

Trait/DNA-marker associations (in polyploids)

Isabel Roldán-Ruiz August 2012

Institute for Agricultural and Fisheries Research



Plant Sciences Unit www.ilvo.vlaanderen.be Agriculture and Fisheries Policy Area



Contents

- 1. Defining the target
- 2. Identifying suitable germplasm
- 3. Mapping population structures
- 4. Phenotypic evaluation
- 5. Genotyping (as discussed in lecture 1)
- 6. Strategies to identifying marker/trait associations
 - Linkage analysis (in polyploids)
 - Bulked Segregant Analysis
 - QTL analysis
 - Association Mapping
- 7. Developing markers for application marker validation

Only at this step the markers are tested in breeding plant materials





Introduction

Remember from Lecture 1:

- DNA-marker technologies offer the possibility to screen plant genomes for sequence polymorphisms
- DNA-marker polymorphisms help us to define specific positions (loci) in the genome
 - \Rightarrow DNA-markers make it possible to follow chromosome pieces from one generation to the next
- How to identify marker/trait associations? How do we identify DNA-markers which are located in or in the neighborhood of the genes of interest?
 - In most cases => construction, phenotyping and genotyping (using DNA-markers) of dedicated populations (= mapping populations)
 - But other possibilities exist; e.g. association mapping





1. Defining the target

- Is the trait of importance to breeding program or to biological research?
- Is a DNA-marker needed?
 - What is the cost of a bioassay relative to a marker assay?
 - Financial aspects / time constraints (can we save one generation of crossing by using DNA-markers?)
 - Is the trait dominant or recessive?
 - Recessive traits may be hard to identify in a bioassay => DNA-marker development is then advantageous
 - Perhaps no alternative to marker development
 - Quarantaine trait (the disease is not present in the country in which breeding is carried out) => no bioassay possible
 - Pyramiding disease resistances to avoid resistance breakdown
 - Map-based cloning of genes high resolution map required
 - Gene deployment where desirable alleles are available for several loci, but only one is really needed. How does one decide on the best one to use?





2. Identifying suitable germplasm

CHOICE OF CONTRASTING PARENTS

We can only find markers linked to a trait of interest if we compare plants with and without the trait

- Among available germplasm search for plants with and without the trait
- Or find an existing mapping population that segregates for the trait
- Remember: only traits that segregate in the population can be mapped

Example: to identify genomic loci linked to <u>plant height</u>, which plants would you select as parents of the mapping population? Could you also map <u>seed production</u> in this mapping population? And <u>disease resistance</u>?

	height	seed production	disease resistance
Plant 1	100	200	++++
Plant 2	105	210	++++
Plant 3	150	205	+
Plant 4	95	200	+++





3. Population structures

DESIGN OF STUDY

- Knowledge on genetics
 - Monogenic / multigenic trait
 - Heritability

Marker development is specially advantageous for complex traits with low heritability in the global population (not specifically in the population used for mapping)

- Decisions about population structure
 - Double haploids (DH)
 - F₂
 - Recombinant inbred lines (RILs) created by single seed descent
 - Near Isogenic Lines (NILs) created by back-crosssing
 - Pair-cross between highly heterozygous parents (CP cross pollinator- population)
- Population size (how many progeny plants do we need?)
 - The higher the population size the higher the 'resolution' of the mapping study
 - To determine chromosome location of a single gene => 50 F_2 is appropriate
 - Map-based cloning => over 1000 progeny plants required
 - Numbers depend on genome size and complexity

Sweetpotato - only one possibility: pair-cross between heterozygous







4. Phenotypic characterization

- Phenotypic evaluation should be possible for single plants
- Quality, yield and traits of low heritability (strong environmental influences) => many replicates and /or multiple field trials
- Association mapping => not based on mapping populations; phenotypic information can be collected from existing programs





Qualitative and quantitative traits

DEFINITION

- From the point of view of the molecular geneticist (interested in the identification of DNA-markers linked to genes involved in agronomic traits) there is one question of special relevance:
 - Some traits (e.g. incompatibility, some disease resistances) are <u>controlled by one or</u> <u>two genes with large effects</u>
 - Differences in expression (phenotype) are 'qualitative'
 - You can say that the traits follow a 'Mendelian' inheritance
 - Most traits of agronomic relevance such as growth or yield are complex, being <u>under</u> the control of a number of genes as well as the environment
 - Differences in expression (phenotype) are '**quantitative**'. No clear discontinuity exists between the phenotype of different genotypes, giving the impression of a continuous distribution

Knowledge about the kind of trait we are dealing with is of crucial relevance for the design of the experimental strategy





2808 wild *Malus* seedlings from 310 populations from Kazakhstan, Russia, China & Turkey Phil Fosline, PGRU, USA





Field for phenotypic evaluation of a mapping family

Qualitative and quantitative traits

DEFINITION

Monogenic trait

- Remember: Mendel studied major gene differences in the garden pea: color of the seed (one locus with alleles Y and g, Y>g)
- But: some factors can affect the expression of Mendel's laws



If we analyze the diversity in a collection of plants for this character, we are only looking at the genetic diversity contained in this genetic locus



Multigenic trait

- Multiple genes (and environment) affect the expression of the trait
- The expression in the population is a 'bell-shaped' curve: there are many genotypes and there are no clear phenotypic differences among them



If we analyze the diversity in a collection of plants for this character, we are looking at the genetic diversity contained in all the genetic loci which influence the trait



Qualitative traits

- Mendel studied major gene differences in the garden pea: colour of the seed (one locus with alleles Y and g, Y>g)
- But: some factors can affect the expression of Mendel's laws
 - Environmental influences can complicate the interpretation of the results
 - Presence of lethal alleles (homozygotes die at early developmental stage)
 - Presence of multiple alleles (more than two alleles at a single locus in the population)
 - Incomplete dominance (heterozygous individuals show an <u>intermediate</u> phenotype) and co-dominance (heterozygous individuals show a <u>distinct</u> phenotype)





Quantitative traits



Linkage mapping of QTLs for seed yield, yield components and developmental traits in pea (*Pisum sativum* L.) Timmerman-Vaughan et al. (2004) New directions for a diverse planet: Proceedings of the 4th International Crop Science Congress Brisbane, Australia, 26 Sep – 1 Oct 2004 | ISBN 1 920842 20 9 | WWW.Cropscience.org.au

If we analyze the diversity of seed size in a collection of pea plants, we are simultaneous looking at the genetic diversity contained in several genetic location

Quantitative traits

DIFFICULTIES

- Different genotypes can display the same phenotype
 - 3 genes each with small equal effect on the phenotype
 - each gene has two co-dominant alleles
 - 3³ = 27genotypes are possible, but we only distinguish 7 phenotypes
- <u>Dominance can vary from gene to gene</u>
 3 genes: 2 co-dominant, 1 dominant
- <u>Environmental variation:</u> the same genotype can display different phenotypes depending on the environment
 - Genotype and environment might interact obscuring genetic effects

genotype j	phenotype	
1. AABBCC	6	
2. AABBCc	5	
3. AABBcc	4	
4. AABbcc	3 🔨	Three different
5. AaBbCc	3 🛹	aenotypes aive the
6. AAbbCc	3 🚩	same phenotype
7. AAbbCC	4	
27 Aabbcc	0	

