

Potential and pitfalls of eukaryotic metagenome skimming:

Motivation

Whole genome shotgun sequencing of multi species communities using only a single library layout is commonly used to assess taxonomic and functional complexities of large and diverse microbial communities. We investigate to what extent such metagenome skimming approaches are applicable for in-depth genomic characterizations of obligate symbiotic communities involving eukaryotes, e.g. lichens. It is still unclear how a eukaryotic species mixture 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

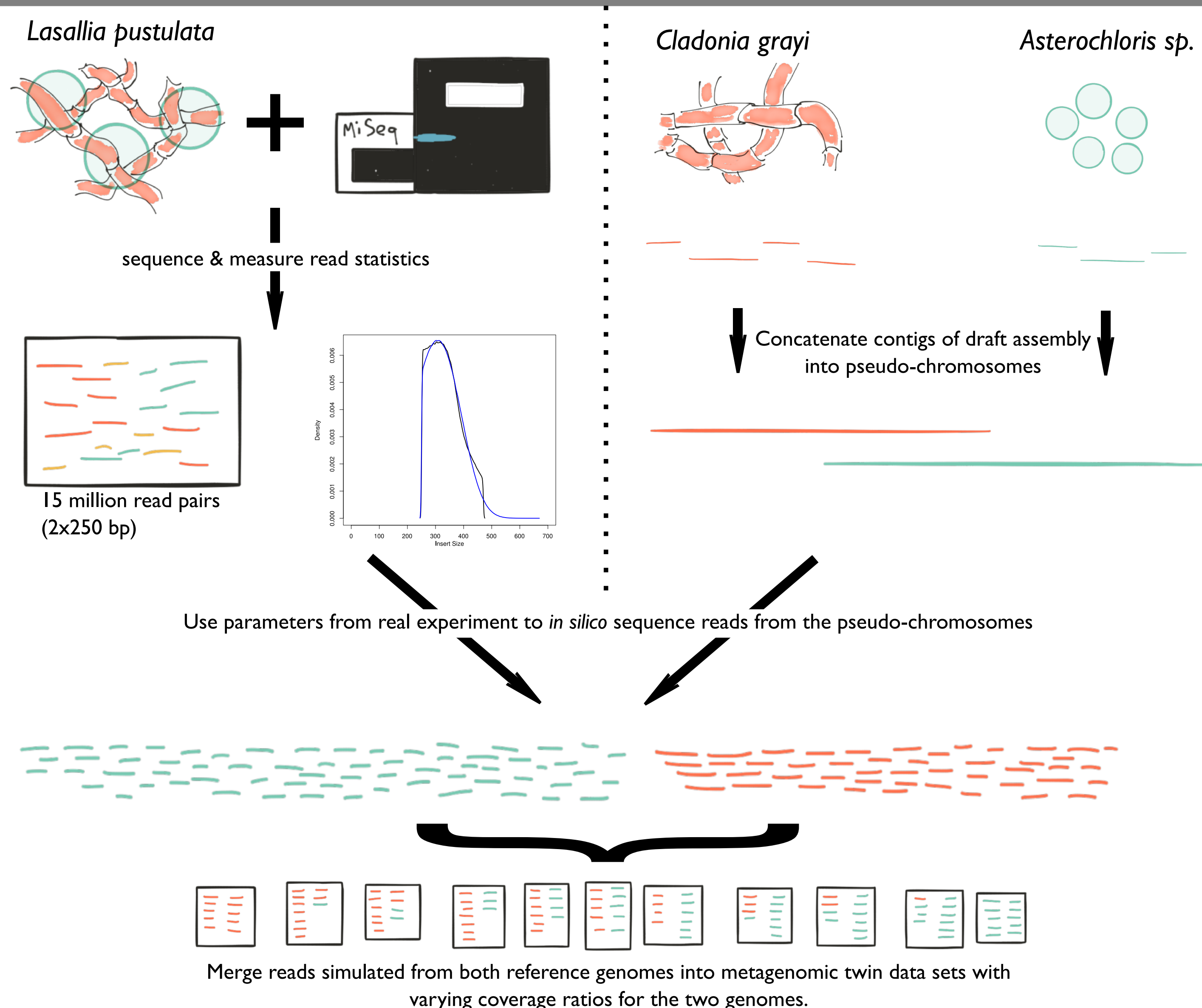
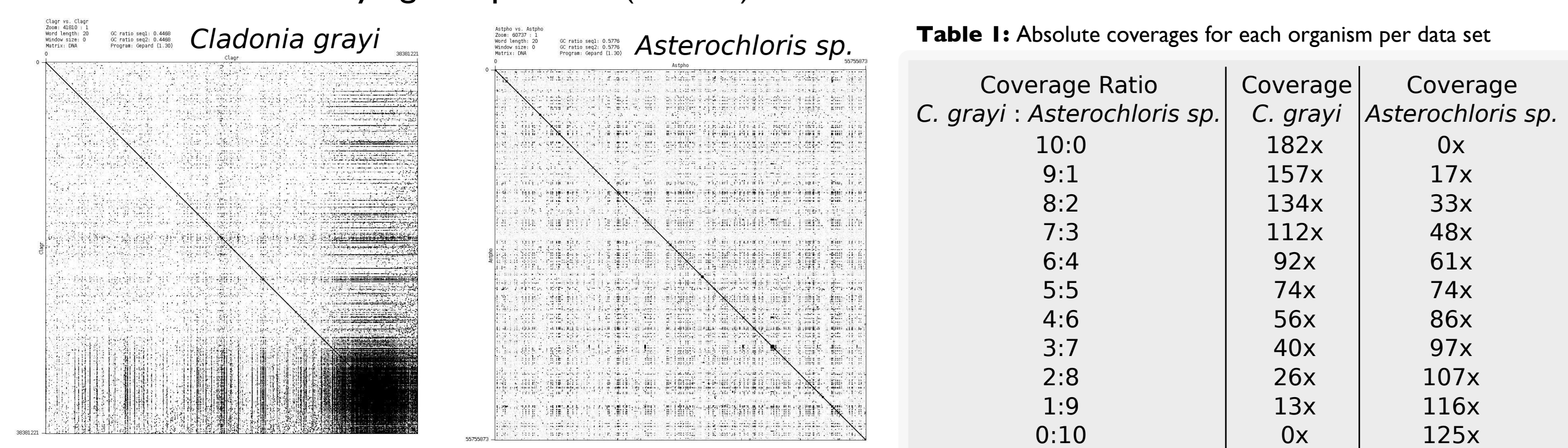


