# Copy Number Variants:

# detection and analysis

Manuel Ferreira & Shaun Purcell

*Boulder, 2009*

# Large chromosomal rearrangements can cause sporadic disease

Down syndrome

Duchenne Muscular Dystrophy (DMD)

DiGeorge-Velo cardiofacial syndrome (VCFS)

…

Lupski 2007 Nat Genet 39: s43

# Large-Scale Copy Number Polymorphism in the Human Genome

Jonathan Sebat,[1] B. Lakshmi,[1] Jennifer Troge,[1] Joan Alexander,[1] Janet Young,[2] Pär Lundin,[3] Susanne Månér,[3] Hillary Massa,[2] Megan Walker,[2] Maoyen Chi,[1] Nicholas Navin,[1] Robert Lucito,[1] John Healy,[1] James Hicks,[1] Kenny Ye,[4] Andrew Reiner,[1] T. Conrad Gilliam,[5] Barbara Trask,[2] Nick Patterson,[6] Anders Zetterberg,[3] Michael Wigler[1]*

...which large duplications and deletions contribute to human genetic ...versity is unknown. Here, we show that large-scale copy number ...(CNPs) (about 100 kilobases and greater) contribute substantially ...tion between normal humans. Representational oligonucleotide ...sis of 20 individuals revealed a total of 221 copy number differ- ...ng 76 unique CNPs. On average, individuals differed by 11 CNPs, ...length of a CNP interval was 465 kilobases. We observed copy ...on of 70 different genes within CNP intervals, including genes ...urological function, regulation of cell growth, regulation of metab- ...veral genes known to be associated with disease.

*Sebat et al 2004 Science 305: 525*

# Detection of large-scale variation in the human genome

A John Iafrate[1,2], Lars Feuk[3], Miguel N Rivera[1,2], Marc L Listewnik[1], Patricia K Donahoe[2,4], Ying Qi[3], Stephen W Scherer[3,5] & Charles Lee[1,2,5]

**We identified 255 loci across the human genome that contain genomic imbalances among unrelated individuals. Twenty-four variants are present in >10% of the individuals that we examined. Half of these regions overlap with genes, and many coincide with segmental duplications or gaps in the human genome assembly. This previously unappreciated heterogeneity may underlie certain human phenotypic variation and susceptibility to disease and argues for a more dynamic human genome structure.**

