

An adsorption model of hybridization behaviour on oligonucleotide microarrays

Conrad J. Burden^{1,2}, Yvonne Pittelkow¹ and Susan R. Wilson¹
Centre for Bioinformation Sciences,
¹Mathematical Sciences Institute and
²John Curtin Curtin School of Medical Research,
Australian National University, Canberra, ACT 0200, Australia.

November 1, 2004

SUMMARY

A physical model of the hybridization of labelled target cRNA to short oligonucleotide probes at the surface of microarray chips is presented. The model is based on competing processes of chemical adsorption and desorption and includes the effects of non-specific hybridization. The adsorption process is modelled as a rate determining nucleation step involving a small number of base pairs, followed by a rapid ‘zipping up’ in which remaining base pairs bind. The model correctly predicts a hyperbolic Langmuir isotherm, consistent with data from the Affymetrix HG-U95A Latin Square spike-in experiment, and explains the differing responses of perfect match and mismatch features at saturation concentrations. A formula relating perfect match and mismatch responses in terms of duplex binding energies is also given.

1. Introduction

Oligonucleotide microarrays are designed to enable the evaluation of simultaneous expression of large numbers of genes in prepared messenger RNA samples. Details of the technology and the design and manufacture of Affymetrix GeneChip arrays, the focus of this paper, can be found in the review of Nguyen et al. (2002) or at the Affymetrix website <http://www.affymetrix.com/technology/index.affx>.

In the manufacture of Affymetrix arrays, single strand DNA probes, 25 bases in length are synthesized base by base onto a quartz substrate using a photolithographic process. They are attached to the substrate via short covalently bonded linker molecules at a separation of the order of 10 nanometres. A microarray chip surface is divided into some hundreds of thousands of regions called features, each about 20 microns square, the single strand DNA probes within each feature being synthesized to a specific nucleotide sequence.

A key step in the process of gene detection with microarrays in the laboratory is the hybridization of RNA target molecules fractionated to lengths of typically 50 to 200 bases onto the single strand DNA probes. The density of hybridized probe-target duplexes in each feature is detected via intensity measurements of fluorescent dye attached to the target RNA molecules. Each gene or EST is represented by a set of 11 to 16 pairs of features using sequences selected for their predicted hybridization properties

and specificity to the target gene. The first element of the pair, termed the perfect match (PM), is designed to be an exact match to the target sequence, while the second element, the mismatch (MM), is identical except for the middle base being replaced by its complement.

In this paper we propose a physical model of the hybridization process based on Langmuir adsorption theory. A number of previous studies have demonstrated the appropriateness of Langmuir adsorption theory for understanding probe-target hybridization at the surface of microarrays. Experimental work includes that of Nelson et al. (2001), Peterson et al. (2001), Peterson et al. (2002) and Dai et al. (2002). Analyses which have sought to match Langmuir adsorption isotherms with data from an Affymetrix spike-in experiment include those of Held et al. (2003), Hekstra et al. (2003), Lemon et al. (2003) and Burden et al. (submitted for publication). The ultimate aim of such work is to establish a functional relationship between measured fluorescence intensities and underlying target concentration parameterized by known physical properties such as probe base sequences. If such a relationship could be established, it would offer the possibility of an absolute measure of RNA target concentration, as opposed to an arbitrarily defined ‘expression measure’. A necessary ingredient in establishing this relationship is a model which accurately describes the physical chemistry of the hybridization process.

The novel aspect of our model is that it acknowledges that hybridization of probes and targets to form duplexes is a two step process: an initial rate-determining nucleation step involving two or three bases, followed by a rapid ‘zipping up’ step in which most, but not necessarily all, of the remaining bases bind (Cantor and Schimmel, 1980). In developing our model we make use of the differing responses of PM and MM features to known spiked-in concentrations of target RNA. The intended purpose of the mismatch feature, as stated in the manufacturer’s webpage, is to allow for the subtraction of signals caused by non-specific cross hybridization (Affymetrix Inc., 2002). In practice, however, there are problems with using the MM signals for this purpose (Irizarry et al., 2003). Rather than trying to interpret MM signals as a measure of non-specific hybridization, we instead view MM features simply as less responsive versions of the PM features. The difference between PM and MM probe signals can then be exploited as the result a single, well controlled change in one of the many parameters influencing the complicated process of hybridization. From this perspective one can obtain powerful insights into the physics and chemistry of hybridization at the microarray surface.

Langmuir adsorption theory is based on an assumption that there are two competing processes driving hybridization: adsorption, i.e. the binding of target molecules to immobilized probes to form duplexes, and desorption, i.e. the reverse process of duplexes dissociating into separate probes and target molecules. Both these processes are determined by chemical rate constants which depend on a number of factors including activation energies and temperature. If the adsorption rate is mainly determined by a nucleation step which, in the majority of cases, is remote from the middle base, adsorption rate constants will differ very little between elements of a (PM, MM) pair. On the other hand, desorption rate constants are affected by duplex binding energies and will be the main cause of difference between PM and MM responses. This picture of hy-

bridization has been borne out in the context of spotted microarrays by the observations of Dai et al. (2002). From the model proposed here we attempt to relate the differing responses of PM and MM features to the difference in binding energies between PM and MM duplexes.

In Section 2 we review current adsorption models including the model of Hekstra et al. (2003) which includes the effects of non-specific hybridization. This model predicts a hyperbolic response function, in good agreement with data from spike-in experiments. However, as we point out in Section 3, it is unable to explain the observed difference between PM and MM signals at saturation concentrations. In Section 4 our modified Langmuir model, which includes the effects non-specific hybridization and nucleation and partial zippering in the forward adsorption process is presented. The model maintains a hyperbolic response function and also correctly predicts that an MM feature will saturate with a lower intensity signal than its PM partner. In Section 5 we derive a relationship between the response curves of PM and MM features and the free energy difference between PM and MM duplexes, and compare the relationship with experimental data. Conclusions are drawn in Section 6.

2. The Langmuir isotherm model

Adsorption models of microarrays generally lead to a hyperbolic response function, or equilibrium Langmuir isotherm, relating RNA target concentration x to a measured equilibrium fluorescence intensity y , namely

$$y = y_0 + b \frac{x}{x + K}. \quad (1)$$

The isotherm is defined by three parameters: y_0 is the measured background intensity at zero target concentration, b is the saturation intensity above background at infinite target concentration, and K is the target concentration required to reach half saturation. The physical origins of these parameters will be discussed in detail below.

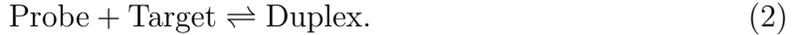
Recently we have carried out an extensive statistical analysis (Burden et al., submitted for publication) of fits of the hyperbolic response function to publicly available data from the Affymetrix Human HG-U95A Latin Square spike-in experiment (http://www.affymetrix.com/support/technical/sample_data). In this experiment genes were spiked in at cyclic permutations of the set of known concentrations, together with a background of cRNA extracted from human pancreas. The data consists of fluorescence intensity values from a set of 14 probesets corresponding to 14 separate genes, each containing 16 probe pairs. For each probeset a set of fluorescence intensity values are obtained for the 14 spiked-in concentrations (0, 0.25, 0.5, 1, 2, 4, ..., 1024) pM. The experiment was replicated three times using microarray chips from different wafers. In common with previous analyses of this data, our study concentrated on data from 12 of the 14 genes, omitting data from two defective genes.

