# Phylogenomics/annotation resources published with the **bread wheat V1 genome papers**

- Resources:
  i. classified homeologs
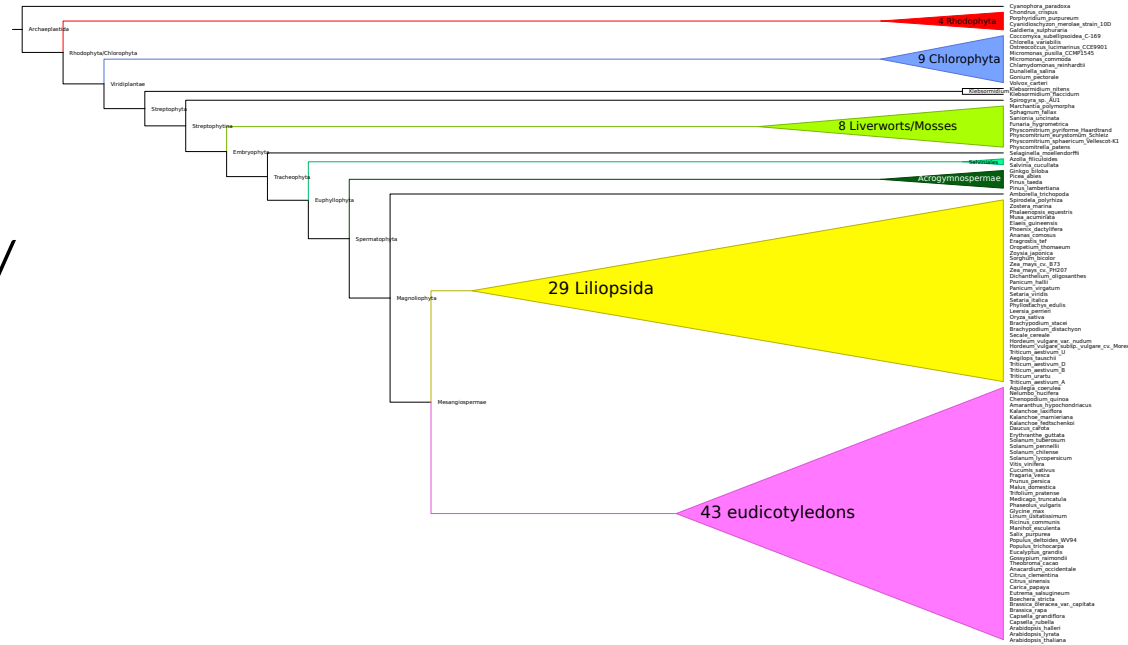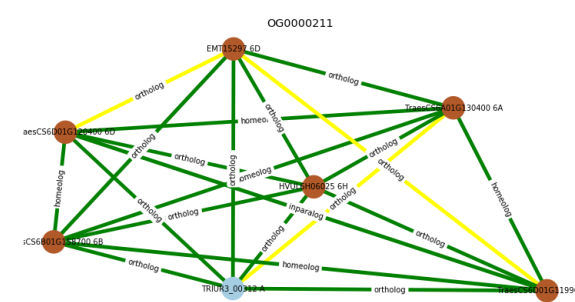  ii. orthology-based GOA|POA|TOA
  iii. wheat Transcription Associated Proteins (TAPs=TF+TR)
  iv. wheat gene families including orthologs from other Viridiplantae
  v. gene families with significant GCNVs
- all downloadable → links also in main paper supplement
  - https://wheat-urgi.versailles.inra.fr/Seq-Repository/Annotations
  - http://dx.doi.org/10.5447/IPK/2018/5
  - http://doi.org/10.1126/science.aar7191
- relevant resources are highlighted on the slides: **resource X**
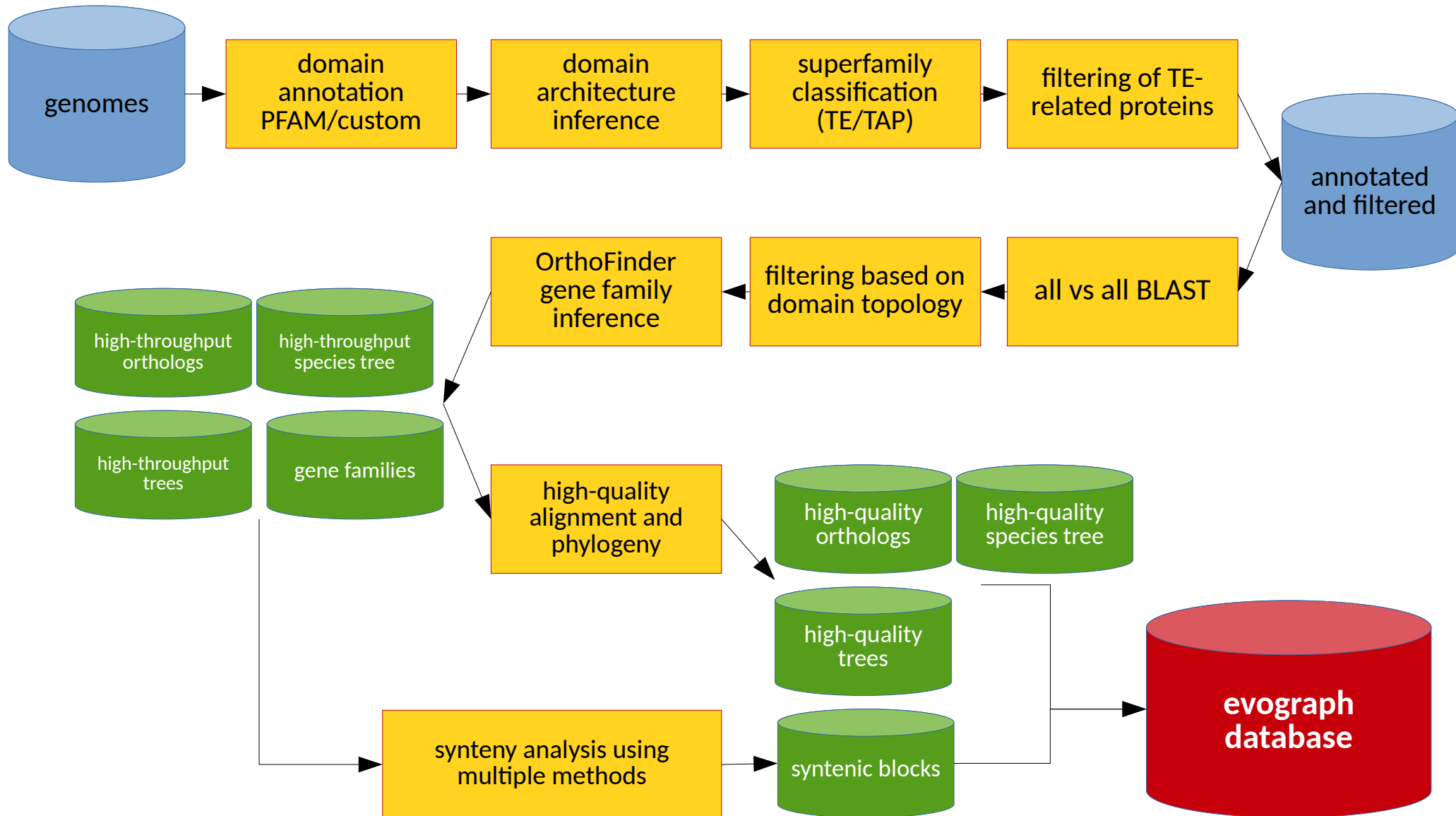
# evograph
## A graph of all plant gene families



