

Поиск консервативных участков в лидерных областях генов в случае известного дерева видов¹

В.А. Любецкий, А.В. Селиверстов

Институт проблем передачи информации РАН,
127994, Россия, Москва,
Большой Каретный переулок, 19,
e-mail: lyubetsk@iitp.ru, slvstv@iitp.ru
Поступила в редколлегию 22.08.2005

Аннотация—Предложен алгоритм поиска консервативных нетранслируемых участков в мРНК с учетом дерева видов. В частности, найдены сигналы у генов, имеющих интроны и транскрибируемых в хлоропластах растений из Embryophyta. У многих хлоропластов гены белков имеют интроны. Поэтому трансляция не может происходить сразу после транскрипции, как это имеет место у бактерий. В тоже время аппарат трансляции у хлоропластов очень близок к таковому у бактерий. В статье описан потенциальный сигнал, который мог бы обеспечивать такую задержку начала трансляции гена с интронами.

1. ВВЕДЕНИЕ

У многих хлоропластов гены белков имеют интроны, или белок транслируется с мРНК, образованной транс-сплайсингом из разных транскриптов. Поэтому трансляция не может происходить сразу после транскрипции, и, соответственно, рибосома не может двигаться вдоль мРНК сразу после РНК-полимеразы, как это имеет место у бактерий. В тоже время аппарат трансляции у хлоропластов очень близок к таковому у бактерий (не требуется экпирование 5' конца мРНК, как у эукариотов, строение рибосом близкое, трансляция мРНК ингибируется хлорамфениколом, как и у бактерий, и т.д.).

В редких случаях время, необходимое для сплайсинга, могло бы обеспечиваться за счет редактирования мРНК: например, за счет преобразования кодона ACG в старт кодон AUG, [1]. Известно, что у *Chlamydomonas reinhardtii* трансляция многих генов (без интронов) регулируется белками, связывающими 5' нетранслируемые области мРНК, [2], [3], [4].

В [5] высказано предположение, что начало трансляции гена с интронами некоторым механизмом задерживается до завершения сплайсинга. В нашей заметке описан потенциальный сигнал, который мог бы обеспечивать такую задержку начала трансляции гена с интронами. В заметке также предлагается модификация алгоритма, описанного в [6] и [7]. Эта модификация тестировалась и показала свою эффективность, в частности, на примере поиска сигналов, т.е. консервативных 5' нетранслируемых участков в мРНК, у генов *atpF* и *petB*, имеющих интроны и транскрибируемых в хлоропластах растений из Embryophyta. Эти растения перечислены в таблице 1.

¹ Работа частично поддержана грантом МНТЦ 2766.

Таксон	Вид	Рибосомальные белки
Embryophyta	<i>Anthoceros formosae</i>	L2 L16 S12 S16
	<i>Marchantia polymorpha</i> (печеночник)	L16 S12
- Tracheophyta	<i>Adiantum capillus-veneris</i> (папоротник)	L2 L16 S12 S16
	<i>Huperzia lucidula</i> (плаун)	L2 L16 S12 S16
	<i>Psilotum nudum</i> (папоротник)	L2 L16 S12
- - Spermatophyta	<i>Pinus koraiensis</i>	L2 L16 S12
	<i>Pinus thunbergii</i> (чёрная сосна)	L2 L16 S12
- - - Magnoliophyta	<i>Amborella trichopoda</i>	L2 L16 S12 S16
	<i>Arabidopsis thaliana</i>	L2 L16 S12 S16
	<i>Atropa belladonna</i>	L2 L16 S12 S16
	<i>Calycanthus floridus</i>	L2 L16 S12 S16
	<i>Cucumis sativus</i> (огурец)	L2 S12 S16
	<i>Epifagus virginiana</i>	L2 L12 L16 S12
	<i>Lotus corniculatus</i> (ледвянец рогатый)	L2 L16 S12 S16
	<i>Nicotiana tabacum</i>	L2 L16 S12 S16
	<i>Nymphaea alba</i> (кувшинка белая)	L2 L16 S12 S16
	<i>Panax ginseng</i> (женьшень обыкновенный)	L2 L16 S12 S16
	<i>Spinacia oleracea</i> (шпинат)	L16 S12 S16
	<i>Oryza nivara</i> (рис)	L2 L16 S12 S16
	<i>Oryza sativa</i> (рис)	L2 L16 S12 S16
	<i>Triticum aestivum</i> (пшеница)	L2 L16 S12 S16
	<i>Zea mays</i> (кукуруза)	L2 L16 S12 S16

Таблица 1. Рибосомальные белки хлоропластов растений из Embryophyta, гены которых подлежат сплайсингу. Знак “-” указывает на принадлежность к подтаксону предыдущего таксона.

