

Mapping the genes underlying phenotypic changes of interest

**Virginie Courtier-Orgogozo
Institut Jacques Monod, Paris**

Tomato shape



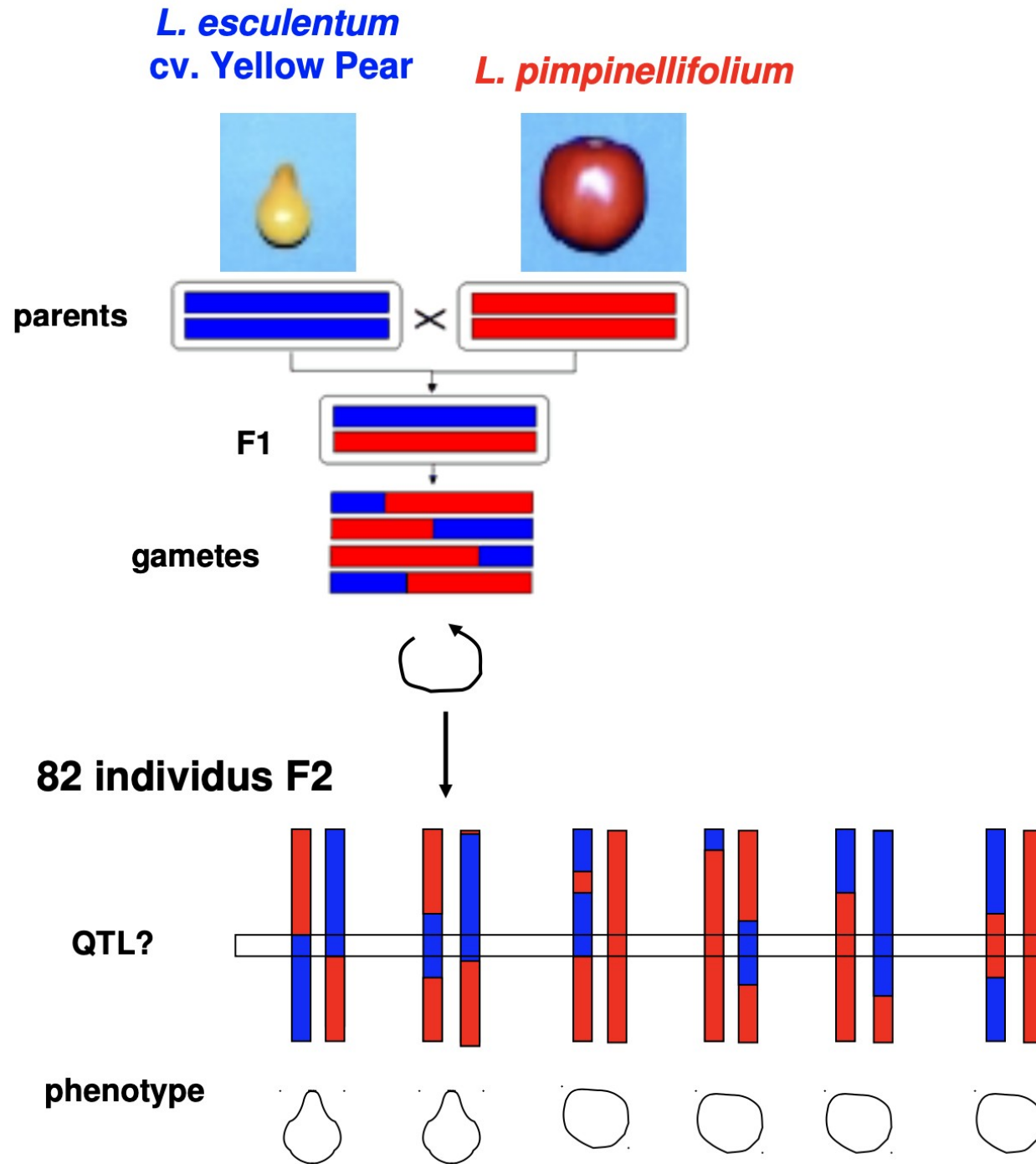
Lycopersicon esculentum



Lycopersicon esculentum cv. Yellow Pear

(Ku et al., 1999; Liu et al., 2002)

QTL mapping

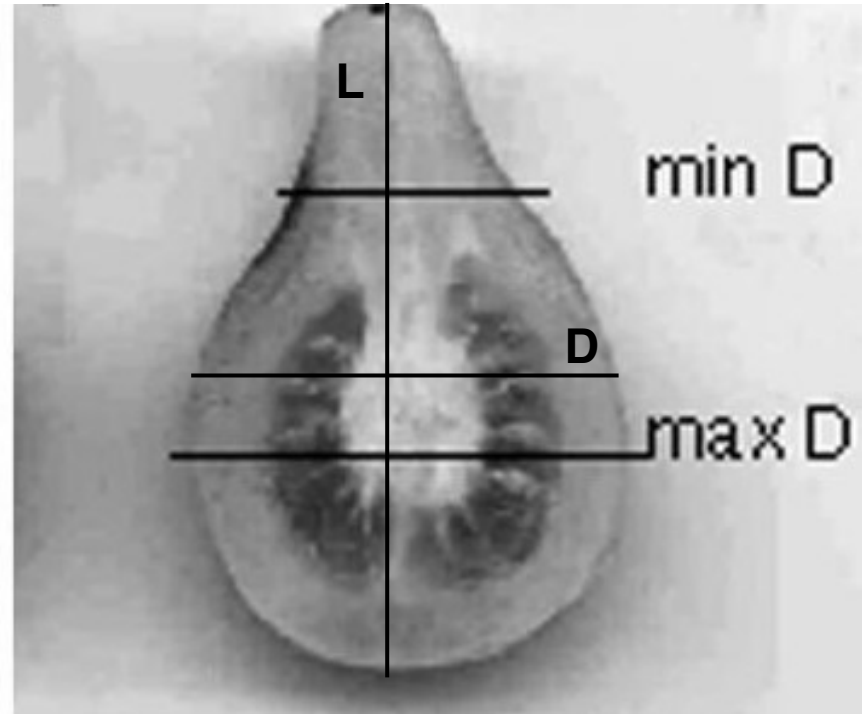
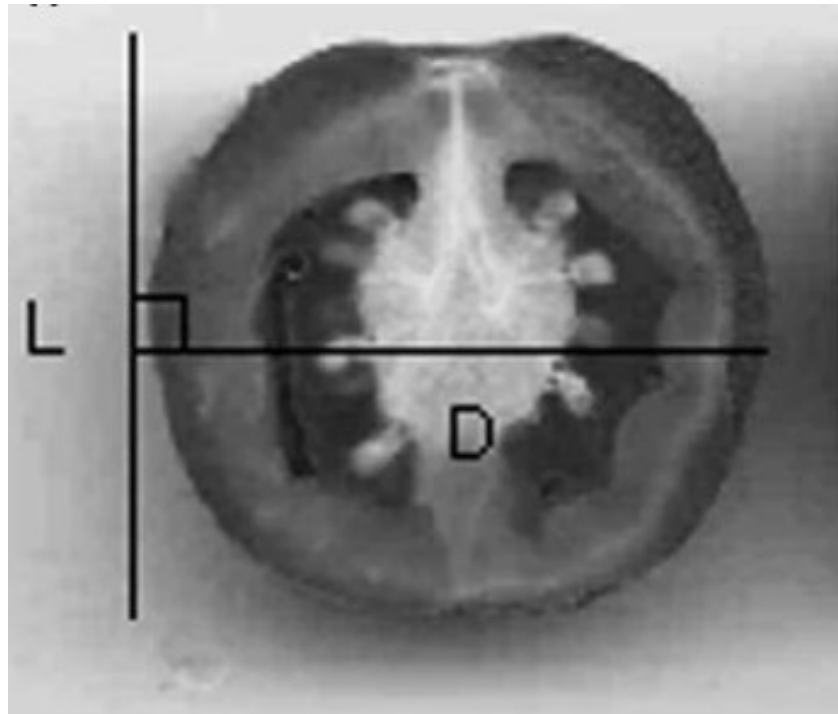


Quantitative measure of the phenotype

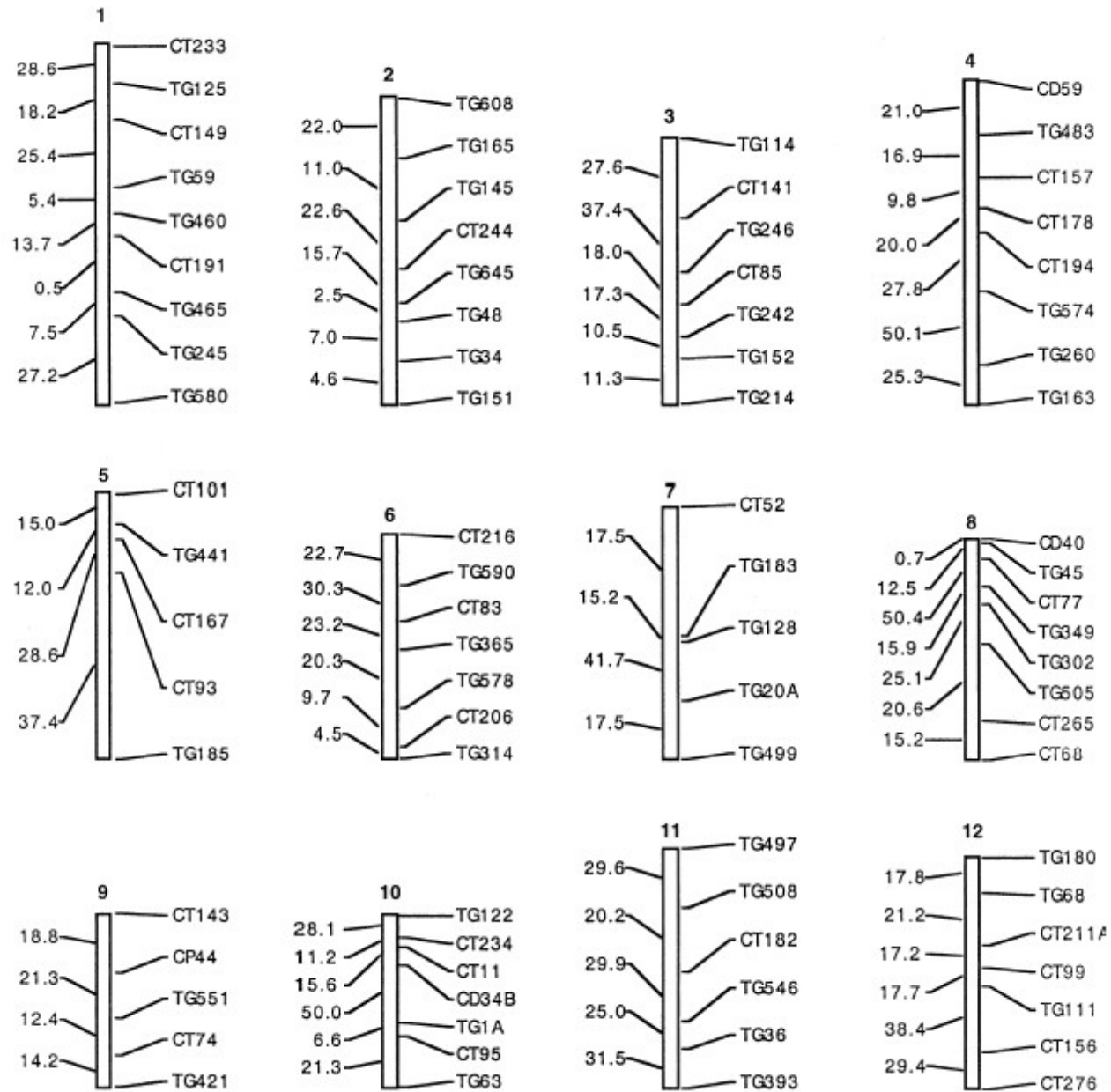
Measure of 2 indexes L/D and D_{min}/D_{max} for 10 fruits per plant

L/D : L = length, D = diameter at equator

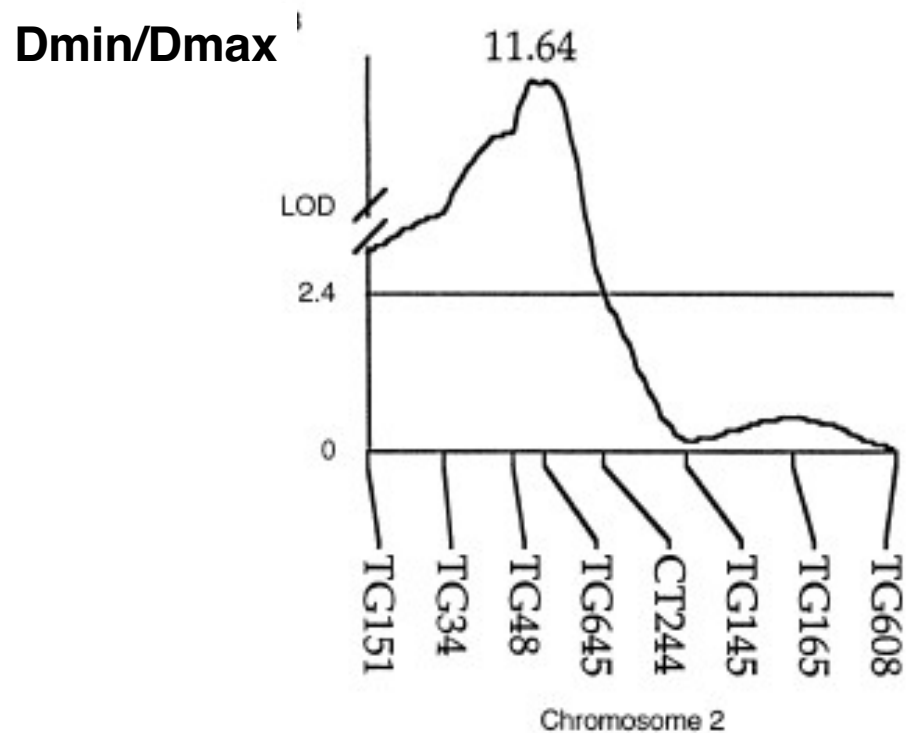
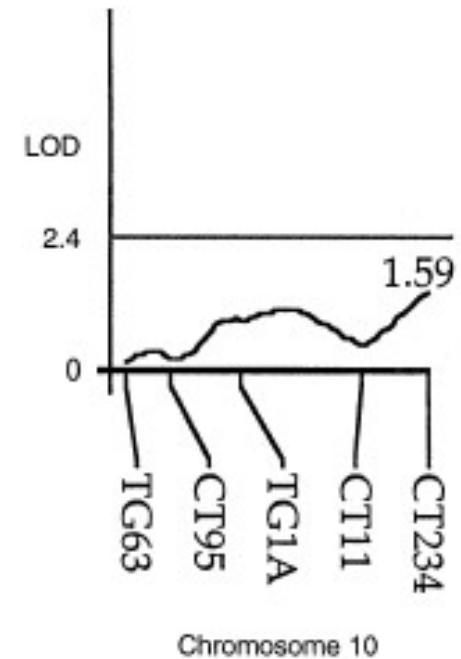
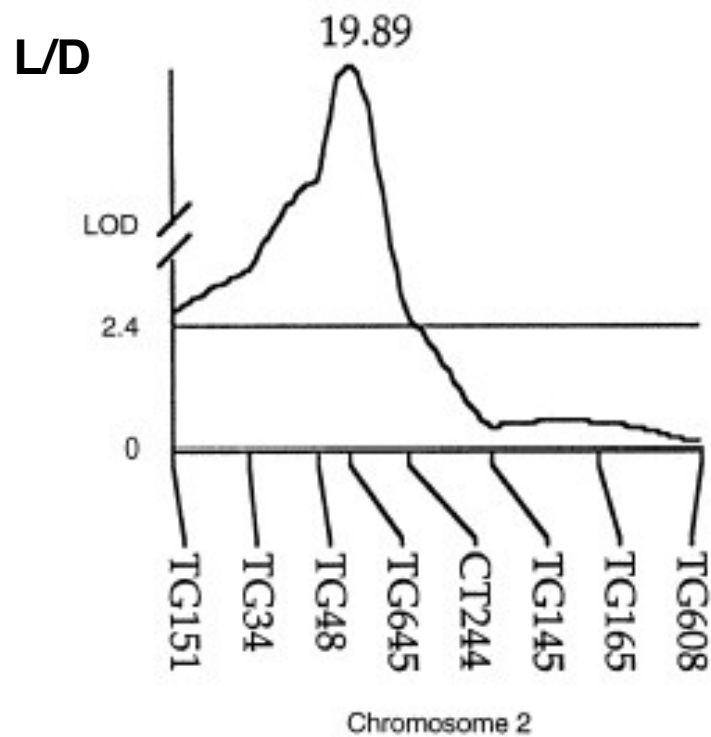
D_{min}/D_{max}



82 molecular markers on the 12 tomato chromosomes



One major locus near marker TG645



responsible for 67% of L/D variance

allele YP = recessive

Two main files

Markers file

```
-start
-Chromosome 1
CF5475      0.4
CF5573      24.7
CT7895      41.0
CT8903      59.0
CF5613      67.7
CT7892      76.0
CT890       89.0
CT233       39.0
Telomere    50.0
-Chromosome 2
CF5671      0
CF5675      10.4
CF5673      34.7
CT789       41.0
CT890       89.0
CT567       115.0
Telomere    130.0
...
```

