

Copy number variation in association studies of human disease

Steven McCarroll

Broad Institute of MIT and Harvard

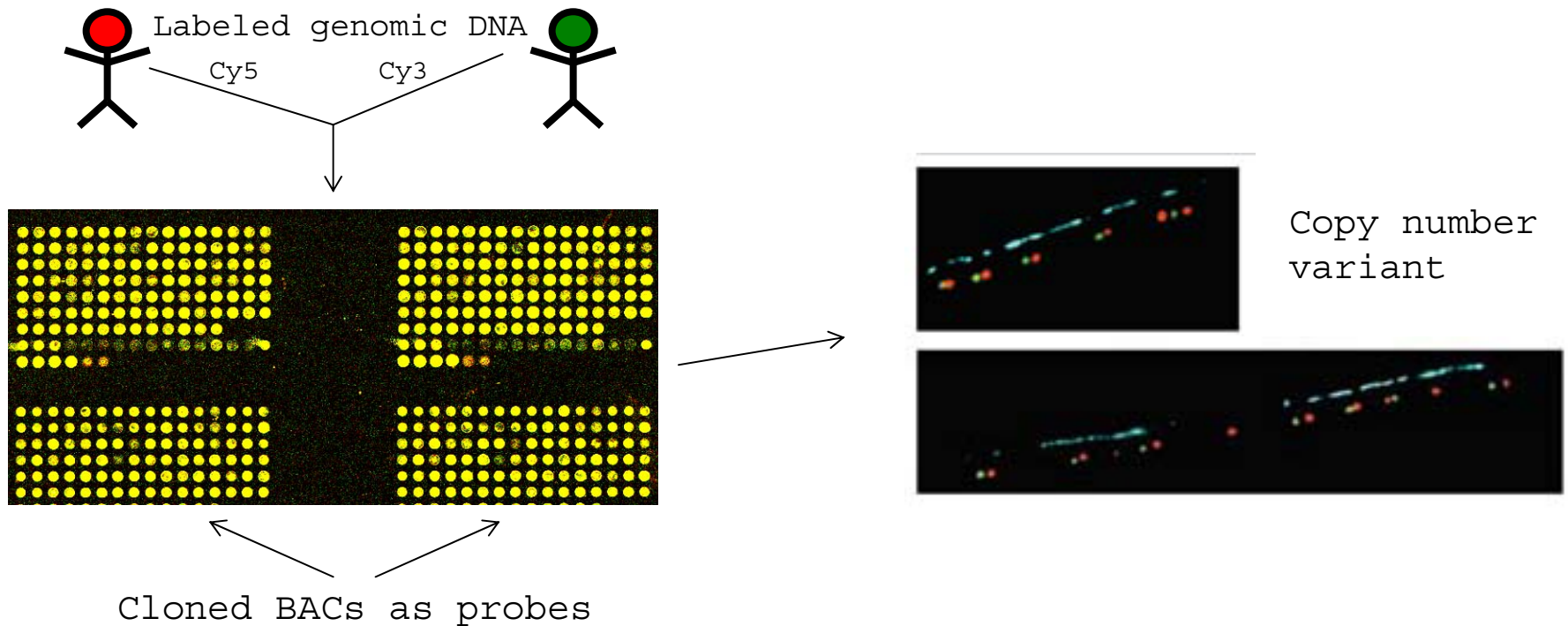
Outline

- Copy number variation
- Re-designing SNP arrays to interrogate copy number variation
- Analysis of copy number variation in whole genome association studies

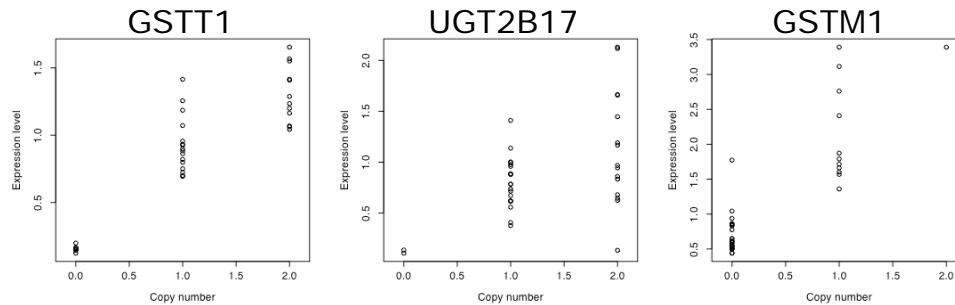
Copy number variation is common and extensive

Sebat et al., *Science* 2004

lafrate et al., *Nature Genetics* 2004



Copy number variation shown to affect gene expression and some disease phenotypes



Deletions - McCarroll et al., 2006
>60 CNVs - Stranger et al., 2007

QuickTime™ and a
TIFF (Uncompressed) decompressor
are needed to see this picture.

HIV progression - Gonzalez et al., 2005
Glomerulonephritis - Aitman et al., 2006
SLE - Faniculli et al., 2007



CCL3L1 CN associated with HIV progression to AIDS

Background:

CCL3L1 is the most potent known ligand for
CC chemokine receptor 5 (CCR5),
the major coreceptor for HIV,
and it is a dominant HIV-suppressive
chemokine

QuickTime™ a
TIFF (Uncompressed) de
are needed to see this



FCGR3 CN associated with systemic autoimmune disease

QuickTime™ and a
TIFF (Uncompressed) decompressor
are needed to see this picture.

Aitman *et al.*, Science 2006

Faniculli *et al.*, Nature Genetics 2007

How broadly does copy number variation influence clinical phenotypes?

- A major outstanding question in the field

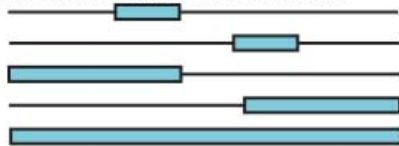
Need for high-resolution maps and genotyping technology

A

BAC probe on which variation has been observed:



Potential variants consistent with the observation:

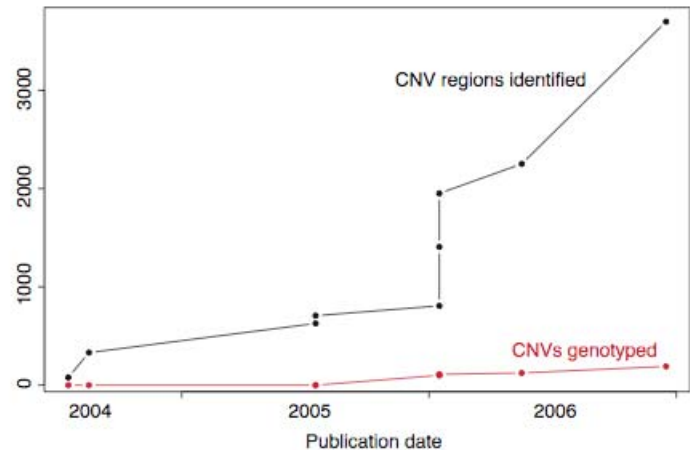


B

Fosmid end pair sequences from which deletion has been inferred:



