

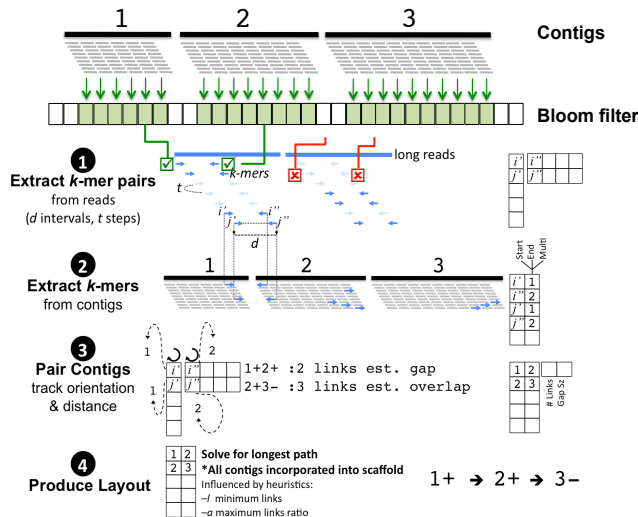
Scalable Algorithms for Long Read Assembly and Scaffolding

René L. Warren,
Chen Yang, Benjamin P. Vandervalk,
Bahar Behsaz, Albert Lagman Steven JM Jones,
Inanç Bırol

Abstract

Routine reconstruction of complex genomes from experimental data is not a solved problem, owing to long repeats that are not resolvable by short reads. Established and emerging long read technologies hold the potential to address present limitations, but their current high errors typically require base correction and/or additional pre-processing before use. Here we present LINKS, a method that exploits the sequence properties of long reads for scaffolding high-quality genome drafts.

Algorithm



References

- Beitzel, K. et al. (2014) Assembling Large Genomes with Single-Molecule Sequencing and Locality Sensitive Hashing. *bioRxiv* doi: <http://dx.doi.org/10.1101/008001>
- Bueter, M. and Provano, W. (2014) SSPACE-LongRead: scaffolding bacterial draft genomes using long read sequence information. *BMC Bioinformatics*, 15, 211.
- Clarke, J. et al. (2009) Continuous base identification for single-molecule nanopore DNA sequencing. *Nat Nanotechnol*, 4, 266-270.
- Goodwin, S. et al. (2015) Oxford Nanopore Sequencing and de novo Assembly of a Eukaryotic Genome. doi: <http://dx.doi.org/10.1101/013490>.
- Genovese, A. et al. (2013) QUAST: quality assessment tool for genome assemblies. *Bioinformatics*, 29, 1072-4.
- Koren, S. and Philippakis, A. M. (2014) One chromosome, one contig: complete microbial genomes from long-read sequencing and assembly. *Curr Opin Microbiol*, 23C, 110-120.
- Quick, J. et al. (2014) A reference bacterial genome dataset generated on the MinION™ portable single-molecule nanopore sequencer. *Genomics*, 3, 22.
- Simpson, J. L. et al. (2009) ABySS: a parallel assembler for short-read sequence data. *Genome Res*, 19, 117-25.
- Warren, R. L. et al. (2007) Assembling millions of short DNA sequences using SSAKE. *Bioinformatics*, 23, 500-1.

Funding

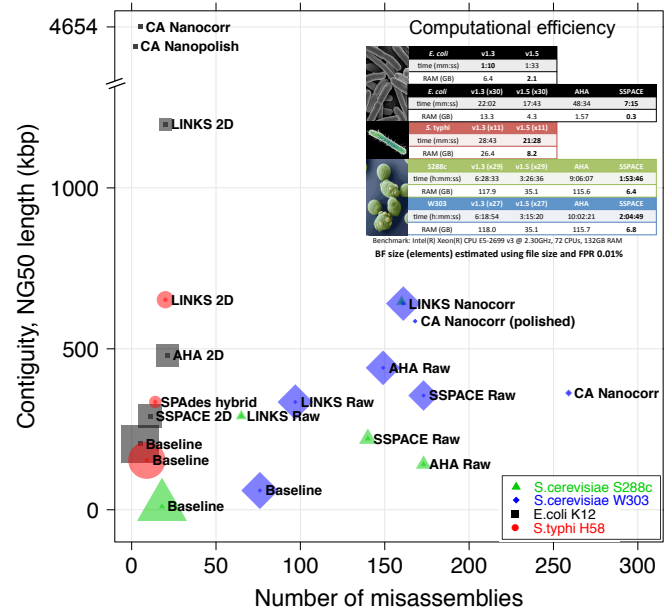


Error Model

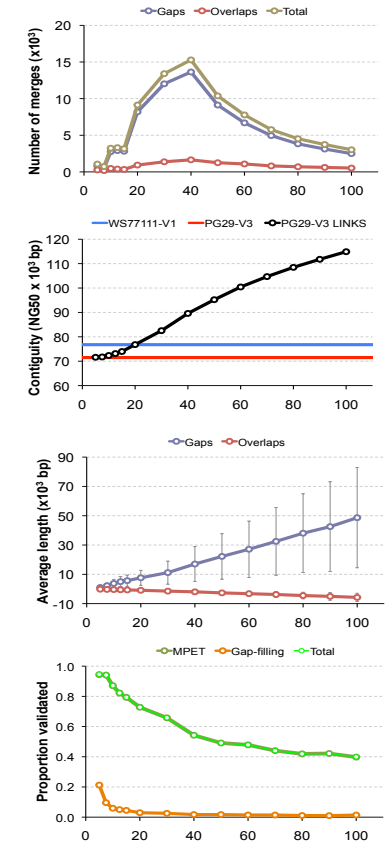
Mismatch: $P_m \sim a_m$ Poisson(l_m) + (1- a_m) Geometric(p)
 Insertion: $P_i \sim \alpha_i$ Weibull(l_i, κ_i) + (1- α_i) Geometric(p_i)
 Deletion: $P_d \sim \alpha_d$ Weibull(l_d, κ_d) + (1- α_d) Geometric(p_d)

	Mismatch			Insertion			Deletion				
	a_m	l_m	p_m	a_i	l_i	κ_i	p_i	a_d	l_d	κ_d	p_d
<i>E. coli</i> R7	0.248	0.480	0.711	0.850	1.004	0.968	0.418	0.870	0.986	1.026	0.403
<i>E. coli</i> R7.3	0.138	0.441	0.476	0.900	1.045	1.473	0.272	0.959	1.059	1.682	0.249
<i>S. cerevisiae</i> R7	0.177	0.499	0.479	0.961	1.024	1.613	0.194	0.891	1.066	1.814	0.207

Performance



Scalability



Iterative scaffolding of the 20 Gbp *Picea glauca* (white spruce) assembly using draft assemblies of two genotypes

Genotype 1 (GCA_000411955.3) 4.2M scaffolds
 Genotype 2 (PRJNA242552) 4.3M scaffolds

- Found 84,529 total merges
- Validated final LINKS assembly with MPET reads and a gap-filling tool (Sealer, Paulino *et al.*, in review)