

Bach User's Guide

ZJU-China iGEM

October 2010

Contents

1.	Introduction	. 3
2.	System Requirement	. 3
3.	Algorithm	. 4
	3.1 Synonymous Substitution	. 4
	3.2 Slow	. 4
	3.3 Fast	. 5
	3.4 Match	. 5
	3.5 Get RiPS	. 6
4.	Using Bach	. 6
	4.1 Getting Started	. 6
	4.2 Data Prep	6
	4.3 Optimization	. 7
	4.4 Calculation of RiPS	8
	4.5 Data Output	10
5.	Copyright	10
6	Annondia	11

1. Introduction

Bach is a gene composer software that optimizes coding sequences and calculates corresponding RiPS to illustrate their translation rates. In sequence optimization, Bach incorporates 4 distinct approaches from which users can freely choose: SYNONYMOUS SUBSTITUTION, SLOW, FAST and MATCH. In the last section of GET RIPS, users are provided with a quantitative view of the translational behavior of the input coding sequence and the four sequences produced by the four methods as well.

Run Bach for the following research needs:

- Heterologous protein expression
- Evaluation of translation rates and behavior
- Design of more predictable and robust genetic circuits

Run Bach for the following benefits:

- Diverse optimization approaches of choice
- Quantitative Analysis of translational behavior
- Free and easy to use

2. System Requirement

Operating System:

- Recommended OS: Windows 7 32-bit Operating System
- Minimum OS: Windows XP 32-bit Operating System

Hardware Configuration:

- No strict limits.
- Recommended configuration:

 $\begin{array}{l} \mathsf{RAM} \ensuremath{\geqslant} \ 1\mathsf{GB} \\ \mathsf{Hard} \ \mathsf{Disk} \ensuremath{\geqslant} \ 100\mathsf{GB} \\ \mathsf{CPU} \ensuremath{\geqslant} \ 1.66\mathsf{GHZ} \end{array}$

3. Algorithm

3.1 Synonymous Substitution

Synonymous Substitution is more of a paraphrase than optimization, the goal of which is to retain and transfer the former performance of translation from the source species into the host species. The codon rank sequence is unchanged, implying the relative availability of aa-tRNA to bind to each codon site along the sequence chain is unchanged.

The input sequence is dealt with the following steps:

- The codon usage table of both the source and host species and the input sequence are loaded.
- For each amino acid, the corresponding n (1≤n≤6) codons are given rank from 1-n according to their relative amount in percentage.
- Each codon in the input sequence is assigned with the rank in the source species, thus producing a new sequence in terms of rank.
- Each codon in the input sequence is then substituted by codon of the same rank in the host species.

Example:

Consider a short coding sequence of AUUAUCAUA, the amino acid sequence of which is III. A is the source species, while B is the host species. The relative amount (in percentage) table of each codon for two species is as follows:

Codon	Amino Acid	Amount in A	Rank in A	Amount in B	Rank in B	
AUU	I 0.47		1	0.21	3	
AUC	I	0.32	2	0.47	1	
AUA	I	0.21	3	0.32	2	

Therefore, AUU which ranked number 1 in A should be substituted with the codon ranked number 1 in B, which is AUC. The new sequence produced by synonymous substitution then is AUCAUAAUU.

3.2 Slow

The output sequence of this phase, denoted by the phase name 'slow', is theoretically the slowest sequence one can obtain without altering the sequence of amino acid. As the codons adopted by the slow sequence correspond to a scarce source/concentration of its aa-tRNA, it's very much likely that the ribosome has to wait longer for the aa-tRNA to come and attach.

The input sequence is transformed with the following procedures:

- The codon usage table of the host species and the input sequence are loaded.
- For codons encoding the same amino acid, the codon with the smallest number/percentage is designated as the slowest codon for that amino acid.
- Each codon in the input sequence is substituted with the slowest codon of the same amino acid.

Example:

Continue to consider the example in Synonymous Substitution. Among the three codons encoding amino acid I, AUU is the slowest codon in species B. Thus all the codons in the sequence encoding amino acid I should be replaced by AUU in this approach, producing the slowest sequence: AUUAUUAUU.

3.3 Fast

Exactly the opposite of SLOW, the FAST section produces the fastest sequence of the same protein one can possibly get. Such coding sequence is usually translated fastest, since all the codons used in the sequence correspond to richest sources of aa-tRNAs, thus speeding up the step of aa-tRNA attachment.

