

Supplement to "Generating high quality assemblies for genomic
analysis of transposable elements"

March 27, 2020

1 Supplementary figures

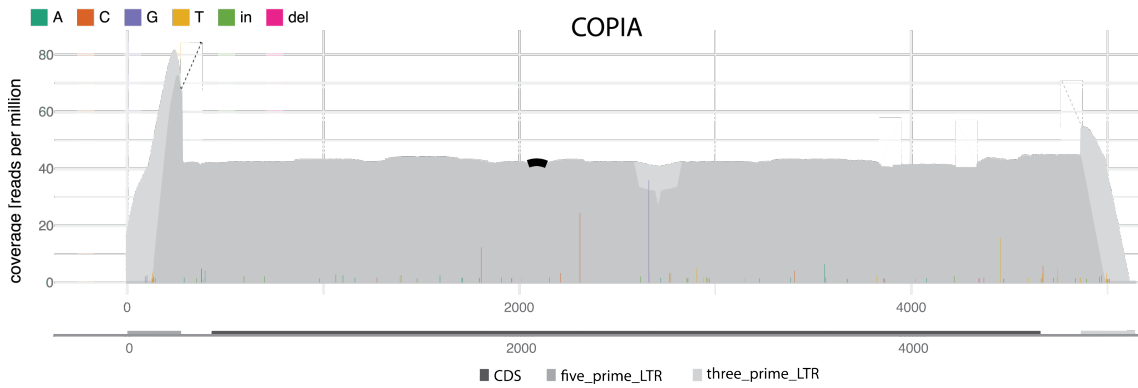


Figure 1: Abundance and diversity for *copia* elements in the *D. melanogaster* strain Canton-S. The coverage (TE abundance in rpm), the position of SNPs (colored lines) and the position of indels (bold arc at the top) are shown. The coverage based on unambiguously (dark grey) and ambiguously (light grey) mapped reads are shown. The plot was generated by DeviateTE (Weilguny and Kofler, 2019) based on Illumina paired end reads mapped to the consensus sequence of *copia* (30 coverage; 2x125bp) and

