Workshop: Introduction to Pretzel



Pretzel

Enabling the seamless connection of plant genotypes to research and breeding knowledge

Duration: 2 hours

Format: Hands-on with guided exercises

Prerequisites:

- An interest in using digital tools for genomic data analysis
- Your own laptop

Learning Objectives:

- Understand what Pretzel does and when to use it
- Load a custom QTL dataset
- Load genotype data and explore allelic diversity
- Search for AGG accessions containing a specific haplotype
- Find a genomic region of interest using BLAST
- How to combine different types of data (genomes, genetic maps, QTLs and genes)







The Australian Grains Genebank Strategic Partnership is a \$30M joint investment between the Victorian State Government and Grains Research and Development Corporation (GRDC) that aims to unlock the genetic potential of plant genetic resources for the benefit of the Australian grain growers. The Partnership is genotypically characterising the AGG plant genetic resource collection and building novel digital tools to improve accessibility and utility of the genebank for industry.

Contents

Introduction to AGG Pretzel	
What is Pretzel?	4
Why use Pretzel?	5
Pretzel documentation	6
Workshop Overview	7
Workshop tasks	8
Starting out in <i>Pretzel</i>	8
Task 1: Log in to <i>Pretzel</i> and complete Data Licence Agreement	8
Pretzel layout	10
Entry points into <i>Pretzel</i>	11
Background information relevant to workshop tasks	12
Upload a QTL to Pretzel	13
Task 2: Prepare a custom QTL upload file	13
Task 3: Upload the Excel file to <i>Pretzel</i>	15
Task 4: Visualise the new QTL dataset	16
Compare positions of QTLs by comparative mapping	21
Task 5: Load a second QTL dataset	21
Task 6: Load anchored markers for comparative mapping	23
Axis Title menu	25
Compare position of a QTL with a reported gene	27
Task 7: Identify a gene of interest using feature search	27
Task 8: Identify the equivalent gene in a different genome based on comparative m	
Explore diversity in a QTL region contributing to plant height	
Task 9: Load genotype data	35
Task 10: Explore diversity in the region of interest	37
Task 11: Sort AGG accessions in the genotype table based on allele calls	40
Task 12: Filter AGG accessions with a defined haplotype of interest	41
Explore diversity around a locus containing a flowering time gene orthologue	46
Task 13: Perform a BLAST search	46

Task 14: Explore gene annotations	50
Visualising genotype calls for specific accessions	56
Task 15: Visualise genotype data in QTL region for specific AGG accessions	57
Additional tasks for a challenge	67
References	68
Useful URLs	69
Appendix	70
Pretzel data upload templates	70
Overview	70
Metadata sheet	71
Dataset sheet	73
Troubleshooting	75
Contact Details	76

Introduction to AGG Pretzel

What is Pretzel?

Pretzel (https://agg.plantinformatics.io/) is an open-source, interactive web application designed to enable users to connect plant genotype data to research and breeder knowledge. It allows users to visualise and interact with genetic and genomic data without needing specialised bioinformatics skills.

As genotype data is generated through the Australian Grains Genebank (AGG) Strategic Partnership, the genotype data and related datasets (such as genomes, BLAST databases and marker positions) are being loaded into *Pretzel*.

Pretzel currently contains genotype data and accessory files for wheat, barley, chickpea, lentil and field pea, with lupin being released shortly. Data for faba bean, mungbean, sorghum, oat and Brassica is planned for release in 2025-26.

Pretzel is a flexible platform built to enable users to analyse genetic and genomic data in many ways. Some of the most common uses include:

- Visualise and explore genetic and genomic datasets, including AGG genotype data
- Search for accessions with specific alleles or haplotypes, including trait-linked markers
- Compare marker order across genetic maps and genome assemblies
- Align QTLs, markers and genes across different genome assemblies
- Investigate pan-genomes
- Explore genes underlying QTLs
- Find alternatives for trait-linked markers
- Compare own private data with genotype data for accessions in the AGG

Why use Pretzel?

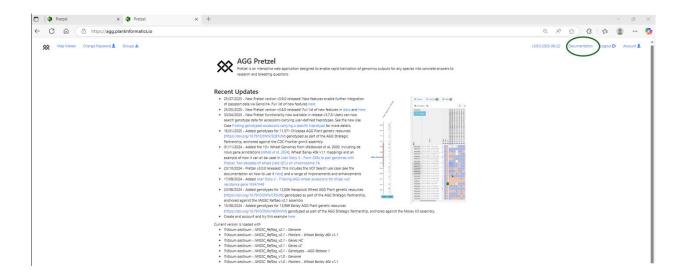
There are many digital tools available, each with their own advantages and disadvantages. *Pretzel* does not replace tools such as genome browsers, data repositories or BLAST servers, but is designed to complement these tools.

What makes *Pretzel* different:

- Highly interactive with broad utility
- Access and interact with genotype data for AGG accessions
- Upload custom, private datasets e.g. genetic maps, QTLs, genes, markers, genotype data
- Group permissions enable private data sharing with specific users
- Integrates different types of data, including legacy datasets
- Connects with external tools e.g. Germinate, Gigwa
- Built with industry via co-design and user feedback
- Online documentation and training available to support industry use

Pretzel documentation

Online documentation is available at https://docs.plantinformatics.io/.



Three types of documentation are available:

- Basic Functions
 - o Describes application layout and core features
- Use Cases
 - o Provides instructions to perform discrete *Pretzel* functions
- User Stories
 - Examples of how *Pretzel* can be used to answer complex research and breeding questions

Workshop Overview

Today's workshop will cover:

- Layout of the *Pretzel* GUI
- Uploading a QTL as a private dataset
- Comparing QTLs anchored to different genomes
- Exploring genotypic diversity in a QTL region
- Exploring annotated genes in a QTL interval
- Performing a feature search for a known gene
- Find similar genes across different genomes
- Finding chickpea accessions from the AGG that contain a haplotype of interest
- Comparing allele calls for specific chickpea AGG accessions
- Performing a BLAST search

Optional exercises if time permits include:

Exploring accession passport data with Genolink

Workshop tasks

Starting out in *Pretzel*

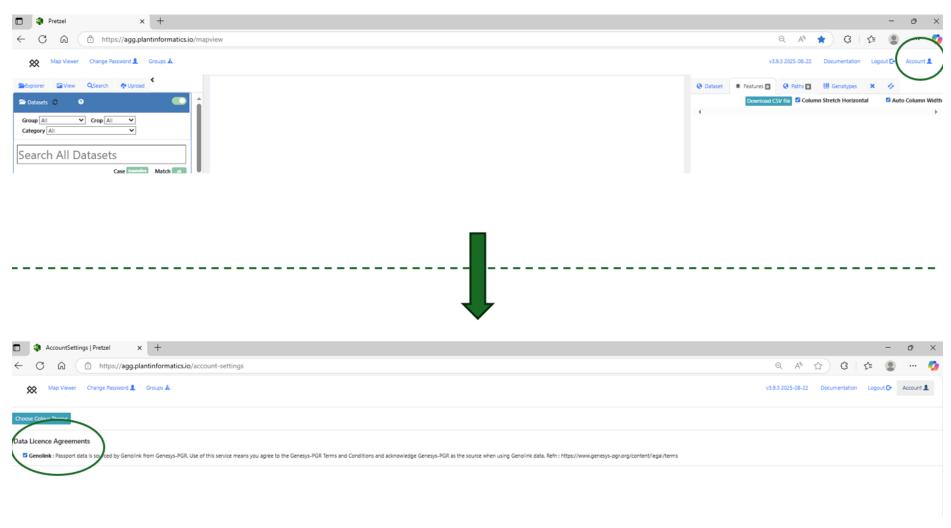
Task 1: Log in to *Pretzel* and complete Data Licence Agreement

If you haven't already requested a *Pretzel* account, click the 'Sign Up' link.

Otherwise, continue to log in with the 'Log In' link.

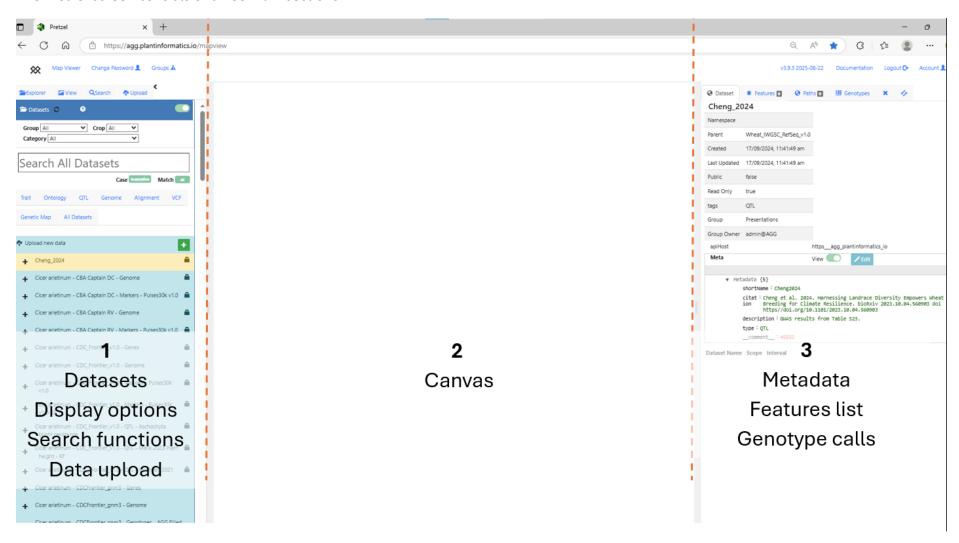


Click on 'Account' then tick the 'Genolink' box to agree to the Genesys-PGR Terms and Conditions .



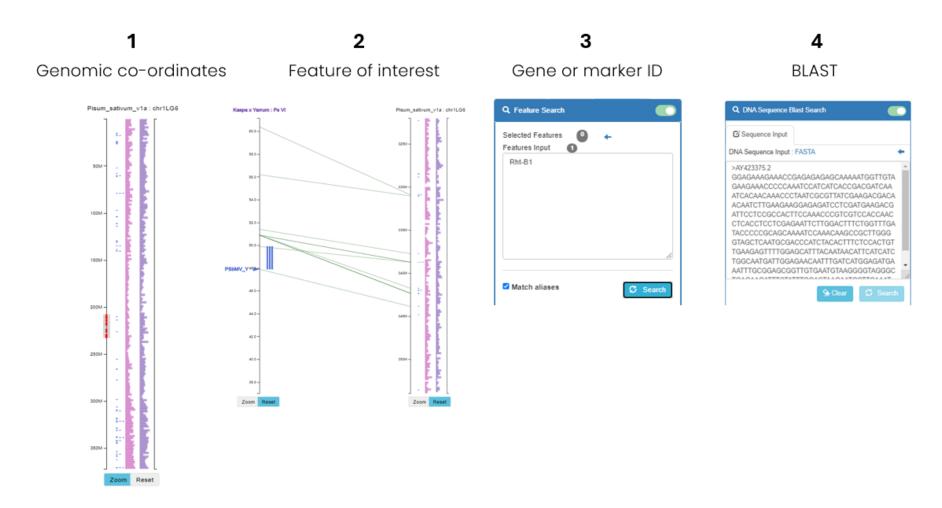
Pretzel layout

The *Pretzel* screen consists of three main sections:



Entry points into *Pretzel*

There are multiple ways for users to figure out which part of the genome to start interacting with.



Background information relevant to workshop tasks

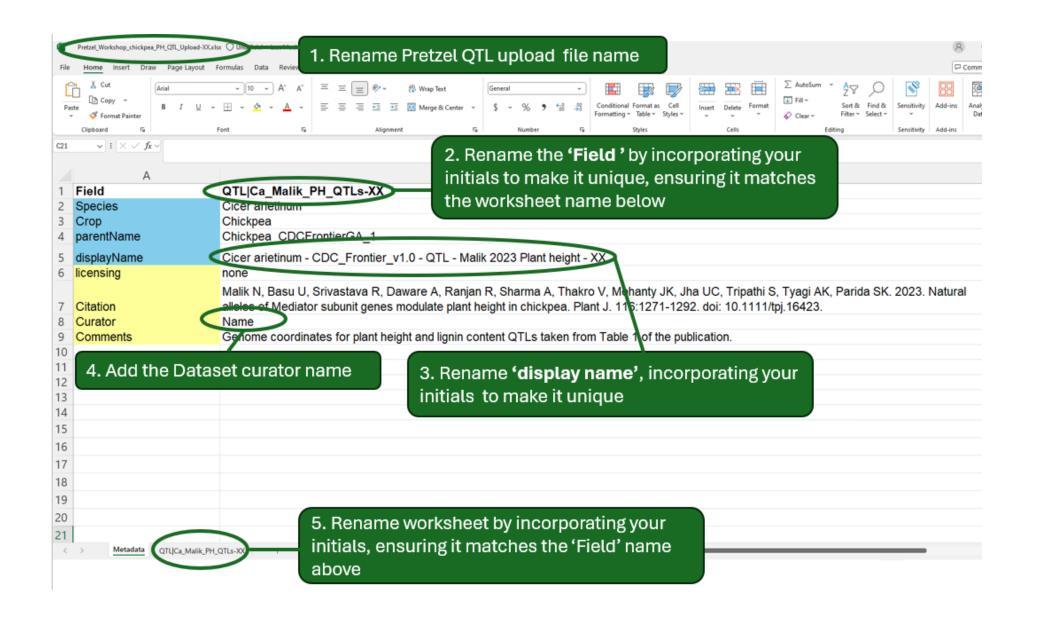
The exercises undertaken in this workshop are associated with genes involved in phenology and plant height in chickpea. Modern domesticated varieties have a narrow genetic base, which can limit their adaptability and resilience. Globally, there is increasing interest in mining existing diversity in crop wild relatives and diverse germplasm to help meet challenges in food security, sustainability and climate adaptation. However, breeders are often hesitant to utilise wild crops and diverse germplasm due to linkage drag and reduced agronomic performance, since it can be difficult to separate beneficial traits from undesirable ones. For example, genes associated with flowering time have implications for adaption to different climates and growing seasons, and genes associated with plant height can have a big impact on plant architecture and yield. Importantly, the presence of some wild alleles could result in flowering at disadvantageous times, or tall plants that can be impacted by lodging or reduced seed set. Having knowledge of which alleles are present is highly advantageous for selecting the right material for its desired purpose. The tasks in this workshop show users how to look at genetic diversity among AGG accessions around QTLs and candidate genes for flowering time and plant height and show how research knowledge can be used to identify chickpea accessions from the AGG that are more likely to contain wild introgressions. These examples demonstrate how users can make more informed selections of PGRs in research and breeding.

Upload a QTL to Pretzel

This task uses data published by Malik et. al. (2023) investigating diversity associated with plant height in chickpea. Information for QTLs on plant height and lignin content was extracted from the publication and entered into an Excel file in advance, with QTL positions reported against the CDC Frontier v1 (kabuli) genome assembly (Varshney et al, 2013). This pre-filled file has been provided to you. Today we will edit this file to ensure the name of the data is unique and then upload it into Pretzel for further investigation. After uploading the custom QTLs, the data will remain private, and you can easily locate the position of the QTLs whenever you log into *Pretzel*.

