

Finding function in complex crop genomes

David Edwards
University of Queensland, Australia
Dave.Edwards@uq.edu.au

Outline

- Second generation DNA sequencing technology
- TAGdb
- Candidate gene discovery
- Brassica repeat identification
- Brassica SNP discovery

Second-generation sequencing (2GS)

- Illumina GAIx and HiSeq2000
 - ↑↑↑ sequence
 - ↓ money
 - ↓ time
 - ↓ read-length
 - ↑ computation



Illumina paired reads



- Illumina GAIIx/HiSeq 2000
- Read length (100 bp)
- Insert size 300 - 500 bp

TAGdb

Welcome to ACPFG Bioinformatics .

This service performs BLAST alignment between a single query and short pair reads of selected species.

Please enter a valid email address

Note: Your result will be sent to the specified address.

Sequence data

Either: Select the sequence file to upload:

Otherwise: Enter a sequence in FASTA format:

Note: Query sequence must be less than 5000 nucleotides.

<http://flora.acpfg.com.au/tagdb>

Species selection

Please choose a query species

- Barley
- Brassica
- Diplotaxis
- Hirschfeldia
- Leptosphaeria
- Lotus
- Nicotiana
- Pongamia



Short paired-read library selection

Please select one or more paired-read libraries to search

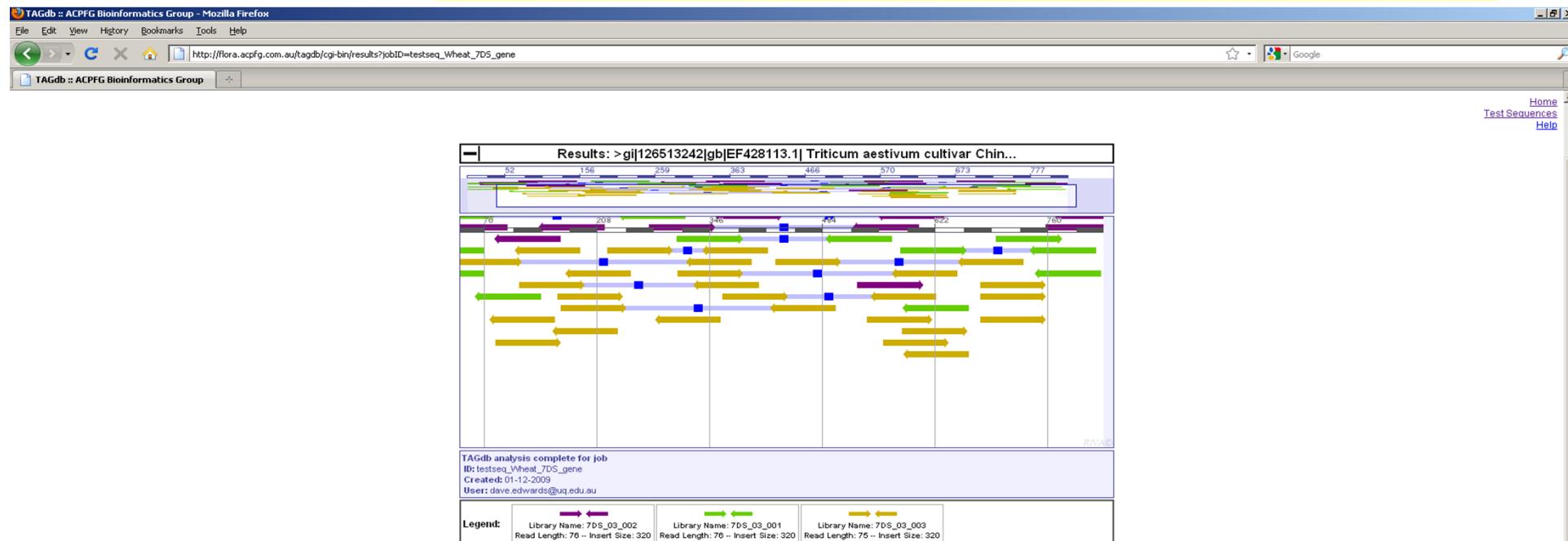
- B. rapa chiifu - 36 - 300 - BrC_03_002
- B. rapa chiifu - 35 - 2700 - BrC_27_001
- B. rapa chiifu - 35 - 2800 - BrC_37_001
- B. rapa chiifu - 35 - 2800 - BrC_37_002
- B. rapa kenshin - 36 - 410 - BrK_03_001
- B. rapa kenshin - 76 - 410 - BrK_03_002
- B. nigra - 76 - 2700 - Bni_37_001
- B. oleracea - 76 - 3000 - Bol_37_001

