

Base pair interactions and hybridization isotherms of matched and mismatched oligonucleotide probes on microarrays

Hans Binder*, Stephan Preibisch, Toralf Kirsten

Interdisciplinary Centre for Bioinformatics, University of Leipzig

* corresponding author: Interdisciplinary Centre for Bioinformatics of Leipzig University, D-4107 Leipzig, Haertelstr. 16-18, binder@izbi.uni-leipzig.de, fax: ++49-341-97-16679

Supplementary Material

S1: Single base contribution of fluorescence emission

The length of the RNA fragments, N_b^{RNA} , typically exceeds the length of the 25meric oligomer probes. Consequently also labelled bases which dangle outside of the probe/target duplex potentially contribute to the measured fluorescence intensity in addition to labels attached to the 25meric target region. Let us denote the number of bases outside of the respective 25meric duplex by N_b^{out} for a RNA fragment of total length $N_b^{\text{RNA}} = N_b + N_b^{\text{out}}$. The respective number of labelled bases inside and outside of the 25mer is $N_p^{\text{F},\text{in}}(\xi^{\text{T}})$ and $N_p^{\text{F},\text{out}}(\xi^{\text{T},\text{out}})$, respectively, where $\xi^{\text{T},\text{out}}$ is the subsequence of the target RNA exceeding the probe on both sides. The fluorescence intensity of a RNA fragment is related to the number of labelled c* and u*, which is given by the number of complementary G and A of the target gene according to

$$N_p^{\text{F},\text{S}} = N_p^{\text{F},\text{in}}(\xi^{\text{T}}) + N_p^{\text{F},\text{out}}(\xi^{\text{T},\text{out}}) = N_p^{\text{u}^*}(\xi^{\text{T}} + \xi^{\text{T},\text{out}}) + N_p^{\text{c}^*}(\xi^{\text{T}} + \xi^{\text{T},\text{out}}) = N_p^{\text{A}}(\xi^{\text{P}} + \xi^{\text{P},\text{out}}) + N_p^{\text{G}}(\xi^{\text{P}} + \xi^{\text{P},\text{out}})$$

if one assumes exclusively WC pairings. The contribution to the sensitivity of a selected probe owing to the number of potentially labelled bases per target, $N_p^{\text{F},\text{S}}$, is (see Eqs. 1 and 7)

$$Y_p^{\text{P},\text{F}} \equiv \log N_p^{\text{F}} - \langle \log N_p^{\text{F}} \rangle_{\text{set}} = \Delta_p^{\text{F}} \cdot \left[\sum_{k=1}^{N_b} \sum_{B=A,G} (\delta(B, \xi_k^{\text{P}}) - f_k^{\text{set}}(B)) - \sum_{k=1}^{N_b} \sum_{B=T,C} (\delta(B, \xi_k^{\text{P}}) - f_k^{\text{set}}(B)) \right] \quad (\text{A1})$$

$$\text{with } \Delta_p^{\text{F}} = \frac{\log N_p^{\text{F}} - \langle \log N_p^{\text{F}} \rangle_{\text{set}}}{\delta N_p^{\text{F},\text{in}}} \quad \text{and} \quad \delta N_p^{\text{F},\text{in}} = N_p^{\text{F},\text{in}} - \langle N_p^{\text{F},\text{in}} \rangle_{\text{set}} = \sum_{k=1}^{N_b} \sum_{B=A,T,G,C} (\delta(B, \xi_k^{\text{P}}) - f_k^{\text{set}}(B))$$

where averaging was performed over the probe set ($\Sigma \equiv \text{set}$).

The coefficient Δ_p^{F} specifies the contribution of fluorescence labelling per potentially labelled base in the considered target sequence of length N_b . Effectively each labelled base pair increases and each nonlabelled pair decreases the sensitivity by Δ_p^{F} . With $\delta N_p^{\text{F}} = \delta N_p^{\text{F},\text{in}} + \delta N_p^{\text{F},\text{out}}$ ($\delta N_p^{\text{F},i} = N_p^{\text{F},i} -$