Genotypes and phenotype(s) file

```
-start individuals markers
Ind_1 0 0 1 1 0 0 0 0 0 1 2 2 2 2
Ind_2 0 0 0 1 0 1 0 0 1 1 1 1 0 0
Ind_3 2 2 2 2 2 1 0 1 1 1 1 0 0 0
Ind_4 0 1 0 0 0 0 0 1 1 1 2 2 1 1
Ind_5 0 1 0 0 0 0 0 1 1 1 1 2 2 2
Ind_6 1 1 1 1 1 1 1 1 1 0 0 0 0 0
Ind_7 1 1 1 1 1 1 1 0 1 n n 1 1 1
Ind_8 2 2 2 1 1 1 1 0 1 1 1 1 1 0
Ind_9 1 1 1 1 1 1 1 0 0 1 1 1 1 1
Ind_1 0 2 2 1 1 1 1 1 0 0 0 1 1 2
-stop individuals markers

-start individuals traits 1 LoverD named
Ind_1      5.5
Ind_2      3.0
Ind_3      4.0
Ind_4      7.0
Ind_5      6.5
Ind_6      5.0
Ind_7      3.5
Ind_8      6.0
```

Simple linear regression for each marker

L/D of individual $i = a + b \cdot x_i + \varepsilon$

$x_i = 0$ if Le/Le, $= 1$ if Le/Lp, $= 2$ if Lp/Lp

a, b = best fit parameters (least square regression)

ε assumed to have a normal distribution

Test $H_0: b = 0$ versus $H_1: b = \text{estimated } b$

Likelihood ratio test statistic

$$D = -2(\ln(\text{likelihood for null model}) - \ln(\text{likelihood for alternative model}))$$

$$= -2 \ln \left(\frac{\text{likelihood for null model}}{\text{likelihood for alternative model}} \right).$$

The **probability distribution** of the **test statistic** can be approximated by a **chi-square distribution** with $(df_1 - df_2)$ **degrees of freedom**, where df_1 and df_2 are the degrees of freedom of models 1 and 2 respectively

Interval mapping

L/D of individual $i = a + b \cdot x_i + e$

x_i = indicator variable specifying the probabilities of an individual being in different genotypes for the tested position, constructed by flanking markers

$x_i = 0$ if Le/Le, $= 1$ if Le/Lp, $= 2$ if Lp/Lp

a, b = best fit parameters (maximum likelihood)

Test $H_0: b=0$ versus $H_1: b=\text{estimated } b$



Interval mapping

L/D of individual $i = a + b.x_i + e$

x_i = indicator variable specifying the probabilities of an individual being in different genotypes for the tested position, constructed by flanking markers

$x_i = 0$ if Le/Le, $= 1$ if Le/Lp, $= 2$ if Lp/Lp

a, b = best fit parameters (maximum likelihood)

Test $H_0: b=0$ versus $H_1: b=\text{estimated } b$

Composite Interval mapping

L/D of individual $i = a + b.x_i + c.y_i + e$

x_i = indicator variable specifying the probabilities of an individual being in different genotypes for the tested position, constructed by flanking markers

$x_i = 0$ if Le/Le, $= 1$ if Le/Lp, $= 2$ if Lp/Lp

$y_i = 0$ if Le/Le, $= 1$ if Le/Lp, $= 2$ if Lp/Lp at marker y

LOD score

$$L/D \text{ of individual } i = a + b \cdot x_i + e$$

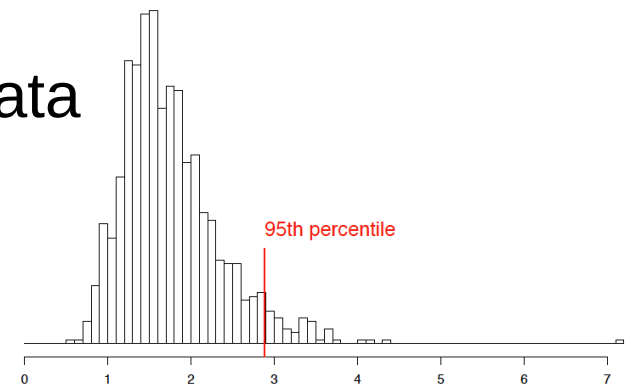
Test $H_0: b = 0$ versus $H_1: b = \text{estimated } b$

$L_0 = \text{pr}(\text{data} \mid \text{no QTL})$ – phenotypes assumed to follow a normal distribution
 $L_1 = \text{pr}(\text{data} \mid \text{QTL at tested position})$

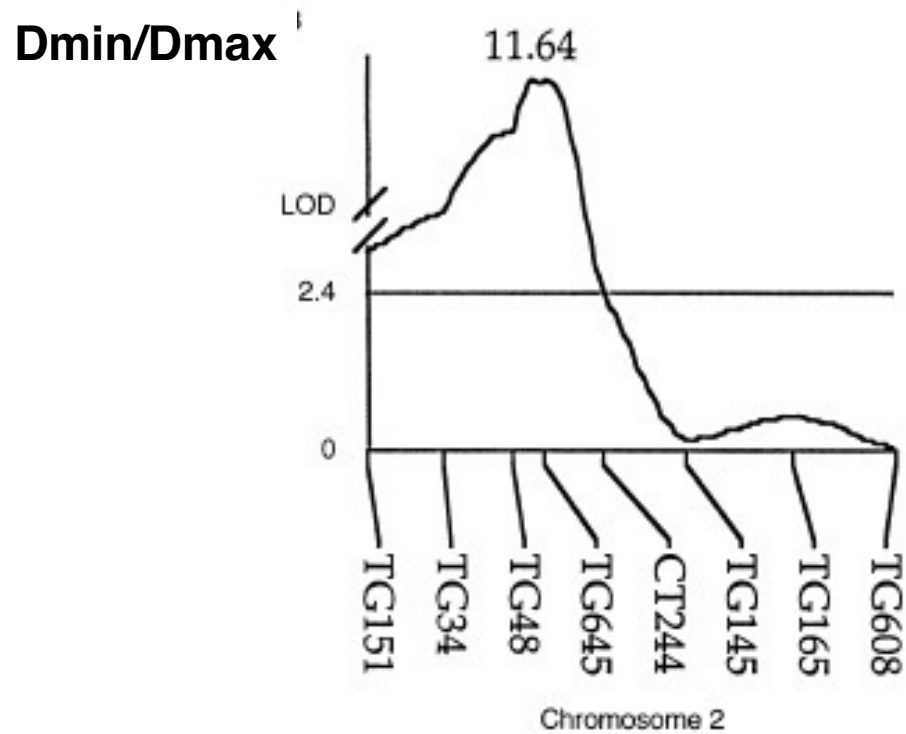
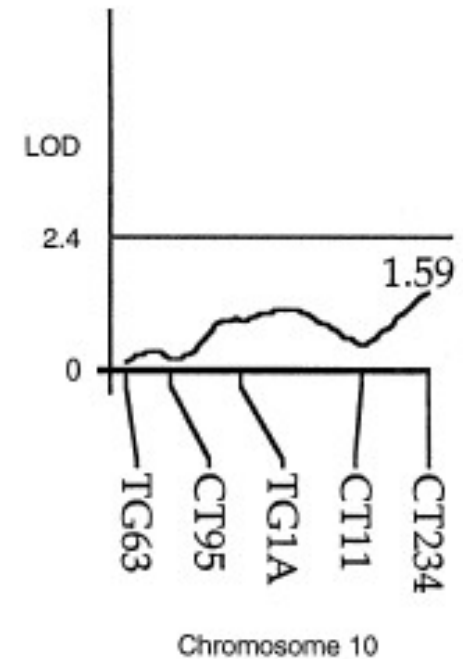
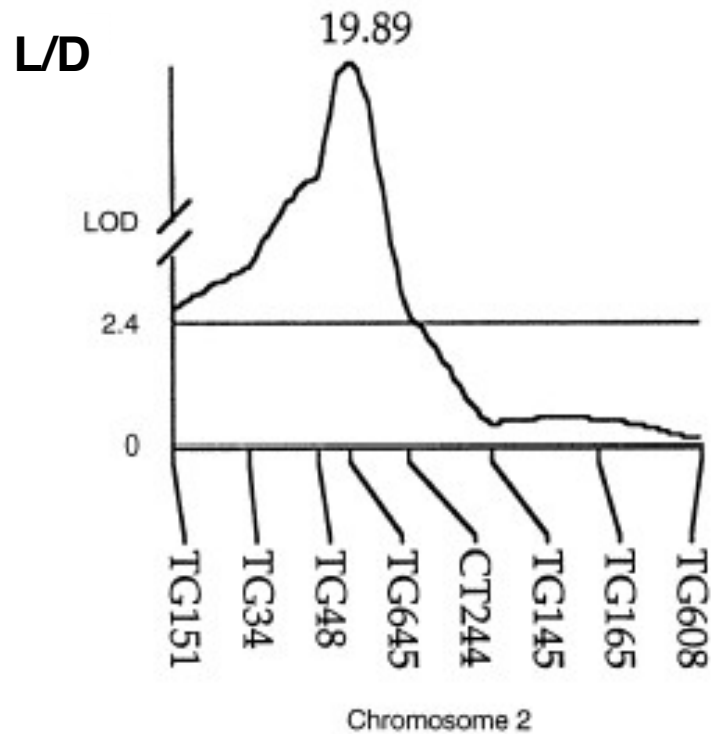
$$LOD = -\log \frac{L_0}{L_1}$$

Significance threshold

- 10,000 permutations of phenotype/genotype data
 - random distribution of LOD scores
 - 1% or 5% significance threshold



One major locus near marker TG645



responsible for 67% of L/D variance

allele YP = recessive

BAC library (Bacterial Artificial Chromosomes)

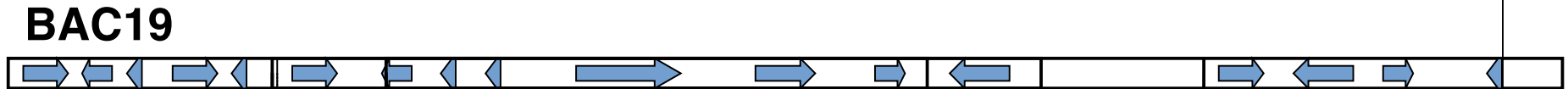
contains genomic DNA fragments of 100-350kb



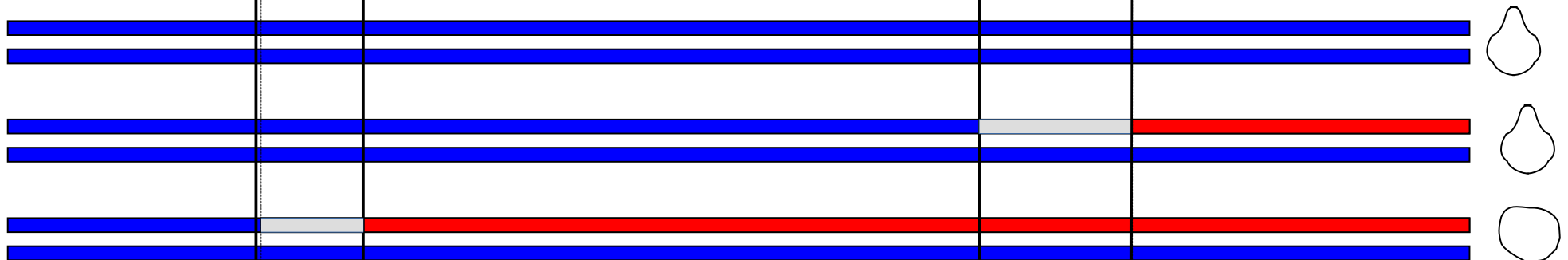
Screen of the library with marker TG645

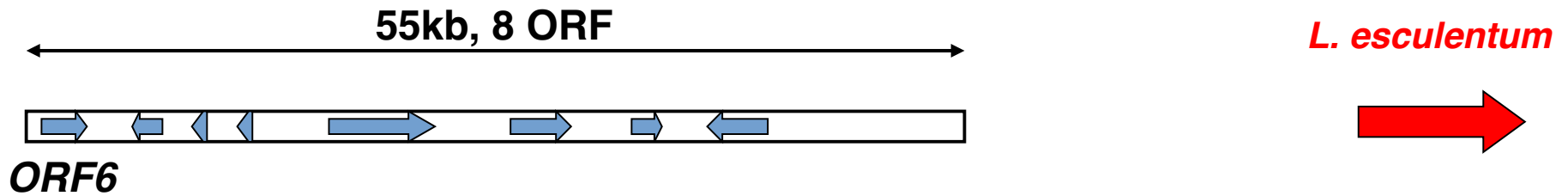


BAC19 containing 105kb, 17 ORF (open reading frame) TG645



Design of new molecular markers to genotype the previously obtained recombinant tomato plants





Sequencing of the region in the 2 tomato varieties

1 SNP (single nucleotide polymorphism)
 et 1 indel (insertion-deletion) of 2bp
 in non-coding regions

L. esculentum
 cv. Yellow Pear



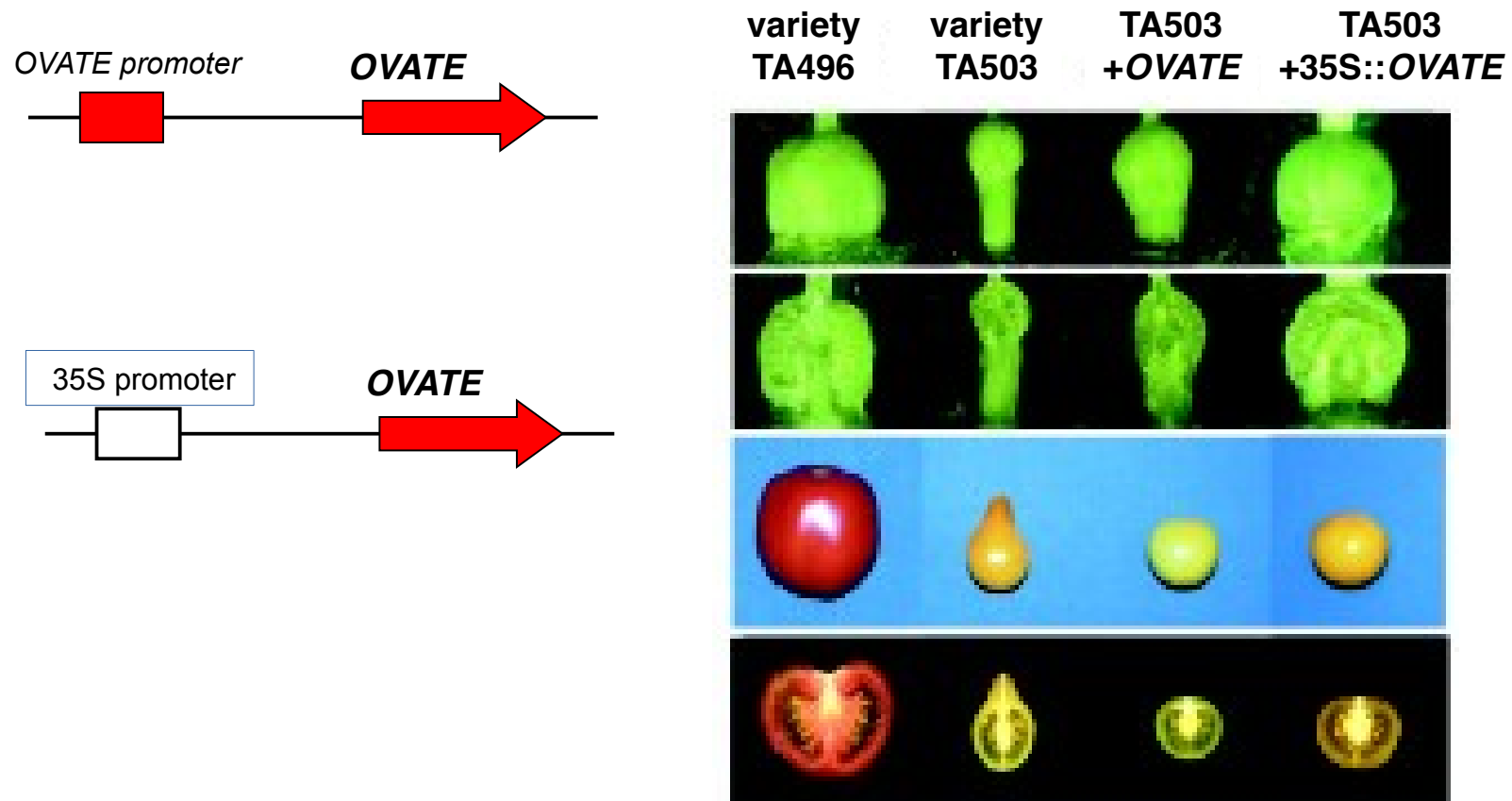
1 SNP in *ORF6* : G496T, stop codon stop,
 truncated protein with last 75 amino acids
 missing

Hypothesis: the causing gene is *ORF6* = *OVATE*

The causing gene is *OVATE/ORF6*

Same mutation in 3 other pear tomato varieties

Complementation of the mutation by transgenesis



OVATE = protein with NLS (nuclear localization signal), unknown function, expressed in developing fruits