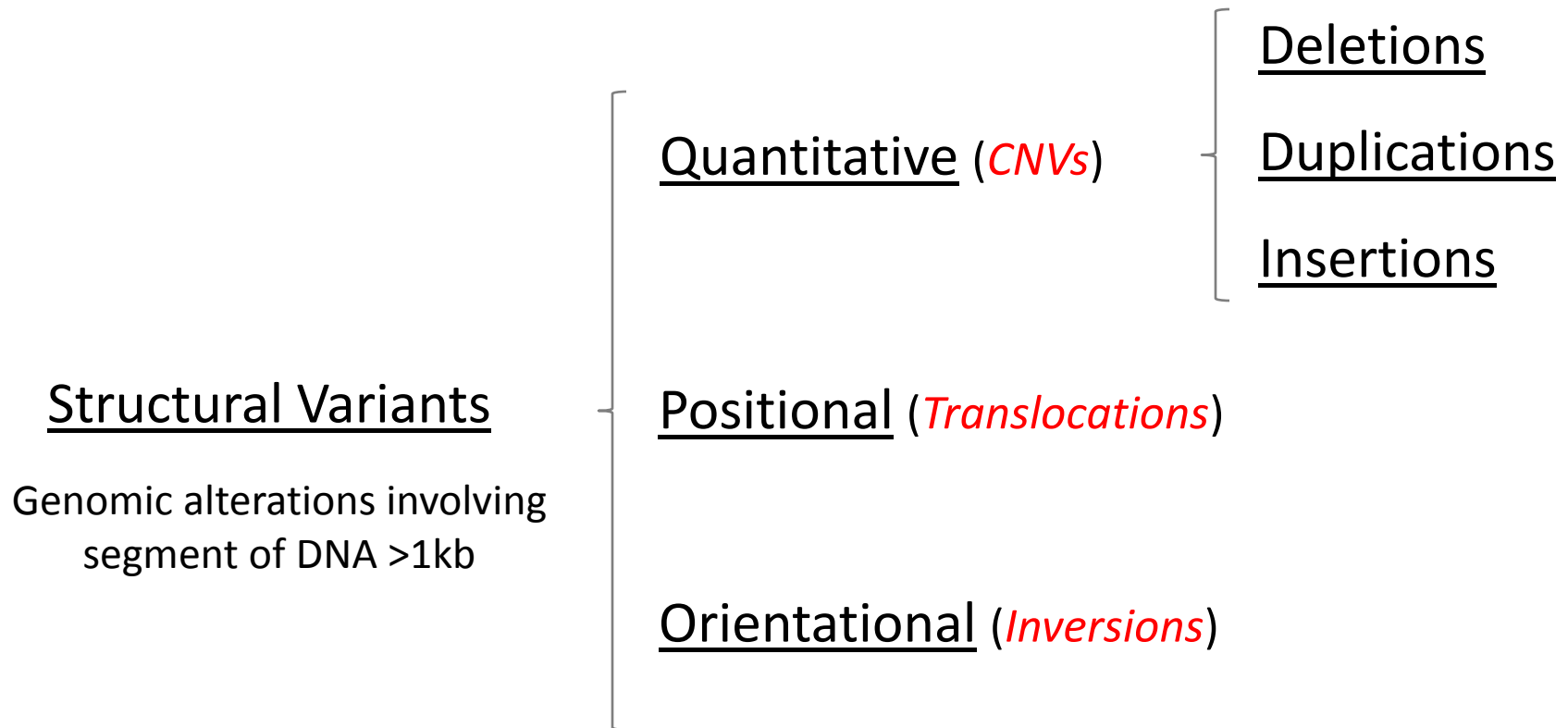*Iafrate et al 2004 Nat Genet 36: 949*

# Outline

1. What is a Copy Number Variant (CNV)

2. Genome-wide detection of CNVs

3. Association analysis of CNVs

4. Online databases

# 1. What is a CNV?

# What is a CNV?

## 1. Classes of structural variants



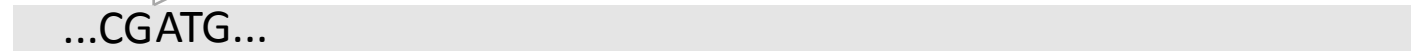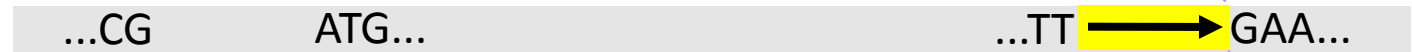Copy Number Polymorphism (CNP) is a *CNV* that occurs in >1% population
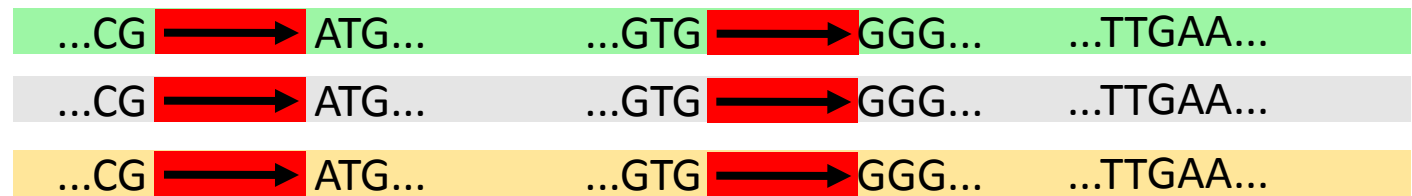
Sequence variation
Single nucleotide
- Base change – substitution – point mutation
→ Insertion-deletions ("indels")
- SNPs – tagSNPs

Structural variation

2 bp to 1,000 bp
- Microsatellites, minisatellites
→ Indels
- Inversions
- Di-, tri-, tetranucleotide repeats
- VNTRs

1 kb to submicroscopic
→ Copy number variants (CNVs)
→ Segmental duplications
- Inversions, translocations
→ CNV regions (CNVRs)
- Microdeletions, microduplications

