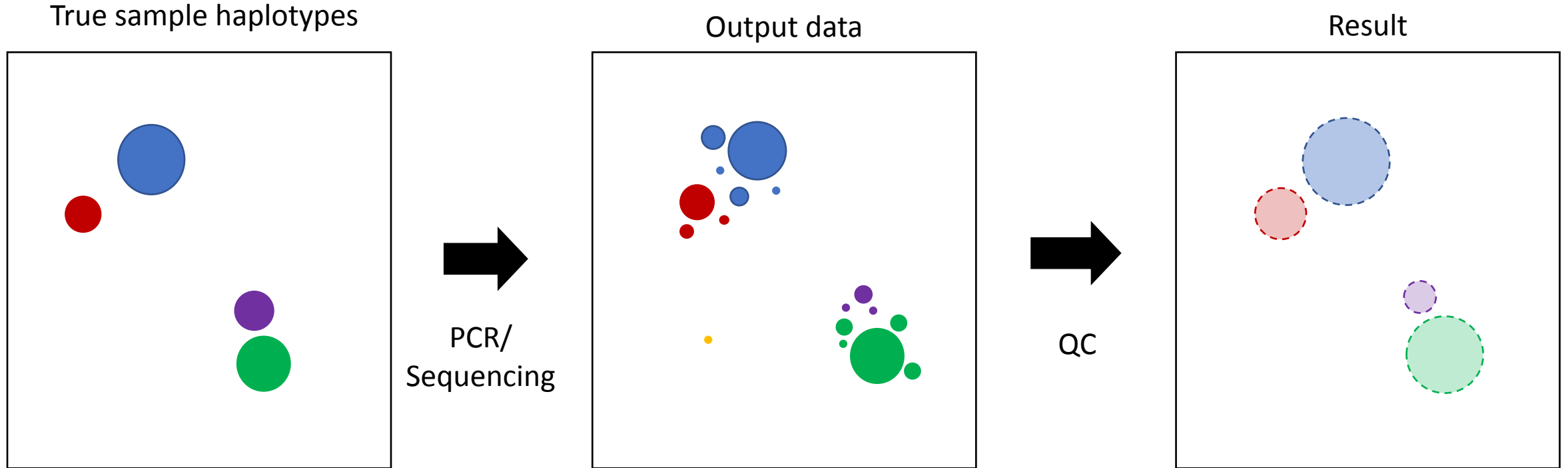


# **Amplicon sequencing pipeline validation**

2017-2018

# Amplicon sequencing analysis

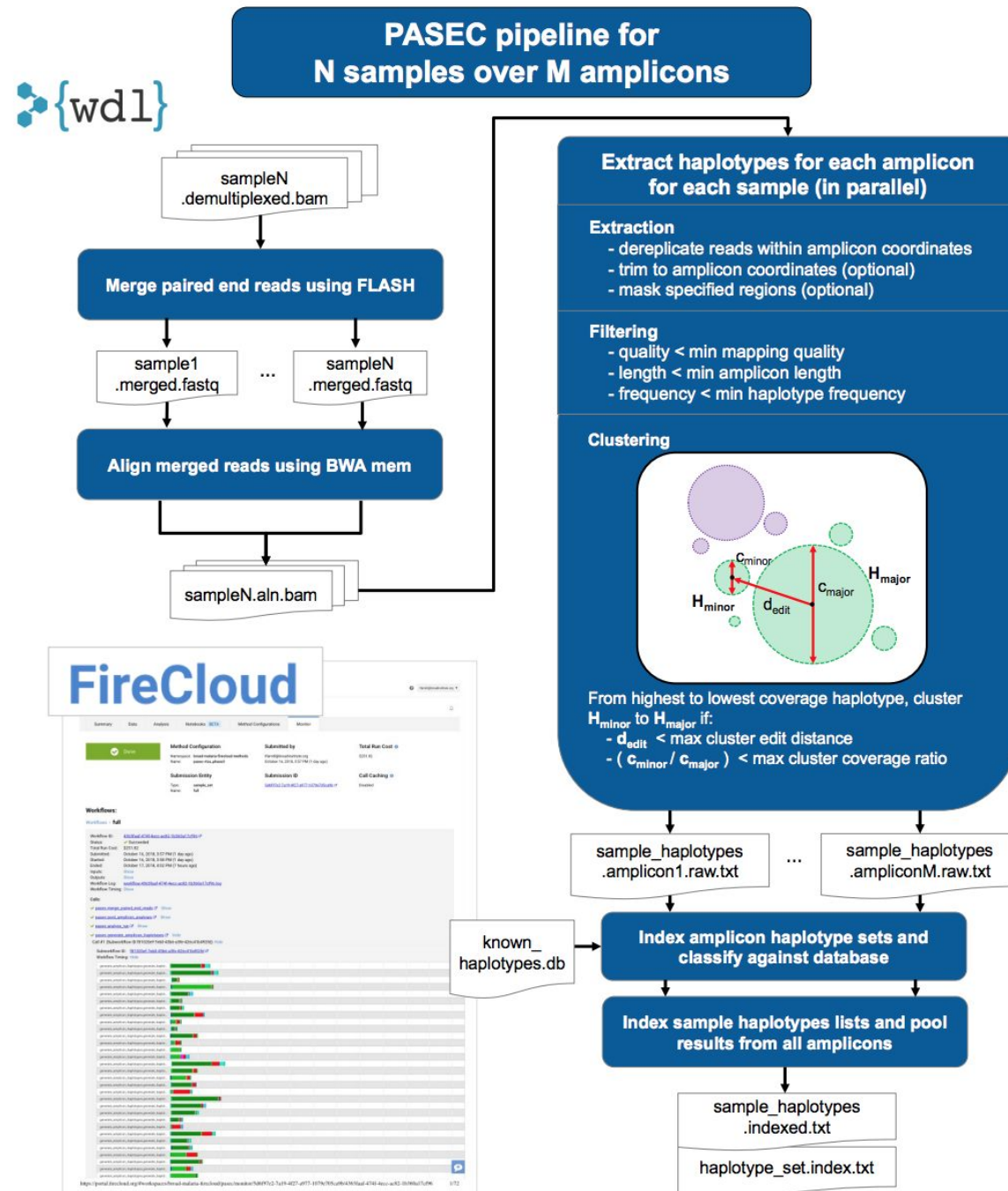


**Task:** How to eliminate technical variation without compromising biological variation?

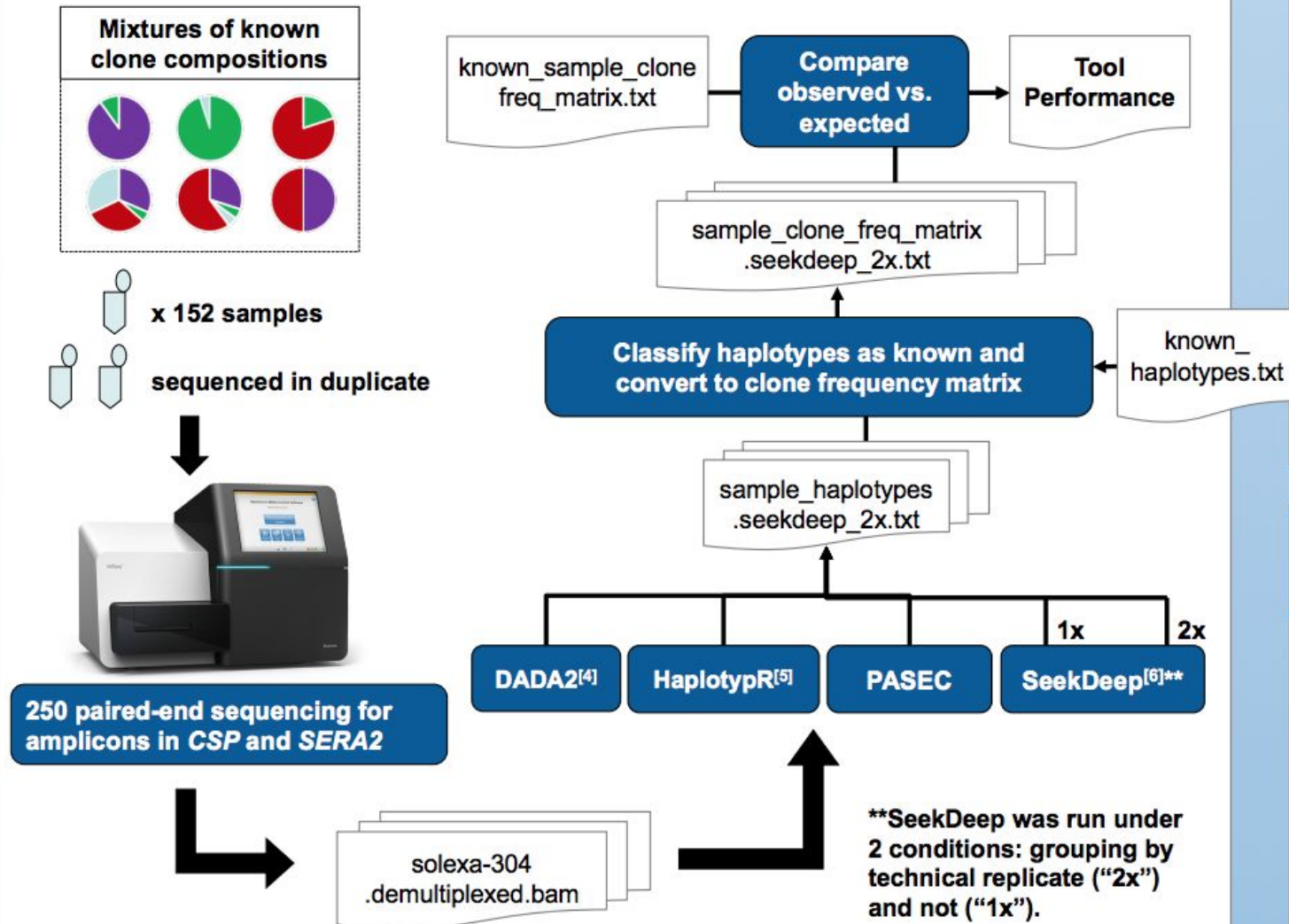
