



TGAC 
The Genome Analysis Centre™

 **BBSRC**

 **EEDA**
East of England Development Agency

**Greater Norwich
Development
Partnership**



Bread Wheat Chromosome-based Survey Sequencing Initiative Update

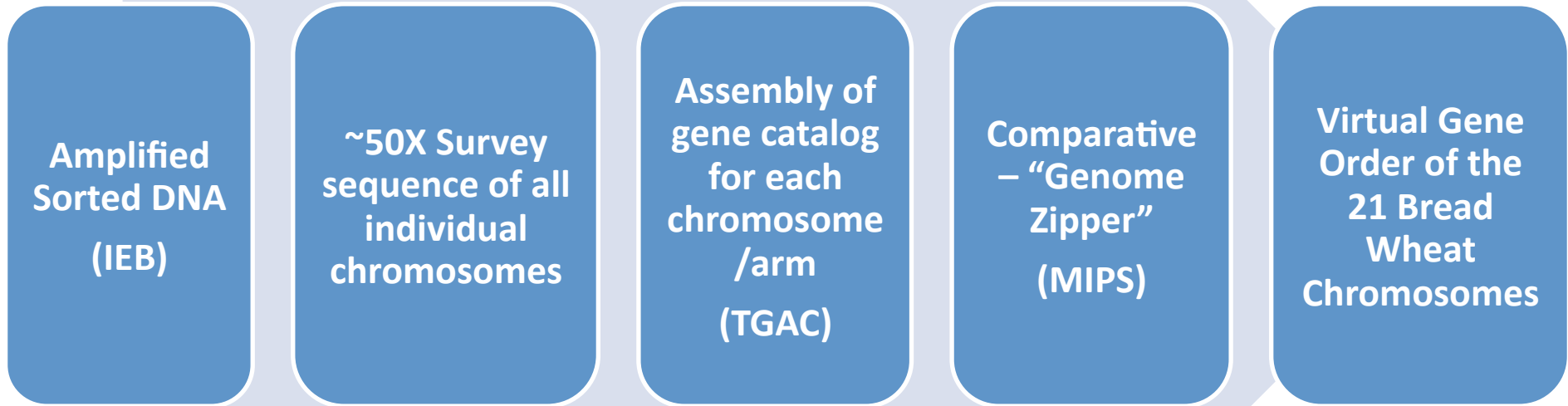
Jane Rogers

**IWGSC Workshop
International Triticeae Mapping Initiative - 2012**

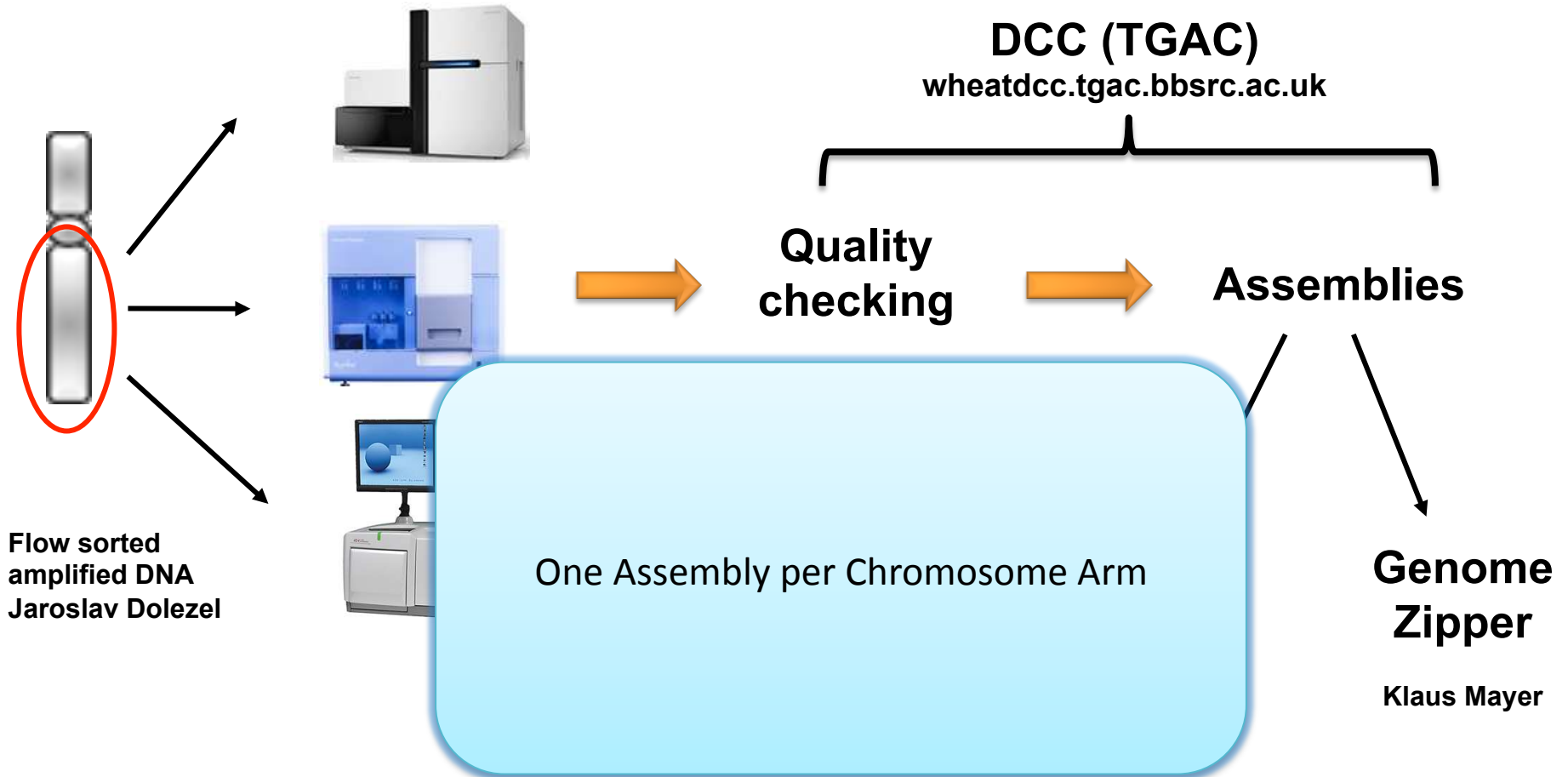
Sequencing Survey Initiative



