Molecular marker redundancy check and construction of a high density genetic map of tetraploid cotton

Jean-Marc Lacape

CAACGAAATACTCCAAAATGCAGTCAGCAGCAGCAACAGCAACAGCAACAGCAACTCTCAGCAT CAACGAAATACTCCAAAATGCAGTCAGCAGCAGCAACAGCAACAGCAACAGCAACTCTCAGCAT CAACAAAATACTCCAAAACGCAGTCAGCAGCAGCAGCAGCAGCAACAGCAAC----TCTCAGCAT







Raleigh ICGI, October 10, 2012



ICGI/2008 Anyang

Numerous marker projects (several 10s thousands) SSRs: BNL, CIR, etc.. SNPs, etc ...

Structural Genetics WG priorities

1- Marker nomenclature, synonymy and redundancy

Several genetic/QTL maps (several 10s) Interspecific (*Gh*x*Gb*) Intraspecific (*Gh*x*Gh*) Others ..

2- Map convergence, and building of a consensus map

Genome sequencing D-genome (A- genome) (AD-genome)



3- Convergence geneticphysical

Objectives

Numerous marker projects (several 10s thousands) SSRs: BNL, CIR, etc.. SNPs, etc ...

1- Curate redundancy among SSR markers

Several genetic/QTL maps (several 10s) Interspecific (*Gh*x*Gb*) Intraspecific (*Gh*x*Gh*) Others ..

2- Build a High Density Consensus (**HDC**) map from 6 component maps

Genome sequencing D-genome (A- genome) (AD-genome)



3- Conduct geneticphysical map alignments

Library	No seq. In library
BNL	379
CIR	392
СМ	53
DOW	52
DPL	200
Gh	700
HAU	3383
JESPR	309
MGHES	84
MON	2568
MUCS	617
MUSB	1316
MUSS	554
NAU	3249
NBRI	2233
PGML	308
STV	192
ТМВ	754
total	17343

17343 markers from CMD (18 different
libraries) as of April 2012
pair-wise sequence alignment (Smith-
Waterman algorithm)
similarity threshold 90%
ATACTCCAAAATGCAGTCAGCAGCAGCAACAGCAACAGCAACAGCAACTCTCAGCATTGCACCACCCC ATACTCCAAAATGCAGTCAGCAGCAGCAACAGCAACAGCAACAGCAACTCTCAGCATTGCACCACCCC ATACTCCAAAACGCAGTCAGCAGCAGCAGCAGCAGCAACTCTCAGCATTGCACCACCCC
ACACACACACACACA TACACAAA CCAAGCAA CCATGGATA CAAGAAGCAAAA CACCTAA TGAA CCATGGCTGCTT ACACA CACACACACACA TACACAAA CCAAGCAACCATGGATA CAAGAAGCAAAA CACCTAA TGAACCATGGCTGCTT ACACACACACACACACATA CACAAA CCAAGCAACCATGGATA CAAGAAGCAAAA CACCTAA TGAACCATGGCTGCTT ACACACACACACACACA CCAAGCAACCATGGATA CAAAAAGCAAAA CACCTAA TGAACCATGGCTGCTT ACACACACACACACACA CCAAGCAACCATGGATA CAAAAAGCAAAA CACCTAA TGAACCATGGCTGCTT ACACACACACACACACA CCAAGCAACCATGGATA CAAGAAGCAAAA CACCTAA TGAACCATGGCTGCTT ACACACACACACACACA CCAAGCAACCATGGATA CAAGAAGCAAAA CACCTAA TGAACCATGGCTGCTT ACACACACACACACACA
TTCATTTCATTCAAATCGAAGGCTCTCTCTCTCTCTCTCT
SCATGTAAATATATATATATCGTCTTTTTTAAAAGTTTTTCAAT – GGGTCAATTAAGTAGTAGTAGTAGTAGTAAGTAAGAT – GTTTAAA SCATGTAAATATATATATATCGTCCTTTTTTAAAAAGTTCTTCAAT – GGGGTCAATTAAGTAGTAGTAGTAGTAGTAAGAAGAT – GTTTAAA SCATGTAAATATATATATATCGTCCTTTTTTAAAAAGTTCTTCGAT – GGGTCAATTAAGTAGTAGTAGTAGTA – – – AGTAAGAT – GTTTAAA SCATGTAAATATATATATATCGTCCTTTTTTAAAAGTTCTTCAAT – GGGTCAATTAAGTAGTAGTAGTAGTA – – – AGTAAGATTGTTTAAA GCATGTAAATATATATATATCGTCCTTTTTTAAAAGTTCTTCAAT – GGGTCAATTAAGTAGTAGTAGTAGTA – – – AGTAAGATTGTTTAAA GCATGTAAATATATATATATCGTCCTTTTTTAAAAGTTCTTCAAT – GGGTCAATTAAGTAGTAGTAGTAGTA – – – AGTAAGAT – GTTTAAA GCATGTAAATATATATATATATCGTCCTTTTTTAAAAGTTCTTCAAT – GGGTCAATTAAGTAGTAGTAGTA – – – AGTAAGAT – GTTTAAA

