

Sequencing the un rearranged human Immunoglobulin heavy chain locus (IgH) from hydatidiform mole

SFAF
Santa Fe, USA
René Warren – June 2010



Immunoglobulin (Ig)

- Loci comprises V,(D),J gene segments that recombine to make heavy chain (IgH) or light chain (IgK, IgL) messages
- Combinatorial rearrangements and “somatic hypermutation” (addition/removal of non-templates bases) creates astonishing theoretical antibody diversity 10^{14}
- Germline variation can confer disease susceptibility
 - Defective IgK V allele fail to protect Navajos natives against *H.influenzae* type B (Feeney, 1996)
 - Homozygous IgH-V deletions can cause nephritis (Lupus) (Cho, 2003)

IgH reference shortcomings

- hg19-IgH consensus is 1.28 Mbp from Matsuda et al. 1998
 - IgH reference is a mosaic from diploid libraries from various sources
 - Sample source from lymphoid origins
 - lymphoblastoid cell lines
 - peripheral blood cells
 - Ig and TCR loci rearranged in lymphoid tissue
 - V(D)J somatic rearrangements that occur in lymphoid tissue and captured by sequencing predicted to include bases not encoded by the human genome and rearranged Ig loci
 - hg19-IgH not a haploid reference

The diagram illustrates the IgH locus structure. It features a central horizontal line representing the genome with various genomic features labeled along it. Above the line, several grey boxes represent different rearranged segments, each with a unique identifier (e.g., 3-41, (I)40-1, 7-40, (II)38-1, 3-38, 3-37, 3-36, 3-35, 7-34-1, 4-34, 3-32, (II)31-1, 4-31, 4-30-2, (II)30-1, 3-30-5, 3-30-3, 4-30-2, 4-30-1). Below the line, other genomic features are labeled: G3, G1, EP1, A1, GP, G2, G4, E, A2, 3', Enhancer, D7-27, J1 to ~60, K1 to ~1000 kb, and 1 250 kb. The diagram also includes a scale bar indicating distances of 600 kb, 800 kb, and 1 250 kb.

Hydatidiform mole



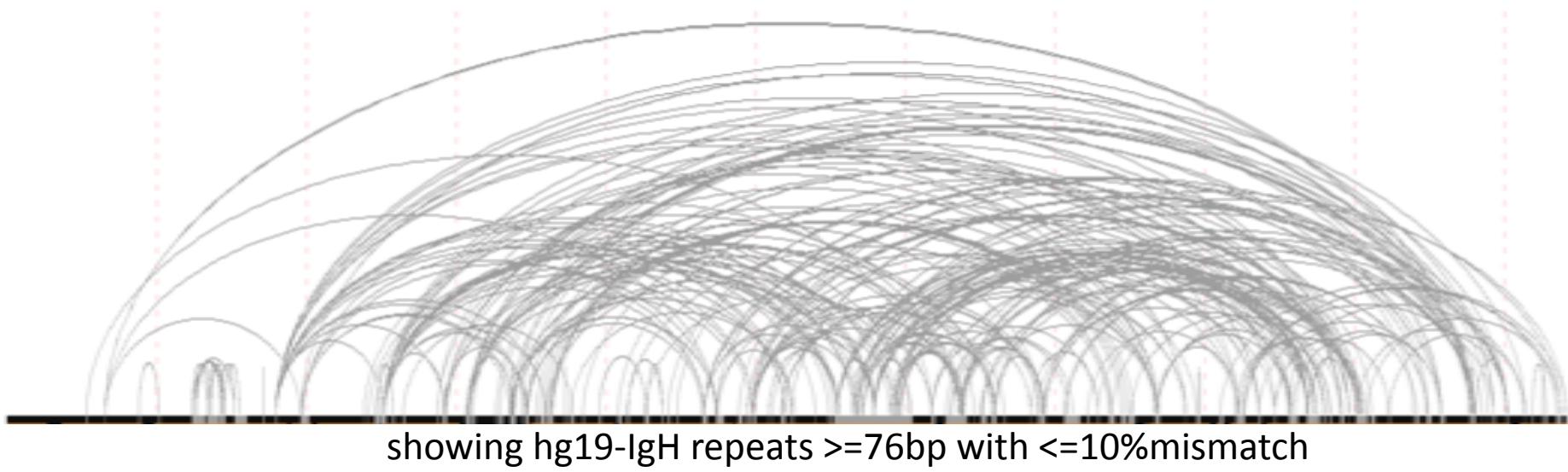
- Pregnancy abnormality (human) generated by fertilization of an empty egg
- Haploid, non-lymphoid
 - Source of unarranged Ig gene segments
- BAC physical map exists (BCGSC) and resources available [BACPAC] (CHORI)



