

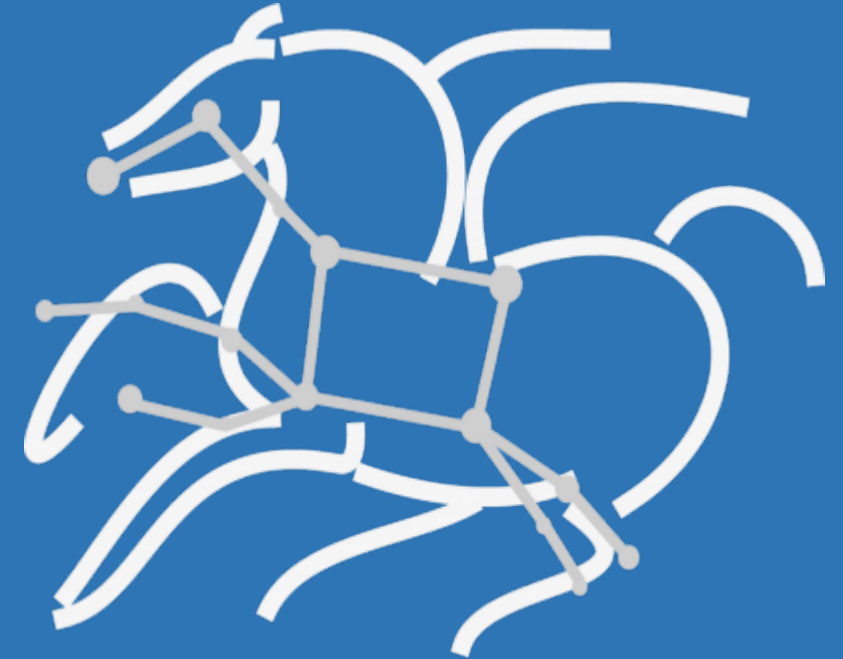


U.S. DEPARTMENT OF
ENERGY



End to End Workflow Monitoring Panorama 360

George Papadimitriou
georgpap@isi.edu



USC Viterbi
School of Engineering
Information Sciences Institute

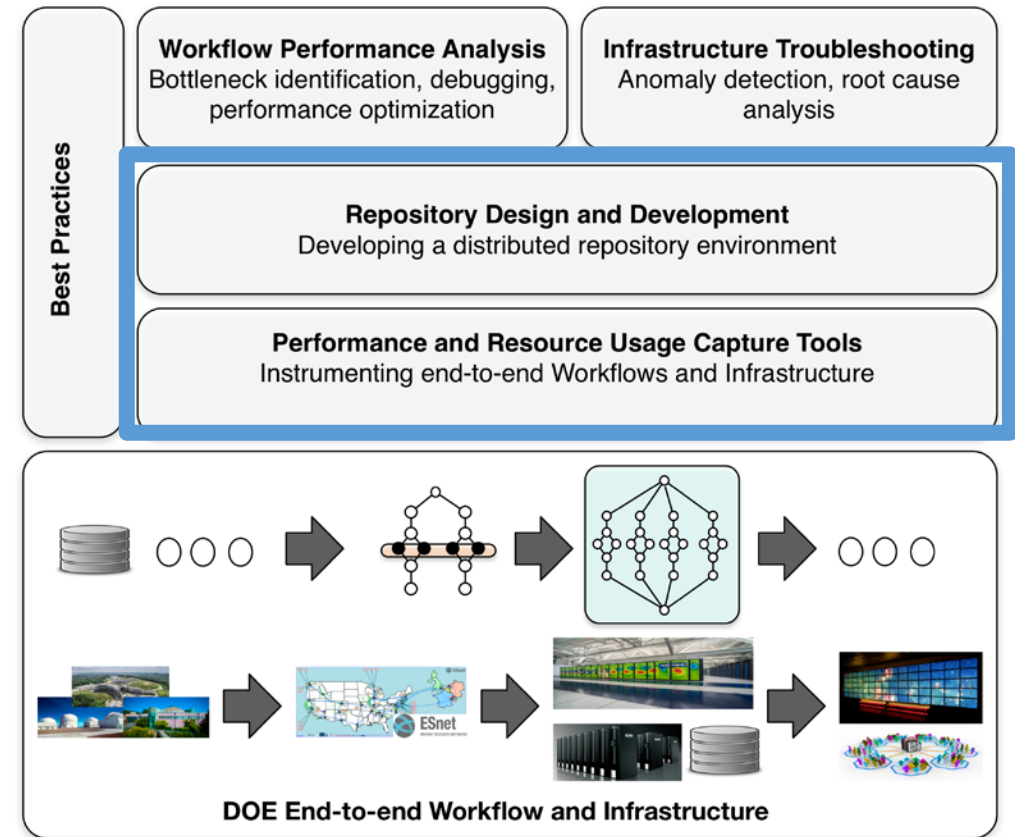
<http://pegasus.isi.edu>

Panorama 360: Project Overview

- Leverage the Pegasus WMS to structure, execute and monitor workflow execution
- Characterize performance: instrument data capture, summarize data, and publish results
- Create an open access **common repository** for storing end-to-end workflow performance and resource data captured using a variety of tools

**Open for external contributors*

- Apply and develop ML techniques for **workflow performance analysis** and **infrastructure troubleshooting**
- Record findings, distill **best practices**, and share and refine them with other program teams



Panorama 360 Overview.

