

计算识别水稻基因组中microRNA 基因^{*}

金伟波^{a, b}, 吴方丽^a, 孔 栋^a, 郭蔼光^{a, b}, 杨淑慎^a

(西北农林科技大学 a 生命科学学院, b 陕西省农业分子生物学重点实验室, 陕西杨凌 712100)

[摘要] 一些低丰度的m RNA 和组织特异性m RNA 往往很难发现, 利用生物信息学方法, 根据已知水稻m RNA 的各种属性, 设计m RNA 前体预测程序——PreM iRFind, 从全基因组范围内预测出1 375 条m RNA 前体, 然后利用已知拟南芥的m RNA 对这1 375 条前体进行同源检索, 找出166 条m RNA 前体, 再用水稻已知的m RNA 对这166 条序列进行检索, 发现其中153 条与已知的水稻microRNA 完全相同, 对剩余的13 条m RNA 前体预测序列用m Fold 作进一步结构分析, 最后得到10 条新的候选m RNA 基因序列。

[关键词] 水稻基因组; RNA fold; microRNA

[中图分类号] S511; Q522

[文献标识码] A

[文章编号] 1671-9387(2006)12-0097-04

MicroRNA (m RNA) 是存在于动植物体内的起调节 mRNA 稳定性及翻译作用的一类 ncRNA^[1-5]。最早发现的m RNA lin-4 和let-7 可通过与靶m RNA 3'末端形成碱基配对来抑制翻译, 很多新发现的m RNA 也以相似的机制发挥作用^[6]。不过m RNA 大家族里也有一些成员是通过RNAi途径使靶m RNA 降解达到基因阻抑目的, 类似于具有相似长度的 siRNA 的功能。

随着人们对m RNA 认识的加深, 更多的人开始相信, m RNA 是早期RNA 家族在进化上的遗留物, 其之所以没有被更为先进的蛋白质调控机制所代替, 可能主要是因为m RNA 的表达比通常的蛋白质编码基因转录物的表达要迅速得多, 并且m RNA 作为调节子不会受到翻译过程的影响, 所以有利于对靶基因进行更加迅速和有效的调控^[7]。m RNA 的发现, 丰富了人们对蛋白质表达调控的认识, 同时也是对中心法则中RNA 作用的重要补充, 将促使人们对生物发育调控进行更深入的探索。

近10年来, 有关m RNA 的研究工作已经取得了突飞猛进的进展。到目前为止, 在Rfam 中登陆的m RNA 已有1 000 多条。但Bentwich 等^[8]认为, 生物体内的m RNA 远大于此数, 他们通过计算并结合microarray 技术, 从人基因组中又发现了89 条新的m RNA, 自此在人类基因组中总共发现321 条m RNA, 但他们认为人类基因组编码的m RNA 至

少有800 条。目前在Rfam 中登陆的水稻microRNA 有123 条, 而其中通过实验确证的却只有44 条。由此推断, 在水稻基因组中也应存在大量尚未被发现的m RNA。本研究采用生物信息方法, 找出了10 条新的水稻m RNA, 现将研究结果报道如下。

1 材料与方法

1.1 材料

本研究所用的水稻基因组和EST 数据从NCBI 库(<http://www.ncbi.nlm.nih.gov/>) 中下载得到; 133 条水稻m RNA 和76 条拟南芥m RNA 均来自于Rfam 数据库(<http://microRNA.sanger.ac.uk/sequences/index.shtml>), 用Vienna 软件包中的RNAfold 程式预测二级结构; RNAstructure 为m fold 的windows 版, 下载于生物软件网(<http://www.bio-soft.net/>)。

1.2 方法

1.2.1 m RNA 属性的选取 为了区别m RNA 与非m RNA, 本研究选取了m RNA 的前体结构、平均长度、发夹结构中茎的长度、GC 含量、环长及与拟南芥的相似度等几个属性, 用于后面的预测。

1.2.2 水稻m RNA 前体的预测 根据上述属性, 对已知的133 条水稻m RNA 进行统计分析, 以得到的结果为参数, 设计出用于探寻水稻前体m RNA 的程序——PreM iRFind, 其工作流程见图1。

* [收稿日期] 2005-11-26

[基金项目] 国家转基因研究与产业化开发专项(JY03A-11-01); 陕西省农业分子生物学重点实验室项目

[作者简介] 金伟波(1977-), 男, 浙江温岭人, 在读博士, 主要从事植物分子生物学研究。

[通讯作者] 郭蔼光(1943-), 女, 陕西西安人, 教授, 博士生导师, 主要从事植物分子生物学研究。E-mail: guoaguang@yahoo.com.cn

