

PRIP-TR-58

December 9, 2000

## Fully Automatic Grid Fitting for Genetic Spot Array Images Containing Guide Spots

*Norbert Brändle, Hilmar Lapp<sup>1</sup> and Horst Bischof*

### Abstract

In the domain of biotechnology array-based methods are used to gain rapid access to genetic information based on the signals of the individual array elements (spots). For an automated analysis of the spots it is necessary to fit a grid to the spots in the digital image in order to represent the array distortions that may occur in the course of the experiment. In order to make the grid fitting problem tractable in a certain class of experiments spot arrays contain a sub-grid of guide spots with a known signal characteristic. We present an automatic grid fitting method for spot array images containing guide spots. Our approach uses simple image processing methods and takes into account prior knowledge inherent in the imaging process.

---

<sup>1</sup>Hilmar Lapp is with Novartis Research Institute, Inflammatory Diseases, Brunner Str. 59, A-1230 Vienna, AUSTRIA. E-mail: [Hilmar.Lapp@pharma.novartis.com](mailto:Hilmar.Lapp@pharma.novartis.com)

# Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Notation and Definitions</b>	<b>6</b>
2.1	Spot Array Image Representation . . . . .	6
2.2	Spot Array Representation . . . . .	7
2.3	Prior Knowledge . . . . .	9
<b>3</b>	<b>Guide Spot Detection</b>	<b>11</b>
3.1	Spot Detection with Matched Filter . . . . .	11
3.2	Guide Spot Location Amplification (GSLA) . . . . .	12
3.3	Local Maximum Search . . . . .	14
3.3.1	Region-of-Interest and Prior Guide Spot Locations . . . . .	15
3.3.2	Maximum Search Algorithm . . . . .	17
3.4	Summary . . . . .	21
<b>4</b>	<b>Grid Fitting</b>	<b>22</b>
4.1	Alignment of Detected Guide Spot Locations . . . . .	23
4.1.1	Global Rotation Estimation . . . . .	23
4.1.2	Global Translation Estimation . . . . .	26
4.1.3	Alignment of locations to grid nodes . . . . .	30
4.2	Consistent Spot Grid . . . . .	31
4.2.1	Parameterization of the Guide Spot Grid . . . . .	31
4.2.2	Robust Grid Parameter Fitting . . . . .	32
4.2.3	Field Parameter Correction . . . . .	33
4.2.4	Abortion Criterion . . . . .	33
4.2.5	Location initialization of regular spots . . . . .	34
4.3	Summary . . . . .	34
<b>5</b>	<b>Experimental Results</b>	<b>35</b>
5.1	Visual Examples . . . . .	35
5.2	Evaluation of Image Quality versus Grid Fitting Success . . . . .	41
<b>6</b>	<b>Conclusion and Outlook</b>	<b>43</b>

# 1 Introduction

Rapid access to genetic information is central to the revolution taking place in molecular genetics [6, 10]. Biological systems read, store and modify genetic information by molecular recognition. Because each DNA strand carries with it the capacity to recognize a uniquely complementary sequence through base pairing ( $A \leftrightarrow T$ ,  $C \leftrightarrow G$ ) [11], the process of recognition, or *hybridization*, is highly parallel, as every nucleotide in a large sequence can in principle be queried at the same time. Thus, hybridization can be used to efficiently analyze large amounts of nucleotide sequence. The primary approaches include array-based technologies that can identify specific expressed gene products on high density formats, including filters, microscope slides and micro-chips [6].

Common to all array-based approaches is the necessity to analyze digital images of the array. The images look like a grid with small, bright round spots. Figure 1 shows a typical image generated in the course of an oligonucleotide fingerprint (ONF) experiment [13][12]: The intensity of every spot corresponds to the amount of label remaining after hybridizing a liquid containing the labelled probes and subsequently washing off probe not bound to the genetic material. For details about the physical imaging process refer to [7]. The grid of the spots is generated by a robot in the following way: A rectangular piece of synthetic material, the so-called *membrane*, is mounted on a fixed plane. The membrane is divided into rectangular *fields*. A robot arm carries a matrix of needles (see Fig. 2a). Every needle has a small amount of liquid stuck to it containing the genetic material to be spotted.

The DNA material is spotted onto the membrane in subsequent spotting cycles. In each cycle, the spotting robot first loads from one well-plate (see Fig. 2b) the needle matrix with liquid containing the material to be spotted. Afterwards, it sets the matrix down onto the membrane at a certain position within a field, thereby transferring the liquid to spots on the membrane. For each field, this cycle is repeated with a set of defined offsets in  $x$ - and  $y$ -direction. As the largest offset is smaller than the distance between two needles, this procedure yields a spot grid of higher density than the needle matrix. The spots in one field originating from the same needle form a pattern of a rectangular *block*.

The ultimate image analysis goal is to automatically assign a quantity to every array element (spot) giving information about the hybridization signal (*quantification*). For a successful quantification of the hybridization signals it is necessary to assign to every array element a location in the spot array image. We call this location assignment *grid fitting*. The grid fitting has to cope with the following problems:

- **Distortions:** The spot grid in the image is only approximatively regular and rectangular due to anisotropic shrinking and expansion of the membrane, inaccuracies of the spotting robot, optical distortions, bent needles and other factors.
- **Rotations:** The grid need not be aligned with the image coordinates because the membrane is put manually into the imaging device.
- **Outliers:** Not every node of the grid necessarily contains luminescent spots because the intensity of the hybridization signal to be measured can be low or even zero.

When a large number of outliers is expected from the hybridization experiment (in ONF images for example only 5–20% of all spots can be expected to hybridize) so-called *guide spots* are

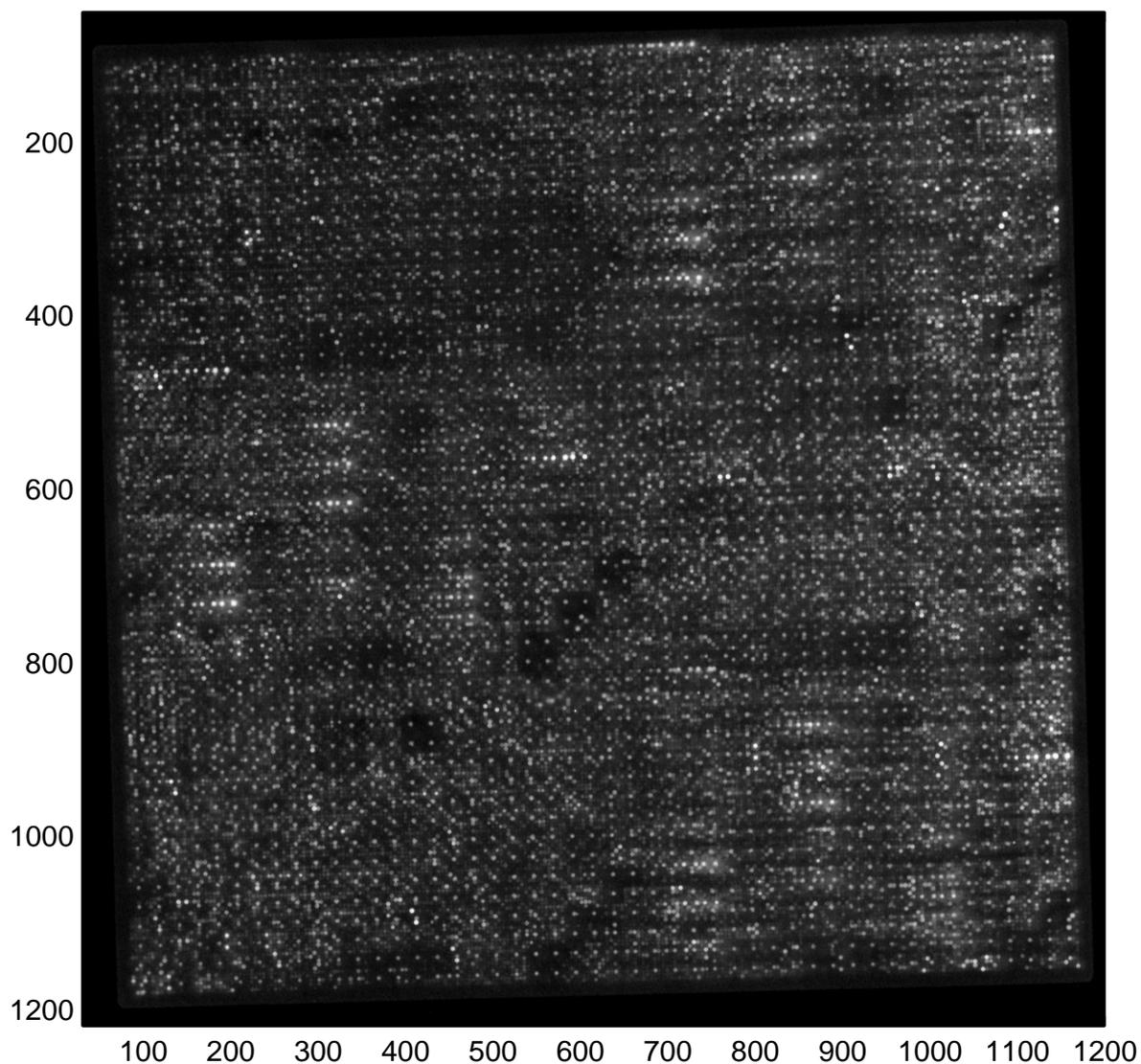
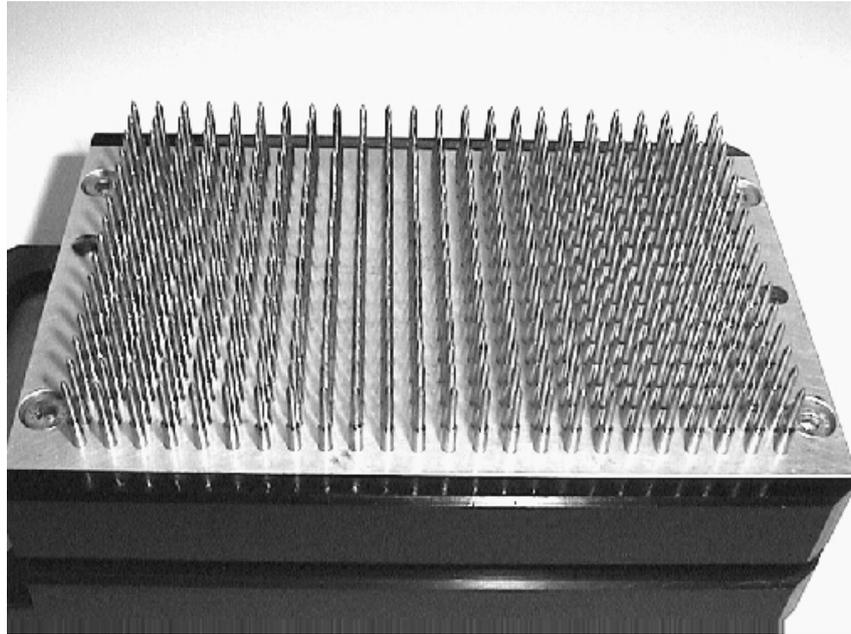
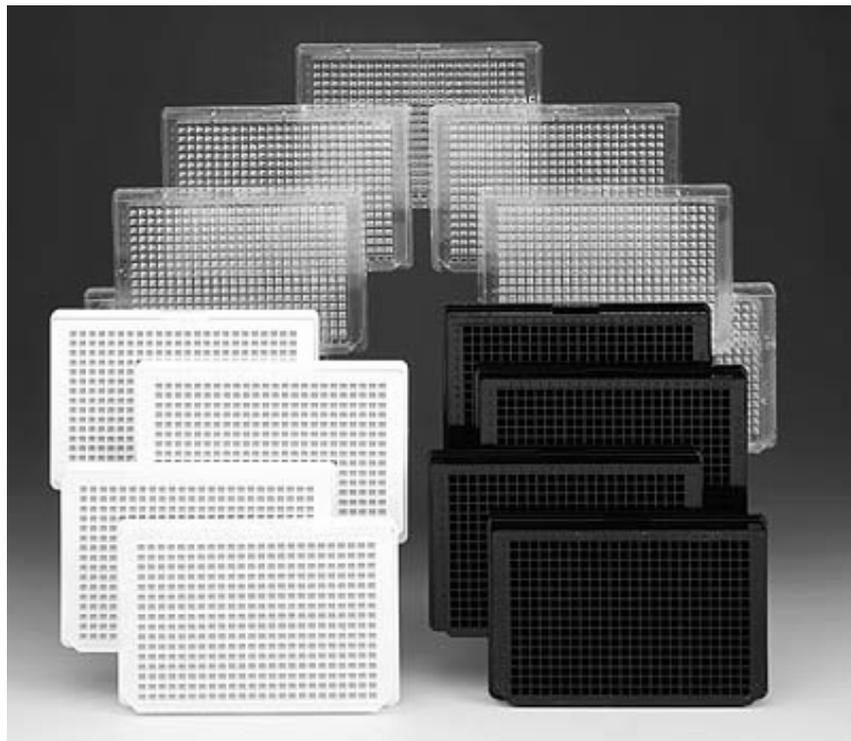


Figure 1: Oligonucleotide Fingerprint (ONF) Image. The intensity of every spot corresponds to the amount of label remaining after hybridizing a liquid containing the labelled probes and subsequently washing off probe not bound to the genetic material. The ultimate image analysis goal is to automatically assign a quantity to every array element (spot) giving information about the hybridization signal (*quantification*). For a successful quantification of the hybridization signals it is necessary to assign to every grid node an image location (grid fitting).



(a) Needle matrix



(b) Well-plates

Figure 2: (a) A needle matrix is carried by a robot arm. Every needle has a small amount of liquid stucked to it containing the genetic material to be spotted. The needle matrix is set down onto a membrane at a certain position thereby transferring the liquid to spots on the membrane. (b)The genetic material is spotted onto the membrane in subsequent spotting cycles. In each cycle, the spotting robot first loads from one well-plate the needle matrix with liquid containing the material to be spotted.

used: A spot at a certain position within a block contains DNA material which has always a strong hybridization signal irrespective of the hybridization probe.

One solution of the grid fitting problem was already presented by Hartelius [8] in his Ph.D. thesis. It involves an image rotation in order to align the grid with the image coordinates, block finding and spot finding via Markov Random Fields (MRF) and Simulated Annealing [17]. His algorithm is computationally demanding and is semi-automatic since the user must provide the locations of the corner nodes of the grid. In this report we present a fully automatic grid fitting approach for spot images containing guide spots. Our approach takes into account different constraints inherent in the imaging process and involves simple image processing operators. Quantification is a problem in its own right and is described in [1] and [4].

This report is structured as follows: Section 2 introduces the notations and definitions necessary to describe the grid fitting in a formal manner. Section 3 deals with the detection of potential guide spot locations. Section 4 describes the mapping of the potential guide spot locations to the nodes of the guide spot grid. Section 5 presents experimental results for images of different quality originating from different hybridization experiments.

## 2 Notation and Definitions

This section introduces the notation and definitions of the main modeling concepts of the spot image and the spot array. The general notation of sets is shown in table 1.

Not.	Description
$\mathcal{A}, \mathcal{B}, \mathcal{C}, \dots$	general sets
$\text{card}(\mathcal{A}), \text{card}(\mathcal{B}), \dots$	cardinality of sets
$\emptyset$	empty set
$\mathbb{N}$	set of cardinal numbers (inclusive zero)
$\mathbb{Z}$	set of integer numbers
$\mathbb{R}$	set of real numbers
$\mathbb{R}^2$	two-dimensional plane.
$\mathbb{R}^{M \times N}$	set of $M \times N$ matrices over $\mathbb{R}$ $M, N \geq 1$
$\mathbf{A}, \mathbf{B}, \mathbf{C}, \dots \in \mathbb{R}^{M \times N}$	$M \times N$ matrices over $\mathbb{R}$
$\mathbf{A}^T, \mathbf{B}^T, \dots \in \mathbb{R}^{N \times M}$	transpose of $M \times N$ matrices over $\mathbb{R}$
$\mathbf{a}, \mathbf{b}, \mathbf{c}, \dots \in \mathbb{R}^d$	vectors (points)
$\circ(\cdot)$	Rounding operator to the nearest integer

Table 1: General notations used in this report.

### 2.1 Spot Array Image Representation

This section deals with the representation of the scanned spot array images.

#### Definition 1 (Spot array image)

A (scanned)  $M \times N$  spot array image is represented as a matrix  $\mathbf{S} \in \mathbb{Z}^{M \times N}$  with pixel coordinates  $(m, n)$  and pixel intensities  $\mathbf{S}[m, n]$ :

$$\mathbf{S} = \begin{bmatrix} \mathbf{S}[1, 1] & \mathbf{S}[1, 2] & \dots & \mathbf{S}[1, N] \\ \mathbf{S}[2, 1] & \mathbf{S}[2, 2] & \dots & \mathbf{S}[2, N] \\ \dots & \dots & \dots & \dots \\ \mathbf{S}[M, 1] & \mathbf{S}[M, 2] & \dots & \mathbf{S}[M, N] \end{bmatrix} \quad (1)$$

Two alternative Cartesian coordinate systems are used:

- a *spatial coordinate system*  $(x, y)$  with the origin at the upper left corner with  $x$  increasing to the right and  $y$  increasing downward.
- a *central coordinate system*  $(x, y)$  with the origin at the image center with  $x$  increasing to right and  $y$  increasing upward.

## 2.2 Spot Array Representation

In the following we define the representation of the spot array in the spot array image. The different dimensions of a spot array are illustrated in Fig 3.

### Definition 2 (Grid)

A *grid*  $\mathcal{G}$  is a set of nodes in  $\{1 \dots I_G\} \times \{1 \dots J_G\}$ , with  $I_G$  as the number of grid rows and  $J_G$  as the number of grid columns.

The grid row extraction relation  $R_G$  extracts  $J_G$  grid nodes belonging to one row of a grid  $\mathcal{G}$ :

$$R_G(\{1 \dots I_G\}) = \{(i, j) \mid (i, j) \in \mathcal{G} \wedge i = k\} \quad \forall k \in \{1 \dots I_G\} \quad (2)$$

Similarly, the column extraction relation  $C_G$  extracts  $I_G$  grid nodes belonging to one column of a grid  $\mathcal{G}$ :

$$C_G(\{1 \dots J_G\}) = \{(i, j) \mid (i, j) \in \mathcal{G} \wedge j = l\} \quad \forall l \in \{1 \dots J_G\}. \quad (3)$$

For a successful quantification of the hybridization signals it is necessary to assign to every node  $(i, j) \in \mathcal{G}$  an element of a set of image locations with sub-pixel accuracy.

### Definition 3 (Grid Fitting)

*Grid fitting* is the location assignment function  $L : \mathcal{G} \times \mathbf{S} \rightarrow \mathcal{L}$ , where  $\mathcal{L} = \{(x, y) \mid x, y \in \mathbb{R}, x, y \geq 0\}$ . The location of a grid node  $(i, j) \in \mathcal{G}$  will sometimes be simply denoted as  $L((i, j))$ .

In order to cope with inaccuracies of the spotting robot the grid is usually divided into subunits to represent the nature of the spotting cycles.

### Definition 4 (Field)

A *field*  $\mathcal{F}_{pq} \subset \mathcal{G}$  is a subunit of a grid  $\mathcal{G}$  and is a set of nodes in  $\{1 \dots I_F\} \times \{1 \dots J_F\}$  with  $I_F$  as the number of field rows and  $J_F$  as the number of field columns. A grid  $\mathcal{G}$  is partitioned into  $F_1 * F_1$  fields such that

$$\bigcup_{p=1}^{F_1} \bigcup_{q=1}^{F_1} L(\mathcal{F}_{pq}, \mathbf{S}) = L(\mathcal{G}, \mathbf{S}) \quad (4)$$

and

$$L(\mathcal{F}_{pq}, \mathbf{S}) \cap L(\mathcal{F}_{rs}, \mathbf{S}) = \emptyset, \quad \forall (p, q) \neq (r, s). \quad (5)$$

The row extraction and column extraction relations  $R_{\mathcal{F}}$  and  $C_{\mathcal{F}}$  are defined similarly to Eqns (2) and (3).

