

Building the Reference Assembly

Deanna M. Church
Senior Director of Genomics and Content
Personalis, Inc



 @deannachurch

Short Course in Medical Genetics 2014

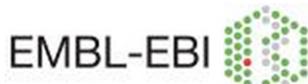
The Genome Reference Consortium consists of:



The Wellcome Trust Sanger Institute



The Genome Institute at Washington University



The European Bioinformatics Institute



The National Center for Biotechnology Information

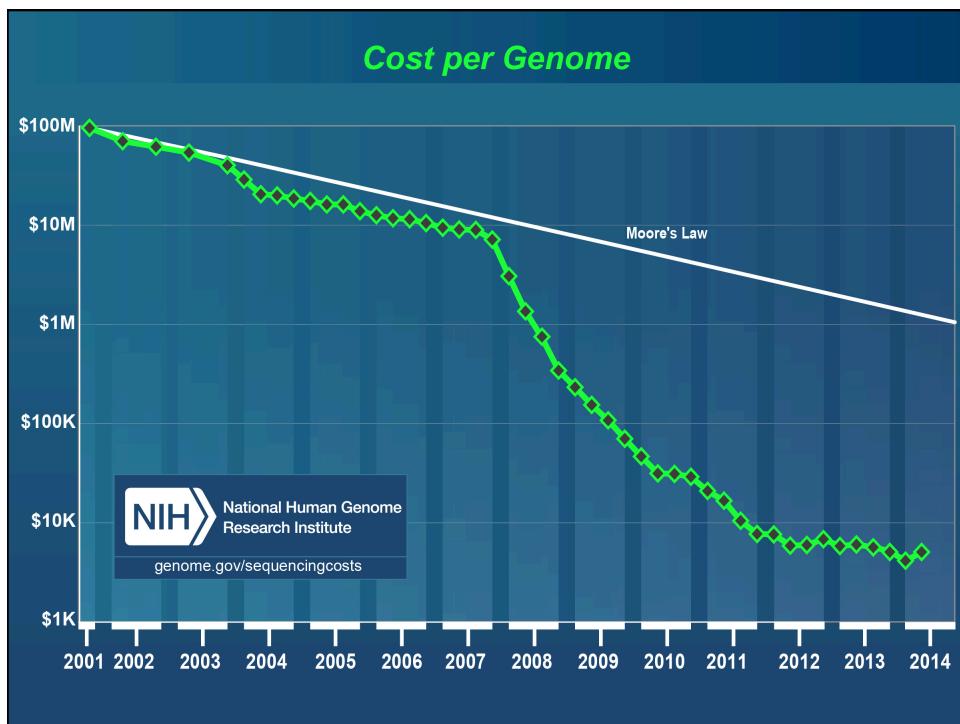
Valerie Schneider, NCBI

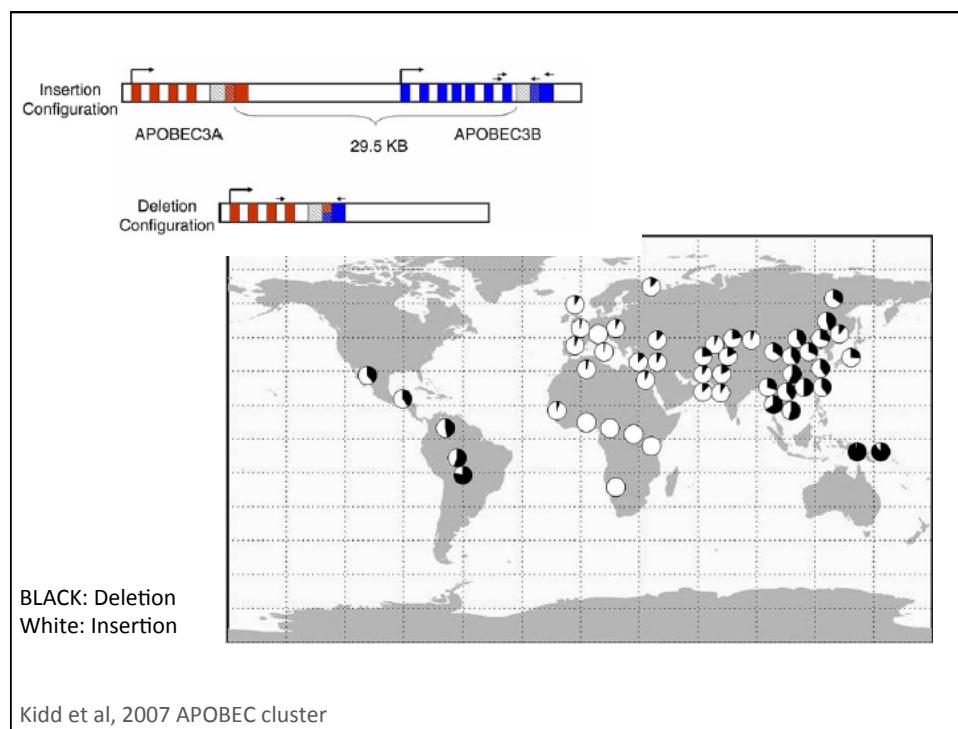
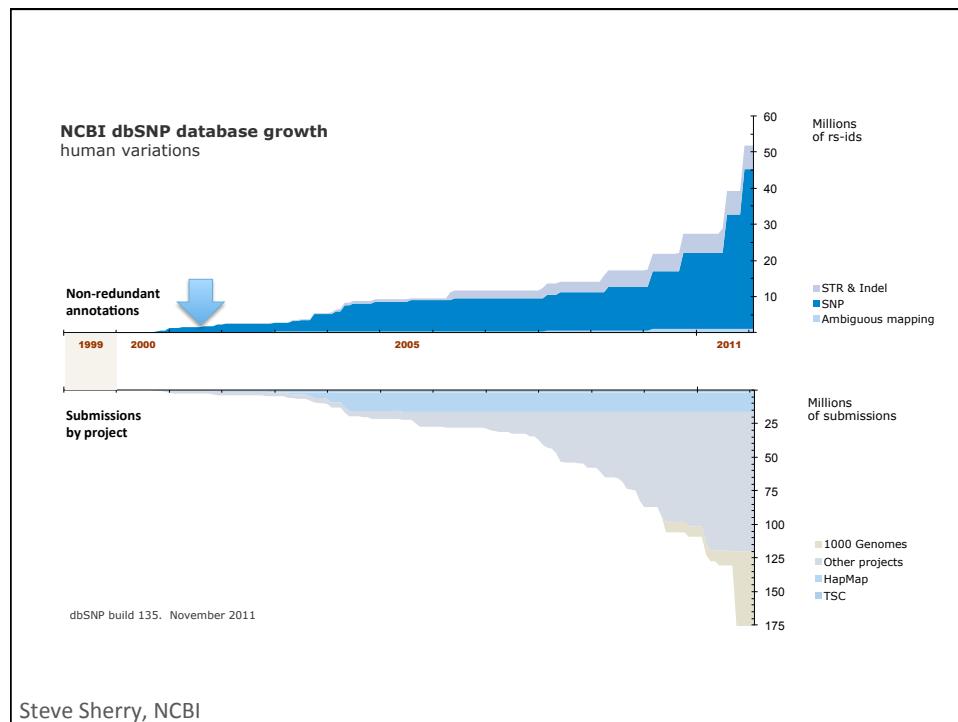
<http://genomereference.org>

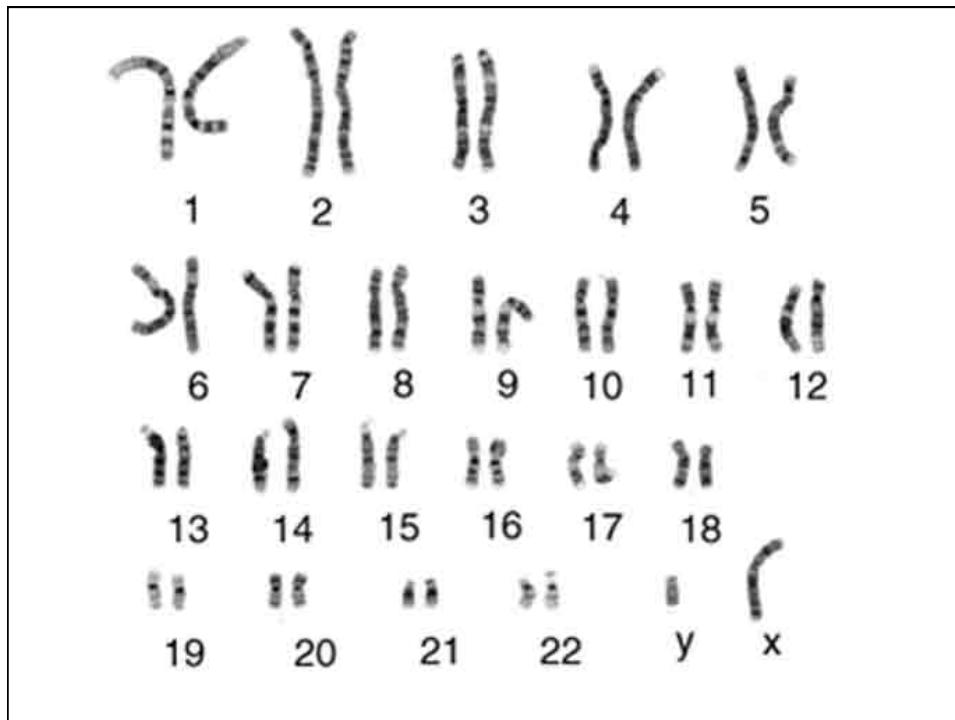
HGP Goals

Area	Goals 1993-98	Status as of Oct. 1998	Goals 1998-2003
Genetic map	Average 2- to 5-cM resolution	1 cM map published Sept. 1994	Completed
Throughput: 500 Mb/year Cost: < \$0.25 per base Variation: 100,000 SNPs mapped			
Gene identification	Develop technology	30,000 ESTs mapped	Full-length cDNAs
Functional analysis	Not a goal	-	Develop genomic-scale technologies
Model organisms	<i>E. coli</i> : complete sequence <i>C. elegans</i> : most of sequence <i>Drosophila</i> : begin sequencing Mouse: map 10,000 STSs	Published Sept. 1997 Released Apr. 1996 80% complete 9% done 12,000 STSs mapped	- Complete Dec. 1998 Sequence by 2002 Develop extensive genomic resources Lay basis for finishing sequence by 2005 Produce working draft before 2005

Collins FS et al, 1998







Reference assembly history

PERSPECTIVE

Against a Whole-Genome Shotgun

Philip Green¹

Department of Molecular Biotechnology, University of Washington, Seattle, Washington 98195

PERSPECTIVE

Human Whole-Genome Shotgun Sequencing

James L. Weber^{1,3} and Eugene W. Myers²

¹Center for Medical Genetics, Marshfield Medical Research Foundation, Marshfield, Wisconsin 54449;

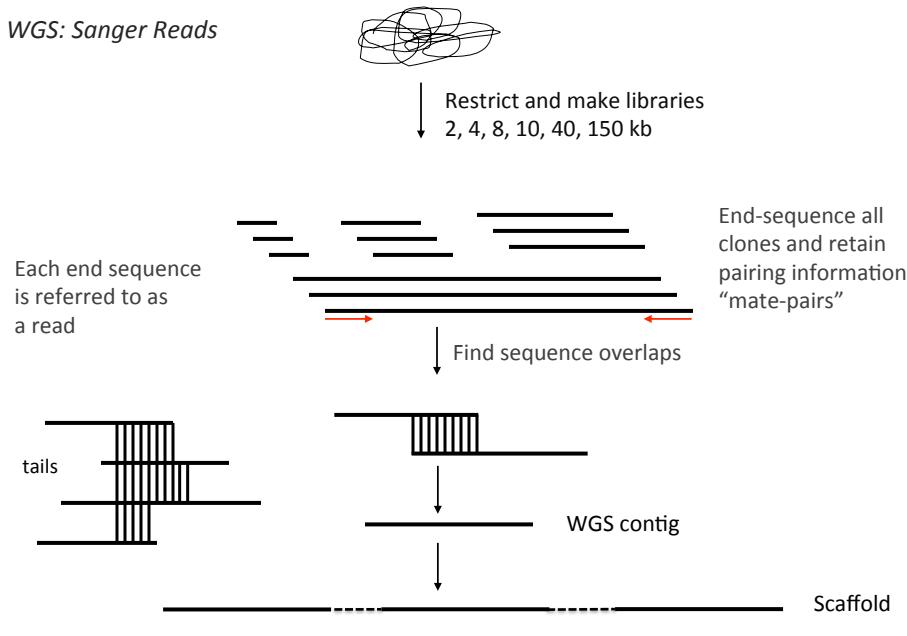
²Department of Computer Science, University of Arizona, Tucson, Arizona 85721

