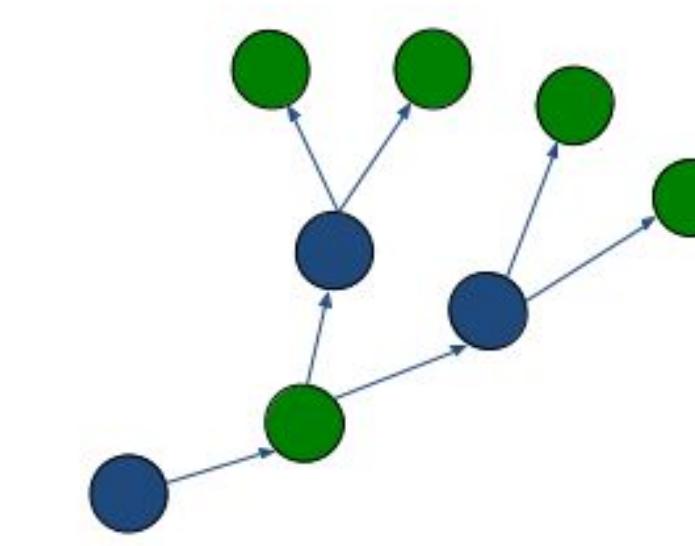




# CSSI Element: DataSwarm: TaskVine: A User-Level Framework for Data Intensive Scientific Applications

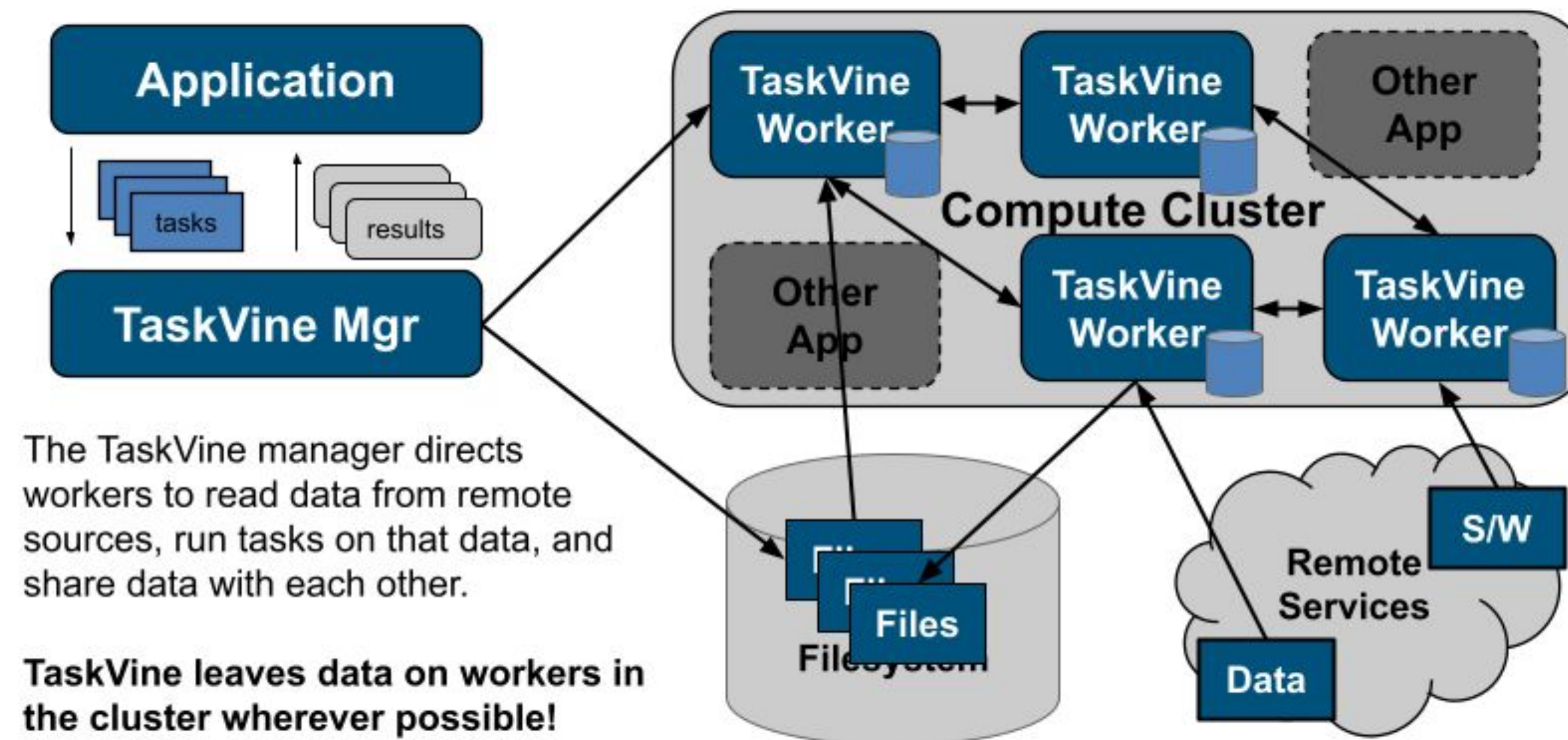
PI: Douglas Thain, University of Notre Dame, Award #: 1931348



# TaskVine

<http://cctools.readthedocs.io>

