# Intertwining of transposable elements and non-coding RNAs in plant genomes

Douglas Silva Domingues

w/ Daniel Longhi Fernandes Pedro,  Alexandre Rossi Paschoal

São Paulo State University, Institute of Biosciences at Rio Claro, Brazil
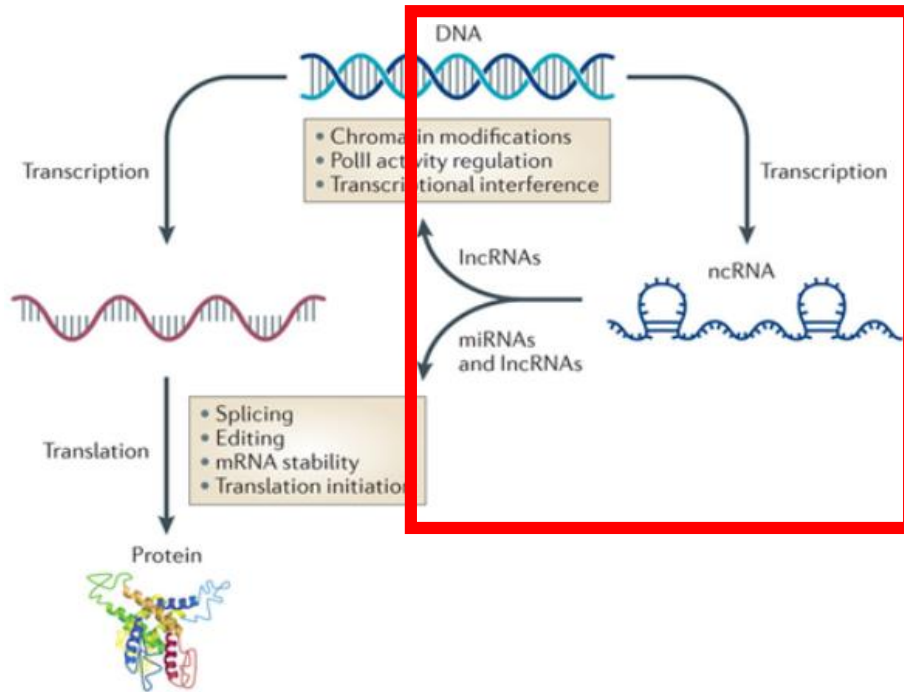**Federal Technology University of Paraná, Graduation Program in Bioinformatics, Brazil**

# KEEP CALM AND STUDY PLANTS

# Topics

- ncRNAs – Non-coding RNAs

- TEs – Transposable Elements

- PlaNC-TE: a comprehensive knowledgebase of non-coding RNAs and transposable elements in plants
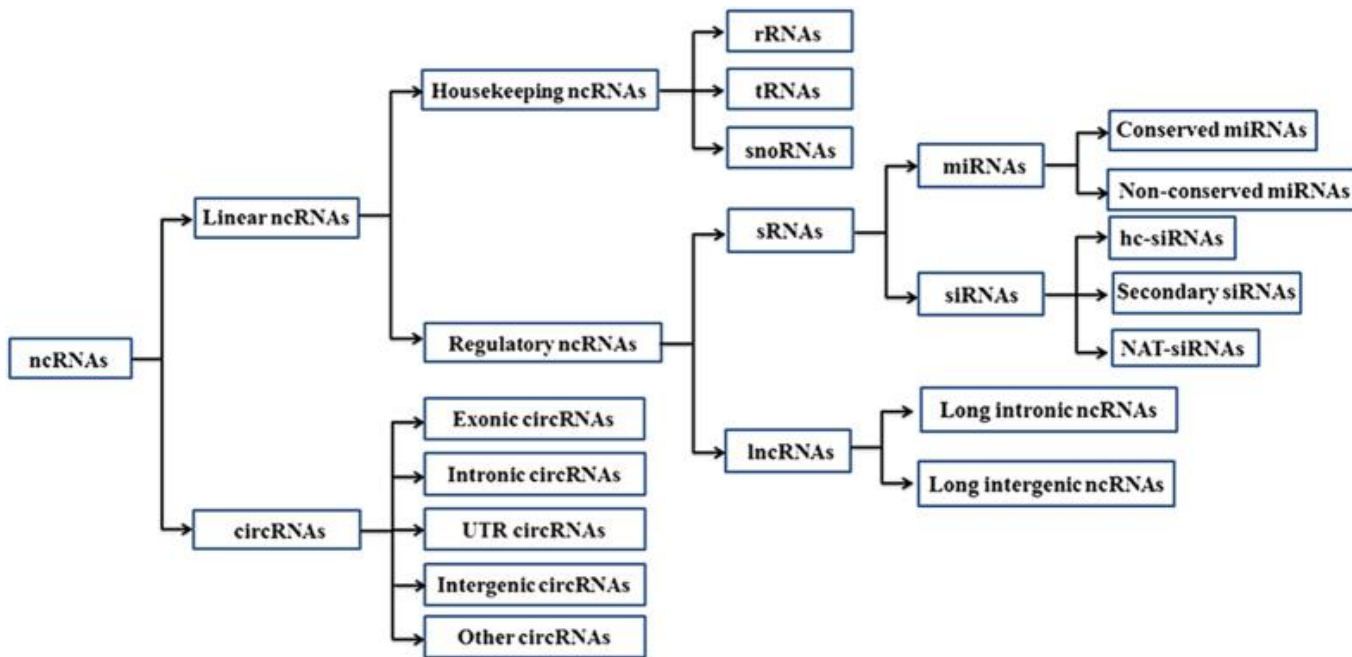
- Next steps