XE



Quantitative traits

DIFFICULTIES

• Epistasis

 Non-allelic interaction: the effect of the genotype at one locus depends on the genotype at another locus

No epistasis

				Genotypic difference	
	aa	Aa	AA	AA - aa	The change in phonetype when AA
bb	1	2	3	3 – 1 = 2	is replaced by as is the same
Bb	2	3	4	4 - 2 = 2	whatever the genotype at the B/b
BB	3	4	5	5 – 3 = 2	locus
Epis	stasis				
				Genotypic difference	The change in phenotype when AA
	aa	Aa	AA	AA - aa	is replaced by aa varies
bb	1	4	9	9 – 1 = 8	considerable depending on the B/b
Bb	4	9	16	16 – 4 = 12	genotype
BB	9	16	25	25 – 9 = 16	





6. Strategies to identify marker/trait associations

STRATEGIES

Qualitative traits

Linkage map construction + mapping of the trait

- The chromosomal location of a 'phenotype' or a 'mutation' is determined by identifying nearby genetic markers which are co-transmitted from parent to progeny with the phenotype
- Requires extensive phenotyping and genotyping of many plants of a mapping (=segregating) population
- Typical experiment: parents and >100 offspring plants + hundreds to thousands of molecular markers

Bulk Segregant Analysis (=BSA)

- Allows to reduce the amount of genotyping work
- Typical experiment: parents and 2 bulks + hundreds to thousands of molecular markers

Quantitative traits

Quantitative Trait Loci-mapping (=QTL-mapping)

- Requires a detailed linkage map (hundreds of molecular markers), using large progenies (typically 200 or more)
- Requires sophisticated statistical tools

Linkage Disequilibrium mapping (for quantitative and qualitative traits)

- Uses 'natural' or breeding populations to map traits by means of association analysis
- Mostly used in humans, recently extended to plants
- Requires sophisticated statistical tools; extremely difficult in highly heterozygous crops







Institute for Agricultural and Fisheries Research



Plant Sciences Unit www.ilvo.vlaanderen.be Agriculture and Fisheries Policy Area



EXAMPLE

Two phenotypic traits segregating in a progeny

In corn, colored aleurone (in the kernel) is due to one dominant allele R. The recessive allele (of the same locus) r, produces colorless aleurone. The plant color is controlled by another locus, with two alleles (Y and y). Y is dominant and results in green color, while y results in yellow color. In a cross between a plant with colored aleurone and green color, and a plant that is homozygous recessive for both traits, the following progeny was obtained:

Colored + green	88	88 R_Y_
Colored + yellow	12	12 R_yy
Colorless + green	8	8 rrY_
Colorless + yellow	92	92 rryy

Parent plant: <u>RrYy</u>	or	RryY	or	RRYY	
RY (parental)	ry = RrYy	88	(12-	+8)/(88+12+8+92) = 10% recomb.	
Ry (recombinant)	ry = Rryy	12		Π	
rY (recombinant)	ry = rrYy	8		₹,È	
ry (parental)	ry = rryy	92	Мар	o distance between loci R and Y =	=
•				10 cM	





EXAMPLE





Adapted from Paterson (1996). Explanations in next slide

ILVO

EXAMPLE

Backcross:

Hybrid	X	Recurrent parent
AcD		<u>aCd</u>
aCd		aCd

Compariso	on of A and C:		
Gametes		Progeny	Plant number
Hybrid	Rec. parent		
Ac (P)	aC	<u>Ac</u> aC	
aC (P)	aC	<u>aC</u> aC	
AC (R)	aC	<u>AC</u> aC	1, 13, 18
ac (R)	aC	<u>ac</u> aC	3, 10, 19
Comparise	on of A and D:		
Gametes		Progeny	Plant number
Hybrid	Rec. parent		
AD (P)	ad	<u>AD</u>	
		ad	
ad (P)	ad	<u>ad</u>	
		ad	
Ad (R)	ad	<u>Ad</u>	1, 13
		ad	
aD (R)	ad	<u>aD</u> ad	3, 10





DEFINITION

- Graphical representation of the genome of an organism
- A <u>linear map</u> of the relative positions of <u>genes</u>, <u>molecular</u> <u>markers and phenotypic markers</u> along a chromosome. Distances are established by <u>linkage analysis</u>, which determines the frequency at which two loci become separated during chromosomal recombination
- A genetic linkage map can be compared to <u>a road map</u>. Just as mile posts guide the motorist along a linear highway, molecular tools enable the geneticist to establish specific 'genetic markers' (DNA-markers) at defined places along each linear chromosome.





HISTORY

- Visual markers have been used in genetic studies since 1910
- Earliest linkage maps contained a few morphological markers mostly representing genes for which mutants were available
- At the DNA-marker level the amount of variation is so large that in a particular mapping population thousands of markers segregate
- Linkage maps were being used before it was known that DNA was the hereditary material, but <u>DNA-based technology</u> <u>revolutionized genetic mapping in plants</u>. This is because large numbers of genetic markers can be found by studying differences in the DNA molecule itself.





GENETIC BASIS

- "<u>Genetic linkage</u>", or co-transmission from parent to progeny of genetic markers (or genes or markers and genes) which are close together on the same chromosome, provides a means for determining the order of DNA markers along the chromosome
- By using in the analysis <u>hundreds to thousands of markers</u>, we can build up a 'linkage map' describing the relationships among the markers
- At the end we have a <u>schematic representation</u> of each chromosome, with the relative position of the DNA-markers
 - Only DNA-markers which are polymorphic in the studied population can be mapped => usually 'unrelated' genotypes are used as parents to construct the segregating population, to increase the probability that they will carry different alleles at many marker loci





TOOLS REQUIRED

- <u>A mapping population</u>. Different kinds of segregating populations can be used to construct a linkage map
- Hundreds to thousands of <u>DNA-markers</u>
- Appropriate statistical tools and software





STEPS

- 1. Creation of mapping population
- 2. Detection of polymorphic loci
- 3. Segregation analysis of single markers
- 4. Estimation of recombination frequencies
- 5. Identification of linkage groups
- 6. Calculation of map distances





1. Creation of mapping population

- 2. Detection of polymorphic loci
- 3. Segregation analysis of single markers
- 4. Estimation of recombination frequencies
- 5. Identification of linkage groups
- 6. Calculation of map distances





FEATURES OF DIFFERENT MAPPING POPULATIONS







MAPPING POPULATIONS

- Backcross population: constructed by crossing the F1 hybrid to one of the parents (the "recurrent" parent). Only alleles derived from the "donor" (non-recurrent parent) segregate
- F2 population: can be constructed by selfing the F1 hybrid
- Double haploid (DH) population: made by regenerating plants from single (haploid) pollen grains produced by the F1, and inducing chromosome doubling
- Recombinant Inbred (RI) population: the progeny of an F2 cross are self-pollinated during several generations, by applying 'single-seeddescent'
- Near-Isogenic (NI) population: the F1 hybrid is backcrossed to one of the parents (the "recurrent" parent) for various generations





