

Predicting genomic coverage by probabilistic binning and extrapolation

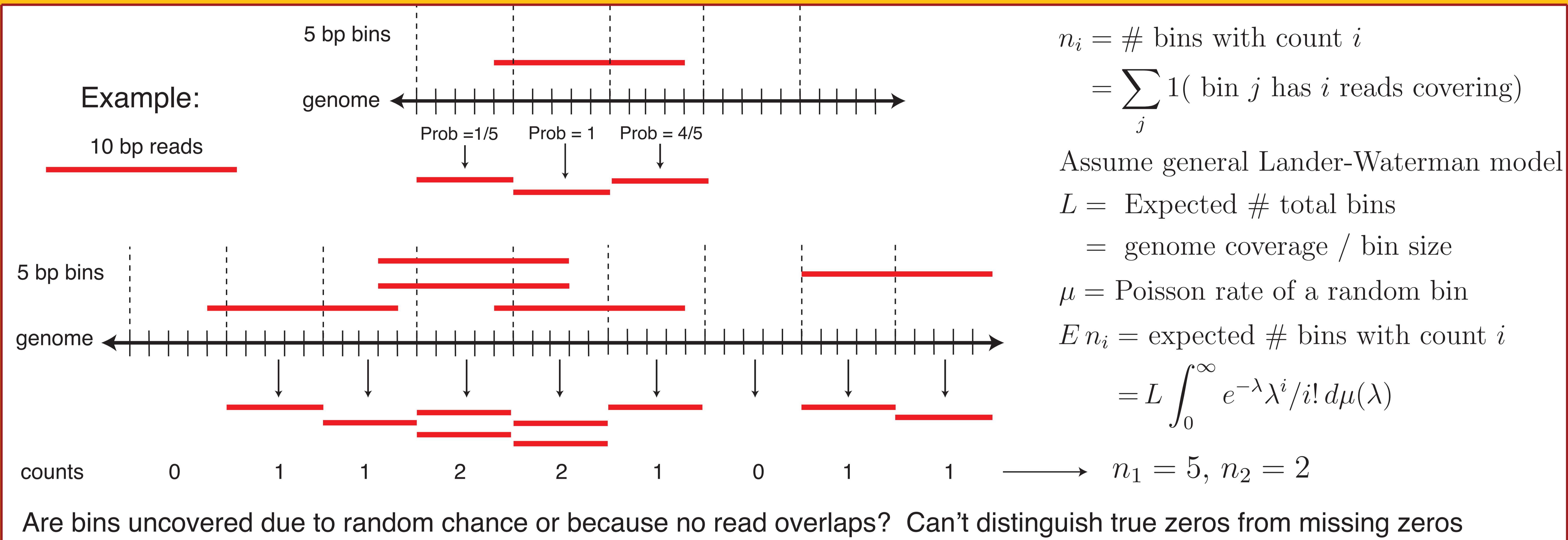


USC University of Southern California

Timothy Daley¹ and Andrew Smith²

¹Department of Mathematics and ²Department of Molecular and Computation Biology, University of Southern California

Probabilistic binning of reads:



Genomic coverage \approx bin size \cdot # covered bins:

Extrapolating coverage \longleftrightarrow Extrapolating covered bins

Daley & Smith *Nature Methods* (2011):

- Predicting number of distinct reads from additional sequencing.
- Treat distinct bins like distinct reads.
- Apply Good & Toulmin's non-parametric empirical Bayes solution (*Biometrika* 1959) and rational function approximation

$\Delta(t) = \#$ new covered bins from sequencing an additional tN reads

$$= E(\text{bins uncovered in additional experiment}) - E(\text{initial } \# \text{ uncovered bins})$$

$$= L \int_0^\infty e^{-\lambda} (1 - e^{-(t-1)\lambda}) d\mu(\lambda)$$

$$= \sum_{i=1}^\infty (-1)^{i+1} (t-1)^i E n_i \approx \frac{p_0 + p_1(t-1) + \dots + p_P(t-1)^P}{1 + q_1(t-1) + \dots + q_Q(t-1)^Q}$$

