

THE ASSEMBLATHON

N50 must die?

Genome assembly workshop, Santa Cruz, 3/15/11

twitter: [@assemblathon](#) web: assemblathon.org

Should N50 die in its role as a frequently used measure of genome assembly quality? Are there other measures out there? Are there other measures that are better? Read on to find out more.



Overview

Two parts of the talk. First a brief recap of what the Assemblathon is all about and who is in it, then a look at some of the metrics we have been working with.

UC Davis

Genome Center

Ian Korf

Dawei Lin

Keith Bradnam

Aaron Darling

Joseph Fass

Ken Yu

Vince Buffalo

UC Santa Cruz

Haussler Lab

David Haussler

Dent Earl

Benedict Paten

John St. John

Ngan Nguyen

Two groups involved in the evaluation of the Assemblathon. David Haussler's group will talk next.

The challenge

Assemble genome of species 'A'

112 MB diploid synthetic genome

Reference genome of species 'B' available

2 paired-read + 2 mate pairs libraries

Goals

What makes a good *assembler*?

Goals

What makes a good *assembly*?

This is less obvious. The best assembler in the world doesn't amount to much if other errors are made when putting the final assembly together. One of the checks that we put into the Assemblathon (E. coli contamination of known genome), helps differentiate between good assemblers and good assemblies.

Goals

What makes a good *assembly metric*?

Can we move beyond a reliance on N50? Can we see which measures work well together and which are revealing different aspects of a good genome assembly? Are there easy-to-calculate metrics that perform just as well as complex, computationally-involved metrics?

Assemblathon by the numbers



Assemblathon by the numbers

17 teams

Assemblathon by the numbers

62 assemblies

21 by UC Davis

If you want to develop a good assembly metric, it is useful to have both good and bad assemblies to measure. Some of our 21 assemblies are intentionally bad, e.g. not using mate pair information. Bad assemblies should score badly, so these assemblies provide an extra check on how useful our metrics are.

Assemblathon by the numbers

> 100 metrics

Many of which are clearly correlated to each other. Altogether, the metrics tackle the ‘what makes a good assembly?’ question from many different angles. Some are simple statistical descriptions of assemblies, and others are new measures that use more complex calculations.

Assemblathon by the numbers

7,949,257,310 bp



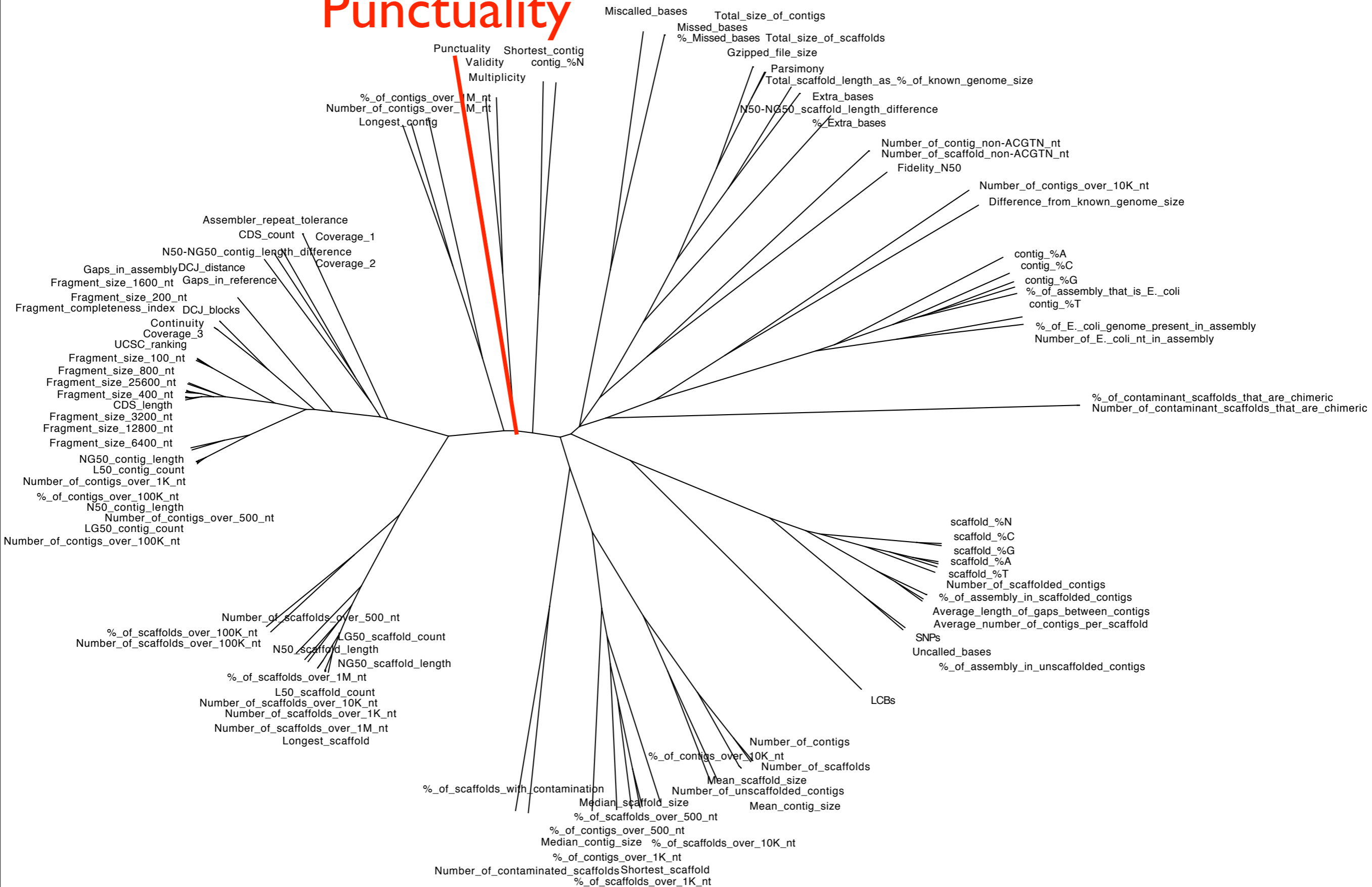
Metrics

What are we measuring?

