

Supplemental Material: sample quality

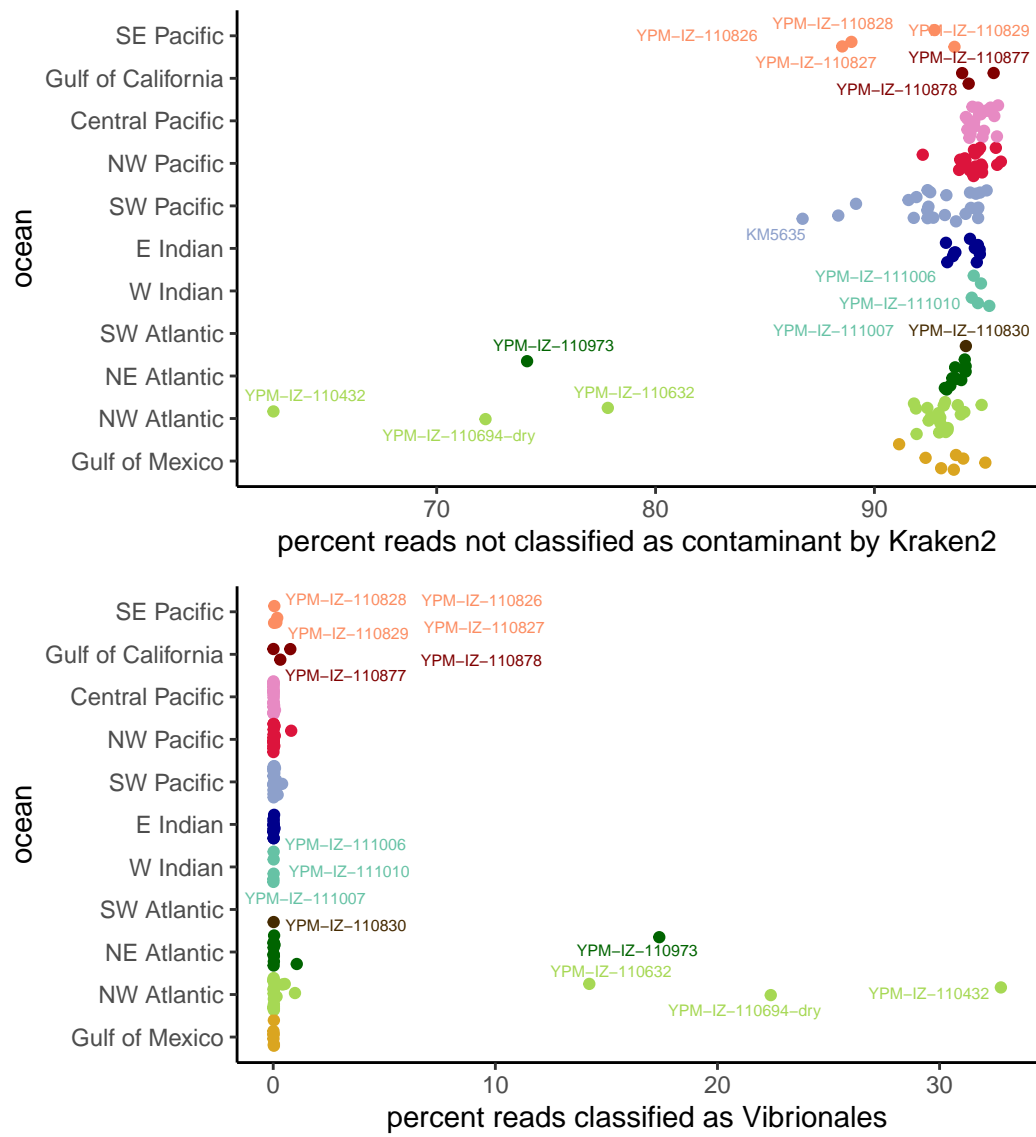
Sample quality

Kraken2: cross-species contamination

Kraken2 attempts to classify reads based on k-mers, against a standard database of known contaminants. **Kraken2** identifies several samples as having a large percentage of “classified” reads. These reads are largely classified as Vibrionales bacteria. Based on these, we will exclude the following samples as potentially highly contaminated:

- YPM-IZ-110432
- YPM-IZ-110973
- YPM-IZ-110632

Sample YPM-IZ-110694-dry was a resequenced sample from dried tissue. For main analyses we will use the other sample, YPM-IZ-110694, and not the dried specimen.



