

Web-based Tools for the Integration of Genomic Data

Mario Caccamo (TGAC)

The Genome Analysis Centre

A new facility to provide critical mass and excellence in genomics specialised in **animal**, **microbial** and **plant** research:

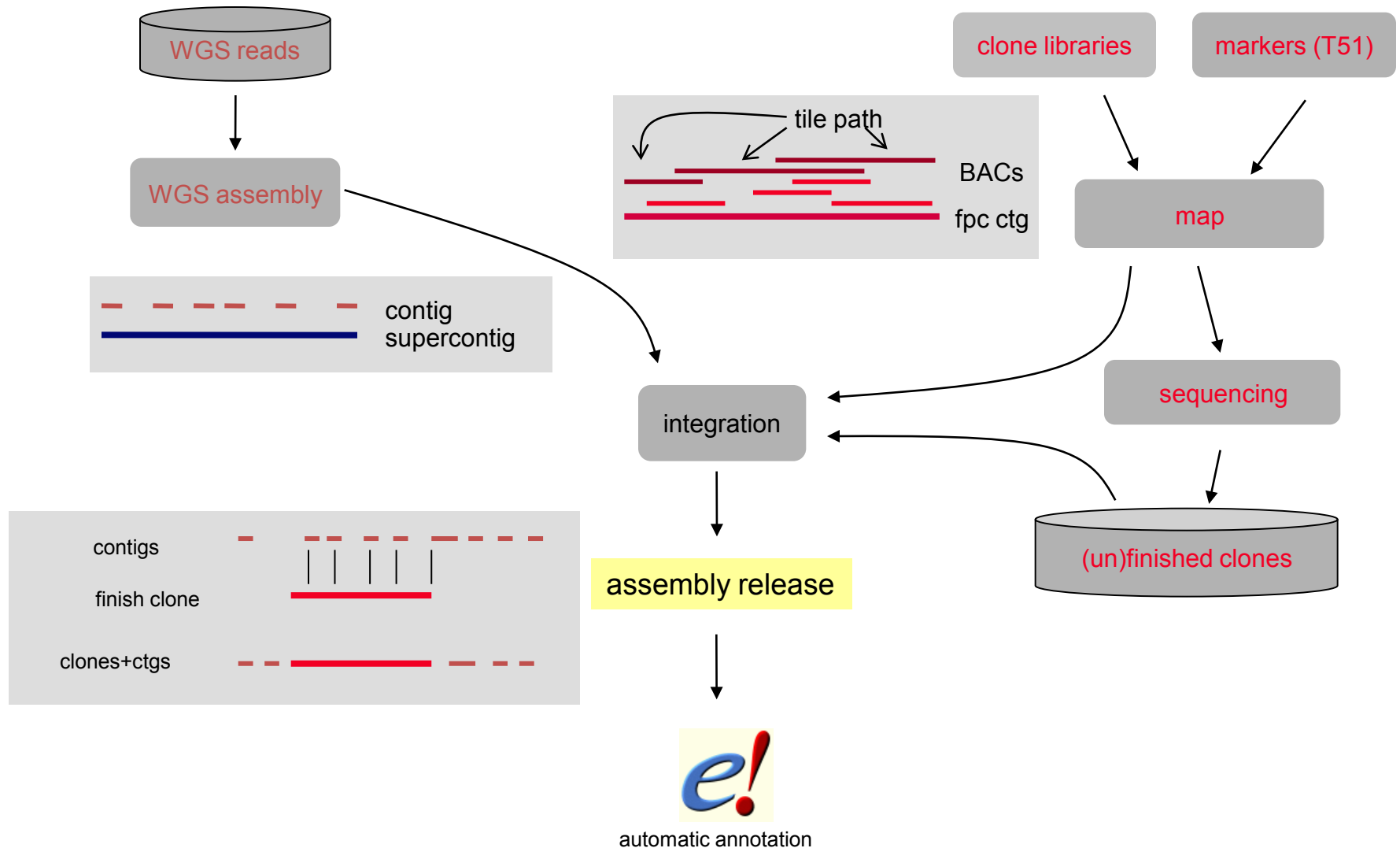
- high throughput sequencing
- new technology platforms
- bioinformatics
- impact through innovation and enterprise