Many scientific applications are expressed as high-throughput **workflows** that consist of large graphs of data assets and tasks to be executed on large parallel and distributed systems. A challenge in executing these workflows is **managing data**: both datasets and software must be efficiently distributed to cluster nodes; intermediate data must be conveyed between tasks; output data must be delivered to its destination. Scaling problems result when these actions are performed in an uncoordinated manner on a shared filesystem. **TaskVine** is a system for exploiting the aggregate local **storage and network capacity** of a large cluster. TaskVine tracks the lifetime of data in a workflow – from archival sources to final outputs-- making use of local storage to distribute, and re-use data wherever possible.



```
# Example of key TaskVine operations
import taskvine as vine
m = vine.Manager(9123)

# File objects are first class citizens.
file = m.declareFile("mydata.txt")
buffer = m.declareBuffer("Some literal data")
url = m.declareURL("https://nd.edu/data.tar.gz")
temp = m.declareTemp();

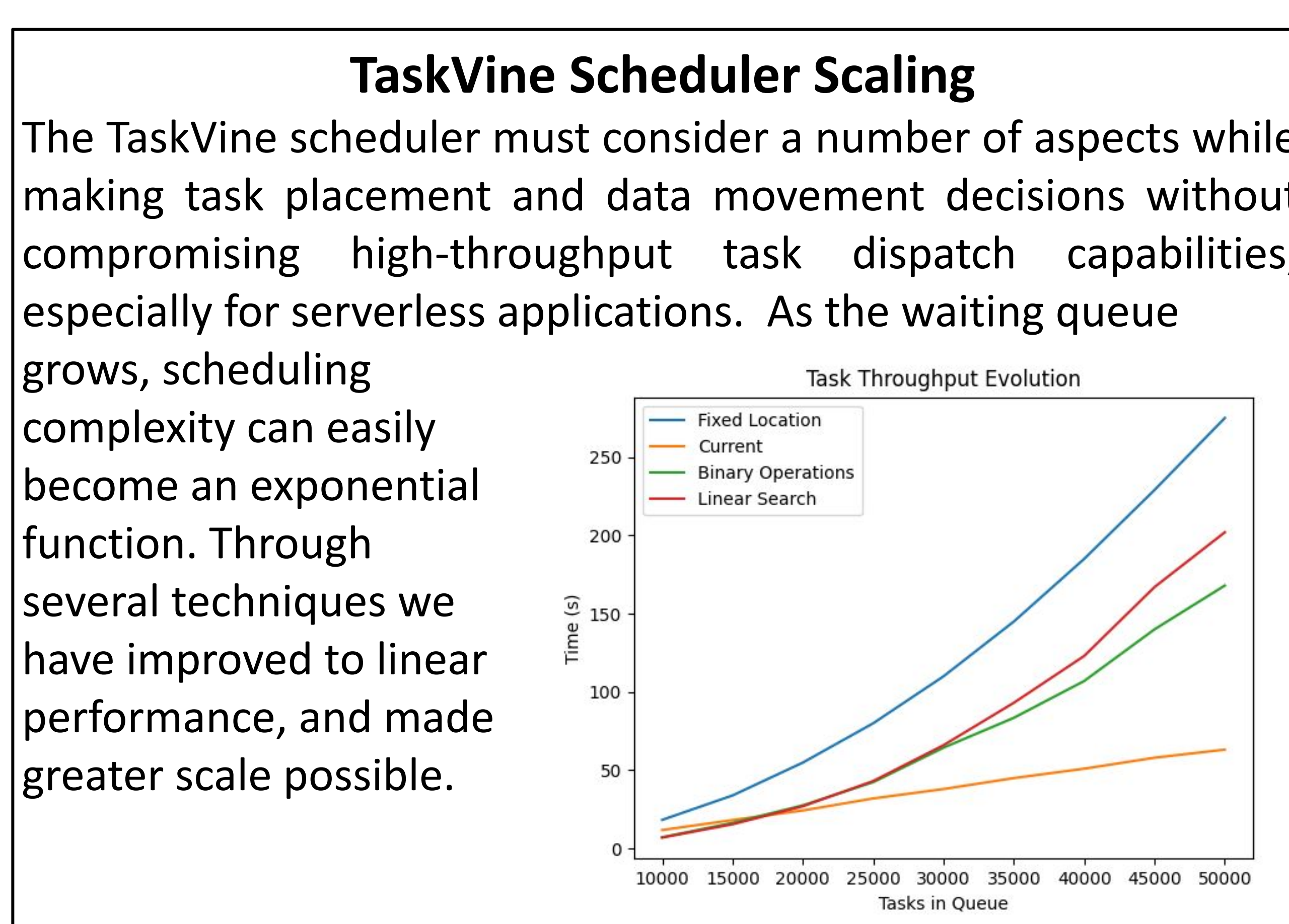
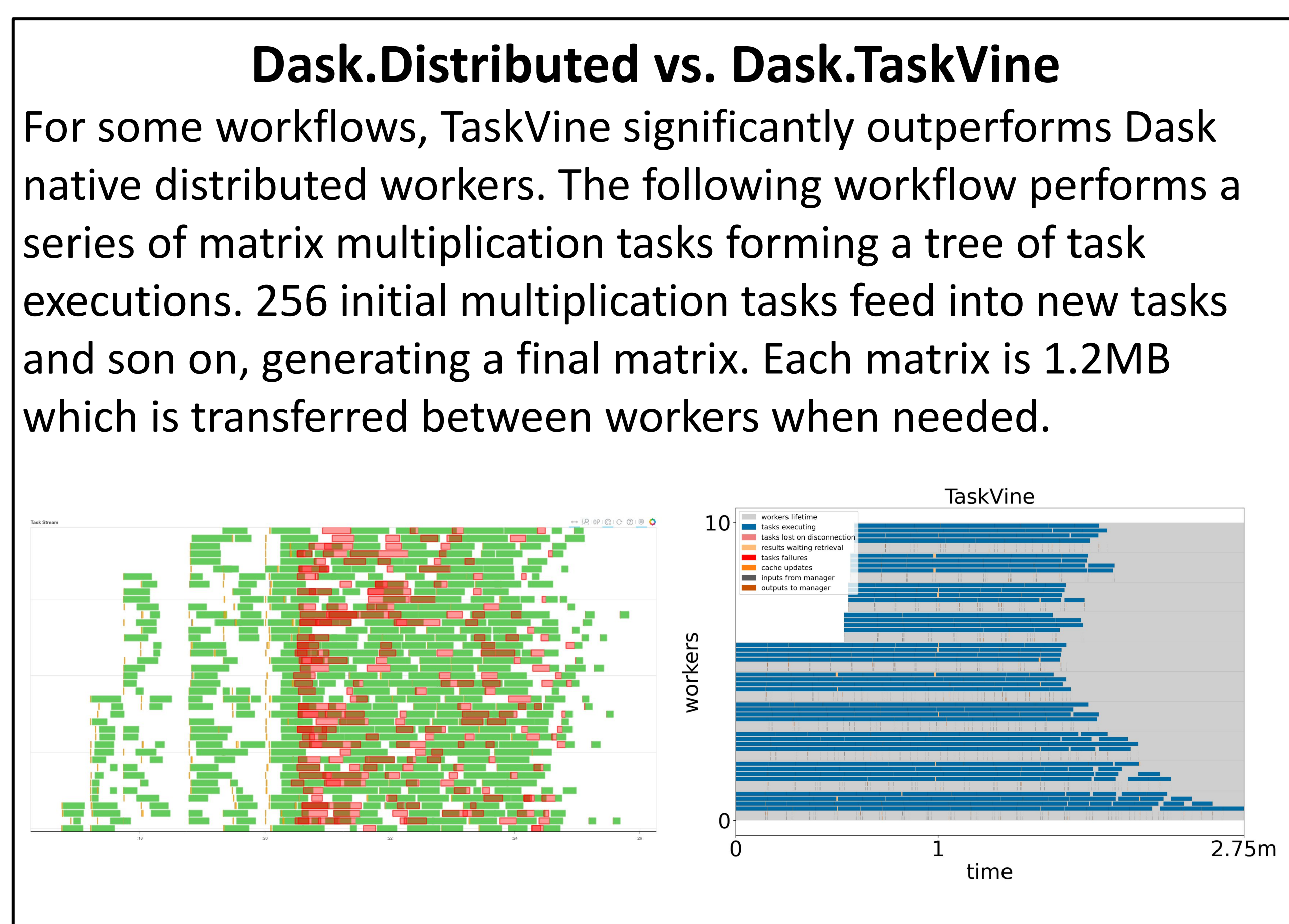
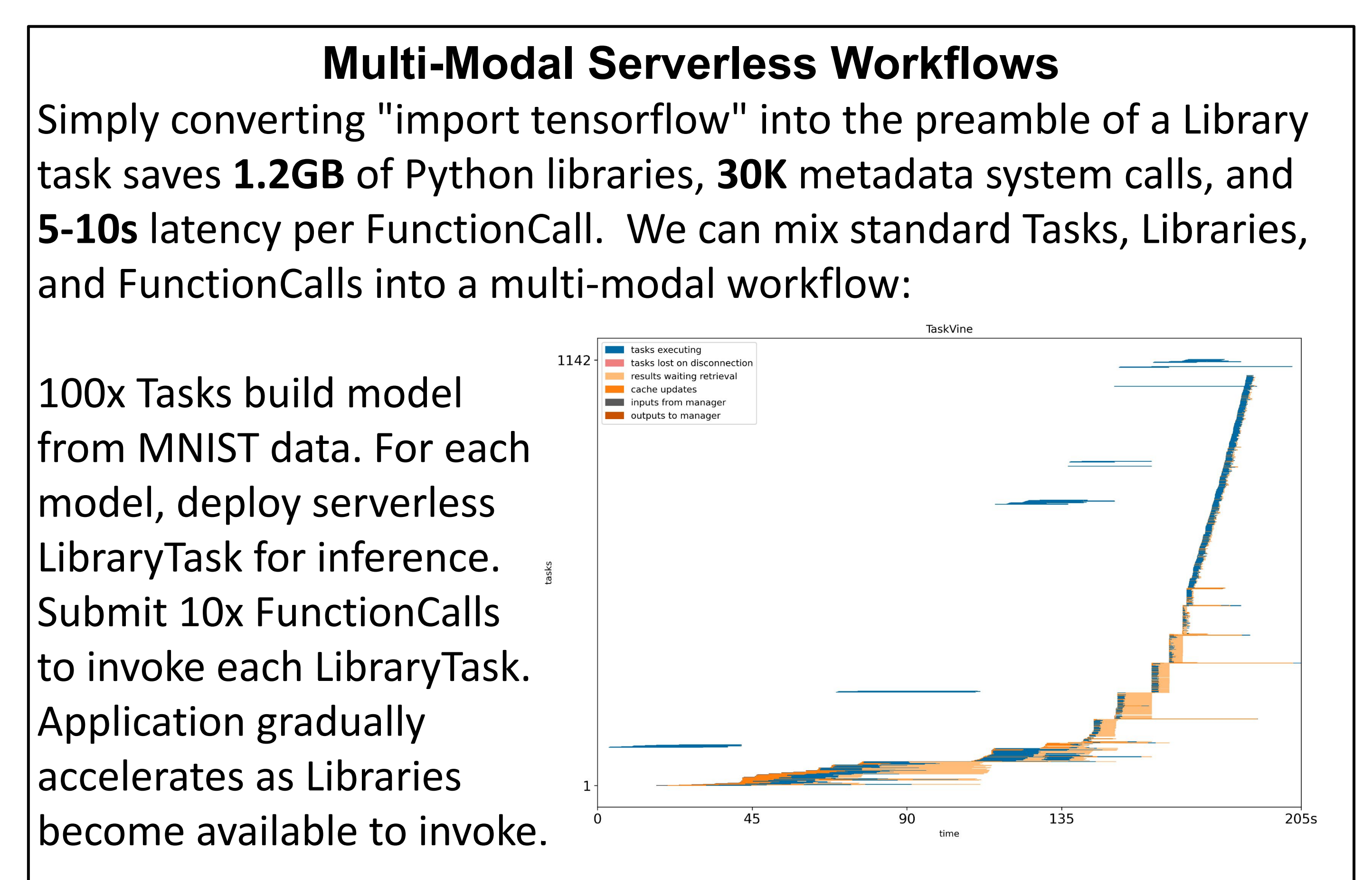
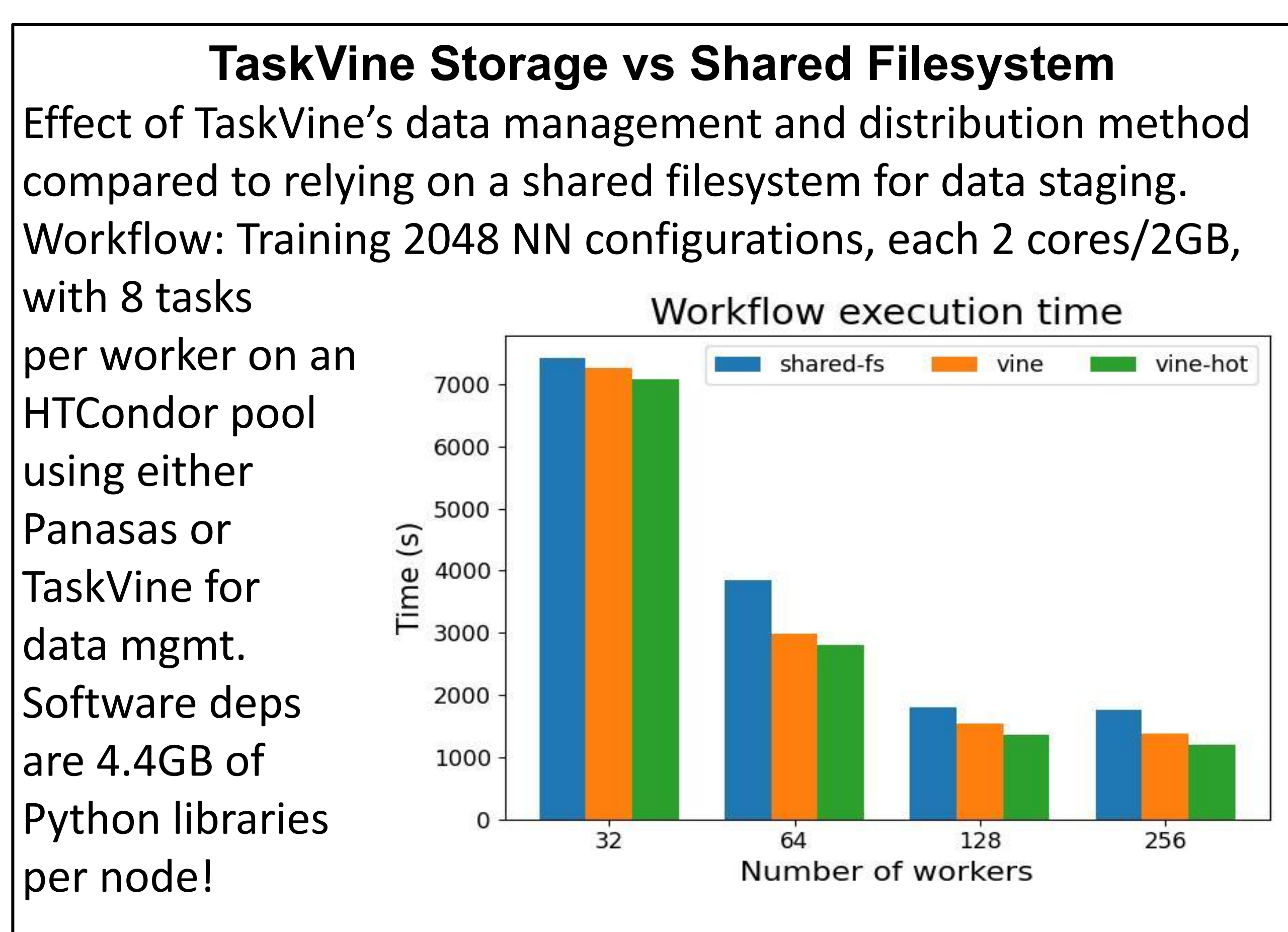
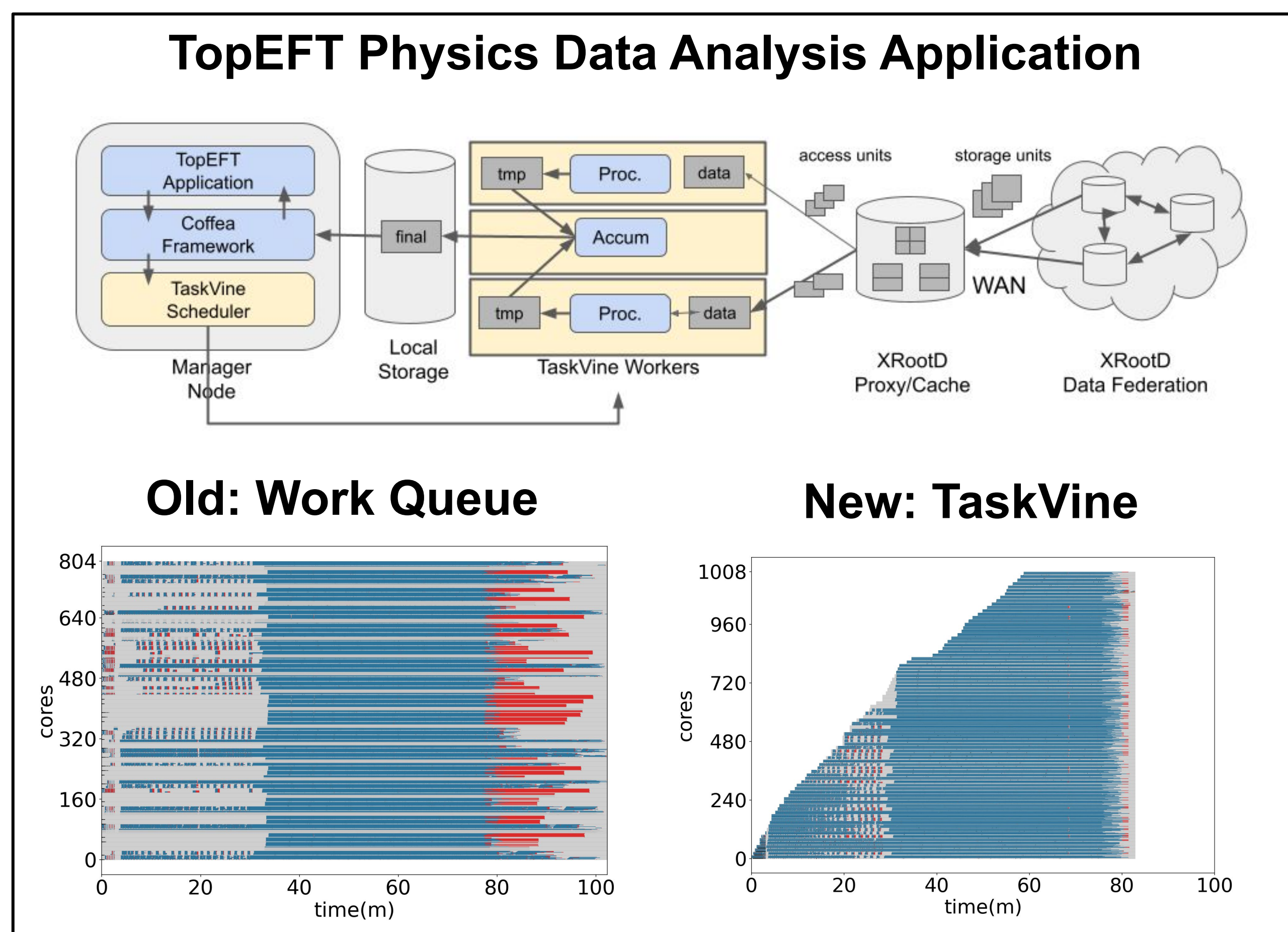
