



# A Peek into Workflow I/O

**Jakob Lüttgau**, Shane Snyder, Phil Carns, Justin M. Wozniak, Julian Kunkel, Thomas Ludwig

**BoF Analyzing Parallel I/O, SC'18**

November 13, 2018 / Dallas, TX

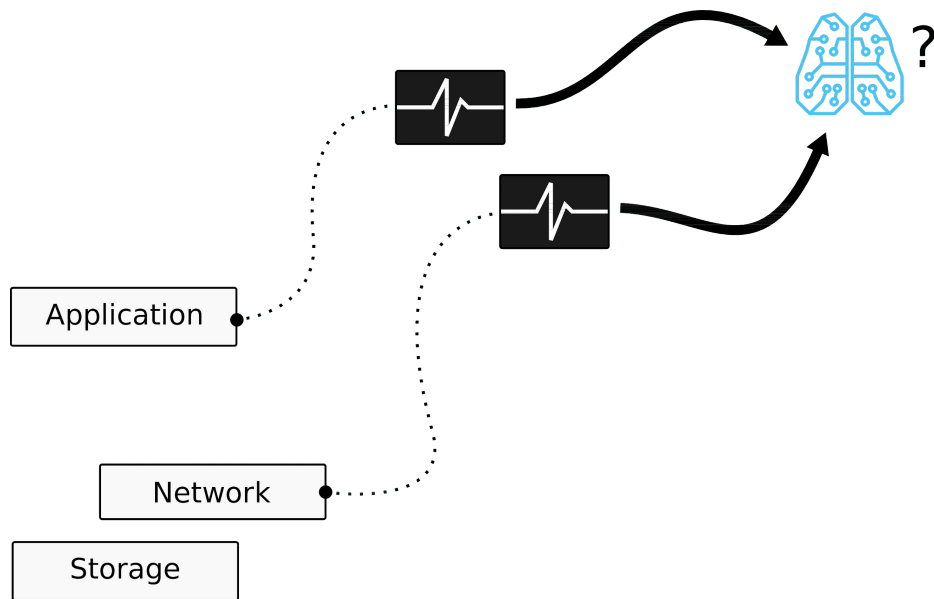
# Workflows and HPC

Workflows offer ...

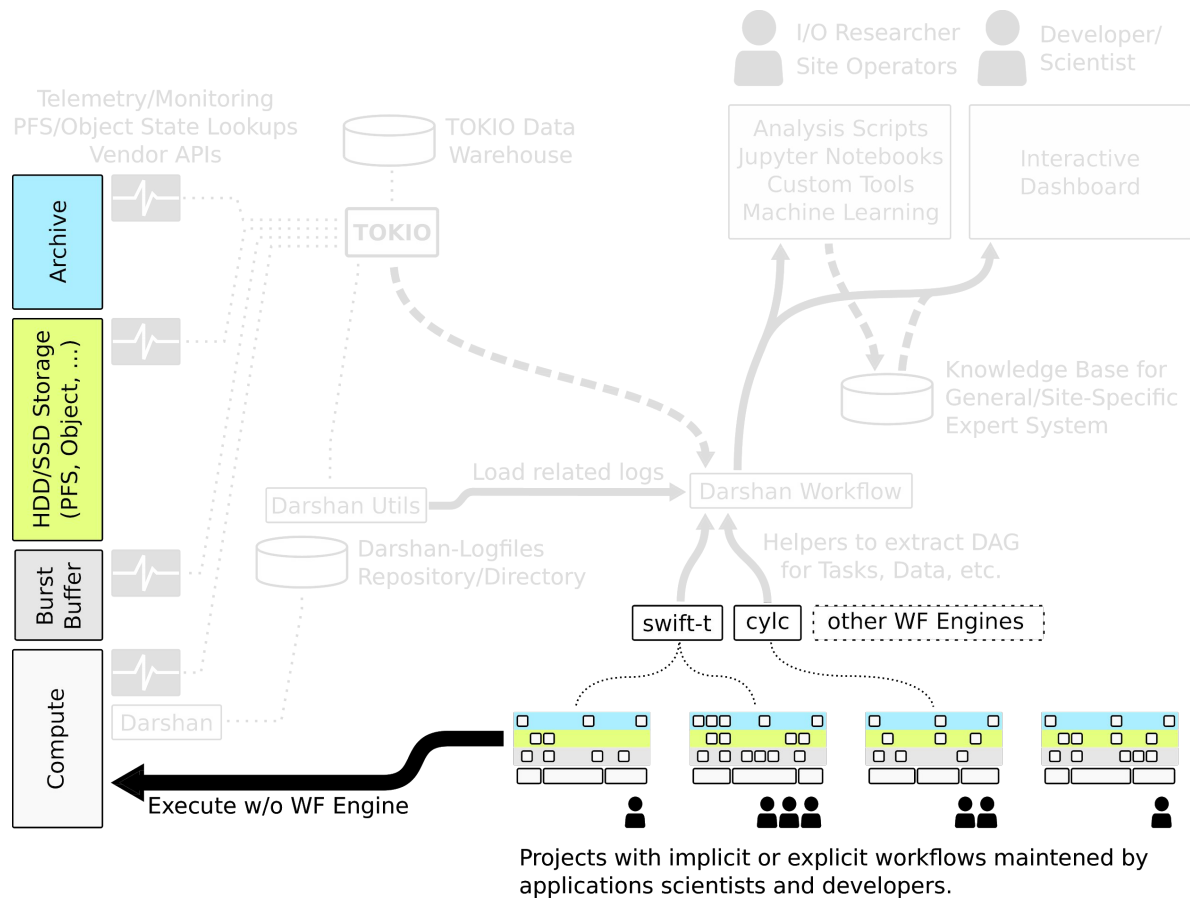
- ... anticipatable future activity
- ... implicit intent to be discovered
- ... explicit intent description

