

TRANSLATION REGULATION OF INTRON-CONTAINING GENES IN CHLOROPLASTS

ALEXANDER SELIVERSTOV* and VASSILY LYUBETSKY†

*Institute for Information Transmission Problems
Russian Academy of Sciences, B. Karetny
19, Moscow 127994, Russia
*slstv@iitp.ru
†lyubetsk@iitp.ru*

Received 18 October 2005
Revised 14 February 2006
Accepted 16 February 2006

The paper provides a short description of the originally developed algorithm for searching of the conservative protein–RNA binding sites. The algorithm is applied to analyze chloroplast genes. The candidate protein–RNA binding sites were detected upstream of *atpF*, *petB*, *clpP*, *psaA*, *psbA*, and *psbB* genes in many chloroplasts of algae and plants. We suggest that some of these sites are involved in suppressing translation until splicing is completed.

Keywords: Translation; chloroplasts; multiple alignment.

1. Introduction

The gene expression in chloroplasts of algae and plants is regulated by binding of nuclear-encoded proteins to the chloroplast mRNA.¹ These proteins are involved in editing, translation and maintaining stability of the chloroplast mRNA. The detailed analysis of regulatory sites was published for alga *Chlamydomonas reinhardtii* and some plants.^{1–3} For example, protein binding to the *psbA* 5′-untranslated region in *C. reinhardtii* leads to activation of translation.³

Here, we briefly describe an algorithm for identification of protein–RNA binding sites. The algorithm was used to identify regulatory sites in 5′-untranslated regions of protein-coding chloroplast genes.

Many chloroplast protein-coding genes contain introns. Thus, their translation should not start immediately after transcription. However, the translation machinery of chloroplasts closely resembles that of bacteria, particularly, in the ribosome being able to follow the RNA-polymerase on the mRNA strand. If the ribosome arrives at the end of an exon before splicing starts, it will preclude splicing. To avoid this, in some rare cases the AUG start codon is derived from ACG by editing of mRNA, which prevents translation from starting immediately.² RNA editing is

known for chloroplasts of higher plants and is absent, e.g. in liverwort *Marchantia polymorpha*.

Our algorithm detects candidate protein–RNA binding sites upstream of *atpF*, *petB*, *clpP*, *psaA*, *psbA*, and *psbB* genes in many chloroplasts.

We suggest that some of these sites are involved in suppressing translation until splicing is completed. This conjecture is also consistent with the presence of conserved sites upstream of these genes in multiple alignments (given below) and is supported by experiment in published evidence.³

2. Materials

Chloroplast genomes were obtained from GenBank (NCBI). The dataset contained 5'-untranslated intergenic regions from chloroplast genomes of algae and plants: *Euglena gracilis*, *Odontella sinensis*, *Guillardia theta*, *Cyanidioschyzon merolae*, *Cyanidium caldarium*, *Gracilaria tenuistipitata*, *Porphyra purpurea*, *Chlamydomonas reinhardtii*, *Nephroselmis olivacea*, *Chaetosphaeridium globosum*, *Mesostigma viride*, *Anthoceros formosae*, *Marchantia polymorpha*, *Adiantum capillus-veneris*, *Huperzia lucidula*, *Psilotum nudum*, *Pinus thunbergii*, *Amborella trichopoda*, *Arabidopsis thaliana*, *Atropa belladonna*, *Calycanthus floridus*, *Cucumis sativus*, *Epifagus virginiana*, *Lotus corniculatus*, *Nicotiana tabacum*, *Nymphaea alba*, *Oryza nivara*, *Oryza sativa*, *Panax ginseng*, *Spinacia oleracea*, *Triticum aestivum*, *Zea mays*.

The only nonphotosynthetic species in this group is *E. virginiana*. In *P. nudum*, the *psbB* gene is annotated *psbT*.

3. Methods

Consider a dataset of leader regions upstream of orthologous genes and a corresponding species tree. A set of shallow phylogenetic subtrees (groups) is selected in the species tree. For each of these groups, the algorithm searches for conserved regions of fixed length n (that can be varied) *via* finding cliques in a suitable multipartite graph. The basic idea is as follows. The algorithm finds clusters of very similar sites, called signals or motifs, of a fixed length n for each of these phylogenetic groups. A cluster of sites generates a weight matrix $4 \times n$, where the k th column of the matrix, $1 \leq k \leq n$, contains letter frequencies in the k th site position from the cluster. Further, the algorithm generates clusters of weight matrices for different suitable n . The clusters of matrices are generated accounting for distances between the ancestors of the initial phylogenetic groups in the species tree. The algorithm of finding cliques also constructs these clusters of matrices. In each matrix cluster, the matrices are replaced by corresponding site clusters and thus obtained the clusters of clusters are then combined. The described procedure can be iterated. The algorithm is described in detail in Ref. 4.