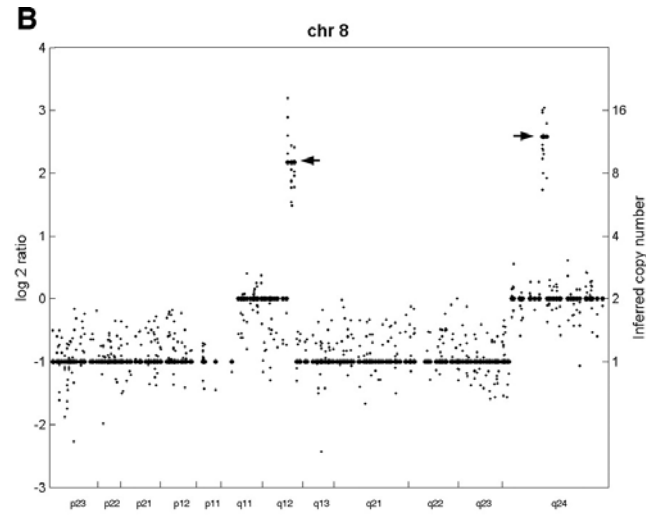
Potential deletions consistent with the observation:



Outline

- Copy number variation
- Re-designing SNP arrays to interrogate copy number variation
- Analysis of copy number variation in whole genome association studies

CN analysis on earlier SNP arrays



Zhao et al., Cancer Res 2004

Little/no coverage of most common germline CNPs

SNP probes weren't optimized for copy number analysis

How could we improve CN analysis on SNP arrays?

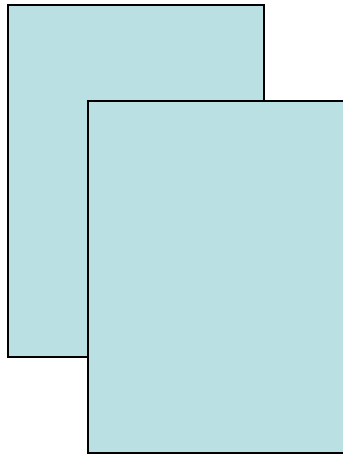
- Optimize probe locations in genome (within constraints of platform)
 - More-uniform coverage
 - Dense coverage in regions of interest (e.g. reported CNV)
- Optimize probe sequences for responsiveness
 - Unconstrained by SNP locations
 - Empirical probe selection

How could we improve CN analysis on SNP arrays?

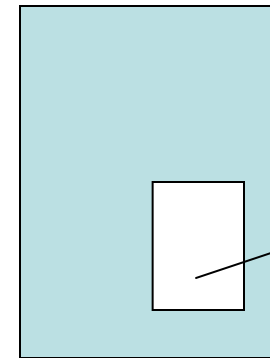
- Optimize probe locations in genome (within constraints of platform)
 - More-uniform coverage
 - Dense coverage in regions of interest (e.g. reported CNV)
- Optimize probe sequences for responsiveness
 - Unconstrained by SNP locations
 - Empirical probe selection

---> an “expression” probe for DNA

Probe reduction made it possible to explore this idea



500K chip set



Room for new stuff

5.0 and 6.0 arrays

Selection of copy number probes (6.0 array)

- Selected
 - 800 K probes across genome -> more-uniform coverage
 - 140 K probes in “target regions”
 - Reported CNV regions
 - Regions of interest to cancer research

Copy number data on 6.0 array

- 1.8 M point measurements across the genome
 - 900 K from SNP probe sets
 - 940 K from individual copy number probes

Outline

- Copy number variation
- Re-designing SNP arrays to interrogate copy number variation
- Analysis of copy number variation in whole genome association studies

Copy number analyses

	<u>Inherited CNPs</u>	<u>De novo CNVs</u>
Mechanism	Inheritance	<i>De novo</i> mutation
# transmitted per meiosis	~ 100s *	~ 0.01 * (may be more in cases)
Distribution	Shared by many people	Unique or rare
Size	Vast majority < 100 kb	Up to several Mb

(* estimate limited by current methods for ascertainment)



	<u>Inherited CNPs</u>	<u>De novo CNVs</u>
Mechanism	Inheritance	<i>De novo</i> mutation
# transmitted per meiosis	~ 100s *	~ 0.01 * (may be more in cases)
Distribution	Shared by many people	Unique or rare
Size	Vast majority < 100 kb	Up to several Mb

(* estimate limited by current methods for ascertainment)

Finding copy number variants *de novo*

- In each sample or sample-pair, find genomic regions across which a series of probes indicates reduced (or increased) intensity
- Evidence for CNV must be strong (genome-wide significance)
- Various existing approaches
 - HMMs
 - Binary segmentation
 - Smith-Waterman

Birdseye (Josh Korn)

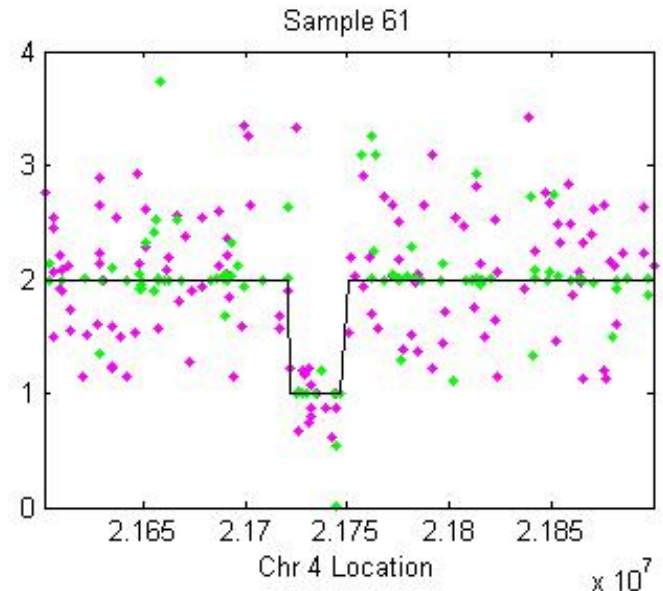
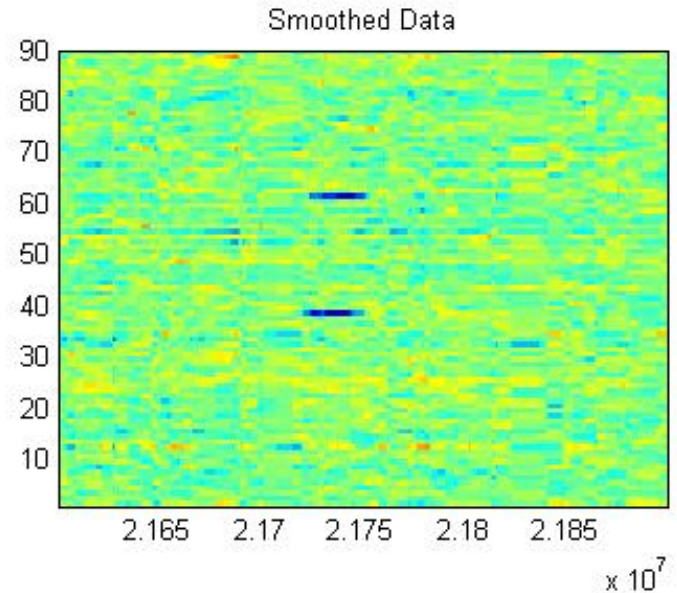
Sophisticated HMM

