

# Gene Similarity-based Approaches for Determining Core-Genes of Chloroplasts

Bassam AlKindy\*<sup>†</sup>, Christophe Guyeux\*, Jean-François Couchot\*, Michel Salomon\*, Jacques M. Bahi\*

\*FEMTO-ST Institute, UMR 6174 CNRS, DISC Computer Science Department,  
University of Franche-Comté, France

<sup>†</sup>Department of Computer Science, University of Mustansiriyah, Baghdad, Iraq  
{bassam.al-kindy, christophe.guyeux, jean-francois.couchot, michel.salomon, jacques.bahi}@univ-fcomte.fr

**Abstract**—<sup>1</sup> In computational biology and bioinformatics, the manner to understand evolution processes within various related organisms paid a lot of attention these last decades. However, accurate methodologies are still needed to discover genes content evolution. In a previous work, two novel approaches based on sequence similarities and genes features have been proposed. More precisely, we proposed to use genes names, sequence similarities, or both, insured either from NCBI or from DOGMA annotation tools. Dogma has the advantage to be an up-to-date accurate automatic tool specifically designed for chloroplasts, whereas NCBI possesses high quality human curated genes (together with wrongly annotated ones). The key idea of the former proposal was to take the best from these two tools. However, the first proposal was limited by name variations and spelling errors on the NCBI side, leading to core trees of low quality. In this paper, these flaws are fixed by improving the comparison of NCBI and DOGMA results, and by relaxing constraints on gene names while adding a stage of post-validation on gene sequences. The two stages of similarity measures, on names and sequences, are thus proposed for sequence clustering. This improves results that can be obtained using either NCBI or DOGMA alone. Results obtained with this “quality control test” are further investigated and compared with previously released ones, on both computational and biological aspects, considering a set of 99 chloroplastic genomes.

**Index Terms**—Chloroplasts, Clustering, Quality Control, Methodology, Pan genome, Core genome, Evolution

## I. INTRODUCTION

The idea motivating the importance of identifying core genes is to understand the shared functionality of a given set of species. We introduced in a previous work [1] two methods for discovering core and pan genes of chloroplastic genomes using both sequence similarity and alignment based approaches. To determine these core and pan genomes for a large set of DNA sequences, we propose in this work to improve the alignment based approach by considering a novel sequence quality control test. More precisely, we focus on the following questions considering a collection of 99 chloroplasts: how can we identify the best core genome (an artificially designed set of coding sequences as close as possible to the real biological one) and how to deduce scenarii regarding their gene loss.

The term Chloroplast comes from the combination of plastid and chloro, meaning that it is an organelle found in plant and eukaryotic algae cells which contains chlorophyll. Chloroplasts may have evolved from *Cyanobacteria* through endosymbiosis and since their main objective is to conduct photosynthesis, these fundamental tiny energy factories are present in many organisms. This key role explains why chloroplasts are at the basis of most trophic pyramids and thus responsible for evolution and speciation. Moreover, as photosynthetic organisms release atmospheric oxygen when converting light energy into chemical energy and simultaneously produce organic molecules from carbon dioxide, they originated the breathable air and represent a mid to long term carbon storage medium. Consequently, exploring the evolutionary history of chloroplasts is of great interest and therefore further phylogenetic studies are needed.

An early study of finding the common genes in chloroplasts was realized in 1998 by *Stoebe et al.* [2]. They established the distribution of 190 identified genes and 66 hypothetical protein-coding genes (*yxf*) in all nine photosynthetic algal plastid genomes available (excluding non photosynthetic *Astasia tonga*) from the last update of plastid genes nomenclature and distribution. The distribution reveals a set of approximately 50 core protein-coding genes retained in all taxa. In 2003, *Grzebyk et al.* [3], have studied the core genes among 24 chloroplastic sequences extracted from public databases, 10 of them being algae plastid genomes. They broadly clustered the 50 genes from *Stoebe et al.* into three major functional domains: (1) genes encoded for ATP synthesis (*atp* genes); (2) genes encoded for photosynthetic processes (*psa* and *psb* genes); and (3) housekeeping genes that include the plastid ribosomal proteins (*rpl* and *rps* genes). The study shows that all plastid genomes were rich in housekeeping genes with one *rbcLg* gene involved in photosynthesis.