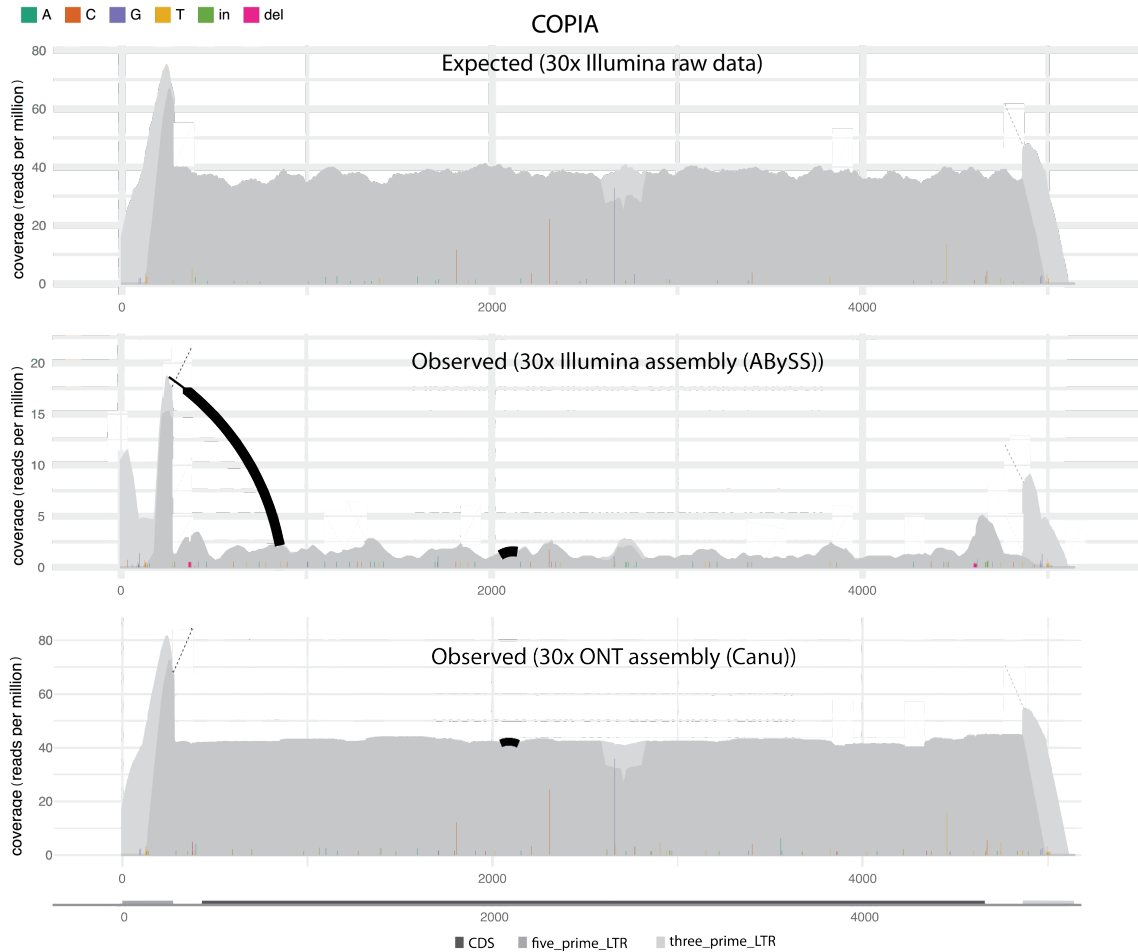


Figure 2: Expected and observed abundance and diversity of *copia* elements in Canton-S. Expected values are based on Illumina raw reads aligned to the consensus sequence of *copia*. Observed values are shown for assemblies based on short (ABySS) and long (Canu) reads. The normalised coverage is shown for ambiguously (light grey) and unambiguously (dark grey) mapped reads. The positions of SNPs (colored lines) and the position of indels (bold arcs) are shown. Note that both, the expected coverage (TE abundance) and diversity (SNPs and indels) of *copia*, are best reproduced by the long-read based assembly (Canu).