### Definition 5 (Block)

A *block*  $\mathcal{B}_{ij} \subset \mathcal{F}$  is a subunit of a field  $\mathcal{F}$  and is a set of nodes in  $\{1 \dots I_B\} \times \{1 \dots J_B\}$  with  $I_B$  as the number of block rows and  $J_B$  as the number of block columns. A field  $\mathcal{F}$  is partitioned into  $I_w * J_w$  blocks such that

$$\bigcup_{i=1}^{I_w} \bigcup_{j=1}^{J_w} L(\mathcal{B}_{ij}, \mathbf{S}) = L(\mathcal{F}, \mathbf{S}) \quad (6)$$

and

$$L(\mathcal{B}_{ij}, \mathbf{S}) \cap L(\mathcal{B}_{kl}, \mathbf{S}) = \emptyset \quad \forall (i, j) \neq (k, l). \quad (7)$$

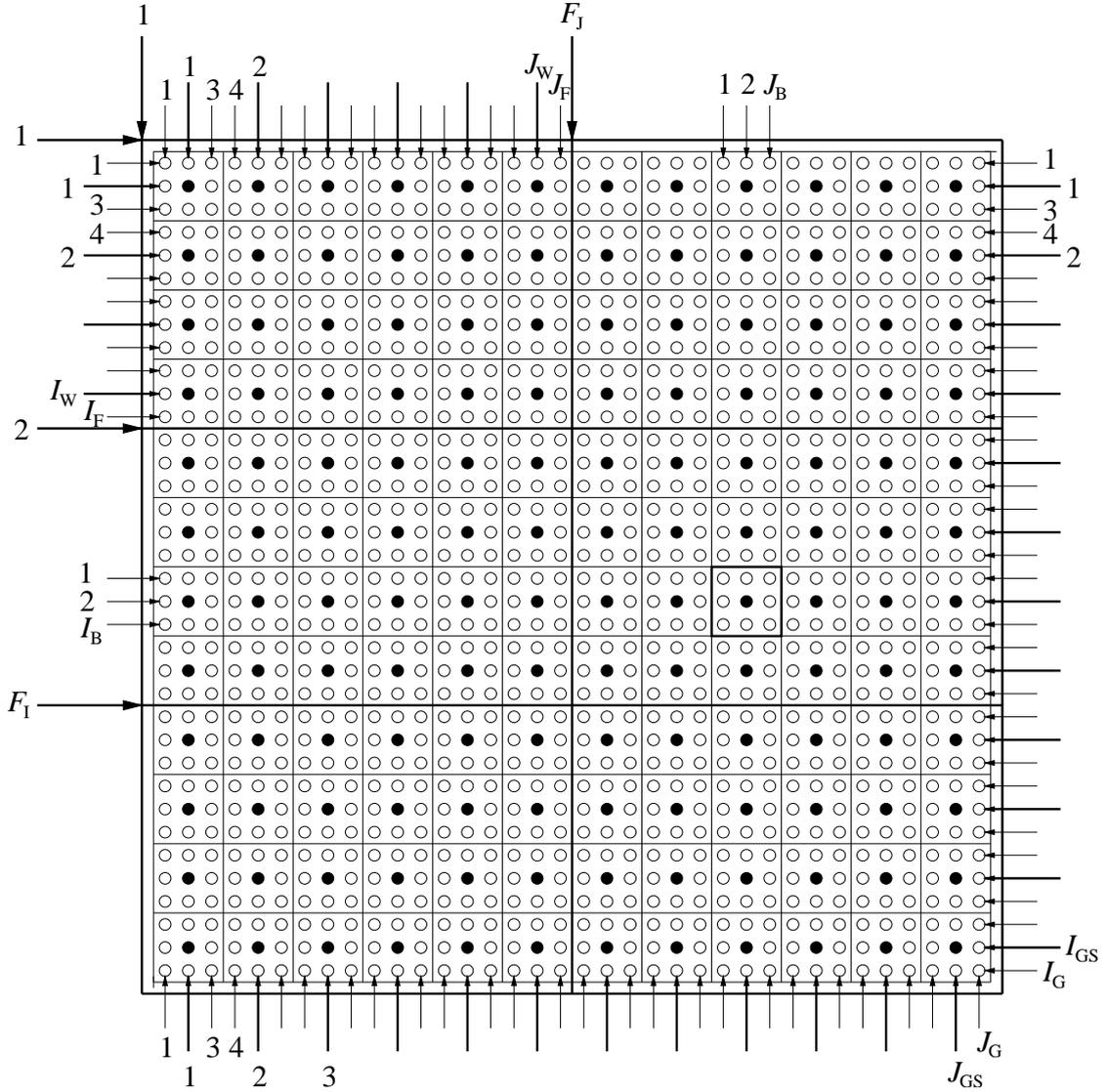


Figure 3: An example grid and its subunits: The grid consists of  $I_G \times J_G = 36 \times 36$  spots and  $F_1 \times F_2 = 3 \times 2$  fields. One field consists of  $I_F \times J_F = 12 \times 18$  spots and of  $I_W \times J_W = 4 \times 6$  blocks. One block consists of  $I_B \times J_B = 3 \times 3$  spots. Guide spots are centered in the blocks and marked black. The guide spots define a  $I_{GS} \times J_{GS} = 12 \times 12$  guide spot grid.

The dimensions of a field counted in blocks corresponds to the dimensions of the well plate. The following equations hold:

$$I_F = I_W * I_B \quad \text{and} \quad J_F = J_W * J_B. \quad (8)$$

When a large number of outliers is expected from the hybridization experiment so-called *guide spots* are used: A spot at a certain position within a block  $\mathcal{B}$  contains DNA material which has always a strong hybridization signal irrespective of the hybridization probe. We denote *regular spots* as spots not belonging to the class of guide spots.

**Definition 6 (Guide spot grid)**

A *guide spot grid*  $\mathcal{G}^*$  is a set of nodes in  $\{1 \dots I_{\text{GS}}\} \times \{1 \dots J_{\text{GS}}\}$ , with  $I_{\text{GS}}$  as the number of guide spot grid rows and  $J_{\text{GS}}$  as the number of guide spot grid columns, where

$$I_{\text{GS}} = F_I * I_w \quad \text{and} \quad J_{\text{GS}} = F_J * J_w \quad (9)$$

The row extraction and column extraction relations  $R_{\mathcal{G}^*}$  and  $C_{\mathcal{G}^*}$  are defined similarly to Eqns (2) and (3).

**Definition 7 (Guide spot field)**

A *guide spot field*  $\mathcal{F}^*_{pq} \subset \mathcal{G}^*$  is a subunit of a guide spot grid  $\mathcal{G}^*$  and is a set of nodes in  $\{1 \dots I_w\} \times \{1 \dots J_w\}$ . A guide spot grid  $\mathcal{G}^*$  is partitioned into  $F_I * F_J$  guide spot fields such that

$$\bigcup_{p=1}^{F_I} \bigcup_{q=1}^{F_J} \mathbf{L}(\mathcal{F}^*_{pq}, \mathbf{S}) = \mathbf{L}(\mathcal{G}^*, \mathbf{S}) \quad (10)$$

and

$$\mathbf{L}(\mathcal{F}^*_{pq}, \mathbf{S}) \cap \mathbf{L}(\mathcal{F}^*_{rs}, \mathbf{S}) = \emptyset, \quad \forall (p, q) \neq (r, s). \quad (11)$$

The row extraction and column extraction relations  $R_{\mathcal{F}^*}$  and  $C_{\mathcal{F}^*}$  are defined similarly to Eqns (2) and (3).

## 2.3 Prior Knowledge

The prior knowledge used for grid fitting consists of the information about a theoretical grid of a perfect experiment without any distortions. We therefore derive theoretical distances in pixels between grid nodes.

**Definition 8 (Theoretical spot distance)**

The *theoretical horizontal spot distance*  $S_x \in \mathbb{R}$  and the *theoretical vertical spot distance*  $S_y \in \mathbb{R}$  are the distances in pixels between two adjacent spots in the horizontal and vertical direction computed as

$$S_x = \frac{N_x}{\Delta x} \quad \text{and} \quad S_y = \frac{N_y}{\Delta y} \quad (12)$$

with  $N_x$  and  $N_y$  as the spot (needle) distances on the filter in millimeters and  $\Delta x$  and  $\Delta y$  as the scanner resolution in millimeters per pixel. Note that  $S_x$  and  $S_y$  are not rounded to integer values.

**Definition 9 (Theoretical block distance)**

The *theoretical horizontal block distance*  $B_x \in \mathbb{R}$  and the *theoretical vertical block distance*  $B_y \in \mathbb{R}$  are the distances in pixels between two adjacent blocks in the horizontal and vertical direction computed as

$$B_x = S_x * J_B \quad \text{and} \quad B_y = S_y * I_B \quad (13)$$

Table 2 provides an overview of the notation introduced in this section.

<b>Not.</b>	<b>Description</b>	<b>Not.</b>	<b>Description</b>
$I_B$	# block rows	$J_B$	# block columns
$I_W$	# well plate rows # guide spot field rows	$J_W$	# well plate columns # guide spot field columns
$I_F$	# field rows $I_F = I_W * I_B$	$J_F$	# field columns $J_F = J_W * J_B$
$F_1$	# fields in vertical direction	$F_1$	# fields in horizontal direction
$I_G$	# grid rows $I_G = F_1 * I_F$	$J_G$	# grid columns $J_G = F_1 * J_F$
$I_{GS}$	# guide spot grid rows $I_{GS} = F_1 * I_W$	$J_{GS}$	# guide spot grid columns $J_{GS} = F_1 * J_W$
$\Delta y$	vertical scanner resolution	$\Delta x$	horizontal scanner resolution
$N_y$	vertical spot distance [mm]	$N_x$	horizontal spot distance [mm]
$S_y$	vertical spot distance (theoretical) [pixel] $S_y = N_y / \Delta y$	$S_x$	horizontal spot distance (theoretical) [pixel] $S_x = N_x / \Delta x$
$B_y$	vertical block distance (theoretical) [pixel] $B_y = S_y * I_B$	$B_x$	horizontal block distance (theoretical) [pixel] $B_x = S_x * J_B$
$M$	vertical spot array image size	$N$	horizontal spot array image size

Table 2: Overview of the notation of the spot array image and the spot array. The left part of the table describes “vertical” entities, the right part of the table contains “horizontal” entities.

### 3 Guide Spot Detection

The output of the grid fitting procedure must provide approximate image locations for every spot. The guide spot grid  $\mathcal{G}^*$  in a spot array image is a reliable “safety” grid which spans the complete spot grid  $\mathcal{G}$ . It is therefore natural to focus the grid fitting efforts to the guide spot grid fitting, i.e. the fitting of the guide spot locations to the nodes of the guide spot grid. Consequently, the initial step of grid fitting consists of a detection of the guide spot locations. The other spot locations can be entirely inferred from the guide spot locations, because only local distortions are expected from the experiment. In the following we describe a guide spot detection approach which includes two filter operations and a local maximum search. In a first step, the spot array image is filtered with a matched filter describing the shape of the spots. Since the shape of many regular spots resembles the shape of the guide spots, the guide spot locations are amplified with the help of a nonlinear filter which considers the matched filter responses at the theoretical guide spot neighborhoods. The resulting response image is expected to contain maximum values at the guide spot locations. The guide spots can therefore be detected by a local maximum search.

#### 3.1 Spot Detection with Matched Filter

The first step towards a detection of guide spots is an amplification of their locations. We want to find signals in the spot array image which resemble the shape of the guide spots. One possible way is to use a matched filter (MF), which is a digital filter whose shape matches the shape of the signal one is trying to find [14]. The MF is optimal with respect to Gaussian noise. In order to find a random signal with non-zero mean in white noise, the filter should be matched to the mean of the signal. The MF for guide spots is constructed by forming an average template from a number of guide spots in the following way: The image patches  $\mathbf{G}_k$ ,  $k \in \{1 \dots G\}$  with the dimension  $M_{\text{MF}} \times N_{\text{MF}}$  contain the intensity values of the  $G$  guide spots which are manually selected by the user. The matched filter dimensions  $M_{\text{MF}}$  and  $N_{\text{MF}}$  should cover the guide spot extension for a given image resolution and spotting geometry. We formally define the matched filter dimensions as a function of the theoretical spot distances  $S_y$  and  $S_x$  defined in (12):

$$M_{\text{MF}} = \begin{cases} \circ(S_y) + 2 & \text{for } \circ(S_y) = 2k + 1 \\ \circ(S_y) + 1 & \text{otherwise} \end{cases} \quad (14)$$

and

$$N_{\text{MF}} = \begin{cases} \circ(S_x) + 2 & \text{for } \circ(S_x) = 2k + 1 \\ \circ(S_x) + 1 & \text{otherwise} \end{cases} \quad (15)$$

with  $k \in \mathbb{N}$  and  $\circ(\cdot)$  as the rounding operator. Irrespective of its parity, the theoretical spot size rounded to the next integer is increased to the next higher odd number. This is done since the guide spots may have very strong hybridization signals and might exceed the theoretical spot size.

The  $G$  images patches  $\mathbf{G}_k$  containing guide spots are formally rearranged as  $D$ -dimensional vectors  $\mathbf{g}_k$  by lexicographical ordering, with  $D = M_{\text{MF}} \cdot N_{\text{MF}}$ . The vectors  $\mathbf{g}_k$  are normalized as

$$\tilde{\mathbf{g}}_k = \mathbf{g}_k - \mu_{\mathbf{g}_k} \cdot \mathbf{1}, \quad (16)$$

where  $\mathbf{1}$  is a  $D \times 1$  vector of ones and  $\mu_{\mathbf{g}_k}$  is the mean intensity value of the the image patch defined as

$$\mu_{\mathbf{g}_k} = \frac{1}{D} \sum_{i=1}^D \mathbf{g}_k[i]. \quad (17)$$

The matched filter  $\mathbf{m}$  is constructed by averaging the  $G$  normalized examples  $\tilde{\mathbf{g}}_k$ :

$$\mathbf{m} = \frac{1}{G} \sum_{k=1}^G \tilde{\mathbf{g}}_k \quad (18)$$

Filtering of the spot array image  $\mathbf{S}$  with the matched filter  $\mathbf{m}$  results in a response image  $\mathbf{R}^M$  which is constructed as follows: If  $\mathbf{s}_{[m,n]}$  denotes an  $M_{\text{MF}} \times N_{\text{MF}}$  image patch around a pixel  $\mathbf{S}[m,n]$  rearranged as a  $D$ -dimensional vector, the image patch  $\mathbf{s}_{[m,n]}$  is first normalized to the local intensity mean value  $\mu_{\mathbf{s}_{[m,n]}}$  similarly to Equation (16) as

$$\tilde{\mathbf{s}}_{[m,n]} = \mathbf{s}_{[m,n]} - \mu_{\mathbf{s}_{[m,n]}} \cdot \mathbf{1}. \quad (19)$$

The matched filter response value is then the dot product

$$\mathbf{R}^M[m,n] = \tilde{\mathbf{s}}_{[m,n]} \cdot \mathbf{m} \quad (20)$$

corresponding to the similarity or statistical covariance between the image patch and the matched filter. High response values in  $\mathbf{R}^M$  indicate potential guide spot locations. Figure 4 shows a filtering example. Figure 4b is the matched filter response image  $\mathbf{R}^M$  of the image in Figure 4a. Note that some spots in this example are very similar to the guide spots with regard to shape and intensity. Consequently, regular spots can have higher response values than guide spots. We want to detect guide spot locations based on maximum response values. Hence an additional filtering step is necessary in order to get rid of the high response values in  $\mathbf{R}^M$  which do not belong to guide spots.

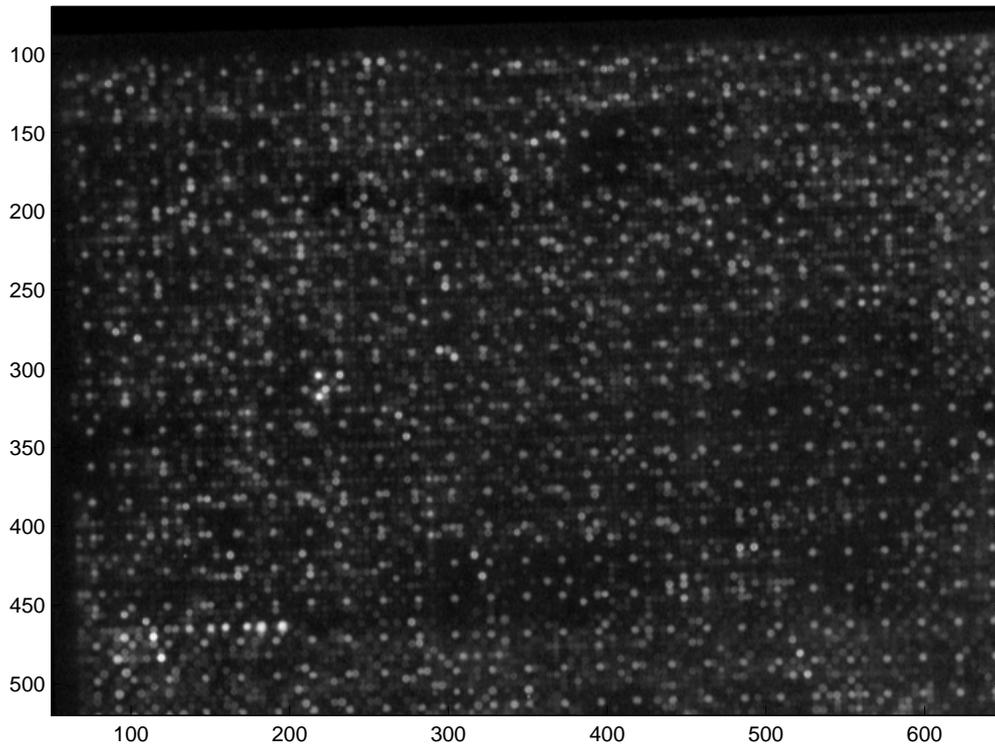
### 3.2 Guide Spot Location Amplification (GSLA)

We saw in Fig. 4 that the shape of the guide spots can be very similar to the shape of the regular spots. It is therefore not guaranteed that all high values in the matched filter response image  $\mathbf{R}^M$  indicate guide spot locations. The main idea to overcome this problem is to amplify the locations of potential guide spots by considering the MF response values at the theoretical guide spot neighborhood locations. Since a guide spot location is part of a grid, its grid neighborhood locations must also have high matched filter response values. If this is not the case, it is likely that the location is not a guide spot.

