

SUPPLEMENTAL MATERIAL: ESTIMATING ERROR MODELS FOR WHOLE GENOME SEQUENCING USING MIXTURES OF DIRICHLET-MULTINOMIAL DISTRIBUTIONS

Steven H. Wu¹, Rachel S. Schwartz^{1,3}, David J. Winter¹, Donald F. Conrad⁴, Reed A. Cartwright^{1,2,*}

¹The Biodesign Institute, and ²School of Life Sciences, Arizona State University, Tempe, AZ, 85281, USA; ³Department of Biological Sciences, The University of Rhode Island, Kingston, RI 02881, USA; ⁴Department of Genetics, Department of Pathology and Immunology, Washington University School of Medicine, Saint Louis, MO, 63110, USA

*To whom correspondence should be addressed. Tel: +1 480-965-9949; Email: cartwright@asu.edu

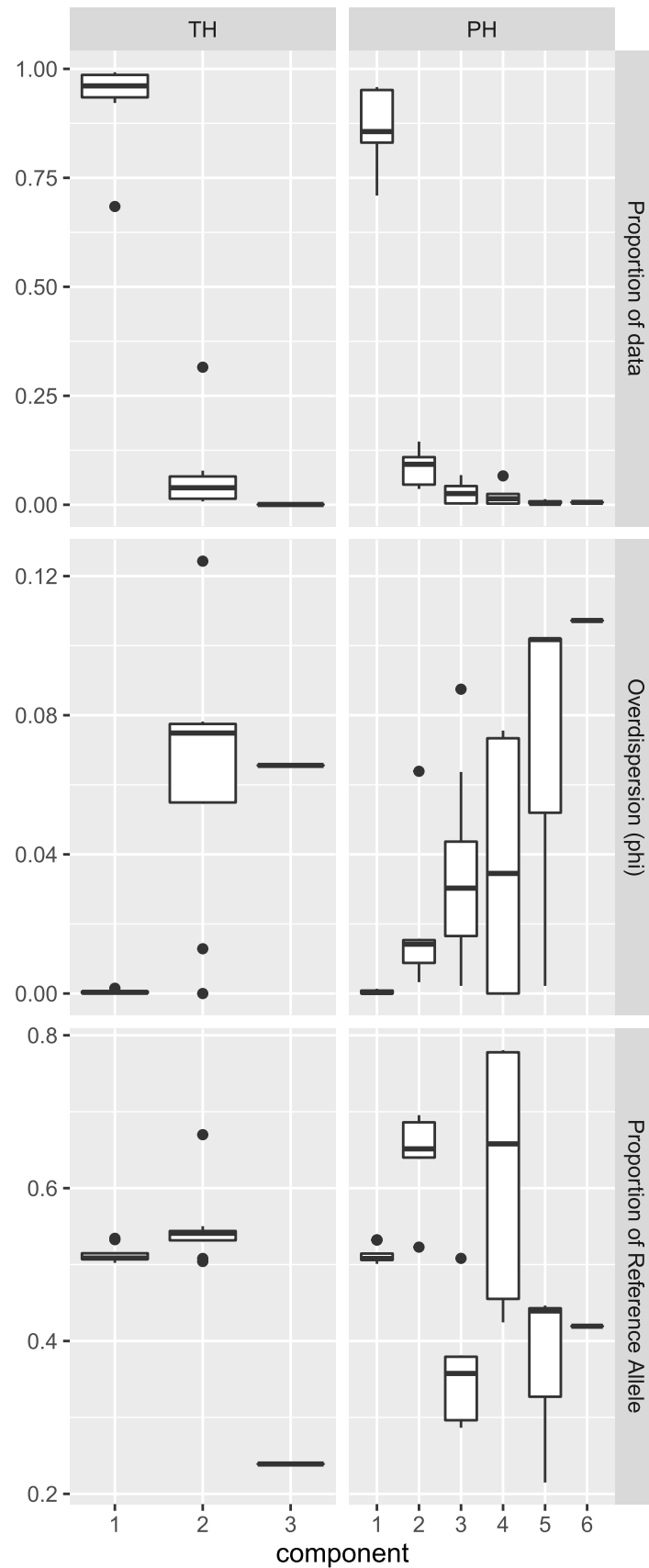


Figure S1: **Comparison of the parameter distributions between PH and TH datasets in the CEU data.** Boxplots summarize the proportion of data in each component, the overdispersion parameter for each component, and the proportion of reference allele in each component across different years and different chromosome regions. Most data is found in the major component of the model, but this is more true of the TH data. The overdispersion of the data increases as the proportion of data in the component decreases, with the major component having almost no overdispersion. The major component contains approximately equal proportions of the reference and alternate allele, with other components having widely disparate proportions of each allele. Thus, the major component, comprising the majority of the data is similar to a binomial with no bias.

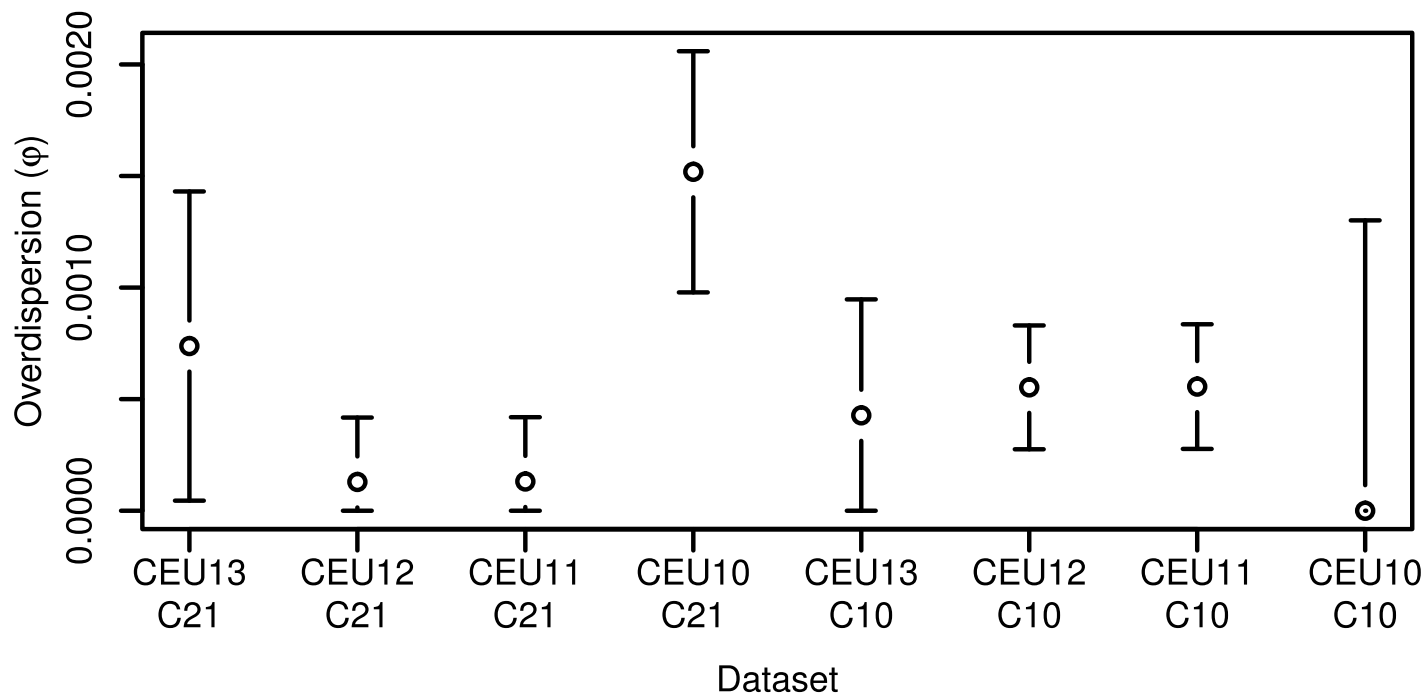


Figure S2: Confidence interval for φ from the major component for different TH datasets. Most of the estimated φ overlap in chromosome 21 and all of them in chromosome 10. Confidence intervals calculated from the EM algorithm.