Uses CN and SNP probes
together

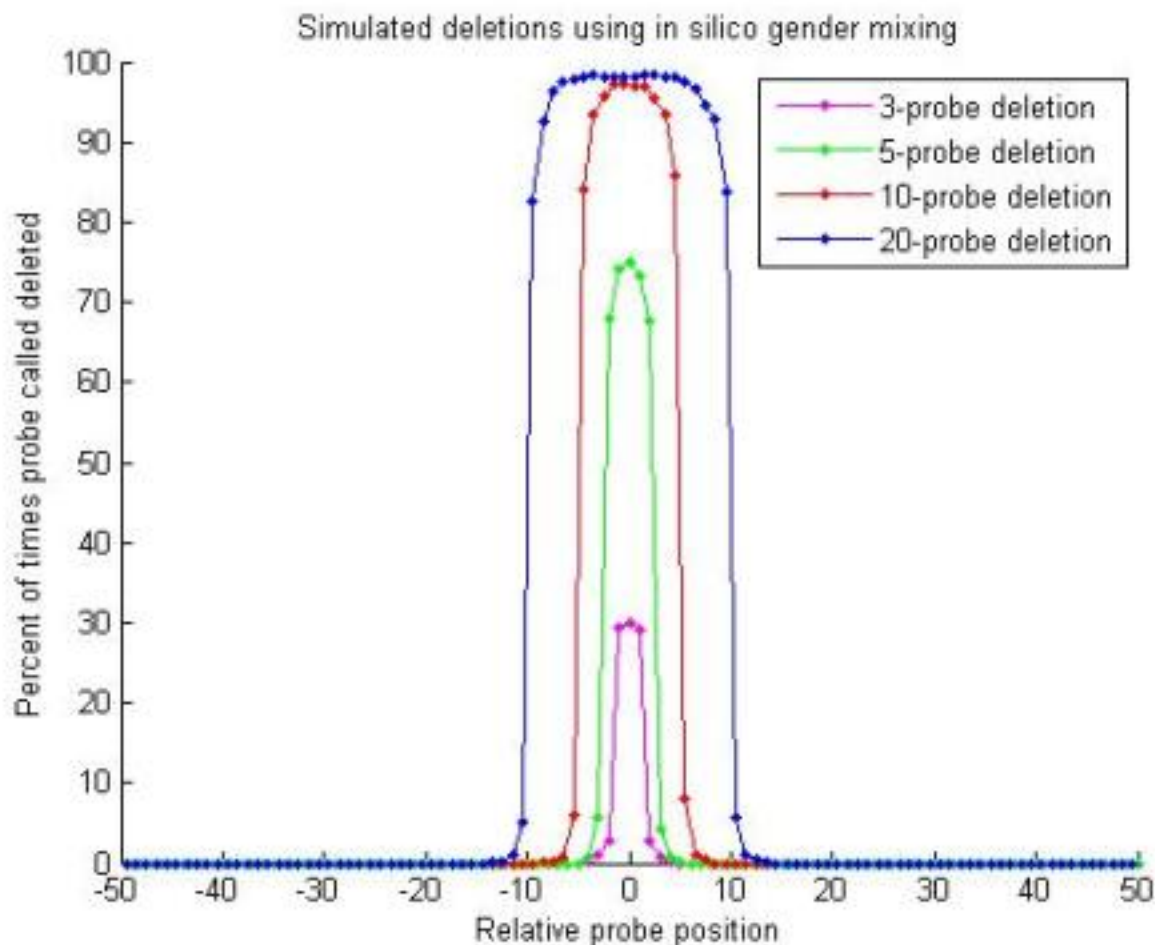
Models the response properties
of each probe

Discovers variants *de novo*

QuickTime™ and a
TIFF (Uncompressed) decompressor
are needed to see this picture.



Performance on in silico gender mixing experiment



Potential association analyses for rare and *de novo* CNVs

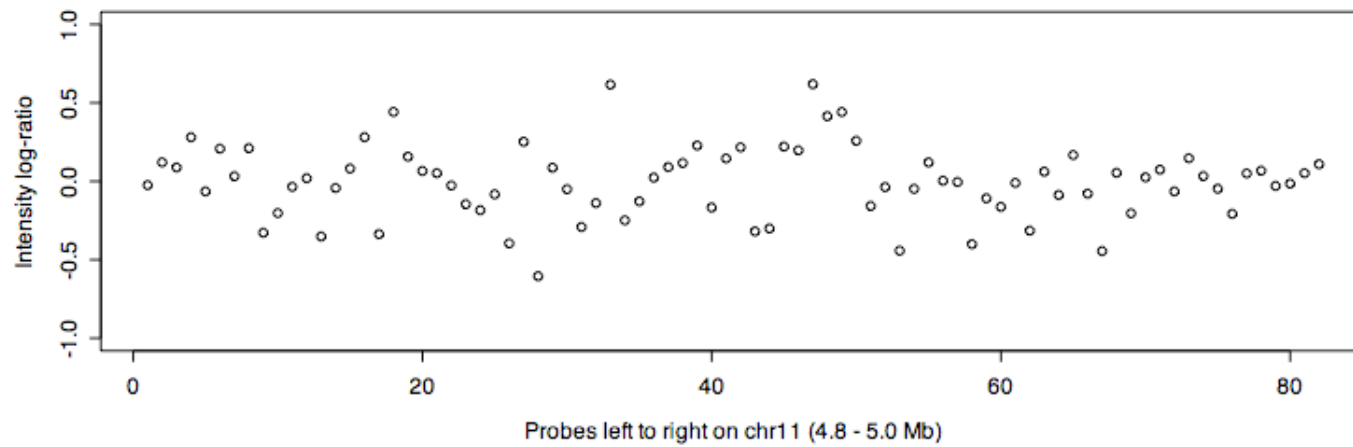
- Pileup at a particular genomic location
 - IMR, Sharp *et al.*, Nat Genet 2006
- Other ways of summing CNVs
 - Across entire genome (sporadic autism, Sebat *et al.*, Science 2006)
 - Critical to define criteria in advance
 - *Post hoc* criteria invented to fit observed data are unlikely to replicate in independent study
 - Precedent in analyses that group rare sequence variants (e.g. Cohen *et al.*, Science 2004, NEJM 2006)
 - Discussion in McCarroll and Altshuler, Nat Genet 2007



