

# QRC

## *Scaffolding with Linked Reads*

**René Warren**

RECOMB-Seq, May 2017



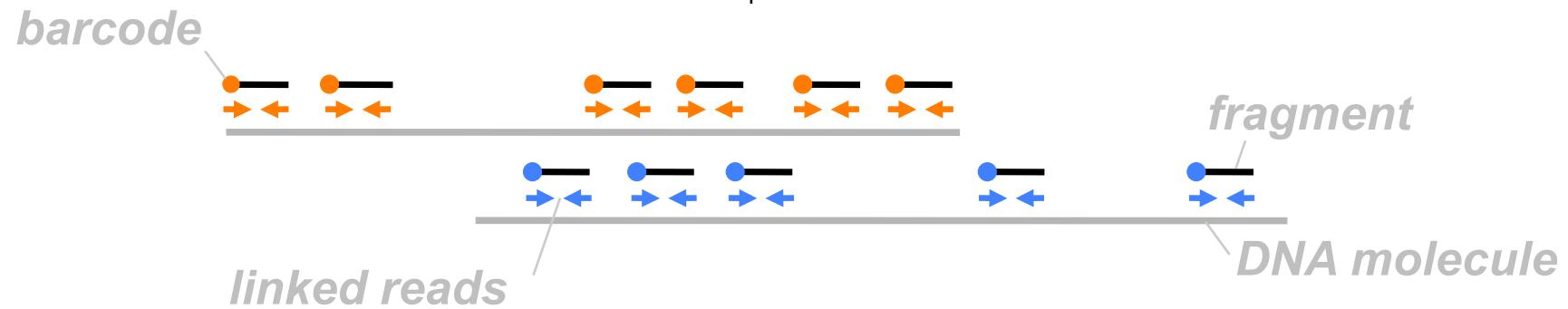
CANADA'S MICHAEL SMITH  
GENOME  
SCIENCES  
CENTRE  
[www.bcgsc.ca](http://www.bcgsc.ca)

# Genome scaffolding

- **Order and orient contigs**
  - Separated by gaps or overlaps
  - Goal : one scaffold, one chromosome
  - Recover complete genes
- **10X Genomics' Chromium : co-located linked reads**
  - DNA molecules : ↗ physical, ↘ individual sequence coverage
- **Bioinformatics Technologies**
  - Applicable to linked reads:
    - *fragScaff* [HiC data](#) (Adey 2014, Mostovoy 2016) > Alignment-based
    - *Architect* [Moleculo read cloud](#) (Kuleshov 2016) >
    - *Supernova* [Chromium](#) (Weisenfeld 2016) *de novo* assembler

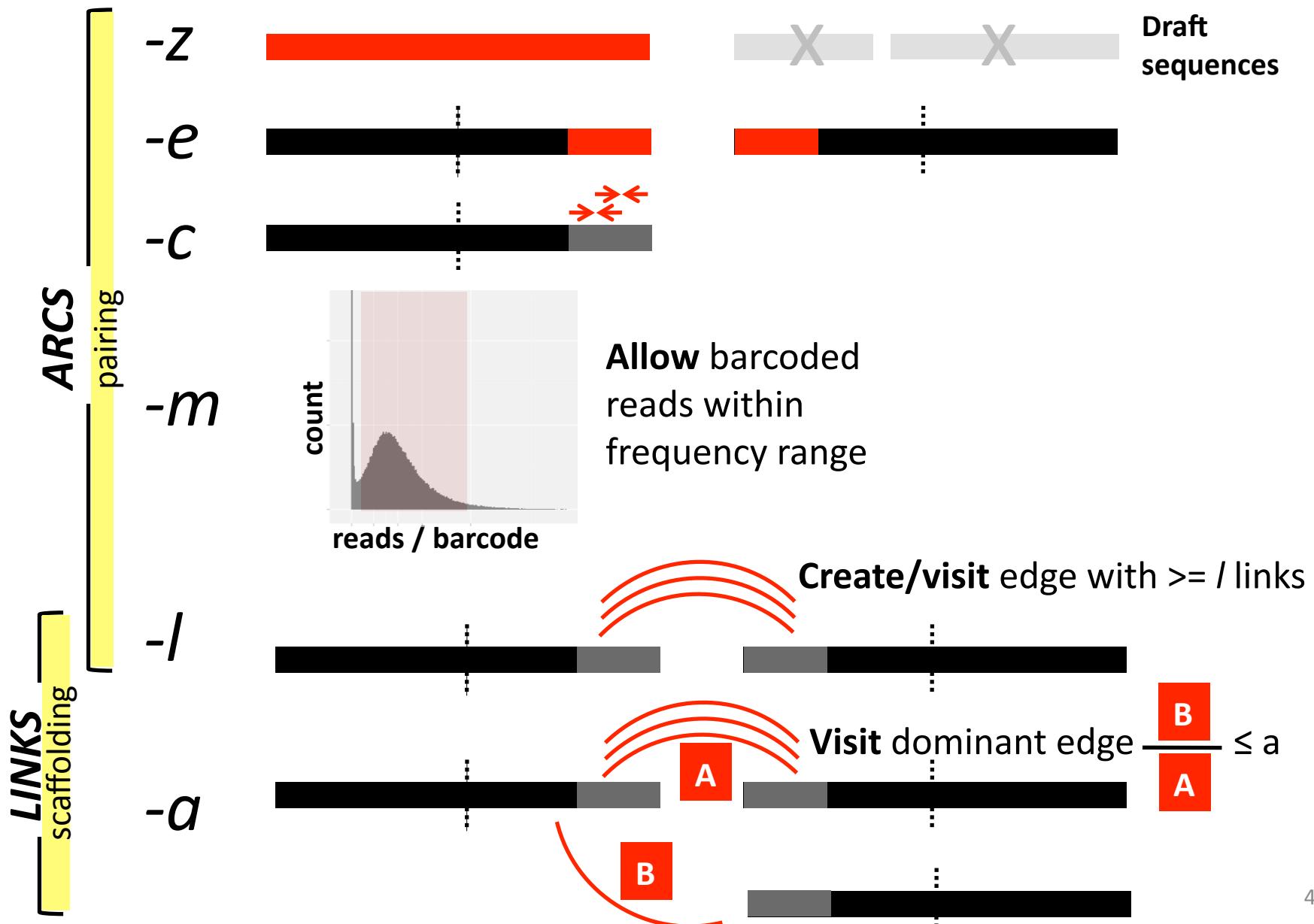
# arcs

## Assembly Roundup by Chromium Scaffolding



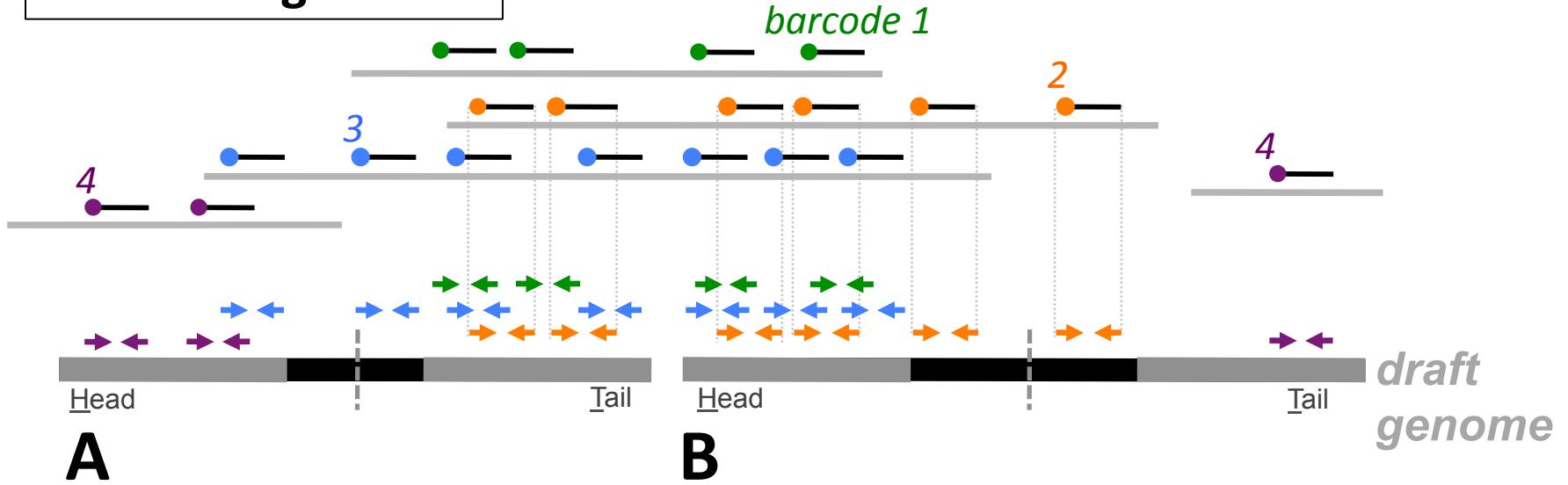
- Captures long-range info
- Identifies co-located scaffolds in genome draft
- Improves *de novo* assembly

