
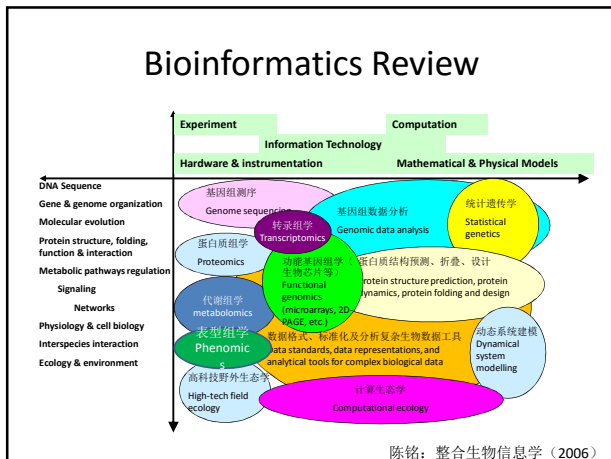
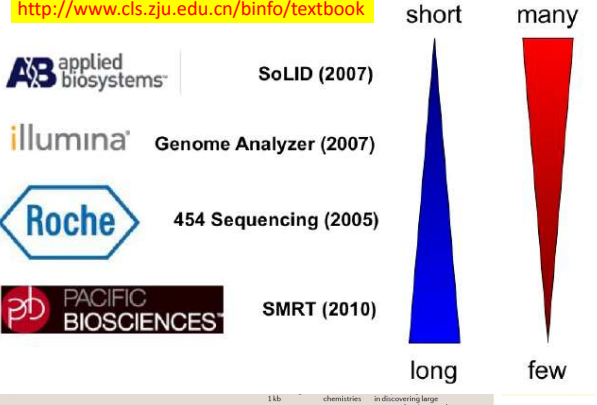
 **Bioinformatics**
生物信息学

2. NGS & Transcriptomics

陈铭 (mchen@zju.edu.cn)
2015年9月21日

<http://www.cls.zju.edu.cn/binfo/textbook>



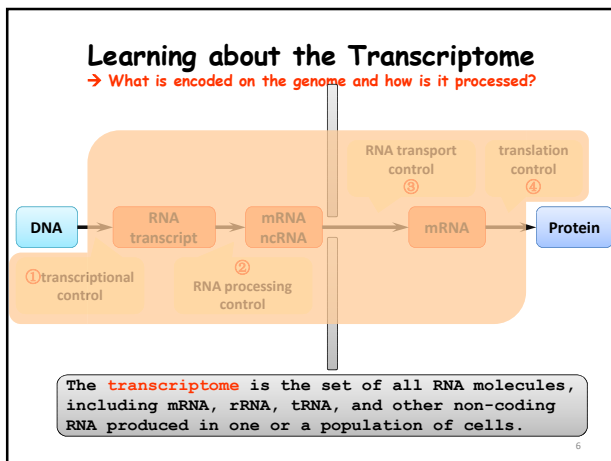
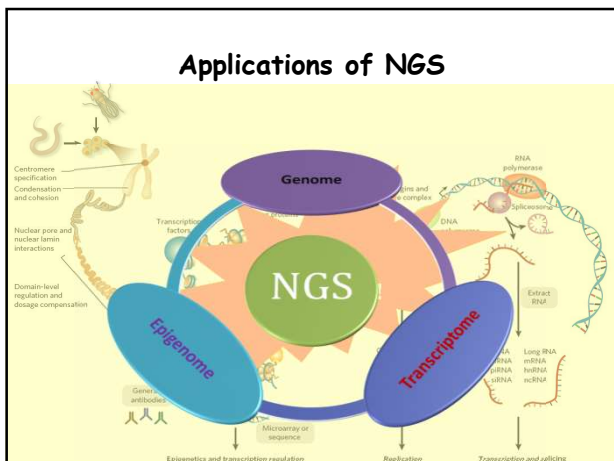
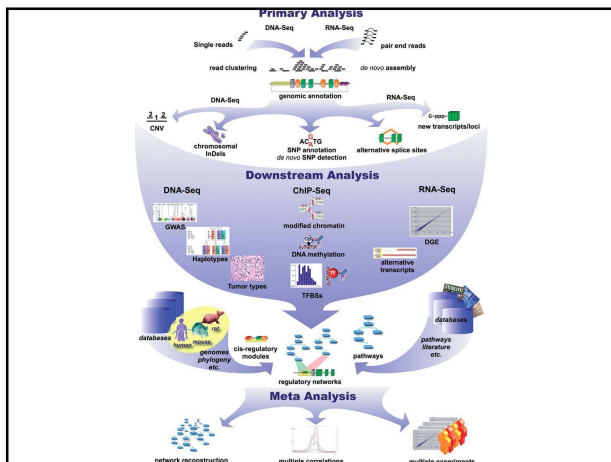
short **many**

long **few**

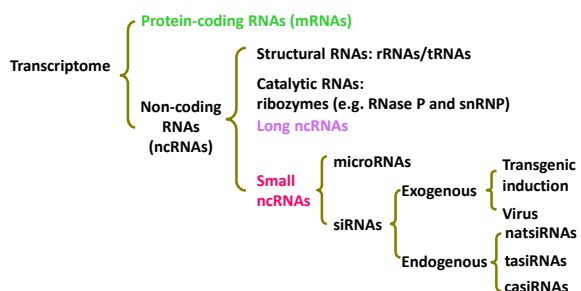
1 kb chemistries in discovering large structural variants and haplotype blocks.

*Average read lengths: *Fragment run: 100-200 bp; *Tag: Fragment: CA; Genome Analyzer: GS; Genome Sequencer: 454; NGS: next-generation sequencing; GS: pyrosequencing; SL: reversible terminator; SLL: sequencing by ligation; SOLiD: support.

Nature Rev. Genet. 2010 11:31-46.



Transcriptome

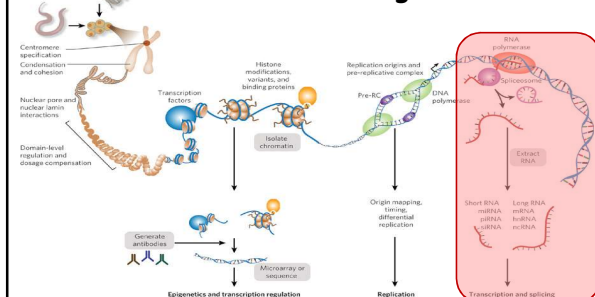


Gene expression data analysis

- Gene expression biology
- Measuring gene expression level
- Identifying differentially expressed genes
- Advanced analysis
- Hickory gene expression data analysis
- Training topics

Part 1: Gene expression biology

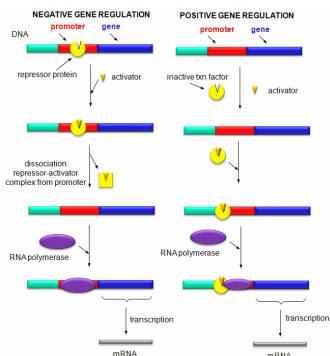
Sequence-based Functional Elements on Central Dogma



Gene expression is the process by which information from a gene is used in the synthesis of a functional gene product. These products are often proteins, but in non-protein coding genes such as rRNA, tRNA or snRNA, the product is a functional RNA.

Nature, 2009 Jun 18;459(7249):927-30.

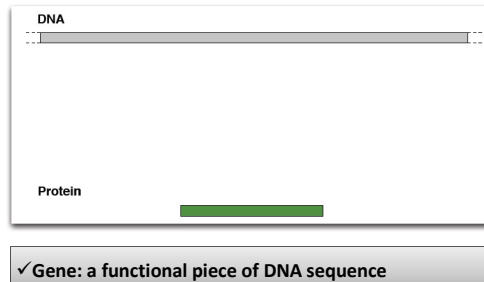
How can gene expression be regulated at the transcriptional level?



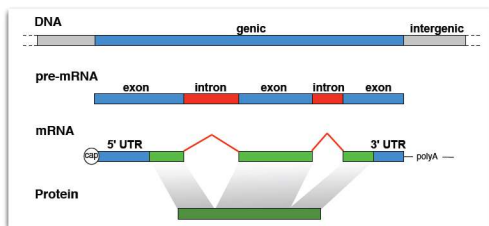
- Chromatin domains
- **Transcription**
- Post-transcriptional modification
- RNA transport
- Translation
- mRNA degradation

- physiological status (nutrition, environment)
- sex and age
- various tissues and cell types
- response to stimuli (drugs, signals, toxins)
- health and disease

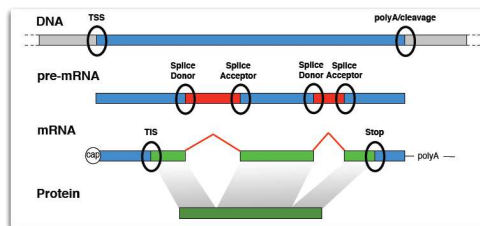
Protein-coding gene



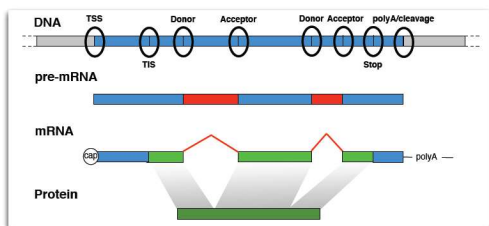
Computational Gene Finding



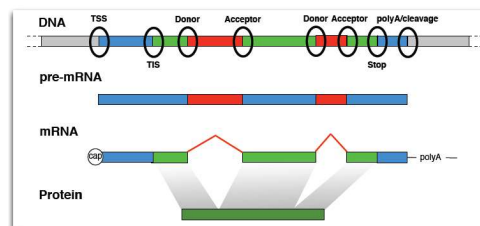
Predict signals used during processing



Predict signals used during processing



Computational Gene Finding



✓ Predict the correct corresponding label sequence with labels "intergenic", "exon", "intron", "5' UTR", etc

