# Module 1

# Understanding and Exploring Genome Assemblies

Kerstin Howe
Wellcome Trust Sanger Institute
zfish-help@sanger.ac.uk

## Genome assembly generation

- The Genome Reference Consortium

- Generating Assemblies

- Working With Assemblies

## Browsing genome assemblies

- Genome browsing

- The Ensembl gene set

- Guided examples

- [Make your own data visible: BED files]

## Genome assembly generation

- The Genome Reference Consortium

- Generating Assemblies

- Working With Assemblies

# Genome Reference Consortium
# genomereference.org

International consortium looking after
human, mouse and zebrafish reference assemblies

- maintaining reference assemblies

- improving reference assemblies

- adding variation

- All issues documented on website

**GRCz10** released August 2014

**GRCz11** released May 2017

Within GRC:

- Handover from Sanger Institute to ZFIN after GRCz11

- Curation now only in reaction to user enquiries!

# Generating zebrafish assemblies

- Restriction analysis (Fingerprint) and clone contig building (FPC map)

- Meiotic and RH maps

- Fill in WGS contigs

- Check and adjust with additional data (e.g. BioNano maps, Strandseq)
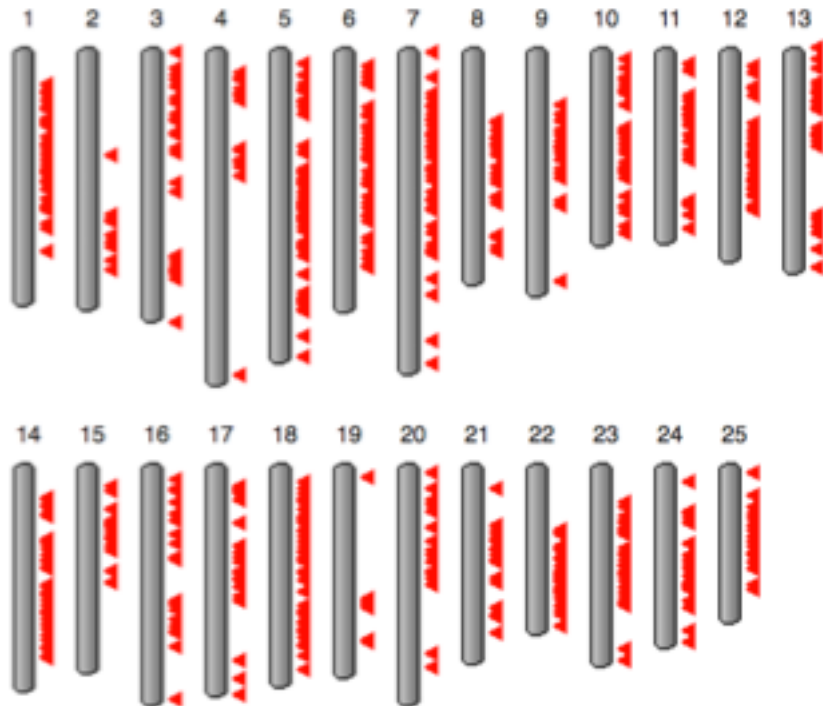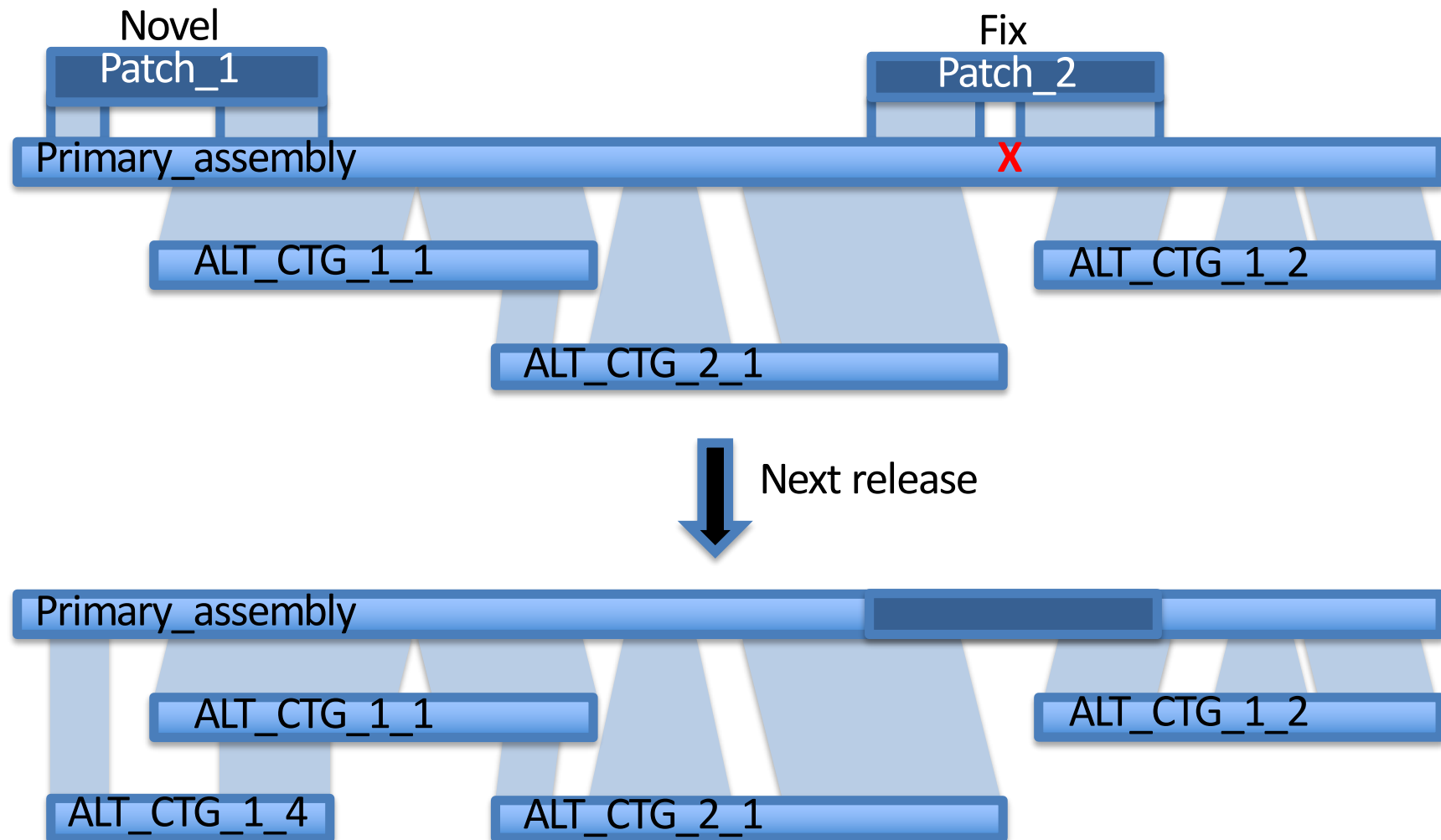
**Individual components**



**Scaffolds**

**chromosomes**

4

**Another layer: ALT_REF_LOCI**

- Representation of variation extracted from surplus clone sequence

- > 5kb indel compared to reference

**ALT_REF_LOCI**



930 ALT_REF_LOCI for GRCz11

# Generating Assemblies

**...and another layer: Patches**

# Working With Assemblies

You might want to

- Find a clone that covers a feature

- Confirm that a region is correctly positioned

- Find out whether a gap could be closed

- Check whether a gene is really duplicated

**Help is at hand**

- zfish-help@sanger.ac.uk

- Ensembl

- **gEVAL**    geval.sanger.ac.uk

**gEVAL**          **geval.sanger.ac.uk**

- genome evaluation browser
- Numerous assembly versions
- alignments of BAC/FOS ends, markers, optical maps, cDNAs, other genome assemblies, etc. to check consistency in the assembly
- reports GRC investigations
- offers 'punchlists' denoting issues with an assembly
- extensive documentation
- featured GRCz11 from release date