Evolution of morphology in threespine sticklebacks



marine



Paxton Lake, Canada

Gasterosteus aculeatus

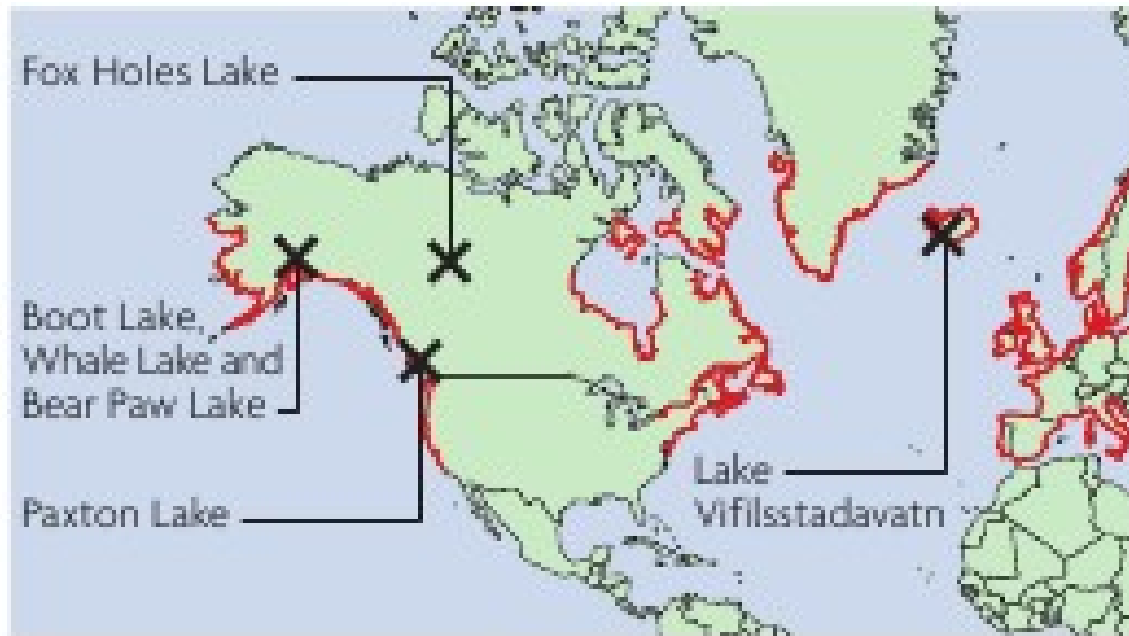
(Peichel et al., 2001 ; Shapiro et al, 2004 ; Chan et al. 2010)

Marine fishes with robust pelvis = ancestral

**Freshwater fishes with reduced pelvic structures = derived,
independently at least 20 times**

- limited calcium availability
- absence of gape-limited predatory fishes
- predation by grasping insects

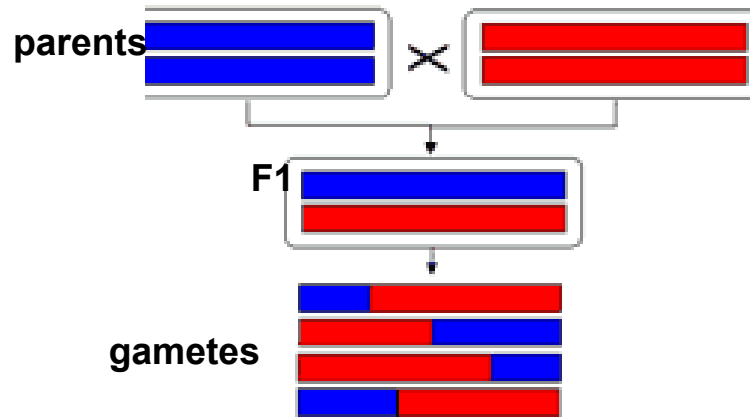
Last glacier retreat = 10 000 – 20 000 years ago



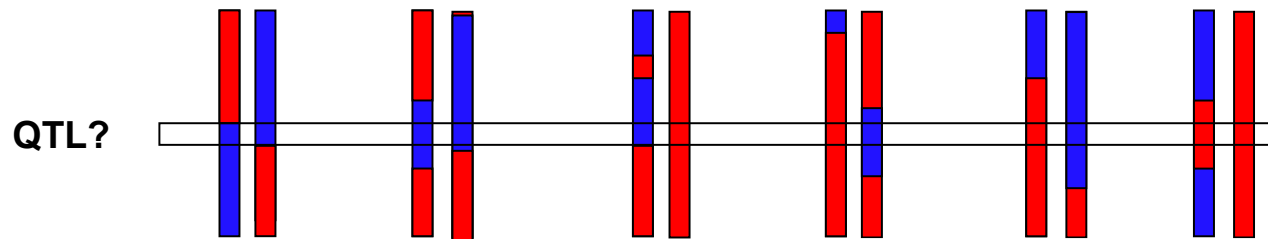
QTL mapping

lake

marine



375 F2 individuals

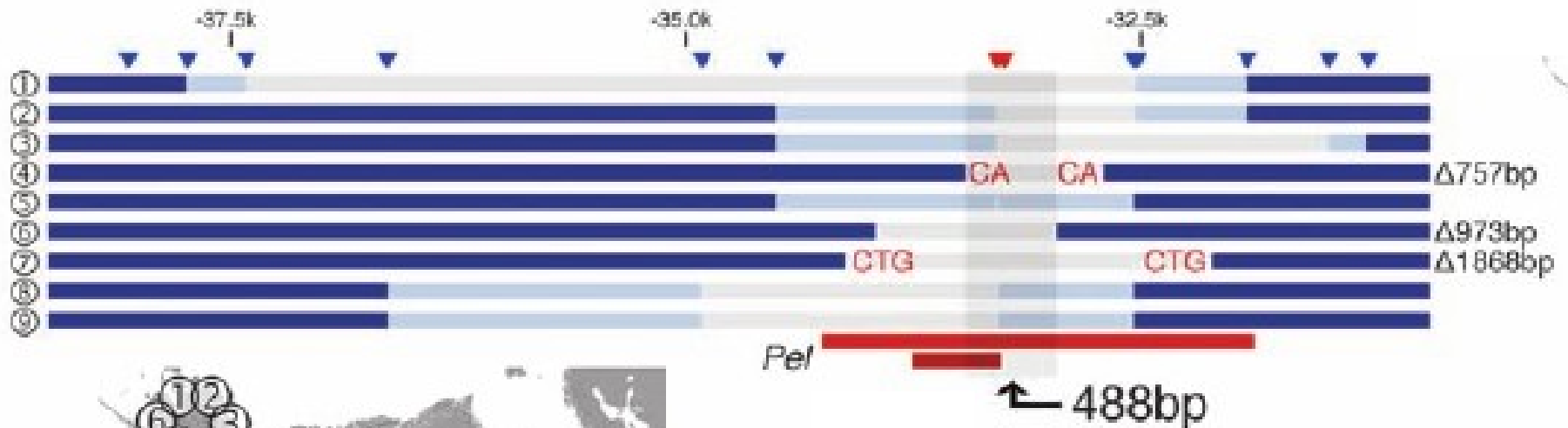


(Shapiro et al., 2004)

Several independent deletions in the cis-regulatory region of *Pitx1*

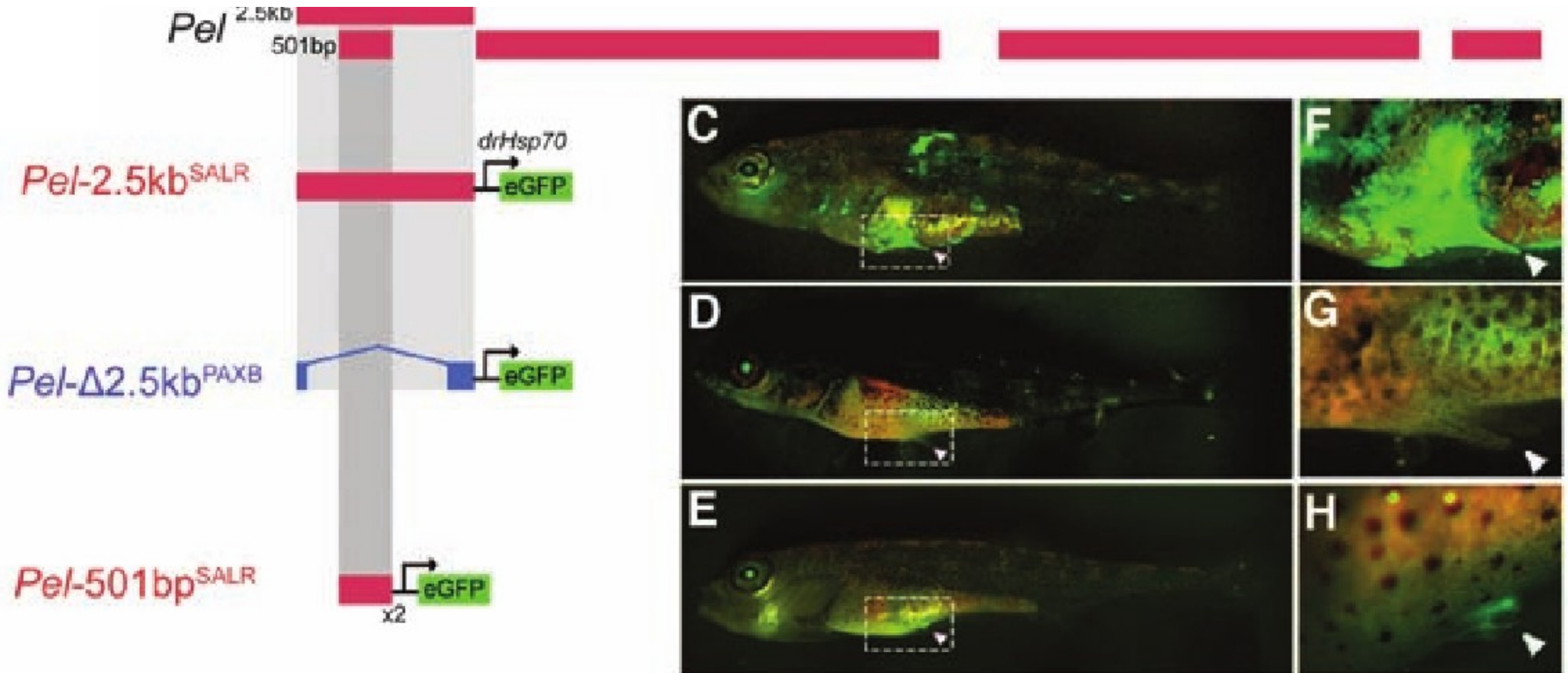
Region sequenced in two lake populations: a 2-kb deletion in one
and a 757-bp deletion in the other one

SNP genotyping in 13 populations with reduced pelvis
and in 21 populations with complete pelvis

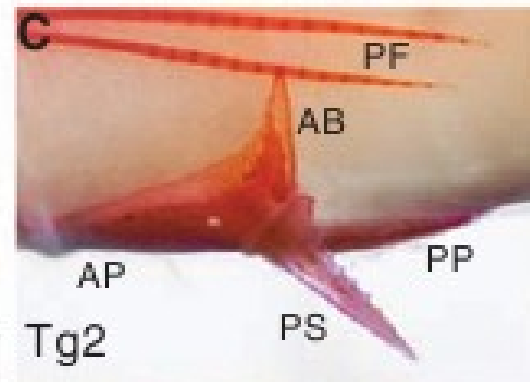


9 different deletions

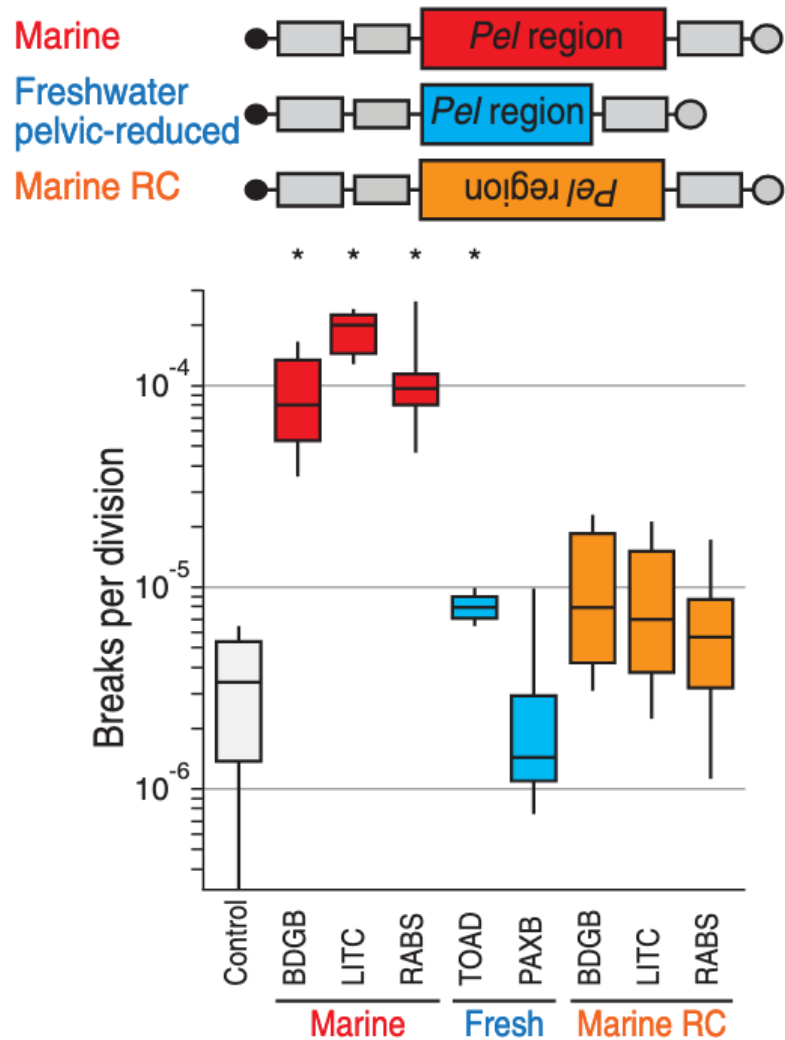
Test of *Pitx1* cis-regulatory regions



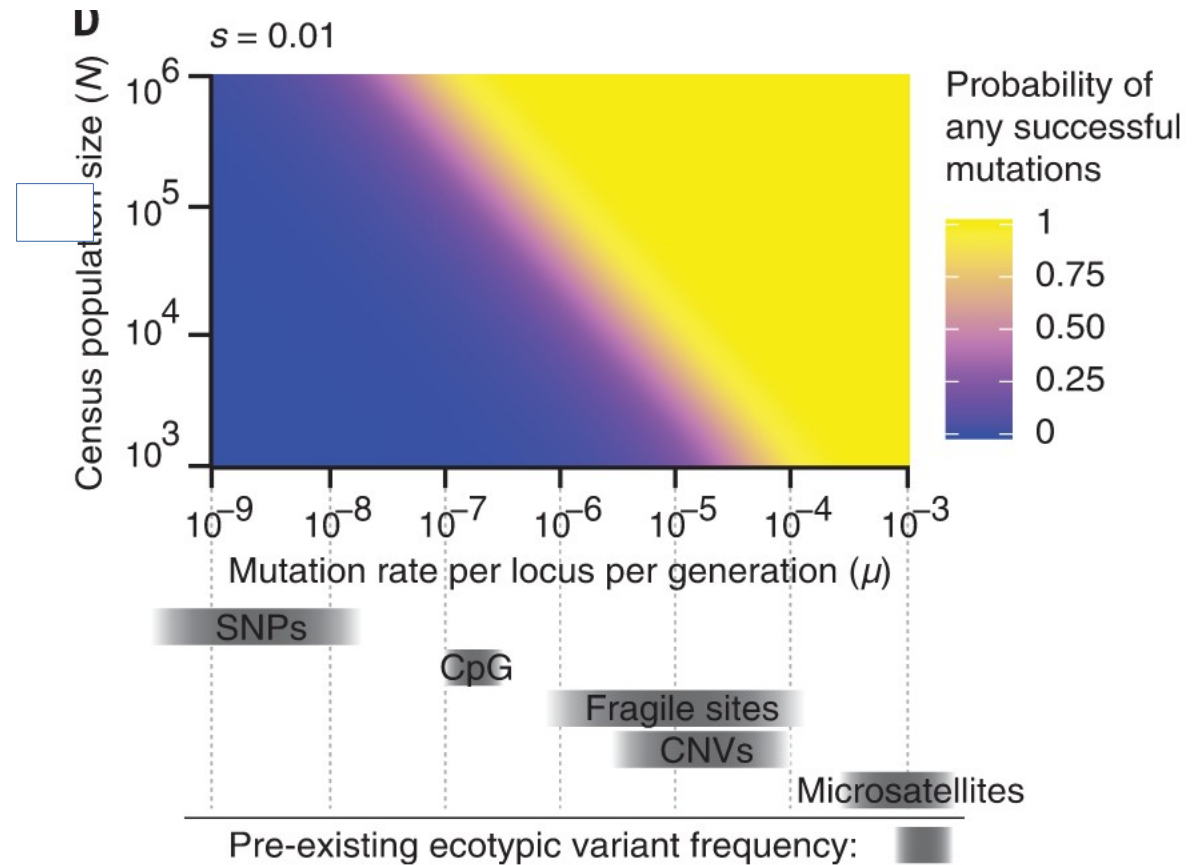
Rescue of a pelvis in freshwater individuals