Part 2: Measuring gene expression level

Quantitate gene expression level method

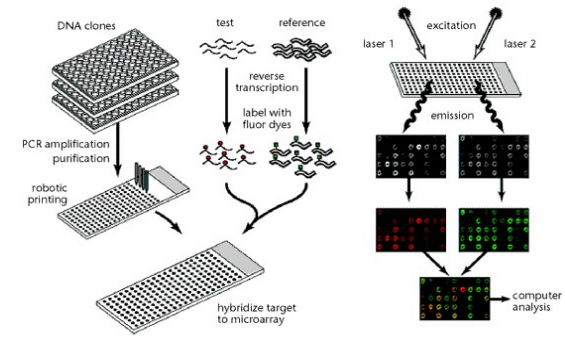
- ✓ Experiment-based approaches:
 - a) RT-PCR
 - b) Northern blot
- ✓ Hybridization-based approaches :
 - a) Microarrays/chip;
 - b) genomic tiling microarrays.
- ✓ Sequence-based approaches:
 - a) EST: Expression Sequence Tag (~400 bp, 20-7000 bp)
 - b) tag-based methods:
 - ✓ CAGE: cap analysis of gene expression (~14-20 bp, 5' ends)
 - ✓ SAGE: serial analysis of gene expression (~14-20 bp, 3' ends)
 - ✓ MPSS: massively parallel signature sequencing (17-20 bp)
- ✓ Next-generation Sequencing-based method:
 - RNA-Seq

Nat Methods. 2008 Jul;5(7):585-7.
Annu Rev Genomics Hum Genet. 2009;10:135-51.
Nat Rev Genet. 2009 Jan;10(1):57-63.

Advantages and disadvantages

- ✓ **Experiment-based approaches:**
 - Low throughput
 - expensive
- ✓ **Hybridization-based approaches :**
 - based on genome sequence;
 - cross-hybridization (high background levels);
 - limited dynamic range of detection (<1000-fold);
 - normalization problems(across different experiments).
- ✓ **Sequence-based approaches:**
 - a) **EST: Expression Sequence Tag (~400 bp, 20-7000 bp)**
 - low throughput;
 - expensive;
 - not quantitative.
 - b) **tag-based methods:**
 - based on expensive Sanger sequencing technology;
 - high throughput;
 - more precise;
 - a portion the short tags cannot be uniquely mapped
- ✓ **Next-generation Sequencing-based method: RNA-Seq**
 - Can be used to detect transcripts of any genome.
 - Low background, highly accurate
 - Large dynamic range of expression levels (~10000-fold)
 - High levels of reproducibility(both for technical and biological replicates)
 - Requires less RNA sample (cloning steps)
 - Lower cost

Microarray schema



From Duggan et al. *Nature Genetics* 21, 10 – 14 (1999)

RNA-seq technologies

➤ Commercially available sequencing technologies used for transcriptome sequencing applications (Sep 15, 2008).

Sequencing platform	ABI3730xl Genome Analyzer	Roche (454) FLX	Illumina Genome Analyzer	ABI SOLiD	HelixScope
Sequencing chemistry	Automated Sanger sequencing	Pyrosequencing on solid support	Sequencing-by-synthesis with reversible terminators	Sequencing by ligation	Sequencing-by-synthesis with virtual terminators
Template amplification method	In vivo amplification via cloning	Emulsion PCR	Bridge PCR	Emulsion PCR	None (single molecule)
Read length	700-900 bp	200-300 bp	32-40 bp	35 bp	25-35 bp
Sequencing throughput	0.03-0.07 Mlb/h	13 Mlb/h	25 Mlb/h	21-28 Mlb/h	83 Mlb/h
Company	Applied Biosystems	Roche Applied Science	Illumina	Applied Biosystems	Helicos
Web site	http://www.appliedbiosystems.com	http://www.roche-applied-science.com	http://www.illumina.com	http://www.appliedbiosystems.com	http://www.helicobio.com

Annu Rev Genomics Hum Genet. 2009;10:135-51.

RNA-Seq: Advantages

- ◆ Sequencing length: 30 - 400bp.
- ◆ Advantages:
 - can be used to detect transcripts of any genome.
 - **low background, highly accurate**
 - **large dynamic range of expression levels (~10000-fold)**
 - **high levels of reproducibility** (both for technical and biological replicates)
 - requires less RNA sample (cloning steps)
 - lower cost

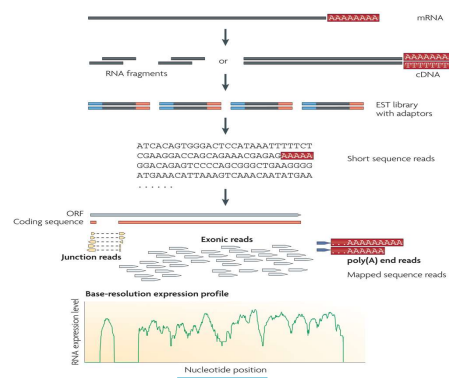
RNA-Seq: Advantages

➤ RNA-Seq v.s. other transcriptomics methods

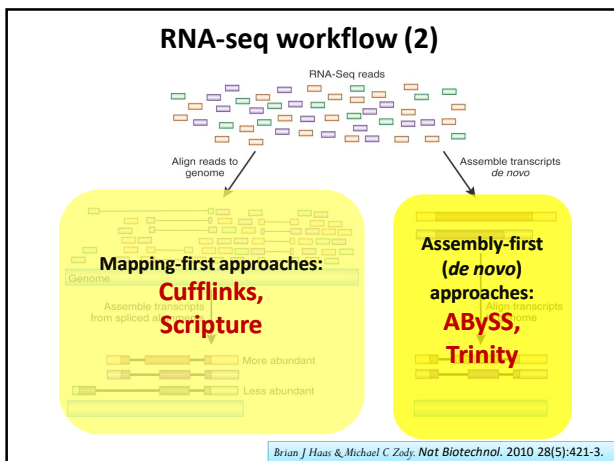
Technology	Tiling microarray	cDNA or EST sequencing	RNA-Seq
Technology specifications			
Principle	Hybridization	Sanger sequencing	High-throughput sequencing
Resolution	From several to 100 bp	Single base	Single base
Throughput	High	Low	High
Reliance on genomic sequence	Yes	No	In some cases
Background noise	High	Low	Low
Application			
Simultaneously map transcribed regions and gene expression	Yes	Limited for gene expression	Yes
Dynamic range to quantify gene expression level	Up to a few-hundredfold	Not practical	>8,000-fold
Ability to distinguish different isoforms	Limited	Yes	Yes
Ability to distinguish allelic expression	Limited	Yes	Yes
Practical issues			
Required amount of RNA	High	High	Low
Cost for mapping transcriptomes of large genomes	High	High	Relatively low

Nat Rev Genet. 2009 Jan;10(1):57-63.

RNA-seq workflow (1)



Zhong Wang et al



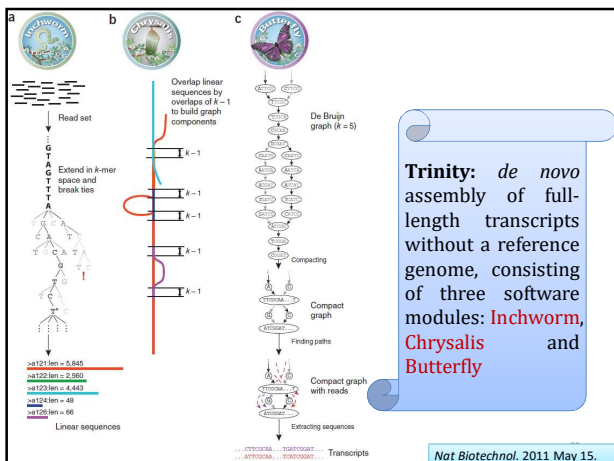
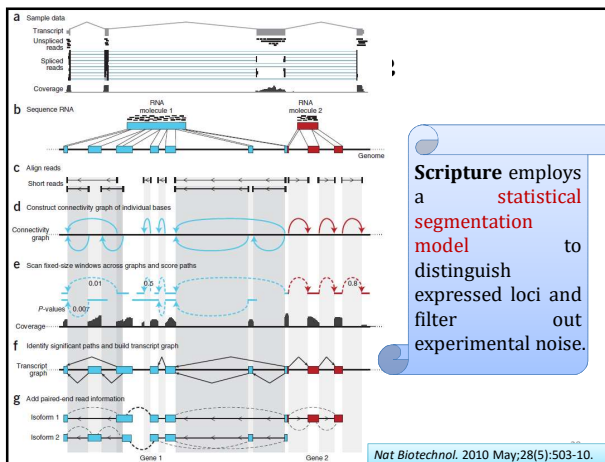
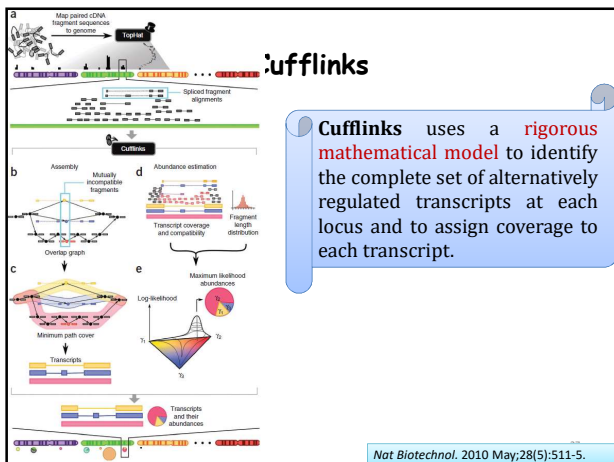
Gene expression level measurement for RNA-seq