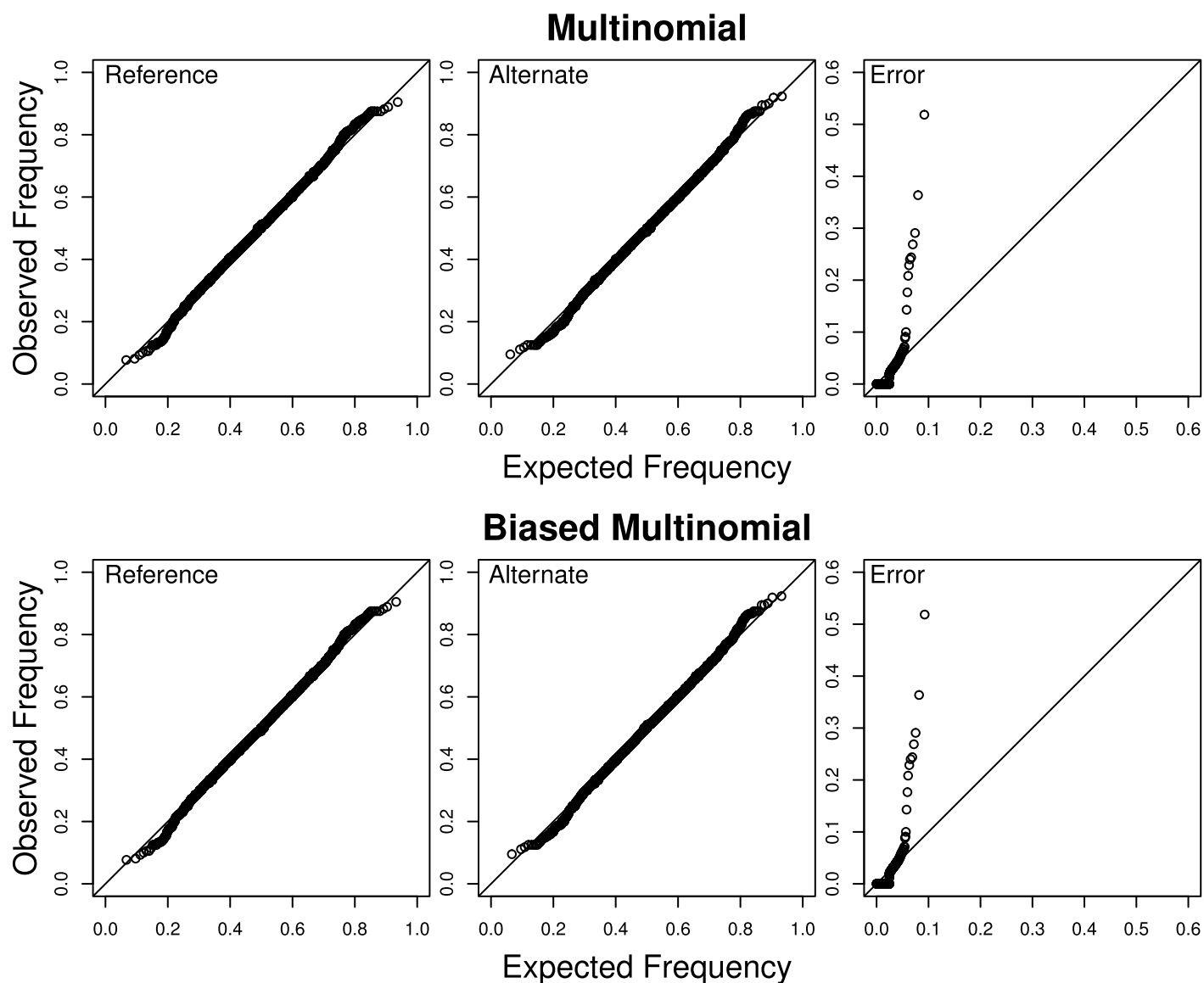
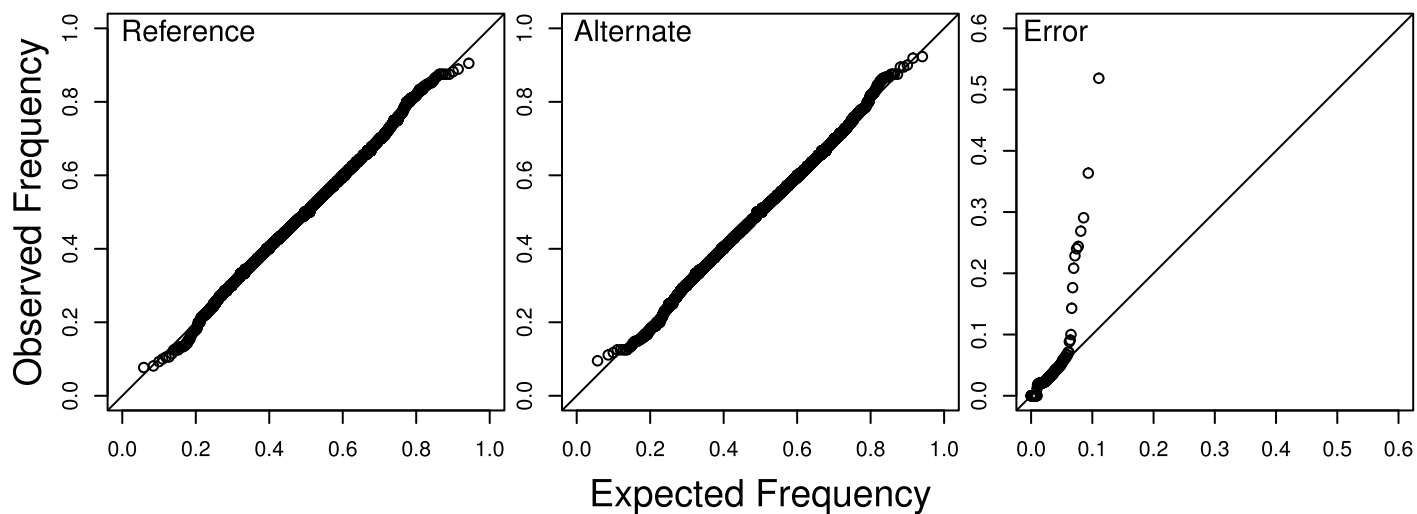
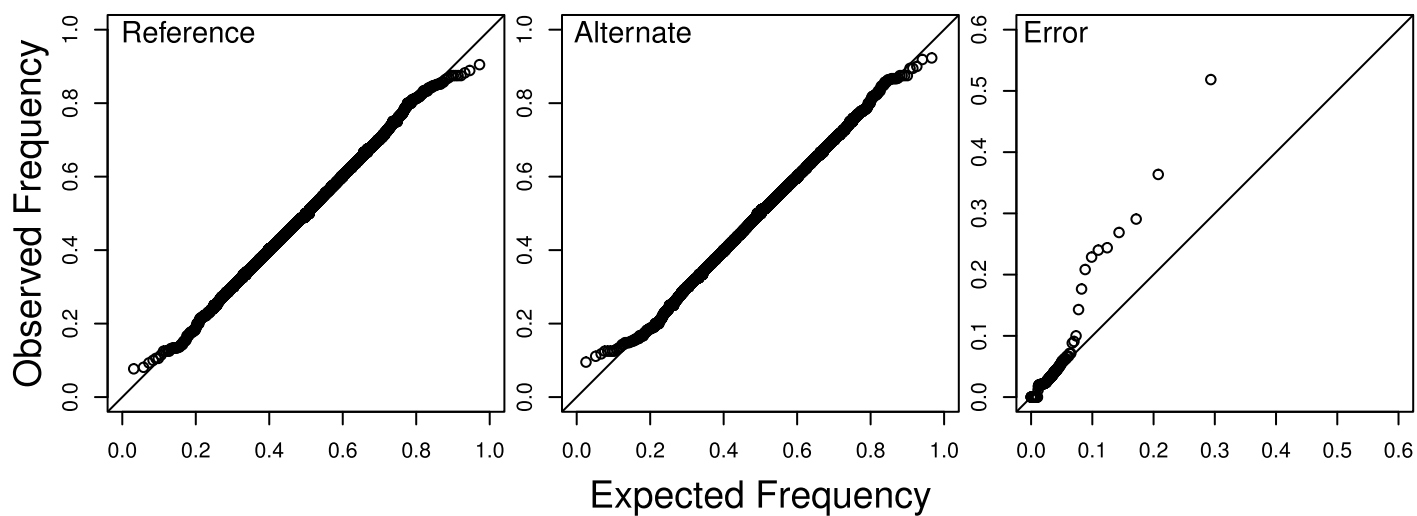


