



Exploiting Locality in Scientific Workflow System: A Cross-Layer Solution

Dong Dai¹, Robert Ross², Dounia Khaldi³, Yonghong Yan⁴, Matthieu Dorier², Neda Tavakoli¹, and Yong Chen¹



¹Department of Computer Science Texas Tech University, ²Mathematics and Computer Science Division Argonne National Laboratory

³Institute for Advanced Computational Science, Stony Brook University, ⁴Computer Science Department, Oakland University



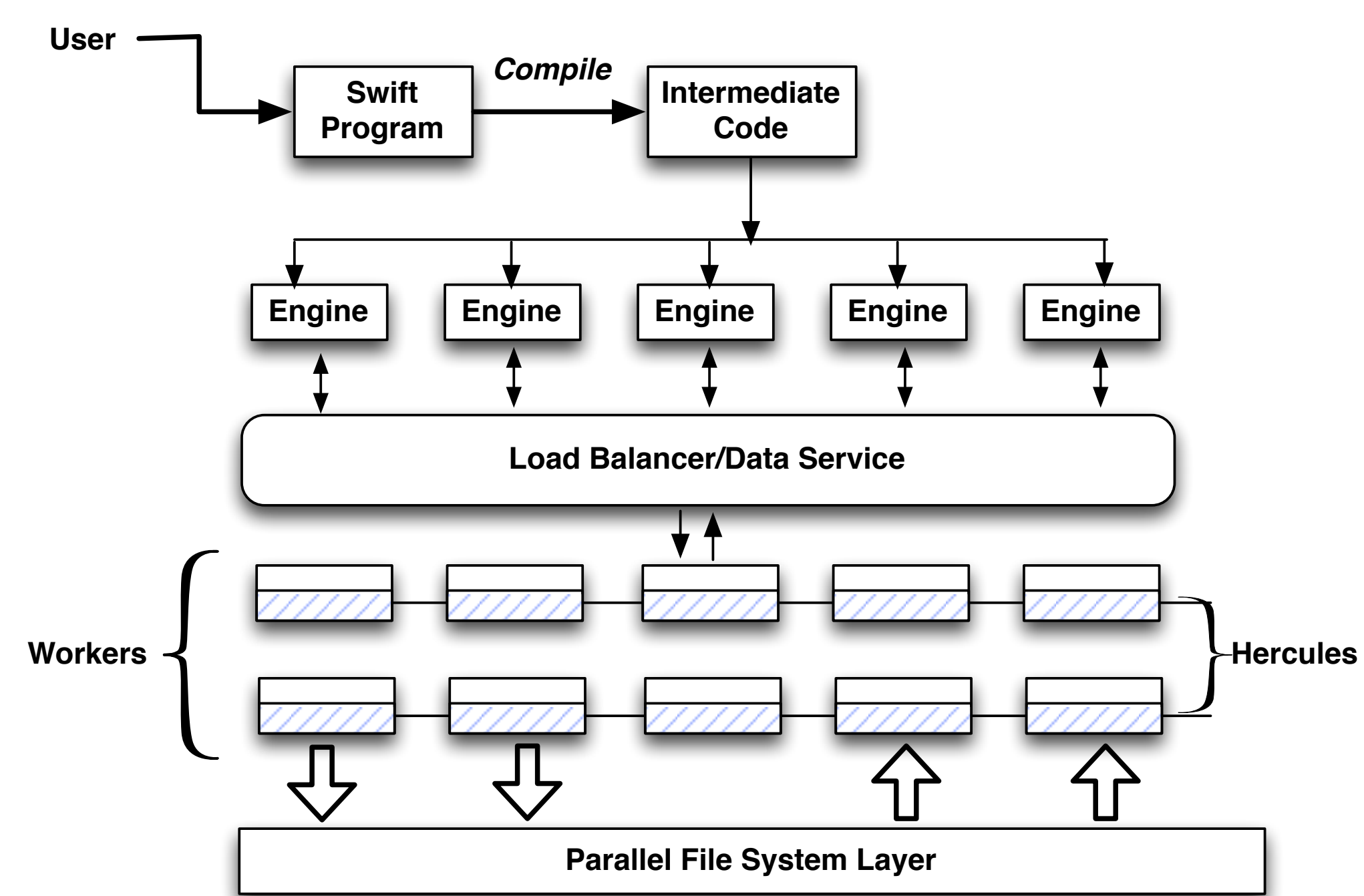
1. MOTIVATION

- Workflows are one of the solutions for complex scientific app.
- Workflows run upon a parallel file system (PFS)
- This leads to significant I/O overhead:
 - Input/Output data needs to be transferred
 - Intermediate data needs to be stored in parallel file system (PFS)
 - Checkpoints need to be stored in PFS
- One solution: **Exploit the Data Locality**
 - If the data is stored in compute nodes, near the workflow tasks, we can obtain locality.
 - Build a files system in compute nodes
 - Use their local resources (memory, SSDs)

2. BACKGROUND

In this research, we take Swift/T and Hercules as the basis to exploit data locality.

- Swift/T is a workflow system.
- Hercules is a compute-node side storage.



3. CHALLENGES OF EXISTING SOLUTION: SWIFT/T WITH HERCULES

```

1 file fs[] = glob("in/*.txt");
2 file outs[];
3
4 app (file out) sort (file in)
5 {
6     "/usr/bin/sort" "-o" out in;
7 }
8