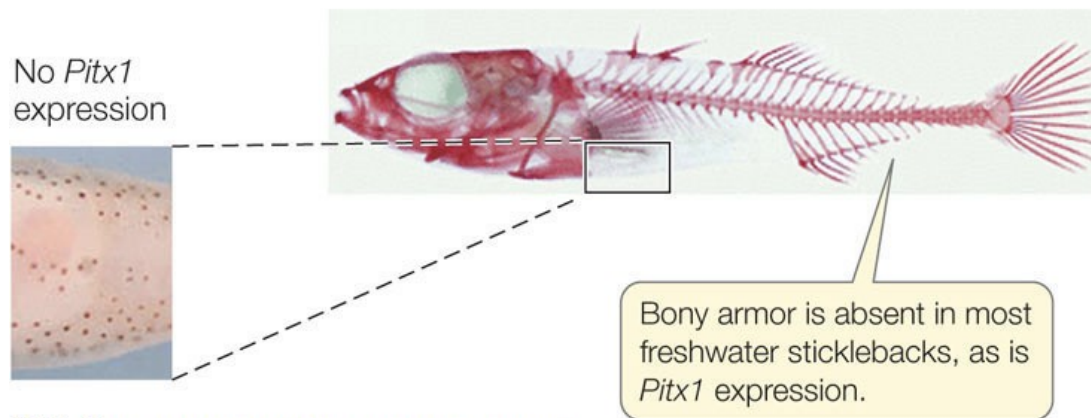
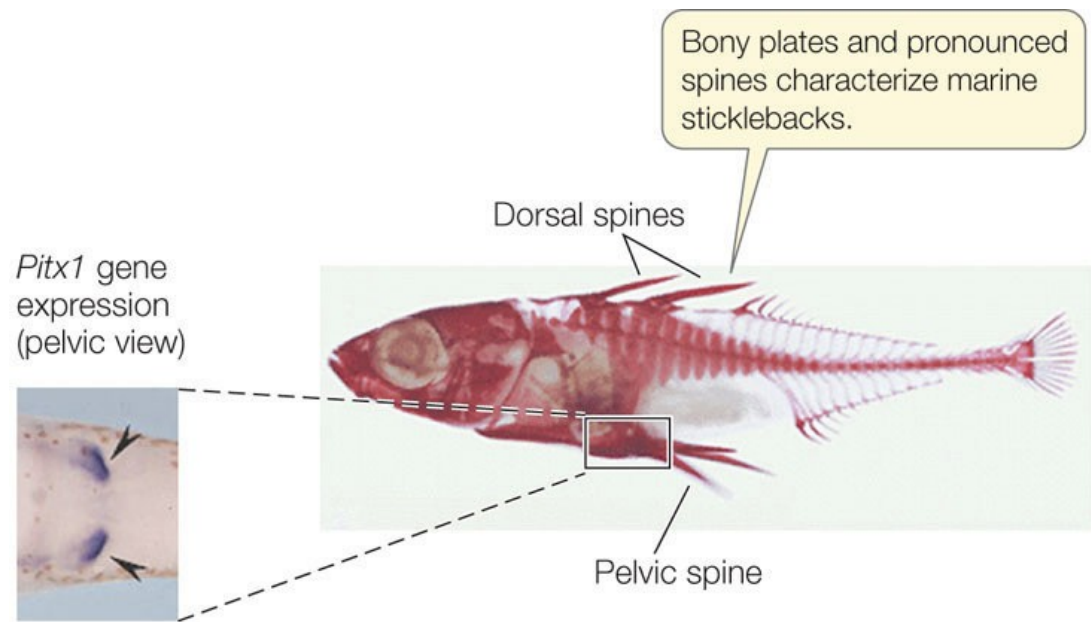


Pitx1 is in a fragile DNA region



Control = artificial chromosome without test region

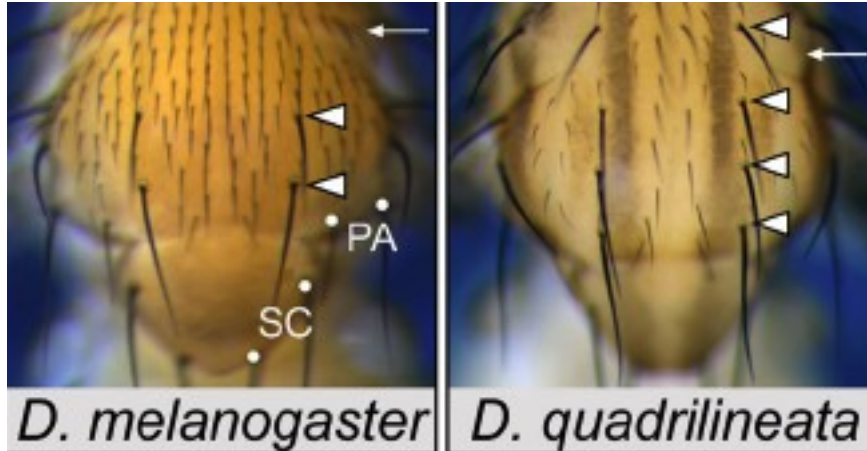




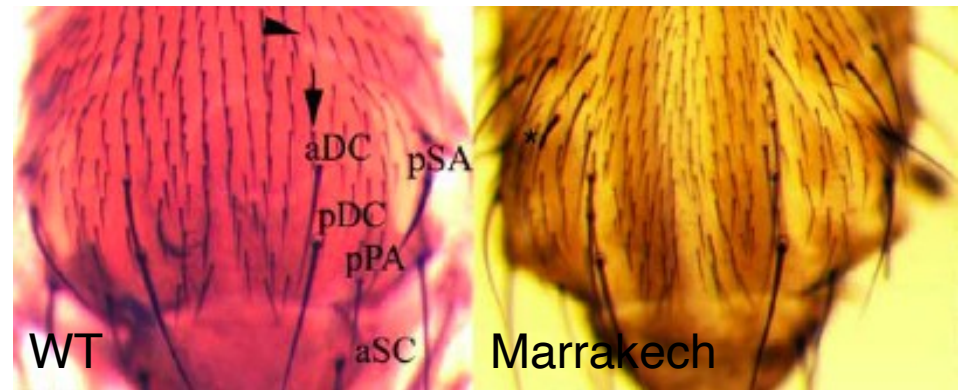
14.21: Courtesy of Mike Shapiro and David Kingsley.

Evolution of extra bristles

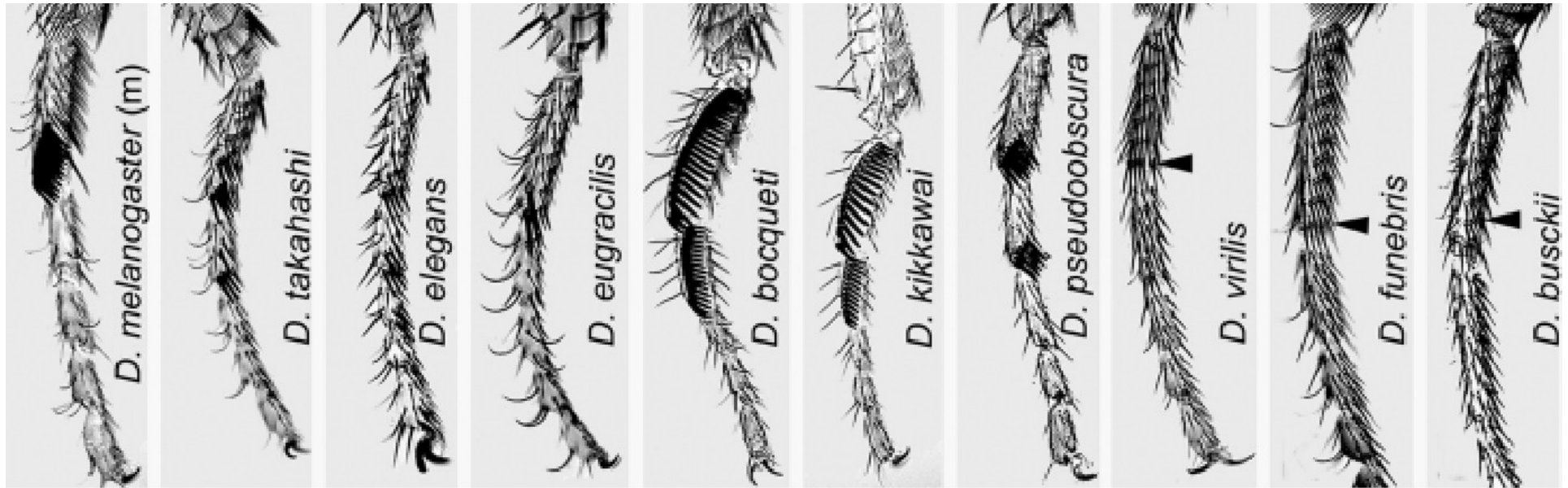
Interspecific change
in *D. quadrilineata*



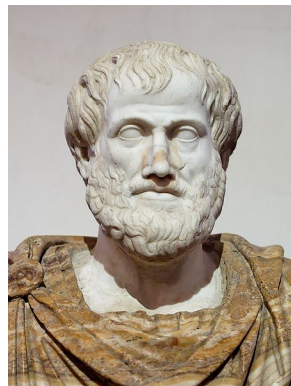
Intraspecific change
in *D. melanogaster*



Finding genetic rules on bristle evolution



Randsholt and Santamaria 2008



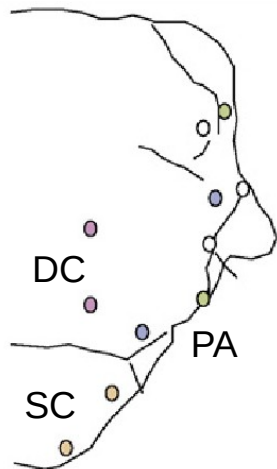
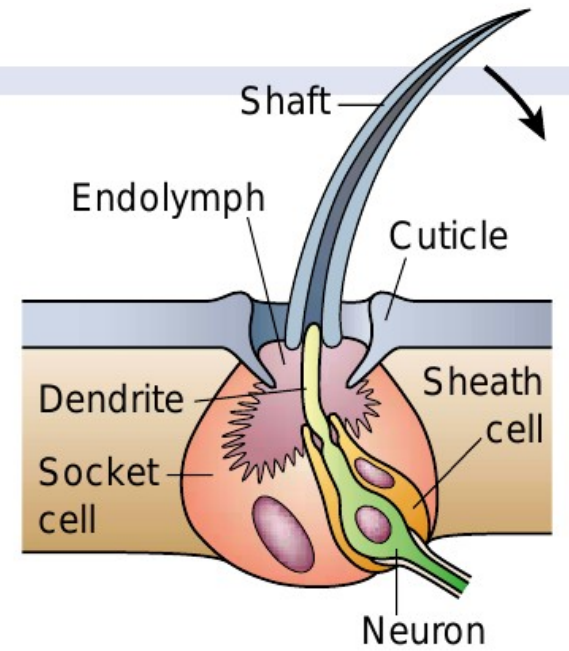
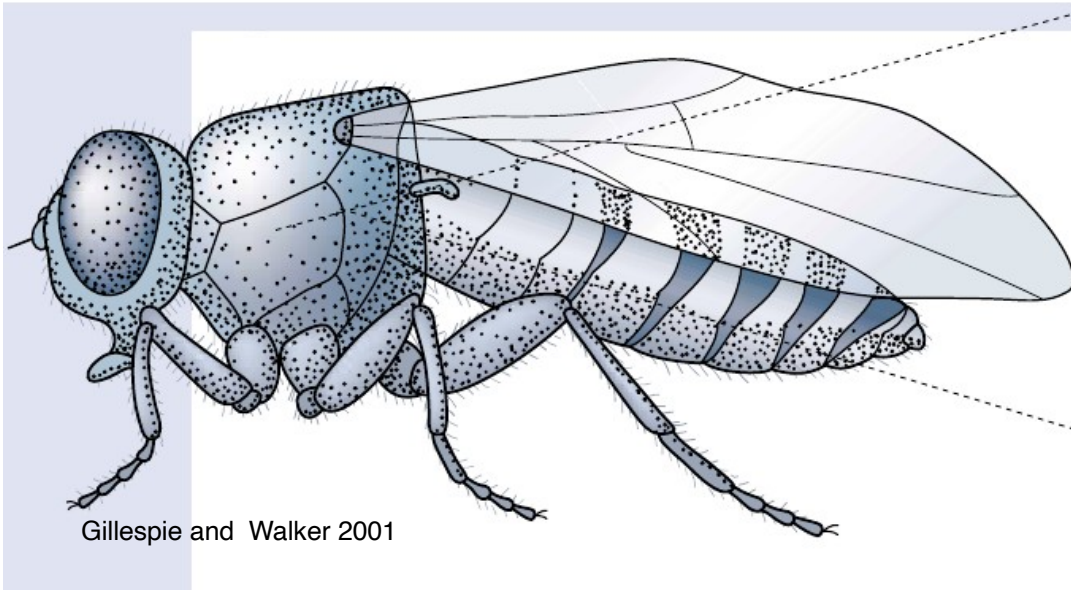
Aristotle, *Historia animalium*, book I, 2, 300BC

color, type, orientation
shape and size
presence/absence
position

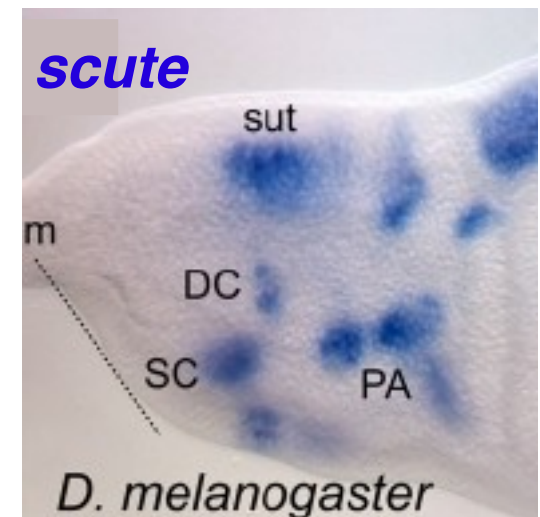
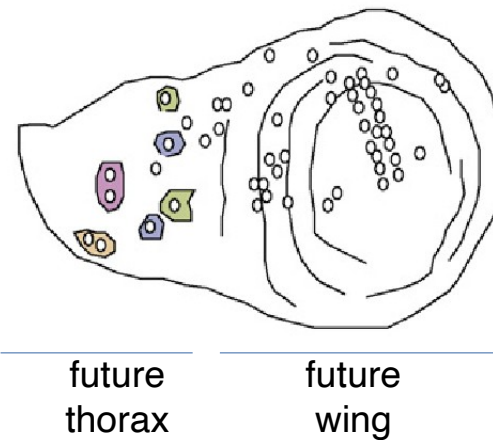
CRE mutations
in *achaete-scute*

Stern and Orgogozo 2009

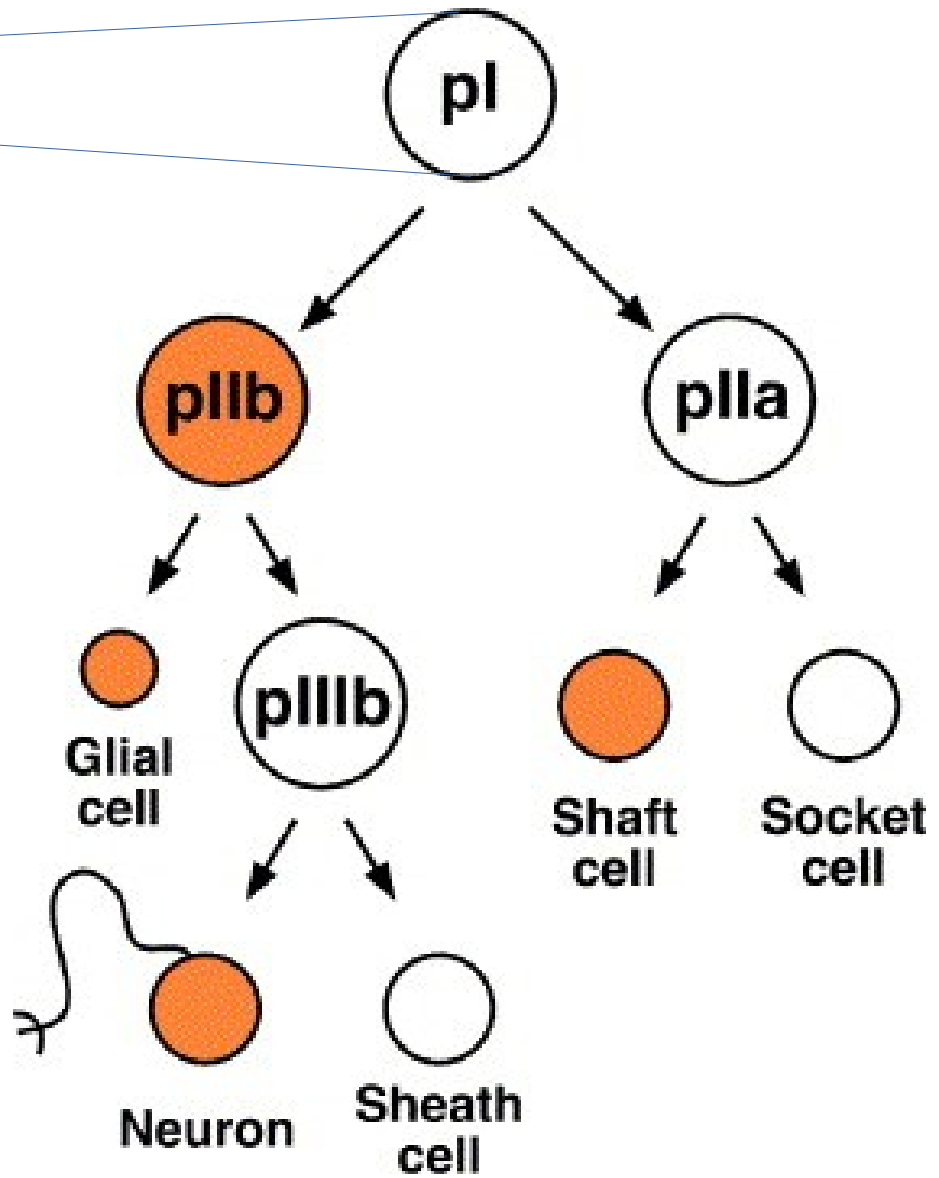
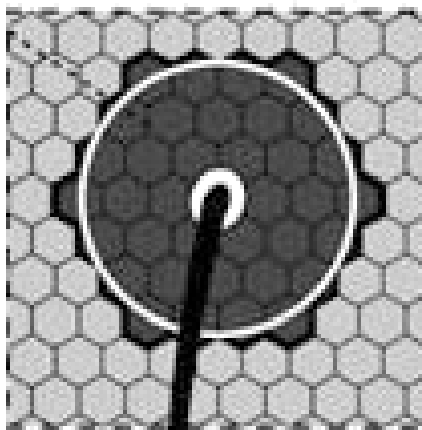
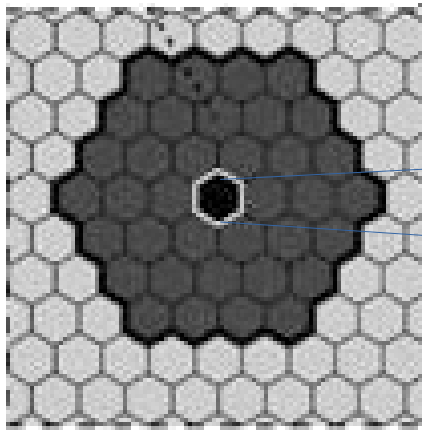
Bristle development