The steps for such optimization is similar to the one in SLOW:

- The codon usage table of the host species and the input sequence are loaded.
- For codons encoding the same amino acid, the codon with the largest number/percentage is designated as the fastest codon for that amino acid.
- Each codon in the input sequence is substituted with the fastest codon of the same amino acid.

Example:

Still take use of the example in the previous two sections. In this method, all codons encoding amino acid I should be replaced by the fastest codon in species B, which is AUC, thus resulting in the output sequence: AUCAUCAUC.

3.4 Match

This section aims to produce the sequence that matches the codon bias of the host species the best without largely altering the input sequence. The product of MATCH serves beyond the scope of elongation process and may be subjected to less failure when expressed in the host organism.

Computation has been made according to the following procedures:

- The codon usage table of the host species and the input sequence are loaded.
- The input sequence is converted into a new sequence of amino acid.
- According to the protein sequence and the codon bias of the host species, the number of each codon is calculated by multiplying the number of its corresponding amino acid in the sequence and the percentage of the codon of each amino acid in the host species.
- The difference between the assumed number of each codon obtained from the previous step and the original number of each codon in the input sequence is calculated.
- The replace of codon in the input sequence then proceeds from 5'-3' according the difference obtained above, so that the output sequence has the proportion of each codon of the same amino acid approximately the same with that of the host species.

Example:

Let's now consider a longer sequence AUUAUUAUUAUUAUUAUUAUUAUUAUAUAAUAAGGG, in which

there're 10 amino acid I. The relative amount of each codon in species B is still the same as the previous examples. To match the codon bias in species B, we can first calculate the assumed amount of each codon encoding amino acid I in the sequence. The amount of AUU = $10*0.21\approx2$; the amount of AUC = $10*0.47\approx5$; the amount of AUA= $10*0.32\approx3$. The difference of each codon can then be calculated as -4 for AUU, 4 for AUC and 0 for AUA. The output sequence is AUUAUUAUCAUCAUCAUCAUCAUAUAUAAUAAGGG, in which the number of each codon encoding amino acid I agrees with the assumed value calculated from the codon bias of the host.

3.5 Get RiPS

One of the brilliance of Bach lies in GET RIPS, in which RiPS (Ribosome per sec) of each sequence can be calculated. RiPS is a newly emerged term in synthetic biology to denote translation rates. The computation is based on the mathematical model built by ZJU-CHINA 2010 iGEM team.

Please refer to <u>http://2010.igem.org/Team:ZJU-China/Modeling</u> for a detailed description of the underlying model and algorithm.

4. Using Bach

4.1 Getting Started

To start running Bach, double click the icon of Bach. The window for input will appear.

4.2 Data Prep

In the software package downloaded, there's a file 'Species Data' in which data of more than 1200 species is stored. The codon usage table of each species is stored in each distinct text file in 'Species Data', named by the species Latin name. Click Load a Source Species button and load the source species (the original host of the gene or the species where the gene is first discovered) from the 'Species Data' file in the software package. Then select it from the list above the Load a Source Species button.

Choosing the host species is similar. Click Load a Host Species button and load the host species (the host you want the gene to be expressed) from the 'Species Data' file in the software package. Then select it from the list above the Load a Host Species button.

Note that for both the source species and host species you've loaded, the data and history will be kept, which means you don't have to load the same species again in the future. If you want to clean the history of loading, you can manage to do so by deleting the database file in your software package. Then, you'll have to load each species data again.

If you either are unsatisfied or simply can't find the species you desire in the 'Species Data' we've

provided, you can make your own data file and load it into the software, as long as the form of the data file is exactly the same as those we provide. More information on how to make eligible data file can be found in the Appendix.

٩	Bach	Load a Sour	ce Species						? 🛛
	Source Species	查找范围(I):	🚞 Test Sequen	ce	*	6	1 🖻	 -	
	Input Species Codon: Saccharomyces_cerev v Load a Source Species Input Sequence:	ました (1) (1) (1) (1) (1) (1) (1) (1) (1) (1)	Sample DNA E Sample prote Sample RNA	in					
	<sequence file="" name=""></sequence>	"美国" 没有的文档							
	Data Export >>	家の実施							
	Host Species Input Species Codon:	《 网上邻居							
	Load a Host Species		文件名 (2): 文件类型 (2):	Sample DNA Text File(*.txt)			~		打开 (<u>(</u>) 取消
	CMD History NO.2: Add rankJoin info success. NO.1: Add host species codon success. NO.0: Add source species codon success.								

