

Design and Analysis of Gene Expression Microarray Experiments

M. Kathleen Kerr

Gary A. Churchill

The Jackson Laboratory

OUTLINE

Experimental Design for Microarrays

Analysis of Variance for Microarrays

Identifying Differentially Expressed Genes

Corning Data Analysis

Synteni Data Analysis

Chu et al. Data Analysis

“Statistical procedure and experimental design are only two different aspects of the same whole, and that whole comprises all the logical requirements of the complete process of adding to natural knowledge by experimentation”

-Fisher 1935

“It is possible, and indeed it is all too frequent, for an experiment to be so conducted that no valid estimate of error is available”

-Fisher 1935

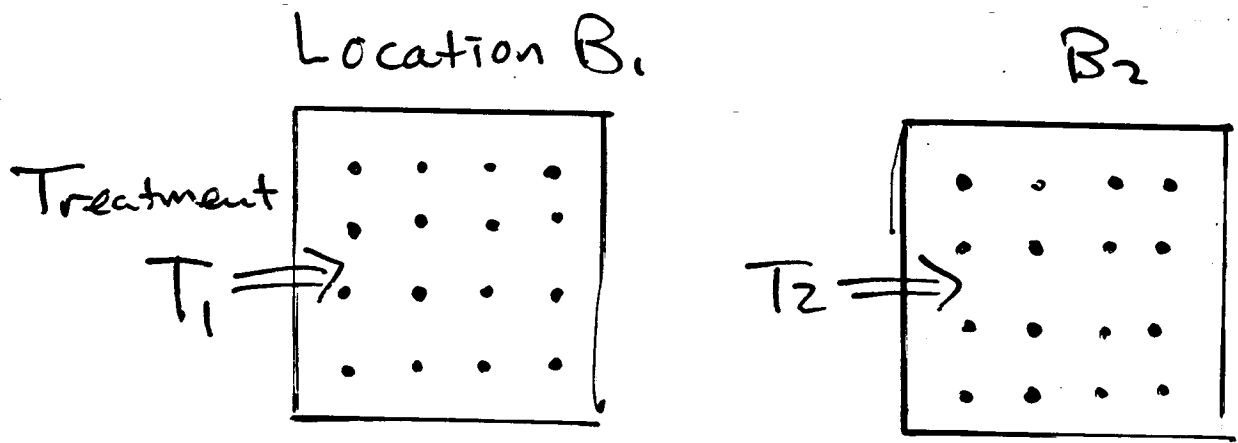
"A well-designed experiment will usually allow its conclusions to be easily obtained, whereas no computations, however industriously or ingeniously performed, can produce entirely satisfactory conclusions from an ill-designed one.

Considerable tact is needed in discussion of these matters."

Finney 1953

The “design” of an experiment is

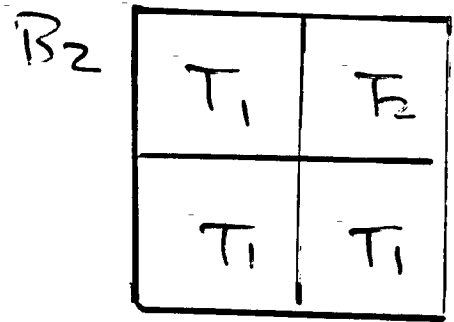
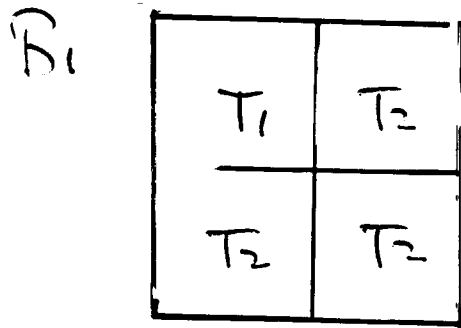
1. the set of treatments selected for comparison
2. the specification of the units to which the treatments will be applied
3. the rules by which the treatments are allocated to the experimental units
4. the specification of the measurements to be made.



Confounding

$$\text{Var}(T_1 - T_2) = 2\sigma^2$$

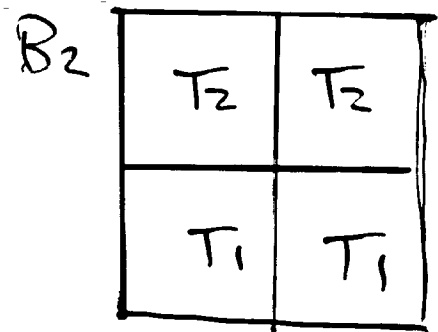
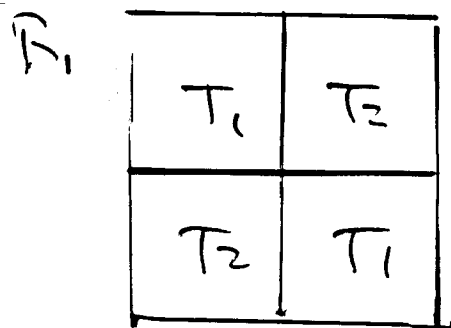
$$\text{error } dF = 0$$



Partial Confounding

$$\text{Var}(T_1 - T_2) = \frac{2}{3}\sigma^2$$

$$\text{error } dF = 6$$



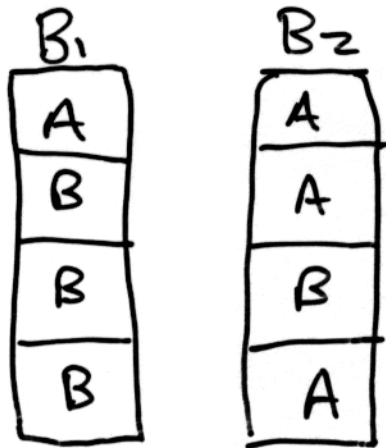
Orthogonality

$$\text{Var}(T_1 - T_2) = \frac{1}{2}\sigma^2$$

$$\text{error } dF = 6$$

Relative Efficiency

unbalanced



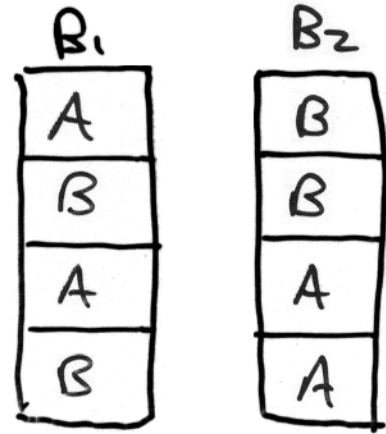
$$X = \begin{bmatrix} 1 & 1 \\ -1 & 1 \\ -1 & 1 \\ -1 & 1 \end{bmatrix}$$

$$XX' = \begin{bmatrix} 8 & 0 & 0 \\ 0 & 8 & -4 \\ 0 & 4 & 8 \end{bmatrix}$$

$$X'X = \begin{bmatrix} 25 & 0 & 0 \\ 0 & 167 & .083 \\ 0 & .083 & .167 \end{bmatrix}$$

$$\sigma^2(T_A T_B) = 4 \times .167 = 0.667$$

balanced



$$X = \begin{bmatrix} 1 & 1 \\ -1 & 1 \\ -1 & 1 \\ -1 & 1 \end{bmatrix}$$

$$XX' = \begin{bmatrix} 8 & 0 & 0 \\ 0 & 8 & 0 \\ 0 & 0 & 8 \end{bmatrix}$$

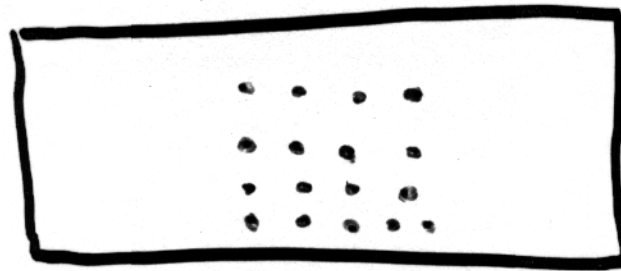
$$X'X = \begin{bmatrix} 25 & 0 & 0 \\ 0 & 25 & 0 \\ 0 & 0 & 125 \end{bmatrix}$$

$$4 \times 125 = 0.5$$

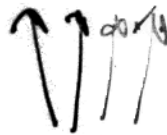
Spotted cDNA Microarrays

Treat.


Ctrl

array on
glass sub



cDNAs

represent genes

Design Factors in a Microarray Experiment:

A_i arrays $i = 1, \dots, K$

D_j dyes $j = 1, 2$

G_g genes $g = 1, \dots, N$

V_k varieties $k = 1, \dots, K$

S_r spots r

Genes

- main effects G

average expression level
dye incorporation
hybridization efficiency

- array by gene $A \times G$

→ variations in spot size

- dye by gene $D \times G$

→ the little green spot

- variety by gene $V \times G$

→ variety specific gene effect

VG_{ig} VG_{jg} = change in expression

Structure of Design Space

- $G \perp A, D, V$
- $AG \perp DG$
- AG, VG partial confounded
 DG, VG can be \perp or partial confounded
- High order interactions
 $A * D * C$ $A * V * G$
 $D * V * C$ $A * D * V * G$
are confounded with L.O.Ts

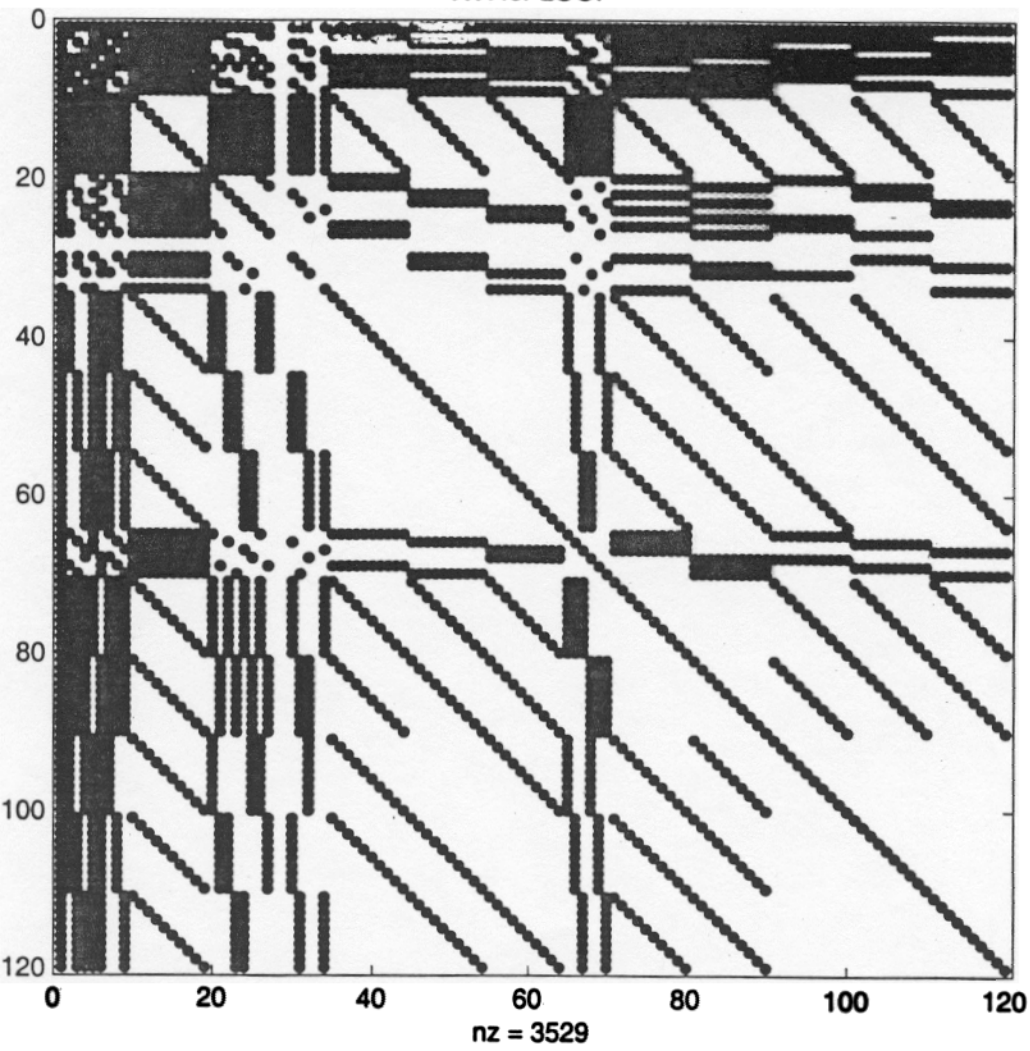
Arrays, Dyes and Varieties

- Interaction effects

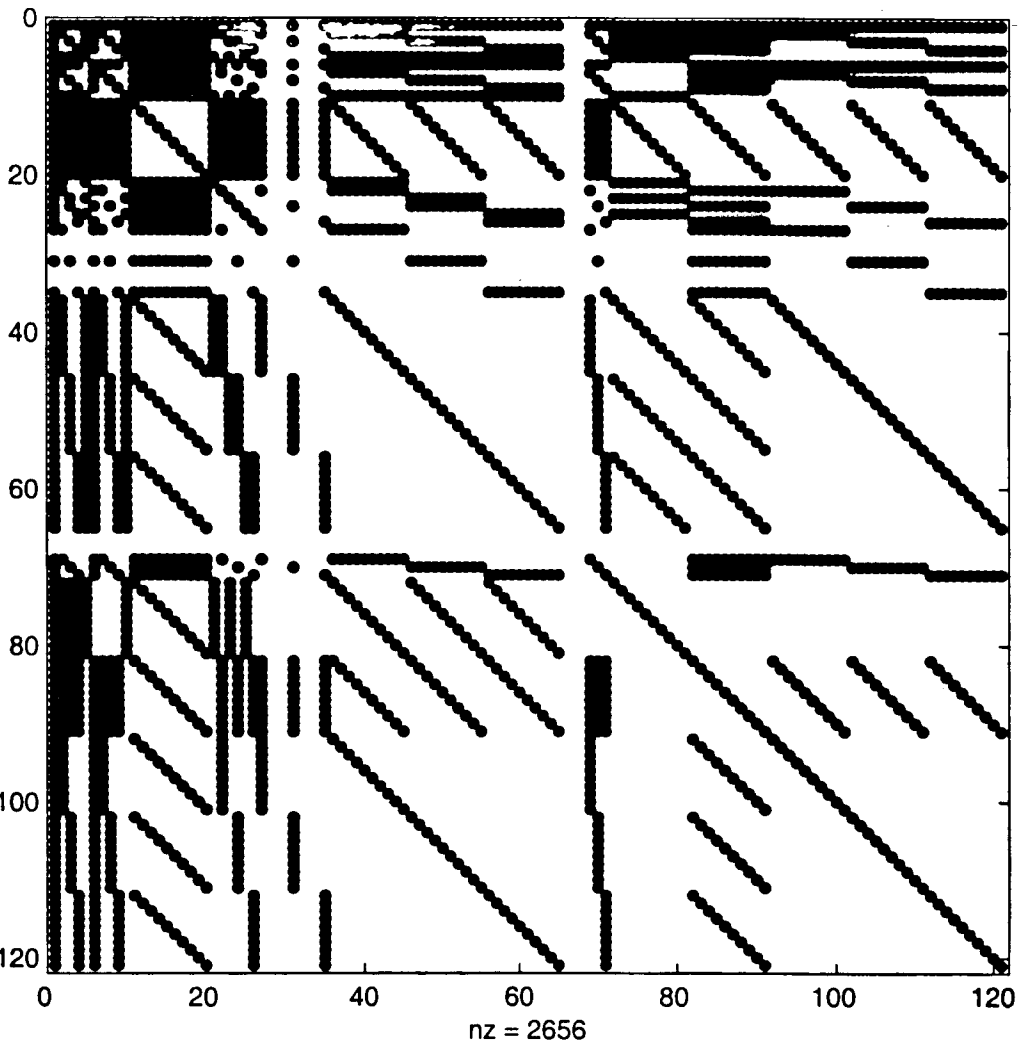
$A * D$, $A * V$, $D * V$, $A * D * V$
are confounded with the
main effects A, D, V

A and D are always balanced

- D and V can be balanced
- A and V cannot be balanced
if $k > 2$

$X^T X$ for LOOP

$X^T X$ for REF



A simple ANOVA model includes only the factor main effects and the effects of interest, variety \times gene:

$$y_{ijk_g} = \mu + A_i + D_j + V_k + G_g + (VG)_{kg} + \epsilon_{ijk_g}. \quad (4.1)$$

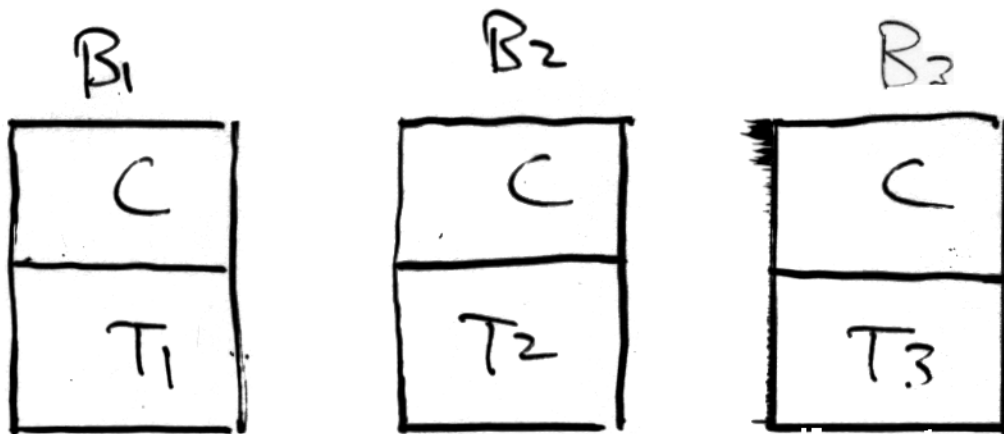
A more plausible model accounts for spot-to-spot variation by including array \times gene effects:

$$y_{ijk_g} = \mu + A_i + D_j + V_k + G_g + (VG)_{kg} + (AG)_{ig} + \epsilon_{ijk_g}. \quad (4.2)$$

Another possibility is to further account for genes interacting with dyes:

$$y_{ijk_g} = \mu + A_i + D_j + V_k + G_g + (VG)_{kg} + (AG)_{ig} + (DG)_{jg} + \epsilon_{ijk_g}. \quad (4.3)$$

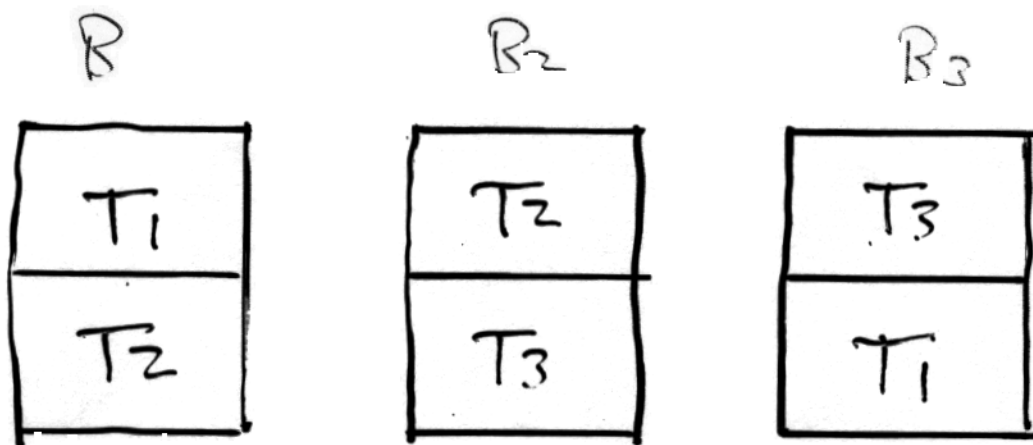
We assume that there is independent, additive error $\epsilon_{ijk_g} \sim F$, where F is a distribution with mean 0 and variance σ^2 .



Comparison to a common control

$$\text{Var}(T_A \ T_B) = 4\sigma^2$$

$$\text{error } dF = 0$$

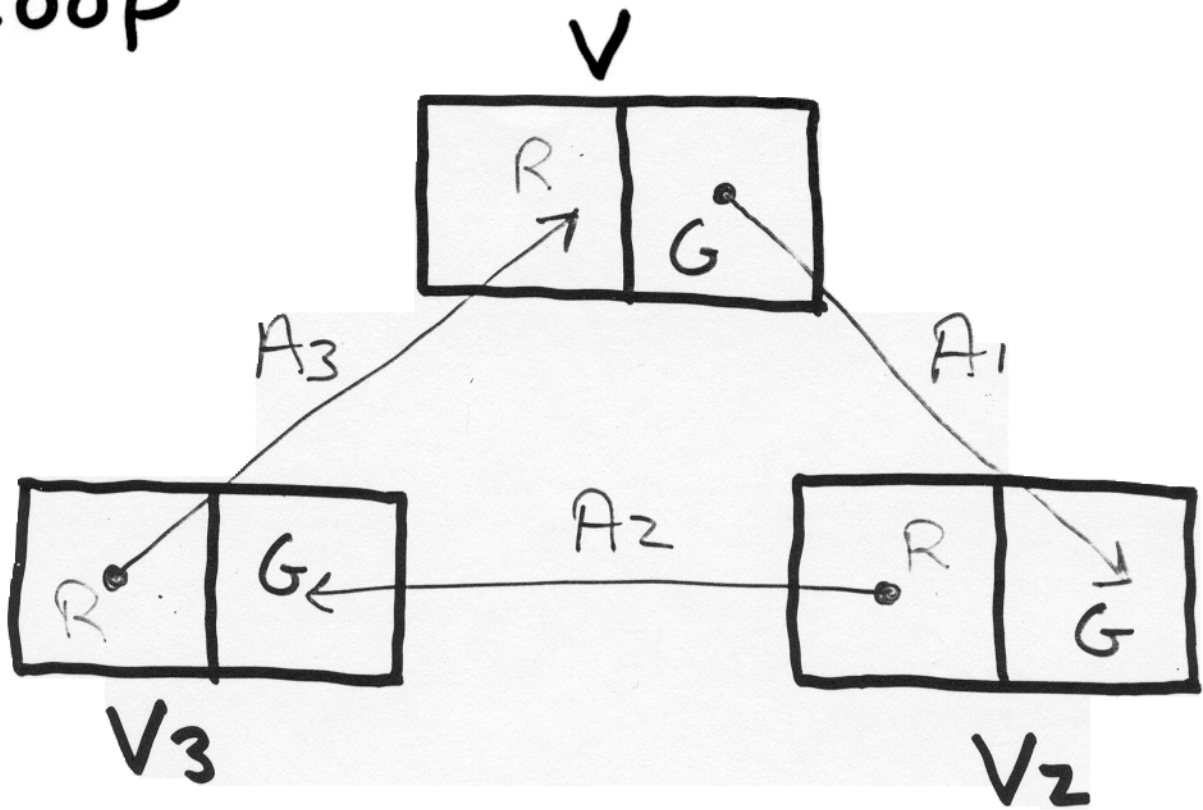


Balanced incomplete blocks

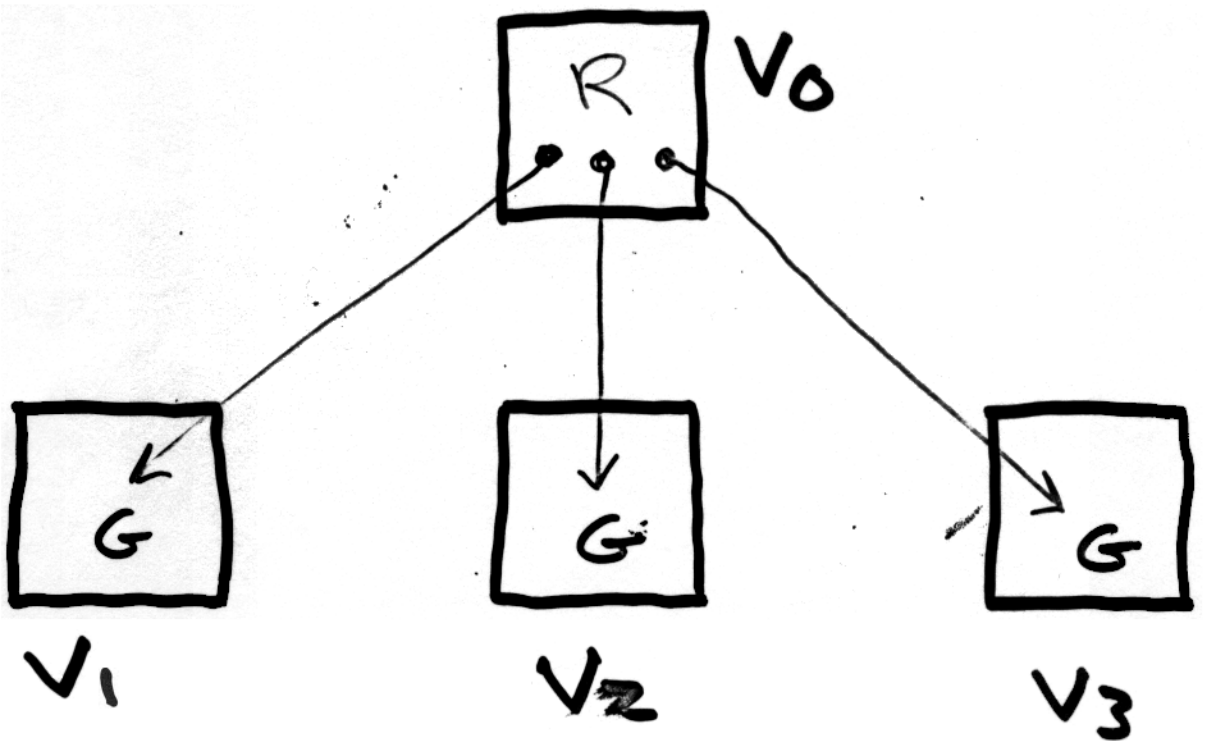
$$\text{Var}(T_A \ T_B) = 4/3 \sigma^2$$