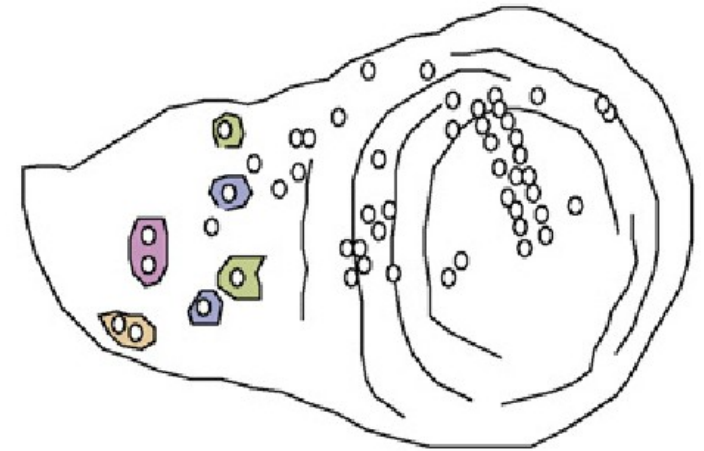
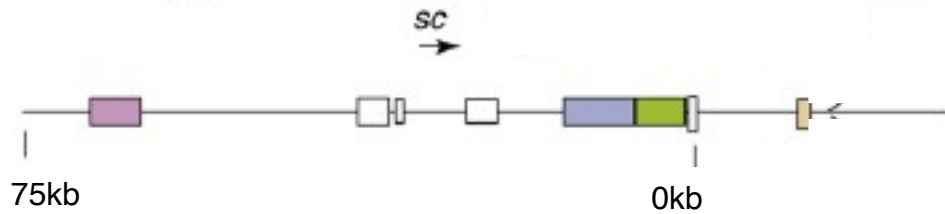
Simpson 2007



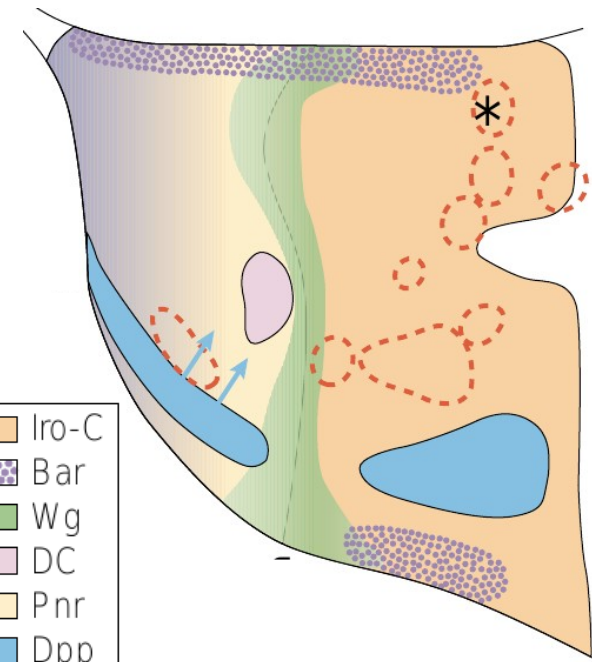
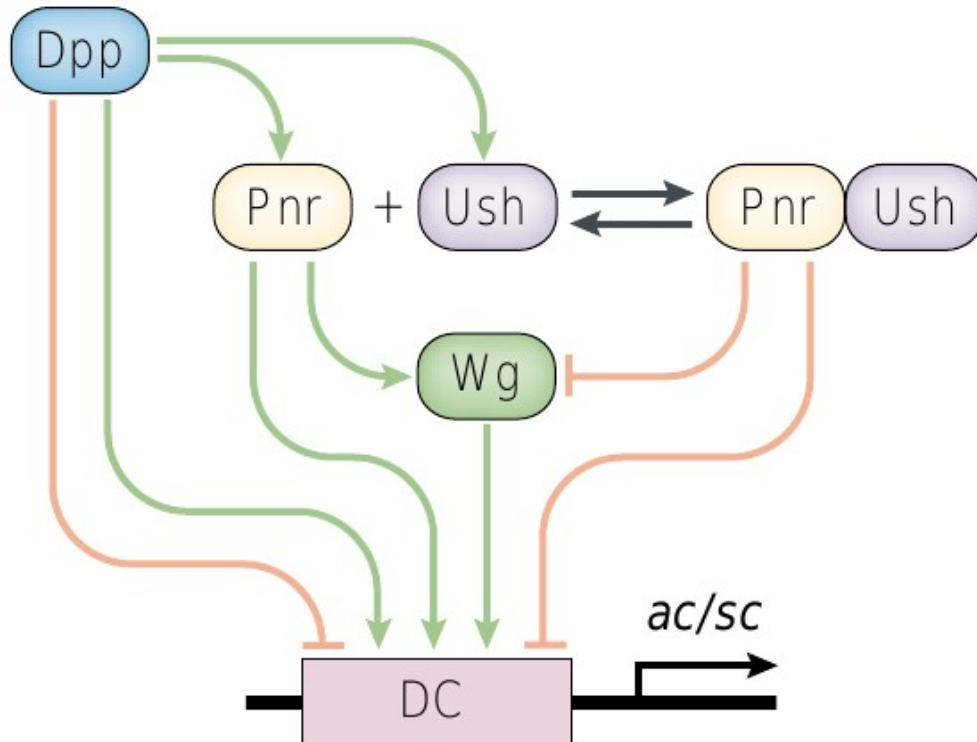
Marcellini 2006 - PloS Biology



scute cis-regulatory elements are “master switches”



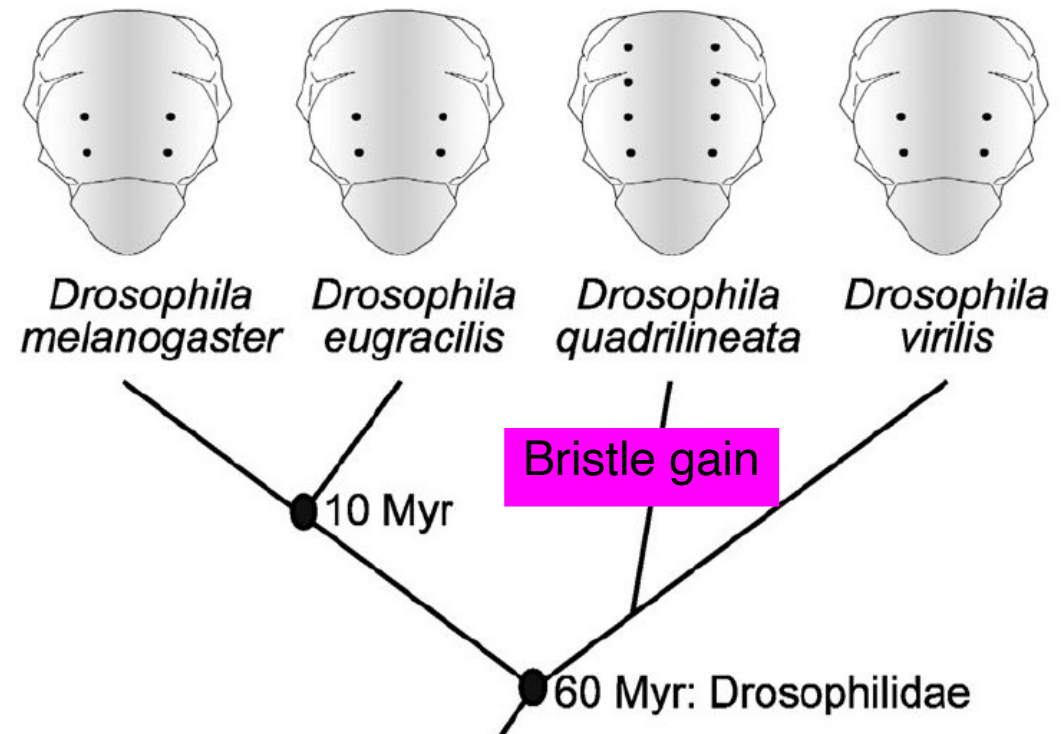
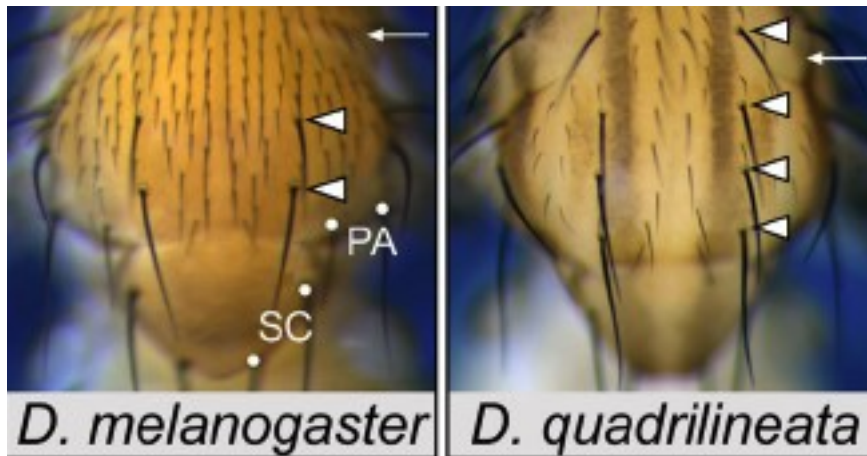
Simpson 2007



- Iro-C
- Bar
- Wg
- DC
- Pnr
- Dpp
- Ush

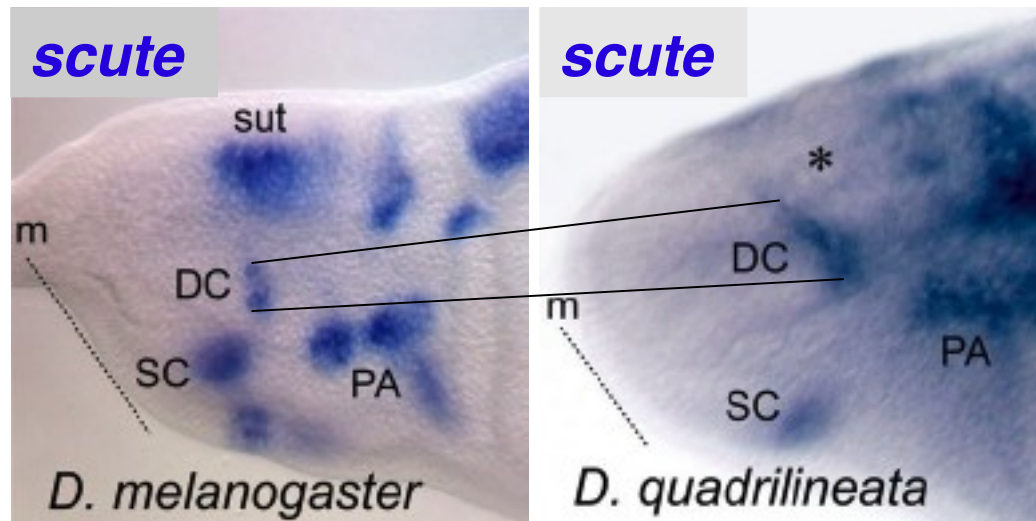
Gómez-Skarmeta 2003

Extra bristles in *D. quadrilineata*



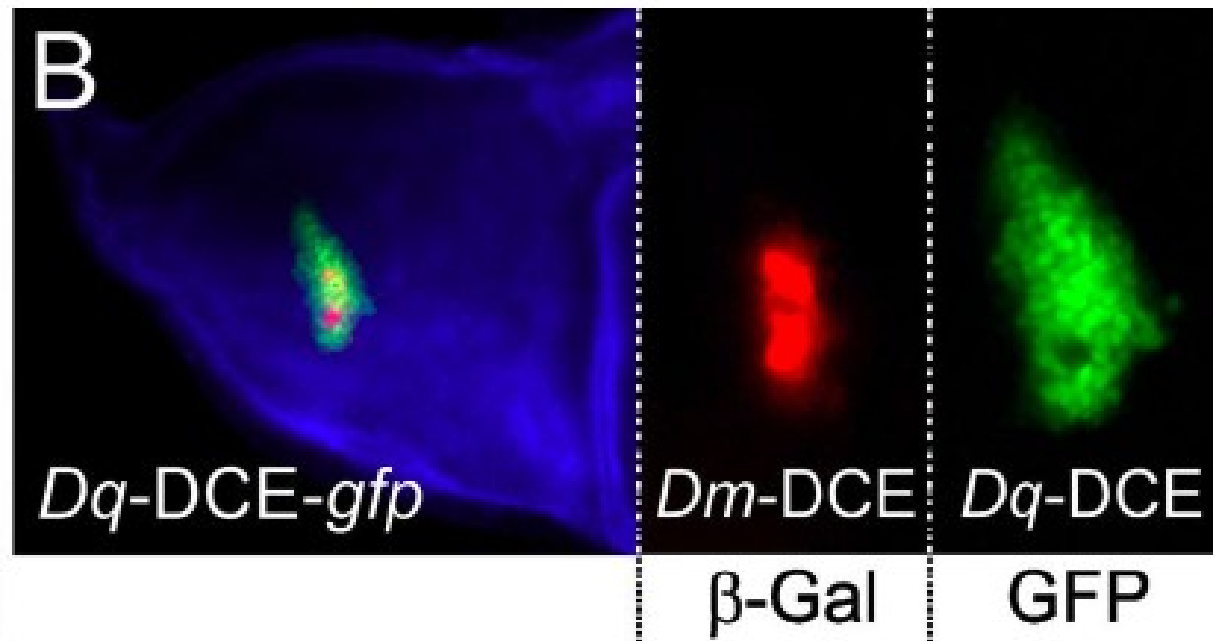
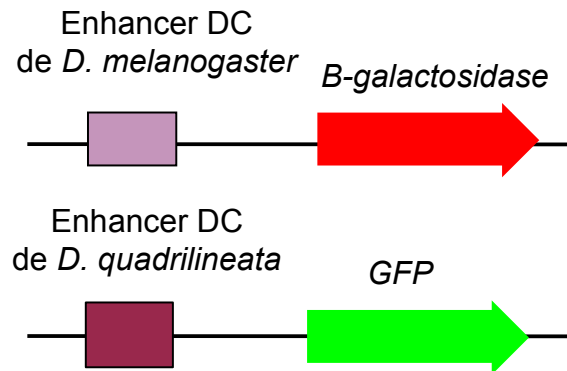
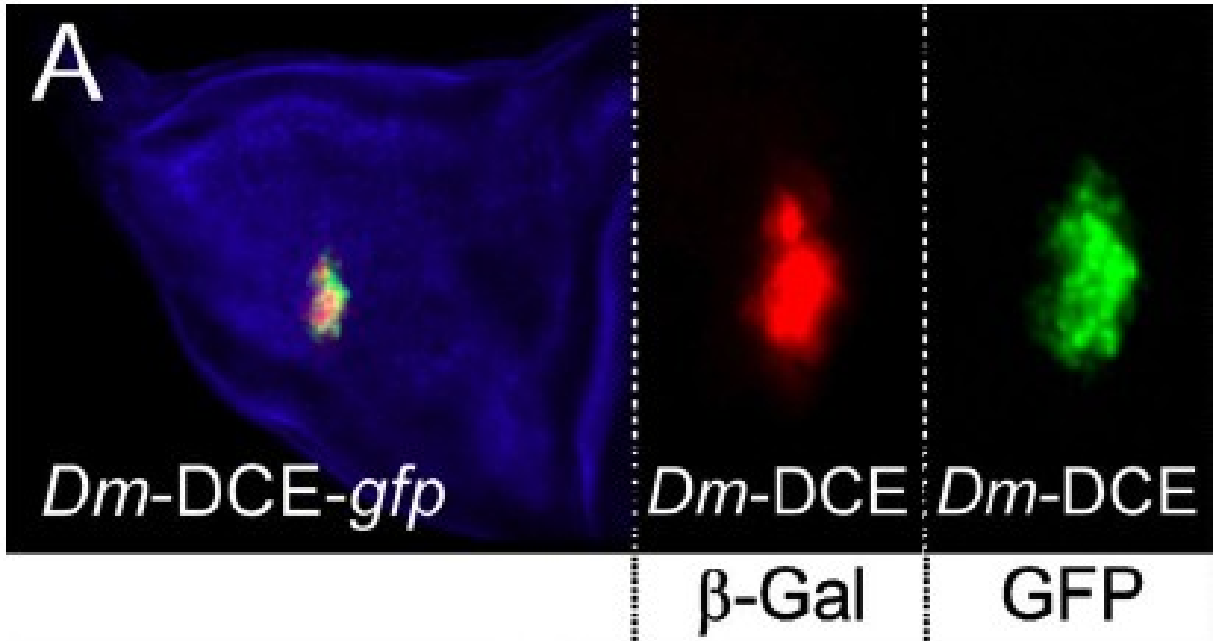
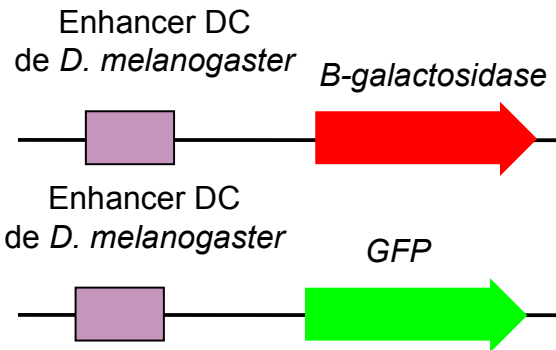
Extra bristles in *D. quadrilineata* correlate with larger *scute* expression domain

