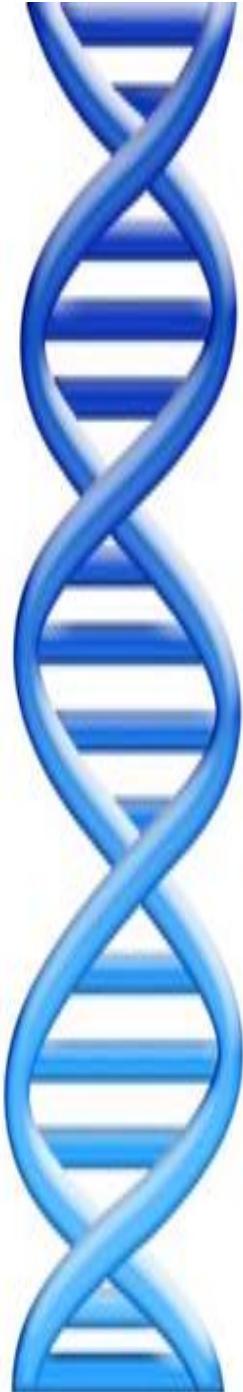


Computational Genomics

Michael Schatz

Oct 3, 2011
Frontiers in Genomics





Outline

Part I: Schatz Lab Overview

Part 2: Sequence Alignment

Part 3: Genome Assembly

Part 4: Parallel & Cloud Computing

A Little About Me



Computational Biology

"Computer science is no more about computers than astronomy is about telescopes."
Edger Dijkstra

Computer Science = Science of Computation

- Compute solutions to problems, designing & building systems
- Computers are very, very dumb, but we can instruct them
 - Build complex systems out of simple components

Computational Biology = Thinking Computationally about Biology

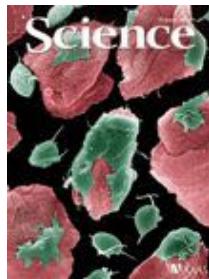
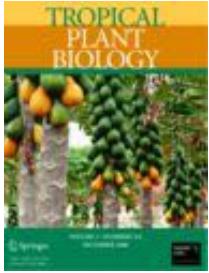
- Analysis: Make more powerful instruments, analyze results
- Design: experimental protocols, procedures, systems

Computational Genomics

1. Alignment
2. Assembly
3. Expression
4. Comparative Genomics

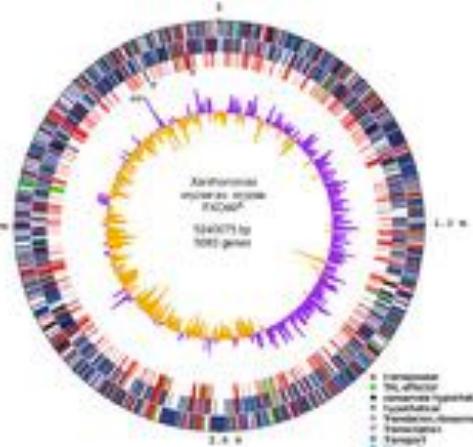
Computational Thinking

1. Algorithm
2. Data structure
3. Computational Analysis
4. Computational Modeling

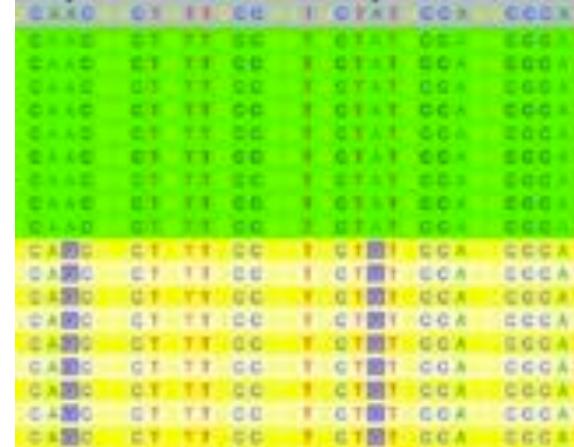


Genomics & Quantitative Biology

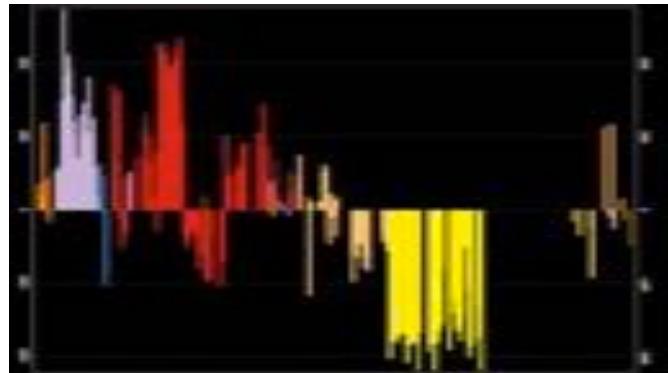
Genome Assembly



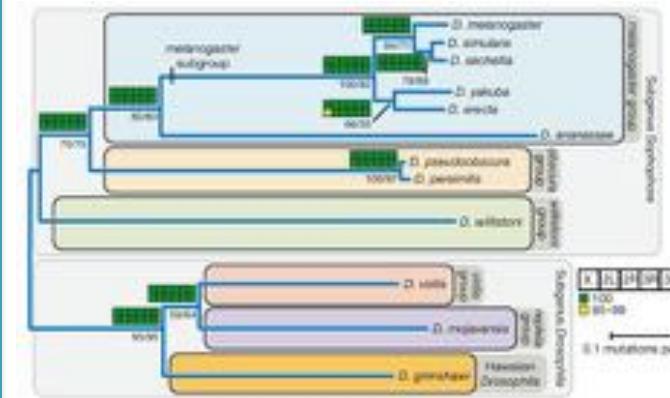
Mutations & Disease

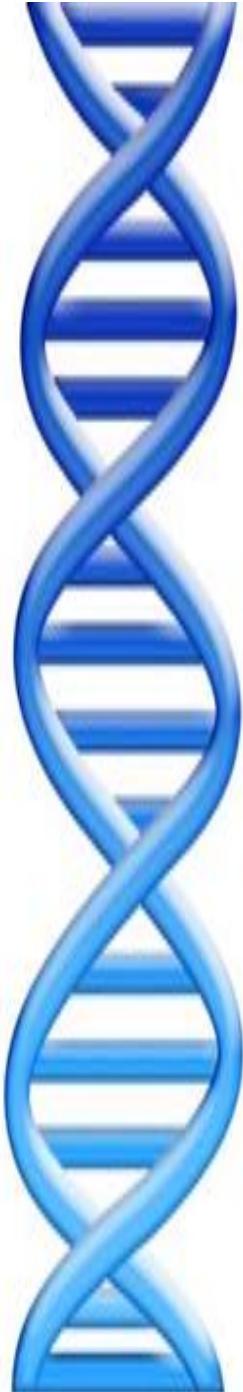


Differential Analysis



Phylogeny & Evolution





Outline

Part I: Schatz Lab Overview

Part 2: Sequence Alignment

- Exact Matching
- Suffix Arrays
- Bowtie and the BWT

Part 3: Genome Assembly

Part 4: Parallel & Cloud Computing

Searching for GATTACA

- Where is GATTACA in the human genome?
- Strategy I: Brute Force

I	2	3	4	5	6	7	8	9	10	11	12	13	14	15	...
T	G	A	T	T	A	C	A	G	A	T	T	A	C	C	...
G	A	T	T	A	C	A									

No match at offset 1

Searching for GATTACA

- Where is GATTACA in the human genome?
- Strategy I: Brute Force

I	2	3	4	5	6	7	8	9	10	11	12	13	14	15	...
T	G	A	T	T	A	C	A	G	A	T	T	A	C	C	...
	G	A	T	T	A	C	A								

Match at offset 2

Searching for GATTACA

- Where is GATTACA in the human genome?
- Strategy I: Brute Force

I	2	3	4	5	6	7	8	9	10	11	12	13	14	15	...
T	G	A	T	T	A	C	A	G	A	T	T	A	C	C	...
		G	A	T	T	A	C	A	...						

No match at offset 3...

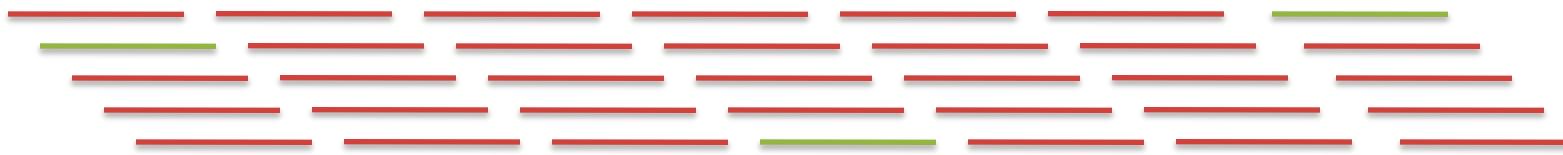
Searching for GATTACA

- Where is GATTACA in the human genome?
- Strategy I: Brute Force

I	2	3	4	5	6	7	8	9	10	11	12	13	14	15	...
T	G	A	T	T	A	C	A	G	A	T	T	A	C	C	...
								G	A	T	T	A	C	A	

No match at offset 9 <- Checking each possible position takes time

Brute Force Analysis



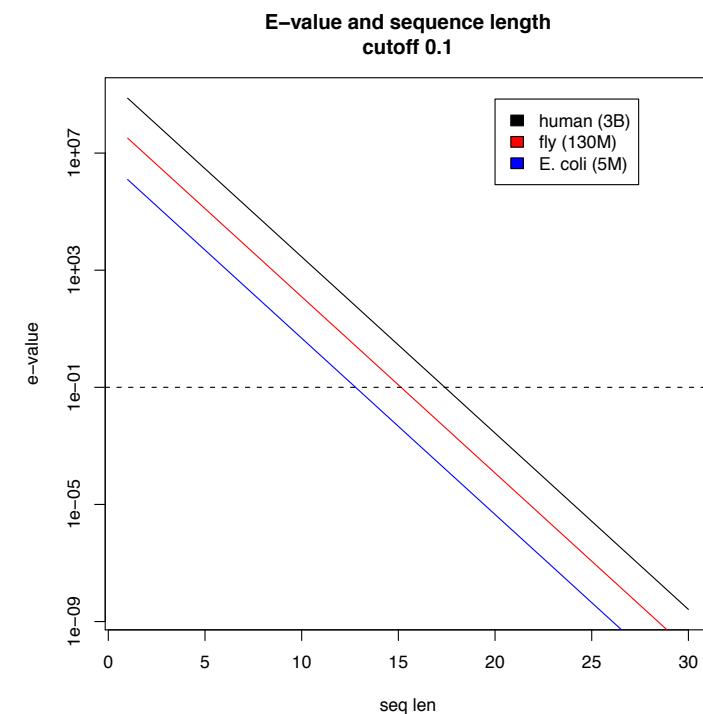
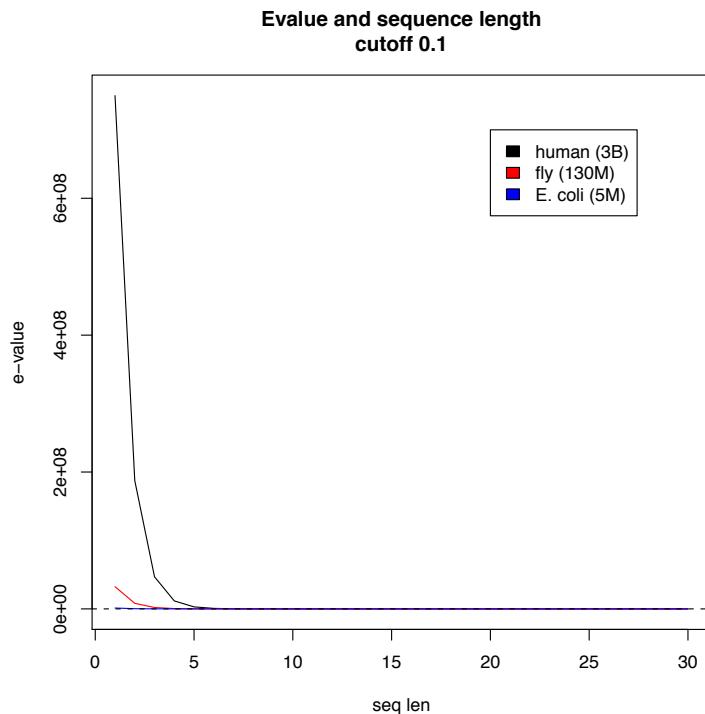
- Brute Force:
 - At every possible offset in the genome:
 - Do all of the characters of the query match?
- Analysis
 - Simple, easy to understand
 - Genome length = n [3B]
 - Query length = m [7]
 - Comparisons: $(n-m+1) * m$ [2IB]
- Overall runtime: $O(nm)$
 - [How long would it take if we double the genome size, read length?]
 - [How long would it take if we double both?]

Expected Occurrences

The expected number of occurrences (e-value) of a given sequence in a genome depends on the length of the genome and inversely on the length of the sequence

- 1 in 4 bases are G, 1 in 16 positions are GA, 1 in 64 positions are GAT, ...
- 1 in 16,384 should be GATTACA
- $E=n/(4^m)$

[183,105 expected occurrences]
[How long do the reads need to be for a significant match?]



Brute Force Reflections

Why check every position?

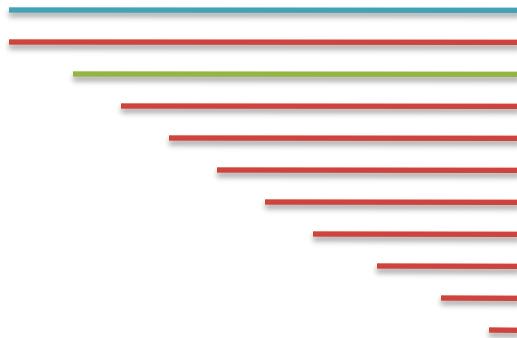
- GATTACA can't possibly start at position 15 [WHY?]

I	2	3	4	5	6	7	8	9	10	11	12	13	14	15	...
T	G	A	T	T	A	C	A	G	A	T	T	A	C	C	...
								G	A	T	T	A	C	A	

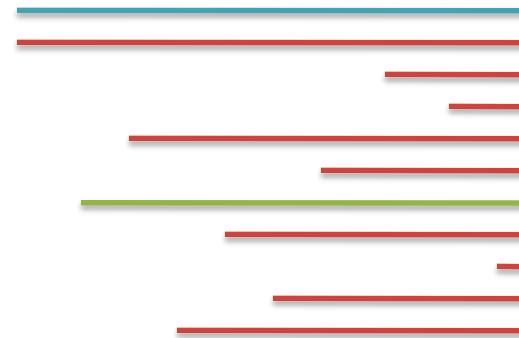
- Improve runtime to $O(n + m)$ [3B + 7]
 - If we double both, it just takes twice as long
 - Knuth-Morris-Pratt, 1977
 - Boyer-Moyer, 1977, 1991
- For one-off scans, this is the best we can do (optimal performance)
 - We have to read every character of the genome, and every character of the query
 - For short queries, runtime is dominated by the length of the genome

Suffix Arrays: Searching the Phone Book

- What if we need to check many queries?
 - We don't need to check every page of the phone book to find 'Schatz'
 - Sorting alphabetically lets us immediately skip 96% (25/26) of the book *without any loss in accuracy*
- Sorting the genome: Suffix Array (Manber & Myers, 1991)
 - Sort every suffix of the genome



Split into n suffixes



Sort suffixes alphabetically

[Challenge Question: How else could we split the genome?]

Searching the Index

- **Strategy 2: Binary search**
 - Compare to the middle, refine as higher or lower
- Searching for GATTACA
 - Lo = 1; Hi = 15;

#	Sequence	Pos
1	ACAGATTACC...	6
2	ACC...	13
3	AGATTACC...	8
4	ATTACAGATTACC...	3
5	ATTACC...	10
6	C...	15
7	CAGATTACC...	7
8	CC...	14
9	GATTACAGATTACC...	2
10	GATTACC...	9
11	TACAGATTACC...	5
12	TACC...	12
13	TGATTACAGATTACC...	1
14	TTACAGATTACC...	4
15	TTACC...	11

Searching the Index

- **Strategy 2: Binary search**
 - Compare to the middle, refine as higher or lower
- Searching for GATTACA
 - $\text{Lo} = 1; \text{Hi} = 15; \text{Mid} = (1+15)/2 = 8$
 - Middle = Suffix[8] = CC

#	Sequence	Pos
1	ACAGATTACC...	6
2	ACC...	13
3	AGATTACC...	8
4	ATTACAGATTACC...	3
5	ATTACC...	10
6	C...	15
7	CAGATTACC...	7
8	CC...	14
9	GATTACAGATTACC...	2
10	GATTACC...	9
11	TACAGATTACC...	5
12	TACC...	12
13	TGATTACAGATTACC...	1
14	TTACAGATTACC...	4
15	TTACC...	11

Lo
→

Hi
→

Searching the Index

- **Strategy 2: Binary search**
 - Compare to the middle, refine as higher or lower
- Searching for GATTACA
 - $\text{Lo} = 1; \text{Hi} = 15; \text{Mid} = (1+15)/2 = 8$
 - Middle = Suffix[8] = CC
=> Higher: Lo = Mid + 1

#	Sequence	Pos
1	ACAGATTACC...	6
2	ACC...	13
3	AGATTACC...	8
4	ATTACAGATTACC...	3
5	ATTACC...	10
6	C...	15
7	CAGATTACC...	7
8	CC...	14
9	GATTACAGATTACC...	2
10	GATTACC...	9
11	TACAGATTACC...	5
12	TACC...	12
13	TGATTACAGATTACC...	1
14	TTACAGATTACC...	4
15	TTACC...	11

Searching the Index

- Strategy 2: Binary search
 - Compare to the middle, refine as higher or lower
- Searching for GATTACA
 - Lo = 1; Hi = 15; Mid = (1+15)/2 = 8
 - Middle = Suffix[8] = CC
=> Higher: Lo = Mid + 1
 - Lo = 9; Hi = 15;

#	Sequence	Pos
1	ACAGATTACC...	6
2	ACC...	13
3	AGATTACC...	8
4	ATTACAGATTACC...	3
5	ATTACC...	10
6	C...	15
7	CAGATTACC...	7
8	CC...	14
9	GATTACAGATTACC...	2
10	GATTACC...	9
11	TACAGATTACC...	5
12	TACC...	12
13	TGATTACAGATTACC...	1
14	TTACAGATTACC...	4
15	TTACC...	11

Lo

Hi

Searching the Index

- **Strategy 2: Binary search**
 - Compare to the middle, refine as higher or lower
- Searching for GATTACA
 - $Lo = 1; Hi = 15; Mid = (1+15)/2 = 8$
 - Middle = Suffix[8] = CC
=> Higher: $Lo = Mid + 1$
 - $Lo = 9; Hi = 15; Mid = (9+15)/2 = 12$
 - Middle = Suffix[12] = TACC

#	Sequence	Pos
1	ACAGATTACC...	6
2	ACC...	13
3	AGATTACC...	8
4	ATTACAGATTACC...	3
5	ATTACC...	10
6	C...	15
7	CAGATTACC...	7
8	CC...	14
9	GATTACAGATTACC...	2
10	GATTACC...	9
11	TACAGATTACC...	5
12	TACC...	12
13	TGATTACAGATTACC...	1
14	TTACAGATTACC...	4
15	TTACC...	11

Lo

Hi

Searching the Index

- **Strategy 2: Binary search**
 - Compare to the middle, refine as higher or lower
- Searching for GATTACA
 - Lo = 1; Hi = 15; Mid = $(1+15)/2 = 8$
 - Middle = Suffix[8] = CC
=> Higher: Lo = Mid + 1
 - Lo = 9; Hi = 15; Mid = $(9+15)/2 = 12$
 - Middle = Suffix[12] = TACC
=> Lower: Hi = Mid - 1
 - Lo = 9; Hi = 11;

#	Sequence	Pos
1	ACAGATTACC...	6
2	ACC...	13
3	AGATTACC...	8
4	ATTACAGATTACC...	3
5	ATTACC...	10
6	C...	15
7	CAGATTACC...	7
8	CC...	14
9	GATTACAGATTACC...	2
10	GATTACC...	9
11	TACAGATTACC...	5
12	TACC...	12
13	TGATTACAGATTACC...	1
14	TTACAGATTACC...	4
15	TTACC...	11

Lo
→

Hi
→

Searching the Index

- **Strategy 2: Binary search**
 - Compare to the middle, refine as higher or lower
- Searching for GATTACA
 - Lo = 1; Hi = 15; Mid = $(1+15)/2 = 8$
 - Middle = Suffix[8] = CC
=> Higher: Lo = Mid + 1
 - Lo = 9; Hi = 15; Mid = $(9+15)/2 = 12$
 - Middle = Suffix[12] = TACC
=> Lower: Hi = Mid - 1
 - Lo = 9; Hi = 11; Mid = $(9+11)/2 = 10$
 - Middle = Suffix[10] = GATTACC

#	Sequence	Pos
1	ACAGATTACC...	6
2	ACC...	13
3	AGATTACC...	8
4	ATTACAGATTACC...	3
5	ATTACC...	10
6	C...	15
7	CAGATTACC...	7
8	CC...	14
9	GATTACAGATTACC...	2
10	GATTACC...	9
11	TACAGATTACC...	5
12	TACC...	12
13	TGATTACAGATTACC...	1
14	TTACAGATTACC...	4
15	TTACC...	11

Lo

Hi

Searching the Index

- **Strategy 2: Binary search**
 - Compare to the middle, refine as higher or lower
- Searching for GATTACA
 - Lo = 1; Hi = 15; Mid = $(1+15)/2 = 8$
 - Middle = Suffix[8] = CC
=> Higher: Lo = Mid + 1
 - Lo = 9; Hi = 15; Mid = $(9+15)/2 = 12$
 - Middle = Suffix[12] = TACC
=> Lower: Hi = Mid - 1
 - Lo = 9; Hi = 11; Mid = $(9+11)/2 = 10$
 - Middle = Suffix[10] = GATTACC
=> Lower: Hi = Mid - 1
 - Lo = 9; Hi = 9;

#	Sequence	Pos
1	ACAGATTACC...	6
2	ACC...	13
3	AGATTACC...	8
4	ATTACAGATTACC...	3
5	ATTACC...	10
6	C...	15
7	CAGATTACC...	7
8	CC...	14
9	GATTACAGATTACC...	2
10	GATTACC...	9
11	TACAGATTACC...	5
12	TACC...	12
13	TGATTACAGATTACC...	1
14	TTACAGATTACC...	4
15	TTACC...	11

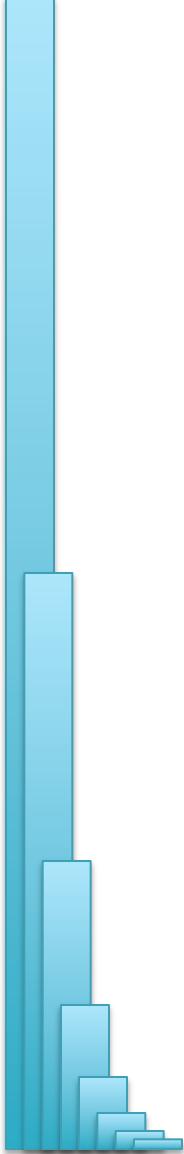
Lo
Hi

Searching the Index

- **Strategy 2: Binary search**
 - Compare to the middle, refine as higher or lower
- Searching for GATTACA
 - Lo = 1; Hi = 15; Mid = $(1+15)/2 = 8$
 - Middle = Suffix[8] = CC
=> Higher: Lo = Mid + 1
 - Lo = 9; Hi = 15; Mid = $(9+15)/2 = 12$
 - Middle = Suffix[12] = TACC
=> Lower: Hi = Mid - 1
 - Lo = 9; Hi = 11; Mid = $(9+11)/2 = 10$
 - Middle = Suffix[10] = GATTACC
=> Lower: Hi = Mid - 1
 - Lo = 9; Hi = 9; Mid = $(9+9)/2 = 9$
 - Middle = Suffix[9] = GATTACA...
=> Match at position 2!



#	Sequence	Pos
1	ACAGATTACC...	6
2	ACC...	13
3	AGATTACC...	8
4	ATTACAGATTACC...	3
5	ATTACC...	10
6	C...	15
7	CAGATTACC...	7
8	CC...	14
9	GATTACA GATTACC...	2
10	GATTACC...	9
11	TACAGATTACC...	5
12	TACC...	12
13	TGATTACAGATTACC...	1
14	TTACAGATTACC...	4
15	TTACC...	11



Binary Search Analysis

- Binary Search

Initialize search range to entire list

$mid = (hi+lo)/2$; $middle = suffix[mid]$

 if query matches middle: done

 else if query < middle: pick low range

 else if query > middle: pick hi range

Repeat until done or empty range

[WHEN?]

- Analysis

- More complicated method

- How many times do we repeat?

- How many times can it cut the range in half?

- Find smallest x such that: $n/(2^x) \leq 1$; $x = \lg_2(n)$

[32]

- Total Runtime: $O(m \lg n)$

- More complicated, but **much** faster!

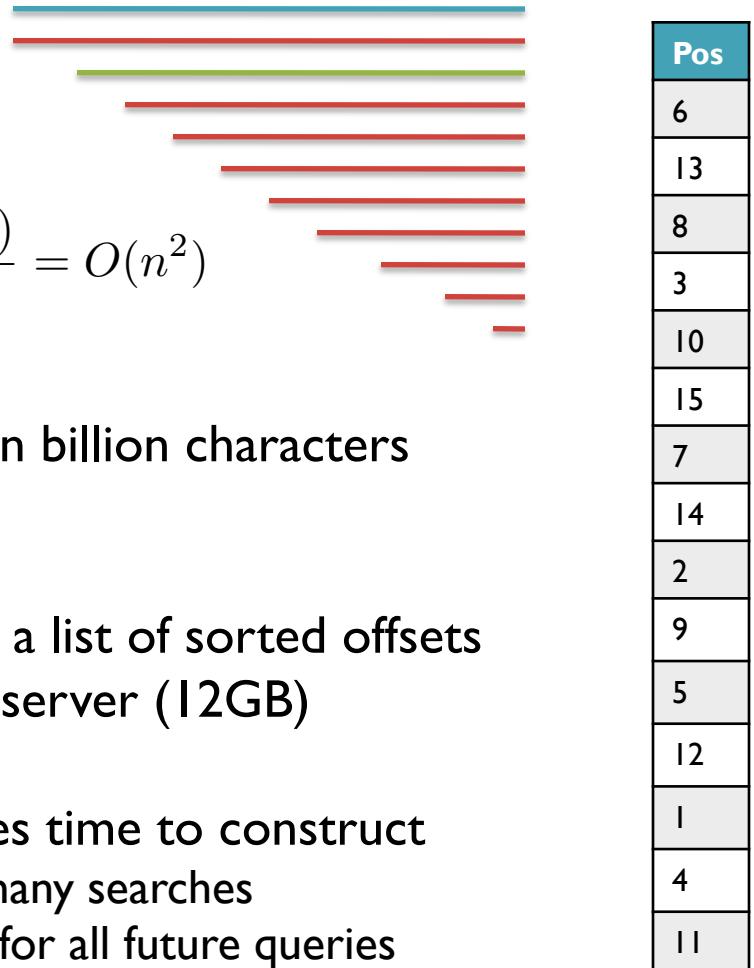
- Looking up a query loops 32 times instead of 3B

[How long does it take to search 6B or 24B nucleotides?]

Suffix Array Construction

- How can we store the suffix array?
[How many characters are in all suffixes combined?]

$$S = 1 + 2 + 3 + \dots + n = \sum_{i=1}^n i = \frac{n(n+1)}{2} = O(n^2)$$



- Hopeless to explicitly store 4.5 billion billion characters
- Instead use implicit representation
 - Keep 1 copy of the genome, and a list of sorted offsets
 - Storing 3 billion offsets fits on a server (12GB)
- Searching the array is very fast, but it takes time to construct
 - This time will be amortized over many, many searches
 - Run it once "overnight" and save it away for all future queries

TGATTACAGATTACC

Sorting

Quickly sort these numbers into ascending order:

14, 29, 6, 31, 39, 64, 78, 50, 13, 63, 61, 19

[How do you do it?]