- scalable, iterative phylogenomics workflow: *evograph*
  - targeting Plantae diversity
  - combine multiple iterations
  - broad taxon sampling ↔ family/ pangenome

- specific run for the bread wheat analysis
  - focussing on grasses → Triticeae
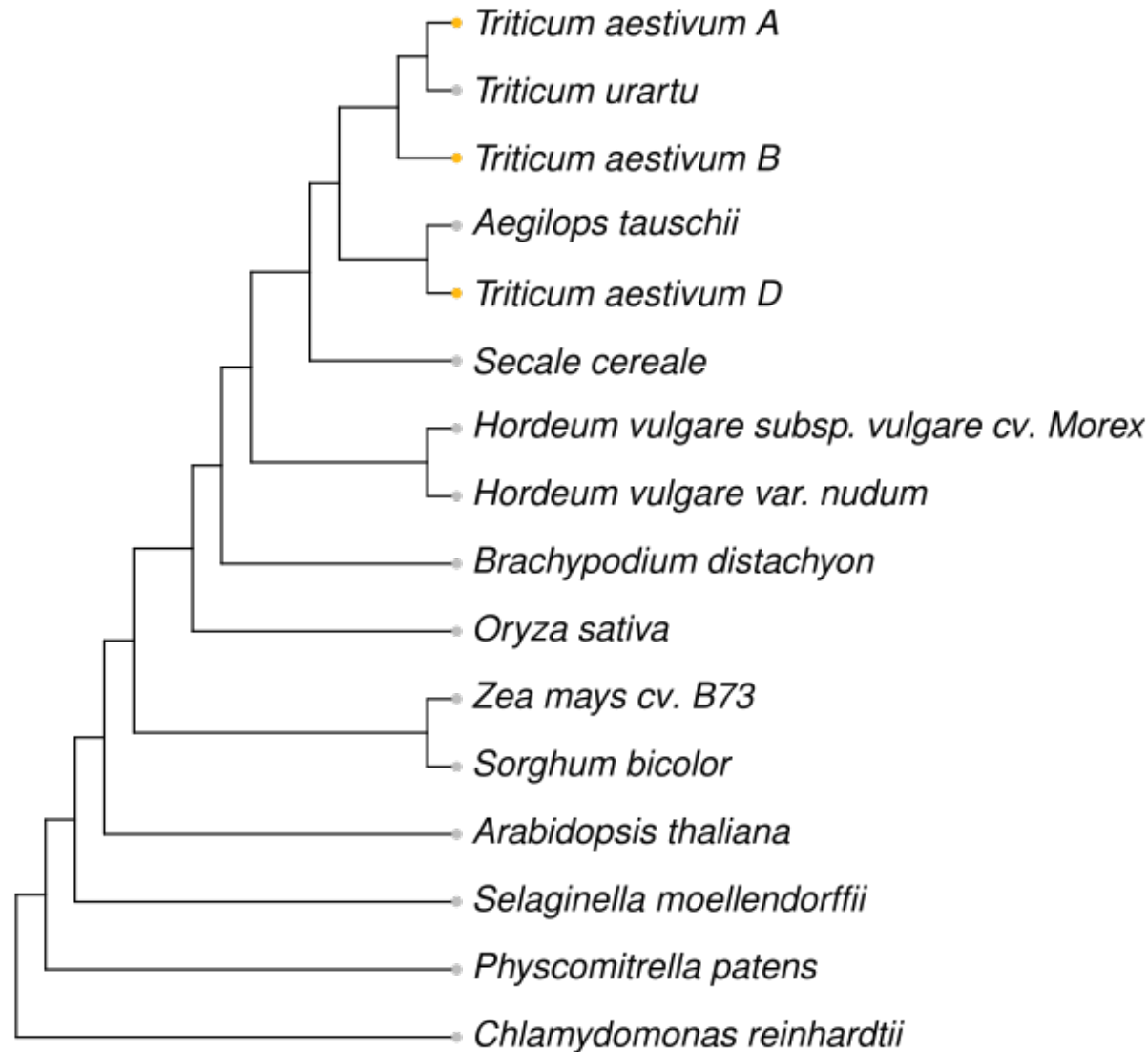  - with representative Viridiplantae as outgroups
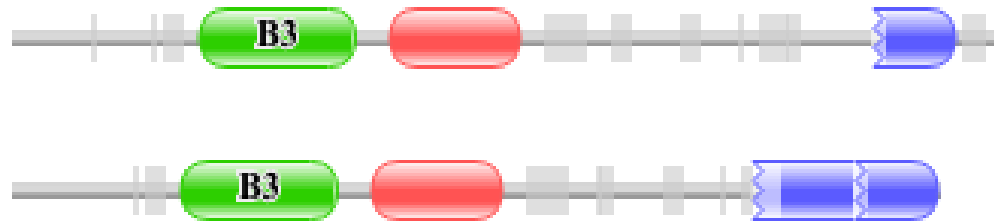
# evograph
## iterative workflow

