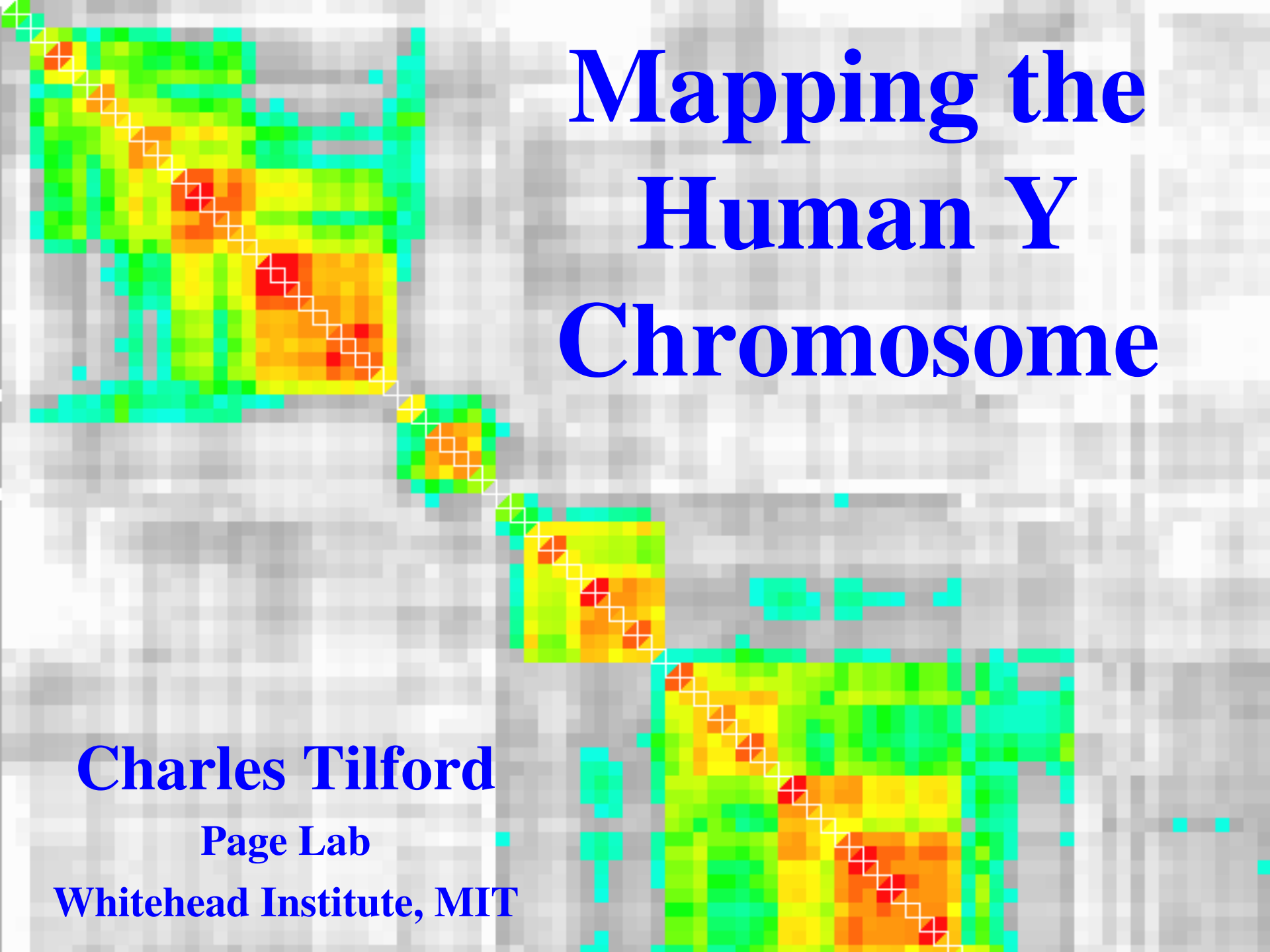


Mapping the Human Y Chromosome

Charles Tilford

Page Lab

Whitehead Institute, MIT



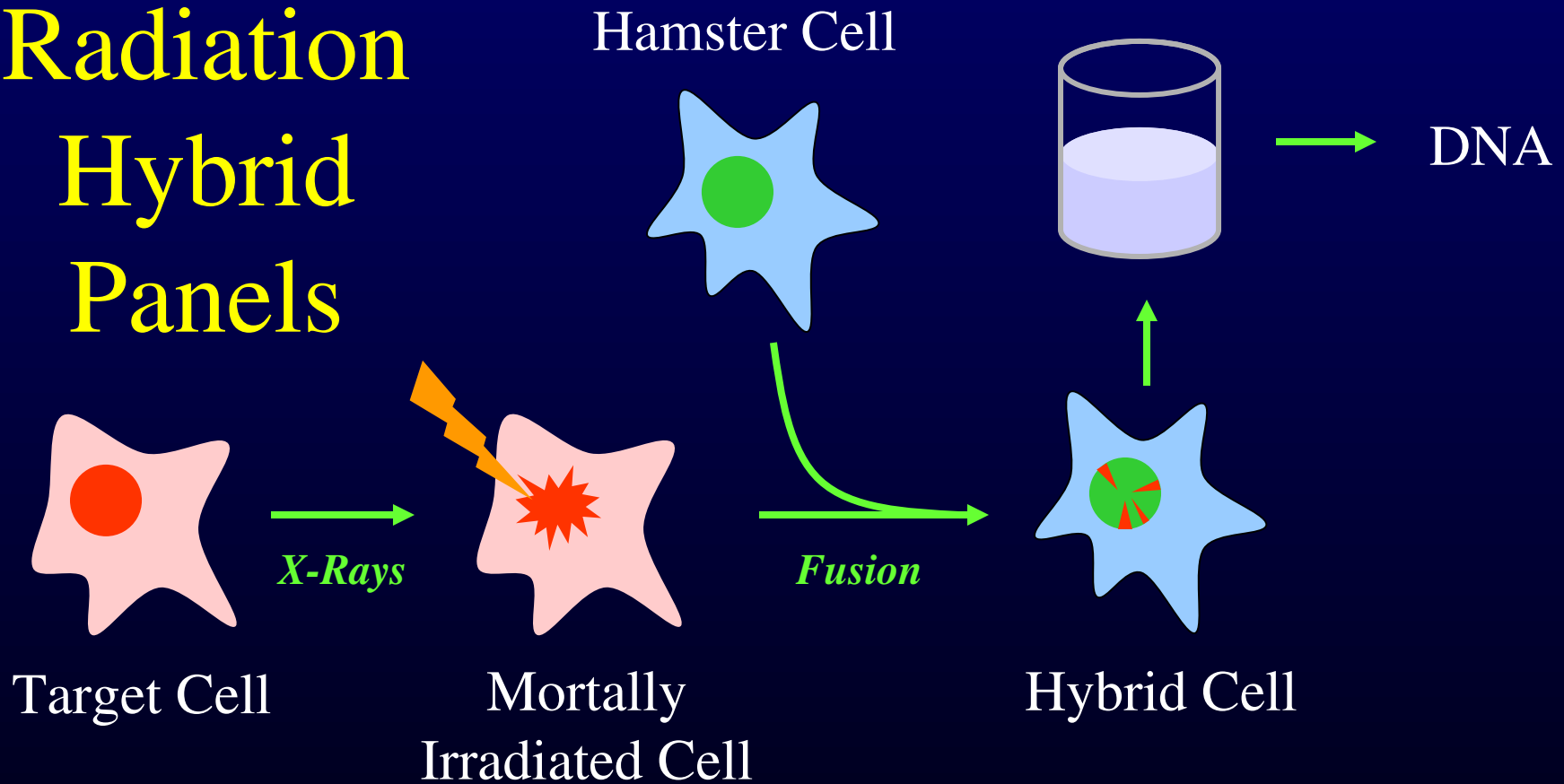
The Human Y Chromosome



Size/Mb:	2.5	8	1	15		30	0.35
Region:	PA					Heterochromatin	PA
	Euchromatic						
	Non-Recombining Region (NRY)						

- Male sex determining, passed through male lineage only
- Only the pseudoautosomal regions (PA) recombine with X
- Heterochromatin has 10kb sequence complexity
- NRY contains 20+ genes in two broad categories:
X-Y homologous, or Y-specific multicopy families.

