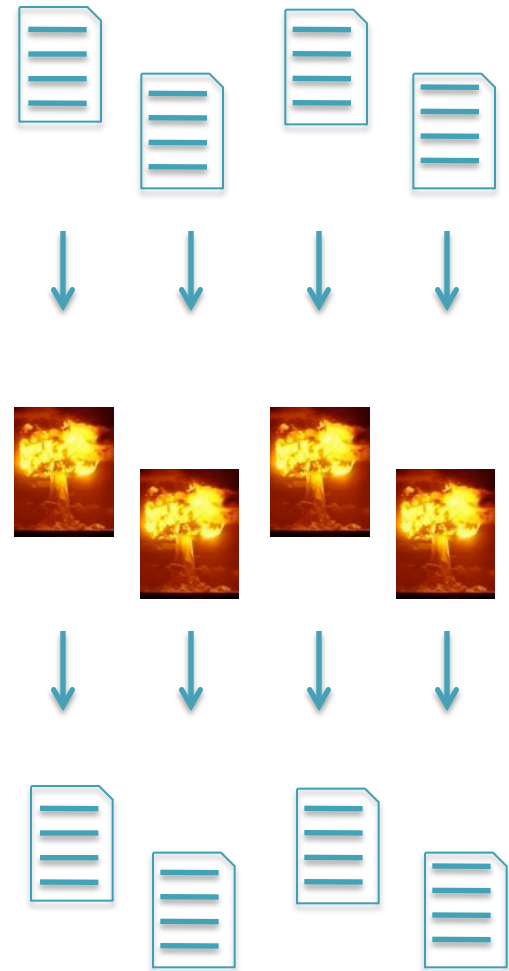
## Tightly Coupled



Iterative MapReduce  
Constant Sync

# I. 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

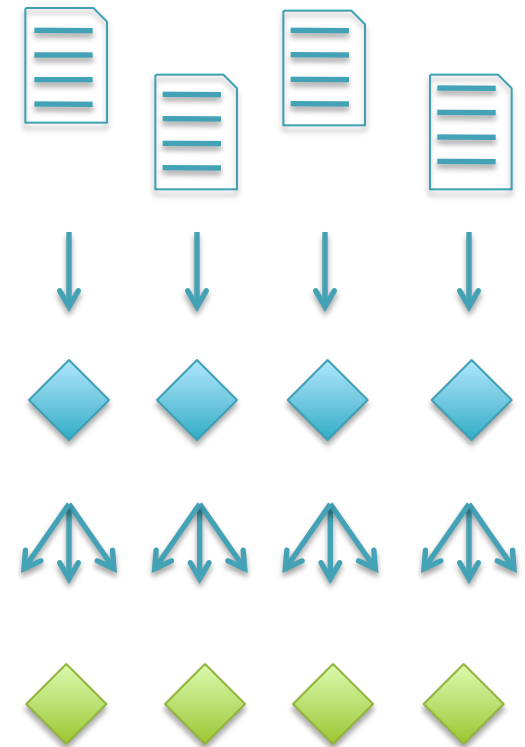


# Elementary School Dance



## 2. Loosely Coupled

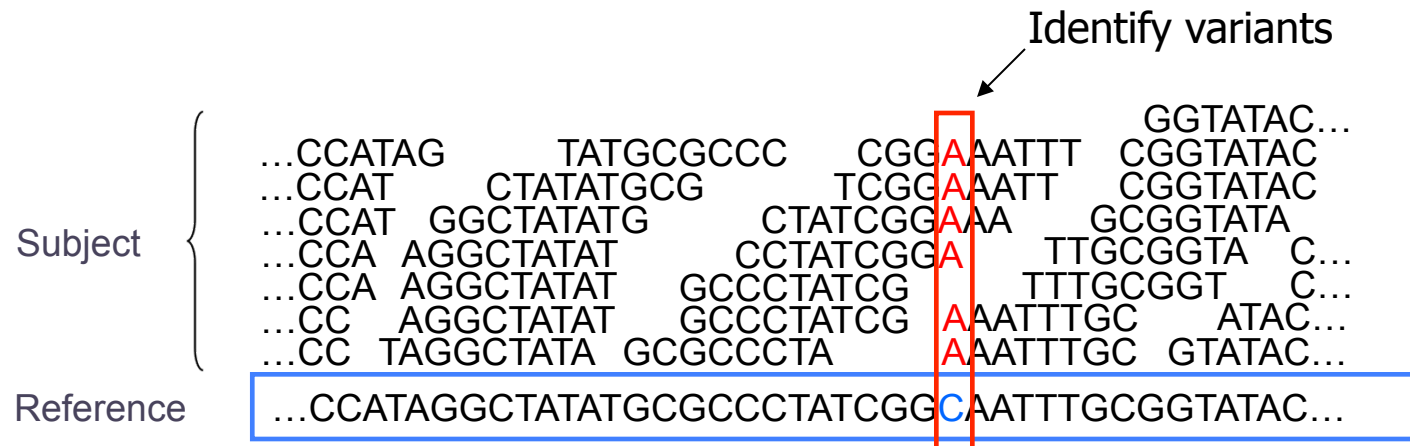
- Divide and conquer
  - Independently process many items
  - Group partial results
  - Scan partial results into final answer
- Challenges
  - Batch computing challenges
  - + Shuffling of huge datasets
- Technologies
  - Hadoop, Elastic MapReduce, Dryad
  - Parallel Databases