# How can we better understand and react to I/O behavior of workflows in HPC?

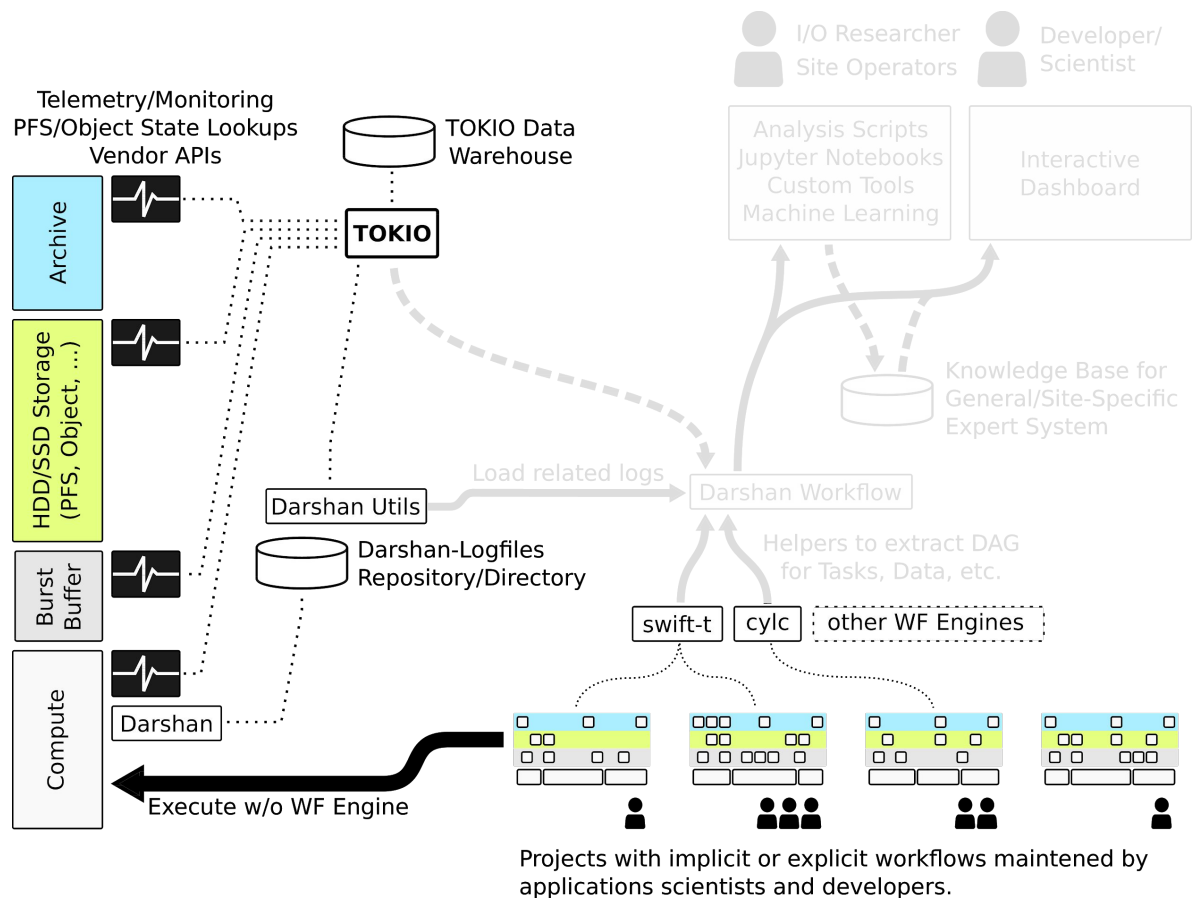
- Explore convenient toolchain for researchers and site operators.
- Considerations for user facing components to improve communicating advice.



