

Improved gene annotation of the fungal wheat pathogen *Zymoseptoria tritici* based on combined Iso-Seq and RNA-Seq evidence



Nicolas Lapalu



Lucie Lamothe



Nicolas Lapalu¹, Lucie Lamothe¹, Yohann Petit¹, Anne Genissel¹, Camille Delude², Alice Feurtey^{3,4}, Leen N. Abraham³, Dan Smith⁵, Robert King⁵, Alison Renwick⁶, Mélanie Appertet², Justine Sucher², Andrei S. Steindorff⁷, Stephen B. Goodwin⁹, Gert H.J. Kema¹¹, Igor V. Grigoriev^{7,8}, James Hane⁶, Jason Rudd⁵, Eva Stukenbrock¹⁰, Daniel Croll³, Gabriel Scalliet², Marc-Henri Lebrun¹

¹Université Paris-Saclay, INRAE, UR1290 BIOGER, Palaiseau, France

²Syngenta Crop Protection AG, Stein, Switzerland

³University of Neuchâtel, Neuchâtel, Switzerland

⁴ETH Zurich, Zurich, Switzerland

⁵Dept of Protecting Crops and the Environment, Rothamsted Research, UK

⁶Centre for Crop and Disease Management, Curtin University, Perth, Australia

⁷U.S. Department of Energy Joint Genome Institute, Berkeley, USA

⁸Department of Plant and Microbial Biology, University of California Berkeley, Berkeley, USA

⁹USDA-Agricultural Research Service, West Lafayette, USA

¹⁰Environmental Genomics, Max Planck Institute for Evolutionary Biology, Germany

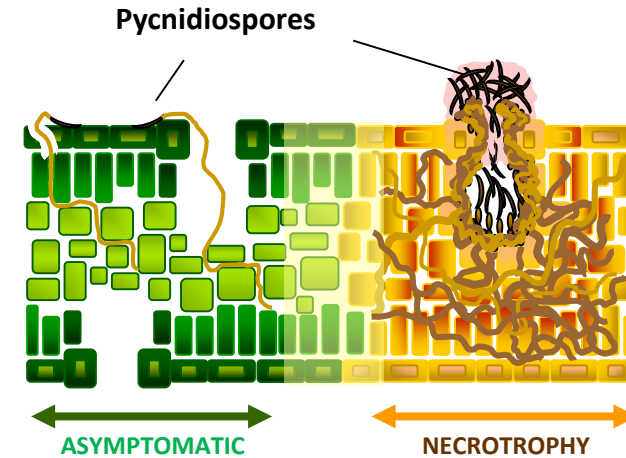
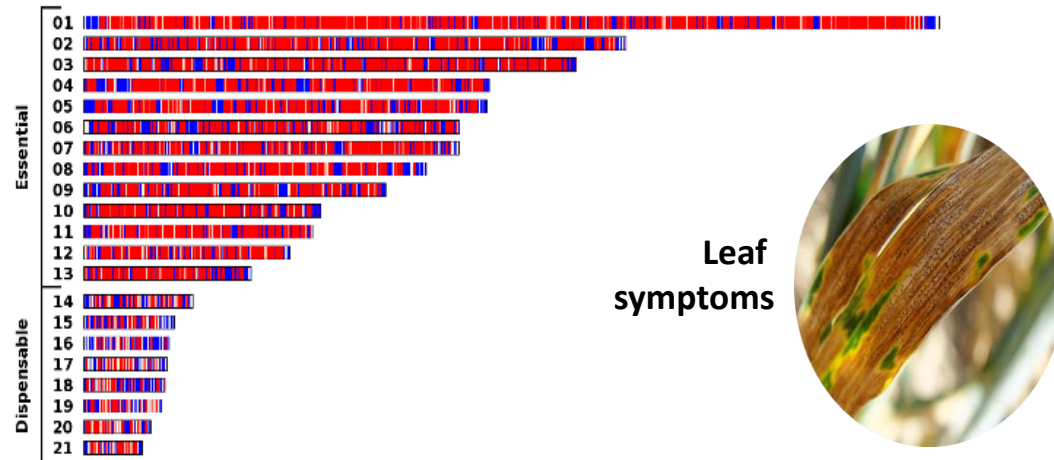
¹¹Wageningen University and Research, The Netherlands

Lapalu et al. 2023. bioRxiv. doi: <https://doi.org/10.1101/2023.04.26.537486>

Zymoseptoria tritici, a fungal pathogen of wheat

Fully sequenced genome 21 chromosomes,
39.7 Mb, JGI: 10.952 genes, 17-20 % transposons

