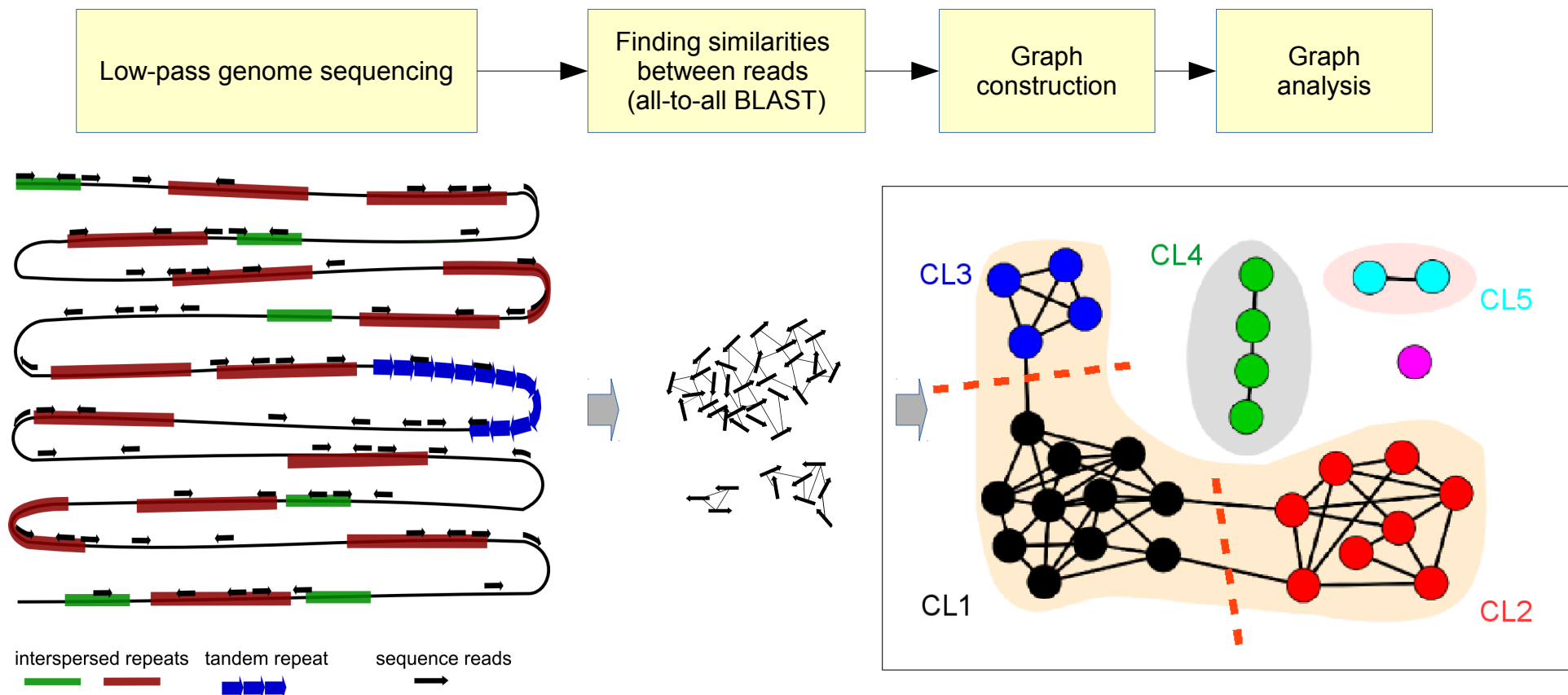# Design of sequencing and repeat analysis experiments

*Platform and coverage*

- Illumina is preferred, use **paired-end reads**

- avoid coverage > 1x (→ similarity hits from single/low copy sequences) , **0.1-0.5x is optimal**
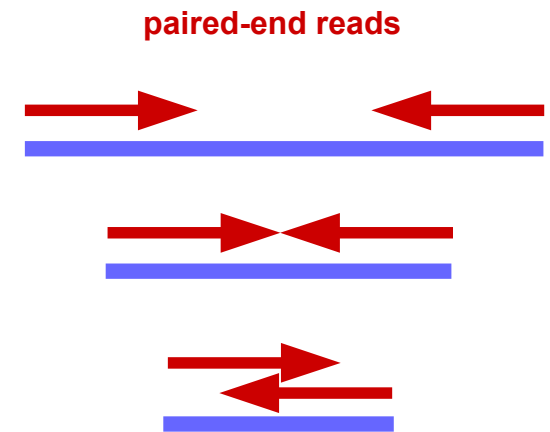
# Design of sequencing and repeat analysis experiments

*Platform and coverage*

- Illumina is preferred, use **paired-end reads**

- avoid coverage > 1x (→ similarity hits from single/low copy sequences) , **0.1-0.5x is optimal**