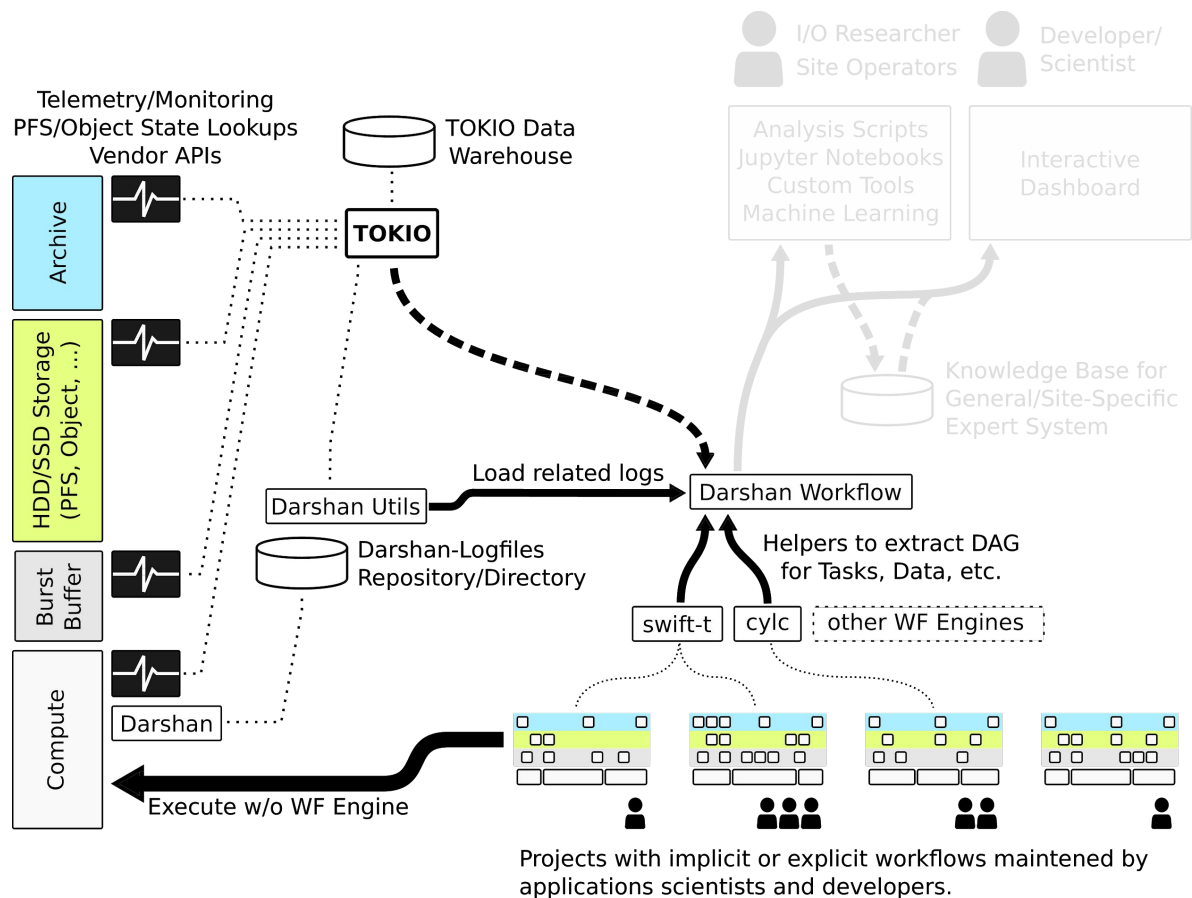
# Architecture for Augmenting I/O in Workflows



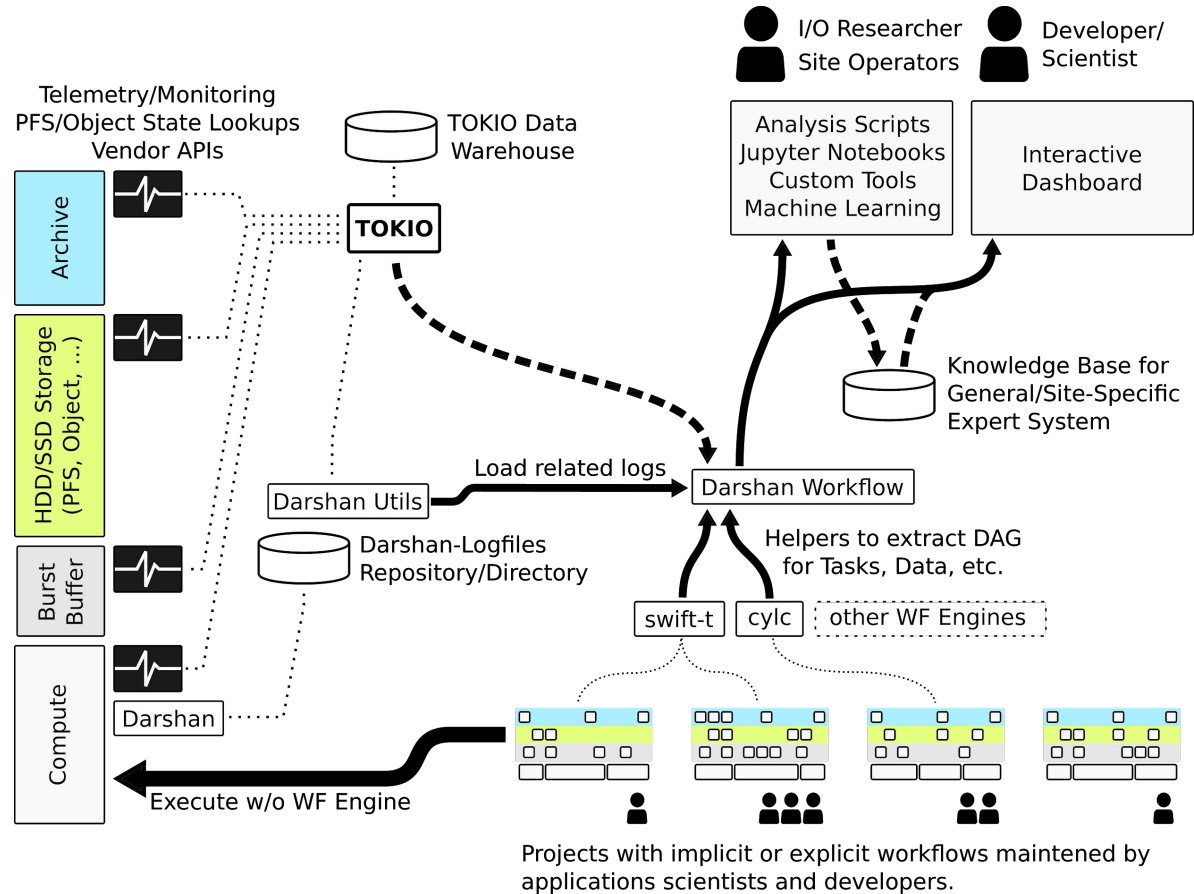
# Architecture for Augmenting I/O in Workflows



