



## ‘Extremotaxis’: Computing with a bacterial-inspired algorithm

Dan V. Nicolau Jr.<sup>a,\*</sup>, Kevin Burrage<sup>b,d</sup>, Dan V. Nicolau<sup>c</sup>, Philip K. Maini<sup>a,d</sup>

<sup>a</sup> Centre for Mathematical Biology, Mathematical Institute, University of Oxford, Oxford OX1 3LB, United Kingdom

<sup>b</sup> Institute for Molecular Biosciences, University of Queensland, St. Lucia 4072, Australia

<sup>c</sup> Department of Electrical and Electronic Engineering, Liverpool University, Brownlow Hill, L693GJ Liverpool, United Kingdom

<sup>d</sup> Oxford Centre for Integrative Systems Biology, University of Oxford, Oxford OX13LB, United Kingdom

### ARTICLE INFO

#### Article history:

Received 31 May 2007

Received in revised form 29 October 2007

Accepted 23 May 2008

### ABSTRACT

We present a general-purpose optimization algorithm inspired by “run-and-tumble”, the biased random walk chemotactic swimming strategy used by the bacterium *Escherichia coli* to locate regions of high nutrient concentration. The method uses particles (corresponding to bacteria) that swim through the variable space (corresponding to the attractant concentration profile). By constantly performing temporal comparisons, the particles drift towards the minimum or maximum of the function of interest. We illustrate the use of our method with four examples. We also present a discrete version of the algorithm. The new algorithm is expected to be useful in combinatorial optimization problems involving many variables, where the functional landscape is apparently stochastic and has local minima, but preserves some derivative structure at intermediate scales.

© 2008 Published by Elsevier Ireland Ltd.

### 1. Introduction

The correspondence between living systems and computers has been stressed in recent years (e.g. Bray, 1995; Nicolau and Nicolau, 2006). Many biological processes can be thought of as processes of constrained optimization. Therefore, the mechanism or mechanisms used by a biological system to carry out a function is analogous to an algorithm or set of algorithms; the biological system is then an unconventional computer; and an instance of such a process taking place is analogous to a computational run. Of course, there are enormous differences between ‘biological computing’ and ‘classical computing’. Biological computations are massively parallel, feature a large degree of stochasticity and (intrinsic and extrinsic) noise – which, aside from being unavoidable, also plays a direct role in the computation – and, rather than being able to compute individual functions with high precision, deal instead with problems of a ‘systems engineering’ flavour, such as the control of very large systems, in the presence of non-linear constraints.

Because living systems are adapted to the environments in which they exist and therefore to the computational tasks required for survival, these natural computing paradigms are expected to be successful for dealing with problems similar to those confronting biosystems (Nicolau and Nicolau, 2006). An increasing number of algorithms are based on or inspired by biological strategies. These

include neural networks (Basheer and Hajmeer, 2000), evolutionary computing (Eiben and Smith, 2003), DNA computing (Adleman, 1994), particle swarm optimization (Call et al., 2007), computing with bio-agents (Nicolau et al., submitted for publication), ant optimization algorithms (Dorigo and Blum, 2005) and others. Increasingly these methods have been successfully applied to a spectrum of problems ranging from pattern identification and matching to aerodynamics engineering problems (Obayashi, 1997).

Chemotaxis, the process by which organisms direct their movements according to certain chemicals in their environment, is crucial for many biological functions. Bacteria such as *Escherichia coli* use chemotaxis to find food (for example, glucose) by swimming towards the highest concentration of food molecules, or to flee from poisons (for example, phenol). In multicellular organisms, chemotaxis is critical to development as well as normal function (Wadhams and Armitage, 2004). By analogy with the process of finding the maximum of a function (represented by the attractant concentration profile in space), chemotaxis is a algorithm for optimization. This computational facet of chemotaxis has already been noted by several authors (Bremermann, 1974; Muller et al., 2002). Recently, Vergassola et al. (2007) proposed a chemotaxis-inspired search method in the absence of gradients.

In this paper, we present a biocomputation approach that is based on the “run-and-tumble” chemotactic mechanism of the bacterium *E. coli*. This method is essentially a general-purpose search algorithm that can be used to optimize a function or set of functions. The method bears some resemblance to particle swarm optimization (PSO) in that the potential solutions (the particles) move

\* Corresponding author.

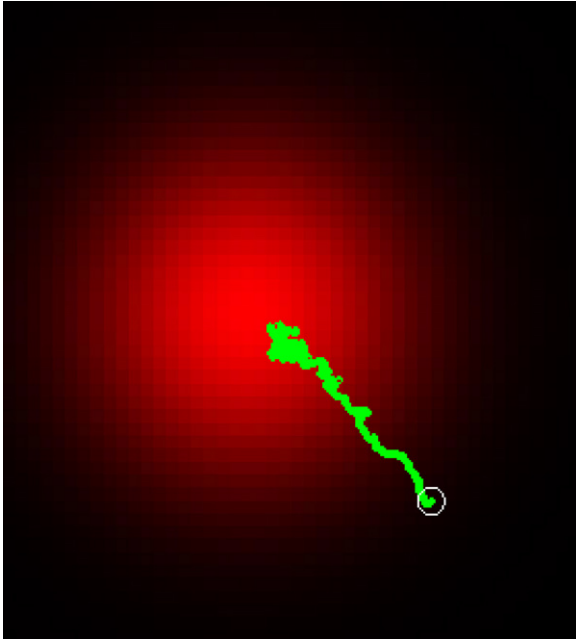
E-mail address: [nicolau@maths.ox.ac.uk](mailto:nicolau@maths.ox.ac.uk) (D.V. Nicolau Jr.).

through the function space. However, unlike PSO or, for example, ant colony optimization, it does not use inter-particle communication and does not bias the trajectories based on the best solutions found over time, relying instead completely on the chemotactic drift property of bacteria to converge locally (not as a swarm) to solutions. We illustrate the potential of the method by applying it to four representative optimization problems of different types. We also present a discrete version of the algorithm.