In Fig 1 we show fits of Eq. (1) to fluorescence intensity data from the 16 PM and MM features corresponding to one of the 12 genes. These fits were estimated using a generalized linear model assuming the data at each spike-in concentration for each probe

sequence to be drawn from a Gamma distribution. The behaviour of this gene is typical of all 12 genes considered. Our findings are summarised as follows:

1. Measured fluorescence values can be approximated by a Gamma distribution with mean given by Eq. (1) and constant coefficient of variation, here ≈ 0.17 ,
2. The equilibrium isotherm Eq. (1) tracks fold changes from both PM and MM probes over the range of spiked-in concentrations from $< 1\text{pM}$ to $> 1000\text{pM}$,
3. All three parameters y_0 , b and K are probe sequence dependent (in contrast with the findings of Held et al. (2003)),
4. MM features invariably saturate at a lower asymptotic intensity $y_0 + b$ than their PM counterparts.

The Langmuir adsorption model is based on an assumption that the hybridization process is driven by competing processes of adsorption and desorption



Herein we shall always use the word ‘probe’ to indicate single strand DNA immobilised on the microarray, ‘target’ to indicate RNA in solution and ‘duplex’ to indicate a bound probe-target pair.

2.1 *Naive Langmuir model without non-specific hybridization*

We begin with a description of the simplest version of the Langmuir adsorption theory (Dai et al., 2002; Held et al., 2003). Let the fraction of probes within a feature which have formed duplexes by time t since the beginning of hybridization be $\theta(t)$. In the simplest form of the adsorption model the forward adsorption reaction is assumed to occur at a rate $k_f x(1 - \theta(t))$, proportional to target concentration x and fraction $(1 - \theta(t))$ of unoccupied probe sites. The backward desorption reaction is assumed to occur at a rate $k_b \theta(t)$, proportional to the fraction of occupied probe sites. k_f and k_b are the forward and backward reaction rates respectively. The target concentration x is assumed not to change significantly during hybridization. The fraction of occupied probe sites $\theta(t)$ at time t is therefore given by the differential equation

$$\frac{d\theta}{dt} = k_f x(1 - \theta) - k_b \theta. \quad (3)$$

Setting $d\theta/dt = 0$ gives the equilibrium isotherm

$$\theta = \frac{x}{x + K_S}, \quad (4)$$

where $K_S = k_b/k_f$. Setting y to be the measured fluorescence intensity and assuming the intensity above the physical background value a at zero concentration to be proportional to θ , we arrive at the Langmuir isotherm

$$y = a + b_S \theta = a + b_S \frac{x}{x + K_S}, \quad (5)$$

which takes the form of Eq. (1). The subscript S indicates parameters appropriate to ‘specific’ hybridization, as opposed to non-specific hybridization, discussed below.

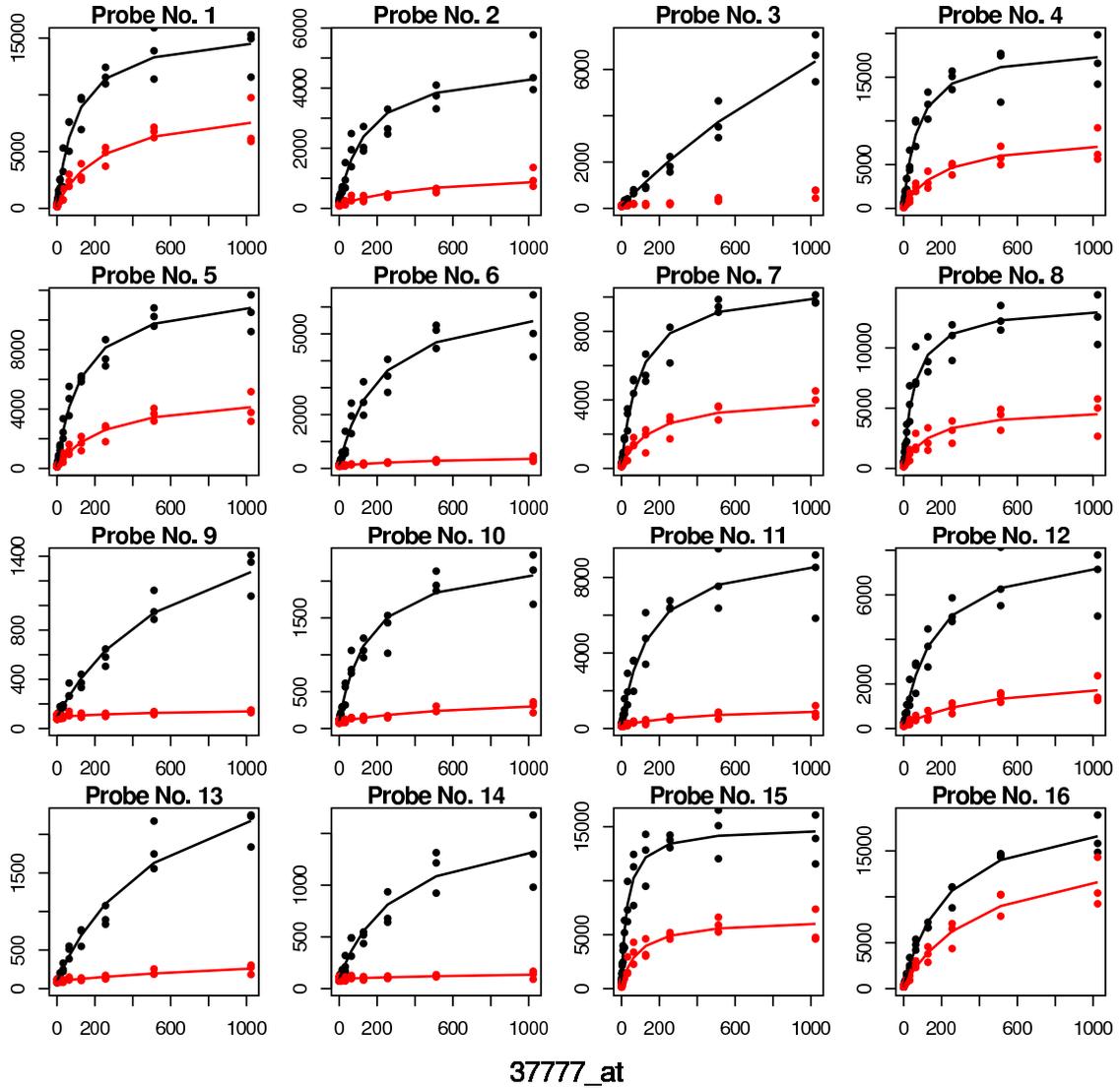


Figure 1. Fits of Eq. (1) to fluorescence intensity data for the 16 PM (black) and 16 MM (red) features of the gene *37777_at* probeset of the Affymetrix spike-in experiment. Concentrations (horizontal axes) are in picomolar and fluorescence intensities (vertical axes) are in the arbitrary units used in Affymetrix .cel files. The fit to MM probe No. 3 gave unphysical negative values to the parameters K and b and is not shown.

2.2 Naive Langmuir model with non-specific hybridization

Hekstra et al. (2003) determine changes to the Langmuir parameters caused by non-specific hybridization from a second target species. It is straightforward to extend their results to any number of non-specific species. Herein we define ‘specific’ to mean PM specific. All other hybridization will be referred to as ‘non-specific’.

