

Collecting and Processing Earth Science Data Metrics at NASA ESDIS

Jianfu Pan¹, Nelson Casiano¹, Stephan Klene¹, James Smith¹, Lalit Wanchoo², Durga Kafle²

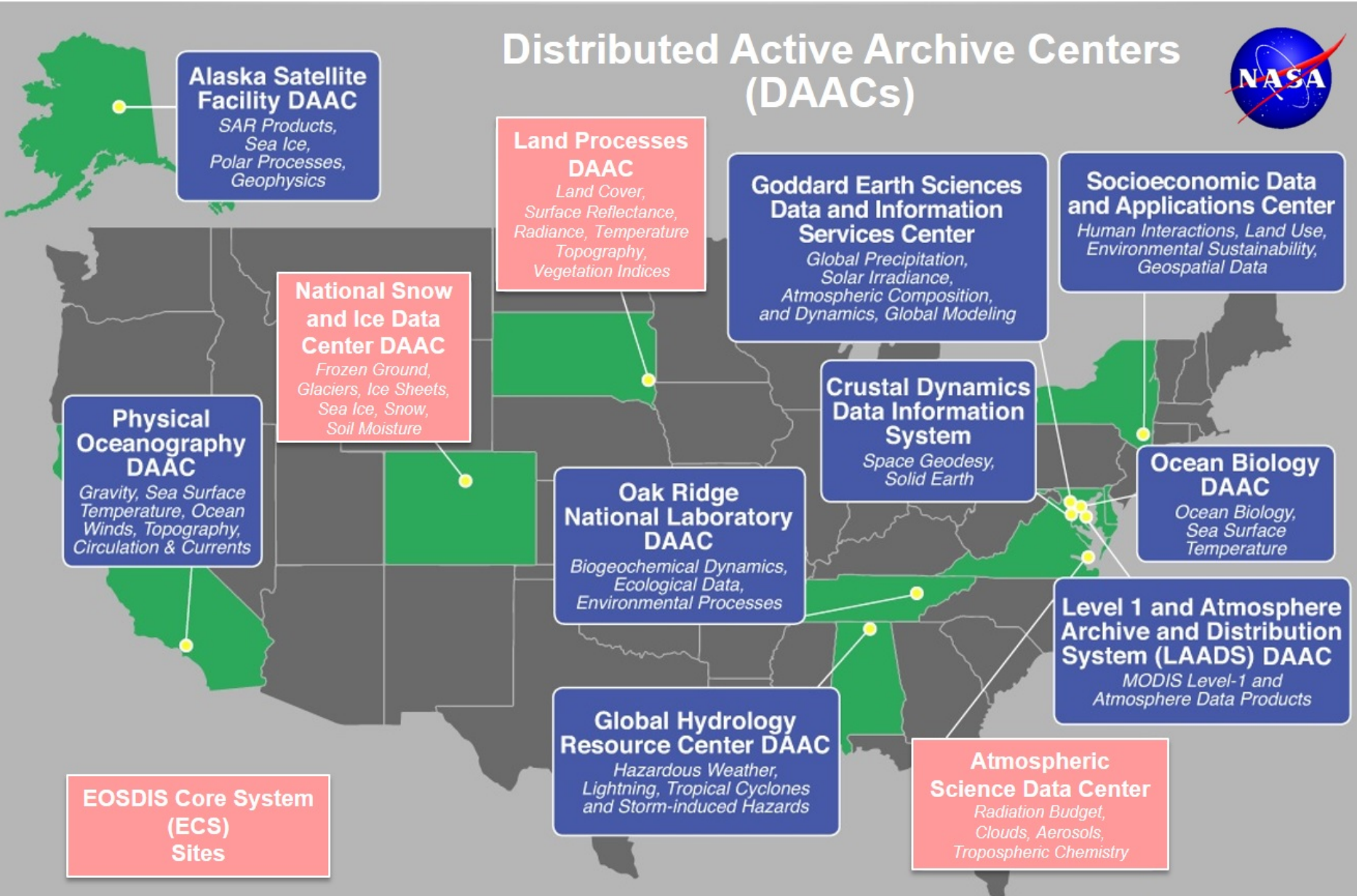
¹ SSAI, Inc., Lanham, MD 20706; ² Adnet Systems, Inc., Lanham, MD 20706

INTRODUCTION

NASA's Earth Science Data and Information System (ESDIS) project manages the science systems of the Earth Observing System Data and Information System (EOSDIS), which is responsible for processing, archiving, and distributing a wide range of Earth Science data products. These data activities occur at twelve EOSDIS Distributed Active Archive Centers (DAACs) located throughout the United States (Figure 1), with each DAAC focusing on specific science discipline areas within the Earth Sciences. DAACs receive and/or produce data products of various missions (satellite, aircraft, field campaign, model, in-situ, and others), then ingest them into their product system. Science data products are then distributed to end users. DAACs are also required to provide metrics or activity logs of ingest, archive and distribution to ESDIS Metrics System (EMS).