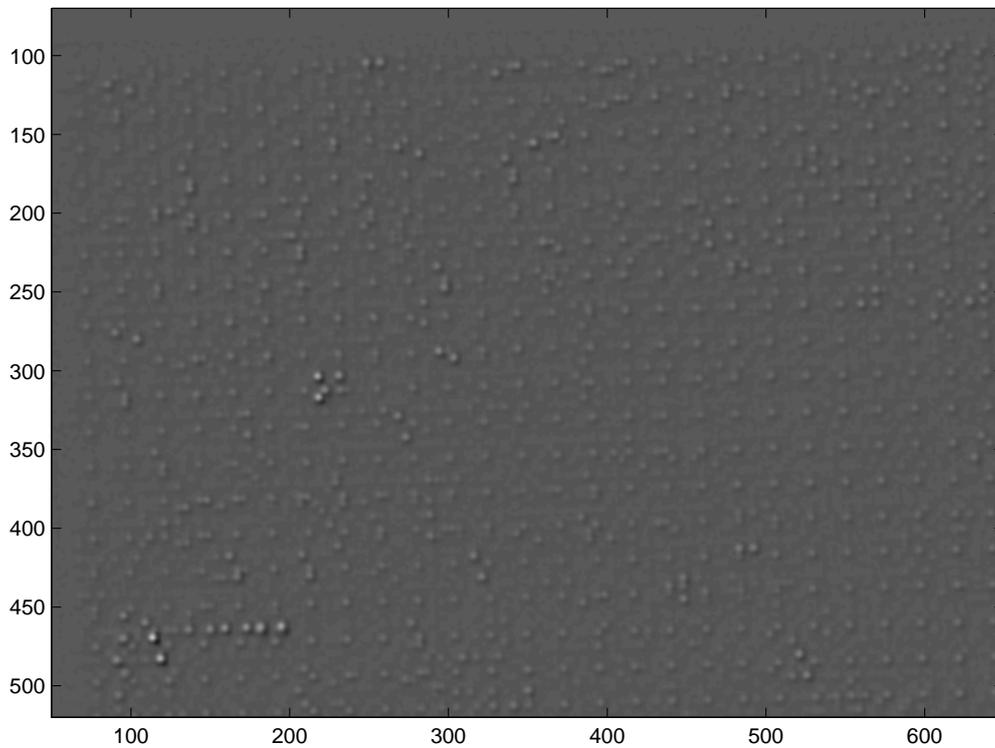
We formally define the set  $\mathcal{T}_{[m,n]}$  which includes the response value  $\mathbf{R}^M[m,n]$  and the response values of the theoretical guide spot neighborhood locations of  $(m,n)$  in  $\mathbf{R}^M$  as

$$\mathcal{T}_{[m,n]} = \{\mathbf{R}^M[m+k, n+l] \mid k \in \{0, \circ(B_y), \circ(-B_y)\} \wedge l \in \{0, \circ(B_x), \circ(-B_x)\}\}. \quad (21)$$

with  $B_y$  and  $B_x$  as the theoretical block distances defined in (13). Figure 5 illustrates the neighborhood set  $\mathcal{T}_{[m,n]}$  for the center black pixel assuming theoretical block distances  $B_y = B_x = 15$ . The GSLA response value  $\mathbf{R}^A[m,n]$  is determined as



(a) Part of a  $1300 \times 1286$  spot array image



(b) Matched filter response image of (a)

Figure 4: Digital filtering with a matched filter. A high pixel value at the matched filter response image (b) indicates a high similarity to the matched filter of the guide spots and is proportional to the probability of a guide spot location in (a). Responses of regular spots can be stronger than the guide spots responses.

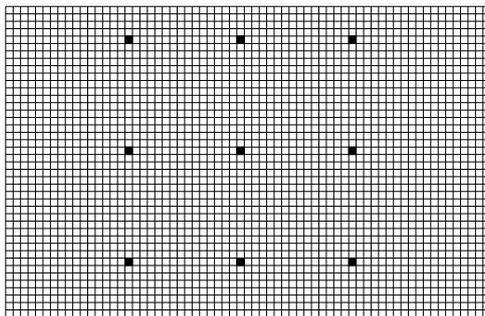


Figure 5: Example for guide spot location amplification (GSLA): The black pixels are considered in the computation of the GSLA response value for the center black pixel (assuming theoretical block distances of 15 pixels). The median of the intensities at the 9 locations is taken as the response value of the center pixel. If the center black pixel is a guide spot location, the grid neighborhood locations will have high response values and the median of the response values will be high.

$$\mathbf{R}^A[m, n] = \text{median}(\mathcal{T}_{[m, n]}). \quad (22)$$

If a MF response value at a location  $(m, n)$  and all the response values in the theoretical neighborhood are high it is likely that  $(m, n)$  is a guide spot location. This is not the case for locations where the neighborhood response values are low. Note that the median value is more robust for the guide spot location amplification than the mean value. In case of a regular spot (or an artifact) with a very high MF impulse response compared to the theoretical neighbors, the mean value measure would propagate high values to the theoretical grid neighbors and generate new local grid structures. Figure 6 shows the GSLA response  $\mathbf{R}^A$  of the MF response image in Fig. 4b. The locations of the guide spots are strongly amplified. Nevertheless, it is still possible that false maximum response values survive, especially if a set of bright regular spots establish a local grid. A drawback of the median is that the grid *corner nodes* will be suppressed since only at most three – and therefore less than 50 % – of the theoretical neighborhood positions will have a high response value.

### 3.3 Local Maximum Search

After having applied the GSLA filter to the matched filter response image it is very likely that maximum values in the GSLA response image  $\mathbf{R}^A$  indicate guide spot locations. The strategy of the maximum search is to reliably extract a set of guide spot locations  $\mathcal{L}^* \subset \mathcal{L}$  from  $\mathbf{R}^A$ . By reliable we mean that an image coordinate  $(x, y)$  in  $\mathcal{L}^*$  should not belong to a regular spot location even at the risk of not detecting all guide spot locations (the lacking guide spot locations are determined in a subsequent step). Put into other words, we want to minimize the detection of false positives for guide spot locations. Since we assume one guide spot per block  $\mathcal{B}$ , local maximum values are searched in non-overlapping windows having approximately the dimensions of the theoretical block size. Since it is by no means guaranteed that the maxima are centered in the search windows, it is necessary to perform the local maximum search in at least two passes. As mentioned above, even after the application of the GSLA filter it is

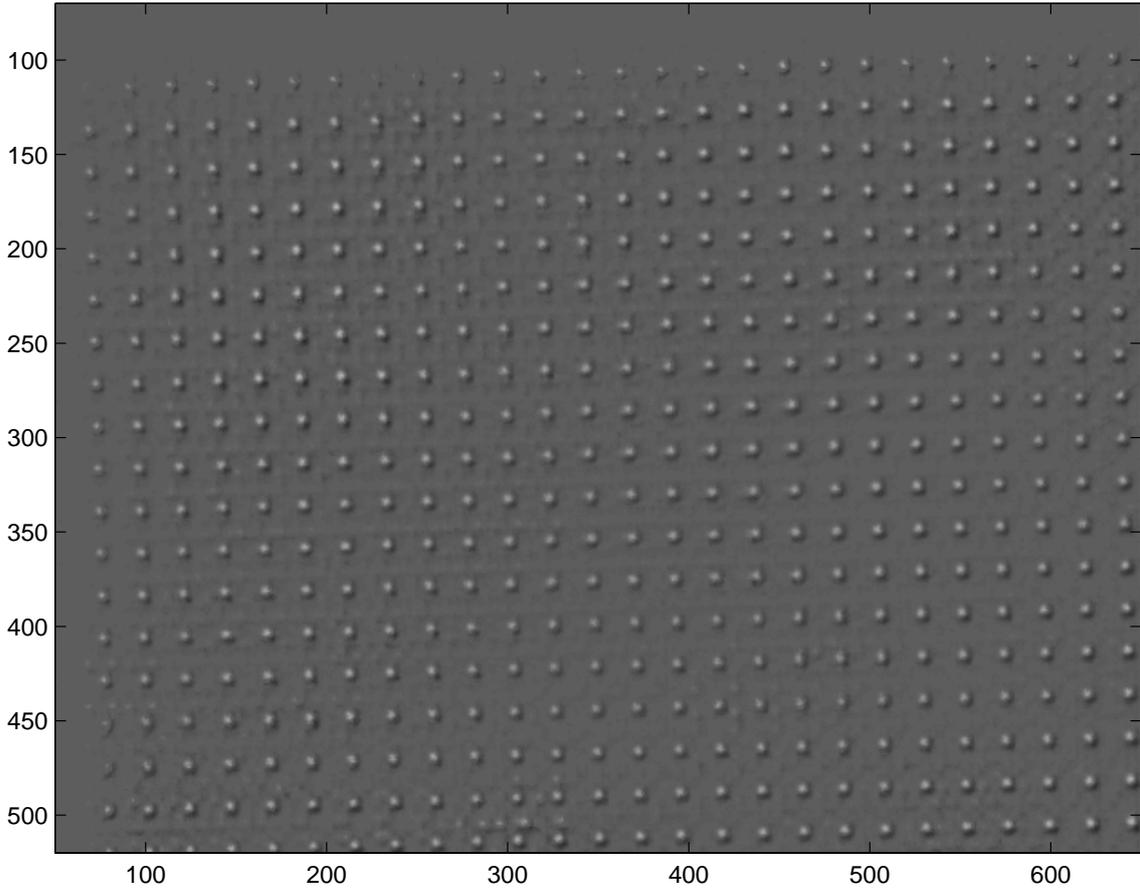


Figure 6: Guide spot location amplification (GSLA) filtered response image of the matched filter response image in Fig. 4b. High intensity values indicate that locations of the theoretical guide spot neighborhood also have high intensities. The guide spot location at the upper left grid corner is missing because the median of 8 neighborhood positions is computed. A corner location only has 3 guide spot neighbors.

still possible that maximum response values not belonging to guide spots could have survived. These maxima can be removed in a third pass by considering again the theoretical guide spot grid neighborhood. Before these three passes of maximum search are performed, it is necessary to specify a region of interest (ROI) for the maximum search. Otherwise, the maximum search would deliver local maxima in image regions containing no hybridization information.

### 3.3.1 Region-of-Interest and Prior Guide Spot Locations

When performing a maximum search for guide spot locations in local search windows, it is necessary to constrain the search area to a region-of-interest (ROI). Without an ROI, even very low local maximum values of the spot array image border (the image region not covering the physical spot filter) would be marked as potential guide spot locations and unnecessarily complicate the grid fitting. The ROI is closely related to the concept of the *prior guide spot locations*:

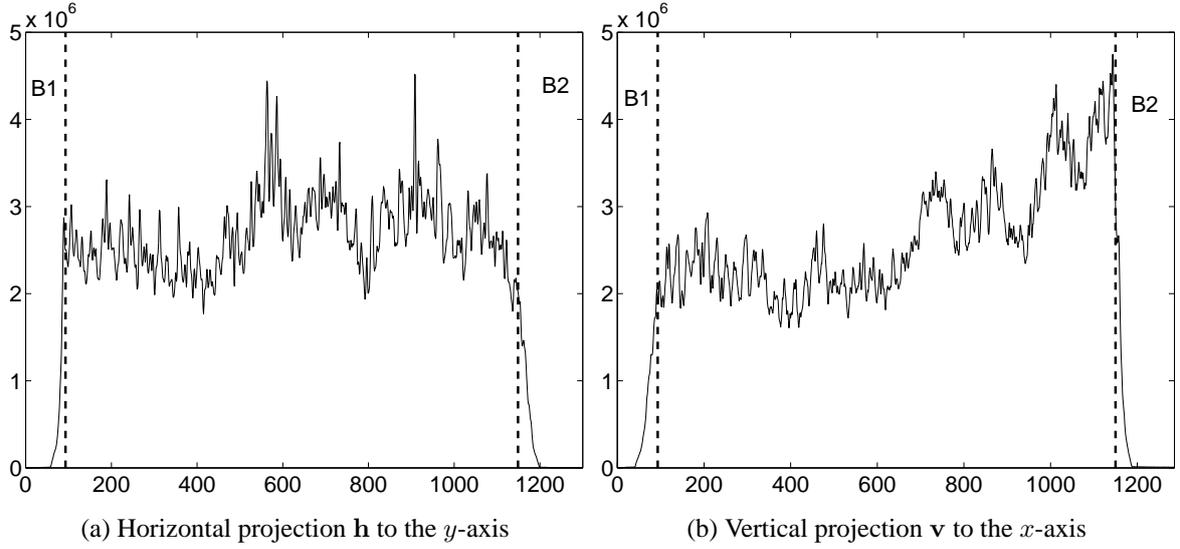


Figure 7: Projections of the spot array image intensities in Fig. 1. The regions  $B1$  and  $B2$  are considered to belong to the image border (the part of the image not containing the physical filter) and are determined with the help of the prior knowledge about the image size and the spotting geometry. The projections form the basis of the region-of-interest for the maximum search.

The prior guide spot locations  $L_p(\mathcal{G}^*, \mathbf{S})$  are initial estimates for the locations of the guide spots based on the prior knowledge of the theoretical block distances  $B_x$  and  $B_y$ . The prior guide spot locations can be obtained as follows: The image border of the spot array image  $\mathbf{S}$  has significantly lower intensity values than the filter region. It is therefore possible to estimate the locations of the prior corner guide spot locations with the help of *projections* of the intensity values of  $\mathbf{S}$  to the axis of the spatial image coordinate system. A projection of a two-dimensional function  $f(x, y)$  is a line integral in a certain direction. The line integral of  $f(x, y)$  in the vertical direction is the projection of  $f(x, y)$  onto the  $x$ -axis; the line integral in the horizontal direction is the projection of  $f(x, y)$  onto the  $y$ -axis. Expressed in terms of an  $M \times N$  image matrix  $\mathbf{S}$  we have

- an  $M \times 1$  horizontal projection vector  $\mathbf{h}$  in which the  $i$ -th element corresponds to the sum of the components of the  $i$ -th row vector of  $\mathbf{S}$ .
- an  $N \times 1$  vertical projection vector  $\mathbf{v}$  in which the  $i$ -th element corresponds to the sum of the components of the  $i$ -th column vector of  $\mathbf{S}$ .

Figure 7 shows the horizontal and vertical projection of the spot array image intensities in Fig. 1. The  $x$ -axis of Fig. 7a corresponds to the components of the horizontal projection vector  $\mathbf{h}$  (with  $x$  as the row number or  $y$ -coordinate of the spot array image). The  $y$ -axis of Fig. 7a corresponds to the projection values of the intensity values normalized by the image width. Figure 7b shows the vertical projection. Projections of the image border have significantly lower projection values than projections of the spot filter area. It is also expected that the components of  $\mathbf{h}$  and  $\mathbf{v}$  belonging to the background form two connected regions  $B1$  and  $B2$  (see Fig. 7). The expected

number of components  $M_b$  belonging to the horizontal projections of the image background is easily determined with the help of the theoretical vertical spot distance  $S_x$ , the grid height  $I_G$  and the spot array image height  $M$ :

$$M_b = M - S_y \cdot I_G. \quad (23)$$

The number  $N_b$  of vertical projections belonging to the background is determined in a similar manner. The location of the upper left prior guide spot location  $L_p((1, 1))$  is determined with the following algorithm:

1. Sort the components of the horizontal projection  $\mathbf{h}$  by their projection values and label the  $M_b$  lowest projection values. Do the same for the  $N_b$  lowest projection values of the vertical projection  $\mathbf{v}$ . We expect that the labeled components with the lowest projection values correspond to the background regions  $B1$  and  $B2$  in Fig. 7.
2. The first (“leftmost”) unlabeled components of  $\mathbf{h}$  and  $\mathbf{v}$  plus an offset of half the theoretical block distances  $\circ(B_y/2)$  and  $\circ(B_x/2)$ , respectively, are an estimate of the upper left prior guide spot location  $L_p((1, 1))$ . This estimate corresponds to the left dashed line in Fig. 7a and b.
3. The locations of the other grid nodes  $L_p((i, j))$  have horizontal pixel distances  $B_x$  and vertical pixel distances  $B_y$ . The right dashed lines in Fig. 7 mark the rightmost prior guide spot location.

Figure 8 shows the prior guide spot locations for the spot array image in Fig. 1 superimposed to the GSLA response image  $\mathbf{R}^A$  of Fig. 6. Note that the spot grid  $\mathcal{G}$  in the spot array image is significantly rotated. The prior guide spot grid is not rotated, so that the deviations from the true guide spot locations are balanced (keep in mind that Fig. 8 shows only a  $450 \times 600$  part of the  $1300 \times 1286$  image). In order to determine the rectangular ROI for the maximum search we must add a tolerance area to the area covered by the prior guide spot locations. We define the corner points of the rectangular ROI as the prior grid corner locations translated by a tolerance offset vector  $\mathbf{t} = (\pm B_x \pm B_y)^T$  (the sign depends on the corner of the grid).

### 3.3.2 Maximum Search Algorithm

We expect that the guide spot locations correspond to maximum values in the ROI of the GSLA response image. Since we assume one guide spot per block  $\mathcal{B}$ , the maximum values are searched in non-overlapping windows of approximately the theoretical block size. The maxima are searched in three passes.

**Pass 1** The initial guide spot location set  $\mathcal{L}^*$  consists of the locations of the maximum value in every *non-overlapping*  $M_L \times N_L$  window. The window dimensions  $M_L$  and  $N_L$  are the next smaller odd number of the theoretical block dimension  $B_y$  and  $B_x$ :

$$M_L = \begin{cases} \circ(B_y) - 2 & \text{if } \circ(B_y) = 2k + 1 \\ \circ(B_y) - 1 & \text{otherwise} \end{cases} \quad (24)$$

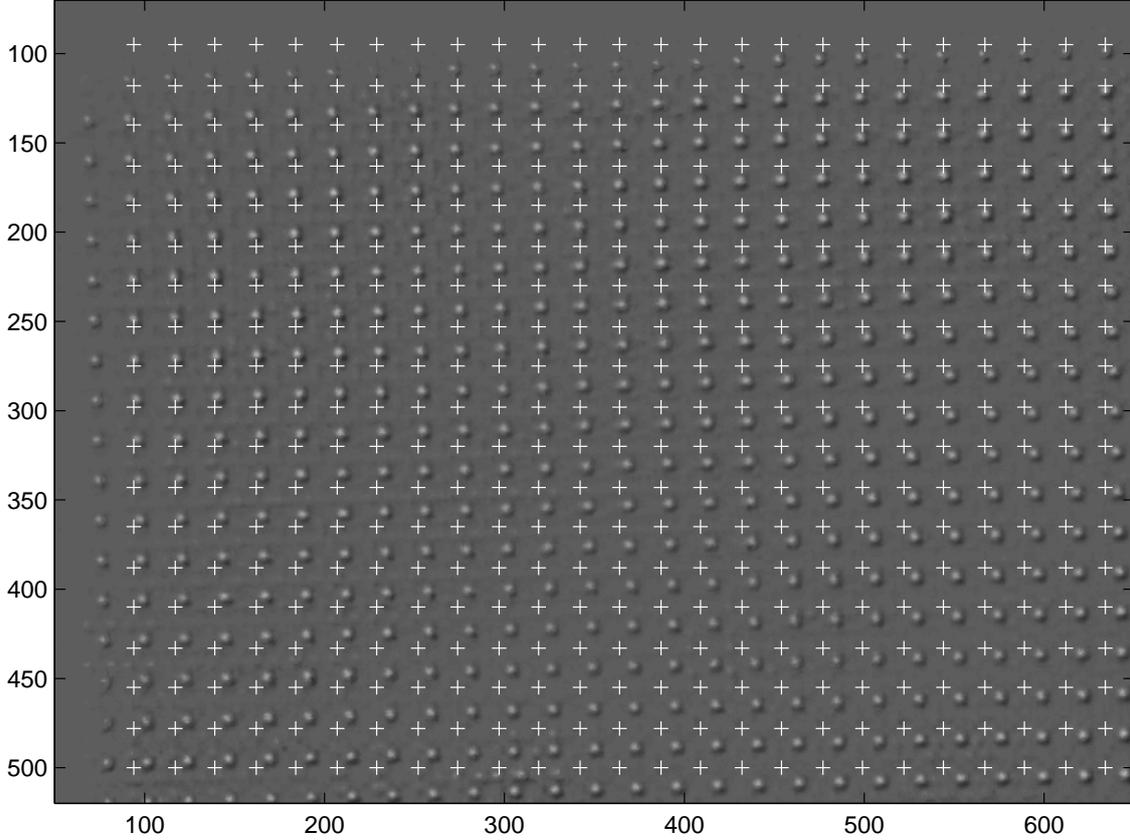


Figure 8: Prior locations for the guide spots in Fig. 1 based on the projections of Fig. 7. The white crosses mark the locations of the prior guide spot grid nodes. Together with a tolerance of approximately one block size these locations span the ROI for the maximum search. The prior guide spot locations are transformed in a later step in order to assign the detected guide spot locations to guide spot grid nodes.

and

$$N_L = \begin{cases} \circ(B_x) - 2 & \text{if } \circ(B_x) = 2k + 1 \\ \circ(B_x) - 1 & \text{otherwise} \end{cases} \quad (25)$$

for  $k \in \mathbb{N}$ . This window size avoids that two guide spot locations fall into one search window and therefore avoids the cancellation of potential guide spot locations. Figure 9a illustrates this step. The black pixels indicate the maxima in the  $M_L \times N_L$  blocks. Note the two close maxima in the second block row.

