

Potential and pitfalls of eukaryotic metagenome skimming:

Motivation

Metagenomic sequencing with only a single library layout is used to quickly and cheaply assess the taxonomic and functional complexity of large and diverse microbial communities. We investigate to what extent such metagenome skimming approaches are applicable for the in-depth characterizations of genomes found in obligate symbiotic communities of eukaryotes, e.g. lichens. It is still unclear how a eukaryotic species mixture, with larger and more repeat-rich genomes, influences different *de novo* assembly paradigms, such as *de Bruijn Graph* based methods or *Overlap Layout* based assemblers and how to optimize assembly parameters as k-mer or overlap sizes.

1. *in silico* Sequencing

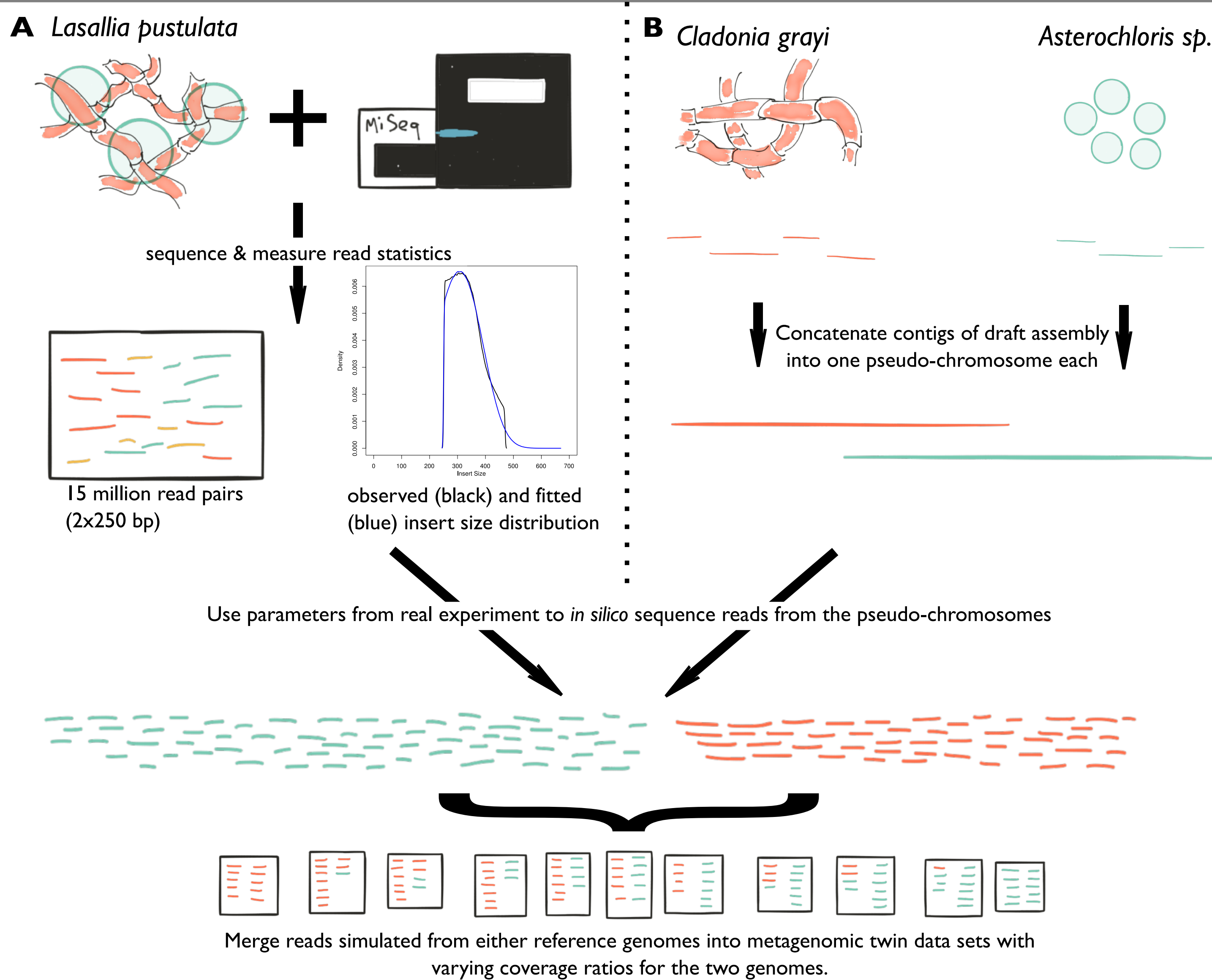


Figure 1: Workflow for generating twin data sets, resembling a real sequencing data set with respect to insert size distribution, read number and read length.