Cross-contamination and duplicated samples

In addition to YPM-IZ-110694-dry, we sequenced a second biological replicat (YPM-IZ-110474-2 of YPM-IZ-110474) and we generated two technical replicates (YPM-IZ-110269-A and -B of YPM-110269). We used these replicates as quality control for library preparation and sequencing, and as a reference for identifying duplicated samples.

We assesses similarity by calculating relatedness scores using the King-robust metric in PLINK2. Using a cutoff of 0.35 to detect potential duplicate samples, we identify the following pairs:

- YPM-IZ-106941 and YPM-IZ-106944
- YPM-IZ-110878 and YPM-IZ-110879
- YPM-IZ-111013 and YPM-IZ-111015
- YPM-IZ-111018 and YPM-IZ-111019

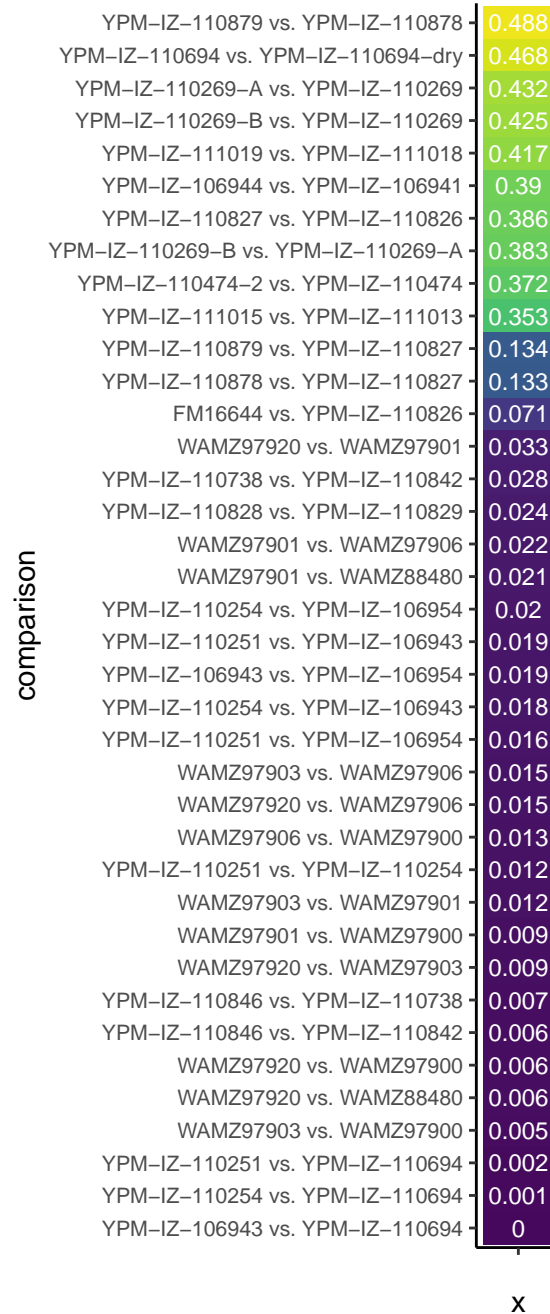
The latter sample in each pair was subsequently excluded.

We also identified two pairs of samples that had moderately high relatedness (0.07 and 0.13, respectively) with samples from improbably long distances. In both cases, DNA was extracted in the same batch, suggested cross contamination.

