

<Summary>

Development of new generation sequencers enabled genome sequencing feasible for every organism in a laboratory. A typical data flow of de novo sequencing includes (1) assembly of sequence reads, (2) estimation of open reading frames, (3) annotation of proteins, and (4) finding RNA genes. The annotation is normally performed by BLASTP searches against several different databases. However, it is usually hard to find a plausible annotation by just looking at the results of BLASTP searches.

Here I propose a potentially automatic method of annotation that exploits automatic protein clustering using the software GCLUST, which estimates proper similarity threshold for each list of homologs using 'entropy-optimized organism count' method (Sato 2009). The software has been used to construct a homolog database including both prokaryotic and eukaryotic proteins (<http://gclust.c.u-tokyo.ac.jp/>). For use in genome annotation, we need de novo clustering including many genomes of related organisms as well as genomes of representative organisms. Application of protein clustering in the annotation in *Arthrospira platensis* was the first successful case (Fujisawa et al. 2010). I present here results of protein clustering of total predicted proteins in two draft genomes of cyanobacteria along with total predicted proteins of 41 cyanobacteria available at NCBI. For each of the resultant protein clusters, an alignment and a phylogenetic tree were also prepared for assistance in functional annotation. The quality of alignments was evaluated by counting ill-aligned proteins (missing N- or C-terminus, or insertion/deletion), which was 4-13% of total predicted proteins in most cyanobacterial genomes. Annotation may be automated by extracting significant key words already assigned for member proteins of clusters or by comparison with reference protein clusters.

1. Introduction

Current way of genome sequencing

- DNA isolation from bacterial cells
- Library construction
- Sequencing (454 etc)
- Assembly (newbler, MIRA)
- Annotation pipeline (MiGAP: Sugawara & Kurokawa labs.)
 - ORF estimation (MetaGeneAnnotator)
 - RNA genes estimation (tRNAscan-SE, RNAmmer)
 - BLAST (COG, RefSeq, TrEMBL) : **to be improved by all-against-all BLASTP and automatic clustering with Gclust → de novo ortholog clusters**
 - Automatic annotation
 - N-terminal correction

5. Prospects

For annotation of better quality

- Proposal:
 - Annotation of a new genome should be performed in the framework of related genomes
- This process is assisted by automatic clustering of all proteins in related genomes
- N-terminal correction will give better alignment
- Consensus annotation may be easy

3. Automatic protein clustering with Gclust software

Entropy-optimized organism count (EEOC) method

Problem: All previous clustering of proteins used a simple threshold value such as $E = 1.0 \times 10^{-6}$ (criterion may be other parameters), but similarity of proteins is very different in different protein families: eg. PsaA: 10^{-150} , PsbO: 10^{-45} , PsbL: 10^{-10} . Use of a single threshold should produce unusually large clusters containing unrelated proteins and divergent paralogs.