Figure 1: Workflow for generating twin data sets, which resemble a real sequencing data set in terms of insert size distribution, read number and read length.

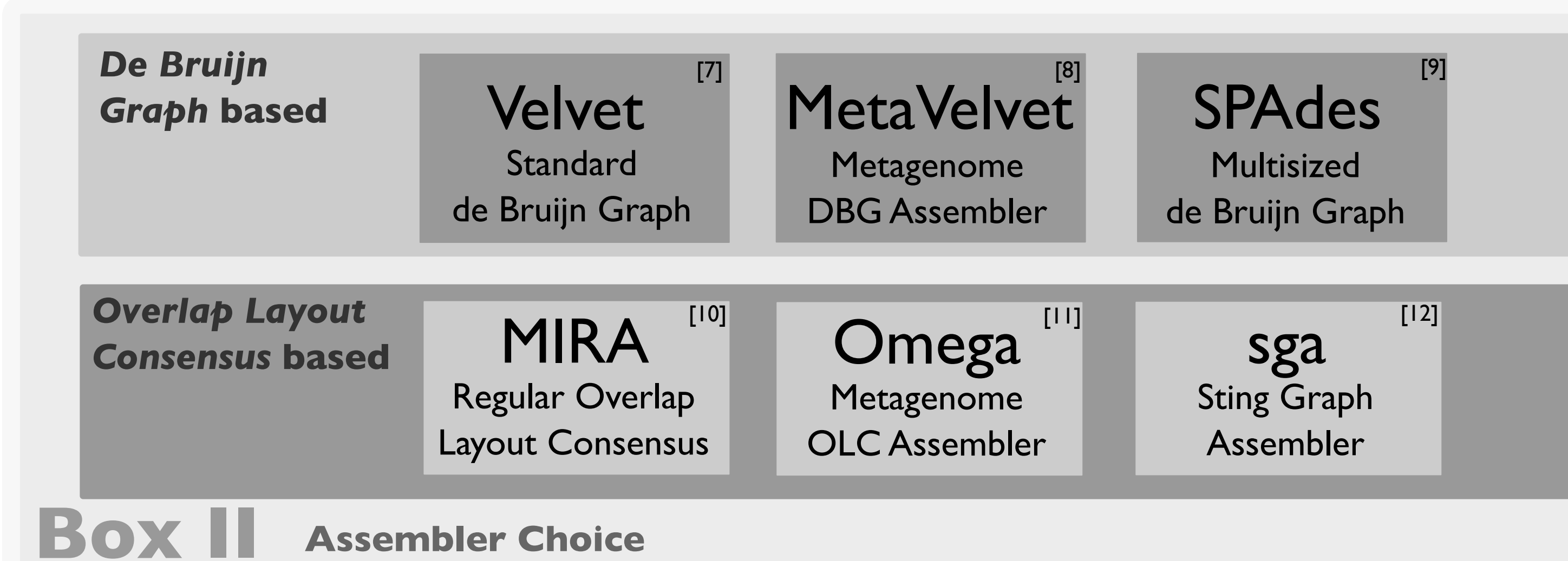
DNA from a thallus of *Lasallia pustulata* was sequenced using Illumina MiSeq technology, yielding 15 million read pairs with a length of 250 bp. To estimate the insert size distribution, we overlapped read pairs using FLASH [1] and fitted a censored Weibull distribution to the observed insert sizes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

