

基于比较基因组学的玉米 ESTs 定位方法

张祖新^{1,2}, 张绍鹏^{1,2}, 郑用琰²

(1. 河北农业大学农学院, 保定 071001; 2. 华中农业大学作物遗传改良国家重点实验室, 武汉 430070)

摘要:描述了以水稻基因组数据和玉米与水稻的比较遗传图谱为桥梁, 基于水稻和玉米间存在的标记和序列水平上的广泛的共线性, 对大量的玉米 ESTs 初步定位于玉米连锁群上新方法, 为对 ESTs 开展进一步的基因组学研究和基因克隆提供参考信息。对 139 条玉米 ESTs 的定位发现, 96 条玉米 ESTs (69%) 可在水稻基因组中找到同源序列, 77 条 ESTs (55%) 可使用该策略进行定位, 证实了该方法的可行性和有效性。

关键词:表达序列标签 (ESTs); 比较作图; 基因组; 分子标记

中图分类号: Q943

文献标识码: A

文章编号: 0253-9772(2006)03-0339-06

A Strategy Based on Comparative Genomics to Align ESTs of Maize

ZHANG Zu-Xin^{1,2}, ZHANG Shao-Peng^{1,2}, ZHENG Yong-Lian²

(1. College of Agronomy, Hebei Agricultural University, Baoding 071001, China;

2. National Key Laboratory of Crop Genetic Improvement, Huazhong Agricultural University, Wuhan 430070, China)

Abstract: In this study, a new strategy to locate ESTs on maize linkage groups was described. In the strategy, the rice (*Oryza sativa* L.) genomic sequence database of was employed to locate maize ESTs on rice linkage groups, and then to locate on maize linkage group by comparative genetics mapping between rice and maize genome. The aligned ESTs information should available for further study on genomics and gene cloning. As an example, 139 ESTs of maize were assayed, and 96 maize ESTs (69%) were homologous with rice genomic sequence, 55% (77/139) ESTs were located on maize linkage groups based on the strategy, indicating that the locating approach of ESTs is feasible and available.

Key words: expression sequence tags (ESTs); comparative mapping; genome; molecular marker

大规模的 ESTs 测序计划和功能基因组学的研究使玉米 ESTs 数据日益积累, 如何建立这些 ESTs 与特定生理过程或性状表现的联系已是当前急需解决的问题。BAC 筛选结合 FISH 杂交是定位 ESTs 的有效策略, 而基于遗传转化的基因功能获得或敲除是了解 ESTs 功能的最直接方法, 这对于研究少量的 ESTs 也是可行的。比较基因组学为大规模的 ESTs 研究提供了新的策略, 如基于 ESTs 功能咨询的方法可以将一些 ESTs 归于某一功能类群, 而多

种 *in silico* 的方法可以将 ESTs 归入可能的功能网络^[1]。虽然这些方法仍是对 ESTs 功能的推测, 同时也无法建立 ESTs 与性状表现的直接联系, 但为 ESTs 的研究提供了有益的参考信息。

对分子标记图谱的比较研究证实, 玉米、高粱和甘蔗的基因组间具有大量的共线性, 且高粱和甘蔗的基因组的组成更为相似, 高粱和甘蔗的许多基因组区域与玉米中的 2 个不同区域同源^[2~5]。Ahn 等^[6]研究表明, 在 RFLP 标记水平上, 水稻和玉米近

收稿日期: 2005-07-28; 修回日期: 2005-11-21

基金项目: 农业部“引进国际先进农业科学技术”项目(编号: 2003-Q03-1-22) [Granted by the Introducing International Advanced Agriculture Science and Technology Project of the Ministry of Agriculture (No. 2003-Q03-1-22)]

作者简介: 张祖新(1964—), 男, 博士, 教授, 研究方向: 遗传学。Tel: 0312-7528408; E-mail: nxzxx@mail.hebau.edu.cn