GCATGTAAATATATATATCGTNCTTTTTAAAAGTTCTTCAATTGGGTCAATTAAGTAGTAGTAGTA---AGTAAGAT-GTTTAAA

► at 90% sequence similarity threshold, **5741 markers**,

Lib	BNL	CIR	СМ	DOW	DPL	Gh	HAU	JESPR N	MGHES	MON	MUCS	MUSB	MUSS	NAU	NBRI	PGML	STV	тмв
BNL	9	10	1		4	8	4	4		16	2			4	1			3
CIR		6				7	2	3		8				2	1	1		40
СМ			5		1	3		28		3								1
DOW				1	1	3		1										
DPL					1	2	5			2				5	3			3
Gh						103	5	29	1	27	1	1		4	10			1
HAU							646	3	49	244	145	1	206	872	136	55	46	3
JESPR								14	4	11				4	2			1
MGHES	5								3	40	6		10	52	3	2		
MON										183	14	2	22	182	56	6	8	6
MUCS											54		162	114	28	22		
MUSB												168		12	1			
MUSS													5	224	26	19		
NAU														593	123	50	115	
NBRI															146	5	4	1
PGML																16	2	3
STV																	12	
тмв																		35

Level of redundancy within libraries



Level of redundancy between libraries

Lib	BNL	CIR	СМ	DOW	DPL	Gh	HAU	JESPR	MGHES	MON	MUCS	MUSB	MUSS	NAU	NBRI	PGML	STV	тмв
BNL	9	10	1		4	8	4	4		16	2			4	1	<u> </u>		3
CIR		6				7	2	3		8				2	1	1		40
СМ			5		1	3		28		3								1
DOW				1	1	3		1										
DPL					1	2	5			2				5	3			3
Gh						103	5	29	1	27	1	1		4	10			1
HAU							646	3	49	244	145	1	<mark>206</mark>	872	<mark>136</mark>	55	46	3
JESPR								14	4	11			7	4	2			1
MGHES	5								3	40	6		10	52	3	2		
MON										183	14	2	22	182	56	6	8	6
MUCS											54		<u> 162</u>	114	28	22		
MUSB			- 1		. 111							168		12	1			
MUSS			etw	veer		orar	les	,					5	224	26	19		
NAU	r	nat	che	s ar	nor	าตร	t							593	123	50	115	
NBRI		$\mathbf{n} \mathbf{n}$	+ lik	vori		3-									146	5	4	1
PGML		nos	ot III.	лап	e 5											16	2	3
STV																	12	
ТМВ																		35

	No seq.	Total
Library	In	cross library
	library	(pair-wise)
BNL	379	57
CIR	392	74
CM	53	37
DOW	52	5
DPL	200	26
Gh	700	102
HAU	3383	1776
JESPR	309	90
MGHES	84	167
MON	2568	647
MUCS	617	494
MUSB	1316	17
MUSS	554	669
NAU	3249	1763
NBRI	2233	400
PGML	308	165
STV	192	175
TMB	754	62
total	17343	6726



- our sequence-based redundancy check aimed at recovering <u>identical-by-sequence</u> SSR markers (cross-PCR-amplification), thus recovering as « redundants »:
- homoeologs, paralogs and other copies
- sequences with mutiple SSRs
- alleles etc ...