Click Load a Sequence button and load the sequence file for optimization. In the software package, there's a 'Test Sequence' file for tryout if you don't have any gene sequence at hand. If you do have a sequence, please make your sequence file eligible for Bach to recognize:

- State the nature of the input sequence (DNA, RNA or protein) in the first line.
- Paste your sequence in the second line. The sequence can be acknowledged and extracted by Bach as long as there're no characters other than numbers and English letters.
- As Bach can't distinguish between coding sequences and non-coding sequences, please make sure the sequence you've submitted is from the coding region, otherwise Bach wouldn't be of much help in the optimization and calculation of RiPS.

4.3 Optimization

Click >> button and the sequence optimization results are displayed on the right. You can switch to view sequences optimized by different methods by clicking on the different tags in the upper window. If you don't have any idea what these sequences stand for or how these sequences are obtained, you can refer to the previous section of this manual on Algorithm.

In each window, there's a description to better guide you through each window and find what you need. The box in the middle displays the input sequence you've loaded. The optimized sequence is displayed in the box underneath. The sequence is displayed in the unit of codon to better illustrate the change of codon and the number on the left helps you to better track the length of the sequence.

🗞 Bach					
Source Species Input Species Codon: Saccharomyces_cerev Load a Source Species	Synonymous Substitution Slow Fast Match Get RiPS Description Each codon in the input sequence is substituted with the codon of the same rank for the output. Each codon of the same amino acid is first ranked in both the source and the host species respectively.				
Input Sequence: I:\Bach\Test Sequence\Sample DNA.txt Load a Sequence Data Export <<	the source species Input Sequence: DOI CGC UGU GGU ACA CGC UGU GGG ACC GCU ACG GCC UGU AUG UGG UGG O45 AUG AAG CCA AUA UUG AAA CCC ACG GCA UGG UGC CAA UGA AUC GUC O40 UGA CCG AUG AUC CGC GCU GGC UAC CGG CGA UGA UGA UCA UCU 135 CGC GAA UGG UGC AGC GCG AUC GUA AUC ACC CGG UGA AUC ACG ACG UGA 135 CGA UGU GCC UGG GGA UG AAU CAG GCC ACG GCC AA AUC ACG ACG CGC UAA 135 CGA UGU ACC CGG GG GA GG CGA CCC CGG CUA AUC ACG ACG CGC UGA 210 AUC GCC UGG GGA UCA AAU CUG UCG AUC CUU CCC GCC UGC AGU 2315 CGA UGU ACG CGC GGG GUG AUG AAG ACC AGC CCU UCC CGG CUG UGC				
Host Species Input Species Codon: Escherichia_coli Load a Host Species CMD History No.0: Read Sequence file success.	the host species Output match sequence: Output match sequence: O1 CGA UGU GGU ACA CGA UGU GCU ACU GCA ACG GCG UGU AUG UGG UGG O45 AUG AGA CCC AUC CUG AAA CCA ACG GCC UGG UGC CAG UGA AUA GUG O90 UGA CCC AUG AUA CGA GGA UAC AGG AGA UGA GCU AAC GCU UAA 135 CGA GAA UGG UGC UGC GCU AUA GUC AIDA CUA AGG GU CUA AUCU AAC 180 GGU CGA UGG GGC AUG AAU CAA GCG ACG GCU CUU AUA ACG ACG CGA 2225 UGU AUA GCA GGC UCU AAU CUC UCG AUA UGG CCC AUC UUA ACG AGU 2315 AGA UGU ACG CGA GCU UGG AUG AAG ACU UCC CCU UCA AGG CUC UGC				

4.4 Calculation of RiPS

To get RiPS of each sequence, click the tag Get RiPS. For a more precise calculation of RiPS, you need to first change the values of parameters.