FEATURES OF DIFFERENT MAPPING POPULATIONS

	# possible genotypes per locus (diploids)	# generations to make	Replication?
DH	2	2	yes
Backcross	2	2	clonal only
F2	3	2	clonal only
RI	2	6 to 8	yes
Cross-pollinator	4	1	clonal only





- 1. Creation of mapping population
- 2. Detection of polymorphic loci
- 3. Segregation analysis of single markers
- 4. Estimation of recombination frequencies
- 5. Identification of linkage groups
- 6. Calculation of map distances





DETECTION OF POLYMORPHIC LOCI







DETECTION OF POLYMORPHIC LOCI







DETECTION OF POLYMORPHIC LOCI







DETECTION OF POLYMORPHIC LOCI

- Dominant versus co-dominant markers
- Null-alleles can cause dominance



e.g. B allele is too long to amplify

- \Rightarrow No visible amplification product
- \Rightarrow No difference between AB and AA => A-





- 1. Creation of mapping population
- 2. Detection of polymorphic loci
- 3. Segregation analysis of single markers
- 4. Estimation of recombination frequencies
- 5. Identification of linkage groups
- 6. Calculation of map distances





SEGREGATION ANALYSIS OF SINGLE MARKERS

Expected genotype frequencies are defined per locus, according to the population structure and kind of marker (dominant/co-dominant)

Locus is AB in F1 => Expected segregation in F2 is 1:2:1 Locus is A- in F1 => Expected segregation in F2 is 3:1

Locus	Parent 1	Parent 2	F1		F2 individuals							
				1	2	3	4	5	6	7	8	9
1	AA	BB	AB	AA	AB	AB	AB	BB	BB	AA	AB	AB
2	AA		A-		A-	A-			A-	A-	A-	A-





SEGREGATION ANALYSIS OF SINGLE MARKERS

 Check for deviation from expected segregation => X² test

$$\chi^2 = \Sigma(O-E)^2$$
E

O = observed frequency

E = Expected frequency





SEGREGATION ANALYSIS OF SINGLE MARKERS



Genotypes	Expected	Observed	(O-E) ² /E
АА	25%	20%	1.00
AB	50%	40%	2.00
BB	25%	40%	9.00
		$\Sigma =$	12.00

 $\chi^2 = 12.00$ with 2 df

⇒p<0.005

⇒Hypothesis of 1:2:1 segregation is rejected






• e.g. Phaseolus vulgaris





Example

RIL

/// JM20D	emo.loc - Kladblol	¢											
Bestand	Bewerken Opm	aak Beeld ⊢	lelp										
; arabi	dopsis recom	binant inb	red lines	popula	ation							<u>_</u>	
name = popt = nloc = nind =	JM20Demo RI8 178 101											E	
CHS aabb- bbbbb b er	b-bba baabb abbab bbbbb	aa-ba aba bba bbb	b- bb-bb bbbba	bbbbb b bbbaa b	obbb- bbbbb oaa-a bbbba	bbbaa bab-a							
baaaa abbba b	ıbba babba ı aabbb bbb	bab-b -ab bbbbb -bb	-b aaabb a ab abbbb a	abbab b aab-b b	obbbb bbb-b obb-b babba	aabba abbba							
g10086 aabaa aaaba b	-aaaa bbbba abba- babb-	bbaba aaa bbbba baa	bb bbaaa ab abbbb a	aabab a aabba b	aaaab aabaa bbaaa b-baa	aaaba baaaa							
g17288 baaaa aabbb a	babba babaa aabbb abbaa	aabbb aab bbbbb abb	ab aabba ab abbbb I	abbab b babab a	obbbb bbbbb ababb abbbb	abbbb abbba							
g17311 baaaa -abaa b	u b-bab abaab u -aaba aabba	aaabb baa babaa aab	aa aabaa bb bbbba l	aabba a baaaa b	aaab- aabba babbb bbbba	aaaab aaaab							
g2368 bbbba -abaa	u b-aba bbbbb u aaabb abbab	baabb aba aba-a baa	bb abbba ab baaab	abbaa a abb-a a	aabba ababa aaaaa aaaab	babab aaaaa							
g2616 abbaa -aabb	b-aba bbbaa bbab aabbb	abbba aaa abbbb aab	bb bbaba l ba abbbb a	bbbab k aab-b k	oaba- aabab obabb bbaab	babbb bbaab							
g2778 abbba -aaab	b-aba bbbaa -aaab baaba	aaabb aaa abb-b aab	b- ab-ab bb aaaaa l	abbab a baa-a b	abab- baaaa D-aba babab	abbba abbba							
g3088 aaaaa -bbba	b-aab bbaaa bb-bb baabb	bbb-a -ba bbaaa baa	b- abaaa ab aaabb a	abbab a abb-a b	aabbb ababa D-aaa babaa	baaaa abaaa							
g3713 babaa -bbaa b	a-aab bbaab bbabb bbabb	bbbbb aba bbaab baa	a- bbaaa ab bbabb l	ababb a bab-a b	ababb ababa baaaa babaa	abaaa abaab							
g3715 aabab -bbab b	b-bba baabb abaab bbbbb	abbba bba abbaa abb	ba bbaba l ab babba i	bbbbb b abb-a b	obbba abbbb baaaa bbaba	bbbab babab							
g3786 bbbbb -baba	a-bab ababa -aabb bbabb	abbba aab bbaba baa	bb aaabb l aa babbb i	baabb a aab-a b	aabb- bbaba o-aab bb-aa	abbaa baaaa				_		-	
*									JM20Demo.loc - Kladblok	1		Þ	
	🧱 📰 🛸	🧾 👩 Micro	soft PowerPoi		JoinMap 4		🔀 Microsoft Excel - De	//// JM20Demo.loc - Kla			NL < 🎗 💽 🔱	间 🛃 🍖 21:57	