✓ RPKM : Reads per kilobase per million mapped reads.

$$RPKM = \frac{\text{Total exon reads}}{\text{mapped reads (millions)} \times \text{exon length (KB)}}$$

1kb transcript with 1000 alignments in a sample of 10 million reads (out of which 8 million reads can be mapped) will have RPKM = $1000 / (1 * 8) = 125$

✓ FPKM : Fragments Per Kilobase of exon per Million fragments mapped (for paired-end sequencing).



Program	Website	Publications
BLAST	http://www.ncbi.nlm.nih.gov/blast/	1990, J. Mol. Biol.
BLAT	http://www.soe.ucsc.edu/~kent/src/	2002, Genome Research
Cross_match	http://www.phrap.org/phredphrapconsed.html	***
ELAND	http://www.illumina.com/	***
TopHat	http://ccb.umd.edu/	2009, Bioinformatics
Novocall		
Mosaik		
Bowtie		2009, Genome Biology
BWA		2009, Bioinformatics
MAQ		2008, Genome Research
SOAP/SOAP2		2009, Bioinformatics
ZOOM		2009, Bioinformatics
PerM		2009, Bioinformatics
BWT-SW		2008, Bioinformatics
RMAP	http://tulai.cshl.edu/map/	2008, BMC Bioinformatics
SHRIMP	http://compbio.cs.toronto.edu/shrimp/	2009, PLoS Computational Biology
SeqMap	http://biogbbs.stanford.edu/~jiangh/SeqMap/	2008, Bioinformatics
MOM	http://mom.csbc.vcu.edu/	2009, Bioinformatics
ProbMatch	http://www.cs.wisc.edu/~jignesh/probmatch/	2009, Bioinformatics
Exonerate	http://www.ebi.ac.uk/~guy/exonerate/	2005, BMC Bioinformatics
SSAHA2	http://www.sanger.ac.uk/Software/analysis/SSAHA2/	2001, Genome Research
Edena	http://www.genomic.ch/edena	2008, Genome Research
VCAKE	http://sourceforge.net/projects/vcake/	2007, Bioinformatics
Euler-SR	***	2007, Genome Research

Part 3: Identifying differentially expressed genes

Statistical methods for finding differentially expressed genes

- Comparing two independent groups
 - a) T-test } Normal distribution
 - b) Linear regression model } Normal distribution
 - c) Wilcoxon rank sum test } Any distribution
 - d) SAM } Any distribution
- Comparing more than two groups
 - a) F-test } Normal distribution
 - b) Linear regression model } Normal distribution
 - c) Wilcoxon rank sum test } Any distribution
 - d) SAM } Any distribution
- Software: R language (Bio-conductor)



T-test

- ✓ Suppose we want to find genes that are differentially expressed between different conditions/phenotypes, e.g. two different tumor types.

Tumor sample	1	1	1	1	2	2	2	2
gene1	X1	X2	X3	X4	Y1	Y2	Y3	Y4
gene2								
gene3								

$$S^2 = \frac{1}{n_1 + n_2 - 2} \left(\sum_{i=1}^{n_1} (X_i - \bar{X}_1)^2 + \sum_{i=1}^{n_2} (Y_i - \bar{Y}_2)^2 \right)$$

$$T(X, Y) = \frac{\bar{X} - \bar{Y}}{s \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$$

- Need check normal assumption
- More arrays in each group more confidence in results

- ✓ After a test statistic is computed, it is convenient to convert it to a p-value. $P \text{ value} = P(t > T(X, Y))$

Linear regression model

- ✓ Expression of gene x is made of a baseline expression level (from control group), plus the group effect.

$$Y = Y_0 + \beta Z$$

Y_0 : baseline exp. Level; β : group effect; Z : group variable (0 for control obs., 1 for group obs.)

- ✓ P-value can be used to test group effect.

ANOVA Table

	d.f.	Sum Sq	Mean Sq	F statistic	p-value
Group	1	29.4115	29.4115	31.323	0.000512
Residuals	8	7.5119	0.939		

- ✓ Results – one p-value per gene

Linear regression model

- ✓ Expression of gene x : baseline expression level, group effect and patient age group

$$Y = Y_0 + \beta Z + \gamma W$$

Y_0 : baseline exp. Level;

β : group effect;

Z : group variable (0 for control obs., 1 for group obs.)

γ : age effect

W : age variable (0 for 0-15, 1 for 16-29, 2 for 30+)

- ✓ ANOVA table:

	d.f.	Sum Sq	Mean Sq	F statistic	p-value
Treatment	1	20.6848	20.6848	25.9737	0.000263
Age	2	27.2838	13.6419	17.13	0.000305
Treatment:Age	2	0.5526	0.2763	0.3469	0.713707
Residuals	12	9.5565	0.7964		

- ✓ Results: a list of p-values

Wilcoxon rank sum test

- ✓ Non-parametric test for equality of two distributions.
- ✓ Compute the ranks of observations in the pooled sample.

Observations: 0:3 0:5 0:8 0:9 1:3 2:4

Ranks: 1 2 3 4 5 6

Groups: 1 1 1 2 2 2

- ✓ The test statistic is a function of the sum of ranks in group 1; here, $R_1 = 6$.
- ✓ For small sample sizes, the null distribution of the test statistic can be computed exactly. For large sample size, a normal approximation is used.
- ✓ Advantage: Non-parametric, robust against outliers

➤ SAM

- ✓ Does not assume normal distribution.
– Instead, p-values computed via permutation
- ✓ The SAM ('significance analysis of microarrays') test statistic is

$$S = \frac{R_g}{c + SE_g}$$

R_g be the mean log ratio of the expression levels of one gene;

SE_g be its standard error;

constant c can be taken to be the 90th percentile SE value.

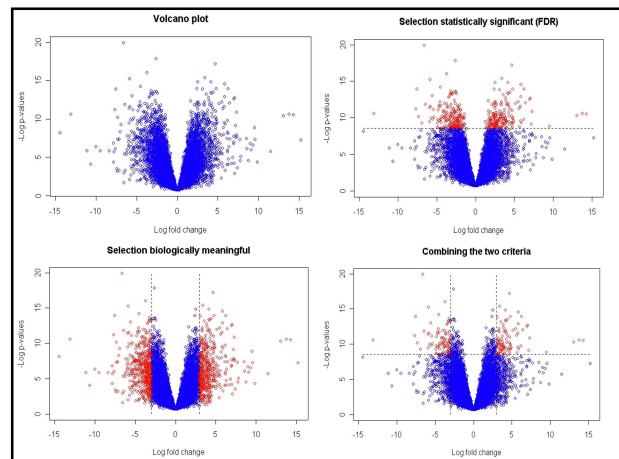
- ✓ One p-value per gene

➤ Multiple testing: the problems

- ✓ Type I: or false-positive error occurs when we declare a gene to be differentially expressed when in fact it is not.
- ✓ Type II: or false-negative error occurs when we fail to detect a differentially expressed gene.
- ✓ The available methods to address the problems:
 - Family-wise error-rate control:** One approach to multiple testing is to control the family-wise error rate (FWER), which is the probability of accumulating one or more false-positive errors over a number of statistical tests.
 - False-discovery-rate control:** An alternative approach to multiple testing considers the false-discovery rate (FDR), which is the proportion of false positives among all of the genes initially identified as being differentially expressed - that is, among all the rejected null hypotheses.

➤ P-value vs. Fold change

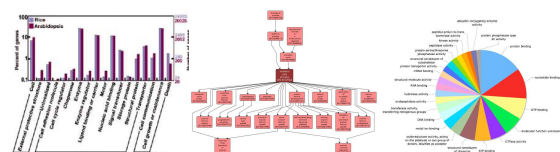
- ✓ P-values measure distance in terms of probability.
– Statistical significance
- ✓ Fold changes: measure distance in arbitrary scale.
The simplest method for identifying differentially expressed genes is to evaluate the log ratio between two conditions (or the average of ratios when there are replicates) and consider all genes that differ by more than an arbitrary cut-off value to be differentially expressed.
– Biological meaning
- ✓ Differentially expressed gene selection: Need combination of these two.



Part 4: Advanced analysis

GO analysis

- ✓ The Gene Ontology, or GO, is a major bioinformatics initiative to unify the representation of gene and gene product attributes across all species.
- ✓ Tools: AmiGO (http://amigo.geneontology.org/cgi-bin/amigo/blast.cgi?session_id=6985amigo1343799107)
OBO-Edit (<http://oboedit.org/>)
WEGO (<http://wego.genomics.org.cn/cgi-bin/wego/index.pl>).
- ✓ Inputs: FASTA file, GO number list... ..
- ✓ Outputs: Histogram, Interactive GO graph, Pie Charts... ..



