

大豆 EST 资源的 SSR 信息分析

陈相艳¹, 李 伟², 戴海英², 张礼凤²

(¹山东省农业科学院办公室, 山东 济南 250100; ²山东省农业科学院作物所, 山东 济南 250100)

摘要:微卫星或简单重复序列 (simple sequence repeats, SSR) 存在于表达序列标签 (expressed sequence tags, ESTs) 中。为了在大豆中开发 EST-SSR 功能性标记, 利用生物信息学对 NCBI 公共数据库中的 394 370 条大豆 ESTs 序列进行 EST-SSRs 特征分析。剔除冗余序列, 得到全长为 51 332.21 kb 的无冗余 EST 80 735 条。在这些序列中搜索出 7 754 个 SSR, 分布于 6 674 条 EST 中, 出现频率是 8.27%。这些 EST-SSR 的平均长度为 15.26 bp, 平均分布频率是 1/6.62 kb。在 1~6 bp 的重复基元中, 三核苷酸重复基元的 SSRs 出现频率最高 (43.72%), 其次是二核苷酸 (37.34%)、单核苷酸 (15.92%)。AG/CT 和 AAG/CTT 是二、三核苷酸中的优势重复基元, 分别占二、三核苷酸重复的 62.83% 和 25.22%。本研究为开发多态性大豆微卫星标记提供了候选序列。

关键词:大豆; EST; SSR; 频率; 特性

中图分类号: S565.1

文献标识码: A

文章编号: 1000-9841(2009)03-0394-06

Analysis of SSR Information in EST Resource of Soybean (*Glycine max*)

CHEN Xiang-yan¹, LI Wei², DAI Hai-ying², ZHANG Li-feng²

(¹Office, Shandong Academy of Agricultural Sciences, Jinan 250100, Shandong; ²Crop Institute, Shandong Academy of Agricultural Sciences, Jinan 250100, Shandong, China)

Abstract: Expressed sequence tags (ESTs) are important resources for development of new SSR markers. In this study, 394 370 ESTs of soybean (*Glycine max* L.) in the database of NCBI were downloaded and analyzed. After the preprocessing, resulting in 80 735 non-redundant ESTs with total length about 51 332.21 kb. Totally 7 754 SSRs distributed in 6 674 ESTs were detected, accounting for 8.27% of the non-redundant ESTs. The average length and distribution distance of the EST-SSRs were about 15.26 bp and 6.62 kb, respectively. Dinucleotide and trinucleotide repeats with similar frequency are the main types, accounting for 81.06% of all the SSRs. AG/CT and AAG/CTT are the most frequent motifs, accounting for 62.83% and 25.22% in the dinucleotide and trinucleotide repeats, respectively. These EST-SSRs will help to develop new SSR markers with high polymorphism for soybean.

Key words: Soybean; EST; SSR; Frequency; Characteristics

大豆是重要的经济作物, 在中国驯化栽培已有上千年的历史^[1], 是我国重要的油料作物, 然而产量低一直是限制大豆生产发展的重要因素。传统育种在大豆的生产中做出了重大贡献, 但受育种周期长、随机性大、农艺性状难以综合等因素的限制, 很难在产量和品质上有大的突破。随着分子生物学和基因工程的发展, 大量植物功能基因的分离克隆、表达调控方式和分子机制的阐明, 以及分子水平遗传改良和分子标记辅助育种手段的不断完善, 为大豆分子育种提供了契机。

在众多分子标记中, SSR (simple sequence repeat) 具有在基因组中分布随机、信息量高、多态性强、共显性和孟德尔遗传等优点, 已在大豆起源、遗传图谱构建、遗传多态性分析、品种鉴定和分子标记辅助育种等方面广泛应用^[2-3]。但 SSR 标记的开发通常需要经过文库构建, 包括 SSR 克隆的筛选、测序、引物设计和 PCR 扩增与分析等步骤, 要投入大量的人力物力。