Format: SourceName - ReadLength - InsertSize - LibraryName

Start

Marshall, D.J., et al. (2010) *Plant Methods*. 6:19

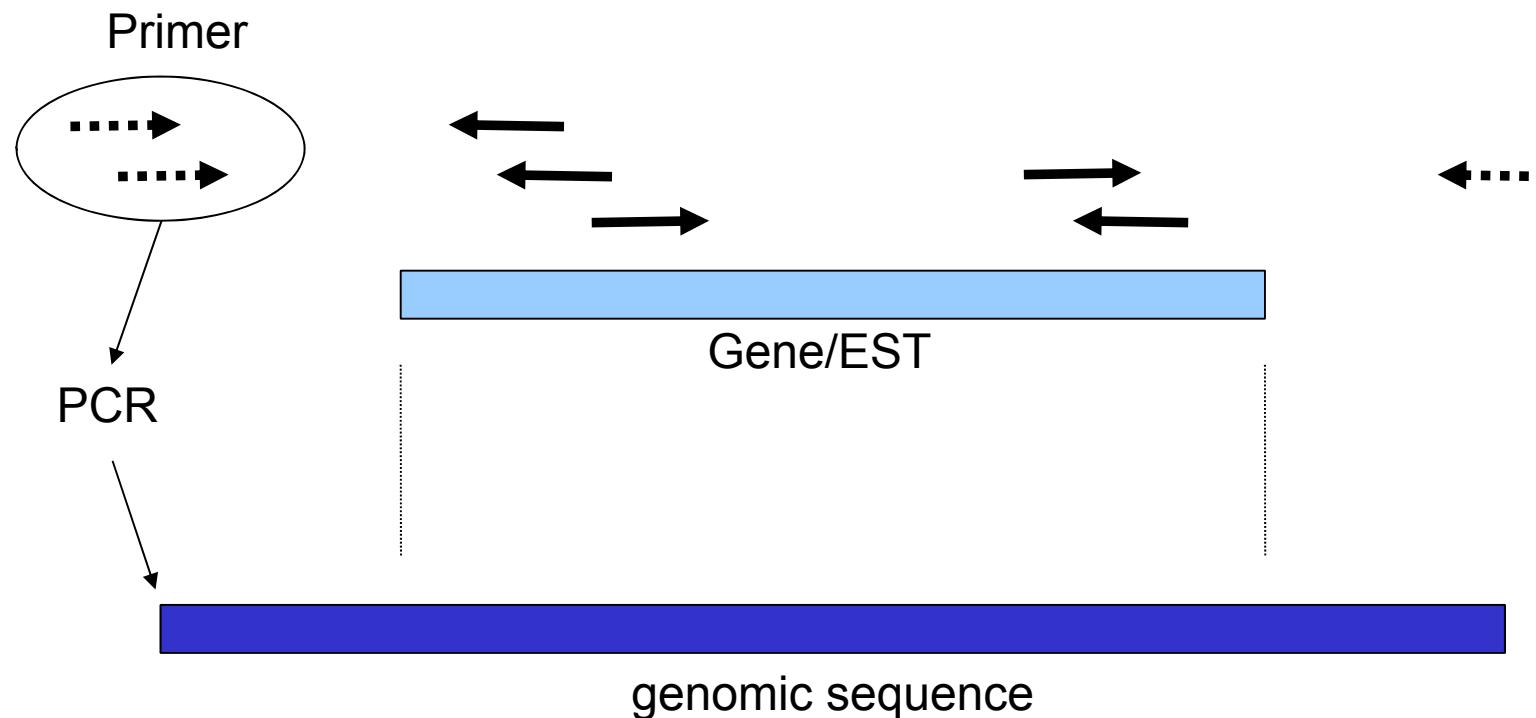
TAGdb output



Tag alignment information (download fasta)						
ID	Start	End	Tag sequence	Orient	Library	Type
1	330	405	TTCTATGCTTTCCTTTCGGGAACATTGATTTGATGCTTCTATTACATGTTGTTCATGTTG	1	7DS_03_002	P
1	553	628	TAGGGACAGAAAATTATCAACTGAGTTACAATTAACTTACTTACCGAGTGAGATACTCCCTAAGGTTGGATCG	-1	7DS_03_002	P
2	17	92	GAACAGGACCGTGTCCACCGCTGGAGCTCAAGCCGTCATGGTGGCCAGCAGGCCAGGGTTAGGTTGGCGGC	1	7DS_03_002	P
2	273	348	AAAGAAAACAAGCATAGAAAATCATGATCAACACCAAGAGAGCTAGCACAGGCCATATGCAAGATGGAGCATTAGT	-1	7DS_03_002	P
3	354	429	TGGGGAAACATTGATTTGATGCTCTATACATGTTGCTATGCAAGCAGCAGCAATTAA	1	7DS_03_002	P
3	558	633	ATATCTAGGGACAGAGAAAATTATCAACTGAGTTACAATTAACTTACCGAGTGAGATACTCCCTAAGGTTGGATCT	-1	7DS_03_002	P
4	272	347	GACTAAATGCTCATCTTCATATGGCTCTGCTCTGCTCTGGTGTGATCATGTTCTATGCTTGTCTT	1	7DS_03_002	P
4	527	602	GTTCACAAATTAGTAGCTTACCGTAGGAGATCTCCCTAACCTTCAAGGTTGGATCGCTGGACTGGACATCTGGCTAC	-1	7DS_03_002	P
5	76	151	ATATAGGCAATTAGATAGTGACTGTGACGTAGCTACGGAGGTGAGAAGGTTCTACCTCATGGCCCAACCTCAACCC	-1	7DS_03_002	A
5	-	-	TCCCTCGTATCTCTAGCTTACGGCCGGTGGATCTATACAAAGAGGAAGGGGGATGGCCGGAGGGGAGAGGCC	-	7DS_03_002	-
6	675	750	CTTGACAGATACTCCCGGTACAACTGGTGCATCTTGGCCAGGAGGGGGATGTCCTACGAGAGCCCTGTCGGC	1	7DS_03_002	A
6	-	-	CTATATACTAGTAACACTAGCCACGGTGGATAATTGCTGACTTGGGGGGGGTGGGGACGGGGGGCATGGTAGA	-	7DS_03_002	-
7	142	217	CACAGTACTATATATAGAGCTGATCATGCTGGGTGGGGAGACTTTCAAGGCTAAACATATAGGGAT	-1	7DS_03_002	A
7	-	-	TATCTAGCTAGCAGGGCGGTGATCTATACAAAGAGGAAGGGGGATGGCCGGGGAGAGAGACCGCGTGGT	-	7DS_03_002	-
8	527	602	GTAGACCCAGATGCTCAAGTCCAACGGATCCACCTTAAAGGAGGTATCTCACTGGTAAGTACTAAATTGTAAC	1	7DS_03_002	A
8	-	-	GACAGGGGGCGAGGGTTGAGAGCTGGCAAGTCCCTGGTGTGAAGTCTGGGCCACCGGGAGGGTACACG	-	7DS_03_002	-

Done

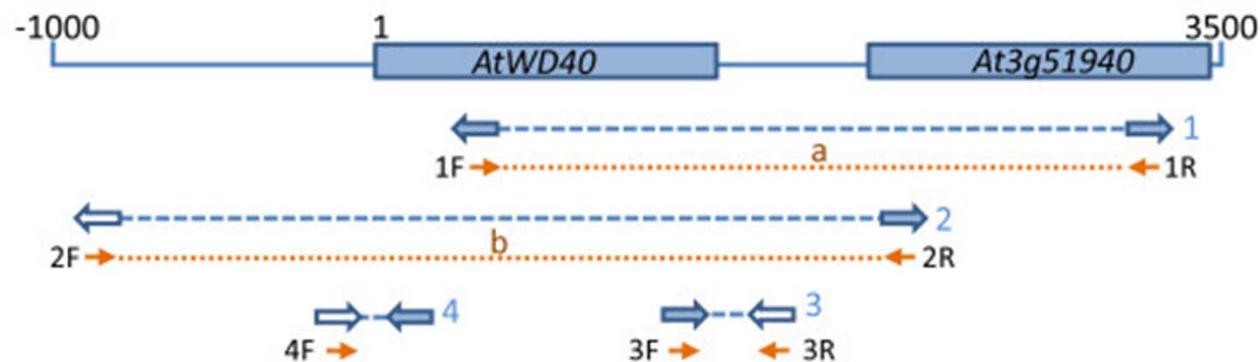
TAGdb – Gene discovery



Known 
(canola)

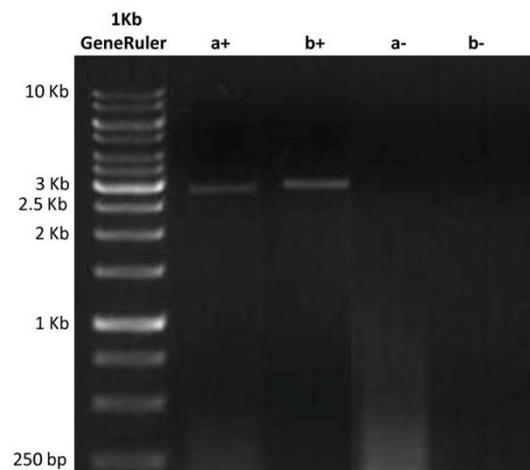
Unknown  →
(wild Brassica)