*Read length vs. fragment length*

- sequenced fragments should be > 2 x read length



**paired-end reads**

*Consider eventual bias in template preparation*
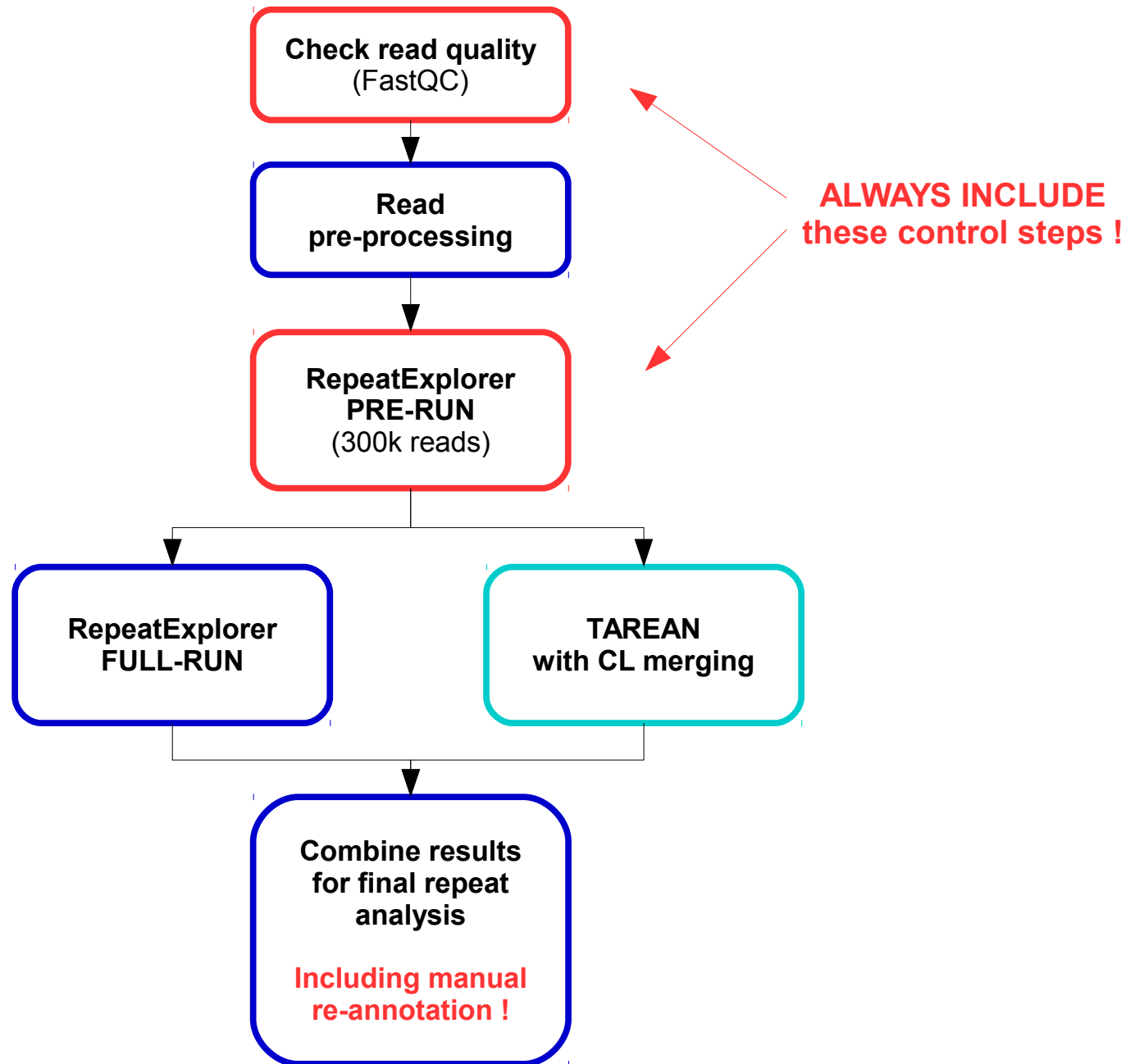
- PCR-based   x   PCR-free kits

- avoid transposon-based kits

*Spend some more money to make <u>experimental replicates for quantification</u>*

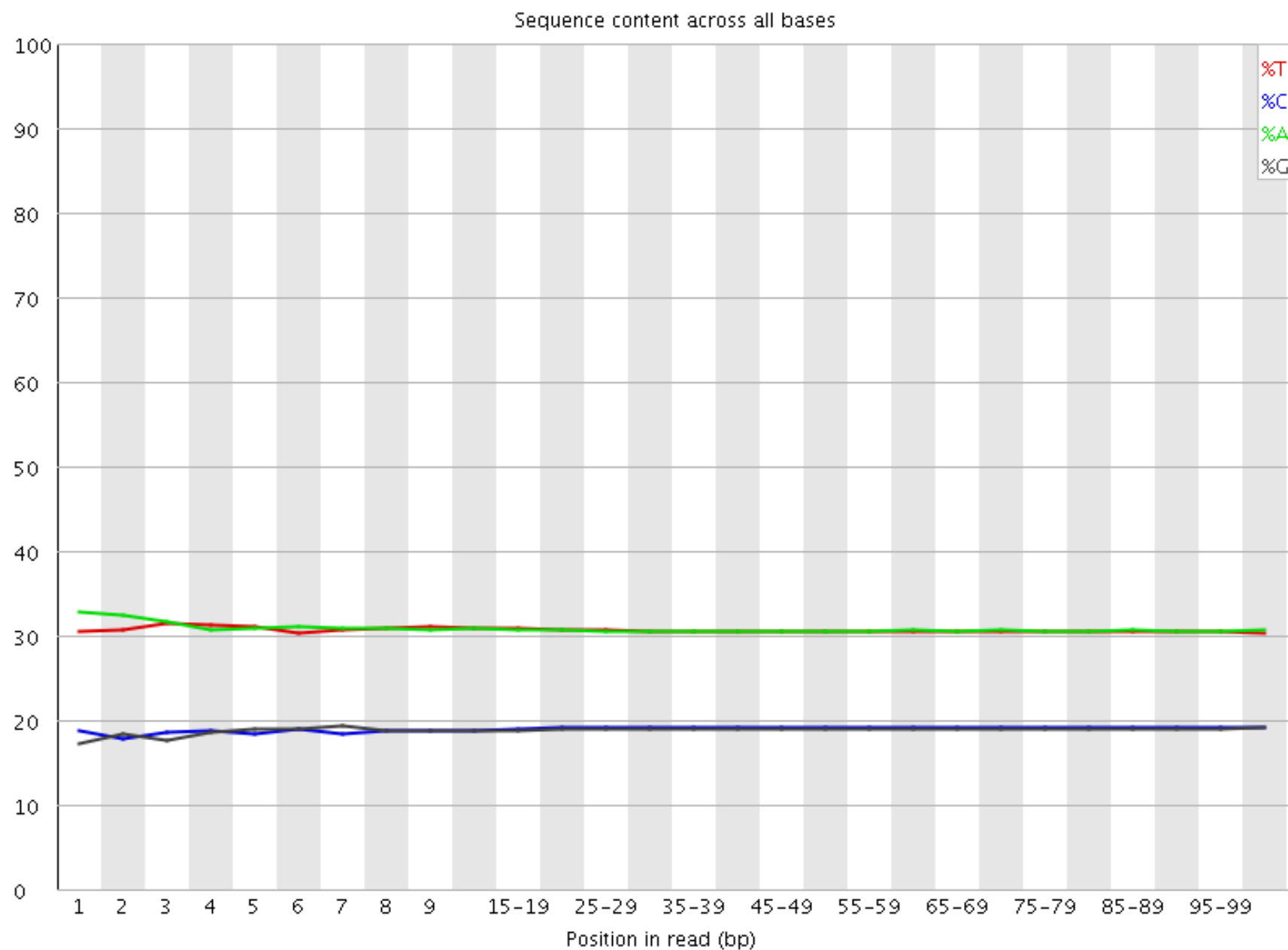- repeating library preparations (not re-sequencing from the same library)