In situ hybridization



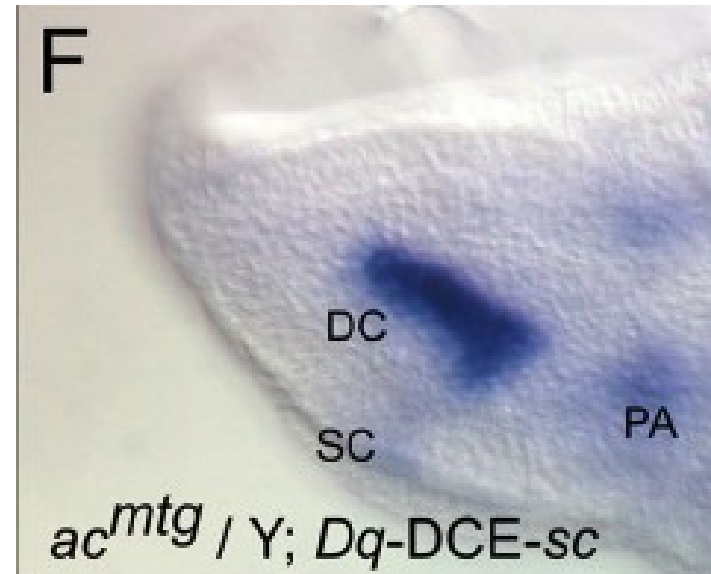
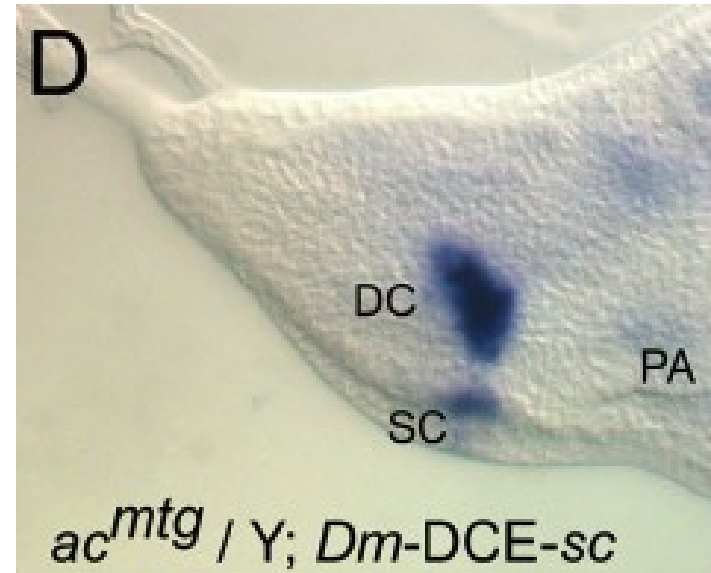
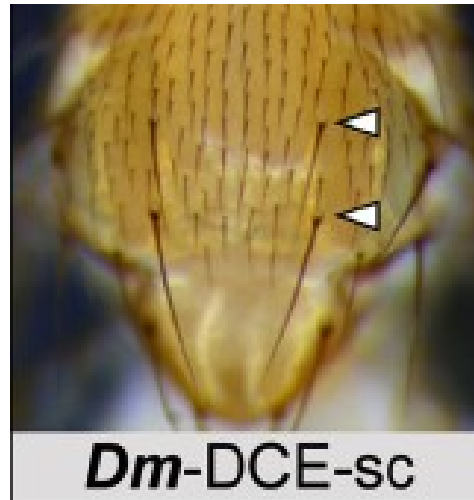
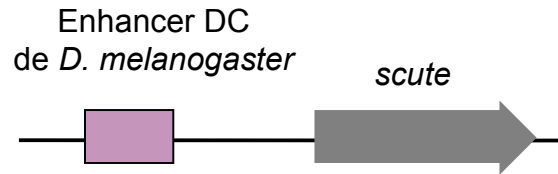
Test for a cis-regulatory change (1)

D. melanogaster
transgenics



Test for a cis-regulatory change (2)

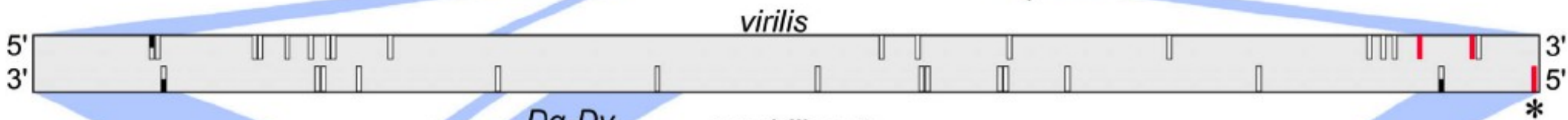
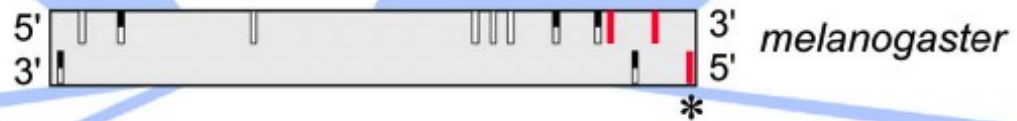
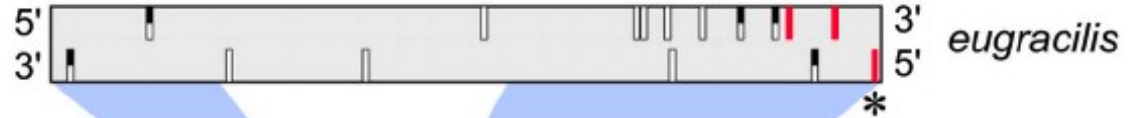
D. melanogaster
transgenics



Alignment of the DC region

A

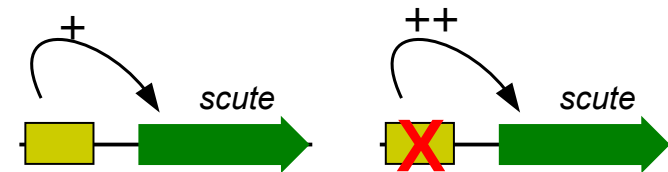
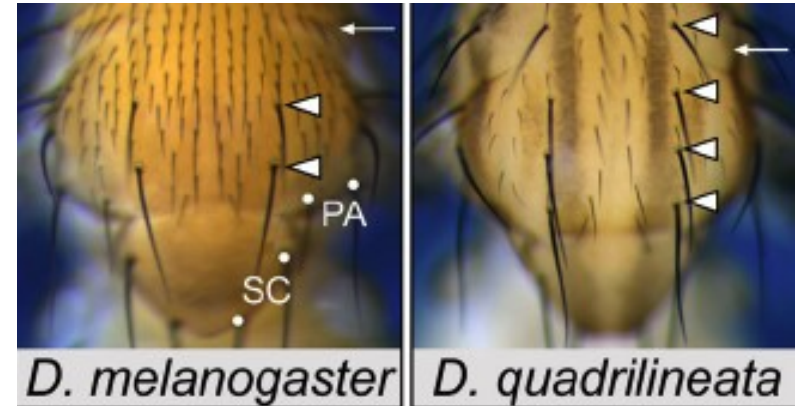
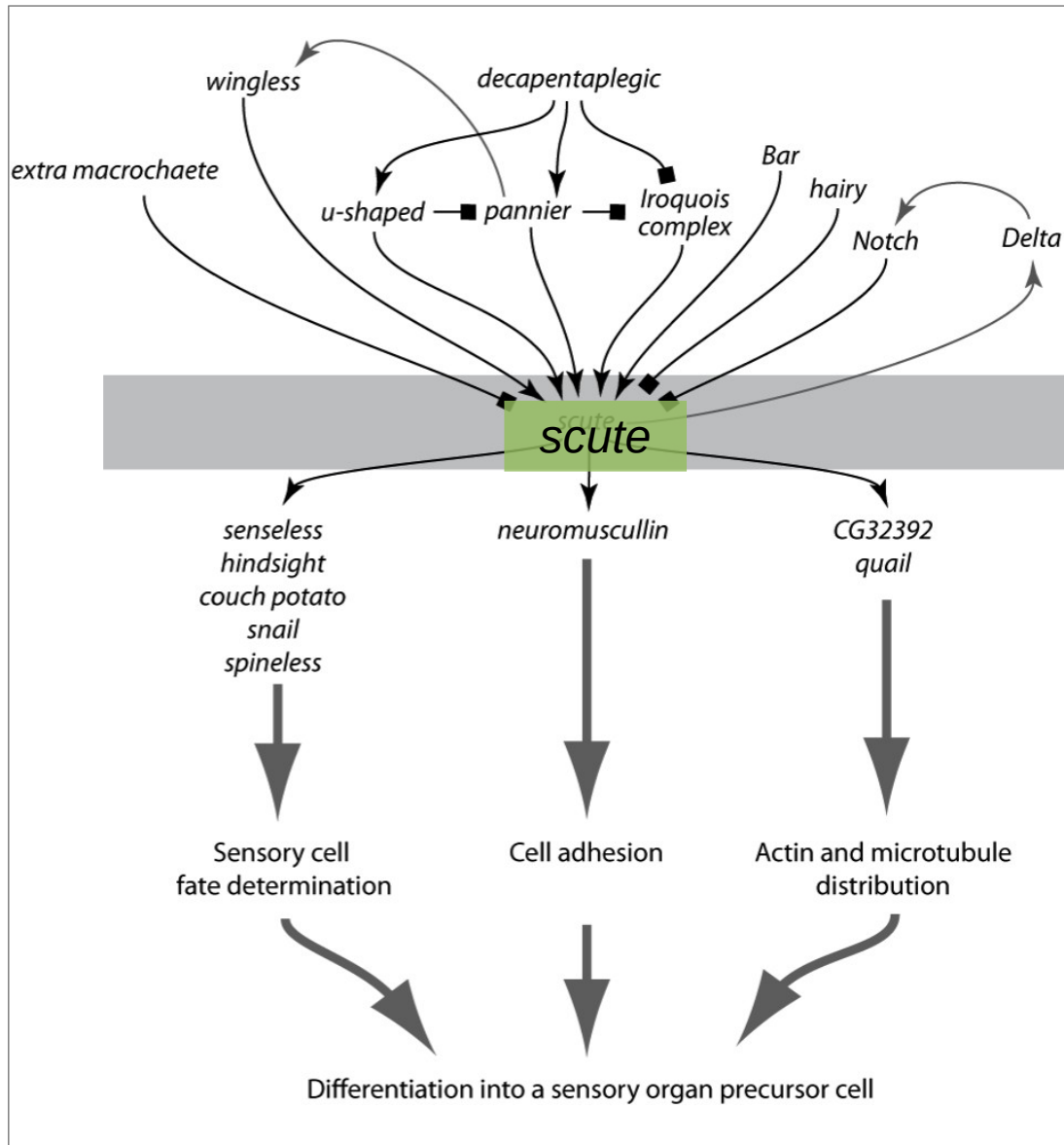
500 bp



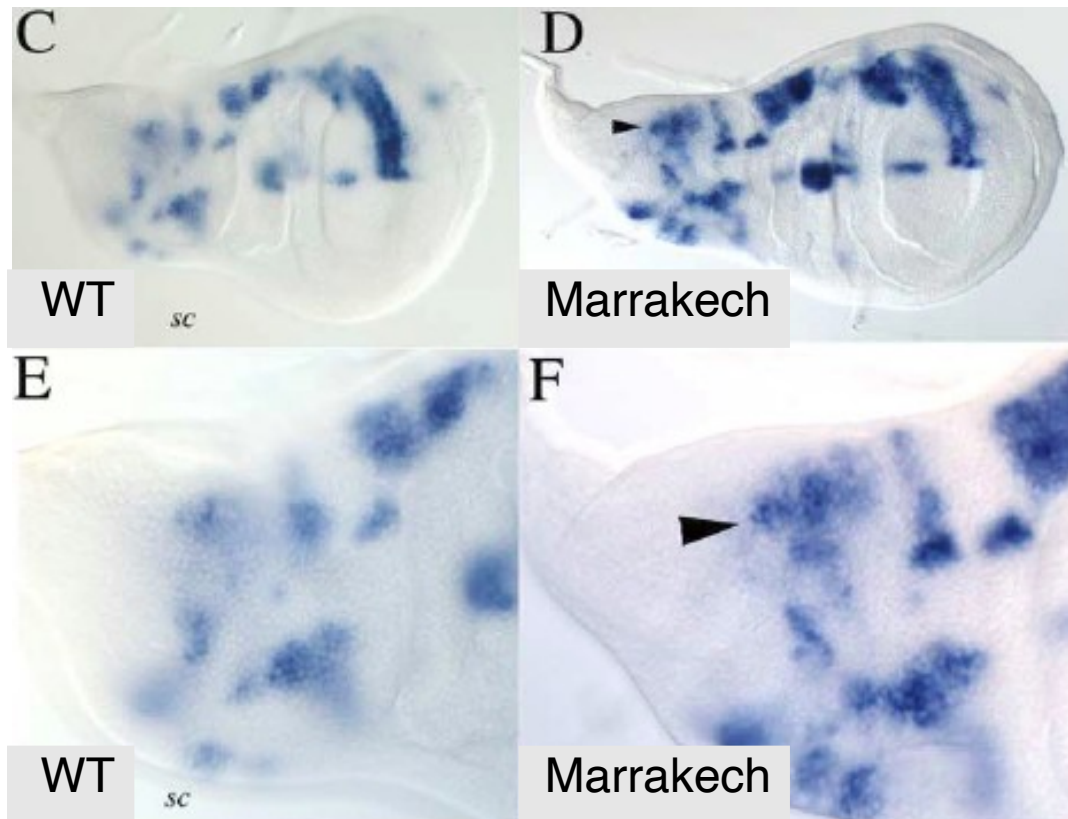
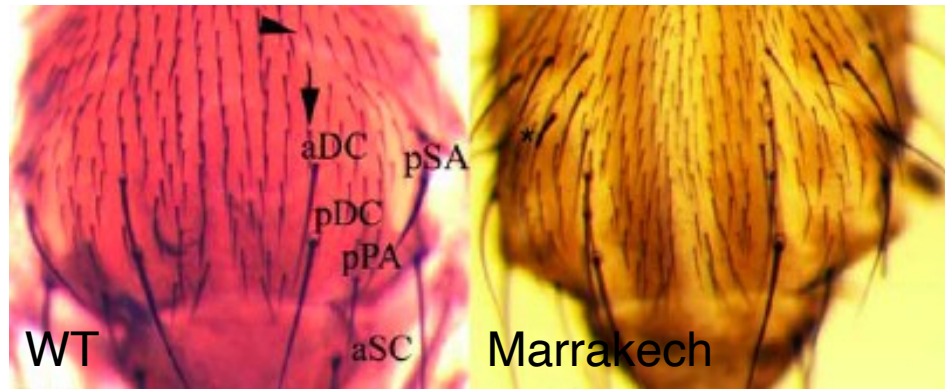
Dq-Dv *quadrilineata*