TAGdb

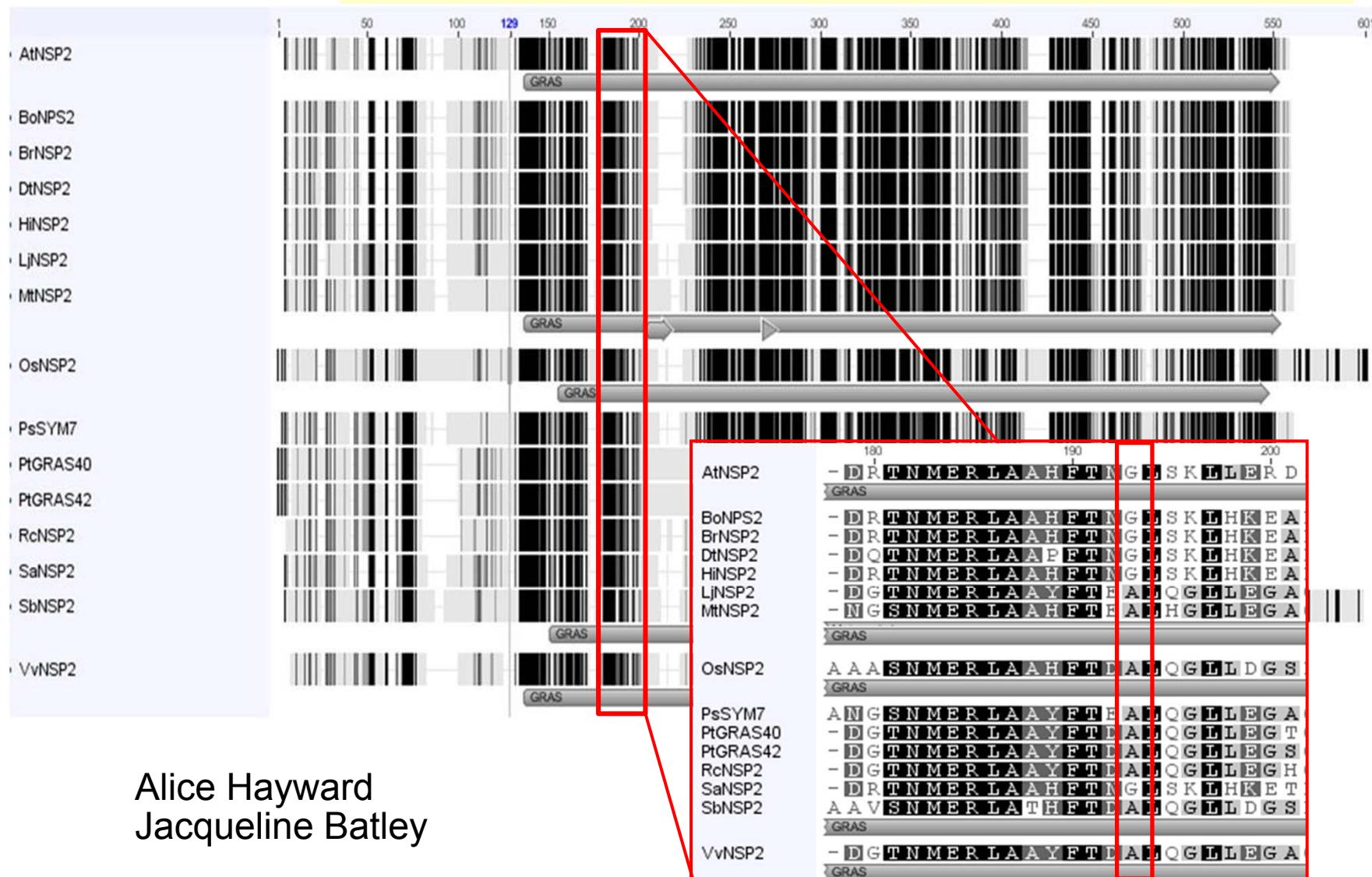


→ ← : Primers
→ ← : PCR products

: Brassica large-insert tag pairs 250 bp
 : Brassica small-insert tag pairs
 : Tags lacking significant similarity to the reference in TAGdb



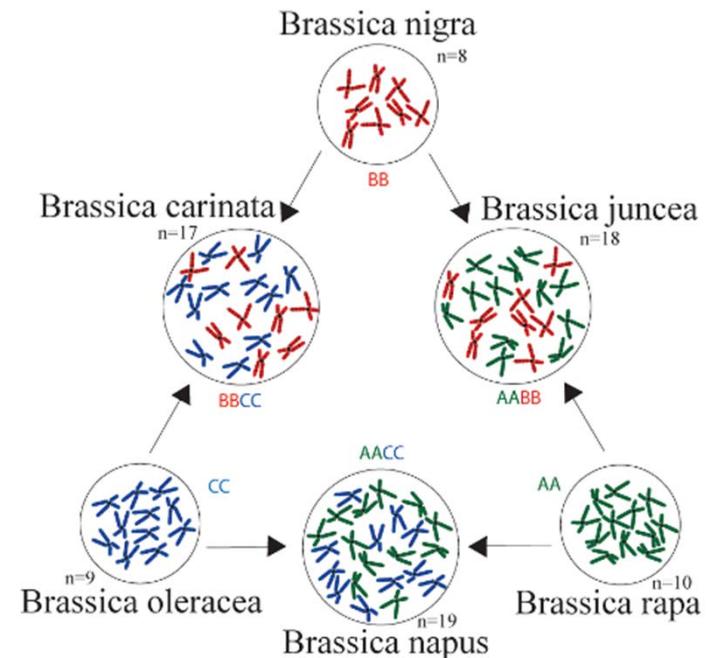
Sequencing SYM genes in *Brassicas*: NSP2



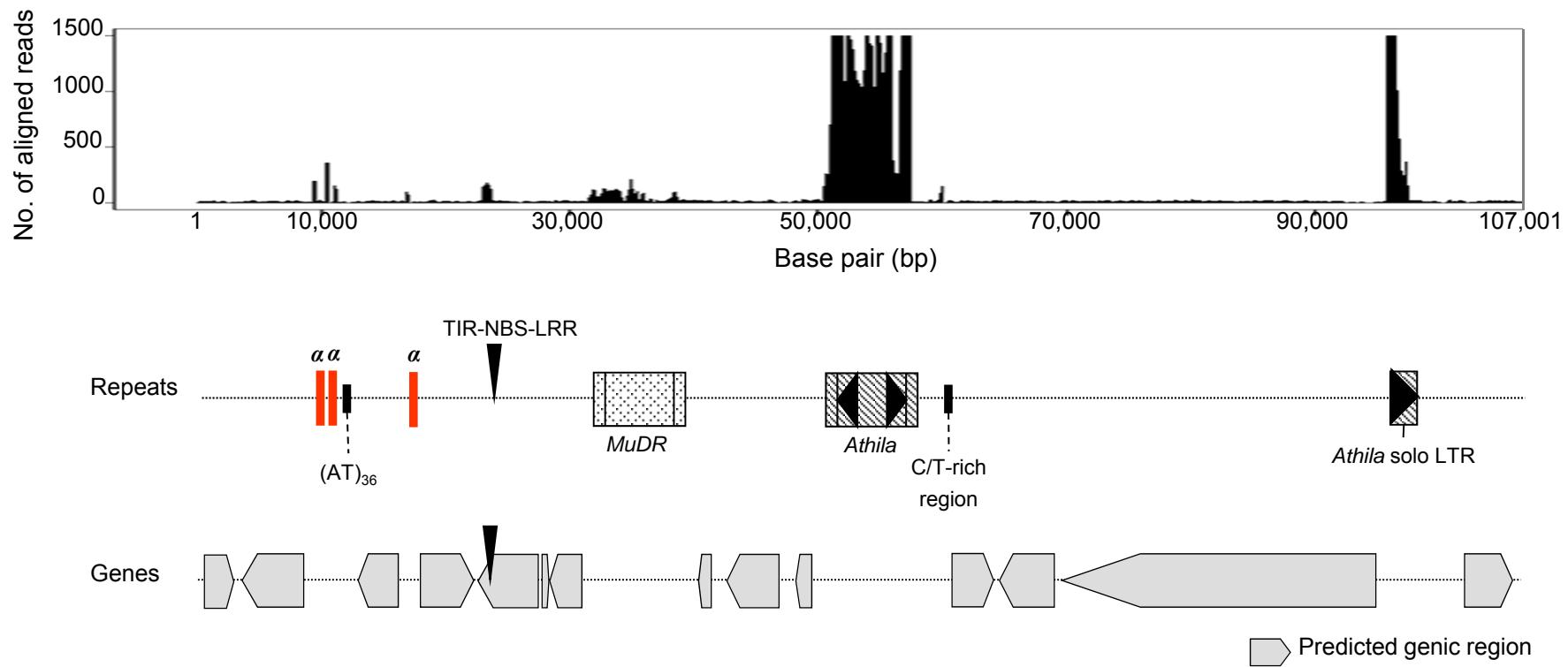
Alice Hayward
Jacqueline Batley

Brassica genomes

- Illumina GAIIx and Hi-Seq data for:
 - *B. rapa* BA, XA
 - *B. oleracea* BC
 - 8 *B. napus* cultivars
 - Wild Brassica species