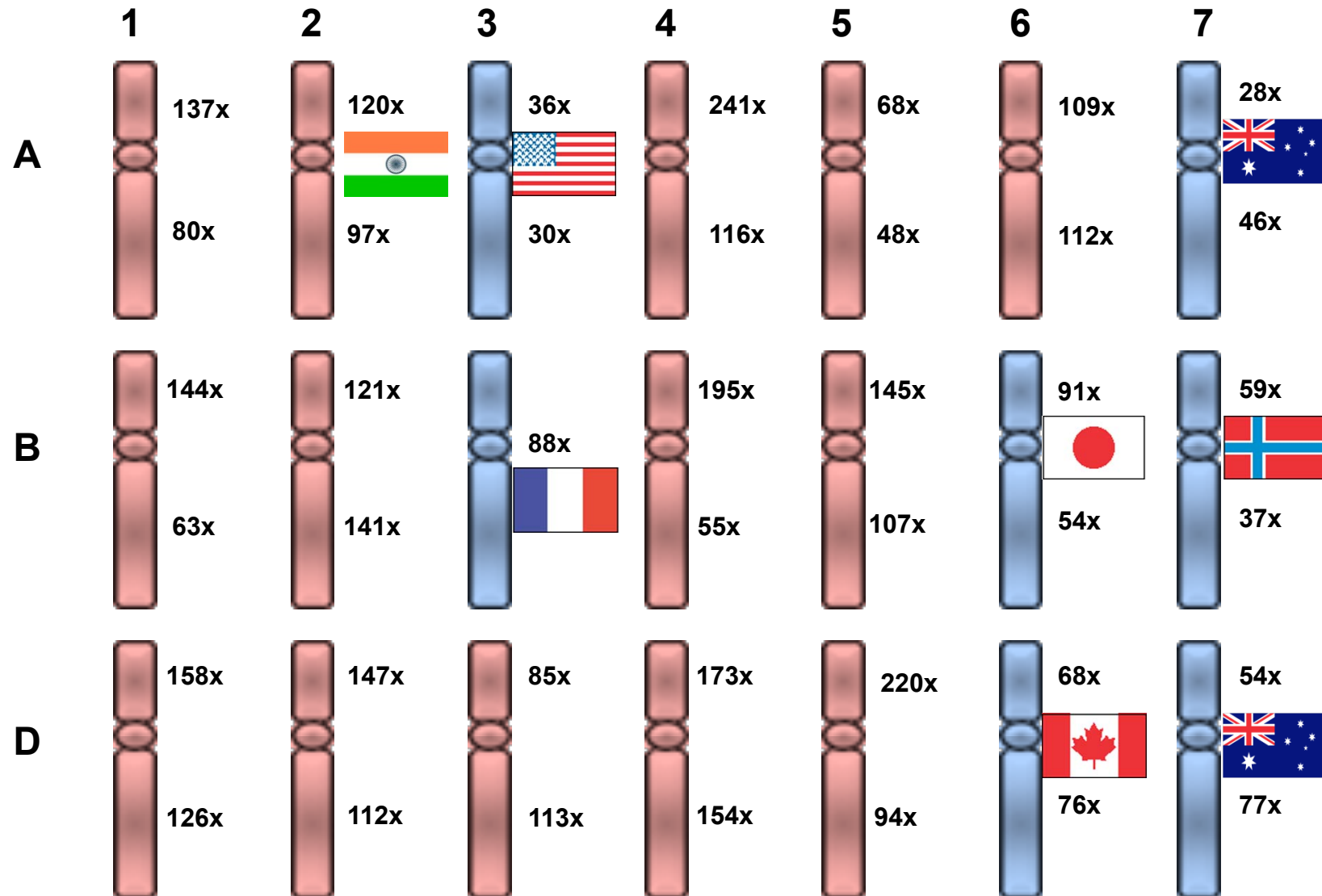
HelmholtzZentrum münchen
German Research Center for Environmental Health



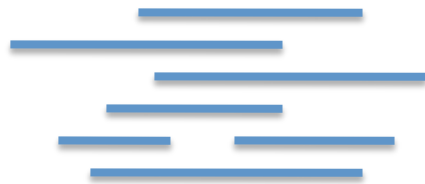
Project Overview



Illumina Sequence Coverage



Stage 1: Sequence Assembly



Set of contigs



Remove contigs
< 200 bp

Assessment of
assemblies

Resource

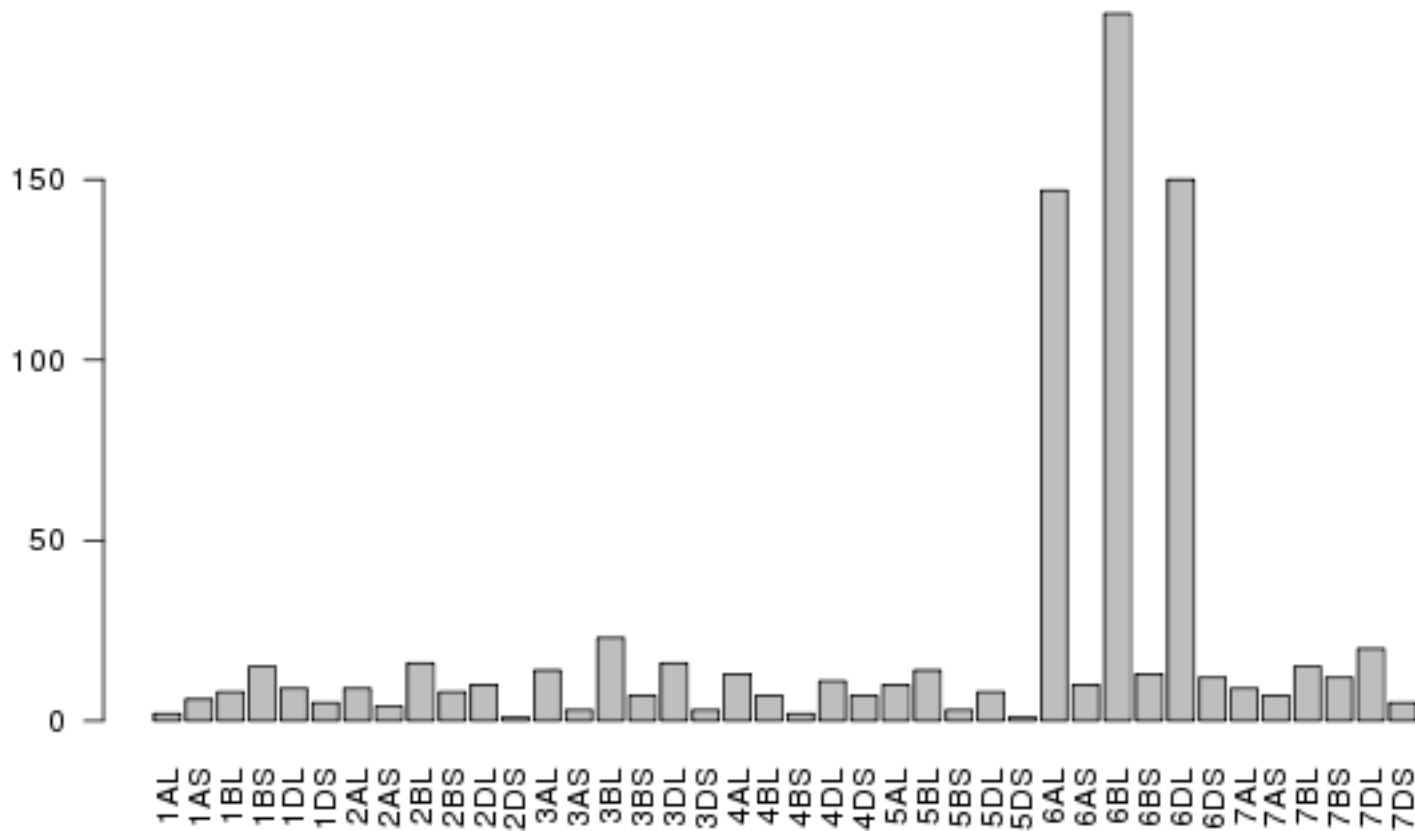
ABYSS: A parallel assembler for short read sequence data