Data Sources: Application and Infrastructure

- **Pegasus Stampede** events regarding the workflow and its status
- **Pegasus Kickstart Online** collects resource usage traces with frequency as low as 1 second in real-time
- **Darshan** collects file access statistics (eg. POSIX, MPI-IO) during the execution
- **Globus** collects transfer statistics and general information about the transfer (throughput, file transfer errors etc.)



DARSHAN
HPC I/O Characterization Tool



Pegasus Kickstart Online Traces

- In the Panorama branch `pegasus-kickstart` supports fine-grained monitoring capabilities by invoking a helper tool called `pegasus-monitor`.
- `pegasus-monitor` can pull resource usage statistics for workflow running tasks within a predefined time interval.
- Minimum interval is limited to **1 second**
- The statistics are being read from the `/proc` entry of the running task and among them include:
 - Number of processes and threads
 - stime and utime
 - Bytes read and written
 - iowait

```
4  {
5    "site": "condorpool",
6    "@version": "1",
7    "xformation": "individuals",
8    "bread": 0,
9    "ts": 1552869746.009779,
10   "rss": 25004,
11   "procs": 5,
12   "pid": 29304,
13   "bwrite": 74547200,
14   "threads": 6,
15   "exe": "/var/lib/condor/execute/dir_29120/pegasus-4.9.2panorama/bin/pegasus-kickstart",
16   "iowait": 0,
17   "vm": 311864,
18   "brecv": 0,
19   "bsend": 0,
20   "event": "kickstart.inv.online",
21   "wf_uuid": "0a114f53-f208-44da-9b5c-9c30edcf49c6",
22   "dag_job_id": "individuals_ID0000001",
23   "utime": 8.4,
24   "rank": 0,
25   "task_id": "ID0000001",
26   "syscr": 1088,
27   "hostname": "workers1-1",
28   "rchar": 1376495,
29   "wchar": 7964288,
30   "condor_job_id": "14988.0",
31   "syscw": 47652,
32   "stime": 1.57,
33   "type": "kickstart",
34   "wf_label": "1000genome-20190317T235728Z",
35   "@timestamp": "2019-03-18T00:42:27.355Z"
36 }
```

Darshan

- Darshan is an HPC lightweight application-level I/O profiling tool that captures statistics on the behavior of HPC I/O operations.
- It captures data for each file opened by the application, including I/O operation counts, common I/O access sizes, cumulative timers, etc.
- I/O behavior is captured for POSIX IO, MPI-IO, HDF5, and Parallel netCDF data interface layers.
- It also captures a set of job-level characteristics, such as the number of application processes, the job's start and end times, and the job unique identification provided by the scheduler.
- *In Panorama we only expose accumulated performance data from the STDIO and POSIX modules*
- Reference: <https://www.mcs.anl.gov/research/projects/darshan/>

```
4 {
5   "@version": "1",
6   "darshan_log_version": "3.10",
7   "STDIO_module_data": {
8     "agg_perf_by_slowest": 0.94878,
9     "shared_files": {
10       "time_by_cumul_io_only": 0.000728,
11       "time_by_open": 0,
12       "time_by_slowest": 0.004884,
13       "time_by_open_lastio": 62.774539,
14       "time_by_cumul_meta_only": 0
15     },
16     "total_bytes": 848131,
17     "agg_perf_by_open": 0.954247,
18     "agg_perf_by_open_lastio": 0.012713,
19     "unique_files": {
20       "slowest_rank_io_time": 0.847622,
21       "slowest_rank_meta_only_time": 0.004305,
22       "slowest_rank": 0
23     },
24     "agg_perf_by_cumul": 0.953428
25   },
26   "uid": "1003",
27   "@timestamp": "2018-07-18T19:22:36.788Z",
28   "type": "stampede",
29   "start_time": 1531941678,
30   "exe": "/shared/software/NAMD_2.12_Linux-x86_64-MPI/namd2 equilibrate_200.conf",
31   "monitoring_event": "darshan.perf",
32   "end_time_ascii": "Wed Jul 18 19:22:22 2018",
33   "xwf_id": "e8c30de5-6bb2-4b38-a000-2b719e688e38",
34   "end_time": 1531941742,
35   "nprocs": 8,
36   "run_time": 65,
37   "event": "stampede.task.monitoring",
38   "sched_id": "10469.0",
39   "job_id": "namd_ID0000002",
40   "POSIX_module_data": {
41     "agg_perf_by_slowest": 18.470956,
42     "shared_files": {
43       "time_by_cumul_io_only": 0,
44       "time_by_open": 0,
45       "time_by_slowest": 0,
46       "time_by_open_lastio": 0,
47       "time_by_cumul_meta_only": 0
48     },
49     "total_bytes": 403752,
50     "agg_perf_by_open": 18.470956,
51     "agg_perf_by_open_lastio": 18.470956,
52     "unique_files": {
53       "slowest_rank_io_time": 0.020846,
54       "slowest_rank_meta_only_time": 0.019316,
55       "slowest_rank": 0
56     },
57     "agg_perf_by_cumul": 18.470956
58   },
59   "compression_method": "ZLIB",
60   "start_time_ascii": "Wed Jul 18 19:21:18 2018",
61   "job_inst_id": 7,
62   "jobid": "1547",
63   "metadata": {
64     "h": "romio_no_indep_rw=true;cb_nodes=4",
65     "lib_ver": "3.1.6"
66   }
67 }
```