近年来, 许多作物开展了大规模 cDNA 测序工作, GenBank 中 EST 数据高速增长, 大大增加了基于

收稿日期: 2008-12-11

基金项目: 山东省自然科学基金资助项目 (Y2006D20)。

作者简介: 陈相艳 (1973-) 女, 硕士, 现主要从事大豆育种工作; 李伟为共同第一作者。

通讯作者: 张礼凤, 副研究员。E-mail: zhanglifeng9639@sina.com。

ESTs 的 SSR 标记开发能力。目前,EST-SSR 的开发在小麦^[4]、大麦^[5]、黑麦^[6]、水稻^[7]、马铃薯^[8]、葡萄^[9]和棉花^[10]等作物中已有报道,并广泛应用于基因组研究和分子育种。目前,大豆连锁图谱上 SSR 分子标记已有 1014 个,但图谱上大于 5 cM 的区间有 138 处,其中大于 10 cM 的区间有 26 处^[11-12]。因而有必要寻找新的 SSR 位点,创建更为密集分子标记图谱,为大豆重要农艺性状基因的标记和克隆奠定基础。

截止 2008 年 6 月 25 日,在 GenBank 中释放的大豆 EST 共有 394 370 条。对现有大豆 EST 中的 SSR 信息进行了全面分析,明确了大豆 EST-SSR 的发生频率和分别特点,为进一步开发大豆 SSR 标记和探索其在大豆遗传育种中的应用奠定基础。

1 材料和方法

1.1 大豆 EST 来源

大豆 EST 来自 NCBI(美国国家生物技术信息中心)数据库(<http://www.ncbi.nlm.nih.gov/>),共计 394 370 条(截止 2008 年 6 月 25 日)。

1.2 EST 前处理

采用 EST-trimmer 软件(http://pgrc.ipk-gatersleben.de/misa/download/est_trimmer.pl)去除 5' 端或 3' 端 50 bp 的 polyT 或 polyA 以及长度小于 100 bp 的 EST 序列;长度大于 700 bp 的 EST 则保留其 5' 端 700 bp。

1.3 聚类去冗余

前处理后的 ESTs 通过软件 Phrap^[13-14] 进行片段重叠群分析和聚类。拼接时设定的初始装配参数为:最小匹配碱基数(minmatch)为 30,最小分值(minscore)为 30。对错误拼接序列设置比较高的装配参数再次进行拼接,判别,共进行了 3 次。

1.4 SSR 筛选

用 MISA 软件(MicroSATellite)(<http://www.pgrc.ipk-gatersleben.de/misa>)对聚类后的 ESTs 进行 SSR 搜索。筛选标准为:单核苷酸重复的次数在 16 次或 16 次以上,二核苷酸重复的次数在 6 次或 6 次以上,三至六核苷酸重复的次数在 5 次或 5 次以上。同时,也筛选中间被少数碱基(间隔小于 10 或等于 10)打断的(interrupted)不完全重复的 SSR^[15]。

2 结果与分析

2.1 大豆 EST 中出现 SSR 的频率

对经过处理后得到的 80 735 条无冗余 EST 序列进行搜索,共检出含有 SSR 的序列 6 674 条,发生频率(含有 SSR 的 EST 数目与总 EST 数目的比值)为 8.27%。其中,5781 条含单个 SSR,893 条含有 2 个或 2 个以上的 SSR。共检出 7 754 个 SSR,占无冗余 EST 的 9.60%。在这 7 754 个 SSR 中,完全重复 SSR 为 7481 个,不完全重复 SSR 有 273 个。从分布情况看,大豆 EST 中平均每 6.62 kb 就出现 1 个 SSR,但不同重复类型间差异很大(表 1)。

