

Rapid Indexing and Searching of Genomes

Eli Ben-Michael

Cynthia Chen

Biostatistics Big Data Summer Institute

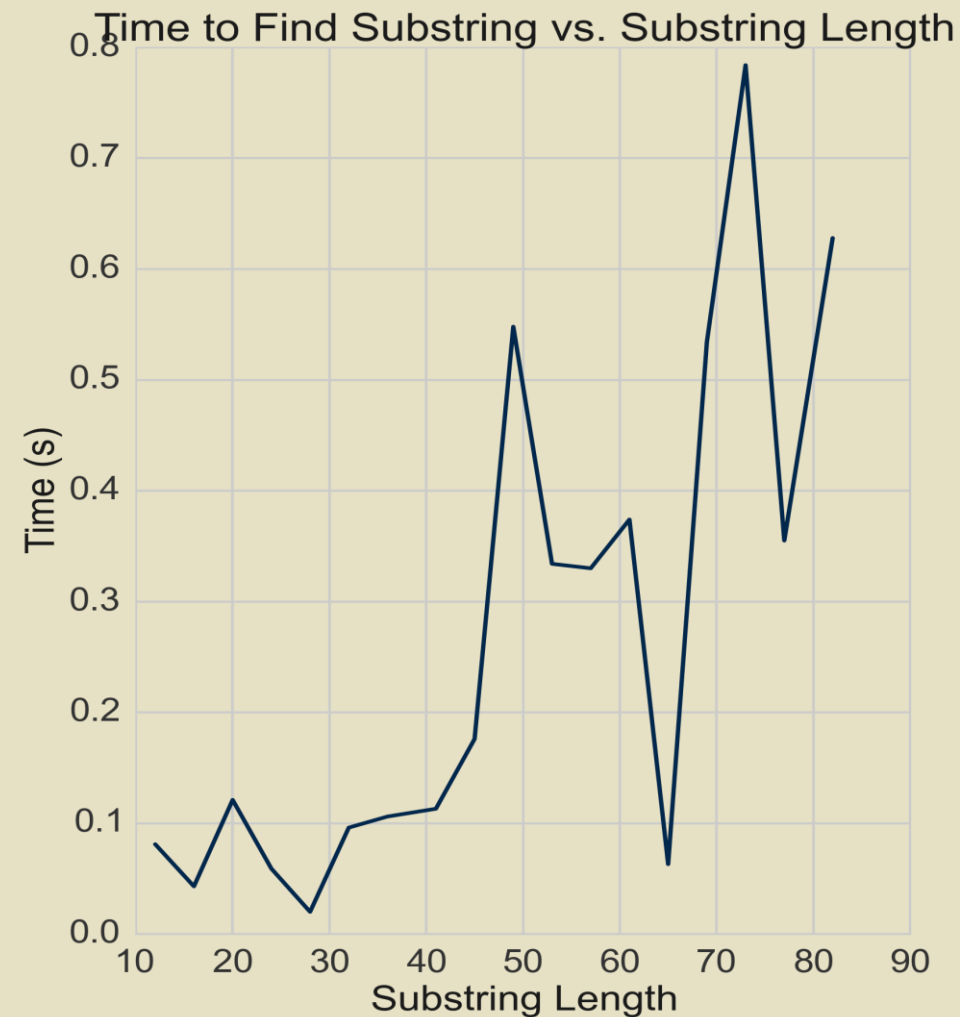
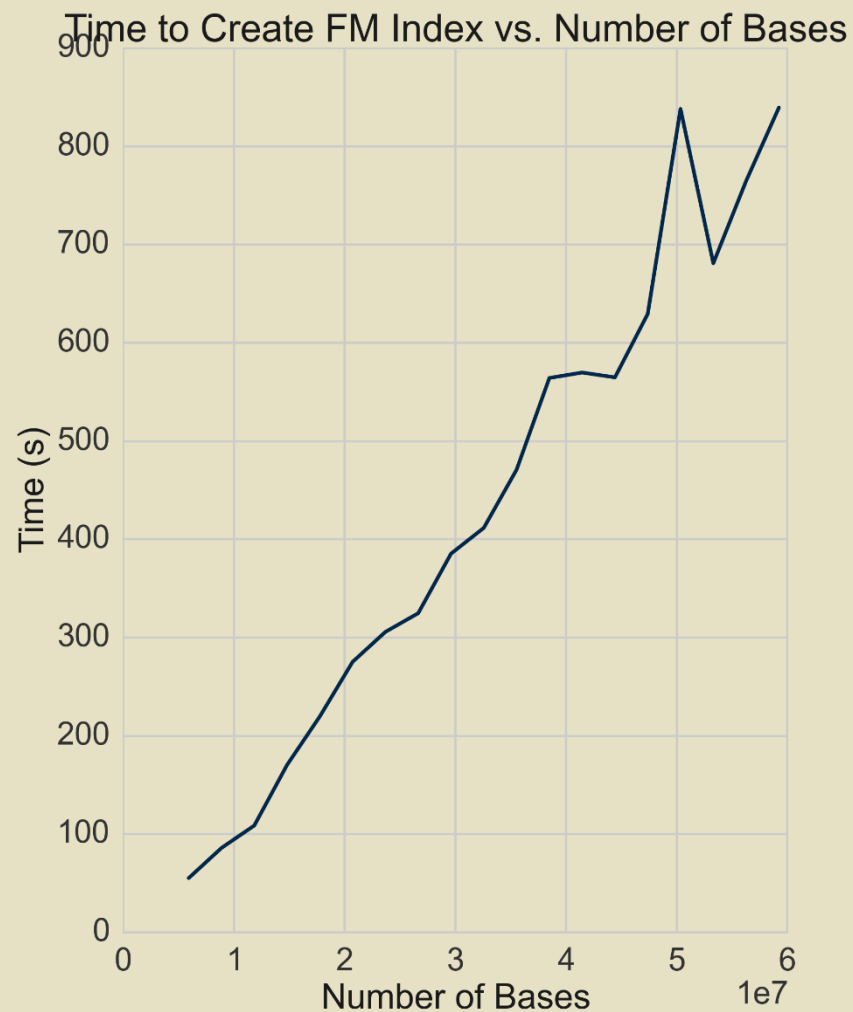
Symposium on Big Data, Human Health, and Statistics