Figure S3: **Mixtures of Dirichlet-multinomials provide the best fits to genomic data: CEU13 Chr10 TH.** QQ plots evaluating the fits of several model to this dataset. The quantiles of the observed read count frequencies are calculated from the datasets, and the quantiles of the expected read count frequencies are estimated from the fitted model. A model that fits the data well produces points that fall along the diagonal.

Dirichlet-Multinomial



Mixture of 2 Dirichlet-Multinomials



Mixture of 3 Dirichlet-Multinomials

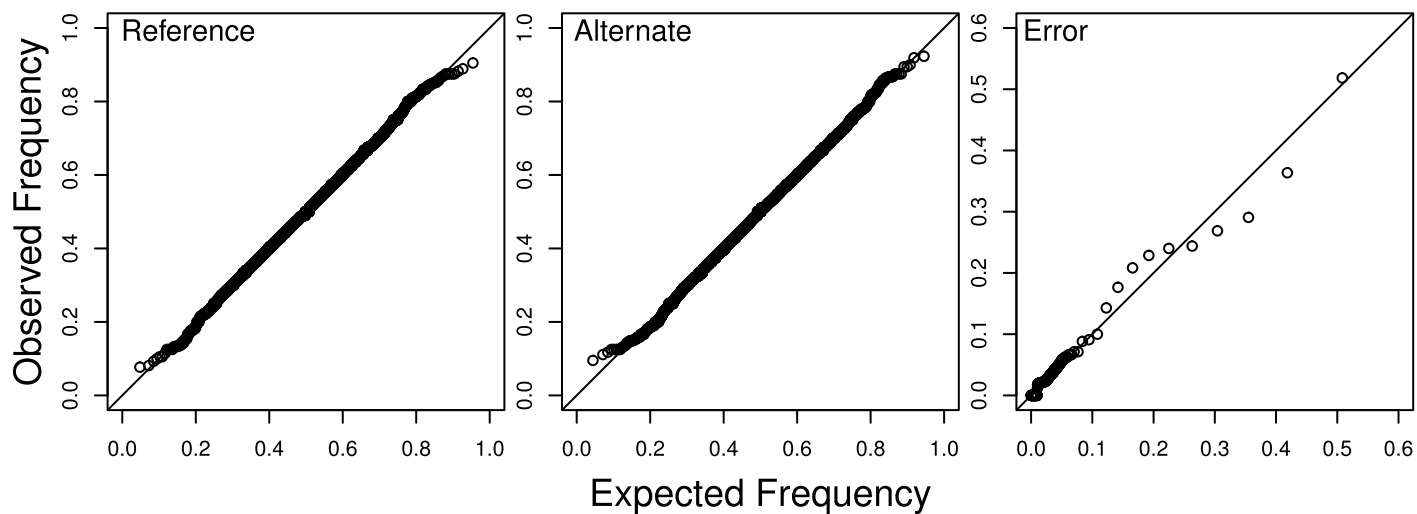
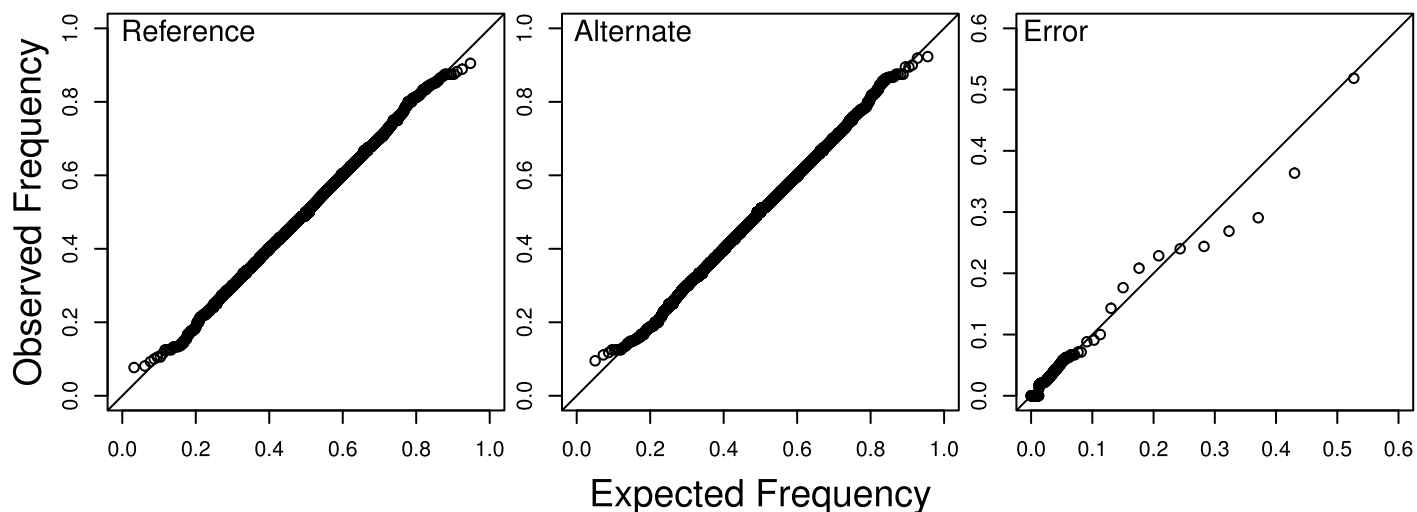
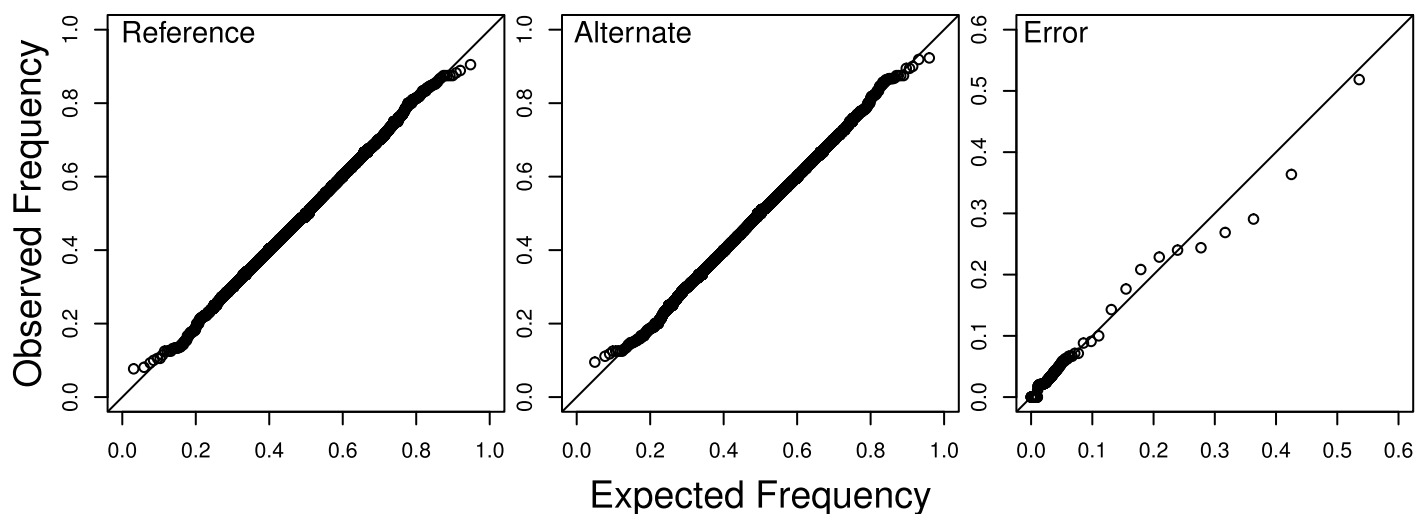


