

	The Mapping Problem and the Assembly Problem
The Mapping	g Problem
INPUT:	m reads $S_1, \ldots, S_m$ of length l and an approximate reference genome R.
QUESTION:	What are the positions $x_1,, x_m$ along R where each read $S_1,, S_m$ matches, respectively?
The Assembl	y Problem
<b>INPUT:</b>	m reads $S_1, \ldots, S_m$ of length l.
<b>QUESTION:</b>	What is the sequence of the full genome?
The crucial di assembly is reference sequence d	fference between the problems of mapping and s that in case of assembly we do not have a genome, and we must assemble the full irectly from the reads.

	The Mapping Problem
The Short Rea	d Mapping Problem on reference genome R
INPUT:	m reads $S_1, \ldots, S_m$ of length 1 and an approximate reference genome R.
QUESTION:	What are the positions $x_1,, x_m$ along R where each read $S_1,, S_m$ matches, respectively?
For example: af existing seq	ter sequencing the genome of a person we want to map it to an <b>uence of the human genome.</b>
• The new sam – natural va – mismatch – sequencir	nple will not be 100% identical to the reference genome: ritation in the population es or gaps og errors
<ul> <li>repetitive</li> <li>Humans are</li> <li>different =</li> <li>two slight</li> </ul>	regions e diploid organisms alleles on the maternal and paternal chromosomes thy different reads mapping to the same location (some with mismatches)
• Viruses: hig mapping pr	gh mutation rates, many variations haplotypes => hybrid oblem



# Applying the Burrows-Wheeler transform BW(T) to the text T = "the next text that i index.": 1. First, we generate all cyclic shifts of T. 2. Next, we sort these shifts lexicographically. a define the character '.' as the minimum and we assume that it appears exactly once, as the last symbol in the text. a followed lexicographically by ' (space) followed by the English letters according to their natural ordering. Call the resulting matrix M. The transform BW(T) is defined as the sequence of the last characters in the rows of M. Note that, the last column is a permutation of all characters in the text since each character appears in the last position in exactly one cyclic shift.







The following holds for BW(T):

- 1. # occurrences of char c in T = # occurrences of char c in BW(T) (BW(T) permutation of the T).
- 2. The first column of the matrix M can be obtained by sorting BW(T) lexicographically.
- 3. Determine the number of occurrences of the substring 'xt' in T:
  - **BW**(**T**) is the last column of the lexicographical sorting of the shifts.
  - The character at the last position of a row appears in the text T immediately prior to the first character in the same row (each row is a cyclical shift).
  - => consider the interval of 't' in the first column, and check how many of these rows have an 'x' at the last position.







Burrows-who	eeler Transform
Given $\overline{BW(T)}$ also the second column of	can be derived:
• 'xt' appears <i>twice</i> in the text, and <i>thre</i> 'x'.	ee rows start with an
• <i>Two of the three</i> must be followed by lexicographical sorting determines w	y a 't', where the which 'x'.
<ul> <li>The third 'x' is followed by a '.' (see a follow the first 'x' in the first column lexicographically than 't'</li> </ul>	first row) => '.' must since '.' is smaller
• The second and third occurrences of	'v' in the first column
are therefore followed by 't'.	haracters at the second
<ul> <li>The second and third occurrences of are therefore followed by 't'.</li> <li>Note: We can use the same process to recover the c column for each interval, and then the third, etc.</li> <li>. the next text that I index</li> </ul>	haracters at the second $1^{\text{st}}$ x
<ul> <li>The second and third occurrences of are therefore followed by 't'.</li> <li>Note: We can use the same process to recover the c column for each interval, and then the third, etc</li></ul>	haracters at the second
<ul> <li>The second and third occurrences of are therefore followed by 't'.</li> <li>Note: We can use the same process to recover the c column for each interval, and then the third, etc</li></ul>	haracters at the second $1^{st} x$ $2^{nd} x$
<ul> <li>The second and third occurrences of are therefore followed by 't'.</li> <li>Note: We can use the same process to recover the c column for each interval, and then the third, etc.         <ul> <li>the next text that I index</li> <li>t text that i index. the next text</li> <li>t that i index. the next text</li> </ul> </li> </ul>	haracters at the second $1^{st} x$ $2^{nd} x$ $3^{rd} x$
<ul> <li>The second and third occurrences of are therefore followed by 't'.</li> <li>Note: We can use the same process to recover the c column for each interval, and then the third, etc.         <ul> <li>the next text that I index</li> <li>t text that i index. the next text that i index. the next text that i index. the next text that i index. the next</li> </ul> </li> </ul>	haracters at the second $1^{st} x$ $2^{nd} x$ $3^{rd} x$
<ul> <li>The second and third occurrences of are therefore followed by 't'.</li> <li>Note: We can use the same process to recover the c column for each interval, and then the third, etc the next text that I index the next text that I index. the next that i index. the next that i index. the next text that i index.</li> </ul>	haracters at the second $1^{st} x$ $2^{nd} x$ $3^{rd} x$







