

PopPsiSeq Summary

Charlie Soeder

11/20/2018

Contents

| | | |
|----------|-------------------------------------------|----------|
| 1 | Introduction | 1 |
| 2 | Materials, Methods, Data, Software | 1 |
| 2.1 | Reference Genomes | 1 |
| 2.2 | Sequenced Reads | 1 |
| 2.2.1 | Pre-processing | 2 |
| 2.3 | Mapped Reads | 4 |
| 2.3.1 | Read & Alignment Quality | 5 |
| 2.3.2 | Depth & Breadth of Coverage | 5 |
| 3 | Results | 7 |
| 4 | References | 7 |
| 4.1 | Software | 7 |
| | Bibliography | 8 |

1 Introduction

Explain motivation, overview of PsiSeq and PsiSeq2

Population-based approach, rather than ancestral

2 Materials, Methods, Data, Software

2.1 Reference Genomes

The droSim1 and droSec1 reference genomes were downloaded in FASTA format from UCSC Genome Browser.

2.2 Sequenced Reads

A backcross and introgression experiment was performed, in which simulans females were mated with sechellia males, and the hybrid offspring were selected for avoidance of morinda odorants. The offspring were sequenced after 15 rounds of backcrossing and introgression (Earley and Jones 2011). One sample was sequenced in this experiment; a follow-up experiment generated three more samples with two replicates each. As a control, several wild-type sechellia sequences were downloaded from NCBI:

```
## Warning: package 'bindrcpp' was built under R version 3.4.4
```

Table 1: Sequenced Experimental Samples

| name | paired | experimental | source |
|------------|--------|--------------|--------|
| SRR5860570 | TRUE | control | NCBI |

| name | paired | experimental | source |
|-----------|--------|--------------|----------------|
| SRR303333 | FALSE | selection | EarlyJones2011 |
| 10A | TRUE | selection | EarlyJones2013 |

For population-wide data, wild *D. simulans* and *D. sechellia* flies were captured and sequenced by Daniel Matute:

Table 2: Number of Sequenced Samples per Species

| species | sample_count |
|-----------------------------|--------------|
| <i>drosophila sechellia</i> | 3 |
| <i>drosophila simulans</i> | 3 |

2.2.1 Pre-processing

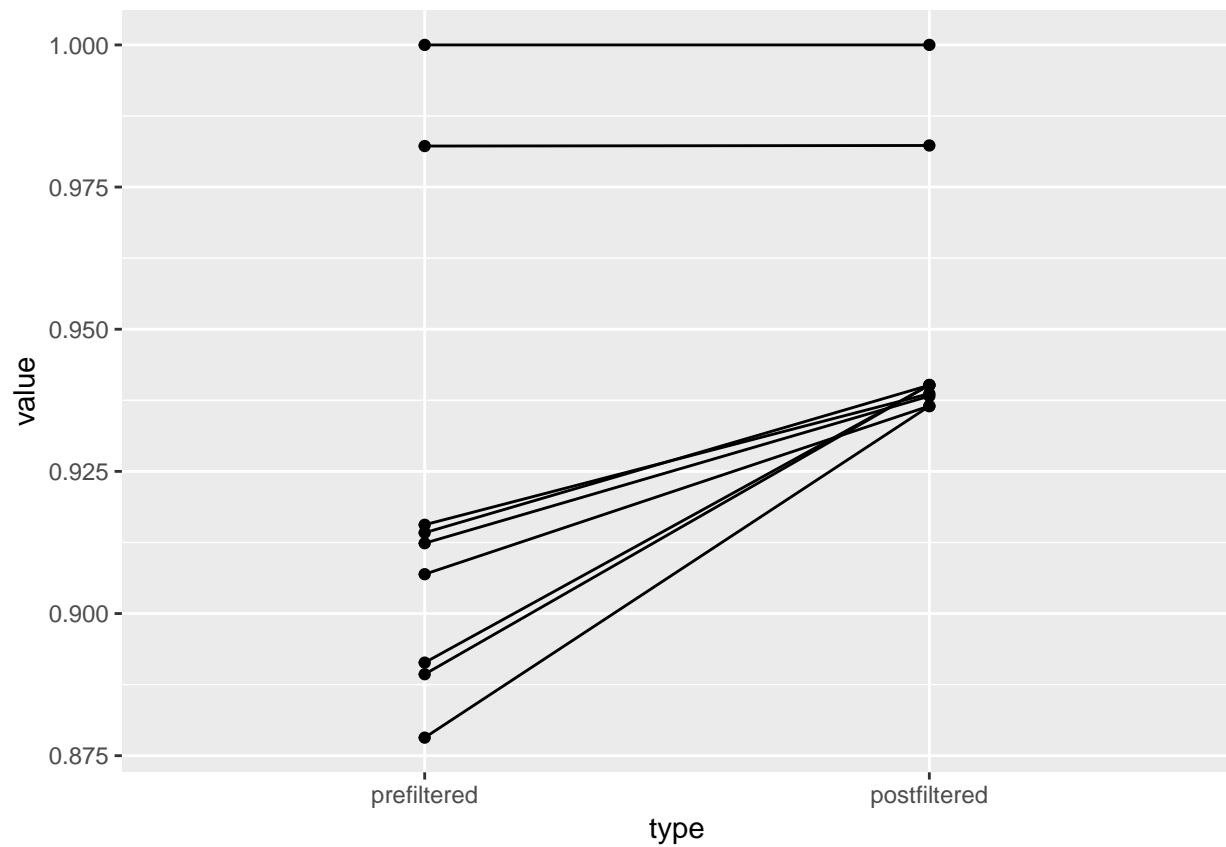
These reads were preprocessed with FASTP (S. Chen et al. 2018) for quality control and analytics.

