



# PRUNE: A Preserving Run Environment for Reproducible Scientific Computing



UNIVERSITY OF  
NOTRE DAME

# Reproducibility

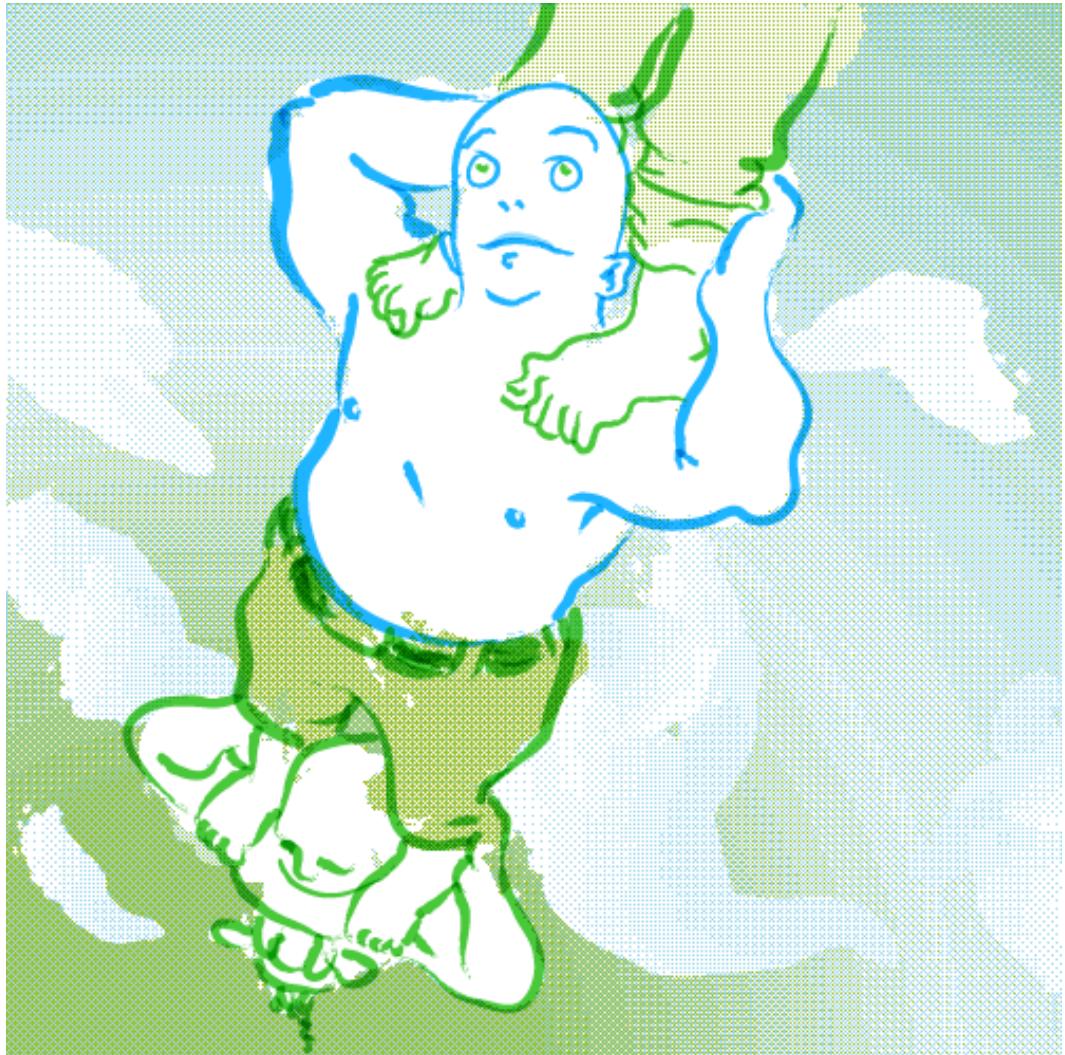
---



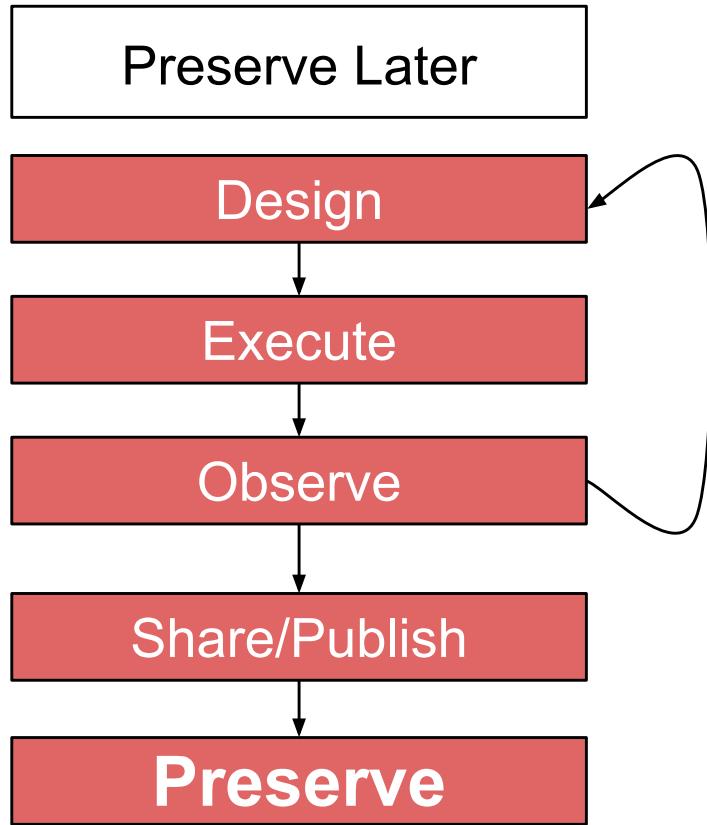
- "[An article about computational science in a scientific publication is not the scholarship itself, it is merely advertising of the scholarship. The actual scholarship is the complete software development environment and the complete set of instructions]"  
–Jon Claerbout

