

Genome Assembly

Michael Schatz

Sept 6, 2023

Lecture 3: Applied Comparative Genomics



Assignment I

Due end of day on Sept 6 (right before midnight)

The screenshot shows a GitHub repository page for 'appliedgenomics2023/assignment1'. The left sidebar displays a file tree with 'main' and 'assignments/assignment1' branches. The 'assignments/assignment1' branch contains files like README.md, TAIR10.chrom.sizes, ce10.chrom.sizes, chr22.fa.gz, dm6.chrom.sizes, ecoli.chrom.sizes, hg38.chrom.sizes, tomato.chrom.sizes, wheat.chrom.sizes, yeast.chrom.sizes, lectures/01.introduction.pdf, LICENSE, and README.md. The main content area is titled 'Assignment 1: Chromosome Structures'. It includes assignment details: 'Assignment Date: Wednesday, August 30, 2023' and 'Due Date: Wednesday, Sept. 6, 2023 @ 11:59pm'. The 'Assignment Overview' section states: 'In this assignment you will profile the overall structure of the genomes of several important species and then study human chromosome 22 in more detail. As a reminder, any questions about the assignment should be posted to Piazza.' The 'Question 1: Chromosome structures [10 pts]' section asks students to download chromosome size files for various species and create a table. It lists 8 species with links to their descriptions. The 'Question 2. Coverage simulator [20 pts]' section provides instructions for simulating sequencing coverage and includes a snippet of pseudocode.

Assignment Date: Wednesday, August 30, 2023
Due Date: Wednesday, Sept. 6, 2023 @ 11:59pm

Assignment Overview

In this assignment you will profile the overall structure of the genomes of several important species and then study human chromosome 22 in more detail. As a reminder, any questions about the assignment should be posted to [Piazza](#).

Question 1: Chromosome structures [10 pts]

Download the chromosome size files for the following genomes (Note these have been preprocessed to only include main chromosomes):

1. *Arabidopsis thaliana* ([TAIR10](#)) - An important plant model species [\[info\]](#)
2. *Tomato* ([Solanum lycopersicum v4.00](#)) - One of the most important food crops [\[info\]](#)
3. *E. coli* ([Escherichia coli K12](#)) - One of the most commonly studied bacteria [\[info\]](#)
4. *Fruit Fly* ([Drosophila melanogaster, dm6](#)) - One of the most important model species for genetics [\[info\]](#)
5. *Human* ([hg38](#)) - us :) [\[info\]](#)
6. *Wheat* ([Triticum aestivum, IWGSC](#)) - The food crop which takes up the largest land area [\[info\]](#)
7. *Worm* ([Caenorhabditis elegans, ce10](#)) - One of the most important animal model species [\[info\]](#)
8. *Yeast* ([Saccharomyces cerevisiae, sacCer3](#)) - an important eukaryotic model species, also good for bread and beer [\[info\]](#)

Using these files, make a table with the following information per species:

- Question 1.1. Total genome size
- Question 1.2. Number of chromosomes
- Question 1.3. Largest chromosome size and name
- Question 1.4. Smallest chromosome size and name
- Question 1.5. Mean chromosome length

Question 2. Coverage simulator [20 pts]

- Question 2.1. How many 100bp reads are needed to sequence a 1Mbp genome to 3x coverage?
- Question 2.2. In the language of your choice, simulate sequencing 3x coverage of a 1Mbp genome with 100bp reads and plot the histogram of coverage. Note you do not need to actually output the sequences of the reads, you can just uniformly randomly sample positions in the genome and record the coverage. You do not need to consider the strand of each read. The start position of each read should have a uniform random probability at each possible starting position (1 through 999,901). You can record the coverage in an array of 1M positions. Overlay the histogram with a Poisson distribution with lambda=3. Also overlay the distribution with a Normal distribution with a mean of 3 and a standard deviation of 1.73 (which is the square root of 3). Here is the pseudocode for the simulator:

```
num_reads = calculate_number_of_reads(genomesize, readlength, coverage)
## use an array to keep track of the coverage at each position in the genome
genome_coverage = initialize_array_with_zero(genomesize)
for (i = 0; i < num_reads; i++)
```

<https://github.com/schatzlab/appliedgenomics2023/tree/main/assignments/assignment1>

On ChatGPT...

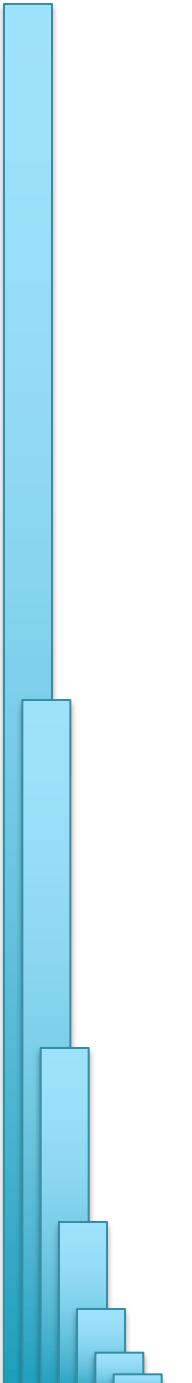
ChatGPT 3.5 output:

... Notably, Professor Michael Schatz, a pioneer in genomics research at Johns Hopkins University, has made significant contributions to the development of high-throughput computational tools and techniques for genome assembly and analysis (Schatz et al., 2019). These advances are crucial for microbiologists working with Acetobacter to harness the wealth of genomic data available, advancing our understanding of their evolution and potential biotechnological applications....

References:

...

4. Schatz, M. C., Maron, L. G., & Stein, J. C. (2019). Assembly and annotation of genomes from long-read sequencing. *Current Opinion in Plant Biology*, 48, 69-75.



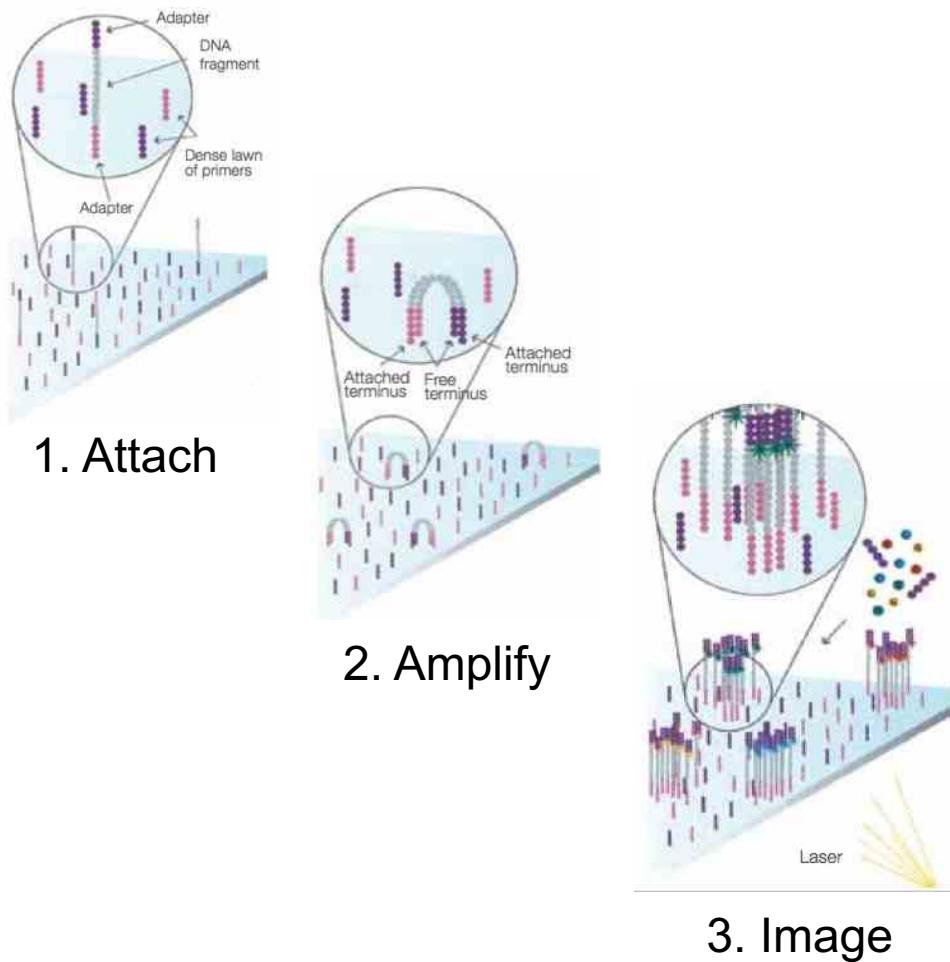
Part I: Recap and Illumina Sequencing

Second Generation Sequencing



Illumina NovaSeq 6000
Sequencing by Synthesis

>3Tbp / day
(JHU has 4 of these!)

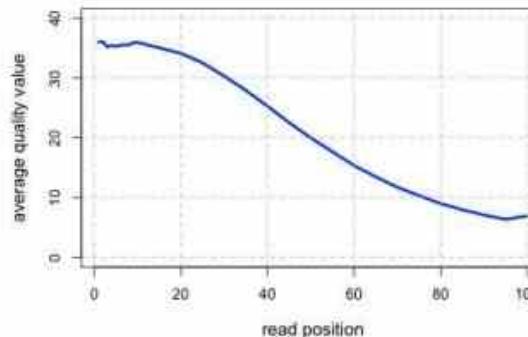


Metzker (2010) Nature Reviews Genetics 11:31-46
<https://www.youtube.com/watch?v=fCd6B5HRaZ8>

Illumina Quality

QV	p_{error}
40	1/10000
30	1/1000
20	1/100
10	1/10

$$Q_{\text{sanger}} = -10 \log_{10} p$$



```

SSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSS.....  

.....XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX.....  

.....IIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIII.....  

.....JJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJ.....  

LLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLL.....  

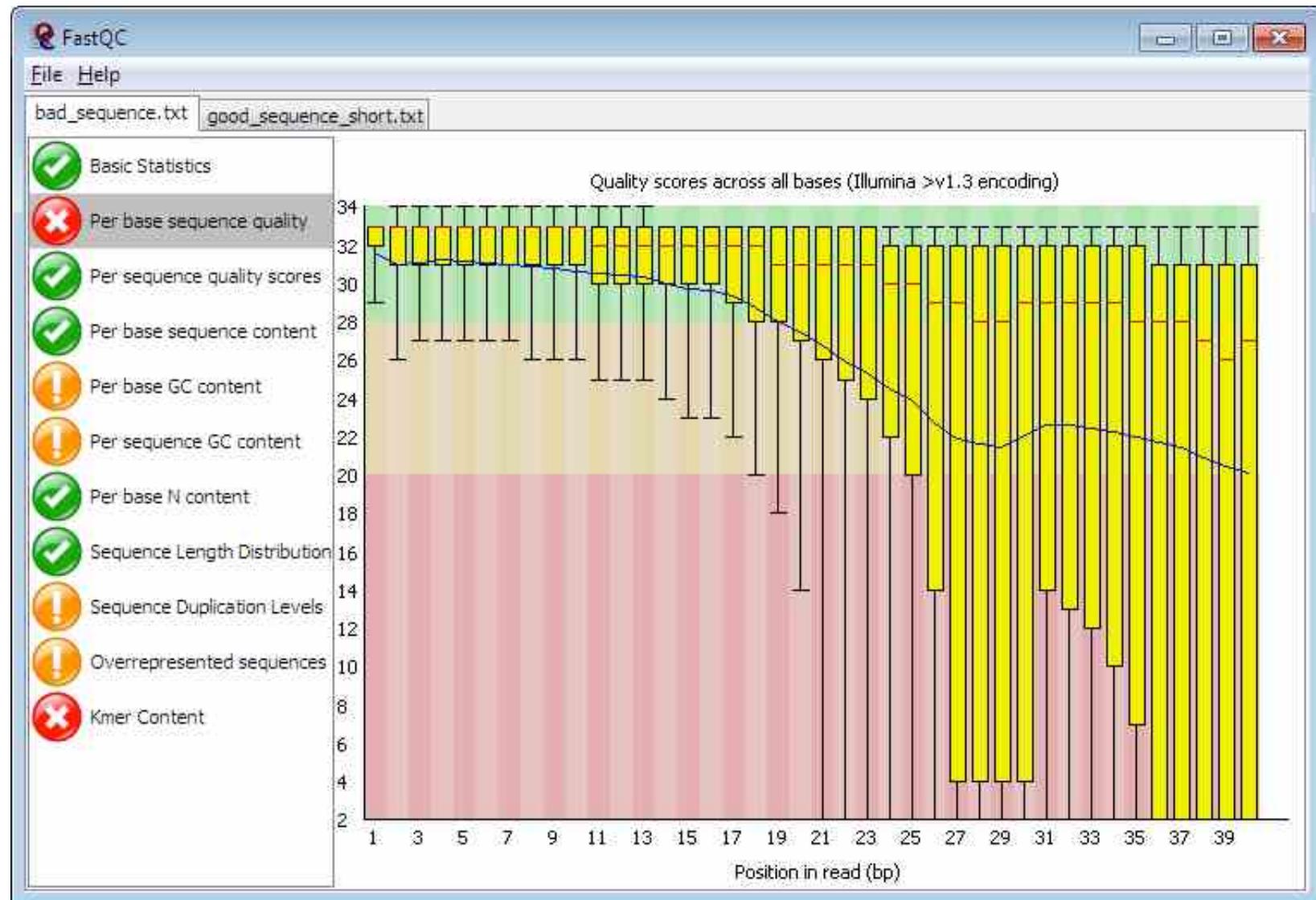
! "#$%&'()*+,.-./0123456789:;=>?@ABCDEFGHIJKLMNPQRSTUVWXYZ[\]^_`abcdefghijklmnopqrstuvwxyz{|}-
|           |           |           |           |  

33          59          64          73          104          126

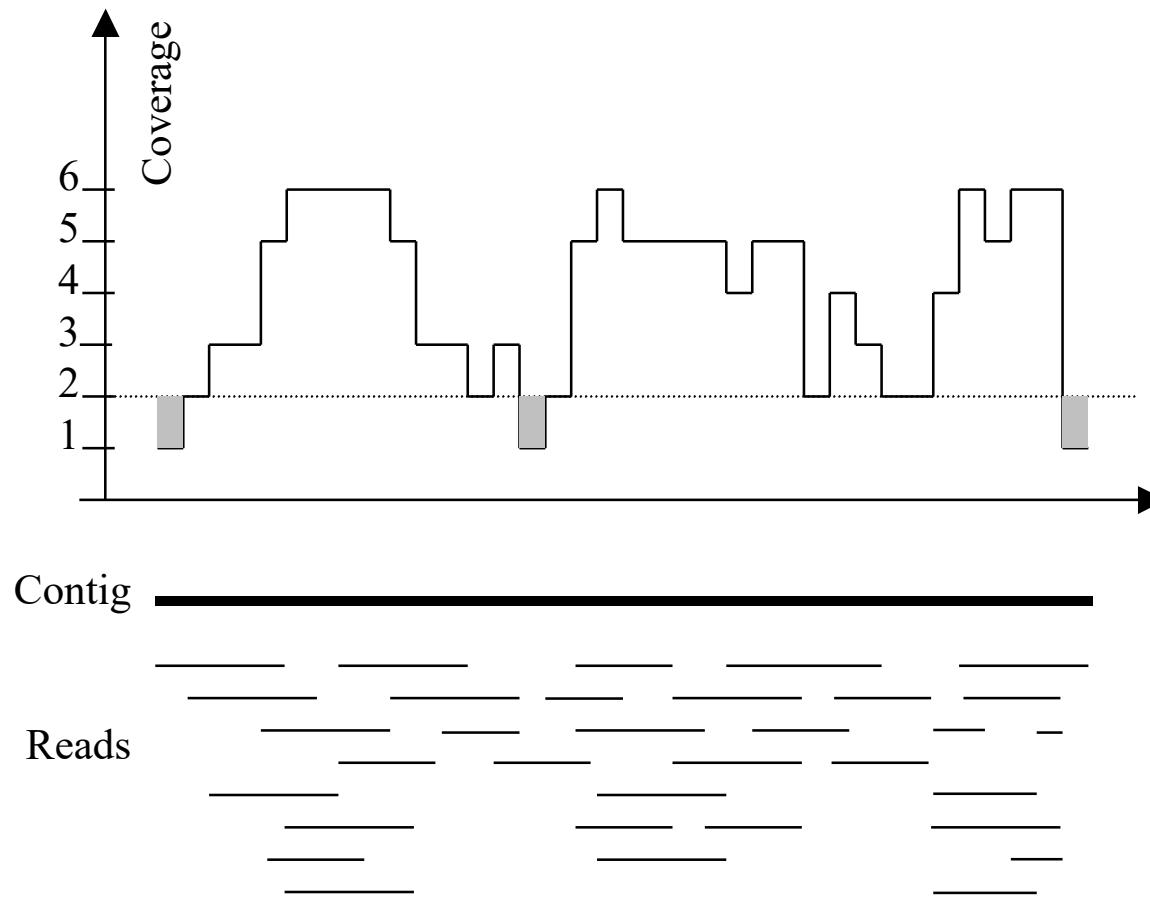
```

