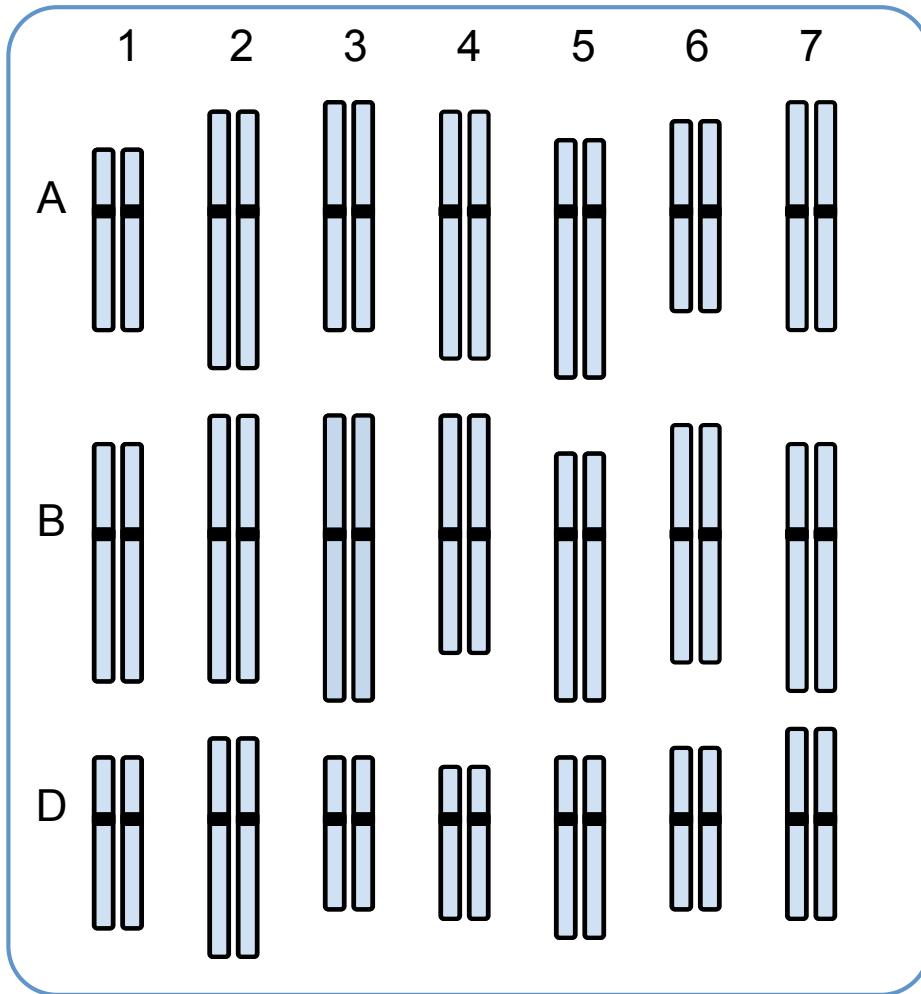# Challenges of the gene and repeat annotation of the wheat genome