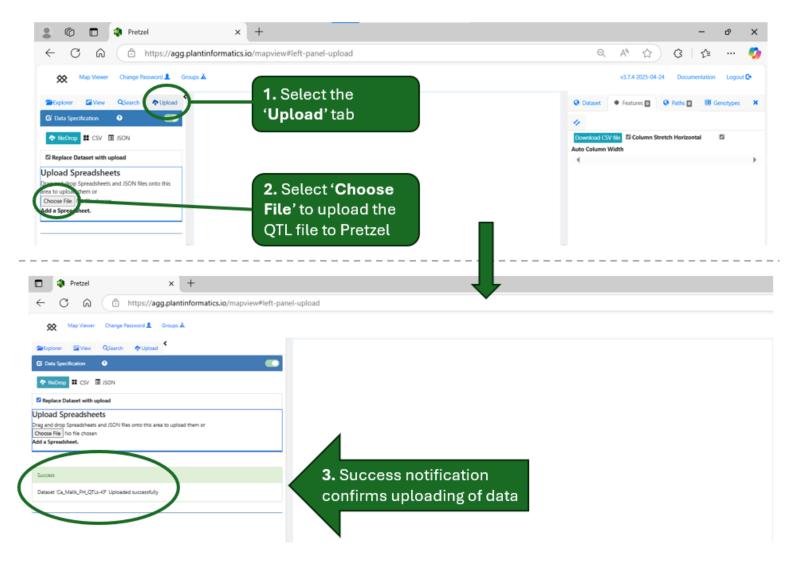
Task 2: Prepare a custom QTL upload file

- 1. Rename the *Pretzel* QTL upload file called '**Pretzel_Workshop_chickpea_PH_QTL_Upload-XX.xlsx**', replacing 'XX' with your initials e.g. Pretzel Workshop chickpea PH QTL Upload-KF.xlsx
- 2. Open the upload file and rename the sheet 'QTL|Ca_Malik_PH_QTLs-XX', replacing 'XX' with your initials e.g. QTL|Ca_Malik_PH_QTLs-KF.
- 3. Update the dataset name beside the 'Field' column in the 'Metadata' sheet.
- 4. Update the 'displayName' field 'Cicer arietinum CDC_Frontier_v1.0 QTL Malik 2023 Plant height XX', replacing 'XX' with your initials.
- 5. Add your name to the 'Curator' field.
- 6. Save the QTL upload file.



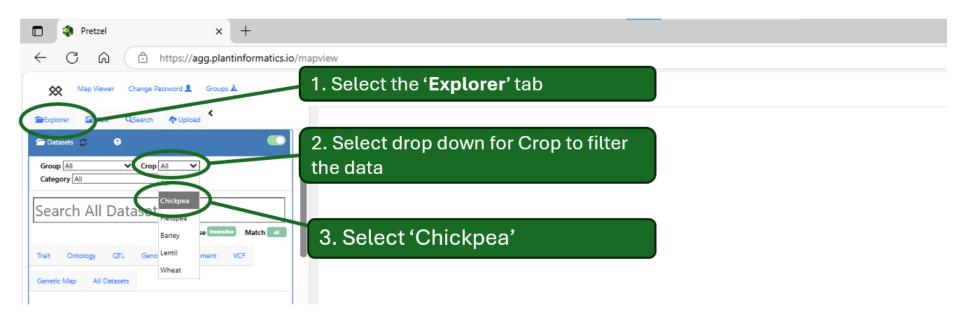
Task 3: Upload the Excel file to Pretzel

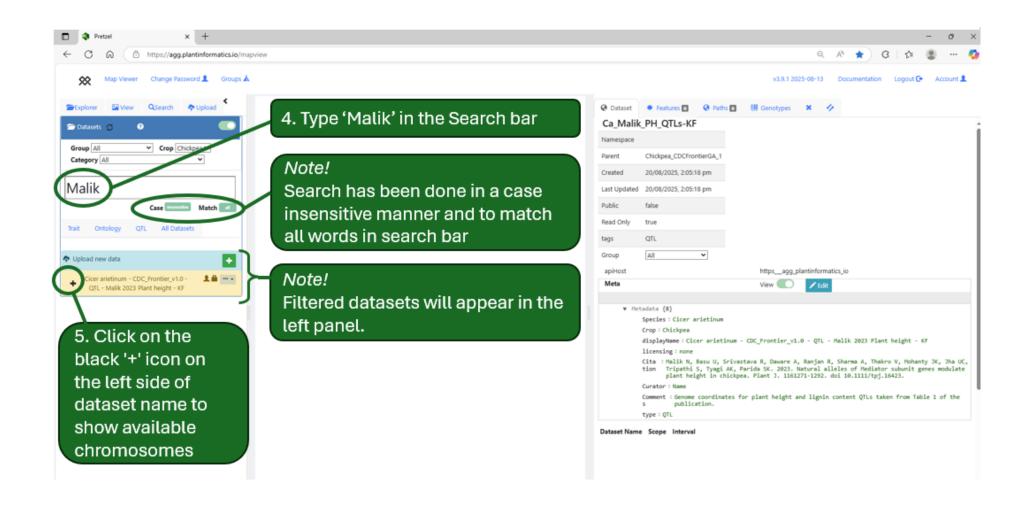
By following the steps below, the QTL will be uploaded as a private dataset.

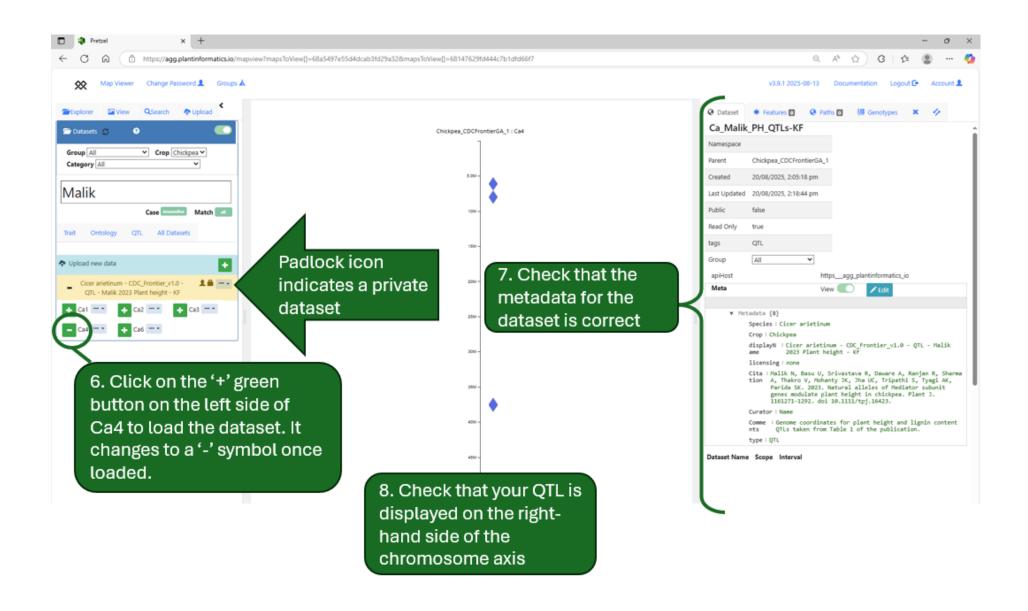


Task 4: Visualise the new QTL dataset

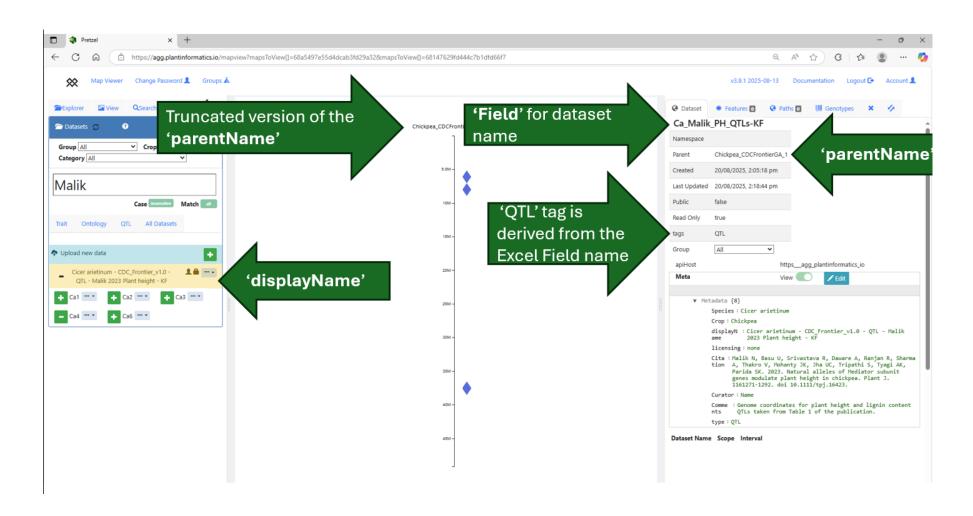
Next we will load the new custom dataset so we can see the location of the QTL on chromosome 4, then adjust the display to colour-code the QTLs according to the type of traits that were defined in the Excel file.



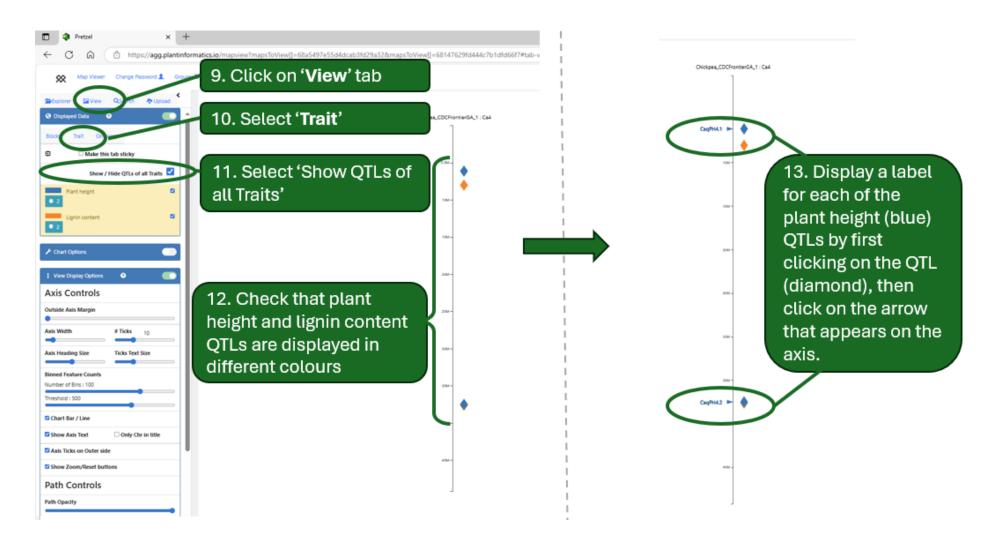




The screenshot below shows how different fields contained in the Excel QTL upload file are incorporated into the *Pretzel* Explorer panel, canvas (i.e. on the chromosome axis), and in the Dataset's metadata section. Refer to the Appendix for more information on the Excel template file for uploading custom QTLs.



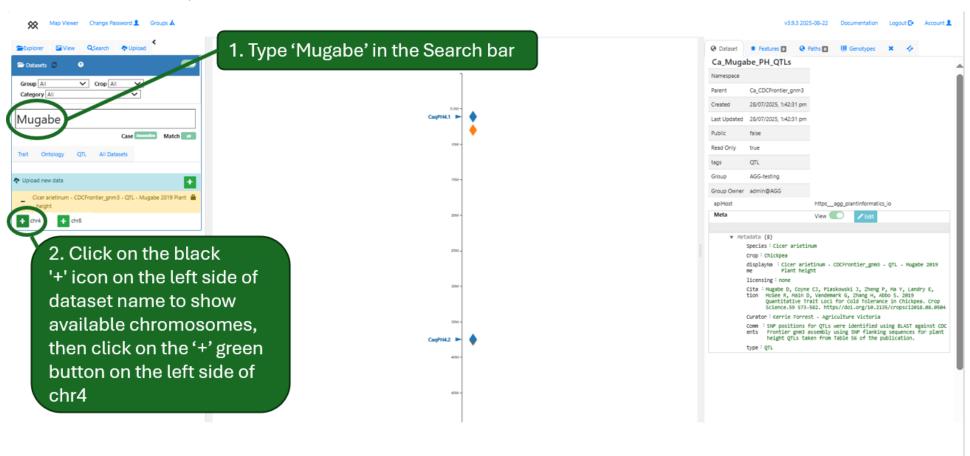
Next the QTL display will be adjusted to colour code the QTLs according to trait type.

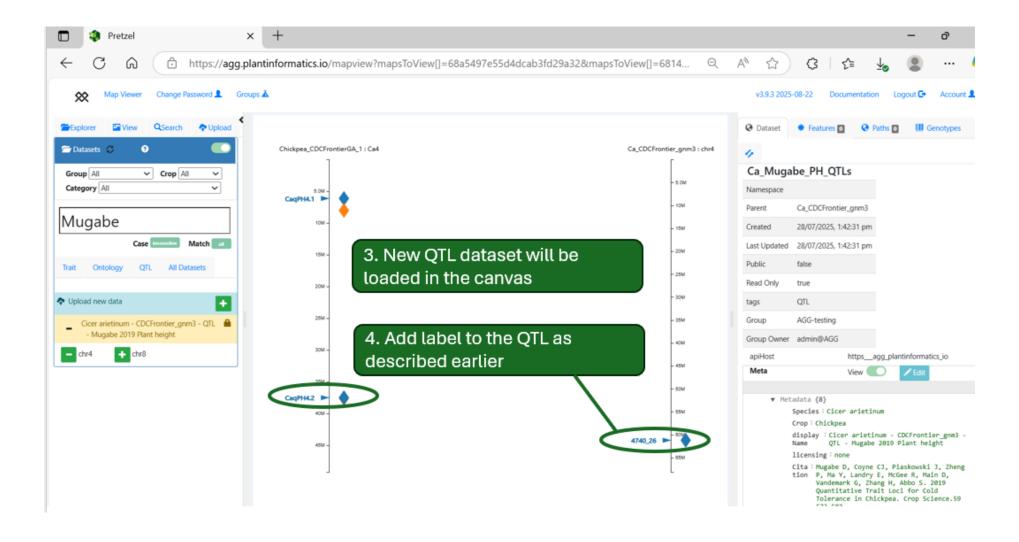


Compare positions of QTLs by comparative mapping

Next, we will load a QTL for plant height on chromosome 4 reported in an independent publication by Mugabe et al. (2019) to see whether any of the QTLs from Malik et al (2023) loaded above are detecting the same locus. Since QTLs from the two different publications have been anchored to different genome assemblies, this will be accomplished through a comparative mapping approach using the positions of Illumina Infinium™ Pulses 30K v1.0 BeadChip markers that have been anchored to both assemblies and loaded into *Pretzel*.