表 1 SSR 在大豆无冗余 EST 中的出现频率

Table 1 Occurrence of SSRs in non-redundant soybean ESTs

类型 Type	数目 Number	各类型的比例 Proportion in all SSRs/%	频率 Frequency/%	平均距离 Average distance/kb
单核苷酸 Mononucleotide	1235.00	15.92	1.32	41.56
二核苷酸 Dinucleotide	2895.00	37.34	3.09	17.73
三核苷酸 Trinucleotide	3390.00	43.72	3.61	15.14
四核苷酸 Tetranucleotide	148.00	1.91	0.16	346.84
五核苷酸 Pentanucleotide	28.00	0.36	0.03	1833.29
六核苷酸 Hexanucleotide	58.00	0.75	0.06	885.04
总计 Total	7754.00	100.00	8.27	6.62

大豆的 EST-SSR 种类十分丰富,一至六核苷酸重复类型都能看到,但各种类型出现的频率相差很大(表 2)。主要集中在一至三核苷酸重复上,占总 EST-SSR 的 96.98%,其中,三核苷酸重复最为常

见,占总 SSR 的 43.72%,二核苷酸重复占总 SSR 的 37.34%,而四至六核苷酸重复所占比例较小,仅占总 SSR 的 3.02%。

续表 2

重复基元 SSR Motif	重复次数 Number of repeats													合计 Total
	5	6	7	8	9	10	11	12	13	14	15	>16		
AATCCC/AGGTT	1	1	0	0	0	0	0	0	0	0	0	0	0	2
AATGAT/ACTATT	1	0	0	0	0	0	0	0	0	0	0	0	0	1
ACACCT/ATGTGG	1	0	0	0	0	0	0	0	0	0	0	0	0	1
CACTC/AGTGTG	0	1	0	0	0	0	0	0	0	0	0	0	0	1
ACATCT/AGATGT	1	0	0	0	0	0	0	0	0	0	0	0	0	1
ACCAGC/CCTGCT	1	1	0	0	0	0	0	0	0	0	0	0	0	2
ACCAGT/ATGGTC	1	0	0	0	0	0	0	0	0	0	0	0	0	1
ACCCGC/CCTGGG	0	1	0	0	0	0	0	0	0	0	0	0	0	1
ACCCGT/ATGGGC	1	0	0	0	0	0	0	0	0	0	0	0	0	1
ACCCTC/AGTGGG	1	0	0	0	0	0	0	0	0	0	0	0	0	1
ACCTCC/AGGTGG	1	0	0	0	0	0	0	0	0	0	0	0	0	1
ACGAGT/ATGCTC	1	0	0	0	0	0	0	0	0	0	0	0	0	1
ACGATC/AGTGTCT	1	0	0	0	0	0	0	0	0	0	0	0	0	1
ACGGAG/CCTCTG	1	0	0	0	0	0	0	0	0	0	0	0	0	1
ACGGGC/CCCCTG	0	1	0	0	0	0	0	0	0	0	0	0	0	1
ACTCCT/AGGATG	1	0	0	0	0	0	0	0	0	0	0	0	0	1
AGAGCC/CGGTCT	1	0	0	0	0	0	0	0	0	0	0	0	0	1
AGAGGT/ATCTCC	1	0	0	0	0	0	0	0	0	0	0	0	0	1
AGCGGT/ATCGCC	2	1	0	0	0	0	0	0	0	0	0	0	0	3
合计 Total	2035	1691	889	543	388	250	168	105	78	52	44	1511	7754	

3 讨论

对 NCBI 数据库中的 394 370 条大豆 ESTs 序列进行了聚类分析,共得到无冗余的 EST 序列 80 735 条,从中搜索出 7 754 个 SSR,占无冗余 EST 序列总数的 9.60%。这个比率远高于甘蔗(2.9%)^[16]、水稻(4.7%)^[17]、普通小麦(5.4%)^[18]等物种 EST 数据库中筛选的 SSR 比率,但低于油菜(15.58%)^[15]和花生(11.27%,私人通讯),与白菜(10.34%)^[19]相近。这种差异可能是物种间的真实 SSR 信息差异或搜寻 SSR 时所用长度最低标准不同造成的。