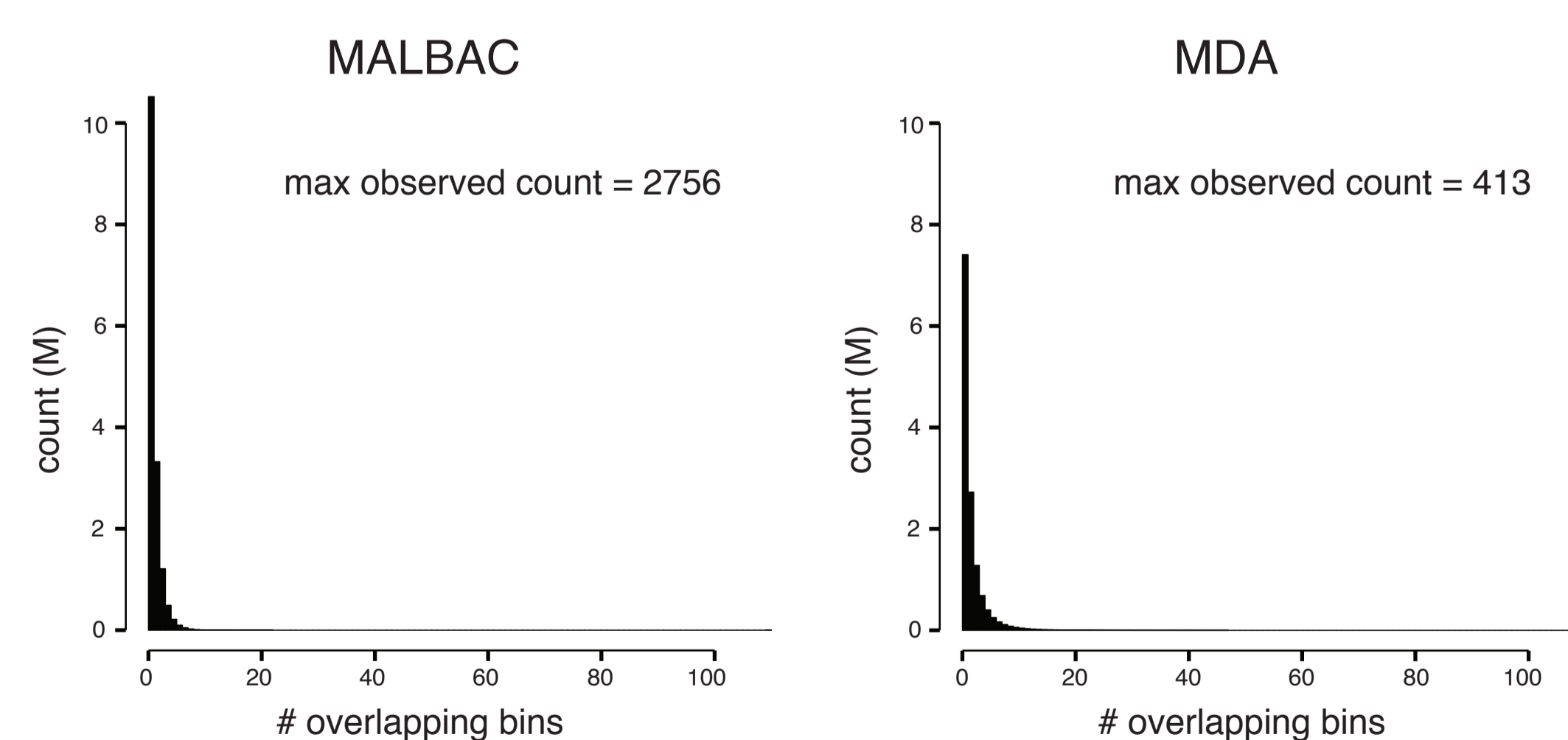
Results: extrapolating using 2% of the data

Data: 2 Single Cell libraries comparing library construction methods (Zong *et al.* Science 2012). MALBAC (Zong *et al.* Science 2012) & MDA (Dean *et al.* Gen. Res. 2001)