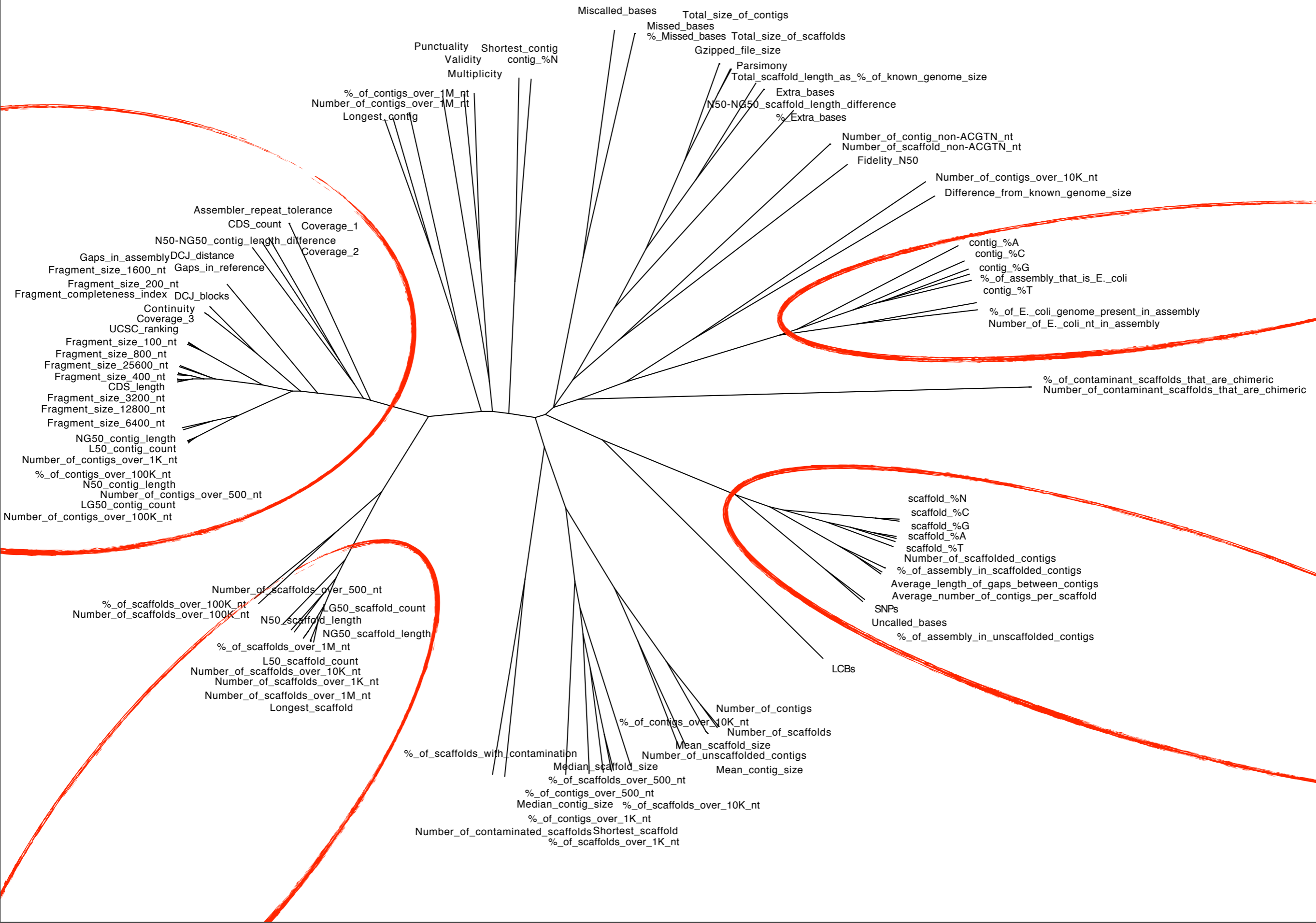
Number of scaffolds	N50-NG50 scaffold length difference	Fragment size 100 nt	% of assembly in scaffolded contigs
Total size of scaffolds	contig %A	Fragment size 200 nt	% of assembly in unscaffolded contigs
Total scaffold length as % of known genome size	contig %C	Fragment size 400 nt	Average number of contigs per scaffold
Longest scaffold	contig %G	Fragment size 800 nt	Average length of gaps between contigs
Shortest scaffold	contig %T	Fragment size 1600 nt	Number of contigs
Number of scaffolds > 500 nt	contig %N	Fragment size 3200 nt	Number of scaffolded contigs
% of scaffolds > 500 nt	contig %non-ACGTN	Fragment size 6400 nt	Number of unscaffolded contigs
Number of scaffolds > 1K nt	Number of contig non-ACGTN nt	Fragment size 12800 nt	Total size of contigs
% of scaffolds > 1K nt	Gzipped file size	Fragment size 25600 nt	Longest contig
Number of scaffolds > 10K nt	Difference from known genome size	Fragment completeness index	Shortest contig
% of scaffolds > 10K nt	Punctuality	Assembler repeat tolerance	Number of contigs > 500 nt
Number of scaffolds > 100K nt	Number of E. coli nt in assembly	Coverage 1	% of contigs > 500 nt
% of scaffolds > 100K nt	% of assembly that is E. coli	Coverage 2	Number of contigs > 1K nt
Number of scaffolds > 1M nt	Number of contaminated scaffolds	Coverage 3	% of contigs > 1K nt
% of scaffolds > 1M nt	% of scaffolds with contamination	Continuity	Number of contigs > 10K nt
Mean scaffold size	Number of contaminant scaffolds that are chimeric	Fidelity N50	% of contigs > 10K nt
Median scaffold size	% of contaminant scaffolds that are chimeric	Validity	Number of contigs > 100K nt
N50 scaffold length	% of E. coli genome present in assembly	Multiplicity	% of contigs > 100K nt
L50 scaffold count	Extra bases	Parsimony	Number of contigs > 1M nt
NG50 scaffold length	% Extra bases	UCSC ranking	% of contigs > 1M nt
LG50 scaffold count	Missed bases	LCBs	Mean contig size
scaffold %A	% Missed bases	DCJ distance	Median contig size
scaffold %C	CDS count	DCJ blocks	N50 contig length
scaffold %G	CDS length	SNPs	L50 contig count
scaffold %T	Number of scaffold non-ACGTN nt	Miscalled bases	NG50 contig length
scaffold %N	N50-NG50 contig length difference	Uncalled bases	LG50 contig count
scaffold %non-ACGTN			Gaps in reference
			Gaps in assembly

These are all the UC Davis metrics that we have been working with. Some are simple, others quite complex. Many are closely related to each others.

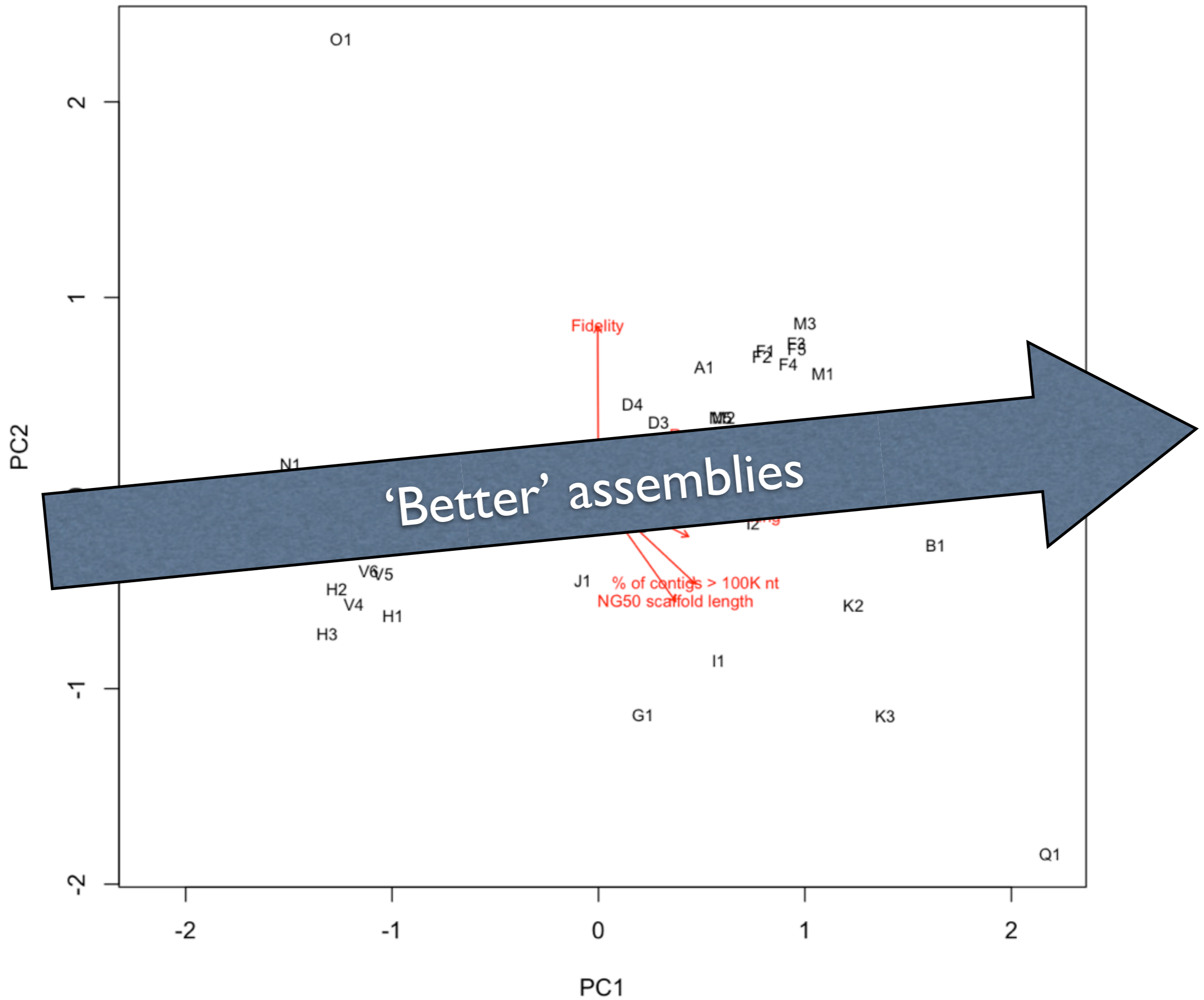
Punctuality



You can see how well any two metrics agree and from the strength of the correlation, produce a distance, and therefore produce a tree. One metric was based on how long it took groups to submit an assembly. Thankfully, this doesn't seem to correlate with any quality metrics. So groups who submitted later, did not seem to gain any obvious advantage.



Many metrics are clearly unrelated to each other, but some metrics group very closely together. Not always in the most expected fashion. E.g. % of assembly that is E. coli is highly correlated with contig %G



Principal components analysis can be used to see if assemblies can be grouped in any meaningful way. Overall, the first principal component (based on 9 selected metrics) seems to do a good job at ordering assemblies based on their overall quality.

