

From fragments to the whole: progress in wheat genome annotation

LUCA VENTURINI

Computational Biologist, Swarbreck group



Earlham
Institute

Decoding Living Systems



A timeline of wheat annotations

Draft of the first-generation wheat reference genome. Almost **one million gene models** annotated, of which **124,201** are classified as **high-confidence**

The **TGACv1 annotation** identifies 217,907 loci, of which **114,247** are high confidence.

Release of a hybrid assembly of *Ae tauschii*, increasing contiguity by 100X.

2012

Release of a draft gene centric assembly; the number of wheat genes is estimated to be around **96,000**.

2013

Drafts of the genomes of *T. urartu* and *Ae. tauschii*; annotated with **34,879** and **43,150 gene models**

2014

Chromosome 3B is released and annotated with **8055** genes, of which **5099** are coding

2016

Release of the IWGSC pseudo-molecule annotation expected later in the year

2017

The CSS assembly: a first step towards a reference sequence

Draft of the first-generation wheat reference genome. Almost **one million gene models** annotated, of which **124,201** are classified as **high-confidence**

2013

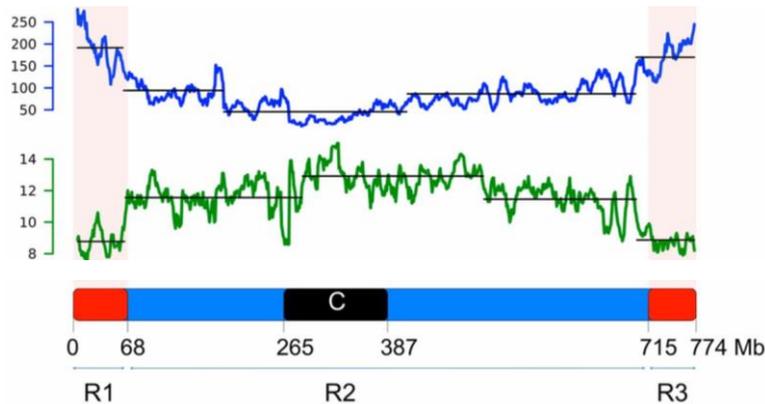
The first draft of the wheat genome, published in 2014, allowed for the first time to anchor a gene catalogue onto its chromosomal locations. It allowed for the first time to analyse genes in terms of their genomic context.

However, the assembly was still not optimal:

- Genes were fragmented (over 1 million low and high confidence genes)
- Out of the 133,090 high confidence genes, 124,201 (93%) could be assigned to a genomic location. However, only 44% were identified as likely full length.

The importance of a contiguous assembly: chr3B

- Having a pseudomolecule, it was possible to define loci completely, and categorise them
- Each gene is linked to its genomic context, allowing to analyse how genes segregate together
- Each gene is linked to long-distance markers, helping in GWAs and breeding



Gene density (in blue) and expression density (in green) along the chr3B pseudomolecule.

	All	Full genes	Pseudogenes
No. of genes	7264	5326	1938
Average size (bps) of coding sequences (\pm standard deviation)	1095 \pm 807	1187 \pm 821	840 \pm 710
Average number of exons (\pm standard deviation)	4.2 \pm 4.4	4.4 \pm 4.6	3.6 \pm 3.8
Gene density (kb^{-1})	107	145	400
No. of expressed genes	5185	4125	1060
% genes with alternative splicing	61	63	56

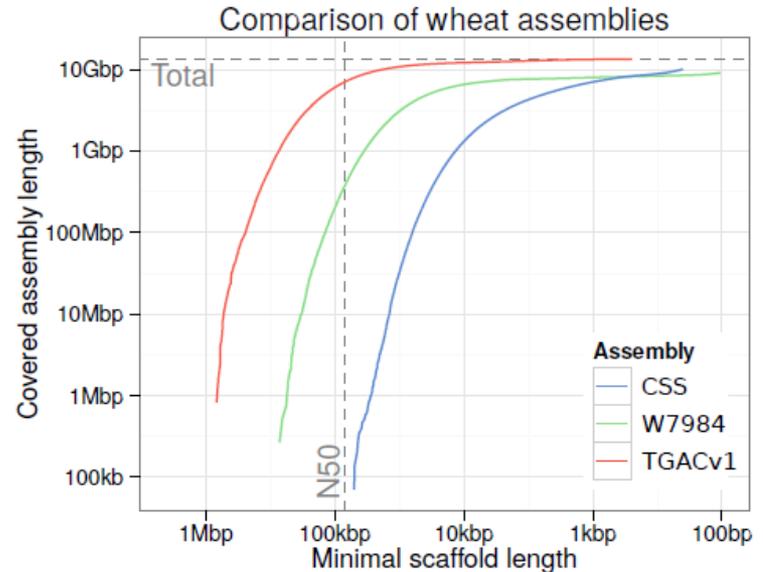
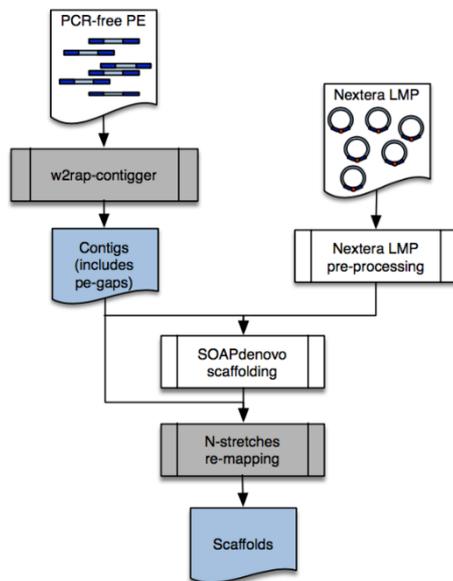
Choulet et al, Science 2014

The TGACv1 assembly

The Earlham Algorithm Development team
Bernardo Clavijo, Gonzalo Garcia Accinelli
Jon Wright



<https://github.com/bioinfologics/w2rap>



- At EI, we developed novel library preparation methods and novel algorithms to assemble genomes *de novo* quickly and reliably.
- Our assembly captured 60% more of the genomic content, compared with the previous best effort
- The method is completely open, reproducible, and fast.

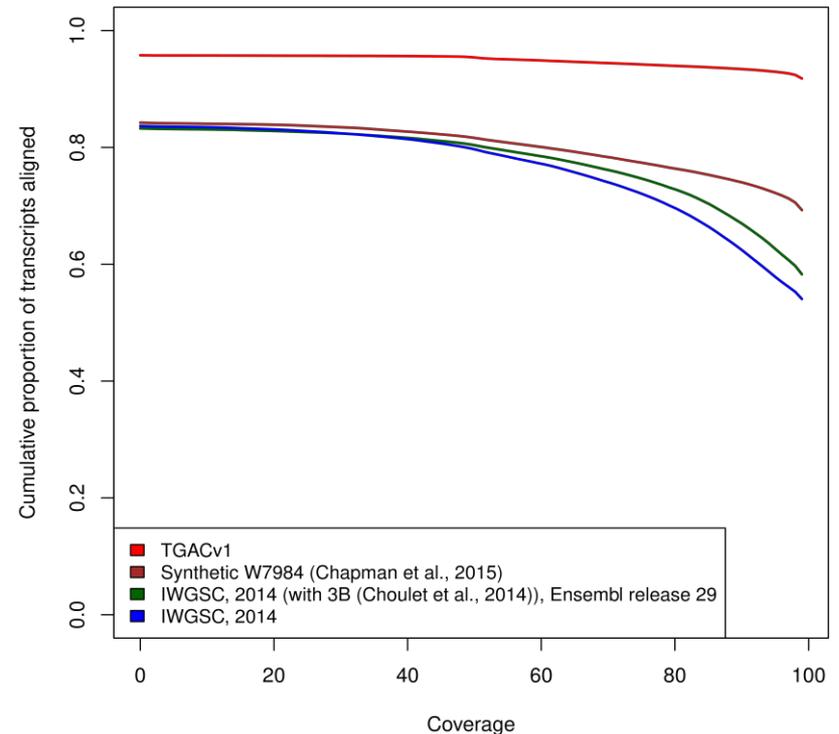
Generation of high-quality input transcript data for genome annotation

