

Identifying Gene Knockout Strategy Using Bees Hill Flux Balance Analysis (BHFBA) for Improving the Production of Ethanol in *Bacillus Subtilis*

Yee Wen Choon¹, Mohd Saberi Mohamad¹, Safaai Deris¹, Rosli M. Illias²,
Lian En Chai¹, and Chuii Khim Chong¹

¹ Artificial Intelligence and Bioinformatics Research Group, Faculty of Computer Science and Information Systems, Universiti Teknologi Malaysia, Skudai, 81310 Johor, Malaysia

{ywchoon2, lechai2, ckchong2}@live.utm.my, {saberis, safaaai}@utm.my

² Department of Bioprocess Engineering, Faculty of Chemical Engineering,

Universiti Teknologi Malaysia, Skudai, 81310 Johor, Malaysia

r-rosli@utm.my

Abstract. *Bacillus subtilis* strains can be manipulated to improve product yield and growth characteristics. Optimization algorithms are developed to identify the effects of gene knockout on the results. However, this process is often faced the problem of being trapped in local minima and slow convergence due to repetitive iterations of algorithm. In this paper, we proposed Bees Hill Flux Balance Analysis (BHFBA) which is a hybrid of Bees Algorithm, Hill Climbing Algorithm and Flux Balance Analysis to solve the problems and improve the performance in predicting optimal sets of gene deletion for maximizing the growth rate and production yield of desired metabolite. *Bacillus subtilis* is the model organism in this paper. The list of knockout genes, growth rate and production yield after the deletion are the results from the experiments. BHFBA performed better in term of computational time, stability and production yield.

Keywords: Bees Algorithm, Hill Climbing, Flux Balance Analysis, *Bacillus subtilis*, Optimization.

1 Introduction

Microbial strains optimization has become popular in genome-scale metabolic networks reconstructions recently as microbial strains can be manipulated to improve product yield on desired metabolites and also improve growth characteristics [1]. Reconstructions of the metabolic networks are found to be very useful in health, environmental and energy issues [2]. The development of computational models for simulating the actual processes inside the cell is growing rapidly due to vast numbers of high-throughput experimental data.

Many algorithms were developed in order to identify the gene knockout strategies for obtaining improved phenotypes. The first rational modeling framework (named OptKnock) for introducing gene knockout leading to the overproduction of a desired

metabolite was developed by Burgard *et al.*, 2003 [3]. OptKnock identifies a set of gene (reaction) deletions to maximize the flux of a desired metabolite with the internal flux distribution is still operating such that growth is optimized.

OptKnock is implemented by using mixed integer linear programming (MILP) to formulate a bi-level linear optimization that is very promising to find the global optimal solution. OptGene is an extended approach of OptKnock which formulates the *in silico* design problem by using Genetic Algorithm (GA) [4]. Meta-heuristic methods are capable in producing near-optimal solutions with reasonable computation time, furthermore the objective function that can be optimized is flexible. SA is then implemented to allow the automatic finding of the best number of gene deletions for achieving a given productivity goal [5]. However, the results are not yet satisfactory.

A hybrid of BA and FBA was proposed by Choon *et al.*, 2012 [6], it showed a better performance in predicting optimal gene knockout strategies in term of growth rate and production yield. Pham *et al.*, 2006 [7] introduced Bees Algorithm (BA), is a typical meta-heuristic optimization approach which has been applied to various problems, such as controller formation [8], image analysis [9], and job multi-objective optimization [10]. BA is based on the intelligent behaviours of honeybees. It locates the most promising solutions, and selectively explores their neighbourhoods looking for the global maximum of the objective function. BA is efficient in solving optimization problems according to the previous studies [7, 10]. However, due to the dependency of BA on random search, it is relatively weak in local search activities [11]. Hence, BHFBA is proposed to improve the performance of BAFBA as Hill climbing algorithm is a promising algorithm in finding local optimum. This paper shows that BHFBA is not only capable in solving larger size problems in shorter computational time but also improves the performance in predicting optimal gene knockout strategy than previous works. In this work, we present the results obtained by BHFBA in two case studies where *B. subtilis* (*Bacillus subtilis*) iBsu1103 model is the target microorganisms [12]. In addition, we also conduct a benchmarking to test performance of the hybrid of Bee algorithm and Hill climbing algorithm.

