



Swati Ranade¹, Jason Chin¹, Steve Kujawa¹, Paul Peluso¹, John Harting¹, Lawrence Hon¹, Richard Hall¹, Primo Baybayan¹, Kevin Eng¹, Julie Karl², Michael Graham², Roger Wiseman², David O'Connor², Daniel E. Geraghty³, Chul-Woo Pyo³, Jason Howard⁴, Erich Jarvis⁴ & Jonas Korlach¹

¹PacBio, 1380 Willow Road, Menlo Park, CA; ²University of Wisconsin, Madison, WI; ³Fred Hutchinson Cancer Research Center, Seattle, WA; ⁴Howard Hughes Medical Institute & Duke University, Durham, NC

Abstract

Fine mapping of causal immune-gene variants for their known association with cancer, drug-induced hypersensitivity and autoimmune diseases, are of paramount importance¹. As population-scale medical genomics initiatives' collection of short read whole genome and exome sequencing (WGS)/(WES) data and SNP array data for Genome Wide Association Studies (GWAS) has increased, the development of imputation software for analysis of immune gene variants has been rampant.

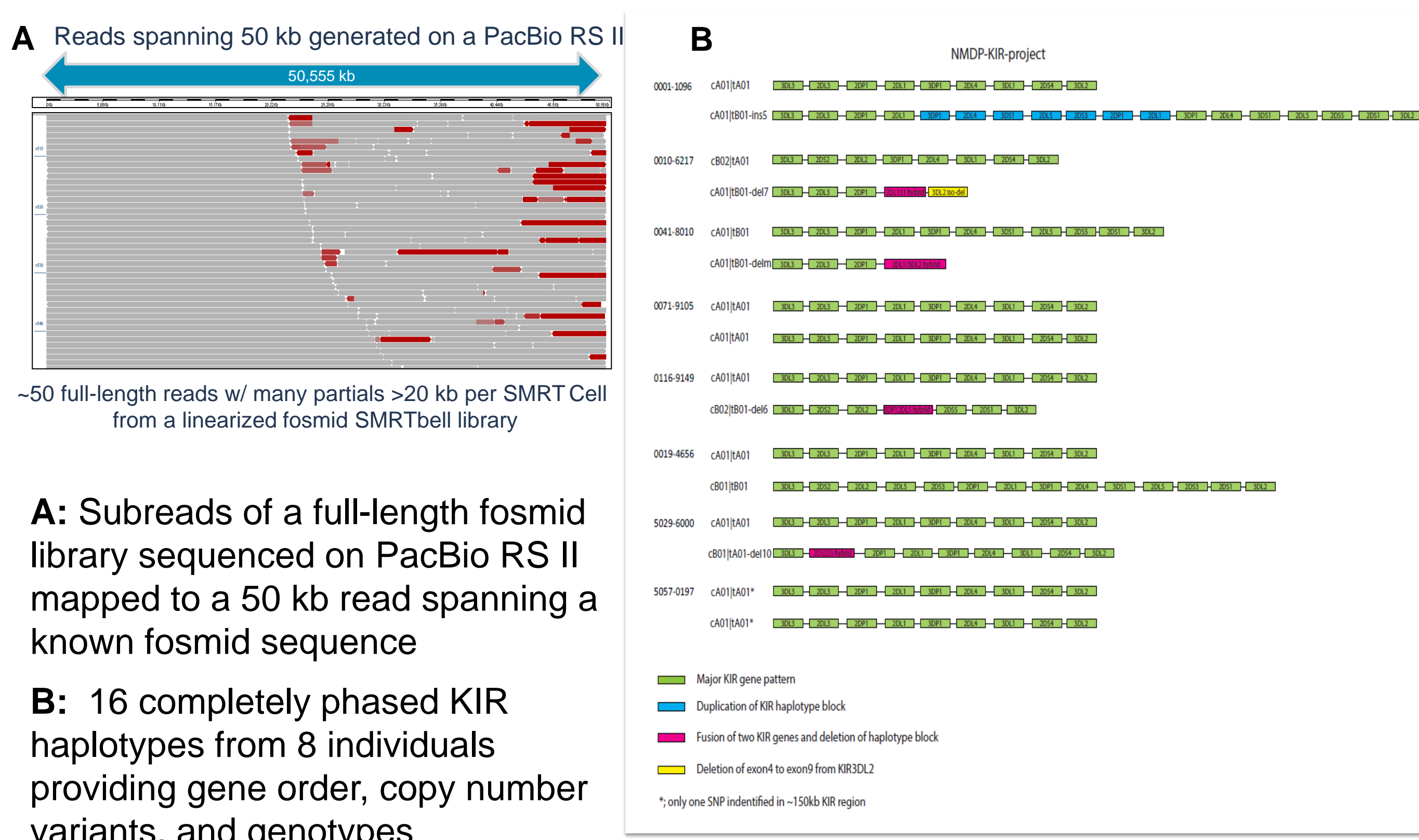
Despite increased efforts, fine mapping has yielded few discoveries¹. This may be because extraction of haplotype information in these complex immune regions characterized by excessive copy number variants (CNVs), high-density hyper-polymorphic genes, pseudo genes, conserved extended haplotypes (CEH) with enormous linkage disequilibrium (LDs), and long stretches of repetitive DNA is challenging for any reference-reliant, short-read sequencing or GWAS methods².

