# Toward Understanding I/O Behavior in HPC Workflows

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

**PDSW-DISC, SC'18**
November 12, 2018 / Dallas, TX

# Overview
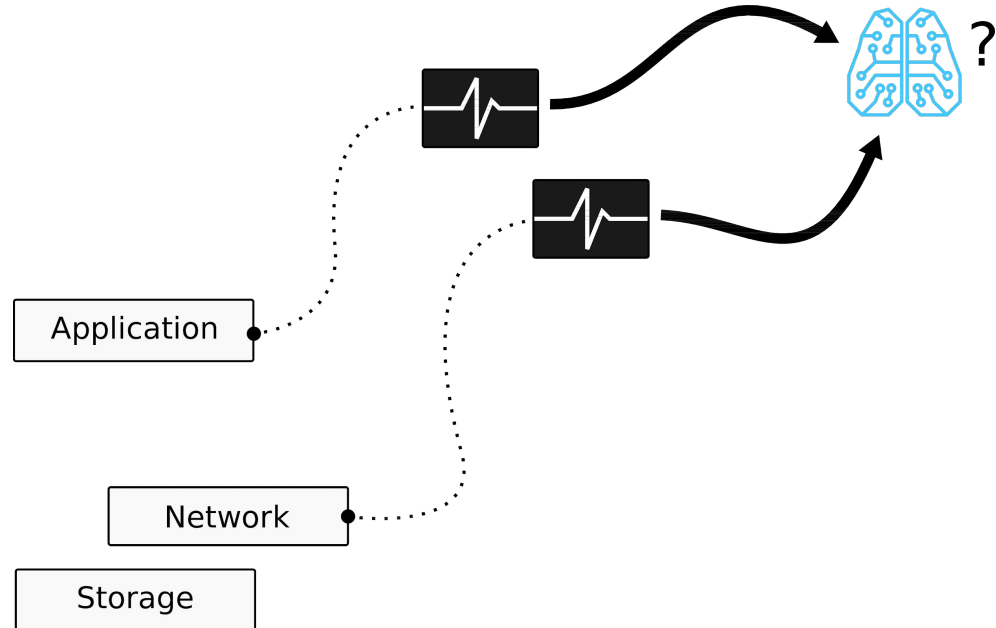
# Trying to add a missing link so we can move closer to realizing smarter systems...

Require **new interfaces** to preserve information about **structure of data**.

How to anticipate **user intentions** and **I/O behavior** of applications ?

Require **tools to observe and record** system activity as a basis to gain insight

# Workflows
a HPC Storage
Perspective?

Workflows offer ...

... anticipatable future activity

... implicit intent to be discovered

... explicit intent description

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

Cylc, Swift-k, Fireworks

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

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

A large **integrated** (MPI) application with many different tasks within the application. With **exascale** in mind and also closer to **in situ** enabled workflows.
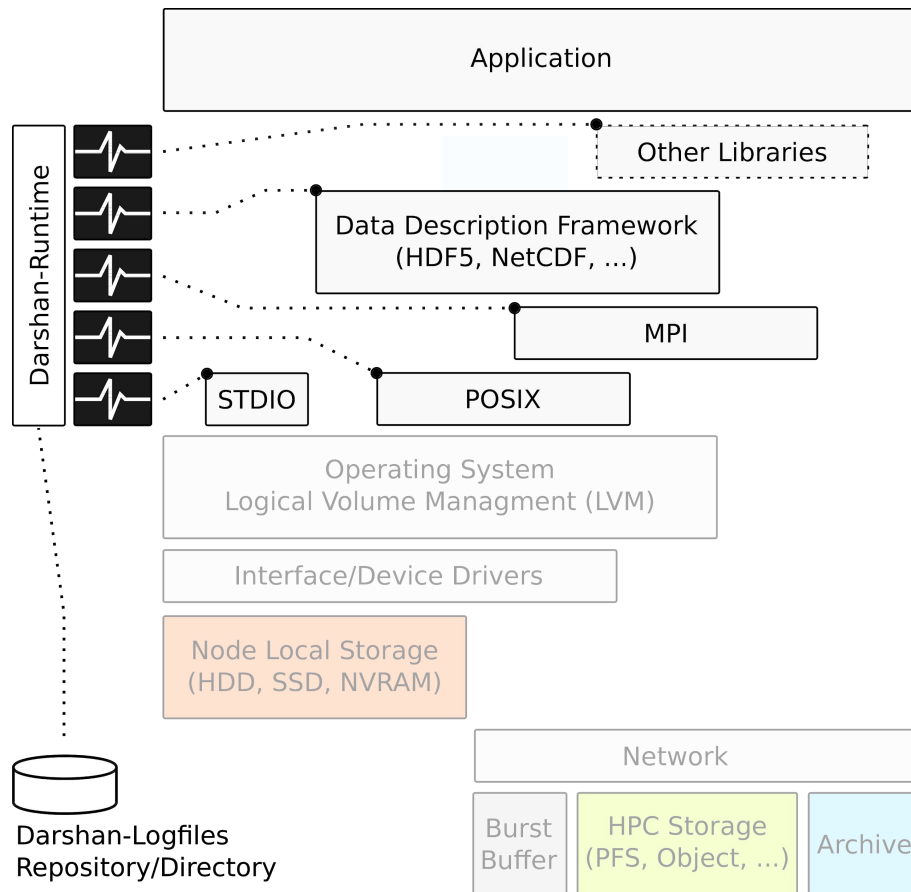
*Usually, a coarse granular dependency graph.*

*Closer to a programming language.*

# Holistic
# I/O Monitoring
# for HPC

Tracking at the Application/Library Layer

Total Knowledge of  I/O in Data Centers

# Darshan: Instrumentation at Library/Application Layer



Application

Other Libraries

Darshan-Runtime

Data Description Framework
(HDF5, NetCDF, ...)