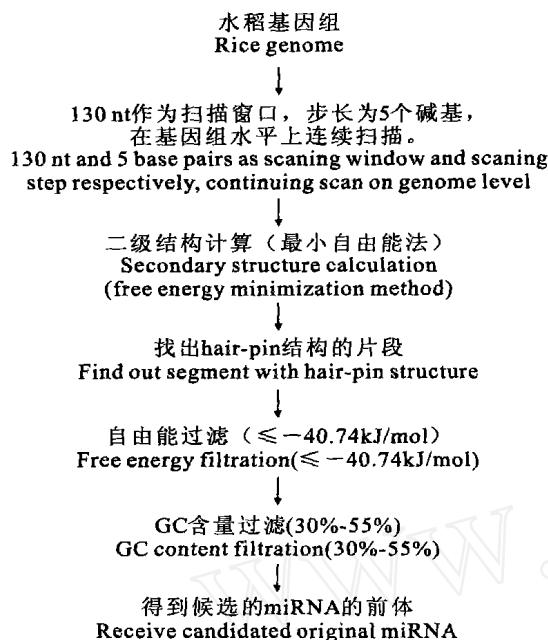


图1 PreM iRFind 工作流程图

Fig. 1 Flow chart of PreM iRFind

1.2.3 水稻m RNA 的预测 以拟南芥已知的m RNA 为quest, 对上述获得的前体m RNA 进行blast 检索, 所获片段被认为是最后预测的m RNA, 如检索到的序列有重复, 只取其前体形成二级结构自由能最小的序列, 其他的作为冗余片段丢弃, 最后再去除已发表的水稻m RNA, 剩下的就是新的m RNA。

2 结果与分析

2.1 水稻m RNA 的特征

对已知的133条水稻m RNA 进行分析, 结果

表1 预测出的水稻m RNA 的序列及其在基因组上的位置

Table 1 Sequences of new m RNAs and their locations in the rice chromosome

m RNA 编号 m RNA number	染色体 Chromosome	前体片段 起始位置 position of original sequence	链位置 Chain location	候选的成熟m RNA 序列 Candidated maturity m RNA sequence	同源的拟南芥 m RNA Originated m RNA of <i>A. thaliana</i>
m iR01	1	38139509~ 38139603	-	U GCCAAA GGA GAUUU GCA UAC	ath-M iR 399
m iR02	1	5718080~ 5718148	+	U GU CAU CUU CAU CAU CAU CU G	ath-M iR 414
m iR03	2	12956021~ 12956100	+	UU GGA GAA GA GA GU GA GCA CA	ath-M iR 156
m iR04	3	30345921~ 0346041	-	A GCU GCCA GCAU GAU CUA ACU	ath-M iR 167
m iR05	4	32655912~ 32656028	+	UU GU AU AA GU GAA GU GU UU G	ath-M iR 395
m iR06	4	31838759~ 31838870	-	A GU GU UU GGGGAA CU CU CGA	ath-M iR 395
m iR07	5	2832917~ 2833037	+	U CCAAAGGGAU CGCAAU GU CU	ath-M iR 393
m iR08	7	28469789~ 28469879	-	GAA GCA GGGCA CGU GCAU GCA	ath-M iR 164
m iR09	8	1701243~ 1701325	+	CU CCCU GU AU GCCA CU CAU CU	ath-M iR 160
M iR10	12	25510435~ 25510546	+	U GCCA GCAU GAU CUA GCU CU G	ath-M iR 167

注: + . 有意义链; - . 无意义链。

Note: + . Sense strand; - . Anti-sense strand

发现水稻前体m RNA 的平均长度约为131 nt, 所有前体均具有发夹结构, 发夹环长度大于4 nt, GC 含量为30% ~ 55%; 成熟的m RNA 由20或21 nt 个碱基组成。

2.2 水稻前体m RNA 的搜索

以130 nt 作为窗口长度, 5 nt 作为步长, 从水稻每条染色体的5 端开始滑动窗口, 调用RNA fold 预测子片段的二级结构(zuker), 判断其是否有发夹结构, 没有发夹结构的丢弃, 有发夹结构的, 若环长> 4 nt, 茎长 20 nt, 并且GC 含量为36% ~ 80%, 则保存此片段作为候选前体片段, 不符合的丢弃; 然后向下滑动5 nt, 重复上述工作。经过这一步后, 得到约1 万条序列, 去除冗余后还剩余1 375 条候选前体片段。