Clustering gene expression data

Algorithms:
 a) K-means
 b) Hierarchical clustering
 c) K-median
 d) Bi-clustering

Tools and software:
 a) R language,
 b) Clustal,
 c) Mev.

If two genes are related (have similar functions or are co-regulated), their expression profiles should be similar (e.g. low Euclidean distance or high correlation).

Pathway mapping and analysis

Identify up-/down-regulated genes

KO ID mapping KEGG

Co-expression network reconstruction

Algorithms:
 a) PCC
 b) Weighted PCC
 c) Multiple rank (MR)

Visualization software:
 Cytoscape

GO enrichment analysis
Function model analysis

Gene regulatory network reconstruction

Gene expression data Discretization
 ✓ Equal Width Discretization
 ✓ Equal Frequency Discretization
 ✓ Kmeans Discretization
 ✓ Column Kmeans Discretization
 ✓ Bkmeans Discretization

Gene regulatory network reconstruction
 ✓ Greedy search
 ✓ K2
 ✓ Aracne
 ✓ Matlab
 ✓ ...

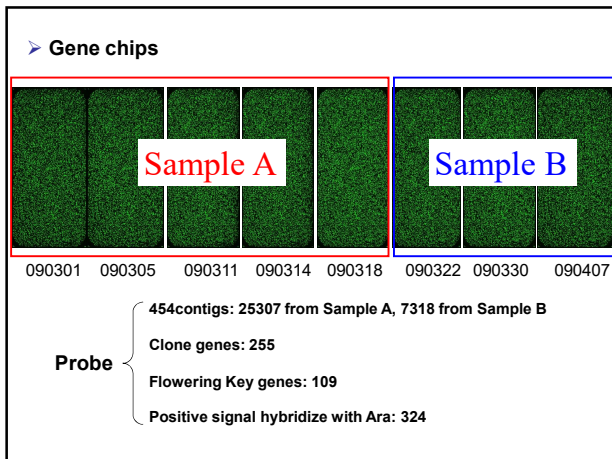
Part 4: gene expression data analysis: Examples

Hickory gene expression data analysis

Materials and Methods

➤ 454 sequencing

	Sample A	Sample B
Read number	431,759	444,905
Avg. read length	332	332
contig	25339	26935
Specific gene	4951	5887
ORF number	15085	16387



Methods

- Flowering network construction of Arabidopsis based on literatures.**
 - Key word: flowering floral ect.
 - The total number of literatures: About 1500.
 - Flowering genes: 436 (Common name, Locus ID).
 - Flowering construction and visualization based on Cytoscape software.
- 454 sequencing analysis.**
 - Contig assemble: CAP3 software (Sample A, Sample B and All)
 - Blast analysis against Arabidopsis: Blast software (Contigs->Ara. genes).
 - Result filter: Identity percent: 80%, E-value: 1e-5, Coverage: 70%.

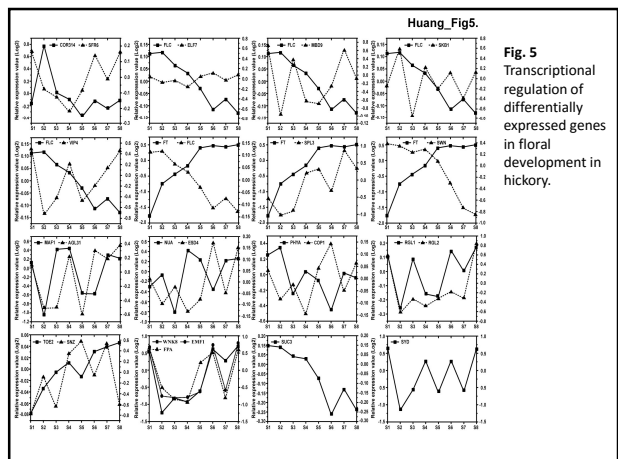
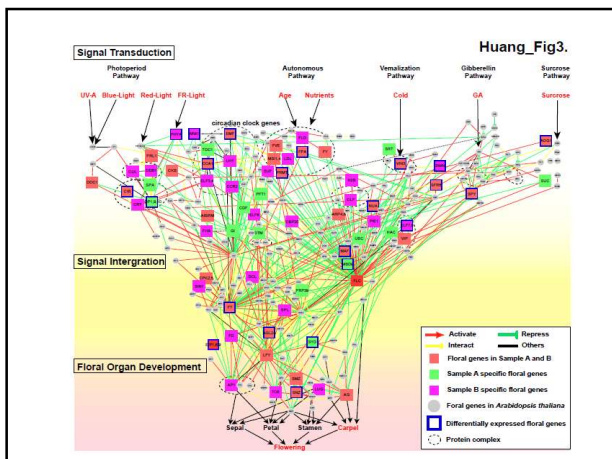
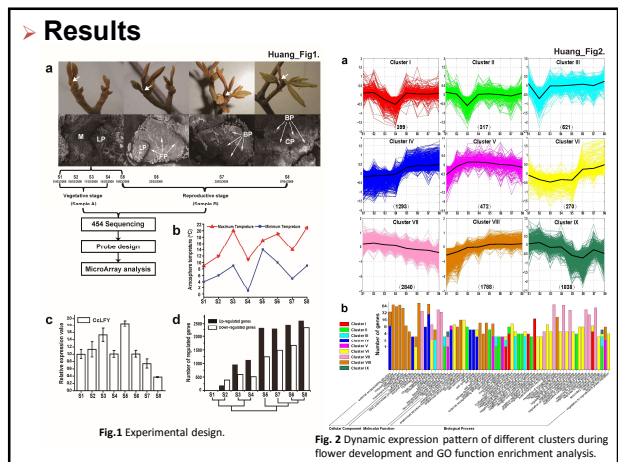
Methods

- Differentially expressed gene analysis.**

Constraint conditions:
Fold change:4, Num(fc): 1. Signal value: except all A's
- Gene expression pattern analysis.**

Software: MeV software.
Algorithm: K-means.
- GO Enrichment analysis**
- Co-expression network reconstruction for flowering genes.**