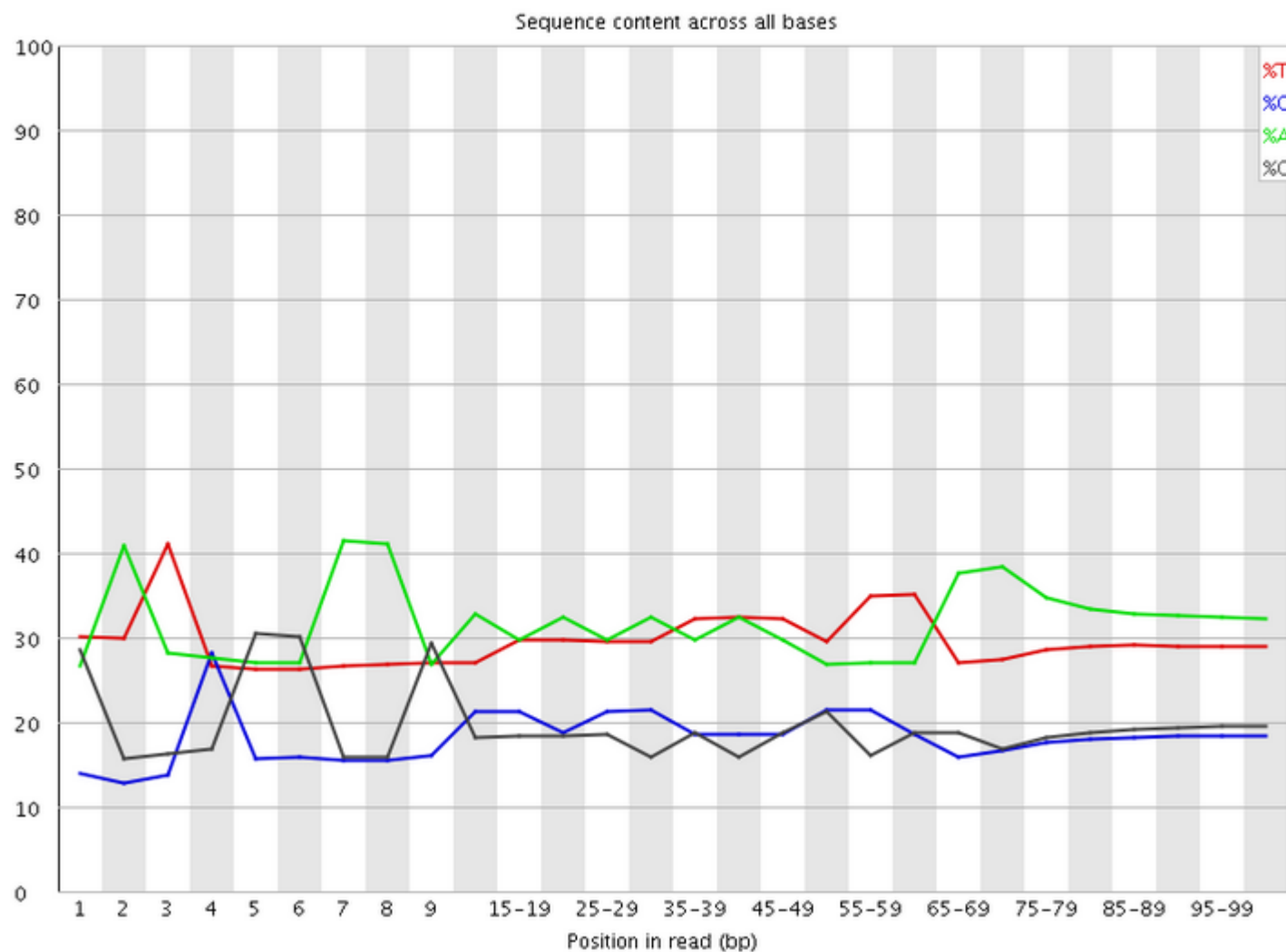**CHECK YOUR INPUT DATA !  (and do not trust public sequence archives)**

# Analysis workflow

```
         ┌─────────────────────┐
         │  Check read quality  │
         │      (FastQC)        │
         └─────────────────────┘
                    │
                    ▼
         ┌─────────────────────┐
         │        Read         │
         │   pre-processing    │
         └─────────────────────┘
                    │
                    ▼
         ┌─────────────────────┐
         │   RepeatExplorer    │
         │      PRE-RUN        │
         │     (300k reads)    │
         └─────────────────────┘
              │          │
              ▼          ▼
   ┌──────────────┐  ┌──────────────┐
   │ RepeatExplorer│  │    TAREAN    │
   │   FULL-RUN   │  │ with CL merging │
   └──────────────┘  └──────────────┘
              │          │
              └────┬─────┘
                   ▼
         ┌─────────────────────┐
         │  Combine results    │
         │  for final repeat   │
         │      analysis       │
         │                     │
         │  Including manual   │
         │   re-annotation !   │
         └─────────────────────┘
```