- YPM-IZ-110826 (Chile) with sample FM16644 (Canary Islands)
- YPM-IZ-110827 (Chile) with samples YPM-IZ-110878 and YPM-IZ-110879 (Gulf of California).

The former in each pair was subsequently excluded.

No other samples indicated a relatedness above the predicted score for a close relative.



Mapping statistics

We examine sample quality by mapping to our assembled genome. Samples with low mapping percentages indicate potential contamination or low-quality. We can count the number of total

reads, reads that map to the reference genome, and reads that are properly paired. From this we can detect a significantly lower mapping percentage for the sample:

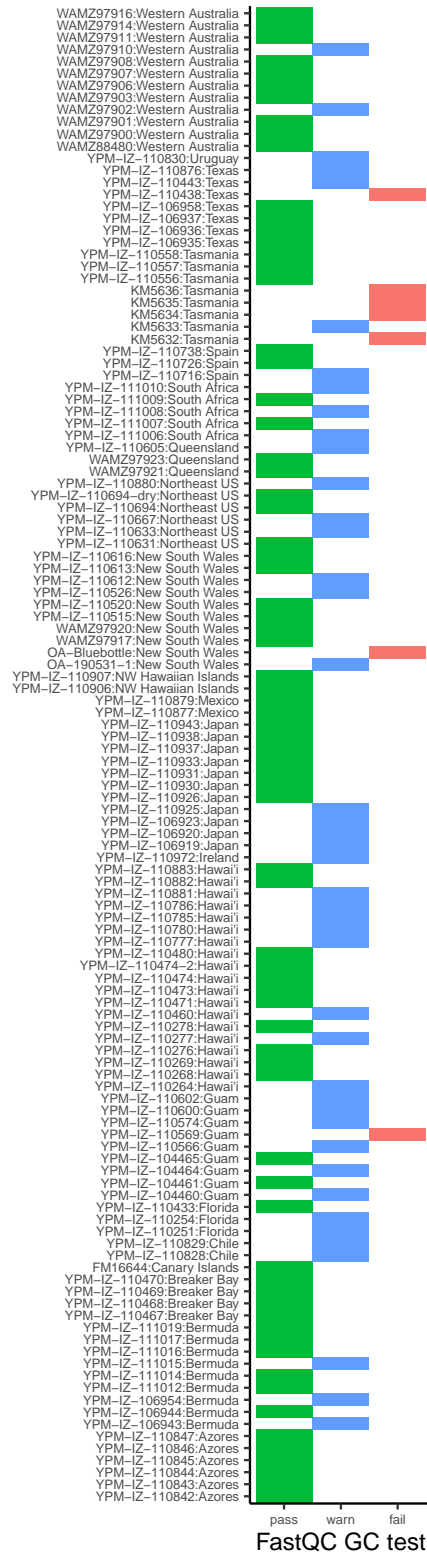
- YPM-IZ-110574

Tests using *in silico* PCR to extract 18S from this sample further indicated contamination from arthropod DNA. This sample will likewise be excluded.

Comparing statistics on mapped reads vs. properly paired reads shows no significant outliers in terms of quality.

FastQC: GC content

FastQC flagged no sequences as having low sequence quality, all were retained. But **FastQC** flagged several samples as having GC contents distinct from the expected distribution. Based on these results we will flag eight samples that fail the GC test and exclude those from our strict analyses (but retain them in general analyses presented in the manuscript).