Starting FASTQ files contained a total of 400M reads; after QC, this dropped to 379M.

Table 3: Read Count and Percent Retention

| type | minimum | average | maximum |
|-------------------|---------|---------|---------|
| prefiltered | 1.48 M | 44.5 M | 85.4 M |
| postfiltered | 1.4 M | 42.1 M | 81.1 M |
| percent retention | 93.3 | 95.6 | 100 |

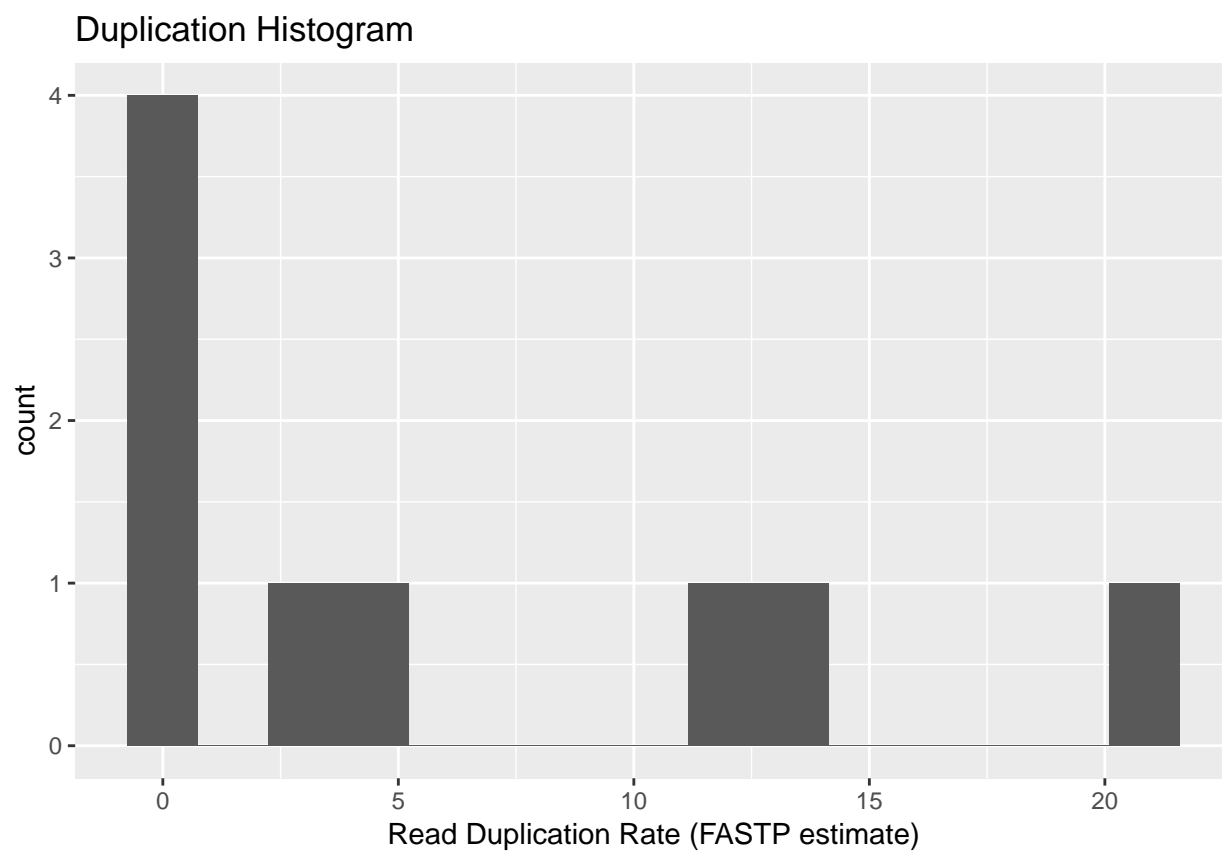
Filtration also increased the read quality, as seen in the increase in the fraction of reads with an average quality score > 30 :



