

邓继忠, 林伟森, 甘四明,等. 一种二倍体片段测序中 SNP 检测系统的构建[J]. 华南农业大学学报,2016,37(3):115-120.

## 一种二倍体片段测序中 SNP 检测系统的构建

邓继忠<sup>1</sup>,林伟森<sup>1</sup>,甘四明<sup>2</sup>,黄华盛<sup>1,3</sup>,李梅<sup>2</sup>,金济<sup>1</sup>,何明昊<sup>1</sup>

2 中国林业科学研究院 热带林业研究所,广东 广州 510520; 3 广东科技学院,广东 东莞 523083)

摘要:【目的】开发基于模式识别方法的二倍体片段测序中单核苷酸多态性(Single nucleotide polymorphism, SNP)自动检测系统,提高检测的准确性。【方法】采用 LabWindows/CVI 9.0 开发平台,结合 Matlab 函数库编程,以二倍体 PCR 片段测序的. abl 或. scf 格式文件作为源数据,首先分离出碱基 G、A、T 和 C,进行一维离散小波滤波,再对各碱基处的波形进行典型特征提取,最后运用基于反向传播神经网络的分类器完成 SNP 识别和判断。【结果】系统界面友好、运行稳定。SNP 等级分为 6 级,允许用户对可疑的 SNP 进行人工修正,对尾叶桉 Eucalyptus urophylla 的 26 个测序序列 143 个 SNP 的测试中检测准确率、假阳性率和假阴性率均明显优于之前的类似软件。【结论】本文所构建的 SNP 自动检测系统准确性高,不需参考序列,可用于二倍体 PCR 片段测序中 SNP 的高效检测。

关键词:二倍体; 测序; SNP 检测; 模式识别; 尾叶桉

中图分类号:TP391.4; Q523.8

文献标志码:A

文章编号:1001-411X(2016)03-0115-06

# Development of an automatic system for SNP detection in diploid fragment sequencing

DENG Jizhong<sup>1</sup>, LIN Weisen<sup>1</sup>, GAN Siming<sup>2</sup>, HUANG Huasheng<sup>1,3</sup>, LI Mei<sup>2</sup>, JIN Ji<sup>1</sup>, HE Minghao<sup>1</sup> (1 College of Engineering, South China Agricultural University, Guangzhou 510642, China;

- 2 Research Institute of Tropical Forestry, Chinese Academy of Forestry, Guangzhou 510520, China;
  3 Guangdong University of Science and Technology, Dongguan 523083, China)
- Abstract: [Objective] This study aims to develop a pattern-recognition based system for automatic single nucleotide polymorphism (SNP) detection in diploid fragment sequencing and improve the detection accuracy. [Method] The LabWindows/CVI 9.0 platform and Matlab environment were combined for analyzing abl or sef files generated in diploid PCR fragment sequencing. Firstly, four bases G, A, T and C were separated for eliminating noise through one-dimensional discrete wavelet filtering, following with extraction of typical features of each base position (peak) from a fluorescence curve. A classifier based on back-propagation neural network was then used for SNP recognition and diagnosis. [Result] This established system was characterized by friendly interface, stable operation and manual modification accessibility. It classified the SNP reliability into six grades. Performance test with 143 SNPs of 26 sequencing fragments from Eucalyptus urophylla demonstrated that our system outperformed three previously reported software packages in detecting accuracy, false positive and false negative rates. [Conclusion] Our system has a high rate of accuracy without the need for a reference sequence. It could be used for efficient SNP detection in diploid PCR fragment sequencing.

收稿日期:2015-08-26 优先出版时间:2016-04-15

优先出版网址; http://www.cnki.net/kcms/detail/44.1110.s.20160415.1554.002.html

作者简介:邓继忠(1963—),男,副教授,博士,E-mail: jz-deng@ scau. edu. cn;通信作者:甘四明(1970—),男,研究员,博士, E-mail:Siming\_Gan@ 126. com

基金项目:863 计划(2013AA102705);国家自然科学基金(31270702, 31070592)

