



OPEN

## PlantPathMarks (PPMdb): an interactive hub for pathways-based markers in plant genomes

Morad M. Mokhtar<sup>1</sup>, Achraf El Allali<sup>1✉</sup>, Mohamed-Elamir F. Hegazy<sup>2</sup> & Mohamed A. M. Atia<sup>3✉</sup>

Over the past decade, the problem of finding an efficient gene-targeting marker set or signature for plant trait characterization has remained challenging. Many databases focusing on pathway mining have been released with one major deficiency, as they lack to develop marker sets that target only genes controlling a specific pathway or certain biological process. Herein, we present the PlantPathMarks database (PPMdb) as a comprehensive, web-based, user-friendly, and interactive hub for pathway-based markers in plant genomes. Based on our newly developed pathway gene set mining approach, two novel pathway-based marker systems called pathway gene-targeted markers (PGTMs) and pathway microsatellite-targeted markers (PMTMs) were developed as a novel class of annotation-based markers. In the PPMdb database, 2,690,742 pathway-based markers reflecting 9,894 marker panels were developed across 82 plant genomes. The markers include 691,555 PGTMs and 1,999,187 PMTMs. Across these genomes, 165,378 enzyme-coding genes were mapped against 126 KEGG reference pathway maps. PPMdb is furnished with three interactive visualization tools (Map Browse, JBrowse and Species Comparison) to visualize, map, and compare the developed markers over their KEGG reference pathway maps. All the stored marker panels can be freely downloaded. PPMdb promises to create a radical shift in the paradigm of the area of molecular marker research. The use of PPMdb as a mega-tool represents an impediment for non-bioinformatician plant scientists and breeders. PPMdb is freely available at <http://ppmdb.easyomics.org>.

Today, an enormous amount of released biodata collections stimulates the development of computerized applications worldwide<sup>1</sup>. Inspiration drives scientists to develop intelligent approaches to extract hidden knowledge behind these massive amounts of data. Biodata mining attempts to find novel, reliable, useful, and meaningful insights from these vast amounts of data. Over the last decades, the scope of biodata mining has expanded from genome-mining to phenome-mining approaches<sup>2</sup>. Its applications outside of basic research have become extensive to study plants' genetic diversity and improve economic crop breeding programs. Generally, plants produce various chemical compounds involving nutritional or medicinal benefits<sup>3</sup>. Our understanding of plant natural products' biological pathways is still insufficient; however, scientists assume that genomic and metabolomic information can provide clues about unidentified enzymes and reactions involved in a particular biosynthesis process<sup>4</sup>.

Plants have undergone complicated evolutionary events that have resulted naturally or are due to human-made plant breeding experiments and finally lead to polyploidy or genome duplication<sup>3</sup>. This duplication drives the availability of genes either as homologous or paralogous. Due to the phylogenetic relationship of homologous genes, they continue to preserve their core functions. Gene products are often mapped to many pathways revealing the same catalyst properties under different reaction conditions or showing the same reaction but in distinct subcellular locations<sup>5</sup>. Such an association provides a mechanism to study genes' roles; their expression may be regulated spatially or temporally due to evolutionary implications on plants, gene products, gene functions, and, finally, developmental stages<sup>1</sup>.

A biological pathway is broadly defined as a series of actions among cell molecules that lead to a particular product or cellular change. It can also trigger novel molecule assembly, turn genes on and off, or drive a cell to move. Any biological pathway comprises a cluster of genes that exhibit identical or complementary biological

<sup>1</sup>African Genome Center, Mohammed VI Polytechnic University, Ben Guerir, Morocco. <sup>2</sup>Chemistry of Medicinal Plants Department, National Research Centre, Giza 12622, Egypt. <sup>3</sup>Molecular Genetics and Genome Mapping Laboratory, Genome Mapping Department, Agricultural Genetic Engineering Research Institute (AGERI), Agriculture Research Center (ARC), Giza 12619, Egypt. ✉email: [achraf.elallali@um6p.ma](mailto:achraf.elallali@um6p.ma); [matia@ageri.sci.eg](mailto:matia@ageri.sci.eg)

functions<sup>6</sup>. Nowadays, pathway analysis is a flourishing research area in systems biology. It relies on extracting knowledge from raw data generated from high-throughput sequencing technologies by building a model that describes and summarizes underlying biological processes. These high-throughput technologies mostly provide a list of differentially expressed genes (DEGs) between a control and a case of interest. Unfortunately, generated DEG data are usually removed from their biological context, causing the resultant genes to require further validation against biological pathways to prove their biological significance<sup>7</sup>. Notably, each of these genes may be a part of many pathways, which usually start with the expression of a gene of a particular ligand and is terminated with an alteration in the concentration of a specific signaling micro-molecule<sup>8</sup>.

The dilemma of plant trait/metabolite characterization and improvement based on single or few genes targeting marker selections remains challenging, and many approaches have been developed to tackle these issues<sup>9</sup>. For decades, massive studies have aimed to determine a few differentially expressed marker genes participating in certain natural products' biosynthesis processes or conferring a particular phenotype<sup>10–12</sup>. Nevertheless, these methods may not be accurate and may lack in capturing an in-depth snapshot of biological processes. Therefore, studies focusing on gene sets belonging to a particular pathway to investigate and explain phenotypic changes have gained increased popularity. These methods are more acceptable, as they assess the significance of a group of predefined genes with shared biological functions to explain specific phenotypes<sup>13–15</sup>.

In contrast to traditional gene-based methods, pathway-based methods combined with biological knowledge can help gain a better understanding of functional insights into phenotypic differences. Pathway-mining approaches attempt to rectify such drawbacks by utilizing the available biological knowledge base about the structures and operations of biological pathways accompanying computational methods<sup>5</sup>. Therefore, instead of focusing on specific gene-based marker development, our proposed pathway gene set mining (PGM) approach is critically more appropriate and promises to create a radical shift in the paradigm of molecular markers as a research area. The dramatically increasing number of sequenced plant genomes is supposed to reveal more basic architectural principles of biosynthetic pathways for generating chemical diversity of natural products. Based on released knowledge on interactions between genes and their phenotypic reflections, many publicly available databases have been developed for storing this pathway mining information<sup>16</sup>.

After numerous plant genomes are fully sequenced, scientists often search for an efficient scheme that can analyze their data and answer certain biological questions<sup>3</sup>. Common types of analyses include pathway enrichment analysis and the analysis of gene ontology and gene–gene interactions. Over the last decade, many powerful resources for studying metabolic pathways in plants, such as KEGG<sup>17</sup>, Gramene<sup>18</sup>, Plant Reactome<sup>19</sup>, MapMan<sup>20</sup>, MetaCyc<sup>21</sup>, Plant Metabolic Network<sup>22</sup>, and BioCyc<sup>23</sup>, have been developed. These databases try their best to integrate the most available information when they were initially launched. Their construction process includes manual curation and computational efforts, regular screening of newly released publications, and updating of interactors. Despite extraordinary efforts, no single database is fully equipped with all the preferred information and may not accommodate all the necessary bioinformatics tools. Notably, all these resources have the significant limitation of missing genome-scale marker sets that target only the genes responsible for a specific pathway or that are involved in/control a particular biological process at the level.