- Take home message

# What are ncRNAs?

# non-coding RNAs (ncRNAs)

## Sequences that are not translated into protein

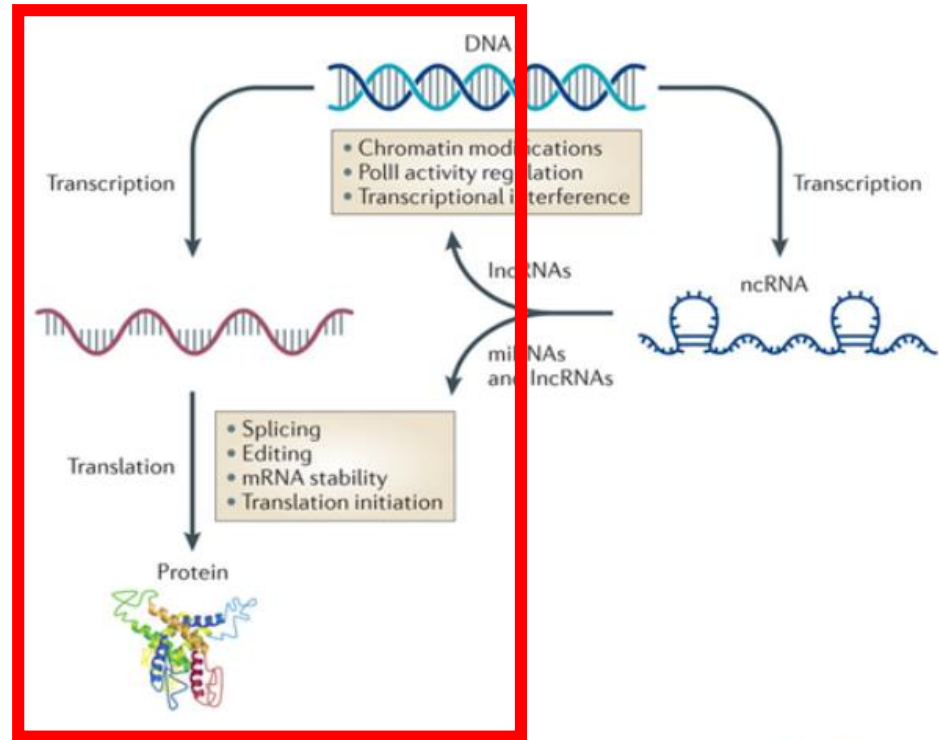Nature Reviews | Drug Discovery
Wahlestedt, 2013

