

Scalable Genomic Assembly through Parallel de Bruijn Graph Construction for Multiple K-mers

Kanak Mahadik, Chris Wright, Milind Kulkarni,
Saurabh Bagchi, Somali Chaterji

Sequence assembly

- Reads : fragments of genome read by sequencing instruments
- Reconstruct the original genome from the reads
- *Approximate*, no single best assembly tool*
 - Repeating patterns in original genome
 - Uneven sampling of original genome
 - Errors in reads

* Salzberg, S. L., Phillippy, A. M., Zimin, A., Puiu, D., Magoc, T., Koren, S., & Marçais, G. (2012). GAGE: A critical evaluation of genome assemblies and assembly algorithms. *Genome research*, 22(3), 557-567.

De bruijn Graph (DBG)

SEQUENCE:

A A T G C C

READS:

A A T

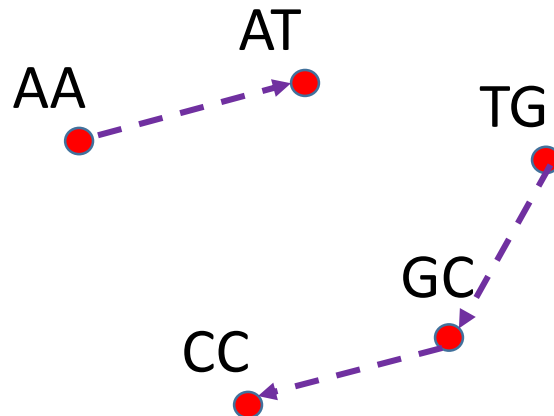
T G C

G C C

Object graph of DBG is undirected. Vertices are consecutive k -mers in a read & edges are $(k-1)$ nucleotides of the k -mer u are the last $(k-1)$ nucleotides of the k -mer v are the first $(k-1)$ nucleotides of the k -mer w . Constructed by traversing the DBG to identify maximal paths in the graph. (All vertices have in-degree and out-degree equal to 1, except at the start and end of the path)

K-mer set for $k=2$

{AA, AT, TG, GC, CC}



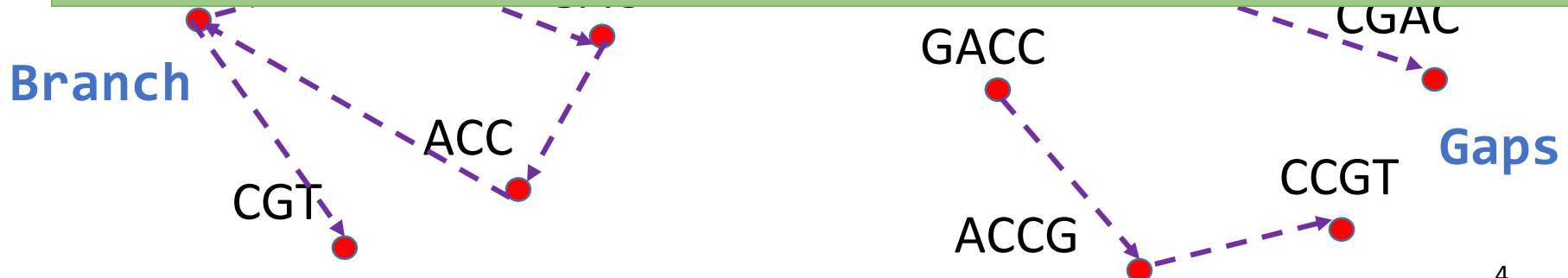
Contig set $k=2$
{AAT, TGCC}

Significance of K-value

SEQUENCE:

K-mer set for k=3

- Smaller k-value
 - Probability of extraction of valid k-mer from a read is higher
 - branched DBG
- Larger k-value
 - Can resolve repeats of greater length
 - fragmented DBG
- Both stop contig extension → reduced contig lengths



Iterative de bruijn graph tools

- Tools use multiple k-values sequentially : IDBA-UD, SPAdes
- Small k-value graph traversed to update larger graph
- Graph “**accumulated**” and “**updated**” at each iteration in a range
- IDBA-UD tool : medium-fast and memory efficient, SPAdes tool : slow and memory inefficient*