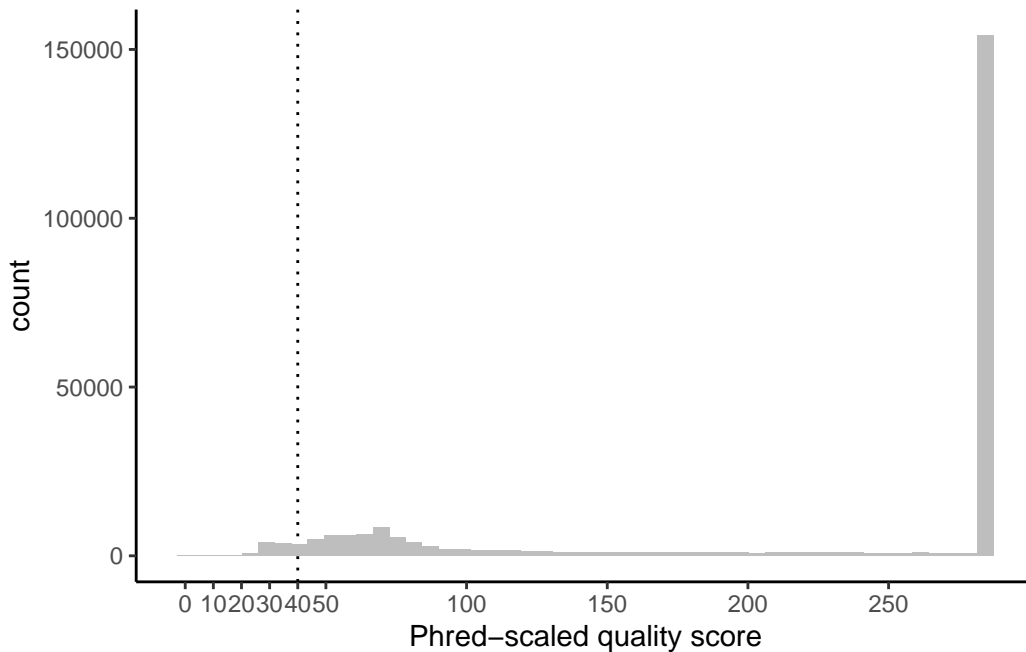


Establishing appropriate filters

Site quality

Most sites have high phred quality scores (e.g. a Phred encoded score of 40 indicates a 1 in 10,000 chance of an erroneous call). Based on this, we set the following:

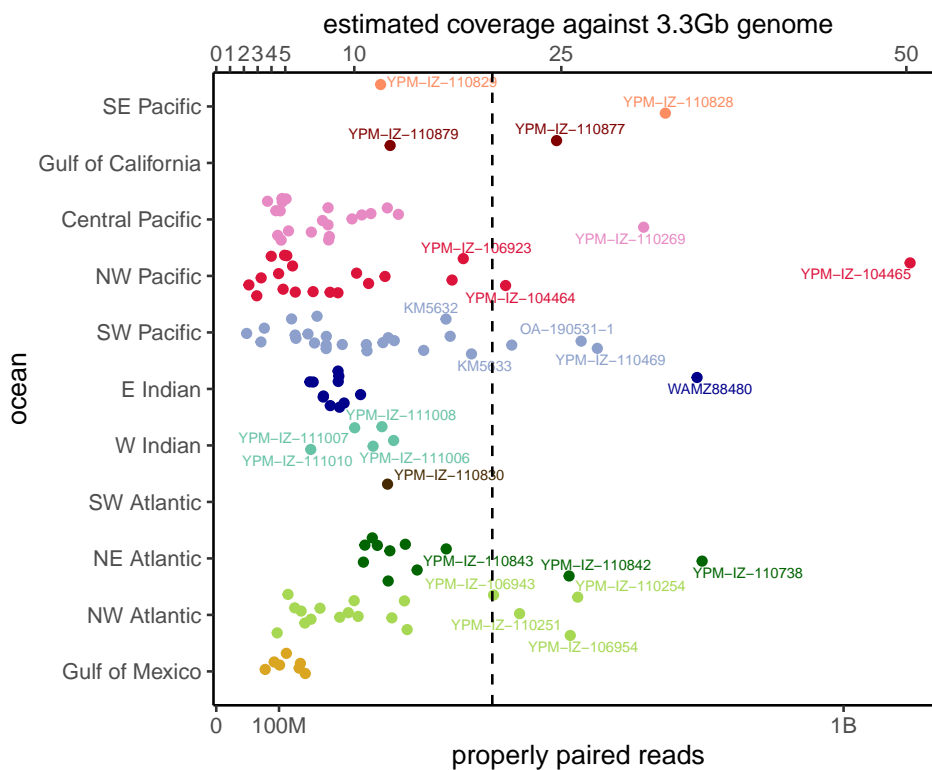
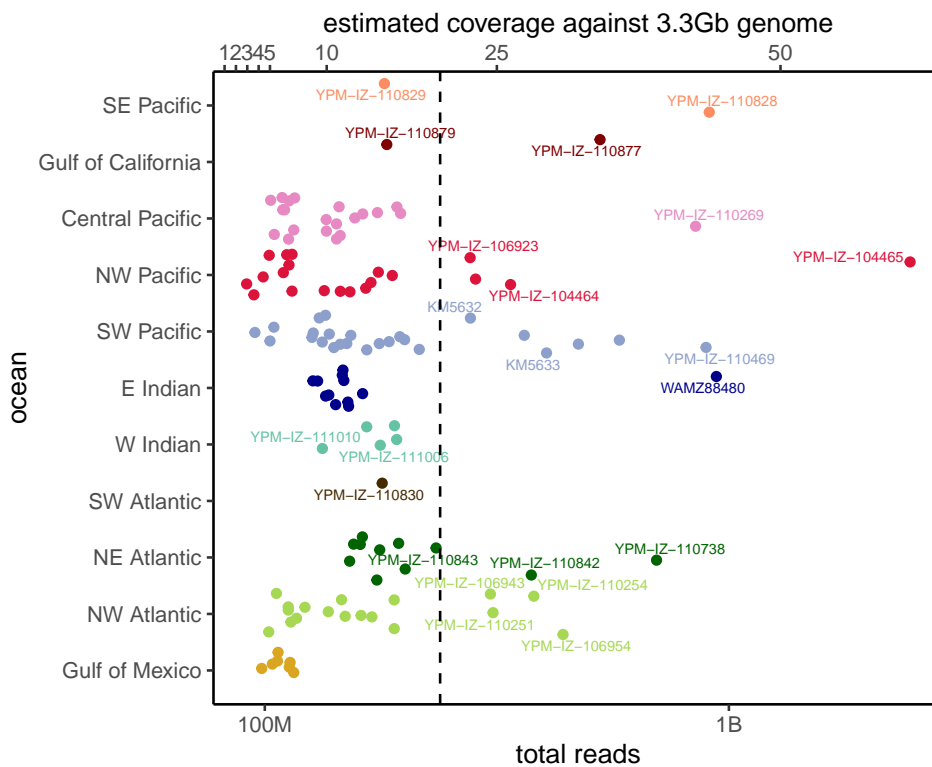
- minimum site quality score = 40