www.tgac.bbsrc.ac.uk

Zebrafish Genome Project

whole genome shotgun sequencing

clone mapping and sequencing

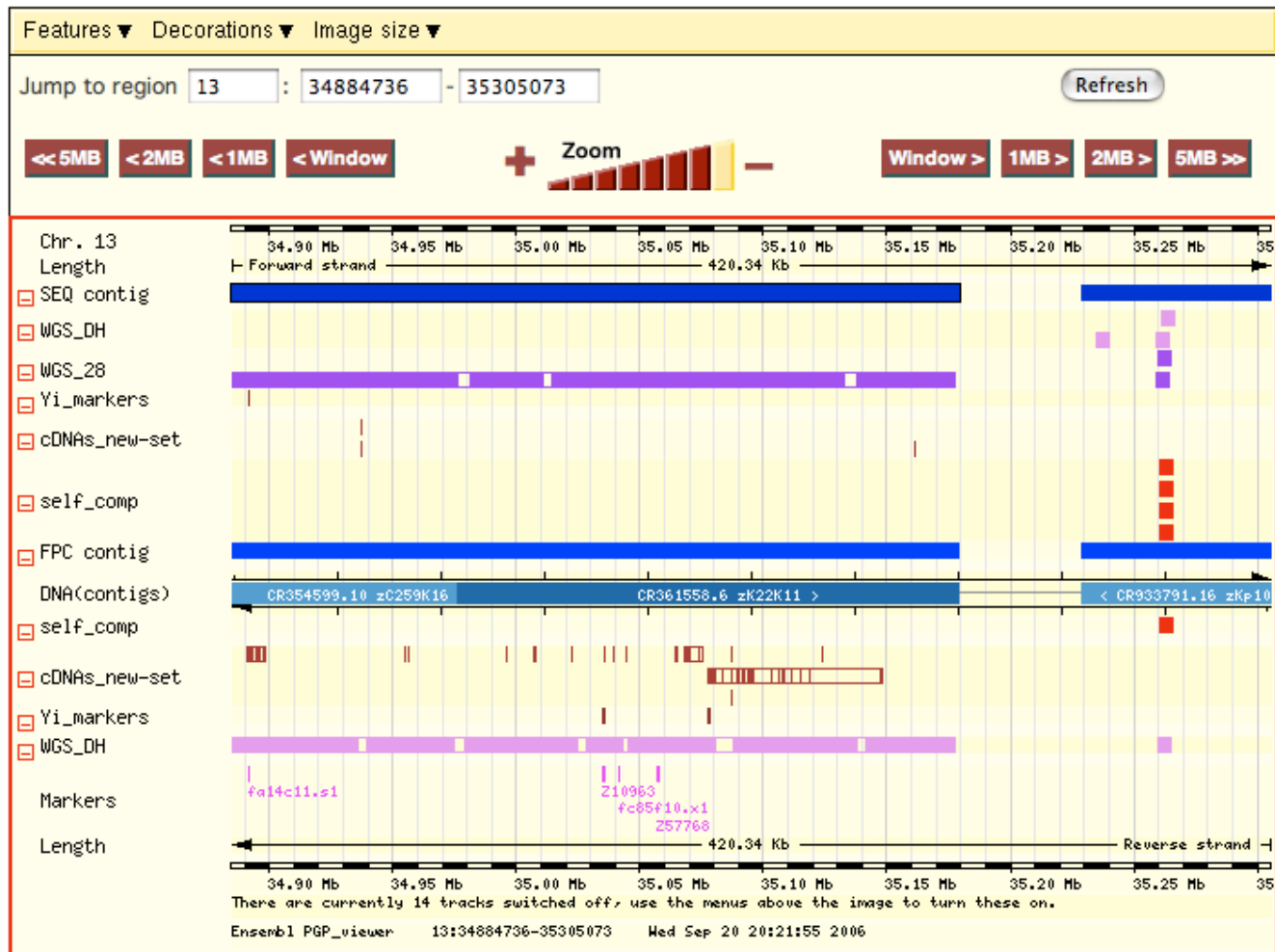


Sequence Analysis Tools

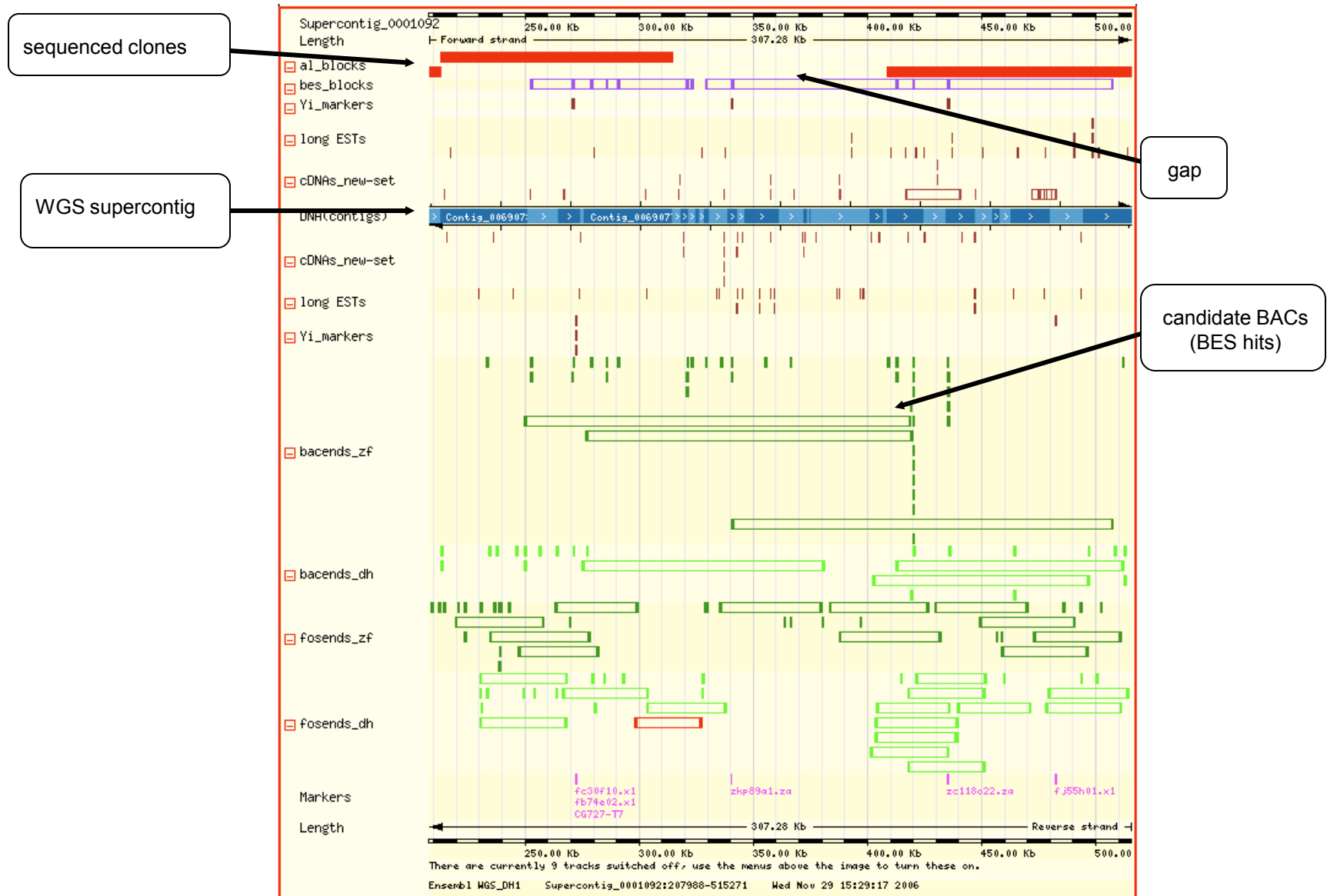
Design and implementation of tools to analyse, integrate and visualise the available resources