2 Bees-Hill Flux Balance Analysis (BHFBA)

In this paper, we propose BHFBA in which BAFBA is only applied to identify optimal gene knockout strategies recently. Fig. 1 shows the flow of BAFBA while Fig. 2 shows our proposed BHFBA. The important steps are explained in the following subsections. Both BHFBA and BAFBA are using binary variables rather than continuous variables. The main difference between BHFBA and BAFBA is the neighbourhood search part, BHFBA improves the operation by combining hill climbing algorithm into BAFBA.

2.1 Model Pre-processing

The model is pre-processed through several steps based on biology assumptions as well as computational approaches to reduce the search space as while as increase the accuracy. Lethal reactions such as the genes that are found to be lethal *in vivo*, but not *in silico*, should be removed to improve the quality of the results. The results are

invalid if a lethal reaction is deleted. The following are the details of computational pre-processing steps to the model [5].

a. Fluxes that are not associated with any genes, such as the fluxes related to external metabolites and exchange fluxes that represent transport reaction should not be involved in the process. These fluxes do not have a biological meaning thus they should not be knocked out.

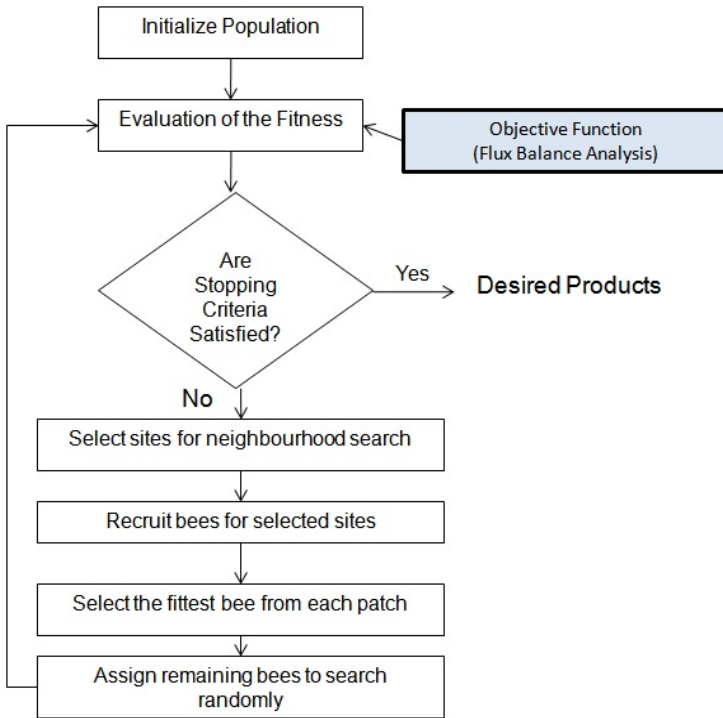
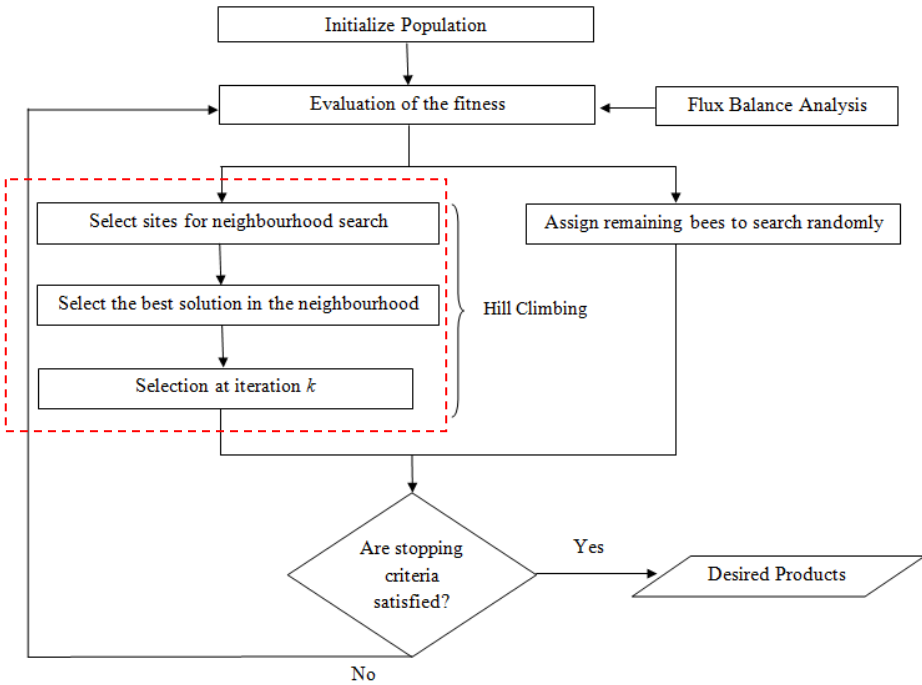