**geval.sanger.ac.uk**

**Browsing genome assemblies**

- Genome browsing

- The Ensembl gene set

- Guided examples

- [Make your own data visible: BED files]

# Genome Browsers

- NCBI Map Viewer          http://www.ncbi.nlm.nih.gov/mapview/

- UCSC Genome Browser      http://genome.ucsc.edu/

- Ensembl Genome Browser   http://www.ensembl.org/



**Caveat: check the assembly version!**

# Ensembl : What annotation is available?

- **Genes**  protein coding genes
  gene/transcript/peptide models (coding and non-coding)

- **Comparative data**  orthologues and paralogues, gene trees
  protein families
  whole genome alignments
  syntenic regions

- **Variation data**  Single Nucleotide Polymorphisms (SNPs), indels, phenotypes,
  population frequencies, variant effect prediction (VEP), etc.

- **Regulatory data**  e.g. regulatory elements from ENCODE

- **IDs**  crossreferences to other databases

- **aligned data**  cDNAs, RNAseq, peptides, micro array probes, BAC clones, etc.

- **Cytogenetic bands, markers, repeats** etc.

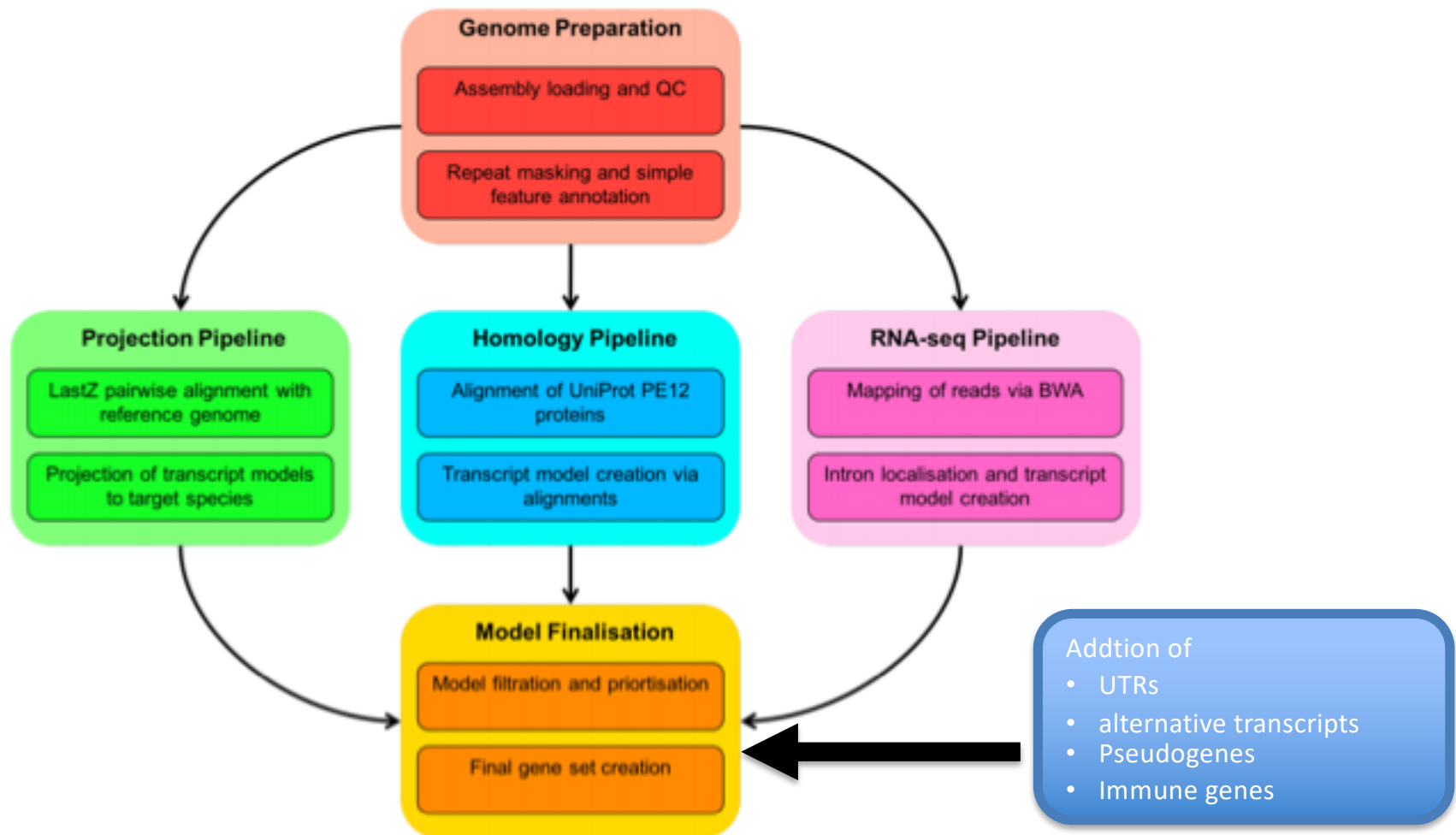- **External resources**  e.g. GRC trackhub, mapped next-gen reads