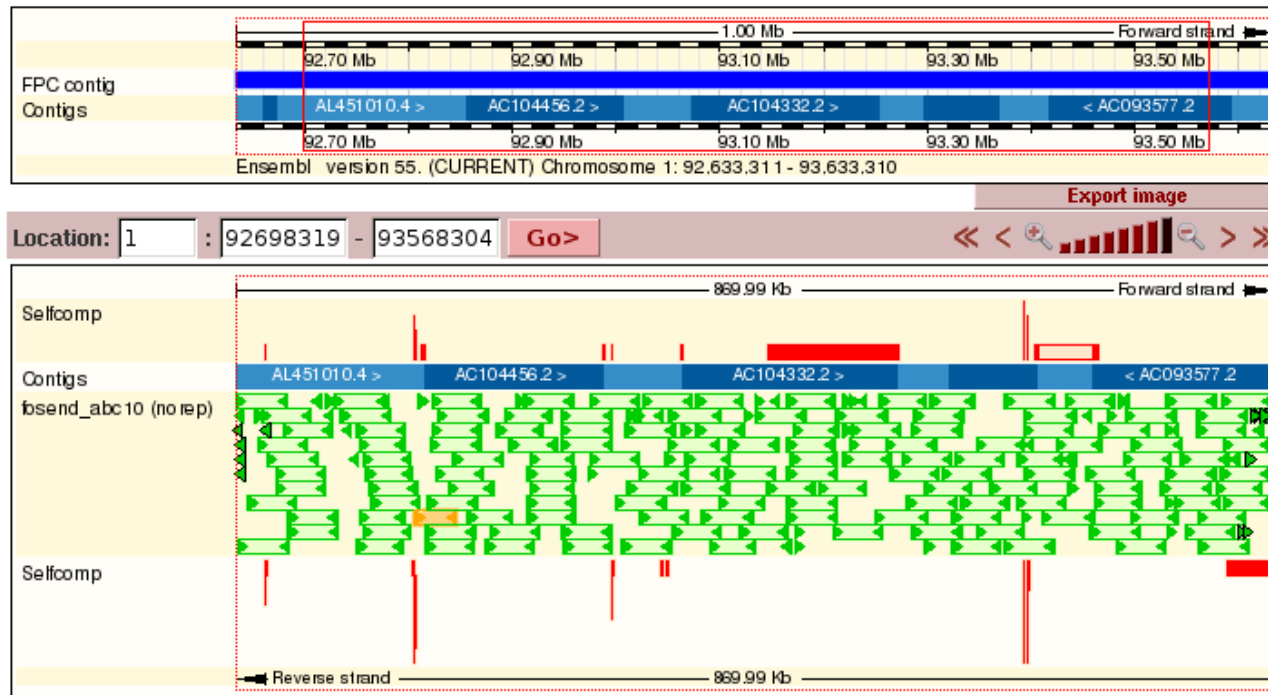
- PGP-viewer and WGS browsers
- “Punchlists”
- Data Analysis
eg. analysis of gaps, missing cDNAs, integration of WGS data

PGP Viewer



Bridging Gaps with WGS Sequences





Mapped 1 time



Mapped multiple times



Wrong direction (<<, <>, >>)

