

Microsoft® Research

Faculty Summit

10
YEAR ANNIVERSARY

Computational Challenges in Analyzing Complex Traits

Jun Zhu

Institute of Bioinformatics
Zhejiang University

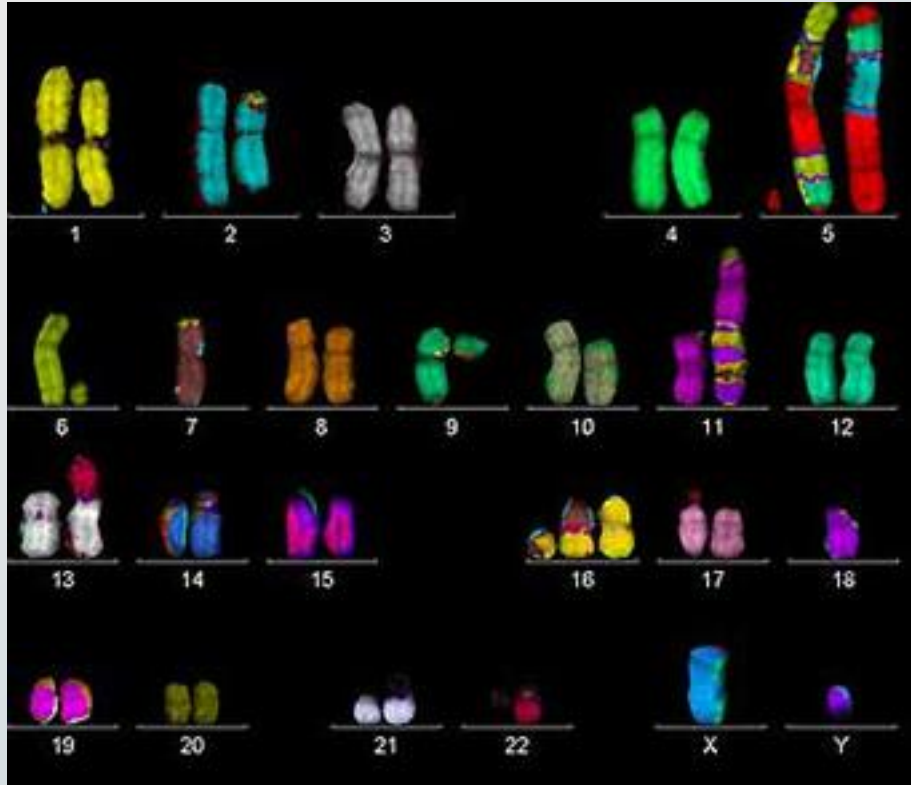
Mendelian Traits

- Phenotypes controlled by singular genes
- No epistasis
- No GE interaction

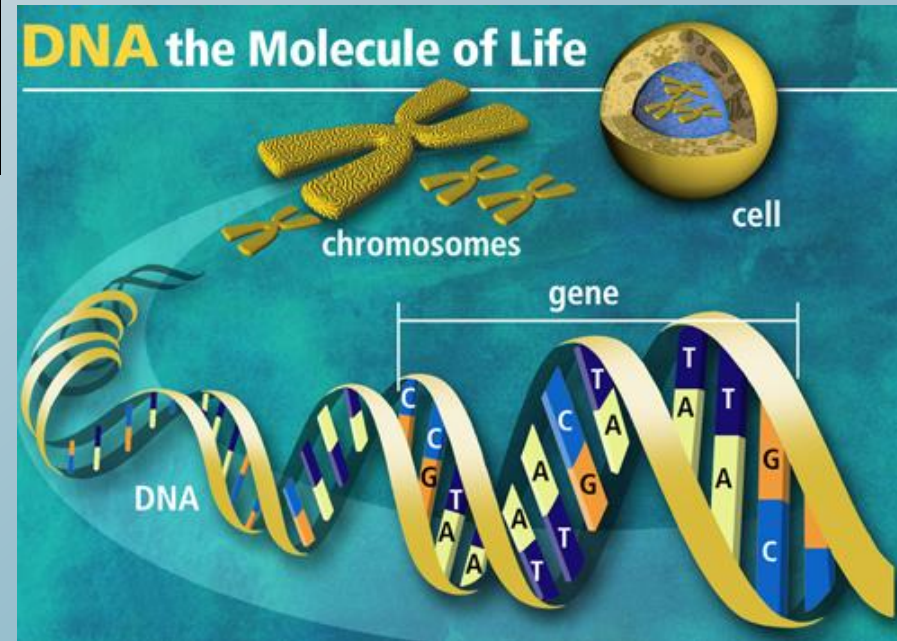
Complex Traits

- Phenotypes controlled by multiple genes
- Epistasis (gene-gene interaction)
- Gene-environment interaction
- Genetic pleiotropy and heterogeneity
- Low heritability
- Limited statistical power

