

Genome Assembly Toolkit

LINKS scaffold graph (*E. coli* K12)

René L Warren

ISMB, July 2016



2015

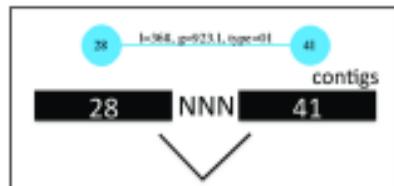
RESEARCH

Open Access

LINKS: Scalable, alignment-free scaffolding of draft genomes with long reads



René L. Warren¹, Chen Yang, Benjamin P. Vondervisck, Bahar Behsaz, Albert Lagman, Steven J. M. Jones and Inanc Birol



Vondervisck et al. *BMC Medical Genomics* 2015, **8**:Suppl 10:S1
<http://www.biomedcentral.com/1755-8289/8/S1/S1>



RESEARCH

Open Access

Konnecter v2.0: pseudo-long reads from paired-end sequencing data

Benjamin P. Vondervisck, Chen Yang, Zhuyi Xue, Karthika Raghavan, Justin Ong, Hamid Mohammadi, Shaun D. Jackman, Readman Chiu, René L. Warren, Inanc Birol¹

From IEEE International Conference on Bioinformatics and Biomedicine (BIBM 2014)
Belfast, UK, 2-5 November 2014

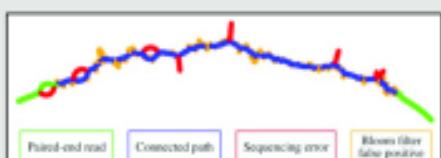


Figure 1 A connecting path between two non-overlapping paired-end sequencing reads within a de Bruijn graph. Konnecter

Paulino et al. *BMC Bioinformatics* (2015) 16:200
DOI 10.1186/s12859-015-0663-4



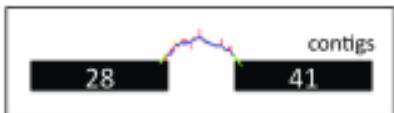
SOFTWARE

Open Access

Sealer: a scalable gap-closing application for finishing draft genomes

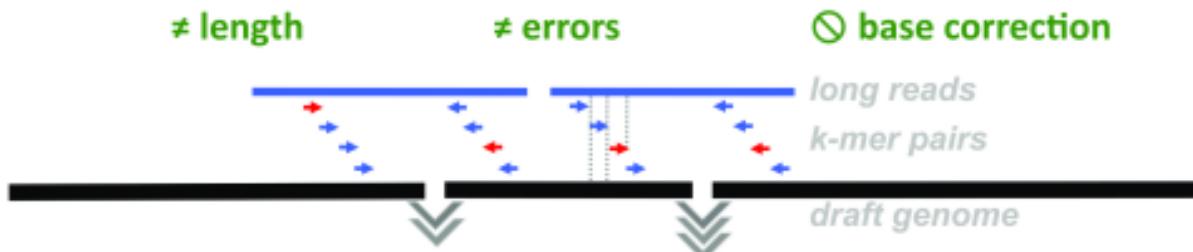


Daniel Paulino¹, René L. Warren¹, Benjamin P. Vondervisck², Anthony Raymond¹, Shaun D. Jackman¹ and Inanc Birol^{1,2*}





- **Scaffolder** : order & orient sequences
- ***k-mer* based** : no alignments
- **Vast *k-mer* space** : no fragment length limitation
- **Versatile** : long-reads, draft sequences, MPET