Jared T. Simpson,¹ Kim Wong, Shaun D. Jackman, Jacqueline E. Schein,
Steven J.M. Jones, and İnanç Birol²

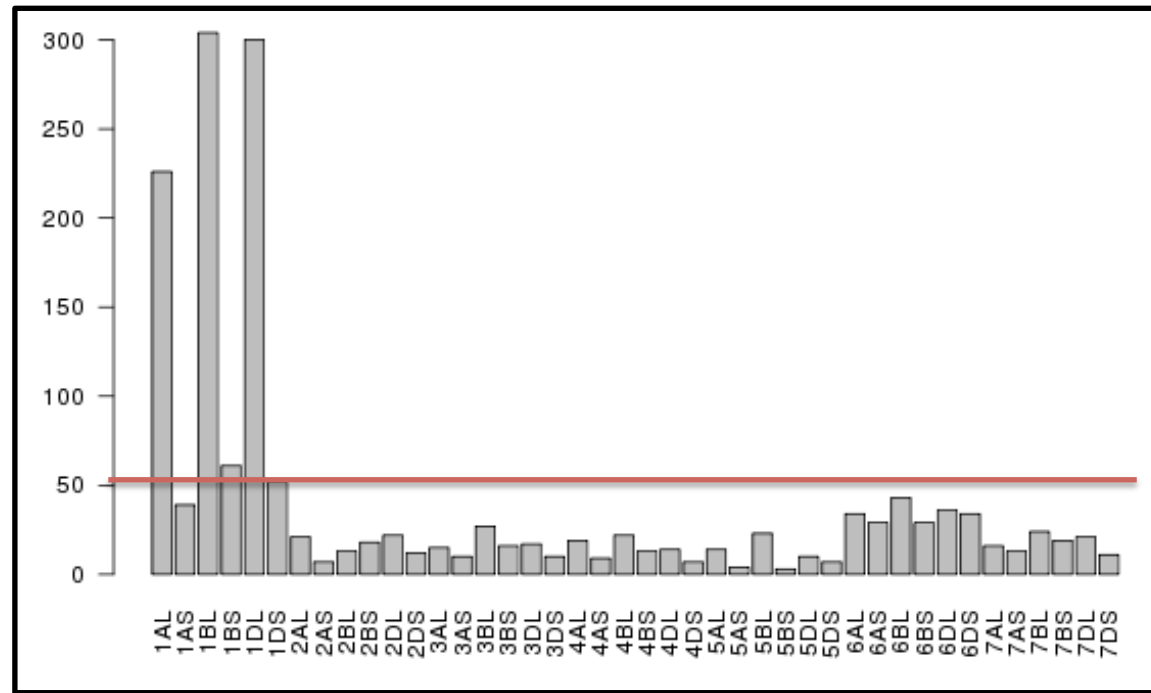
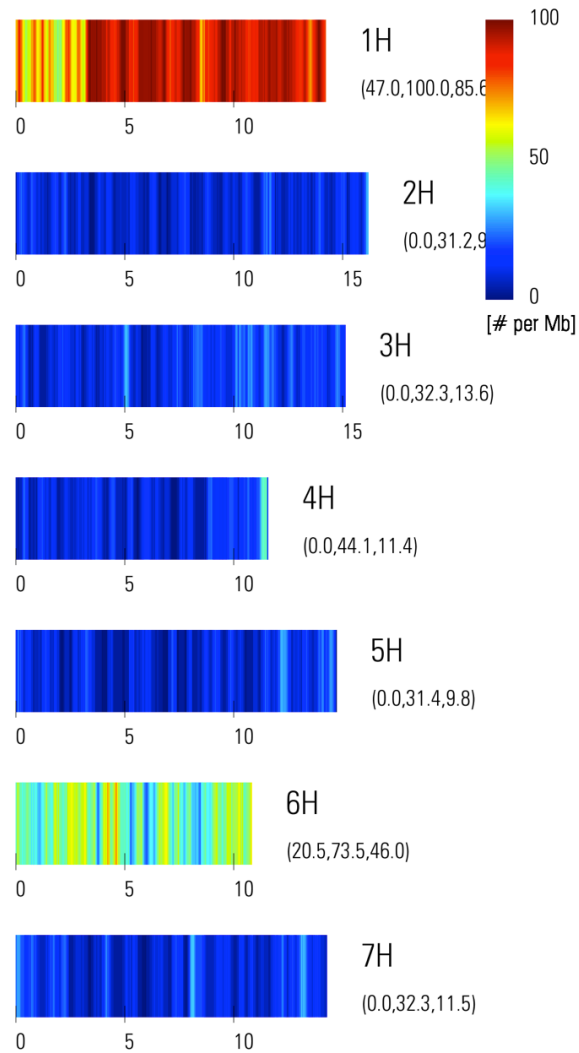
Genome Sciences Centre, British Columbia Cancer Agency, Vancouver, British Columbia V5Z 4E6, Canada

