# Human Genetic Variation and Association Studies

Jacques Fellay, MD
SNF Professor

Global Health Institute
EPFL School of Life Sciences

Institute of Microbiology
University of Lausanne / CHUV

# We Are All Different

# We Are All Different



About 2% of people have two copies of APOE4 and are very likely to succumb to Alzheimer's disease
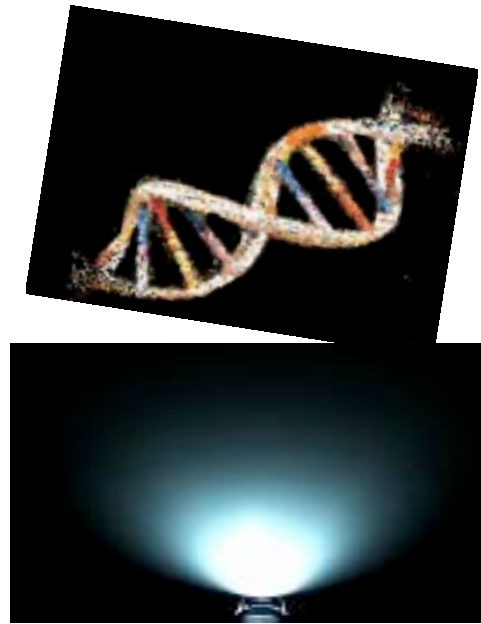
# We Are All Different



About 2% of people have two copies of APOE4 and are very likely to succumb to Alzheimer's disease

About 1% of us have two copies of a small deletion in CCR5 and are largely immune to infection by the HIV virus

# We Are All Different



About 2% of people have two copies of APOE4 and are very likely to succumb to Alzheimer's disease
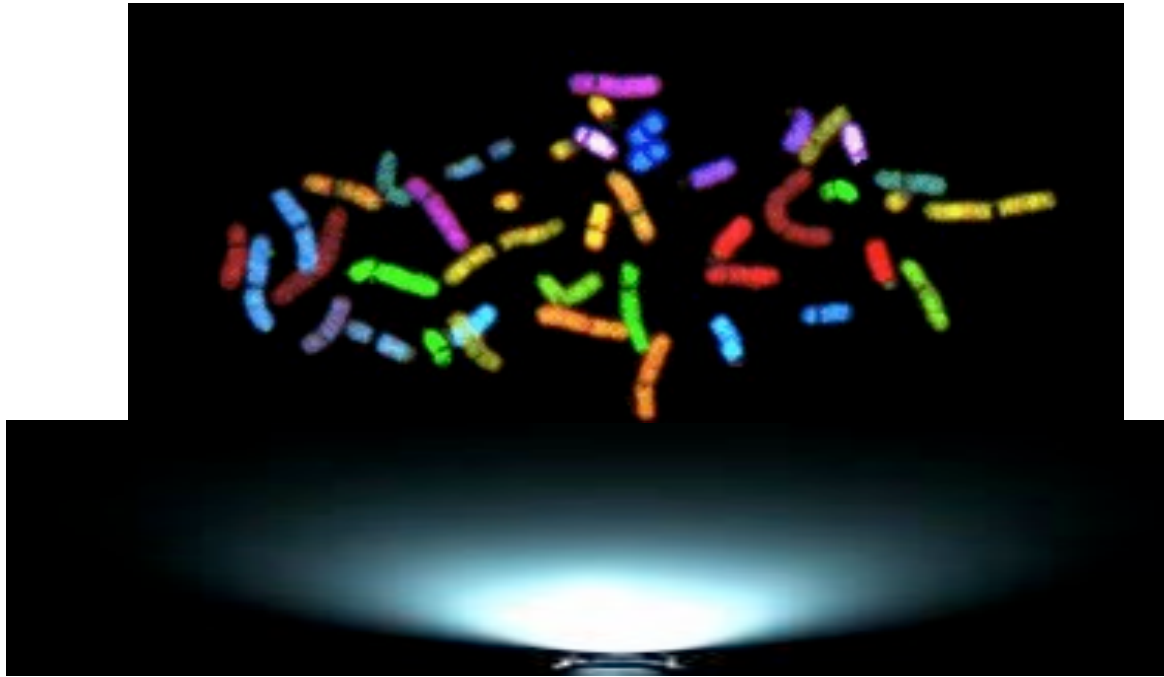
About 1% of us have two copies of a small deletion in CCR5 and are largely immune to infection by the HIV virus

And about 7% do not make any functional CYP2D6 enzyme and therefore codeine provides no pain relief