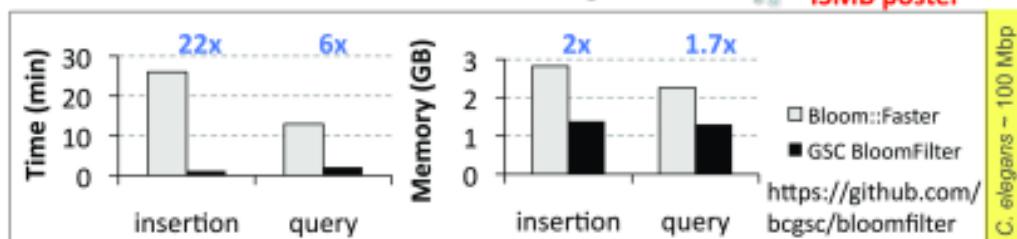
New Features

v1.6 : Custom Bloom filter

- ⌚ *ntHash* recursive *k*-mer hashing : fast



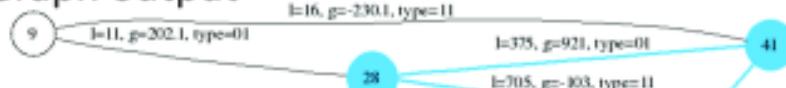
Mohamadi et al.
ISMB poster



- ⌚ Build large genomes filters

Species	Genome (Gb)	Time	RAM (GB)
human	3	: 40m	26
bullfrog	6	: 2h	37
spruce	20	: 5h	178

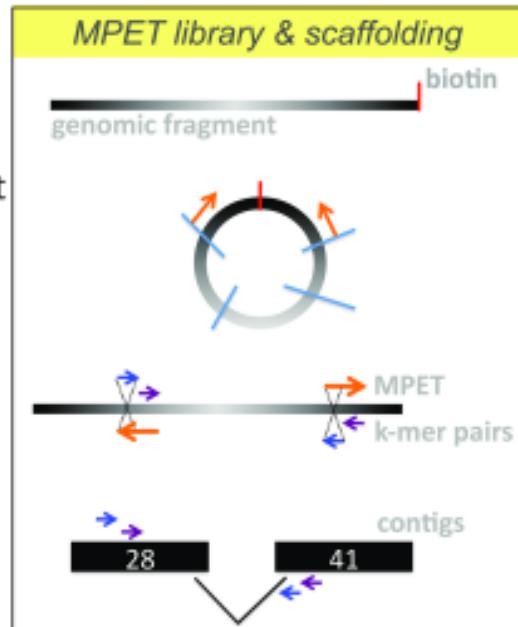
- ⌚ Graph output



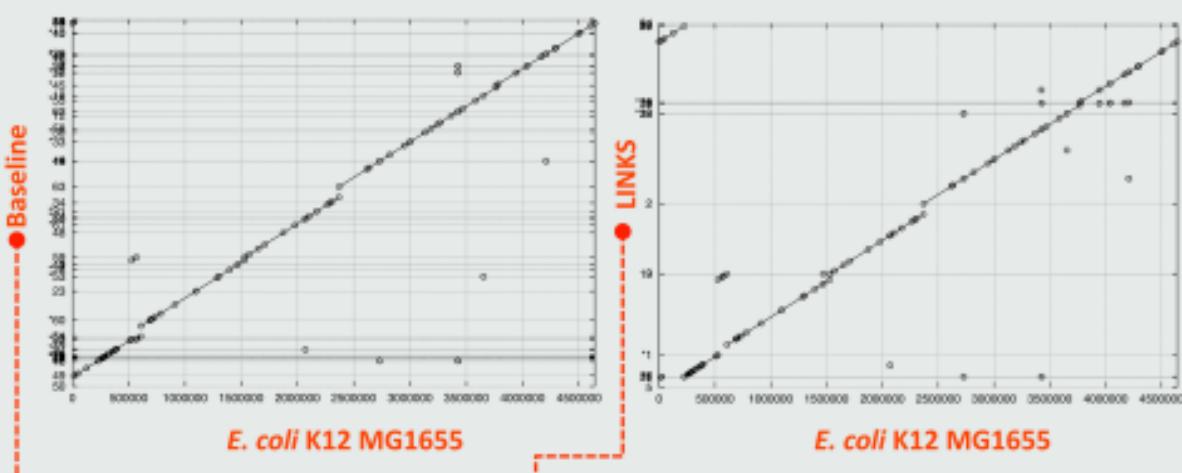
New Features

v1.7 : MPET Scaffolding

- Long-range pairing
- Massive depth, linkage support
- [Still] no alignments
- ↗ Scaffolding checkpoints
- ↗ Support for zip/gzip