Gene sets involved in a particular pathway often work together in an orchestral pattern to reflect a specific biological function. Herein, we proposed a novel PGM pipeline/approach to develop a specific marker set for each pathway at the genome-scale level in plant genomes.

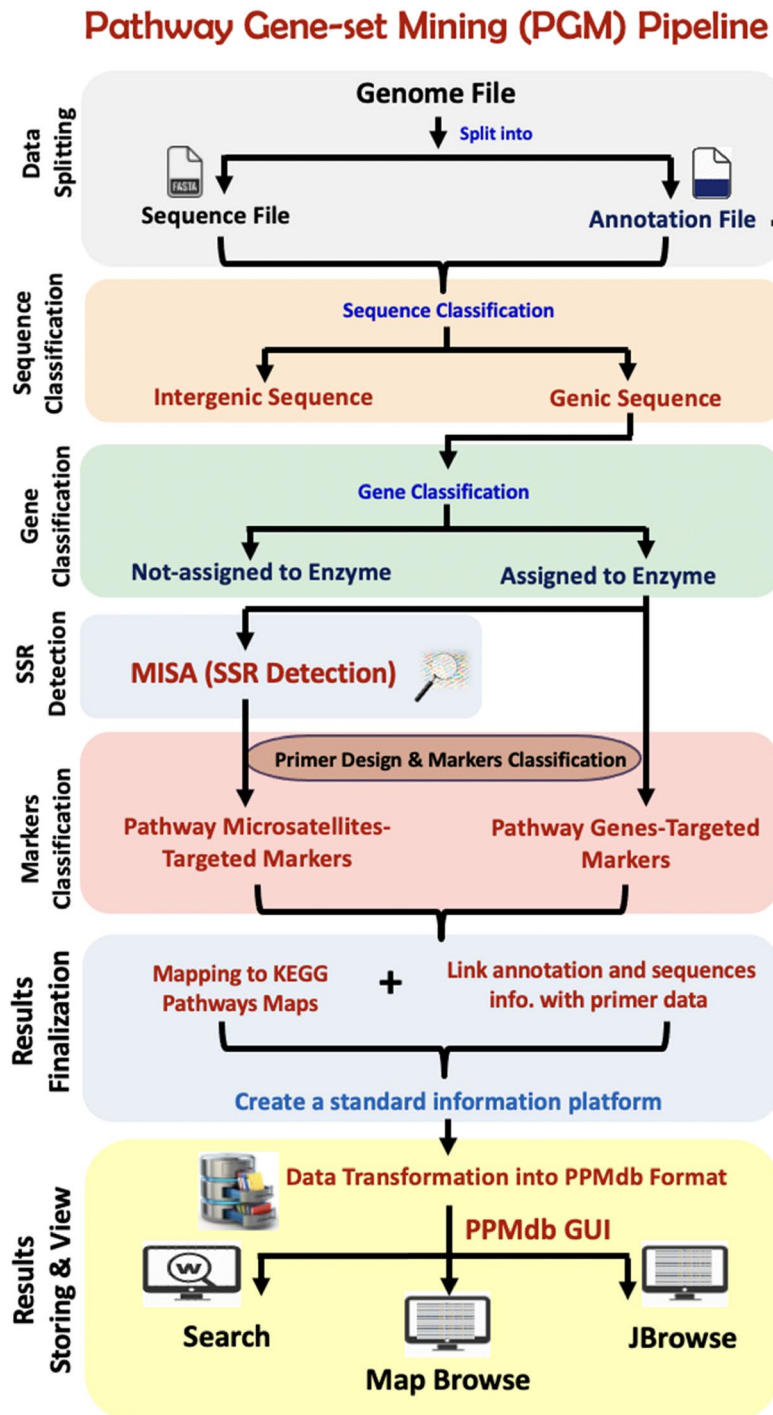
Microsatellites, also known as simple sequence repeats (SSRs), are a set of one to six nucleotides that are repeated tandemly and classified either into perfect, imperfect, or composite repeats<sup>24,25</sup>. Due to their unique characteristics (hypervariability, multi-allelic nature, high reproducibility, Mendelian inheritance, and overall high abundance in plant genomes), as well as their essential roles in many biological functions (modulating transcription factor binding, tuning knobs of gene expression levels, organizing chromatin, nucleosome positioning, acting as recombination hot spots, and facilitating unusual structural conformations), microsatellites have become one of the most preferred choices among all genetic markers because they afford a molecular basis for plant genome evolution and demonstrate fast acclimatization to many abiotic stresses<sup>26</sup>. Therefore, we proposed a mining schema of microsatellites within gene sets involved in all plant genome pathways as a robust and functional approach to develop a novel marker system called pathway microsatellite-targeted markers (PMTMs).

Ultimately, this study aims to develop a comprehensive hub for two novel pathway-based marker systems coined pathway gene-targeted markers (PGTMs) and pathway microsatellite-targeted markers (PMTMs) in plant genomes. Furthermore, the developed markers were mapped against reference KEGG pathway maps interactively. The use of these powerful portal represents an impediment for non-bioinformatician plant scientists and breeders.

## Materials and methods

**Data collection.** To build a comprehensive pathway-based marker database in plants, we retrieved complete genome sequences and annotations of 82 various plant species, including 72 eudicots and 10 monocots. The sequences and annotations were downloaded as GenBank files from NCBI's FTP site (<ftp://ftp.ncbi.nlm.nih.gov/>). Furthermore, we downloaded all KEGG gene annotations available on the KEGG database of these plant genomes (approximately 2,296,000 genes) in addition to all KEGG pathway reference maps (<https://www.genome.jp/kegg/>).

**Data analysis and database construction.** To construct the PlantPathMarks (PPMdb) core database, we integrated two powerful tools called M<sup>I</sup>croS<sup>A</sup>tellite identification (MISA)<sup>27</sup> and primer<sup>3</sup><sup>28</sup> into our in-house-developed “Perl and Shell” scripts to strap all standalone bioinformatics analysis steps in one main pipeline called the PGM pipeline. This pipeline was built to develop two novel marker systems called PGTM and PMTM.



**Figure 1.** The workflow of pathway gene-set mining (PGM) pipeline.

This pipeline was implemented to construct our PPMdb sub-databases through seven main steps: (a) splitting the retrieved genomic data into sequence and annotation files, (b) classifying genomic sequence into genic or intergenic, (c) sub-classifying genes into enzyme coding or non-enzyme coding, (d) mining microsatellites on enzyme-coding genes only (this step exclusive to PMTM), (e) designing primers and classifying markers into PMTM or PGTM, (f) mapping all designed primers against KEGG reference maps<sup>17</sup> besides linking associated information (sequence and annotation) to all pathway-mapped markers, and (g) data transforming into the PPMdb database format and implementing all datasets into the PPMdb web interface (Fig. 1).