- S - Sanger Phred+33, raw reads typically (0, 40)
- X - Solexa Solexa+64, raw reads typically (-5, 40)
- I - Illumina 1.3+ Phred+64, raw reads typically (0, 40)
- J - Illumina 1.5+ Phred+64, raw reads typically (3, 40)
with 0=unused, 1=unused, 2=Read Segment Quality Control Indicator (bold)
(Note: See discussion above).
- L - Illumina 1.8+ Phred+33, raw reads typically (0, 41)

FASTQC: Are my data any good?



Typical sequencing coverage



Imagine raindrops on a sidewalk

We want to cover the entire sidewalk but each drop costs \$1

If the genome is 10 Mbp, should we sequence 100k 100bp reads?

Poisson Distribution

The probability of a given number of events occurring in a fixed interval of time and/or space if these events occur with a known average rate and independently of the time since the last event.

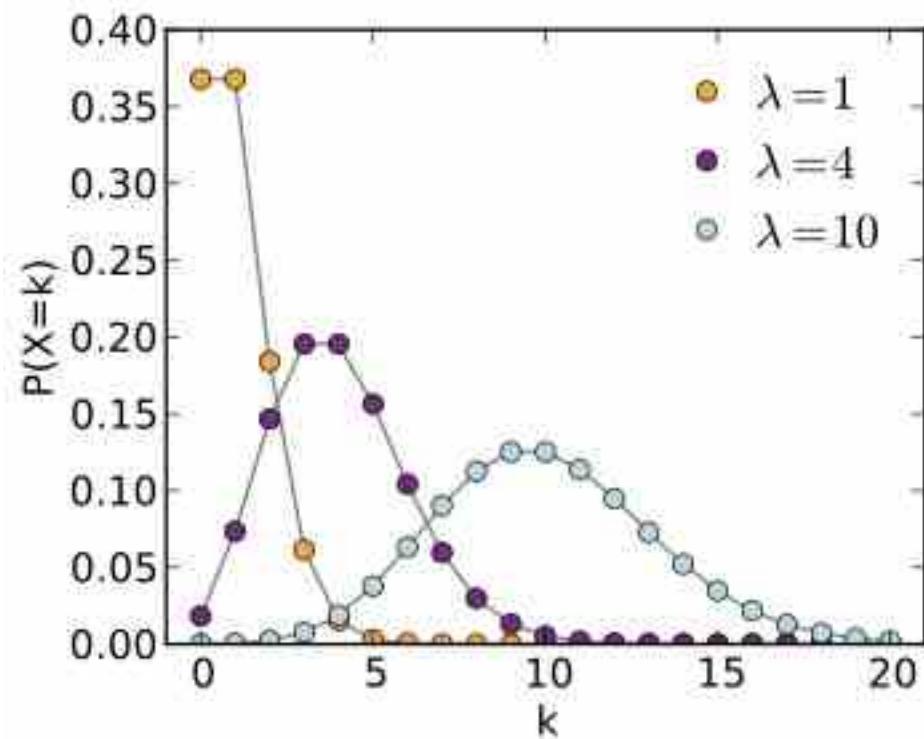
Formulation comes from the limit of the binomial equation

Resembles a normal distribution, but over the positive values, and with only a single parameter.

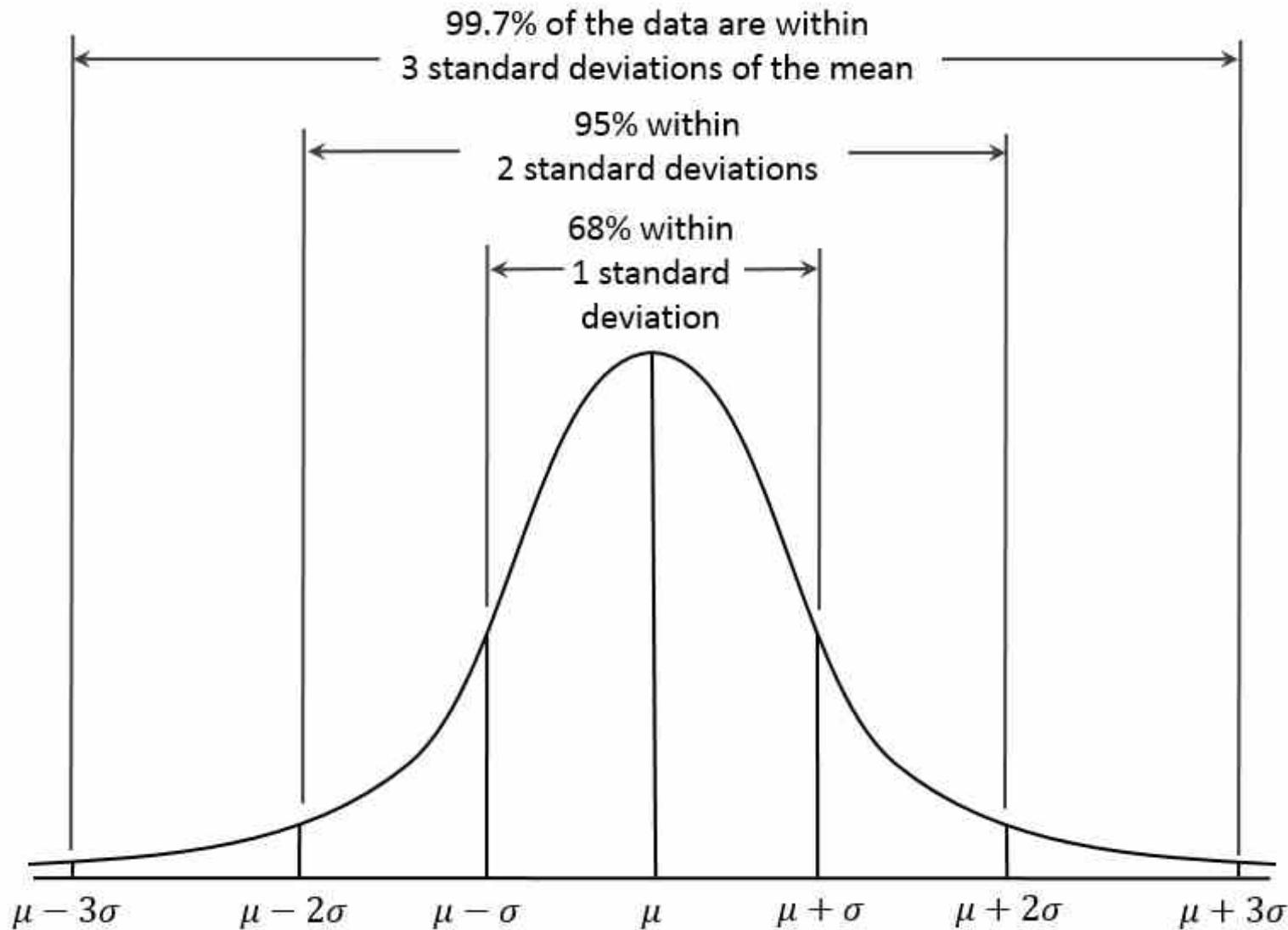
Key properties:

- ***The standard deviation is the square root of the mean.***
- ***For mean > 5, well approximated by a normal distribution***

$$P(k) = \frac{\lambda^k}{k!} e^{-\lambda}$$



Normal Approximation



Can estimate Poisson distribution as a normal distribution when $\lambda > 10$

K-mers and K-mer counting

GATTACATACACATTGGATG

GAT ACA ACA ATT GAT

ATT CAT CAC TTG ATG

TTA ATA ACA TGG

TAC TAC CAT GGA

Kmers:

- Divide a string into substrings of length k
- Notice every position is covered k times
- Notice there are $G - k + 1$ kmers from a string of length G

Computation: Very easy to compute, exact matches, represent 32mers in 64 bits

Biological: The “atomic unit” of a sequence, creates a fingerprint of a genome/read

Transcription Factors

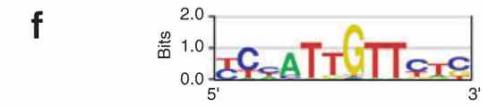
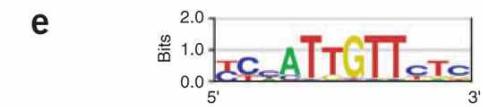
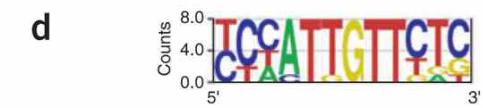
A transcription factor (or sequence-specific DNA-binding factor) is a protein that controls the rate of transcription of genetic information from DNA to messenger RNA, by binding to a specific DNA sequence.

- Transcription factors work alone or with other proteins in a complex, by promoting (as an activator), or blocking (as a repressor) the recruitment of RNA polymerase to specific genes.
- A defining feature of transcription factors is that they contain at least one DNA-binding domain (DBD)
- Figure (a) Eight known genomic binding sites in three *S. cerevisiae* genes. (b) Degenerate consensus sequence. (c,d) Frequencies of nucleotides at each position. (e) Sequence logo (f) Energy normalized logo using relative entropy to adjust for low GC content in *S. cerevisiae*.

a HEM13 CCCATTGTTCTC
HEM13 TTTCTGGTTCTC
HEM13 TCAATTGTTTAG
ANB1 CTCATTGTTGTC
ANB1 TCCATTGTTCTC
ANB1 CCTATTGTTCTC
ANB1 TCCATTGTTCGT
ROX1 CCAATTGTTTTG

b YCHATTGTTCTC

c A 002700000010
C 464100000505
G 000001800112
T 422087088261



Bob Crimi

What are DNA sequence motifs?

D'haeseleer (2006) Nature Biotechnology 24, 423 – 425 doi:10.1038/nbt0406-423

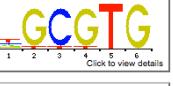
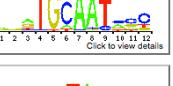
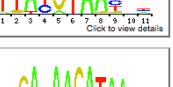
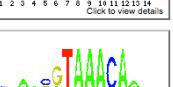
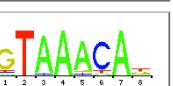
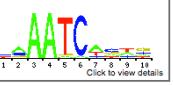
Transcription Factors Database

The JASPAR database Michael

jaspar.genereg.net/cgi-bin/jaspar_db.p?rm=browse&db=core&tax_group=vertebrates

SEARCH Name AND Species AND Class SEARCH ?