$$\text{error } dF =$$

Loop



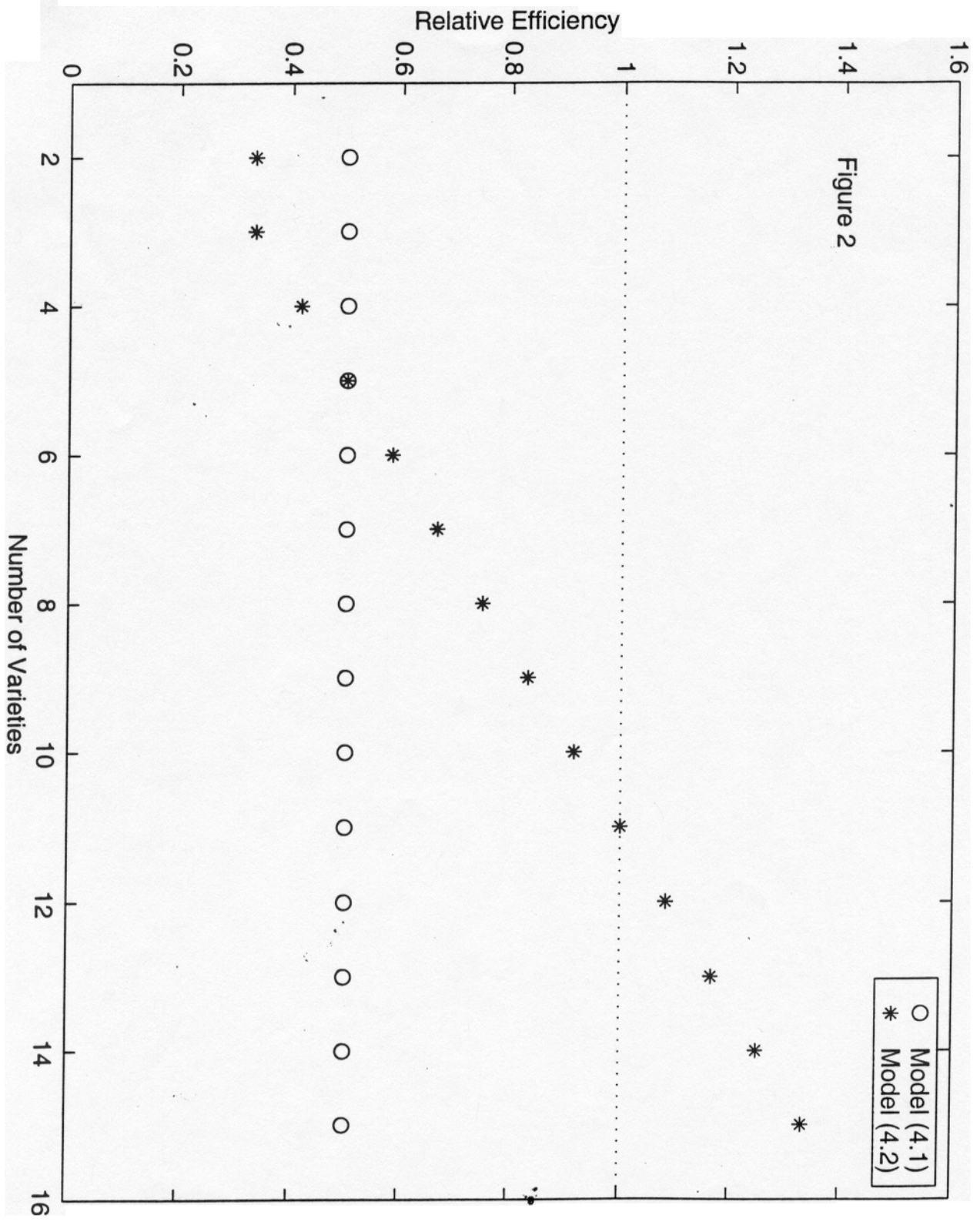
Reference



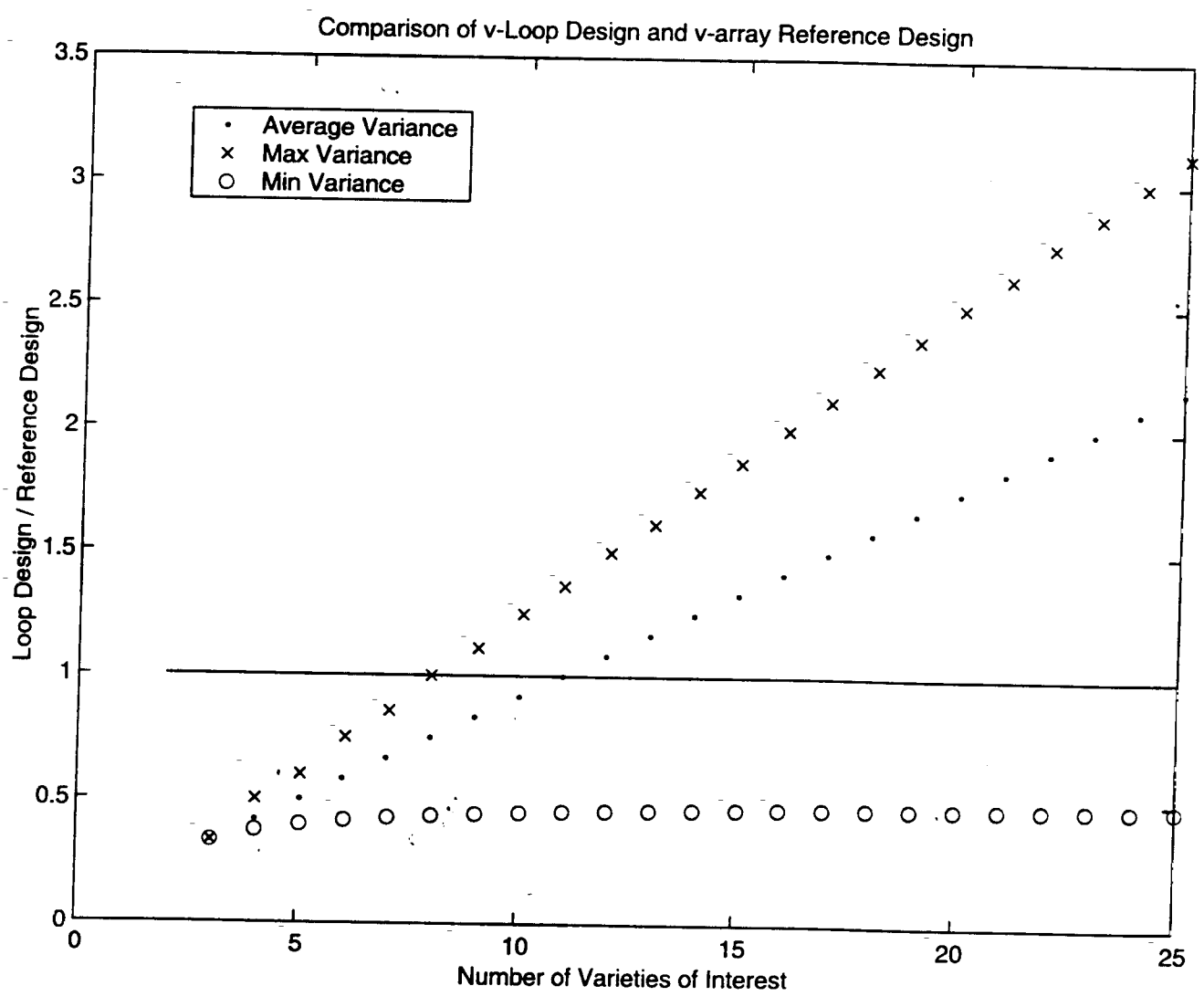
Source	Reference df	Loop df
A,D,T	$2K-1$	$2K-1$
G	$N-1$	$N-1$
TG	$K(N-1)$	$(K-1)(N-1)$
AG	$(K-1)(N-1)$	$(K-1)(N-1)$
residual	0	$N-1$
total	$2KN-1$	$2KN-1$

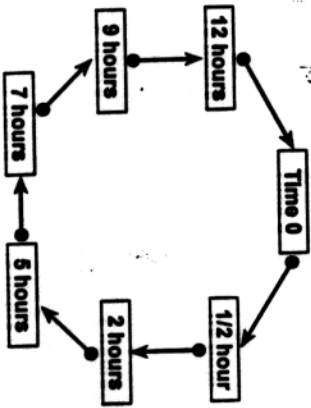
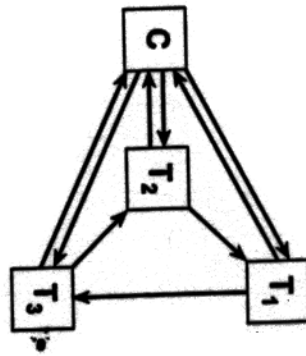
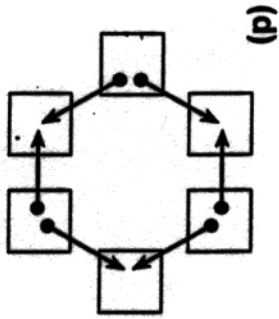
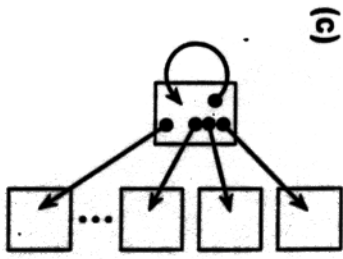
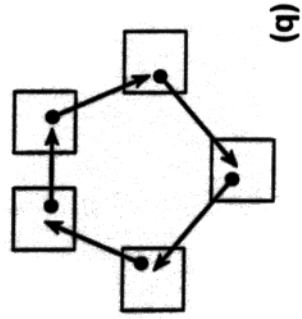
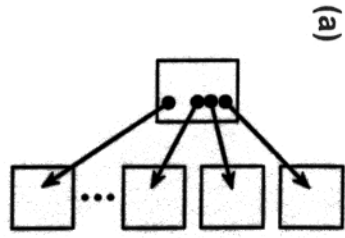
ANOVA table for designs with
K arrays and N genes.

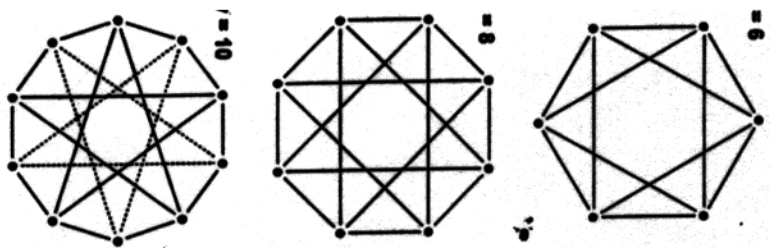
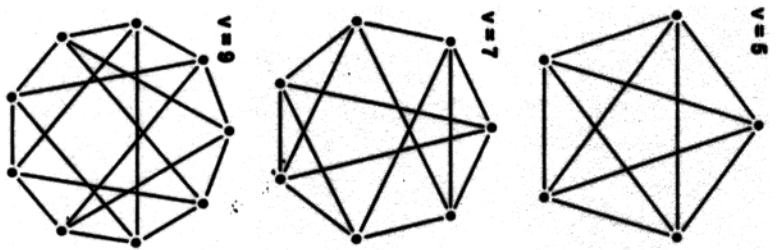
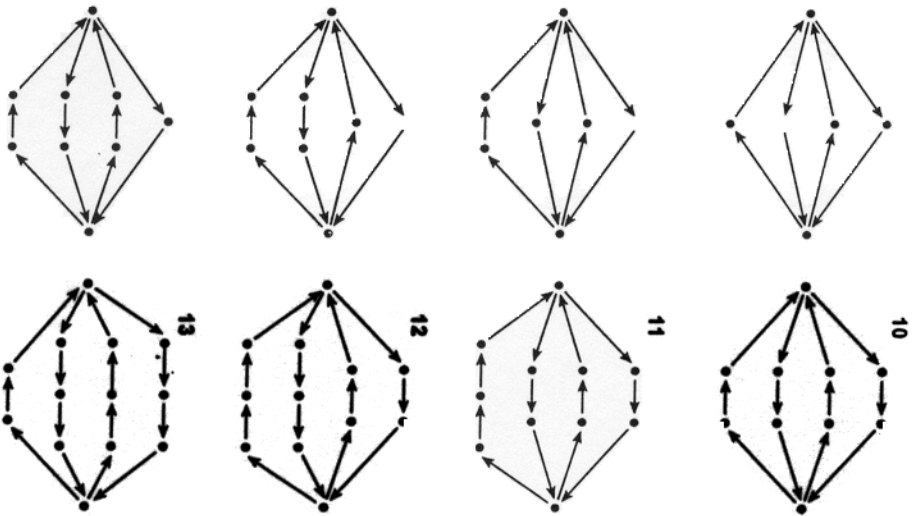
Loop vs Reference Avg. Variance



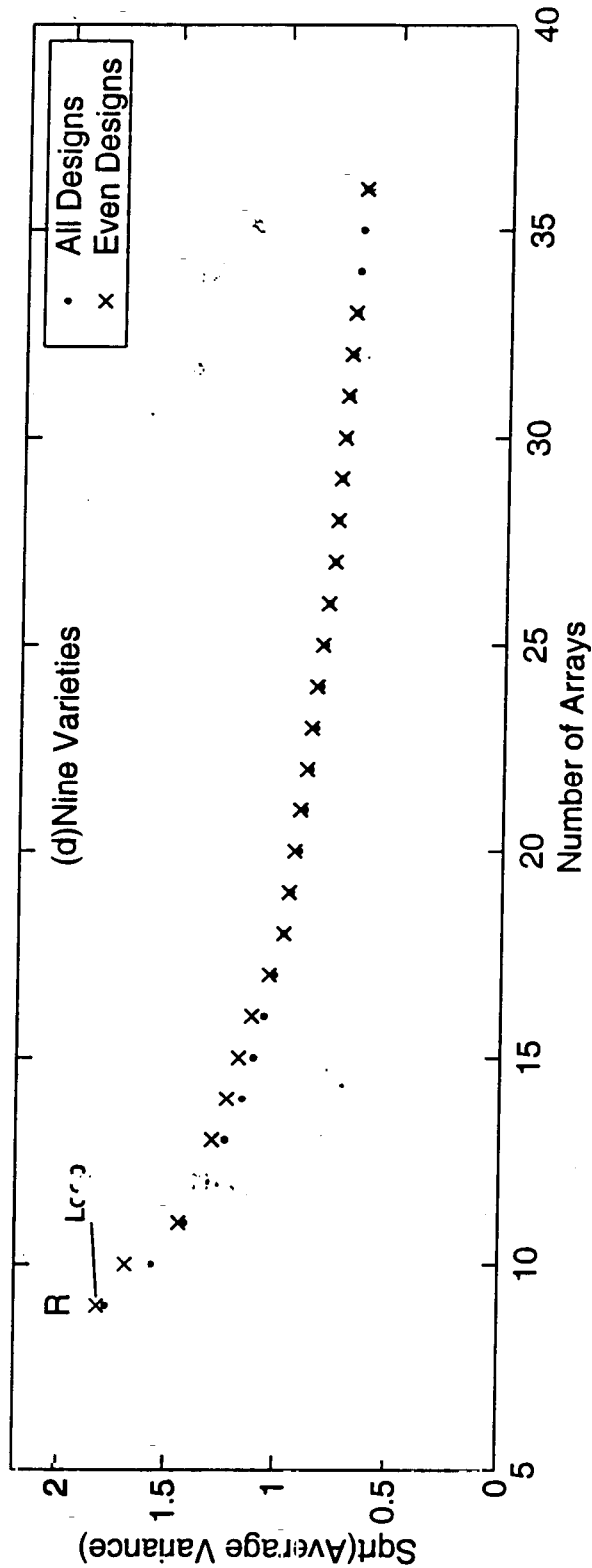
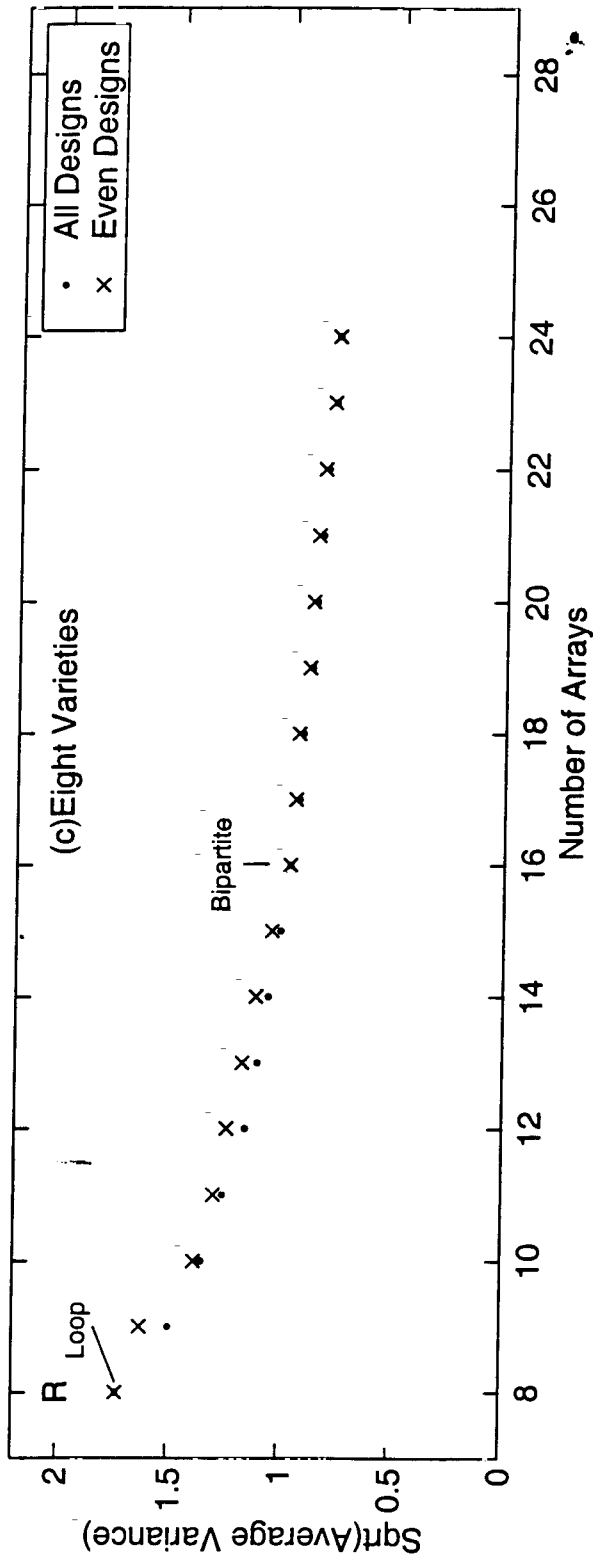
Comparison with k







Optimal Designs



Factor A

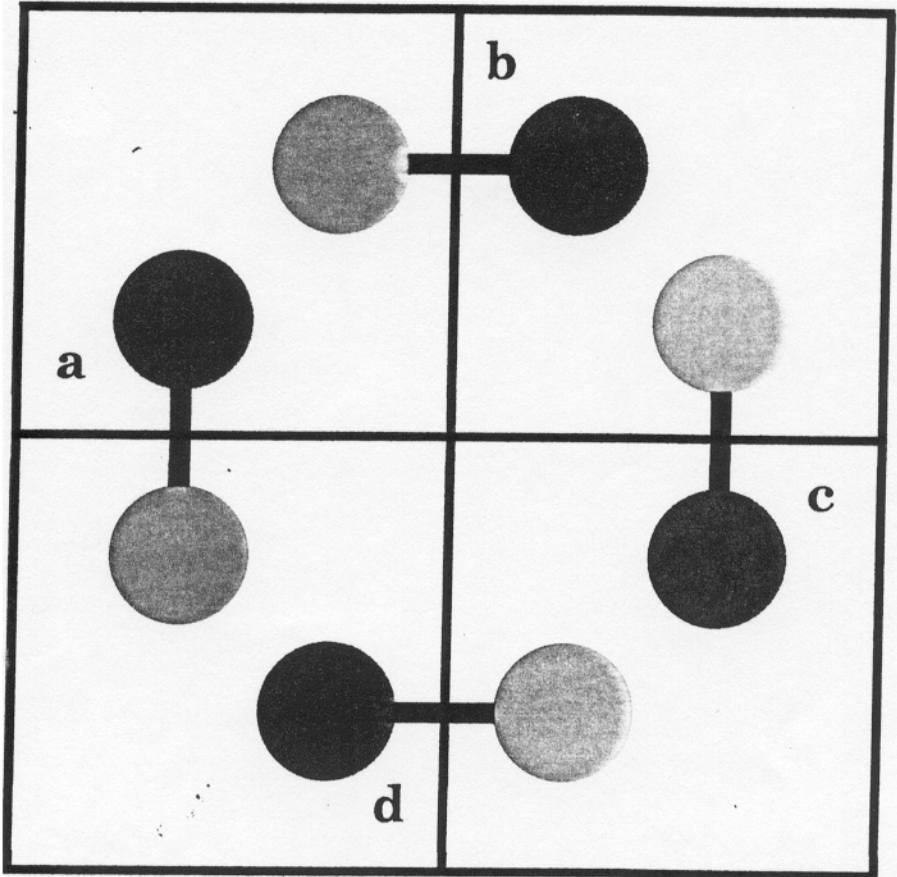
0

1

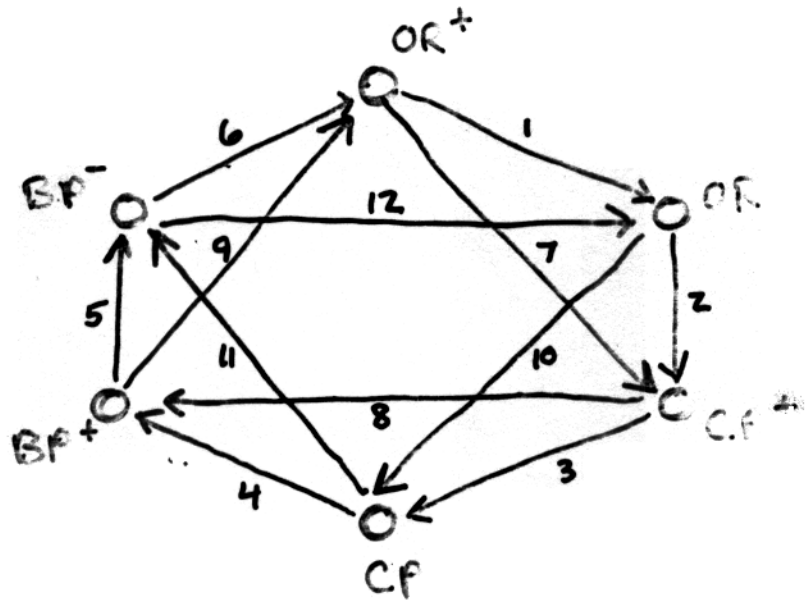
0

Factor B

1



Gary Churchill, Ph.D.
The Jackson Laboratory



i	j	k	l
background	genoty	array	cyr
BP	+/+	2	
CP	/+		
OR		2	

pairwise comparisons

$$OR^+ \quad OR \quad \frac{\gamma_7 \gamma_7}{\gamma_6 \gamma_9} \quad \frac{\gamma_2 \gamma_{11}}{\gamma \gamma_{10}}$$

contrasts

$$(CP^+ \quad BP \quad OR^+) \quad (CP \quad BP \quad OP^-)$$

multiple df comparisons

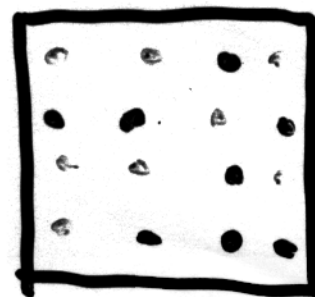
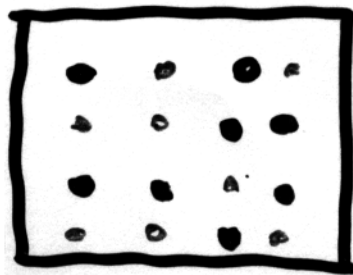
$$BP^- \quad CP \quad OP^-$$

Concluding Remarks

- The key concept that we left out:

Randomization

- The critical question that no one asked
IF the design is so bad, why do we get the right answers?



"One characteristic common to all biological material is that it varies."