## 2. Methods

We begin by briefly describing the chemotactic swimming pattern of *E. coli*, on which our algorithm is based. *E. coli* is a common intestinal bacterium, cylindrical in shape and roughly  $2 \mu\text{m}$  long and  $1 \mu\text{m}$  wide. Each cell is equipped with approximately six flagella, each with a rotary motor at its base, embedded in the cell wall. The flagella are randomly distributed on the cell membrane. The rotary motor can turn clockwise and counter-clockwise at different times and is reversible. When all the motors turn in concert in a counter-clockwise direction, the flagella form a bundle that propels the cell forward in a “run”. Runs are not perfectly rectilinear due to rotational Brownian motion that perturbs the cell direction by roughly  $0.5\sqrt{t}$ , where  $t$  is in seconds. If one or more of the motors reverse direction and turn clockwise, the bundle becomes unstable and the cell turns in place (“tumbles”) in a random fashion and with negligible displacement. This serves to reorient the cell; the orientation is not perfect and there is some persistence of direction after a tumble (the mean angle between the direction before and after a tumble is  $63^\circ$ ) (Locsei, 2007). We omit this property in the present work, assuming that the reorientation is perfect.

*E. coli* cells use the system of motors and flagella to execute chemotactic swimming towards regions of high nutrient concentration (or away from toxins) as follows. Due to the high stochasticity of the environment and its small size, the bacterium cannot accurately measure an attractant gradient across its body. In lieu of computing a spatial gradient directly, a simple biochemical memory mechanism is used to perform temporal comparisons. During swimming, the bacterium monitors the concentration of chemoattractant (e.g. serine or aspartate) in the environment, comparing the average concentration measured over the last second with that measured over the previous 3 s. If the comparison indicates that the attractant concentration has increased, the cell is more likely to continue a straight-line run, while if it indicates that conditions have deteriorated, it is more likely to reorient by performing a tumble. In this way, the bacterium performs a biased random walk, leading it (in a stochastic fashion) up a chemoattractant gradient. In the absence of any such gradient, both the run and tumble times are exponentially distributed with means of 1.0 and 0.1 s, respectively (Locsei, 2007).



**Fig. 1.** Finding the global maximum in a Gaussian gradient field. Five hundred consecutive iterates of the colony centre (starting at the lower right) are shown in green; higher values of the attractant are shown as shades of red. The circle indicates the point at which the 250 particles are initialized.

We propose to use an analogous strategy to locate regions in a multi-dimensional space where a continuous (respectively discrete) function takes a global maximum (or minimum) value. We define such a “bacterial optimizer”  $B$  as a set of  $n$  particles ( $b_1, b_2, \dots, b_n$ ), each possessing an  $m$ -dimensional position vector function  $\mathbf{p}_i(t)$  and a velocity vector function  $\mathbf{v}_i(t)$  such that  $\mathbf{p}_i \in \mathbb{R}^m$ ,  $\mathbf{v}_i \in \mathbb{R}^m$ ,  $i = 1, \dots, n$  and  $t \in \mathbb{N}$ . This is the continuous version of the algorithm (a discrete version is described below). Let  $f: \mathbb{R}^m \rightarrow \mathbb{R}$  be the objective function and let  $\hat{\mathbf{x}}_i = \max_{t \geq 0} \mathbf{x}_i(t)$  and  $\hat{\mathbf{g}} = \max_{\mathbf{x}} f(\mathbf{x}_i)$ ,  $i = 1, \dots, n$ . Let  $U[x, y]$  be a random number between  $x$  and  $y$ , independently sampled from the uniform distribution and let  $N(\mu, \sigma)$  be a random number independently sampled from the normal distribution with mean  $\mu$  and standard deviation  $\sigma$ . The algorithm proceeds as follows (the meanings of the functions  $T$  and  $A$  and the various parameters are described after the algorithm):

- (i) Initialise  $\mathbf{p}_i$  and  $\mathbf{v}_i$  for all  $i$ . A simple choice is  $\mathbf{p}_i = U(a_j, b_j)$ ,  $j = 1, \dots, m$  and  $\mathbf{v}_i = \mathbf{r}_1$  where  $\mathbf{r}_1$  is a  $m \times 1$  vector each of whose entries  $\mathbf{r}_{1,j} = U[-1, 1]$ .  $a_j$  and  $b_j$  are the limits of the search domain in each dimension.
- (ii)  $\hat{\mathbf{g}} = \min_x f(\mathbf{x}_i)$ ,  $i = 1, \dots, n$
- (iii) While not converged:

For  $1 \leq i \leq n$ :

$\mathbf{x}_i \leftarrow \mathbf{x}_i + \varpi \mathbf{v}_i + \beta \mathbf{r}_r$ , where  $\mathbf{r}_r$  is a  $m \times 1$  vector, each of whose entries  $\mathbf{r}_{r,j} = N(0, 1)$

If  $f(\mathbf{x}_i) \geq f(\hat{\mathbf{g}})$ :  
 $\hat{\mathbf{g}} \leftarrow \mathbf{x}_i$

$Pr(tumble) = T(A_i(t))$

If  $r \leq Pr(tumble)$ , where  $r = U[0, 1]$ :

$\mathbf{v}_i = \mathbf{r}_{tumble}$ , where  $\mathbf{r}_{tumble}$  is a  $m \times 1$  vector each of whose entries  $\mathbf{r}_{tumble,j} \in U[-1, 1]$

End If

End For

$t \leftarrow t + 1$

End While

Convergence can be decided either by setting an upper limit on the number of iterations  $t_{\max}$  or by setting an acceptable value for  $f(\hat{\mathbf{g}})$ .

The tumbling probability function  $T$  is calculated as follows:

$$T(t, i) = \begin{cases} p_w, & A_i(t) < 0 \\ p_b, & A_i(t) \geq 0 \end{cases} \quad (1)$$

where

$$A_i(t) = \sum_{\tau=0}^{\min(w_r + w_d, t)} f(\mathbf{x}_i(t - \tau)) M(\tau) \quad (2)$$

and  $M: \mathbb{R} \rightarrow \mathbb{R}$  is a memory comparison function, which can take any number of forms but which we define for simplicity here as

$$M(\tau) = \begin{cases} \frac{1}{w_r}, & 0 < \tau \leq w_r \\ -\frac{1}{w_d}, & w_r < \tau \leq w_d \\ 0, & \text{otherwise} \end{cases} \quad (3)$$