Example

RIL

	Info	Data Loci	Indiv	dual	s L	ocus	Genot	t. Fre <u>q</u> .	Individual G	not. Freq. Similarity of Loci Similarity of Individuals Groupings (text) Groupings (tree)	
no	S/n /	Nr Locus	а	h	0 1	d	-	X2	Df Signif.	Classification	
oing i	1	1 CHS	24	0	65	0 (12	18.89	1 ******	[a:b]	
roup 1	2	2 er	29	0	61	0 0	11	11.38	1 *****	[a:b]	
oup 2	3	3 g10086	56	0	41	0 (4	2.32	! 1 -	[aːb]	
up 3	4	4 g17288	38	0	63	0 0	0	6.19	1 **	[a:b]	
up 5	5	5 g17311	57	0	40	0 (4	2.98	1 *	[aːb]	
12	6	6 g2368	54	0	43	0 0	4	1.25	1 -	[a:b]	
n 1	7	7 g2616	38	0	58	0 0	5	4.17	1 **	[a:b]	
in 2	8	8 g2778	51	0	41	0 0	9	1.09	1 -	[a:b]	
ID 3	9	9 g3088	51	0	42	0 (8	0.87	1 -	[a:b]	
p 3	10	10 g3713	47	0	50	0 0	4	0.09	1 -	[a:b]	
n 5	11	11 g3715	36	0	62	0 0	3	6.90	1 ***	[a:b]	
	12	12 g3786	45	0	49	0 0	7	0.17	1 -	[a:b]	
	13	13 g3829	43	0	45	0 (13	0.05	1 -	[a:b]	
	14	14 g3837	38	0	60	0 0	3	4.94	1 **	[a:b]	
	15	15 g3843	40	0	57	0 0	4	2.98	1 *	[a:b]	
	16	16 g3845	55	0	41	0 0	5	2.04	1 -	[a:b]	
	17	17 g4014	50	0	39	0 0	12	1.36	1 -	[a:b]	
	18	18 g4026	65	0	33	0 0	3	10.45	1 ****	[a:b]	
	19	19 g4028	47	0	51	0 0	3	0.16	1 -	[a:b]	
	20	20 g4117	52	0	47	0 0	2	0.25	1 -	[a:b]	
	21	21 g4133	35	0	61	0 (5	7.04	1 ***	[a:b]	
	22	22 g4514	35	0	61	0 0	5	7.04	1 ***	[a:b]	
	23	23 g4523	52	0	43	0 0	6	0.85	1 -	[a:b]	
	24	24 g4532	33	0	60	0 (8	7.84	1 ***	[a:b]	
	25	25 g4552	69	0	29	0 (3	16.33	1 *******	[a:b]	
	26	26 g4553	33	0	62	0 0	6	8.85	1 ****	[8:0]	
	27	27 g4560	34	0	60	0 0	7	7.19	1 ***	[a:D]	
	28	28 g4564-a	56	0	43	0 0	2	1.71	1 -	[8:0]	
	29	29 g4564-b	50	0	41	0 0	10	0.89	1-		
	30	30 g4/08	51	0	45	0 0	5	0.38	1-		
	31	31 94/11	4/	0	48	0 0	6	0.01	1 -	[a.b]	
	32	32 g4/15-b	38	0	55	0 0	8	3.11	1 ^	[8:0]	
	33	33 g4/15a	48	0	48	0 0	5	0.00	1-	[8:0]	
	34	34 g6837	55	0	44	0 0	2	1.22	1	[a:0]	





SEGREGATION ANALYSIS OF SINGLE MARKERS

- Deviation from Mendelian segregation is called segregation distortion
- Causes:
 - selection
 - competition of gametes
 - non-random mating
 -
- Consequences:
 - innacurate estimation of recombination frequencies between markers
 - Innacurate estimation of QTL positions and effects
 - ...

Usually only markers that fit expected genotype proportions are kept for further steps





- 1. Creation of mapping population
- 2. Detection of polymorphic loci
- 3. Segregation analysis of single markers
- 4. Estimation of recombination frequencies
- 5. Identification of linkage groups
- 6. Calculation of map distances





ESTIMATION OF RECOMBINATION FREQUENCIES

A diploid plant with two chromosomes (2n=2x=4)







ESTIMATION OF RECOMBINATION FREQUENCIES

Parents: AB/AB (donor) x ab/ab (recurrent)

F1:

AB/ab

gametes	AB	Ab	aB	ab
ab	ab/AB	ab/Ab	ab/aB	ab/ab

BC1 genotypes	ab/AB	ab/Ab	ab/aB	ab/ab
frequencies	(1-r)/2	r/2	r/2	(1-r)/2





ESTIMATION OF RECOMBINATION FREQUENCIES

Estimation in a BC population:

gametes	AB	Ab	aB	ab
ab	ab/AB	ab/Ab	ab/aB	ab/ab

Genotype	Observed numbers	Expected freq
ab/AB	N1	(1-r)/2
ab/Ab	N2	(r/2)
ab/aB	N3	(r/2)
ab/ab	N4	(1-r)/2

r = (N2+N3)/(N1+N2+N3+N4)





ESTIMATION OF RECOMBINATION FREQUENCIES

- Loci located in close (physical) proximity => small chance of recombination
- The number of crossovers between two loci

=> estimate the map distance (centiMorgans)

 Recombination events can be recognized only the basis of haplotypes

=> determine haplotypes in mapping populations





ESTIMATION OF RECOMBINATION FREQUENCIES

LOD-Score

- Maximum likelihood (ML)
 - Generic name of statistical approaches in which one aims to find the parameter value (the value of <u>r</u> in our case) that maximizes the likelihood of the data. The question is: given the observed genotype values, what can we say about r?
 - L(r | x) = likelihood of r given the observed vector of data x
 - The MLE (maximum likelihood estimate) is the value of r that is most likely to have produced the data observed
- Likelihood Ratio test (LRT)
 - The ratio between the likelihood of r taking value r₁ (MLE), and the likelihood of r under the null hypothesis r₀ (r₀ =0.5; no linkage)
 - $L1 / L0 = L(r_1 | x) / L(r_0 | x)$
- LOD score (Z)
 - $Z = \log_{10} (L(r_1 | x) / L(r_0 | x))$
 - The alternative hypothesis is 10^z times more likely than the null hypothesis (Z = 3 means that the alternative hypothesis is 10³ times more likely than the null hypothesis)
 - In general, when Z = 3 or higher linkage is considered significant





Example

RIL

Info Data	a Loci	Individuals Loc	us Genot. Freg. I	lividual Genot. Fre	q. Similarity of Loc	i Similarity of Individua	Is Groupings (text) Grouping	ngs (tree)
S/n Nr1	Locus1	Nr2 Locus2	Similarity	-		1	1	- · · · ·
1 3	a10086	28 g4564-a	0.960					
2 3	g10086	54 m326	0.950					
3 3	q10086	119 w279	0.950					
4 4	q17288	118 w277	0.960					
5 7	g2616	60 m506	0.990					
6 14	g3837	41 m217	0.970					
7 16	q3845	72 w112	0.950					
8 16	q3845	135 w330	0.950					
9 18	g4026	109 w230	0.950					
10 18	q4026	127 w304	0.950					
11 18	q4026	150 w408	0.950					
12 20	q4117	48 m249	0.990					
13 23	q4523	64 m583	0.950					
14 28	q4564-a	34 g6837	0.960					
15 28	q4564-a	54 m326	0.970					
16 28	q4564-a	119 w279	0.990					
17 28	q4564-a	124 w30	0.970					
18 33	q4715a	59 m488	0.980					
19 33	q4715a	145 w372	0.960					
20 34	q6837	54 m326	0.950					
21 34	q6837	79 w128	0.960					
22 34	q6837	119 w279	0.970					
23 34	g6837	124 w30	0.970					
24 38	m105	84 w142	0.950					
25 47	m247	133 w323	0.950					
26 54	m326	119 w279	0.960					
27 58	m457	166 w456	0.950					
28 59	m488	145 w372	0.980					
29 69	w103	167 w462	0.980					
30 70	w106	105 w208	0.960					
31 70	w106	162 w441	0.950					
32 71	w111	74 w116	0.960					
33 71	w111	86 w15	0.980					
34 71	w111	96 w192	0.980					