2/3 的基因组具有共线性。禾本科植物间,控制重要农艺性状的 QTL 也具有共线性^[7]。基于分子标记的比较基因组学研究在宏观水平上反应了基因组间的相似性。而对 $\alpha 1\text{-sh2}$ 基因区域的 DNA 序列分析表明,玉米、高粱和水稻均含这两个基因,且转录方向相同,表现出序列水平上的微共线性;但在玉米中,这 2 个基因被 140 kb 的片段所分开,而在高粱和水稻中,这 2 个基因间仅有 19 kb 的片段^[8,9]。同时外显子序列的同源性比内含子序列的同源性要高得多^[10]。玉米和水稻基因组间的同源性为基于水稻的基因组信息来研究玉米基因组研究提供了理论依据。

玉米和水稻的亲缘关系较近,并且作为模式植物的水稻,其基因组测序计划已经完成^[11,12],同时,随着分子标记的不断开发和利用,玉米和水稻的分子标记图谱进一步精细。Yuan 等(2003)使用 13 251 个水稻标记,将 2 464 个水稻 BAC/PAC 序列排列在水稻遗传图谱上,进而在水稻 BAC/PACs 中搜索 1 259 个已定位的玉米标记,在高严谨条件下,350 个玉米标记被排列于水稻连锁群上,且这些资料可供进一步的研究使用 (www.tigr.org/tdb/e2k1/osa1/maize/description.shtml)^[13]。因此,利用水稻基因组序列数据和玉米与水稻间的比较图谱,构建一套以水稻基因组为桥梁的比较基因组学研究方法,将玉米 ESTs 定位于玉米染色体上成为可能。

1 实验原理和方法

1.1 网站及工具软件

实验所借助的基因组数据库网站和软件主要为:(1)玉米基因组数据库 www.maizeGDB.org;(2)禾本科植物基因组数据库 www.gramene.org,它提供了水稻基因组数据、水稻与玉米基因组的比较图谱数据及相关软件;(3)美国国家生物技术信息中心 GenBank 数据库 www.ncbi.nlm.nih.gov,该数据库中基于 NCBI/Blast 软件的核酸序列同源性比对分析,可对功能未知的 ESTs 进行功能预测。

1.2 实验基本过程

1.2.1 实验原理

许多研究证实,玉米和水稻间在标记水平上具有广泛的共线性,同时也具有序列水平(特别是 ESTs 水平)上的微共线性,因而,玉米 ESTs 则可在 Gramene 网站的水稻基因组数据库中搜索到同源的水稻 ESTs 和基因组 DNA 序列。由于水稻基因组序

列已经完成,与玉米同源的 DNA 序列则可找到其所在的染色体位置、BAC 以及与该序列紧密连锁的分子标记等。通过玉米和水稻的比较遗传图谱以及玉米与水稻间共有的标记,进而可以将这些标记重新定位到玉米的连锁群上。由于水稻 DNA 序列与这些分子标记连锁,玉米 ESTs 与水稻 DNA 同源,因而,通过水稻基因组数据库作为桥梁,就可以将玉米 ESTs 确定在玉米某一连锁群的片段上(图 1)。

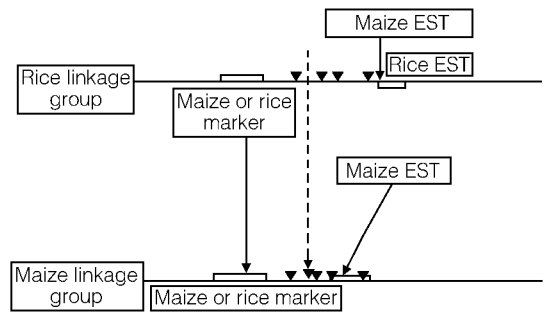


图 1 基于比较基因组学的玉米 ESTs 定位的示意图

Fig. 1 An illustration on the process for maize ESTs location on maize linkage groups based on comparative genomics strategy