Algorithm: Mutual Rank (MR) (2008, NAR)
- Real time quantitative PCR**



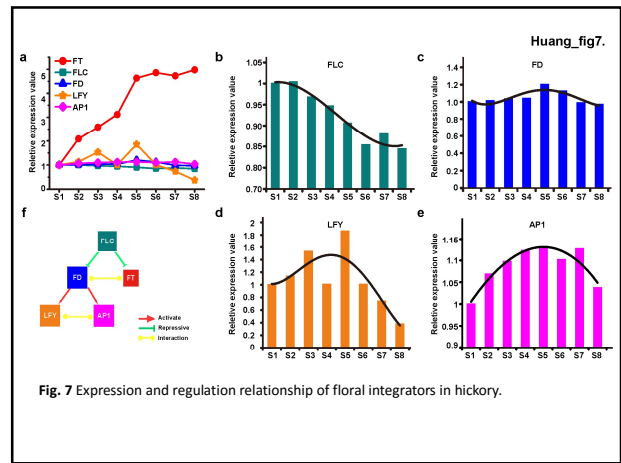
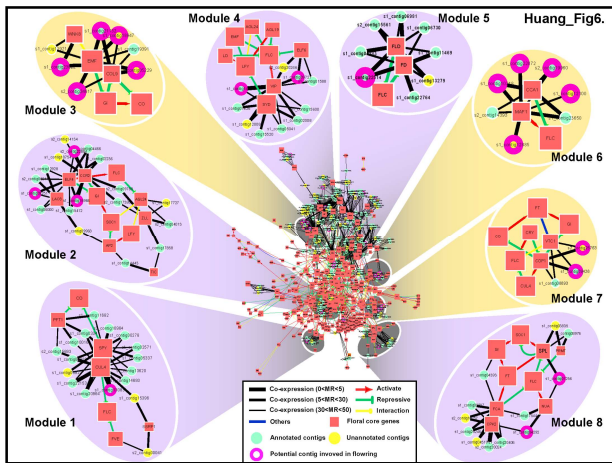
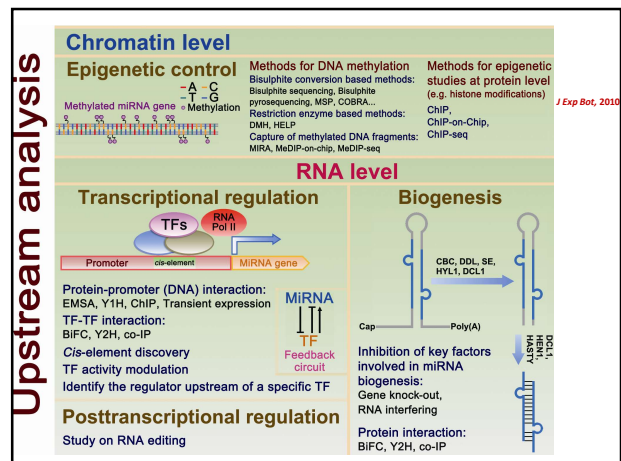
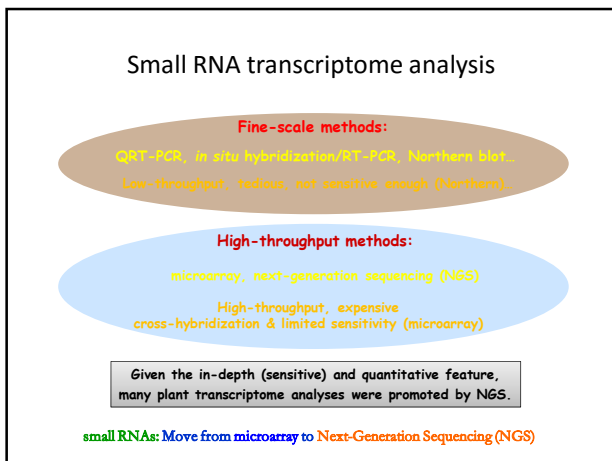
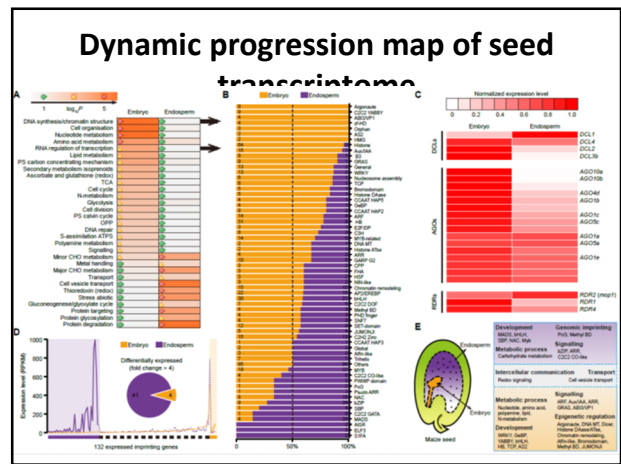
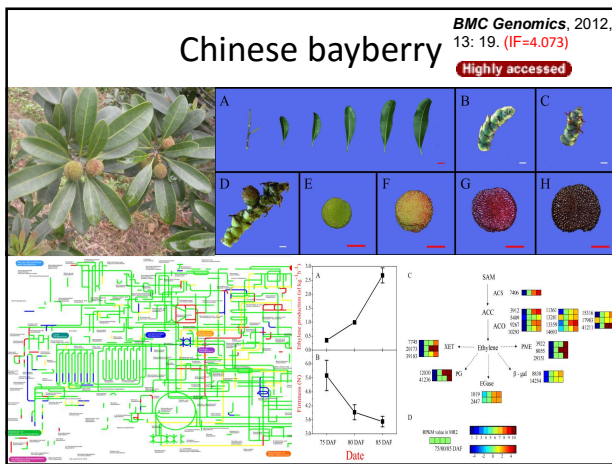
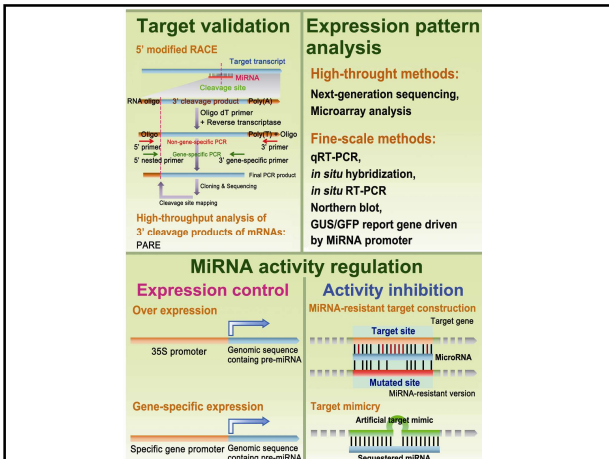


Fig. 7 Expression and regulation relationship of floral integrators in hickory.





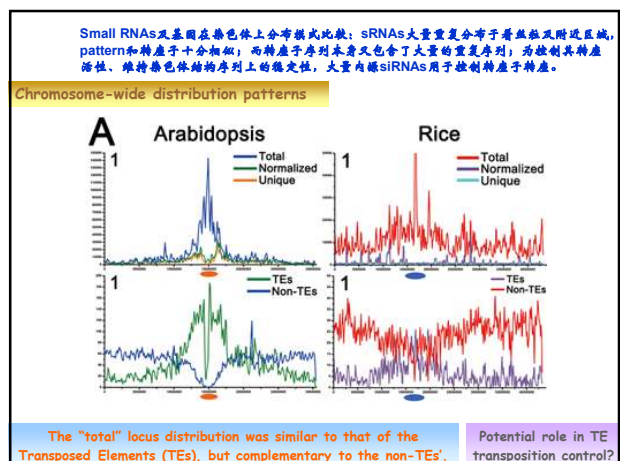
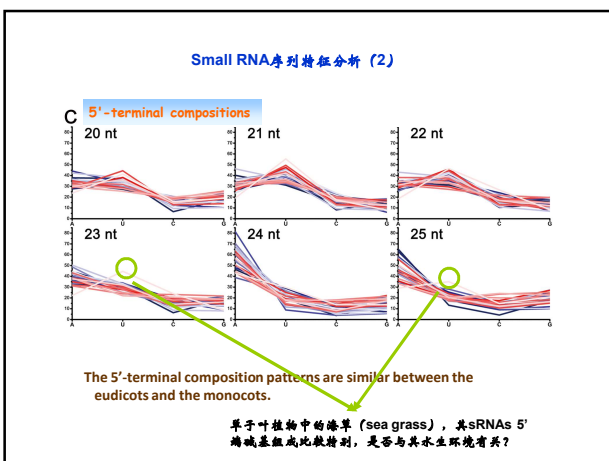
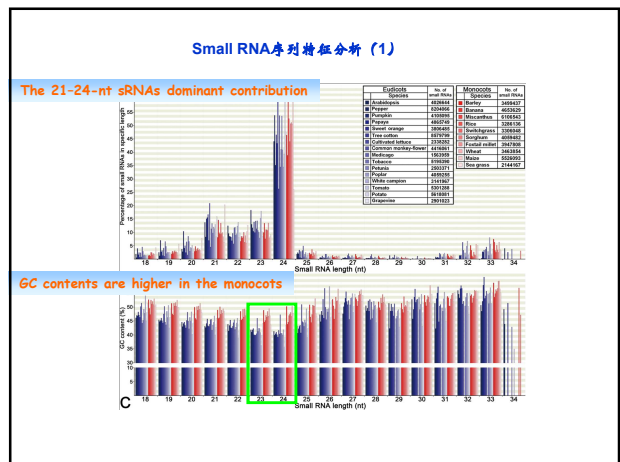
Research work of my lab

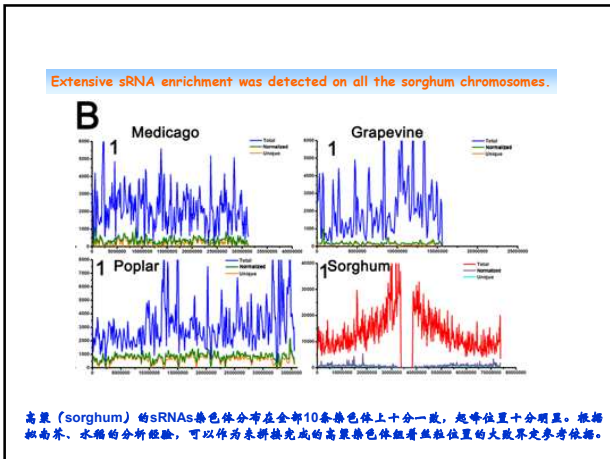
- Plant ncRNAs
- biogenesis,
- characteristics,
- expressions,
- interactions,
- regulations,
- even dynamic functions, 3D...

Small RNAs in angiosperms: sequence characteristics, distribution and generation

Eudicots (16)	Monocots (10)
Arabidopsis 拟南芥	水稻 rice
Tomato 西红柿	玉米 Maize
Medicago 苜蓿	大麦 barley
Pepper 辣椒	香蕉 banana
Pumpkin 南瓜	柳枝稷 Switchgrass
Sweet orange 甜橙	高粱 Sorghum
Tree cotton 木棉	小麦 wheat
Cultivated lettuce 莴苣	海草 Sea grass
Common monkey-flower 猴面花	芒草 Miscanthus
Tobacco 烟草	谷子 Foxtail millet
Petunia 矮牵牛花	
Poplar 白杨	
White campion 白花耧子草	
Potato 土豆	
Grapevine 葡萄	
Papaya 木瓜	