To determine core chloroplast genomes for a given set of photosynthetic organisms, bioinformatics investigations using sequence annotation and comparison tools are required, and therefore various choices are possible. The purpose of our research work is precisely to study the impact of these choices on the obtained results. A general presentation of the approaches we propose is provided in Section II. A closer examination of the approaches is given in Section III. Section III-A will present coding sequences clustering method based on sequence

<sup>1</sup>This paper was submitted to IEEE International Conference on Bioinformatics and Biomedicine (BIBM 2014) on 15/07/2014, Accepted 05/09/2014. Oral presentation was on 05/11/2014.

similarity, while Section III-B will describe quality test method based on quality genes. The paper ends with a discussion based on biological aspects regarding the evolutionary history of the considered genomes, leading to our methodology proposal for core and pan genomes discovery of chloroplasts, followed by a conclusion section summarizing our investigations.

## II. AN OVERVIEW OF THE PIPELINE

Instead of considering only gene sequences taken from NCBI or DOGMA [4], an improved quality test process now takes place as shown in Figure 1. It works with gene names and sequences, to produce what we call “quality genes”. Remark that such a simple general idea is not so easy to realize, and that it is not sufficient to only consider gene names provided by such tools. Providing good annotations is an important stage for extracting gene features. Indeed, gene features here could be considered as: gene names, gene sequences, protein sequences, and so on. We will subsequently propose methods that use gene names and sequences for extracting core genes and producing chloroplast evolutionary tree.

Real genomes were used in this study, which cover eleven types of chloroplast families (see [1] for more details). Furthermore, two kinds of annotations will be considered in this document, namely the ones provided by NCBI on the one hand, and the ones by DOGMA on the other hand.

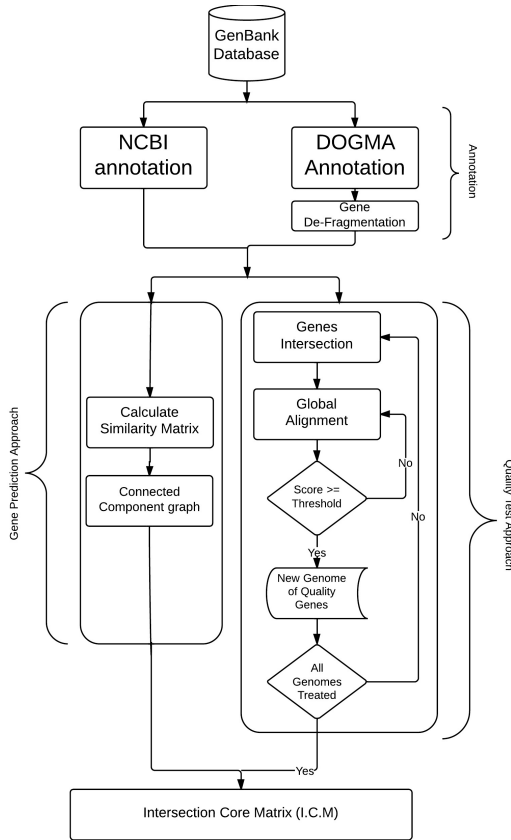


Fig. 1: An overview of the pipeline

## III. CORE GENES EXTRACTION

To make this document self contained, we recall the same definition with a fast revision of similarity based method.

### A. Similarity-based approach

Basically, this method starts with annotated genomes either from NCBI or DOGMA and uses a distance  $d : N = \{A, T, C, G\}^* \times \{A, T, C, G\}^* \rightarrow [0, 1]$  on genes coding sequences to group similar alleles in a same cluster.