Globus

- Globus is a **research data management service**, built on top of gridftp.
- It can be used to **transfer files** for your own computations or **share files** with the community.
- For every transfer request Globus creates logs containing transfer statistics, such as:
 - *Request and Completion time*
 - *Source and Destination*
 - *Transfer rate*
 - *Number of failures*
- Reference: <https://www.globus.org>

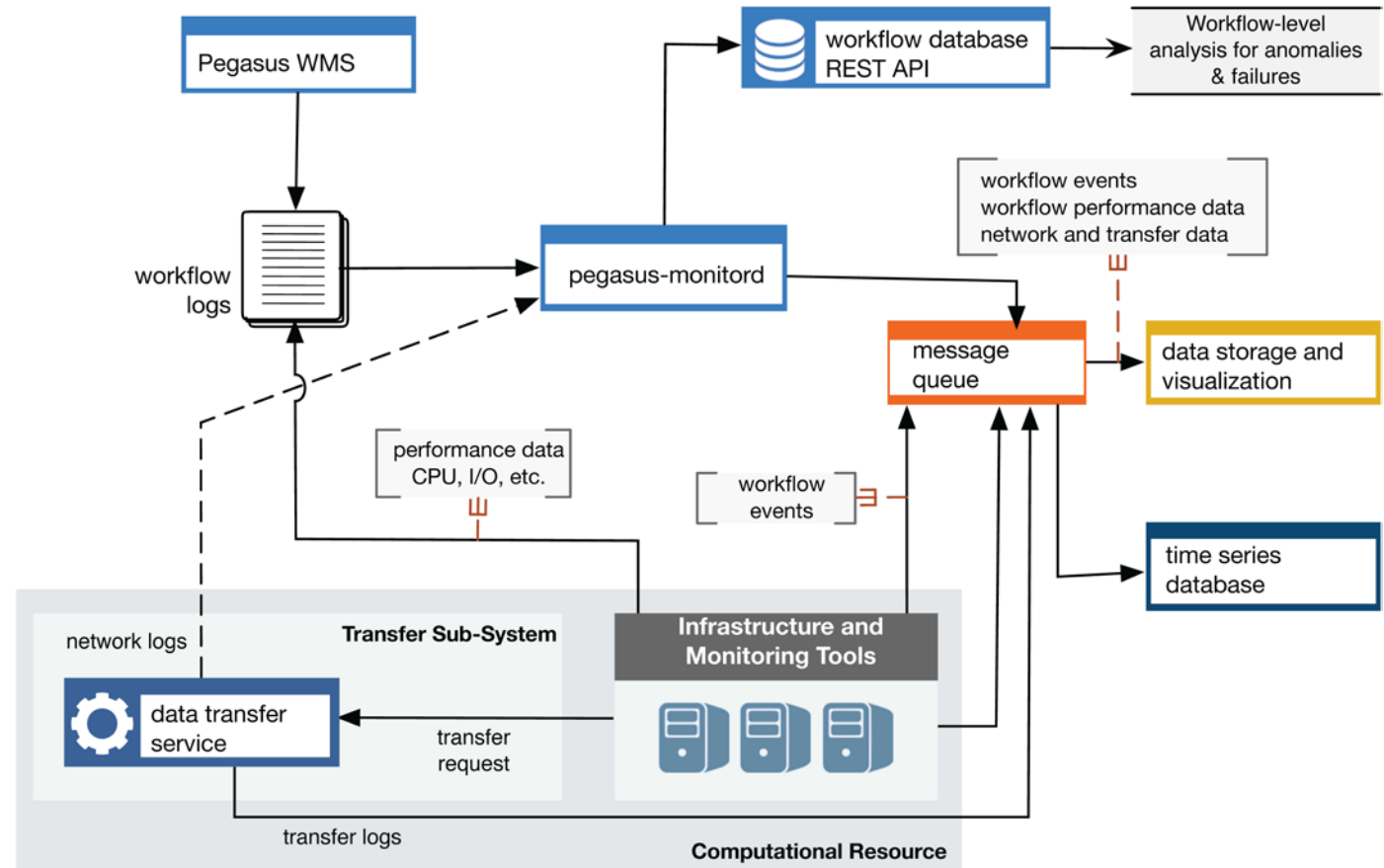
```
4 {
5   "@version": "1",
6   "label": "0a114f53-f208-44da-9b5c-9c30edcf49c6 - stage_in_local_local_2_2",
7   "subtasks_pending": 0,
8   "files_transferred": 2,
9   "sync_level": null,
10  "destination_endpoint_display_name": "georgepap#exo-master",
11  "symlinks": 0,
12  "nice_status_details": null,
13  "completion_time": "2019-03-18 00:40:26+00:00",
14  "source_endpoint_display_name": "georgepap#exo-data",
15  "canceled_by_admin": null,
16  "effective_bytes_per_second": 7919,
17  "delete_destination_extra": false,
18  "wf_uuid": "0a114f53-f208-44da-9b5c-9c30edcf49c6",
19  "deadline": "2019-03-19 00:40:20+00:00",
20  "fatal_error": null,
21  "task_id": "68df23b8-4916-11e9-9e69-0266b1fe9f9e",
22  "nice_status_expires_in": null,
23  "nice_status": null,
24  "recursive_symlinks": "ignore",
25  "username": "georgepap",
26  "owner_id": "b821e4e5-52df-41b7-a27c-4881691d259a",
27  "bytes_transferred": 32248,
28  "subtasks_expired": 0,
29  "status": "SUCCEEDED",
30  "history_deleted": false,
31  "type": "TRANSFER",
32  "@timestamp": "2019-03-18T00:40:31.932Z",
33  "subtasks_failed": 0,
34  "subtasks_total": 4,
35  "preserve_timestamp": false,
36  "is_paused": false,
37  "directories": 0,
38  "verify_checksum": false,
39  "transfer_events": [{}],
40  "DATA_TYPE": "task",
41  "request_time": "2019-03-18 00:40:21+00:00",
42  "files_skipped": 0,
43  "subtasks_retrying": 0,
44  "source_endpoint": "georgepap#exo-data",
45  "faults": 0,
46  "event": "transfer.inv.go",
47  "dag_job_id": "stage_in_local_local_2_2",
48  "subtasks_succeeded": 4,
49  "subtasks_canceled": 0,
50  "destination_endpoint": "georgepap#exo-master",
51  "files": 2,
52  "command": "API 0.10",
53  "encrypt_data": false,
54  "canceled_by_admin_message": null,
55  "source_endpoint_id": "1baf5918-f813-11e8-8cda-0a1d4c5c824a",
56  "nice_status_short_description": null,
57  "destination_endpoint_id": "00a7def6-f813-11e8-8cda-0a1d4c5c824a",
58  "bytes_checksummed": 0
59 }
```