The TGACv1 annotation was based on the following sets of data:

- Two billion Illumina reads from public available datasets
- 800 million long, strand-specific Illumina dataset
- Over one million and a half PacBio full length cDNAs
- Protein models from six different species

	# of sequences	Notes
Public Illumina reads	2,409,760,971	
Internal Illumina dataset	824,241,135	Strand-specific, 250bp PE
IsoSeq	1,509,322	
Protein models	316,385	Six different species

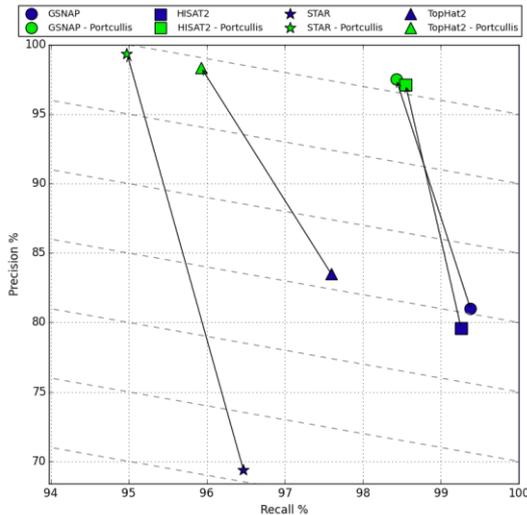
The high proportion of PacBio transcripts aligned to the TGACv1 assembly indicates an excellent representation of the gene space.





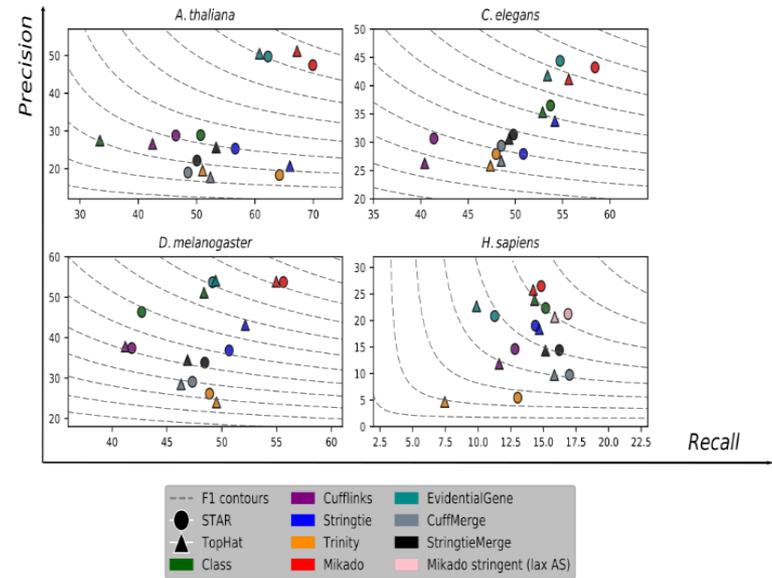
<https://github.com/maplesond/portcullis>

- It allows to distinguish between real and artifactual splicing junctions
- Especially important in datasets with deep sequencing and with polyploidy
- Machine-learning based

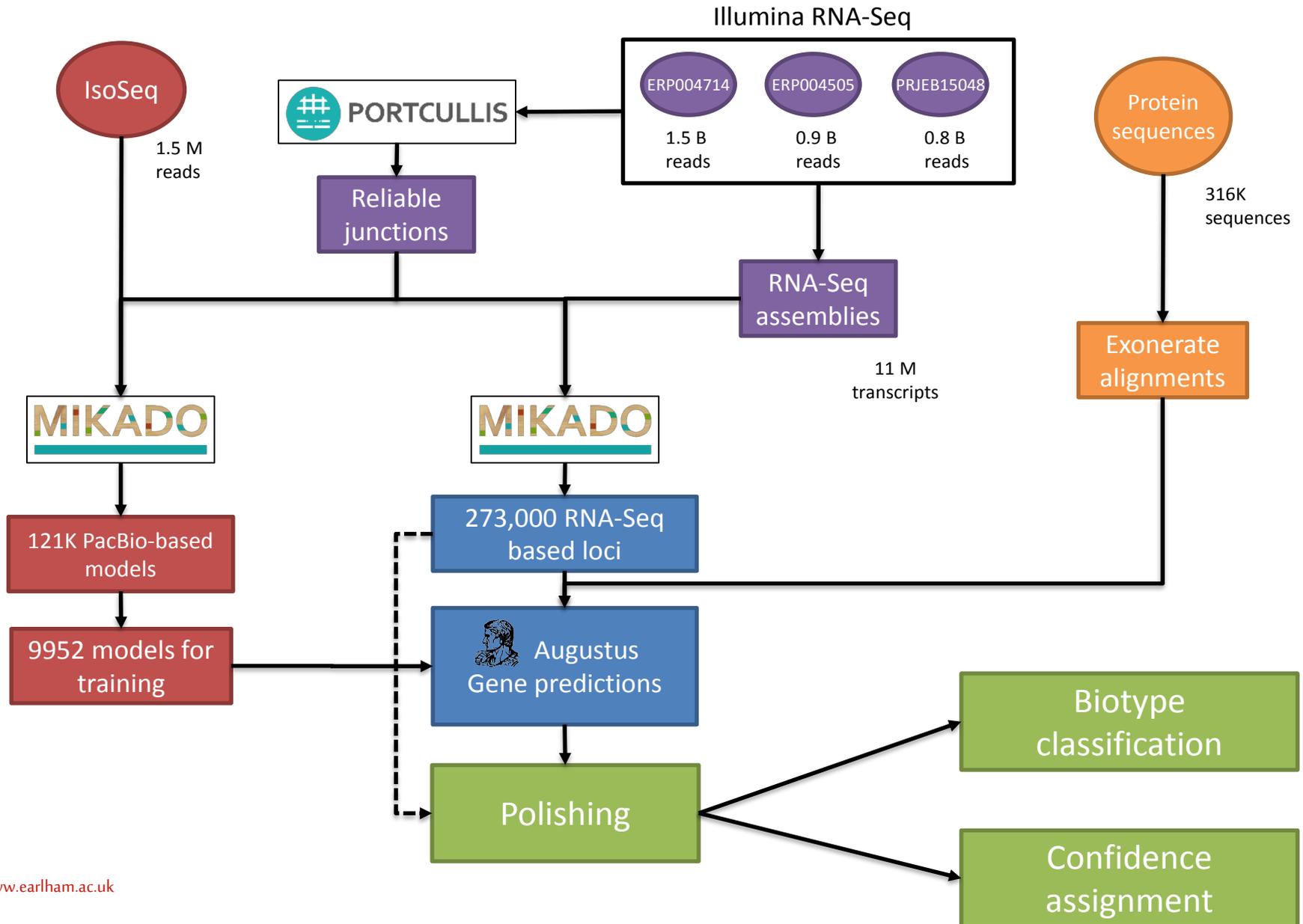


<https://github.com/luventurini/mikado>

- Scores transcripts on the basis of intrinsic and extrinsic features (eg. Portcullis junctions)
- Robustly integrate multiple RNA-Seq assemblies
- Detects and resolves chimeric transcripts