Haploid hemi-biotrophic ascomycete
Life cycle: sexual and asexual reproduction



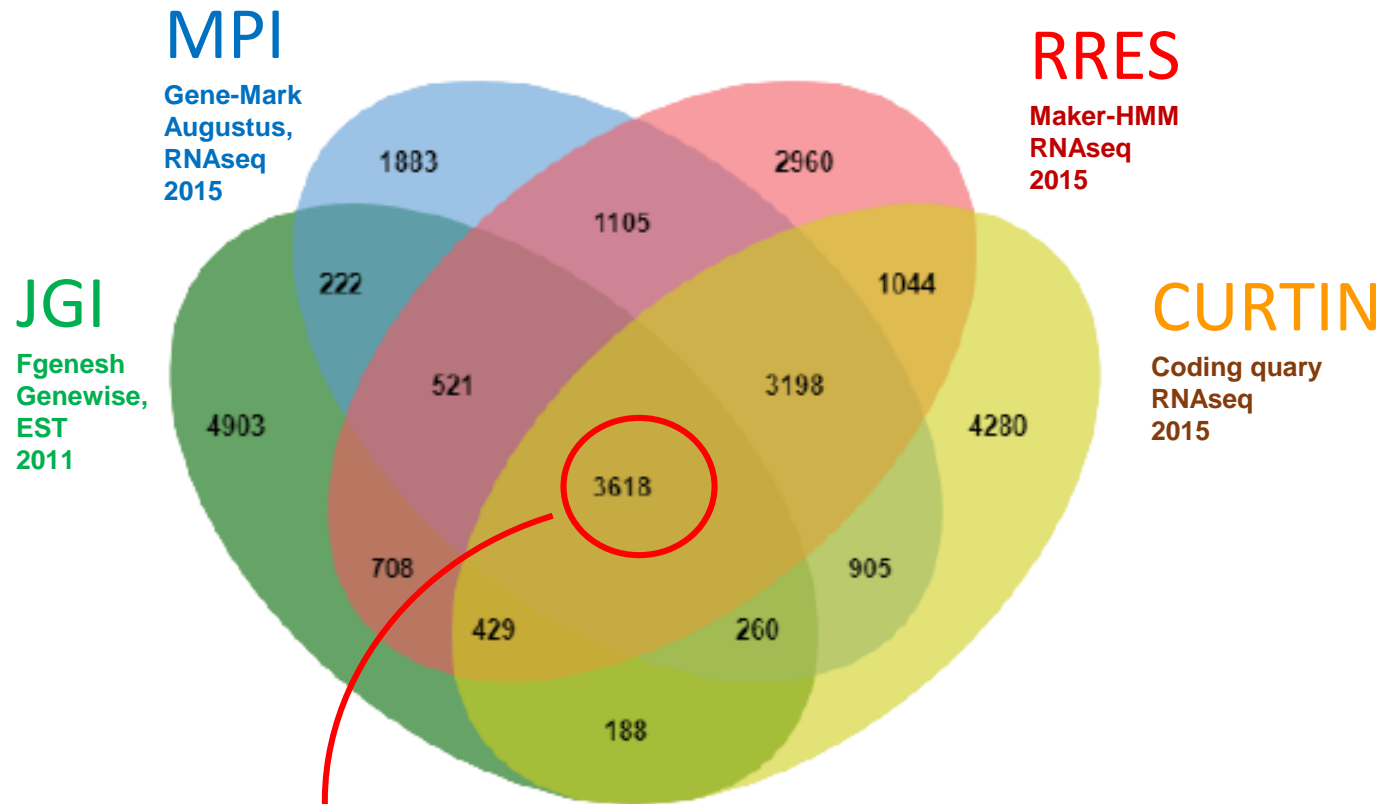
Genetics crosses, genome sequences, genome annotations
Transformation, reverse genetics using Ku70 deficient strains, cellular biology, expression vectors
Large infection RNAseq datasets,
GWAS, populations genomics, experimental evolution

Improvement of *Z. tritici* gene annotation

Why do we need to improve *Z. trititi* gene annotation ?

Why do we need to improve *Z. tritici* gene annotation ?

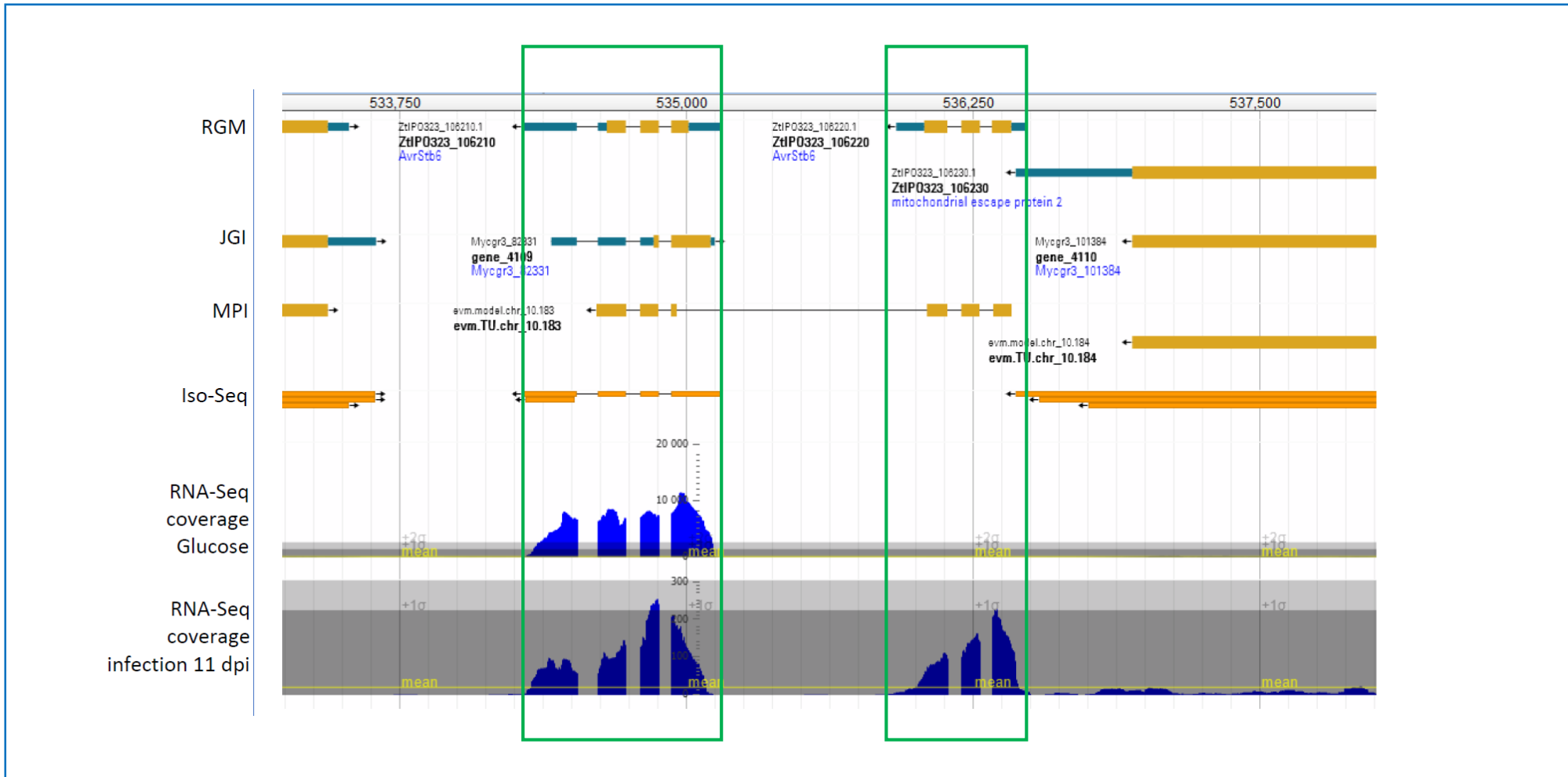
Comparison of existing IPO323 *Z. tritici* gene annotations



Among 16.000 metagenes (loci),
only 3618 (22 %) displayed identical gene structures