DNA from a thallus of *Lasallia pustulata* was sequenced using Illumina MiSeq technology, yielding 15 million read pairs with a length of 250 bp. To estimate the insert size distribution, we overlapped read pairs using FLASH [1] and fitted a censored Weibull distribution to the observed insert size distribution (Figure 1, A).

The scaffolds of the genomes of *Cladonia grayi* [2] and *Asterochloris sp.* [3] were each concatenated to create a contiguous pseudo-chromosome, respectively (Figure 1, B). Both were checked for repeat content & self-similarity using Repeatmasker [4] (Box I) and Gepard [5] (Fig 2).

Using the pseudo-chromosomes as templates, we simulated reads using ART [6], parameterized using the values from Fig. 1, A. The reads were mixed into 11 twin data sets by mixing fungal and algal reads in varying ratios (Table I).

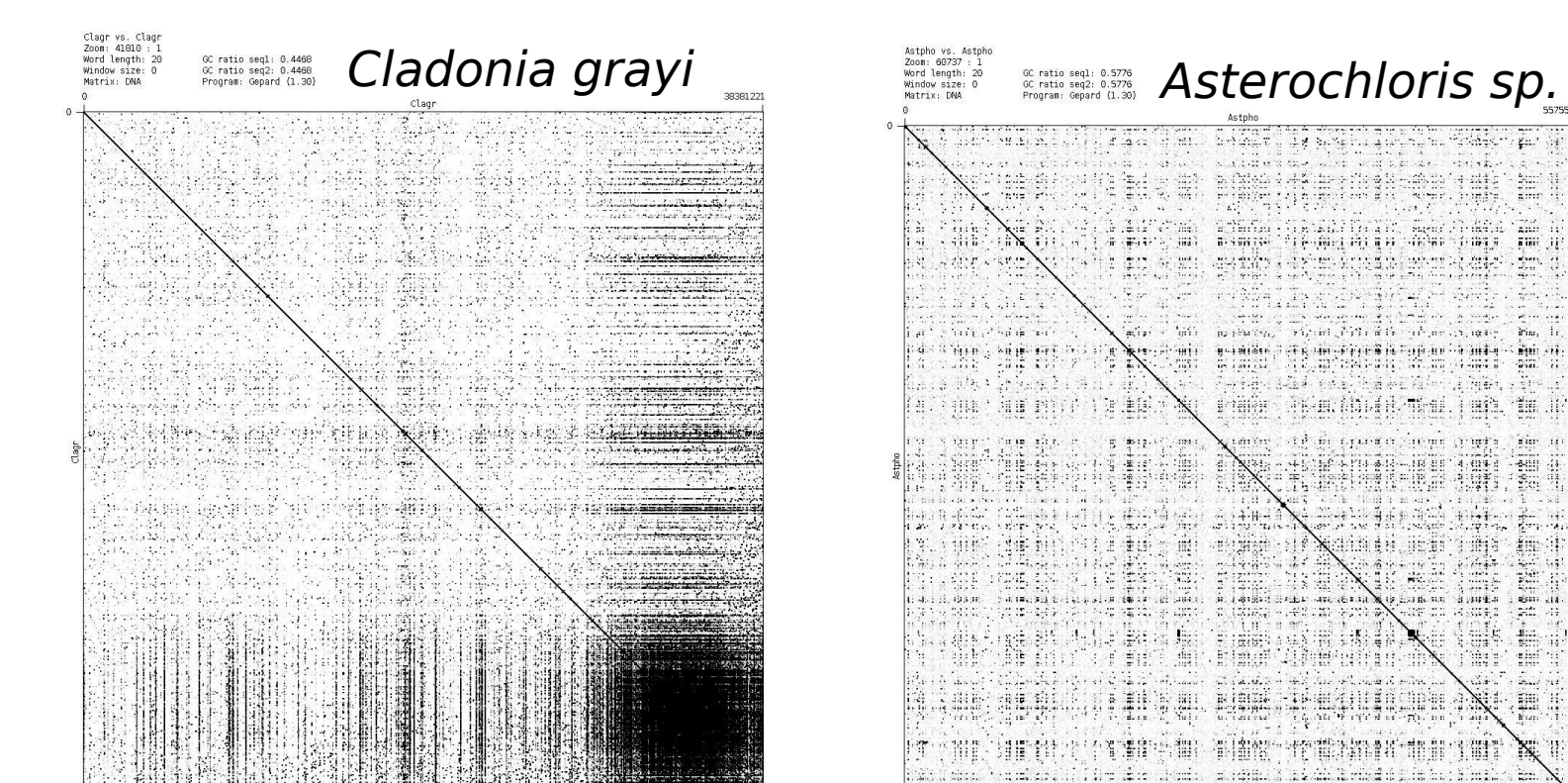
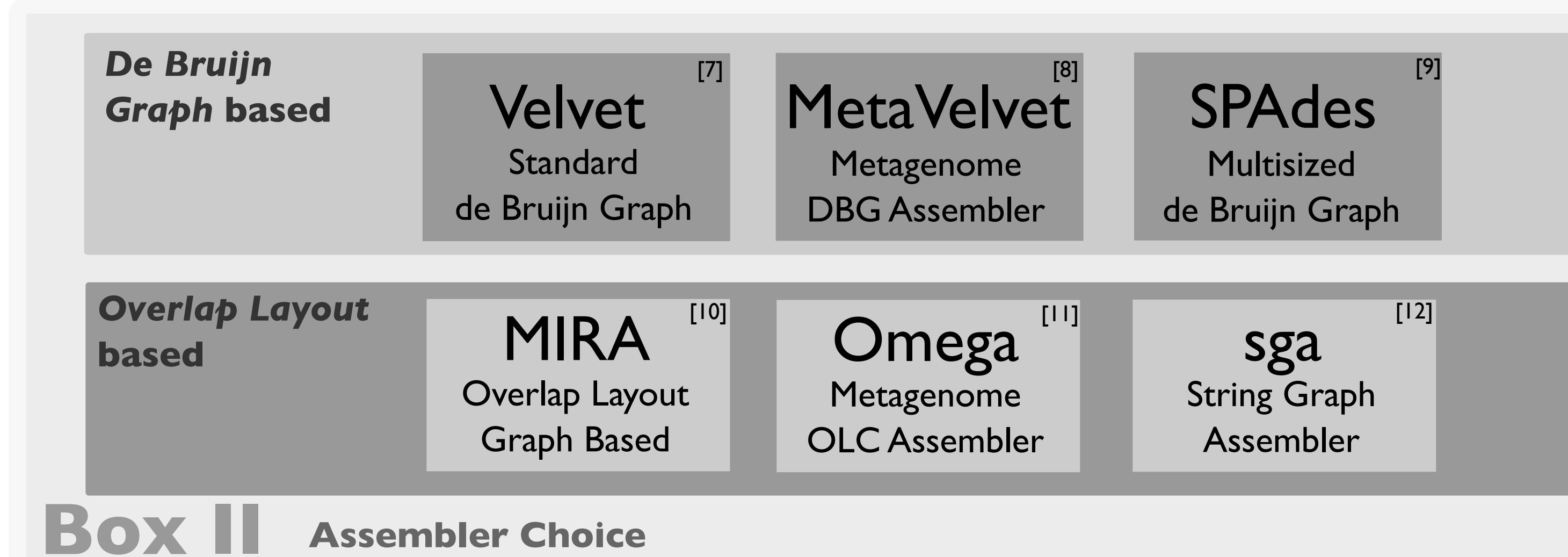


Figure 2: Dotplot of the pseudo-chromosomes of *Cladonia grayi* and *Asterochloris sp.*

2. Assembler Selection & Optimisation



For Omega, sga, Velvet & MetaVelvet we explored the parameter space (overlap size and k-mer size respectively) and use the maximization of the N50 size as the acceptance objective.