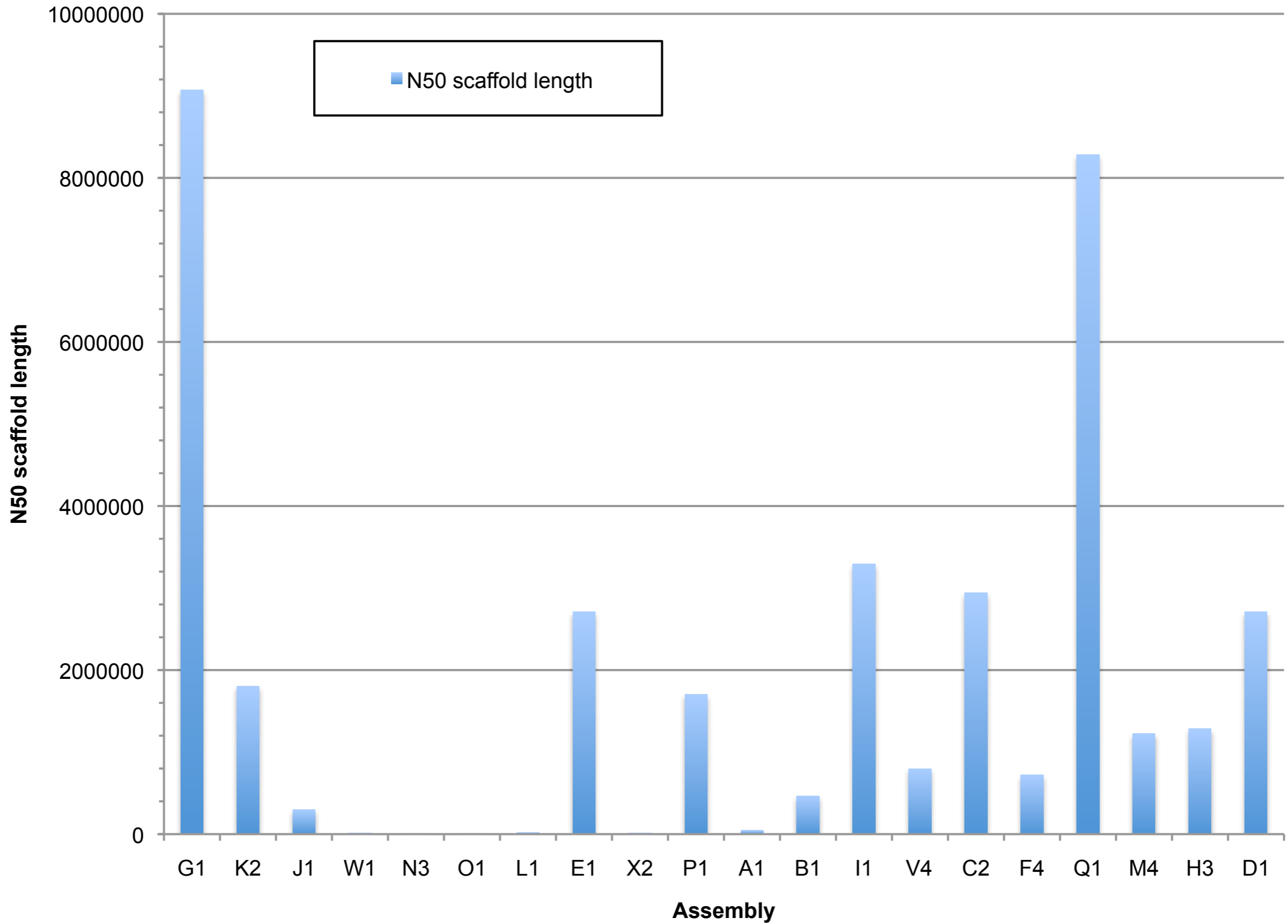
N50

Scaffold/contig length at which you have covered
50% of total assembly length

NG50

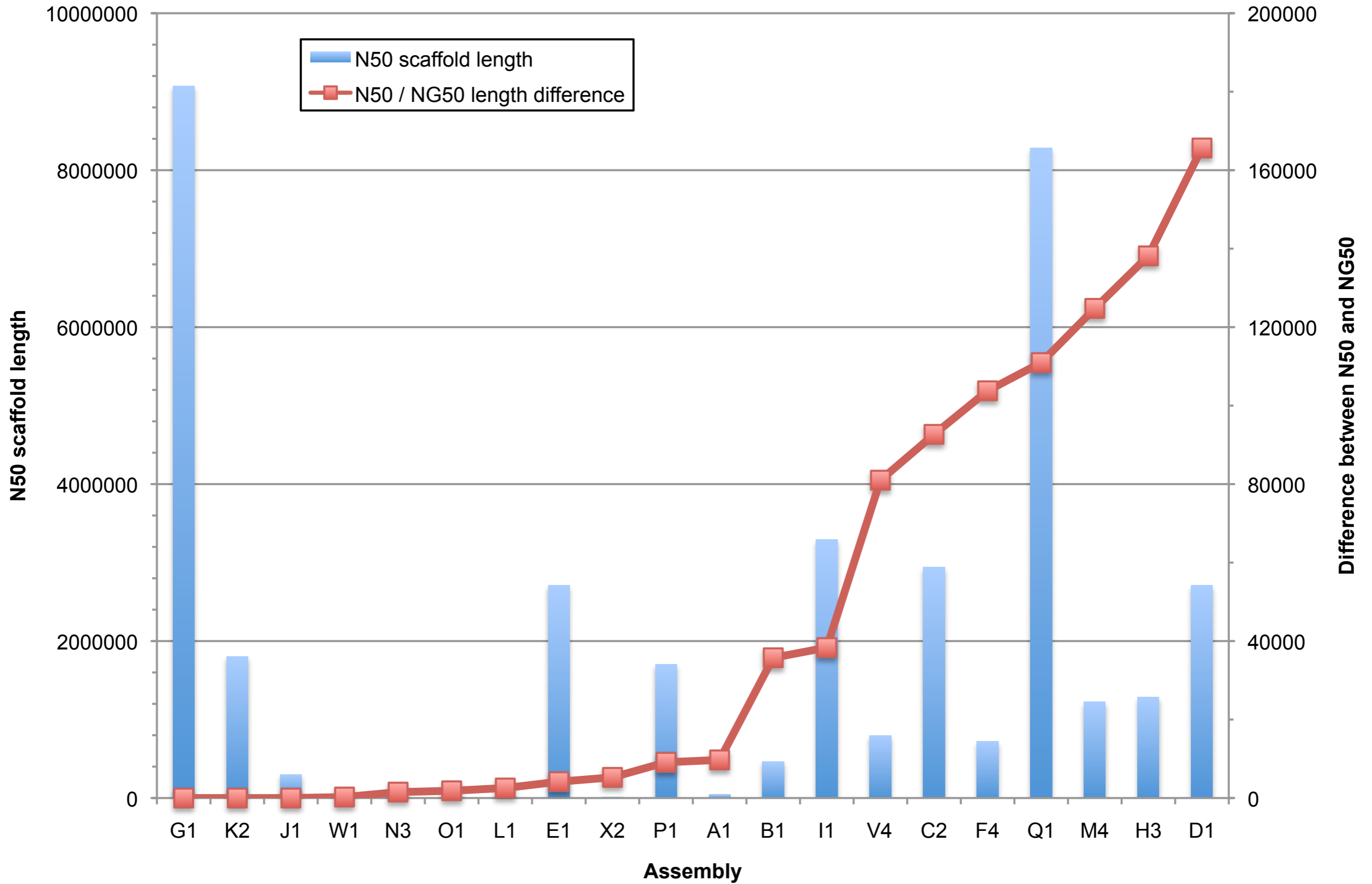
Scaffold/contig length at which you have covered
50% of total *genome* length