# Verify and Extend

- Don't re-invent the wheel
- Stand on the shoulders of giants

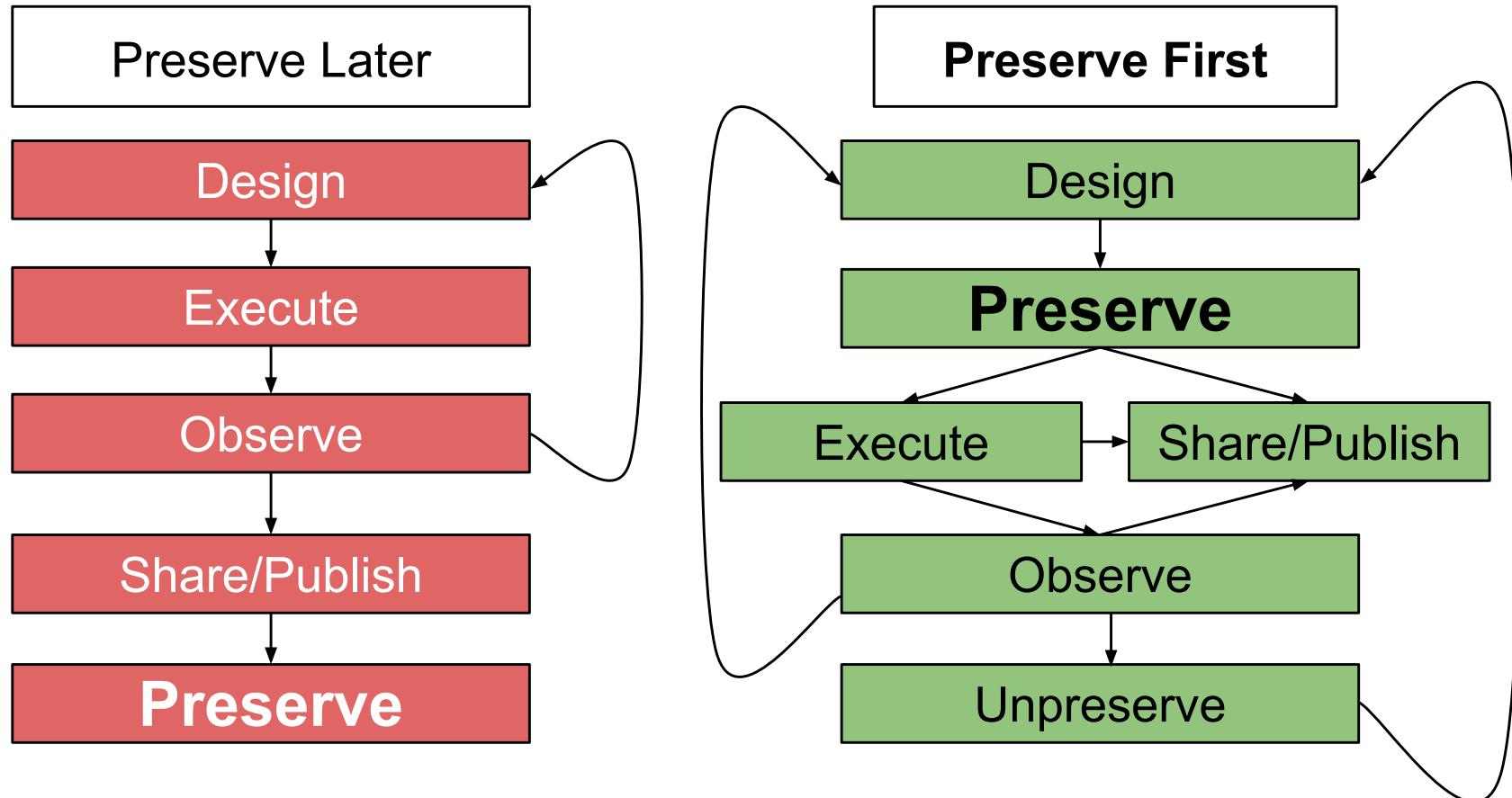


