

Models for Microarray Analysis: Sequence Effects and RNA Degradation

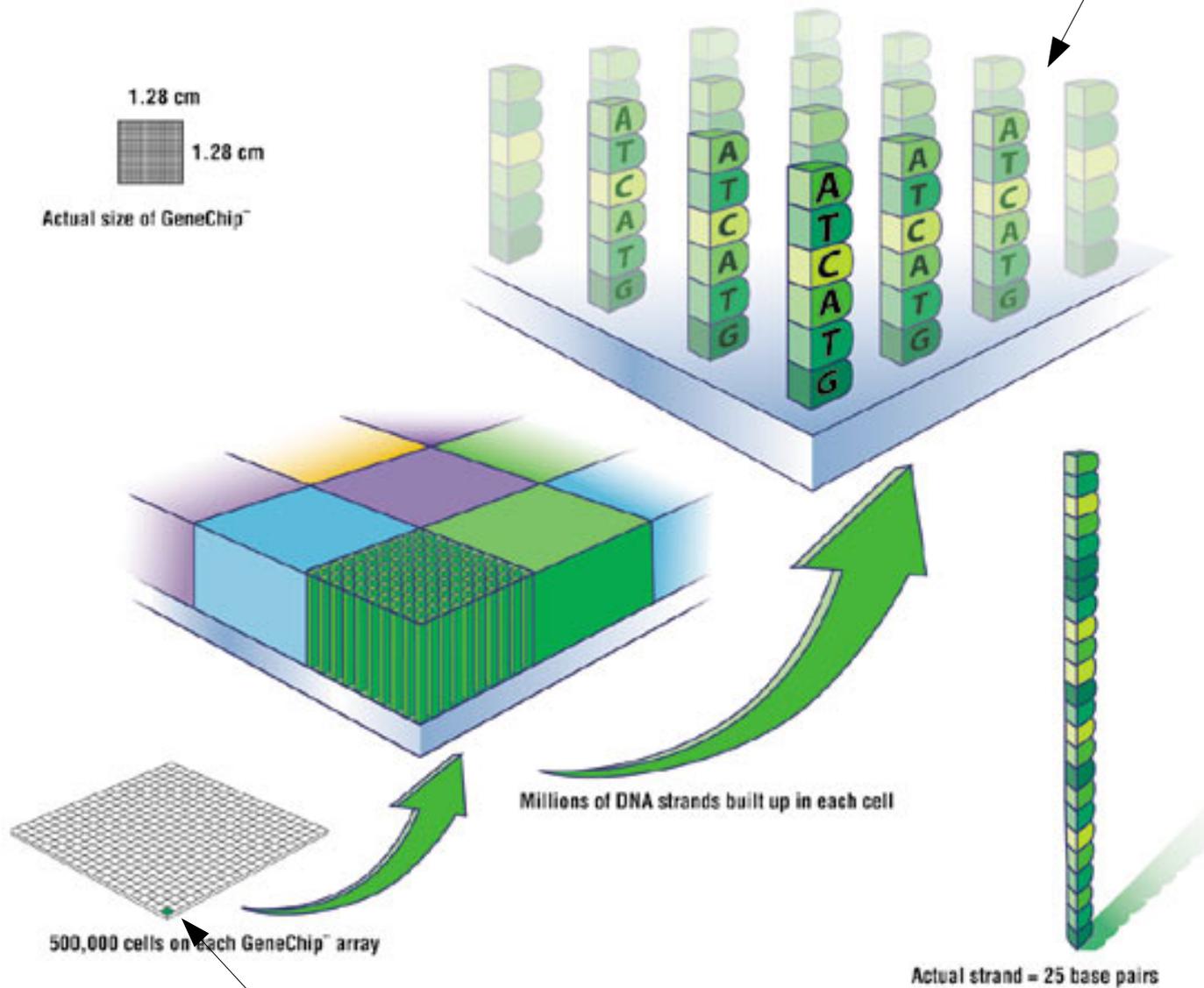
„All models are wrong, and increasingly
you can succeed without them.”

(Peter Norvig, Research Director Google Inc.)

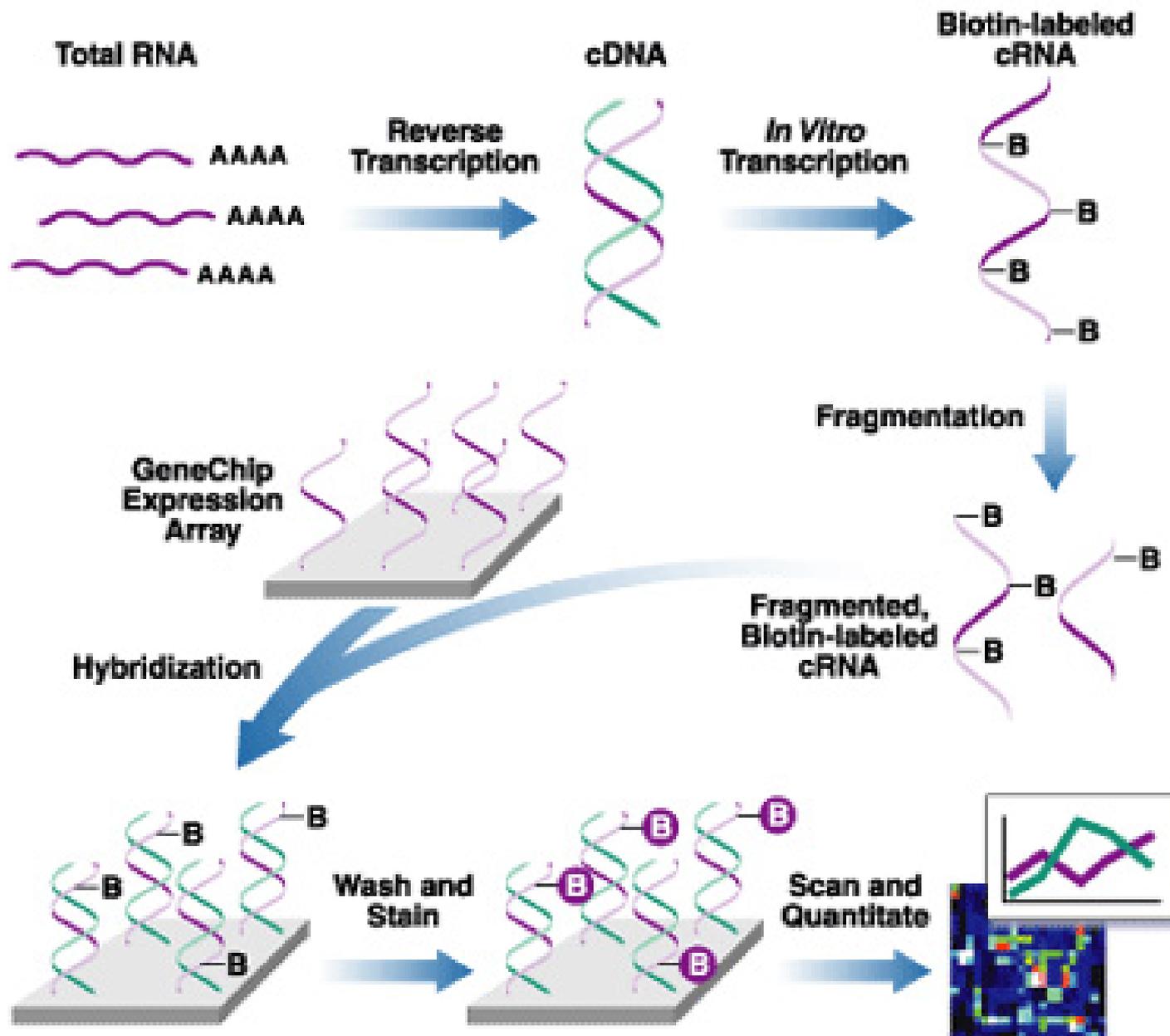
Affymetrix Microarrays



Probe Sequence $\xi_p \in \{A, C, G, T\}^{25}$

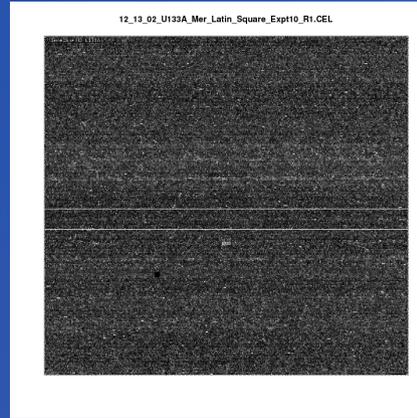


Probe Position on the Chip X_p, Y_p

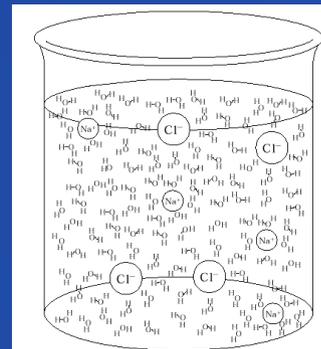


Probe Intensity I_p

Probe 1 I_p
Probe 2 I_p
Probe 3 I_p
...
Probe 11 I_p



„Microarray Correction”