当 4 种不同碱基随机组合时,若 EST-SSR 数目足够大且无偏倚性,将可能产生 2 种单核苷酸、4 种二核苷酸、10 种三核苷酸、33 种四核苷酸、102 种五核苷酸和 350 种六核苷酸基本重复基元类型^[20]。对大豆 EST-SSR 分析的结果显示,不同重复基元出现与否及其频率高低表现出明显的偏倚性。二、三核苷酸重复基元最为常见,占大豆总 SSR 的 81.06%。GC 重复基元在多数植物中很难见到,如拟南芥、杏树、桃树、水稻、玉米等^[21]植物中都没发现。在大豆 EST-SSR 中发现了 GC 重复基元,只是出现频率很低(0.008/100 kb)。在大多数植物中,四、五和六核苷酸重复基元虽有出现,但类型很少,表现出明显的偏倚性,但在大豆中,四核苷酸重复基

元出现类型很多,高达 21 种,占理论类型(33 种)的 63.64%。

EST 是功能基因的一部分序列,从 EST 中开发出的 EST-SSR 标记理论上可为功能基因提供“绝对”的标记^[22]。从研究的结果看,大豆 EST 中的 SSR 不但出现频率高,而且类型丰富,为进一步开发大豆 EST-SSR 标记奠定了基础。

参考文献

- [1] 徐豹,郑惠玉,路琴华,等.大豆起源地的三个新论据[J].大豆科学,1986,5(2):123-130. (Xu B, Lu Q H, Zhao S W, et al. Three new evidences of the original area of soybean[J]. Soybean Science, 1986, 5(2): 123-130.)
- [2] 王彪,邱丽娟.大豆 SSR 技术研究进展[J].植物学通报,2002, 19(1):44-48. (Wang B, Qiu L J. Current advance of simple sequence repeats in soybean[J]. Chinese Bulletin of Botany, 2002, 19(1):44-48.)
- [3] 宋启建.大豆 SSR 分子标记的创制及其应用[J].大豆科学, 1999, 18(3):248-254. (Song Q J. A review of development and application of simple sequence repeat in soybean[J]. Soybean Science, 1999, 18(3):248-254.)
- [4] Peng J H, Nore L, Lapitan V. Characterization of EST-derived microsatellites in the wheat genome and development of eSSR markers [J]. Functional & Integrative Genomics, 2005, 5:80-96.
- [5] Thiel T, Michalck W, Varshney R K, et al. Exploiting EST databases for the development and characterization of gene-derived SSR-markers in barley (*Hordeum vulgare* L.) [J]. Theoretical and Ap-