These examples come from looking at only the tiniest fraction of our genome

It is now possible to scan the *whole genome* to find the genetic determinants of key differences amongst people

# Overview

➢ Types of human genetic variation

➢ Mapping approaches

 ▪ GWAS

 ▪ Sequencing

➢ Real life examples

# Different forms of genetic variation

# Different forms of genetic variation

-single nucleotide variants
    -3-4 million per individual

-multiple nucleotides variants
    -greater content than single site changes
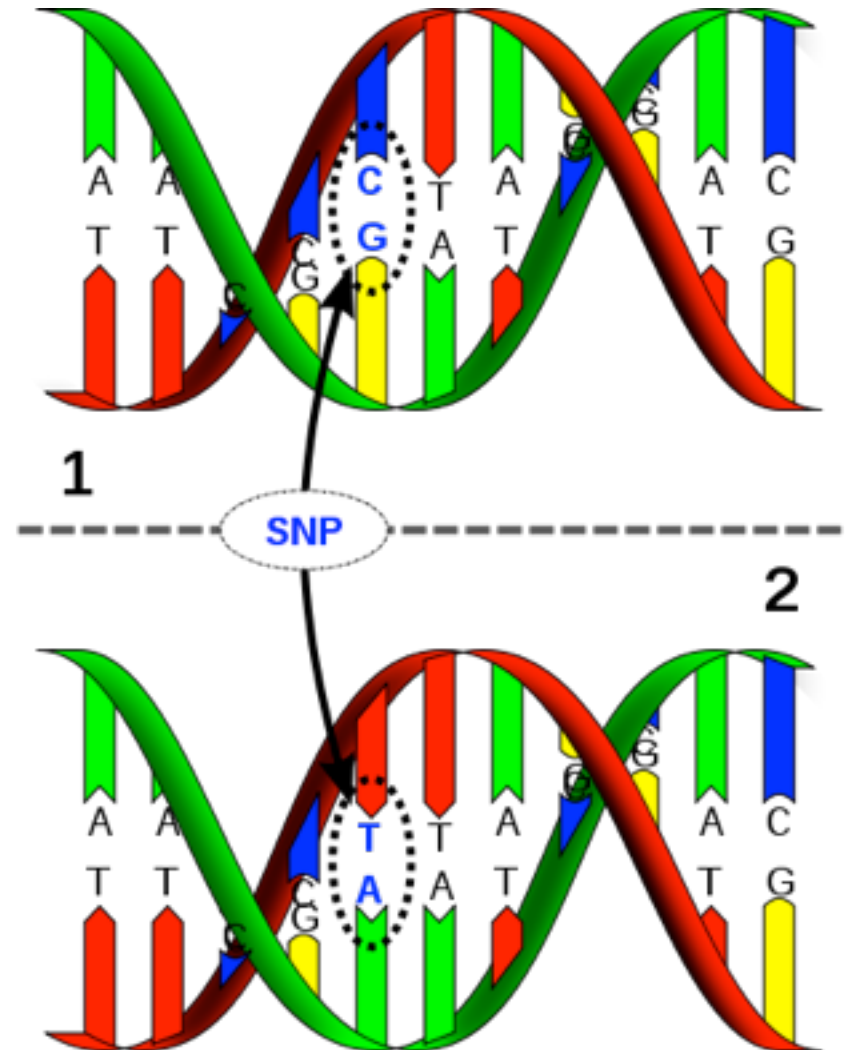
# Different forms of genetic variation