For effective marker development, we used a straightforward strategy to boost data analysis and marker generation processes. Before we began, we established in-house “Perl and Shell” scripts to convert all retrieved genome sequence and annotation data and unify input and output files used for the development of the PPMdb database.

Initially, the data retrieved from the 82 plant genomes were dissected to classify each genome into genic and intergenic according to the genome's annotation information. This classification step was achieved with the aid of the gene-finding format (GFF) file for each genome. The gene ID list for each plant was further searched against the KEGG database (<https://www.genome.jp/kegg/genes.html>) to sub-classify genes into enzyme coding or non-enzyme coding. Thereafter, perfect and compound microsatellite motifs were identified using the MISA tool<sup>27</sup> according to the following parameters: mono- ( $\geq 9$ ), di- ( $\geq 4$ ), tri- ( $\geq 3$ ), tetra- ( $\geq 3$ ), penta- ( $\geq 2$ ), and hexa-nucleotide ( $\geq 2$ ). The parameter was set as  $\geq 2$  repeats interrupted by  $\leq 20$  bp for the characterization of compound microsatellite motifs<sup>29</sup>.

PMTM and PGTM primers were then designed for all enzyme-coding genes for all plant genomes using the Primer3 software<sup>28</sup> according to the following criteria: a) optimal primer length of 20 bp; b) optimal melting temperature of 55 °C; c) product size range of 90–500 bp for PMTM primers and 100–1000 bp for PGTM primers; and d) a 50% G/C content. Then, e) a unique primer ID was assigned.

All developed PMTM and PGTM primers were mapped against the KEGG pathway reference maps (<https://www.genome.jp/kegg/pathway.html>) by using in-house “Perl and Shell” scripts to integrate and localize our developed PMTM and PGTM primers and build a map set of marker-enriched pathways for each plant. Each map contains a pathway ID, pathway image, IDs of mapped enzymes (highlighted), PMTM/PGTM markers associated with mapped enzymes, annotation information of each marker, and other information related to this primer (Tm, GC%, length, etc.). All developed maps were visualized in an attractive user form using JQuery ImageMapster Plugin<sup>30</sup>.

The JBrowse tool<sup>31</sup> was also integrated into our PPMdb database as a powerful web application for genome analysis and visualization. Herein, it was used to map and browse the identified microsatellite motifs and developed PMTMs and PGTM and their linked information for each organism.

All generated data of PMTMs, PGTM, and maps were further processed to build a standard information platform for all marker types before integrating them into the PPMdb SQL database. With the aid of the in-house “Perl and Shell” scripts, we batch processed all developed markers, maps, and any associated data and converted them into a consistent format. Additionally, we uploaded these scripts to online open-source housing website (GitHub) to make it available for all users at the following link (<https://github.com/MoradMMokhtar/PlantPathMars-Scripts.git>). The PPMdb database was developed with aid of the LPPM (Linux + Perl + PHP + MySQL) web application platform; finally, JavaScript, CSS, and HTML languages were used to design a user-friendly interactive web interface.

## Features and utilities

**Database interface.** PPMdb presents an interactive, user-friendly portal well equipped with many features to enable users to search and download PMTMs and PGTM across 82 plant genomes. PPMdb provides users with a navigation bar designed to help access the PPMdb database sections and tools in a responsive and convenient way. The PPMdb data can be straightforwardly browsed and retrieved via nine interactive pages: Home (PPMdb Quick-Access), Database Search, Map Browse, JBrowse, Statistics, Data Resources, Bulk Download, Species Comparisons and Manual. Under these pages, plants are taxonomically grouped or sorted alphabetically to improve the convenience of exploration and selection.

The Home page introduces PPMdb as a comprehensive database along with its PGM pipeline by describing the PPMdb database sections, in addition to providing users with a “PPMdb Quick-Access; Analyzed Genomes in PPMdb Database Section” option for all analyzed plant genomes. In the “PPMdb Quick-Access” section, plant genomes analyzed in the PPMdb database are categorized according to plant type/class (dicot plants or monocot plants). Under each group, interactively, the plants are sorted alphabetically to enable straightforward access and searches within each plant genome independently through “Map Browse” or “JBrowse” quick-access links (Fig. 2).

The Database Search page provides users with a searching utility within two main sections: (a) the Search section, which enables users to obtain results by selecting the following types of interest keywords: pathway name (e.g., Glycolysis), marker type (e.g., PGTM), and organism name (e.g., *Arabidopsis lyrata*), from the available drop-down menus to effortlessly access and retrieve all marker and map data stored in the PPMdb database. (b) The KEGG Pathway Maps Overview section, which provides users the necessary information for each pathway, such as pathway ID (hyperlinked to KEGG ref. pathways), pathway class (e.g., metabolism), pathway sub-class (e.g., carbohydrate metabolism), and pathway map. The search results of this section are designed in a user-friendly style containing essential descriptive information for each marker (e.g., enzyme ID, marker type, repeat type and sequence [exclusive to PMTM], primer sequence, primer annealing temperature, primer position within CDS, product size, JBrowse view link, NCBI gene accession number and CDS sequence, and gene information) (Fig. 3).

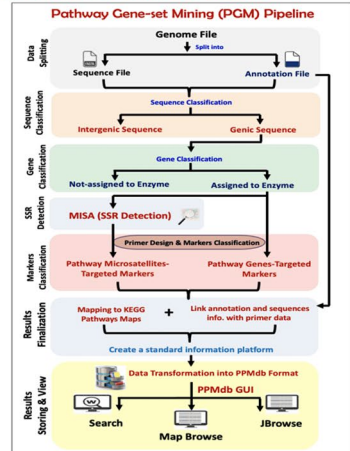
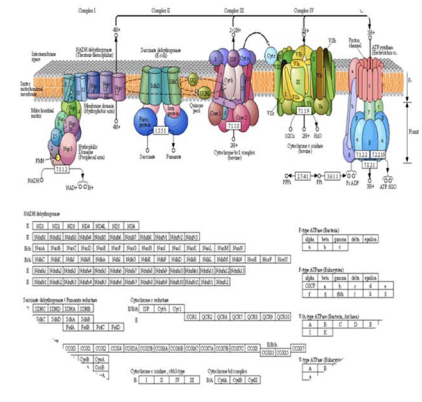
In addition, general statistic tables are presented to summarize the total number of designed primers, number of PGTM primers, number of PMTM primers, and number of failed designed PMTM primers within the selected pathway of a particular plant. All search results can be downloaded freely for each plant straightforwardly as a CSV file (Fig. 3).

The Map Browse page offers users with searching utility in two styles inside the same page: (a) selecting the plant of interest from the page-side plant list for a direct access of pathways and markers, (b) selecting the plant of interest from drop-down menus to directly access and obtain all markers and maps. In both styles, the search results are presented in a convenient, visualization-supported, and well-dissected manner involving essential information for each pathway, such as pathway ID (hyperlinked to KEGG ref. pathways), pathway class (e.g., metabolism), pathway sub-class (e.g., carbohydrate metabolism), pathway map browse (hyperlinked to interactive pathway map localizing all developed PMTM and PGTM markers on the pathway map), pathway description, and pathway map (Fig. 4).