# Parameters



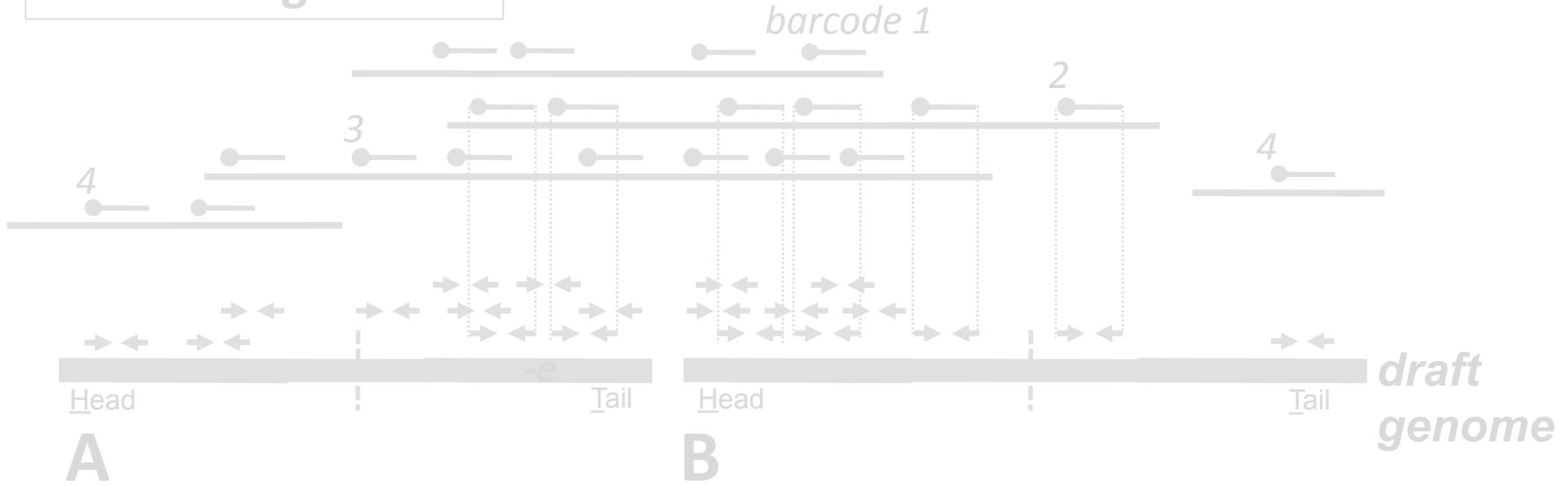
# Method

## 1. Read alignments



# Method

## 1. Read alignments



## 2. Binomial test & tally

Sequence A barcode	Head	Tail	Sequence B barcode	Head	Tail
1	0	2	T-H	1	2
2	0	2	T-H	2	0
3	1	2	T-H	3	0
4	2	0	H-T	4	1

Is read distribution in 5' (head) or 3' (tail) end different from a uniform distribution?



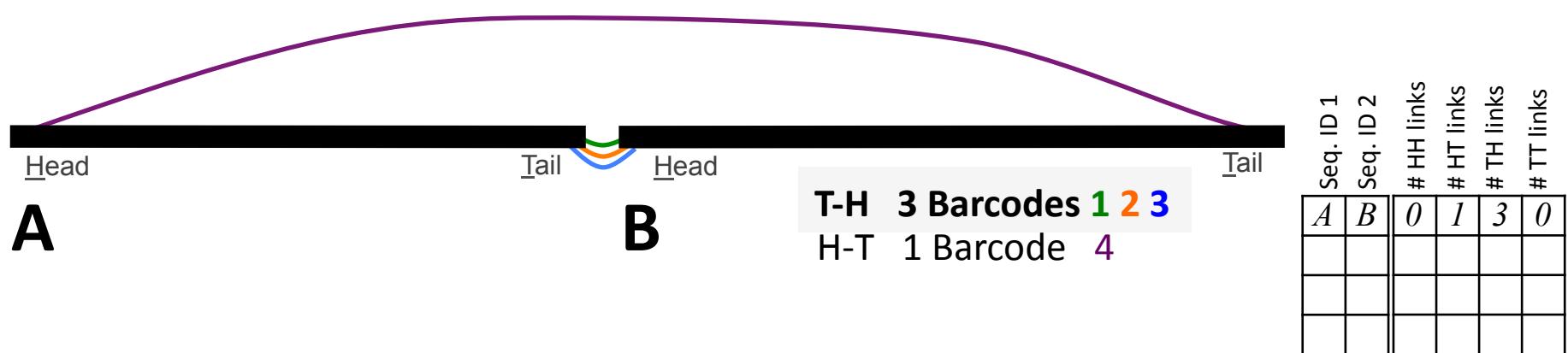
Barcode	(Seq. ID, H/T), count
1	(A, T): 2 (B, H): 2
2	(A, T): 2 (B, H): 0
3	(A, T): 2 (A, H): 1 (B, H): 3
4	(A, H): 2 (B, T): 1

