

学校编码: 10384

分类号 \_\_\_\_\_ 密级 \_\_\_\_\_

学号: 23220071152866

UDC \_\_\_\_\_

厦门大学

硕士学位论文

基于 Tiling Array 的植物基因结构分段与  
表达差异分析

Segmentation and Expression Analysis of Plant Tiling  
Array Data

武姗姗

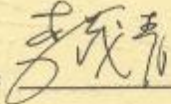
指导教师姓名: 吉国力 教授

专业名称: 系统工程

论文提交时间: 2010 年 5 月

论文答辩日期: 2010 年 月

学位授予日期: 2010 年 月

答辩委员会主席: 

评阅人: \_\_\_\_\_

2010 年 5 月

厦门大学博硕士学位论文摘要库

## 厦门大学学位论文原创性声明

本人呈交的学位论文是本人在导师指导下,独立完成的研究成果。本人在论文写作中参考其他个人或集体已经发表的研究成果,均在文中以适当方式明确标明,并符合法律规范和《厦门大学研究生学术活动规范(试行)》。

另外,该学位论文为( )课题(组)的研究成果,获得( )课题(组)经费或实验室的资助,在( )实验室完成。(请在以上括号内填写课题或课题组负责人或实验室名称,未有此项声明内容的,可以不作特别声明。)

声明人(签名): 武世明

2010年5月30日

厦门大学博硕士学位论文摘要库

## 厦门大学学位论文著作权使用声明

本人同意厦门大学根据《中华人民共和国学位条例暂行实施办法》等规定保留和使用此学位论文，并向主管部门或其指定机构送交学位论文（包括纸质版和电子版），允许学位论文进入厦门大学图书馆及其数据库被查阅、借阅。本人同意厦门大学将学位论文加入全国博士、硕士学位论文共建单位数据库进行检索，将学位论文的标题和摘要汇编出版，采用影印、缩印或者其它方式合理复制学位论文。

本学位论文属于：

1. 经厦门大学保密委员会审查核定的保密学位论文，  
于 年 月 日解密，解密后适用上述授权。

2. 不保密，适用上述授权。

（请在以上相应括号内打“√”或填上相应内容。保密学位论文应是已经厦门大学保密委员会审定过的学位论文，未经厦门大学保密委员会审定的学位论文均为公开学位论文。此声明栏不填写的，默认为公开学位论文，均适用上述授权。）

声明人（签名）：武珊珊

2010年5月30日

## 摘要

在对植物基因芯片数据分析的相关研究领域中，基于 Tiling Array 芯片技术的植物基因结构分段和基因表达差异分析是当前生物信息学的研究热点之一。作为迄今为止分辨率最高的基因芯片类型，Tiling Array 芯片技术已经广泛应用于生物学和医学的各个研究领域，成为目前研究人类基因组和其他各种模式生物基因组复杂性的最强有力的工具。

本文是基于实验室课题组与美国迈阿密大学李教授 (Dr.Q. Quinn Li) 实验室的项目合作对植物选择性多聚腺苷酸化相关问题的研究成果，应用计算机和数学算法对植物基因结构分段和基因表达差异问题进行分析和研究。论文首先探讨了 Tiling Array 芯片的设计模式及其独特的 PM-MM 探针选取方法以及在此基础上发展而来的一些 Tiling Array 经典信号识别算法。然后针对本课题的实际情况，经过合理的实验设计、试验流程，设计了一个完整的基于拟南芥 Tiling Array 芯片数据分析的数据处理流程，一步一步描述了从拿到芯片原始探针浓度数据后应该如何进行数据处理。应用了包括 R 语言、Bioconductor 平台、matlab 生物信息工具箱等一系列适合分析 Tiling Array 芯片数据的计算机分析平台、软件和数学算法，对拟南芥 Tiling Array 芯片进行了数据处理，特别是在拟南芥数据预处理和标准化方面做了重点介绍，如为消除非特异性杂交干扰影响的 VSN 背景校正、为消除不同样本之间系统偏差所进行的 Quantile 芯片间归一化、DNA 参考标准化、为从生物学角度上更好的解释及使数据满足特定分布的数据转换、探针过滤等步骤，并给出了对应算法所进行的芯片数据处理结果。最后通过 R、SQL Server 数据库将构造后的探针数据导入到 Matlab 分段识别算法中进行分段识别并运用统计学分析方法进行差异表达基因的挑选。实验结果表明，本文基于 Tiling Array 芯片数据的拟南芥基因分段和基因表达差异分析方法是可行和有效的。

**关键词：**Tiling Array；拟南芥；基因分段；差异分析

厦门大学博硕士学位论文摘要库

## Abstract

Segmentation and expression analysis of plant Tiling Array data is an interesting but challenging problem. As an important branch of gene chips, which has the highest resolution of all gene chips by far, Tiling Array chip technology has been widely used in all research areas of biology and medicine, and is the most powerful tool for the study of the human genome and the complexity of other genomes of various model organisms.