EMS collects information on ingest, archive and distribution, along with product metadata from all DAACs. EMS daily processing converts the information received into various metrics. The metrics help both ESDIS and DAAC management in resource planning, understanding user behaviors, and identifying important data products that are popular to user communities.

Figure 1. NASA EOSDIS Distributed Active Archive Centers (DAACs).
(source: <https://earthdata.nasa.gov/eosdis/daacs>)



EMS PROCESSING SYSTEM

Each day, DAACs collect logs of ingest, archive, and distribution of science data products and send them to EMS for processing. One example of the distribution logs is the Apache log produced by an Apache web server where end users access and download the data. Other forms of records include what is called customer logs, which are constructed by DAACs to include selected metrics fields before being sent to EMS.

Figure 2 illustrates the architecture of the EMS log processing system. DAACs send logs collected to a dedicated file server serving as an interface between DAACs and EMS. EMS picks up logs and other files from the file server for processing. The interface file server is also used by EMS to store certain specialized reports for DAACs, while more standardized reporting is done through the Oracle Application Express (Apex/HTMLDB) platform, a web-based Oracle application and reports development tool. All processed metrics are permanently stored in the EMS database. The Apex/HTMLDB platform is used to make metrics reports, and is accessible by all DAACs.

EMS processing runs two instances in a semi-parallel fashion, with steps involving database processing can only be single-threaded. In addition, it also uses an independent preprocessing module, called pre-staging, to preprocess files before being further processed by the EMS workflow. Within an EMS processing instance, a typical workflow consists of the steps in Table 1, executed in that order.

Table 1. Workflow steps in EMS processing system

▪ <i>GetFiles</i>	Gets log files stored on the interface server for processing, normally in batches.
▪ <i>StageFiles</i>	First major processing step for hostname and IP geolocation lookups.
▪ <i>LoadFiles</i>	Wrapper to Oracle processing for validation, user profiles lookup, and standardization of field values. Creates output files on the database server.
▪ <i>FileXfer</i>	Transfers outputs created on the database server to local machine for further processing.
▪ <i>ProductSearch</i>	Maps distribution records to a science products for product usage metrics.
▪ <i>DistSummary</i>	Consolidates distribution records and saves the final metrics results to the EMSdatabase.



Data Providers (DAACs)

- Data ingest/Archive & Archive delete/Distribution logs
- User profiles
- Product metadata, product search terms

Data Provider – EMS Interface File Server (fs1)