# Junior High Dance



# Short Read Mapping



- 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

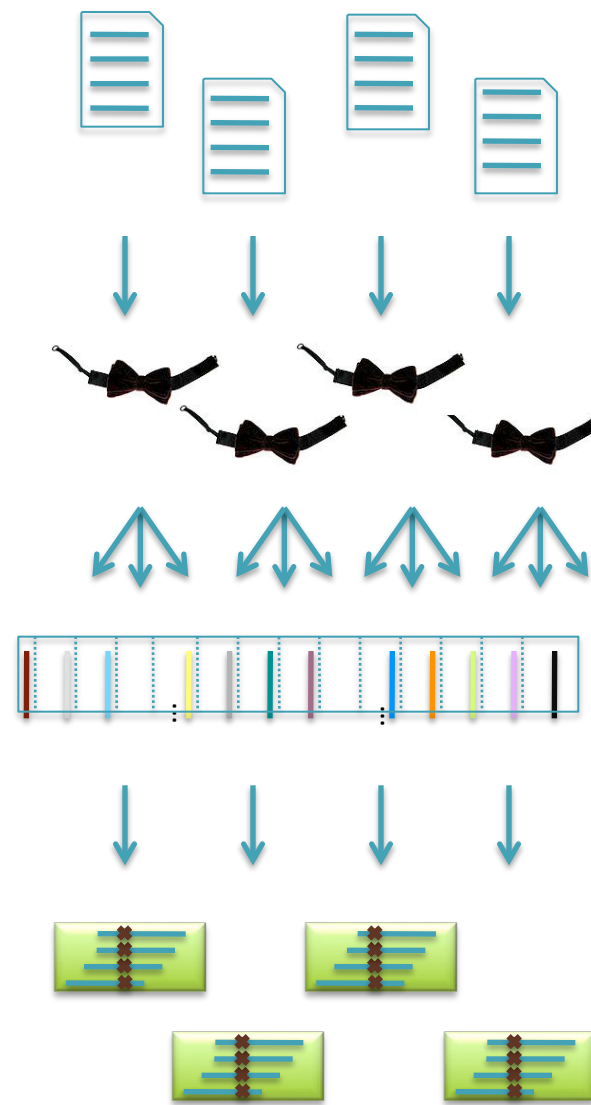




# Crossbow

<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>

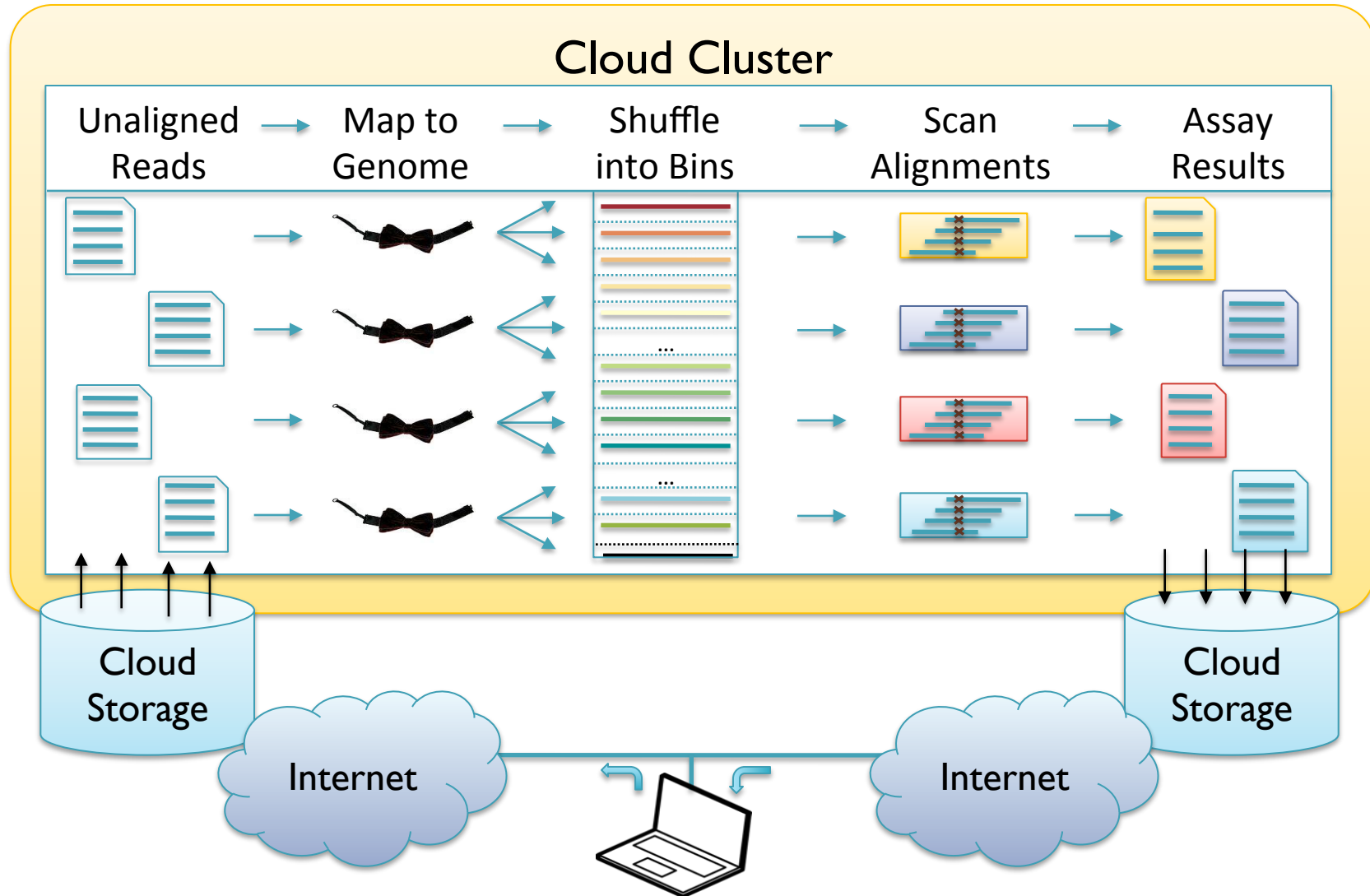
	Asian Individual Genome		
<b>Data Loading</b>	3.3 B reads	106.5 GB	\$10.65
<b>Data Transfer</b>	1h :15m	40 cores	\$3.40
<b>Setup</b>	0h : 15m	320 cores	\$13.94
<b>Alignment</b>	1h : 30m	320 cores	\$41.82
<b>Variant Calling</b>	1h : 00m	320 cores	\$27.88
<b>End-to-end</b>	4h : 00m		\$97.69

Discovered 3.7M SNPs in one 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

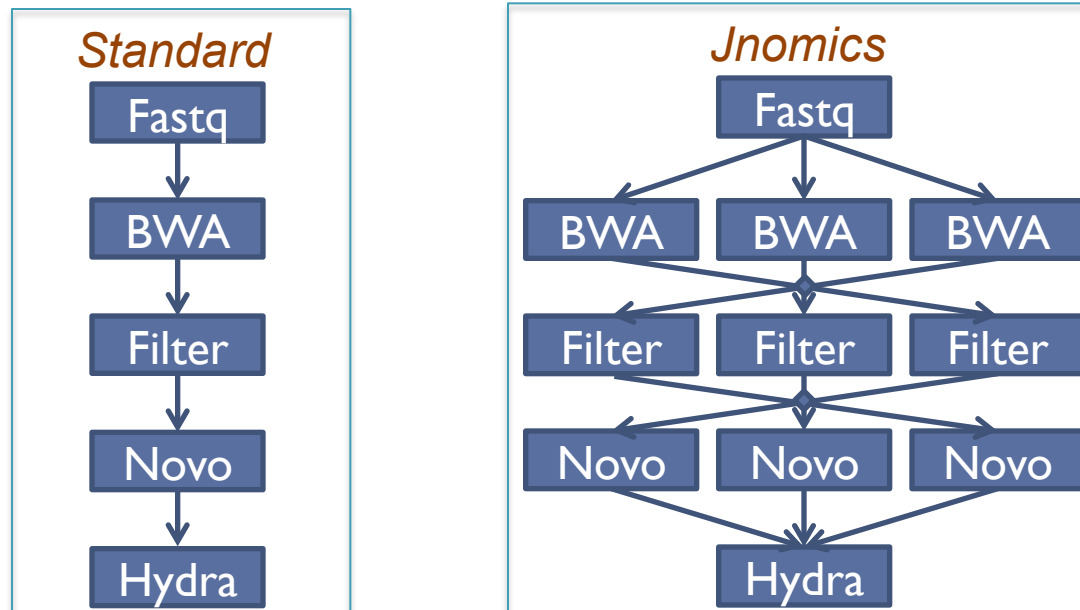


**Cloud Computing and the DNA Data Race.**

Schatz, MC, Langmead B, Salzberg SL (2010) *Nature Biotechnology*. **28**:691-693