**Pass 2** For every detected maximum location  $(x, y) \in \mathcal{L}^*$ , select the location  $(x', y')$  with the maximum response value in an  $M_L \times N_L$  window around  $(x, y)$ . If  $(x', y') \neq (x, y)$ , remove  $(x, y)$  from  $\mathcal{L}^*$  and add  $(x', y')$  to  $\mathcal{L}^*$ . After this pass it is guaranteed that two guide spot locations have plausible distances. In Figure 9b the grey-shaded area illustrates the plausible distance criterion. If we assume that in the second row the first local maximum value is higher than the second maximum value, the second maximum will be canceled.

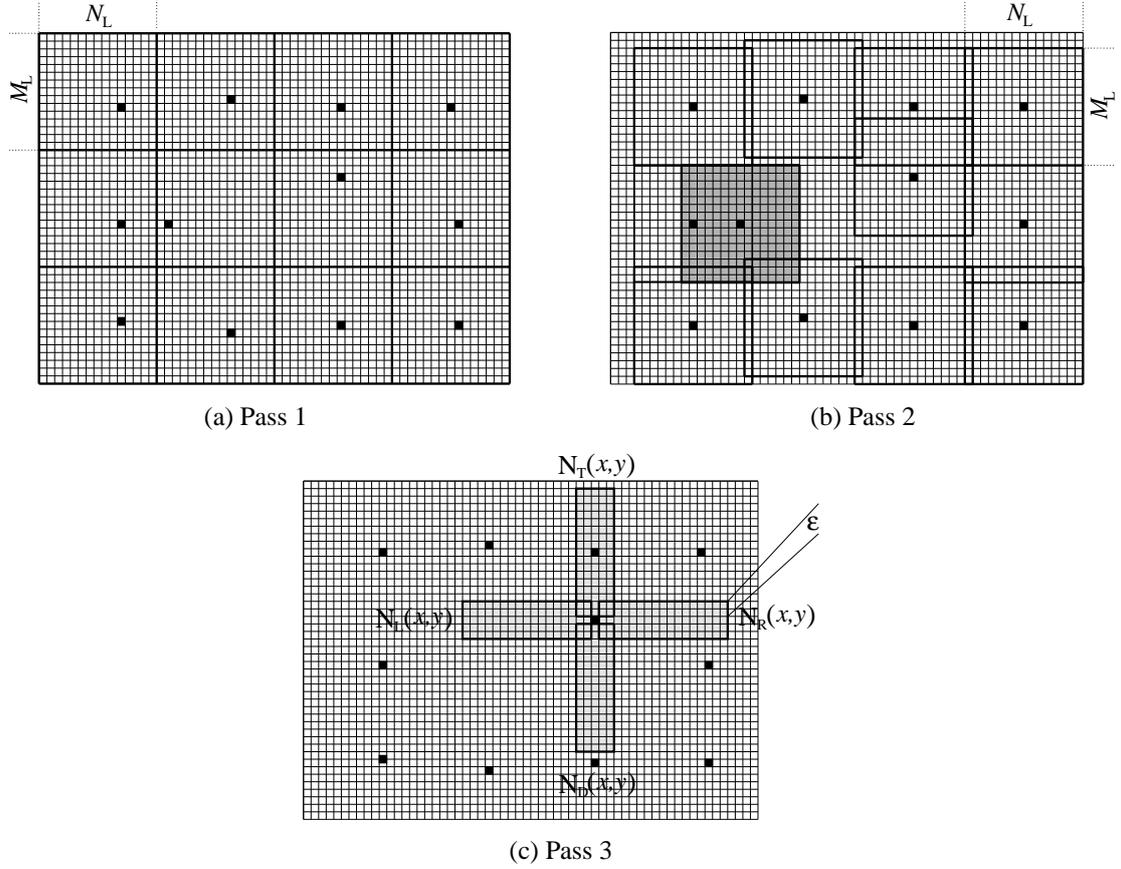


Figure 9: Principle of 3-pass maximum search. Maxima are first searched in non-overlapping  $M_L \times N_L$  windows which are smaller than the block size. Black pixels in (a) indicate local maxima. The Maximum search is repeated in pass 2 in order to get plausible guide spot distances. In (b) the lower maximum from the first pass will be canceled in the shaded area. The third pass tries to remove locations not belonging to the guide spot grid. A location in is removed if it does not have a left and a right or a top and a bottom neighbor.

**Pass 3** For every maximum location  $(x, y)$ , define a left neighborhood location set  $N_L(x, y)$ , a right neighborhood location set  $N_R(x, y)$ , an upper neighborhood location set  $N_U(x, y)$  and a lower neighborhood location set  $N_D(x, y)$  within a tolerance  $\varepsilon$ . Fig. 9c sketches the neighborhood location sets for a tolerance  $\varepsilon = 2$ . For real spot array images, it is useful to chose  $\varepsilon$  as the theoretical spot distance  $S_x$  or  $S_y$ . The neighborhood sets should help eliminate false positives that do not lie within the grid of the guide spots. Hence every location  $(x, y)$  in  $\mathcal{L}^*$  that has *no left and no right neighbor* or *no upper and no lower neighbor* is removed from  $\mathcal{L}^*$ . Formally, we define the conditions

$$N_L(x, y) \cap \mathcal{L}^* = \emptyset \quad \wedge \quad N_R(x, y) \cap \mathcal{L}^* = \emptyset \quad (26)$$

and

$$N_U(x, y) \cap \mathcal{L}^* = \emptyset \quad \wedge \quad N_D(x, y) \cap \mathcal{L}^* = \emptyset. \quad (27)$$

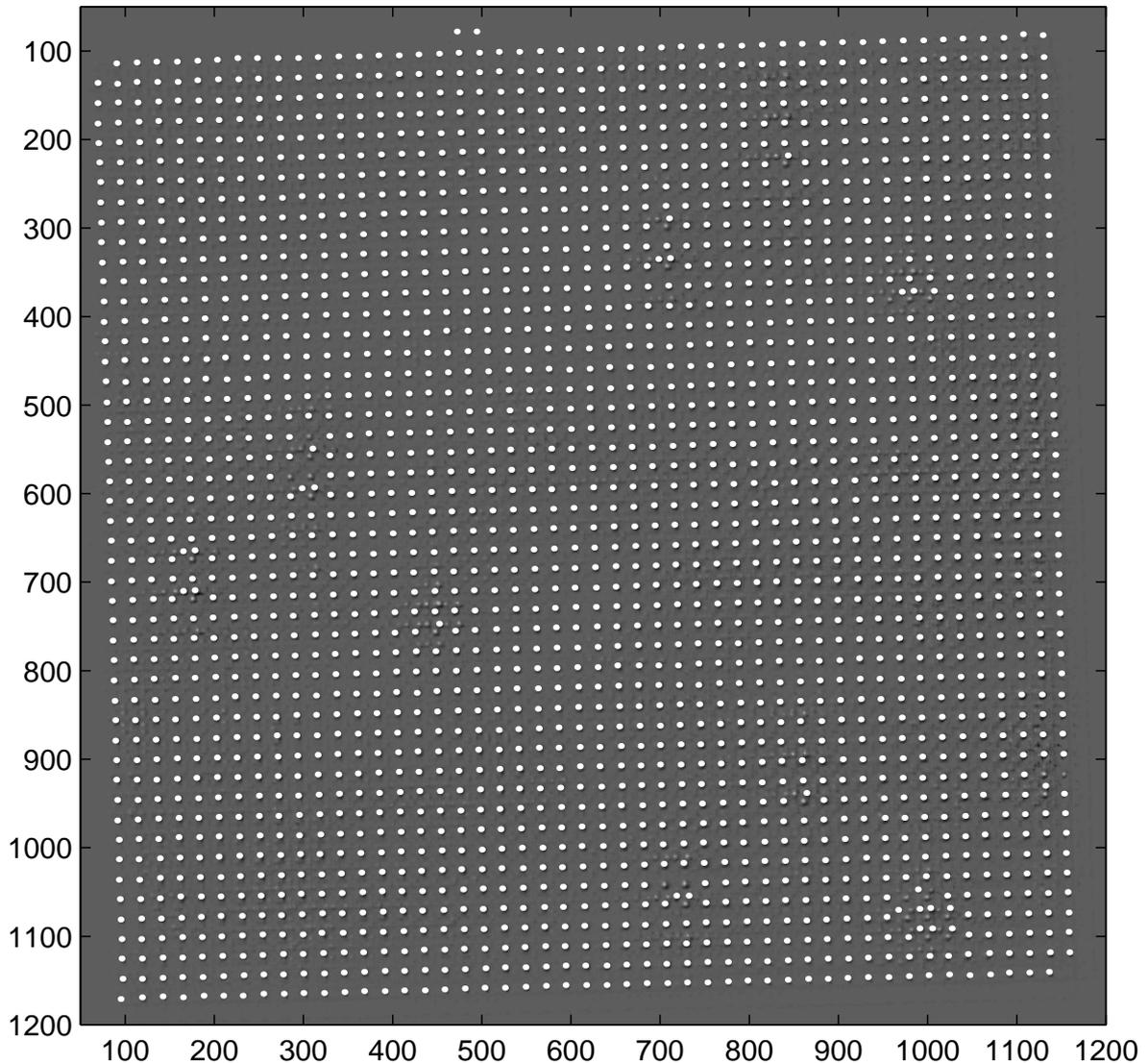


Figure 10: Detected guide spot locations for the spot array image of Fig. 1 superimposed as white dots to the GSLA response image in Fig. 6. There are still false positives, for example at the top of the spot array. The task of the grid fitting is to detect these false positives, to assign the true positives to grid nodes and to determine the locations of the lacking guide spots.

If condition (26) or condition (27) holds for a detected guide spot location  $(x, y) \in \mathcal{L}^*$ , remove  $(x, y)$  from  $\mathcal{L}^*$ .

Figure 10 shows the guide spot detection results for the spot array image in Fig. 1 superimposed as white dots to the GSLA response image of Fig. 6. One can see the following problems:

- Two corner locations (the upper left and the lower right) are not detected. This happens because the GSLA filter suppresses grid corner locations.
- Two detected locations (on the top of the grid) are outside the guide spot grid. The

ROI must include this area because of the rotation of the grid. If there are at least two neighbored locations outside the grid, the third pass of the maximum search cannot erase these locations because they are embedded into the grid structure.

- Some locations within the guide spot grid are false positives. This occurs when many strong regular spots in contiguous blocks span a local grid.

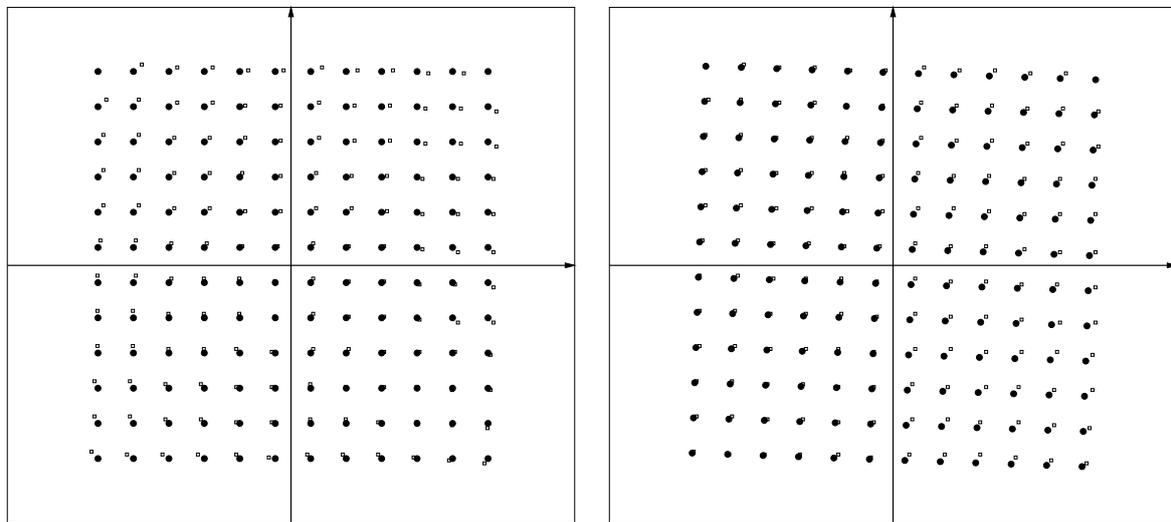
### 3.4 Summary

The goal of the guide spot detection step is to extract a set of locations which are likely to belong to guide spot centers in the spot array image. The guide spot locations are first searched with the help of a matched filter whose signal resembles the shape of the guide spots. The matched filter is built manually from a series of template guide spots and is normalized to its mean value. Since the regular spots and guide spots may have similar shape and intensity, the guide spot locations are amplified with a second nonlinear filter. This guide spot location amplification (GSLA) filter amplifies locations that lie within the guide spot structure and suppresses isolated regular spots. The set of potential guide spot locations is determined by a search of local maxima in non-overlapping windows of approximately the block size. Since local maxima outside the area covered by the physical filter do not make sense, the search is constrained to a region-of-interest (ROI). The ROI is spanned by the prior guide spot locations which are determined with the help of the horizontal and vertical projections of spot array image intensity values to the  $x$ - and  $y$ -axis of the image coordinate system. The maximum search also tries to eliminate locations that do not lie within the grid structure. The guide spot detection does not guarantee that all the detected guide spot locations are true positives. It is the task of the grid fitting to remove the false positives, to assign the true positives to the nodes of the guide spot grid and to determine the guide spot locations that were not detected.

## 4 Grid Fitting

The grid fitting procedure tries to span a grid from the set of detected guide spot locations. These potential guide spot locations which were found with the digital filters and the maximum search as described in Sect. 3 are only an unordered set with no information about the position on the guide spot grid. The first main task of the grid fitting is therefore to map the detected locations to the correct nodes of the guide spot grid. We also call this mapping the *alignment* of the detected guide spot locations. The idea to align the detected locations is to transform the prior guide spot locations of Sect. 3.3.1 in a way that they can serve as reference locations for the detected locations. As illustrated in Fig. 11, a reference location only has to be in the near of the detected guide spot location. We assume that the nonlinear distortions of the guide spot grid during the imaging process of a spot array image are small enough such that a *linear transformation* of the prior guide spot locations is sufficient. The nonlinearities of the guide spot grid are captured by the local search for the detected guide spot locations. In order to apply a linear transformation, we must extract information about two features of the grid: First, we must estimate the global grid rotation. Second, we must know how to translate (or shift) the prior guide spot locations.

The second main task of grid fitting is the construction of a consistent grid. A consistent grid is a grid in which the *false negatives*, i.e. the guide spot locations that were not detected, are fitted. Furthermore, a consistent grid does not contain *false positives*, i.e. locations, which



(a) Prior guide spot locations (circles)

(b) Transformed guide spot locations (circles)

Figure 11: Alignment of detected guide spot locations: The black circles in (a) indicate the prior guide spot locations as defined in Sect. 3.3.1, the squares indicate the detected guide spot locations. The prior guide spot locations are transformed such that they define a reference grid for the detected guide spot locations. A local search in the neighborhood of the reference locations then tries to map the detected guide spot locations to the correct grid nodes.

do not belong to guide spots. The main idea is to robustly fit straight lines for every column and every row of every field. A lacking guide spot location can then be determined by the intersection of the corresponding straight lines. Once a consistent guide spot grid is available, the locations of the other (regular) spots are initialized. We provide a detailed description of the solutions to all the grid fitting tasks in this section.

## 4.1 Alignment of Detected Guide Spot Locations

The first main task of the grid fitting is to map the detected guide spot locations  $\mathcal{L}^*$  to the correct nodes of the guide spot grid, i.e. we want to find the locations  $L_D(\mathcal{G}^*)$  of the correct guide spot grid. This is accomplished by computing a set of reference grid locations  $L_R(\mathcal{G}^*)$  being *in the near* of the detected guide spot locations  $L_D(\mathcal{G}^*)$ . A search for a detected location  $L_D((i, j))$  in a neighborhood of  $L_R((i, j))$  should then provide us with the correct location mapping. A reference location  $[x_R \ y_R]^T$  in the central coordinate system is computed from a prior guide spot location  $[x_P \ y_P]^T$  as follows:

$$\begin{bmatrix} x_R \\ y_R \end{bmatrix} = \begin{bmatrix} \cos \theta_G & \sin \theta_G \\ -\sin \theta_G & \cos \theta_G \end{bmatrix} \begin{bmatrix} x_P \\ y_P \end{bmatrix} + \begin{bmatrix} t_x \\ t_y \end{bmatrix}, \quad (28)$$

where  $\theta_G$  is an estimate of the global grid rotation angle and  $\mathbf{t} = [t_x \ t_y]^T$  is an estimate for the translation vector of the rotated prior guide spot locations. In the next two sections we describe how the parameters  $\theta_G$  and  $\mathbf{t}$  can be estimated.

### 4.1.1 Global Rotation Estimation

A rotation angle  $\theta_G$  which estimates the rotation of the spot array can be computed by generalizing the concept of the horizontal and vertical projections (Fig. 7) to projections in directions of an arbitrary angle  $\theta$ . The theoretical framework of such projections is given by the the Radon transform [9]. In general, the Radon transform of a function  $f(x, y)$  is the line integral of  $f$  parallel to the  $y'$ -axis:

$$R_\theta(x') = \int_{-\infty}^{\infty} f(x' \cos \theta - y' \sin \theta, x' \sin \theta + y' \cos \theta) dy' \quad (29)$$

where

$$\begin{bmatrix} x' \\ y' \end{bmatrix} = \begin{bmatrix} \cos \theta & \sin \theta \\ -\sin \theta & \cos \theta \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix}. \quad (30)$$

Figure 12 illustrates the geometry of the Radon transform. Suppose the projections are applied to a GSLA image (Fig. 6) in which only the guide spot locations are expected to have high response values. When projecting these response values of the GSLA image along different directions  $\theta$ , only the projection for the correct grid rotation will go through all the guide spot locations. The rotation angle  $\theta_G$  is therefore determined by the projection with the *maximum projection values* for the guide spot locations.

The *discrete* Radon transform  $\mathbf{R}^T$  is defined as an  $R \times C$  matrix. The number of rows  $R$  corresponds to the number of directions for which the projection is computed.  $R$  depends on the angle resolution  $\Delta\theta$  of the projection angles and the maximum rotation angle  $\theta_M$ :

$$R = \frac{2\theta_M + 1}{\Delta\theta} \quad (31)$$

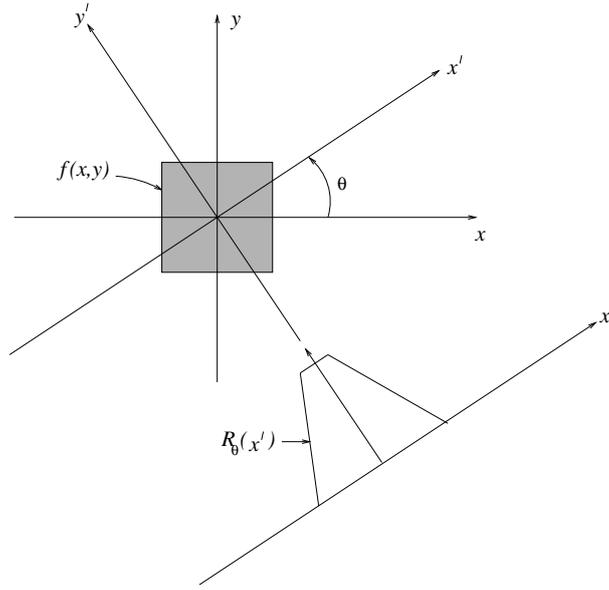


Figure 12: Geometry of the Radon transform. A function  $f(x, y)$  is projected along the  $y'$ -axis of a coordinate system rotated by the angle  $\theta$ .