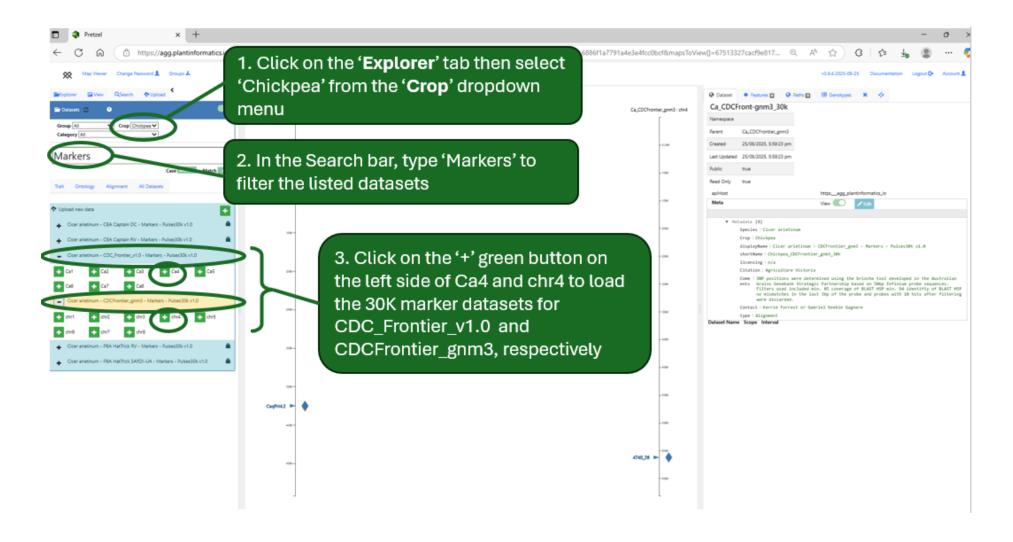
Task 5: Load a second QTL dataset



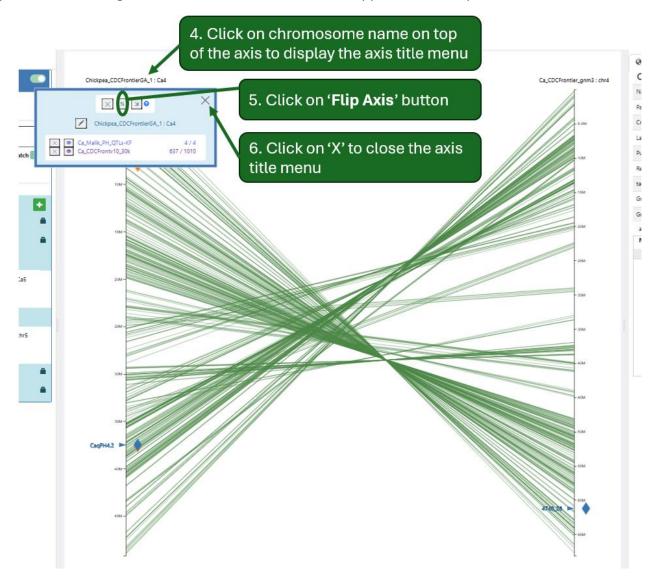


Based on the current display, the QTLs **CaqPH4.2** and **4740_26** towards the ends of the chromosomes might be detecting the same locus. We will investigate this further next.

Task 6: Load anchored markers for comparative mapping

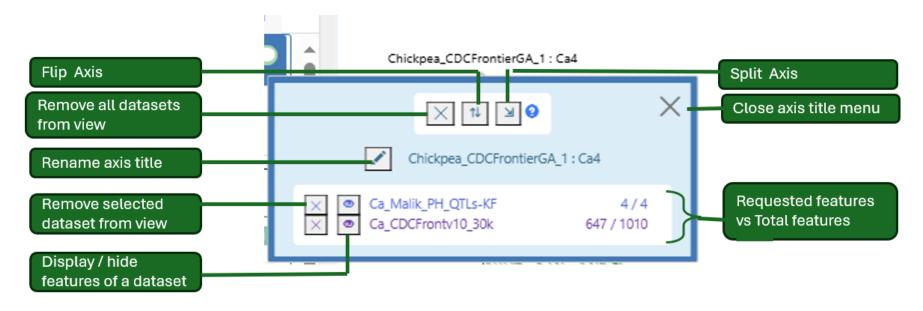


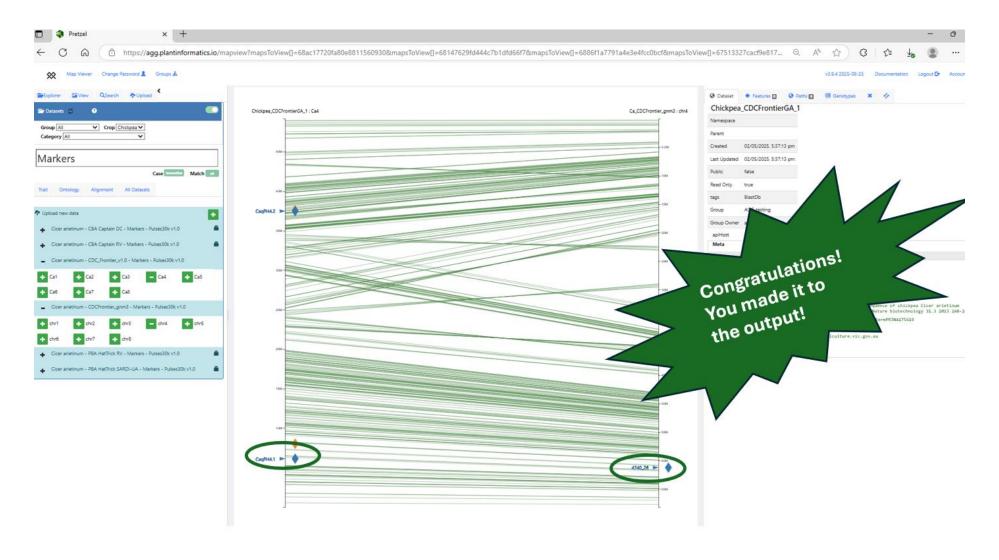
A green line drawn between the axes indicates that a marker with the same name is anchored to each of the two different versions of chromosome 4 displayed. In this case, marker ordering between the two versions of chromosome 4 is reversed, as indicated by the many green lines crossing over. One of the axes needs to be flipped. We will flip CDC Frontier v1.



Axis Title menu

The information below describes functionality that becomes available when the Axis Title menu is opened, and information about the loaded datasets.



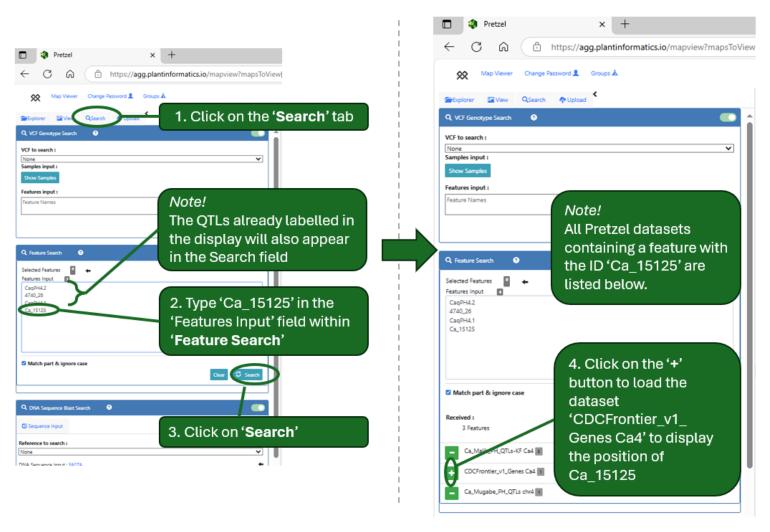


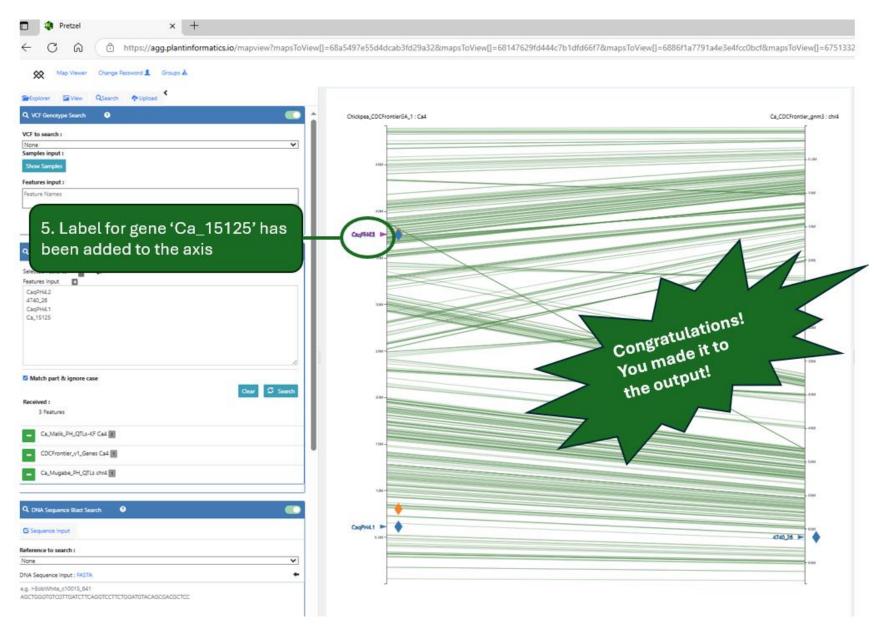
We have now successfully compared the positions of QTLs from two different publications, one set anchored to CDC Frontier v1 (Varshney et al. 2013) and the other anchored to an independent genome assembly CDC Frontier gnm3 (Cook, 2022). After taking into account the different orientations of the chromosome sequences, it is now clear that QTLs **CaqPH4.1** and **4740_26** are likely detecting the same locus, while **CaqPH4.2** is detecting a different locus.

Compare position of a QTL with a reported gene

In this example, we will search for a gene called CaMED23 reported to be involved in plant height and compare its position to the published QTL regions being investigated. The gene ID 'Ca_15125' in the CDC Frontier v1 genome annotation was reported in Malik et. al. (2023).

Task 7: Identify a gene of interest using feature search

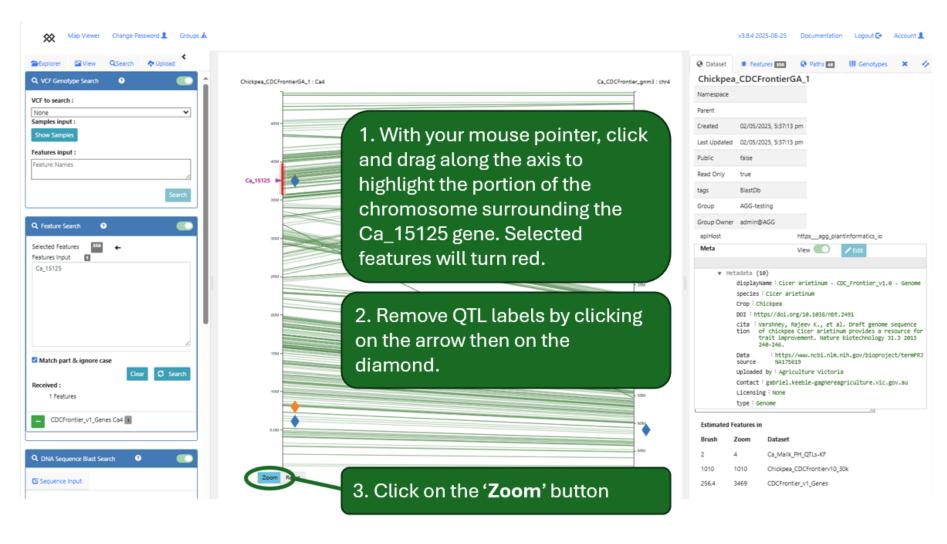


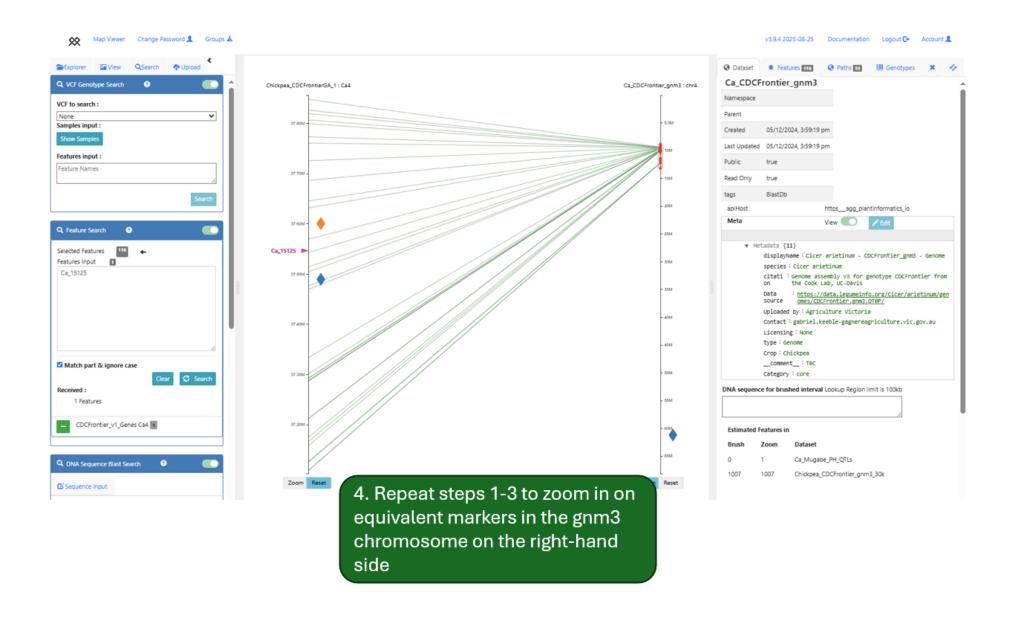


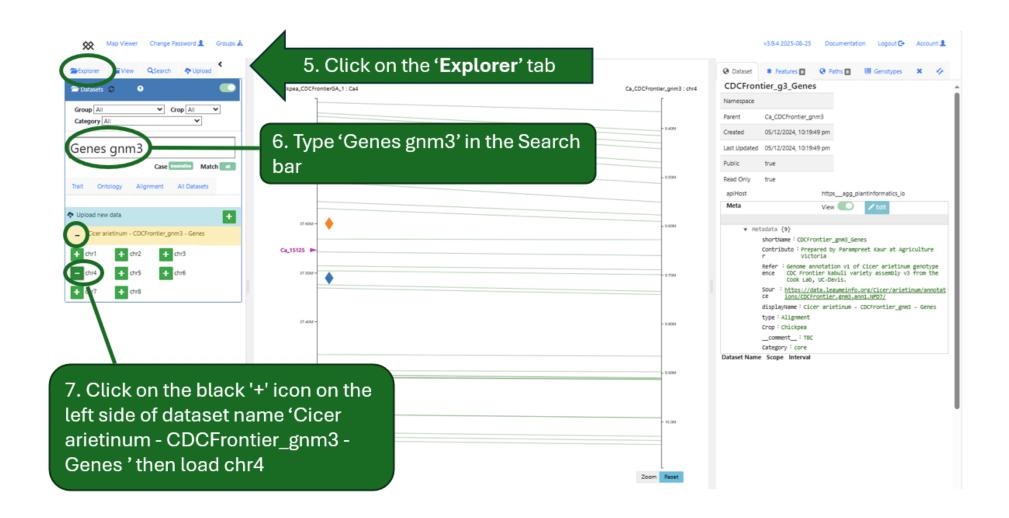
The position of gene Ca_15125 in chromosome Ca4 of the CDC Frontier v1 assembly is now indicated with a purple arrow and label. The gene is located very close to QTL CaqPH4.2 from Malik et. al. (2023).

