



Exploring in-network data caching - ESnet-US CMS collaboration study

Alex Sim, Katherine Zhang, Ellie Copps, John Wu at LBNL
Chin Guok, Inder Monga at ESnet
Diego Davila, Frank Wuerthwein, Edgar Hernandez at UCSD

- **We know about**

- Data volume increase in experiments and simulations
- Data volume moving through network also increases
- Network bandwidth requirement gets higher

- **Observation**

- Significant portion of the popular dataset is transferred multiple times to different users as well as to the same user

- **Data sharing**

- Reduce the redundant data transfers
- Save network traffic volume, consequently.
- Lower data access latency
 - **Overall application performance is expected to be improved**

Pilot experiment

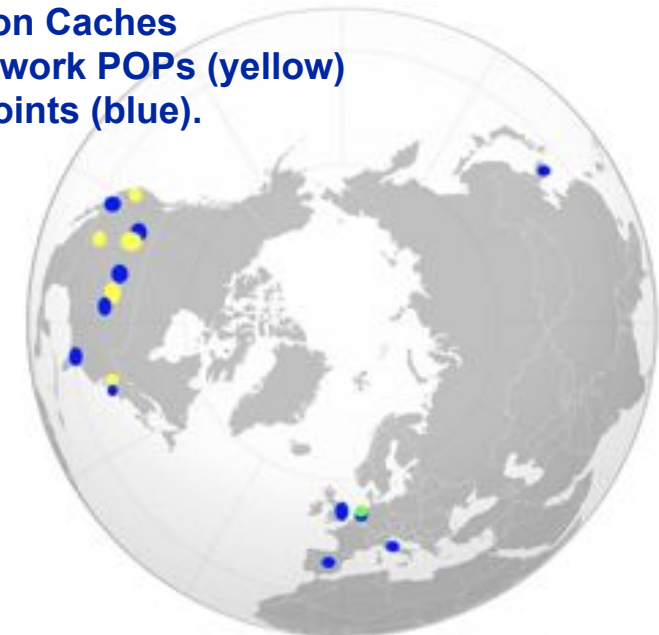
- **In-network temporary data cache for data sharing**
- **Collaboration with UCSD, US CMS, Open Science Grid (OSG)**
 - ESnet cache node as a part of SoCal Petabyte scale cache for US CMS
 - Petabyte scale cache is deployed/operated by UCSD and Caltech
 - ESnet: Provide a storage host to US CMS
 - Installed a temporary storage cache node in ESnet network
 - Prepared a server, which is physically connected to nersc-tb1 but having its routing original on sunn-cr55
 - Monitor the network utilization of the node
 - UCSD: Operation support for the ESnet node
 - Deploy/operate the OSG/Stashcache(Xcache) software stack and monitoring
 - Application-level monitoring
- **Goals**
 - Study how network cache storage helps network traffic performance and the overall application performance
 - Accumulate experience on how the US DOE scientific experiments and simulations share data among their users

- **Diverse science relevant to DOE HEP & NP**
 - Regional storage repo and data caching are one of the hot topics at HSF/WLCG meetings
 - At present, there are caches in production for ATLAS, CMS, and OSG, all based on XRootD
 - **OSG cache use dominated by Dune, LIGO, Virgo, MINERVA, DES, NOVA, and a liquid XENON detector R&D for future dark matter and neutrinoless double beta decay experiments. Electron Ion Collider R&D is in planning**

Collaboration	Working Set	Data Read	Reread Multiplier
DUNE	25 GB	131 TB	5400
LIGO (private)	41.4 TB	3.8 PB	95
LIGO (public)	4.3 TB	1.5 PB	318
MINERVA	351 GB	116 TB	340
DES	268 GB	17 TB	66
NOVA	268 GB	308 TB	1200
RPI_Brown	67 GB	541 TB	8300

OSG Data Federation Caches are deployed at network POPs (yellow) and compute endpoints (blue).

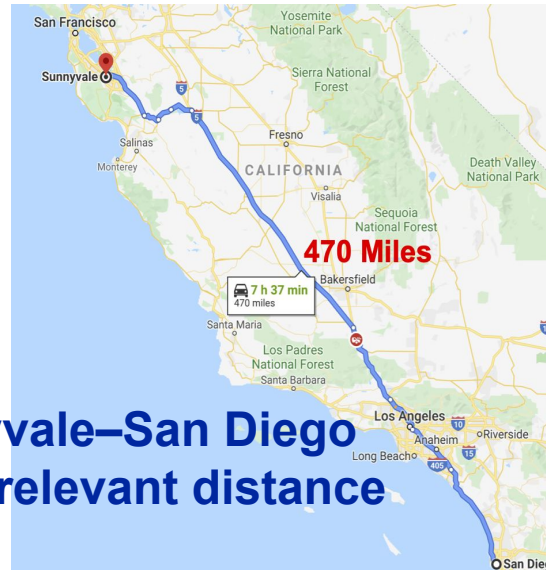
- Cache at institution
- Cache in the backbone
- Future Deployments



Data read from OSG data federation caches in 6 month period 3/2020-8/2020

Application use case with CMS

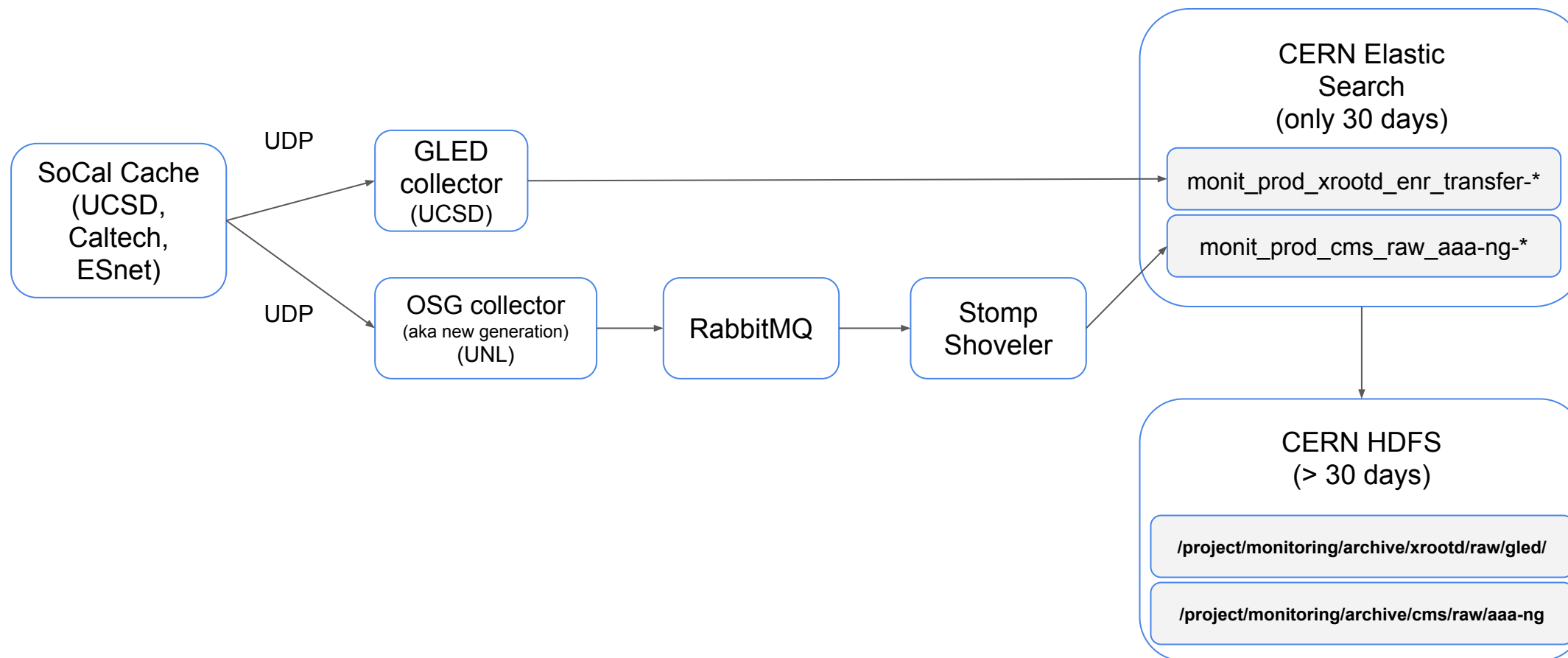
- **R&D Towards HL-LHC**
 - High-Luminosity-LHC: the LHC performance to increase the potential for discoveries after 2025
 - All processing done via buffers
 - All analysis done via caches
- **High level assumptions of annual volumes and use**
 - 384 PB of RAW } Mostly kept on Tape => accessed a couple times per year
 - 240 PB of AOD } Mostly kept on disk => heavily re-used by many researchers
 - 30 PB of MINI }
 - 2.4 PB of NANO }
- **Petabyte scale cache for CMS in CA**
 - Deployed/Operated by UCSD and Caltech
 - To gain experience with MiniAOD reuse
 - Includes the ESnet cache node
 - 500 miles distance for a distributed cache is a socio-politically very relevant distance scale