Radiation Hybrid Panels



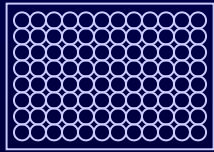
- Each hybrid has a full hamster genome, and some random fraction of the target genome
- X-ray intensity determines fragment size, which determines STS linkage range and ordering resolution

Radiation Hybrid Theory

- Fragments are independently retained/rejected by host cell. Each hybrid will have a unique subset of fragments.
- Unlinked markers will associate randomly, since they will always be on different fragments.
- Linked markers with "zero" separation will associate absolutely, and with decreasing frequency as the distance between them (and thus the probability of a radiation-induced break occurring) increases

Radiation Hybrids in Action

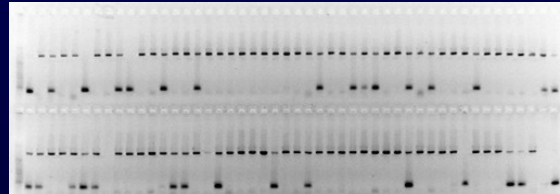
RH Panel



PCR



Agarose Gel



Scoring



Vector

10100100110...

*Generate
Vector for
each STS*



```
sY14      0000101000011000000...  
sY137     1000010010011000010...  
sY63      0000000000000000010...  
sY140     1000020010001000010...  
sY217     0000000000001000010...  
SP18-15   0000000000010001000...  
SP18-07   0000000000010000000...  
SP18-09   0000000000010000000...  
...
```

*Pairwise
Comparison*



Theta Matrix
(*Large*)



Distance and
Linkage Estimates

STS Order

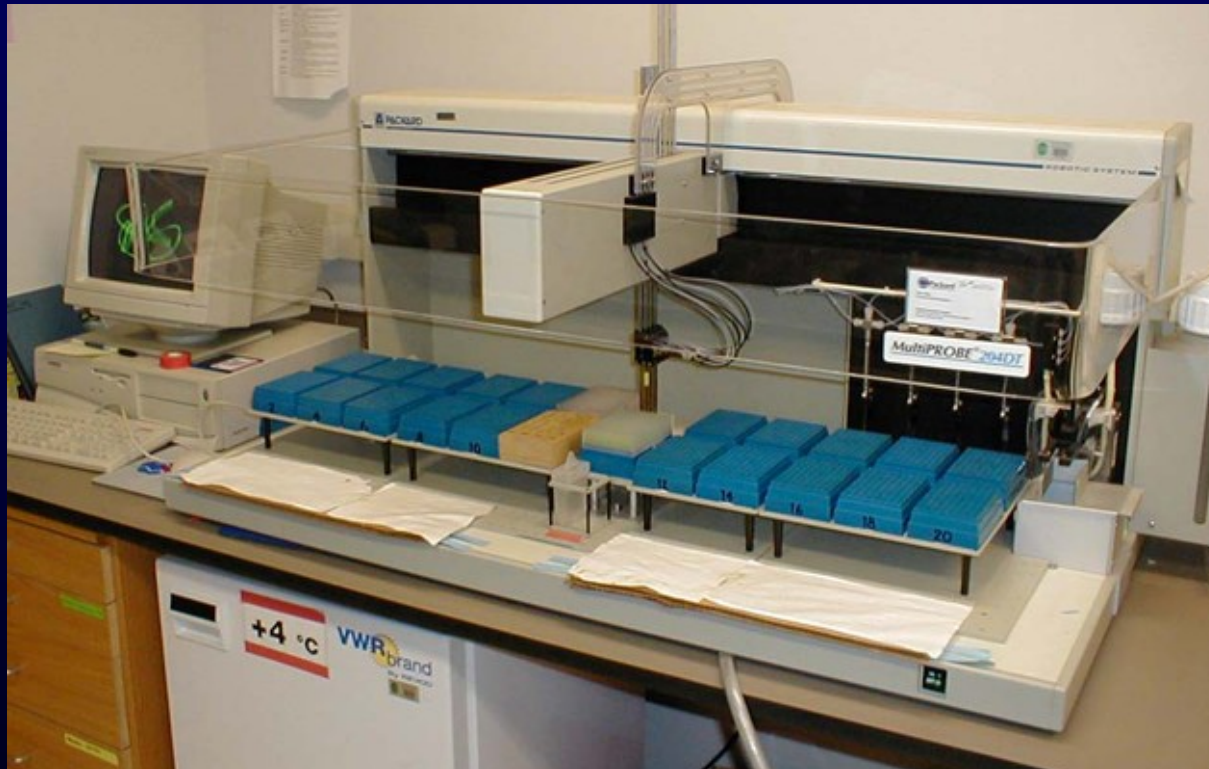
*Distance
Minimization*



Problem: Scale of Project

- Over 1000 markers to be mapped
- Each marker requires at least one plate of PCR, two plates for markers that provide useful map information
- PCR reactions are visualized with agarose gels
- Vector data needs to be extracted, stored and analyzed