	Burrows-Wheeler Transform
Proof:	
1. Follo cycli	we directly from the fact that each row in $M$ is a cal shift.
2. Let 2	$X_j$ denote the j-th occurrence of character X in L, and let $\alpha$ be the character following $X_j$ in the text
	and $\beta$ the character following $X_{i+1}$ .
The begi	n, since $X_j$ appears above $X_{j+1}$ in L, $\alpha$ appears at the nning of a row above the row that starts with $\beta$ .
The hence	rows are lexicographically ordered, ee $\alpha$ must be equal or lexicographically smaller than $\beta$ .
Now	$\alpha \leq_{\text{lexicographically}} X \beta$ holds.
Hen char be a	ce, as the rows are lexicographically ordered, if acter $X_j$ appears in F it is followed by $\alpha$ , and thus will bove $X_{j+1}$ which is followed by $\beta$ .
Hen char be a Thu	ce, as the rows are lexicographically ordered, if acter $X_j$ appears in F it is followed by $\alpha$ , and thus will bove $X_{j+1}$ which is followed by $\beta$ . s proofing the Lemma.



















	De Bruijn Graphs
Definition	
A k-dimensiona representing	l de Bruijn graph of n symbols is a directed graph overlaps between sequences of symbols.
It has n <sup>k</sup> vertices Note: the sar	s, consisting of all possible k-tuples of the given symbols. ne symbol may appear multiple times in a tuple.
If we have the so $V = \{ (a_1, \dots, (a_1, \dots, (a_1, \dots, (a_n, \dots, (a_$	et of symbols $A = \{a_1,, a_n\}$ then the set of vertices is: , $a_1, a_1$ , $(a_1,, a_1, a_2)$ ,, $(a_1,, a_1, a_n)$ , , $a_2, a_1$ ,,, $(a_n,, a_n, a_n)$ }
If a vertice w ca by one place a directed ed	n be expressed by shifting all symbols of another vertex v to the left and adding a new symbol at the end, then v has ge to w.
Thus the set of $c$ E = {( (v <sub>1</sub> , v <sub>2</sub> )	directed edges E is: $v_2,, v_k$ , $(w_1, w_2,, w_k)$ $ v_2 = w_1, v_3 = w_2,, v_k = w_{k-1}$ , and w. new}

















## Other Assembly Algorithms

- HMM based
- Majority based
- Etc.
- Long reads: string graphs

## Other Assembly Algorithms

K.R. Bradnan et al. Assemblathon 2: evaluating *de novo* methods of genome assembly in three vertebrate species (http://gigascience.biomedcentral.com/articles/10.1186/2047-217X-2-10, 2013)

"Many current genome assemblers produced useful assemblies, containing a significant representation of their genes and overall genome structure.

However, the high degree of variability between the entries suggests that there is still much room for improvement in the field of genome assembly and that:

approaches which work well in assembling the genome of one species may not necessarily work well for another."

### Other Assembly Algorithms

Salzberg SL, Phillippy AM, Zimin A, Puiu D, Magoc T, Koren S, Treangen TJ, Schatz MC, Delcher AL, Roberts M, et al. Gage: A critical evaluation of genome assemblies and assembly algorithms. Genome Res. 2012; 22(3):557–67.

Three conclusions:

**1.Quality:** data quality, rather than the assembler itself, has a dramatic effect on the quality of an assembled genome

2. Variability: the degree of contiguity of an assembly varies enormously among different assemblers and different genomes

3.Correctness: the correctness of an assembly also varies widely and is not well correlated with statistics on contiguity.