Let the concentration of non-specific species i be z_i , $i = 1, 2, \dots$, the corresponding forward and backward reaction rates be k_{fi} and k_{bi} respectively, and the fraction of probe sites occupied by the i th non-specific species be ϕ_i . The kinetic equations are then

$$\frac{d\theta}{dt} = k_f x (1 - \theta - \sum_j \phi_j) - k_b \theta, \quad (6)$$

$$\frac{d\phi_i}{dt} = k_{fi} z_i (1 - \theta - \sum_j \phi_j) - k_{bi} \phi_i. \quad (7)$$

At equilibrium, setting $d\theta/dt = d\phi/dt = 0$ and defining $K_S = k_b/k_f$, $K_i = k_{bi}/k_{fi}$ gives

$$\begin{aligned} \theta &= \frac{x/K_S}{1 + x/K_S + \sum_i z_i/K_i} \\ \phi_i &= \frac{z_i/K_i}{1 + x/K_S + \sum_j z_j/K_j}. \end{aligned} \quad (8)$$

Introducing proportionality constants b_i for the non-specific hybridizations analogous to the b_S specified previously, the measured fluorescence intensity is given by

$$y(x) = a + b_S \theta + \sum_i b_i \phi_i \quad (9)$$

$$= y_0 + b \frac{x}{x + K}, \quad (10)$$

where

$$y_0 = a + A, \quad (11)$$

$$b = b_S - A, \quad (12)$$

$$K = K_S B, \quad (13)$$

and

$$A = \frac{1}{B} \sum_i \frac{b_i z_i}{K_i}, \quad (14)$$

$$B = 1 + \sum_i \frac{z_i}{K_i}. \quad (15)$$

The effect of non-specific hybridization is to maintain the hyperbolic form Eq. (1) while amending the isotherm parameters y_0 , b and K . The purpose of Eqs. (11) to (15) is to

relate the estimated isotherm parameters to the underlying physical parameters: a (the physical background value in the absence of any hybridization), b_S and b_i (proportionality constants relating the incremental change in measured intensity to an incremental change in duplex fraction for the ‘specific’ and ‘non-specific’ hybridizations respectively) and k_f , k_b , k_{fi} and k_{bi} (chemical reaction rate constants), given a set of non-specific background target concentrations z_i . The parameters b_S and b_i are a measure of the amount of fluorescent light emitted per hybridized target molecule. Fluorescent dye is bound only to the target molecules (in fact only to U and C bases), so b_S and b_i can only be functions of specific and non-specific target sequences, and not probe sequences.

3. Inconsistency of the naive model with observed PM/MM saturation intensities

In this section we argue that, if the Hekstra model with non-specific hybridization leading to the form given by Eqs. (6) to (15) is assumed, we are led inescapably to a conclusion that the PM and MM intensity measurements for a given probe pair must saturate at the same asymptotic intensity value, in contradiction with the observed fits to the Affymetrix spike-in experiment.

Consider two neighbouring features on a microarray, one PM and one MM, their probe sequences differing only by the middle base. Note that, in this paper, we define the word ‘specific’ to mean those target RNA which are exact complements to the PM sequence, even when dealing with the MM feature. In what follows this definition will prove useful given that, for most probe pairs, the dominant part of the MM signal at high spike-in concentrations in the Affymetrix experiment appears to come from hybridization of spiked-in target RNAs complementary to the PM sequence (see Fig. 1 for instance). Parameters relating to the PM and MM features will be indicated by superscripts PM and MM respectively.

Although the sum occurring in Eqs. (14) and (15) will be over the same set of non-specific targets for PM as for MM, one can expect $A^{\text{MM}} \neq A^{\text{PM}}$ since in general $K_i^{\text{MM}} \neq K_i^{\text{PM}}$. Considering the asymptotic intensities at high concentration, however, Eqs. (10) to (12) imply that, under the Hekstra model, the non-specific hybridization effects cancel out:

$$\begin{aligned} y^{\text{MM}}(\infty) &= y_0^{\text{MM}} + b^{\text{MM}} = a + b_S, \\ y^{\text{PM}}(\infty) &= y_0^{\text{PM}} + b^{\text{PM}} = a + b_S. \end{aligned} \tag{16}$$

An essential step in this argument is the claim that the parameters a and b_S do not differ between intensity measurements from neighbouring PM and MM features. For the physical background a this is clearly a reasonable assumption: physical properties of the chip in the absence of any hybridization, such as reflectance, are unlikely to vary significantly over a distance of a few microns. For the parameter b_S the argument is more subtle. Recall that b_S is a function of the amount of fluorescent light emitted per hybridized specific target RNA molecule, and as such is a function of the target sequence only, and not the probe sequence. By our current definition of ‘specific’ this is the *same sequence* for PM and MM, and so b_S is common to both. The Hekstra model formulated above then necessarily entails that $y_0^{\text{MM}} + b^{\text{MM}} = y_0^{\text{PM}} + b^{\text{PM}}$, in gross contradiction with

the values of y_0 , and b obtained by fitting the spike-in data.

The source of the problem is that the equilibrium solution of the model leading to the rate equations (6) and (7) entails that at sufficiently high specific target concentration, all probes form duplexes: as $x \rightarrow \infty$, $\theta \rightarrow 1$. That is, all probes in the feature form duplexes if saturated with enough specific target, even in the case of the MM feature. In practice however, this is clearly not true: for instance if the hybridization were carried out at the duplex melting temperature, then by definition θ would be one half.

This problem has been recognised previously in the context of the naive Langmuir model without non-specific hybridization by Peterson et al. (2002), who explain their experimental data by invoking a Sips isotherm to explain a lower MM response curve at high target concentrations. The Sips isotherm (Sips, 1948) is an empirical response curve which can be shown to correspond to an adsorption model in which chemical reaction rates are drawn from a pseudo-Gaussian distribution. Peterson et al.’s experimental results are indeed a good fit to the Sips isotherm, however their experiment differs from the conditions of the hybridization of Affymetrix chips in one important aspect, namely the hybridization temperature. The Peterson experiment was carried out at a hybridization temperature of 20°C, while Affymetrix microarrays are hybridized at 45°C, which is much closer to the duplex melting temperature. Furthermore, Peterson et al. found that heating the hybridization buffer to 37°C and then cooling back to 20°C almost completely removed any difference in equilibrium saturation intensities between PM and MM probes at a temperature well below the melting temperature.

To determine whether the hyperbolic or Sips isotherm is more appropriate for the Affymetrix spike-in data we have carried out a statistical analysis comparing the fits of the MM data to both isotherms. Our results, summarized in the Appendix, show that for the Affymetrix spike-in data the extra parameters involved in invoking the Sips isotherm are not significant, and that a hyperbolic response function adequately describes the data. We conclude that, at a hybridization temperature of 45°C, the asymptotic response of microarrays at high spike-in concentration is determined by a thermodynamic reduction in saturation efficiency which is more pronounced for MM than for PM features.

Thermodynamic effects manifest themselves through the temperature dependence of chemical rate constants. However, we have seen above that, irrespective of the values of chemical rate constants, a single step Langmuir model erroneously predicts that features will saturate at 100% efficiency, even in the presence of non-specific hybridization. In the next section we propose an adsorption model which takes into account that hybridization is a two step process involving an initial rate determining nucleation step. Our modified Langmuir model correctly predicts that MM probes should saturate at a lower efficiency than PM probes.

4. A solution to the saturation problem

The naive model leading to Eq. (3) assumes two processes, adsorption and desorption. Since the probes are immobilized and spatially separated, it is reasonable to assume that the desorption process should be an independent process from one probe to another. The assumption that the backward desorption should occur at a rate proportional to the

fraction θ of occupied probes therefore seems quite sound.

The assumption that the forward, adsorption process should occur at a rate that rises linearly with RNA target concentration, however, is questionable. Because the target molecules are mobile, and need to find their way to the probes, and given that duplex formation takes a certain time, at a sufficiently high concentration one could expect congestion at the chip surface. One way to model this is to take into account the time taken for a probe-target duplex to form. Textbook descriptions of duplex formation (see, for instance, Cantor and Schimmel (1980), pages 1215 to 1219) imply a two step process: a slow rate determining step in which an initial two or three base pairs form, followed by a fast ‘zipping-up’ step involving some, though not necessarily all, of the remaining base pairs. The model we propose here assumes that only zipped-up duplexes survive the washing of the chip and contribute to the measured fluorescence intensity.