# Jnomics: Cloud-scale genomics

Matt Titmus, James Gurtowski, Michael Schatz



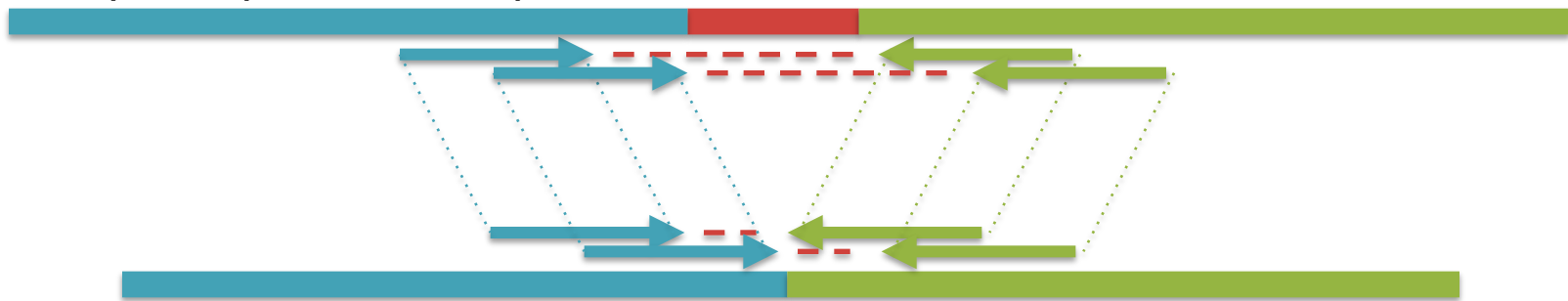
- Rapid parallel execution of NGS analysis pipelines
  - FASTX, BWA, Bowtie, Novoalign, SAMTools, Hydra
  - Sorting, merging, filtering, selection, of BAM, SAM, BED, fastq
  - Population analysis: Clustering, GWAS, Trait Inference

**Answering the demands of digital genomics**

Titmus, M.A., Schatz, M.C.. (2011) *Under Review*

# Jnomics Structural Variations

Sample Separation: 2kbp



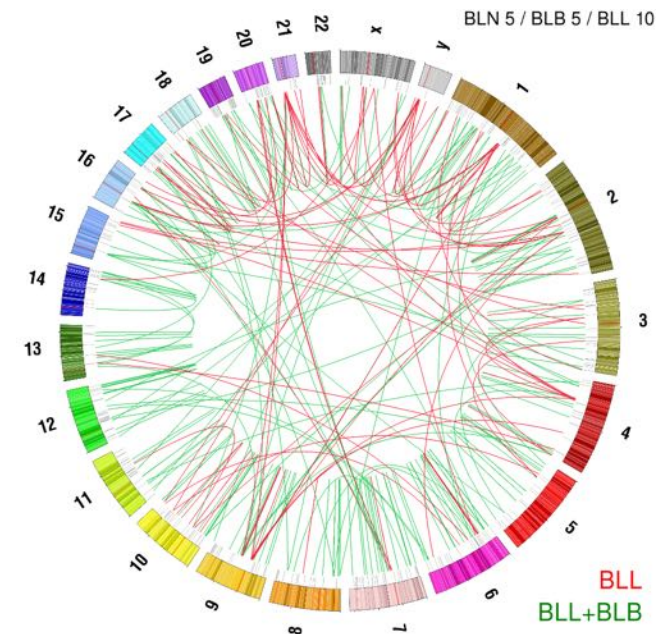
Mapped Separation: 1kbp

## Discordant Pair Analysis

- Identify clusters of pairs too close or too far away indicating a SV

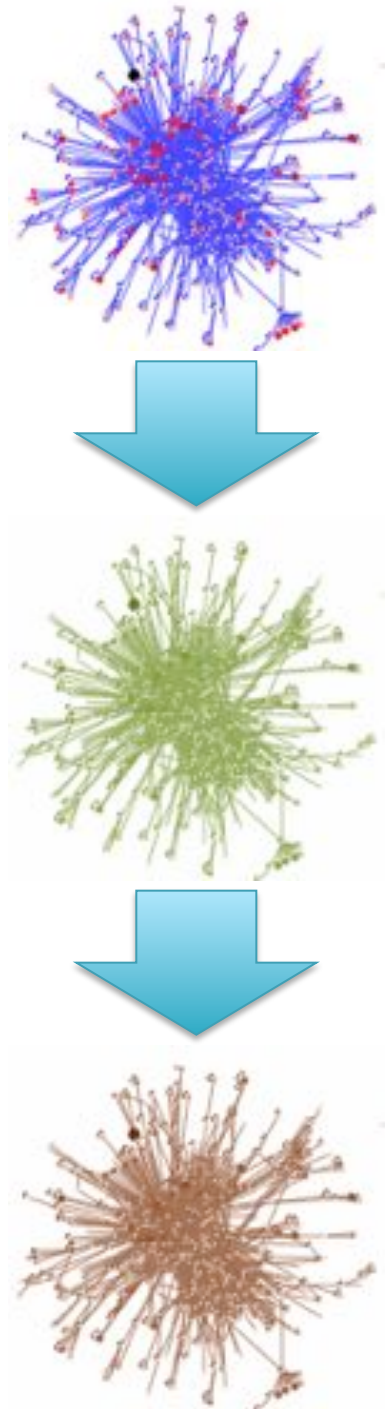
## Circos plot of high confidence SVs specific to esophageal cancer sample