-single nucleotide polymorphisms (or SNPs)
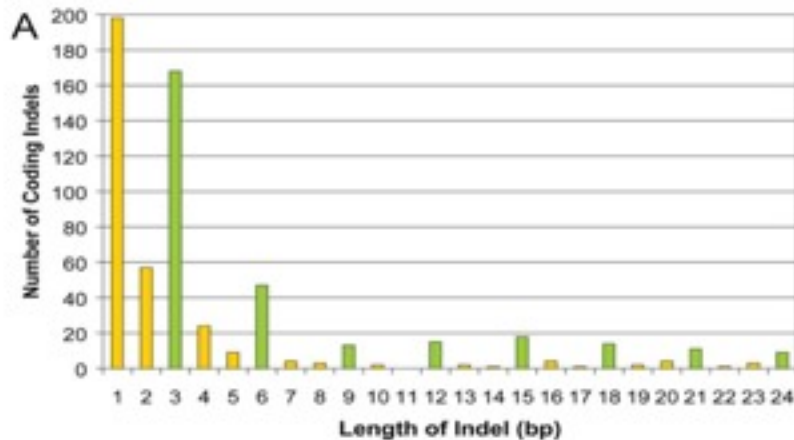
functional?
-missense
-non-sense
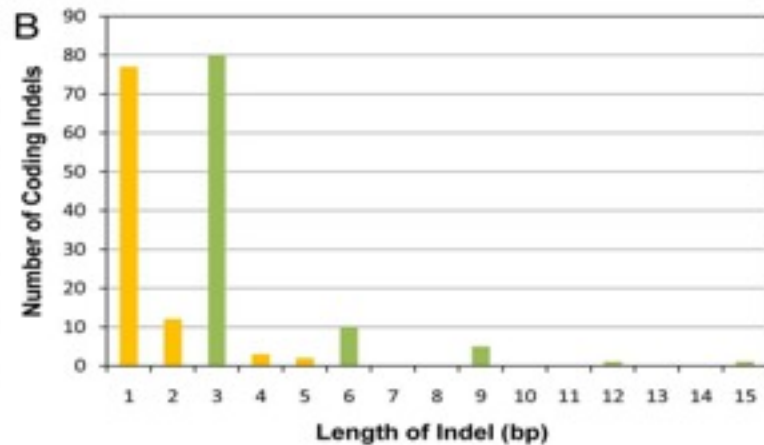-splice site

non-functional?
-silent
-intronic
-intergenic

# Different forms of genetic variation

-single nucleotide polymorphisms (or SNPs)
-small insertions or deletions (indels)
    -coding/non-coding



A — J. C. Venter's Genome (Figure from Ng et al. PLoS Genetics 4(8): e1000160)

B — This study (individual genome average)

# Different forms of genetic variation

-single nucleotide polymorphisms (or SNPs)
-small insertions or deletions (indels)
-short tandem repeats/microsatellites
      -repeat of 2, 3, 4 or more nucleotides
      -10-100x
      -highly polymorphic
      -error during replication (slippage)

# Trinucleotide repeat diseases

14 known diseases

   -9 due to glutamine repeats (CAG trinucleotide)
   -neurodegenerative disease (polyglutamine disease)
   -neuronal decay
   -spinocerebellar ataxias and Huntington's disease

# Trinucleotide repeat diseases

14 known diseases
>   -9 due to CAG trinucleotide=Glutamine
>   -neurodegenerative disease (polyglutamine disease)
>   -neuronal decay
>   -Spinocerebellar ataxias and Huntington's disease

Huntington's disease trinucleotide repeats
>   -tract of <28                =normal
>   -tract of 28 to 35           =intermediate
>   -tract of 36-40              =reduced penetrance/affected
>   -tract of >40                =full penetrance/affected

anticipation – tract expands with successive generations leading to earlier age of onset and more severe disease.

# Different forms of genetic variation

-single nucleotide polymorphisms (or SNPs)
-small insertions or deletions (indels)
-short tandem repeats/microsatellites

-retrotransposons (RNA intermediate)