It is <u>necessary</u> to account-for marker/locus redundancy for (1) across- map alignments and (2) consensus map construction

Six selected interspecific Gh x Gb genetic maps (only map positions, **no raw codings**)

Map code	Parents	Gene- ration	References	standardize marker names
				(BNL119=BNL0119)
				verify orientation
T3	TM-1 x 3-79	RIL	Yu et al (2012)	
GV	Guazuncho 2 x VH8	BC ₁ -RIL	Lacape et al (2009)	acronyms T3,
PK	Palmeri x K101	F_2	Rong et al (2004)	GV, PK, TH, E3 and
TH	TM-1 x Hai 7124	BC_1	Guo et al (2008)	СН
E3	Emian22 x 3-79	BC_1	Yu et al (2011)	
CH	CRI 36 x Hai 7124	F_2	Yu et al (2007)	

Six maps were saturated with combination of different markers types (SSRs as predominant)

Map code	Gene- ration	Pop size		сM	No	. of m	apped markers			
			Total		RFLP	AFLP	SSR	SNP	Others	
T3	RIL	186	2072	3380			1825	247		
GV	BC1-RIL	215	1745	3637	190	715	781		59	
PK	F2	57	2584	4448	2459		124		1	
TH	BC1	138	2247	3541		71	1865	10	301	
E3	BC1	141	2316	<u>4419</u>			2311		5	
CH	F2	186	1080	4418		93	690		297	
Total			12044		2649	879	<mark>7596</mark>	257	663	

Six maps with good level of connections: as many as 10086 bridges (same chromosome) altogether



Following redundancy check, groups of suspected redundant markers (2357 clusters with up to 14 SSR) were assigned a collective name 'CLU' (cluster)



Additional evidence of redundancy from map data

► CLU1057



Redundancy check (replacing marker names by a collective 'CLU' name) resulted in:

additionnal new correspondances: 2 redundant markers mapped on 2 different maps, including erroneous paralogs

numerous spurious correspondances generated by colocalized redundant markers on the same map

need for a curation of duplicates both internally and

between maps

<u>all inversions</u>
<u>between all</u>
<u>possible pairwise</u>
<u>comparisons</u>
<u>had to be solved</u>



Marker redundancy and map curation resulted in a decrease of the number of bridges (5578 between instead of 10086)



		After ma	ap curatio	n	Map integration (using <i>Biomercator v3</i>)
					$\lambda / (the O rescale O) / areal TO$
Map		No unique	No loci for	map	with 2 maps, GV and 13,
code		loci	integratio	n	given «higher confidence»,
	Total				a HDC map was build in <u>2</u>
T3	2072	677	1977		<u>steps</u> :
GV	1745	778	1671		
PK	2584	1352	2559		► Step 1: integration of GV
TH	2247	796	1892		and 13, <u>no reference</u> (result=GVT3)
E3	2316	480	1814		
CH	1080	33	417		Step 2: iterative
Total	12044	4116	10330		projections of PK, TH, E3
					and CH with GVT3 as fixed
					reference

Step 1 *(Biomercator):* integration of GV+T3, 2 maps of 1st priority (258 bridge loci) ► **GVT3** (3538 cM, 3374 loci)