# Bread wheat phylogenomics study

- 12 grass genomes
  - 9 grass species
  - A/B/D subgenomes as individual taxa
  - 2 barley varieties
- 4 Viridiplantae outgroups
- using OrthoFinder on filtered BLAST links
  - DendroBLAST phylogenies
- custom species tree
  - based on >12k subfamily trees [ASTRAL]
- dated chronogram of grasses used for comparative analysis
  - dated based on previous molecular clock estimates and fossil dates



*Triticum aestivum A*
*Triticum urartu*
*Triticum aestivum B*
*Aegilops tauschii*
*Triticum aestivum D*
*Secale cereale*
*Hordeum vulgare subsp. vulgare cv. Morex*
*Hordeum vulgare var. nudum*
*Brachypodium distachyon*
*Oryza sativa*
*Zea mays cv. B73*
*Sorghum bicolor*
*Arabidopsis thaliana*
*Selaginella moellendorffii*
*Physcomitrella patens*
*Chlamydomonas reinhardtii*

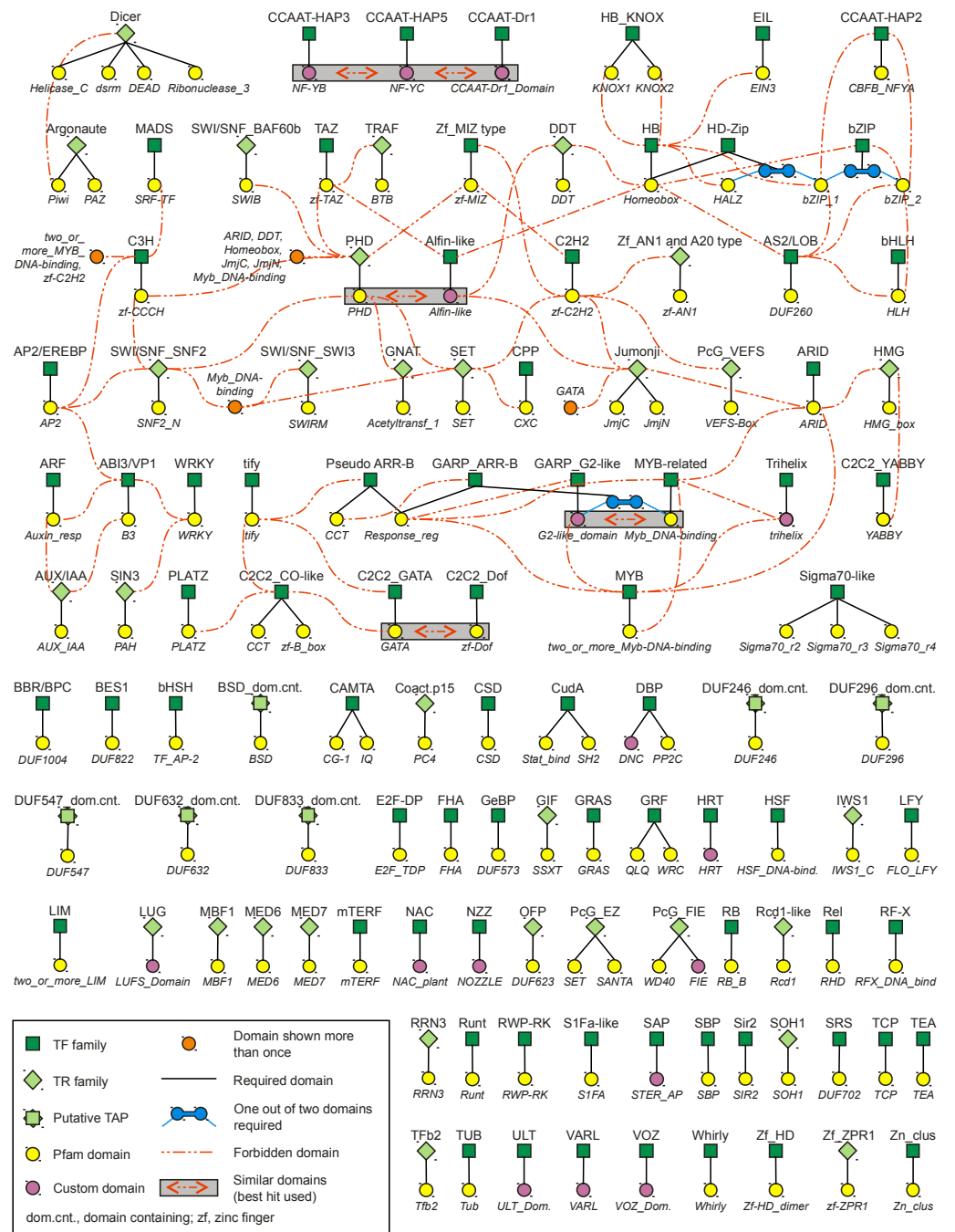# Protein **domain architectures** define gene families

- e.g. these represent 2 distinct families of ARFs:



- subfamily = comprising only orthologs and inparalogs
- gene family here:
  - 1-few subfamilies ideally traced to one ancestor gene in LCA of Viridiplantae
- inference of domain architectures
  - PFAM and custom HMMs
  - filtering of all vs. all BLAST links used as OrthoFinder inputs
    - hits must cover at least domain architecture
  - domain presence/abscence rules:
    - TE classification
    - TF superfamily

resource iv

# Using protein domain presence/absence to classify proteins e.g. **genome wide annotation of wheat Transcription Associated Proteins (TAPs) = TFs + TRs**

resource iii

| genome | genes | TAP super families |
|--------|-------|--------------------|
| A | 3025 | 104 |
| B | 3211 | 103 |
| D | 3163 | 103 |
| U | 158 | 45 |
| **TAPs** | **9557** | |

Lang, D., B. Weiche, G. Timmerhaus, S. Richardt, D.M. Riaño-Pachón, L.G.G. Corrêa, R. Reski, B. Mueller-Roeber, S.A. Rensing (2010): Genome-wide phylogenetic comparative analysis of plant transcriptional regulation: a timeline of loss, gain, expansion and correlation with complexity Genome Biology and Evolution 2, 488-503.