To address these questions, we performed an *in silico* study, simulating a genome skimming experiment of a lichen. We show that the quality of genome reconstructions from metagenome skimming data depends essentially on assembler choice, but also on the parameter optimisation strategy used. In the worst case optimising standard assembly metrics can lead to the excluding of complete genomes. Reconciling the expectations from the *in silico* study with the outcome of a real-world metagenome skimming of the lichen *Lasallia pustulata* indicates an even larger biodiversity, causing the underrepresentation of one symbiont in the shotgun library.

3. Assembly Results

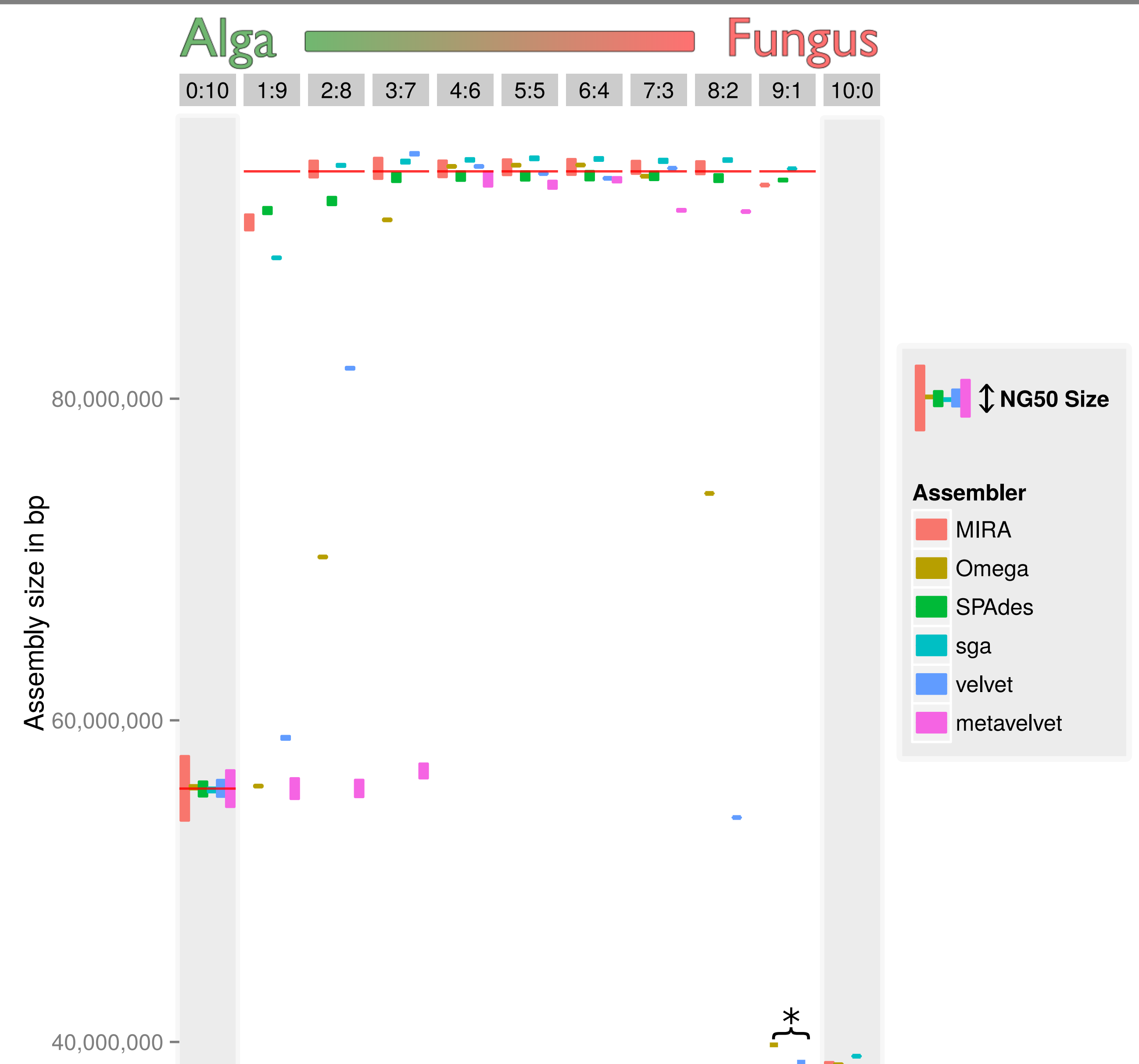


Figure 3: Assembly results for the 11 data sets and the different assemblers. Bars are centered at total assembly length, red lines are reference lengths. Height of the bars shows the NG50 size. For the assemblies with the asterisk the total assembly length was less than 50% of the reference length. A default height was used in those instances.