4. Notation

In the following sections, R denotes A or G, Y denotes C or U, W denotes A or U, S denotes C or G, D denotes A, G or U, N denotes any nucleotide and * denotes any nucleotide or gap.

5. Results

In many chloroplasts, the algorithm found long conserved binding sites containing conserved helices upstream of the genes *atpF*, *petB*, *clpP*, *psaA*, *psbA*, and *psbB*. A summary of the results is given in Table 1.

Table 1. Occurrence of introns and predicted sites upstream of chloroplast genes *atpF*, *petB*, *clpP*, *psaA*, *psbA*, and *psbB*. Notation: “+” — candidate protein binding site present; “-” — no candidate binding site; “s” — introns present; “n” — no gene homolog in the species; “&” — helices in the site; “E” — editing of the start codon.

Filum	Species	<i>atpF</i>	<i>clpP</i>	<i>petB</i>	<i>psaA</i>	<i>psbA</i>	<i>psbB</i>
Euglenozoa	<i>Euglena gracilis</i>	-s	-	-s	-s	-s	-s
Bacillariophyta	<i>Odontella sinensis</i>	-	-	-	+&	+&	-
Cryptophyta	<i>Guillardia theta</i>	-	-	-	+&	+&	-
Rhodophyta	<i>Cyanidioschyzon merolae</i>	-	-	-	-	+&	-
	<i>Cyanidium caldarium</i>	-	-	-	-	-	-
	<i>Porphyra purpurea</i>	-	-	-	+&	+&	+
	<i>Gracilaria tenuistipitata</i>	-	-	-	-	+&	-
Chlorophyta	<i>Chlamydomonas reinhardtii</i>	-	-	-	-s	+&s	-
	<i>Nephroselmis olivacea</i>	-	-	-	+&	+&	+
Streptophyta	<i>Chaetosphaeridium globosum</i>	-	+&s	-s	+&	+&	+
	<i>Mesostigma viride</i>	-	-	-	+&	-	-
Anthocerophyta	<i>Anthoceros formosae</i>	+s	+&s	+&s	+&	+&	+
Hepatophyta	<i>Marchantia polymorpha</i>	+s	+&s	+&s	+&	+&	+
Lycopodiophyta	<i>Huperzia lucidula</i>	+s	+&s	+&s	+&	+&	+
Pteridophyta	<i>Adiantum capillus-veneris</i>	+sE	+&s	-sE	+&	+&	+
Psilophyta	<i>Psilotum nudum</i>	+s	+&s	+&s	+&	+&	+
Pinophyta	<i>Pinus thunbergii</i>	+s	+&	+&s	+&	+&s	+
Magnoliophyta (eudicotyledons)	<i>Amborella trichopoda</i>	+s	+&s	+&s	+&	-	+
	<i>Arabidopsis thaliana</i>	+s	+&s	+&s	+&	+&	+
	<i>Atropa belladonna</i>	+s	+&s	+&s	+&	+&	+
	<i>Calycanthus floridus</i>	+s	+&s	+&s	+&	+&	+
	<i>Cucumis sativus</i>	+s	+&s	+&s	+&	+&	+
	<i>Epifagus virginiana</i>	n	+&s	n	n	n	n
	<i>Lotus corniculatus</i>	+s	+&s	+&s	+&	+&	+
	<i>Nicotiana tabacum</i>	+s	+&s	+&s	+&	+&	+
	<i>Nymphaea alba</i>	+s	+&s	+&s	+&	+&	+
	<i>Panax ginseng</i>	+s	+&s	+&s	+&	+&	+
	<i>Spinacia oleracea</i>	+s	+&s	+&s	+&	+&	+
Magnoliophyta (Liliopsida)	<i>Oryza nivara</i>	+s	+&s	+&s	+&	+&	+
	<i>Oryza sativa</i>	+s	+&s	+&s	+&	+&	+
	<i>Triticum aestivum</i>	+s	+&s	+&s	+&	+&	+
	<i>Zea mays</i>	+s	+&s	+&s	+&	+&	+

5.4. Photosystem I

In many chloroplasts, the 5'-untranslated regions of *psaA* gene (photosystem I P700 apoprotein A1) are found to possess a conservative motif adjacent upstream of the start codon AUG with the consensus

WGUURGYRRGUYUYUYU*UAUN*****
 NUYGUCYGRARAGAGGAGRA*CUCR.