Fig. 1. BAFBA Flowchart

b. Essential genes that cannot be deleted from the microorganism's genome need to be removed. The search space for optimization is reduced due to that these genes should not be considered as targets for deletion. A linear programming problem is defined by setting the corresponding flux to 0, while maximizing the biomass flux for each gene in the microorganism's genome. If the biomass flux result from the Linear Programming algorithm is zero (or near zero) then the gene is marked as essential. This biological meaning of this fact is that the microorganism is unable to survive without this gene. This process does not suggest any changes to the model like the previous one, but provides favorable information for the optimization algorithms. With the help of biologists, the list of essential genes can be manually edited to include genes that are known to be essential *in vivo*, but not *in silico*.

c. Given the constraints of the linear programming problem, the fluxes need to be removed if the fluxes cannot exhibit values different from 0. Two linear programming are solved for every reaction in the model: the first is to define the flux over that

reaction as the maximization target, while the second is to set the same variable as minimization target. If the objective function is 0 for both problems, then the variable is removed from the model.



Note: Red-dotted box is Hill Climbing algorithm which is newly hybridized into BA.

Fig. 2. BHFBA Flowchart

2.2 Bee Representation of Metabolic Genotype

One or more genes can be discovered in each reaction in a metabolic model. In this paper, each of those genes is represented by a binary variable indicating its absence or presence (0 or 1), these variables form a ‘bee’ representing a specific mutant that lacks some metabolic reactions when compared with the wild type (Fig. 3.)

2.3 Initialization of the Population

The algorithm starts with an initial population of n scout bees. Each bee is initialized as follows: assume that a reaction with n genes. Bees in the population are initialized by setting present or absent status to each gene randomly. Initialization of the population is done randomly so that all bees in the population have an equal chance of being selected. The result might not truly reflect the population if it is done with bias setting.

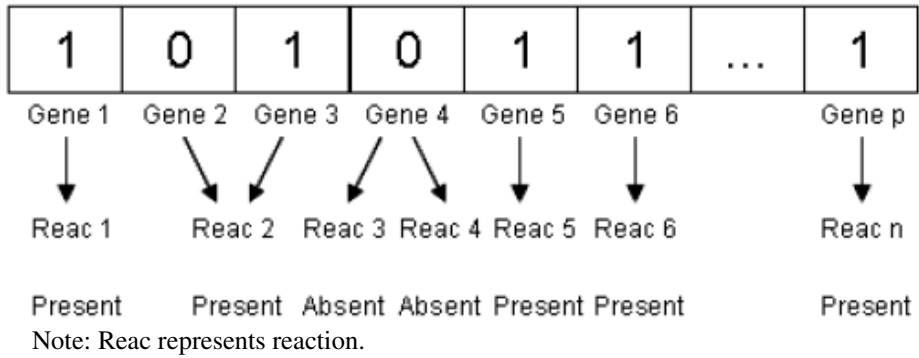


Fig. 3. Bee representation of metabolic genotype

2.4 Scoring Fitness of Individuals

Each site is given a fitness score that determines whether to recruit more bees or should be abandoned. In this work, we used FBA to calculate the fitness score for each site and the equation is as follow:

Maximize Z, where

$$Z = \sum_i c_i v_i = \mathbf{c} \cdot \mathbf{v} \tag{1}$$

where \mathbf{c} = a vector that defines the weights for of each flux.