We present detailed views of human and non-human Killer Immune Receptor (KIR) and Major Histocompatibility Complex (MHC) genes and/or haplotypes, captured by SMRT Sequencing. High accuracy, long reads and uniform coverage from SMRT Sequencing allowed allele-level segregation of sequences with full phasing across SNP-poor regions. Using either whole-genome or targeted sequencing approaches, we were able to carry out *de novo* analysis, including avian species.

Human KIR Sequencing

DE NOVO ASSEMBLY OF COMPLETE KIR HAPLOTYPES

- Unambiguous assembly of KIR regions has been difficult to impossible.
- Long reads in combination with a recombineering sample-preparation approach for isolating a tiling of targeted fosmid containing contiguous genomic regions (~35 kb – 50 kb)³ allowed *de novo* assembly of individual KIR haplotypes with high confidence⁴.
- 16 such haplotypes were assembled from 8 individuals (available in NCBI)
- A 9 kb partial reference, available for homozygous sample 0071-9105 in the IMGT database, was 100% concordant with the PacBio sequence.



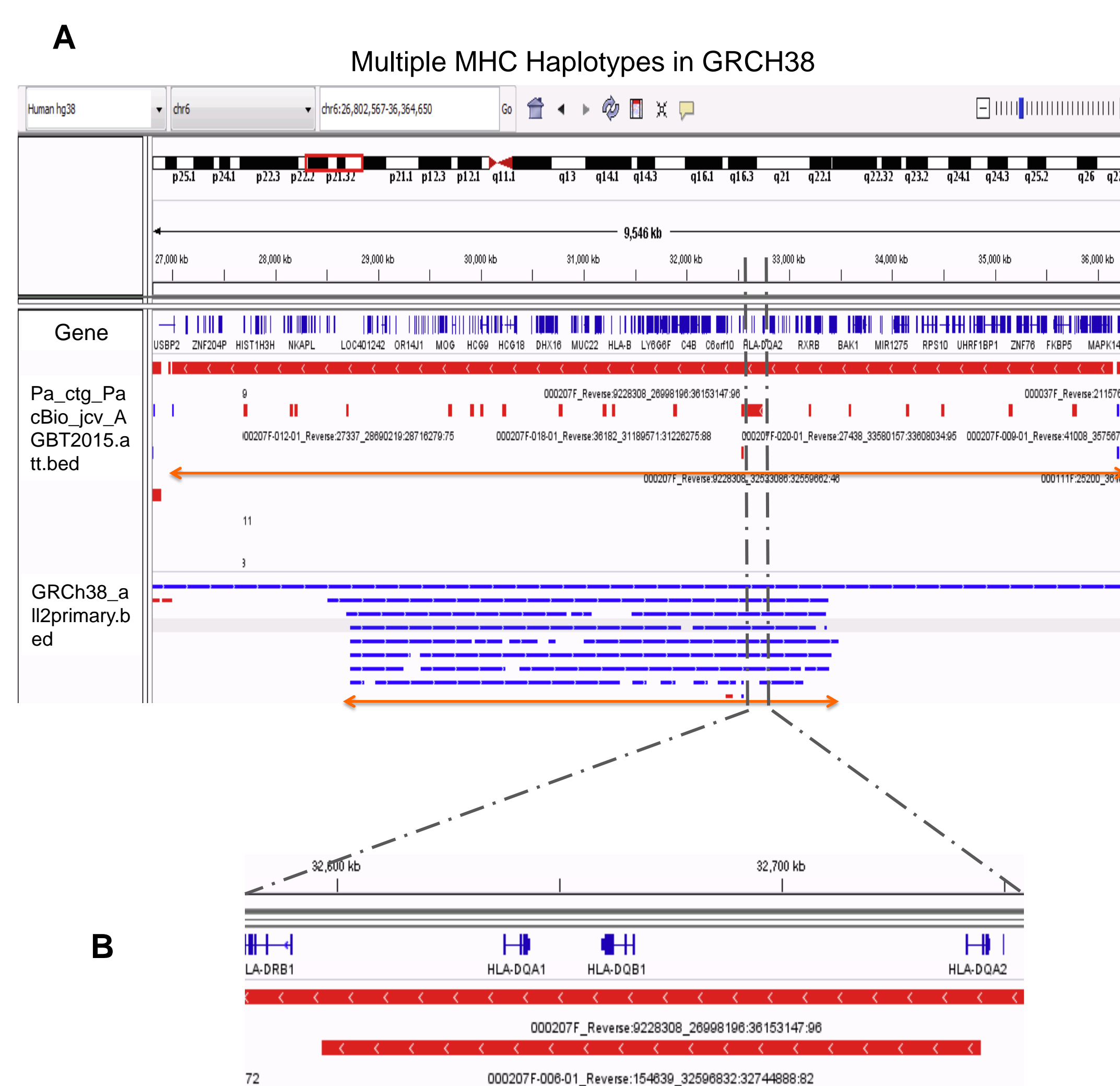
References

- Trowsdale J, and Knight JC. (2013) Major histocompatibility complex genomics and human disease. *Annual Review Genomics Human Genetics*, 14, 301-323.
- Brandt, D. Y. C., et al., (2015) Mapping bias overestimates reference allele frequencies at the HLA genes in the 1000 Genomes Project phase I data. *5(5)*, 931-941.
- Pyo, Chul-Woo et al., *AGBT* (2015) Complete resequencing of extended genomic regions using fosmid target capture and single molecule real-time (SMRT) long read sequencing technology.
- Hall, Richard J. et al., *AGBT* (2015) Assembly of complete KIR haplotypes from a diploid individual by the direct sequencing of full-length fosmids.
- Chin, J. et al. *AGBT* (2015) Toward comprehensive genomics analysis with *de novo* assembly.
- Ranade, S. et al. *AGBT* (2015) Access full spectrum of polymorphisms in HLA class I & II genes, without imputation for disease association and evolutionary research.
- DNA oligo hybridization with Nimblegen SeqCap EZ System ([Application Note](#))
- Zhang et al. *Science* (2015) Comparative genomics reveal insights into avian genome evolution and adaptation. *Science* 346(6215), 1311-1320.