In case of the single species data sets almost all assemblers reconstruct the two genomes over their full length (Figure 3, column 0:10 & 10:0), however with varying NG50 sizes. For the alga many assemblers exceeded the NG50 size of the original draft assembly. For the fungus, repeats hindered such an extension with the present WGS library layout (Figure 4).

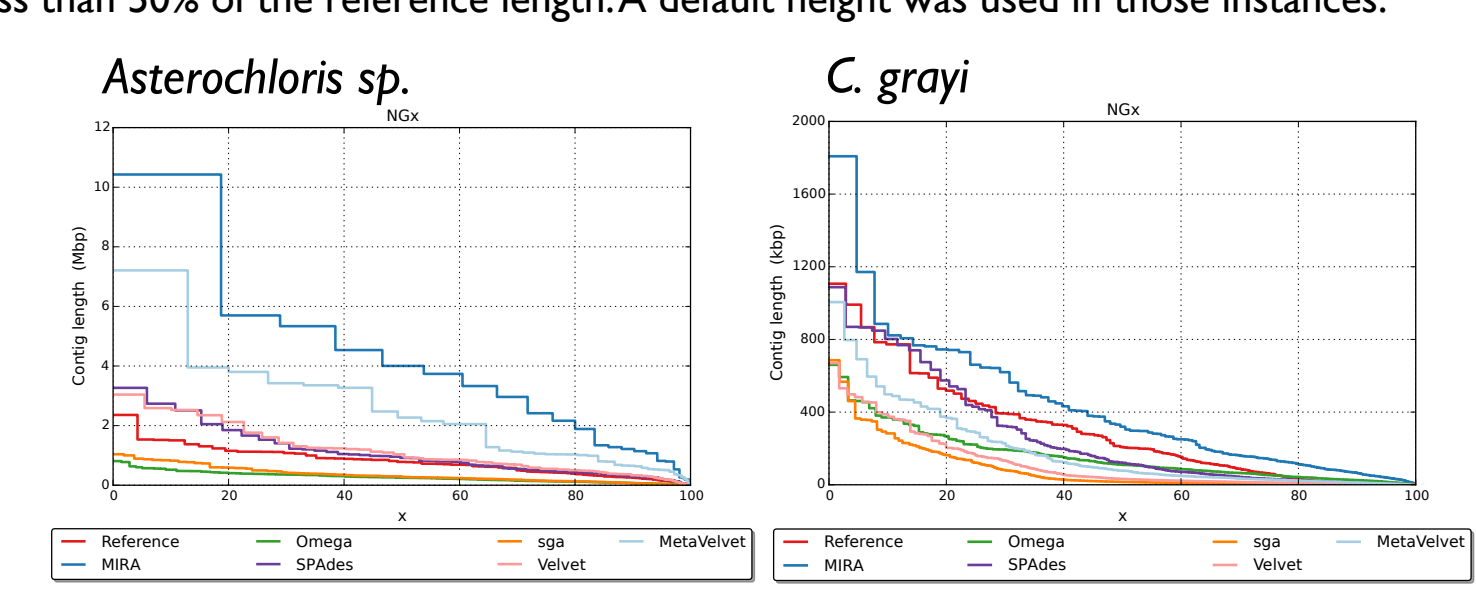


Figure 4: NGx distributions for *Asterochloris sp.* & *C. grayi*

For the the mixed species data, completeness of the genome reconstructions depends heavily on assembler choice and coverage ratios. *MIRA* and *SPAdes* perform best across all data sets. In contrast, *Omega*, *Velvet* and *MetaVelvet* fail to assemble large parts of the low coverage genome once coverage ratios become extreme (Fig. 3, 1:9 - 3:7, 9:1).

