Graph analysis of fragmented long-read bacterial genome assemblies

Pierre Marijon, Rayan Chikhi, Jean-Stéphane Varré

Inria, University of Lille

Introduction

Introduction: de novo assembly problem, solved ?

Assembly of 3rd generation sequencing data

- requires correction (not my problem today)
- solves almost all genomic repetitions

Assembly graph of the *E. coli* genome¹:



¹One chromosome, one contig [Koren and Phillippy, 2015]

Introduction: de novo assembly problem, solved ?

Assembly of 3rd generation sequencing data

- requires correction (not my problem today)
- solves almost all genomic repetitions

Assembly graph of the *E. coli* genome¹:



But in reality ...

¹One chromosome, one contig [Koren and Phillippy, 2015]

 $\ensuremath{\text{NCTC}}$: 3000 bacteria cultures sequenced with PacBio, and assembled with $\ensuremath{\text{HGAP}}^2$

599 / 1136 (34 %) assemblies are not single-contig (as of Feb 2019)

Species	Strain	Sample	Runs	Automated Assembly	Manual Assembly	Manual Assembly Chromosome Contig Number	Manual Assembly Plasmid Contig Number	Manual Assembly Unidentified Contig Number
Achromobacter xylosoxidans	NCTC10807	ERS451415 C	ERR550491 C ERR550506 C ERR550507 C	Pending	EMBL 0	1	0	0
Budvicia aquatica	NCTC12282	ERS462988	ERR581162 C	Pending	EMBL 0	2	0	0
Campylobacter jejuni	NCTC11351	ERS445056 C	ERR550473 C* ERR550476 C*	Pending	EMBL 0	1	0	0
Cedecea neteri	NCTC12120	ERS462978 2	ERR581152 2 ERR581168 2 ERR597265 2	Pending	EMBL 0	7	1	0
Citrobacter amalonaticus	NCTC10805	ERS485850 2	ERR601566 C [*] ERR601575 C [*]	Pending	EMBL 0	1	2	0
Citrobacter freundii	NCTC9750 2	ERS485849 C	ERR601559 C ERR601565 C	Pending	EMBL 0	1	0	0
Citrobacter koseri	NCTC10849 C	ERS473430 C	ERR581173 C	Pending	EMBL 0	1	1	0
Corynebacterium diphtheriae	NCTC11397 2	ERS451417 C	ERR550510 C	Pending	EMBL 0	1	0	0
Cronobacter sakazakii	NCTC11467 2	ERS462977 2	ERR581151 C ERR581167 C	Pending	EMBL 0	4	3	0

²[Chin et al., 2013]

 $\ensuremath{\mathsf{NCTC}}$: 3000 bacteria cultures sequenced with PacBio, and assembled with $\ensuremath{\mathsf{HGAP}}^2$

599 / 1136 (34 %) assemblies are not single-contig (as of Feb 2019)

Species	Strain	Sample	Runs	Automated Assembly	Manual Assembly	Manual Assembly Chromosome Contig Number	Manual Assembly Plasmid Contig Number	Manual Assembly Unidentified Contig Number
Achromobacter xylosoxidans	NCTC10807	ERS451415 C	ERR550491 C ERR550506 C ERR550507 C	Pending	EMBL 0	1	0	0
Budvicia aquatica	NCTC12282	ERS462988 C	ERR581162 C	Pending	EMBL 0	2	0	0
Campylobacter jejuni	NCTC11351	ERS445056 C	ERR550473 C* ERR550476 C*	Pending	EMBL 0	1	0	0
Cedecea neteri	NCTC12120 2	ERS462978 🖙	ERR581152 C ERR581168 C ERR597265 C	Pending	EMBL 0	7	1	0
Citrobacter amalonaticus	NCTC10805	ERS485850 C	ERR601566 C* ERR601575 C*	Pending	EMBL 0	1	2	0
Citrobacter freundii	NCTC9750 2	ERS485849 C	ERR601559 C ERR601565 C	Pending	EMBL 0	1	0	0
Citrobacter koseri	NCTC10849 C	ERS473430 C	ERR581173 C	Pending	EMBL 0	1	1	0
Corynebacterium diphtheriae	NCTC11397 2	ERS451417 C	ERR550510 C	Pending	EMBL 0	1	0	0
Cronobacter sakazakii	NCTC11467 2	ERS462977 C	ERR581151 C ERR581167 C	Pending	EMBL 0	4	3	0

Assembly problem is solved for many bacteria but not for all.

²[Chin et al., 2013]

KNOT: Knowledge Network Overlap exTraction

KNOT: A synthetic example

- **Dataset**: *Terriglobus roseus* synthetic pacbio, 20x coverage (LongISLND³)
- Assembly tools: Canu⁴

tig 1





³[Lau et al., 2016] ⁴[Koren et al., 2017]

KNOT: A synthetic example

- **Dataset**: *Terriglobus roseus* synthetic pacbio, 20x coverage (LongISLND³)
- Assembly tools: Canu⁴

tig 1 tig 4 tig 8

Can we recover missing edges between contigs?