High Throughput PCR



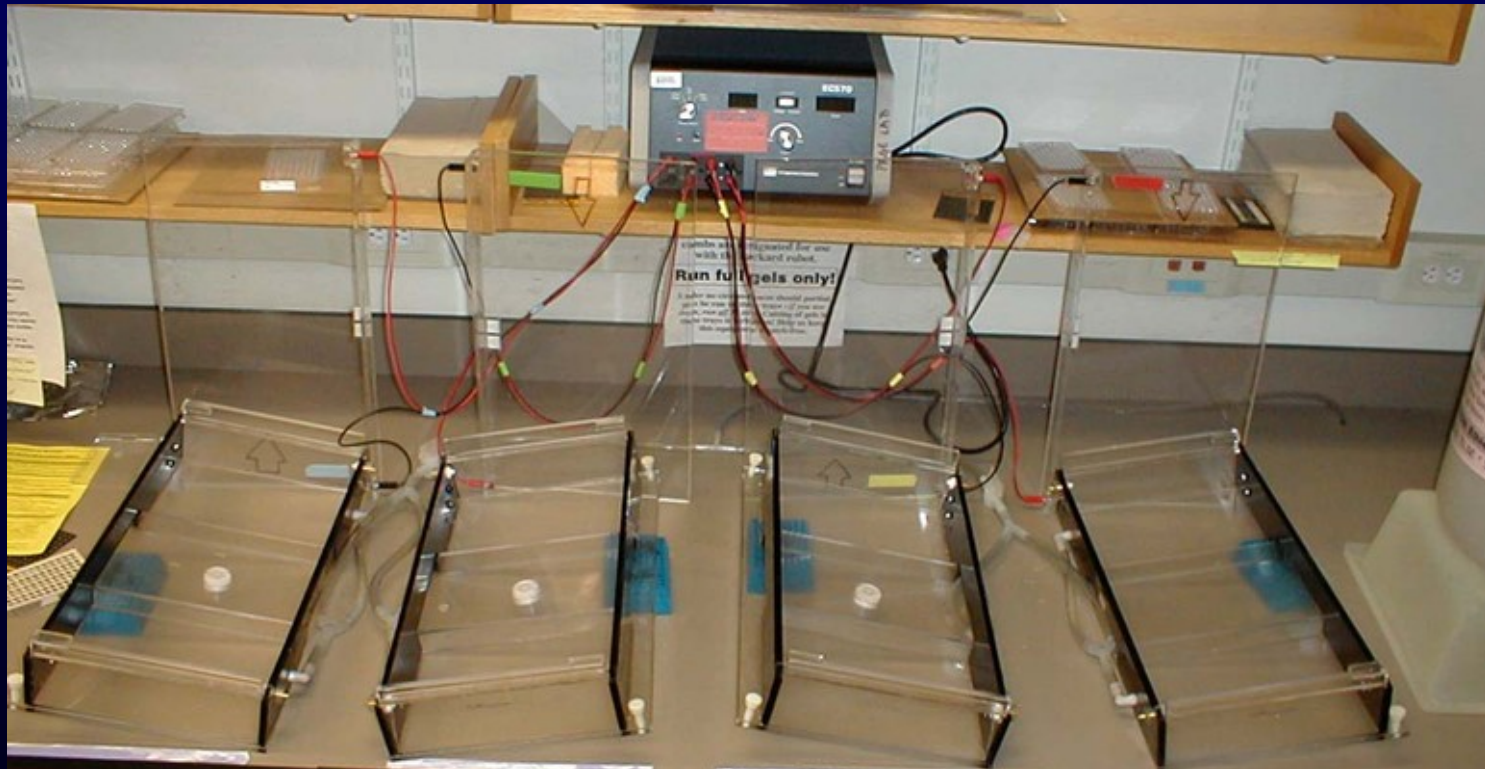
- Packard Multiprobe adapted for 20 x 96-well plates
- Adds template DNA, primers, PCR master mix, mineral oil and loading dye

High Throughput PCR



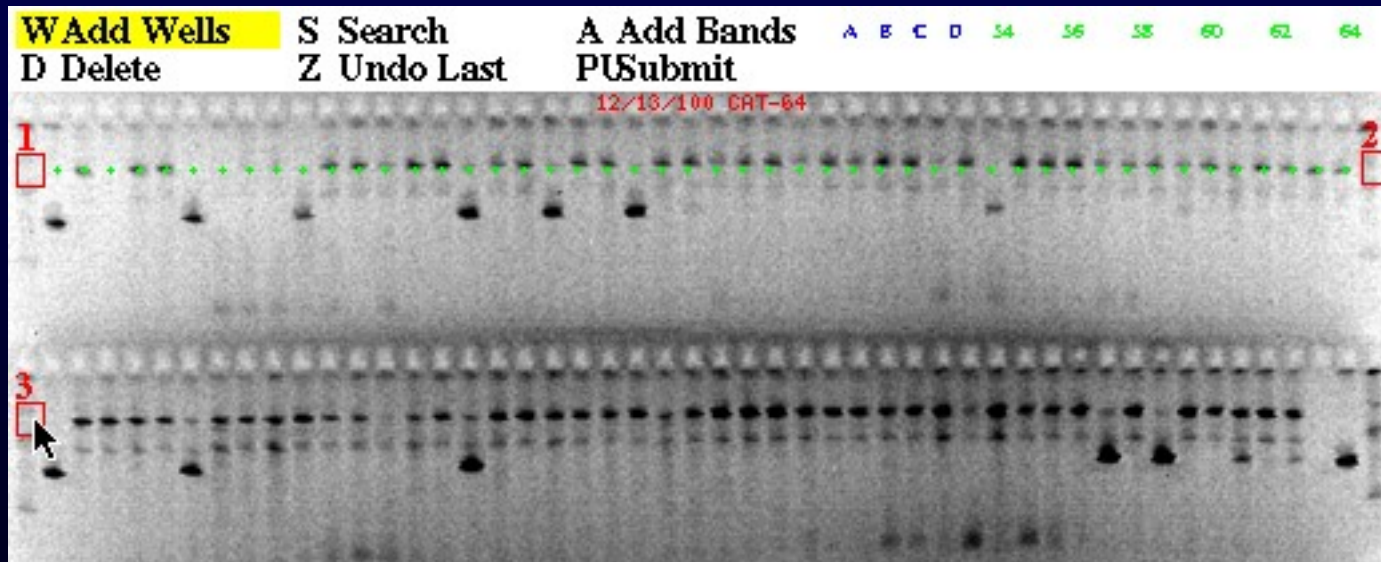
- Small, inexpensive, basic-function thermal cyclers (Hybaid Omn-E)

High Throughput PCR



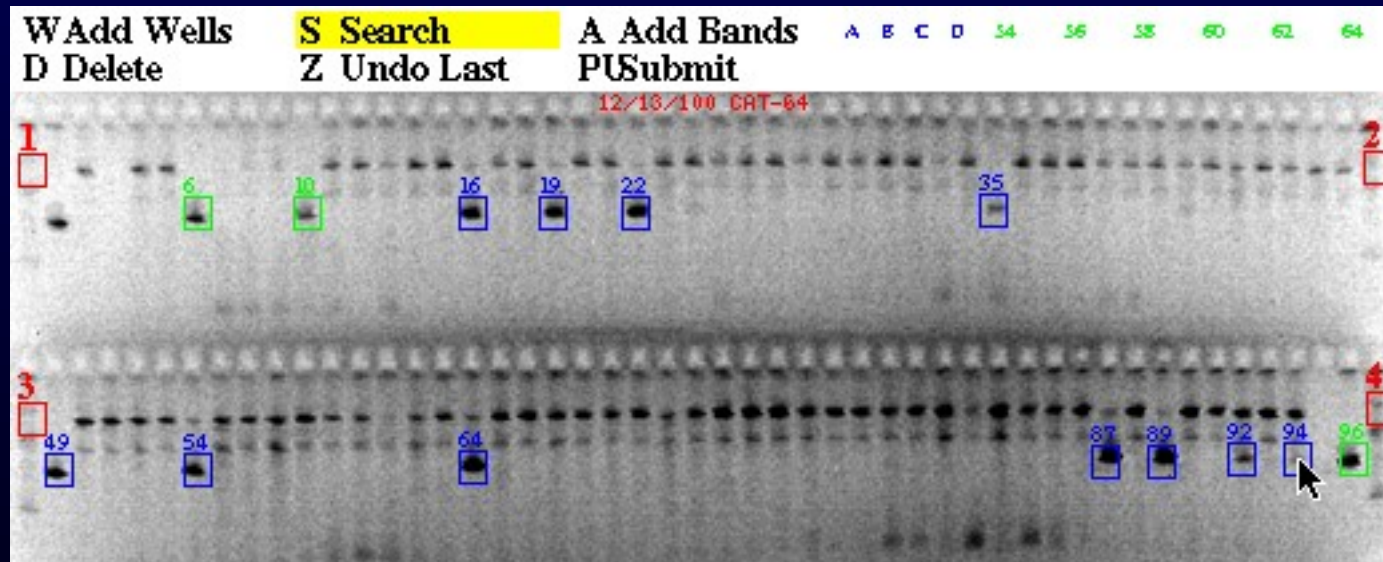
- Four Centipede gel boxes, each with 10 x 50-well combs
- Buffer plumbing system for rapid "dry" gel loading