Duplicate reads were also detected; these will be filtered during alignment:

Table 4: Percentage Duplication

| minimum | average | median | maximum |
|---------|---------|--------|---------|
| 0.1 | 6 | 2.6 | 20.9 |



2.3 Mapped Reads

Reads were first mapped to a reference genome using the BWA SAMPE/SE algorithm. Then, the alignment file was filtered for uniqueness (ie, a read must be aligned optimally with no alternative or runner-up hits, “XT:A:U.X0:i:1.X1:i:0”), mapping/sequencing quality (“-q 20 -F 0x0100 -F 0x0200 -F 0x0300 -F 0x04”), and deduplication.

2.3.1 Read & Alignment Quality

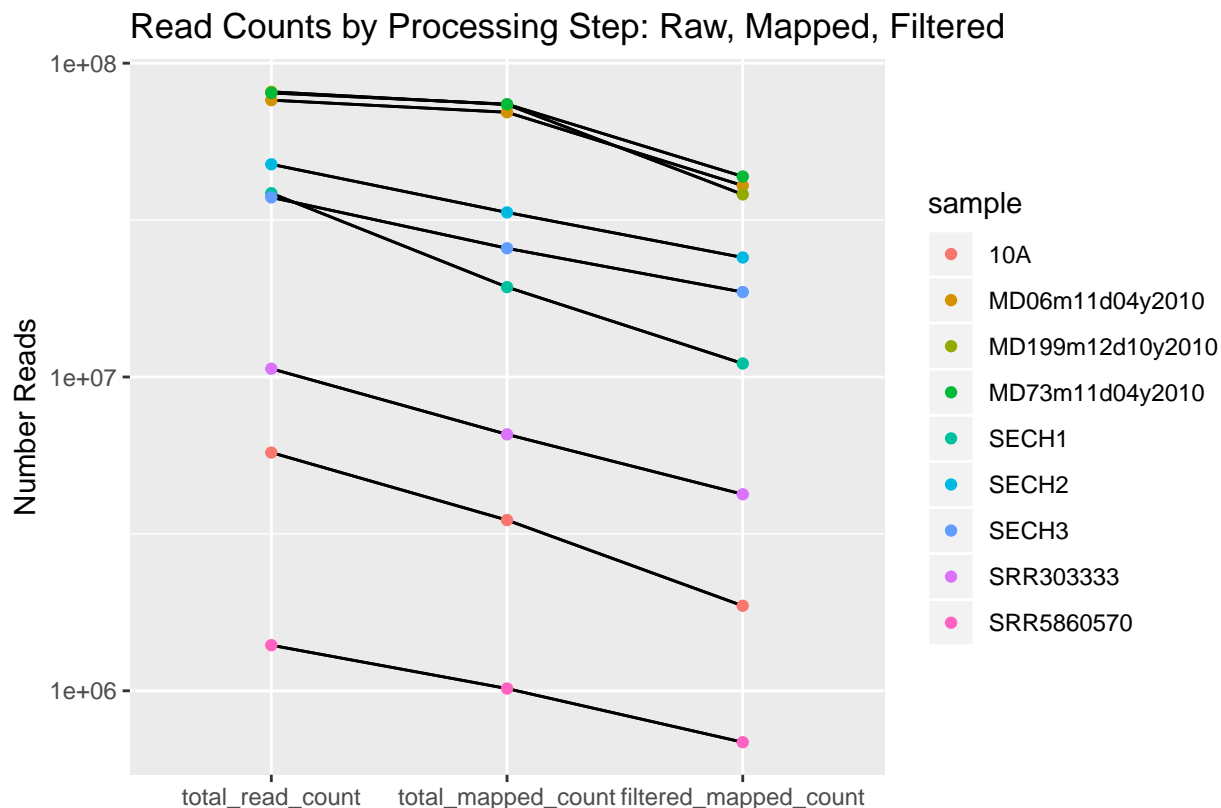