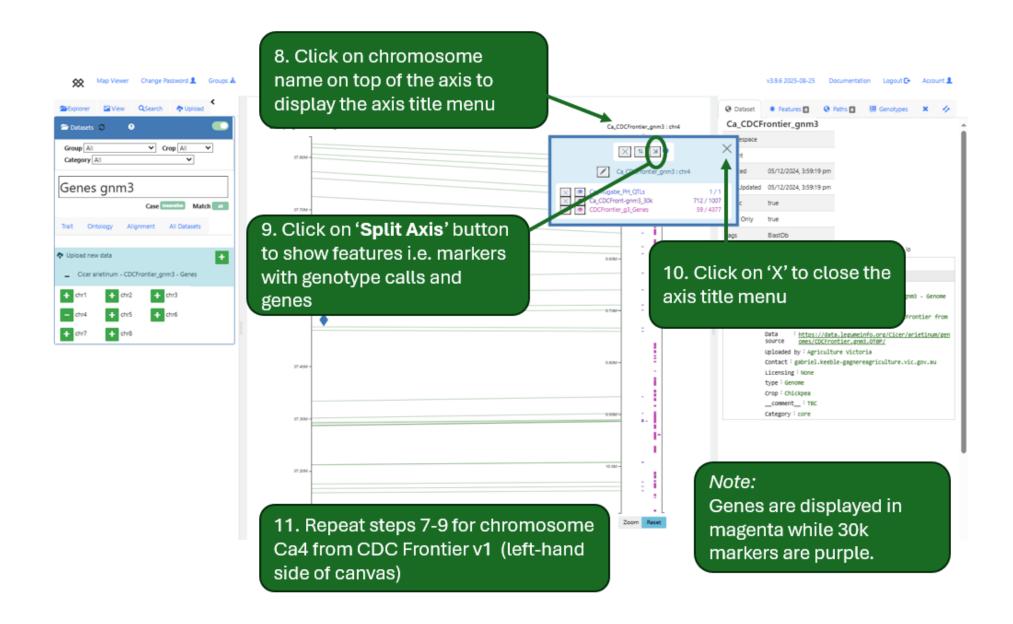
Task 8: Identify the equivalent gene in a different genome based on comparative mapping

Next we will identify the annotated gene for CaMED23 in a newer chickpea genome, CDC Frontier gnm3 (Cook, 2022). This will be accomplished through comparative mapping using Infinium™ Pulses 30K v1.0 BeadChip (Illumina) markers that have been anchored to both assemblies.

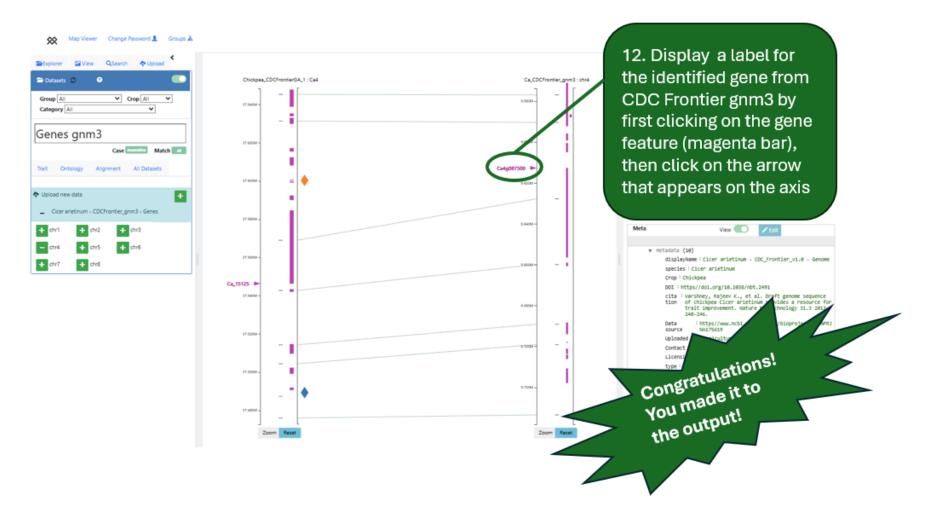






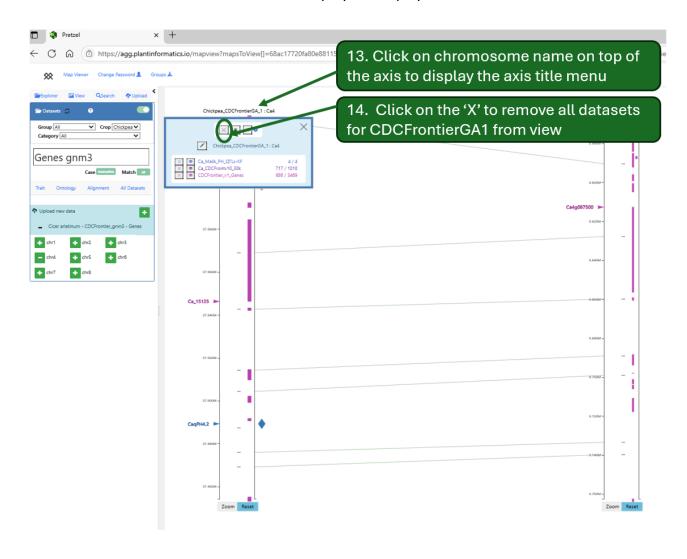


Continue to zoom in on both axes as required until you can see the start and stop position of Ca_15125 in CDC Frontier v1 and the equivalent gene in CDC Frontier gnm3. Next, we will add a label to the candidate gene in CDC Frontier gnm3 genome annotation to see the gene ID.



By comparing the relative positions of annotated genes in the two genome assemblies, we have successfully identified gene ID **Ca4g087500** in CDC Frontier gnm3 as the equivalent gene of **Ca15125** in v1.

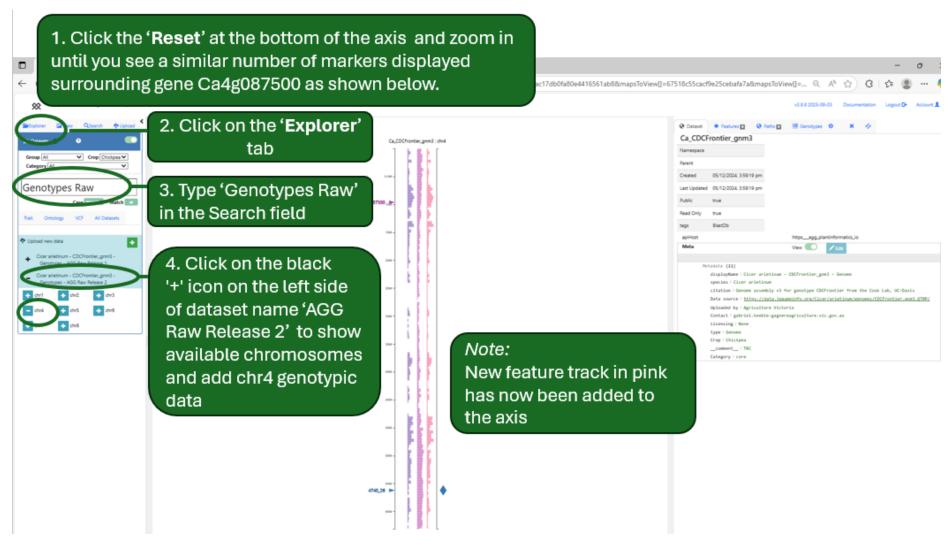
We will now close the Ca4 chromosome for CDC Frontier v1 to simplify the display for downstream tasks.

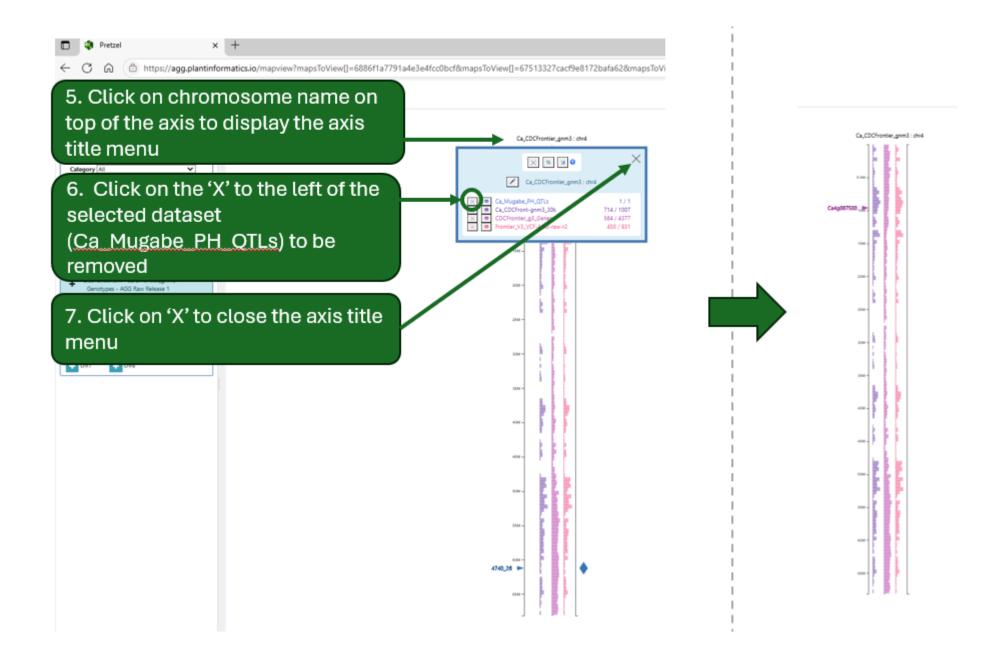


Explore diversity in a QTL region contributing to plant height

Next, we will load a vcf file containing genotypic data for chickpea accessions in the AGG and explore diversity of allele calls around the QTL CaqPH4.2 and annotated genes that contribute to plant height on chromosome 4.

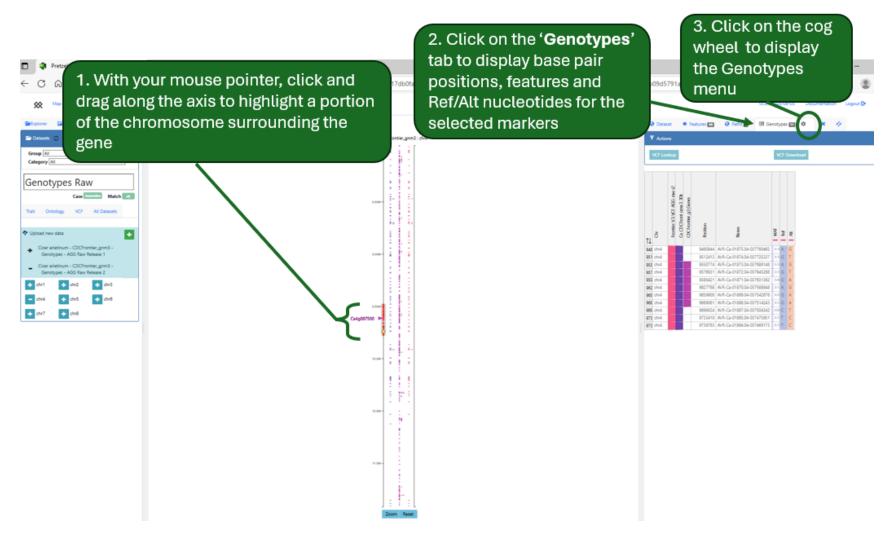
Task 9: Load genotype data





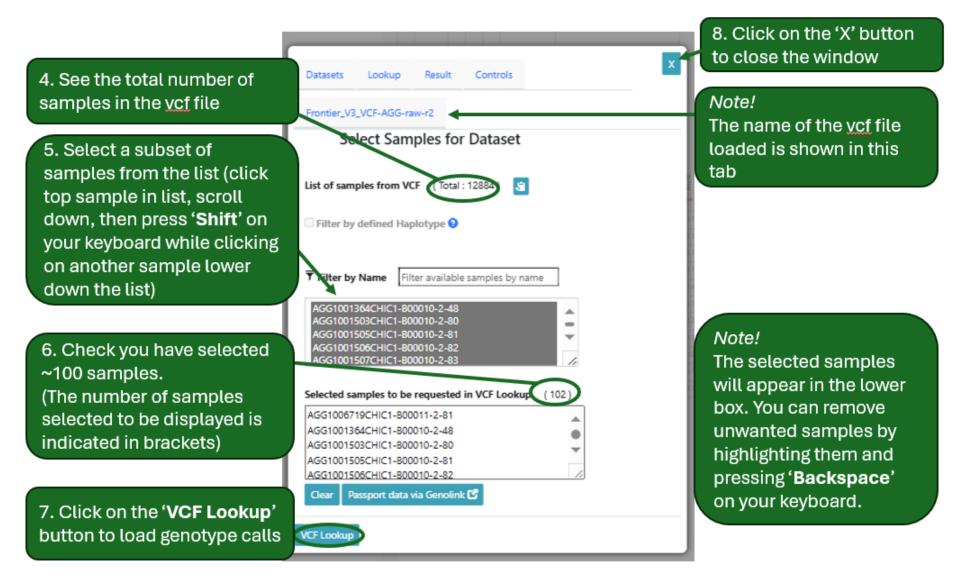
Task 10: Explore diversity in the region of interest

This task will involve zooming in and selecting markers that you wish to view genotype calls for surrounding the candidate gene.

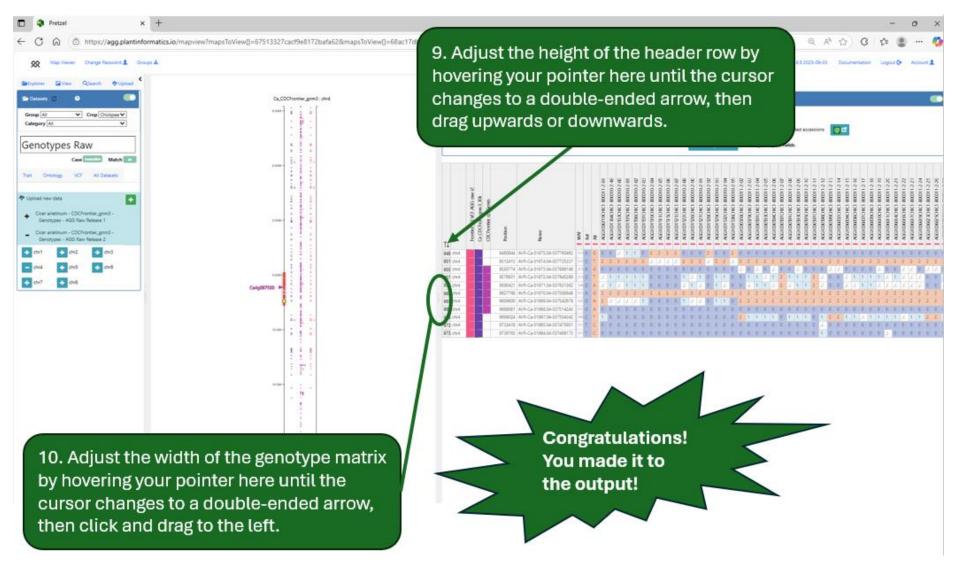


Markers that overlap annotated genes are indicated by the coloured (magenta) boxes in the 'CDCFrontier_g3_Genes' column. Note that while the nucleotides in the 'Ref' and 'Alt' columns are defined, the actual genotype calls are not yet loaded.

After displaying the Genotypes menu, you should select only a subset of samples to display e.g. 100. Attempting to load thousands of samples will result in slow performance.