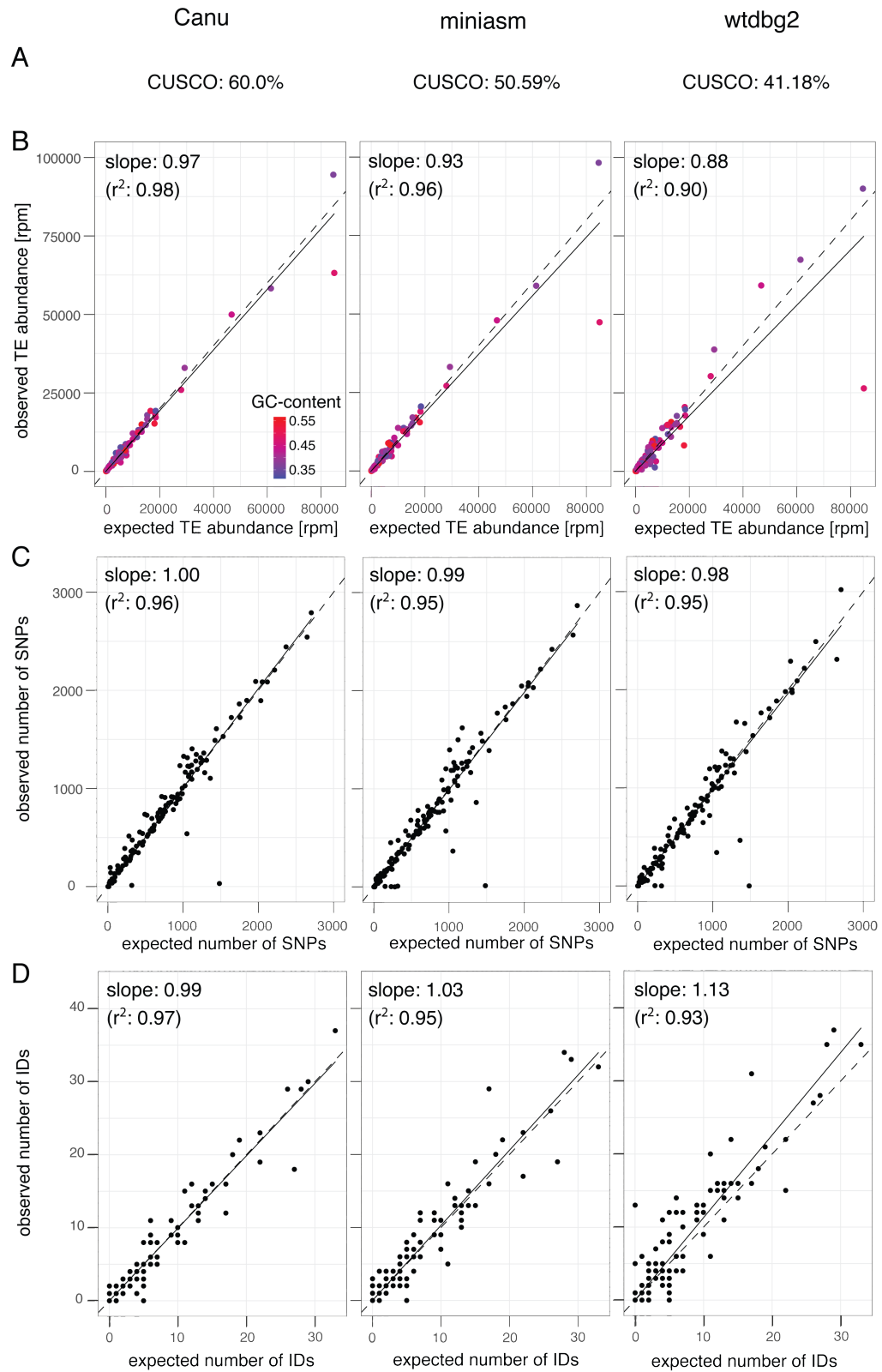


Figure 3: Influence of the assembly algorithm on the quality of assemblies of the *D. melanogaster* strain Canton-S. Assemblies are based on 30x coverage with ONT reads.

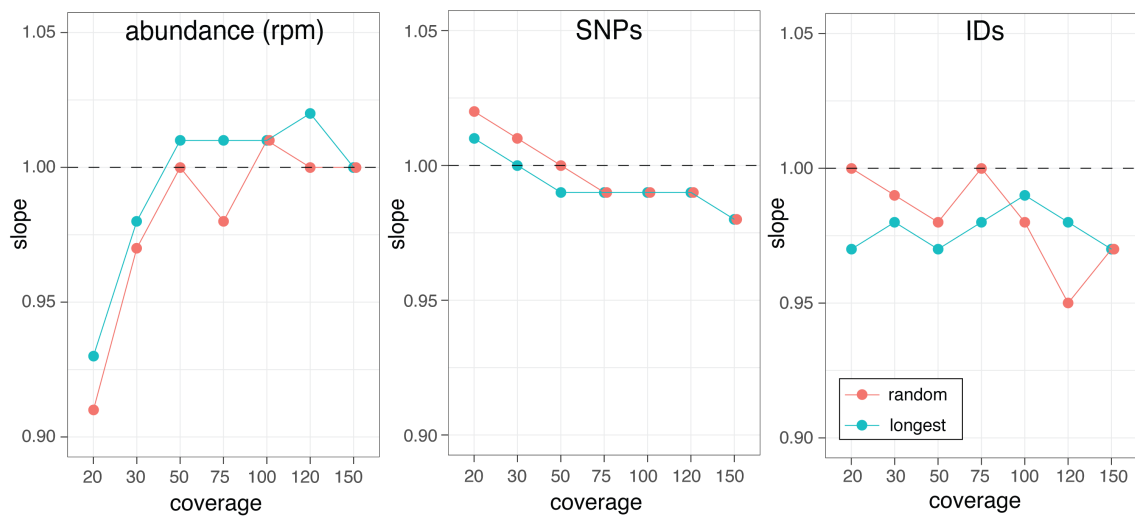


Figure 4: Assembly quality with different subsets of reads. Either random reads (random) or the longest reads (longest) were used to generate assemblies with Canu. The assembly quality is assessed using three of our TE-centered quality metrics (abundance, SNPs, IDs). The dashed lines indicate the optimal representation of TEs.

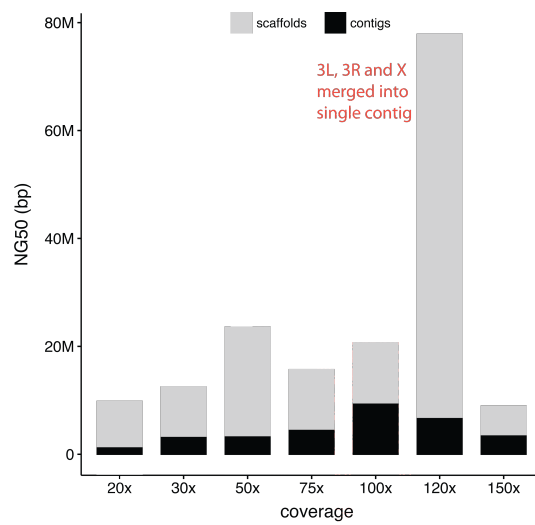


Figure 5: NG50 values of Canton-S assemblies generated with Canu and different subsamples of the longest reads. Values are shown before (contigs) and after Hi-C based scaffolding.