**ALWAYS INCLUDE
these control steps !**

## Per base sequence content



Sequence content across all bases

# ⊗ Per base sequence content

**BAD !**



Sequence content across all bases

| | %T |
| | %C |
| | %A |
| | %G |

Position in read (bp)

# ⊗ Overrepresented sequences

| Sequence | Count | Percentage | Possible Source |
|---|---|---|---|
| GATCGGAAGAGCACACGTCTGAACTCCAGTCACACTGATATATCTCGTAT | 1743087 | 10.335846049950064 | TruSeq Adapter, Index 5 (97% over 37bp) |
| GATCGGAAGAGCACACGTCTGAACTCCAGTCACACTGATATATATCGTAT | 36128 | 0.21422536344576945 | TruSeq Adapter, Index 5 (97% over 37bp) |
| CGGAAGAGCACACGTCTGAACTCCAGTCACACTGATATATCTCGTATGCC | 35955 | 0.21319953893635515 | TruSeq Adapter, Index 5 (97% over 34bp) |
| GATCGGAAGAGCACACGTCTGAACTCCAGTCACACTGCTATATCTCGTAT | 22960 | 0.13614410830145224 | TruSeq Adapter, Index 5 (97% over 37bp) |

# FASTX-Toolkit: Compute quality statistics -> Draw nucleotides distribution chart
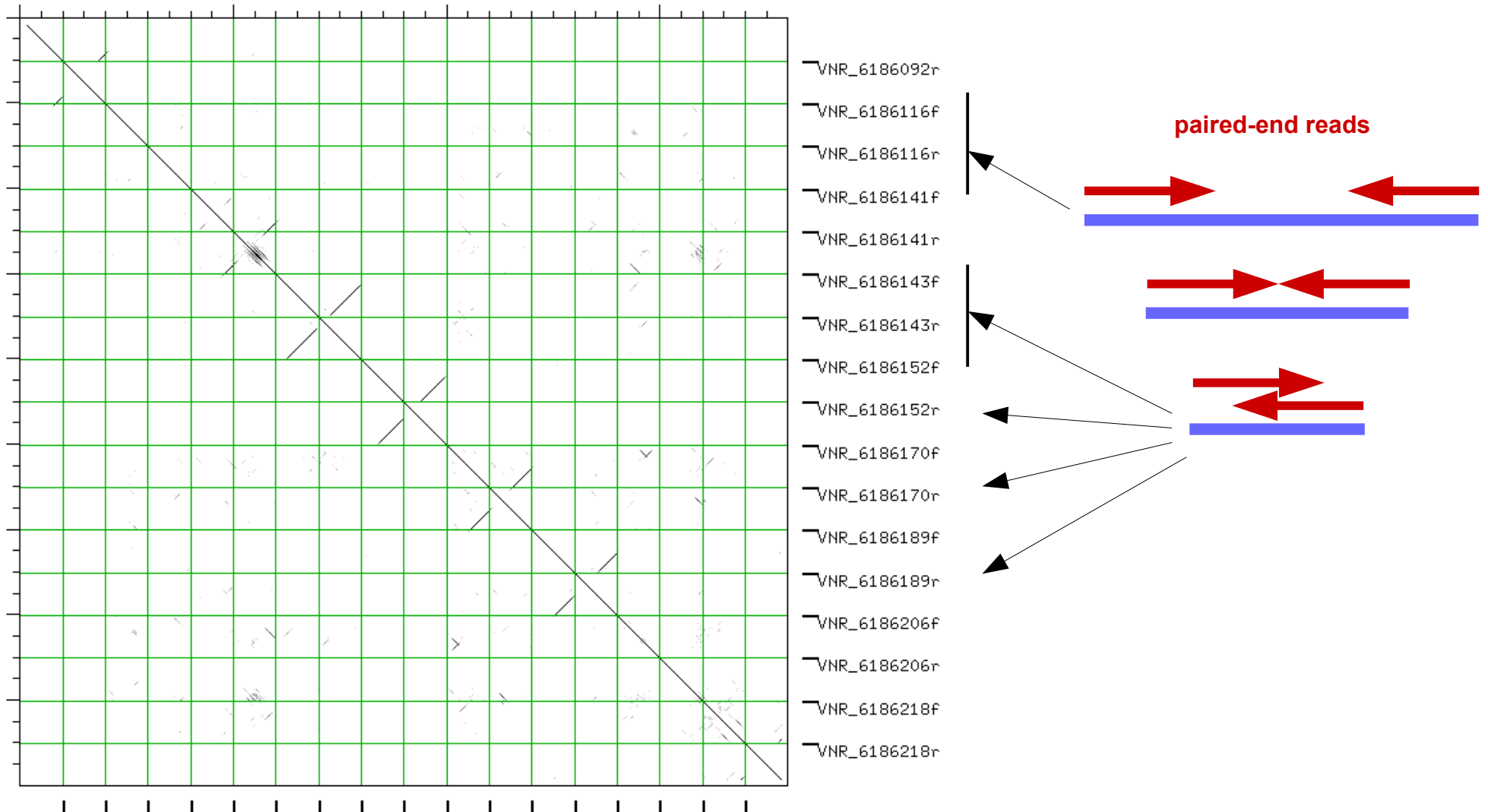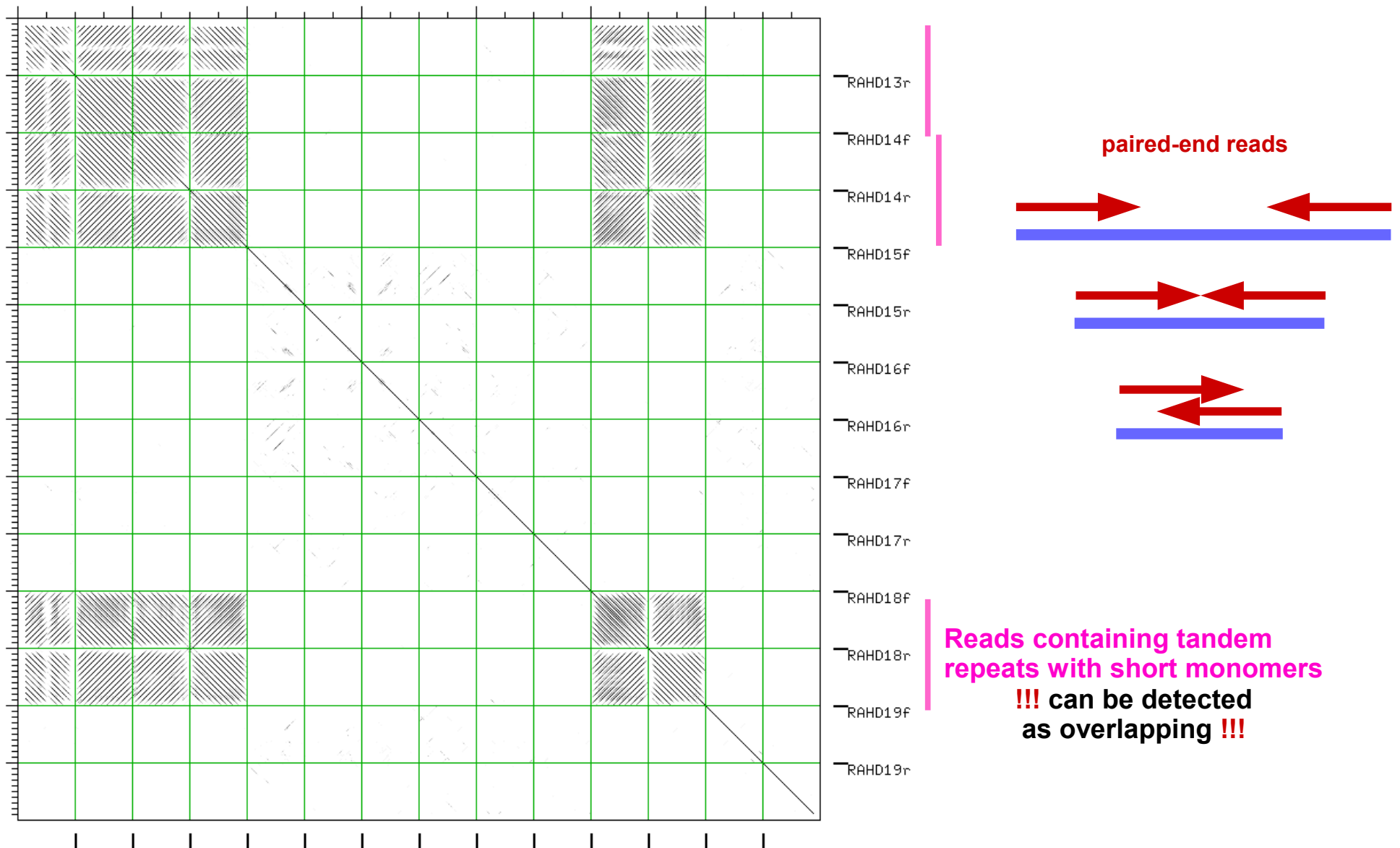(alternative to FastQC)