MPET Scaffolding



<i>E. coli</i> K12	n:500	NGA50
Baseline contig	61	180 kbp
LINKS v1.7	27	300 kbp

Trimmed 4kb MPET: 260-fold physical coverage

<i>H. sapiens</i>	NA19238 NGA50 *	GIAB HG004 NGA50 ^
Baseline contig	56 kbp	48 kbp
LINKS v1.7	473 kbp	1.5 Mbp

*Trimmed 2kb MPET: 150-fold physical coverage

^Trimmed 6kb MPET: 175-fold physical coverage

10h @ 400 GB RAM

New Features

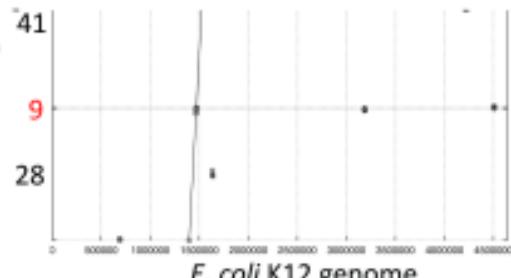
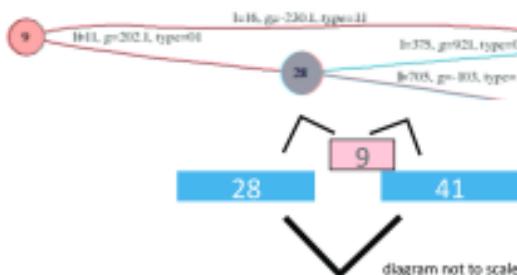
v1.8 : Native k-mer length *intervals*

↪ Single-run, Δi

<i>E. coli</i> K12 MinION full2D ONT	paper (v1.5) 30 iterations	v1.8 1 iteration 15 intervals*	Fold change
NG50 length (bp)	633,147	1,152,663	1.8
Memory (GB)	4.3	4.3	-
Time (mm:ss)	17:21	02:54	5.5

*1-10:1kb, 11-20:2kb

↪ Prioritizes proximal contigs





LINKS

BBT



KOLLECTOR



Konnector
&

RAILS

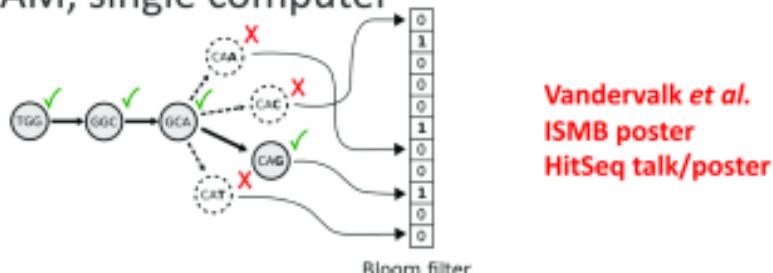
Sealer

ABySS

de novo assembly of large genomes

Human Genome First:

- ∅ Then: parallel, short read *k*-mer assembler
 - MPI to aggregate memory (Simpson et al, 2009)
- ∅ Now: de Bruijn graph Bloom filter representation
 - 1/10th RAM, single computer



BBT

Bio Bloom Tools

Sequence classification with Bloom filters

general purpose filtering

contaminant screening

pathogen discovery

BioBloom-

Maker

- Build filters : Re-usable
- Customizable : Flexible
- Multi-filter : Concurrent