Table 5: Read Counts During Alignment & Filtration

| measure | minimum | average | median | maximum |
|-----------------------|---------|---------|--------|---------|
| filtered_mapped_count | 686 k | 20.3 M | 18.6 M | 43.5 M |
| total_mapped_count | 1.02 M | 34.1 M | 25.7 M | 74 M |
| total_read_count | 1.4 M | 42.1 M | 38.5 M | 81.1 M |

The fraction of reads retained at each point:

Table 6: Percentage of Reads Retained at Each Step

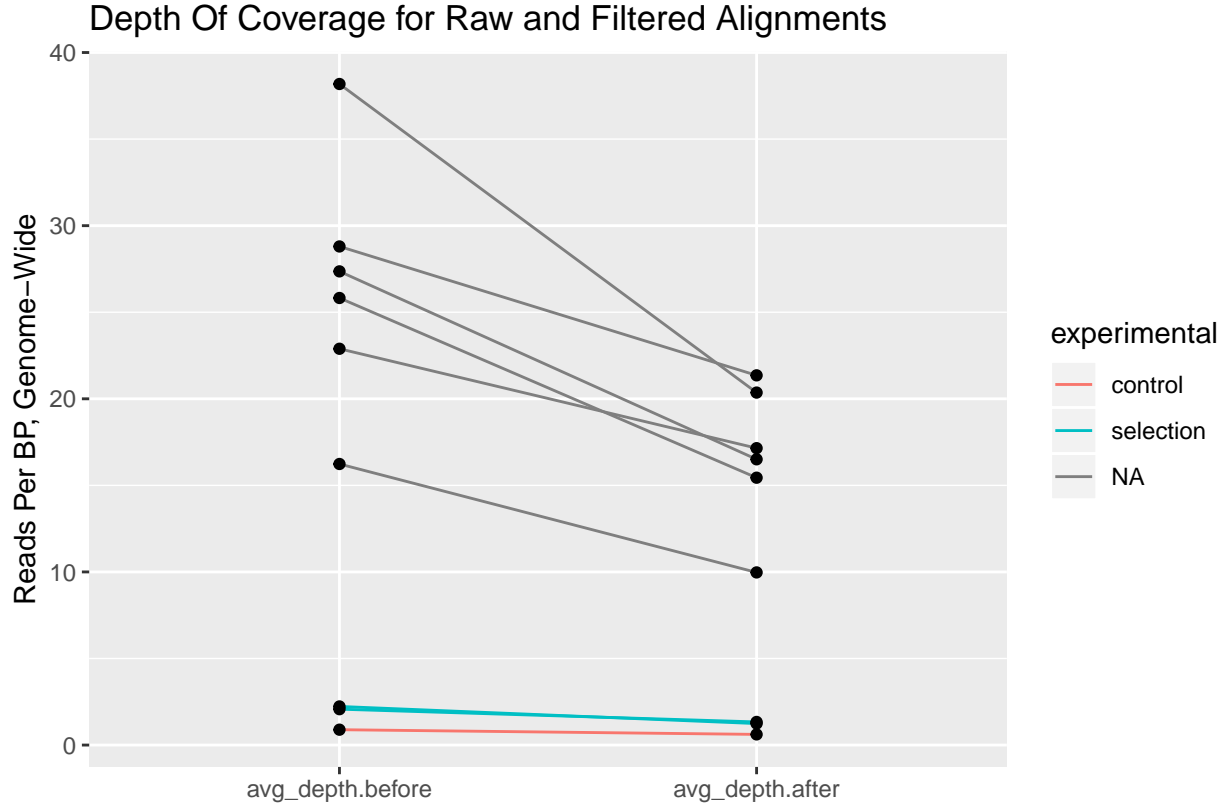
| measure | minimum | average | median | maximum |
|-------------------|---------|---------|--------|---------|
| filter_retention | 51.7 | 61.7 | 58.8 | 72.5 |
| mapping_retention | 50.2 | 73.3 | 70.2 | 92.1 |