# Amplicon seq error correction tools

- PASEC (Early *et al*, **Malaria J** 2019 ; Neafsey *et al*, **NEJM** 2015)
  - Clusters based on distance and coverage
  - Manually mask difficult-to-sequence regions (homopolymers, etc)
- SeekDeep (Hathaway *et al*, **Bioinformatics** 2017)
  - Iteratively clusters based on weighted-distance, where weight is a function of difference type (mismatch or indel) and base quality
  - Derives power from duplicate PCRs
- DADA2 (Callahan *et al*, **Nature Methods** 2016)
  - Clusters based on error model prediction

$$p_A(i \rightarrow j) = \frac{1}{1 - p_{\text{pois}}(c_i \lambda_{ij}, 0)} \sum_{c'=c_j}^{\infty} p_{\text{pois}}(c_i \lambda_{ij}, c')$$

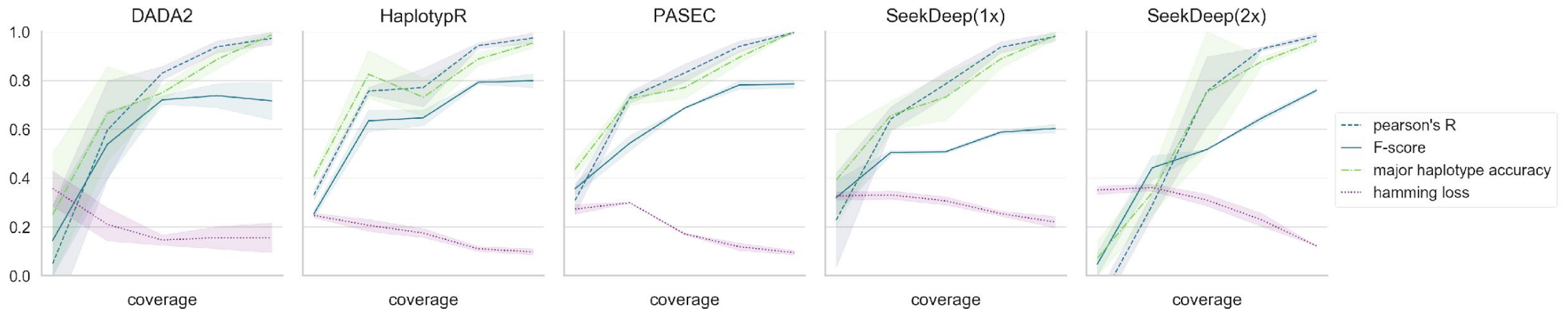
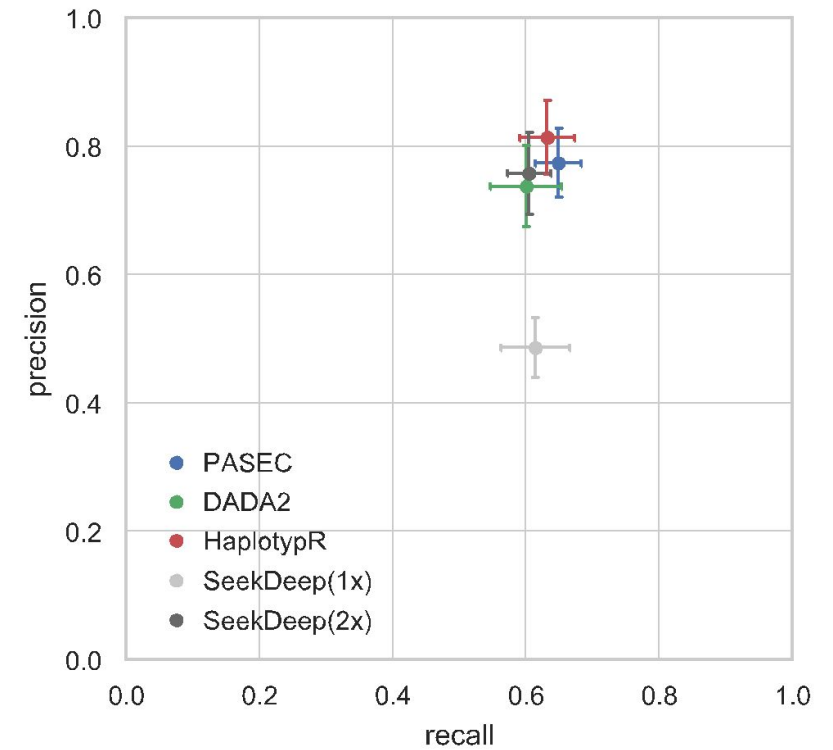


