

SUPPORTING APPENDIX: A re-analysis of the Choe *et al* Affymetrix GeneChip control dataset

Alan R Dabney, John D Storey*

Address: Department of Biostatistics, University of Washington, Seattle, WA 98195, USA

Email: adabney@u.washington.edu, jstorey@u.washington.edu;

*Corresponding author

Simulation details. Let $\mathbf{X}_{i0} = (X_{i0A}, X_{i0B})^T$ be the RNA amounts for gene i , $i = 1, 2, \dots, m$, on a randomly-selected individual under conditions A and B. Suppose that \mathbf{X}_{i0} is distributed according to $N_2(\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$, the bivariate normal distribution with mean $\boldsymbol{\mu}_i = (\mu_{iA}, \mu_{iB})^T$ and covariance matrix

$$\boldsymbol{\Sigma}_i = \sigma_i^2 \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix}. \quad (1)$$

Equivalently, we can write

$$\mathbf{X}_{i0} = \boldsymbol{\mu}_i + \boldsymbol{\epsilon}_i, \quad (2)$$

where $\boldsymbol{\epsilon}_i \sim N_2(\mathbf{0}, \boldsymbol{\Sigma}_i)$, $i = 1, 2, \dots, m$. Let \mathbf{X}_i be the observed expression from a microarray. We write

$$\mathbf{X}_i = \boldsymbol{\mu}_i + \boldsymbol{\epsilon}_i + \boldsymbol{\phi}_i, \quad (3)$$

where $\boldsymbol{\phi}_i = (\phi_{iA}, \phi_{iB})^T \sim N_2(\mathbf{0}, \tau_i^2 \mathbf{I}_2)$ represents hybridization variability, $i = 1, 2, \dots, m$; by \mathbf{I}_2 , we mean the 2×2 identity matrix.

Suppose we form three technical replicates on a single individual. Then, the three observations can be written as

$$\mathbf{X}_{ijT} = \boldsymbol{\mu}_i + \boldsymbol{\epsilon}_i + \boldsymbol{\phi}_{ij}, \quad (4)$$

$i = 1, 2, \dots, m$, $j = 1, 2, 3$. If, on the other hand, we sample three independent individuals, we obtain

$$\mathbf{X}_{ijI} = \boldsymbol{\mu}_i + \boldsymbol{\epsilon}_{ij} + \boldsymbol{\phi}_{ij}, \quad (5)$$

$i = 1, 2, \dots, m$, $j = 1, 2, 3$. Note that, in terms of Figure 1 in Choe *et al*, ϵ represent variability introduced in columns 1-4, while ϕ represents the variability introduced between column 4 and 5.

The simulations summarized in Figures 2–4 were generated as follows. For $m = 15000$ genes, baseline means μ_{i0} were generated from the $N(0, 1)$ distribution. Then, for 6000 alternative genes, differential expression was created by sampling δ_i from the $N(0, 1)$ distribution. The means for each alternative gene were formed as $\mu_{iA} = \mu_{i0} + \delta_i$ and $\mu_{iB} = \mu_{i0}$. The means for the remaining 9000 null genes were formed as $\mu_{iA} = \mu_{iB} = \mu_{i0}$. The standard deviations σ_i and τ_i were taken to be 0.3 and 0.2, respectively. The correlation parameter ρ was chosen as 0.85. In simulation one, three technical replicates were formed. In simulation two, three biological replicates were formed. The above scenario was simulated 30 times.