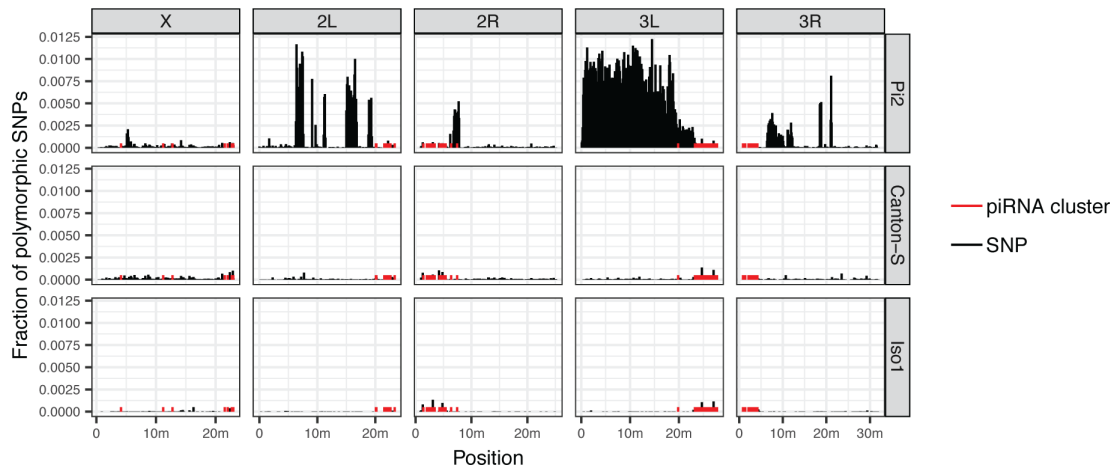


Figure 6: Location of piRNA clusters (red) and of regions with segregating polymorphisms (black) for several *D. melanogaster* strains. Segregating polymorphisms are shown for 100kb windows.

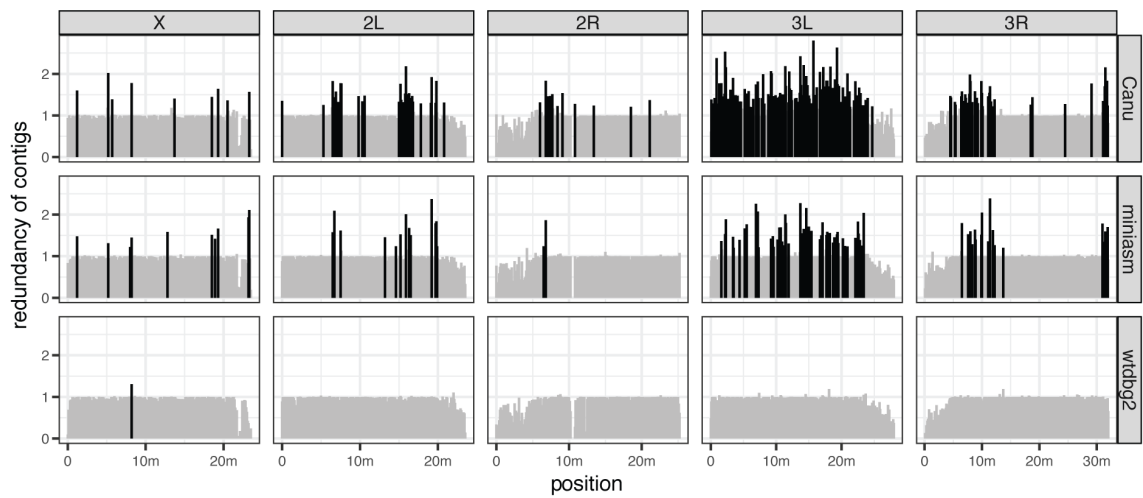


Figure 7: Origin of redundant contigs for assemblies generated by three different algorithm (right panel). Non-overlapping 1kb subsequences of an assembly were aligned to the reference. The average coverage per 100kb window is shown. Coverages > 1.2 indicate redundant contigs (shown in black).

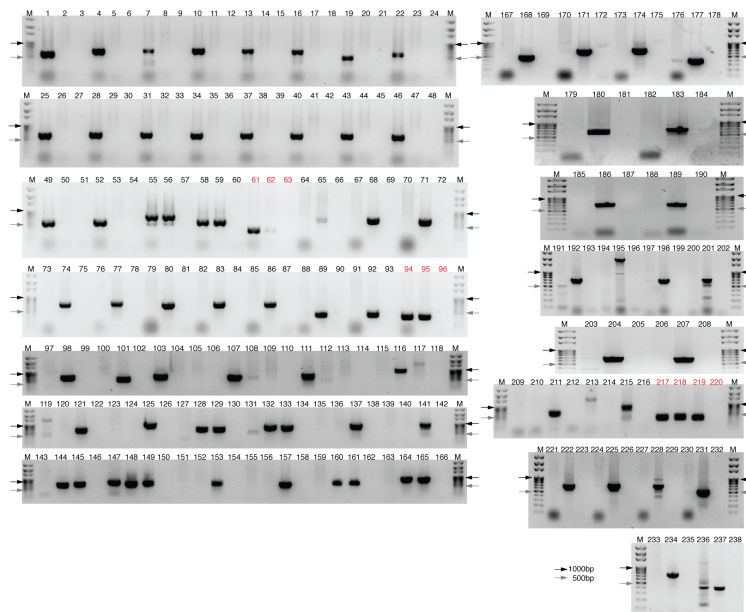


Figure 8: PCR validation of polymorphic TE insertions in piRNA clusters. Numbers above lanes refer to entries in supplementary tables 6 ; Arrows indicate the position of the 1000bp and 500bp size markers. Positive controls (*RpL32*) are labeled in red.