Data Management - JavaGel



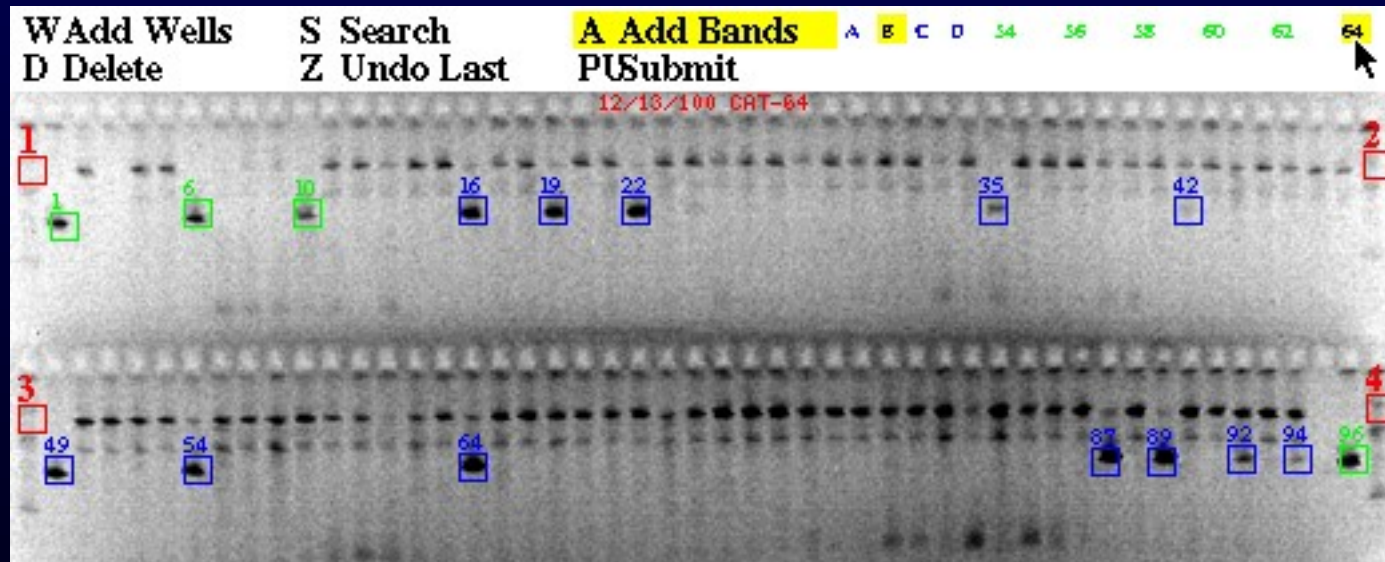
- Java applet for image annotation and data retrieval
- User clicks on the four edge marker lanes

Data Management - JavaGel



- In search mode, clicking on the faintest band finds all other bands of that size with the same or better S/N ratio

Data Management - JavaGel



- Quality and annealing temperature information are entered



11/11/2019

12/13/00

12/13/100 CAT-64

Con 12 15 18 31 68

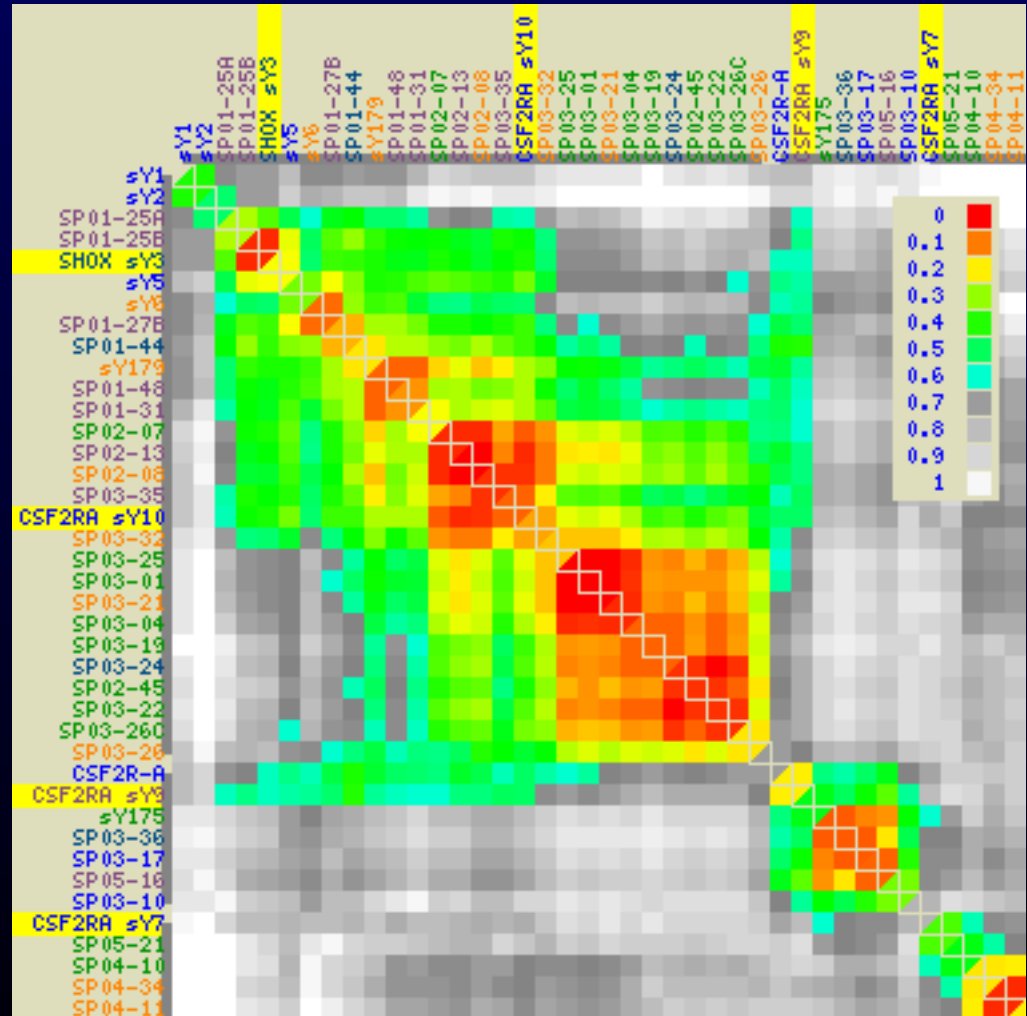
45 50 60 83 85 88 90 Con

78 Negative = 87%	A1 (Human) = 1	sY89 1	0.418	SP34-35	0.682	sY892	0.050
12 Positive = 13%	B1 (Hamster) = 0	sY892	0.050	sY89 1	0.418	sY914	0.000
0 Ambiguous = 0%	B3 (Human 20%) = 1	sY914	0.000	sY892	0.050	sY912	0.180
	B5 (Human 5%) = 1	sY912	0.180	sY914	0.000	BE508K05	0.310
	G12 (Blank) = 0	BE508K05	0.310	sY912	0.180	sY920	0.269
	H12 (Human) = 1						

