

Genome Assembly

Sijia Huo Sean Kelly Gregory Raskind

University of Michigan

University of Illinois

14 July 2017

Background

Genomics
BDSI

Sijia Huo,
Sean Kelly,
Gregory
Raskind

In our cells, DNA carries genetic information:

- Each DNA strand contains a sequence of A, C, G, and Ts
- The human genome is around 3 billion base pairs long

Types of DNA Variation:

- Single base variations
- Copy number variations

```
>HSBGPG
GGCAGATTCCCCCTAGACCCGCCCCGCACCATGGTCAG
TGGGCACAGCCCCAGAGGGTATAAACAGTGCTGGAGGC
AGTCCTGAGCAGCAGCCAGCGCAGCCACCGAGACAC
CTCGCCCTATTGGCCCTGGCCGCACTTTGCATCGCTG
CCACCTCCCCTCAGGCCGCATTGCAGTGGGGGCTGAG
CACCTCTTCTCACCCCTTTGGCTGGCAGTCCCTTTGC
AGGCTCAATCCATTTGCCCCAGCTCTGCCCTTGCAAG
AGCTGCCCGAGACGCAGGGGAAGGAGGATGAGGGCCC
ACCAGGCTCCCTTTCTTTGCAAGTGCGAAGCCGAGC
GTGCAAGGTATGAGGATGGACCTGATGGGTTCTTGAC
CCCTCAGTCTCATTCCCCCACTCCTGCCACCTCCTGT
GCCTGCTCCCCACCTGATCCTCCCAACCCAGAGCCA
CTCCACAGCCTTTGTGTCCAAGCAGGAGGGCAGCGAG
GCTACCTGTATCAATGGCTGGGGTGAGAGAAAAGGCA
```

```
@H3GFVCCXX150415:8:2224:9627:35467/1
GGGAATTTTAACTGGCAAACTCAGAACTCCATCCAAAC.
+
AAAFF<<FFAFAFAAAFFFAFAAAAAFAAFAFAFAFA
@H3GFVCCXX150415:8:2224:8957:23407/1
CATACTTGATGGTCTCAGATATGTGTGGATTTTGAATT
+
<FAFAFFAAAA7AAFAAAFAAAAAFAFAFFFAAFFFF.
@H3GFVCCXX150415:8:2224:8907:25745/1
GTTAATTAAGCCCTTTACGAATGGACTAGATGTACCT
+
AAAFFAF/FFAAAA<FAA7FFAAFAFAF7FAAFAAF
@H3GFVCCXX150415:8:2224:8825:55175/1
GCACCTGTGTCAACAACCTGAGAGTGGCCTTGAGTTGC
+
AAFAAAFAAAFFFAFAFAFAAAFAAAFFFAFAFAA
@H3GFVCCXX150415:8:2224:7780:21245/1
```

(a) FASTA file stores reference genome

(b) FASTQ file stores sequenced reads

Data Source: National Center for Biotechnology Information

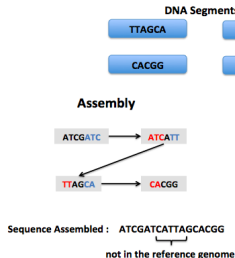
Mapping vs Assembly

Genomics BDSI

Sijia Huo,
Sean Kelly,
Gregory
Raskind

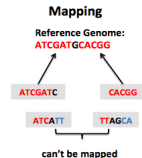
Assembly

- Reconstructs genome from sequenced reads
- Particularly useful to detect large genetic variants. (i.e. free from bias towards reference genome)
- Slow, memory-intensive, and hard to implement



Mapping

- Match to reference genome to locate variation
- Fast, easy to understand
- Cannot handle sequenced reads containing complex variation



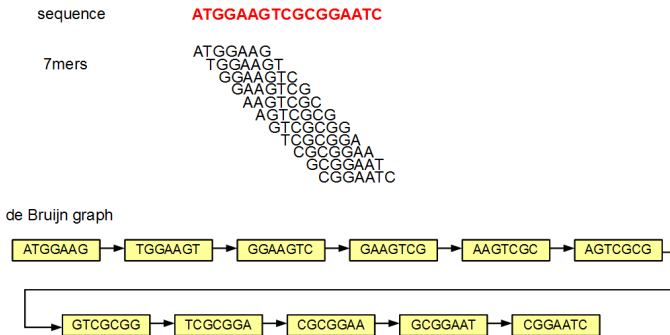
Constructing De Bruijn Graph of Reference Genome