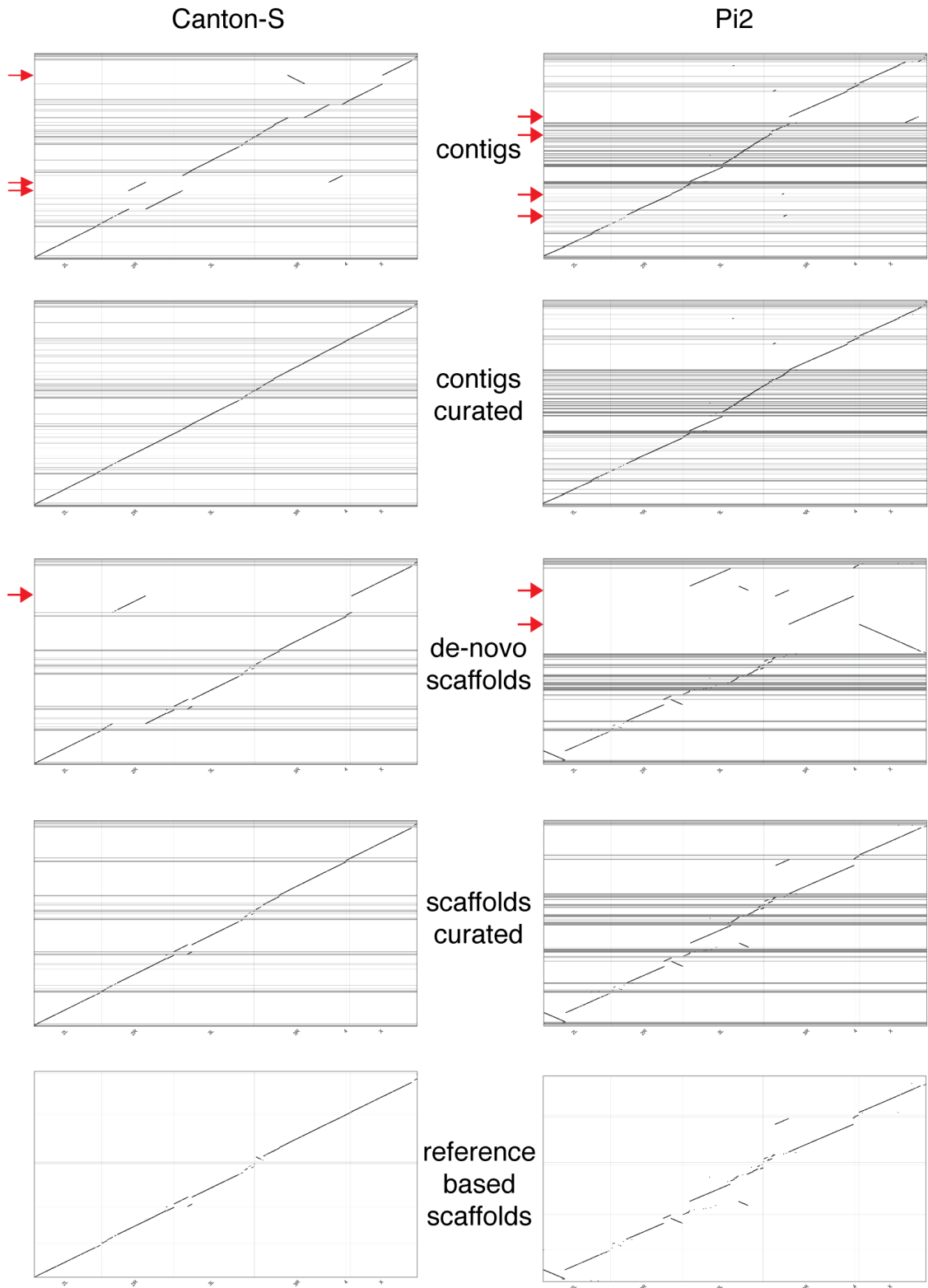


Figure 9: Manual curation steps of the final assemblies of Pi2 and Canton-S. Misassemblies (red arrows) were manually broken up at each step.