MPI

STDIO

POSIX

Operating System
Logical Volume Managment (LVM)

Interface/Device Drivers

Node Local Storage
(HDD, SSD, NVRAM)

Network

Burst Buffer

HPC Storage
(PFS, Object, ...)

Archive

Darshan-Logfiles
Repository/Directory

```
$ 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

# Toward Understanding Workflow I/O

**Combine workflow descriptions with monitoring information from Darshan/TOKIO, etc.**
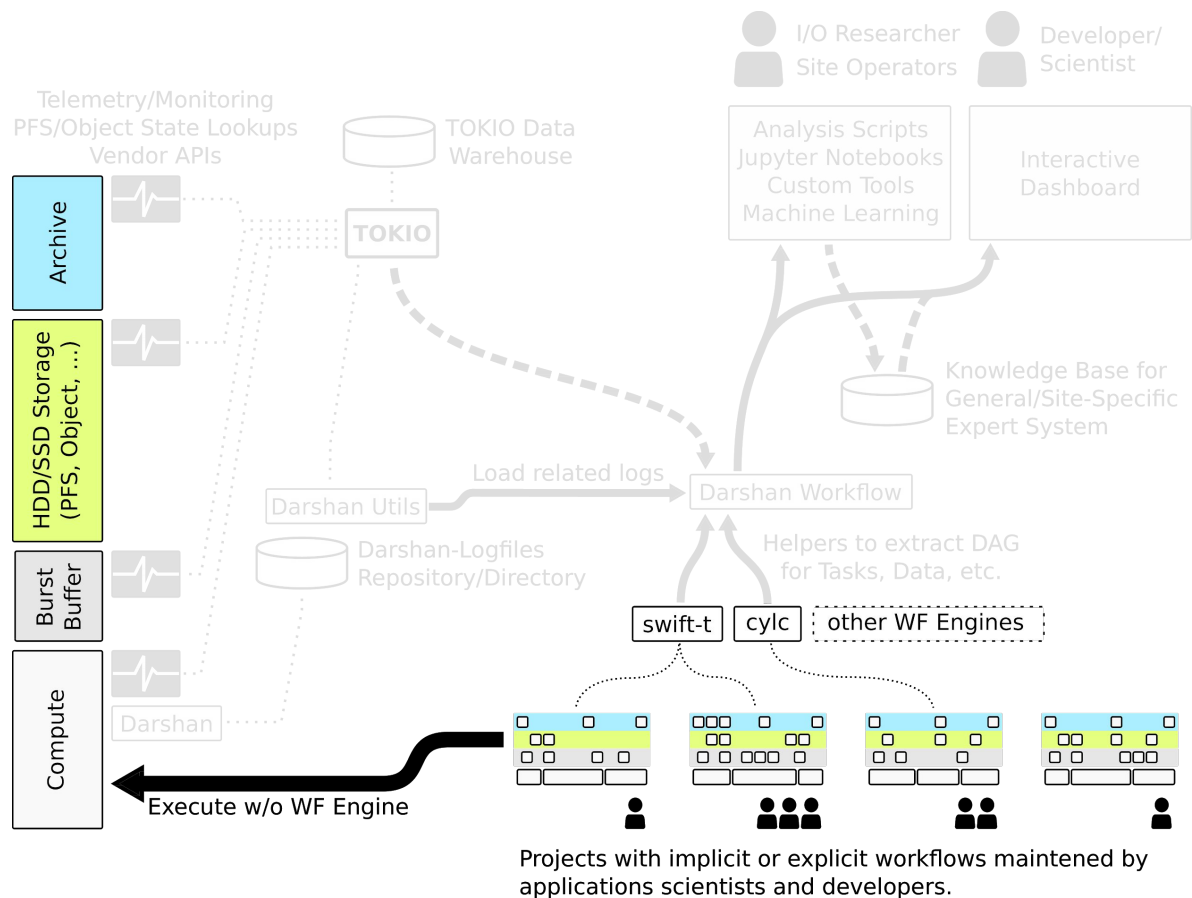
Benefits:

      Insight useful for operating decisions and system design

      Communication with users, relatable to their scientific process
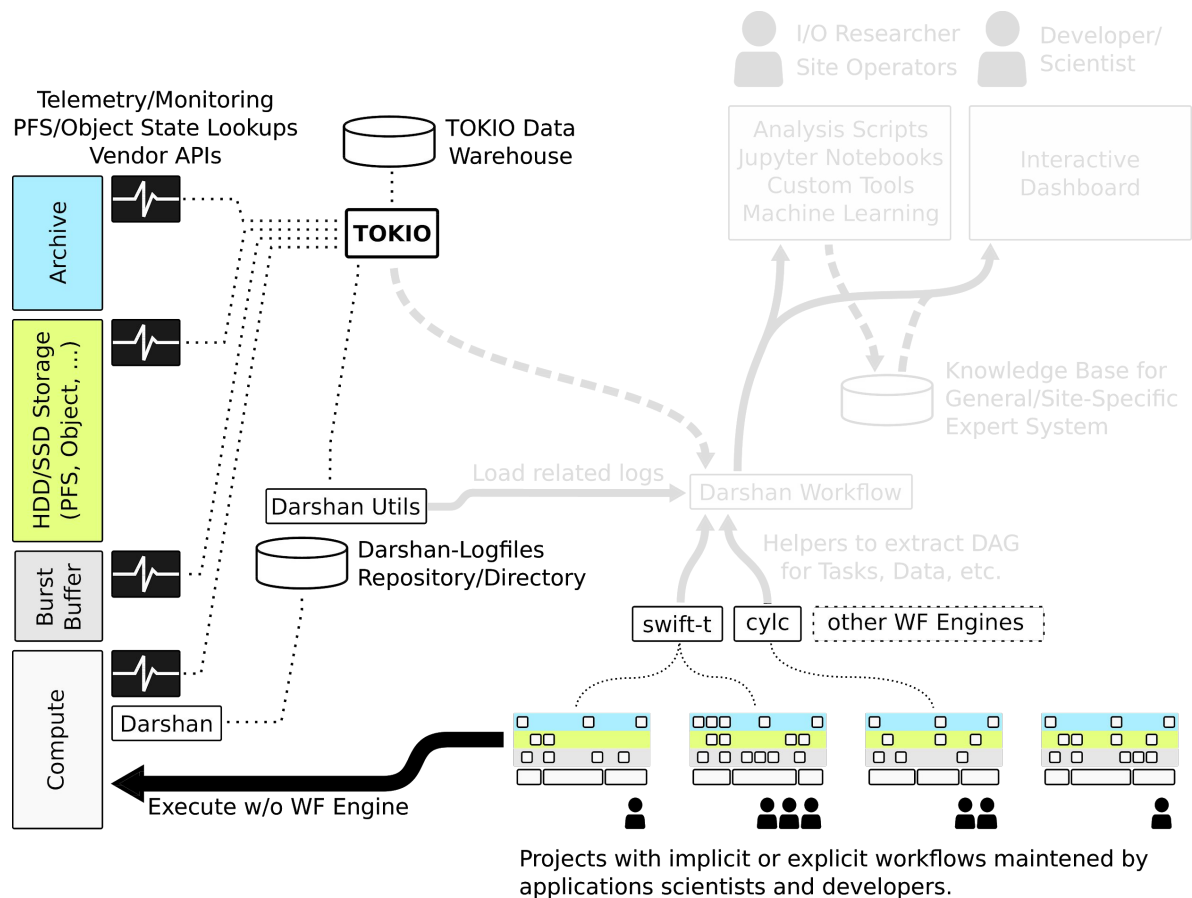
      Source of information for smarter systems

Requirements:

      Support multiple workflow engines as communities use different tools across difference sites

      Explore convenient toolchain for researchers and operators

      User facing component to communicate advice
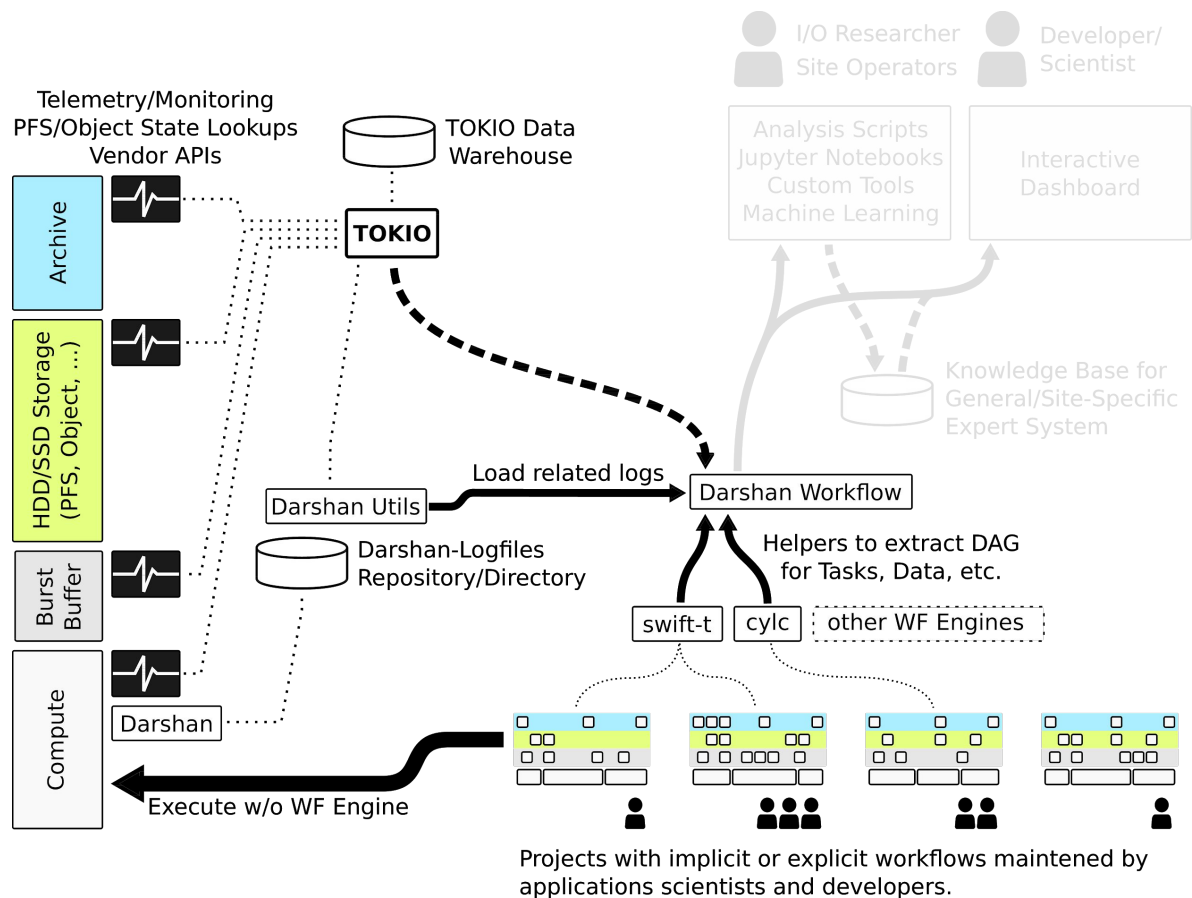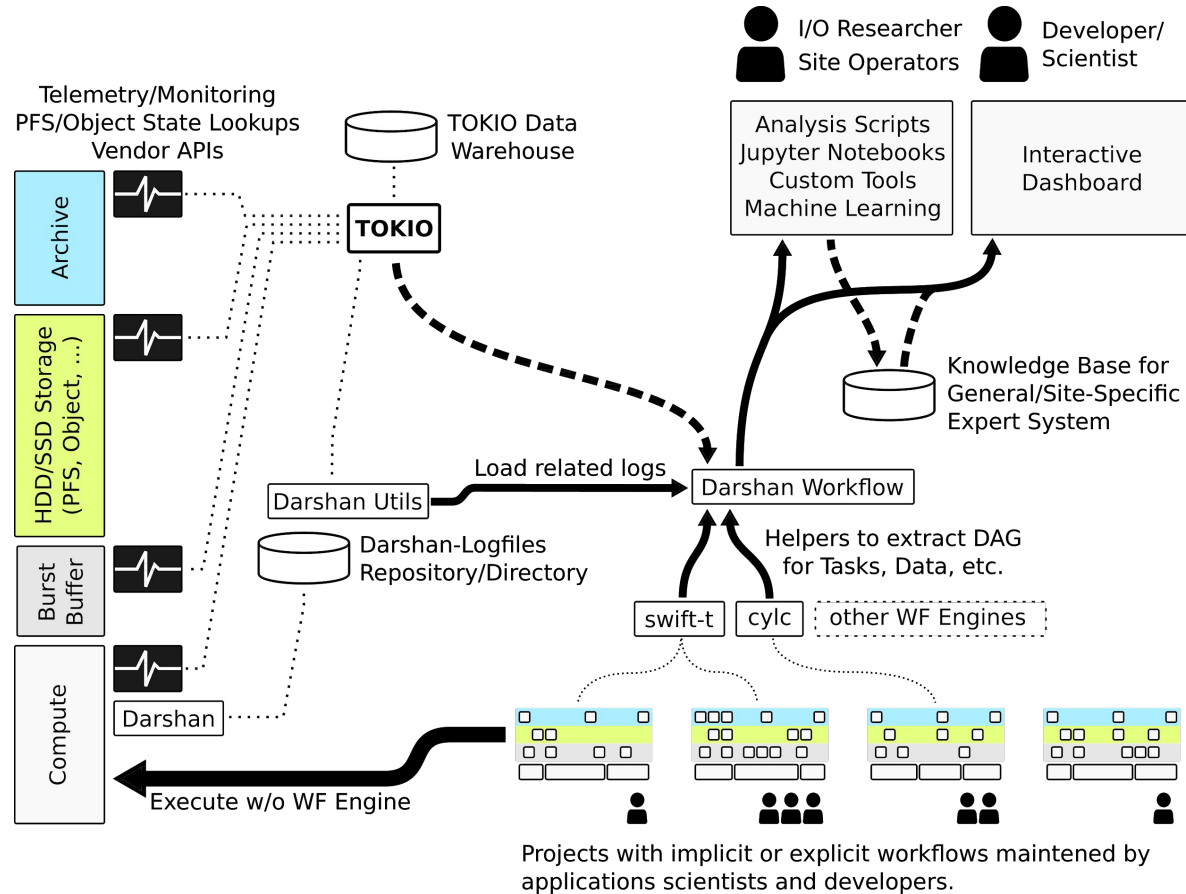
# Architecture for Augmenting I/O in Workflows



I/O Researcher
Site Operators

Developer/
Scientist

Telemetry/Monitoring
PFS/Object State Lookups
Vendor APIs

TOKIO Data
Warehouse

Analysis Scripts
Jupyter Notebooks
Custom Tools
Machine Learning

Interactive
Dashboard

Archive

TOKIO

HDD/SSD Storage
(PFS, Object, …)

Knowledge Base for