Genome Research, May, 1997

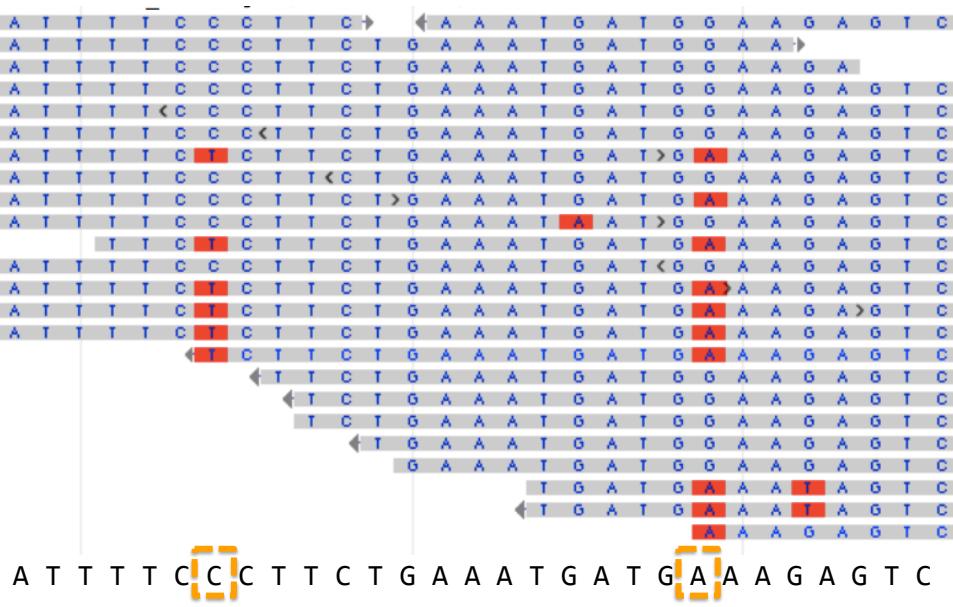


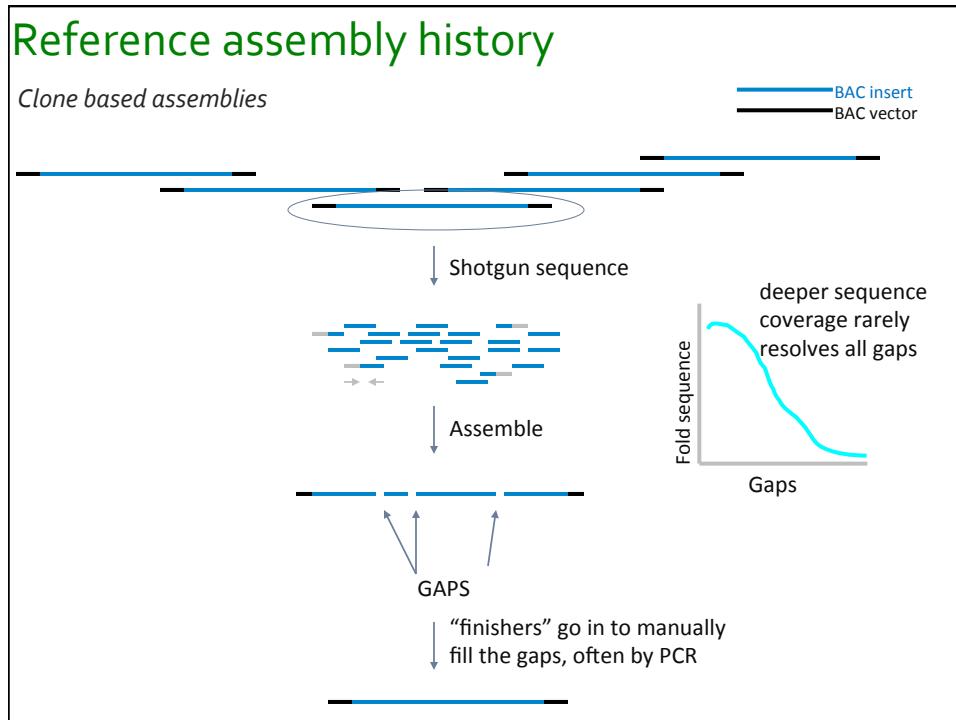
Reference assembly history

WGS: Sanger Reads



Reference assembly history

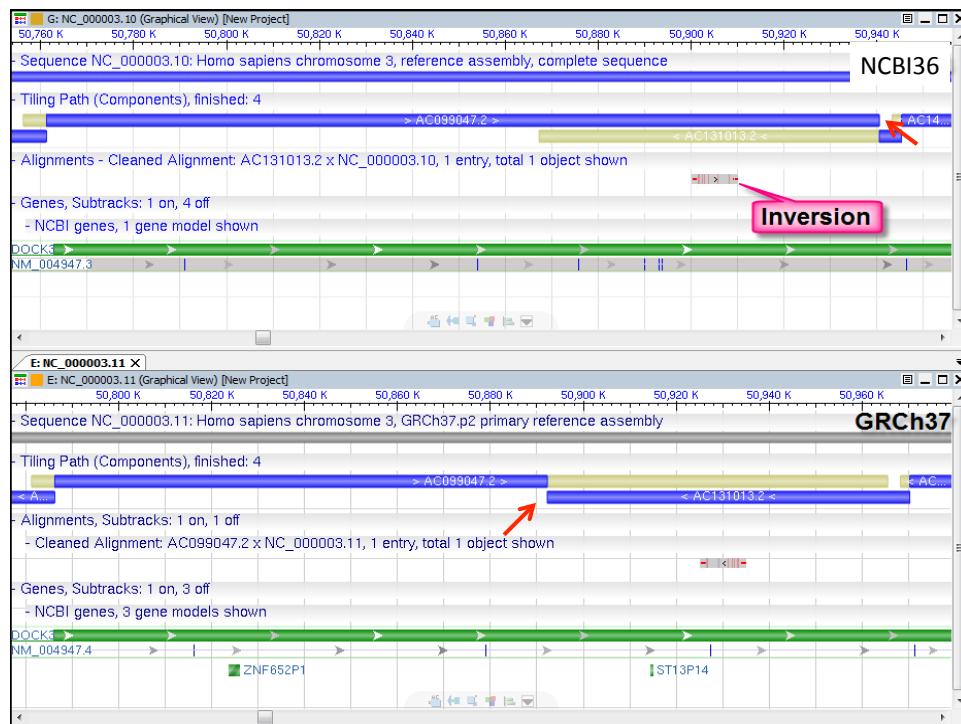
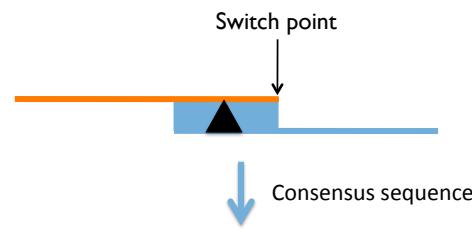




Reference assembly history

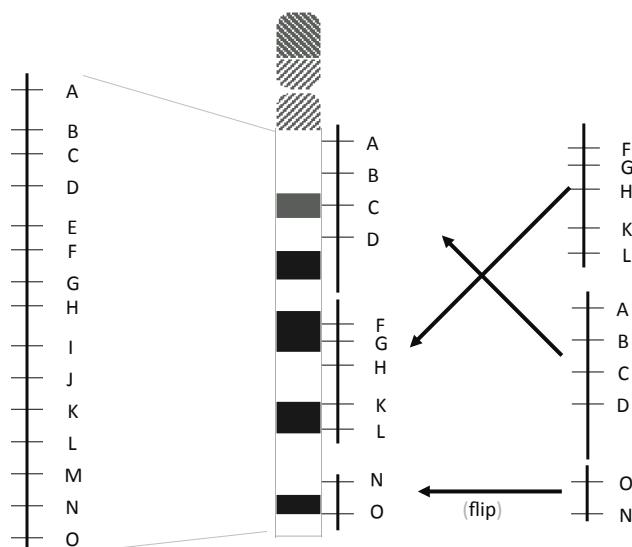
Build sequence contigs based on contigs defined in TPF.

- ✗ Check for orientation consistencies
- ✗ Select switch points
- ✗ Instantiate sequence for further analysis



Reference assembly history

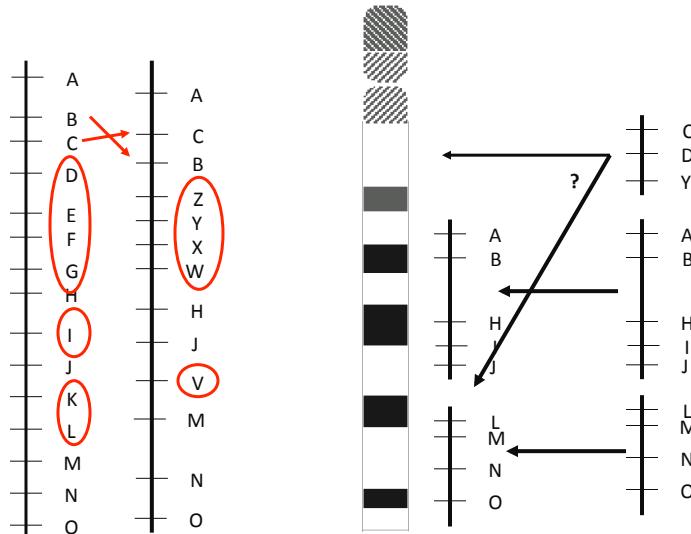
Ideally...

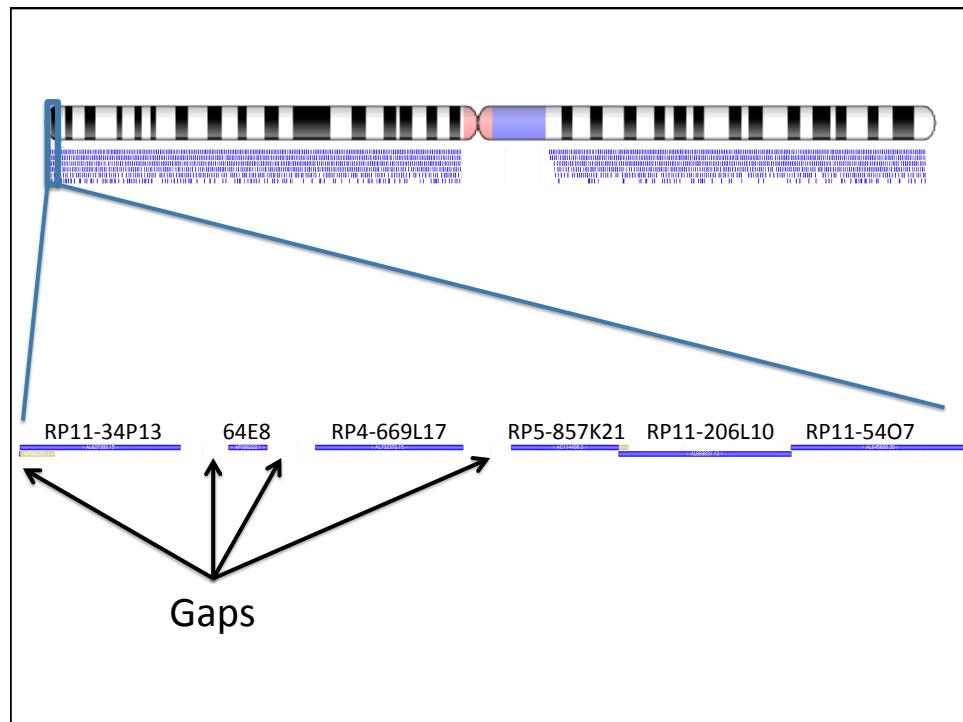
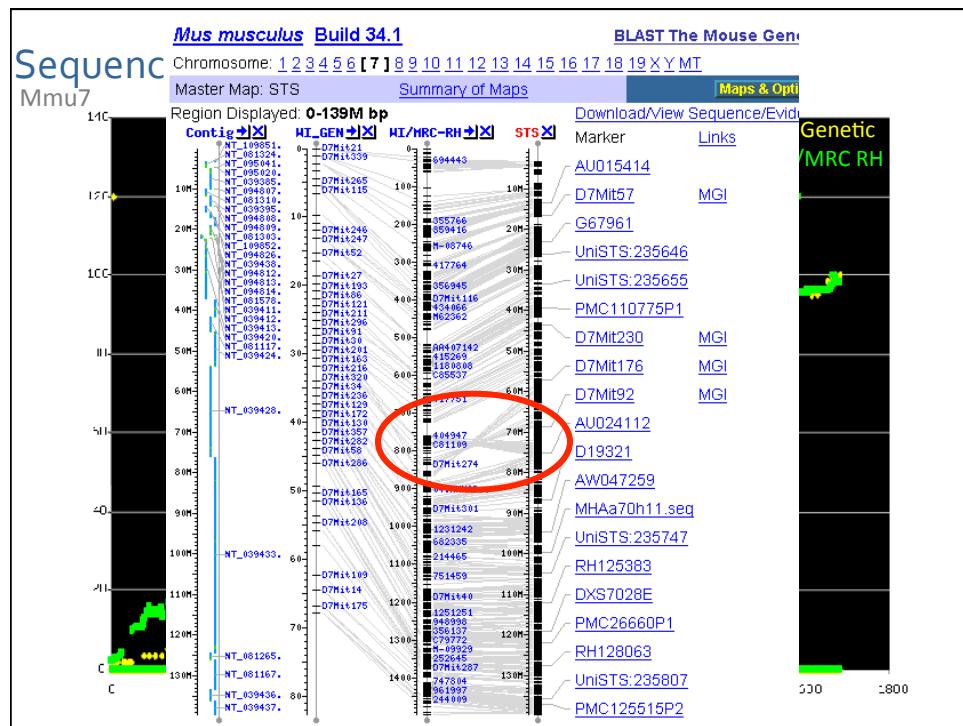