4.1 *Modified Langmuir model without non-specific hybridization*

We first present our modified model without non-specific hybridization. We extend the single step Langmuir model behind Eq. (3) to a modified model which acknowledges that the forward, duplex forming, reaction involves two steps: a slow rate determining step in which the first two or three base pairs form, followed by a fast step involving some or all of the remaining base pairs. Denoting the probe and target molecules by P and T respectively, the partially formed duplex after the rate determining step by PT*, and the completed target-probe duplex by PT, the hybridization process can be modelled by the chemical reactions illustrated in Fig. 2.

Let the target concentration be x , the fraction of probes in a feature which have formed a zipped up duplex PT be θ and the fraction which have formed an initiated duplex PT* be ζ . The fraction of free single strand probes is then $1 - \theta - \zeta$. Reading off the chemical rate equations from Fig. 2 gives

$$\frac{d\zeta}{dt} = k_1x(1 - \theta - \zeta) - (k_{-1} + k_2)\zeta, \quad (17)$$

$$\frac{d\theta}{dt} = k_2\zeta - k_b\theta. \quad (18)$$

At equilibrium, $d\zeta/dt = 0$, giving

$$\zeta = \frac{k_1x(1 - \theta)}{k_{-1} + k_2 + k_1x}, \quad (19)$$

and $d\theta/dt = 0$, giving

$$k_b\theta = k_2\zeta = \frac{k_f x(1 - \theta)}{1 + k_f x/k_2}, \quad (20)$$

where we have defined

$$k_f = \frac{k_1k_2}{k_{-1} + k_2}. \quad (21)$$

Note that, for x small, Eq. (20) becomes the equilibrium form of the naive Langmuir model of Eq. (3), that is, k_f is the forward reaction rate constant in the limit of low target

concentration. At higher x , Eq. (20) says that the forward reaction rate is corrected by an effective Michaelis-Menton type law as the target concentration increases. Solving for θ gives

$$\theta = \lambda \frac{x}{x + K_S} \quad (22)$$

where we have further defined

$$\lambda = \frac{k_2}{k_2 + k_b}, \text{ and } K_S = \lambda \frac{k_b}{k_f}. \quad (23)$$

The important point to note is that, unlike the naive Langmuir model, at saturation concentration $x \rightarrow \infty$ the feature is not completely covered by zipped up probe-target duplexes. Assuming the washing process removes the unformed duplexes PT^* , the measured fluorescence intensity is given in the absence of non-specific hybridization by the analogue of Eq. (5),

$$y = a + b_S \theta = a + b_S \lambda \frac{x}{x + K_S}, \quad (24)$$

where $0 < \lambda < 1$ is the fraction of probes occupied by duplexes at saturation. As pointed out by Dai et al. (2002), the experimental evidence is that the main difference between PM and MM reaction rates is through the backward reaction rate, viz. $k_b^{MM} > k_b^{PM}$. From Eq. (23), this implies $\lambda^{PM} > \lambda^{MM}$, and that, ignoring non-specific hybridization, the asymptote $a + b_S \lambda$ of the isotherm is higher for PM than for MM. We next proceed to include the effects of non-specific hybridization within the modified model.

4.2 Modified Langmuir model with non-specific hybridization

Following the notation of Section 2 we consider competition between a specific target species T and a number of non-specific species T_i , and set

$$\begin{aligned} x &= \text{concentration of specific target RNA species T} \\ z_i &= \text{concentration of non-specific target RNA species } T_i \\ \theta &= \text{fraction of feature covered by specific duplexes PT} \\ \phi_i &= \text{fraction of feature covered by non-specific duplexes } PT_i \\ \zeta &= \text{fraction of feature covered by specific unzipped duplexes } PT^* \\ \eta_i &= \text{fraction of feature covered by non-specific unzipped duplexes } PT_i^* \\ \Phi &= \theta + \sum_j \phi_j + \zeta + \sum_j \eta_j \end{aligned}$$

so that $1 - \Phi$ is the fraction of the feature covered with single strand unmatched DNA probes. We also assume chemical reactions analogous to those shown in Fig. 2 for each non-specific species i with reaction rates k_{1i} , k_{-1i} , k_{2i} and k_{bi} . The complete set of rate

equations is

$$\begin{aligned}
\frac{d\zeta}{dt} &= k_1x(1 - \Phi) - (k_{-1} + k_2)\zeta, \\
\frac{d\eta_i}{dt} &= k_{1i}z_i(1 - \Phi) - (k_{-1i} + k_{2i})\eta_i, \\
\frac{d\theta}{dt} &= k_2\zeta - k_b\theta, \\
\frac{d\phi_i}{dt} &= k_{2i}\eta_i - k_{bi}\phi_i.
\end{aligned} \tag{25}$$

At equilibrium $d\theta/dt = d\phi_i/dt = d\zeta/dt = d\eta_i/dt = 0$, and the rate equations solve to give the isotherms

$$\begin{aligned}
\theta &= \frac{\lambda x/K_S}{1 + x/K_S + \sum_j z_j/K_j}, \\
\phi_i &= \frac{\lambda_i z_i/K_i}{1 + x/K_S + \sum_j z_j/K_j},
\end{aligned} \tag{26}$$

where

$$K_S = \lambda \frac{k_b}{k_f}, \quad K_i = \lambda_i \frac{k_{bi}}{k_{fi}}, \tag{27}$$

$$\lambda = \frac{k_2}{k_2 + k_b}, \quad \lambda_i = \frac{k_{2i}}{k_{2i} + k_{bi}}, \tag{28}$$

$$k_f = \frac{k_1 k_2}{k_{-1} + k_2}, \quad k_{fi} = \frac{k_{1i} k_{2i}}{k_{-1i} + k_{2i}}. \tag{29}$$

Assuming as before that only zipped up duplexes PT and PT_i contribute to the measured fluorescence intensity, we again have

$$\begin{aligned}
y &= a + b_S \theta + \sum_i b_i \phi_i \\
&= y_0 + b \frac{x}{x + K},
\end{aligned} \tag{30}$$

where now

$$y_0 = a + A \tag{31}$$

$$b = \lambda b_S - A \tag{32}$$

$$K = K_S B, \tag{33}$$

and

$$A = \frac{1}{B} \sum_i \frac{b_i \lambda_i z_i}{K_i} \tag{34}$$

$$B = 1 + \sum_i \frac{z_i}{K_i}, \tag{35}$$

and λ , K_S , λ_i and K_i are given in terms of the rate constants by Eqs. (27) to (29).

We see that the hyperbolic form of Eq.(1) is retained in Eq.(30), but with parameters amended to account for non-specific hybridization and the effects of congestion at the chip surface affecting the saturation asymptote. All non-specific hybridization effects are contained in the parameters A and B . The parameter λ is again the fraction of the feature covered by the specific PT duplexes at saturation, and if $k_b^{\text{MM}} > k_b^{\text{PM}}$, Eq. (28) implies $0 < \lambda^{\text{MM}} < \lambda^{\text{PM}}$. In other words, the asymptote $y(\infty) = y_0 + b = a + \lambda b_S$ is lower for MM than for PM isotherms.

4.3 Modified Langmuir model with non-specific hybridization and partial zipping

Finally we describe the full version of our model which allows for only partial zipping of duplexes. In the previous section we treated the zipped-up duplex PT as a single configuration. A more realistic model is to consider a distribution of configurations PT_α , $\alpha = 1, 2, \dots$, where the index α labels all possible partial zipperings, i.e. configurations in which only bases m to n , where $1 \leq m < n \leq 25$, are bound (Deutsch et al., 2004).