Categorizer

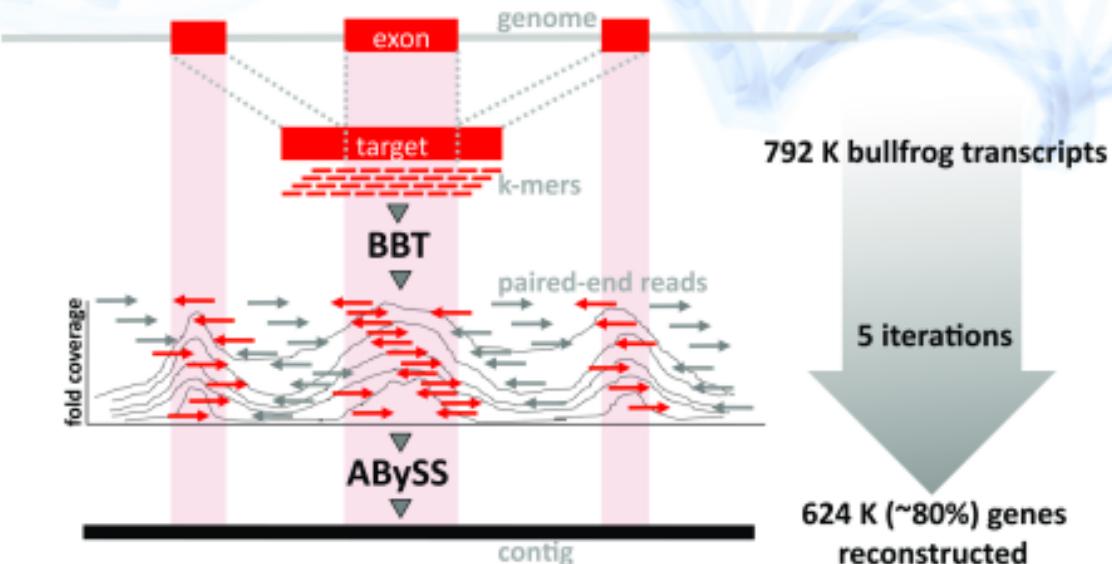
- Bins sequences : Analysis summary



Chu et al. 2014
HitSeq
poster/talk

KOLLECTOR

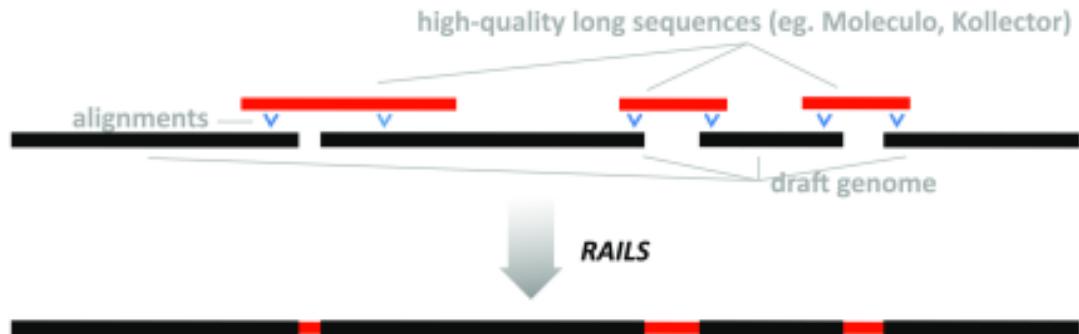
Targeted *de novo* assembly of gene loci from RNA



RAILS

Radial Assembly Improvement by Long Sequence Scaffolding

- ♂ Scaffolding and gap-filling
- ♂ Uses LINKS algorithm



Sealer

Automated genome finishing

⦿ Gap-filler (resolve Ns) **Application of Konnector**

⦿ Bloom filter de Bruijn graph (Scalable)