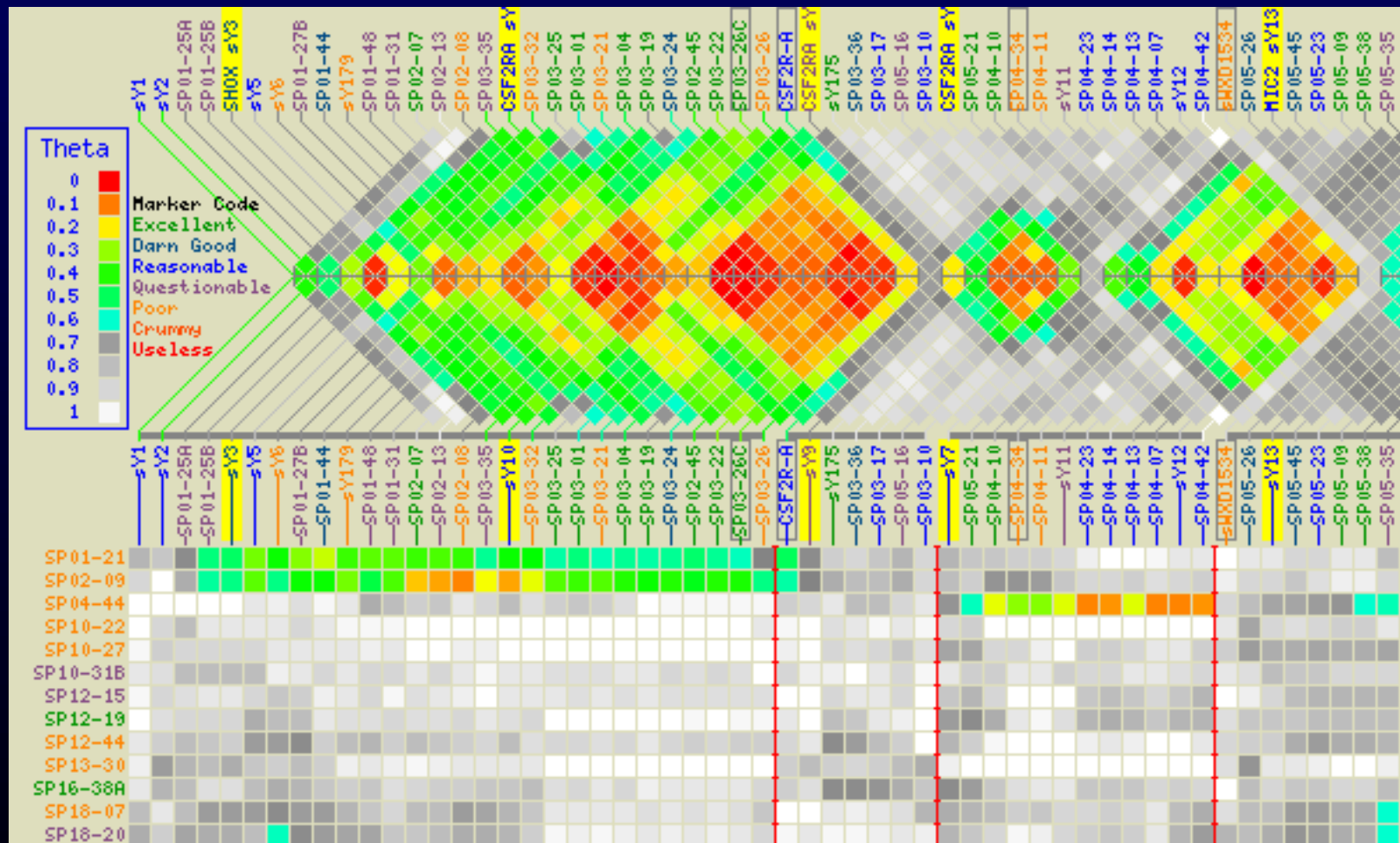
Controls, best hits

Graphical View of Theta Matrix

- Theta values for likely linkage represented by spectrum
- Likely unlinked values (>0.6) shown as gray scale
- Color of marker name indicates quality of STS
- Diagonal values (technically 0) are split between flanking markers



Diagonal View of Matrix



Ordering Markers

- User needs to grossly order markers due to multiple Y-specific repeat families
- Algorithmic ordering minimizes subdiagonal distance:
 - Exhaustive ordering within a small (6-8 marker) sliding window
 - "Random cost" ordering generates excellent order with much larger marker sets
 - Wang et al. A fast random cost algorithm for physical mapping. Proc Natl. Acad. Sci. USA 1994 (91) 11094-11098

RH "Simulator"

Radiation Hybrid Simulator			
<i>Mouse-over any field for help on the functioning of each setting</i>			
Model Statistics Generation		Help / FAQ View Statistics for 49 simulations. View Vectors for 1387 markers. View previous Output 0 Mb.	
Panel: A*B	Execute	kb / cR = 5.3	Retention Frequency 12
Scan from <input type="text"/> kb to <input type="text"/> kb in <input type="text"/> steps.	<input checked="" type="checkbox"/> Log Spacing	Number of Hybrids 90	
Load File Click "Execute" to load	Custom Marker Order:	Local Vectors: All Combinations Grid is: Ordered	
<input type="button" value="Browse..."/>	sY1 sY2 SP01-25A	Discard: LOD >3 Ambiguous >6 Distance >350 kb. Show <input type="text" value="1"/> next-best matches. <input checked="" type="checkbox"/> Discard extreme matches <input checked="" type="checkbox"/> Calculate Linkage map Estimate 11.6 seconds to calculate neighbors, 2.2 hours for all combinations.	
<input type="checkbox"/> Parse on-the-fly	Compare to Sequence		

- Novel method for determining linkage likelihood
- Developed due to concern that existing algorithms would be less likely to find linkage in marker pairs with differing copy number :

Triploid 010100010100110010011000100101

Haploid 01000001000000000000110000000000

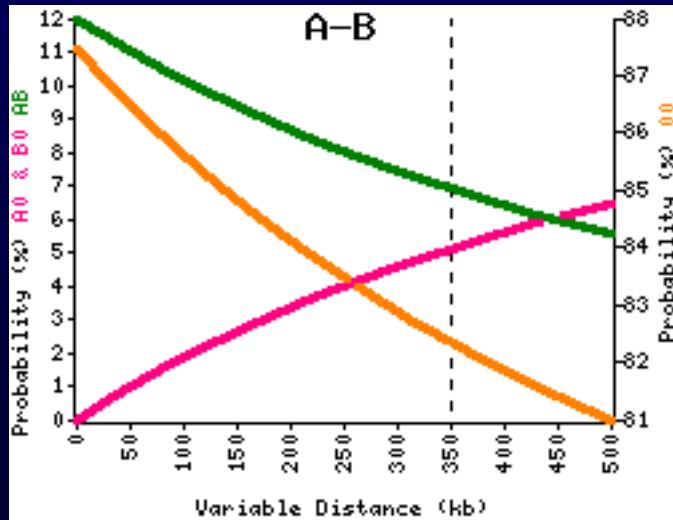
RH Simulator Theory

- System is provided with "models" of marker arrangement:
 - A-B A-B A-B A-B A-B B B *etc.*
 - Each model provides linkage information only, not order
- For each model, the probability is calculated for a hybrid being one of four classes in a pairwise comparison:

Marker-A	0	0	0	1	0	0	1	1	0	0	0	1	0	0	1	0	1	0	0	0	1	0	0	0	1	0	0	0	...	
Marker-B	0	0	1	0	0	1	0	0	1	0	0	0	1	1	0	0	1	0	1	1	0	0	0	0	0	0	0	1	0	...
				A0				B0							AB											00				

- These probabilities depend on:
 - The model, the distance between linked markers, the retention frequency (RF) and the average kilobases per centiRay. The last two values need to be observed from a subset of actual data.