5

# ncRNAs Classification



Liu et al. 2017

# Central dogma of molecular biology.

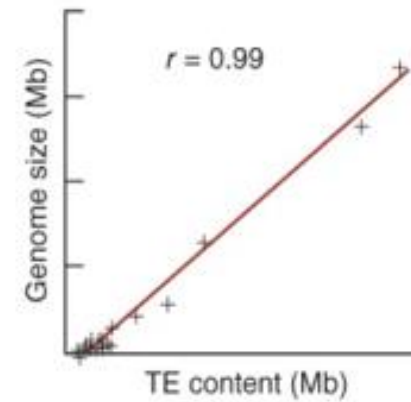DNA -> RNA -> Protein
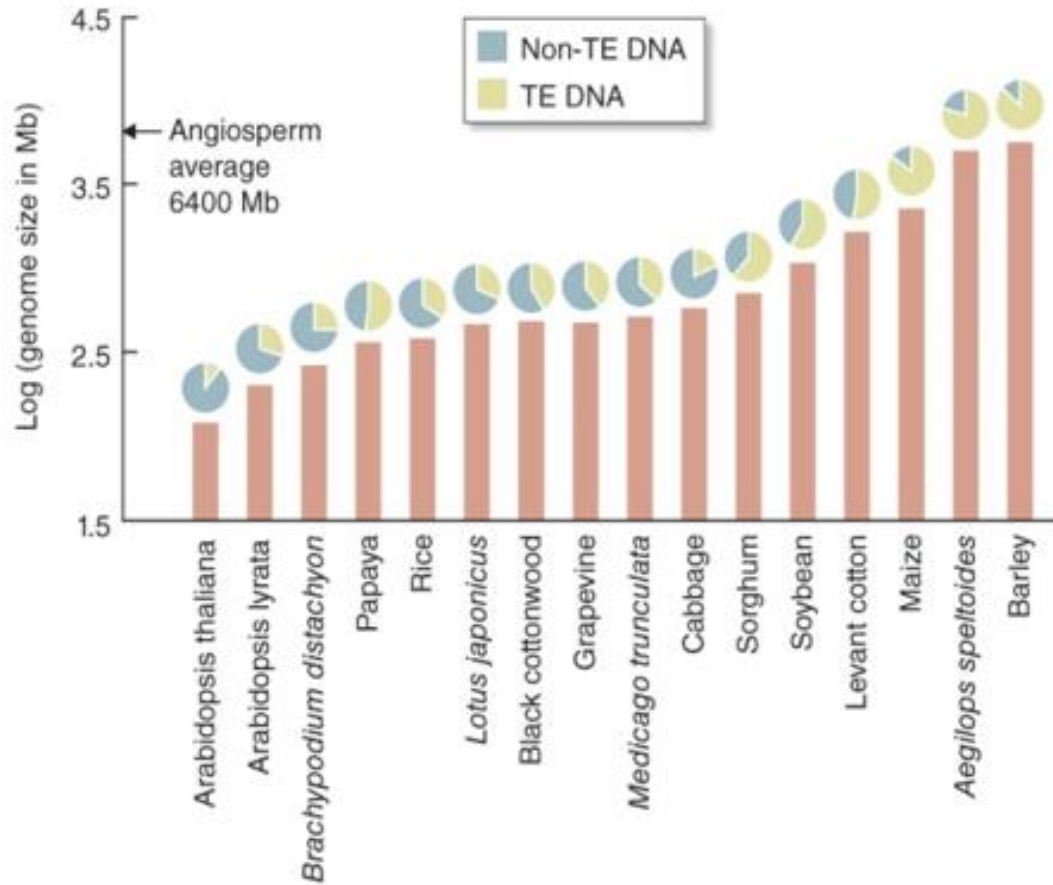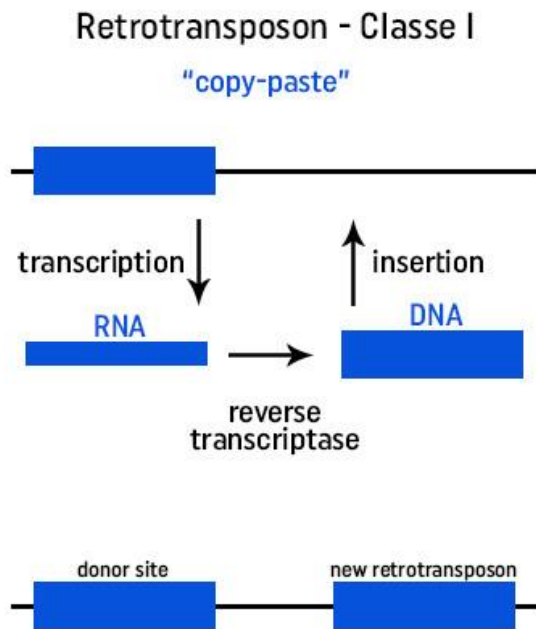


Wahlestedt, 2013

KEEP CALM.