1.2.2 操作步骤

以下面 1 条 EST 序列为例说明实验具体的操作步骤:

(1) 同源序列咨询。进入 Gramene 网站的 Blast 页面,将该序列贴入查询框,各项参数选择默认值,进行 Blast 咨询。Results Summary 表格显示水稻数据库中一条与咨询序列同源的水稻序列(chr3:28 545 877~28 545 939 bp),称为匹配序列(当有多条匹配序列时,按照同源程度从高到低在 Sequence 项中排序,选择 Score 值高、E-value 值低的序列)。Alignment 提供了两条序列具体的比对情况,点击匹配序列链接,进入到该序列的详细注释页面。

(2) 展示水稻序列的注释。即点击 Chromosome、Overview、Detailed View 和 Basepair View,逐步获得更为详细的信息。在 Detailed View 的 Zoom 梯状图中选择 200 kb,表示显示以匹配序列为中心的 200 kb 范围内的序列;Detailed View 随后显示对该序列的注释信息,包括 Rice markers、Rice EST、Maize markers、Maize EST 等。Markers 项中的标记有 C944、C51151S、R1538 和 R1770 等,选择与匹配序列连锁最近的 C51151S,点击进入 C51151S 标记的信息界面。

(3) 标记在水稻遗传图谱中位置的确定。选择 Map Positions 中 Rice-GR TIGR Assm IRGSP Seq

```

TCCGACCAGA ACTCTCTAGA GTGGATGCC GTTTTCCTCG TCATACTGCT GCTGGGCGGT
CTGCAGCACC CGACCATCGC CGCCGGTCTG GCGGCATCT ACATCGTTGC GAGGTTCTTC
TACTTCAAGG GATACGCCAC CGGC GTTCCG GACAACCGTC TCAAGATTGG GGGGCTCAAC
TACTTGCGGT TGCTGGGGCT GATCATTTC ACAGCATCTT TCGGCATCAA CTGCTCATC
AGGGAGATGC TCTAAACTTT GGGTTTATAT GTGAGCTCCA GGCTGGGTTA ACCGCACAAG
TTACTCATGG TCTGAGCATC AAAACTTTGT GATACTAGTC GTTTGGGGTT TTATAGTGG
GTTATCGTTT CTTTCTTTGT CTGTGTGTTG TTAGCAAAG AAGTGTGCGT CTGCGTCTTA
TAGCTGTGAT GGAATAAGCA GACAGATGGG TTGAAGCGTG AATCGTACAG

```

图 2 一条有待定位的玉米 EST 序列

Fig.2 A Located EST sequence of maize

2005, 点击 view on map (有时一个标记会在一个图谱中有多处定位, 点一个即可), 显示 C51151S 在所选项谱中的位置。同时通过 Options Menu 中的 Feature Types 调节图谱中所要显示的标记。

(4) 比较图谱的绘制。在 Comparative map 的下拉菜单中选择玉米遗传图谱 IBM2 neighbors 2004 的所有染色体, 再点击 Redraw map, 显示水稻 chr3 和玉米染色体的对应关系比较图(图 2)。由于共线性过

于密集, 可以返回上一界面, 通过 Comparison Menu 中的 Map Start 和 Map End 缩小显示范围, 以便看到更清晰共线性。图中得到在水稻上的 *umc1431* 和 *csu554b* 分别在玉米的第 1 和第 5 染色体上找到共线标记 *umc107a* 和 *csu554b*。而与玉米 EST 同源的水稻 DNA 序列紧密连锁的标记 C51151S 与玉米第 5 染色体上的 *csu554b* 具有共线性, 因而推测该 EST 可能在玉米第 5 染色体上的 *csu554b* 标记附近区域(图 3)。