**Key words**: diploid; sequencing; SNP detection; pattern recognition; Eucalyptus urophylla

单核苷酸多态性 (Single nucleotide polymorphism, SNP) 是指生物体基因组中单个碱基的颠换或转换导致的 DNA 变异, 是最普遍和最广泛的基因组变异形式, 如拟南芥  $Arabidopsis\ thaliana\ 3$  个品系间存在超过 82 万个 SNP, 平均约 1.4 kb 就存在 1个 SNP<sup>[1]</sup>。SNP 目前已广泛应用于生物、医学和农学等研究领域<sup>[23]</sup>。

传统的 PCR 片段测序中,测序厂家所提供的软 件只能识别各波峰位置的单个碱基,而对二倍体个 体内存在 SNP 的双峰处的低峰碱基不能有效检测。 已有一些第三方软件可用于二倍体物种 PCR 片段测 序的 SNP 自动检测,如 Mutation Surveyor (http:// www. softgenetics. com/MutationSurveyor. html) novoS-NP<sup>[4]</sup>和 PolyPhred<sup>[5]</sup>等。但是,这些软件均需参考序 列,具有局限性,不能有效用于一些序列的分析,如 表达序列标签(Expressed sequence tag, EST)和候选 基因的测序中的内含子区域,而且操作上较为繁琐。 这些软件的检测准确性均偏低[6]。此外, Poly-Phred<sup>[5]</sup>还需要至少8个测序文件以保证足够的准确 性,不利于少量样品的检测。针对二倍体 PCR 片段 测序中如何进行 SNP 的自动、有效检测,本文基于模 式识别的方法构建了无需参考序列进行 SNP 自动检 测的计算机系统,并验证了系统的准确性。

## 1 材料与方法

#### 1.1 系统开发平台

本研究基于 LabWindows/CVI (以下简称 CVI) 平台,以 CVI 9.0 作为系统前台进行主界面的设计和基本功能菜单的实现。CVI 是基于 C 语言的虚拟仪器软件开发平台,其功能强大,具有灵活的交互式编程方法和丰富的用户接口资源<sup>[7-8]</sup>,且运行速度快、界面控件丰富。

采用 MATLAB 2012 函数库编程实现模式识别的分类器构建。MATLAB 具有编程简单、仿真能力强和易于扩展移植等优点<sup>[9]</sup>,提供内部神经网络工具箱,且可外接用户自行开发的模式识别算法或软件包,有利于用户进行各种模式识别算法的测试与建模。本文主要通过 MATLAB 实现模式识别的分类器结构的建立,利用反向传播神经网络(Back-propagation neural network, BPNN)进行权值和阈值的训练,再将训练好的 BPNN 结构移植至 CVI 环境用作SNP 分类器。

系统的运行环境为 Windows XP Professional、 http://xuebao.scau.edu.cn Windows 7 或 Window 8。

#### 1.2 系统功能模块和界面的设计

系统主要包括测序数据导入、碱基分离、噪声处理、SNP识别、数据显示与存储、人工校正等功能模块。测序数据导入可打开扩展名为. abl 或. scf 的测序文件。碱基分离是根据测序数据中 G、A、T 和 C 4种碱基的不同荧光颜色而将它们分别提取出来。噪声处理通过一维离散小波变换滤除噪声,便于后续的数据特征提取。SNP识别主要根据测序峰图的波形特征,结合二倍体内 SNP 表现为双峰的现象,提取各碱基位置的典型波形特征作为 BPNN 的输入向量,输出 SNP属性分数,再参照判别标准给出 SNP等级。数据显示可在运行窗口直观显示测序峰图和序列以及 SNP的位置、双峰碱基、属性分数和等级,识别结果可存为. txt 文件。人工校正允许用户对误判和漏判的 SNP进行删除、添加和更改碱基等手动操作。系统的工作流程见图 1。

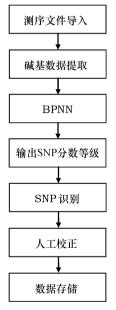


图 1 SNP 检测流程图

Fig. 1 A flow chart of the SNP detection process

系统界面包括主操作区、原始峰图显示区、SNP 识别结果区和人工校正区等4个功能区域(图2)。

主操作区位于系统界面的上部,如图 2 所示,包含面板选项卡、菜单栏、快捷键和信息提示栏。面板选项卡可提供多个检测面板,各检测面板的序列分析数据相互独立,以便进行多个测序文件的检测。菜单栏包括文件、编辑、滤波选择、识别方法和帮助等一级菜单,各一级菜单可下拉显示二级菜单。文件的二级菜单包括原始测序文件导入、SNP 检测结

