

SHED: Streaming Heterogeneous Event Data Tracking with Provenance

Christopher J. “CJ” Wright¹, Songsheng Tao¹, and Simon J. L. Billinge^{*1, 2}

¹ Department of Applied Physics and Applied Mathematics, Columbia University, New York, NY 10027 ² Brookhaven National Laboratory, Upton, NY 11973

DOI: [10.21105/joss.03119](https://doi.org/10.21105/joss.03119)

Software

- [Review](#) ↗
- [Repository](#) ↗
- [Archive](#) ↗

Editor: [Arfon Smith](#) ↗

Reviewers:

- [@remram44](#)

Submitted: 29 October 2020

Published: 17 March 2021

License

Authors of papers retain copyright and release the work under a Creative Commons Attribution 4.0 International License ([CC BY 4.0](#)).

Summary

The SHED package provides a python based software framework for real time analysis of streaming data with provenance. It is highly flexible, allowing analysis pipelines to be built for any variety of experimental data that is arriving in a time series. It allows straightforward serialization of the pipeline, and the analysis results, for storage in databases with provenance information that allows analyses later to be pulled from the databases, adapted if desired, and rerun.

Statement of Need

The accelerating rate of data coming from modern scientific experiments enable completely new classes of experiments to be carried out such as *in situ* and *operando* measurements of time-dependent phenomena. Many current and planned world class experimental facilities such as x-ray synchrotron sources produce mountains of data. Real-time processing of these large volumes of data can be made possible by streaming data processing, where the data is processed and reduced as it is being collected. While many streaming data processing systems exist ([Silva et al., 2007](#))([Davison et al., 2014](#))([Picard, 2018](#))([Carbone et al., 2015](#)) few handle the heterogeneity which is prevalent in many experimental environments. For example, x-ray synchrotrons support thousands of users per year doing a wide array of bespoke experiments and measurements. Additionally, few streaming systems track the provenance of the processed data, creating reproducibility and discoverability concerns.

SHED is a Python package for handling heterogeneous streaming data and the tracking of the provenance of the produced results. The SHED API translates data coming from a streaming experiment into python objects which are passed into data processing pipelines written using the `streamz` library. It expects data in the form of the Bluesky Event Model schema, which is a flexible open-source document model suitable for scientific time-series experiments. The Bluesky event model ([Koerner et al., 2020](#)) is used by nearly all of the beamlines at the National Synchrotron Light Source-II (NSLS-II) as well as at other facilities.

SHED also provides translation from the python objects flowing through a `streamz` pipeline back into the Event Model enabling the use of data visualization and storage tools written for the Event Model. In addition to the translation features SHED passively tracks the provenance of produced data, capturing the data processing pipeline, data unique IDs, and the order of data insertion into the pipeline. The pipeline is captured in a way which is both human and machine readable, enabling searching and comparison of analyzed data sets. An extendable

^{*}Corresponding author

39 system also captures information about the software environment in which the pipeline ran,
40 including the conda environment, allowing in principle data analysis campaigns to be serial-
41 ized into a database then later recovered and rerun, or adapted and then rerun and stored
42 independently.

43 SHED is designed to provide a bridge between the data collected from high throughput ex-
44 periments if they are, or can be put, into event model form (such as data from NSLS-II
45 experiments), and the NSLS-II's data visualization and storage tools. This system has been
46 used to write x-ray scattering data processing pipelines currently used in production at the
47 NSLS-II XPD and PDF beamlines (28-ID-2 and 28-ID-1, respectively). The design and some
48 applications of shed are described in greater detail in (Wright, 2020).