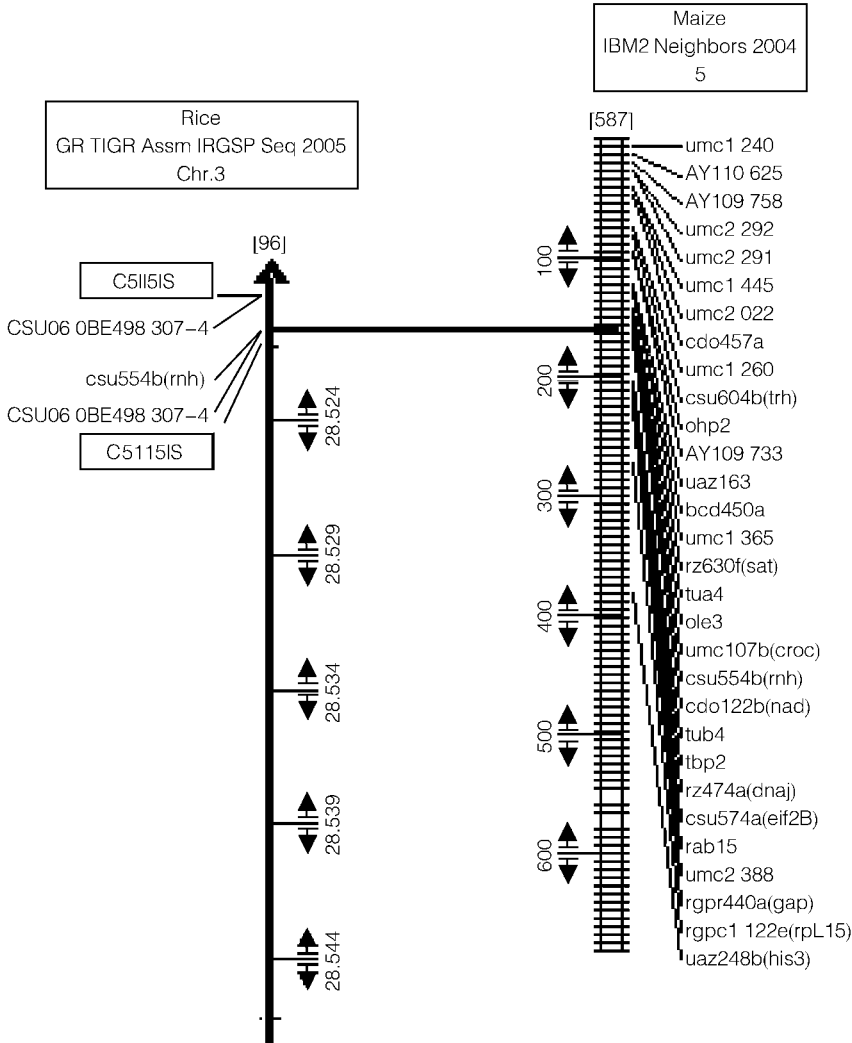


图 3 基于比较图谱的玉米 ESTs 定位的实例

Fig.3 An example to locate maize ESTs by comparative mapping

2 实例

以 139 条在玉米 S-(*Rf3*)和 S-(*rf3*)花粉中差异表达的 ESTs 为例^[14],来说明基于水稻基因组数据和玉米与水稻的比较基因组遗传图谱的玉米 ESTs 定位策略。

利用玉米 ESTs 与水稻基因组序列进行 BLASTn,在 E-value $<10^{-5}$ 、配对长度至少 50 bp 的条件下,96 条(约 69%)玉米 ESTs 在水稻染色体上找到了同源序列,而有 43 条 ESTs(约 31%)在水稻染色体上找不到匹配序列。如果降低同源配对的严谨度,E-value $<10^{-5}$ 同时配对长度大于 30 bp 的条件下,玉米 ESTs 在水稻染色体上找到的同源序列则会大大增加,81% 以上的玉米 ESTs 与水稻序列同源。19% 的 ESTs 未找到同源水稻序列。从总体上来看,玉米和水稻在 ESTs 序列水平上确实存在着广泛的同源性和进化上的保守性。尽管如此,与水稻基因组比较,在标记水平上玉米基因组存在大量的倒位和重排。