Cellular growth is defined as the objective function Z, vector \mathbf{c} is used to select a linear combination of metabolic fluxes to include in the objective function, \mathbf{v} is the flux map and i is the index variable (1, 2, 3, ..., n). After optimizing the cellular growth, mutant with growth rate more than 0.1 continues the process by minimizing and maximizing the desired product flux at fixed optimal cellular growth value. Hence, we can enhance yield of our desired products at fixed optimal cellular growth. Production yield is the maximum amount of product that can be generated per unit of substrate. The following shows the calculation for production yield:

$$\text{Production yield} = \frac{(\text{production rate}_{\text{production}})}{(\text{consumption rate}_{\text{substrate}})} \tag{2}$$

(mmol/mmol)(gm/gm)

where mmol = millimole and gm is gram.

We used Biomass-product coupled yield (BPCY) as the fitness score in this work, the calculation for BPCY is as follow:

$$\text{BPCY} = \text{product yield} * \text{growth rate} \text{ (mmol(mmol*hr}^{-1}\text{))(gm (gm * hr}^{-1}\text{))} \tag{3}$$

where mmol is millimole, hr is hour and gm is gram.

2.5 Neighbourhood Search (Hill Climbing Algorithm)

The algorithm carries out neighbourhood searches in the favored sites (m) by using Hill climbing algorithm. Hill climbing is an iterative algorithm that starts with an arbitrary solution to a problem, then attempts to find a better solution by incrementally changing a single element of the solution. In this paper, the initial solution is the m favored sites from the population initialized by BA. The algorithm starts with the solution and makes small improvements to it by adding or reducing a bee to the sites. User defined the value of initial size of patches (ngh) and uses the value to update site (m) which is declared in the previous step to search in neighbourhood area. In this paper, m is equal to 15 and ngh is equal to 30, the values are obtained by conducting a small number of trials with the range of 10 to 25 and 20 to 35 respectively. This step is important as there might be better solutions than the original solution in the neighbourhood area.

2.6 Randomly Assigned and Termination

The remaining bees in the population are sent randomly around the search space to scout for new feasible solutions. This step is done randomly to avoid overlooking the potential results that are not in the range. These steps are repeated until either the maximum loop value is met or the fitness function has converged. At the end, the colony generates two parts to its new population – representatives from each selected patch and other scout bees assigned to perform random searches.

3 Results and Discussion

In this work, we use *E.Coli* and *B.subtilis* models to test on the operation of BAFBA. The *E.Coli* model is a small-scale model of the central metabolism of *E. coli* [12]. It is a modified subset of the iAF1260 model, and contains 134 genes, 95 reactions, and 72 metabolites. We use *E.coli* core model in this work because this model is useful for testing new constraint-based analysis methods, since the results of most constraint-based calculations are easier to interpret on this smaller scale. The second model is *B.subtilis* iBsu1103 model [13] which includes 1437 reactions associated with 1103 genes. We pre-process this model and the size is reduced to 763 reactions. The experiments are carried out by using a 2.3 GHz Intel Core i7 processor and 8 GB DDR3 RAM computer.

Table 1 and Table 2 summarize the results obtained from BHFBA for succinic acid production from *E.coli* and ethanol production from *B.subtilis*. Succinic acid is one of the intermediates of the TCA cycle and is a chemical to be used as a feedstock for the synthesis of a wide range of other chemical with several industrial applications. Besides, as a metabolite from the central carbon metabolism, succinic acid represents a good case study for identifying metabolic engineering strategies. Ethanol is a volatile, flammable, colourless liquid, and it is a promising biofuel. Ethanol is currently used as an alternative fuel for gasoline worldwide. As shown from the results, this method has produced better results to the previous works in term of growth rate and BPCY meanwhile potential genes which can be removed are identified [5][10].