6, 13, 14, 29, 31, 39, 64, 78, 50, 63, 61, 19

6, 13, 14, 29, 31, 39, 64, 78, 50, 63, 61, 19

6, 13, 14, 19, 29, 31, 39, 64, 78, 50, 63, 61

6, 13, 14, 19, 29, 31, 39, 64, 78, 50, 63, 61

6, 13, 14, 19, 29, 31, 39, 64, 78, 50, 63, 61

6, 13, 14, 19, 29, 31, 39, 50, 64, 78, 63, 61

6, 13, 14, 19, 29, 31, 39, 50, 61, 64, 78, 63

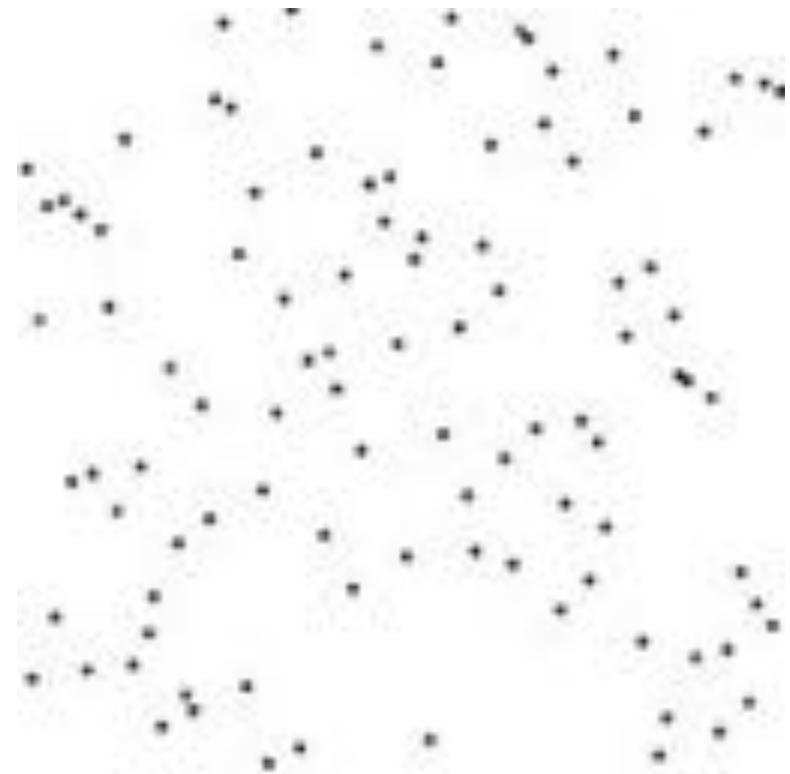
6, 13, 14, 19, 29, 31, 39, 50, 61, 63, 64, 78

6, 13, 14, 19, 29, 31, 39, 50, 61, 63, 64, 78

6, 13, 14, 19, 29, 31, 39, 50, 61, 63, 64, 78

6, 13, 14, 19, 29, 31, 39, 50, 61, 63, 64, 78

6, 13, 14, 19, 29, 31, 39, 50, 61, 63, 64, 78



Selection Sort Analysis

- Selection Sort (Input: list of n numbers)

```
for pos = 1 to n
```

```
    // find the smallest element in [pos, n]
```

```
    smallest = pos
```

```
    for check = pos+1 to n
```

```
        if (list[check] < list[smallest]): smallest = check
```

```
    // move the smallest element to the front
```

```
    tmp = list[smallest]
```

```
    list[pos] = list[smallest]
```

```
    list[smallest] = tmp
```

- Complexity Analysis

$$T = n + (n - 1) + (n - 2) + \dots + 3 + 2 + 1 = \sum_{i=1}^n i = \frac{n(n + 1)}{2} = O(n^2)$$

- Outer loop: pos = 1 to n

- Inner loop: check = pos to n

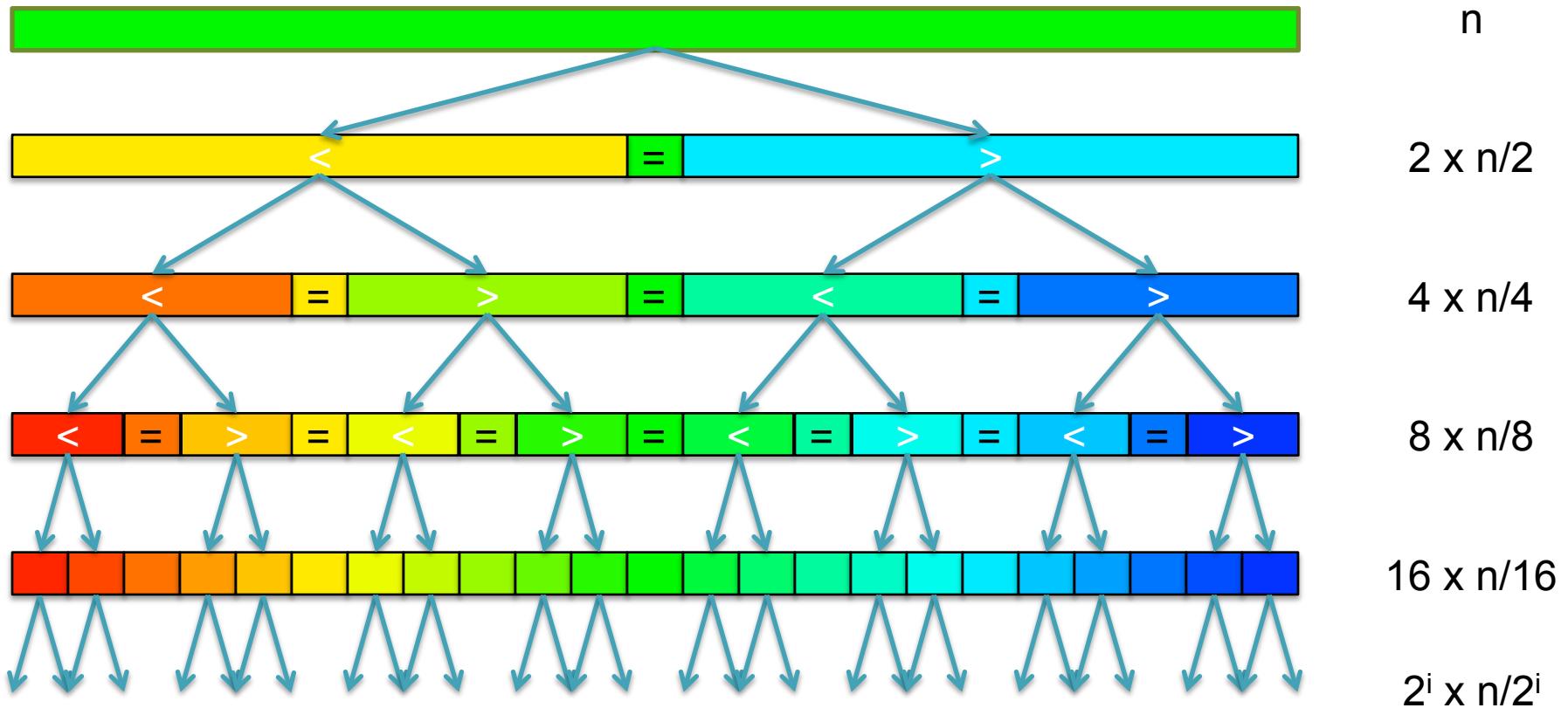
- Running time: Outer * Inner = $O(n^2)$

[4.5 Billion Billion]

[Challenge Questions: Why is this slow? / Can we sort any faster?]

Divide and Conquer

- Selection sort is slow because it rescans the entire list for each element
 - How can we split up the unsorted list into independent ranges?
 - Hint 1: Binary search splits up the problem into 2 independent ranges (hi/lo)
 - Hint 2: Assume we know the median value of a list



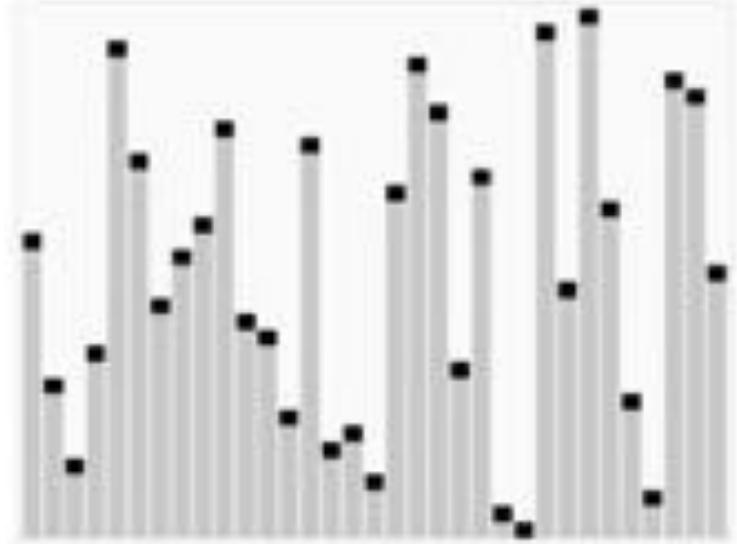
[How many times can we split of n items a list in half?]

QuickSort Analysis

- QuickSort(Input: list of n numbers)
// see if we can quit
if (length(list)) <= 1): return list

// split list into lo & hi
pivot = median(list)
lo = {}; hi = {};
for (i = 1 to length(list))
 if (list[i] < pivot): append(lo, list[i])
 else: append(hi, list[i])

// recurse on sublists
return (append(QuickSort(lo), QuickSort(hi)))



<http://en.wikipedia.org/wiki/Quicksort>

- Complexity Analysis (Assume we can find the median in $O(n)$)

$$T(n) = \begin{cases} O(1) & \text{if } n \leq 1 \\ O(n) + 2T(n/2) & \text{else} \end{cases}$$

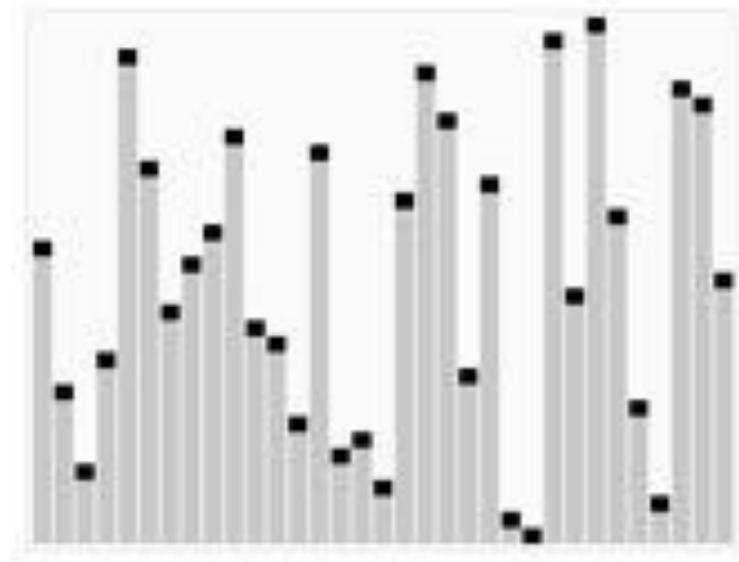
$$T(n) = n + 2\left(\frac{n}{2}\right) + 4\left(\frac{n}{4}\right) + \dots + n\left(\frac{n}{n}\right) = \sum_{i=0}^{\lg(n)} \frac{2^i n}{2^i} = \sum_{i=0}^{\lg(n)} n = O(n \lg n) \quad [\sim 94B]$$

QuickSort Analysis

- QuickSort(Input: list of n numbers)
// see if we can quit
if (length(list)) <= 1): return list

// split list into lo & hi
pivot = median(list)
lo = {}; hi = {};
for (i = 1 to length(list))
 if (list[i] < pivot): append(lo, list[i])
 else: append(hi, list[i])

// recurse on sublists
return (append(QuickSort(lo), QuickSort(hi)))



<http://en.wikipedia.org/wiki/Quicksort>

- Complexity Analysis (Assume we can find the median in $O(n)$)

$$T(n) = \begin{cases} O(1) & \text{if } n \leq 1 \\ O(n) + 2T(n/2) & \text{else} \end{cases}$$

$$T(n) = n + 2\left(\frac{n}{2}\right) + 4\left(\frac{n}{4}\right) + \dots + n\left(\frac{n}{n}\right) = \sum_{i=0}^{\lg(n)} \frac{2^i n}{2^i} = \sum_{i=0}^{\lg(n)} n = O(n \lg n) \quad [\sim 94B]$$

In-exact alignment

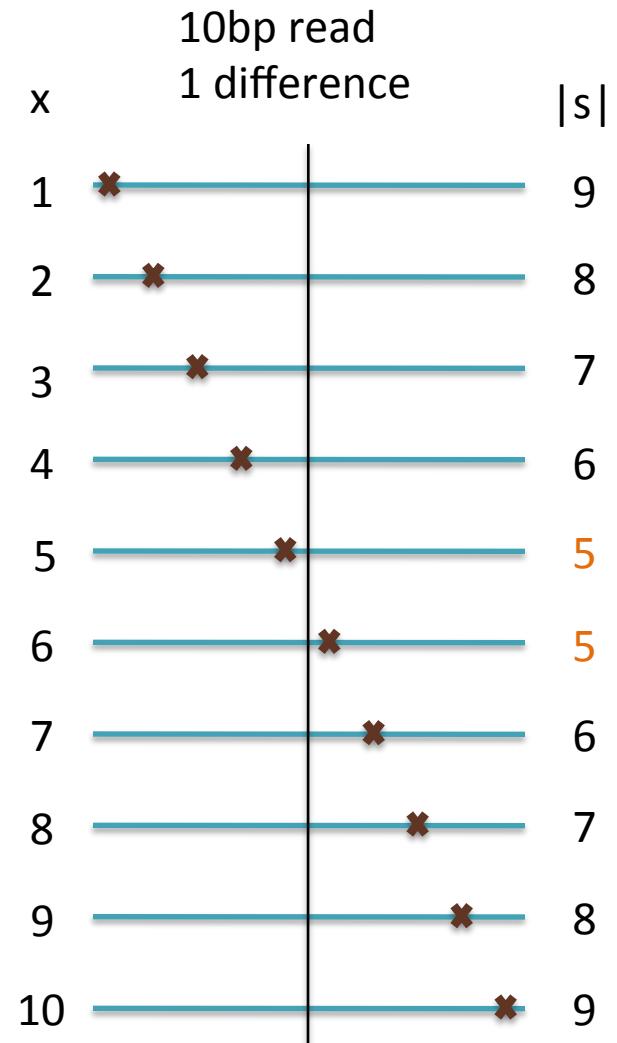
I	2	3	4	5	6	7	8	9	10	11	12	13	14	15	...
T	G	A	T	T	A	C	A	G	A	T	T	A	C	C	...

- Where is GATTACA *approximately* in the human genome?
 - And how do we efficiently find them?
- It depends...
 - Define 'approximately'
 - Hamming Distance, Edit distance, or Sequence Similarity
 - Ungapped vs Gapped vs Affine Gaps, Global vs Local
 - Algorithm depends on the data characteristics & goals
 - Smith-Waterman: Exhaustive search for optimal alignments
 - BLAST: Hash-table based homology searches
 - Bowtie: BWT alignment for short read mapping

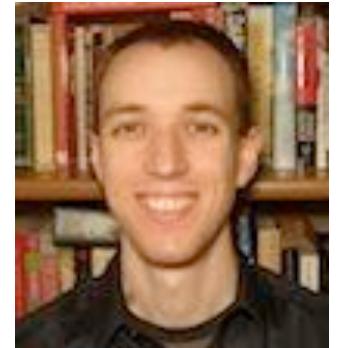
Seed-and-Extend Alignment

Theorem: An alignment of a sequence of length m with at most k differences **must** contain an exact match at least $s=m/(k+1)$ bp long
(Baeza-Yates and Perleberg, 1996)

- Proof: Pigeonhole principle
 - 1 pigeon can't fill 2 holes
- Seed-and-extend search
 - Use an index to rapidly find short exact alignments to seed longer in-exact alignments
 - BLAST, MUMmer, Bowtie, BWA, SOAP, ...



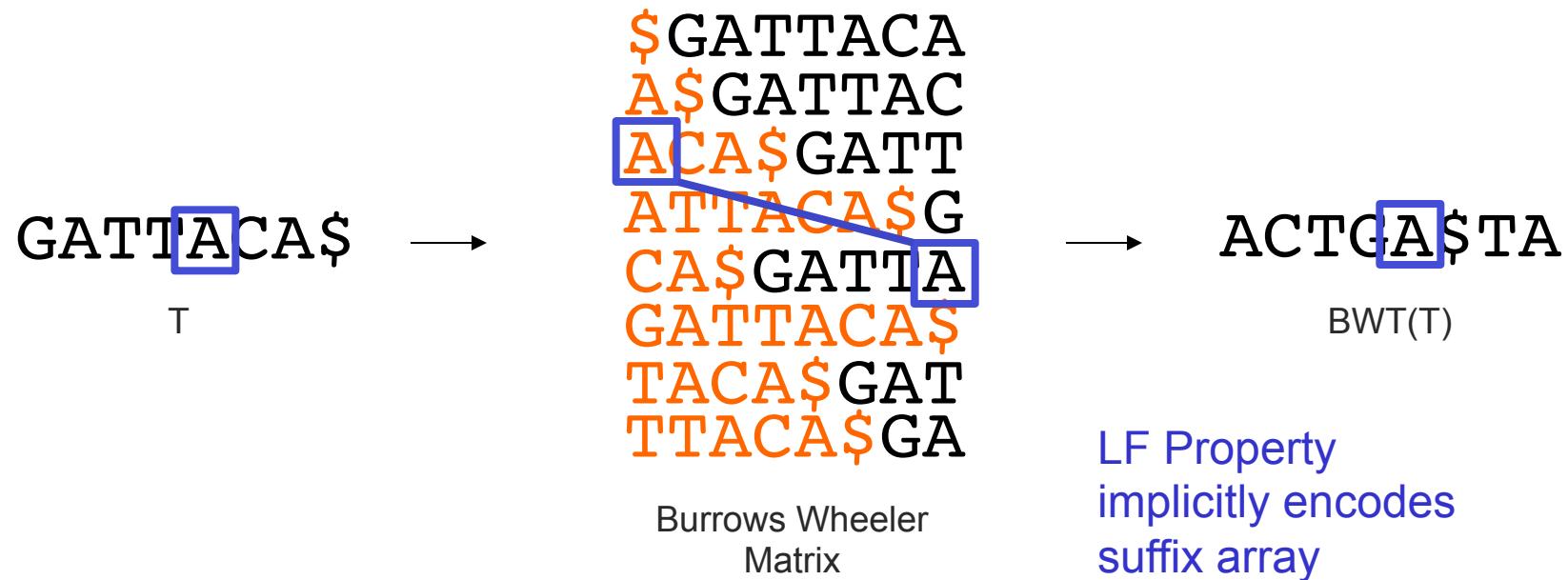
[How could you use seed-and-extend with a suffix array?]



Bowtie: Ultrafast and memory efficient alignment of short DNA sequences to the human genome

Slides Courtesy of Ben Langmead
langmead@umiacs.umd.edu

Burrows-Wheeler Transform



- Suffix Array is fast to search, but much larger than genome
 - BWT is a reversible permutation of the genome based on the suffix array
 - Core index for Bowtie (Langmead et al., 2009) and most recent short read mapping applications

A block sorting lossless data compression algorithm.

Burrows M, Wheeler DJ (1994) *Digital Equipment Corporation, Palo Alto, CA* 1994, Technical Report 124

Bowtie algorithm

Reference



BWT(Reference)

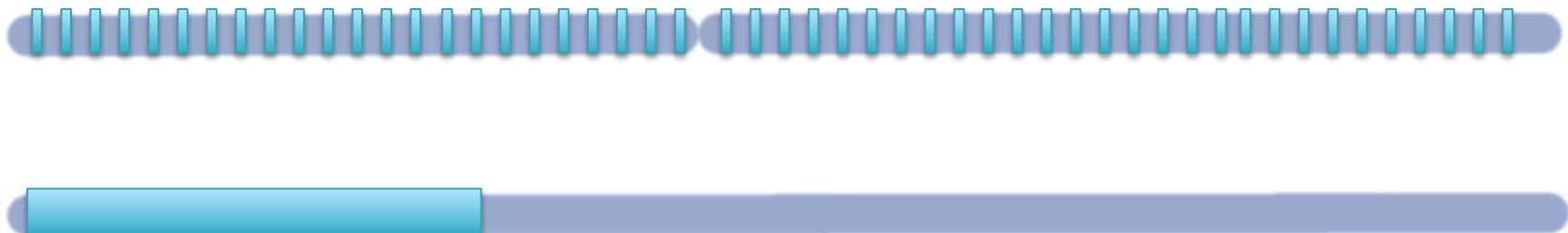
Query:

AATGATAACGGCGACCAACCGAGATCTA



Bowtie algorithm

Reference



BWT(Reference)

Query:

AATGATACGGCGACCAACCGAGATCTA



Bowtie algorithm

Reference



BWT(Reference)

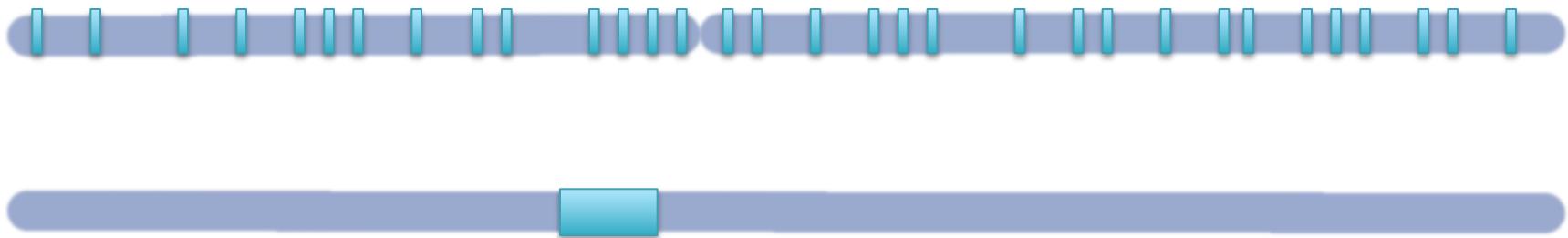
Query:

AATGATAACGGCGACCAACCGAGATC TA



Bowtie algorithm

Reference



BWT(Reference)

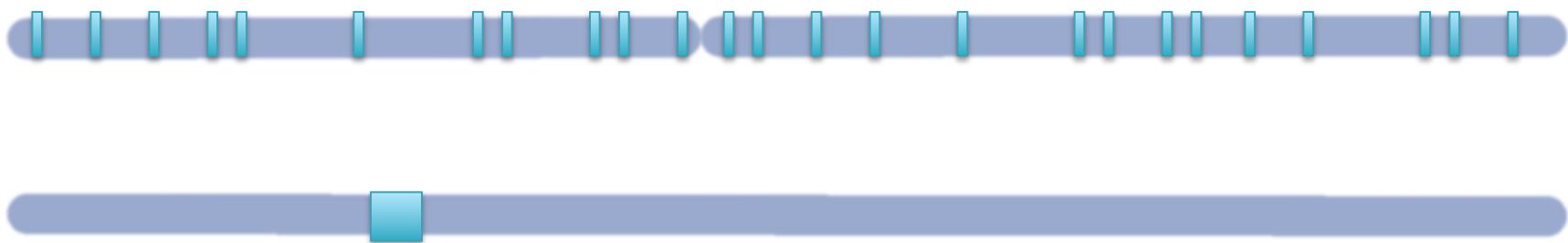
Query:

AATGATAACGGCGACCAACCGAGATCTA



Bowtie algorithm

Reference



BWT(Reference)

Query:

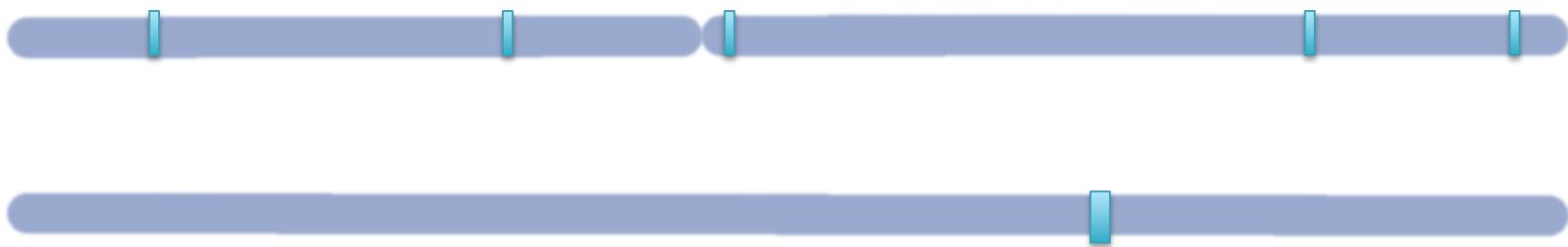
AATGATAACGGCGACC

CACCGAGATCTA



Bowtie algorithm

Reference



BWT(Reference)

Query:

AATGATACGGCGACCAACCGAGATCTA



Bowtie algorithm

Reference



BWT(Reference)

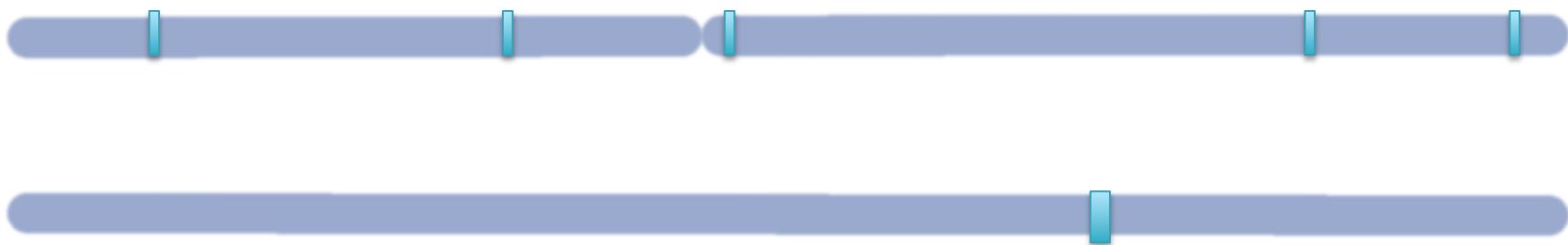
Query:

AATGATACGGCGACCAACCGAGATCTA



Bowtie algorithm

Reference



BWT(Reference)

Query:

AATG T TACGGCGACCACCGAGATCTA



Bowtie algorithm

Reference



BWT(Reference)

Query:

AATG **T** TACGGCGACCAACCGAGATCTA

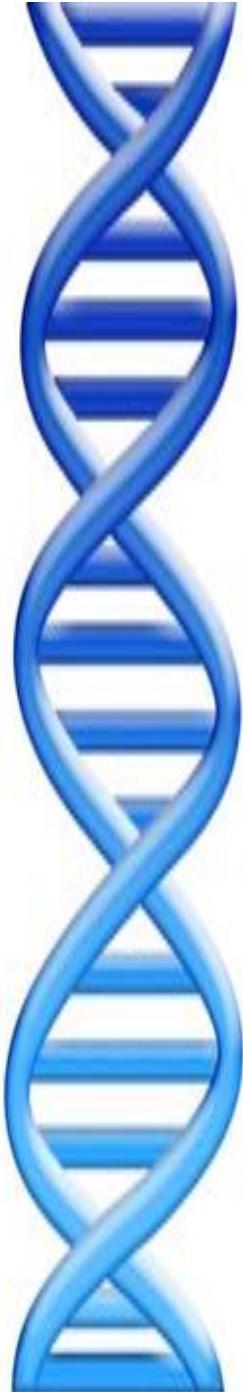


Part I: Summary

- Short Read Mapping: Seed-and-extend search of the BWT
 - If we fail to reach the end, back-track and resume search
 - The beginning of the read is used as high confidence seed
 - 100s of times faster than competing approaches, works entirely in RAM
- Algorithms choreograph the dance of data inside the machine
 - Algorithms add provable precision to your method
 - A smarter algorithm can solve the same problem with much less work
- Computational Techniques
 - **Binary search:** Fast lookup in any sorted list
 - **Divide-and-conquer:** Split a hard problem into an easier problem
 - **Recursion:** Solve a problem using a function of itself
 - **Indexing:** Focus on just the important parts
 - **Seed-and-extend:** Anchor the problem using a portion of it

Break





Outline

Part 1: Schatz Lab Overview

Part 2: Sequence Alignment

Part 3: Genome Assembly

- Assembly by analogy
- Coverage, read length, and repeats
- Contiging & Scaffolding
- Assembly Forensics

Part 4: Parallel & Cloud Computing

Shredded Book Reconstruction

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

It was the best of times, it was the worst of times, it was the age of wisdom, it was the age of foolishness, ...