Assessing the assemblies for purity

Alignment of bin-mapped wheat ESTs to the repeat masked assemblies (Qi *et al*, 2004)



Contamination in assemblies? (1BL)



Aligning bin-mapped wheat ESTs to contigs (TGAC)

Aligning contigs to the virtual barley genome
(Mihaela Martis, MIPS)

Contamination Assessment

11 assemblies contaminated

Action taken

- Regenerate flow-sorted chromosomes from wheat DNA
- Remake the libraries
- Resequence
- Reassemble

This improved 7 assemblies, leaving 4 problematic ones (1AL, 1BL, 5DL, 6BS)
TGAC developed a kmer-based cleaning approach to generate clean assemblies for these arms – see Jon Wright's presentation on Monday

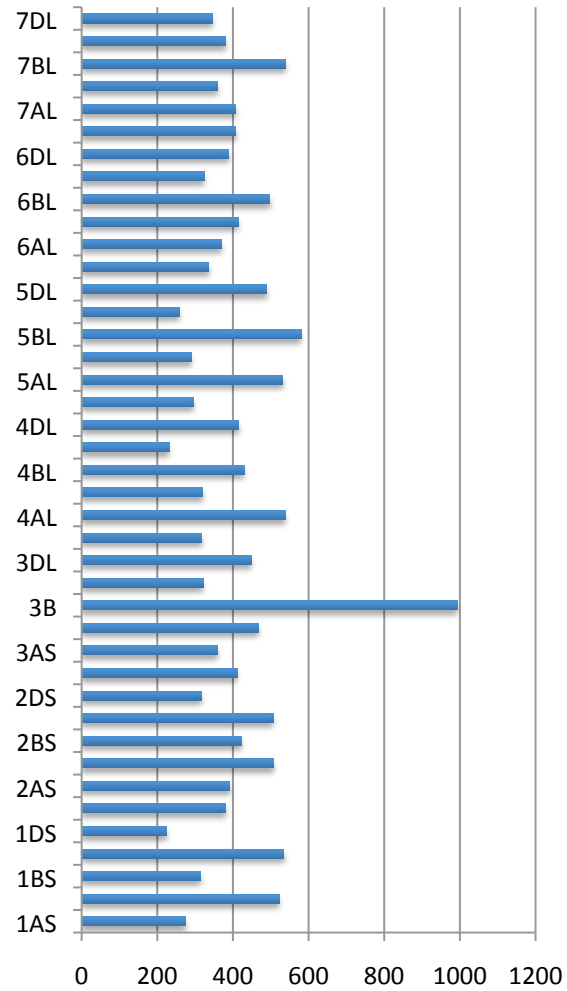
Summary Statistics for whole genome

Assembly statistics

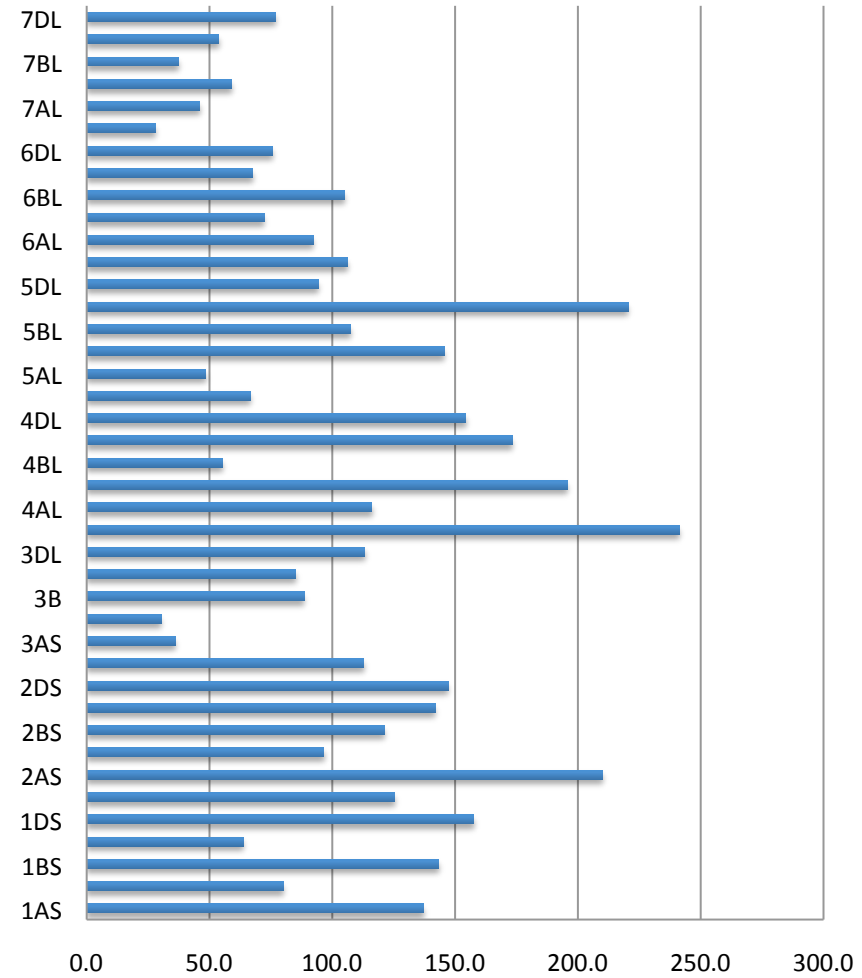
- Average GC content = **45.0%**
- N50 contig length (after filtering) = **2.4 kbp**
- Estimated gene count (based on hits to a barley gene set provided by IBGSC)
 - **1,526** (average per short arm)
 - **2,460** (average per long arm)
 - Total **83,977**