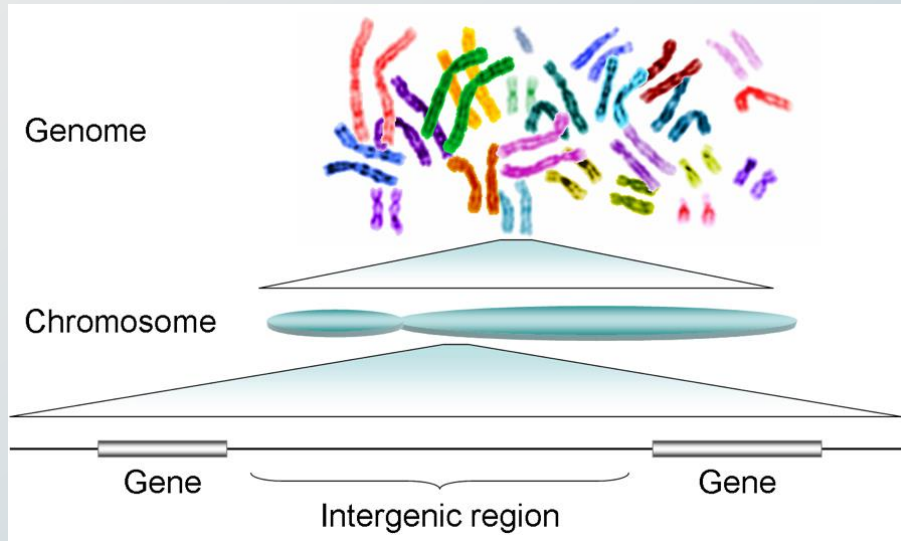
● Human Chromosomes



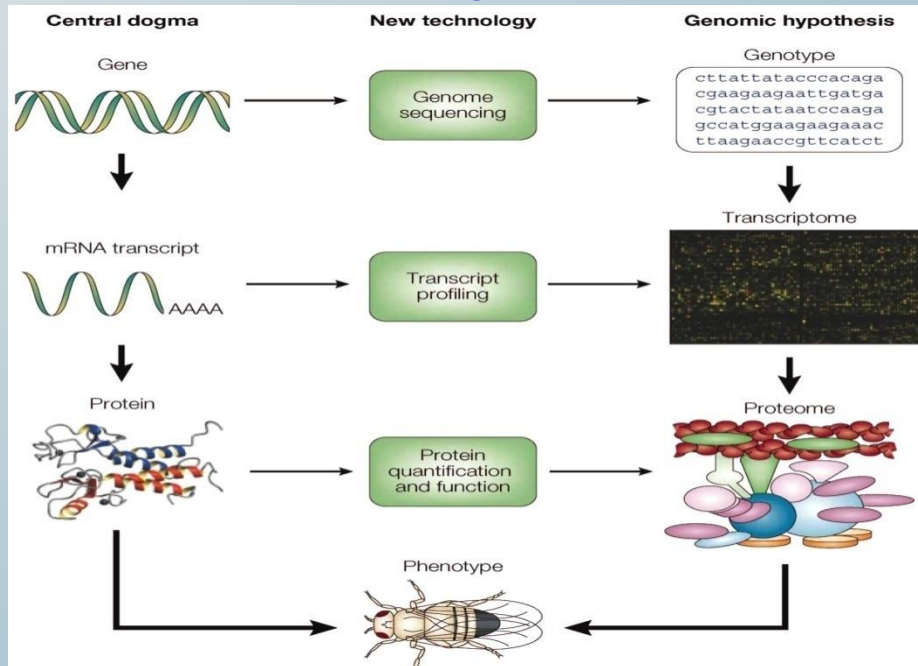
● Chromosome to DNS



● From Genome to Genes



● From Gene to Phenotype



Partition of Phenotypic Variation

Phenotype

$$y = \mu + E + G + GE + e$$



Genetic Effect: A, D, I

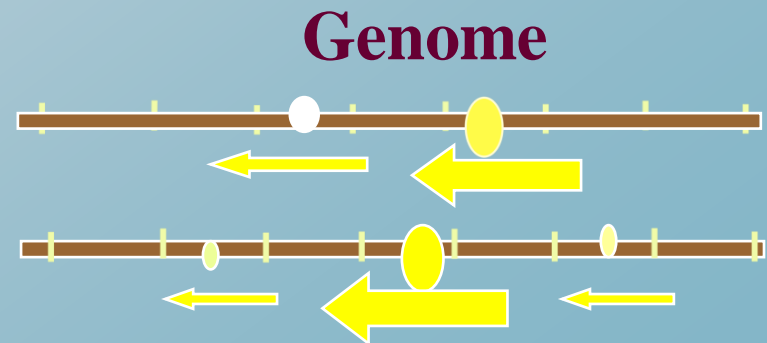
GE: AE, DE, IE

Macro Env, Micro Env

Genetic Effects



QTL
Position &
Effects



Partitioning of Phenotypic Variation

$$y = \mu + E + G + GE + e$$

- **Gene Effects (G):**
Additive (A), Dominance (D)
Gene-Gene Interaction (AA, AD, DD)
- **Environment Effects (E):**
Location, Weather, Treatment
- **Gene-Environment Interaction (GE)**
AE, DE, AAE, ADE, DDE
- **Random Error (e)**