在 96 条与水稻序列具有同源性的 ESTs 中,77 条 ESTs(占 ESTs 总数的 55%)在其同源序列两侧各 100 kb 区间内,存在与之紧密连锁的共线性标记,在玉米的第 1 至第 10 条染色体上的分布分别为

17、7、17、8、17、10、8、15、5 和 8 条(表 1),这些标记在玉米连锁群上的分布如图 3,其中玉米第 1、3、5 和 8 染色体上所定位的 ESTs 较多,而第 2 和第 9 染色体上所定位的 ESTs 较少。由于所有的 ESTs 是来源于 *Rf3/rf3* 的差异所引起的表达水平的差异,已知 *Rf3/rf3* 位于玉米第 2 染色体的长臂 2.09 bin^[15,16],但第 2 染色体上却只定位了 7 条 ESTs,其中 5 条位 2L 上的 2.07~2.08 bin,主要是与信号传导有关的 Calmodulin-binding protein 和 protein phosphatase 2C。类似的,第 1 染色体上的 ESTs 主要分布在 1.05 bin,1.07~1.08 bin,第 8 染色体上主要分布在 8.01~8.02 bin,第 9 染色体上分布在 9.02 bin 等(图 4)。基于生物信息学的功能分析发现,这些 ESTs 可能涉及到代谢、细胞结构、信号传导、转录与翻译、物质转运及细胞凋亡等生理生化过程^[14]。这些说明,玉米基因组结构上,功能相似的基因也存在成簇排列并协同表达趋势。另外,有 19 条玉米 ESTs 虽然在水稻基因组序列中有同源序列,由于玉米连锁遗传图谱上的标记密度不如水稻,在水稻基因组 200 Kb 的区域内所找到的分子标记却与玉米连锁群上的标记没有共线性,因此,无法将这些 ESTs 通过比较图谱进行定位。

表 1 定位于玉米 10 条染色体上的 ESTs 数及其与 ESTs 连锁的标记

Table 1 The number of located ESTs on 10 maize linkage groups and markers link with ESTs

染色体号 Chr.	定位的 ESTs 数 Number of mapped ESTs	与 ESTs 连锁的标记 Linked markers with ESTs
1	17	<i>csu503, rs2, tub1, AY109506, rgpc250, umc1383, rz403, gln2, uaz151, csu663b, umc2239, umc1571, kn1, ids1, rgpc1122c, csu675a, csu256, tua1, umc1571, AY110296, umc1245, umc107a, umc1571, umc1571, rs2, umc1446, umc1579</i>
2	7	<i>rz474c, rgpc1122a, csu154a, umc1555, umc1464, mmp33</i>
3	15	<i>umc1968, abp1, rz382a, umc1973, AW258116, rz390c, umc1608, rgpc385b, umc1300, csu899a, umc1347, umc2275, rz527a, AY105849, umc1504, umc1608</i>
4	8	<i>csu324a, umc2360, prh1, gpc1, rz900a, az145, uaz228c, rz143b, rgp1102, rpd3, umc2360, csu554b, rz474a, ppp1, ivr2, csu663a, AY105029, rgpc1122b, gpc440b, csh1c, csu587b, umc1692, xet1, uaz190, ua4, AY109606, mc1447, umc2294, su774, bcd1072, az219, mc2386</i>
5	17	<i>umc1250, mir1, rz390d, dh1, gpc2, rz143a, mc1996, umc2162, umc1314, csu360, uaz220, umc1760, mdh6, rgpc1122b, rgpr440b, ij1, umc1456, asg8, AY109703, csu847b, umc1125, AY103821, csu179c, csu675b, AY110056, rz390a, rgpc131b, AY106269, AY109853, su368</i>
8	6	<i>umc1728, RZ382, csu891, AY110056, AY105457, AY106269</i>
9	5	<i>rgpr3235a, wx1, umc1636, rz2a, prc1, rz593</i>
10	8	<i>rgpc1122d, rgpr440c, umc1115, tip5, rz900c, por2, AY109698, gdcp1</i>