HelmholtzZentrum münchen
German Research Center for Environmental Health

Daniel.Lang@helmholtz-muenchen.de

PGSB Plant Genome and Systems Biology

# Homolog relationships /
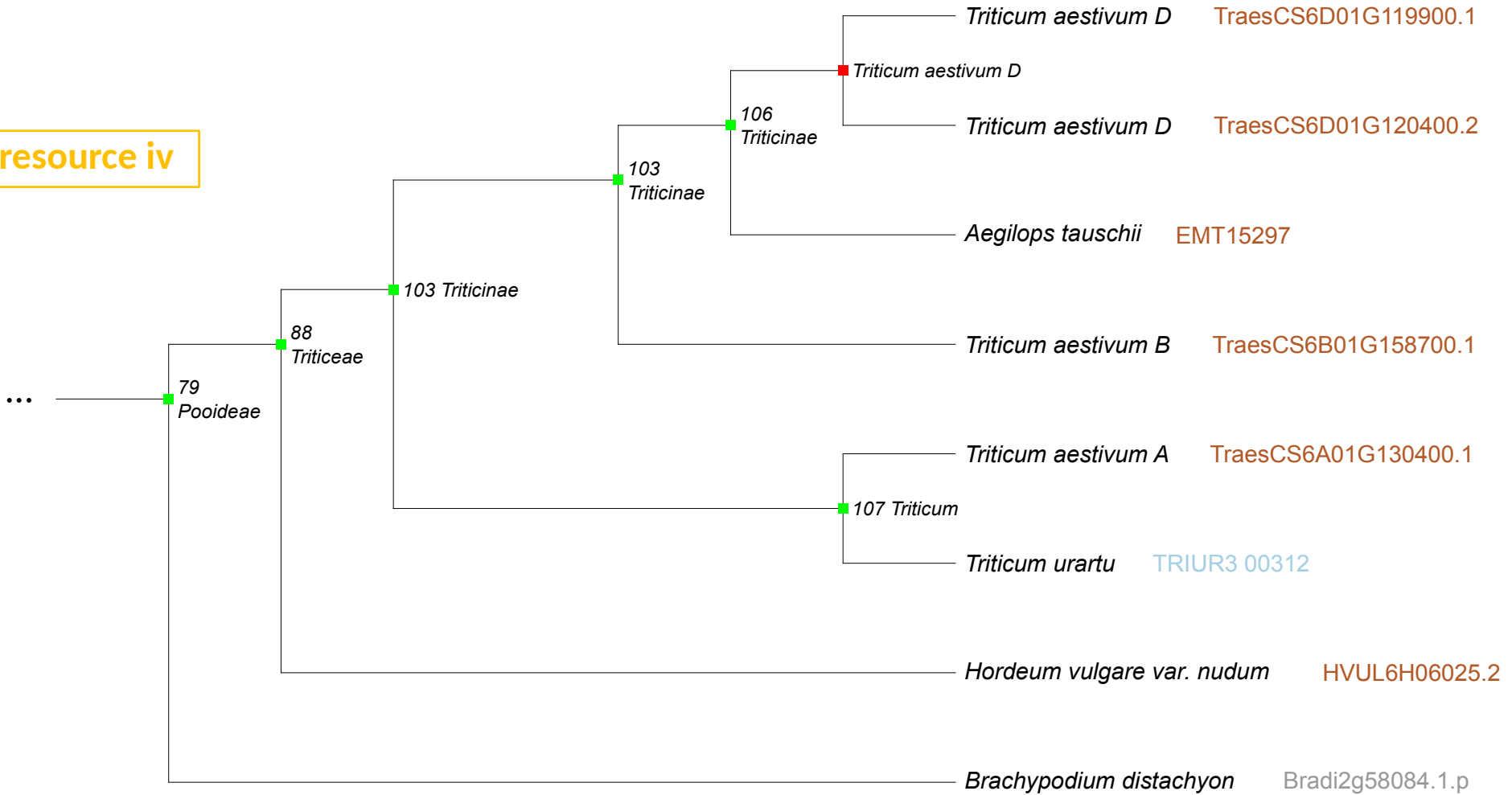**ortholog and homeolog** inference

- reconciliation of gene trees with the species using Species Overlap algorithm

  - infer **speciation** and **duplication** events → **ortholog**, **inparalog** and **outparalog** relationships

- wheat subgenomes treated as independent taxa

  - subgenome orthologs → **homeologs**

- use colinearity, chromosomal ancestry and position to assess translocations/transduplications

| class | pairs |
|---|---|
| homeolog | 123,588 |
| inparalog | 1,797,522 |
| ortholog | 2,555,680 |
| outparalog | 59,046,632 |

**resource i**

# Example of an annotated **reconciled gene tree** of a Pooideae-specific subfamily

resource iv

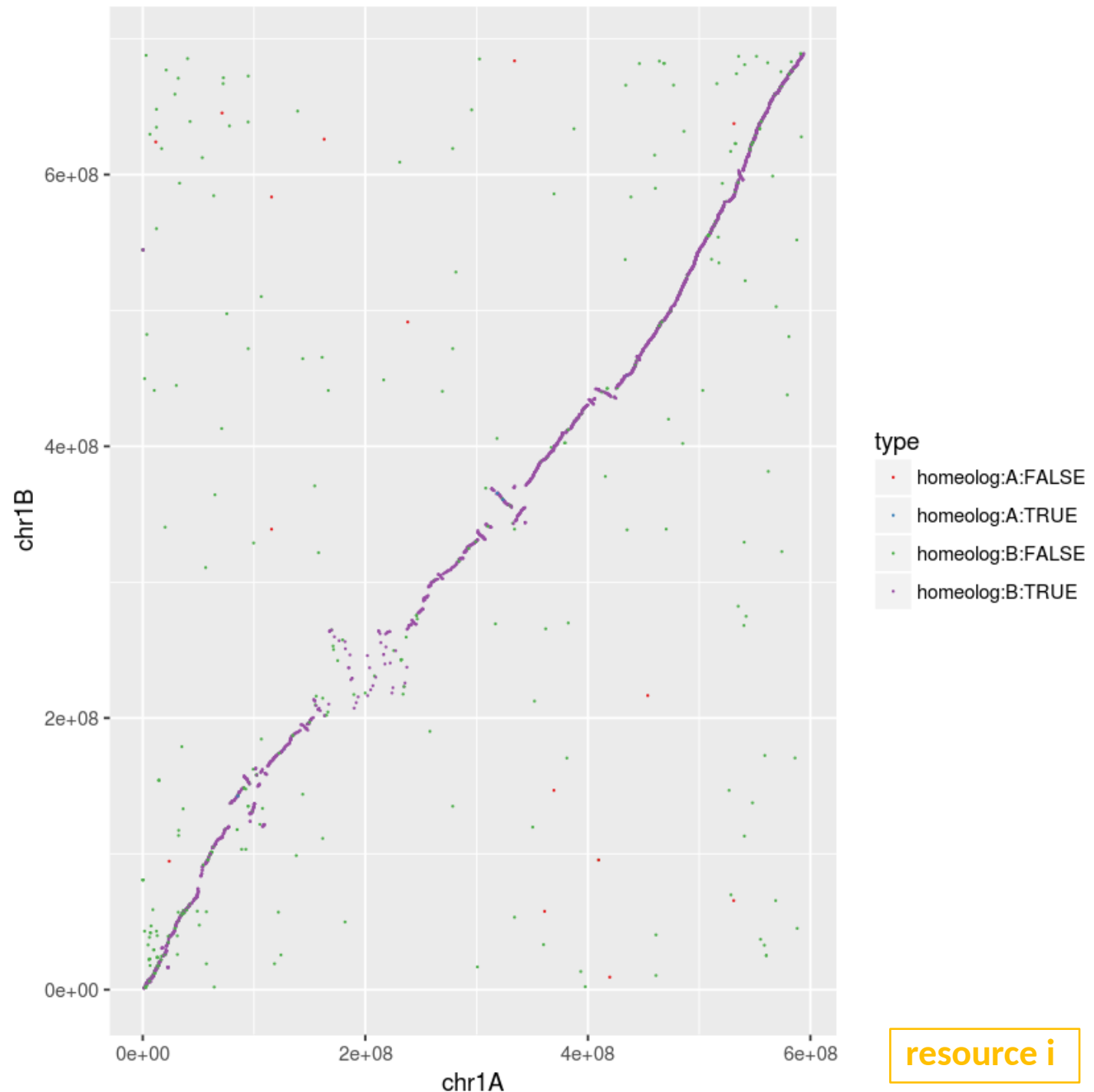# Orthology-based **ontology annotation** for wheat V1