2. МАТЕРИАЛЫ И МЕТОДЫ

Геномы хлоропластов взяты из ГенБанка.

В качестве набора исходных последовательностей брались лидерные области перед ортогологичными генами. В общем случае можно искать сигнал – консервативные участки в таком наборе с помощью алгоритма поиска клики в многодольном графе, который описан в [6] и [7]. Но когда исходные последовательности принадлежат геномам, для которых известно дерево видов, этот алгоритм может быть удачно усовершенствован, как это описано ниже. При тестировании так усовершенствованного алгоритма был получен результат, показанный в таблицах 2 и 3.

Исходный набор последовательностей в алфавите {A, U, C, G} разбивается на группы, состоящие из близкородственных по дереву видов последовательностей, и в каждой группе по отдельности ищется сигнал – набор попарно близких сайтов фиксированной длины n (которая возрастала до тех пор, пока сигнал не теряется в заметной части последовательностей из такой группы). Таким образом, для каждой филогенетической группы находится один или несколько сигналов. Каждый сигнал определяет свою весовую матрицу – матрицу размера $4 \times n$, в которой указаны частоты встречаемости каждой буквы в каждом столбце сигнала. Затем отбираются цепочки таких матриц, у которых имеются попарно близкие подматрицы одного (по возможности, большого) размера $4 \times m$, и эта цепочка определяет набор последовательностей уже из разных филогенетических групп и соответствующий этому набору общий сигнал. Этот шаг алгоритма выполняется с помощью того же исходного алгоритма поиска клики в многодольном графе.

Наконец, выполняется сборка односайтовых сигналов в многосайтовые, которые уже, вообще говоря, разделены неконсервативными участками. И затем строится множественное выравнивание, согласованное с найденными сигналами и с исходным деревом видов.

3. РЕЗУЛЬТАТЫ ТЕСТИРОВАНИЯ

3.1. АТФ синтаза

Ген *atpF*, кодирующий одну из субъединиц АТФ синтазы, имеет интроны в генах хлоропластов из Embryophyta (кроме *Epifagus virginiana* – не фотосинтезирующего паразита). У *Adiantum capillus-veneris* старт кодон AUG получается из кодона ACG при редактировании мРНК. В остальных случаях старт кодоном является RUG, где R – один из нуклеотидов А или G. В лидерных областях гена *atpF* у видов из Embryophyta найдены по два консервативных участка, консенсусы которых соответственно равны AAUNAAAAA и UAUCUAUAAGAGGAGANNA, второй участок расположен правее и вблизи старта трансляции гена, он перекрывает предполагаемый сайт связывания рибосомы, таблица 2. Здесь N означает произвольный нуклеотид. Между этими консервативными участками расположены AU богатые неконсервативные участки с длинами от 19 до 61 нуклеотидов, неконсервативные даже в Magnoliophyta. В таблице 2 показано соответствующее множественное выравнивание 5' лидерных областей, непосредственно примыкающих к старт кодону гена. Вероятно, некоторый белок *предотвращает* инициацию трансляции, перекрывая область связывания рибосомы, см. дискуссию.