-Finney 1953

"If I had to replicate my experiments, I could only do half as much"

-Botstein 1999

Synteni Data Analysis

Latin Square Design

Dye	Array	
	1	2
Red	Liver	Muscle
Green	Muscle	Liver

Analysis of Variance

Source	df	SS	MS
Array	1	92.34	92.34
Dye	1	0.74	0.74
Tissue	1	2.97	2.97
Gene	1285	1885.89	1.47
Gene x Tissue	1285	1357.28	1.06
Residual	2570	242.76	0.09
Corrected Total	5143	3581.99	

Synteni Data Analysis

Reference Sample Design

Dye	Array	
	1	2
Red	Placenta	Placenta
Green	Liver	Muscle

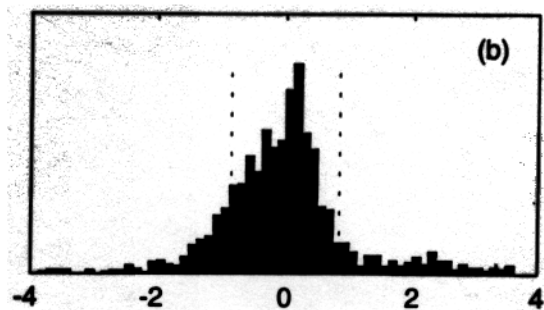
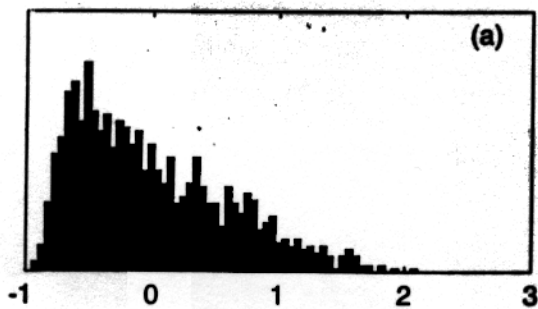
Analysis of Variance

Source	df	SS	MS
Array, Tissue	3	761.97	253.99
Gene	1904	3394.17	1.78
Genex Tissue	3808	1264.43	0.33
Residual	1904	55.21	0.03
Corrected Total	7619	5475.78	

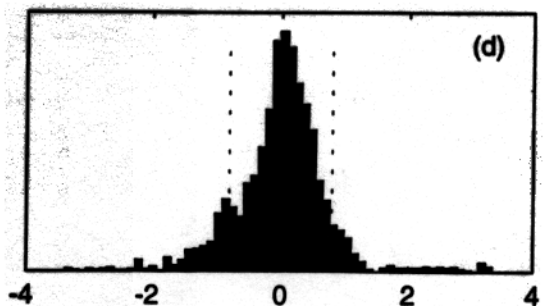
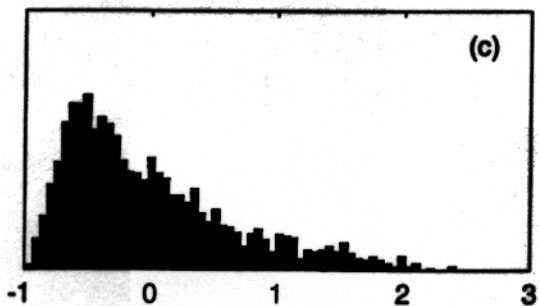
G effects

$V \times G$ effects

Loop

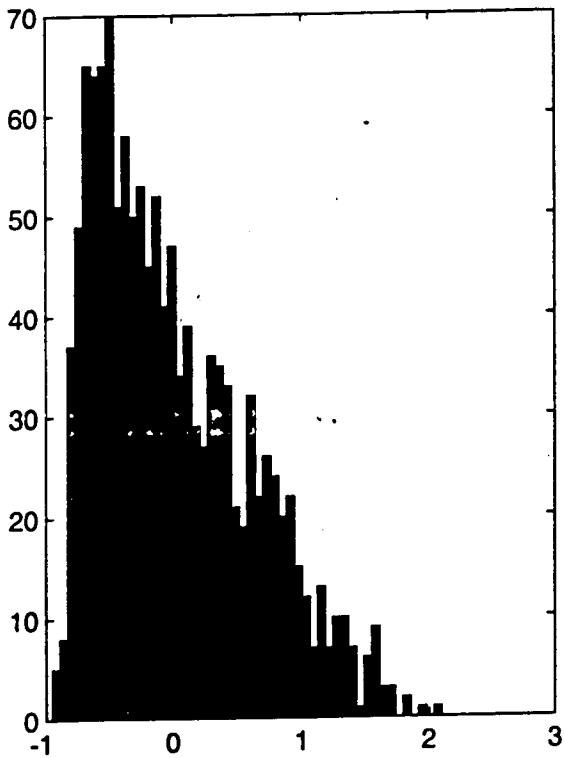


ReF

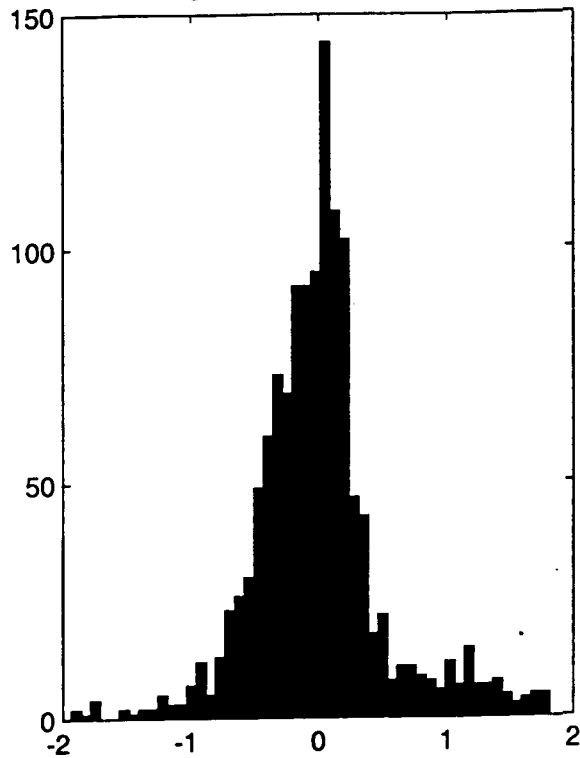


SYNTENI LATIN SQUARE EXPERIMENT

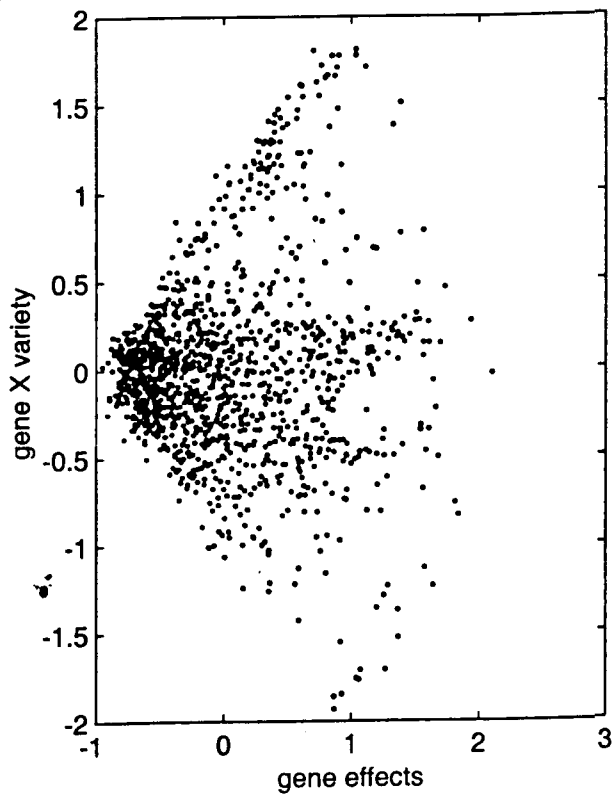
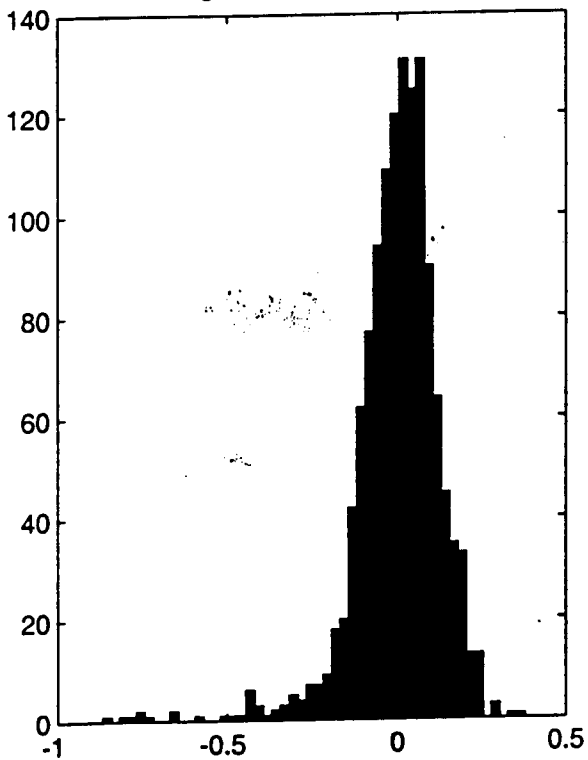
gene effects

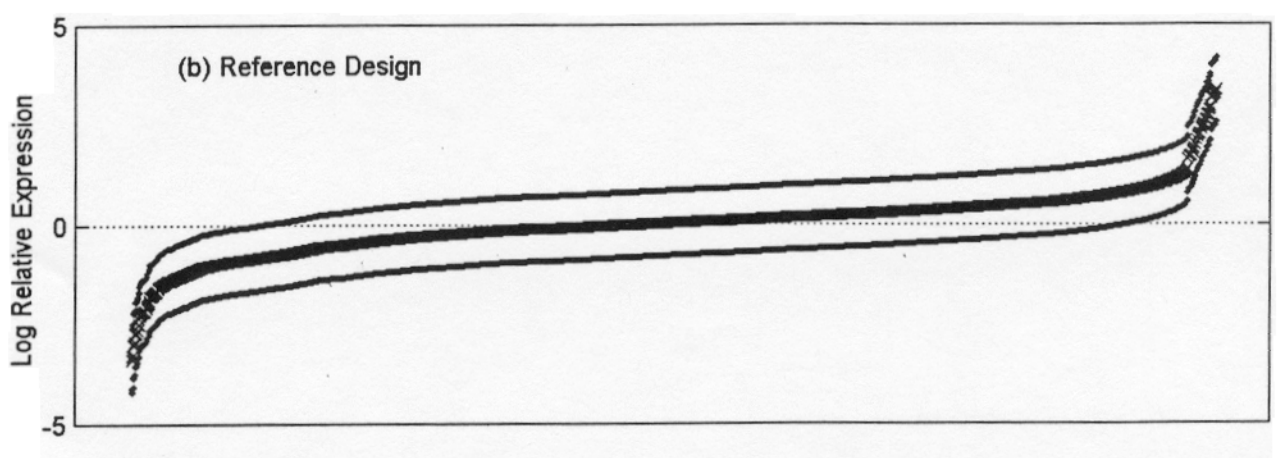
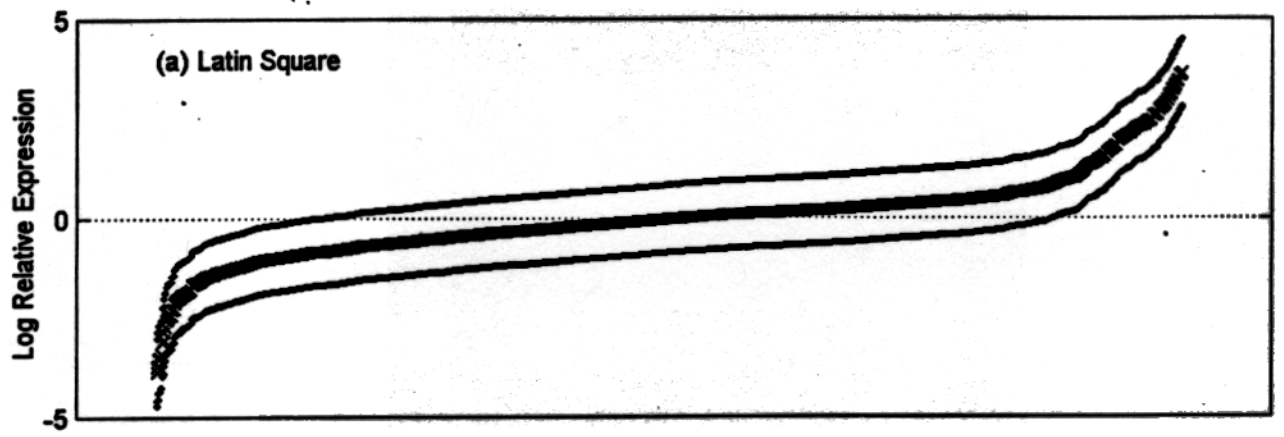


gene X variety effects

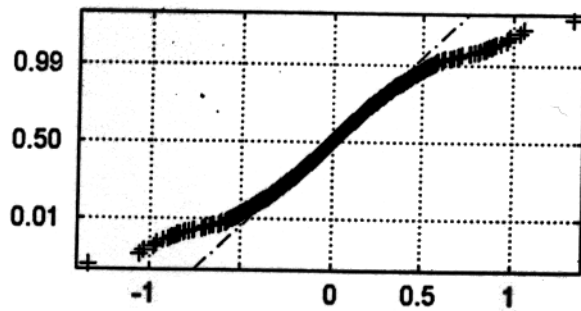


gene X dye effects

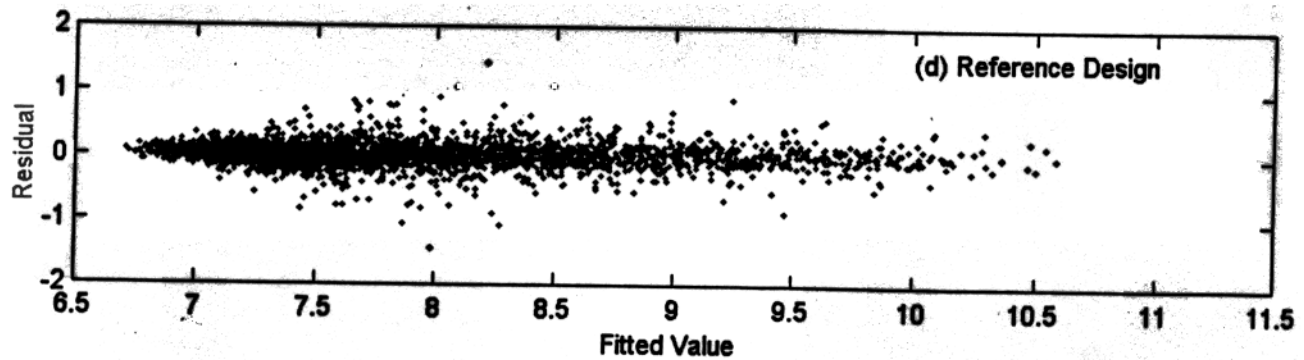
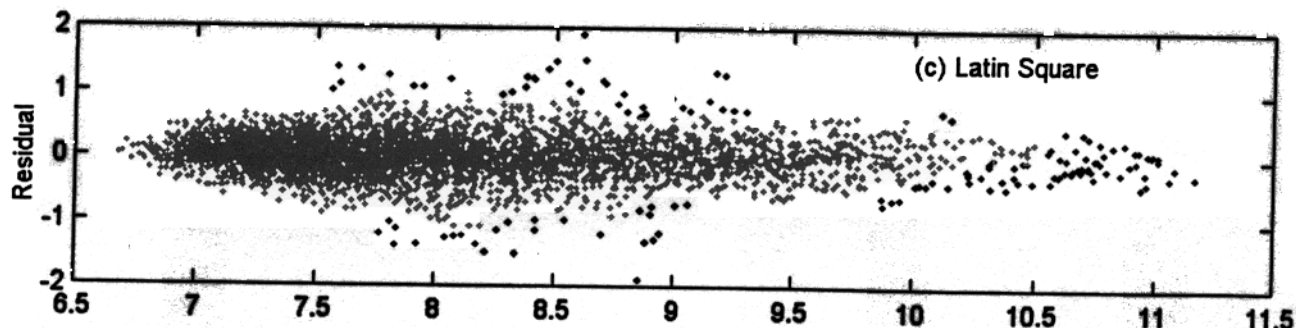
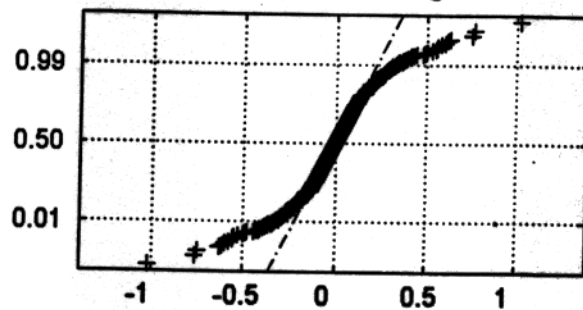




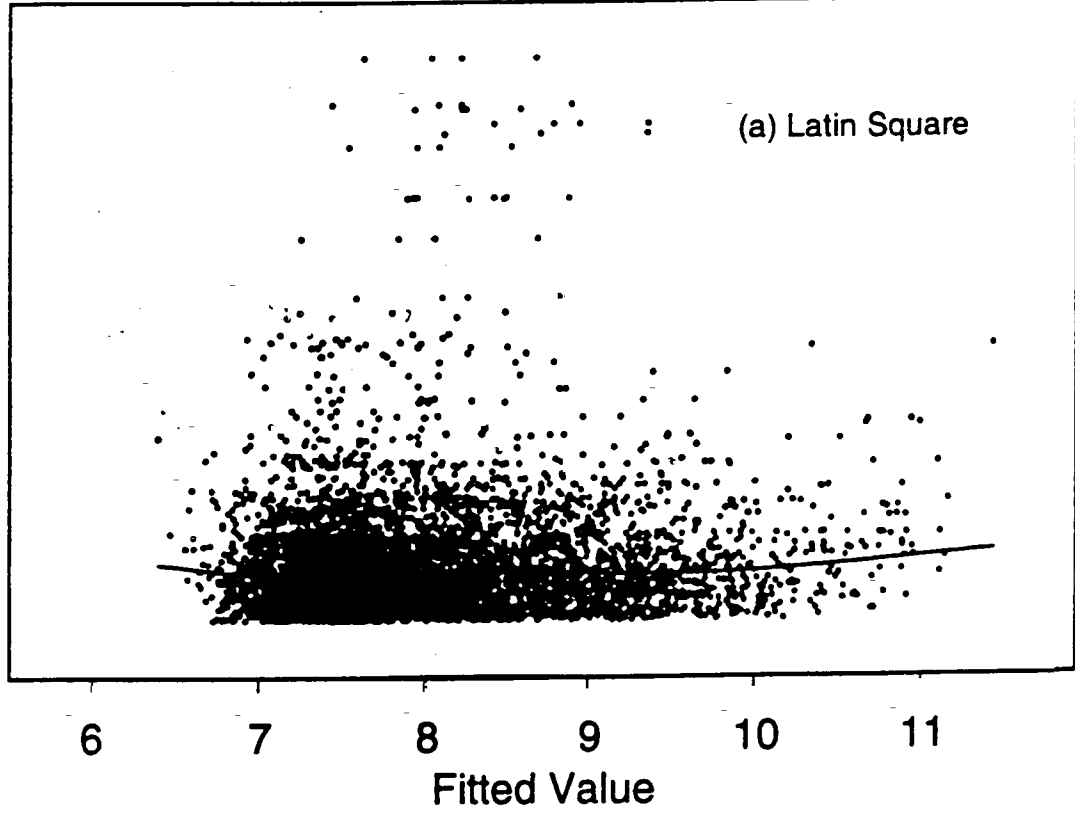
(a) Latin Square



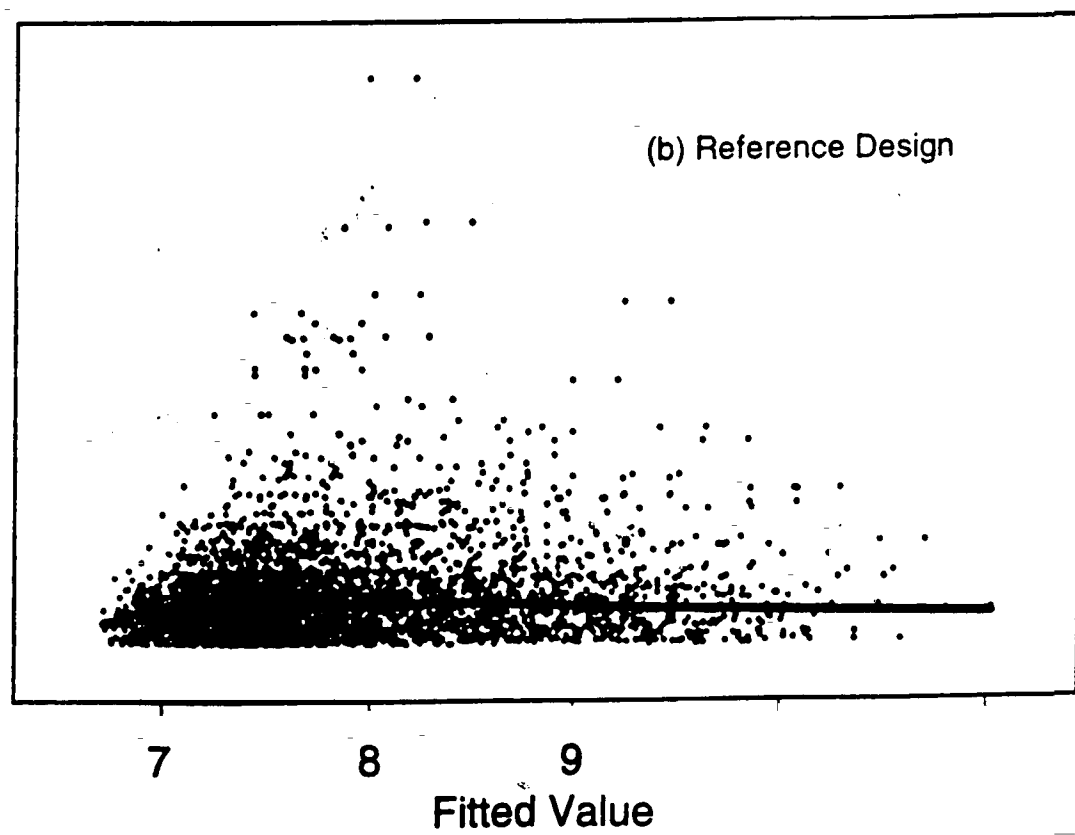
(b) Reference Design



Absolute Residual

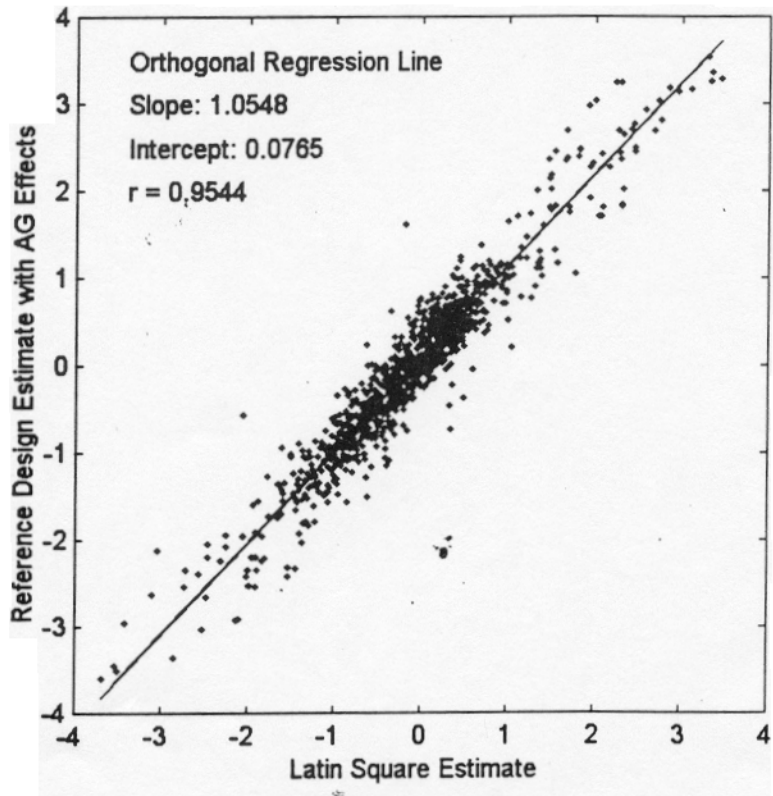
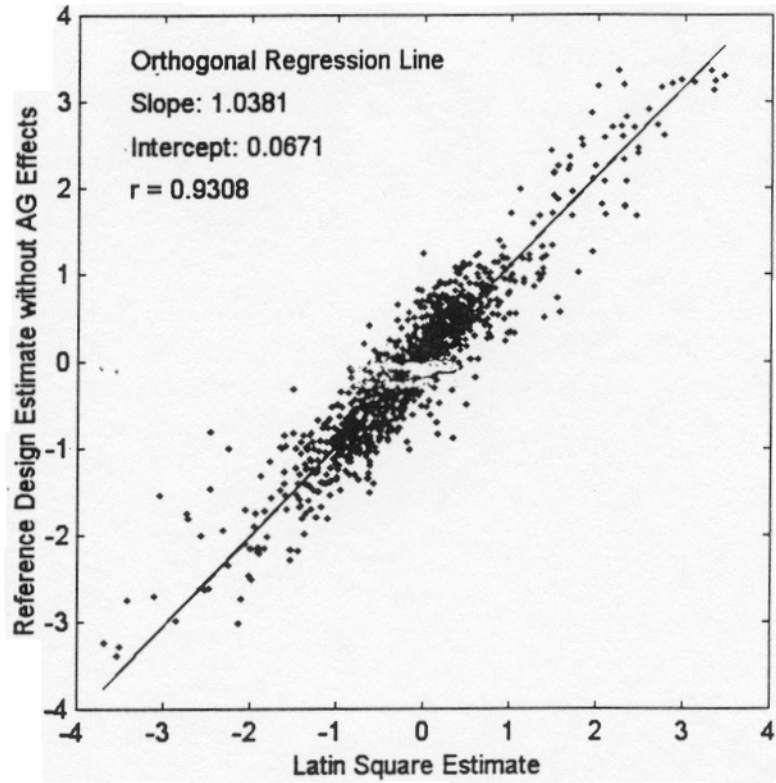