Data Sources: Problems

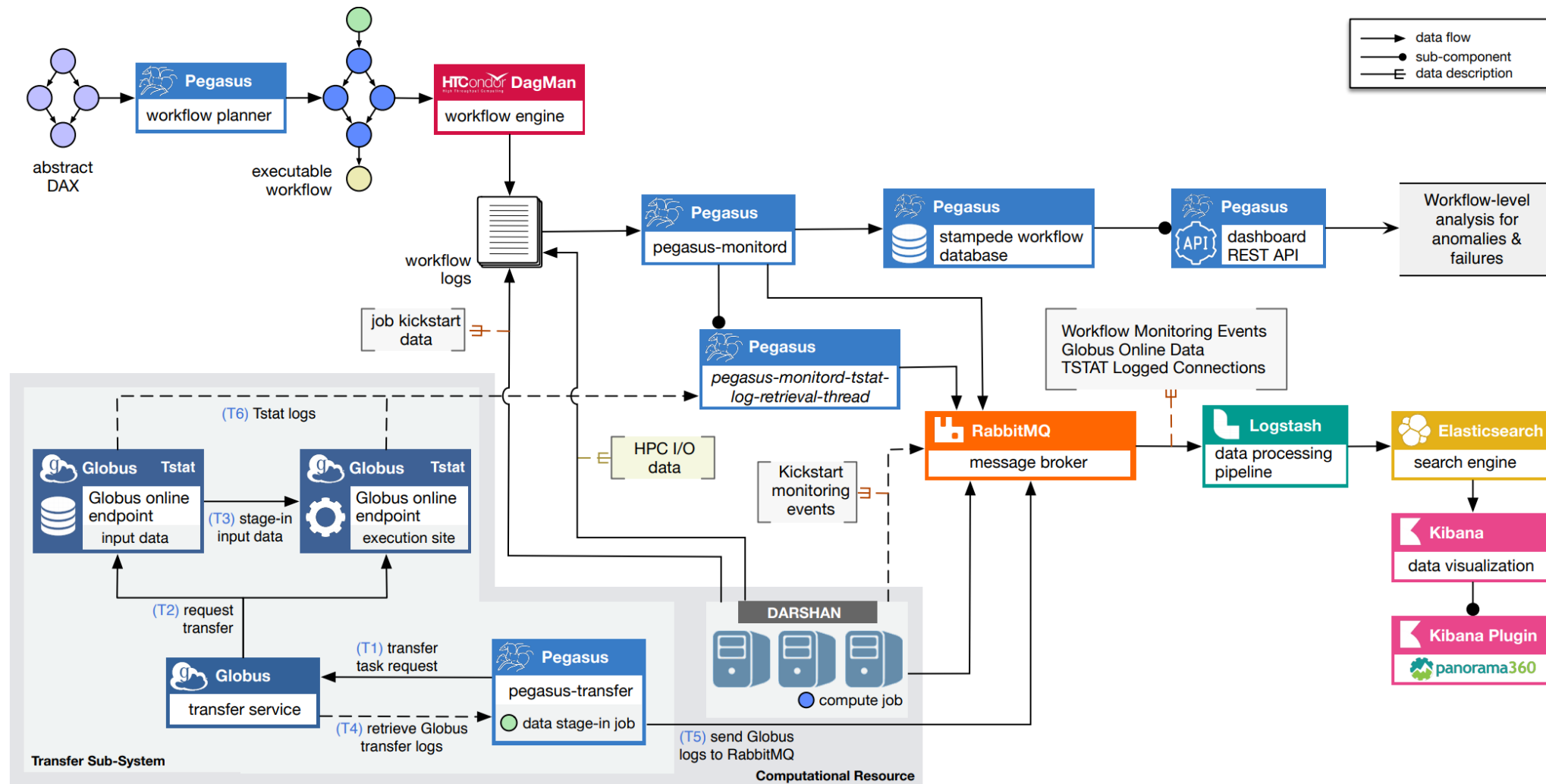
- They are scattered across multiple locations (Eg. execution site, cloud service, pegasus logs)
- They don't contain metadata about the workflow, and it's very hard to locate and match them in the future
- Captured data don't have a common format
 - Pegasus Kickstart logs are in XML format
 - Pegasus Stampede events are in JSON format
 - Pegasus Kickstart online logs are in JSON format
 - Globus logs are in JSON format
 - Darshan logs are in binary format