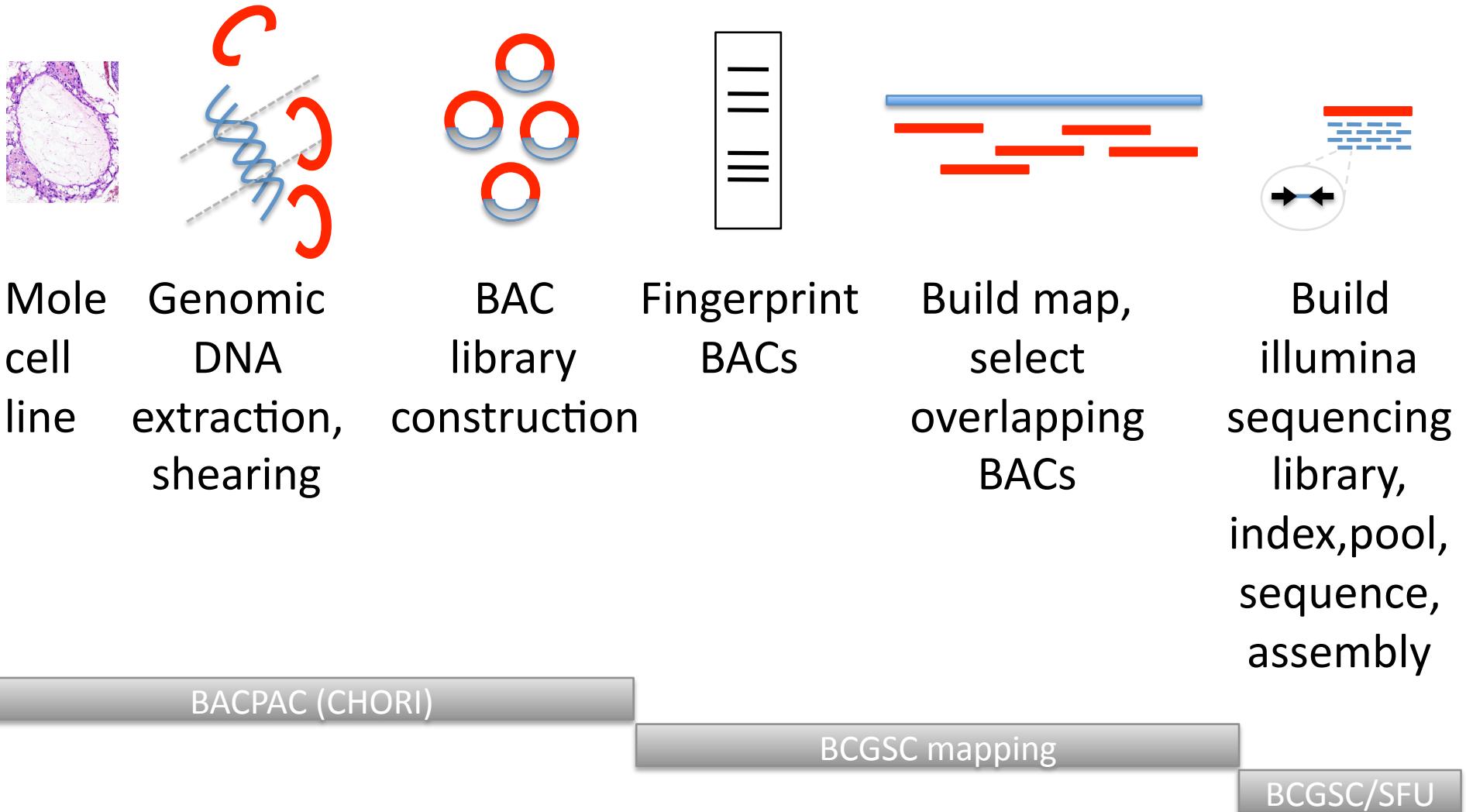
1.0 mm

Sequencing challenges

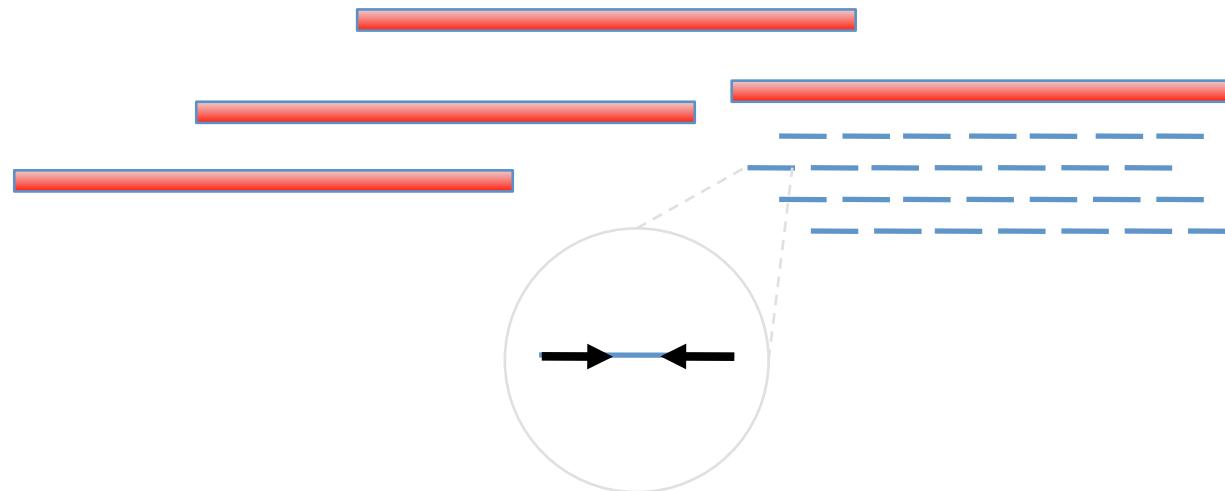
- Human genome interspersed repeats
- Ig gene segments (VDJ) share high sequence homology
- IgH most segmentally duplicated locus in the human genome



Approach



High-throughput sequencing



- IgH, IgK and IgL tiling path BAC sequences sheared
- BAC Sequence library built & indexed for HT sequencing
 - hexamer barcode (Illumina)
- BAC libraries pooled
- Sequenced onto 1 lane Illumina GAIIx instrument
 - 27.6M x 76bp PET, 200nt insert size
 - ~1950-fold average sequence coverage

De novo short read assembly

- Using in-house assembler SSAKE (v3.5)

- 1st short read assembler
 - continuously improving



- Handles errors, mate pairs, reads of different lengths
 - Used in conjunction with base quality trimming
- Suited for IgH sequencing from BACs
 - Certain short read assemblers mishandle BAC overlap sequence redundancy as large repeats

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

What's new in SSAKE v3.5?