General/Site-Specific
Expert System

Load related logs

Darshan Utils

Darshan Workflow

Burst
Buffer

Darshan-Logfiles
Repository/Directory

Helpers to extract DAG
for Tasks, Data, etc.

swift-t    cylc    other WF Engines

Compute

Darshan

Execute w/o WF Engine

Projects with implicit or explicit workflows maintained by
applications scientists and developers.

# Architecture for Augmenting I/O in Workflows



Telemetry/Monitoring
PFS/Object State Lookups
Vendor APIs

TOKIO Data
Warehouse

TOKIO

Archive

HDD/SSD Storage
(PFS, Object, ...)

Burst Buffer

Compute

Darshan

Execute w/o WF Engine

Darshan Utils

Darshan-Logfiles
Repository/Directory

Load related logs

Darshan Workflow

I/O Researcher
Site Operators

Developer/
Scientist

Analysis Scripts
Jupyter Notebooks
Custom Tools
Machine Learning

Interactive
Dashboard

Knowledge Base for
General/Site-Specific
Expert System

Helpers to extract DAG
for Tasks, Data, etc.

swift-t    cylc    other WF Engines

Projects with implicit or explicit workflows maintained by
applications scientists and developers.

# Architecture for Augmenting I/O in Workflows



Telemetry/Monitoring
PFS/Object State Lookups
Vendor APIs

TOKIO Data Warehouse

**TOKIO**

Archive

HDD/SSD Storage
(PFS, Object, ...)

Burst Buffer

Compute

Darshan

Darshan Utils

Load related logs

Darshan-Logfiles