- concept: transfer of ontology term annotation from orthologs
  - rich data from Gene Ontology consortium, TAIR, Gramene/Planteome
- Gene Ontology (**GO**), Plant Ontology (**PO**), Trait Ontology (**TO**)
- evidence quality:
  - transfer: IEA =IEA, all other: ISO (inferred by sequence orthology)
  - pfam2GO/domain architecture: ISM (inferred from sequence model)
- sequence similarity-based annotation using AHRD (only IEA):
  - 120,185 genes
  - **1,820** GO terms
- orthology-based yields more and more specific annotations:

| ontology | nassoc | nfamilies | nseq | nterms |
|---|---|---|---|---|
| **PO** | 3,622,724 | 12,631 | 79,530 | 446 |
| **GO** | 1,560,102 | 15,983 | 113,815 | **7,408** |
| **TO** | 8,173 | 227 | 1,060 | 279 |

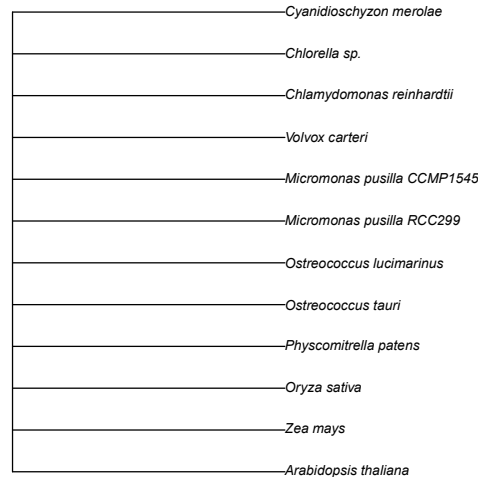| evidence_code | count |
|---|---|
| ISO | **4,989,419** |
| IEA | 177,154 |
| ISM | 155,326 |

resource ii

# Genome-scale resolution and **classification of homeologs:**

syntenic (TRUE) ↔ translocated/-duplicated (FALSE)
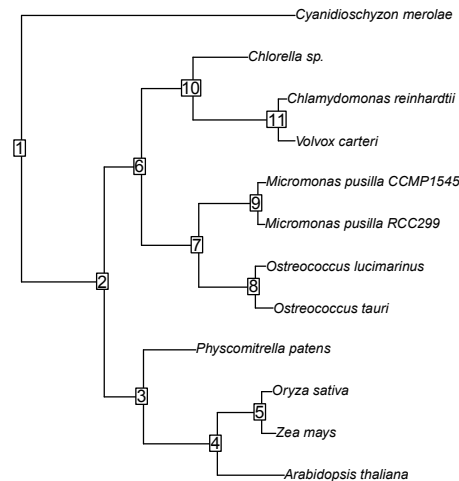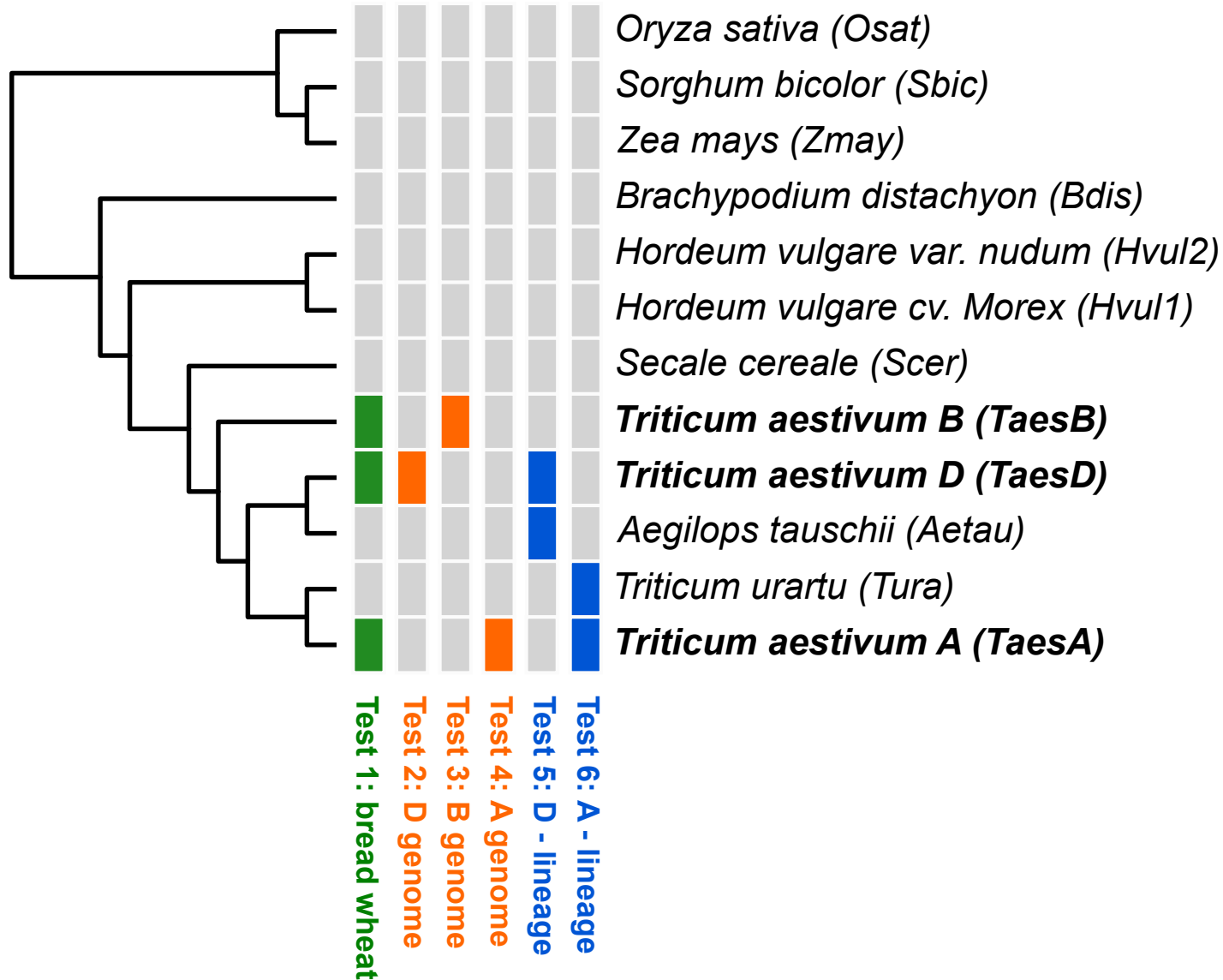
Any **comparative analysis** among species is affected by phylogenetic relationships