Absolute Residual

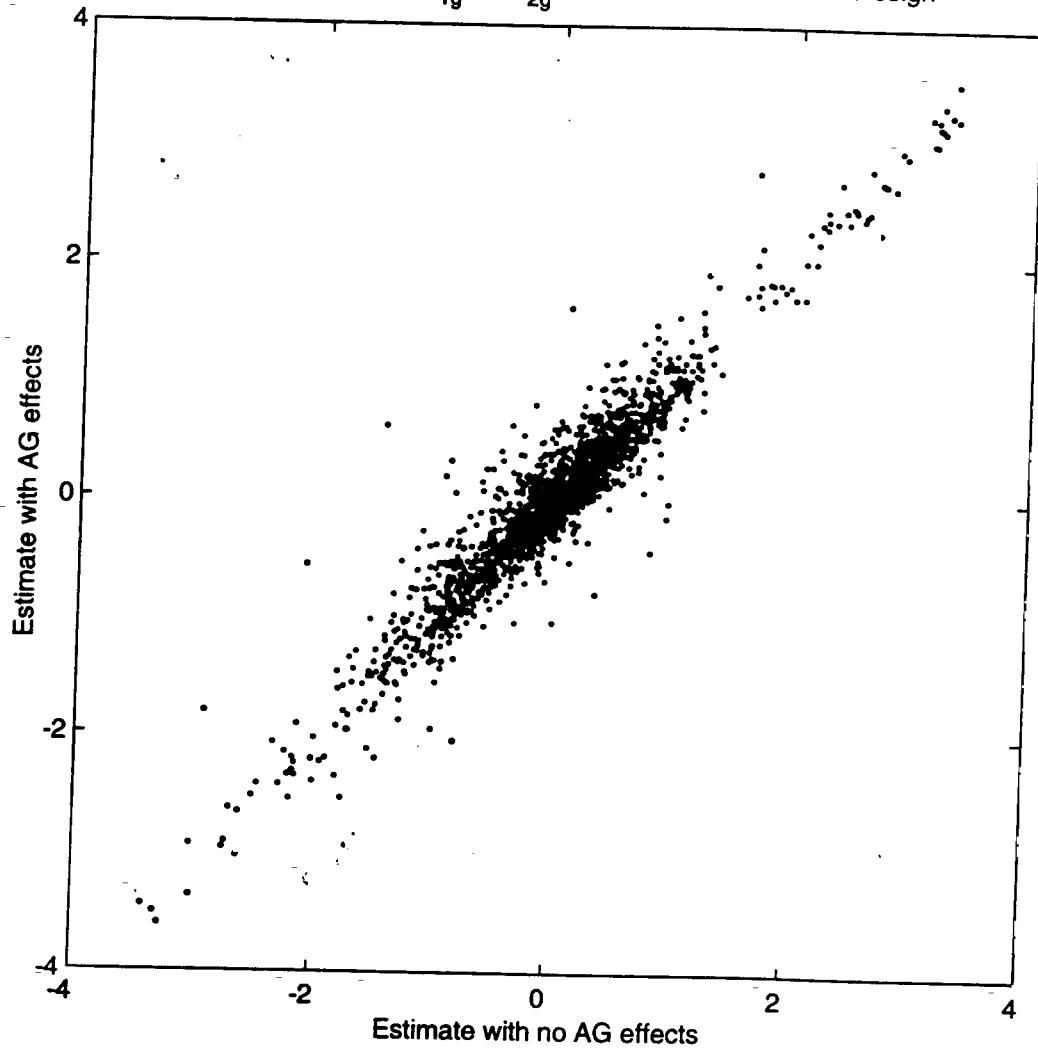


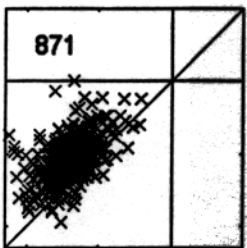
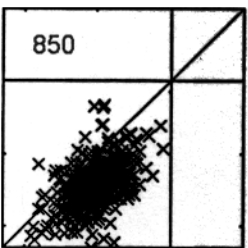
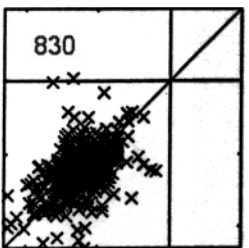
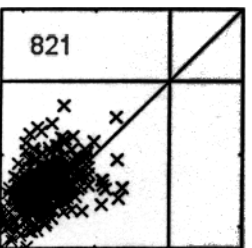
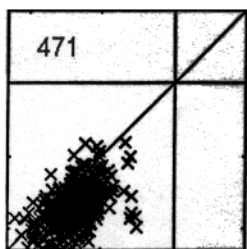
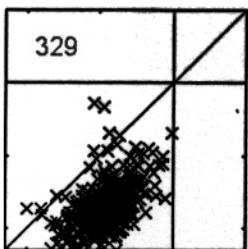
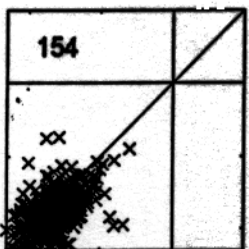
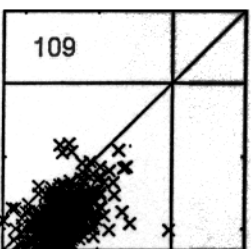
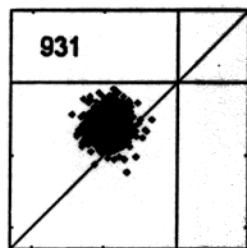
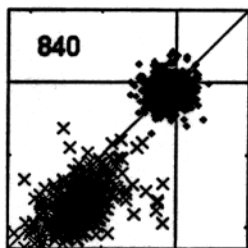
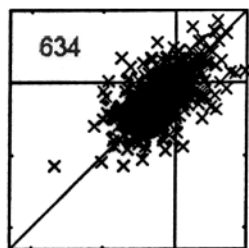
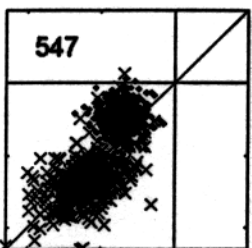
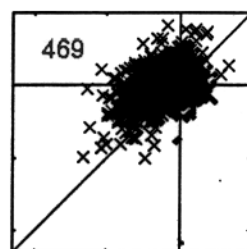
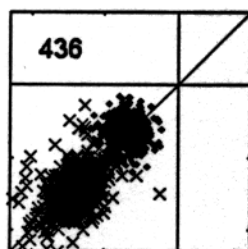
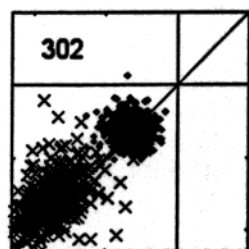
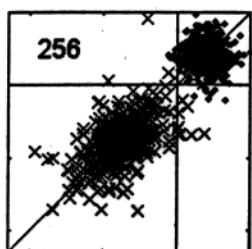
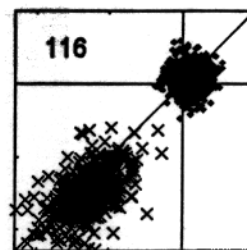
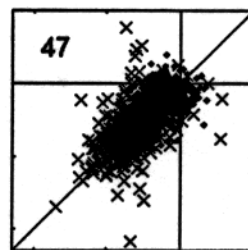
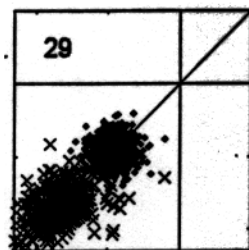
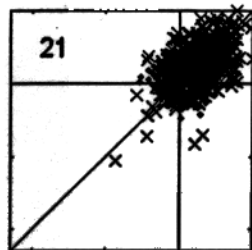
Concordance of the Liver Muscle Differences

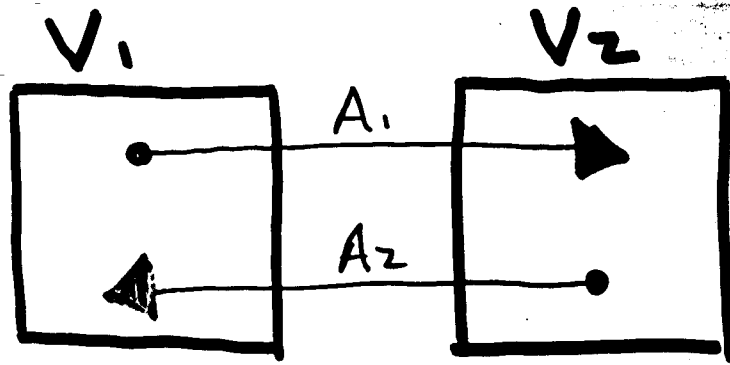
Latin Square Design	Reference Sample Design		
	L < M	N C	L > M
L < M	55	26	0
N C	58	792	46
L > M	0	6	84



Comparison of $(VG)_{1g} - (VG)_{2g}$ Estimates for Reference Design







The latin square is a half-fraction of the theoretical complete factorial design. It has a simple aliasing structure.

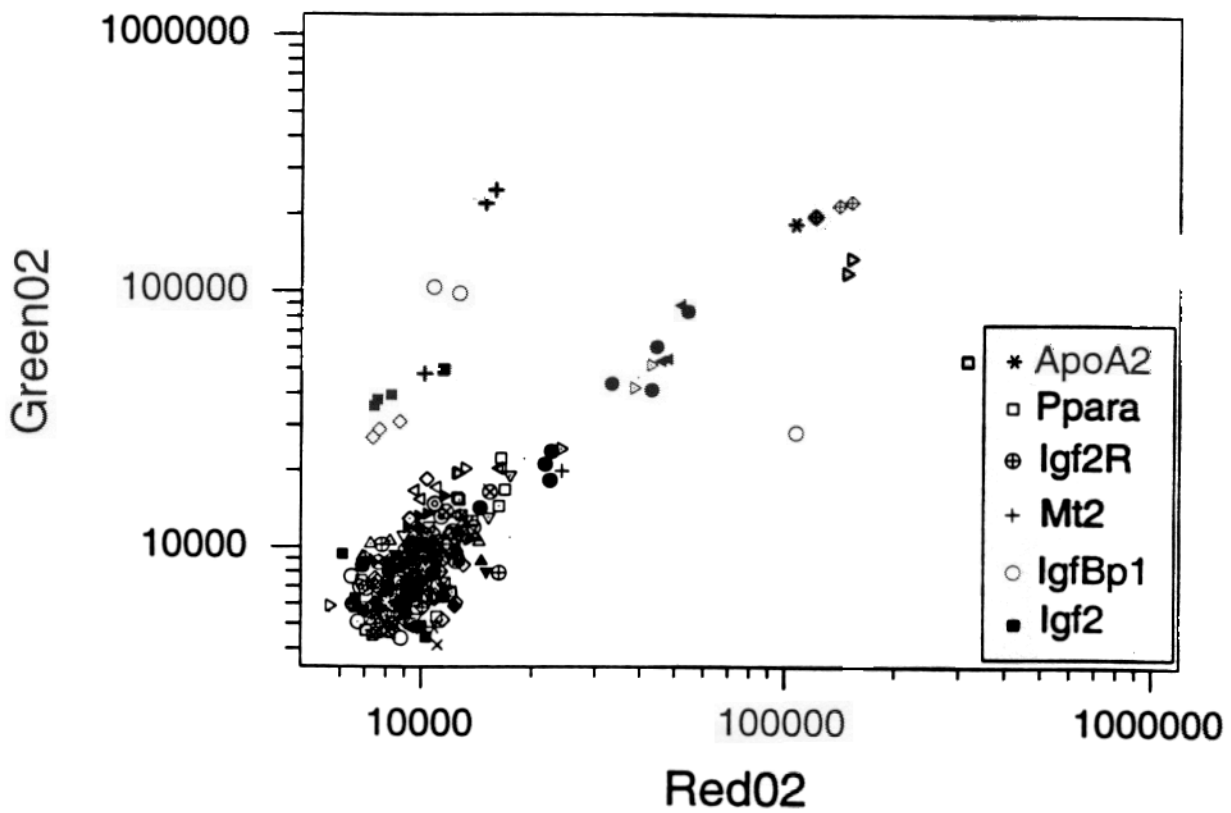
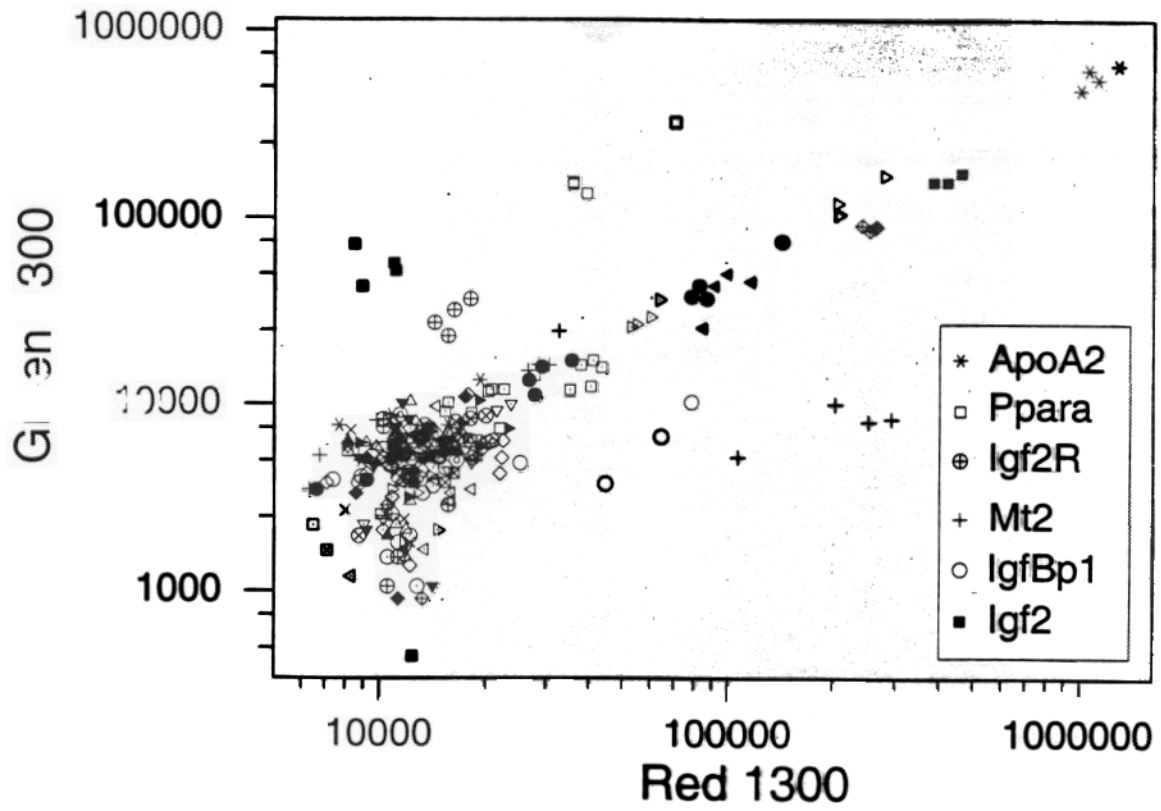
mean	~	ADV
A	~	DV
D	~	AV
V	~	AD
G	~	ADVG
VG	~	ADG
AG	~	DVG
DG	~	AVG

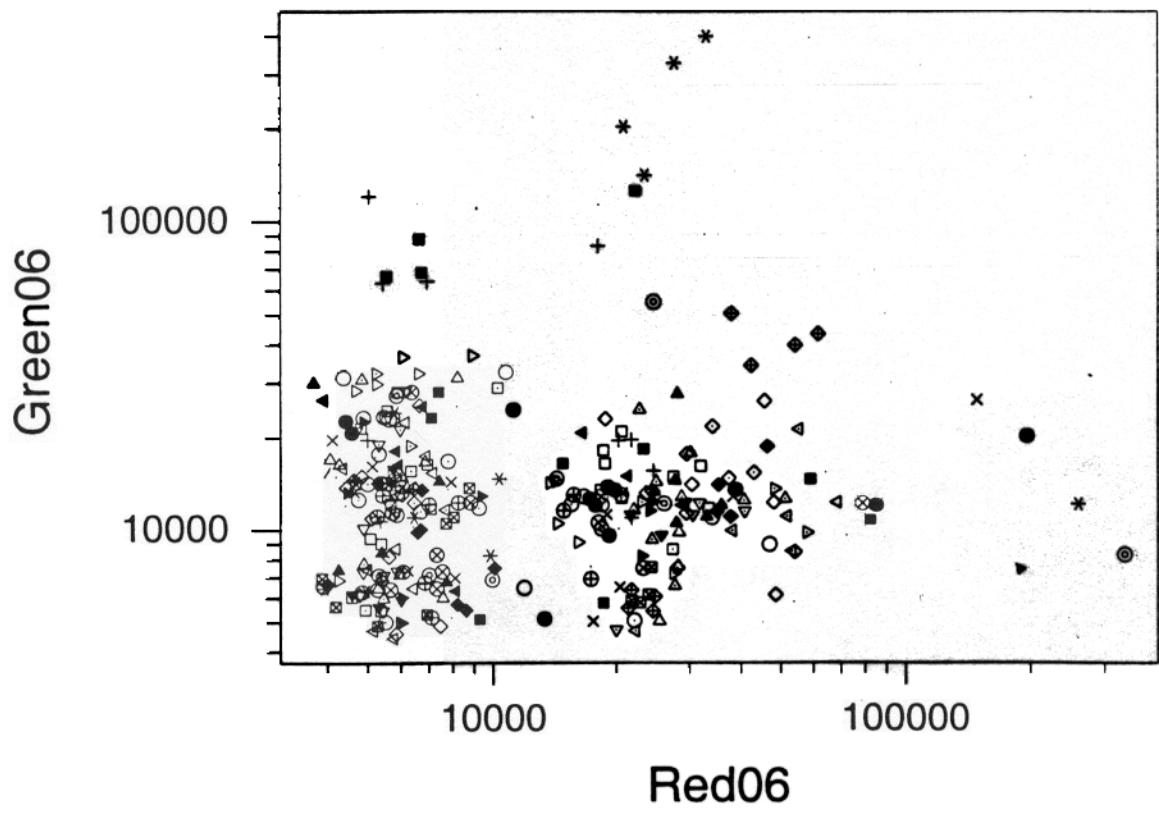
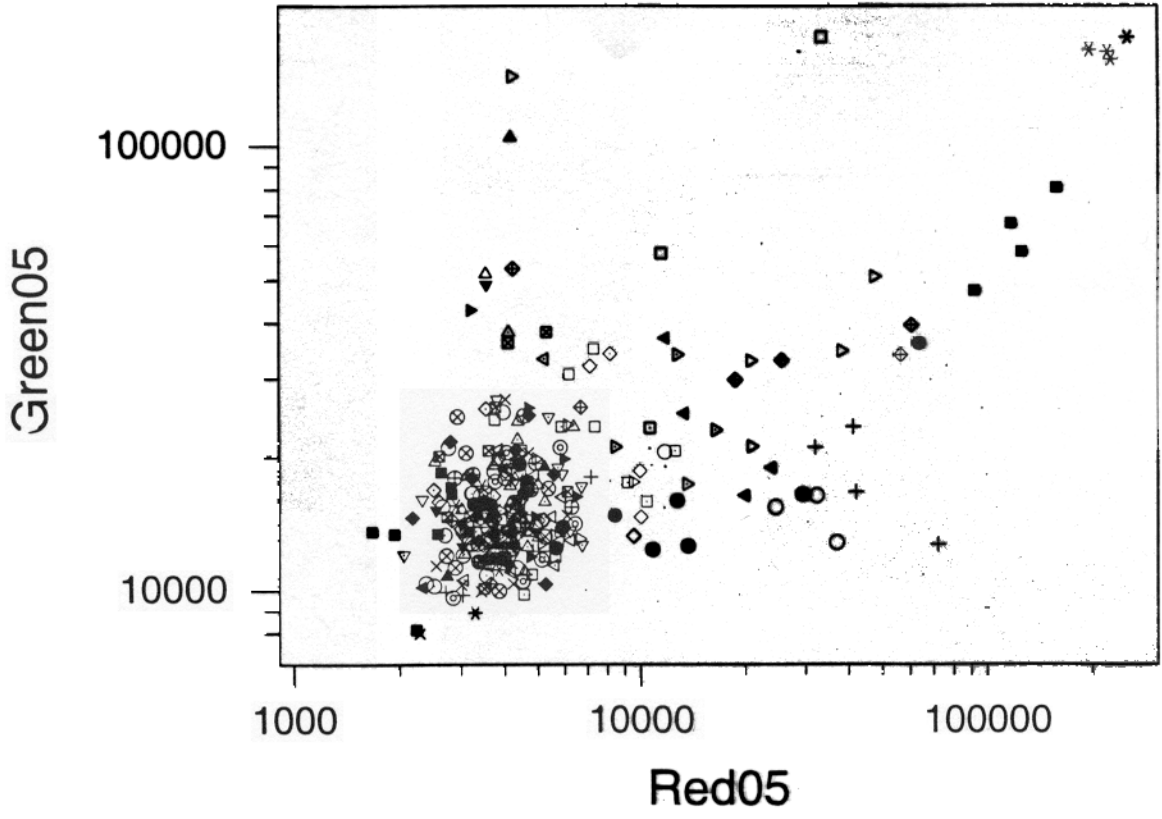
Corning Latin Square Experiment

Source	SS	df	MS
Array	3.11	1	3.11
Dye	79.81	1	79.82
Variety	53.06	1	53.06
Gene	954.60	83	11.30
Variety × Gene	70.84	83	0.85
Spot × Gene	165.36	332	0.50
Dye × Gene	30.79	83	0.37
Residual	30.38	759	0.0400
Adj Total	1387.96	1343	

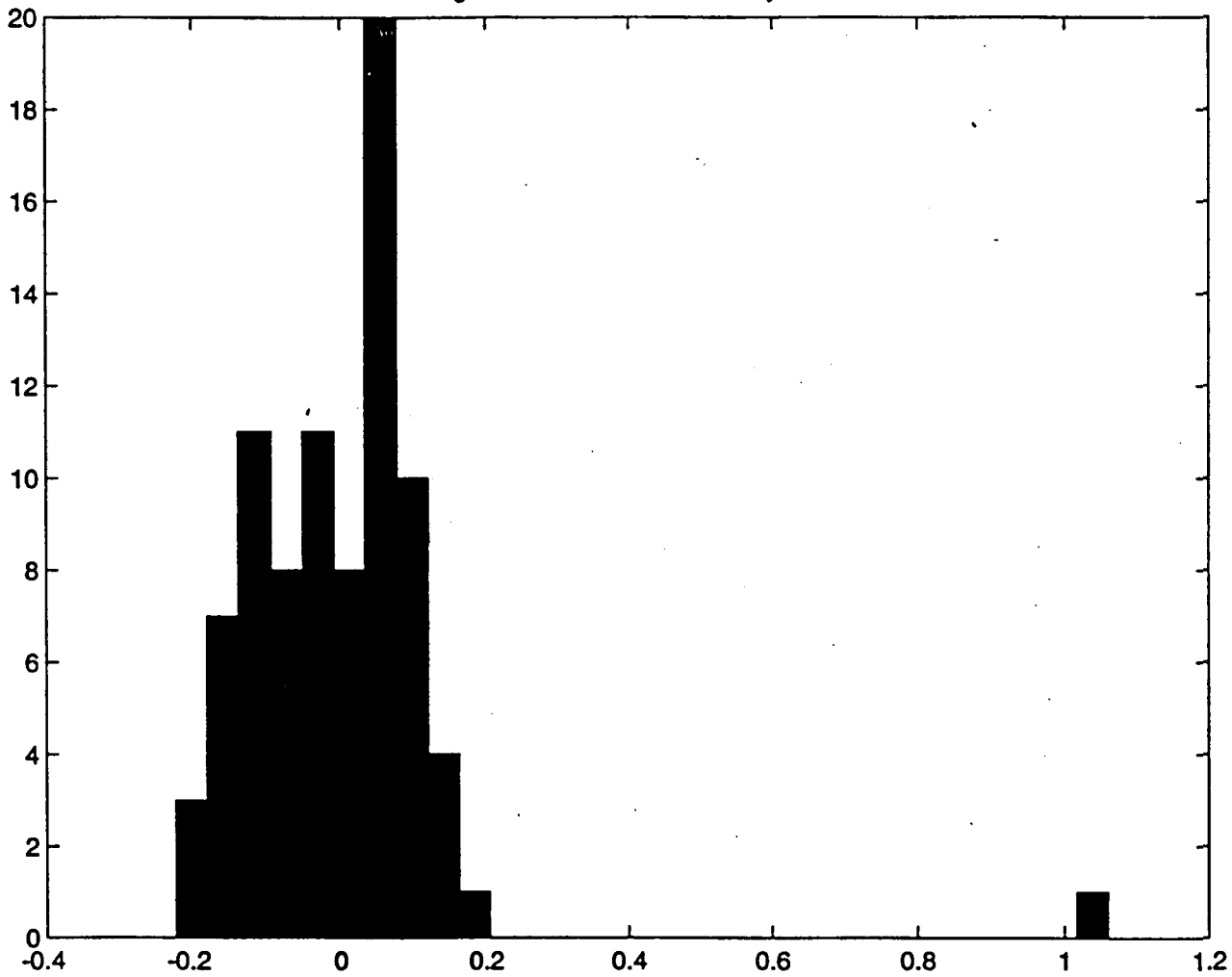
Remove Gene 98/Igf2:

Source	SS	df	MS
Array	3.33	1	3.33
Dye	87.35	1	87.35
Variety	53.99	1	53.99
Gene	949.77	82	11.58
Variety × Gene	69.87	82	0.85
Spot × Gene	164.69	328	0.50
Dye × Gene	12.60	82	0.15
Residual	30.25	750	0.0403
Adj Total	1371.86	1327	

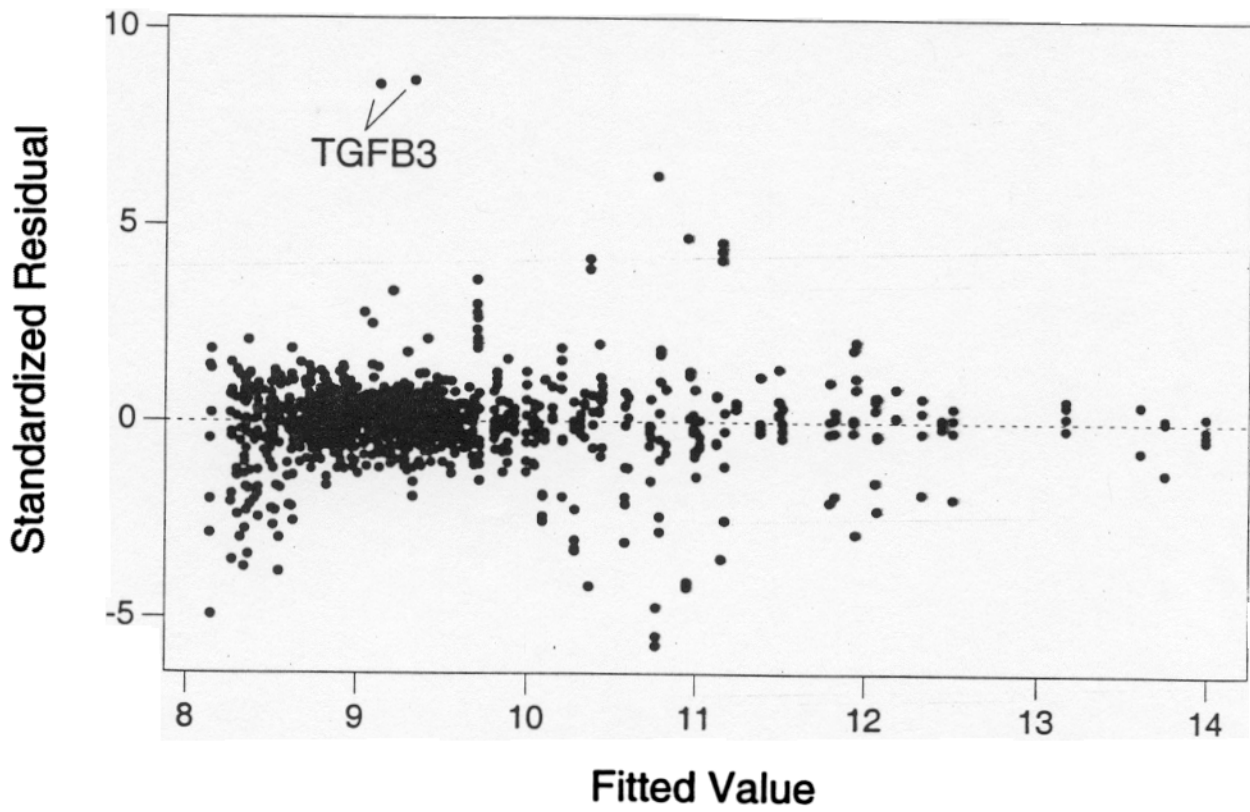




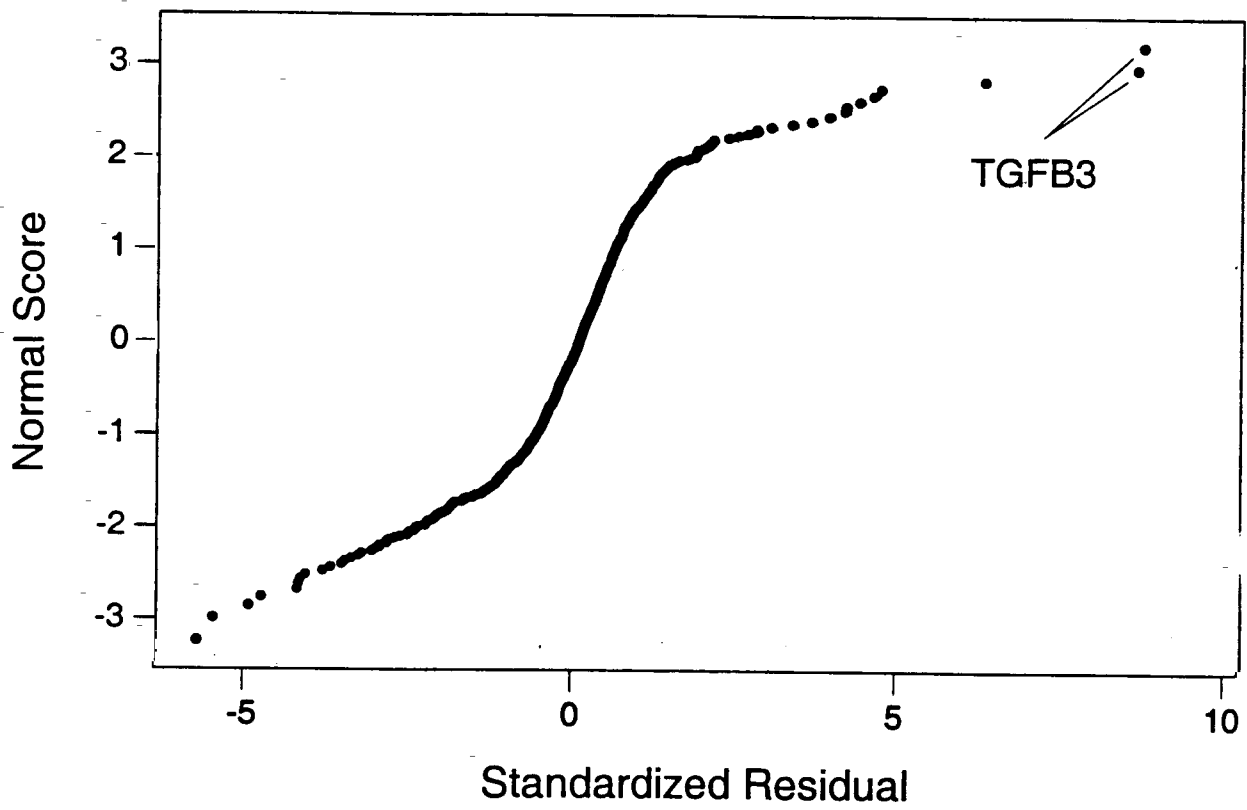
Coming - Distribution of Gene X Dye Effects

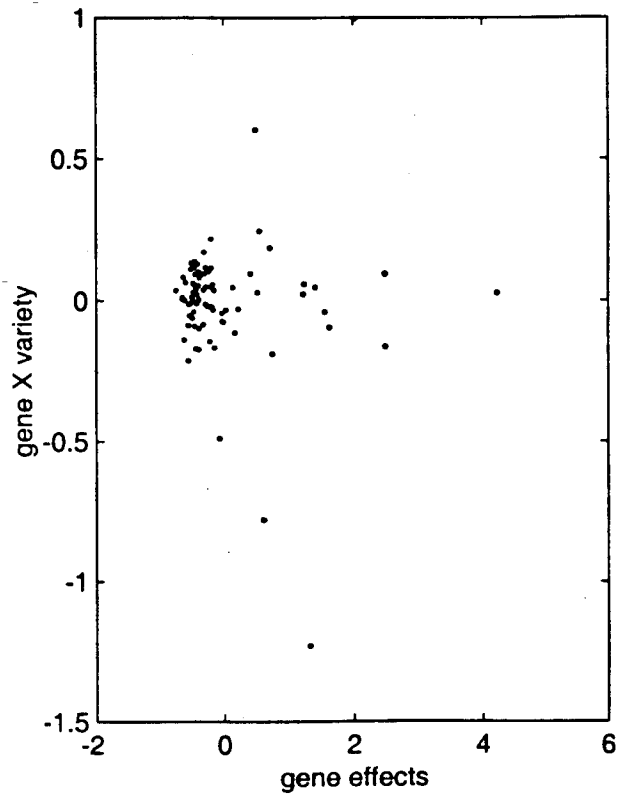
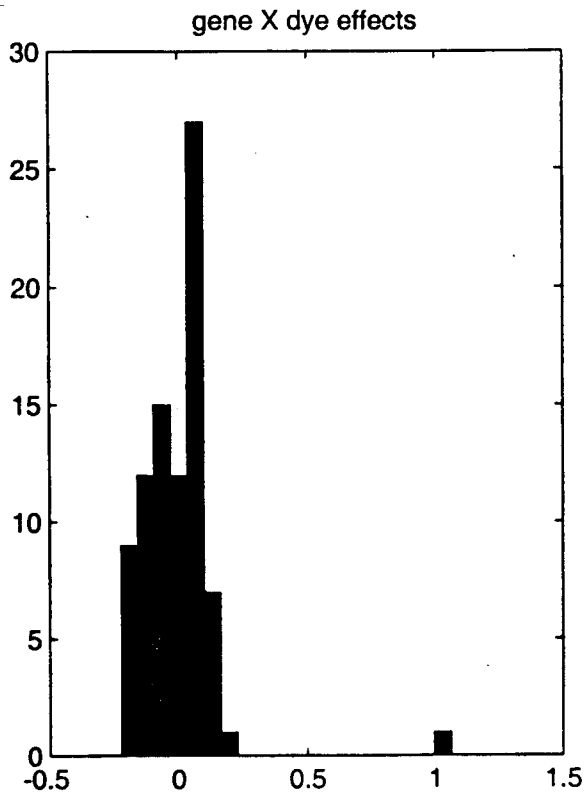
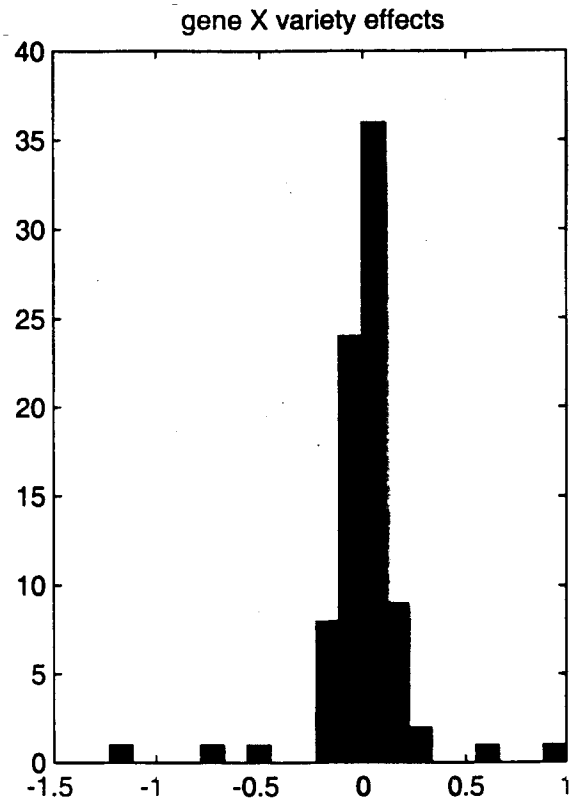
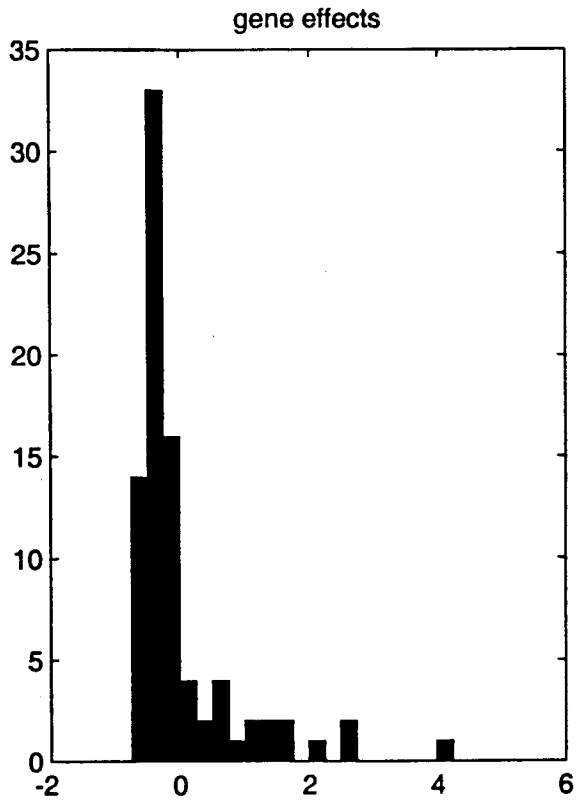


Residuals Versus the Fitted Values



Normal Probability Plot of the Residuals





The Question: Is gene g differentially expressed across varieties k_1 and k_2 ?

$$\text{Is } (VG)_{k_1g} = (VG)_{k_2g}?$$

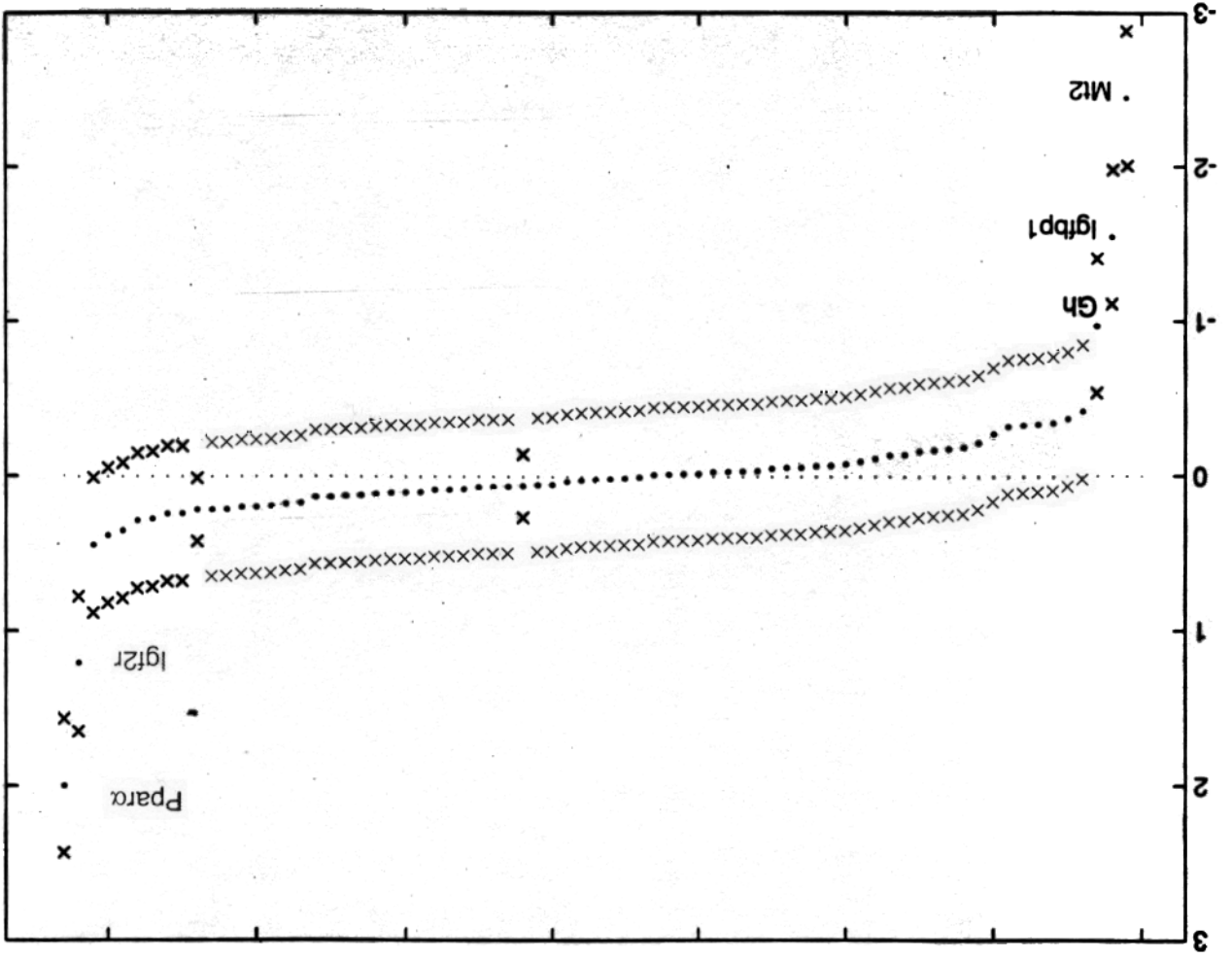
Parametric Bootstrap

Produce bootstrap datasets $\log(y_{ijkgr})^* =$

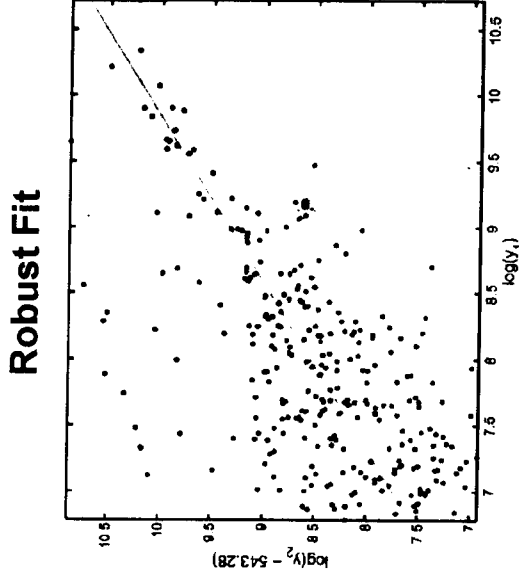
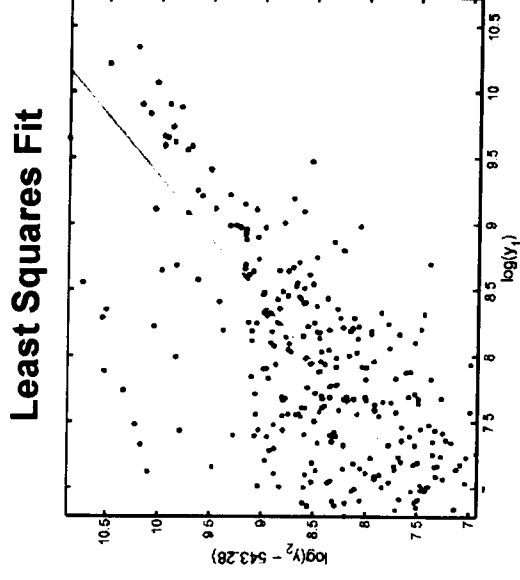
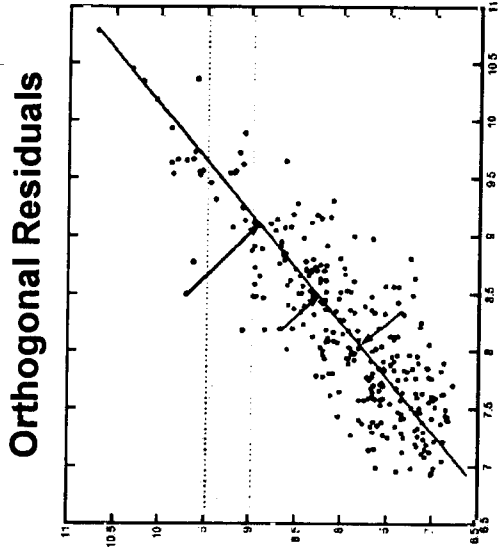
$$\hat{\mu} + \hat{A}_i + \hat{D}_j + \hat{V}_k + \hat{G}_g + (\widehat{VG})_{kg} + \epsilon_{ijkgr}^*$$

The ϵ_{ijkgr}^* are drawn independently from the empirical distribution of residuals from the original fit of the model.

Form 99% confidence intervals using the percentiles of the bootstrap distribution of $(\widehat{VG})_{Tg_0}^*$ $(\widehat{VG})_{Cg_0}^*$ for each gene.



Fitting a Regression Line to the Data

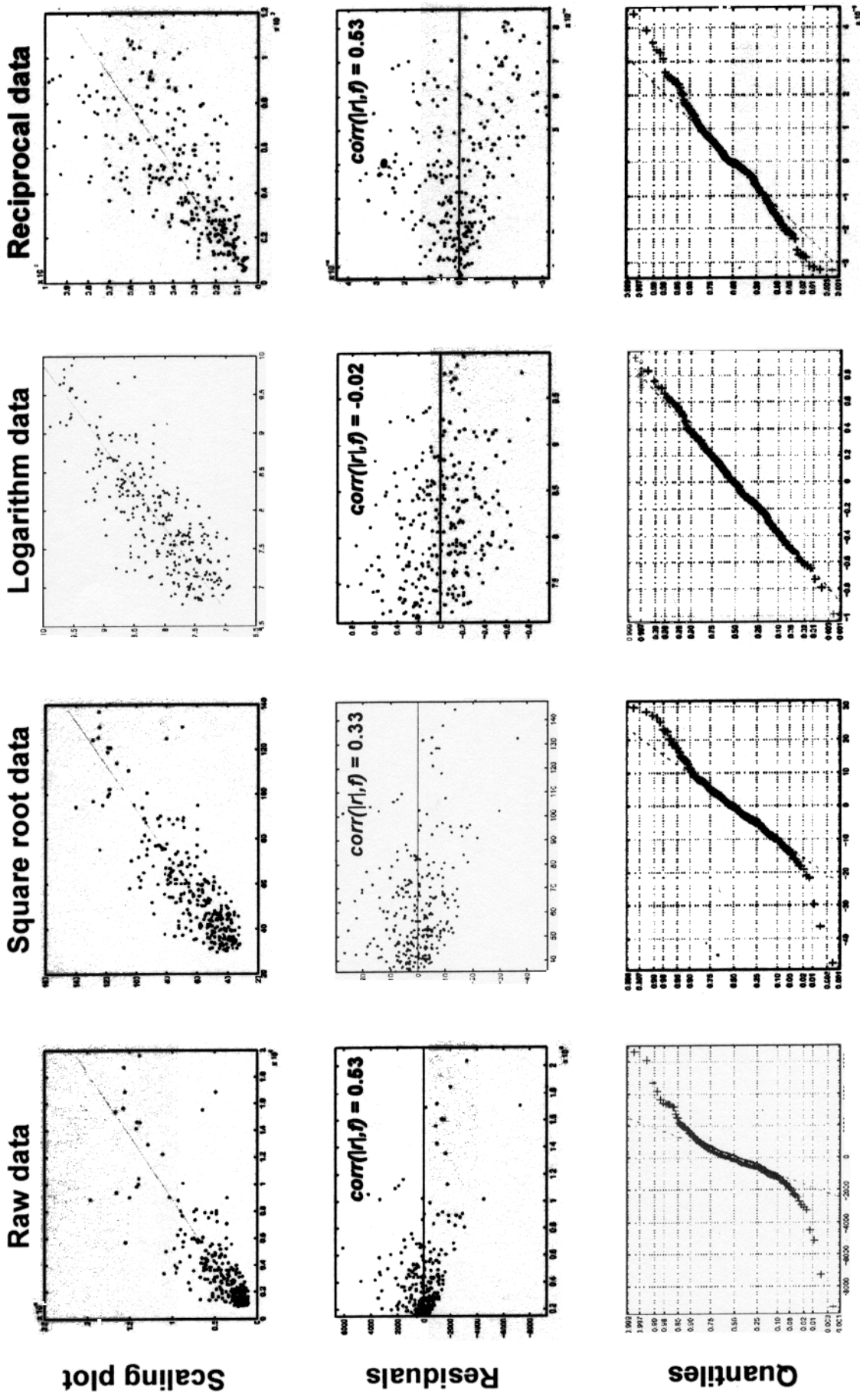


Orthogonal residuals provide a symmetric treatment of the two fluorescent intensities y_1 and y_2 .