Microscopic to subchromosomal
→ Segmental aneusomy
- Chromosomal deletions – losses
- Chromosomal insertions – gains
- Chromosomal inversions
- Intrachromosomal translocations
- Chromosomal abnormality
→ Heteromorphisms
- Fragile sites

Whole chromosomal to whole genome
- Interchromosomal translocations
- Ring chromosomes, isochromosomes
- Marker chromosomes
→ Aneuploidy
→ Aneusomy

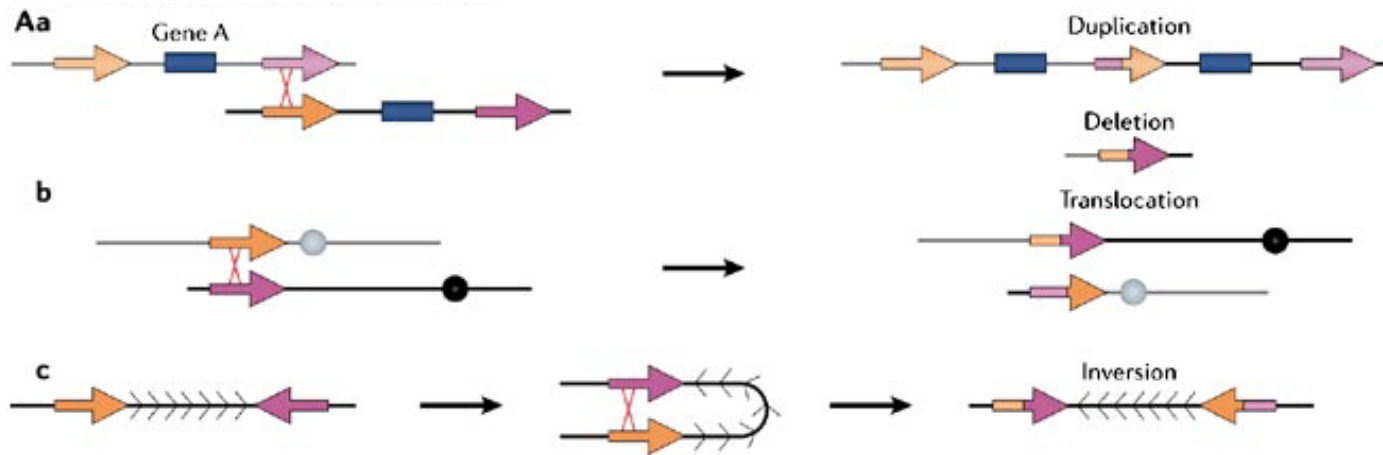→ Term defined or discussed in Box 1

Molecular genetic detection

Cytogenetic detection

Scherer 2007 Nat Genet 39: s7

# What is a CNV?

## 2. Origins of CNVs

*(A) Non-allelic homologous recombination*



*(B) Non-homologous end joining*
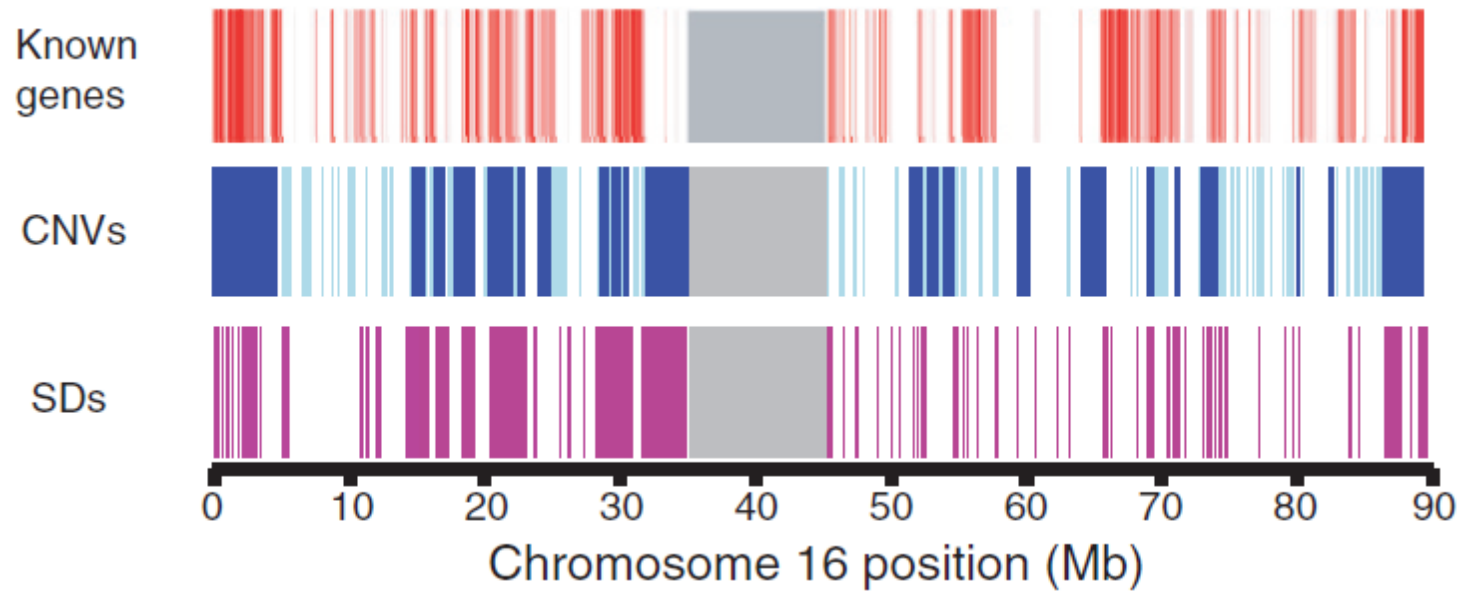
*(C) Tandem repeat sequences*

*(D) Retrotransposons*

# What is a CNV?

## 3. CNVs are abundant in the genome

| Human vs Human | SNPs | CNVs |
|---|---|---|
| Base pairs | 2.5 Mb | 4 Mb |
| | 1/1,200 bp | 1/800 |
| % genome | 0.08% | 0.12% |

# What is a CNV?

## 4. CNVs significantly overlap with known genes



Cooper et al 2007 Nat Genet 39: s22

# What is a CNV?

## 5. CNVs influence gene expression

83.6%  17.7%



Stranger et al 2007 Science 315: 848

# What is a CNV?

**6. In healthy individuals, most CNVs are inherited…**

>99%  <u>inherited</u>

Rare CNVs   10%

Common CNVs   90%
>1% population

<1%  <u>de novo</u>

McCarroll 2008 Hum Mol Genet 17: R135

McCarroll et al 2008 Nat Genet 40: 1166