- Red: SVs specific to tumor
- Green: SVs in both diseased and tumor samples



# 3. 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



# High School Dance

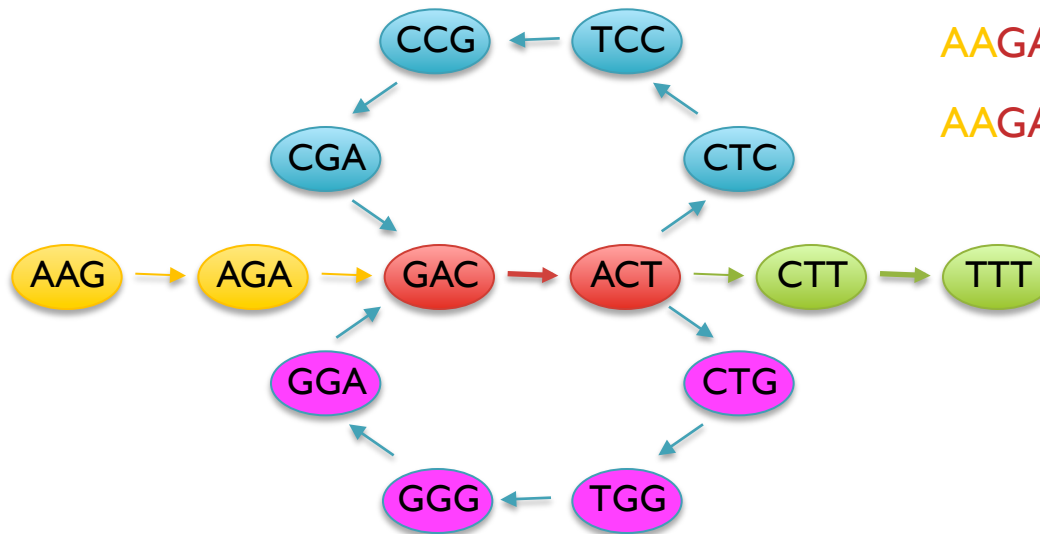


# Short Read Assembly

## Reads

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

## de Bruijn Graph



## Potential Genomes

AAGACTCCGACTGGGACTTT

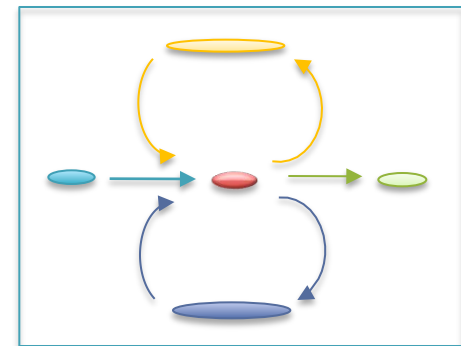
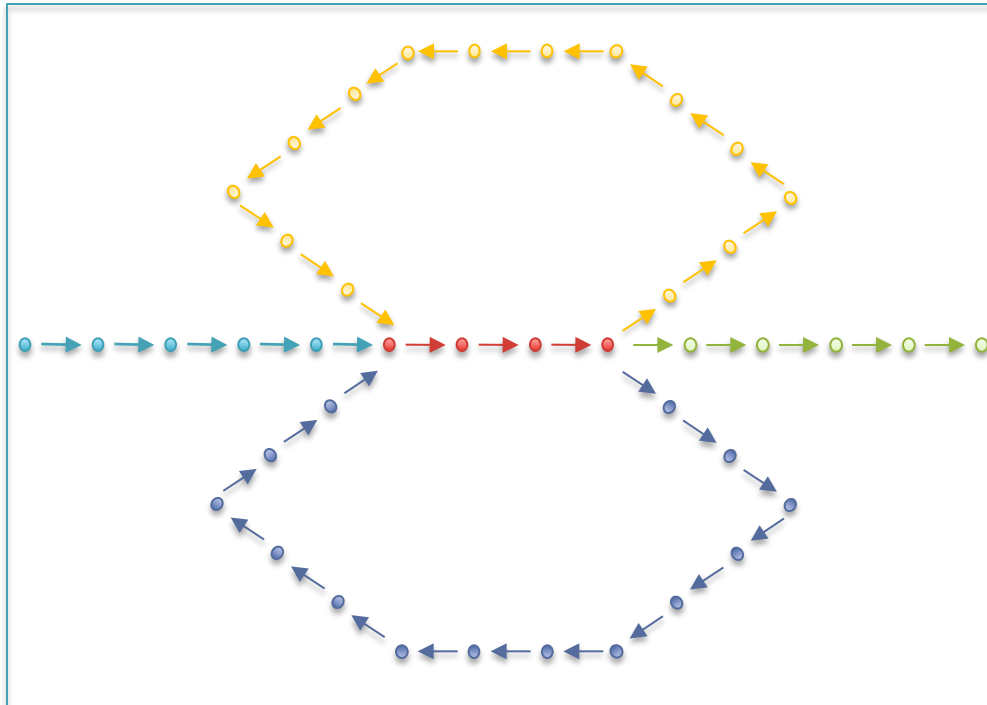
AAGACTGGGACTCCGACTTT

- 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



# Graph Compression

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

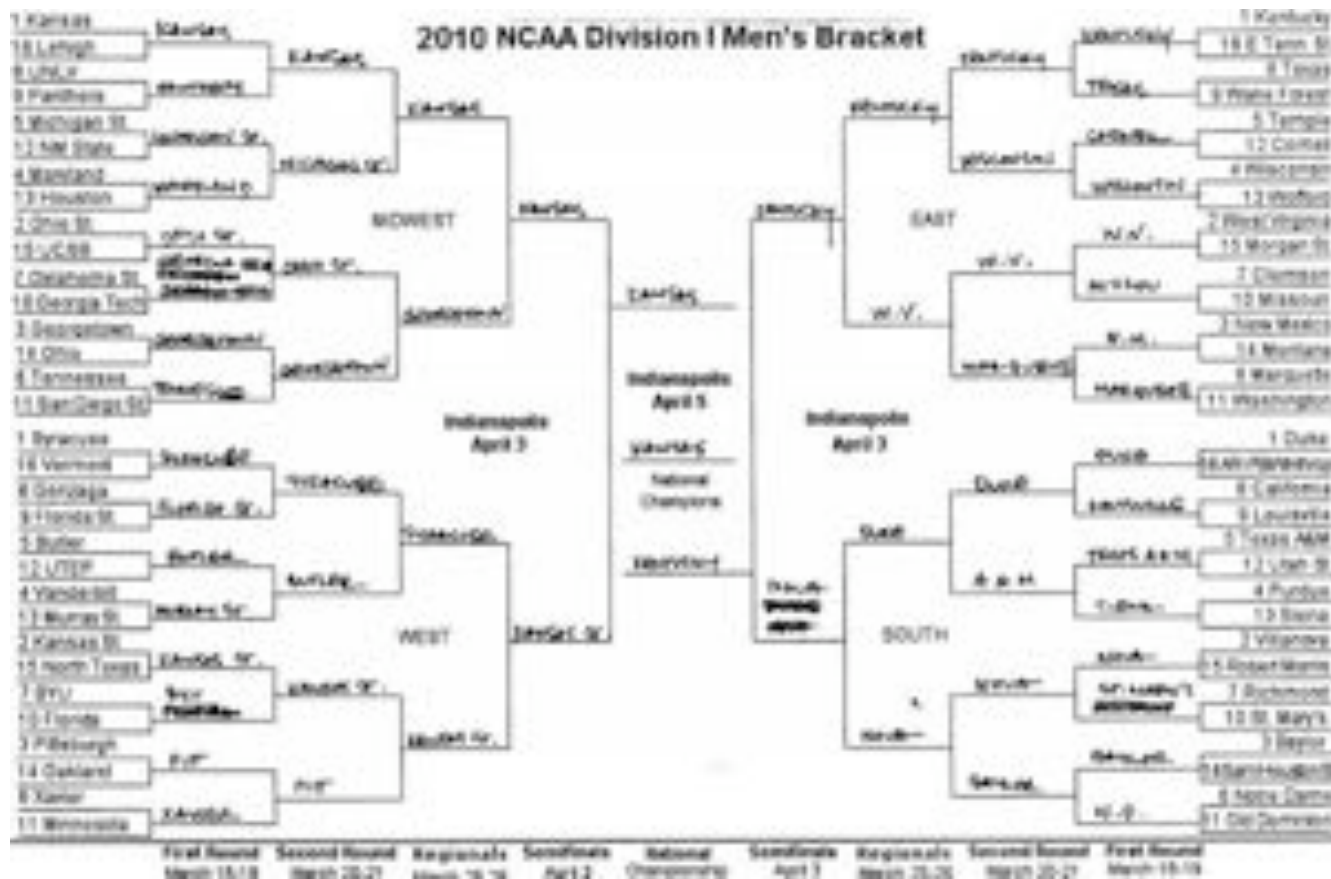


**Design Patterns for Efficient Graph Algorithms in MapReduce.**

*Lin, J., Schatz, M.C. (2010) Workshop on Mining and Learning with Graphs Workshop (KDD-2010)*

# Warmup Exercise

- Who here was born closest to Feb 7?
  - You can only compare to 1 other person at a time



