



Fig S11 (preceding page). Effect of Corr-ITSPCA iterations on the RoCA sector inference using binary synthetic data. The parameters used in the simulation study (see S2 Text for details) were $M = 500$ residues and $r = 5$ non-overlapping units S_i of size 12%, 10%, 8%, 6%, and 4% of M , respectively for $i = 1, \dots, 5$. The corresponding ℓ_i were set to equally spaced values between $\ell_1 = 6$ and $\ell_5 = 4$, and to model the phylogenetic effect, we set $\ell_0 = 8$. The ratio of the number of samples to the number of residues was fixed at $\frac{N}{M} = 4$ and the simulation was run for 500 Monte Carlo realizations. (A) The maximum percentage mismatch PM_{\max} between sectors formed using the PCs estimated at a particular Corr-ITSPCA iteration and the corresponding units. PM_{\max} decreases as the number of iterations increases, demonstrating that the iterative procedure in Corr-ITSPCA helps to accurately predict the true units. Here, the intermediate iteration corresponds to half of the total number of iterations Corr-ITSPCA took to converge in each Monte Carlo realization. (B-D) Illustration of the convergence of RoCA sectors to the corresponding units in a single Monte Carlo realization for the simulation setting of (A). Snapshot of (B) PCs representing the true units (here, the first PC represents phylogeny while the subsequent five PCs represent the five units), (C) PCs of the sample correlation matrix (constructed using the phylogeny-filtered MSA) used in forming PCA sectors, and (D) PCs at first, fourth, and eight (last) iteration of the Corr-ITSPCA method.