ANNOTATION PIPELINE



Gene classification and confidence assignment

- We classified each model in two complementary ways:
- Verifying their support using known protein sequences (protein or homology rank)
- Verifying their concordance with RNA-Seq data (transcript or structural rank)

Rank level	Protein	Transcript
1	Over 80% homology	Full support from PacBio models
2	60-80% homology	Full support from Illumina models
3	30-60% homology	Structural congruence greater than 50%
4	Lower than 30% homology	Structural congruence lower than 50%
5	No homology with known proteins	No transcriptomic support

Confidence rankings of coding transcripts

Transcript count	protein rank	transcript rank
66404	P1	T1
43423	P1	T2
20937	P1	T3
10013	P1	T4
21469	P1	T5
3461	P2	T1
3545	P2	T2
3392	P2	T3
2084	P2	T4
6213	P2	T5
1813	P3	T1
4521	P3	T2
3995	P3	T3
3406	P3	T4
12210	P3	T5
781	P4	T1
3116	P4	T2
2846	P4	T3
2494	P4	T4
7484	P4	T5
2079	P5	T1
4638	P5	T2
3944	P5	T3
2915	P5	T4
12364	P5	T5

	High confidence
	Low confidence

Results: transcripts and gene classification

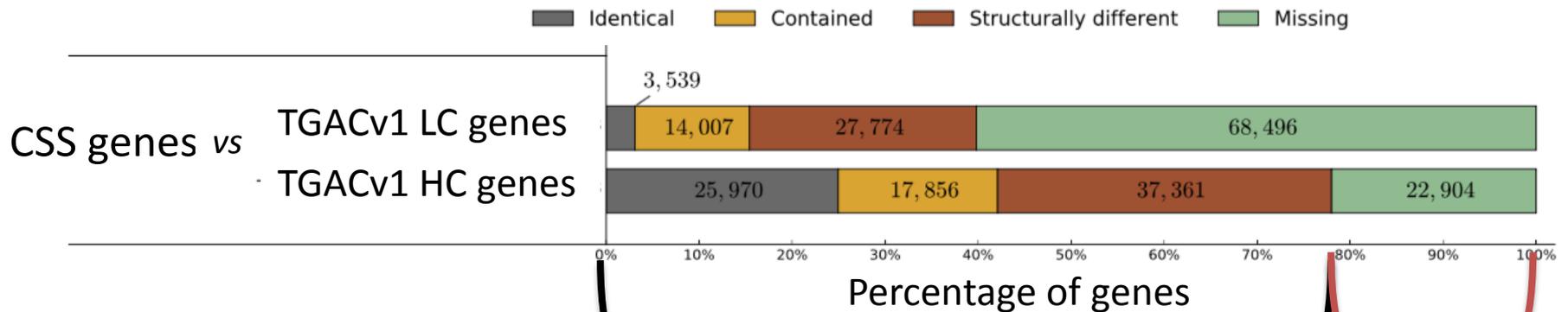
	All TGAC Models	mRNA HC	mRNA LC	ncRNA HC	ncRNA LC	Repeat-associated
Genes	217,907	104,091	83,217	10,156	9,933	10,510
Transcripts	273,739	154,798	85,778	11,591	10,438	11,134
Transcripts per gene	1.26	1.49	1.03	1.14	1.05	1.06
Transcript mean cDNA size (bp)	1,766.12	2,119.52	1,304.53	1,368.24	1,083.98	1,462.71
Exons per transcript	4.48	5.83	2.8	2.58	2.76	2.27
Exon mean size (bp)	394.15	363.73	465.27	530.25	392.24	644.09
Transcript mean CDS size (bp)	1,165.52	1,361.82	839.97	-	-	891.05
Mono-exonic transcripts	60,322	19,034	30,479	3,061	3,044	4,704
	22.04%	12.30%	35.53%	26.41%	29.16%	42.25%
Genes with alternative splicing	32,616	28,608	2,033	1,037	460	478
	14.97%	27.48%	2.44%	10.21%	4.63%	4.55%

The final set of TGACv1 annotations comprises 217,907 loci, of which 104,091 are classified as high-confidence protein coding genes. Compared with the previous CSS assembly, therefore, our annotation displays:

- A similar number of high confidence genes
- A much decreased number of low-confidence genes, many of which will probably be characterised as pseudogenes in the future
- The explicit characterization of long non-coding RNAs, which were absent from previous catalogues.

A more comprehensive and accurate wheat annotation

- We aligned the CSS/3B (IWGSC) gene models to the TGACv1 assembly and compared against the TGACv1 gene models.



78% of TGACv1 genes overlapped with the full set of CSS genes (LC + HC).

- 6665 (29%) are fully supported by PacBio
- 19 810 (86%) have cross species protein similarity support

Data availability

The reference sequence and annotation can be retrieved at:

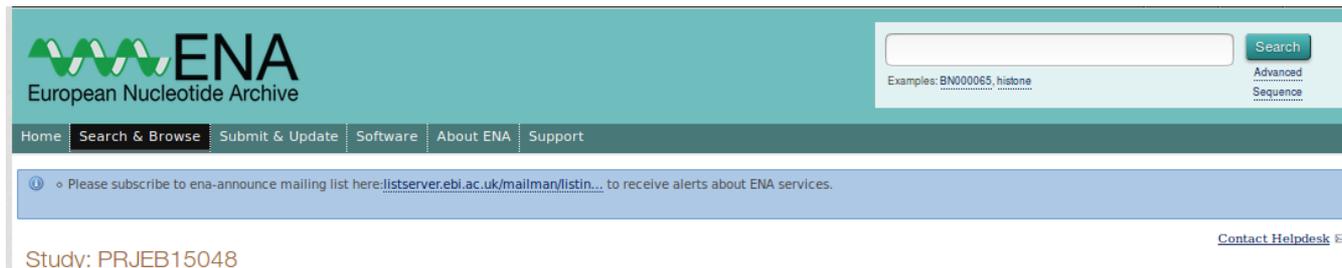


http://opendata.earlham.ac.uk/Triticum_aestivum/TGAC/v1/annotation/

http://plants.ensembl.org/Triticum_aestivum/Info/Index

The RNA sequencing reads can be downloaded from ENA (project PRJEB15048):

<http://www.ebi.ac.uk/ena/data/view/PRJEB15048>



An improved assembly and annotation of the allohexaploid wheat genome identifies complete families of agronomic genes and provides genomic evidence for chromosomal translocations.

Bernardo J. Clavijo^{1*}, Luca Venturini^{1*}, Christian Schudoma¹, Gonzalo Garcia Accinelli¹, Gemy Kaithakottil¹, Jonathan Wright¹, Philippa Borrill², George Kettleborough¹, Darren Heavens¹, Helen Chapman¹, James Lipscombe¹, Tom Barker¹, Fu-Hao Lu², Neil McKenzie², Dina Raats¹, Ricardo H. Ramirez-Gonzalez¹, Aurore Coince¹, Ned Peel¹, Lawrence Percival-Alwyn¹, Owen Duncan³, Josua Trösch³, Guotai Yu², Dan Bolser⁴, Guy Namaati⁴, Arnaud Kerhornou⁴, Manuel Spannagl⁵, Heidrun Gundlach⁵, Georg Haberer⁵, Robert P. Davey^{1,6}, Christine Fosker¹, Federica Di Palma^{1,6}, Andrew Phillips⁷, A. Harvey Millar³, Paul J. Kersey⁴, Cristobal Uauy², Ksenia V. Krasileva^{1,6,8}, David Swarbreck^{1,6+}, Michael W. Bevan²⁺ and Matthew D. Clark^{1,6+}.