It was the best of times, it was the worst of times, it was the age of wisdom, it was the age of foolishness, ...

It was the best of times, it was the worst of times, it was the age of wisdom, it was the age of foolishness, ...

It was the best of times, it was the worst of times, it was the age of wisdom, it was the age of foolishness, ...

It was the best of times, it was the worst of times, it was the age of wisdom, it was the age of foolishness, ...

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

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

Greedy Reconstruction

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]

[Any ideas on how to proceed?]

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

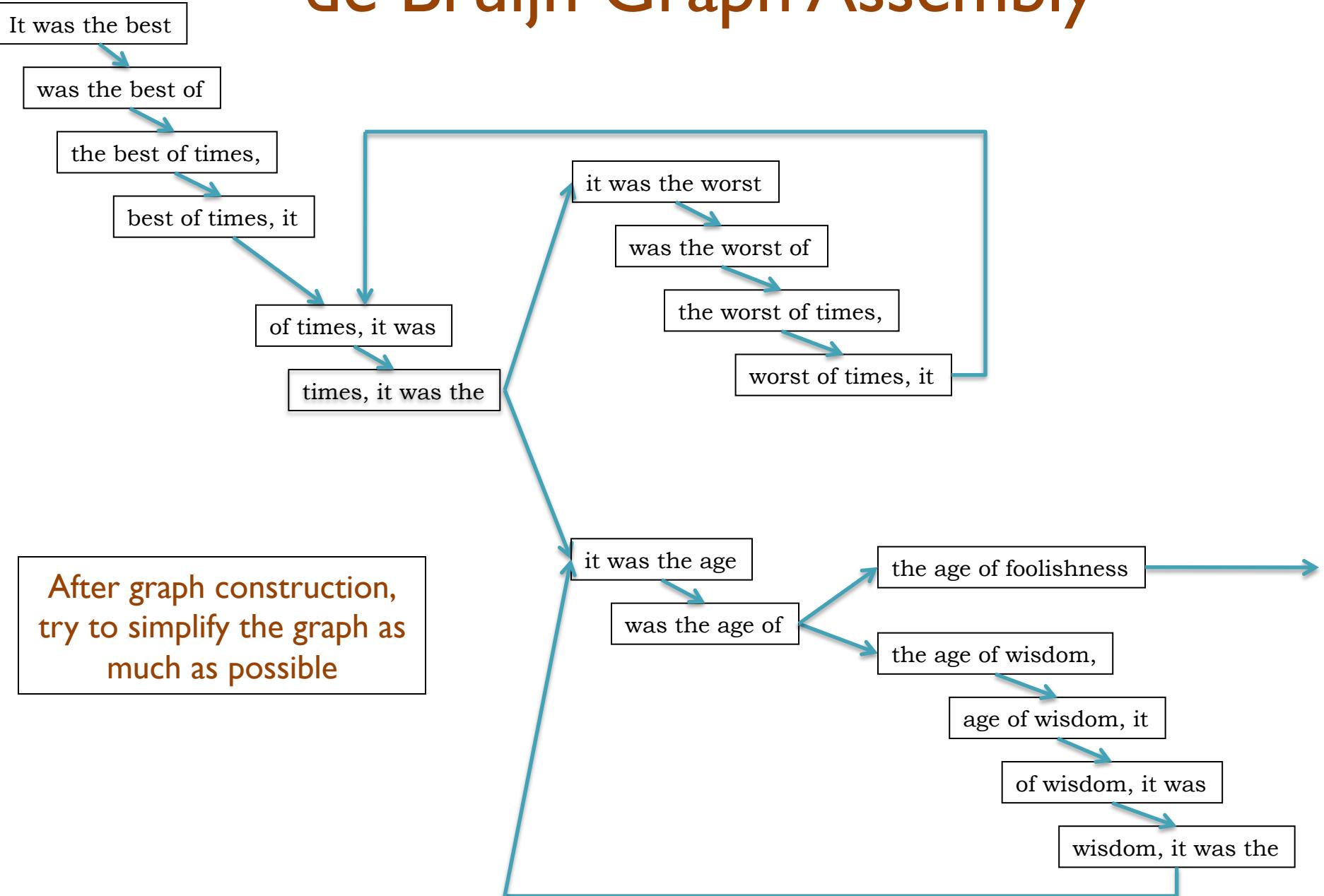
- Locally constructed graph reveals the global sequence structure
 - Overlaps between sequences implicitly computed

de Bruijn, 1946

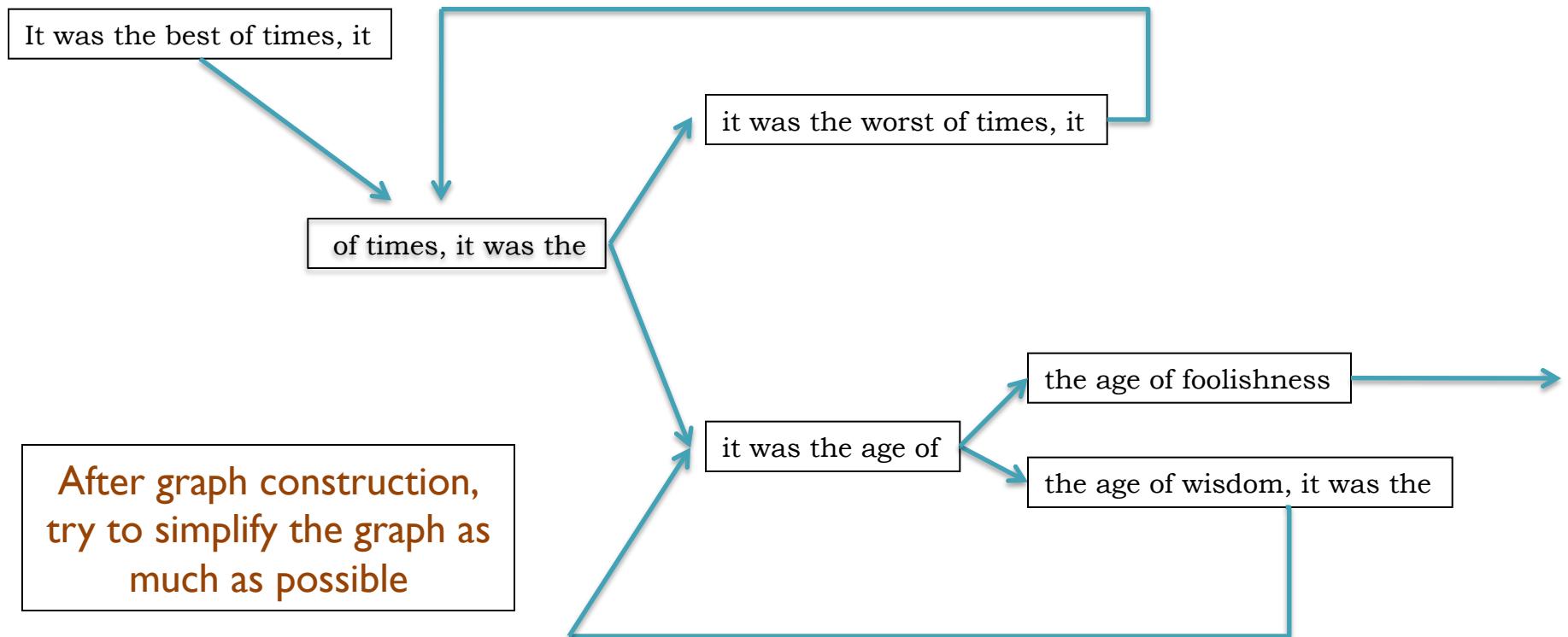
Idury and Waterman, 1995

Pevzner, Tang, Waterman, 2001

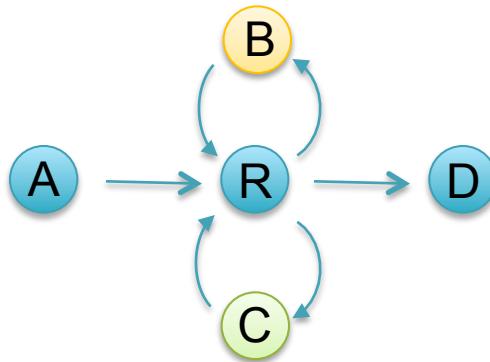
de Bruijn Graph Assembly



de Bruijn Graph Assembly



Counting Eulerian Tours



ARBRCRD
or
ARCRBRD

Generally an exponential number of compatible sequences

- Value computed by application of the BEST theorem (Hutchinson, 1975)

$$\mathcal{W}(G, t) = (\det L) \left\{ \prod_{u \in V} (r_u - 1)! \right\} \left\{ \prod_{(u,v) \in E} a_{uv}! \right\}^{-1}$$

L = $n \times n$ matrix with $r_u - a_{uu}$ along the diagonal and $-a_{uv}$ in entry uv

$r_u = d^+(u) + 1$ if $u=t$, or $d^+(u)$ otherwise

a_{uv} = multiplicity of edge from u to v

Assembly Complexity of Prokaryotic Genomes using Short Reads.
Kingsford C, Schatz MC, Pop M (2010) *BMC Bioinformatics*.

Milestones in Genome Assembly

Science Vol. 207 February 20, 1980
articles

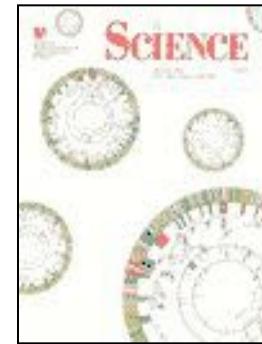
Nucleotide sequence of bacteriophage Φ X174 DNA

E.Sanger, G.M.Air, R.G.Bonell, N.E.Browne, A.R.Coulson, J.C.Fiddes,
C.A.Hinchliffe, B.P., P.M.Sherrard* & M.Smith*

*MRC Laboratory of Molecular Biology, 40 Cambridge Heath Road, London E1 2AD, UK.

An Φ X174 genome of approximately 54,000 bp was sequenced by the chain-terminating method using the rapid and simple 'one and seven' method. The minimum detection limit of the technique required for the production of the same or better quality of sequence data is approximately 10% of the genome. The sequence of the genome of Φ X174 is given for the positions and 871 bp. Two pairs of genes are coded by the same regions of Φ X174 using different reading frames.

The genome of bacteriophage Φ X174 is a single-circled, non-recombinant molecule of double-stranded DNA. The code of these genes is discussed by Sanger et al. (1980).



1977. Sanger et al.
1st Complete Organism
5375 bp

1995. Fleischmann et al.
1st Free Living Organism
TIGR Assembler. 1.8Mbp



1998. C.elegans SC
1st Multicellular Organism
BAC-by-BAC Phrap. 97Mbp



2000. Myers et al.
1st Large WGS Assembly.
Celera Assembler. 116 Mbp



2001. Venter et al., IHGSC
Human Genome
Celera Assembler/GigaAssembler. 2.9 Gbp



2010. Li et al.
1st Large SGS Assembly.
SOAPdenovo 2.2 Gbp

Like Dickens, we must computationally reconstruct a genome from short fragments

Current Applications

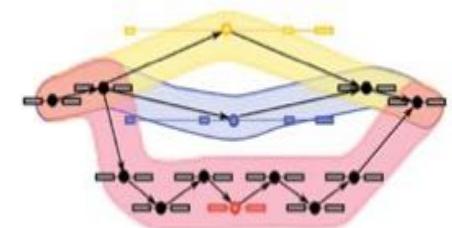
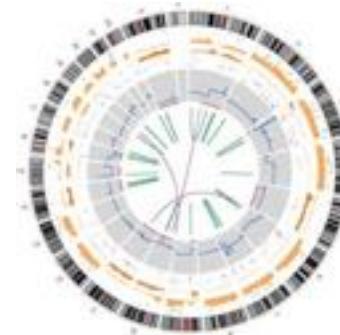
- Novel genomes



- Metagenomes

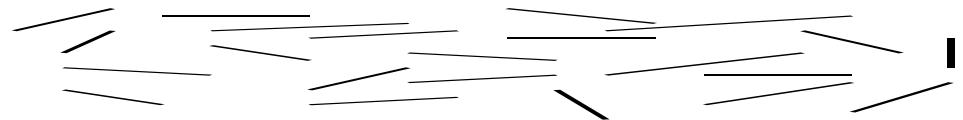


- Sequencing assays
 - Structural variations
 - Transcript assembly
 - ...



Assembling a Genome

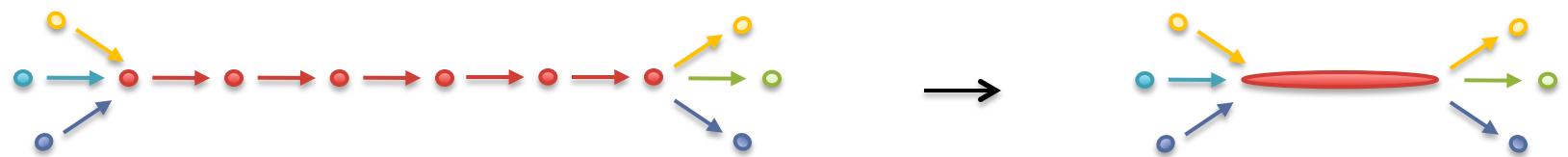
1. Shear & Sequence DNA



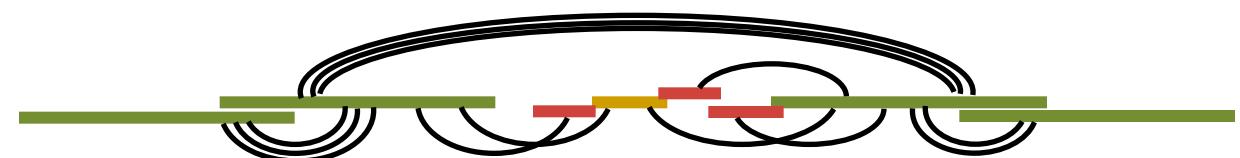
2. Construct assembly graph from overlapping reads

...AGCCTAGACCTACA**GGATGCGCGACACGT**
GGATGCGCGACACGTCGCATATCCGGT...

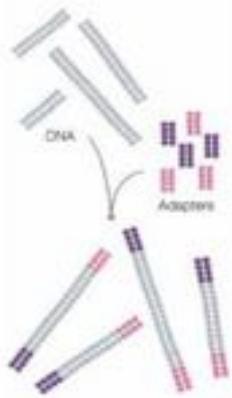
3. Simplify assembly graph



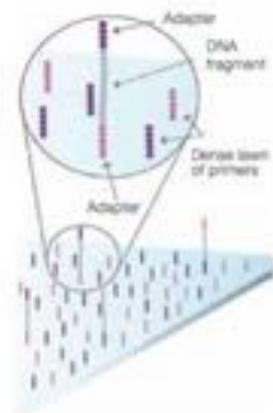
4. Detangle graph with long reads, mates, and other links



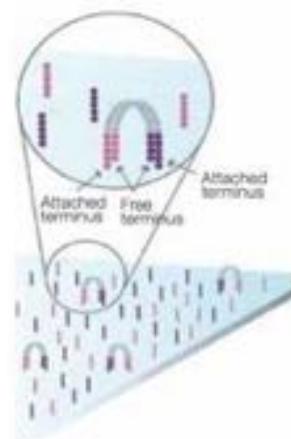
Illumina Sequencing by Synthesis



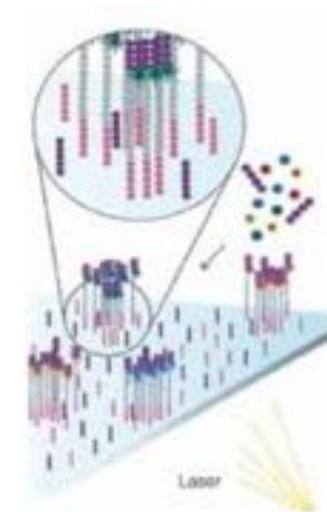
1. Prepare



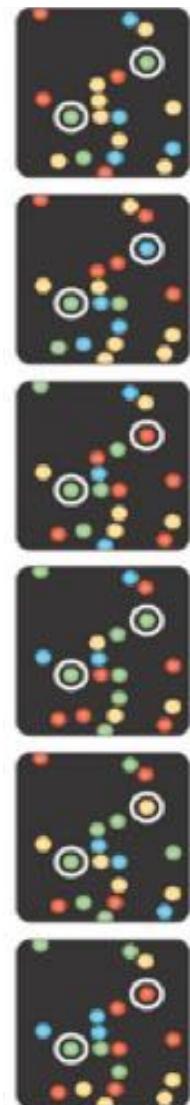
2. Attach



3. Amplify



4. Image



5. Basecall

Metzker (2010) Nature Reviews Genetics 11:31-46

http://www.illumina.com/documents/products/techspotlights/techspotlight_sequencing.pdf

Paired-end and Mate-pairs

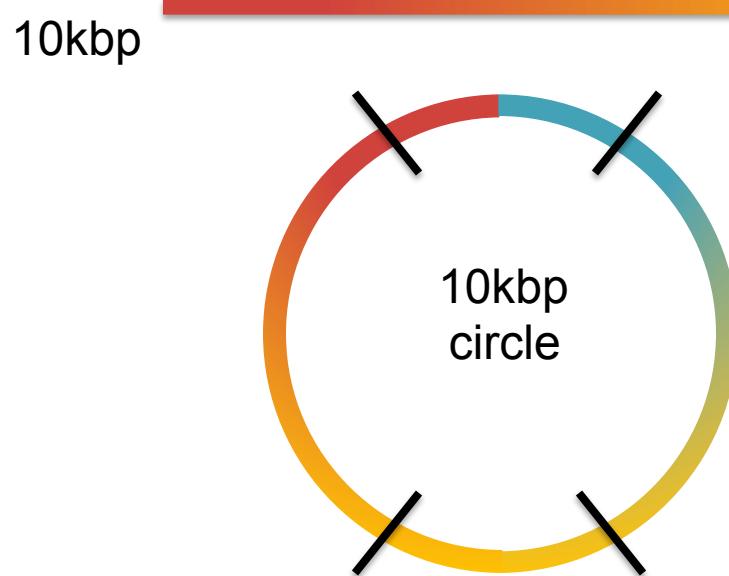
Paired-end sequencing

- Read one end of the molecule, flip, and read the other end
- Generate pair of reads separated by up to 500bp with inward orientation



Mate-pair sequencing

- Circularize long molecules (1-10kbp), shear into fragments, & sequence
- Mate failures create short paired-end reads



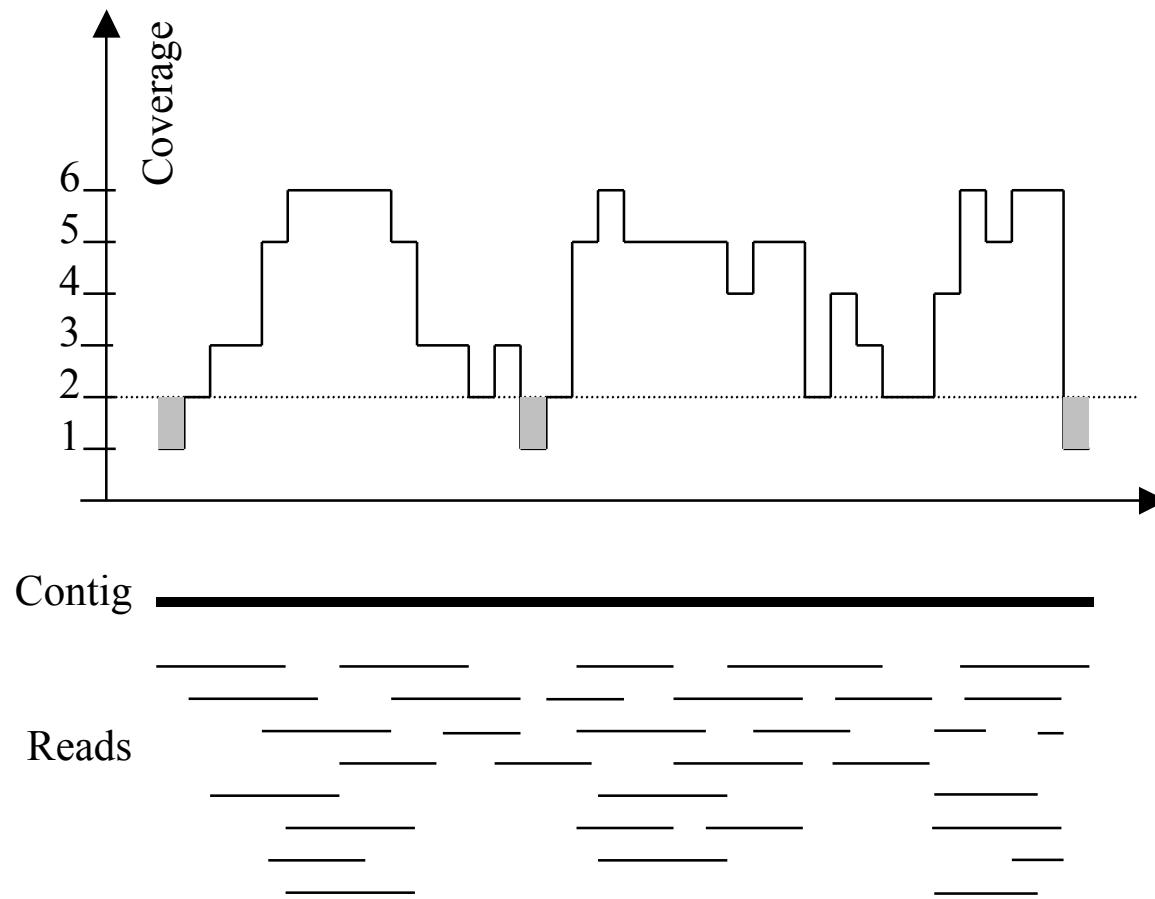
2x100 @ ~10kbp (outies)



2x100 @ 300bp (innies)



Typical genome coverage

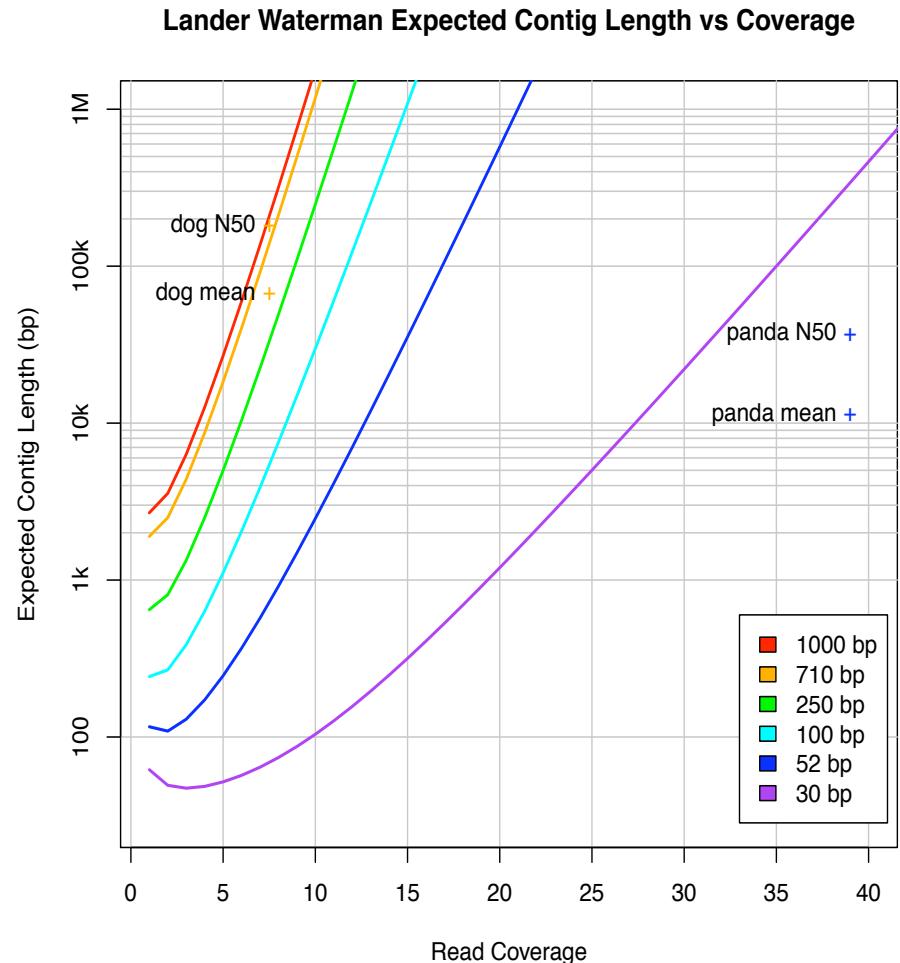


Imagine raindrops on a sidewalk

Coverage and Read Length

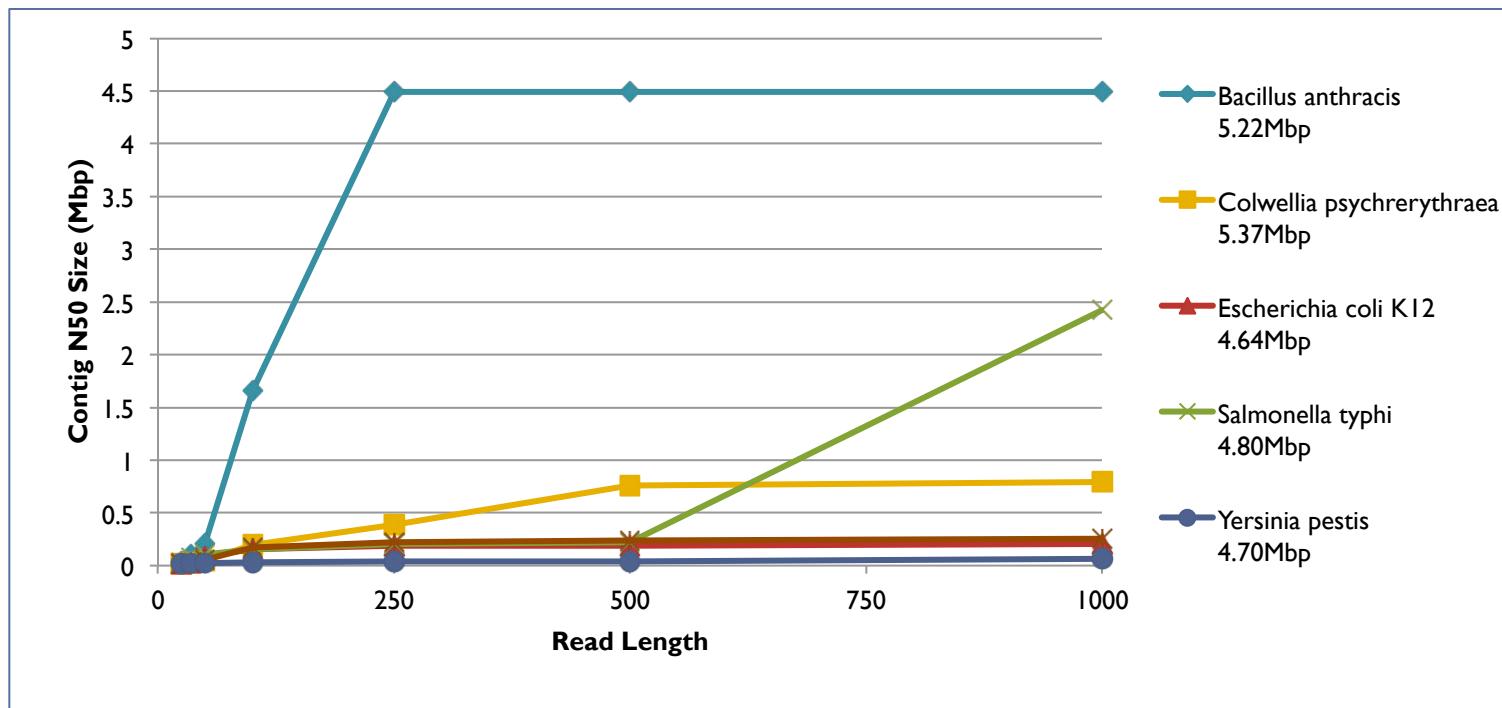
Idealized Lander-Waterman model

- Reads start at perfectly random positions
- Poisson distribution in coverage
 - Contigs end when there are no overlapping reads
- Contig length is a function of coverage and read length
 - Effective coverage reduced by overlap
 - Short reads require much higher coverage to reach same expected contig length



Assembly of Large Genomes using Second Generation Sequencing
Schatz MC, Delcher AL, Salzberg SL (2010) *Genome Research*. 20:1165-1173.

