

研究报告

Research Report

五种豆科植物物种基于 De Novo 序列装配和注释的比较研究

Sagar S. Patel¹, Dipti B. Shah¹, HetaKumar J. Panchal¹

1. G. H. Patel Post Graduate Department of Computer Science and Technology, Sardar Patel University, Vallabh Vidyanagar, Gujarat-388120, India
2. Gujarat Agricultural Biotechnology Institute, Navsari Agricultural University, Sana, Gujarat- 395007, India

通信作者: sp389@psu.ac.in | 作者

计算分子生物学, 2014年, 第3卷, 第12期 | doi: 10.51707/cmhb.v3i12.0012

本文档以英文发表, 符合 Computational Molecular Biology 1.0 规范, 版权所有人授权的许可协议, 采用 Creative Commons Attribution License 协议进行授权。如中文再次发表与传播, 只要对原内容有恰当的引用, 版权所有人允许并同意第三方可无条件的使用与传播。如果您对中文文档理解有困难, 请联系英文文档。

建议最佳引用格式:

Patel et al., 2014. Comparative study of five Legume species based on De Novo Sequence Assembly and Annotation. Computational Molecular Biology, Vol. 3, No. 12, doi: 10.51707/cmhb.v3i12.0012

摘要 豆科植物是世界热带和亚热带地区的一种重要的油料作物。最近, 名为 RNA-seq 的新一代测序技术为转录组分析提供了强有力的方法。这项研究是集中在 RNA 序列对五种豆科植物, 分别是来自 NCBI 数据库的花生 *SRRI22866*, 鹰嘴豆 *SRRI27764*, 菜豆 *SRRI28304*, 豌豆 *SRR061197* 和豌豆 *SRR019101*。比较研究侧重于各种重要特征如: 用 NSD, 序列组功能注释产生 reads, 用已知的蛋白质和基因进一步搜索; 其中, 许多基因是根据 GO 功能分类并注释过 KEGG 数据库和序列注释快速集中, 这些数据库将有助于基因发现和功能研究, 并且在当前研究中报道的大多数物种特征与这五种豆科植物的有价值的遗传资源。

关键词 De Novo 装配; 生物信息学; 豆科植物; 序列组合和注释

Comparative study of five Legume species based on De Novo Sequence

Assembly and Annotation

Sagar S. Patel¹, Dipti B. Shah¹, HetaKumar J. Panchal¹

1. G. H. Patel Post Graduate Department of Computer Science and Technology, Sardar Patel University, Vallabh Vidyanagar, Gujarat-388120, India
2. Gujarat Agricultural Biotechnology Institute, Navsari Agricultural University, Sana, Gujarat- 395007, India

通信作者: sp389@gmail.com | Author

Abstract Legume species are an important oilseed crop in tropical and subtropical regions of the world. Recently, next-generation sequencing technology, termed RNA-seq, has provided a powerful approach for analyzing the Transcriptome. This study is focus on RNA-seq of five legume species which are *Arachis hypogaea* L. (The peanut) of *SRRI22866*, *Cicer arietinum* L. of *SRRI27764*, *Phaseolus vulgaris* L. of *SRRI28304*, *Trigonella foenum-graecum* L. of *SRR061197* and *Vicia sativa* L. of *SRR019101* from NCBI database. Comparative study focuses on various important features like: reads generated with NSD, sequence assembly contigs which is further searched with known proteins and genes; among these, how many genes were associated with gene ontology (GO) functional categories and sequences mapped to pathways by searching against the Kyoto Encyclopedia of Genes and Genomes pathway database (KEGG). These data will be useful for gene discovery and functional studies and the large number of transcripts reported in the current study will serve as a valuable genetic resource of these five legume species.

Keywords De Novo assembly; Bioinformatics; Legume species; Sequence assembly and annotation

介绍

新一代测序方法——高通量RNA测序(转录组)正越来越多地应用在植物的检测和定量已知

或新型转录物的选择技术。这种转录组分析方法快速和简单的, 因为它不需要cDNA的克隆。这些cDNA的直接测序可以深度产生短reads。测序后, 得到的reads可以组装成基因组规模的转录配置文件。它是一种更全面和有效的方法来测量转录组组成, 获得RNA表达模式, 并发现新的外显子和基因(Montaravi et al., 2008; Wang et al., 2009);