hash tracks head-tail tally  
for each barcode

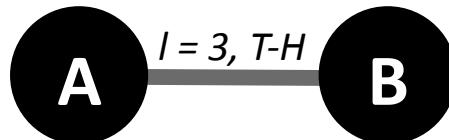
# Method

## 3. Pair sequences

For any two sequences, tally head & tail barcode support



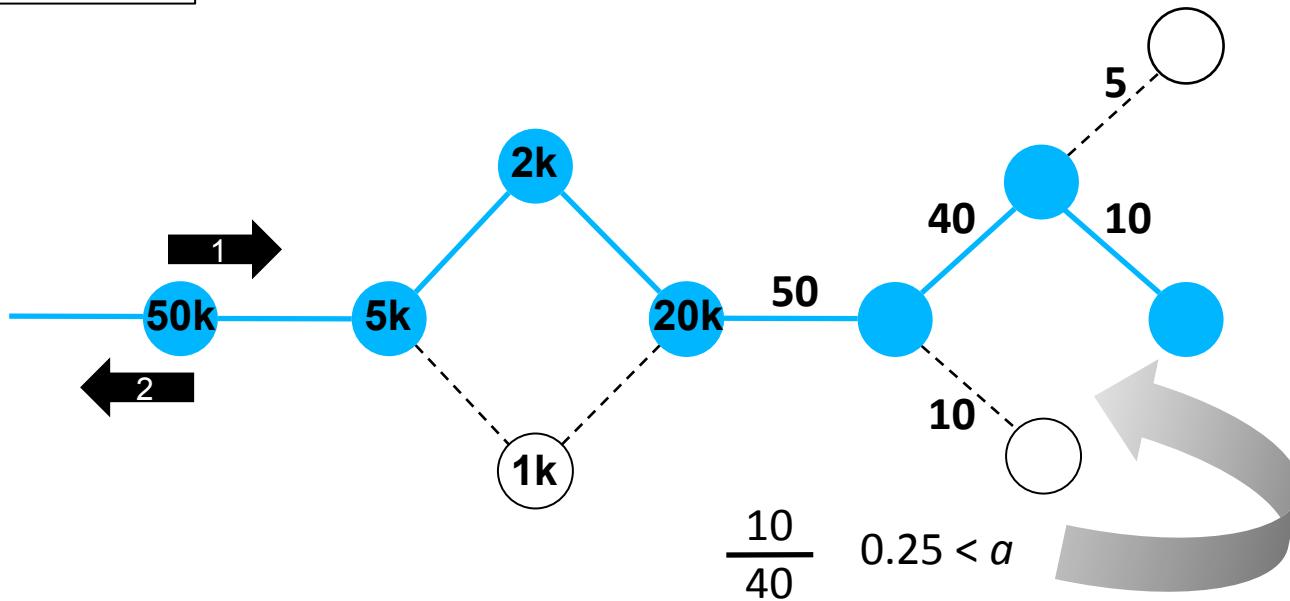
## 4. Output graph



Apply heuristics favoring most supported H-T combo and build graph

# Method

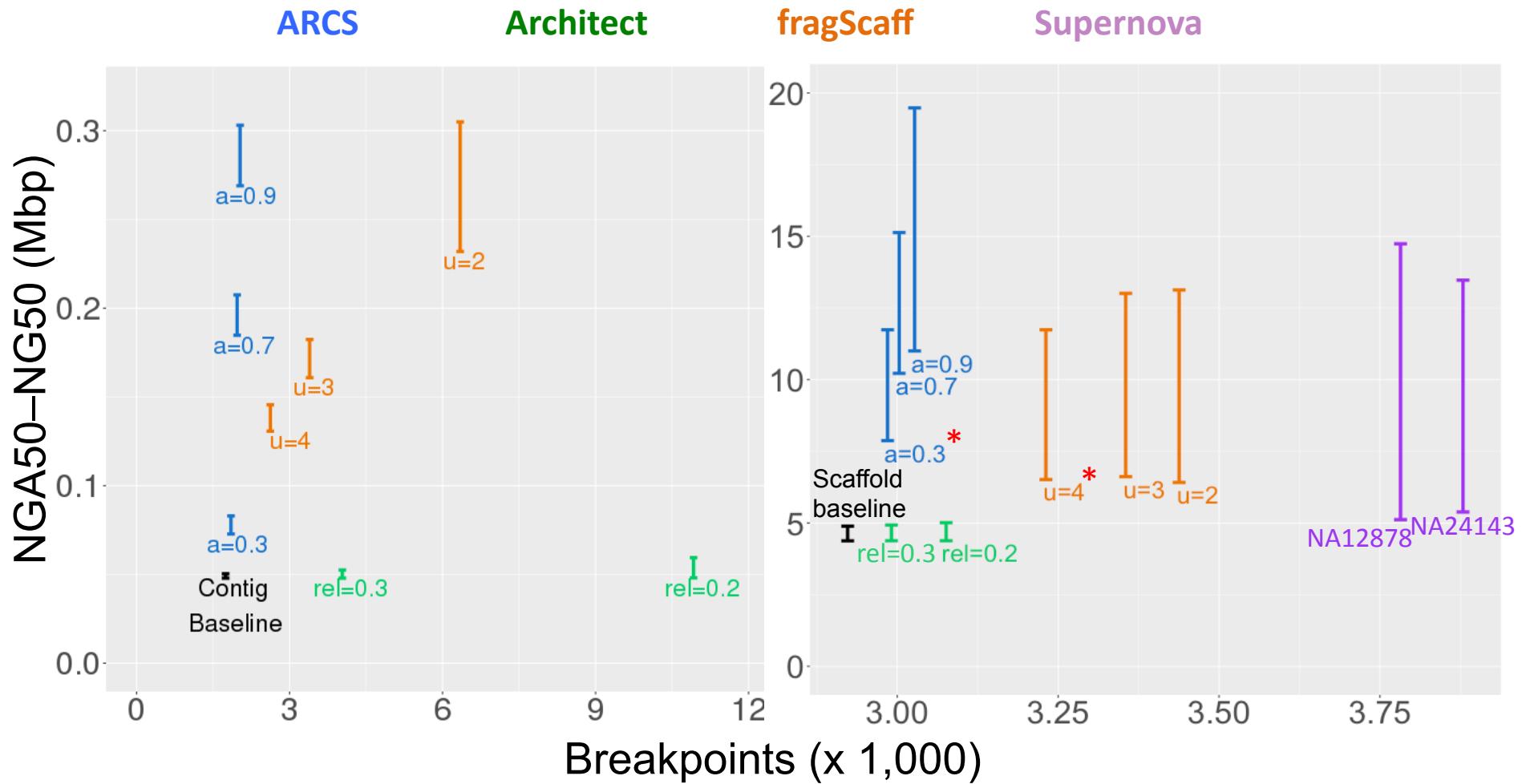
## 5. Scaffolding



Currently utilizes the *LINKS* (Warren *et al.* 2015) scaffolder

# Scaffolding human drafts

comparison to GRCh38 human genome reference



# Scaffolding human drafts

Resources

Tool	Time	Mem
ARCS	<b>0:55</b>	<b>3.4</b>
fragScaff*	1:59	14.1
Architect	6:07	9.6

\*64 threads

h:mm

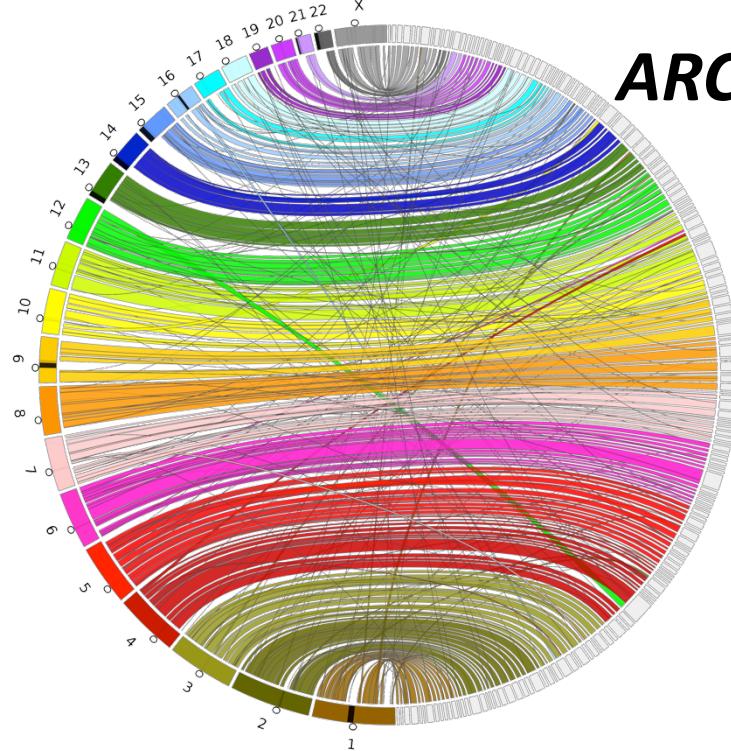
GB

Misassemblies

Assembly	Intra	Inter
Baseline	1090	399
ARCS	<b>1228</b>	<b>434</b>
fragScaff	1911	510