Cooperating with Dr.Q.Quinn Li at the Department of Botany, Miami University, and based on the research results of differentially expressed gene, this thesis is to research on the problem of segmentation and differentially expression gene using computer and mathematic algorithm. In this article, we describe a whole data processing framework based on *Arabidopsis thaliana* Tiling Array, introducing how to deal with the raw Tiling Array probe data. First we discussed the Tiling Array chip design patterns and unique PM-MM probe selection method, and some classical signal recognition algorithm of Tiling Array. Then according to our actual condition, we designed a analysis procedure for *Arabidopsis thaliana* Tiling Array data, including the application of a series of computer analysis softwares, such as R, Bioconductor platform, matlab bioinformatics toolbox, and applied them to the actual array data to do VSN background correction to eliminate the interference effects of non-specific hybridization, Quantile normalization to eliminate systematic bias between different samples, DNA reference normalization, Log transformation and probe filtering and gives the corresponding results of each algorithm. Finally, we conducted a segmentation analysis of the data, selecte the differentially expressed genes, using the sub-recognition algorithms and statistical analysis, and then compared different algorithms' results. Experimental results showed that the methods we used on *Arabidopsis* Tiling Array data are feasible and effective.

**Key Words:** Tiling Array; *Arabidopsis Thaliana*; Segmetation; Expression Analysis



厦门大学博硕士学位论文摘要库

# 目 录

<b>第一章 绪论</b> .....	<b>1</b>
1.1 引言 .....	1
1.2 本文研究背景及意义 .....	2
1.3 本文结构框架 .....	4
<b>第二章 Tiling Array 芯片技术与应用</b> .....	<b>7</b>
2.1 基因芯片技术概述 .....	7
2.1.1 基因芯片的产生和发展 .....	7
2.1.2 基因芯片基本原理 .....	8
2.1.3 基因芯片相关技术 .....	9
2.2 Tiling Array 芯片介绍 .....	10
2.2.1 Tiling Array 芯片基本概念 .....	10
2.2.2 Tiling Array 芯片设计原理 .....	11
2.2.3 Tiling Array 芯片实际应用研究 .....	13
2.3 Tiling Array 经典信号识别算法 .....	14
2.3.1 滑窗算法 (sliding window, SW) .....	14
2.3.2 基于 HMM 的信号识别算法 .....	14
2.3.3 支持向量机 (SVM) .....	16
2.3.4 三种算法比较及结论 .....	17
2.3.5 隐马尔可夫-支持向量机算法 (HM-SVM) .....	18
2.4 Tiling Array 主要数据分析软件 .....	19
2.4.1 Affymetrix GCOS 系统 .....	19
2.4.2 R 平台及 Bioconductor .....	20
2.4.3 MATLAB 生物信息工具箱 .....	21
<b>第三章 拟南芥基因结构分段与表达差异分析</b> .....	<b>23</b>
3.1 实验对象及数据文件 .....	23

3.1.1	模式植物拟南芥 .....	23
3.1.2	本文相关数据文件 .....	24
<b>3.2</b>	<b>Tiling Array 芯片数据预处理 .....</b>	<b>26</b>
3.2.1	背景校正 .....	28
3.2.2	芯片间标准化 .....	30
3.2.3	DNA 参考标准化 .....	32
3.2.4	数据的对数转换 .....	34
3.2.5	探针过滤 .....	35
<b>3.3</b>	<b>基因结构分段分析 .....</b>	<b>39</b>
3.3.1	实验目的 .....	39
3.3.2	实验预期 .....	40
<b>3.4</b>	<b>表达差异基因挑选及分析 .....</b>	<b>41</b>
3.4.1	研究差异表达基因的意义 .....	41
3.4.2	算法介绍 .....	41
<b>第四章</b>	<b>实验结果分析与讨论 .....</b>	<b>45</b>
4.1	环境搭建及构造变量 .....	45
4.2	输出结果与分析 .....	48
4.1.1	分段结果显示 .....	48
4.1.2	差异表达基因 .....	50
<b>第五章</b>	<b>总结与展望 .....</b>	<b>55</b>
5.1	全文总结 .....	55
5.2	后续展望 .....	56
<b>附 录</b>	<b>.....</b>	<b>57</b>
<b>参考文献</b>	<b>.....</b>	<b>59</b>
<b>致 谢</b>	<b>.....</b>	<b>63</b>