Sunnyvale–San Diego is the relevant distance scale

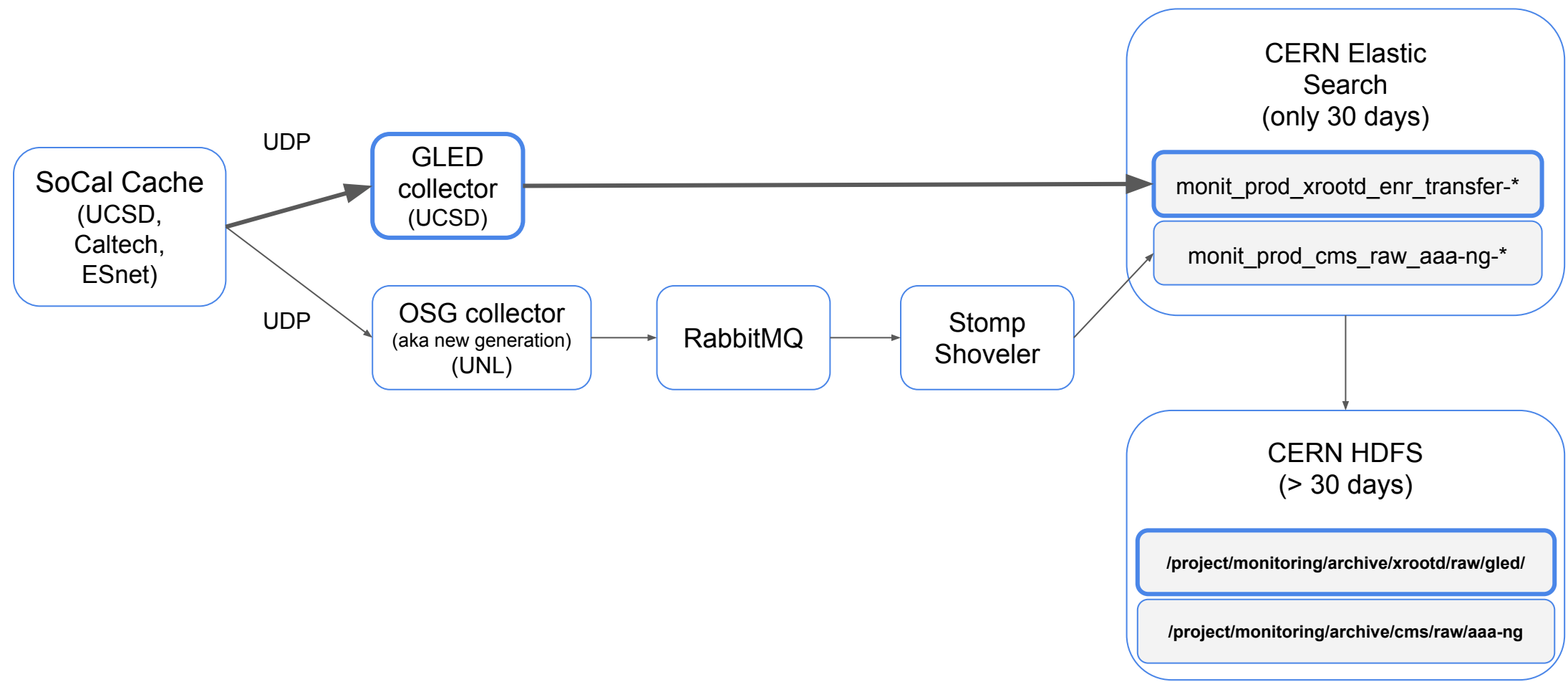


Current monitoring data paths



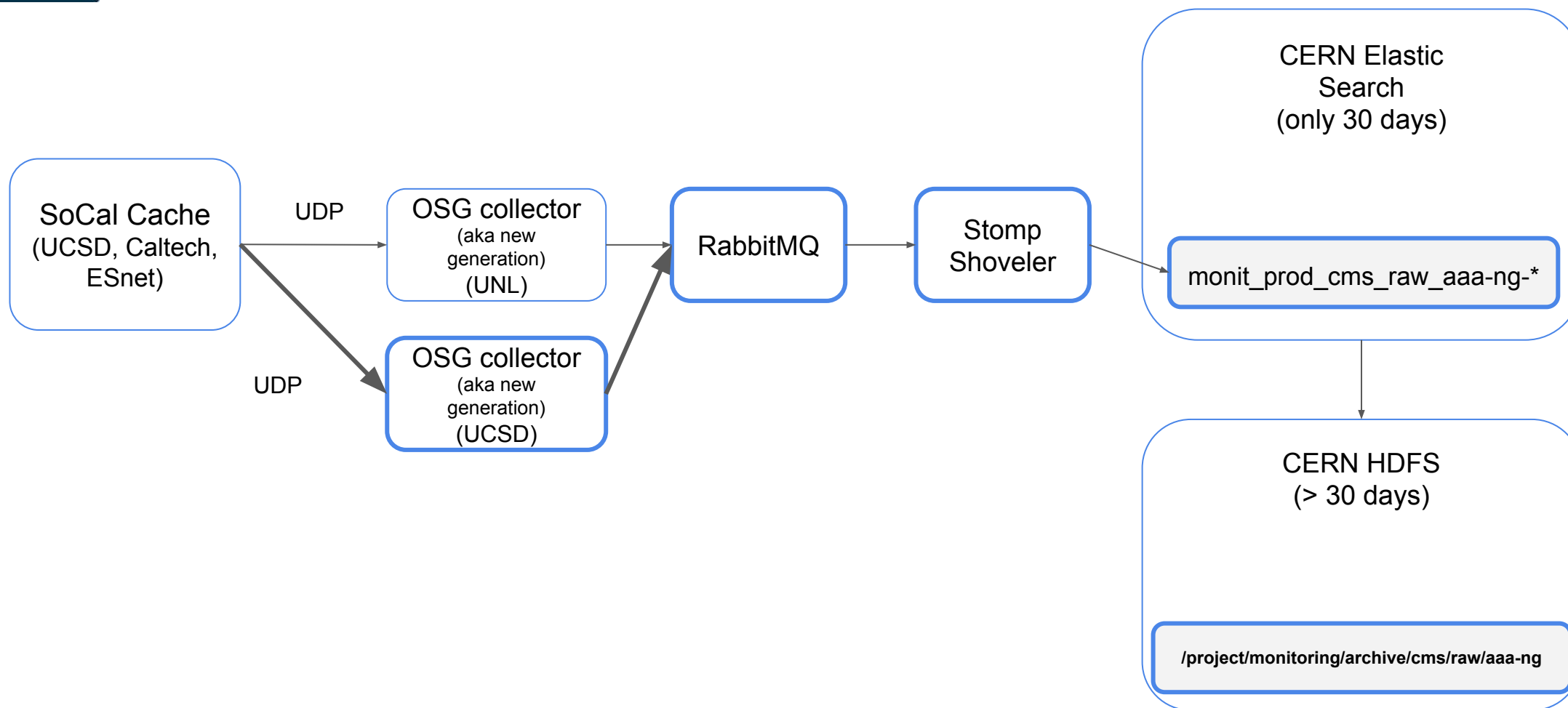
- The XRootD servers send monitoring data to 2 different collectors
- The data ends up in 2 different indexes at ES and kept there for a maximum of 30 days
- Periodically, the ES data gets stored in HDFS for long-term storage

Current used path



Given the proximity of the GLED collector to the SoCal cache, the chances for UDP packets being dropped are less on this path

Coming next



In order to deprecate the GLED collector without the risk of losing data we will deploy a New Generation collector at UCSD so that it is close to the SoCal cache



Acknowledgements & Resources

- **Acknowledgements**

- Alex, John, Katherine, Ellie at SDM, CRD, LBNL
- Adam, Anne, Chin, Eli, Eric, George, Goran, Inder, Kate, Yatish at ESnet
- Diego Davila, Dima Mishin, Edgar Hernandez, Frank Wuerthwein, Michael Sinatra at UCSD
- ESnet Infrastructure team: System build, config, and management
- ESnet Engineering team: Network connectivity

- **Resources**

- Hardware: 40TB storage and 40Gbps networking capability
- Expected network utilization: about 10-20 Gbps

- **In operation since May 2020**

- **This work is supported by the U.S. Department of Energy, Office of Science, Office of Advanced Scientific Computing Research, under Contract No. DE-AC02-05CH11231, and also used resources of the National Energy Research Scientific Computing Center (NERSC).**

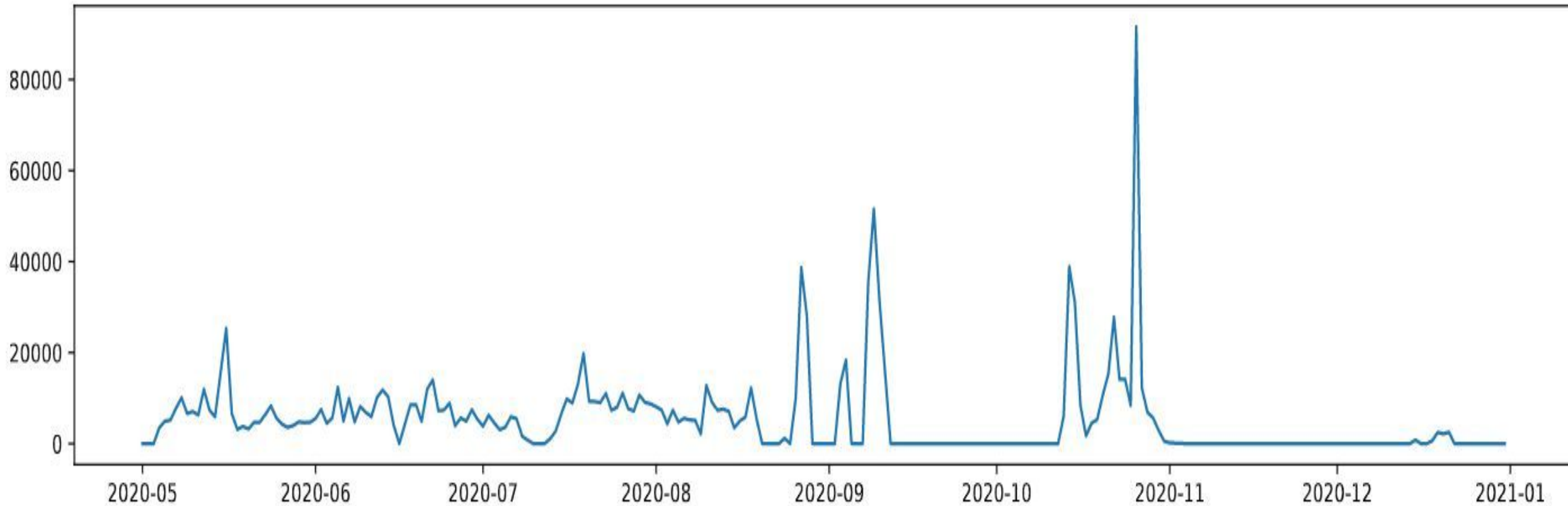
Summary statistics for data accesses

	Number of accesses	Data transfer size (GB)	Shared data access size (GB)
May 4-31, 2020	189,984	30,150.50	47,986.56
June 2020	215,452	40,835.23	55,929.47
July 2020	205,478	33,399.81	66,457.35
Aug 2020	203,806	30,819.80	68,723.19
Sep 2020	165,910	10,153.97	38,036.19
Oct 2020	306,118	22,723.93	45,614.91
Nov 2020	276	3.33	47
Dec 2020	8514	1236.81	4523
Total (May-Oct)	1,286,748	168,083.27	322,747.67
Daily average	9,674.79	1,263.8	2,426.67

- Data transfer size (= first time data access size, cache misses): From remote sites to the local node cache
- Shared data access size (= repeated data accesses, cache hits, network bandwidth savings): From the local node cache to the application, excluding the first time accesses (data transfers)
- Total number of active days until the end of Oct 2020: 133