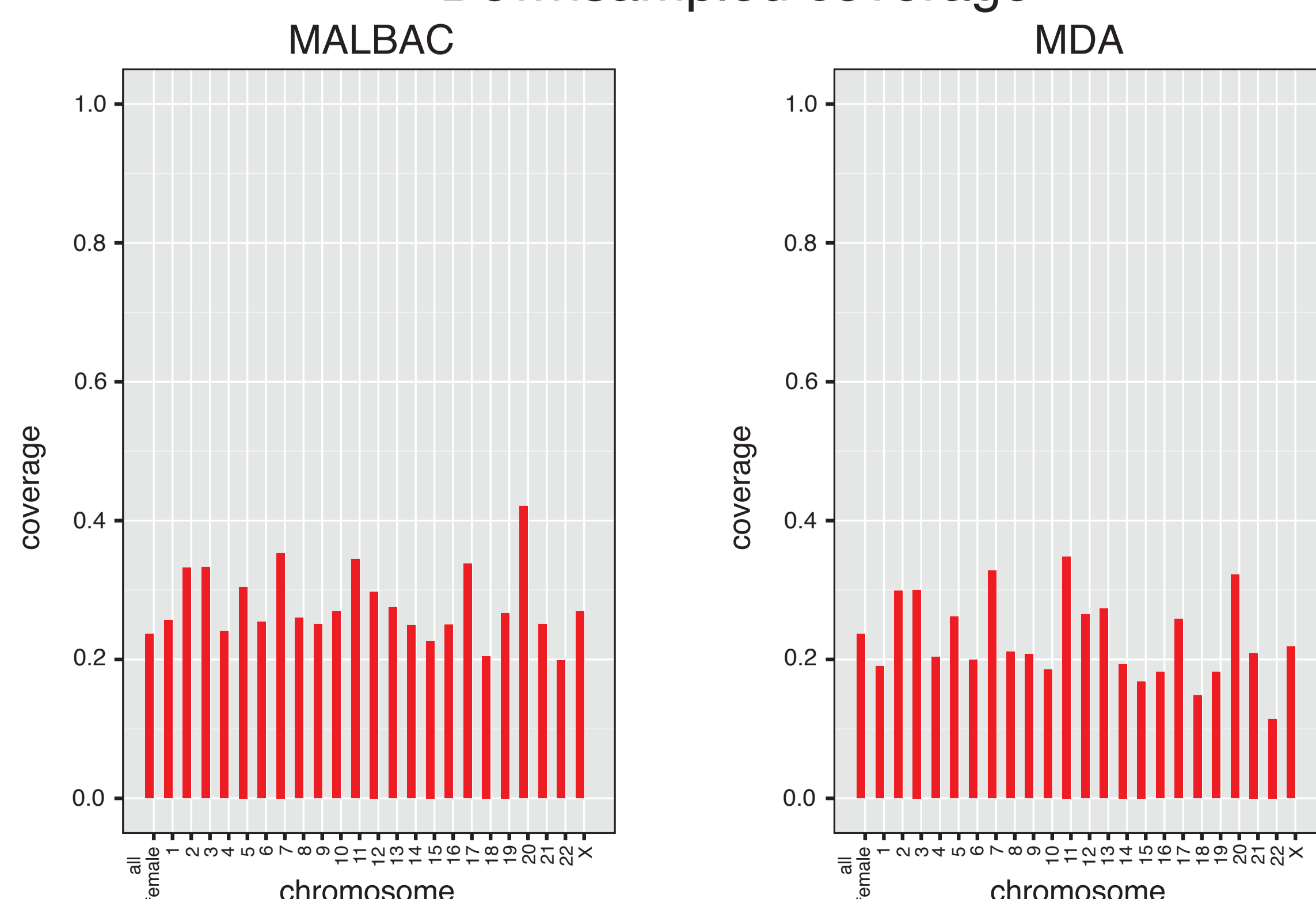
Library	MALBAC	MDA
SRA	SRX202978	SRX204160
Total reads	563.7M	534.8M
Reads mapping to female chroms (bwa)	449.6M	482.1M

Downsample 2% of reads with Picard:

50bp covered bin count histograms



Downsampled coverage



Extrapolated vs. Observed coverage