CSS Assembly Analysis

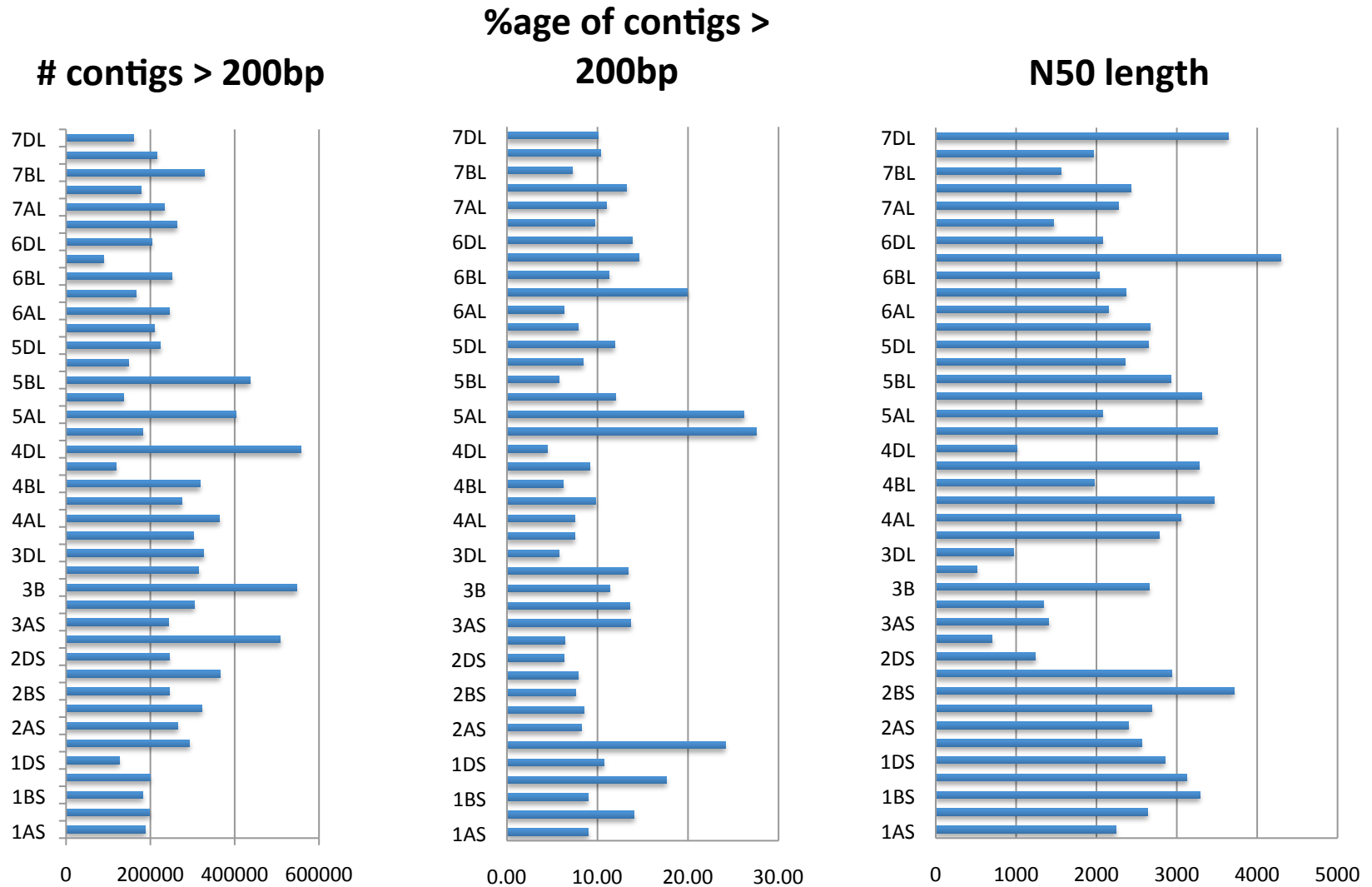
Chromosome size



Chromosome arm coverage



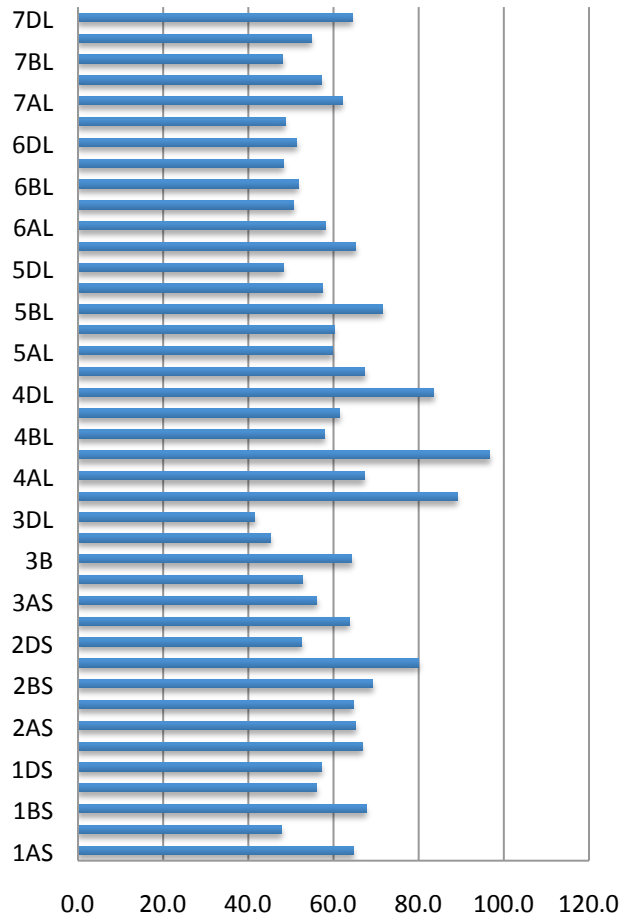
CSS Assembly Analysis – assembled contigs