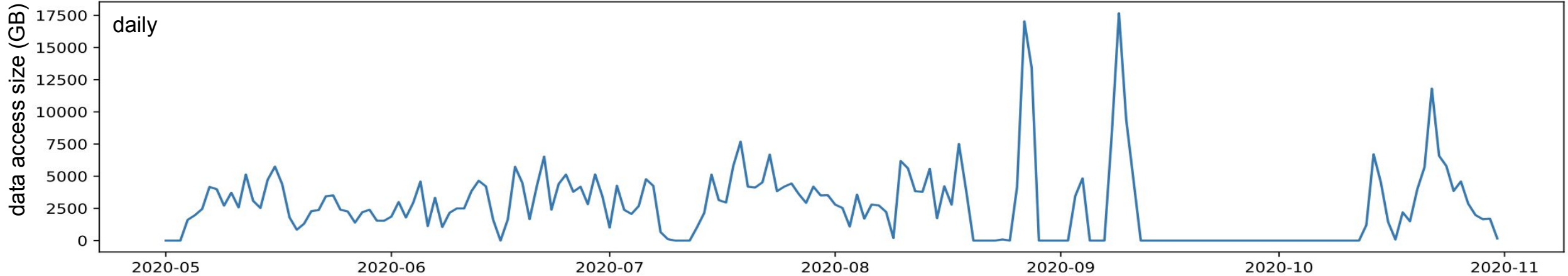
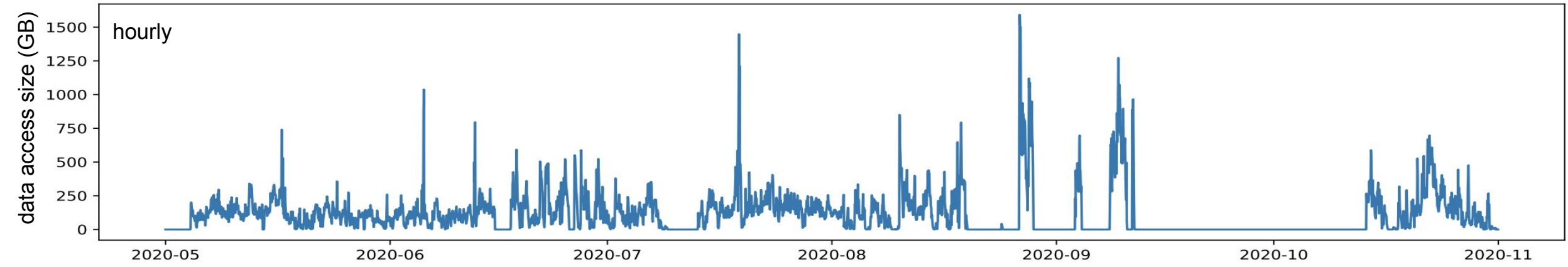
Number of data accesses over time (daily)



- Monitoring issues from Sep to Dec that logs from many days were lost
- Number of total data accesses per day during 5/2020-10/2020, total count=1,286,748



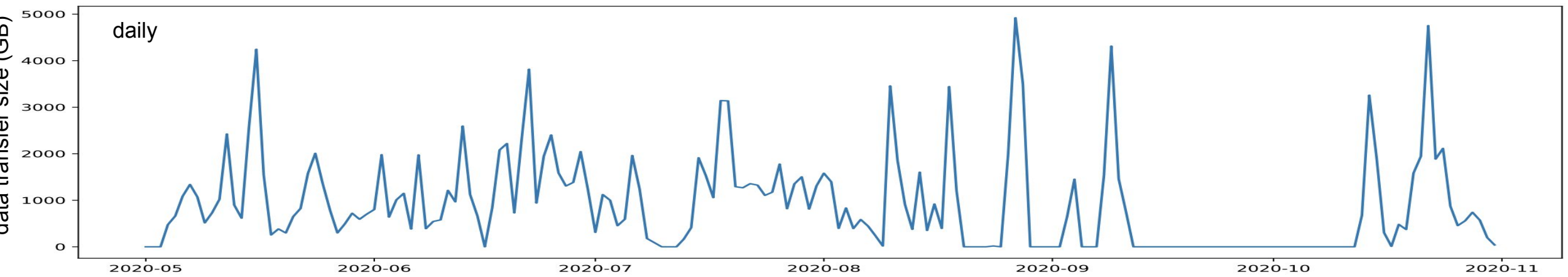
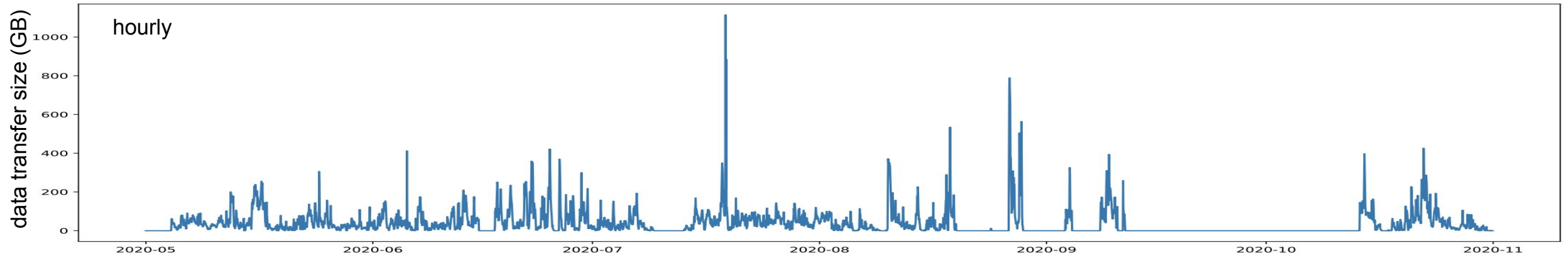
Total data access size over time - from Xcache to applications



- Total data access size (per hour, per day) during 5/2020-10/2020, total size=490.831TB
- Total data access includes the first time and repeated accesses
- **Total reads from cache**



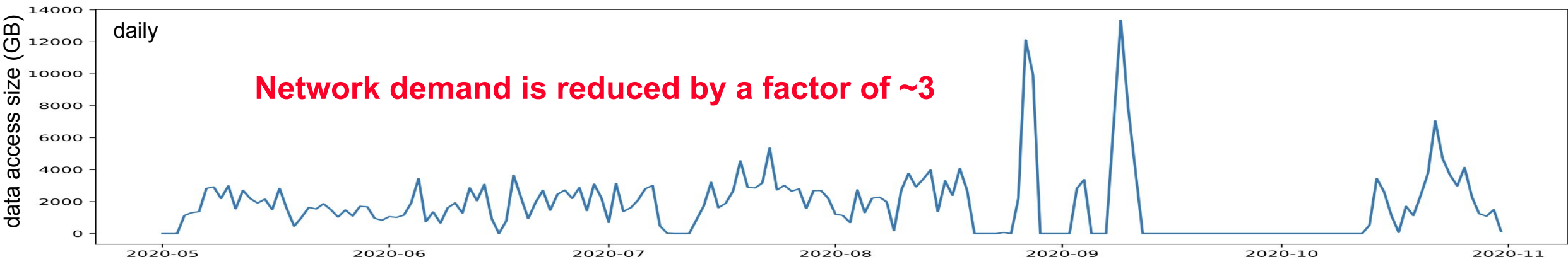
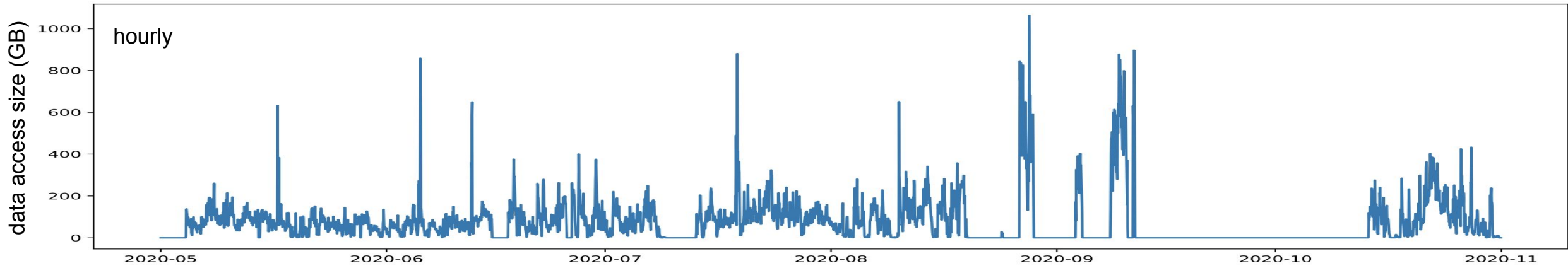
Data transfer size over time - from remote sites to Xcache