**contributed equally to this work*

+corresponding authors

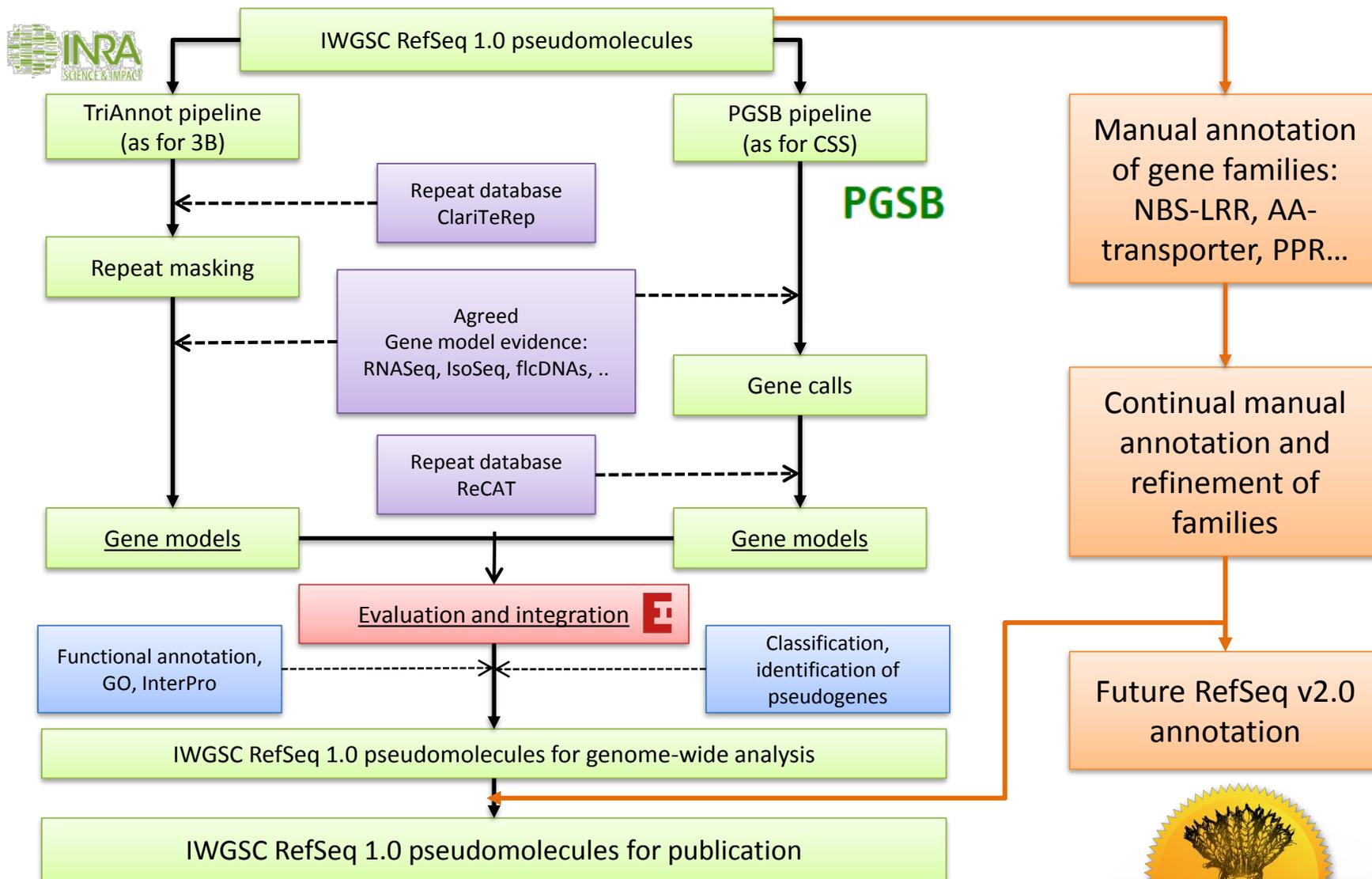
Manuscript available on [biorXiv](https://www.biorxiv.org/), and accepted for publication in Genome Research

IWGSC RefSeq v1.0 Gene Prediction Strategy

- Coordination and oversight – Jane Rogers, IWGSC
- Two annotation teams:
 - INRA-GDEC – Frédéric Choulet, Hélène Rimbart, Philippe Leroy
 - PGSB – Sven Twarzdiok, Klaus Mayer, Manuel Spannagl
- Evaluation and integration team
 - Earlham Institute – David Swarbreck, Luca Venturini, Gemy Kaithakottil



IWGSC RefSeq v1.0 Annotation Approach



Characteristics of the two annotations

	PGSB		INRA			
	PGSB HC	PGSB All	INRA HC	INRA LC	INRA Pseudo	INRA all
Number of genes	104,696	205,643	65,884	41,342	73,044	180,270
Number of transcripts	297,971	432,097	65,884	41,342	73,044	180,270
Number of monoexonic genes	24,231	88,313	23,677	18,683	34,501	76,861
Average transcripts per gene	2.85	2.10	1.00	1.00	1.00	1.00
Average CDS length	1,384	1,145	1,110	1,231	766	998
CDS exons per transcript	6.32	4.96	4.08	4.04	2.84	3.57

The pipelines underlying the two annotations utilise different approaches and input data:

- The **PGSB** annotation:
 - is an alignment-driven approach (based on protein and wheat transcriptome data)
 - Captures splice variants
 - Classifies genes as high and low based on homology
- The **INRA** annotation:
 - is an evidence-guided approach (utilising gene predictions and aligned evidence)
 - provides an annotation of the UTR of transcripts
 - Identifies pseudogenes

Each approach will have specific strengths and weaknesses.