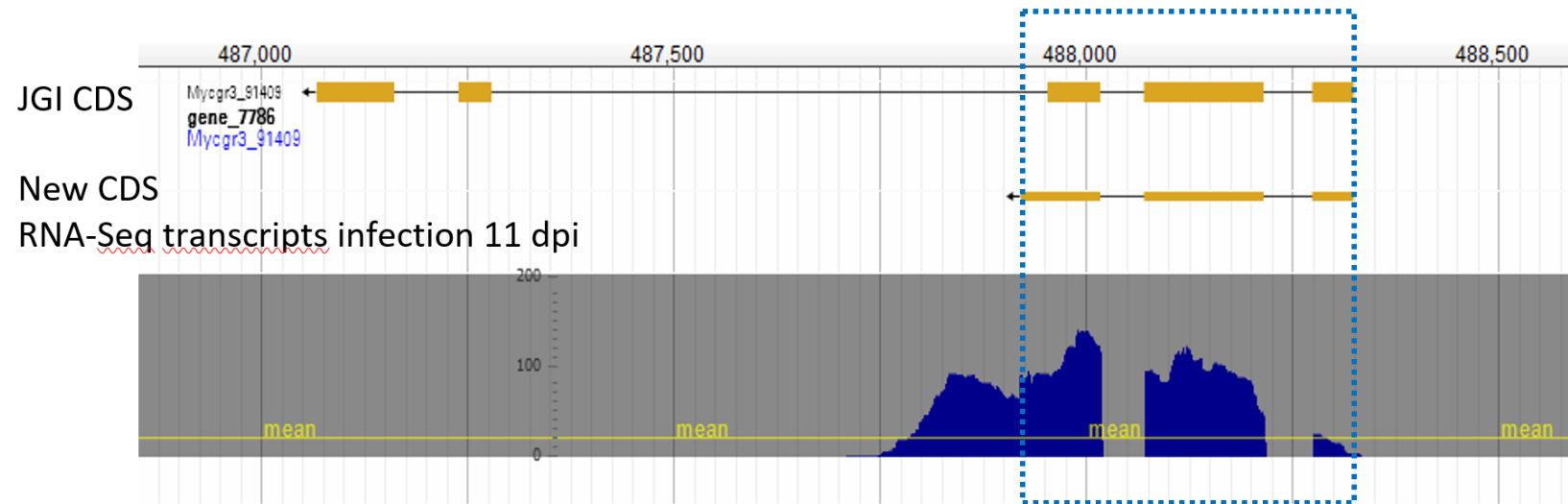
Why do we need to improve *Z. trititi* gene annotation ?

***Avr-Stb6* encoding genes were not predicted by any previous annotations**



Why do we need to improve *Z. trititi* gene annotation ?

Zt-NIP1 encoding gene was not predicted accurately by previous JGI annotation (NCBI)



Zt-Mycgr3-91409

MAPIFTYAVAALAFQAQSAAYAVVYAARCKFGNPLVQNNRITRAVCDLTNEHTTKDGSWHYVEVDNECKYLAGDNPRDQPGWAVFVKYCPEQNAAADKSRARGEQTCGFVRTPVDDVSRAESATTPVCEGD