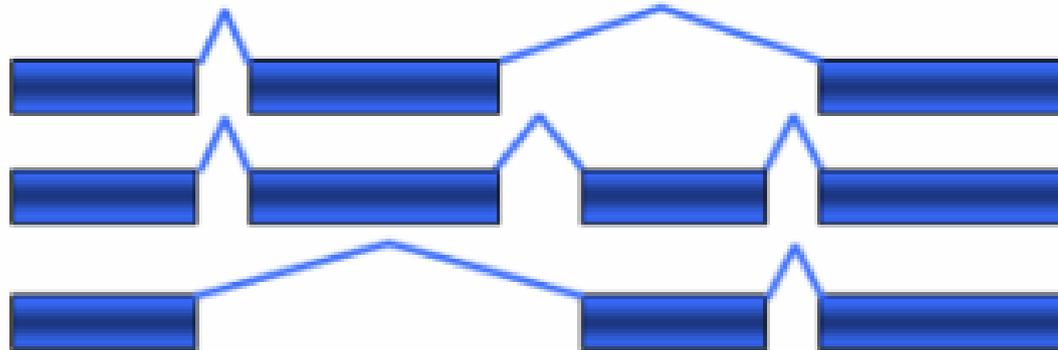


mRNA
Concentration
(Expression Measure)

Genomic locus



mRNA transcripts



Probe locations in
3' focused arrays



Probeset

Probes are designed to match sequences towards the 3' end in expression arrays

Microarray Correction Methods

Statistical Methods

MAS5 VSN
RMA FARMS
PLIER

gcRMA

Hook

Use of Physico-chemical Models

The „Hook“-Method

Saturation

$$I_p = \frac{L_p}{1 + L_p/I^{max}} + I^{min}$$

The „Hook“-Method

Saturation

$$I_p = \frac{L_p}{1 + L_p/I^{max}} + I^{min}$$

Two-species: specific
and non-specific
binding [S]/[NS]

$$L_p = L_p^N + L_p^S$$

The „Hook“-Method

Saturation

$$I_p = \frac{L_p}{1 + L_p/I^{max}} + I^{min}$$

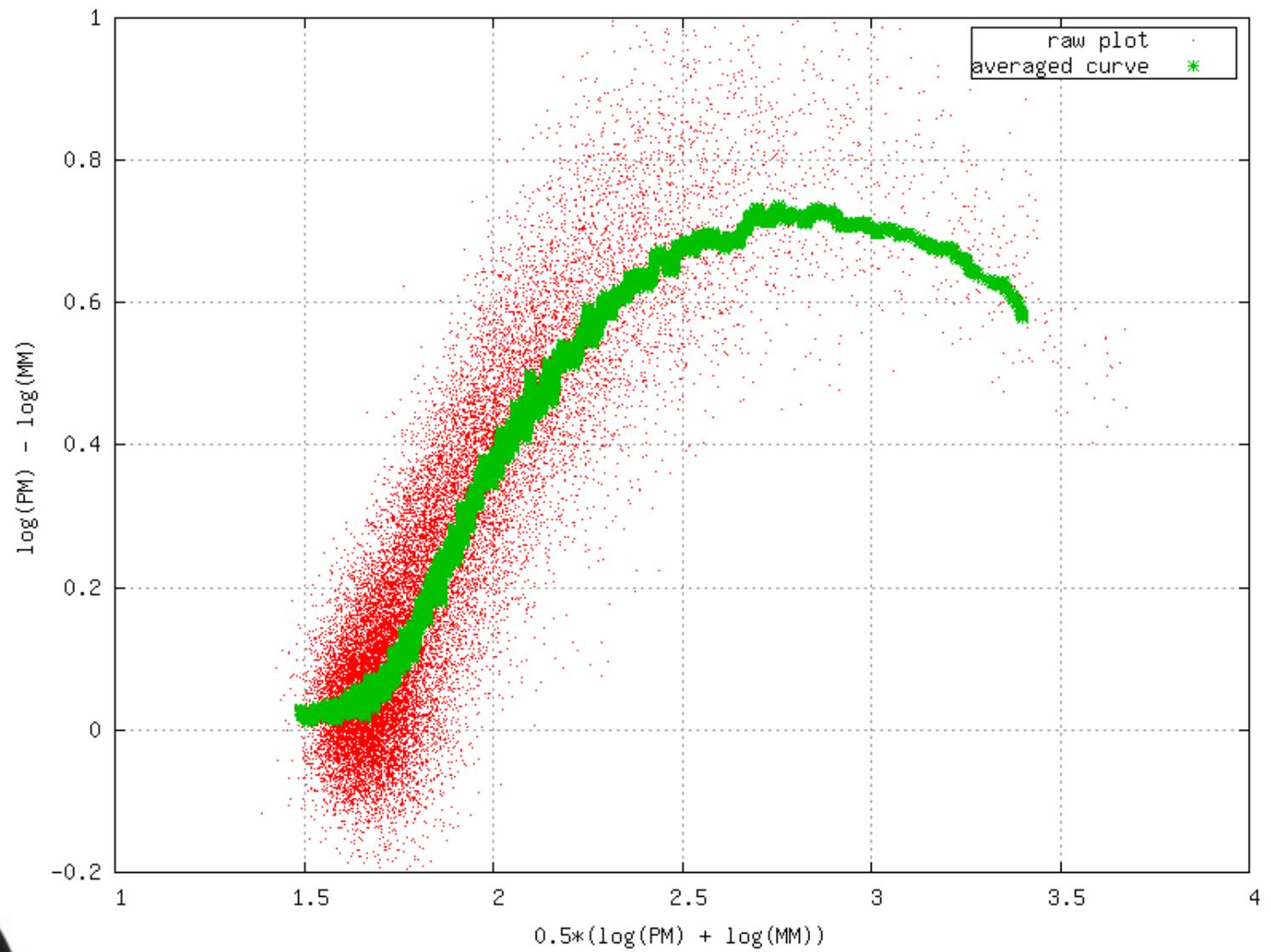
Two-species: specific
and non-specific
binding [S]/[NS]

$$L_p = L_p^N + L_p^S$$

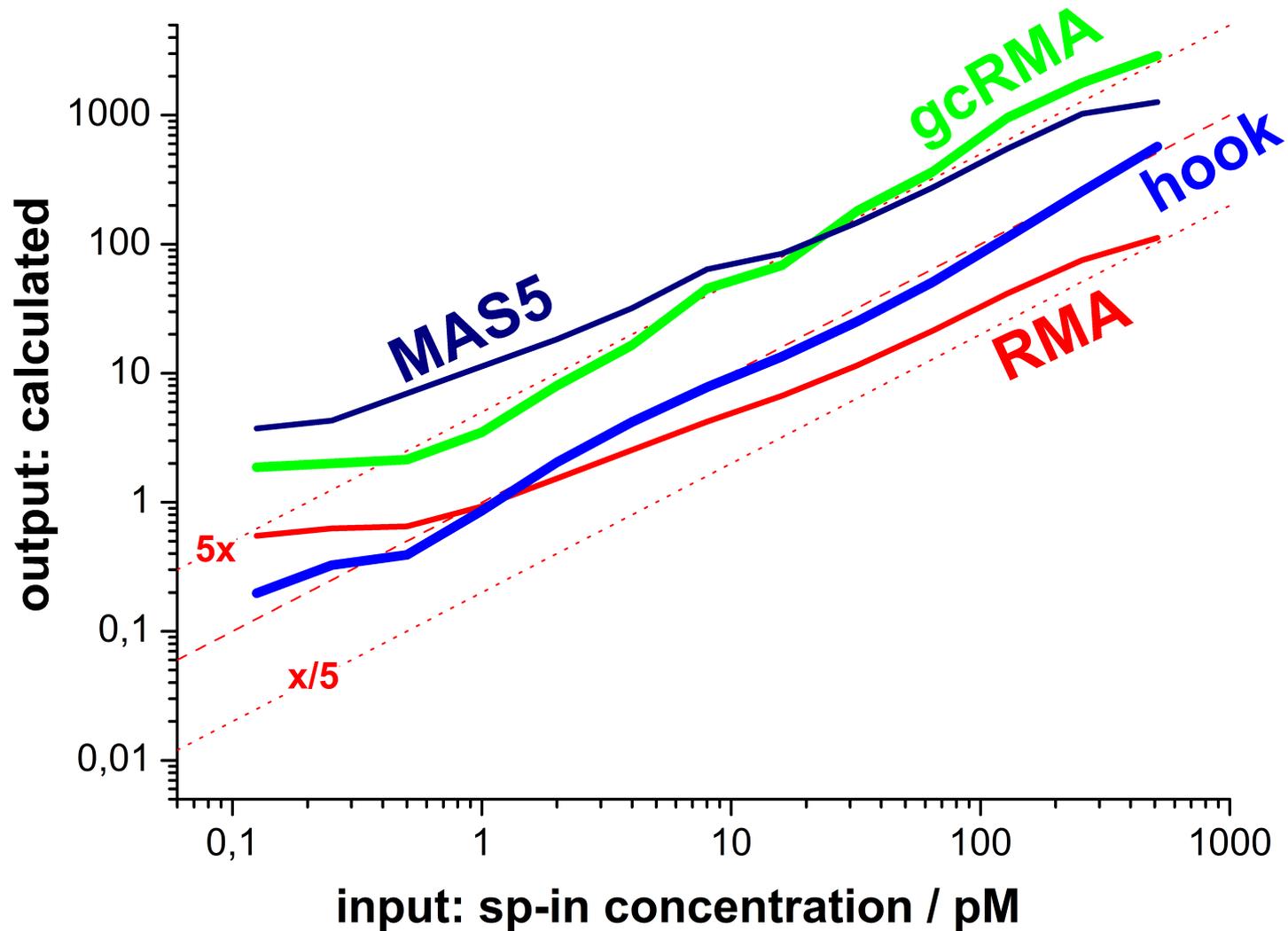
Sequence effects

$$\log L^S = \log[S] + \delta A(\xi)$$

$$\log L^N = \log[N] + \delta A(\xi)$$



The „Hook”



