

甜瓜 EST 序列中微卫星的分布特征

胡建斌, 刘颖, 王兰菊, 李建吾*

河南农业大学园艺学院, 郑州 450002

摘要: GenBank 中 35 547 条甜瓜 EST 经去冗余处理后, 得到总长度为 250.3 Mb 的无冗余 EST 34 438 条。这些序列中有 2 813 个微卫星简单重复序列 (Simple sequence repeat, SSR), 分布于 2 107 条 EST 中, 出现频率为 8.16%, 平均分布距离为 8.90 kb。三核苷酸重复是主导重复类型, 占 SSR 总数的 47.14%; 其次是二核苷酸和单核苷酸重复, 分别占 SSR 总数的 20.72% 和 16.99%。AAG/TTC 是优势重复基元, 占微卫星总数的 29.26%, AG/CT 和 A/T 分别占 14.61% 和 16.25%。在所有的 SSR 中, 重复次数为 4~10 次的占 70.32%, 长度为 12~20 bp 的占 51.12%。并对这些 SSR 的多态性潜能进行了评价。

关键词: 甜瓜; EST 序列; 微卫星; 特征

Characteristics of Microsatellites in Melon (*Cucumis melo* L.) EST Sequences

HU Jian-Bin, LIU Ying, WANG Lan-Ju, LI Jian-Wu*

College of Horticulture, Henan Agricultural University, Zhengzhou 450002, China

Abstract: 34 438 melon expressed sequence tags (ESTs) were obtained after the removal of redundant sequences from 35 547 melon ESTs deposited in GenBank. In these ESTs, 2 813 microsatellites (simple sequence repeat, SSR) were tested and distributed in 2 107 ESTs. The occurrence frequency of SSR was 8.16% and the average distribution distance was 8.90 kb. Trinucleotide repeats were dominant repeat types and accounted for 47.14%. Dinucleotide and mononucleotide repeats accounted for 20.72% and 16.99%, respectively. AAG/TTC was most frequent repeat motif and accounted for 29.26% in all SSRs. AG/CT and A/T accounted for 14.61% and 16.25%, respectively. Among all SSRs, SSRs which repeated 4 to 10 times accounted for 70.32%, and SSRs which was 12 to 20 bp long accounted for 51.12%. The potential polymorphism of these SSRs was also analyzed.

Key words: melon (*Cucumis melo* L.); EST sequence; microsatellite; characteristic

微卫星或简单重复序列 (simple sequence repeat, SSR) 广泛存在真核生物的基因组中, 因其数量丰富、多态性高、多等位性、共显性等特点 (Powell 等 1996; Varshney 等 2005a), SSR 标记被公认为目前遗传学中最令人信赖的分子标记之一。然而, SSR 引物的开发通常需构建文库、筛选、测序等工作, 费时费力。公共数据库中数量不断增加的表达序列标签 (expressed sequence tag, EST) 极大地增强了对基于 EST 的 SSR 标记开发能力。迄今, 在多种农作物中, 大量的 EST-SSR 标记已经得到开发, 并广泛用于遗传图谱的构建、基因发掘、遗传进化及比较基因组的研究 (李永强等 2004; 李小白等 2006)。

甜瓜 (*Cucumis melo* L.) 为葫芦科 (Cucurbitaceae) 黄瓜属植物, 是一种世界性的园艺作物, 栽培历史悠久, 种质资源丰富, 世界各国都非常重视其遗传学研究。甜瓜基因组学的研究进展较快, 现已构建了 12 张遗传图谱, 一些控制重要经济性状和抗病性

的基因或 QTL 已定位在图谱中 (苏芳等 2007), 为甜瓜基因结构和功能的研究奠定了基础。然而, 现已发表的甜瓜遗传图谱主要由同工酶、RAPD、RFLP、AFLP 等标记构成, 与小麦、水稻和棉花等农作物的图谱相比, 甜瓜遗传图谱所包含的 SSR 标记较少。另一方面, RAPD、AFLP、基因组 SSR 等分子标记或是扩增非编码区域, 或是随机在基因组中扩增, 得到的位点一般与目标性状基因的距离较远, 以致分子标记在应用上与其目标有一定的偏差。EST-SSR 来源于编码区 DNA, 代表基因表达信息, 这使得直接鉴定决定重要表型性状的等位基因的工作成为可能 (Kumpatla 和 Mukhopadhyay 2005)。到现在为止, 甜瓜基因组测序已得到 3 万多条 EST, 为 EST-SSR 标记的开发提供了丰富的资

