

Agriculture et Agroalimentaire Canada

Genotype-by-sequence in oat

Nick Tinker, AAFC Ottawa



Acknowledgements

Coauthors: Jesse Poland Eric Jackson Shiaoman Chao Gerard Lazo Becky Oliver



Tinker Lab: Charlene Wight Kyle Gardner Phil Couroux Jiro Hattori Biniam Hizbai Benazir Marquez Xiaomei Luo Arsh Singh

CORE:

Rick Jellen, Marty Carson, Howard Rines, Don Obert, Joe Lutz, Irene Shackelford, Abraham Korol, Aaron Beattie, Åsmund Bjørnstad, Mike Bonman, Jean-Luc Jannink, Mark Sorrells, Gina Brown-Guedira, Jennifer Mitchell Fetch, Steve Harrison, Catherine Howarth, Amir Ibrahim, Fred Kolb, Mike McMullen, Paul Murphy, Herb Ohm, Brian Rossnagel, Weikai Yan, Kelci Miclaus, Jordan Hiller, Jeff Maughan, Rachel Redman, Joe Anderson, Emir Islamovic And others









Imagine the complete oat genome!



← (this genome stretches across China) \rightarrow

Imagine every complete oat genome !



State-of-the-art DNA sequencing

- Make 100's of copies
- Put them through a shredder
- Try to put them back together



e.g. 150x coverage = 20 billion pieces x 100 bp = \$100 K



Until then....

- Focus on differences among varieties
 - That's what we care about the most
- Order differences by linkage (count recombinations)
- Associate with phenotype (also by linkage)
 - Mapping populations (2-parents, lots of kids)
 - 'Natural' populations (unknown family structure)



Single Nucleotide Polymorphism (SNP)

• The most common genetic difference



- SNP = SNP no matter how you find it !
 - "Old" non-sequence-based methods (AFLP, DArT)
 - Discover by sequence / assay by design
 - Discover and assay by sequencing (GBS)

SNP – discover by sequence, assay by design



CORE – Illumina SNP array



- 20 varieties
- 9 million reads
- 18,000 templates



- 25 varieties
- 4 million reads
- 12,000 templates



...TGATCGCTA [G/T] CTGGCATGGCT......

- 80,000 predicted SNPs
- 4600 tested SNPs
- 2300 validated SNPs (Golden Gate)
- 6000 SNPS in progress (Infinium)
- (we estimate 4000 will work)



Genome Studio Software (example SNPs)







Genotype by sequence (GBS) - concept

- Discover and assay SNPs by direct sequencing
- Similar to SNP discovery for planned assay
 But <u>much</u> larger numbers of and sequences
- Based on subset of genome (enzyme / amplify)





GBS details

Development of High-Density Genetic Maps for Barley and Wheat Using a Novel Two-Enzyme Genotyping-by-Sequencing Approach

Jesse A. Poland^{1,2}*, Patrick J. Brown³, Mark E. Sorrells⁴, Jean-Luc Jannink^{4,5}

1. Digest with *Pstl* & *Mspl*



2. Ligate sequencing adapters + variety-specific barcode



GBS - details 3 billion sequenced fragments

5. Trim barcode, trim to 64 bases, keep track of variety

- 6. Identify all unique tags, count in each variety (Tag x Taxa)
- 7. Match tag pairs, call SNPs (across full data set)



GBS - caveats

- Missing data 10% to 70%
 - Depends on sequencing depth ("plexity")
 - Depends on how many SNPs you call
 - e.g. to get 95% complete, I could only call 2400 SNPs
- Bioinformatics is under development
 - Different methods give different SNP sets
- Massive storage and computing requirements
- Data sets too large for some software
 - More markers <u>and</u> more samples
 - Have not yet managed a consensus map

DArT vs SNP vs GBS

	DArT	SNP	GBS
Assay cost per sample	\$50	\$68	\$20 *
Markers across all taxa	1500	4000 *	40 k → 100 k
Markers per population	300	800	4 k → 20 k
Missing data	>5 %	< 1%	10% → 50%
Co-dominant	0%	> 25%	100% *
Genes / orthology	20%	100%	5%
Duplicate loci (map inconsistently)	> 5%	< 3 %	?

Are they useful ?

- Already have more GBS data than anything else:
 - Not just more loci..... more varieties too
 - DArT: 350 diversity + 4 bi-parental populations
 - SNP: 108 diversity + 6 bp-pop (400)
 - GBS: 738 diversity + 8 bp-pop (700) + 16 iso-lines
- 10 x more likely to find [marker QTL] ?
- But missing data ...
 - Mapping difficulty ?
 - Association artefacts ?
 - MAS predictions ?