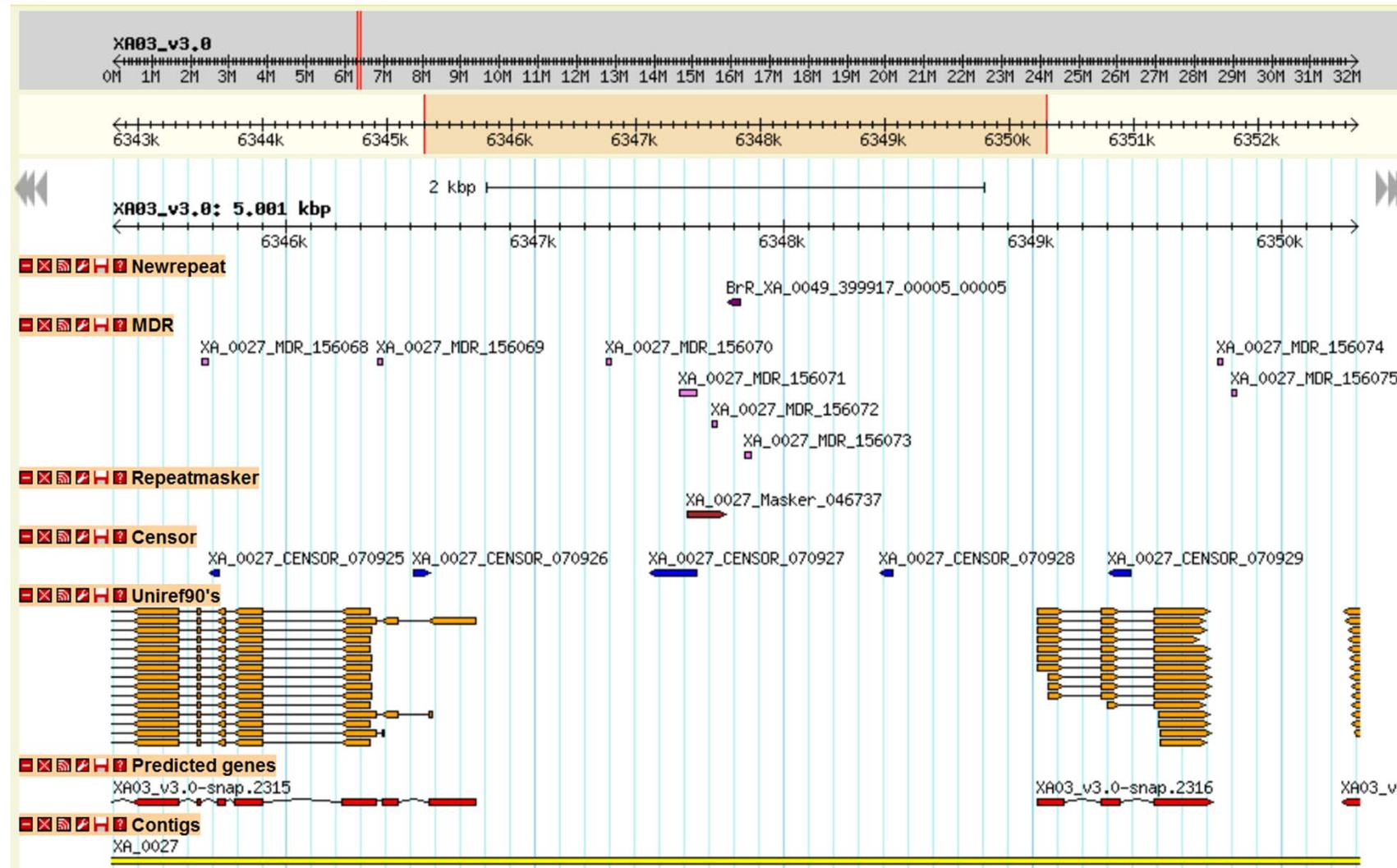


Repeat detection and annotation



High-covered regions of short reads and their corresponding annotation.

www.brassicagenome.net



Applications Places System 3.19 22 °C Fri Oct 8, 1:27 PM uqlsmits

hiruko brassica_private - Mozilla Firefox

File Edit View History Bookmarks Tools Help

http://localhost/gbrowse/hiruko/?name=Chr1_XA_V3.0%3A1..10000

Most Visited Getting Started Latest Headlines Acpfg Bioinformatics DAWGPAWS Google Vertalen blat-genomic-inhou... hiruko brassica_priv... 7DS: 470 bp from 7... localhost / localhost... GBrowse Administra... (Untitled) hiruko brassica_priv... +

File Help

hiruko brassica_private

Browser Select Tracks Upload and Share Tracks Preferences

Search

Landmark or Region: Irr Search Download Track Data Configure... Go Examples: Chr1_BA_V4.0:1..10000, Chr1_XA_V3.0:1..10000.

Data Source: hiruko brassica_private Scroll/Zoom: << < > >> Show 1.434 kbp + - Flip

The following 382 regions match your request.

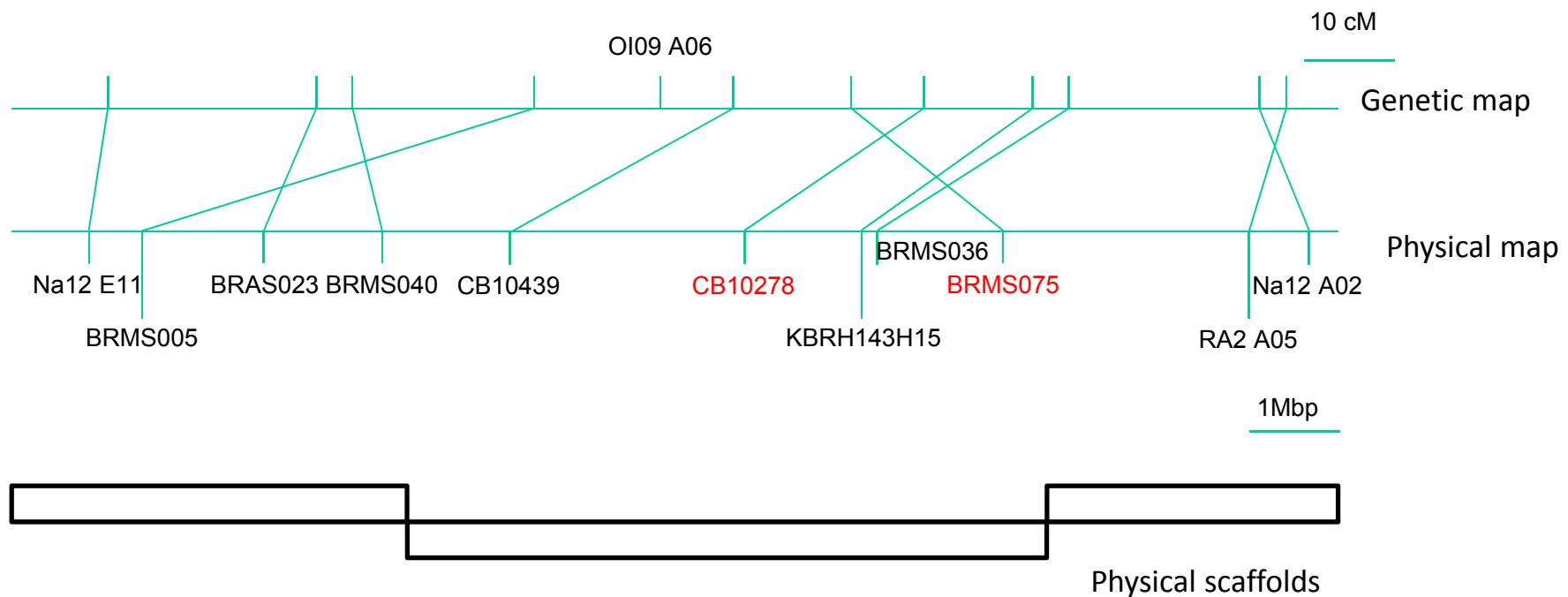
Name	Type	Description	Position	Match Score
UniRef90_Q8GUQ5	protein:blastp	Brassinosteroid LRR receptor kinase n:4 Tax:Solanum ReplID:BRI1_SOLLC	Chr1_BA_V4.0:227116..230562	10
UniRef90_Q8GUQ5_1	cds:blastp	Brassinosteroid LRR receptor kinase n:4 Tax:Solanum ReplID:BRI1_SOLLC	Chr1_BA_V4.0:227116..230562	10
Chr1_BA_V4.0.snap.416	gene:SNAP	Probable LRR receptor-like serine/threonine-protein kinase At4g36180 n:2 Tax:Arabidopsis ReplID:Y4361_ARATH	Chr1_BA_V4.0:1284157..1288031	10
UniRef90_COLGS2_1	cds:blastp	Probable LRR receptor-like serine/threonine-protein kinase At4g36180 n:2 Tax:Arabidopsis ReplID:Y4361_ARATH	Chr1_BA_V4.0:1284211..1284750	10
UniRef90_COLGS2	protein:blastp	Probable LRR receptor-like serine/threonine-protein kinase At4g36180 n:2 Tax:Arabidopsis ReplID:Y4361_ARATH	Chr1_BA_V4.0:1284211..1288028	10
UniRef90_COLGS2_2	cds:blastp	Probable LRR receptor-like serine/threonine-protein kinase At4g36180 n:2 Tax:Arabidopsis ReplID:Y4361_ARATH	Chr1_BA_V4.0:1284826..1285635	10
UniRef90_COLGS2_3	cds:blastp	Probable LRR receptor-like serine/threonine-protein kinase At4g36180 n:2 Tax:Arabidopsis ReplID:Y4361_ARATH	Chr1_BA_V4.0:1285759..1287093	10
UniRef90_COLGS2_4	cds:blastp	Probable LRR receptor-like serine/threonine-protein kinase At4g36180 n:2 Tax:Arabidopsis ReplID:Y4361_ARATH	Chr1_BA_V4.0:1287516..1288028	10
UniRef90_HP0000162920_4	cds:blastp	disease resistance protein (TIP_NBS_LRR class)	Chr1_BA_V4.0:1212407..1212421	10