The scaffolds of the draft genomes of *Cladonia grayi* [2] and *Asterochloris sp.* [3] were concatenated to create contiguous pseudo-chromosomes (Figure 1, right). These were analysed using Repeatmasker [4] (Box I) and Gepard [5] (Fig 2.)

From the pseudo-chromosomes reads were simulated with the parameters estimated from *L. pustulata*, using ART [6], and mixed into 11 twin data sets of varying composition (Table I).



2. Assembler Selection & Optimisation



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

To investigate how those paradigms perform on eukaryotic species mixtures, we performed an *in silico* study based on a lichen system, mimicking a genome skimming study. 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. Reconciling the expectations from the *in silico* study with the outcome of a real-world metagenome skimming of the lichen *Lasallia pustulata* indicates methodological problems causing the underrepresentation of one symbiont in the shotgun library.

3. Assembly Results

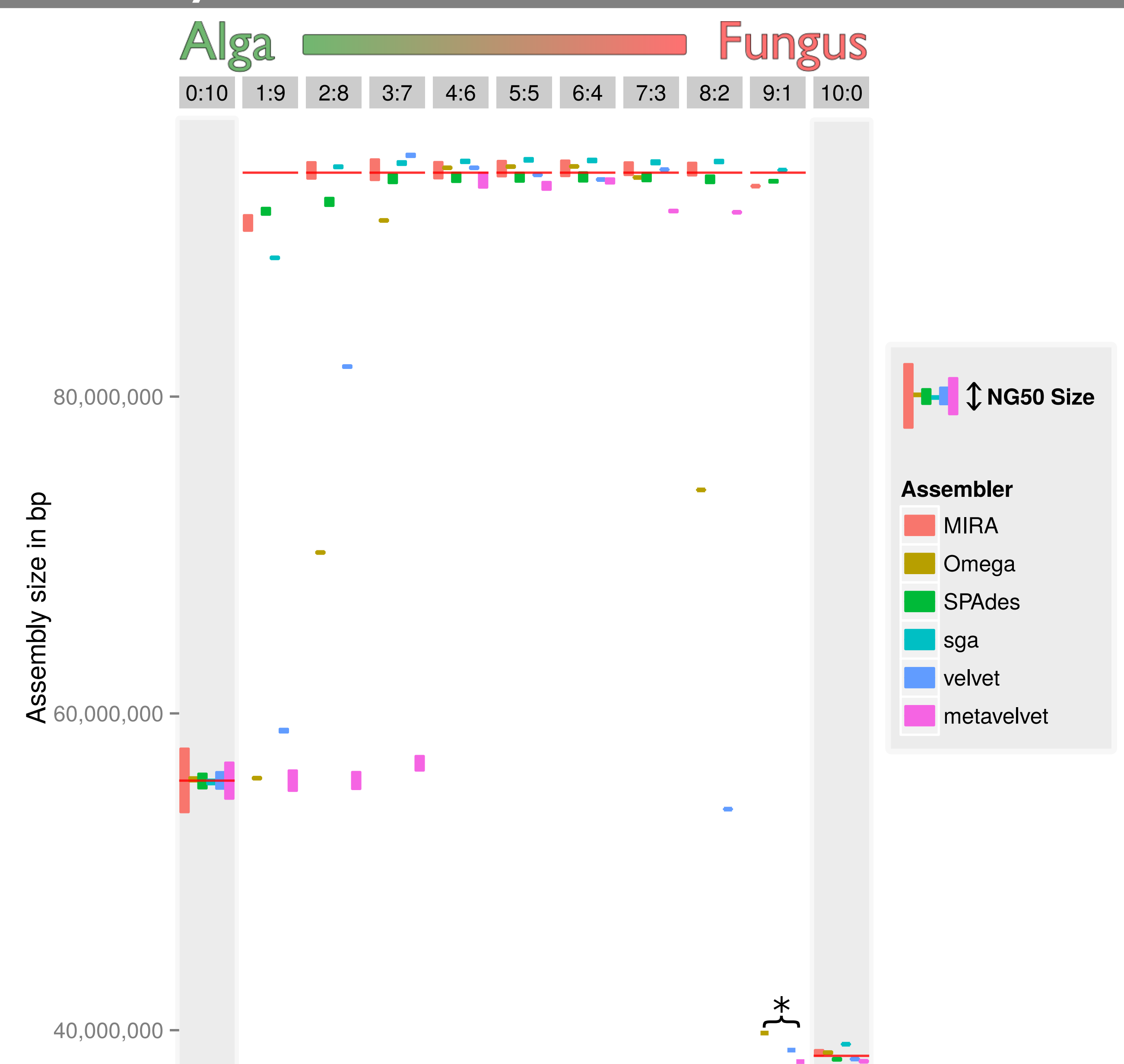
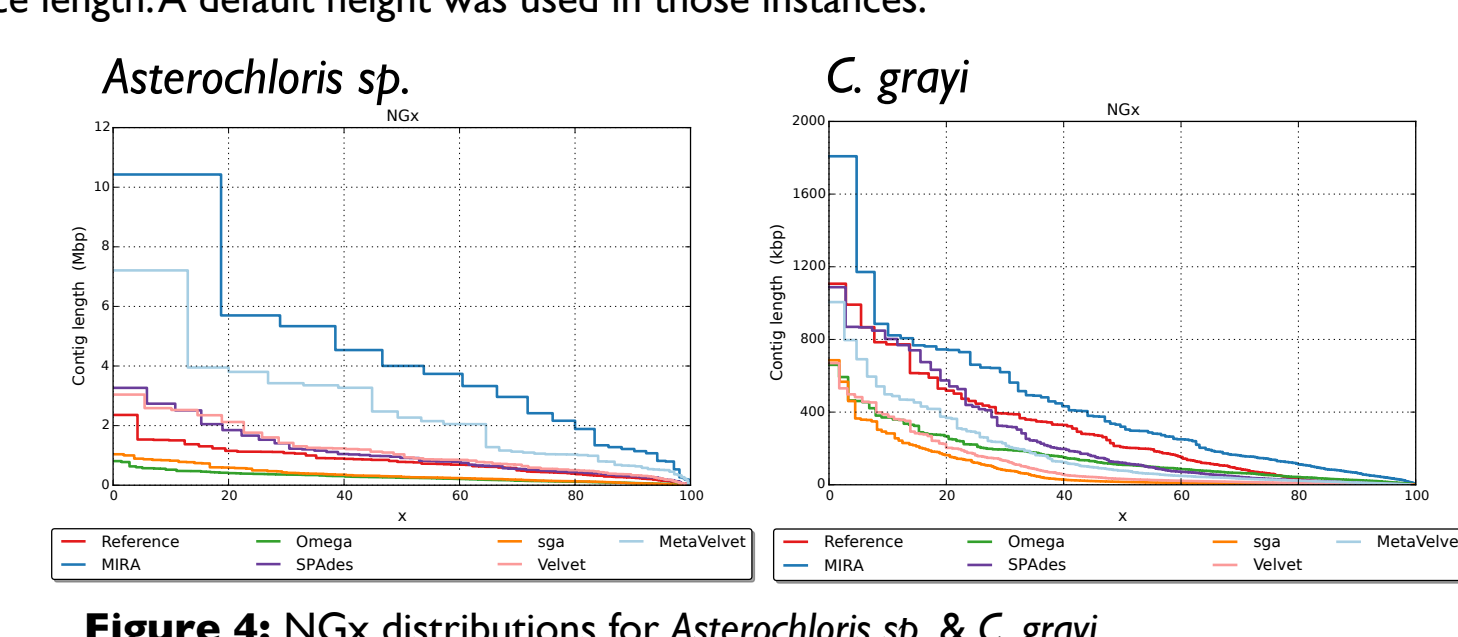