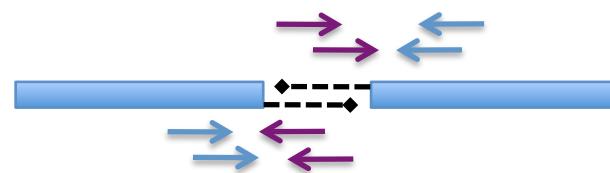
① PET usage at run-time

- helps resolve repeats during contig extension



② Fill gaps between adjacent contigs in a scaffold

- fill gaps with missing/misplaced mates



E.coli K12 short read test set (SRA)

12.2M x 36bp mate pairs from Illumina GA :: 6.4M pairs and 2.9M unpaired quality trimmed

Metric	SSAKE v3.4	SSAKE v3.5	
time	0:37:05	0:39:18	Fewer contigs & scaffolds providing higher coverage and higher scaffold contiguity
# contigs	896	856	
N50 (bp)	12,433	11,946	
largest (bp)	50,675	48,852	
bases (bp)	4,569,095	4,569,002	
% coverage	96.3 (880)	97.0 (842)	
% accuracy	99.98	99.98	
# merged contigs (fill gaps)	NA	563	run-time PET usage
N50 (bp)	NA	21,664	
largest (bp)	NA	70,162	
bases (bp)	NA	4,568,849	
% coverage	NA	93.4 (544)	
% accuracy	NA	99.88	