2.3 成熟水稻m RNA 的预测

为了进一步在m RNA 前体上预测成熟片段, 本研究用已知的拟南芥成熟m RNA 作为quest, 对上述预测出的1 375 条候选片段进行blast 分析, 要求匹配长度大于17 nt, 结果得到166 条符合条件的候选片段。然后再用已知的水稻m RNA 对这166 条片断进行检索, 以去除已发表的m RNA, 最后得到13 条m RNA 片断。又用RNA structure 对这13 条序列进行重折叠, 去除3 条结构不符片段, 最终得到了10 条新的候选m RNA 前体, 其二级结构见图2, 序列及其前体在染色体上的位置见表1。再将这10 条m RNA 递交到水稻EST 库进行blast 检索验证, 其中plus/plus 匹配的有m iR02, m iR05 和m iR07 3 条。

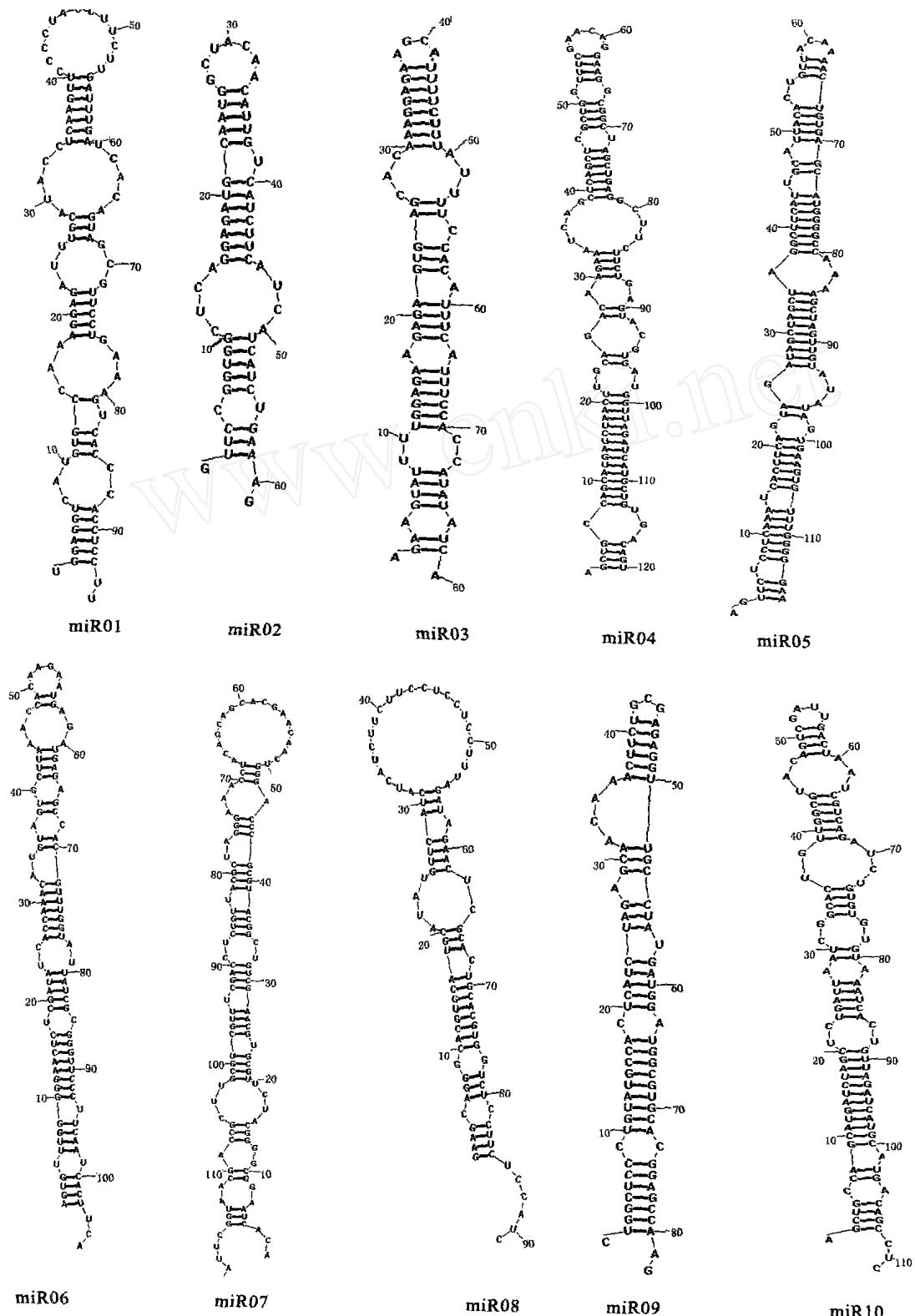


图2 预测出的水稻m iRNA 前体的二级结构

Fig. 2 Putative secondary structures of selected m iRNA precursors

3 讨 论

m iRNA 在基因调控中非常重要^[9]。Lew is 等研究认为, 人类有 1/3 的基因由 m iRNA 所调控。因

此, 尽快找出所有的 m iRNA 并研究其功能, 对进一步了解基因表达调控具有重要意义。但由于其片段较短小, 在实验水平快速识别 m iRNA 具有很大困

难^[11]。因为目前在实验水平上检测microRNA,主要是分离18~28 nt的小片段RNA,然后再通过克隆和测序的手段来获得,但由于实验本身的问题,研究者们克隆到的mRNA仅是表达丰度较高的少数mRNA,而大批的低丰度mRNA却很难通过实验手段分离到。因此,利用生物信息学的方法来预测mRNA意义深远。通过计算的方法预测mRNA能在短时间内识别出大量的mRNA,但同时也会产生大量的假阳性序列,因此如何提高识别准确率是mRNA预测中亟待解决的问题,也是生物信息学领域中普遍存在的问题。

本研究采用全基因组扫描结合同源检索来识别mRNA的方法是一种高效、准确的方法。在本研究中总共识别出166条候选片段,而在用已知的水稻mRNA进行筛选后只剩下13条,说明其中包含了大量的已知水稻mRNA,这也从另一个侧面反映了本研究方法的识别率较高。最后将得到的10条microRNA对水稻EST库进行blast检索,发现有3条是Plus/Plus匹配,说明这3条microRNA在水稻生长阶段肯定能转录出来。但其他7条mRNA是否为真正的mRNA,还需要进一步的研究。

[参考文献]

- [1] Lagos-Quintana M, Rauhut R, Lendeckel W, et al Identification of novel genes coding for small expressed RNAs[J]. Science, 2001, 294(5543): 853-858
- [2] Lau N C, Lai M P, Weinstein E G, et al An abundant class of tiny RNA s with probable regulatory roles in *Caenorhabditis elegans*[J]. Science, 2001, 294(5543): 858-862
- [3] Lee R C, Ambros R, et al An extensive class of small RNAs in *Caenorhabditis elegans*[J]. Science, 2001, 294(5543): 862-864
- [4] Reinhart B J, Weinstein E G, Rhoades M, Bartel B, et al MicroRNAs in plants[J]. Genes Dev, 2002, 16(13): 1616-1626
- [5] Llave C, Kasschau K D, Rector M A, et al Endogenous and silencing associated small RNAs in plants[J]. Plant Cell, 2002, 14(7): 1605-1619
- [6] Bartel D P. MicroRNAs: genomes, biogenesis, mechanism, and function[J]. Cell, 2004(116): 281-297