The meaning of each of the parameters is as follows.  $\varpi$  is a speed factor for the particles in the functional search space, since the elements of  $\mathbf{v}_i$  are bounded by  $-1$  and  $1$ .  $\beta$  is a strictly positive parameter that if greater than  $0$  ensures that the runs are not perfectly straight and simulates rotational Brownian motion during a run.  $p_w$  and  $p_b$  are probabilities of tumbling if conditions have improved ( $A(t) \geq 0$ ) and deteriorated ( $A(t) < 0$ ), respectively. In practice, the probability of tumbling must be larger if conditions have deteriorated than if they have improved, so we have  $p_w > p_b$ .  $w_r$  and  $w_d$  are the number of iterations (window lengths) over which the recent and distant past are averaged, respectively. A balance must be struck between accuracy (using longer window lengths) and fast response time to improving or deteriorating conditions (leading to the use of shorter window lengths). Because

*E. coli* compares (roughly speaking) the last second of its life with the previous 3 s during chemotactic swimming (Strong et al., 1998), this would suggest a simple rule of  $w_d \approx 3w_r$ . In applying our algorithm to different functions, this is likely to vary with the nature and properties of the function in question.

An issue particular to the use of iterative algorithms (a large class of which the new algorithm is a member) is *stalling*, a phenomenon characterised by many iterations without any improvement (Li et al., 2006). Occasionally, this is caused by having found the global optimum but usually it is caused by the algorithm either having ‘passed by’ the optimum or having become stuck in one or more local minima. One way to deal with stalling is to perturb the algorithm in some way in order to escape the local minimum or to explore new regions. A potential implementation of this strategy for ‘extremotaxis’ is to force all the particles (bacteria) to tumble if no global improvement has been seen for some time. This can be done by modifying the probability of tumbling, calculated in the main loop of the algorithm, as follows:

$$Pr(\text{tumble}) = \begin{cases} T(A_i(t)), & \varepsilon < \varepsilon_{\text{critical}} \\ 1, & \varepsilon \geq \varepsilon_{\text{critical}} \end{cases} \quad (4)$$

where  $\varepsilon$  is the number of iterations since the last improvement in the global optimum found (this must be recorded over the course of the computation) and  $\varepsilon_{\text{critical}}$  is the critical number of iterations required to trigger a forced ‘colony’ tumble.

It is also possible to modify our algorithm so that it can be applied to discrete problems. The key change is to restrict the elements of the velocity vector  $\mathbf{v}_i$  to positive values smaller than 1 and to treat these as probabilities of the entries in  $\mathbf{x}_i$  changing state. For example, for a problem in which the variables can only take binary values (0 or 1), an element of  $\mathbf{v}_i$  equal to 0.1 means a 10% probability that the corresponding element of  $\mathbf{x}_i$  will change state at the next iteration of the algorithm. Additionally, the speed  $w$  should be set to 1 in order for the probabilities to be guaranteed to be between 0 and 1. In problems where the variables can take a number of discrete values (for example where they can take any positive integer value), each element of  $\mathbf{x}_i$  could, of course, be incremented by 1 or decremented by 1; therefore, the ‘direction’ of the increment should be chosen at random. Finally, it may be desirable to use increments greater than 1 (this would correspond to using a greater speed, in the continuous version of the algorithm). If this is done, then in order to avoid equal-sized increments at each point where a variable changes (and thus miss intermediate values), the size of the jump should be sampled from a suitable probability distribution with a mean of  $\xi$ , where  $\xi$  is the average increment size.

One issue if using this discrete version of the algorithm is that an appropriate fitness function may need to be more carefully chosen (or would be more difficult to find) than in the continuous version. In the latter, the fitness is evaluated simply as the value of the function at the point in  $m$ -dimensional space represented by the position vector. However, in the case of a discrete function, if the number of values that the function can assume is small or if these values are not consecutive (or both), this may not be appropriate. This is because the algorithm relies on ‘tumbles’ being more likely when the fitness is relatively inferior and less likely when it is close to the desired value. Therefore, using the value of the function as the fitness function may result in completely stochastic behaviour. We illustrate the issue using a simple (NP-complete) problem: Boolean satisfiability. Here it is required to determine, for a Boolean expression in  $n$  variables, what set of values for the variables (if any) makes the expression TRUE. Using a fitness function that simply takes the values 1 (for true) and 0 (for false) would not be advantageous in this case, reducing effectively to a random search through the function space (which would require exponential time proportional to  $2^n$ ). A more appropriate choice of fitness function may be

$$f(\mathbf{x}) = \max_{\text{all } z} c(\mathbf{z}) - c(\mathbf{x}) \quad (5)$$

where  $c(\mathbf{x})$  is the number of clauses in the Boolean expression that evaluate to 0. In this way, the algorithm would favour position vectors  $\mathbf{x}$  that result in a smaller number of such clauses evaluating to 0, thus in some sense being closer to finding a combination of variables that will cause the Boolean function to evaluate to TRUE. Of course, other possible fitness functions exist and in general it is to be expected that the choice of fitness function would vary with the problem and the determination of a suitable such function would present difficulties for some discrete problems.

### 3. Results

We implemented our algorithm using MATLAB and applied it to three different optimization problems. The first of these is trivial: finding the maximum of a two-dimensional Gaussian function. The second is finding the global minimum of a difficult two-dimensional function with many local minima. The third is concerned with how  $n$  particles should be distributed on a sphere so as to minimize the potential energy of the system.

#### 3.1. Finding the maximum of a Gaussian function

In order to demonstrate the operation of our algorithm, we first applied it to the Gaussian function:

$$f(x, y) = \frac{1}{2\pi\sigma^2} e^{-((x-x_0)^2 + (y-y_0)^2)/(2\sigma^2)}. \quad (6)$$

The Gaussian is an attractive first choice for a two reasons. Firstly, because the fundamental solution of the diffusion equation is a Gaussian, we might expect naturally occurring attractant gradients to take this form and therefore, due to adaptation, we might expect bacterial chemotaxis to be efficient at finding the global maximum of such a function (corresponding, in vivo, to, for example, finding the point of maximum nutrient concentration in a local environment). Secondly, it possesses a continuous and smooth gradient that is everywhere non-zero. It is, nonetheless, a non-trivial function.

