# Improving the Layout of Oligonucleotide Microarrays: Pivot Partitioning

Sérgio A. de Carvalho Jr.[1,2,3]    Sven Rahmann[1,2]

[1]Algorithms and Statistics for Systems Biology, Genome Informatics, Technische Fakultät, Universität Bielefeld, Germany

[2]International NRW Graduate School in Bioinformatics and Genome Research

[3]Graduiertenkolleg Bioinformatik

Introduction
○○○○○

Conflict Index
○○○○

Pivot Partitioning
○○○○

Summary
○○

# Outline

Outline

Introduction
○●○○○

Conflict Index
○○○○

Pivot Partitioning
○○○○

Summary
○○

# High-Density Oligonucleotide Microarrays



Source: Affymetrix, Inc.

**Introduction**
○○●○○

Conflict Index
○○○○

Pivot Partitioning
○○○○

Summary
○○

## Probe Synthesis: Photolitographic Masks



Source: Affymetrix, Inc.

- Probes are synthesized on the chip in a series of steps
- Each step appends a particular nucleotide to selected regions
- Selection occurs by exposure to light

**Introduction**
○○○●○

Conflict Index
○○○○

Pivot Partitioning
○○○○

Summary
○○

## Deposition Sequence and Probe Embeddings

| $p_1$ ACT | $p_2$ CTG | $p_3$ GAT |
|---|---|---|
| $p_4$ TCC | $p_5$ GAC | $p_6$ GCC |
| $p_7$ TAC | $p_8$ CGT | $p_9$ AAT |

$\mathcal{S}$ = ACGTACGTACGT
$\mathcal{E}_1$ = ------------
$\mathcal{E}_2$ = ------------
$\mathcal{E}_3$ = ------------
$\mathcal{E}_4$ = ------------
$\mathcal{E}_5$ = ------------
$\mathcal{E}_6$ = ------------
$\mathcal{E}_7$ = ------------
$\mathcal{E}_8$ = ------------
$\mathcal{E}_9$ = ------------

**Introduction**
○○○●○

Conflict Index
○○○○

Pivot Partitioning
○○○○

Summary
○○

# Deposition Sequence and Probe Embeddings



$$\mathcal{S} = \text{ACGTACGTACGT}$$
$$\mathcal{E}_1 = \text{A-----------}$$
$$\mathcal{E}_2 = \text{-----------}$$
$$\mathcal{E}_3 = \text{-----------}$$
$$\mathcal{E}_4 = \text{-----------}$$
$$\mathcal{E}_5 = \text{-----------}$$
$$\mathcal{E}_6 = \text{-----------}$$
$$\mathcal{E}_7 = \text{-----------}$$
$$\mathcal{E}_8 = \text{-----------}$$
$$\mathcal{E}_9 = \text{A-----------}$$

**Introduction**
ooooeo

Conflict Index
oooo

Pivot Partitioning
oooo

Summary
oo

# Deposition Sequence and Probe Embeddings



$$\mathcal{S} = \text{ACGTACGTACGT}$$
$$\varepsilon_1 = \text{A-----------}$$
$$\varepsilon_2 = \text{-C----------}$$
$$\varepsilon_3 = \text{------------}$$
$$\varepsilon_4 = \text{------------}$$
$$\varepsilon_5 = \text{------------}$$
$$\varepsilon_6 = \text{------------}$$
$$\varepsilon_7 = \text{------------}$$
$$\varepsilon_8 = \text{-C----------}$$
$$\varepsilon_9 = \text{A-----------}$$

**Introduction**
oooo●o

Conflict Index
oooo

Pivot Partitioning
oooo

Summary
oo

## Deposition Sequence and Probe Embeddings



$\mathcal{S}$ = ACGTACGTACGT
$\mathcal{E}_1$ = A-----------
$\mathcal{E}_2$ = -C----------
$\mathcal{E}_3$ = --G---------
$\mathcal{E}_4$ = ------------
$\mathcal{E}_5$ = --G---------
$\mathcal{E}_6$ = --G---------
$\mathcal{E}_7$ = ------------
$\mathcal{E}_8$ = -CG---------
$\mathcal{E}_9$ = A-----------

**Introduction**
ooooo●o

Conflict Index
oooo

Pivot Partitioning
oooo

Summary
oo

# Deposition Sequence and Probe Embeddings



$$\mathcal{S} = \text{ACGTACGTACGT}$$
$$\varepsilon_1 = \text{A}-----------$$
$$\varepsilon_2 = -\text{C}----------$$
$$\varepsilon_3 = --\text{G}---------$$
$$\varepsilon_4 = ---\text{T}--------$$
$$\varepsilon_5 = --\text{G}---------$$
$$\varepsilon_6 = --\text{G}---------$$
$$\varepsilon_7 = ---\text{T}--------$$
$$\varepsilon_8 = -\text{CGT}--------$$
$$\varepsilon_9 = \text{A}-----------$$

