

1 **SciApps: A Cloud-Based Platform for Analyses and Distribution of the MaizeCODE data**

2

3 Liya Wang¹, Zhenyuan Lu¹, Melissa delaBastide¹, Peter Van Buren¹, Xiaofei Wang¹, Cornel
4 Ghiban¹, Michael Regulski¹, Jorg Drenkow¹, Xiaosa Xu¹, Carlos Ortiz Ramirez², Cristina F.
5 Marco¹, Jason Williams¹, Alexander Dobin¹, Kenneth D. Birnbaum², David P. Jackson¹, Robert
6 A. Martienssen¹, William R. McCombie¹, David A. Micklos¹, Michael C. Schatz^{1,3}, Doreen H.
7 Ware^{1,4}, Thomas R. Gingeras¹

8

9 ¹Cold Spring Harbor Laboratory, Cold Spring Harbor, NY, ²New York University, New York, NY,
10 ³Johns Hopkins University, Baltimore, MD. ⁴USDA-ARS Robert W. Holley Center for Agriculture
11 and Health, Ithaca, NY, United States

12

13 MaizeCODE is a project aimed at identifying and analyzing functional elements in the maize
14 genome. In its initial phase, MaizeCODE assayed up to five tissues from four maize strains
15 (B73, NC350, W22, TIL11) by RNA-Seq, Chip-Seq, RAMPAGE, and small RNA sequencing. To
16 facilitate reproducible science and provide both human and machine access to the MaizeCODE
17 data, we developed SciApps, a cloud-based portal, for analysis and distribution of both raw data
18 and analysis results. Based on the SciApps workflow platform, we generated new components
19 to support the complete cycle of the MaizeCODE data management. These include publicly
20 accessible scientific workflows for reproducible and shareable analysis of various functional
21 data, a RESTful API for batch processing and distribution of data and metadata, a searchable
22 data page that lists each MaizeCODE experiment as a reproducible workflow, and integrated
23 JBrowse genome browser tracks linked with workflows and metadata. The SciApps portal is a
24 flexible platform that allows integration of new analysis tools, workflows, and genomic data from
25 multiple projects. Through metadata and a ready-to-compute cloud-based platform, the portal
26 experience improves access to the MaizeCODE data and facilitates its analysis.

27

28 **Keywords:** bioinformatics, functional annotations, cloud computing, ENCODE, workflows

1 INTRODUCTION

2

3 Maize is one of the most biologically, socially, and economically important crop plants. Following
4 the sequencing of its genome (Jiao et al., 2017; Schnable et al., 2009), the next critical step in
5 understanding maize biology will involve identifying and deciphering functional sequence
6 regions. Modeled on the Encyclopedia of DNA Elements (ENCODE) Project for the human
7 genome (ENCODE Consortium, 2004), the MaizeCODE project is an integrated and multi-
8 disciplinary project aimed at revealing the functional regions of the maize genome by identifying
9 loci that are transcribed, methylated, or bound by specific modified histones and transcription
10 factors in various tissues. In addition, MaizeCODE is designed to store, collate, display, and
11 disseminate the data to the wider community of plant biologists worldwide.

12

13 To curate, process, and distribute the ENCODE data, the ENCODE Data Coordination Center
14 (DCC) group established the ENCODE portal (Sloan et al., 2016), which relies on both rich
15 metadata and commercial cloud resources through the DNAnexus platform
16 (<https://www.dnanexus.com>). Within this platform, standard processing pipelines for human
17 genome analysis are constructed to ensure consistent and reproducible processing of primary
18 sequence data. However, both the ENCODE DCC and end users are required to cover the cost
19 of the DNAnexus service and commercial cloud resources. In order to provide a cost-free data
20 processing platform for academic users, the MaizeCODE DCC group decided to leverage two
21 NSF-funded resources, XSEDE (Towns et al., 2014) at the Texas Advanced Computing Center
22 (TACC) for computing power and the CyVerse Data Store (Goff et al., 2011) for cloud-based
23 data storage.

24

25 To automate bioinformatics analysis over both the XSEDE/TACC cloud and CyVerse Data
26 Store, we developed a bioinformatics workflow platform called SciApps (Wang et al., 2018). In

1 the work described here, we further improved SciApps by adding a RESTful API for automating
2 batch processing of the MaizeCODE data and metadata management, a searchable
3 MaizeCODE data page powered by a relational database, several analysis workflows, and
4 Genome Browser tracks automatically generated from unique workflow identifiers via the
5 RESTful API. The SciApps platform has been used to support both MaizeCODE DCC and end
6 users to process/reprocess and manage multi-omics data through either the GUI or the API.

7

8 **METHODS**

9

10 **Overview of the entire MaizeCODE data management cycle**

11

12 To improve accessibility, reproducibility, reusability, and interoperability, data generated by the
13 MaizeCODE Consortium members are uploaded to a cloud-based data storage system, the
14 CyVerse Data Store (Goff et al., 2011). The Data Store, which is built on top of iRODS (a rule-
15 oriented data system) ([Moore and Rajsekar 2010](#)), supports data virtualization, sharing, bulk
16 uploads/downloads, and collaborations. Once uploaded, the experimental metadata are
17 attached to raw data files to facilitate the reuse of data, as well as submission of data to the
18 NCBI Short Read Archive (SRA). The metadata is later retrieved via the Terrain API
19 (<https://github.com/cyverse-de/terrain>) to automate batch analyses via SciApps (Wang et al.,
20 2018). SciApps provides a ready-to-compute cloud-based platform for automating complex
21 analyses constructed using modular applications (or apps). Previously, SciApps was operated
22 via a web GUI, but since that time we have developed SciApps RESTful APIs to support batch
23 processing of MaizeCODE and other data. In addition, we have integrated over 20 new apps for
24 ground-level analyses such as quality control (QC), alignment to the reference genome, filtering,
25 quantification (e.g. for gene expression), and peak calling (if needed). Replicates (and controls if
26 available) of each assay are organized as a single experiment (or workflow with a unique ID),

1 which represents an entity that chains together raw data, analysis results, experimental
2 metadata, and computational provenance. SciApps extracts experimental metadata and
3 attaches them to a specific workflow so that users can access them directly on the SciApps
4 portal. Genome browser tracks are automatically generated and displayed within an integrated
5 version of JBrowse (Skinner et al., 2009) by looping through the list of experiments/workflows
6 via the SciApps RESTful API.