Repeats and Read Length



- Explore the relationship between read length and contig N50 size
 - Idealized assembly of read lengths: 25, 35, 50, 100, 250, 500, 1000
 - Contig/Read length relationship depends on specific repeat composition

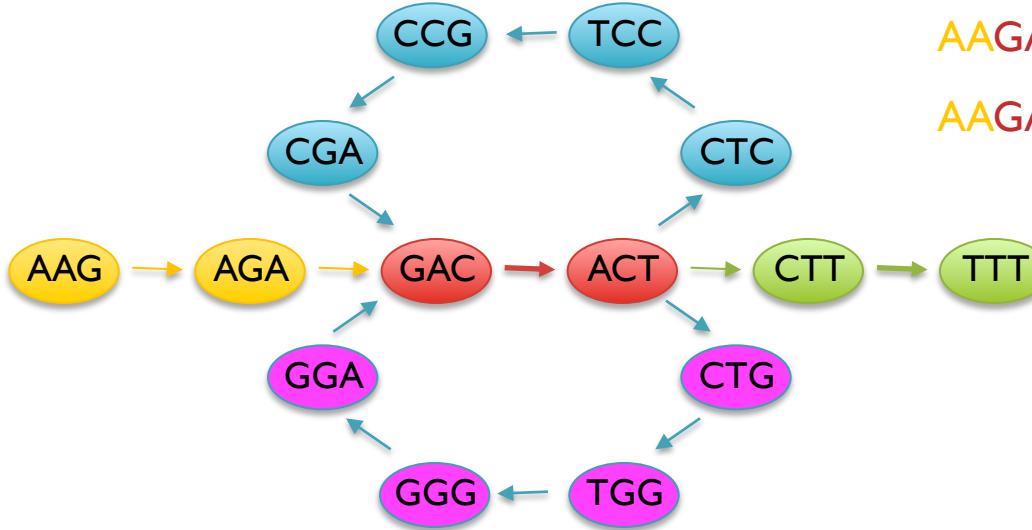
Assembly Complexity of Prokaryotic Genomes using Short Reads.
Kingsford C, Schatz MC, Pop M (2010) *BMC Bioinformatics*. 11:21.

Short Read Assembly

Reads

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

de Bruijn Graph



Potential Genomes

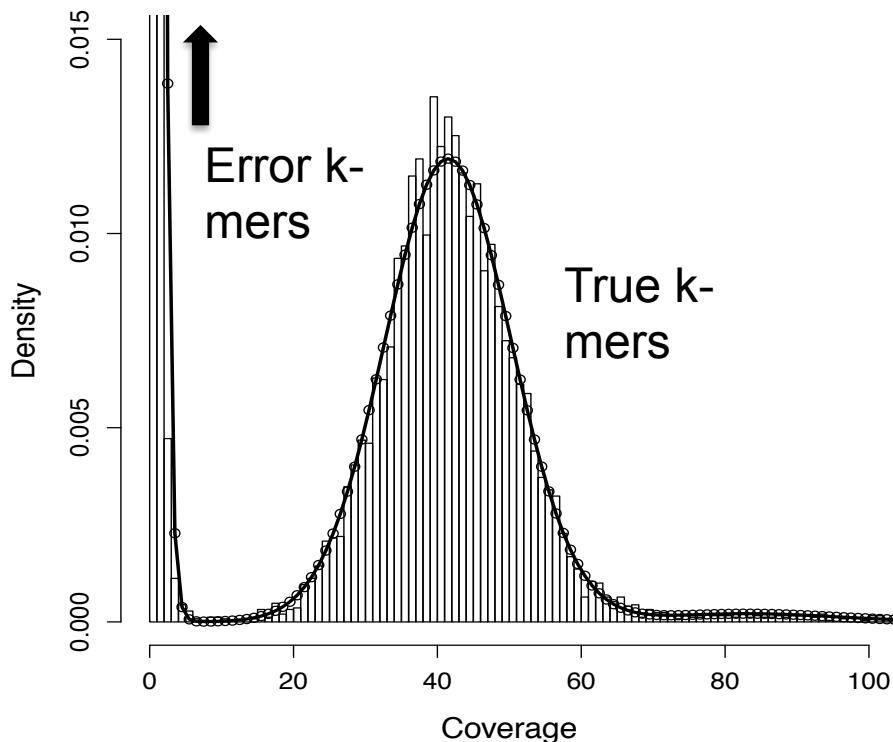
AAGACTCCGACTGGGACTTT
AAGACTGGGACTCCGACTTT

- 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

Error Correction with Quake

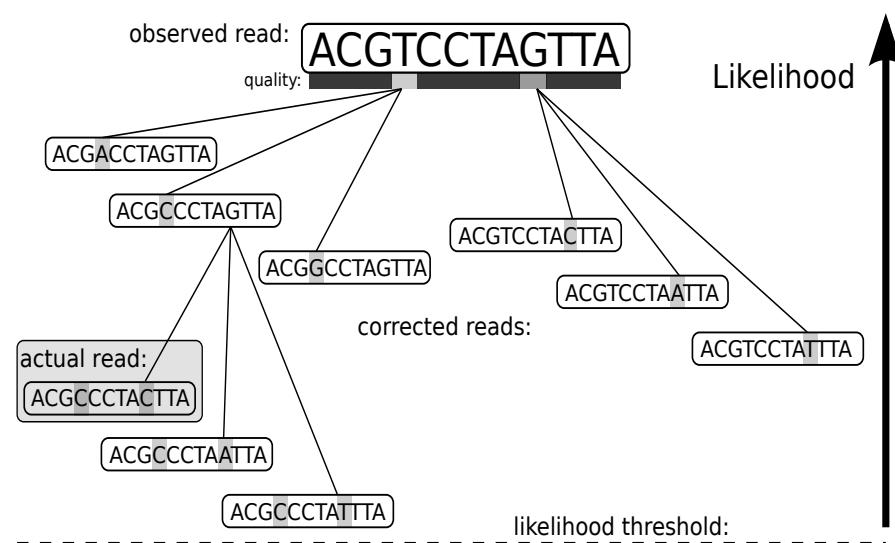
I. Count all “Q-mers” in reads

- Fit coverage distribution to mixture model of errors and regular coverage
- Automatically determines threshold for trusted k-mers



2. Correction Algorithm

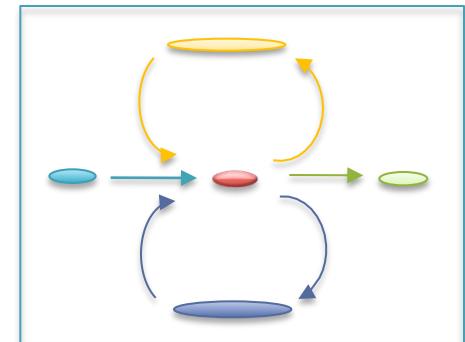
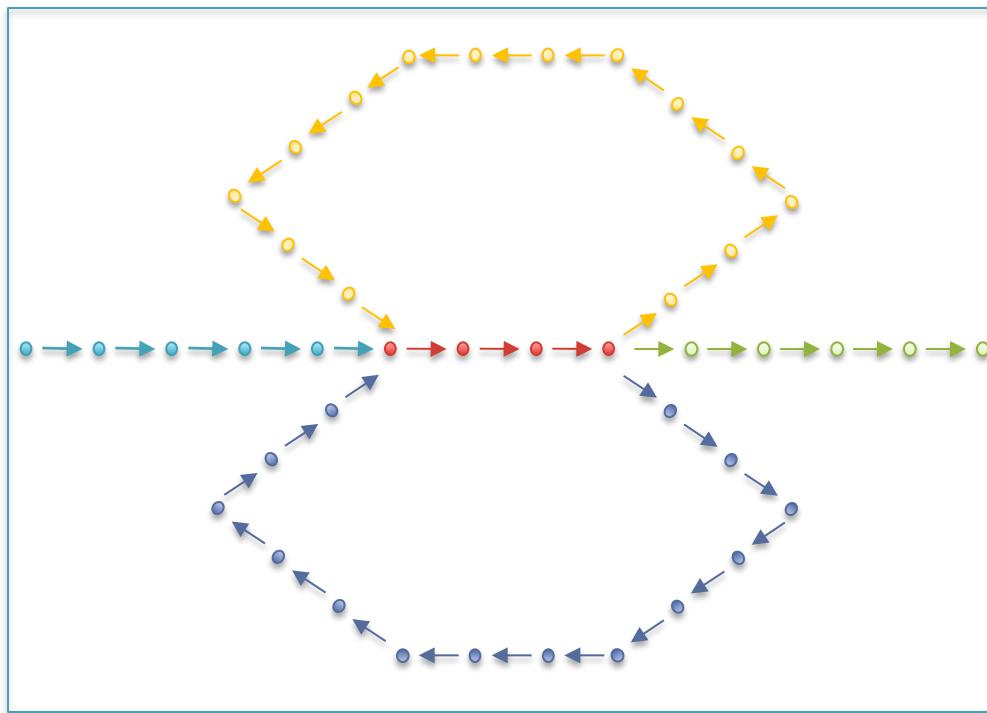
- Considers editing erroneous kmers into trusted kmers in decreasing likelihood
- Includes quality values, nucleotide/nucleotide substitution rate



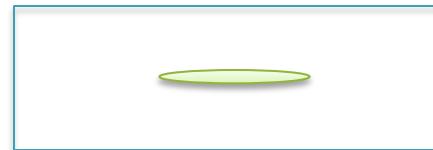
Quake: quality-aware detection and correction of sequencing reads.
Kelley, DR, Schatz, MC, Salzberg SL (2010) *Genome Biology*. 11:R116

Graph Compression

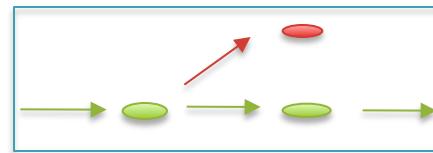
- After construction, many edges are unambiguous
 - Merge together compressible nodes
 - Error correction reduces number of nodes, number of false edges, and allows for longer word size



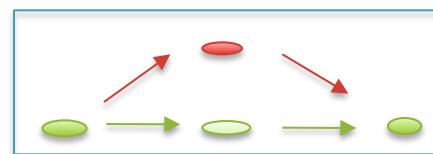
Node Types



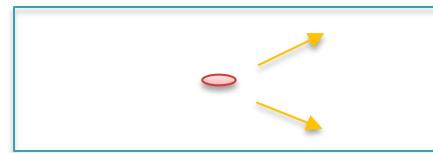
Isolated nodes (10%)



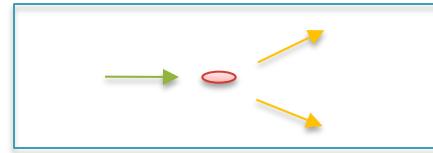
Tips (46%)



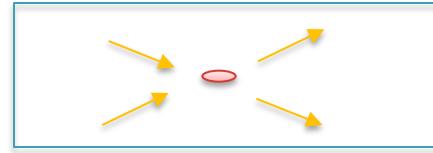
Bubbles/Non-branch (9%)



Dead Ends (.2%)



Half Branch (25%)

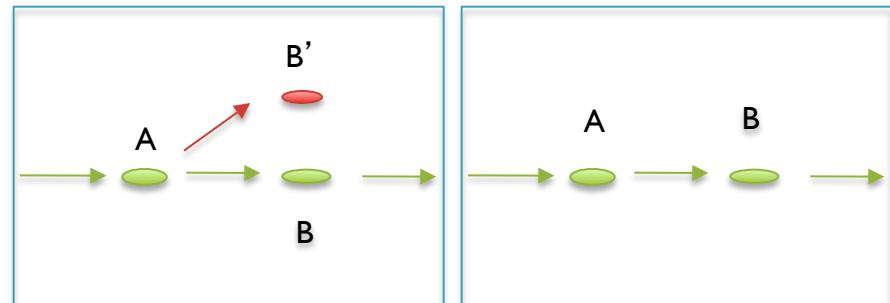


Full Branch (10%)

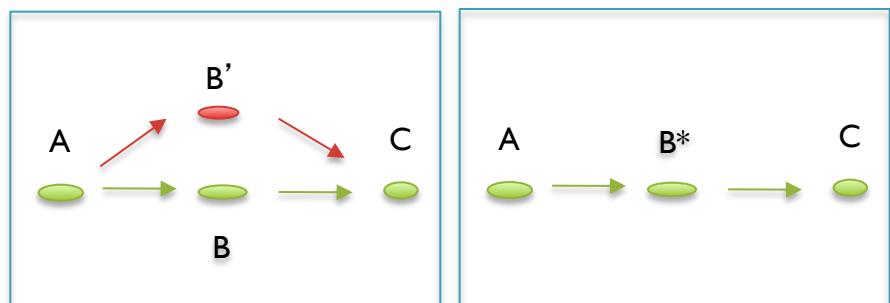
(Chaisson, 2009)

Graph Correction

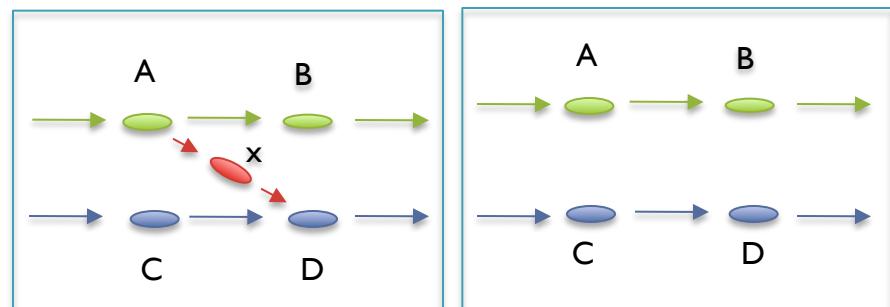
- Errors at end of read
 - Trim off ‘dead-end’ tips



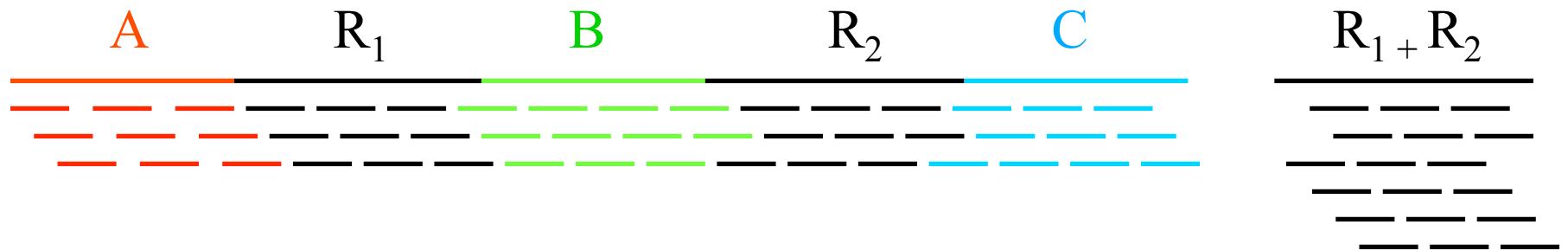
- Errors in middle of read
 - Pop Bubbles



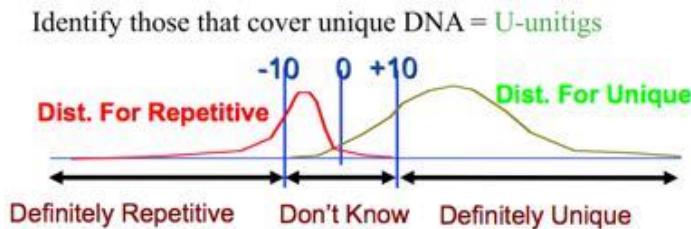
- Chimeric Edges
 - Clip short, low coverage nodes



Coverage Evaluation

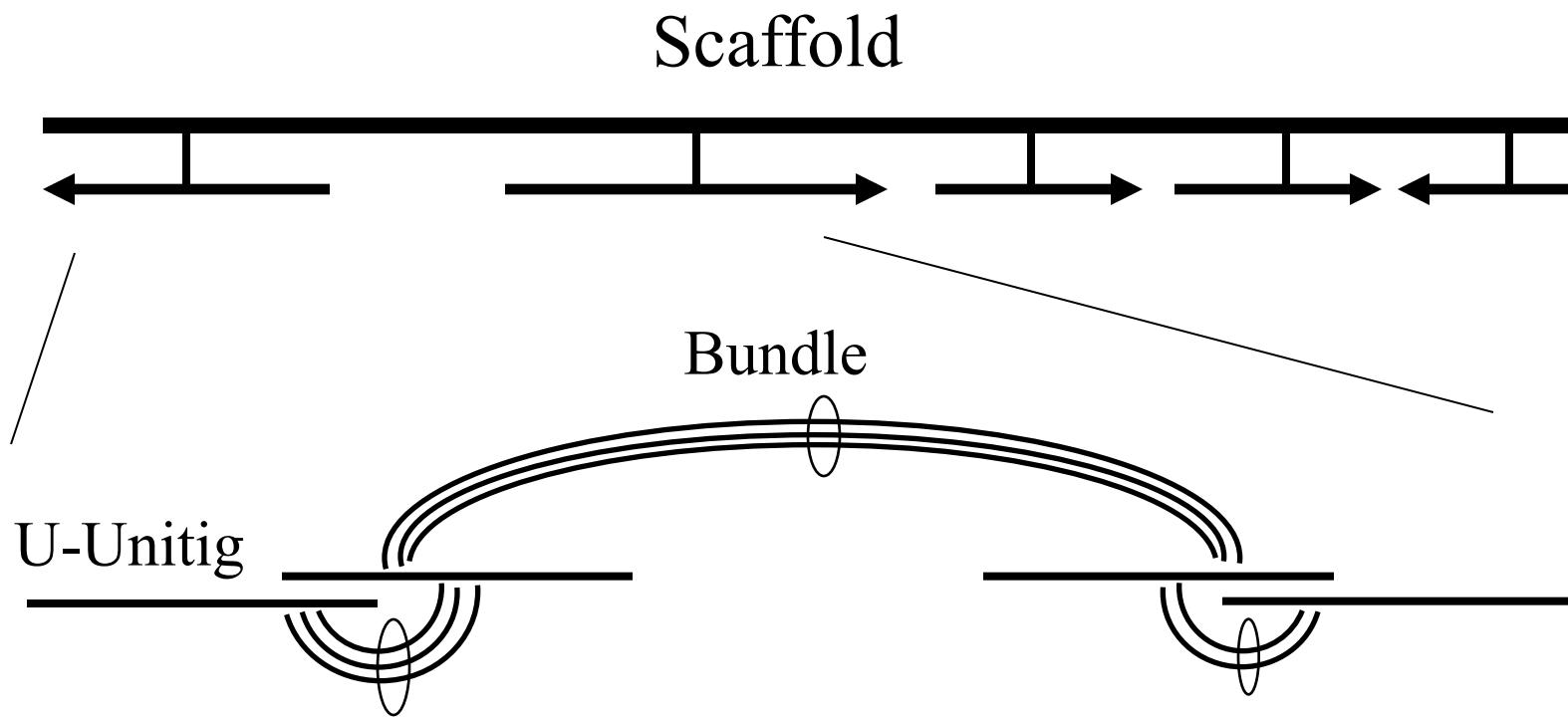


- If n reads are a uniform random sample of the genome of length G , we expect $k = n \Delta/G$ reads to start in a region of length Δ .
 - If we see many more reads than k (if the arrival rate is $> A$) , it is likely to be a collapsed repeat
 - Requires an accurate genome size estimate



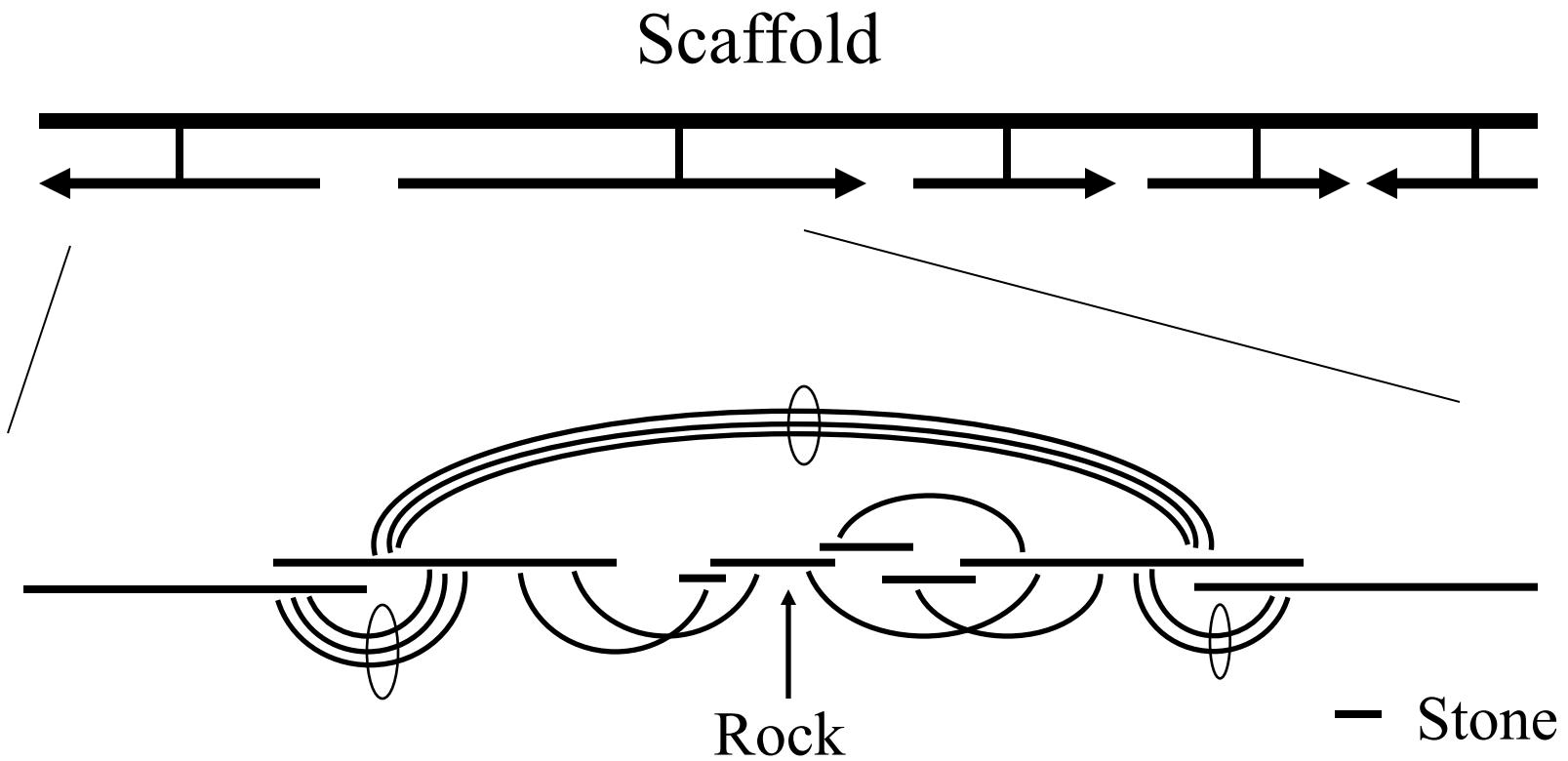
$$A(\Delta, k) = \ln\left(\frac{\Pr(1 - \text{copy})}{\Pr(2 - \text{copy})}\right) = \ln\left(\frac{\frac{(\Delta n/G)^k e^{-\Delta n}}{k!}}{\frac{(2\Delta n/G)^k e^{-2\Delta n}}{k!}}\right) = \frac{n\Delta}{G} - k \ln 2$$

Initial Scaffolding



Create an initial scaffold of basic contigs (“unitigs”) whose coverage indicates they are not repetitive ($A\text{-stat} > 5$).

Repeat Resolution



Then add in remaining repetitive contigs based on their mate relationships allowing repetitive sequences to be placed multiple times.