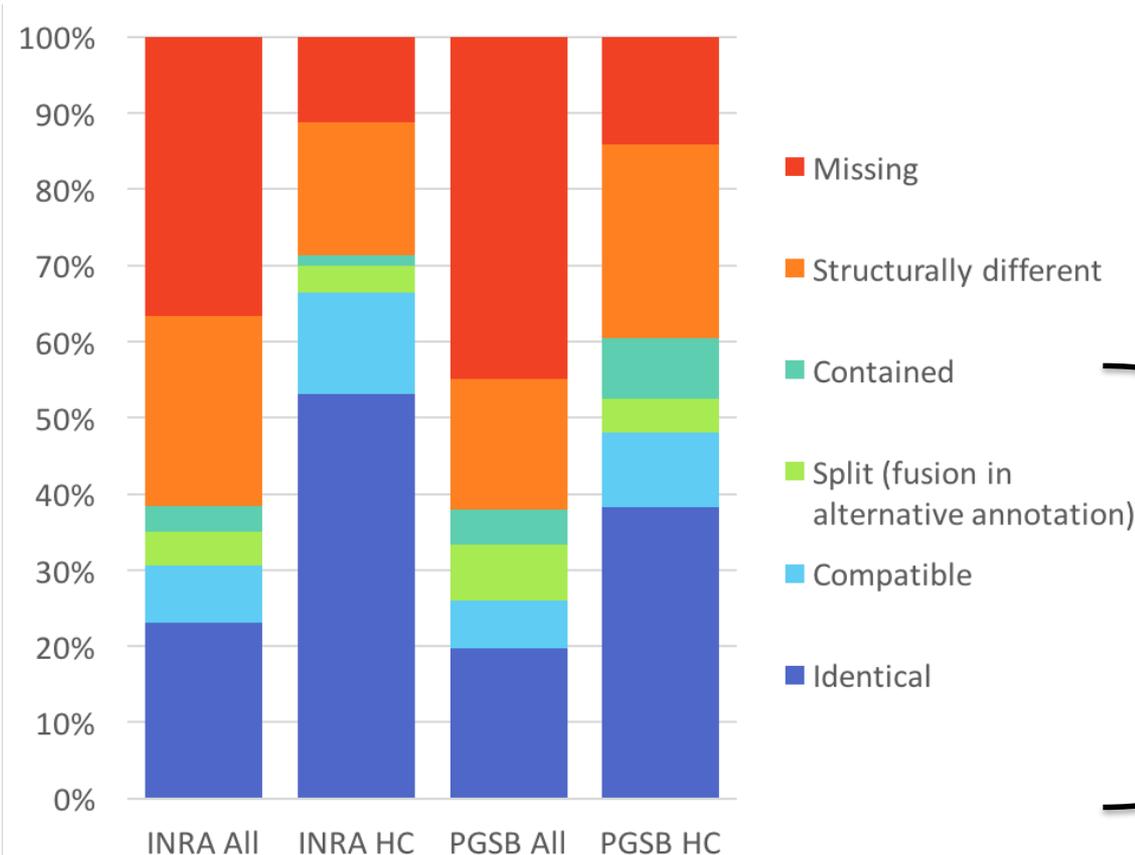


*preliminary results



Majority of high confidence genes are identified in both annotations

*preliminary results



Over 80% of the genes defined as high confidence by INRA/PGSB overlapped a gene in the alternative annotation

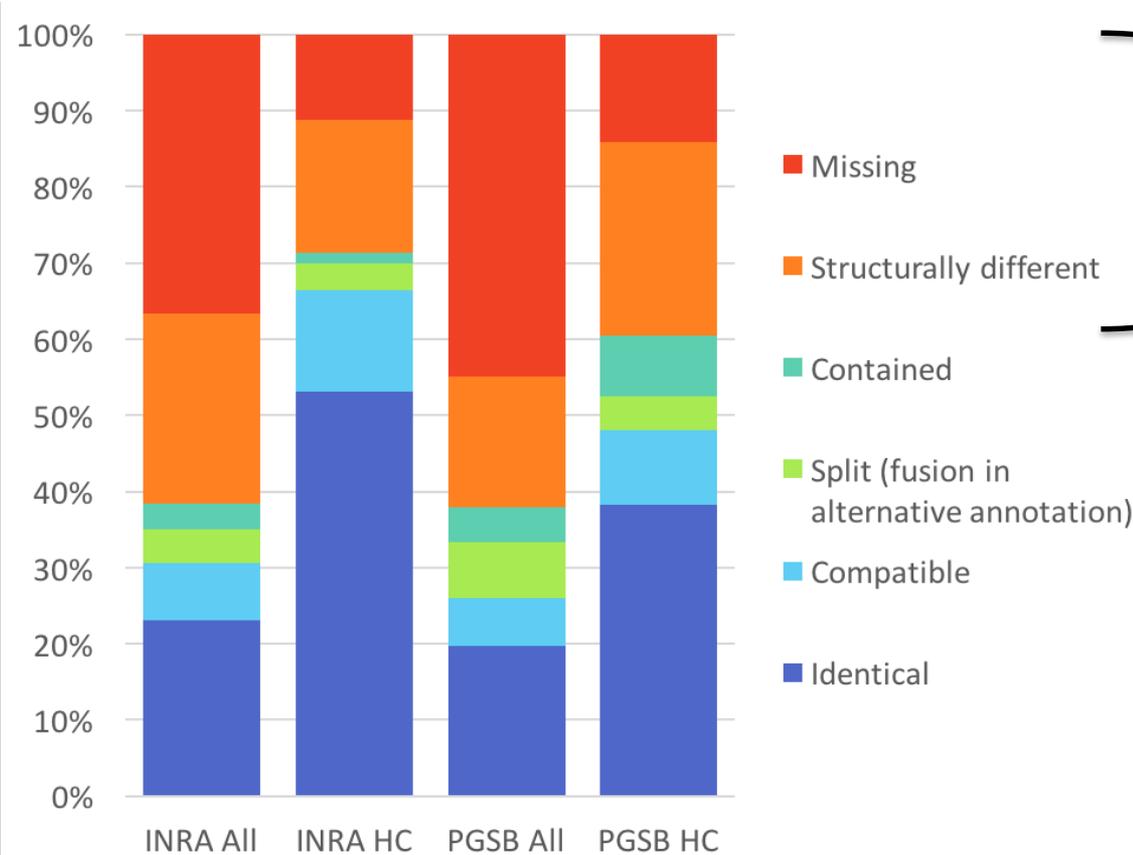
A substantial percentage of genes show highly similar structures for the high-confidence gene sets (in [blue](#) here) ..



*preliminary results



Some low confidence genes are present in only one annotation



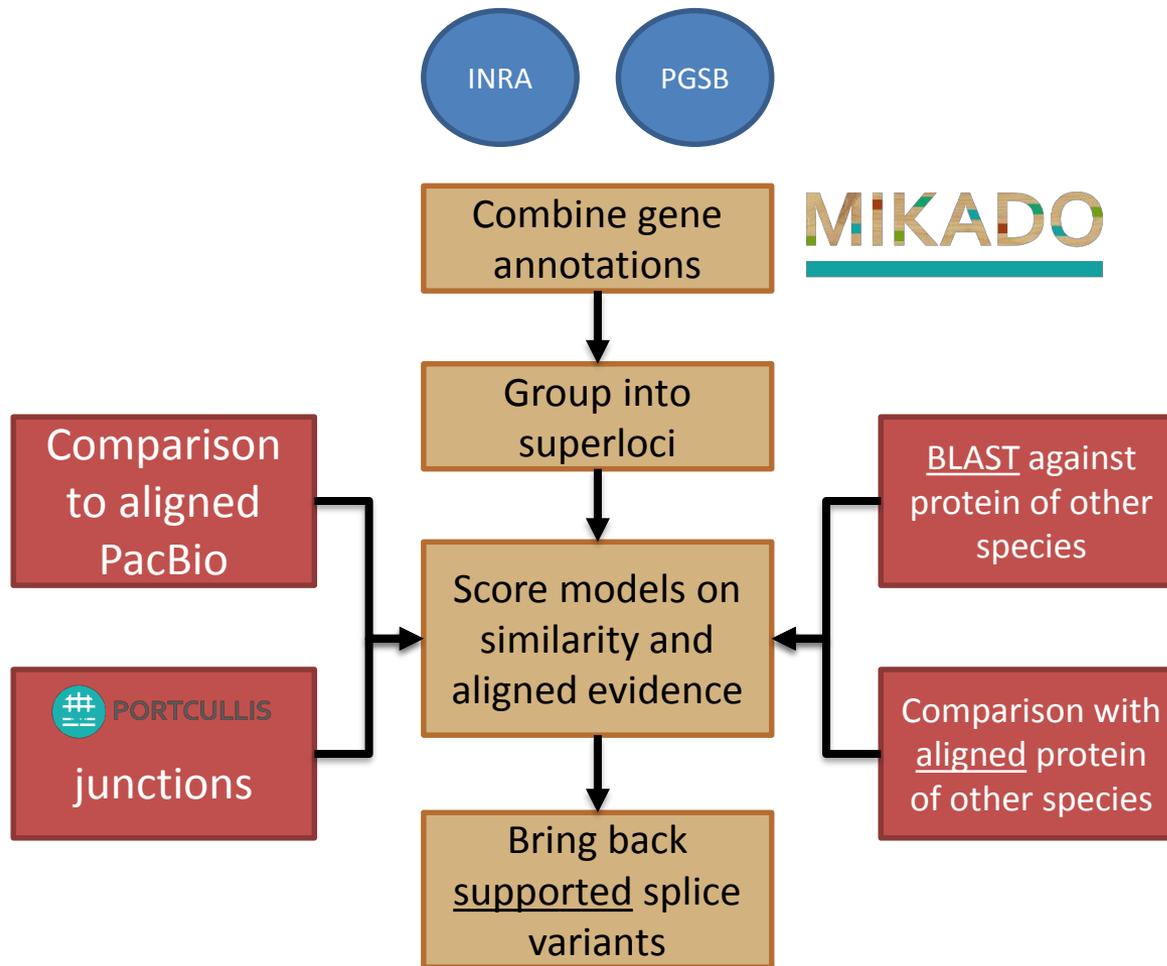
.. however, many low-confidence genes are present in only one of the two annotations (red fractions).