**Introduction**
○○○●○

Conflict Index
○○○○

Pivot Partitioning
○○○○

Summary
○○

## Deposition Sequence and Probe Embeddings

| $p_1$ ACT | $p_2$ CTG | $p_3$ GAT |
|---|---|---|
| $p_4$ TCC | $p_5$ GAC | $p_6$ GCC |
| $p_7$ TAC | $p_8$ CGT | $p_9$ AAT |

$$
\begin{aligned}
S &= \text{ACGTACGTACGT} \\
\varepsilon_1 &= \text{A----C-T----} \\
\varepsilon_2 &= \text{-C-----T--G-} \\
\varepsilon_3 &= \text{--G-A--T----} \\
\varepsilon_4 &= \text{---T-C---C--} \\
\varepsilon_5 &= \text{--G-A----C--} \\
\varepsilon_6 &= \text{--G--C---C--} \\
\varepsilon_7 &= \text{---TAC------} \\
\varepsilon_8 &= \text{-CGT--------} \\
\varepsilon_9 &= \text{A---A------T}
\end{aligned}
$$

**Introduction**
○○○●○
Conflict Index
○○○○
Pivot Partitioning
○○○○
Summary
○○

## Deposition Sequence and Probe Embeddings

| $p_1$ | $p_2$ | $p_3$ |
|---|---|---|
| ACT | CTG | GAT |
| $p_4$ | $p_5$ | $p_6$ |
| TCC | GAC | GCC |
| $p_7$ | $p_8$ | $p_9$ |
| TAC | CGT | AAT |

$\mathcal{S}$ = ACGTACGTACGT
$\varepsilon_1$ = A----C-T----
$\varepsilon_2$ = -C-----T--G-
$\varepsilon_3$ = --G-A--T----
$\varepsilon_4$ = ---T-C---C--
$\varepsilon_5$ = --G-A----C--
$\varepsilon_6$ = --G--C---C--
$\varepsilon_7$ = ---TAC------
$\varepsilon_8$ = -CGT--------
$\varepsilon_9$ = A---A------T
$\varepsilon'_9$ = A---A--T----

Left-most embedding!

**Introduction**
○○○●○

Conflict Index
○○○○

Pivot Partitioning
○○○○

Summary
○○

## Deposition Sequence and Probe Embeddings

| $p_1$ | $p_2$ | $p_3$ |
|-------|-------|-------|
| ACT | CTG | GAT |
| $p_4$ | $p_5$ | $p_6$ |
| TCC | GAC | GCC |
| $p_7$ | $p_8$ | $p_9$ |
| TAC | CGT | AAT |

$$
\begin{aligned}
\mathcal{S} &= \text{ACGTACGTACGT} \\
\mathcal{E}_1 &= 100001010000 \\
\mathcal{E}_2 &= 010000010010 \\
\mathcal{E}_3 &= 001010010000 \\
\mathcal{E}_4 &= 000101000100 \\
\mathcal{E}_5 &= 001010000100 \\
\mathcal{E}_6 &= 001001000100 \\
\mathcal{E}_7 &= 000111000000 \\
\mathcal{E}_8 &= 011100000000 \\
\mathcal{E}_9 &= 100010000001 \\
\mathcal{E}'_9 &= 100010010000
\end{aligned}
$$

# Problem: Unintended Illumination



- Untargeted spots can be accidentally activated
  - Diffraction of light
  - Internal reflection
- Production of defective probes
- More likely near the borders between masked and unmasked spots: border conflict

## Border Length Minimization Problem (Hannenhalli et al., 2002)

- Find arrangement (and embeddings) with minimum number of border conflicts

Introduction
○○○○○

Conflict Index
●○○○

Pivot Partitioning
○○○○

Summary
○○

# Outline

Introduction
○○○○○

Conflict Index
○●○○○

Pivot Partitioning
○○○○

Summary
○○

# Conflict Index: Motivation

- Border Length measures the quality of a particular mask
  - We are more interested in a per-probe measure
- Practical considerations need to be taken into account:
  - a) Stray light might activate probes that are as far as three cells away from the targeted spot
  - b) Imperfections produced in the middle of a probe are more harmful than in its extremities

## Conflict Index of a probe $p$

$$\mathcal{C}(p) := \sum_{t=1}^{T} \Big( \omega(p, t) \sum_{p'} \delta(p, p', t) \Big),$$

where $\delta(p, p', t)$ are distance-dependent weights (a) and $\omega(p, t)$ are position-dependent weights (b) defined as follows.