NEW Zt-Mycgr3-91409-2

MAPIFTYAVAALAFQAQSAAYAVVYAARCKFGNPLVQNNRITRAVCDLTNEHTTKDGSWHYVEVDNECKYLAGDNPRDQPGWAVFVKYCTYYKGPDA

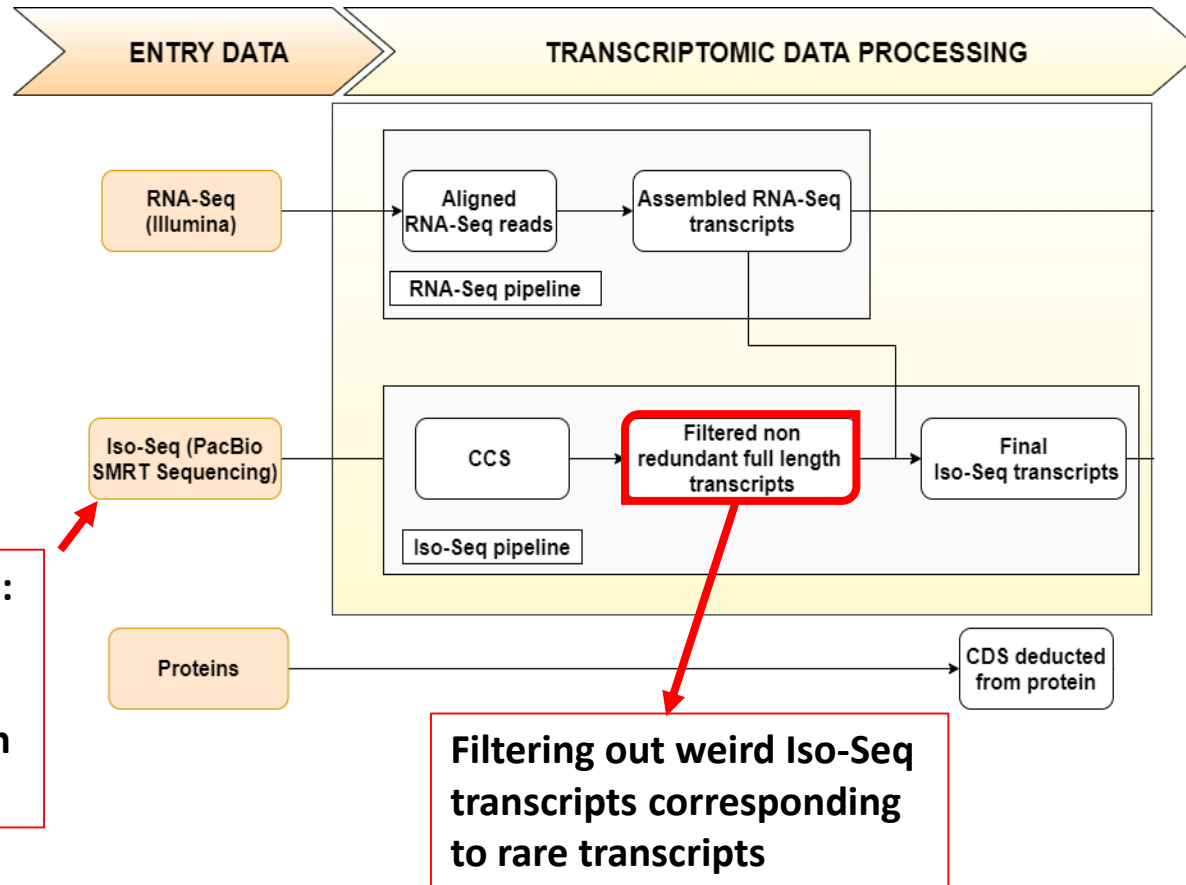
Errors in database searches

Improvement of *Z. tritici* gene annotation

**Search for a modifiable strategy
to improve *Z. trititi* gene annotation**

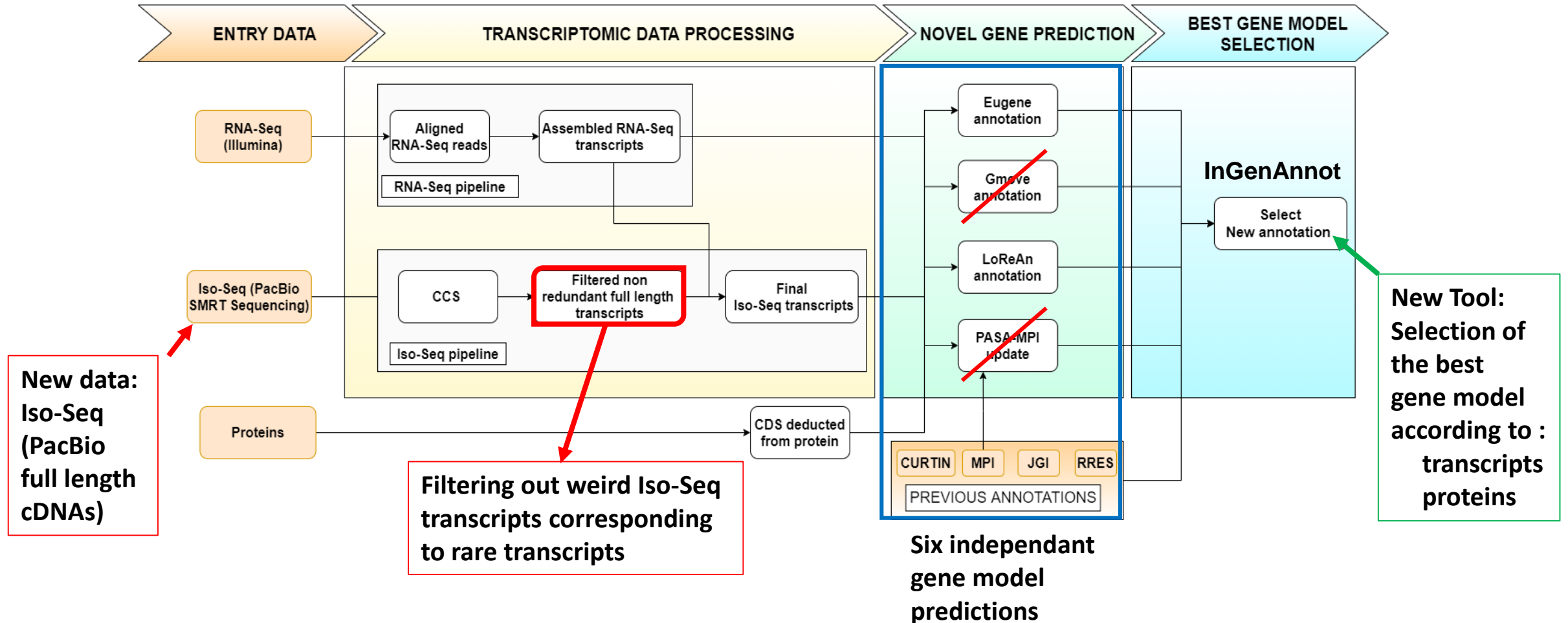
Improvement of *Z. tritici* gene annotation

Workflow for gene model prediction and selection



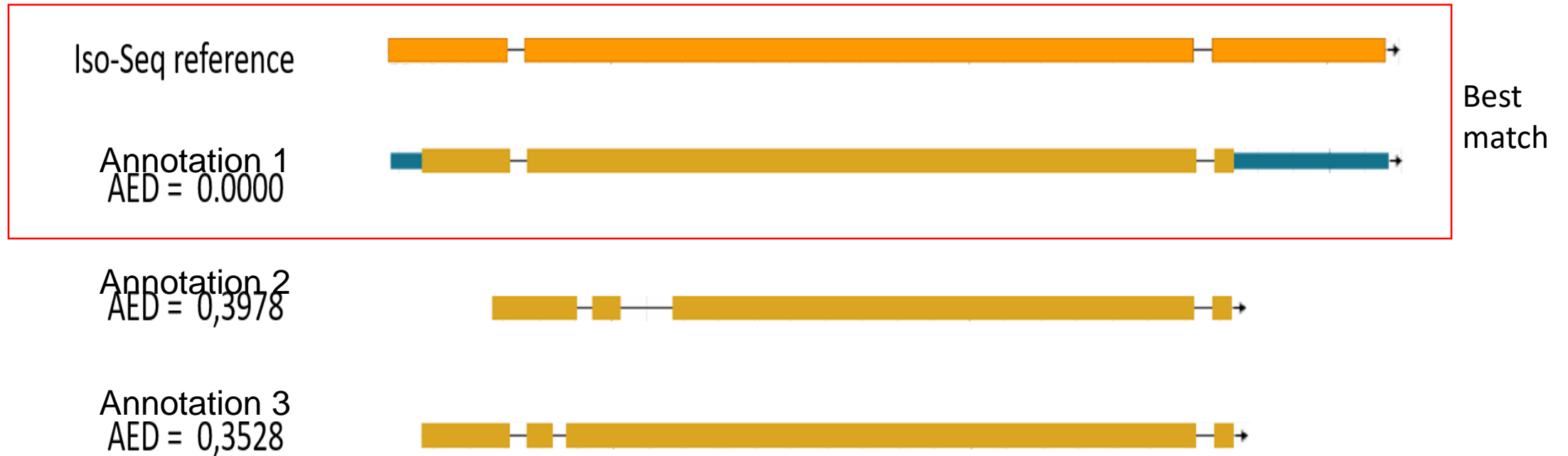
Improvement of *Z. tritici* gene annotation