# Architecture for Augmenting I/O in Workflows



# Architecture for Augmenting I/O in Workflows



# Data-Intensive Exascale Workflow: Climate Modeling



ICON is a climate model used by Researchers at Max-Planck and by the German Weather Service (DWD).  
CDO is a pre/post-processing tool (climate operators) for NetCDF files.  
ParaView is a popular visualisation toolkit built on top of VTK.





```
pq@mcswl209:~/ANL/darshan-workflow/demo/workflow-cylc-examle
File Edit View Search Terminal Tabs Help

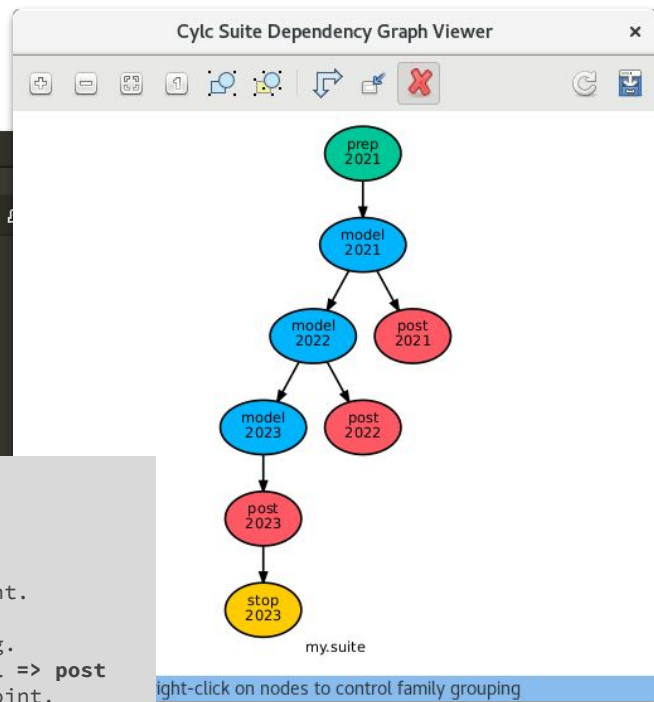
pq@mcswl209:~/ANL/darshan-workflow/demo/workflow-cylc-examle
$ ./02-visualize.sh

pq@mcswl209:~/ANL/darshan-workflow/demo/workflow-cylc-examle
$ ./run.sh
REGISTER my.suite: /home/pq/ANL/darshan-workflow/demo/workflow-cylc-examle/suites/test
my.suite | A first Cylc suite. | ~/ANL/darshan-workflow/demo/workflow-cylc-examle/suites/test
Valid for cylc-UNKNOWN
/home/pq/ANL/darshan-workflow/develop/testbed/install/software/darshan
libdarshan.a libdarshan.so libdarshan-stubs.a libdarshan-util.a libdarshan-util.so Number pkgconfig Tex

The Cylc Suite Engine [UNKNOWN]
Copyright (C) 2008-2018 NIWA

This program comes with ABSOLUTELY NO WARRANTY;
see 'cylc warranty'. It is free software, you
are welcome to redistribute it under certain
conditions; see 'cylc conditions'.

2018-07-17T17:02:05-05 INFO - Suite starting: server=mcswl209.mcs.anl.
2018-07-17T17:02:05-05 INFO - Cylc version: UNKNOWN
2018-07-17T17:02:05-05 INFO - Run mode: live
2018-07-17T17:02:05-05 INFO - Initial point: 2021
2018-07-17T17:02:05-05 INFO - Final point: 2023
2018-07-17T17:02:05-05 INFO - Cold Start 2021
2018-07-17T17:02:05-05 INFO - [prep.2021] -submit-num=1, owner@host=lo
2018-07-17T17:02:06-05 INFO - [prep.2021] -(current:ready) submitted a
2018-07-17T17:02:06-05 INFO - [prep.2021] -job[01] submitted to localh
2018-07-17T17:02:06-05 INFO - [prep.2021] -health check settings: subm
2018-07-17T17:02:06-05 INFO - [prep.2021] -(current:submitted)> starte
2018-07-17T17:02:06-05 INFO - [prep.2021] -health check settings: exec
2018-07-17T17:02:07-05 INFO - [prep.2021] -(current:running)> succee
2018-07-17T17:02:08-05 INFO - [model.2021] -submit-num=1, owner@host=l
2018-07-17T17:02:09-05 INFO - [model.2021] -(current:ready) submitted
2018-07-17T17:02:09-05 INFO - [model.2021] -job[01] submitted to local
2018-07-17T17:02:09-05 INFO - [model.2021] -health check settings: sub
2018-07-17T17:02:09-05 INFO - [model.2021] -(current:submitted)> start
2018-07-17T17:02:09-05 INFO - [model.2021] -health check settings: exe
2018-07-17T17:02:10-05 CRITICAL - [model.2021] -(current:running)> fai
2018-07-17T17:02:10-05 CRITICAL - [model.2021] -job[01] failed
2018-07-17T17:02:11-05 WARNING - suite stalled
2018-07-17T17:02:11-05 WARNING - Unmet prerequisites for stop.2023:
2018-07-17T17:02:11-05 WARNING - * post.2023 succeeded
2018-07-17T17:02:11-05 WARNING - Unmet prerequisites for model.2022:
2018-07-17T17:02:11-05 WARNING - * model.2021 succeeded
2018-07-17T17:02:11-05 WARNING - Unmet prerequisites for post.2021:
2018-07-17T17:02:11-05 WARNING - * model.2021 succeeded
```



```
[scheduling]
initial cycle point = 2021
final cycle point = 2023
[[dependencies]]
[[[R1]]] # Initial cycle point.
graph = prep => model
[[[R/P1Y]]] # Yearly cycling.
graph = model[-P1D] => model => post
[[[R1/P0Y]]] # Final cycle point.
graph = post => stop
```

```
[runtime]
[[prep]]
script = mpiexec -np 1 ./prep
[[model]]
script = mpiexec -np 4 ./model
[[post]]
script = mpiexec -np 1 ./post
```