Coverage

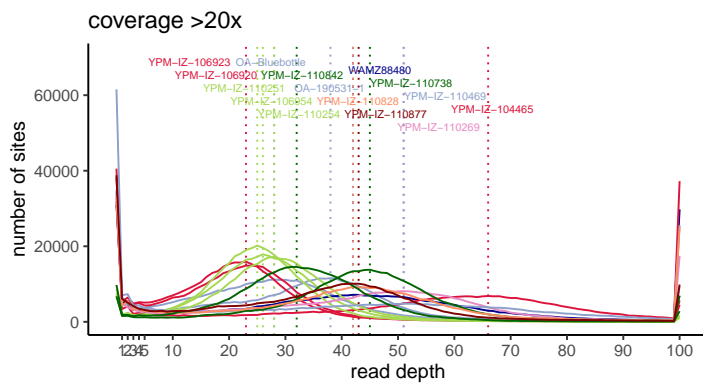
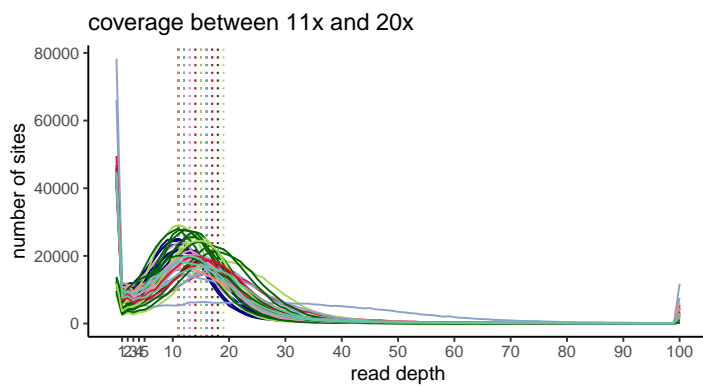
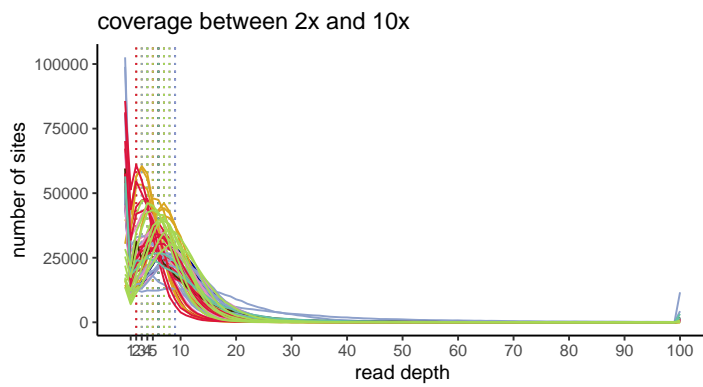
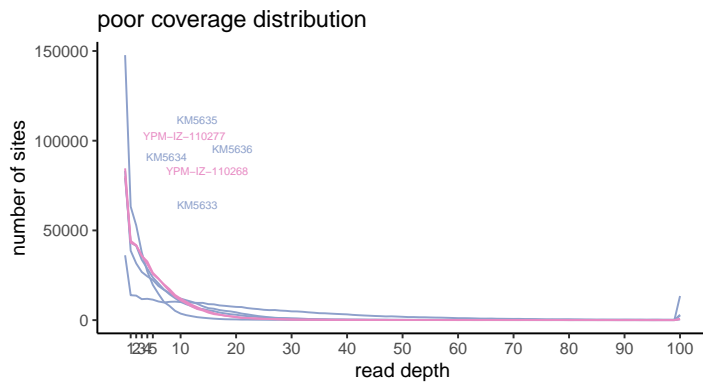
We can estimate coverage by dividing the number of basepairs sequenced for each sample by the length of the reference genome, 3.3Gb. Samples were sequenced to variable depths, with estimated coverage ranging from 5-60x.

Samples with greater than 20x target coverage (dashed line) were considered high coverage, and were selected for genome size estimation and other high-coverage analyses.