JASPAR matrix models:

TOGGLE	ID	name	species	class	family	Sequence logo
■	MA0004.1	Arnt	Mus musculus	Basic helix-loop-helix factors (bHLH)	PAS domain factors	
■	MA0006.1	Ahr:Arnt	Mus musculus	Basic helix-loop-helix factors (bHLH);Basic helix-loop-helix factors (bHLH)	PAS domain factors::PAS domain factors	
■	MA0019.1	Ddit3::Cebpa	Rattus norvegicus	Basic leucine zipper factors (bZIP);Basic leucine zipper factors (bZIP)	C/EBP-related::C/EBP-related	
■	MA0025.1	NFIL3	Homo sapiens	Basic leucine zipper factors (bZIP)	C/EBP-related	
■	MA0029.1	Mecom	Mus musculus	C2H2 zinc finger factors	Factors with multiple dispersed zinc fingers	
■	MA0030.1	FOXF2	Homo sapiens	Fork head / winged helix factors	Forkhead box (FOX) factors	
■	MA0031.1	FOXD1	Homo sapiens	Fork head / winged helix factors	Forkhead box (FOX) factors	
■	MA0038.1	Gfi1	Rattus norvegicus	C2H2 zinc finger factors	More than 3 adjacent zinc finger factors	
■	MA0040.1	Foxq1	Rattus norvegicus	Fork head / winged helix factors	Forkhead box (FOX) factors	

ANALYZE selected matrix models:

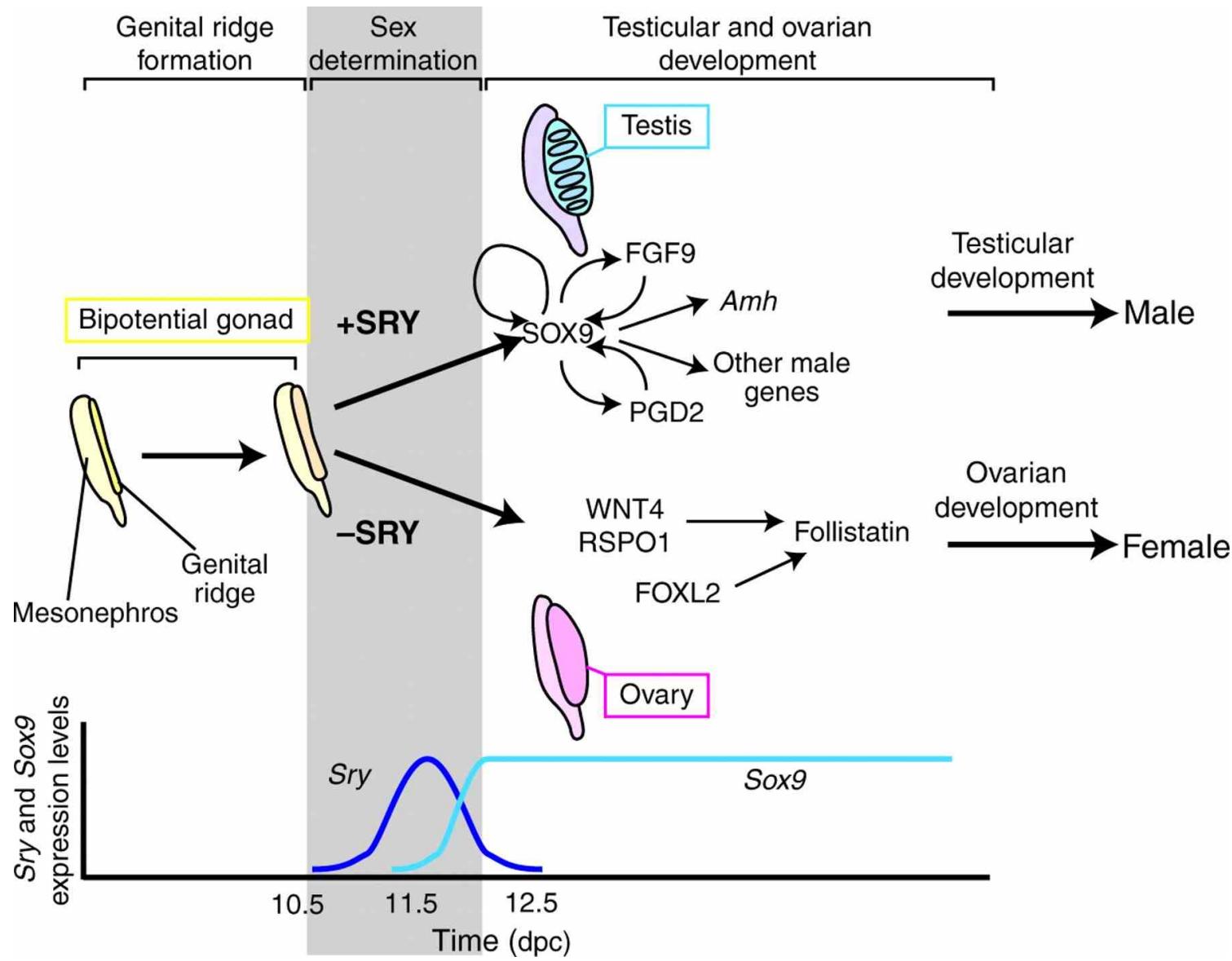
CLUSTER selected models using STAMP

Create RANDOM matrix models based on selected models
Number of matrices: 200 Format: Raw **RANDOMIZE**

Create models with PERMUTED columns from selected:
Type: Within each matrix Format: Raw **PERMUTE**

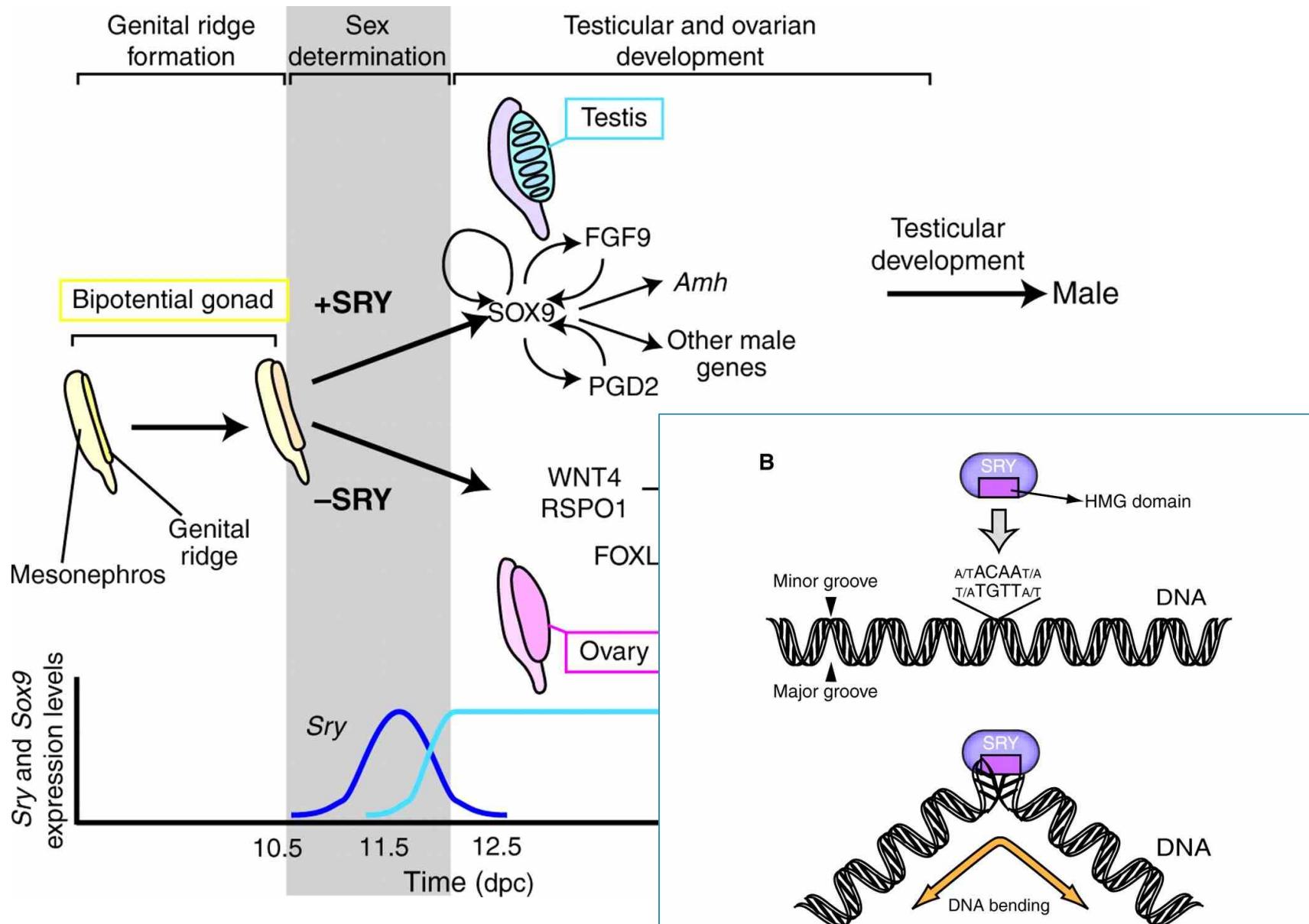
SCAN this (fasta-formatted) sequence with selected matrix models
Relative profile score threshold: 80 % **SCAN**

JASPAR 2014: an extensively expanded and updated open-access database of transcription factor binding profiles
Anthony Mathelier (2014) Nucleic Acids Res. 42 (D1): D142-D147. DOI: <https://doi.org/10.1093/nar/gkt997>



SRY: The master switch in mammalian sex determination

Kashimada and Koopman (2010) Development 137: 3921-3930; doi: 10.1242/dev.048983



SRY: The master switch in mammalian sex determination

Kashimada and Koopman (2010) Development 137: 3921-3930; doi: 10.1242/dev.048983

K-mers and K-mer counting

GATTACATACACATTGGATG

1 : 7 (**ATG**, **TGG**, ...)

2 : 4 (**GAT**, **CAT**, **ATT**, **TAC**)

3 : 1 (**ACA**)

How long should k be?

K=1 : Too short, every base is present

K=2 : Too short, every pair of bases will be present

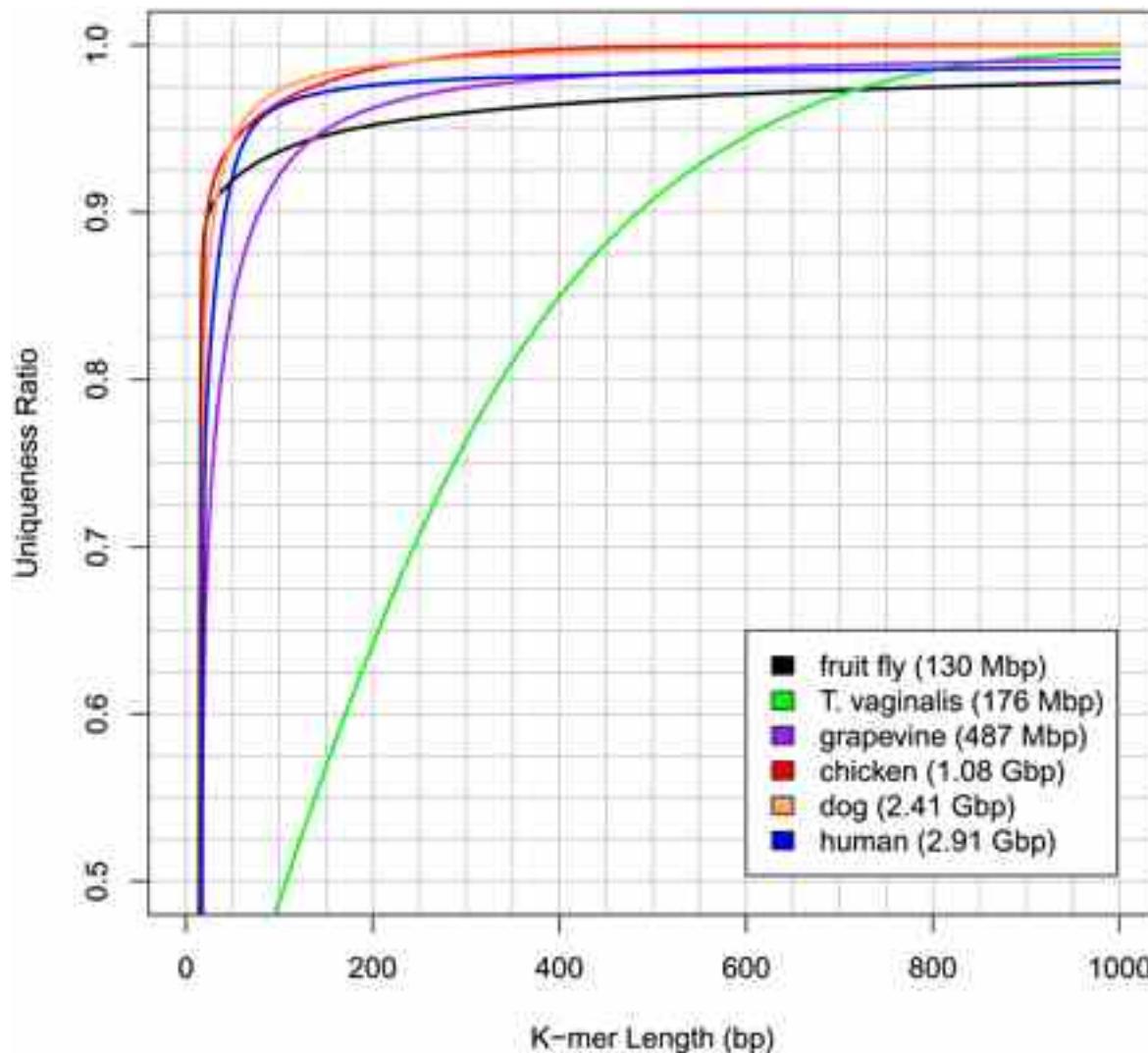
Pick k so that $G/(4^k) \ll 1$

$$k = \log_4(G)$$

At least 15 for human, often a bit longer

But not too long or could loose resolution

K-mer Uniqueness



Assembly of large genomes using second-generation sequencing
Schatz et al. (2010) Genome Research. doi: 10.1101/gr.101360.109

Question?