Bioinformatics, 2010





Small RNAs derived from gene models

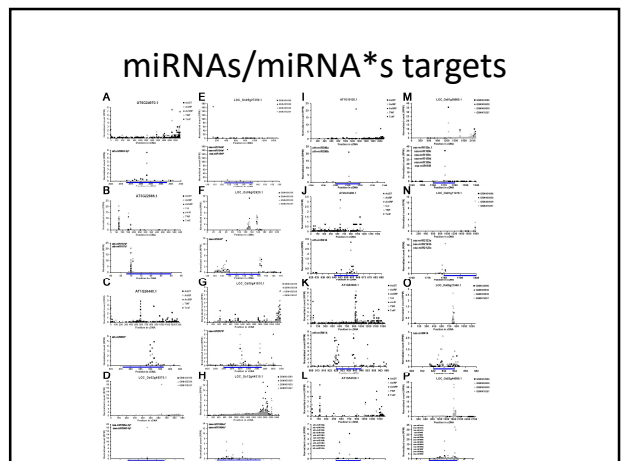
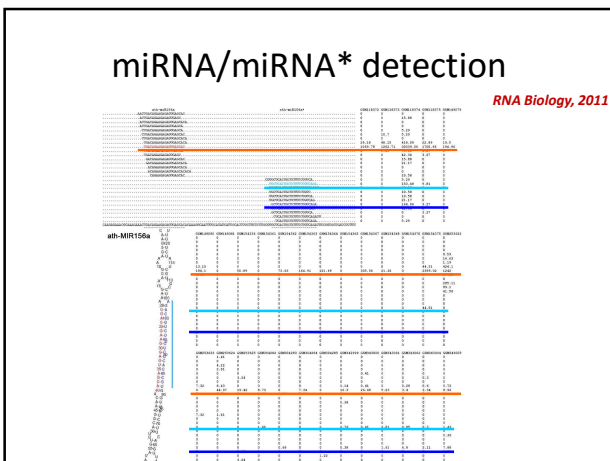
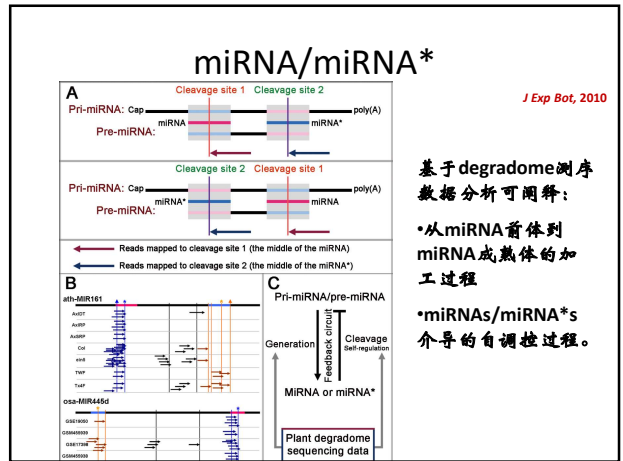
Bioinformatics, 2010

Species	Major division (percentage) ^a	Subdivision (percentage) ^b	No. of sRNA loci analyzed (total/unique)
Arabidopsis	Intergenic loci (Total ^a : 80.48%; Unique ^b : 79.30%)	-	9,008,884/2,641,530
	Intragenic ^c loci (Total ^a : 19.04%; Unique ^b : 20.14%)	5' UTR ^d (Total ^b : 0.79%; Unique ^b : 1.65%) 3' UTR ^d (Total ^b : 1.38%; Unique ^b : 3.63%) Exons ^d (Total ^b : 83.21%; Unique ^b : 79.85%) Introns ^d (Total ^b : 7.37%; Unique ^b : 9.19%) Others ^d (Total ^b : 7.05%; Unique ^b : 5.68%)	
	Other loci ^e (Total ^a : 0.49%; Unique ^b : 0.56%)	-	
	Intergenic loci (Total ^a : 80.30%; Unique ^b : 85.24%)	-	
	Intragenic ^c loci (Total ^a : 19.31%; Unique ^b : 14.42%)	5' UTR ^d (Total ^b : 0.72%; Unique ^b : 1.77%) 3' UTR ^d (Total ^b : 1.76%; Unique ^b : 7.12%) Exons ^d (Total ^b : 56.30%; Unique ^b : 39.74%) Introns ^d (Total ^b : 37.75%; Unique ^b : 46.08%) Others ^d (Total ^b : 3.47%; Unique ^b : 5.29%)	
Rice	Other loci ^e (Total ^a : 0.38%; Unique ^b : 0.35%)	-	22,147,409/1,529,832

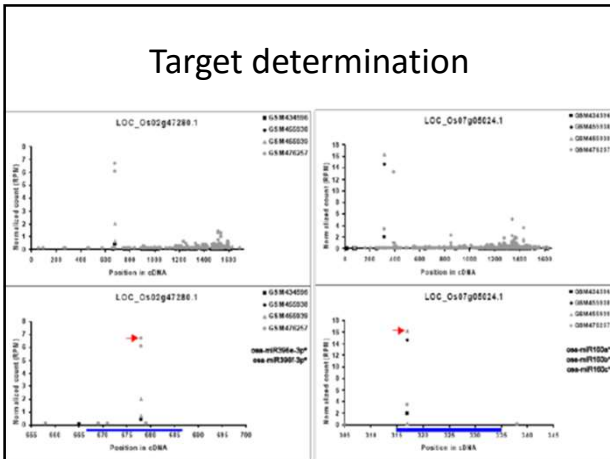
Plant microRNA knowledge base

Nucleic Acids Res, 2011

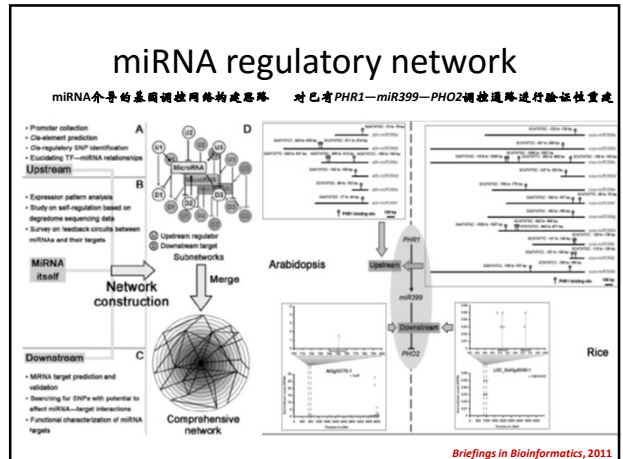
包含4个主要功能模块



Target determination



miRNA regulatory network

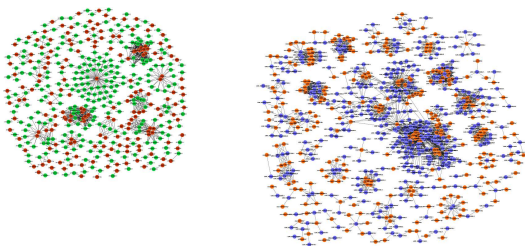


miRNA/miRNA* regulatory network

基于miRNA target lists, miRNA* target lists和co-regulated target lists构建network

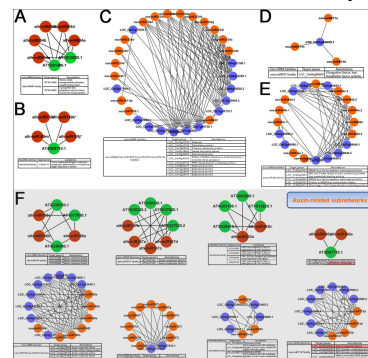
Arabidopsis (all targets)

Rice (all targets)



基于降解组测序数据, 对拟南芥、水稻中已知miRNAs的靶基因预测和大规模鉴定; 利用sRNA高通量测序数据, 基于表达量鉴定了所有miRNAs对应的miRNA*, 并对miRNA*的潜在的靶基因进行了预测鉴定; 最终构建了由miRNAs/miRNA*介导的基因调控网络

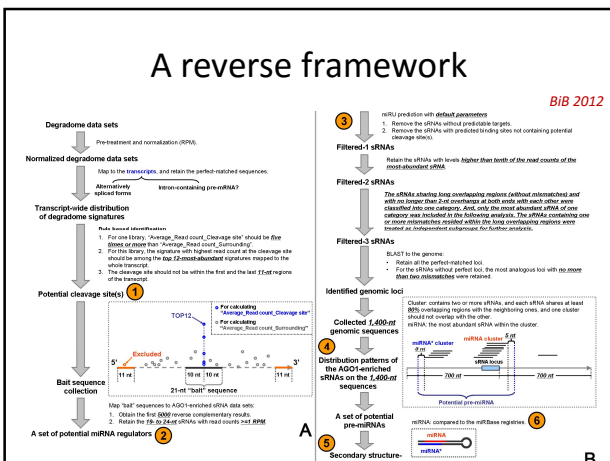
Subnet analysis



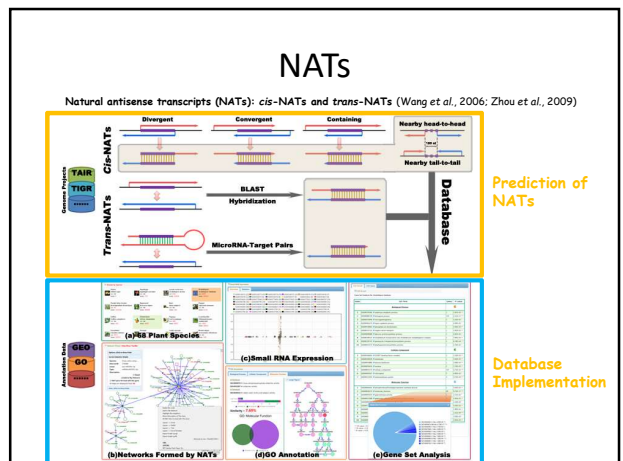
发现一些可能具有重要生物学功能 (参与重金属胁迫应答、植物抗病相关) 的子网络

生长素信号相关子网络在拟南芥、水稻中具有高度保守性, 但也发现了一些物种特异的调控关系 (红色背景)

A reverse framework



NATs



NATs Generated Small RNAs

sRNA loci are enriched in the overlapping regions of trans-NATs, but not for cis-NATs.