# Methods



# Tool performance comparison

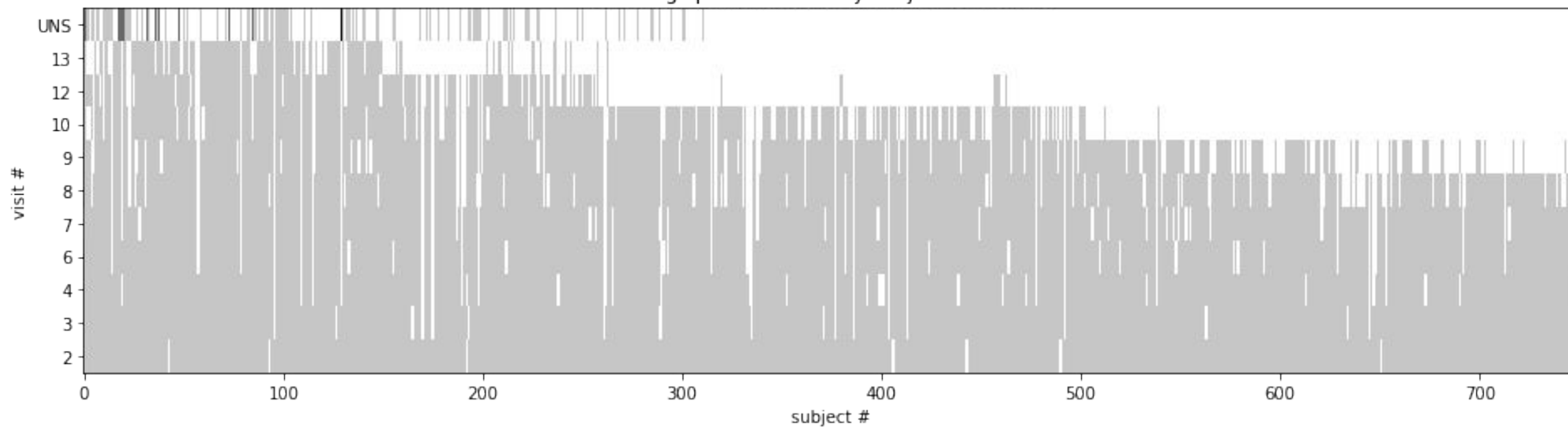
<https://www.biorxiv.org/content/early/2018/10/25/453472>.