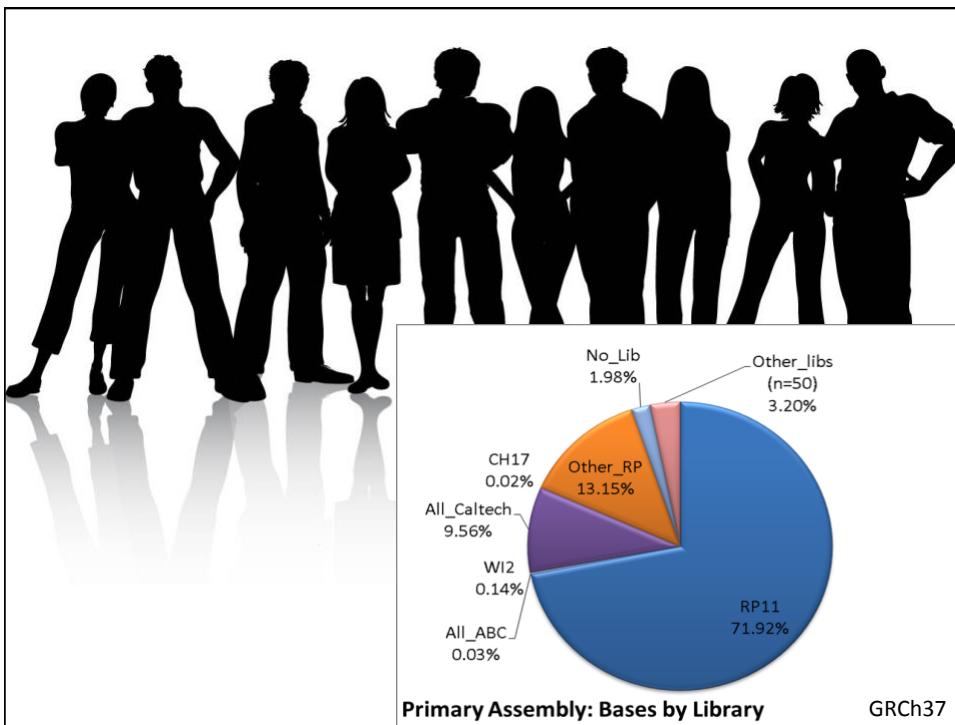
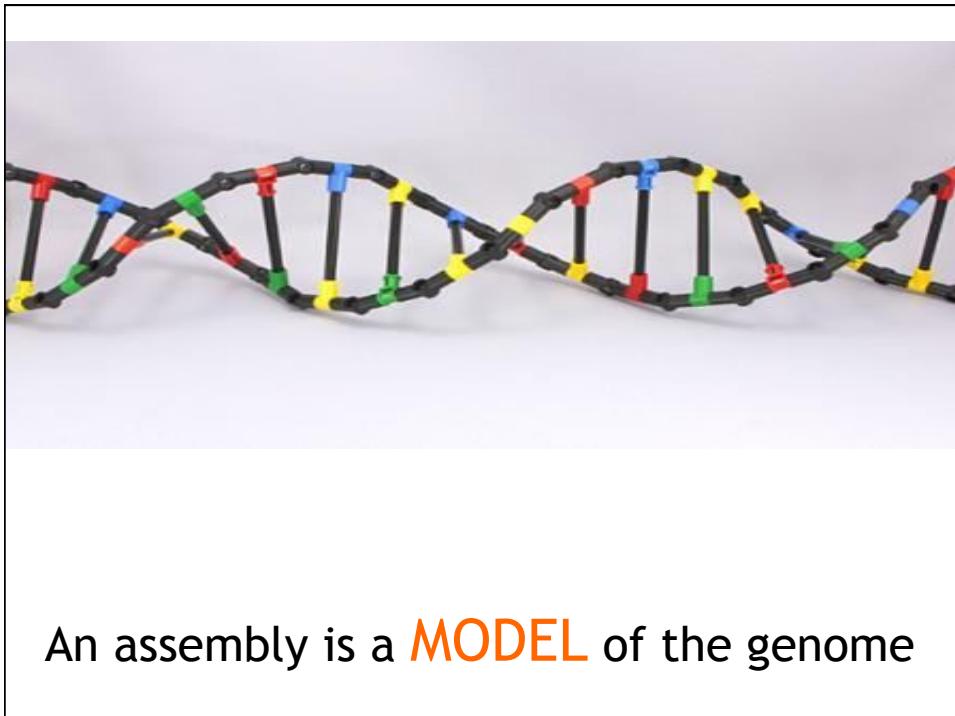
Non-sequence based Map

Reference assembly history

More like...







GRCh37 Chr. 7

NC_000007.13: 99M..99M (51b+)

GRCh37 genome-wide recombination rate from Phase 2 HapMap estimated from phased haplotypes

Association Results

NCBI Genes

CYP3A5

GRCh37 allele (C) is non-functional

Issue Report for HG-1321

Category: Variation

Affects GRCh37

version(s):

Report type: RefSeq Report

Description: The reference may be representing the non-protein-coding allele of GenID:1577 (CYP3A5)

Last updated: 2012-12-03

Status: Open

Experiment type: na

http://www.ncbi.nlm.nih.gov/projects/genome/assembly/grc/issue_detail.cgi?id=HG-1321

GRCh37 chr. 9

NC_000009.11: 95M..95M (78b+)

RefSeq Alignments

ASPN

GRCh37 D14 allele (3nt insertion) associated with osteoarthritis

Issue Report for HG-1012

Category: Variation

Affects GRCh37

version(s):

Report type: RefSeq Report

Description: The reference genome is representing the D14 (less common) allele of NM_017680.4

Last updated: 2013-03-01

Status: Awaiting External Info

Experiment type: na

Assembly Information

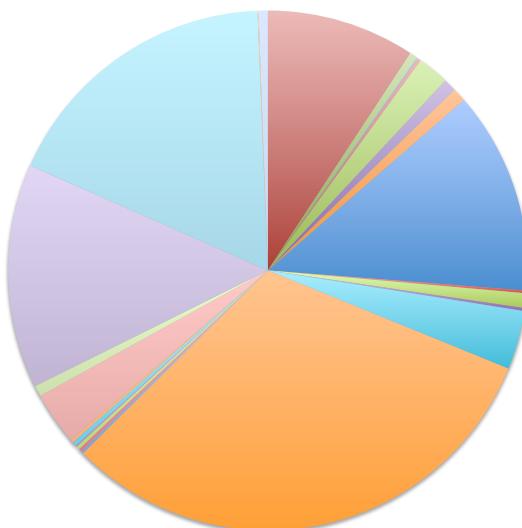
Select a placement below to display it in the Sequence Viewer.

GRCh37.p12 chr9:95,118,957-95,271,297 (View Regions: Ensembl | NCBI | UCSC)

NCBI36 chr9:94,158,778-94,311,118 (View Regions: Ensembl | NCBI | UCSC)

http://www.ncbi.nlm.nih.gov/projects/genome/assembly/grc/issue_detail.cgi?id=HG-1012

Center sequence distribution: NCBI36



AECOM BCM Beijing CGM CHGC CMGWCH CSHL GBF GS GTC IIGB-CNR
IMB JGI JST Keio MPIIMG RIKEN SC SDSTDC SHGC TIGR Tokai
UOKNOR UTSW UUGC UWGC UWMSC WIBR WUGSC YMGC unknown

Finding the data

ends of each insert. Centres differed in the fluorescent labels employed and in the degree to which they used dye-primers or dye-terminators. The sequence detectors included both slab gel- and capillary-based devices. Detailed protocols are available on the web sites of many of the individual centres (URLs can be found at http://www.nhgri.nih.gov/genome_hub.html). The extent of



404 - Not Found

The link you followed or typed is either incomplete, outdated, contains disallowed characters, or we've made a mistake!

- Please try our *Powered-by-Google* search engine above (your best bet).
- Navigate to what you want directly from the [Home page](#).
- Please [send us feedback](#) if you think there is a problem with our site.

Our apologies for any inconvenience!

GRC Genome Reference Consortium

Church et al., 2011 PLoS

GRC Home Data Help Report an Issue Contact Us Credits Curators Only

Human | Mouse | Zebrafish

The Genome Reference Consortium
Putting sequences into a chromosome context.

The original model for representing the genome assemblies was to use a single, preferred tiling path to produce a single consensus representation of the genome. Subsequent analysis has shown that for most mammalian genomes a single tiling path is insufficient to represent a genome in regions with complex allelic diversity. The GRC is now working to create assemblies that better represent this diversity and provide more robust substrates for genome analysis.

The GRC has started the submission of GRCz10 to GenBank. We will provide an update on this website when the submission is complete. If you have questions or concerns about this [let us know](#).

Attending Genome Informatics 2014? Register now for the GRC Assembly Workshop!

Transitioning to GRCh38? Try the [NCBI Remapping Service](#), which uses the same assembly-assembly alignments used by the GRC.

Subscribe to the [grc-announce](#) email list to receive email notification for all GRC assembly updates.

GRC Blog

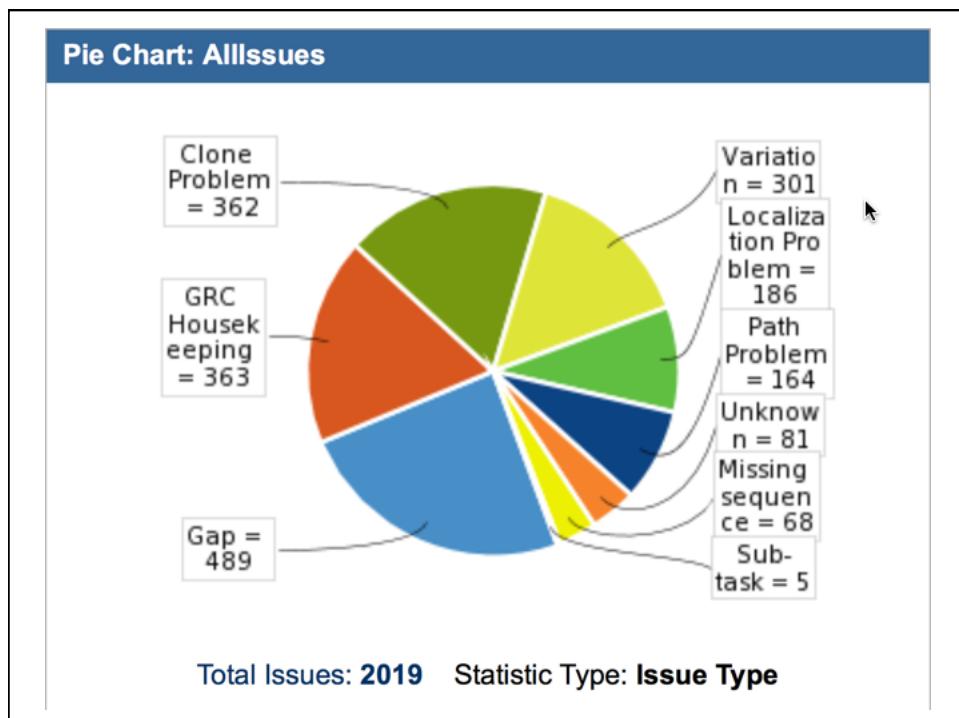
- Chromosome 9 peri-centromeric assembly improvement Apr 04, 2014
- GRCh38: Incorporating Modeled Centromere Sequence Jan 14, 2014
- Centromeres are specialized chromatin structure... [see all](#)

Resolved Issues

- Zebrafish (ZG-6653) Jun 25, 2014
- CH211-274B20 has been removed from the Chr25 TPF. It has been re-sequenced and matches to Chr25 but is deleted so
- Zebrafish (ZG-6652) Jul 2, 2014
- MGH markers localise contig ctg10380 to chromosome 4 in the gap tracked in JIRA ticket ZG-3934

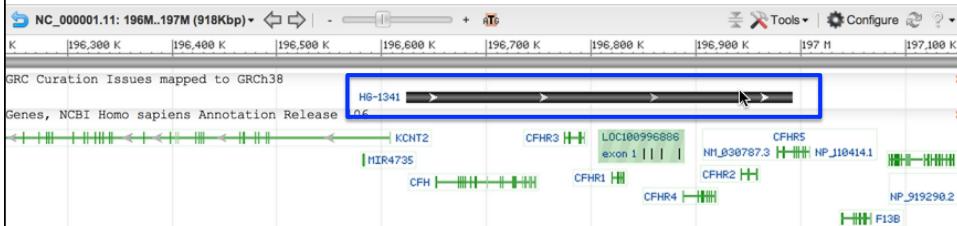
[see all](#)

<http://genomeref.org>

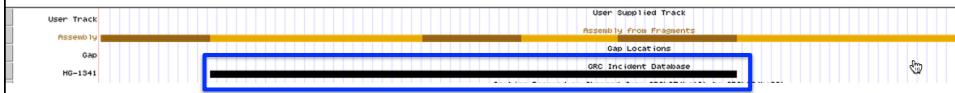


Finding Curated Regions: Browsers

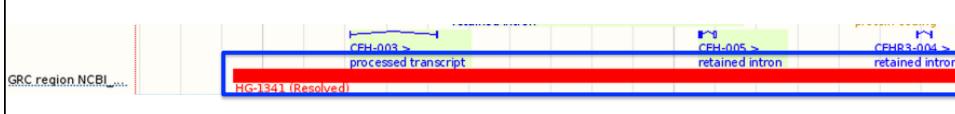
NCBI: GRC Curation Issues mapped to \$assembly_name



UCSC: GRC Incident Database



Ensembl: GRC_region_NCBI37



Finding Curated Regions: GRC web browser

Issue Report for HG-1341

Category: Unknown

Affects version(s): GRCh37

Report type: TPF Analysis

Description: Clones have been selected for sequencing that align near the AL049744.8 - AL139418.9 region of Chr. 1. These may be used to improve the reference or to represent an alternate haplotype.