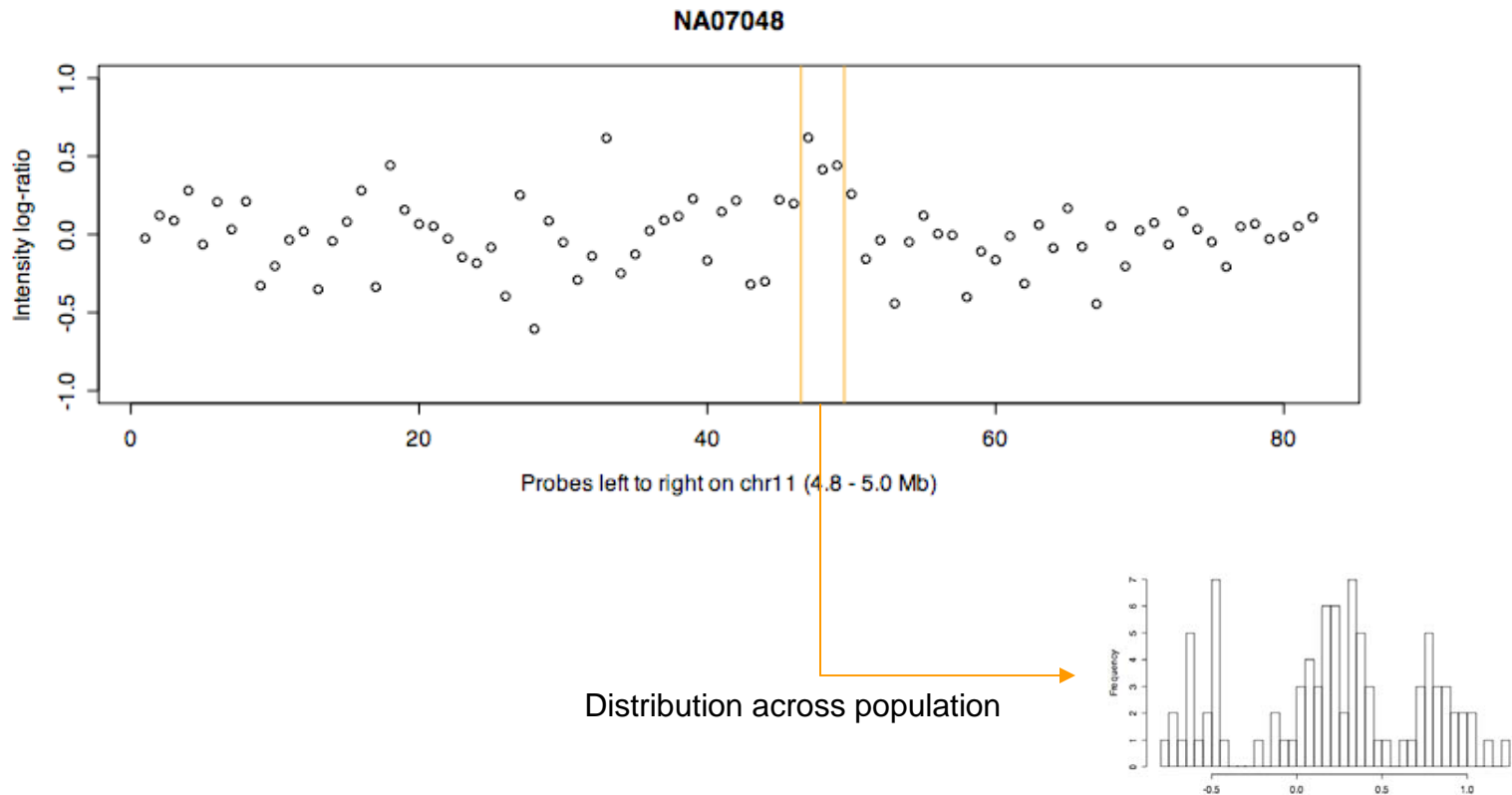
	<u>Inherited CNPs</u>	<u>De novo CNVs</u>
Mechanism	Inheritance	<i>De novo</i> mutation
# transmitted per meiosis	~ 100s *	~ 0.01 * (may be more in cases)
Distribution	Shared by many people	Unique or rare
Size	Vast majority < 100 kb	Up to several Mb

(* estimate limited by current methods for ascertainment)

NA07048



CNP analysis using prior information about CNP locations



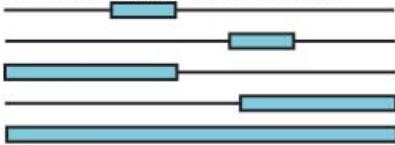
Need for high-resolution maps and genotyping technology

A

BAC probe on which variation has been observed:



Potential variants consistent with the observation:

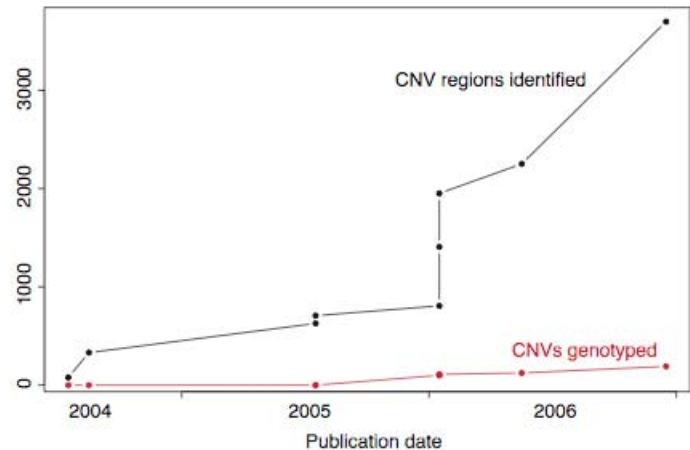


B

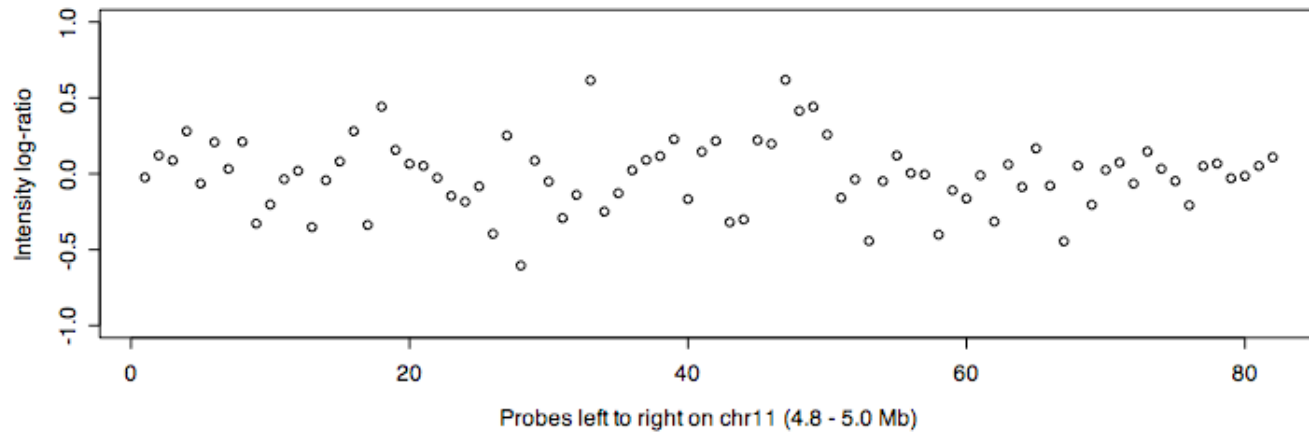
Fosmid end pair sequences from which deletion has been inferred:



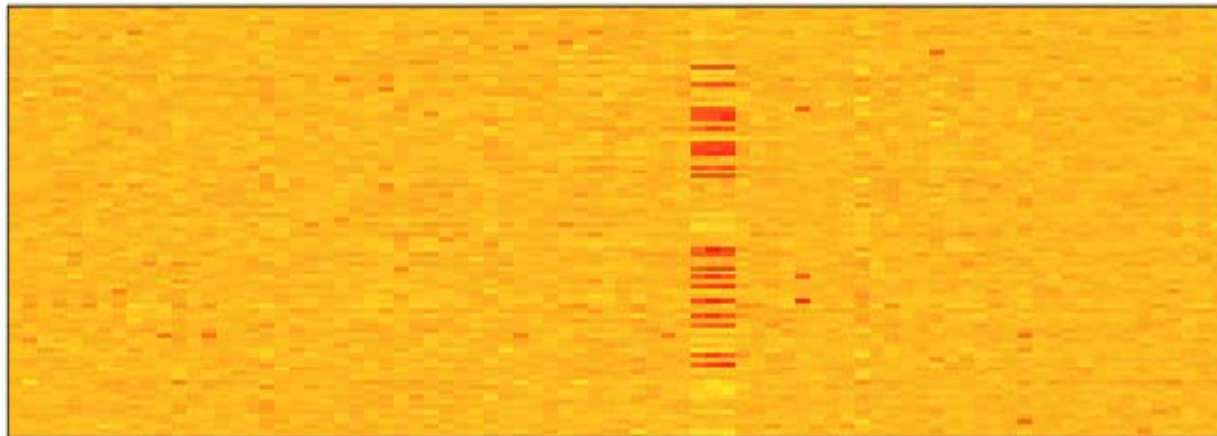
Potential deletions consistent with the observation:



NA07048



Across
90 samples

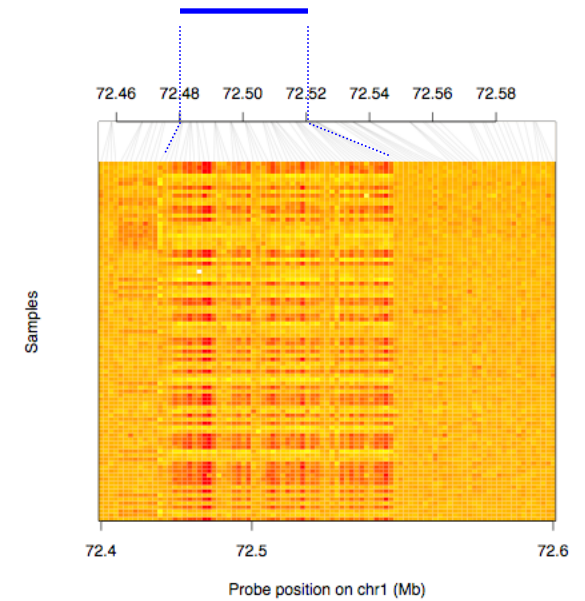
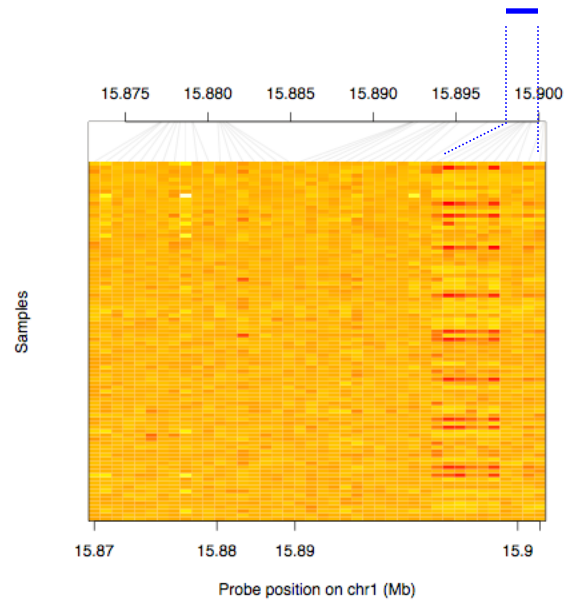
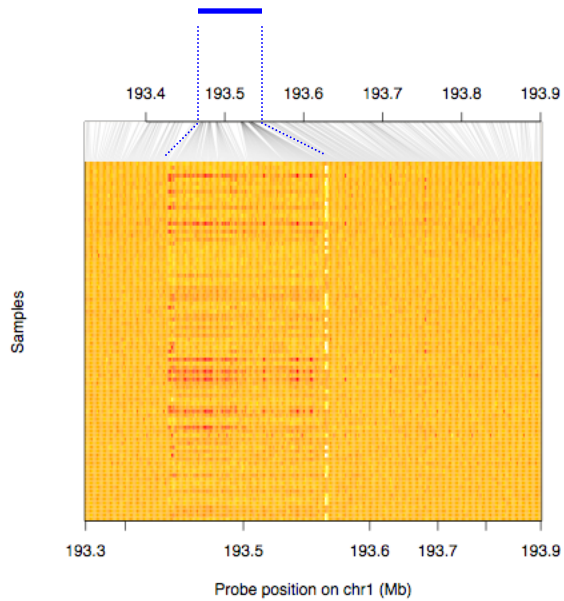