收稿日期: 2014年12月24日

接受日期: 2014年12月24日

发表日期: 2014年12月25日

Copyright © 2014 BioPublisher

Issam Freni Shengwuac | Vol.3 | No.12 | 1-7

使用各种装配工具, 基因的功能注释和用各种生物信息学工具携带的途径分析来组装转录组的测序数据。 本研究报告的大量转录本可以作为描述五种豆类物种的宝贵的遗传资源。

高通量短read测序是基因组学公开的最新测序技术之一。例如, 在Illumina-基因组分析仪上的平均单次运行可以得到超过3000至4000万个单端序列(~35 nt)。然而, 输出结果可以轻松超过传统Sanger测序的长度设计的基因组分析系统, 甚至获得较小体积数据的454(Roche)测序技术。通常, 短read测序初期使用是局限于与参考基因组几乎相同的基因组数据的低配。全基因组表达水平的转录组分析全是短read测序的理想应用。传统上, 这种分析包括互补DNA(cDNA)文库构建, EST的Sanger测序和微阵列分析。与传统的Sanger方法相比, 新一代测序已经成为增加测序深度和覆盖范围, 同时减少时间和成本的可行方法(L.J Collins et al.)。

1 方法

1.1 序列检索

本研究的重点是五种豆科植物的de novo测序和序列注释, 分别是来自NCBI数据库的花生 *SR1212866*, 鹰嘴豆 *SR627764*, 菜豆 *SR1243064*, 葫芦巴 *SR0066197* 和豌豆 *SR403901*。用于 de novo 转录组分析, 从 Illumina HiSeq 2000 平台和 LS454-454GS FLX 平台来源的NCBI SRA中下载的原始数据(<http://trace.ncbi.nlm.nih.gov/Traces/sra/>)。使用NCBI SRA TOOL KIT将原始序列转化为 fastq 文件格式用于进一步注释 (<http://trace.ncbi.nlm.nih.gov/Traces/sra/trace.cgi?view=software>)。

1.2 NGS QC 工具包

NGS QC工具包, 它是用于高质量数据的质量检查和过滤的应用程序, 此工具包是独立的应用程序, 可从<http://www.nipgr.res.in/ngsquality.html> 免费获得。该工具包包括用Roche 454和Illumina平台生成测序数据的用户容易掌握的QC工具以及用于消除QC(序列格式转换器和修整工具)和分析(统计工具的附加工具, 提供了各种选项能于用户定义的QC参数。该工具包预计对利于较好下游分析的NGS

数据的QC非常有用(Patel RK, et al.)。

1.3 通过 CLC GENOMICS WORKBENCH 运行 De novo 测序

用于分析, 比较和下一代测序数据可视化的全面和用户容易掌握的分析包。这个软件包用 de novo 测序工具的默认参数对序列进行 de novo 测序 (<http://www.cbio.com/products/clc-genomics-workbench/>)。

1.4 BLASTX

对测序文件进一步注释, 第一步是从序列中鉴定出翻译的蛋白质序列。改变几个参数后在NCBI上进行BLASTX比对, 如选择非冗余蛋白质数据库(nr)作为数据库; 双链叶植物在生物选择和算法参数的最大目标序列设置为10, 期望阈值设置为6。

1.5 Blast2GO

Blast2GO(<http://www.blast2go.com/b2gohome>)是一个用于(新)序列的功能注释和注释数据分析的 ALL in ONE工具。基于蛋白质数据库注释的结果, Blast2GO被用于获得基于GO项的ungenes的功能分类, 转录序列根据三个GO term分类, 如分子功能, 细胞过程和生物过程(Nees et al., 2011; Shi et al., 2011; Wang et al., 2010)。用 WEGO(<http://www.wego.genomics.org.cn/>)工具对所有ungenes进行GO功能分类, 并在宏观层面了解该物种的基因功能的分布。用 KEGG 数据库 (<http://www.genome.jp/kegg/pathway.html>)注释这些ungenes的途径。

1.6 SSR 挖掘

我们使用MicroSatellite(MISA)(<http://pgrc.ipk-gatersleben.de/misa>)进行微卫星挖掘, 统计输出其产生有用转录组的信息。

1.7 植物转录因子

PlantTFcat, 在线植物转录因子和转录调节因子分类和分析工具, 用于鉴定序列中的植物转录因子 (<http://plmigrm.nobk.org/PlantTFcat/>)。