Repository/Directory

I/O Researcher
Site Operators

Developer/
Scientist

Analysis Scripts
Jupyter Notebooks
Custom Tools
Machine Learning

Interactive
Dashboard

Knowledge Base for
General/Site-Specific
Expert System

Darshan Workflow

Helpers to extract DAG
for Tasks, Data, etc.

swift-t    cylc    other WF Engines

Execute w/o WF Engine

Projects with implicit or explicit workflows maintained by
applications scientists and developers.

# Architecture for Augmenting I/O in Workflows



Telemetry/Monitoring
PFS/Object State Lookups
Vendor APIs

TOKIO Data Warehouse

TOKIO

Archive

HDD/SSD Storage (PFS, Object, …)

Burst Buffer

Compute

Darshan

Darshan Utils

Darshan-Logfiles Repository/Directory

Load related logs

Execute w/o WF Engine

I/O Researcher Site Operators

Developer/ Scientist

Analysis Scripts
Jupyter Notebooks
Custom Tools
Machine Learning

Interactive Dashboard

Knowledge Base for General/Site-Specific Expert System

Darshan Workflow

Helpers to extract DAG for Tasks, Data, etc.

swift-t     cylc     other WF Engines

Projects with implicit or explicit workflows maintained by applications scientists and developers.

# Case Study
# & Demonstration

Example Workflow

Research Perspective

User Perspective

Start
Outer Loops
Inner Loops
check()
then / else
f()
g()
sum()

Task
Data
Spawn Task
Data wait/write