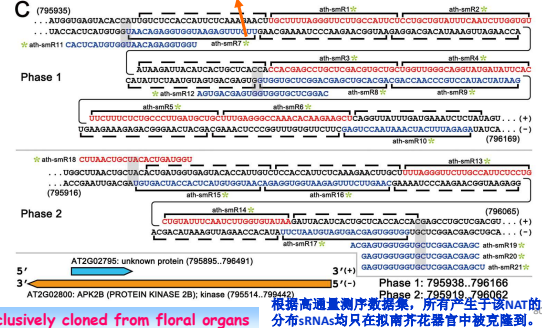
Species	Cis-NATs			
	Overlap [†] [total/unique] [†]	All [†] [total/unique] [†]	Average score [†] [total/unique] [†]	P-value [‡] [total/unique] [†]
Arabidopsis	38.89/7.11	10.62/5.63	3.10/1.95	<0.0001/0.0448
Poplar	8.42/11.19	5.42/2.68	2.61/5.26	0.4525/0.1348
Papaya	7.05/3.85	4.66/2.33	1.99/1.97	0.0094/0.0011
Rice	3.28/1.13	4.62/0.58	1.62/2.31	0.0011/<0.0001
Maize	13.33/1.73	11.68/1.19	1.32/2.24	0.0453/<0.0001
Sorghum	8.13/3.64	8.11/2.54	1.69/2.17	0.9836/0.0727

Species	Trans-NATs			
	Overlap [†] [total/unique] [†]	All [†] [total/unique] [†]	Average score [†] [total/unique] [†]	P-value [‡] [total/unique] [†]
Arabidopsis	169.65/60.06	48.62/19.00	3.74/3.51	<0.0001/<0.0001
Poplar	159.94/9.19	23.80/2.63	8.63/5.48	<0.0001/<0.0001
Grapevine	35.25/0.74	17.87/0.47	2.39/1.95	<0.0001/<0.0001
Papaya	26.84/7.52	20.14/1.13	1.56/1.42	<0.0001/0.2838
Medicago	61.37/5.00	28.49/1.74	3.17/4.53	<0.0001/<0.0001
Rice	210.30/6.23	17.33/2.65	14.06/7.03	<0.0001/<0.0001
Maize	116.44/6.97	18.97/1.61	7.13/6.15	<0.0001/<0.0001
Sorghum	344.77/5.17	64.09/2.39	10.22/3.37	<0.0001/<0.0001

Organ specific - Arabidopsis

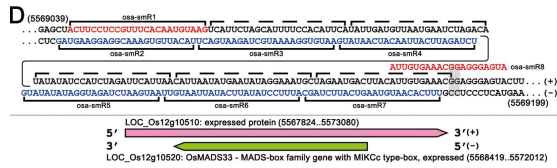
Phase-distributed sRNA in the overlapping region of a cis-NAT in Arabidopsis

标记星号的是在基因组上仅有一个完全匹配位点的sRNAs



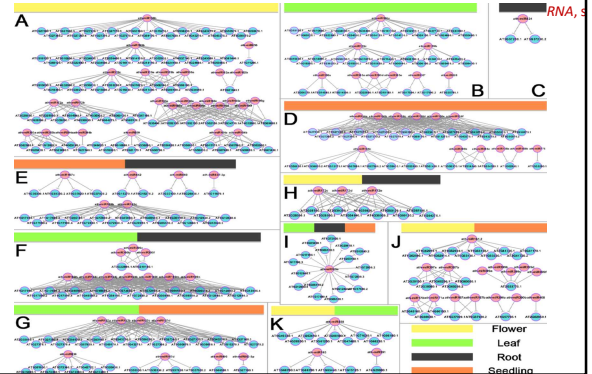
Organ specific - rice

Phased sRNA in the overlapping region of a cis-NAT in rice



Organ-specific regulatory role?

Organ-specific miRNAs in Arabidopsis



PlantNATsDB

Plant Natural Antisense Transcripts DataBase

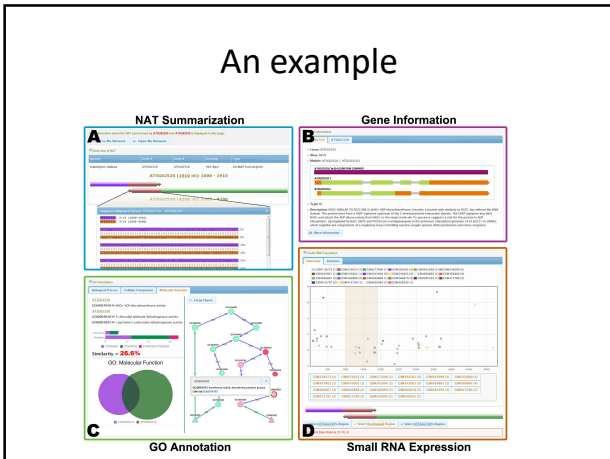
NAR, 2012

Statistics

PlantNATsDB predicted 2,066,720 NATs from 69 plant species

No.	ID	Scientific name	MicroRNAs ^{a,b}	Genes	Cis-NATs ^b	Trans-NATs (MicroRNA-Target Pairs)	All-NATs
1	aco	Allium cepa	NA	4063 (10)	NA	5 (NA)	5
2	aco	Aquilegia coerulea	45 (45)	13556 (610)	NA	772 (631)	772
3	aly	Arabidopsis lyrata	375 (373)	32670 (12527)	918	19636 (15686)	20554
4	ath	Arabidopsis thaliana	243 (243)	33239 (13875)	3005	16915 (12648)	19920
5	bdt	Brachypodium distachyon	19 (19)	25532 (6007)	36	110526 (3747)	110562
6	bna	Brassica napus	48 (48)	50542 (20723)	NA	46668 (738)	46668
7	brv	Beta vulgaris	NA	4785 (249)	NA	192 (NA)	192
8	can	Capsicum annuum	NA	14727 (2138)	NA	6119 (NA)	6119
9	oca	Coffea canephora	NA	7511 (202)	NA	163 (NA)	163
10	ocl	Citrus clementina	5 (5)	32287 (2238)	NA	3665 (111)	3665
11	cpa	Carica papaya	1 (1)	25536 (4001)	180	4047 (14)	4227
12	cre	Chlamydomonas reinhardtii	85 (84)	15935 (8761)	1450	28051 (4919)	29501
13	osa	Cucumis sativus	NA	32775 (6104)	1471	16014 (NA)	17485
14	asi	Citrus sinensis	64 (59)	26081 (3392)	NA	8385 (893)	8385
15	osa	Euphorbia esula	NA	10727 (103)	NA	96 (NA)	96
16	ari	Ectocarpus siliculosus	NA	9122 (387)	NA	340 (NA)	340
17	tax	Festuca arundinacea	15 (14)	10617 (295)	NA	229 (78)	229

An example



Small RNAs derived from gene models

Species	Major division (percentage) ^a	Subdivision (percentage) ^b	No. of sRNA loci analyzed (total/unique)
Arabidopsis	(Total: 80.48%; Unique ^a : 79.30%)	Intergenic loci	~1.8%
		5' UTR ^c (Total: 0.79%; Unique ^b : 1.65%)	
		3' UTR ^c (Total: 1.38%; Unique ^b : 3.63%)	
		Intragenic ^d loci (Total: 19.04%; Unique ^b : 20.14%)	
		Exons ^e (Total: 83.21%; Unique ^b : 79.85%)	
Introns ^f (Total: 7.37%; Unique ^b : 9.19%)			
Others ^g (Total: 7.05%; Unique ^b : 5.68%)			
Rice	(Total: 19.31%; Unique ^a : 14.42%)	5' UTR ^c (Total: 0.72%; Unique ^b : 1.77%)	~6.6%
		3' UTR ^c (Total: 1.76%; Unique ^b : 7.12%)	
		Exons ^e (Total: 56.30%; Unique ^b : 59.74%)	
		Introns ^f (Total: 37.75%; Unique ^b : 46.08%)	
		Others ^g (Total: 3.47%; Unique ^b : 5.29%)	
Other loci ^h	(Total: 0.38%; Unique ^a : 0.35%)		

Identification of intronic long hairpins

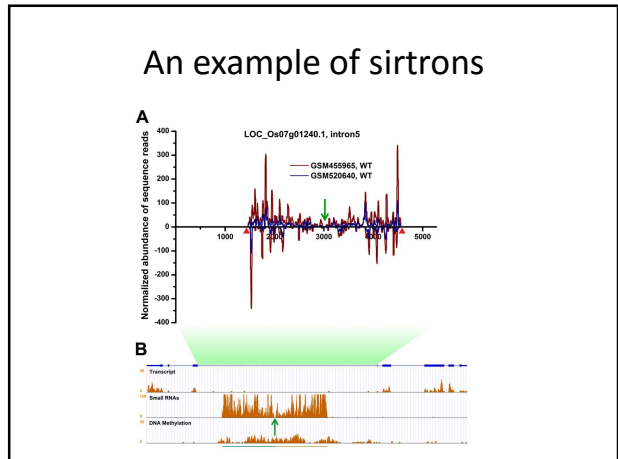
RNA, 2011

Table 1. A list of 21 *IR*-introns with significant numbers of sRNAs^a from the sense strand.