Brief introduction describing the PPMdb sections & statistics

## PlantPathMarks (PPMdb) A hub for Plant Pathways-Based Markers



- Valuable Resources
- NCBI
  - SSRome
  - VIODB
  - MISAweb
  - Primer3web
  - EMBL-EBI
  - Phytozome
  - TAIR
- Useful Links
- CPAN
  - MySQL
  - KEGG

The PlantPathMarks (PPMdb) portal is an extensive website providing interactive and convenient access and mapping to thousands of pathways-based primers/markers in 82 plant genomes. Here, we also developed a novel approach called the Pathway Gene-set Mining (PGM). Besides, new marker systems were developed based on their pathways-based mapping of developed markers. The systems were named Pathway Microsatellites-Targeted Markers (PMTM) and Pathway Genes-Targeted Markers (PGTM); as a kind of annotation-based marker. In the PlantPathMarks hub, 9894 marker-panels (comprising 2,690,742 markers) were developed across eighty-two plant genomes. This huge number of markers includes 1,999,187 PMTM and 691,555 PGTM markers

A

**"PPMdb Quick-Access" Section**

Analyzed Genomes in PPMdb Database

Plant Class	Genome	Map Browse	JBrowse	Image
Dicot Plants	Arabidopsis lyrata	Map Browse	JBrowse	
	Arabidopsis thaliana	Map Browse	JBrowse	
Monocot Plants	Arachis duranensis	Map Browse	JBrowse	
	Arachis ipaensis	Map Browse	JBrowse	
	Beta vulgaris	Map Browse	JBrowse	
	Brassica napus	Map Browse	JBrowse	

The analyzed plant genomes are categorized according to plant class (dicot or monocot)

Users can directly access all KEGG Pathways Maps and their mapped markers under this plant genome

Users can directly access and browse all marker types mapped across this plant genome in a user-friendly style containing essential descriptive information for each marker

B

Figure 2. Screenshot of the PPMdb homepage (A) Database description and (B) PPMdb Quick-Access section.

The "JBrowse" page enables the users to visualize and map all identified microsatellite motifs besides all developed PMTM and PGTM markers against the Refseq CDS for each plant genome by selecting the plant of interest from the page-side plant list. The mapped microsatellite motifs and developed PMTM/PGTM markers were linked to its essential information (Fig. 5).

The 'Statistics' page was designed to provide users a primary indication of the taxonomy ID, number of mapped pathways, number of mapped enzymes, number of designed primers, number of PGTM primers, number of PMTM primers, and number of failed designed PMTM primers for each plant species stored in the PPMdb database.

The Data Resources page provides users with hyperlinks of all types/formats of the data utilized to build the PPMdb database. The page includes hyperlinks of taxonomy ID, Genbank file, Fasta file, GFF file, and KEGG gene annotation for each plant species analyzed within the PPMdb.

**PlantPathMarks (PPMdb)**  
A hub for Plant Pathways-Based Markers

Home Database Search Map Browse JBrowse Statistics Data Resources Bulk Download Species Comparisons Manual

**Search Section**

- \* Please select the Pathway Name from Pathway Name drop down menu.
- \* Please select Marker type (PGTM, PMTM) from Marker Type drop down menu.

Pathway Name: Glycolysis / Gluconeogenesis | Marker Type: PGTM | Organism Name: Arabidopsis lyrata | Search

**KEGG Pathway Maps Overview**

Metabolism

Genetic Information Processing

Environmental Information Processing

Cellular Processes

Organismal Systems

**Pathway Name: Glycolysis / Gluconeogenesis**

Pathway ID: k000010

Class & Sub-Class: Class: Metabolism Sub-Class: Carbohydrate metabolism

Description: Glycolysis is the process of converting glucose into pyruvate and generating small amounts of ATP (energy) and NADH (reducing power). It is a central pathway that produces important precursor metabolites: six-carbon compounds of glucose-6P and fructose-6P and three-carbon compounds of glyceraldehyde-3P, glyceralate-3P, phosphoenolpyruvate, and pyruvate [MD-M0000]. Acetyl-CoA: another important precursor metabolite; is produced by oxidative decarboxylation of pyruvate [MD-M00307]. When the enzyme genes of this pathway are examined in completely sequenced genomes; the reaction steps of three-carbon compounds from glyceraldehyde-3P to pyruvate form a conserved core module [MD-M0000], which is found in almost all organisms and which sometimes contains open structures in bacterial genomes. Gluconeogenesis is a synthesis pathway of glucose from noncarbohydrate precursors. It is essentially a reversal of glycolysis with minor variations of alternative paths [MD-M0000].

**Pathway Name: Citrate cycle (TCA cycle)**

Pathway ID: k000020

Class & Sub-Class: Class: Metabolism Sub-Class: Carbohydrate metabolism

Description: The citrate cycle (TCA cycle; Krebs cycle) is an important aerobic pathway for the final steps of the oxidation of carbohydrates and fatty acids. The cycle starts

**Statistics**

Total Number of designed primers: 407  
Number of PGTM primers: 133  
Number of PMTM primers: 186  
Number of Failed designed PMTM primers: 19

Download Search Result

**Enzyme Id: K02841**

Marker Type	PGTM
Repeat type & sequence	Repeat type: Repeat sequence:
Repeat start & end within CDS	Repeat start: Repeat end:
Forward primer	Primer sequence: GGACACACTGATGAGCCAA Tm(°C): 59.969 GC(%): 55 Size: 20
Reverse primer	Primer sequence: TTTCGGTAAAGCCTGACC Tm(°C): 59.962 GC(%): 55 Size: 20
Primer start, end within sequence and product size	Start: 429 End: 504 Product size (bp): 424
JBrowse View	JBrowse
NCBI Gene Accession Number & CDS Sequence	Accession Number: g333021 CDS Sequence: g333021
Gene Information	Gene Information

**Enzyme Id: K00844**

Marker Type	PGTM
Repeat type & sequence	Repeat type: Repeat sequence:
Repeat start & end within CDS	Repeat start: Repeat end:
Forward primer	Primer sequence: TGGGTTCATTCTCGAAGC Tm(°C): 59.866 GC(%): 50 Size: 20
Reverse primer	Primer sequence: GATTCGGCCTTTGGATGC Tm(°C): 59.878 GC(%): 55 Size: 20
Primer start, end within sequence and product size	Start: 427 End: 497 Product size (bp): 491
JBrowse View	JBrowse
NCBI Gene Accession Number & CDS Sequence	Accession Number: g333024 CDS Sequence: g333024
Gene Information	Gene Information

**Callouts:**

- Users can directly select pathways belonging to one of this five classes.
- The pathways categorized under these classes based on their molecular interaction, reaction and relation networks.
- Search Section enables users to obtain results by selecting the following types of interest keywords: pathway name, marker type, and organism name, from the available drop-down menus
- The search results containing essential descriptive information for each marker (e.g., enzyme ID, marker type, repeat type & sequence, primer sequence, primer annealing temperature, primer position within CDS, product size, JBrowse view link, NCBI gene accession number & CDS sequence, and gene information)