49 Figures

50 Example data processing pipeline with SHED nodes.

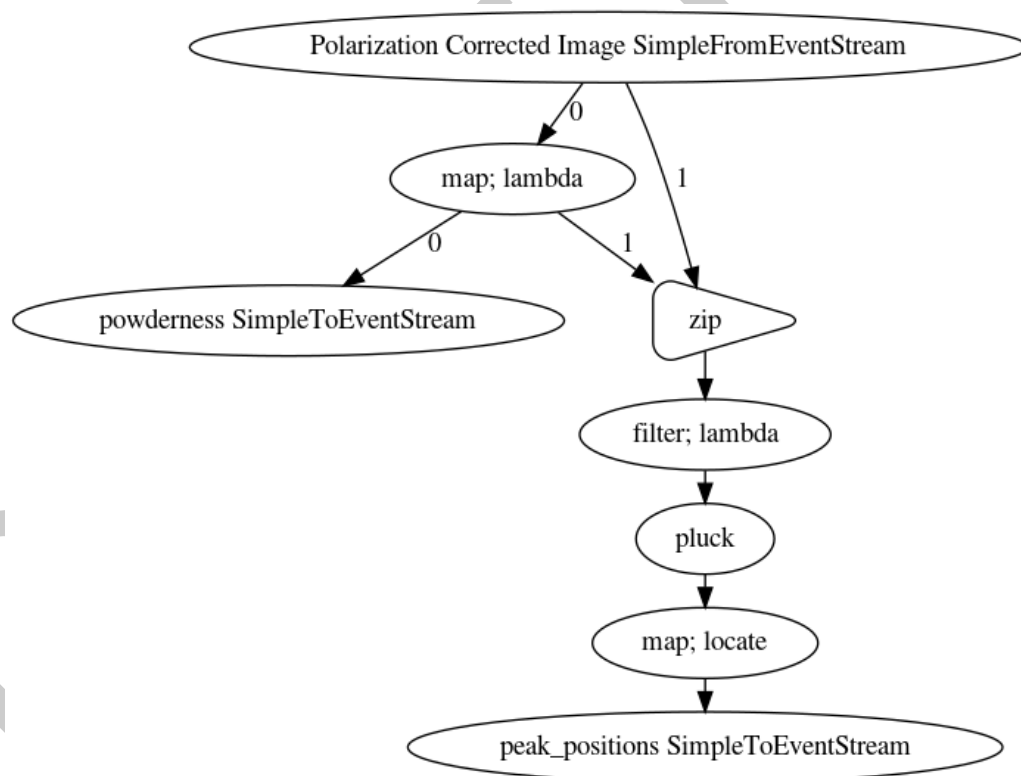


Figure 1: example pipeline

51 Acknowledgments

52 This work was supported as part of GENESIS: A Next Generation Synthesis Center, an Energy
53 Frontier Research Center funded by the U.S. Department of Energy, Office of Science, Basic
54 Energy Sciences under Award Number DE-SC0019212

References

- 56 Carbone, P., Katsifodimos, A., Ewen, S., Markl, V., Haridi, S., & Tzoumas, K. (2015). Apache
57 flink: Stream and batch processing in a single engine. *Bulletin of the IEEE Computer*
58 *Society Technical Committee on Data Engineering*, 36(4).
- 59 Davison, A. P., Mattioni, M., & Samarkanov, D. (2014). Sumatra: A Toolkit for Reproducible
60 Research. In *Implementing Reproducible Research*. <https://doi.org/10.1201/b16868-9>
- 61 Koerner, L. J., Caswell, T. A., Allan, D. B., & Campbell, S. I. (2020). A python instrument
62 control and data acquisition suite for reproducible research. *IEEE Transactions on In-*
63 *strumentation and Measurement*, 69(4), 1698–1707. [https://doi.org/10.1109/TIM.2019.](https://doi.org/10.1109/TIM.2019.2914711)
64 [2914711](https://doi.org/10.1109/TIM.2019.2914711)
- 65 Picard, R. (2018). *Hands-on reactive programming with python: Event-driven development*
66 *unraveled with RxPY*. Packt Publishing Ltd.
- 67 Silva, C. T., Freire, J., & Callahan, S. P. (2007). Provenance for Visualizations: Repro-
68 ducibility and Beyond. *Computing in Science Engineering*, 9(5), 82–89. [https://doi.org/](https://doi.org/10.1109/MCSE.2007.106)
69 [10.1109/MCSE.2007.106](https://doi.org/10.1109/MCSE.2007.106)
- 70 Wright, C. J. (2020). *Towards real time characterization of grain growth from the melt* [PhD
71 thesis, Columbia University]. <https://doi.org/10.7916/d8-5tdx-bb07>

DRAFT