收稿 2008-11-27 修定 2009-01-08

资助 国家“863”计划 (2007AA10Z100) 和河南农业大学博士启动基金 (30400247)。

* 通讯作者 (E-mail: lijw555@sohu.com; Tel: 0371-63554959)。

源。本文从 GenBank 公布的 EST 中查找甜瓜 SSR, 分析这些 SSR 在甜瓜转录组中的特点和分布规律, 从而为 EST-SSR 标记的开发建立基础。

材料与方法

以“melon mRNA”为关键词采用 FASTA 格式从 GenBank/dbEST (<http://www.ncbi.nlm.nih.gov/entrez>) 下载 35 547 条甜瓜 (*Cucumis melo* L.) EST 序列(以 2008 年 5 月 20 日 NCBI 公布的数据为准), 它们主要来自甜瓜的果实、叶片、根等组织。

采用 EST-trimmer (http://pgrc.ipk-gatersleben.de/misa/download/est_trimmer.pl)、cross-match (www.phrap.org) 等网络软件除去 EST 中重复、polyA/T “尾巴”、载体等冗余序列。用 MISA (<http://pgrc.ipk-gatersleben.de/misa/misa.html>) 在无冗余 EST 中搜索 SSR, 并结合手工查寻。搜索条件为: 含有单、二和三核苷酸基元的最小重复数分别为 20、8 和 5, 四核苷酸或四核苷酸以上的 SSR 最小重复数均为 4。分别统计甜瓜不同组织中 SSR 的数目和出现频率。为了便于统计和分析 SSR 序列, 将重复基元所有循环序列和互补序列视为一类。如 ACT、CTA、TAC、TGA、GAT 和 ATG 均归为一类。

实验结果

1 甜瓜 EST-SSR 的频率和分布密度

从 NCBI 网站下载的 35 547 条甜瓜 EST 序列, 主要来源于甜瓜的根(13 853 条)、果实(10 585 条)、子叶(5 664 条)、叶片(3 212 条)和韧皮部(1 800 条), 另外还有少量来自悬浮细胞和愈伤组织。经净化处理(除去重复、polyA/T、载体等)后共获得 34 438 条无冗余的 EST, 序列总长度为 250.3 Mb。按照查找标准, 共发现 2 107 条至少含有 1 个 SSR 的 EST 序列, 占无冗余 EST 总数的 6.12%, 表明甜瓜 EST 中微卫星含量较为丰富。在 2 107 条 EST 中, 含有单个 SSR 的 EST 为 1 458 条, 含有 2 个或 2 个以上 SSR 的 EST 为 649 条, 其中还有 69 条序列出现两个 SSR 串联。共检出 2 813 个精确重复的 SSR, 占无冗余 EST 的 8.16%, 即甜瓜基因组中 EST-SSR 的出现频率。

甜瓜 EST-SSR 的优势重复基元为单、二和三核苷酸, 三者共占 EST-SSR 总数的 84.85%, 其中

又以三核苷酸重复所占的比例最大(47.14%), 二核苷酸重复次之(20.72%), 单核苷酸重复最少(16.99%)。基元长度大于等于 4 的重复序列所占的比例较小, 共计 15.15%。从 SSR 的分布密度来看, 甜瓜 EST 中平均每 8.90 kb 就出现 1 个 SSR, 但不同重复单元出现的平均距离各不一致, EST-SSR 出现的频率越高, 其平均距离越小(表 1)。甜瓜不同组织中 SSR 的分布密度不尽相同, 韧皮部中 SSR 密度最大(1/6.14 kb), 子叶次之(1/8.44 kb), 果实中最少(1/10.53 kb)。