Last updated: 2013-07-18

Status: Resolved

Experiment type: Clone Sequencing

Resolution: This CH17 pathway does not meet the criteria to replace the reference or for alt_loci creation.

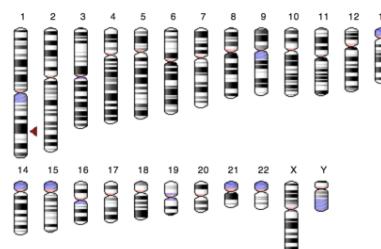
Fix version(s): GRCh38

Assembly Information

Select a placement below to display it in the Sequence Viewer.

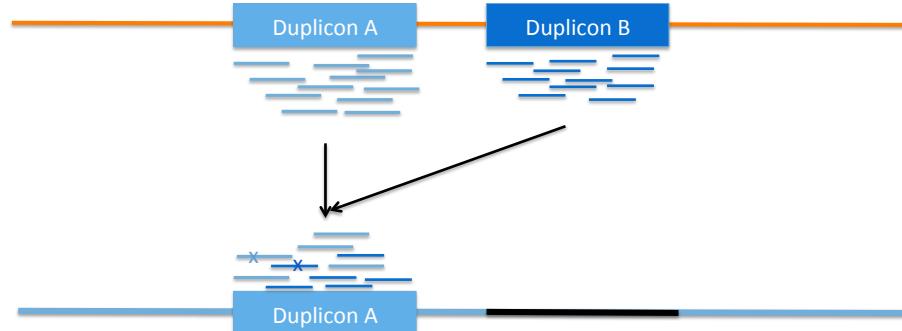
GRCh38 chr1:196,623,707-196,992,672 (View Regions: [NCBI](#) | [UCSC](#))

GRCh37.p13 chr1:196,592,837-196,961,802 (View Regions: [Ensembl](#) | [NCBI](#) | [UCSC](#))



Missing Sequence: Gaps

Sample genome

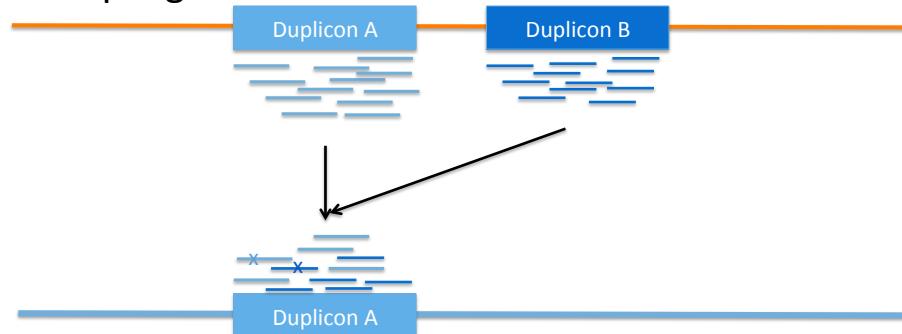


G>A (allelic difference – true variant)
G>C (paralogous sequence variant- false positive)

May or may not detect increased coverage depending on sequencing depth
and library quality (easier to find with new technologies than with old, low through technologies)

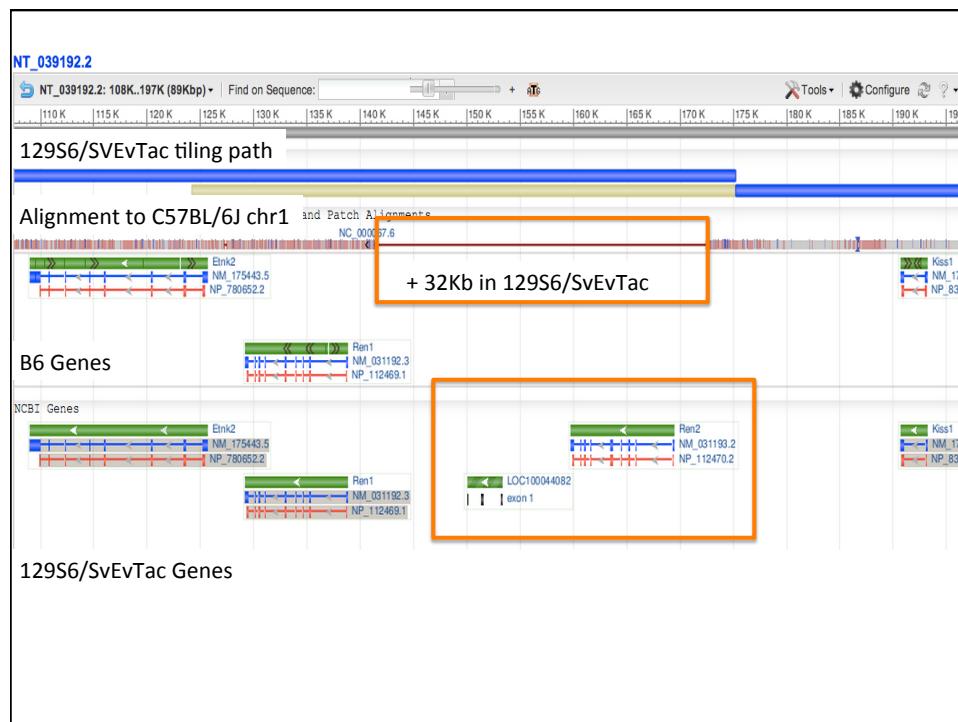
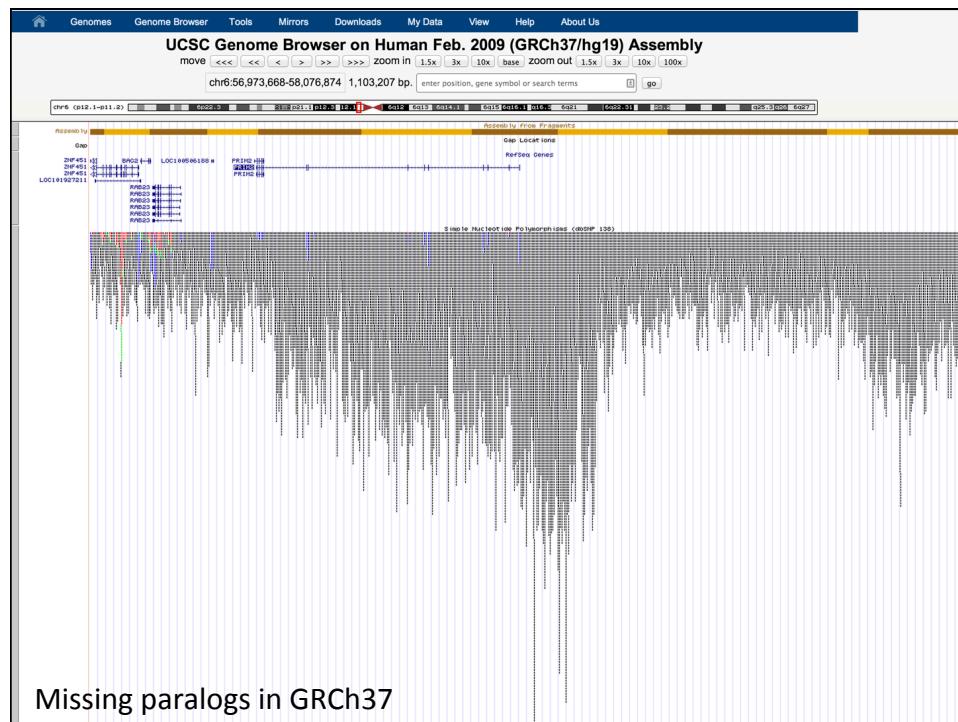
Missing sequence: collapse

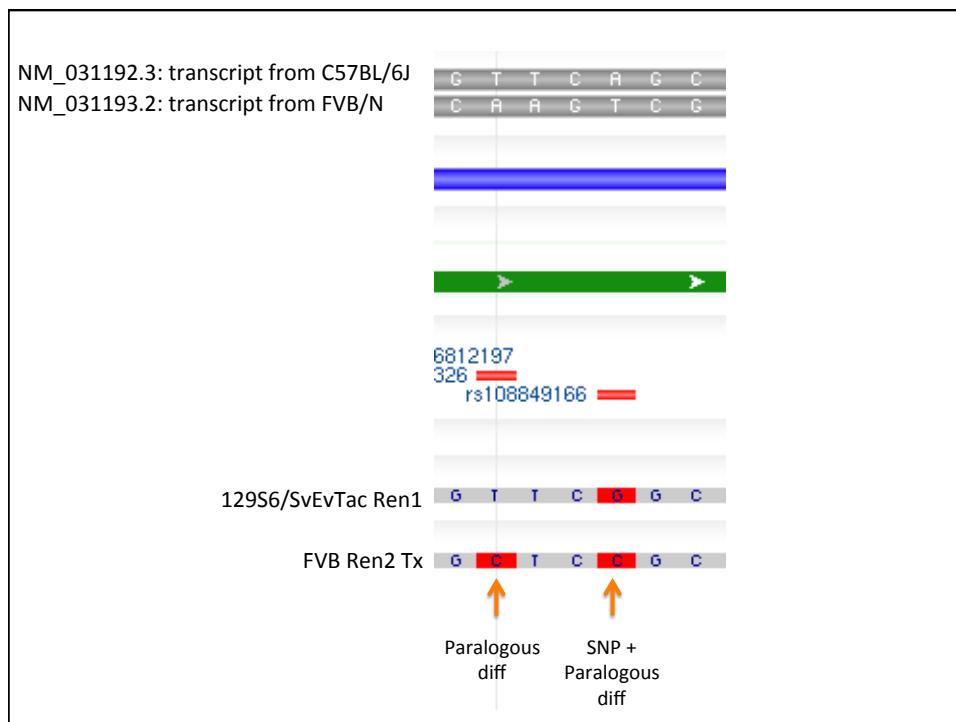
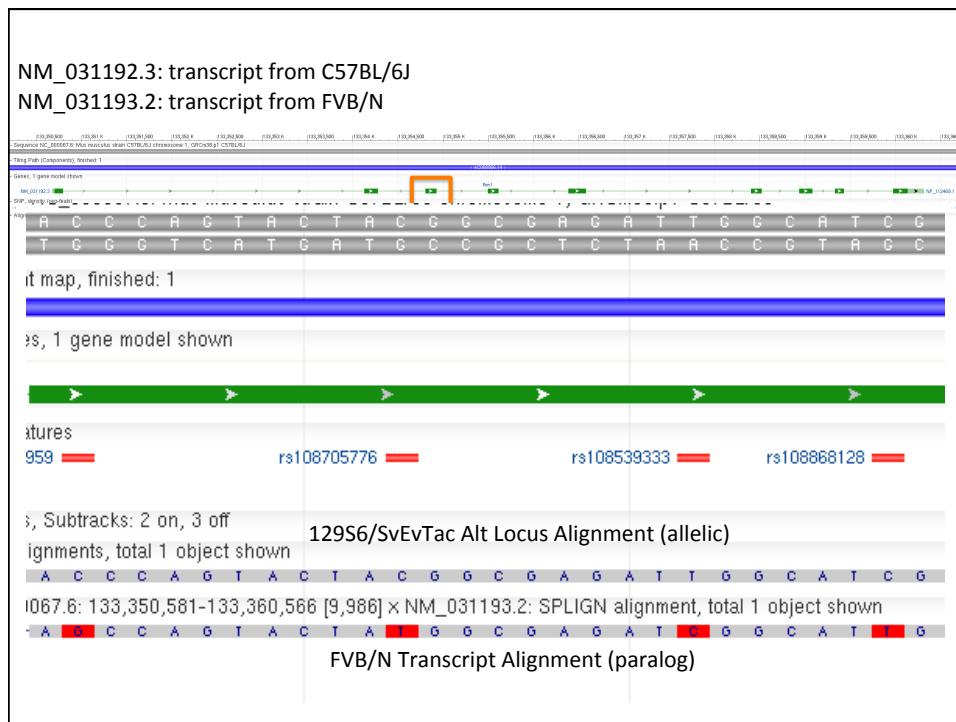
Sample genome

