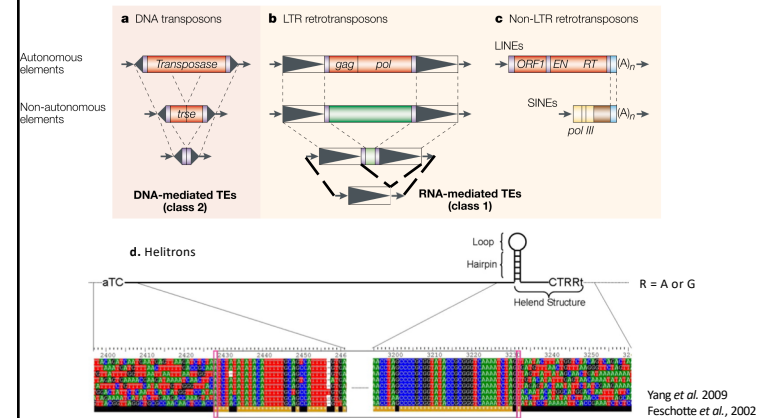
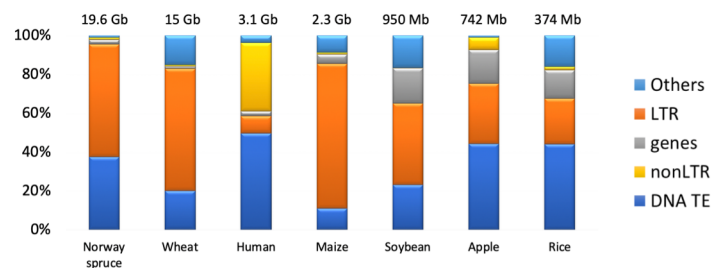


Transposable Element Annotation

Annotating Transposable Elements



TEs are prevalent in eukaryotic genomes



Approaches to TE Annotation

- Copy-number based
 - 👍: sensitively identify repetitive sequences
 - 🙋: not knowing what are the sequences, could be duplicated genes
- Homology based
 - 👍: reuse prior knowledge; quick
 - 🙋: limited by prior knowledge; many TEs sequences are not conserved
- Structural based
 - 👍: codable; independent of database
 - 🙋: limited by knowledge of sequence structures; high FDR

Structural Features of TEs

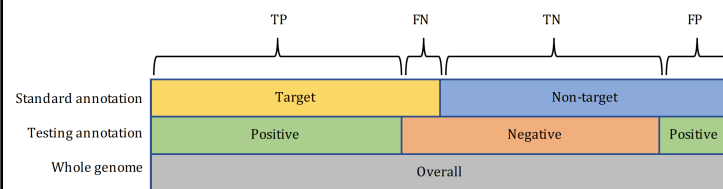
| Class | Superfamily | Traditional name | Systematic classification | TSD | Element size | TIR | Terminal or sequence size |
|----------|-----------------------|------------------|---------------------------|------------|----------------|----------------|---------------------------------------|
| | | | | | Autonomous | Non-autonomous | |
| Class I | L1R-Copia | RLC | 5 bp | 4–11 kb | 202 bp to 4 kb | 85 bp to 3 kb | TG...C/G |
| | L1R-Gypsy | RLG | 5 bp | 5–20 kb | 2–5 kb | 120 bp to 6 kb | TG...CA |
| | LINE | R1L/R1 | Variable | Up to 9 kb | NA | NA | ...AAAA |
| | SINE | RST/RSL | Variable | NA | 80–500 bp | NA | ...AAAA |
| Class II | Au/Dn/Alu | DFA | 8 bp | 3–6 kb | 110 bp–3 kb | 5–22 bp | C/TA...TA |
| | Enj/Spm/dSpm/CACTA | DTC | 3 bp | 6–21 kb | 200 bp–6 kb | 12–28 bp | CACTA/G...C/TAGTGTG |
| | MuDR/Mustator/Mu/MULE | DTM | 7–11 bp, mostly 9 | 4–16 kb | 120 bp–3 kb | 6–800 bp | G/C...G/C |
| | P1/Harbinger/Tourist | DTH | TNA | 3–7 kb | 80 bp–3 kb | 14–60 bp | GGG/CC...GG/CCC GG/AGCA TGC/TCC |
| | T1/Mariner/Stowaway | DTT | TA | 3–7 kb | 80 bp–3 kb | 11–120 bp | CTCCCTC...GGAGGG |
| | Helicon | DHH | None | 5–17 kb | 150 bp–20 kb | None | TC...CTR |