³[Lau et al., 2016] ⁴[Koren et al., 2017]

Not even a repetition problem..

Dotplot of T. roseus genome against itself.



Length of the tandem repeat is 460 kbp. The repetition explains only one of the two contig breaks.

An assembly graph can be defined as :

- nodes \rightarrow reads
- edges \rightarrow overlaps

KNOT: A synthetic example

An assembly graph can be defined as :

- nodes \rightarrow reads
- edges \rightarrow overlaps

Overlap graph (constructed by Minimap2 ⁵), reads are colored by Canu contig.

⁵[Li, 2018]

KNOT: A synthetic example

An assembly graph can be defined as :



Overlap graph (constructed by Minimap2 ⁵), reads are colored by Canu contig.

⁵[Li, 2018]



KNOT: definition of an Augmented Assembly Graph

The AAG is an undirected, weighted graph:

nodes: contigs extremities

edges:

- between extremities of a contig (weight = 0),
- paths found between contigs (weight = path length in bases)

KNOT: definition of an Augmented Assembly Graph

The AAG is an undirected, weighted graph:

nodes: contigs extremities

edges:

- between extremities of a contig (weight = 0),

- paths found between contigs (weight = path length in bases)



Plain links are paths compatible with true order of contigs, dotted links are other paths.

We classify paths based on their length (in base pairs):



In prokaryotes, most repetitions are < 10 kbp $^{\rm 6}$

⁶[Treangen et al., 2009]

KNOT: Hamilton walk

AAG's are generally complete graphs. We can enumerate all their Hamilton walks.

The weight of a walk is the of sum of all edge weights.

KNOT: Hamilton walk

AAG's are generally complete graphs. We can enumerate all their Hamilton walks.

The weight of a walk is the of sum of all edge weights.

Supposedly: We assume that **lowest-weight walk** is the true genome.



We selected 38 datasets from NCTC3000, where **Canu**, **Miniasm** and **Hinge** didn't produce the expected number of chromosomes (*i.e. unsolved assemblies*).

- 19 datasets were manually solved by NCTC
- 17 remained fragmented
- 2 with no assembly attempt by NCTC

Across 38 datasets:

Mean number of	
Canu contigs	4.32
Edges in AAG	32.67
Theoretical max. edges in AAG	41.83
Distant edges	28.64
Adjacency edges	4.02
Dead-ends in Canu contigs	4.94
Dead-ends in AAG, adjacency edges	2.70

Across 38 datasets:

Mean number of	
Canu contigs	4.32
Edges in AAG	32.67
Theoretical max. edges in AAG	41.83
Distant edges	28.64
Adjacency edges	4.02
Dead-ends in Canu contigs	4.94
Dead-ends in AAG, adjacency edges	2.70

Almost half of the missing paths in contigs graph are recovered.

KNOT: Hamilton walk



KNOT: Hamilton walk



Generally, the true contig ordering is a low-weight Hamiltonian walk

Graph analysis of fragmented long-read bacterial genome assemblies

Summary:

- Bacterial assembly is not solved for all datasets
- Build and analyse Augmented Assembly Graph can help

Future:

- Assembly graph between contig
- Biological validation (we search collaboration)
- Application to larger genome/metagenome
- Performance improvement (path search step)



https://gitlab.inria.fr/pmarijon/knot



Teng tensors.							
Barelbooks gath							
				- Constanting -			
1000 100 1000 1							
			1. (2)	14 M			

13

References i

 Chin, C.-S., Alexander, D. H., Marks, P., Klammer, A. A., Drake, J., Heiner, C., Clum, A., Copeland, A., Huddleston, J., Eichler, E. E., Turner, S. W., and Korlach, J. (2013).
 Nonhybrid, finished microbial genome assemblies from long-read SMRT sequencing data. Nature Methods, 10(6):563–569.

 Koren, S. and Phillippy, A. M. (2015).
 One chromosome, one contig: complete microbial genomes from long-read sequencing and assembly.
 Current Opinion in Microbiology, 23:110–120.

References ii

 Koren, S., Walenz, B. P., Berlin, K., Miller, J. R., Bergman, N. H., and Phillippy, A. M. (2017).
 Canu: scalable and accurate long-read assembly via adaptive

k-mer weighting and repeat separation.

Genome Research, 27(5):722–736.

- Lau, B., Mohiyuddin, M., Mu, J. C., Fang, L. T., Asadi, N. B., Dallett, C., and Lam, H. Y. K. (2016).
 LongISLND:in silicosequencing of lengthy and noisy datatypes. Bioinformatics, 32(24):3829–3832.
- Li, H. (2018).

Minimap2: pairwise alignment for nucleotide sequences. Bioinformatics, 34(18):3094–3100.



Treangen, T. J., Abraham, A.-L., Touchon, M., and Rocha, E. P. (2009).

Genesis, effects and fates of repeats in prokaryotic genomes. *FEMS Microbiology Reviews*, 33(3):539–571.