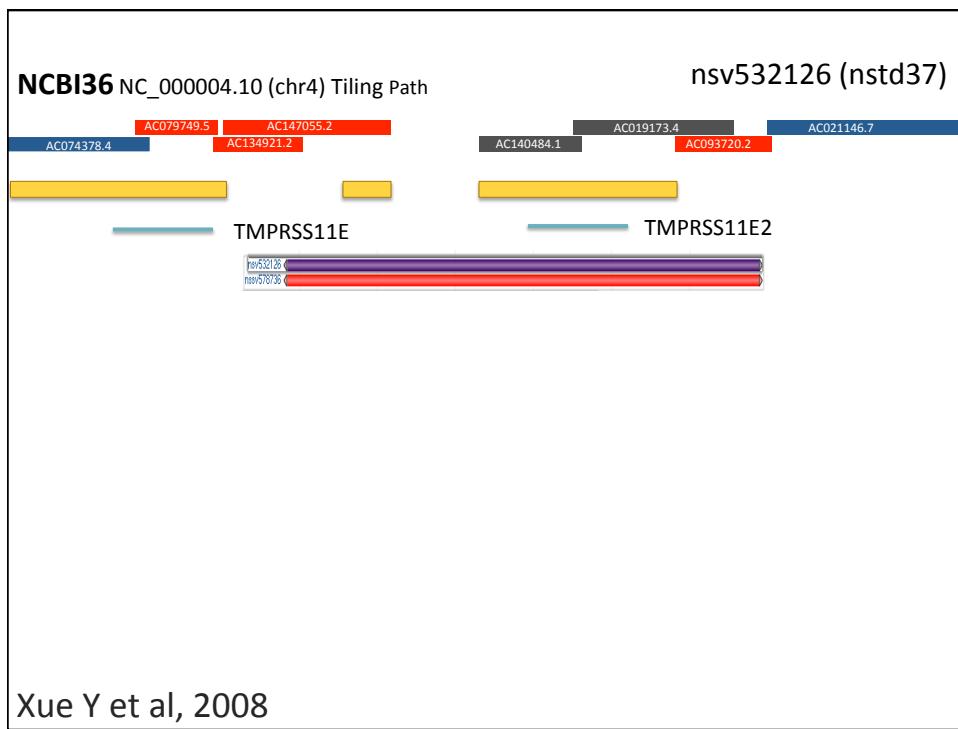
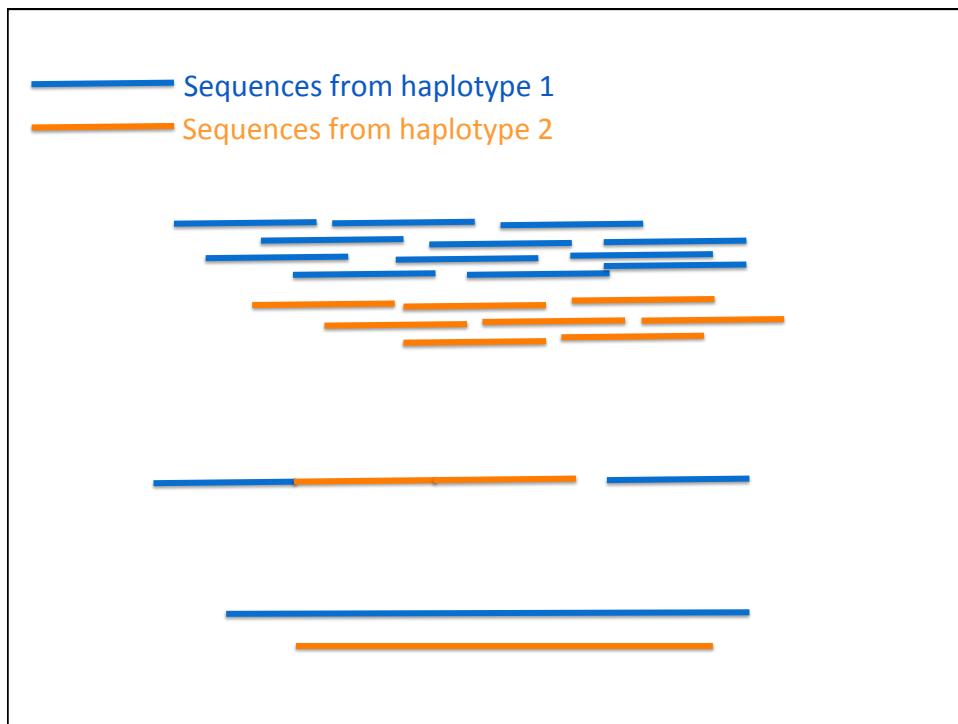


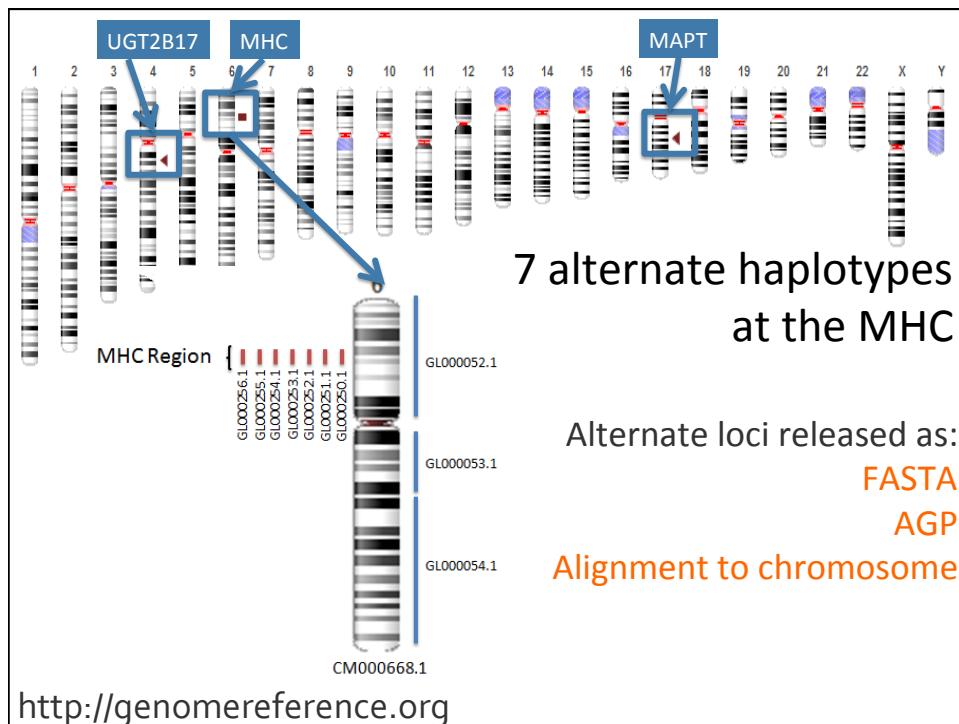
G>A (allelic difference – true variant)
G>C (paralogous sequence variant- false positive)

May or may not detect increased coverage depending on sequencing depth
and library quality (easier to find with new technologies than with old, low through technologies)





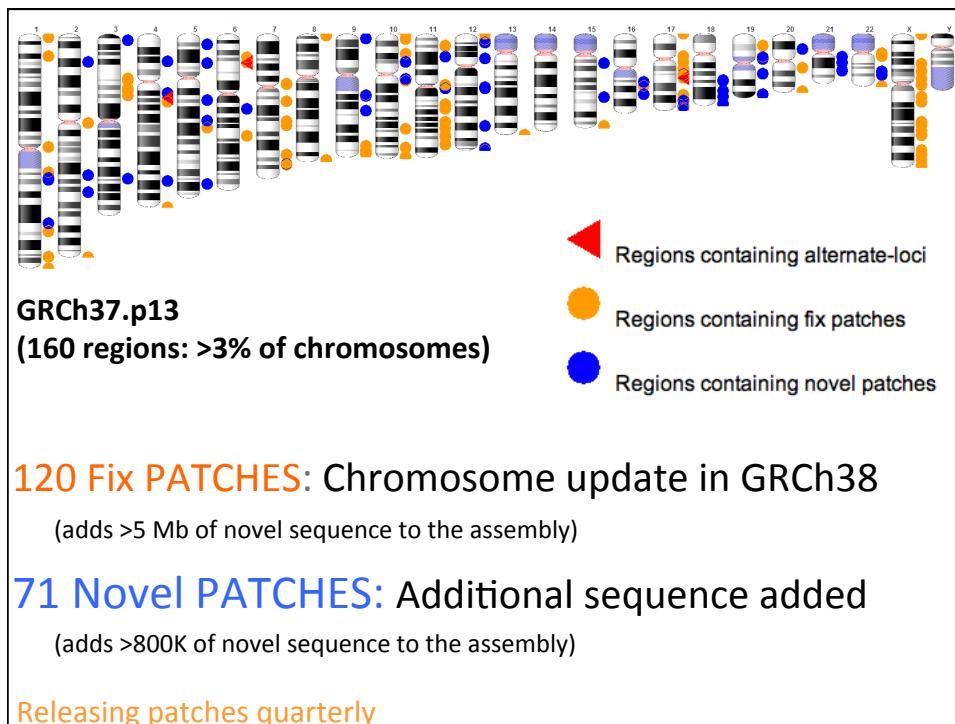




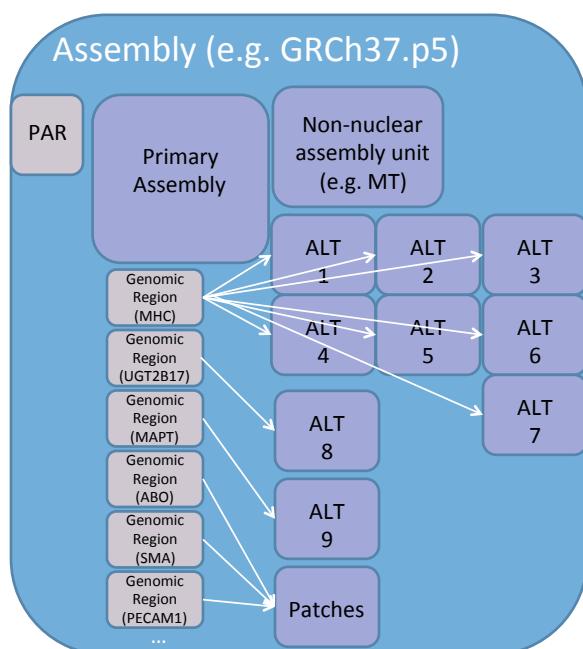
Updating the assembly

Oh No! Not a new version of the human genome!