### Other Assembly Algorithms

Scaffolding and completing genome assemblies in real-time with nanopore sequencing

By Minh Duc Cao, Son Hoang Nguyen, Devika Ganesamoorthy, Alysha G. Elliott, Matthew A. Cooper & Lachlan J. M. Coin

Nature Communications 8, Article number: 14515 (2017)

"Long read sequencing technologies, for example Pacific Biosciences' (PacBio) and Oxford Nanopore MinION sequencing, allow users to generate reads spanning most repetitive sequences, which can be used to close gaps in fragmented assemblies."





#### Viral Quasispecies Assembly Viruses HIV, Zika, and Ebola: Goal: a viral quasispecies assembly presenting all of the viral haplotypes, and • their abundance rates. • Problems: -The number of different strains is usually unknown. Different strains can differ by only minor amounts of distinguishing mutations. - Abundance rates can be as low as the sequencing error rates - Reference genomes representing high-quality consensus genome sequences can be obsolete (as a result of great diversity and high mutation rates)





Dhred quality scores ()	are defined as a property which is log-	arithmically related to th	have calling error probabilities $\mathcal{P}^{[2]}$
$\Omega = -10 \log_{10} P$	are defined as a property willer is log-	anamically related to th	to base-canning error probabilities r.
$Q = -10 \log_{10} r$			
-9			
$P = 10^{\frac{3}{10}}$			
For example, if Phred as	ssigns a quality score of 30 to a base,	the chances that this b	ase is called incorrectly are 1 in 1000.
Phred quality sco	res are logarithmically linked to er	ror probabilities	
Phred Quality Score	Probability of incorrect base call	Base call accuracy	
10	1 in 10	90%	
20	1 in 100	99%	
30	1 in 1000	99.9%	
40	1 in 10,000	99.99%	
	1 in 100 000	99.999%	
50	1 11 100,000		









Table 1. Characteristics of benchmarking data set	Table 1.	Characteristics of benchmarking d	ata set
---	----------	-----------------------------------	---------

Data set	Virus type	Genome length (bp)	Average coverage	Strain count	Strain abundance	Pairwise divergence
600× HIV mix	HIV-1	9478-9719	600×	5	20%	1%-6%
5-strain HIV mix	HIV-1	9478-9719	20.000×	5	5%-28%	1%-6%
10-strain HCV mix	HCV-1a	9273-9311	20.000×	10	5%-19%	6%-9%
3-strain ZIKV mix	ZIKV	10.251-10.269	20.000×	3	16%-60%	3%-10%
15-strain ZIKV mix	ZIKV	10.251-10.269	20.000×	15	2%-13%	1%-10%
Lab mix	HIV-1	9478-9719	20,000×	5	10%-30%	1%-6%

For each benchmark, we specify virus type, genome length, average coverage, strain count, relative abundance, and pairwise divergence. For the 600× HIV mix, the strains were homogeneously distributed with a relative abundance of 20% each.



## Bibliography

- Ben Langmead, Cole Trapnell, Mihai Pop and Steven L. Salzberg. Ultrafast and memory-effcient alignment of short DNA sequences to the human genome. Genome Biology, 2009.
- [2] Michael Burrows and David Wheeler. A block sorting lossless data compression algorithm. Technical Report 124, Digital Equipment Corporation, 1994.
- [3] Brent Ewing and Phil Green. Base-calling of automated sequencer traces using phred. II. Error probabilities. Genome Research, 1998.
- [4] Paulo Ferragina and Giovani Manzini. Opportunistic data structures with applications. FOCS '00 Proceedings of the 41st Annual Symposium on Foundations of Computer Science, 2000.
- [5] Heng Li, Jue Ruan and Richard Durbin. Mapping short DNA sequencing reads and calling variants using mapping quality scores. Genome Research, 2008.
- [6] Daniel R. Zerbino and Ewan Birney. Velvet: algorithms for de novo short read assembly using de Bruijn graphs. Genome Research, 2008.
- [7] Ron Shamir, Computational Genomics Fall Semester, 2010 Lecture 12: Algorithms for Next Generation Sequencing Data, January 6, 2011 Scribe: Anat Gluzman and Eran Mick
- [8] Minh Duc Cao, Son Hoang Nguyen, Devika Ganesamoorthy, Alysha G. Elliott, Matthew A. Cooper & Lachlan J. M. Coin. Scaffolding and completing genome assemblies in realtime with nanopore sequencing. Nature Communications 8, Article number: 14515, 2017.