      LINE -long interspersed repetitive elements
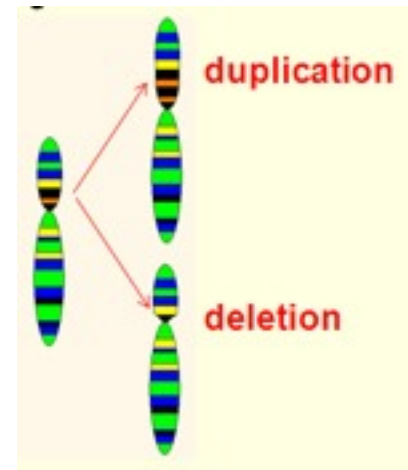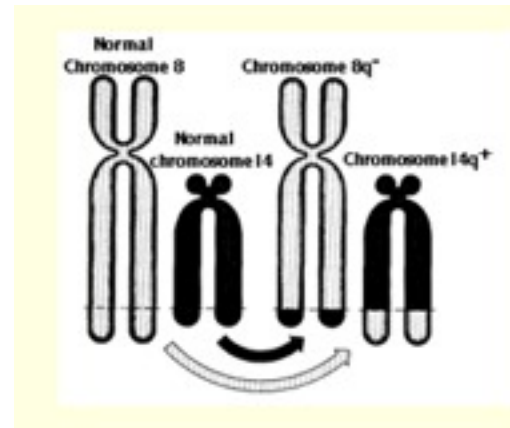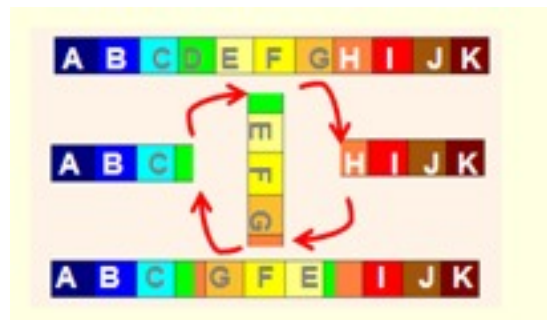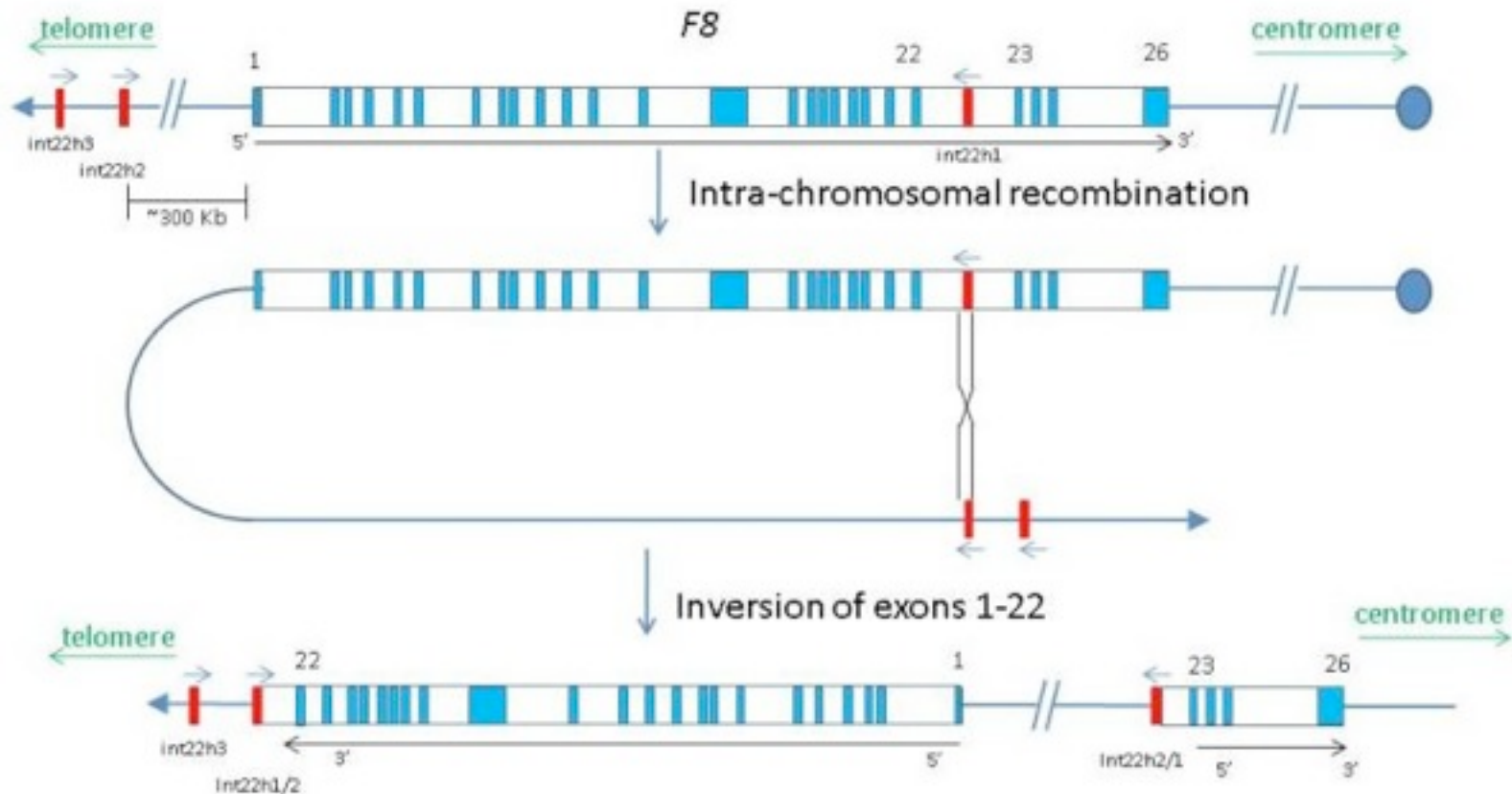               -makes copies
               -17% of genome

      SINE -short interspersed repetitive elements
               -alu sequence (around 300 bp)
               -over 10% of genome

      LTR   -long terminal repeats
               -8% of genome
               -still active?