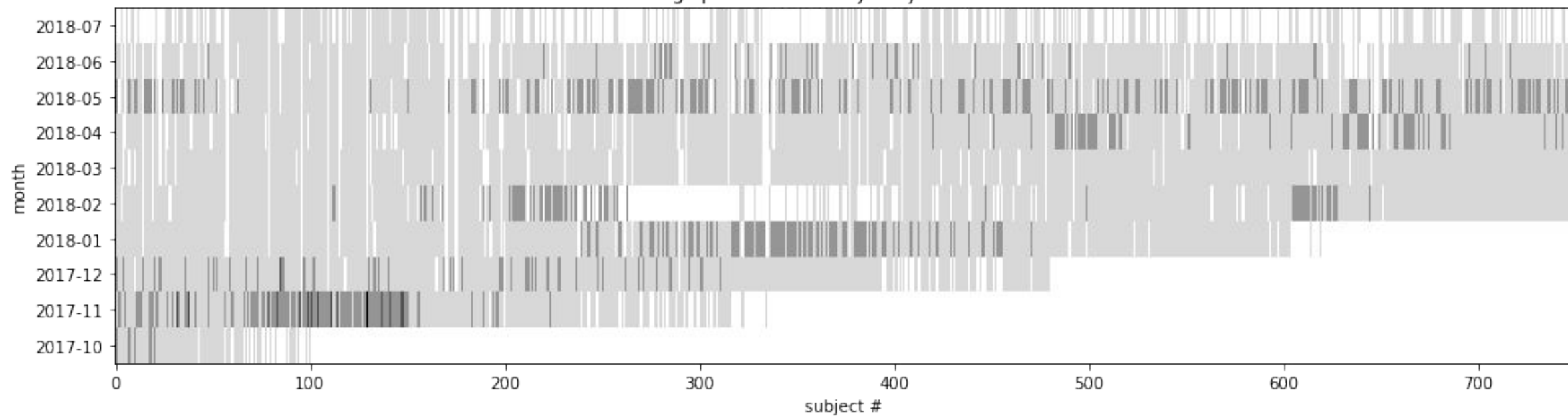
# **Phase IV anti-malarial vaccine clinical trial project**

2018-2019

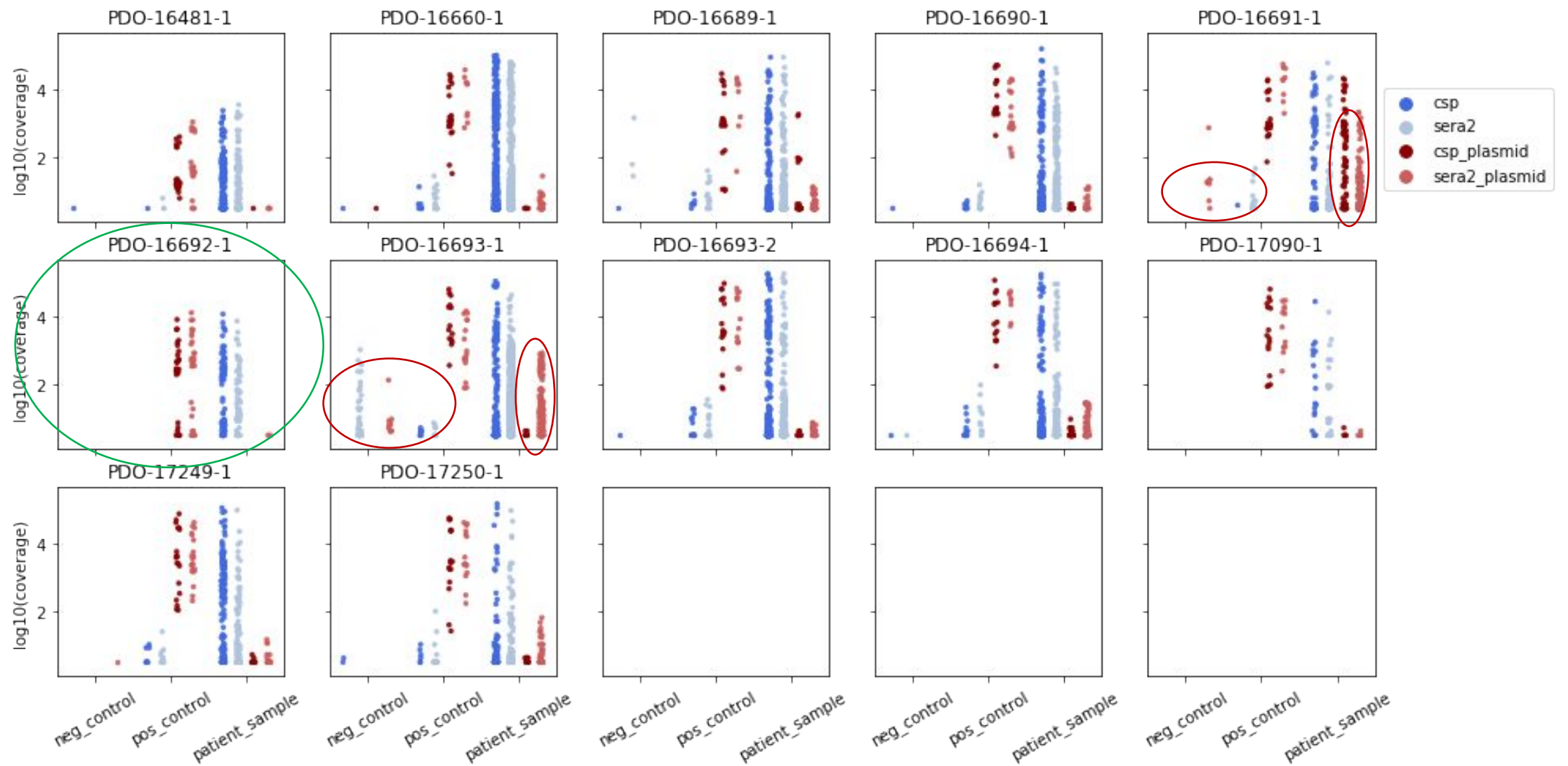
missing specimen data by subject and visit #