Hans Binder and Stephan Preibisch: "Hook"-calibration of GeneChip-microarrays:
Theory and algorithm; *Algorithms for Molecular Biology*, 2008

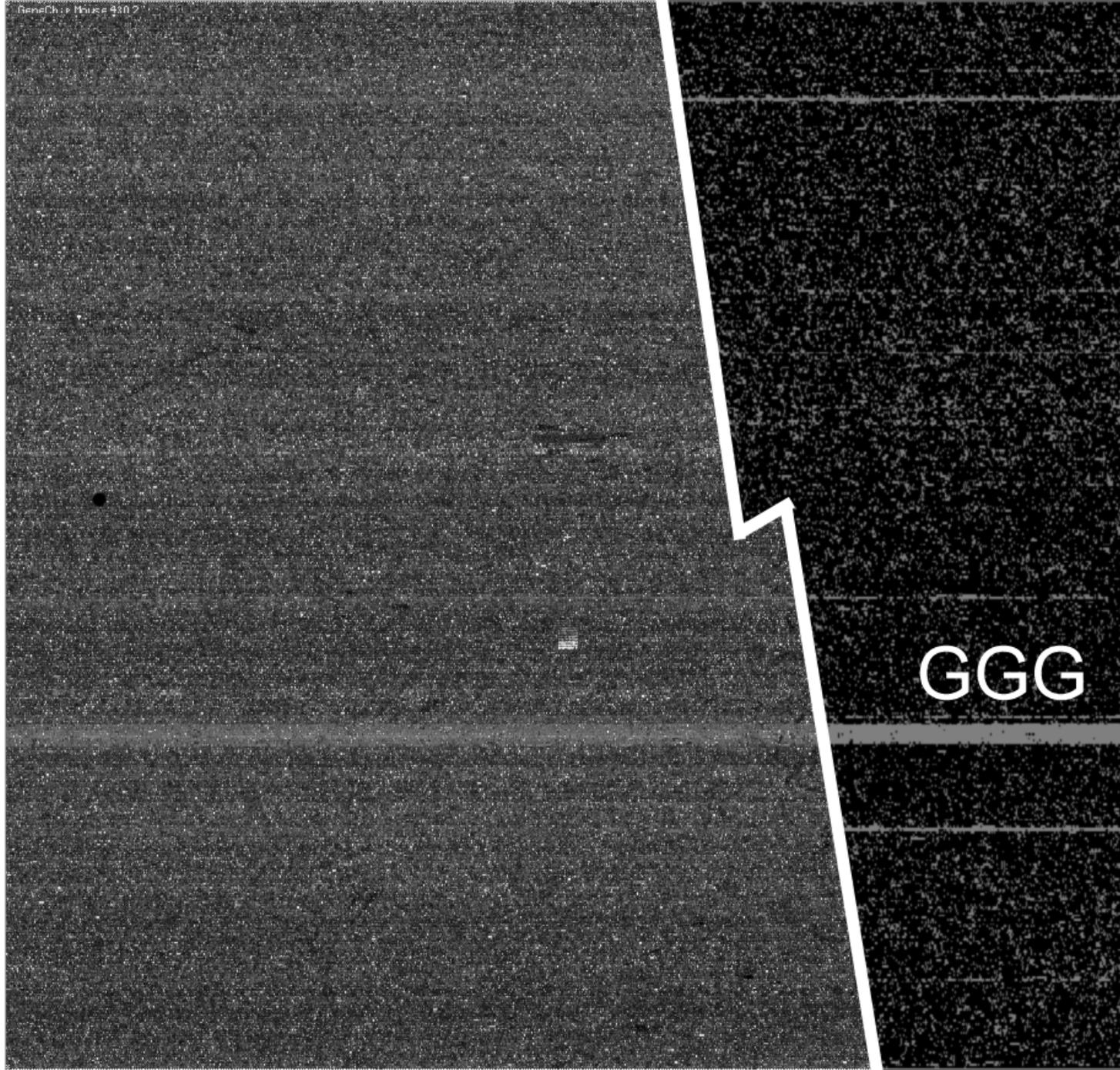
http://www.izbi.de/downloads_links/programs/hook.php

GeneChip Mouse 430 2'



GeneCh: Mouse 4302

GGG



Sequence Effects

$$\log L^N = \log[N] + \delta A(\xi)$$

$$\delta A(\xi) = \sum_{k=1}^{25-r+1} \sigma_k(\xi^{k, k+r-1})$$

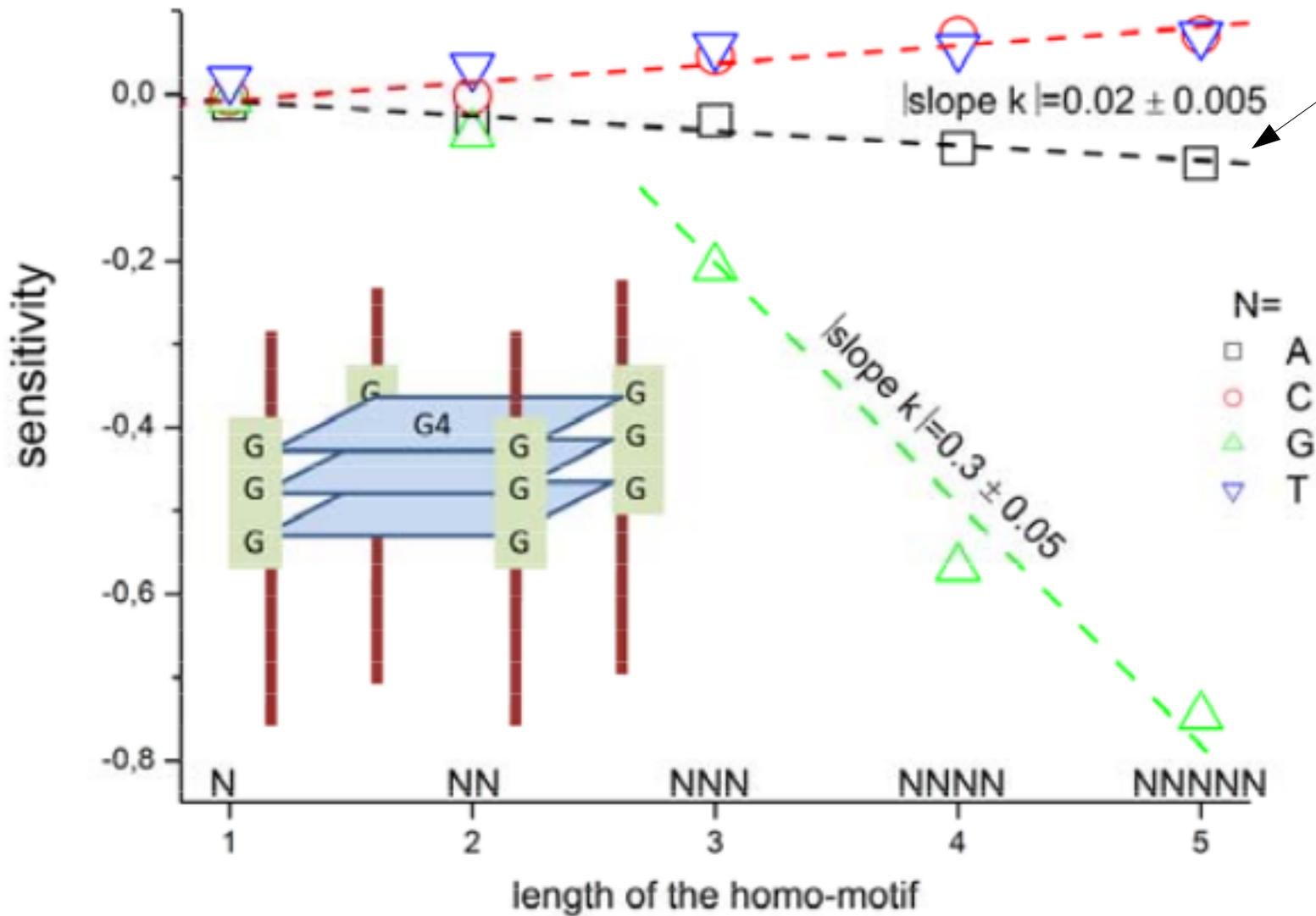
Example:

$r = 2$ (Nearest-Neighbor-Model), Sequence $\xi = \text{GTGACCGTTATCCA}$

$$\delta A(\xi) = \sigma_1(GT) + \sigma_2(TG) + \cdots + \sigma_{24}(CA)$$

Stacks of G-tedrads

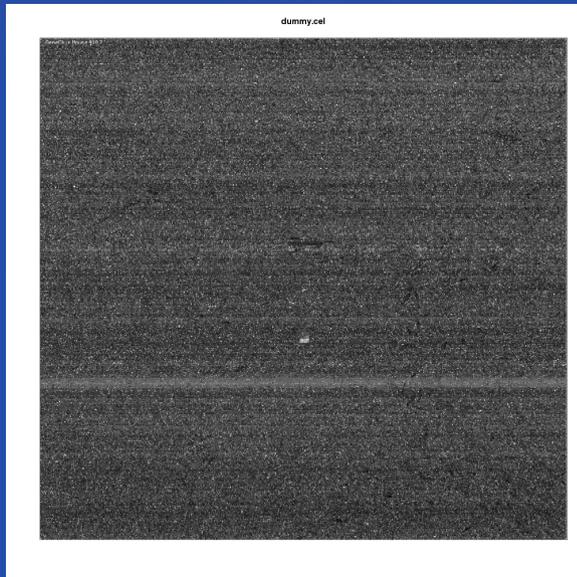
$$\sum_{k=1}^{20} \sigma_k(AAAAA)$$



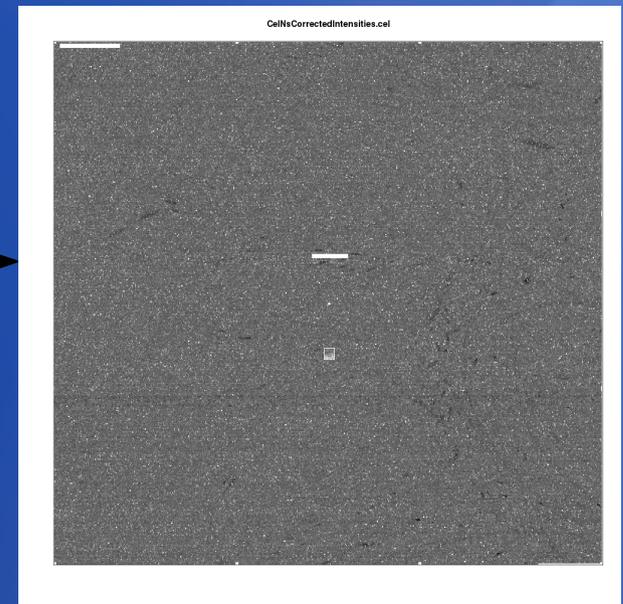
Correcting with the NN+GGG Sequence Affinity Model

$$\delta A(\xi_p) = \sum_{i=1}^{24} \sigma_i(\xi_{p,i\dots i+1}) + \sum_{i=1}^{23} \tilde{\sigma}_i(\xi_{p,i\dots i+2})$$

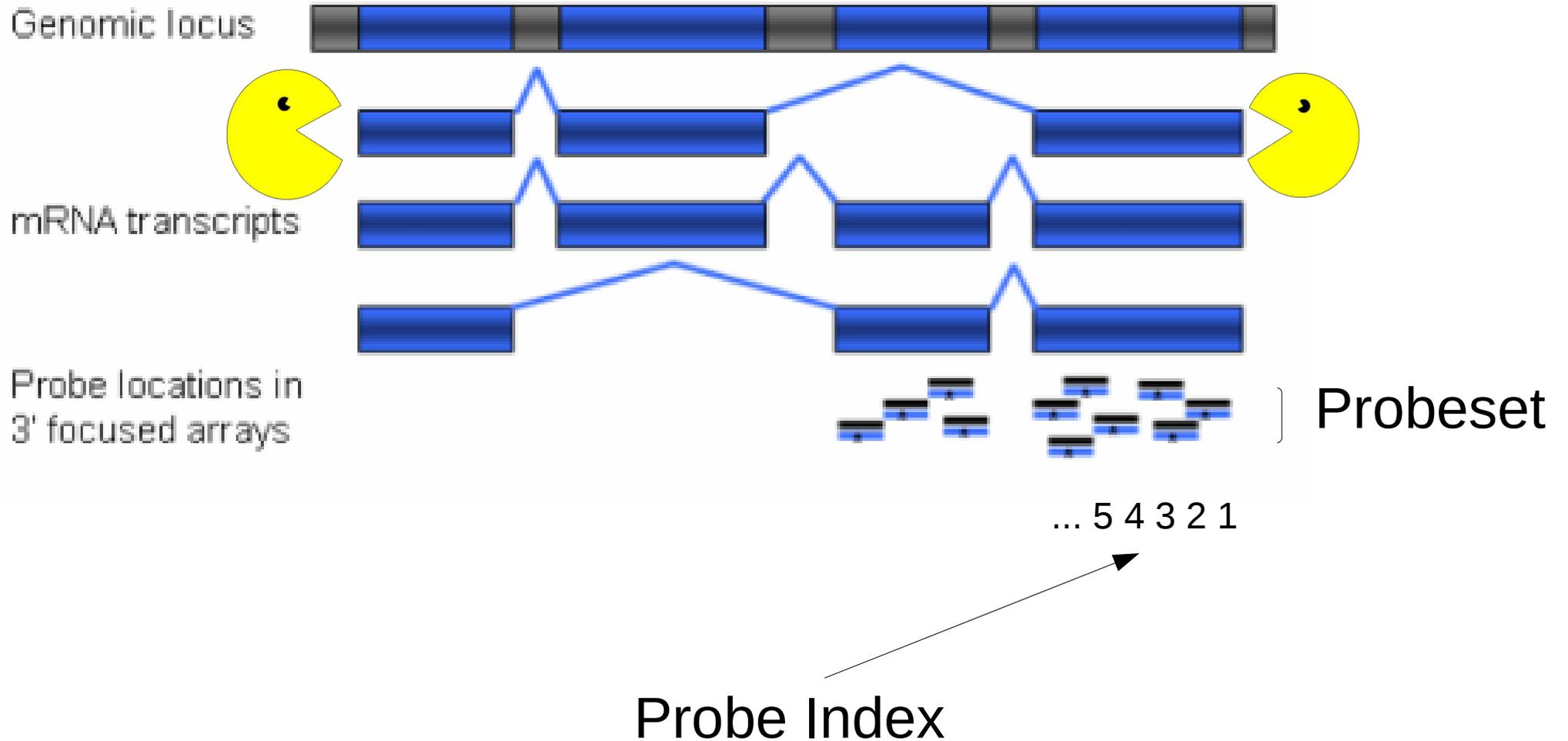
$$\tilde{\sigma}_i(b) = 0 \text{ for } b \neq GGG$$