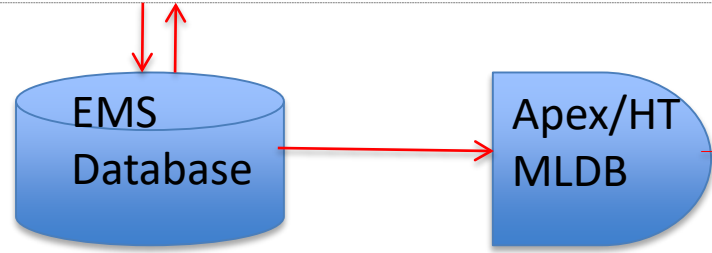
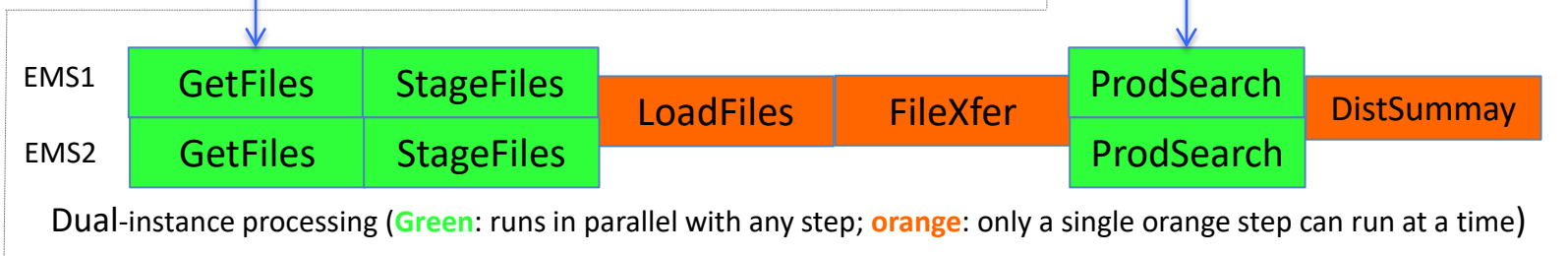
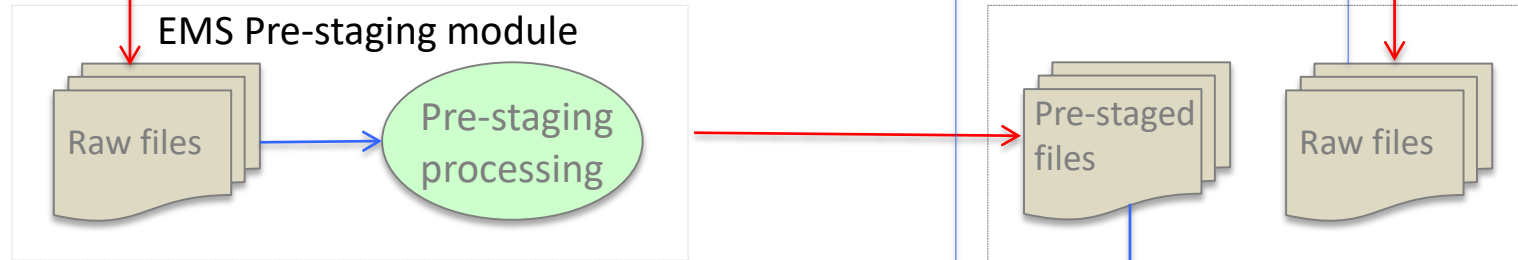
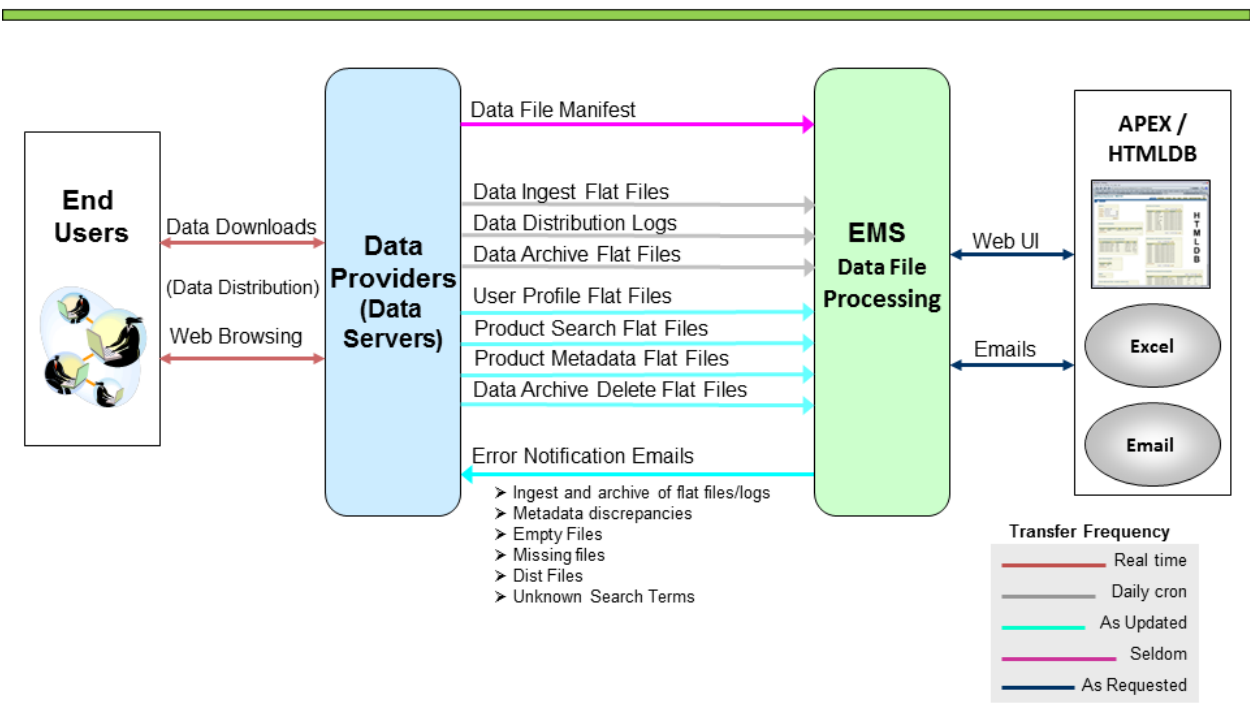


Figure 2. Architecture of EMS processing system

Reports and notifications

Figure 3 is a schematic diagram illustrating communication between DAACs, also known as Data Providers, and EMS. DAACs send all logs and other necessary informational files, such as product metadata, to EMS as flat files. These are plain text files with delimited fields that are normally pre-defined and registered in EMS. The fields and formats follow a mutually agreed upon Interface Control Document (ICD) between ESDIS and DAACs. While flat file format is used to exchange logs and file-based reports, Data Providers use the APEX/HTMLDB reporting tool to create metrics reports and monitor file processing status (Figure 4). From time to time, Data Providers also receive email notifications from EMS for announcements and processing errors if any.