果保存等。编辑的二级菜单包括 X 轴和 Y 轴缩放以及隐藏/显示碱基 G、A、T 和 G。滤波选择的二级菜单包括小波变换等方法。识别方法的二级菜单包括 BPNN 等多种方法。快捷键工具栏列出常用功能,包括文件打开、SNP 检测(按默认的滤波和模式识别方法)、增加面板、关闭当前面板、峰图和结果显示窗口 X 轴和 Y 轴缩放等。信息提示栏显示文件路径等。

原始峰图显示区位于界面右中区域,见图 2,不同颜色代表碱基 G、A、T 和 C 显示测序的峰图,可通

过主操作区的快捷按钮选择性地显示一种或几种碱基类型,也可进行 X 轴和 Y 轴的缩放。

SNP识别结果区位于界面下部区域,见图 2,界面右下区域显示 SNP识别结果,包括碱基位置、碱基类型、双峰 SNP、SNP分数和 SNP等级等。界面左下区域列表总结 SNP检测的结果,点击某一行可与右下窗口相应位置的 SNP进行快速定位。

人工校正区位于界面左中区域(图 2),允许用户对任一位置碱基进行更改、删除和增加等操作。

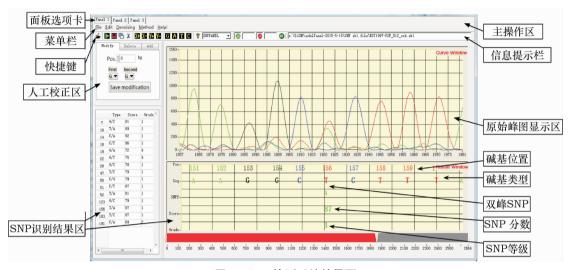


图 2 SNP 检测系统的界面

Fig. 2 A typical interface of the SNP detection system

#### 1.3 系统主要功能的实现

1.3.1 测序数据的导入和碱基分离 本系统可导入. ab1 或. scf 格式的测序文件,根据文件内容分离出各种碱基。以. ab1 文件为例说明。按照美国 Applied Biosystems 公司的 ab1f 文件格式, . ab1 文件是由一个指定的文件夹存储文件信息,包括文件名称、碱基类型和碱基数量等。文件 7 ~ 34 字节是'DirEntry'结构体,指向文件夹所在的位置,其结构为:

Struct DirEntry { sInt32 name; sInt32 number; sInt16 elementtype; sInt16 elementsize; sInt32 numelements; sInt32 datasize; sInt32 dataoffset; sInt32 datahandle; },

其中,number 显示通道编号(对应 4 种碱基), elementsize 显示碱基所占空间大小,numelements 显示碱基数量,dataoffset 显示碱基位置的偏移量。根据这些信息即可将各通道的碱基和荧光数据提取出来。

1.3.2 滤波的实现 分离出的碱基荧光信号通常 含有噪声,影响 SNP 分析的准确性,可通过小波分解 进行噪声滤除。小波分解具有良好的时频分析特 性,在时域和频域内都具有突显信号局部特征的能力,并且能够在多种分辨率下观测数据,是噪声滤除的有效方法<sup>[10-11]</sup>。

本系统调用了 CVI 中 Signal Processing Toolkit 7.0.2 工具箱所含的小波工具包的函数,一维离散小波的滤波过程如下:分离 4 个通道的碱基数据;读取各通道的碱基数据长度,获得滤波参数;调用 SptAllocCoeffWFBD 函数,分配小波滤波工具包的滤波参数结构;调用 SptReadCoeffWFBD 函数,读取小波工具包参数;调用 SptDiscreteWaveletTransform (filterData, dataLenght, analysisFilter, SptSymmetric, NULL, NULL, scales, shift, dwt, &outputSize, lengths)函数,进行一维离散小波分解,滤除后的数据存储于变量 dwt 中,用于后续分析。