Robust regression reduces the influence of differentially expressed genes of the Fitted line.

Selection of the Scaling Function:

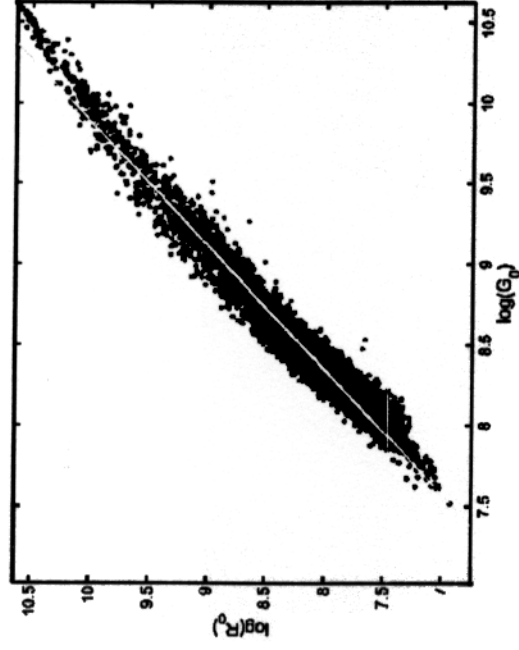
Placenta-Placenta Self Comparison on GEM Array



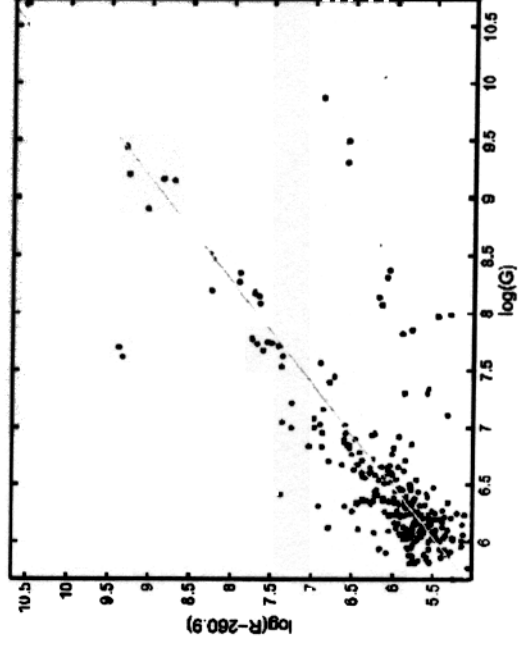
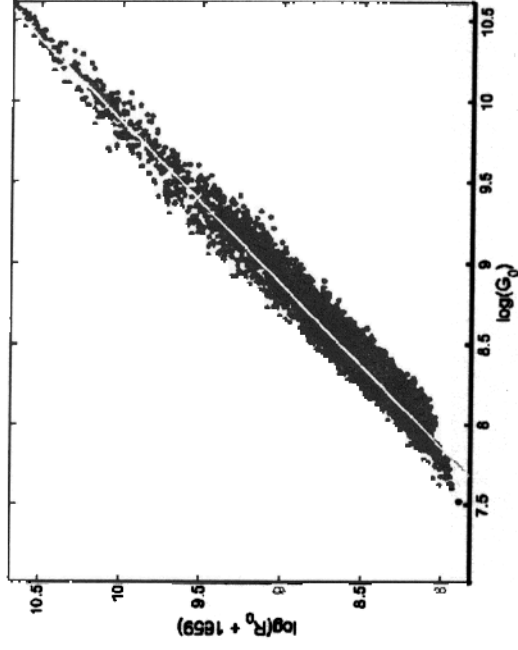
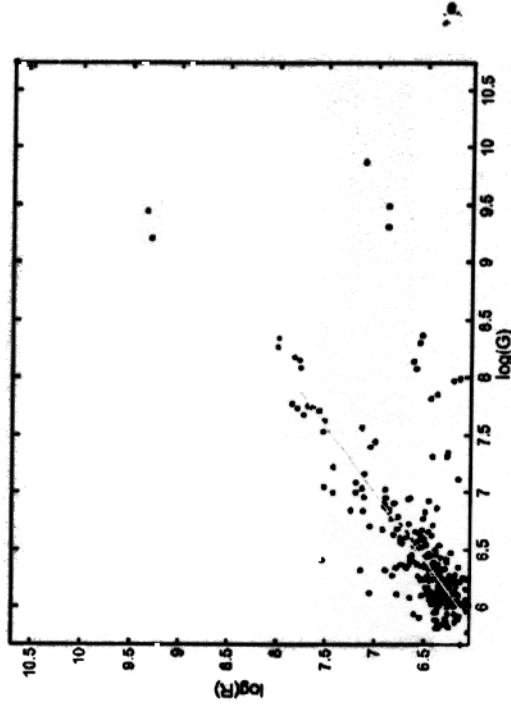
Conclusion: The residuals from a logarithm scaling transformation are approximately normal and have minimal correlation with fitted values.

The Shifted Logarithm Scaling Function

Stanford data

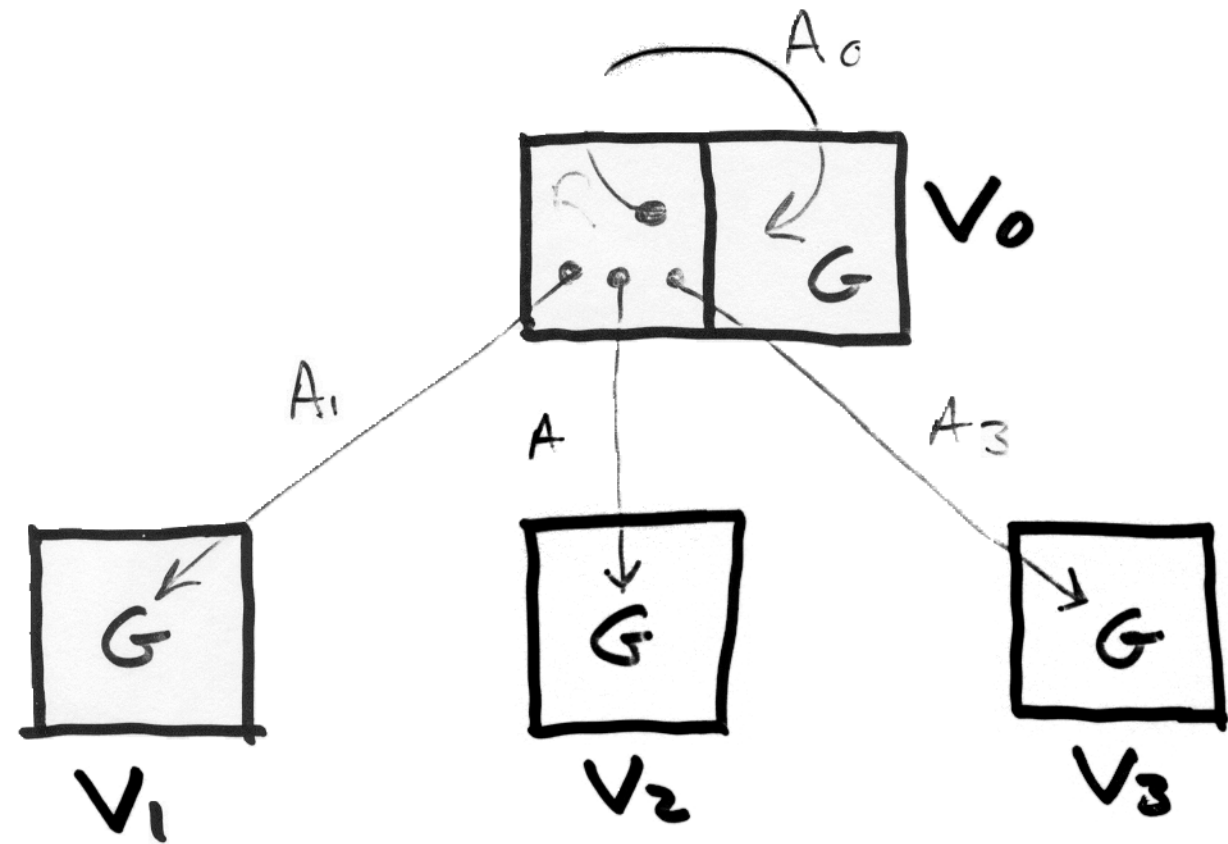


Corning data



Conclusion: Shifting the raw measurement improves linearity on the logarithm scale.

Augmented Reference Design

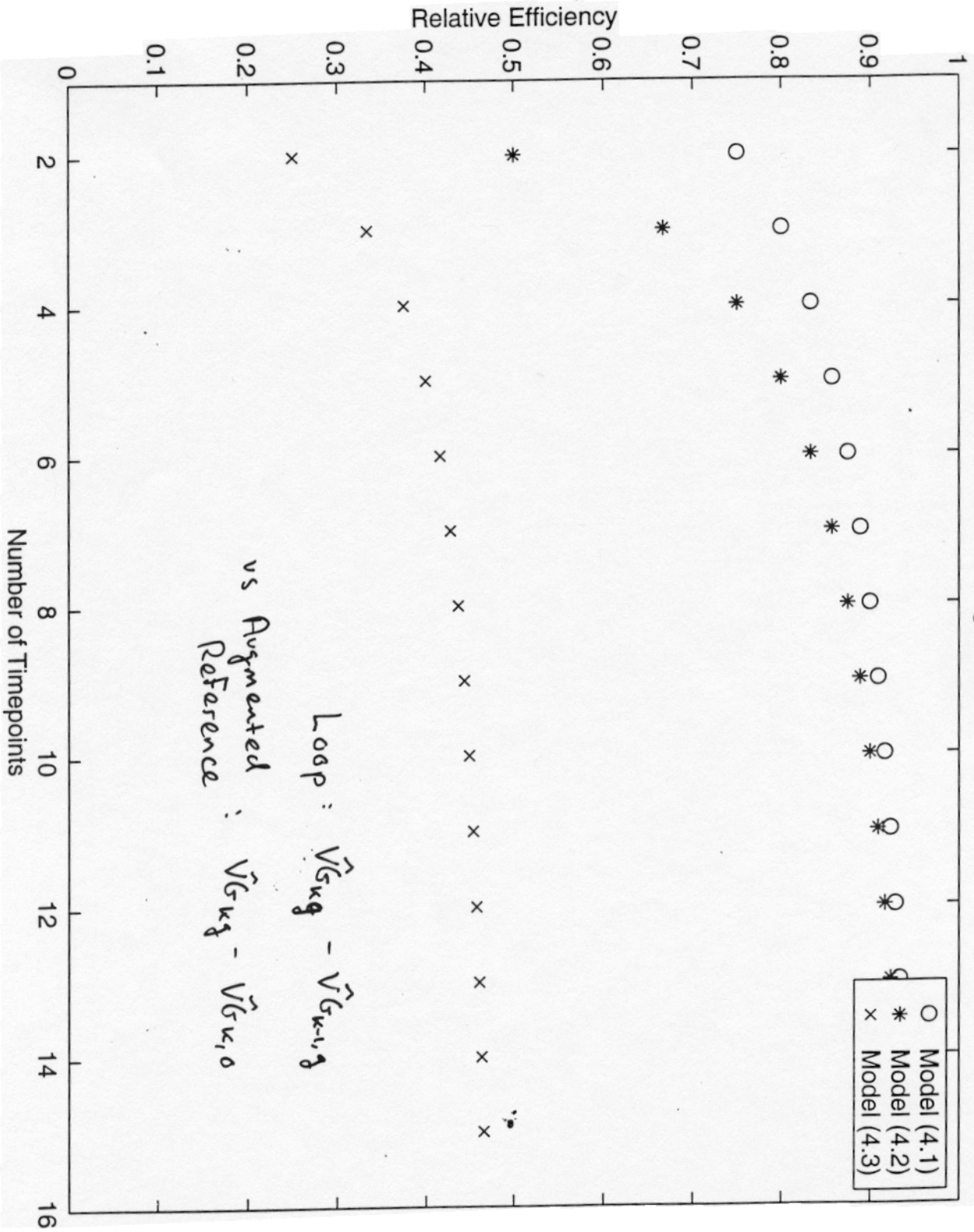


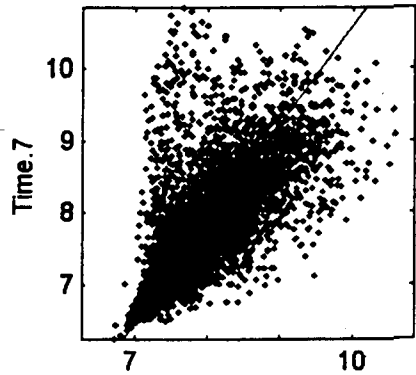
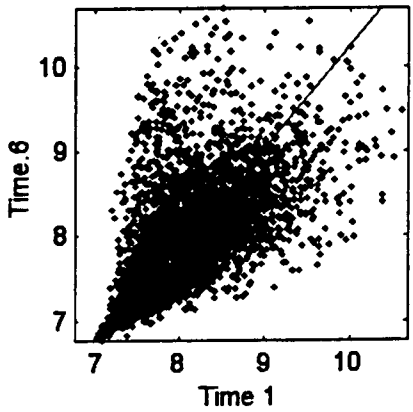
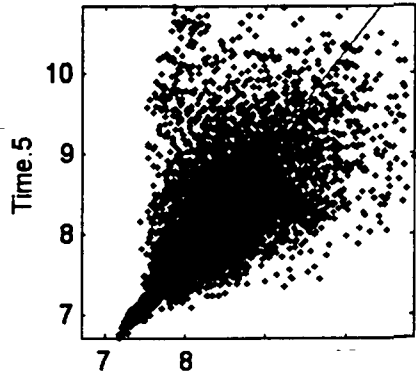
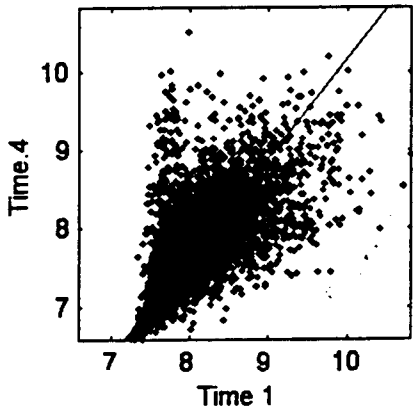
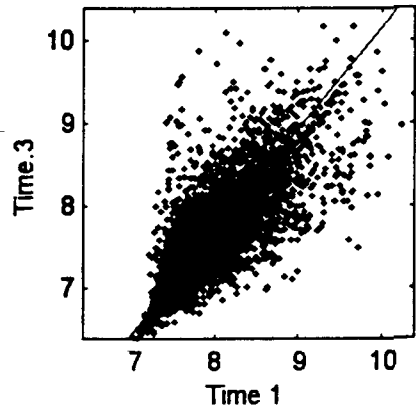
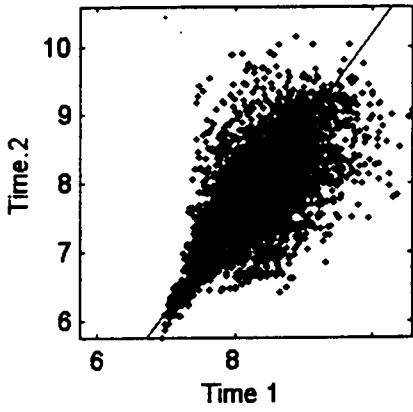
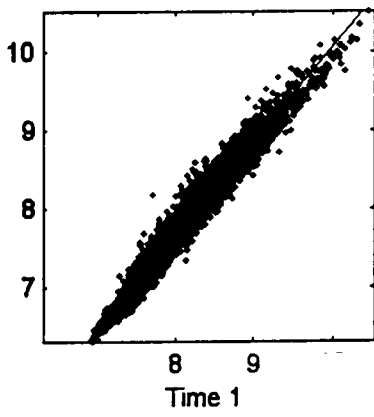
Residual dF

Unbalanced \Rightarrow partial confounding of VG DG

Residuals are all from one array

Figure 4





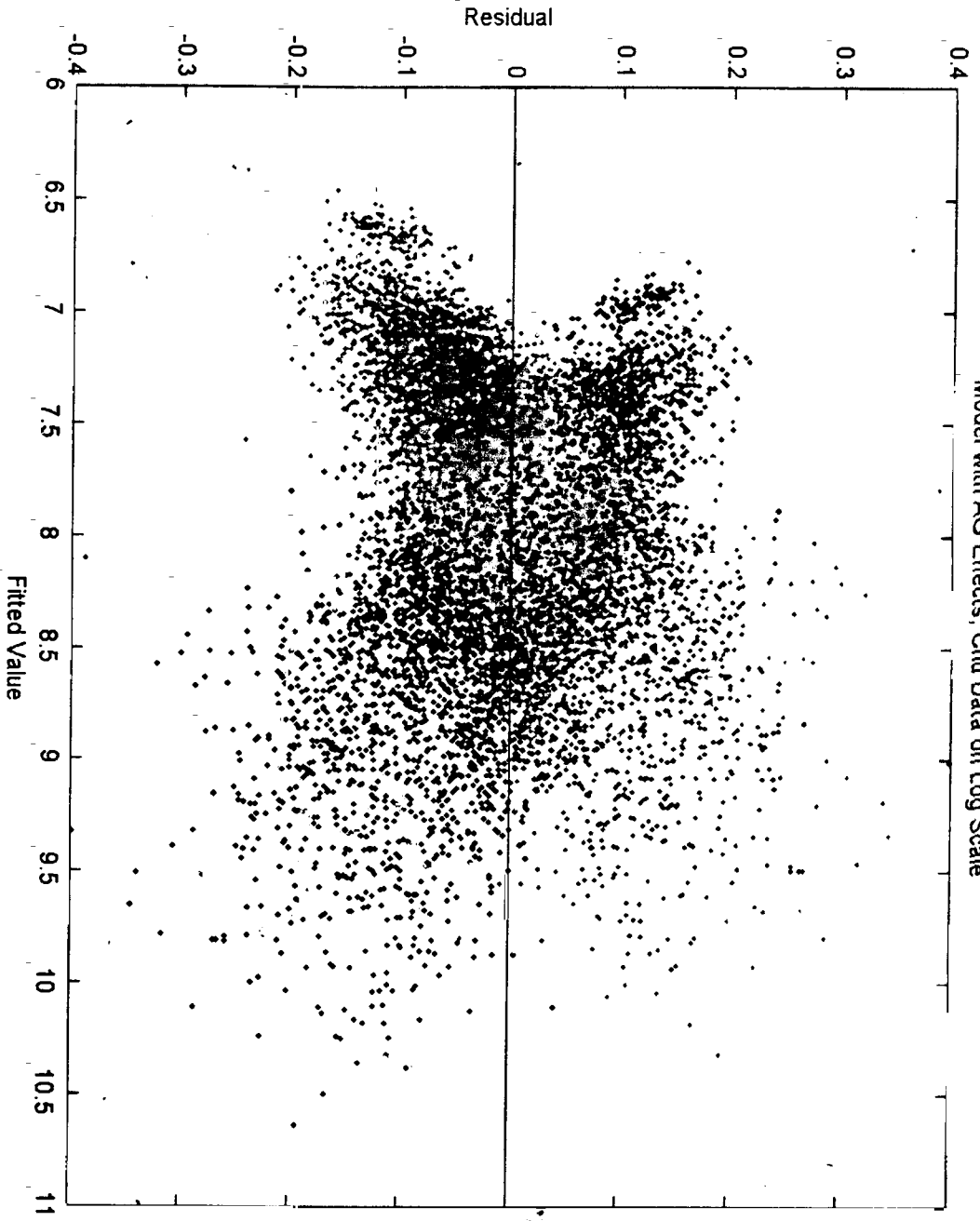
Chu, et al Analyses

Source	SS	df	MS
Array,Dye,Variety	4871.16	13	374.70
Gene	18992.28	6117	3.10
→ VG	8250.36	36702	0.22
Residual	<u>2743.69</u>	42819	0.0641
Adj Total	34857.50	85651	

Source	SS	df	MS
Array,Dye,Variety	4871.16	13	374.70
Gene	18992.28	6117	3.10
→ VG,AG	10909.11	73404	0.15
Residual	<u>84.94</u>	6117	0.0139
Adj Total	34857.50	85651	

Source	SS	df	MS
Array,Dye,Variety	4871.16	13	374.70
Gene	18992.28	6117	3.10
→ VG,DG	8596.02	48936	0.18
Residual	<u>2398.03</u>	36702	0.0653
Adj Total	34857.50	85651	

Model with AG Effects, Chu Data on Log Scale



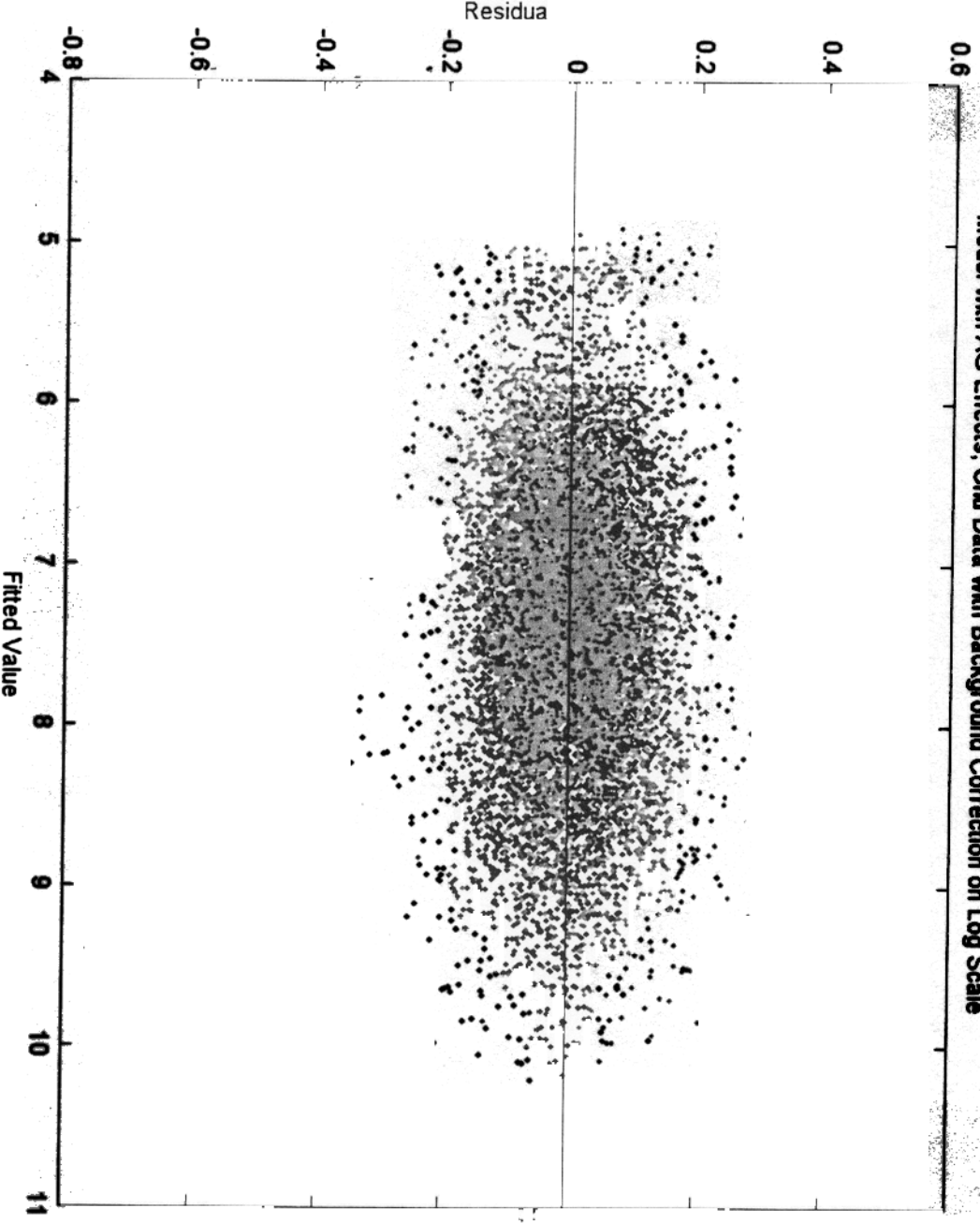
Chu Data, Background Corrected

Source	SS	df	MS
Array,Dye,Variety	6896.24	13	530.48
Gene	48329.71	6117	7.90
VG	16681.88	36702	0.45
Residual	<u>6314.46</u>	42819	0.1475
Adj Total	78222.28	85651	

