UCLA

Publications

Title

Too big to share? Scaling up knowledge transfer workflows from little science to big science

Permalink

https://escholarship.org/uc/item/3rc3p5qp

Authors

Randles, Bernadette M. Sands, Ashley E. Borgman, Christine L.

Publication Date

2016-04-01

Copyright Information

This work is made available under the terms of a Creative Commons Attribution-NonCommercial-NoDerivatives License, available at https://creativecommons.org/licenses/by-nc-nd/4.0/

Too Big to Share?

Scaling Up Knowledge Transfer Workflows in Computational Sciences Bernadette M. Randles, Ashley E. Sands, Christine L. Borgman UCLA, Information Studies

Problem:

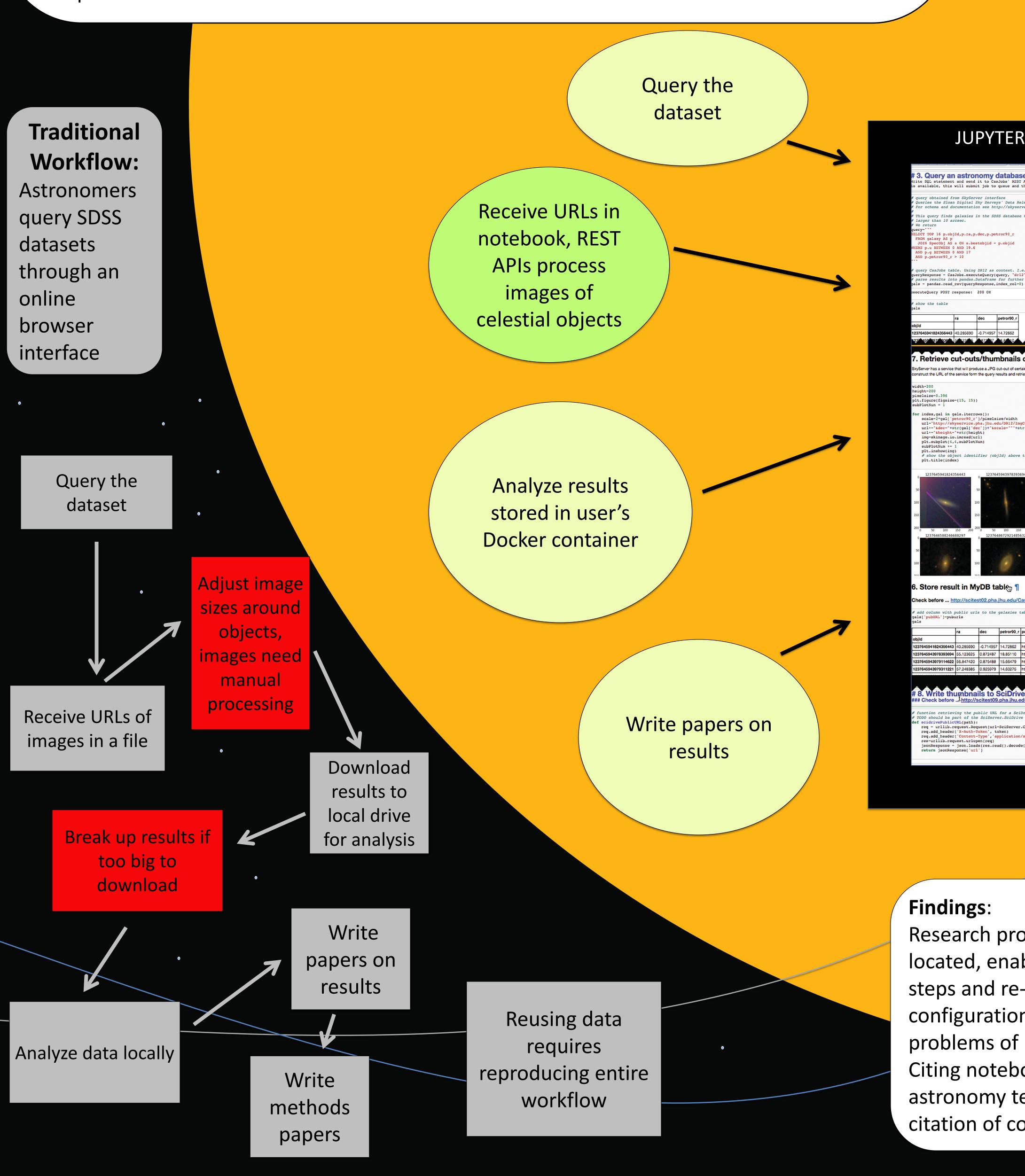
Academic researchers experience computational and storage challenges when retrieving and analyzing up to petabyte-scale datasets. Researchers desire infrastructures and tools that support scalable retrieval, storage, and analysis of research data. Adding to these challenges, funders may require the resulting datasets to be cited on a more granular level, be reproducible and publically available, further complicating the traditional research process of documentation.

Case Study:

We compare a traditional astronomy research workflow of querying the Sloan Digital Sky Survey dataset for celestial object images to a new method using a Jupyter notebook. The Jupyter notebooks are the astronomer's main interface tool to remotely access datasets, providing a live computational research document¹ containing functionality to access, analyze, and display results of dataset queries or simulations. This particular (beta) instantiation of the Jupyter notebook platform resides on the server or server cluster where the data lives, keeping the analysis close to the data². Another key element of the new workflow is the Docker container, a lightweight virtual space on the server enabling separate, storage spaces for users to perform larger computations and queries.

What are Jupyter Notebooks?

Jupyter notebooks are (i)python-based, open source, browser-based documents that support rich text elements, graphs, images, and executable code; supporting over 50 computer languages such as R, Python, and C++. Jupyter notebook sociotechnical benefits include a dedicated user community. Drawbacks emerge in resistance to switch tools and work habits.



Generate public URL to link to the notebook Reusing data entails accessing

New Workflow:

the notebook

through public

URL

Queries to find celestial objects from the SDSS dataset are in the notebook, with textual comments, code, and results including images all contained in one place.

Findings:

JUPYTER NOTEBOOK

Research processes, code, and results are colocated, enabling others to retrace procedural steps and re-run the same code. This configuration shows potential to address problems of reproducibility and discoverability. Citing notebooks may solve some of the astronomy team's difficulties surrounding citation of code and datasets.

- 1 Millman, K. J., & Pérez, F. (2014). Developing Open-Source Scientific Practice. *Implementing Reproducible Research*, 149.
- 2 Raicu, I., Zhao, Y., Foster, I. T., & Szalay, A. (2008, June). Accelerating large-scale data exploration through data diffusion. In Proceedings of the 2008 international workshop on Data-aware distributed computing (pp. 9-18). ACM.

Thank you to the Alfred P. Sloan Foundation ("If Data Sharing is the Answer, What is the Question?" Award# 2015-14001. Christine L. Borgman, UCLA, Principal Investigator).

Thank you to the Johns Hopkins University IDIES, SciServer and SDSS teams.

Thank you to the members of the UCLA Knowledge Infrastructures Team, which include: Peter T. Darch, Irene Pasquetto, and Milena Golshan. http://knowledgeinfrastructures.gseis.ucla.edu @UCLA_KI (and yes, Jupiter does have rings!)