# The Ensembl gene set

- **Bimonthly releases, updated gene set ~ every 6 months**
- **New 'Genebuild' with every new assembly**

The following archives are available for this page:

- Ensembl 91: Dec 2017 (GRCz10)
- Ensembl 90: Aug 2017 (GRCz10) - patched/updated gene set Jun 2017
- Ensembl 89: May 2017 (GRCz10)
- Ensembl 88: Mar 2017 (GRCz10)
- Ensembl 87: Dec 2016 (GRCz10) - patched/updated gene set Nov 2016
- Ensembl 86: Oct 2016 (GRCz10) - patched/updated gene set Jul 2016
- Ensembl 85: Jul 2016 (GRCz10)
- Ensembl 84: Mar 2016 (GRCz10) - patched/updated gene set Jan 2016
- Ensembl 83: Dec 2015 (GRCz10)
- Ensembl 82: Sep 2015 (GRCz10)
- Ensembl 81: Jul 2015 (GRCz10)
- Ensembl 80: May 2015 (GRCz10) - gene set updated May 2015
- Ensembl 79: Mar 2015 (Zv9)
- Ensembl 78: Dec 2014 (Zv9)
- Ensembl 77: Oct 2014 (Zv9)
- Ensembl 76: Aug 2014 (Zv9)
- Ensembl 75: Feb 2014 (Zv9) - patched/updated gene set Feb 2014
- Ensembl 74: Dec 2013 (Zv9) - patched/updated gene set Jul 2013
- Ensembl 67: May 2012 (Zv9) - patched/updated gene set Mar 2012
- Ensembl 54: May 2009 (Zv8) - gene set updated Apr 2009
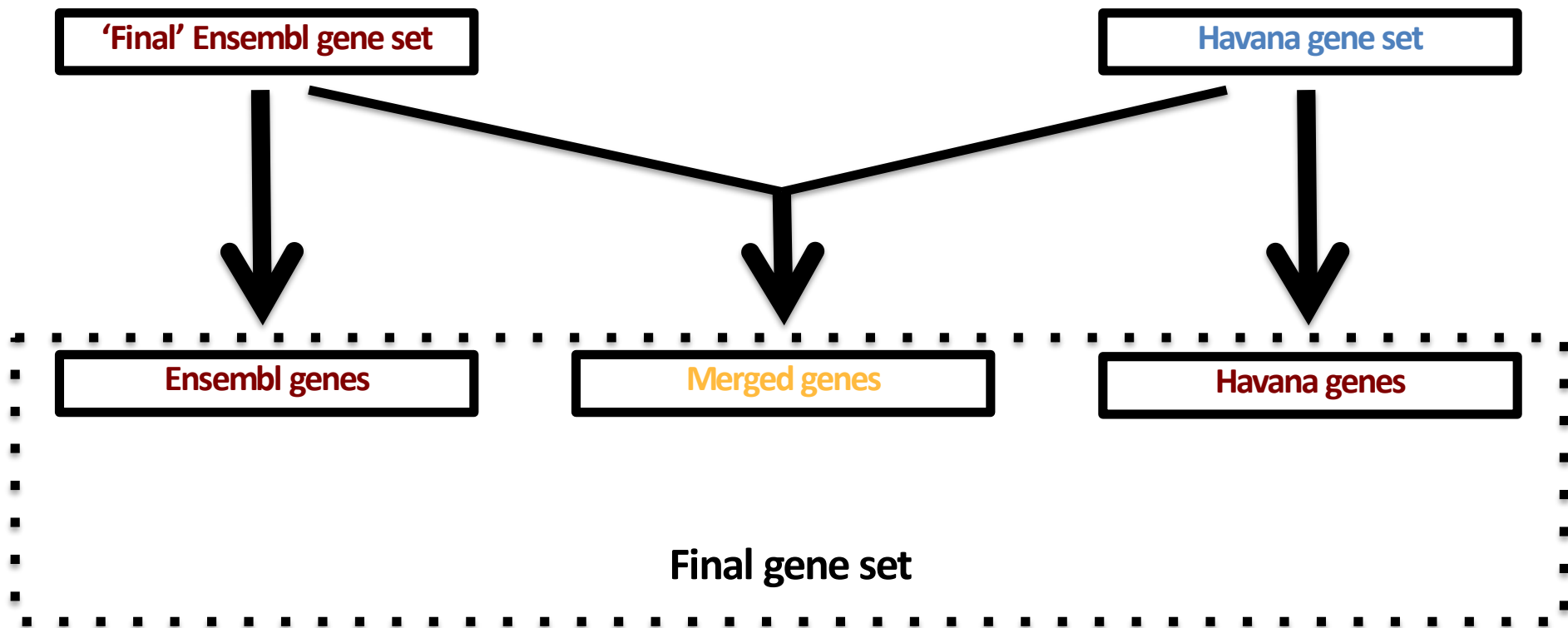
- **Genes are built on evidence, no gene is predicted on sequence alone!**