Human MHC Sequencing

DE NOVO PHASED MHC COMPLEXES

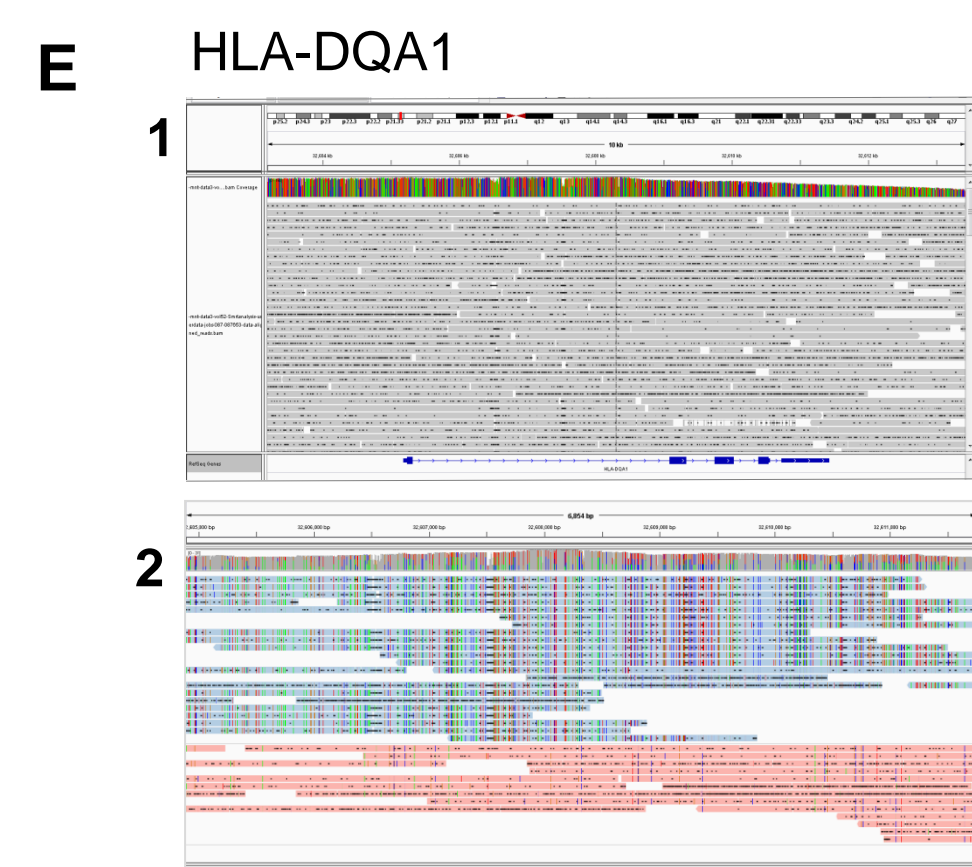
- The MHC locus is difficult to extract from traditional human genome assemblies
- Inspected contigs corresponding to MHC in PacBio *de novo* human genome assembly, (85x coverage, FALCON assembler)⁵
- The 4 Mb MHC region of chromosome 6 was located within a single 9 Mb contig
- Haplotigs separated by diploid aware software 'Falcon', unzips and identifies 20 significant structural variations between the homologous MHC chromosomes



PHASING MHC GENES BY TARGETED SEQUENCING

- MHC regions from the above sample were interrogated with two orthogonal targeted sequencing approaches
- 8 HLA genes were amplified and sequenced using GENDX NGS-go⁶ workflow for PacBio⁶
- Roche NimbleGen's SeqCap EZ Enrichment strategy⁷ was used to capture and sequence MHC regions
- HLA amplicon and whole genome data were 100% concordant, with a single difference in homopolymeric region of one allele

Locus	Alleles 1	Alleles 2	% cDNA Identity to IMGT-HLA Database	Gene Length Allele 1 / Allele 2	% Base Identity with Whole-Genome Assembly
A	HLA-A*01:01:01:01	HLA-A*02:01:01:01	100%	3138 / 3152	100%*
B	HLA-B*13:02:01	HLA-B*37:01:01	100%	3394 / 3395	100%
C	HLA-C*06:02:01:01	Homozygous	100%	3424 / -	100%
DRB1	HLA-DRB1*07:01:01:01	HLA-DRB1*10:01:01:01	100%	3904 / 3654	100%
DQB1	HLA-DQB1*02:02:01	Homozygous	100%	4114 / -	100%
DPA1	HLA-DPA1*01:03:01:02	Homozygous	100%	4857 / -	100%
DPB1	HLA-DPB1*04:01:01	Homozygous	100%	5867 / -	100%
DQA1	HLA-DQA1*01:05	HLA-DQA1*02:01	100%	5746 / 5657	100%



