# Creating a FAIR data ecosystem for incorporating single-cell genomics data into agricultural G2P research

**Muskan Kapoor**, Enrique Sapena Ventura, Alexey Sokolov, Galabina Yordanova, Nicholas J. Provart, Irene Papatheodorou, Nancy George, Doreen Ware, Sunita Kumari, Timothy Tickle, Lance Daharsh, James Koltes, Benjamin Cole, Marc Libault, Christine Elsik, Wesley Warren, Tony Burdett, Peter Harrison, and Christopher Tuggle

## ABSTRACT

The agriculture genomics community has numerous data submission standards but limited knowledge describing and storing single-cell (e.g., scRNAseq) data. Other single-cell genomics infrastructure efforts, such as the Human Cell Atlas Data Coordination Platform (HCA DCP), have resources that could benefit our community. For example, the HCA DCP is integrated with Terra, a cloud-native workbench for computational biology developed by Broad, Verily, and Microsoft that houses tools for scGenomics analysis. In Aim-1 we describe a pilot-scale project to determine if our current metadata standards for livestock and crops can be used to ingest scRNAseq datasets in a manner consistent with HCA DCP standards and if established resources (e.g., Terra) can be used to analyze the ingested data. Currently, the most comprehensive data ingestion portal for high throughput sequencing datasets from plants, fungi, protists, and animals (including humans) at the EMBL-European Bioinformatics Institute, Annotare, ensures that sufficient metadata are collected to enable re-analysis and dissemination via the Single Cell Expression Atlas (SCEA). Another EMBL-EBI portal limited to animal datasets, the FAANG portal, provides bulk and scRNAseq data access which uses a semi-automated process to submit and validate files using the HCA DCP metadata and data validation service. Once incorporated, datasets are used to augment the DCP resource for the scientific community. These files are also incorporated using EMBL-EBI's HCA DCP ingestion service and then transferred to Terra for further analysis. In Aim-2 we test and develop prototype tools to visualize the output of scRNAseq analyses on genome browsers and comparing across tissues and cell populations.

## GOAL: IMPROVING SC DATA ANALYSIS INFRASTRUCTURE



Fig 1. (A) Current Status describing the pilot scale project when the data and meta-data file is created and transferred to a computing environment without the help of data wranglers or curators
(B) Future Vision: a more detailed explanation for transferring data to a computing environment that can further be utilized by the agricultural community.

## GENERAL WORKFLOW



Fig 2. The general route of meta-data flow in Animal, Plant, and Public Cell Atlas.

## PLANT SIDE WORKFLOW



Fig 3. Ingestion of Plant Side, the Single-cell data from public archives follows a route through three scripts. The data can be visualized in the SCEA portal itself and analyzed in the GALAXY through web API retrieval.

## RESULTS- AIM 1

## ANIMAL SIDE WORKFLOW



Fig 4A. Ingestion of Animal Side workflow, the meta-data from the FAANG data portal was validated and transferred to HCA- DCP ingestion service. The red box indicates that we are generating a new spreadsheet because of the schema difference in both portals.



Fig 4B. Ingestion in HCA-DCP; First result is one of the entities in the API, showing the transformation in the format. The second result is an image of the UI of validated entities
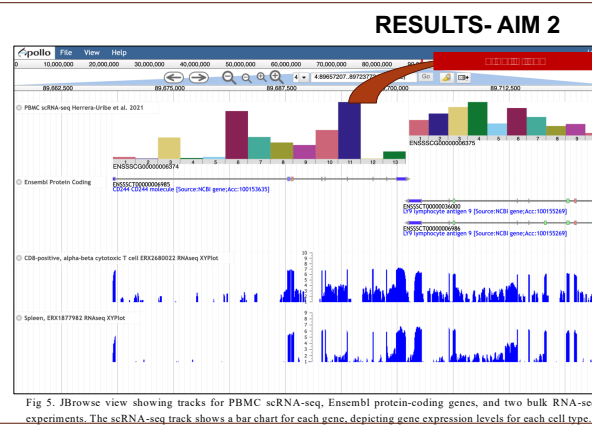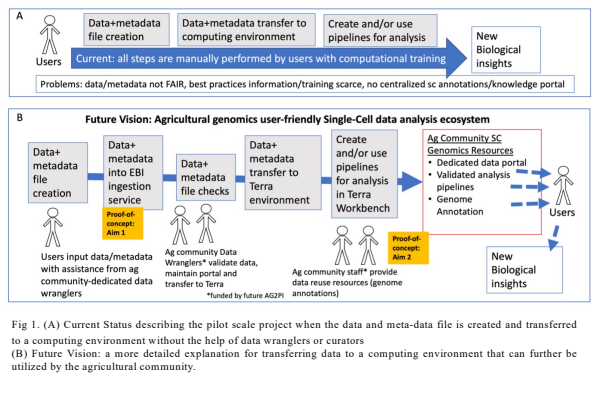
## RESULTS- AIM 2