Figure S3 (Continued): **Mixtures of Dirichlet-multinomials provide the best fits to genomic data: CEU13 Chr10 TH.**

Mixture of 4 Dirichlet–Multinomials



Mixture of 5 Dirichlet–Multinomials



Mixture of 6 Dirichlet–Multinomials

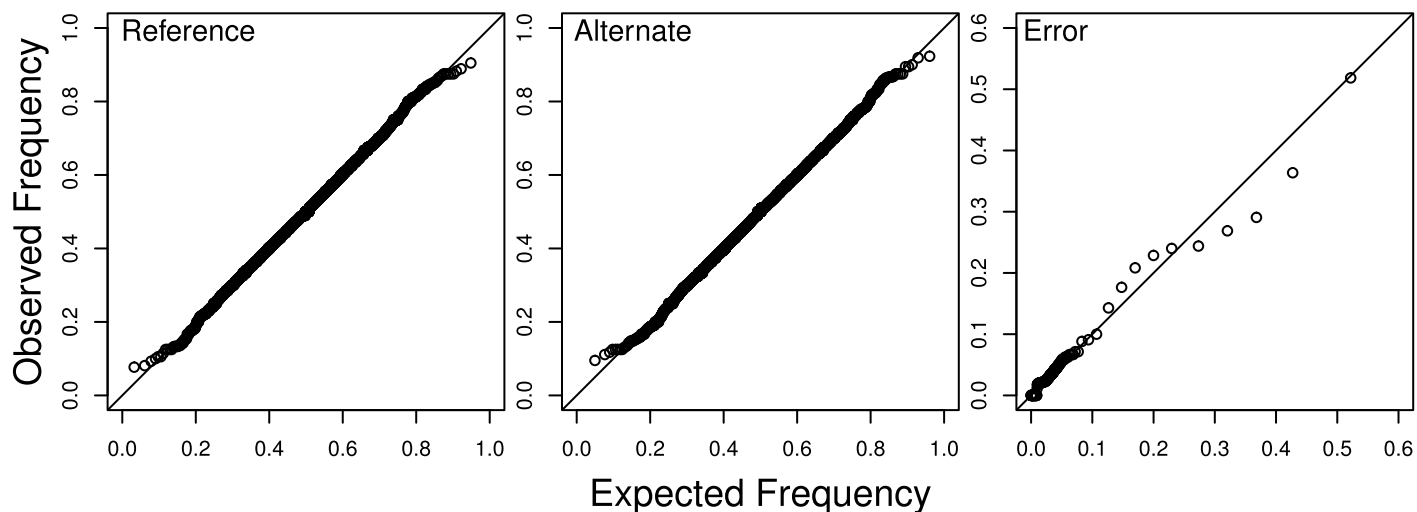
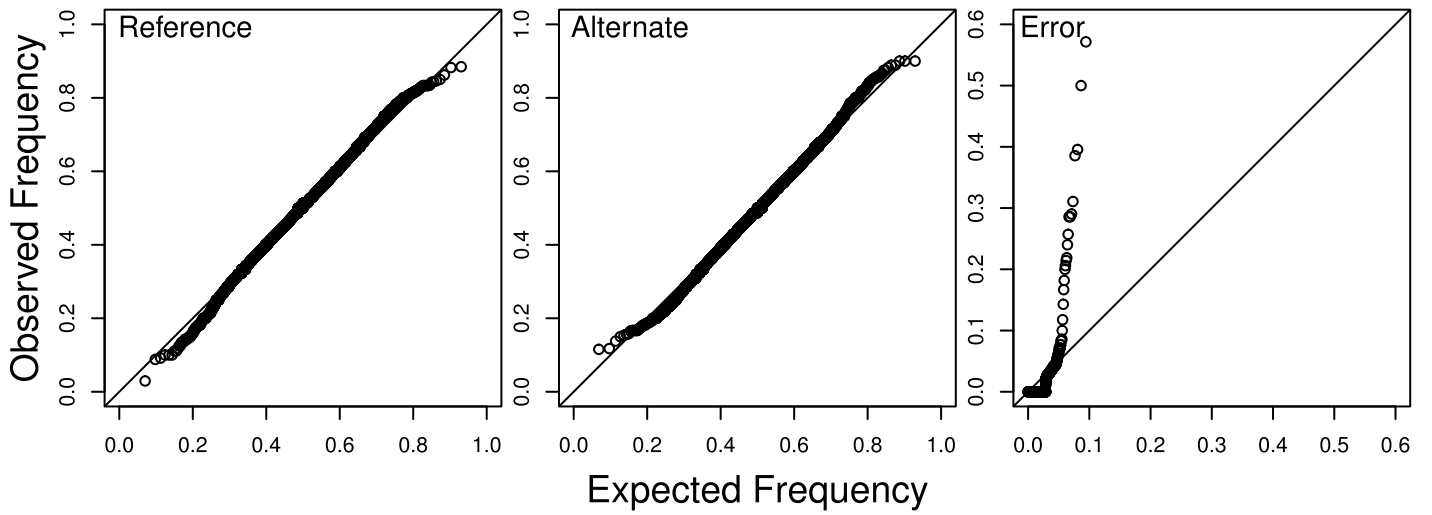
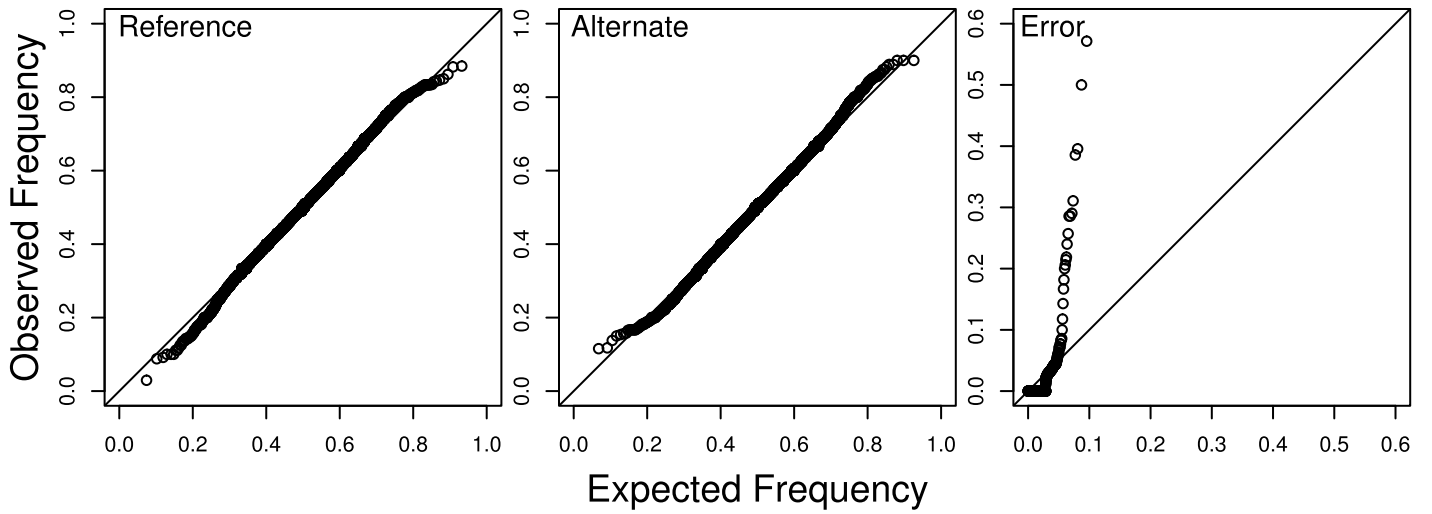