N50 vs NG50

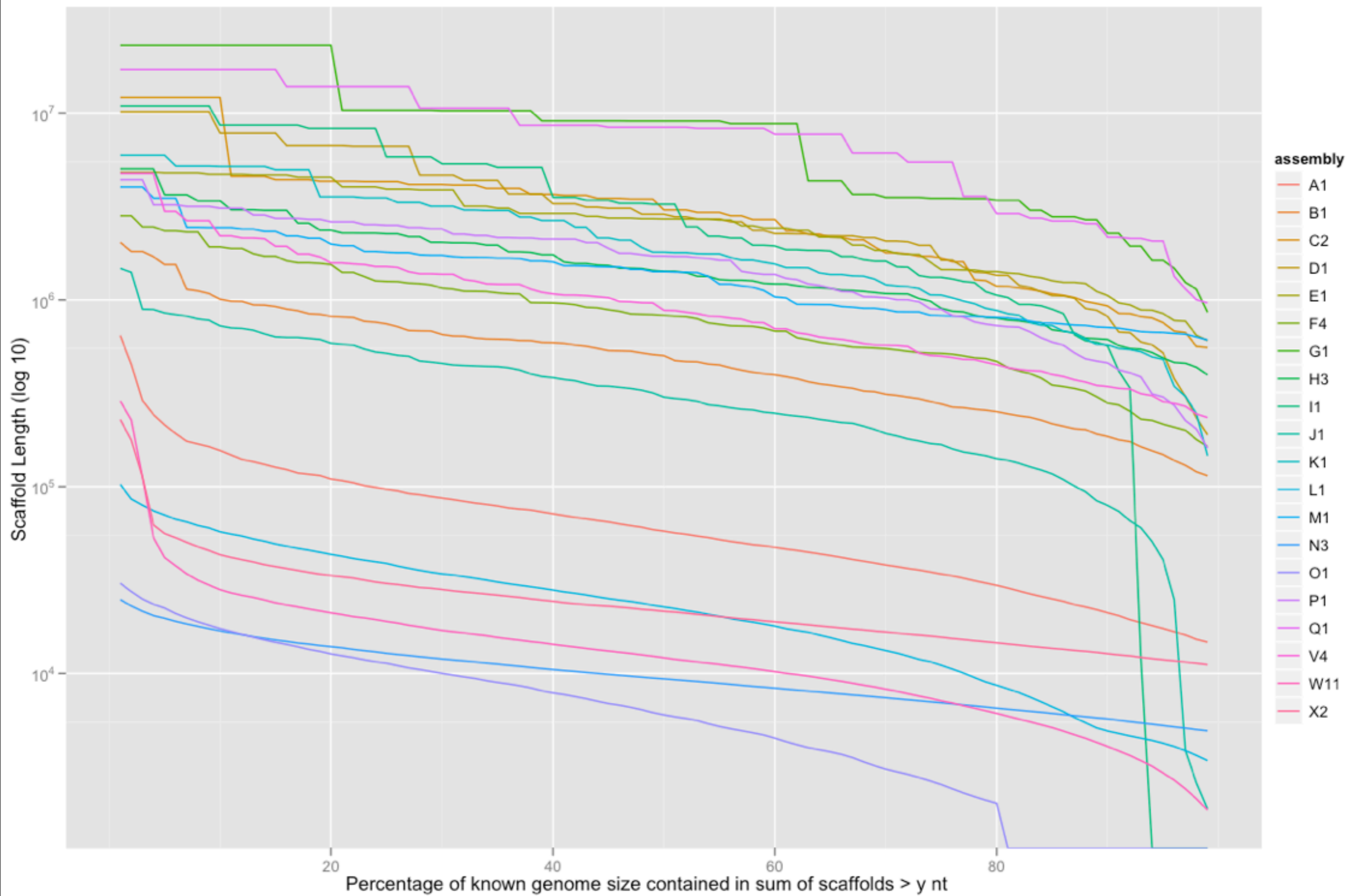


Huge variation in assemblies when looking at N50. G1 and Q1 are clear winners by this metric. Only showing best assembly from each team.

N50 vs NG50



Second series (in red) shows difference between N50 and NG50. Assemblies are ordered in increasing N50/NG50 difference. Some assemblies have no difference, some have large differences (over 160 kbp).



This NG(X) graph shows the values not just of NG50, but all values from NG1 through to NG99. Y-axis is on a log scale. This graph allows you compare all assemblies in a visual manner. Total area under the curve could be used as another assembly metric.

Fragment analysis

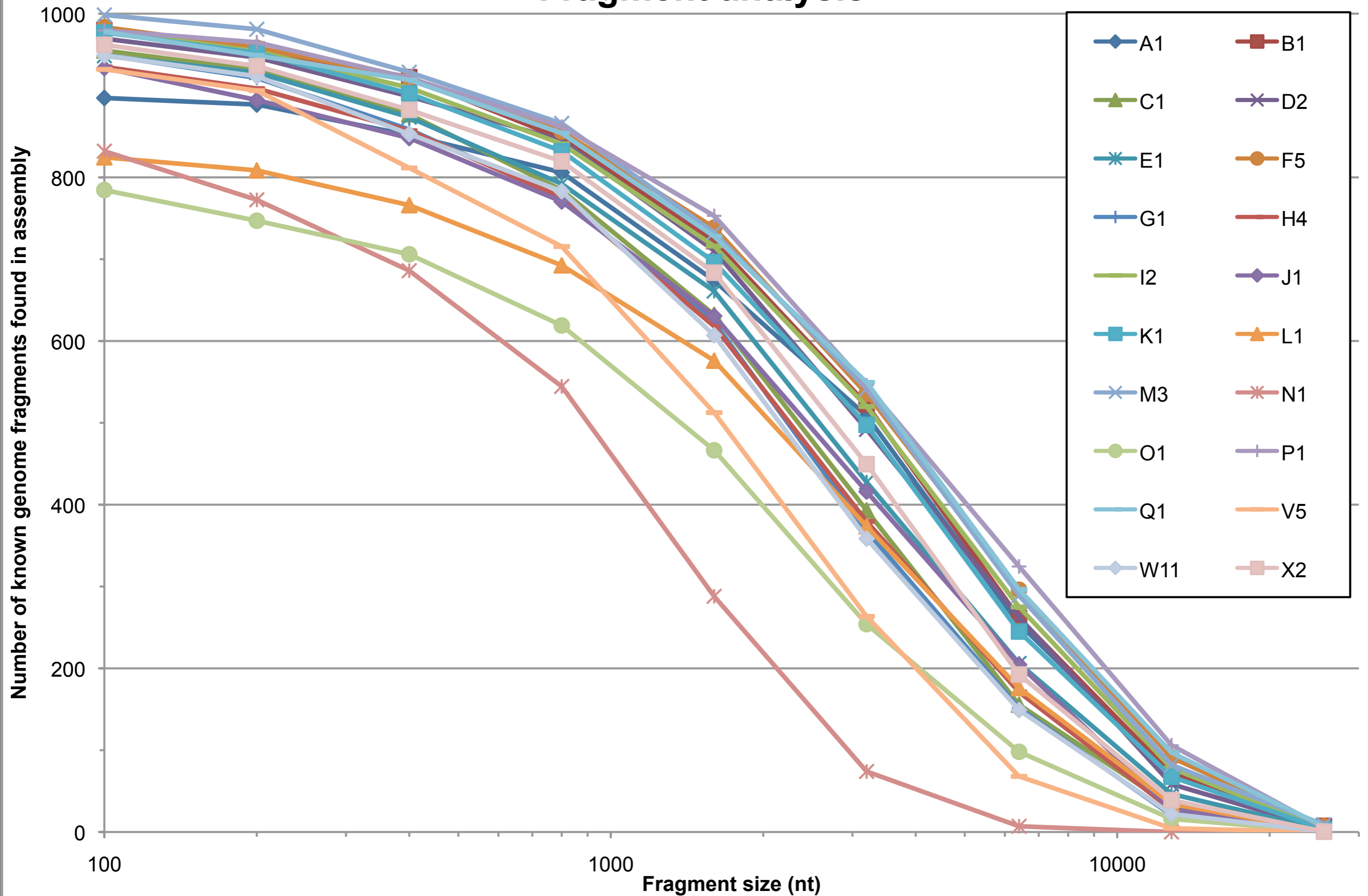
Count how many randomly chosen fragments from species A genome can be found in assembly