# scaffolds	270	224	N50 doubled, #contigs halved after automated gap closure
N50 (bp)	55,882	69,464	
largest (bp)	134,207	221,647	
Satisfied pairs (%)	99.75	99.79	

IgH assemblies

Metric	SSAKE v3.4	SSAKE v3.4 data cleaned*	SSAKE v3.4 stringency up	SSAKE v3.5 pair@RT/fill gap
# contigs^	1,817	1,293	936	725
N50 (bp)	1,444	2,553	3,181	4,704
largest (bp)	14,322	19,427	19,422	19,421
bases (bp)	1,390,273	1,222,204	1,128,894	1,119,786
# merged contigs (fill gaps)	na	na	na	508
N50 (bp)	na	na	na	8,313
largest (bp)	na	na	na	31,411
# scaffolds	720	689	431	390
N50 (bp)	9,199	9,785	11,535	12,585
largest (bp)	61,664	53,487	49,960	50,660
Satisfied pairs (%)	95.8	96.3	97.2	97.9

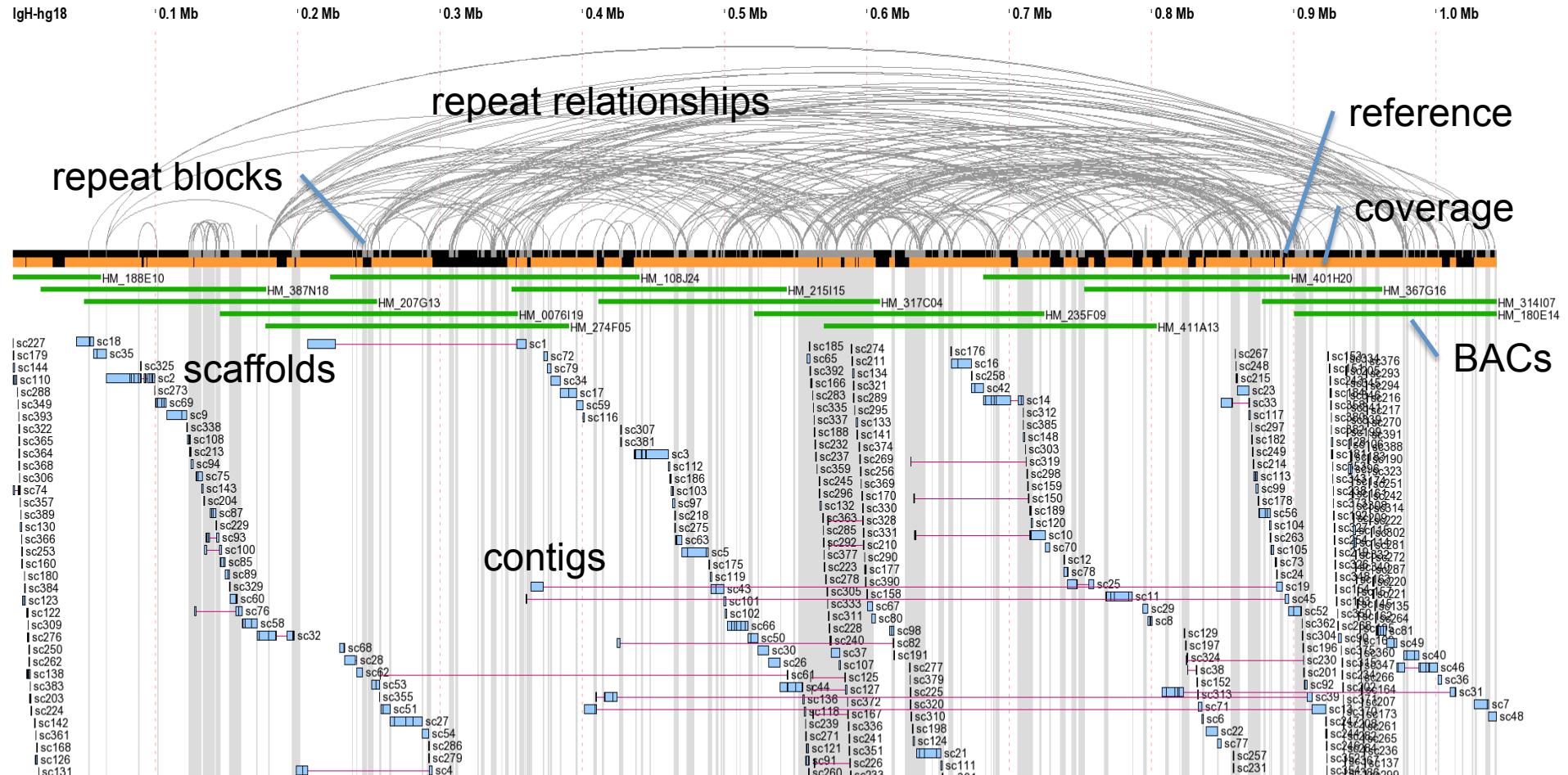
*vector and *E.coli* reads removed, segregate IgH reads on HQ indexes (>=Q20)