Probes left to right on chr11 (4.8 - 5.0 Mb)

High-resolution map of CNPs

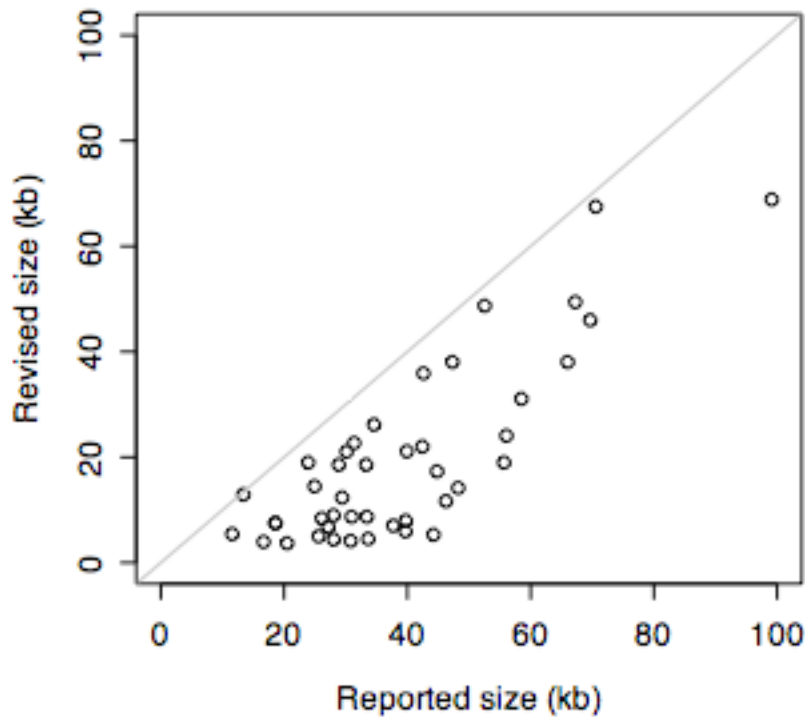
Reported CVNR:

Revised CVNR:

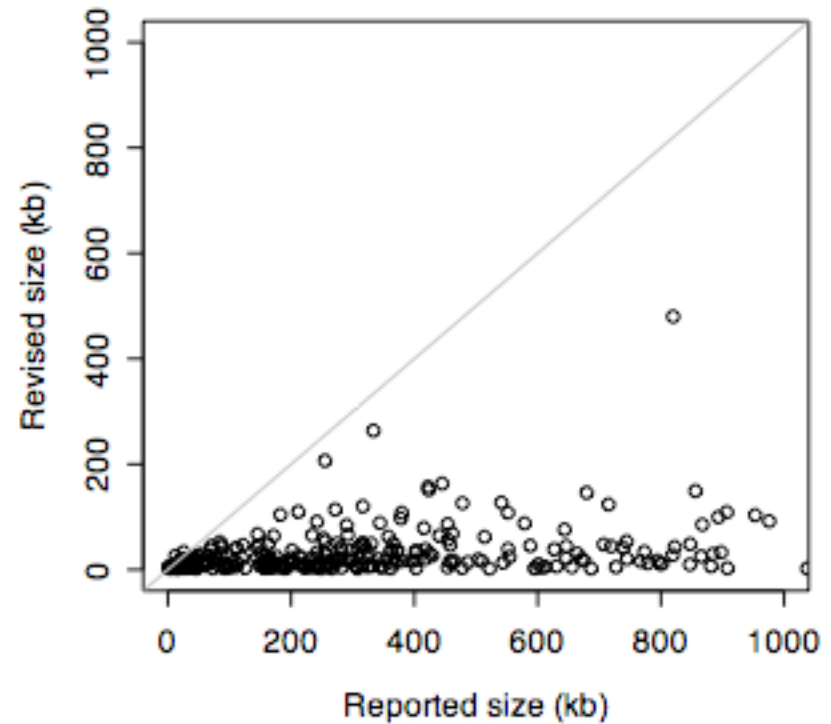


Most CNVs smaller than reported

Deletions (Tuzun et al.)



CNPs (Redon et al.)

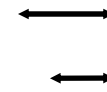
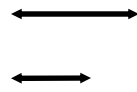


Compared to earlier survey, more CNVs, but less of genome affected

- Recent survey (Redon *et al.*, 2006)
 - 70 loci at different copy number levels between any two people
 - Mean size 228 kb
 - Cover 16 MB of genome
- This survey
 - 250 loci at different copy number levels between any two people
 - Mean size 20 kb
 - Cover 5 MB of genome