Partitioning of Phenotypic Variation

$$y = \mu + E + G + GE + e$$

- **Classical quantitative genetics**

$$G = A + D + (AA + AD + DD)$$

$$GE = AE + DE + (AAE + ADE + DDE)$$

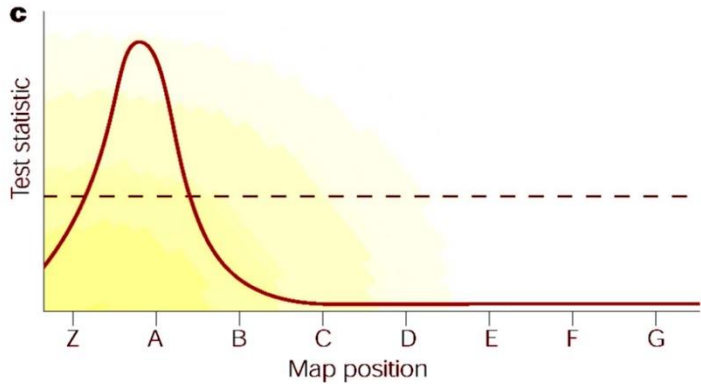
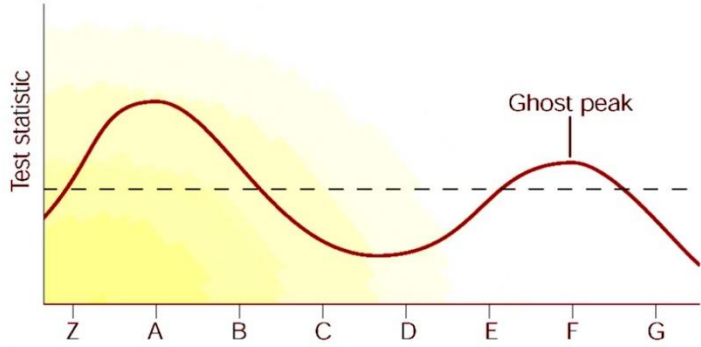
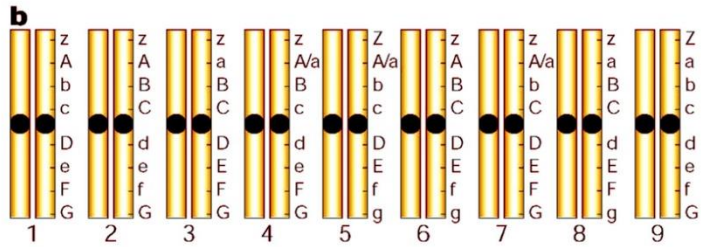
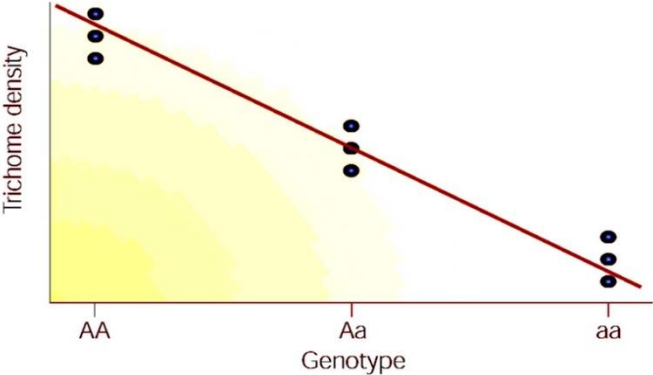
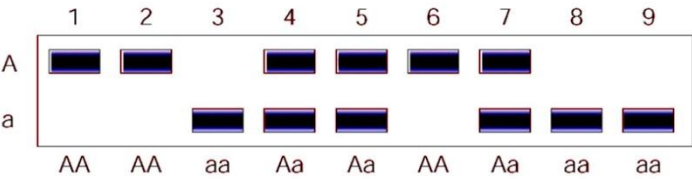
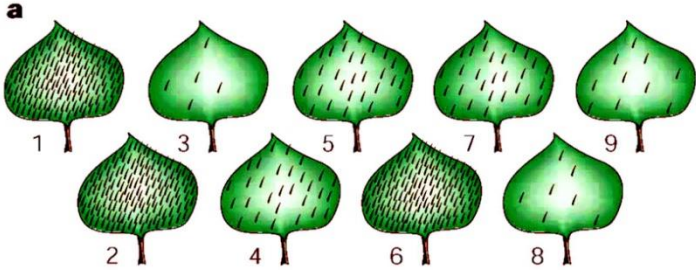
- **Molecule quantitative genetics**