We would love to generate
longer and longer reads with this technology

What can we do?

Illumina Hacking

BIOINFORMATICS ORIGINAL PAPER

Vol. 29 no. 12, 2013, pages 1492–1497
doi:10.1093/bioinformatics/btt178

Genome analysis

Advance Access publication May 22, 2013

Assembling the 20Gb white spruce (*Picea glauca*) genome from whole-genome shotgun sequencing data

Inanc Birol^{1,2,3,*}, Anthony Raymond¹, Shaun D. Jackman¹, Stephen Pleasance¹, Robin Coop¹, Greg A. Taylor¹, Macaire Man Saint Yuen⁴, Christopher I. Keeling⁴, Dana Brand¹, Benjamin P. Vandervalk¹, Heather Kirk¹, Pawan Pandoh¹, Richard A. Moore¹, Yongjun Zhao¹, Andrew J. Mungall¹, Barry Jaquish⁵, Alvin Yanchuk⁶, Carol Ritland^{4,6}, Brian Boyle^{7,8}, Jean Bousquet^{7,8}, Kermit Ritland⁶, John MacKay^{7,8}, Jörg Bohlmann¹, Steven JM Jones^{1,2,9}

¹Genome Sciences Centre, British Columbia Cancer Agency, Vancouver, BC V5Z 4S6, Canada; ²University of British Columbia, Vancouver, BC V6H 3W1, Canada; ³Computational Biology, Burnaby, BC V5A 1S6, Canada; ⁴Michael Smith Laboratories, University of British Columbia, Vancouver, BC V6W 9CZ, Canada; ⁵Department of Forest Sciences, University of British Columbia, Vancouver, BC V6W 9CZ, Canada; ⁶Institute for Systems and Integrative Biology, Université Laval, Québec, QC G1V 0A6, Canada; ⁷Department of Wood and Forest Sciences, Université Laval, Québec, QC G1V 0A6, Canada; ⁸Molecular Biology and Biochemistry, Simon Fraser University, Burnaby, BC V5A 1S6, Canada; ⁹Associate Editor: Michael Brudno

ABSTRACT

White spruce (*Picea glauca*) is a dominant conifer of the boreal forests of North America, and providing genomic resources for this commercially valuable tree will help improve forest management and conservation efforts. Seeing and assessing the large and highly repetitive spruce genome throughput pushes the boundaries of the current technology. Here we describe a whole-genome shotgun sequencing strategy using two Illumina sequencing platforms and an assembly approach using the ABYSS software. We report a 20 giga base pairs draft genome in 4.9 million scaffolds, with a scaffold N50 of 23,956 bp. We demonstrate how recent improvements in the sequencing technology, especially increasing read lengths and paired end reads from longer fragments have a major impact on the assembly contiguity. We also note that scalable bioinformatics tools are instrumental in providing rapid draft assemblies.

Availability: White spruce genome sequencing and assembly data are available through NCBI's Accession# ALW2D100000000 (PID: PRJNA8345). <http://www.ncbi.nlm.nih.gov/bioproject/8345>.

Contact: ibiro@csbc.ca

Supplementary information: Supplementary data are available at *Bioinformatics* online.

Received on March 20, 2013; revised on April 10, 2013; accepted on April 11, 2013

1 INTRODUCTION

The assembly of short reads to develop genomic resources for non-model species remains an active area of development (Schatz *et al.*, 2012). The feasibility of the approach and its scalability to

*To whom correspondence should be addressed.

© The Author 2013. Published by Oxford University Press.
This is an Open Access article distributed under the terms of the Creative Commons Attribution Non-Commercial License (<http://creativecommons.org/licenses/by-nc/3.0/>), which permits non-commercial re-use, distribution, and reproduction in any medium, provided the original work is properly cited.

large genomes was demonstrated by Simpson *et al.* (2009) using human and was later extended to other organisms such as the mouse (Li *et al.*, 2012). High quality results, as demonstrated by Ladurner *et al.* (2013), can be successfully applied to numerous genomes (Chan *et al.*, 2011; Chu *et al.*, 2011; Godet *et al.*, 2012; Swart *et al.*, 2012).

Estimated at 20 giga base pairs (Gbp), the genome of the pine (*Pinus taeda*) family present unique challenges in the sequencing and assembly of the genome. In the pine (*Pinus taeda*) genome, we found that the size of the genome (ca. 4.9 Gbp) and the number of genes (ca. 20,000) are similar to the human genome. The genome is highly repetitive, containing approximately 80% repeats (Swart *et al.*, 2012). The genome is also highly methylated (ca. 90%) (Liu *et al.*, 2012).

We addressed the data representation and sequencing multiple whole-genome HiSeq 2000 and MiSeq sequences (CA, USA). Compared with localized sequencing approaches of isolating ~10-kb DNA sequencing fragments in high throughput (CA, USA), a shotgun only sequencing strategy data effectively covering the genome can be an order of magnitude less expensive than sequencing the entire genome.

In this work, we demonstrate that at this scale remains viable and provides a high quality genome assembly.

2 METHODS

2.1 Sample collection

Apical shoot tissues were collected in April 2006 from a single white spruce (*Picea glauca*, genotype PG29) tree at the Kalamalka Research Station of the British Columbia Ministry of Forests and Ranges, Vernon, British Columbia, Canada. Genomic DNA was extracted from whole tissue by BioS&T (<http://www.biost.com/>), Montreal, QC, Canada) using an organic extraction method yielding 300 µg of high quality purified nucleic acid.

2.2 Library preparation and sequencing

DNA quality was assessed by spectrophotometry and gel electrophoresis before library construction. DNA was sheared to 45 µm using an E210 shearing system (Covaris, Woburn, MA) and size fractionated using a 1.5% agarose gel (for libraries with 200 bp insert size) or 4% agarose gel (for libraries with 500 bp insert size). DNA size fractions were excised and eluted from the gel slices overnight at 4°C in 300 µl of elution buffer [5 µl [vol/vol] Laemmli buffer (3 mM Tris-HCl (pH 7.5), 0.2 mM EDTA]/7.5 µM ammonium acetate] and were purified using a Spin-Filter Titer Plate (Fisher Scientific) and then quantified using a Nanodrop (Nanodrop). DNA was eluted and purified using a QIAquick PCR Purification kit (Qiagen) and ligated with Illumina PE adaptors (Illumina Inc.). This involved DNA end repair and formation of 3' adenine overhangs using the Klenow fragment of DNA polymerase I (3'-5' exonuclease minus) and ligation to Illumina PE adaptors (with 5' overhangs). Adapter-ligated products were purified on QIAquick spin columns (Qiagen) and amplified using Phusion Hot Start Polymerase (NEB) and 10 PCR cycles with the PE primer and 2.0 µM Q5 PCR primer. After each cycle of PCR, the purified adapter-ligated products were run on 8% PAGE gels. DNA quality was assessed and quantified using an Agilent DNA 1000 series II assay (Agilent) and Nanoplot spectrophotometer (Nanodrop). DNA was subsequently diluted to 8 µM. The final concentration was confirmed using a Quant-iT dsDNA HS assay kit and Qubit fluorometer (Invitrogen).

mate pair (MPET, a.k.a. jumping libraries) were constructed using 4-kb gap sequencing (Illumina Nextera Mate-Pair library construction protocol and reagent (FC-132-1001)). The genomic DNA sample was simultaneously fragmented and tagged with a bovine containing mate pair junction adapter, which left a short sequence gap in the fragmented DNA. The gap was filled by a strand displacement reaction using a polymerase to ensure that all fragments are flush and ready for circularization. After an AMPure bead clean-up, size selection was done on a 0.6% agarose gel and the fragments were purified using a ZymoClean Gel Purification Kit. The fragments were then end-repaired and A-tailed following the

protocol and ligated to indexed TruSeq adaptors. The final library was created by three PCR amplification using AMPure bead clean-up. Library quality and size were assessed by Agilent DNA 1000 series II assay and KAPA Library Quantification protocol. The two fractions were pooled for sequencing paired end 100 bp using Illumina HiSeq2000.

The construction of the 12 kb mate pair libraries was achieved by a hybrid 454/Illumina procedure. Briefly, 50 µg of genomic DNA (fragmented for 20 cycles at speed code 12 using a Hydroshear, Marburg, MA) equipped with a large excess of modular 18 kb mate pair end adaptors (Illumina) were extracted. Biotinylated circularization adaptors (Titanium Paired-end Adapter set (454 Life Sciences, Roche CT)) were added to ends of the gel-extracted fragments. Recombination of the ends was performed with Cre recombinase (Engle, Ipswich, MA), and linear molecules removed were removed with Plasmid Safe (Epiconic, Madison, WI). Molecules were fragmented again and 12 kb end adaptor-tagging was performed with the GS Rapid Library Preparation Kit (Illumina, San Diego, CA). Biotinylated end adaptors (Illumina) were ligated to the repaired/A-tailed ends. Biotinylates were enriched using Streptavidin-coupled Dynabeads (Life Technologies, Grand Island, NY) and amplified by PCR using Illumiprimer.

Recombinant bacterial artificial chromosome (BAC) was performed using DNA from the same genome on a Titanium with 6 kb paired-end libraries at the Plat-Forme Génomique of the Institute for Systems and Integrative Biology, Laval, Quebec, QC. A single paired-end prepared on a standard assembly at speed code 12. 10-kb fragments were isolated from GS-FLEX library and were ligated to the mate-pair end adaptors using 8% PAGE gel. DNA quality was assessed and quantified using an Agilent DNA 1000 series II assay (Agilent) and Nanoplot spectrophotometer (Nanodrop). DNA was subsequently diluted to 8 µM. The final concentration was confirmed using a Quant-iT IT dDNA HS assay kit and Qubit fluorometer (Invitrogen).

A new reagent tray that opens the snap-hood latches cartridge together (Supplementary Figs S1B and S2), give the reagent tubes, yet allowing the tube to be put without damage to its components (Supplementary Fig. 40), the stock length-dependent reagent containers allow a maximum of ~650 cycles in total. To maximize the potential of kit approach, a new reagent tray with 70 ml wells was placed in a modified clamshell style.

Assembling the 20 Gb white spruce genome

Assembling the genome was performed using the ABYSS algorithm (Simpson *et al.*, 2009), which captures a representation of read-to-read overlaps by a distributed de Bruijn graph and uses parallel computations to build the target genome. The modular nature of the tool allowed us to execute a large number of tests, tune the memory parameters for a successful run, train the assembly parameters for a successful assembly and quantify the utility of long reads for large genome assemblies. To the best of our knowledge, the ABYSS algorithm is unique in its ability to enable genome assemblies of this scale using whole-genome shotgun sequencing data.

After sequencing, the raw data were assembled using ABYSS (Simpson *et al.*, 2009) and the resulting contig map was visualized using CIRCOS (Krzywinski *et al.*, 2009). The contig map was then used to generate a genome browser using Galaxy (Galaxy Team, 2012).

Annotations were performed using the RAST server (Overbeek *et al.*, 2009) and the results were visualized using Circos (Krzywinski *et al.*, 2009).

The genome was annotated using the RAST server (Overbeek *et al.*, 2009) and the results were visualized using Circos (Krzywinski *et al.*, 2009).

The genome was annotated using the RAST server (Overbeek *et al.*, 2009) and the results were visualized using Circos (Krzywinski *et al.*, 2009).

Assembling the 20 Gb white spruce (*Picea glauca*) genome from whole-genome shotgun sequencing data

Inanc Birol^{1,2,3,*}, Anthony Raymond¹, Shaun D Jackman¹, Stephen Pleasance¹, Robin Coop¹, Greg A Taylor¹, Macaire Man Saint Yuen⁴, Christopher I Keeling⁴, Dana Brand¹, Benjamin P Vandervalk¹, Heather Kirk¹, Pawan Pandoh¹, Richard A Moore¹, Yongjun Zhao¹, Andrew J Mungall¹, Barry Jaquish⁵, Alvin Yanchuk⁶, Carol Ritland^{4,6}, Brian Boyle^{7,8}, Jean Bousquet^{7,8}, Kermit Ritland⁶, John MacKay^{7,8}, Jörg Bohlmann¹, Steven JM Jones^{1,2,9}

¹British Columbia Cancer Agency, Genome Sciences Centre, Vancouver, BC V5Z 4S6

²University of British Columbia, Department of Medical Genetics, Vancouver, BC V6H 3N1

³Simon Fraser University, School of Computing Science, Burnaby, BC V5A 1S6

⁴University of British Columbia, Michael Smith Laboratories, Vancouver, BC V6T 1Z4

⁵British Columbia Ministry of Forests, Lands and Natural Resource Operations, Victoria, BC V8W 9C2

⁶University of British Columbia, Department of Forest Sciences, Vancouver, BC V6T 1Z4

⁷Université Laval, Institute for Systems and Integrative Biology, Québec, QC G1V 0A6

⁸Université Laval, Department of Wood and Forest Sciences, Québec, QC G1V 0A6

⁹Simon Fraser University, Department of Molecular Biology and Biochemistry, Burnaby, BC V5A 1S6

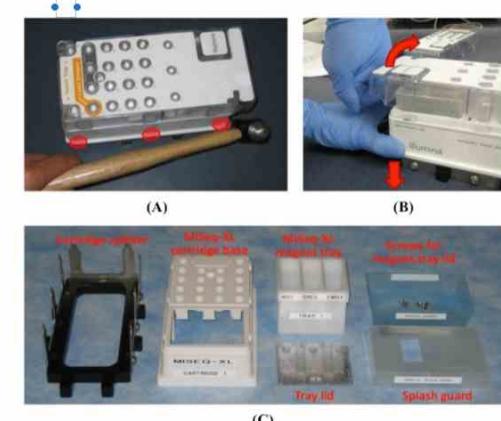


Figure S1. Modification of the MiSeq cartridge. MiSeq reagent cartridge was modified to allow for longer read lengths. (A, B) Opening of the clamshell style cartridge. (C) Contents of the modified cartridge. This was initially used to combine two PE150 kits for PE300 runs. When Illumina introduced the P250 kit, the same apparatus was used to enable PE500 runs.