# 2. Detection of CNVs

# Detection of CNVs

**A. Using intensity data from whole-genome arrays**

SNPs $\longrightarrow$ *Genotype known common variants*

CNVs

(A)  Genotype *known common* variants

(B)  *Identify* and *genotype new*, potentially rarer variants

# Detection of CNVs

*(A) Genotype known common CNVs using whole-genome arrays*

Nimblegen

array-CGH, CNV only, test vs reference

custom or whole-genome (up to 2,1M probes)

Affy 6.0

>940,000 CNV non-polymorphic probes

High-density in ~5,600 CNV regions in DGV +
extended to whole-genome

Illumina 1M

36,000 CNV non-polymorphic probes

covering ~4,000 CNV regions in DGV

# Detection of CNVs

S

...CG ⟶ ATG...

| Ind | Genotype Mat/Pat | Copy number at S | | Amount of DNA at S |
|---|---|---|---|---|
| 1 | S/S | 2 | ...CG ⟶ ATG...<br>...CG ⟶ ATG... | ⬆⬆ |

# Detection of CNVs



*Non-polymorphic probes*

# Detection of CNVs

(B) *Identify and genotype new, potentially rarer CNVs from whole-genome array data (CGH, Affymetrix/Illumina)*
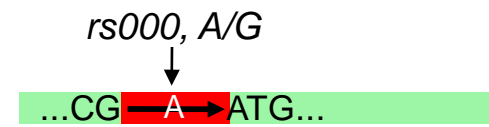
*Example:* rs1006737 **A**/**G**

        **AA**
        **AG**
        **GG**

... AGCCCGAA**A**TGTTTTCAGA...    *probe 1*

... AGCCCGAA**G**TGTTTTCAGA...    *probe 2*

# Detection of CNVs

*rs000, A/G*

...CG A ATG...

| Ind | Genotype *Mat/Pat* | Pattern | Copy number for: | | |
|---|---|---|---|---|---|
| | | | A | G | Total |
| 1 | A/G | ...CG A ATG... <br> ...CG G ATG... | 1 | 1 | 2 |

# Detection of CNVs

...CG ——A→ ATG...