Table 1. Comparison between different methods for production of Succinic acid in *E.coli*

Method	Growth Rate (1/hr)	BPCY	List of knockout genes
BHFBA	0.7988	0.93656	PTAr**, RPE, SUCD1i
BAFBA [10]	0.62404	0.66306	FUM, PTAr**, TPI**
SA + FBA [5]	N/A	0.39850	ACLD19*, DRPA, GLYCDx, F6PA, TPI**, LDH_D2, EDA, TKT2, LDH_D-
OptKnock [3]	0.28	N/A	ACKr, PTAr**, ACALD*

Note: The shaded column represents the best result. N/A – Not Applicable. * Common genes for all methods. ** Common genes in either 2 methods. BPCY is in gram (gram-glucose.hour)⁻¹.

Table 1 shows that BHFBA performed better than the previous works with growth rate 0.7988 and BPCY 0.93656. Knocking out succinate dehydrogenase (SUCD1i) interrupts the formation from succinic acid to fumarate. Without the conversion from succinic acid to fumarate, production yield of succinic is improved. Next, phosphotransacetylase (PTAr) is removed, according to Burgard *et al.*, 2003[3], the mutants can grow anaerobically on glucose by producing lactate. In the next step, ribulose-5-phosphate-3-epimerase (RPE) is suggested to knockout. This knockout involves the inflow reaction of ammonium. As stated in Bohl *et al.*, 2010 [14], the utilization of nitrate as electron acceptor and ammonium source under anaerobic conditions can improve succinate production. Figure 4 shows the comparison among the methods in term of growth rate and BPCY.

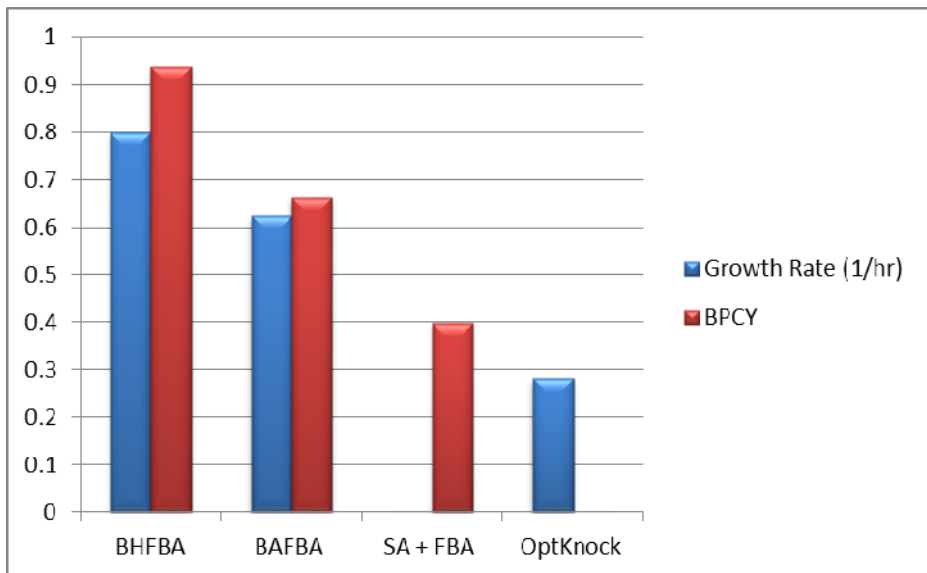

Fig. 4. Growth rate and BPCY comparison among available methods

Table 2. Comparison between different methods for production of ethanol in *B.subtilis*

Method	Growth Rate (1/hr)	BPCY	List of knockout genes
BHFBA	122.9089	1.15680e+05	ALAD_L, GPDH, LDH_L
BAFBA	122.8861	1.1154e+05	ALAD_L, LDH_L, XYLI1, inosose 2,3-dehydratase

Note: The shaded column represents the best result. N/A – Not Applicable. * Common genes for all methods. ** Common genes in either 2 methods. BPCY is in gram (gram-glucose.hour)⁻¹.