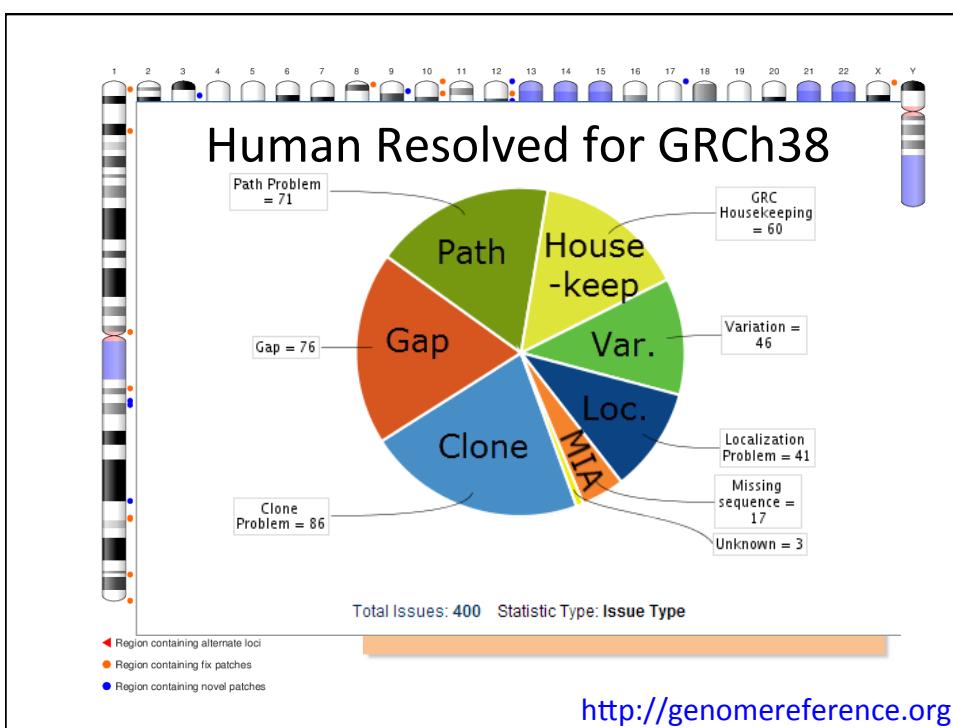
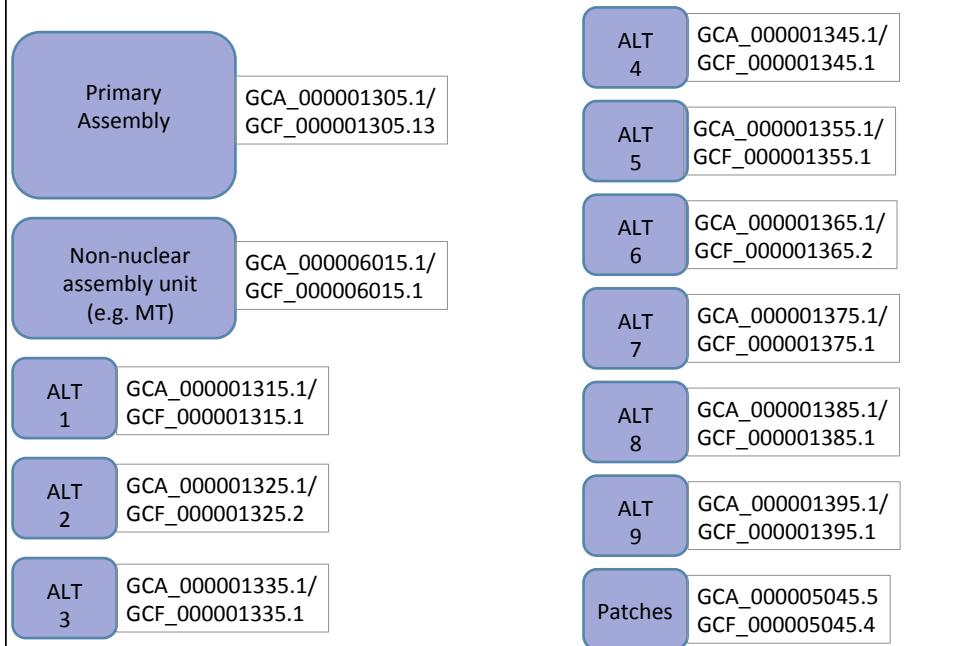


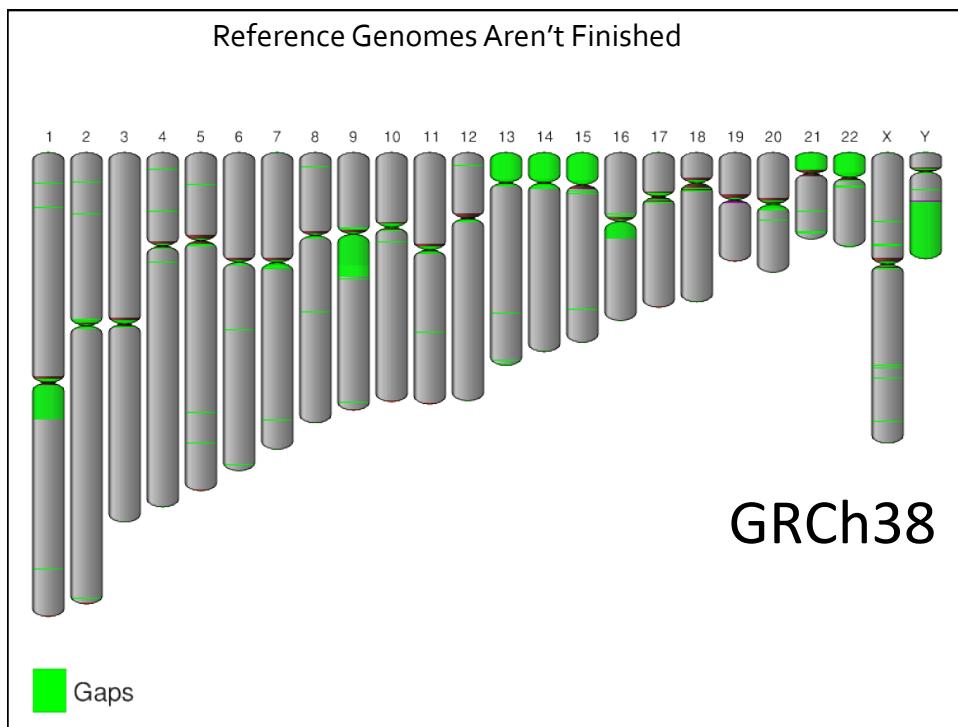


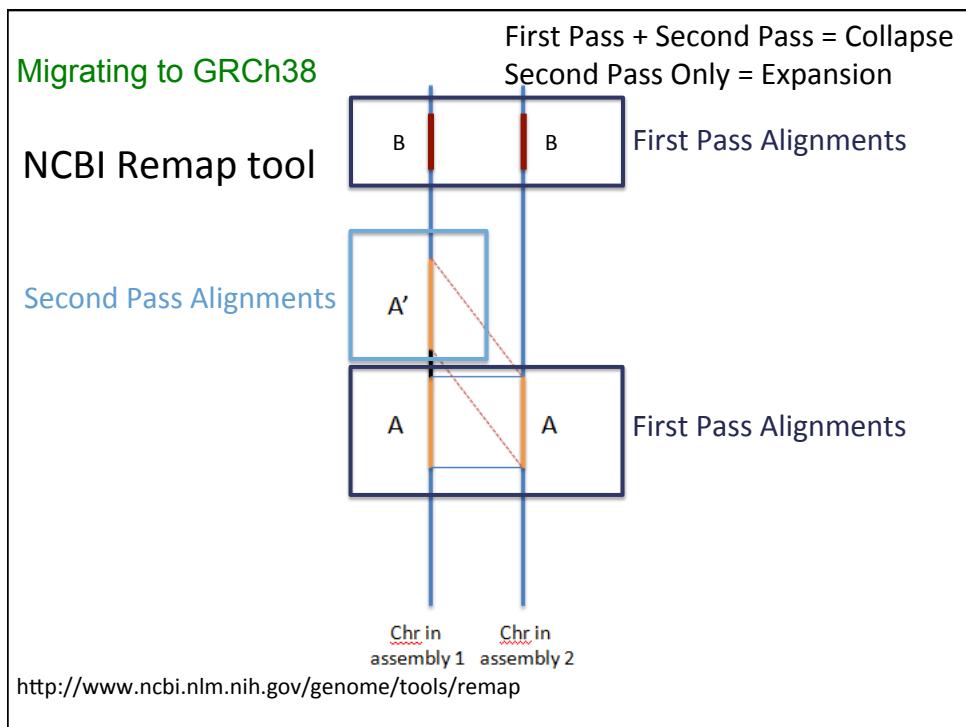
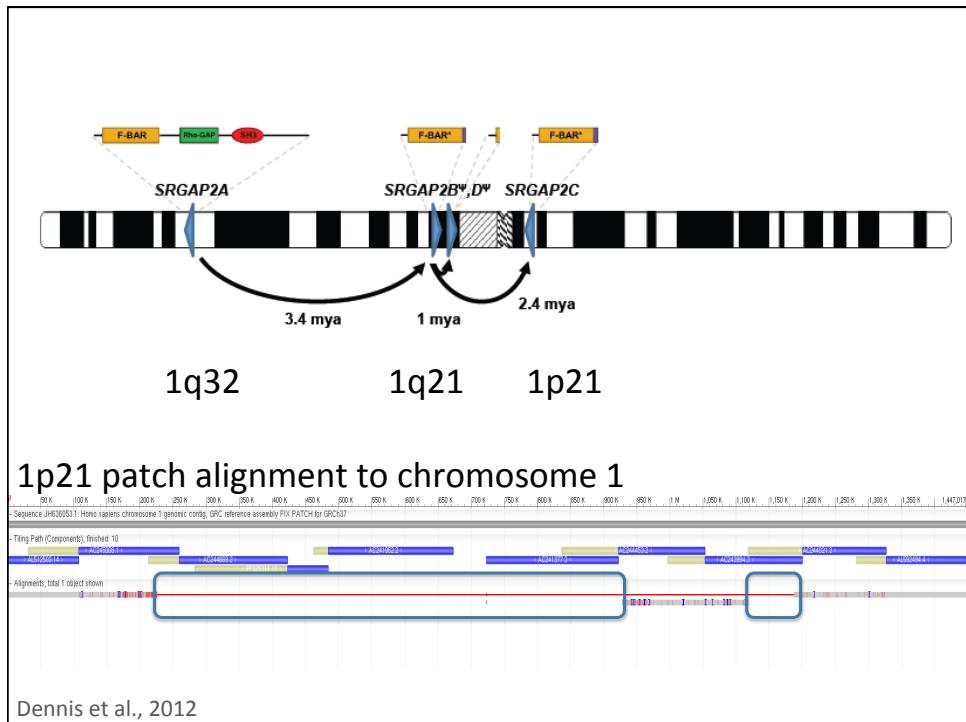
Data Model

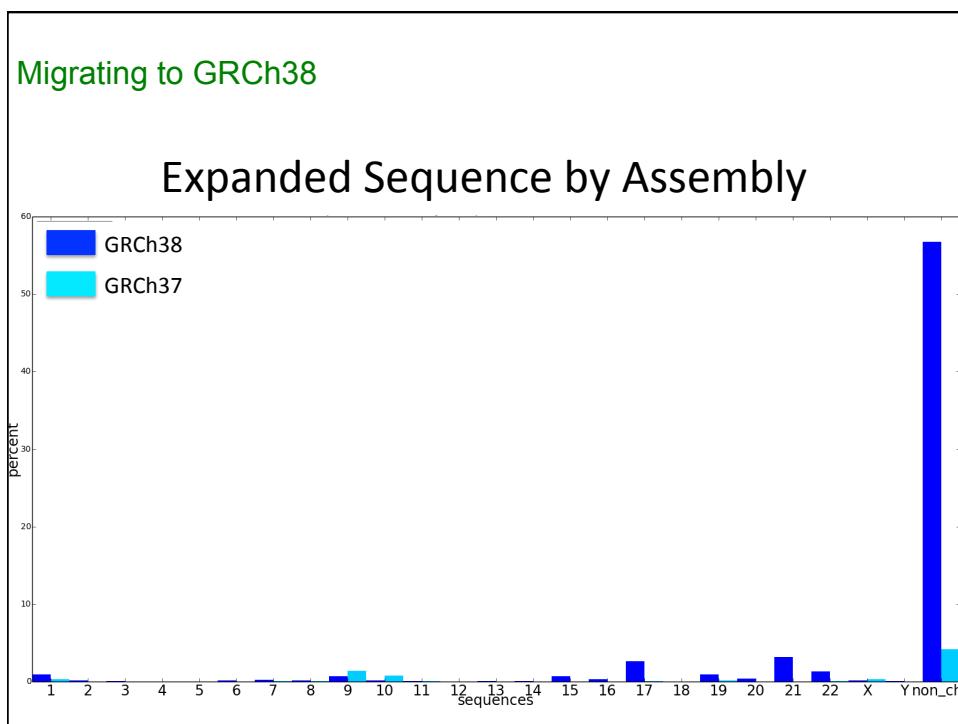
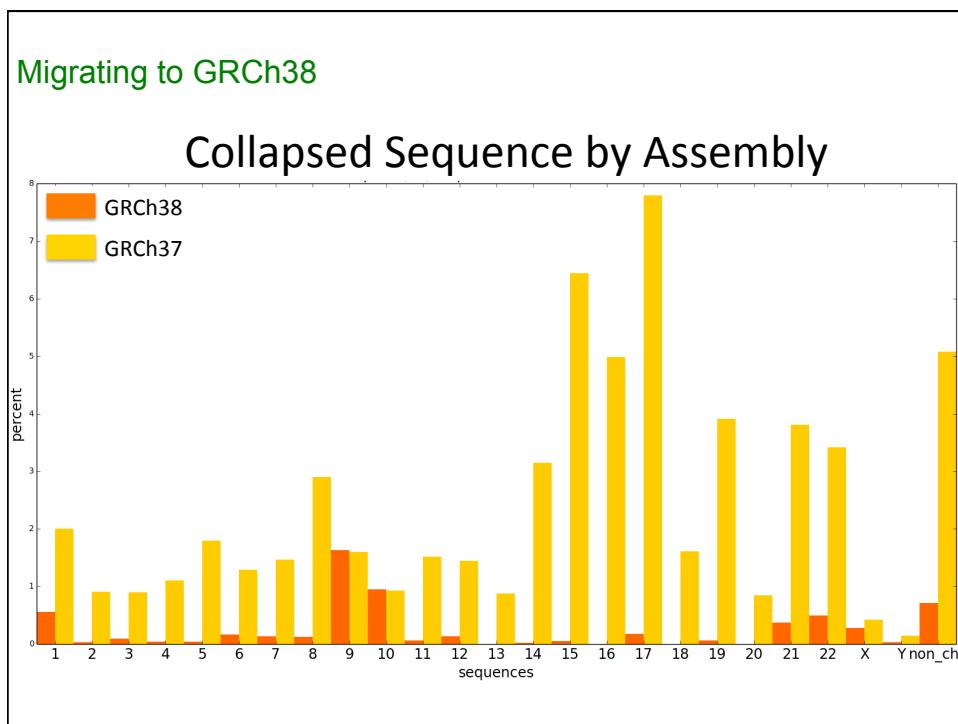