Frédéric Choulet
GDEC, INRA, UCA, Clermont-Ferrand, France
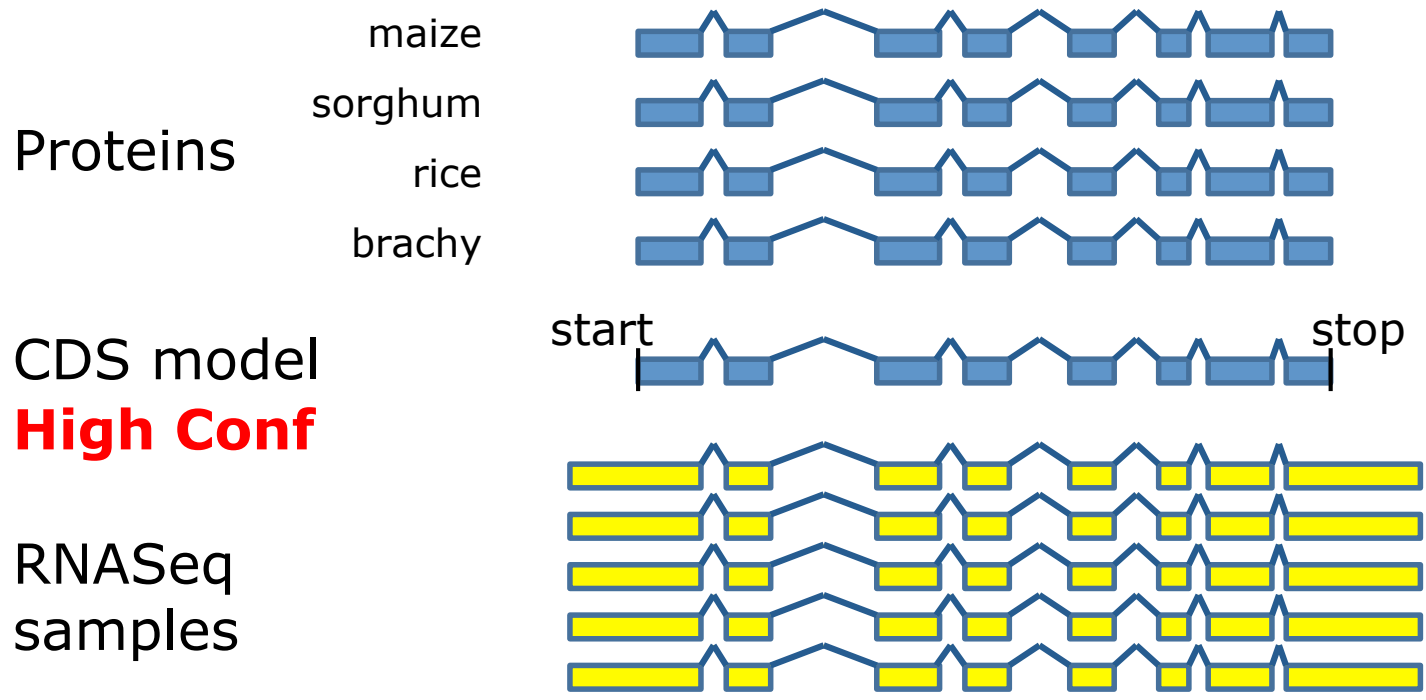
IWGSC RefSeq v1
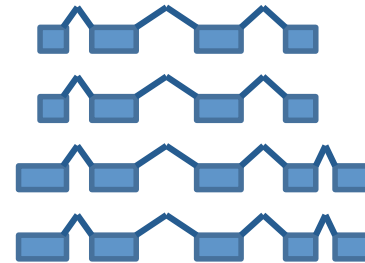21 pseudomolecules

## ❑ **TriAnnot – gene modeling**

➢ *Hélène Rimbert*
➢ *Philippe Leroy*

- o Uses 4 gene-modeling approaches
  - Augustus, FGENESH (ab-initio)
  - SIMSearch
  - BLASTX-Exonerate

- o **Scoring**
  - based on alignment with best hit in *Poaceae* proteomes
  - score = `QueryCoverage*2 + HitCoverage + %Identity`
                `+/- penalties`

➔ *1 CDS per locus*

- o Assign **confidence index** (of the gene structure)
  - **HC/LC:** agreement with mapped evidence
  - **Pseudogenes**

- o **Filtering out** doubtful predictions (=discard models with no similarity with RNASeq/IsoSeq OR proteins)

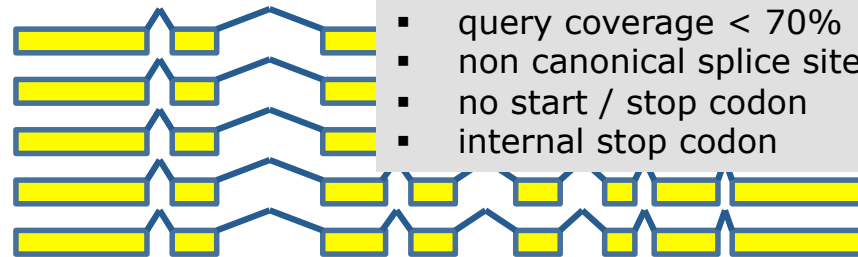- o Cufflinks (RNASeq)
➔ *splicing variants + UTRs*