The parameters on the left in the RiPS Parameters box relate to the Monte Carlo simulation itself and are species independent:

- Walpha: The parameter denoting the initiation rate of translation, the value of which may be influenced by the RBS sequence and other nonsense sequences in both the 5' and 3' tails of the mRNA.
- T: Time for Monte Carlo simulation in terms of 0.001s. The more times Monte Carlo run, the more accurate value of RiPS you'll get. Values more than 100 for T is highly recommended, otherwise the calculation results will display a much more random fashion.
- Length: The length of ribosome occupation in terms of codon. The value of 11 is recommended according to published researches. The range of the parameter is set as 1≤ Length≤Total Length.
- Look Index: The site along the sequence where RiPS is calculated. Look Index of 20 indicates that RiPS is calculated by counting the number of ribosomes passing the 20th codon (5'->3') per second. As the first 0-20 codons constitute a zone which performs quite differently from the latter part of the sequence, Look Index of values from 20-n is highly recommended, where n is the total length of the sequence in terms of codon displayed in the box on the right.

The parameters on the right in the Advanced Parameters box may be species dependent. The default values have been obtained from experiments for Escherichia coli. If you are setting E.coli as the host species or want to calculate the RiPS of each sequence when expressed in E.coli, no

alteration has to be made about the default parameters. If not, there may be a need to reevaluate these parameters before calculating RiPS. Furthermore, if you are not satisfied with the default value of parameters in either case, you can always change it based on your own research or other more reliable sources you'd like.

Advanced parameters are as follows:

- Wg: Rate constant for transition from 4-5, in which tRNA is shifted and site A is occupied by EF-G in the GTP bound form.
- Wp: Rate constant for transition from 2-1, when falsely bound aa-tRNA-EF-Tu complex dissociates from site A.
- Wh1: Rate constant for transition from 2-3, denoting the probability of the GTP part of EF-Tu hydrolized to GDP.
- Wh2: Rate constant for transition from 5-6, representing the hydrolysis of GTP to GDP and released. As transition from 6-7-1 can be negligible compared with transition from 5-6, wh2 is used in the model to indicate the transition rate from 5-1.
- K2: Rate constant for transition from 3-4, when the phosphate group, a product of the hydrolysis, leaves and releases the EF-Tu.
- Wbeta: The dissociation rate of ribosome from mRNA at the stop codon.
- Codon total length: The length of the selected sequence in terms of codon. The value of total length casts limits on values of other parameters such as Length and Look Index, thus helping users to set reasonable values for other parameters.

Click Data Confirm button to confirm the values of all the parameters you want to use for the calculation of RiPS. Then choose the sequence you wish to compute RiPS. Upon clicking on each sequence tag, the corresponding sequence will be instantly displayed in the box beneath. Click RiPS Calculation to finally compute.

A green bar of progress will appear at the bottom of the window to show how long you have to wait for the result. When the bar is completed, a small box will pop out, on which the value of RiPS is shown. Upon closing the small box, the result will also be shown under the bar.

💊 Bach	
Source Species Input Species Codon: Saccharomyces_ceret Load a Source Species	Synonymous Substitution Slow Fast Match Get RiPS Description For each coding sequence obtained in the previous phases, RiPS (Ribosome per sec) can be calculated to denote each translation rate. Note that the default value of all parameters only work for E coli as the host species.
Input Sequence: I:\Bach\Test Sequence\Sample DNA.txt Load a Sequence Data Export << Heat Sequence	RiFS Parameters Walpha: Bach T: 100 Length: 11 LookIndex: 20 Choose sequence for RiPS calculation: 157 Data Confirm Stronymous
Input Species Codon: Escherichia_coli Load a Host Species CMD History NO.3: Read Sequence file success. NO.1: Add host species codon success. NO.1: Add host species codon success. NO.0: Add source species codon success.	Calculate Sequence: OD1 CGC UGU GGU ACA CGC UGU GCG ACC GCU ACG GCC UGU AUG UGG UGG OD5 AUG AAG CCA AUA UUG AAA CCC ACG GCA UGG UGC CAA UGA AUC GUC OD90 UGA CCG AUG AUC CGC GCU GGC UAC CGG CCA UGA GCG AAC GGC UAA ISO GGU CGC UGG GGA AUG AAU CAC GCC ACG GCG CUA AUC ACG ACG CGC RiPS calculation accomplish bar: RiPS Calculation RIPS Outcome:

4.5 Data Output

By clicking the Data Export button on the left, the optimization results (containing newly composed sequences from different methods) can be exported into a text file.