- Data transfer size (per hour, per day) during 5/2020-10/2020, total size=168.083TB
- Corresponds to **cache misses**



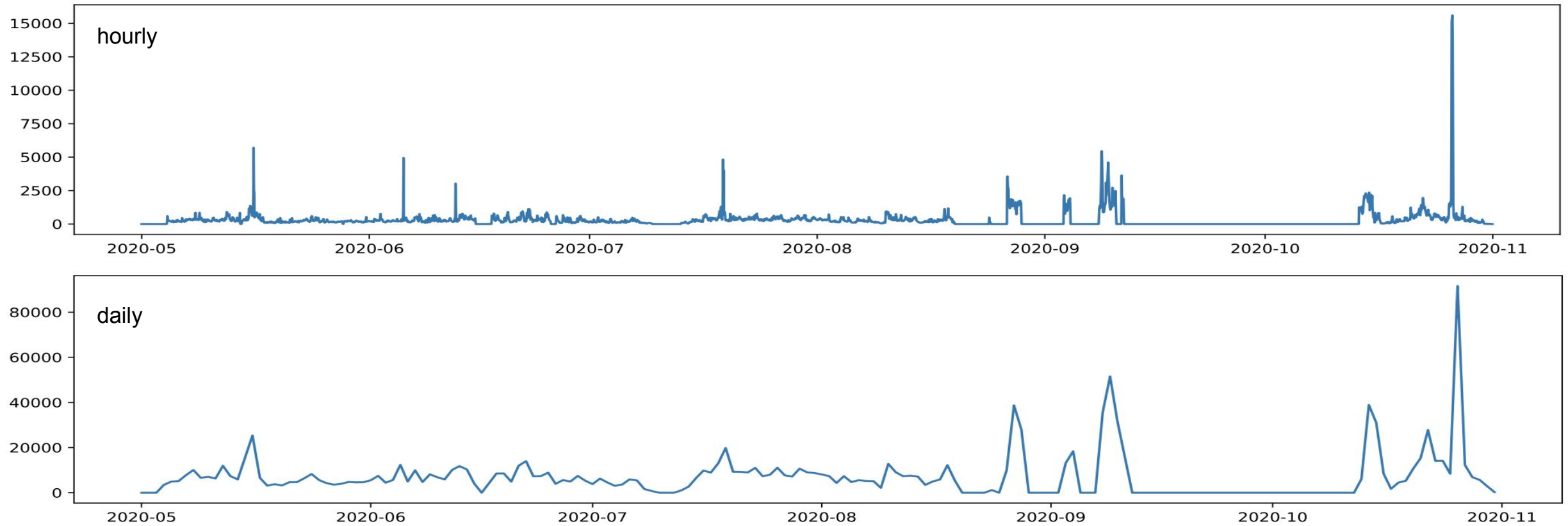
Shared data access size over time - from Xcache to applications



- Shared data access size (per hour, per day) during 5/2020-10/2020, total size=322.748TB
- Shared data access size = network bandwidth savings with only repeated accesses
- Corresponds to **cache hits**



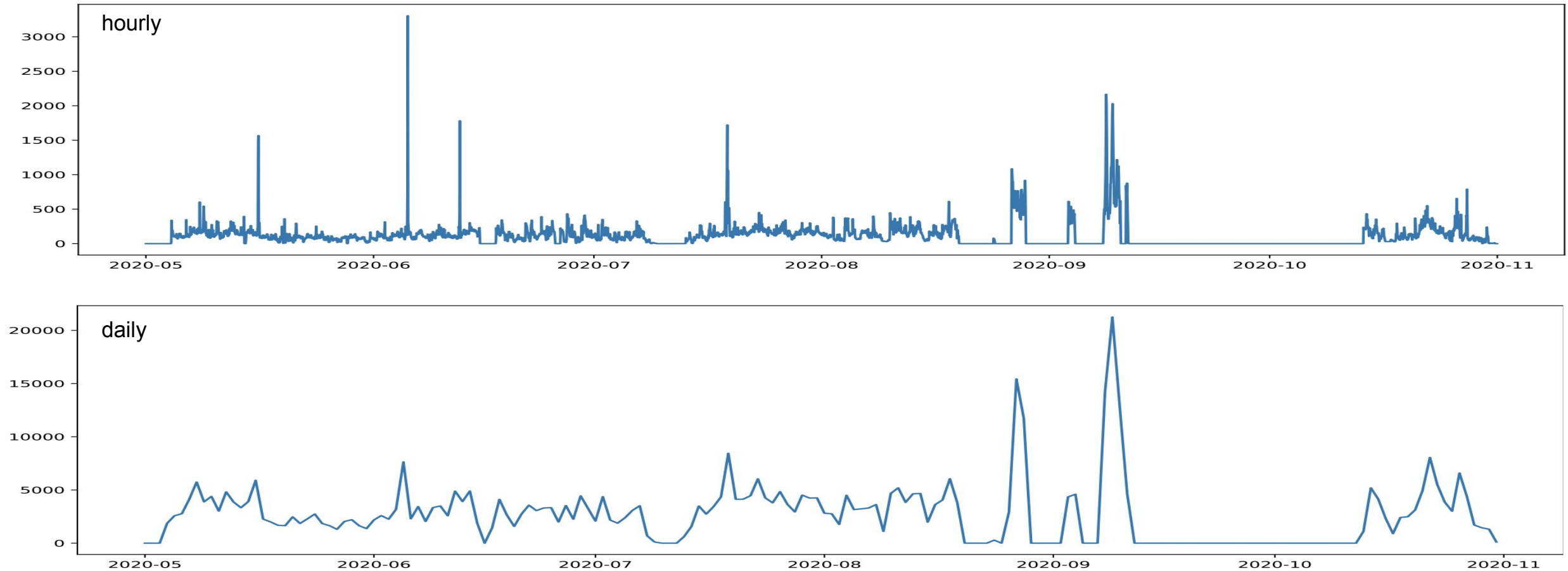
Number of data accesses over time



Number of total data accesses (per hour, per day) during 5/2020-10/2020, total count=1,286,748