Note that also negative rotation angles must be considered, hence the factor 2 in (31). The rotation angle  $\theta(r)$  belonging to a row index  $r \in \{1 \dots R\}$  is given by

$$\theta(r) = r\Delta\theta - \theta_M - 1. \quad (32)$$

The number of columns  $C$  is defined as the size of the spot array image diagonal:

$$C = \lceil \sqrt{M^2 + N^2} \rceil. \quad (33)$$

Figure 13 shows a part of a discrete Radon Transform for a GSLA image with an angle resolution  $\Delta\theta = 0.2^\circ$  and a maximum rotation angle  $\theta_M = 10^\circ$ . The global grid rotation estimate  $\theta_G = \theta(r)$  is in the row  $r$  of the Radon transform  $\mathbf{R}^T$  having the highest projection values. Let

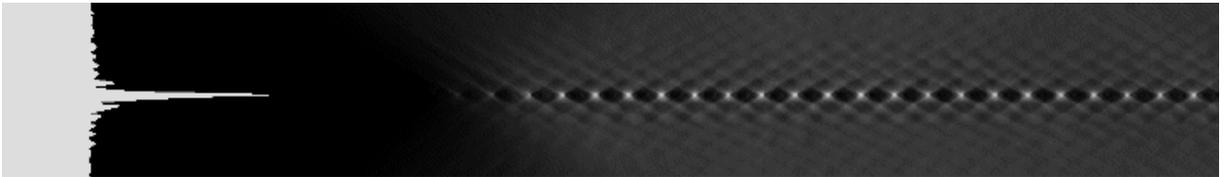


Figure 13: Part of a discrete Radon transform. Every row is a projection of the GSLA response values (Fig. 6) along a specific direction. The estimate for the grid rotation is given by the row with the highest projection values for the  $J_{GS}$  guide spot locations. Since the projection values between the guide spot locations are low, the row with the highest median of the  $J_{GS}$  highest projection values is determined as the grid rotation. The median values are illustrated as horizontal bars at the left hand side of the image.

$\mathcal{M}_r$  be the set of the  $J_{GS}$  (number of guide spots in a row) highest projection values in the row  $r$  of  $\mathbf{R}^T$ . The estimated rotation angle is then found in the row with the highest median of  $\mathcal{M}_r$ :

$$r = \arg \max\{\text{median}(\mathcal{M}_r)\}. \quad (34)$$

In Fig. 13, the values  $\text{median}(\mathcal{M}_r)$  for every row are visualized as horizontal bars. Note that the row with the true rotation angle has the highest projection values but large intervals of low projection values. The other rows have lower projection values, but they are distributed over the columns. A simple horizontal addition of the projection values would therefore not lead to reliable rotation estimations.

The computation of a row of the Radon transform  $\mathbf{R}^T$  involves a transformation of all the pixels of the spot array image  $\mathbf{S}$ . In order to increase the efficiency, one should compute the Radon transform in a hierarchical manner, determine a reasonable maximum rotation angle  $\theta_M$  and an angle resolution  $\Delta\theta$ .

**Hierarchical Radon transform.** It is useful to compute  $\mathbf{R}^T$  in a *hierarchical* manner with increasing angle resolutions  $\Delta\theta$ :

1. Start at the initial resolution, for example  $\Delta\theta_0 = 1.0^\circ$ , and compute a Radon transform for the maximum rotation angle between  $-\theta_M$  and  $+\theta_M$  in order to get a coarse rotation estimate  $\theta_{G_0}$ .
2. Double the angle resolution to  $\Delta\theta_1 = \Delta\theta_0/2 = 0.5^\circ$ . As the accuracy of the initial angle estimate  $\theta_{G_0}$  is  $\pm\Delta\theta_0$ , or  $\pm 1^\circ$ , it is sufficient to compute 5 projections for the angle set  $\{\theta_{G_0} - 2\Delta\theta_1, \theta_{G_0} - \Delta\theta_1, \theta_{G_0}, \theta_{G_0} + \Delta\theta_1, \theta_{G_0} + 2\Delta\theta_1\}$ .
3. Repeat step 2 until the desired angle resolution  $\Delta\theta$  is reached. Every step except of the first one only requires the computation of 5 projections.

Table 4.1.1 illustrates the speedup that can be gained using the hierarchical approach.

Non-Hierarchical Radon transform				Hierarchical Radon transform			
Iteration	$\Delta\theta$	$\theta_G$	Projections	Iteration	$\Delta\theta$	$\theta_G$	Projections
0	$0.125^\circ$	$1.625^\circ$	54	0	$1.0^\circ$	$2^\circ$	9
				1	$0.5^\circ$	$1.5^\circ$	5
				2	$0.25^\circ$	$1.75^\circ$	5
				3	$0.125^\circ$	$1.625^\circ$	5
Total number of projections:			54	Total number of projections:			24

Table 3: Speedup for the hierarchical Radon transform with a maximum rotation angle of  $\theta_M = 4^\circ$ . For an angle resolution of  $\Delta\theta = 0.125^\circ$ , the non-hierarchical Radon transform needs 54 projections according to (31). The hierarchical approach gradually refines an initially coarse angle resolution ( $\theta_G$  shows the intermediate results for the grid rotation angle). Every refinement of the initial results only needs 5 projections, the total number of projections for the desired angle resolution therefore sums up to only 24 projections.

**Maximum Rotation Angle.** For some spot array images it is feasible to compute a maximum possible rotation which depends on the size of the physical filter and the size of the digital spot array image. However, the pixel dimensions of some spot array images are much larger than the pixel dimensions of the physical filter they comprise. As a consequence, a filter could have theoretically any rotation  $\theta_G$  in the digital image. We found empirically that no filter was rotated more than  $3^\circ$ . In order to have an additional tolerance we set  $\theta_M = 4^\circ$ .

**Angle Resolution.** The necessary angle resolution  $\Delta\theta$  depends on the size of the  $M \times N$  spot array image  $S$ . The bigger the image, the more orientations a straight line can have in the digital image. In the central coordinate system with the origin at the image center, the straight line with the minimal possible rotation has a  $\Delta x$  of  $N/2$  pixels and a  $\Delta y$  of 1 pixel. Since we are dealing with digital images with  $M > 1000$  and  $N > 1000$ , there is barely a difference between a straight lines with  $\Delta y = 1$  and a straight line with  $\Delta y = 2$ . We therefore fix the minimal  $\Delta y$  to 2 pixels and have

$$\Delta\theta = \frac{180}{\pi} \arctan\left(\frac{2}{N/2}\right) = \frac{180}{\pi} \arctan\left(\frac{4}{N}\right). \quad (35)$$

Applying the estimated rotation  $\theta_G$  to the set of prior guide spot locations  $L_p(\mathcal{G}^*)$  yields a set  $L_\theta(\mathcal{G}^*)$  of rotated prior guide spot locations.

#### 4.1.2 Global Translation Estimation

The second step on the way to define the reference locations  $L_R(\mathcal{G}^*)$  for the detected guide spot locations  $L_D(\mathcal{G}^*)$  is an estimation of how the rotated prior guide spot locations  $L_\theta(\mathcal{G}^*)$  need to be shifted. Formally, we must determine a global translation vector  $\mathbf{t}$  with the following property:

$$L_R(\mathcal{G}^*) = L_\theta(\mathcal{G}^*) + \mathbf{t} \approx L_D(\mathcal{G}^*). \quad (36)$$

It is sufficient to determine the location of one grid node of the reference grid  $L_R(\mathcal{G}^*)$ : If, for example, the translation  $\mathbf{t}_{UL}$  to the reference location  $L_R((1, 1))$  of the upper left node of  $\mathcal{G}^*$  is known, the shift for all the other grid nodes is the same. Formally, the translation vector  $\mathbf{t}_{UL}$  is given by

$$\mathbf{t}_{UL} = \widehat{L}_D((1, 1)) - L_\theta((1, 1)), \quad (37)$$

i.e. we subtract the rotated prior guide spot location of the upper left node of  $\mathcal{G}^*$  from an estimate for the detected guide spot location of the upper left corner node of  $\mathcal{G}^*$ . Note that the upper left location  $L_D((1, 1))$  is not yet available, since the determination of the grid of detected guide spot locations  $L_D(\mathcal{G}^*)$  is the task of the grid fitting itself. We must therefore determine an estimate  $\widehat{L}_D((1, 1))$  for the location  $L_D((1, 1))$  of the upper left corner node from the set  $\mathcal{L}^*$  of detected guide spot locations. In order to be more consistent, we also estimate the location of the lower right corner node  $(I_{GS}, J_{GS})$  of  $\mathcal{G}^*$  and define the translation vector  $\mathbf{t}_{LR}$  as

$$\mathbf{t}_{LR} = \widehat{L}_D((I_{GS}, J_{GS})) - L_\theta((I_{GS}, J_{GS})). \quad (38)$$

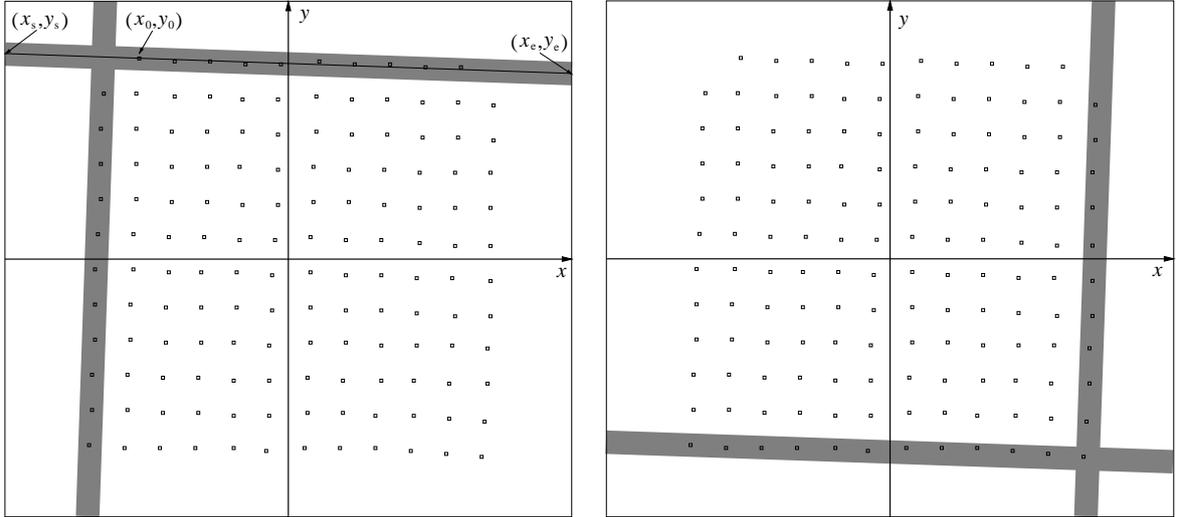
The translation vector in (36) is then determined as

$$\mathbf{t} = \text{mean}(\mathbf{t}_{UL}, \mathbf{t}_{LR}). \quad (39)$$

The location estimate  $\widehat{L}_D((1, 1))$  can be computed in the following way: try to extract from the set  $\mathcal{L}^*$  of detected guide spot locations the locations  $L_D(\mathcal{R}_{\mathcal{G}^*}(1))$  belonging to the first row and the locations  $L_D(\mathcal{C}_{\mathcal{G}^*}(1))$  belonging to the first column of  $\mathcal{G}^*$ .  $\widehat{L}_D((1, 1))$  is then the intersection point of the straight lines fitted to the locations of the first row and first column, respectively. Similarly, extract the locations  $L_D(\mathcal{R}_{\mathcal{G}^*}(I_{GS}))$  belonging to the last row and the locations  $L_D(\mathcal{C}_{\mathcal{G}^*}(J_{GS}))$  belonging to the last column of  $\mathcal{G}^*$  and intersect the corresponding fitted straight lines in order to estimate the location of the lower right grid node of  $\mathcal{G}^*$ .

**Extracting the First Row.** In the following, we describe an algorithm to extract the locations  $L_D(\mathcal{R}_{\mathcal{G}^*}(1))$  belonging to the first row of  $\mathcal{G}^*$  (see also Fig. 14):

1. Select as initial location  $(x_0, y_0) \in \mathcal{L}^*$  the detected guide spot location with the biggest  $y$ -coordinate in the central coordinate system.
2. Determine the end points  $(x_s, y_s)$  and  $(x_e, y_e)$  of a digital straight line crossing the initial location  $(x_0, y_0)$ . The straight line has a slope corresponding to the estimated global grid rotation angle  $\theta_G$  (Sect. 4.1.1). The straight line reaches from the first pixel column of the  $M \times N$  spot image  $\mathbf{S}$  to the last pixel column of the spot image.
3. Define an area  $\mathcal{A}$  in which to look for the detected locations belonging to the first row. The search area  $\mathcal{A}$  includes locations of the digital straight line between the end points  $(x_s, y_s)$  and  $(x_e, y_e)$ . These locations belonging to the digital straight line are determined



(a) Search area for first row and first column

(b) Search area for last row and last column

Figure 14: Determination of the points belonging to the first/last row/column. A search area  $\mathcal{A}$  (one of the four shaded bars) is spanned by a straight line with a slope corresponding to the estimated grid rotation angle going to the point with the smallest/biggest  $y/x$ -coordinate. The width of the search area corresponds to the theoretical block size.

by a modified Bresenham algorithm [3, 16]. Add also all locations to the search area  $\mathcal{A}$  belonging to straight lines between  $(x_s, y_s \pm \tau)$  and  $(x_e, y_e \pm \tau)$ , where  $\tau$  is a pixel tolerance parameter with  $\tau \in \{1 \dots \circ (B_y/2)\}$ . The height of the search area therefore corresponds to the theoretical block size.

4. The locations belonging to the first grid row of  $\mathcal{G}^*$  correspond to the intersection of the set of detected guide spot locations  $\mathcal{L}^*$  with the locations of the search area  $\mathcal{A}$ :

$$L_D(\mathbf{R}_{\mathcal{G}^*}(1)) = \mathcal{A} \cap \mathcal{L}^*. \quad (40)$$

The location set is only accepted if the number of locations in the intersection (40) is at least 50 % of the theoretical guide spot grid width  $J_{GS}$ . We therefore have the plausibility condition

$$\text{card}(L_D(\mathbf{R}_{\mathcal{G}^*}(1))) \geq J_{GS}/2. \quad (41)$$

Condition (41) is necessary to cope with outliers, which for example can be seen in Fig. 10. If inequality (41) does not hold, the extraction restarts with step 1. Before, all the locations previously assigned to  $L_D(\mathbf{R}_{\mathcal{G}^*}(1))$  are removed from the set of detected guide spot locations  $\mathcal{L}^*$ .

The locations  $L_D(\mathbf{R}_{\mathcal{G}^*}(I_{GS}))$  belonging to the nodes of the last row are determined in the same way except that the initial point  $(x_0, y_0)$  is chosen as the point with the smallest  $y$ -coordinate in the central coordinate system. Likewise, the locations belonging to the nodes of the first and last column  $L_D(\mathbf{C}_{\mathcal{G}^*}(1))$  and  $L_D(\mathbf{C}_{\mathcal{G}^*}(J_{GS}))$  are determined in a similar way, except that the digital straight lines reach from the first row to the last row of the spot image.

**Fitting Straight Lines.** Once the  $R$  locations  $(x_k, y_k) \in L_D(\mathbf{R}_{\mathcal{G}^*}(1))$ ,  $k \in \{1 \dots R\}$  belonging to the first row are determined, we can fit to the location data the parameters  $a_r$  and  $b_r$  of a straight line model

$$y(x) = y(x; a_r, b_r) = a_r + b_r x. \quad (42)$$

The fitting is performed in the standard least squares sense, i.e. we want to minimize the sum of the squares of the residuals  $e_k$  between the location data points  $(x_k, y_k)$  and the model:

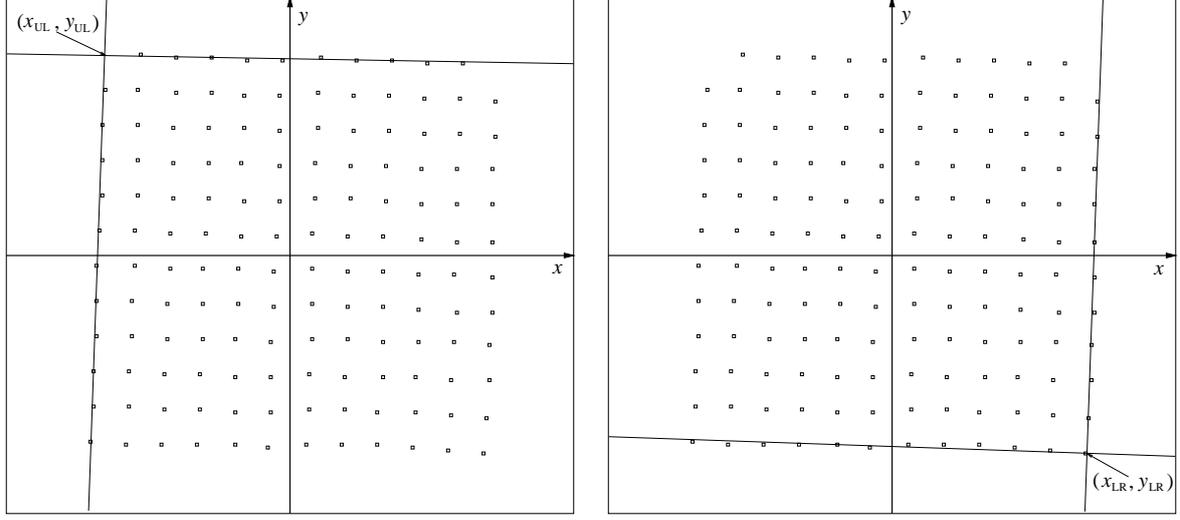
$$\sum_{k=1}^R e_k^2 = \sum_{k=1}^R (y_k - a_r - b_r x_k)^2 \rightarrow \min. \quad (43)$$

The optimal solution is given by [15, 5]

$$a_r = \mu_y - b_r \mu_x \quad \text{and} \quad b_r = \frac{\sigma_{xy}}{\sigma_x^2}. \quad (44)$$

The mean values  $\mu_x$  and  $\mu_y$  are computed as

$$\mu_x = \frac{1}{R} \sum_{k=0}^R x_k \quad \text{and} \quad \mu_y = \frac{1}{R} \sum_{i=k}^R y_k, \quad (45)$$



(a) Determination of upper left location

(b) Determination of lower right location

Figure 15: Determination of the location of the upper left grid node and the location of the lower right grid node of the guide spot grid. The locations are determined by the intersection of straight line models whose parameters are fitted to the locations belonging to the first/last row/column.

with  $(x_i, y_i) \in L_D(\mathbf{R}_{G^*}(I_{GS}))$ . The variance  $\sigma_x^2$  of the  $x$ -coordinates is given by

$$\sigma_x^2 = \frac{1}{R-1} \sum_{k=0}^R (x_k - \mu_x)^2, \quad (46)$$

and the covariance  $\sigma_{xy}$  between the  $x$ -coordinates and the  $y$ -coordinates is given by

$$\sigma_{xy} = \frac{1}{R} \sum_{k=0}^R (x_k - \mu_x)(y_k - \mu_y). \quad (47)$$

The parameters for the straight line of the last row are computed in the same way. As for the parameter computation of the straight line of the first and last column, the  $x$ - and  $y$ -coordinates are swapped. This is necessary because an unrotated grid would result in an infinite slope. A straight line of the form  $y' = a_c + b_c x'$  in the swapped  $x', y'$ -coordinate system has the form  $y = -a_c/b_c - 1/b_c x$  in the original coordinate system (provided that  $b_c \neq 0$ ).