2 结果与讨论

2.1 序列比较

(表1)。

2.2 NGS QC 工具包

通过去除接头和其他污染的材料用该工具过滤序列，然后用该工具检查序列的质量，最终用 de novo 序列组的高质量过滤序列文件(表2)。

2.3 De novo 序列组装

CLC GENOMICS WORKBENCH 7 考虑用于 De novo 序列组装，使用默认参数比如 Mismatch Cost = 2, Insertion Cost = 3, Deletion Cost = 3, Length Fraction = 0.5, Similarity Fraction = 0.8, Word size = 21。由本软件产生的序列平均长度和其他细节列于表3。

表 1 物种序列比较

Species	SRR Number	Reads	%GC Content	Platform
<i>Arachis hypogaea</i> L.	SRR1212866	7.3 M spots	48.5	Illumina HiSeq 2000
<i>Cicer arietinum</i> L.	SRR0627764	36 M spots	41.8	Illumina
<i>Phaseolus vulgaris</i> L.	SRR1283084	20.4 M spots	46.4	Illumina HiSeq 2000
<i>Trigonella foenum-graecum</i> L.	SRR066197	627,117 spots	45.2	454 GS FLX
<i>Vicia sativa</i> L.	SRR403901	12.4 M spots	42.4	Illumina HiSeq 2000

表 2 NGS QC 工具包结果

Species	Total number of reads (Original File)	Total number of High Quality reads (High Quality Filter file)	Total number of HQ bases (HQ bases File)	Total number of High Quality bases (Original High Quality Filter file)	Percentage of HQ reads
<i>Arachis hypogaea</i> L.	7300624	7216150	365031200	360807500	98.84%
<i>Cicer arietinum</i> L.	1942297463	1942030113	1060985	1040469	99.99%
<i>Phaseolus vulgaris</i> L.	20444892	13418627	1042589492	684193777	65.63%
<i>Trigonella foenum-graecum</i> L.	627117	609277	146335656	141577237	97.15%
<i>Vicia sativa</i> L.	12427455	12131939	608945295	594465011	97.62%

表 3 序列测量长度

Species	N50	Minimum	Maximum	Average	Conti (Contigs)
<i>Arachis hypogaea</i> L.	448	199	6635	425	10824
<i>Cicer arietinum</i> L.	1239	179	8439	805	34678
<i>Phaseolus vulgaris</i> L.	293	187	5386	302	6999
<i>Trigonella foenum-graecum</i> L.	470	86	3231	445	7256
<i>Vicia sativa</i> L.	588	197	6080	503	22748

2.4 BLASTX 和 Mast2GO 的功能注释

2.4.1 BLASTX 使用 10.6 的 E 值阈值进行 BLASTX 以将序列与非冗余序列数据库比对。BLAST 结果的各种统计信息列于表 4。

2.4.2 酶代码(EC)分类 表 5 是酶序列分类，进一步分为六类，氧化还原酶，转移酶，水解酶，裂解酶，异构酶和连接酶。 2.4.3 基因本体(GO)分类 为了将各种豆科植物的转录序列进行功能分

类, 分配GO terms组成转录序列, 转录序列按照 其分布在分子功能, 生物过程和细胞组分的三个主要GO功能类别进行分组(<http://www.geneontology.org>), 要类别(表6).

表4 BLAST 结果的比较

Species	Without Results	Blat/Without Hits	Blat/With Results	Blat/With Results	Mapping/Unmatched Sequences	Total Sequences
<i>Arachis hypogaea</i> L.	60	688	4789	568	4719	10824
<i>Cicer arietinum</i> L.	3492	3996	25459	786	945	34678
<i>Phaseolus vulgaris</i> L.	102	2601	1988	629	1679	6999
<i>Trigonella foenum-graecum</i> L.	167	2656	1983	192	2258	7256
<i>Vicia sativa</i> L.	0	1114	15482	500	7652	22748

表5 EC分类

Species	Oxidoreductases	Transferases	Hydrolases	Lyases	Isomerase-ses	Ligases	Total
<i>Arachis hypogaea</i> L.	301	614	431	78	51	71	1546
<i>Cicer arietinum</i> L.	51	92	76	20	5	4	248
<i>Phaseolus vulgaris</i> L.	110	232	147	23	20	31	563
<i>Trigonella foenum-graecum</i> L.	148	149	179	34	38	27	575
<i>Vicia sativa</i> L.	429	927	718	80	82	100	2336