7
8 Together, the SciApps portal supports the complete cycle of MaizeCODE data management,
9 and the level of automation it provides greatly decreases the chances of human errors in data
10 organization, analysis, and distribution.

11

12 **Processing and accessing the MaizeCODE data with the SciApps RESTful API**

13

14 The cloud-based architecture of SciApps (Wang et al., 2018) enables highly scalable processing
15 of MaizeCODE data on the XSEDE/TACC cloud. Both intermediate and final results are
16 archived in the CyVerse Data Store, where the raw data are also hosted. As discussed above,
17 each SciApps workflow captures experimental metadata and computational provenance along
18 with the raw and processed data. Batch processing is supported through the RESTful API; the
19 API endpoints are provided in **Table 1**.

20

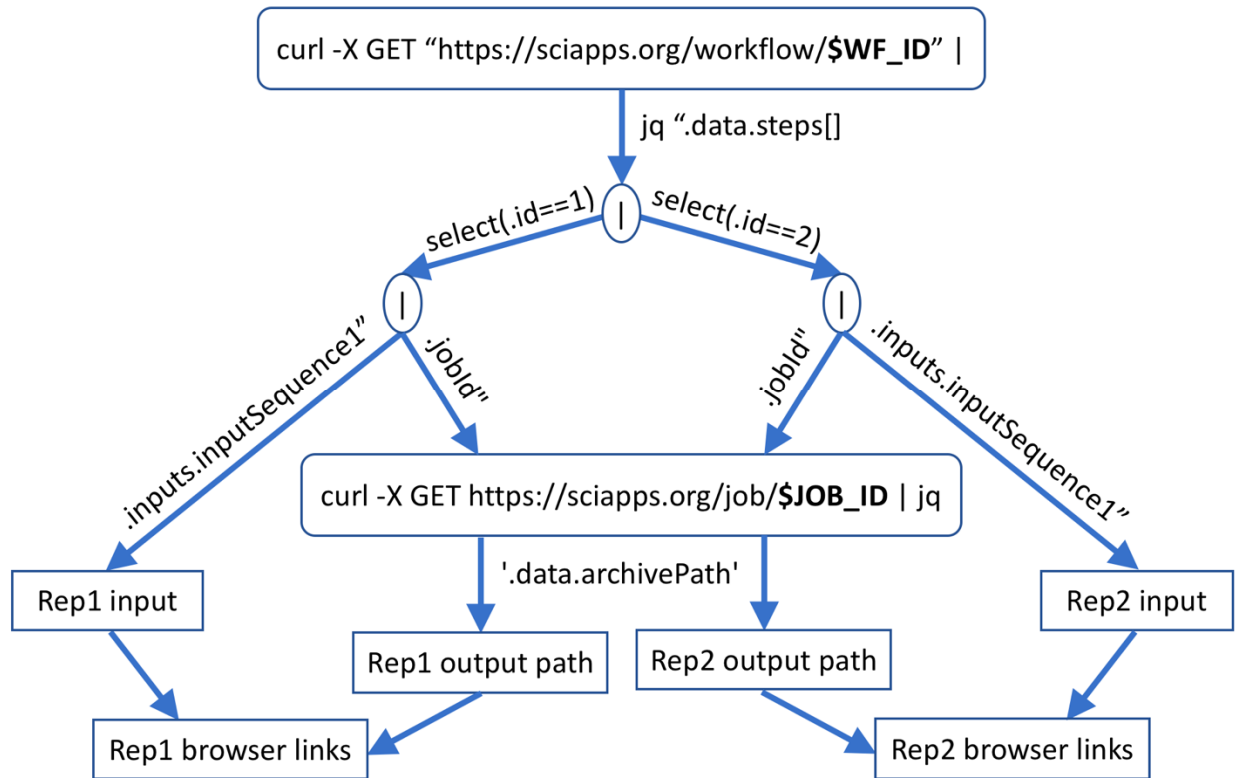
21 **Table 1.** SciApps release 1.0 RESTful API

Endpoint	Method	Description
/job	GET	List all jobs
/job/new/{id}	POST	Run a new job
/workflow/build	POST	Build a workflow from jobs

/workflowJob/new	POST	Generate a workflow JSON
/workflow/new	POST	Save a new workflow
/job/{id}	GET	Return the job JSON
/job/{id}/delete	GET	Delete the job
/workflowJob/run/{id}	GET	Run a new workflow
/workflow/{id}/metadata	GET	Get the workflow metadata
/workflow/{id}/update	POST	Update the workflow with metadata etc
/workflow	GET	List all workflows
/apps/{id}	GET	Return the application JSON
/workflow/{id}/delete	GET	Delete the workflow
/workflow/{id}	GET	Return the workflow JSON
/apps	GET	List all integrated apps

1
2 The analysis workflow for a specific assay is typically built interactively within the SciApps GUI
3 using one data set as a template. Once the workflow is captured, it can then be easily and
4 automatically applied to analyze other genomes and tissues. Alternatively, users may also build
5 workflows entirely programmatically with a series of analysis job IDs via the API. Experimental
6 metadata are retrieved via the CyVerse Terrain API, and then attached to the workflow via the
7 SciApps API at runtime. The API supports the MaizeCODE DCC for automatically processing a
8 large amount of data and also supports retrieval of results and metadata by end users. For
9 example, genome browser tracks can be automatically generated given a workflow ID by the
10 following steps (**Figure 1**): 1. Retrieve job IDs and inputs with the workflow endpoint, given a
11 workflow ID; 2. Retrieve the output path with the job endpoint, given a job ID; 3. Construct the
12 browser-ready link with the retrieved information. To simplify the process, the MaizeCODE DCC
13 encodes the genome, tissue, and replicate information into the input raw data file path, which is
14 also accessible through the workflow metadata endpoints. SciApps also names the output

