

ABSTRACT

WORKFORCE

DEVELOPMEN1

& EDUCATION

Today's scientific projects and simulations often require repeated transfer of large data volumes between the storage system and the client. This increases the load on the network, leading to congestion. In order to mitigate these effects, regional data storage cache systems are used to store data locally. This project examines the XCache storage system to closely analyze data trend patterns in the data volume and data throughput performance, while also creating a model for predicting how caches could potentially impact network traffic and data transfer performance overall. The results of the data access patterns demonstrated that traffic volume was reduced by an average factor of 2.35. The hourly and daily prediction models also showed low error values, reinforcing the learning methods used in this effort.

BACKGROUND

- High-Luminosity Large Hadron Collider aims to increase performance by 2025
- Study uses data from Southern California Petabyte Scale Cache (SoCal Repo) of 24 XCache nodes
- SoCal Repo handles data from Large Hadron Collider (LHC), which is expected to produce 30x more data in 2028 than 2018

• Study period: July 2021 - June 2022 Table 1. Cummers, of nodes in the regional each

Table T. Summary of hodes in the regional cache					
	UCSD	Caltech	ESnet		
Number of Nodes	12	11	1		
Disk Capacity	24 TB each	98 - 288 TB	40 TB		
Network Connections	10 Gbps each	40 Gbps each	40 Gbps		

Table 2: Summary data access from July 2021 to June 2022						
	Number of accesses	Data transfer size (TB)	Shared data size (TB)	Shared data Percentage		
Total	8,021,922	8,210.78	4,499.44	35.40%		
Daily Average	22,283	22.81	12.46			

DATA











Figure 7: Hourly average throughput performance per access with 24-hour moving average: (a) cache hits, (b) cache misses.



Л	activation function	dropout rate	# of training epoch
	tanh	0.04	50
	relu	0.1	50
	tanh	0.04	50
	relu	0.1	50

DISCUSSION

- Daily cache utilization shows different patterns over time, and streaming jobs were observed between Oct 2021 and Jan 2022.. (Fig. 1)
- Avg traffic volume reduction for the whole period is 1.55, and avg traffic volume reduction from July 2021 to Sept 2021 is 2.35. (Fig. 2)
- LSTM models for hourly data tend to have less errors than models for daily data, primarily due to the fact that there are more data points. (Fig. 4 & 6)
- LSTM model with moving average fits better with less extreme values. (Fig. 5 & 7)
- Hourly data plot for cache misses shows different patterns compared to the hourly data for cache misses with 24-hour moving average (Fig. 6b & 7b), indicating that 24-hour moving average reduced extreme values, thus there are smaller RMSE values.

CONCLUSIONS

- General in-network regional cache could supplement the existing repository and benefit wider user community, reduce the redundant data transfers, and save network traffic bandwidth.
- Cache utilization and network throughput performances are predictable by LSTM models.
- Prediction can be better with more data records.
- Further studies: data access patterns of the different regional repositories, and longer term network requirements

URTHER READING



ACKNOWLEDGMENTS

Thanks to my mentors and collaborators, Inder Monga, Chin Cuok, Alex Sim, John Wu, Frank Würthwein, Diego Davila, Harvey Newman and Justas Balcas. This work was supported in part by the Office of Advanced Scientific Computing Research, Office of Science, of the U.S. Dept. 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. Dept. of Energy, Office of Science, Office of Workforce Development for Teachers and Scientists (WDTS) under the Science Undergraduate Laboratory Internship (SULI) program.

Office of Science