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

## Augmenting Workflow I/O with Darshan/TOKIO

Scientific discovery increasingly depends on complex workflows consisting of multiple phases and sometimes millions of parallelizable tasks or pipelines. Typical workflows on HPC systems routinely require the pre-processing, generation by simulation and post-processing of data. Unfortunately, most workflow models focus on the scheduling and allocation of resources for tasks while the impact on storage systems remains a secondary objective.

By combining a workflow description (e.g., from a workflow engine like Swift or Cylc) with log data or telemetry information we can gain insight on the I/O behavior of a complete workflow and then optimize applications, middleware and systems accordingly.

In [1]: `1 import darshan.workflow`

In [2]: `1 # Load a workflow report  
2 wf = darshan.workflow.load('data-workflow.json')  
3 wf.data`

Out[2]: 

```
{ 'nodes': [{ 'id': 'model.2021',  
            'label': 'model.2021',  
            'x': 140.5,  
            'y': -410.0,  
            'data': {},  
            'group': 'report'},  
  { 'id': 'model.2022',  
    'label': 'model.2022',  
    'x': 91.5,  
    'y': -320.0,  
    'data': {},  
    'group': 'report'},  
  { 'id': 'post.2021',  
    'label': 'post.2021',  
    'x': 189.5,  
    'y': -320.0,  
    'data': {}},  
  { 'id': 'model.2023',  
    'label': 'model.2023',  
    'x': 42.5,
```

## Inspecting a Workflow from Darshan

Darshan/TOKIO-workflow is build to ease interactive and automatic analysis of workflows with a focus on the I/O perspective. As such it explores a variety of convenience methods and common visualisations such as the `task_summary()` and `show_graph()` methods.

In [3]: `1 wf.task_summary()`

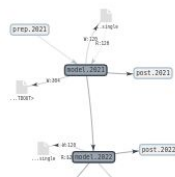
```
model.2021  Records: 2    Layers: POSIX, STDIO  
model.2022  Records: 2    Layers: POSIX, STDIO  
post.2021   Records: 0    Layers:  
model.2023  Records: 2    Layers: POSIX, STDIO  
post.2022   Records: 0    Layers:  
post.2023   Records: 0    Layers:  
stop.2023   Records: 0    Layers:  
prep.2021   Records: 0    Layers:
```

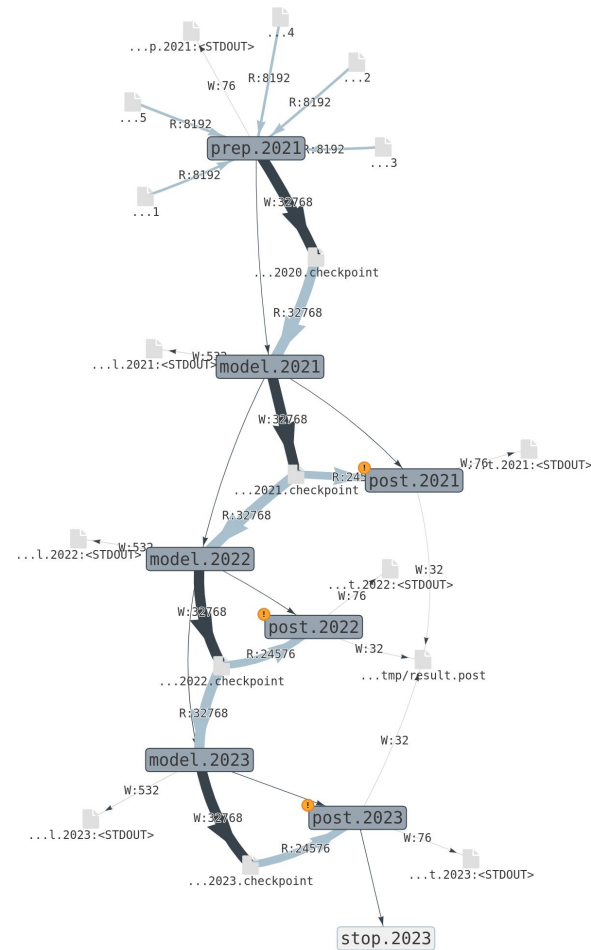
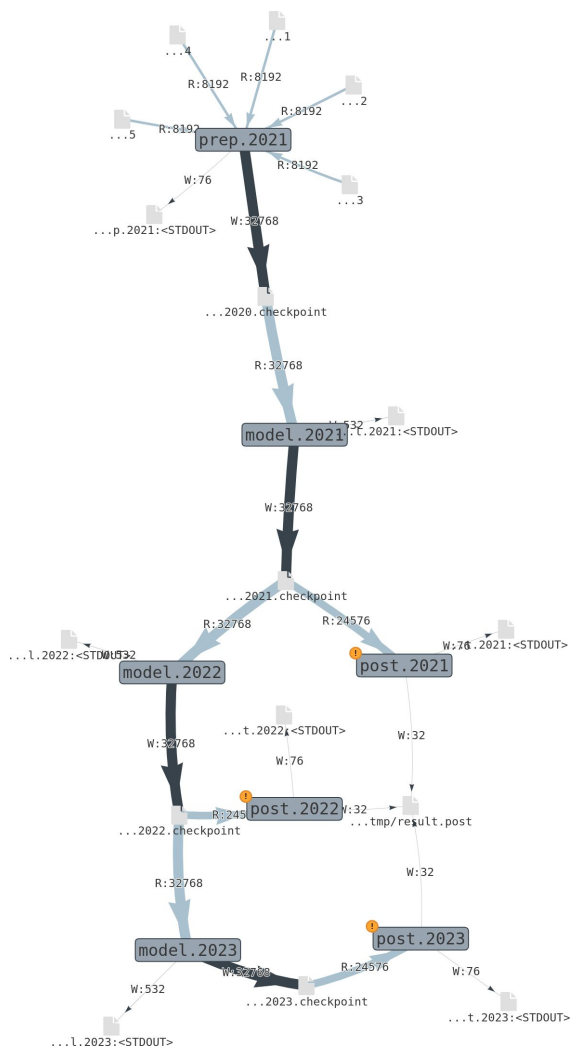
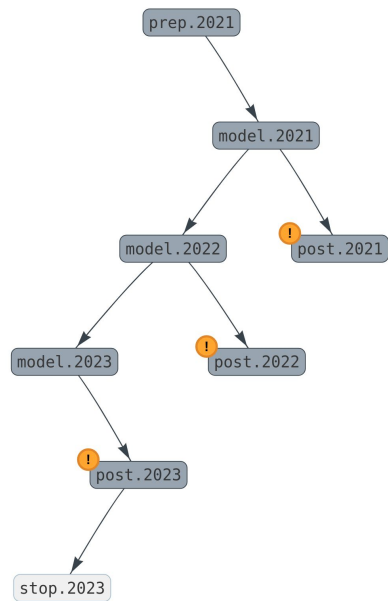
## Displaying and Interacting with Workflow