1 filename based on the input filename with the output ID (defined by the app) as the prefix. As
2 shown in **Figure 1**, once the input filename, input path, and output path to cloud storage are
3 retrieved by calling the API, the output file path can be constructed to build the browser links.
4



5
6

7 **Figure 1.** Flowchart of generating genome browser tracks via cURL for a MaizeCODE workflow.
8 Each workflow has two replicates. **\$WF_ID** is the workflow ID, and **\$JOB_ID** is the job ID of a
9 step of the workflow. As an example, the cURL for retrieving the input filename of replicate 1 is:
10 `curl -X GET --header "Accept: application/json" "https://www.sciapps.org/workflow/$WF_ID" | jq`
11 `".data.steps[] | select(.id==1) | .inputs.inputSequence1".`

12

1 Accessing the MaizeCODE experiments as reproducible workflows

2 The MaizeCODE data page can be accessed under 'Data' from the top navigation bar of
3 SciApps (**Figure 2**). Keyword search is supported to allow the user to narrow down the list of
4 experiments to a specific genome or tissue or assay in real time. Once an experiment is
5 selected, the user can access the metadata, workflows, and ground-level analysis results of the
6 experiments, starting from raw sequence data. With the 'Relaunch' tab, user can reproduce the
7 entire analysis with one click or apply the same analysis workflow to new data. Using the 'Share'
8 tab, the analysis can be shared with others. Users can load the results to the History panel and
9 subject them to further analysis using the modular apps. Because all results are archived in the
10 cloud, downstream analyses can be completed quickly, e.g., differential expression analysis
11 between two tissues can be completed in a few minutes, rather than hours when starting from
12 the raw sequence data.

13

The screenshot displays the MaizeCODE web interface. At the top, there is a navigation bar with 'Home', 'Data', 'Workflow', 'Tools', 'Help', and 'Login'. Below this, a search bar is labeled 'Keyword search'. The main content area is divided into three panels:

- Left Panel (Modular apps):** A sidebar with a 'Search Apps' field and a list of categories: Alignment, Annotation, Assembly, Calculation, Clustering, Comparison, Conversion, Data handling, and Mapping. Under 'Alignment', several tools are listed, including BWA_index-0.7.17, BWA_mem-0.7.17, Bowtie2-2.3.2, HISAT2_align-2.0.5, HISAT2_index-2.0.5, STAR_align-2.5.3, and STAR_index-2.5.3.
- Middle Panel (List of workflows/experiments):** A table with columns 'Name' and 'Description'. Below the table, a 'Metadata' panel is open, showing details for a selected experiment (MC_B73_cr). The metadata includes: age (5 DAG), assay (Long RNA-Seq (LG)), cultivar (B73), design_description (After total RNA isolation, 1 ug was used for the library prep...), forward_read_length (150), fragment_size (523), growth_protocol (Growth chamber: Maize growth conditions for tissue collection), instrument_model (NextSeq 500), library_input (Total RNA), and library_layout (Paired).
- Right Panel (History):** A panel titled 'History' showing a list of jobs. The first job is '1: MCma-0.0.1' with a list of files: multiloq_report.html, rsem_cn_rep1_R1.txt, sig_l_cn_rep1_R1.txt, sig_l_cn_rep1_R1.bw, star_cn_rep1_R1.bam, star_cn_rep1_R1.bam.bai, and str_cn_rep1_R1.gif. Below it, '2: MCma-0.0.1' and '3: MCma-0.0.1' are also listed with their respective file lists.

Red arrows point to various features: 'Reproduce the analysis' points to the 'Relaunch' button; 'Load results' points to the 'Load' button; 'Keyword search' points to the search bar; 'Visualization' points to the 'Visualize' button; and 'Loaded results' points to the 'Load' button in the history panel. A 'Modular apps' label points to the left sidebar.