Genomics
BDSI

Sijia Huo,
Sean Kelly,
Gregory
Raskind

Create process to store and reassemble a reference genome

Method: Break DNA string into 'kmers' and store in a De Bruijn Graph

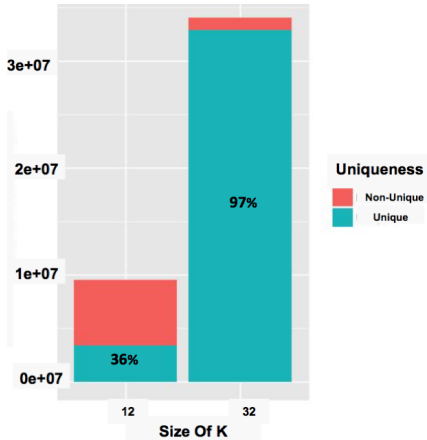


<http://www.homolog.us/Tutorials/Tut-Img/Set1/fig5.png>

Impact of K on complexity of De Bruijn Graph

Genomics
BDSI

Sijia Huo,
Sean Kelly,
Gregory
Raskind



Methodology: Functions

Genomics
BDSI

Sijia Huo,
Sean Kelly,
Gregory
Raskind

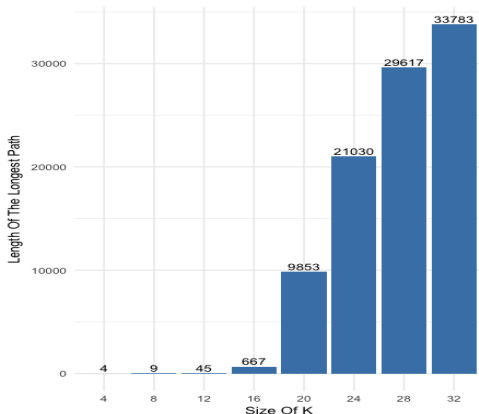
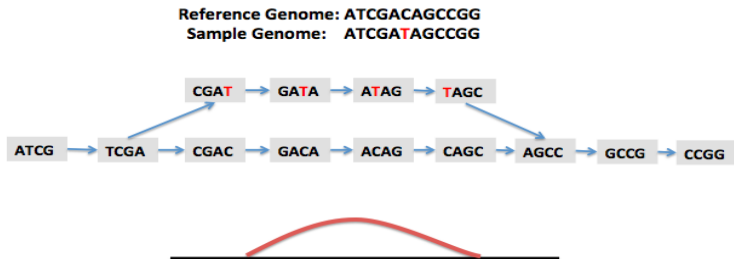


Figure: Comparison of The Lengths Of The Longest Paths

DNA Variation in a De Bruijn Graph

Genomics
BDSI

Sijia Huo,
Sean Kelly,
Gregory
Raskind



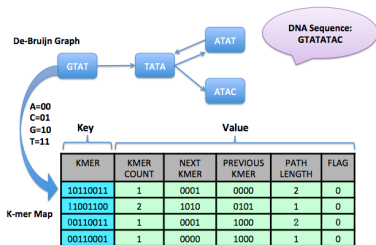
Given an individual's sequenced DNA, we want to locate where their genome varies from the reference

Method: Create De Bruijn graph for sequenced reads to store only variations in genome

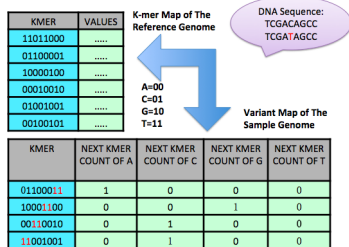
Storage Structure

Genomics
BDSI

Sijia Huo,
Sean Kelly,
Gregory Raskind



(a) Kmer Graph Structure



(b) Variant Graph Structure

Figure: Data Structures Of The Graphs

Implemented in C++

Future Plans

Genomics
BDSI

Sijia Huo,
Sean Kelly,
Gregory
Raskind

Cleaning/Pruning Variant Graph

- More accurate and usable
- Reduces storage size

FM indexing

- Allows compression of input text while still permitting fast substring queries

References

Genomics
BDSI

Sijia Huo,
Sean Kelly,
Gregory
Raskind

For Images:

- Iqbal, Zamin, Mario Caccamo, Isaac Turner, Paul Flicek, and McVean Gil. "De Novo Assembly and Genotyping of Variants Using Colored De Bruijn Graphs." Nature Genetics (2012): Nature Publishing Group. Web.

For information on sequencing and assembly methods

- Paolo Ferragina and Giovanni Manzini (2000). "Opportunistic Data Structures with Applications". Proceedings of the 41st Annual Symposium on Foundations of Computer Science. p.390.
- <http://genome.cshlp.org/content/18/5/821.long#sec-17>
- http://cortexassembler.sourceforge.net/index_cortex_var.html
- Umich Biostats lectures

Acknowledgments

Genomics
BDSI

Sijia Huo,
Sean Kelly,
Gregory
Raskind

University of Michigan and Bhramar Mukherjee
Hyun Min Kang, Goncalo Abecasis, Greg Zajac
NSF