9 app (void signal) mpitask (file f[], string mpi)
10 {
11     "/usr/bin/mpiexec" mpi f;
12 }
13
14 foreach v,i in fs
15 {
16     file y<sprintf("out/%i",i)> = sort(v);
17     outs[i] = y;
18 }
19
20 mpitask(outs, mpiprogram);

```

Take left program as an example, it is hard to exploit data locality in current systems.

- Hercules does not support explicit control over locations. Scientists have to estimate location and tell Swift.
- Swift compiler does not collect metadata regarding data and tasks, like data size or task complexity to help schedule in runtime system.
- Swift runtime does not feedback to storage systems to direct data organization or movements.

4. A CROSS-LAYER SOLUTION TO EXPLOIT DATA LOCALITY

• Location-Aware File System Extension

- Extend POSIX APIs to be location-aware
 - * Add a new mode argument (`S_LOC`) for the `O_CREAT` flag of the `OPEN` function.
 - * Store location information of a file in its `xattr` attributes. Obtained using `getxattr()`.
- Build distributed location service
 - * Store file in *real-loc*, but the location metadata is stored in *orig-hash-loc*.

• Hint-Assist Workflow Compiler

- * `@size` indicates the size of an existing file.
- * `@task` indicates the key parameters for a task.
- * `@compute-complex` indicates the computation cost of a task.
- * `@input-output-ratio` indicates the output size of a task regarding to its inputs.