14

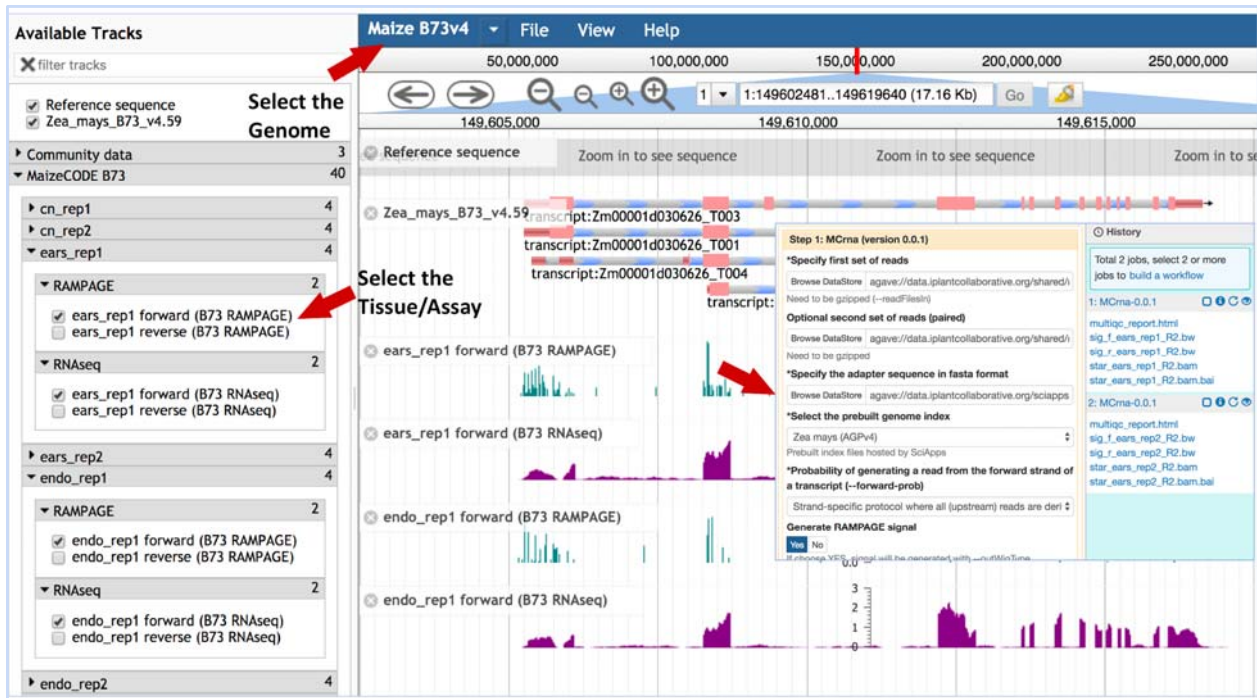
15

16 **Figure 2.** Web browser interface of the MaizeCODE data page. In the middle panel page, a list
17 of workflows/experiments is presented. Above the list, several action buttons are available:

1 'Relaunch' the analysis, 'Visualize' the graphic diagram of the workflow (with URLs for the raw
2 sequence files from the input file node), 'Load' the results to the History panel, 'Share' the
3 analysis with others, and display the experimental 'Metadata'. User can perform a keyword
4 search for a specific dataset (e.g., B73 ears RNA-Seq). In the right panel, SciApps displays the
5 history of the selected datasets; the visualization (eye) icon opens a panel where users can
6 generate links to visualize the results in a web browser (e.g., a QC report) or genome browser
7 (e.g., alignments or signal tracks). The left panel shows a list of modular apps that can be
8 launched to perform a variety of downstream analyses with the loaded results.

9 **Accessing the MaizeCODE data as Genome Browser tracks**

10 Once the analysis is completed, genome browser tracks are automatically generated given the
11 workflow ID by calling the SciApps API for an integrated version of JBrowse. The browser tracks
12 can be accessed under the 'Tools' menu within the top navigation bar. As shown in **Figure 3**,
13 tracks are organized by genome, tissue, replicate, and assay. Checking the box next to each
14 track will load it into the browser. The SciApps workflow ID is embedded, so clicking on a track
15 brings up the workflow 'Relaunch' interface, which can be used to reproduce the track signal if
16 needed. In this interface, the user can also check the parameters used for the analysis, as well
17 as additional results in the History panel. At the bottom of the interface, a diagram button
18 visualizes the workflow diagram, and a metadata button displays the experimental metadata
19 associated with the workflow. From the results, user can also generate additional browser track
20 links through the visualization (eye) icon. For example, this can be used to verify the signal track
21 with the alignment files (in the BAM format). As mentioned earlier, the results can also be used
22 to perform a downstream analysis on the same interface. Finally, the browser tracks are
23 available as a JSON file for integration into other platforms (e.g., the JSON file for B73 is
24 available at https://data.sciapps.org/view2/data2/B73/v4/apollo_data/trackList.json.



2 **Figure 3.** Genome browser tracks for the MaizeCODE data. JBrowse is used to hold the
3 MaizeCODE signal tracks, which are organized in the following order: genome, tissue, replicate,
4 and assay. Clicking on each track brings up the workflow 'Relaunch' interface.

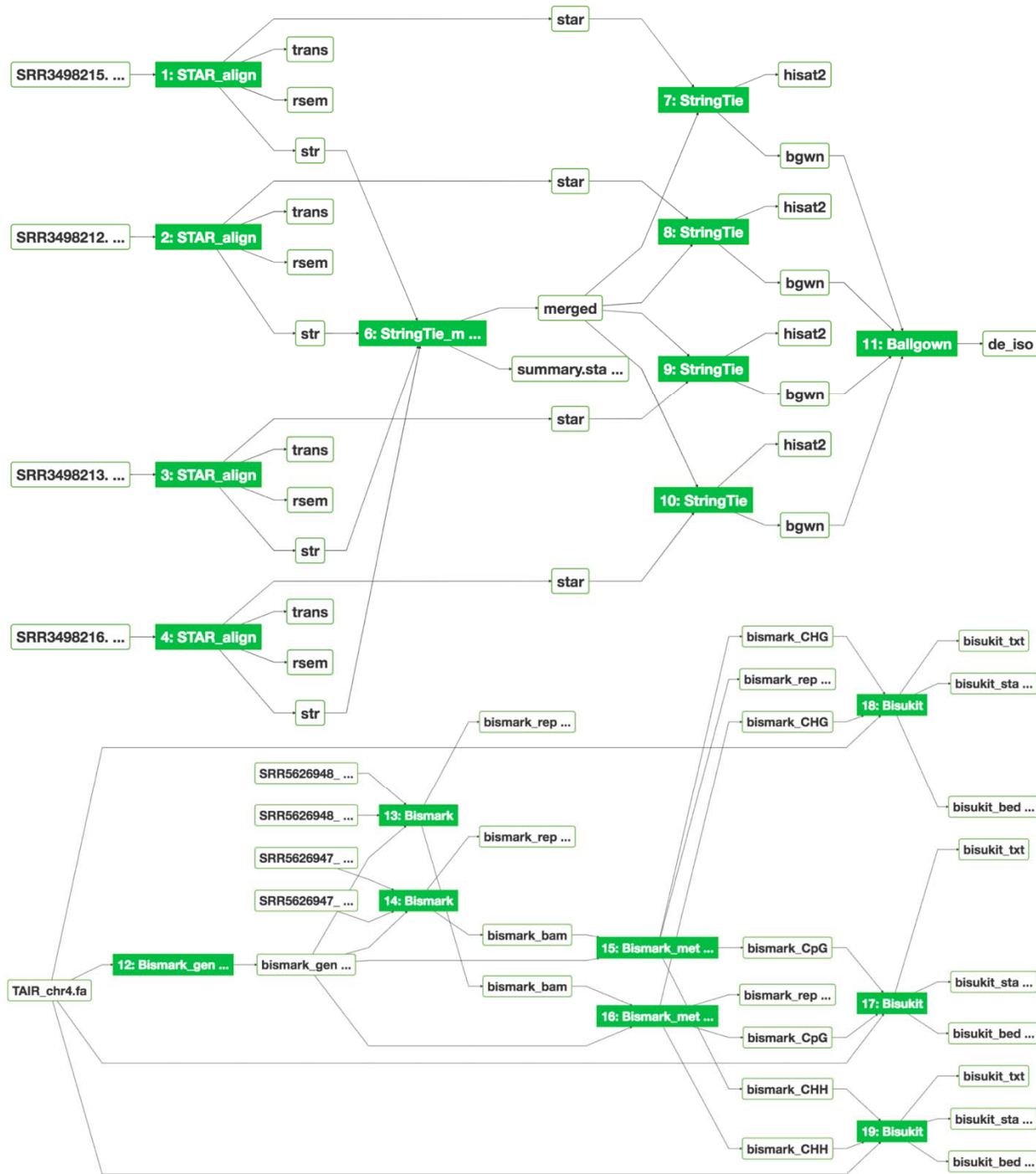
5 **Accessing the raw reads on CyVerse Data Store**