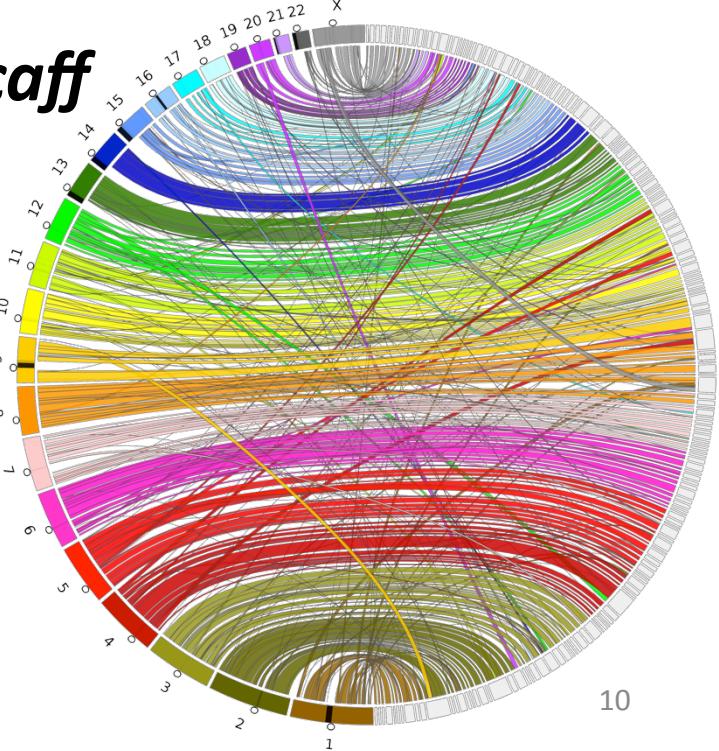
Inconsistent contig order

*ARCS*



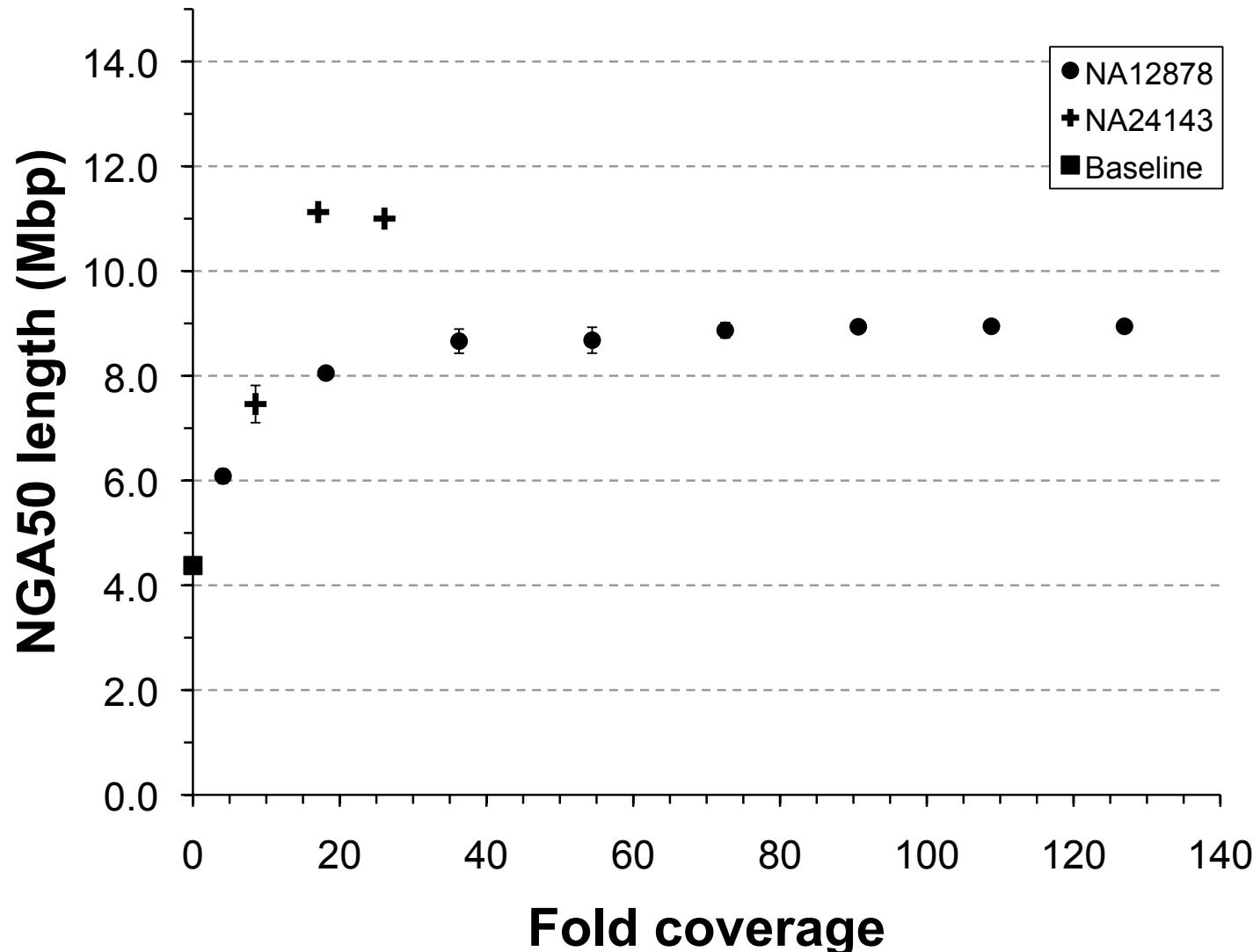
NA24143  
1kbp+ GRCh38  
alignments

*fragScaff*



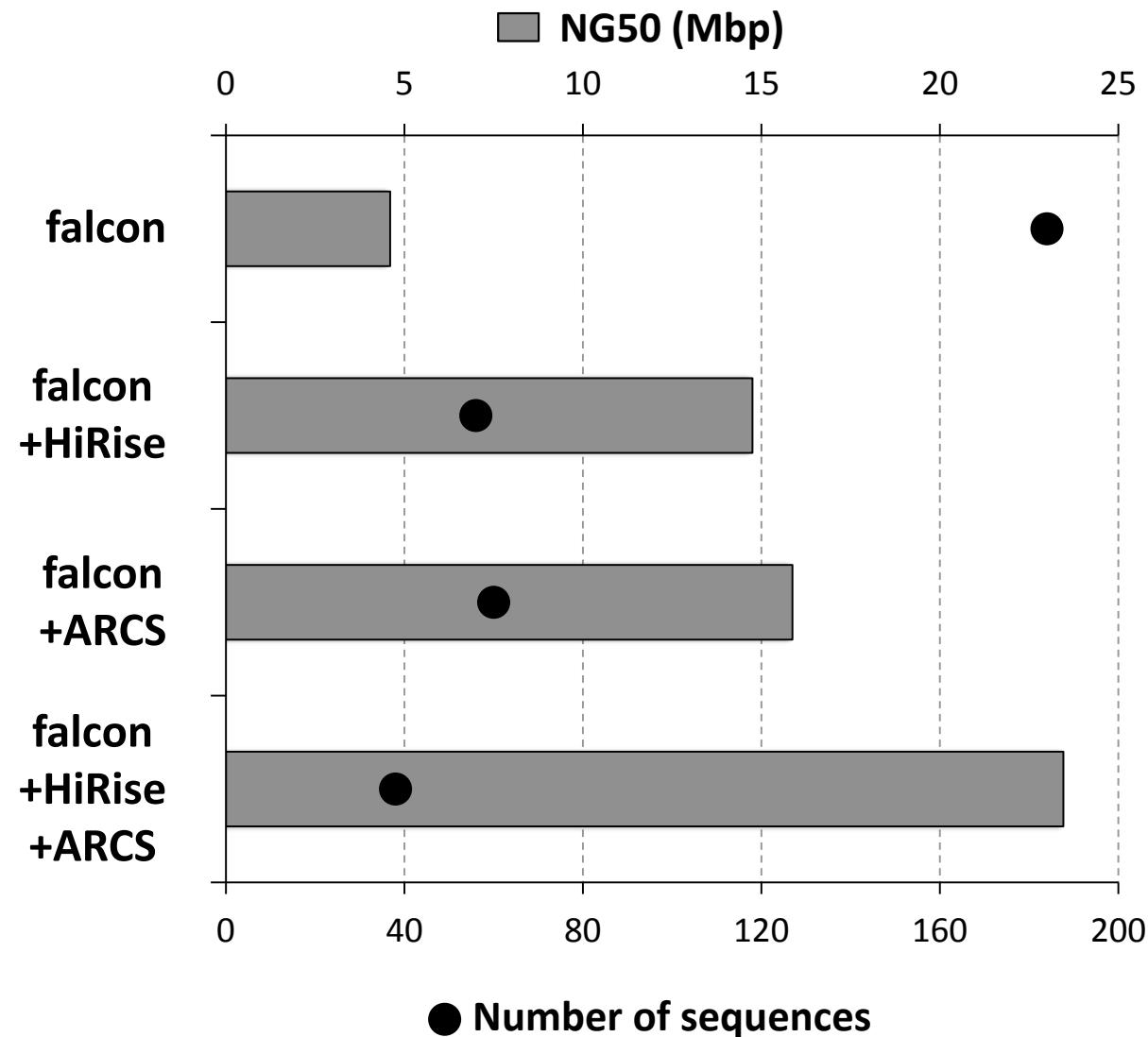
# Scaffolding human drafts

works with low (20X) Chromium sequence coverage



# Retrospective scaffolding

long-read (pacbio) NA24143 sequencing



# Spruce genome scaffolding



## Economically important lumber

- **Aims**

- Improve wood quality, pest resistance