Define the forward zipping up reaction rate $\text{PT}^* \rightarrow \text{PT}_\alpha$ to be $p_\alpha k_2$ where $\sum_\alpha p_\alpha = 1$, and the backward decay rate for the process $\text{PT}_\alpha \rightarrow \text{P} + \text{T}$ to be $k_{b\alpha}$. The third equation in the set Eq. (25) is then replaced by

$$\frac{d\theta_\alpha}{dt} = p_\alpha k_2 \zeta - k_{b\alpha} \theta_\alpha, \quad \alpha = 1, 2, \dots \quad (36)$$

the remaining three equations being unchanged. It is straightforward to see that, in the equilibrium regime in which all time derivatives are set to zero, we recover all the results of the previous section provided we define the total fraction covered by zipped duplexes to be $\theta = \sum_\alpha \theta_\alpha$ and define an effective backward decay rate k_b by

$$\frac{1}{k_b} = \sum_\alpha \frac{p_\alpha}{k_{b\alpha}}. \quad (37)$$

More specifically, the hyperbolic isotherm solution of Eqs. (27) to (35) remains intact.

Note also that, with k_f as defined by Eq.(29), one easily obtains from the equilibrium version of Eqs. (25) and (36) that

$$k_{b\alpha} \theta_\alpha = p_\alpha k_f x (1 - \Phi). \quad (38)$$

That is, $p_\alpha k_f$ is the effective forward reaction rate at equilibrium for the process $\text{P} + \text{T} \rightarrow \text{PT}_\alpha$. Standard physical chemistry then gives the following relationship between the effective equilibrium forward and backward reaction rates and the free energies ΔG_α of the bound duplexes PT_α relative to the unbound state $\text{P} + \text{T}$:

$$\frac{k_f}{k_b} = \sum_\alpha \frac{p_\alpha k_f}{k_{b\alpha}} = C \sum_\alpha e^{-\Delta G_\alpha / (RT)}, \quad (39)$$

where T is absolute temperature, $R = 1.987 \text{ calK}^{-1}\text{mol}^{-1}$ is the gas constant, and the pre-exponential factor C is a measure of the frequency of encounters of reactants,

independent of their energy (Atkins, 2000). Our convention is such that $\Delta G < 0$ for a bound state.

From Eqs. (31) and (32) we also have

$$\lambda = \frac{y_0 + b - a}{b_S}. \quad (40)$$

Combining Eqs. (27), (33) and (40) then gives

$$\frac{k_f}{k_b} = \frac{B(y_0 + b - a)}{Kb_S}. \quad (41)$$

Eqs. (39) and (41) then imply the following relationship between the isotherm parameters y_0 , b and K , the sequence specific parameters b_S and ΔG_α and the non-specific hybridization parameter B :

$$\frac{B(y_0 + b - a)}{Kb_S} = C \sum_{\alpha} e^{-\Delta G_\alpha/(RT)}. \quad (42)$$

5. Comparison between PM and MM

A test of the theory proposed in the previous section is that it should be quantitatively consistent with fits to the Affymetrix spike-in data. The existence of MM features provides a convenient set of data in which only a single aspect of the experimental conditions, namely the central base of a probe sequence, has been altered and the resultant change in response measured. Below we compare the predictions of our modified Langmuir model with the observed differences between PM and MM gene expression measurements. As before we consider a neighbouring (PM, MM) pair of features and use the superscripts PM and MM to indicate quantities belonging to elements of the pair.

Because, as argued in Section 3, the parameters a and b_S are common to PM and MM, we immediately have from Eq. (42)

$$\frac{(y_0^{\text{MM}} + b^{\text{MM}} - a)/K^{\text{MM}}}{(y_0^{\text{PM}} + b^{\text{PM}} - a)/K^{\text{PM}}} = \frac{B^{\text{PM}} Z^{\text{MM}}}{B^{\text{MM}} Z^{\text{PM}}}, \quad (43)$$

where

$$Z^{\text{PM}} = \sum_{\alpha} e^{-\Delta G_\alpha^{\text{PM}}/(RT)}, \quad Z^{\text{MM}} = \sum_{\alpha} e^{-\Delta G_\alpha^{\text{MM}}/(RT)}. \quad (44)$$

The left hand side of Eq. (43) contains only the parameters y_0 , b and K which were estimated in our previous statistical analysis, and the physical background a which might be estimated for instance by taking the lowest fitted value of y_0 over the data set.

To estimate the right hand side, we split Z^{PM} into three pieces

$$Z^{\text{PM}} = U + V + W, \quad (45)$$

where U = a sum over duplex configurations bound within the set of bases not including the middle base, i.e. within the set of bases 1 to 12 or the set of bases 14 to 25; V =

a sum over duplex configurations bound over a stretch which terminates at the middle base; and W = a sum over the remaining duplex configurations, i.e. those which straddle the middle base. Binding energies ΔG of RNA/DNA duplexes can be estimated using a nearest neighbour stacking model (Sugimoto et al., 1995). In this model the duplex free binding energy is calculated as the sum of an initiation energy plus a contribution from each nearest neighbour pair of bases along the length of the duplex, while mismatches can be accounted for by including a contribution from the triplet straddling the mismatched base (Sugimoto et al., 2000). A consequence of the stacking model is that the free energy of a duplex configuration which is not bound at the middle base will be the same for a PM probe as for an MM probe. A further consequence is that the free energy of a configuration straddling the middle base will differ between PM and MM by a fixed amount $\Delta\Delta G = \Delta G_{\alpha}^{\text{MM}} - \Delta G_{\alpha}^{\text{PM}}$ for all such configurations α . Also, there will be no MM configurations which terminate at the middle base. It can be shown from these considerations that

$$Z^{\text{MM}} = U + We^{-\Delta\Delta G/(RT)}. \quad (46)$$

With some rearrangement, Eq. (43) then becomes

$$\begin{aligned} & \ln \frac{K^{\text{MM}}}{y_0^{\text{MM}} + b^{\text{MM}} - a} - \ln \frac{K^{\text{PM}}}{y_0^{\text{PM}} + b^{\text{PM}} - a} \\ &= \frac{\Delta\Delta G}{RT} + \ln \left[\left(1 + \frac{U + V}{W}\right) \left(1 + \frac{Ue^{\Delta\Delta G/(RT)}}{W}\right)^{-1} \right] + \ln \frac{B^{\text{MM}}}{B^{\text{PM}}}. \end{aligned} \quad (47)$$

As a reasonable first approximation, one might presume that $\Delta\Delta G$ should only depend on the middle letter of the probe sequence. If this were the case, a log-log plot of $K/(y_0 + b - a)$ for MM against PM over a range of probe sequences would have slope 1 and offset dependent only on the middle letter, plus corrections for the final two terms in Eq (47). The first correction term is due to partial zippering and the second from non-specific hybridization.

In Fig. 3 is shown a plot of $\ln[K/(y_0 + b - a)]$ for MM against PM. The parameters y_0 , b and K are estimated from fits of the Langmuir isotherm Eq. (1) to data from the 12 genes of the Affymetrix spike-in experiment. The physical background a is chosen to have a value of 85, which is slightly less than the minimum over all fits of y_0 , though the plot changes little over the range $60 < a < 100$. Each point corresponds to a different PM probe sequence, and colours indicate the middle base of the PM sequence. The four coloured lines are linear regressions to the model

$$X - Y = \beta_1(X + Y) + \beta_0 \quad (48)$$

where (X, Y) are the coordinates of the plotted points and β_0 , β_1 are fitted parameters. The fitting function is designed to be linear and symmetric with respect to the X and Y axes, that is, it is a standard least squares fit with respect to axes set at 45 degrees to the original axes. The lines are very close to parallel, as predicted, but the slope is far from 1. The fitted slopes $(1 - \beta_1)/(1 + \beta_1)$ are

$$\text{A: 1.498, C: 1.493, G: 1.515, T: 1.530.} \quad (49)$$

The dashed lines join constant values of the dimensionless free energy difference $\Delta\Delta G/RT$. We observe an increase in the offset as we move from the more strongly bound duplexes (small $K/(y_0 + b - a)$) to the weakly bound duplexes (large $K/(y_0 + b - a)$).