表 1 甜瓜中 EST-SSR 的数量、比例和平均距离

Table 1 Number, percentage and mean distance of EST-SSRs in melon

重复类型	数目	所占比例 /%	出现频率 /% ¹⁾	平均距离 /kb ²⁾
单核苷酸	478	16.99	1.39	52.37
二核苷酸	583	20.72	1.69	42.94
三核苷酸	1326	47.14	3.85	18.88
四核苷酸	115	4.09	0.33	217.67
五核苷酸	180	6.40	0.52	139.07
六核苷酸	131	4.66	0.38	191.08
总计	2813	100.00	8.16	8.90

1) 出现频率 = 检出的 SSR 数目 / 无冗余 EST 总数;

2) 平均距离 = 无冗余 EST 总长度 / SSR 总数。

2 甜瓜 EST-SSR 中基元类型及比例

在搜索到的 2 813 个甜瓜 EST-SSR 中共观察到 69 种重复基元, 一至六核苷酸重复分别有 2、3、10、14、22 和 18 种。不同类型的基元的出现频率不一致, 存在明显的偏倚性。单核苷酸重复基元中 A/T 占绝对优势, 所占比例为 16.25%, C/G 极少(18 个, 占 0.64%); 二核苷酸重复基元中 AG/CT 比例最高, 达 14.61%, AT 次之(3.63%), AC/GT 较少(2.49%), 而 CG 则没有出现; 三核苷酸重复基元中 AAG/TTC 比例高达 29.26%, 除 CCG/CGG 外(15 个, 占 0.53%), 其他各类三核苷酸重复基元的比例均在 1%~5% 之间; 四核苷酸重复基元中比例最高的是 AAAG/CTTT (1.78%); 五核苷酸重复基元 AAAAG/CTTTT 比例最高(2.13%); 六核苷酸中 AAAAAG/CTTTTT 和 ACCACG/CTGGTG 的比例分别为 0.64% 和 0.60% (表 2)。

表2 甜瓜 EST 中主要重复基元

Table 2 The major repeat motifs in melon ESTs

重复类型	重复基元	数量	所占比例 /%	出现频率 /%
单核苷酸	A/T	457	16.25	1.33
	C/G	18	0.64	0.05
二核苷酸	AG/CT	411	14.61	1.19
	AT/AT	102	3.63	0.30
	AC/GT	70	2.49	0.20
三核苷酸	AAG/CTT	823	29.26	2.39
	AGG/CCT	116	4.12	0.34
	AAC/GTT	97	3.45	0.28
	AGC/CGT	67	2.38	0.19
	AGT/ATC	54	1.92	0.16
	AAT/ATT	39	1.39	0.11
	ACC/GGT	40	1.42	0.12
	ACT/ATG	38	1.35	0.11
	ACG/CTG	37	1.32	0.11
	CCG/CGG	15	0.53	0.04
	四核苷酸	AAAG/CTTT	50	1.78
AAAT/ATTT		30	1.07	0.09
AAGG/CCTT		10	0.36	0.03
AAAC/GTTT		6	0.21	0.02
AATG/ACTT		5	0.18	0.01
五核苷酸	AAAAG/CTTTT	60	2.13	0.17
	AAAAT/ATTTT	17	0.60	0.05
	AATAC/ATGTT	16	0.67	0.05
	AAGAG/CTCTT	15	0.53	0.04
	AAAAC/GTTTT	12	0.43	0.03
	AAAGG/CCTTT	9	0.32	0.03
	AACCT/ATTGG	9	0.32	0.03
	AATCCAGGTT	8	0.28	0.02
	AAACC/GGTTT	5	0.18	0.01
	AATTC/AAGTT	5	0.18	0.01
	六核苷酸	AAAAAAG/CTTTTT	18	0.64
ACCACG/CTGGTG		17	0.60	0.05
AAGAGG/CCTTCT		10	0.36	0.03
AAAAAT/ATTTTT		7	0.25	0.02
AGGTAT/ATATCC		6	0.21	0.02
AACACG/CTTGTT		6	0.21	0.02
AAGATG/ACTTCT		5	0.18	0.01
AACCTC/AGTTGG		5	0.18	0.01
AAGACG/CTTCTG		5	0.18	0.01