N50 size

Def: 50% of the genome is in contigs larger than N50

Example:

1 Mbp genome

Contigs: 300k, 100k, 50k, 45k, 30k, 20k, 15k, 15k, 10k,

N50 size = 30 kbp

($300k+100k+50k+45k+30k = 525k \geq 500\text{ kbp}$)

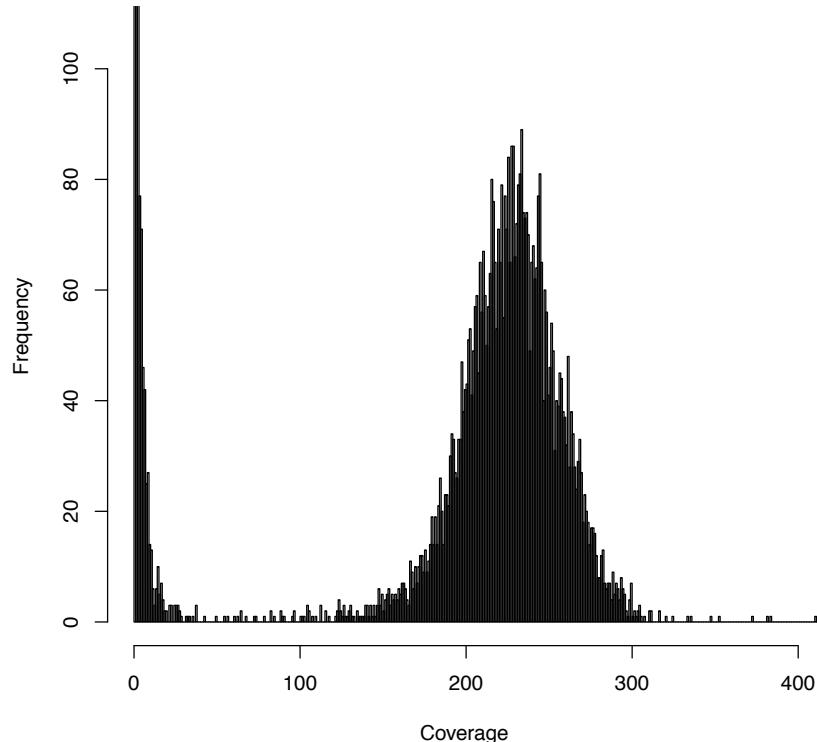
Note:

N50 values are only meaningful to compare when base genome size is the same in all cases

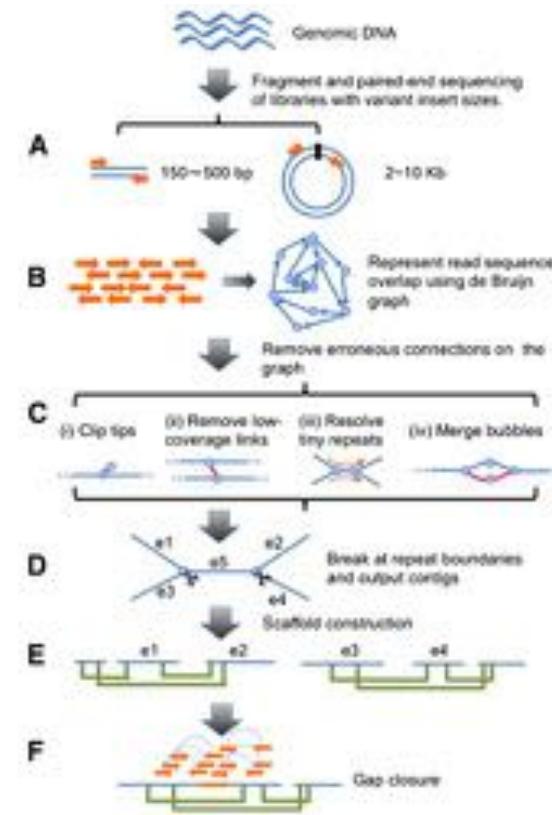
Illumina Sequencing & Assembly

Quake Results

2x76bp @ 275bp
2x36bp @ 3400bp



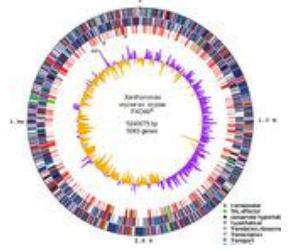
SOAPdenovo Results



Validated	51,243,281	88.5%
Corrected	2,763,380	4.8%
Trim Only	3,273,428	5.6%
Removed	606,251	1.0%

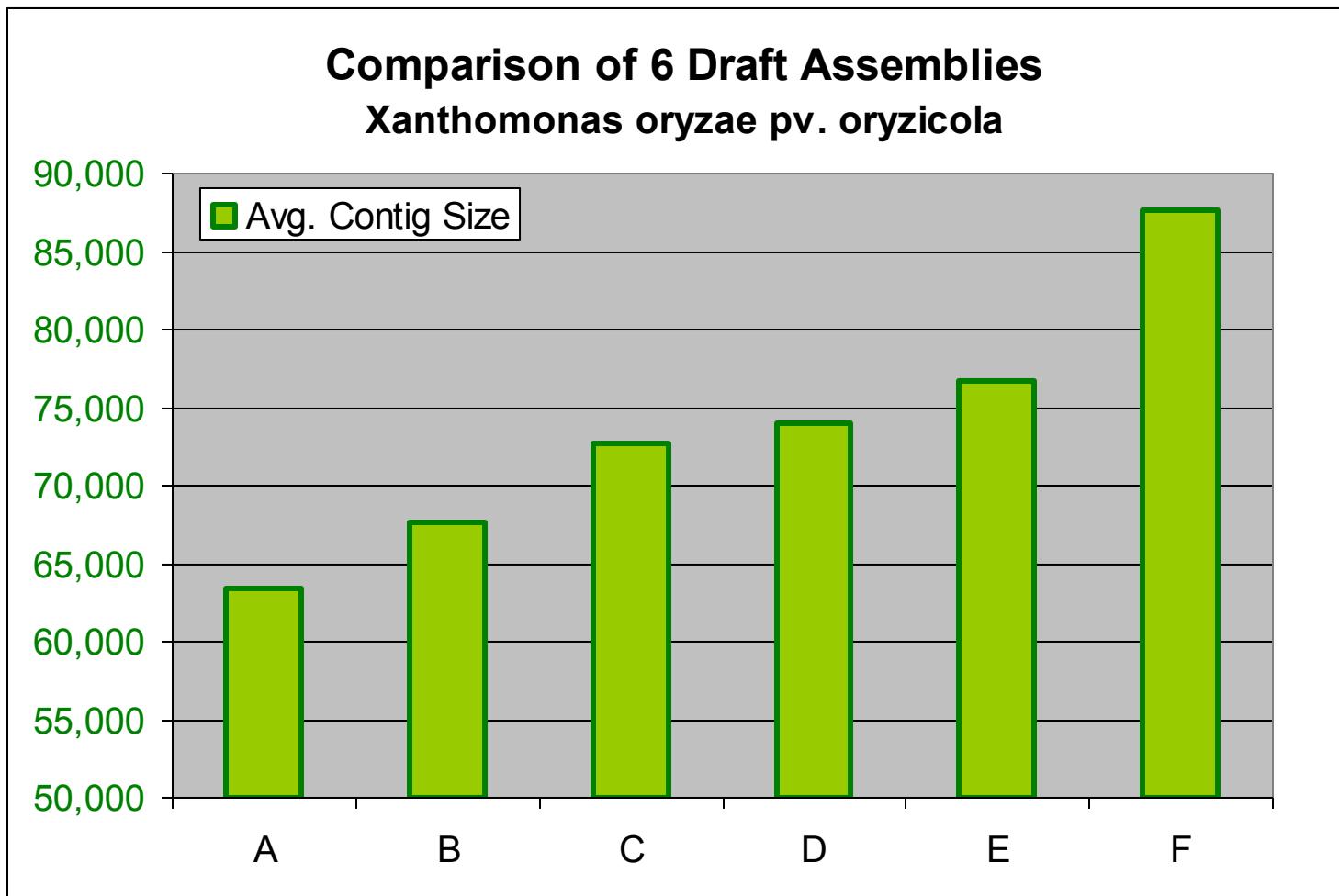
	# ≥ 100bp	N50 (bp)
Scaffolds	2,340	253,186
Contigs	2,782	56,374
Unitigs	4,151	20,772

Assembly realities

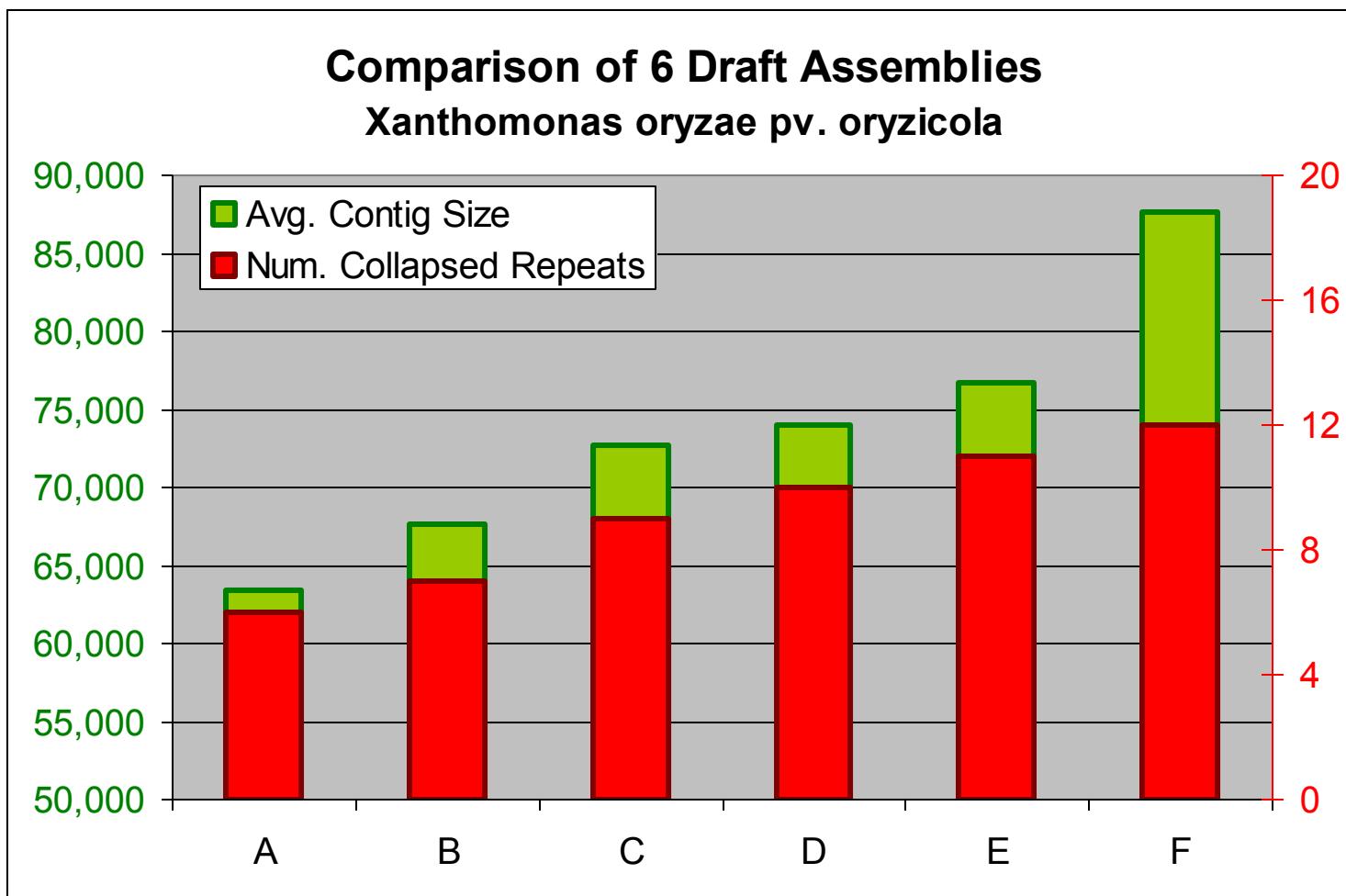


- Contigs are never as large as predicted
 - High coverage is a necessary but not sufficient condition
 - Error correction is required for good assembly
 - Sequencing is basically random, but sequence composition is not
- Repeats control the quality of the assembly
 - Assemblers break contigs at ambiguous repeats
 - Highly repetitive genomes will be highly fragmented
- Assemblers make mistakes
 - Mis-assemblies confuse all downstream analysis
 - Tension between overlap error rate and repeat resolution

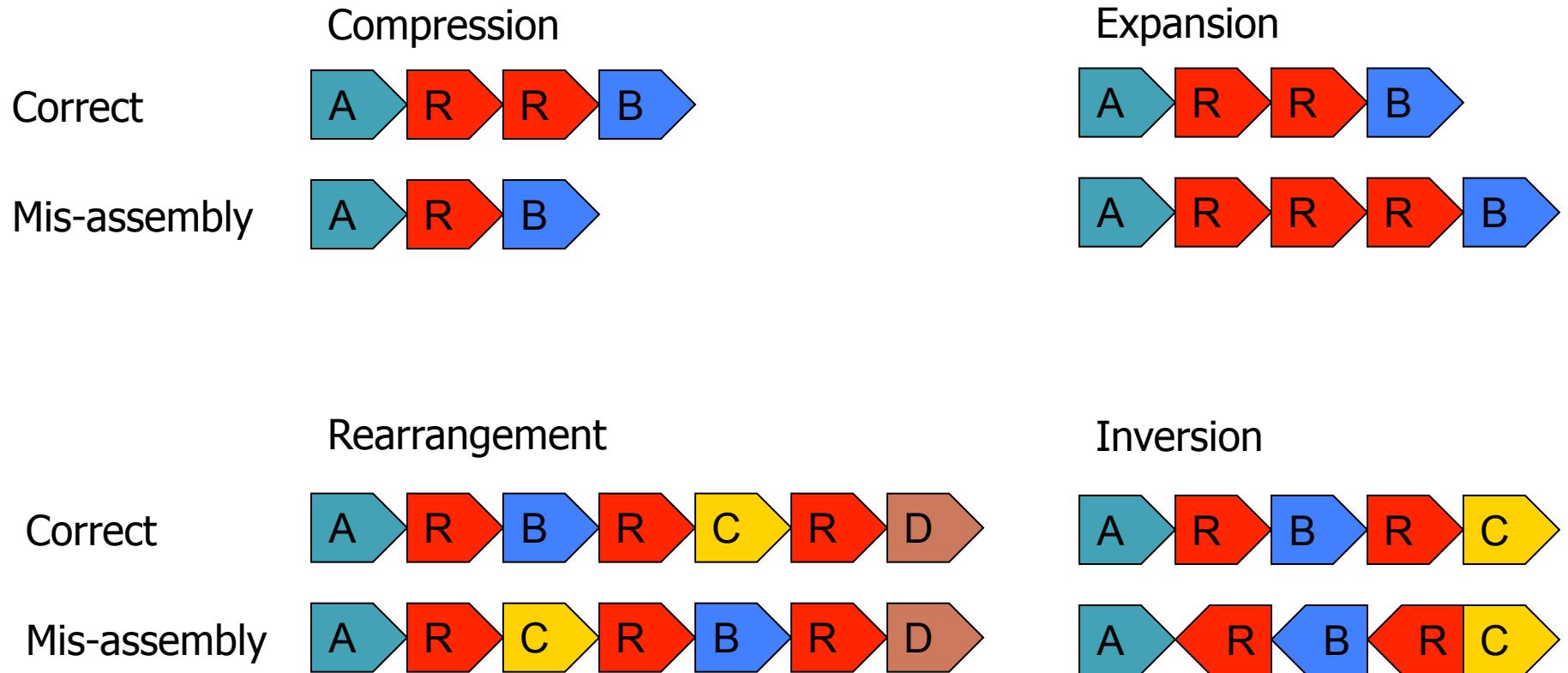
Assembly Evaluation



Assembly Evaluation



Mis-assembly Types



Basic mis-assemblies can be combined into more complicated patterns:
Insertions, Deletions, Giant Hairballs

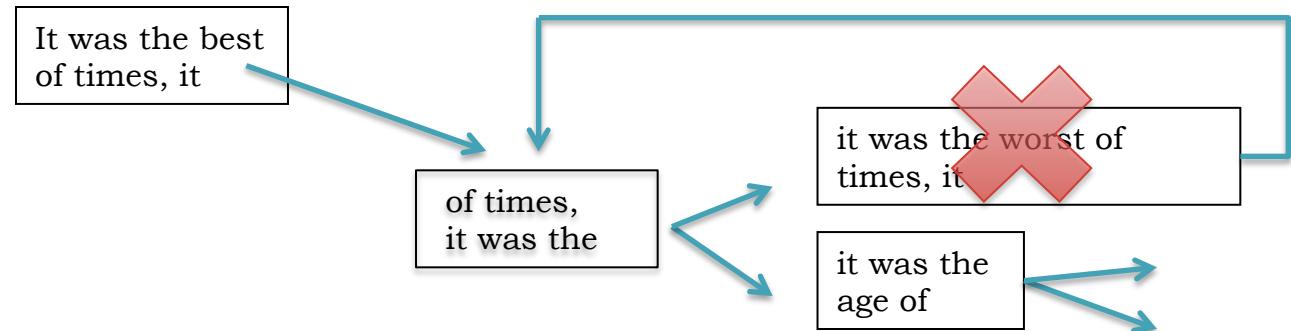
Assembly Forensics



Automatically scan an assembly to locate misassembly signatures for further analysis and correction

Assembly-validation pipeline

1. Evaluate Mate Pairs & Libraries
2. Evaluate Read Alignments
3. Evaluate Read Breakpoints
4. Analyze Depth of Coverage

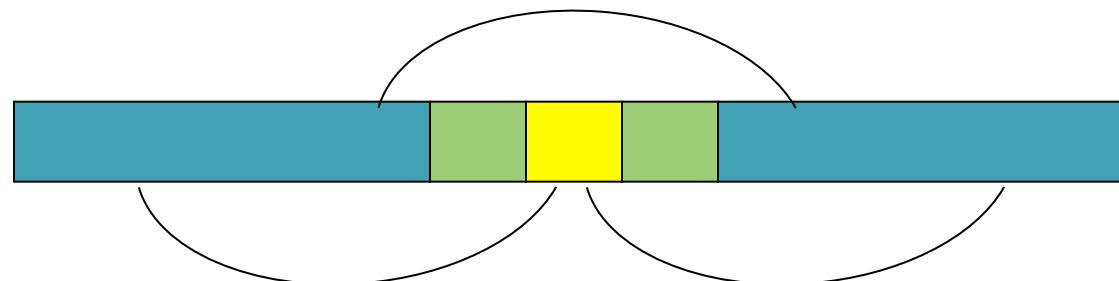


Genome Assembly forensics: finding the elusive mis-assembly.
Phillippy, AM, Schatz, MC, Pop, M. (2008) *Genome Biology* 9:R55.

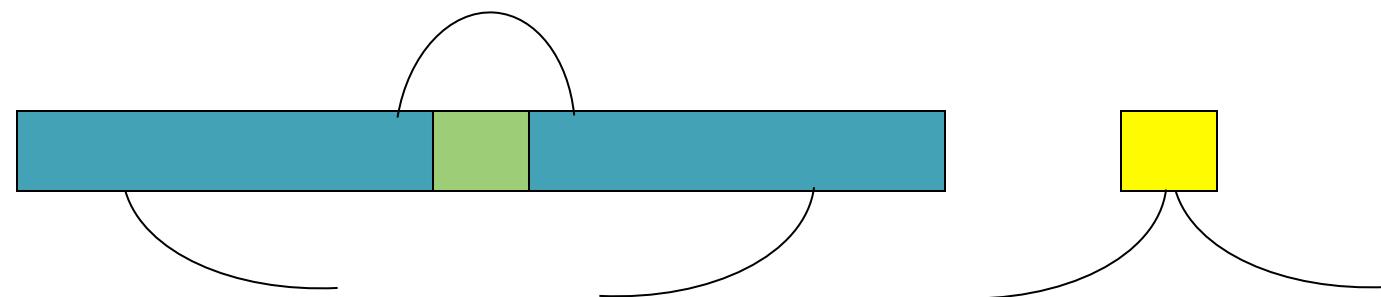
Mate-Happiness: asmQC

- Excision: Skip reads between flanking repeats

- Truth



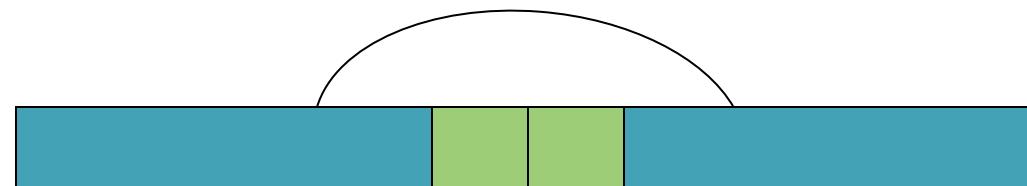
- Misassembly: Compressed Mates, Missing Mates



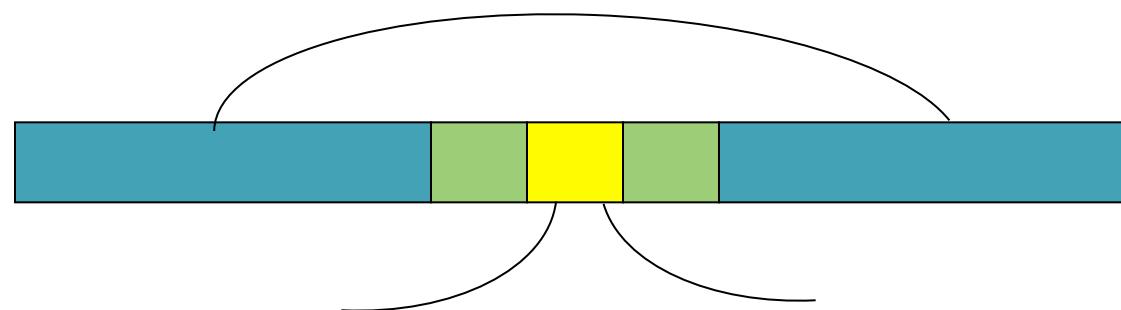
Mate-Happiness: asmQC

- Insertion: Additional reads between flanking repeats

- Truth



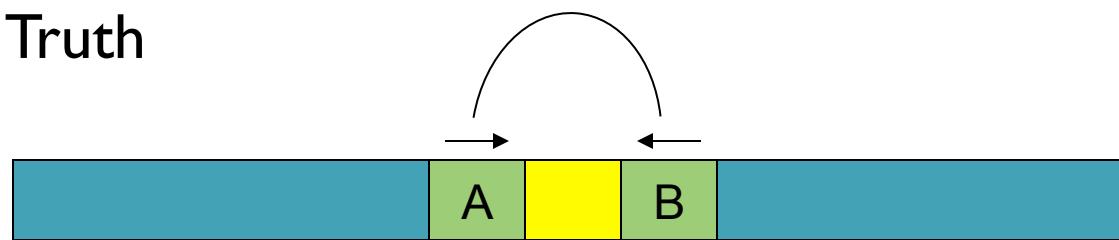
- Misassembly: Expanded Mates, Missing Mates



Mate-Happiness: asmQC

- Rearrangement: Reordering of reads

- Truth

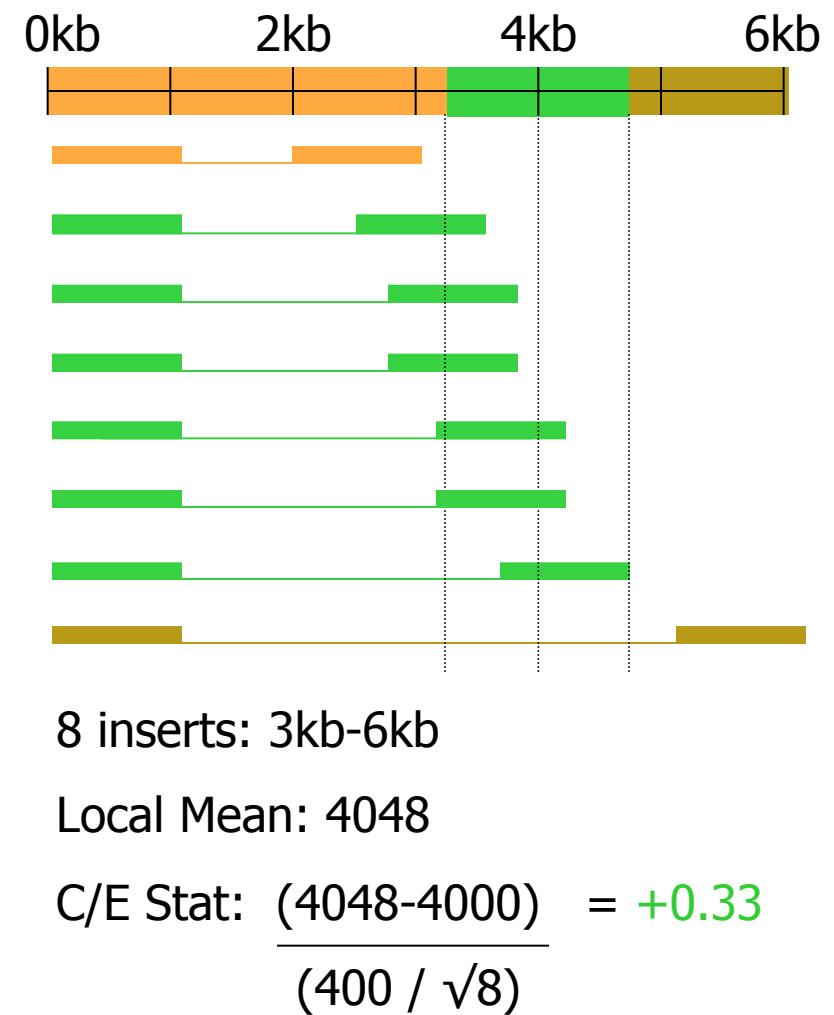
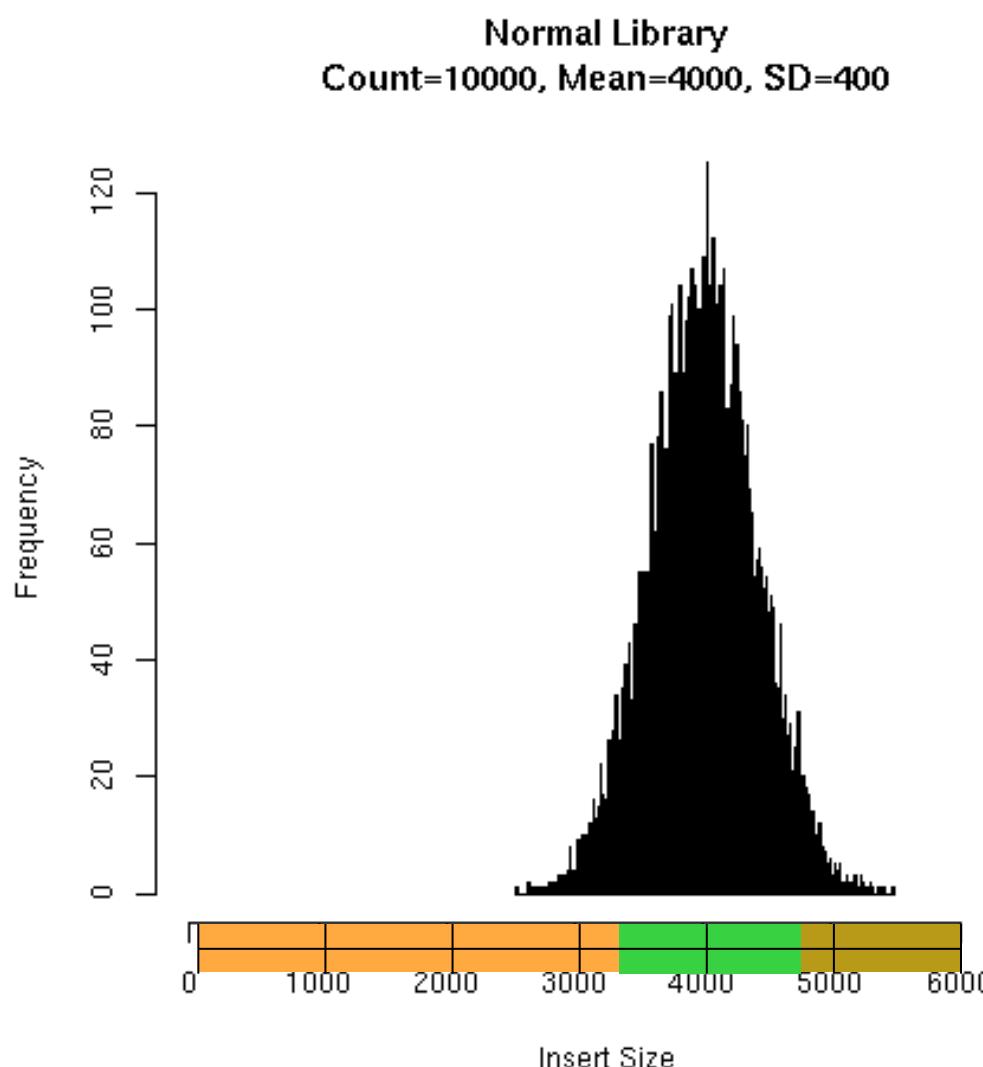


- Misassembly: Misoriented Mates



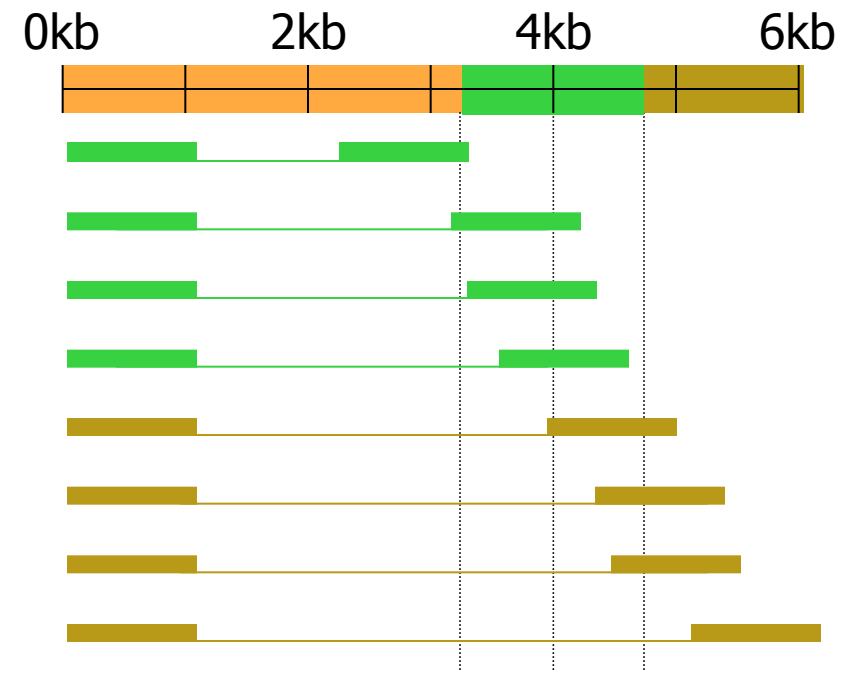
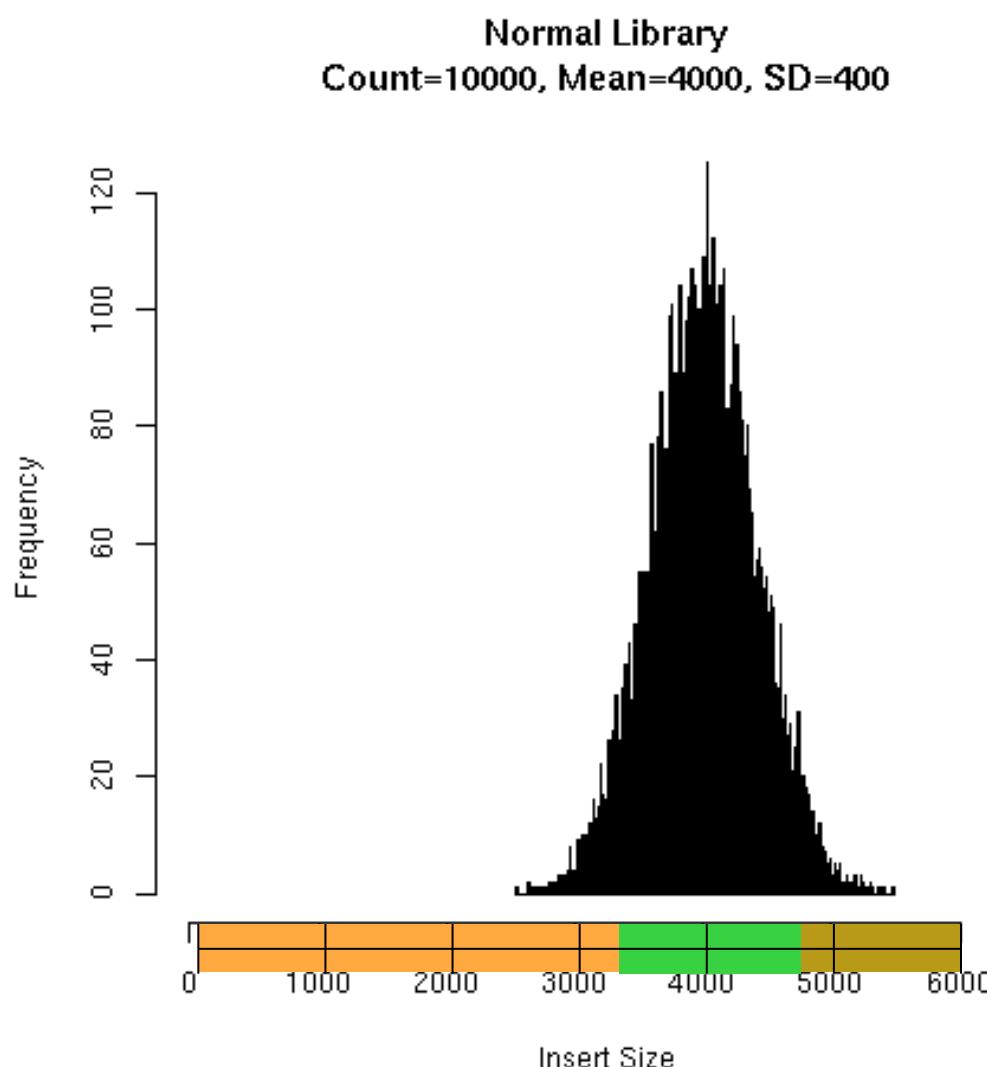
Note: Unhappy mates may also occur for biological or technical reasons.