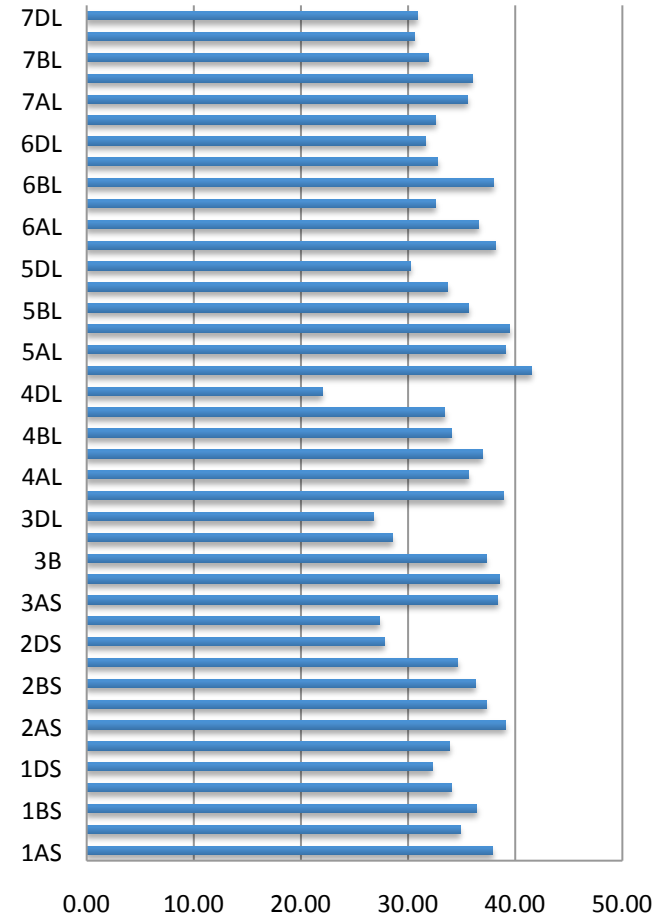
CSS Assembly Analysis

- chromosome coverage estimates

%age of arm represented



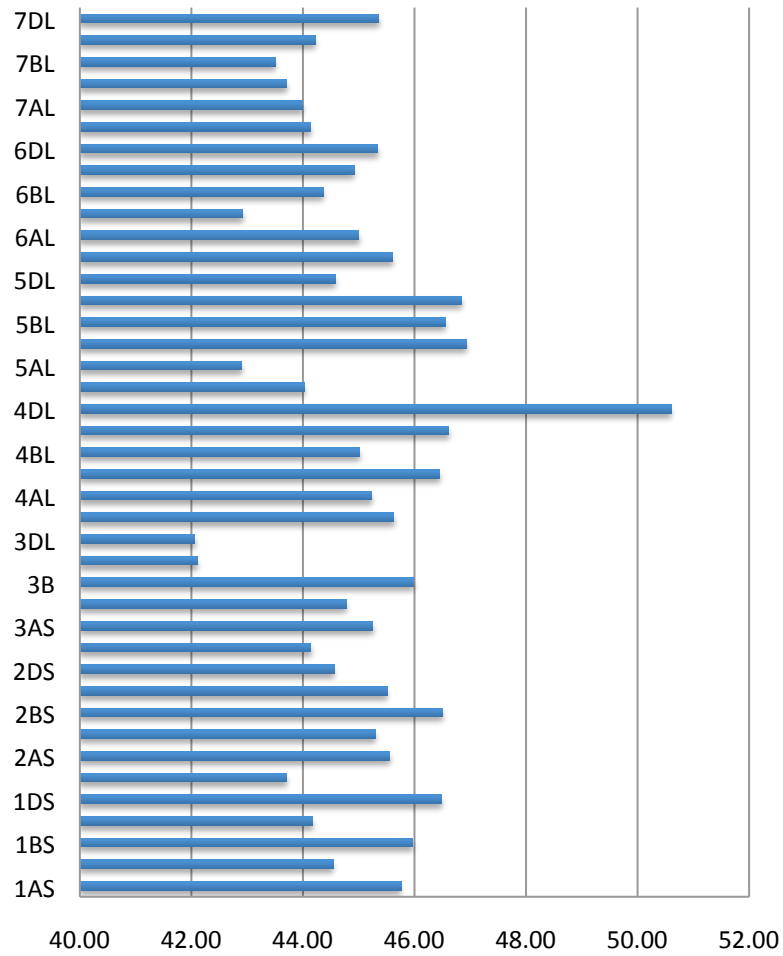
%age of bases masked



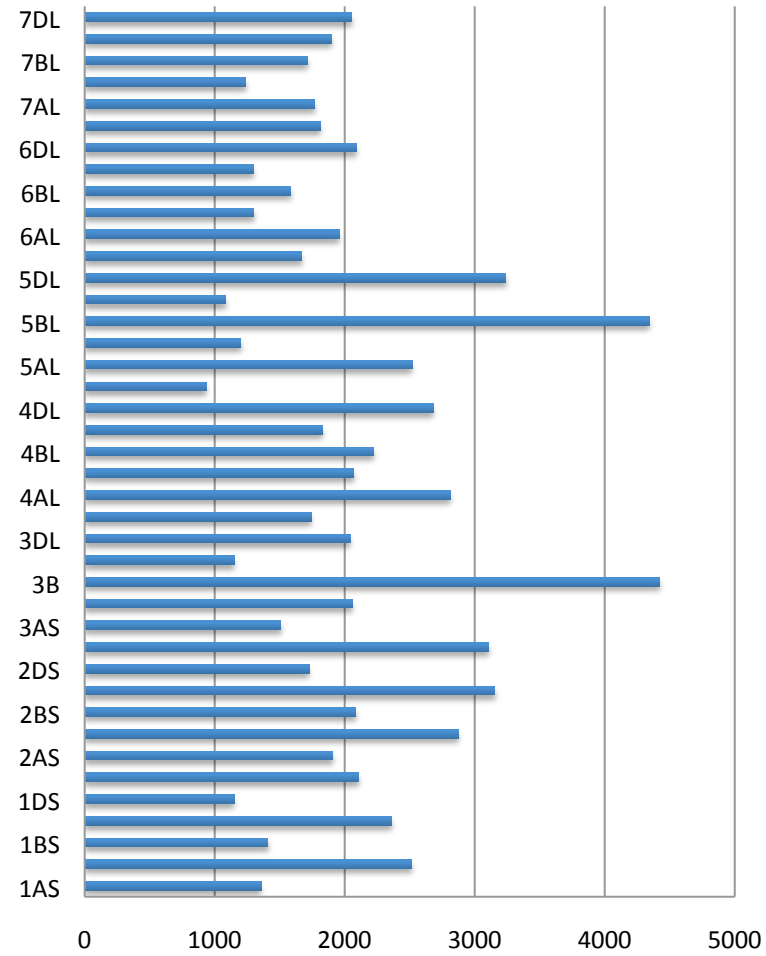
CSS Assembly Analysis