Fig 5. JBrowse view showing tracks for PBMC scRNA-seq, Ensembl protein-coding genes, and two bulk RNA-seq experiments. The scRNA-seq track shows a bar chart for each gene, depicting gene expression levels for each cell type.
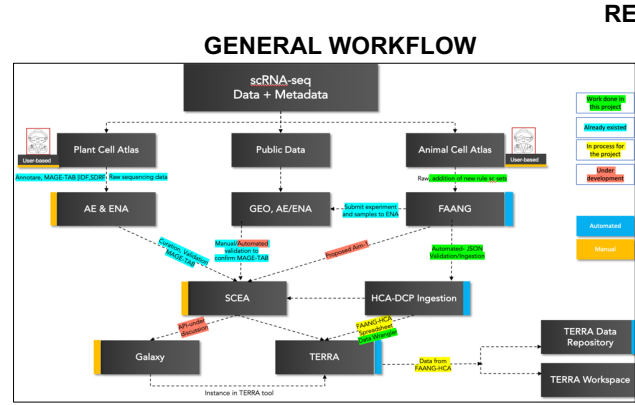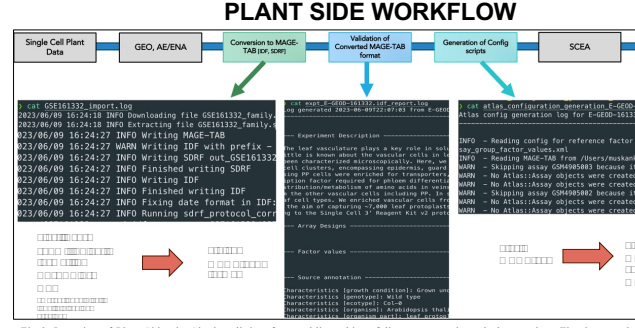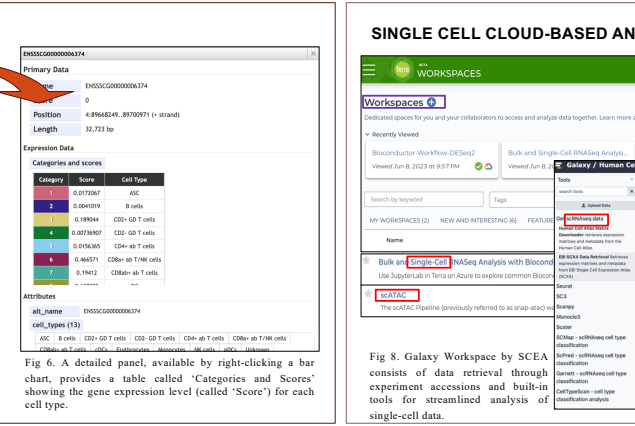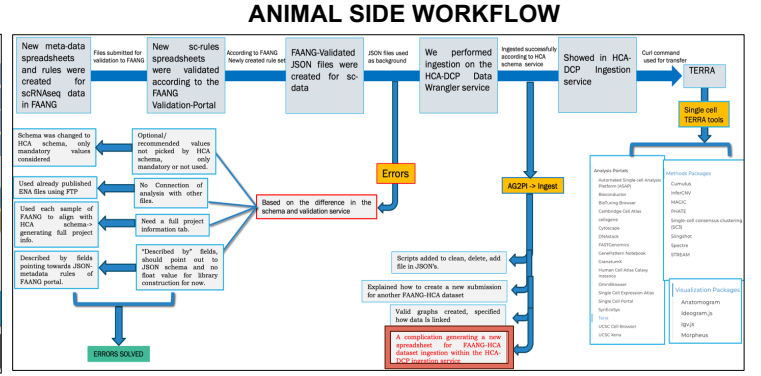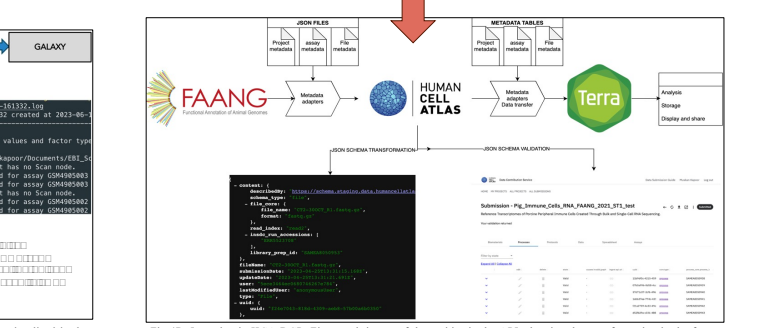


Fig 6. A detailed panel, available by right-clicking a bar chart, provides a table called 'Categories and Scores' showing the gene expression level (called 'Score') for each cell type.

## SINGLE CELL CLOUD-BASED ANALYSIS TOOLS- FUTURE VISION



Fig 7. TERRA Workspace for exploring single-cell data, and in-built pipelines for standardizing the workflows. Data with standard metadata can be retrieved using the curl command for simplified meta-analyses.



Fig 8. Galaxy Workspace by SCEA consists of data retrieval through experiment accessions and built-in tools for streamlined analysis of single-cell data.

## CONCLUSIONS & FUTURE SCOPE

- We are building upon existing tools to develop a scientist-friendly data resource and analytical ecosystem and facilitate single cell-level genomic analysis across agricultural species. This is being accomplished through data ingestion, storage, retrieval, re-use, visualization, and comparative annotation.

- On the animal side, we intend to develop an automated pipeline to transfer the meta-data from the FAANG portal to the SCEA portal.

- On the Plant side, we intend to make the pipeline to SCEA more automated, in a similar fashion to the FAANG data portal functionality.

## CONTACT INFORMATION

Muskan Kapoor

Graduate Research Assistant

Bioinformatics & Computational Biology, Animal Science

Email: muskan@iastate.edu ; cktuggle@iastate.edu

**IOWA STATE UNIVERSITY**
Bioinformatics & Computational Biology Graduate Program