- 1. Creation of mapping population
- 2. Detection of polymorphic loci
- 3. Segregation analysis of single markers
- 4. Estimation of recombination frequencies
- 5. Identification of linkage groups
- 6. Calculation of map distances





IDENTIFICATION OF LINKAGE GROUPS

- Use pairwise recombination frequency of all marker pairs
- Define groups of markers with a high likelihood to segregate together

=> at increasing stringency levels of a test for linkage=> using a LOD score as threshold





IDENTIFICATION OF LINKAGE GROUPS

- Different LOD thresholds => different numbers of linkage groups
 => LOD score low => low number of linkage groups
 => LOD score high => high number of linkage groups
- LOD values between 4 and 7 usually the 'best' grouping
- Ideally the number of linkage groups equals the haploid chromosome number over a wide range of LOD values

=> BUT if too few markers

=> more linkage groups than haploid chromosome number





Example

RIL

IM20Demo	Data Loci wea	K Linkages Strong	g Linkages	Maximi	linkages Suspect Linkages Sta <u>r</u> t Order Fi <u>x</u> ed Orders
Grouping 1	S/n Nr1 Locus1	Nr2 Locus2	Rec. Freq.	LOD	
Group 1	1 5 g17311	12 g3786	0.4990	0.00	
Group 2	2 5 g17311	13 g3829	0.4990	0.00	-
Group 3	4 5 g17311	33 g4715a	0.4990	0.00	
Group 4	6 5 g17311	45 m235	0.4990	0.00	
Group 5	7 5 g17311	50 m253	0.4990	0.00	
Grouping 2	8 5 g17311	59 m488	0.4990	0.00	
Group 1	9 5 g17311	68 w100	0.4990	0.00	
Group 2	10 5 g17311	71 w111	0.4990	0.00	
Group 3	11 5 g17311	73 w113	0.4990	0.00	
Group 4	12 5 g17311	74 w116	0.4990	0.00	
Group 5	13 5 g17311	86 w15	0.4990	0.00	
IM20Demo	14 5 g17311	88 w163	0.4990	0.00	
	16 5 g1/311	96 w192	0.4990	0.00	
	17 5 g17311	102 w203	0.4990	0.00	
	21 5 g17311	115 w265	0.4990	0.00	
	25 5 g17311	139 w348	0.4990	0.00	
	26 5 g17311	145 w372	0.4990	0.00	
	28 5 g17311	155 w423a	0.4990	0.00	
	29 5 g17311	163 w443	0.4990	0.00	
	30 5 g17311	175 w62	0.4990	0.00	
	31 5 g17311	176 w63	0.4990	0.00	
	32 12 g3786	25 g4552	0.4990	0.00	
	33 12 g3786	52 m315	0.4990	0.00	
	34 12 g3786	62 m532	0.4990	0.00	
	35 12 g3786	69 w103	0.4990	0.00	
	37 12 g3786	87 w157	0.4990	0.00	
	38 12 g3786	94 w185	0.4990	0.00	
	39 12 g3786	97 w193	0.4990	0.00	
	40 12 g3786	105 w208	0.4990	0.00	
	45 12 g3786	158 w425	0.4990	0.00	
	46 12 g3786	162 w441	0.4990	0.00	
	47 12 g3786	167 w462	0.4990	0.00	
	48 13 g3829	18 g4026	0.4990	0.00	
	49 13 g3829	25 g4552	0.4990	0.00	· · · · · · · · · · · · · · · · · · ·
	10.				
	# (24			-	
	Microsoft Pow	/erPoi 😽 🎎 Joir	nMap 4		🛛 Microsoft Excel - De 🔰 JM20Demo.loc - Kla NL 🔍 😼 🍓 22:2

- 1. Creation of mapping population
- 2. Detection of polymorphic loci
- 3. Segregation analysis of single markers
- 4. Estimation of recombination frequencies
- 5. Identification of linkage groups
- 6. Calculation of map distances





CALCULATE MAP DISTANCES

- Once we know which markers belong to the same linkage group and which is the pair-wise recombination between all of them, we can order the markers along the linkage group
- To do this, we try to find the marker order that minimizes the number of crossovers
- Several approaches available in mapping programs

$$A = B = C = 0.10$$

$$FBC = 0.07$$

$$FAC = 0.15$$





$C\mathsf{ALCULATE} \mathsf{MAP} \mathsf{DISTANCES}$

- The <u>distance along a genetic map is derived from the frequency of</u> <u>recombination</u> between genetic markers. 1 cM (*centiMorgan*) corresponds to an average of 1 recombinant in 100 gametes (1% recombination)
 - Usually, the distance is adjusted for the possibility of '<u>double recombination</u>', which makes individuals to appear to represent a parental type, but in fact containing two recombination events
 - The probability of double recombination is proportional to the square of the recombination distance between two loci.
- <u>The precision</u> with which genetic distance is measured, is directly related to the <u>number of individuals</u> which is studied (if no recombinants found a sample of 20 progeny plants => recombination fraction = 0; but analyzing 80 additional individuals 1 recombinant can appear => recombination fraction = 1).
 - Typically a primary genetic map is constructed based on 50-100 individuals, permitting to detect recombination between markers 1-3 cM appart.
 - If higher precision is required, more individuals should be analyzed





CALCULATE MAP DISTANCES

- Map function translates recombination frequencies into map distances
 - Haldane function

$X = \frac{1}{2} \ln (1-2r)$

- Takes into account double strand , three and four strand crossovers
- BUT assumes that 2-3-4 strand crossovers does not interfere with single crossovers





CALCULATE MAP DISTANCES

- Map function translates recombination frequencies into map distances
 - Kosambi function
 - Takes into account interference

$X = \frac{1}{4} \ln \left[\frac{1+2r}{1-2r} \right]$





C ALCULATE MAP DISTANCES

The genetic distance is only loosely related to the physical distance

=> influenced by - genetic

- epigenetic
- environmental factors
- For example, repetitive DNA elements are relatively inert in recombination, recombination is also reduced around the centromers





C ALCULATE MARKER ORDER

- There are many possibilities, but one of the computationally 'cheapest' algorithms is the greedy algorithm (stepwise build-up of the map by adding one marker at a time):
 - Start with two markers A and B
 - Add third marker C at the three different positions which are possible: CAB, ACB, ABC
 - The three orders are compared and the best fitting one is chosen (ACB)
 - Then add a fourth marker D at the four different positions which are possible: DACB, ADCB, ACDB, ACBD
 - The four orders are compared and the best fitting one is chosen
- JoinMap uses this algorithm, but with several refinements:
 - The order in which the markers are added is not random
 - After a marker has been added a 'local reshuffling' is applied to prevent that the previous sequence will not be changed any more
- There are many possibilities to use ad 'fitting criterion' (minimization of the sum of adjacent distances, maximization of the sum of adjacent LOD scores, weighted least squares)





Example

RIL









MAPPING POPULATIONS IN POLYPLOIDS

2	2	yes
2	2	clonal only
3	2	clonal only
2	6 to 8	yes
4	1	clonal only
400!!!!!	1	clonal only
-	2 3 2 4 400!!!!!	2 2 3 2 2 6 to 8 4 1 400!!!!! 1



MAPPING POPULATIONS IN POLYPLOIDS

Why is mapping in polyploids more complicated than in diploids?