Figure S3 (Continued): **Mixtures of Dirichlet-multinomials provide the best fits to genomic data: CEU13 Chr10 TH.**

Multinomial



Biased Multinomial



Dirichlet-Multinomial

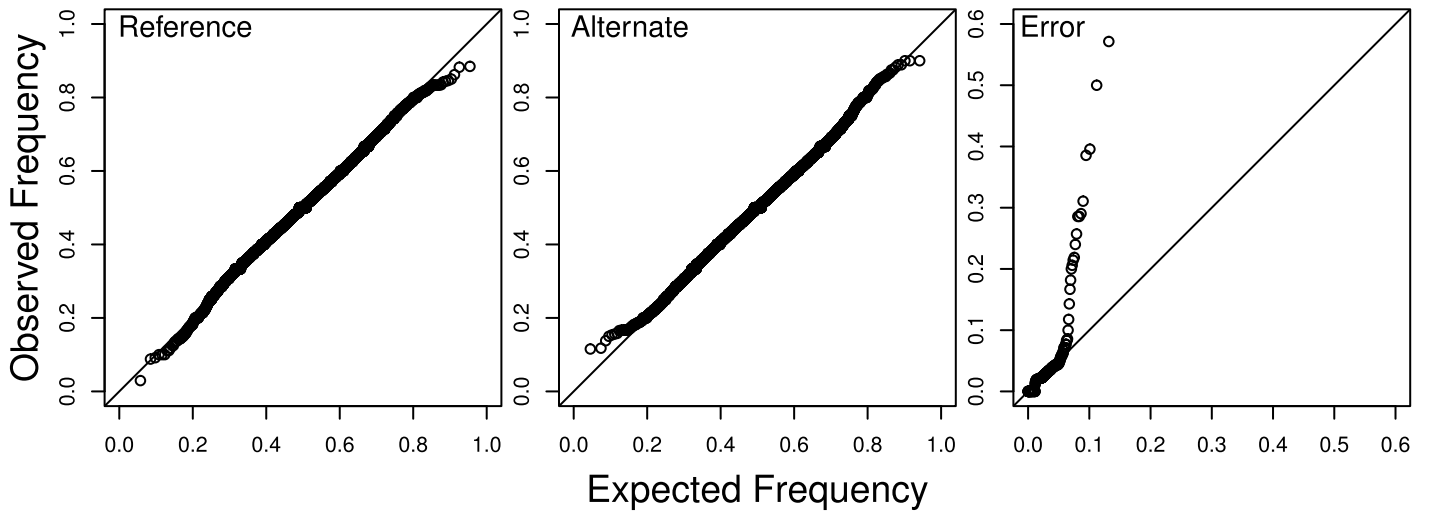
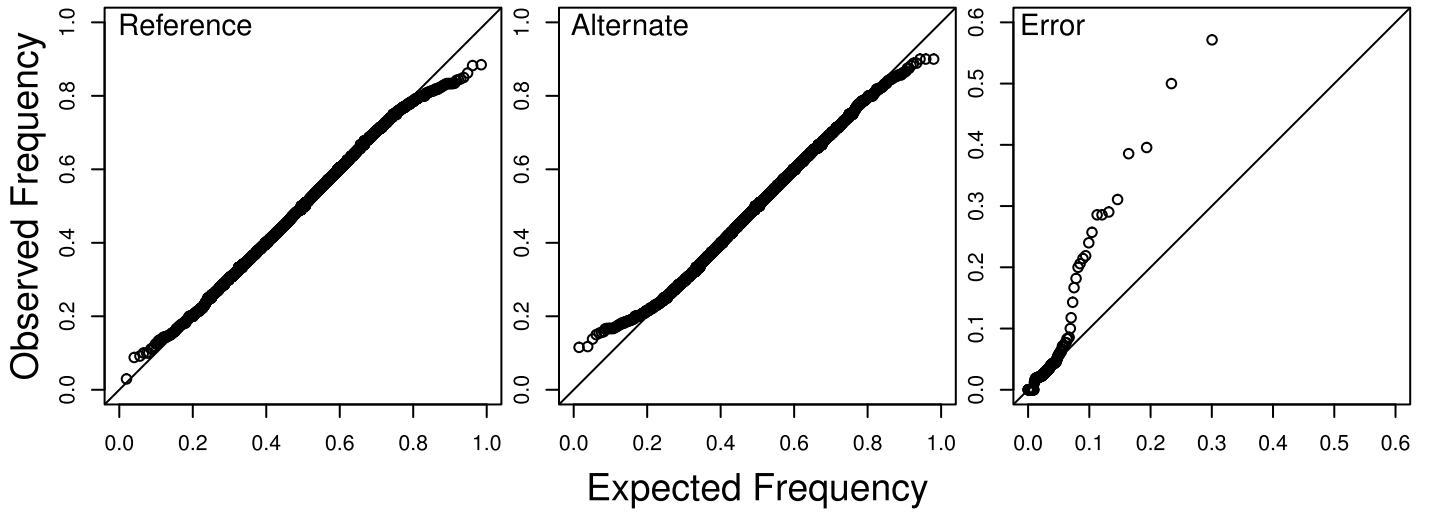
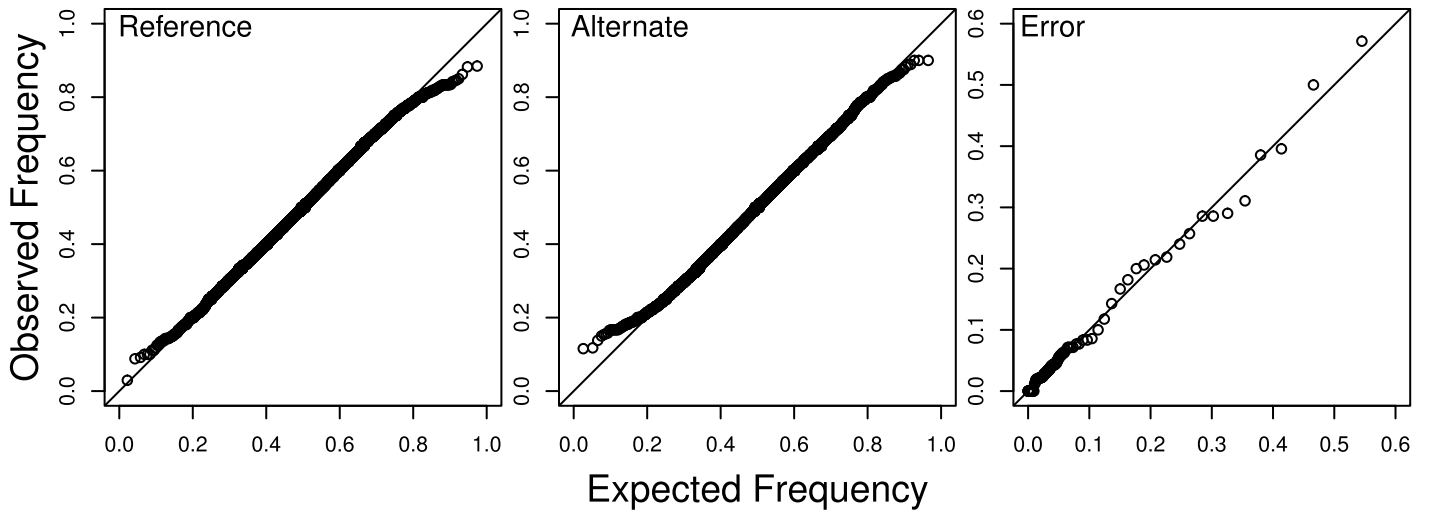