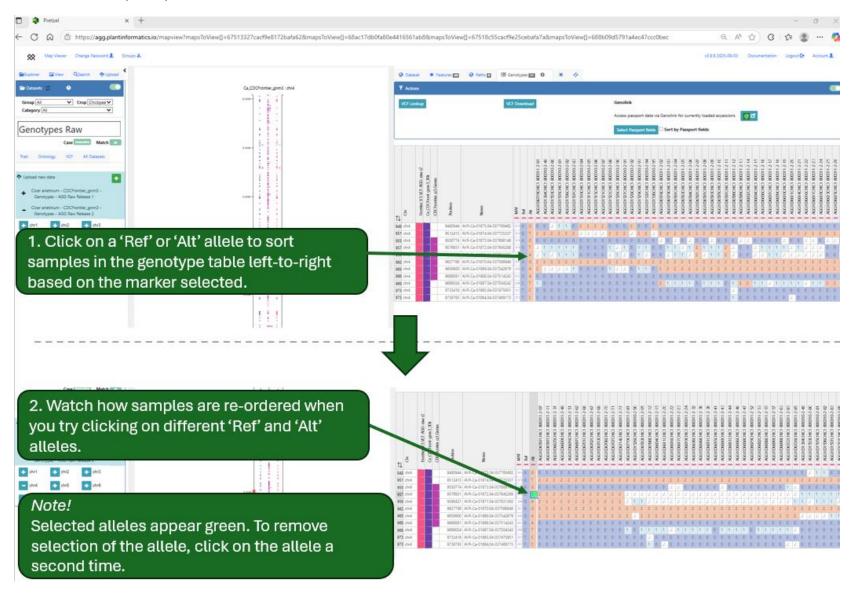


Genotypes will now be displayed in the 'Genotypes' table, similar to the example below. Markers are listed from top to bottom based on order along the chromosome. Samples are ordered left to right. In this case, you can see samples are labelled with AGG sample identifiers.



Task 11: Sort AGG accessions in the genotype table based on allele calls

In this task, you will learn how to sort accessions from left to right in the 'Genotypes' table, based on the allele you select for a chosen marker. You can sort samples by the 'Ref' or 'Alt' columns.

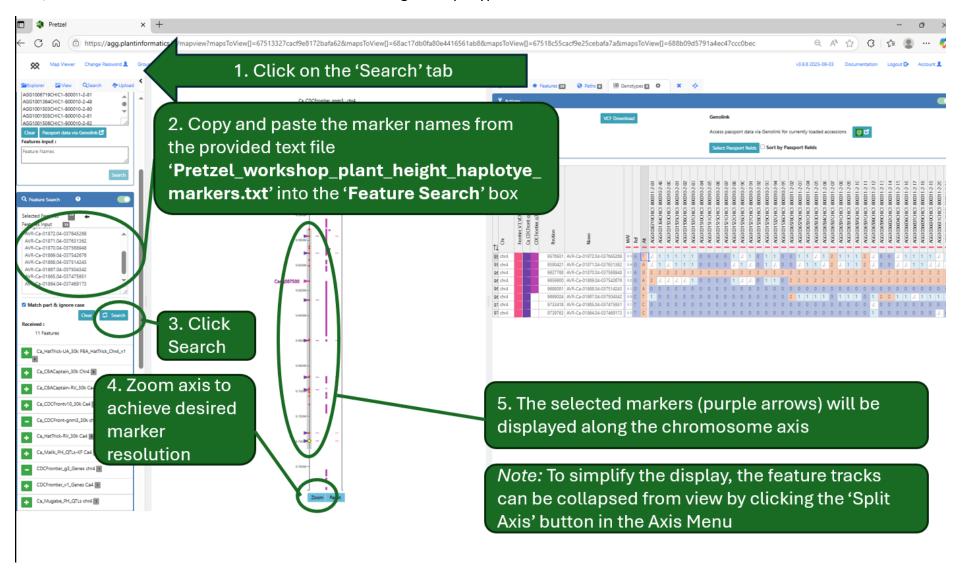


Task 12: Filter AGG accessions with a defined haplotype of interest

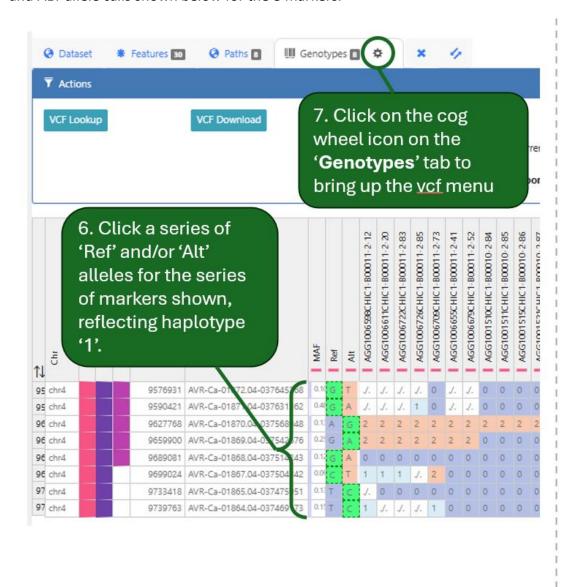
Using data from a research project, haplotypic diversity in this region containing a locus for plant height was explored in wild chickpea (*Cicer reticulatum* and *Cicer echinospermum*) as well as cultivated chickpea (*Cicer arietinum*). It was discovered that some haplotypes (i.e. series of genotype calls) were common among wild chickpea accessions while other haplotypes were common among cultivated chickpea accessions. Some of the most frequently observed haplotypes are shown in the table below. Our hypothesis is that we can identify accessions from the AGG with wild chickpea introgressions around the plant height gene by filtering accessions to display based on specific haplotypes.

				Haplotypes predominant in wild			Haplotypes predominant in cultivated					
Marker Name	Position	REF	ALT	1	2	3	4	5	6	7	8	9
AVR-Ca-01872.04-037645268	9576931	G	Т	REF	REF	REF	REF	REF	REF	REF	REF	REF
AVR-Ca-01871.04-037631362	9590421	G	Α	REF	REF	REF	REF	REF	ALT	REF	ALT	./.
AVR-Ca-01870.04-037568948	9627768	Α	G	ALT	ALT	ALT	ALT	ALT	ALT	REF	ALT	ALT
AVR-Ca-01869.04-037542676	9659900	G	Α	ALT	ALT	ALT	ALT	ALT	ALT	REF	ALT	./.
AVR-Ca-01868.04-037514243	9689081	G	Α	REF	REF	REF	REF	REF	REF	REF	ALT	REF
AVR-Ca-01867.04-037504342	9699024	С	Т	REF	REF	REF	ALT	REF	REF	REF	REF	REF
AVR-Ca-01865.04-037475951	9733418	Т	С	ALT	REF	REF	REF	./.	REF	REF	REF	REF
AVR-Ca-01864.04-037469173	9739763	Т	С	ALT	REF	REF	REF	ALT	REF	REF	REF	REF

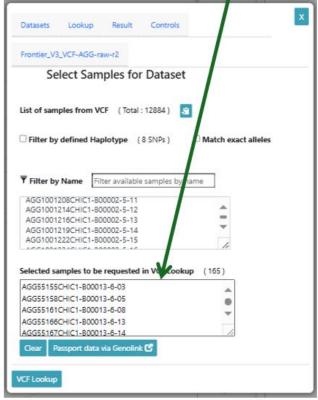
First, we will do a Feature search of the markers IDs defining the haplotypes above to make them easier to find in the chromosome.



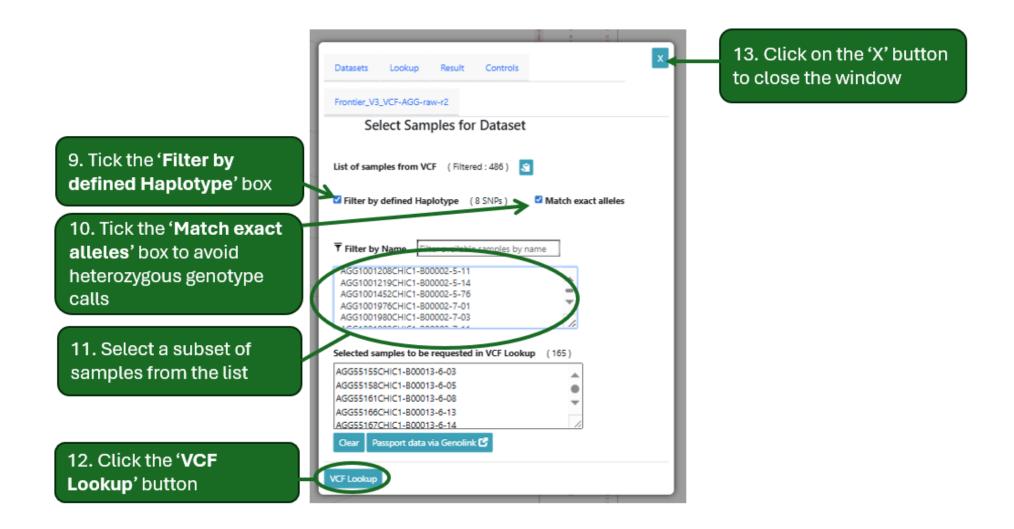
In this case, we will seek to identify chickpea samples in the AGG that have the wild haplotype '1' by selecting the specified combination of REF and ALT allele calls shown below for the 8 markers.

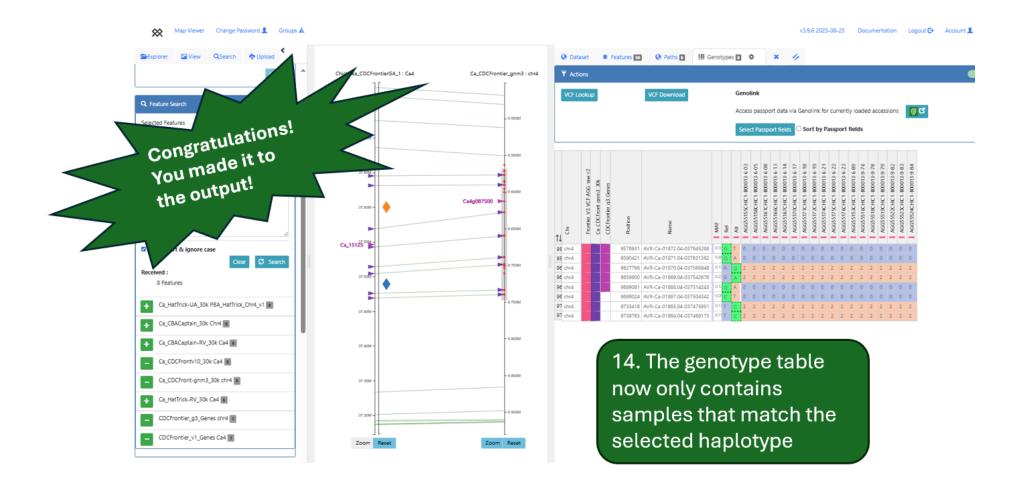


8. If samples are still listed in the lower box from the previous VCF Lookup, click into the box, select all samples with 'Ctrl' + 'A', then clear them with the 'Backspace' button on the keyboard.



IMPORTANT: select only a subset of the filtered samples. Loading thousands of samples will result in slow *Pretzel* performance.





We have now successfully identified chickpea accessions in the AGG with the wild haplotype '1'.

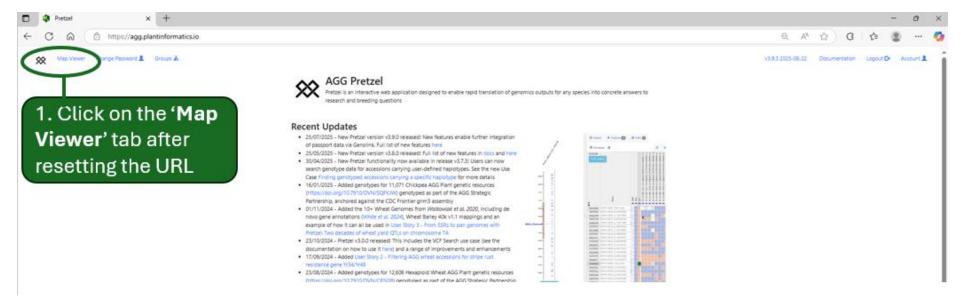
Explore diversity around a locus containing a flowering time gene orthologue

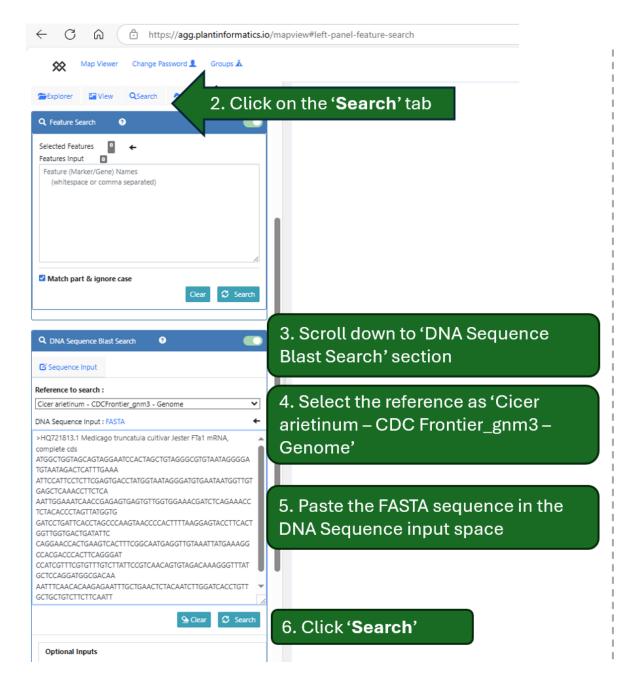
The next task will demonstrate how you can use BLAST to find a region of interest in the genome to investigate in *Pretzel*. Here, we will aim to find a gene in chickpea that is orthologous to Medicago gene MtFTa1 that was reported by Laurie et. al. (2011) to be involved in flowering time.

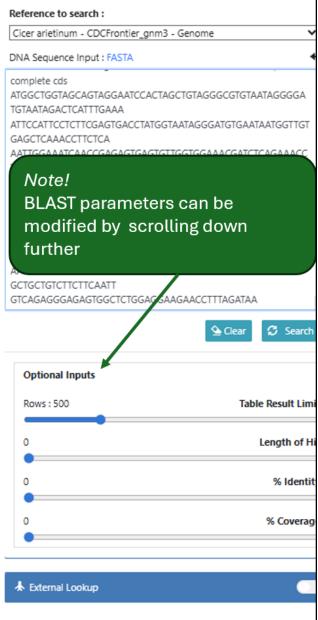
Task 13: Perform a BLAST search

Start by restarting your *Pretzel* session (e.g. re-enter the URL https://agg.plantinformatics.io/), since this is the quickest way to remove loaded datasets from view. Then, perform BLAST Search using the nucleotide sequence in the provided file 'Flowering_time_HQ721813.txt'.

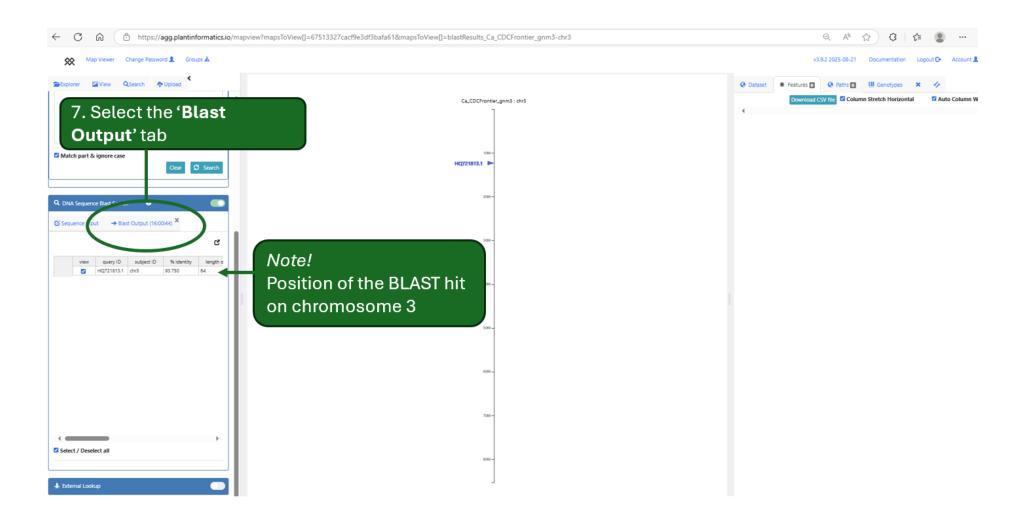
Alternatively, you can copy the sequence from NCBI with the link https://www.ncbi.nlm.nih.gov/nuccore/HQ721813.1?report=fasta.