For a given threshold  $T \in [0, 1]$  and a similarity measure  $d$ , the method builds the *similarity* undirected graph where vertices are alleles and s.t. there is an edge between  $g_i$  and  $g_j$  if we have  $d(g_i, g_j) \leq T$ .

Each connected component (CC) of this graph defines a class of the DNA sequences and is abusively called a “gene”, whereas all its nodes (DNA sequences) are the “alleles” of this gene. Let  $\pi$  the function that maps each sequence into its representative gene. Each genome  $G = \{g_1^G, \dots, g_{m_G}^G\}$  is thus mapped into the set  $\{\pi(g_1^G), \dots, \pi(g_{m_G}^G)\}$  where duplicated genes are removed.

Consequently, the core genome (resp., the pan genome) of two genomes  $G_1$  and  $G_2$  is defined as the intersection (resp., as the union) of their projected genomes. The intersection (resp. the union) of all the projected genomes constitutes the core genome (resp. the pan genome) of the whole species.

Let us now consider the 99 chloroplastic genomes introduced earlier. We use in this case study either the coding sequences downloaded from NCBI website or the sequences predicted by DOGMA. Each genome is thus constituted by a list of coding sequences. In this illustration study, we have evaluated the similarity between two sequences by using a global alignment. More precisely, the measure  $d$  introduced in the first approach is the similarity score provided after a Needleman-Wunch global alignment, by the *emboss* package released by EMBL [5].

The number of genes in the core genome and in the pan genome have been computed. Obtained results from various threshold values are represented in Table I. Remark that when the threshold is large, the pan genome is large too. No matter the chosen annotation tool, this first approach suffers from producing too small core genomes, for any chosen similarity threshold, compared to what is usually expected by biologists.

TABLE I: Size of core and pan genomes w.r.t. the similarity threshold

Threshold(%)	Method 1			
	NCBI		DOGMA	
	core	pan	core	pan
50	1	163	1	118
55	<b>5</b>	<b>692</b>	2	409
60	2	1032	2	519
65	1	1454	2	685
70	0	2000	1	1116
75	0	2667	1	1781
80	0	3541	0	2730
85	0	4620	0	3945
90	0	5703	0	5181
95	0	7307	0	7302
100	0	8911	0	10132

### B. Quality test approach

Let us present our new approach. In this one, we propose to integrate a similarity distance on gene names into the pipeline.

Each similarity is computed between a name from DOGMA and a name from NCBI, as shown in Figure 2.

Genome: NC_001713.1.fasta		Threshold= 60		Mode:osneedle				
Genes in NCBI: 138		Genes in Dogma: 155						
Common Genes: 119		NCBI: 86.23%		Dogma: 76.77%				
No	Gene	Len_NC	Len_Do	N.Start	N.Stop	D.Start	D.Stop	Score
1.	RPOC2	4446	4442	ATG	TGA	TGA	AAT	99.90
2.	YCF4	546	542	ATG	TAA	TGC	TTA	99.30
3.	YCF3	540	536	ATG	TAA	ATG	TTT	99.30
4.	RPOB	4140	4136	ATG	TAA	TGA	TTT	99.90
5.	RPOA	939	935	ATG	TAA	ATG	TAA	99.60
6.	ATPH	249	245	ATG	TAA	TGG	GGT	98.40
7.	YCF47	225	221	ATG	TGA	TGC	GGA	98.20
8.	YCF46	1494	1490	ATG	TAA	ATG	TAA	99.70
9.	YCF45	1368	1364	ATG	TAA	ATG	TAA	99.70
10.	ATPG	471	467	ATG	TAA	TGG	ATT	99.20
11.	ATPF	540	536	ATG	TGA	TGG	TTA	99.30
12.	YCF41	342	338	ATG	TAA	ATG	TAA	98.80
13.	ATPD	564	560	ATG	TAA	TGA	ATT	99.30
14.	YCF24	1461	1457	ATG	TAA	TGA	GGT	99.70
15.	RPL14	366	362	ATG	TAA	ATG	TCT	98.90
16.	RPL16	414	410	ATG	TAA	ATG	AAT	99.00
17.	RPL11	426	422	ATG	TAA	ATG	AGA	99.10
18.	RPL12	384	380	ATG	TAA	ATG	AAA	99.00
19.	RPL13	420	416	ATG	TAA	ATG	TAT	99.00
20.	RPL18	408	404	ATG	TAA	ATG	ATT	99.00
21.	RPL19	363	359	ATG	TAA	TGT	CGA	98.90
22.	THIG	786	782	ATG	TAG	TGA	CAT	99.50
23.	RPS2	690	686	ATG	TAA	TGG	AAA	99.40
24.	YCF39	960	956	ATG	TAA	ATG	AAT	99.60
25.	RPL6	540	536	ATG	TAA	ATG	AAA	99.30
26.	RPL4	648	644	ATG	TGA	ATG	TGG	99.40
27.	RPL5	669	713	TTG	TAA	ATG	AGA	92.70
28.	RPL2	828	824	ATG	TAA	ATG	TTC	99.50
29.	RPL3	624	587	GTG	TAA	ATG	AAA	94.10
30.	RPL1	693	689	ATG	TAA	ATG	TCT	99.40
31.	FTSH	1935	1931	ATG	TAA	TGA	AAC	99.80