Revising gene content of CNVs

Earlier reports
of CNVs



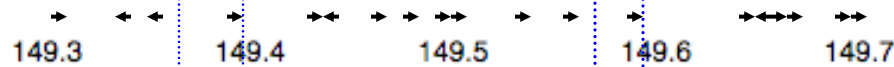
Tuzun *et al.*, 2005

McCarroll *et al.*, 2006



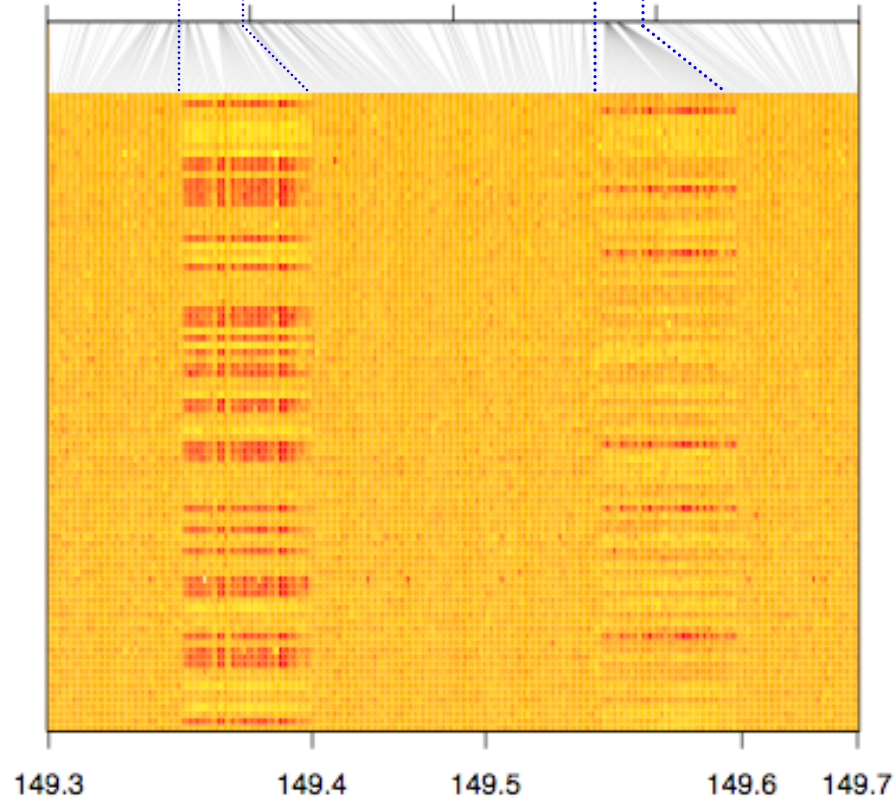
Redon *et al.*, 2006

Genes



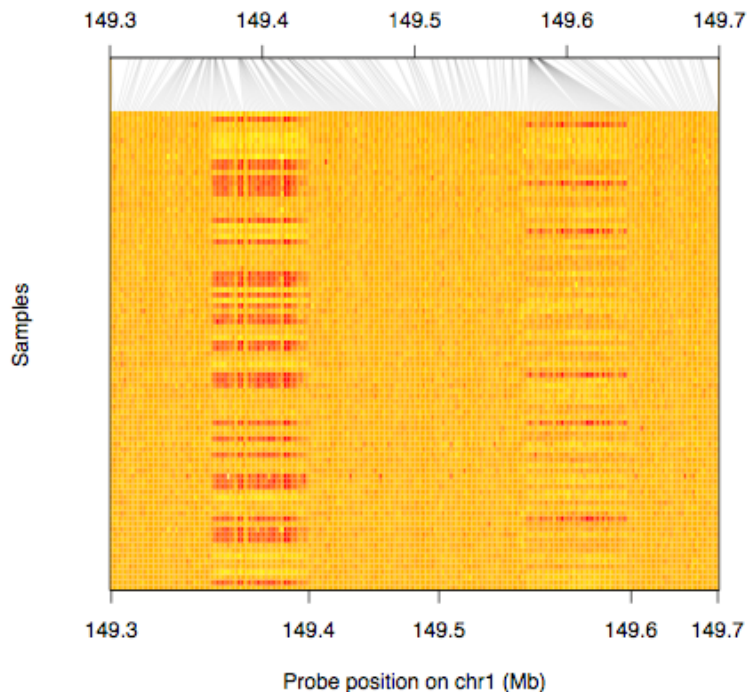
6.0 Array

Samples



Probe position on chr1 (Mb)

Most CNPs aren't directly interrogated by SNP assays (or by earlier platforms)

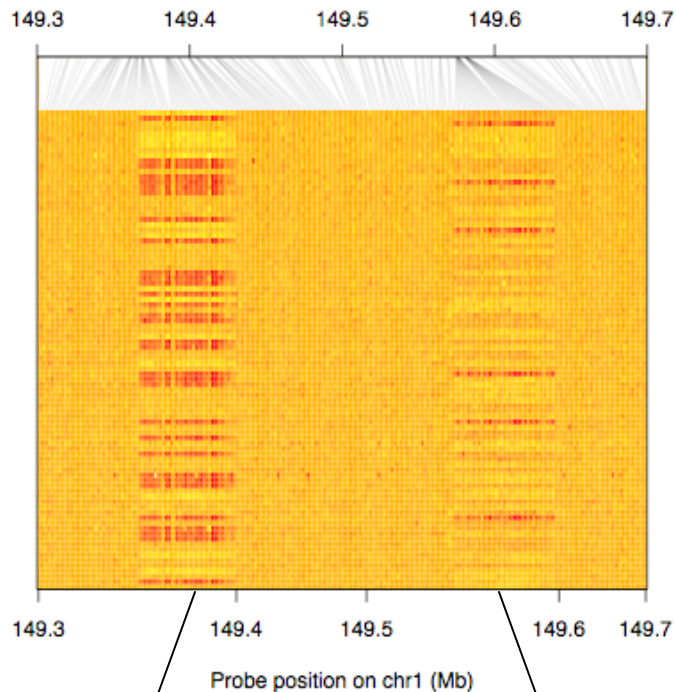


Affymetrix 500K
Illumina 300K

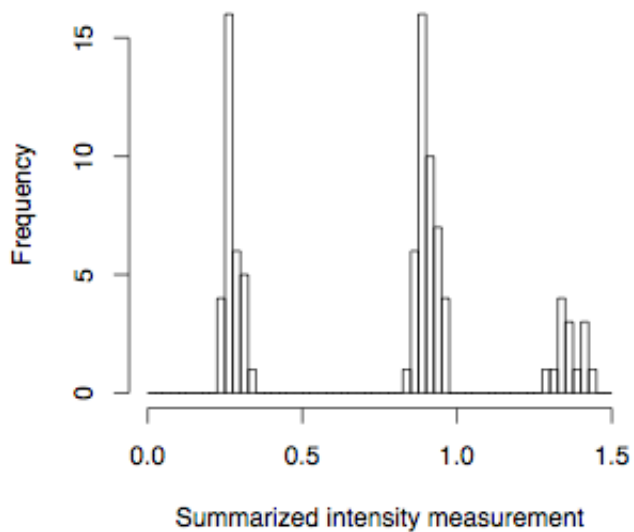
<u>SNP assays in "CNV region" (DGV)</u>	<u>SNP assays in actual CNVs within region</u>
110	0
24	0