Paired-end and Mate-pairs

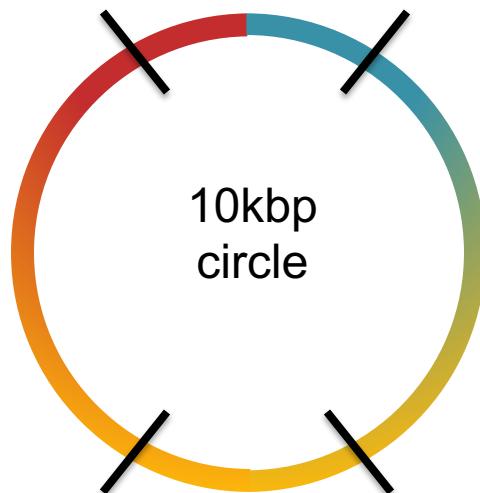
Paired-end sequencing

- 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



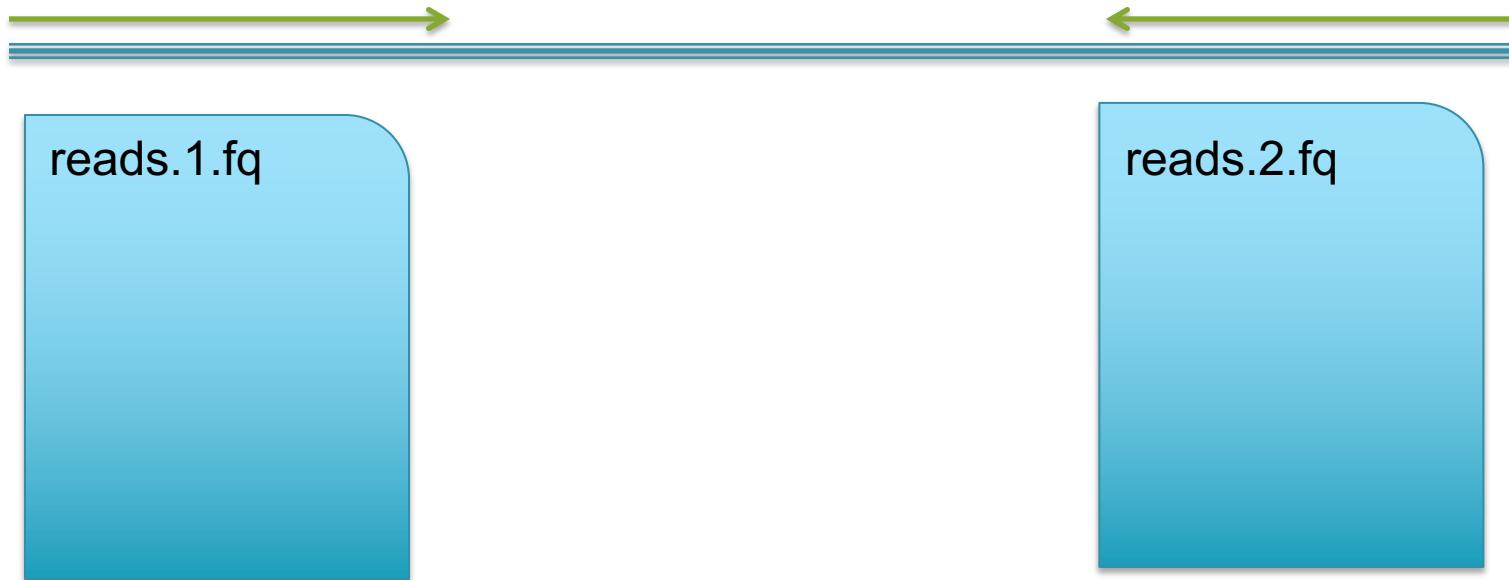
2x100 @ ~10kbp (outies)



2x100 @ 300bp (innies)

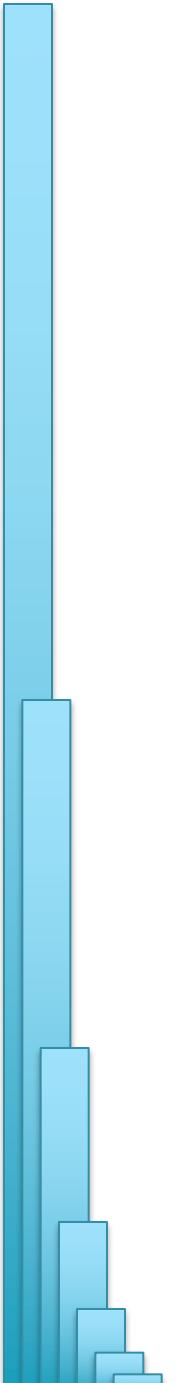


FASTQ Files

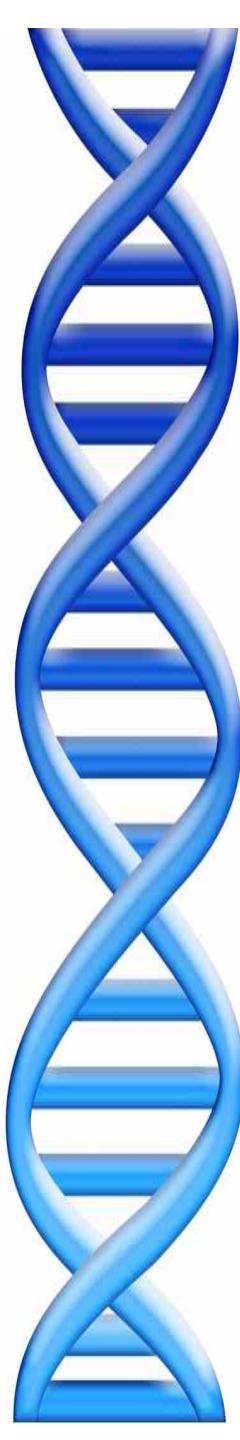


```
@SEQ_ID  
GATTGGGGTTCAAAGCAGTATCGATCAAATAGTAAATCCATTGTTCAACTCACAGTTT  
+  
! ' *((( (**+) ) % % + +) ( % % % ) . 1 *** - + * ' ) ) **55CCF>>>>CCCCCCCC65
```

@Identifier
Sequence
+Separator
Quality Values
...



Part 2: De novo genome assembly



Outline

1. Assembly theory

- Assembly by analogy

2. Practical Issues

- Coverage, read length, errors, and repeats

3. Whole Genome Alignment

- MUMmer recommended

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 copies x 138,656 words / 5 words per fragment = 138k 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]

Model the assembly problem as a graph problem

How long will it take to compute the overlaps?

de Bruijn Graph Construction

- $G_k = (V, E)$
 - V = Length- k sub-fragments
 - E = Directed edges between consecutive sub-fragments
 - Sub-fragments overlap by $k-1$ words

Fragments $|f|=5$

It was the best of

was the best of times

Sub-fragment $k=4$

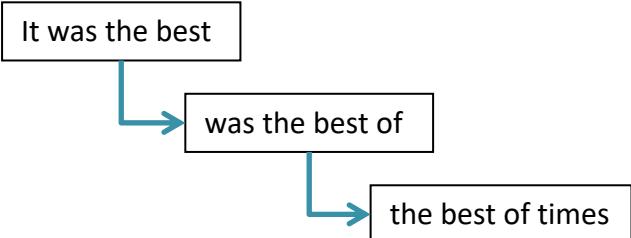
It was the best

was the best of

was the best of

the best of times

Directed edges (overlap by $k-1$)



- Overlaps between fragments are implicitly computed

How to pronounce:

https://forvo.com/word/de_bruijn/

de Bruijn, 1946

Idury et al., 1995

Pevzner et al., 2001

de Bruijn Graph Assembly

It was the best

was the best of

the best of times,

best of times, it

of times, it was

times, it was the

it was the worst

was the worst of

the worst of times,

worst of times, it

it was the age

was the age of

the age of foolishness

the age of wisdom,

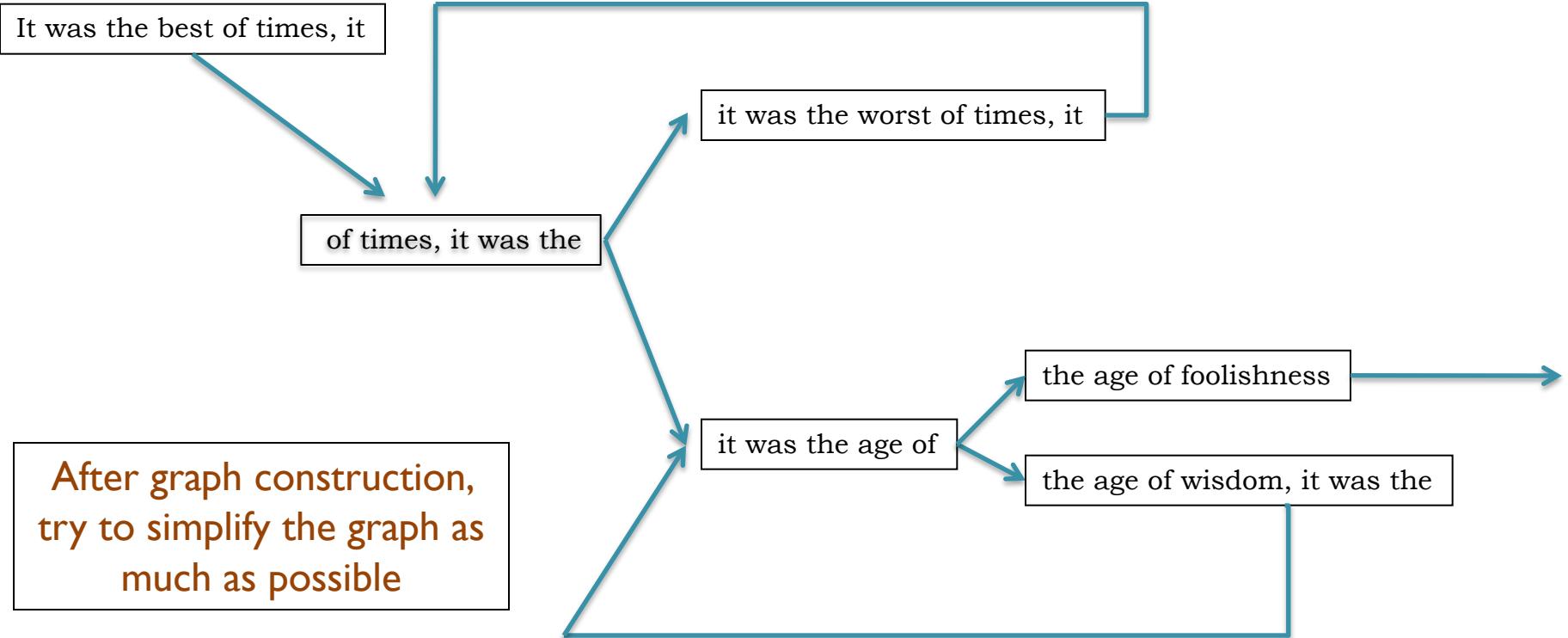
age of wisdom, it

of wisdom, it was

wisdom, it was the

After graph construction,
try to simplify the graph as
much as possible

de Bruijn Graph Assembly



The full tale

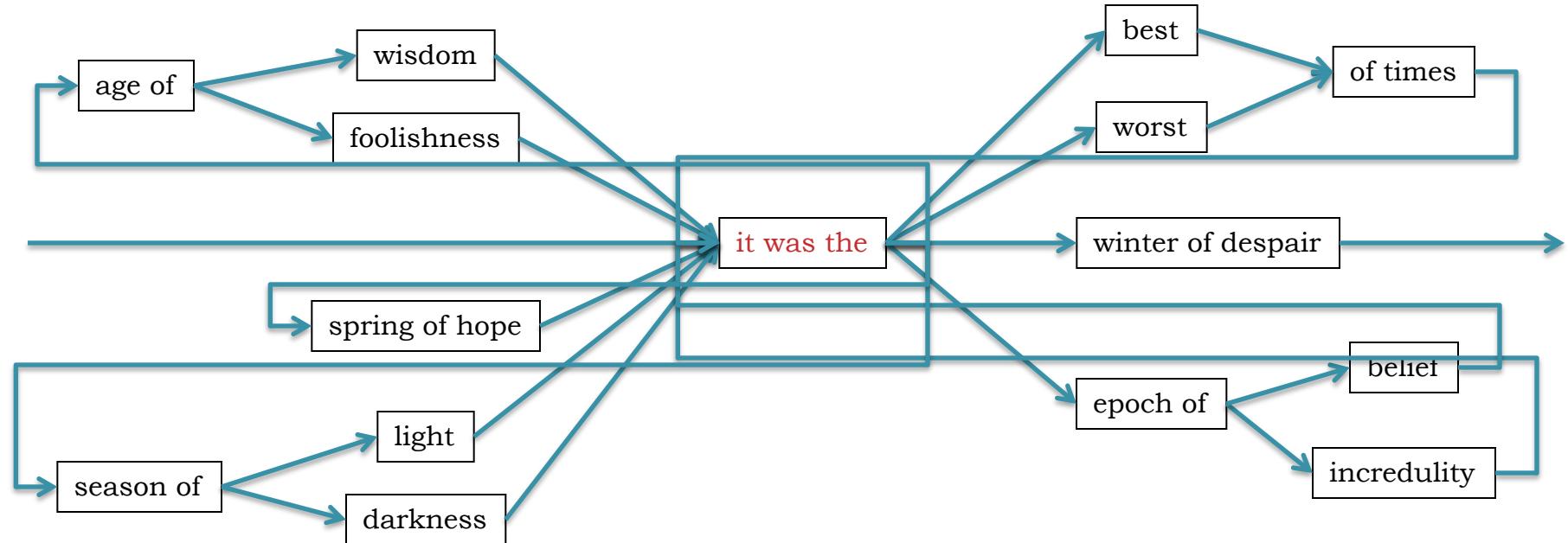
... it was the best of times it was the worst of times ...