Fig. 2: Part of the implementation of the second method, comparison of the common genes from NCBI and DOGMA.

The proposed distance is the Levenshtein one, which is close to the Needleman-Wunsch, except that gap opening and extension penalties are equal. The same name is then set to sequences whose NCBI names are close according to this edit distance. The risk is now to merge genes that are different but whose names are similar (for instance, ND4 and ND4L are two different mitochondrial genes, but with similar names). To fix such a flaw, the sequence similarity, for intersected genes in a genome, is compared too in a second stage (with a Needleman-Wunsch global alignment) after selecting a genome accession number, and the genes correspondence is simply ignored if this similarity is below a predefined threshold. We call this operation, which will result in a set of quality genes, a quality test. A result from this quality test process is a set of quality genes. These genes will then constitute the quality genomes. A list of generated quality genomes based on specific threshold will construct the intersection core matrix to generate the core genes, core tree, and phylogenetic tree after choosing an appropriate outgroup.

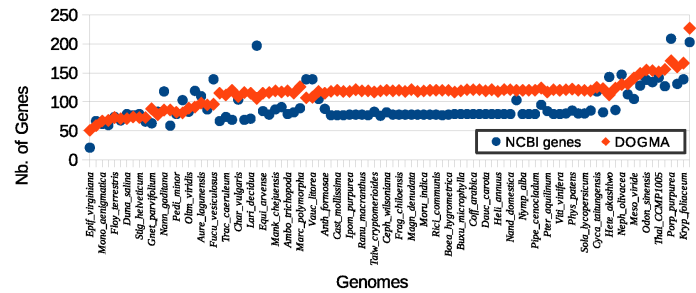
It is important to note that DNA sequence annotation raises a problem in the case of DOGMA: contrary to what happens with gene features in NCBI, genes predicted by DOGMA annotation may be fragmented in several parts. Such genes are stored in the Gene-Vision file format produced by DOGMA, as each fragment is in this file with the same gene name. A gene whose name is present at least twice in the file is either a duplicated gene or a fragmented one. Obviously, fragmented genes must be defragmented before the DNA similarity computation stage (remark that such a defragmentation has already been realized on NCBI website). As the orientation of each fragment is given in the Gene-

Vision output, this defragmentation consists in concatenating all the possible permutations (in the case of duplication), and only keeping the permutation with the best similarity score in comparisons with other sequences having the same gene name, if this score is larger than the given threshold.