PLANTS

HAVE

PROTEIN, TOO.

# What are TEs?

# Major componentes in plant genomes and relevant to genome size!



(a)

# Transposable Elements – Classes



Retrotransposon - Classe I
"copy-paste"

transcription ↓    insertion ↑
RNA    →    DNA
reverse transcriptase

donor site    new retrotransposon

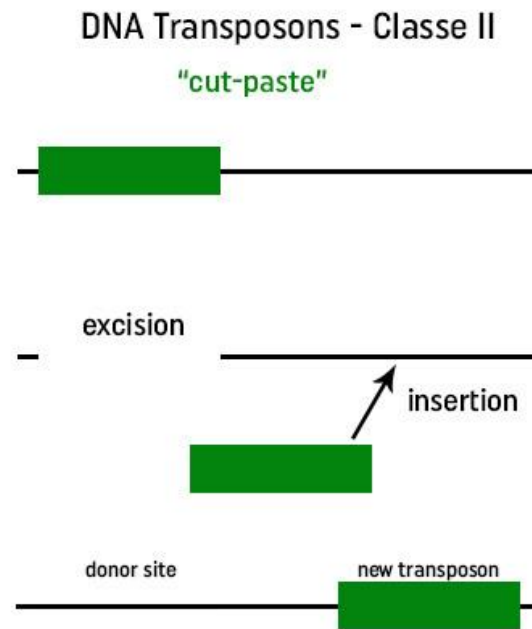DNA Transposons - Classe II
"cut-paste"

excision
insertion

donor site    new transposon

Sequences that can change their position within a genome.
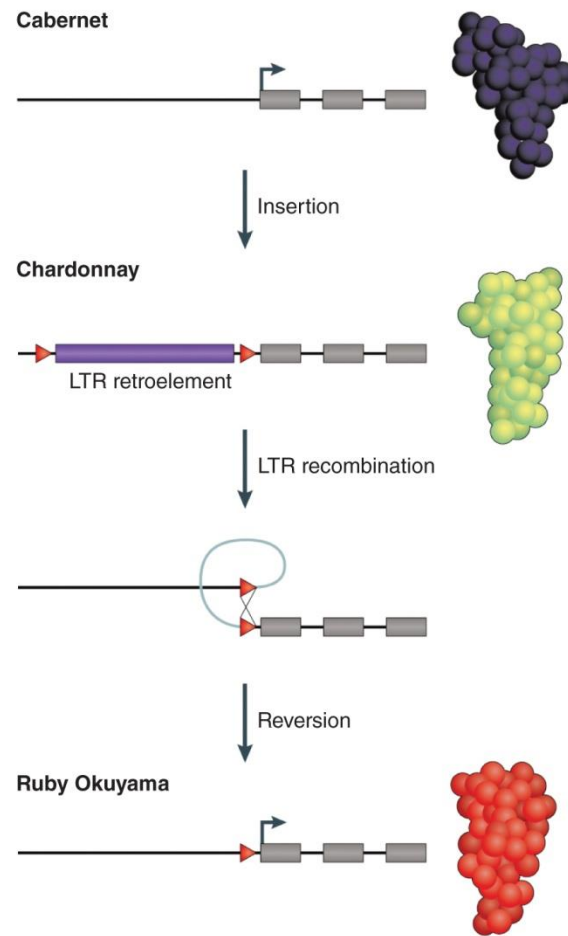
Own

11

**Figure 3.10** Control of fruit color in grapes by a retrotransposon. Cabernet grapes have a fully functional pigment gene (exons indicated by gray boxes). Insertion of a retrotransposon just upstream of the gene blocks pigment production and leads to green Chardonnay grapes. The LTRs of the element can recombine and remove most of the transposon, but one LTR is left, causing reduced transcription of the locus in Ruby Okoyama grapes. Lisch (2013). Reproduced with permission of Macmillan Publishing Ltd.