**OK**



**Adaptor contam.**

# Read length vs. fragment length

- sequenced fragments should be > 2 x read length

# *Read length vs. fragment length*

- sequenced fragments should be > 2 x read length



RAHD13r

RAHD14f

RAHD14r

RAHD15f

RAHD15r

RAHD16f

RAHD16r

RAHD17f

RAHD17r

RAHD18f

RAHD18r

RAHD19f

RAHD19r

**paired-end reads**

**Reads containing tandem repeats with short monomers**
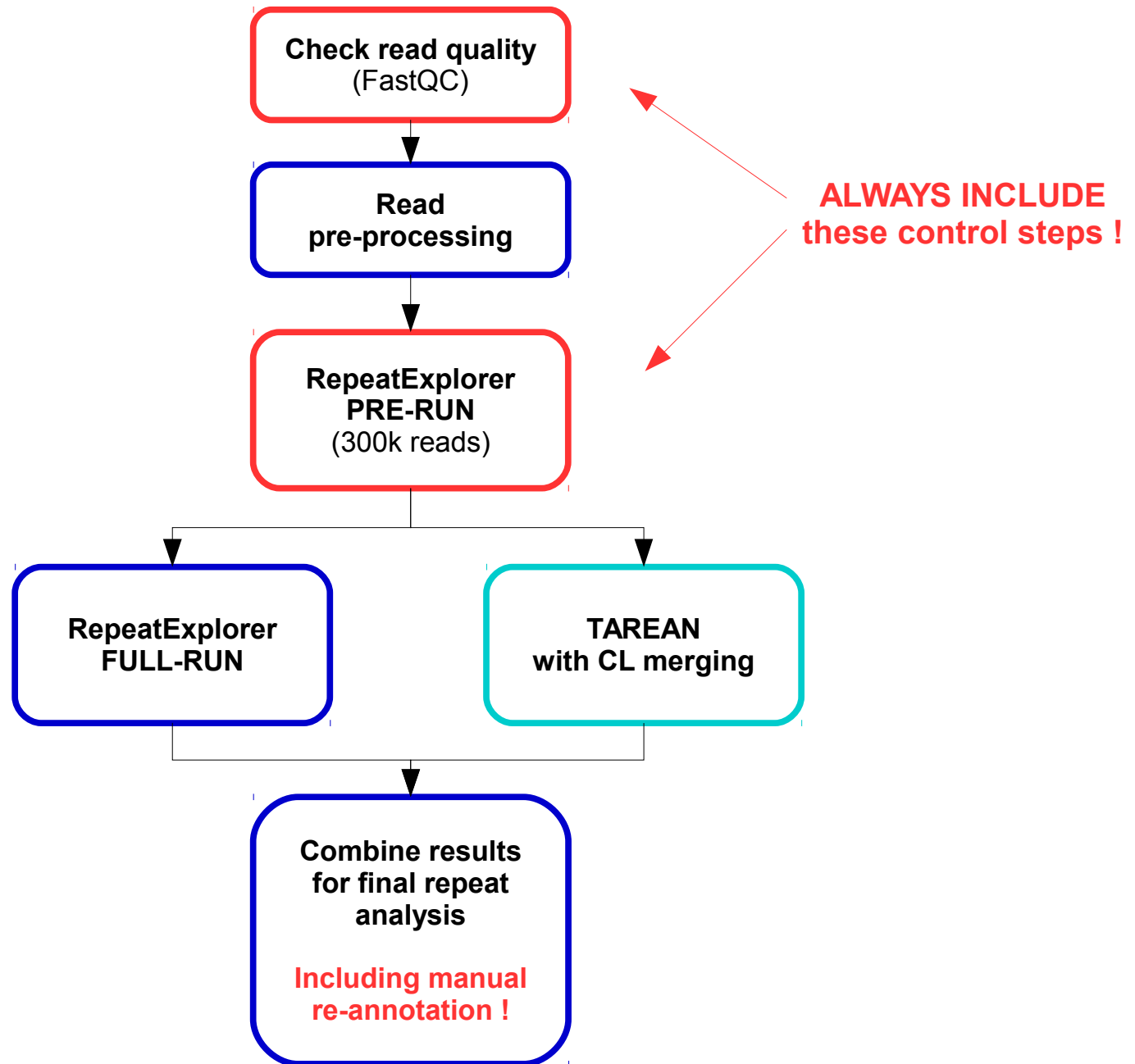**!!! can be detected as overlapping !!!**

# Adela (*Adelhaide kratzmarii*)

- Fictional plant from a popular Czech movie

- We reconstructed its genome ;-)

- It turned out to be very small but made only by repeats !

- EXCELLENT for training at this workshop

- Do not expect such nice results with your real data
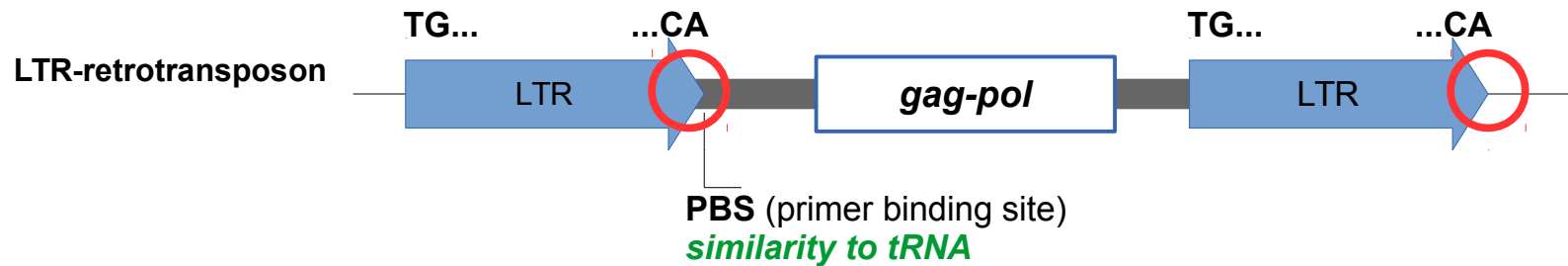
- On the other hand, your real model will not eat you