Select Tracks | Select Tracks

Done

hiruko b... uqlsmi... uqlsmi... uqlsmi... Inbox (7... uqlsmi... brassica... Jupiter uqlsmi... Save Sc... Compos... 13

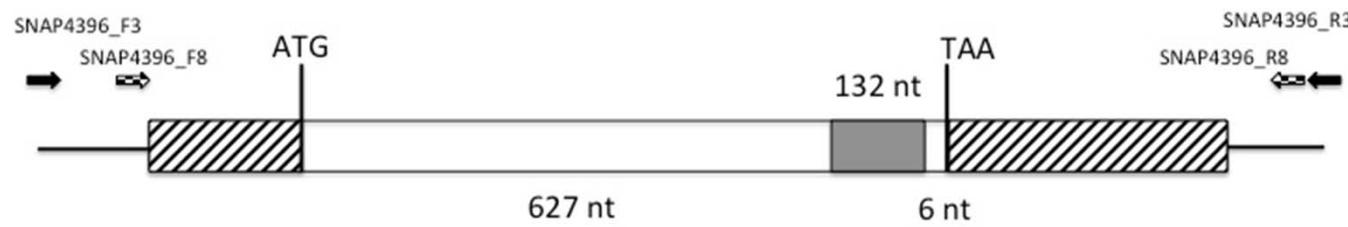
Candidate gene discovery



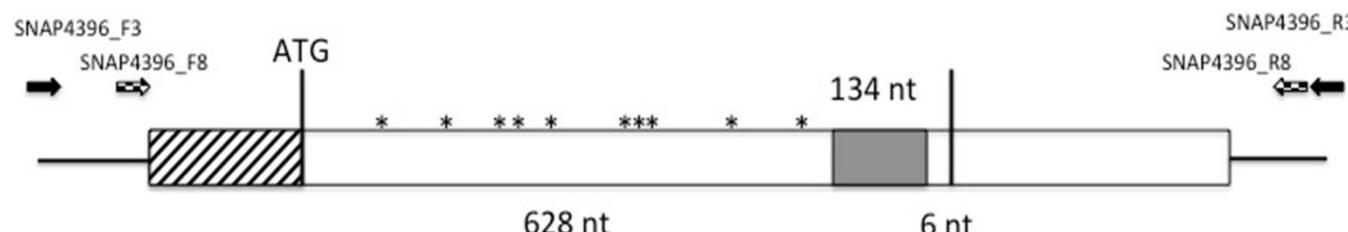
Candidate gene discovery

- Sequencing confirmed the presence of a large (402 nt) insertion in the 3' region of Ag Spectrum