Sampling the Genome



Near 0 indicates overall happiness

CE Statistic: Expansion



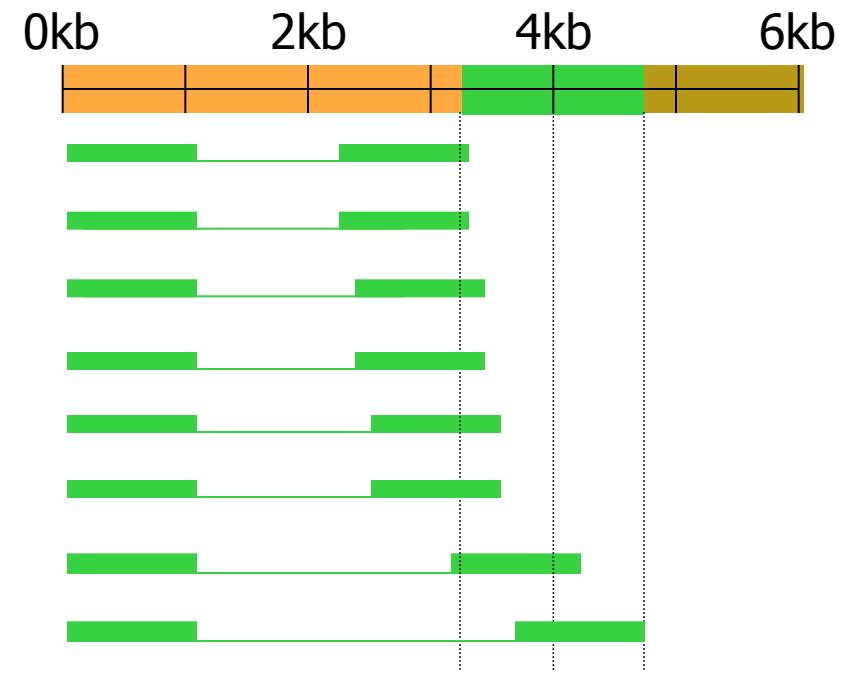
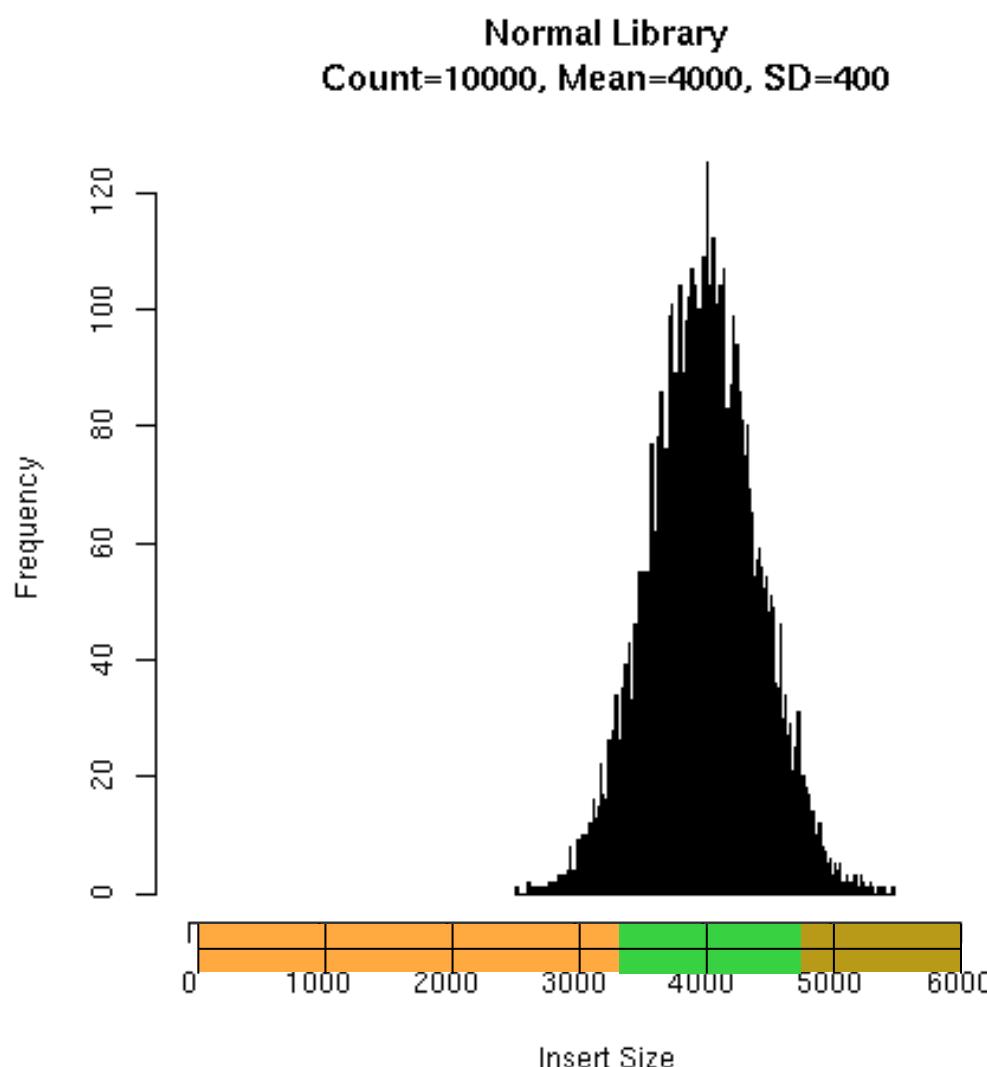
8 inserts: 3.2kb-6kb

Local Mean: 4461

$$\text{C/E Stat: } \frac{(4461 - 4000)}{(400 / \sqrt{8}}) = +3.26$$

C/E Stat ≥ 3.0 indicates Expansion

CE Statistic: Compression



Local Mean: 3488

$$\text{C/E Stat: } \frac{(3488 - 4000)}{(400 / \sqrt{8}} = -3.62$$

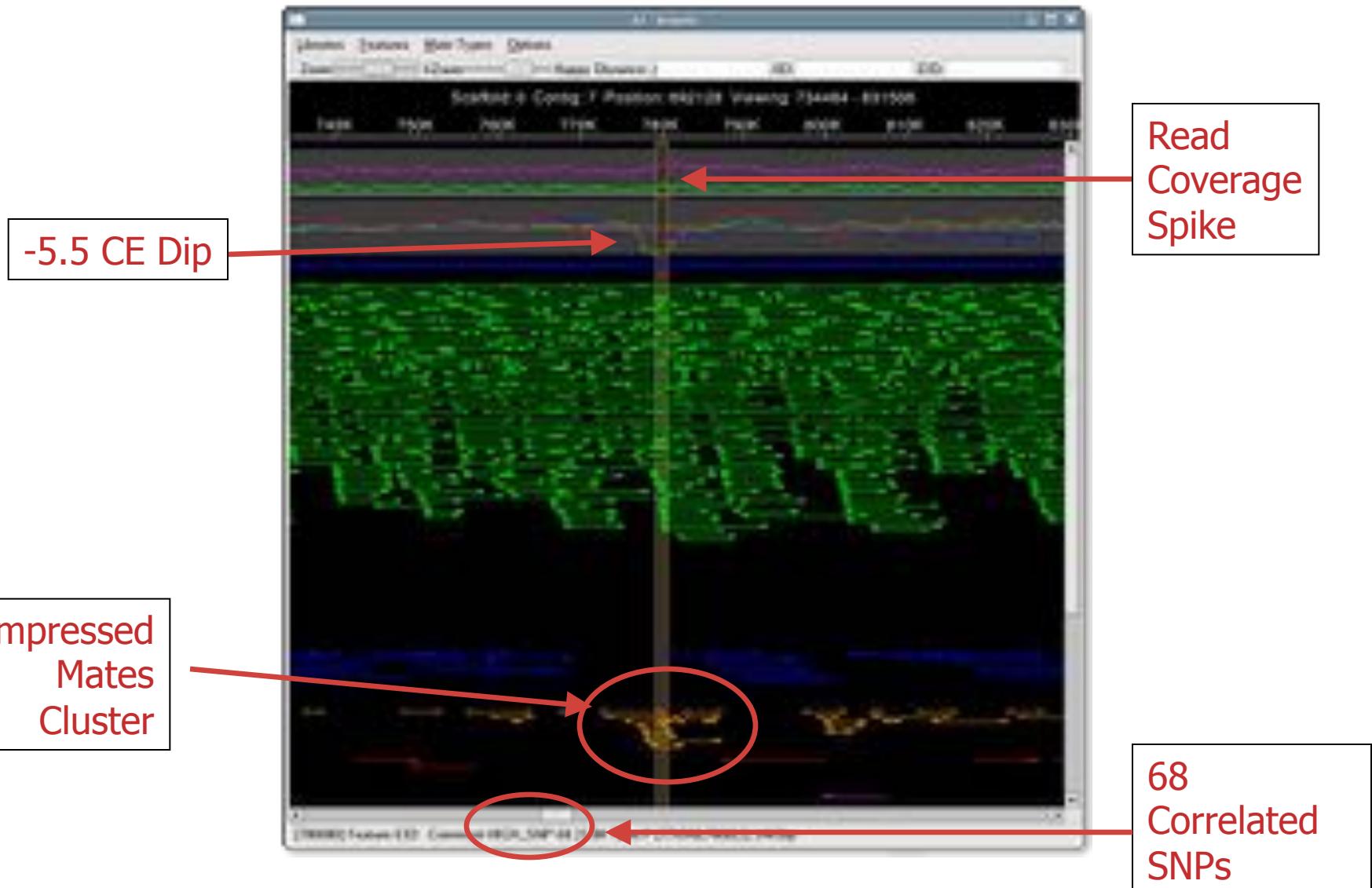
C/E Stat ≤ -3.0 indicates Compression

Read Alignment

- Multiple reads with same conflicting base are unlikely
 - 1x QV 30: 1/1000 base calling error
 - 2x QV 30: 1/1,000,000 base calling error
 - 3x QV 30: 1/1,000,000,000 base calling error
 - Regions of correlated SNPs are likely to be assembly errors or interesting biological events
 - Highly specific metric
 - AMOS Tools: analyzeSNPs & clusterSNPs
 - Locate regions with high rate of correlated SNPs
 - Parameterized thresholds:
 - Multiple positions within 100bp sliding window
 - 2+ conflicting reads
 - Cumulative QV \geq 40 (1/10000 base calling error)

A G C
A G C
A G C
A G C
A G C
A G C
C T A
C T A
C T A
C T A

Collapsed Repeat



Hawkeye: a visual analytics tool for genome assemblies.

Schatz, MC, Phillippy, AM, Shneiderman, B, Salzberg, SL. (2007) Genome Biology 8:R34.

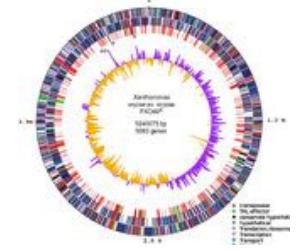
Validation Accuracy

Figure 1

¹Because of administrative and economic difficulties most countries require a minimum fee for the collection of personal information with others.

bottom-up, passive long-term memory of non-existent concepts (Siegler, 1996) and subsequent influence over decisions (Siegler, 1996) are given. A link from these subsections to the following discussion of cross-cultural differences in decision-making is provided.

Assembly Summary



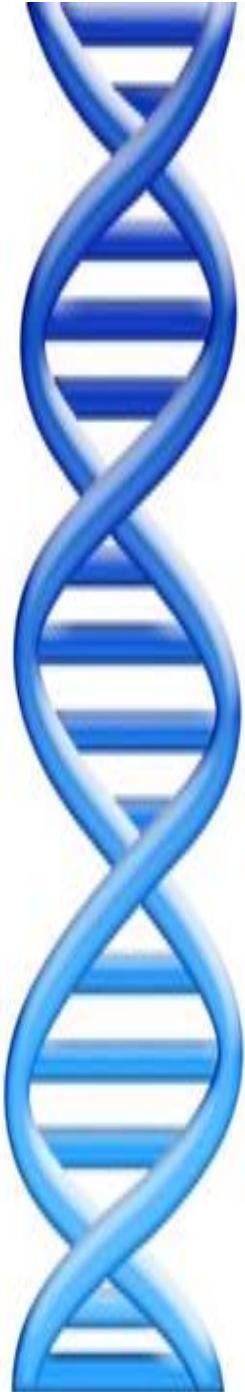
Assembly quality depends on

1. **Coverage**: low coverage is mathematically hopeless
2. **Repeat composition**: high repeat content is challenging
3. **Read length**: longer reads help resolve repeats
4. **Error rate**: errors reduce coverage, obscure true overlaps

- Assembly is a hierarchical, starting from individual reads, build high confidence contigs/unitigs, incorporate the mates to build scaffolds
 - Extensive error correction is the key to getting the best assembly possible from a given data set
- Watch out for collapsed repeats & other misassemblies
 - Globally/Locally reassemble data from scratch with better parameters & stitch the 2 assemblies together

Break





Outline

Part 1: Schatz Lab Overview

Part 2: Sequence Alignment

Part 3: Genome Assembly

Part 4: Parallel & Cloud Computing

- Milestones in DNA Sequencing
- Hadoop & Cloud Computing
- Sequence Analysis in the Clouds

Milestones in DNA Sequencing

1970

1980

1990

2000

2010

核酸酸 1977 年卷 107 号

107

articles

Nucleotide sequence of bacteriophage ϕ X174 DNA

E. Sanger, G. M. Air*, B. G. Barrell, N. L. Botting†, A. R. Coulson, Z. C. Fiddes,
C. A. Hoadley‡, P. M. Slocombe§ & M. Smith*

*MRC Laboratory of Molecular Biology, 20 Cambridge Embankment, London WC2R 3PD, U.K.

In 1974 a sequence for the genome of bacteriophage ϕ X174 of approximately 5,375 nucleotides has been determined using the rapid and simple 'gap and vector' method. The sequence identifies some of the features responsible for the production of the proteins in the other three genes of the genome, including lacunase and haemocyanin genes (for the proteins and RNA), and two genes of unknown function (for the genes of protein and RNA). Two pairs of genes are coded by the same region of ϕ X174 using different reading frames.

The genome of bacteriophage ϕ X174 is a single-stranded, linear DNA of approximately 5,375 nucleotides coding for six known proteins. The code of base pairs on bacteriophage is nearly identical^{1,2} to λ -DNA^{3,4} (Table 1). Class I, II and III code for structural proteins of the virus capsid, and Class IV has twelve functional genes coding for a total mean protein

content 10% of the nucleic acid sequence in the genome and, in certain enzymes, are also involved in that a protease fragment can be formed and measured. Only one reading frame is found. The complementary strand contains no sequence which would disrupt the otherwise reading frame dictated by the sequence of the gene of protein⁵. (See also 1.2.3.1.)

In the early sequencing techniques used, protein synthesis from DNA, conversion into living phage⁶, and labourious measurement of each individual protein component, or part of the otherwise functioning virion, was very slow and difficult. The introduction of the 'gap and vector' technique⁷ and the 'double-stranded phage'⁸ have greatly simplified the work. The double-stranded phage technique⁹ facilitated the physical separation of the double-stranded phage. This was achieved by a system of two nested filters, giving rise to two main fractions: 'bottom', containing proteins and 'bottom', difficult to pass over and so

1977

Sanger et al.

Ist Complete Organism
Bacteriophage ϕ X174
5375 bp



Radioactive Chain Termination

5000bp / week / person

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

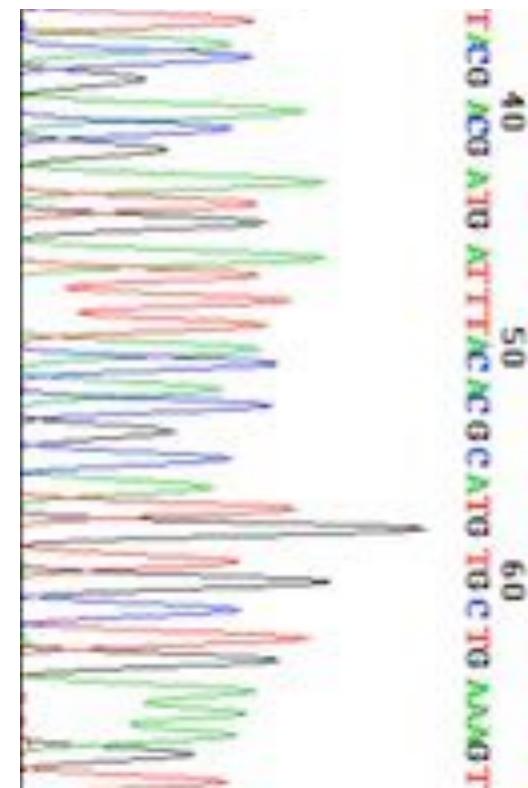
Milestones in DNA Sequencing



1987

Applied Biosystems markets the ABI 370 as the first automated sequencing machine

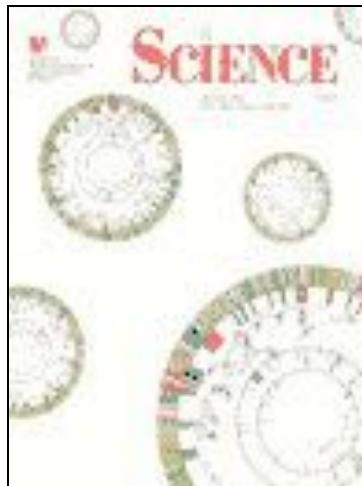
http://commons.wikimedia.org/wiki/File:370A_automated_DNA_sequencer.jpg



Fluorescent Dye Termination
350bp / lane x 16 lanes =
5600bp / day / machine

<http://www.answers.com/topic/automated-sequencer>

Milestones in DNA Sequencing



1995
Fleischmann *et al.*
1st Free Living Organism
TIGR Assembler. 1.8Mbp



2000
Myers *et al.*
1st Large WGS Assembly.
Celera Assembler. 116 Mbp



2001
Venter *et al.*,
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 DNA Sequencing



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



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

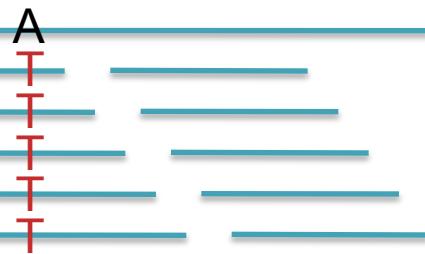


2008
ABI / Life Technologies
SOLiD Sequencing
Current Specs (5500xl):
5B 75bp reads / run =
30Gbp / day

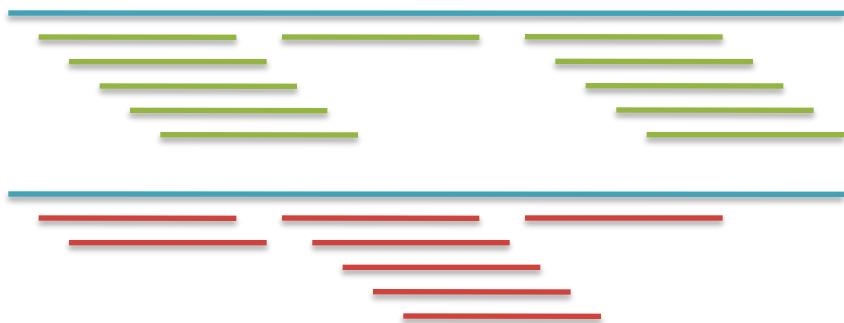
Second Generation Sequencing Applications



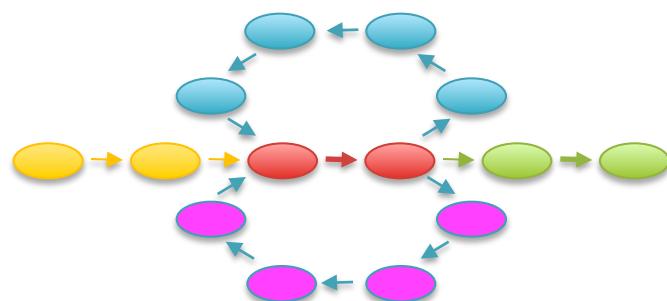
Alignment & Variations



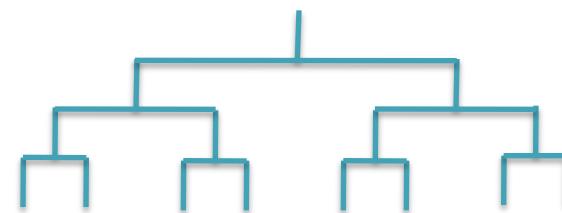
Differential Analysis



De novo Assembly



Phylogeny & Evolution



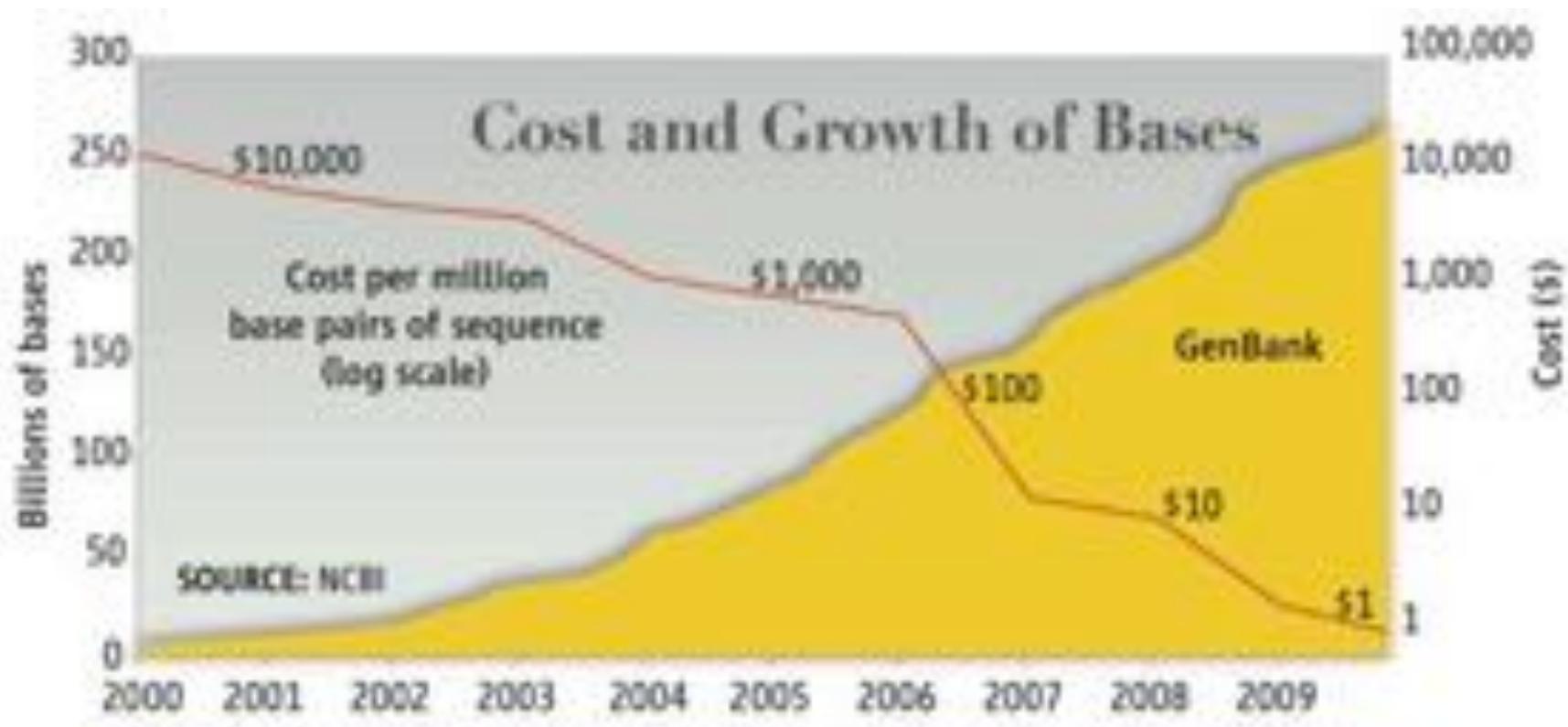
Sequencing Centers



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

DNA Data Tsunami

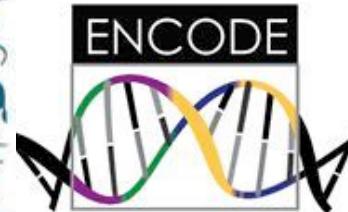
*Current world-wide sequencing capacity exceeds 13Pbp/year
and is growing at 5x per year!*



"Will Computers Crash Genomics?"

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

Genomics and Parallel Computing



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



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 2.0¢ / core / hour
- Simple Storage Service (S3)
 - Scalable data storage
 - 15¢ / GB monthly fee
- Plus many others



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



Amazon Machine Images



- A few Amazon sponsored images
 - Suse Linux, Windows
- Many Community Images & Appliances
 - CloudBioLinux: Genomics Appliance
 - Crossbow: Hadoop, Bowtie, SOAPsnp
 - Galaxy: CloudMan
- Build your own
 - Completely customize your environment
 - Your results could be totally reproducible

Getting Started

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

The screenshot shows a web browser displaying the 'Amazon Elastic Compute Cloud Getting Started Guide (API Version 2013-06-06)' at <http://docs.amazonwebservices.com/AWSEC2/latest/GettingStartedGuide/>. The page has a blue header bar with the Amazon logo. On the left, there's a sidebar titled 'Get Started with EC2' containing links: 'Sign up for EC2', 'Launch an instance', 'Associate an IAM user', 'Ssh access', 'Connect to your Windows instance', 'Terminate your instance', 'Delete your instance', 'Where will my IP be from now?', 'Post-Instance Feedback', and 'Leave this guide'. The main content area has a title 'Get Started with EC2' and a paragraph explaining what EC2 is. It includes a flowchart showing the process: Sign up for EC2 → Launch instance → Connect to Linux instance or Connect to Windows instance → Terminate instance. Below the flowchart is a note about getting started and a 'Get Started' button.