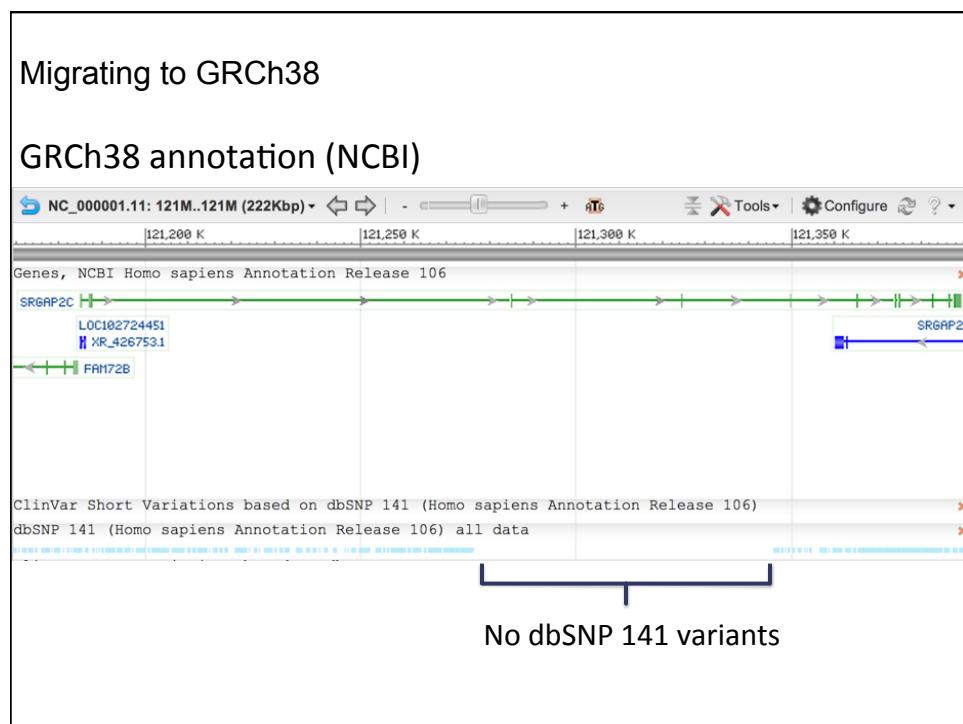
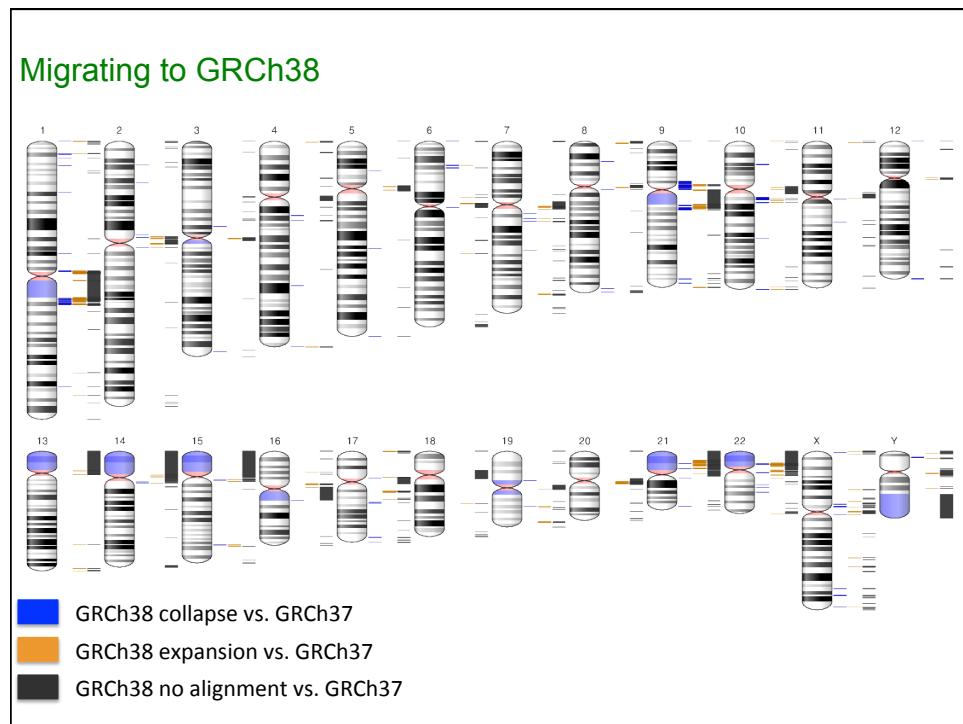
Data Model

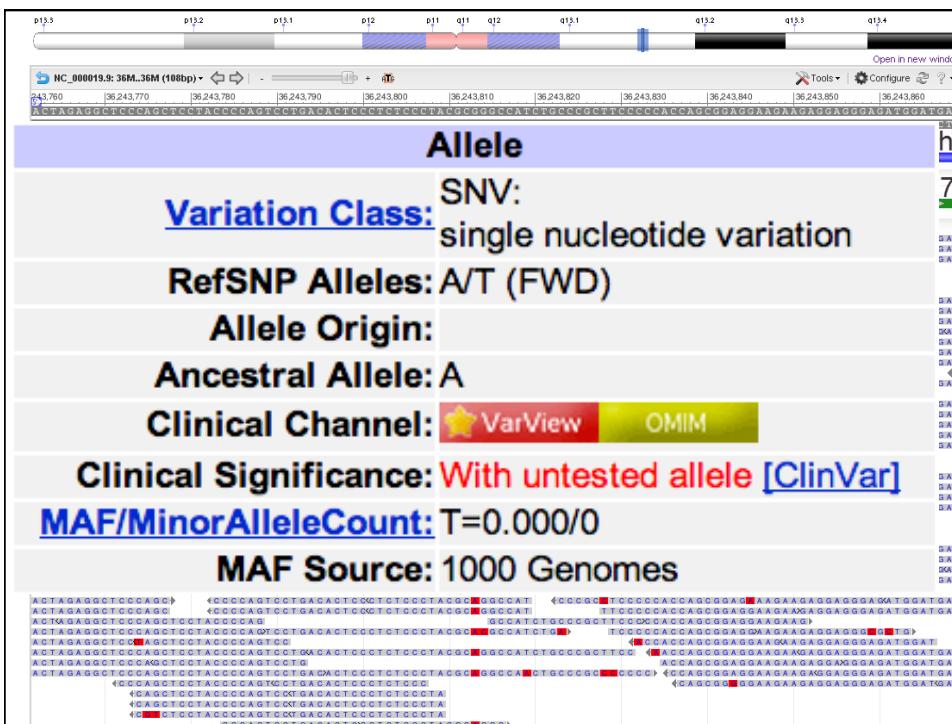
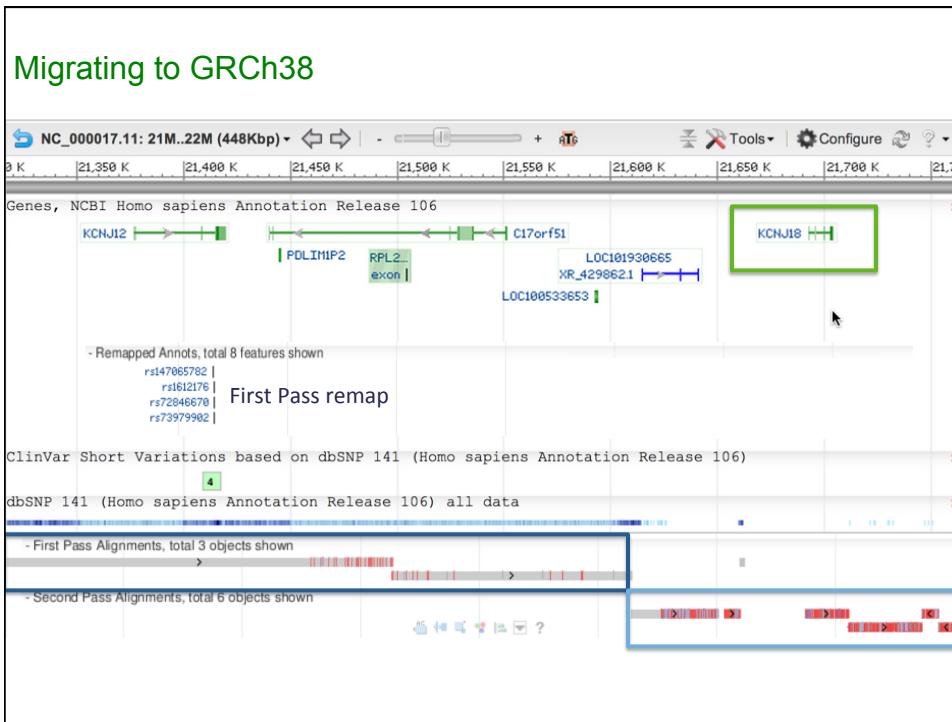


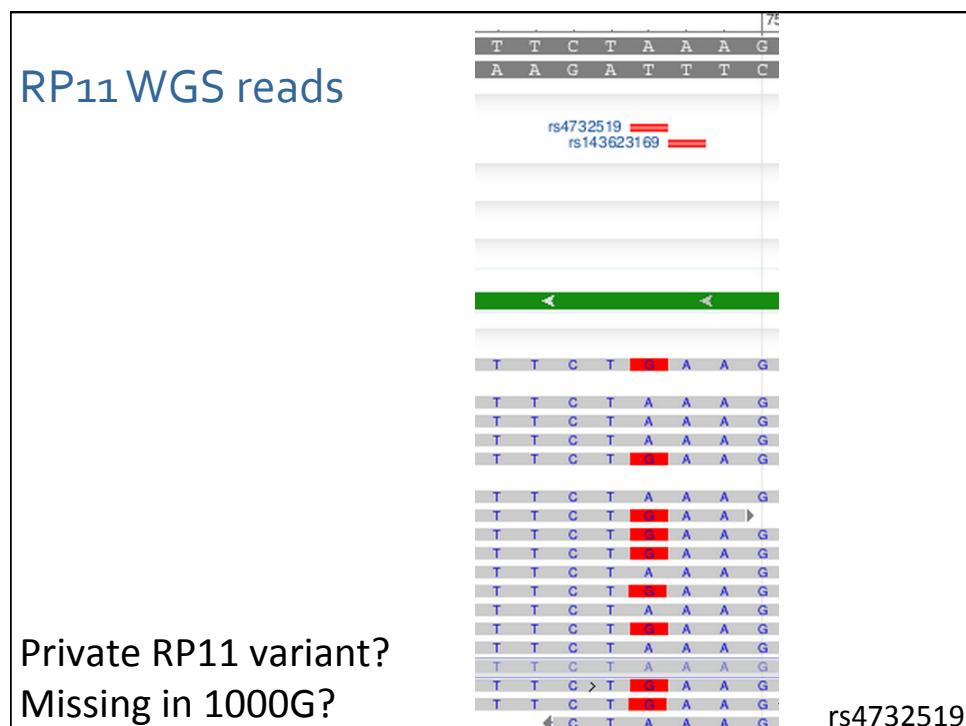
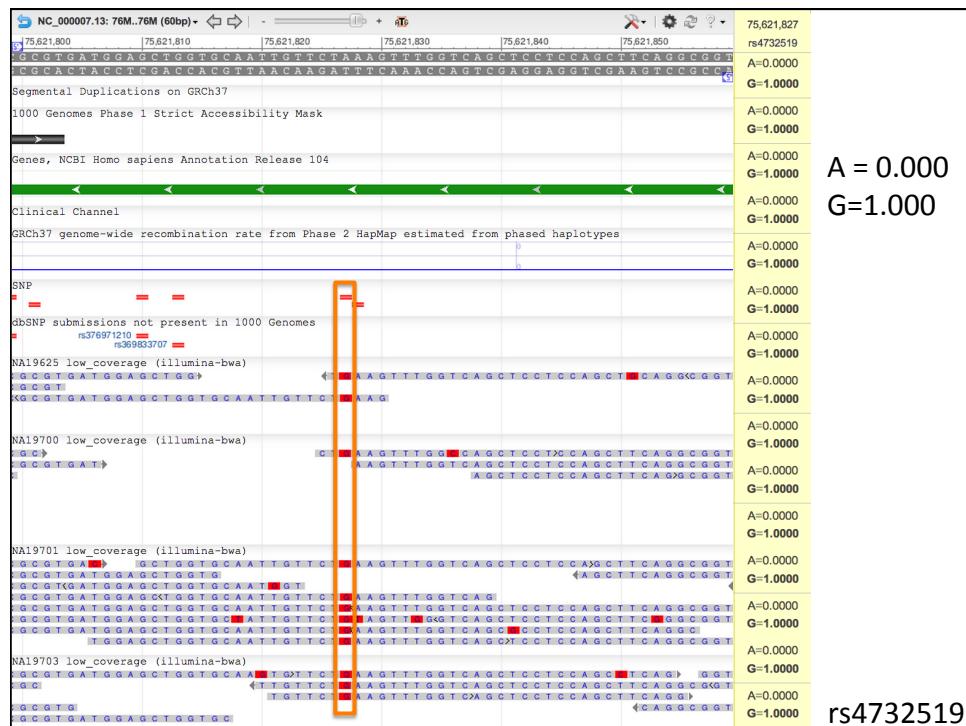