Data Collection: End-to-End Workflow Execution Monitoring

- **Pegasus** apart from planning and running the workflow, orchestrates the data collection
- A **message queueing system** is used, to decouple the publishers from the datastore
- Flexible **search** and **visualization engines** are used to explore the data



Data Collection: Architecture Overview

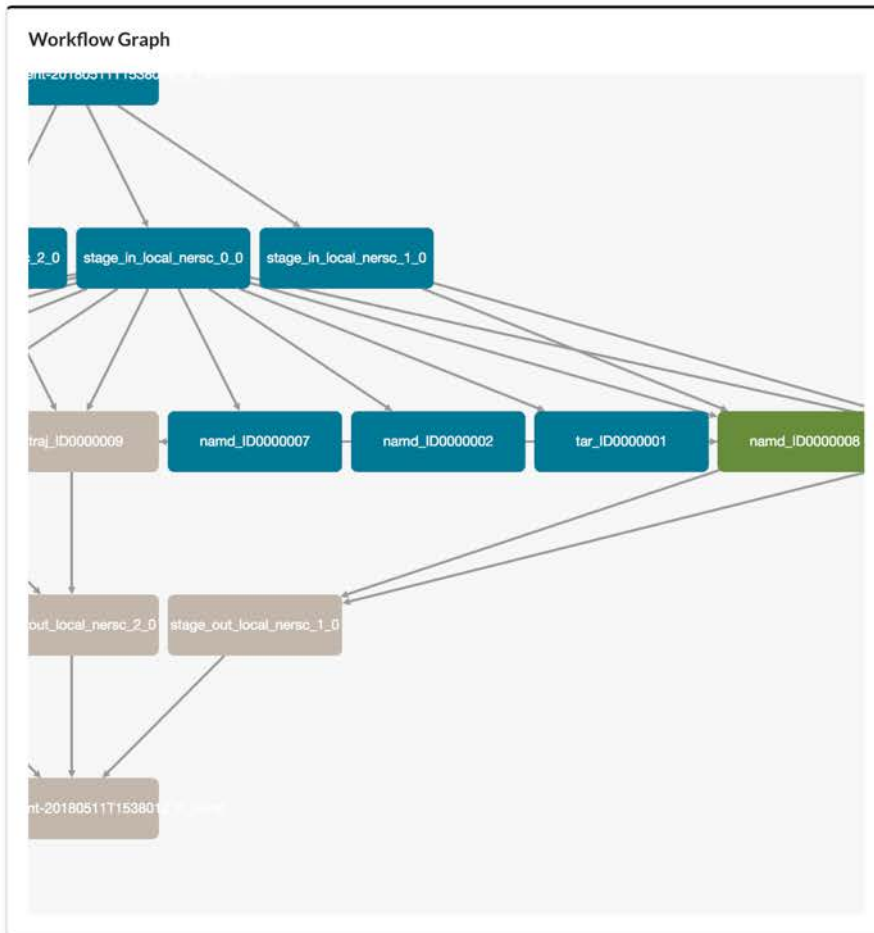


Data Collection: Tool enhancements and new tools

- **pegasus-monitord**: extended with JSON output format, the ability to pickup job related monitoring messages, and publish to amqp endpoints
- **pegasus-transfer**: extended to support Globus transfers and the ability to publish statistics in json format to amqp endpoints
- **pegasus-darshan**: wrapper to darshan-parser, that pushes darshan logs in JSON format to **pegasus-monitord**