- **Challenges**

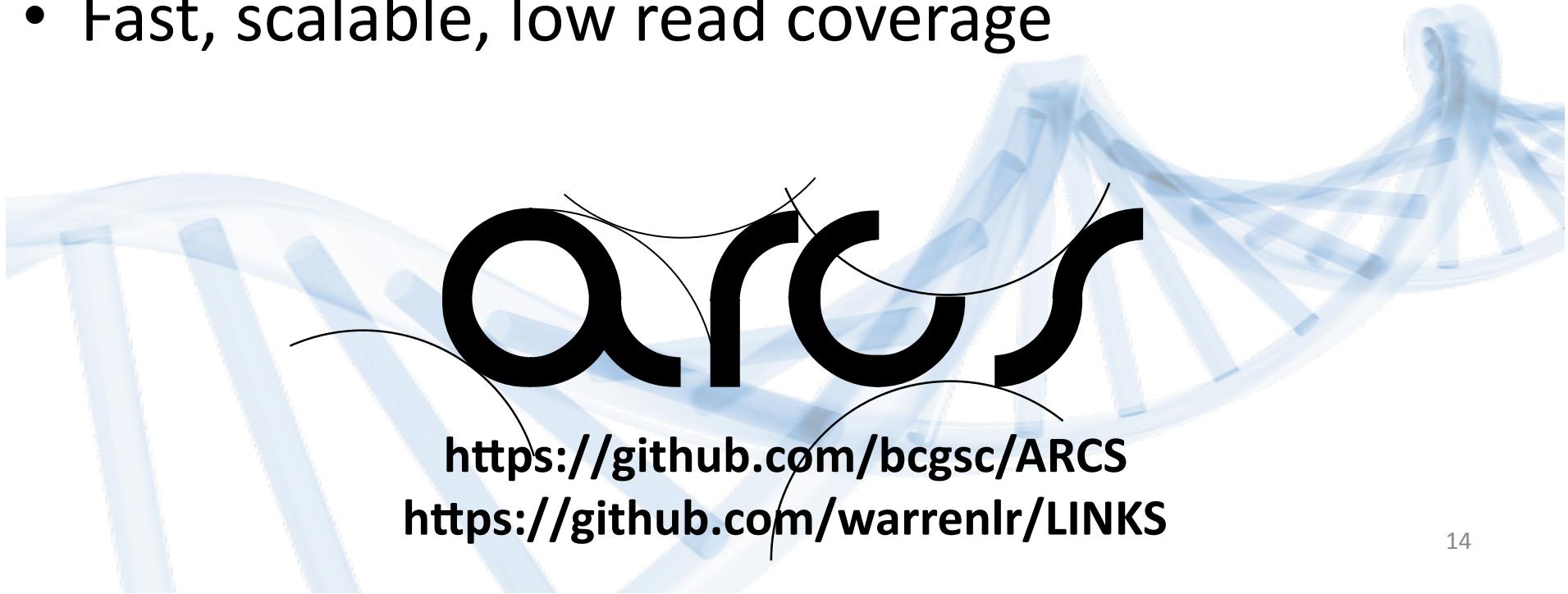
- Large 20 Gbp genome, drafts fragmented

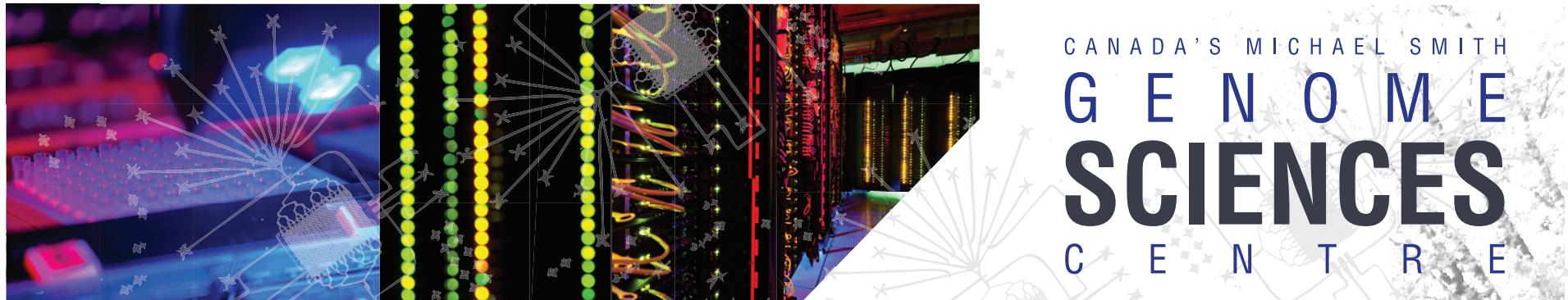


- 8 Chromium lanes (~50-fold coverage)
- 3-fold NG50 increase to 315 kbp (from 87 kbp)
- 9 additional complete genes recovered
  - CEGMA (42%) and BUSCO (31%)

# Summary

- Stand-alone, linked-read genome scaffolder
  - retrospective scaffolding
- *Mbp*-range, accurate drafts from short reads
- Fast, scalable, low read coverage