# Accepted philosophy



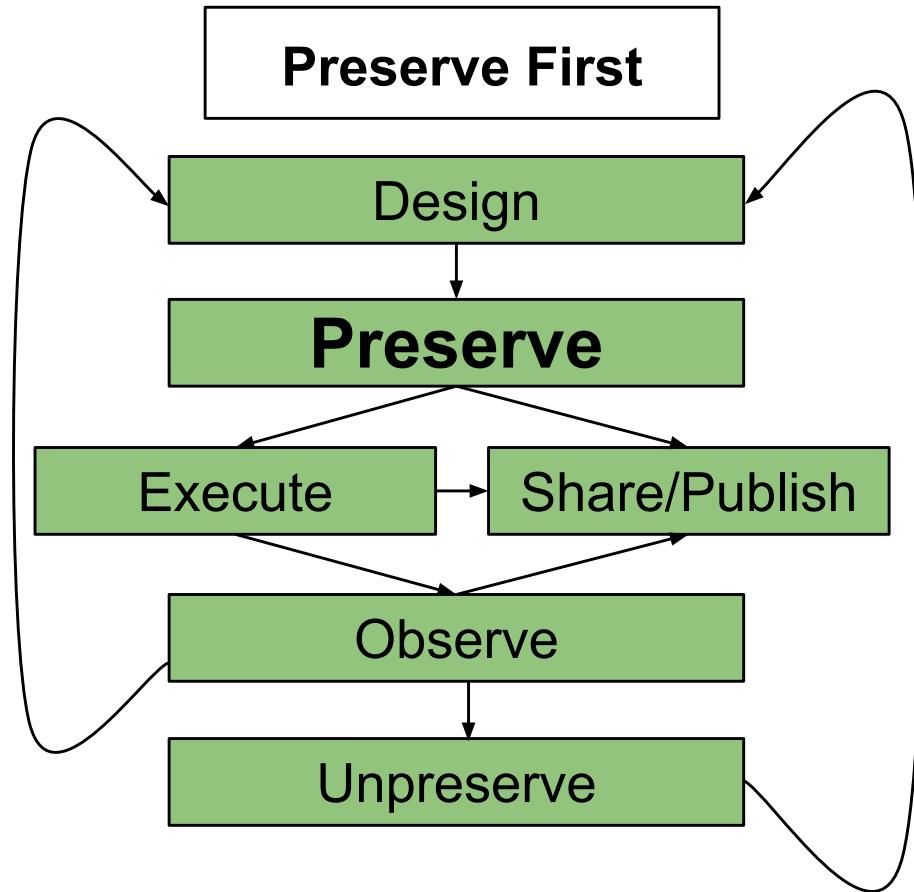
- Libraries
- Hardware
- Network
- System Administrators
- Remote Collaborators
- Graduated Students