Genetic evolution is partly predictable

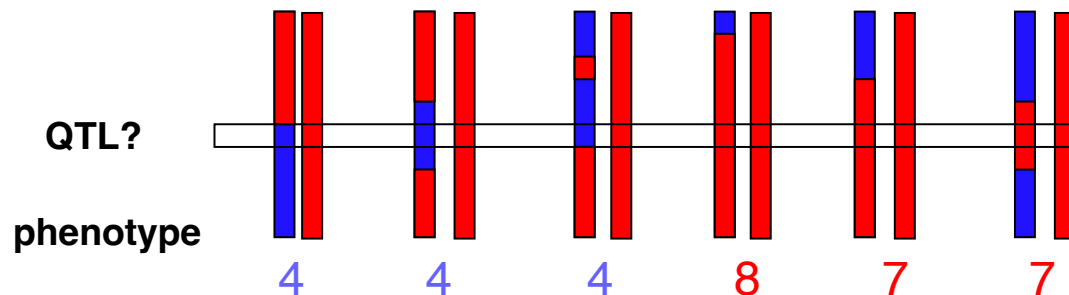
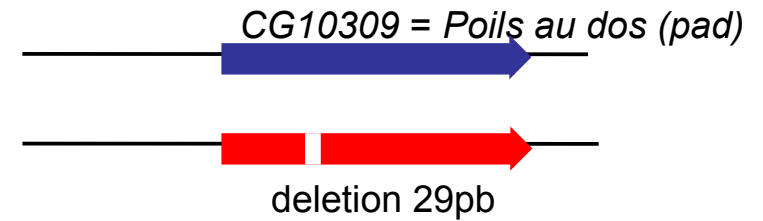
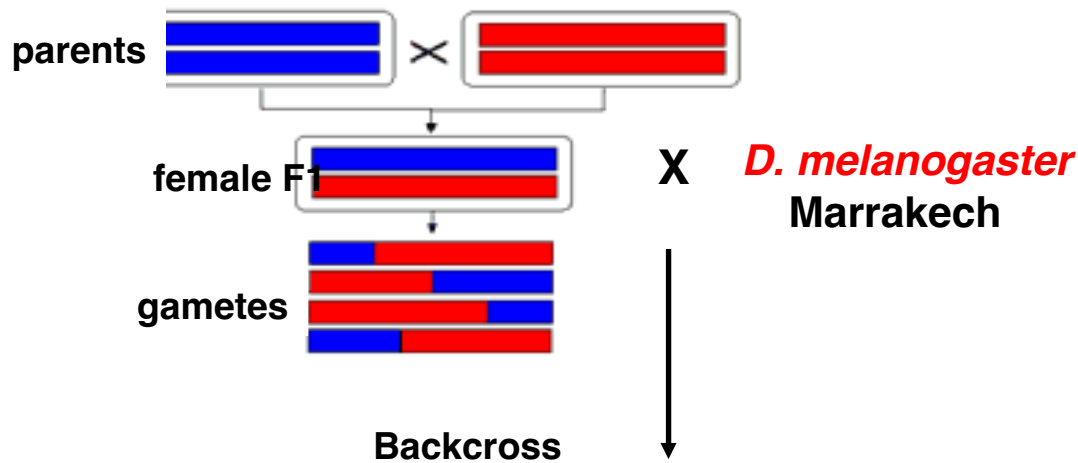
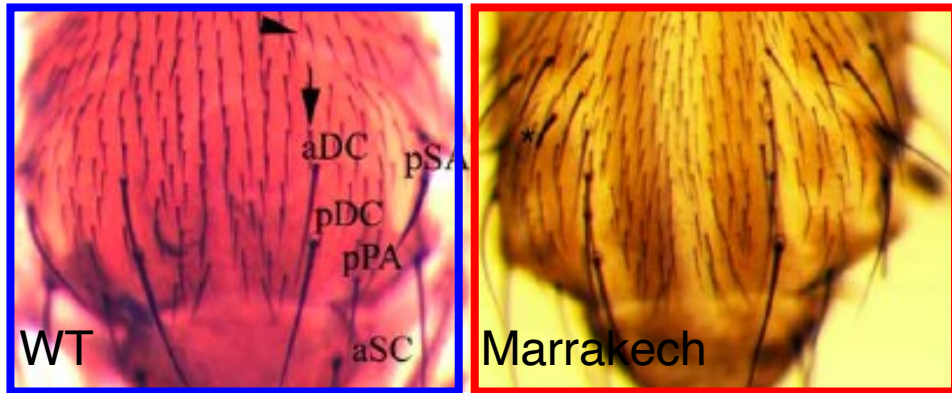


Extra bristles in *D. melanogaster*-Marrakech correlate with larger *scute* expression domain

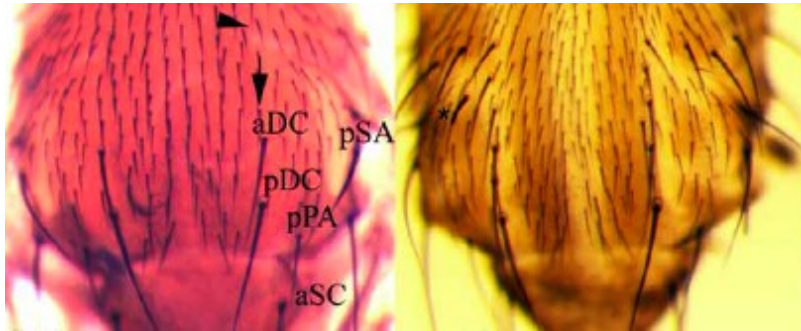


scute

Extra bristles in *D. melanogaster*-Marrakech due to mutation(s) in *poils-au-dos*

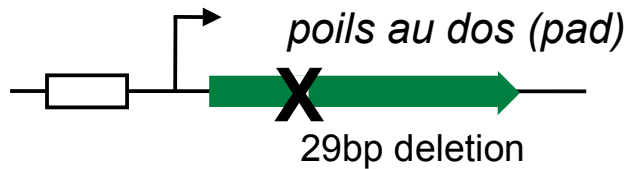


Short-term evolution...



D. melanogaster

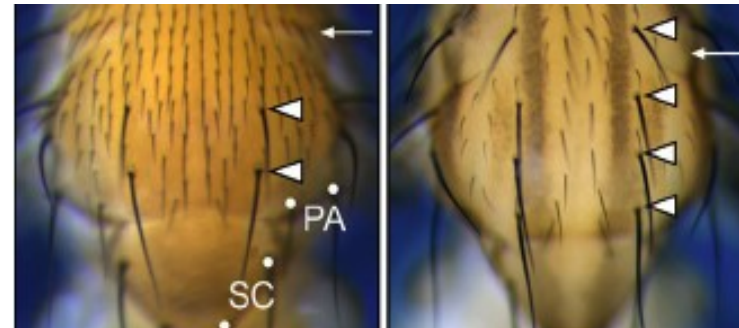
D. melanogaster variant



null mutation in coding region
change in thorax and wing

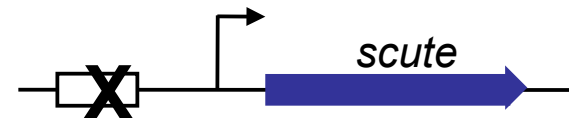
(Gibert et al., 2005)

...versus long-term evolution



D. melanogaster

D. quadrilineata



cis-regulatory mutation
change in the thorax only

(Marcellini et al. 2006)

Methods to identify the genes and the mutations responsible for phenotypic evolution



Various methods

Genetic

which chromosome (ex: autosomal versus sex)

QTL mapping

Genetic association studies

Complementation tests

General biology

General knowledge of the genes involved in the phenotype

Similarity with a known phenotype

Correlation with a change in gene expression level/pattern

Final test of protein activity

in vitro in *E. coli*, by transgenesis in the studied species or the closest model organism (ex: *beta-defensin* of dogs tested in mouse)

Final test of cis-regulatory regions

- with reporter constructs, transgenesis, comparison of both regions
- comparison of allele expression levels in hybrids (pyrosequencing)

Two types of approaches



no a priori, fewer bias
long and tedious
rarely ends with identification of the gene

only with strains/species which produce fertile
hybrids

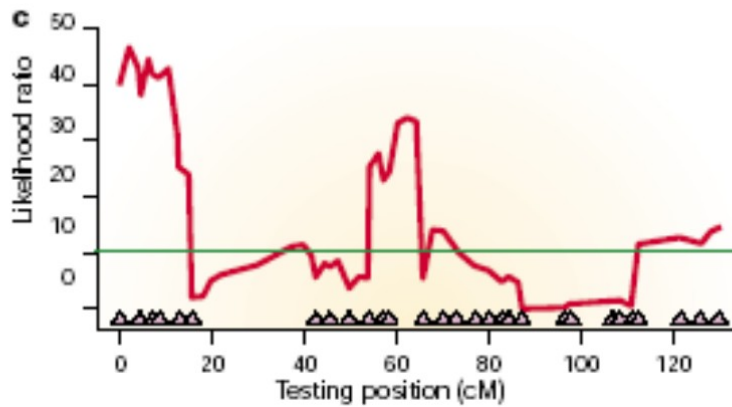
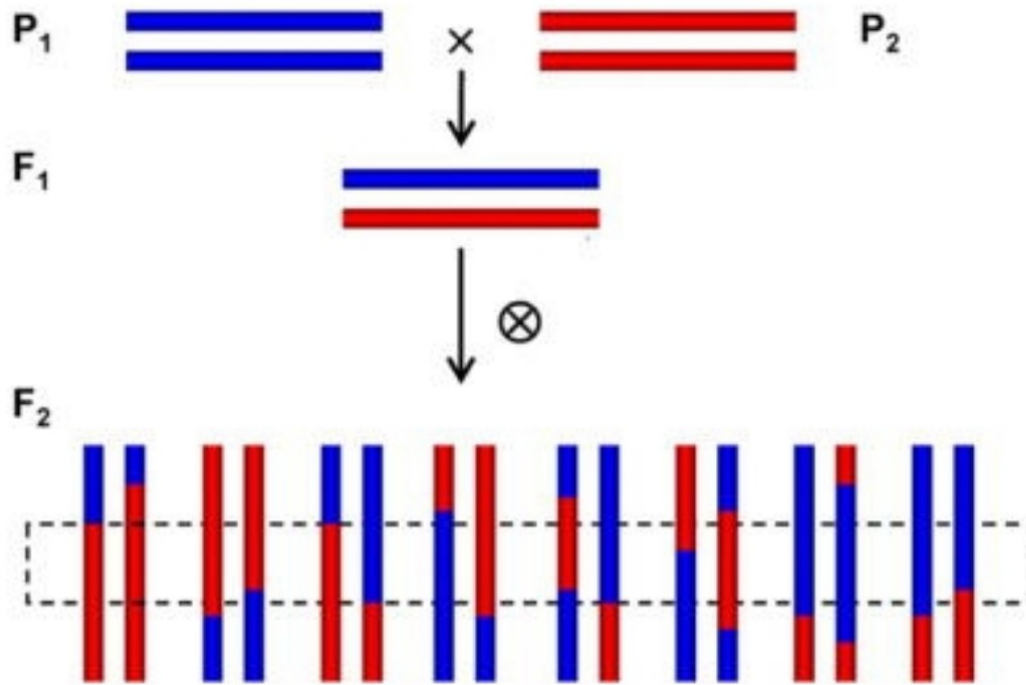
Based on an a priori idea
can be fast and efficient

will only find known genes

In both cases, genes with small effect are more difficult to identify

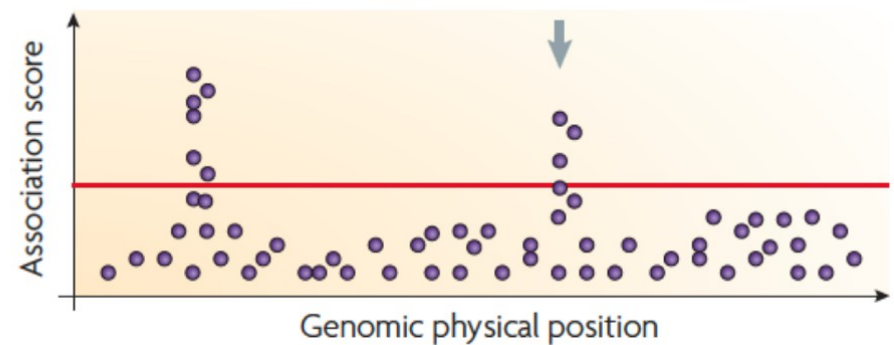
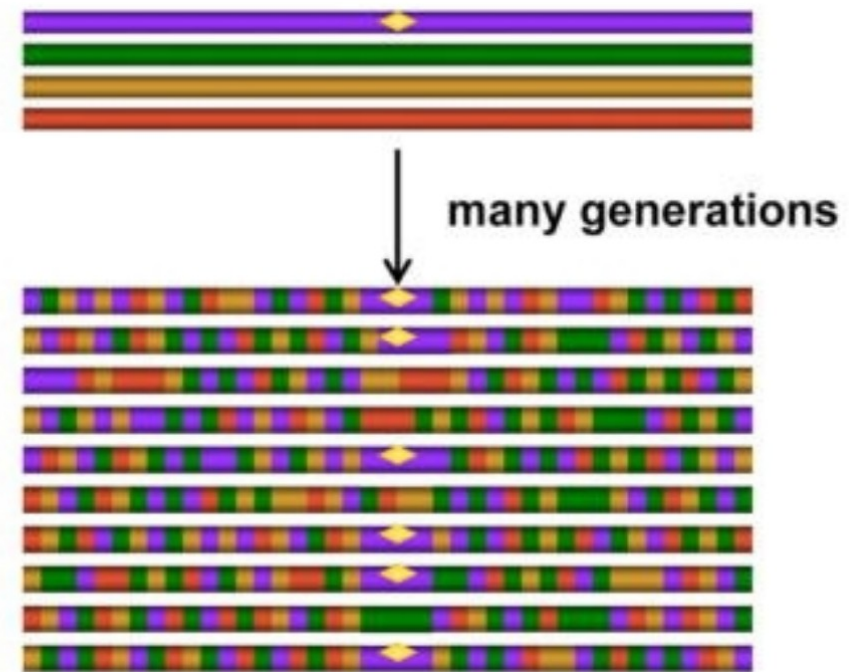
Linkage Mapping