1.3.3 SNP 识别分类器的构建 利用 BPNN 构建 SNP 识别分类器。BPNN 是一种反向传播且能修正 误差的多层映射模式,由输入层、隐含层和输出层及 各层神经元之间连接组成<sup>[12]</sup>。SNP 识别可以视作从 波形特征输入到 SNP 类型输出的非线性映射问题。

基于文献[13],选择了SNP处双峰的波峰距离、 高度比值和起伏度比值作为特征参数,归一化处理 http://xuebao.scau.edu.cn 后作为 BPNN 输入层的特征向量,即输入层的神经元数为3。BPNN 输出层为 SNP 类别的诊断结果,可能是双峰 SNP 位点或单峰非 SNP 位点,即输出层神经元为1;输出层中设定范围0~100为1个位点的属性分数,并结合属性分数和周边杂峰将 SNP 可能性分为6个等级,等级越小越可能为 SNP。

MATLAB 神经网络工具箱中调用 BPNN 构建 SNP 分类器的过程具体如下:调用 net = newff[input, output, 10, ('tansig', 'purelin'), 'trainlm'],初始建立适于 SNP 识别的 BPNN 分类器;调用 net. trainParam. goal = 0.01,设置误差,当网络精度达到误差则停止训练,训练好的权值和阈值存于 MATLAB 中 net 结构体的 net. IW、net. LW 和 net. b 元素;调用 net = train (net, input, output), CVI 中重新塑造 BPNN 即可用于 SNP 检测。

#### 1.4 数据处理

因 SNP 检测可能存在假阳性和假阴性,人工校正的碱基删除、添加和更改的结果可直接显示在当

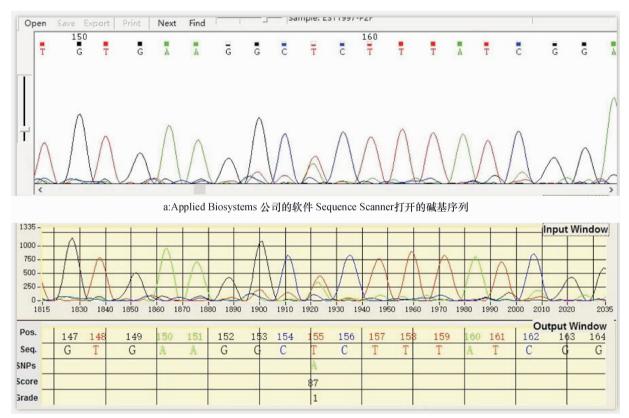
前面板。SNP 检测结果的保存格式为. txt 文件。

利用尾叶桉 Eucalyptus urophylla 26 个测序文件的 143 个 SNP 进行系统测试,并与 Mutation Surveyor (http://www.softgenetics.com/MutationSurveyor.html)、novoSNP<sup>[4]</sup>和 PolyPhred<sup>[5]</sup>的结果进行比较,分析软件的有效性。所用 26 个测序文件为尾叶桉SNP 开发的 EST 重测序文件<sup>[14]</sup>。

## 2 结果与分析

#### 2.1 系统主要功能的实现

系统导入测序文件后能正确地进行碱基判读, 判读结果与测序结果较为一致。图 3 显示了本系统 对 1 个序列的判读与测序结果一致。但对测序质量 较差或者杂峰较多的区域,测序结果亦不准确,与本 系统的判读结果可能有异。图 3B 表明本系统能可 靠地识别双峰碱基位置的 SNP,155 bp 处的 SNP 判 读为 T/A,分值为87,等级为1。并且,其他较低的杂 峰被有效排除,未被误判为 SNP。



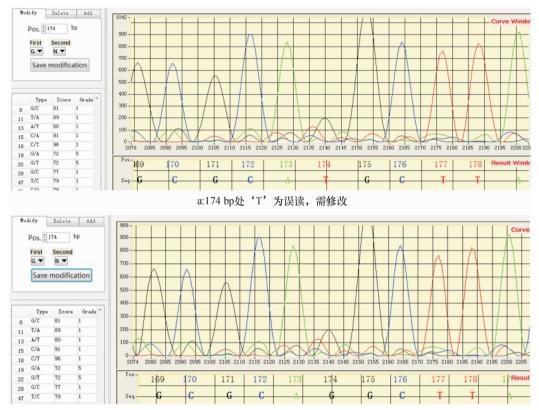
b:本系统判读的序列,155 bp 处双峰、SNP 为 T/A

图 3 本系统导入测序文件后的判读序列与测序比较

Fig. 3 Base sequence obtained in our system being compared to that of sequencer

人工校正可对特定位置的第 1 和/或第 2 碱基进行修改、删除和添加操作。图 4 显示了对 174 bp处误读的'T'修改为'G',表明系统可有效地进行人工校正。