Amazon Elastic Compute Cloud Getting Started Guide (API Version 2013-06-06)

Amazon

Get Started with EC2

Amazon Elastic Compute Cloud (Amazon EC2) is a web service that enables you to launch and manage Linux/Windows server instances in Amazon's datacenters. You can get started with Amazon EC2 by following the tasks shown in the following diagram. You'll primarily use the EC2 Management Console, a plain web-based user interface.

```
graph LR; A[Sign up for EC2] --> B[Launch instance]; B --> C[Connect to Linux instance]; C --> D[Connect to Windows instance]; D --> E[Terminate instance]
```

Free guide walks you through launching and connecting to your first Amazon EC2 instance. To start, click the following 'Get Started' button.

[Get Started](#)

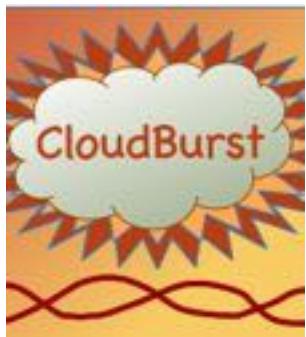
Hadoop MapReduce

<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)
 - 946 PB 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*

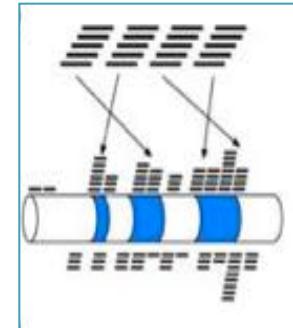
<http://cloudburst-bio.sf.net>

(Schatz, 2009)

Myrna

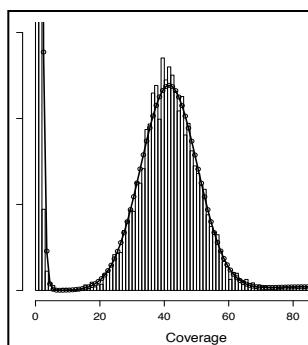
Cloud-scale differential gene expression for RNA-seq

*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.ccb.umd.edu/software/quake/>

(Kelley, Schatz,
Salzberg, 2010)

Genome Indexing

Rapid Parallel Construction
of Genome Index

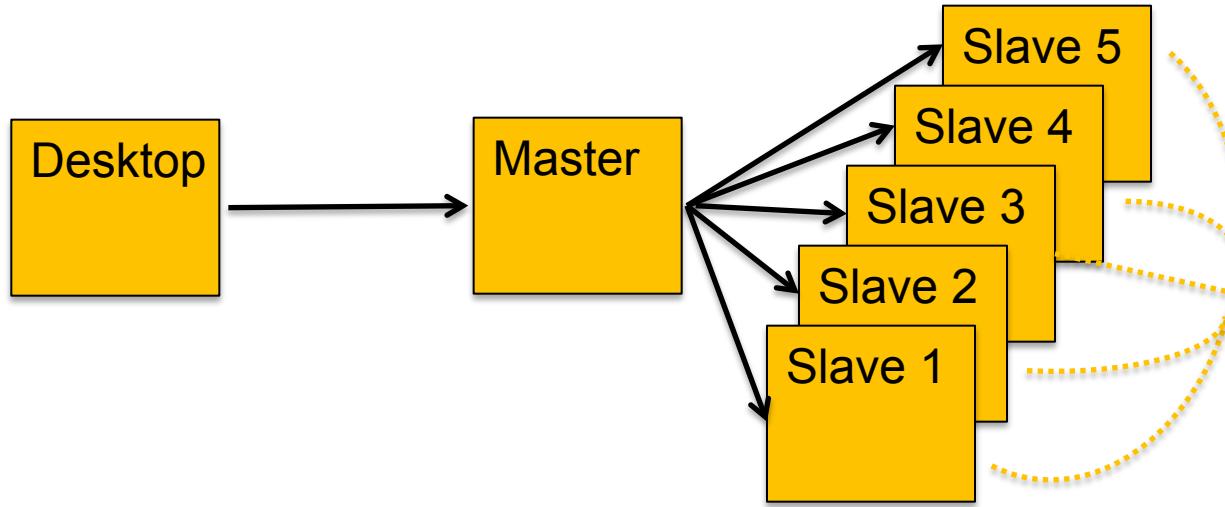
*Construct the BWT of
the human genome in 9 minutes*

\$GATTACA
A\$GATTAC
ACA\$GATT
ATTACA\$G
CA\$GATTA
GATTACA£
TACA\$GA T
TTACA\$GA

(Menon,
Bhat, Schatz, 2011*)

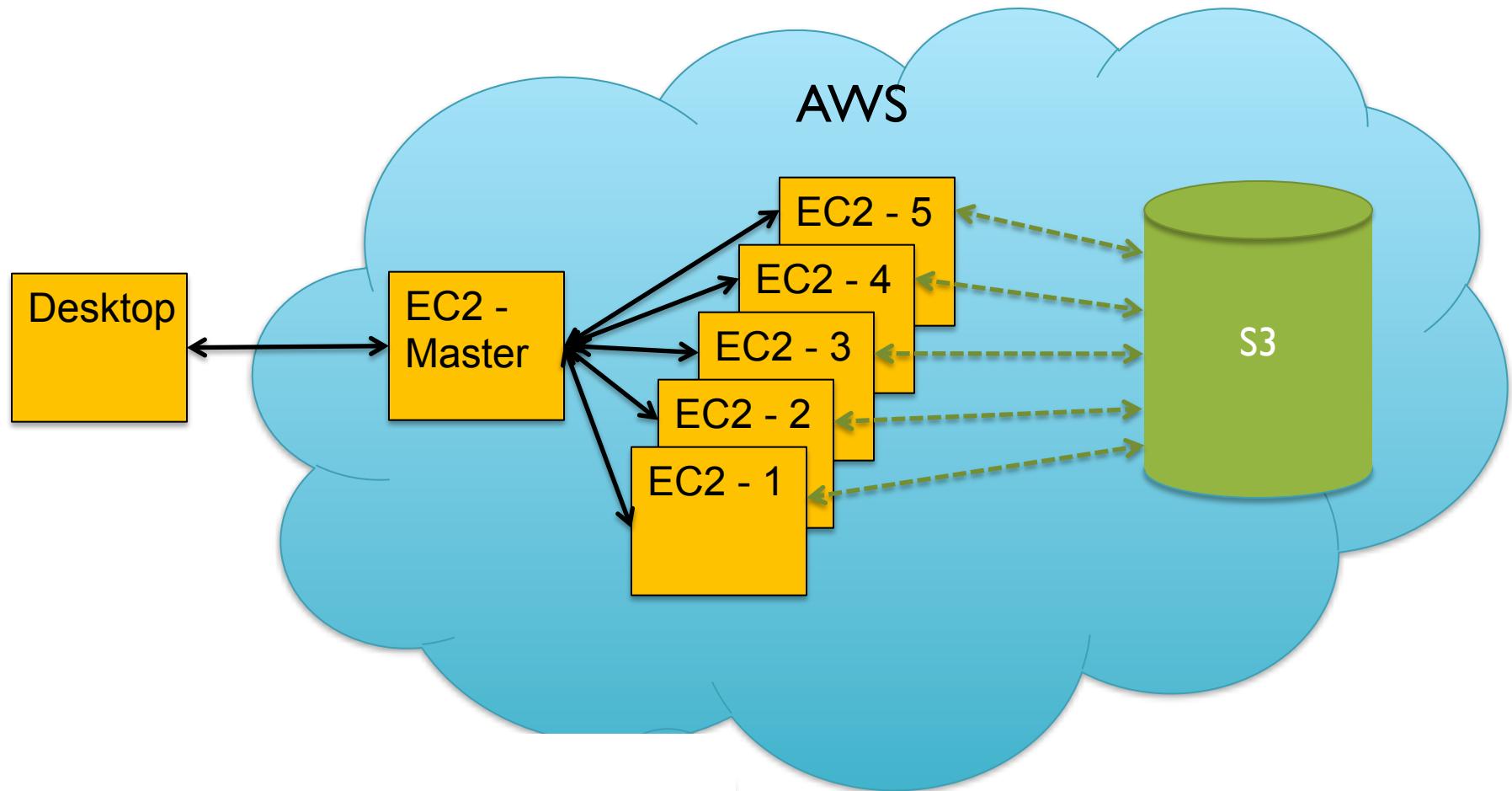
[http://code.google.com/p/
genome-indexing/](http://code.google.com/p/genome-indexing/)

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
 - Provides many disks in addition to many cores

Hadoop on AWS



If you don't have 1000s of machines, rent them from Amazon

- After machines spool 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

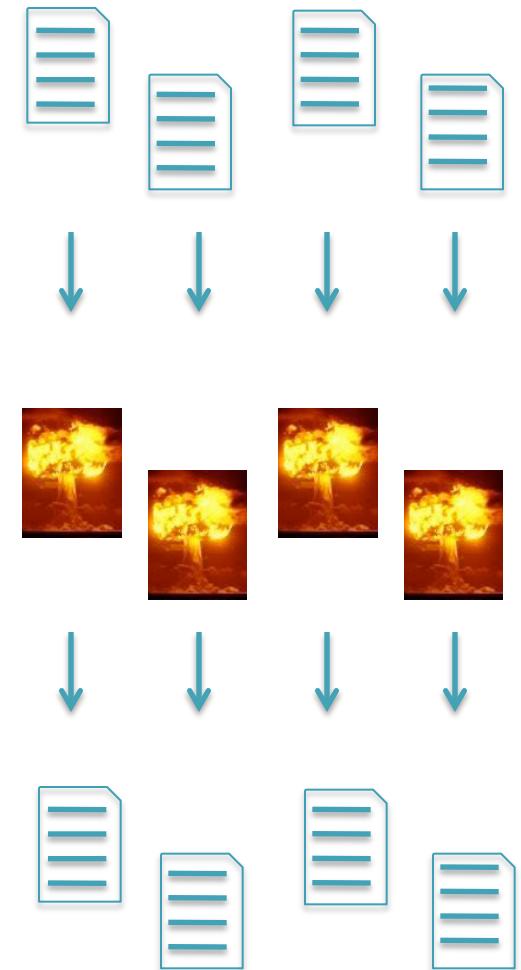
Tightly Coupled



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

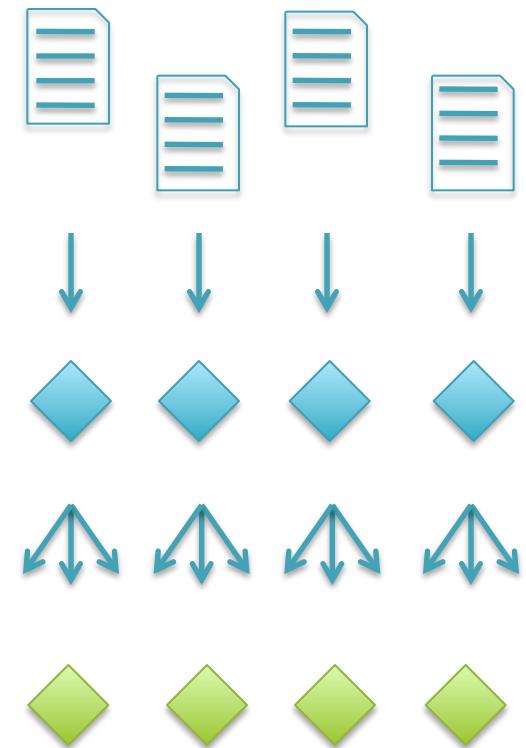


Elementary School Dance



2. Loosely Coupled

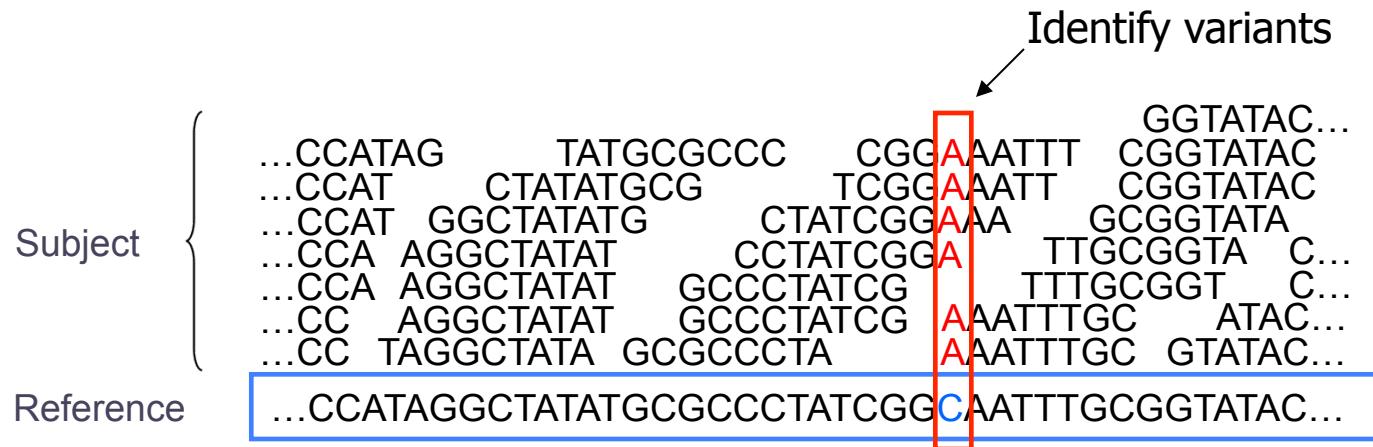
- Divide and conquer
 - Independently process many items
 - Group partial results
 - Scan partial results into final answer
- Challenges
 - Batch computing challenges
 - + Shuffling of huge datasets
- Technologies
 - Hadoop, Elastic MapReduce, Dryad
 - Parallel Databases



Junior High Dance



Short Read Mapping



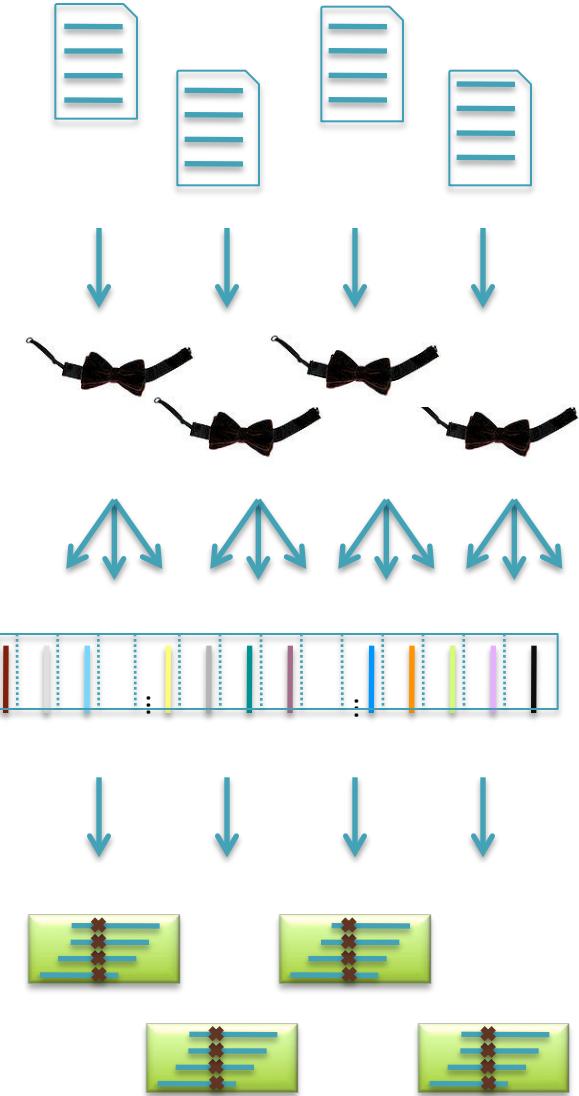
- Given a reference and many subject reads, report one or more “good” end-to-end alignments per alignable read
 - Find where the read most likely originated
 - Fundamental computation for many assays
 - Genotyping RNA-Seq Methyl-Seq
 - Structural Variations Chip-Seq Hi-C-Seq
 - Desperate need for scalable 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>

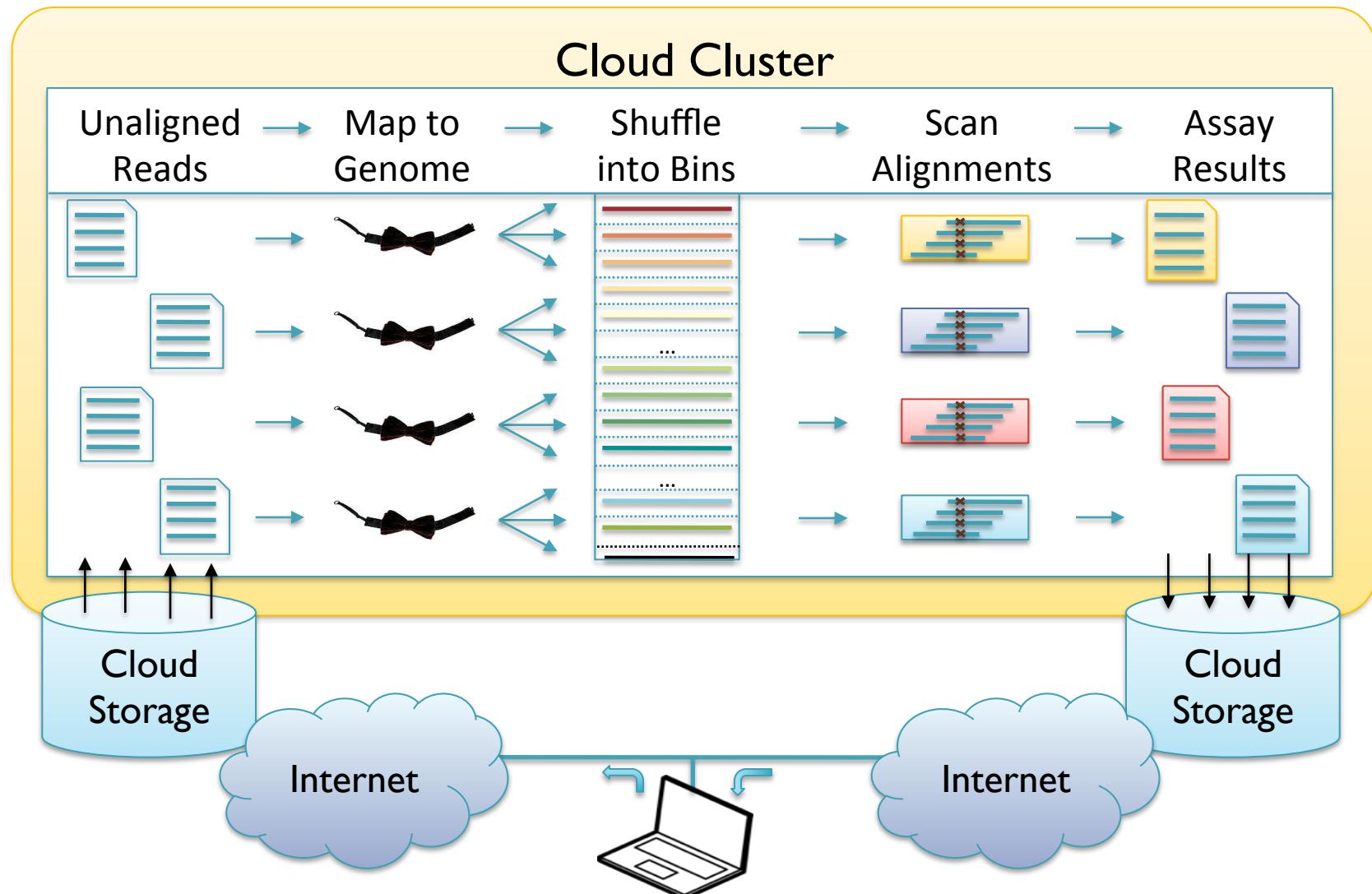
Asian Individual Genome			
Data Loading	3.3 B reads	106.5 GB	\$10.65
Data Transfer	1h :15m	40 cores	\$3.40
Setup	0h :15m	320 cores	\$13.94
Alignment	1h :30m	320 cores	\$41.82
Variant Calling	1h :00m	320 cores	\$27.88
End-to-end	4h :00m		\$97.69

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

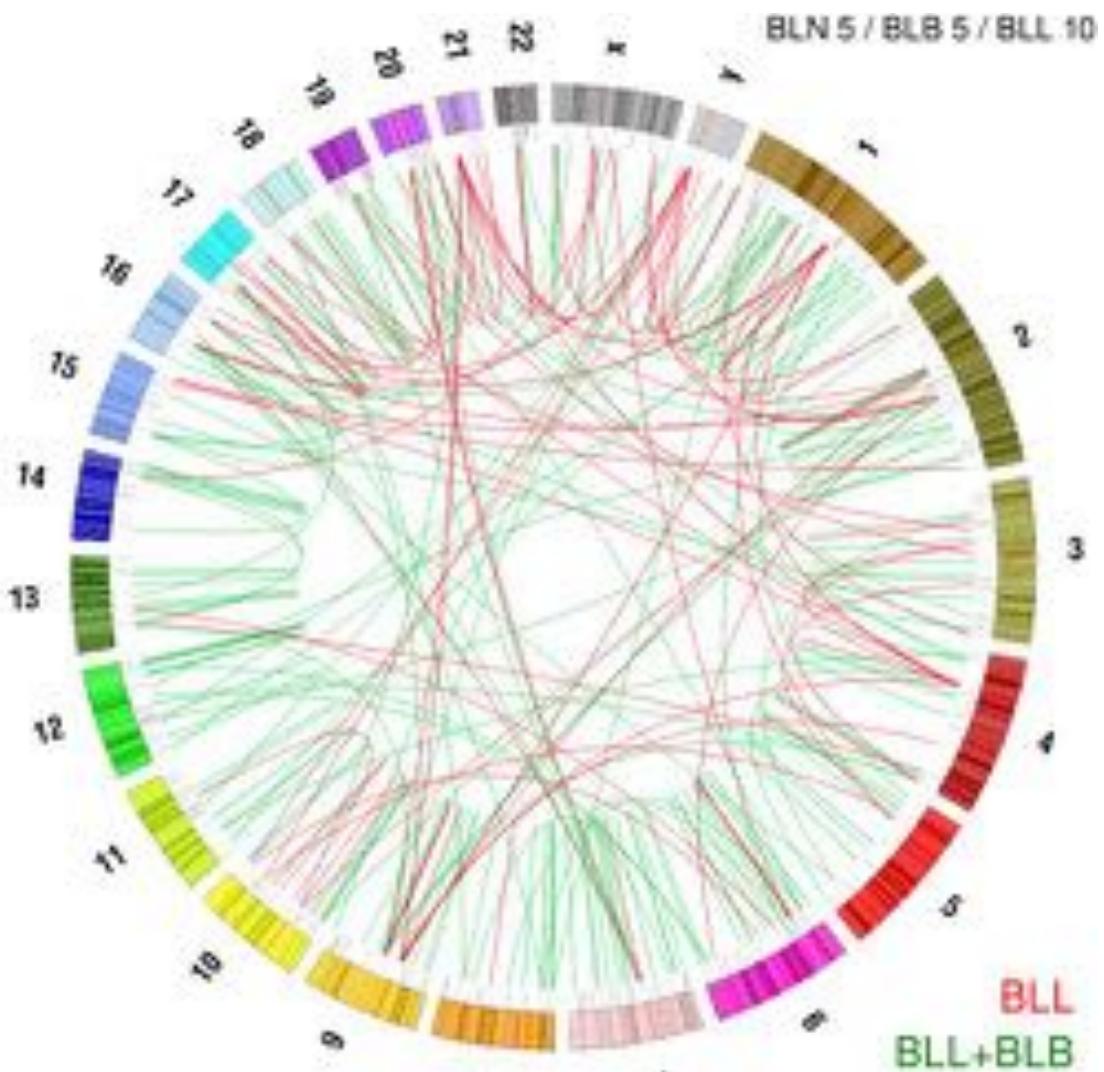
Jnomics Structural Variations

Circos plot of high confidence SVs specific to esophageal cancer sample

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

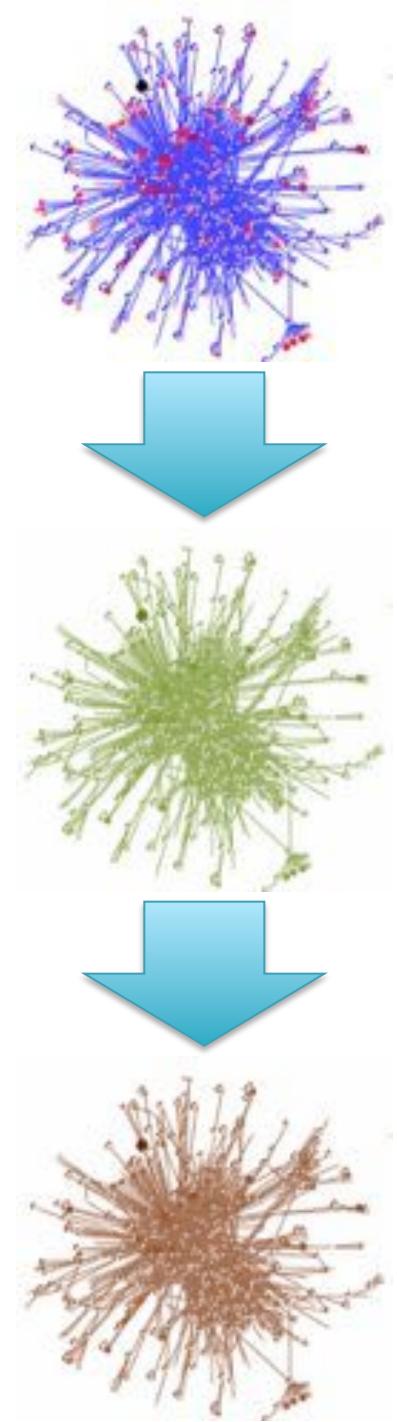
Detailed analysis of disrupted genes and fusion genes in progress

- Preliminary analysis shows many promising hits to known cancer genes



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

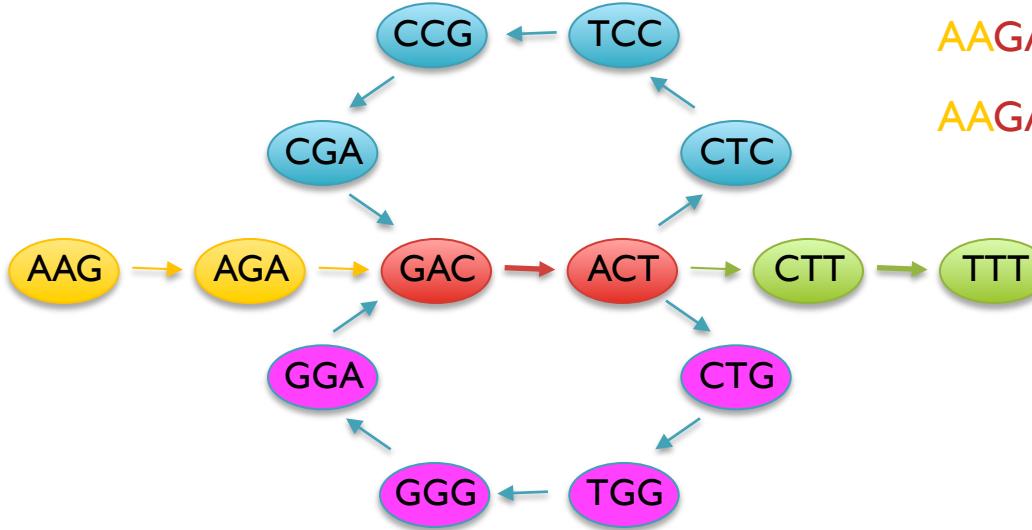


Short Read Assembly

Reads

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

de Bruijn Graph



Potential Genomes

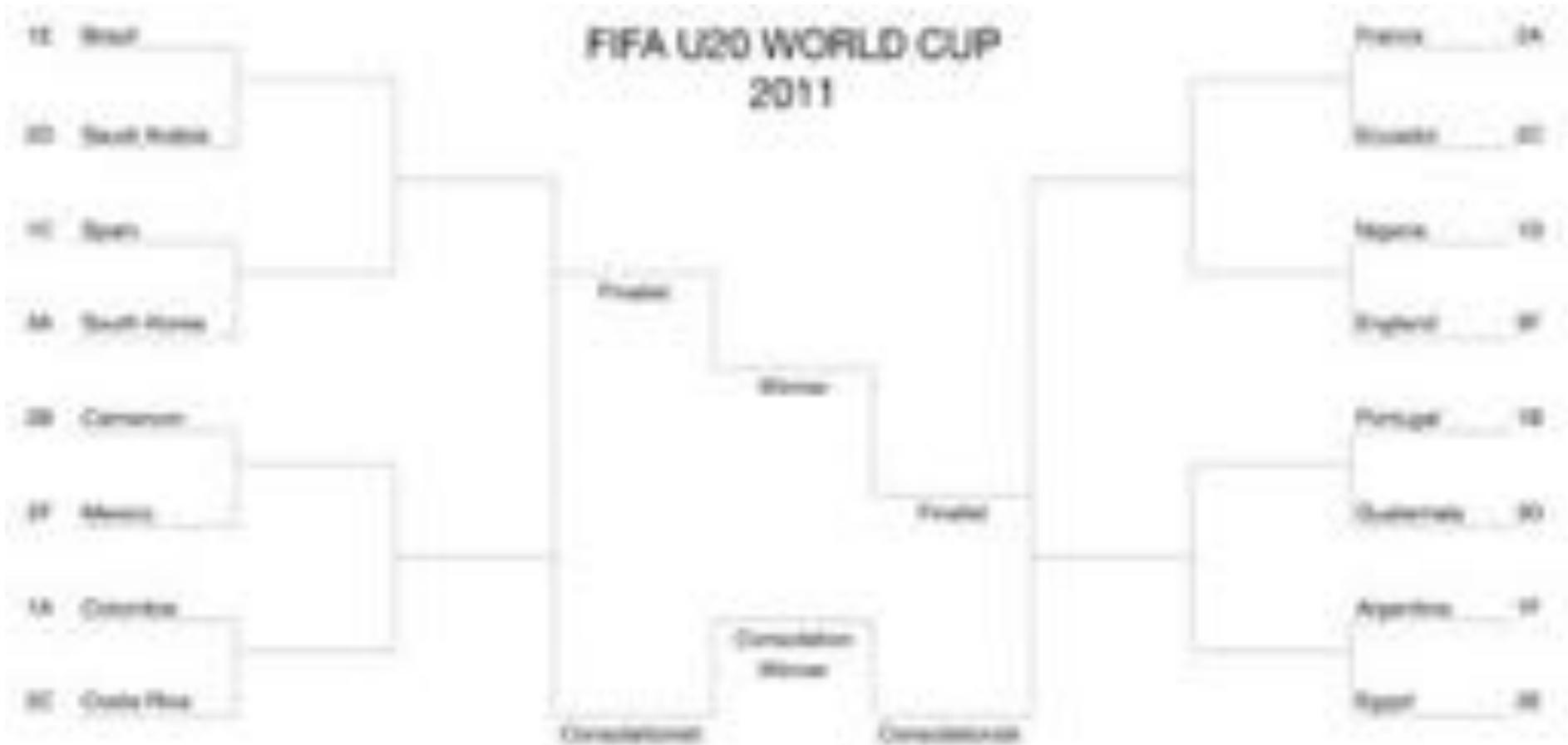
AAGACTCCGACTGGGACTTT
AAGACTGGGACTCCGACTTT

- 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

Warmup Exercise

Who here was born closest to Oct 3?