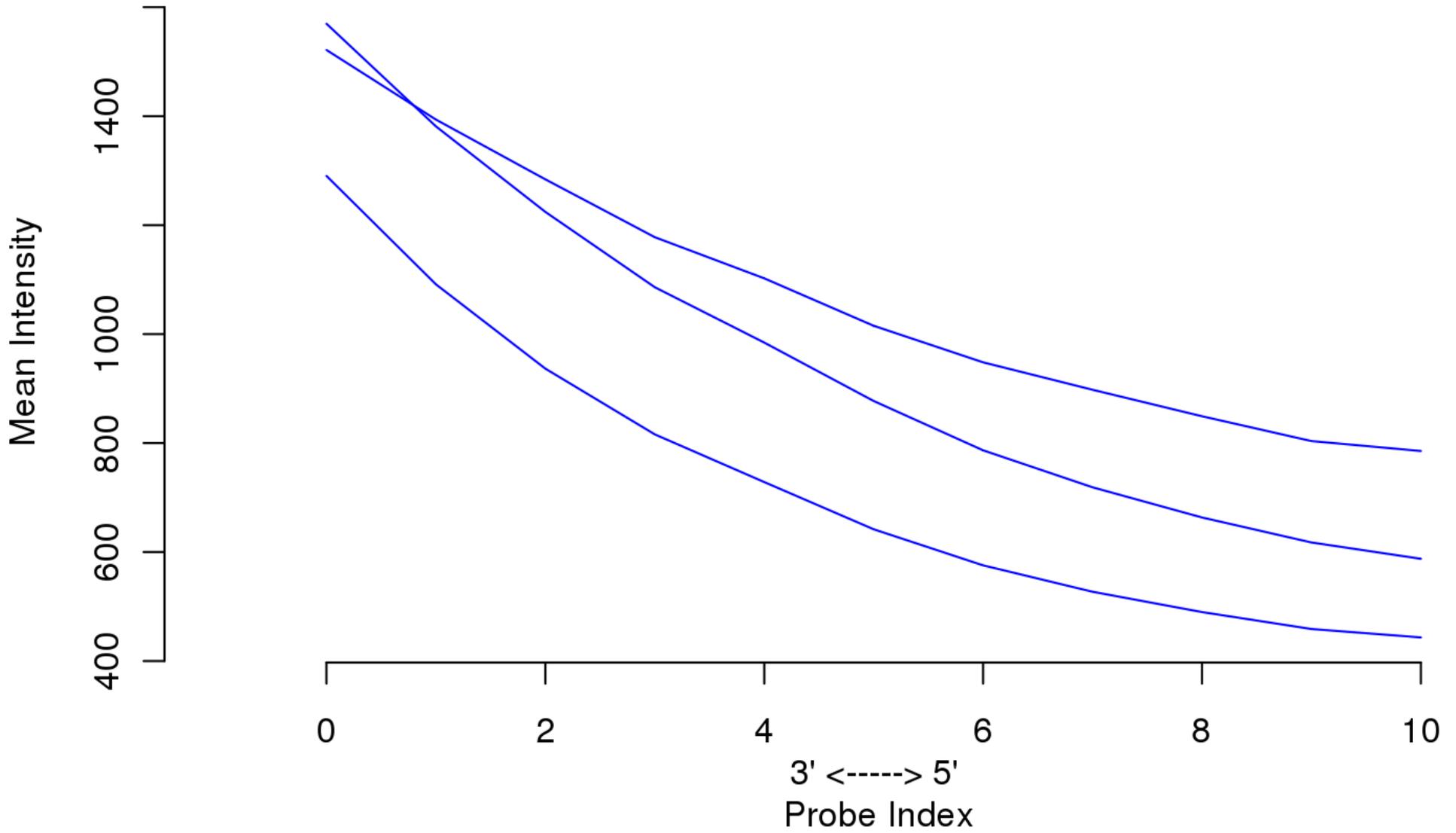
Sequence Correction

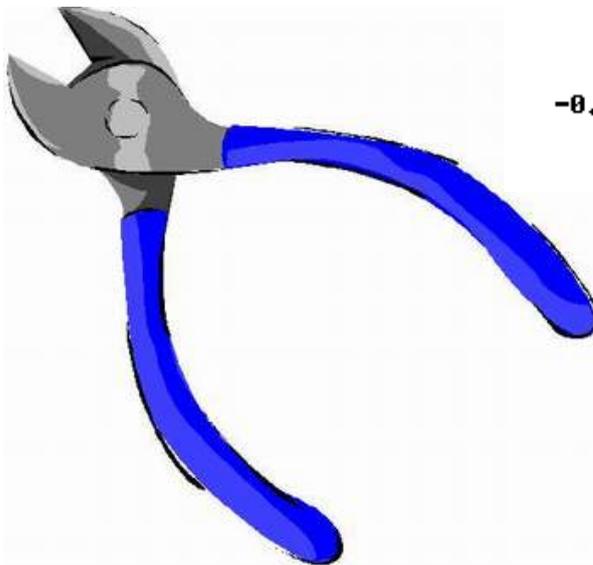
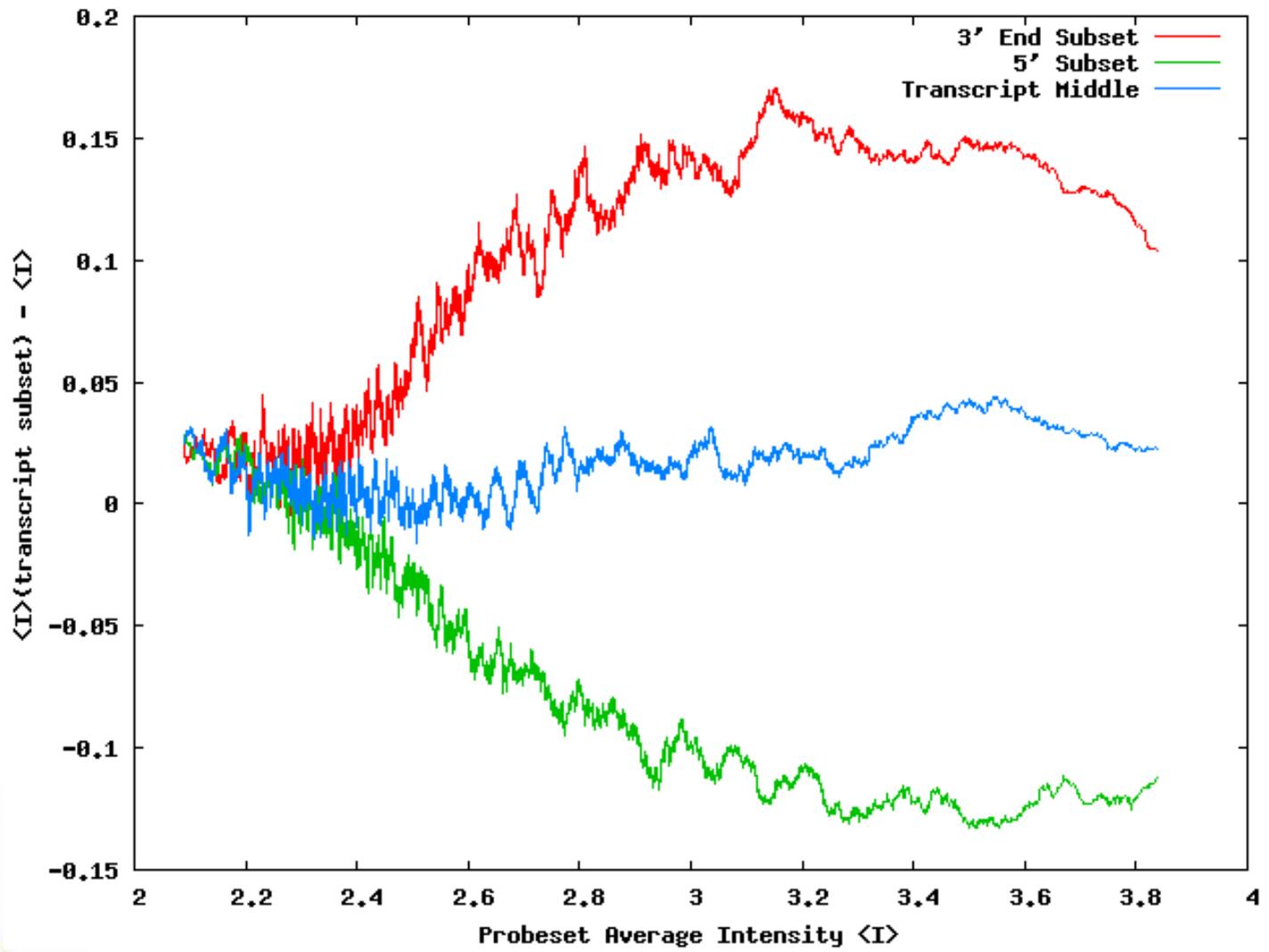