Simulated MAS using GBS

- Scenario:
 - 4 target loci, simulate with random GBS loci
 - Discover markers by association in odd # lines
 - Predict genotypes in even # lines

Marker	f(A)	f(B)	f(H)	Chr
avjp100014	568	782	20	1C
avjp100585	339	966	18	16A
avjp100261	150	1188	9	9D
avjp108144	1196	143	8	7C_17A

Predictive markers (target locus avjp100014)

204 "good"

141 "bad"



Predictive markers (target locus avjp100261)

62 "good"

287 "bad"



Simulated MAS using GBS

- Surprise!
 - Many loci in (near) perfect LD with a "QTL"
 - Bias in genome regions sampled by GBS ?
 - Too much good data ?
- Test predictions of array-based SNPs
 - Independently discovered
 - Most based on expressed genes
 - BUT.... SNP data only available for 108 varieties

Simulated QTLs based on Illumina SNPs

- 4 random Illumina cDNA SNPs (targets)
- Fit model with 54 odd # lines
- Test model with 54 even # lines

Marker	f(A)	f(B)	Chr
GMI_ES01_c10033_104	56	52	1C
GMI_ES15_c10388_464	57	51	18D
GMI_ES01_c796_180	84	24	19A
GMI_ES15_c8064_341	25	83	3C

Predictive markers (locus GMI_ES01_c10033_104)

27 "good"

27 "bad"

	1 1																																	
GMI ES CC9711 432	11	1 1 1	1	11	11	1 1	1	1 1	1 1	11	1 1	1	1 1	1	1 1	1	3 0	0	3 3	3 3	3 3	3	3 3	3	3 3	3	3 3	3 3	3 3	3	3 3	3	3 3	3
GMI ES01 c16239 143	11																3 3	3																
GMI ES01 c8817 242	11																																	
GMI ES01 lrc28711 407	11	0 0 1			10																												3 0	3
GMI_ES02_c2898_314	10	111		1 1	0 1		0										3 0	3	3 0	3 (0 0		3 0	0									3 3	3
GMI_ES17_c2941_220	11			1 1	1 1		1										3 3	3	3 3	3 3	3 3		3 3	3										
avjp8646	10									11	0 1	1	0 1	0								3	0 3							3	0 3			
avjp86831	0 1	110	0 1	0 1	0 1		1	0 0	2 (0 1	0 1	0	0 0	0	0 1		3 0	3				3	3 3				3 0	3	0 3	3	з з			
avjp109287	1 1	0 1 1	l 1	1 0	0 2		0	1 0	2 (0 1	0 1	1	0 1	0	0 1	0	3 3							2			3 3	3 3	3 3				3 2	3
avjp24905	11	1 1 1		1 1	1 1		1	11	1 1	11	0 1	1	0 1	1	1 1	0	3 0	3	3 0	3		3	0 3	0		3	0	0		0			3 3	
avjp44794	2 1	110	1	10	1 1	2 1		1 1	2	0	0 1	0	0 1	0	0 1	0	3 0	3	0 3	3 (0 0	3	0 0	3		0				3	0 3			
avjp95770	1 1		0	1 0				1 1	2 (0 0	0 1			0	0 1	0	3 0	3		0	0 0	3	0 3		3 0) 0	0		3 0	3	0 3		3 0	3
avjp60667	01	1 1 0	0 0				. 1	0 1	1 (0 1	0 1			0	0 0	0 1	0 0	3		0	3 0	0	0 3		3 0	3			3 0	3	0 0	0		
avjp2085	1 1	1 0 0	0 0	10	10	1 0	1	0 1	1 (0 1	0 1	0	0 1	0	0 0	0 0	3 0	3		0	0 0	0				3	0 3			3	0 3			
avjp43761	01	110	1	1 0		1 0	1	1 0	2	0	0 1	1	0 1	0	0 1	0					3 0	0	3 0	3		0	0			3	0 3			
avjp94341	10						1	0 1		0	0 1			0	0 0	0 0	3 0	3	3 0	3						0	0			3	0 3			
avjp45798	10	<mark>1 1</mark> C	0 0	0 1	10		1	1 0	1 1	11	0 1	1	0 1	0		0	3 0	0	03	3 (0 0		3 0	3	0 3	6 З	0	3-3	0 0	0	0 0	0		
avjp90887	01	111	0	1 1	0 1		0	0 1	2 (0 1	0 1	1	0 1	0	0 1	0	3 0	0																
avjp1651	1.1	<mark>1 1</mark> C	1	1 0	0 0		. 1	0 0		0	0 1	1	0 1	0	0 1	l 1	3 0	3	03			3	0 0	3		6 З	0	0	0 3	3	0 3	0		0
avjp24599	10	110	0 0	1 1	0 1				1 (0 0	0 1	0	0 0	1	0 1	. 0	3 0	0		0		3	0 0	3		0	0 3	0	0 3	0	0 0	0		0
avjp45799	10	110	0 0	0 1	10	1 1	1	1 0	1 1	l 1	0 1	1	0 1	0	1 1	. 0	3 0	0	0 3	3 (0 0	3	3 0	3	0 3	3	0 3	3 3	0 0	0	0 0	0		
avjp19552	10	0 <mark>1</mark> 0	0 0	0 1	1 1	1 0	1	0 0	1 (0 0	0 0	1	1 0	1	0 0	0 0	3 0	3		0	3 0	3	0 0	3		0		0	0 0	3	0 0	3		
avjp93294	00	00	0 0	0 0	10	1 0	0	0 0	1 (0 0	0 0	1	0 1	0	0 1	0	3 0	3	3 3	3	3 0	3	0 0	3	3 3	, З		0	2 3	3	0 3	3	3 3	3
avjp106786	1 1	110	1	0 0	10	1 0	1	1 1	0	0 0	0 1	0	0 1	0	1 1	0	3 0	0	0 0	0	D 3	0	0 0	0	0 0) 0		3-3	0 0	3	0 0	3	0 3	0
avjp13850	10	0 0 1	0	0 0	0 0	1 0	0	0 1	0	0 1	0 0	1	0 0	0	0 0	0 1	3 0	0		3 (D 3	3	3 3	0	0 0) (3_3	3 0	0	0 0	13	0 0	0
avjp31268	$1 \ 1$	0 <mark>1</mark> 0	0 0	1 0		1 0	0	0 1		0	0 1	1	0 1	0	1 0	0 0	3 3	3	3 3	3	3 3	0	0 3	3	3 3	0		0		0	3 3	3		3
avjp17734	10		11	10	1 1	1 0	0	1 1	1 1	11	0 1	1	0 1	0	0 1	0	0 0	3	3 0	3	3 0	0	0 3	0	0 3	0	3 3	0		0	3 0	0		0
avjp38924	01	111	l 1	0 1	1 0	1 0	1	11	0	0	0 0	0	0 1	0	0 1	0	3 3	3	3 3	3 (0 0	0	0 3	0	0 3	0	3 0	3		0	0 0	0		
avjp63997	10	L 0 C	0 0	0 1	10	1 0	0	0 0	3 (01	0 1	0	0 1	0	0 1	. 1	3 0	3	0 0	0	3 0	0	3 0	3	0 3	0		3_3		0	0 0	0		3
avjp64036	1 1	00	0 0	0 1	01	1 0	0	0 1		0	0 1	0		0	1 (0 0	3 0	3		3	0 0	0	0 0	3	3 0) ()		0	3 3	3	0 3	0		0
avjp69451	10	110	0 0	0 0	10	1 0	1	11		0	0 1	0	1 1	1	0 1	. 0	3 0	3	3 3	3	0 0	3	0 0	3	3 0	3		0	0 0	0	0 3			3
avjp77615	10	200	01	0 0	0 0	2 2	0	0 0	1 1	ι 1	0 1	1	0 1	0	1 1	0	3 3	0	0 3	0 (D 3	3	0 0	3	0 0) 3		0	0 0	0	0 3	3		0
avjp83537	10	100	0 0	0 0	10	1 0	0	1 0	1 (0 0	0 1	1	0 1	0	0 0	0 0	3 0	3	3 3	3 (0 0	0	0 0	3	3 3	0	3 3	0	0 0	3	0 0	3	3-3	3