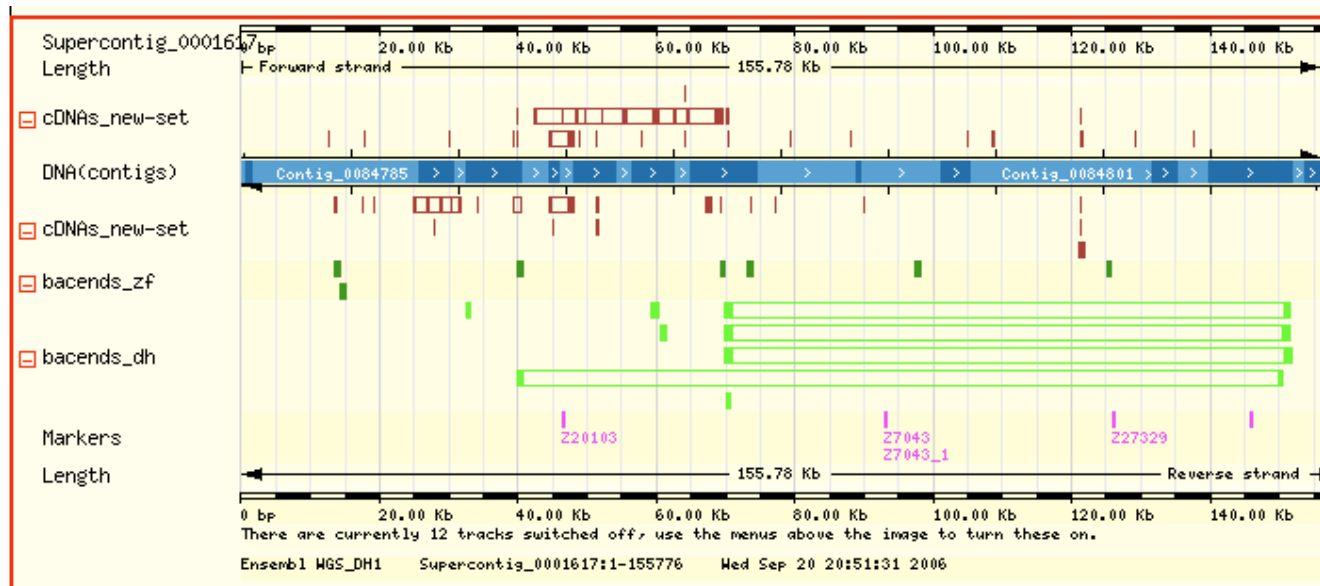


Wrong distance

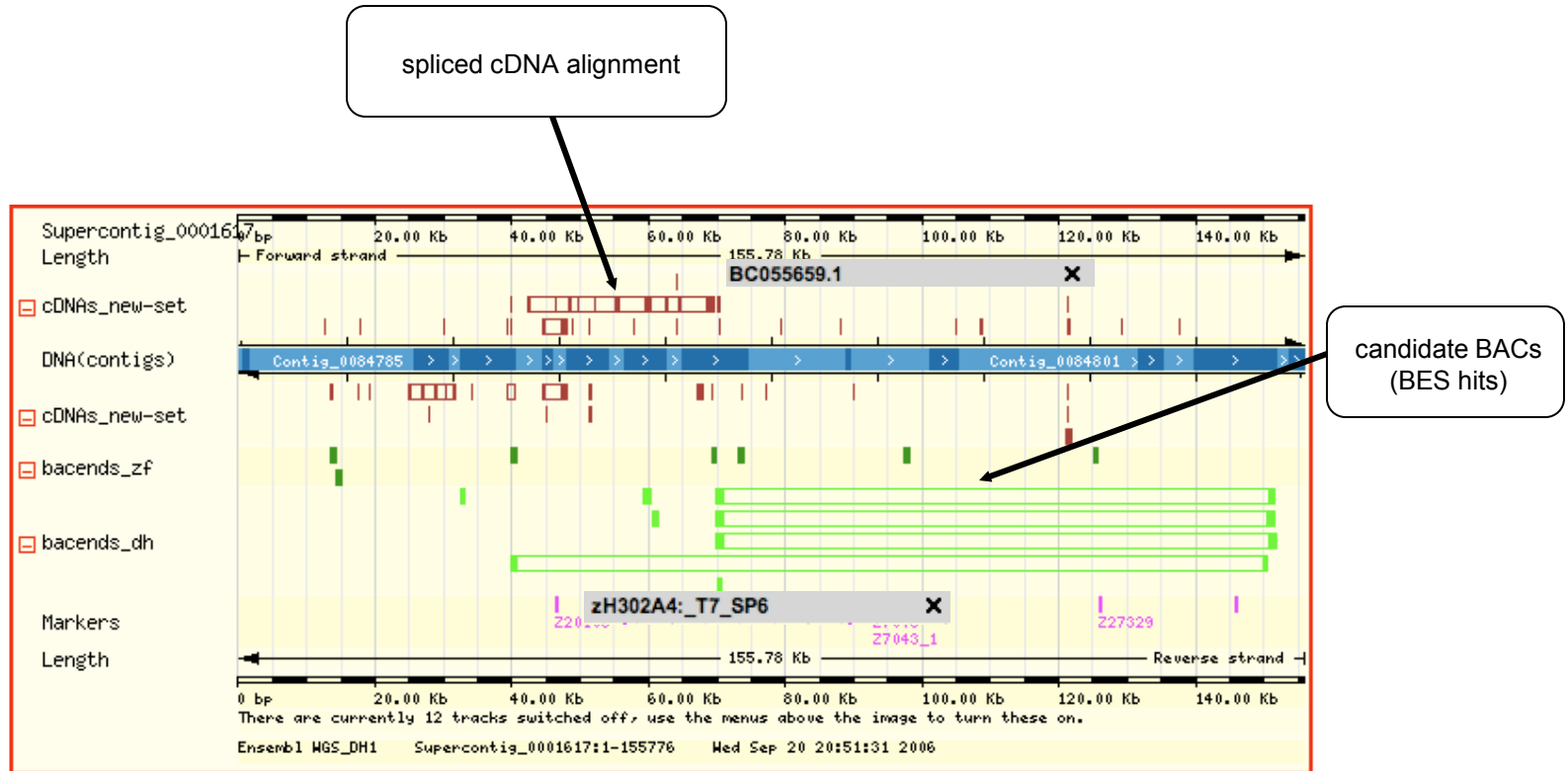


Spanner parter in the vicinity

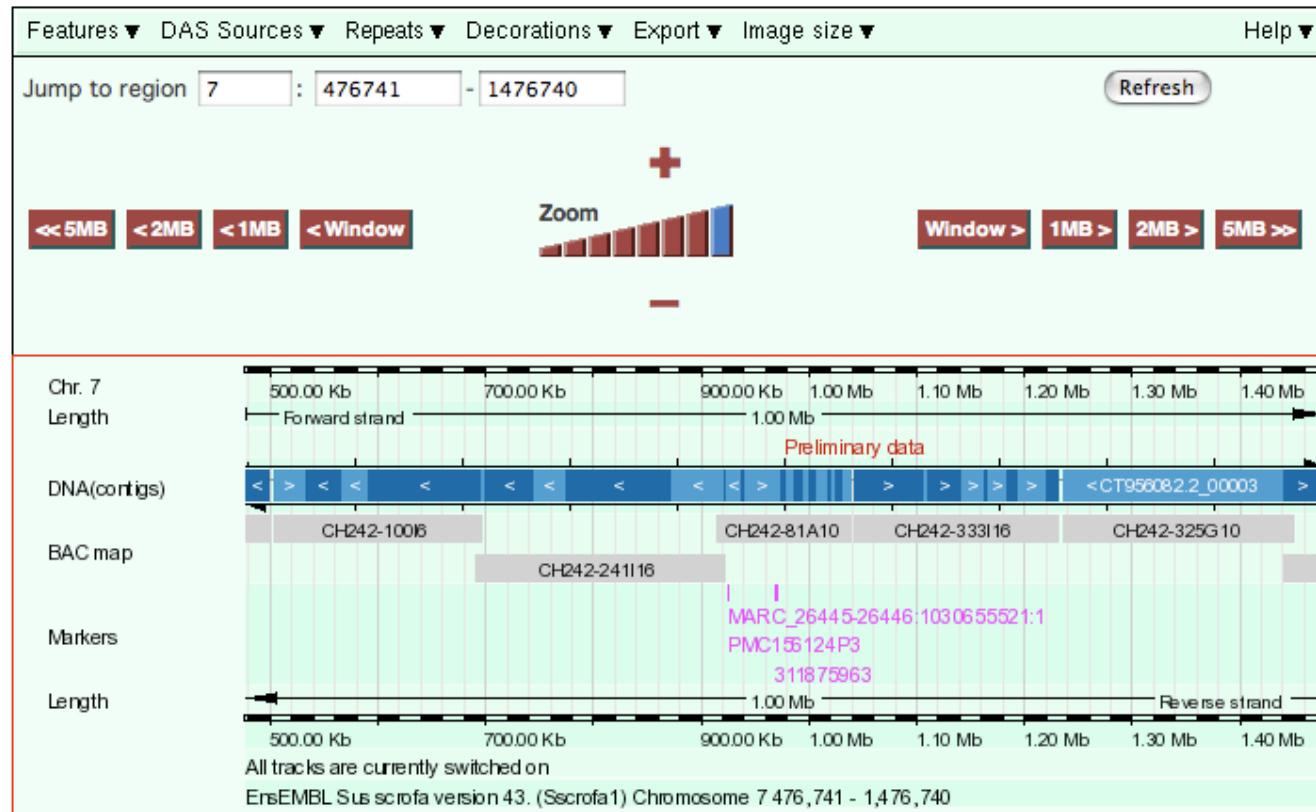
Region not covered by the physical map



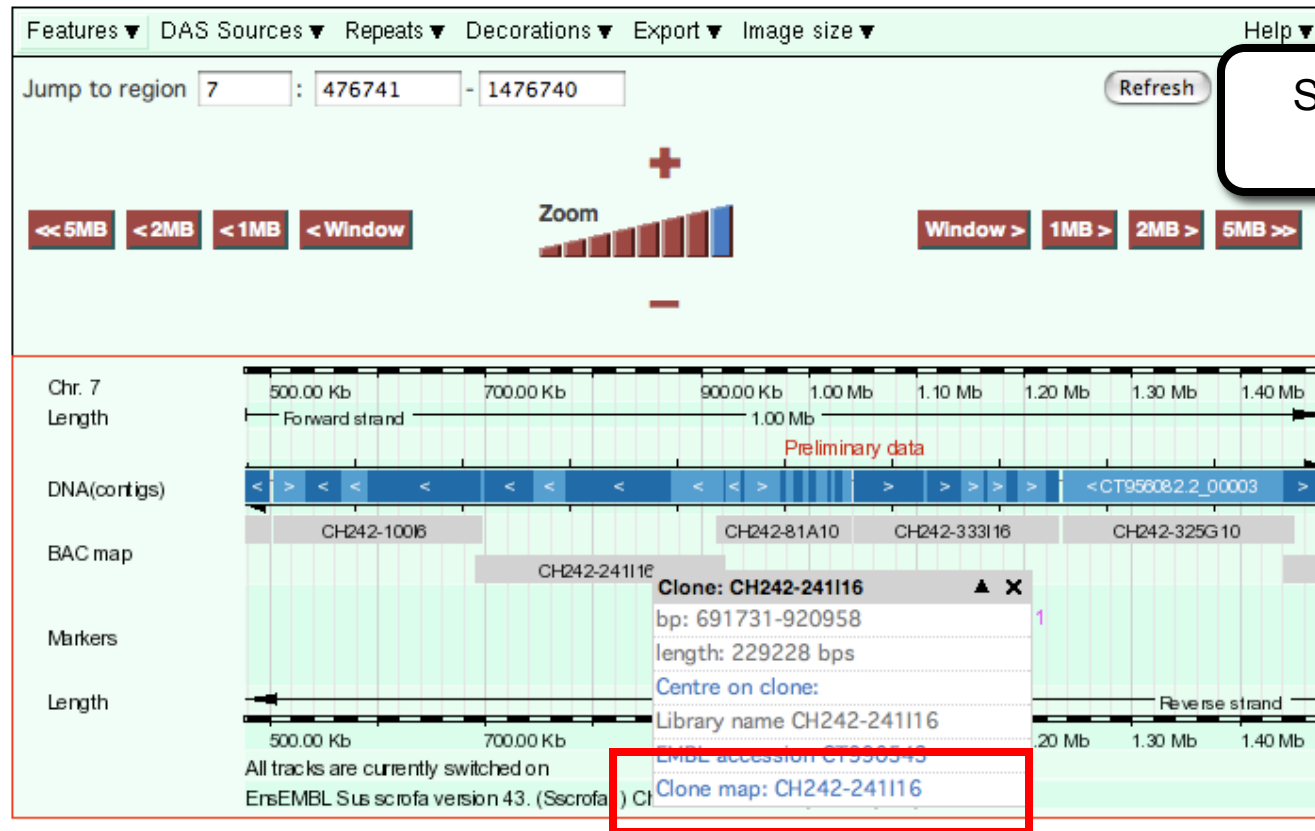
cDNAs in Regions not Covered by the MAP



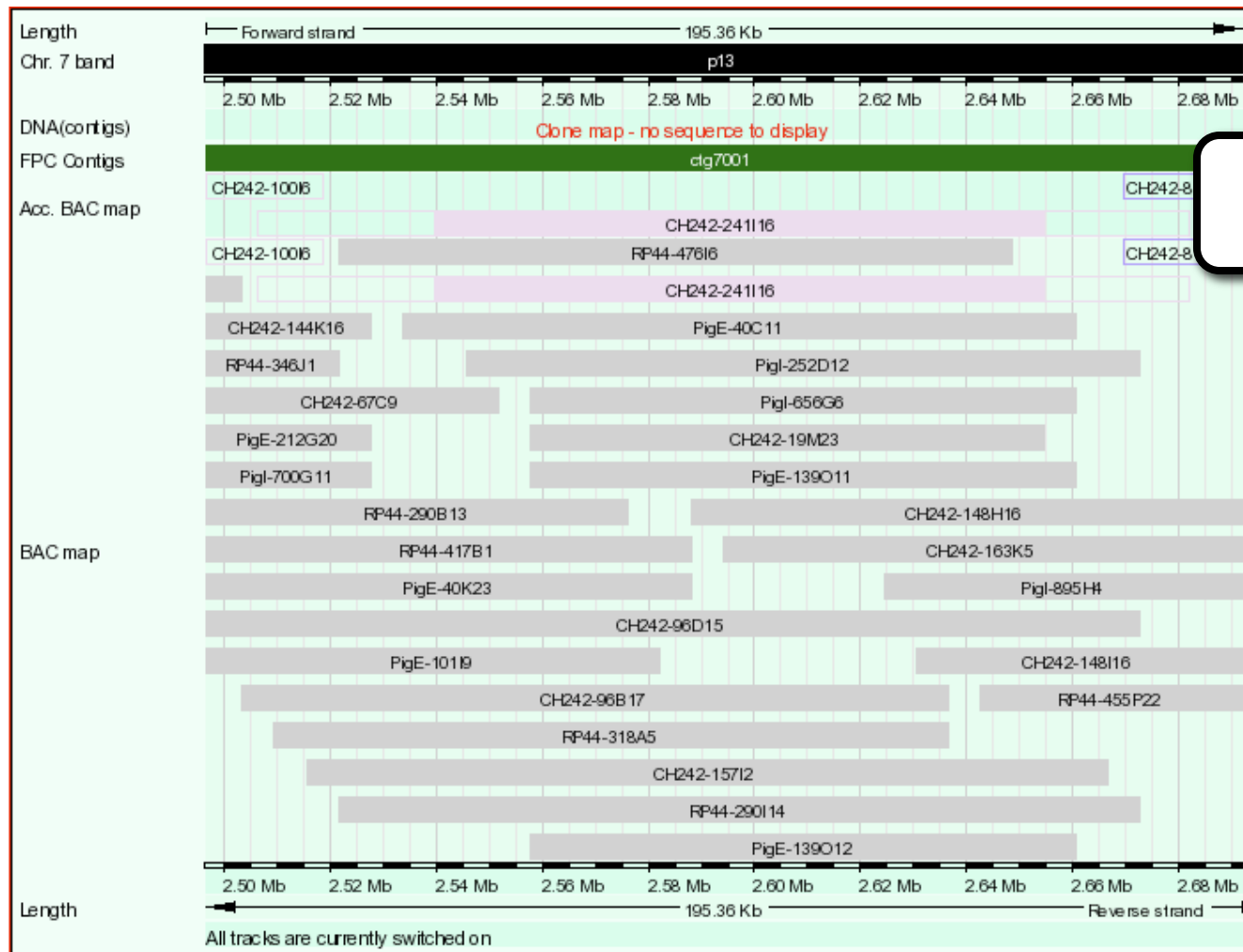
Tilepath



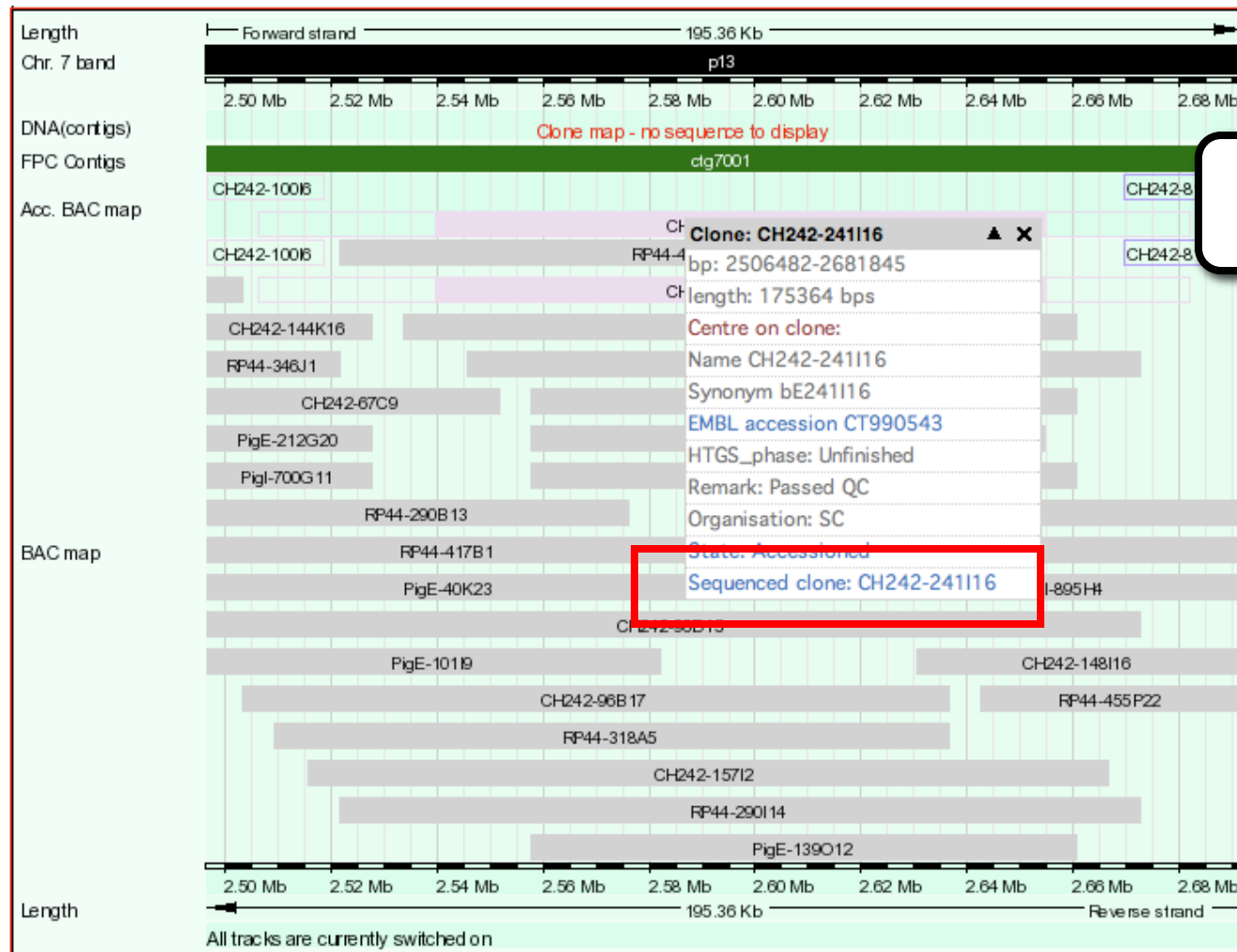
From the Sequence to the Map

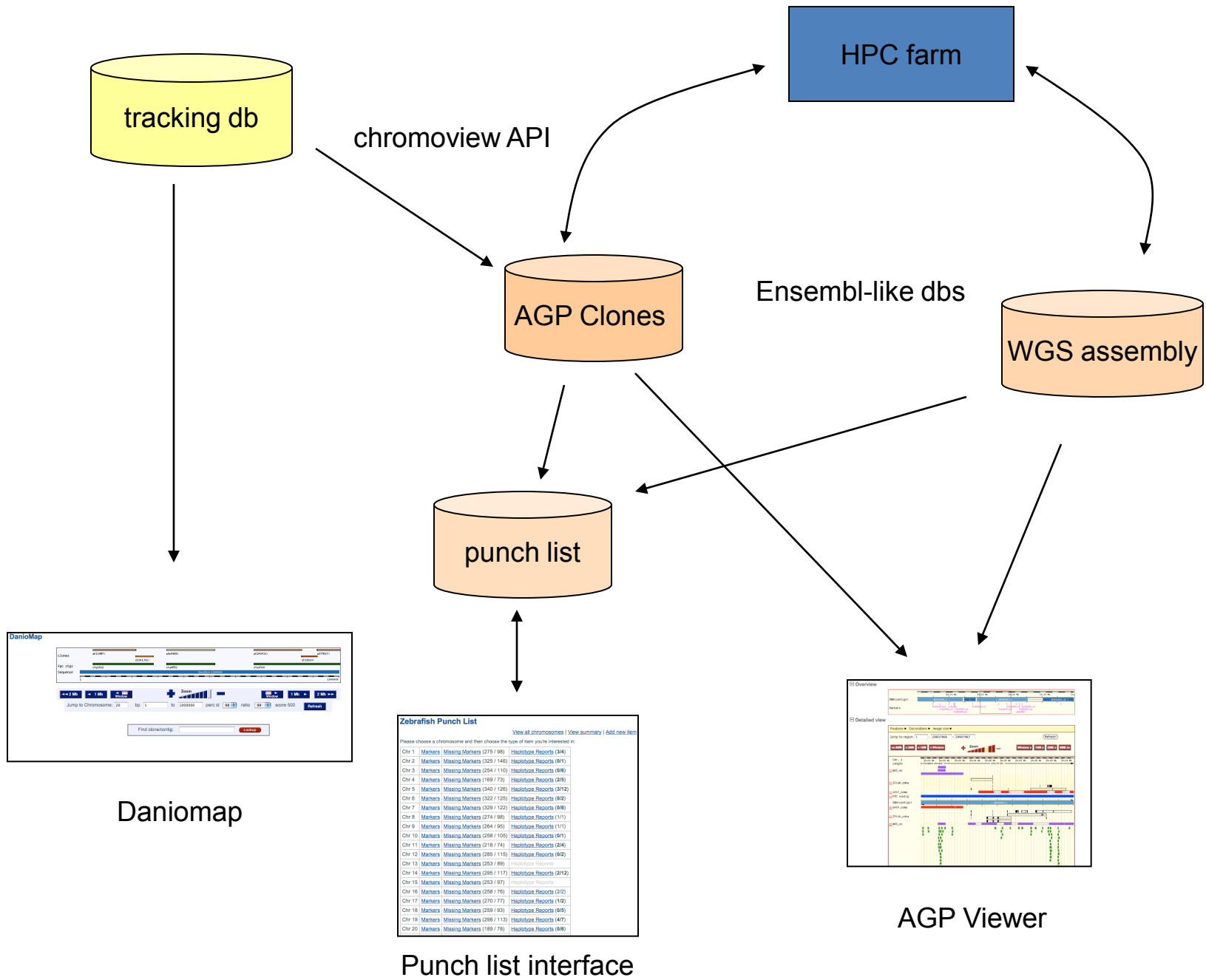


From the Sequence to the Map



From the Map to the Sequence





```
$slice=$slice_adaptor->fetch_by_region("chromosome",20);  
  
my $clone_projs = $slice->project("clone");  
  
foreach my $clone_proj (@$clone_projs){  
  
    #get all bacends for $clone_proj  
  
}
```

Punch Lists

Zebrafish Punch List

[View all chromosomes](#) | [View summary](#) | [Add new item](#) | [Add new haplotype report](#) | [Add new tandem report](#) | [Logout](#)

Chr: | Search: in