表6 GO 分类

Species	Molecular Function	Biological Process	Cellular Components	Total
<i>Arachis hypogaea</i> L.	4512 (4.3%)	2352 (22%)	2467 (23%)	10471
<i>Cicer arietinum</i> L.	916 (40%)	734 (22%)	654 (28%)	2304
<i>Phaseolus vulgaris</i> L.	1727 (47%)	1168 (31%)	829 (22%)	3724
<i>Trigonella foenum-graecum</i> L.	2792 (28%)	4407 (43%)	2980 (29%)	10179
<i>Vicia sativa</i> L.	7026 (37%)	5815 (31%)	5920 (32%)	18761

图1是WEGO工具的输出结果, 它表明, 在分子功能类别中, 编码与催化活性相关的结合蛋白和蛋白质的基因是最富集的, 与代谢过程和细胞过程

相关的蛋白质富集在生物过程类别中。关于细胞成分类别, 细胞和细胞部分是最高度表示的类别, 我们发现所有其他豆类种类相同, 所以我们只引用此图说明WEGO工具。

许多基因在KEGG数据库中用不同的途径注释(<http://www.genome.jp/kegg/pathway.html>), 进一步的比较结果显示在表中。许多转录物包括各种途径, 如代谢途径, 植物-病原体相互作用途径, 脂肪

酸代谢途径和脂肪酸生物合成。

表7 KEGG结果

Species	Genes	KEGG Pathway
<i>Arachis hypogaea</i> L.	568	109
<i>Cicer arietinum</i> L.	786	78
<i>Phaseolus vulgaris</i> L.	629	89
<i>Trigonella foenum-graecum</i> L.	192	87
<i>Vicia sativa</i> L.	500	122

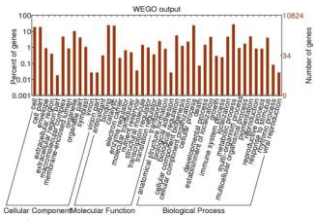


图1 花生 WEGO 工具结果
Figure 1 WEGO Tool Result of *Arachis hypogaea* L.

2.5 SSR 检测 的详细信息。最大部分的SSR是单核苷酸SSR，随后是三核苷酸SSR和二核苷酸SSR。虽然在转录物中只鉴定了一小部分四核苷酸SSR、五核苷酸SSR和六核苷酸SSR，但在大多数物种中该数目是相当显著的。

表8 转录组中鉴定出的 SSR 统计

SSR Mining	Species				
	<i>Arachis hypogaea</i> L.	<i>Cicer arietinum</i> L.	<i>Phaseolus vulgaris</i> L.	<i>Tilganella foeniculum-guacum</i> L.	<i>Vicia sativa</i> L.
Total number of sequences examined:	10824	14678	6999	7256	22748
Total size of examined sequences (bp):	4605095	27932177	2110290	3226271	11444673
Total number of identified SSRs:	742	5228	1405	3107	1150
Number of SSR containing sequences:	649	4391	1304	2191	1055
Number of sequences containing more than one SSR:	74	681	86	747	92
Number of SSRs present in compound formation:	48	337	64	747	48
Distribution to different repeat type classes:					
Mono-nucleotide	265	2019	1218	2589	362
Di-nucleotide	164	1271	87	235	243
Tri-nucleotide	299	1818	90	243	529
Tetra-nucleotide	10	78	7	28	10
Penta-nucleotide	2	17	2	10	3
Hexa-nucleotide	2	25	1	2	3

2.6植物转录因子

此外,通过与已知的转录因子基因家系的序列比较来鉴定转录因子编码转录物。表9中的结果显示,确定了转录因子基因分布在与家族中,并且表9和图2中描述的是Trigonella foenum-graecum L.的植物转录因子结果。转录因子编码转录物在各种已知蛋白质家族中的整体分布与早期预测的其他豆类非常相似(Lhankh et al., 2009)。

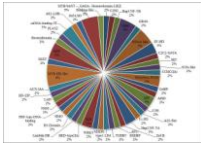


图2 扇形图显示转录因子结果
Figure 2 Plant Transcription Factor Result of Trigonella foenum-graecum L.