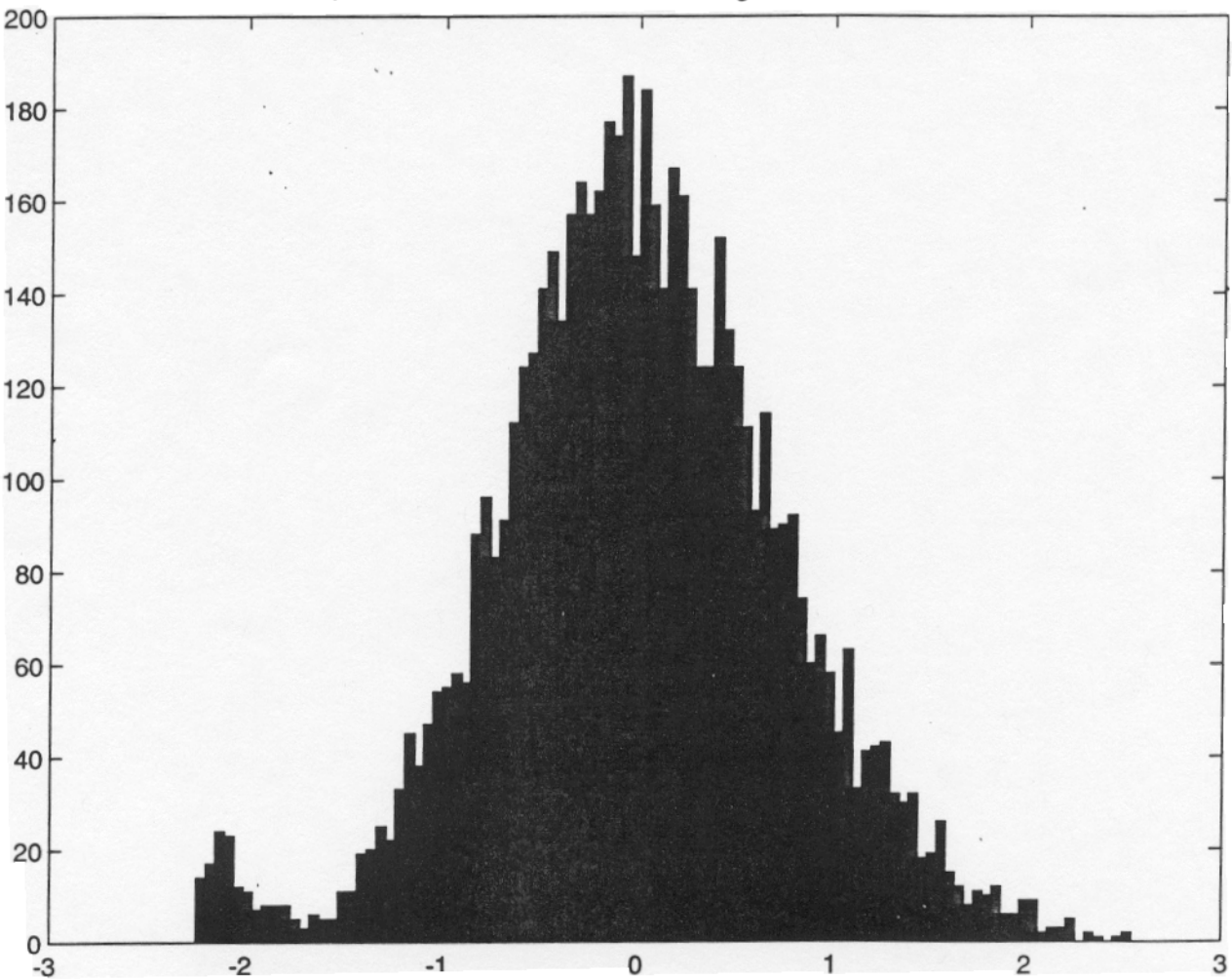
Source	SS	df	MS
Array,Dye,Variety	6896.24	13	530.48
Gene	48329.71	6117	7.90
<u>VG,AG</u>	22907.16	73404	0.31
Residual	<u>89.18</u>	6117	0.0146
Adj Total	78222.28	85651	

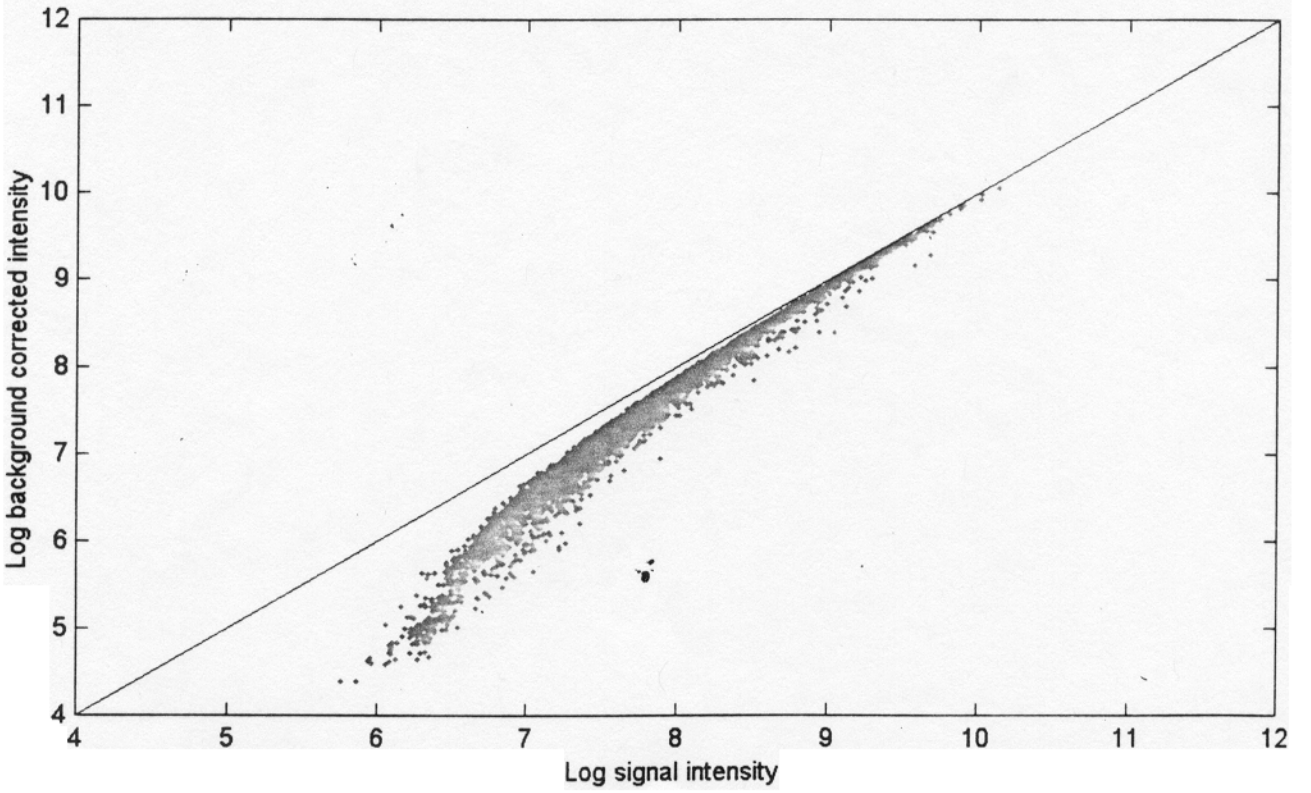
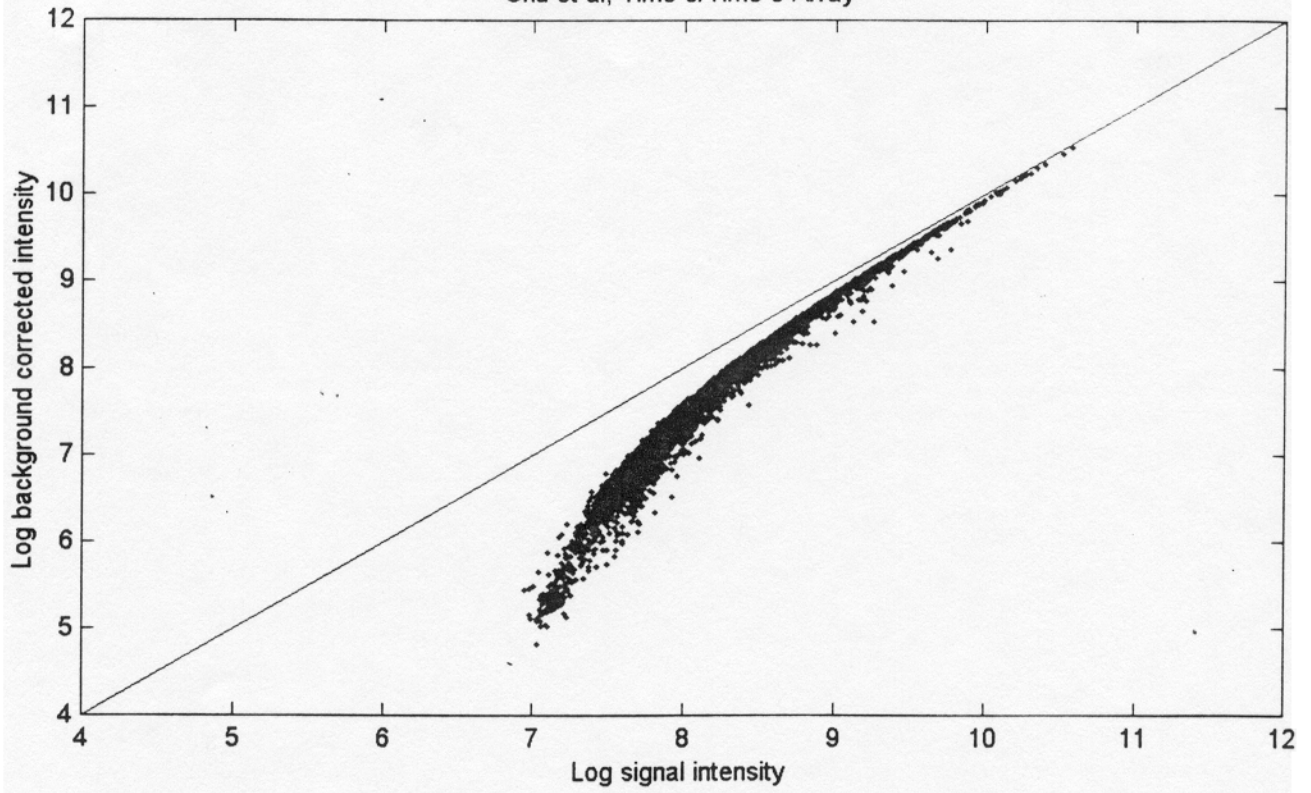
Source	SS	df	MS
Array,Dye,Variety	6896.24	13	530.48
Gene	48329.71	6117	7.90
VG,DG	17215.01	48936	0.35
Residual	<u>5781.33</u>	36702	0.1575
Adj Total	78222.28	85651	

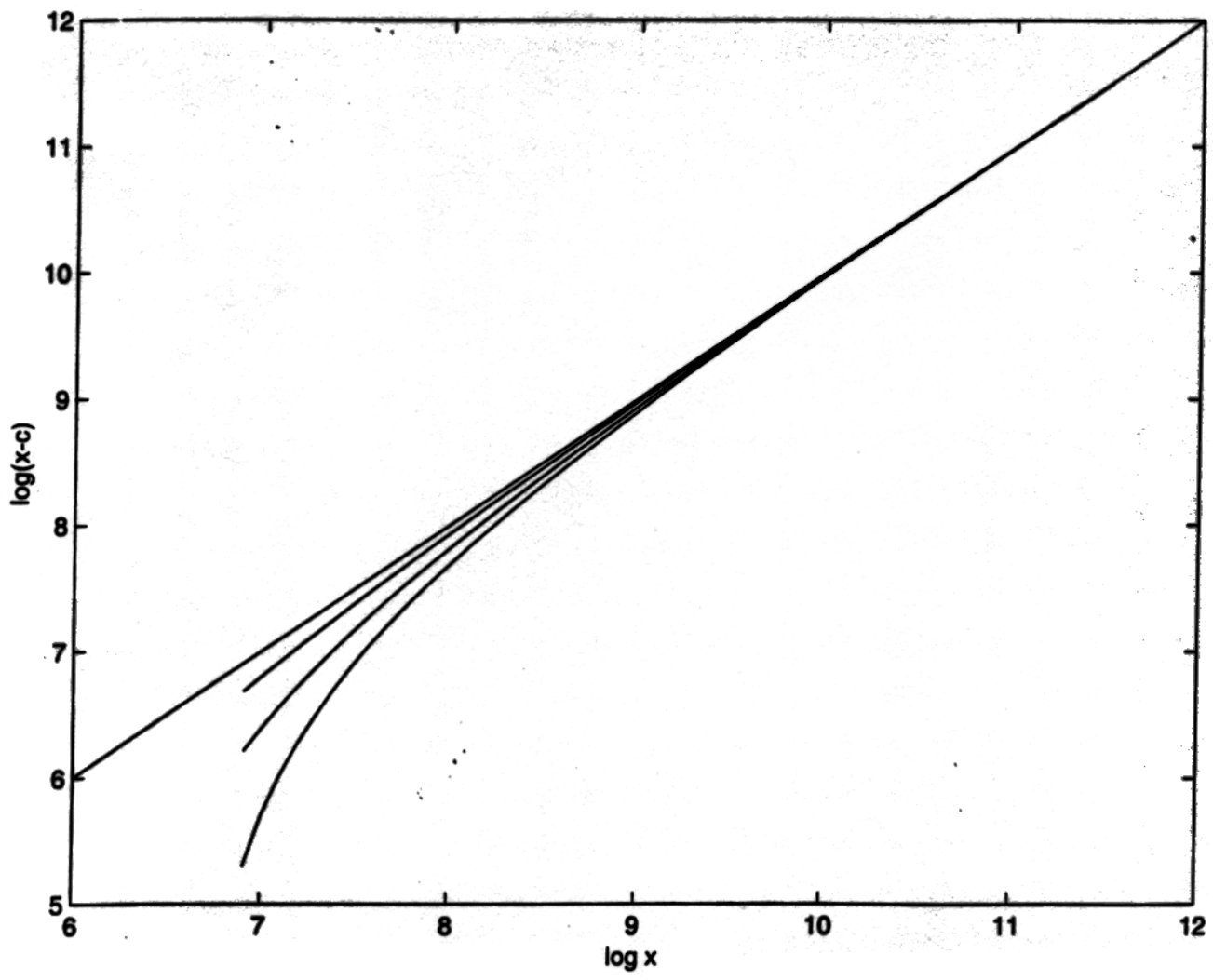
Model with AG Effects, Chu Data with Background Correction on Log Scale



Chu, Distribution of Gene Effects, Background Corrected Data







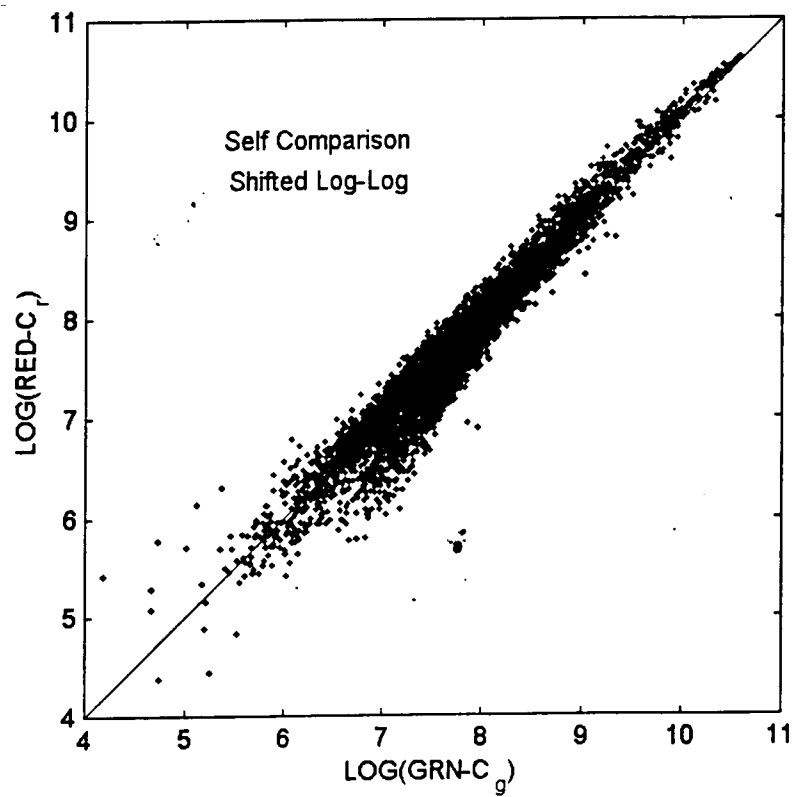
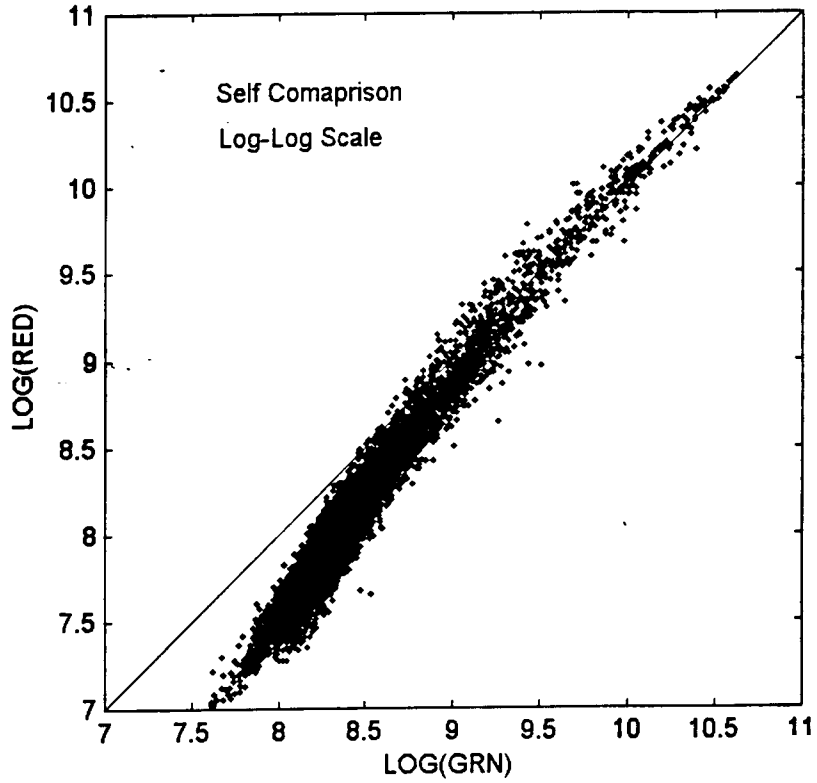
4.

Chu Data with Median Background Correction

Source	SS	df	MS
Array, Dye, Variety	7353.57	13	565.66
Gene	43768.29	5827	7.51
VG	20081.60	34962	0.57
Residual	8855.74	40789	0.2171
Adj Total	80059.21	81591	

Source	SS	df	MS
Array, Dye, Variety	7353.57	13	565.66
Gene	43768.29	5827	7.51
VG, AG	28804.40	69924	0.41
Residual	132.95	5827	0.0228
Adj Total	80059.21	81591	

Source	SS	df	MS
Array, Dye, Variety	7353.57	13	565.66
Gene	43768.29	5827	7.51
VG, DG	20688.91	46616	0.44
Residual	8248.43	34962	0.2359
Adj Total	80059.21	81591	



Model with AG Effects

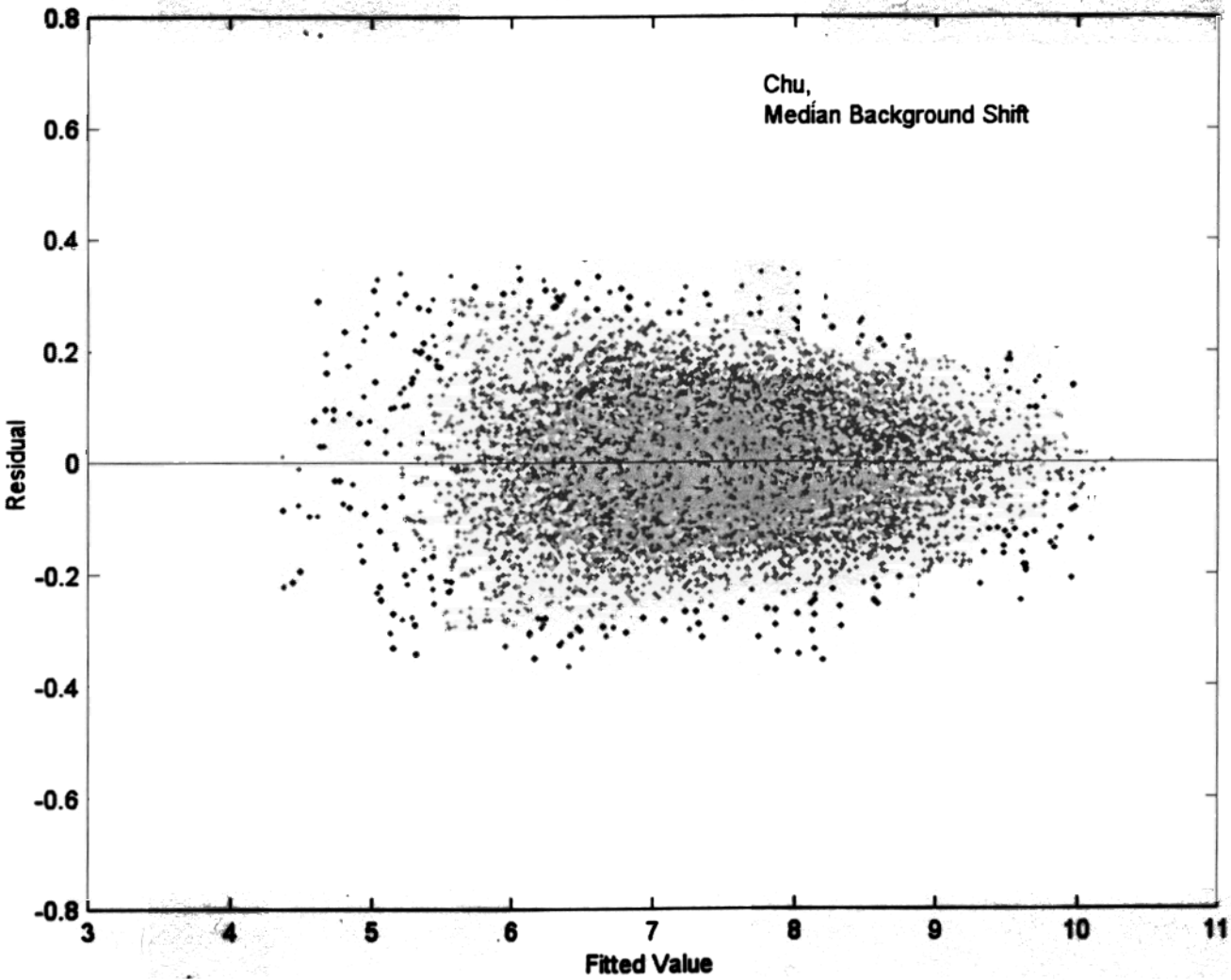
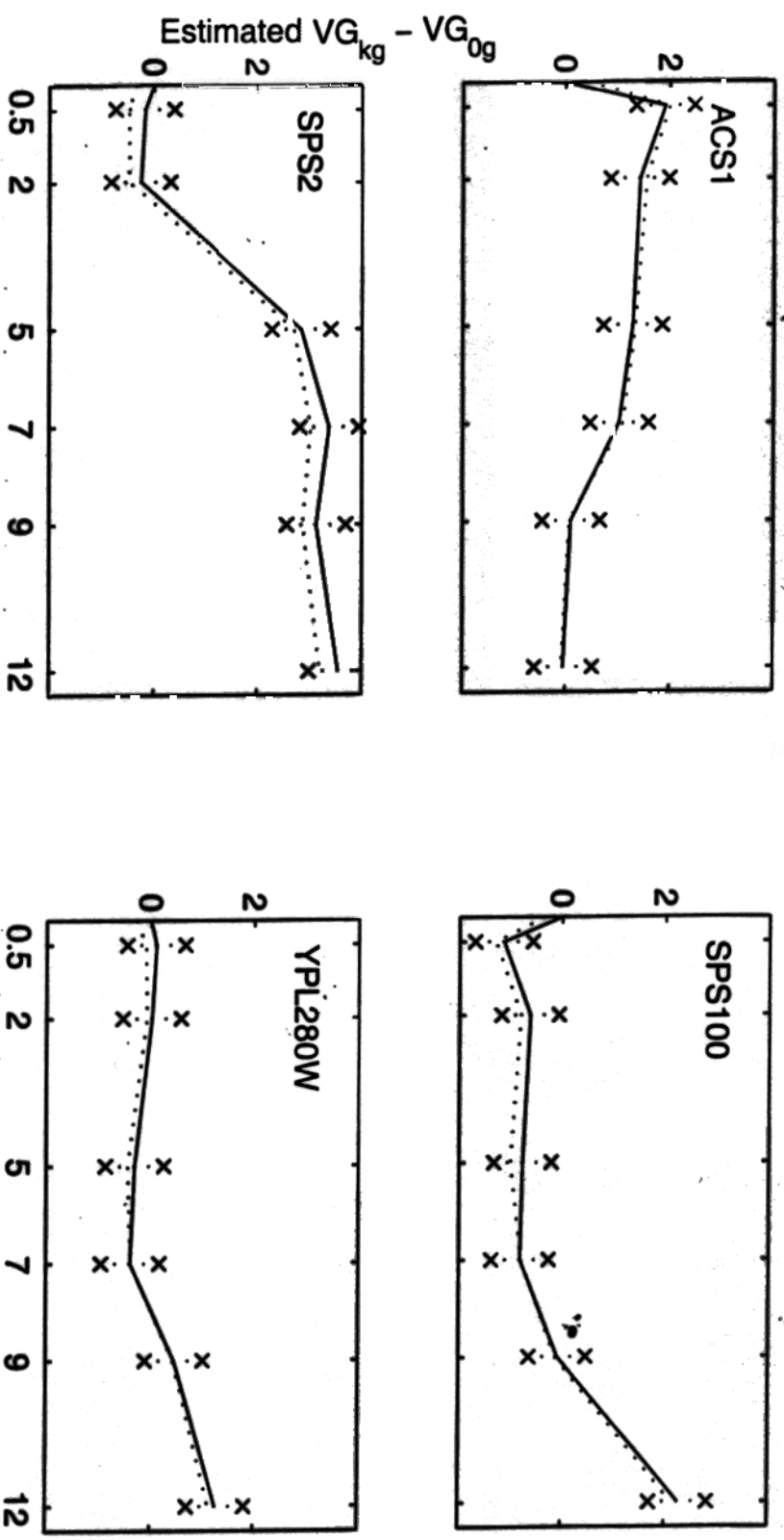