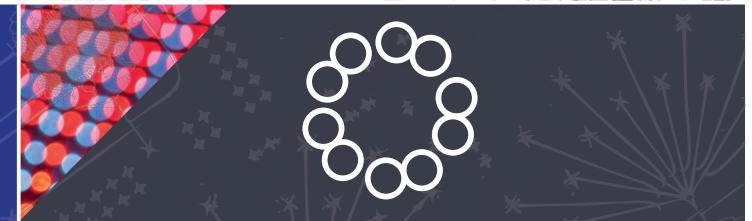




CANADA'S MICHAEL SMITH  
**GENOME**  
**SCIENCES**  
CENTRE

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

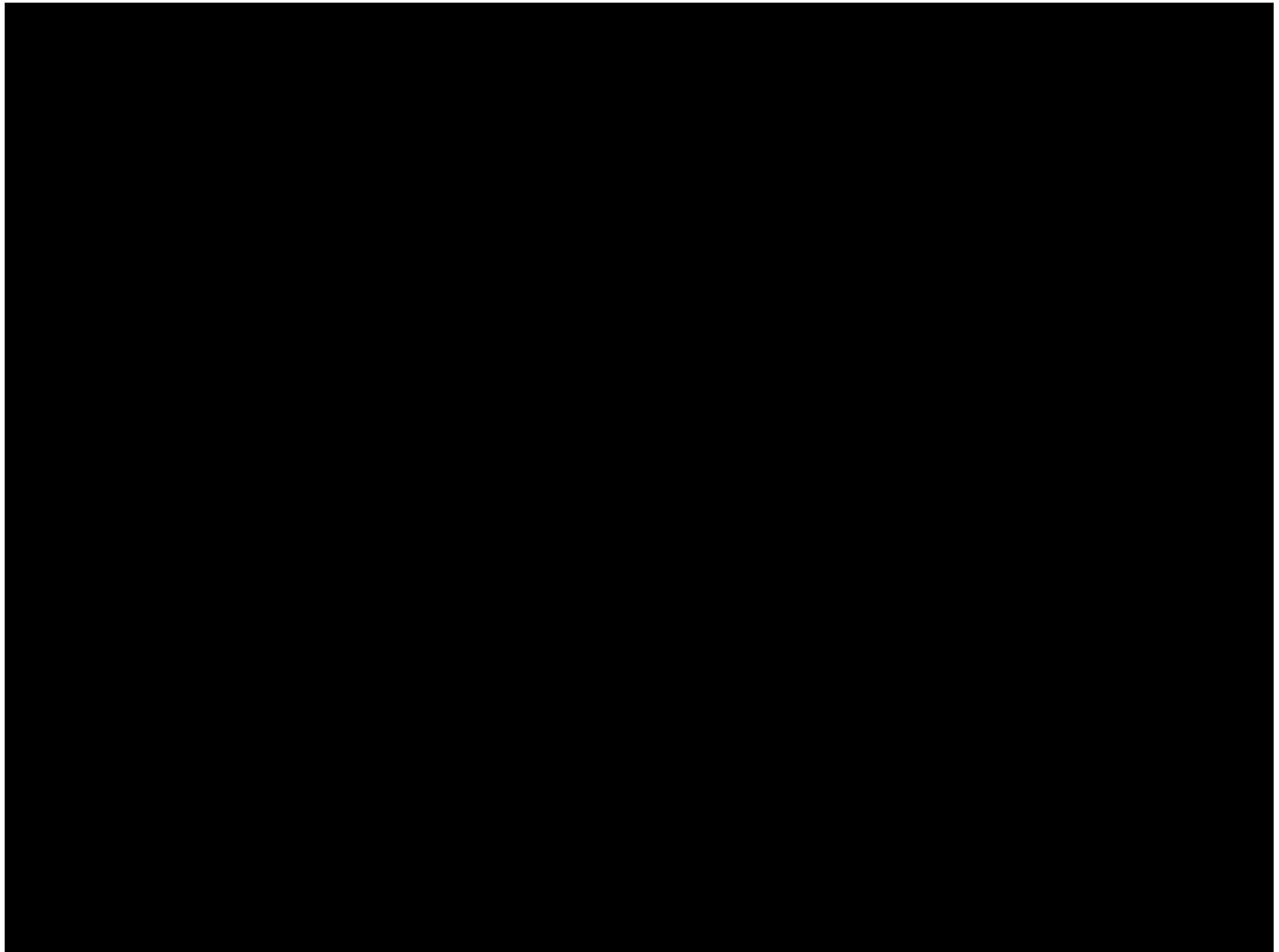
# Thank You!