Workflow for gene model prediction and selection



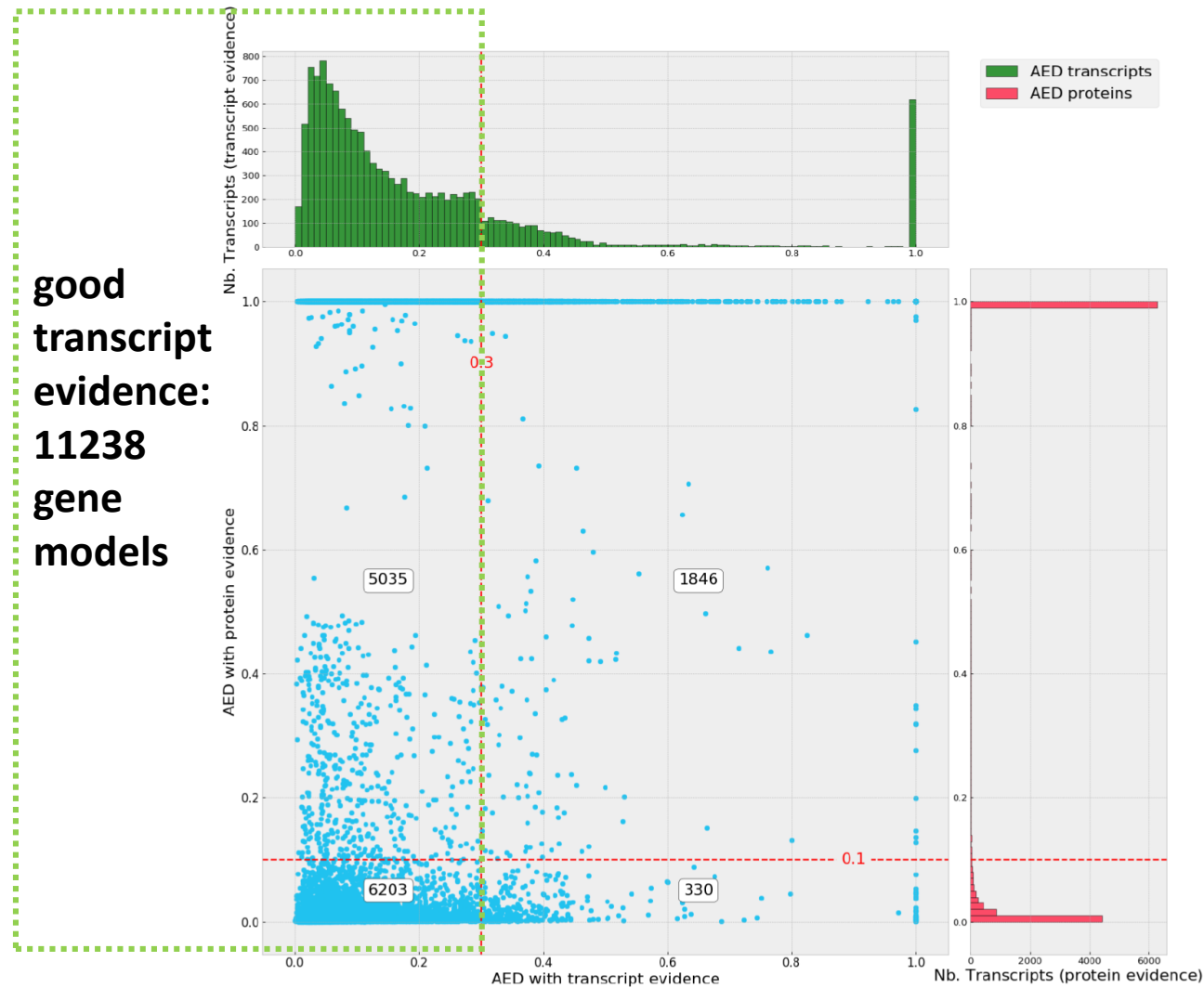
Improvement of *Z. tritici* gene annotation

Selection of the best gene models according to AED scores
(AED = value integrating evidences from transcripts and proteins)



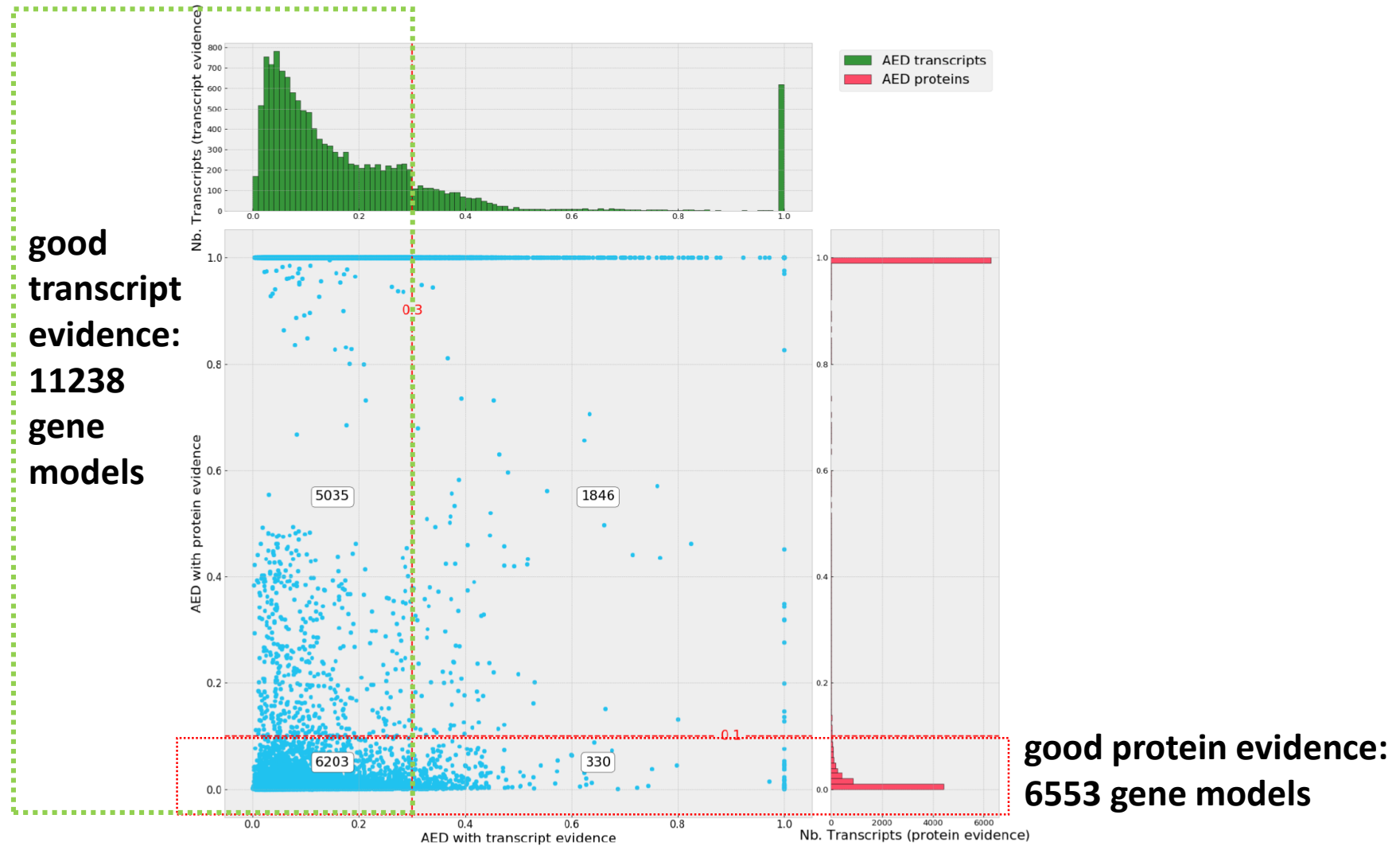
Improvement of *Z. tritici* gene annotation