– base composition and gene content

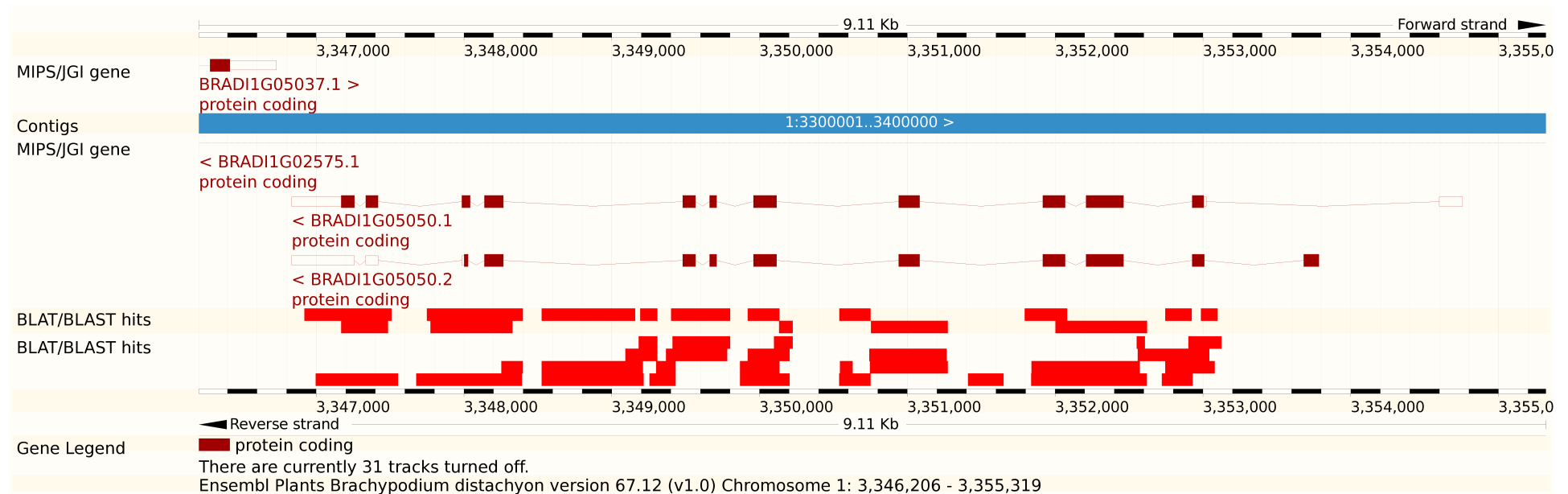
GC content



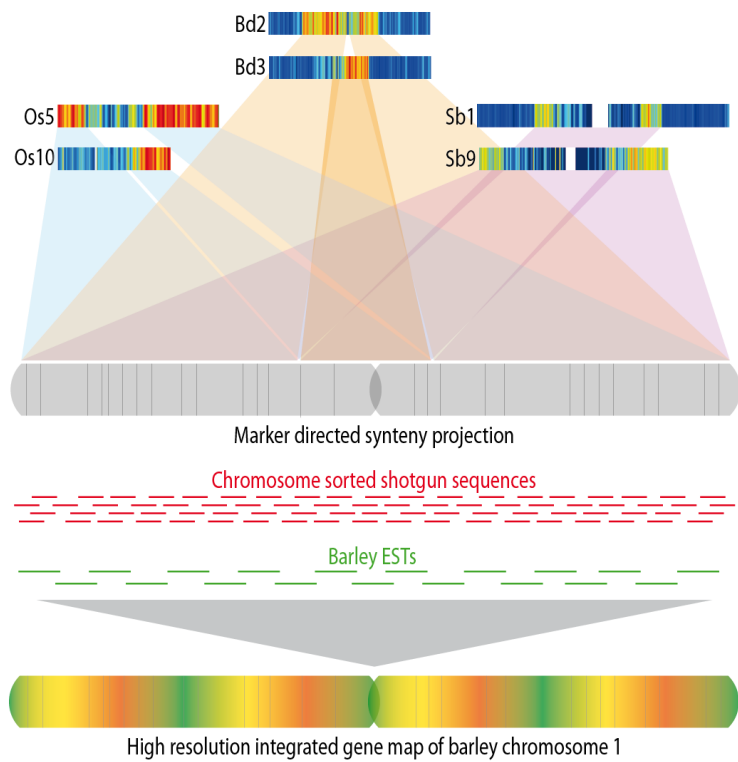
Gene content



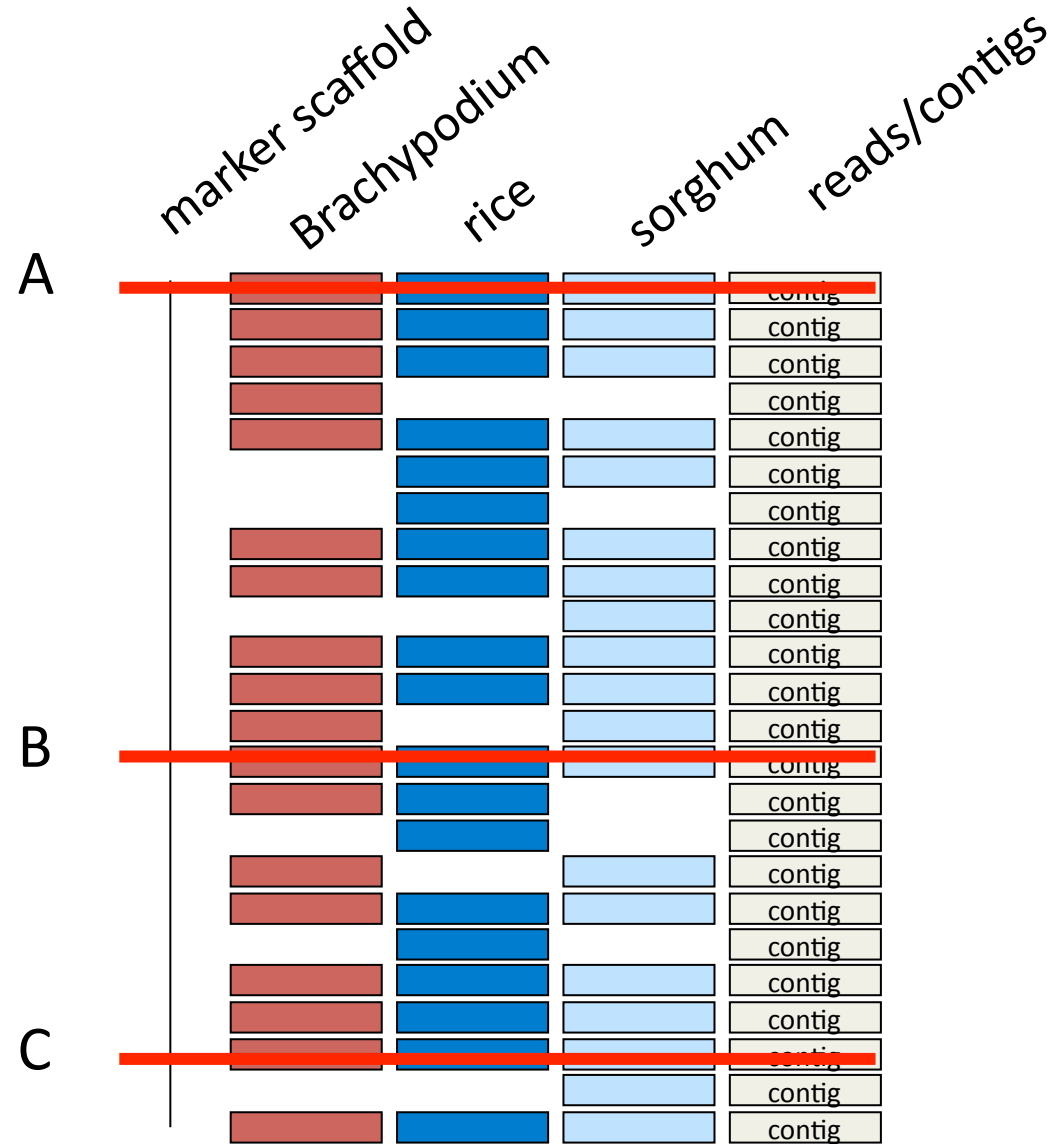
Alignment of assembled contigs with Brachypodium – an Ensembl browser view