***Sarah Yeo | Lauren Coombe | Justin Chu | Ben Vandervalk  
Shaun Jackman | S. Austin Hammond | Inanc Birol***

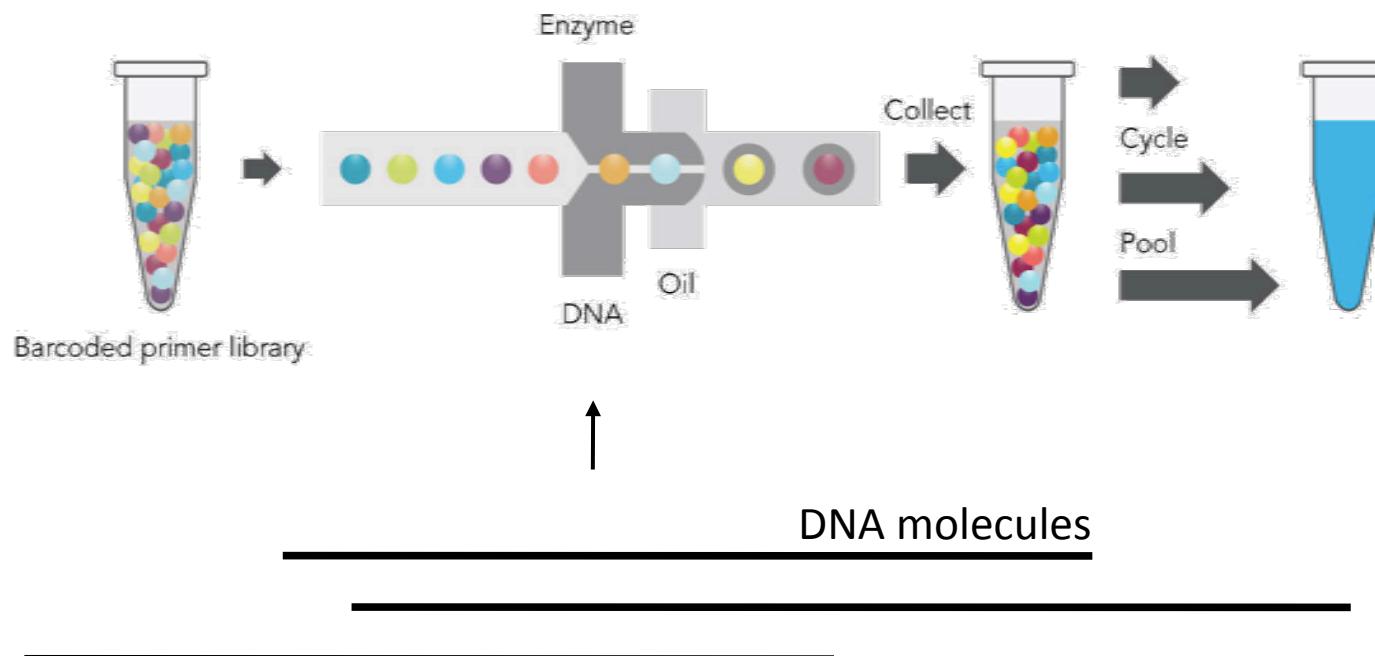
*Funding*



**BCCA CANCER RESEARCH CENTRE**

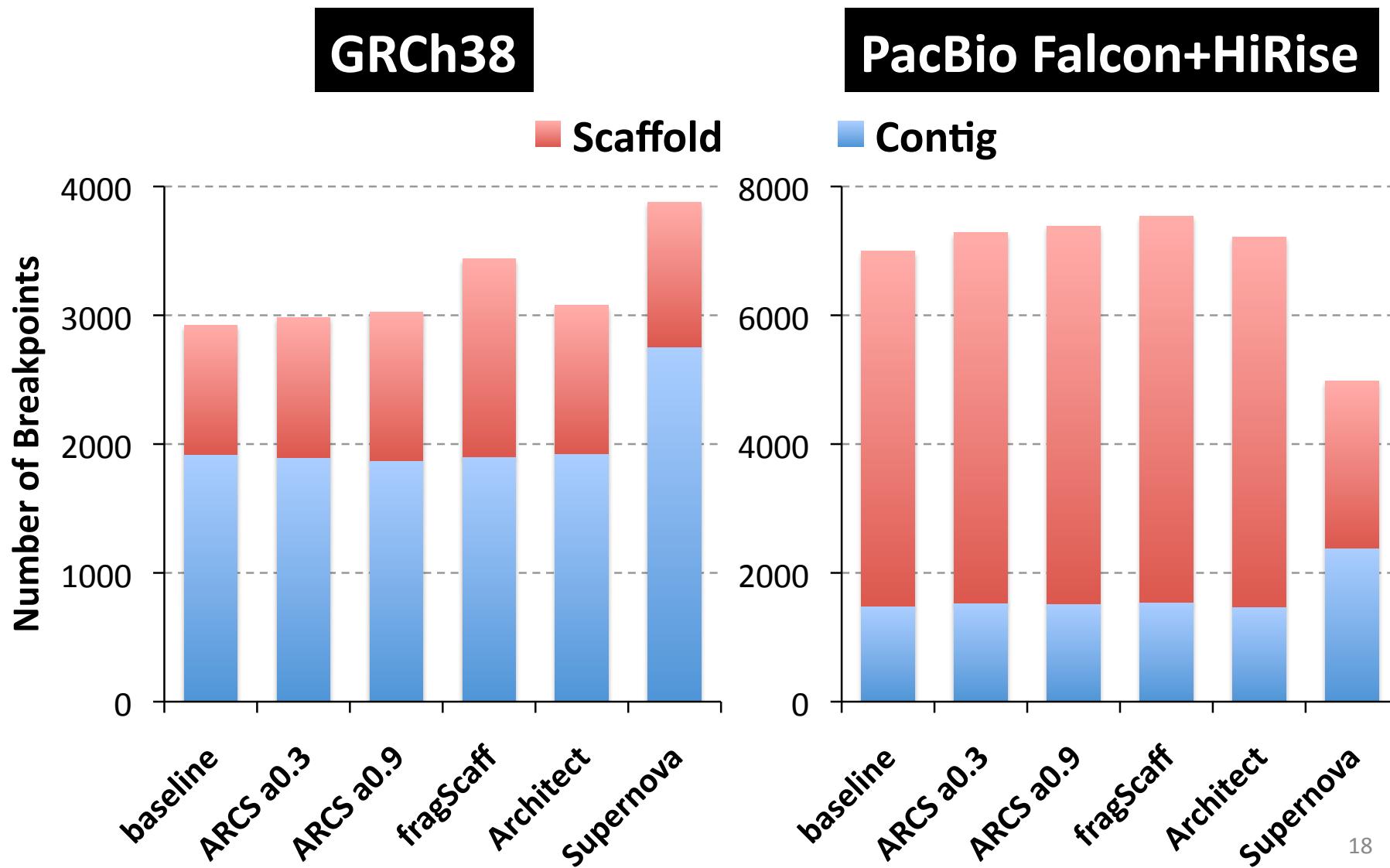


# Overview of 10X Genomics



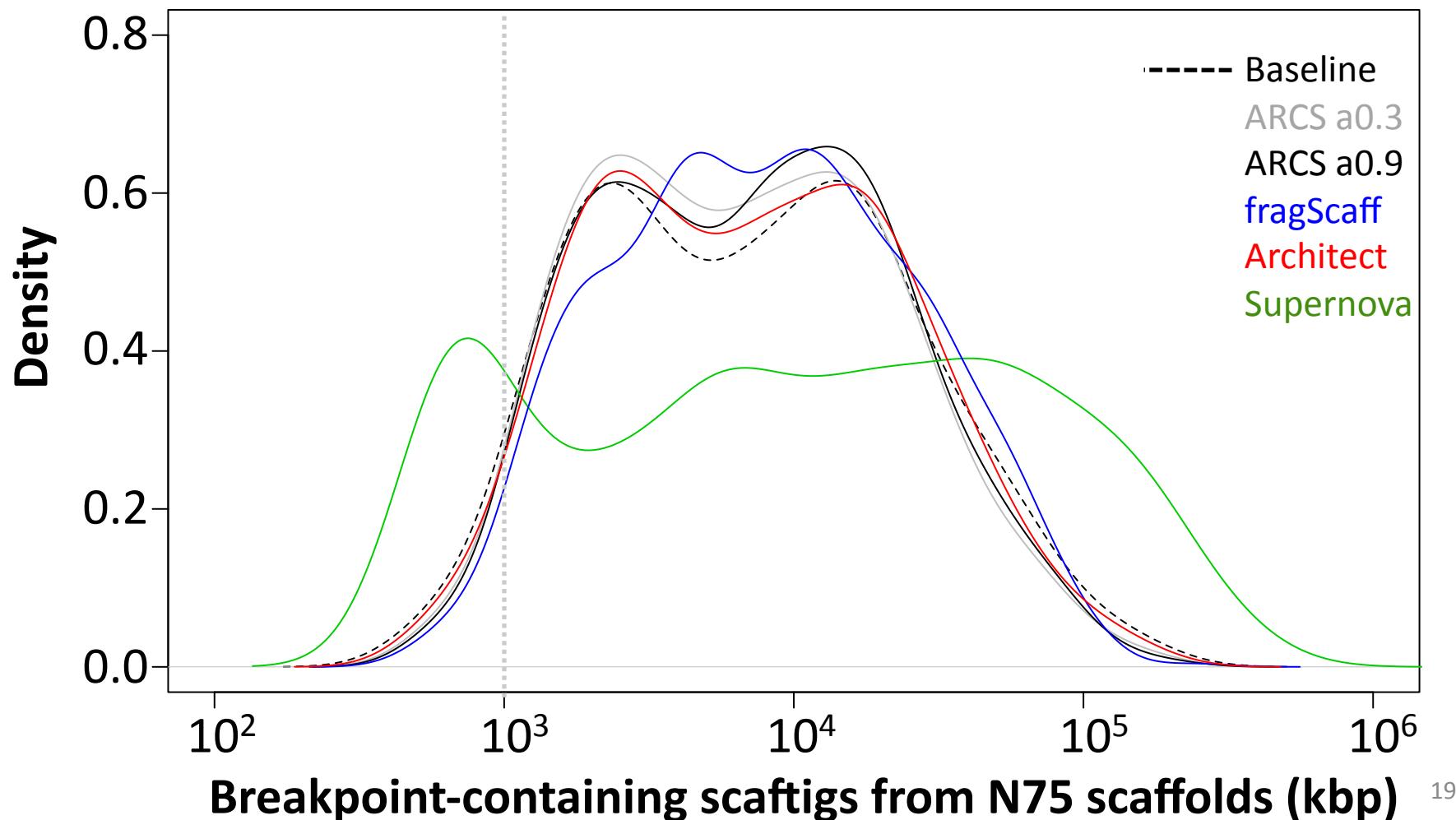
# Scaffolding NA24143 assembly drafts

## breakpoint analysis



# Scaffolding NA24143 assembly drafts

ScafTigs	n:500	N50
Baseline	210,391	33,905
Supernova	69,299	115,158



# Scaffolding human drafts

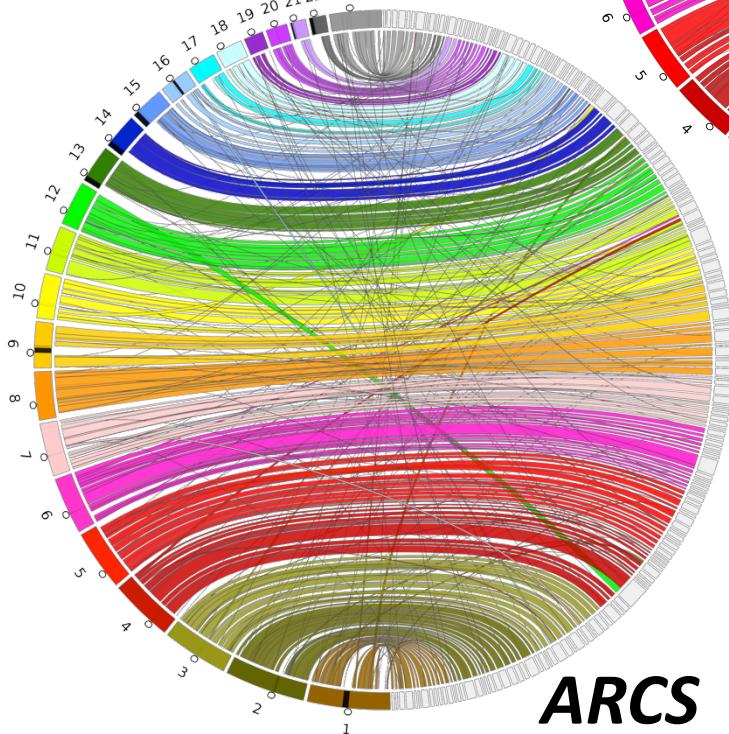
## Resources

Tool	Time	Mem
ARCS	<b>0:55</b>	<b>3.4</b>
Architect	6:07	9.6
fragScaff*	1:59	14.1