**Intersecting straight lines.** After having determined the straight line parameters  $a_r$  and  $b_r$  for the first row and  $a_c$  and  $b_c$  for the first column, the intersection point can be determined with the following equation:

$$a_r + b_r x = -\frac{a_c}{b_c} - \frac{1}{b_c} x \quad \text{for } b_c \neq 0. \quad (48)$$

Note that the right hand side of the equation is the straight line of the first column transformed to the coordinate system of the straight line of the first row. After some manipulation we have the

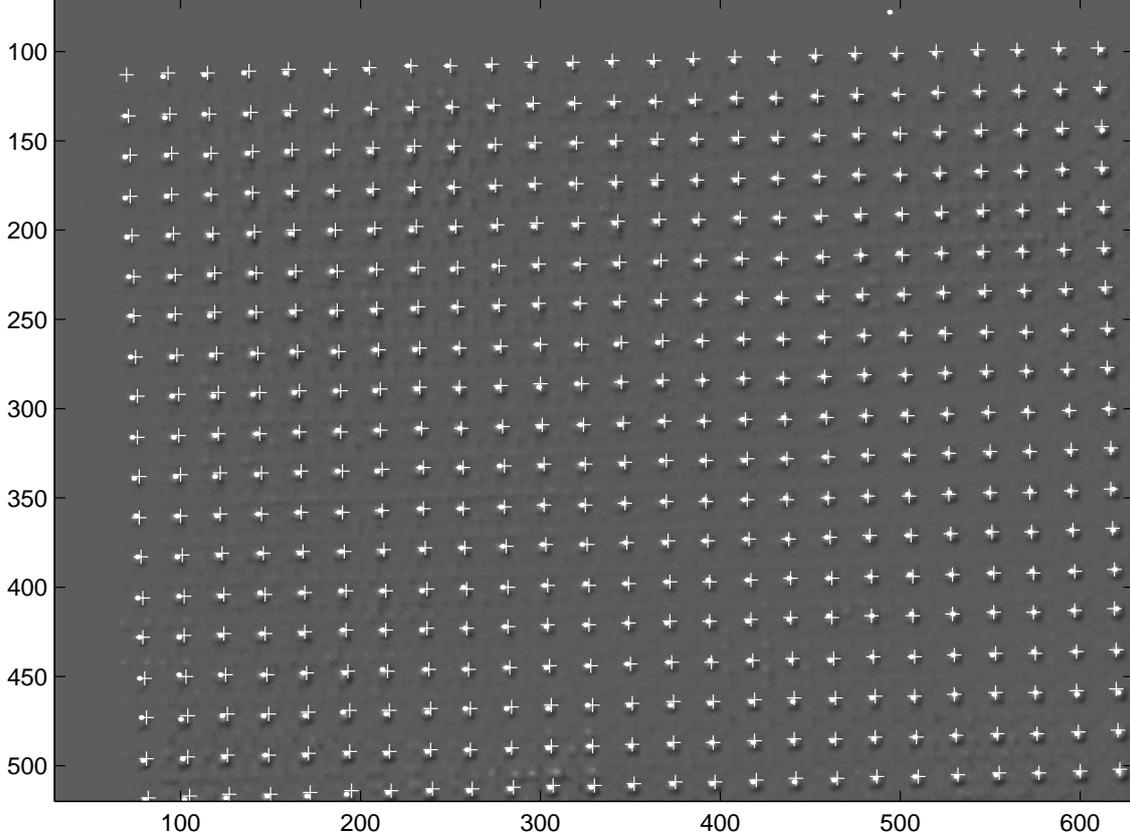


Figure 16: Reference guide spot locations. The grid of crosshairs illustrates the reference guide spot locations, which are the prior guide spot locations of Fig. 8 rotated by the estimated grid rotation angle  $\theta_G$  and translated by the translation vector  $\mathbf{t}$ . The detected guide spot locations (illustrated as white dots) can be easily assigned to the corresponding nodes of the guide spot grid  $\mathcal{G}^*$ .

following coordinates for the intersection point  $(x_{UL}, x_{UL}) = \widehat{L}_D((1, 1))$  of the upper left corner:

$$x_{UL} = \begin{cases} \left( -a_r - \frac{a_c}{b_c} \right) / \left( b_r - \frac{1}{b_c} \right) & \text{for } b_c \neq 0 \\ a_c & \text{for } b_c = 0 \end{cases} \quad (49)$$

$$y_{UL} = a_r + b_r x_{UL}. \quad (50)$$

The intersection point  $(x_{LR}, y_{LR}) = \widehat{L}_D((I_{GS}, J_{GS}))$  for the lower right corner is computed in a similar manner to (49) and (50). The intersection points  $(x_{UL}, x_{UL})$  and  $(x_{LR}, x_{LR})$  are illustrated in Fig. 15. The final translation vector  $\mathbf{t}$  is the mean vector according to (39).

### 4.1.3 Alignment of locations to grid nodes

Having an estimate of the rotation angle  $\theta_G$  and the translation vector  $\mathbf{t} = [t_x \ t_y]^T$ , we are able to linearly transform the prior guide spot locations  $[x_p \ y_p]^T \in \mathbf{L}_P(\mathcal{G}_P^*)$  to the reference guide spot

locations  $[x_r \ y_r]^T \in \mathbf{L}_R(\mathcal{G}_p^*)$  according to (28). The reference guide spot locations  $\mathbf{L}_R(i, j)$  are expected to be near the detected guide spot locations  $\mathcal{L}^*$  (see Fig. 16). The detected guide spot locations can therefore be assigned to guide spot grid nodes  $(i, j)$  by investigating a rectangular  $M_L \times N_L$  location window  $\mathbf{W}(\mathbf{L}_R(i, j))$  with  $\mathbf{L}_R(i, j)$  in the window center:

$$\mathbf{L}_D(i, j) = \mathbf{W}(\mathbf{L}_R(i, j)) \cap \mathcal{L}^*. \quad (51)$$

The dimensions  $M_L$  and  $N_L$  of  $\mathbf{W}$  correspond to the size of maximum search window as defined in (24) and (25).

## 4.2 Consistent Spot Grid

The algorithm for the detection of guide spots in Sect. 3 tries to minimize the number of false positives. Nevertheless, it cannot be guaranteed that all the detected locations in  $\mathcal{L}^*$  belong to guide spots (see also Fig. 10). We must therefore find a way to eliminate the false positives after the grid alignment of Sect. 4.1. Additionally, in order to have a consistent grid, we must eliminate the false negatives, i.e. compute the locations of the guide spots that were not detected. Both tasks - the elimination of false positives and false negatives - can be accomplished by the robust fitting of straight lines to the rows and columns of the fields. Finally, once a consistent guide spot grid  $\mathcal{G}^*$  is available, we must initialize the locations of the regular spots.

### 4.2.1 Parameterization of the Guide Spot Grid

It is possible that during the assignment of the detected locations to the grid nodes (Sect. 4.1) no location can be found in the search window (Eqn. (51)). Formally, this fact is expressed as

$$\mathbf{L}_D(i, j) = \mathbf{W}(\mathbf{L}_R(i, j)) \cap \mathcal{L}^* = \emptyset. \quad (52)$$

We call a node  $(i, j)$  as *invalid* if condition (52) holds. If a node  $(i, j)$  is invalid we must estimate  $\mathbf{L}_D(i, j)$  from a subset of the aligned guide spot locations  $\mathbf{L}_D(\mathcal{G}^*)$ .

One possible approach is to fit straight lines to every row and column of the guide spot grid  $\mathcal{G}^*$  and estimate the lacking location by the intersection of the straight line of row  $i$  and the straight line of column  $j$ . This would be equivalent to the computations of Sect. 4.1.2, in which the locations of the upper left and the lower right node are estimated. However, different fields  $\mathcal{F}_{pq}^*$  can be significantly shifted. It is therefore a good idea to perform straight line fits on the fields  $\mathcal{F}_{pq}^*$  of the guide spots instead of the whole guide spot grid  $\mathcal{G}^*$ . The parameter set  $\mathcal{P}_{pq}$  for a guide spot field  $\mathcal{F}_{pq}^*$  is defined by

$$\mathcal{P}_{pq} = \{((a_{r_i}, b_{r_i}), (a_{c_j}, b_{c_j})) \mid 1 \leq i \leq I_F, 1 \leq j \leq J_F\}, \quad (53)$$

with the parameters  $a$  and  $b$  as defined in (42). If a node  $(i, j)$  of a guide spot field is not valid, the parameters  $((a_{r_i}, b_{r_i}), (a_{c_j}, b_{c_j}))$  modeling field row  $i$  and field column  $j$  can be used to perform straight line intersection as introduced in (49) and (50).

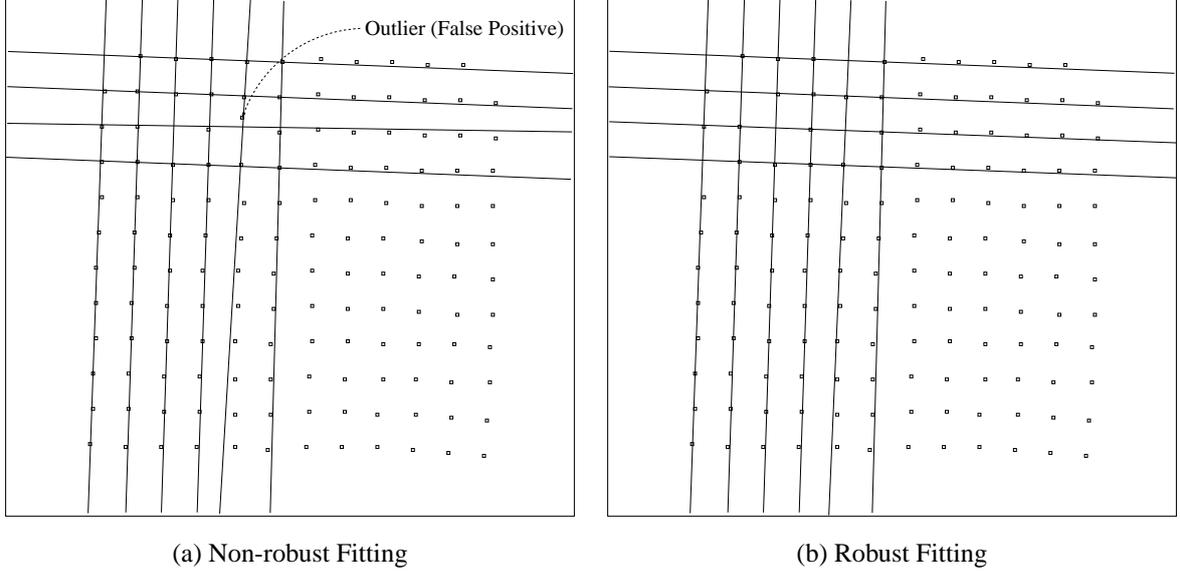


Figure 17: Example for the parameterization of the upper left  $4 \times 6$  guide spot field grid. In order to infer lacking guide spot locations (false negatives) from neighborhood locations, straight lines are fitted to every row and every column of the guide spot field grid. (a) An outlier (false positive) biases the straight line models such that intersections in which the biased straight line models are involved would be inaccurate. (b) After the removal of the outlier (position (2, 5)), the corresponding straight line models are no longer biased. Outliers are detected with the help of residual values from the model to the data.

#### 4.2.2 Robust Grid Parameter Fitting

The field row and field column parameters as defined in (53) must be fitted in a *robust* manner because false positives (outliers) can bias the fitted straight lines (see Fig. 17). One possible approach to perform robust fitting is to consider the residuals  $e_k$  of the fitted straight line model and the data points (see (43)). The algorithm for the fitting of a straight line belonging to row  $i$  of a guide spot field is outlined as follows (the application to column parameters is straightforward):

1. **Initialization.** Fit the parameters  $a_{r_i}$  and  $b_{r_i}$  according to (44) if  $K_0 \geq 2$ , where  $K_0$  is the number of valid nodes in row  $i$  (see (52)). If  $K_0 < 2$  go to 5.
2. **Large residual removal.** Sort the residuals  $e_k$  of the initial fit and mark the nodes with the  $\rho$  highest residuals as invalid, where  $\rho$  is a percentage  $P_1$  of the valid guide spots  $K_0$ .
3. **Refitting.** Re-fit the parameters  $a_{r_i}$  and  $b_{r_i}$  if still at least 2 nodes in the row  $i$  are valid; if fewer than 2 nodes are valid, go to 5.
4. **Outlier removal.** Sort the residuals  $e_k$ . If the largest residual is below a threshold  $t$  [pixel], the fitting is finished and the parameters are valid; otherwise mark the corresponding node as invalid. If more than 50 % of the initial number of valid nodes  $K_0$  are valid, go to 3. Else go to 5.
5. **Algorithm abortion.** Mark the *parameters* as invalid and return.

A reasonable choice is to set  $P_1 = 10$ , meaning that the straight lines are re-fitted at least once with 10 % fewer nodes than at the first fit (the 10 % with the highest residuals). After the re-fit, the absolute values of the residuals are regarded. The parameters are re-fitted

- if the distance between a model point and a data point exceeds a threshold  $t$  [pixels] and
- if not more than 50 % of the initial valid nodes have already been marked as invalid.

The threshold  $t$  is a percentage  $P_2$  of the theoretical block sizes  $B_x$  and  $B_y$ , respectively, as defined in (13). It is reasonable to also set  $P_2 = 10$ . The 50 % limit is set because an original model with more than 50 % outliers cannot be restored with this kind of fitting procedure.

### 4.2.3 Field Parameter Correction

After the fitting of the guide spot field rows and columns, the parameters are checked:

1.  $a_{r_i}$  and  $b_{r_i}$  might be marked as invalid because of the failure of the straight line fitting algorithm
2. the absolute value of the slope  $b_{r_i}$  might differ too much ( $0.5^\circ$ ) from the absolute mean value of all the valid slopes of the field rows. This situation may occur since it is not guaranteed that the node with the highest residual value must necessarily be the outlier. Sometimes it is therefore possible that - even after the repeated refitting - false positives do survive.

If 1) or 2) holds for the parameters  $a_{r_i}$  and  $b_{r_i}$  of a row, they are estimated with the help of the parameters of the nearest neighborhood row  $a_{r_n}$  and  $b_{r_n}$ . Since we do not expect large variation between neighbored slopes we can set  $b_{r_i} = b_{r_n}$ . The intercept  $a_{r_i}$  can be estimated by setting

$$a_{r_i} = a_{r_n} - (n - i)B_y, \quad (54)$$

i.e. subtracting the theoretical block distances from the neighbor depending on how far away the neighbor on the grid is.

### 4.2.4 Abortion Criterion

After the robust determination of the field parameter sets (53), the entire guide spot grid  $\mathcal{G}^*$  is tested for consistency. If – due to a very bad image quality – the final guide spot grid is not plausible, the distance between the locations of at least two nodes must be too large. Formally, there must exist at least one node  $(i, j) \in \mathcal{G}^*$  for which the following condition holds:

$$\|\mathbf{L}((i, j)) - \mathbf{L}((i + 1, j + 1))\| > t \quad (55)$$

where  $t$  is the residual threshold entity of Sect. 4.2.2. If (55) holds, the grid fitting is aborted and the user is notified.

### 4.2.5 Location initialization of regular spots

Locations belonging to regular spots in a block  $\mathcal{B}$  are inferred from the guide spot location with the help of the prior knowledge of the theoretical spot distances  $S_x$  and  $S_y$  and the block rotation. The rotation  $\theta_{\mathcal{B}}$  of a block  $\mathcal{B}$  in row  $i$  of a field  $\mathcal{F}$  is given by the slope  $b_{r_i}$  of the field grid parameter set (53) as follows:

$$\theta_{\mathcal{B}} = \arctan(b_{r_i}). \quad (56)$$

Given the guide spot  $(i, j) \in \mathcal{B}$  and its location  $\mathbf{L}((i, j)) = [x_G \ y_G]^T$ , the regular spot locations  $\mathbf{L}((m, n))$  with  $(m, n) \in \mathcal{B}$  and  $(m, n) \neq (i, j)$  are initialized as follows:

$$\mathbf{L}((m, n)) = \begin{bmatrix} x_G \\ y_G \end{bmatrix} + \begin{bmatrix} \cos \theta_{\mathcal{B}} & \sin \theta_{\mathcal{B}} \\ -\sin \theta_{\mathcal{B}} & \cos \theta_{\mathcal{B}} \end{bmatrix} \begin{bmatrix} (m - i)S_x \\ (n - j)S_y \end{bmatrix} \quad \forall (m, n) \neq (i, j). \quad (57)$$

Equation (57) is valid for blocks with guide spots residing at an arbitrary position of the block. The regular spot locations (57) are used as initial estimates of the center of a parametric spot model.

## 4.3 Summary

The goal of the grid fitting step is to map potential guide spot locations detected in Sect. 3 to the guide spot grid nodes and to initialize the locations of the regular spots. The prior guide spot locations introduced in Sect. 3.3.1 are rotated and translated in order to serve as reference locations for the detected guide spot locations. The global rotation is estimated with the help of a Radon transform which is based on the projection of the image intensity values along different directions. In order to save computation time, the Radon transform can be computed in a hierarchical manner, in which the angle resolutions are successively refined. The global translation of the prior guide spot locations is determined with the help of the upper left and lower right guide spot grid corner points. These corner points are determined with the help of straight lines which are fitted in the least squares sense to the points belonging to the first/last row/column. In order to get a consistent grid, false negatives (lacking guide spots) and false positives (outliers) are removed with the help of robust fitting of straight lines to the field rows and column. If the residual values from the straight line model to the data are too large, the corresponding nodes are marked as invalid which is equivalent to the removal of false positives. A lacking guide spot location can be estimated with the help of the intersection of the corresponding straight lines. Sometimes the parameters of the straight line model have to be estimated from the nearest neighbor. This is the case for images with very bad quality or for image regions in which a large majority of the spots has a strong hybridization signal.

## 5 Experimental Results

The quality of the grid fitting cannot be assessed with a simple location distance measure, since *no ground-truth data* is available for the spot array images. We haven chosen two other ways to demonstrate the effectiveness of the grid fitting algorithm presented in this report. We first show five examples of image types originating from different hybridization experiments, having different quality, resolution and size. We then show that the grid fitting success is correlated with the image quality.