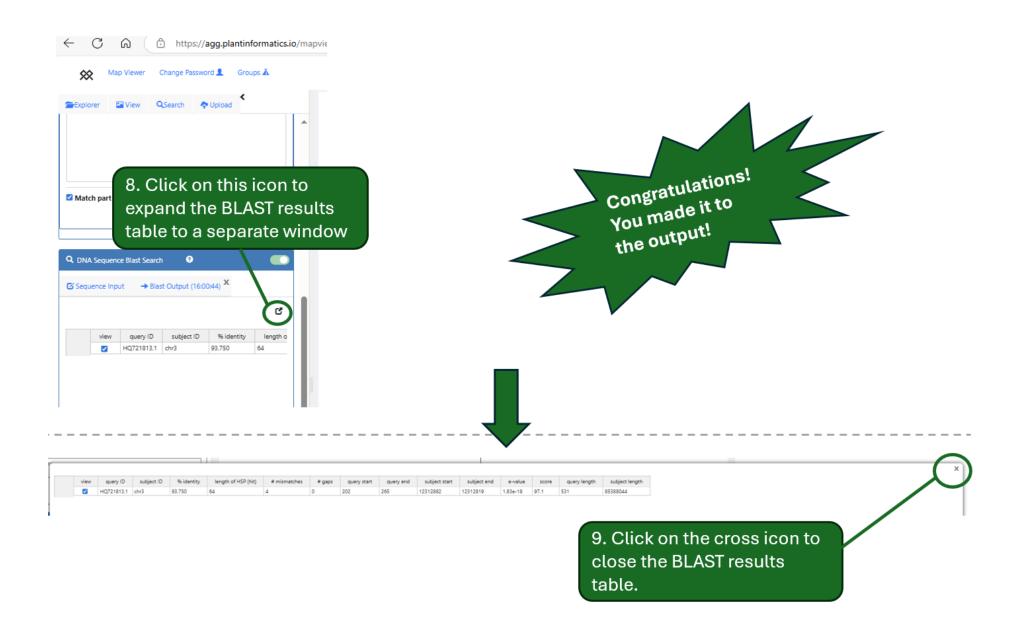






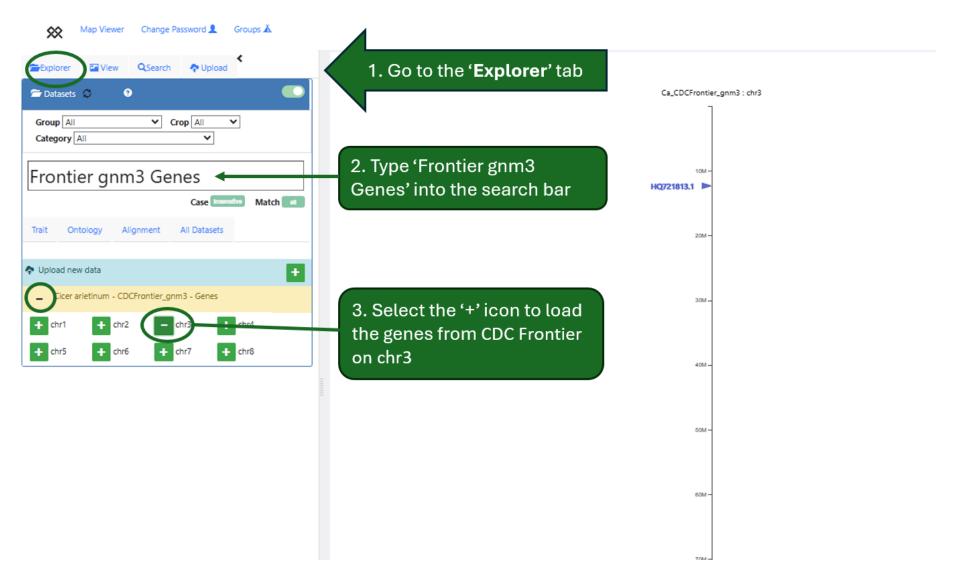
When you click on the BLAST output sheet, each chromosome with a BLAST hit is automatically loaded, with the position of the BLAST hits marked with an arrow and label. Note that the label is taken from the FASTA sequence identifier.

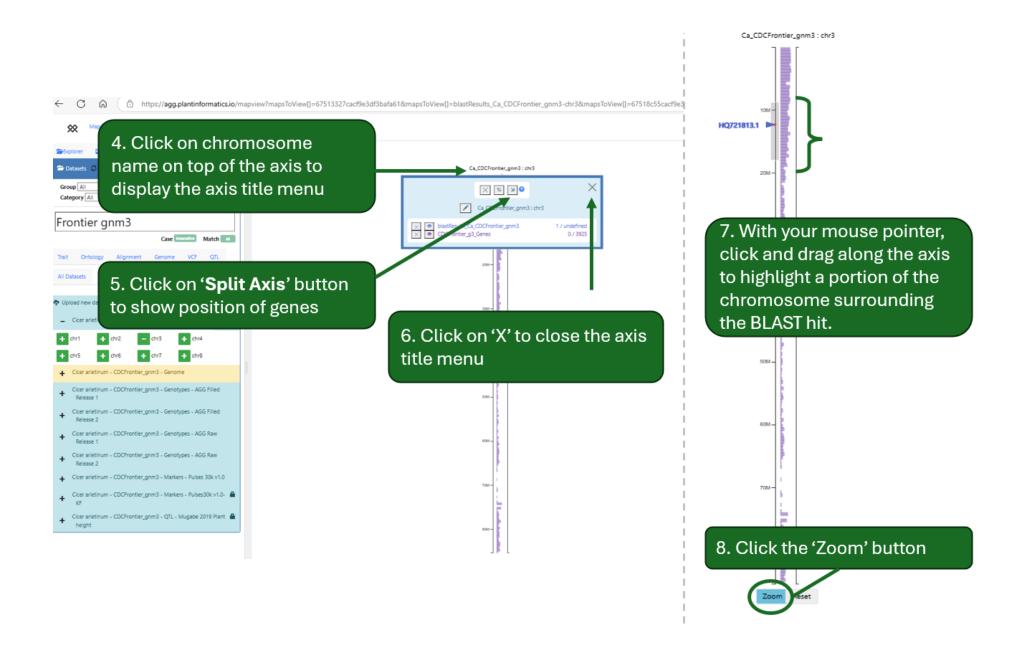




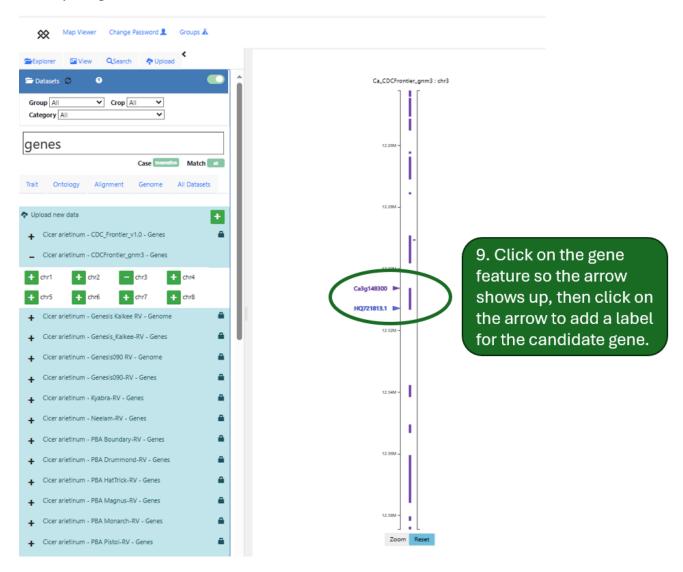
Task 14: Explore gene annotations

In this example, we will find a chickpea gene that corresponds to the Medicago MtFTa1 gene, then review chickpea gene annotations in the Feature Table to find a candidate gene contributing to flowering time in chickpea.

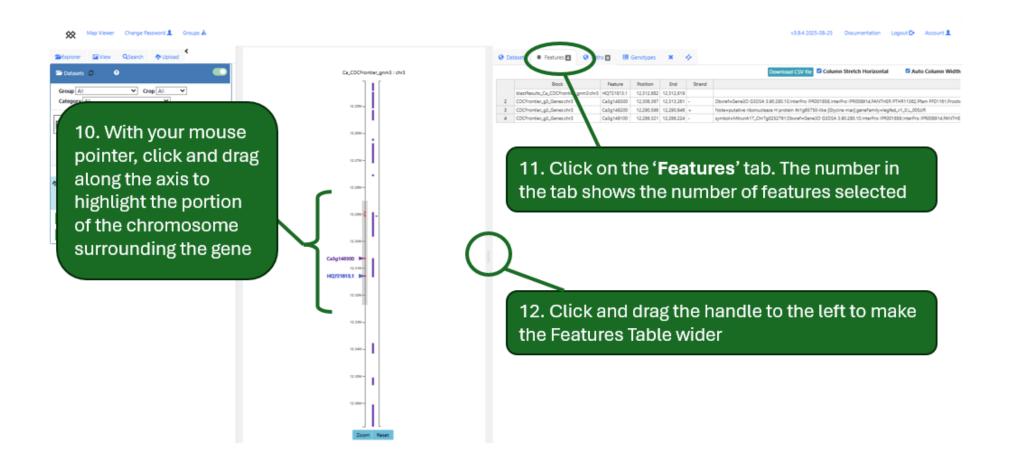


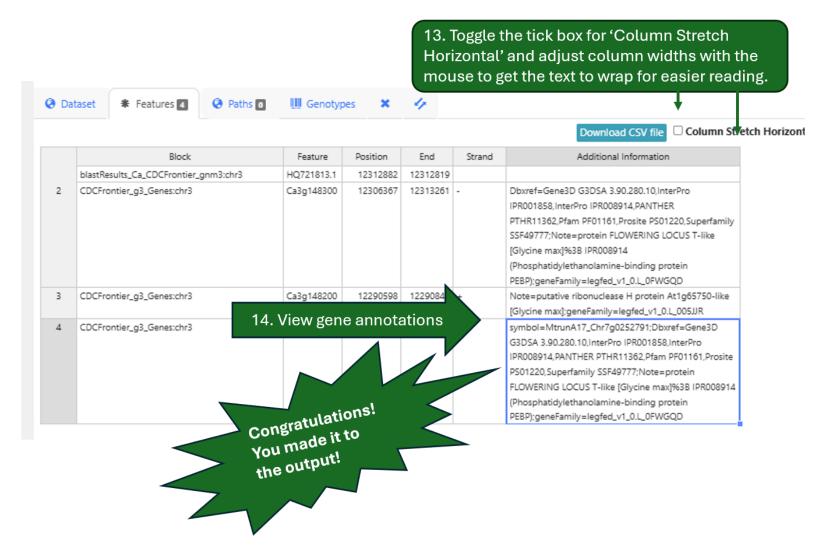


Repeat the zoom step above until you can see the position of the BLAST hit with the closest annotated gene, as shown below. Then, we can identify the gene ID with a label.



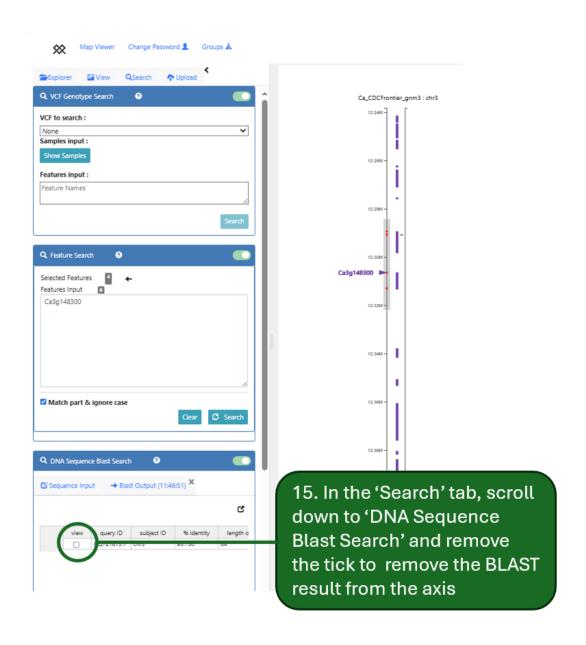
You can now see that the BLAST hit for the Medicago MtFTa1 gene sequence (**HQ721813**) matches the position of chickpea gene **Ca3g148300** from the CDC Frontier gnm3 assembly.





We can see the Prosite information describes gene Ca3g148100 as 'FLOWERING LOCUS T-like', suggesting we have correctly identified the chickpea gene in the CDC Frontier gnm3 gene annotation file that is orthologous to Medicago MtFTa1.

We can now close the BLAST result from view.

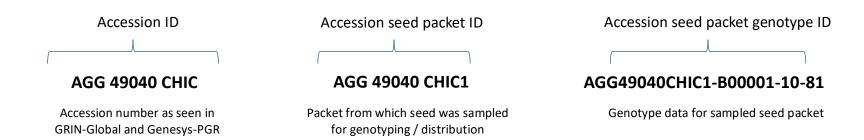


Visualising genotype calls for specific accessions

Using data from a research project, haplotypic diversity surrounding a gene for flowering time was explored in wild chickpea (*Cicer reticulatum* and *Cicer echinospermum*) as well as cultivated chickpea (*Cicer arietinum*). It was discovered that some haplotypes (i.e. series of genotype calls) were common among wild chickpea accessions while other haplotypes were common among cultivated chickpea accessions. Some of the most frequently observed haplotypes are shown in a table below in Task 15. Our hypothesis is that we can differentiate accessions in the AGG that carry wild introgressions versus cultivated alleles around the flowering time gene using haplotype information.

In this exercise, we will search for some known varieties of chickpea and some hybrids (wild introgression lines) to assess whether they contain haplotypes typically seen in wild or cultivated chickpea around a flowering time gene. AGG sample identifiers are provided in the electronic file 'Pretzel_workshop_Task15_AGG_chickpea_sample_list.txt'.