Figure S4: **Mixtures of Dirichlet-multinomials provide the best fits to genomic data: CEU13 Chr21 TH.**

Mixture of 2 Dirichlet-Multinomials



Mixture of 3 Dirichlet-Multinomials



Mixture of 4 Dirichlet-Multinomials

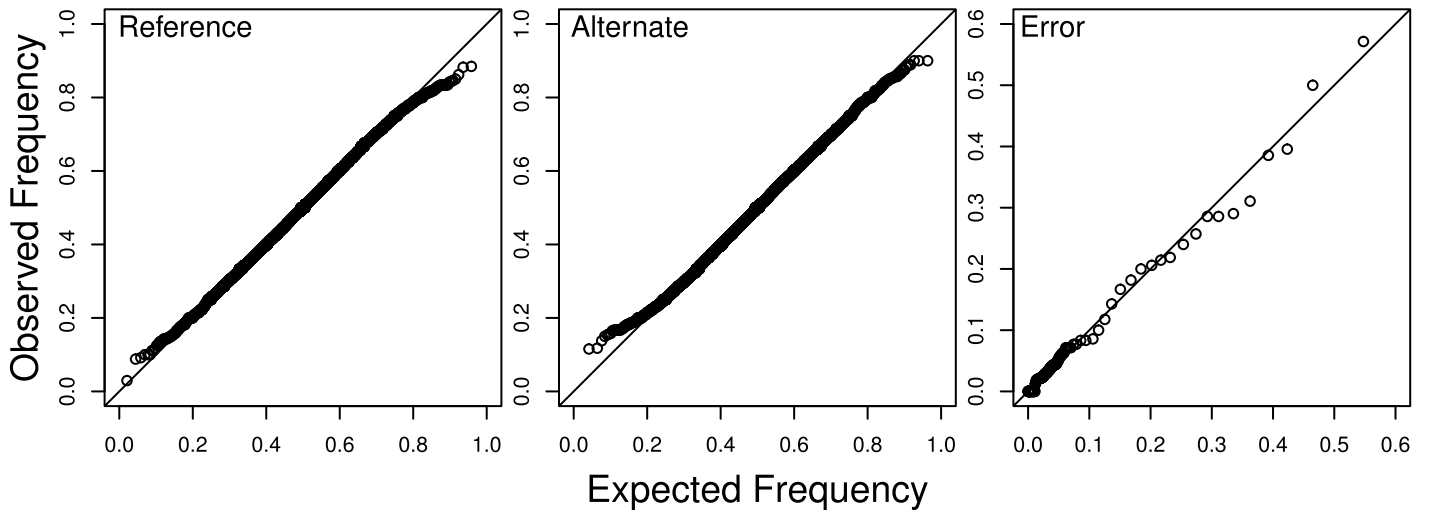
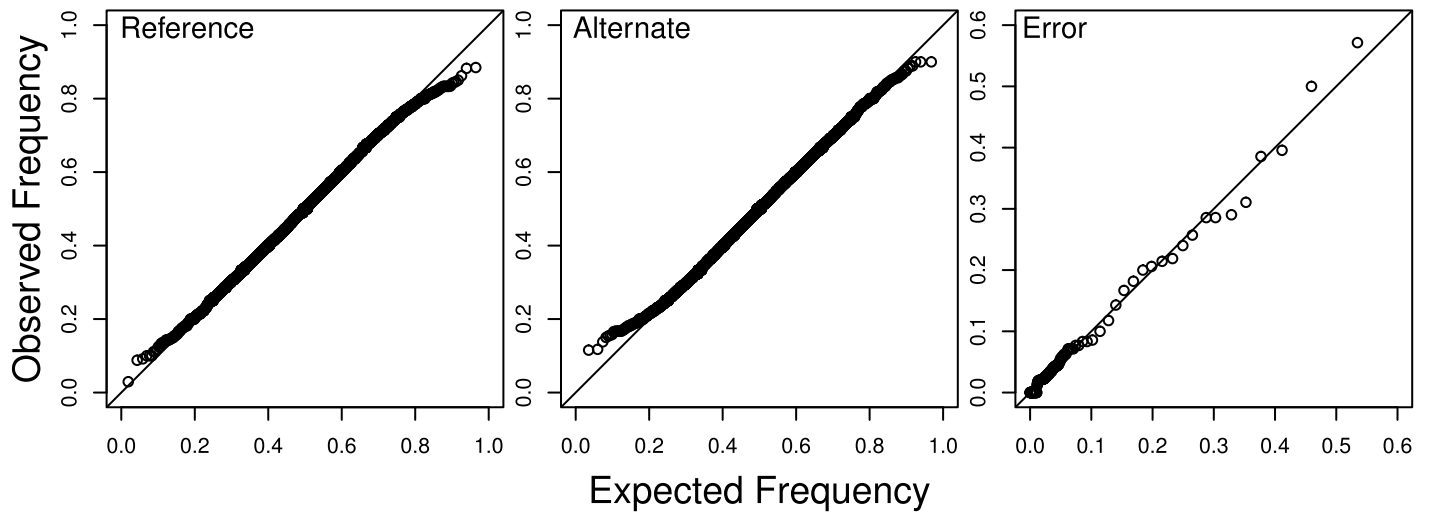


Figure S4 (Continued): **Mixtures of Dirichlet-multinomials provide the best fits to genomic data: CEU13 Chr21 TH.**