Skipton



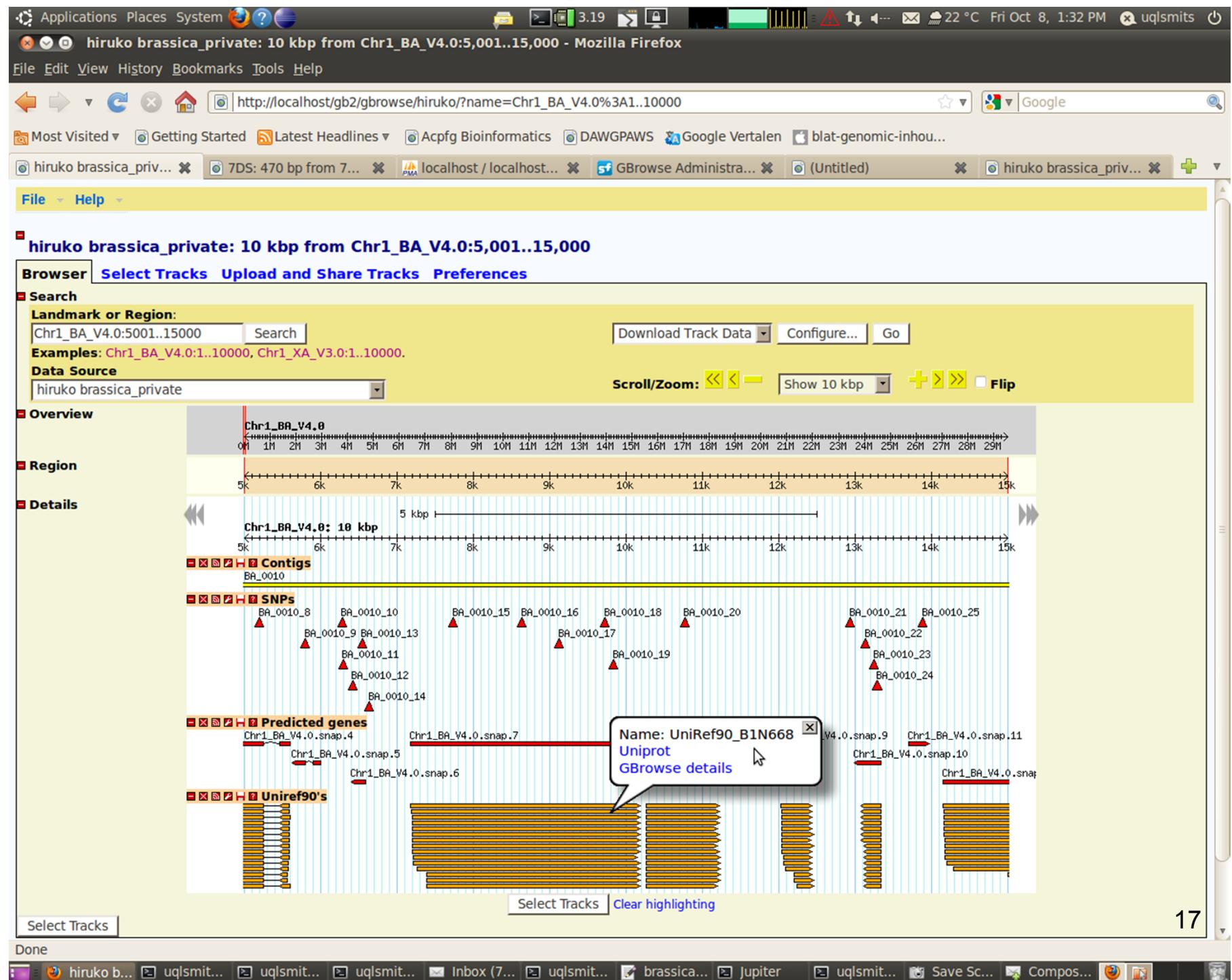
Ag Spectrum



Insertion in 3' UTR
Stop codon at position
44 – truncated protein

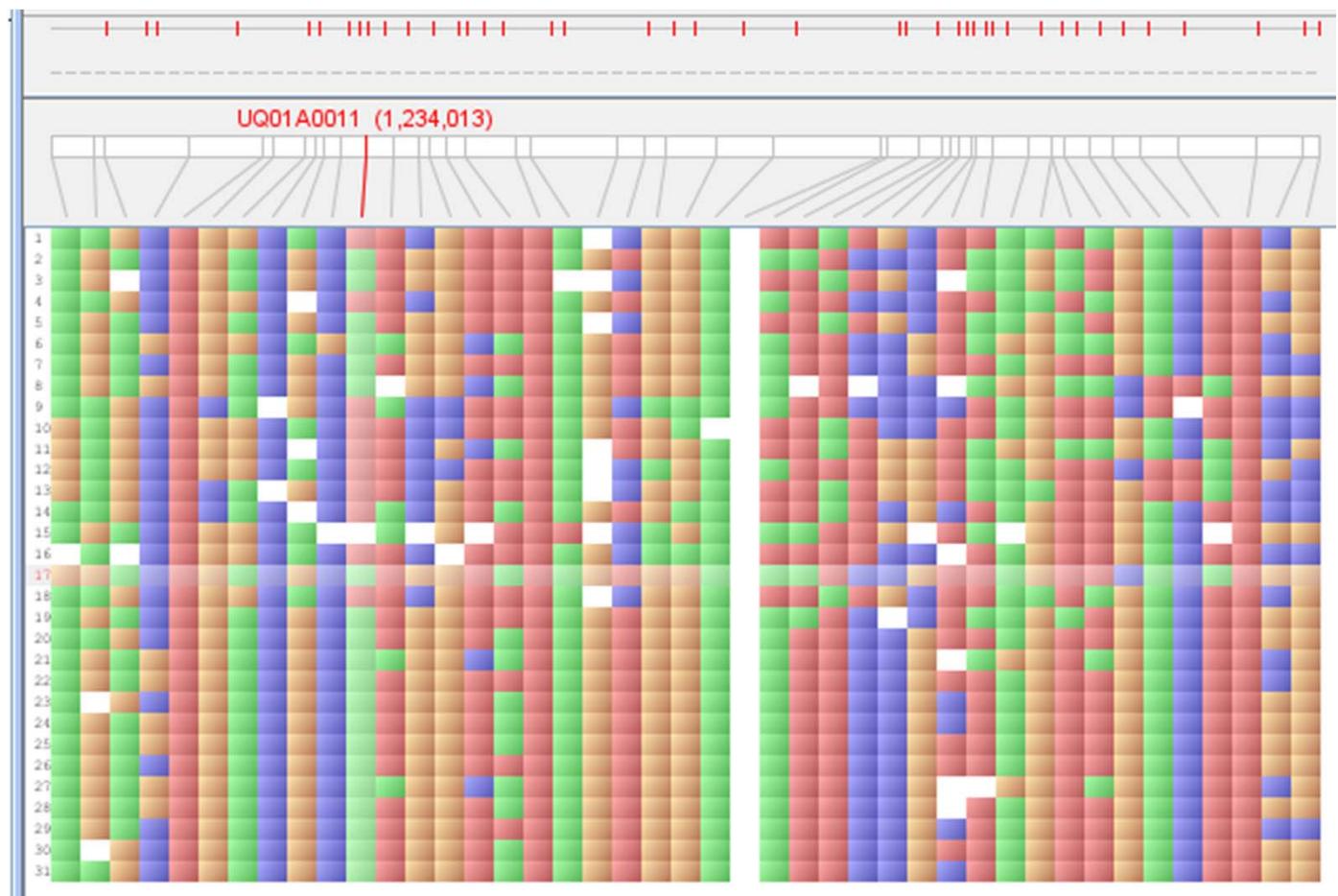
B. napus SNP discovery

- Illumina paired end sequence from 8 *B. napus* cultivars
- Map reads to reference using SOAP
- Identify varietal genomic SNPs using custom algorithm
 - input: BAM files
 - output: text, goldengate, GFF3, and VCF
- Identified > 1 million SNPs
- Validated accuracy > 96%