Apart from this overall change in offset, the ordering of the offset for the four middle bases is as expected. The main difference in free energies should be between the strong H-bond nucleotides (C,G) and weak H-bond nucleotides (A,T), and indeed the C and G fits come out above the A and T curves. On top of the H-bond shift, we expect a difference between the smaller pyrimidine nucleotide (C,T) and larger purine nucleotide (G,A) free energy difference because replacing a small nucleotide with a larger one will make the duplex less stable, whereas replacing a larger one by a smaller one will simply leave a dangling bond. From the plot we see that the C fit is shifted upwards relative to the G fit, and the T shifted by a roughly equal amount upwards from A.

By comparison, we show in Fig. 4 the plot obtained from the equivalent exercise for the naive Langmuir model, in which the analogue of Eq. (47) is simply (Hekstra et al., 2003)

$$\ln K^{\text{MM}} - \ln K^{\text{PM}} = \frac{\Delta\Delta G}{RT}. \quad (50)$$

The slope is closer to 1, but the four curves are not parallel, that is, there seems to be no clear dependence on the middle letter. Furthermore, the points are more broadly scattered, and some points correspond to negative values of $\Delta\Delta G$, entailing the unlikely prediction that there are probes for which the mismatch duplex is more stable than the perfect match.

Returning to Fig. 3, we address whether the two correction terms in Eq. (47) can account for the increase in offset from left to right across the plot. The first of these, the effect of partial zippering, can be estimated from the RNA/DNA stacking parameters measured by Sugimoto et al. (1995). Using an algorithm described in (Deutsch et al., 2004) to calculate sums over binding configurations, one finds that for the sequences occurring in the Affymetrix spike-in experiment, U/W and V/W both lie between about 10^{-4} and 10^{-10} with a median of about 10^{-7} . Expanding the log, one sees that the first correction term is first order in U/W and V/W , which is too small to account for the slope of the fits in Fig. 3. The second correction term, dependent on the ratio of non-specific hybridization parameters B defined by Eq. (35) is unfortunately harder to estimate. Ideally what is needed is a spike-in experiment without non-specific hybridization, in which case B^{MM} and B^{PM} could be measured from the change in the parameter K in the presence of non-specific RNA via Eq. (33). In the absence of such data we leave open the possibility that the final term in Eq. (47) can explain the slope of the fits in Fig. 3, or alternatively imply a more appropriate dependence than the linear form (48) employed above.

6. Discussion and Conclusions

An understanding of the physical processes driving hybridization is essential if the design of expression measures is to advance to a point where target concentration can be measured in absolute terms. In this paper we have presented a model of hybridization at the surface of an oligonucleotide microarray based on Langmuir adsorption theory, which

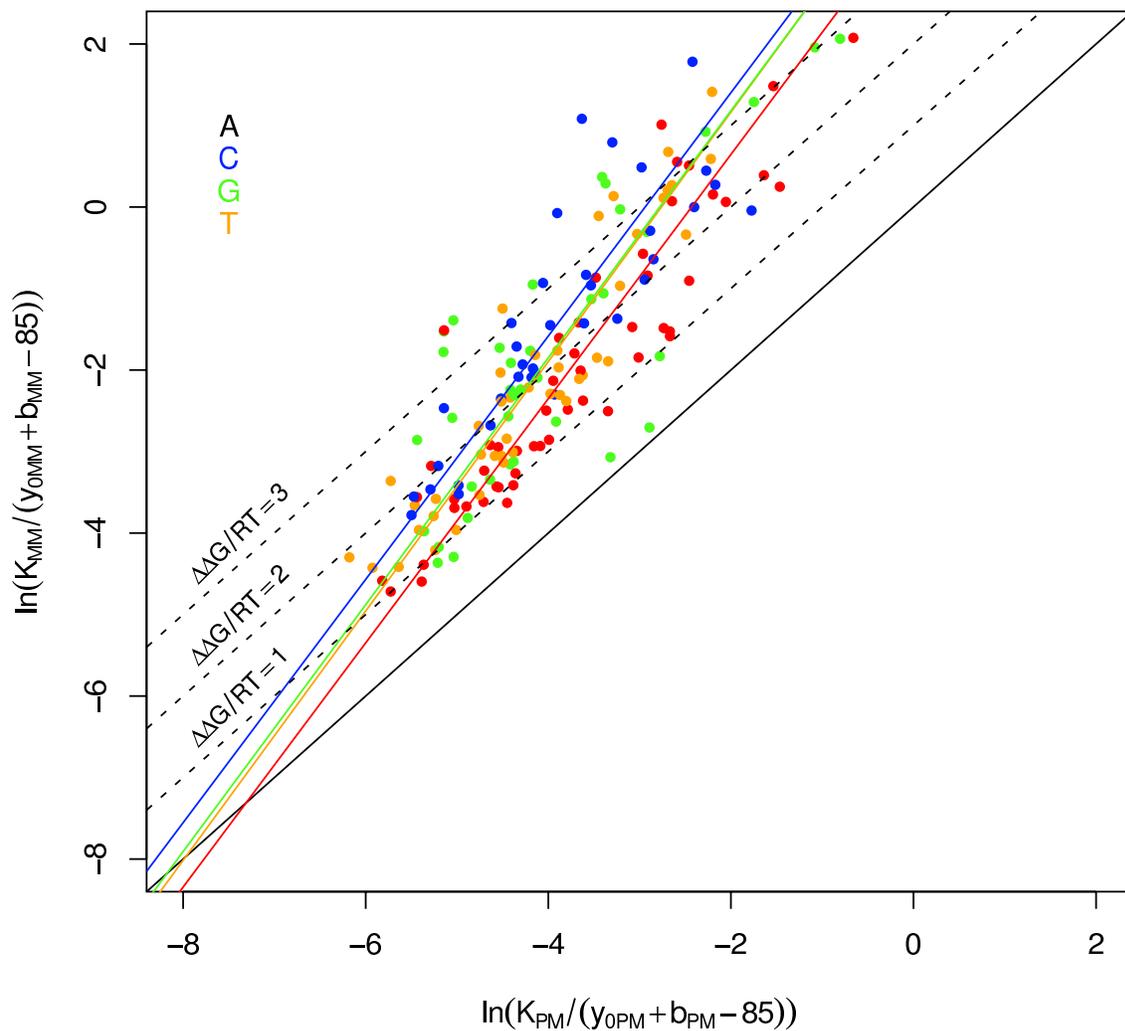


Figure 3. Plot comparing the two terms on the left hand side of Eq. (47). Colours indicate the middle base of the PM sequence. Our modified Langmuir theory predicts that probe sequences with the same difference in free binding energy between MM and PM duplexes will lie along the lines parallel to the diagonal with possible corrections for partial zippering and non-specific hybridization. At a hybridization temperature of 45°C, $\Delta\Delta G/RT = 1$ equates to a binding energy of 0.632 kcal mol⁻¹.