```
int X = 50, Y = 50;
int A[][];
int B[];

foreach x in [0:X-1] {
  foreach y in [0:Y-1] {
    if (check(x, y)) {            // mask a region which gets computed
      A[x][y] = g(f(x), f(y));    // compute result for this cell (a physics process)
    } else {
      A[x][y] = 0;                // default for skipped cells
    }
  }
  B[x] = sum(A[x]);              // compute some aggregate metric
}
```

http://swift-lang.org

prep
2021

model
2021

model
2022

post
2021

model
2023

post
2022

post
2023

stop
2023

my.suite

ight-click on nodes to control family grouping

pq@mcswl209:~/ANL/darshan-workflow/demo/workflow-cylc-example

File  Edit  View  Search  Terminal  Tabs  Help

pq@mcswl209:~  ✕  pq@mcswl209:~/ANL/darshan-workflow/demo/workflow-cylc...  ✕

```
$ ./02-visualize.sh

pq@mcswl209:~/ANL/darshan-workflow/demo/workflow-cylc-example
$ ./run.sh
REGISTER my.suite: /home/pq/ANL/darshan-workflow/demo/workflow-cylc-example/suites/test
my.suite | A first Cylc suite. | ~/ANL/darshan-workflow/demo/workflow-cylc-example/suites/test
Valid for cylc-UNKNOWN
/home/pq/ANL/darshan-workflow/devel/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)> succeede
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/

# Perspective for **I/O Research** and **Site Operating**?

**Interactive** Tools/Dashboards to ease navigating **overwhelming  amounts of log data**, with "algebra"-like semantics for convenient aggregation of multiple tasks, data objects or pipelines.

Python Library for use in, e.g.,  **jupyter notebooks**, to draft/prototype/provide **templates** for more sophisticated and **reproducible** analysis.

JavaScript Packages (NPM) for visualisation/tools allowing easy **reuse** in custom tools , jupyter notebooks (widget plugins), and dashboards (e.g., Grafana).