Mixture of 5 Dirichlet–Multinomials



Mixture of 6 Dirichlet–Multinomials

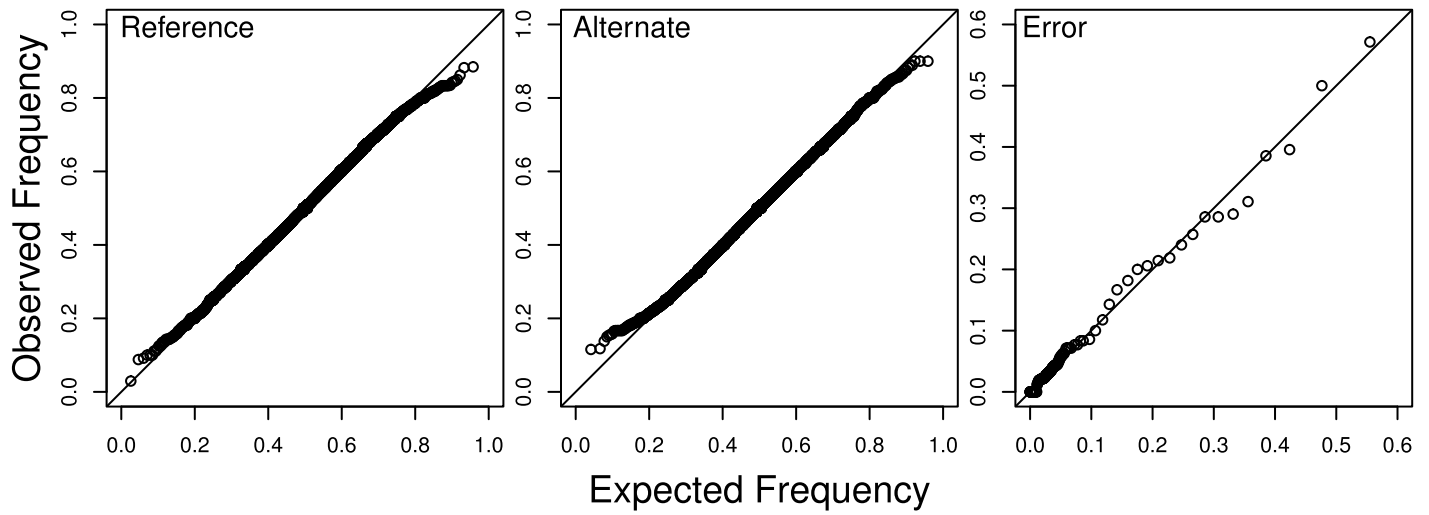
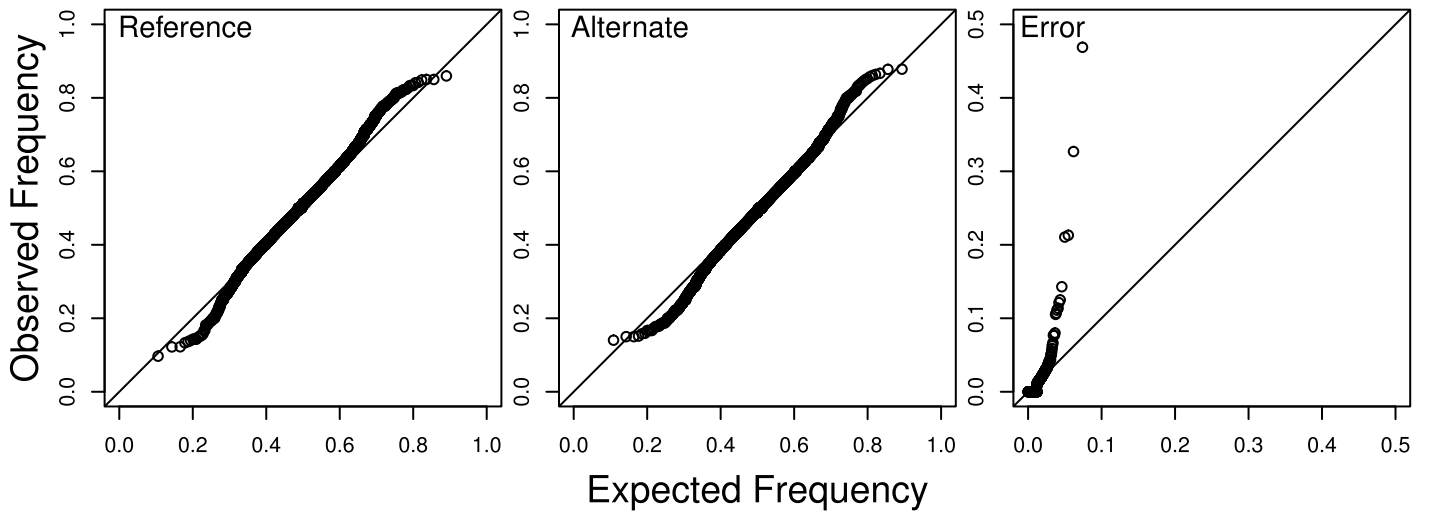
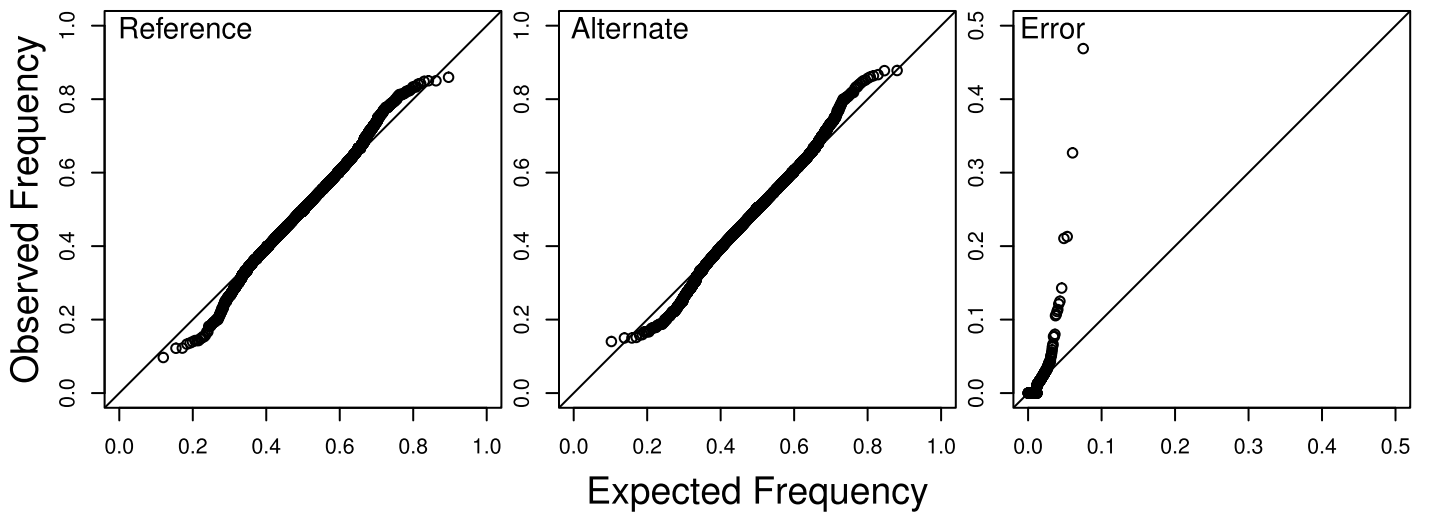


Figure S4 (Continued): **Mixtures of Dirichlet-multinomials provide the best fits to genomic data: CEU13 Chr21 TH.**