表9 植物转录因子结果
Table 9 Plant Transcription Factor Result

Species	At least different families
<i>Arachis hypogaea</i> L.	70
<i>Citrus aurantium</i> L.	97
<i>Phaseolus vulgaris</i> L.	43
<i>Trigonella foenum-graecum</i> L.	45
<i>Vicia sativa</i> L.	82

参考文献

Collins J. L., Biggs J. P., Weckel C. and July S., 2008. An approach to transcriptome analysis of non-model organisms using short-read sequences. *Genome Informatics* 21:3-14
http://dx.doi.org/10.1142/9781848852251_0001

Jian Zhang, Shao Liang, Haidi Duan, Jin Wang, Siqing Chen, Zongshu Cheng, Qiang Zhang, Xuanqing Liang and Yurong Li, 2012. De novo assembly and Characterisation of the Transcriptome during seed development, and generation of gssic-SSR markers in Peanut (*Arachis hypogaea* L.). *BMC Genomics* 2012:13-90
<http://dx.doi.org/10.1186/1471-2164-13-90>

Lhankh, M., Joshi, T., Benedito, V.A., Xu, D., Udwari, M.K., and Sharkey, G., 2009. Legume Transcription Factor Genes: What makes legumes so special? *Plant Physiology* 151: 991-1001
<http://dx.doi.org/10.1104/pp.109.144105>

Motzavi, A., Williams, B.A., McCue, K., Scheffler, L., and Wald, B., 2008. Mapping and quantifying mammalian transcriptomes by RNA-Seq. *Nat Methods* 5(7): 621-8
<http://dx.doi.org/10.1038/nmeth.1226>

Nice, R.W., Solt, M., and Barrett, S.C.H., 2011. De novo sequence assembly and characterization of the floral transcriptome in cross and self-fertilizing plants. *BMC Genomics* 12: 298
<http://dx.doi.org/10.1186/1471-2164-12-298>

Patel RK, Jain M., 2012. NGS QC Toolkit: A Toolkit for Quality Control of Next Generation Sequencing Data. *PLoS ONE* 7(2): e30619. doi:10.1371/journal.pone.0030619
<http://dx.doi.org/10.1371/journal.pone.0030619>

Issam Fenzi Shengwuac. | Vol.3 | No.12 | 1-7

3 结论

本研究侧重于NCBI数据库中五种不同豆科植物的de novo测序和分析。通过四合一[Cillumina和454测序进行RNA-seq分析,转录组测序使得能够对生物体进行各种功能性基因组学研究。虽然已经开发了用于快速测序和表观转录组的几种高通量技术,但是表达的序列数据仍然不能用于许多生物体,包括许多作物植物。在这项研究中,我们对五种不同豆科植物进行了de novo功能注释。不考虑任何参考物种具有显著的非冗余数据集的成绩单没有考虑任何具有显著的非冗余数据集的参考物种。根据五种植物数据的详细分析得到了几个重要的特征如GC含量,豆科植物和其他植物物种的保守基因,通过GO term区分功能类别和通过MISA工具鉴定SSR。值得注意的是,对花生、鹰嘴豆、蚕豆、葫芦巴和豌豆这五种不同豆科植物的比较研究将有利于进一步的基因组研究,因为它包括每个物种完整注释的有用信息。

致谢

感谢 Prof. (Dr) P.V. Vepari, Director, GDCST, Sardar Patel University, Vallabh Vidyamagar 为研究工作提供设备。

- Rohini Garg, Ravi K. Patel, Akhlesh K. Tyagi, and Mahesh Jain. 2011. De Novo Assembly of Chickpea Transcriptome Using Short Reads for Gene Discovery and Marker Identification. *DNA RESEARCH* 18, 53-61; doi:10.1093/dnares/dqr028
<http://dx.doi.org/10.1093/dnares/dqr028>
- Shi, C.Y., Yang, H., and Wei, C.L., 2011. Deep sequencing of the *Caracalla siemensis* transcriptome revealed candidate genes for major metabolic pathways of tea-specific compounds. *BMC Genomics* 12: 131
<http://dx.doi.org/10.1186/1471-2164-12-131>
- Vadva K., Ghosh A., Kumar N, Chandhary S, Sivastava N, Kauria K, Tiwari T and Chikara K., 2012. De novo transcriptome sequencing in *Trigonella foenum-graecum* to identify genes involved in the biosynthesis of dsosmin. *The Plant Genome*
<http://dx.doi.org/10.3835/plantgenome2012.08.0021>
- Wang, X.W., Luo, J.B., Li, J.M., Bao, Y.Y., Zhang, C.X., and Liu, S.S., 2010. De novo characterization of a whitefly transcriptome and analysis of its gene expression during development. *BMC Genomics* 11: 400
<http://dx.doi.org/10.1186/1471-2164-11-400>
- Wang, Z., Gerstein, M., and Snyder, M., 2009. RNA-Seq: a revolutionary tool for transcriptomics. *Nat Rev Genet.* 10(1): 57-63
<http://dx.doi.org/10.1038/nrg2584>