- correct for phylogenetic dependence with **phylogenetic comparative (PC) methods**
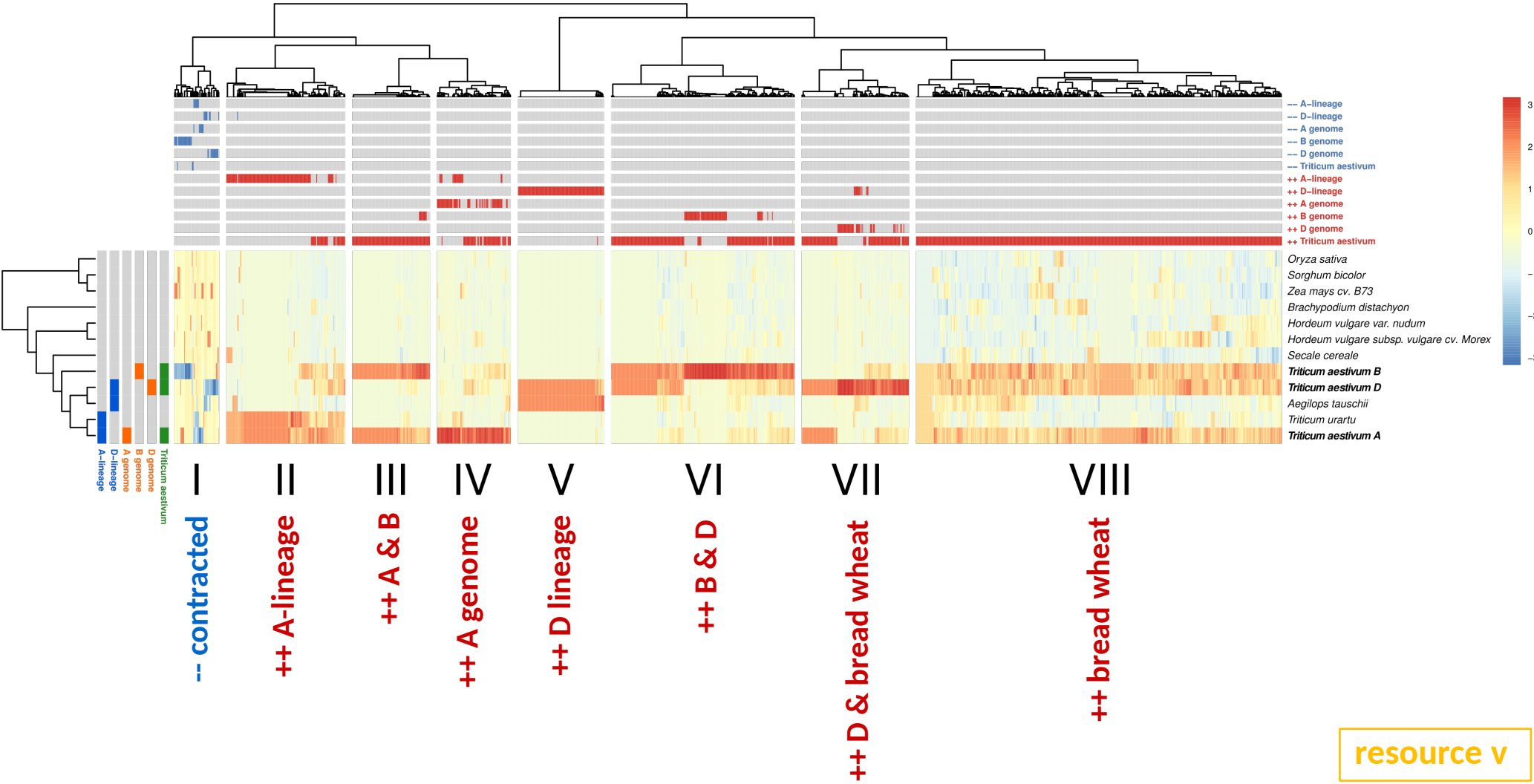  - developed for morphological traits
  - e.g. Felsenstein's PIC 1985