RNA Degradation



RNA degradation plot



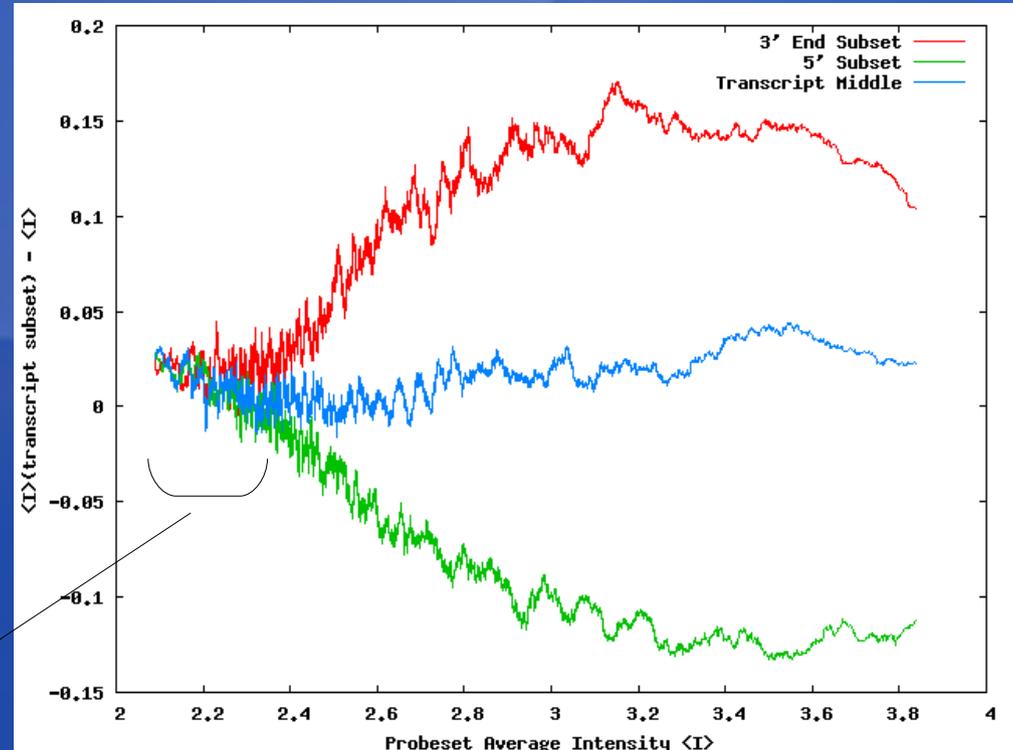
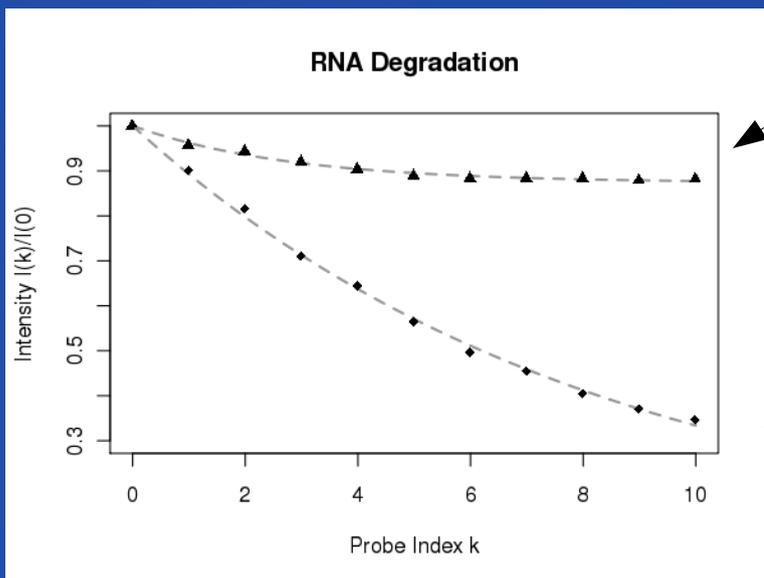