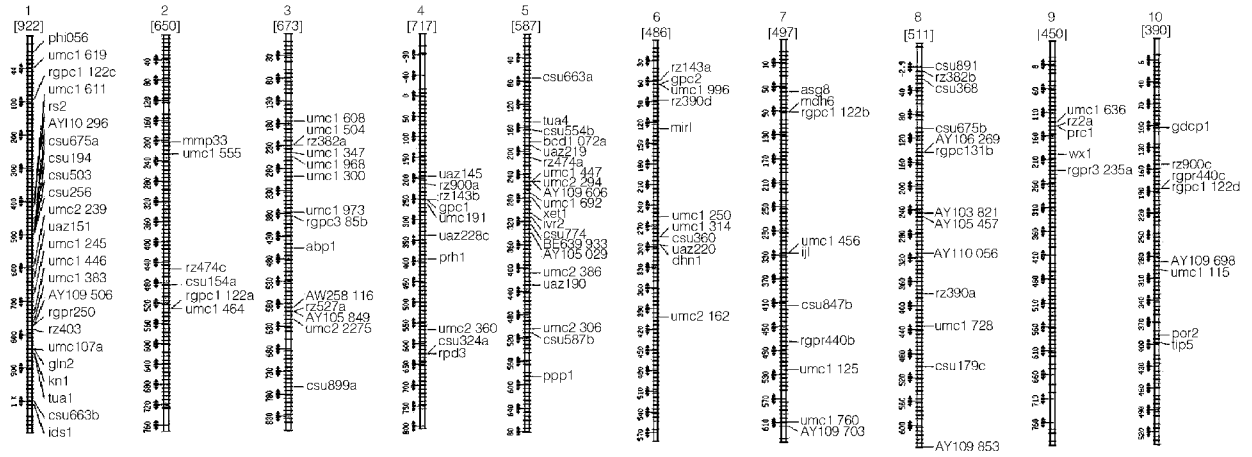


图 4 与玉米 ESTs 连锁的标记在玉米连锁群上的分布

Fig.4 The distribution of marker linked with maize ESTs on maize linkage groups

3 讨论

比较基因组研究主要是利用相同的 DNA 分子标记在相关物种之间进行遗传或物理作图,比较这些标记在不同物种基因组中的分布特点,揭示染色体或染色体片段上的基因及其排列顺序的相同或相似性,并由此对相关物种的基因组结构和起源进化进行分析。由于不同基因组间,基因编码区的同源性远高于非编码区,因而,基于 ESTs 的比较基因组研究则可以研究基因组功能区域的相似程度,指导基因组功能的预测^[10]。对 27 294 条非冗余的玉米 ESTs 与拟南芥的 ESTs 比较发现,62%~68% 的玉米 ESTs 可以与编码拟南芥蛋白质的 ESTs 匹配^[17]。Salse 等^[18]将 1 411 条玉米 unigene 与水稻 ESTs 和蛋白质进行比较,发现 74% 的玉米基因与水稻基因同源,47% 的可以确定为潜在的直向同源基因。而我们对 139 条非冗余的 ESTs 研究也发现,基于严谨条件的不同,69%~81% 的玉米 ESTs 与水稻基因具有同源性,这些说明,水稻和玉米在 ESTs 上具有很高的同源,利用水稻的 ESTs 信息可以有效地指导相应的玉米基因的定位和克隆。但在本实验中,对 139 条 ESTs 进行定位,其中 43 条 ESTs 不能在水稻中找到符合要求的同源序列。这也反映出利用 ESTs 进行同源对比的不足。导致这种情况的可能原因有:1) ESTs 长度不够;2) 5'-UTR 和 3'-UTR 序列的存在掩盖了编码区的真实的同源性;3) 长期进化过程中变异的积累,使少数玉米和水稻基因的同源性降低;4) 少数 ESTs 可能为玉米或水稻中特