... it was the age of wisdom it was the age of foolishness ...

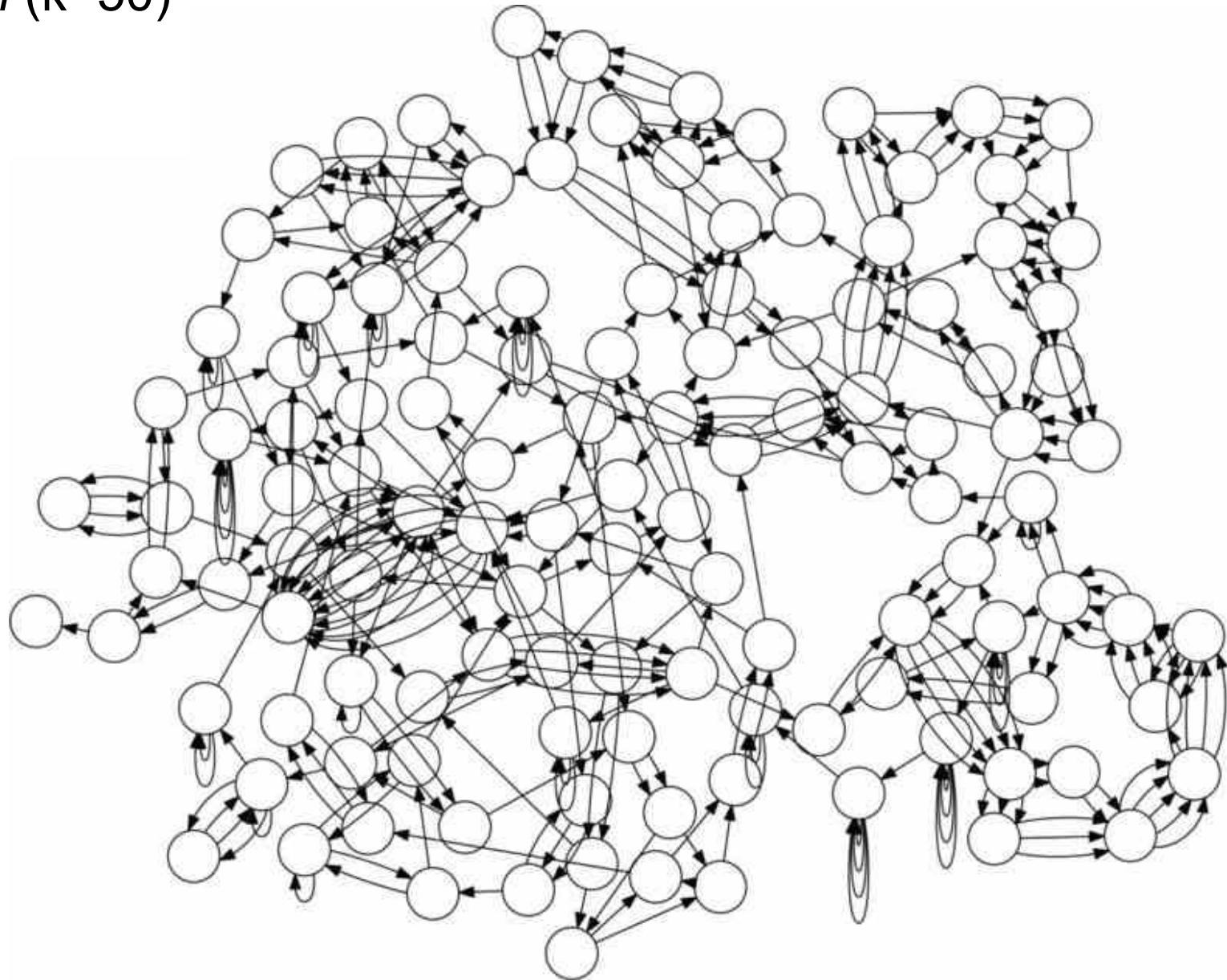
... it was the epoch of belief it was the epoch of incredulity ...

... it was the season of light it was the season of darkness ...

... it was the spring of hope it was the winter of despair ...



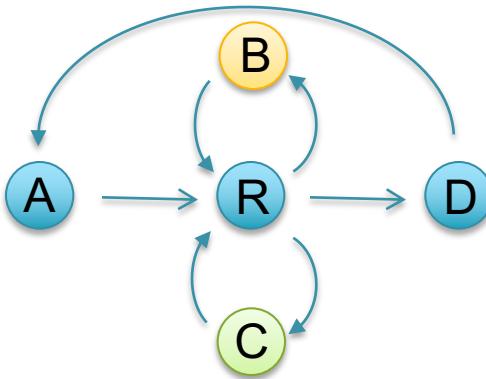
E. coli (k=50)



Reducing assembly complexity of microbial genomes with single-molecule sequencing

Koren et al (2013) Genome Biology. 14:R101 <https://doi.org/10.1186/gb-2013-14-9-r101>

Counting Eulerian Cycles



ARBRCRD
or
ARCRBRD

Generally an exponential number of compatible sequences

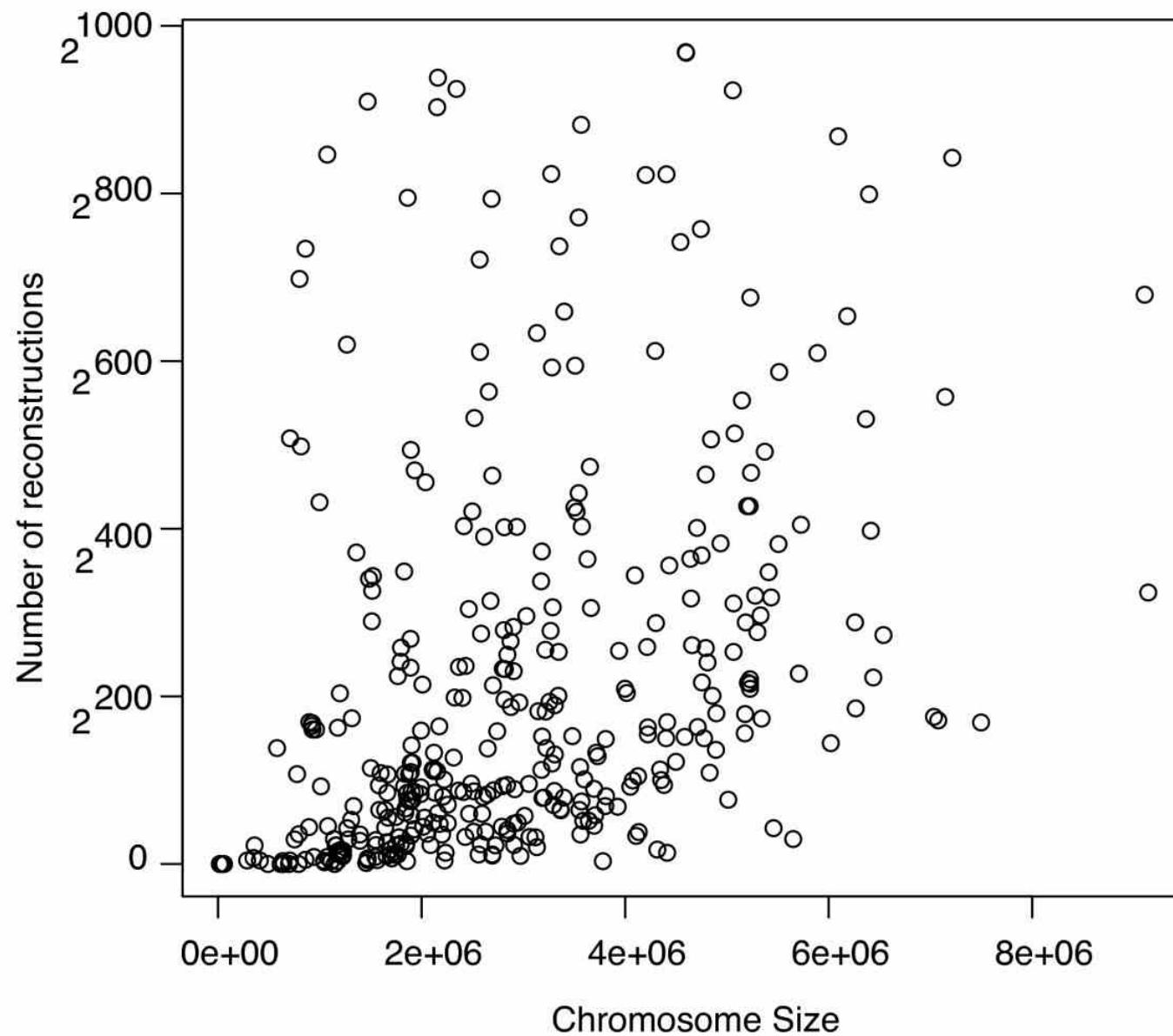
- Value computed by application of the BEST theorem (Hutchinson, 1975)

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

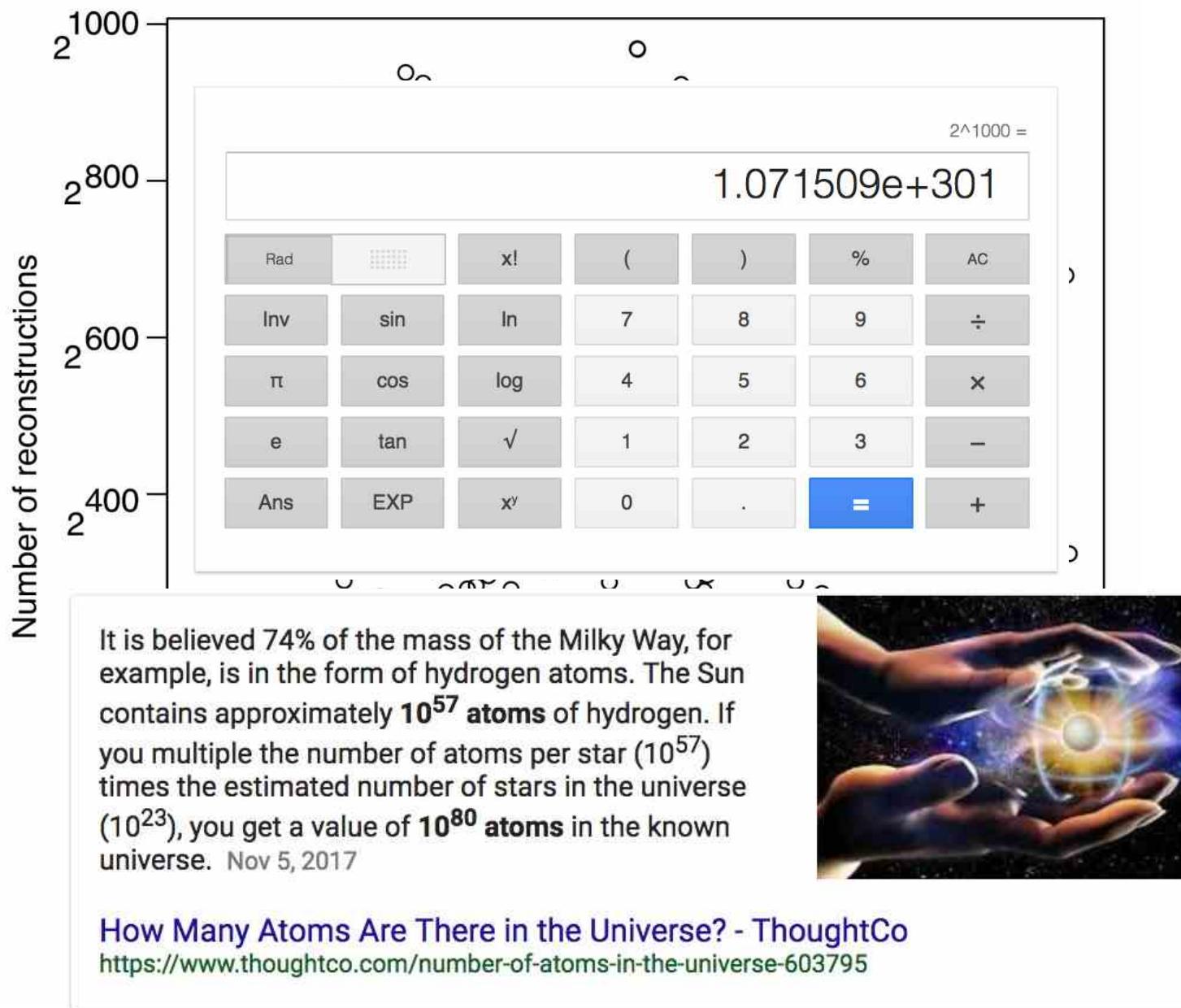
L = $n \times n$ matrix with $r_u - a_{uu}$ along the diagonal and $-a_{uv}$ in entry uv

$r_u = d^+(u) + 1$ 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*.



Assembly Complexity of Prokaryotic Genomes using Short Reads.

Kingsford C, Schatz MC, Pop M (2010) *BMC Bioinformatics*.