^contigs >=200bp analyzed

RT: run-time

SC: scaffolding

Comparison to a reference



WGA alignment to reference

Scale 1:500

Identity >= 90 Blast Alignment >= 90 Alignment Size >= 100 Repeat Size >= 76 Repeat Base Mismatch <= 10

Showing Scaffolds with 1+ contigs

360 scaffolds, 625 contigs satisfy alignment filters

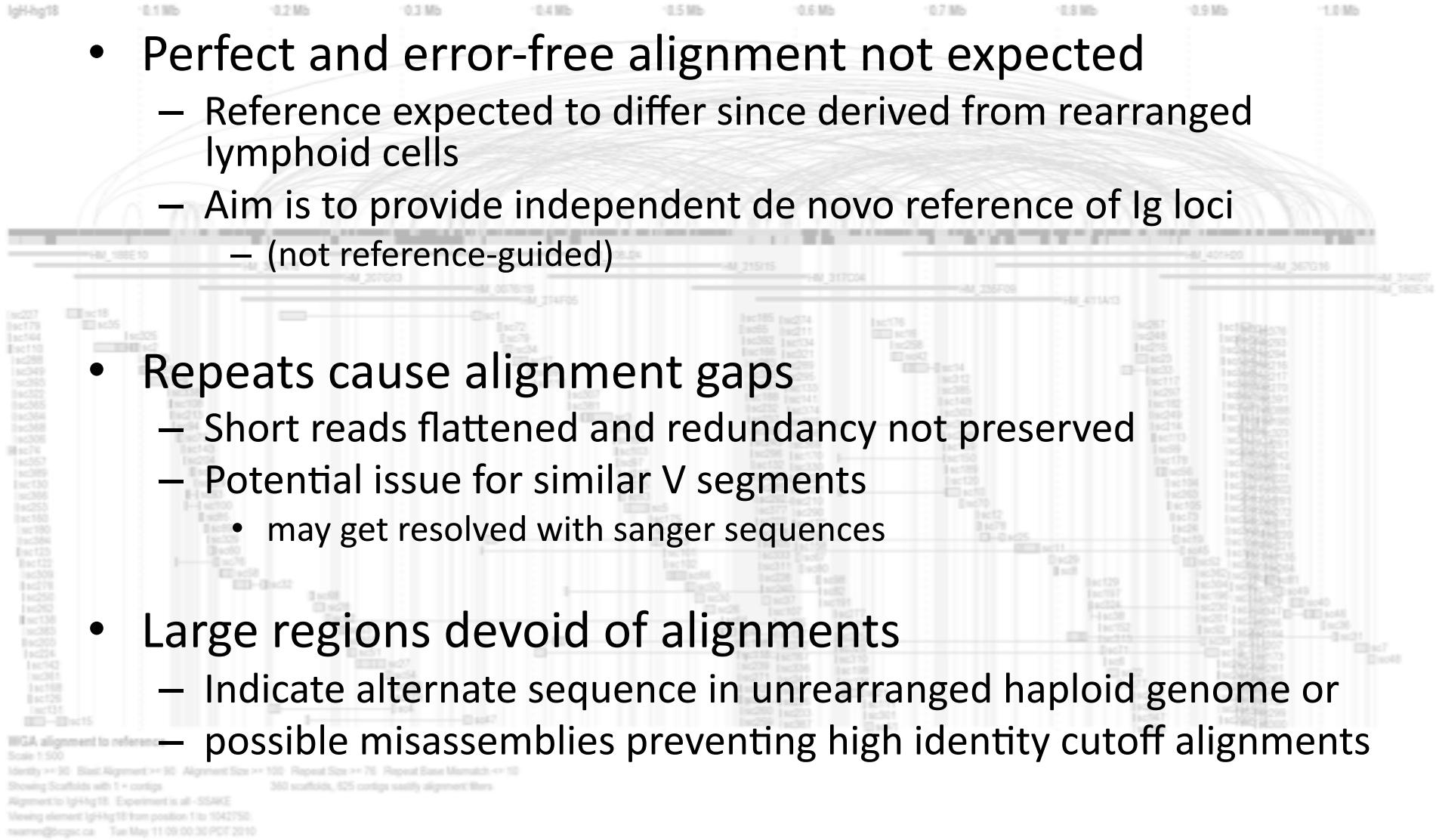
Alignment to IgH-hg18. Experiment is all -SSAKE

Viewing element IgH-hg18 from position 1 to 1042750.

rwarren@bcgsc.ca Tue May 11 09:00:30 PDT 2010

80% coverage by contigs >=90% seq.id.

Comparison to a reference



In the works

- 5kb plasmid library from minimum BAC tiling set
 - light sequence coverage from plasmid clones
- Additional SSAKE developments
 - Support for paired-end sanger sequences
 - improve [short+long range] contiguity / fix local misassemblies
 - Support for libraries of different insert sizes
- Phase I: illumina as a framework to produce draft sequence
- Phase II: Assemblies of illumina+sanger data set and directed finishing
 - Primer walking
- Goal is to get a fully finished IgH sequence

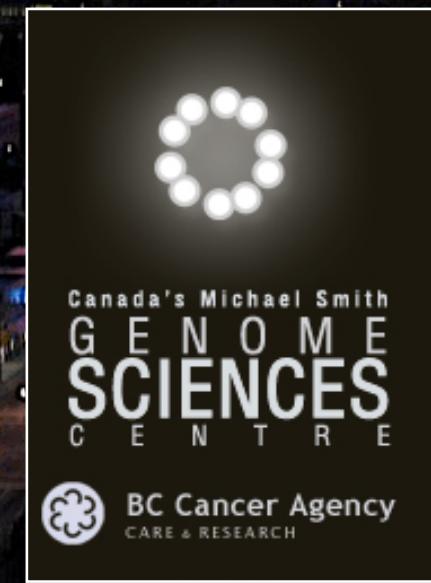
Acknowledgements

BC Genome Sciences Centre

- Rob Holt
- Jacquie Schein and BCGSC mapping
- BCGSC sequencing group

Simon Fraser University

- Felix Breden
- Corey Watson



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