Jupyter Notebooks can be extended with custom widget, which allows to turn them into versatile tools for custom but powerful tools in I/O analysis. This is especially true for workflows, which generate a lot of log data so that interactive tools make data exploration more convenient.

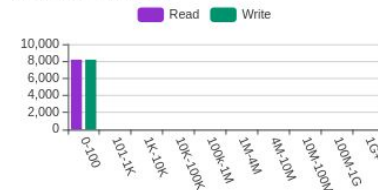
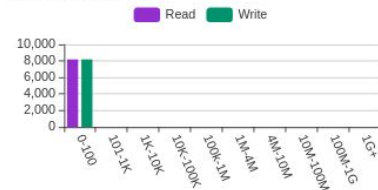
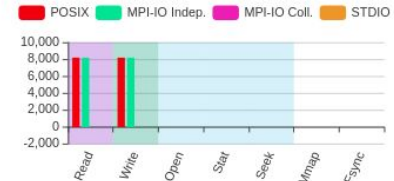
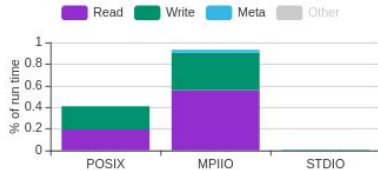
In [4]: `1 import darshan.ipynwidgets.example as wfgraph  
2 import json  
3 wf.show_graph()  
4 graph = wfgraph.HelloWorld()  
5 graph`

T000: GRAPH IpythonWidget





jobid: 19911	uid: 1000	nproc: 4	runtime: 1 seconds
--------------	-----------	----------	--------------------



Layer	Access Size	Count
POSIX	50013	14
POSIX	50007	14
POSIX	49986	9
POSIX	49998	9
MPI-IO	1020	212
MPI-IO	512	34

type	number of files	avg. size	max size
total opened	4	1.8M	7.2M
read-only files	0	0	0
write-only files	1	16K	16K
read/write files	3	2.4M	7.2M
created files	4	1.8M	7.2M

```
Task: model.2021
model.2021: darshan-io
Task: model.2022
model.2022: darshan-io
Task: post.2021
post.2021: darshan-logs
Task: model.2023
model.2023: darshan-io
Task: post.2022
post.2022: darshan-logs
Task: post.2023
post.2023: darshan-logs
Task: prep.2021
prep.2021: darshan-logs
Files

model.2021: /tmp/2021.0
model.2021: /tmp/2021.1
model.2021: <STDOUT>
model.2022: /tmp/2022.0
model.2022: /tmp/2021.1
model.2022: <STDOUT>
post.2021: /tmp/result.p
post.2021: /tmp/2021.ch
post.2021: <STDOUT>
model.2023: /tmp/2023.0
model.2023: /tmp/2022.1
model.2023: <STDOUT>
post.2022: /tmp/result.p
post.2022: /tmp/2022.ch
post.2022: <STDOUT>
post.2023: /tmp/result.p
post.2023: /tmp/2023.ch
post.2023: <STDOUT>
prep.2021: /tmp/file-raw
prep.2021: /tmp/file-raw
prep.2021: /tmp/file-raw
prep.2021: /tmp/2020.ct
prep.2021: /tmp/file-raw
prep.2021: /tmp/file-raw
prep.2021: <STDOUT>
```

# Lessons Learned

## Requirements for Workflow Engines

- Expose Context / DAGs of Workflows
- Data/(File) Notions
- Reflection in Execution Runtime?