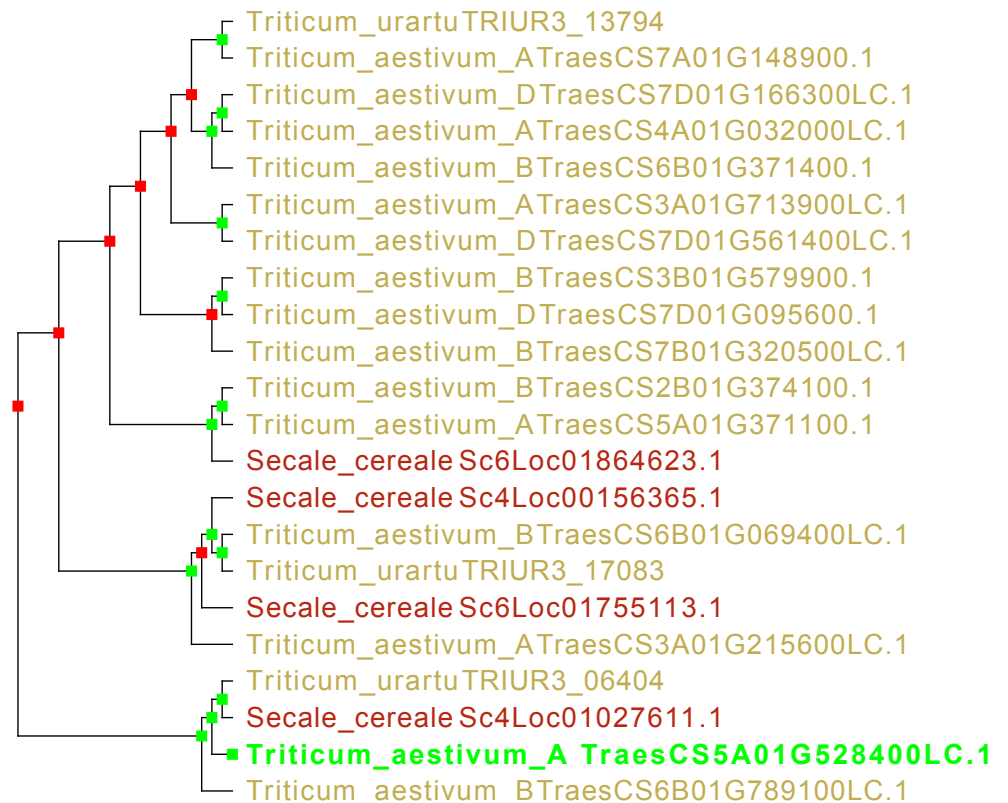# Phylogenetic one-way ANOVA testing for **Gene Copy Number Variations** (**GCNVs**) in these sets

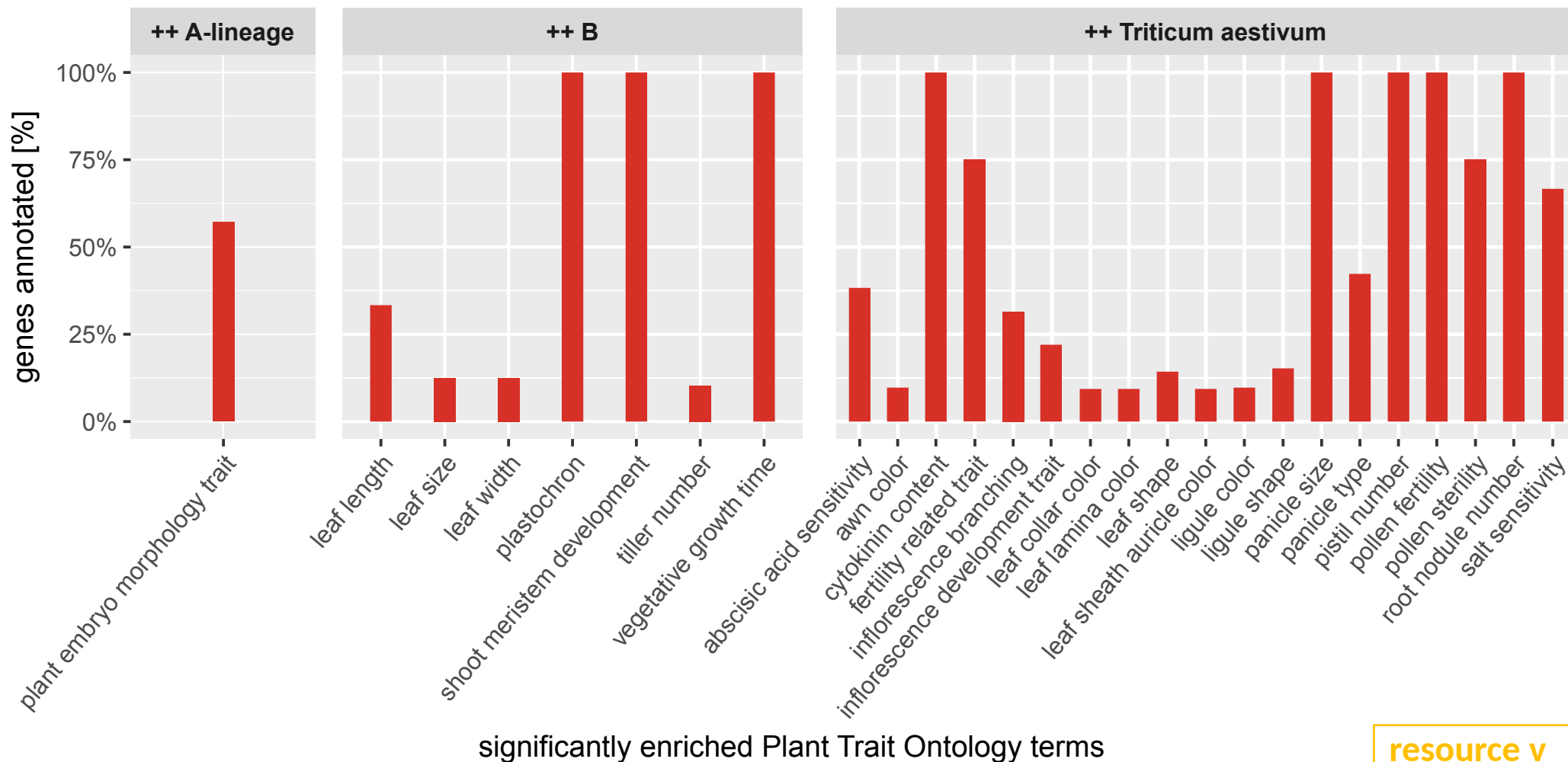# **GCNVs** in gene families shaping the traits that make bread wheat

# Wheat **GCNVs**: rich resource to find novel/interesting **traits<->genes**!

- **all expert-curated gene families** are in expansions e.g. prolamins, NLRs, HMW glutenin...