<i>Anthoceros formosa</i>	ggaaugaauaaugscuuaagcuuuccuauuuuaaguaggauuggau
<i>Marchantia polymorpha</i>	uaaaugaaaaaacugugaguaaacaacacucca*****
<i>Adiantum capillus-veneris</i>	uuaauaaauaaugugucaaauucugcuaaggcuggauc*****
<i>Hyperzia lucidula</i>	auaaggaaaaaacuaguaaacuu*****
<i>Psilotum nudum</i>	uuaauaaaagaaaaacuuguc*****
<i>Pinus thunbergii</i>	caaaauaagaaaaaacaacaaaaaucugua*****
<i>Pinus koraiensis</i>	aggaucaaaaaagaaaagaacaaaaaugaug*****
<i>Amborella trichopoda</i>	uaaaauaaaaauaagauauggggugaagugaucaaaaaaga*****
<i>Arabidopsis thaliana</i>	agaauaaaaaaaaggacagaguuccu*****
<i>Atropa belladonna</i>	auaaauaaaagagggggcgaagugcucaaaaaaga*****
<i>Calycanthus floridus</i>	ucaaucaaaaaaaggggggcgaagugaucaaaaaagaac****
<i>Cucumis sativus</i>	aaaauaaaaaaaauagaaagaauaga*****
<i>Lotus corniculatus</i>	caaaaaaaauaaaagaaaaaucuacaaaaauagga*****
<i>Nicotiana tabacum</i>	agaauaaaauaaaagagggggcgaagugcucaaaaaaga*****
<i>Nymphaea alba</i>	gugaucaaaaaaggaauucuuuuuuuguauuuug*****
<i>Panax ginseng</i>	uaauuaaaaaaagaacaggggcgaaguuaucaaaaaagaac****
<i>Spinacia oleracea</i>	ucaauaaaaaaauaaaauucuuuugaaguagcaaac auu****
<i>Oryza nivara</i>	gcauuaaaaaacsgaucaaaaagggcgagcgaaguaagugau***
<i>Oryza sativa</i>	gcauuaaaaaacsgaucaaaaagggcgagcgaaguaaaguga***
<i>Triticum aestivum</i>	ccgaucaaaaagggcgagcgaaguaagugaucgaaa*****
<i>Zea mays</i>	ccgaucaaaaagggcgagcgaaguaaguuaucaaaaaagga*****
консенсус	**AAUNAAAAA*****

Таблица 2. Выравнивание 5' лидерных областей перед иницирующим кодоном гена *atpF* у видов из Embryophyta. Здесь N обозначает произвольный нуклеотид, а * нуклеотид или делецию.

<i>Anthoceros formosa</i>	aaggaagagacaуacuaagacuаааgааccuaugaugggagagagagu
<i>Marchantia polymorpha</i>	*****uааuuuuсaаuaаuааааacgааааааagaggacagc****
<i>Adiantum capillus-veneris</i>	*****gаааuugccсааааacгуааааасуucgaggagggааааgааu*
<i>Huperzia lucidula</i>	*****ggаuааuааассугuааugggagааааgаu*
<i>Psilotum nudum</i>	*****aаuаgаuаgucаuuаugggagagguauu
<i>Pinus thunbergii</i>	*****gаасаuаucсуuаucuaugaggggagagcgu**
<i>Pinus koraiensis</i>	*****uаgаасаuаucсуuаucuaugaggggcgagcau**
<i>Amborella trichopoda</i>	****acucсguuugguuуuuаgucсуаucуаgааgaggagagau**
<i>Arabidopsis thaliana</i>	*****uuuuuаuаguuuаgсuаgааgaggagauuаu**
<i>Atropa belladonna</i>	*****acucугucгуuсgаuuuuuuаgucуаucуаuааgaggagaucаu**
<i>Calycanthus floridus</i>	**cugcgcauuuuguuаgccсуаucуаuucуаuааgaggааagcau**
<i>Cucumis sativus</i>	*****uaаuuаguuuuаucуаuаааagggаuсаu**
<i>Lotus corniculatus</i>	*****аucуаuаааgаgааuucгуuuуаucсуаuаaggagagaucаu**
<i>Nicotiana tabacum</i>	*****acucугucгуuсgаuuuuuuаgucуаucуаuааgaggagaucаu**
<i>Nymphaea alba</i>	*****uuаgucсуаuucсуаucсуаааgаuаaggagagagcau**
<i>Panax ginseng</i>	**ucугucгуuuuuuuuuuuuuuаgucуаucуаuааaggagagaucаu**
<i>Spinacia oleracea</i>	**gаааuааuасаасgаuuuuuuuguuuаucуаuааaggagagaucаu**
<i>Oryza nivara</i>	*cgааааасуuugucсуuuугucгucсуаucуаuааaggagagagcau**
<i>Oryza sativa</i>	ucgааааасуuugucсуuuугucгucсуаucуаuааaggagagagcau**
<i>Triticum aestivum</i>	*****aacуuugucсуuuугucгucсуаucуаuааaggagagagcau**
<i>Zea mays</i>	*****aacуuucуucсуuuугucгucсуаucуаuааaggagagagcau**
консенсус	*****UаUCUаUаAGAGGаGаNNа***

Продолжение таблицы 2.