$\langle N_p^{F,i} \rangle_{\text{set}}$, $i = \text{in, out}$) and $\langle \delta N^F \rangle_{\text{set}} = 0$ one obtains the following approximation for Δ_p^F in the limit of small $\delta N_p^F / \langle N_p^F \rangle_{\text{set}} \ll 1$, which is justified for sequence lengths $N_b^{\text{RNA}} > 20$,

$$\Delta_p^F = [\log(1 + \delta N_p^F / \langle N_p^F \rangle_{\text{set}}) - \langle \log(1 + \delta N_p^F / \langle N_p^F \rangle_{\text{set}}) \rangle_{\text{set}}] / \delta N_p^{F,\text{in}}$$

$$\approx (\ln 10 \cdot \delta N_p^{F,\text{in}} \cdot \langle N_p^F \rangle_{\text{set}})^{-1} (\delta N_p^F - \langle \delta N^F \rangle) = (1 + \delta N_p^{F,\text{out}} / \delta N_p^{F,\text{in}}) / (\ln 10 \langle N_p^F \rangle_{\text{set}}).$$

The binominal distribution $B(N^F, N_b^{\text{tot}}, p) = \binom{N_b^{\text{tot}}}{N^F} p^{N^F} (1-p)^{N_b^{\text{tot}} - N^F}$ specifies the probability to find N_p^F

potentially labelled nucleotides among a total sequence length of the target fragment of N_b^{RNA} nucleotides where $p \approx 0.5$ is the probability for a uracyl or a cytosine at any position of the target sequence. After substitution of the set average by the overall mean of the number of labels per target by $\langle N_p^F \rangle_{\text{set}} \approx \langle N_p^F \rangle_{\text{binom}} = p \cdot N_b^{\text{RNA}}$ one gets the relative fluorescence contribution per sequence position

$$\Delta_p^F \approx \left(1 + \frac{\delta N_p^{F,\text{out}}}{\delta N_p^{F,\text{in}}} \right) \cdot \Delta_0^F \quad \text{with} \quad \Delta_0^F = (\ln 10 \cdot p \cdot N_b^{\text{RNA}})^{-1} \approx \frac{2}{\ln 10 \cdot N_b^{\text{RNA}}} \quad (\text{A2})$$

Its mean value,

$$\Delta^F \equiv \langle \Delta_p^F \rangle_{\text{chip}} \approx \left(1 + \sqrt{\frac{N_b^{\text{out}}}{N_b}} \right) \cdot \Delta_0^F = \left(1 + \sqrt{\frac{N_b^{\text{RNA}}}{N_b} - 1} \right) \cdot \Delta_0^F, \quad (\text{A3})$$

provides the average contribution per considered base within a probe sequence of length N_b as a function of the total length of the RNA fragment, N_b^{RNA} . The mean incremental contributions are approximated in Eq. A3 by the standard deviation of the binominal distribution according to $\langle \delta N \rangle_{\text{set}} \approx p \cdot N^{0.5}$.

Bases in the probe sequence referring (B=A,G) and not-referring (B=T,C) to labels in the complementary target sequence add and subtract the constant contribution Δ^F to the sensitivity, respectively.

S2: Signal and sensitivity error of single Affymetrix GeneChips

The weighting factor for the least squares fits of the positional dependent sensitivity models is given by the variance of the experimental sensitivity data, $\omega_p^2 \approx \text{var}(Y_p)$ ¹. It can be estimated for each probe from chip replicates using standard error analysis. The SB sensitivity contributions are partly obtained from least square fits of the sensitivity data of single chips. We therefore developed a method, which estimates $\text{var}(Y_p)$ for each individual chip using selected probe intensities.

The variance of the sensitivity can be directly related to the variance of the respective signal intensity according to Eq. 6, $\text{var}(Y_p) \approx \text{var}(\log(I_p^P)) + \text{var}(\langle \log(I_p^P) \rangle) \approx \text{var}(\log(I_p^P)) (1 + (N_{\text{probe}} - 1)^{-1}) \approx \text{var}(\log(I_p^P))$ where $N_{\text{probe}} = 11 - 20$ is the number of probes per probe set. For the estimation of the chip-specific value of ω_p^2 we make use of the fact that a considerable number of PM and MM probes are present as replicates on each Affymetrix[®] chip. We identified repeated probes by comparison of all sequences present on the chip. For example, the human HG U133 chip contains 3463 probes in duplicate (2x), 725 in triplicate (3x), 186 fourfold (4x), 77 fivefold (5x), 37 sixfold (6x), 7 sevenfold (7x), 2 ninefold (9x) and one each for 12x, 16x and 20x. We calculated the variance, $\text{var}_{\text{exp}}(\log(I_p^P))$ (P=PM, MM) and log-averaged mean intensity, $\langle I_p^P \rangle = \exp(\langle \ln I_p^P \rangle_{\text{replicate}})$, for each of these groups of replicates for a selected chip.