Model Probability Curves



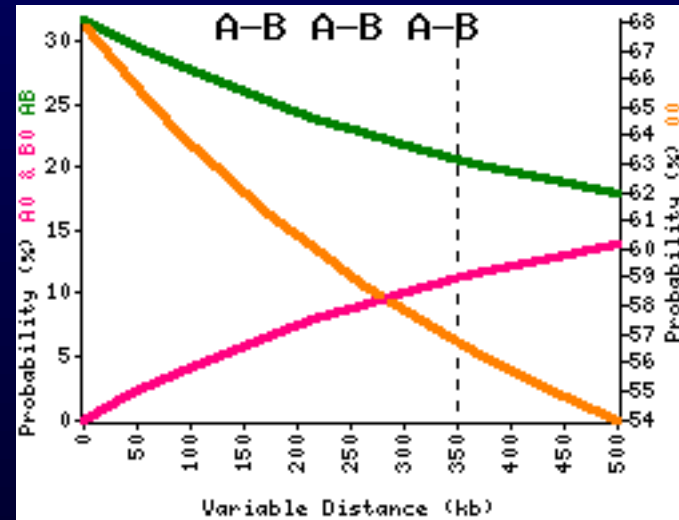
Linked Haploid

$A0 = A^+ B^-$
 $B0 = A^- B^+$

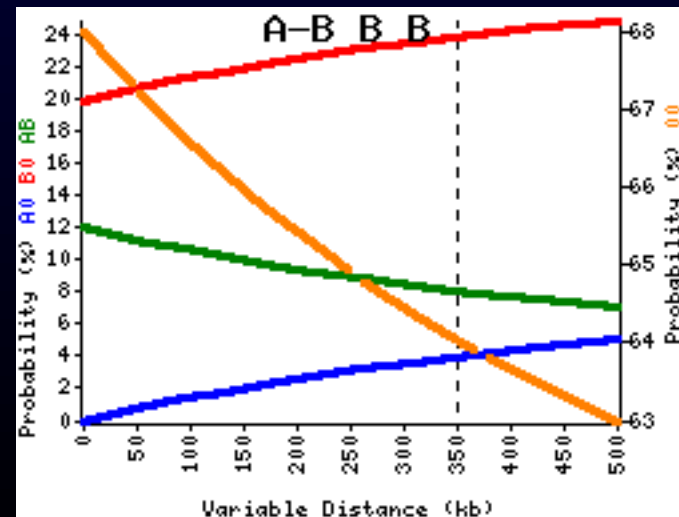
Shown in pink when probabilities are equal

$AB = A^+ B^+$

$00 = A^- B^-$ shown on right axis



Linked Triploid



Haploid - Triploid

Overall Probability

- Given a vector pair, the number of hybrid pairs that are 00, AB, A0 & B0 are summed (for example, 79, 8, 1 & 2)
- Using the probability curves, the overall probability that a model would generate that set of 4 sums can be calculated for each distance. The distance with the best (greatest) probability is then reported for the model.
- These "model probabilities" can then be compared to each other to determine the likelihood of one model over the others

Ambiguous Data

- Markers are typically tested twice, on different days
- On average, about 1-3% of hybrids show conflicts between duplicate trials (i.e. one trial is positive, the other negative)
- Conflicts result from error (often gel loading, but also PCR failure) or marginal assays with "consistent inconstancy"
- Such ambiguous data are represented as "2" in the vector
- Program will analyze all possible vector pairs when ambiguous data are encountered, and report the average results with standard deviations

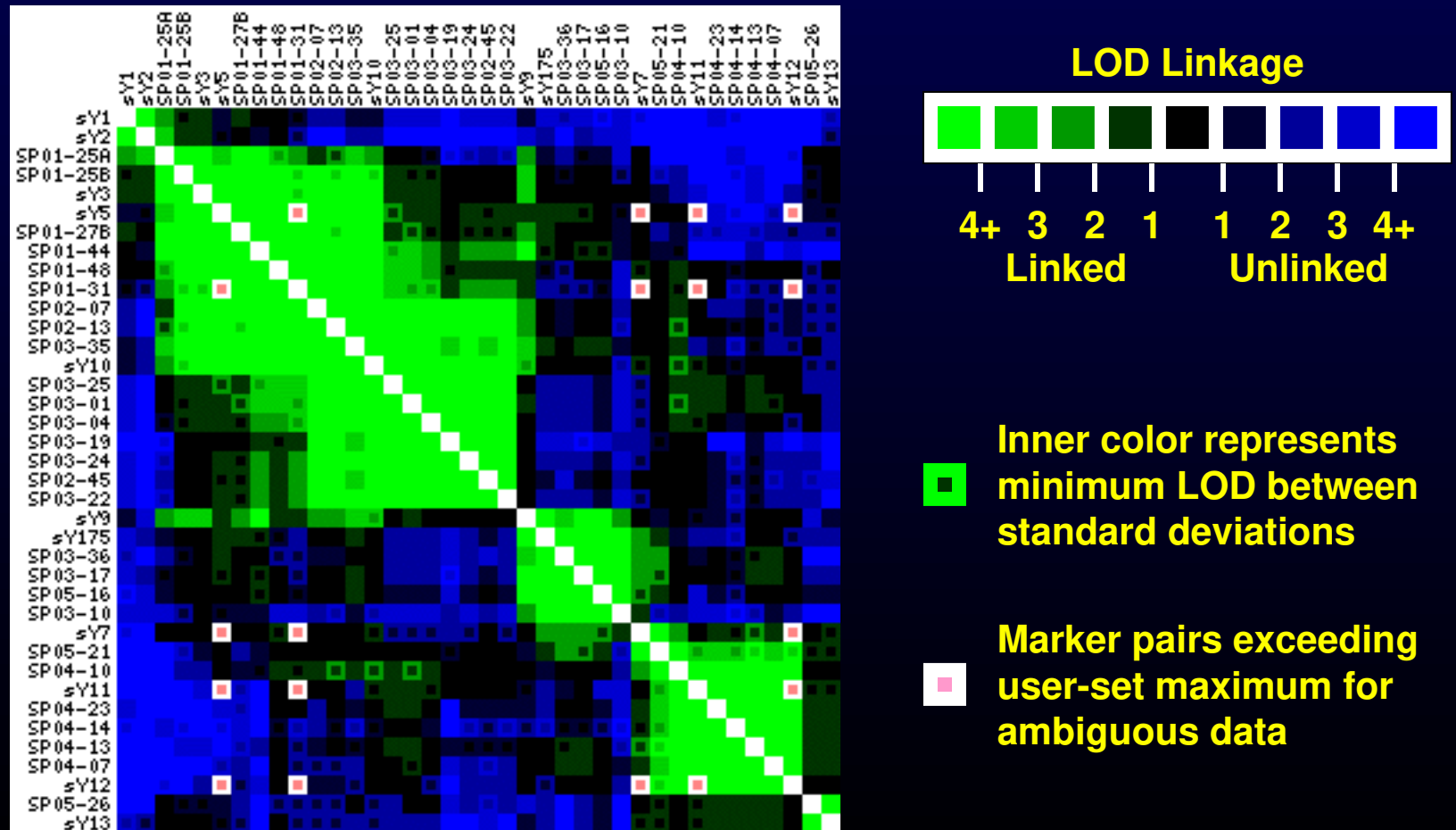
Sample Model Prediction Output

	Model	$-\log_{10}(P)$	Distance
A = sY1 B = sY2	A-B B	1.95 ± 0.12 at	0 ± 0 kb
	A-B B B	3.04 ± 0.16 at	0 ± 0 kb
	A-B B B B	5.44 ± 0.30 at	0 ± 0 kb
	A-B A-B	5.85 ± 0.34 at	256 ± 29 kb
	A A-B B	6.16 ± 0.19 at	0 ± 0 kb
	A A-B B B	6.97 ± 0.14 at	0 ± 0 kb
	A B B	8.22 ± 0.55	Unlinked Model
	A-B B B B B	8.37 ± 0.41 at	0 ± 0 kb
	A-B A-B A-B	8.62 ± 0.47 at	194 ± 22 kb
	A A-B B B B	9.13 ± 0.25 at	0 ± 0 kb
	A B	9.32 ± 0.53	Unlinked Model
	A B B B	9.41 ± 0.59	Unlinked Model
	B B-A A A	10.43 ± 0.26 at	0 ± 0 kb
			etc...