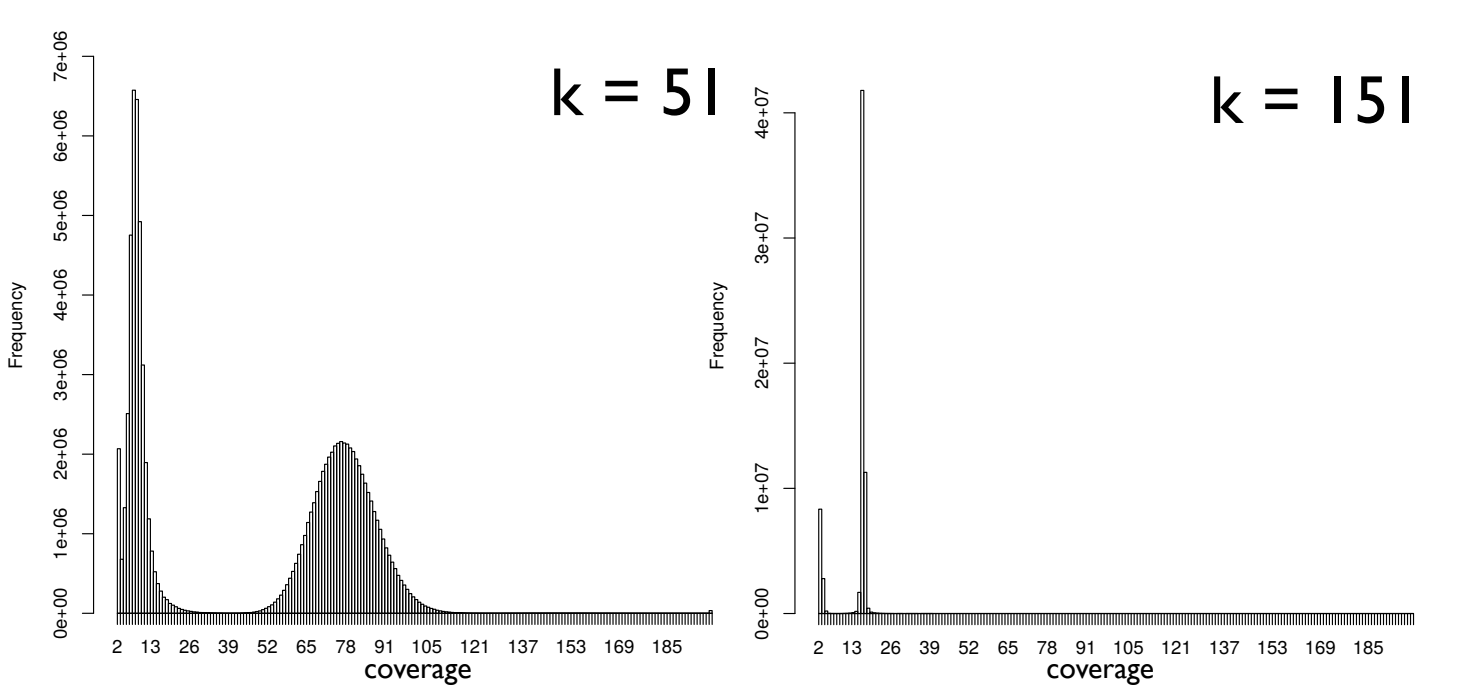
Figure 3: Assembly results over the 11 data sets. Bars are centered at total assembly length (red lines are reference lengths). Height of the bars shows the NG50 size. For the assemblies with the asterisk the total assembly length was less than 50% of the reference length. A default height was used in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

For 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 were able to exceed the NG50 size of the original draft genome. For the fungus, repeats hindered such an extension with the present WGS library layout (Figure 4).