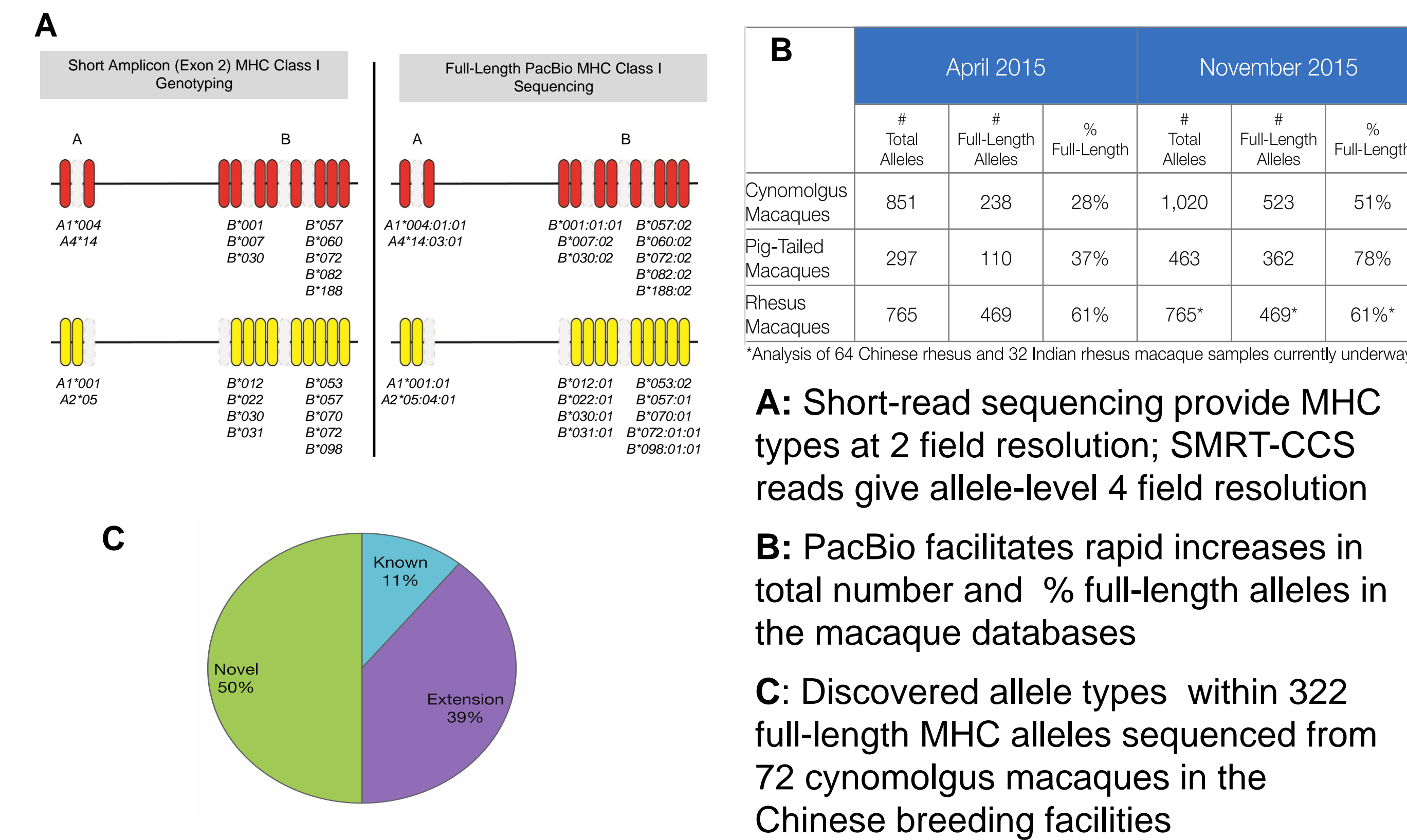
E1: Nimblegen SeqCap EZ MHC target-enriched PacBio reads aligned to HLA-DQA1 (hg19).