$$G = \sum a + \sum d + (\sum \sum aa + \sum \sum ad + \sum \sum dd)$$

$$GE = \sum ae + \sum de + (\sum \sum aae + \sum \sum ade + \sum \sum dde)$$

Method of Mapping QTL

Box 2 | Quantitative trait loci mapping methods



Presentation of QTL

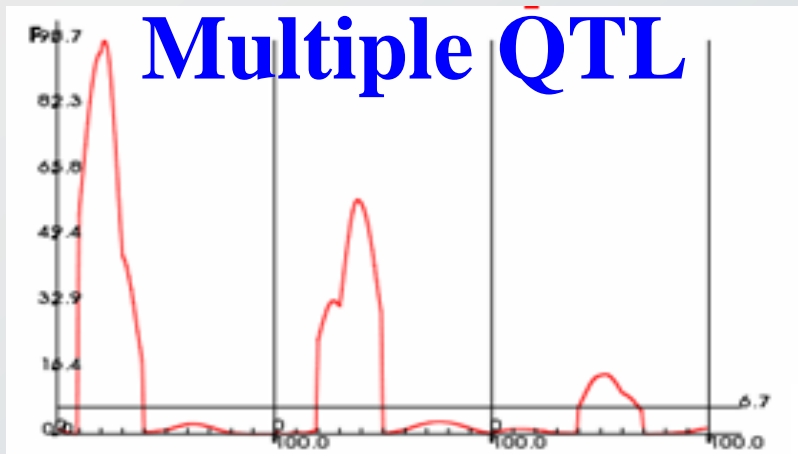


Fig. 1. Individual 3 QTL

QTL × Env

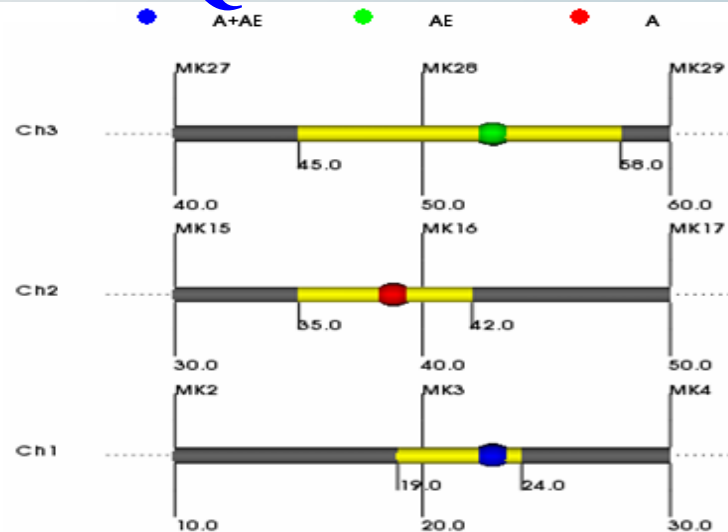


Fig. 2. Differential expression of 3 QTL

QTL Network

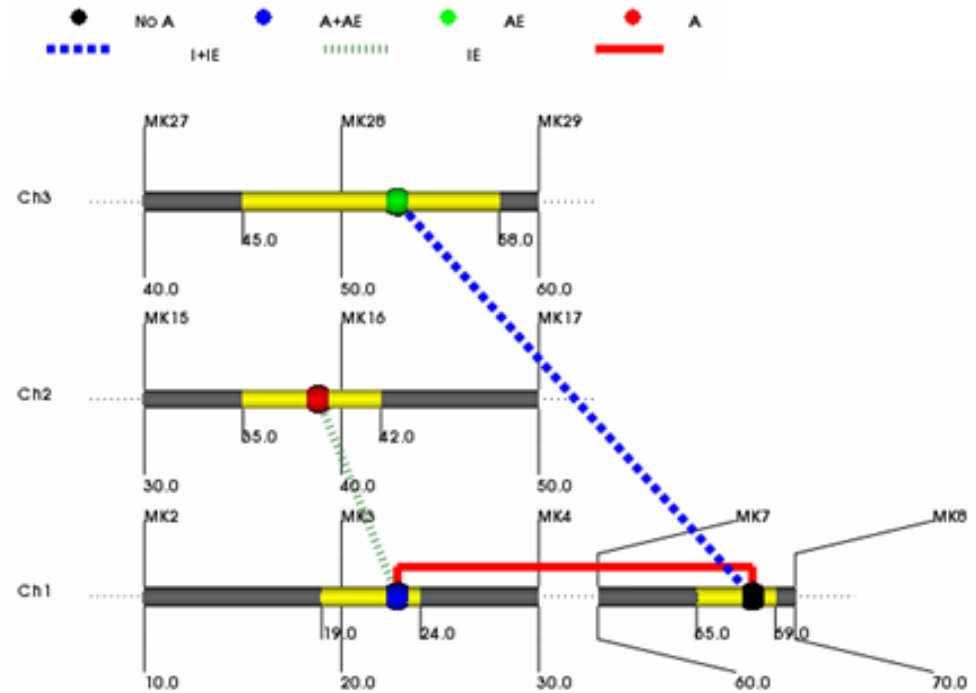
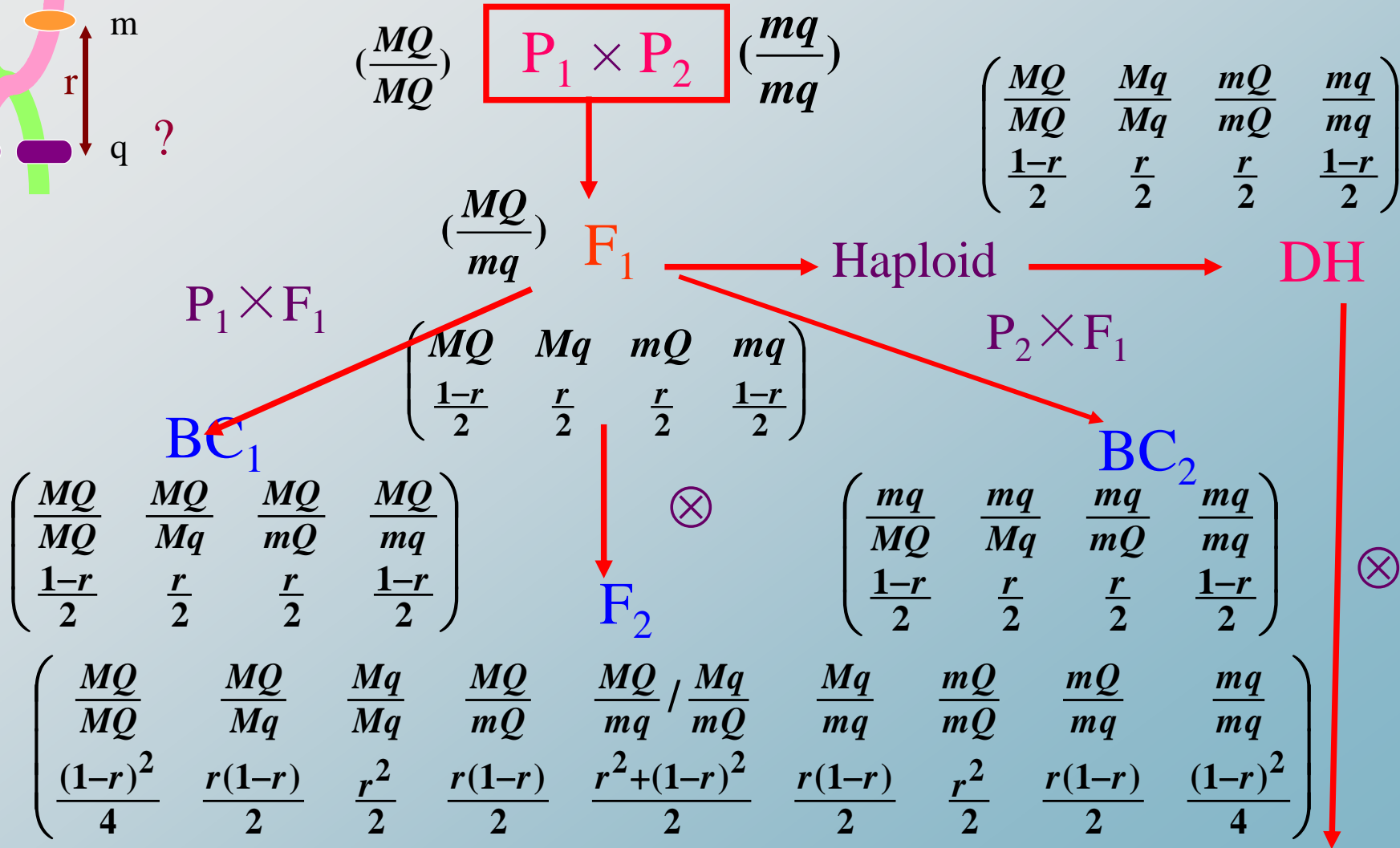
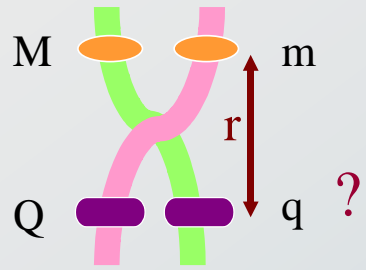