Please choose a chromosome and then choose the type of item you're interested in:

| Chr | Markers | Missing Markers (627 / 627) | Manual Haplotype Reports | Automatic Haplotype Reports | Failed AGP Joins | Failed PCR Joins | Tandem Reports |
|--------|--------------------|-----------------------------|-----------------------------|-----------------------------|------------------|------------------|----------------|
| Chr 1 | Ma | | Zebrafish Punch List | | | | |
| Chr 2 | Ma | | | | | | |
| Chr 3 | Ma | | | | | | |
| Chr 4 | Ma | | | | | | |
| Chr 5 | Ma | | | | | | |
| Chr 6 | Ma | | | | | | |
| Chr 7 | Ma | | | | | | |
| Chr 8 | Ma | | | | | | |
| Chr 9 | Ma | | | | | | |
| Chr 10 | Ma | | | | | | |
| Chr 11 | Ma | | | | | | |
| Chr 12 | Ma | | | | | | |

Chromosome 1 - Manual Haplotype Reports

[Add new haplotype report to chromosome 1](#)

Clones are clustered according to AGP as of 2006-09-15.

Manual Haplotype Report | Chromosome 1 | UNRESOLVED | Added by mc2, 2005-12-08 17:15 | [Add Additional Clone\(s\)](#)