### 5.1 Visual Examples

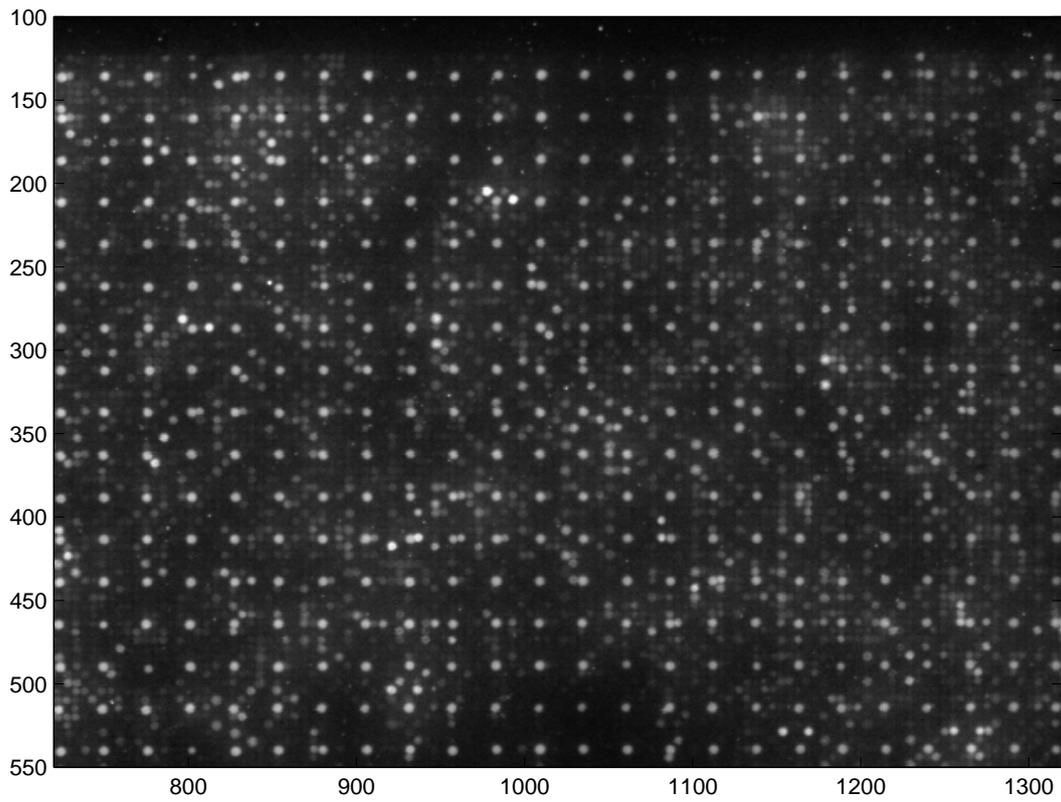
We present five examples of images for which the grid fitting was successful. In order to demonstrate the different scanning resolutions of the images, the image parts in Fig. 18- 22 are of the same size and in scale, meaning that the spots of a high-resolution image are displayed larger than the spots of a low-resolution image.

Figure 18a shows a part of a  $1596 \times 1482$  ONF image (Sect. 1) which has been scanned at a resolution of  $175\mu m$  at the Max Planck Institute for Molecular Biology (MPIMG) Berlin. The guide spots at the center of the  $5 \times 5$  blocks in Fig. 18a are bright and clearly identifiable. Figure 18b shows the same image with the computed guide spot locations superimposed as cross-hairs. For the sake of overview, the initialized locations of the regular spots are not shown – they are simply derived from the guide spot locations as demonstrated in (57). Please note that the locations need not necessarily be right in the center of the (guide) spot: They are just the initializations for the center of a parametric spot model.

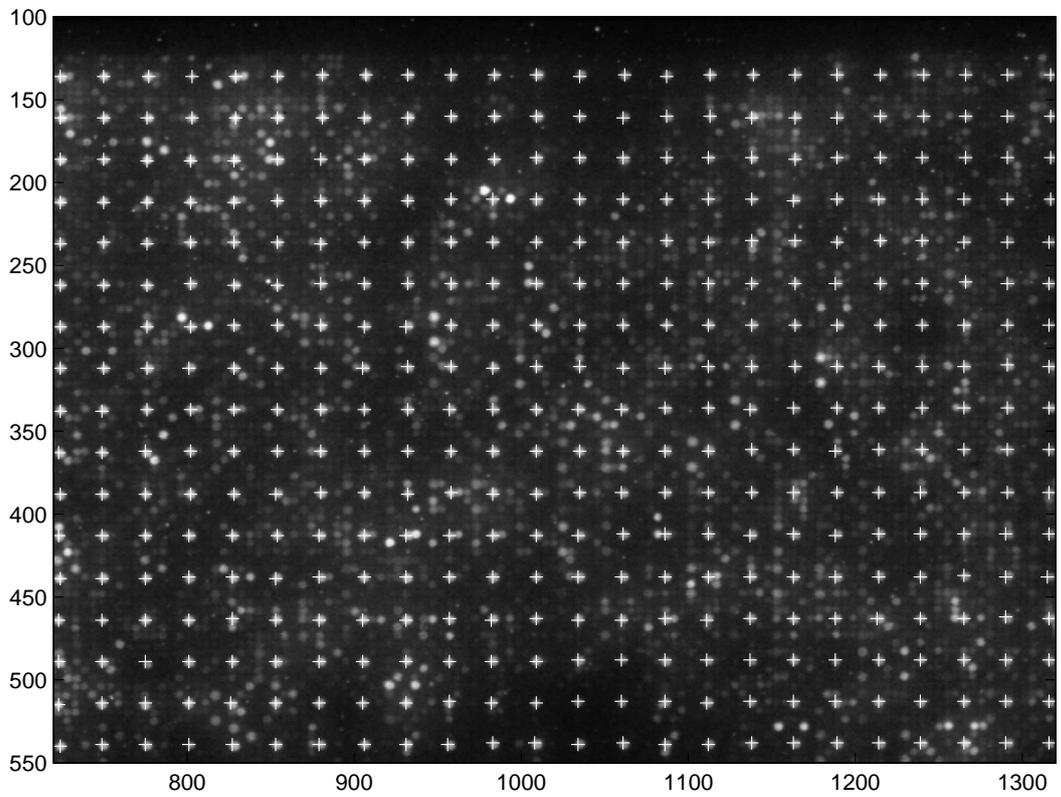
Figure 19a shows a part of a  $1300 \times 1200$  ONF image which has been scanned at a resolution of  $200\mu m$  at the Novartis Forschungsinstitut (NFI) Vienna. The signals of the hybridization signals of the guide spots at the center of the  $5 \times 5$  blocks are relatively low in comparison with the signals of the hybridized regular spots, for example at the lower right corner of the filter. Furthermore, there are regions in which the signal-to-noise ratio is very low. It can be seen in Fig. 19b that our algorithm is able to restore the guide spot grid. In this example, the border between two fields can be noticed by comparing guide spot columns 4 and 5 of Fig. 19b: There is a leap in the  $y$ -coordinates indicating a field shift and therefore justifying the parameterization (53) of the fields.

Figure 20a shows a part of a  $1300 \times 1586$  image originating from hybridizations of complex cDNA samples. It has been scanned at a resolution of  $200\mu m$  at the NFI Vienna. The grid is nearly full, but the guide spots at the center of the  $5 \times 5$  blocks are brighter than the majority of the regular spots. Note the visible vertical field shift in pixel row 800 of the image. The correct grid fitting output can be seen in Fig. 20b.

Figure 21a shows a part of a  $1300 \times 1486$  image originating from colony filter hybridizations. The image has been scanned at a resolution of  $200\mu m$  at the NFI Vienna. The upper left field is clearly noticeable. The intensities of the of the guide spots at the center of the  $5 \times 5$  blocks differ significantly: They are very high at the first row and first column of the guide spot grid and are partly not distinguishable from the regular spots in regions within the field. The dark rectangular region around pixel (375, 250) in the image indicates that a needle was lacking (broken) on the needle matrix (Fig. 2a). Due to the parameterization (53) of the fields, such lacking guide spot grid information can be easily restored, as is demonstrated in Fig. 21b.

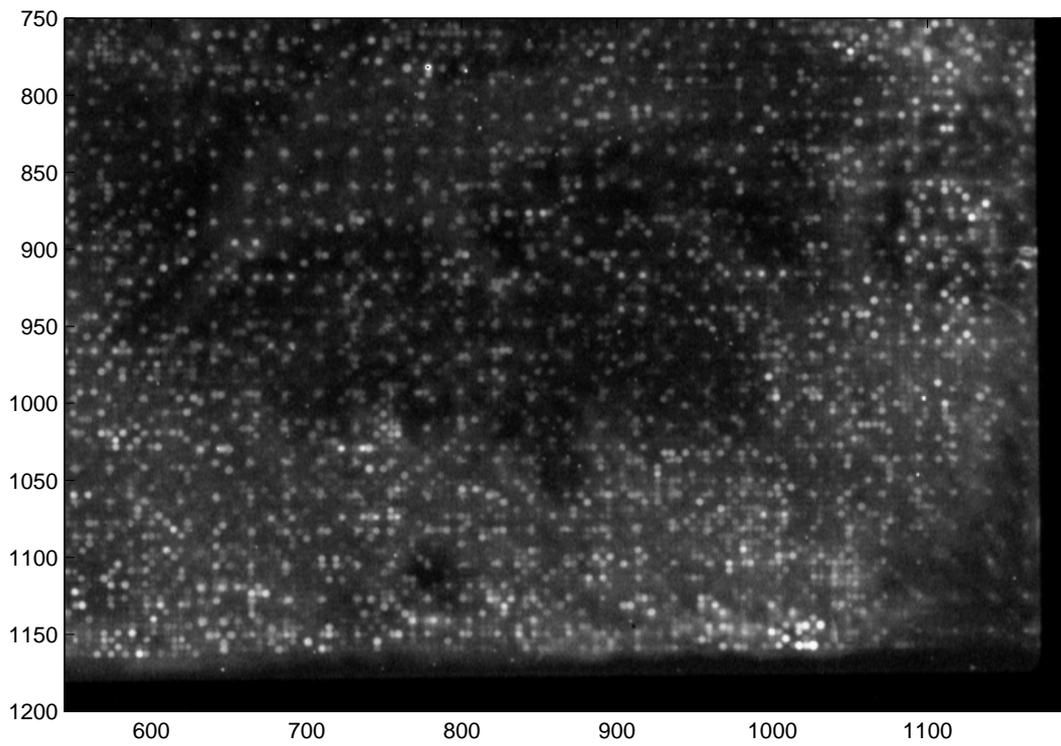


(a) Intensity Data

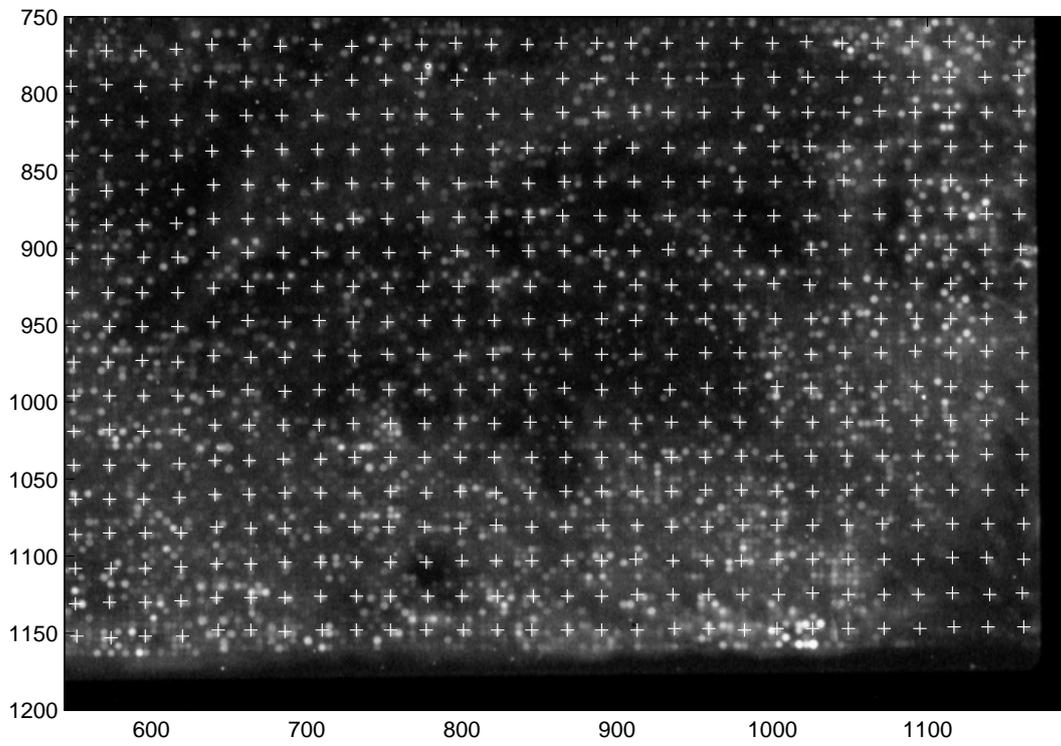


(b) Intensity Data + Guide Spot Locations

Figure 18: Part of a  $1596 \times 1482$  ONF Image with  $175\mu m$  resolution scanned at MPIMG Berlin.

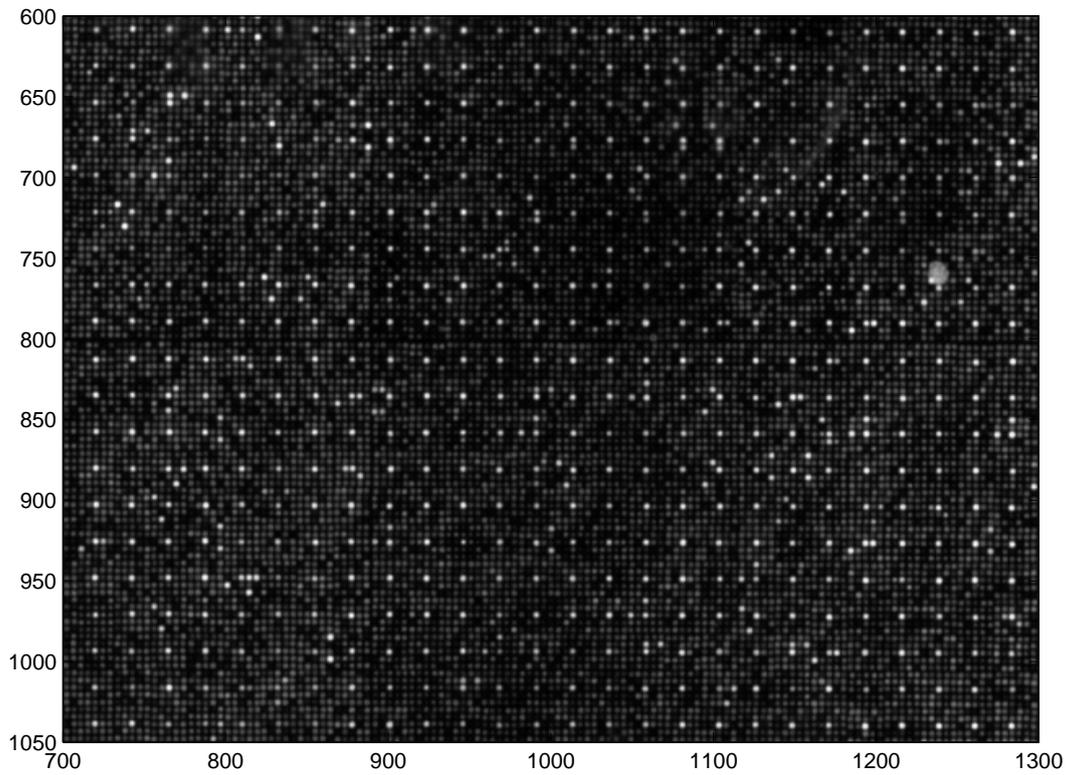


(a) Intensity Data

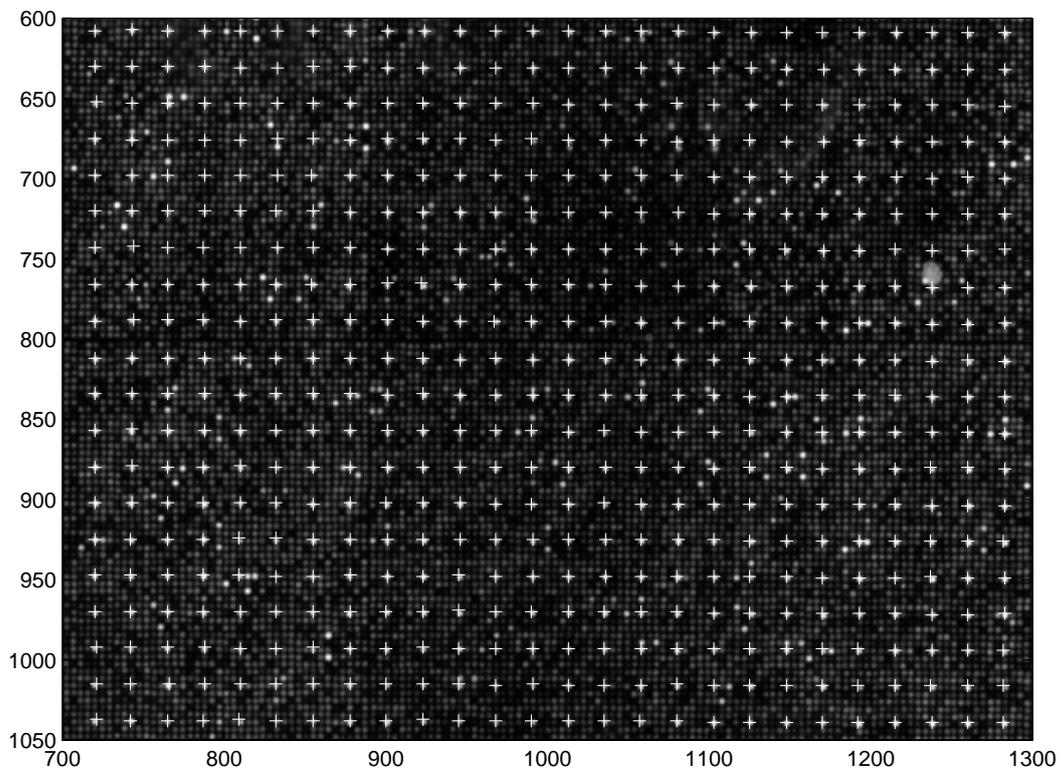


(b) Intensity Data + Guide Spot Locations

Figure 19: Part of a  $1300 \times 1286$  ONF Image with  $200\mu m$  resolution scanned at NFI Vienna.