- 1. Lage number of genotypes expected (up to 400 if 12 alleles segregating)
- 2. Genotype of an individual is not always readily inferred (even using codominant markers)
 - Result of genotyping at one locus = ABC
 - AAAABC, ABBBCC, AABBCC,?
- 3. The type of ploidy (allopolyploidy or autopolyploidy) of many crops is unclear
 - Wheat is an allopolyploid (mapping is as in a diploid with many chromosomes)
 - Sweetpotato is an autopolyploid with hexasomic inheritance (Cervantes-Flores et al 2008)







MAPPING POPULATIONS IN POLYPLOIDS

Chromosome segregation in polyploids:



Figure: De Baker (2012)







SEGREGATION ANALYSIS OF SINGLE MARKERS

Kreigner et al (2003) - hexasomic inheritance in sweetpotato

Table 1. Expected segregation ratios (presence:absence) for the inheritance of a dominant marker in hexaploid sweetpotato, according to four cytological hypotheses (Jones 1967)

Marker dose	Hypothesis I Autohexaploid	(hexasomic)	Hypothesis II Tetradiploid (te somic, disomic	and III etra-disomic, tetra- :)	Hypothesis IV Allohexaploid (disontic)	
Simplex	Aaaaaa	1:1	Aaaa aa	1:1	Aa aa xa 1:1	
			aaaa Aa	1:1		
Duplex	AAaaaa	4:1	AAaa aa	5:1 ²⁾	Aa Aa aa 🗙 3:1	
			Aaaa Aa	3:1 3)	AA aa aa 1:0	
			aaaa AA	1:0 1)		
Triplex	AAAaaa	19:1	AAAa aa	_	Aa Aa 7:1	
			AAaa Aa	11:1	AA Aa aa 1:0	
			Aaaa AA	1:0		
Quadruplex	AAAAaa	1:0	AAAA aa	1:0	AA Aa Aa 1:0	

1) disomic inheritance

²⁾ tetrasomic inheritance

3) tetra-disomic inheritance

Using this table, it is possible to check the segregation of individual dominant markers in sweetpotato







$\label{eq:segregation} Segregation \ \text{analysis} \ \text{of single markers}$

- Deviation from Mendelian segregation is called <u>segregation</u> <u>distortion</u>
- Causes:
 - selection
 - competition of gametes
 - non-random mating
 - non-random pairing of chromosomes (in autopolyploids)
 -
- Consequences:
 - innacurate estimation of recombination frequencies between markers
 - Innacurate estimation of QTL positions and effects







ESTIMATION OF RECOMBINATION FREQUENCIES

For tetraploids, suited linkage mapping program available: Tetraploidmap (<u>http://www.bioss.ac.uk/knowledge/tetraploidmap/</u>)

- Handles dominant and co-dominant markers in all possible configurations, taking account of null alleles
- Free to download
- Approach in hexaploids (Cervantes-Flores et al 2008):
- Double pseudotestcross strategy
 - Reconstruct individual chrosomosomes in each parent independently
- Simplex AFLP markers (=SDF) are present in one parent in a single copy => segregate 50%/50% (present:absent) in progeny
 - A00000 x 000000 => A000000 / 0000000
- Mapping in this case is similar to diploid situation and can be carried out using standard mapping programs such as JoinMap or MapMaker
- Limitation: using simplex markers, impossible to identify homologous chromosomes







ESTIMATION OF RECOMBINATION FREQUENCIES

Simplex markers in coupling allow to reconstruct individual chromosomes in each parent separately

Table 2 Marker pair configurations and expected phenotypic frequencies used in this study

Hexasomic inheritance marker-pair configuration ^a	Phe	notype probabilities
Simplex/simplex coupling	AB	1/2(1 - r)
AB/00/00/00/00/00 \times 00/00/00/00/00/00/00	Α	1/2r
	В	1/2r
	0	1/2(1 - r)
Simplex/duplex coupling	AB	1/2 - 1/5r
AB/0B/00/00/00/00 \times 00/00/00/00/00/00	Α	1/5r
	в	3/10 + 1/5r
	0	1/5 – 1/5r
Duplex/duplex coupling	AB	$4/5 - 2/5r + 1/5r^2$
AB/AB/00/00/00/00 × 00/00/00/00/00/00	Α	$2/5r - 1/5r^2$
	в	$2/5r - 1/5r^2$
	0	$1/5 - 2/5r + 1/5r^2$
Simplex/triplex coupling	AB	1/2 - 1/20r
AB/0B/0B/00/00/00 × 00/00/00/00/00/00	Α	1/20r
	В	9/20 + 1/20r
	0	1/20 - 1/20r
Duplex/triplex coupling	AB	$4/5 - 1/4r + 1/20r^2$
AB/AB/0B/00/00/00 × 00/00/00/00/00/00	A	$1/10r - 1/20r^2$
	В	$3/20 - 1/10r - 1/20r^2$
	0	$1/20 - 1/10r + 1/20r^2$



a Distribution of alleles of two loci in a base chromosome group (chromosomes are separated by "/"). "A": presence of band at locus

A, "B": presence of band at locus B, "0": absence of band, (Source: Kriegner et al. 2003; Ripol et al. 1999)







One linkage group per chromosome (90) of each parent => in total ~180 linkage groups

How to know which linkage groups are homologous?







ESTIMATION OF RECOMBINATION FREQUENCIES

Multiplex (duplex and triplex) markers allow to establish homology relationships among the 6 parental chromosomes

Table 2 Marker pair configurations and expected phenotypic frequencies used in this study

Hexasomic inheritance marker-pair configuration ^a	Phe	notype probabilities
Simplex/simplex coupling	AB	1/2(1 - r)
AB/00/00/00/00/00 × 00/00/00/00/00/00	Α	1/2r
	В	1/2r
	0	1/2(1 - r)
Simplex/duplex coupling	AB	1/2 - 1/5r
AB/0B/00/00/00/00 × 00/00/00/00/00/00	А	1/5r
	в	3/10 + 1/5r
	0	1/5 – 1/5r
Duplex/duplex coupling	AB	$4/5 - 2/5r + 1/5r^2$
AB/AB/00/00/00/00 \times 00/00/00/00/00/00/00	Α	$2/5r - 1/5r^2$
	в	$2/5r - 1/5r^2$
	0	$1/5 - 2/5r + 1/5r^2$
Simplex/triplex coupling	AB	1/2 - 1/20r
AB/0B/0B/00/00/00 \times 00/00/00/00/00/00	Α	1/20r
	в	9/20 + 1/20r
	0	1/20 - 1/20r
Duplex/triplex coupling	AB	$4/5 - 1/4r + 1/20r^2$
AB/AB/0B/00/00/00 \times 00/00/00/00/00/00/00	Α	$1/10r - 1/20r^2$
	В	$3/20 - 1/10r - 1/20r^2$
	0	$1/20 - 1/10r + 1/20r^2$



^a Distribution of alleles of two loci in a base chromosome group (chromosomes are separated by "/"). "A": presence of band at locus

A, "B": presence of band at locus B, "0": absence of band. (Source: Kriegner et al. 2003; Ripol et al. 1999)









*

Software available

Diploids and allopolyploids:

- JoinMap (http://www.kyazma.nl)
- MapMaker
- Carte Blanche (for extremely large datasets)....