## Requirements for Monitoring Solutions

- Pick up context to allow associations
- Support User-Specific Metadata with record
- API to interact with monitoring toolkit
- Allow counters per MPI Communicator

## Requirements for Application Developers

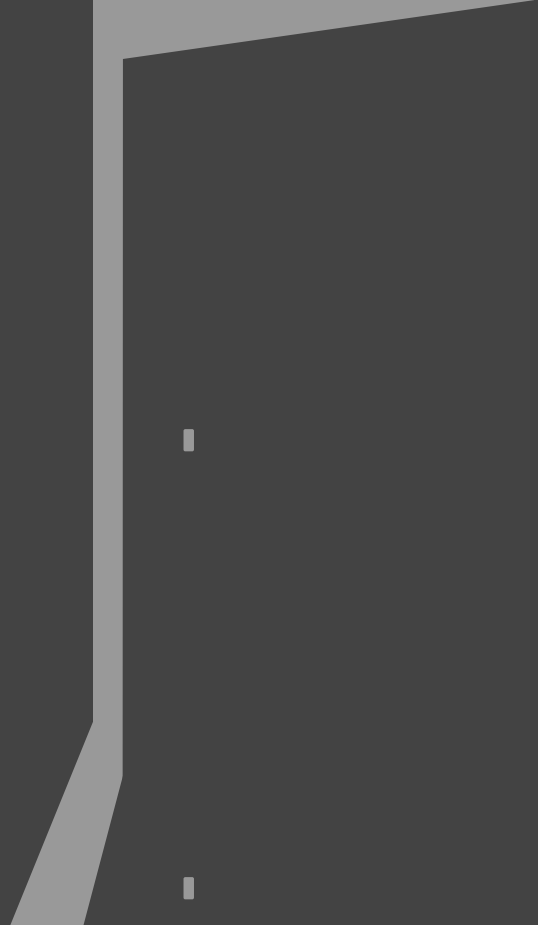
- Make intent explicit: use Libs/DSL (e.g. HDF5)
- Enable Darshan with at least a subset of runs
- Donate traces and logs for a training body.

# Thank you! Questions?

luettgau@dkrz.de

# Thank you! Questions?

luettgau@dkrz.de



# Workflow Engines: Swift, Cylc, Tigres, etc.

Cylc, Swift-k, Fireworks

**Job centric**, with tasks and data targets. Tasks are distributed and possibly run on **remote systems**. Data products might be moved between sites.

*Usually, a coarse granular dependency graph.*

Swift-t, Tigres, Spark/RDD Lineage, QDO

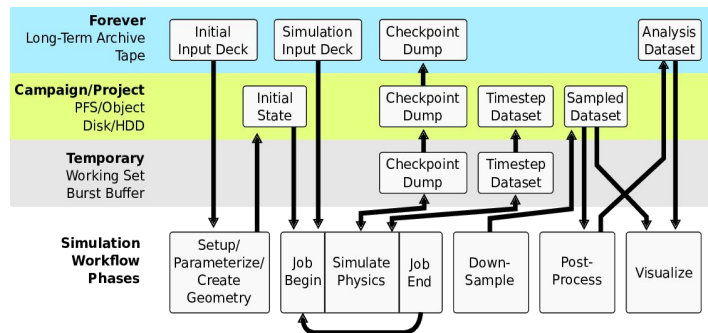
A large **integrated** (MPI) application with many different tasks within the application. With **exascale** in mind and also closer to **in situ** enabled workflows.

*Closer to a programming language.*



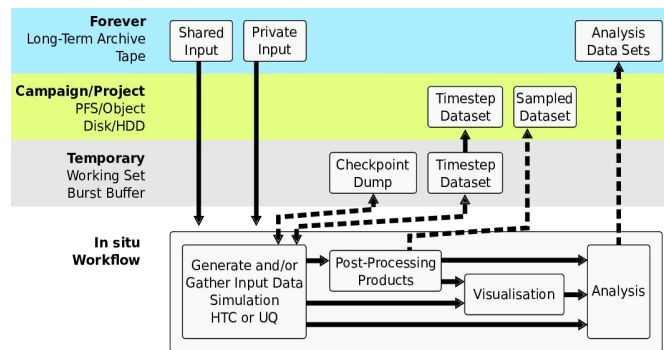
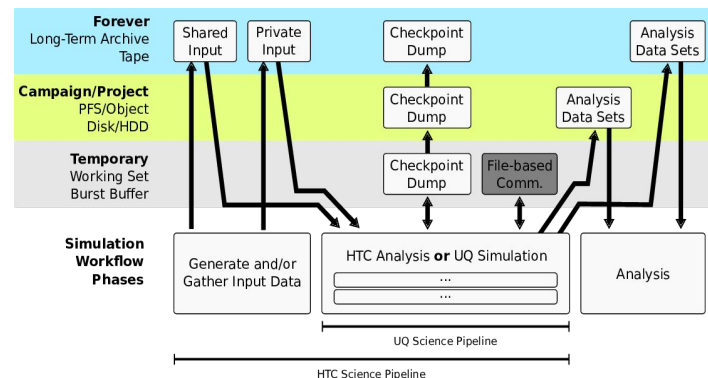
# Common Scientific Workflows in HPC