# Proposed philosophy

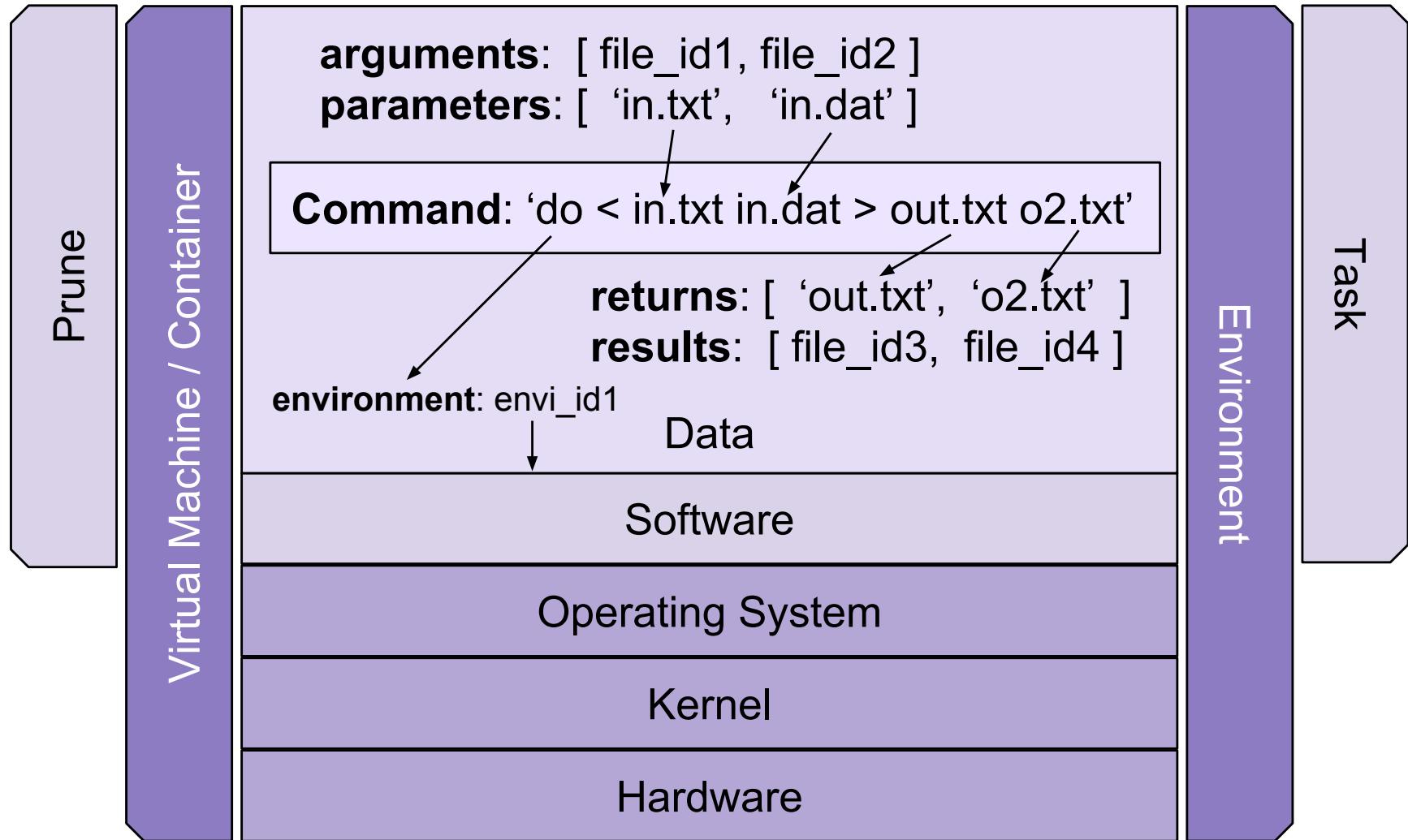


# Differences

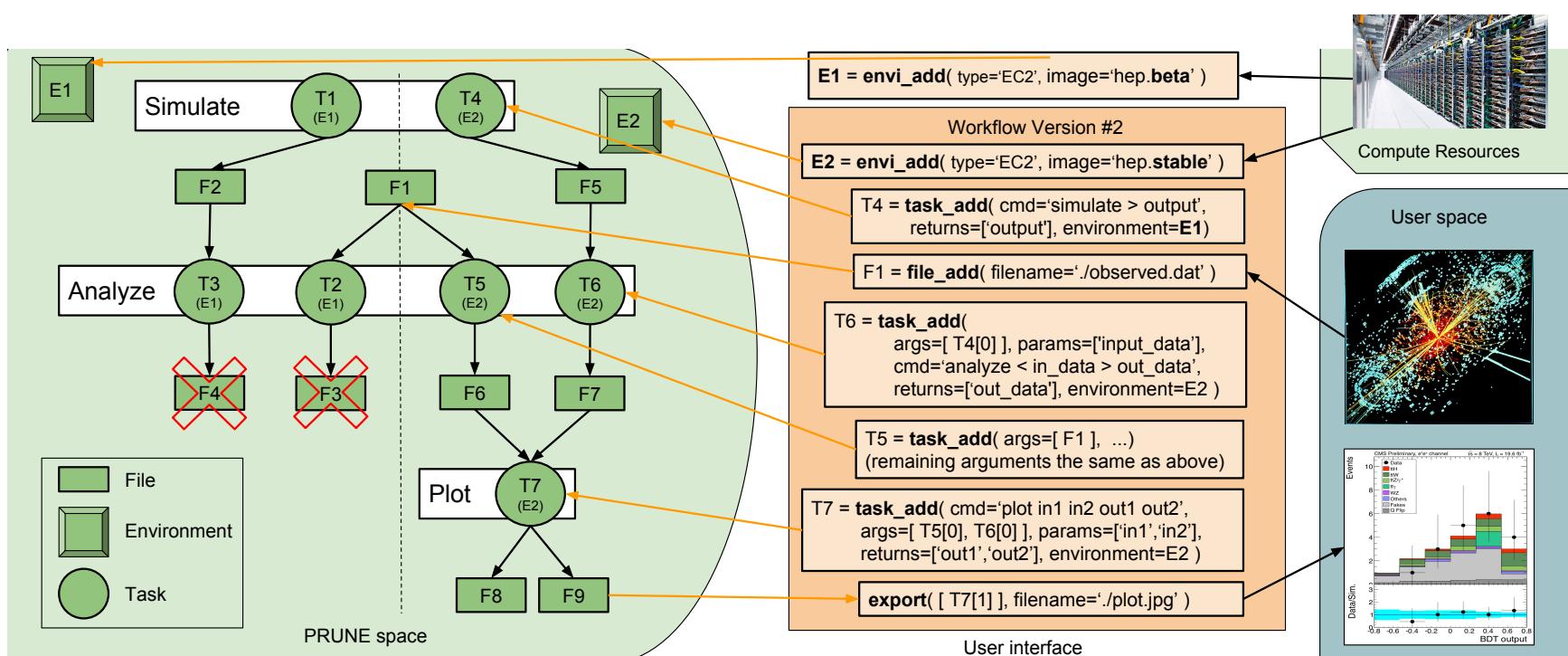
- ~~Git: User decides when to preserve~~
- Preserve ALL specification changes
- ~~Git: Code Commits separate from Code Execution~~
- System Manages ALL computation
- Remove unneeded code later on