(Zhao, D., A. Ferguson, N. Jiang, BBA Gene Regulatory Mechanisms 2016)

Many Different TE Annotation Software

- General TE annotator
 - RepeatModeler, RepeatScout, RepeatMasker, RepBase, Red, RECON, CENSOR, PILER, REPET, GenericRepeatFinder (GRF)
- LTR retrotransposons
 - LTR_STRUC, LTR_MINER, LTR_FINDER, LTRharvest, MGEScan_LTR/LTRrho, LTR Annotator, LTR_retriever, LtrDetector, GRF
- Non-LTR retrotransposons
 - SINE-Finder, SINE_scan, SINEBase, MGEScan_nonLTR, TSDFinder
- MITE
 - MITE-Hunter, detectMITE, IRF, TIRvish, GRF, miteFinderII, P-MITE, TIR-Learner, MITE-Tracker, MUSTv2, MITE Digger, RSPB
- Helitron
 - HelSearch2, HelitronFinder, HelitronScanner

How to Pick a Software?



$$\text{Sensitivity} = P(\text{positive}|\text{target}) = \frac{TP}{TP + FN}$$

$$\text{Specificity} = P(\text{negative}|\text{non_target}) = \frac{TN}{TN + FP}$$

$$\text{Accuracy} = P(\text{true_classification}) = \frac{TP + TN}{TP + TN + FP + FN}$$

$$\text{Precision} = P(\text{target}|\text{positive}) = \frac{TP}{TP + FP}$$

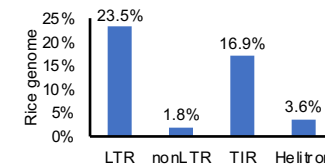
$$F_1 = \frac{2 * \text{Sensitivity} * \text{Precision}}{\text{Sensitivity} + \text{Precision}} = \frac{2TP}{2TP + FP + FN}$$

$$\text{FDR} = P(\text{non_target}|\text{positive}) = \frac{FP}{TP + FP}$$

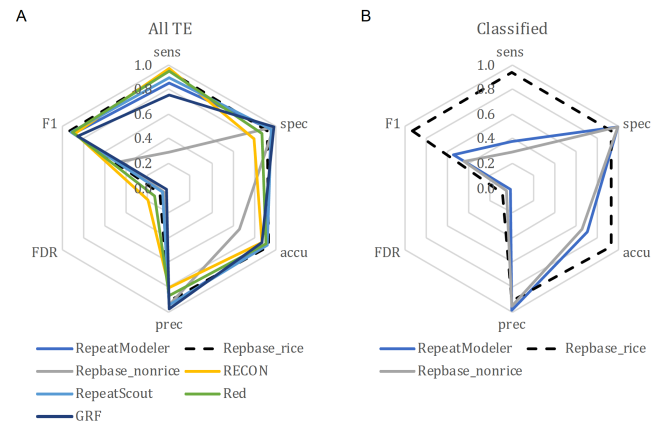
Ou et al. (2019) Genome Biol.

Subject genome – *Oryza sativa* (rice)

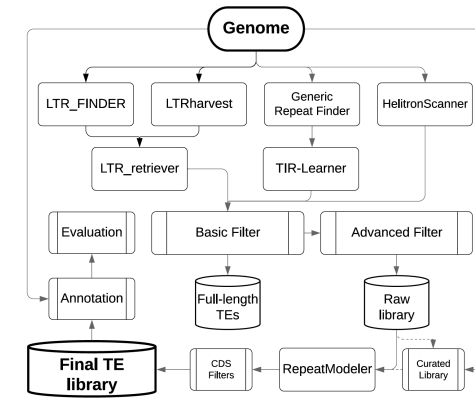
- High-quality TE space assembly (LAI = 20, “gold” quality)
- Harbors wide range of TEs in reasonable abundance (46%)
- Small genome (375 Mb) allows quick analyses
- Has a manually curated TE library.