6 The raw sequence data is deposited into the CyVerse Data Store via iCommands
7 (<https://docs.irods.org/4.2.1/icommands/user/>), with metadata attached before submission to the
8 NCBI SRA. From there, users can access the raw data in several ways. Within SciApps, the
9 input file node of the graphic diagram for a workflow/experiment is linked to the raw sequence
10 file. Clicking on the input node will open the CyVerse Data Common landing page in a web
11 browser. The metadata attached to the raw sequence file is also displayed on the same page.
12 The user can further navigate through all released raw data from the landing page
13 (<http://datacommons.cyverse.org/browse/iplant/home/shared/maizecode/released/>); the
14 SciApps workflow ID is attached as metadata to the raw data files if it has been processed. The

1 user can use the ID to load the workflow on the SciApps portal. For batch downloading of raw
2 sequence files through the GUI or the command line, we recommend CyberDuck
3 (<https://cyberduck.io/>) or iCommands, respectively.

4 **Analysis with reproducible workflows**

5 Bioinformatics applications (or apps) are integrated into SciApps as modular components that
6 can be chained with other apps into an automated workflow. Individual apps are built with
7 Singularity images (Kurtzer et al., 2017) from BioConda recipes (Grüning et al., 2018) or directly
8 from Dockerfiles to ensure reproducibility across different cloud resources. To support the
9 analysis of MaizeCODE data, over 20 software tools are integrated. **Figure 4** shows two
10 publicly accessible workflows for differential expression analysis and cytosine methylation
11 analysis, building on the popular STAR (Dobin et al., 2013)/RSEM (Li and Dewey,
12 2011)/StringTie (Pertea et al., 2015) and Bismark (Krueger and Andrews, 2011) pipelines,
13 respectively. These workflows can be constructed either with the SciApps GUI or through the
14 API. The user can retrieve the inputs, metadata, results, and provenance of the software used in
15 the analysis with a unique workflow ID. The interactive graph, along with the platform guide
16 (<https://cyverse-sciapps-guide.readthedocs-hosted.com/en/latest/index.html>), helps users to
17 understand how multiple apps are used together to analyze a specific assay. For MaizeCODE
18 data, the graph is also helpful for visually inspecting the input–output relationship. Additionally,
19 the user can check the input data (through the input file node) and relaunch each individual step
20 of the analysis, or even the entire analysis, via the web interface or API.



1

2 **Figure 4.** Graphical workflow diagrams for differential expression analysis (top) and MethylC-
 3 seq analysis (bottom). The interactive graph demonstrates the relationships among input–output
 4 files, displays provenance of the software tools, and provides real-time job status updates of

1 new analyses via the color of the app node (green: completed; blue: running; yellow: pending;
2 red: failed).