3.2. Цитохром b6

В 5' лидерной области гена *petB* (цитохром b6) хлоропластов у всех видов из Embryophyta, кроме *Adiantum capillus-veneris* и *Epifagus virginiana*, найден консервативный участок, примыкающий к старт кодону AUG с консенсусом равным GGTAGTTCGAYCGYGGAATT*YTTT**GTTTNGTATTTYYGGAAT, таблица 3. Здесь Y означает один из нуклеотидов C или U, и N произвольный нуклеотид, а * возможная вставка не более чем одного нуклеотида. Здесь обнаружена консервативная спираль, показанная подчеркиванием в множественном выравнивании в таблице 3.

Перед геном *petB* на расстоянии 3–10 нуклеотидов от старта трансляции этого гена нет высоко консервативного AG-богатого участка, характерного для сайта связывания рибосомы. Найденный нами консервативный участок имеет большую длину и перед ним расположен слабо консервативный участок. Поэтому можно думать, что наш консервативный участок перекрывает сайт связывания рибосомы. И тогда некоторый белок, связываясь с этим консервативным участком в 5' лидерной области гена *petB* с учетом вторичной структуры, в отличие от предыдущего случая активирует инициацию трансляции.

У *Epifagus virginiana* ортолог гена *petB* отсутствует. У *Adiantum capillus-veneris* соответствующая 5' лидерная область не выравнивается с другими видами из Embryophyta. У всех видов из Embryophyta ген *petB* содержит интроны.

<i>Anthoceros formosa</i>	***uuuucccagug*gugguaguuu <u>aaucgugcaacua</u> cugaaaaaaaaggauuuuugaaau
<i>Marchantia polymorpha</i>	ucauuuuuuuaauuuu*agguaguuu <u>aaauuguguaauua</u> *uuaa**auucaaggauuu*uugaaau
<i>Huperzia lucidula</i>	auugauc <u>ccuuccuuu</u> *gguaguuu <u>aaucguguaauu</u> *cuga***aucaaaggauuuuagaaau
<i>Psilotum nudum</i>	*ucauaaaaaaagac*gaggcagu <u>ugaucacgca</u> aaauuuuu***auuuauugauguuuugaaau
<i>Pinus thunbergii</i>	*agcuuauc <u>uuguuc</u> **cacuagu <u>uuugaucguguaauua</u> cuuuu**cucaaaggauuuuuggaau
<i>Pinus koraiensis</i>	*agcuuauc <u>uuguuc</u> **caauagu <u>uuugaucguguaauua</u> cuuuu**cucaaaggauuuuuggaau
<i>Amborella trichopoda</i>	uuggguuu <u>cuagguu</u> *a*ggguagu <u>ucgaccgugca</u> auuuccuuu***guuucgguauuuuccggaau
<i>Arabidopsis thaliana</i>	***ccuauuc <u>uccuu</u> **ugguagu <u>ucgaccgca</u> aaauuuuuuuuugcguuugaauuuuuccggaau
<i>Atropa belladonna</i>	*cauuc <u>uaauuu</u> cuuuu*ugguagu <u>ucgaucgugga</u> auuuu <u>cuuu</u> ***guuucgguauuuuccggaau
<i>Calycanthus floridus</i>	*uuuuc <u>uagcc</u> cauuc*ugguagu <u>ucgaccgugga</u> auuuccguu***guuucgguauuuuccggaau
<i>Cucumis sativus</i>	*uuagcc <u>uacuc</u> uuuuuugguagu <u>ucgaucgugga</u> auuuuauuu***uuucgguauuuuccggaau
<i>Lotus corniculatus</i>	*cauuc <u>cuuuuu</u> uuuu*ugguagu <u>ucgaucgugga</u> acuuu <u>cuuu</u> ***guuucgguauuuuccggaau
<i>Nicotiana tabacum</i>	*cauuc <u>uaauuu</u> cuuuu*ugguagu <u>ucgaucgugga</u> auuuu <u>cuuu</u> ***guuucgguauuuuccggaau
<i>Nymphaea alba</i>	ucauc <u>ucauuc</u> cuguu**ugguagu <u>ucgaccgca</u> aaauuu <u>cuuu</u> ***guuucgguauuuuccggaau
<i>Panax ginseng</i>	*cagcc <u>cauuc</u> uaauuu*ugguagu <u>ucgaccgca</u> aaauuu <u>cuuu</u> ***guuucgguauuuuccggaau
<i>Spinacia oleracea</i>	*uaauuu <u>cauucc</u> uuu*ugguagu <u>ucgaucgugga</u> auuuu <u>cuuu</u> ***cuuucgguauuuuccggaau
<i>Oryza nivara</i>	*cauuu <u>cuagaca</u> uuc*ugguagu <u>ucgaccgugga</u> auu <u>uuuug</u> **guuucgguauuc <u>ucugga</u> au
<i>Oryza sativa</i>	*cauuu <u>cuagaca</u> uuc*ugguagu <u>ucgaccgugga</u> auu <u>uuuug</u> **guuucgguauuc <u>ucugga</u> au
<i>Triticum aestivum</i>	*cauuu <u>cuagau</u> uuu*augguagu <u>ucgaccgca</u> aaauuuuuuu***guuucgguauuc <u>ucugga</u> au
<i>Zea mays</i>	acauuu <u>cuagaca</u> uuc*ugguagu <u>ucgaccgugga</u> auu <u>uuuu</u> ***guuuugguauuc <u>ucugga</u> au
консенсус	****YY*****UYU*UGGUAGUUCGAYCGYGGAUU*YUUU***GUUUNNGUAUUUYYGAAU