出现频率小于 0.01% 的重复基元未列出。

3 甜瓜 EST-SSR 重复次数和长度

SSR 重复次数的变异引起位点长度的变化是产生 SSR 多态性的主要原因。对 2813 个甜瓜 EST-SSR 以及主要重复类型(二核苷酸和三核苷酸重复)进行分类统计, 结果表明, 随着重复次数的增加, SSR 数量迅速减少(图 1)。甜瓜 EST-SSR 按重复次数可分为 3 个区间, 即 4 次至 10 次重复为第一个区间,

11 次至 20 次重复为第二区间, 20 次重复以上为第三区间。统计结果发现, 甜瓜 EST-SSR 主要分布在第一区间, 这一区间共有 1 978 个微卫星, 约占全部微卫星的 70.32%, 一至六核苷酸重复基元均有分布, 其中三核苷酸重复比例最大; 第二区间有 437 个微卫星, 约占总数的 15.54%, 主要为二核苷酸和三核苷酸重复; 第三区间有 398 个微卫星, 主要是

单核苷酸和部分二核苷酸重复, 占 14.15%。

重复基元长度变化是 EST-SSR 位点多态性的主要表现形式。甜瓜 EST-SSR 长度分布情况见图 2。由于搜索标准的严格度, 部分 SSR 被过滤掉(特别是部分长度小于 20 bp 的单核苷酸重复), 因此 EST-SSR 长度分布是不连续的, 如长度为 13 bp、

17 bp 和 19 bp 的 SSR 没有出现。总的来说, 大部分甜瓜 EST-SSR 长度集中在 12~20 bp 范围内(1466 个 SSR, 占 52.12%), 几乎全为二核苷酸和三核苷酸重复; 其次是 21~30 bp (920 个 SSR, 占 32.71%), 这一区间三核苷酸所占比例最大; 30 bp 以上的 SSR 数量相对较少(427 个 SSR, 占 15.18%), 主要由单