Fig. 1 shows a typical simulation of the algorithm. A colony of 100 particles is initially distributed at random points chosen from  $x_{\text{init}} \in [-3, 3]$  and  $y_{\text{init}} \in [-3, 3]$ . We also randomly choose  $x_0 \in [-3, 3]$  and  $y_0 \in [-3, 3]$ . The initial velocity vectors are also chosen at random such that directions are uniformly distributed in  $[-\pi, \pi]$  and the magnitude of each direction vector is 0.02. We also set  $\sigma = 1$  for simplicity. The figure shows the first 500 consecutive iterations, with the maximum value found converging to the true maximum at  $(x_0, y_0)$ .

#### 3.2. Optimizing a difficult two-dimensional function

We next applied our algorithm to Problem 4 of the 100-digit challenge (Trefethen, 2002; Strang, 2005). This problem asks for the minimum of the function

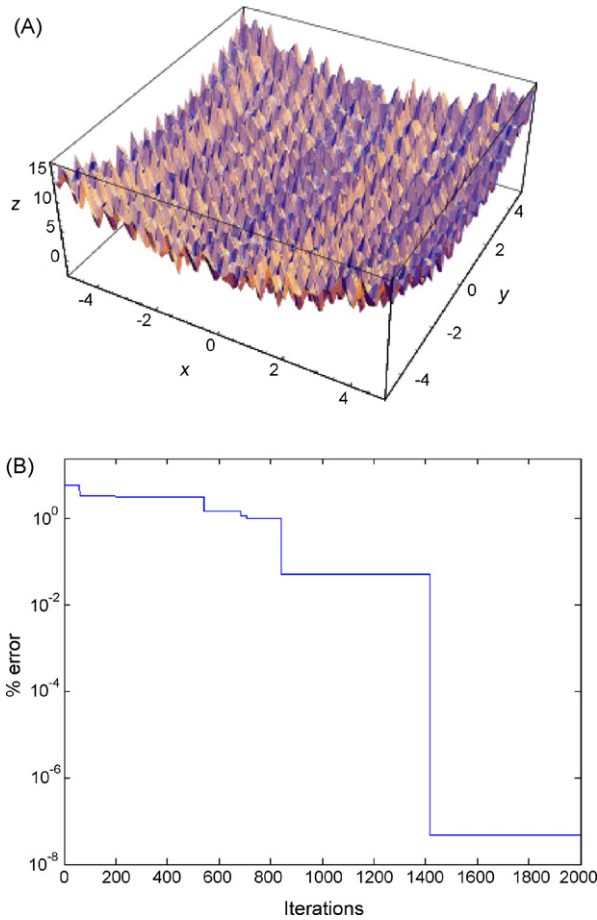
$$f = e^{\sin 50x} + \sin(60e^y) + \sin(70 \sin x) + \sin(\sin(80y)) - \sin(10(x+y)) + \frac{1}{4}(x^2 + y^2). \quad (7)$$

It is made difficult by the presence of many local minima that are very close to the global minimum—the latter is approximately  $f_{\text{min}} \approx -3.30686864747523728$  and occurs at  $(x, y) \approx (-0.0244030796943785, 2.10612427155358)$ . Fig. 2 A illustrates the difficulty, with a graph of the function showing the behaviour near the global minimum.

With the bacterial algorithm, setting up the problem consists of placing a number (250 in our computations) of particles at random in the two-dimensional function space near the minimum ( $-5 \leq x \leq 5$ ,  $-5 \leq y \leq 5$  are appropriate intervals) and randomising their directions. Again, the behaviour of the algorithm is good—Fig. 2 B shows the percentage difference of the algorithm’s best estimate from  $f_{\text{min}}$  over the course of a computational run. The 2000 iterations require, for 250 particles, only 2–3 s of computer time (on a 1-GHz desktop machine running MATLAB) to find the minimum with 8-digit accuracy (9-digit accuracy requires tens of seconds on the same machine). In the simple implementation presented here, the memory function and other parameters such as the particle velocity, directional persistence, etc. have not been optimized; this time would be reduced by some (unknown) factor if these steps were taken.

#### 3.3. Finding the minimum-energy configuration of particles on a sphere

Lastly, we applied our algorithm to a difficult  $n$ -body configuration problem: how to distribute  $n$  particles on a sphere so as to



**Fig. 2.** (A) A graph of the function in Eq. (7), showing the complex behaviour near the global minimum. (B) Finding the minimum of this function using taxis. The vertical axis shows the percentage difference between the best estimate,  $f(\hat{g})$  and the (known) global minimum.

minimise the potential energy

$$E = \sum_{i \neq j} \frac{1}{d(i, j)} \quad (8)$$

where  $d(i, j)$  is the (great-circle) distance on the sphere between particles  $i$  and  $j$ . This problem is computationally difficult because, similarly to the  $n$ -body problem and to protein folding, the potential energy space to be searched grows very rapidly with the number of particles—at each iteration of a search algorithm, all pairwise distances must be re-evaluated. Additionally, because a small difference in even a single distance can make a large difference to the sum, a brute force search will fail due to the fine required partition of the search space. Approximations to optimal configurations for this problem are known (Hardin and Sloane, 1995) for various numbers of particles.

To apply the taxis algorithm to the problem, first we convert the coordinates of the particles to spherical coordinates (latitude and longitude). For two particles  $i$  and  $j$ , let  $\phi_i$  and  $\phi_j$  be the latitudes and  $\lambda_i$  and  $\lambda_j$  the longitudes. Then the great-circle distance on a sphere of unit radius is

$$d(i, j) = \arccos[\sin \phi_i \sin \phi_j + \cos \phi_i \cos \phi_j \cos(\lambda_i - \lambda_j)]. \quad (9)$$

Each bacterium in the computation represents one possible solution, i.e. one arrangement of particles. The  $n$ -dimensional location vector of each computational agent  $i$  is then  $\mathbf{p}_i =$

**Table 1**

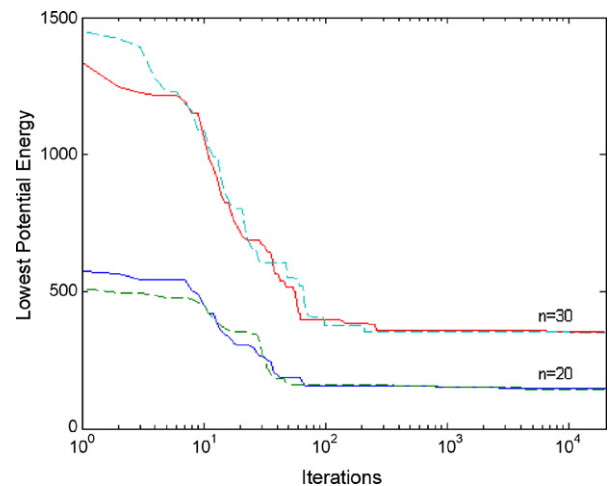
Lowest energy particle arrangements found with 'taxis' compared with previously published values