**Individuals with duplication(s)**

*ie. total CN > 2*

*Normalized intensity of allele G*

*Normalized intensity of allele A*

**Individuals with deletion(s)**

*ie. total CN < 2*

*Polymorphic probe in CNV region*

# Detection of CNVs



**Birdseye**
**Affy 5.0, 6.0**
Korn et al 2008 Nat Genet 40: 1253

**PennCNV**
**Affymetrix and Illumina**
Wang et al 2007 Genome Res 17: 1665

*Combine information across probes to identify new CNVs*

| For example... | Cases | Controls |
|---|---|---|
| 100kb deletion chr. 2 | 10/5,000 | 1/5,000 |

## Detecting CNVs through GWAS arrays is challenging…

Lots of confounders: DNA quality, concentration, source, batch effects, date effects.

Arrays have poor resolution for CNVs (>100kb).

Genotype calling is computationally demanding, as it requires analysis of very large 'raw' cell files.

Genotype calling software often platform specific, not very user friendly.
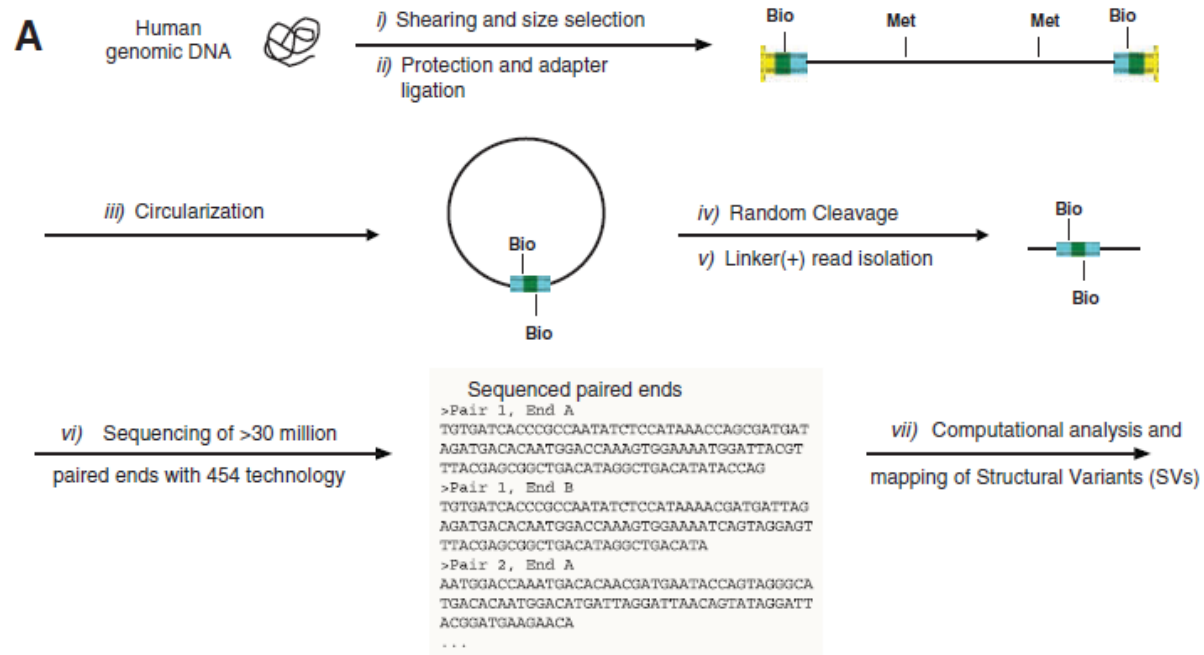
# Detection of CNVs

## B. Identifying CNVs through genotyping errors

▷ <u>Mendelian Inconsistencies</u>

▷ <u>Failure Hardy-Weinberg equilibrium</u>

Conrad et al 2006 Nat Genet 38: 75

McCarroll et al 2006 Nat Genet 38: 86

# Detection of CNVs

## C. Targeted or whole-genome sequencing



Korbel et al 2007 Science 318: 420

# Summary so far…

CNVs are abundant, often overlap genes, can influence gene expression and most are inherited in healthy individuals

Known and new CNVs can be identified and genotyped in large-scale studies using whole-genome genotyping arrays, such as the 6.0 and 1M. Low resolution (>100Kb) & low signal/noise ratio.