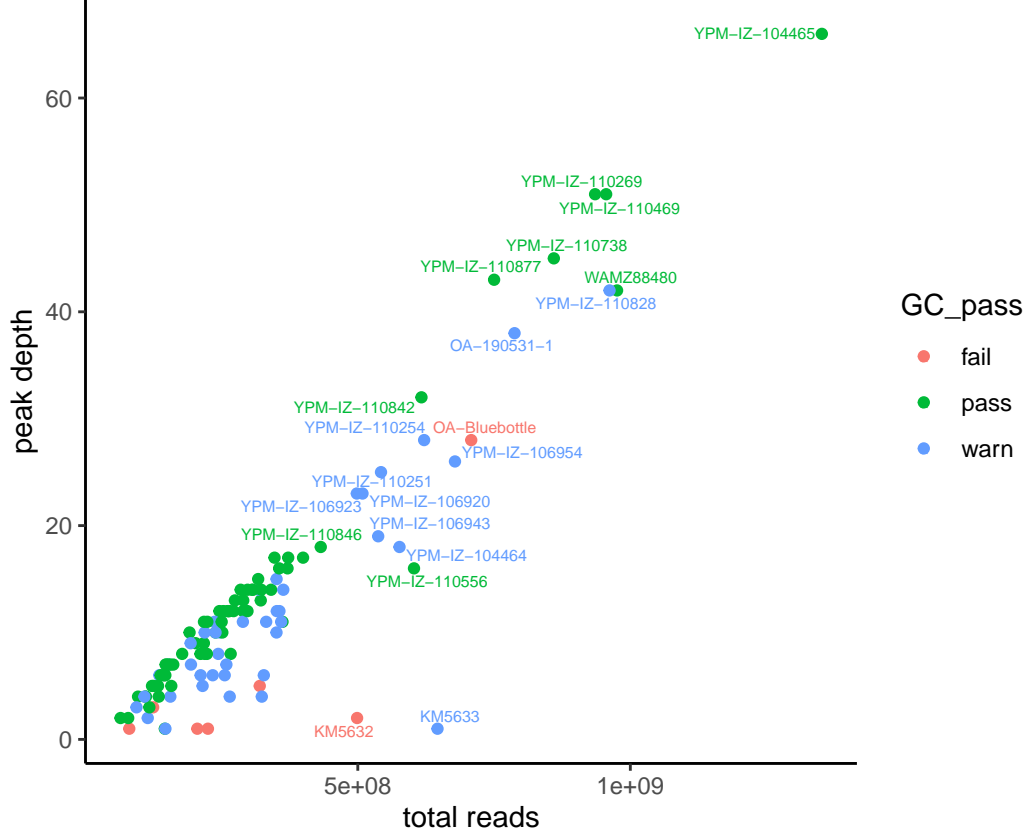
We use **angsd** to estimate realized depth across samples and sites using a subset of genomic regions. We categorize samples into those with poor coverage distributions, those with low coverage (between 2 and 10x), moderate coverage (11 and 20x) and high coverage (>20x). Poor coverage samples were the following:

- KM5633
- KM5634
- KM5635
- KM5636
- YPM-IZ-110277
- YPM-IZ-110268



From the previous distributions we calculated peak depth, and then compared to total reads. We observed that GC distribution, as assessed using **FastQC** impact the relationship between input reads and realized depth.

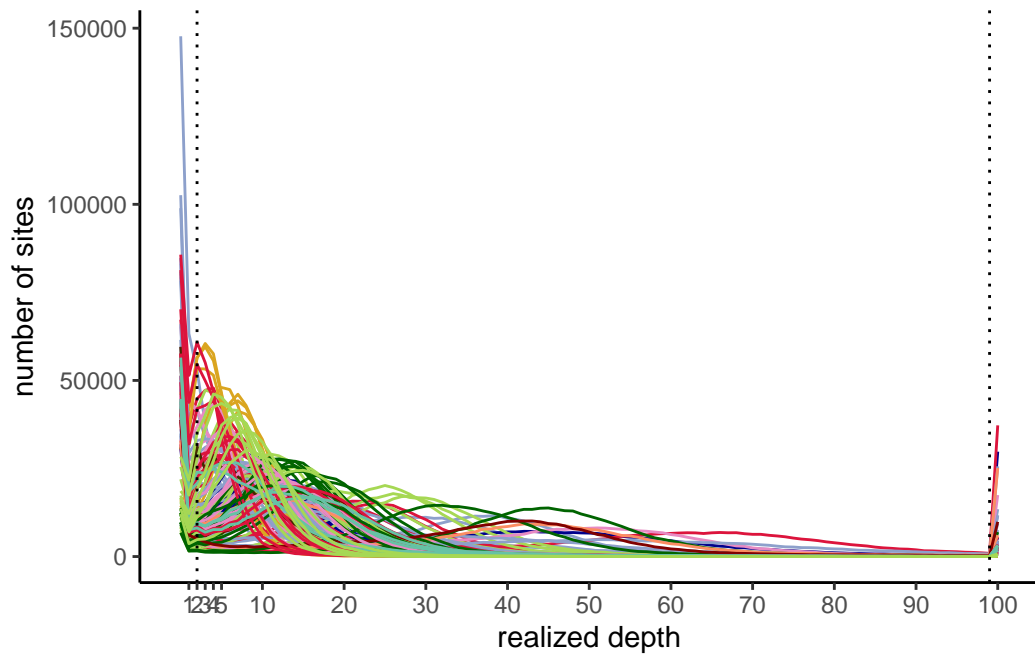
To address this, we run an additional robustness test using only samples that pass the GC test.