Figure 3. Communication between DAACs (Data Providers) and EMS



Summary

Options

Provider: ALL

Startdate: 2019-10-31

Enddate: 2019-11-07

Ancillary: Include Ancillary Data

Query

EMS Processing Status

Date	Processing Status
07-NOV-19	PROCESSING

1 - 1

Successful Distributions By Provider to Public Users

1 - 24

Provider ↑	MBs	Files	Granules	Orders	Requests	Hosts
AMSR2NRT	30,029,892	11,662	0	0	0	41
AMSR2NRT2	574,547	147	0	0	0	49
ASF	37,303,146,918	717,672	0	0	0	3,775
CDDIS	6,136,248,240	33,104,920	0	0	0	14,992
GESDISC	91,764,853,210	10,780,047	0	0	0	4,288
GESDISCNRT	700,560,672	200,554	0	0	0	200
GHRC	438,383,600	1,594,428	0	0	0	1,029
ISSLISNRT	2,573,450	14,498	0	0	0	18
LARCANGE	531,063,096	3,311	0	3,311	0	1
LARCECS	26,236,675,669	109,379	109,376	2,341	2,341	29
LARCNRT	42,218	444	444	222	222	1
LARCNRT2	6,078	96	96	0	0	1
LARCSVC	20,327,380,883	164,477	0	495	0	41
LPDAAC	73,873,982,125	5,259,468	5,184,680	13,019	12,718	18,149
MODAPS	235,259,108,902	6,650,671	0	0	0	7,352
MODAPSNRT	18,078,651,390	1,436,163	0	0	0	11,446
MOPITTNRT	18,652,667	1,204	0	0	0	11
NSIDC	18,830,545,487	2,115,235	1,932,583	2,289	6,238	576
NSIDCV0	268,655,013	148,974	0	0	0	179
OBDAAC	25,721,067,821	597,925	0	0	0	1,812

Daily File Processing Missing Files Provider Effective Dates

Daily File Processing

Query

Provider: GESDISC

Rangetype: By Filename Date

Filetype: ALL

Datasource: ALL

Startdate: 2019-10-31

Enddate: 2019-11-07

Query

Daily File Processing

row(s) 1 - 31 of 78 [Next](#)

Provider ↑	Filename	Processing Date	Total Recs	Failed Recs
GESDISC	20191031_GESDISC_ArchDel_s4pa.fit	2019-11-01 08:51:39 AM	13,968	0
GESDISC	20191031_GESDISC_Arch_s4pa.fit	2019-11-01 08:51:11 AM	33,855	0
GESDISC	20191031_GESDISC_DistCustom_Push.fit	2019-11-01 08:52:42 AM	8,601	0
GESDISC	20191031_GESDISC_DistHTTP_httpairs.fit	2019-11-01 01:37:02 AM	582,428	0
GESDISC	20191031_GESDISC_DistHTTP_httpaura.fit	2019-11-01 02:20:13 AM	507,734	0
GESDISC	20191031_GESDISC_DistHTTP_httpdisc.fit	2019-11-01 03:23:10 AM	694,031	0
GESDISC	20191031_GESDISC_DistHTTP_httpreason.fit	2019-11-01 03:26:34 AM	396,802	0
GESDISC	20191031_GESDISC_Ing_s4pa.fit	2019-11-01 08:50:57 AM	20,699	0
GESDISC	20191031_GESDISC_Meta_s4pa.fit	2019-11-01 12:35:37 AM	4,577	0
GESDISC	20191031_GESDISC_UsrProf_s4pa.fit	2019-11-01 12:36:16 AM	1,407	0
GESDISC	20191031_GESDISC_searchExp_s4pa.fit	2019-11-01 12:35:46 AM	17,040	0
GESDISC	20191101_GESDISC_ArchDel_s4pa.fit	2019-11-02 08:51:15 AM	13,169	0
GESDISC	20191101_GESDISC_Arch_s4pa.fit	2019-11-02 08:50:46 AM	29,887	0
GESDISC	20191101_GESDISC_DistCustom_Push.fit	2019-11-02 08:52:04 AM	8,489	0
GESDISC	20191101_GESDISC_DistHTTP_httpairs.fit	2019-11-02 01:20:16 AM	460,091	0
GESDISC	20191101_GESDISC_DistHTTP_httpaura.fit	2019-11-02 02:43:42 AM	547,755	0
GESDISC	20191101_GESDISC_DistHTTP_httpdisc.fit	2019-11-02 03:32:05 AM	604,104	0

Figure 4. Screenshots of APEX/HTMLDB reporting tool.