# Different forms of genetic variation

-single nucleotide polymorphisms (or SNPs)
-small insertions or deletions (indels)
-short tandem repeats/microsatellites
- retrotransposons (RNA intermediate)
-copy number variants (CNVs)
      -deletions/duplications

CNVs



duplication

deletion

# Different forms of genetic variation

-single nucleotide polymorphisms (or SNPs)
-small insertions or deletions (indels)
-short tandem repeats/microsatellites
-transposable elements
-copy number variants (CNVs)
-large structural variation
      -inversions
      -translocations

# Factor VIII gene inversions

-severe hemophilia A
-40% of individuals have a large 400kb inversion

# Variants may be…

-common (>1%)
-rare (<1%)

-single family
-small region
-one population
-all populations

# International HapMap Project

- Identify SNPs from 270 individuals
  - CEU: CEPH (Utah residents with ancestry from northern and western Europe) (30 trios)
  - CHB: Han Chinese in Beijing, China (45 individuals)
  - JPT: Japanese in Tokyo, Japan (45 individuals)
  - YRI: Yoruba in Ibadan, Nigeria (30 trios)

# International HapMap Project

- Establish haplotypes
- Identify Tagging SNPs

# Finding gene variants

# Mapping strategies

- Map based
  - Use a set of markers spread throughout the genome designed to capture most regions/ common variants
- Complete resequencing

# Mapping using genetic variation

linkage analysis using microsatellites



-genotype 300-400 markers for affected and unaffected individuals

-10 cM resolution = 10 Mb



○  *PTPN11* mutation -

⊛  *PTPN11* mutation +

■  **Affected**

⊡  **Clinically unaffected, but mutation confirmed**

# Mapping using genetic variation

Genome-Wide Association Studies (GWAS), using SNPs

Goal: identify common variation associated with a specific phenotype/trait

# GWAS Basic Strategy

- Large sample size (1000's)
- Well defined phenotype
- Case/control or continuous phenotype
- Whole-genome genotyping
- Appropriate statistics

       -correction for multiple testing

       -population stratification

# Whole Genome Genotyping

-10 million common variants in our genome

-Genotyping chips can assay over 2.5 million SNVs

-Excellent coverage of common genetic variants for most populations

# Illumina Infinium Assay

# SNP output

# Copy Number Variation (CNVs)

-Many known common insertions/deletions identified throughout the genome

-Whole-genome genotyping arrays can detect larger CNV events (>50kb)

# Common vs Rare Variation

- GWAS generally only have the power to detect common SNP variation


- Rare variation?

  - Genome Sequencing

# Transition to Sequencing

Moore's Law

sequencing throughput/cost

| 2007 | 2008 | 2009 | 2010 | 2011 |

single run 1Gb                                                single run 650 Gb

# Next Generation Sequencing Setup

-Huge amount of data (terabytes!)

-Analysis computationally intensive

-Dedicated IT infrastructure

# Sequencing Approach

1. Identify subjects

   -Extreme phenotype or family based
2. Sequence (50-100 individuals)
3. Align to reference and call variants
4. Compare to 100's of sequenced controls

5. Follow-up genotyping in **larger** cohorts!

# Sequencing Approach



Cirulli ET& Goldstein DB. *Nature Reviews Genetics*, 2010

# Filtered FastQ sequence

Data from a single cluster/read (75bp)

@G:1:1:11:1079#0/1

TGATTGATTCCATTCCATTCCATTCCATTTCATTCCATTGCAATCCCTTCCAATCCATTCCATTCCATTCCATTC

+G:1:1:11:1079#0/1
`Xa^YO\_^a_`__`a__^a^a^_a``^_\'\\]``[XUGXXXXXWUTWWVWUSTXXPUWYYRVWYYYXZYXYWZ

*a finished genome will have over 1 billion reads

# Analysis

- Alignment to reference genome
  - 3 billion bases
- Call variants
  - Single nucleotide variants (SNVs)
  - Small insertion/deletions (indels)
  - Structural Variants (SV/CNV)

# Summary of a single human genome

| SNVs | 3.5 million |
|---|---|
| Premature stop | 120 |
| Stop loss | 25 |
| Non-synonymous | 11,000 |
| Synonymous | 11,000 |
| Essential splice site | 100 |

| indels | 610,000 |
|---|---|
| Frameshift | 500 |
| In-frame | 900 |

# Whole genome vs. exome sequencing

Exome
    -Coding regions
    -Cheaper/Faster
    -Uneven capture of both alleles
    -Incomplete capture of target region
    -Bias towards known biology