Fig. 3. Differential expression of QTL network

Populations for Mapping QTL



Experimental Design

(Observations $n = 250 \times 12 \times 3 = 9,000$)

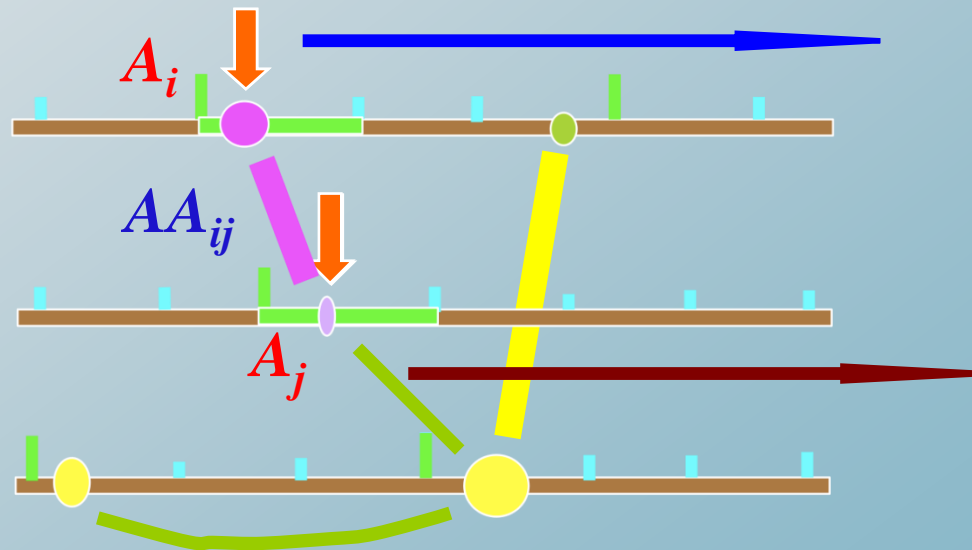
- **Mapping Population**
250 DH lines
- **Environments**
4 Location, 3 Years
- **Replications**
3 Blocks

Mixed-model-based Composite Interval Mapping

(Wang & Zhu et al. 1999, TAG, V99)

Mapping QTL with A+AA and QE Interaction (DH, RIL)

$$y = \mu + A_i + A_j + AA_{ij} + E + A_iE + A_jE + AA_{ij}E + G_M + G_{MM} + G_ME + G_{MM}E + \varepsilon$$



Mixed-Linear Model Approaches

$$\begin{aligned} \mathbf{y} &= \mathbf{X}\mathbf{b} + \mathbf{U}_E \mathbf{e}_E + \mathbf{U}_{QE} \mathbf{e}_{QE} \\ &\quad + \mathbf{U}_M \mathbf{e}_M + \mathbf{U}_{MM} \mathbf{e}_{MM} + \mathbf{U}_{ME} \mathbf{e}_{ME} + \mathbf{U}_{MME} \mathbf{e}_{MME} + \mathbf{e}_\varepsilon \\ &= \mathbf{X}\mathbf{b} + \sum \mathbf{U}_u \mathbf{e}_u \\ &\sim N(\mathbf{X}\mathbf{b}, \mathbf{V} = \sum_{u=1}^m \sigma_u^2 \mathbf{U}_u \mathbf{R}_u \mathbf{U}_u^T) \end{aligned}$$

The Likelihood Function for QTL Mapping Model

$$L(\mathbf{b}, \mathbf{V}) = (2\pi)^{-\frac{n}{2}} |\mathbf{V}|^{-\frac{1}{2}} \exp\left[-\frac{1}{2} (\mathbf{y} - \mathbf{X}\mathbf{b})^T \mathbf{V}^{-1} (\mathbf{y} - \mathbf{X}\mathbf{b})\right]$$

Estimation of QTL Main Effects

$$\hat{\mathbf{b}} = (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \mathbf{y}$$