Crosses in the lab

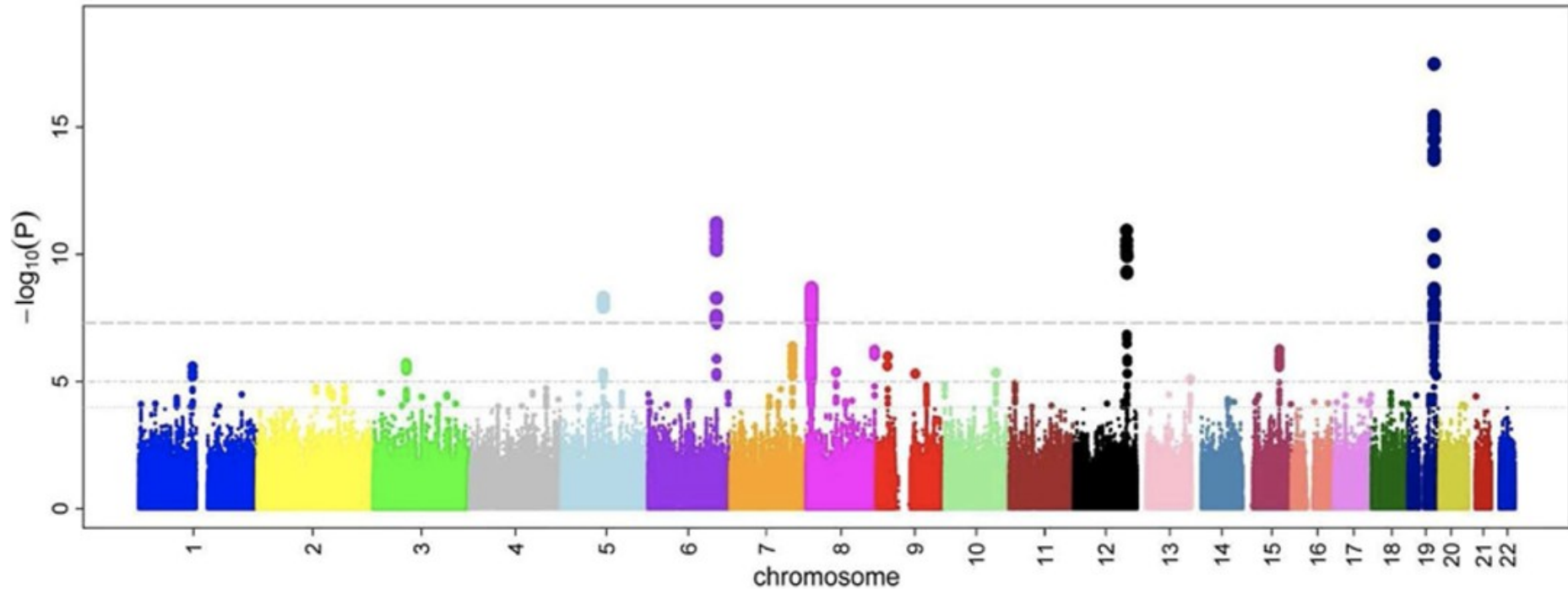


Association Mapping

Past crosses in natural populations

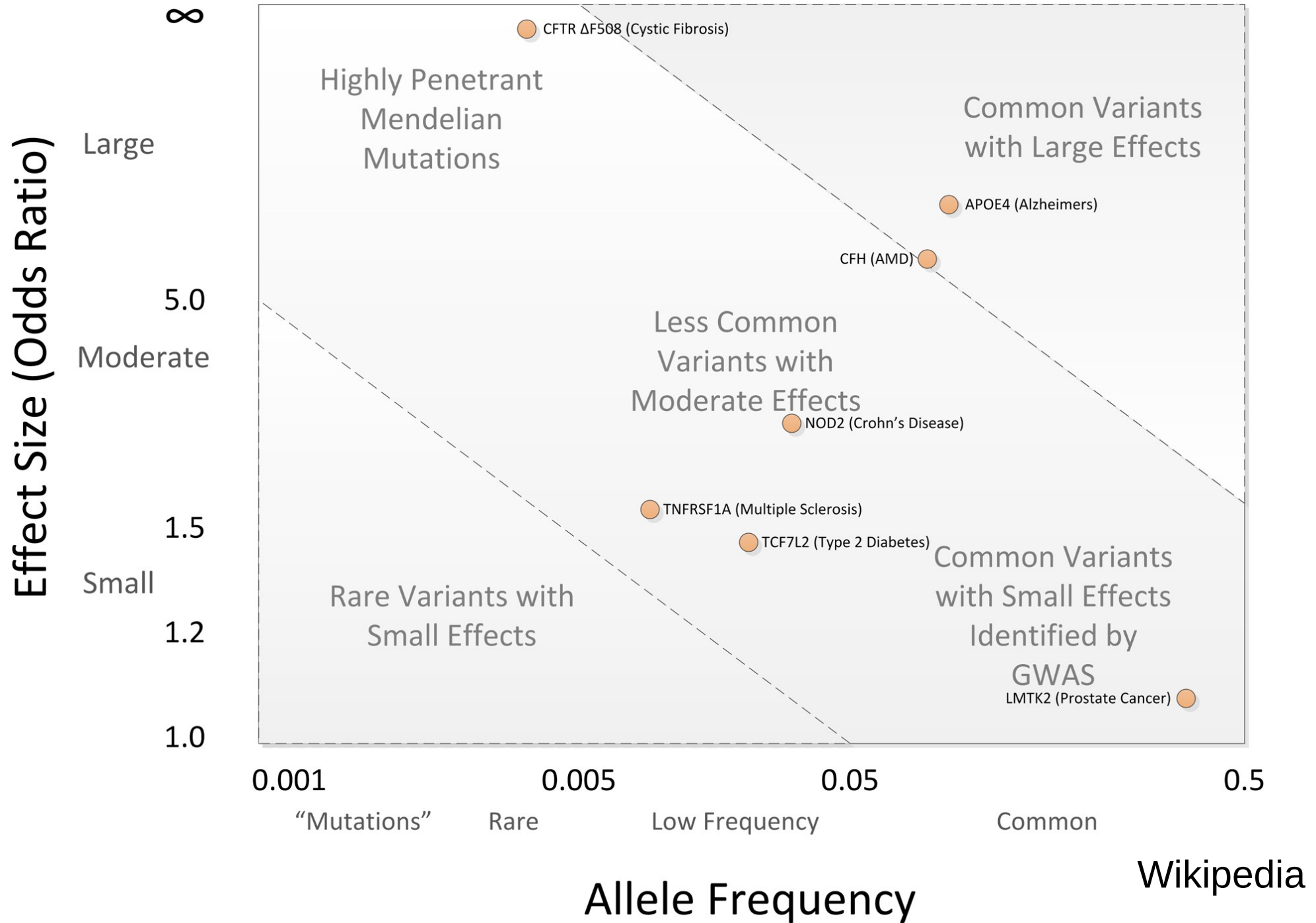


Genome-wide association study (GWAS)

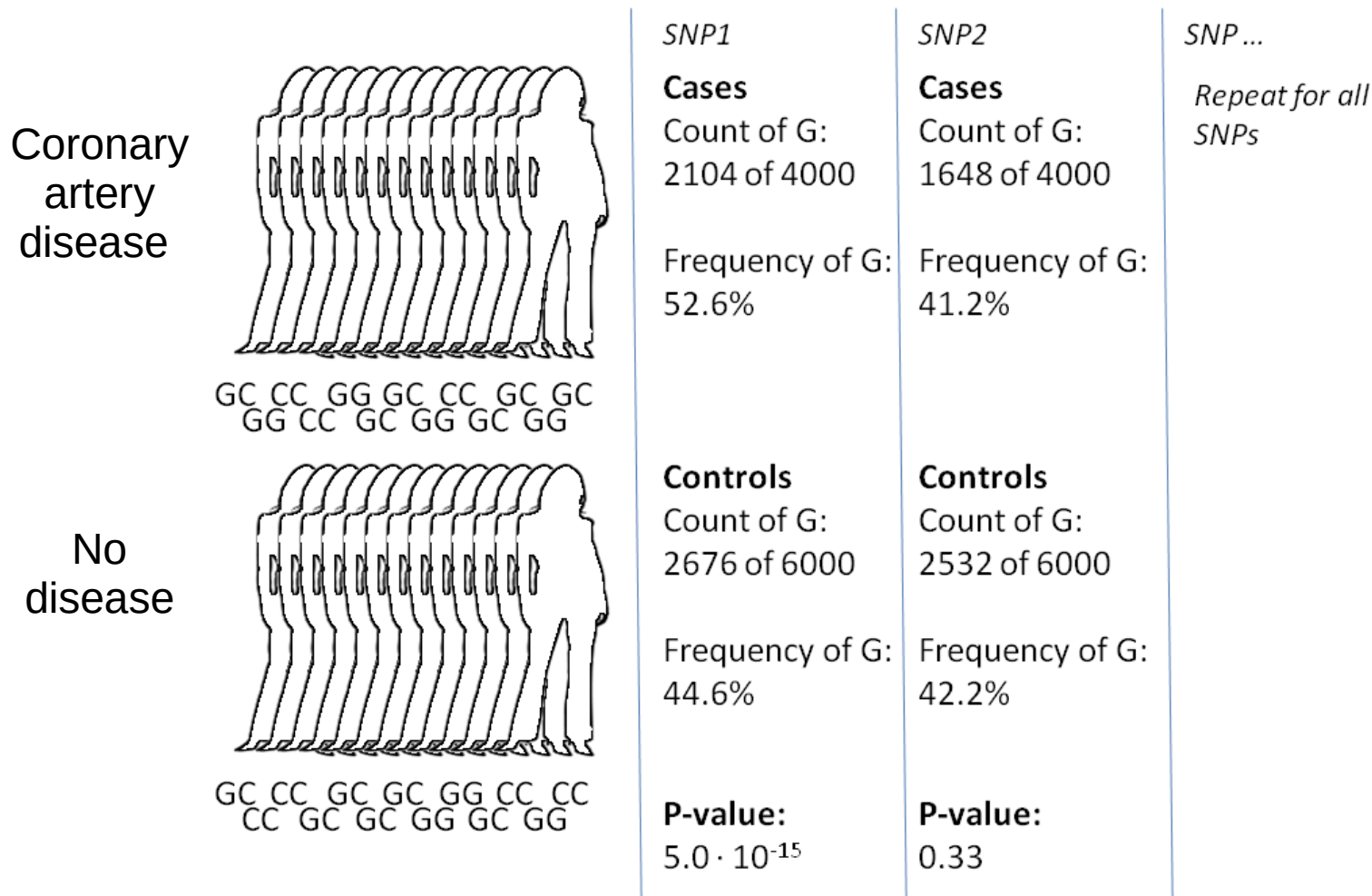


Manhattan plot depicting several strongly associated risk loci. Each dot represents a **SNP**, with the X-axis showing genomic location and Y-axis showing **association level**.

GWAS typically identify common alleles

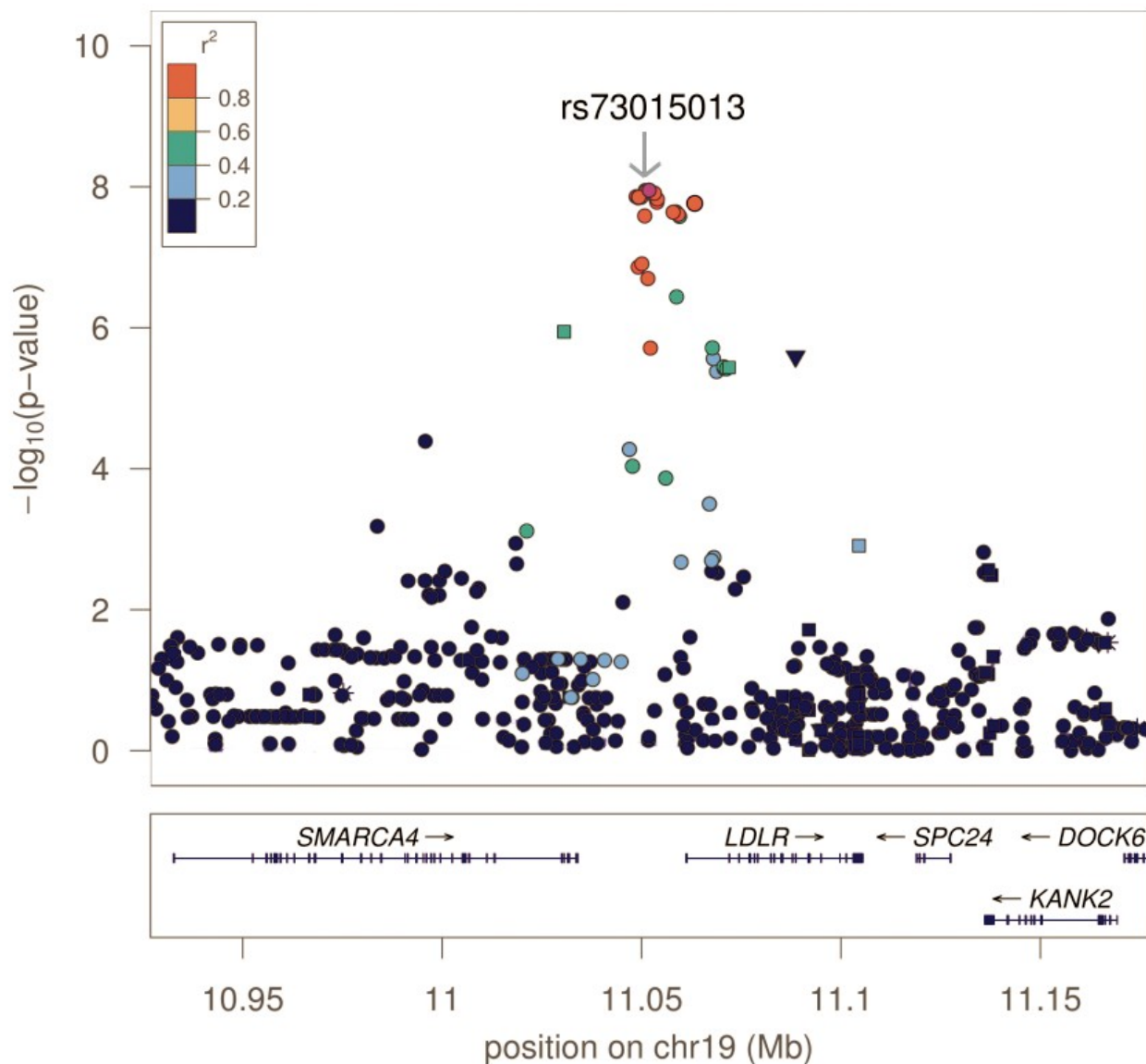


Methodology of a case-control GWA study



The allele count of each measured SNP is evaluated, in this case with a chi-squared test, to identify variants associated with the trait in question.

Regional association plot



Association to LDL-cholesterol levels.

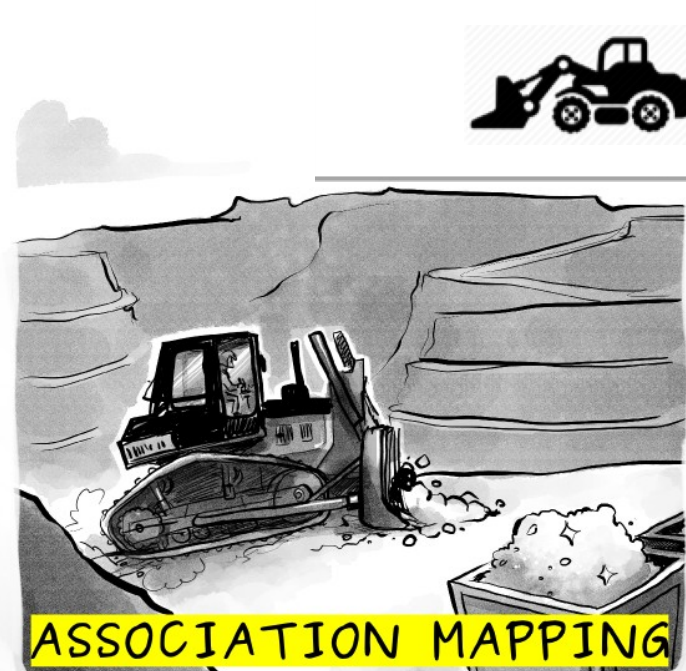
The haploblock structure is visualized with colour scale and the association level is given by the left Y-axis.

The dot representing the rs73015013 SNP (in the top-middle) has a high Y-axis location because this SNP explains some of the variation in LDL-cholesterol.

THREE APPROACHES to FIND the GOLDEN LOCI of EVOLUTION



REVERSE GENETICS
From genes to traits

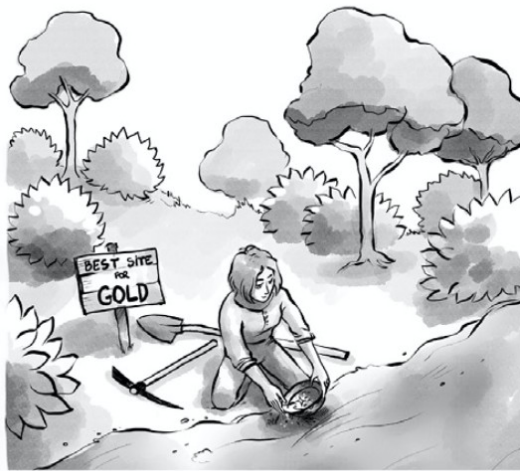


FORWARD GENETICS
From traits to genes

Little Ascertainment Bias, but Micro-Evolution only

Experimental Evidence

3 categories, each with biases



Candidate Gene

Reverse Genetics:

looking for sequence differences and trait effects based on previous studies of a given gene

Experimental Principle



Linkage Mapping

Forward Genetics:

trait mapping in hybrids obtained from laboratory crosses, using recombination over a few generations



Association Mapping

Forward Genetics:

statistical SNP/character state association in large cohorts, using recombination over many generations

Example

66 cases of color variation associated to *MC1R* coding mutations in vertebrates

F2 crosses between melanic and amelanic phenotypes in cavefish : identification of *MC1R* and *Oca2* alleles in distinct cave populations

GWAS of human pigmentation (skin, hair, eyes): identification and confirmation of causal variants at >15 genes including *Oca2* p.His615Arg in Eastern Asia

Ascertainment Bias on Locus Identification

High

Low to Intermediate

(depending on resolution / cross size)

Low

Molecular Type Bias

Favors identification of **coding mutations**

Little molecular bias

Can miss structural variants (short read genotyping)

Trait Type Bias

Favors traits with small molecular targets, large-effect size

Amenable to dissection of complex traits with small-effect size (large crosses, multiparental families)

Most common approach for complex traits with small-effect size

Taxonomic Breadth

Large

Narrow, limited to interfertile lineages (populations or sister species)

Very narrow, limited to polymorphic or intermixing populations

QTL Mapping

4 steps: crosses, genotyping, phenotyping, statistical analysis

Crosses

Backcross with one line
Backcross in both directions
F2
Crosses for several generations
Introgression lines
Recombinant Inbred Lines
...

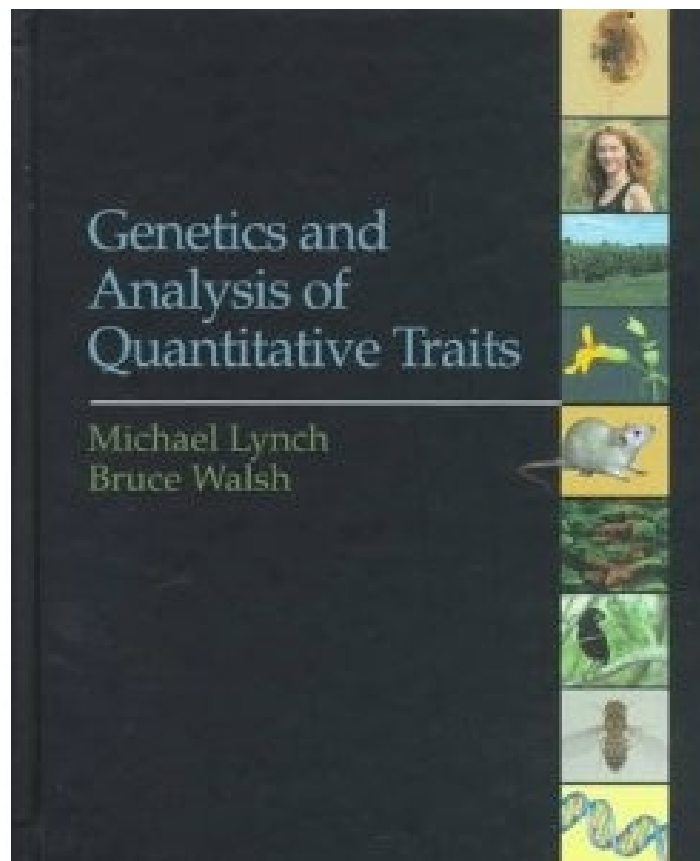
Always try to maximize the number of recombination events

Markers

yes-no PCR
PCR length polymorphism
Pyrosequencing
Probe hybridization
Microarray
RADseq
High-throughput sequencing

How many markers?

theory



practice