Use BLASTN, 95% identity over 95% length

Repeat at different fragment sizes

Repeat for both species A haplotypes

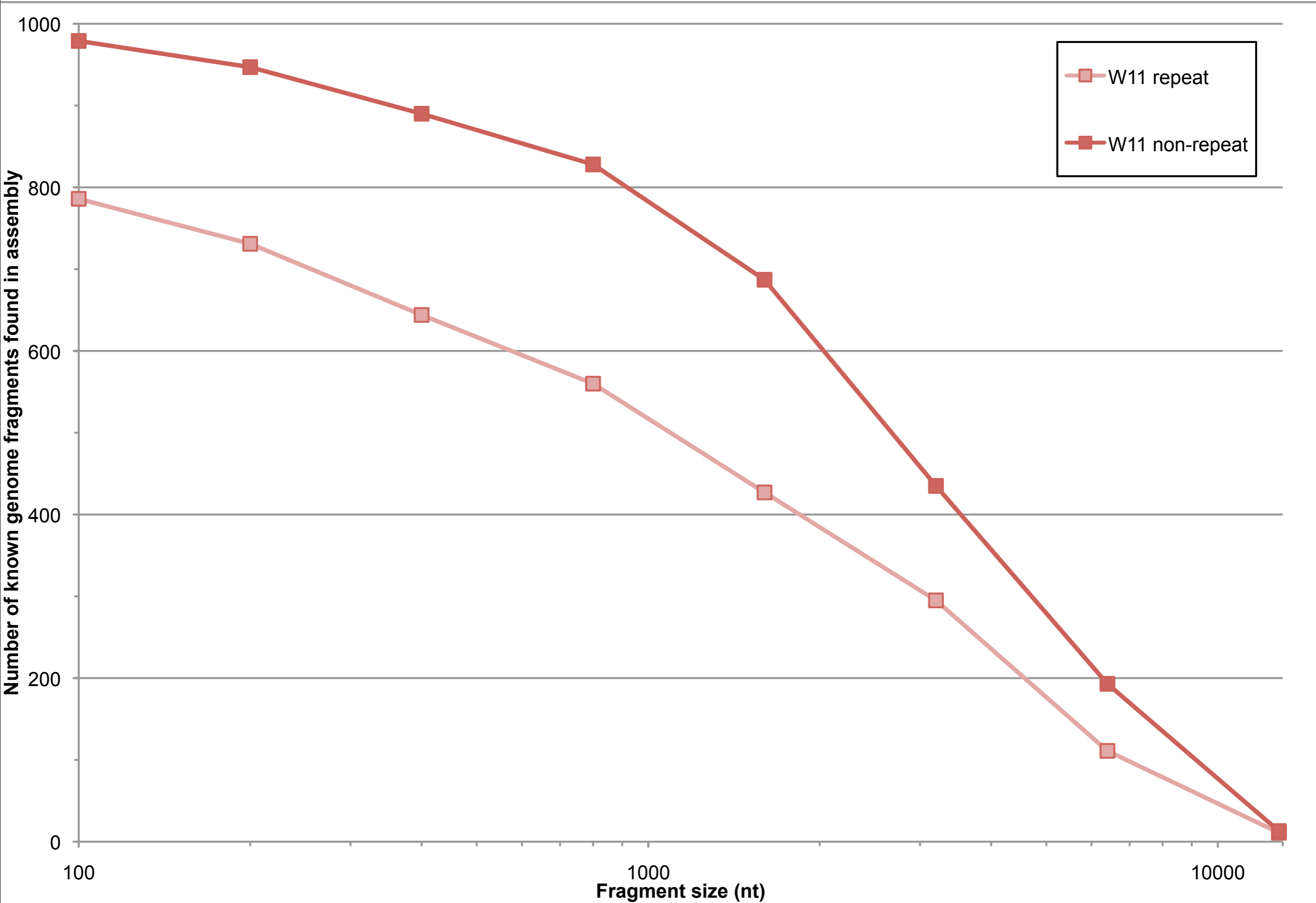
Fragment analysis



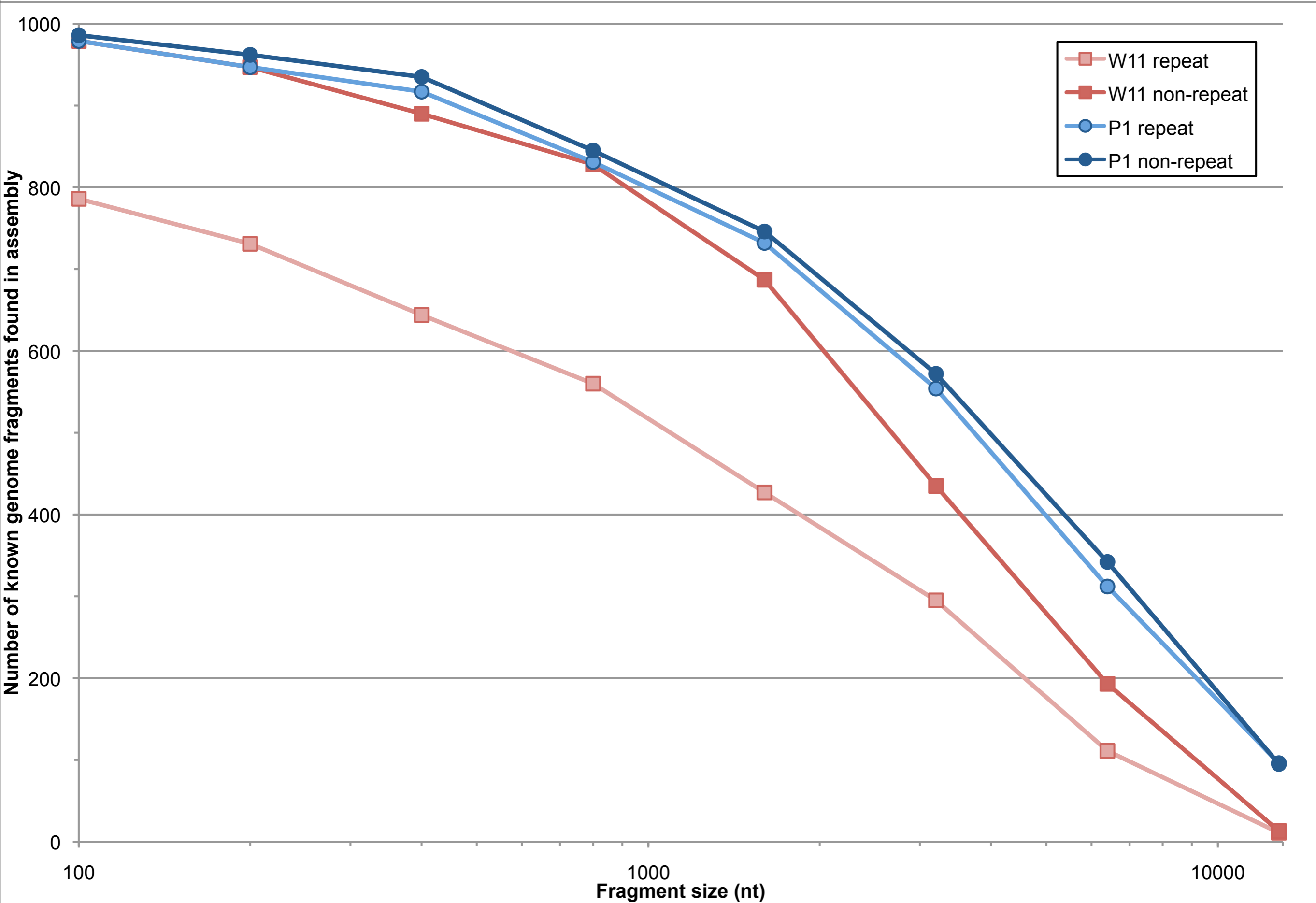
x-axis is log scale. Assemblies do relatively well at containing near-perfect regions of the known genome as long as fragments are short. Longer fragments can not be found because of errors and the problem of dealing with haplotype differences. The total area under the curve serves as a useful metric.

Repeat analysis

Choose fragments that either overlap
or don't overlap a known repeat



The (intentionally bad) W11 assembly does a poor job at coping with repeats, and the two lines on the above graph are far apart. The assembly contains fewer fragments of the species A genome that overlap known repeats. The sum difference between the repeat and non-repeat lines also makes for a useful metric.

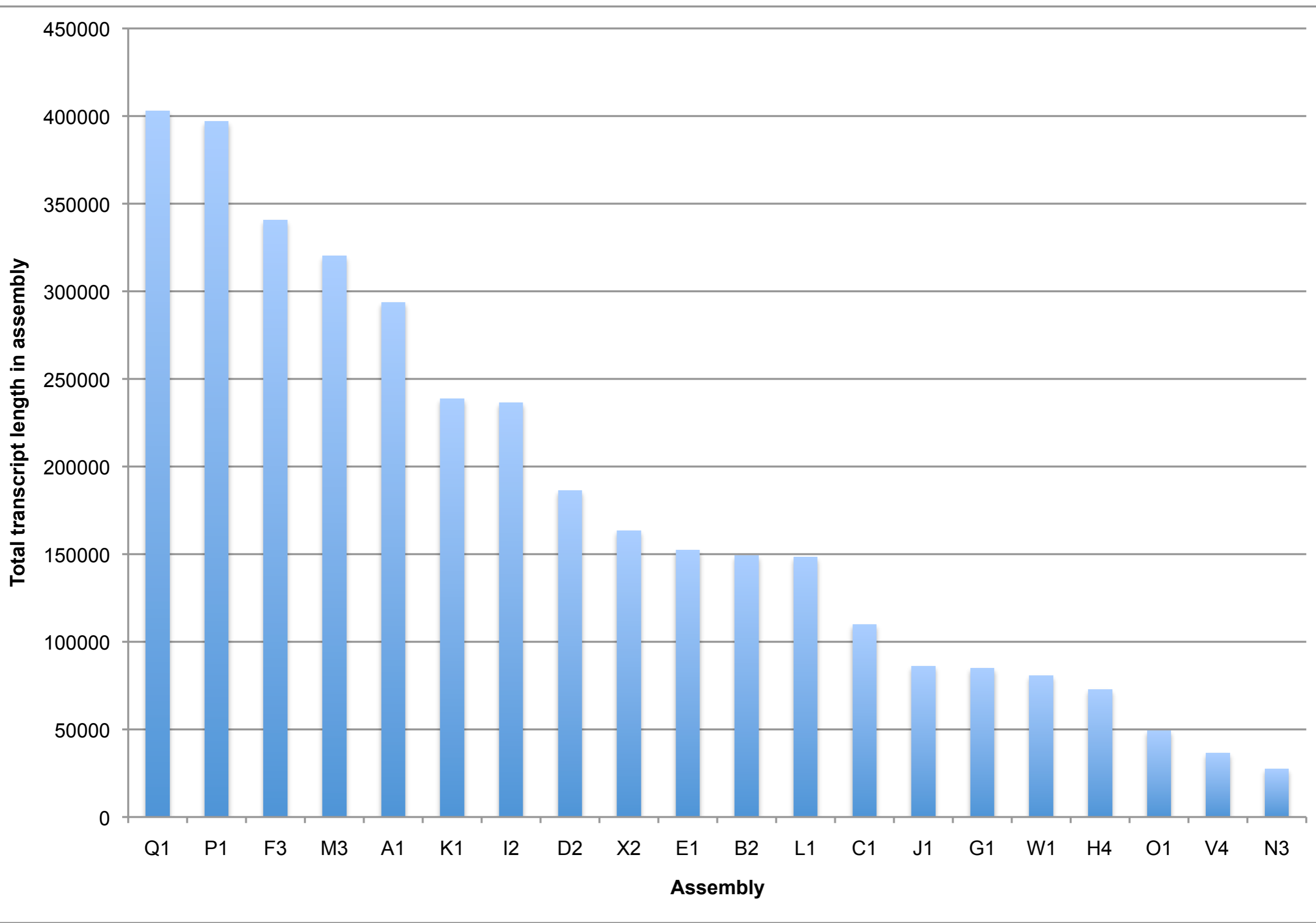


In contrast, the P1 assembly does a very good job at coping with repeats, and the two lines on the above graph are very close together.