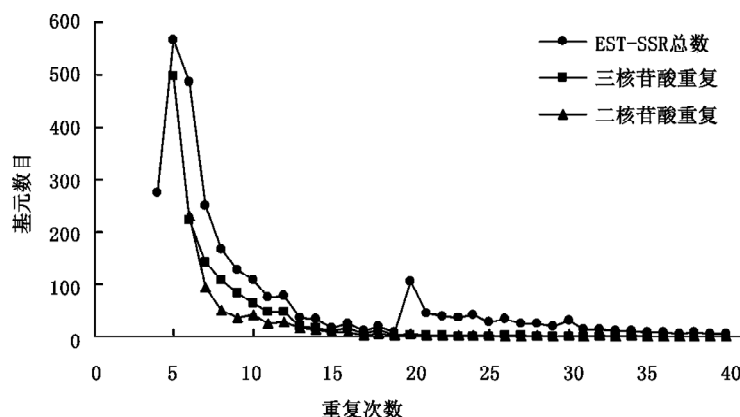


图1 甜瓜 EST-SSR 重复次数分布

Fig.1 Distribution of EST-SSR repeat numbers in melon

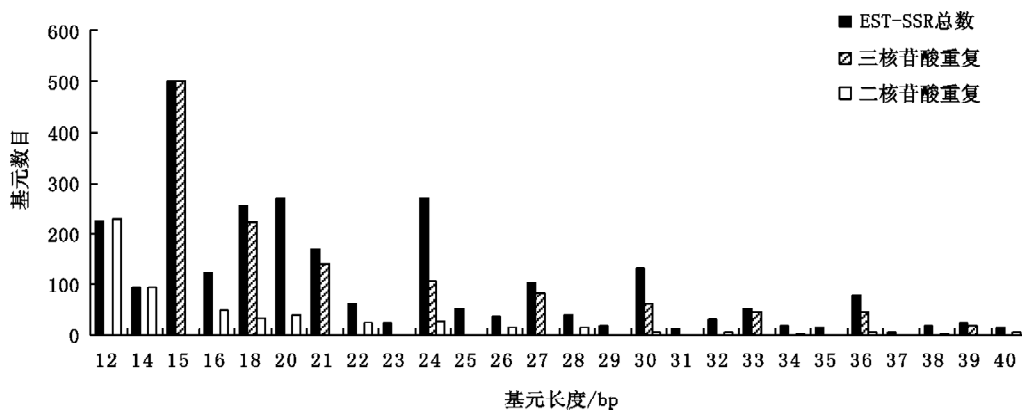


图2 甜瓜 EST-SSR 长度分布

Fig.2 Distribution of EST-SSR length in melon

核苷酸和三核苷酸重复组成。

讨 论

公共数据库中数量不断增加的 EST 信息极大地增强了 SSR 标记的开发能力。国际葫芦科基因组计划已全面启动, 甜瓜作为葫芦科研究作物之一, 已经产生 3 万多条 EST 序列, 这为 EST-SSR 标记的开发提供了宝贵的资源。但现已开发的甜瓜 EST-SSR 标记引物却不足 30 对(Danin-Poleg 等

2001; Kong 等 2007), 这与 EST 数量极不相配, 因此甜瓜 EST-SSR 标记还有进一步开发的潜力。明确甜瓜基因组中 EST-SSR 的分布特征是开发新的 SSR 标记的前提。本文发现, 甜瓜 EST-SSR 基元类型丰富、出现频率高、密度大, 平均距离仅为 8.17 kb, 高于拟南芥(14.9 kb)、小麦(15.6 kb)、棉花(20.0 kb)和番茄(11.1 kb)等作物(李永强等 2004), 这说明大规模开发 EST-SSR 标记的可能性较大。例如, 雷蒙德氏棉 EST 中 SSR 的平均距离是 14.8

kb, 目前已从 58 906 条非冗余雷蒙德氏棉 EST 中成功开发 1 554 对 EST-SSR 引物(王长彪等 2006)。Kong 等(2007)从 5747 条甜瓜 EST (2006 年 7 月 31 日前 NCBI 公布的数据)中发现了 383 个 SSR 位点, 主要是二核苷酸重复(49.9%)和三核苷酸重复(43.6%), EST-SSR 出现频率是 1/4.7 kb, 高于本研究结果 (1/14.8 kb), 这主要是由于搜索标准不同(最小重复数为 5)所致。本文中甜瓜 EST-SSR 中主要重复类型为单、二和三核苷酸(约占总数的 85%), 这与大多数植物中的报道相似(李永强等 2004; Kumpatla 和 Mukhopadhyay 2005)。甜瓜 EST-SSR 中出现频率最高的基元是 AAG/CTT (29.26%), 这与柑橘(14.2%) (Jiang 等 2006)、棉花(26.13%) (王长彪等 2006)和拟南芥(29.00%) (Cardle 等 2000)中的报道相似, 进一步验证了 Gao 等(2003)认为 AAG/TTC 是双子叶植物中优势重复基元的观点。甜瓜 EST 中 AAG/TTC 的高频率出现可能与其作为三联体密码编码相应蛋白质时的高频率使用有关。

SSR 位点多态性主要是因重复基元数量和基元碱基数不同所产生的简单序列长度多态性和随机扩增微卫星多态性。一般认为, SSR 位点的变异频率与基元重复数存在一定的正相关, 即重复次数越多 SSR 产生变异的可能性越大(Schlötterer 2000)。而本文结果表明, 甜瓜 70.32% 的 EST-SSR 是低重复数基元(4 至 10 次重复), 只有 29.68% 的 EST-SSR 重复数在 10 以上。从这个角度讲, 本研究所发掘的 EST-SSR 中, 仅有少部分具有多态性潜能。而 Xu 等(2000)认为微卫星变异依赖于等位基因的长度, 等位基因的重复序列扩张和收缩的频率在总体上是相等的, 即重复序列长度存在一个阈值, 长度在阈值以下的 SSR 倾向扩张, 而长度在阈值以上的 SSR 倾向收缩, 大多数真核生物阈值约为 20 bp。甜瓜 EST-SSR 中, 长度在 20 bp 以下的占 42.48%, 长于 20 bp 的占 47.88%, 这些位点因未达到或超过阈值长度而倾向于扩张或收缩。因此, 按照 Xu 等(2000)观点, 我们所发掘的甜瓜 EST-SSR 位点大部分具有多态性潜能。在 Kong 等(2007)所开发的 22 对甜瓜 EST-SSR 引物中, 约 1/4 SSR 位点长度小于或等于 20 bp, 但在品种间均能表现出多态性, 充分说明 SSR 多态性潜能与其长度(或重复次数)没有直接关系。