Department of Biostatistics, University of Michigan

FM Indexes

- Genome is long
 - ~ 3 billion base pairs
 - How can we do efficient search on the genome?
 - How can we compress the genome?
- FM Indexing
 - Places characters with similar right contexts together
 - Efficient $O(1)$ search
 - Good for compression (Huffman Coding)

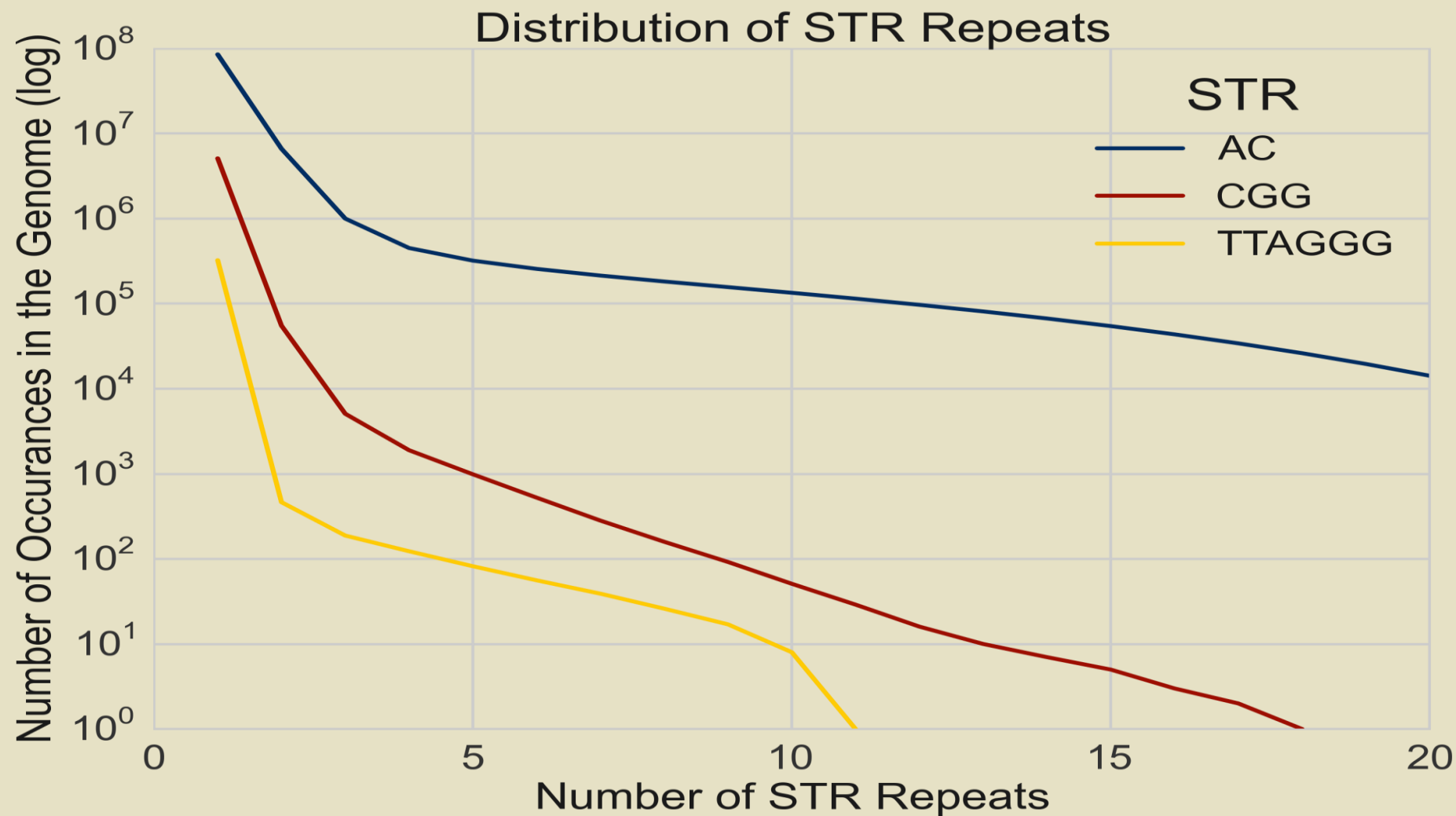
Results (1)



Short Tandem Repeats (STRs)

- STRs
 - E.g.
 - “AC”: Association with breast cancer
 - “CGG”: Repeat varies based on gender
 - “TTAGGG”: Telomere sequence
 - 2-13 nucleotides repeated consecutively on a strand of DNA
 - Important to understanding genetic variation in humans
 - Little is known about the frequencies

Results (2)



Acknowledgments

- Prof. Goncalo Abecasis
- Prof. Hyun Min Kang
- Sayantan Das
- Prof. Bhramar Mukherjee
- University of Michigan School of Public Health