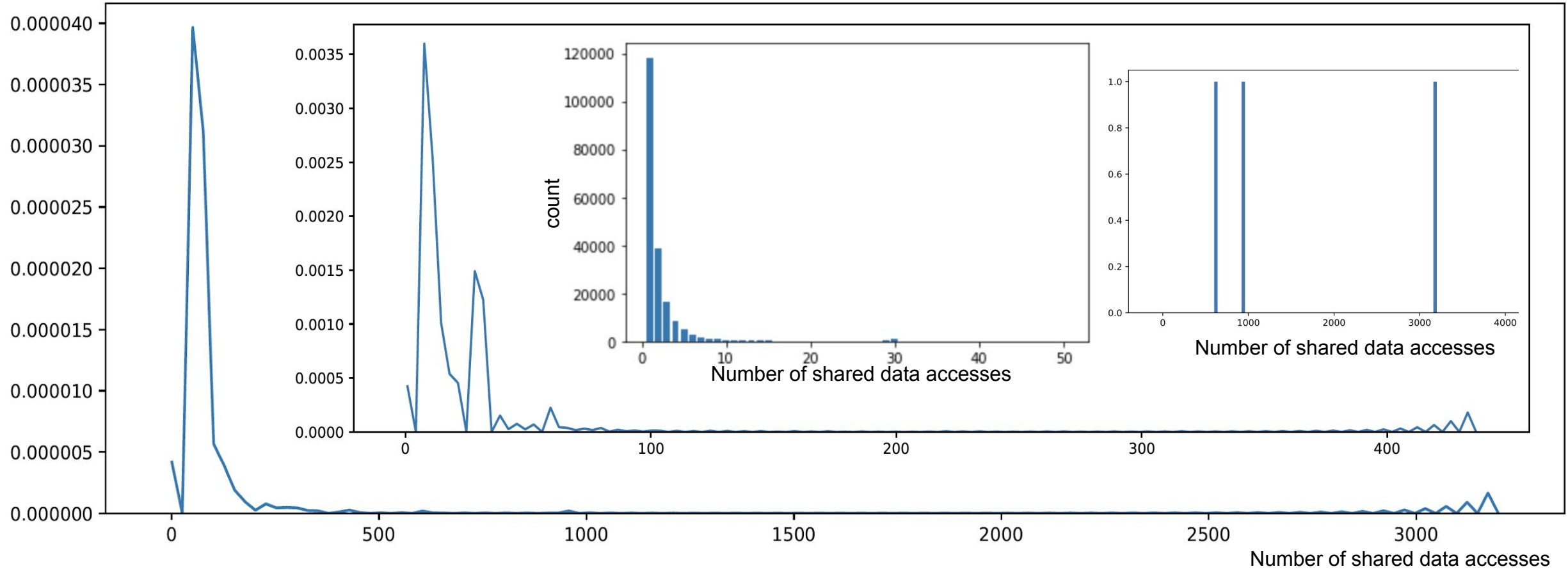


Number of shared data access over time



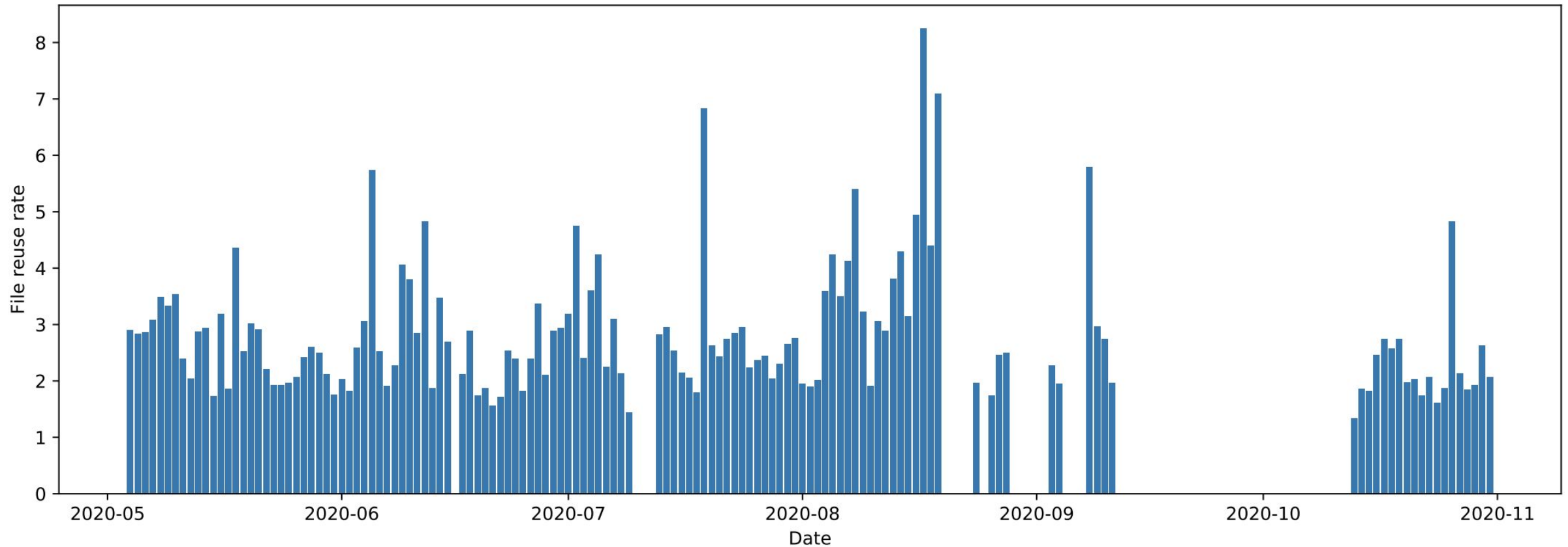
Number of shared reads (per hour, per day) during 5/2020-10/2020, total count=490,944

Distribution of the shared data accesses



- Distribution of the shared data access count during 05/2020-10/2020
total shared access count=490,944, unique file count=198,940
- Distribution of the shared data access count (≤ 500), total shared access count=486,182, unique file count=198,937
- Density plot of shared data access count (> 500), total shared access count=4,762, unique file count=3