Vandervalk et al, 2015
Paulino et al, 2015

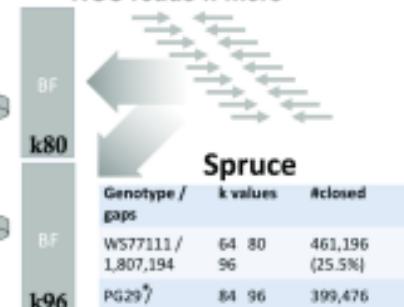


1 Extract flanking sequences

2 Assemble Konnector k-sweep

3 Generate new draft

Build Bloom filters
(Konnector)
NGS reads k-mers

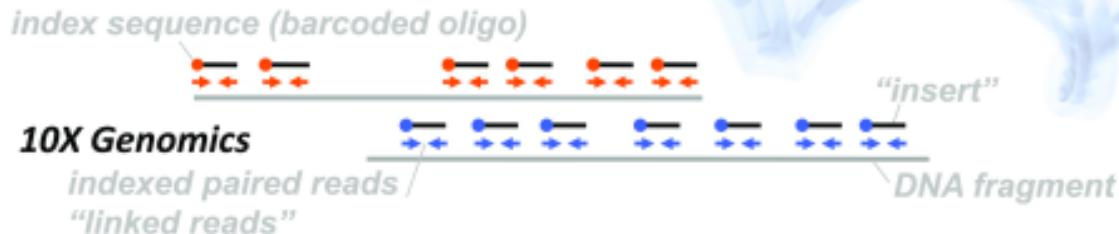


*4.5B Illumina MiSeq/HiSeq2000 reads

Peak memory: 44 GB RAM

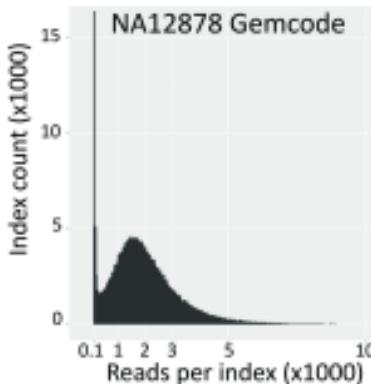
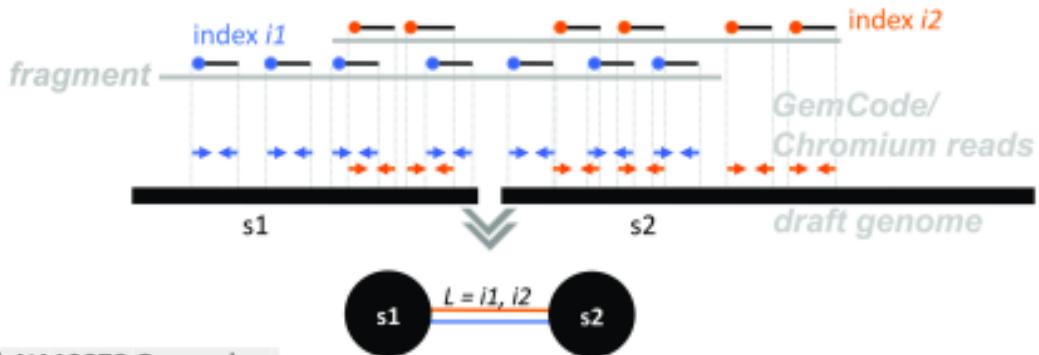


Assembly Roundup by *Chromium Scaffolding* *gemCode*



- ⌚ Capture long-range info
- ⌚ Identify co-located scaffolds in genome draft
- ⌚ Improve *de novo* assembly