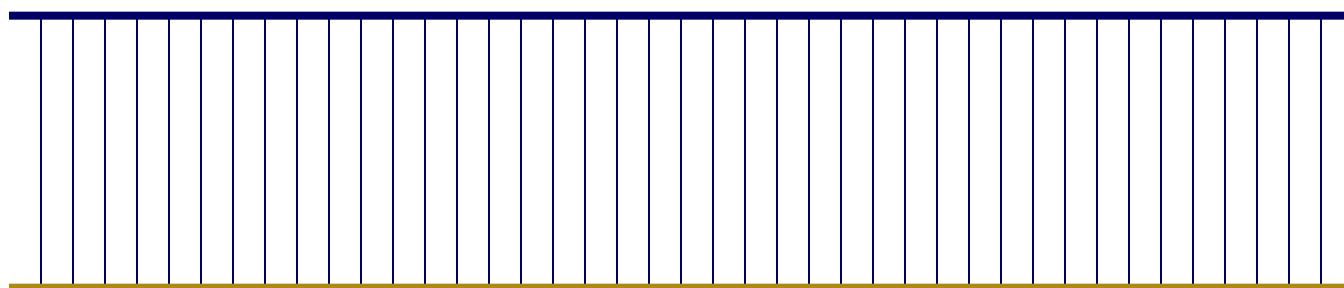
Whole Genome Alignment with MUMmer

Slides Courtesy of Adam M. Phillippy
NHGRI

Goal of WGA

- For two genomes, A and B , find a mapping from each position in A to its corresponding position in B

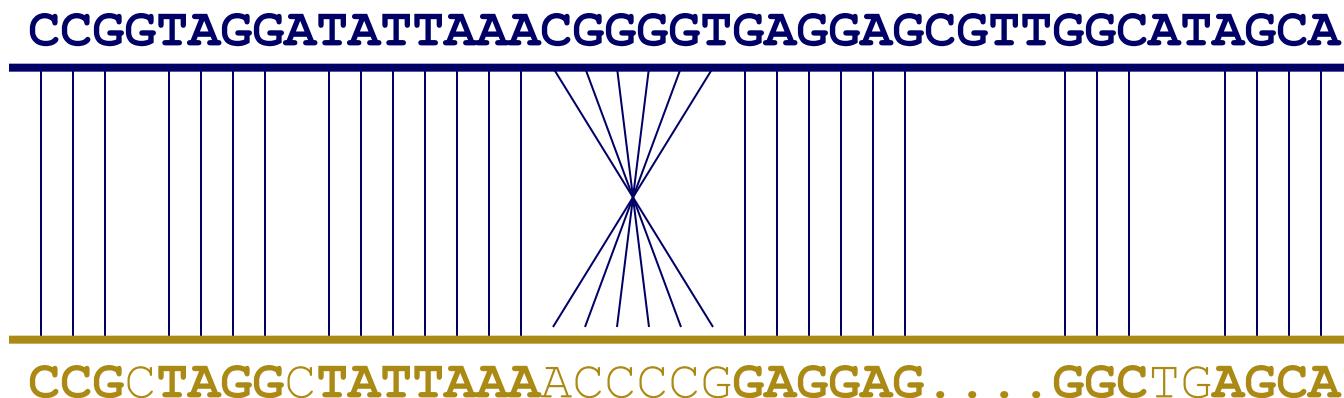
CCGGTAGGCTATTAAACGGGGTGAGGAGCGTTGGCATAGCA



CCGGTAGGCTATTAAACGGGGTGAGGAGCGTTGGCATAGCA

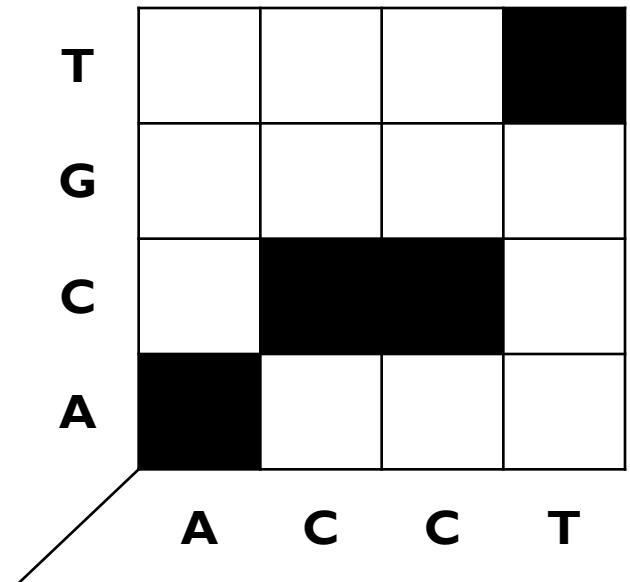
Not so fast...

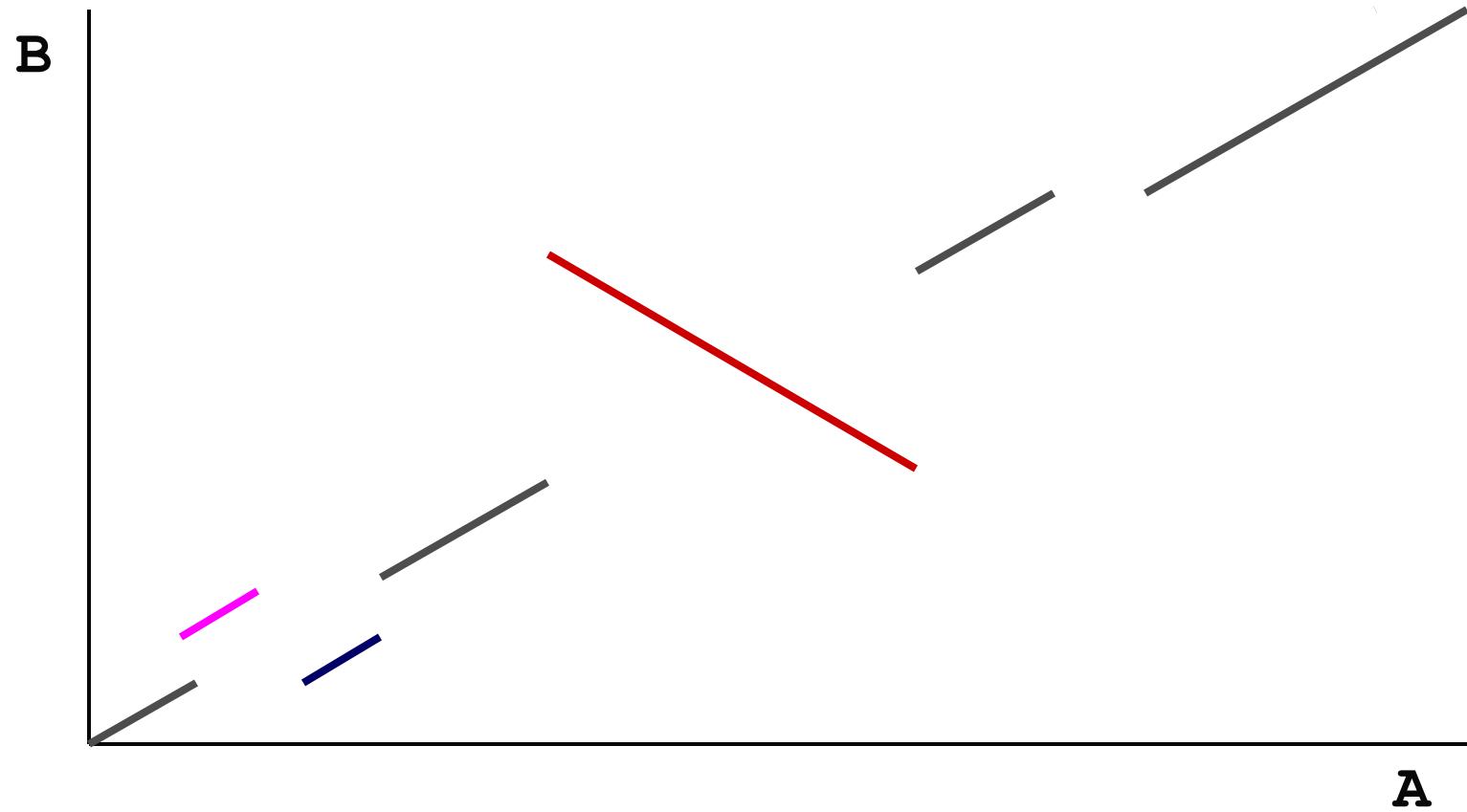
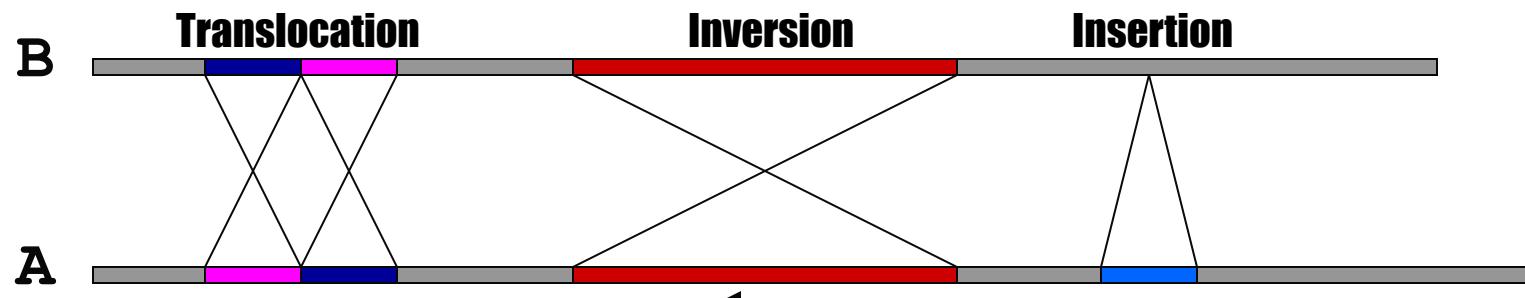
- Genome A may have insertions, deletions, translocations, inversions, duplications or SNPs with respect to B (sometimes all of the above)



WGA visualization

- How can we visualize *whole genome* alignments?
- With an alignment dot plot
 - $N \times M$ matrix
 - Let i = position in genome A
 - Let j = position in genome B
 - Fill cell (i,j) if A_i shows similarity to B_j
 - A perfect alignment between A and B would completely fill the positive diagonal