**WILEY** Blackwell

# Transposable Elements – Hierarchy



WICKER et al. 2007

Class | Order | Superfamily



**LTR**: Long Terminal Repeat
**non-LTR**: non-Long Terminal Repeat

# Why ncRNA:TEs?



Roberts et al. 2014



Adapted from: Maiti et al. 2012

# Why ncRNA:TEs?



Roberts et al. 2014

Adapted from: Maiti et al. 2012



Qin et al., 2015

# - ncRNA and TEs: known but ignored at in large-scale analyses

Genome Analysis

## Mammalian microRNAs derived from genomic repeats

### Neil R. Smalheiser and Vetle I. Torvik

University of Illinois at Chicago, UIC Psychiatric Institute, MC 912, 1601 W. Taylor Street, Chicago, IL 60612 USA

BIOINFORMATICS

## Dual coding of siRNAs and miRNAs by plant transposable elements

JITTIMA PIRIYAPONGSA and I. KING JORDAN

School of Biology, Georgia Institute of Technology, Atlanta, Georgia 30332-0230, USA

## Expression and diversification analysis reveals transposable elements play important roles in the origin of Lycopersicon-specific lncRNAs in tomato

Xin Wang, Guo Ai, Chunli Zhang, Long Cui, Jiafa Wang, Hanxia Li, Junhong Zhang and Zhibiao Ye

Key Laboratory of Horticultural Plant Biology, MOE, and Key Laboratory of Horticultural Crop Biology and Genetic improvement (Central Region), MOA, Huazhong Agricultural University, Wuhan Hubei 430070, China
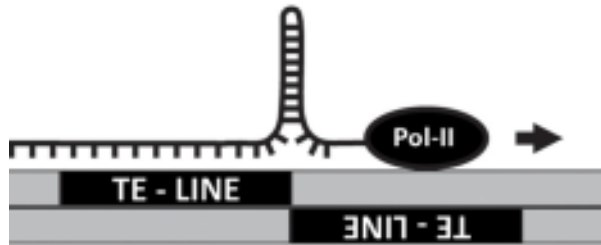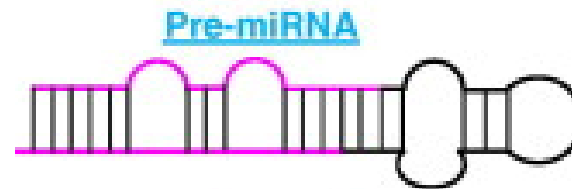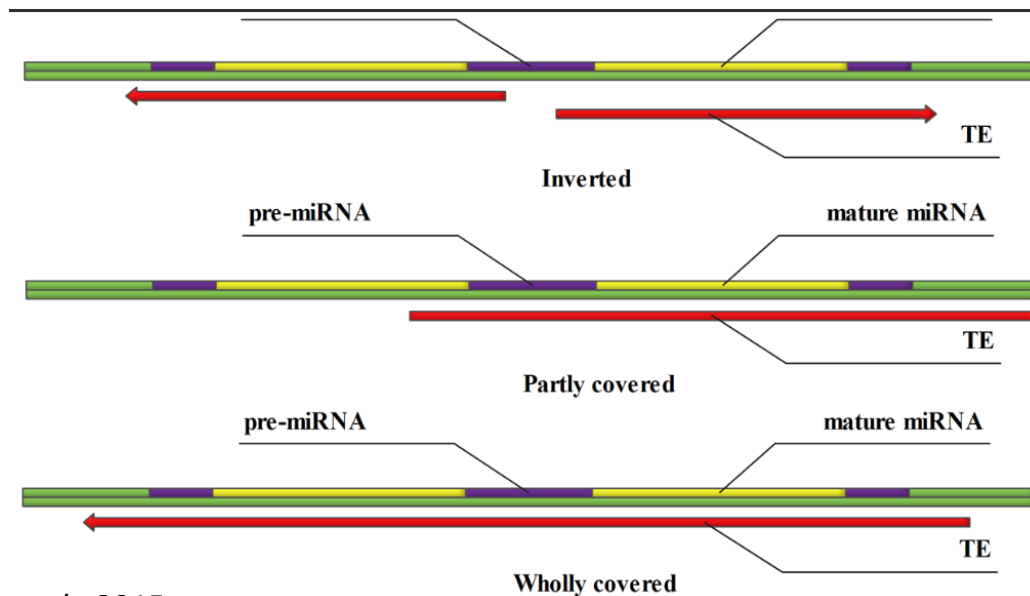
## Burgeoning evidence indicates that microRNAs were initially formed from transposable element sequences

Justin T Roberts, Sara E Cardin, and Glen M Borchert*

HYPOTHESIS

## The RIDL hypothesis: transposable elements as functional domains of long noncoding RNAs

RORY JOHNSON[1,2,3,4] and RODERIC GUIGÓ[1,2,3]