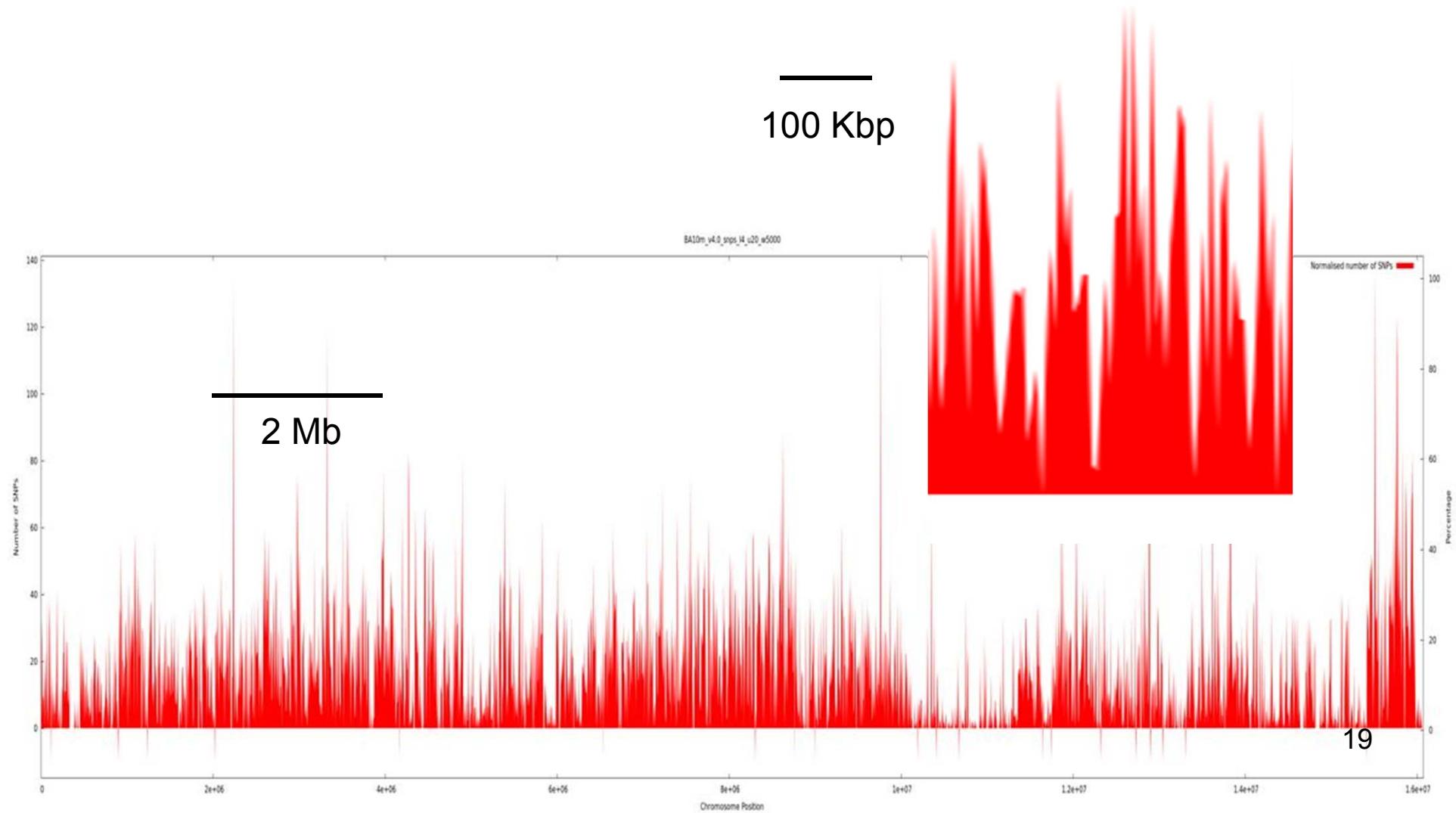


B. napus 6K infinium

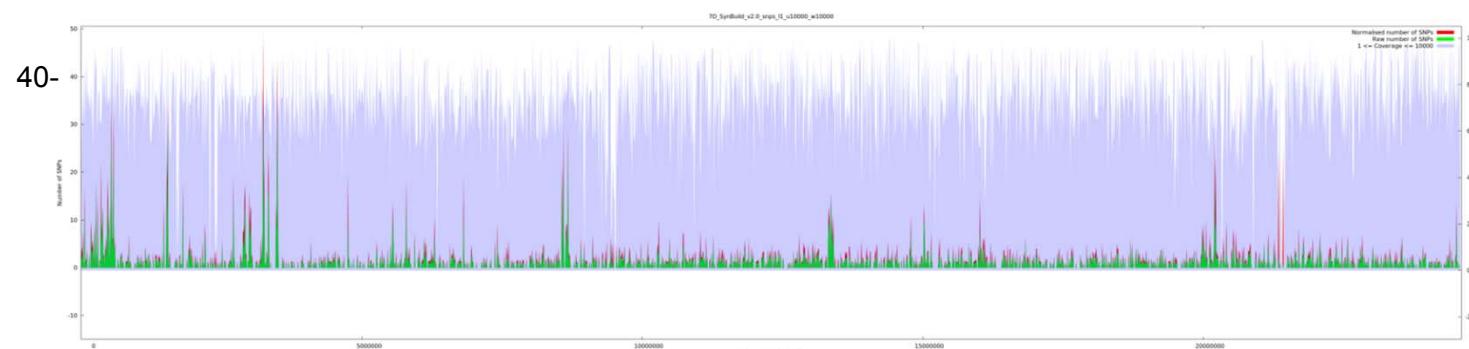
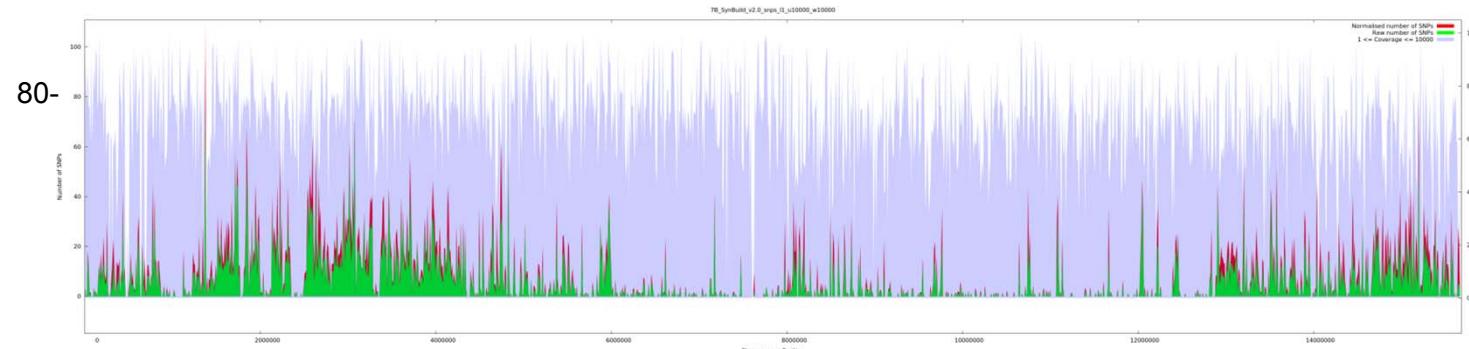
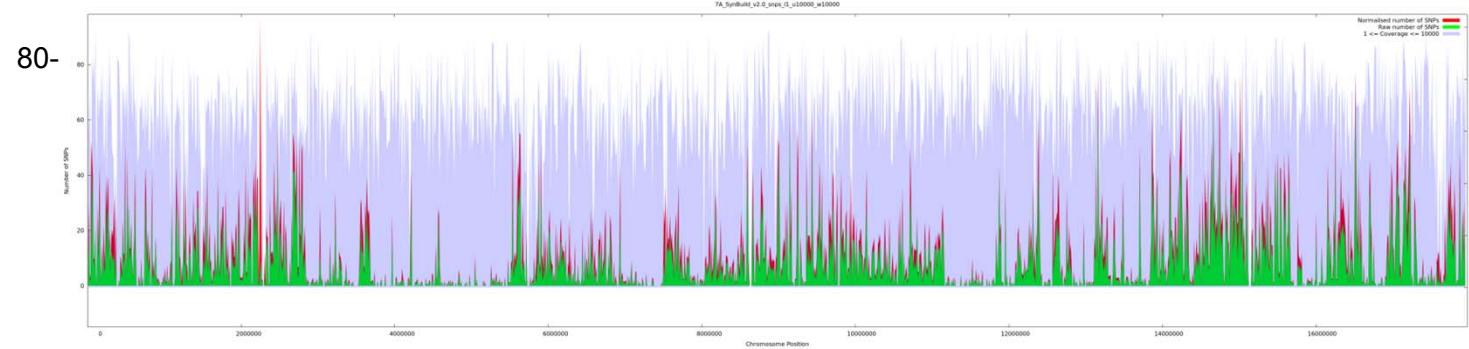
- 5306 genome wide SNPs
- Genotype >2000 lines