„Tongs“-Plot

RNA Degradation Depends on Probe Intensity

$$\langle I \rangle(k) = \langle I \rangle(0) * D(k)$$

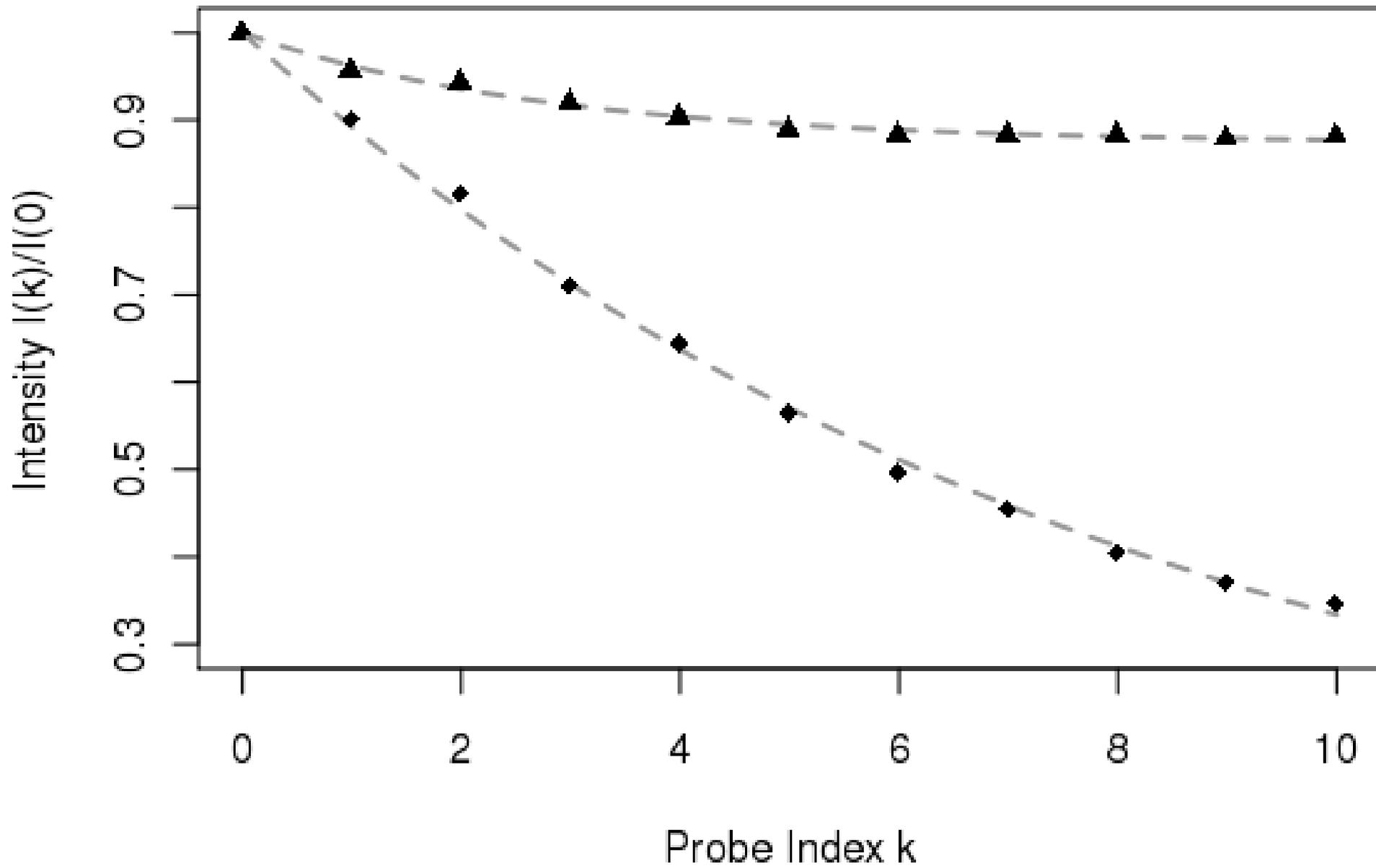
$$D(k) = (1 - d_1)e^{-d_2 k} + d_1$$



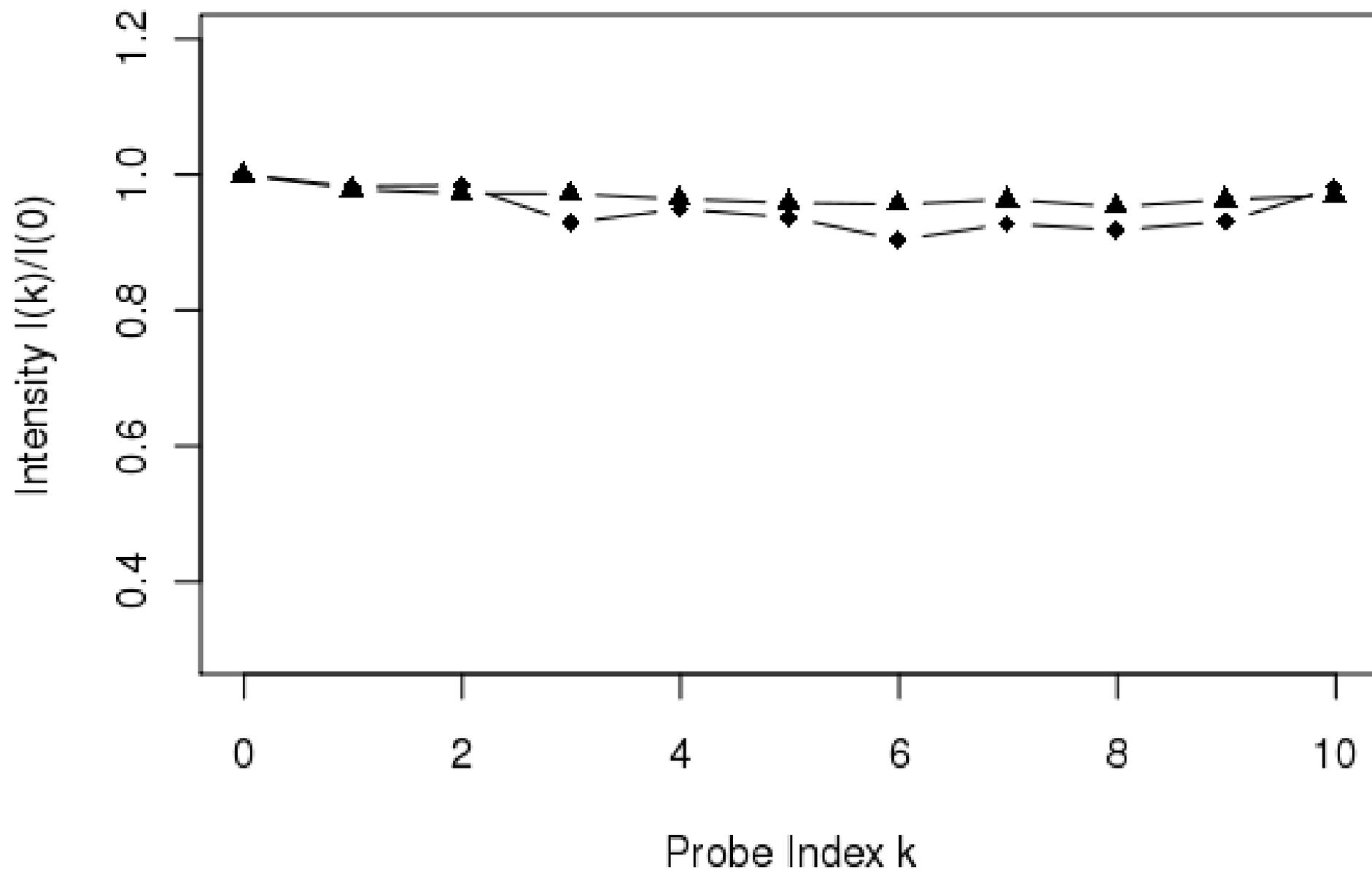
Non-specific binding

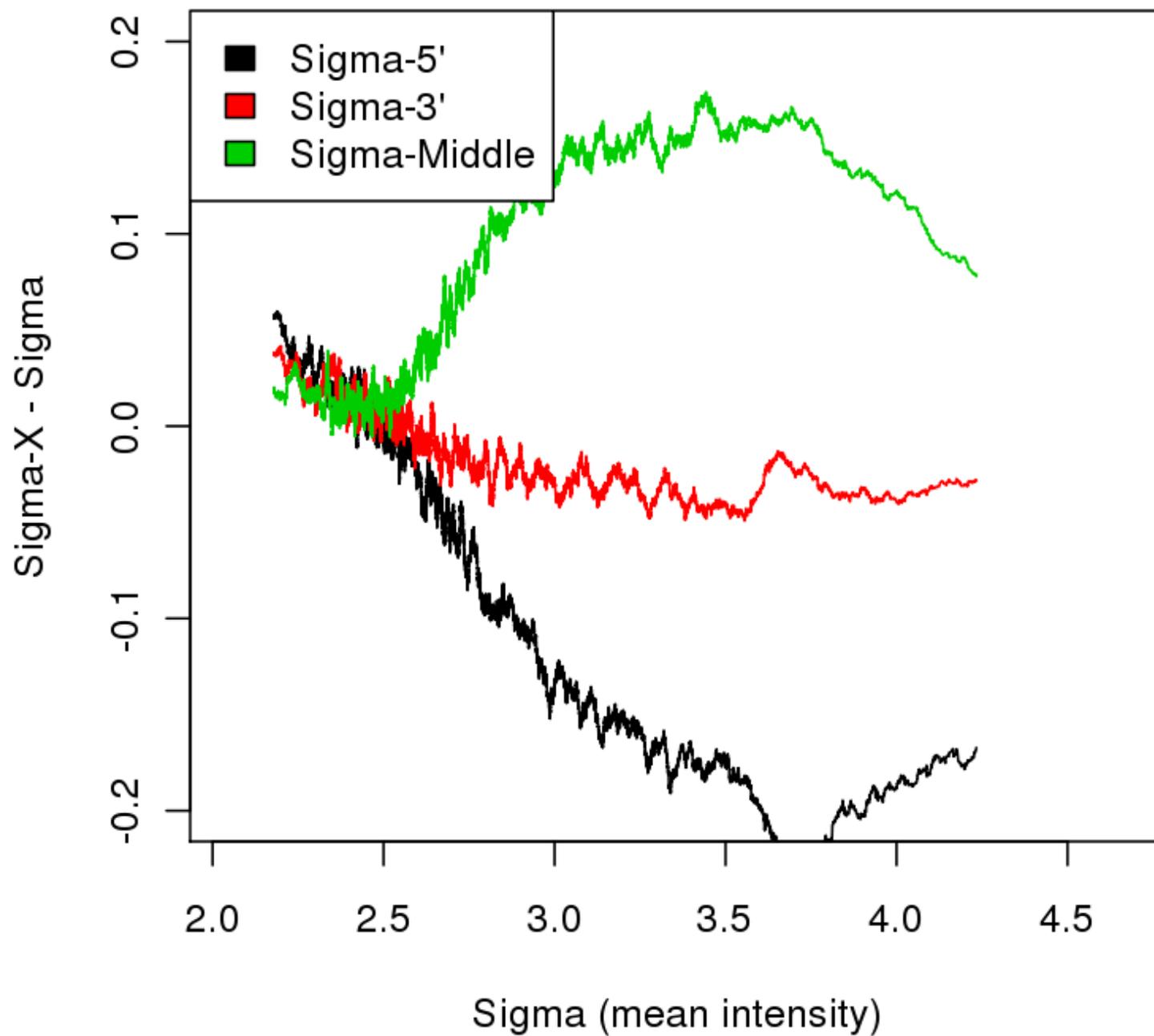
Specific binding

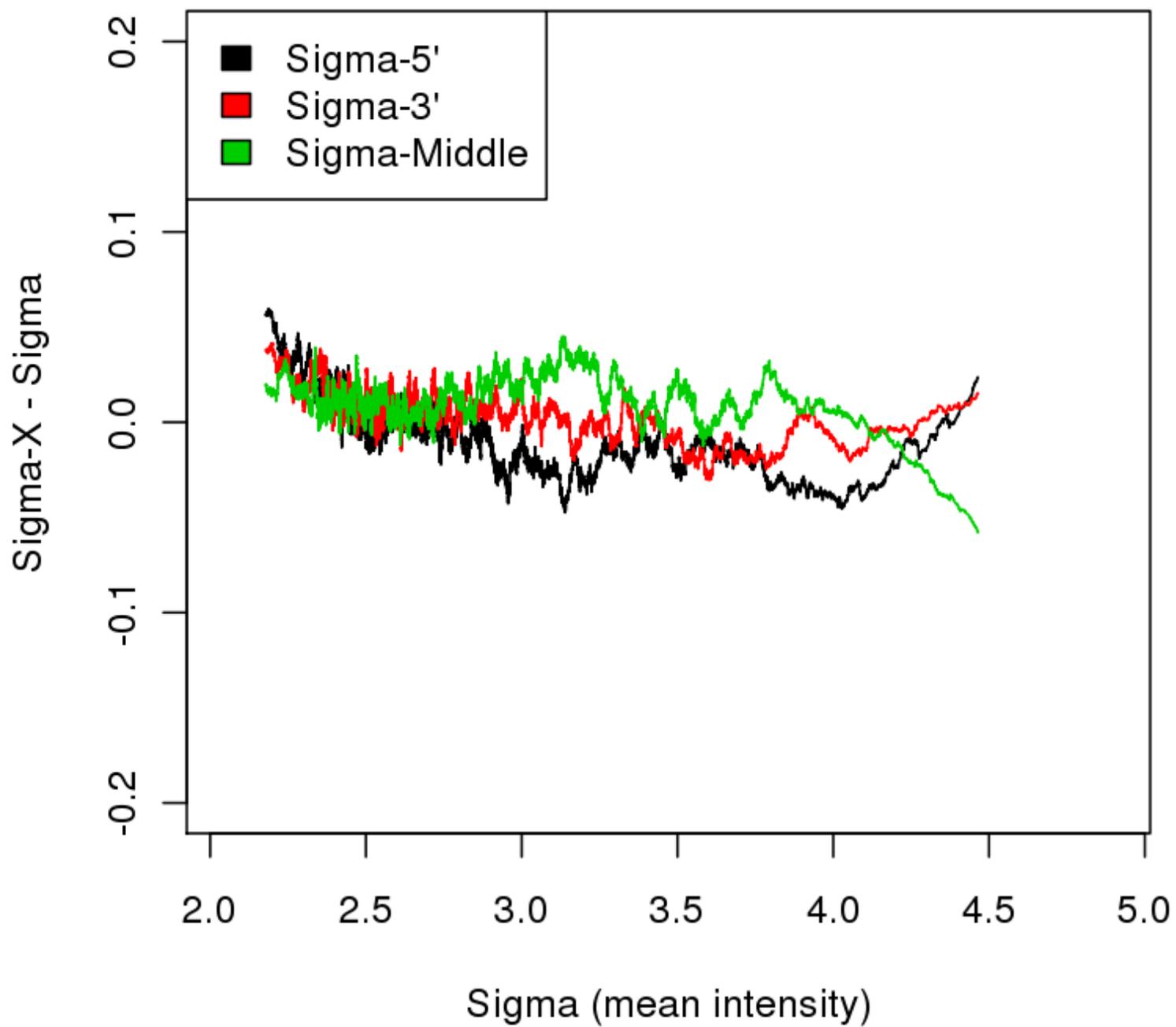
RNA Degradation



RNA Degradation (Corrected)