From CNVs to *genotypes*

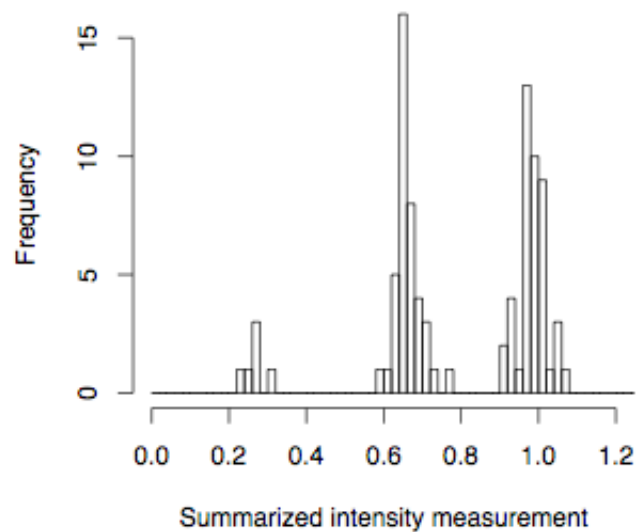
Samples



CNP 66
chr1, 149.369–149.399 Mb

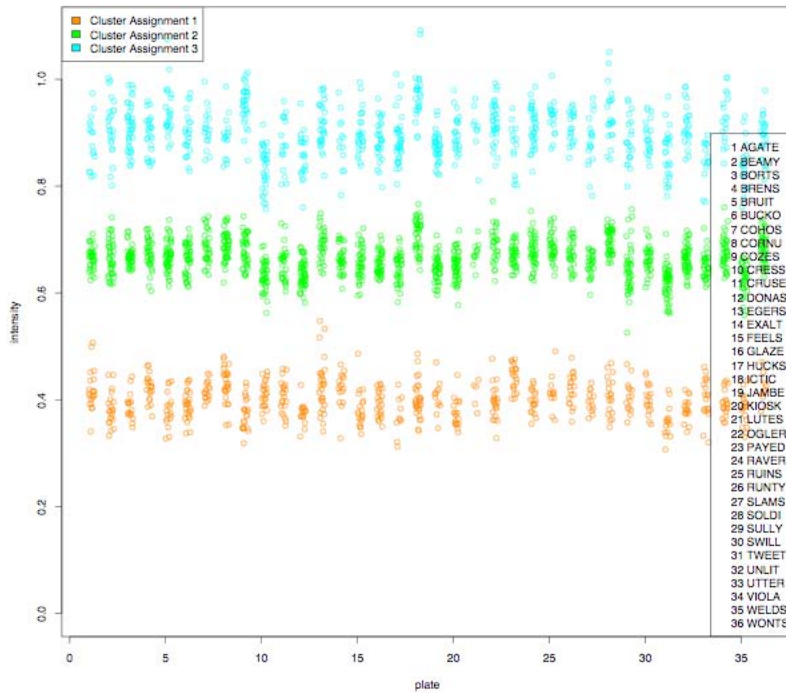


CNP 67
chr1, 149.574–149.582 Mb

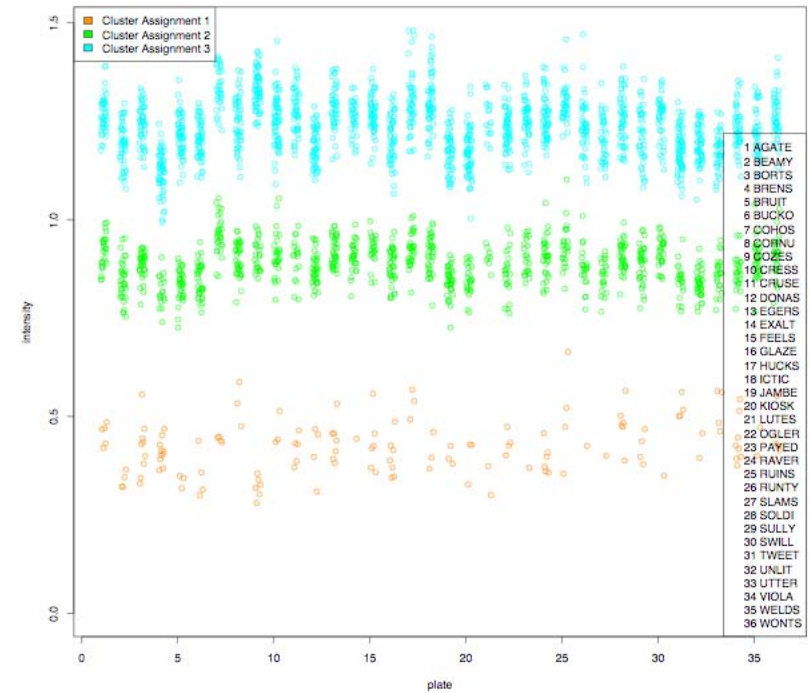


CNP genotyping in large studies

CNP 992 #probes [6]
chr [6] pos [103.859]
Uncertainty 95% [7.7e-06] 90% [3.2e-07]



CNP 1324 #probes [3]
chr [8] pos [72.378]
Uncertainty 95% [1.4e-05] 90% [1.0e-06]



Evaluating the accuracy of CNP genotypes by inheritance

Minor allele frequency	Population	# of CNVs	Mendel failures per trio
1% - 5%	CEU	136	0.008
	YRI	209	0.004
5% - 15%	CEU	219	0.007
	YRI	364	0.008
15% - 50%	CEU	340	0.008
	YRI	387	0.008

- Accomplished in CANARY
 - Genotype samples for several hundred common, inherited CNPs

How much does copy number variation influence clinical phenotypes?

- A major outstanding question in the field
- Enabled by new array technology, will soon be possible to learn



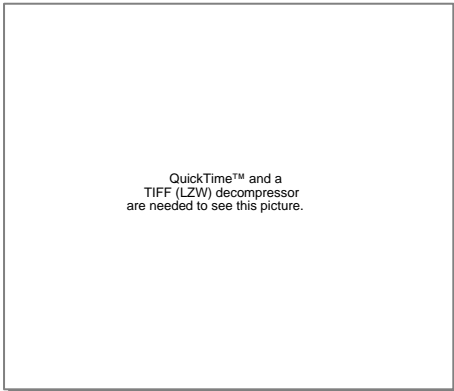
Finny Kuruville



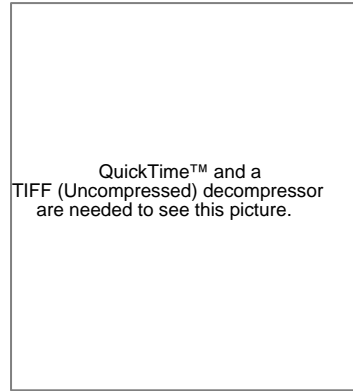
Josh Korn

Alec Wysoker
Jim Nemesh
Paul de Bakker
Casey Gates
Marcia Nizarri

Simon Cawley
Steve Lincoln
Keith Jones



Stacey Gabriel



Mark Daly

QuickTime™ and a
TIFF (Uncompressed) decompressor
are needed to see this picture.

David Altshuler

Acknowledgements

Affymetrix

Simon Cawley

Steve Lincoln

Keith Jones

Earl Hubbell

Teresa Webster

Sean Walsh

Rui Mei

Xiaojun Di

Geoff Yang

Hajime Matsuzaki

Guoying Liu

Alan Williams

Harley Gorrell

Chuck Sugnet

Fan Shen

Michael Shapero



Finny Kuruville



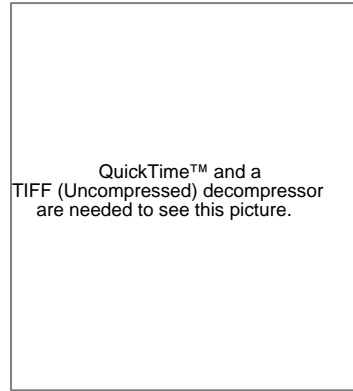
Josh Korn

Alec Wysoker
Jim Nemesh
Paul de Bakker
Casey Gates
Marcia Nizarri

Simon Cawley
Steve Lincoln
Keith Jones



Stacey Gabriel



Mark Daly

QuickTime™ and a
TIFF (Uncompressed) decompressor
are needed to see this picture.

David Altshuler