Genome
    -Complete sequence
    -Expensive/Throughput
    -IT issues (10 fold more data)

# GWAS Examples

# Host genetic differences contribute to variation in response to HIV

➢ Susceptibility to infection

➢ Natural history of disease

○ Viral load

○ Immunological progression

○ AIDS events / death

# Host genetics of HIV disease

# Genome-wide research demands

- precise phenotype

- careful selection of patients

- efficient genotyping

- powerful analysis

# Phenotype: HIV viral load at set point

# Patients / Cohorts

- ~2500 white patients
- High quality viremia data
- Genetic consent

# Genotyping

## WG chips:

### 500K to >1 mio single nucleotide polymorphisms (SNPs)



FIGURE 5: HUMAN1M-DUO GENOMIC COVERAGE

- CEU (mean 0.95 median 1.0)
- CHB+JPT (mean 0.95 median 1.0)
- YRI (mean 0.85 median 1.0)

The Human1M-Duo BeadChip content covers the majority of HapMap common variation in three distinct populations. Graphs are based on the HapMap release 23 data set of > 2.3 million common SNPs.

# Analysis
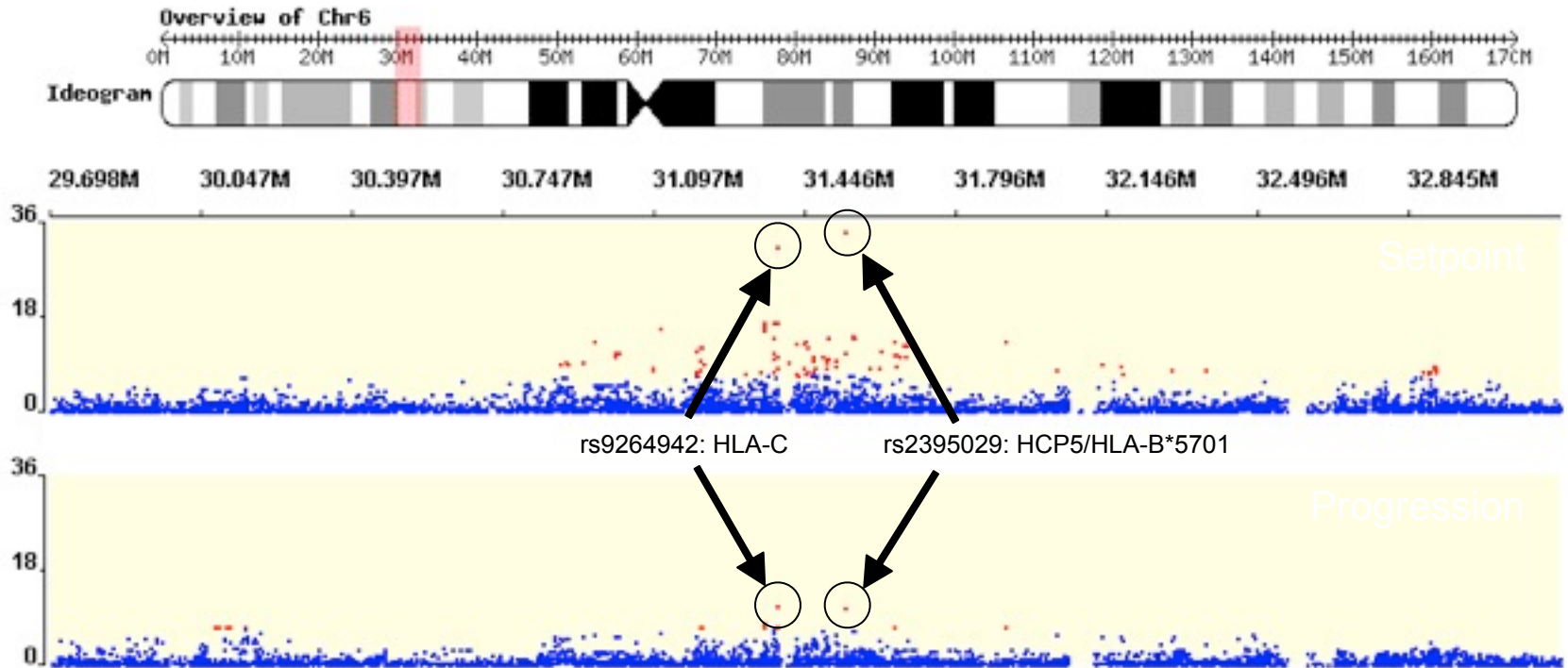
# Analysis

# Analysis

# Genetics of HIV-1 control: results



CHAVI_SETPOINT_FINAL_linear.assoc.linear

# Genetics of HIV-1 control: results



*Fellay et al. Science 2007*

# Genetics of HIV-1 control: results



rs2395029: HCP5/HLA-B*5701

*Fellay et al. Science 2007*

# Genetics of HIV-1 control: results



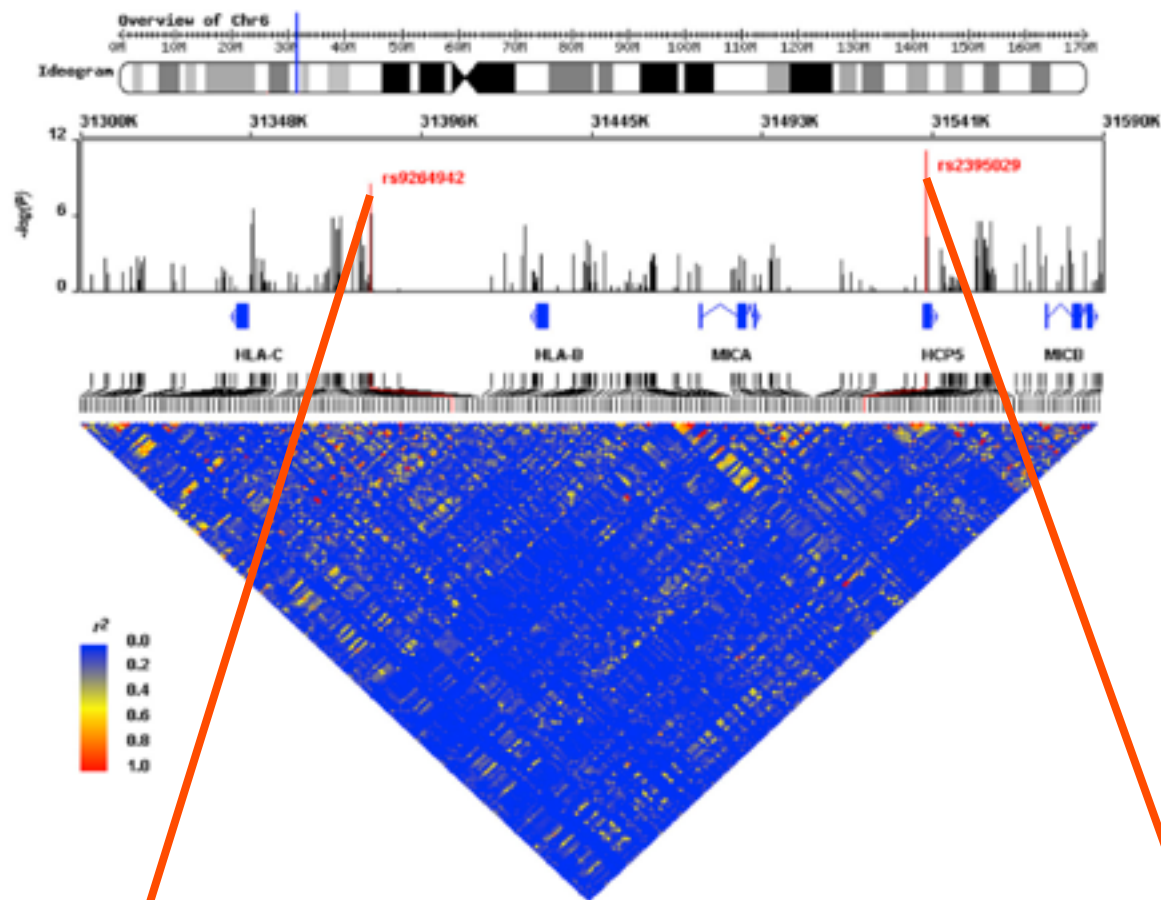*Fellay et al. Science 2007*

# Genetics of HIV-1 control: results



*Fellay et al. Science 2007*

| Gene & SNP | P-value for association with HIV-1 viral load at setpoint N=2362 | P-value for association with protection against progression (CD4 <350) N=1071 |
|---|---|---|
| HCP5 / HLA-B*5701 rs2395029 | 4.5E-35 | 1.2E-11 |
| HLA-C rs9264942 | 5.9E-32 | 7.4E-12 |
| ZNRD1 / RNF39 rs9261174 | 1.1E-04 | 3.8E-08 |
| CCR5 Δ32 het rs333 | 1.7E-10 | 2.6E-06 |