# The Ensembl gene build



From: http://www.ensembl.org/info/genome/genebuild/2018_03_zebrafish_genebuild.pdf

# The Ensembl gene build : Gene merge

**Getting the optimum gene set by combining automated and manual annotation**

# The Ensembl gene build : Gene merge

# Ensembl gene names

Genes with ZFIN record      **yes** ⟶    nrf1

                                      zgc:345

⬇ No ZFIN record

Genes with human ortholog    **yes** ⟶    TXNL4B

                                 PROZ (1 of 2)

⬇ no                              C1H13orf34

Relation to other databases    **yes** ⟶    5S_rRNA.386 (RFAM)

                                 dre-mir-734.1 (miRBase)

⬇ no

Gene names based on component names ⟶    CU856394.1

                                             CABZ01063868.2

# Access to Genome Annotation

- Release web site                     http://www.ensembl.org/
- Archive                               http://archive.ensembl.org

- BioMart                 http://www.ensembl.org/biomart/martview
- Downloads        http://www.ensembl.org/info/data/ftp/index.html

- Perl API            http://www.ensembl.org/info/docs/index.html
- REST API

# Help and Information

- Zebrafish specific help  zfish-help@sanger.ac.uk / zfinadmin@zfin.org

- Zebrafish genome project at Sanger Institute
  www.sanger.ac.uk/science/data/zebrafish-genome-project

- Zebrafish at the GRC    genomereference.org

- Zebrafish genome assembly evaluation geval.sanger.ac.uk

- Ensembl helpdesk        helpdesk@ensembl.org

- View animated tutorials
  www.ensembl.org/info/website/tutorials/index.html

- Mailing lists:             announce@ensembl.org
                             ensembl-dev@ensembl.org

# Guided examples

- A stroll through Ensembl - abbreviated

- Introduction to BioMart - abbreviated

- Making your own data visible (in your own time)