Gene finding

How many genes are present in each assembly?

The Species A genome contains 176 genes. For many end users of a genome assembly the genes might be all they care about, and it doesn't matter how long contigs are...as long as they contain a full-length gene.

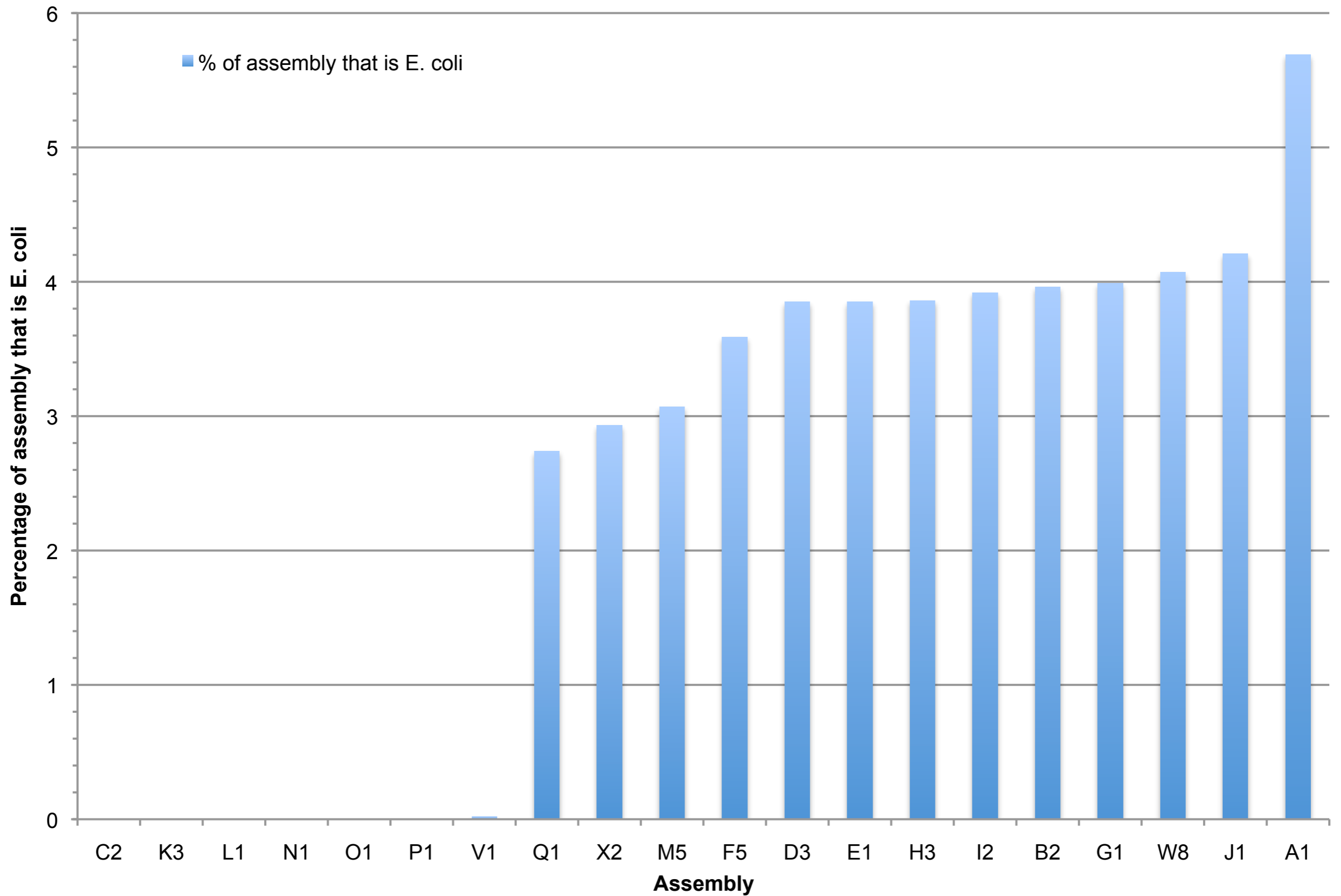


This shows the sum length of all 176 genes that can be found in each assembly.

Contamination

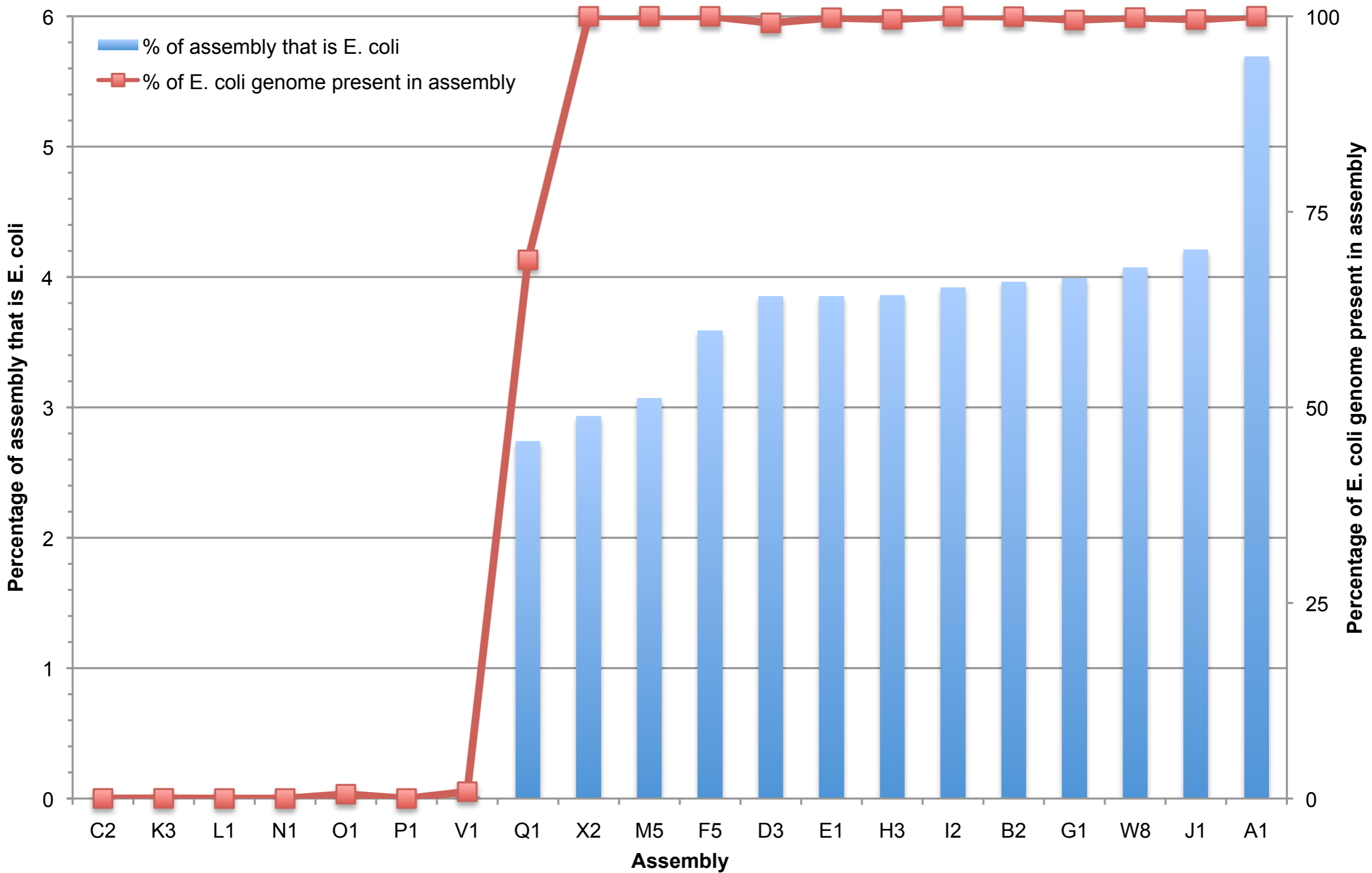
`“all libraries will contain
some bacterial contamination”`

Bacterial contamination



Assemblies either filtered out all contamination or they ended up with the E. coli genome sequence.