SNP density across A01



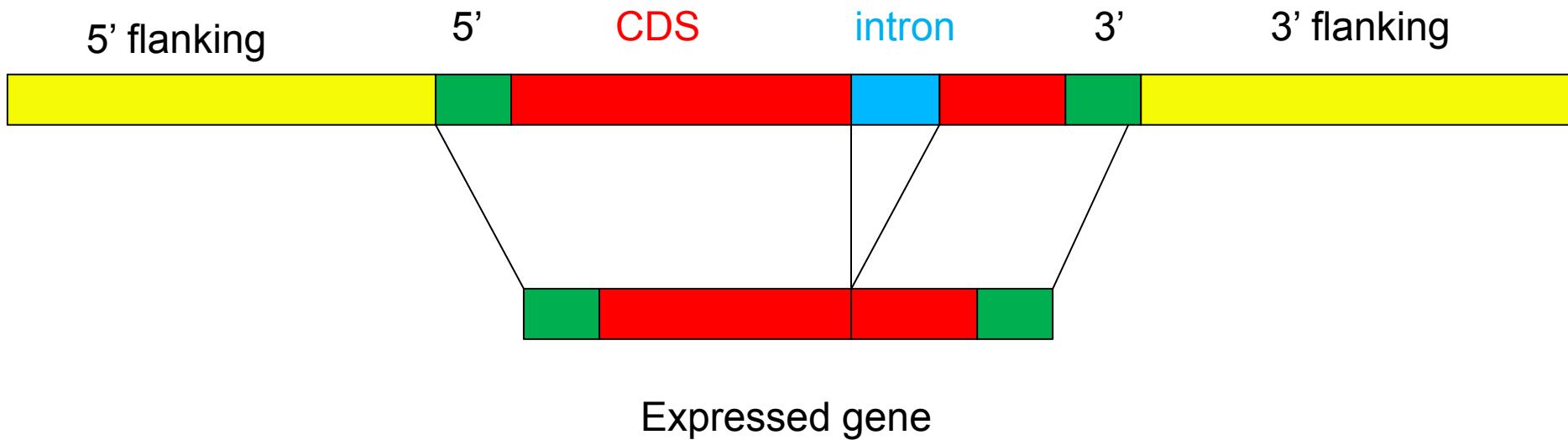
Wheat genomic SNP discovery



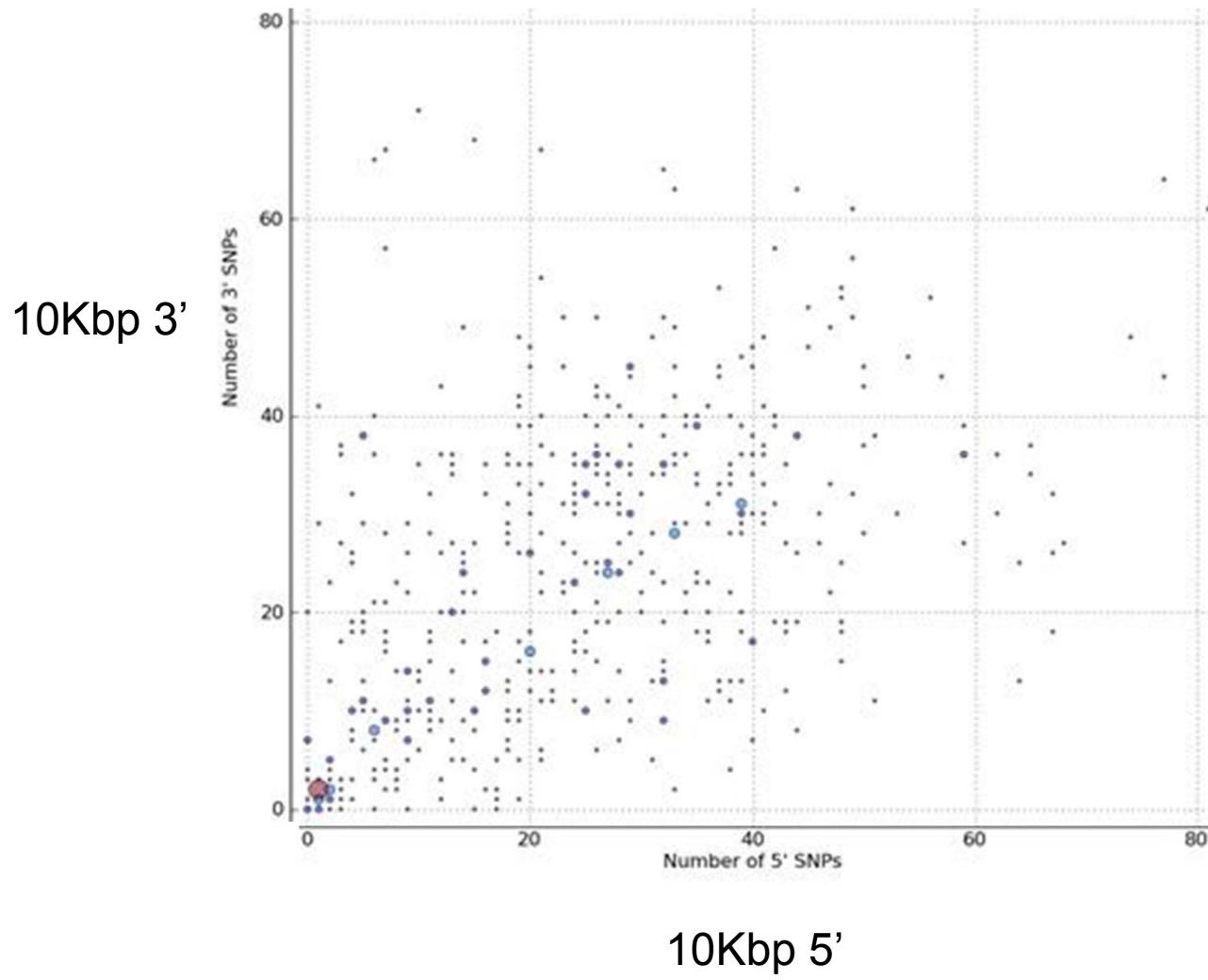
SNP density

- Causes of SNP density variation
 - Constraints on variation tolerance in expressed genes
 - Breeding and selection, reducing SNP density in regions and fixing alleles
 - Diverse crossing/introgression, increasing variation in regions

SNP density around genes



SNP density around genes



Which genes have low SNP density

Example genes:

Q6NLD5	Ethylene-responsive transcription factor ERF015
P49592	Protein Dr1 homolog
O49550	Dof zinc finger protein DOF4.5
O95780	zinc finger protein 682
Q9FJS2	Homeobox-leucine zipper protein HDG5
Q8GXT3	Transcription factor bHLH123
Q9SAH7	Probable WRKY transcription factor 40
Q9SZ69	Zinc finger A20 and AN1 domain-containing stress-associated protein 7
Q39081	Transcription factor CAULIFLOWER
Q7XJK6	Agamous-like MADS-box protein AGL36
Q9S7L5	Ethylene-responsive transcription factor ERF018
Q1PDN3	Heat stress transcription factor A-6a
Q9SJ41	Zinc finger CCCH domain-containing protein 18
Q9FX25	Auxin response factor 13
Q5RJC5	Zinc finger CCCH domain-containing protein 67
Q8L500	Two-component response regulator-like APRR9
Q8GZ13	Transcription factor BEE 1
Q8L7A4	Probable ADP-ribosylation factor GTPase-activating protein AGD11
Q9SIB4	WUSCHEL-related homeobox 3
Q9SGJ6	Dehydration-responsive element-binding protein 1E
Q38828	Auxin-responsive protein IAA10
Q9LV52	Heat stress transcription factor C-1
Q38827	Auxin-responsive protein IAA9
Q6J9Q2	Ethylene-responsive transcription factor ERF086

Gene classes:

1. transcription (7.33)
2. ion transport (3.42)
3. signal transduction / membrane receptors (2.87)
4. F-box protein (2.73)
5. transposable element (2.27)
6. response to auxin / stimuli (2.03)
7. chromatin assembly, DNA damage/repair (1.93)
8. auxin transporter(1.91)
9. developmental processes (1.81)
10. immune response (1.71)

SNP density

- Causes of SNP density variation
 - Constraints on variation tolerance in expressed genes
 - Breeding and selection, reducing SNP density in regions and fixing alleles
 - Diverse crossing/introgression, increasing variation in regions

Finding function

- Which alleles have been selected for in different germplasm?
- Have all favourable alleles been fixed?
- Have unfavourable alleles been dragged along for the ride?
 - What is the impact of linkage drag?
- Can this information be used for selecting parents/progeny?
- Can this be applied for breeding better canola?

Acknowledgements



Australian Government

Australian Research Council



THE UNIVERSITY
OF QUEENSLAND
AUSTRALIA



Grains
Research &
Development
Corporation



Paul Berkman
Kenneth Chan
Chris Duran
Michael Imelfort
Kaitao Lai
Hong Lee
Edmund Ling
Michal Lorenc
Sahana Manoli
Pradeep Ruperio
Jiri Stiller

Dominic Eales
Lars Smits

Jacqueline Batley
Alice Hayward
Emma Campbell
Jessica Dalton-Morgan
Reece Tollenaere

Harsh Raman

Bart Lambert
Benjamin Laga

Contact:
Dave.Edwards@uq.edu.au