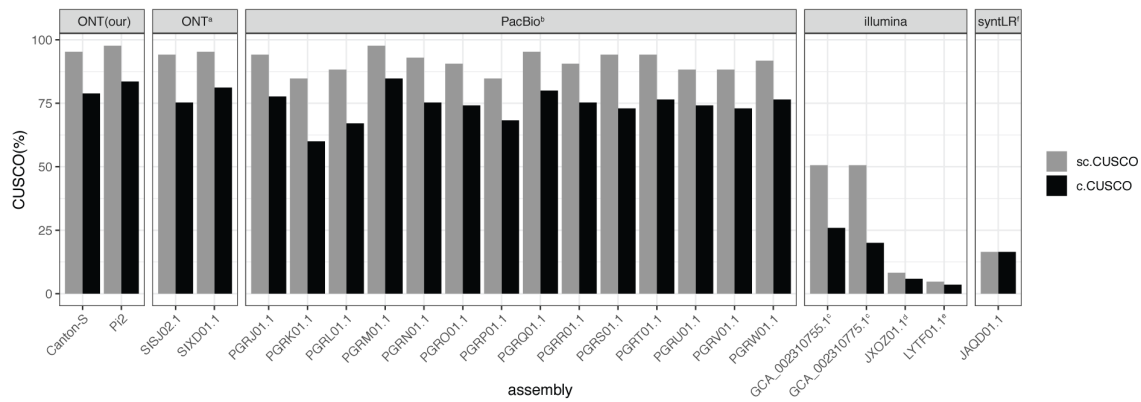


Figure 10: CUSCO values for our assemblies and publicly available assemblies of different *D. melanogaster* strains. We used assemblies from NCBI databases with following accession numbers: ^aWGS: SIXD01000000 and SISJ02000000 (Ellison and Cao, 2020) for ONT; ^bBioproject: PRJNA418342 (Chakraborty et al., 2019) for PacBio; ^cGenbank: GCA_002310755.1 and GCA_002310775.1 (Anreiter et al., 2017), ^dWGS: JXOZ01000000 (Vicoso and Bachtrog, 2015), ^eWGS: LYTF01000000 (Singhal et al., 2017) for illumina; ^fWGS: JAQD01000000 (McCoy et al., 2014) for illumina synthetic long reads.