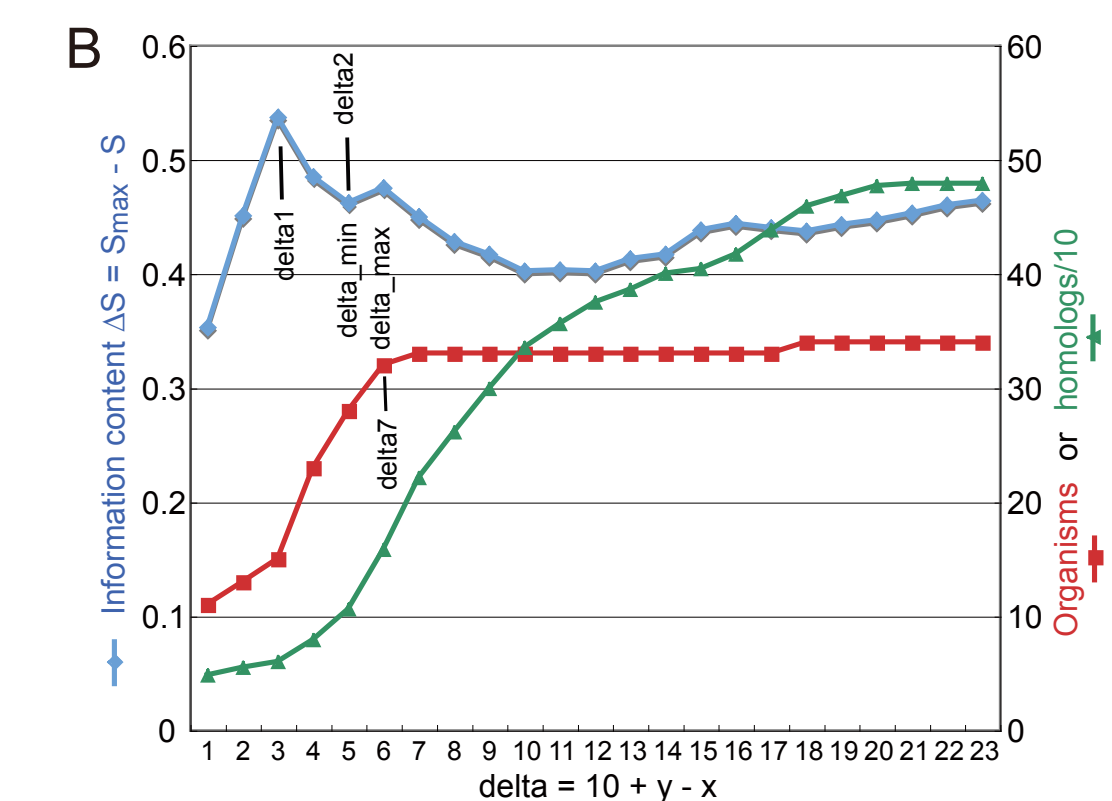
Solution: For each list of protein similarity data (may include two variables), organisms are counted from the top. This is justified by the fact that orthologs are usually found near the top, and then, paralogs, then partially similar proteins. However, there are many different cases. Entropy of distribution is useful in estimating a proper threshold. In the actual implementation in Gclust, the two measures are considered to obtain orthologs and highly related paralogs.

A xy-table

overlap sc	0	0.1	0.2	0.3	0.4	0.5	0.6	0.7	0.8	0.9	1
thr	0	0	0	0	0	0	0	0	0	0	1
1.00E-99	0	0	0	0	0	0	0	0	0	0	4
1.00E-90	0	0	0	0	0	0	0	0	0	1	4
1.00E-80	0	0	0	0	0	0	0	0	0	0	1
1.00E-70	0	0	0	0	0	0	0	0	0	0	1
1.00E-60	0	0	0	0	0	0	0	0	0	0	1
1.00E-50	0	0	0	2	5	3	16	19	29	18	1
1.00E-45	0	0	0	0	1	1	10	10	13	10	0
1.00E-40	0	0	0	1	2	4	11	16	11	3	2
1.00E-35	0	0	0	1	2	5	8	2	6	1	0
1.00E-30	0	0	0	1	1	1	5	3	5	2	1
1.00E-25	0	0	0	0	2	1	4	1	1	0	0
1.00E-20	0	0	0	0	2	1	4	2	3	0	0
1.00E-16	0	0	0	4	12	7	7	0	0	0	0
1.00E-13	0	0	0	2	1	8	6	2	0	0	0
1.00E-10	0	0	1	1	1	1	3	2	0	0	0
1.00E-08	0	0	3	4	1	2	0	0	0	0	0
1.00E-06	0	0	7	2	1	0	0	0	0	0	0
1.00E+00	1	0	0	0	0	0	0	0	0	0	0

Best local maximum
Finally optimized border
Line for $\Delta S = 5$
 $-5 = x - 10$

In the xy table, consider an upper triangular area defined by $(x=10, y=0)$ and $(y - \Delta S) = (x - 10)$. Let P_i be the number of proteins in division i , and the number of divisions in the triangle be n . The entropy of distribution is $S = - \sum P_i \cdot \log P_i$. Maximal S is defined by $S_{max} = - \sum P_i \cdot \log P_i = \log n$. We use $\Delta S = S_{max} - S$ for estimating a local maximum.



Other merits of the Gclust software:

1. Automatic domain identification based on homology region data of BLAST. This is used for excluding multidomain proteins.
2. Automatic identification of transit peptides based on domain identification.
3. Comparison of both prokaryotic and eukaryotic proteins in a single dataset. This is the reason why Gclust is the only software that can analyze proteins of endosymbiotic origin.
4. A powerful replacement of COG. Gclust clusters are suitable for annotation of data from new generation sequencer.