The differences between the two annotations reflect the challenge of annotating a transcriptionally complex polyploid species as well as differences in the datasets utilised and the annotation approaches (eg. sensitivity for detecting pseudogenes)

*preliminary results



An integrative approach to improve wheat annotation



By utilising an approach that selects gene models from across the two annotations we can exploit the strength of each annotation pipeline.

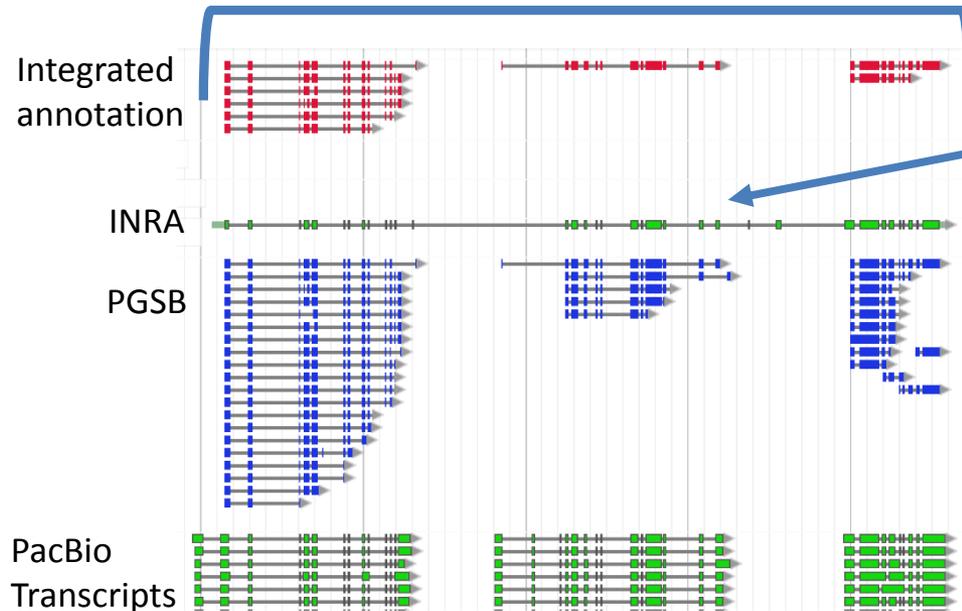
Gene models were selected based on the extent of evidence support and intrinsic gene characteristics

Including support from:

- PacBio transcripts
- High quality RNA-Seq assemblies
- Aligned proteins
- Validated junctions
- Sequence homology with known proteins

Integrating both annotations allows us to identify and resolve errors

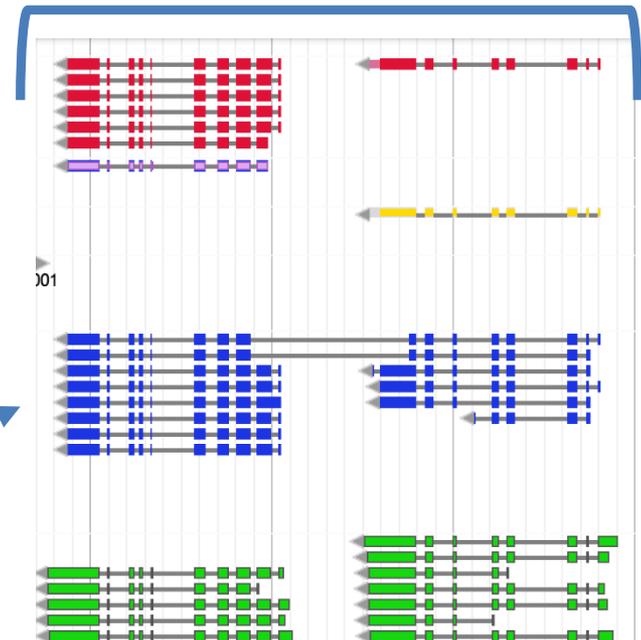
By “cherry picking” from the two annotations we can resolve issues where genes were incorrectly fused in one of the original annotations.



Correctly defined as 3 genes in the integrated gene set

Incorrect INRA fusion

Correctly defined as 2

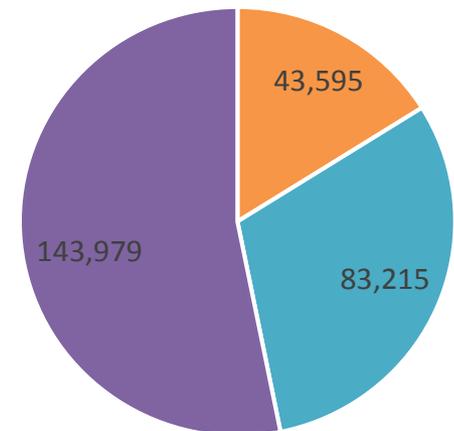


Incorrect PGSB fusion



Preliminary IWGSC v1.0 gene annotation results

	High Confidence	Low Confidence	All
Number of genes	110,790	158,793	270,789
Number of transcripts	137,056	162,011	299,076
Number of monoexonic genes	31,660	104,364	136,263
Average transcripts per gene	1.24	1.02	1.10
Average CDS length	1,323.77	604.3	934.01
CDS exons per transcript	5.27	1.86	3.42
Average CDS exon length	251.29	325.15	273.02
Average CDS intron length	476.92	799.06	538.81



- Derived from both annotations
- Derived from INRA
- Derived from PGSB

Both PGSB and INRA contributed significantly to the final integrated annotation

The integrated annotation combines the two annotation, removing redundancy and retaining the best models from both datasets.

The number of high confidence genes is similar to previous estimates, and to those found in the TGACv1 annotation.