Proteins

maize
sorghum
rice
brachy

CDS model
**High Conf**

start                                                                        stop

RNASeq
samples

Proteins

CDS model
**Pseudogene**

start                              stop

RNASeq
samples

- query coverage < 70%
- non canonical splice sites
- no start / stop codon
- internal stop codon

## ❑ **To do**

PGSB, Munich
➢ S. Twardziok
➢ M. Spannagl
➢ K. Mayer

IE, Norwich
➢ D. Swarbreck
➢ L. Venturini

o Protein-coding

- **Combining** 2 sets of gene predictions
    - TriAnnot
    - PGSB pipeline

- Integrate manually curated gene families

- Function assignment

o ncRNAs??????

o NTRs (novel transcribed regions)
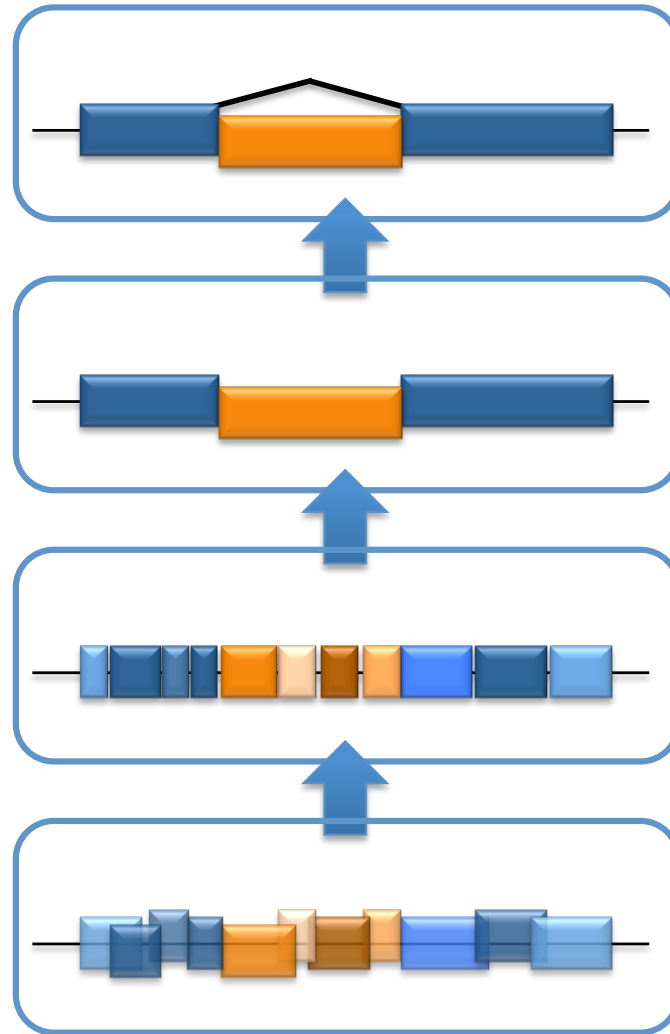
# ❑ **Homeologs / Orthologs / Paralogs**
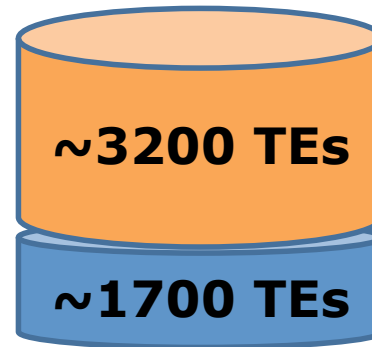
<span style="color:red">work in progress…</span>