Introns	Length (bp)	No. of sRNAs ^a	% sRNAs from s ¹	Length (bp)	5' arm	Paired stem regions ^d	3' arm	Identity (%)	sRNA density ^e
LOC_Os07g01240.1.intron_3	5275	39969	67.7	978	2012-2091	3027-4000	95	16.98	
LOC_Os01g66379.1.intron_2	10049	5824	64.3	906	4091-5001	5163-6084	93	3.241	
LOC_Os07g31069.1.intron_6	6540	8108	78.2	865	2403-2176	3536-4401	94	3.221	
LOC_Os12g13440.1.intron_1	4436	2553	64.2	811	1478-2253	2589-3426	93	1.443	
LOC_Os11g1760.1.intron_1	675	778	67.1	184	1-184	445-628	90	1.285	
LOC_Os07g35600.1.intron_2	8625	3107		62	4873	87	1.188		
LOC_Os03g34359.1.intron_2	9177	1432		13	8272	96	1.101		
LOC_Os01g1614.1.intron_1	5284	1043		45	5627	92	1.096		
LOC_Os09g17730.1.intron_1	4168	721		13	2373	91	0.759		
LOC_Os02g15039.1.intron_8	5898	2107		21	4096	97	0.536		
LOC_Os03g1370.1.intron_1	3641	911		92	3066	81	0.328		
LOC_Os01g1270.1.intron_3	1224	341		33	966	93	0.274		
LOC_Os08g37700.1.intron_2	601	134	79.9	181	65-245	373-553	95	0.185	
LOC_Os04g32600.1.intron_27	2331	137	85.4	635	711-1362	1438-2080	82	0.089	
LOC_Os09g9910.1.intron_7	576	72	69.4	208	33-242	312-519	93	0.072	
LOC_Os02g12570.1.intron_4	439	32	71.9	160	13-172	185-344	86	0.072	
LOC_Os01g27100.1.intron_3	678	36	63.9	191	18-208	364-554	85	0.068	
LOC_Os04g28420.1.intron_9	581	62	80.6	213	68-284	349-562	92	0.063	
LOC_Os10g3275.1.intron_7	678	56	76.8	195	192-386	412-607	90	0.062	
LOC_Os01g1604.2.intron_8	18327	122	59.9	766	8993-9769	10232-10938	86	0.044	
LOC_Os02g10280.1.intron_4 ^f	499	24	83.3	182	100-285	311-404	87	0.041	

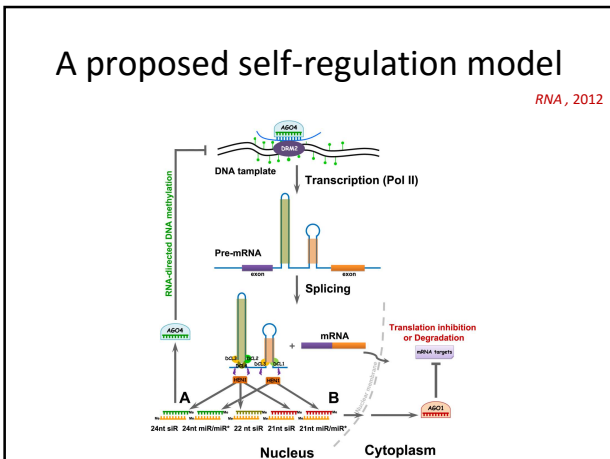
sirtrons

An example of sirtrons

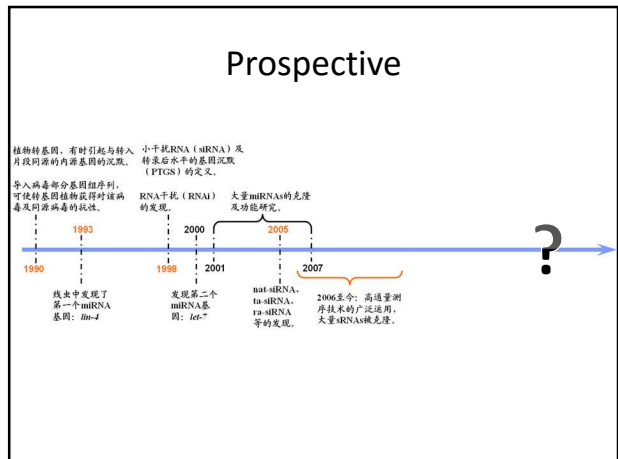


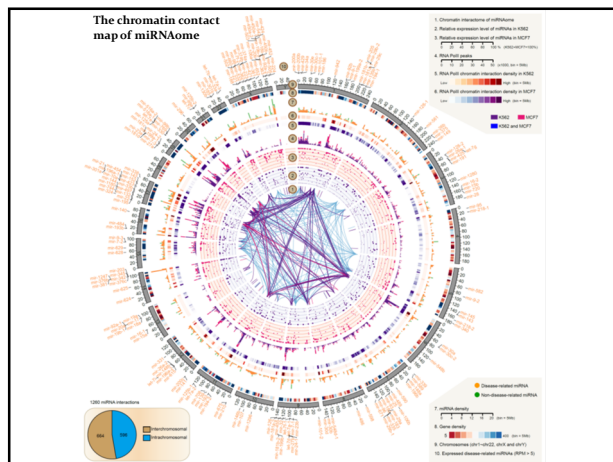
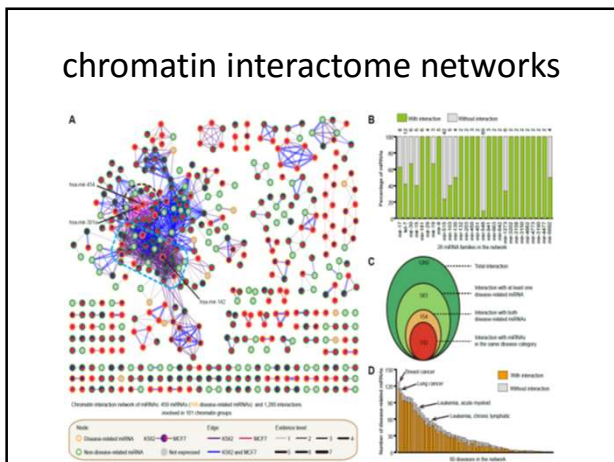
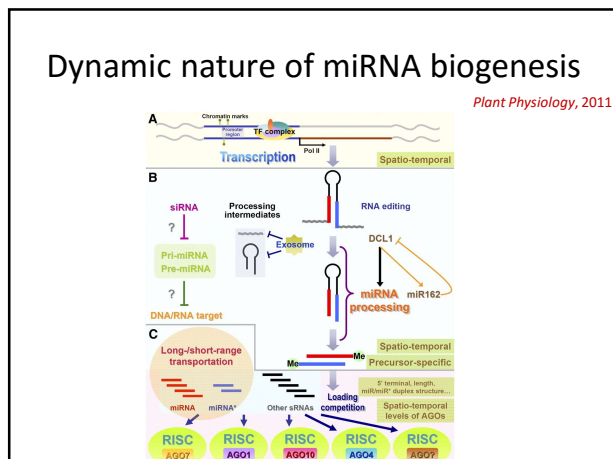
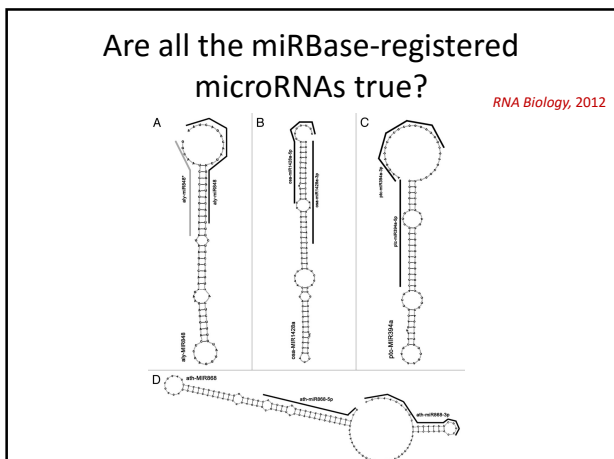
A proposed self-regulation model

RNA, 2012

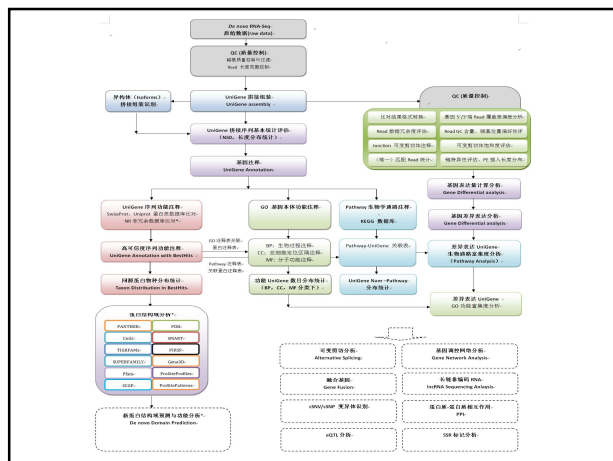


Prospective





Part 4: Practice



➤ **Training Topics**

One: Using R language find differentially expressed genes.

Two: Do the GO analysis for identified differentially expressed genes.

Three: Do the pathway analysis for the differentially expressed genes.

<http://www.cls.zju.edu.cn/binfo/links.htm>