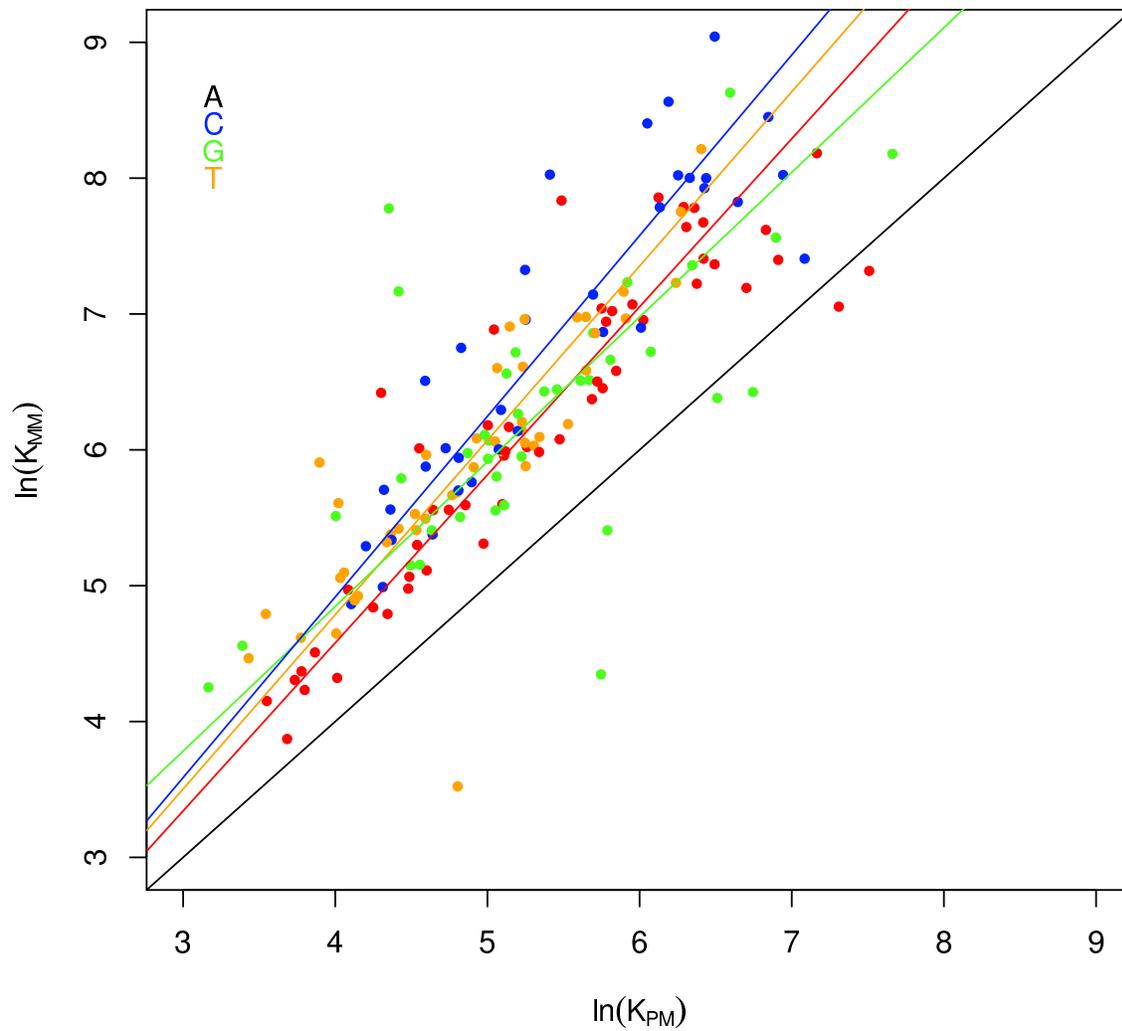


Figure 4. Plot comparing the two terms on the left hand side of Eq. (50). Colours indicate the middle base of the PM sequence.

takes into account 1) the effects of non-specific hybridization, 2) that DNA/RNA duplex formation proceeds via a slow initiation step followed by a rapid zipping-up step, and 3) a thermal distribution of partially zipped up RNA/DNA duplexes. We find that the hyperbolic form Eq. (1) of the Langmuir isotherm is maintained, in agreement with fits to the Affymetrix spike-in experiment. In Sections 4.2 and 4.3 we give the isotherm parameters in terms of more fundamental parameters driving the physical processes. One success of this model over previous models is that it is able to account for the clear difference in response of PM and MM features at saturation concentrations.

We have also made some advance towards establishing a relationship between the respective adsorption isotherm parameters of PM and MM probes. Knowing how these parameters are related can only improve the interpretation of microarray data, given that current expression measures either erroneously attempt to subtract the MM signal as a measure of nonspecific binding or ignore the MM signal completely. According to our model, Eq. (47) relates isotherm parameters of the MM and PM response, to the difference in free binding energy $\Delta\Delta G$ between a PM and an MM DNA/RNA duplex plus a possible correction arising from non-specific hybridization from non-specific binding. We have attempted, with limited success, a comparison of the fitted hyperbolic isotherms to Eq. (47) assuming $\Delta\Delta G$ to be a function of the middle base. We are unfortunately unable to determine the magnitude of the non-specific hybridization correction term in Eq. (47) with the available data. Further progress could be made if data were available from a repeat of the spike-in experiment without non-specific hybridization from a complex background. Furthermore, there is experimental evidence (Sugimoto et al., 2000) that $\Delta\Delta G$ is strongly dependent not only on the middle base, but on the complete triplet of bases straddling the middle base. However there is insufficient data from the current spike-in experiment to test Eq. (47) against the full set of 64 possible triplets.

Ultimately one would wish to be able to determine isotherm parameters solely from probe sequences. We see the current work as leading to a model not dissimilar in flavour to the positional dependent nearest neighbour model of Zhang et al. (2003), or the recent thermodynamic partial zippering model of Deutsch et al. (2004). Examination of more extensive spike-in experiments will be the subject of future research.

Appendix

In this appendix we carry out a statistical analysis of fits to the Langmuir isotherm, Eq. (1), and the Sips isotherm

$$y = y_0 + b \frac{x^\gamma}{x^\gamma + K^\gamma}, \tag{51}$$

to determine which model is the better fit to the MM data of the Affymetrix spike-in experiment. The method used is described in detail in an earlier paper which compares fits of the PM data to a number of isotherm models (Burden et al., submitted for publication).

The stochastic component of the fluorescence intensity y is assumed to be drawn from a gamma distribution. The data is fitted using the generalized linear model formalism as defined in McCullagh and Nelder (1989), in which the negative log likelihood of the fit, or

Table 1

Comparisons of fits to Langmuir and Sips isotherms. Δr is the decrease in residual degrees of freedom for each gene and ΔD_{scaled} is the corresponding scaled decrease in deviance.

Gene	Δr	ΔD_{scaled}	omitted probes
37777_at	14	6.43	3, 9
684_at	12	3.62	3, 5, 7, 8
1597_at	12	14.56	9, 11, 14, 15
38734_at	9	10.11	1, 3, 4, 9, 11, 12, 6
39058_at	5	11.34	1, 2, 3, 5, 6, 7, 9, 10, 12, 14, 16
36311_at	13	3.46	7, 8, 14
1024_at	16	15.19	
36202_at	15	6.18	6
36085_at	15	7.29	13
40322_at	16	39.58	
1091_at	14	3.83	1, 2
1708_at	12	2.36	11, 12, 13, 14
All genes	153	133.80	

deviance, is minimised over the parameters y_0 , b , K and, in the case of the Sips isotherm, also γ . To compare fits to the Langmuir and Sips models with r_L and r_S residual degrees of freedom and deviances D_L and D_S respectively, we use the scaled deviance

$$\Delta D_{\text{scaled}} = (D_L - D_S) \frac{r_S}{D_S}. \quad (52)$$

Note that $r_L > r_S \gg 1$. To evaluate the null hypothesis, $\gamma = 1$, ΔD_{scaled} can be compared with a chi-squared distribution with $\Delta r = r_L - r_S$ degrees of freedom (McCullagh and Nelder, 1989).