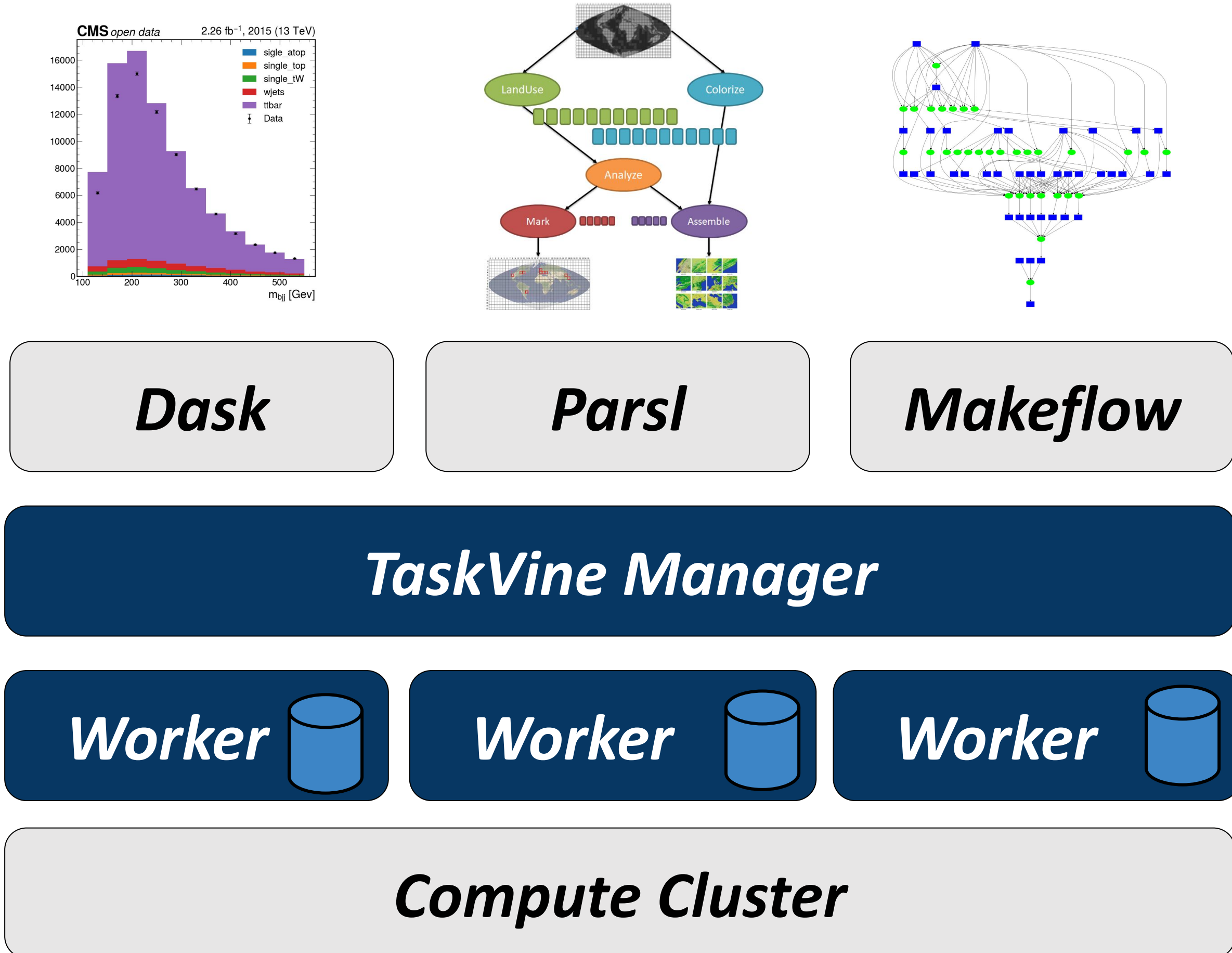
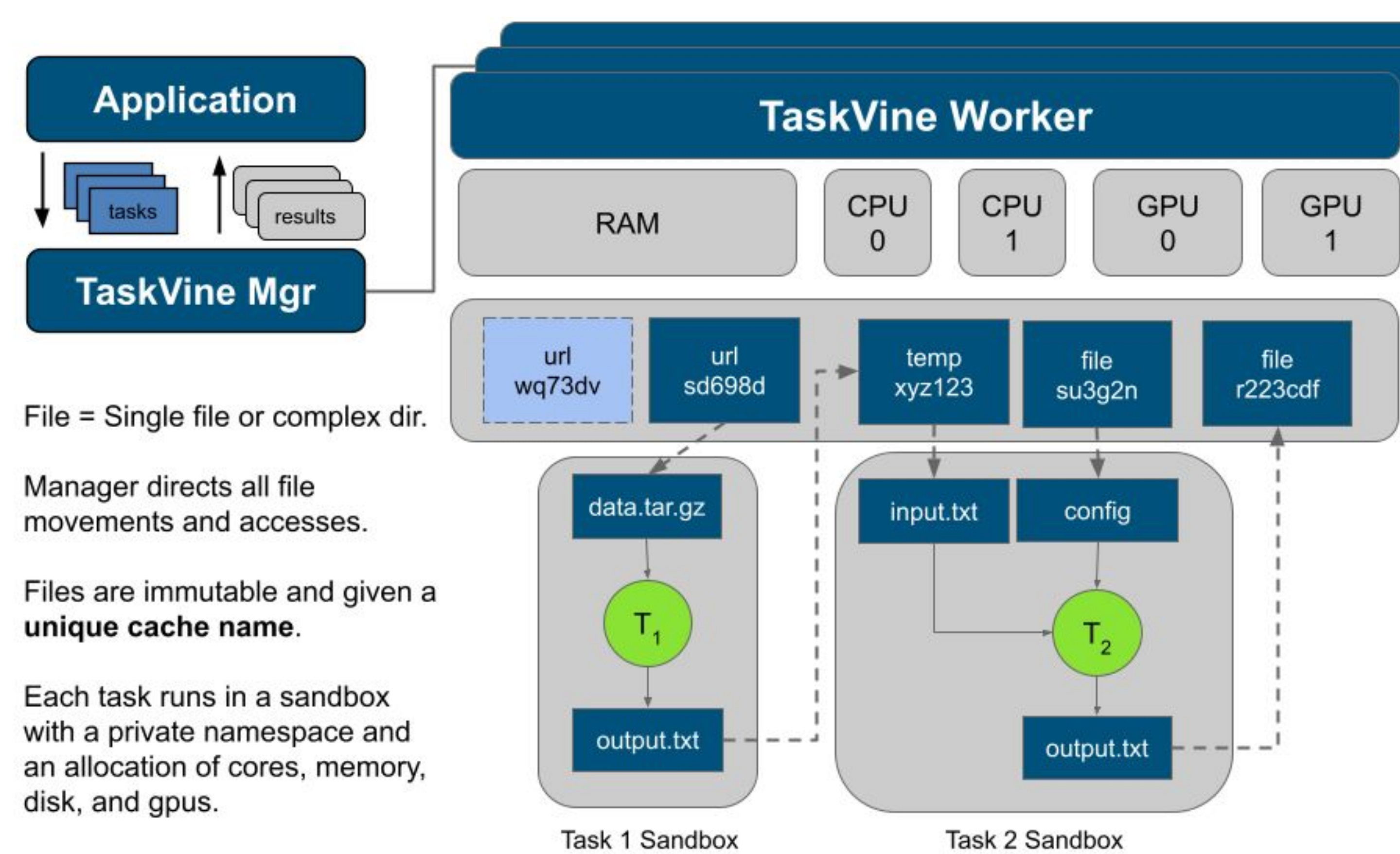
# Perform a standard transformation on a file
data = m.declareUntar( url )
software = m.declarePoncho( package )

# Submit a standard executable task.
task = vine.Task("mysim.exe -p 50 input.data -o output.data")

t.add_input(url, "input.data")
t.add_output(temp, "output.data")
t.set_cores(4)
t.set_memory(2048)
t.set_disk(100)
taskid = m.submit(t)

# Submit a Python function execution
t = vine.PythonTask(
    simulate_func, molecule, parameters)
taskid = m.submit(t)

# Wait for any task to complete.
t = m.wait()
print( t.output )
```



Barry Sly-Delgado, Thanh Son Phung, Colin Thomas, David Simonetti, Andrew Hennesse, Ben Tovar, and Douglas Thain, **"TaskVine: Managing In-Cluster Storage for High Throughput Data Intensive Workflows"**, WORKS Workshop at Scupercomputing, November 2023.

## Cooperative Computing Lab at the University of Notre Dame