Table 3. Comparison between average computational time of BHFBA and BAFBA for 1000 iterations

Method	Computation Time (s)
BHFBA	7028
BAFBA	22515

Table 2 shows the results of BHFBA and previous works. BHFBA obtained a better growth rate and BPCY that are 122.9089 and 1.15680e+05 respectively. In the experiment of Kim *et al.*, 2012, deletion of NADH-dependent glycerol-3-phosphate dehydrogenase 1 (GPDH) showed a slight improvement in ethanol yield. As stated in Kim *et al.* (2012), lactate dehydrogenase (LDH_L) plays a key role in fermentative metabolism in metabolic engineering of *B.subtilis* for ethanol production. The deletion of LDH_L inhibits the conversion from pyruvate to lactate therefore more pyruvate is decarboxylated to acetaldehyde and further converted to ethanol.

Table 3 shows the computational time comparison between BHFBA and BAFBA for 1000 iterations. The average computational time of BHFBA improved 69% of the BAFBA result for 1000 iterations.

In addition, since BA and Hill Climbing algorithm is a new hybrid algorithm. Hence, we conducted a benchmarking to test performance of a hybrid of BA and Hill Climbing algorithm (BH). As BA is looking for the maximum, the functions are inverted before the algorithm is applied. The De Jong, Martin & Gaddy, and Griewangk functions are used in this paper. Table 4 shows the mathematical representation of the functions. Table 5 shows mean and standard deviation (STD) of the three functions, De Jong, Martin & Gaddy, and Griewangk, tested on both original BA and BH.

Table 4. Mathematical representation of De Jong and Beale functions

Name	Mathematical representation
De Jong	$\max F = (3905.93) - 100(x_1^2 - x_2)^2 - (1 - x_1)^2$
Martin & Gaddy	$\min F = (x_1 - x_2)^2 + ((x_1 + x_2 - 10) / 3)^2$
Griewangk	$\min F = 1 / (0.1 + (\sum(x(1,i)^2 / 4000)) - \sum(\cos(x(1,i) / \sqrt{i}) + 1))$

As seen from the results, both BHFBA and BH performed better than other algorithms. It can be concluded that the capability of Hill Climbing algorithm in finding local optimum improved the performance of the original BA. The original BA with the problem of repetitive iterations of the algorithm in local search where each bee keep searching until the best possible answer is reached. Our proposed BHFBA solved the problem by implementing Hill Climbing algorithm into the local search

part. Hill Climbing algorithm is a powerful local search algorithm which attempts to find a better solution by incrementally changing a single element of the solution until no further improvements can be found, the search process is recorded so the process is not repeated. Furthermore, one of the advantages of Hill Climbing algorithm is it can return a valid solution even if it is interrupted at any time before it ends.

Table 5. Obtained fitness value of both De Jong and Beale functions

Function	Mean		STD	
	BA	BH	BA	BH
De Jong	3.91e+03	3.90e+03	0.000504	4.79e-13
Martin & Gaddy	11.1083	11.1111	0.002797	0
Griewangk	-0.5263	-0.5263	5.76765E-09	0

4 Conclusion and Future Works

In this paper, BHFBA is proposed to predict optimal sets of gene deletion to maximize the production of desired metabolite. BHFBA improves the performance of BAFBA as Hill climbing algorithm is a promising algorithm in finding local optimum. Experimental results on *B.subtilis* iBsu1103 model obtained from literature showed that BHFBA is effective in generating optimal solutions to the gene knockout prediction, and is therefore a useful tool in Metabolic Engineering [12]. In the future, to improve the performance of BHFBA we are interested in applying an automated pre-processing function in BHFBA to refine the genome-scale metabolic model. We are also interested in applying other fitness functions in BHFBA such as minimization of metabolic adjustment (MOMA) and regulatory on/off minimization (ROOM) to further improve the performance of BHFBA. Besides that, BA employs many tunable parameters which are difficult for the users to determine so it is important to find ways to help the users choose suitable parameters.

Acknowledgement. Institutional Scholarship MyPhD provided by the Ministry of Higher Education of Malaysia finances this work. We would like to thank Malaysian Ministry of Science, Technology and Innovation for supporting this research by an e-science research grant (Grant number: 06-01-06-SF1029). This research is also supported by UTM GUP research grant that was sponsored by Universiti Teknologi Malaysia.