We were able to obtain fits with positive parameter values to both the Langmuir and Sips isotherms for about 80% of the probes. For most of the remaining cases the MM response was too small to provide a useful fit (see probes 3 and 9 in Fig. 1 and Table 1 for instance). Results for the scaled deviance are shown in Table 1. The total deviance of 133.8 lies at the 13th percentile of a chi-squared distribution with 153 degrees of freedom, showing no reason to consider a more complex model than the Langmuir isotherm. Finally, a histogram of the fitted values of the Sips parameter, Fig. 5, shows that the Sips parameter is symmetrically distributed about $\gamma = 1$, as expected if the Langmuir isotherm is the more accurate model.

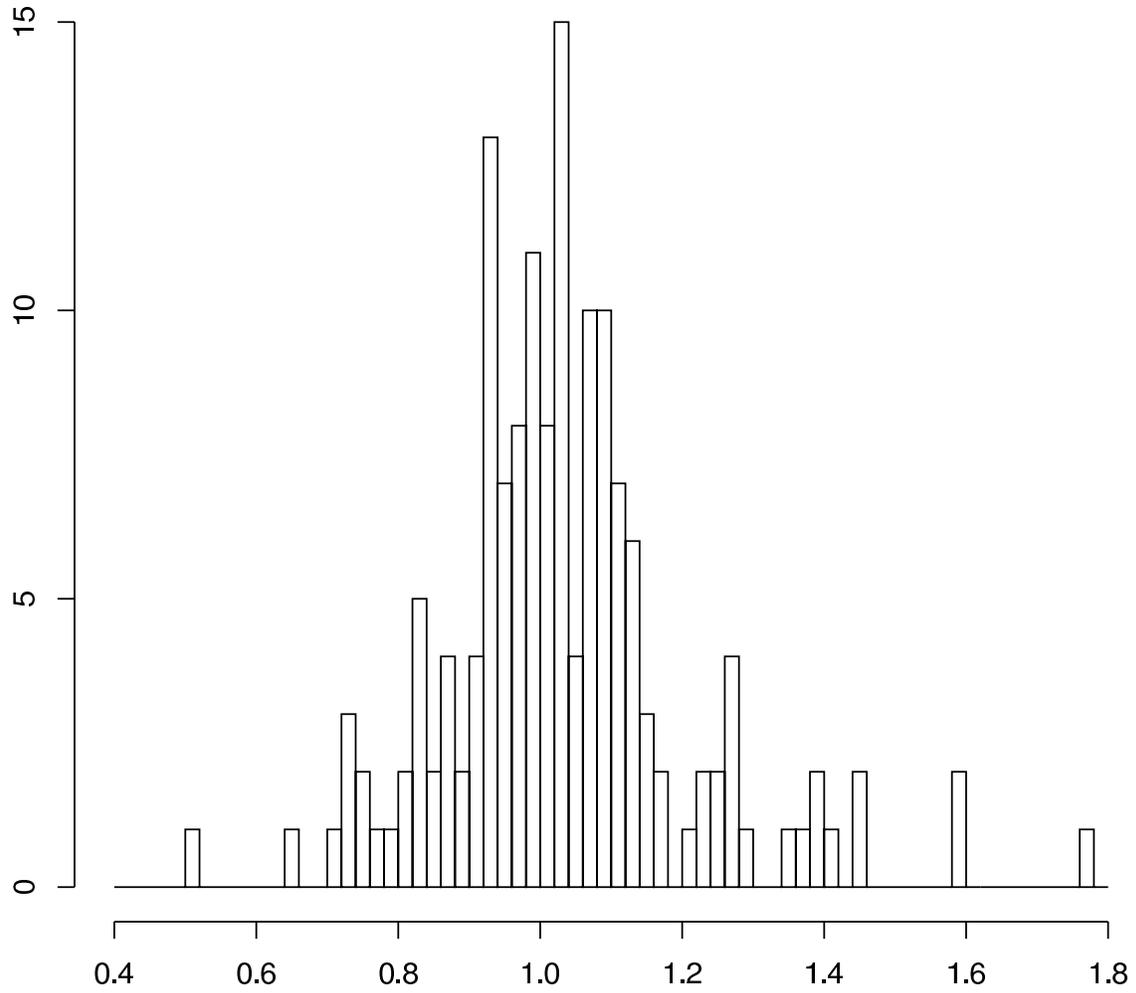


Figure 5. Histogram of fitted values of the Sips parameter γ for the MM data.

REFERENCES

- Affymetrix Inc. (2002). Statistical algorithms description document. Available at <http://www.affymetrix.com/support/technical/whitepapers.affx>.
- Atkins, P. W. (2000). *Physical Chemistry*. W. H. Freeman and Co., New York, NY, USA, 6th edition.
- Cantor, C. R. and Schimmel, P. R. (1980). *Biophysical Chemistry, Part 3: The behaviour of biological macromolecules*. W. H. Freeman and Co., San Francisco, CA, USA, 1st edition.
- Dai, H., Meyer, M., Stepaniants, S., Ziman, M. and Stoughton, R. (2002). Use of hybridization kinetics for differentiating specific from non-specific binding to oligonucleotide microarrays. *Nucleic Acids Research* **30**, e86.
- Deutsch, J. M., Liang, S. and Narayan, O. (2004). Modeling of microarray data with zippering. *arXiv* pages q-bio.BM/0406039.
- Hekstra, D., Taussig, A. R., Magnasco, M. and Naef, F. (2003). Absolute mRNA concentrations from sequence-specific calibration of oligonucleotide arrays. *Nucleic Acids Research* **31**, 1962–1968.
- Held, G. A., Grinstein, G. and Tu, Y. (2003). Modeling of DNA microarray data by using physical properties of hybridization. *Proceedings of the National Academy of Science* **100**, 7575–7580.
- Irizarry, R. A., Hobbs, B., Collin, F., Beazer-Barclay, Y. D. and et al. (2003). Exploration, normalization and summaries of high density oligonucleotide array probe level data. *Biostatistics* **4**, 249–264.
- Lemon, W. J., Liyanarachchi, S. and You, M. (2003). A high performance test of differential gene expression for oligonucleotide arrays. *Genome Biology* **4**, R67.1 – R67.11.
- McCullagh, P. and Nelder, J. A. (1989). *Generalized Linear Models*. Chapman and Hall, London, UK, 2nd edition.
- Nelson, B. P., Grimsrud, T. E., Liles, M. R., Goodman, R. M. and Corn, R. M. (2001). Surface plasmon resonance imaging measurements of DNA and RNA hybridization adsorption onto DNA microarrays. *Analytical Chemistry* **73**, 1–7.
- Nguyen, D. V., Arpat, A. B., Wang, N. and Carroll, R. J. (2002). DNA microarray experiments: Biological and technological aspects. *Biometrics* **58**, 701–717.
- Peterson, A. W., Heaton, R. J. and Georgiadis, R. M. (2001). The effect of surface probe density on DNA hybridization. *Nucleic Acids Research* **29**, 5163–5168.
- Peterson, A. W., Wolf, L. K. and Georgiadis, R. M. (2002). Hybridization of mismatched or partially matched DNA at surfaces. *Journal of the American Chemical Society* **124**, 14601–14607.
- Sips, R. (1948). On the structure of a catalyst surface. *Journal of Chemical Physics* **16**, 490–495.
- Sugimoto, N., Nakano, M. and Nakano, S. (2000). Thermodynamics - structure relationship of single mismatches in RNA/DNA duplexes. *Biochemistry* **39**, 11270–11281.
- Sugimoto, N., Nakano, S., Katoh, M., Matsumura, A., Nakamuta, H. and Ohmichi,

- T. (1995). Thermodynamic parameters to predict stability of RNA/DNA hybrid duplexes. *Biochemistry* **34**, 11211–11216.
- Zhang, L., Miles, M. F. and Aldape, K. D. (2003). A model of molecular interactions on short oligonucleotide microarrays. *Nature Biotechnology* **21**, 818–821.