Method

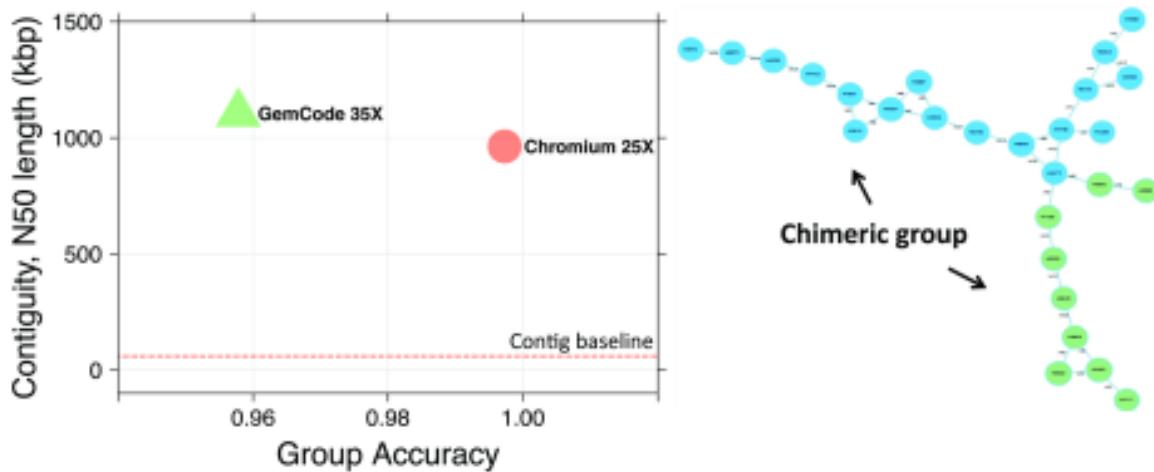


🔍 Pair sequences, Build graph

- Allow indexed reads within frequency range
- Link 2 sequences \geq reads map per index
- Create edge with $\geq L$ links