3

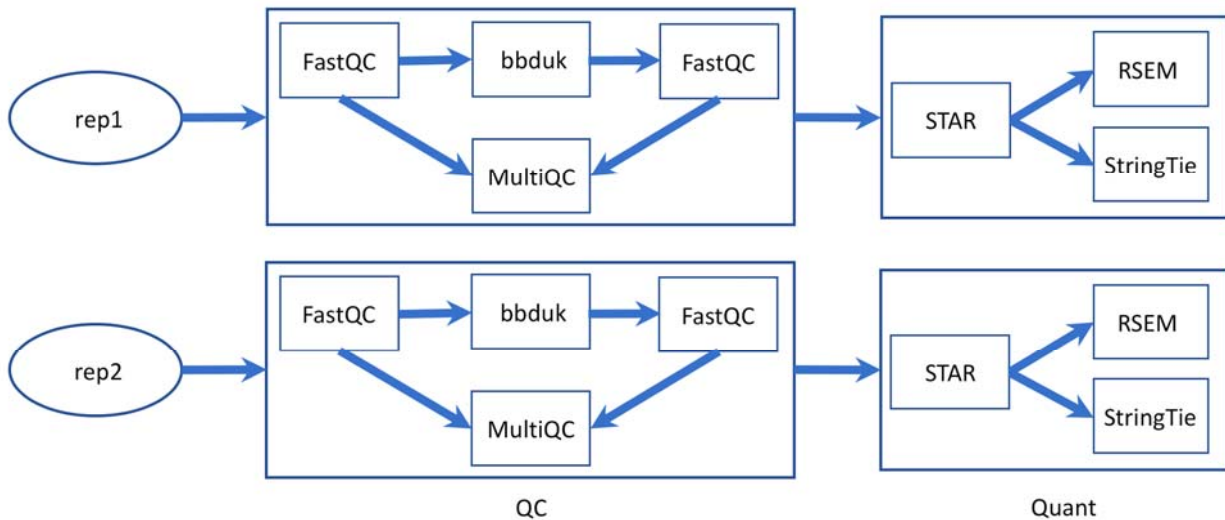
4 **RESULTS AND DISCUSSION**

5 A large variety of software is needed to process the MaizeCODE data. For each experiment
6 (consisting of two replicates), a workflow with a unique ID is provided via the SciApps platform.
7 One major goal of SciApps is to empower anyone in the community to easily repeat an entire
8 analysis, or use a workflow with alternative parameters for each step if so desired. A second
9 major goal is to empower community members to process their own comparable data sets using
10 protocols validated by MaizeCode. In the following sections, we describe how RNA-seq data is
11 processed, how the results can be visualized, and how the primary analysis results can be used
12 for differential expression analysis.

13 **Processing the RNA-seq data**

14 Besides the UCSC genome analysis tools (for format conversion and generating browser track
15 signals), the major software used in MaizeCODE RNA-seq data analysis are bbduk
16 (<https://jgi.doe.gov/data-and-tools/bbtools/bb-tools-user-guide/bbduk-guide/>) for trimming low
17 quality reads and adapter sequences, FastQC and MultiQC (Ewels et al., 2016) for visually
18 checking read quality, STAR (Dobin et al., 2013) for alignment, RSEM (Li and Dewey, 2011) for
19 quantifying gene expression, and StringTie (Pertea et al., 2015) for transcriptome assembly. All
20 tools are integrated into SciApps individually as separate apps, and also combined as a single
21 app, MCrna, for rapid batch processing of RNA-seq data without requiring intermediate results
22 to be transferred between the TACC and CyVerse cloud. **Figure 5** shows the relationship
23 among these analysis tools within the MCrna app, which is used to process both replicates of an

1 experiment in parallel.



2

3 **Figure 5.** MaizeCODE MCrna app for processing RNA-seq data from two replicates.

4

5 For each replicate, the MultiQC software outputs a quality report for the sequence data, before
6 and after trimming, in an interactive HTML format. This report can be accessed via the
7 visualization (eye) icon in the History panel (next to each loaded replicate, as shown in **Figure**
8 **2**). As with the HTML format, text, image, and other web browser-compatible files can be
9 visualized by clicking the icon. For files that can be displayed on a Genome Browser (e.g., BAM,
10 bigwig, etc), the user can also generate browser-ready links by clicking the same icon. These
11 links address the cloud storage system from the CyVerse project, so they can be displayed on
12 Genome Browsers hosted by different portals. If the user clicks on the output file name (from the
13 History panel), they will be directed to the CyVerse Data Commons landing page, where they
14 can preview or download the results. For files over a few GBs, we recommend that the user
15 download their data using either iCommands or CyberDuck, using the CyVerse Data Store path
16 available in the file URL.

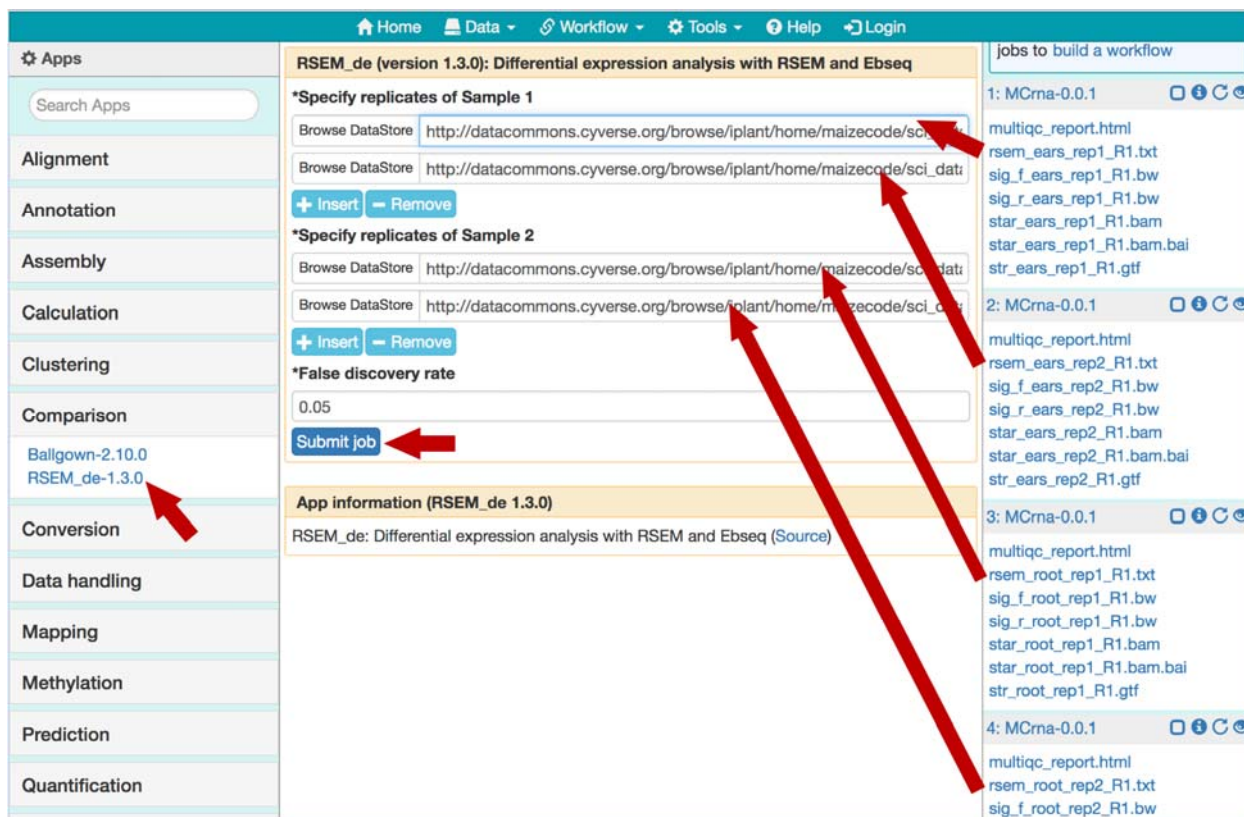
17 **Automated differential expression analysis**