- plied Genetics, 2003, 106:411-422.
- [6] Hackauf B, Wehling P. Identification of microsatellite polymorphisms in an expressed portion of the rye genome [J]. Plant Breeding, 2002, 121:17-25.
- [7] Cho Y G, Ishii T, Temnykh S, et al. Diversity of microsatellites derived from genomic libraries and GenBank sequences in rice (*Oryza sativa* L.) [J]. Theoretical and Applied Genetics, 2000, 100:713-722.
- [8] Feingold S, Lloyd J, Norero N, et al. Mapping and characterization of new EST-derived microsatellites for potato (*Solanum tuberosum* L.) [J]. Theoretical and Applied Genetics, 2005, 111 (3):456-466.
- [9] Decriq V, Fave M G, Hagen L, et al. Development and transferability of apricot and grape EST microsatellite markers across taxa [J]. Theoretical and Applied Genetics, 2003, 106:912-922.
- [10] Han Z G, Guo W Z, Song X L, et al. Genetic mapping of EST-derived microsatellites from the diploid *Gossypium arboreum* in allotetraploid cotton [J]. Molecular Genetics and Genomics, 2004, 272:308-327.
- [11] Cregan P B, Jarvik T, Bush A L, et al. An integrated genetic linkage map of the soybean genome [J]. Crop Science, 1999, 39:1464-1490.
- [12] Song Q J, Marek L F, Shoemaker R C, et al. A new integrated genetic linkage map of the soybean [J]. Theoretical and Applied Genetics, 2004, 109:122-128.
- [13] Gordon D, Abajian C, Green P. Consed: a graphical tool for sequence finishing [J]. Genome Research, 1998, 8:195-202.
- [14] Gordon D, Desmarais C, Green P. Automated finishing with auto-finish [J]. Genome Research, 2001, 11:614-625.
- [15] 李小白, 张明龙, 崔海瑞. 油菜 EST 资源的 SSR 信息分析 [J]. 中国油料作物学报, 2007, 29(1):20-25. (Li X B, Zhang M L, Cui H R. Analysis of SSR information in EST resource of oilseed rape [J]. Chinese Journal of Oil Crop Sciences, 2007, 29 (1):20-25.)
- [16] Cordeiro G M, Casu R, McIntyre C L, et al. Microsatellite markers from sugarcane (*Saccharum ssp.*) ESTs cross transferable to erianthus and sorghum [J]. Plant Sciences, 2001, 160:1115-1123.
- [17] Kantety R V, Rota M L, Matthews D E, et al. Data mining for simple sequence repeats in expressed sequence tags from barley, maize, rice, sorghum and wheat [J]. Plant Molecular Biology, 2002, 48:501-510.
- [18] Gupta P K, Rustgi S, Sharma S, et al. Transferable EST-SSRs markers for the study of polymorphism and genetic diversity in bread wheat [J]. Molecular Genetics and Genomics, 2003, 270:315-323.
- [19] 忻雅, 崔海瑞, 卢美贞, 等. 白菜 EST-SSR 信息分析与标记的建立 [J]. 园艺学报, 2006, 33(3):549-554. (Xin Y, Cui H R, Lu M Z, et al. Data mining for SSRs in ESTs and EST-SSR marker development in Chinese cabbage [J]. Acta Horticulturae Sinica, 2006, 33(3):549-554.)
- [20] Rota L R, Kantety R V, Yu J K. Nonrandom distribution and frequencies of genomic and EST-derived microsatellite markers in rice, wheat and barley [J]. BMC Genomics, 2005, 6:23.
- [21] Gao L F, Tang J F, Li H W, et al. Analysis of microsatellites in major crops assessed by computational and experimental approaches [J]. Molecular Breeding, 2003, 12:245-261.
- [22] 王长彪, 郭旺珍, 蔡彩平, 等. 雷蒙德氏棉 EST-SSRs 分别特征及开发与利用 [J]. 科学通报, 2006, 51(3):316-320. (Wang C B, Guo W Z, Cai C P, et al. Characterization, development and exploitation of EST-derived microsatellites in gossypium raimondii ulbrich [J]. Chinese Science Bulletin, 2006, 51(3):316-320.)

欢迎订阅 2009 年《大豆科学》

《大豆科学》是由黑龙江省农科院主管主办的大豆领域专业学术性期刊,是中国自然科学核心期刊,中国科学引文数据库来源期刊及国内外多家权威数据库收入期刊源。主要刊登有关大豆的遗传育种,品种资源,生理生态,耕作栽培,病、虫、杂草防治,营养施肥,生物技术,食品加工,药理研究和工业用途等方面的科研报告,学术论文,国内、外研究进展评述,研究简报,学术活动简讯、新品种介绍等。《大豆科学》主要面向从事大豆科学研究的科技工作者,大专院校师生、各级农业技术推广部门的技术人员及科技种田的农民。

国内外公开发行,双月刊,16 开本,每期 180 页。国内每期订价:10.00 元,全年 60.00 元,邮发代号:14-95。国外每期订价:10.00 美元(包括邮资),全年 60 美元。国外由中国国际图书贸易总公司发行,北京 399 信箱。国外代号:Q5587。

本刊热忱欢迎广大科研及有关企事业单位刊登广告,广告经营许可证号:2301030000004。

地 址:哈尔滨市南岗区学府路 368 号《大豆科学》编辑部

邮 编:150086

电 话:0451-86668735

E-mail: dadoukx@sina.com