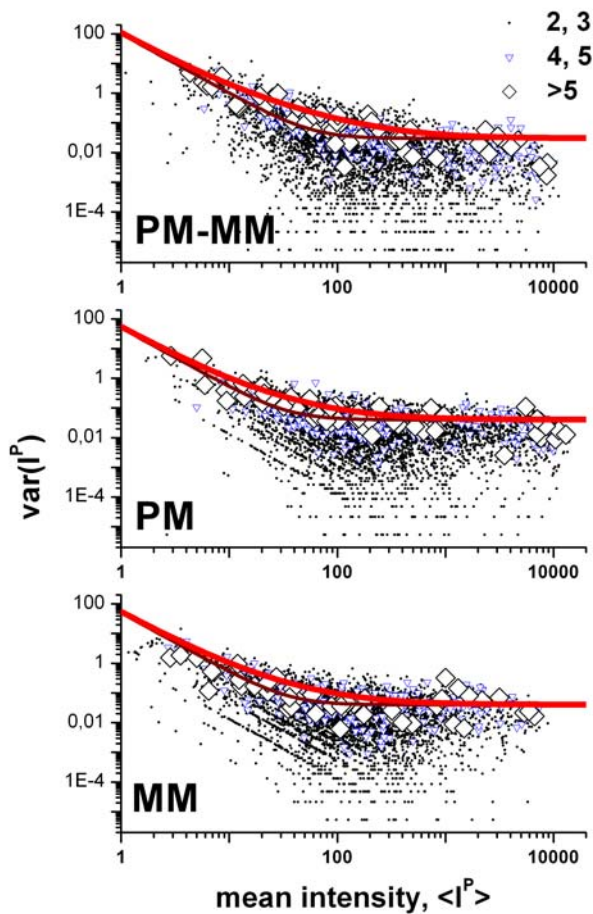


Figure S2: Log-log plot of the variance of the intensities of replicate PM and MM probes present on a HG U133 chip as a function of the mean intensity $\langle I^P \rangle$ averaged over probes present in duplicate and triplicate (see legend in the Figure: 2,3...small points), four- and fivefold (4,5...small triangles) and more than fivefold (>5...rhombes). The lines are calculated according to the error model (Eq. A5) with $a/b/c = 0.04/5/50$ (thick line) and $0.04/0/50$ (thin line) for P=PM and MM. The panel above shows the respective analysis of log-intensity differences, PM-MM. In this case the variance is given by $\text{var}(\log I^{\text{PM-MM}}) = \text{var}(\log I^{\text{PM}} - \log I^{\text{MM}})$ and the squared mean intensity by $\langle I^{\text{PM+MM}} \rangle^2 = \langle I^{\text{PM}} \cdot I^{\text{MM}} \rangle^{-1} = \exp[\langle \ln(I^{\text{PM}} + I^{\text{MM}}) \rangle]$. The thick and thin lines are calculated with $a/b/c = 0.03/10/100$ and $0.03/0/100$, respectively.

The uncorrected signal intensity can be rewritten according to Eq. 1 (in the original article) in a simplified version as $I_p^{P*} \approx \langle F_{\text{chip}} \cdot N^F \cdot c_{\text{RNA}} \rangle + e_F \cdot \exp[\langle \ln K_p^P \rangle + e_G] + \langle \beta_p^P \rangle + e_B$ where the angular brackets, $\langle \dots \rangle$, denote means over replicated probes. The e_i ($i=F, G, B$) are error terms and β_p^P is the optical background of each probe, which is not related to hybridization. With $\langle I_p^P \rangle \approx \langle F_{\text{chip}} \cdot N^F \cdot c_{\text{RNA}} \cdot K_p^b \rangle$ one obtains the background-corrected intensity

$$I_p^P \approx I_p^{P*} - \langle \beta_p^P \rangle \approx (\langle I_p^P \rangle + e_F) \cdot \exp(e_G) + e_B \quad (\text{A4})$$