Prediction of QTL-Environment Interaction Effects

$$\hat{\mathbf{e}}_u = \sigma_u^2 \mathbf{U}_u^T \left[\mathbf{V}^{-1} - \mathbf{V}^{-1} \mathbf{X} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^+ \mathbf{X}^T \mathbf{V}^{-1} \right] \mathbf{y}$$

Challenges in Computation for Inverses of Big Matrix \mathbf{V} (9000 × 9000)

- Rice Map Distance = 2031 cM
- Step of QTL Searching = 1 cM
- Steps of Two-dimension Search
= $2031 \times 2030 / 2 = 2.06$ Millions
- Detecting QTLs for 10 Traits Needs to
Calculate 20.6 Millions Inverses of \mathbf{V}

QTLNetwork version 2.0

QTLNetwork - Project1

Project(P) Edit(E) View(V) Windows(W) View Angle(U) Help(H) Setting(S)

Chromosome: ALL Trait: trait1

Project1

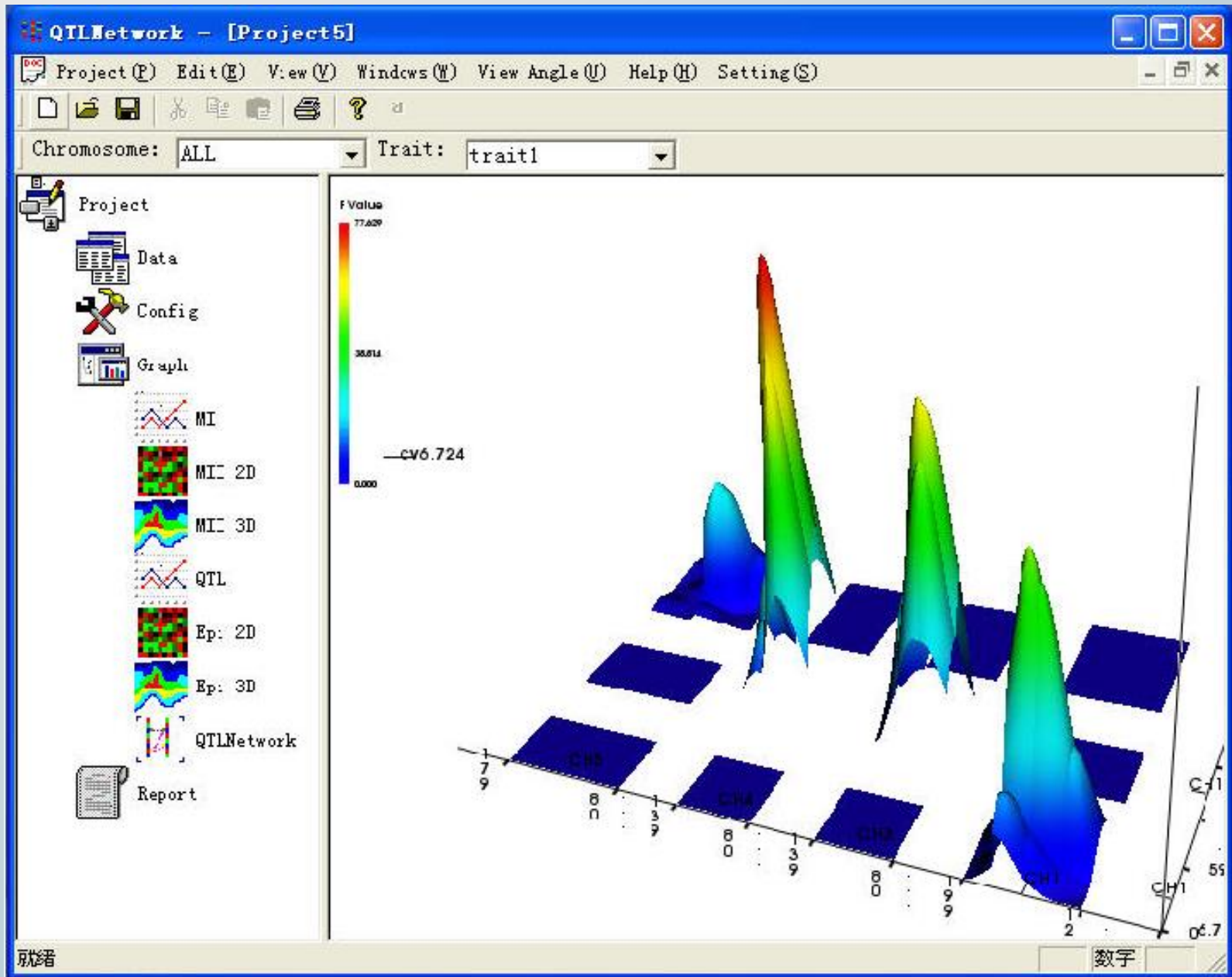
- Project
 - Data
 - Map File
 - Data File
 - Config
 - Graph
 - MI
 - MII 2D
 - MII 3D
 - QTL
 - Epi 2D
 - Epi 3D
 - QTLNetwork
 - Report

Population DH
_ Genotypes 200
_ Observations 400
_ Environments yes
_ Replications no
_ TraitNumber 1
_ TotalMarker 33
_ MarkerCode P1=A P2=B F1=H F1P1=C F1P2=D