Introduction
○○○○○

Conflict Index
○○●○

Pivot Partitioning
○○○○

Summary
○○

## Conflict Index: Definition

Conflict Index of a probe $p$

$$\mathcal{C}(p) := \sum_{t=1}^{T}\Big(\omega(p, t) \sum_{p'} \delta(p, p', t)\Big)$$

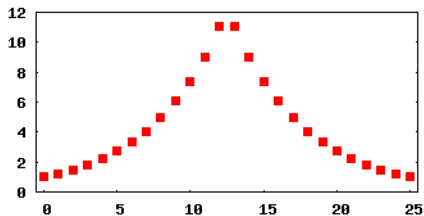| 0.06 | 0.08 | 0.10 | 0.11 | 0.10 | 0.08 | 0.06 |
|------|------|------|------|------|------|------|
| 0.08 | 0.13 | 0.20 | 0.25 | 0.20 | 0.13 | 0.08 |
| 0.10 | 0.20 | 0.50 | 1.00 | 0.50 | 0.20 | 0.10 |
| 0.11 | 0.25 | 1.00 | $p$  | 1.00 | 0.25 | 0.11 |
| 0.10 | 0.20 | 0.50 | 1.00 | 0.50 | 0.20 | 0.10 |
| 0.08 | 0.13 | 0.20 | 0.25 | 0.20 | 0.13 | 0.08 |
| 0.06 | 0.08 | 0.10 | 0.11 | 0.10 | 0.08 | 0.06 |

a) Distance-dependent weights $\delta(p, p', t)$

$$\delta(p, p', t) := \left\{ \begin{array}{ll} (d(p, p'))^{-2} & \text{if } p' \text{ is unmasked at step } t, \\ 0 & \text{otherwise,} \end{array} \right.$$

where $d(p, p')$ is the Euclidean distance between the spots of $p$ and $p'$.

Introduction
○○○○○

Conflict Index
○○○●

Pivot Partitioning
○○○○

Summary
○○

# Conflict Index: Definition

### Conflict Index of a probe $p$

$$\mathcal{C}(p) := \sum_{t=1}^{T} \Big( \omega(p, t) \sum_{p'} \delta(p, p', t) \Big)$$



### b) Position-dependent weights $\omega(p, t)$

$$\omega(p, t) := \left\{ \begin{array}{ll} c \cdot \exp\left(\theta \cdot \lambda(p, t)\right) & \text{if } p \text{ is masked at step } t, \\ 0 & \text{otherwise,} \end{array} \right.$$

where $c > 0$ and $\theta > 0$ are constants,

$$\lambda(p, t) := 1 + \min(b_{p,t}, \ell_p - b_{p,t}),$$

$b_{p,t}$ denotes the number of nucleotides synthesized up to and including step $t$, and $\ell_p$ is the length of probe $p$.

# Outline

Introduction
00000

Conflict Index
0000

Pivot Partitioning
0●00

Summary
00

# Previous Work: Place and Re-embed

- The microarray layout problem has been traditionally approached in two phases:
  1) Placement of probes given a fixed embedding
  2) Re-embedding of probes given a fixed placement

## Placement: Row-epitaxial (Kahng *et al.*, 2003)

- Essentially greedy
- Spots are filled in a pre-defined order
  - Select probe from a list $Q$ such that conflicts with filled spots are minimized
- Restrict the maximum size of $Q$

# Previous Work: Place and Re-embed

## Re-embedding: several algorithms

- All based on the Optimum Single Probe Embedding (OSPE)
- OSPE re-embed a probe optimally in regards to its neighbors
- Difference is in the order in which re-embeddings take place

$\mathcal{S}$ = ACGTACGTACGT
$\varepsilon_1$ = A----C-T----
$\varepsilon_2$ = -C-----T--G-
$\varepsilon_3$ = --G-A--T----
$\varepsilon_4$ = ---T-C---C--
$\varepsilon_5$ = ????????????
$\varepsilon_6$ = --G--C---C--
$\varepsilon_7$ = ---TAC------
$\varepsilon_8$ = -CGT--------
$\varepsilon_9$ = A---A------T

## Optimum Single Probe Embedding (OSPE)

- Dynamic Programming
- Originally developed for border length minimization
- Now extended for conflict index minimization

Introduction
ooooo

Conflict Index
oooo

Pivot Partitioning
oooo●

Summary
oo

# Previous Work: Partitioning

More recently, a partitioning algorithm was proposed

- Divide the problem into smaller sub-problems
- Each sub-problem is treated as a separate placement
- Reduce run-time; may improve placement

Partitioning: Centroid-based Quadrisection (Kahng *et al.*, 2003)

- To do...

Introduction
○○○○○

Conflict Index
○○○○

Pivot Partitioning
○○○○

Summary
●○

# Summary

- Conflict Index: new model for evaluating microarray layouts
- Pivot Partitioning: new partitioning algorithm
  - Faster and better selection of pivots
  - Improved assignment of probes to regions
  - First to combine placement and re-embedding

Introduction
○○○○○

Conflict Index
○○○○

Pivot Partitioning
○○○○

**Summary**
○●

# Thanks!



- Prof. Dr. Jens Stoye
- Prof. Dr. Robert Giegerich
- AG Genominformatik
- Graduiertenkolleg Bioinformatik
- Graduate School in Bioinformatics and Genome Research
- ...and thank you for your attention!

More info on

http://gi.cebitec.uni-bielefeld.de/assb/chiplayout