EVOLUTION OF EMS PROCESSING SYSTEM

Early version of the EMS processing system used a simple pass-through multi-step processing workflow of the steps described earlier. With the increase in the data provided by the data providers over years, the system took significant time to process all of the data on daily basis and required to be upgraded to improve the efficiency of processing. Figure 5 shows monthly total number of records processed by EMS over years, for ingest, archive, and distribution. It seems that the amount of data distributed doubles every 3 to 4 years.

The first major upgrade introduced parallel processing, taking advantages of the two most time consuming steps, StageFiles and LoadFiles, being run on separate servers with StageFiles on an application server and most work in loadFiles being database processing on a database server. Without the parallelization, one server would be idle for nearly half of the time, waiting for the other step to complete before moving along in the workflow. More recently, we also implemented a prestaging process that can run either independently or as part of the EMS processing workflow. The system architecture with the prestaging is illustrated in Figure 6. The prestaging module preprocesses files and can run independently. The preprocessed files can then skip several time-consuming steps when being merged into the EMS process work flow.

These system upgrades have greatly improved the processing performance. As shown in Figure 7, EMS was running above 20 hours each day for its daily processing before around 2015 for several years, making the system highly stressed and very challenging for down times needed by hardware maintenance and occasional file reprocessing. The upgraded system has been steadily running with about 10 hours of processing time, with just three largest files are preprocessed.

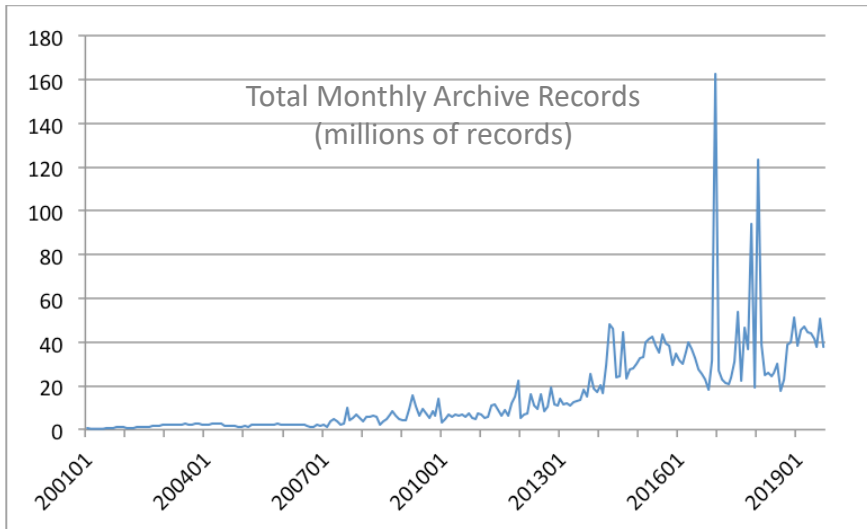
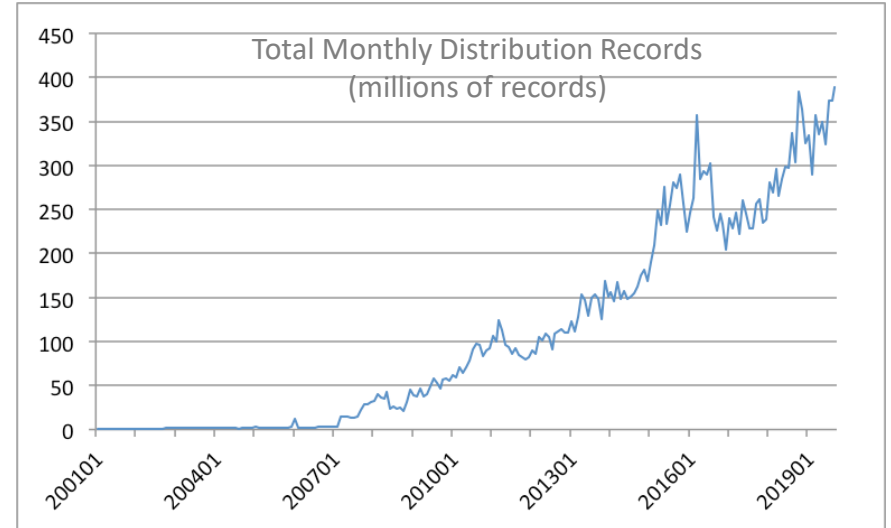
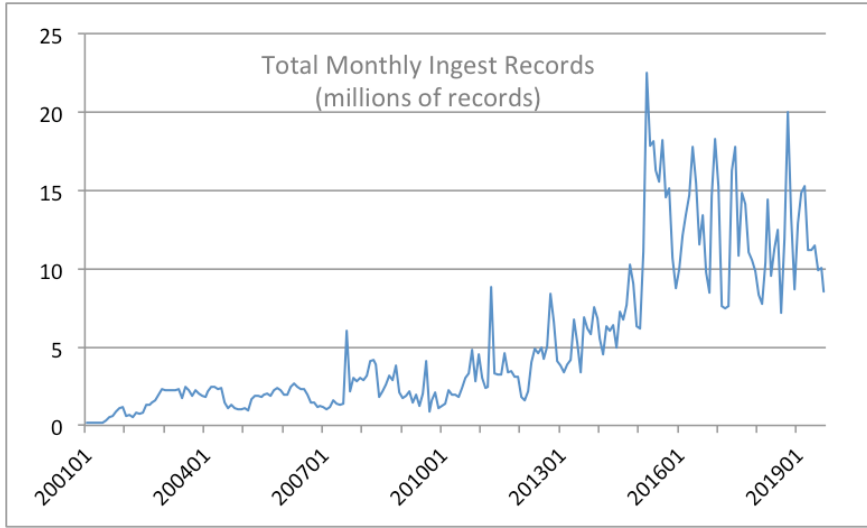