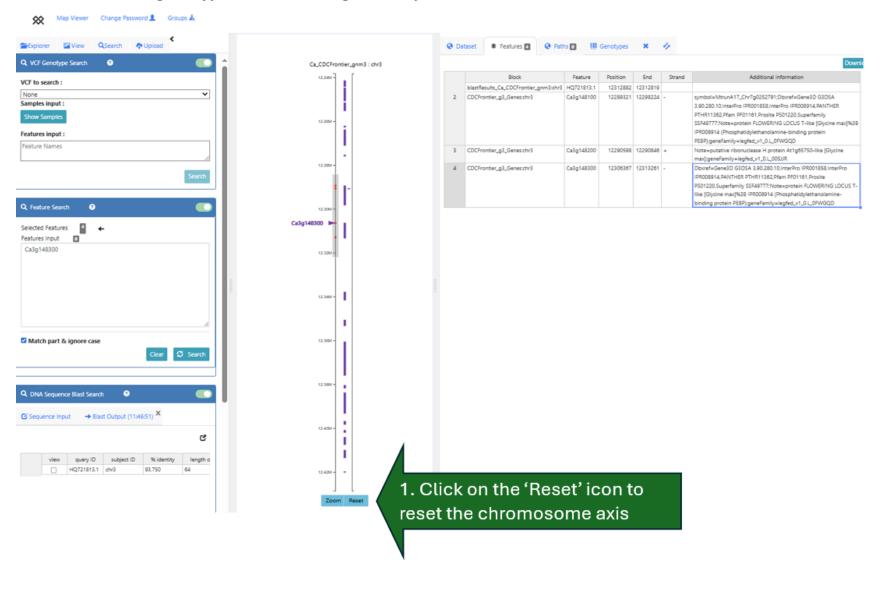
The Genotype ID provided uniquely identifies the DNA sample that was genotyped in the AGG Strategic Partnership for that accession.

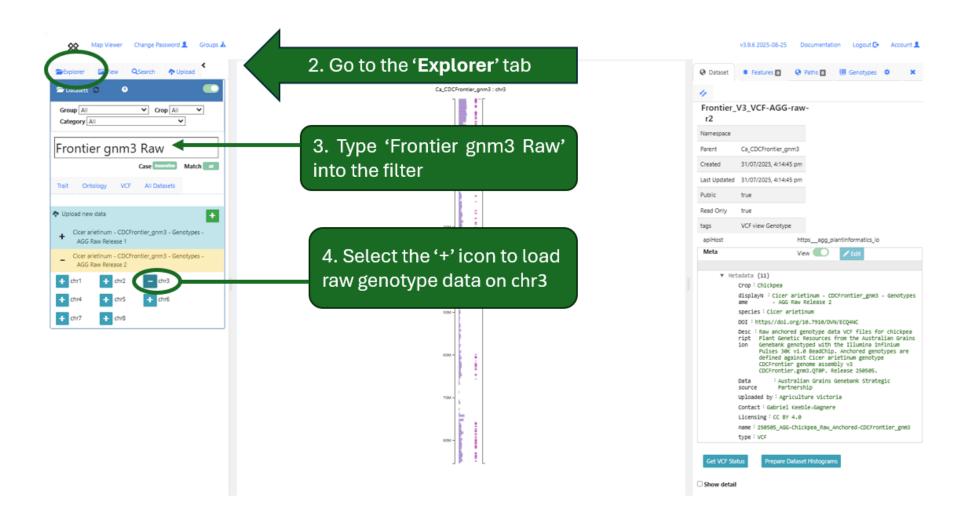


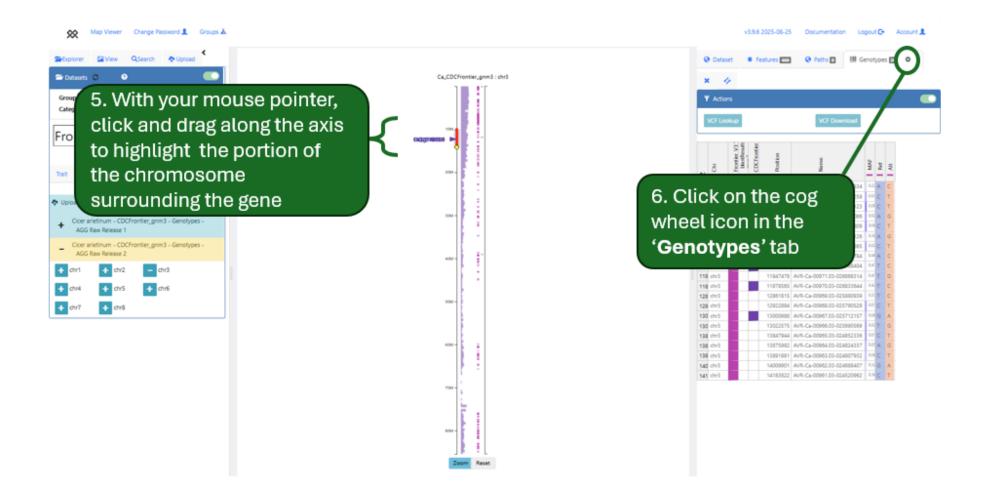
Accession IDs can be found from several sources:

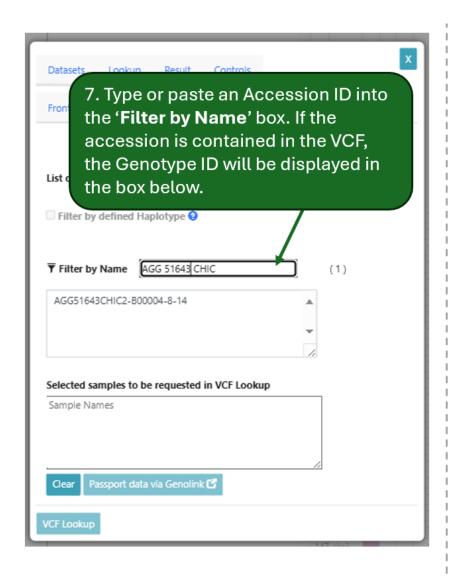
- AGG GRIN-Global database (https://ausgenebank.agriculture.vic.gov.au/gringlobal/search)
- Genolink (<u>https://genolink.plantinformatics.io/</u>)
- Genesys-PGR (<u>https://www.genesys-pgr.org/</u>)
- Publications

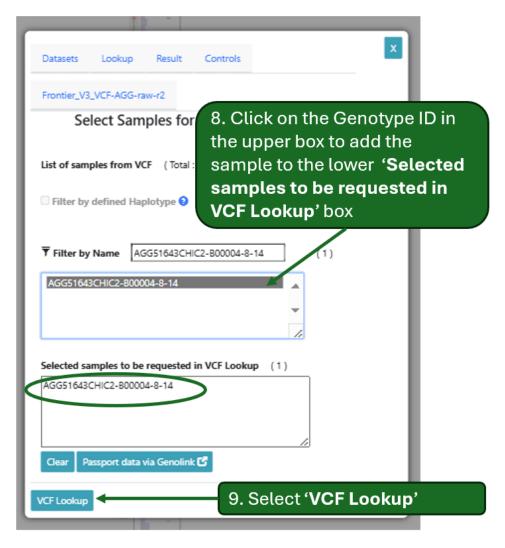
Task 15: Visualise genotype data in QTL region for specific AGG accessions

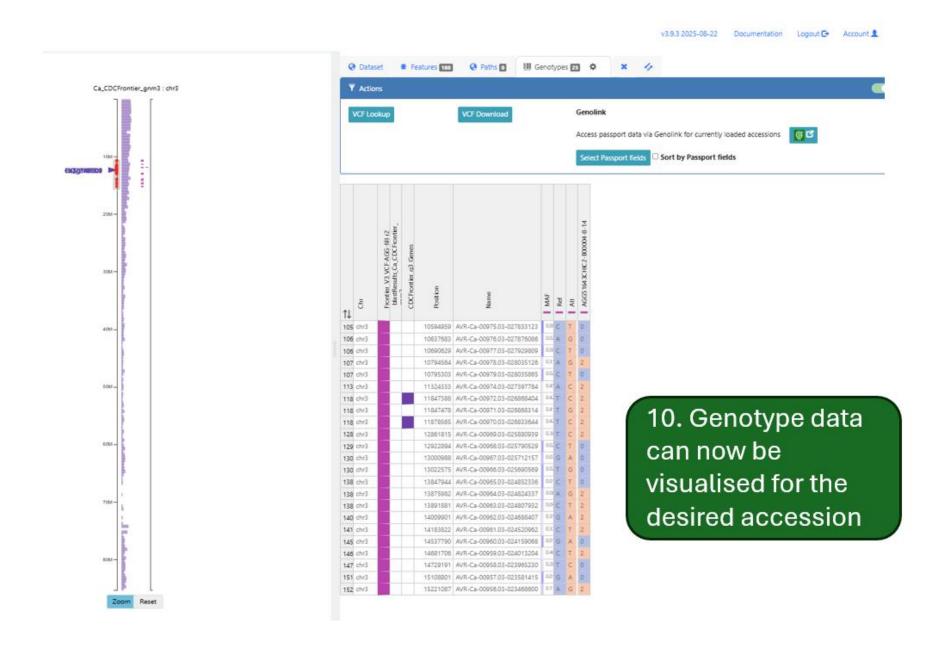




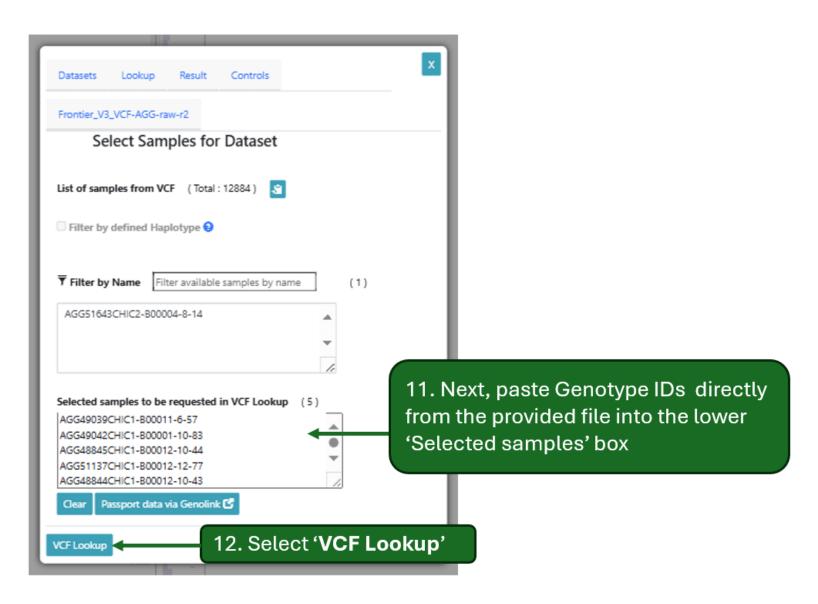




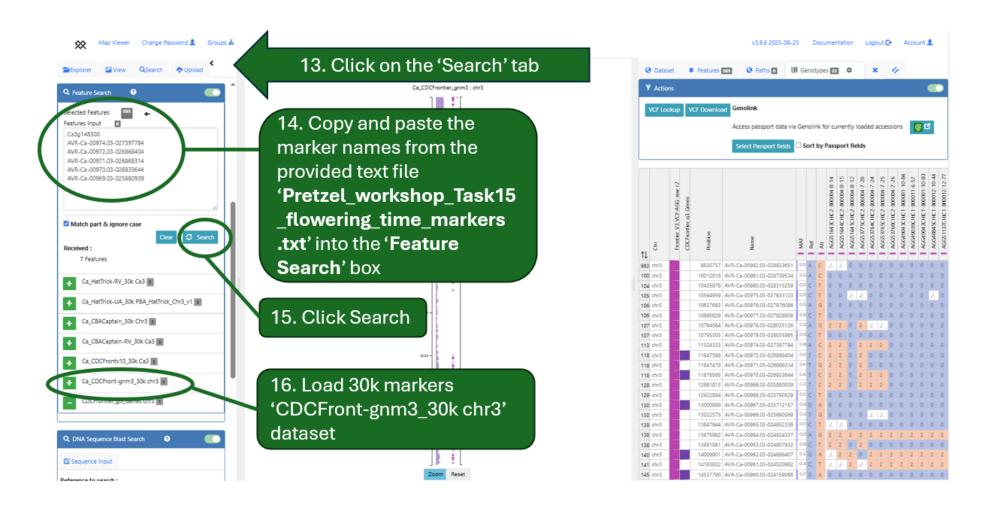


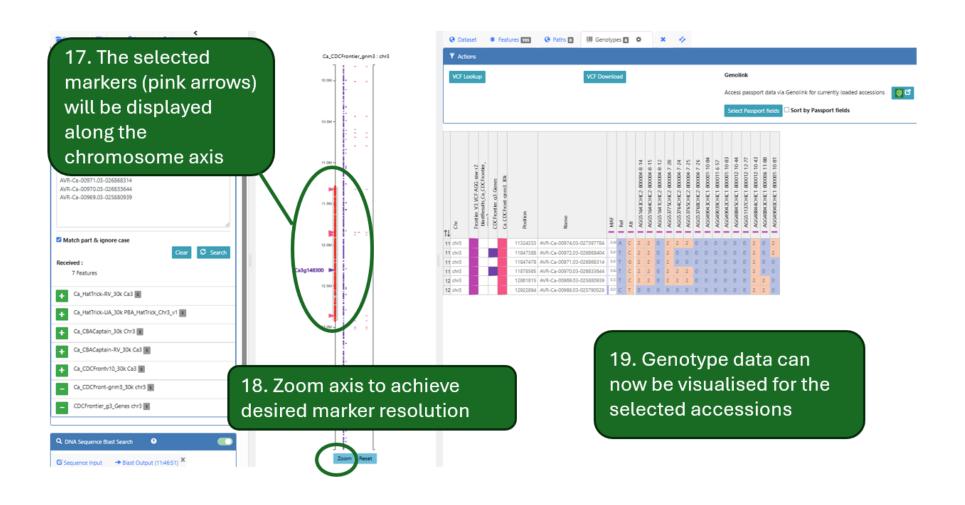


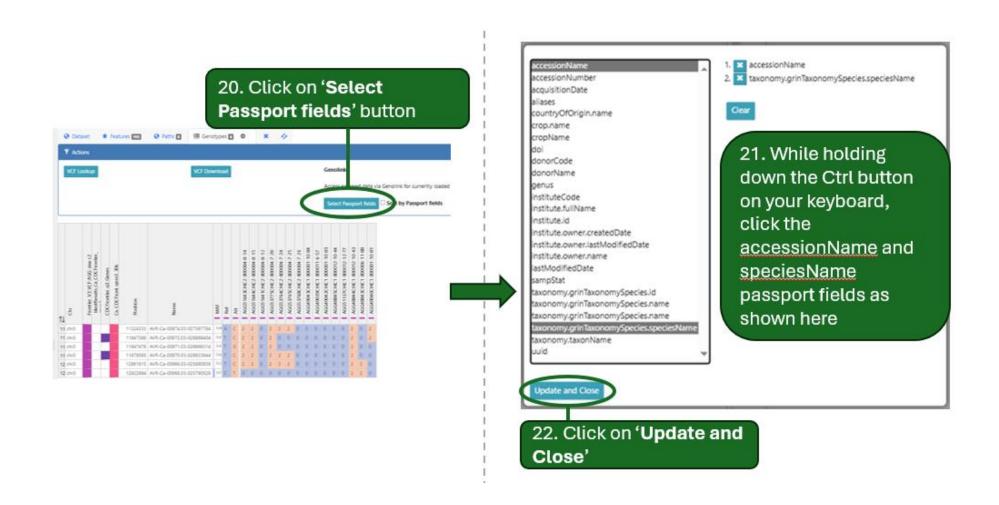
Since we already have the Genotype IDs for chickpea samples of interest, we can directly copy and paste them into the lower box. Copy and paste Genotype IDs from the electronic file 'Pretzel_workshop_Task15_AGG_chickpea_sample_list.txt' as shown below.



In order to find the markers making up the haplotype of interest more easily, we will undertake a feature search to highlight the position of the markers on the chromosome.