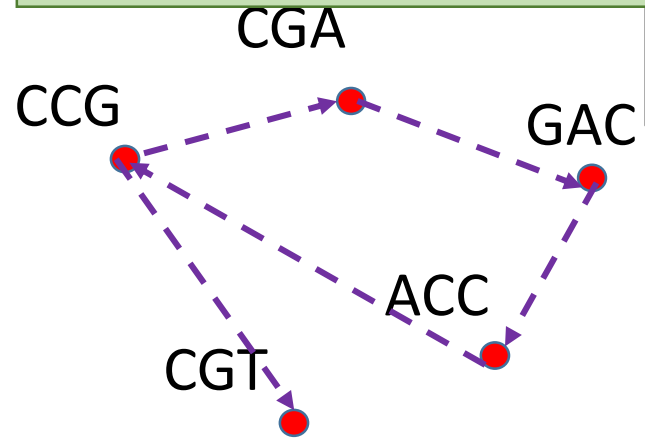
IDBA-UD iterative process

SEQUENCE:
C C G A C C G T
READS:
C C G A C **X**

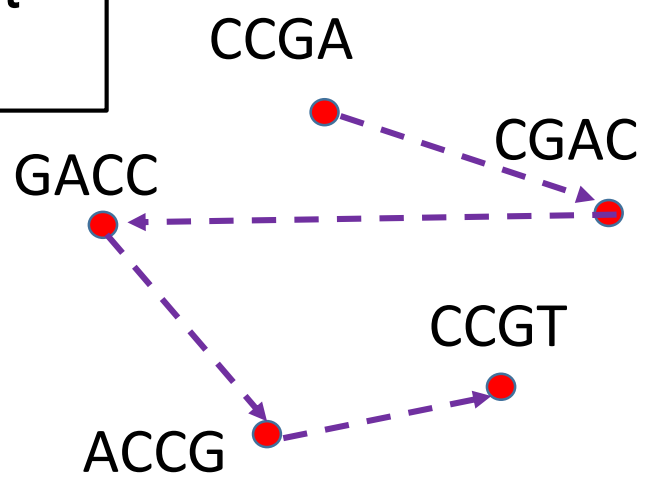
Graphs $G_{k=3}$ updated to represent reads
 are $G_{k=4}$ based on information from reads
 Set $k=4$ and Contig Set $C_{k=3}$

K-mer set for $k=3$
 {CCG, CGA, GAC, ACC, CCT}

- **Dependence** of graph at iteration "i+1" on graph at iteration "i", contigs of graph at iteration "i" and read set at iteration "i"
- Forces **sequential** operation on a chain of k-values



NEW READ SET
ACCGT



Problems with IDBA-UD

- Iterative graph construction iterations dominate the overall execution time(>90%)
- Iterative de bruijn assemblers are **sequential**
- **No parallelism** for **long chain of k-values**
- Scalability is limited to **single node**
- Insufficient for **large** datasets

ScalaDBG Insight

SEQUENCE:

C C G A C C G T

READS:

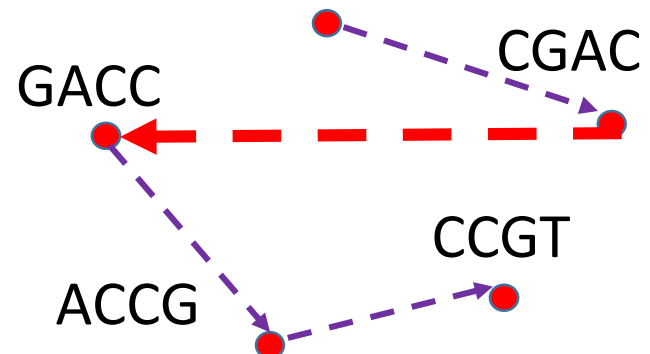
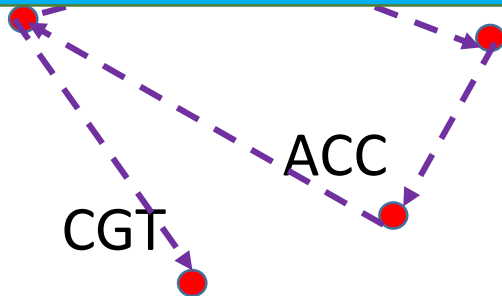
C C G A C

K-mer set for k=3

{CCG,CGA,GAC,ACC,CGT}

K-mer set for k=4

- Build (partial) graphs independently in parallel
- Push or “patch” information from lower k-value graph to higher k-value graph



ScalaDBG Parallel Patch

$$k1 < k2 < k3 < k4$$

$$G_{k1} = \text{Build}(k1)$$

$$G_{k2} = \text{Build}(k2)$$

$$G_{k3} = \text{Build}(k3)$$

$$G_{k4} = \text{Build}(k4)$$

- Number of serialized patching steps grow *logarithmically* with number of k-values
- Well known tree reduction parallel pattern