Visualization: Detailed Workflow and Job Characteristics

Workflow Dashboard Panorama 360 Workflow Performance Explorer



Workflow Characteristics

Workflow Label
refinement-20180511T153801Z

Workflow ID
f6d04639-6bd2-4e77-a7f7-d7490e697db7

Start Time
5/11/2018, 8:38:09 AM

Running

Time series

Job Characteristics

Job ID
namd_ID0000002

Executable
/global/project/projectdirs/m2187/papajim/sns-wrappers/namd_wrapper

Type
compute

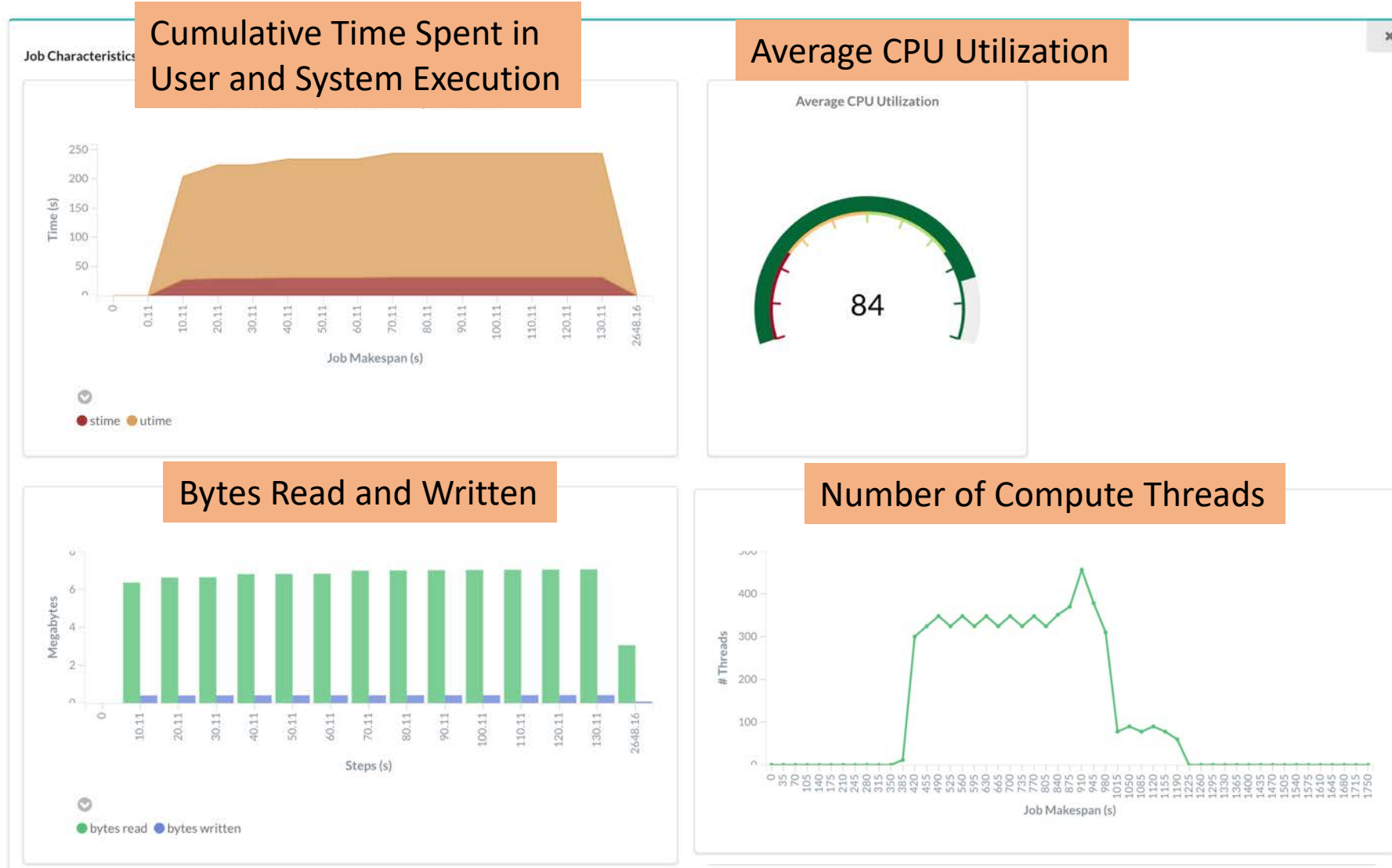
Duration
12.817s

Remote CPU Time
0.508s

Exit Code
0

Time series

Visualization: Time Series Data of Workflow Performance



Repository: Organization

ElasticSearch Index	Description
panorama_transfer	Globus logs
panorama_kickstart	Pegasus-Kickstart online traces
panorama_stampede	Workflow Events and Darshan logs