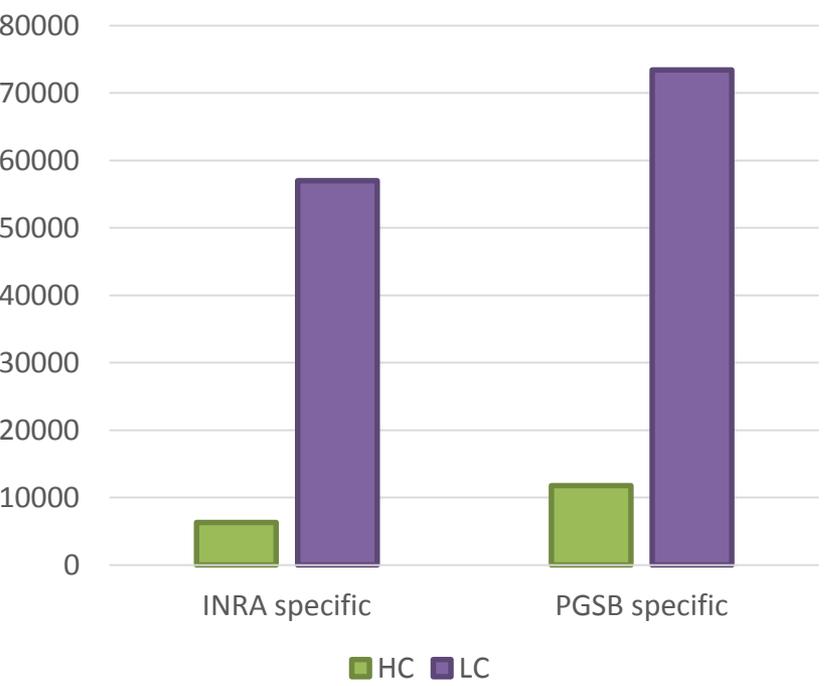


*preliminary results

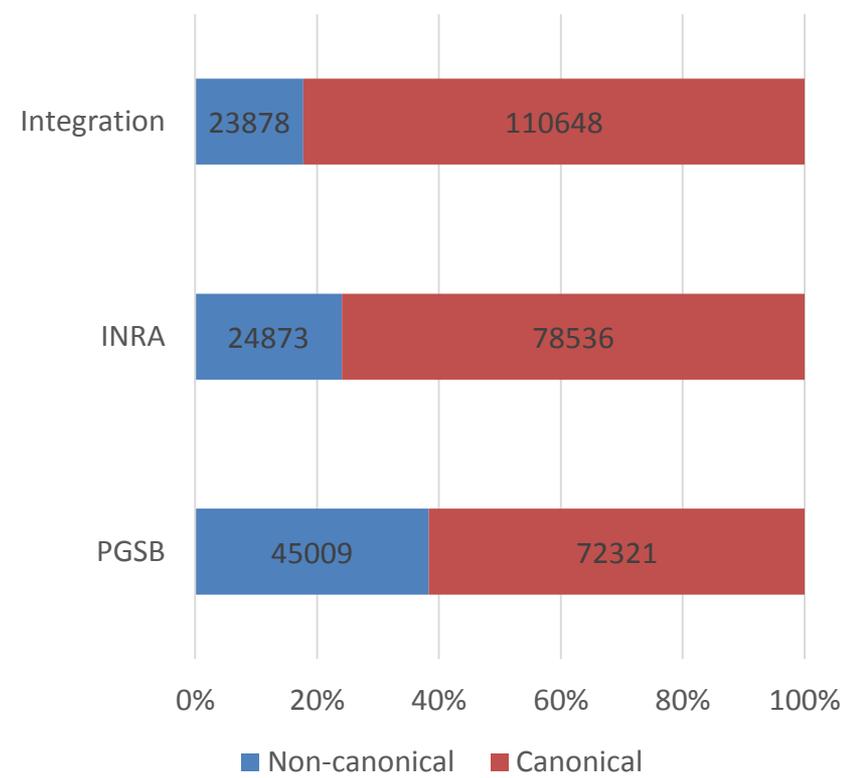


The integration provides a more comprehensive representation of genes and potential pseudogenes

- The combined annotation contains 6320 high-confidence genes that were absent in PGSB, and 11,780 high-confidence genes that were not represented in the INRA annotation.



The integration removes many potential incorrect splice variants

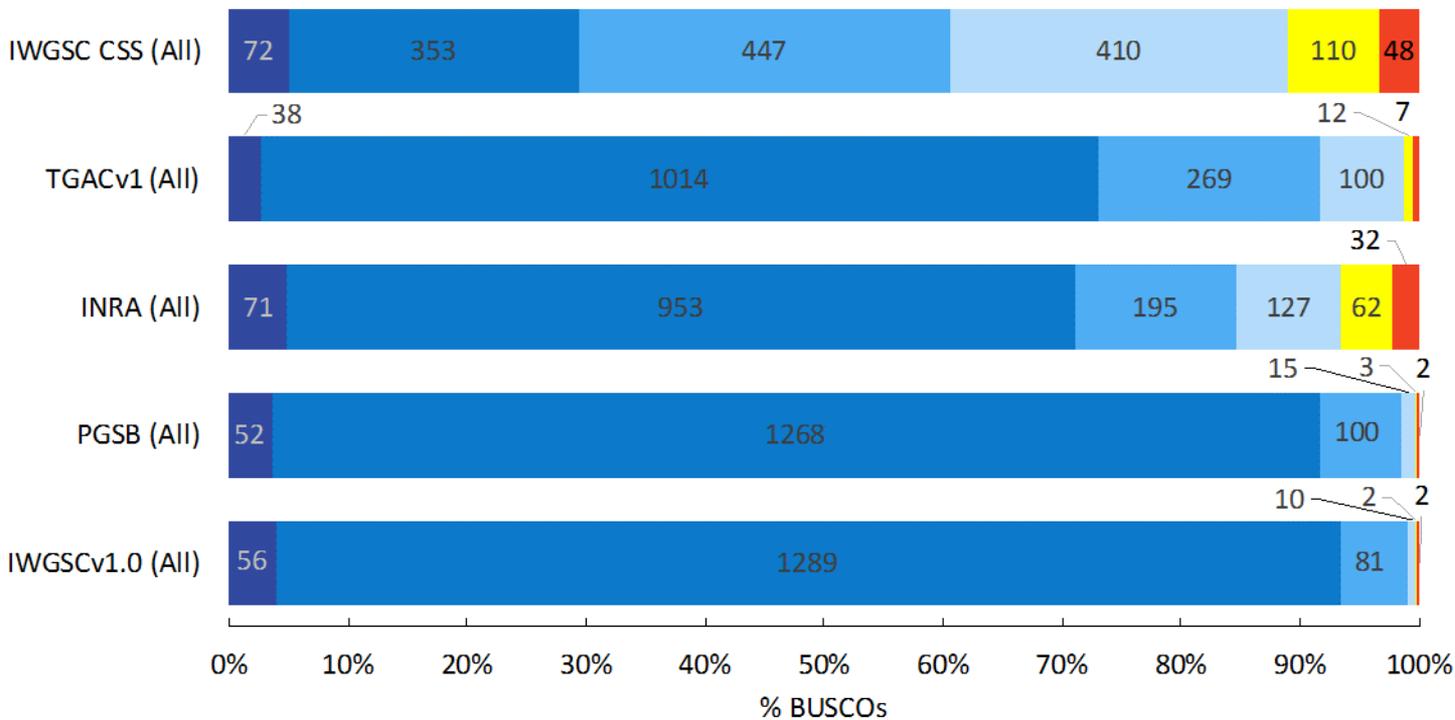
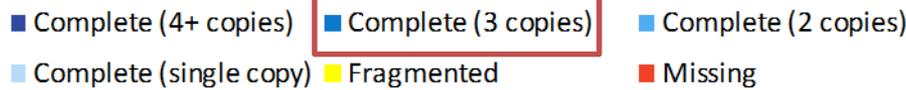


*preliminary results



As annotations improve, more of the triplet highly-conserved homeologs are captured

BUSCO Assessment results



An important quality check for gene annotations is to verify that they contain expected genes that are conserved across lineages – and when the organism is polyploidy, that the correct number of copies are present.

For wheat, the expectation is to find **three copies** for each of these highly conserved genes.