New *Z. tritici* annotation with 13,414 gene models fitting to transcript ($0.3 < \text{AED}$) and/or protein ($0.1 < \text{AED}$) evidence



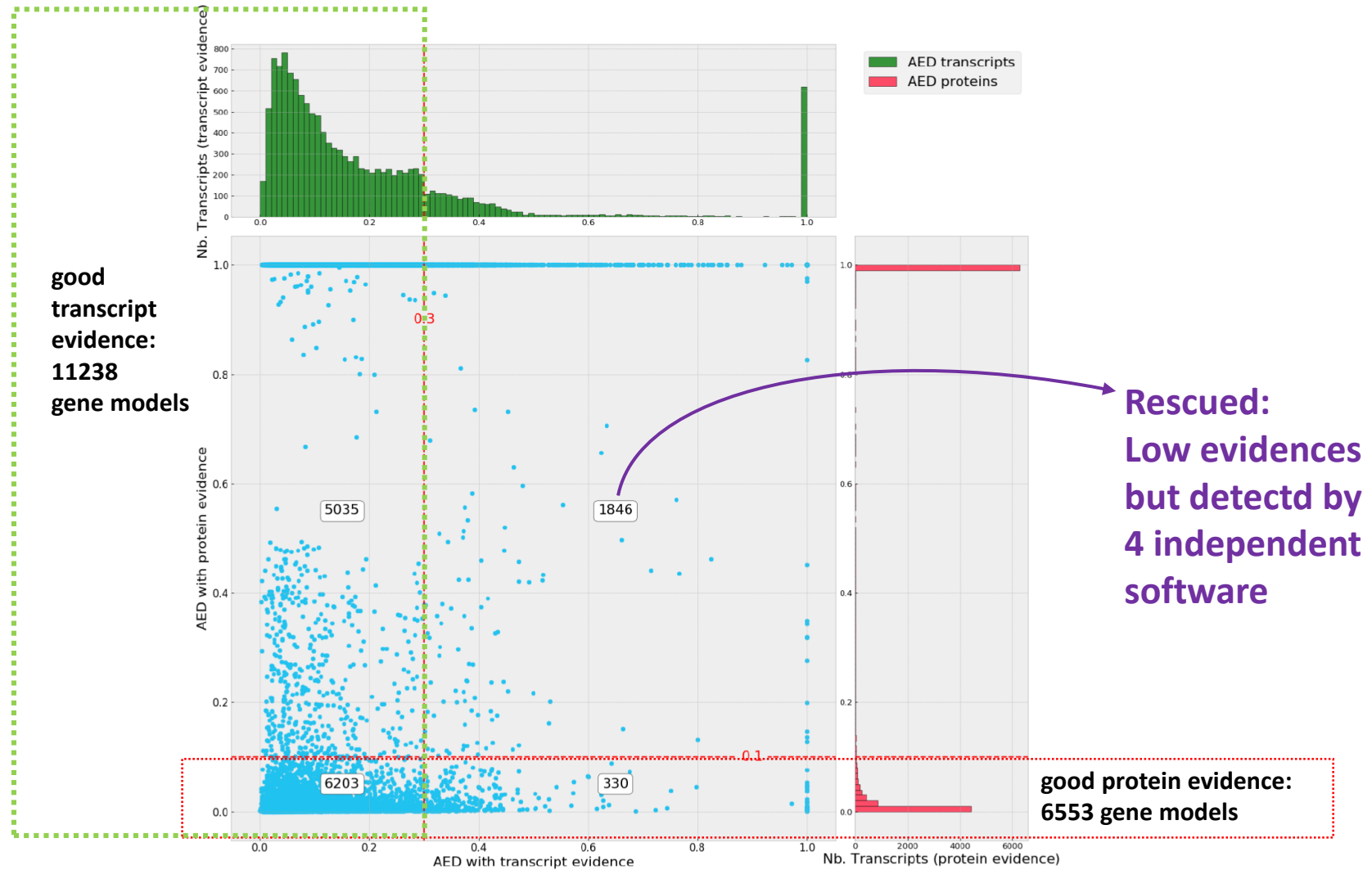
Improvement of *Z. tritici* gene annotation

New *Z. tritici* annotation with 13,414 gene models fitting to transcript ($0.3 < \text{AED}$) and/or protein ($0.1 < \text{AED}$) evidence