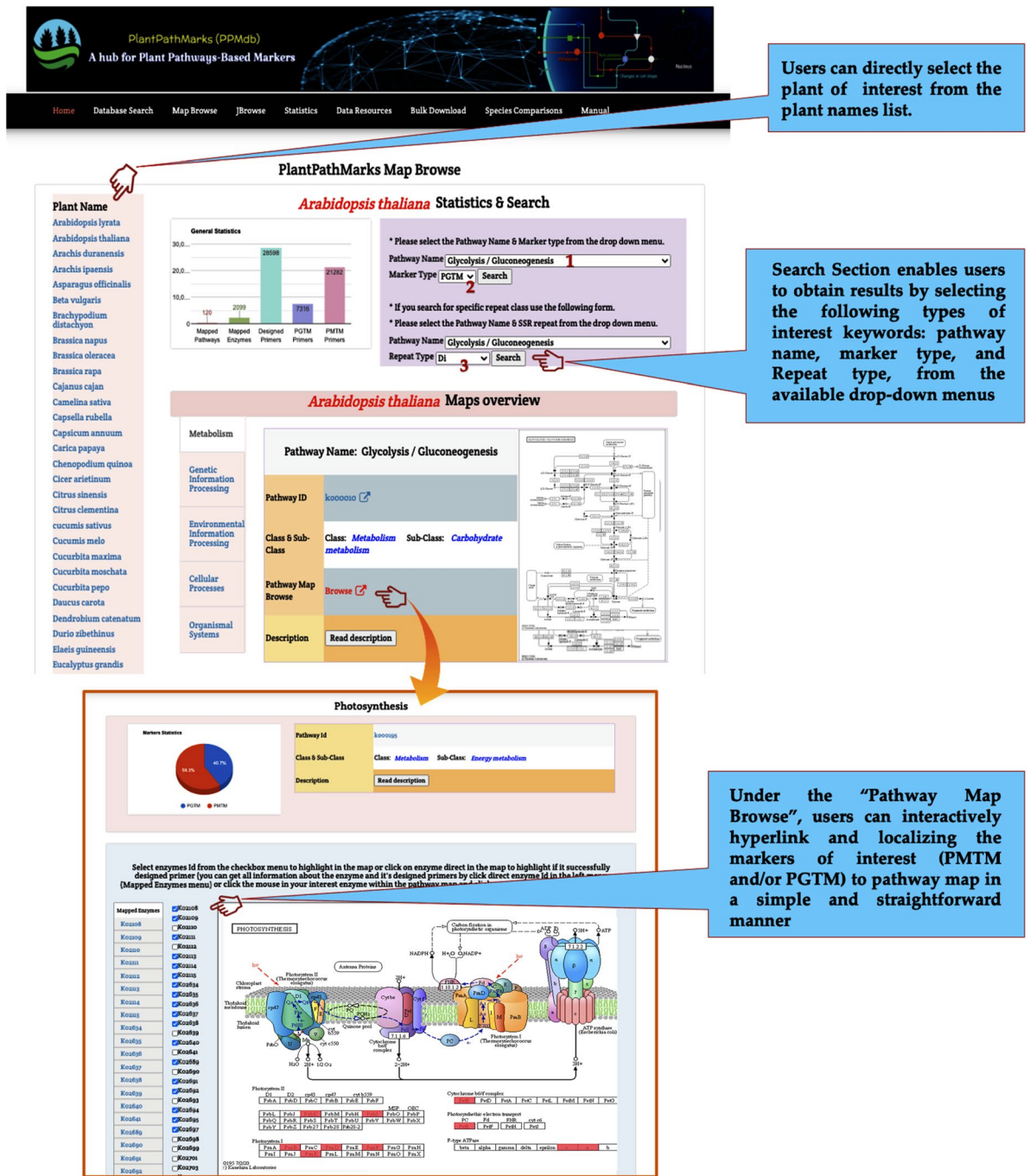
Figure 3. Screenshot of the PPMdb “Database Search” page layout.

The Bulk Download page enables users to download all developed primers and their secondary information under a specific pathway in the organism of interest through the selection of organism name (e.g., *Arabidopsis lyrata*), pathway name (e.g., Glycolysis/Gluconeogenesis), and marker type (e.g., PGTM), from the drop-down menus straightforwardly in a convenient manner.

The Species Comparisons page offers users a powerful utility to compare two or three organisms in a particular pathway. The comparison results are provided under three main sections: general statistics, marker statistics, and distribution of the various SSR classes. The general statistics section layouts the numbers of mapped pathways, mapped enzymes, total designed primers, number of designed PGTM and PMTM primers under a particular pathway between the organisms of interest. Furthermore, the page offers a simple comparison of the distribution of various SSR classes within the pathway of interest. All comparisons are offered in a simple and effective visualization style.

### Statistics and discussion

**PPMdb database statistics.** As of January 2021, the PPMdb was launched and consisted of 2.7 million pathway-based markers distributed over 9894 marker panels developed across 82 plant genomes. Across these genomes, 165,378 enzyme-coding genes were mapped against 126 KEGG reference pathway maps. The SSR mining of all enzyme-coding genes identified 3,471,782 SSR motifs, including 2,844,501 perfect motifs and



**Figure 4.** Screenshot of the PPMdb database “Map Browse” section including page layout and Pathway markers mapping layout.

627,281 compound motifs across the 82 plant genomes. Moreover, a total of 691,555 PGTMs and 1,999,187 PMTMs were developed (Fig. 6). This massive number of developed pathway-based markers have been mapped and supplemented with all essential information to offer users a modern version of markers called “pathway gene set markers”. All developed molecular markers for the 82 plant species are saved in separate backend tables for each plant. These datasets are searchable and can be downloaded conveniently via the PPMdb website. The statistical records of all analyzed plant genomes within the PPMdb, including the number of mapped pathways and designed primers, are summarized in Table 1.

**Insights.** Advances in plant omics promise to transform the molecular markers research area, in which the main challenge will not be the development of novel markers rather than the optimum selection and validation of a group of useful functional markers from the big collection of candidates<sup>32</sup>. For many decades, genetic diver-



JBrowse

**Arabidopsis thaliana**

**Plant Name**

- Arabidopsis lyrata
- Arabidopsis thaliana
- Arachis duranensis
- Arachis ipaensis
- Asparagus officinalis
- Beta vulgaris
- Brachypodium distachyon
- Brassica napus
- Brassica oleracea
- Brassica rapa
- Cajanus cajan
- Camelina sativa
- Capsella rubella
- Capsicum annuum
- Carica papaya
- Chenopodium quinoa
- Cicer arietinum
- Citrus sinensis