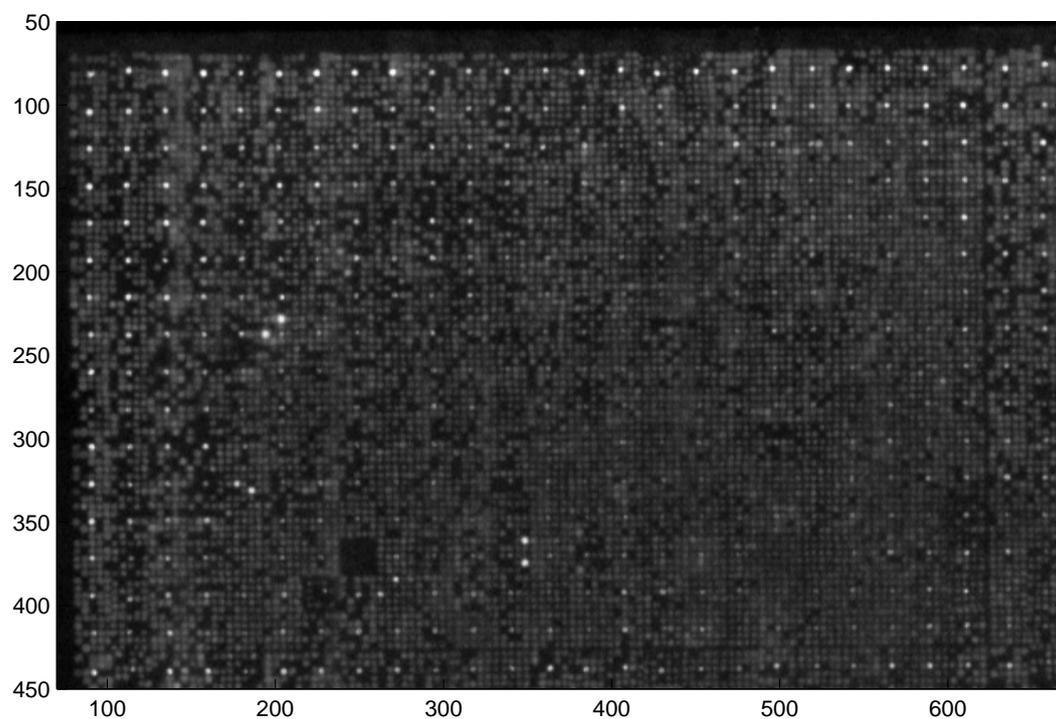


(a) Intensity Data

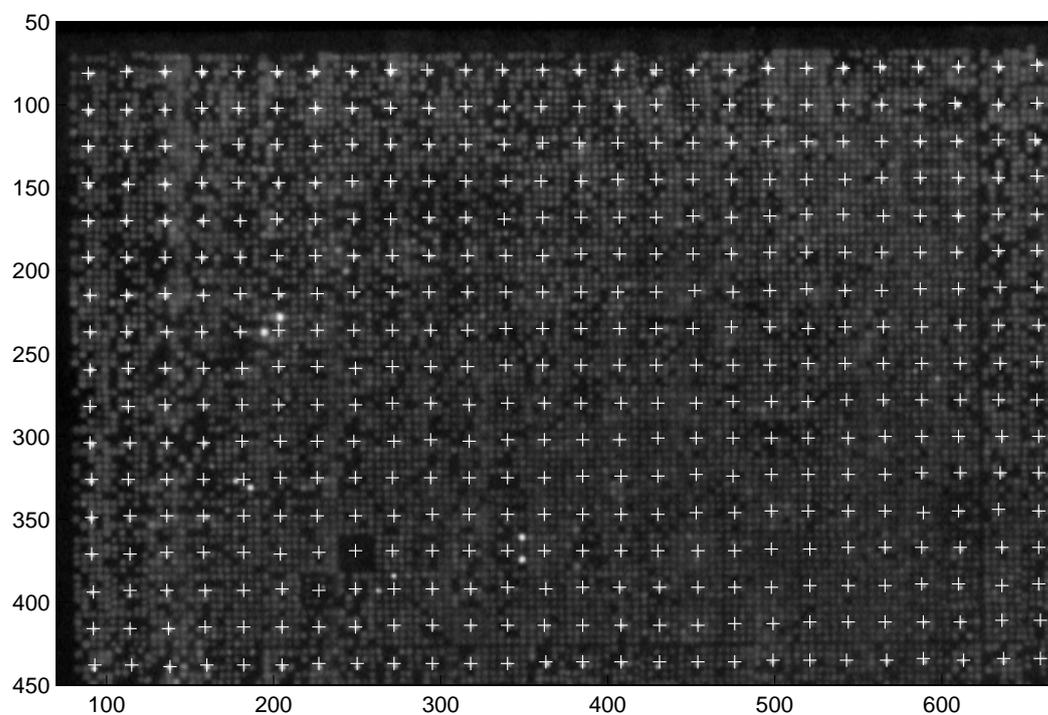


(b) Intensity Data + Guide Spot Locations

Figure 20: Part of a  $1300 \times 1586$  ComplexHyb Image with  $200\mu m$  resolution scanned at NFI Vienna.

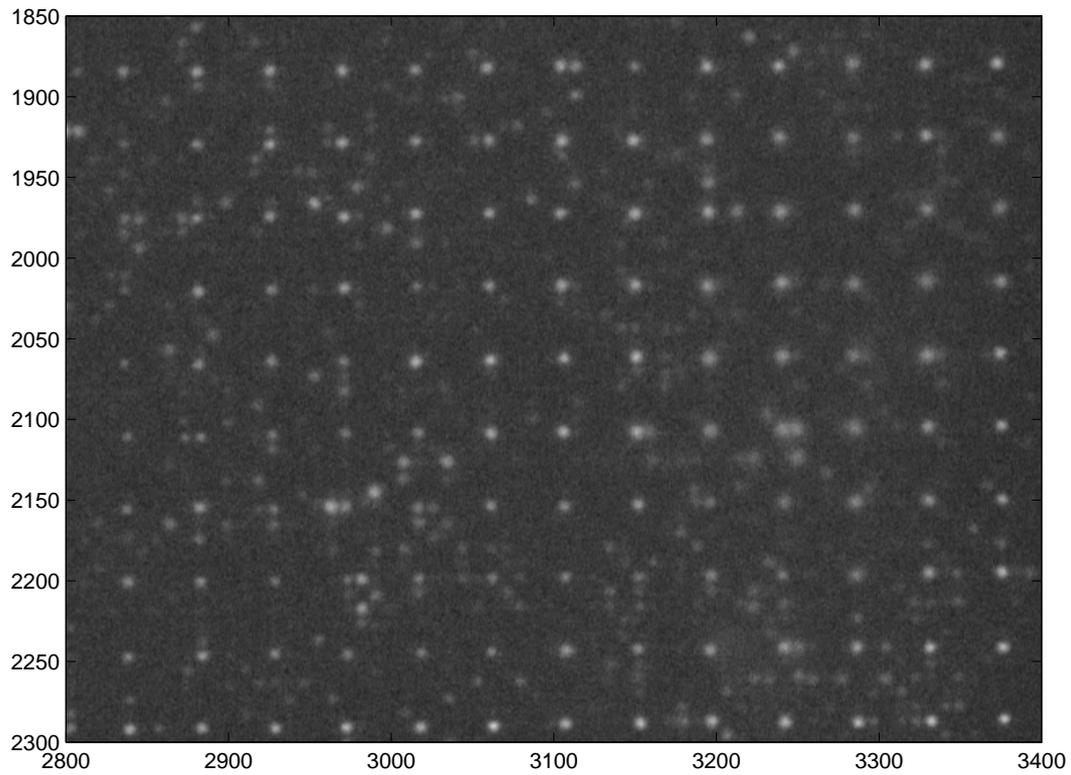


(a) Intensity Data

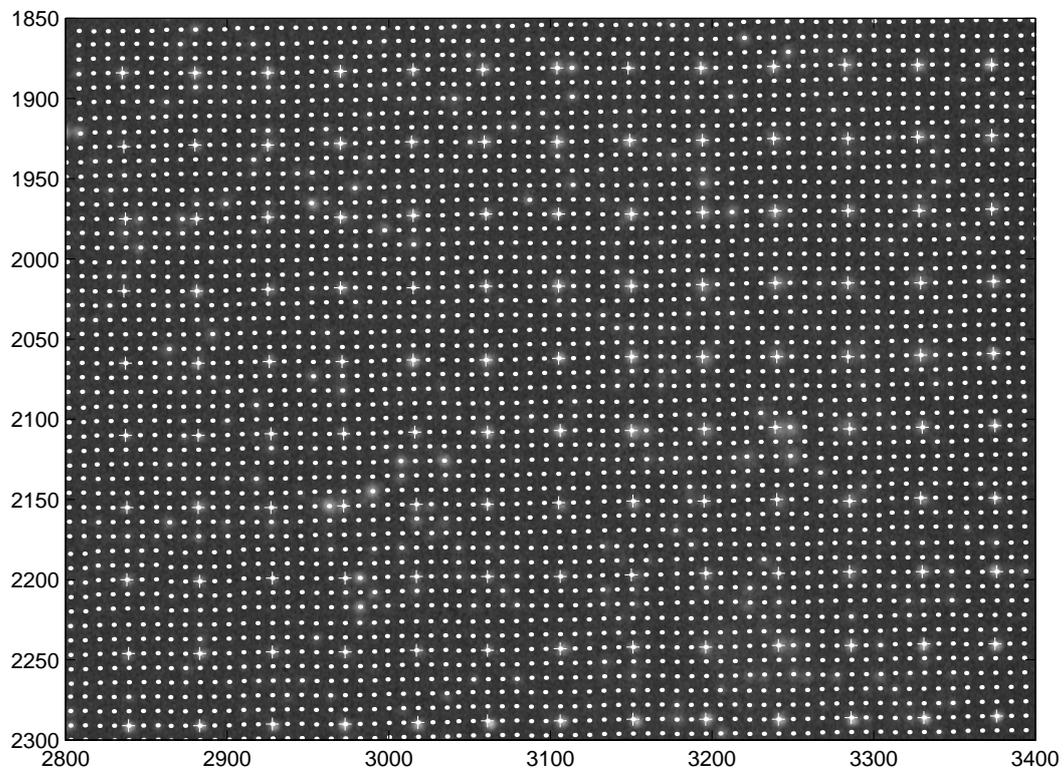


(b) Intensity Data + Guide Spot Locations

Figure 21: Part of a  $1300 \times 1486$  ColonyHyb Image with  $200\mu m$  resolution scanned at NFI Vienna.



(a) Intensity Data



(b) Intensity Data + spot locations

Figure 22: Part of a  $2400 \times 3544$  ComplexHyb Image with  $100\mu m$  resolution scanned at NFI Vienna.

Fig. 22a shows a part of a  $2400 \times 3544$  image originating from hybridizations of complex cDNA samples. It was scanned at a resolution of  $100\mu m$  at the NFI Vienna. It is an example for a low signal-to-noise ratio hybridization image. Due to the high resolution of the image, Fig. 22b also shows the computed locations of the regular spots superimposed as dots.

Figure	Image Name	Experiment	Resolution [ $\mu m$ ]	Size
18	o163_04260_A1	ONF	175	$1596 \times 1482$
19	o100_295107_y1	ONF	200	$1300 \times 1286$
20	coctail2_c64102_y1	ComplexHyb	200	$1300 \times 1586$
21	990401ptaa_b111dk_y1	ColonyHyb	200	$1300 \times 1486$
22	u266-1-99031_c64120_x1	ComplexHyb	100	$2400 \times 3544$

Table 4: Overview of the image examples demonstrated in Fig. 18–22.

Table 5.1 summarizes the image information and indicates the names of the image files.

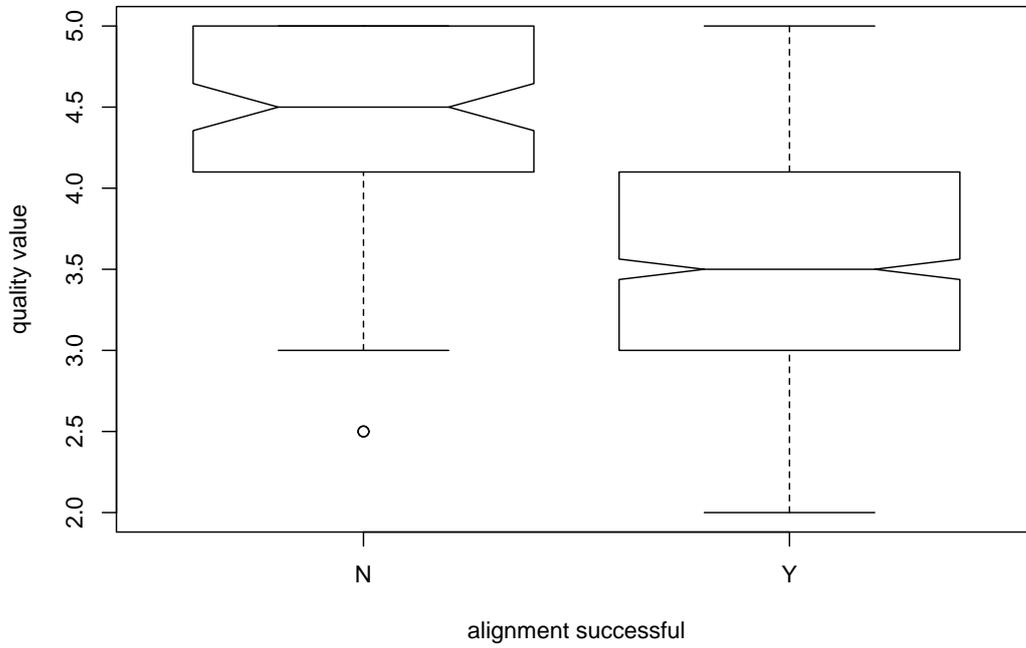
## 5.2 Evaluation of Image Quality versus Grid Fitting Success

The largest test set consisted of ONF images from two cDNA-libraries named w08 (855 images) and w09 (885 images). These ONF images had been scanned at the Novartis Research Institute Vienna at a resolution of  $\Delta x = \Delta y = 200\mu m$ . The image quality of every image in w08 and w09 had been rated by a human with numbers between 1 (very good) and 5 (very bad). Hence it is at least possible to investigate the correlation between the image quality and the success of the grid fitting (abortion criterion (55)). Figure 23 shows both for w08 and w09 images the Box plot of image qualities for which the grid fitting fails (left hand side) and the Box plot of image qualities for which the grid fitting was successful (right hand side). Figure 23 can be interpreted as follows:

- As expected, the grid fitting is successful for images with good quality. However, the algorithm can also cope with some images the quality of which was rated as very bad.
- Most of the images for which the grid fitting fails are rated as bad. There are, however, some outliers of images with good quality which do not meet the expectations.

The good quality images for which the grid fitting failed have one common feature: They violate the assumption stated on page 16 that the image border of the spot array image has significantly lower intensity values than the filter region. The background values of non-hybridized regions are as dark as the image border regions. As a consequence, the regions  $B1$  and  $B2$  in Fig. 7 do not contain all the lowest projection values and the prior guide spot locations are wrong. Furthermore, whole rows of guide spots will not be detected because of the wrong region-of-interest defined by the prior guide spot locations, resulting in a wrong guide spot grid. These rare cases can be handled by trying several different prior guide spot locations. There should be not too much computational overhead for this solution, since the most expensive computations (matched filter, GSLA filter and Radon transform) must not be recomputed.

**w08: Image quality evaluation vs grid alignment success**



**w09: Image quality evaluation vs grid alignment success**

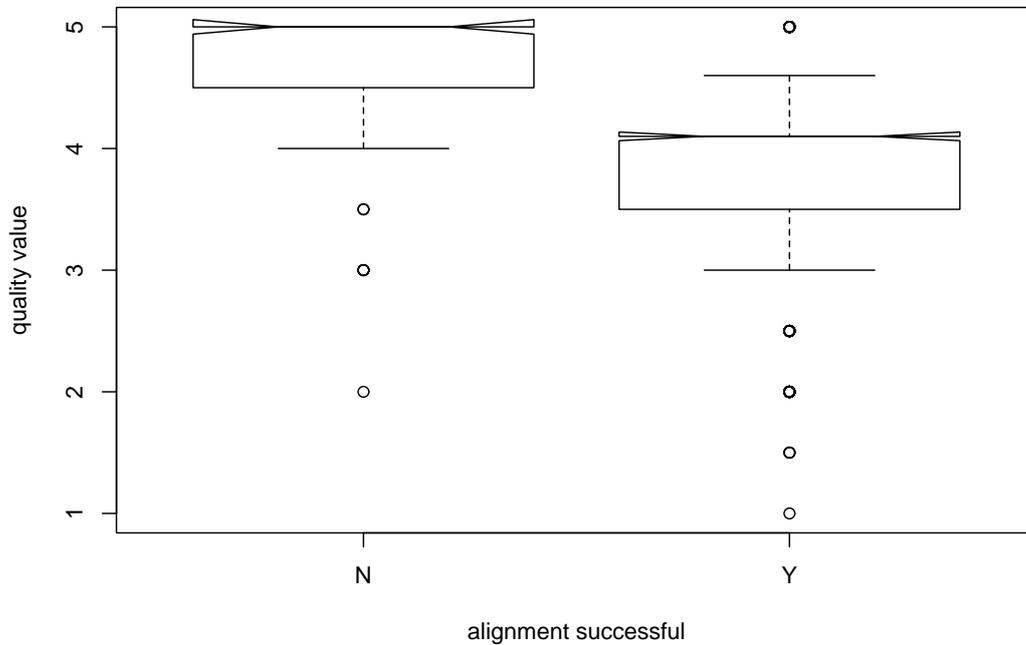


Figure 23: Box plots of image qualities versus grid fitting success. The grid fitting is successful for images with good quality. The algorithm can also cope with some images the quality of which was rated as very bad. Most of the images for which the grid fitting fails are rated as bad. There are, however, some outliers of images with good quality which do not meet the expectations

## 6 Conclusion and Outlook

We have presented a grid fitting approach for genetic spot array images containing guide spots. In the first main step, the guide spots are detected with a matched filter, a filter which amplifies the guide spot locations and a maximum search. In the second step, the detected guide spots are used to span a consistent guide spot grid. The fitted spot locations of our algorithm are considered as spot center initializations for parametric spot models [1, 4].

The crucial point of our algorithm is the region-of-interest (Sect. 3.3.1) which is determined with the help of a horizontal and a vertical projection. Experimental results have shown that our prior assumption of the spot image intensity distribution (much darker image border than filter area) does not always hold. In these cases, our algorithm possibly rejects spot images of good quality. The problem can be solved as follows:

- The horizontal projection is replaced with the GSLA projection of the correct angle which is given by the Radon transform. Most peaks in the projection will then result from guide spots.
- The vertical projection is replaced with the GSLA projection of the correct angle, preferably computed with an additional Radon transform.
- The two projections are considered simultaneously in order to solve an energy maximization problem, where the energy contributions result from the intersections of the straight lines belonging to guide spot location hypotheses.
- Non-linearities of the straight line intersections are considered by a local maximum search in the matched filter response image.

The approach outlined above is applicable with minor changes to spot array images which do not contain guide spots.

## References

- [1] N. Brändle, H-Y. Chen, H. Bischof, and H. Lapp. Robust parametric and semi-parametric spot fitting for spot array images. In *ISMB-2000 8th Intl. Conference on Intelligent Systems for Molecular Biology, August 20–23*, page 46, 2000. 5, 43
- [2] N. Brändle, H. Lapp, and H. Bischof. Automatic Grid Fitting for Genetic Spot Array Images Containing Guide Spots. In *8th Intl. Conf. on Computer Analysis of Images and Patterns, Ljubljana, Slovenia, September 1–3*, pages 357–366, 1999.
- [3] J. E. Bresenham. Algorithm for Computer Control of a Digital plotter. *IBM Systems Journal*, 4(1):25–30, 1965. 28
- [4] H-Y. Chen, N. Brändle, H. Bischof, and H. Lapp. Robust spot fitting for genetic spot array images. In IEEE Signal Processing Society, editor, *ICIP-2000 Intl. Conference on Image Processing, September 10–13, Vancouver, Canada, 2000*. 5, 43

- [5] Vladimir S. Cherkassky and Filip M. Mulier. *Learning from Data : Concepts, Theory, and Methods*. John Wiley and Sons, 1998. 28
- [6] M. Chee *et al.* Accessing Genetic Information with High-Density DNA Arrays. *Science*, 274:610–614, 1996. 2, 2
- [7] R. J. Johnston *et al.* Autoradiography using storage phosphor technology. *Electrophoresis*, 11:355–360, 1990. 2
- [8] K. Hartelius. *Analysis of Irregularly Distributed Points*. PhD thesis, Institute of Mathematical Modelling, Technical University of Denmark, 1996. 5
- [9] Anil K. Jain. *Fundamentals of Digital Image Processing*. Prentice-Hall, 1986. 23
- [10] E. Lander. The new genomics: Global views of biology. *Science*, 274:536–539, 1996. 2
- [11] Benjamin Lewin. *Genes VI*. Oxford University Press, 1997. 2
- [12] S. Meier-Ewert, J. Lange, H. Gerst, R. Herwig, A. Schmitt, and J. Freund *et al.* Comparative gene expression profiling by oligonucleotide fingerprinting. *Nucleic Acids Research*, 26(9):2216–2223, 1998. 2
- [13] S. Meier-Ewert, E. Maier, A. R. Ahmadi, J. Curtis, and H. Lehrach. An automated approach to generating expressed sequence catalogues. *Nature*, 361:375–376, 1993. 2
- [14] Shree K. Nayar and Tomaso Poggio, editors. *Early Visual Learning*. Oxford University Press, 1996. 11
- [15] William H. Press, Saul A. Teukolsky, William T. Vetterling, and Brian P. Flannery. *Numerical Recipes in C*. Cambridge University Press, 1992. 28
- [16] A. Watt. *3D Computer Graphics*. Addison-Wesley, 1994. 28
- [17] Gerhard Winkler. *Image Analysis, Random Fields and Dynamic Monte Carlo Methods*. Springer Verlag, 1995. 5