Reference

Sato, N. (2009) Gclust: trans-kingdom classification of proteins using automatic individual threshold setting. *Bioinformatics* 25: 599-605

2. Revision of functional categories

Problems of functional categories of COG and a revision

COG is based on clustering of proteins of 61 organisms, which are divided into 14 groups.

Only ortholog groups shared by ≥ 3 groups are assigned COG number

Category 'Energy production and conversion' does not include photosynthesis

Difficulty in assignment of transporter and DNA-binding proteins to a correct functional category

4(B) Improvement of annotation

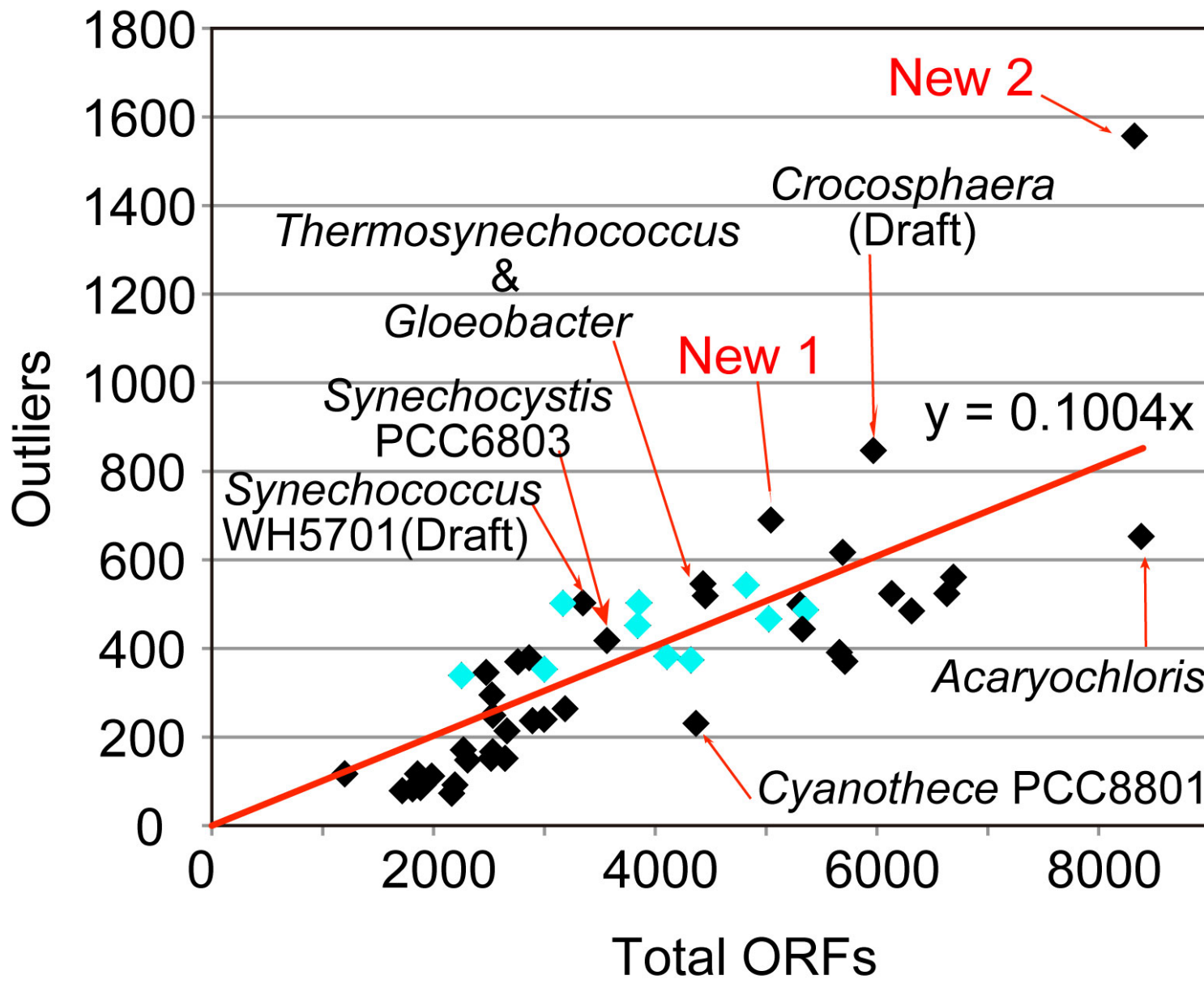
CyanoClust is a database of ortholog groups of cyanobacteria