🔗 Connected components = groups

Grouping, human NA12878



⌚ Accuracy of groups assessed with BBT

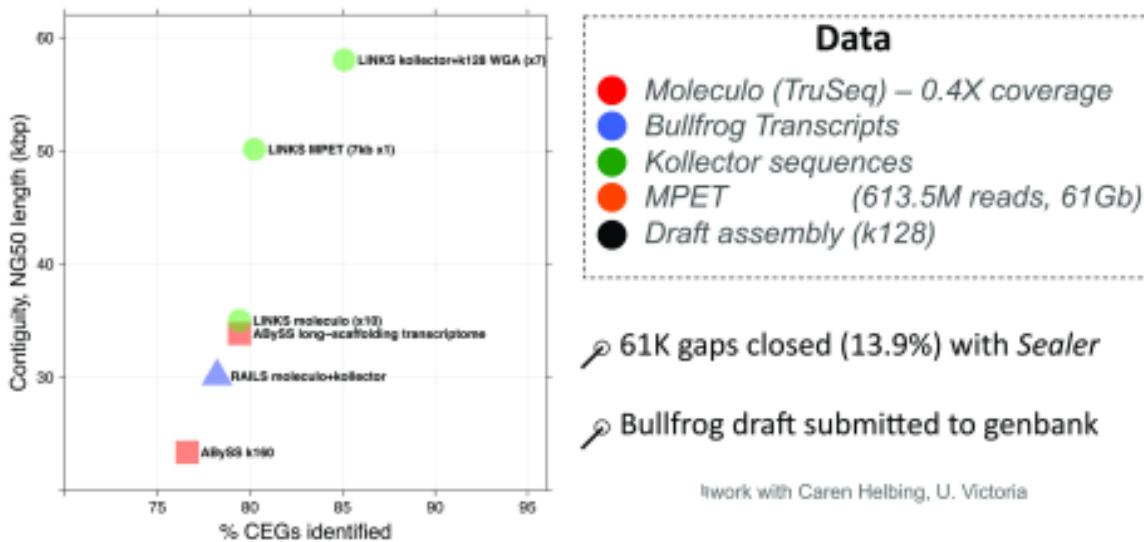
Application to large genomes



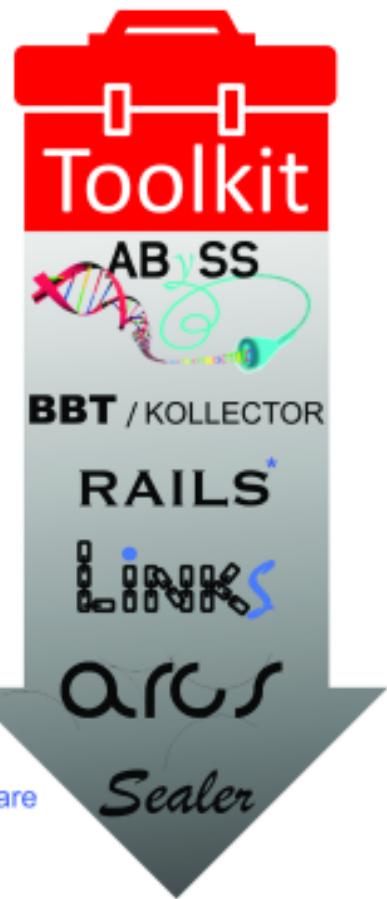
Image source: https://en.wikipedia.org/wiki/American_bullfrog

Scaffolding the bullfrog genome[¶]

Assemblies	ABySS k160	RAILS	Long-Scaffolding	LINKS
#Scaffold merges	NA	36,137	NA	177,506
NG50 length (bp)	23,361	30,085	33,847	58,021
Core Eukaryotic Genes	77%	78%	79%	85%

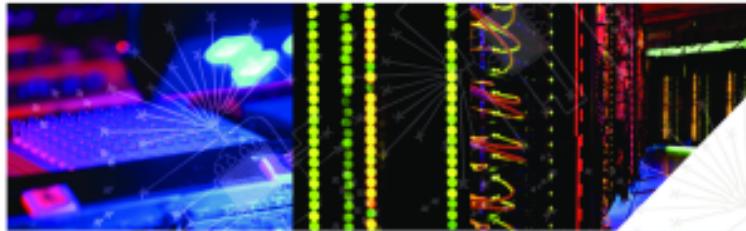


Genome Assembly



<http://www.bcgsc.ca/platform/bioinfo/software>

*<ftp://ftp.bcgsc.ca/supplementary/RAILS>



CANADA'S MICHAEL SMITH
GENOME
SCIENCES
CENTRE

A leading international centre for genomics and bioinformatics research committed to advancing knowledge of cancer-related diseases, improving human health through disease prevention, diagnosis and therapeutic approaches, and realizing social and economic benefits of genomics research.

54 TERABASES SEQUENCED • A HUMAN GENOME EVERY 17 MINUTES • HIGH-PERFORMANCE COMPUTING



AFFILIATIONS BC Cancer Research Center • BC Cancer Agency • BC Cancer Foundation • Genome BC • Simon Fraser University • University of British Columbia • Genome Sciences Institute

Thank You!

Austin Hammond | Ben Vandervalk | Chen Yang | Daniel Paulino | Erdi Kucuk Golnaz Jahesh | Hamid Mohamadi | Inanc Birol | Justin Chu | Lauren Coombe René Warren | Sarah Yeo | Shaun Jackman | Tony Raymond