*Bonferroni threshold for genome-wide significance:  5E-08*

# Genes and viral load to predict disease progression
## (no progression after 5 years, %)

# Sequencing example

# Host genetics of HIV-1 control

# Sequencing of extreme HIV progressors

1. Rapid Progressors
   - Known date of seroconversion
   - CD4 <350 in less than 3 years
   - [Severe PHI = immediate CD4 depletion without spontaneous recovery]

2. Controllers
   - VL >50 cp/ml
   - Excluding HLA-B*57, B*27 and B*5801

# Exome sequencing in 31 rapid progressors and 10 controllers

## Overview of genetic variation

- Mean coverage: 72x

- Total SNVs: 101057
- Novel SNVs: 40385 (40%)
- Total indels: 12149
- Novel indels: 5330 (45.5%)

SNVs

indels

**SNVs legend:**
- NON_SYNONYMOUS_CODING 35.07%
- STOP_GAINED 0.41%
- STOP_LOST 0.04%
- upstream 1.26%
- intronic_EXON_BOUNDARY 5.23%
- 5' UTR 2.78%
- 3' UTR 3.13%
- SYNONYMOUS_CODING 24.59%
- intronic 27.40%
- intergenic 0.09%

**indels legend:**
- CODING_DISRUPTED_FRAMESHIFT 35.
- CODING_DISRUPTED_OTHER 8.93%
- TRANSCRIPT_INCLUDED 0.00%
- 5PRIME_UTR 3.43%
- 3PRIME_UTR 2.87%
- INTRONIC_EXON_BOUNDARY 8.86%
- UPSTREAM 1.70%
- DOWNSTREAM 1.23%
- INTRONIC 37.78%
- INTERGENIC 0.10%

# Analysis of genetic variants

➢ Single variant analysis (ATAV)
  ○ Case-control comparison of single variants (SNVs and indels) using Fisher's exact tests for allelic, dominant, recessive, and genotypic models, plus Cochran-Armitage trend test

➢ Ranking of putatively functional variants (SVA)
  ○ listing of homozygous or heterozygous variants observed mostly (or only) in cases, ranked by numbers

➢ Gene prioritization (SVA and ATAV)
  ○ Case-control comparison of genes carrying key functional variants, using Fisher's exact tests with assessment of genome-wide significance by permutations

# How does it work?

# Single variant analysis

| variant | RS | gene | function | RP | VC | Ctrls | P_value |
|---|---|---|---|---|---|---|---|
| 19_59711073_T | - | **LAIR2 / CD306** | STOP_GAINED | 1/4/26 | 0/0/10 | 0/19/208 | 0.05 (genotypic:RPvsCtrls) |
| 5_86731030_G | rs2230641 | **CCNH** | NS | 2/12/17 | 0/0/10 | 10/59/160 | 0.009 (allelic:RPvsVC) |
| 16_55617854_C | rs28438857 | **NLRC5** | NS | 3/7/21 | 0/1/9 | 0/59/167 | 0.002(recessive:RPvsCtrls) |
| 11_60533649_A | rs12360861 | **CD6** | NS | 3/8/19 | 0/3/7 | 2/48/160 | 0.02 (trend:31vsCtrls) |
| 4_74921673_INS_T | - | **CXCL6** | FRAME SHIFT | 0/2/29 | 0/0/10 | 0/2/226 | 0.018 (trend:RPvsCtrls) |
| 1_26517124_A | - | **CD52** | NS | 0/3/28 | 0/0/10 | 0/0/229 | 0.0006 (trend:RPvsCtrls) |
| 1_12108645_G | rs2230625 | **TNFRSF8 / CD30** | NS | 0/3/28 | 0/3/6 | 0/5/210 | 0.0009 (trend:VCvsCtrls) |
| 1_158052037_T | rs61823162 | **FCRL6** | STOP_GAINED | 3/3/25 | 0/1/9 | 7/52/170 | 0.06 (genotypic:RPvsCtrls) |

# Next steps

1. **More "extreme" samples:**

   ➜         MACS:    25 rapid progressors
                   25 controllers

2. **More sequence:**

   ➜              Whole genome sequencing

# From GWAS to sequencing, and beyond…

- Only a limited amount of the genetic basis for much phenotypic variation has been located by GWAS
  - Still much 'missing heritability'

- Rare and/or causal variants will be identified by sequencing

- Data integration and systems approaches represent the next frontier