# Analysis workflow

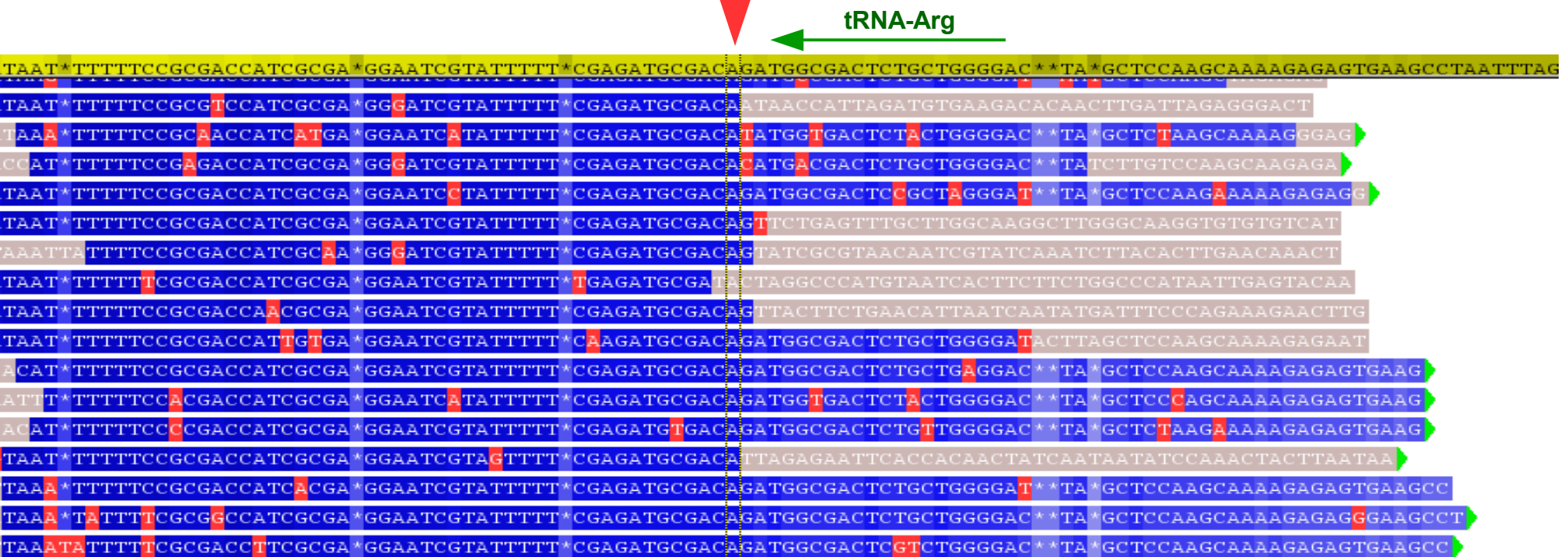# Insertion sites of mobile elements
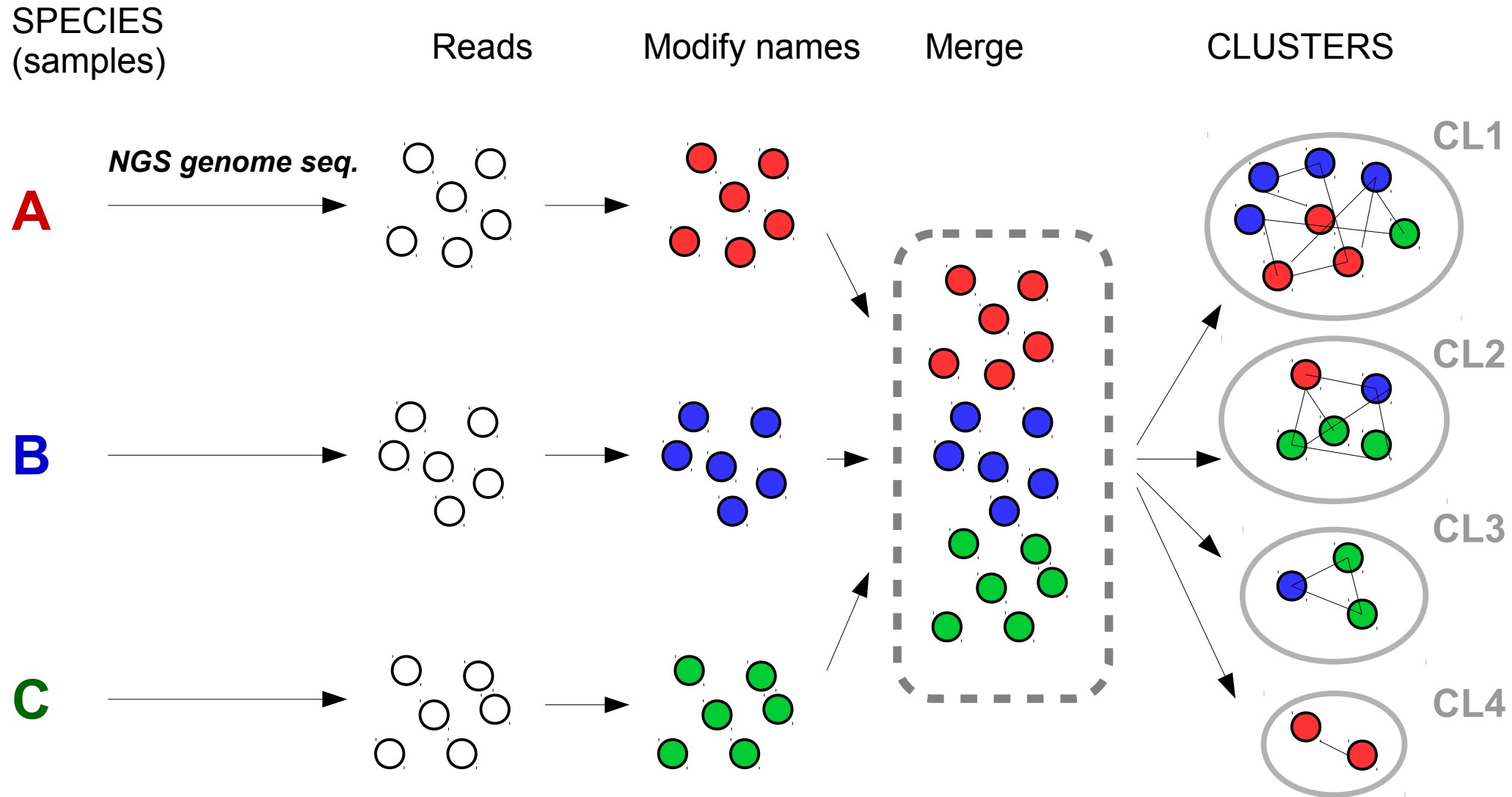
## A new tool for detection of LTR / PBS sites