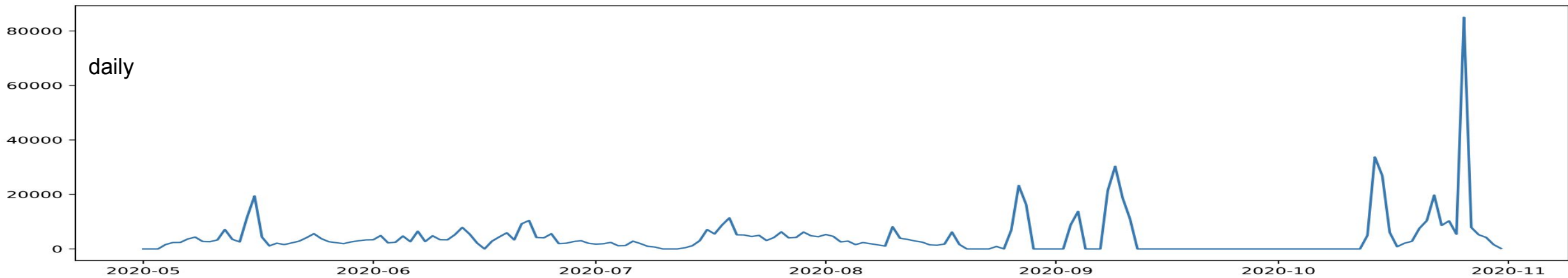
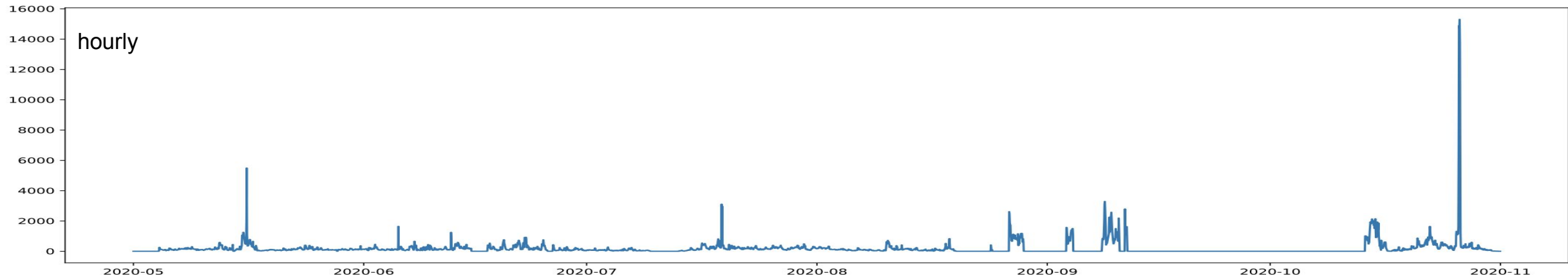
Data file reuse rates (daily)



- Data file reuse rate = (sum of reuses) / (total number of unique files)
 - Sum of reuses = all shared data access counts of the day (cache hits)
 - Total number of unique files = number of unique files for the day