## Augmenting Workflow I/O with Darshan/TOKIO

Scientific discovery increasingly depends on complex workflows consisting of multiple phases and sometimes millions of parallelizeable 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 behaivior 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 convienience methods and common visualisations such as the `task_summary()` and `show_graph()` methods.

```
In [3]:    1   wf.tasks_summary()
```
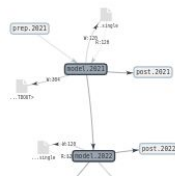
```
model.2021      Records: 2     Layers: POSIX, STDIO
model.2022      Records: 2     Layers: POSIX, STDIO
post.2021       Records: 0     Layers:
model.2023      Records: 2     Layers: POSIX, STDIO
post.2022       Records: 0     Layers:
post.2023       Records: 0     Layers:
stop.2023       Records: 0     Layers:
prep.2021       Records: 0     Layers:
```
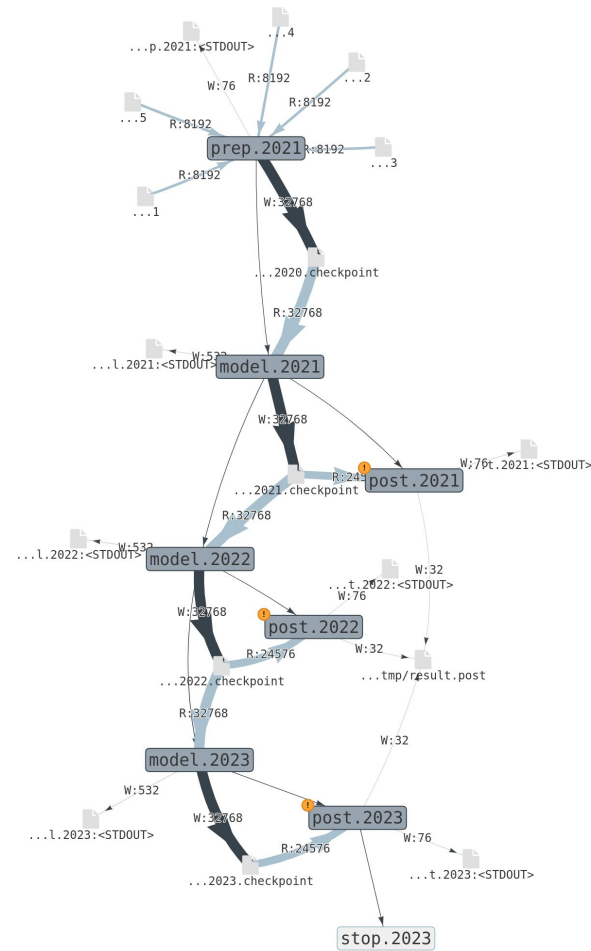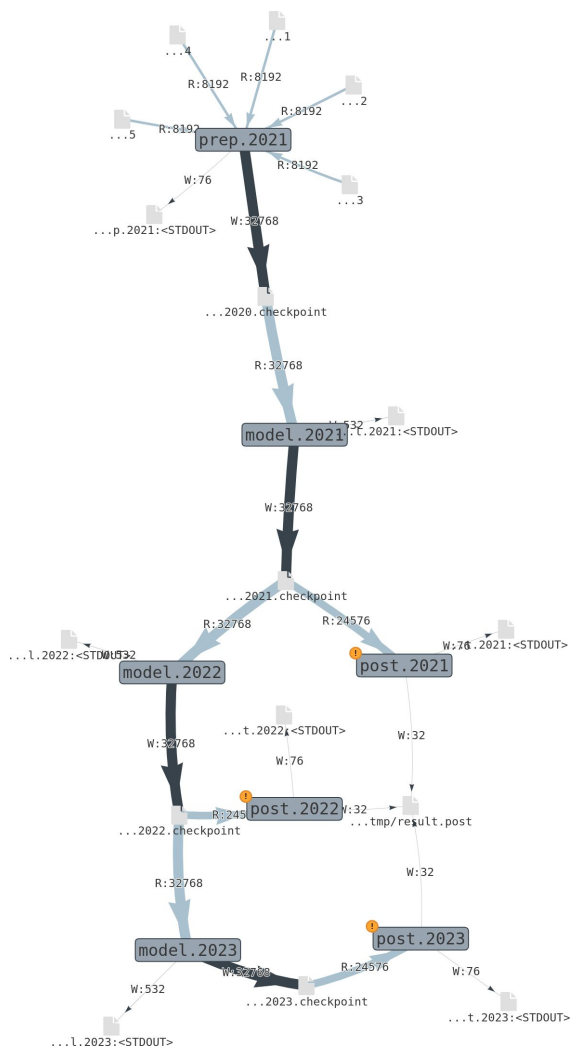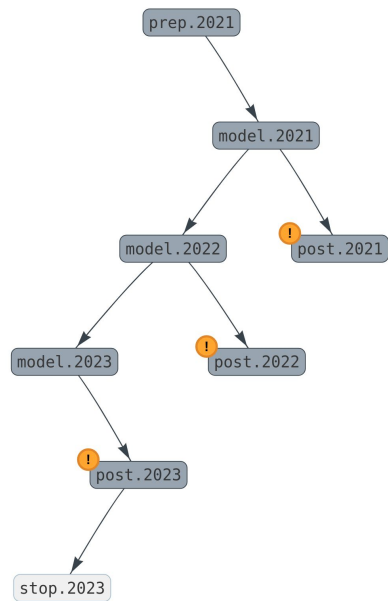
### Displaying and Interacting with Workflow

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 convienient.