Bacterial contamination



The second series (in red) shows how much of the E. coli genome was present in each assembly. In most cases E. coli contamination produced contigs that consisted of just bacterial sequence. But some assemblies also had chimeric contigs with a mix of E. coli and species A sequence.

Mauve analysis

Uses whole genome alignment to reveal:

Miscalled bases

Uncalled bases

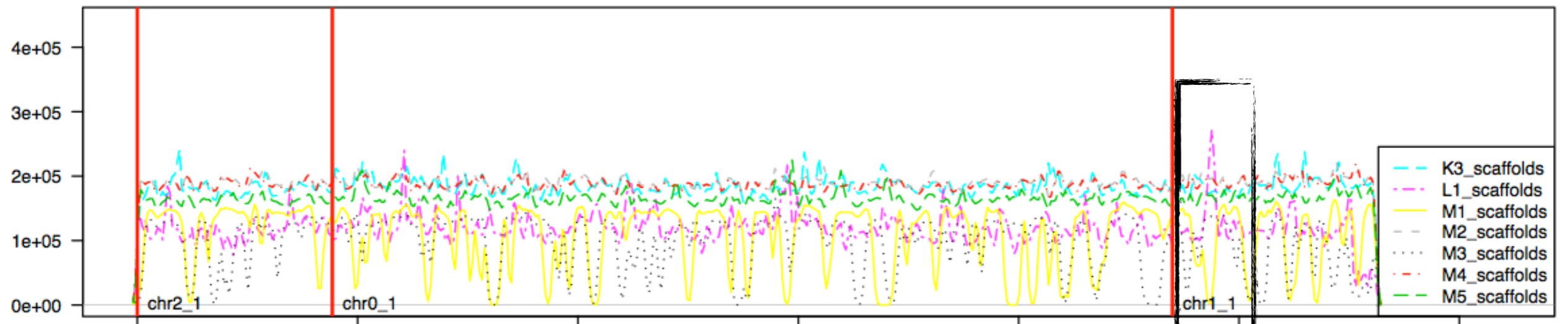
Missing bases

Extra bases

Misassemblies

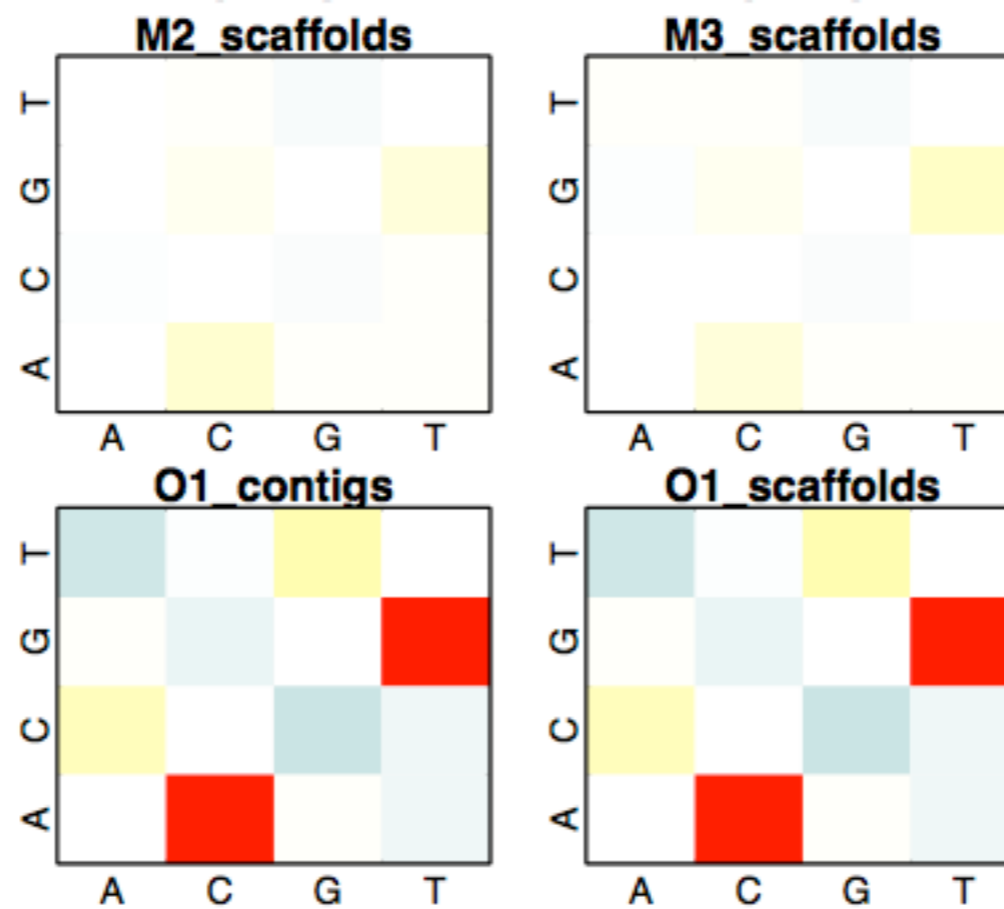
Double cut & join distance

Density of miscalled bases



Each panel shows analysis of 7 assemblies. Red lines denote boundaries between the 3 chromosomes of species A. Black box shows an area where many assemblies had a peak of miscalled bases.

Direction of miscalled bases



Most miscalled bases had no bias, but some assemblies had miscalled bases that were more likely to be a specific other base as this heat map shows.

BWA analysis

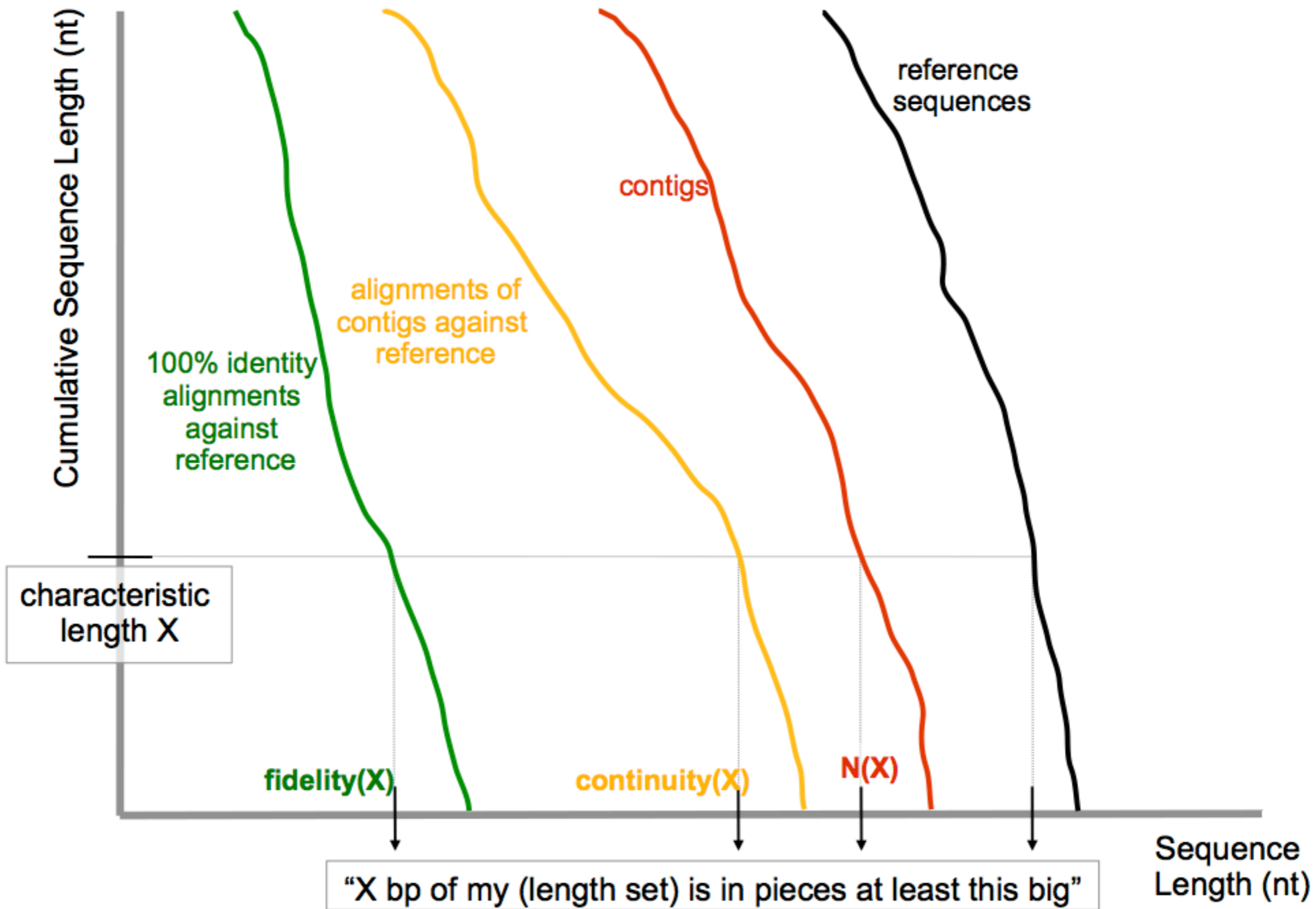
Align contigs to genome to reveal:

What fraction align?