Example of a Full TE Annotation Pipeline



Ou et al. (2019) *Genome Biol.*



Ou et al. (2019) Genome Biol.

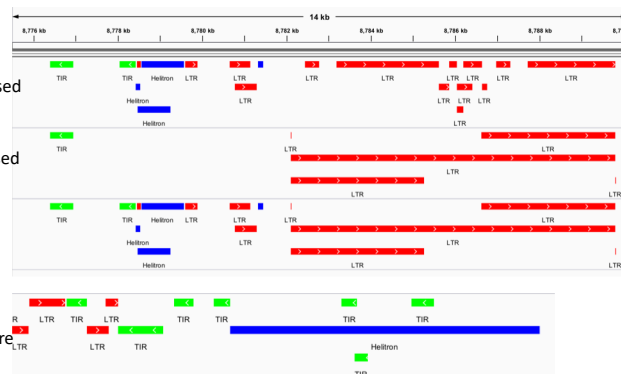
Example of a Full TE Annotation Pipeline

Homology-based

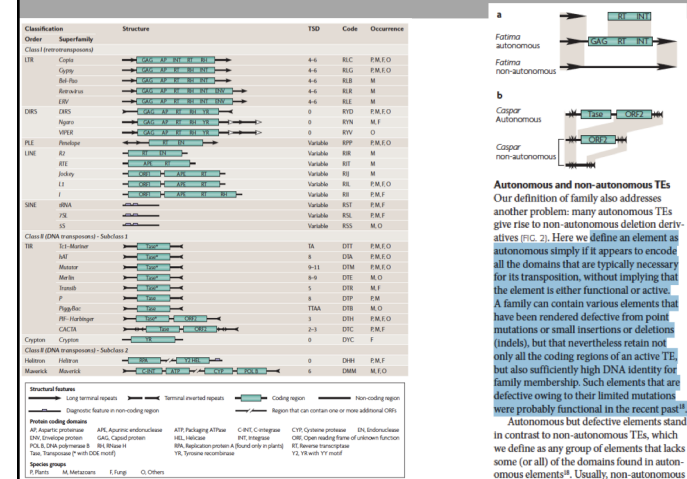
Structural-based

Combined

Nesting/capture



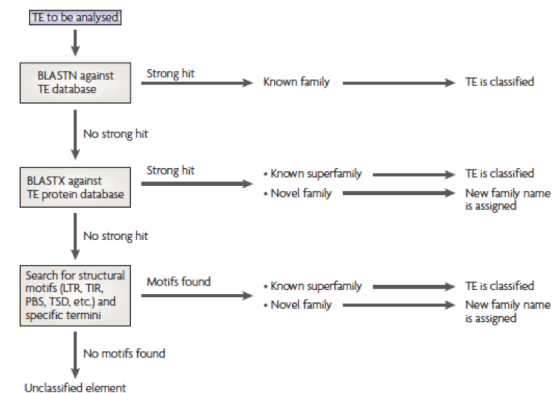
Autonomous vs. Non-Autonomous



80-80-80 Rule

- Used to define family and subfamily
- TE Family – group of TEs that have high DNA sequence similarity in their coding region/internal domain, or in their terminal repeat region
- **80%** sequence similarity in at least **80%** of aligned sequence when considering segments of **80bp** or longer

Superfamily and Family Classification



TE Annotation Takeaways

- Like gene annotation – there are rules...but they aren't as defined as they are for genes
- Annotations are done using structure and homology
- Different programs used to annotate different TE orders each with different accuracy levels
- Family classification important for downstream biological interpretation