missing specimen data by subject # and month

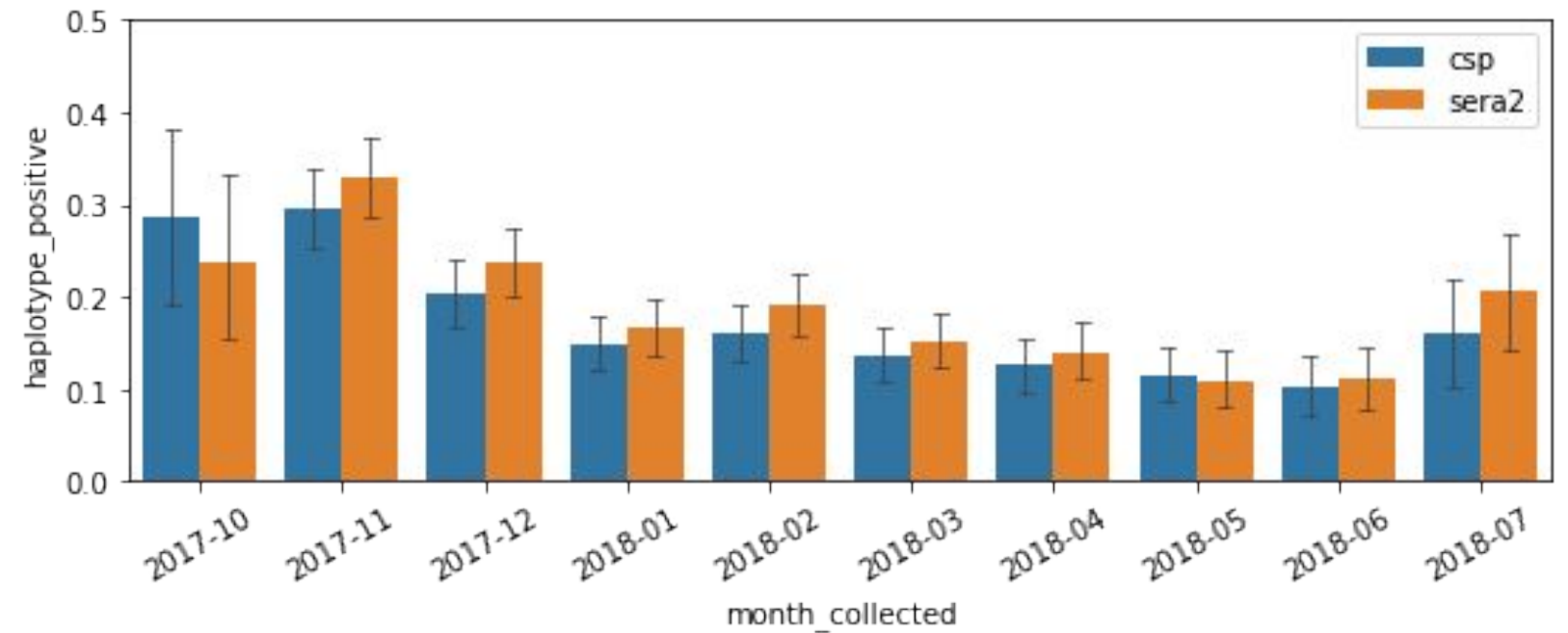
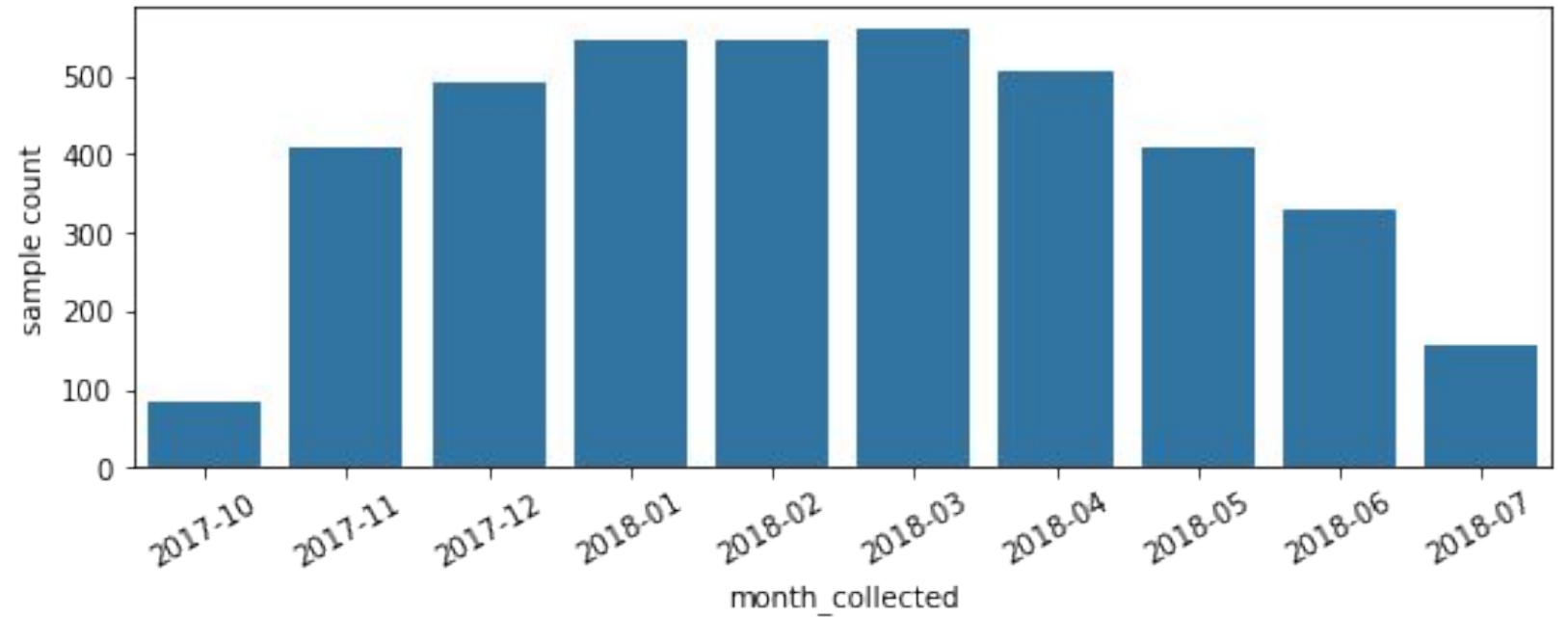






Amplicon (blue) and positive control plasmid (red)  $\log_{10}(\text{coverage})$  by PDO (each box, with PDO-version as title) and sample type (x-axis). Red circles indicate contamination (amplicons in controls or plasmids in negative controls or patient samples). The green circle indicates one of the cleanest (but lowest coverage) runs. One take away, when we see cross-contamination of the plasmids, there is usually also amplicon cross-contamination. Also of note, PDO-16693-2 (reworked PDO-16693) is a good quality PDO.

# Seasonality?



# How to define new infections?