Number of particles	Lowest energy	
	Hardin et al. (1996)	'Taxis' algorithm
5	6.4746915	6.3395412
10	32.7169495	30.2804821
15	80.6702441	77.9961307
20	150.8815683	142.582206
25	243.8127603	238.313036
30	359.6039459	351.460535

Computations were carried out on a 3.4-GHz IBM machine with 1 GB RAM, running MATLAB<sup>TM</sup> version 7.1. Each run consisted of 40,000 iterations, corresponding to 30 s total for the  $n = 30$  (slowest) case.

$\{(\phi_1, \lambda_1); (\phi_2, \lambda_2), \dots, (\phi_n, \lambda_n)\}$  and the velocity vector is  $\mathbf{v}_i = \{\theta_1, \theta_2, \dots, \theta_n\}$ , where  $\theta_j$  is the bearing of the  $j$ th particle, with  $\theta_j \in [-\pi, \pi]$ . The location and velocity vectors are initially chosen at random for each computational agent (we used 20 in our simulations). Note that with this definition, a tumble corresponds to all the particles in one potential solution reorienting.

At each step of the calculation, the optimum arrangement among the  $k$  agents (the one with the smallest potential energy) is recorded and represents the best arrangement found up to that point. Using larger values of  $k$  increases the probability of rapid convergence and decreases the probability of the entire system becoming stuck in local minima, but increases the running time in proportion to  $k$ . Table 1 presents the results of this algorithm (left column) compared with the values given by Hardin and Sloane (1995). Remarkably, 'taxis' seems to find more optimal arrangements than those previously known. Fig. 3 shows the convergence of the system to these values for different numbers of particles (in all results shown, the computations were stopped after 20,000 iterations). The algorithm converges quickly and reproducibly in all cases. Adding a tumbling perturbation to the system every 100–1000 iterations to combat stalling does not have a significant effect on the performance. However, since the values found are better than those known (Hardin and Sloane, 1995) and probably very close to optimal, this is not a major consideration here.



**Fig. 3.** The convergence of a system of agents to the lowest potential-energy arrangement of  $n$  particles on a sphere (dotted line show the results with additional perturbation after 1000 iterations of stalling).



### 3.4. Detection of microarray features

One of the major difficulties of microarray technology relates to the processing of large and, importantly, error-loaded images of the dots on the chip surface (Qin et al., 2005). Whatever the source of these errors, those obtained in the first stage of data acquisition (segmentation) are passed down to the subsequent processes, with deleterious results.

The interpretation of the microarray data starts with the integration of the signal compared with the background on the area of individual dots – segmentation – the results being further used for elaborate clustering methods (Qin et al., 2005). It follows that the incorrect demarcation of the circular dot features will propagate throughout the whole microarray data processing and will add to the other sources of variability of microarray data: biological variability, technical variability and labeling (Zakharkin et al., 2005). Several methods have been proposed to de-noise data (Adjero et al., 2006). These include adaptive split and merge algorithm (Barra, 2006), polynomial-hyperbolic spot shape model in combination with the Box–Cox transformation (Ekström et al., 2004), spectral embedding (Higgs et al., 2006), noise-resistant algorithms (Novikov and Barillot, 2006), background extraction (O'Neill and Magoulas, 2003), just to name a few recent contributions. We tested the performance of ‘extremotaxis’ on both computer-generated (artificial) microarray-like spots and real microarray images. Thinking about the intensity of an image as different levels of nutrient concentration in space allows us to use this algorithm to identify regions of high intensity. The form of the ‘attractant’ (image

intensity profile) will affect the motion of the model bacteria. We initially suppose for simplicity that the intensity of a spot can be well described by a Gaussian function (see first example, above) and that the spots are arranged in a rectangular array of boxes, each spot of different randomly chosen size (different standard deviation of the Gaussian) and each placed at a different location within its home box (different spatial means of the Gaussian).

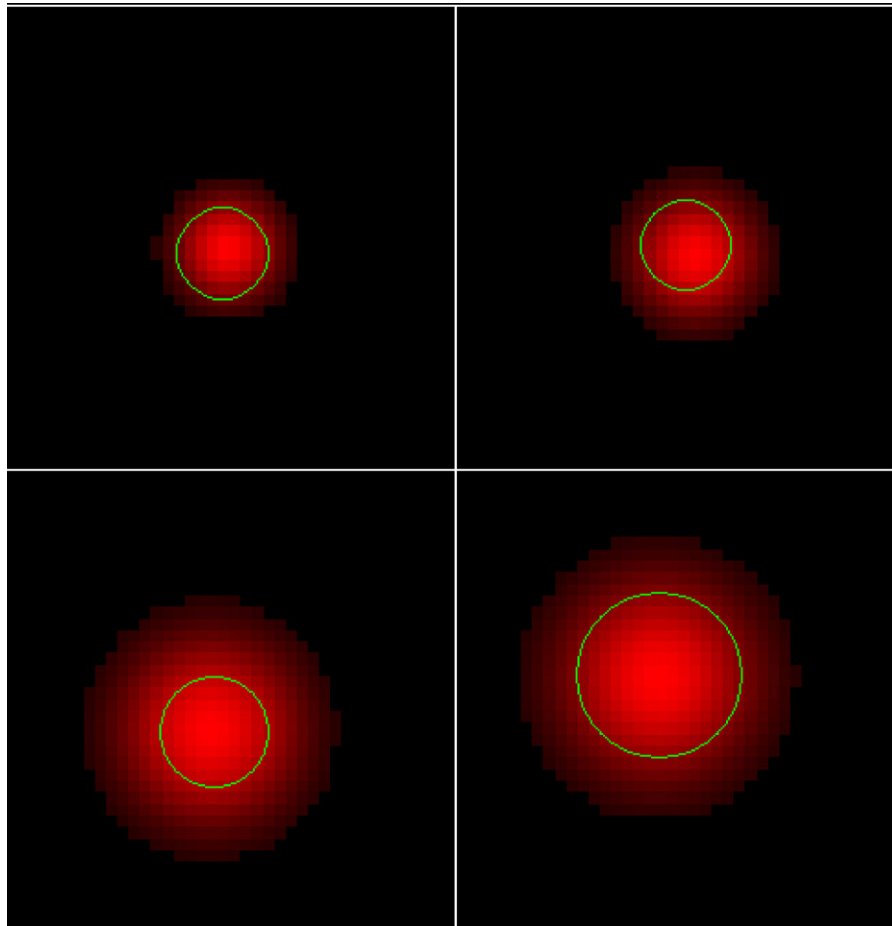
To locate the centre of a spot that is located randomly inside the simulation area, we can compute, at periodic intervals, the “centre of mass” of the model bacteria as follows:

$$C_x = \frac{\sum_{i=1}^n x_i}{n} \quad C_y = \frac{\sum_{i=1}^n y_i}{n} \quad (10)$$

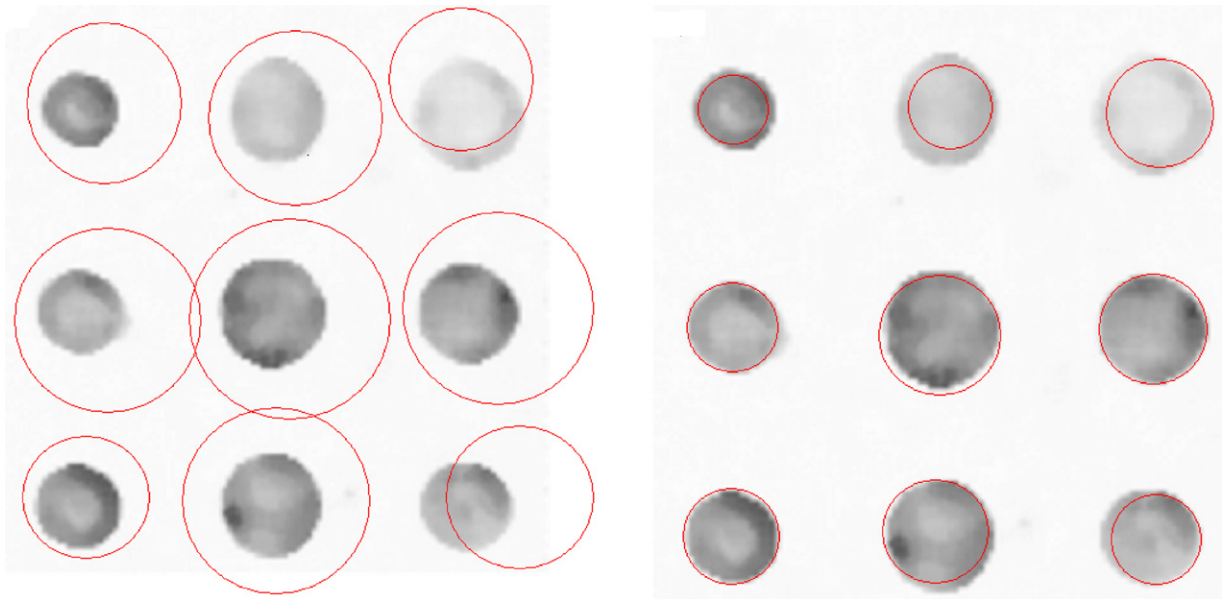
where  $x_i$  and  $y_i$  are the  $x$  and  $y$  coordinates of bacterium  $i$ . As bacteria aggregate around the intensity centre of the image (the centre of the spot), the point  $(c_x, c_y)$  approaches the true centre of the intensity spot. We can also obtain a measure of the dimensions of the spot by computing the standard deviation of the bacterial positions:

$$\sigma_{\text{spot}} = \sqrt{\sum_{i=1}^n (x_i - c_x)^2 + \sum_{i=1}^n (y_i - c_y)^2 / n} \quad (11)$$

Gaussian intensity distributions lead to good performance of the algorithm, possibly because in the natural habitat of bacteria,



**Fig. 4.** Spot regions and centres on a  $2 \times 2$  model microarray located using the bacterial algorithm. Each spot is a truncated Gaussian with  $\theta = 0.15$  (see text).



**Fig. 5.** Performance of the algorithm on a real  $3 \times 3$  microarray image after 100 (left) and 1000 iterations (right) (determined spot areas shown by solid circles).

the nutrient is dispersed through diffusive processes, which are typically characterised by Gaussian or Gaussian-like distributions (since the fundamental solution of Fick's equation takes this form). However, this form may not be a realistic model for a microarray spot since, for example, we expect the edges of the spot to be much sharper than the long tail of a Gaussian distribution. We can simulate this by truncating the Gaussians using a threshold of intensity  $\theta$ , so that for all values of intensity below  $\theta$  we simply set the intensity to 0. Eq. 10 can then be modified so that only the model bacteria located on voxels whose intensity is non-zero will be taken into account. Eq. 11 can be similarly modified. Typical results on a  $2 \times 2$  array are shown in Fig. 4.

Finally, we measured the performance of our algorithm on a real microarray image chosen from (Rhodes, 2005). Fig. 5 shows the determined spot sizes and centres.

Importantly, the algorithm presented here has potential even if the image is not segmented, as we have assumed it to be. In this more general case, the model bacteria would still converge on the spots present in the image, but it would not be possible to identify the number of spots or set bounds on their locations *a priori*. To solve this problem, one would need to detect when a cluster of bacteria has formed, this in turn requiring an algorithm dedicated to this task. One possibility is to mimic "quorum sensing", a phenomenon whereby bacteria not only consume but also release a chemoattractant into the environment, to which they are in turn attracted. Other bacteria are then attracted to this region and in such a way a stable cluster is formed. It would then be possible to identify a cluster formed by model bacteria when the mean "chemoattractant" distributed over a set of voxels exceeds a known threshold. This will form the subject of future work.