Find winner among 64 teams in just 6 rounds

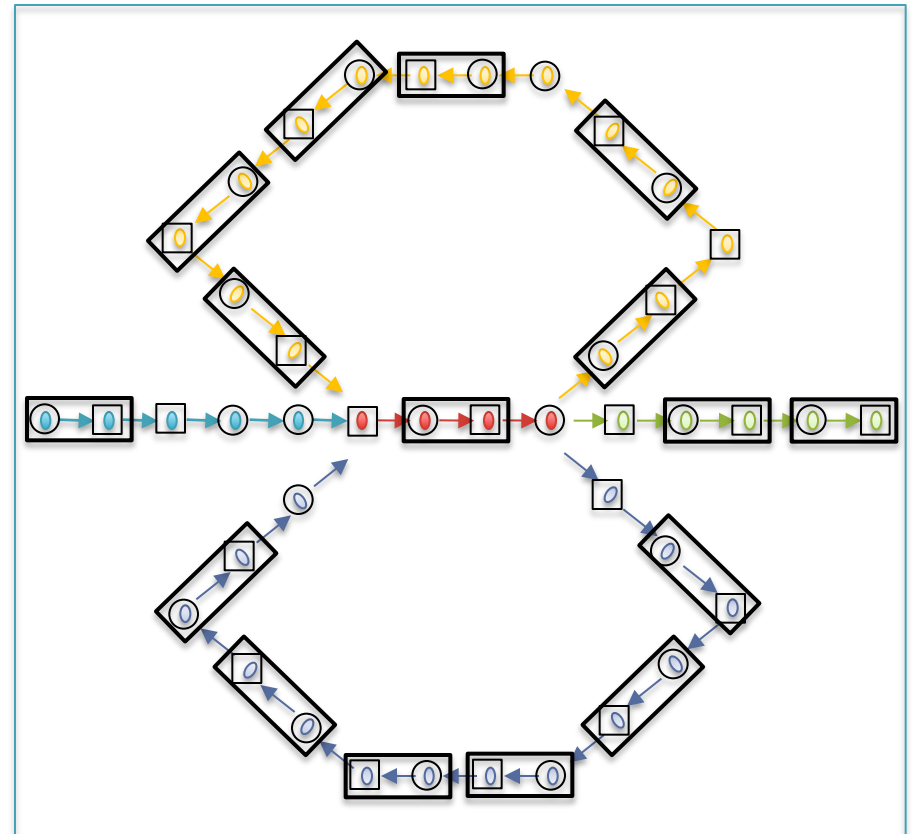
# Fast Path Compression

## Challenges

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

## Randomized List Ranking

- Randomly assign  $\textcircled{H}$  /  $\square T$  to each compressible node
- Compress  $\textcircled{H} \rightarrow \square T$  links



Initial Graph: 42 nodes

## Randomized Speed-ups in Parallel Computation.

Vishkin U. (1984) *ACM Symposium on Theory of Computation*. 230-239.

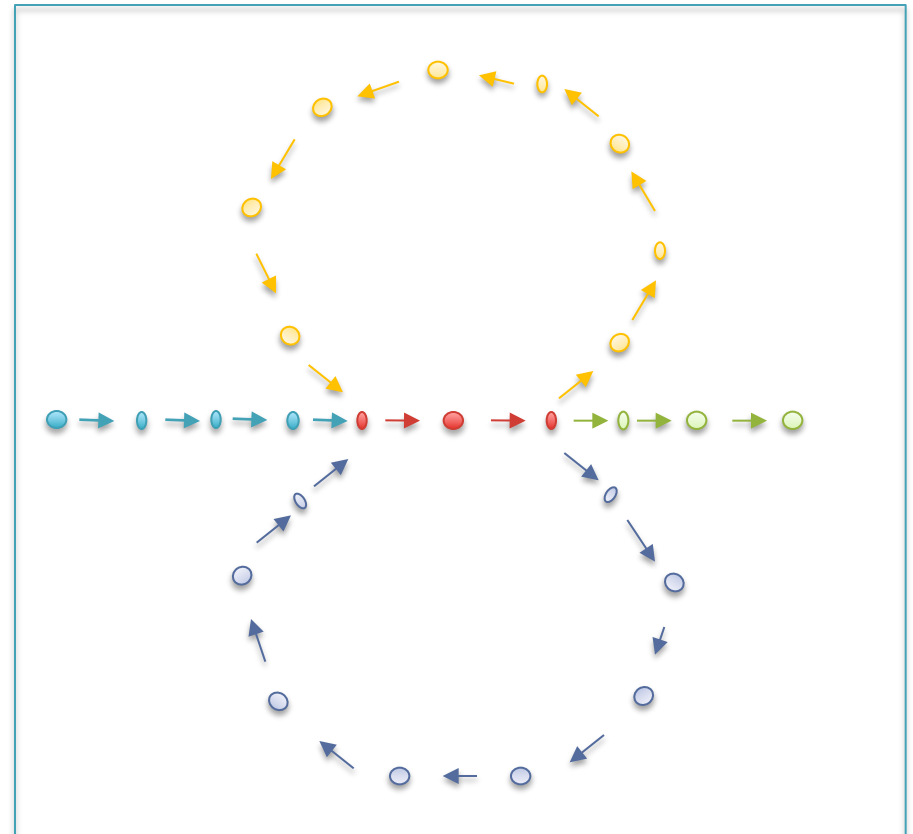
# Fast Path Compression

## Challenges

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

## Randomized List Ranking

- Randomly assign  $\textcircled{\text{H}}$  /  $\boxed{\text{T}}$  to each compressible node
- Compress  $\textcircled{\text{H}} \rightarrow \boxed{\text{T}}$  links



Round 1: 26 nodes (38% savings)

## Randomized Speed-ups in Parallel Computation.

Vishkin U. (1984) *ACM Symposium on Theory of Computation*. 230-239.

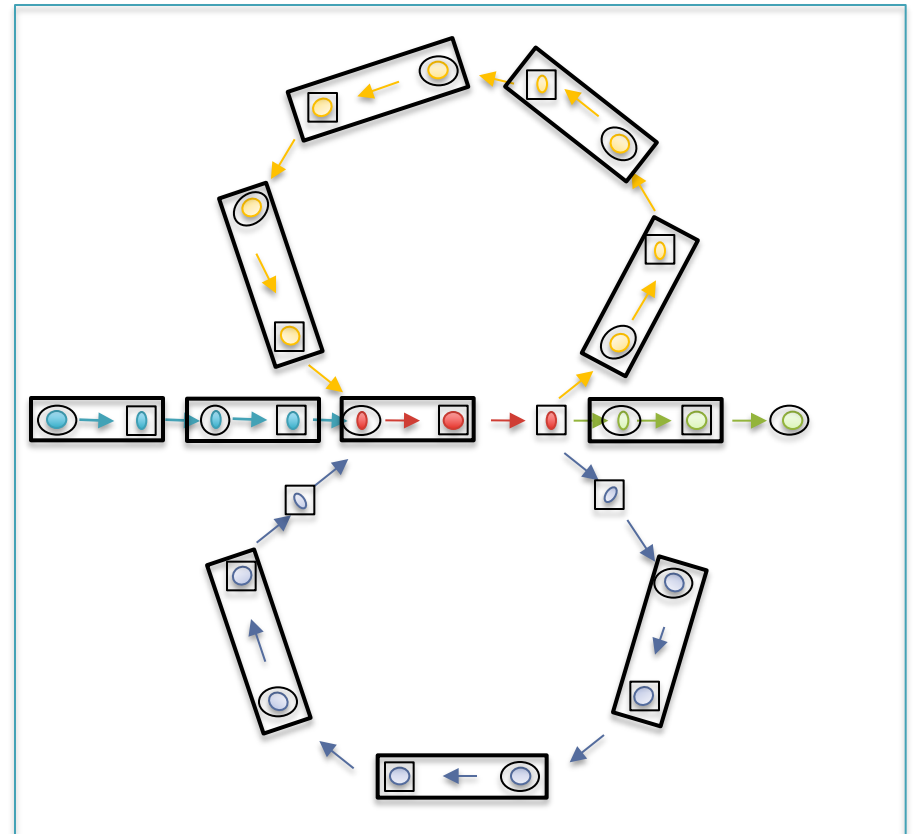
# Fast Path Compression

## Challenges

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

## Randomized List Ranking

- Randomly assign  $\textcircled{H}$  /  $\square T$  to each compressible node
- Compress  $\textcircled{H} \rightarrow \square T$  links



Round 1: 26 nodes (38% savings)

## Randomized Speed-ups in Parallel Computation.

Vishkin U. (1984) *ACM Symposium on Theory of Computation*. 230-239.