☒ Mark Resolved ☐ Haplotype ☐ Redundant ☐ Clone Swap ☐ Duplication ☐ Large Overlap ☐ Repeat ☐ Certificate ☐ Not Haplotype

☒ Mark Pending

Clone Clusters:

| Chr 1 ctg15 | Chr 1 ctg15 | Chr 1 ctg15 |
|---|---|---|
| zC59A4 BX004756.7 Analysed View in PGP Viewer View in DanioMap View in Vega Position: 3994244 | zC25H13* CR956435.1* Cleared for library making View in PGP Viewer View in DanioMap | zC105H17 CT033808.5 Analysed View in PGP Viewer View in DanioMap Position: 4544143 |
| | zC274B7 BX323090.7 Analysed View in PGP Viewer View in DanioMap View in Vega Position: 4372790 Failed join report: [1] | |

Comments:

2005-12-08 17:15 mc2:
zC59A4 and zC274B7 share a fragment of gene dmd (check Ensembl) but they don't overlap in the current AGP

2005-12-12 13:56 dje:
zC274B7 is non dH but overlap with zC59A4 is haplotypic. A new zH is in place to try and resolve this.

2006-07-13 10:10 gkl:
zC105H17 also involved looks haplotypic wrt to zC59A4 - see otter

2006-08-23 16:59 sb2:
Haplotype issues with this area of ctg15, bZ86K19 (92% dh) / zK9K8 (94% dh) / zC59A4 (93% dh) look haplotypic to zK111N6 (3% dh) / zC105H17 (11% dh) / zC274B7 (42% dh)

[Add Comment](#)

Genome Reference Consortium

<http://genomereference.org>

Genome Reference Consortium

[GRC Home](#)

[Human](#)

[Mouse](#)

[Help](#)

[Report an Issue](#)

[Contact Us](#)

[Curators Only](#)

The Genome Reference Consortium

At the time the human reference was initially described, it was clear that some regions were recalcitrant to closure with existing technology. What was not as clear was the degree to which structural variation affected our ability to produce a truly representative genome sequence at some loci. It is now apparent that some regions of the genome are sufficiently variable that they are best represented by multiple sequences in order to capture all of the sequence potentially available at these loci.