- Observation - likelihoods are insufficient to choose a single model with high certainty.

Linkage Likelihood with Models

- By subtracting the best linked and unlinked models, a likelihood of linkage for STS pairs can be generated.

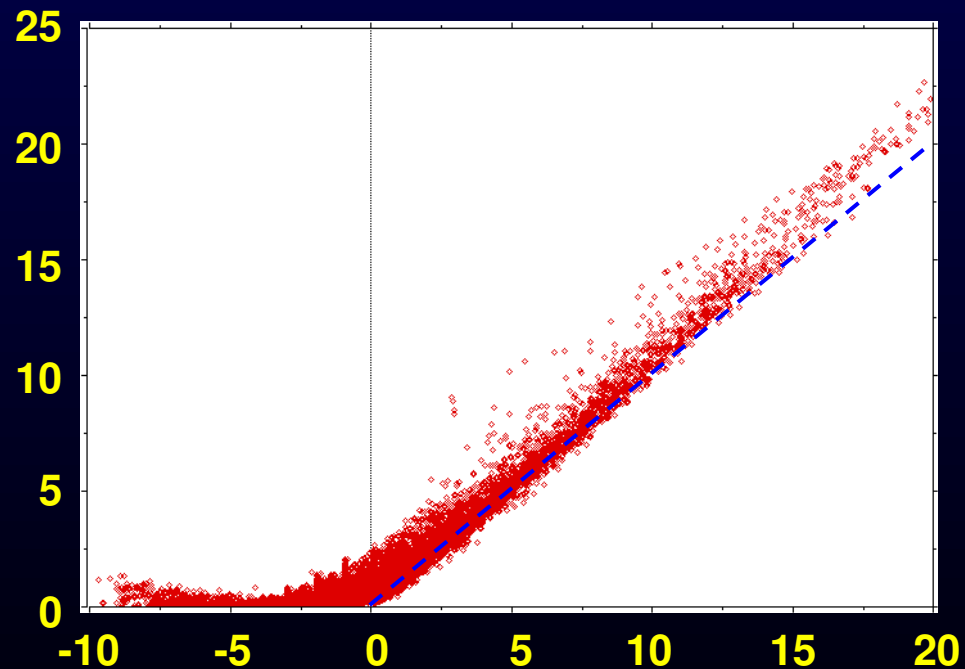


Comparison to Standard Methods

- LOD linkage results are similar to those obtained by standard algorithms, but tend to be lower.

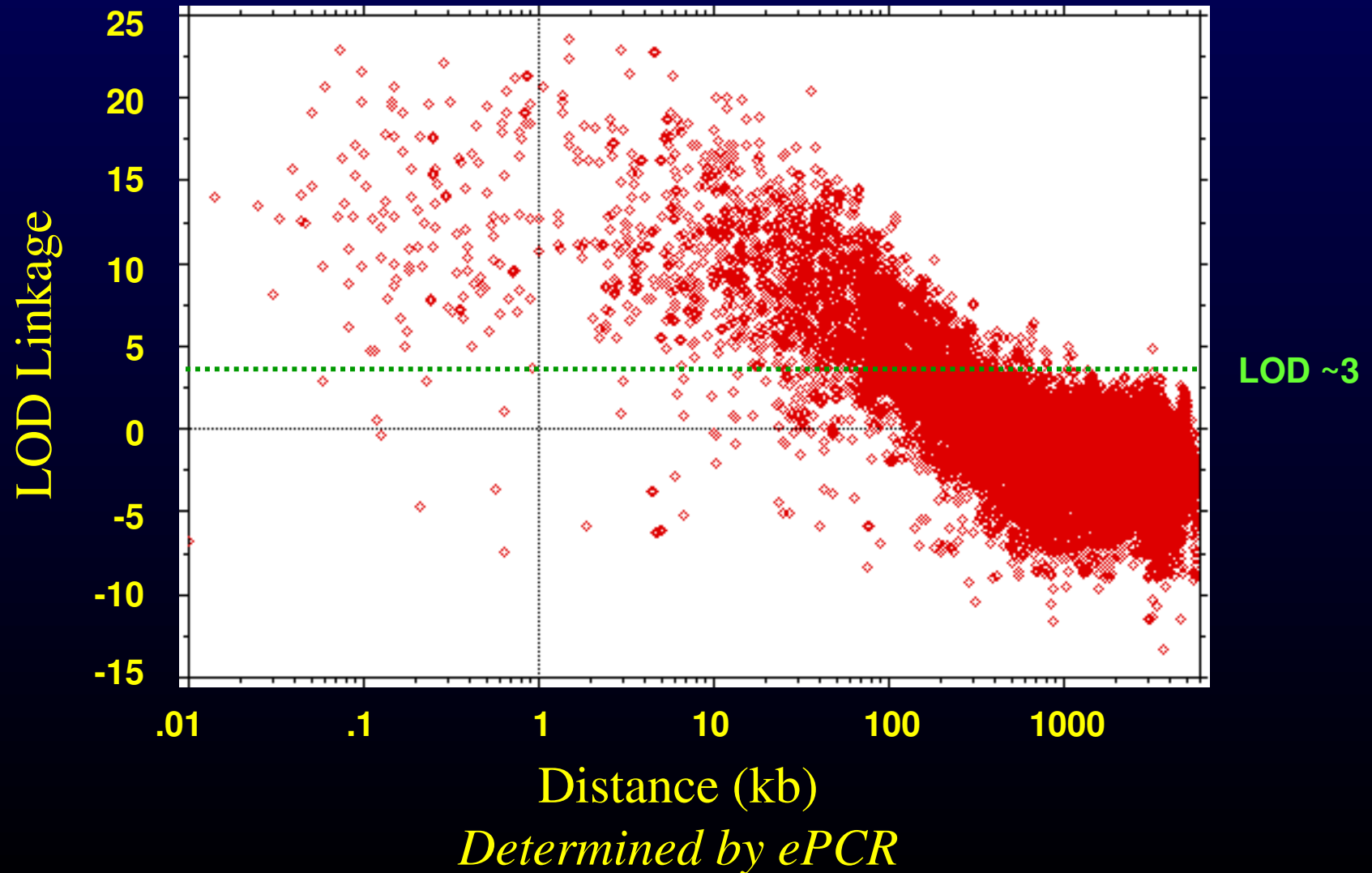
LOD Linkage as per:

Jones HB. Pairwise analysis of radiation hybrid mapping data. *Ann. Hum. Genet.* 60, 351-357 (1996)



LOD Linkage as per Simulation

LOD score vs. Distance



Conclusions

- Model simulation provides a stringent assay for determining linkage between RH vectors
- Likelihoods of copy number based on retention frequency are too similar to allow reliable choice of one model over all others when comparing markers.
- Final map has 700+ markers, with 11 gaps. Sequence comparison shows reasonable order within contiguous segments.