➢ *Romain De Oliveira*

o 5 species:
wheat, rice, Brachypodium, maize, sorghum

o 3 tested approaches:
  • orthoMCL *(Li et al. Genom Res 2003)*
  • Silix *(Miele et al. BMC Bioinfo 2011)*
  • OMA *(Altenhoff et al. NAR 2015)*
=> define **homeologs, orthologs, gene families**
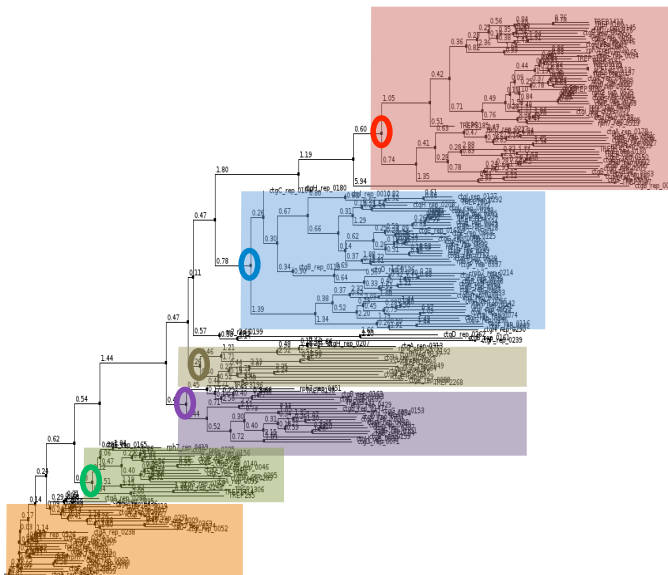
# ❏ TE modeling with CLARI-TE



*Daron et al. Genome Biol 2015*

# ☐ ClariTeRep

~3200 TEs — Choulet et al. 2010    CACTA++

~1700 TEs — TREP

=> **MCL** denovo clustering + manual curation

RLG_famc1.1

RLG_famc1.2

RLG_famc1.3
RLG_famc1.4
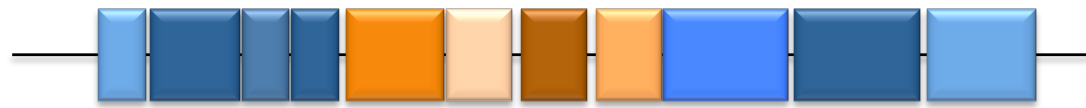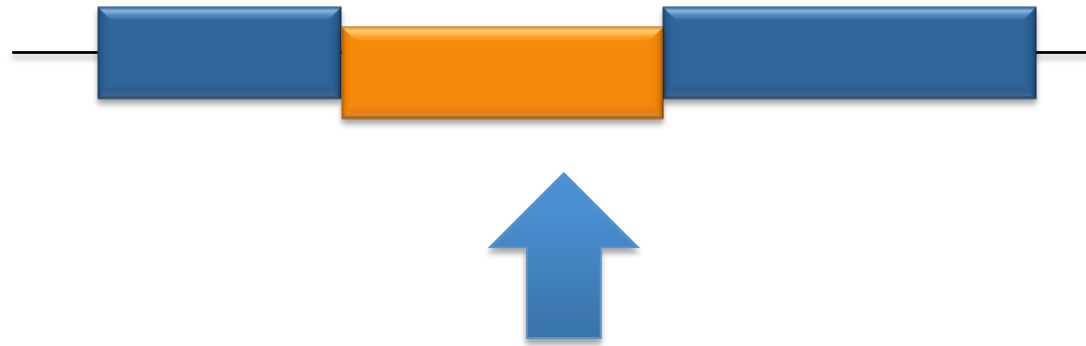RLG_famc1.5

**ClariTeRep**

~500 classified families

*download @ github.com/.../CLARI-TE*

- CLARI-TE uses new classification for defragmentation

- CLARI-TE uses new classification for better defragmentation
- CLARI-TE uses info from LTR coordinates