数据存储的. txt 文件见图 5,第 1 行为文件名,数字代表显示碱基位置,识别序列上面一行为碱基序列,下面一行显示识别的 SNP 的第 2 碱基。测序前 20 个碱基一般是较杂乱的峰, SNP 可靠性低。



b:校正后 174 bp 处更正为'G'

图 4 人工校正的有效性

Fig. 4 Efficient manual correction of a misidentified base

2.2 SNP 检测的有效性

对于测试的尾叶桉 26 个测序文件,对比另外 3

个软件,本系统具有较高的准确性和误检率。表1

显示了不同软件在高和低的标准下对 143 个 SNP 的识别准确率、漏判(假阴性)率和误判(假阳性)

率,不管是在高还是低的标准下,本文所开发的系

统均有最高的 SNP 判读准确率以及最低的假阴性

率和假阳性率, PolyPhred 除外, 因其未检测到任何

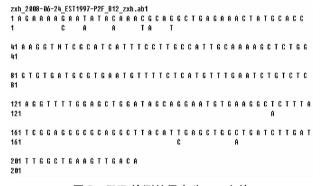


图 5 SNP 检测结果存为. txt 文件

Fig. 5 SNP detection result saved as a . txt file

表 1 不同软件识别 SNP 的结果对比 $^{1)}$ 

 $SNP_{\circ}$ 

Tab. 1 Comparison of software performance in SNP detection

判别	SNP	DiSNPindel		novoSNP		Mutation Surveyor		PolyPhred	
标准	准确性	SNP 数量	判别率/%	SNP 数量	判别率/%	SNP 数量	判别率/%	SNP 数量	判别率/%
高	准确	138	96.5	12	8.4	30	21.0	0	0
	漏判	5	3.5	131	91.6	113	79.0	143	100
	误判	45	24.6	17	58.6	31	50.8	0	0
低	准确	141	98.6	34	23.8	31	21.7	0	0
	漏判	2	1.4	109	76.2	114	78.3	143	100
	误判	75	42.9	89	72.4	31	50.0	0	0

<sup>1)</sup> DiSNPindel、novoSNP、Mutation Surveyor 和 PolyPhred 高的标准分别为等级 2、分值 13、中敏感性和等级  $2^{[6]}$ ,低的标准分别为等级 4、分值 6、高敏感性和等级 4 准确 SNP 的判别率 = 准确 SNP 数量/(准确 SNP 数量 + 漏判 SNP 数量)/×100%,漏判 SNP 的百分比 = 漏判 SNP 数量/(准确 SNP 数量 + 漏判 SNP 数量/(误判 SNP 数量 + 准确 SNP 数量)×100%。

## 3 结论与讨论

本文基于 CVI 9.0 开发平台,结合 MATLAB 函数库,运用特征提取和模式识别等技术构建了二倍体片段测序的 SNP 自动检测系统。本系统开发中结合了 CVI 运行速度快和用户接口资源丰富以及MATLAB 函数库多样和用户可自行扩充函数等优点,降低了系统开发的难度,并且检测系统界面友好、易于扩充。系统各模块功能的完整实现和测试的有效性也证明相关技术的合理运用。

高通量的新一代测序技术已经兴起,但传统的PCR 片段测序仍非常重要<sup>[15]</sup>。测序厂家提供的软件只能识别各序列位置的最高峰所对应的碱基,因此双峰位置的低峰碱基的识别需要第三方软件。另外,测序也常出现杂峰,且杂峰可能高于碱基峰而导致测序厂家的软件不能正确读序。本文针对这一问题研发了SNP自动检测系统,可以有效实现双峰判读和杂峰滤除等功能,为二倍体PCR 片段测序的SNP 判读提供了有力工具。