Predictive markers (locus GMI_ES15_c8064_341)

41 "good"

13 "bad"



Case study with naked (hulless) oat

- Other genes (N2,N3,N4) or will we just find N1 ?
 - N1 was Mapped in Terra x Marion
 - Poorly placed by comparative mapping
- GBS data:
 - Diversity lines (*rare trait)
 - I only had phenotypes for 100 covered + 20 naked
 - 100 Terra x Marion progeny
 - 8 pairs of naked / covered iso-lines
 - OT253/Marion, from F6 heterozygotes
 - Developed by Solomon Kibite

Look at all pairwise linkages and LD's

$N1 \leftrightarrow GBS$ -tag



N1 predictions

- LD analysis based on partial data:
 - "training set" (the only ones I had data for at the time)
 - 20 naked / 100 covered
- Based on "avjp23455" we predicted:
 - 12 more naked lines (remaining 600 = covered)
 - All were correct!

Training set	
95Ab13050	OT253
Boudrias	Provena
Bullion	Racoon
Gehl	Salomon
IOI033	Shadow
IOI108	Terra
IOI114	VAO-44
IOI150	Vao-48
IOI191	VAO-51
Navaro	VAO58

Predicted Naked Lines

98Ab7265 FL03184-FLID-B-S1 FL04178-FLID-B-S-2 HLA05AB1-34 IL02-10836 IL03-7936 LA02012-S-B-139-S2-B-S2-B-S2 LA0210SBSBSBSB-S1 LA03066SBS-L1 Lennon Nudist Zuton

Conclusions

- GBS will be exceptionally efficient for tagging genes
 - Bi-parental AND association mapping
- Excellent potential for MAS
- SNP array or custom assays are better for
 - Genome analysis, comparative genomics
 - Critical genotyping, well-characterized targets
- Work required
 - Streamline informatics
 - Build GBS/SNP based consensus map
 - Evaluate consistency of map position
 - Evaluate genomic selection

American Oat Workers Conference Ottawa Canada, 2014 Sun. July 13 to Wed. July 16, 2014

Ottawa Marriott Hotel 100 Kent Street, Ottawa Ontario K1P 5R7 Canada



web: **aowc.ca**