Mapping the genes underlying phenotypic changes of interest

Virginie Courtier-Orgogozo Institut Jacques Monod, Paris

Tomato shape



Lycopericon 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 Dmin/Dmax for 10 fruits per plant L/D : L= length, D = diameter at equator Dmin/Dmax



82 molecular markers on the 12 tomato chromosomes



One major locus near marker TG645



Two main files

Markers file

-start	
-Chromosome	1
CF5475	0.4
CF5573	24.7
СТ7895	41.0
СТ8903	59.0
CF5613	67.7
СТ7892	76.0
СТ890	89.0
СТ233	39.0
Telomere	50.0
-Chromosome	2
CF5671	0
CF5675	10.4
CF5673	34.7
СТ789	41.0
СТ789 СТ890	41.0 89.0
CT789 CT890 CT567	41.0 89.0 115.0
CT789 CT890 CT567 Telomere	41.0 89.0 115.0 130.0

Genotypes and phenotype(s) file

-start	t :	inc	liv	7ic	lua	als	s r	naı	cke	ers	5					
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	1	1	1	2	2	1	1	1		
Ind_5	0	1	0	0	0	0	1	1	1	1	2	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	iı	nd	LV	Ĺdι	ıa	Ls	ma	ar}	دeı	cs						
-start	5 3	inc	liv	7ic	dua	als	5 t	ra	ait	S	1	L	ove	erD	name	d
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.xi + ε xi = 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 Ho: b = 0 versus H1: b = estimated b

Likelihood ratio test statistic

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

The probability distribution of the test statistic can be approximated by a chi-square distribution with (df1 – df2) degrees of freedom, where df1 and df2 are the degrees of freedom of models 1 and 2 respectively

Interval mapping

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

xi = indicator variable specifying the probabilities of an individual being in different genotypes for the tested position, constructed by flanking makers xi = 0 if Le/Le, = 1 if Le/Lp, = 2 if Lp/Lp

a,b = best fit parameters (maximum likelihood)

Test Ho: b=0 versus H1: b=estimated b



Interval mapping

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

xi = indicator variable specifying the probabilities of an individual being in different genotypes for the tested position, constructed by flanking makers xi = 0 if Le/Le, = 1 if Le/Lp, = 2 if Lp/Lp

a,b = best fit parameters (maximum likelihood)

Test Ho: b=0 versus H1: b=estimated b

Composite Interval mapping

L/D of individual i = a + b.xi + c.xi + e

xi = indicator variable specifying the probabilities of an individual being in different genotypes for the tested position, constructed by flanking makers xi = 0 if Le/Le, = 1 if Le/Lp, = 2 if Lp/Lp

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

LOD score

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

Test Ho: b = 0 versus H1: b = estimated b

Lo = pr (data | no QTL) – phenotypes assumed to follow a normal distribution L1 = pr (data | QTL at tested position)

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

Significance threshold

10,000 permutations of phenotype/genotype data

- $\rightarrow\,$ random distribution of LOD scores
 - \rightarrow 1% or 5% significance threshold



One major locus near marker TG645







Sequencing of the region in the 2 tomato varieties

L. esculentum cv. Yellow Pear



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

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



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



⁽Shapiro et al., 2004)

Quantitative measurement of the phenotype



1000 microsatellite markers

26 linkage groups



One major locus at the end of linkage group 7



Major locus responsible for 65% of the variance

One major locus at the end of linkage group 7 A few minor loci





One major locus at the end of linkage A few minor loci group 7 Pitx1 Linkage group 1 Linkage group 7 3 b Linkage group 2 C, 5 8 4 LOD score 6 80 Pelvic spine 3 Pelvic girdle 70 2 Ascending branch Tbx4 \mathbf{Z} Asymmetry [L/(R+L)] -0---

0

0

20

40

60

Map position (cM)

80

10 20 30 40 50 60 70 80 90 100 0 100 Map position (cM)



60

50

40

30

20

10

0

0

20

60

40

Map position (cM)

80

Pitx1, responsible for the phenotypic change?

Pitx1 null mutations in mice (pelvis reduction, stronger on right side)

QTL mapping

Same coding sequence in lake and marine forms

Pitx1 expressed at stage 29 in marine individuals but not in marine individuals

Pitx1, responsible for the phenotypic change?