# TE content – IWGSC RefSeq v1

**4.0M** TEs (from 8.8M RepeatMasker matches)

**14%**

|  |  |  |  | **A** | **B** | **D** |
|---|---|---|---|---|---|---|
| • All TEs |  |  | **84.8%** | 86 | 85 | 83 |
| • ClassI | • | Gypsy | 45.2% | 51 | 47 | 41 |
|  | • | Copia | 16.1% | 17 | 16 | 16 |
|  | • | Others | 3.1% |  |  |  |
|  | • | LINE | 0.9% |  |  |  |
|  | • | SINE | 0.0% |  |  |  |
| • ClassII | • | CACTA | **15.0%** | 13 | 16 | 19 |
|  | • | Mutator | 0.4% |  |  |  |
|  | • | Mariner | 0.1% |  |  |  |
|  | • | Harbinger | 0.1% |  |  |  |
|  | • | hAT | 0.0% |  |  |  |
|  | • | Helitron | 0.0% |  |  |  |
| • Unk | • | Unk | 0.7% |  |  |  |

5 Gb —

4 Gb —

━ exons

☐ non-TE

3 Gb —

CACTA

2 Gb —

Copia

1 Gb —

Gypsy

0 Gb —
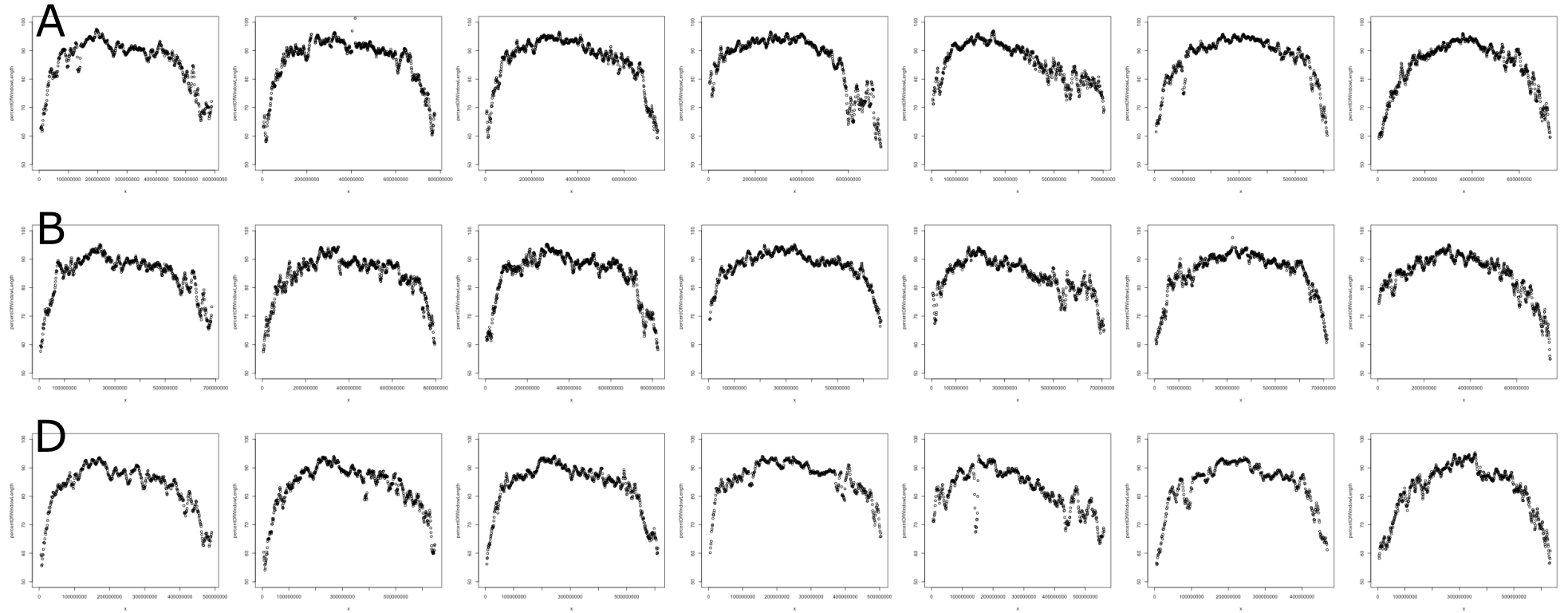
o   TE distribution

- Cereba+Quinta (centrom. gypsy families) distribution
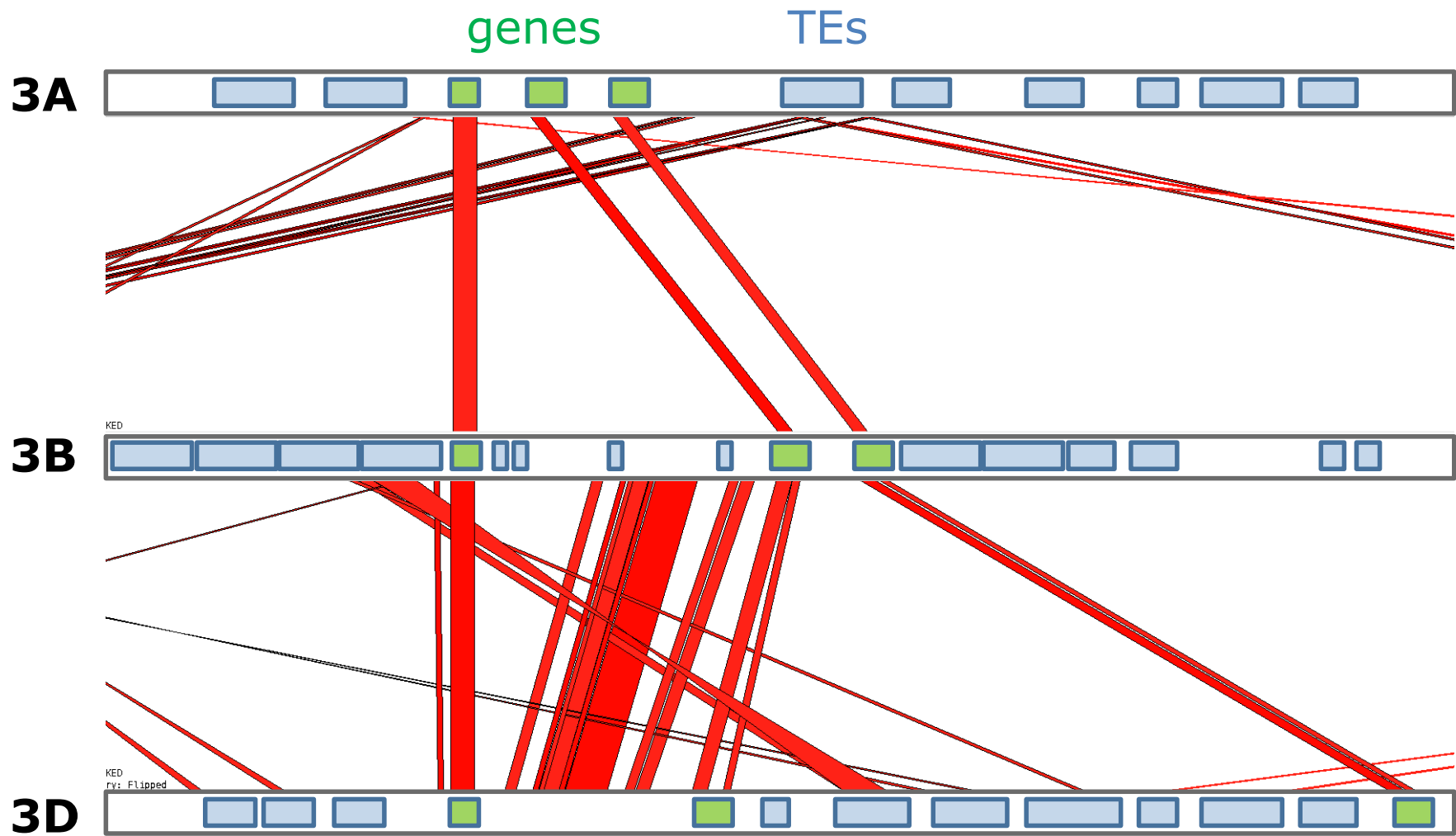
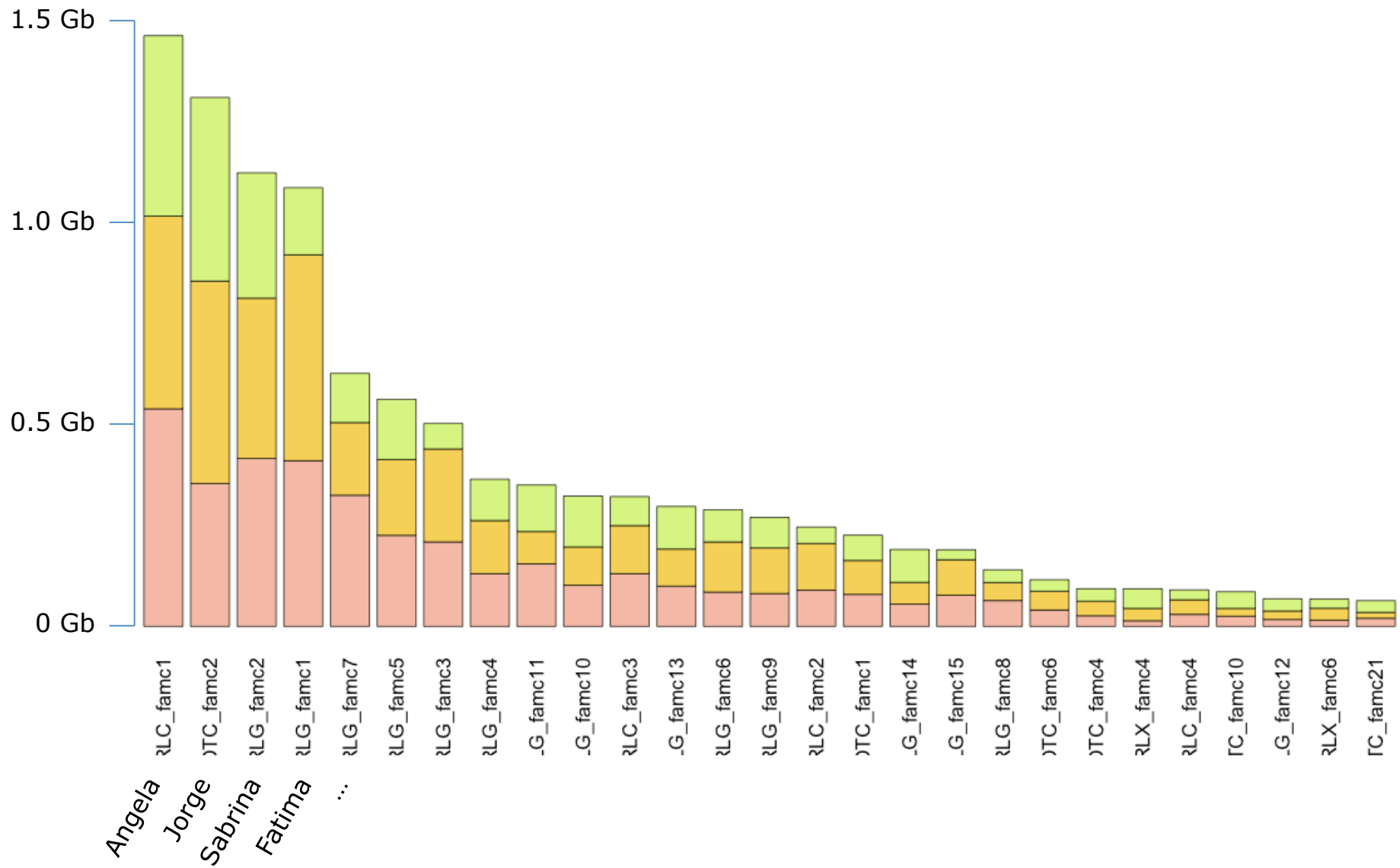A