Repository: Open access data



<https://data.panorama.isi.edu>



<https://kibana.panorama.isi.edu>

How to Deploy: Prerequisites

- HTCondor 8.6+:
 - <https://research.cs.wisc.edu/htcondor/downloads/>
- Pegasus Panorama:
 - Compile from source: <https://github.com/pegasus-isi/pegasus/tree/panorama>
 - Pre-compiled binaries: <http://download.pegasus.isi.edu/pegasus/4.9.2panorama/>
- Docker 17.02+:
 - <https://docs.docker.com/install/>
- Docker Compose:
 - <https://docs.docker.com/compose/>

How to Deploy: Monitoring Backend (RabbitMQ, ELK Stack)

- On a **host** that has **Docker** and **Docker Compose** installed, clone <https://github.com/Panorama360/data-collection-arch>
- Change to the cloned directory and execute the following command:

docker-compose up -d

- **Example:**

```
georgpap@iris:~/GitHub/panorama360/data-collection-arch$ docker-compose up -d
Creating network "panorama_net" with driver "bridge"
Creating panorama-rabbitmq      ... done
Creating panorama-elasticsearch ... done
Creating panorama-kibana        ... done
Creating panorama-logstash      ... done
georgpap@iris:~/GitHub/panorama360/data-collection-arch$ █
```

How to Deploy: Checking Services (RabbitMQ, ELK Stack)

Now the **host** should have **RabbitMQ**, **Elasticsearch**, **Logstash**, and **Kibana** running as **Docker** containers with their service ports exposed. Try to access them...

- **RabbitMQ**: <http://<hostname or ip>:15672>
- **Elasticsearch**: <http://<hostname or ip>:9200>
- **Logstash**: <http://<hostname or ip>:9600>
- **Kibana**: <http://<hostname or ip>:5601>

How to Deploy: Enabling Stampede Events

- In order to get **pegasus-monitord** to publish **all** of its events to an AMQP endpoint in JSON format, **3 properties** must be specified in the workflow's properties file (eg. "pegasus.properties").
 - **pegasus.monitord.encoding** = json
 - **pegasus.catalog.workflow.amqp.url** =
amqp://[username:password]@<hostname>[:port]/<exchange_name>
 - **pegasus.catalog.workflow.amqp.events** = stampede.*

- **Example:**

```
19 # Monitord Events
20 pegasus.monitord.encoding=json
21 pegasus.catalog.workflow.amqp.url=amqp://panorama:panorama@amqp.isi.edu:5672/panorama/monitoring
22 pegasus.catalog.workflow.amqp.events = stampede.*
```

- More about stampede events: https://pegasus.isi.edu/documentation/stampede_wf_events.php

How to Deploy: Enabling Transfer Events

- In order to get pegasus-transfer to publish transfer statistics from the Globus Transfer Service to an AMQP endpoint in JSON format, 2 profiles must be specified in the workflow's sites catalog (eg. "sites.xml"), under the site where pegasus-transfer is going to be invoked (eg. "local").
 - `env.PEGASUS_TRANSFER_PUBLISH = 1`
 - `env.PEGASUS_AMQP_URL =`
`amqp://[username:password]@<hostname>[:port]/<exchange_name>`
- **Example:**

```
5 <!-- The local site contains information about the submit host -->
6 <site handle="local">
7   <directory type="shared-scratch" path="${PWD}/work/scratch">
8     <file-server operation="all" url="go://georgepap#exo-master/${PWD}/work/scratch"/>
9   </directory>
10  <directory type="local-storage" path="${PWD}/work/outputs">
11    <file-server operation="all" url="go://georgepap#exo-master/${PWD}/work/outputs"/>
12  </directory>
13  <!-- These profiles tell pegasus-transfer to publish stats to an AMQP endpoint -->
14  <profile namespace="env" key="PEGASUS_TRANSFER_PUBLISH">1</profile>
15  <profile namespace="env" key="PEGASUS_AMQP_URL">amqp://panorama:panorama@amqp.isi.edu:5672/panorama/monitoring</profile>
16 </site>
```

How to Deploy: Enabling Kickstart Online Traces

- In order to get pegasus-kickstart to publish traces of resource usage statistics to an AMQP endpoint in JSON format, 2 profiles must be specified in the workflow's sites catalog (eg. "sites.xml") under the compute site.
 - `pegasus.gridstart.arguments = -m <interval in seconds>`
 - `env.KICKSTART_MON_URL =`
`rabbitmq://[USERNAME:PASSWORD]@<hostname>[:port]/api/exchanges/<exchange_name>/publish`

- **Example:**

```
23 <site handle="condorpool" arch="x86_64" os="LINUX">
24   <!-- These profiles tell Pegasus that the site is a plain Condor pool -->
25   <profile namespace="pegasus" key="style">condor</profile>
26   <profile namespace="condor" key="universe">vanilla</profile>
27   <!-- These profiles tell pegasus-kickstart to publish stats to an AMQP endpoint -->
28   <profile namespace="pegasus" key="gridstart.arguments">-m 10</profile>
29   <profile namespace="env" key="KICKSTART_MON_URL">rabbitmq://panorama:panorama@amqp.isi.edu:15672/api/exchanges/panorama/monitoring/publish</profile>
30 </site>
```

- Alternatively if we want to customize the monitoring interval per computational task we can specify the profile in the workflow's transformation catalog (eg. "tx.txt")

How to Deploy: Enabling Kickstart Online Traces (MPI Jobs)

- Usually MPI jobs are not launched by Pegasus-Kickstart. Thus, adding the `gridstart.arguments` profile doesn't have any effect.
- We can work around this by using a wrapper script for the MPI job, that invokes directly `pegasus-monitor`.
- We still need to specify **KICKSTART_MON_URL** in the `sites catalog`.