Autotetraploids:

• Tetraploidmap

Autohexaploids:

• Software for diploids, but using strategy described above







Use of linkage maps

- Locate genes affecting trait of interest (statistical associations between DNA-marker variants and trait)
 - Genetically simple traits (monogenic); disease resistances in plants are frequently controlled by one or a few genes
 - Genetically complex traits, involving many genes (quantitative trait loci, QTL), and affected by the environment; economically important traits such as yield or stress resistances
- Fine-mapping of genes
- Compare the genomes of different species
- Select DNA-markers evenly spread over the genome to perform a diversity study







BSA-analysis

Institute for Agricultural and Fisheries Research



Plant Sciences Unit www.ilvo.vlaanderen.be Agriculture and Fisheries Policy Area


BSA

DEFINITION

<u>B</u>ulked: it makes use of bulked samples => saves time and money <u>S</u>egregant: it makes use of a segregating population

=> but it does not require map information!!! <u>Analysis: it screens the whole genome</u>







DEFINITION

- Involves comparing two pooled DNA samples of individuals from a segregating population originating from a single cross
- Within each bulk the individuals are identical for the loci of interest but are arbitrary for all other loci
- Markers that are polymorphic between the pools are markers putatively linked to the loci involved in the trait of interest









DEFINITION

- BSA assumes that markers adjacent to (=linked to) the target locus will be in linkage disequilibrium (LD*) among themselves and with respect to the trait
 - Recombination will not have randomised these markers with respect to the targeted locus
 - As genetic distance increases, more recombinants will be present in each bulk, culminating in 50% recombinants (= no linkage disequilibrium) and therefore resulting in no differences between the bulks

- * non-random association of alleles at different genomic loci
- * non-random association of a marker-allele at a given locus with a phenotype
- * related to Mendel's second Law







DEFINITION

At a locus not linked to the trait, using dominant markers







BSA

DEFINITION







STEPS

- 1) Create a segregating population from a single cross (for example, F2 or BC progeny)
- 2) Phenotype the progeny and identify individuals with extreme trait-phenotypes
- 3) Construct DNA bulks of the individuals displaying the most extreme trait-phenotypes
- 4) Genotype the parents and the bulks using hundreds to thousands of DNA-markers
- 5) Identify those markers which distinguish the bulks and the parents







BSA

EXAMPLE SWEETPOTATO

Ukoskit et al. (1997) Identifying a randomly amplified polymorphic DNA (RAPD) marker linked to a gene for root-know nematode resistance in sweetpotato

Fenotype screening

• F1 progeny of a single cross => log total nematode number => qualitative trait



Fig 1. Frequency distribution of F_1 single-cross progenies for log total nematode number.

• Segregation in progeny fits 4:1 ratio => RRrrrr x rrrrrr







BSA

EXAMPLE SWEETPOTATO

Ukoskit et al. (1997) Identifying a randomly amplified polymorphic DNA (RAPD) marker linked to a gene for root-know nematode resistance in sweetpotato

Genotype screening

• R bulk, S bulk, R parent, S parent



Fig 2. Bulk segregant analysis of $OPI5_{1500}$. (a) Lane 2 = susceptible parent, lane 3 = resistant parent, lane 4 = resistant bulk, lane 5 = susceptible bulk. (b) All susceptible progenies (* = recombinant progenies, susceptible with marker). (c) Example of resistant progenies. Arrows indicate the polymorphic band at 1500 bp. The first lane of each picture shows 100-bp molecular mass marker.

• 728 polymorphic bands screened; one linked with resistance; linkage is not complete (recombinants present); recombination fraction = 0.2421







QTL-analysis

Institute for Agricultural and Fisheries Research



Plant Sciences Unit www.ilvo.vlaanderen.be Agriculture and Fisheries Policy Area



QUANTITATIVE TRAIT LOCUS (QTL) MAPPING

- A <u>QTL</u> is the location of a gene (or set of genes) that affects a trait that is measured on a <u>quantitative</u> scale. Examples of quantitative traits are plant height, fruit size, grain yield...
- These traits are typically affected by more than one gene, and also by the environment
- Mapping QTL is not as simple as mapping a single gene that affects a qualitative trait
- Tools required:
 - Many polymorphic molecular markers ordered in a linkage map
 - Variation for the trait in the segregating population
 - Detailed and accurate phenotypic data of the segregating population
 - Appropriate statistical tools





QUANTITATIVE TRAIT LOCUS (QTL) MAPPING

 Also based on the assumption that markers adjacent to (=linked to) the targeted loci will be in linkage disequilibrium with respect to the trait

but

- QTL-analysis tries to identify <u>simultaneously</u> the chromosomal location of all the genetic factors affecting the trait
- In a QTL analysis we infer the QTL genotypes in order to estimate the QTL effects and locations from associations with known markers
- A QTL is described by
 - 1. Its chromosomal location
 - 2. The magnitude of its phenotypic effect
 - 3. The effect of gene dosage at the locus
 - 4. Its interactions with other QTLs





QUANTITATIVE TRAIT LOCUS (QTL) MAPPING

Chromosomal location:

Derived, for example, from associations between marker genotypes and trait phenotypes

Phenotypic effects:

In most studies it is found that a few genes explain large proportions of the phenotypic variance, with increasing numbers of genes explaining progressively smaller fractions of phenotypic variance

Effects of gene dosage:

The "gene action" at QTLs is determined by the same principles (additivity, dominance, recessivity...) as are employed for monogenetic traits

Epistasis (= interactions between QTLs)

The collective activities of mapped QTLs explain only a portion of the phenotypic difference between parents. This is usually explained by the importance of epistasis, or nonlinear interactions between unlinked loci





QUANTITATIVE TRAIT LOCUS (QTL) MAPPING

- χ^2 test of association or simple linear regression (no linkage map required)
 - **Pro:**
 - robust to violations of normality in phenotypic data
 - Con:
 - Cannot extract all the information from the data
 - Large progenies required
 - Bad resolution of QTL position obtained
 - Multiple testing problem

Currently only used for initial data exploration

- Single marker analysis
 - Tests for differences in the means of the genetic marker classes
 - Rough estimation of the location of a QTL





QUANTITATIVE TRAIT LOCUS (QTL) MAPPING



Possible marker genotypes and their phenotypic means:

MQ/MQ and MQ/Mq: $\mu_{MM} = (1-r_{MQ}) \mu_1 + r_{MQ} \mu_{BC1}$ MQ/mQ and MQ/mq: $\mu_{Mm} = r_{MQ} \mu_1 + (1-r_{MQ}) \mu_{BC1}$

 r_{MQ} = recombination freq between marker and QTL

 μ_1 = average trait value of the plants with marker genotype MM

 μ_{BC1} = average trait value of the BC population





QUANTITATIVE TRAIT LOCUS (QTL) MAPPING

• Expected difference in average trait values between the two marker classes (MM and Mm) is:

```
\mu_{MM} - \mu_{Mm} = (1-2r_{MQ})\delta
\delta = difference in trait values between P1 and P2
```

- If the trait and the marker are unlinked, r_{MQ} = 0.5: μ_{MM} μ_{Mm} = 0
- What we have to do then is test the hypotheses:

```
H_0: \mu_{MM} - \mu_{Mm} = 0H_1: \mu_{MM} - \mu_{Mm} \neq 0
```

- Standard t-test to compare the means will determine whether the marker genotype has any effect on the trait.
- If we conclude that H₁ is true => the marker has an effect on the trait, and we can conclude that the marker *M* is located in the neigbourhood of QTL *Q*.





QUANTITATIVE TRAIT LOCUS (QTL) MAPPING

Single marker analysis

Pro:

• Robust to violations of normality in phenotypic data

Con:

- The order of the M and the QTL on the genetic map remains unknown (MQ or QM?)
- Low power if the markers are far apart
- Large progenies required
- It is not possible to distinguish between size of a QTL effect and its position (relative to the marker): a marker close to a QTL of small effect will give the same 'signal' as a marker some distance from a QTL of large effect
- t-testing is not appropriate due to deviations from normality. Non-parametric tests such as Kruskal-Wallis are more appropriated
- Multiple testing problem





QUANTITATIVE TRAIT LOCUS (QTL) MAPPING

- <u>Interval mapping</u> (requires linkage map)
 - Estimates the position of a QTL between two markers
 - Systematically searches the genome by calculating a test statistic at each position of the genome
 - Originally based on the maximum likelihood estimates
 - Intervals between adjacent markers along a chromosome are scanned and the LOD of there being one versus no QTL at a particular point is estimated
 - The higher the value of the LOD score, the more likely the data are if there was a QTL present compared to the situation when there is no QTL
 - A LOD-profile is constructed along the chromosome, and the maxima in this profile which exceed a specified significance level, indicate likely sites of a QTL
 - Pro:
 - By taking into account several markers simultaneously, it allows accurate estimation of QTLpositions
 - It is possible to separate the recombination and the size of the gene effect
 - Precision and power are increased by the use of extra information from a second marker
 - Con:
 - The effect of other QTLs present in the genome is neglected, and only the QTLs with the biggest phenotypic QTL on a chromosome, estimated positions and effect may be biased
 - Multiple testing problem





QUANTITATIVE TRAIT LOCUS (QTL) MAPPING



LOD = Log of the odds $Z = \log_{10} (L1/L0)$ Interpretation= the alternative hypothesis (L1) is 10^{Z} times more likely than the null hypothesis (L0)





QUANTITATIVE TRAIT LOCUS (QTL) MAPPING

- <u>Composite Interval mapping</u> (CIM) or Multiple QTL mapping (MQM)
 - Markers located nearby putative QTLs identified by e.g. IM, are used as cofactors in an approximate multiple-QTL model. At each testing point the effect of one or more co-factors is included.
 - By entering QTLs identified by IM (with the biggest effects) as co-factors, the effects of these QTLs is absorbed, increasing the power to identify additional QTLs.

The most used method at present







QTL MAPPING (CIM)

Fig. 5 Tanzania linkage groups showing significant QTL for β -carotene content, starch content, and drymatter content. Significance at 5, 1 and 0, 1% is indicated by *, ** and ***. respectively. The markers most significantly associated with the trait. according to ANOVA. analysis, are shown in bold. OTL are shown as vertical bars on the right side of the respective linkage group. Shaded boxes OTL are named and according to the trait they affect (dm = drymatter content; sta = starch content: and caro = carotene content). OTL with positive effect are indicated by normal text, and those with negative effect are indicated by underlined text

Cervantes-Flores et al (2011)





X

7. Developing markers for application

DNA-MARKER VALIDATION

- When we map one trait in one segregating population, we are analyzing the effects of the genetic factors carried by the parents of the mapping population, but we do not know anything about the situation in the 'global population' (outside the mapping population)
 - For example, in a mapping population we find that allele 1 at locus A (A1) is associated with disease resistance. We do not know yet whether in a different plant (not used for mapping) allele 1 at locus A will also be associated with disease resistance. In many cases it will be a different allele
 - A different genetic locus might be more relevant in other plants
- Test marker trait association in alternative populations => estimate reliability of marker in predicting phenotype
- Identify polymorphisms between lines used in breeding programs (perhaps different alleles present than in the mapping population)
- Develop a palette of suitable markers wih associated polymorphism data
 - E.g. 10 markers within 10 cM of trait for Marker Assisted Backcross selection (MAB)
 - Provide protocols and polymorphism data to breeder





7. Developing markers for application

DNA-MARKER VALIDATION



Association analysis (= linkage disequilibrium mapping) or association mapping, is a population-based survey used to identify trait-marker relationships based on linkage disequilibrium (Flint-Garcia et al. 2003. Annu. Rev. Plant Biol. 54, 357-374)



- F2 family and AM population: co-heritance (= co-segregation) of functional polymorphisms and neighboring DNA markers
- F2 family: only a few opportunities for recombination (only a few generations) => low mapping resolution
- AM population: historical recombination and natural genetic diversity exploited => high mapping resolution





ADVANTAGES OF AM

- Increased mapping resolution
- Reduced research time
- Greater allele number
- Results readily applicable in breeding germplasm







GENOME-WIDE / CANDIDATE-GENE



Genome-wide association mapping

It is a comprehensive approach to systematically search the genome for causal genetic variation. A large number of markers are tested for association with various complex traits, and prior information regarding candidate genes is not required. It works best for a research consortium with complementary expertise and adequate funding.

Candidate-gene association mapping

Candidate genes are selected based on prior knowledge from mutational analysis, biochemical pathway, or linkage analysis of the trait of interest. An independent set of random markers needs to be scored to infer genetic relationships. It is a low cost, hypothesis-driven, and trait-specific approach but will miss other unknown loci. Q, K or both can be included in the analysis, depending on the genetic relationship of the association mapping population and the divergence for the trait examined

- Q = population structure
- K = relative kinship
- E = residual variance





REQUIREMENTS

- Natural diversity
 - AM exploits natural diversity and its results are applicable to a wide germplasm base
 - As QTLs are mapped in collections of breeding lines, landraces, or samples from natural populations, AM offers great potential for future trait improvement
- Genomic technology
 - Technology at low cost/data point (SNPs) new sequencing technologies
 - Information on location and function of genes involved (for candidate-gene approach)
- Methodology for data analysis
 - To date, few efforts for AM tools specifically adapted to plants
 - Specificities of plants: geographical origins, local adaptation, and breeding history