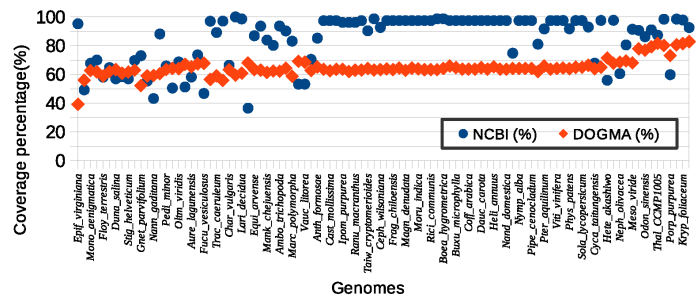
#### IV. IMPLEMENTATION

All algorithms have been implemented using Python language version 2.7, on a personal computer running Ubuntu 12.04 32bit with 6 GByte memory, and a quad-core Intel core i5 processor with an operating frequency of 2.5 GHz.

##### A. Construction of quality genomes



(a) Sizes of genomes based on NCBI and DOGMA annotations. We can see in this figure that the number of genes with DOGMA is larger than with NCBI, because the former generates more tRNAs and rRNAs genes than NCBI.



(b) Percentage of genes coverage between NCBI and DOGMA. The former outperforms the latter, as almost all genes in NCBI genomes have been covered with common genes, while most of DOGMA genes are ignored. However, correlation of them with NCBI (after quality test) is 0.6731, while it is 0.9664 with DOGMA, this latter being thus more accurate than NCBI.

Fig. 3: Original and coverage sizes between NCBI and DOGMA genomes based on threshold of 60%

To produce a core tree and genomes based on quality control approach<sup>2</sup>, we need to know what are the common genes that share almost the same name and sequence from different annotation tools. Figure 3a shows the original amount of genes based on two different annotation tools, their correlation is equal to 0.57. A two steps quality test routine is then launched to produce “quality genomes” and to enlarge the correlation: (1) select all common genes based on gene names and (2) check the similarity of sequences, which must be larger than a predefined threshold. Figure 3b presents the genes coverage percentage between NCBI and DOGMA. Remark that, gene differences between such annotation tools can affect the final

<sup>2</sup>see <http://members.femto-st.fr/christophe-guyeux/en/chloroplasts>

core genome. More precisely, The number of *tRNAs* and *rRNAs* genes are very high in the case of DOGMA annotation, while they are very low in the case of NCBI. There are also some unnamed or badly named *ORFs* genes in the case of NCBI. These genes may improve the final core genome, if their functionality are well defined.

### B. Core and pan genomes

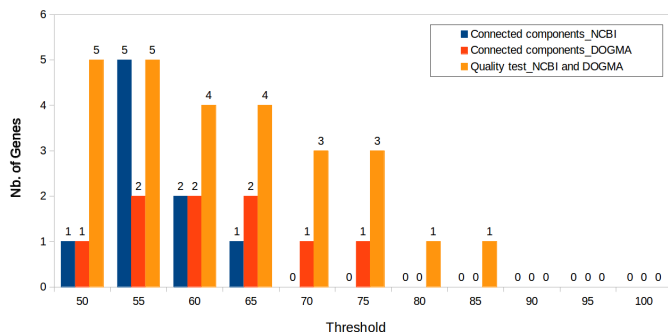


Fig. 4: Amount of core genes from each method w.r.t threshold. Note that a maximal number of core genes does not mean good core genomes: we are looking for genes meeting biological requirements.

The number of *Core genes*, illustrated in Figure 4, represents the amount of genes in the computed core genome. The main goal is to find the largest number of core genes that is compatible with biological background related to chloroplasts. From the first approach with a threshold of 60%, we have obtained 2 genes for 99 genomes with NCBI and DOGMA, whereas 4 genes for 98 genomes have been found using the second approach. In the case of second approach, we have ignored one genome for *Micromonas pusilla* under the accession (NC\_012568.1) from our sample, because we have a few amount of quality genes or none that could have been generated from its correspondents. With the second approach, zero gene in rooted core genome means that we have two or more subtrees of organisms that are completely divergent among each other. Unfortunately, for the first approach with NCBI annotation, the core genes within NCBI cores tree did not provide true biological distribution of the genomes. Conversely, in the case of DOGMA annotation, the distribution of genomes is biologically relevant. The NCBI under performance may be explained by broken subcores due to an artificially low number of genes in some genomes intersection, which could be explained by coding sequence prediction or annotation errors, or by very divergent genomes. More precisely, *Micromonas pusilla* (accession number NC\_012568.1) is the only genome who totally destroys the final core genome with NCBI annotations, for both gene features and gene quality methods.