# Contents

<b>Chapter 1: Introduction .....</b>	<b>1</b>
<b>1.1 Introduction .....</b>	<b>1</b>
<b>1.2 Research Background and Significance.....</b>	<b>2</b>
<b>1.3 Structure .....</b>	<b>4</b>
<b>Chapter 2: Technology and Application of Tiling Array .....</b>	<b>7</b>
<b>2.1 Overview of Gene Chip Technology .....</b>	<b>7</b>
2.1.1 Emergence and Development of Gene Chips.....	7
2.1.2 Basic Principles of Gene Chips .....	8
2.1.3 Related Technologies .....	9
<b>2.2 Introduction of Tiling Array.....</b>	<b>10</b>
2.2.1 Concept of Tiling Array .....	10
2.2.2 Design Principles of Tiling Array .....	11
2.2.3 Application of Tiling Array .....	13
<b>2.3 Classical Signal Recognition Algorithm of Tiling Array .....</b>	<b>14</b>
2.3.1 Sliding Window Algorithm .....	14
2.3.2 HMM Algorithm.....	14
2.3.3 SVM Algorithm .....	16
2.3.4 Comparison of the three Algorithm .....	17
2.3.5 HM-SVM Algorithm.....	18
<b>2.4 Major Data Analysis Software of Tiling Array .....</b>	<b>19</b>
2.4.1 Affymetrix GCOS System .....	19
2.4.2 R and Bioconductor .....	20
2.4.3 MATLAB Bioinformatics Toolbox.....	21
<b>Chapter 3: Segmentation and Expression Analysis of Arabidopsis Thaliana.....</b>	<b>23</b>

<b>3.1</b>	<b>The Sources of Samples and Data Files .....</b>	<b>23</b>
3.1.1	Arabidopsis thaliana.....	23
3.1.2	Related Data Files.....	24
<b>3.2</b>	<b>Pre-processing of Tiling Array .....</b>	<b>26</b>
3.2.1	Background Correction.....	28
3.2.2	Normalization Between Arrays .....	30
3.2.3	DNA Reference Normalization .....	32
3.2.4	Log2 Transformation .....	34
3.2.5	Probe Selection.....	35
<b>3.3</b>	<b>Gene Structure Analysis .....</b>	<b>39</b>
3.3.1	Experiment Purpose.....	39
3.3.2	Expected result .....	40
<b>3.4</b>	<b>Differential Expression Gene Analysis.....</b>	<b>41</b>
3.4.1	Research Significance.....	41
3.4.2	Algorithm Introduction .....	41
<b>Chapter 4: Experiment Result and Analysis .....</b>		<b>45</b>
4.1	Environment and Variables.....	45
4.2	Results and Analysis .....	48
4.2.1	Segmentation results .....	48
4.2.2	Differential Expression Gene .....	50
<b>Chapter 5: Conclusion and Expectation.....</b>		<b>55</b>
5.1	Conclusion of the Whole Thesis .....	55
5.2	Expectation of Problem .....	56
<b>Appendix .....</b>		<b>57</b>
<b>Reference .....</b>		<b>59</b>
<b>Acknowledgement.....</b>		<b>63</b>

厦门大学博硕士学位论文摘要库

# 第一章 绪论

## 1.1 引言

我们正处在一个激动人心的时代，科技的进步已使人类可以窥视生命的秘密。2000年6月26日，美、英、日、法、德、中6个国家16个研究中心联合宣布完成覆盖人的大部分基因组、准确率超过90%的DNA序列人类基因组“工作框架图”，这部由30亿个字符组成的人类遗传密码本已活生生的摆在了我们面前。于此同时，来自其它生物的基因组信息源源不断从自动测序仪中涌出，堆集如山，浩如烟海。这些海量的生物信息是用特殊的“遗传语言”——DNA的四个碱基字符（A、T、G和C）和蛋白质的20个氨基酸字符一一写成。从这一时刻开始，人类历史进入了一个崭新的时代——后基因组（post-genome）时代<sup>[1]</sup>。