- You can only compare to 1 other person at a time



Find winner among 16 teams in just 4 rounds

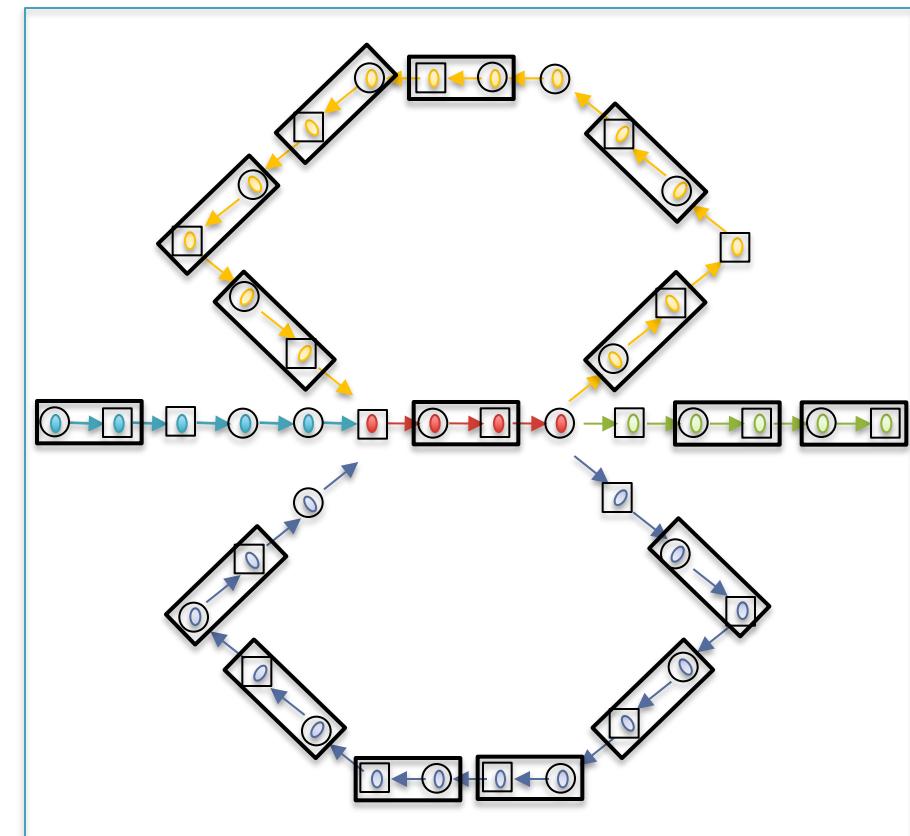
Fast Path Compression

Challenges

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

Randomized List Ranking

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



Initial Graph: 42 nodes

Randomized Speed-ups in Parallel Computation.

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

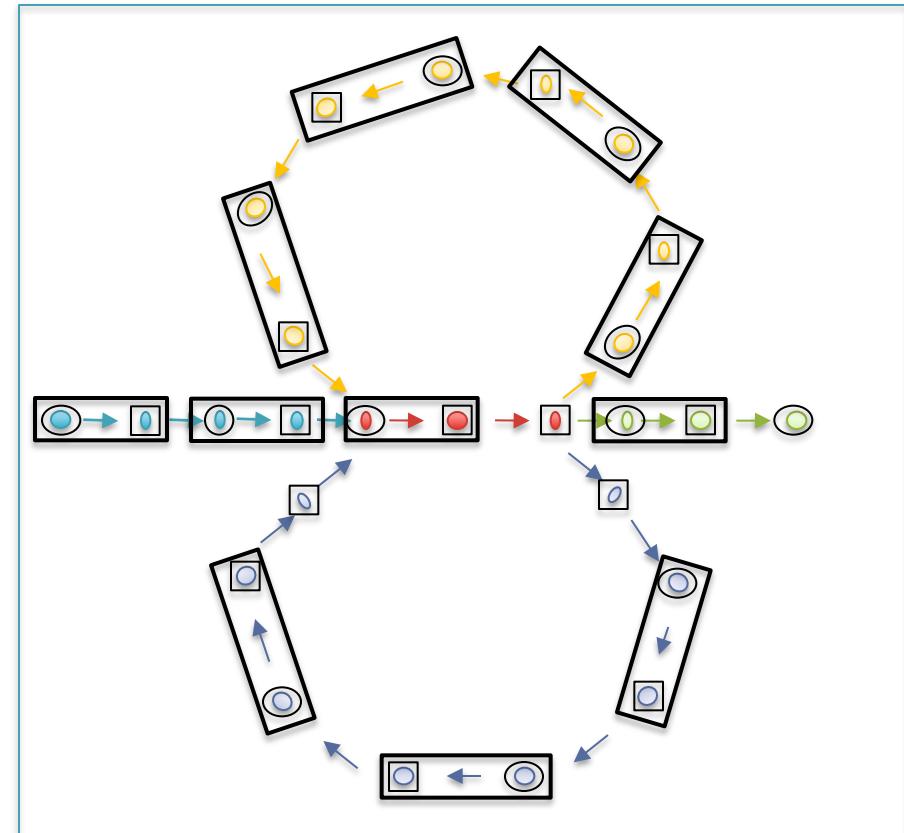
Fast Path Compression

Challenges

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

Randomized List Ranking

- Randomly assign H / T to each compressible node
- Compress $H \rightarrow T$ links



Round 1: 26 nodes (38% savings)

Randomized Speed-ups in Parallel Computation.

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

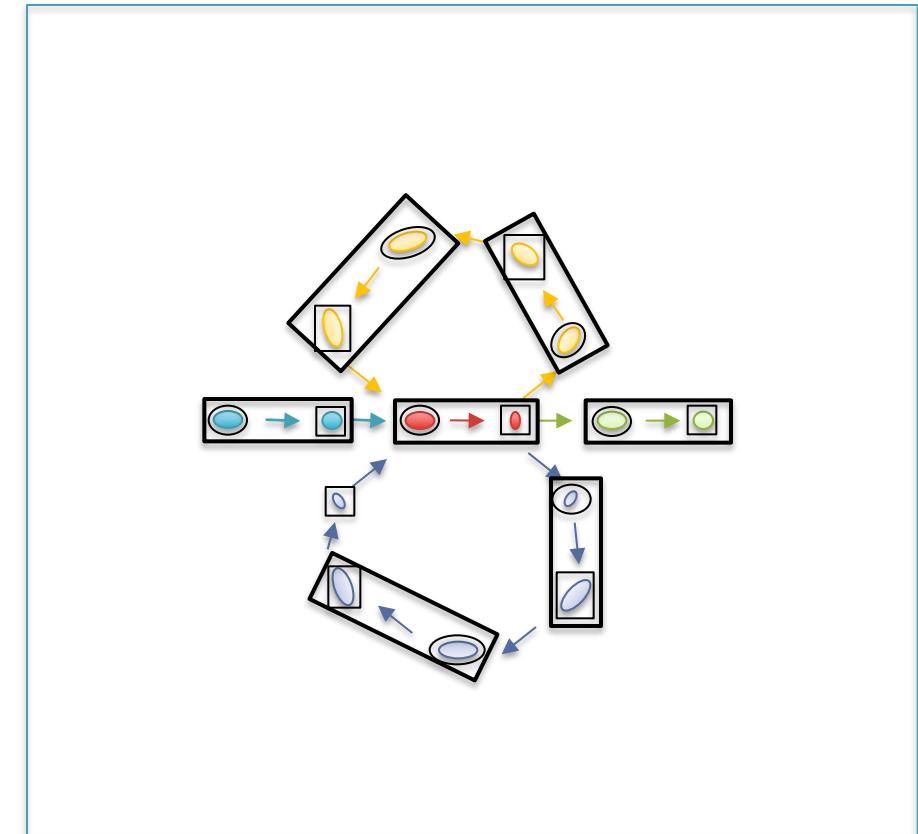
Fast Path Compression

Challenges

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

Randomized List Ranking

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



Round 2: 15 nodes (64% savings)

Randomized Speed-ups in Parallel Computation.

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

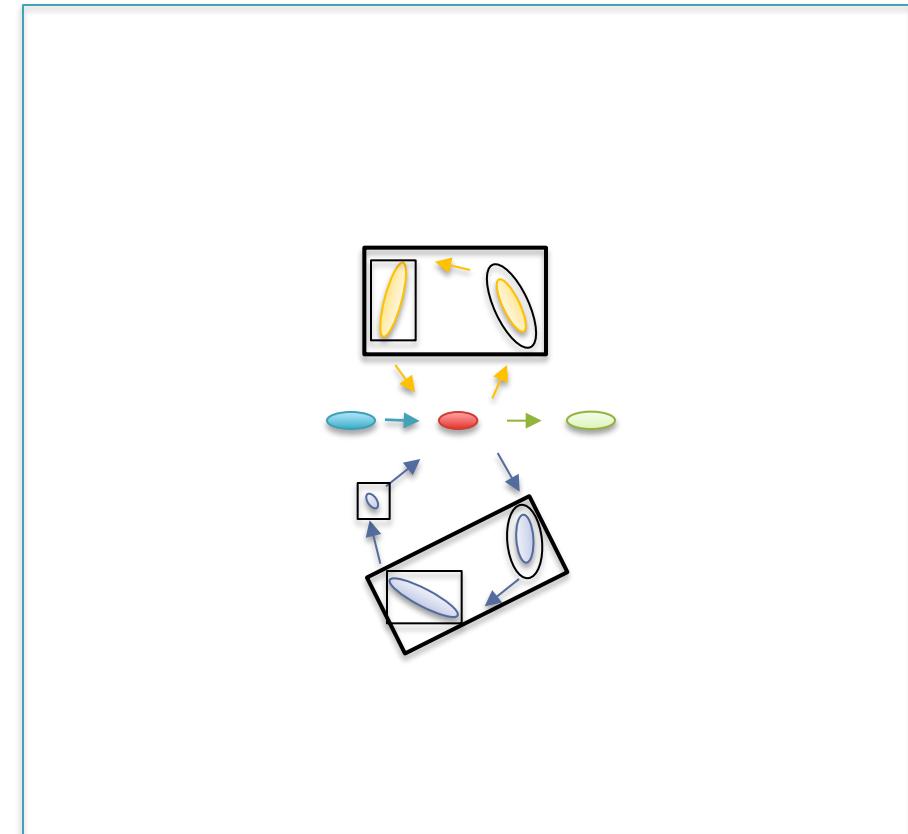
Fast Path Compression

Challenges

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

Randomized List Ranking

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



Round 2: 8 nodes (81% savings)

Randomized Speed-ups in Parallel Computation.

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

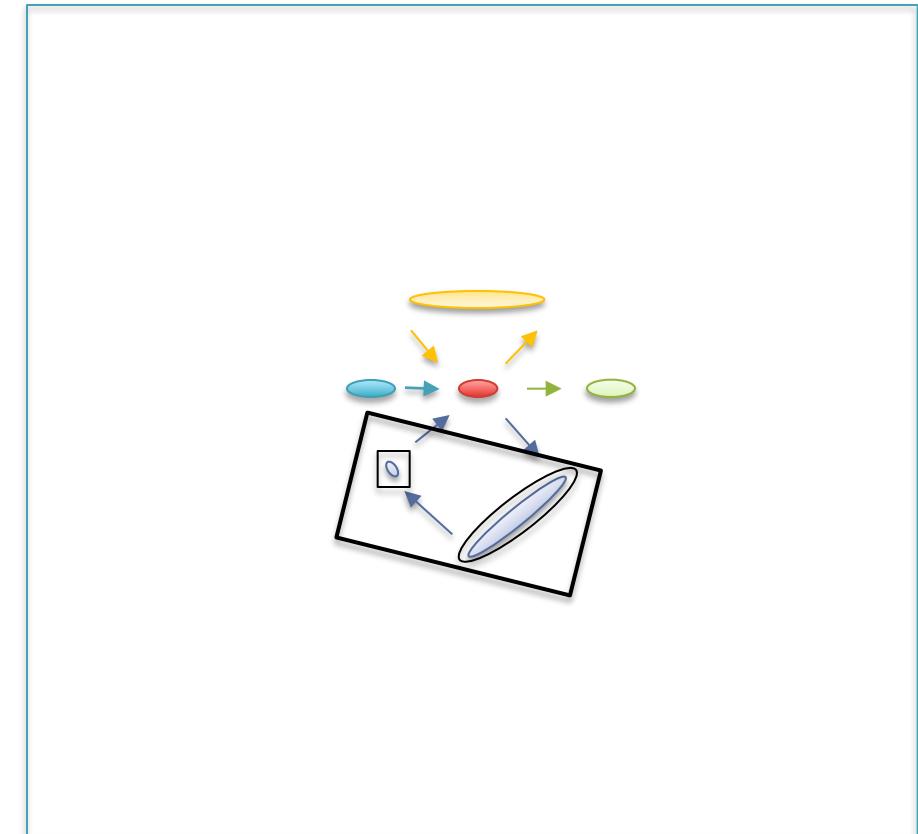
Fast Path Compression

Challenges

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

Randomized List Ranking

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

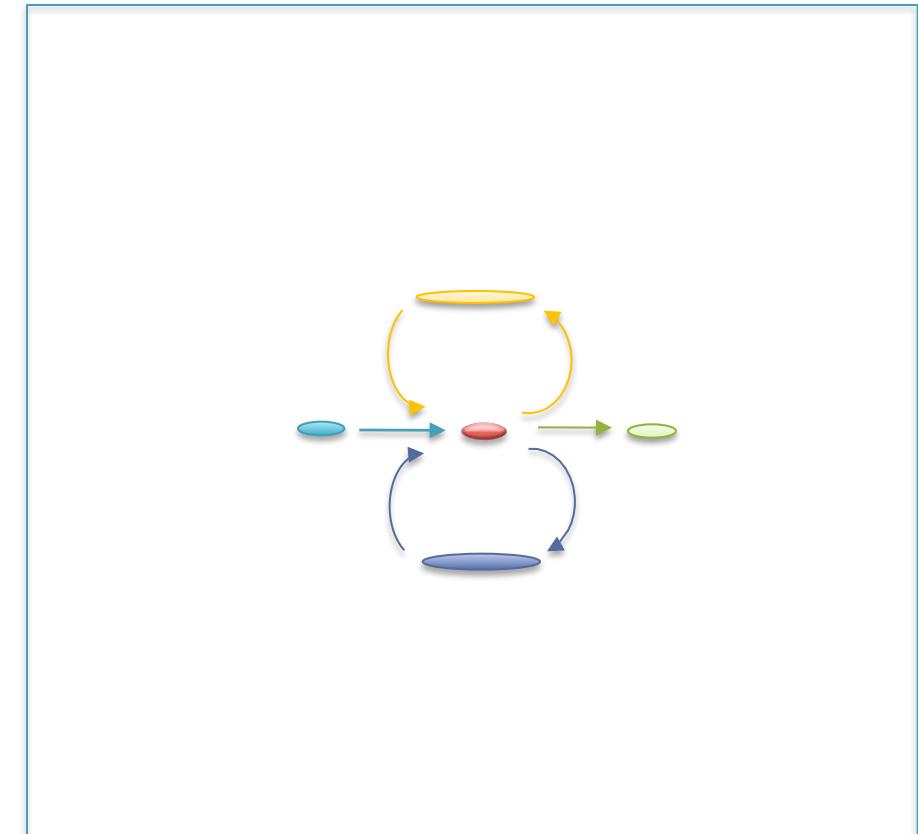
- Nodes stored on different computers
- Nodes can only access direct neighbors

Randomized List Ranking

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

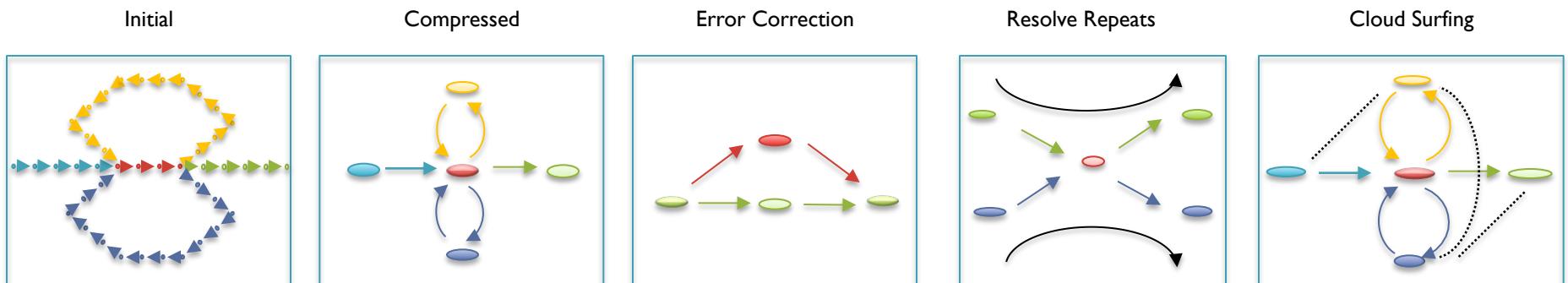
Contrail

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



De novo bacterial assembly

- *Genome: E. coli K12 MG1655, 4.6Mbp*
- *Input: 20.8M 36bp reads, 200bp insert (~150x coverage)*
- *Preprocessor: Quake Error Correction*



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.*

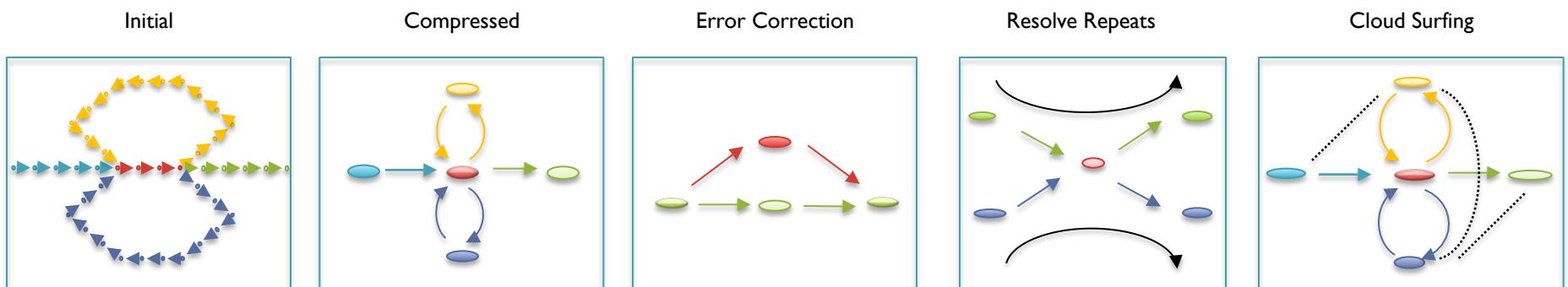
Contrail

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



De novo Assembly of the Human Genome

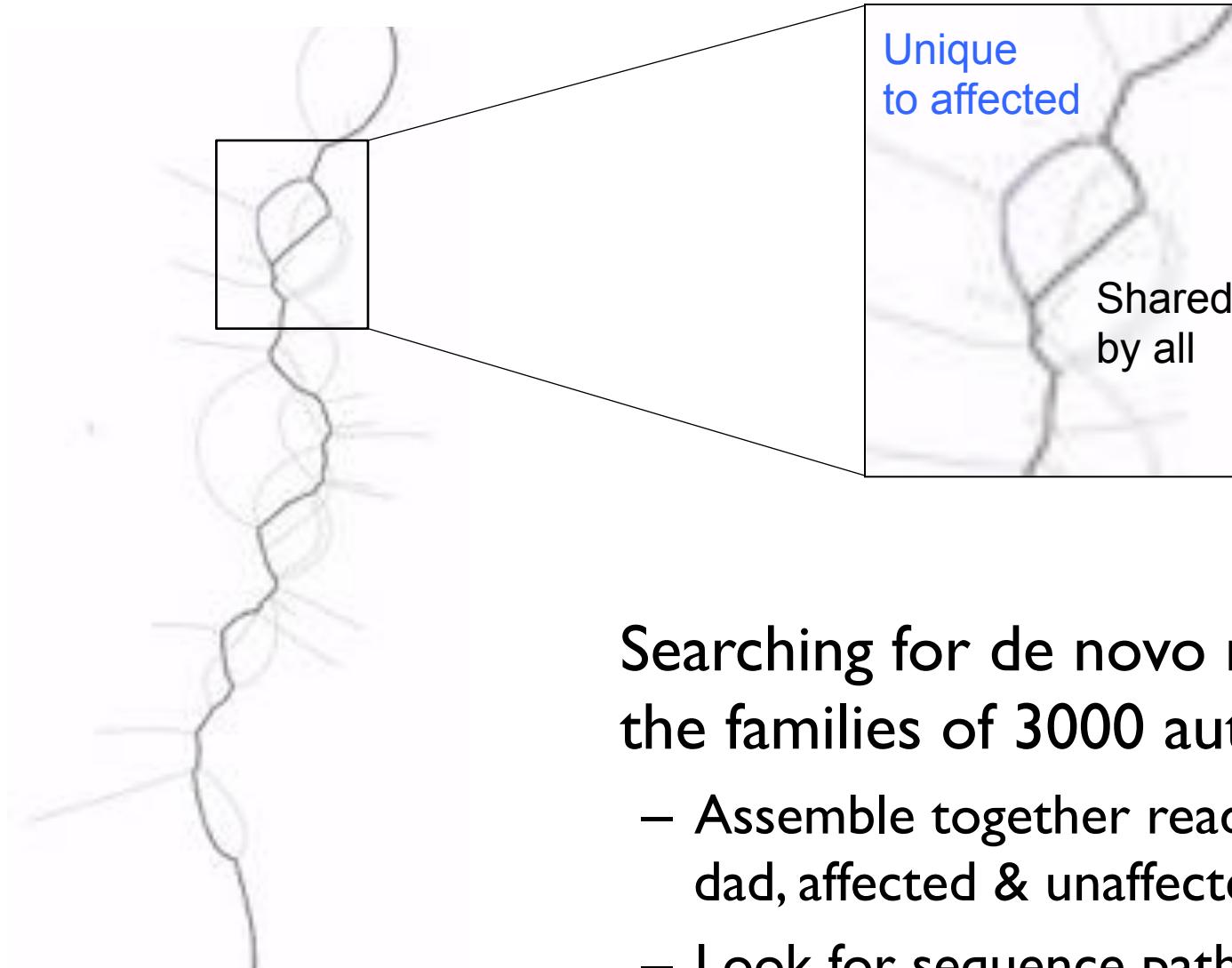
- *Genome:* African male NA18507 (SRA000271, Bentley *et al.*, 2008)
- *Input:* 3.5B 36bp reads, 210bp insert (~40x coverage)



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.*

De novo mutations and de Bruijn Graphs



Searching for de novo mutations in the families of 3000 autistic children.

- Assemble together reads from mom, dad, affected & unaffected children
- Look for sequence paths unique to affected child

COLEC12
C->A



Summary

- We are entering the digital age of biology
 - Next generation sequencing, microarrays, mass spectrometry, microscopy, ecology, etc
 - Parallel computing may be our only hope for keeping up with the pace of advance
- Modern biology requires (is) quantitative biology
 - Computational, mathematical, and statistical techniques applied to analyze, integrate, and interpret biological sensor data
- Don't let the data tsunami crash on you
 - Study, practice, collaborate with quantitative techniques

WATSON SCHOOL of BIOLOGICAL SCIENCES

Since opening in 1999, the WSBS has become a leading PhD program in the biological sciences, one whose fresh approach is quickly being emulated by other programs across the country.

- An innovative Ph.D. program designed for exceptional students
 - Approximately four years from matriculation to Ph.D. degree award
 - A first year with course work and laboratory rotations in separate phases
 - Emphasis on the principles of scientific reasoning and logic
- Learn more: <http://www.cshl.edu/gradschool>

Acknowledgements

Schatzlab

Mitch Bekritsky

Matt Titmus

Hayan Lee

James Gurtowski

Anirudh Aithal

Rohith Menon

Goutham Bhat

CSHL

Dick McCombie

Melissa Kramer

Eric Antonio

Mike Wigler

Zach Lippman

Doreen Ware

Ivan Iossifov

JHU

Steven Salzberg

Ben Langmead

Jeff Leek

NBACC

Adam Phillip

Sergey Koren

Univ. of Maryland

Mihai Pop

Art Delcher

Jimmy Lin

David Kelley

Dan Sommer

Cole Trapnell



Thank You!

<http://schatzlab.cshl.edu>

@mike_schatz