The constant F_{chip} depends on the yield of labelling (fraction of labelled uracils and cytosines), on the number of oligos per spot and on the efficiency of the detector and of the imaging system (see ref. ² for details). Consequently the first error term, e_F , considers effects such as variations of the labelling efficiency, of the number of oligos per probe spot and of their density, of the RNA concentration and the noise of the detector and of the imaging system. The exponential term, e_G , can be rationalized as the error of the free energy of duplex formation, $\Delta G_p^b \propto -\ln K_p^P$, which is related, e.g., to incorrect sequences of individual oligos in each probe spot due to imperfect synthesis and/or to non-equilibrium effects of target binding. The last error, e_B , considers the noise of the detector and of the imaging system in the absence of hybridization.

The variance of log transformed and background corrected signal intensity is described to a good approximation by $\text{var}_{\text{mod}}[\log(I_p^P)] \approx a + c / \langle I_p^P \rangle^2$ with $a \approx s_G^2 / (\ln 10)^2$ and $c \approx (s_F^2 + s_B^2)$ if one assumes exclusively normally distributed error terms with mean 0 and variance s_i^2 ($i=F, G, B$). This result agrees with a previously proposed error model of microarray intensity data ³.

Figure S2 compares experimental and theoretical variance data of PM and MM intensities and of their difference in a double-logarithmic scale. The model curves systematically underestimate the experimental variance data in the intermediate intensity range, $100 < I^P < 1000$ (see thin lines in Fig. S2). Considerable better agreement was achieved if one adds a term $\sim \langle I^P \rangle^{-1}$ according to (see thick lines in Fig. 13)

$$\text{var}_{\text{mod}}(\log I_p^P) \approx a + \frac{b}{\langle I_p^P \rangle} + \frac{c}{\langle I_p^P \rangle^2} \quad (\text{A5})$$

The additional term can be tentatively rationalized as non-Gaussian error terms, which contribute to e_F . Here we use Eq. A5 without further specification as an empirical measure to estimate the weighting factor in the sum of squared residuals in the least squares fits as a function of the signal intensity.

The analysis of the intensity difference, $\text{var}(\log I^{\text{PM}} - \log I^{\text{MM}})$, as a function of $\langle I^{\text{PM}} \cdot I^{\text{MM}} \rangle^{-1} = \exp[-\langle \ln(I^{\text{PM}} + I^{\text{MM}}) \rangle]$ provides similar plots as that for PM and MM (Fig. S2, panel above). The respective background error is however increased whereas the signal error decreases compared with the respective error data of PM and MM probes. This result is compatible with a uncorrelated background noise of PM and MM intensities. In this case one expects for the background error of the log-difference of PM and MM probes a standard deviation of $s_B^2(\text{PM-MM}) \approx s_B^2(\text{PM}) + s_B^2(\text{MM}) \approx 2s_B^2(\text{PM})$, where the arguments PM and MM refer to the log-transformed intensities of the respective probes. On the other hand, the signal error term, “a”, of $\text{var}(\log I^{\text{PM}} - \log I^{\text{MM}})$ slightly reduces when

compared with that of the individual PM and MM probes. This result can be explained with correlations between the PM and MM intensities, which are discussed in the paper.

References

- (1) Bevington, P. R. *Data Reduction and Error Analysis for the Physical Sciences*; McGraw-Hill: New York, 1969.
- (2) Binder, H.; Preibisch, S. *Biophys. J.* **2005**, *89*, 337.
- (3) Rocke, D. M.; Durbin, B. *J. Comput. Biol.* **2001**, *8*, 557.