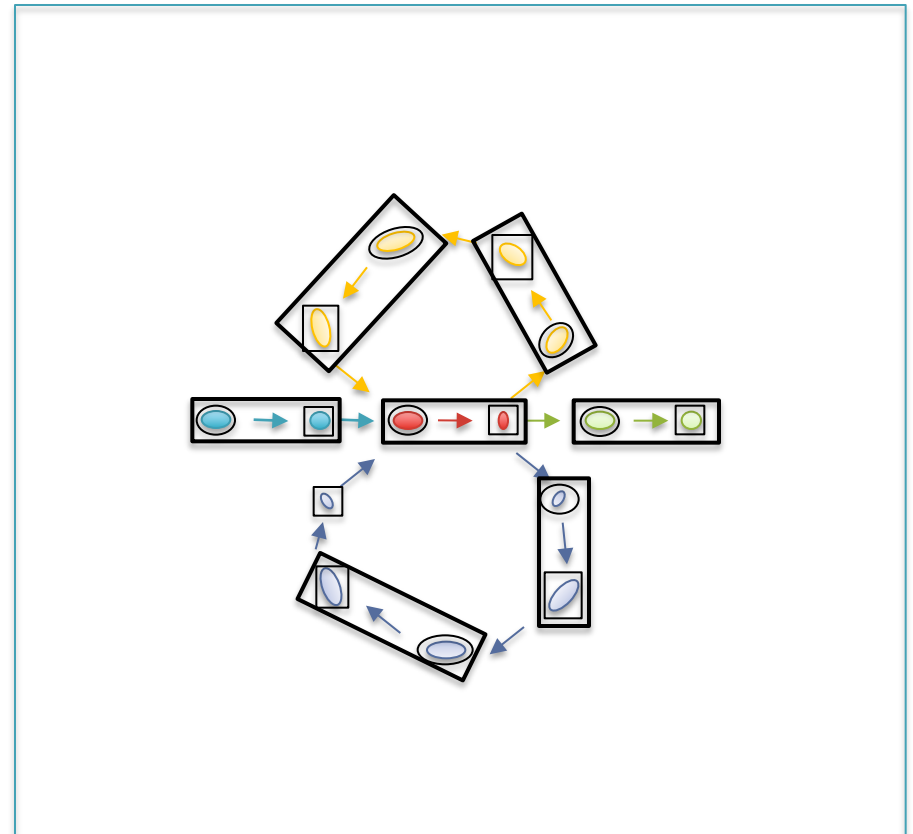
# Fast Path Compression

## Challenges

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

## Randomized List Ranking

- Randomly assign  $\textcircled{\text{H}}$  /  $\boxed{\text{T}}$  to each compressible node
- Compress  $\textcircled{\text{H}} \rightarrow \boxed{\text{T}}$  links



Round 2: 15 nodes (64% savings)

## Randomized Speed-ups in Parallel Computation.

Vishkin U. (1984) *ACM Symposium on Theory of Computation*. 230-239.

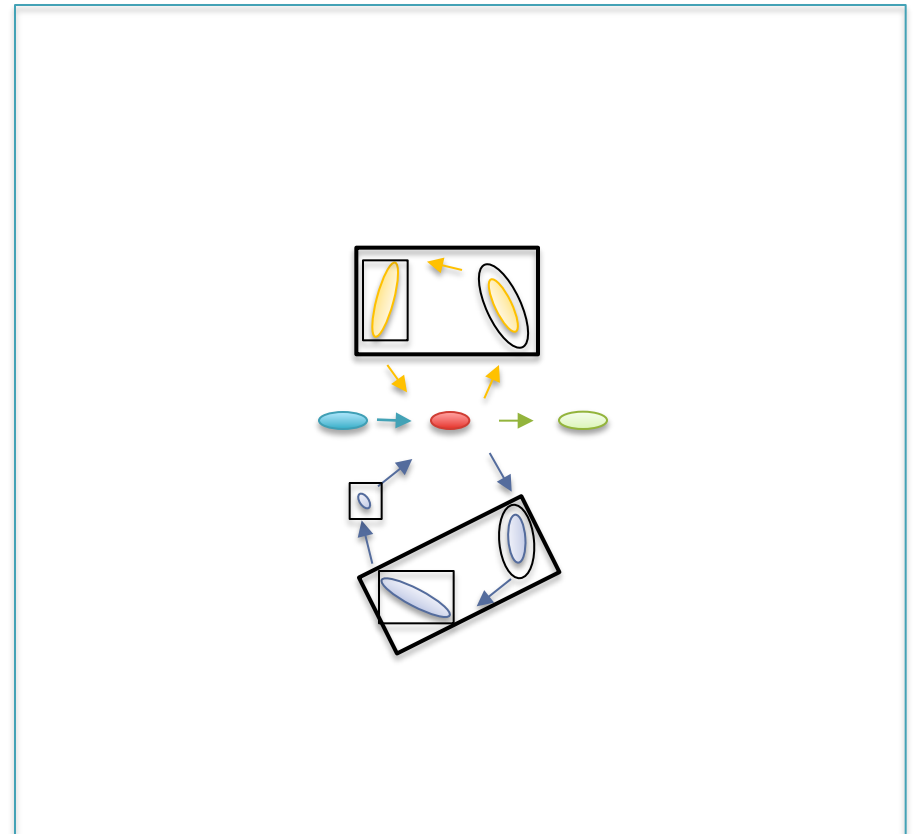
# Fast Path Compression

## Challenges

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

## Randomized List Ranking

- Randomly assign  $\textcircled{\text{H}}$  /  $\boxed{\text{T}}$  to each compressible node
- Compress  $\textcircled{\text{H}} \rightarrow \boxed{\text{T}}$  links



Round 2: 8 nodes (81% savings)

## Randomized Speed-ups in Parallel Computation.

Vishkin U. (1984) *ACM Symposium on Theory of Computation*. 230-239.

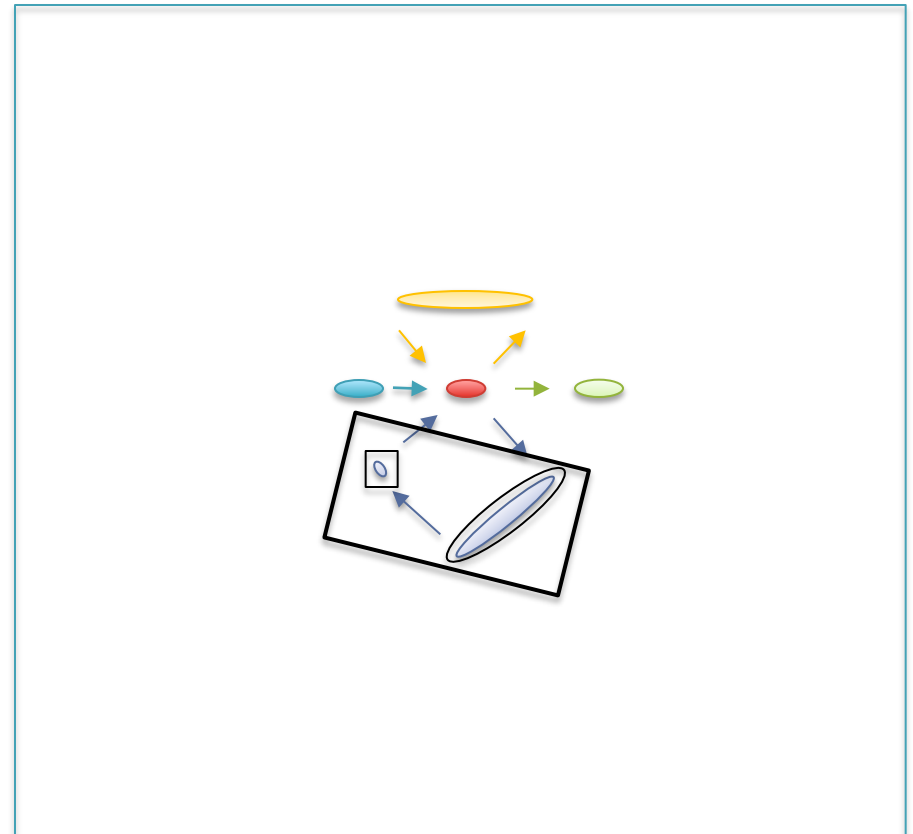
# Fast Path Compression

## Challenges

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

## Randomized List Ranking

- Randomly assign  $\textcircled{\text{H}}$  /  $\boxed{\text{T}}$  to each compressible node
- Compress  $\textcircled{\text{H}} \rightarrow \boxed{\text{T}}$  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