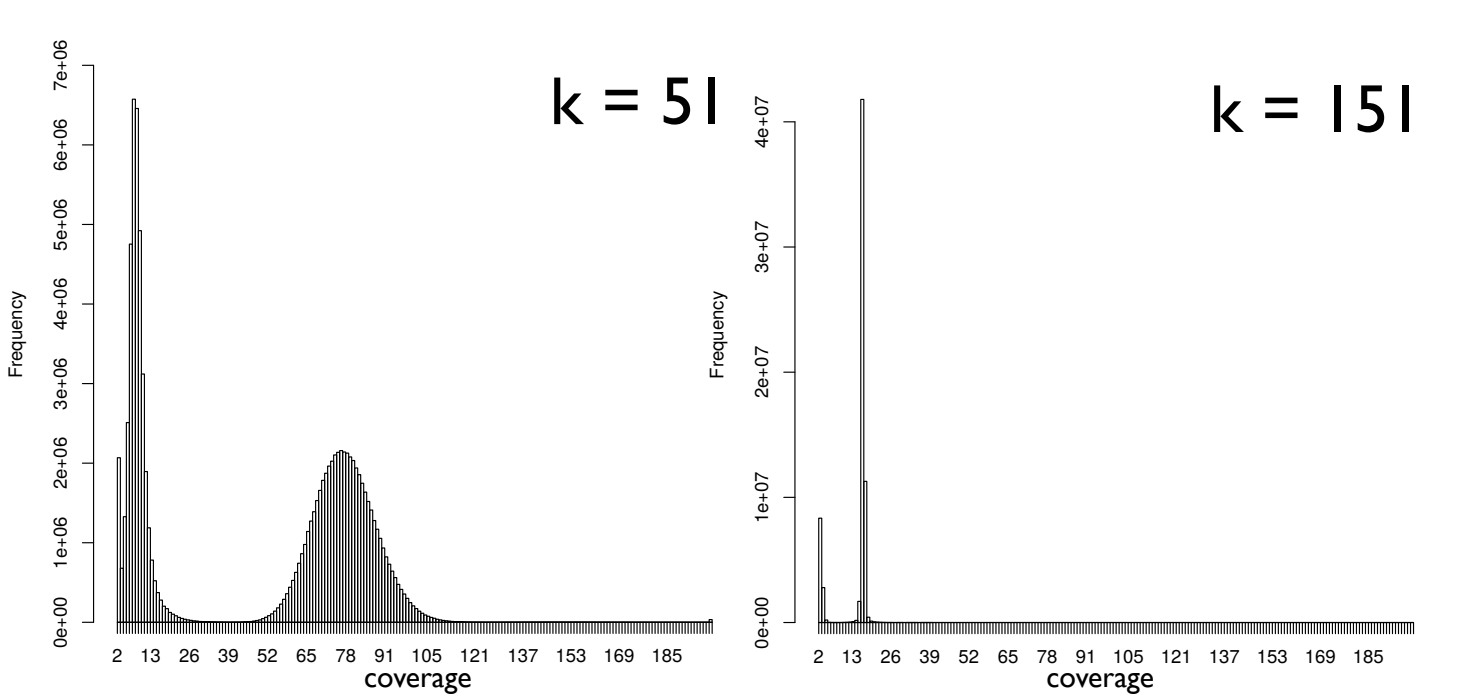
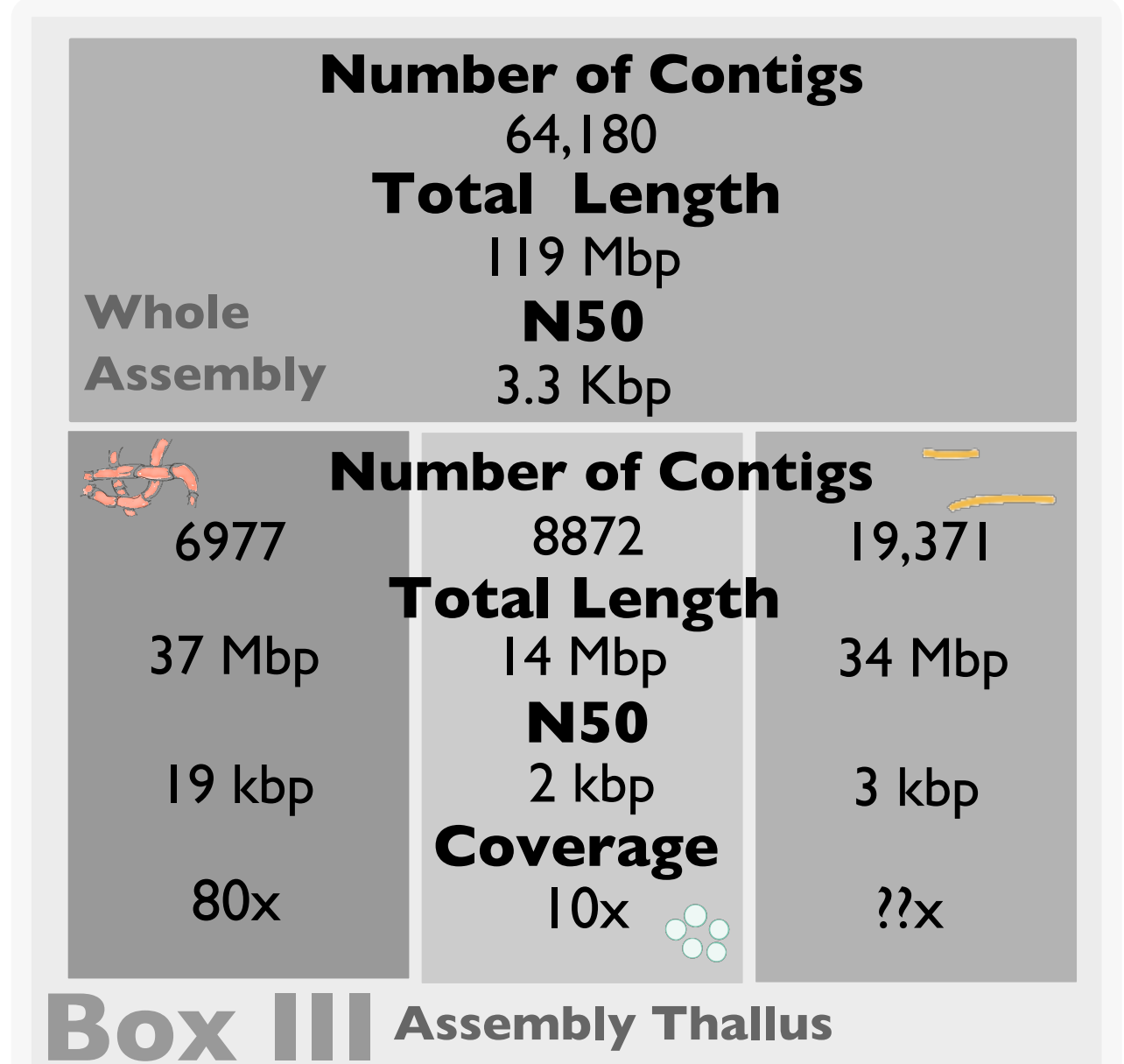


Figure 6: kmer-coverage frequencies for the 1:9 data set.

The k-mer coverage plots explain **the assemblers' sensitivity to biased coverage ratios**. Increasing the value of k reduces the frequency of all k-mers (Fig. 6), causing k-mers from the low coverage genome to overlap with those introduced by the sequencing error. This prevents the formation of typically short contigs, thus increasing the N50 size by not assembling the low coverage genome.

L. pustulata was assembled using MIRA and contigs taxonomically assigned using *MEGAN* [13]. The fungal and algal genome are much more fragmented than expected given the *in silico* study. This is a result of the larger microbial diversity: nearly 1/3 of the assembly is of bacterial origin, which decreases the absolute coverage, especially for the alga (Box III, right column).



Summary

- Twin sets are valuable for guiding strategic decisions during planning of metagenome sequencing and assembly.
- Optimising the N50 can lead to the preclusion of sequences representing the low-coverage genome.
- Assembler performance already varies substantially for single species data.
- The assembly of *L. pustulata* shows that lichen thalli are even more diverse in their composition, driving down absolute coverages for all organisms.
- Mixing data from different species inflates the assembler performance differences, with *MIRA* & *SPAdes* yielding the most contiguous sequences