Stage 2: GenomeZipper + Virtual Gene Map: Syntenic Integration



Klaus Mayer
 Mihaela Martis
MIPS



Stage 3: Survey Sequence Repository at URGI

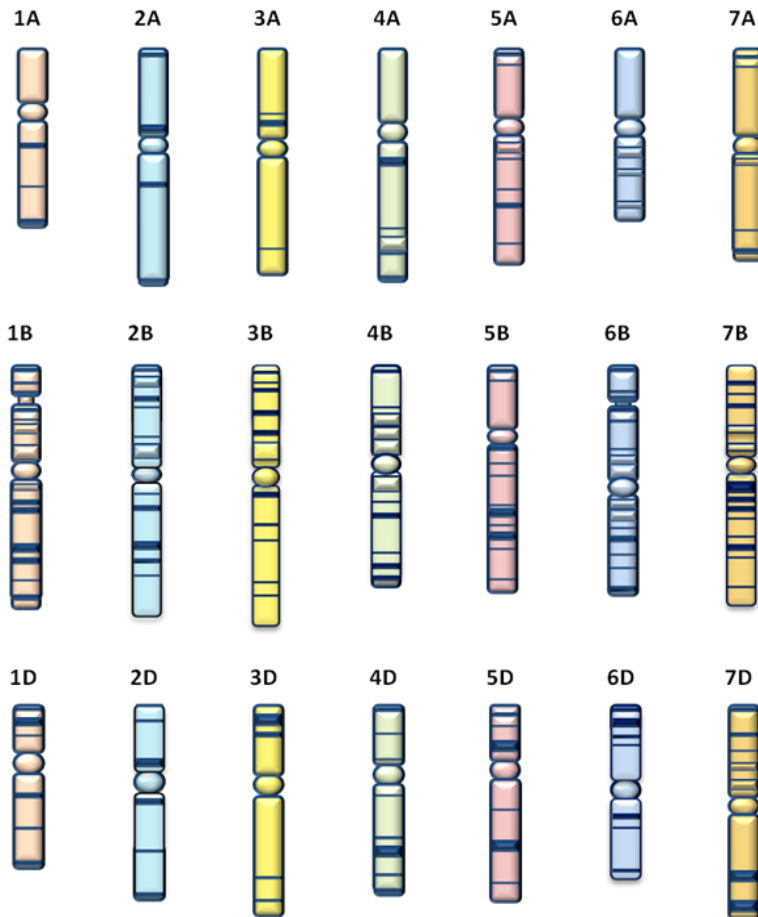
Sequence Repository



International
Wheat Genome
Sequencing
Consortium

Click on a chromosome to access the survey sequence and the viewers.

- [Process](#) to create an account to access the download and BLAST.
- [BLAST](#) direct link (registered access).
- [FAQ](#) section.
- **News:** All chromosomes are now available for download and BLAST.



Click on a chromosome to have access to the survey sequence with **download**, **blast search** and **viewers**.

WSS leaders can **download** and **blast** and non-members can blast and download the contigs matching hits once they agree to a data release policy.

All assemblies are available in the repository

<http://urgi.versailles.inra.fr/Species/Wheat/Sequence-Repository>

Applications of the survey sequence

The survey sequence provides a very fragmented view of the sequences of individual chromosome arms. It does not provide a true representation of the structure of the chromosomes but it does enable:

- Annotation of genes within contigs (intron-exon structure)
- Some limited annotation of 'pseudogenes'
- Analysis of coding variants
- *In silico* mapping of markers (e.g. genetic markers, trait markers) to chromosomes within sub-genomes
- Implement localised synteny studies
- Obtain estimations for
 - coding genes
 - lineage specific genes
 - comparative analysis of homoeologous genes

Next Steps

- Full analysis of data for publication – target submission date autumn 2012
- On publication read data and assemblies will be available from GenBank / EBI / DDBJ repositories in addition to sites with added value, e.g. URGI
- Assembly improvement options include:
 - incorporation of 454 data (chr arm and whole genome)
 - incorporation of BAC end data, where available
 - addition of new data sets , e.g. chr arm mate pair data
- Improving the utility of data, e.g. visualisation options include providing synteny alignments of contigs with other genomes, e.g. Brachypodium and barley

What does the community need?

Acknowledgments

TGAC

Jon Wright

Mario Caccamo

Sarah Ayling

Kirsten McLay

Melanie Febrer

IEB ASCR

Jaroslav Dolezel

Hana Simkova

MIPS

Klaus Mayer

Mihaela Martis

URGI

Michael Alaux

IWGSC leadership

Kellye Eversole

Catherine Feuillet