**Available Tracks**

filter tracks

refseq

Miscellaneous 3

- Microsatellites
- PGTM
- PMTM

Genome Track View Help

0 200 400 600 800 1,000 1,200 1,400 1,600 1,800

ath\_AT1G01060 ath\_AT1G01060:1345..1707

1,375 1,500 1,625

refseq

PMTM

Microsatellites

PGTM

**PMTM**

Primary Data

Type PMTM

Position ath\_AT1G01060:1142..1441

Length 300 bp

Attributes

forward\_primer\_gc% 50

forward\_primer\_sequence TGCAGTCCAACTCCAGCA

forward\_primer\_size 20

forward\_primer\_tm 60

reverse\_primer\_gc% 55

reverse\_primer\_sequence GTTTCTCGCTCACTCCCTGA

reverse\_primer\_size 20

reverse\_primer\_sequence GTTTCTCGCTCACTCCCTGA

reverse\_primer\_size 20

reverse\_primer\_tm 59

seq\_id ath\_AT1G01060

source insilco\_analysis

uniqueID offset-1639586

Region sequence

```
>ath_AT1G01060 ath_AT1G01060:1142..1441 class=PMTM
length=300
TTGCAGTCCAACTCCAGCACTGCAATGGATACGCGTTGAAATCACTCAACCGTTG
AGAAACAACACAGACTCTCGAAGATCAAACTGGCTTGGAAATCTCCAGCTCATCAT
CTGATGATTCAGATGAGACTGGAGTACCAAGCTAAATGCCGACTCAAAAACCAATGATG
ATAAAATGGAGGAGTTGTTGACTGCGCGCTGTGCATGACTCAAACTGCCAGGAAGA
AAAACTTGGGACCGCTCATCTGCTGCTCAAAATACACTTCAGGGAGTGACCCGAA
```

**JBrowse Section** enables users to visualize and map all identified SSR motifs besides all developed PMTM and PGTM markers against the Refseq CDS for the plant genome of interest

The search results containing essential descriptive information for the selected marker (e.g., marker type, repeat type & sequence, primer sequence, primer annealing temperature, primer position, product size, and other information)

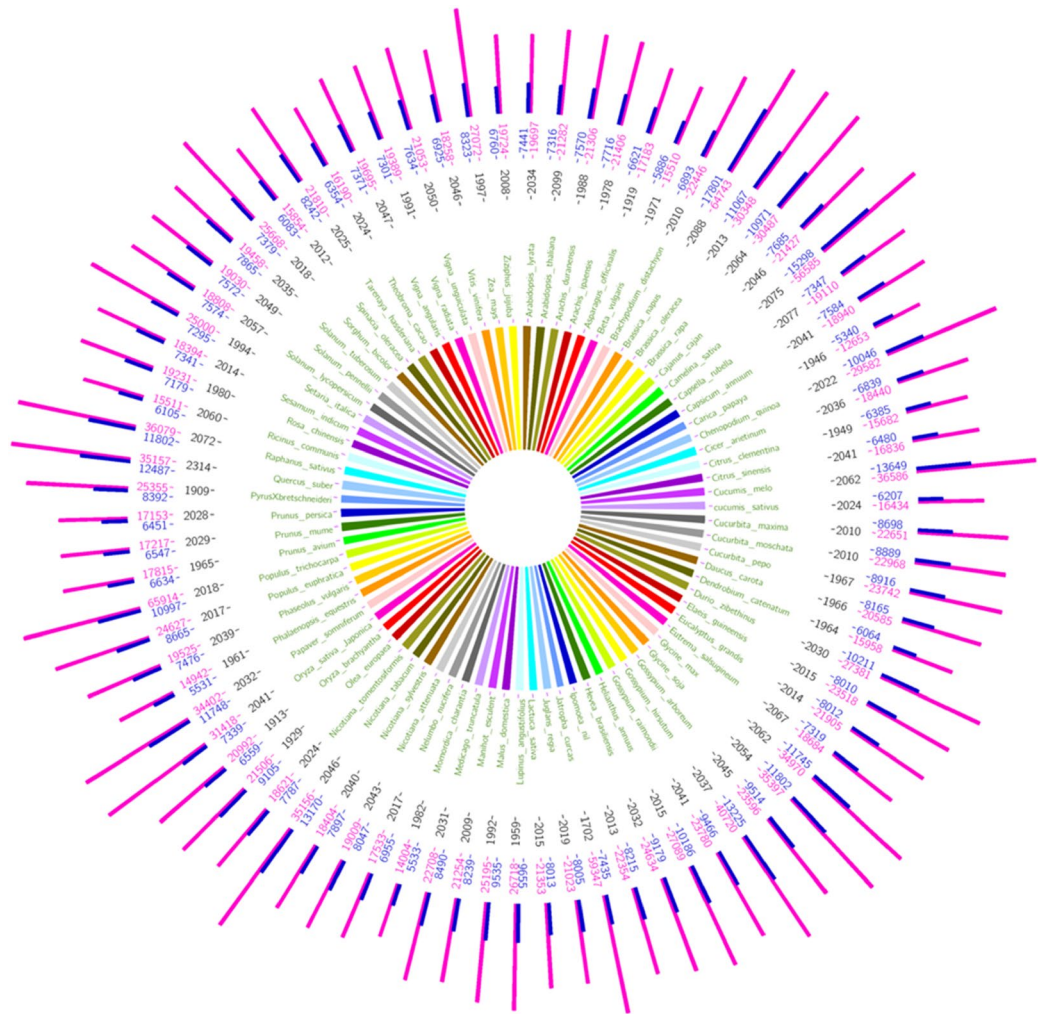
Figure 5. Screenshot of the PPMdb database “JBrowse” page.

ity and molecular breeding studies have involved a few pre-specified candidate markers/genes. This knowledge-based approach was found to run a high risk of missing critical genes related to interest traits.

Genes or gene families involved in a biological pathway are often switched on or off together to reflect a particular biological function or elucidate specific phenotypes<sup>1</sup>. In plants, complex economic traits such as yield, resistance to a particular disease, production of secondary metabolites, etc.,... have thus driven the scientists’ needs for new ‘systems’-based approaches that can illuminate the molecular mechanisms underlying specific trait(s) rather than the effect of distinct genes<sup>3</sup>. Based on this deep understanding, this study’s scientific vision provides a platform for the development of pathways-based markers toward designing future studies that aim to disentangle the causal biological pathways and their phenotypic reflections.

As a part of this vision, microsatellites have been utilized as one of the most common sources of genetic markers, which have served as a keystone for massive genetic studies due to their robust and unique features.





**Figure 6.** Statistical layout of the 82 analyzed plant genomes. The middle part shows the plants’ scientific names (each plant is assigned to a unique color). The black-colored numbers refer to the number of mapped enzymes within each plant genome, where the blue-colored numbers refer to the number of developed PGTM, and the purple-colored numbers refer to the number of developed PMTM.

Microsatellites, which demonstrate their efficiency in broad applications, such as diversity studies, genome mapping, molecular breeding, and molecular phylogeny<sup>33,34</sup>, were successfully implemented in our developed PGM pipeline to develop the PMTM system as a novel class of pathway-targeted markers. In silico microsatellite mining on a plant genome scale is expected to advance our understanding and elucidate the functional impact of microsatellites within biological pathways and, consequently, in the context of systems biology<sup>35</sup>. This study also aimed to develop and map pathway-based genetic marker panels that support and boost molecular breeding programs, genetic diversity, and genetic characterization studies. From this point of view, the value of our developed marker panels will need to be assessed in the context of availability and mining of biological pathway information to generate knowledge that is more actionable rather than more complex.