The motif contains a conservative helix near the putative ribosome-binding site. The corresponding alignment is shown in Fig. 4.

5.5. Photosystem II

In many chloroplasts, the 5'-untranslated regions of the *psbA* gene (photosystem II protein D1) and the *psbB* (photosystem II P680 chlorophyll A apoprotein) gene contain conserved motifs adjacent to the AUG start codon with the consensus

YUUGGGARYYYY*****NAAACYAAG

and

AAAGUNACRUAGU*GUCUAYUUN*****NNAAGGGGURUUU,

respectively. The first motif, upstream of *psbA* contains a conserved helix. These motifs were not found in 5'-untranslated regions of *M. viride* genes, nor in other algae.

<p><i>Odomella sinensis</i> <i>Guillardia theta</i> <i>Porphyra purpurea</i> <i>Anthoceros olivacea</i> <i>Chaetosphaeridium globosum</i> <i>Mesostigma viride</i> <i>Anthoceros formosae</i> <i>Marchantia polymorpha</i> <i>Huperzia lucidula</i> <i>Adiantum capillus-veneris</i> <i>Psilotum nudum</i> <i>Pinus thumbergii</i> <i>Amborella trichopoda</i> <i>Arabidopsis thaliana</i> <i>Atropa belladonna</i> <i>Calycanthus floridus</i> <i>Cucumis sativus</i> <i>Lotus corniculatus</i> <i>Nicotiana tabacum</i> <i>Nymphaea alba</i> <i>Panax ginseng</i> <i>Spinacia oleracea</i> <i>Oryza nivara</i> <i>Oryza sativa</i> <i>Triticum aestivum</i> <i>Zea mays</i> Consensus</p>	<p>cuaaugagaguuucau*aaaau*****UUCgucUCCaaaaGGAGAAgucua auaaagaaagaguuuuagauu*****gcugUCUcaaaagagGAGAccuca uagaaaaaagcguuuu*gaau***ccuugUCUcaagagagGAGAAucuca agccaggaagacuaauu*cauu***CCUCgugugagaGAGGagaucucg uguuuuaagauuuuucuuagc***CUCgUCUgaaaAGAGGAGaaucug uagaggugaguuuuuu*ugug***cCUcaUCUaaaaAGAGGAGaaucucc uguuggcgguuuuuc*caug***CCUCgucugaagaGAGGauaauucg uguugguagguuuuuc*uaug***CCUCgucugaagaGAGGagaaccucg ucuuggcggguuuuuuc*uaug***CCUCgucuggaaaGAGGagaaccucg uguugguagguuuuuc*uauc***CCUCgUCGaaGAGAGGagaguccca ugcuggcagguuuuuc*uaau***CCUCgucucgagaGAGGagaucuca uaauuggcagguuuuucuaauuuuuuagucccgUCCgaaaagaGGAgaa*uuca ucuuggcgggucucuuguaug***uguugUCCggaagaGGAgga*cuca uguuggcggguuuuucuuuguaug***uguugUCCggaagaGGAgga*cuca uguuggcgggucucuuguaug***uguugUCCggaagaGGAgga*cuca uguuggcgggucucuuguaug***uguugUCCggaagaGGAgga*cuca uguuggcgggucucuuguaug***uguugUCCggaagaGGAgga*cuca uguuggcgggucucuuguaug***uguugUCCggaagaGGAgga*cuca uguuggcgggucucuuguaug***uguugUCCggaagaGGAgga*cuca uguuggcgggucucuuguaug***uguugUCCggaagaGGAgga*cuaa uguuggcgggucucuuguaug***uguugUCCggaagaGGAgga*cuaa uguuggcgggucucuuguaug***uguugUCCggaagaGGAgga*cuaa uguuggcgggucucuuguaug***uguugUCCggaagaGGAgga*cuaa WGUURGYRRGUYUYUYU*UAUN*****NUYGUCYGRARAGAGGAGRA*CUCR</p>
---	--

Fig. 4. Alignment of 5'-untranslated regions upstream of the *psaA* start codon. Helices are shown in capitals.

J. Bioinform. Comput. Biol. 2006.04:783-792. Downloaded from www.worldscientific.com
 by INSTITUTE FOR INFORMATION TRANSMISSION OF THE RAS on 10/21/12. For personal use only.

6. Discussion

Conserved motifs detected upstream of the *atpF*, *petB*, *clpP*, *psaA*, *psbA*, and *psbB* genes are likely to be involved in the regulation of translation.

The conserved region upstream of *atpF* contains an AG-rich motif typical for the ribosome-binding sites. This conservative region is considerably longer than typical binding sites. The presence of introns in the genes suggests that translation starts only after completion of splicing.