- Overlap with mapped **QTLs**:
  - IWGSC/Salse collection: 22, families, 25 genes, 23 marker, 62 QTLs
  - comprise at least 3 connections to literature:
    - FAR1-like → Earliness
    - (u)DENN → senescence
    - Sulfate_transp::STAS → patent: enhanced traits in wheat
  - QTLs mapped to orthologs in other species → Plant **T**rait **O**ntology (**TO**)

Triticum_urartu TRIUR3_13794
Triticum_aestivum_A TraesCS7A01G148900.1
Triticum_aestivum_D TraesCS7D01G166300LC.1
Triticum_aestivum_A TraesCS4A01G032000LC.1
Triticum_aestivum_B TraesCS6B01G371400.1
Triticum_aestivum_A TraesCS3A01G713900LC.1
Triticum_aestivum_D TraesCS7D01G561400LC.1
Triticum_aestivum_B TraesCS3B01G579900.1
Triticum_aestivum_D TraesCS7D01G095600.1
Triticum_aestivum_B TraesCS7B01G320500LC.1
Triticum_aestivum_B TraesCS2B01G374100.1
Triticum_aestivum_A TraesCS5A01G371100.1
Secale_cereale Sc6Loc01864623.1
Secale_cereale Sc4Loc00156365.1
Triticum_aestivum_B TraesCS6B01G069400LC.1
Triticum_urartu TRIUR3_17083
Secale_cereale Sc6Loc01755113.1
Triticum_aestivum_A TraesCS3A01G215600LC.1
Triticum_urartu TRIUR3_06404
Secale_cereale Sc4Loc01027611.1
**Triticum_aestivum_A TraesCS5A01G528400LC.1**
Triticum_aestivum_B TraesCS6B01G789100LC.1

resource v

# The **expanded families** are enriched with previously mapped, orthologous **plant traits (TO)**



significantly enriched Plant Trait Ontology terms
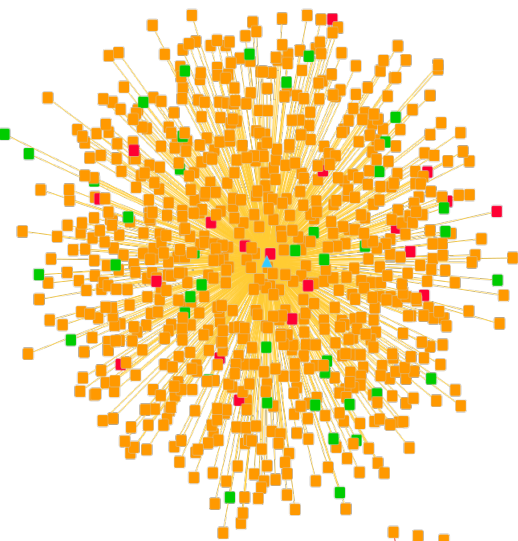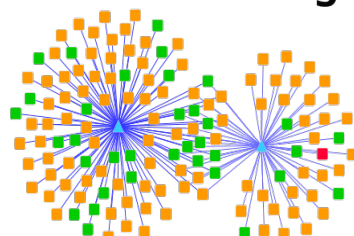
# Clustering of functional networks of enriched ontology terms: **GCNVs** affected broad range of **functions** hinting at subgenome(-lineage)-specific adaptations or subgenomic gene flow

- based on FDR<0.1 sets for gene families consistent domain architectures (>50%)
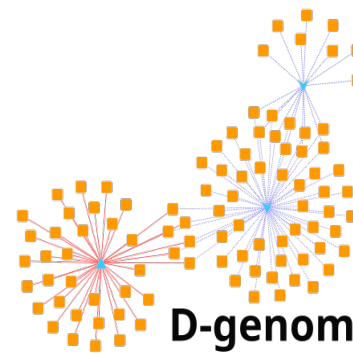- node colors:
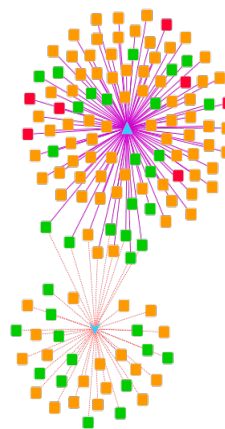  - GO orange
  - PO green
  - TO red



bread wheat expansion

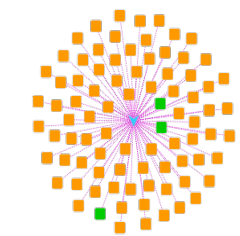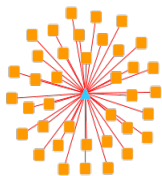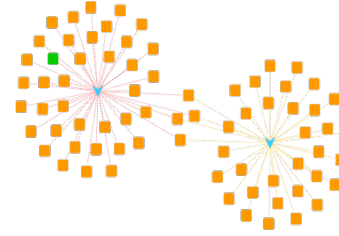B-genome expansion vs D-genome contraction

A-genome & A-lineage expansions

D-genome expansion vs A-genome & A-lineage contraction

D-genome expansion

B-genome contraction

D-genome & bread wheat contractions

resource v

# GCNVs: Defense responses expanded in D-lineage and reduced in A-lineage
## Subgenomic gene flow or lineage-specific adaptation?

# Phylogenomics/annotation resources published with the **bread wheat V1 genome papers**

- Resources:
  i. classified homeologs
  ii. orthology-based GOA|POA|TOA
  iii. wheat Transcription Associated Proteins (TAPs=TF+TR)
  iv. wheat gene families including orthologs from other Viridiplantae
  v. gene families with significant GCNVs
- all downloadable → links also in main paper supplement
  - https://wheat-urgi.versailles.inra.fr/Seq-Repository/Annotations
  - http://dx.doi.org/10.5447/IPK/2018/5
  - http://doi.org/10.1126/science.aar7191
- relevant resources were highlighted on the slides:  resource X

# Thank you!



- everyone @PGSB
- IWGSC cooperators
  - especially the comparative and gene family working groups
- for your attention :-)