## The Role of Transposable Elements in the Origin and Evolution of MicroRNAs in Human

Sheng Qin, Ping Jin, Xue Zhou, Liming Chen, Fei Ma

Published: June 26, 2015 • https://doi.org/10.1371/journal.pone.0131365

## Transposable Elements Are Major Contributors to the Origin, Diversification, and Regulation of Vertebrate Long Noncoding RNAs

Aurélie Kapusta, Zev Kronenberg, Vincent J. Lynch, Xiaoyu Zhuo, LeeAnn Ramsay, Guillaume Bourque, Mark Yandell, Cédric Feschotte

Published: April 25, 2013 • https://doi.org/10.1371/journal.pgen.1003470

# Initial efforts

- First database to organize such information in plants

- PlanTE-MIR DB {10 *ssp*} – *v.1 - 2016*
    - miRNA:TE {152 evidences in 10 genomes}

Funct Integr Genomics (2016) 16:235–242
DOI 10.1007/s10142-016-0480-5

ORIGINAL ARTICLE

## PlanTE-MIR DB: a database for transposable element-related microRNAs in plant genomes

Alan P. R. Lorenzetti[1] · Gabriel Y. A. de Antonio[2] · Alexandre R. Paschoal[2] · Douglas S. Domingues[3]

## PlanTE-MIR DB
Plant Transposable Element-related miRNA Database

Home   About   Search   Download   Team

PlanTE-MIR DB was conceived to deliver positional intersections between annotated Transposable Elements (TEs) and pre-miRNAs in plant genomes. Our findings are useful to researchers whose discoveries relies on the study of miRNA families origin and also can contribute to the knowledge about TE regulation mechanisms.

What if we expand this to all public plant genomes with TE and ncRNA annotation data?

# http://planc-te.cp.utfpr.edu.br

## PlaNC-TE: a comprehensive knowledgebase of non-coding RNAs and transposable elements in plants 🔓

Daniel Longhi Fernandes Pedro, Alan Péricles Rodrigues Lorenzetti,
Douglas Silva Domingues, Alexandre Rossi Paschoal ✉

**OXFORD** ACADEMIC

**DATABASE**
The Journal of Biological Databases and Curation

# Objectives

- Extend PlanTE-MIR to all plant genomes available in Ensembl (53 species)

- Extend to all ncRNA classes available

- Make available a well-organized data
    - Lack of an organized repository of ncRNA:TEs for complete genomes in plants
    - Standardize outputs

- Stimulate studies in TEs and ncRNAs in plant genomes

# PlaNC-TE - *Workflow*



A) Selecting ncRNAs.

B) Filtering TEs.

C) ncRNA:TE analysis

D) PlaNC-TE webpages

21

# PlaNC-TE: A comprehensive knowledgebase of non-coding RNAs and transposable elements in plants.

- Genomic sequence source
    - Ensembl Plants
    - 53 genomes

- Retrieved sequences in ncRNAs and TEs
    - ncRNAs – 58,390 records (53 genomes)
    - TEs – 31,217,630 records (45 genomes)

# PlaNC-TE – *Phylogenetic tree*



ncRNAs  |  TEs  |  % TE overlapping

23

|  | LTR | TIR | LINE | SINE | Unknown | Total |
|---|---|---|---|---|---|---|
| **tRNA** | 2959 | 192 | 1 | 14 | 303 | 3469 |
| **rRNA** | 2962 | 1389 | 25 | 7 | 1082 | 5465 |
| **snRNA** | 1763 | 117 | 14 | 2 | 120 | 2016 |
| **Sense-intronic** | 764 | 20 | – | – | 207 | 991 |
| **Pre-miRNA** | 696 | 190 | 3 | 3 | 94 | 986 |
| **snoRNA** | 529 | 287 | 2 | 2 | 49 | 869 |
| **SRP** | 391 | 70 | – | – | 2 | 463 |
| **Antisense** | 70 | 2 | 1 | – | 16 | 89 |
| **RNase MRP** | 2 | – | – | – | – | 2 |
| **Total** | 10 136 | 2 267 | 46 | 28 | 1873 | 14 350 |

# Overall features

- Overlap features:
  - ~41% of the overlaps are among 4 genomes:
    - *Triticum aestivum; Zea mays; **Oryza sativa; Arabidopsis thaliana***

- Overlap records between ncRNA:TE | **Public data available**

  - Visualization tools (Charts by genome and jBrowse)

  - **14.350 overlaps in 40 genomes**

- Scripts developed in Perl + Bash
  - Automatic updates
  - ZendFramework2, Php7, MySQL, CSS3, HTML5, JavaScript and Debian9.