New functional categories of proteins

This classification is implemented in CyanoClust database. Number of clusters is shown for each sub-category.

Main category	Sub-category	Clusters	Main category	Sub-category	Clusters	
gene expression	genome structure	11	inorganic carbon		6	
	replication	10	ATP synthase		10	
	repair, recombination, modification and nuclease	23	central metabolism and sugar		41	
	transcription	13	lipid		32	
	ribosome, translation	111	porphyrin, heme, cytochrome, pigment		38	
	RNA metabolism	13	phosphorus and sulfur		11	
	sorting	11	nitrogen and amino acid nucleotide		93	
	processing and degradation	14	cofactor biosynthesis		33	
	regulation			other metabolism	24	
	signal transduction	12	cellular structure			
	stress response and chaperon	24				
	metabolism	photosystem	35	extracellular matrix		28
		respiration	15	cell division		13
		hydrogenase	5	transport and membrane		30
			unclassified	unclassified	4	
			hypothetical	hypothetical	248	
			Total		980	

4(A) Improvement of N-termini



The N-termini can be improved by comparing sequences within a cluster. Outliers as defined as proteins having a protruding or lacking N-terminus are about 10% of total ORF of a genome, but are significantly larger than this average in draft sequences. The N-terminus problem exists also in published genomes. Many of them can be improved by inspecting the N-terminal nucleic acid sequence.

Annotation of universally (at least within a phylum) conserved proteins

Each of the 980 conserved protein clusters in 41 cyanobacteria have been given a biologically correct annotation (see above).

This annotation can be transferred to a new cluster constructed for (new + 41 species).

Cluster-based annotation is useful in avoiding 'Inherited strange annotation'

- Some annotations are inherited from those of other genomes based on unreliable homology or improper biological knowledge
- Cluster-based annotation is not susceptible for such inappropriate inheritance of annotation, even though individual annotations (given for original databases) may be variable or sometimes unreliable

An example of consensus-based annotation

Before correction

```
S79_Synpcc7942_0175 ..... MKGLV ..... TFLICLQTLWLGALTAALTAQNPISLVSVLLVLSVRL
Lim_c00231g12 ..... -MGLAIFSSQNIYAVTVKLGAFESIKL
Mae_MAE_31110 ..... -NPLASAIYGGWIMTMAVFAIQNTQVPSLKFQFSIKV
Syn_sll1193 ..... -DLDFGCSGKNPGLGNGNLGNLSEVPLTAGSPFMWERSS
Phks_c00261g3 ..... -MKSTA ..... SLLACSLTAVVAIYFQVQATVLSKFLGYSIKL
Cth4_PCC7424_5099 ..... -MDTLT ..... -NLASILIAWMITLAVFSIQNTTTPVSLKPMFEFPGL
S70_SYNPCC7002_A2518 ..... -MQLTG ..... -RLIAGLLAAVWVVAIYVFSIQNTTTPVSLKPMFEFPGL
Cth1_PCC8801_1340 ..... -MKALT ..... -NLTAIIAIPLWGALIAIYFSIQNTTTPVSLKPMFEFPGL
Cth2_ccc_0816 ..... -MMSMLLT ..... -NVIISVVAIWMGALIAIYFSIQNTTTPVSLKPMFEFPGL
Cnat_g5622 ..... -MKCQT ..... -NVIIVAVMAGWELGALIAIYFSIQNTTTPVSLKPMFEFPGL
Tel_tr0124 ..... -MR ..... -RFLLLMMISSGAILFVQNAKAVSLKFLMWSIQM
Cth5_Cyan7425_1365 ..... -MAALAVPVQNAKAVSLKFLMWSIQM
Amar_AM1_5738 ..... -MHLYGNVPSLQSLVQNAKAVSLKFLMWSIQM
Amax_004086 ..... -MKALP ..... -PFTSLIAVAIAVAIVLQVQNTAVVSLQPMIFPESINI
Apl_NIES39_D06280 ..... -MKALP ..... -PFTSLIAVAIAVAIVLQVQNTAVVSLQPMIFPESINI
Npun_Npun_F6174 ..... -MKLA ..... -PFTSLIVVAIAVAIVLQVQNTAVVSLQPMIFPESINI
Ana_all1363 ..... -MNVGVLRFVVDVDPKMLTA ..... -PFTSLVVAIAVAIVLQVQNTAVVSLKFLMWSIQM
Ava_Ava_1738 ..... -MNACVHLRFVVDVDPKMLTA ..... -PFTSLVVAIAVAIVLQVQNTAVVSLKFLMWSIQM
ruler 1 .....10.....20.....30.....40.....50.....60.....70.....80.....90.....
```