Pitx1 null mutations in mice (pelvis reduction, stronger on right side)

QTL mapping

Same coding sequence in lake and marine forms

Pitx1 expressed at stage 29 in marine individuals but not in marine individuals

BUT

The decrease in *Pitx1* expression levels might have evolved due to mutations in an upstream regulatory gene

Comparison of allele expression in FRIL (with pelvis) Х Х hybrids LITC PAXB (with pelvis) (no pelvis) F₁ larvae (all with pelvis) Pitx1 PAXB Pitx1 LITC Pitx1 Allele-Specific Pitx1 FRIL

Expression

Pyrosequencing





Test of *Pitx1* cisregulatory regions





Rescue of a pelvis in freshwater individuals





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



Pitx1 is in a fragile DNA region



Control = artificial chromosome without test region

(Xie et al., 2019)

Evolution of extra bristles

Interspecific change in *D. quadrilineata*

D. melanogaster D. quadrilineata

Intraspecific change in *D. melanogaster*



Marcellini and Simpson 2005, Gilbert et al. 2005

Finding genetic rules on bristle evolution



Randsholt and Santamaria 2008



color, type, orientation shape and size presence/absence

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

CRE mutations in *achaete-scute*

Stern and Orgogozo 2009

Bristle development









scute cis-regulatory elements are "master switches"







Simpson 2007



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



Genetic evolution is partly predictable





Stern and Orgogozo 2009 Science

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



Gibert et al 2005

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



Gibert et al 2005

Short-term evolution...



...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

Two types of approaches



no a priori, fewer bias long and tedious rarely ends with identification of the gene Based on an a priori idea can be fast and efficient

only with strains/species which produce fertile hybrids

will only find known genes

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

Various methods

<u>Genetic</u> which chromosome (ex: autosomal versus sex) QTL mapping Genetic association studies Complementation tests

<u>General biology</u> 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)

Linkage Mapping

Crosses in the lab

Association Mapping

Past crosses in natural 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





www.Gephebase.org >2000 entries



Suggest an Article





The Database of Genotype-Phenotype Relationships Search Gephebase for genes, phenotypes, taxa, mutations, a

Gephebase compiles genotype-phenotype relationships, i.e. associations between a mutation and a phenotypic variation. Gephebase consolidates data from the scientific literature about the genes and the mutations responsible for phenotypic variation in Eukaryotes (mostly animals, yeasts and plants). We plan to include non Eukaryote species in the future. For now, genes responsible for human disease and for aberrant mutant phenotypes in laboratory model organisms are excluded and can be found in other databases (OMIM, OMIA, FlyBase, etc.). QTL mapping studies that did not identify single genes are not included in Gephebase.

If you use Gephebase for your publication, please cite: Martin, A., & Orgogozo, V. (2013). The loci of repeated evolution: a catalog of genetic hotspots of phenotypic variation. Evolution, 67(5), 1235-1250.

Conference on Gephebase and the loci of evolution (Paris, 2016)

You can retrieve data via HTTP requests through APIs. Below is the list of available APIs. By default, the response data is sent in xml format. For each field, it is possible to only enter a subset of a keyword and still be able to successfully retrieve the desired data. (example: "Bir" for "Birds" will display all data that have the characters "Bir").

Orgogozo et al Nucleic Acid Research 2019

Wrinkled seed: TE insertion (Bhattacharyya 1990)



myostatin coding region (Grobet 1997)



Mc1r coding region (Eizirik 2003)



OVATE coding region (Liu 2002)



FRI coding region (Johanson 2000)



luciferase coding region (Stolz 2003)



anthocyanin-2 coding region (Quattrocchio 1999)



THREE APPROACHES to FIND the GOLDEN LOCI of EVOLUTION



FORWARD GENETICS

From traits to genes Little Ascertainment Bias, but Micro-Evolution only

REVERSE GENETICS From genes to traits Experimental Evidence

3 categories, each with biases



Candidate Gene

Experimental
Principle

Example

Ascertainment Bias on Locus Identification

Molecular Type Bias

Trait Type Bias

Taxonomic Breadth

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

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

High

Favors identification of coding mutations

Favors traits with small molecular targets, large-effect size

Large



Linkage Mapping

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

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

Low to Intermediate (depending on resolution / cross size)

Little molecular bias

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

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



Association Mapping

Forward Genetics:

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

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

Low

Can miss structural variants (short read genotyping)

Most common approach for complex traits with small-effect size

Very narrow, limited to polymorphic or intermixing populations

The wrinkled pea



wrinkled (recessive):

altered starch structure (physiology) and wrinkled seed aspect (morphology)

insertion of a 800 bp TE in the coding region of SBEI = glycosyl hydrolase enzyme gene

probably disrupts the last 61 amino acids of the SBEI protein presumptive null mutation



Class I, retrotransposon mammalian L1



Class II, DNA transposon

Sleeping Beauty (reconstructed from fish) Frog Prince (reconstructed from frogs) Hsmar1 (reconstructed from human) Minos (Drosophila hydei) Tol2 (Oryzias latipes) piggyBac (Trichoplusia ni)



Meta-analysis of the role of transposable elements in phenotypic evolution

It's your turn!

https://tinyurl.com/yyod3zqx