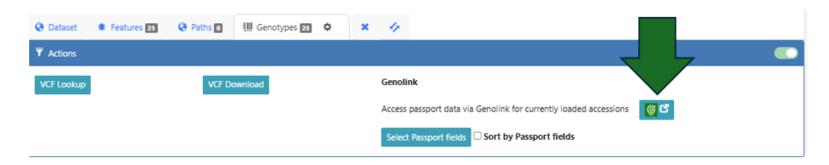


				Haplotypes predominant in wild			Haplotypes predominant in cultivated				
Marker Name	Position	REF	ALT	1	2	3	4	5	6	7	8
AVR-Ca-00974_03-027397784	11,324,333	Α	С	ALT	ALT	ALT	ALT	REF	ALT	REF	ALT
AVR-Ca-00972_03-026868404	11,847,388	Т	С	ALT	REF	ALT	ALT	REF	REF	ALT	ALT
AVR-Ca-00971_03-026868314	11,847,478	Т	G	ALT	REF	ALT	REF	REF	REF	ALT	ALT
AVR-Ca-00970_03-026833644	11,878,565	Т	С	ALT	ALT	ALT	ALT	REF	REF	ALT	ALT
AVR-Ca-00969_03-025880939	12,861,815	Т	С	ALT	ALT	./.	ALT	REF	REF	REF	REF

Based on these results, we can find chickpea hybrid lines that do and do not contain wild-like haplotypes at the flowering time locus. Interestingly, PBA Slasher contains a wild-type haplotype, suggesting it may have wild chickpea in its pedigree.

Additional tasks for a challenge

• Search for chickpea accessions in the AGG that contain wild haplotype 2 around the plant height gene in Task 12, then click on the link to Genolink to view passport data for these accessions derived from Genesys PGR.



• Use Genolink (https://genolink.plantinformatics.io/) to find the AGG barcode for the chickpea accession 'Almaz' that has been genotyped and load AGG genotype data in *Pretzel*. Load genotypes for both plant height (chr X) and flowering time (chr Y) in single genotype table by selecting markers on both chromosomes.

References

Cook D. (2022) Genome assembly for CDC Frontier - CDCFrontier.gnm3.QT0P (https://data.legumeinfo.org/Cicer/arietinum/genomes/CDCFrontier.gnm3.QT0P/)

Laurie RE, Diwadkar P, Jaudal M, Zhang L, Hecht V, Wen J, Tadege M, Mysore KS, Putterill J, Weller JL, Macknight RC. (2011) The Medicago FLOWERING LOCUS T homolog, MtFTa1, is a key regulator of flowering time. Plant Physiol. 156:2207-2224. doi: 10.1104/pp.111.180182

Malik N, Basu U, Srivastava R, Daware A, Ranjan R, Sharma A, Thakro V, Mohanty JK, Jha UC, Tripathi S, Tyagi AK, Parida SK. (2023) Natural alleles of Mediator subunit genes modulate plant height in chickpea. Plant J. 116:1271-1292. doi: 10.1111/tpj.16423

Mugabe D, Coyne CJ, Piaskowski J, Zheng P, Ma Y, Landry E, McGee R, Main D, Vandemark G, Zhang H, Abbo S. (2019) Quantitative Trait Loci for Cold Tolerance in Chickpea. Crop Science. 59: 573-582. doi: 10.2135/cropsci2018.08.0504

Varshney RK, Song C, Saxena RK, Azam S, Yu S, Sharpe AG, Cannon S, Baek J, Rosen BD, Tar'an B, Millan T, Zhang X, Ramsay LD, Iwata A, Wang Y, Nelson W, Farmer AD, Gaur PM, Soderlund C, Penmetsa RV, Xu C, Bharti AK, He W, Winter P, Zhao S, Hane JK, Carrasquilla-Garcia N, Condie JA, Upadhyaya HD, Luo MC, Thudi M, Gowda CL, Singh NP, Lichtenzveig J, Gali KK, Rubio J, Nadarajan N, Dolezel J, Bansal KC, Xu X, Edwards D, Zhang G, Kahl G, Gil J, Singh KB, Datta SK, Jackson SA, Wang J, Cook DR. (2013) Draft genome sequence of chickpea (*Cicer arietinum*) provides a resource for trait improvement. Nat Biotechnol. 31:240-246. doi: 10.1038/nbt.2491

Useful URLs

<u>Pretzel</u>
 https://agg.plantinformatics.io/

 Pretzel documentation https://docs.plantinformatics.io/

Australian Grains Genebank (AGG) website
 https://agriculture.vic.gov.au/crops-and-horticulture/australian-grains-genebank

 Australian Grains Genebank (AGG) online portal - GRIN-Global Database https://ausgenebank.agriculture.vic.gov.au/gringlobal/search

Genesys-PGR
 https://www.genesys-pgr.org/

 Genolink https://genolink.plantinformatics.io/

 AGG Dataverse for accessing genotype data https://dataverse.harvard.edu/dataverse/australiangrainsgenebank

 NCBI Genbank Nucleotide site https://www.ncbi.nlm.nih.gov/nucleotide/

Appendix

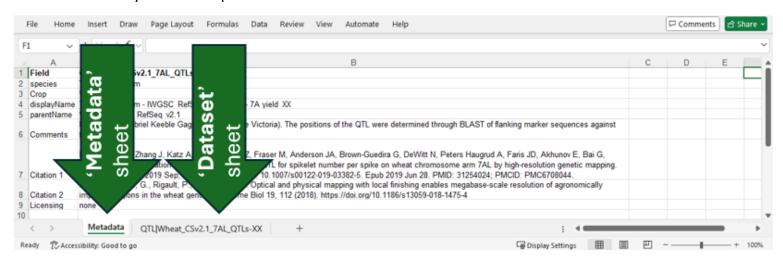
Pretzel data upload templates

Overview

Data can be uploaded into *Pretzel* using several methods, but the Excel upload method is recommended for most users. Three types of data (alignments, genetic maps and QTLs) can be independently uploaded to *Pretzel* as custom datasets, as described below. Note that it is also possible to upload genomes, BLAST databases and vcf files with assistance from the *Pretzel* development team.

Data type	Dataset Name Format	Description
Alignment	Alignment DataName	Refers to a feature with a defined position or interval, except for a QTL. This is commonly used for
		markers or genes, which are typically anchored to a specified genome assembly. Units are in base pairs
		(bp).
Мар	Map DataName	Used for uploading genetic maps, where units for the defined position are in centiMorgans (cM).
QTL	QTL DataName	Used for uploading QTLs. The defined interval positions can relate to either a genome assembly (where
		coordinates are in bp), or a genetic map (where coordinates are in cM),

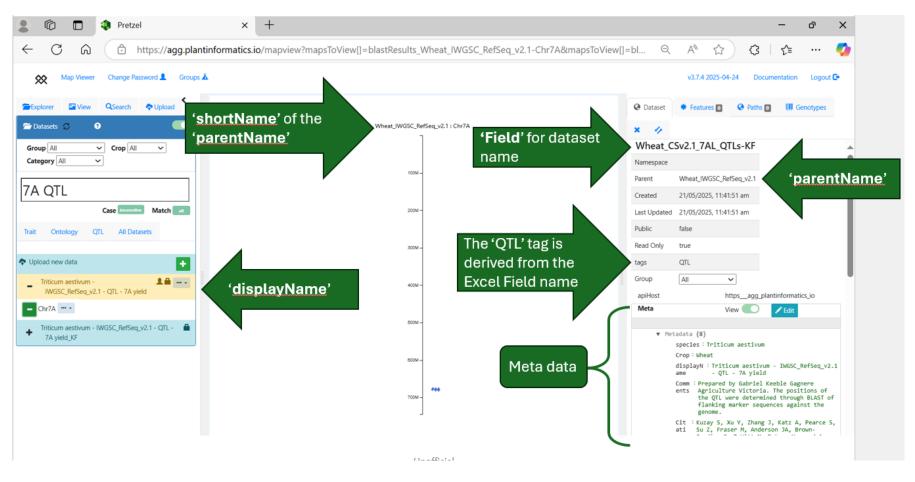
The Excel template file consists of two sheets, the 'Metadata' sheet for defining parameters and providing meta data, and the 'Dataset' sheet that contains data you wish to upload into *Pretzel*.



Metadata sheet

The figure and table below describe the fields in the 'Metadata' sheet for uploading a custom QTL dataset into *Pretzel*. The mandatory fields are used by *Pretzel* and are case sensitive; they must match the relevant dataset in *Pretzel* exactly to link datasets effectively e.g. to link a QTL to a genome assembly. The optional parameters below are used to display meta data in the 'Dataset' panel to align with FAIR (Findable, Accessible Interoperable, Re-usable) data standards. We recommend using the parameters below as best practise. Additional optional columns can be added if desired.

The figure below shows where the mandatory parameters in the 'Metadata' sheet get incorporated into Pretzel.

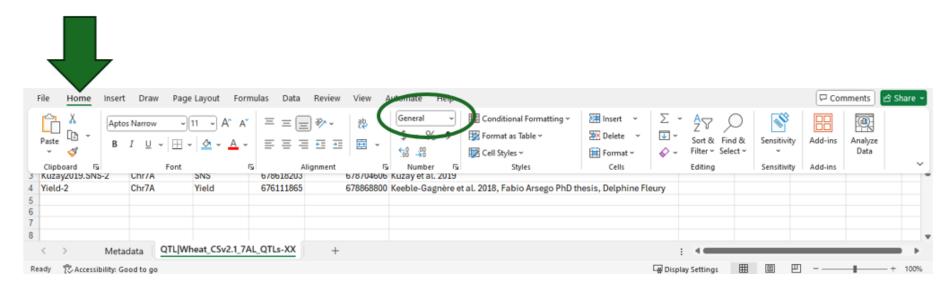


Parameter	Details	Usage			
		QTL	Alignment	Мар	
Field	Name of the custom dataset to be loaded into <i>Pretzel</i> . It is structured as the type of dataset on the LHS, a pipe symbol ' ' and the unique name of the dataset on the RHS. E.g. 'QTL Dataset_XX', 'Alignment Dataset_XX' or 'Map Dataset_XX' This dataset name must be unique, it cannot match any other datasets loaded in <i>Pretzel</i> , including	mandatory	mandatory	mandatory	
	datasets loaded by other users. Note that this dataset name is NOT used in the left-hand side Explorer panel, it is only displayed in the 'Dataset' panel (see Figure below).				
	This must exactly match the dataset name on the RHS of the pipe symbol in the Excel worksheet name i.e. the portion in bold e.g. QTL Wheat_CSv2.1_7AL_QTLs-KF				
	The name of the dataset is limited to 31 characters (including e.g. 4 characters used for 'QTL ') due to limitations set by MS Excel on the sheet name. Take care to ensure there are no spaces.				
Species	Taxonomy e.g. Triticum aestivum	mandatory	mandatory	mandatory	
Crop	Used to filter the datasets in the 'Explorer' panel based on the crop type e.g. Wheat	mandatory	mandatory	mandatory	
displayName	Complete, descriptive name that shows up in the left-hand 'Explorer' panel	mandatory	mandatory	mandatory	
parentName	Name of the genome assembly or genetic map that the feature (e.g. QTL, marker, gene) is being anchored to. This must exactly match the dataset name in <i>Pretzel</i> .	mandatory	mandatory	not applicable	
shortName	Short dataset name displayed at the top of the chromosome axis in <i>Pretzel</i> . Required only for a genome or genetic map dataset.	not applicable	not applicable	mandatory	
platform	Type of technology platform used to generate data e.g. 90k, 40k, SSRs, etc. If multiple types of markers were used e.g. in the case of a consensus map, then the use of 'multiple' is recommended.	not applicable	not applicable	mandatory	
licensing	E.g. CC BY 4.0 (or 'none' if not applicable)	optional	optional	optional	
Citation	Reference to a manuscript, 'personal communication (Researcher name)', 'unpublished' etc.	optional	optional	optional	
DOI	DOI links the uploaded data to the source e.g. manuscript	optional	optional	optional	
Comments	Longer description of the dataset, and/or information on the study, or how the data was compiled. This could include:	optional	optional	optional	
	 parents and filial generation of a mapping population information about the germplasm (e.g. diverse landraces) trait information analysis and/or filtering steps taken 				
Contact	Person / organisation who uploaded the data (e-mail address encouraged)	optional	optional	optional	

Dataset sheet

The fields in the 'Dataset' sheet are very specific to the type of data being uploaded. The names of the mandatory parameters in the LHS column below are case sensitive.

It is important that the numbers entered into the 'Start', 'End' and 'Position' fields are formatted correctly. In excel, highlight the relevant cells and set the format to 'General', ensuring they do not contain any commas or thousand separators e.g. '678380602' not '678,380,602'.



Additional, optional columns can be added as desired to capture other information e.g. human-readable gene descriptions, gene identifiers.

Parameter	Details	Usage				
		QTL	Alignment	Мар		
Name	Name of the QTL or feature	mandatory	mandatory	not applicable		
Marker	Name of the markers in the genetic map	not applicable	not applicable	mandatory		
Chromosome	Name of the chromosome the feature is anchored to. You must check the name of the chromosome in the relevant dataset in <i>Pretzel</i> for the 'parentName'. E.g. if the chromosome in <i>Pretzel</i> is called 'Chr7A', then the use of e.g. '7A' will fail to anchor the feature to the chromosome.	mandatory	mandatory	mandatory		
Trait	Name of the trait being investigated	mandatory	not applicable	not applicable		
Start	Start position of the QTL. If you have anchored the QTL to a chromosome in a genome assembly, this should be the number in basepairs. If you have anchored the QTL to a genetic map, then the number should be in cM.	mandatory	mandatory	not applicable		
End	End position of the QTL. If the QTL is associated with a SNP marker, then the Start and End positions should be the same number.	mandatory	mandatory	not applicable		
Position	Position of the marker in the genetic map	not applicable	not applicable	mandatory		
Ontology	Crop ontology term from https://cropontology.org/ e.g. CO_321:0000902. The data contained in this field enables <i>Pretzel</i> to colour-code QTLs based on crop ontology terms.	optional	not applicable	not applicable		
Reference	Reference to a published manuscript / 'personal communication (Researcher name)', 'unpublished' etc. It can be especially helpful to list a reference for each entry if data has been compiled from multiple sources.	optional	optional	optional		
Comments	Longer description of the dataset, and/or information on the study, or how the data was compiled.	optional	optional	optional		

Troubleshooting

If you encounter a bug and *Pretzel* is no longer responding as expected, reset your session. The easiest way to do this is put your cursor in the URL field of your browser and delete characters back to the main URL: https://agg.plantinformatics.io/, then press 'Enter' on your keyboard.

If any software bugs are causing problems, please report them to the *Pretzel* development team: info@plantinformatics.io.

Contact Details

For requests to set up a group training session:

Kerrie Forrest - Senior Research Scientist

Agriculture Victoria (AgriBio, Bundoora)

kerrie.forrest@agriculture.vic.gov.au

For questions about data management, digital tools and Program 1 of the AGG Strategic Partnership:

Gabriel Keeble-Gagnere - Science Program Leader – Australian Grains Genebank

Agriculture Victoria (AgriBio, Bundoora)

gabriel.keeble-gagnere@agriculture.vic.gov.au

For questions about the Australian Grains Genebank and Program 2 of the AGG Strategic Partnership:

Sally Norton - Genebank Curator and Research Leader of the Australian Grains Genebank

Agriculture Victoria (Horsham)

sally.norton@agriculture.vic.gov.au