E2: PacBio sequences with predicted accuracy >97% phased using SAMtools (blue/pink). Consensus sequence captures larger region than PCR amplicon and useful for validating haplotype structures

Non-Human MHC Sequencing

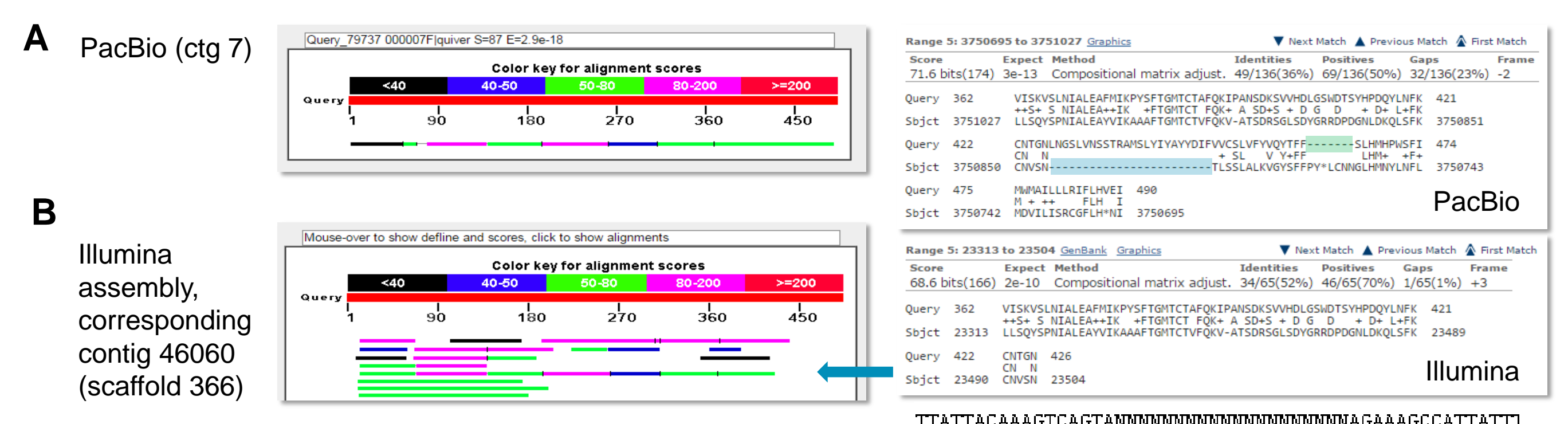
MHC ALLELE-DISCOVERY IN MACAQUES

- Macaque monkey MHC Class I regions are important for infectious disease research and more complex than humans due to gene duplications and variable number of pseudogenes
- Differing number of alleles express per haplotype, but only a few impact disease susceptibility
- SMRT-CCS PacBio reads provide full-length MHC class I cDNA sequences, which is a clear advantage for establishing robust MHC databases for Macaques



RESOLUTION OF AVIAN MHC REGIONS FROM WHOLE-GENOME ASSEMBLIES

- Avian immune regions are also important for understanding their disease resistance mechanism and can ultimately support avian conservation⁸
- Short-read assemblies have limitations for MHC analysis
- PacBio whole genome assembly of Anna's hummingbird using FALCON assembly method resolved several gaps in the short-read assembly



Summary

- SMRT Sequencing provides imputation-free access to complex immune regions.
- Complete views of genotypes and haplotypes of immune regions captured in human and non-human species.
- De novo* allows for reference-free analysis without imputation.
- These methods can be adopted to a multitude of immunology applications, including:

- Establishing databases of gene alleles and/or haplotypes
- Disease association in biomedical research
- Conservation biology

The authors acknowledge and thank everyone involved with the data generation for this poster. We also thank Martin Mayers & Cynthia Vierra-Green, National Bone Marrow Donor Program, Minneapolis, MN, for the 8 cell line samples used in the KIR haplotype sequencing project as well as for guidance provided during the project