We used realized depth to establish the cutoff for minimum and maximum depth across sites.

For **ANGSD** analyses that are robust to low coverage, we did not perform any filtering based on depth, but for **BCF** based analyses (e.g **Fst** calculations), we set the following filters across sites for each sample:

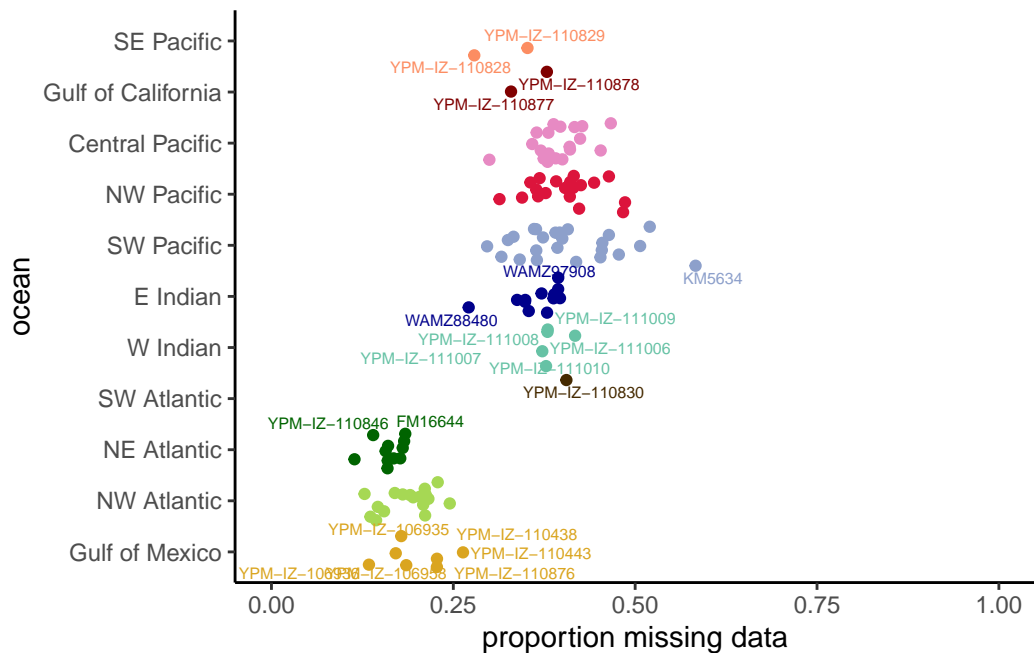
- minimum depth = 2x
- maximum depth = 99x



Missingness

We examined the proportion of missing sites across samples. This was calculated using `vcftools` on a random subsample of sites from across genome regions. From this we can see identify several samples with a high proportion of missing sites, in particular:

- KM5634



We can use the distribution of missing samples across sites to establish an appropriate filter for missingness. Here we set the following:

- tolerate missingness equal to or below 75%

This means we exclude sites for which $>25\%$ of samples are missing data.