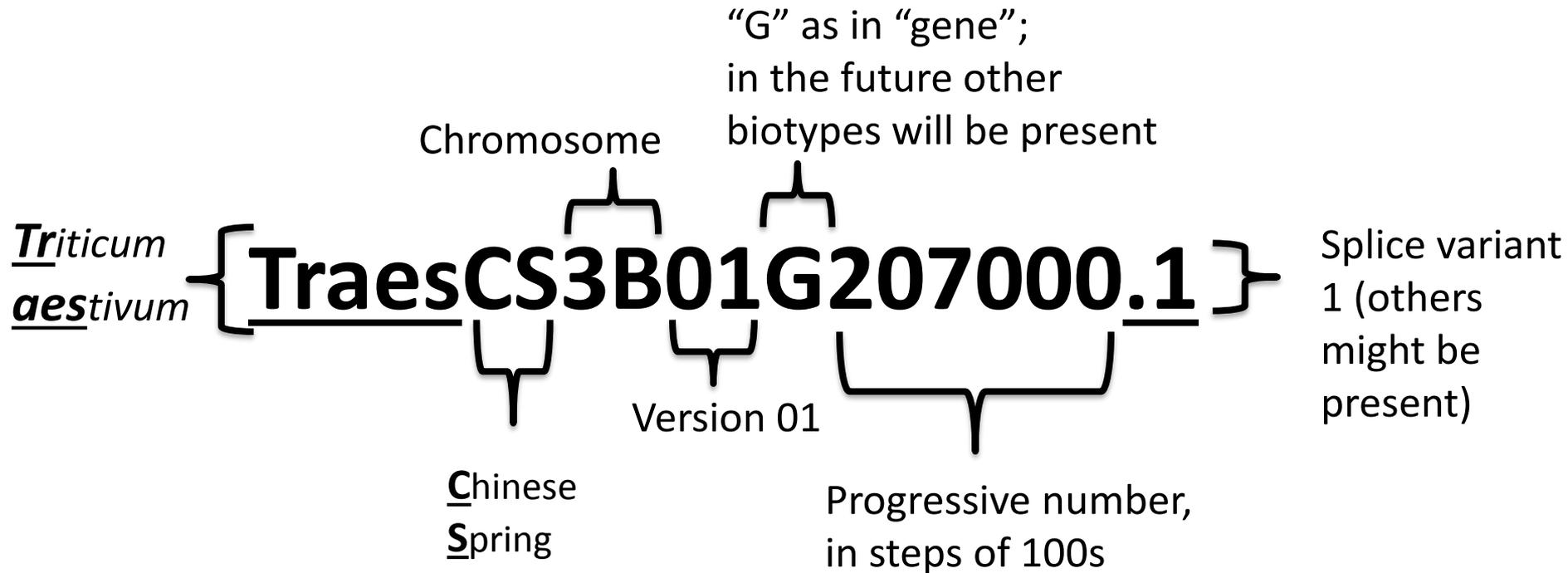


*preliminary results



IWGSC RefSeq v1.0 gene identifiers

For each gene locus, one transcript/isoform/splice variant has been selected as the representative transcript, which is the longest transcript/isoform/splice variant in the respective confidence class.



Conclusions

- More contiguous and complete wheat genome assemblies have enabled the full gene space to be captured
- The forthcoming IWGSC gene annotation will be the most comprehensive wheat annotation to date
- This annotation will be instrumental to perform global analyses such as:
 - Marker positioning for GWAS and breeding
 - Evolution analyses on many gene families
 - Gene regulation across multiple tissues
 - Defining precisely translocation events among different chromosome arms
- IWGSC RefSeq v1.0 focuses mainly on coding genes, however, in future releases we can expect a better definition of other important genomic elements such as:
 - Long and short non-coding RNAs
 - Pseudogenes
 - Expansions of the splicing isoform catalogue
 - Manual revision of specific gene families
- Finally, this assembly and annotation will be the cornerstones upon which future IWGSC projects for high quality functional annotation and for resequencing the breadth of global germplasm diversity



Acknowledgments

IWGSC Leadership: Rudi Appels, Kellye Eversole, Catherine Feuillet, Beat Keller, Jane Rogers

IWGSC Chromosome Leaders:



Etienne Paux, Frédéric Choulet



Bikram Gill



Rudi Appels



Institute of Experimental Botany of the AS CR, v. v. i.

Jaroslav Dolezel, Hana Simkova, Miroslav Valarik, Jan Bartos



Bayer CropScience

Catherine Feuillet
John Jacobs



Hikmet Budak



Nils Stein
Thorsten Schnurbusch



Hirokazu Handa



Universität Zürich UZH

Beat Keller



Curtis Pozniak
Andrew Sharpe



Luigi Cattivelli



Abraham Korol



Kuldeep Singh



Elena Salina



Odd-Arne Olsen



NORTHWEST A&F UNIVERSITY

Song Weining



Nikolai Ravin



Matt Clark



IWGSC RefSeq v1.0 Team Leaders

IWGSC Sequence Repository



Michael Alaux

BAC Libraries



Institute of Experimental Botany of the AS CR, v. v. i.

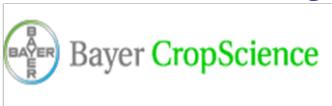
Jaroslav Dolezel, Hana Simkova

BAC Library Pools



Hélène Bergès

BAC WGP Tags



John Jacobs

Genetic Maps



Jesse Poland

RH Mapping



Vijay Tiwari

WGA PIs



Nils Stein



Curtis Pozniak
Andrew Sharpe



Jesse Poland



Frédéric Choulet

NRGene

Gil Ronen



Assaf Distelfeld

TEL AVIV UNIVERSITY

illumina

Mike Thompson



Kellye Eversole
Jane Rogers

Pseudomolecule Team



Frédéric Choulet



Economic Development,
Jobs, Transport
and Resources

Gabriel Keeble-Gagnere



Martin Mascher

Annotation Team



Philippe Leroy
Frédéric Choulet

PGSB

Manuel Spannagl, Klaus Mayer



David Swarbreck

RNASeq



Cristobal Uauy



IWGSC RefSeq v1.0 Annotation Team



- Coordination and oversight – Jane Rogers, IWGSC

- Two annotation teams:

- **INRA** – Frédéric Choulet, Hélène Rimbart, Philippe Leroy



- **PGSB** – Sven Twarzdiok, Klaus Mayer, Manuel Spannagl



- Evaluation and integration team

- **Earlham Institute** – David Swarbreck, Luca Venturini, Gemy Kaithakottil



Thanks to IWGSC Sponsors!



Acknowledgments



- Bernardo Clavijo and his team
 - Gonzalo Garcia Accinelli
 - Jon Wright
- David Swarbreck and his team
 - Gemy George Kaithakottil
 - Daniel Mapleson
- Ksenia Krasileva and her team
 - Christian Schudoma
 - Dina Raats
- Matt Clark and his team
 - George Kettleborough
 - Aurore Coince
 - Ned Peel
 - Lawrence Percival-Alwyn
- Darren Heavens and his team
 - James Lipscombe
- Helen Chapman
- Tom Barker
- Robert Davey
- Christine Fosker
- Federica di Palma



- Cristobal Uauy and his team
 - Philippa Borrill
 - Ricardo H. Ramirez-Gonzalez
- Mike Bevan and his team
 - Fu Hao-Lu
 - Neil Mckenzie
- Guotay Yu



- Harvey Millar and his team
 - Owen Duncan
 - Joshua Troesch



- Paul Kersey and his team
 - Dan Bolser
 - Guy Namaati
 - Arnaud Kerhornou



- Manuel Spannangl
- Heidrun Gundlach
- Georg Haberer



- Andy Phillips