

Scalable Observation System (SOS) for Scientific Workflows Project Overview & Discussion

Chad D. Wood

Supervisors: Prof. Allen Malony and Kevin Huck

"So, where is this talk going?"

To advocate and demonstrate a *run-time system* designed to enable the *characterization* and *analysis* of complex scientific *workflow performance* at scale.

- It is reasonable to want to see "information" during application execution
- Information could come from the application as well as from the environment in which the application is executing
- Application:

Performance, problem-specific data and metadata, ...

Environment:

System state, resource usage, runtime properties, ...

Multiple application components may be running together as a workflow, and higher-level workflow behavior might be interesting



Scientific Workflows





- Multiple components with *data flow*
- Complex interactions with *dynamic behavior*
- Components (or entire flows) may be *parallelized* differently
- Offline episodic performance analysis has limited benefits

UNIVERSITY OF OREGON

- Scalable
- Portable
- Easy to use
- Multi-purpose
- Multiple information sources
- Operates at the time of application (workflow) execution
- Supports in situ access
- Low overhead and low intrusion
- Ability to allocation additional resources to control overhead



- Base on a model of a "global" information space
- Utilize database technology
- Utilize MPI high-performance communication
- Build on launch support in scheduler
- Allow for additional (dedicated) resource allocation
- Flexible publishing interface
- SOS architecture





SOSflow forms a functional overlay.



In situ daemon with its local database

UNIVERSITY OF OREGON





SOS lives side by side with your tasks



Dedicated nodes for aggregate databases

UNIVERSITY OF OREGON



Dedicated nodes for analytics processing





Co-located analytics query modules



Independent ranks of analytics engines

UNIVERSITY OF OREGON



Analytics modules form independent communication channels





SOSflow data is continuous and asynchronous

SOSflow: Data Structure



Every value is conserved, with **its full history** and **evolving metadata**.





```
int main(int argc, char *argv[]) {
int
     var_int
                    = 10;
double var_double = 88.8;
char var_string = "Hello, SOS.";
SOS_pub *pub;
SOS_init( &argc, &argv, SOS_ROLE_CLIENT );
pub = SOS_new_pub("demo");
SOS_pack(pub, "example_int", SOS_VAL_TYPE_INT, (SOS_val) var_int
                                                                       ):
SOS_pack(pub, "example_dbl", SOS_VAL_TYPE_DOUBLE, (SOS_val) var_double );
SOS_pack(pub, "example_str", SOS_VAL_TYPE_STRING, (SOS_val) var_string );
SOS_announce(pub);
SOS_publish(pub);
SOS_finalize();
return (EXIT_SUCCESS);
```





SOSflow: Distributed Asynchronous Runtime (Simple)





SOSflow: Distributed Asynchronous Runtime



NERSC
 Cori
 Edison

LLNL

- CAB
- Catalyst
- University of Oregon
 - ACISS

Software:
 OpenMPI
 MPICH
 Slurm
 PBS

Scalable Observation System for Scientific Workflows

SOSflow: Evaluation



Experimental Setup

- **Explore performance of work-in-progress implementation**
- Synthetic and real-world cases
- What is the latency cost of being async?
- Synthetic Sweep of Parameters
 - □ Iterations: 2 to 10, steps of 2
 - □ Size: 100 to 500 unique values per pub, steps of 100
 - Delay: 0.5 to 1.0 second, each 0.1 second





SOSflow: Evaluation

UNIVERSITY OF OREGON

Real-World Scenario

- TAU Instrumented LULESH on Cori
- □ TAU reports results to SOSflow on a timer
- □ LULESH calls SOSflow API directly at iteration
- □ SOSflow gathers metrics from the OS











- Performance improvements
- Integrate more automatic data gathering for node-level metrics
- Support for deep analytics
- Testing with additional real-world workflows and operating environments