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

TODO: GRAPH IpythonWidget

# Communication with Scientists/Developers

Maintain affinity to scientists perspective
      Stick to relationship of tasks/pipelines used by scientists/developers
      Use intuitiv presentation of data-flow by extending graph of workflow

Interactive to manage complexity
      100s or 1000s of different tasks and files in a workflow
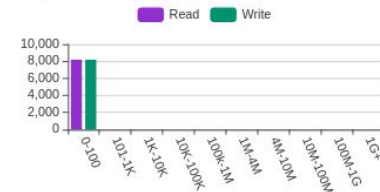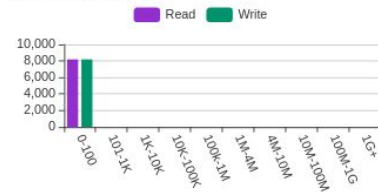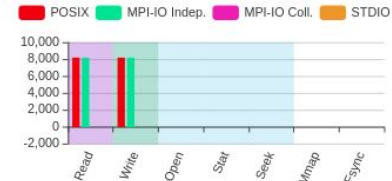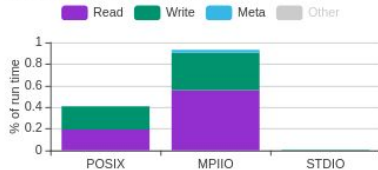      Possibly, millions of log records per task (HTC, UQ)
      Make it easy to aggregate multiple log records

Integration with expert advice
      Human in the loop
      Automatic advisories with machine learning (mid/long-term)

## Workflow Overview



/home/pq/ANL/darshan-workflow/demo/workflow-cylc-example/model

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

### Average I/O cost per process

Legend: Read, Write, Meta, Other



### I/O Operation Counts

Legend: POSIX, MPI-IO Indep., MPI-IO Coll., STDIO



### Access Sizes POSIX

Legend: Read, Write



### Access Sizes MPIIO

Legend: Read, Write



### Most Common Access Sizes
(POSIX and MPI-IO)

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

### File Count Summary
(estimated by POSIX I/O access offsets)

| 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 |

Logfiles:

Task: model.2021
model.2021: darshan-lo
Task: model.2022
model.2022: darshan-lo
Task: post.2021
post.2021: darshan-logs
Task: model.2023
model.2023: darshan-lo
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/2020.c
model.2021: /tmp/2021.c
model.2021: <STDOUT:
model.2022: /tmp/2022.c
model.2022: /tmp/2021.c
model.2022: <STDOUT:
post.2021: /tmp/result.po
post.2021: /tmp/2021.ch
post.2021: <STDOUT>
model.2023: /tmp/2023.c
model.2023: /tmp/2022.c
model.2023: <STDOUT:
post.2022: /tmp/result.po
post.2022: /tmp/2022.ch
post.2022: <STDOUT>
post.2023: /tmp/result.po
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.ch
prep.2021: /tmp/file-raw.
prep.2021: /tmp/file-raw
prep.2021: <STDOUT>

What a **real task**
might look like though…

# Analyzing **Access Patterns**



Input Files

Output Files

In this case diagnostic files otherwise not so clear

# Toward Adaptive I/O Systems

Influence **Job Scheduling** decisions

Support **I/O Middleware**
 Data Placement
 Transformations

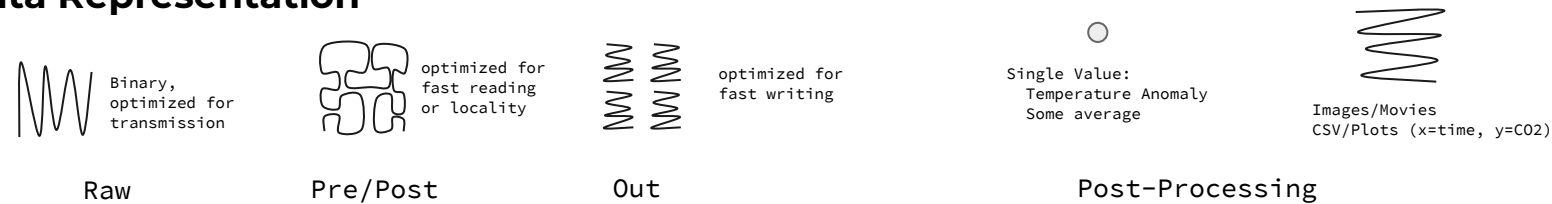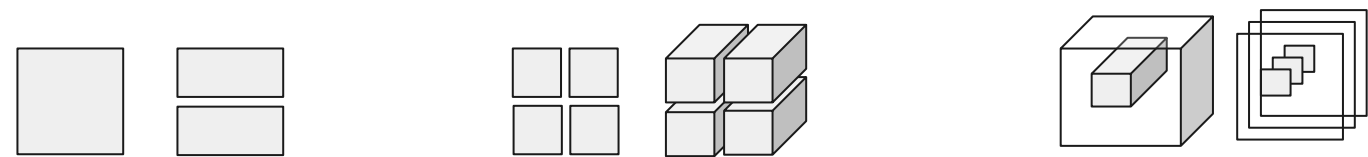# Use Case 1: I/O-Aware Scheduling for Workflows



**I/O in phase**

Job 1
Job 2
Job 3

I/O Activity

Time

**I/O Aware Schedule**

**I/O out of phase**

Job 1
Job 2
Job 3

I/O Activity

Time

# Use Case 2: Benefits for I/O Middleware (1/2)

## Data Representation

Binary, optimized for transmission

optimized for fast reading or locality

optimized for fast writing

Single Value:
Temperature Anomaly
Some average

Images/Movies
CSV/Plots (x=time, y=CO2)

Raw                Pre/Post          Out                            Post-Processing

## Domain Decomposition

## Layout on Storage

# Discussion Summary

**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 instrumentation with a subset of runs

    Collect traces and logs for a training body.