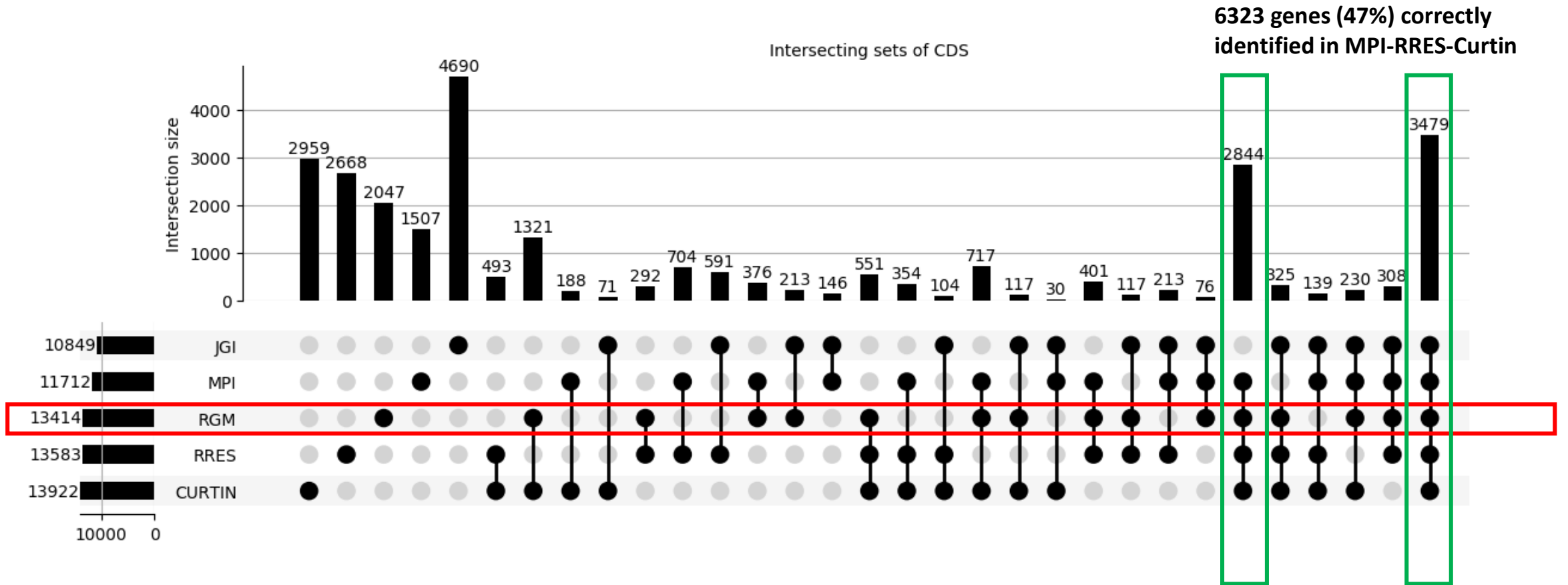
Improvement of *Z. tritici* gene annotation

New *Z. tritici* annotation with 13,414 gene models fitting to transcript ($0.3 < \text{AED}$) and/or protein ($0.1 < \text{AED}$) evidence or **rescued**



Improvement of *Z. tritici* gene annotation

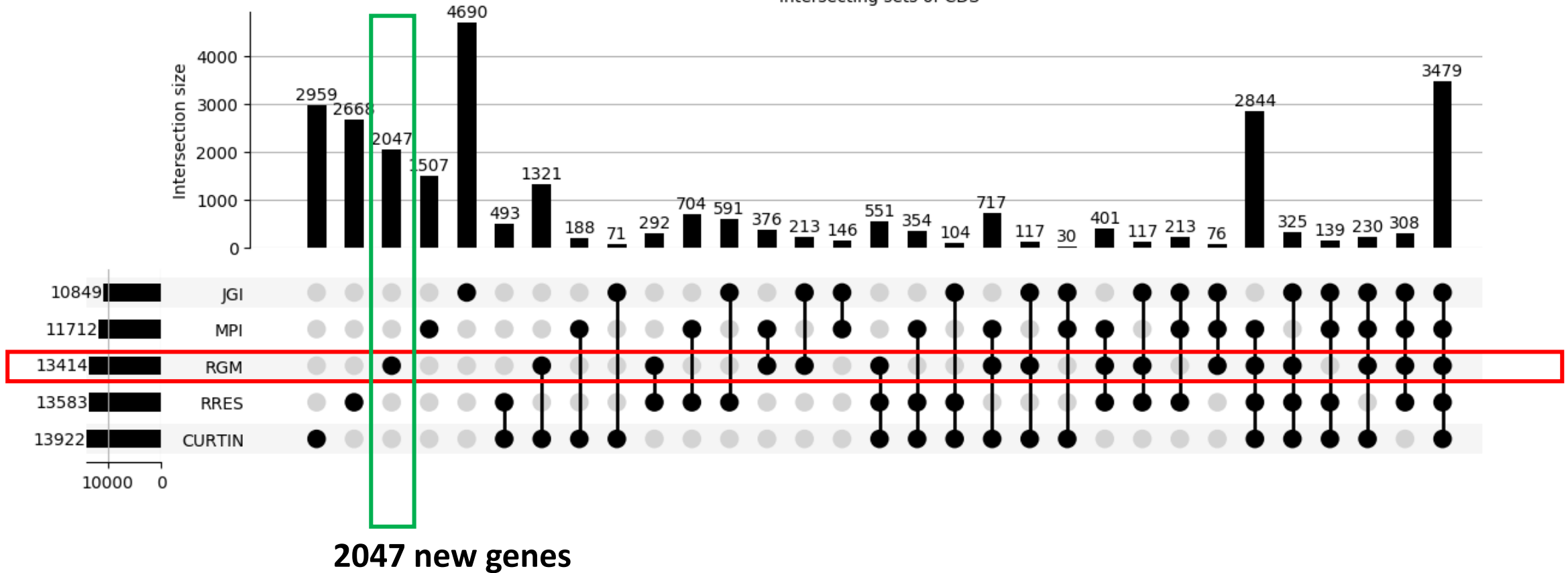
Comparison to previous annotations: % of identical predictions



Improvement of *Z. tritici* gene annotation

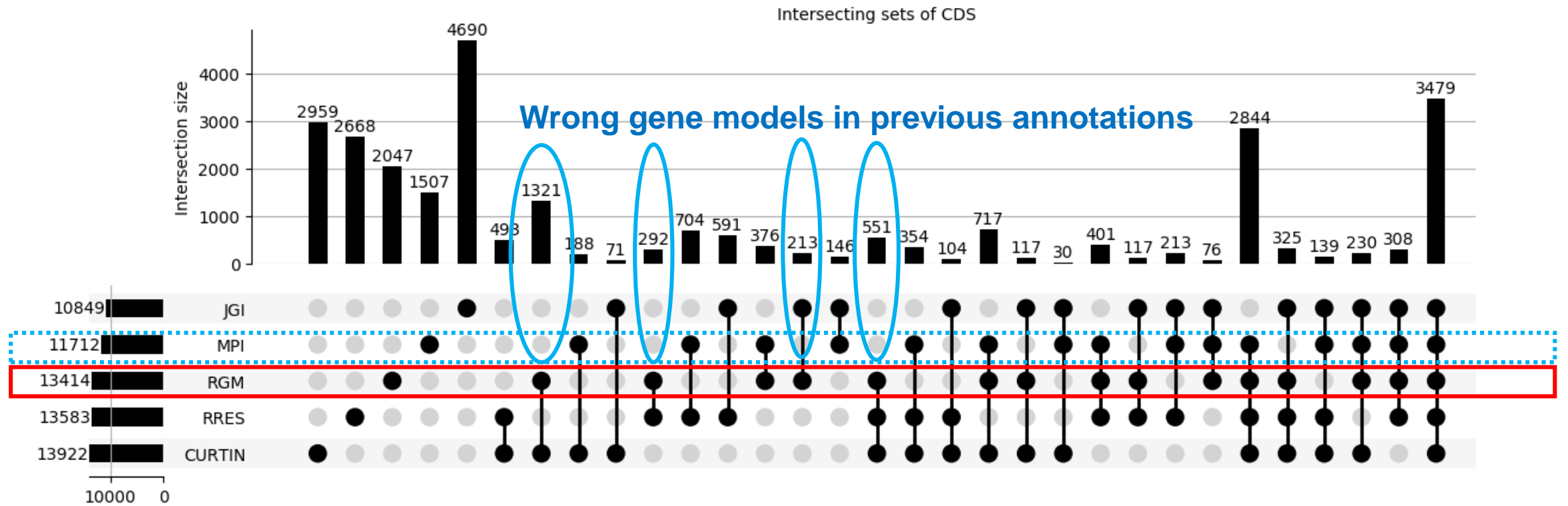
Comparison to previous annotations: new genes ?