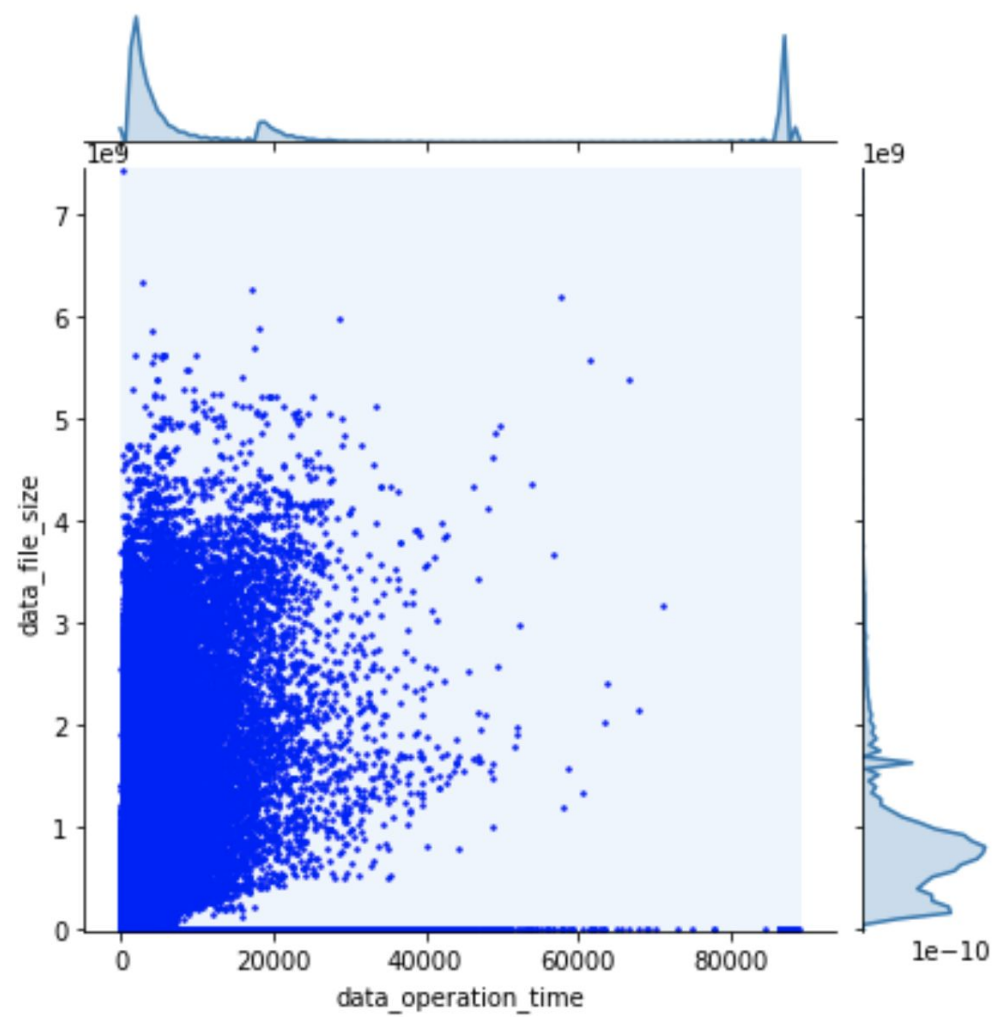


Number of data transfers over time

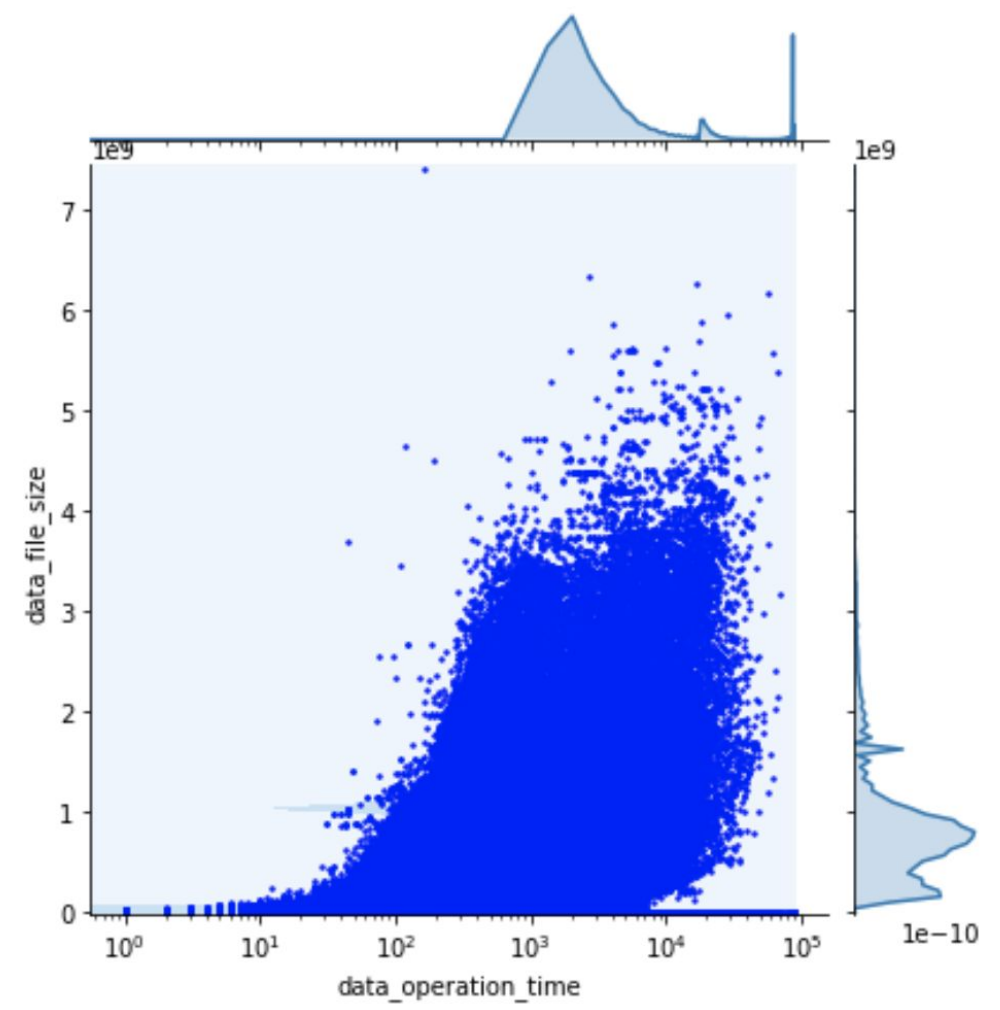


Number of data transfers (per hour, per day) during 5/2020-10/2020, total count=795,804

Transfer duration vs data file size



Transfer Size (bytes) vs. Duration (sec), 05/2020-10/2020



Transfer Size (bytes) vs. Duration (log(sec))

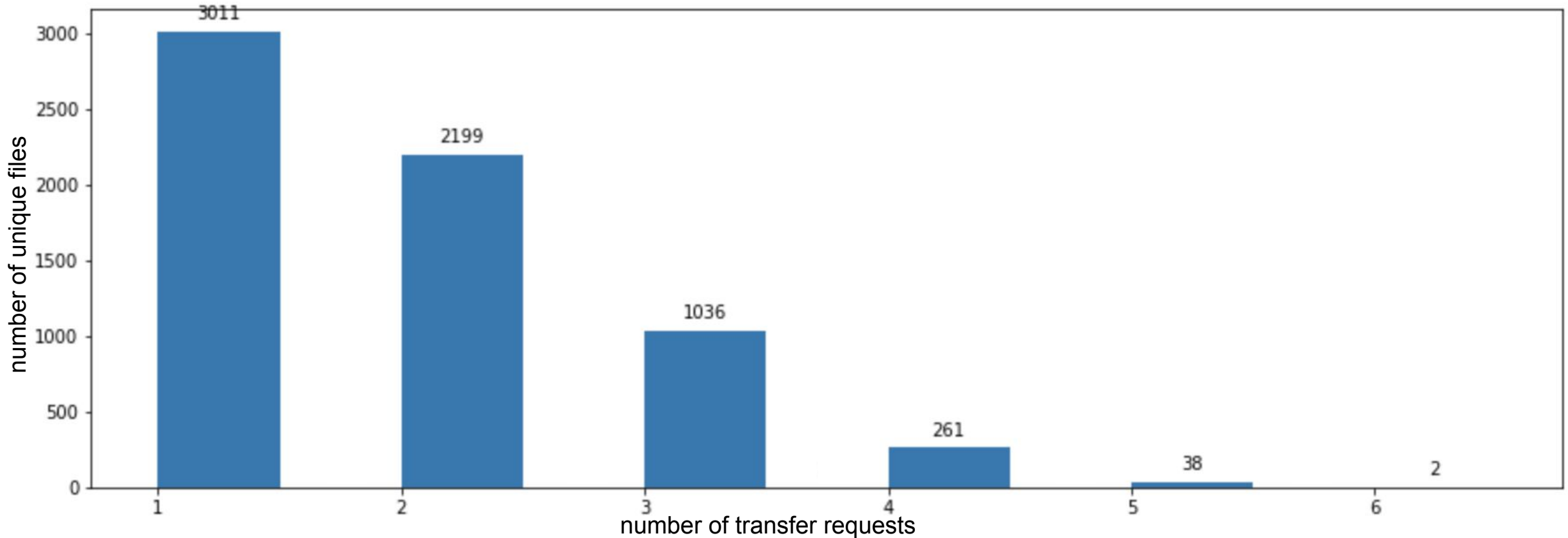


Notes on data transfer spike in 10/2020

	data_user	data_transfer	data_zero_operation_time	metadata_id	data_file_lfn			data_read_bytes		data_file_size		data_operation_time
					count	nunique	count	max	sum	max	sum	mean
0	Alan Malta Rodrigues	False	False	1	1	1	775856549	775856549	775856549	775856549	325.000000	
1	Alan Malta Rodrigues	True	False	1	1	1	0	0	729145901	729145901	621.000000	
2	Alejandro Gomez Espinosa	True	False	8	8	8	0	0	10921350	20772110	83.500000	
3	Alejandro Gomez Espinosa	True	True	1	1	1	0	0	0	0	0.000000	
4	Antoine Lesauvage	True	False	14	14	14	0	0	459441739	726622941	270.285714	
5	Antoine Lesauvage	True	True	2	2	2	0	0	0	0	0.000000	
6	Daina Dirmaite	False	False	45	45	45	13723775	433735867	13723775	433735867	12.444444	
7	Daina Dirmaite	True	False	11763	6547	11763	0	0	14537889	19644221225	6.160928	
8	Daina Dirmaite	True	True	3295	2757	3295	0	0	0	0	0.000000	
9	Devdatta Majumder	True	False	1	1	1	0	0	0	0	1152.000000	

- 12pm-1pm, October 26, 2020, 15058 transfer operations for one user
- 3295 records has 0 file transfer size and 0 transfer time.