与功能类似的软件 Mutation Surveyor(http://www.softgenetics.com/MutationSurveyor.html)、novoS-NP<sup>[4]</sup>和 PolyPhred<sup>[5]</sup>相比,本系统不需参考序列,既为EST 和候选基因的测序中内含子无参考序列的 SNP 检测提供了便利,也减少了输入参考序列的繁琐操作;并且,本文所开发的检测系统具有最高的 SNP 判读准确率以及最低的假阴性率和假阳性率,可用于二倍体PCR 片段测序中 SNP 的高效检测。

实际应用中, Mutation Surveyor 和 novoSNP 可能假阳性率极高,如对15 个基因171 个 SNP 的检测中的假阳性 SNP 分别为 3 728 和 505 个,显著多于正确判读的 SNP<sup>[16]</sup>。另一方面,任何检测系统的准确率都严重依赖于测序质量,实际应用时需要优化测序体系<sup>[4]</sup>。

#### 参考文献:

- [1] OSSOWSKI S, SCHNEEBERGER K, CLARK R M, et al. Sequencing of natural strains of *Arabidopsis thaliana* wit short reads[J]. Genome Res, 2008,18(12):2024-2033.
- [2] 唐立琼,肖层林,王伟平. SNP 分子标记的研究及其应用进展[J]. 中国农学通报, 2012,28(12):154-158.

- [3] 许家磊,王宇,后猛,等. SNP 检测方法的研究进展[J]. 分子植物育种, 2015,13(2):475-482.
- [4] WECKX S, DEL-FAVERO J, RADEMAKERS R, et al. no-voSNP, a novel computational tool for sequence variation discovery [J]. Genome Res, 2005, 15(3):436-442.
- [5] MATTHEW S, JAMES S, ROBERTSON P D, et al. Automating sequence-based detection and genotyping of single nucleotide polymorphisms (SNPs) from diploid samples [J]. Nat Genet, 2006,38(3):375-381.
- [6] DENG J Z, HUANG H S, YU X L, et al. DiSNPindel: Improved intra-individual SNP and InDel detection in direct amplicon sequencing of a diploid[J]. BMC Bioinformatics, 2015,16:343.
- [7] 仇志平,李树军. LabWindows/CVI 虚拟仪器软件在测试 领域中的应用[J]. 计算机工程与设计,2007,28(22): 5544-5548.
- [8] 刘君华. 虚拟程序编程语言 LabWindows/CVI 编程[M]. 北京:电子工业出版社,2001.
- [9] 肖伟,刘忠,曾新勇,等. MATLAB 程序设计与应用[M]. 北京:清华大学出版社,2005.
- [10] BUI T D, CHEN G. Translation-invariant denoising using multiwavelets [J]. IEEE Trans Sig Proc, 1998, 46 (12): 3414-3420.
- [11] PAN Q,ZHANG P,DAI G, et al. Two denoising methods by wavelet transform[J]. IEEE Trans Sig Proc,1999,47(12): 3401-3406.
- [12] MCKEOWN J J, STELLA F, HALL G. Some numerical aspects of the training problem for feed-forward neural nets [J]. Neural Netw, 1997, 10(9):1455-1463.
- [13] 黄华盛. 基于模式识别的二倍体个体内 SNP 和 InDel 自动检测[D]. 广州:华南农业大学,2014.
- [14] YU X,GUO Y,ZHANG X, et al. Integration of EST-CAPS markers into genetic maps of and *Eucalgptus urophylla* and *E. tereticornis* and their alignment with *E. grandis* genome sequence [J]. Silvae Genet, 2012,61(6);247-255.
- [15] STUDER A, ZHAO Q, ROSS-IBARRA J, et al. Identification of a functional transposon insertion in the maize domestication gene tb1 [J]. Nat Genet, 2011, 43 (11):1160-1163.
- [ 16 ] NGAMPHIW C, KULAWONGANUNCHAI S, ASSAWA-MAKIN A, et al. VarDetect: A nucleotide sequence variation exploratory tool [ J ]. BMC Bioinformatics, 2008, 9 (S12):9.

【责任编辑 霍 欢】