S3: Overview of SB free energy parameters of DNA/RNA duolexes

Relations between the positional dependent SB free energy and fluorescence contributions of Watson-Crick (WC) and self complementary (SC) pairings in DNA/RNA oligonucleotide duplexes of the PM and MM microarray probes upon-specific (S) and non-specific (NS) hybridization ^a

Du-plex	Single base contributions			
	probe level	base pair level		
	Probe P= position	PM,MM k≠13	PM k=13	MM k=13
NS	base pairing	WC: B-b ^c	WC: B-b ^c	WC: B ^c -b
	$\epsilon_{0,k}^{P,NS} \approx$ $\epsilon_{0,k}^{PM-MM,NS} \approx$	$\epsilon_{0,k}^{WC}$ 0	$\epsilon_{0,13}^{WC}$ $\epsilon_{0,13}^{WC-WC} \approx -(\log K_0^{PM,NS} - \log K_0^{MM,NS})$	$\epsilon_{0,13}^{WC}$
	$\Delta\epsilon_k^{P,NS}(B) \approx$ $\Delta\epsilon_k^{PM-MM,NS}(B) \approx$	$\Delta\epsilon_k^{WC}(B)$ 0	$\Delta\epsilon_{13}^{WC}(B)$ $\Delta\epsilon_{13}^{WC-WC}(B) = -\Delta\epsilon_{13}^{WC-WC}(B^c) \equiv \Delta\epsilon_{13}^{WC}(B) - \Delta\epsilon_{13}^{WC}(B^c)$ C ≈ T ≈ -G ≈ -A < 0	$\Delta\epsilon_{13}^{WC}(B^c)$
	$\Delta\phi_k^{P,NS}(B) \approx$ $\Delta\phi_k^{PM-MM,NS}(B) \approx$	$\Delta\phi_k^{WC}(B)$ 0	$\Delta\phi_{13}^{WC}(B)$ $\Delta\phi_{13}^{WC-WC}(B) \equiv \Delta\phi_{13}^{WC}(B) - \Delta\phi_{13}^{WC}(B^c) = 2\Delta\phi_{13}^{WC}(B)$ $ \Delta\phi_{13}^{WC-WC}(B) \approx \Delta^F $; G ≈ A ≈ -C ≈ -T > 0	$\Delta\phi_{13}^{WC}(B^c) = -\Delta\phi_{13}^{WC}(B)$
	$\sigma_k^{PM-MM,NS}(B) \approx$	0	$\Delta\epsilon_{13}^{WC-WC}(B) - \Delta\phi_{13}^{WC-WC}(B)$	
S	base pairing	WC: B-b ^c	WC: B-b ^c	SC: B ^c -b ^c
	$\epsilon_{0,k}^{P,S} \approx$ $\epsilon_{0,13}^{PM-MM,S}$	$\epsilon_{0,k}^{WC}$ 0	$\epsilon_{0,13}^{WC}$ $\epsilon_{0,13}^{WC-SC} \equiv \epsilon_{0,13}^{WC} - \epsilon_{0,13}^{SC} \approx -(\log K_0^{PM,S} - \log K_0^{MM,S})$	$\epsilon_{0,13}^{SC}$
	$\Delta\epsilon_k^{P,S}(B) \approx$ $\Delta\epsilon_k^{PM-MM,S}(B) \approx$	$\Delta\epsilon_k^{WC}(B)$ 0	$\Delta\epsilon_{13}^{WC}(B)$ $\Delta\epsilon_{13}^{WC-SC}(B) \equiv \Delta\epsilon_{13}^{WC}(B) - \Delta\epsilon_{13}^{SC}(B^c) \approx \Delta\epsilon_{13}^{WC}(B)$ C > G ≈ T > A	$\Delta\epsilon_{13}^{SC}(B^c)$
	$\Delta\phi_k^{P,S}(B) \approx$ $\Delta\phi_k^{PM-MM,S}(B) \approx$	$\Delta\phi_k^{WC}(B)$ 0	$\Delta\phi_{13}^{WC}(B)$ $\Delta\phi_{13}^{WC-SC}(B) \equiv \Delta\phi_{13}^{WC}(B) - \Delta\phi_{13}^{SC}(B^c) \approx 0$	$\Delta\phi_{13}^{SC}(B^c) = -\Delta\phi_{13}^{WC}(B^c)$
	$\sigma_k^{PM-MM,S}(B) \approx$	0	$\Delta\epsilon_{13}^{WC-SC}(B)$	

^a Single base related free energy (ϵ), fluorescence (ϕ) and sensitivity (σ) contributions to the probe intensities. The index k indicates the position of base B=A,T,G,C along the probe sequence. k=13 refers to the middle base whereas k≠13 refers to all positions outside the middle base. The superscript "c" denotes the complementary base, e.g., for B=A one gets of B^c=T. Single See text.