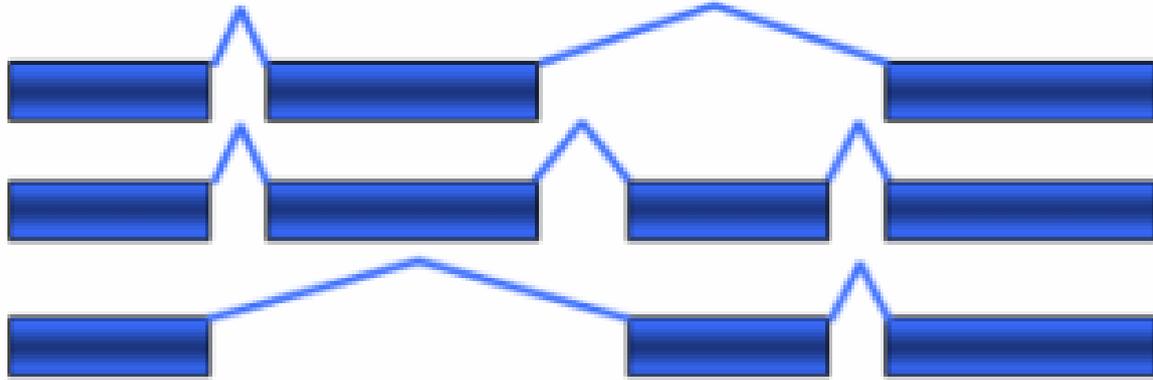


REFSEQ transcript sequences (ESTs)

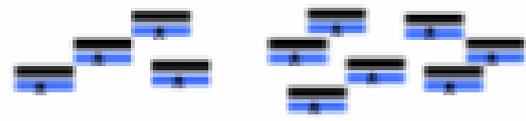
Genomic locus



mRNA transcripts



Probe locations in 3' focused arrays



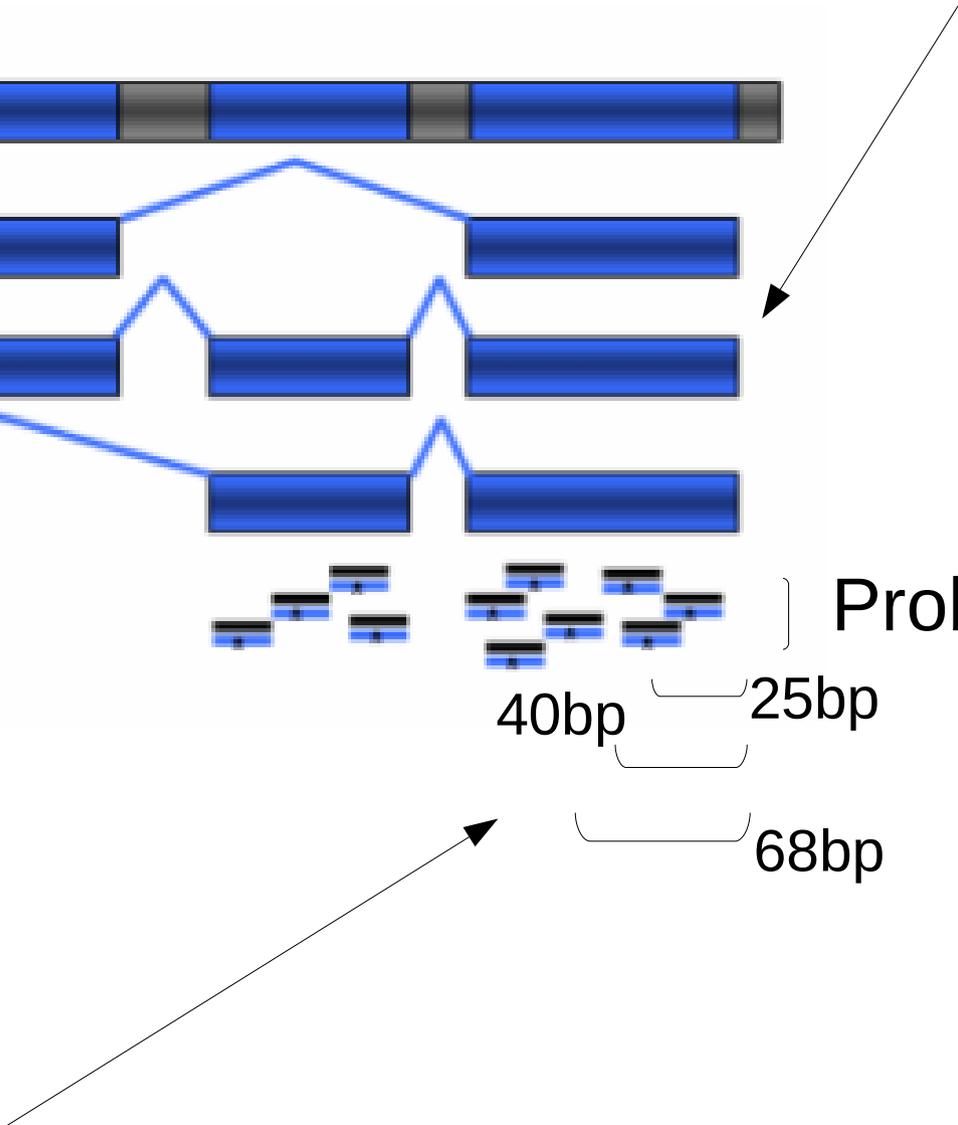
Probeset

40bp

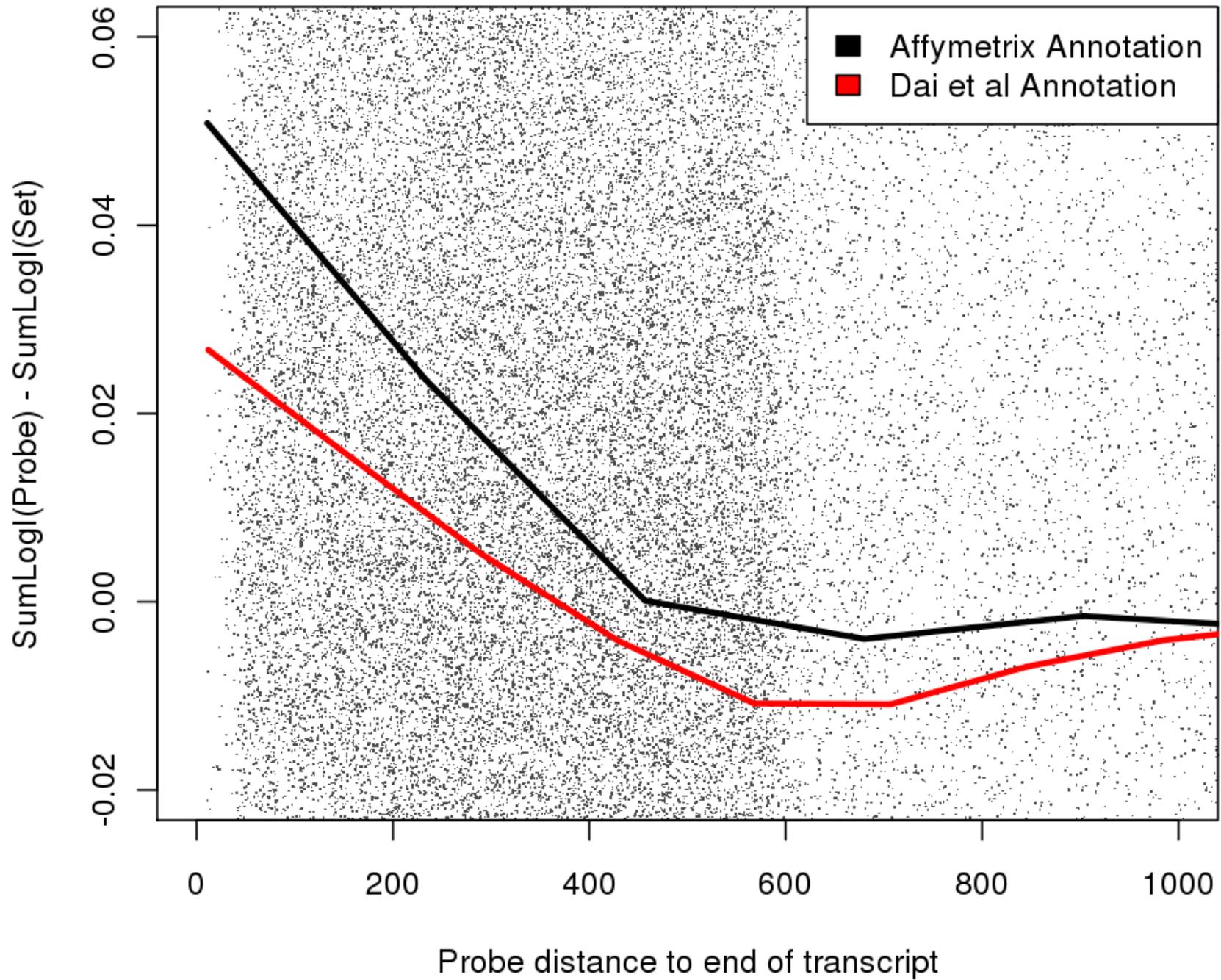
25bp

68bp

Absolute probe distance from 3' end (in base-pairs)



Absolut probe location bias



„All models are wrong, but some are
useful.“

(George Box)