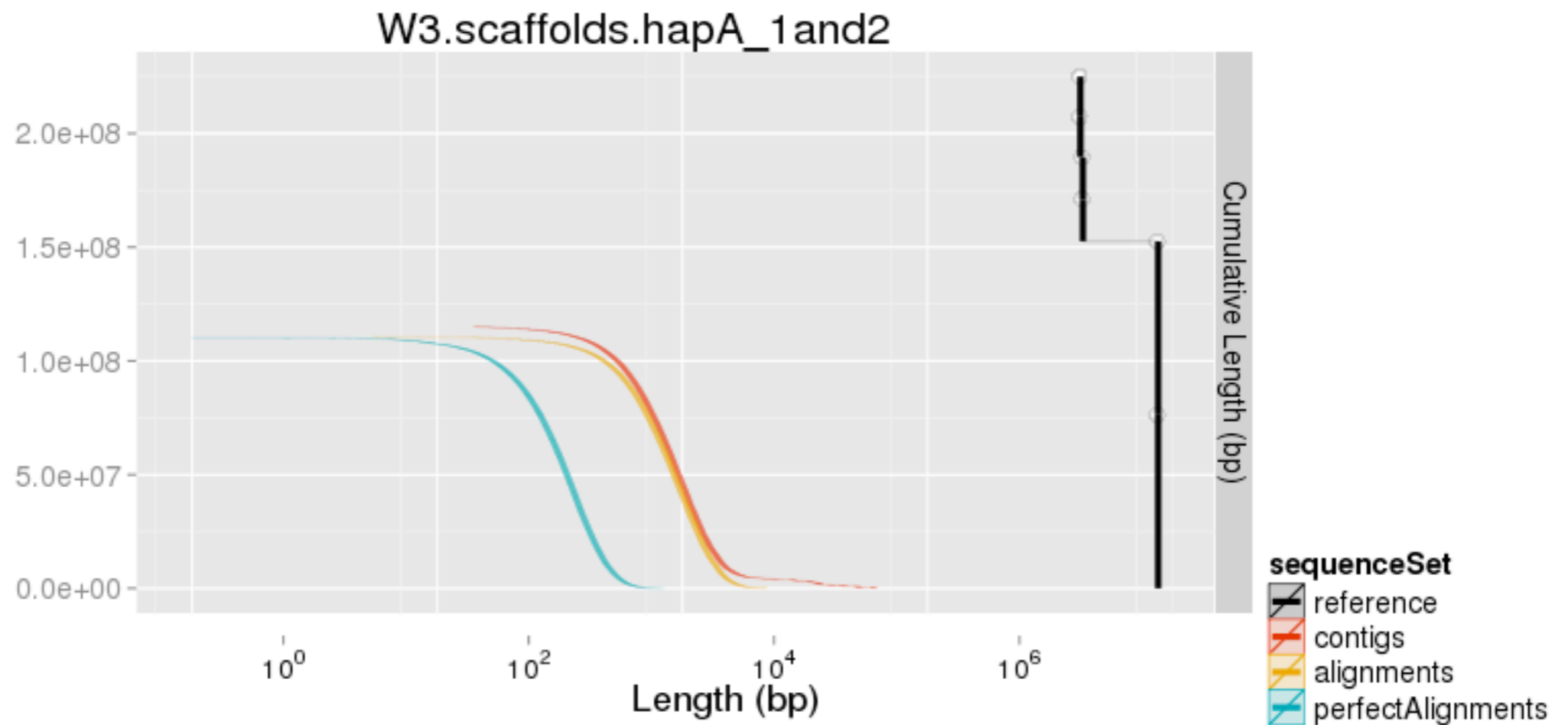
What fraction align *perfectly*?

Lengths of *coverage islands*

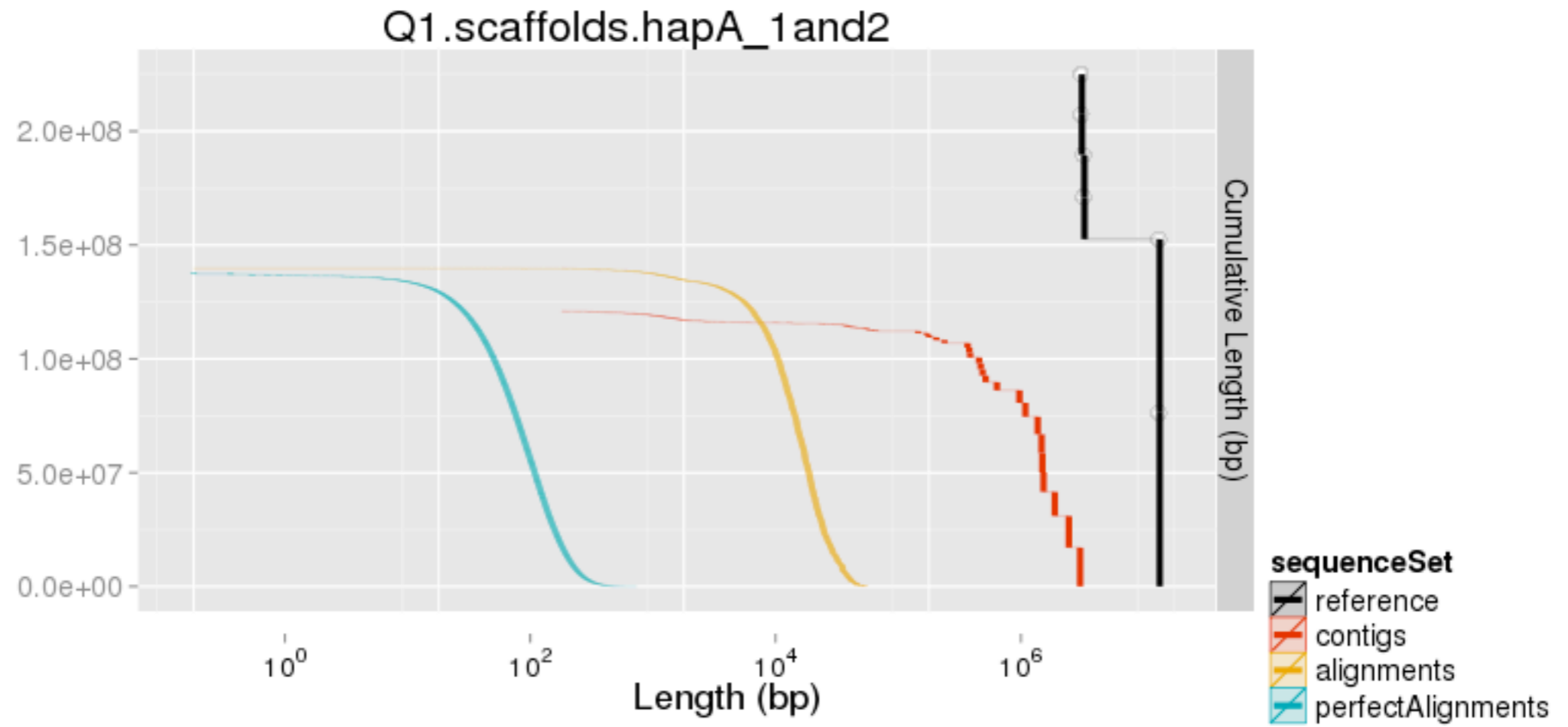
Measures of validity, multiplicity, parsimony,



For any set of sequences (contigs, aligned portion of contigs, or perfectly aligned portion of contigs), can calculate distribution of lengths and calculate measures that are analogous to N50. Continuity N50 = the N50 value for all alignable regions of contigs/scaffolds.



This is a deliberately poor assembly, with short contigs. Most contigs align well (close distance between red and yellow lines) but don't all align perfectly. Height on Y-axis denotes total size of each set of sequence lengths.



A better assembly with much longer scaffolds, and longer lengths of aligned sequence.



Summary

What have we learned?

Easy to produce a lot of data!

Use of any single metric can be misleading

Good assemblies tend to score well
across most, but not all, metrics

Good assembler \neq good assembly

End users may care little about all of this except for wanting to see the sequence of their gene of interest.