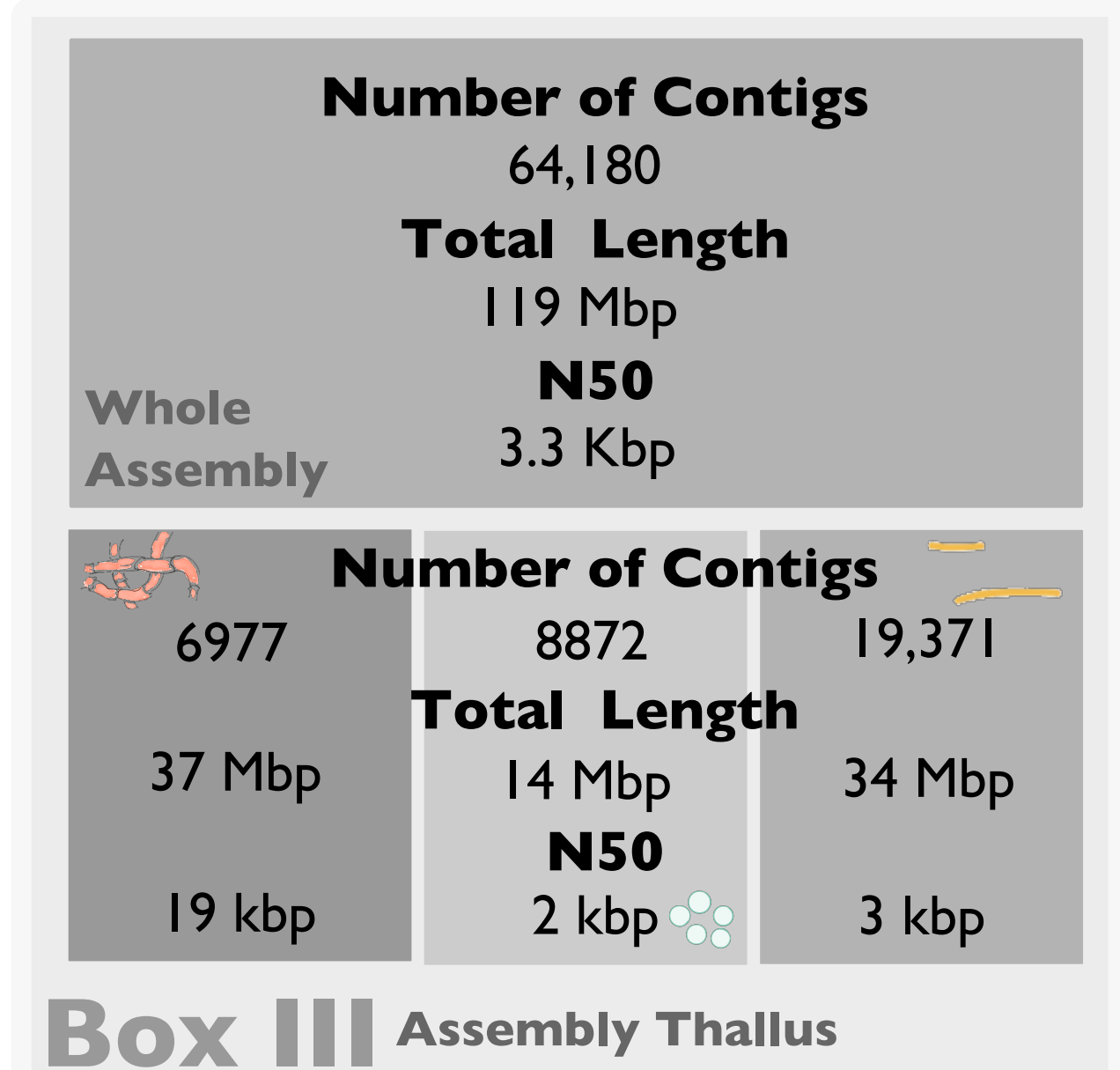


For the the mixed species data, completeness of the genome reconstructions depends heavily on assembler choice and coverage ratios. *MIRA* and *SPAdes* perform best across all data sets. In contrast, *Omega*, *Velvet* and *MetaVelvet* fail to assemble large parts of the low coverage genome once coverage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

The k-mer coverage plots provide an explanation for the sensitivity to biased coverage ratios. Increasing the value of k reduces the frequency of all k-mers (Fig. 6). A high k precludes k-mers from the low coverage genome from the assembly, since their frequency overlaps with that of k-mers introduced by sequencing errors. This prevents the formation of typically short contigs, thus optimizing the N50 size.



The WGS data of the *L. pustulata* thallus was assembled using *MIRA* and the contigs taxonomically assigned using *MEGAN* [13]. Both, fungal and algal, assemblies are much more fragmented than expected given the simulations. This is a result of the bacterial diversity also present in the thallus: nearly 1/3 of the total assembly is of bacterial origin (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.
- Real world data sets from lichen thalli not only contain fungal & algal DNA but also a large bacterial proportion.
- Mixing data from different species inflates the assembler performance differences, with *MIRA* & *SPAdes* performing best.