## What makes a workflow?



SIM

UQ or HTC



in situ

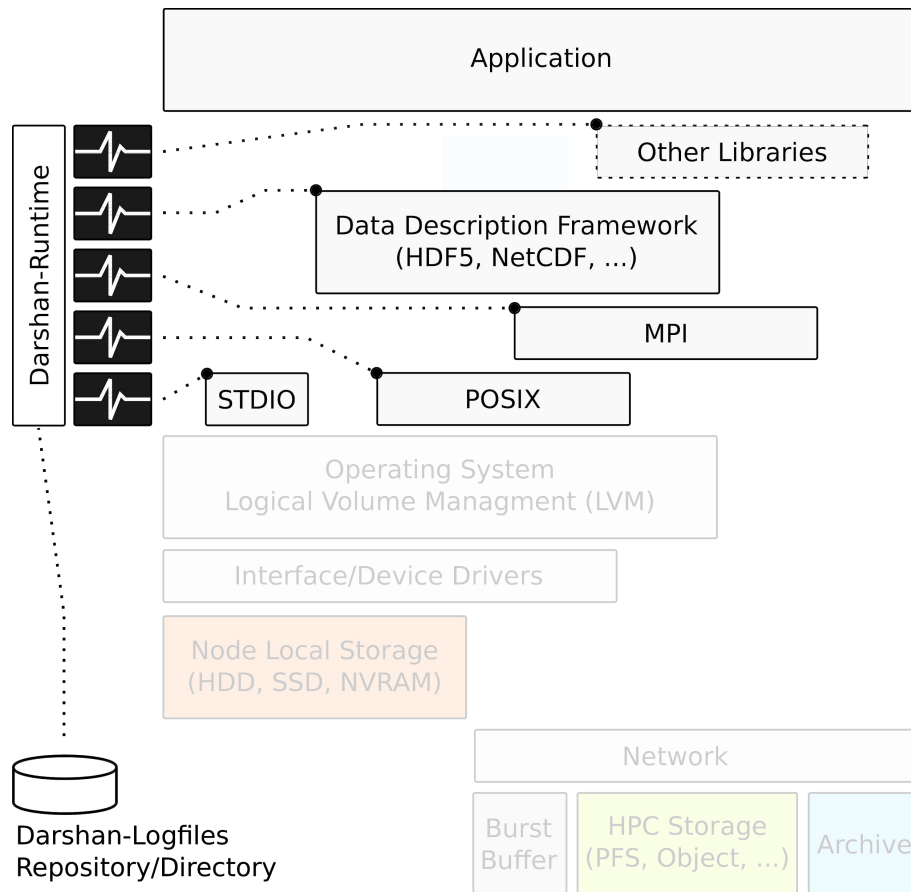
SIM and HTC/UQ are derived figures from [1]. For outlook on workflows refer to [2].

[1] LANL, NERSC, and SNL, "APEX Workflows.", Whitepaper, Mar. 2016

Online: <https://www.nersc.gov/assets/apex-workflows-v2.pdf>

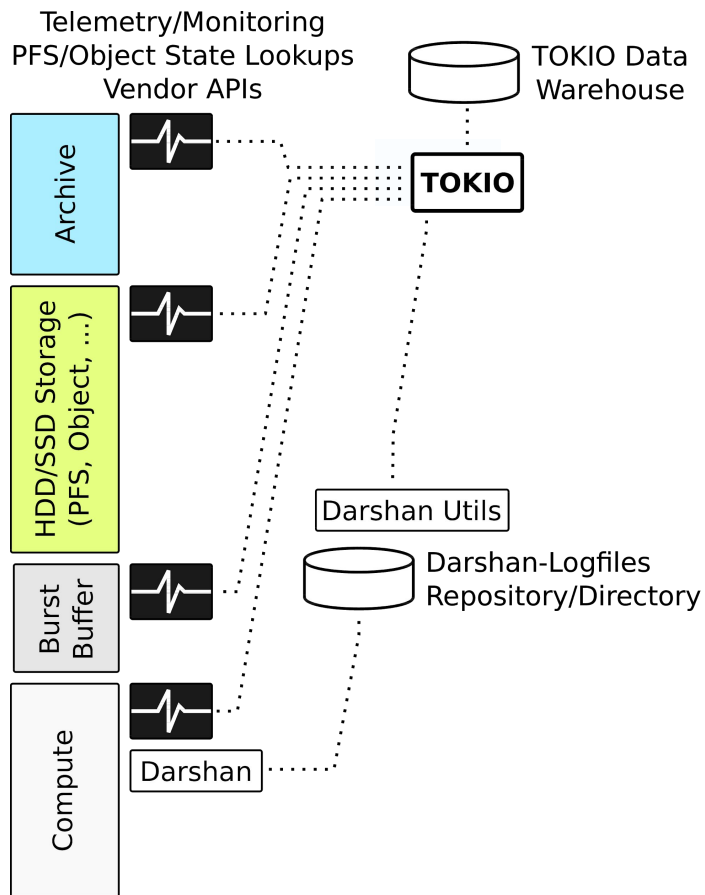
[2] E. Deelman *et al.*, "The future of scientific workflows," *The International Journal of High Performance Computing Applications*, vol. 32, no. 1, pp. 159–175, Jan. 2018.

# Darshan: Instrumentation at Library/Application Layer



```
$ export LD_PRELOAD=libdarshan.so
$ mpiexec -np 4 ./hellompi
```

# TOKIO: Total Knowledge of Input/Output



Comprehensive capture of I/O activity

Support different storage services in data center

May require privileged access in many cases