**Conclusion.** In summary, we present PPMdb as a comprehensive database for pathway-based markers in plant genomes. To our knowledge, PPMdb is the first portal providing unique pathway-targeted marker panels not presented in any previous database. PPMdb classifies developed marker panels based on their biological functions. This advantage allows researchers to deeply focus on the functional roles of utilized markers and may explain many phenotypic variations in the future. Our developed PPMdb is substantially different from similar plant genetic marker databases. The PPMdb will regularly be updated by integrating any newly released plant genomes.

Furthermore, the PPMdb graphical user interface and functionality will always be enhanced and continuously supported with new tools and technologies. Overall, we believe that the PPMdb hub will serve as a starting point or cornerstone for pathway-targeted marker research. In addition, we believe that PPMdb will catch great attention from a wide range of plant scientists in different disciplines, including genetic diversity, species characterization, population genetics, genome mapping, and targeted trait improvement.

Plant Name	No. of mapped pathways	No. of designed primers	Plant Name	No. of mapped pathways	No. of designed primers
<i>Arabidopsis lyrata</i>	121	27,138	<i>Lupinus angustifolius</i>	121	36,373
<i>Arabidopsis thaliana</i>	120	28,598	<i>Malus domestica</i>	121	34,730
<i>Arachis duranensis</i>	121	28,876	<i>Manihot esculent</i>	121	29,493
<i>Arachis ipaensis</i>	121	29,122	<i>Medicago truncatula</i>	120	31,198
<i>Asparagus officinalis</i>	120	23,804	<i>Momordica charantia</i>	120	19,537
<i>Beta vulgaris</i>	120	21,396	<i>Nelumbo nucifera</i>	120	24,488
<i>Brachypodium distachyon</i>	121	29,339	<i>Nicotiana attenuata</i>	121	27,056
<i>Brassica napus</i>	121	82,544	<i>Nicotiana sylvestris</i>	120	26,301
<i>Brassica oleracea</i>	121	41,415	<i>Nicotiana tabacum</i>	120	48,326
<i>Brassica rapa</i>	121	41,458	<i>Nicotiana tomentosiformis</i>	120	26,408
<i>Cajanus cajan</i>	121	29,112	<i>Olea europaea</i>	121	30,611
<i>Camelina sativa</i>	121	71,883	<i>Oryza brachyantha</i>	121	27,551
<i>Capsella rubella</i>	121	26,457	<i>Oryza sativa Japonica</i>	121	38,757
<i>Capsicum annuum</i>	121	26,524	<i>Papaver somniferum</i>	120	46,150
<i>Carica papaya</i>	120	17,993	<i>Phalaenopsis equestris</i>	119	20,473
<i>Chenopodium quinoa</i>	120	39,628	<i>Phaseolus vulgaris</i>	121	27,001
<i>Cicer arietinum</i>	121	25,279	<i>Populus euphratica</i>	121	33,292
<i>Citrus clementina</i>	121	22,067	<i>Populus trichocarpa</i>	121	76,911
<i>Citrus sinensis</i>	121	23,316	<i>Prunus avium</i>	121	24,449
<i>Cucumis melo</i>	120	50,235	<i>Prunus mume</i>	121	23,764
<i>cucumis sativus</i>	120	22,641	<i>Prunus persica</i>	121	23,604
<i>Cucurbita maxima</i>	120	31,349	<i>PyrusXbretschneideri</i>	121	33,747
<i>Cucurbita moschata</i>	120	31,857	<i>Quercus suber</i>	121	47,644
<i>Cucurbita pepo</i>	120	32,658	<i>Raphanus sativus</i>	121	47,881
<i>Daucus carota</i>	121	28,750	<i>Ricinus communis</i>	121	21,616
<i>Dendrobium catenatum</i>	120	22,022	<i>Rosa chinensis</i>	121	26,410
<i>Durio zibethinus</i>	121	37,592	<i>Sesamum indicum</i>	121	25,735
<i>Elaeis guineensis</i>	120	31,528	<i>Setaria italica</i>	121	32,295
<i>Eucalyptus grandis</i>	121	29,917	<i>Solanum lycopersicum</i>	120	26,382
<i>Eutrema salsugineum</i>	121	26,003	<i>Solanum pennellii</i>	121	26,602
<i>Glycine max</i>	121	46,715	<i>Solanum tuberosum</i>	120	27,323
<i>Glycine soja</i>	121	47,199	<i>Sorghum bicolor</i>	121	33,047
<i>Gossypium arboreum</i>	121	33,110	<i>Spinacia oleracea</i>	120	21,937
<i>Gossypium hirsutum</i>	121	53,945	<i>Tarenaya hassleriana</i>	121	30,052
<i>Gossypium raimondii</i>	121	33,246	<i>Theobroma cacao</i>	121	22,544
<i>Helianthus annuus</i>	121	37,275	<i>Vigna angularis</i>	121	27,066
<i>Hevea brasiliensis</i>	121	33,813	<i>Vigna radiata</i>	121	26,690
<i>Ipomoea nil</i>	121	30,569	<i>Vigna unguiculata</i>	121	28,687
<i>Jatropha curcas</i>	118	66,782	<i>Vitis vinifera</i>	121	25,183
<i>Juglans regia</i>	120	29,028	<i>Zea mays</i>	121	35,395
<i>Lactuca sativa</i>	121	29,366	<i>Ziziphus jujuba</i>	121	26,484

**Table 1.** Summarize the number of mapped pathways and the number of designed primers of each plant genome.

## Data availability

PPMdb is an online free access database initiative available at the following link: (<http://PPMdb.easyomics.org>).

Received: 10 March 2021; Accepted: 6 September 2021

Published online: 29 October 2021

## References

1. Kao, P. Y., Leung, K. H., Chan, L. W., Yip, S. P. & Yap, M. K. Pathway analysis of complex diseases for GWAS, extending to consider rare variants, multi-omics and interactions. *Biochim. Biophys. Acta (BBA) Gen. Subj.* **1861**, 335–353 (2017).