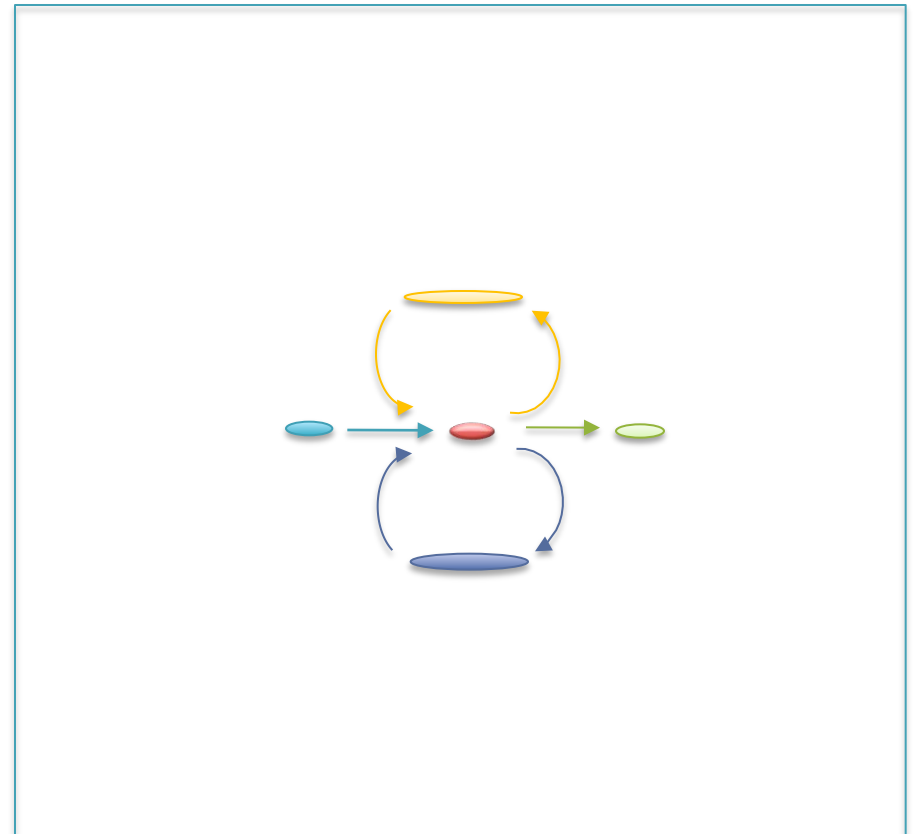
- Nodes stored on different computers
- Nodes can only access direct neighbors

## Randomized List Ranking

- Randomly assign  $\textcircled{\text{H}}$  /  $\boxed{\text{T}}$  to each compressible node
- Compress  $\textcircled{\text{H}} \rightarrow \boxed{\text{T}}$  links

## Performance

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

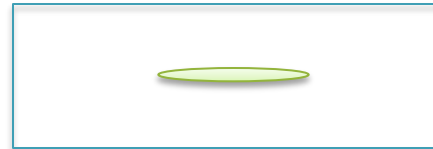


Round 4: 5 nodes (88% savings)

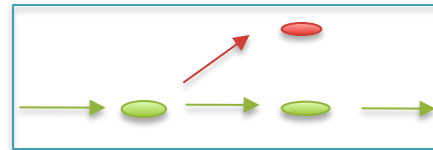
## Randomized Speed-ups in Parallel Computation.

Vishkin U. (1984) *ACM Symposium on Theory of Computation*. 230-239.

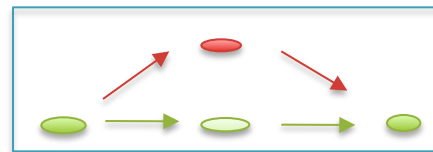
# Node Types



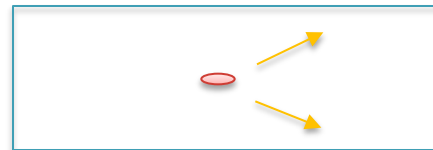
Isolated nodes (10%)



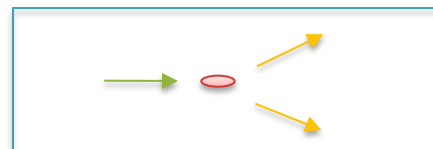
Tips (46%)



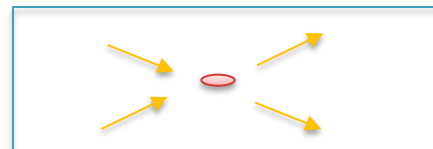
Bubbles/Non-branch (9%)



Dead Ends (.2%)



Half Branch (25%)



Full Branch (10%)

(Chaisson, 2009)

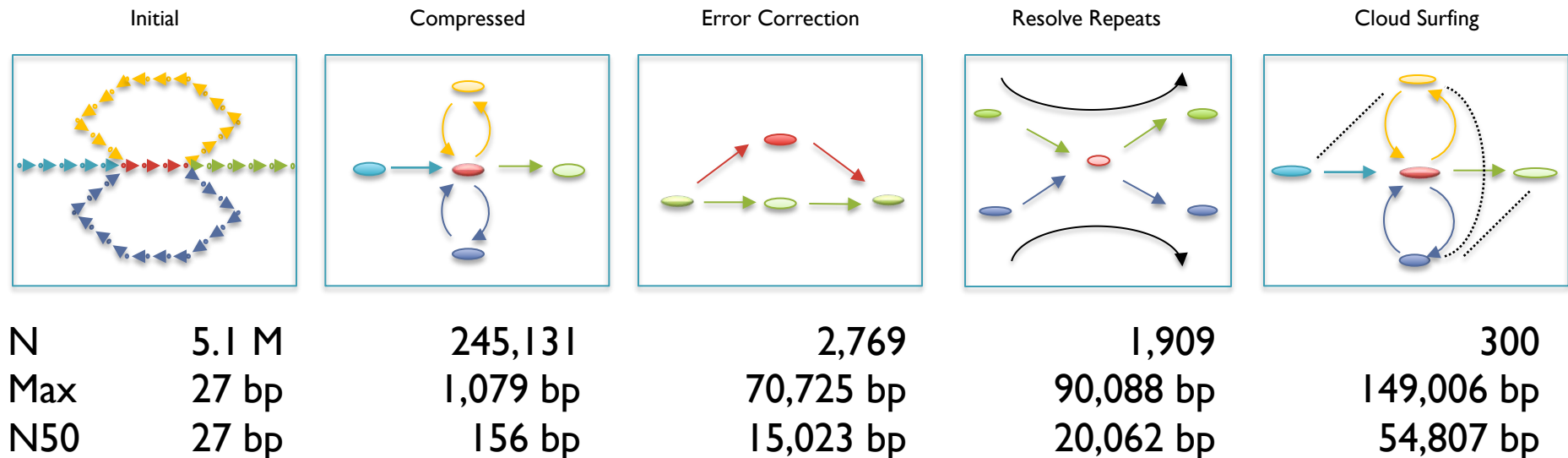
# Contrail

<http://contrail-bio.sourceforge.net>



## De novo bacterial assembly

- *Genome: E. coli* K12 MGI655, 4.6Mbp
- *Input: 20.8M* 36bp reads, 200bp insert (~150x coverage)
- *Preprocessor: Quake* Error Correction