More accurate CNV genotyping maps/arrays/algorithms expected in the next few years.

What are the particular strategies and challenges for association analysis of CNVs?

# 3. Association analysis of CNVs

# Association analysis of CNVs

## 1. Some of the relevant questions

*(A) Are CNPs associated with variation in human traits or diseases?*

*(B) Can we identify rare CNVs associated with large increase in disease risk? Are these de novo or inherited in cases?*

*(C) When considering the whole-genome, do cases have more CNV events then controls, ie. increased burden?*

*(D) How to test SNPs in copy number regions?*

*(E) Are most CNVs tagged by SNPs in genotyping arrays?*

# Example 1: Autism whole-genome CNV analysis

| Sample | 16p11 | Cases | Controls | *P* |
|--------|-------|-------|----------|-----|
| Discovery | Del (600kb) | 5/1,441 | 3/4,234 | 1.1 x 10$^{-4}$ |
| [Affy 500K] | Dup | 7/1,441 | 2/4,234 | |
| Replication 1 (CHB) | Del | 5/512 | 0/434 | 0.007 |
| [array-CGH] | Dup | 4/512 | 0/434 | |
| Replication 2 (deCODE) | Del | 3/299 | 2/18,834 | 4.2 x 10$^{-4}$ |
| [Illumina] | Dup | 0/299 | 5/18,834 | |

COPPER
Birdseye
CNAT

Deletion frequency Iceland

| | | | | del | dup |
|--|--|--|--|-----|-----|
| Autism | 1% | | inherited | 2 | 6 |
| Psychiatric disorder | 0.1% | | de novo | 10 | 1 |
| General population | 0.01% | | unknown | 1 | 4 |

Weiss et al. N Engl J Med 2008; 358: 667

# Example 2: SCZ whole-genome CNV analysis



doi:10.1038/nature07239

nature

LETTERS

**Rare chromosomal deletions and duplications increase risk of schizophrenia**

The International Schizophrenia Consortium*

*Specific loci*

**Cases**

**Controls**

**Chromosome →**

# Specific large (>500kb) rare deletions

# Genome-wide burden of rare CNVs in SCZ

**3,391 patients with SCZ, 3,181 controls**
*Filter for <1% MAF, >100kb*
**6,753 CNVs**

**Cases have greater rate of CNVs than controls**
*1.15-fold increase*
$P = 3 \times 10^{-5}$

**Cases have more genes intersected by CNVs than controls**
*1.14-fold increase*
$= 2 \times 10^{-6}$

**True for *singleton* events (observed only once in dataset)**
*1.45-fold increase*
*(~15% cases versus 11% controls)*
$P = 5 \times 10^{-6}$

**Rate of *genic* CNVs in cases v**
*1.18-fold increase*
$P = 5 \times 10^{-6}$

**CNVs in cases versus controls**
*1.09-fold increase*
$P = 0.16$

**Results invariant to obvious statistical controls**
*Array type, genotyping plate, sample collection site, mean probe intensity*

nature

LETTE...