MarkerBegin

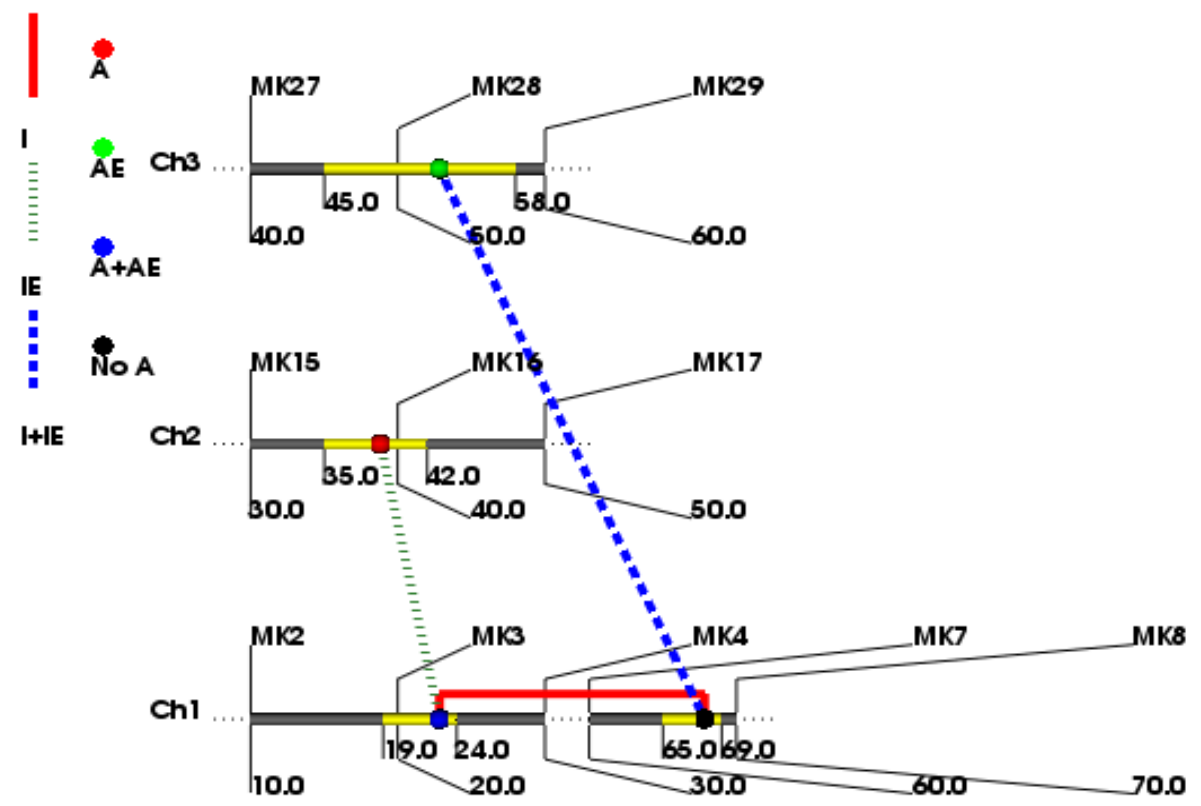
NIndi#	MK1	MK2	MK3	MK4	MK5	MK6	MK7	MK8	MK9	MK10	MK11	MK12	MK13
1	B	B	B	B	B	B	B	B	B	B	B	B	B
2	A	A	A	A	A	A	A	B	A	A	A	A	A
3	B	B	B	B	B	B	B	B	B	A	A	B	A
4	A	A	A	A	A	A	A	A	A	A	A	A	B
5	A	A	A	A	A	A	A	A	A	A	A	B	B
6	B	B	B	B	B	B	B	B	B	B	B	A	A
7	A	B	B	B	A	A	A	A	A	A	A	A	B
8	B	B	B	B	B	B	B	B	B	A	A	A	A
9	A	A	A	A	A	A	B	B	B	B	B	B	B
10	B	A	A	A	A	A	A	B	B	B	B	B	B
11	A	A	A	A	A	A	A	A	A	A	A	A	B
12	B	B	A	A	A	A	A	B	B	B	B	A	A

QTLNetwork 2.0



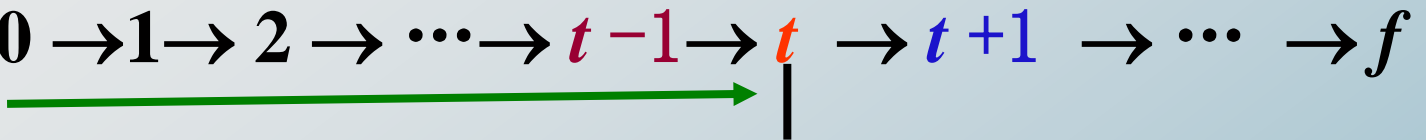
Chromosome: ALL Trait: trait1

- Project
- Data
- Map File
- Data File
- Config
- Graph
- MI
- MII 2D
- MII 3D
- QTL
- Epi 2D
- Epi 3D
- QTLNetwork
- Report



Mapping Developmental QTL for Complex Traits

Time: $0 \rightarrow 1 \rightarrow 2 \rightarrow \dots \rightarrow t-1 \rightarrow t \rightarrow t+1 \rightarrow \dots \rightarrow f$



Unconditional Model for Phenotypic Value at Time t

$$y(t) = \mu(t) + G_Q(t) + E(t) + G_Q E(t) \\ + G_M(t) + G_M E(t) + \varepsilon(t)$$

Analyzing Q & QE Effects from Time $0 \rightarrow t$

Conditional Model for Phenotypic Value at Time t

$$y(t|t-1) = \mu(t|t-1) + G_Q(t|t-1) + E(t|t-1) \\ + G_Q E(t|t-1) + G_M(t|t-1) + G_M E(t|t-1) + \varepsilon(t|t-1)$$