## Assembly of Large Genomes with Cloud Computing.

Schatz MC, Sommer D, Kelley D, Pop M, et al. *In Preparation.*

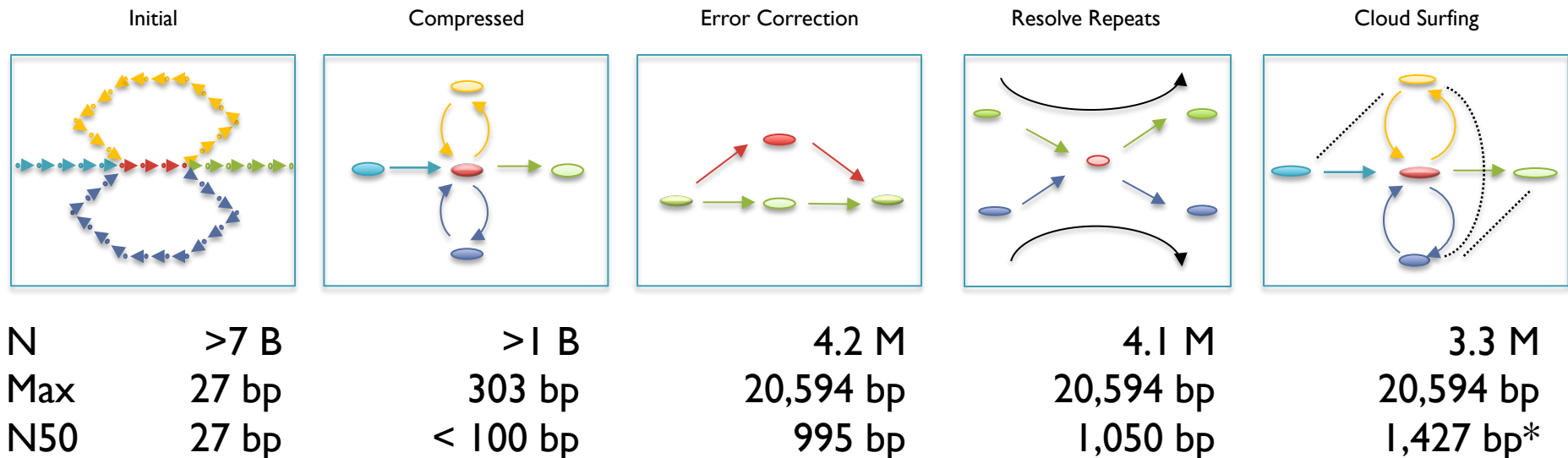
# Contrail

<http://contrail-bio.sourceforge.net>



## De novo Assembly of the Human Genome

- *Genome*: African male NAI8507 (SRA000271, Bentley *et al.*, 2008)
- *Input*: 3.5B 36bp reads, 210bp insert (~40x coverage)

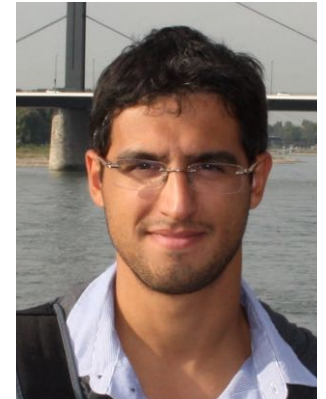


## Assembly of Large Genomes with Cloud Computing.

Schatz MC, Sommer D, Kelley D, Pop M, *et al.* *In Preparation.*

# Scalpel: Haplotype Microassembly

G. Narzisi, D. Levy, I. Iossifov, J. Kendall, M. Wigler, M. Schatz



- Use assembly techniques to identify complex variations from short reads
  - Improved power to find indels
  - Trace candidate haplotypes sequences as paths through assembly graphs



```
Ref:      ...CACAGGATCCACCTTTCTCAAAGACCCAGGATCCTCCTTCCTCGGTGACACTGTATACGTC...
Father:   ...CACAGGATCCACCTTTCTCAAAGACCCAGGATCCTCCTTCCTCGGTGACACTGTATACGTC... [cov:19.5]
Mother_1: ...CACAGGATCCACCTTTCTCAAAGACCCAGGATCCTCCTTCCTCGGTGACACTGTATACGTC... [cov:19.4]
Mother_2: ...CACAGGATCCACCTTT-----CTTGGTGACACTGTATACGTC... [cov:21.5]
Aut_2:    ...CACAGGATCCACCTTTCTCAAAGACCCAGGATCCTCCTTCCTCGGTGACACTGTATACGTC... [cov:28.2]
Aut_1:    ...CACAGGATCCACCTTT-----CTTGGTGACACTGTATACGTC... [cov:33.3]
Sib_1:    ...CACAGGATCCACCTTTCTCAAAGACCCAGGATCCTCCTTCCTCGGTGACACTGTATACGTC... [cov:19.4]
Sib_2:    ...CACAGGATCCACCTTT-----CTTGGTGACACTGTATACGTC... [cov:21.5]
```

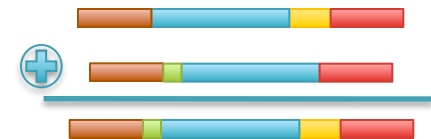
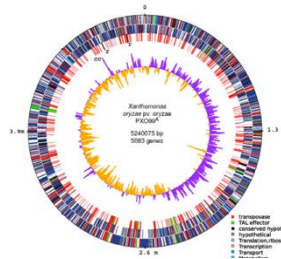
24 bp heterozygous indel at chr5:176026122 GPRIN1

# Summary

I'm focused on the intersection of the most significant biology, biotechnology, and compute technology

Computational research is the key to unlocking the potential of “mega-genomics”

- Explosion in quantitative traits and measurements
- Parallel systems essential for analyzing large data sets
- Algorithms and machine learning to squeeze insight out of diverse data types
- Collaborations and visual informatics systems with biologists to help execute experiments & interpret results

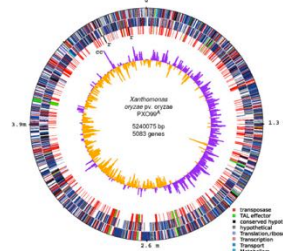


# Acknowledgements



Mitch Bekritsky  
Giuseppe Narzisi

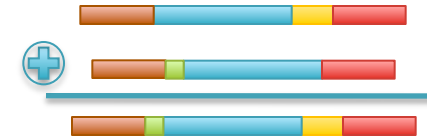
Ivan Iossifov  
Wigler Lab



Hayan Lee  
Matt Titmus  
James Gurtowski

Ware Lab  
McCombie Lab

Adam Phillippy (NBACC)  
Sergey Koren (NBACC)



Paul Baranay (CSHL/ND)

Scott Emrich (ND)  
Steven Salzberg (JHU)  
Mihai Pop (UMD)



# Thank You!

<http://schatzlab.cshl.edu>  
[@mike\\_schatz](#)

Sequencing & Assembly at 10:30a

Genome-wide Analysis at 1:30p

Break-out Discussions at 3:00p