1 One of the key advantages of distributing the MaizeCODE data through SciApps is that it
2 facilitates downstream analysis. In this section, we will show how differential expression
3 analysis can be performed, on either the gene or isoform levels, using the primary analysis
4 results discussed in the last section.

5
6 As shown in **Figure 6**, after loading the results for two experiments from the MaizeCODE data
7 page (**Figure 2**), one for ear tissue and the other for the root tissue, we can launch the
8 RSEM_de app from the 'Comparison' category (or through searching the app panel). For the
9 analysis, users drag and drop output files with names starting with 'rsem_' into the input field for
10 both replicates of each sample. The analysis job can then be submitted to the cloud for running,
11 and the results (i.e., the differentially expressed genes) will be available within a few minutes.
12 Note that the app is flexible in handling different numbers of replicates per sample. Additional
13 input fields can be added using the '+ Insert' button. For the MaizeCODE project, most data sets
14 are generated with two replicates.

15
16 Users can check the results through the History panel or the list of jobs (under the 'Workflow'
17 tab from the top menu). Users can also select jobs from the History panel and save them as
18 new workflows to organize the analysis and/or share it with others. Given that the XSEDE/TACC
19 cloud is a shared resource, and jobs may be queued for several minutes to several hours, we
20 have also established a local cluster (Wang et al., 2015) to quickly process small jobs requiring
21 less than an hour to complete. Powered by the Agave API (Dooley et al., 2012), SciApps treats
22 both the XSEDE/TACC cloud and the local cluster as virtual execution systems, allowing a
23 scientific app to be configured to run on either cloud. As such, by using the CyVerse Data Store
24 as a central storage hub, SciApps workflows can be seamlessly executed on a mixture of
25 XSEDE/TACC cloud and local clusters for efficient yet scalable processing and consumption of
26 MaizeCODE data.

1



2

3 **Figure 6.** Using the RSEM_de app for gene-level differential expression analysis

4

5 To perform differential expression on the isoform level, users can launch the workflow diagram

6 shown in **Figure 4**. From the diagram and the relaunched app forms, the StringTie_Merge app

7 is used to merge assembled transcripts from all replicates to generate a new annotation file.

8 The annotation file is then passed again to the StringTie app to compute gene expression

9 quantification results, which are then input to the Ballgown app (Frazee et al., 2015) to compute

10 differentially expressed isoforms. Again, the user can drag and drop the alignment files and

11 assembled transcripts from the MaizeCODE primary analysis results without repeating the time-

12 consuming alignment and quantification steps. Given that all results are accessible through web

13 URLs, users can also retrieve the data directly to their local server for further interactive

14 analysis. For example, SciApps users can inspect the quantification results of each pair of

1 replicates to confirm that they are consistent, and if they are, proceed with the analysis of
2 differentially expressed genes.

3 **CONCLUSION AND FUTURE WORK**

4 SciApps is a cloud-based platform that provides a data management infrastructure for the
5 MaizeCODE data. The platform supports the management of both experimental and
6 computational metadata; provides a collection of analysis apps covering analysis of multi-omics
7 data sets; and provides both web browser and API access to data, results, metadata, and
8 computational provenance of software tools through a unique workflow ID. Genome browser
9 tracks are also provided to enable visualization of the results using an integrated version of
10 JBrowse.

11
12 SciApps has been designed to integrate cloud resources for scaling and long-term stability.
13 Currently, SciApps has been integrated with the NSF-funded XSEDE/TACC cloud and the
14 CyVerse Data Store to provide academic users with cost-free data storage and computing
15 services. Both resources are integrated as virtual systems via the Agave API, which also
16 supports the integration with commercial cloud platforms like Amazon EC2/S3 and Microsoft
17 Azure. Therefore, SciApps can be scaled for large-scale data management and analysis if
18 needed. Additionally, SciApps can also seamlessly integrate local data and computing
19 resources to complement cloud resources. The successful utilization of our local system
20 suggests that SciApps can facilitate collaborative projects across different institutes for joint data
21 production, analysis, and management with multiple local systems, thereby avoiding high
22 financial costs.

23
24 To process the data sets that are continually being generated by the MaizeCODE project,
25 several new analysis tools have been integrated into SciApps. Future goals include developing

1 new analysis workflows, supporting sophisticated queries against metadata, reanalyzing and
2 distributing published sequencing data sets from raw data, and conducting training and
3 community outreach.

4

5 **IMPLEMENTATION**

6 The major components of the SciApps analysis portal include a web browser user interface, a
7 MySQL database, a workflow engine, and a web API. The user interface was built with the
8 React-Bootstrap front-end framework for data and workflow interactions. The workflow engine is
9 written in Perl and uses the MySQL database to track job status, perform submission of a
10 subsequent job once its inputs are ready, and record the analysis into the database. The web
11 API is written in Perl to complement the web browser user interface, especially for batch
12 processing and metadata handling. SciApps uses the CyVerse Central Authentication Service
13 (CAS) to grant users access to the XSEDE/TACC cloud and CyVerse Data Store. Analysis jobs
14 are submitted to the cloud, and in addition all released raw data, processed results, metadata,
15 and workflows, are made publicly accessible through both the GUI and API with no
16 authentication needed.