5. Copyright

Bach is Copyrignt © ZJU-CHINA iGEM 2010. All rights reserved.

For further queries related to Bach program, please contact us through the email of <u>wendyhu001@gmail.com</u>. Any reports of failure or suggestions on Bach are welcome.

6. Appendix

Species Data Compilation

If you can't find the target species in the Species Data file, you can easily make your own data file by the following procedures:

- Go to Codon Usage Database: <u>http://www.kazusa.or.jp/codon/</u>.
- Search for or choose the codon usage table of your target species. Here we use Saccharomyces cerevisiae as an example.
- Choose '1. Standard' from the Format list under the table.
- Click the Submit button underneath.
- The table is changed into a new format, in which there're four columns of information for each codon. Copy the table into a new text file.
- The text file format should be strictly the same with the example figure below. The number in the first line should be bigger than 1232 and should be distinct for each species. The blank space ' ' in the Latin name of the species in the second line should be replaced with '_'.

📕 Saccharomyces cerevisiae - 记事本	
又件ぜり 編輯ぜり 格式 回り 登者 ぜり 帮助 低り	
1234	~
Saccharomuces cerevisiae	
HULL F 0.59 26.1 (170666) HELLS 0.26 23.5 (153557)	(HAH Y 8.56 18.8 (122728) HGH C 8.63 8.1 (52983)
HILC F 8 41 18 4 (128518) HCC S 8 16 14 2 (92923)	HAC Y 8 44 14 8 (96596) HGC C 8 37 4 8 (31895)
UUA 0.28 26.2 (170884) UCA S 0.21 18.7 (122028)	HAA * 8.47 1.1 (6913) HGA * 8.38 8.7 (4447)
UNC 0 20 27 2 (177573) UNC S 0 10 8 6 (55051)	
	ona - 0.20 0.5 (0012) oaa ii 1.00 10.4 (01107)
CIIII I 8 19 19 9 / 08876\ CCII D 8 91 19 E / 00969\	COULU 8 66 19 6 / 90887) CCU D 8 16 6 6 / 61701)
CUC L 0.13 12.3 (00070) CCC F 0.31 13.5 (00203)	$\begin{array}{cccccccccccccccccccccccccccccccccccc$
CUC L 0.00 5.4 (35545) CCC F 0.15 0.8 (44309)	CHC N 0.30 7.6 (50765) CGC N 0.00 2.0 (10993)
CUH L 0.14 13.4 (87019) CCH P 0.42 18.3 (119641)	GHH Q 0.09 27.3 (178251) GGH K 0.07 3.0 (19502)
CUG L 0.11 10.5 (68494) CCG P 0.12 5.3 (34597)	CHG Q 0.31 12.1 (79121) CGG R 0.04 1.7 (11351)
AUU I 0.46 30.1 (196893) ACU T 0.35 20.3 (132522)	AAU N 0.59 35.7 (233124) AGU S 0.16 14.2 (92466)
AUC I 0.26 17.2 (112176) ACC T 0.22 12.7 (83207)	AAC N 0.41 24.8 (162199) AGC S 0.11 9.8 (63726)
AUA I 0.27 17.8 (116254) ACA T 0.30 17.8 (116084)	AAA K 0.58 41.9 (273618) AGA R 0.48 21.3 (139081)
AUG M 1.00 20.9 (136805) ACG T 0.14 8.0 (52045)	AAG K 0.42 30.8 (201361) AGG R 0.21 9.2 (60289)
GUU V 0.39 22.1 (144243) GCU A 0.38 21.2 (138358)	GAU D 0.65 37.6 (245641) GGU G 0.47 23.9 (156109)
GUC V 0.21 11.8 (76947) GCC A 0.22 12.6 (82357)	GAC D 0.35 20.2 (132048) GGC G 0.19 9.8 (63903) 💻
GUA V 0.21 11.8 (76927) GCA A 0.29 16.2 (105910)	GAA E 0.70 45.6 (297944) GGA G 0.22 10.9 (71216)
GUG V 0.19 10.8 (70337) GCG A 0.11 6.2 (40358)	GAG E 0.30 19.2 (125717) GGG G 0.12 6.0 (39359) 🧓
· · · · · · · · · · · · · · · · · · ·	· · · · · · · · · · · · · · · · · · ·
S.	