GV	BNL1707 (0)	BNL0686 (0) DPL0222 (1,6)		28 (0)CLU2941 (0) (0,6)	T3
c9		TMB0177 (5,2) GA486 (6,3) TMB2916 (7) CLU233 (9,3) TMB0670 (11,3) C9 3_9	c23	(5) 6 (5,3) 4 (6,3) MUSS300 (13,7)	c23
	B5H7 290 (11,2) A1737 (13,1) B4IM4-M (16,2) B3IM55 (16) B3IM55 (16) B6M9 168 (18,1) BNL0686 (19,5) pAR124 (17,2) B3IM55 (16) BCM9 168 (18,1) BNL0686 (19,5) pAR127 (21,8) BCM9 268 (18,2) BSM5 240 (31,3) NAU3155 (32,9) BIM4-M1 (37,5) B4M4-M1 (37,5) B4M4-M1 (37,5) B4M4-M1 (37,5) BASP230 (53,7) CLU15 (50,5) JBSPP230 (53,7) CLU15 (50,5) JBSP230 (53,7) CLU15 (50,5) JBSM1 450 (61,1) BML3647 (64,5) CG11 (66,5) BM13502 (67,6) M92 (70) A1471 (71,9) BIM37M1 (74,9) BML317 (81) pA84M1-M (84,6) NAU0868 (85,1) C11151 (87,4) B43M2-M (90,5) MUSS022 (103,6) BML2590 (94,5) CLU2077 (97,9)	CHB2916 (7) CLU233 (9,3) CLU234 (14,8) BNL162 (17,1) CIR019 (17,9) CLU234 (14,8) BNL162 (17,1) CIR019 (17,9) CLU334 (.12,2,8) JESPE734 (22,8) JESPE734 (25,1) CLU352 (25,9) MUSB0965 (25,9) MUSB0965 (25,9) MUSS100 (50,5) BNL3179 (44,6) CLU352 (42,7) BNL315 (52,6) CLU384 (61,4) BNL315 (52,6) CLU384 (61,4) BNL3140 (66,2) DU5050 (64,1) BNL3140 (66,2) DU5050 (64,1) BNL3140 (66,2) DU5050 (64,1) BNL3140 (66,2) DU5050 (64,1) BNL3147 (75,5) CLU286 (77,6) UCcg107 (79,4) JESPR290 (85,8) CLU397 (86,8) BNL1317 (75,5) CLU297 (86,8) BNL1317 (75,5) CLU296 (77,6) UCcg112 (102,3) UCcg111 (121,3) NAU2354 (123,7) JESPR295 (124) DU50641 (125,5) CLU2700 (115,4) DU50641 (125,5)	C23 GN CLU10 CG13 (C BNL3383 PAR547 R41M1-7 CACC63 BNL269 BNL3511 CG01 (C CG01 (4 (6,3) MUSS300 (13,7) (21,3) TMB0670 (22) 22,5) TMB0157 (22,7) 3 (26,2) BNL3383 (25) (35,3) MUSS0073 (30) (35,3) CLU10 (20,3) (35,3) MUSS0973 (30) (35,3) MUSB0641 (34,6) (37,7) CLU1509 (38,2) 0 (38,5) JBSPR274 (41) TMEB109 (38,2) JBSPR274 (41) TMESB1040 (45,3) CLU15 (53,3) 1 (54,4) MUSB1040 (45,3) 7 (48,8) CLU15 (53,3) 2 (50) CLU33 (63,4) 1 (54,9) CLU33 (63,4) 55,6) CLU33 (67,1) (70,9) TMB0170 (70,9) (70,5) CLU23 (62,1) TMB0170 (70,9) CLU2742 (78,9) (77,7) CLU2742 (78,9) (77,7) CLU2742 (78,9) (87,3) CLU151 (91,1) (9,9) CLU2048 (87,5) (9,1) BNL12590 (100,9) (102,5) CLU2077 (103,6) DPL0518 (105,3) JBSEPR114 (108,7) MUSB1049 (112,3) CLU2269 (114,8) <tr< th=""><th></th></tr<>	
	CG21 (125,2)	BNL0597 (130,4)		38 (119,3) UCcgl13 (125,3 UCcgl13 (125,6 DPL0530 (127,1) JESPR095 (127,6)	
	At60 (131,2)	CLU352 (139,4)	B3M1_31	10 (128,5) BNL3985 (131,2) NAU0864 (133,4) MUCSU72 (134.9)	
GV)	9	9 (T3)	(GV)	23 (T3)	