# PlaNC-TE – *Detailed info*

ncRNAs:TEs overlap

Aegilops tauschii - (276 overlaps) ▾



**TEs (inner)**
- LTR
- TIR
- SINE
- LINE
- Unknown

**ncRNAs (outer)**
- tRNA
- rRNA
- pre_miRNA
- snRNA
- snoRNA
- antisense_RNA
- sense_intronic
- SRP_RNA
- RNase_MRP_RNA

For more information see *download* section.
To export this chart, *click here.*

- First page

- Select genome

- View ncRNA:TE overlaps

# PlaNC-TE – *Examples of nc:TEs*

# PlaNC-TE – *Download*

PlaNC-TE data using Ensembl Plants v.38

| Species | ncRNAs overlapping TEs | TEs overlapping ncRNAs | Overlapped Regions | JBrowse |
|---|---|---|---|---|
| *Aegilops tauschii* (276 overlaps) | *.gff3* \| *.fa* | *.gff3* \| *.fa* | *.gff3* \| *.fa* \| *.tsv (with nts)* | *Click here* |
| *Amborella trichopoda* (219 overlaps) | *.gff3* \| *.fa* | *.gff3* \| *.fa* | *.gff3* \| *.fa* \| *.tsv (with nts)* | *Click here* |
| *Arabidopsis lyrata* (466 overlaps) | *.gff3* \| *.fa* | *.gff3* \| *.fa* | *.gff3* \| *.fa* \| *.tsv (with nts)* | *Click here* |
| *Arabidopsis thaliana* (722 overlaps) | *.gff3* \| *.fa* | *.gff3* \| *.fa* | *.gff3* \| *.fa* \| *.tsv (with nts)* | *Click here* |
| *Beta vulgaris* (3 overlaps) | *.gff3* \| *.fa* | *.gff3* \| *.fa* | *.gff3* \| *.fa* \| *.tsv (with nts)* | *Click here* |
| *Brachypodium distachyon* (180 overlaps) | *.gff3* \| *.fa* | *.gff3* \| *.fa* | *.gff3* \| *.fa* \| *.tsv (with nts)* | *Click here* |
| *Brassica oleracea* (381 overlaps) | *.gff3* \| *.fa* | *.gff3* \| *.fa* | *.gff3* \| *.fa* \| *.tsv (with nts)* | *Click here* |
| *Brassica rapa* (483 overlaps) | *.gff3* \| *.fa* | *.gff3* \| *.fa* | *.gff3* \| *.fa* \| *.tsv (with nts)* | *Click here* |
| *Chlamydomonas reinhardtii* (50 overlaps) | *.gff3* \| *.fa* | *.gff3* \| *.fa* | *.gff3* \| *.fa* \| *.tsv (with nts)* | *Click here* |
| *Chondrus crispus* (2 overlaps) | *.gff3* \| *.fa* | *.gff3* \| *.fa* | *.gff3* \| *.fa* \| *.tsv (with nts)* | *Click here* |
| *Corchorus capsularis* (546 overlaps) | *.gff3* \| *.fa* | *.gff3* \| *.fa* | *.gff3* \| *.fa* \| *.tsv (with nts)* | *Click here* |
| *Cyanidioschyzon merolae* (4 overlaps) | *.gff3* \| *.fa* | *.gff3* \| *.fa* | *.gff3* \| *.fa* \| *.tsv (with nts)* | *Click here* |
| *Galdieria sulphuraria* (1 overlap) | *.gff3* \| *.fa* | *.gff3* \| *.fa* | *.gff3* \| *.fa* \| *.tsv (with nts)* | *Click here* |
| *Leersia perrieri* (220 overlaps) | *.gff3* \| *.fa* | *.gff3* \| *.fa* | *.gff3* \| *.fa* \| *.tsv (with nts)* | *Click here* |
| *Medicago truncatula* (380 overlaps) | *.gff3* \| *.fa* | *.gff3* \| *.fa* | *.gff3* \| *.fa* \| *.tsv (with nts)* | *Click here* |
| *Musa acuminata* (251 overlaps) | *.gff3* \| *.fa* | *.gff3* \| *.fa* | *.gff3* \| *.fa* \| *.tsv (with nts)* | *Click here* |