# Large recurrent microdeletions associated with schizophrenia

Hreinn Stefansson[1]*, Dan Rujescu[2]*, Sven Cichon[3,4]*, Olli P. H. Pietiläinen[5], Andres Ingason[1], Stacy Ste... Ragnheidur Fossdal[1], Engilbert Sigurdsson[6], Thordur Sigmundsson[6], Jacobine E. Buizer-Voskamp[7],
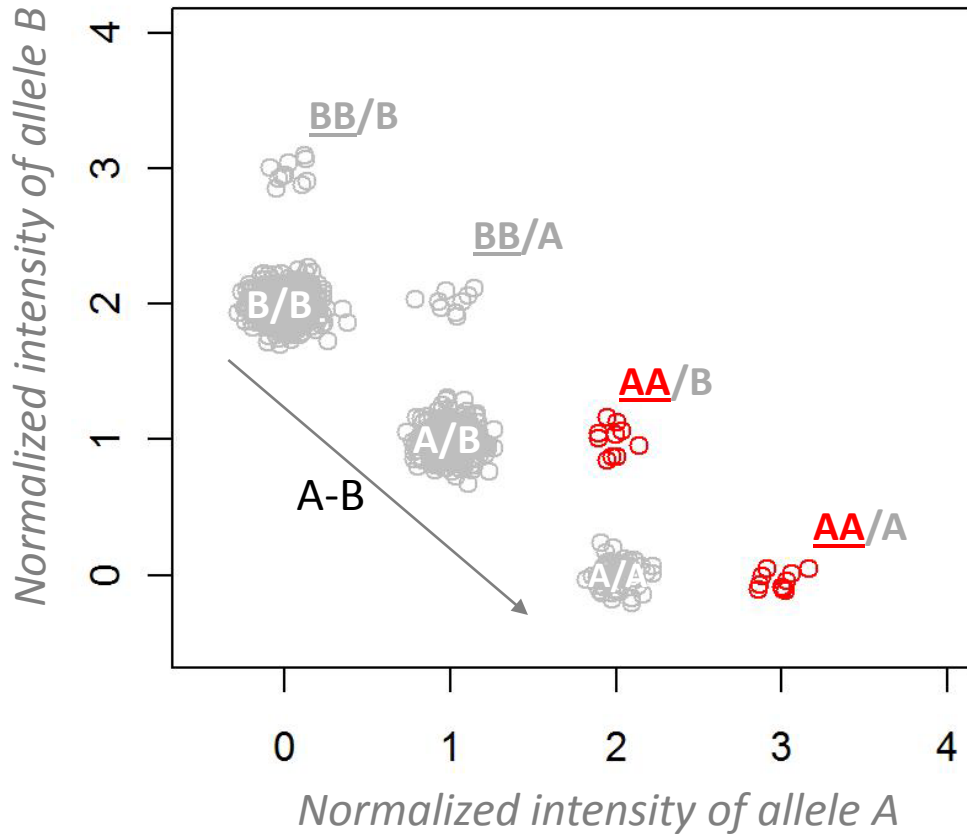


1q21.1 and 15q13.3 also identified
by SGENE consortium

- Two other studies supporting a genome-wide increase in rare CNVs in schizophrenia
  - Walsh et al (2008) *Science*
    - 5% controls, 15% cases, 20% early onset cases
    - neurodevelopmental genes disrupted
  - Xu et al (2008) *Nature Genetics*
    - strong increased de novo rate in sporadic cases; but increased inherited rate also

# Association analysis of CNVs

## 2. Testing SNPs in CNV regions



| | | |
|---|---|---|
| Gene Z reference | **Healthy** | **Healthy** |
| Gene Z with deletion | **Disease** | **Healthy** |
| Gene Z with mutation | **Disease** | **Healthy** |
| Gene Z with deletion and mutation | **Disease** | **Disease** |
| **Individual analysis of SNPs or CNVs** | ✔ | ✘ |

$$y = \beta_1 \cdot \text{SNP} + \beta_2 \cdot \text{CNV}$$

$$\downarrow$$

$$y = \beta_1 \cdot (A - B) + \beta_2 \cdot (A + B)$$

*Allele-specific risk CNV*

Korn et al 2008 Nat Genet 40: 1253

# Association analysis of CNVs

## 3. Testing CNVs through the analysis of SNPs in LD



**Common CNVs**

Coverage limited by lack of SNPs in CNV regions
(poor genotyping)

McCarroll et al 2008 Nat Genet 40: 1166

# 4. Online databases

# Database of Genomic Variants

http://projects.tcag.ca/variation/



Comprehensive summary of structural variation in the
human genome. Healthy control samples

# DECIPHER

https://decipher.sanger.ac.uk/



Database of submicroscopic chromosomal imbalances, from array-CGH data. Focuses on data from patients with developmental delay, learning disabilities or congenital anomalies.