LTR-retrotransposon

**PBS** (primer binding site)
*similarity to tRNA*

## Program output:

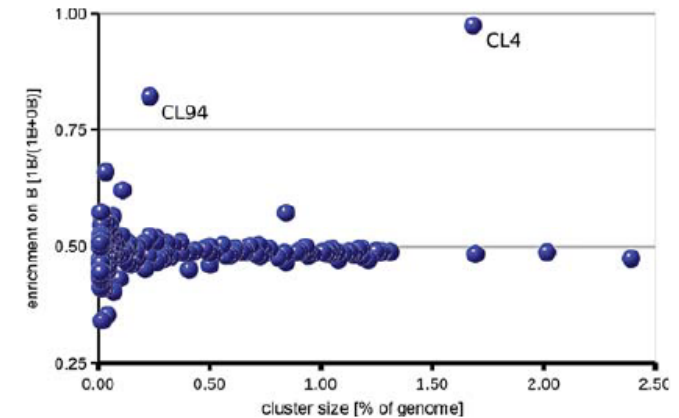| CL | contig | pos. | site | site_depth | out_ma▸ | maske▸ | maske▸ | region_in | region_out | blast to tRNA | % | length | site from | to | tRNA from | to | E-val |
|----|--------|------|------|-----------|---------|--------|--------|-----------|------------|---------------|---|--------|-----------|-----|-----------|-----|-------|
| 19 | 400 | 364 | TGCGA**CA** | 106.6 | 30.4 | 0.0362 | 0.2975 | GCGAGGAA▸ | GATGGCGA▸ | At-chr2.tRNA28-Arg▸ | 100 | 18 | 0 | 0 | 3 | 20 | 23 | 6 | 7E-007 |

(window size 7)

tRNA-Arg

# Comparative analysis - principle

# Comparative analysis

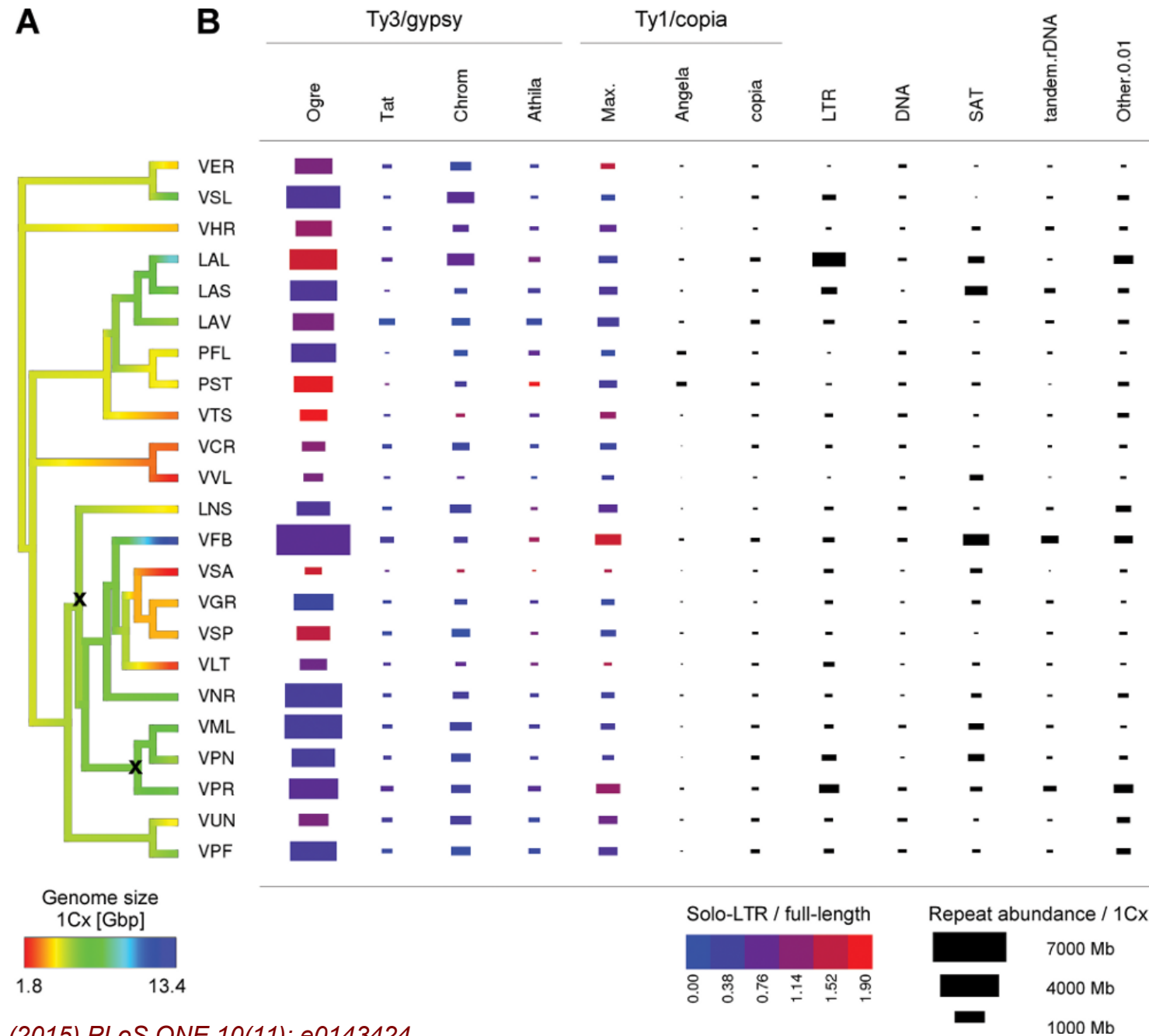*Two samples only* (e.g. genotypes of the same species +/- B chromosomes)

- simultaneous clustering of +B and 0B reads, same genome coverage

- evaluate ratios of +B/0B reads in each cluster



*Multiple samples* (e.g. a set of species differing in genome size)

- comparative clustering      – equal read numbers of genome coverages ?
            - problems with species with big variations in genome sizes
            - problems when analyzing large numbers of samples

- two-step approach

     1./ perform repeat analysis in each species separately

     2./ comparative clustering with reads sampled from (1) – finding "orthologous" repeats

# Comparative study of repeats in 23 species of *Fabeae*

# Comparative study of repeats in 23 species of *Fabeae*



*Macas et al. (2015) PLoS ONE 10(11): e0143424*