GenomeBritishColumbia

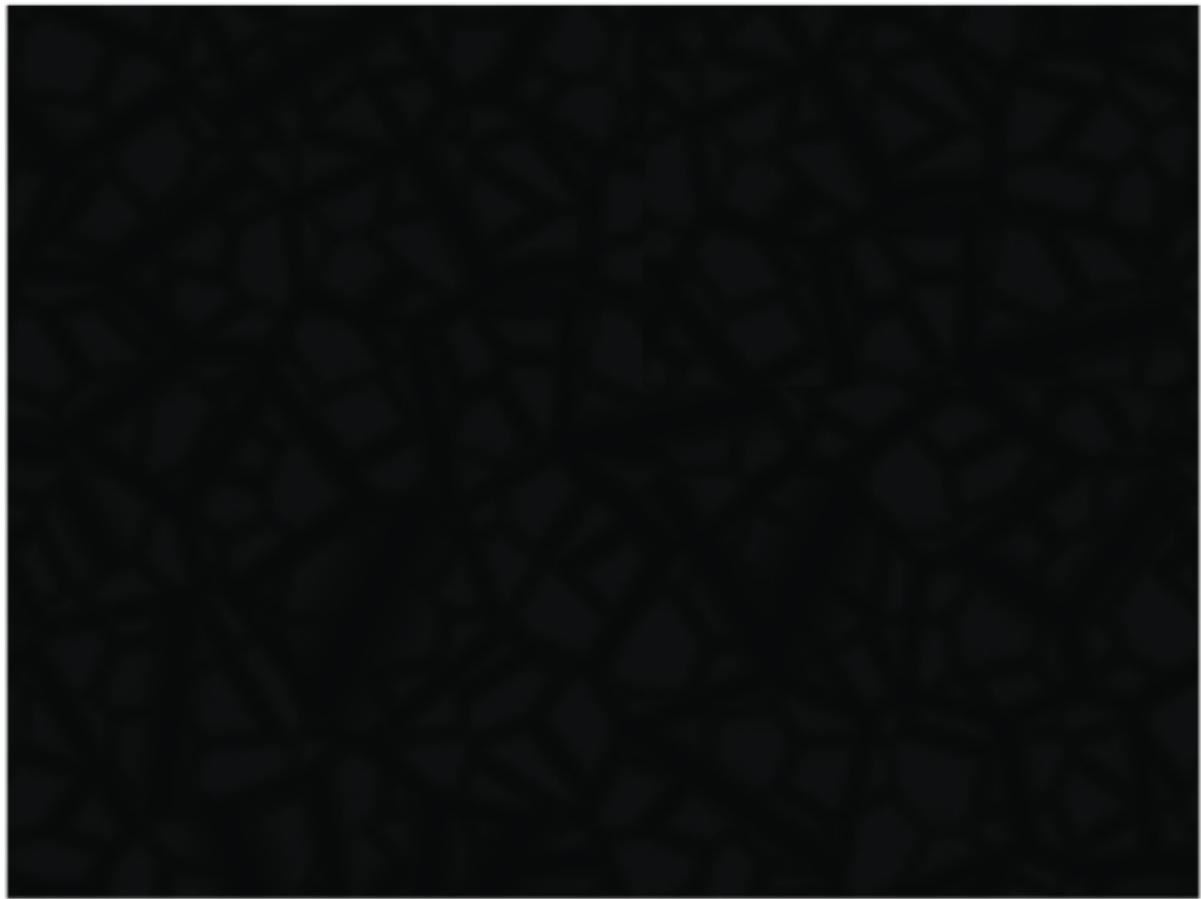


GenomeCanada

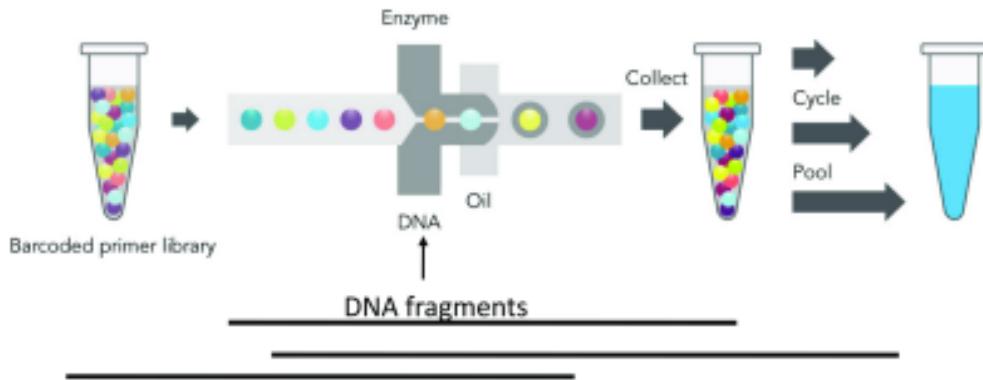


John Jambor Knowledge Fund

BCCA CANCER RESEARCH CENTRE



Overview of 10X Genomics' GemCode Platform



Summary, *E. coli*