# Thank you!
# Questions?
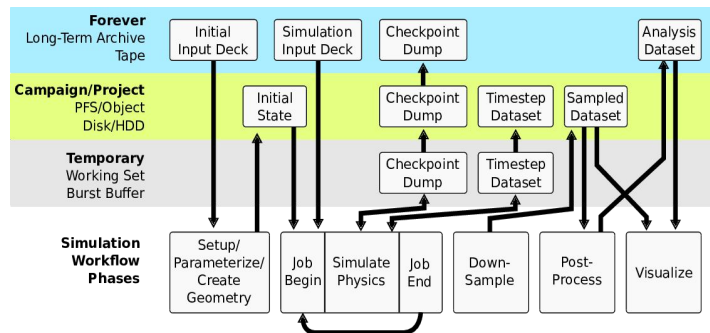
luettgau@dkrz.de

# Disclaimer

# Appendix

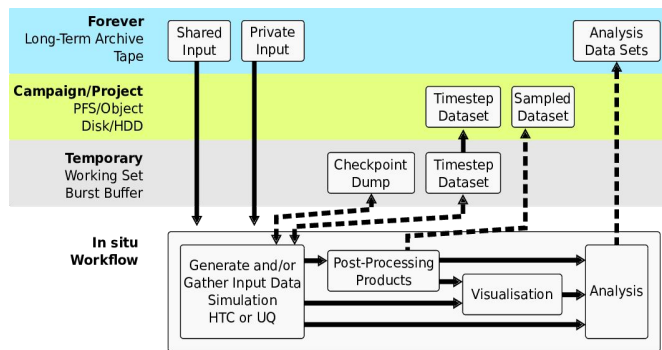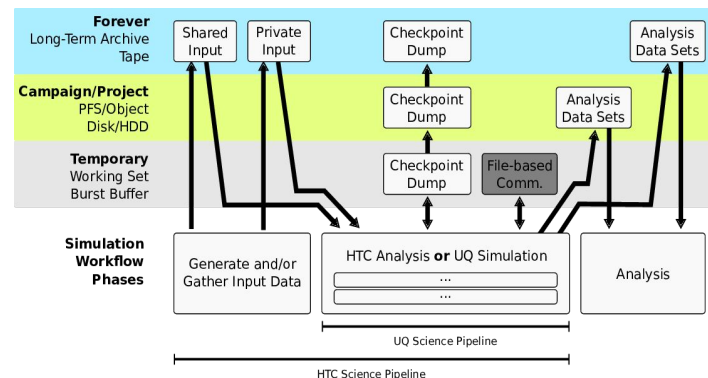Generic HPC Workflows

Example Climate Workflow

# Common Scientific Workflows in HPC
## What makes a workflow?



UQ or HTC

SIM

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.

# 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.

The Cylc Suite Dependency Graph Viewer

```
pq@mcswl209:~/ANL/darshan-workflow/demo/workflow-cylc-example

File  Edit  View  Search  Terminal  Tabs  Help

           pq@mcswl209:~                    ×    pq@mcswl209:~/ANL/darshan-workflow/demo/workflow-cylc...  ×

$ ./02-visualize.sh

pq@mcswl209:~/ANL/darshan-workflow/demo/workflow-cylc-example
$ ./run.sh
REGISTER my.suite: /home/pq/ANL/darshan-workflow/demo/workflow-cylc-example/suites/test
my.suite | A first Cylc suite. | ~/ANL/darshan-workflow/demo/workflow-cylc-example/suites/test
Valid for cylc-UNKNOWN
/home/pq/ANL/darshan-workflow/devel/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)> succeede
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
```
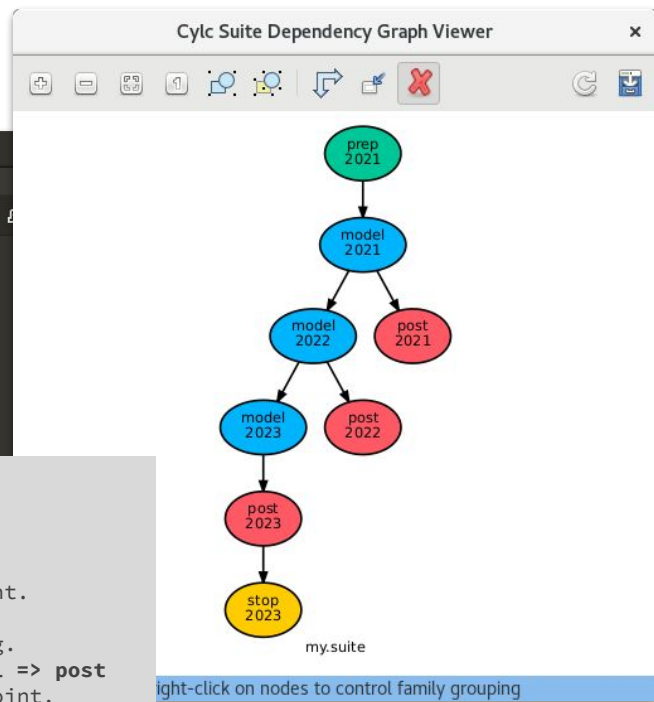
```
[scheduling]
 initial cycle point = 2021
 final cycle point = 2023
 [[dependencies]]
   [[[R1]]]  # Initial cycle point.
    graph = prep => model
   [[[R//P1Y]]]  # Yearly cycling.
    graph = model[-P1D] => model => post
   [[[R1/P0Y]]]  # Final cycle point.
    graph = post => stop

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

https://cylc.github.io/cylc/