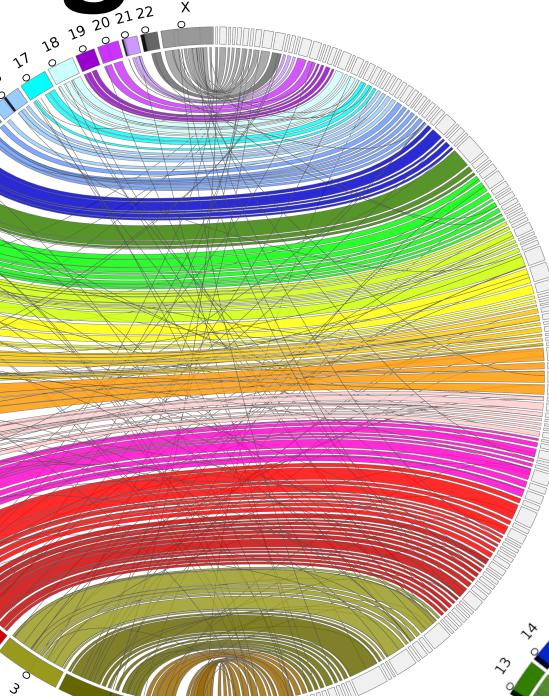
\*64 threads

h:mm

GB



*ARCS*

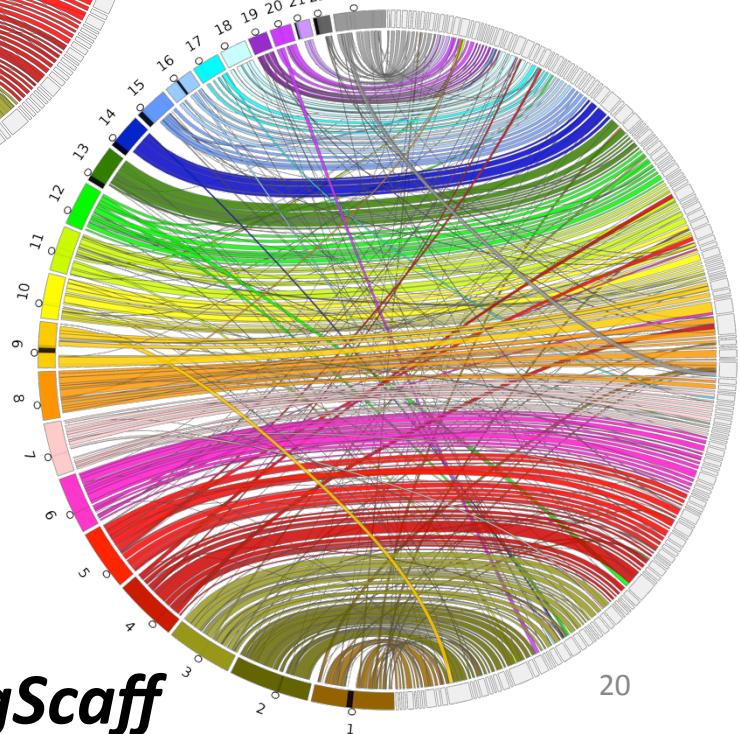


*Supernova*

NA24143  
1kbp+ GRCh38  
alignments

Assembly	Intra	Inter
Baseline	1090	399
ARCS	1228	<b>434</b>
fragScaff	1911	510
Supernova	<b>1198</b>	439

Inconsistent contig order

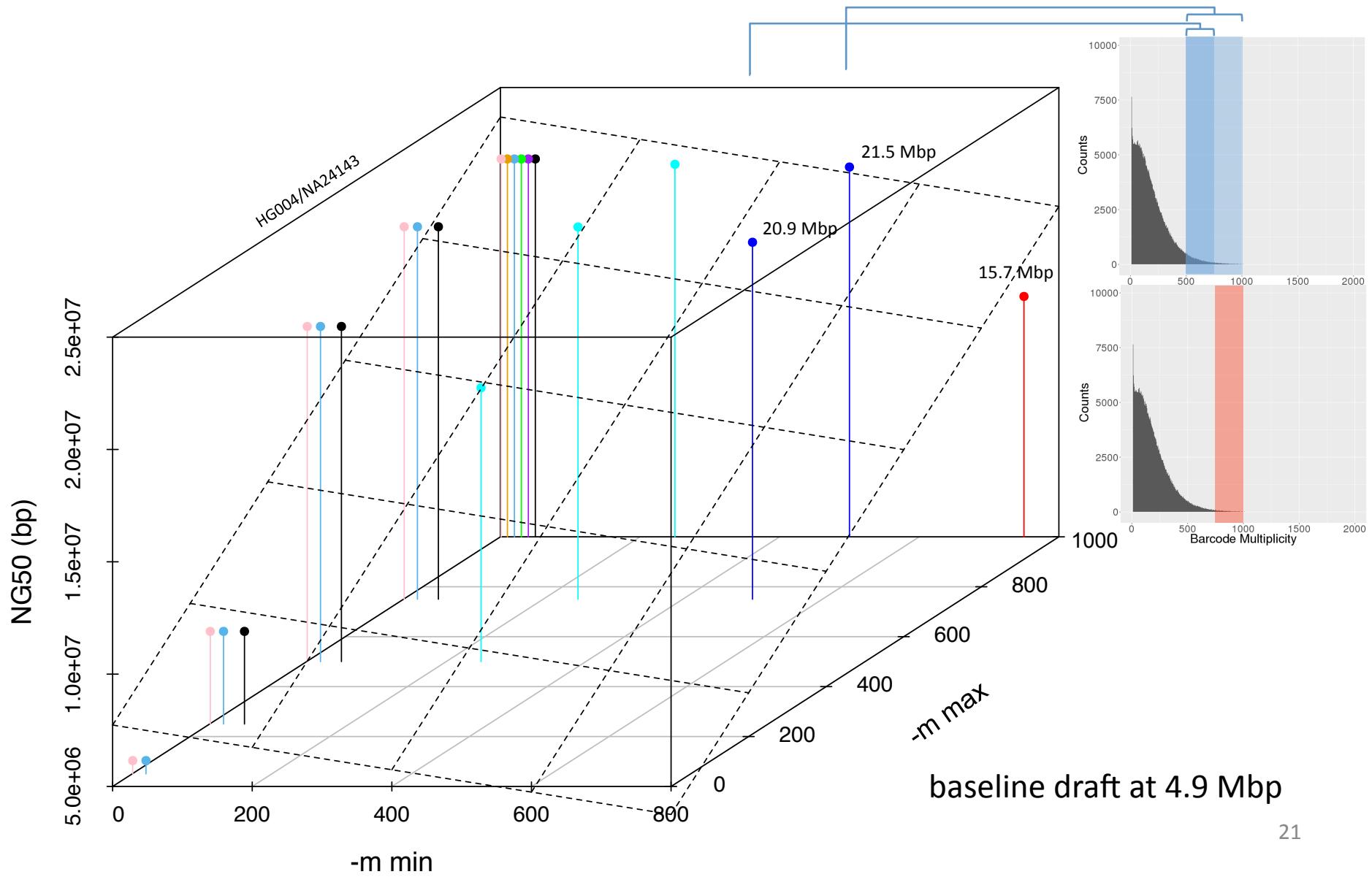


*fragScaff*

20

# Scaffolding human drafts

Effect of  $-m$  barcode multiplicity range parameter



# Spruce genome scaffolding

- 8 lanes Chromium ~50-fold
  - NG50 : 315 kbp (from 87 kbp)
  - Recovers 9 additional CEGs (CEGMA and BUSCO)