2.3.2 Depth & Breadth of Coverage

Depth of coverage, ie, the genome-wide average number of mapped reads per base pair:

Table 7: Depth of Coverage Statistics for Raw and Filtered Alignments

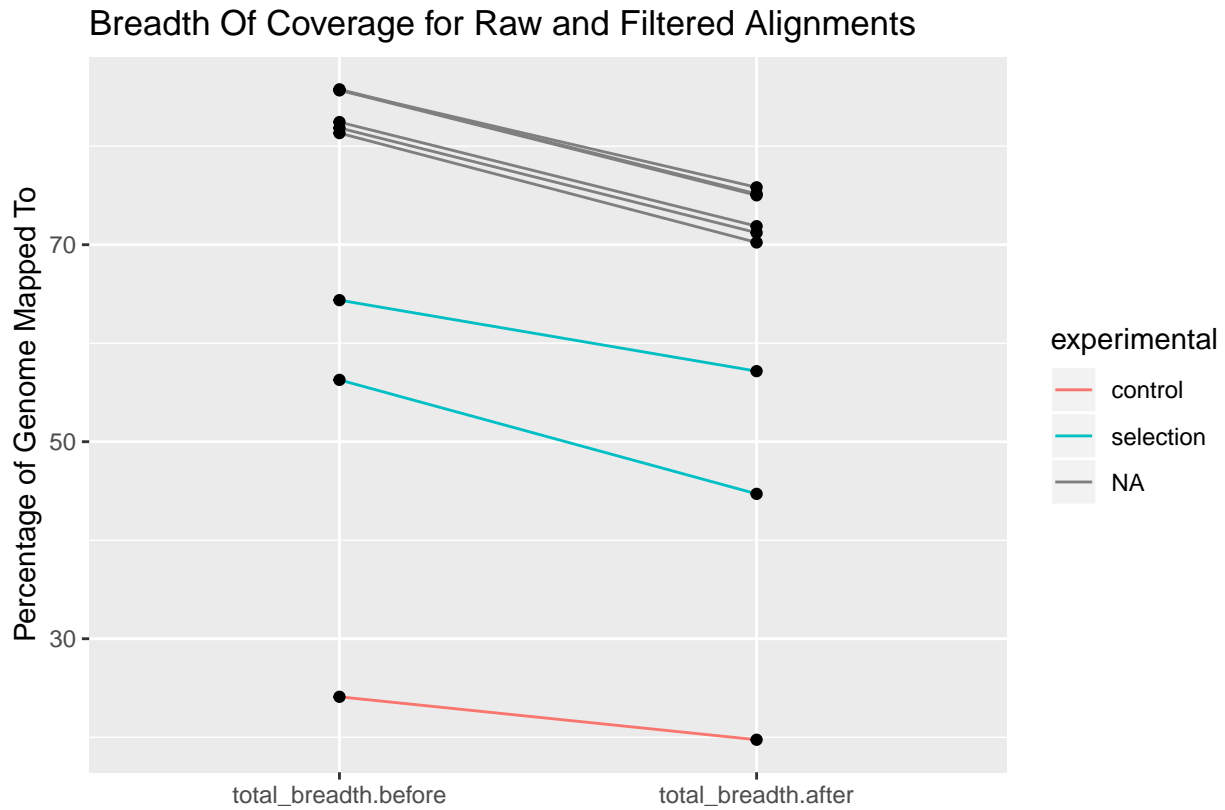
| step | minimum | average | median | maximum |
|-------------------------|---------|---------|--------|---------|
| pre-filtration depth | 889 m | 18.3 | 22.9 | 38.2 |
| post-filtration depth | 616 m | 11.6 | 15.4 | 21.4 |
| depth retention percent | 53.3 | 63.7 | 61.4 | 74.9 |