- **Example:**

```
1  #!/usr/bin/env bash
2
3  cd $PEGASUS_SCRATCH_DIR
4
5  mpirun pegasus-monitor -i 10 /shared/software/NAMD_2.12_Linux-x86_64-MPI/namd2 $@
6
```

How to Deploy: Enabling Darshan Logs (MPI Jobs)

- In case your MPI application wasn't compiled and statically linked with Darshan's library, we need to set a profile in the transformation catalog, adding the path of the library to LD_PRELOAD.
- We launch the application using a **wrapper script**, and as post job steps:
 - Build the darshan log path from the environmental variables
 - Invoke **pegasus-darshan** with the files as input

```
1 tr namd {
2   site local-slurm {
3     pfn "/shared/software/wrappers/namd_wrapper_slurm.sh"
4     arch "x86_64"
5     os "LINUX"
6     type "INSTALLED"
7     profile pegasus "exitcode.successmsg" "End of program"
8     profile pegasus "memory" "2500M"
9     profile globus "jobtype" "single"
10    profile env "LD_PRELOAD" "/shared/software/darshan-3.1.6/lib/libdarshan.so"
11  }
12 }
```

Transformation Catalog

```
1 #!/usr/bin/env bash
2
3  cd $PEGASUS_SCRATCH_DIR
4
5  mpirun pegasus-monitor -i 10 /shared/software/NAMD_2.12_Linux-x86_64-MPI/namd2 $@
6
7  #post job parse darshan output
8  DAY=$(date +%d)
9  DAY=${DAY##0}
10 MONTH=$(date +%m)
11 MONTH=${MONTH##0}
12 YEAR=$(date +%Y)
13
14 darshan_base=/shared/darshan-logs/${YEAR}/${MONTH}/${DAY}
15 darshan_file=${darshan_base}/${SLURM_JOB_USER}_namd2_id${SLURM_JOB_ID}_*.darshan
16
17 for f in $darshan_file; do
18   $PEGASUS_HOME/bin/pegasus-darshan -f "$f"
19 done
20
```

Wrapper Shell Script

How to Deploy: Enabling Darshan Logs (MPI Jobs)

- **pegasus-darshan** will output in stdout a monitoring payload, that will be picked by **pegasus-monitor**, which in its turn will publish it to the AMQP endpoint.
- *This can also be used as a generic way of adding new tools to this architecture.*

```
1  @@@MONITORING_PAYLOAD - START@@@
2  {
3      &quot;monitoring_event&quot;;: &quot;darshan.perf&quot;;,
4      &quot;payload&quot;;: [
5          {
6              &quot;POSIX_module_data&quot;;: {
7                  &quot;agg_perf_by_cumul&quot;;: 14.667417,
8                  &quot;agg_perf_by_open&quot;;: 14.667417,
9                  &quot;agg_perf_by_open_lastio&quot;;: 14.667417,
10                 &quot;agg_perf_by_slowest&quot;;: 14.667417,
11                 &quot;shared_files&quot;;: {
12                     &quot;time_by_cumul_io_only&quot;;: 0.0,
13                     &quot;time_by_cumul_meta_only&quot;;: 0.0,
14                     &quot;time_by_open&quot;;: 0.0,
15                     &quot;time_by_open_lastio&quot;;: 0.0,
16                     &quot;time_by_slowest&quot;;: 0.0
17                 },
18                 &quot;total_bytes&quot;;: 403761,
19                 &quot;unique_files&quot;;: {
20                     &quot;slowest_rank&quot;;: 0.0,
21                     &quot;slowest_rank_io_time&quot;;: 0.026253,
22                     &quot;slowest_rank_meta_only_time&quot;;: 0.023997
23                 }
24             }
25         },
26         &quot;ts&quot;;: 1552878285
27     ]
28 }
29 @@@MONITORING_PAYLOAD - END@@@
```

DEMO



- GitHub:
<https://github.com/Panorama360>
- Website:
<https://panorama360.github.io>



George Papadimitriou

Computer Science PhD Student
University of Southern California

email: georgpap@isi.edu

USC Viterbi
School of Engineering
Department of Computer Science

<https://panorama360.github.io/>



Pegasus

est. 2001

Automate, recover, and debug scientific computations.

Get Started

Pegasus Website

<http://pegasus.isi.edu>

Users Mailing List

pegasus-users@isi.edu

Support

pegasus-support@isi.edu

Pegasus Online Office Hours

<https://pegasus.isi.edu/blog/online-pegasus-office-hours/>

Bi-monthly basis on second Friday of the month, where we address user questions and also apprise the community of new developments