Figure 5. Monthly total number of records processed by EMS since January 2001 (in millions of records)

Figure 6. EMS with Prestaging (prestaging workflow in red). Prestaged results are merged into EMS main processing workflow to achieve final metrics results.

EMS Architecture with Pre-staging Subsystem

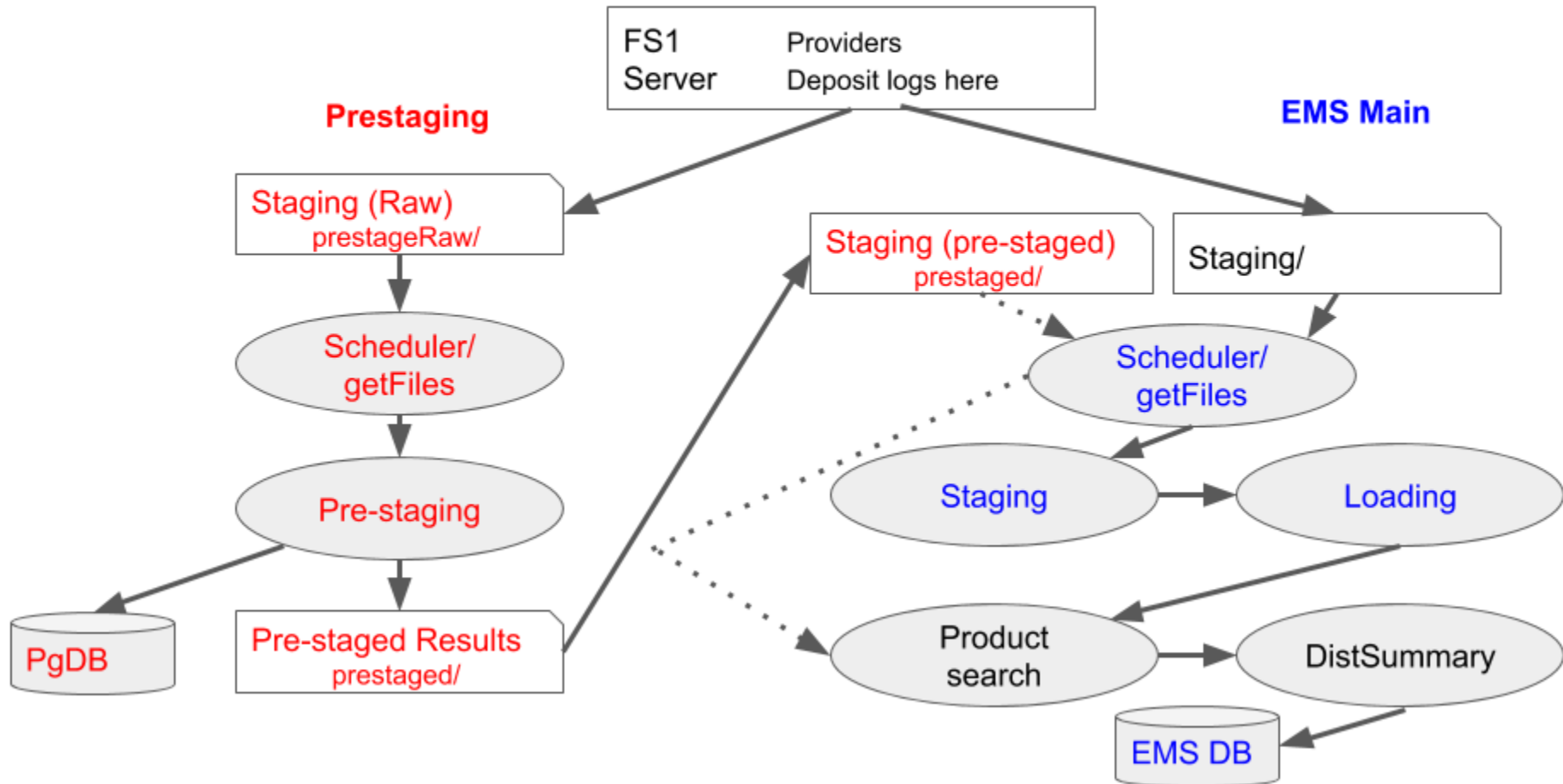
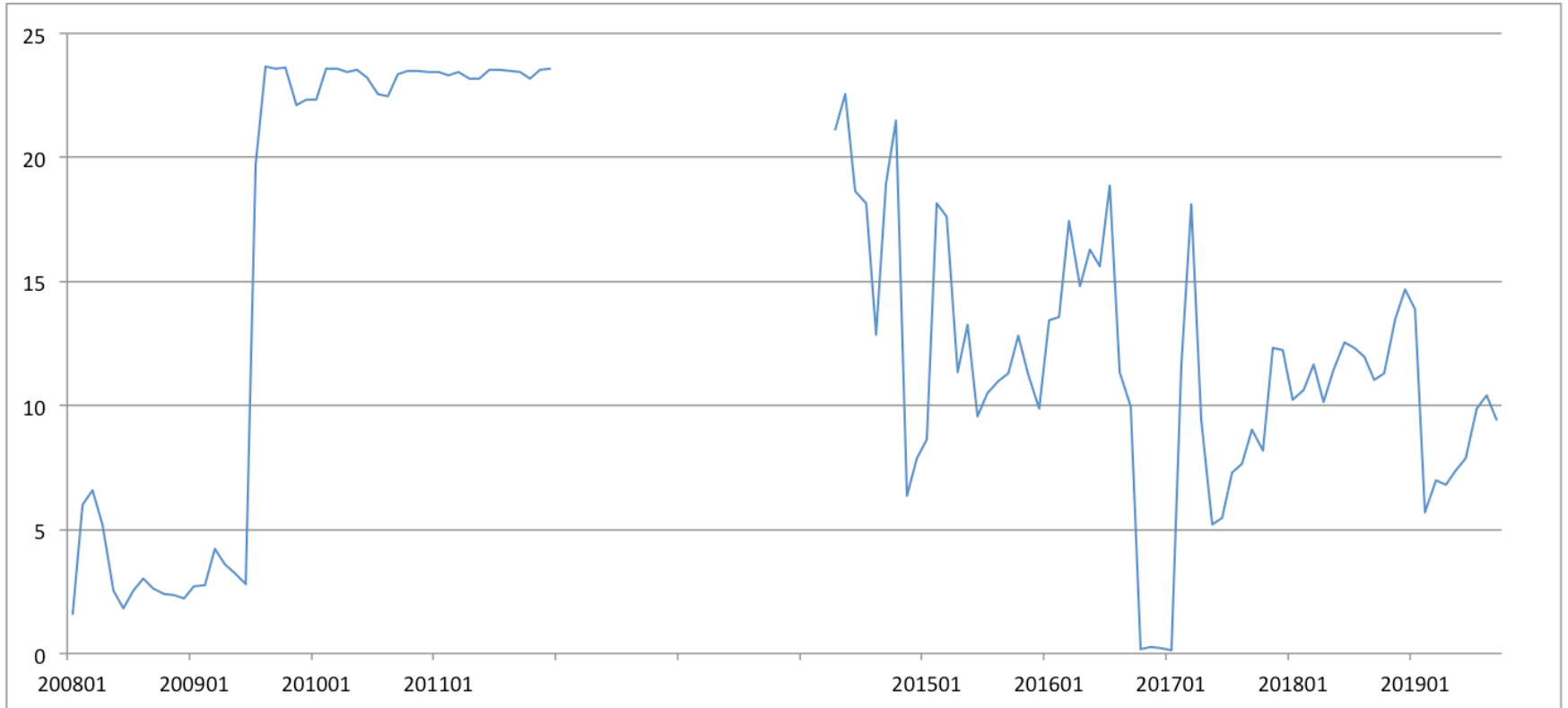


Figure 7. EMS monthly median system daily processing time since January 2008 (hr)
(gap between 2012 and 2014 due to loss of records from system hardware failure)