Breadth of coverage, ie, the percentage of the genome covered by at least one read:

Table 8: Breadth of Coverage Statistics for Raw and Filtered Alignments

| step | minimum | average | median | maximum |
|-------------------------|---------|---------|--------|---------|
| pre-filtration breadth | 24.1 | 71.9 | 81.9 | 85.8 |
| post-filtration breadth | 19.7 | 62.3 | 71.2 | 75.8 |
| breadth retention | 79.4 | 86.0 | 87.2 | 88.8 |



3 Results

4 References

4.1 Software

```
##
## To cite package 'tidyverse' in publications use:
##
##   Hadley Wickham (2017). tidyverse: Easily Install and Load the
##   'Tidyverse'. R package version 1.2.1.
##   https://CRAN.R-project.org/package=tidyverse
##
## A BibTeX entry for LaTeX users is
##
##   @Manual{,
##     title = {tidyverse: Easily Install and Load the 'Tidyverse'},
##     author = {Hadley Wickham},
##     year = {2017},
##     note = {R package version 1.2.1},
##     url = {https://CRAN.R-project.org/package=tidyverse},
##   }
##
## To cite the 'knitr' package in publications use:
##
```

```

## Yihui Xie (2018). knitr: A General-Purpose Package for Dynamic
## Report Generation in R. R package version 1.20.
##
## Yihui Xie (2015) Dynamic Documents with R and knitr. 2nd
## edition. Chapman and Hall/CRC. ISBN 978-1498716963
##
## Yihui Xie (2014) knitr: A Comprehensive Tool for Reproducible
## Research in R. In Victoria Stodden, Friedrich Leisch and Roger
## D. Peng, editors, Implementing Reproducible Computational
## Research. Chapman and Hall/CRC. ISBN 978-1466561595
##
## To see these entries in BibTeX format, use 'print(<citation>,
## bibtex=TRUE)', 'toBibtex(.)', or set
## 'options(citation.bibtex.max=999)'.

##
## To cite package 'yaml' in publications use:
##
## Jeremy Stephens, Kirill Simonov, Yihui Xie, Zhuoer Dong, Hadley
## Wickham, Jeffrey Horner, reikoch, Will Beasley, Brendan O'Connor
## and Gregory R. Warnes (2018). yaml: Methods to Convert R Data to
## YAML and Back. R package version 2.2.0.
## https://CRAN.R-project.org/package=yaml
##
## A BibTeX entry for LaTeX users is
##
## @Manual{,
##   title = {yaml: Methods to Convert R Data to YAML and Back},
##   author = {Jeremy Stephens and Kirill Simonov and Yihui Xie and Zhuoer Dong and Hadley Wickham and},
##   year = {2018},
##   note = {R package version 2.2.0},
##   url = {https://CRAN.R-project.org/package=yaml},
## }
##
## ATTENTION: This citation information has been auto-generated from
## the package DESCRIPTION file and may need manual editing, see
## 'help("citation")'.

```

Bibliography

Chen, Shifu, Yanqing Zhou, Yaru Chen, and Jia Gu. 2018. “Fastp: An ultra-fast all-in-one FASTQ preprocessor.” *Bioinformatics* 34 (17): i884–i890. doi:10.1093/bioinformatics/bty560.

Earley, Eric J., and Corbin D. Jones. 2011. “Next-generation mapping of complex traits with phenotype-based selection and introgression.” *Genetics* 189 (4): 1203–9. doi:10.1534/genetics.111.129445.