$$C_{k1-k2} = \text{Contigs}(G_{k1-k2})$$

$$G'_{k1-k4} = \text{Patch}(G_{k3-k4}, C_{k12})$$

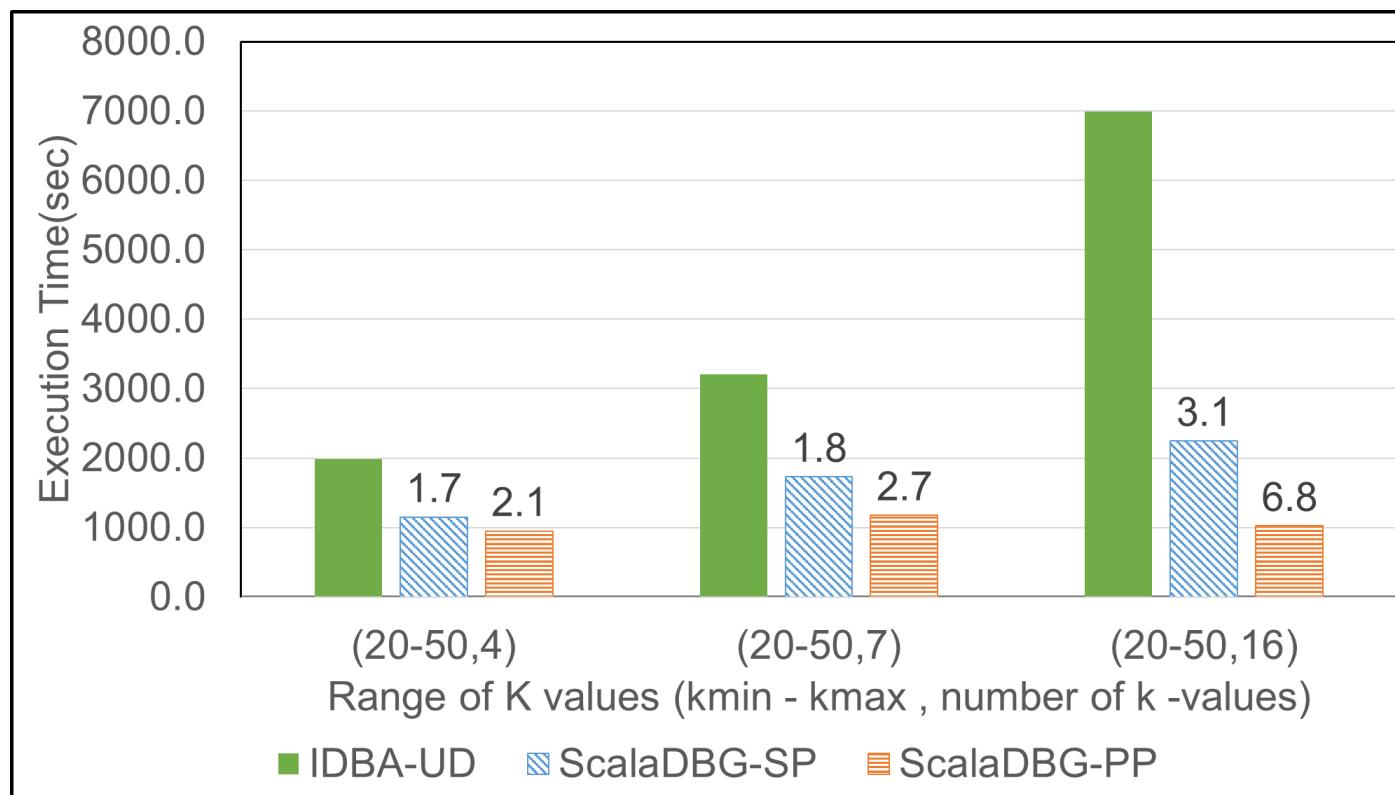
L2

Evaluation - Setup

Name	Read Set Type	Read Length (bp)	# reads
RM1 , RM2	Real, Metagenomic	150	33,140,480
SC – E.coli	Real, Single Cell	100	23,818,596
E. coli	Simulated	75	1,882,418
SC – S.aureus	Real, Single Cell	100	66,997,488
SC – SAR324	Real, Single Cell	100	55,733,218

- Intel Xeon Infiniband cluster, Intel Xeon E5-2670, 2.6 GHz with 16 cores per node, 32 GB of memory
- IDBA-UD (1.1.1)

Performance on SC-SAR 324 dataset



- Speedup increases with number of k-values
- 6.8X of baseline, 2.2X of serial patch for 16 k-values

Accuracy Results Overview

- Used Quast tool to analyze quality
- Difference in metric values not statistically significant
- Metrics Analyzed
 - N50 : Median value of the length of contigs
 - # Contigs : Number of assembled non-gap genomic sequences
 - Max Contig Length
 - NGA50 : Median value of the aligned contig length normalized w.r.t genome length
 - Coverage : Number of aligned bases in the contigs divided by the reference length.
 - #Missassemblies : number of positions in the contigs incorrect

Conclusion

- Technique to exploit parallelism in multi k-value DBG assembly
- Distribute on a cluster of nodes
- Provides modular stages for construction, patching, and contig generation, applicable to other assemblers

Thank You!