原始的生物信息资源挖掘出来后，生命科学工作者面临着严峻的挑战。数以亿计的ACGT序列中包涵着什么信息？基因组中的这些信息怎样控制有机体的发育？基因组本身又是怎样进化的？人们迫切希望可以研究出这些分子数据所具有的丰富内涵，以及其背后隐藏着的尚不为人所知道的生物学知识，这就形成了以基因组为主要研究对象的生物信息学<sup>[2][3]</sup>。《科学》（Science）在2001年2月16日人类基因组专刊上配发了一篇题为“生物信息学：努力在数据的海洋里畅游”的文章，文章写道：“我们身处急速上涨的数据海洋中...，我们如何避免生物信息的没顶之灾呢？一叶轻舟也许可以救命！生物信息学便是我们找到的这样一条“轻舟”。作为一门年青的学科，它充满挑战、机遇且引人入胜。”

一般意义上，生物信息学是研究生物信息的采集、处理、存储、传播、分析和解释等各方面的一门学科，它通过综合利用生物学、计算机科学和信息技术而揭示大量而复杂的生物数据所赋有的生物学奥秘<sup>[4]</sup>。具体而言，生物信息学作为一门新的学科领域，它是把基因组DNA序列信息分析作为源头，在获得蛋白质编码区的信息后进行蛋白质空间结构模拟和预测，然后依据特定蛋白质的功能进行必要的药物设计。基因组信息学、蛋白质空间结构模拟以及药物设计构成了生物信息学的3个重要组成部分。

最初生物信息学仅用于数据的贮存，关注于存储来自基因组测序计划完成的序列数据。随着人类基因组计划的完成，生物信息学也被赋予了更多的内容，人们开始感兴趣如何利用计算方法分析生物数据，如根据核酸序列预测基因结构、功能的算法等<sup>[5]</sup>。生物信息学的研究重点从以前的以 DNA 测序转移到了系统了解基因组内所有基因的生物学功能，即功能基因组（functional genomics）。研究内容主要包括基因功能发现、基因表达分析及突变检测<sup>[6]</sup>。

传统研究功能基因的方法包括减法杂交，差示筛选，cDNA 代表差异分析以及 mRNA 差异显示等，但这些技术不能对基因进行全面系统的分析。对于人类基因组，曾经估计编码蛋白质的基因数大约为 2 万个。除此之外，还有许多不编码蛋白质的 RNA 基因，如 rRNA, tRNA, microRNA, snoRNA 等。全长 cDNA 测序（cDNA sequencing）方法识别了目前已知、高质量的编码蛋白基因，但几乎所有类似的克隆方法都偏向于探测在生物组织中被充分表达的基因。这类技术往往很难深入地探索基因组水平所有的表达信息，也不适用于不同的组织和不同条件下转录信息的提取。在刚刚过去的 5 年间，从微阵列（microarray）基因芯片技术发展而来的新技术 Tiling Array 使高通量、全基因组水平的表达探测在理论上得以开展。随着芯片探针设计密度的不断增加，Tiling Array 在高等真核生物全基因组水平的应用也得以逐渐实现。

## 1.2 本文研究背景及意义

微阵列（microarray）或芯片（chip）技术是近年发展起来的一种新的分子生物学研究工具<sup>[7]</sup>。它利用光导化学合成、照相平板印刷以及固相表面化学合成等技术，在固相表面合成成千上万个寡核苷酸探针，与放射性同位素或荧光物标记的 DNA 或 cDNA 杂交，用于分析 DNA 突变及多态性、DNA 测序、监测同一组织细胞在不同状态下或同一状态下多种组织细胞基因表达水平的差异、发现新的致病基因或疾病相关基因等多个研究领域。

Tiling Array 实验技术，中文译为“嵌合芯片”，是从传统的微阵列（microarray）技术发展而来，是迄今为止分辨率最高的基因芯片类型，Tiling Array 的探针设计几乎涵盖了目标 DNA 的全部序列<sup>[8]</sup>。



Degree papers are in the "[Xiamen University Electronic Theses and Dissertations Database](#)". Full texts are available in the following ways:

1. If your library is a CALIS member libraries, please log on <http://etd.calis.edu.cn/> and submit requests online, or consult the interlibrary loan department in your library.
2. For users of non-CALIS member libraries, please mail to [etd@xmu.edu.cn](mailto:etd@xmu.edu.cn) for delivery details.

厦门大学博硕士论文摘要库