2. Golestan, H. F. S. *et al.* Intelligent mining of large-scale bio-data: Bioinformatics applications. *Biotechnol. Biotechnol. Equip.* **32**, 10–29 (2018).
3. Kanehisa, M. KEGG bioinformatics resource for plant genomics and metabolomics. In *Plant Bioinformatics* (ed. Edwards, D.) 55–70 (Humana Press, 2016).
4. Dharmawardhana, P. *et al.* A genome scale metabolic network for rice and accompanying analysis of tryptophan, auxin and serotonin biosynthesis regulation under biotic stress. *Rice* **6**, 15 (2013).
5. Lee, H. & Shin, M. Mining pathway associations for disease-related pathway activity analysis based on gene expression and methylation data. *BioData Min.* **10**, 1–3 (2017).
6. Bezerianos, A., Dragomir, A. & Balomenos, P. *Computational Methods for Processing and Analysis of Biological Pathways* (Springer, 2017).
7. Mandal, M., Mondal, J. & Mukhopadhyay, A. A PSO-based approach for pathway marker identification from gene expression data. *IEEE Trans. NanoBiosci.* **14**, 591–597 (2015).
8. Allen, L. A. & Felker, G. M. Multi-marker strategies in heart failure: Clinical and statistical approaches. *Heart Fail. Rev.* **15**, 343–349 (2010).
9. Saeys, Y., Inza, I. & Larraaga, P. A review of feature selection techniques in bioinformatics. *Bioinformatics* **23**, 2507–2517 (2007).
10. Mandal, M. & Mukhopadhyay, A. A graph-theoretic approach for identifying non-redundant and relevant gene markers from microarray data using multiobjective binary PSO. *PLoS ONE* **9**, e90949 (2014).
11. Bandyopadhyay, S., Mallik, S. & Mukhopadhyay, A. A survey and comparative study of statistical tests for identifying differential expression from microarray data. *IEEE/ACM Trans. Comput. Biol. Bioinform.* **11**, 95–115 (2013).
12. Lee, E., Chuang, H. Y., Kim, J. W., Ideker, T. & Lee, D. Inferring pathway activity toward precise disease classification. *PLoS Comput. Biol.* **4**, e1000217 (2008).
13. Emmert-Streib, F. & Glazko, G. V. Pathway analysis of expression data: Deciphering functional building blocks of complex diseases. *PLoS Comput. Biol.* **7**, e1002053 (2011).
14. Khatri, P., Sirota, M. & Butte, A. J. Ten years of pathway analysis: Current approaches and outstanding challenges. *PLoS Comput. Biol.* **8**, e1002375–e1002381 (2012).
15. Gu, Z., Liu, J., Cao, K., Zhang, J. & Wang, J. Centrality-based pathway enrichment: A systematic approach for finding significant pathways dominated by key genes. *BMC Syst. Biol.* **6**, 1–13 (2012).
16. Monaco, M. *et al.* Maize metabolic network construction and transcriptome analysis. *Plant Genome* **6**, 1–12 (2013).
17. Kanehisa, M. & Goto, S. KEGG: Kyoto encyclopedia of genes and genomes. *Nucleic Acids Res.* **28**, 27–30 (2000).
18. Gupta, P. *et al.* Gramene database: Navigating plant comparative genomics resources. *Curr. Plant Biol.* **7**, 10–15 (2016).
19. Naithani, S. *et al.* Plant Reactome: A resource for plant pathways and comparative analysis. *Nucleic Acids Res.* **44**, gkw932 (2016).
20. Thimm, O. *et al.* MAPMAN: A user-driven tool to display genomics data sets onto diagrams of metabolic pathways and other biological processes. *Plant J.* **37**, 914–939 (2004).
21. Caspi, R. *et al.* The MetaCyc database of metabolic pathways and enzymes and the BioCyc collection of pathway/genome databases. *Nucleic Acids Res.* **42**, D459–D471 (2014).
22. Zhang, P. *et al.* Creation of a genome-wide metabolic pathway database for *Populus trichocarpa* using a new approach for reconstruction and curation of metabolic pathways for plants. *Plant Physiol.* **153**, 1479–1491 (2010).
23. Karp, P. D. *et al.* The BioCyc collection of microbial genomes and metabolic pathways. *Brief. Bioinform.* **20**, 1085–1093 (2019).
24. Mokhtar, M. M., Adawy, S. S., El-Assal, S. E. & Hussein, E. H. Genic and intergenic SSR database generation, SNPs determination and pathway annotations, in date palm (*Phoenix dactylifera* L.). *PLoS ONE* **11**, e0159268 (2016).
25. Atia, M. A., Osman, G. H. & Elmenofy, W. H. Genome-wide *in silico* analysis, characterization and identification of microsatellites in *Spodoptera littoralis* multiple nucleopolyhedrovirus (SpliMNPV). *Sci. Rep.* **6**, 1–9 (2016).
26. Mokhtar, M. M. & Atia, M. A. SSRome: An integrated database and pipelines for exploring microsatellites in all organisms. *Nucleic Acids Res.* **47**, D244–D252 (2019).
27. Beier, S., Thiel, T., Münch, T., Scholz, U. & Mascher, M. MISA-web: A web server for microsatellite prediction. *Bioinformatics* **33**, 2583–2585 (2017).
28. Untergasser, A. *et al.* Primer3-new capabilities and interfaces. *Nucleic Acids Res.* **40**, 1–12 (2012).
29. Mokhtar, M. M., Hussein, E. H., El-Assal, S. E. & Atia, M. A. VFODB: A comprehensive database of ESTs, EST-SSRs, mtSSRs, microRNA-target markers and genetic maps in *Vicia faba*. *AoB Plants* **12**, plaa064 (2020).
30. <https://github.com/jamietre/imagemapster>. Accessed Jun 2019.
31. Buels, R. *et al.* JBrowse: A dynamic web platform for genome visualization and analysis. *Genome Biol.* **17**, 1–2 (2016).
32. Chaudhary, J. *et al.* Advances in omics approaches for abiotic stress tolerance in tomato. *Biology* **8**, 90 (2019).
33. Esselink, G. D., Nybom, H. & Vosman, B. Assignment of allelic configuration in polyploids using the MAC-PR (microsatellite DNA allele counting—peak ratios) method. *Theor. Appl. Genet.* **109**, 402–408 (2004).
34. Kalia, R. K., Rai, M. K., Kalia, S., Singh, R. & Dhawan, A. K. Microsatellite markers: An overview of the recent progress in plants. *Euphytica* **177**, 309–334 (2011).
35. Quilez, J. *et al.* Polymorphic tandem repeats within gene promoters act as modifiers of gene expression and DNA methylation in humans. *Nucleic Acids Res.* **44**, 3750–3762 (2016).

## Acknowledgements

The authors would like to thank the administration of the Agricultural Genetic Engineering Research Institute (AGERI), as well as the administration of the Agricultural Research Center (ARC), for their continued support. The authors also acknowledge the African Supercomputing Center at Mohamed VI Polytechnic University for providing necessary supercomputing resources to achieve this study.

## Author contributions

Conceptualization, M.A.M.A.; Methodology, M.M., M.A.M.A.; Software, M.M., M.A.M.A., A.E., M.F.H.; Resources, M.M., M.A.M.A., A.E., M.F.H.; Formal Analysis, M.M., A.E.; Visualization, M.M., M.A.M.A., A.E., M.F.H.; Supervision, M.A.M.A.; writing—original draft preparation, M.A.M.A.; Writing—Review and Editing, M.A.M.A., A.E., M.M., M.F.H. All authors have read and agreed to the published version of the manuscript.

## Funding

The author(s) received no specific funding for this work.

## Competing interests

The authors declare no competing interests.

### Additional information

**Correspondence** and requests for materials should be addressed to A.E.A. or M.A.M.A.

**Reprints and permissions information** is available at [www.nature.com/reprints](http://www.nature.com/reprints).

**Publisher's note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



**Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2021