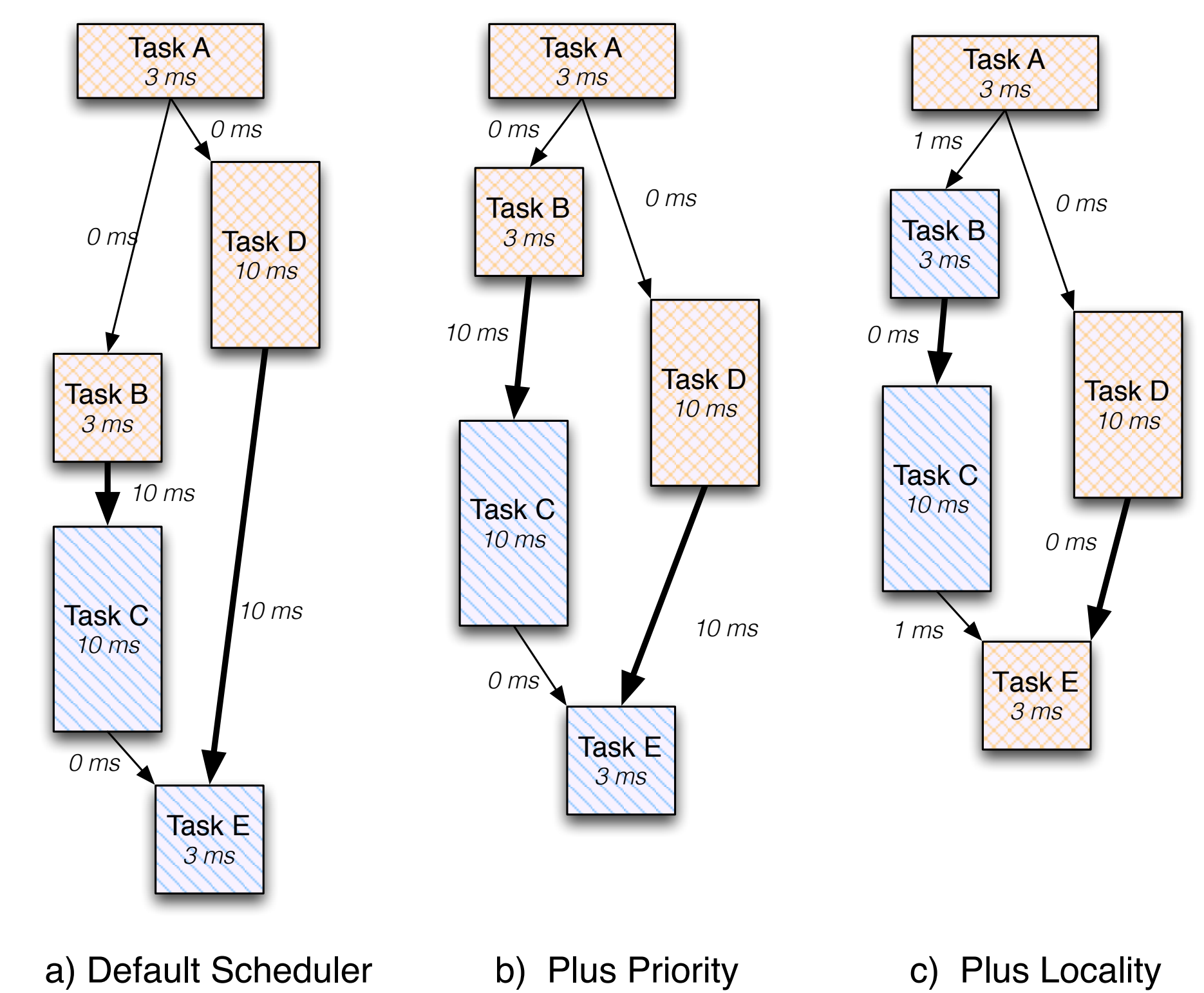
• Locality-Aware Workflow Scheduler

- A heuristic priority for task selection
 - * The task that is far away from the final step in the topological order is assigned higher priority
- Proactive scheduling: schedule “almost-ready” tasks
 - * The scheduling of an almost-ready task will feedback to file system
 - * File system will prefetch intermediate results to target node

5. PRELIMINARY RESULTS

To demonstrate the benefits of the proposed system, we conduct a preliminary evaluation:

- Use a synthetic Swift workflow. Each task is manually controlled to spend a certain time and generate outputs.
- Check performance under three cases: 1) default scheduler; 2) plus priority; 3) plus locality.



6. SUMMARY & FUTURE WORK

• Summary

- We identify the challenges of exploiting data locality in scientific workflow systems even with a compute-node side file system.
- We propose a cross-layer solution to tackle these challenges.

• Future Work

- Full system implementation
- Extend this cross-layer idea into task level
- More detailed performance analysis

REFERENCES

[1] DURO, F. R., BLAS, J. G., ISAILA, F., CARRETERO, J., WOZNIAC, J., AND ROSS, R. Exploiting Data Locality in Swift/T Workflows Using Hercules. In *Proc. NESUS Workshop* (2014).

[2] WOZNIAC, J. M., WILDE, M., AND FOSTER, I. T. Language Features for Scalable Distributed-Memory Dataflow Computing. In *Data-Flow Execution Models for Extreme Scale Computing (DFM), 2014 Fourth Workshop on* (2014), IEEE, pp. 50–53.

ACKNOWLEDGE

This material is based upon work supported by the National Science Foundation under grant CCF-1409946 and CNS-1338078.