17 **AUTHOR CONTRIBUTIONS**

18 LW and ZL designed, implemented, and tested the software. MD and LW managed metadata
19 and submission of raw data to NCBI SRA. LW and PVB designed and maintained the local
20 system consisting of a web server, a data server, and several computing servers. LW, ZL, and
21 XW developed the scientific workflows and participated in testing through the web interface. AD
22 and TG helped in developing the workflow. XW and LW processed the data with the automated
23 workflows. MB, JD, XX, and COR generated the raw data. CG developed the Perl code for
24 interacting with Agave API. CFM, CG, and JW managed the MaizeCODE website. LW, DW,

1 MS, and TG wrote the manuscript. All the authors read and approved the final manuscript.

2

3 **AVAILABILITY AND REQUIREMENTS**

4

5 Project name: SciApps

6 Project home page: <https://sciapps.org/>

7 Source code repository: <https://github.com/warelab/sciapps>

8 Operating system(s): Platform independent

9 Programming language: JavaScript, Perl

10 License: MIT

11 Any restrictions to use the data: Toronto Agreement

12

13 **FUNDING**

14 This work is supported by NSF grant IOS 1445025 for MaizeCODE. SciApps has also been
15 supported by NSF grant DBI-1265383, and USDA-ARS (1907-21000-030-00D).

16

1 **References**

- 2 Dobin, A., Davis, C. A., Schlesinger, F., Drenkow, J., Zaleski, C., Jha, S., et al. (2013). STAR:
3 ultrafast universal RNA-seq aligner. *Bioinformatics* 29, 15–21.
- 4 Dooley, R., Vaughn, M., Stanzione D, T. S., and Skidmore, E. (2012). Software-as-a-service:
5 the iPlant Foundation API. in *5th IEEE Workshop on Many-Task Computing on Grids and*
6 *Supercomputers (MTAGS)*.
- 7 ENCODE Consortium (2004). The ENCODE (ENCyclopedia Of DNA Elements) Project.
8 *Science* 306, 636–640.
- 9 Ewels, P., Magnusson, M., Lundin, S., and Källér, M. (2016). MultiQC: summarize analysis
10 results for multiple tools and samples in a single report. *Bioinformatics* 32, 3047–3048.
- 11 Frazee, A. C., Pertea, G., Jaffe, A. E., Langmead, B., Salzberg, S. L., and Leek, J. T. (2015).
12 Ballgown bridges the gap between transcriptome assembly and expression analysis. *Nat.*
13 *Biotechnol.* 33, 243–246.
- 14 Goff, S. A., Vaughn, M., McKay, S., Lyons, E., Stapleton, A. E., Gessler, D., et al. (2011). The
15 iPlant Collaborative: Cyberinfrastructure for Plant Biology. *Front. Plant Sci.* 2, 34.
- 16 Grüning, B., The Bioconda Team, Dale, R., Sjödin, A., Chapman, B. A., Rowe, J., et al. (2018).
17 Bioconda: sustainable and comprehensive software distribution for the life sciences. *Nature*
18 *Methods* 15, 475–476. doi:10.1038/s41592-018-0046-7.
- 19 Jiao, Y., Peluso, P., Shi, J., Liang, T., Stitzer, M. C., Wang, B., et al. (2017). Improved maize
20 reference genome with single-molecule technologies. *Nature* 546, 524–527.
- 21 Krueger, F., and Andrews, S. R. (2011). Bismark: a flexible aligner and methylation caller for
22 Bisulfite-Seq applications. *Bioinformatics* 27, 1571–1572.

- 1 doi:10.1093/bioinformatics/btr167.
- 2 Kurtzer, G. M., Sochat, V., and Bauer, M. W. (2017). Singularity: Scientific containers for
3 mobility of compute. *PLoS One* 12, e0177459.
- 4 Li, B., and Dewey, C. N. (2011). RSEM: accurate transcript quantification from RNA-Seq data
5 with or without a reference genome. *BMC Bioinformatics* 12, 323.
- 6 Pertea, M., Pertea, G. M., Antonescu, C. M., Chang, T.-C., Mendell, J. T., and Salzberg, S. L.
7 (2015). StringTie enables improved reconstruction of a transcriptome from RNA-seq reads.
8 *Nat. Biotechnol.* 33, 290–295.
- 9 Schnable, P. S., Ware, D., Fulton, R. S., Stein, J. C., Wei, F., Pasternak, S., et al. (2009). The
10 B73 maize genome: complexity, diversity, and dynamics. *Science* 326, 1112–1115.
- 11 Skinner, M. E., Uzilov, A. V., Stein, L. D., Mungall, C. J., and Holmes, I. H. (2009). JBrowse: A
12 next-generation genome browser. *Genome Research* 19, 1630–1638.
13 doi:10.1101/gr.094607.109.
- 14 Sloan, C. A., Chan, E. T., Davidson, J. M., Malladi, V. S., Strattan, J. S., Hitz, B. C., et al.
15 (2016). ENCODE data at the ENCODE portal. *Nucleic Acids Res.* 44, D726–32.
- 16 Townes, J., Cockerill, T., Dahan, M., Foster, I., Gaither, K., Grimshaw, A., et al. (2014). XSEDE:
17 Accelerating Scientific Discovery. in *Computing in Science & Engineering* 5., 62–74.
- 18 Wang, L., Lu, Z., Van Buren, P., and Ware, D. (2018). SciApps: a cloud-based platform for
19 reproducible bioinformatics workflows. *Bioinformatics* 34, 3917–3920.
- 20 Wang, L., Van Buren, P., and Ware, D. (2015). Architecting a distributed bioinformatics platform
21 with iRODS and iPlant Agave API. in *International Conference on Computational Science*
22 *and Computational Intelligence (CSCI)*, 420–423.