#### 4. Discussion

We have so far presented four examples of problems that can be tackled with our method. Clearly, the performance of the algorithm, as is the case with any optimization algorithm, would depend strongly on the nature of the problem under consideration. Why might we, in general, expect computing with taxis to perform well at optimizing certain difficult functions? We can speculate on an

answer to this question. Because bacteria are the oldest motile organisms and because the environments in which they live are complex at different scales of space and time, it might be expected that they be very efficient at solving optimization problems, including through chemotaxis. Recent work (Nicolau et al., submitted for publication) suggests that run-and-tumble is evolutionarily optimal and that, remarkably, this simple algorithm can for example (as a conservative estimate) locate on average more than 92% of the total available nutrient in a Gaussian field. Furthermore, other natural algorithms and biocomputation methods such as neural networks, evolutionary computing, DNA computing and (most similar to taxis computing), particle swarm optimization have been successful. Therefore, there are general reasons to be optimistic about the potential of taxis computing for global optimization.

The question can also be asked in the opposite direction: for what types of functions would the method be expected to perform well? Functions that are difficult to optimize because of the presence of many local extrema are good candidates because they resemble in some sense the natural environments of bacteria. If we think of the presence of many such extrema as an "apparent stochasticity" in the function (from the point of view of a particle walking the functional landscape) then we can draw an analogy between noise in biological environment and the presence of many local minima on this landscape. In other words, local fluctuations in the derivative of a function are analogous to noise in a natural environment, in this sense.

On the other hand, run-and-tumble relies on the presence of gradients to produce a drift towards favourable environments. If the functional landscape is either extremely stochastic or discontinuous, no gradient will be reliably detected and, in the limit, the method reduces to a diffusion-like random local search at a number of random points (equal to the number of bacteria in the system) on the landscape. This may not always be disadvantageous—for example, one can imagine funnel-like landscapes (similar to the postulated energy landscapes of folding protein) that possess smooth gradients on the whole but become very stochastic near the global minimum. In these cases, a combination of gradient-induced global drift and noise-induced random local search may perform well. Nonetheless, taking these ideas together, we expect

the type of function on which taxis computing will perform well relative to other methods to possess local gradients on scales larger than the characteristic velocity of the moving particles.

As mentioned, taxis computing bears some resemblance to particle swarm optimization. Both exhibit some attributes of evolutionary computing: each particle represents a potential solution, these solutions are initially randomly chosen, and the algorithm proceeds by evolving the solutions from iteration to iteration, with each iteration being based on the last. Of course, both methods are based on the concept of a set of particles moving through the problem space.

Two essential differences are that (a) in PSO the particles share information about the best solutions found up to each point in the computational run and (b) in PSO the velocity of each particle in the swarm is changing smoothly while in the model we propose here, the direction of each particle is constant during a run and is randomized by a tumble. The first of these is particularly essential because it means the swarm as a whole may become trapped in local minima. In PSO, at each step the velocity of particle  $i$  is re-evaluated according to the equation (Call et al., 2007):

$$\mathbf{V}_{id} = w\mathbf{V}_{id} + c_1r_1(\mathbf{b}_p - \mathbf{x}_{id}) + c_2r_2(\mathbf{b}_i - \mathbf{x}_{id}) \quad (12)$$

where  $\mathbf{V}_{id}$  is the velocity of the particle,  $\mathbf{x}_{id}$  is the position of the particle,  $\mathbf{b}_p$  is the position of the best solution seen by the swarm as a whole and  $\mathbf{b}_i$  is the position of the best solution found by the particle.  $r_1$  and  $r_2$  are random numbers and  $c_1$ ,  $c_2$  and  $w$  are positive real numbers representing the “weights” of the three terms. Because the particles (a) cooperate amongst themselves and (b) remember and factor in their best solution to date, a sufficiently good local minimum, once found, may trap the particle and in some cases the whole swarm. This cannot happen in taxis computing because these features are not present; instead, each particle relies on the structure of the local environment combined with random reorientations to explore the search space. Although it is possible (though not equally likely as in PSO) that an individual bacterium may become trapped in a local minimum for a time, this cannot happen at the level of the colony. Furthermore, because of the stochastic nature of run-and-tumble, its escape probability from this region will be non-zero and hence the residence time will be finite. Of course, the cooperation property of PSO is often valuable because the swarm as a whole can converge towards a favourable region of the problem space, which can then be searched more efficiently; nonetheless, the reinforcement of solutions already found at both swarm level and individual level means the algorithm has a higher probability of missing the global minimum of functions with properties similar to those described above.

Investigations of the performance of the discrete version of the algorithm will form the subject of future work. One can speculate, however, on the prospects of this method. On the one hand, because chemotaxis relies on the presence of gradients—in the context of computing, a direct and well-behaved relationship between position in  $n$ -space and fitness, we do not expect the method to perform as well or as consistently for discrete functions, for which there is little or no such correlation. For example, the difficulty in solving Boolean satisfiability stems from the property that a change in the state of one single variable (possibly among hundreds or thousands of such variables) will be the difference between the expression evaluating as TRUE or FALSE. On the other hand, a discrete version of PSO (Yang et al., 2004) has been successfully used to solve various discrete problems such as the capacitated vehicle routing problem (CVRP) (Ai-ling et al., 2006). Although in the worst case, taxis computing for discrete problems may reduce to a stochastic random search through the functional landscape (if tumble probability is uncorrelated with changes in the fitness function, or if these changes are very rare), in many cases the method may work well.

This is expected to be the case, for example, when the position vector (i.e. the independent variables) is simply restricted to integer entries, as might happen for an integer optimization problem.

In the numerical results presented here, we have used a memory function of the form in Eq. (3) with  $w_d \approx 3w_r$ , because *E. coli* compare roughly the last second of their lives with the previous 3 s (Strong et al., 1998), and for simplicity. However, the form of the memory function used by a live bacterium is believed to be more complicated (Segall et al., 1986). Even more importantly, the optimal form of the memory function in the context of biocomputation is likely to be (a) different and (b) sensitive to problem or class of problems under consideration. Therefore, future work will explore different memory functions and their performances for different problems. A promising avenue is to “evolve” the memory function *in silico* for a particular class of problems. For example, in recent work we evolved the memory function of a chemotactic bacterium-like organism on a computer, in the presence of a Gaussian attractant distribution, finding that the evolved function resembles the biphasic shape believed to be at work in the chemotactic mechanism of *E. coli*. Presumably, when exposed to different functional landscapes, a “species” of digital organisms equipped with the ability to evolve the memory function will adapt to the function in question, developing an optimal or near-optimal response.