Step 2 (*Biomercator*): Consensus of GV+T3 used as « fixed » framework for projecting markers from 4 other component maps





HDC map represents:

 8254 loci (317/chr)
 6669 « unique » markers, 3450 SSRs (727 as 'CLU'), 1894
 RFLPs
 4070 cM, 2 loci/cM
 a single dense region per chrom.

1292 multi-copy markers, 64% homeoduplicates

ue markers of the HDC map are a sequence (genomic, cDNA)

3- Convergence HDC map / G raimondii

► alignment of the sequences of :

- the 4744 markers mapped on the HDC map and
- the 13 scaffolds (750 Mb) of *G raimondii* (www.phytozome, v2.1, Sept.2012)

▶ using **BBMH**, Best Blast (Bidirectional) Mutual Hit approach, (BLASTN, 1e-5)

► the 3607 single hits (4592 loci) included 2182 hits (2369 loci) for the D chromosomes of HDC map (c14-c26), or 182 anchor loci per chromosome



The sequences originate from different ploidy/evolutionary context (1) from a **tetraploid** (markers on the tetraploid HDC map) and (2) from a **diploid** context (*G raimondii* genome)

3- Convergence HDC map / G raimondii

excellent convergence ratio kb/cM highly variable (average 400 kb/cM); S-shaped curve (high ratio in centromeric regions)



<u>after</u> clustering (5 colinear pairs): possible incongruency
 an internal region of Gr_Chr08 with not hit on genetic map



<u>after</u> clustering (5 colinear pairs) : possible incongruency
 dual correspondance of Gr_Chr09 with c19 and c22



last but not least ...

New data on TropGeneDB/Cmap (www.tropgenedb.cirad.fr)



CLU markers serve as bridges between genetic maps (thanks to aliases/feature names associated with markers)

IN ADDITION :

added (TropGeneDB, Cmap) the *G raimondii* v2.1
 physical map and correspondances/HDC b more QTLs
 (metaQTLs, expression eQTLs hotspots) b created links from markers (Cmap) to a genome browser of *G raimondii*



 (list of) candidate gene(s) in QTL confidence interval tropgenedb.cirad.fr/
 Any request for additional map and QTL data to be included in TropgeneDB, email to: marc.lacape@cirad.fr or chantal.hamelin@cirad.fr

Perspective: exchange with CottonGen

This research is published in *PLosONE:*

A High Density Consensus Genetic Map of Tetraploid Cotton That Integrates Multiple Component Maps through Molecular Marker Redundancy Check

Anna Blenda^{1,2}, David D. Fang³, Jean-François Rami⁴, Olivier Garsmeur⁴, Feng Luo⁵, Jean-Marc Lacape⁴*

1 Department of Genetics and Biochemistry, Clemson University, Clemson, South Carolina, United States of America, 2 Department of Biology, Erskine College, Due West, South Carolina, United States of America, 3 Cotton Fiber Bioscience Research Unit, USDA-ARS-SRRC, New Orleans, Louisiana, United States of America, 4 CIRAD, UMR AGAP, Montpellier, France, 5 School of Computing, Clemson University, Clemson, South Carolina, United States of America

Acknowledgments and Contributors

Acknowledged are

Andrew Paterson and JGI for early access to G raimondii scaffolds
 Johann Joets and Olivier Sosnowski for early access to Biomercator v3

OPEN O ACCESS Freely available online

Contributors:

- Anna Blenda (Clemson Un.)
- David Fang (ARS)
- Feng Luo (Clemson Un.)
- Jean-François Rami (Cirad)
- Olivier Garsmeur (Cirad)
- Chantal Hamelin (Cirad)
- Gaëtan Droc (Cirad)

Thank you, Merci ...