有。不过,69%~81% 的同源性和 55% 的 ESTs 可基于水稻基因组序列成功定位于玉米染色体上,说明基于水稻基因组序列对 ESTs 进行定位的方法是可行的和有效的。

该方法是以水稻基因组数据和已有的禾本科植物遗传图谱为桥梁,基于禾本科植物之间(如水稻和玉米之间)存在的分子标记水平上的广泛的共线性,从而对大量的玉米 ESTs 在玉米连锁群上的分布进行初步定位,为进一步的功能分析和基因克隆提供参考信息。在方法学上,这种尝试不失为一种有益的探索,它为大基因组作物 ESTs 的定位提供了一种新的策略。但标记水平上的共线只能反映宏观水平上的共线性,不足以完全揭示玉米和水稻基因组之间的同源性关系,更不能反映基因组内部大量微小的缺失、插入和移位等分子重排,因而,该方法只是对众多 ESTs 在连锁群上的分布提供了一个可供参考的位置区域,可以利用这些信息从 ESTs 或基因组序列中开发新标记进一步地定位和精细作图等,但利用这些信息进行基因图位克隆会存在潜在的风险。

参考文献(References):

[1] Döhr S, Klingenhoff A, Maier H, Hrabé de Angelis M, Werner T, Schneider R. Linking disease-associated genes to regulatory networks via promoter organization. *Nucleic Acids Research*, 2005, 33(3): 864~872.

[2] Hulbert S H, Richter T E, Axtello J D, Bennetzen J L. Genetic mapping and characterization of sorghum and related crops by means of maize DNA probes. *Proc Natl Acad Sci USA*, 1990,

87:4251~4255.

