

Data Access Pattern Analysis for dCache Storage System

ABSTRACT

dCache is a storage management system with disk space for temporary and permanent file storage and connectivity to a larger high performance storage system (HPSS) for additional file storage. Ideally, all files will be stored locally on the dCache disks when they need to be accessed. Analysis of past dataset access patterns for the files stored on dCache is necessary to determine how dataset popularity varies over time and what factors might affect the fluctuation. Findings suggest that datasets are generally accessed only a few days at a time with an inconsistent day to day popularity measurement, and therefore could only be stored locally for the few days that they are popular.

BACKGROUND

- Particle collision events being studied produce a very high volume of data
- A majority of the data is stored on HPSS tapes, and the rest is stored on dCache disks.

Data is retrieved through a request to the dCache System

dCache searches local cache for the file

Local copy found **in cache:** File can be read from or written to with no HPSS interaction.

Local copy not found in cache: a copy of the file has to be moved for access.

Motivation:

Caching data shortens access time and reduces latency. Forecasting dataset popularity for a future time frame allows the dataset to be cached prior to when it needs to be accessed.



Amanda Pereira¹, Alex Sim², John Wu², Shinjae Yoo³, Hiro Ito³ Drake University¹, Lawrence Berkeley National Laboratory², Brookhaven National Laboratory³

RESEARCH QUESTION

How can future dataset popularity be predicted given information on past dataset popularity?

ANALYSIS METHODS

The analysis of data transfer logs and metadata information are used to address this research question by grouping files into datasets and studying the popularity fluctuation.



ool.root.1

All filepaths with the <u>same 8 digit number</u> are <u>grouped together</u>.





- Greater fluctuation in

SUMMARY

- Search for patterns in popular and unpopular files may help what factors may lead to a high popularity
- For maximum storage efficiency, datasets could be stored in the local cache during the time period in which they are popular, and then moved to the HPSS for storage during the time they are not popular.

FUTURE WORK:

- Develop a machine learning model to next time frame



ACKNOWLEDGMENTS

This work was supported in part by the Office of Advanced Scientific Computing Research, Office of Science, of the U.S. Department of Energy under Contract No. DE-AC02-05CH11231, and also used resources of the National Energy Research Scientific Computing Center (NERSC). This work was also supported in part by the U.S. Department of Energy, Office of Science, Office of Workforce **Development for Teachers and Scientists** (WDTS) under the Science Undergraduate Laboratory Internship (SULI) program.





U.S. DEPARTMENT OF ENERGY