### V. DISCUSSION

According to chloroplast endosymbiotic theory, the primary endosymbiosis has led to three chloroplast lineages among which the two most evolved groups are the chloroplastida

and the rhodophyceae. These chloroplast groups, which respectively consist of *Land plants* and *Green algae*, and *Red algae*, gave rise to secondary plastids when algae cells were engulfed by other heterotrophic eukaryotes through various secondary endosymbioses. Thus *Euglens* [6] come from *Green algae* while *Red algae* gave birth to both *Brown algae* and *Dinoflagellates*.

Now, if we observe the built core trees, in particular the one gained with quality control approach, we can notice that a primary plastid generated by the first endosymbiosis can be found in a single lineage of the chloroplast genome evolution tree: the chloroplastida group corresponds to a lineage, whereas the rhodophyceae group is represented by a second one. The generated core tree is composed by two subtrees, the first one containing the lineages of land plants and green algae and the second one presenting the lineages of brown and green algae. In the tree, some chloroplast lineages such as *Angiosperms* and *green algae* have well biological distributions, while other lineages (*Euglens*, *Dinoflagellates*, and *Ferns*) are badly distributed when compared to their biological history. Indeed, common quality genes from quality control approach are well covered by most NCBI genomes, while a large number of *tRNAs* and *rRNAs* from DOGMA genomes have been lost.

### VI. CONCLUSION

In this research work, we studied two methodologies for extracting core genes from a large set of chloroplastic genomes, and we developed Python programs to evaluate them in practice. A two stage similarity measure, on names and sequences, is thus proposed for DNA sequences clustering in genes, which merges best results provided by NCBI and DOGMA. Results obtained with this “quality control test” are deeply compared with our previous research work, on both computational and biological aspects, considering a set of 99 chloroplastic genomes. Core trees have finally been generated for each method, to investigate the distribution of chloroplasts and core genomes. The tree from DOGMA annotation has revealed the best distribution of chloroplasts regarding their evolutionary history. In particular, it appears to us that each endosymbiosis event is well branched in the DOGMA core tree.

### REFERENCES

- [1] B. Alkindy, J. F. Couchot, C. Guyeux, A. Mouly, M. Salomon, J. M. Bahi, "Finding the Core-Genes of Chloroplasts", *Journal of Bioscience, Biochemistry, and Bioinformatics*, IACSIT press, 4(5):357–364, 2014.
- [2] B. Stoebe, W. Martin, K. V. Kowallik, "Distribution and nomenclature of protein-coding genes in 12 sequenced chloroplast genomes". *Plant Molecular Biology Reporter*, Springer 16(3):243–255, 1998.
- [3] D. Grzebyk, O. Schofield, C. Vetriciani, P. G. Falkowski, THE MESOZOIC RADIATION OF EUKARYOTIC ALGAE: THE PORTABLE PLASTID HYPOTHESIS1. *Journal of Phycology*, 39(2):259–267, 2013.
- [4] Stacia K. Wyman, Robert K. Jansen, Jeffrey L. Boore, "Automatic annotation of organellar genomes with DOGMA", *BIOINFORMATICS*, oxford Press, 20(172004):3252–3255, 2004.
- [5] P. Rice, I. Longden, and A. Bleasby. "EMBOSS: the European Molecular Biology Open Software Suite", *Journal Trends in Genetics*, 16(6):276–277, 2000.
- [6] Geoffrey I. McFadden, "Primary and secondary endosymbiosis and the origin of plastids", *Journal of Phycology*, Wiley Online Library 37(6):951–959, 2001.