After correction

```
S79_Synpcc7942_0175 ..... MKGLV ..... TFLICLQTLWLGALTAALTAQNPISLVSVLLVLSVRL
Lim_c00231g12 ..... -MLRLLLYLCIWMGSMGLAIFSS
Mae_MAE_31110 ..... -MNTIENPLASAIYGGWIMTMAVFAIQNTQVPSLKFQFSIKV
Syn_sll1193 ..... -DLDFGCSGKNPGLGNGNLGNLSEVPLTAGSPFMWERSS
Phks_c00261g3 ..... -MKSTA ..... SLLACSLTAVVAIYFQVQATVLSKFLGYSIKL
Cth4_PCC7424_5099 ..... -MDTLT ..... -NLASILIAWMITLAVFSIQNTTTPVSLKPMFEFPGL
S70_SYNPCC7002_A2518 ..... -MQLTG ..... -RLIAGLLAAVWVVAIYVFSIQNTTTPVSLKPMFEFPGL
Cth1_PCC8801_1340 ..... -MKALT ..... -NLTAIIAIPLWGALIAIYFSIQNTTTPVSLKPMFEFPGL
Cth2_ccc_0816 ..... -MMSMLLT ..... -NVIISVVAIWMGALIAIYFSIQNTTTPVSLKPMFEFPGL
Cnat_g5622 ..... -MKCQT ..... -NVIIVAVMAGWELGALIAIYFSIQNTTTPVSLKPMFEFPGL
Cth5_Cyan7425_1365 ..... -MISICLGLIGIMMAAIAIYFSIQNTAVVSLKFLMWSIQM
Tel_tr0124 ..... -HRPFLMLLMSISSGAILFVQNAKAVSLKFLMWSIQM
Amar_AM1_5738 ..... -MHLLSLLGSLGIAIAVAIVLQVQNTAVVSLQPMIFPESINI
Amax_004086 ..... -MKALP ..... -PFTSLIAVAIAVAIVLQVQNTAVVSLQPMIFPESINI
Apl_NIES39_D06280 ..... -MKALP ..... -PFTSLIAVAIAVAIVLQVQNTAVVSLQPMIFPESINI
Npun_Npun_F6174 ..... -MKIA ..... -PFTSLIVVAIAVAIVLQVQNTAVVSLKFLMWSIQM
Ana_all1363 ..... -MNTAVPFLVVAIAVAIVLQVQNTAVVSLKFLMWSIQM
Ava_Ava_1738 ..... -MNTAVPFLVVAIAVAIVLQVQNTAVVSLKFLMWSIQM
ruler 1 .....10.....20.....30.....40.....
```

Published genomes can also be improved.

Organism	Total ORFs	Outliers	Selected for N-terminal correction	Remaining outliers
Synechocystis sp. PCC 6803	3564	418 (11%)	199	86
Anabaena sp. PCC 7120	6132	524 (8%)	189	70
Arthrospira platensis	6630	524 (7%)	86	38
New cyano 1	5043	690 (13%)	190	85
New cyano 2	8323	1557 (18%)	243	112

Annotation for universally conserved proteins is implemented in CyanoClust

Annotation : respiration NdhE, NADH dehydrogenase kappa subunit