- For the other 11,763 transfer records, 6547 unique files are requested.

Notes on data transfer spike in 10/2020



- 15058 file transfer operations, 12pm-1pm, October 26, 2020 for one user
- 3295 records has 0 file transfer size and 0 transfer time.
- From 11,763 transfer records, 6547 unique files.
- Total transferred file size: 19.64 GB

Summary

- **Demonstrated the capability of a network-based temporary data cache**
- **Shared data caching mechanism**
 - **Reduced the redundant data transfers, saved network traffic volume**
 - **Summary of the 1,286,748 accesses from May 2020 to Oct 2020**
 - **Total 490.831 TB of client data access (first time reads and repeated reads)**
 - **Transferred/cached 168.08 TB (from remote sites to cache)**
 - **Saved 322.748 TB of network traffic volume (repeated reads only)**
 - **Network demand reduced by a factor of ~3**
- **Further studies**
 - **Cache miss rates**
 - **How caches affect each other when one or more of the federated caches are down**
 - **How many time a file needs to be retrieved from remote sites?**
 - **How are the cache misses affecting the application performance?**
 - **Regional cache impacting application performance (local vs remote data access)**
 - **Cache utilization**
 - **How many Xcache installations are good enough?**
 - **What size of each disk cache would be appropriate?**
 - **If the number of physicists using the system doubles, how many more cache deployments are needed?**