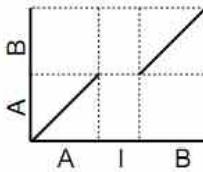




SV Types

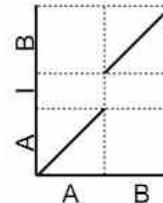
Insertion into Reference

R: AIB
Q: AB



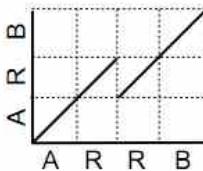
Insertion into Query

R: AB
Q: AIB



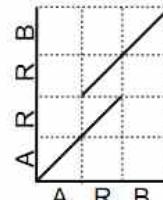
Collapse Query

R: ARRB
Q: ARB



Collapse Reference

R: ARB
Q: ARRB



Collapse Query w/ Insertion

R: ARIRB
Q: ARB

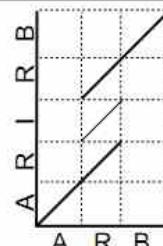
Exact tandem alignment if I=R



Collapse Reference w/ Insertion

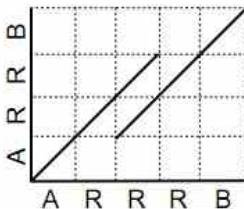
R: ARB
Q: ARIRB

Exact tandem alignment if I=R



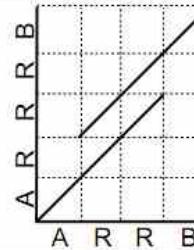
Collapse Query

R: ARRRB
Q: ARRB



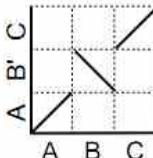
Collapse Reference

R: ARRB
Q: ARRRB



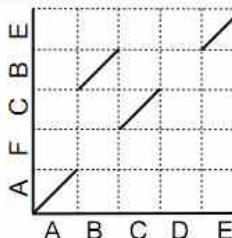
Inversion

R: ABC
Q: ABC



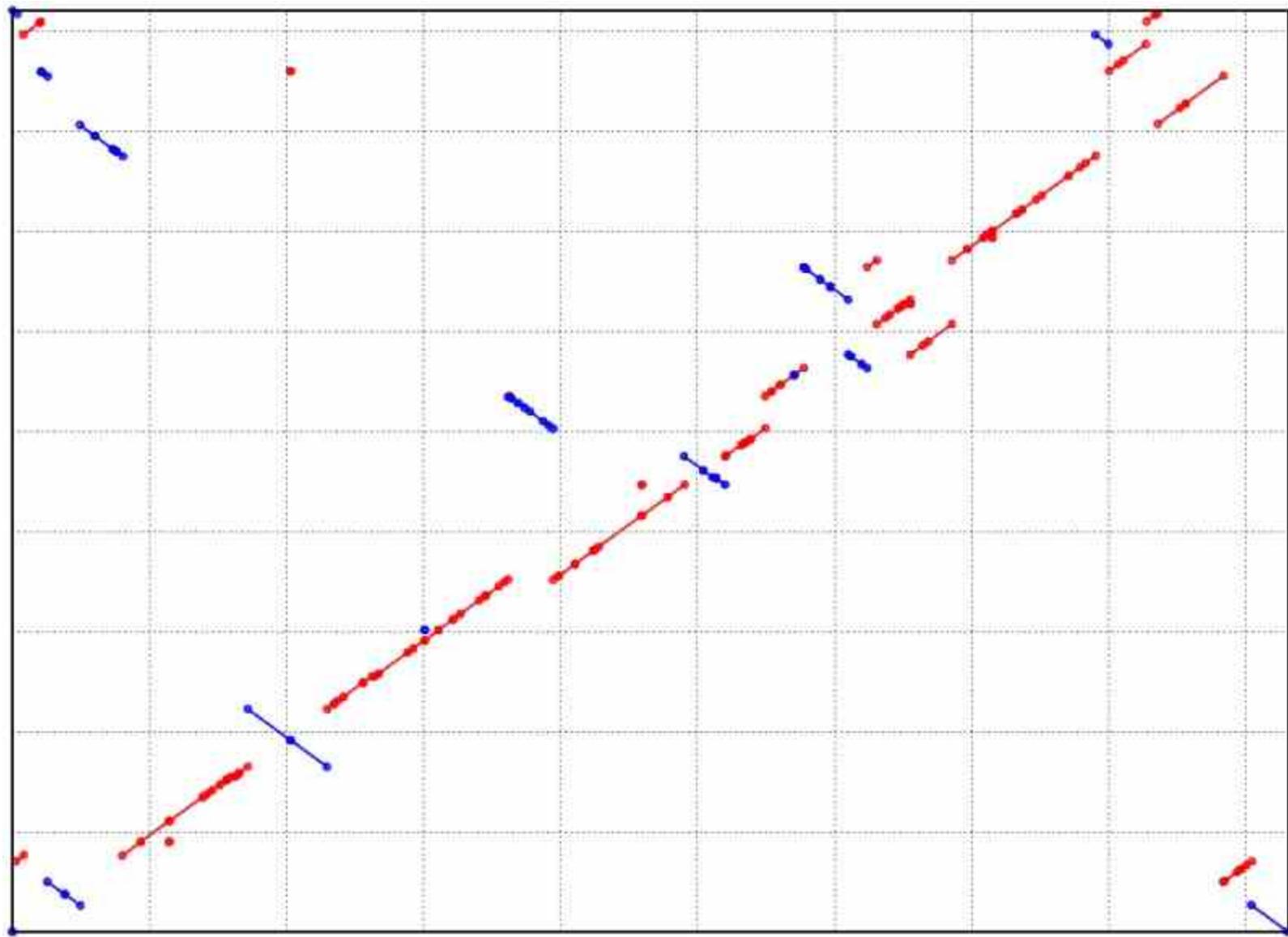
Rearrangement w/ Disagreement

R: ABCDE
Q: AFCBE



- Different structural variation types / misassemblies will be apparent by their pattern of breakpoints

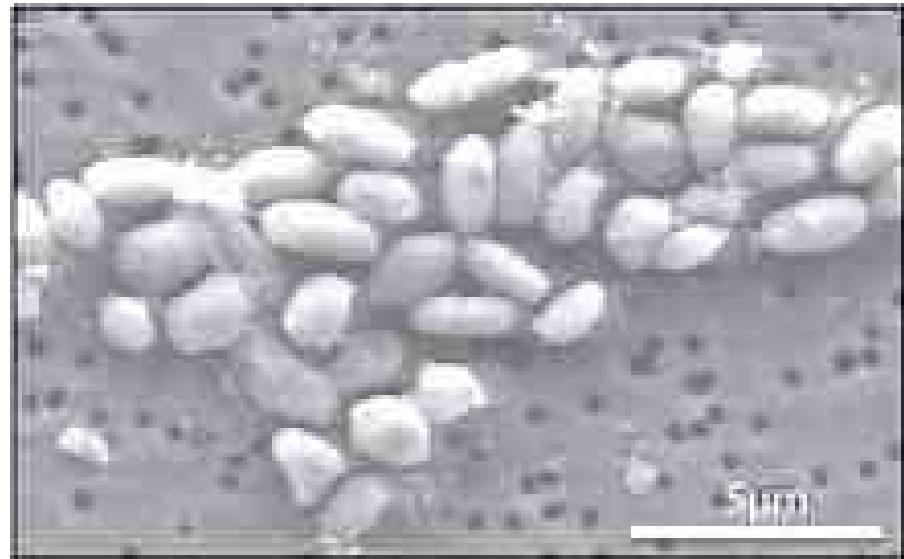
- Most breakpoints will be at or near repeats
- Things quickly get complicated in real genomes



Alignment of 2 strains of *Y. pestis*

<http://mummer.sourceforge.net/manual/>

Halomonas sp. GFAJ-1



Library 1: Fragment

Avg Read length: 100bp

Insert length: 180bp

Library 2: Short jump

Avg Read length: 50bp

Insert length: 2000bp

A Bacterium That Can Grow by Using Arsenic Instead of Phosphorus
Wolfe-Simon et al (2010) *Science*. 332(6034):1163-1166.

Digital Information Storage

Decoding self-referential DNA that encodes these notes.

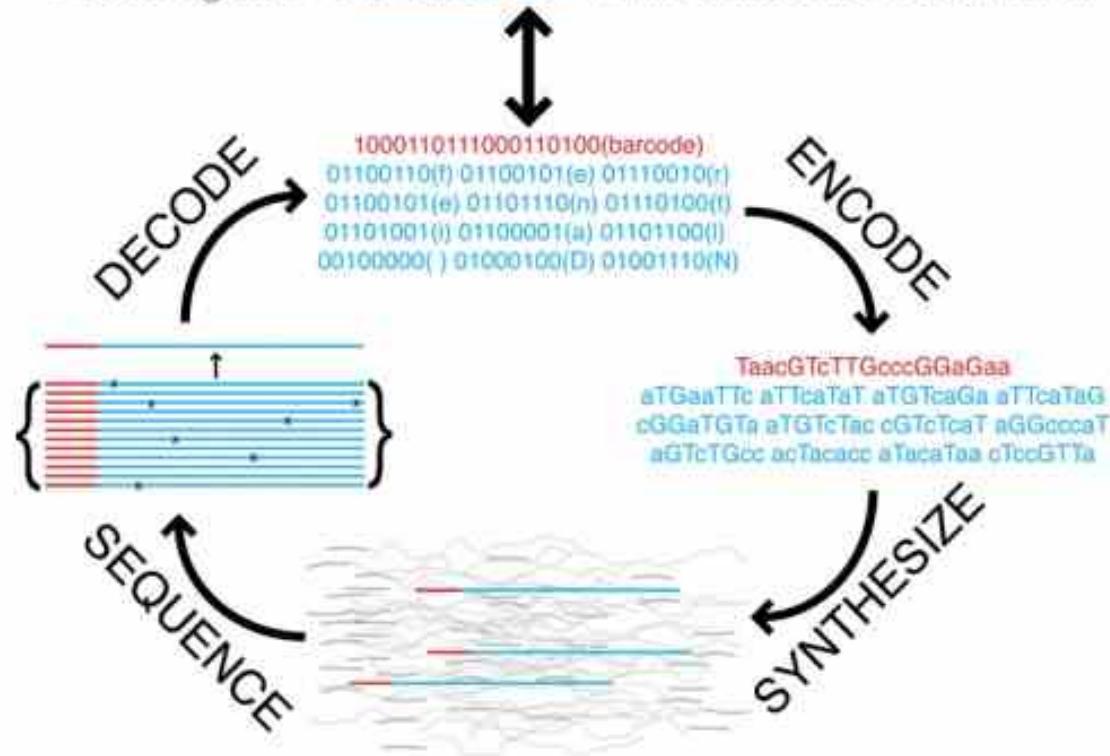


Fig. S1. Schematic of DNA information storage.

Encoding/decoding algorithm implemented in dna-encode.pl from David Dooling.

Next-generation Digital Information Storage in DNA

Church et al (2010) Science. 337(6102)1628

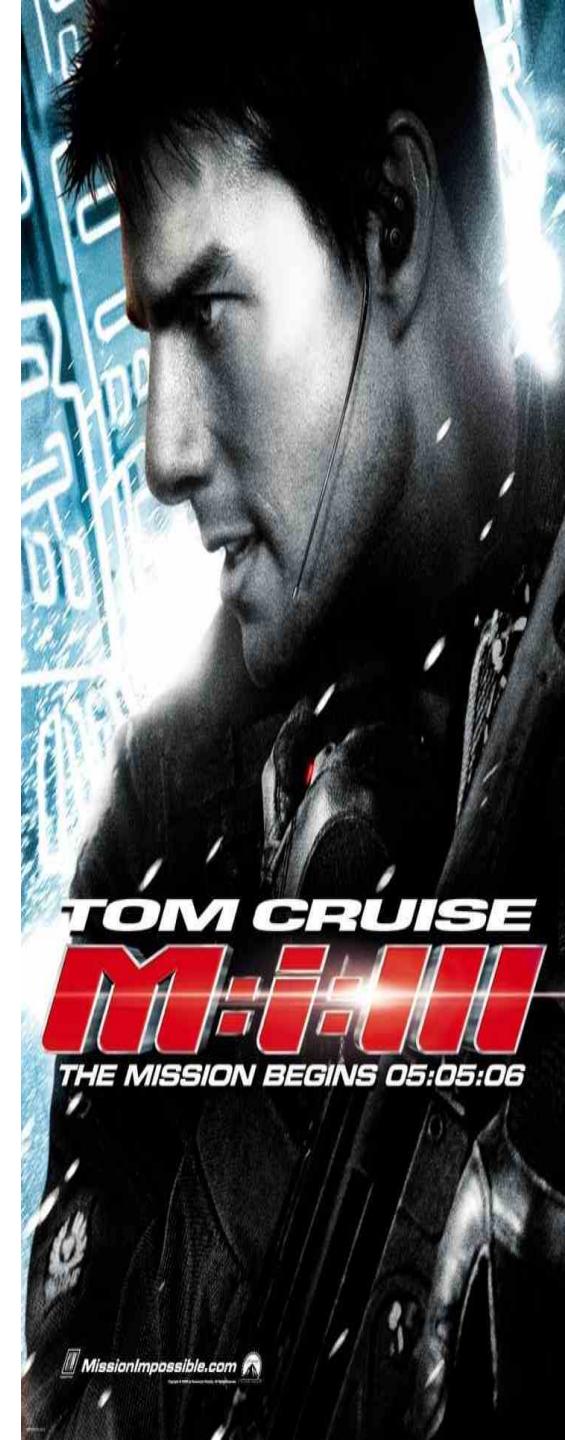
Assignment 2: Genome Assembly

Due Wednesday Sept 13 @ 11:59pm

- 1. Setup Conda/Docker/Ubuntu**
- 2. Initialize Tools**
- 3. Download Reference Genome & Reads**
- 4. Decode the secret message**

1. Estimate coverage, check read quality
2. Check kmer distribution
3. Assemble the reads with *spades*
4. Align to reference with *MUMmer*
5. Extract foreign sequence
6. *dna-encode.pl -d*

<https://github.com/schatzlab/appliedgenomics2023/blob/master/assignments/assignment2/README.md>



Find and decode

nucmer --maxmatch ref.fasta contigs.fasta

--maxmatch Find maximal exact matches (MEMs) without repeat filtering
-p refctg Set the output prefix for delta file

mummerplot --layout --png out.delta

--layout Sort the alignments along the diagonal
--png Create a png of the results

show-coords -rclo out.delta

-r Sort alignments by reference position
-c Show percent coverage
-l Show sequence lengths
-o Annotate each alignment with BEGIN-END/CONTAINS

samtools faidx contigs.fasta

Index the fasta file

samtools faidx contigs.fasta contig_XXX:YYY-ZZZ > msg.fa
dna-decode.py -d -input msg.fa

Assignment 2: Genome Assembly

Due Wednesday Sept 13 by 11:59pm

The screenshot shows a GitHub repository interface. On the left, a sidebar titled "Files" lists the contents of the "main" directory, including "assignments", "assignment1", "assignment2" (which is expanded to show "README.md", "asm.tgz", and "lectures"), "LICENSE", and "README.md". The main pane displays the "README.md" file for "Assignment 2: Genome Assembly". The file content includes:

Assignment 2: Genome Assembly

Assignment Date: Wednesday, September 6, 2023
Due Date: Wednesday, September 13, 2023 @ 11:59pm

Assignment Overview

In this assignment, you are given a set of unassembled reads from a mysterious pathogen that contains a secret message encoded someplace in the genome. The secret message will be recognizable as a novel insertion of sequence not found in the reference. Your task is to assess the quality of the reads, assemble the genome, identify, and decode the secret message. If all goes well the secret message should decode into a recognizable english text, otherwise doublecheck your coordinates and try again. As a reminder, any questions about the assignment should be posted to [Piazza](#).

For this assignment, we recommend you install and run the tools using [bioconda](#). There are some tips below in the Resources section. Note on Mac, we highly recommend you install the x86_64 package even if you are using an M1 chip.

Question 1. Coverage Analysis [20 pts]

Download the reads and reference genome from:
<https://github.com/schatzlab/appliedgenomics2023/blob/main/assignments/assignment2/asm.tgz?raw=true>

Note we have provided both paired-end and mate-pairs reads (see included README for details). Make sure to look at all of the reads for the coverage analysis and kmer analysis, as well as in the assembly.

- Question 1a. How long is the reference genome? [Hint: Try `samtools faidx`]
- Question 1b. How many reads are provided and how long are they? Make sure to measure each file separately [Hint: Try `FastQC`]
- Question 1c. How much coverage do you expect to have? [Hint: A little arithmetic]
- Question 1d. Plot the average quality value across the length of the reads [Hint: Screenshot from `FastQC`]

Question 2. Kmer Analysis [20 pts]

Use `Jellyfish` to count the 21-mers in the reads data. Make sure to use the "-C" flag to count canonical kmers, otherwise your analysis will not correctly account for the fact that your reads come from either strand of DNA.

- Question 2a. How many kmers occur exactly 50 times? [Hint: try `jellyfish histo`]

<https://github.com/schatzlab/appliedgenomics2023/tree/main/assignments/assignment2>

Check Piazza for questions!