- [7] Ruvkun G. Glimpses of a tiny RNA world[J]. Science, 2001, 294(5543): 797-799
- [8] Isaac Bentwich, Amir Avniel, Yael Karov, et al Identification of hundreds of conserved and nonconserved human microRNAs[J]. Nature Genetics, 2005(37): 766 -770
- [9] Bonnet E, Wuyts J, Rouze P, et al Detection of 91 potential conserved plant microRNAs in *A rabidopsis thaliana* and *Oryza sativa* identifies important target genes[J]. Proc Natl Acad Sci USA, 2004(101): 11511-11516
- [10] Lewis B P, Burge C B, Bartel D P. Conserved seed pairing, often flanked by adenines, indicates that thousands of human genes are microRNA targets[J]. Cell, 2005, 120: 15-20
- [11] Wang X J, Reyes J L, Chua N H, et al Prediction and identification of *A rabidopsis thaliana* microRNAs and their mRNA targets[J]. Genome Biol, 2004, 5: 65.

Computational identification of novel family members of microRNA genes in *Oryza sativa* genome

JIN Wei-bo^{a,b}, WU Fang-li^a, KONG Dong^a, GUO Ai-guang^{a,b}, YANG Shu-shen^a

(*a* College of Life Science, *b* Key Laboratory of Agriculture Moleculer Biology in Shaanxi Northwest A & F University, Yangling, Shaanxi 712100, China)

Abstract: MicroRNAs (mRNA) are about 21 nt long non-coding RNAs derived from larger hairpin precursors and play important regulatory roles in both animals and plants. The low abundance of some mRNA s and their time- and tissue-specific expression patterns make experimental mRNA identification difficult. Here a program - PreM iRFind for genome-wide prediction of *Oryza* microRNAs was presented. This method used characteristic features of known *Oryza* mRNA s as criteria to search for mRNA s. Our prediction identified 1 375 pre-mRNA candidates. The 166 pre-mRNA s were reservation as criteria to search for mRNA s conserved between *A rabidopsis* and *Oryza sativa*. 153 pre-mRNA s were the same as known mRNA s and then the structures of the 13 pre-mRNA s were analyzed further. Finally 10 new mRNA s were obtained.

Key words: rice genome; RNA fold; microRNA