Intersecting sets of CDS



Improvement of *Z. tritici* gene annotation

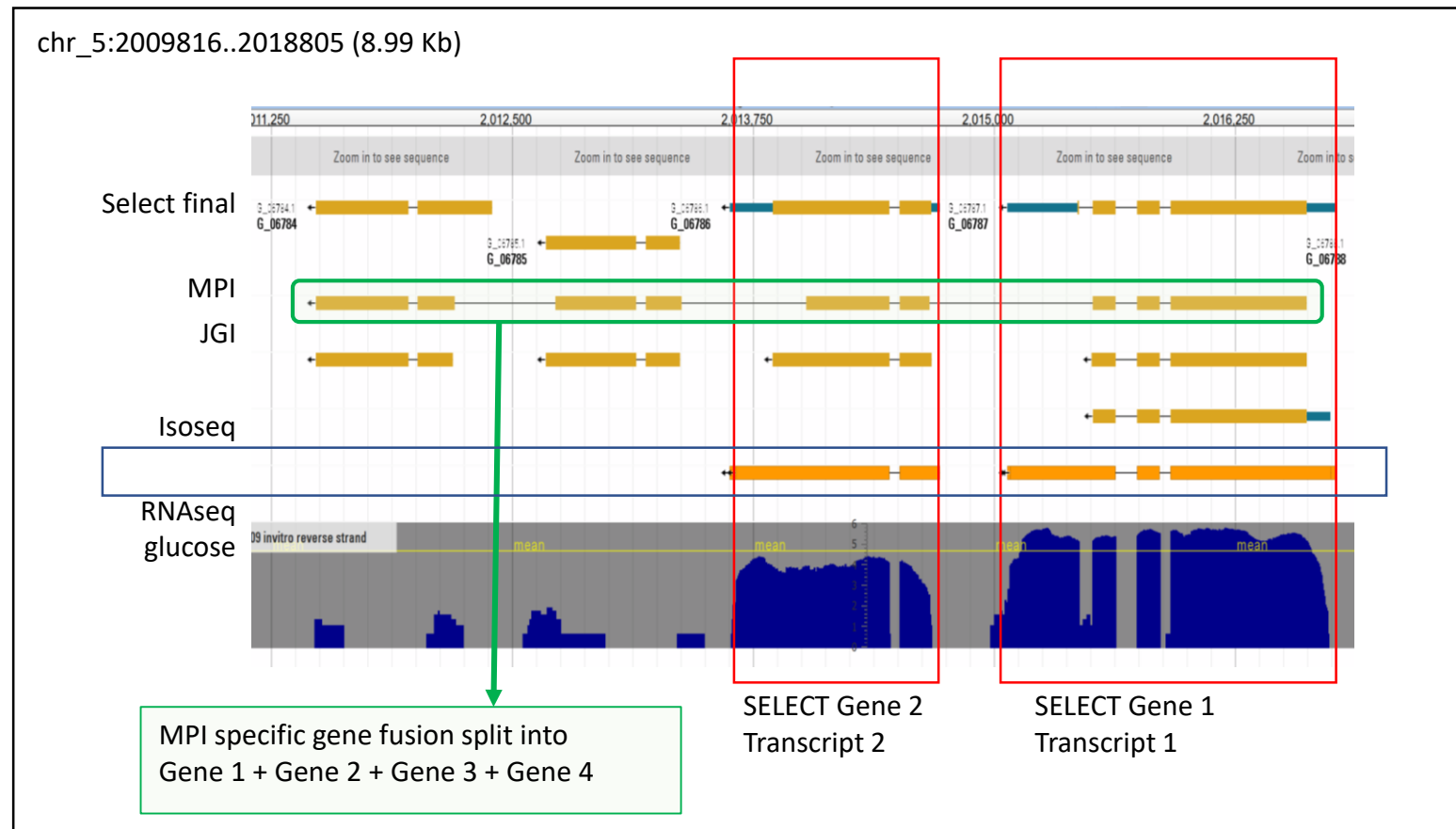
Comparison to previous annotations: systematic errors ?



Improvement of *Z. tritici* gene annotation

Systematic errors in previous gene models: gene fusions

Exemple of a MPI gene model splitted into 4 new genes



Improvement of *Z. tritici* gene annotation

Contribution of each software to the best gene models

	Annotation	Identical	Specific to	% total
Available annotations	JGI	4865	157	36
	MPI	8431	91	63
	RRES	8317	175	62
	CURTIN	9584	506	71
New annotations	Eugene	10224	1603	76
	LoReAn	7769	199	58

**No single software is able to predict all best gene models:
need for many to select the best**

Improvement of *Z. tritici* gene annotation: Conclusions

- **New datasets to improve CDS and gene (UTR) predictions:
Iso-Seq full length transcripts**
- **New tool to select the best gene model according
to transcript and protein evidence using different annotations pipelines
(Ingenannot)**

**Improved *Z. tritici* IPO323 gene annotation with 13,414 gene models
Annotation of 5' and 3' UTRs for 9,856 genes (73 %)**

- **Identification of transcript isoforms
13 % of *Z. tritici* expressed genes have transcript isoforms
(alternative transcriptional start, stop, exon skipping, intron retention)
mostly intron retention**

Improvement of *Z. tritici* gene annotation

What's next

- Add more ab initio software to InGeAnnot
Saturation in the prediction of novel genes ?
- Manual annotation still needed for some gene models
Web site for manual annotation
- Compare our pipeline (InGeAnnot) to other software
Braker3-Tsebra, Deep Learning
- Improve other fungal genomes with InGeannot and PacBio Isoseq data