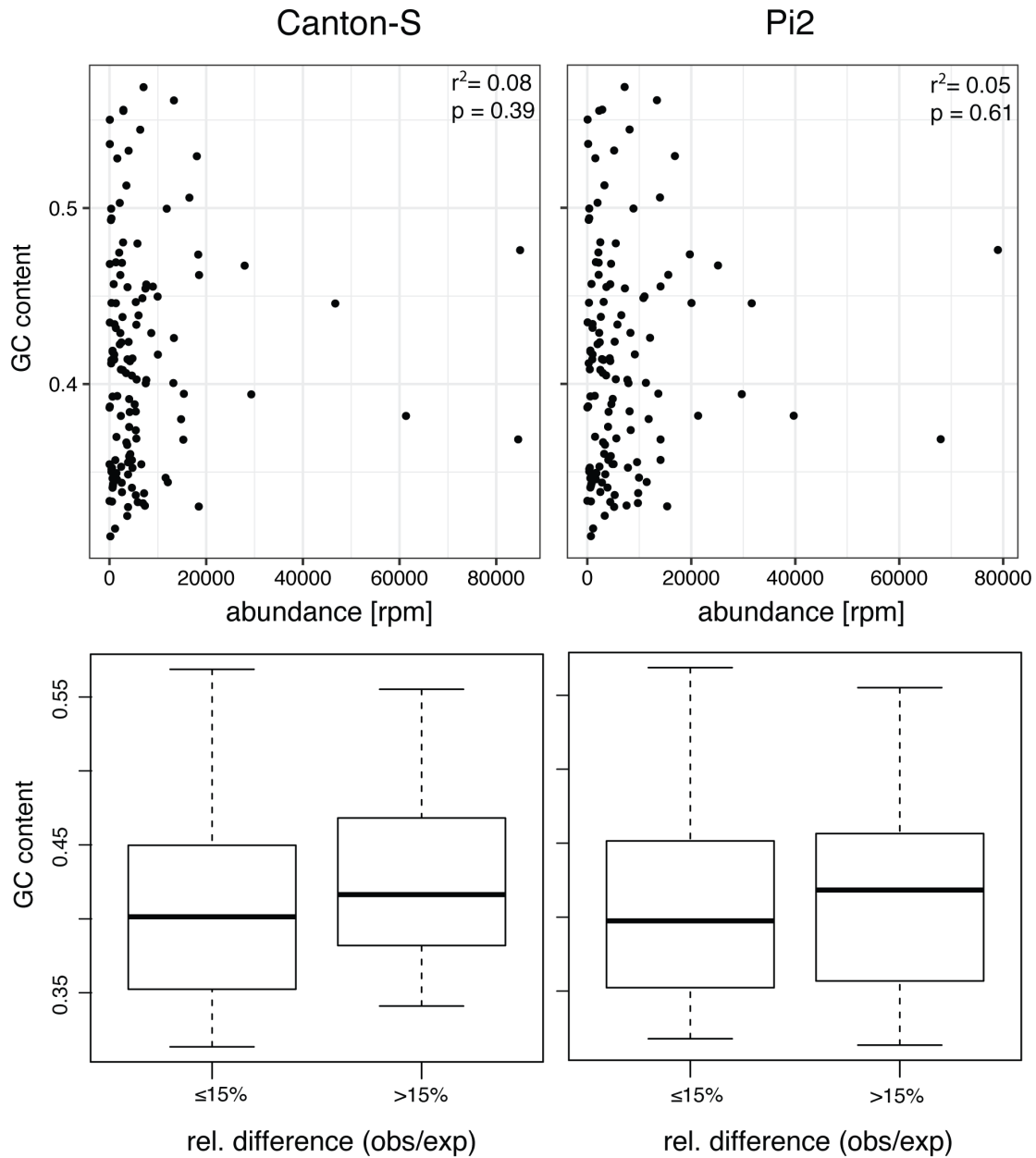


Figure 11: Influence of the GC-content on TE abundance for our assemblies of Canton-S and Pi2. A) Relationship between the the GC-content of a TE and the estimated abundance of a TE. Correlations were calculated with the spearman method. B) GC-content of TEs that strongly deviate from the expected abundance (difference between observed and expected TE abundance $> 15\%$) compared to TEs that do not deviate from expectations (difference $< 15\%$). For both Canton-S and Pi2 the differences between deviating and not-deviating TEs was not significant (Wilcoxon rank-sum test, two-sided: $p_{CS} = 0.109$, $p_{Pi2} = 0.3813$)

2 Supplementary tables

Table 1: Overview of the raw data used for the assemblies; PE paired ends

| | CantonS | Pi2 |
|--------------------------|---------|--------|
| ONT, coverage | 149x | 199x |
| ONT, flow cells | 2 | 3 |
| ONT, mean read length | 7146bp | 8045bp |
| Illumina PE, coverage | 30x | 40x |
| Illumina PE, read length | 125 | 125 |
| Hi-C, coverage | 591x | 260x |

Table 2: Influence of different polishing steps on the quality of Canton-S assemblies (30x coverage with long reads). Assembly quality is estimated with BUSCO and our four TE centered quality metrics (CUSCO, TE abundance, SNPs and IDs in TEs).

| quality | raw | racon | pilon |
|---------|------|-------|-------|
| BUSCO | 76.8 | 85.6 | 98.4 |
| Cusco | 60.0 | 60.0 | 60.0 |
| TE abu. | 0.97 | 0.97 | 0.97 |
| TE SNPs | 0.93 | 0.99 | 1.00 |
| TE IDs | 0.91 | 0.97 | 0.99 |

Table 3: Effect of polishing on BUSCO values. We applied three rounds of polishing with Racon, picked the assembly with the highest BUSCO value (bold) and polished this assembly three times with Pilon, where we again kept the assembly with the highest BUSCO values (bold). In case BUSCO values did not improve between two successive iterations we kept the assembly requiring the fewest polishing steps. The finally used polishing strategy (polishing s.) for each assembly is shown at the bottom (R .. Racon, P.. Pilon). All ONT reads were used for these assemblies.

| | CantonS | | | Pi2 | | |
|--------------|-------------|-------------|-------------|-------------|-------------|-------------|
| | Canu | miniasm | wtdbg2 | Canu | miniasm | wtdbg2 |
| unpolished | 83.0 | 1.1 | 74.9 | 87 | 0.8 | 76.1 |
| 1x Racon | 91.0 | 80.2 | 90.5 | 93.0 | 84.2 | 92.3 |
| 2x Racon | 91.0 | 90.2 | 91.4 | 92.9 | 90.3 | 92.9 |
| 3x Racon | 91.6 | 91.5 | 91.7 | 93.3 | 93.3 | 92.6 |
| 1x Pilon | 98.6 | 98.2 | 98.5 | 98.2 | 98.3 | 96.9 |
| 2x Pilon | 98.8 | 98.6 | 98.9 | 98.4 | 98.5 | 97.7 |
| 3x Pilon | 98.9 | 98.8 | 98.8 | 98.4 | 98.5 | 97.4 |
| polishing s. | 3R,3P | 3R,3P | 3R,2P | 3R,3P | 3R,2P | 3R,2P |