References

1. Feist, A.M., Herrgård, M.J., Thiele, I., Reed, J.L., Palsson, B.Ø.: Reconstruction of biochemical networks in microorganisms. *Nat. Rev. Microbiol.* 7, 129–143 (2009)
2. Chandran, D., Copeland, W.B., Sleight, S.C., Sauro, H.M.: Mathematical modeling and synthetic biology. *Drug Discovery Today Disease Models* 5(4), 299–309 (2008)

3. Burgard, A.P., Pharkya, P., Maranas, C.D.: OptKnock: A bilevel programming framework for identifying gene knockout strategies for microbial strains optimization. *Biotechnol. Bioeng.* 84, 647–657 (2003)
4. Patil, K.R., Rocha, I., Förster, J., Nielsen, J.: Evolutionary programming as a platform for in silico metabolic engineering. *BMC Bioinformatics* 6, 308 (2005)
5. Rocha, M., Maia, P., Mendes, R., Pinto, J.P., Ferreira, E.C., Nielsen, J., Patil, K.R., Rocha, I.: Natural computation meta-heuristics for the in silico optimization of microbial strains. *BMC Bioinformatics* 9, 499 (2008)
6. Pham, D.T., Ghanbarzadeh, A., Koç, E., Otri, S., Zaidi, M.: The bees algorithm – a novel tool for complex optimization problems. In: *Proceedings of the Second International Virtual Conference on Intelligent Production Machines and Systems*, July 3-14 (2006)
7. Choon, Y.W., Mohamad, M.S., Deris, S., Chong, C.K., Chai, L.E., Ibrahim, Z., Omatu, S.: Identifying Gene Knockout Strategies Using a Hybrid of Bees Algorithm and Flux Balance Analysis for in silico Optimization of Microbial Strains. In: *The 9th International Symposium on Distributed Computing and Artificial Intelligence (DCAI 2012)*. University of Salamanca, Spain (2012)
8. Pham, D.T., Darwish, A.H., Eldukhri, E.E.: Optimisation of a fuzzy logic controller using the bees algorithm. *International Journal of Computer Aided Engineering and Technology* 1(2), 250–264 (2006)
9. Olague, G., Puente, C.: The honeybee search algorithm for three-dimensional reconstruction. In: Rothlauf, F., Branke, J., Cagnoni, S., Costa, E., Cotta, C., Drechsler, R., Lutton, E., Machado, P., Moore, J.H., Romero, J., Smith, G.D., Squillero, G., Takagi, H. (eds.) *EvoWorkshops 2006*. LNCS, vol. 3907, pp. 427–437. Springer, Heidelberg (2006)
10. Pham, D.T., Ghanbarzadeh, A.: Multi-objective optimisation using the bees algorithm. Paper. In: *Proceedings of the Third International Virtual Conference on Intelligent Production Machines and Systems*, July 2-13 (2007)
11. Cheng, M.Y., Lien, L.C.: A Hybrid Swarm Intelligence Based Particle Bee Algorithm for Benchmark Functions and Construction Site Layout Optimization. In: *Proceedings of the 28th ISARC*, Seoul, pp. 898–904 (2011)
12. Orth, J.D., Fleming, R.M.T., Palsson, B.Ø.: *Reconstruction and Use of Microbial Metabolic Networks: the Core Escherichia coli Metabolic Model as an Educational Guide*. ASM Press, Washington, DC (2009)
13. Henry, C.S., Zinner, J.F., Cohoon, M.P., Stevens, R.L.: iBsu1103: a new genome-scale metabolic model of *Bacillus subtilis* based on SEED annotations. *Genome Biology* 10, 69 (2009)
14. Bohl, K., de Figueiredo, L.F., Hadicke, O., Klamt, S., Kost, C., Schuster, S., Kaleta, C.: CASOP GS: Computing intervention strategies targeted at production improvement in genome-scale metabolic networks. In: *The 5th German Conference on Bioinformatics 2010*, September 20-22 (2010)
15. Kim, J.W., Chin, Y.W., Park, Y.C., Seo, J.H.: Effects of deletion of glycerol-3-phosphate dehydrogenase and glutamate dehydrogenase genes on glycerol and ethanol metabolism in recombinant *Saccharomyces cerevisiae*. *Bioprocess Biosyst. Eng.* 35, 49–54 (2012)