It was mentioned above that in PSO, the swarm converges to the best solution found to date and that this strategy, while running the risk of missing the global optimum, means that local searches near the best solution found are more efficient—because they are carried out by more particles. In an attempt to introduce this feature into our model by mimicking the natural behaviour of bacteria, one possible variation on the algorithm presented above would also allow the agents to divide (produce offspring) when in a favourable region of the functional landscape. This would maintain the advantage of not swarming to a local minimum while increasing the efficiency of local searches (and, if the agents can also die, reducing the proportion of computational time dedicated to unpromising regions). Finally, future work will also focus on comparing this method with other methods, both of a natural computing flavour and also more classical methods such as steepest descent, random search, etc.

## Acknowledgements

D.V.N. would like to acknowledge the financial support from the Clarendon Fund, the UK Overseas Research Student Award Scheme and the Devorguilla Scholarship from Balliol College, University of Oxford. K.B. gratefully acknowledges support via the Federation Fellowship of the Australian Research Council. P.K.M. was partially supported by a Royal Society-Wolfson Merit award. The authors are indebted to Tjeerd olde Scheper for help with the manuscript.

## References

- Adjeroh, D.A., Zhang, Y., Parthe, R., 2006. On denoising and compression of DNAmicroarray images. *Pattern Recogn.* 39, 2478–2493.
- Adleman, L.M., 1994. Molecular computation of solutions to combinatorial problems. *Science* 266 (11), 1021–1024.
- Ai-ling, C., Gen-ke, Y., Zhi-ming, W., 2006. Hybrid discrete particle swarm optimization algorithm for capacitated vehicle routing problem. *J. Zhejiang Univ. Sci. A* 7 (4), 607–614.
- Barra, V., 2006. Robust segmentation and analysis of DNA microarray spots using an adaptative split and merge algorithm. *Comput. Methods Prog. Biomed.* 8, 174–180.
- Basheer, I.A., Hajmeer, M., 2000. Artificial neural networks: fundamentals, computing, design, and application. *J. Microbiol. Methods* 43 (1), 3–31.
- Bray, D., 1995. Protein molecules as computational elements in living cells. *Nature* 376, 307–312.
- Bremermann, H.J., 1974. Chemotaxis and optimization. *J. Franklin Inst.* 297, 397–404.
- Call, S.T., Zubarev, D.Y., Boldyrev, A.I., 2007. Global minimum structure searches via particle swarm optimization. *J. Comput. Chem.* 28, 1177–1186.

- Dorigo, M., Blum, C., 2005. Ant colony optimization theory: a survey. *Theor. Comp. Sci.* 344 (2–3), 243–278.
- Eiben, A.E., Smith, J.E., 2003. *Introduction to Evolutionary Computing*. Springer, New York.
- Ekstrøm, C.T., Bak, S., Kristensen, C., Rudemo, M., 2004. Spot shape modelling and data transformations for microarrays. *Bioinformatics* 20, 2270–2278.
- Hardin, R.H., Sloane, N.J.A., 1995. *Codes (Spherical) and Designs (Experimental)*, Proceedings of Symposia in Applied Mathematics, Vol. 50.
- Higgs, B.W., Jennifer Weller, J., Solka, J.L., 2006. Spectral embedding finds meaningful (relevant) structure in image and microarray data. *BMC Bioinformatics* 7, 74–87.
- Li, W., Pan, P.Q., Chen, G.T., 2006. Combined projected gradient algorithm for linear programming. *Optim. Method Softw.* 21 (4), 541–550.
- Locsei, J.T., 2007. Persistence of direction increases the drift velocity of run and tumble chemotaxis. *J. Math. Biol.* 55(1), 41–60.
- Muller, S., Marchetto, J., Airaghi, S., Koumoutsakos, P., 2002. Optimization based on bacterial chemotaxis. *IEEE Trans. Evol. Comput.* 6 (1), 16–29.
- Nicolau, D.V., Armitage, J.P., Maini, P.K., submitted for publication. In silico evolution of chemotactic swimming.
- Nicolau, D.V., Nicolau, D.V., 2006. *Biocomputation*. Wiley, New York.
- Novikov, E., Barillot, E., 2006. A noise-resistant algorithm for grid finding in microarray image analysis. *Mach. Vision Appl.* 17, 337–345.
- Obayashi, S., 1997. *Pareto genetic algorithm for aerodynamic design using the Navier-Stokes equations*. Wiley, New York.
- O'Neill, P., Magoulas, G.D., 2003. Improved processing of microarray data using image reconstruction techniques. *IEEE Trans. Nanobiosci.* 2, 176–183.
- Qin, L., Rueda, L., Ali, A., Ngom, A., 2005. Spot detection and image segmentation in DNA microarray data. *Appl. Bioinformatics* 4 (1), 1–11.
- Rhodes, P., 2005. Enhancement of DNA and Microarray Analysis using Image Processing Practices (presentation).
- Segall, J.E., Block, S.M., Berg, H.C., 1986. Temporal comparisons in bacterial chemotaxis. *Proc. Natl. Acad. Sci. U. S. A.* 83 (23), 8987–8991.
- Strang, G., 2005. The SIAM 100-digit challenge—a study in high-accuracy. *Science* 307 (5709), 521–522.
- Strong, S.P., Freedman, B., Bialek, W., Koberle, R., 1998. Adaptation and optimal chemotactic strategy for *E. coli*. *Phys. Rev. E* 57, 4604–4617.
- Trefethen, N., 2002. A hundred-dollar, hundred-digit challenge. *SIAM News* 35 (1).
- Vergassola, M., Villermaux, E., Shraiman, B.I., 2007. Infotaxis as a strategy for searching without gradients. *Nature* 445, 406–409.
- Wadhams, G.H., Armitage, J.P., 2004. Making sense of it all: bacterial chemotaxis. *Nat. Rev. Mol. Cell. Biol.* 5, 1024–1037.
- Zakharkin, S.O., Kim, K., Mehta, T., Chen, L., Barnes, S., Scheirer, K.E., Parrish, R.S., Allison, D.B., Page, G.P., 2005. Sources of variation in affymetrix microarray experiments. *BMC Bioinformatics* 6, 214–225.