These *petB* genes do not have a typical ribosome-binding site but contain a conserved helix that might suggest posttranscriptional modification of the 5'-untranslated regions or binding of a transcription activator. In all plants, the *petB* gene contains introns.

The conserved regions in the 5'-untranslated regions of *clpP* and *psbA* genes were observed upstream of almost all their orthologs, even those lacking introns. The translational regulation of the *psbA* gene has been experimentally observed in *C. reinhardtii*, where transcription is constitutive, but translation is activated at light by a 47 kDa protein, forming a complex with other proteins and mRNA, but not interacting with mRNA directly.³ The complex is inactivated in the dark.

The conserved nature of this region in plants and algae might suggest that the translation regulation machinery for the gene *psbA* preceded the emergence of introns.

Notably, the conserved RNA motifs in the transcripts of *clpP*, *petB*, *psaA* and *psbA* contain helices with conserved flanks likely interacting with protein mediator, which is typical for most regulatory systems.^{5,6}

The long conserved motifs were found upstream of the *psaA* and *psbB* genes, which lack introns in all species containing the motifs. On the other hand, in chloroplasts of *A. capillus-veneris*, all studied 5'-untranslated regions are considerably diverged. Thus, the motif was not found upstream of *petB*, while it was upstream of the other five cases. In the latter situation, however, site trees and species trees differed considerably in the node containing the name of the corresponding species.

Other intron-containing genes in the studied chloroplast genomes were not found to have conserved 5'-motifs, or their 5'-untranslated regions were too short or absent. Two such examples are discussed below.

In studied plants, the upstream regions of the *rbcL* genes encoding a ribulose 1,5-bisphosphate carboxylase/oxygenase subunit contain only a short conserved motif with the consensus ARGGAGGGACYT whose core is the ribosome-binding site. And we have no reason to assign a regulatory role to this motif, as the plant *rbcL* genes lack introns. On the other hand, the *rbcL* genes contain introns in chloroplasts of both algae *E. gracilis* and *C. reinhardtii*, and in the latter case, *rbcL* is regulated by mRNA-binding proteins.³ This seeming discrepancy is not surprising, since in both algae the structure of 5'-untranslated region is completely different from that in studied plants.

A different situation is observed in the case of *ycf3* genes (photosystem I assembly protein Ycf3). It contains introns and a long 5'-untranslated region neither overlapping with other genes in plant chloroplasts, nor was it found to possess conserved motifs.

Acknowledgments

The authors are grateful to M.S. Gelfand for valuable discussion and to L.Y. Rusin for discussion and help. This study was partially supported by ISTC 2766.

References

1. Nickelsen J, Chloroplast RNA binding proteins, *Current Genet* **43**:392–399, 2003.
2. Zerges W, Translation in chloroplasts, *Biochimie* **82**:583–601, 2000.
3. Hauser CR, Gillham NW, Boynton JE, Translation regulation of chloroplast genes, *J Biol Chem* **271**:1486–1497, 1996.
4. Lyubetsky VA, Seliverstov AV, Note on cliques and alignments, *Inf Processes* **4**:241–246, 2004.
5. Vitreschak AG, *et al.*, Riboswitches: the oldest mechanism for the regulation of gene expression? *Trends Genet* **20**:44–50, 2004.
6. Seliverstov AV, *et al.*, Comparative analysis of RNA regulatory elements of amino acid metabolism genes in Actinobacteria, *BMC Microbiol* **5**:54, 2005.



Vassily Lyubetsky obtained his M.Sc., Ph.D. and D.Sc. habilitation degrees in Math from Moscow State University and Russian Academy of Sciences, Russian Federation, in 1968, 1971, and 1991, respectively. His fields of academic expertise are mathematical logic, algebra, and theoretical computer science. Since 1995 he is a faculty professor of Moscow State University at the chair of mathematical logic and algorithmic theory. For over 10 years he is occupying permanent positions in Russian Academy of Sciences as a head of the Laboratory for mathematic methods and models in bioinformatics of the Institute for Information Transmission Problems, and as a principal researcher at the Federal State Enterprise, Institute of Strategic Stability, of the Ministry for Industry and Science of Russian Federation.

He has authored and co-authored over 150 scientific publications, four monographs and over 50 international conference publications.



Alexander Seliverstov received his Diploma in Applied Mathematics from State Classic Academy, Moscow, Russian Federation, in 1999. He is currently in the Laboratory for mathematic methods and models in bioinformatics of the Institute for Information Transmission Problems, Russian Academy of Science.