B

D

# ❏ TE: work in progress

- **Thomas Wicker**
  - explore the unannotated part of the genome
  - chromosome "niche" specificity

- **Heidrun Gunlach**
  - genome-wide characterization of the TE content and distribution

- **Frédéric Choulet**
  - evolutionary dynamics of wheat TEs
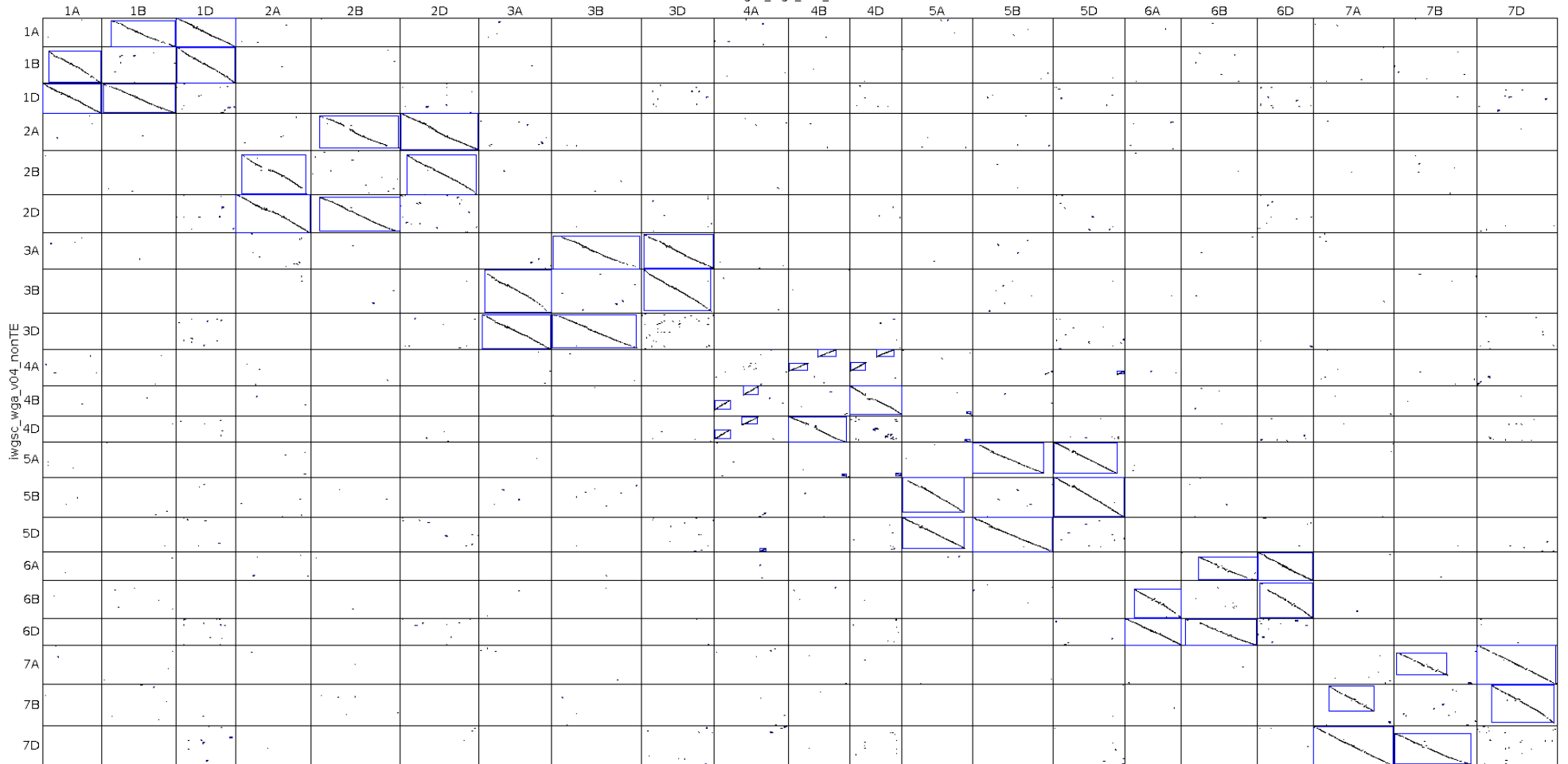  - relationships TE/genes, especially CACTAs

➢ What is the level of TE specificity in A-B-D?

1% (6/505) of TE families specific to 1 subgenome

➔ absence (almost) of subgenome specific TE families!!!

# Acknowledgments

*IWGSC-WGA working group*

*INRA GDEC*
- Hélène Rimbert
- Philippe Leroy
- Romain De Oliveira
- Etienne Paux