Multinomial



Biased Multinomial



Dirichlet-Multinomial

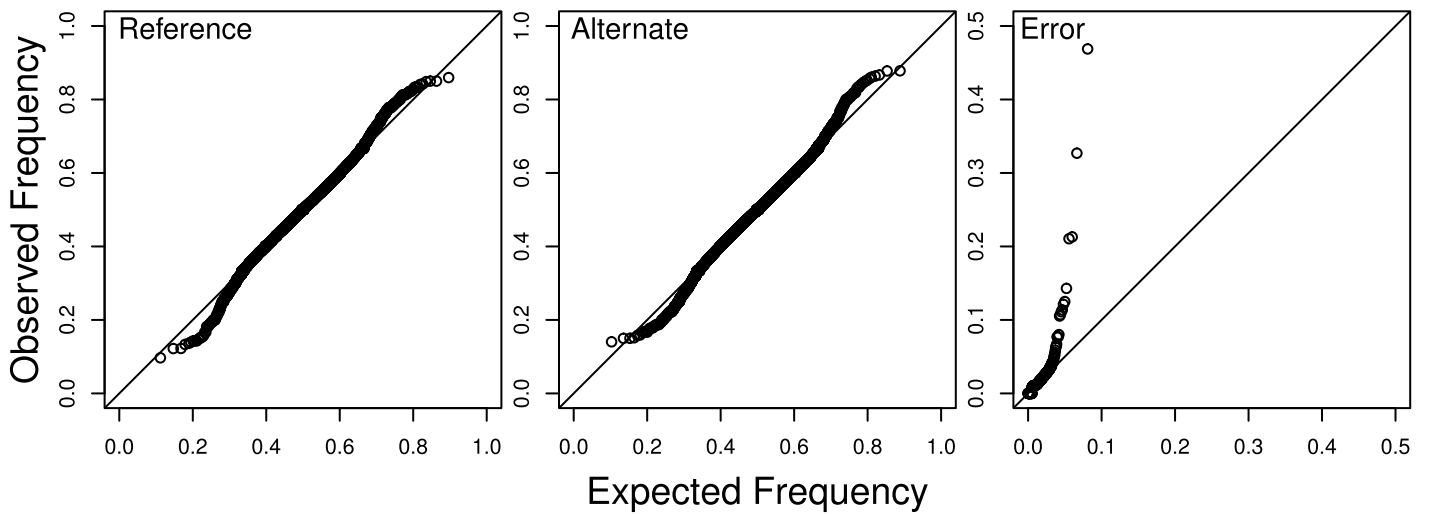
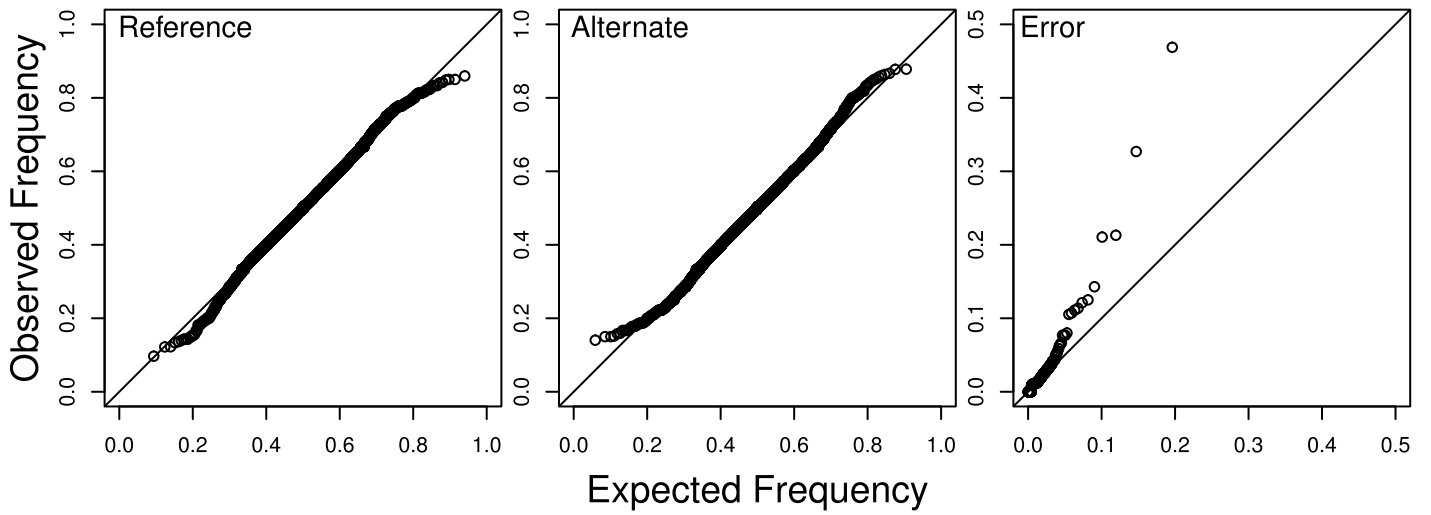
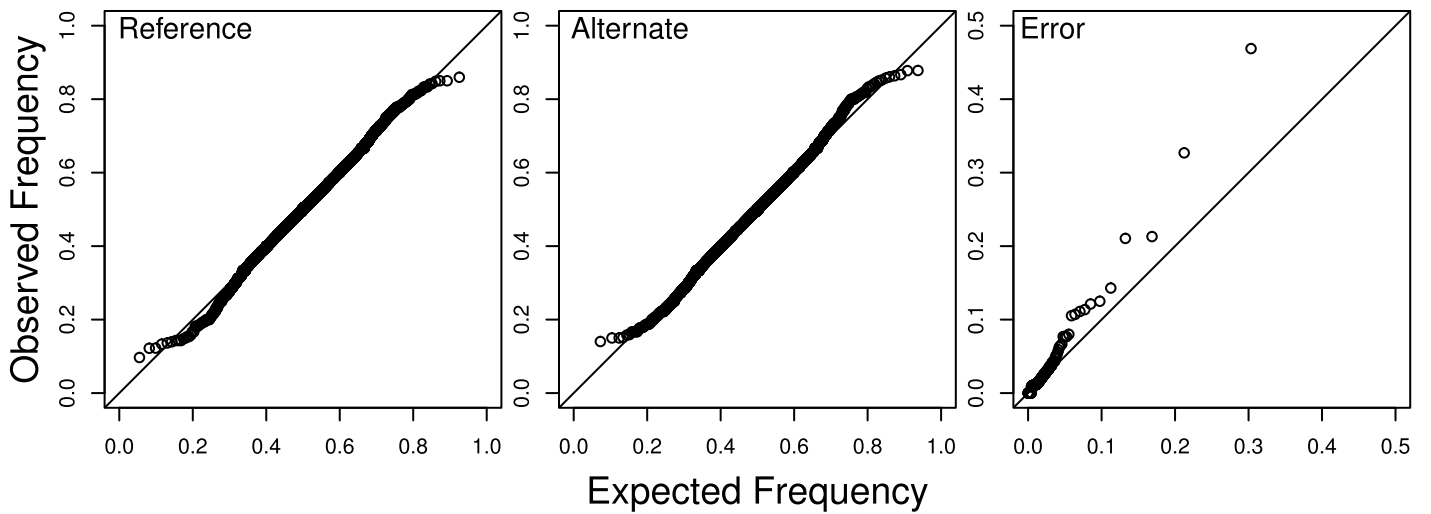


Figure S5: **Mixtures of Dirichlet-multinomials provide the best fits to genomic data: CEU12 Chr10 TH.**

Mixture of 2 Dirichlet–Multinomials



Mixture of 3 Dirichlet–Multinomials



Mixture of 4 Dirichlet–Multinomials

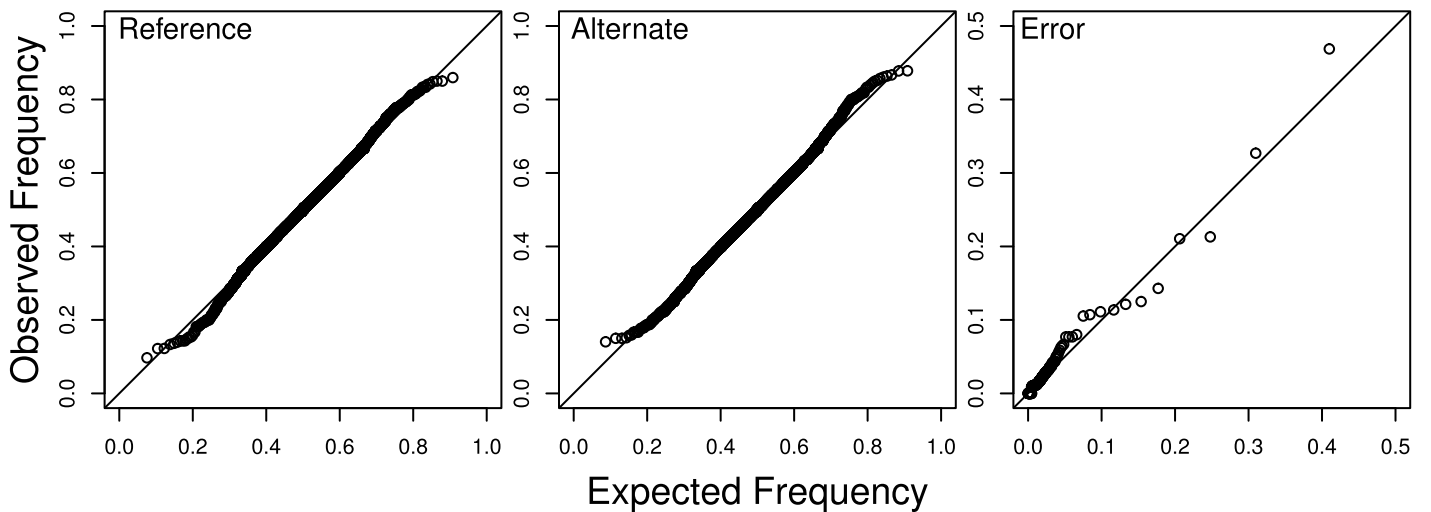


Figure S5 (Continued): **Mixtures of Dirichlet-multinomials provide the best fits to genomic data: CEU12 Chr10 TH.**