Table 4: BUSCO values for assemblies generated with different assemblers and coverages for Canton-S. For each assembly the optimized number of polishing rounds performed with Racon (R) and Pilon (P) are shown (for optimization procedure see supplementary table 3).

| | coverage | 20x | 30x | 50x | 75x | 100x | 120x | 150x |
|---------|-----------|-------|-------|-------|-------|-------|-------|-------|
| Canu | polishing | 3R,3P | 3R,3P | 3R,3P | 2R,3P | 3R,2P | 3R,3P | 3R,3P |
| | BUSCO | 98.4 | 98.4 | 98.6 | 98.6 | 98.6 | 98.8 | 98.9 |
| miniasm | polishing | 3R,3P | 3R,2P | 3R,3P | 3R,2P | 3R,3P | 3R,3P | 3R,3P |
| | BUSCO | 98.4 | 98.4 | 98.6 | 98.6 | 98.6 | 98.8 | 98.8 |
| wtdbg2 | polishing | 1R,3P | 2R,3P | 2R,3P | 3R,3P | 3R,2P | 2R,3P | 3R,2P |
| | BUSCO | 98.8 | 98.6 | 98.7 | 98.7 | 98.7 | 98.6 | 98.9 |

Table 5: Overview of the final assemblies of Canton-S and Pi2. The assembly quality is assessed with classic quality metrics (NG50, BUSCO) as well as our TE centered quality metrics. Misassembled contigs and scaffolds were broken manually (based on dot-plots; supplementary fig. 1). c.CUSCO contig-CUSCO, sc.CUSCO scaffold-CUSCO

| | | CantonS | Pi2 |
|------------|-----------|-----------|-------|
| ONT | Assembler | Canu | Canu |
| | Coverage | 100x | 100x |
| | contigs | 335 | 625 |
| | NG50 | 6.8m | 4.1m |
| | length | 148m | 169m |
| | BUSCO | 82.3 | 85.8 |
| | c.Cusco | 80.00 | 77.65 |
| | TE abu. | 1.02 | 1.04 |
| | TE SNPs | 0.95 | 0.97 |
| | TE IDs | 0.93 | 0.95 |
| pol. | strategy | 2R,2P | 2R,2P |
| | contigs | 335 | 625 |
| | NG50 | 4.6m | 4.1m |
| | length | 149m | 169m |
| | BUSCO | 98.7 | 98.3 |
| | c.Cusco | 81.18 | 83.53 |
| | TE abu. | 1.01 | 1.04 |
| | TE SNPs | 0.99 | 1.01 |
| | TE IDs | 0.99 | 1.01 |
| | Hi-C | scaffolds | 266 |
| NG50 | | 21.4m | 23.6m |
| length | | 149m | 169m |
| BUSCO | | 98.7 | 98.3 |
| sc.Cusco | | 84.71 | 91.76 |
| TE abu. | | 1.01 | 1.04 |
| TE SNPs | | 0.99 | 1.01 |
| TE IDs | | 0.99 | 1.01 |
| ref. scaf. | scaffolds | 15 | 16 |
| | NG50 | 28.2m | 37.2m |
| | length | 149m | 169m |
| | BUSCO | 98.6 | 98.2 |
| | sc.Cusco | 95.29 | 97.65 |
| | TE abu. | 1.01 | 1.04 |
| | TE SNPs | 0.99 | 1.01 |
| | TE IDs | 0.98 | 1.00 |

- hibit widespread allelic heterogeneity and shape variation in complex traits. *Nature Communications*, 10(1):419275.
- Ellison, C. E. and Cao, W. (2020). Nanopore sequencing and Hi-C scaffolding provide insight into the evolutionary dynamics of transposable elements and piRNA production in wild strains of *Drosophila melanogaster*. *Nucleic Acids Research*, 48(1):1–14.
- McCoy, R. C., Taylor, R. W., Blauwkamp, T. A., Kelley, J. L., Kertesz, M., Pushkarev, D., Petrov, D. A., and Fiston-Lavier, A.-S. (2014). Illumina TruSeq synthetic long-reads empower de novo assembly and resolve complex, highly-repetitive transposable elements. *PLoS ONE*, 9(9):e106689.
- Singhal, K., Khanna, R., and Mohanty, S. (2017). Is *Drosophila*-microbe association species-specific or region specific? A study undertaken involving six Indian *Drosophila* species. *World Journal of Microbiology and Biotechnology*, 33(6):103.
- Vicoso, B. and Bachtrog, D. (2015). Numerous Transitions of Sex Chromosomes in Diptera. *PLOS Biology*, 13(4):1–22.
- Weilguny, L. and Kofler, R. (2019). DeviaTE: Assembly-free analysis and visualization of mobile genetic element composition. *Molecular Ecology Resources*, 19(5):1346–1354.