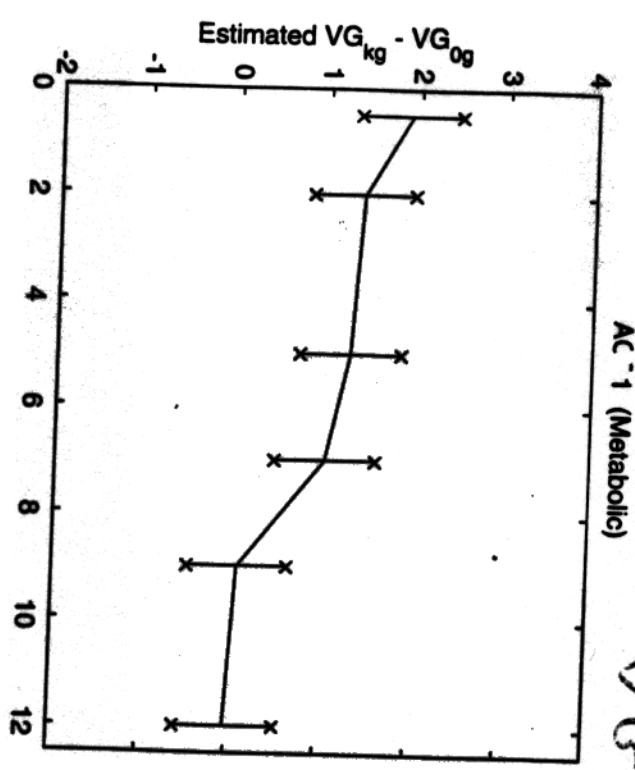
Figure 1



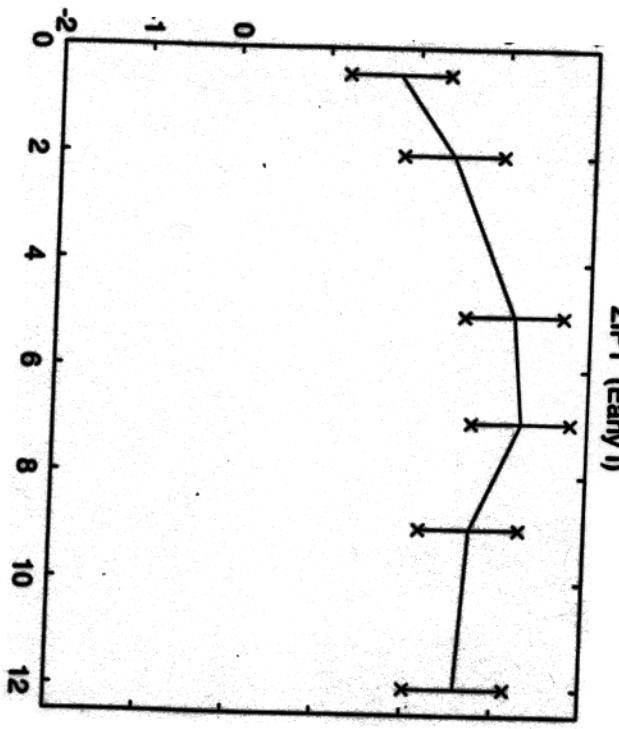
AC - 1 (Metabolic)

$\dot{V}G_{kg}$

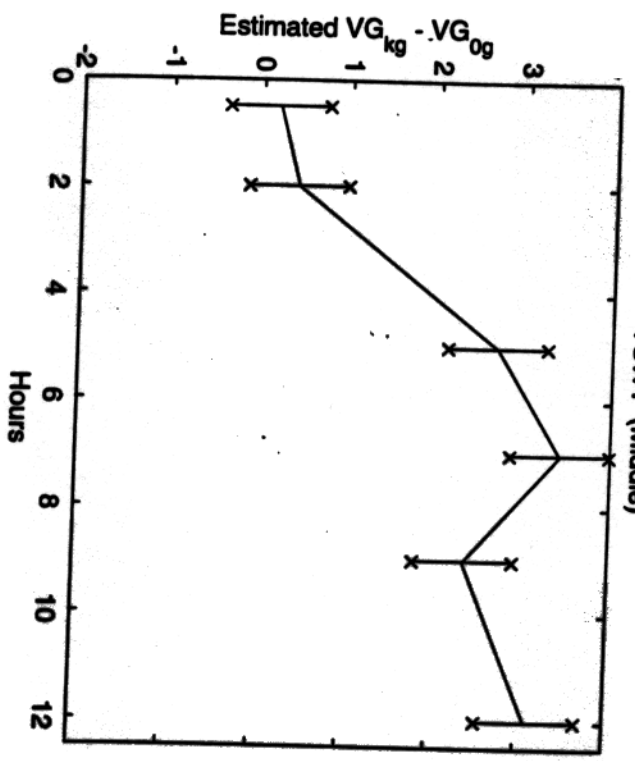
$\dot{V}G_{0g}$



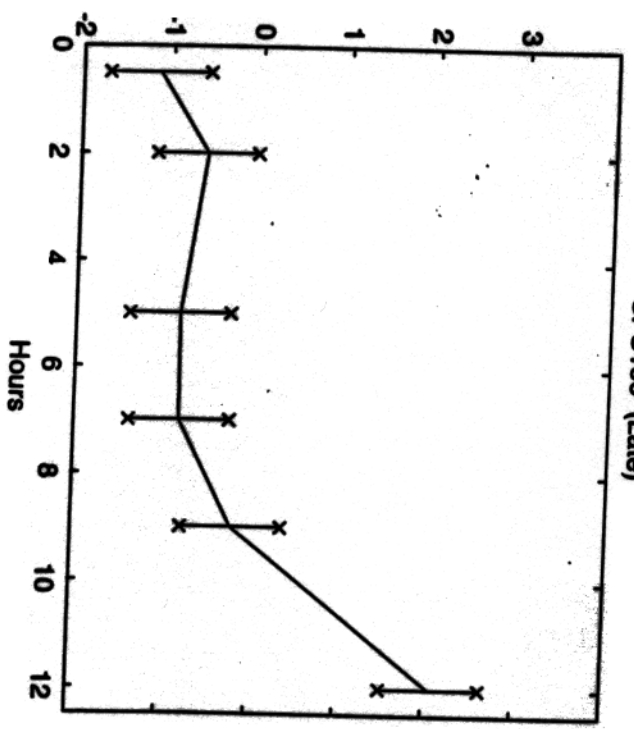
ZIP1 (Early I)



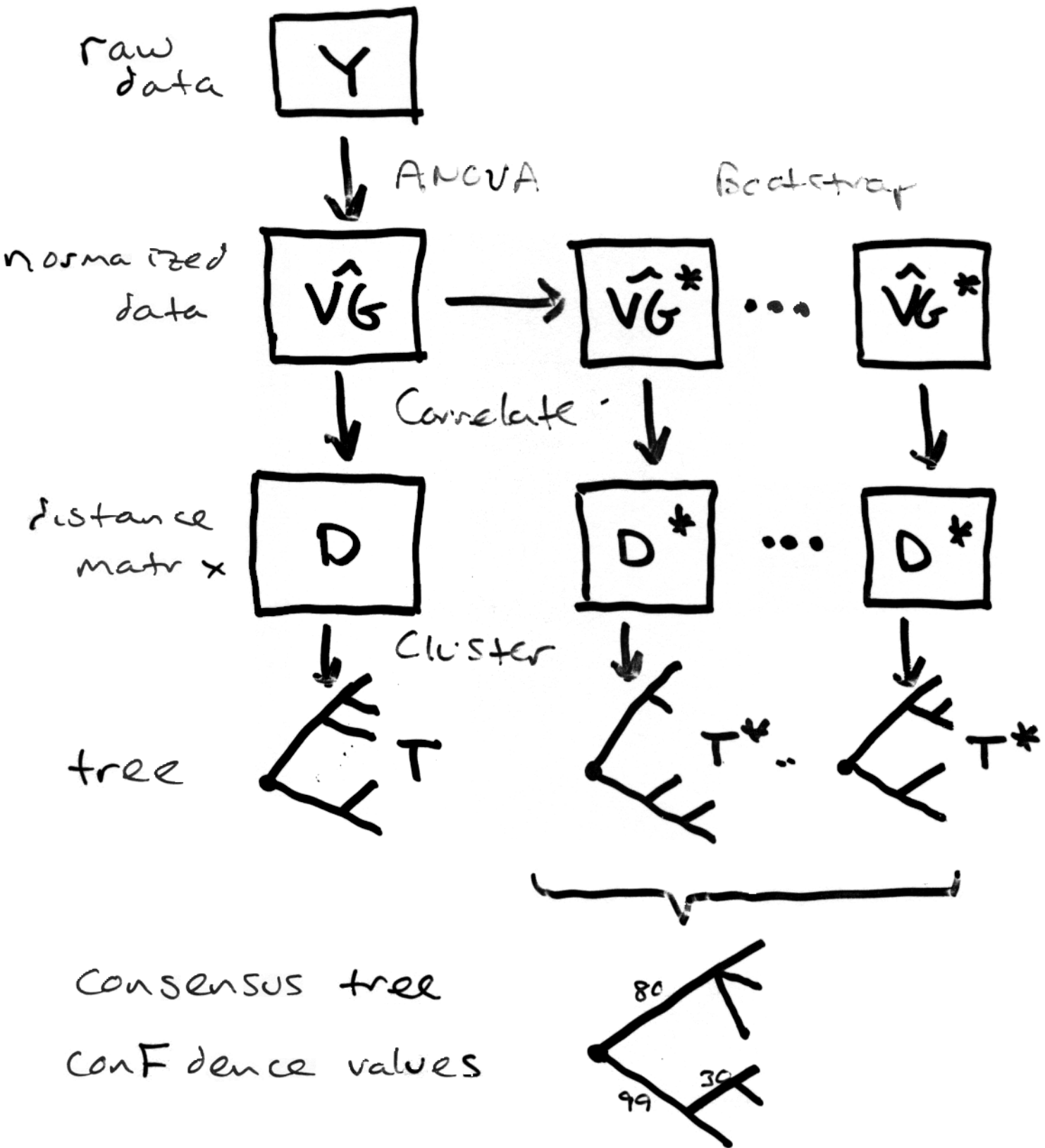
YSW1 (Middle)



SPS100 (Late)



Bootstrapping Cluster Analysis



Target Profiles

Chu et al

Figure 2

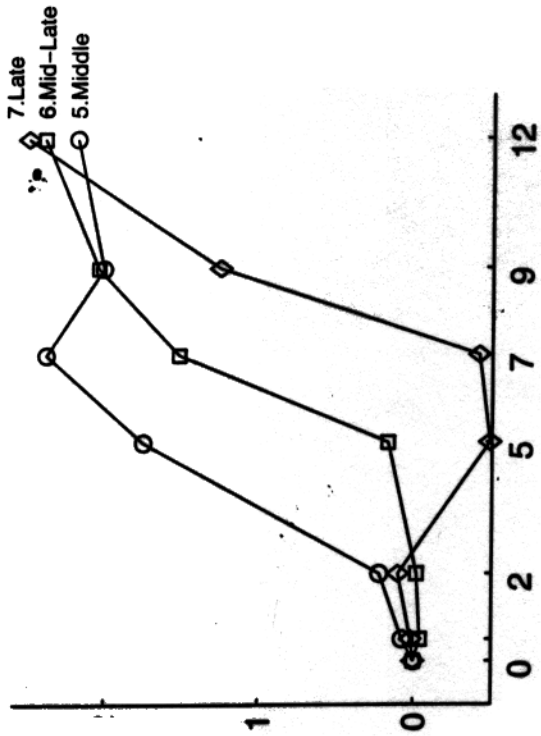
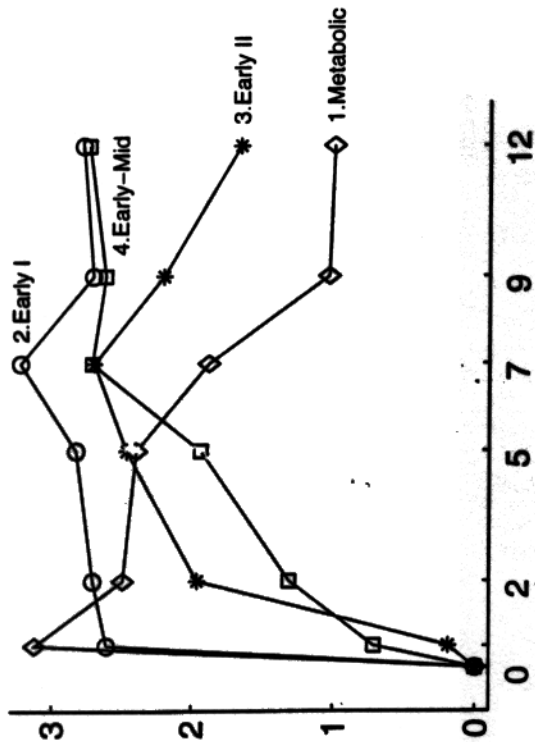
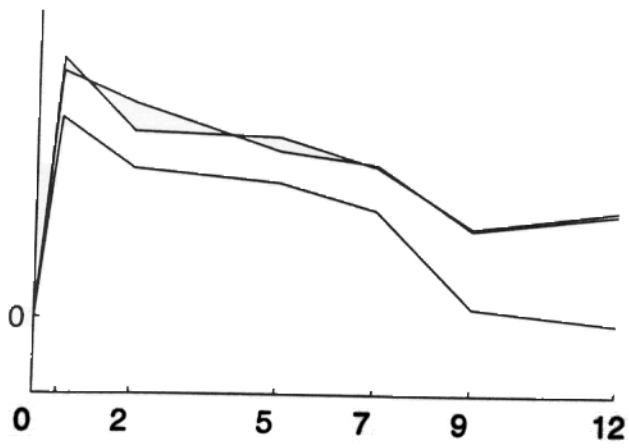
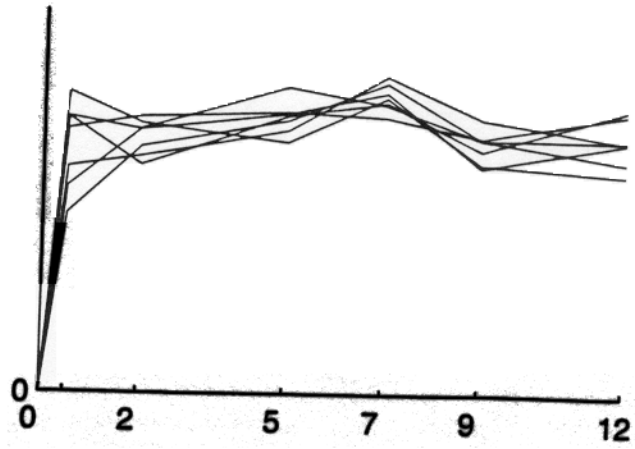


Figure 3

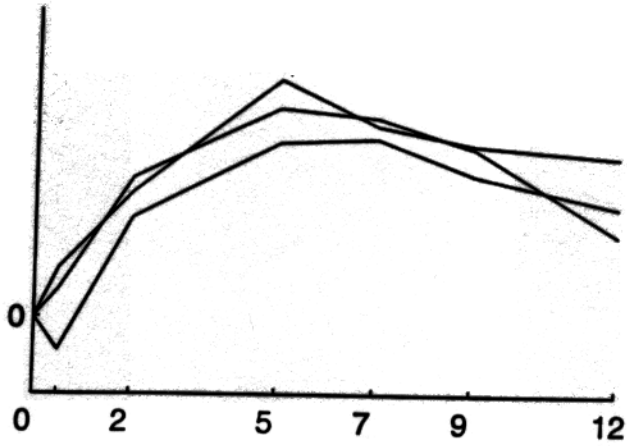
Profile #1,3 95% Stable Genes



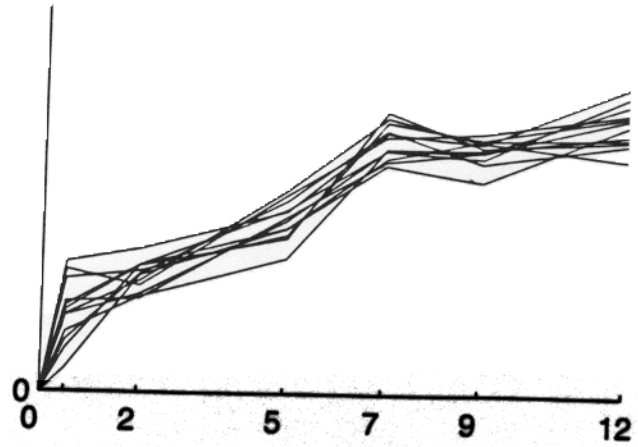
Profile #2,7 95% Stable Genes



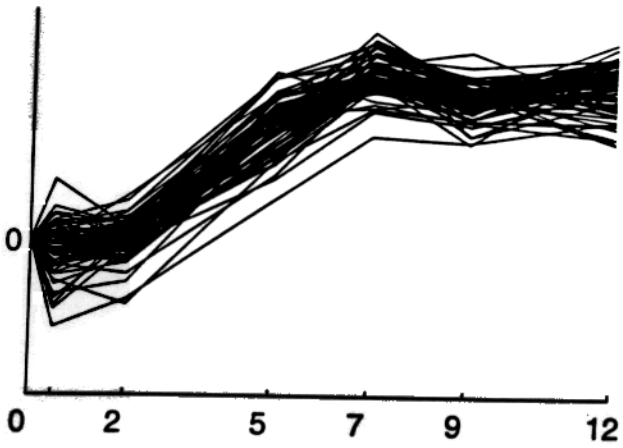
Profile #3,3 95% Stable Genes



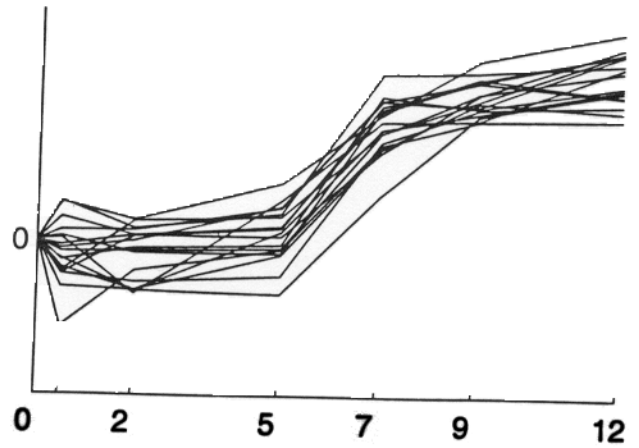
Profile #4,12 95% Stable Genes



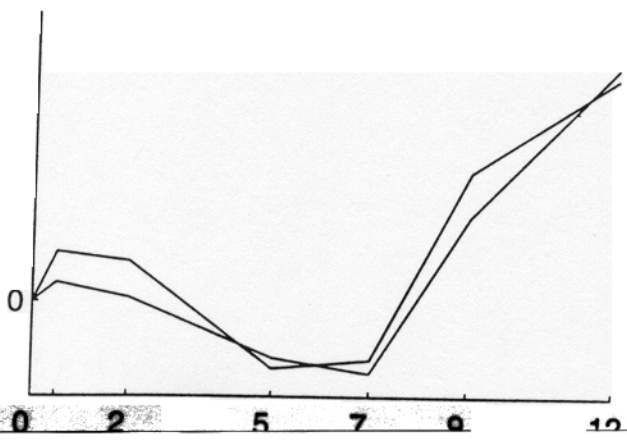
Profile #5,86 95% Stable Genes



Profile #6,17 95% Stable Genes



Profile #7,2 95% Stable Genes



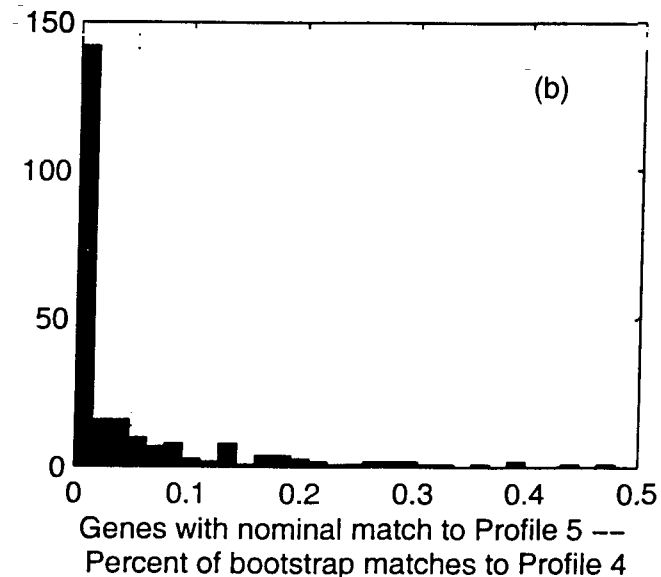
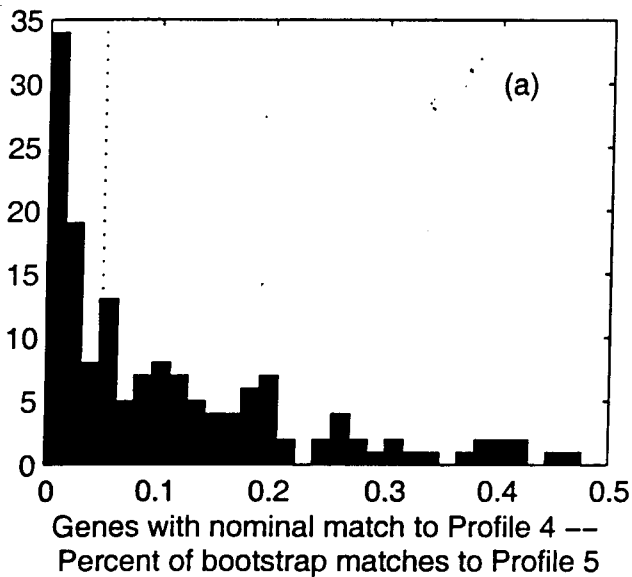
	(a) Chu <i>et al.</i>	(b) Nominal	(c) 95% Stable	(d) 80% Stable
Profile 1	52	65	3	8
Profile 2	61	51	7	11
Profile 3	45	74	3	11
Profile 4	95	151	12	27
Profile 5	158	241	86	120
Profile 6	61	145	17	36
Profile 7	5	15	2	6

Table 2: Number of genes matching to each profile for (a) Chu *et al.* clustering method, (b) modified clustering method with no reliability measure, (c) modified clustering method requiring 95% stability, (d) modified clustering method requiring 80% stability. Column (d) is included because our choice of 95% for stability is somewhat arbitrary.

Correlations among Chu et al Profiles

	Profile					
	2	3	4	5	6	7
1	.65	.19	.03	-.14	-.39	-.41
2		.73	.77	.60	.40	.11
3			.84	.78	.46	.01
4				.95	.84	.44
5					.83	.36
6						.75

Table 3: Pairwise correlations among the seven profiles.



Summary

Experimental Design Concepts For microarrays

- Confounding and partial confounding effects
- Role of replication is 2-fold
- Relative efficiency for comparing designs

OPEN QUESTIONS

Designs for structured experiments

Designs for large surveys

ANOVA as an analysis too

For microarrays

- Normalization of main effects
- Control of spot effects and D_xG_s
- Easy computation / least squares

OPEN QUESTIONS

Scale of measurement

Robustify

- modeling other effects
- Random effects + shrinkage
- gene specific error rates