# PlaNC-TE – *Search & Browse*

Browse by organism

Aegilops tauschii - (276 overlaps) ▼

| TEs | ncRNAs |
|-----|--------|

All - (276 nc-TE overlap) ▼

download as: .gff3 .fa

Show 10 ▼ entries                                                    Search: [          ]

| | Chr/Scaffold ▲ | Class | Type | Start | End | Strand | Length |
|---|---|---|---|---|---|---|---|
| ☐ | C135648261 | LTR/Copia | trep3125 | 1 | 209 | + | 208 |
| ☐ | C140946502 | LTR | LTR_Sb_chr_02_53 | 326 | 448 | - | 122 |
| ☐ | C141036012 | LTR | LTR_Al_scaffold_0002_362 | 408 | 500 | - | 92 |
| ☐ | C141054532 | LTR/Copia | trep3125 | 1 | 151 | + | 150 |
| ☐ | C141088066 | LTR/Copia | trep3125 | 760 | 972 | + | 212 |
| ☐ | C141262992 | LTR | LTR_Al_scaffold_0002_336 | 960 | 1071 | - | 111 |
| ☐ | C141424736 | LTR | LTR_Gm_08_2830 | 519 | 697 | - | 178 |
| ☐ | C141461244 | LTR | LTR_Gm_08_2830 | 1036 | 1214 | - | 178 |
| ☐ | C141470882 | DNA/En-Spm | trep28 | 1499 | 1517 | + | 18 |
| ☐ | C141470882 | LTR | LTR_AC198290.2_10424 | 1474 | 1748 | - | 274 |

Showing 1 to 10 of 276 entries          Previous  1  2  3  4  5  …  28  Next

download as: .gff3 .fa

# PlaNC-TE – *jBrowser*

# http://planc-te.cp.utfpr.edu.br

## PlaNC–TE: a comprehensive knowledgebase of non-coding RNAs and transposable elements in plants 🔓

Daniel Longhi Fernandes Pedro, Alan Péricles Rodrigues Lorenzetti,
Douglas Silva Domingues, Alexandre Rossi Paschoal ✉

OXFORD ACADEMIC

DATABASE
The Journal of Biological Databases and Curation

# So...

- Extend PlanTE-MIR to all plant genomes available in Ensembl (53 species)

# So…

- Extend PlanTE-MIR to all plant genomes available in Ensembl (53 species)

    - But TEs data are available only for 40 genomes!

# So…

- Extend PlanTE-MIR to all plant genomes available in Ensembl (53 species)

    - But TEs data are available only for 40 genomes!

    - If TEs are a major component of genomes, something is wrong!

# Phase 2: Re-annotation of TEs in complete plant genomes

Dataset

8 Genomes

Identification

| 1A | RepeatScout | *PASTEClassifier* |

| 1B | RepeatModeler |

Class I
Class II
Both

| 2 | RepeatMasker | *Repbase Library* | *R.Modeler Library* | *R.Scout Library* |

| 3 | HelitronScanner | MITE-Hunter |

| 4 | LTR_retriever | MGEScan-non-LTR |

Filter

Low complexity
Simple repeat

Annotation

Genome
annotation
for TEs

# Initial analyses raised the number of TE entries in genomes

| Plant genomes | Ensembl Plants | Our approach |
|---|---|---|
| *A. lyrata* | 116,145 | 391,425 |
| *A. thaliana* | 43,442 | 63,879 |
| *B. vulgaris* | 6,295 | 984,280 |
| *B. rapa* | 97,576 | 434,231 |
| *C. sativus* | - | 176,333 |
| *M. acuminata* | 116,189 | 637,112 |
| *P. trichocarpa* | 248,622 | 864,831 |
| *V. vinifera* | 281,476 | 834,298 |

# Take home message

- We still need standardization and better annotation (at least of TEs) in plant genomes

- Up to now, TE annotation is heavily based in alignment: curated datasets can be an starting point for other computational approaches

- Long-term goal: Are any specific characteristic (*feature*) of TE and/or ncRNA that distinguish ncRNA:TE association?

# Team/Acknowledgements



Daniel L. F. Pedro

Dr. Alexandre Rossi Paschoal

Alan P. Rodrigues Lorenzetti

Tharcísio Amorim

Funding

# Thank you!

e-mail: **douglas.domingues@unesp.br**