In order to improve the representation of the reference human genome we have formed the Genome Reference Consortium (GRC). The goal of this group is to correct the small number of regions in the reference that are currently misrepresented, to close as many remaining gaps as possible and to produce alternative assemblies of structurally variant loci when necessary. We will provide mechanisms by which the scientific community can report loci in need of further review. In addition, information about loci currently under review and genome assembly production cycles will be made readily available. The human reference assembly is the cornerstone upon which all whole genome studies are based. It is critical to ensure that we have the best possible view of the genome to facilitate continued progress in understanding and improving human health.

The Genome Reference Consortium consists of:



The Wellcome Trust Sanger Institute



AT WASHINGTON UNIVERSITY The Genome Center at Washington University



The European Bioinformatics Institute



The National Center for Biotechnology Information

GRC News and Updates

GRC in the News

Tue, 29 Jan 2009

The GRC is highlighted in a Nature news feature.

GRCh37 is now available in Map Viewer

Fri, 14 Aug 2009

NCBI has annotated and released the latest version of the public human genome assembly (GRCh37).

[see all](#)

Resolved Issues

Mouse (MG-202)

Jan12, 2010

Switchpoints have been curated to use the sequence of AC102933.6, which agrees with the transcript, through this region of overlap.

Mouse (MG-180)

Dec31, 2009

The switchpoints have been curated to use the sequence of overlapping accession AC138666.6, which does not contain this discrepancy, through the overlap.

[see all](#)

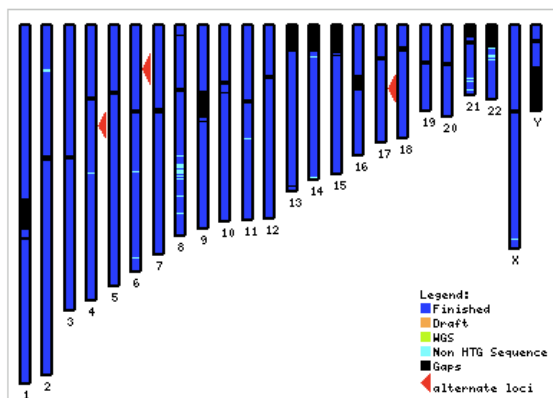
Genome Reference Consortium

Genome Reference Consortium

[GRC Home](#)
[Human](#)
[Mouse](#)
[Help](#)
[Report an Issue](#)
[Contact Us](#)
[Curators Only](#)
[Overview](#)
[Issues under Review](#)
[Assembly Data](#)
[Report a problem](#)

Human Genome Overview

Information concerning continuing improvement of the human genome.



GRCh37: A graphical representation of the latest human assembly. The genome is colored with respect to the genomic component used to build the genome assembly at that location. The red triangles mark regions where alternate loci have been provided.

The most recent assembly for human is GRCh37. This is the first assembly produced by the GRC and is considered the next version of NCBI Build 36 (also known as hg18). Improvements in this assembly include:

- Closure of 25 unspanned gaps found in Build 36
- Resolution of over 150 issues reported as problems in Build 36
- Addition of alternate loci for three complex regions, including the [MHC region](#).
- Standardization of AGPs, including the addition of biological gap information.

GRCh37 is a [haploid assembly](#), constructed from multiple individuals and can be divided into a 'primary assembly' and a set of 'alternate loci'. The [primary assembly](#) represents the assembled chromosomes, plus any [unlocalized](#) or [unplaced](#) sequence that represent the non-redundant, [haploid assembly](#). The [alternate loci](#) represent regions for which there is large scale variation and an alternate tiling path is available for this region. An example of such a region can be found at chromosome 17q21.31, often known as the MAPT locus. This region was described as carrying an inversion polymorphism ([PMID: 15654335](#)) and has been associated with various phenotypes ([PMID: 16718704](#) ; [PMID: 18628315](#)). The version of this region in Build 36 was actually a mosaic of both haplotypes (as tracked in HG-77) and has been resolved in GRCh37 thanks to data described in Zody et al., 2008 ([PMID: 19165922](#)).

Information on alternate loci

| Chromosome region with alternate loci | Length of region | Number of alternate contigs in region | View Region |
|---|------------------|---------------------------------------|----------------------|
| UGT2B17 region (chr4:69,170,077-69,877,175) | 707,099 bp | 1 contig [+] | view |
| MHC region (chr6: 28,477,797-33,448,354) | 4,970,558 bp | 7 contigs [+] | view |
| MAPT region (chr17: 43,384,864-44,913,631) | 1,528,768 bp | 1 contig [+] | view |

GRC News and Updates

GRCh37 now available at UCSC

Fri, 8 May 2009

UCSC has released the latest version of the public human genome assembly.

GRCh37 now available in Pre! Ensembl

Thu, 9 Apr 2009

Ensembl has released the latest release of the public human genome assembly (GRCh37) on their Pre! site.

[see all](#)

Recently Resolved Human Issues

Human (HG-546)

May14, 2009

The gap has been closed by adding CR394530.16 to the TPF

Human (HG-33)

May14, 2009

CR812477 closes the gap

[see all](#)

References

Whole Genome Papers

[The HGP Reference Assembly](#)
[The Venter Genome Assembly](#)

Human Chromosome Papers

[Chr1](#) [Chr2](#) [Chr3](#) [Chr4](#) [Chr5](#) [Chr6](#)
[Chr7](#) [Chr8](#) [Chr9](#) [Chr10](#) [Chr11](#)
[Chr12](#) [Chr13](#) [Chr14](#) [Chr15](#) [Chr16](#)
[Chr17](#) [Chr18](#) [Chr19](#) [Chr20](#) [Chr21](#)
[Chr22](#) [ChrX](#) [ChrY](#)

Conclusions

- Integration of the data
 - speeds up project
 - validation
 - systematic identification of problematic areas
- Scientific Community
 - expertise in relevant areas of the genome
 - more and better eyes to look for issues
 - share ownership of the final product
- Quality
 - lasting resource
 - usability
 - architecture of the sequence