由于 EST 来源于编码区 DNA, EST-SSR 代表了基因表达的信息, 能为功能基因提供“绝对”标

记。此外, EST-SSR 标记具有物种间通用性, EST-SSR 遗传作图将使物种之间连锁信息的转换更快, 实现多个图谱整合, 从而更有利于比较基因组学的研究(Lan 等 2000; Varshney 等 2005b)。现在, 依据本文结果所进行的甜瓜 EST-SSR 引物大规模开发, 已在本实验室全面展开。

参考文献

- 李小白, 崔海瑞, 张明龙(2006). EST 分子标记开发及在比较基因组学中的应用. 生物多样性, 14 (6): 541~547
- 李永强, 李宏伟, 高丽锋, 何蓓如(2004). 基于表达序列标签的微卫星标记(EST-SSRs)研究进展. 植物遗传资源学报, 5 (1): 91~95
- 苏芳, 郭绍贵, 宫国义, 张海英, 许勇(2007). 甜瓜基因组学研究进展. 分子植物育种, 5 (4): 540~547
- 王长彪, 郭旺珍, 蔡彩平, 张天真(2006). 雷蒙德氏棉 EST-SSRs 分布特征及开发与利用. 科学通报, 51 (3): 316~320
- Cardle L, Ramsay L, Milbourne D, Macaulay M, Marshall D, Waugh R (2000). Computational and experimental characterization of physically clustered simple sequence repeats in plants. *Genetics*, 156: 847~854
- Danin-Poleg Y, Reis N, Tzuri G, Katzir N (2001). Development and characterization of microsatellite markers in *Cucumis*. *Theor Appl Genet*, 102: 61~72
- Gao LF, Tang JF, Li HW, Jia JZ (2003). Analysis of microsatellites in major crops assessed by computational and experimental approaches. *Mol Breed*, 12: 245~261
- Jiang D, Zhong GY, Hong QB (2006). Analysis of microsatellites in citrus unigenes. *Acta Genet Sin*, 33 (4): 345~353
- Kong Q, Xiang C, Yu Z, Zhang C, Liu F, Peng C, Peng X (2007). Mining and charactering microsatellites in *Cucumis melo* expressed sequence tags from sequence database. *Mol Ecol Notes*, 7: 281~283
- Kumpatla SP, Mukhopadhyay S (2005). Mining and survey of simple sequence repeats in expressed sequence tags of dicotyledonous species. *Genome*, 48: 985~998
- Lan TH, DelMonte TA, Reischmann KP, Hyman J, Kowalski SP, McFerson J, Kresovich S, Paterson AH (2000). An EST-enriched comparative map of *Brassica oleracea* and *Arabidopsis thaliana*. *Genome Res*, 10: 776~788
- Powell W, Machray GC, Provan J (1996). Polymorphism revealed by simple sequence repeats. *Trends Plant Sci*, 1: 215~222
- Schlötterer C (2000). Evolutionary dynamics of microsatellite DNA. *Chromosoma*, 109: 365~371
- Varshney RK, Graner A, Sorrells ME (2005a). Genic microsatellite markers in plants: features and applications. *Trends Biotechnol*, 23: 48~55
- Varshney RK, Sigmund R, Börner A, Korzun V, Stein N, Sorrells ME, Langridge P, Graner A (2005b). Interspecific transferability and comparative mapping of barley EST-SSR markers in wheat, rye and rice. *Plant Sci*, 168: 195~202
- Xu X, Peng M, Fang Z, Xu X (2000). The direction of microsatellite mutations is dependent upon allele length. *Nat Genet*, 24: 396~399