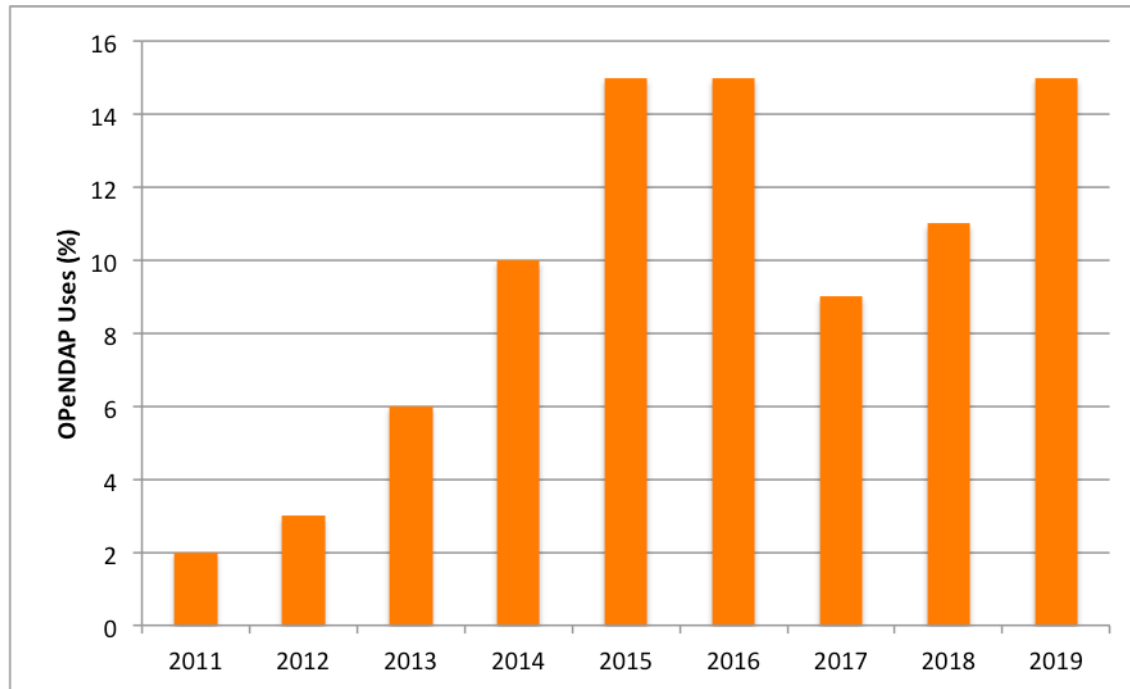
EMS Metrics

EMS started tracking metrics of ingest, archive and distribution of Earth Science data since around 2000. Over the years, DAACs provided various data services making users easier to get to the data, such as OPeNDAP as a data access protocol allowing users to download subset of data, both in geographical region and variables within a data product. EMS incorporated collection of such services protocol in the data provided by the DAACs. Figure 8 shows such OPeNDAP distribution by DAACs to the data users. Table 2 summarizes different categories of metrics that EMS produces.

Table 2. List of typical EMS metrics

Metrics Category	Metrics Item	Description
Data Usage	Volume in bytes	Amount of data distributed to user
	Number of files	Number of files or granules distributed to user
Data Product	Mission	Science mission of the data product
	Instrument	Instrument used to acquire the data measurements
User Profile	Contact email	Occasionally used for notifications from data producers
	Affiliation	User's categorical affiliation such as Government, Education, etc.
	Primary study area	User's primary field of science researches
	Country	Country where user resides
Protocol & Service	FTP / SCP	Traditional file download through FTP, SCP, and alike
	HTTP	Data download through web HTTP protocol
	OPeNDAP	Data download through OPeNDAP protocol
	Giovanni	GES DISC data analysis system that provides analysis results
	Subsetter	Various subset tools mostly for spatial subsetting
	Reformat	Convert data files from one format to another
	Reprojection	Re-project data grids to a different geographic coordinate system

Figure 8. Percentages of OPeNDAP data accesses over years among all successful data distributions.



SUMMARY AND CONCLUSIONS

EMS collects and processes records of ingest, archive and distribution logs of Earth Science data products at NASA DAACs. EMS processing system has been evolving over time to meet the needs to handle ever-increasing amount of logs to process. With the implementation of parallelization and preprocessing module, the performance of EMS system has much improved. Moving forward, EMS will continue to grow to meet the needs to process the growing amount of log data, as well as to handle metrics from data distribution through new technologies such as OPeNDAP and cloud-based data distribution. In near future, EMS is looking to interact with NASA's Common Metadata Repository (CMR) for product metadata. As data distribution is moving to cloud, work is also underway to process metrics from the cloud-based data distribution.

ACKNOWLEDGEMENTS

This study was funded through the Software Engineering Support (SES II) GSFC NASA Contractor No. NNG15CR67C-6002