Таблица 3. Выравнивание 5' лидерных областей перед инициирующим кодоном гена *petB* у видов из Embryophyta. Здесь N обозначает произвольный нуклеотид, а * нуклеотид или делецию.

4. ДИСКУССИЯ

Найденные сигналы – консервативные участки перед генами *atpF* и *petB*, вероятно, связаны с задержкой инициации трансляции до завершения сплайсинга. Сейчас нет достаточных оснований, чтобы решить, основана ли эта регуляция на связывании белка с мРНК. Перед геном *petB* консервативный участок не содержит очевидного сайта связывания рибосомы, но имеется консервативная шпилька, что позволяет предположить: после транскрипции происходит образование вторичной структуры в 5' лидерной области гена и туда присоединяется белок, активирующий инициацию трансляции.

Интроны обнаружены в генах, кодирующих различные белки, экспрессируемые в хлоропластах как растений, так и некоторых одноклеточных эукариотов (*Chlamydomonas reinhardtii*, *Euglena gracilis*). Часто они встречаются в генах белков фотосистем, некоторых субъединиц РНК-полимеразы и NADH-дегидрогеназы, и в генах рибосомальных белков. Количество генов с интронами в хлоропластах примерно одинаковы у всех растений, за исключением *Pinus* spp. и *Epifagus virginiana*, где таких генов меньше.

Хлоропласты у всех видов из Embryophyta, а также у *Chaetosphaeridium globosum* и *Euglena gracilis* содержат интроны в генах белков большой и малой субъединиц рибосомы, таблица 1. Можно думать, что интроны в хлоропластах играют своеобразную регуляторную роль. Если рибосом достаточно много, то рибосомы, транслирующие начальный фрагмент транскрибированной РНК, предотвращают или изменяют сплайсинг, препятствуя образованию новых рибосом. Если рибосом мало, то мРНК длительное время остается открытой для “правильного” сплайсинга и, возможно, редактирования.

СПИСОК ЛИТЕРАТУРЫ

1. Sugiura M., Hirose T., Sugita M. Evolution and mechanism of translation in chloroplasts, *Annu Rev Genet*, 1998, 32, стр. 437–59.
2. Hauser C.R., Gillham N.W., Boynton J.E. Translation regulation of chloroplast genes, *The Journal of Biological Chemistry*, 1996, 271: 3, стр. 1486–1497.
3. Nickelsen J. Chloroplast RNA binding proteins, *Curr Genet*, 2003, 43, стр. 392–399.
4. Herrin D.L., Nickelsen J. Chloroplast RNA processing and stability, *Photosynthesis Research*, 2004, 82, стр. 301–314.
5. Zegers W. Translation in chloroplasts, *Biochimie*, 2000, 82, стр. 583–601.
6. Любецкий В.А., Селиверстов А.В. Некоторые алгоритмы, связанные с конечными группами, *Информационные процессы*, 2003, 3, 1, стр. 39–46.
7. Lyubetsky V.A., Seliverstov A.V. Note on Cliques and Alignments, *Информационные процессы*, 2004, 4, 3, стр. 241–246.