Acknowledgments

Steve Rozen

Programming Consultation

Helen Skaletsky

Statistics / Programming Consultation

Loreall Pooler

Marker Typing

Chad Nussbaum

Mapping Consultation / Reagents

Tom Hudson

Mapping Consultation

Eric Lander

Idea Guru

David Page

Patron and Mentor

RH Simulator Theory I

- System is provided with "models" of marker arrangement:
 - A-B A-B A-B A-B A-B B B *etc.*
 - Each model provides linkage information only, not order
- Two basic levels of probability calculations:
 - Probability of a linked pair (A-B) becoming separated is function of distance between them (d) and "intensity of radiation" ($\lambda = 0.01/\text{kb}_{\text{percR}}$): $\theta = 1 - e^{-\lambda * d}$
 - Probability of retaining a fragment is retention frequency (RF)
 - Both RF and λ are experimentally determined
 - RF and λ may vary across chromosome, but are treated as constant

RH Simulator Theory II

- We can simplify the system by considering only the four possible outcomes for a pairwise comparison of a hybrid:
 - both positive (AB), both negative (00), or A0 or B0
- Models are composed of independent components, calculate probabilities for each of three component types:

	A	B	A-B
pAB	0	0	$(1-\theta) RF + \theta RF^2$
p00	1-RF	1-RF	$(1-\theta)(1-RF) + \theta(1-RF)^2$
pA0	RF	0	$\theta RF(1-RF)$
pB0	0	RF	$\theta RF(1-RF)$

RF = retention frequency, θ = break probability (function of distance)

RH Simulator Theory III

- Find total probabilities for the model as:

$p_{00_{\text{Tot}}}$ Product (p_{00}) for all components

$p_{A0_{\text{Tot}}}$ (Product ($p_{A|0}$) for all components) - $p_{00_{\text{Tot}}}$
where $p_{A|0} = p_{A0} + p_{00}$

$p_{B0_{\text{Tot}}}$ (Product ($p_{B|0}$) for all components) - $p_{00_{\text{Tot}}}$
where $p_{B|0} = p_{B0} + p_{00}$

$p_{AB_{\text{Tot}}}$ $1 - (p_{00_{\text{Tot}}} + p_{A0_{\text{Tot}}} + p_{B0_{\text{Tot}}})$

- Calculate these total probabilities for each model, and for a finite set of distances (generally 12)

RH Simulator Theory IV

- We can now apply these four probabilities to determine the probability that a given model would generate an observed vector pair at a specified distance:

$$P = \frac{p_{00}^w p_{A0}^x p_{B0}^y p_{AB}^z N!}{w! x! y! z!}$$

Where:

w = number of hybrids 00

x = number of hybrids A0

y = number of hybrids B0

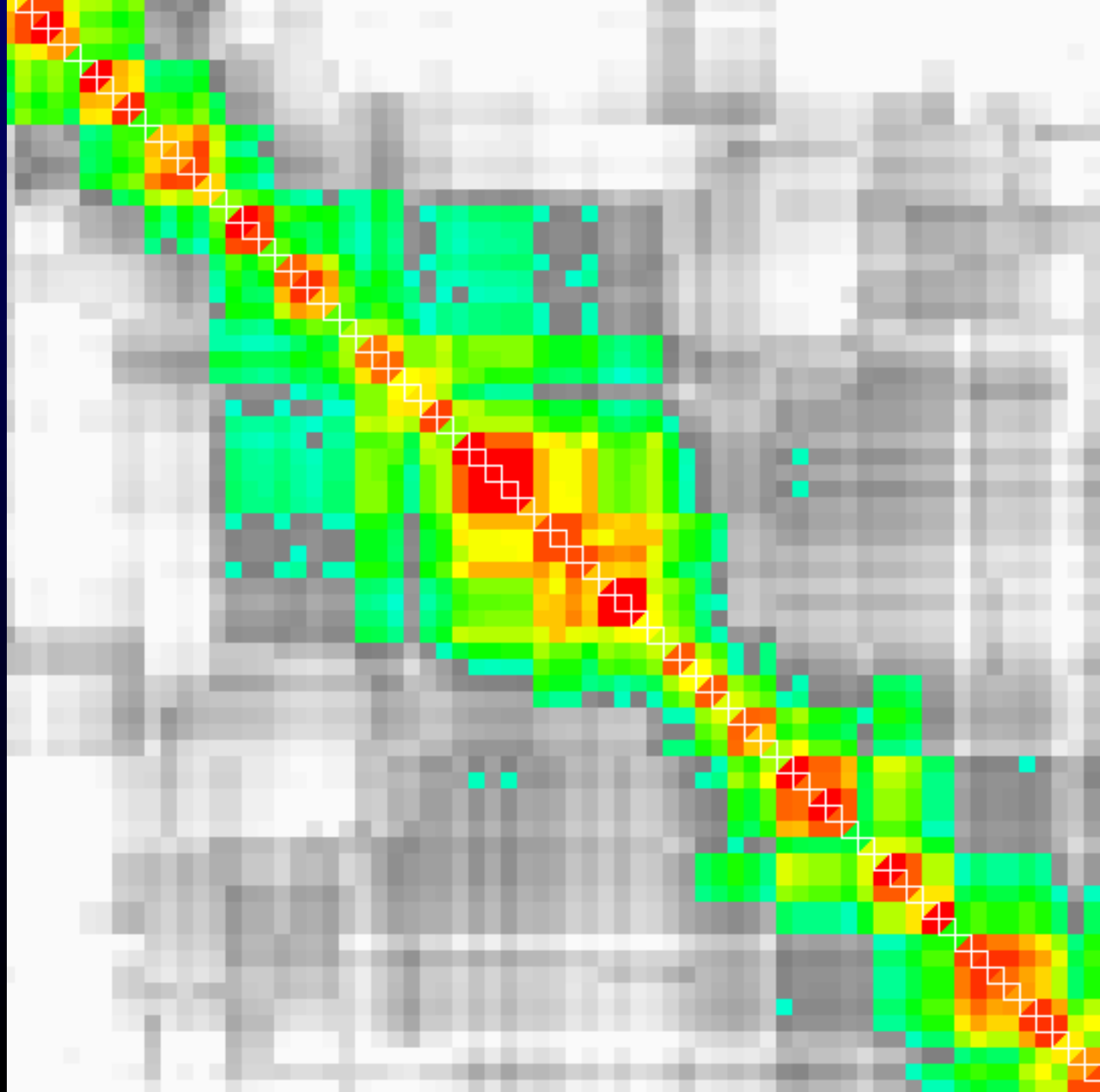
z = number of hybrids AB

N = total number of hybrids

Gap Analysis

LOD=3						
Gap	Left STS	Link	Copy#		Distance / Comment	
		LOD	Lft	Rgt		
----	-----	-----	----	----	-----	
	SP03-26	0.10	1-2	1-2	No Seq Data (not real gap)	
1	SP03-10	2.52	1-2	2-5	No Seq Data	
	SP04-42	2.23	1-3	1-3	No Seq Data (not real gap)	
2	sY23	2.41	1-4	2-5	439 kb	
3	sY1029	-0.41	1-2	2-5	125 kb + Near site of inversion...	
4	sY900	-2.32	1-3	1-4	160 kb +	
5	SP14-22	-1.36	1	1-2	540 kb	
6	sXY701	-2.00	1-2	1-3	270 kb	
7	sY715	1.20	6+	1-2	Centromere	
8	SP30-03	-0.92	1-2	1-2	100 kb	
9	SP37-38	0.69	1-2	1-3	147 kb	
10	SP51-36	-0.91	1-2	1-4	??	
11	sY707	1.93	1-4	1-2	??	

"Well
Behaved"
Markers



"Poorly
Behaved"
Markers