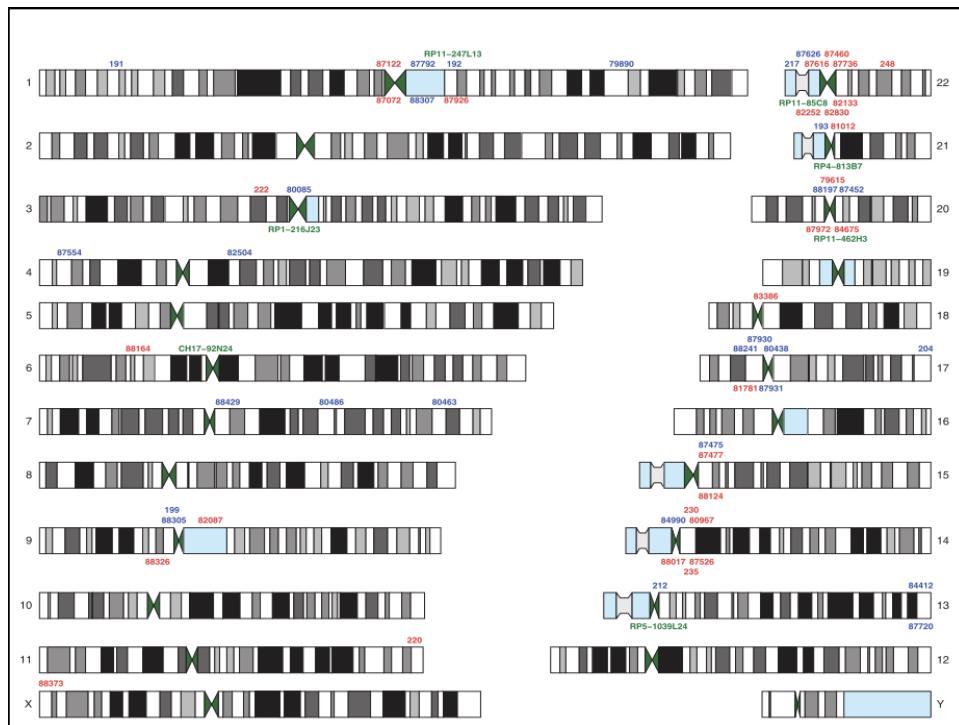
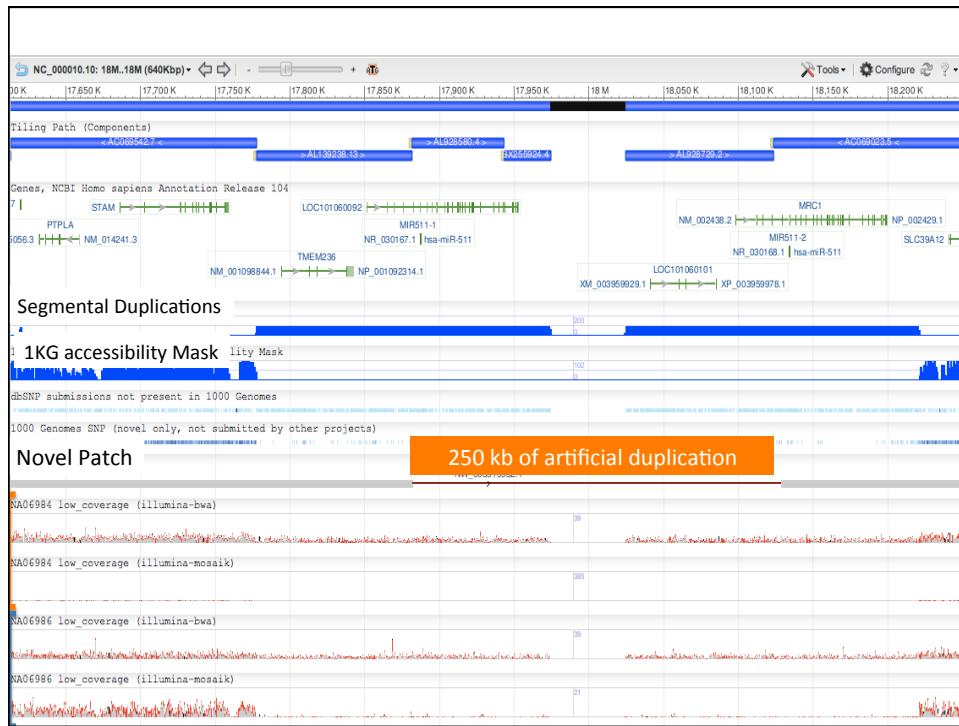


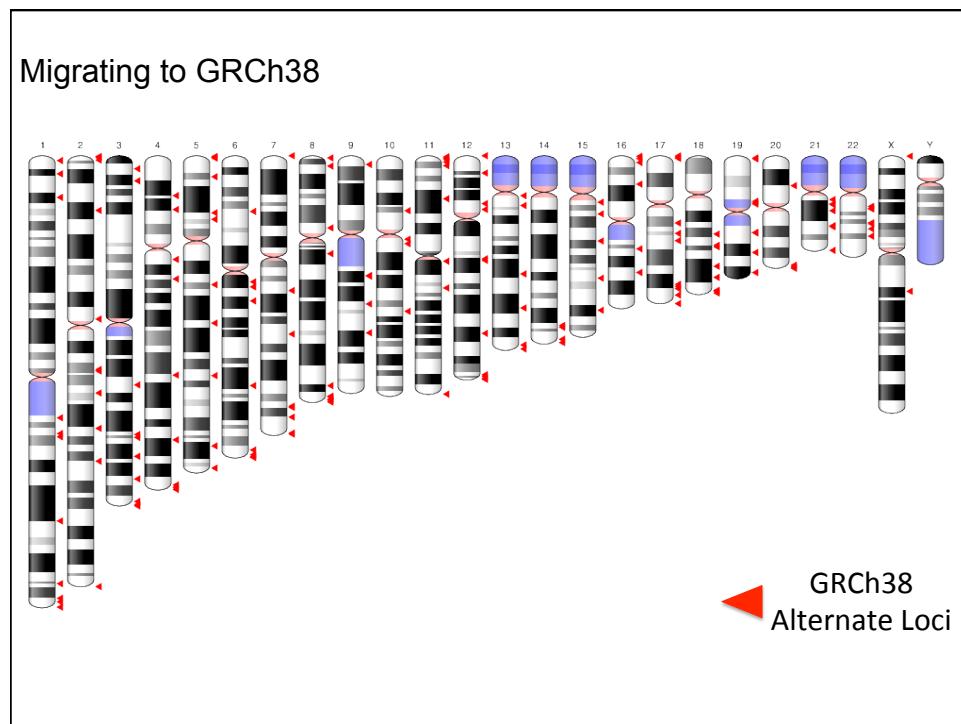
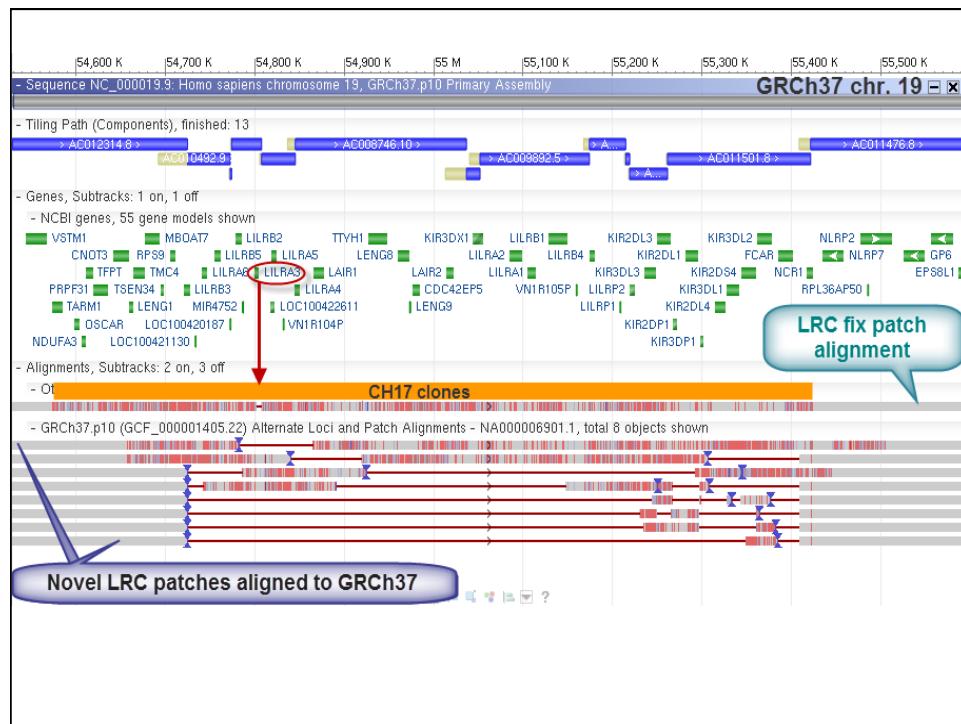


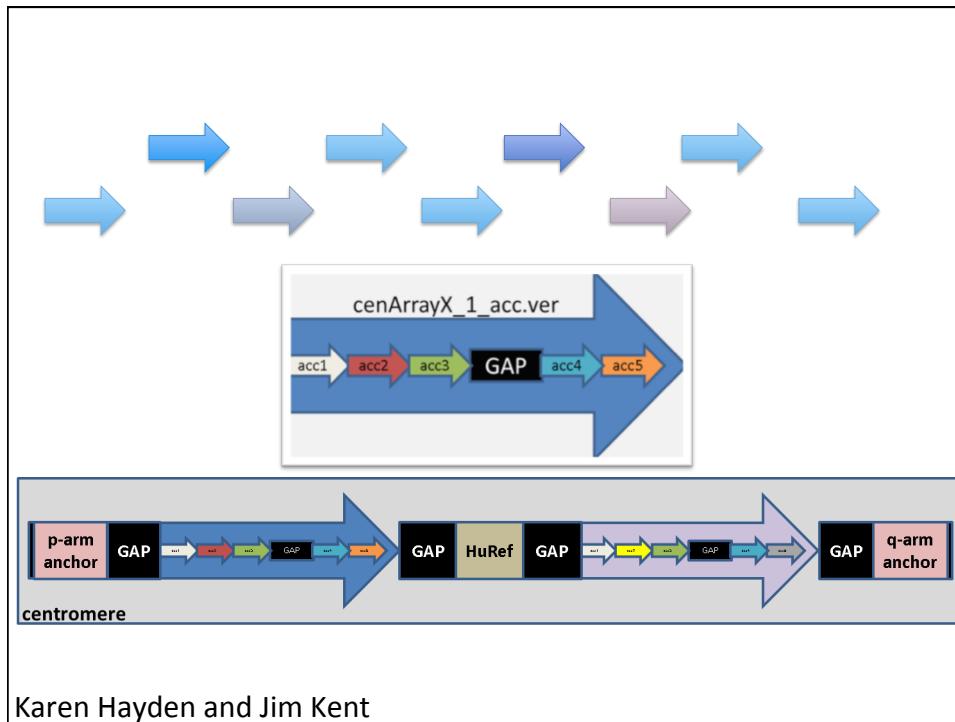












Take home messages

- The human and mouse reference assemblies are not linear single haplotype representations
- Complex variation and segmental duplication makes genome assembly challenging.
 - This is true even using newer long read data
- New assembly will challenge our understanding of annotation