Analyzing Net Q & QE Effects From Time $t-1 \rightarrow t$

Table 2. Chromosomal regions and estimated genetic effects of QTLs for plant height (cm) at different stages in two environments. (Yan, *et al.* 1998, Genetics, V150)

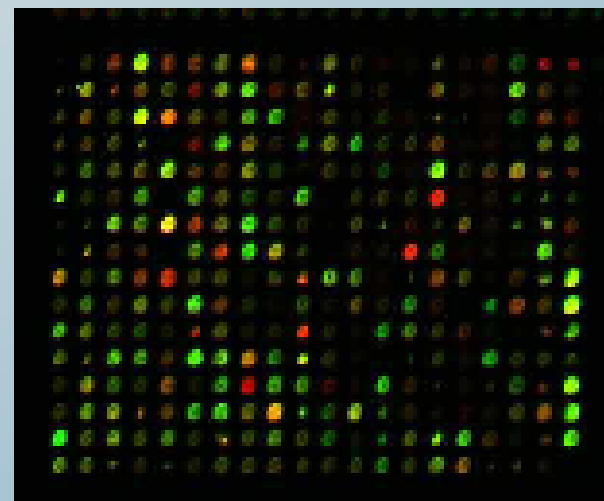
QTL	Maker Interval	Distance (cM)	Days	GE					
				Main effect		in Hangzhou		in Hainan	
				<i>t</i>	<i>t</i> <i>t</i> -1	<i>t</i>	<i>t</i> <i>t</i> -1	<i>t</i>	<i>t</i> <i>t</i> -1
Ph1	RZ730- RZ801	33.1	10	-2.14	-2.14			-1.08	-1.08
			20	-2.47				-0.95	
			30	-3.6				-1.45	-0.52
			40	-4.6				-3	-0.52
			50	-5.06			-1.68	-3.12	1.05
			60	-9.54	-0.55	-3.09	-1.7	-2.64	1.1
			70	-12.57		-4.39		-3.51	
			80	-15.01		-4.47		-4.35	
			90	-16.98		-4.07		-4.36	
Ph2	Amy1A/C -RG95	12.8	10	1.03	1.03			0.63	0.63
			20	1.51				0.76	
			30	1.93				0.93	
			40	2.91		0.99		1.25	
			50	2.77	0.92		-0.74	1.05	-0.74
			60	4.74		2.55			
			70	6.07		2.67			
			80	7.73		2.42		1.93	
			90	7		1.62		1.69	

Challenges in Presentation of G-G and G-E During Developmental Stages

- For G-G interaction presentation, there needs two-dimension display
- When Genes express differently across times and spaces, there needs four-dimension display or three-dimension dynamic display

Experimental Design for Microarray Testing

- **Array Design :**
 - 14112 Genes (84 x 68)/Array
- **Treatment Design (3 Way Factors) :**
 - **Factor G: 14112 Genes**
 - **Factor C: 8 Cancer Cells**
 - **Factor M: 7 Medicine**
- **Replicates: 3**



Combining Analysis

$$y_{ijkl} = G_l + A_i + C_j + M_k + GA_{li} + GC_{lj} + GM_{lk} + \varepsilon_{ijkl}$$

Challenges in Mixed-Model Approaches for Array Analysis

$$\begin{aligned} \mathbf{y} &= \mathbf{X}\mathbf{b}_G + \mathbf{U}_A \mathbf{e}_A + \mathbf{U}_C \mathbf{e}_C + \mathbf{U}_M \mathbf{e}_M \\ &\quad + \mathbf{U}_{GA} \mathbf{e}_{GA} + \mathbf{U}_{GC} \mathbf{e}_{GC} + \mathbf{U}_{GM} \mathbf{e}_{GM} + \mathbf{e}_\varepsilon \\ &= \mathbf{X}\mathbf{b}_G + \sum_{u=1}^7 \mathbf{U}_u \mathbf{e}_u \sim N(\mathbf{X}\mathbf{b}_G, \mathbf{V}_{(n \times n)} = \sum_{u=1}^7 \sigma_u^2 \mathbf{U}_u \mathbf{U}_u^T) \end{aligned}$$

Prediction of Random Effects

$$\begin{aligned} \ddot{\mathbf{e}}_{u(1)} &= \mathbf{U}_u^T \mathbf{V}_{(1)}^{-1} (\mathbf{y} - \mathbf{X}\mathbf{b}) \\ &= \mathbf{U}_u^T \left[\mathbf{V}_{(1)}^{-1} - \mathbf{V}_{(1)}^{-1} \mathbf{X} (\mathbf{X}^T \mathbf{V}_{(1)}^{-1} \mathbf{X})^+ \mathbf{X}^T \mathbf{V}_{(1)}^{-1} \right] \mathbf{Q}_{(1)} \mathbf{y} \end{aligned}$$

Experiment Size

$$n = 14112 \times 8 \times 7 \times 3 = 2.37 \text{ Millions}$$

A Two-step Strategy for Detecting Differential Gene Expression of cDNA Microarray Data

(Lu, Zhu, &, Liu, 2005, Current Genetics. 47: 121–131)

Choosing a subset of potential genes with differential expression

$$y_{ijk(l)} = \mu_{(l)} + A_{i(l)} + V_{j(l)} + D_{k(l)} + \gamma_{ijk(l)}$$

Combining analysis of multiple genes

$$y_{ijkl} = G_l + A_i + C_j + M_k + GA_{li} + GC_{lj} + GM_{lk} + \varepsilon_{ijkl}$$

Experiment Size Reduced to ≈ 300

$$n = 300 \times 8 \times 7 \times 3 \approx 50,400$$

Thank You!