- [3] Dufour P, Deu M, Grivet L, D'Hont A, Paulet F, Bouet A, Lanaud C, Glaszmann J C, Hamon P. Construction of a composite sorghum genome map and comparison with sugarcane, a related complex polyploid. *Theor Appl Genet*, 1997, 94: 409 ~ 418.
- [4] Guimarães C T, Sills G R, Sobral B W S. Comparative mapping of Andropogoneae: *Saccharum* L. (sugarcane) and its relation to sorghum and maize. *Proc Natl Acad Sci USA*, 1997, 94: 14261~14266.
- [5] Ming R, Liu S C, Lin Y R, da Silva J, Wilson W, Braga D, van Deynze A, Wenslaff T F, Wu K K, Moore P H, Burnquist W, Sorrells M E, Irvine J E, Paterson A H. Detailed alignment of *Saccharum* and *Sorghum* chromosomes; comparative organization of closely related diploid and polyploidy genomes. *Genetics*, 1998, 150: 1663~1682.
- [6] Ahn S, Anderson J A, Sorrells M E, Tanksley S D. Homoeologous relationships of rice, wheat and maize chromosomes. *Mol Gen Genet*, 1993, 241: 483~490.
- [7] Paterson A H, Lin Y R, Li Z, Schertz K F, Doebley J F, Pinson S R M, Liu S C, Stansel J W, Irvine J E. Convergent domestication of cereal crops by independent mutations at corresponding genetic loci. *Science*, 2002, 296: 1714~1717.
- [8] Chen M, San Miguel P, Bennetzen J L. Sequence organization and conservation in *sh2/a1*-homologous regions of sorghum and rice. *Genetics*, 1998, 148: 435~443.
- [9] Tikhonov A P, SanMiguel P, Nakajima Y, Gorenstein N M, Bennetzen J L, Avramova Z. Colinearity and its exceptions in orthologous *adh* regions of maize and sorghum. *Proc Natl Acad Sci USA*, 1999, 96: 7409~7414.
- [10] Chen M, SanMiguel P, de Oliveira A C, Woo S S, Zhang H, Wing R A, Bennetzen J L. Microcolinearity in *sh2*-homologous regions of the maize, rice, and sorghum genomes. *Proc Natl Acad Sci USA*, 1997, 94: 3431~3435.
- [11] International Rice Genome Sequencing Project. The map-based sequence of the rice genome. *Nature*, 2005, 436(7052): 793~800.
- [12] Yu J, Hu S, Wang J, Wong G K, Li S, Liu B, Deng Y, Dai L, Zhou Y, Zhang X, Cao M, Liu J, Sun J, Tang J, Chen Y, Huang X, Lin W, Ye C, Tong W, Cong L, Geng J, Han Y, Li L, Li W, Hu G, Huang X, Li W, Li J, Liu Z, Li L, Liu J, Qi Q, Liu J, Li L, Li T, Wang X, Lu H, Wu T, Zhu M, Ni P, Han H, Dong W, Ren X, Feng X, Cui P, Li X, Wang H, Xu X, Zhai W, Xu Z, Zhang J, He S, Zhang J, Xu J, Zhang K, Zheng X, Dong J, Zeng W, Tao L, Ye J, Tan J, Ren X, Chen X, He J, Liu D, Tian W, Tian C, Xia H, Bao Q, Li G, Gao H, Cao T, Wang J, Zhao W, Li P, Chen W, Wang X, Zhang Y, Hu J, Wang J, Liu S, Yang J, Zhang G, Xiong Y, Li Z, Mao L, Zhou C, Zhu Z, Chen R, Hao B, Zheng W, Chen S, Guo W, Li G, Liu S, Tao M, Wang J, Zhu L, Yuan L, Yang H. A draft sequence of the rice genome (*Oryza sativa* L. ssp. *Indica*). *Science*, 2002, 296: 19~92.
- [13] Yuan Q, Ouyang S, Liu J, Suh B, Cheung F, Sultana R, Lee D, Quackenbush J, Buell C R. The TIGR rice genome annotation resource: annotating the rice genome and creating resources for plant biologists. *Nucleic Acids Research*, 2003, 31(1): 229~233.
- [14] Zhang Z X, Tang W H, Zhang F D, Zheng Y L. Fertility restoration mechanisms in S-Type cytoplasmic male sterility of maize (*Zea mays* L.) revealed through expression differences identified by cDNA microarray and suppression subtractive hybridization. *Plant Mol Biol Rep*, 2005, 23(1): 17~38.
- [15] Kamps T L, Chase C D. RFLP mapping of the maize gametophytic restorer-of-fertility locus (*rf3*) and aberrant pollen transmission of the nonrestoring *rf3* allele. *Theor Appl Genet*, 1997, 95: 525~531.
- [16] SHI Yong-Gang, ZHENG Yong-Lian, LI Jian-Sheng, LIU Ji-Lin. Mapping CMS-S restores gene *Rf3* with RFLPs and RAPDs. *Acta Agronomica Sinica*, 1997, 23(1): 1~6.
石永刚, 郑用链, 李建生, 刘纪麟. 玉米 S 组 CMS 育性恢复基因的分子标记定位. *作物学报*, 1997, 23(1): 1~6.
- [17] Brendel V, Kurtz S, Walbot V. Comparative genomics of Arabidopsis and maize: prospects and limitations. *Genome biology*, 2002, 3(3): reviews1005. 1~1005. 6
- [18] Salse J, Piegu B, Cooke R, Delseny M. New in silico insight into the synteny between rice (*Oryza sativa* L.) and maize (*Zea mays* L.) highlights reshuffling and identifies new duplications in the rice genome. *Plant J*, 2004, 38: 396~409.