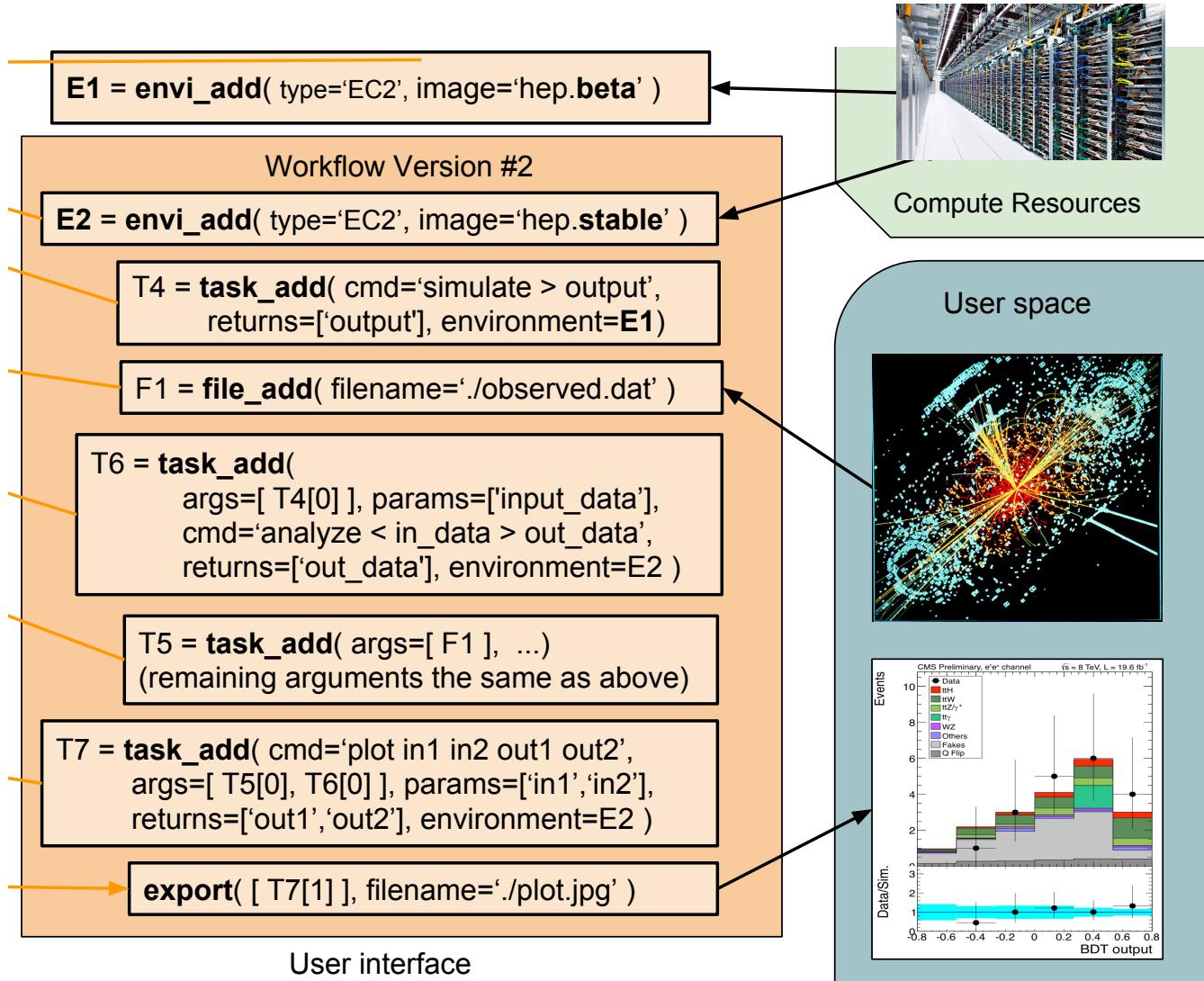
# What to Preserve



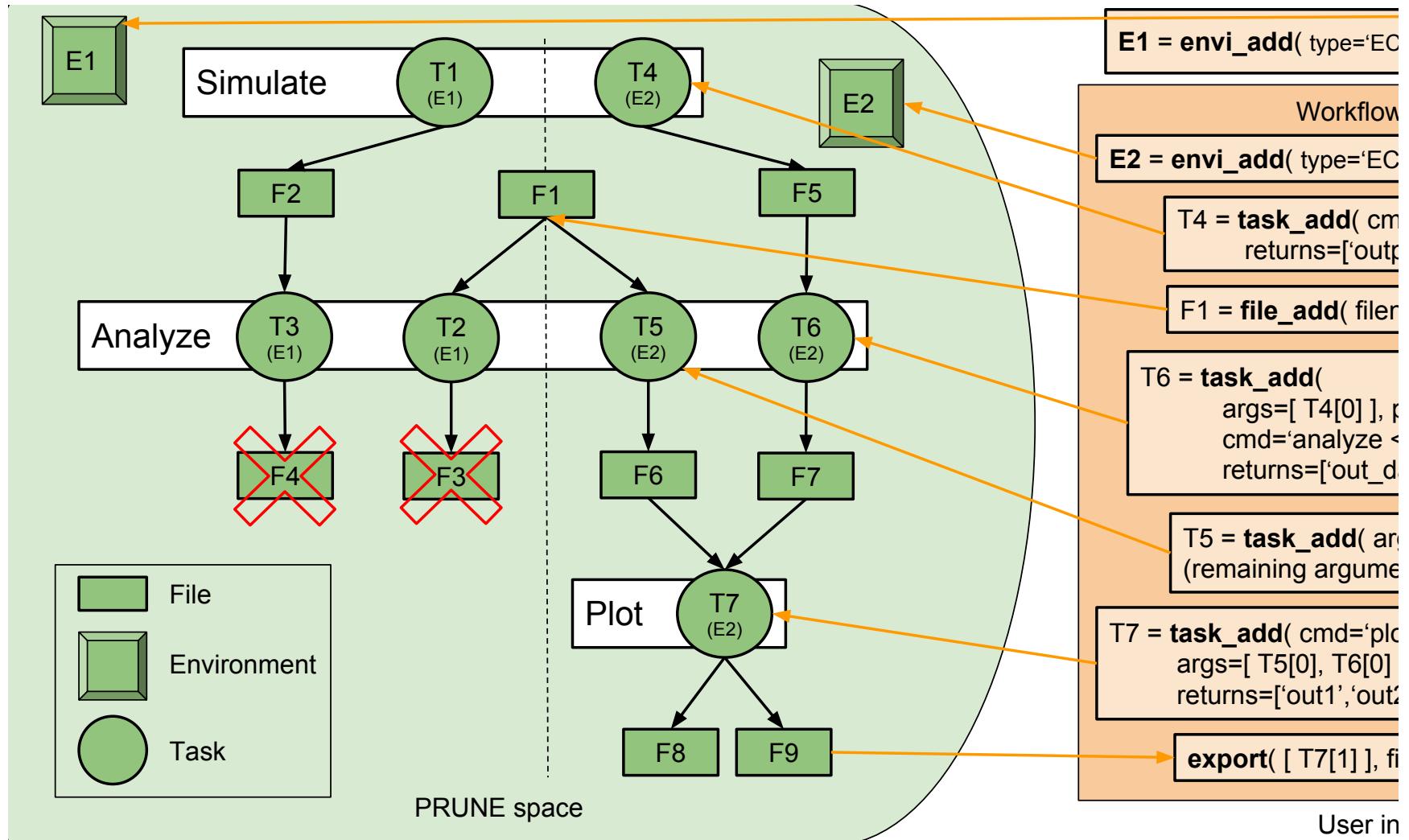
# Overview



# User Interface



# Overview



# Sample code: Merge sort

---

```
#!/usr/bin/env python
from prune import client
prune = client.Connect() #Use SQLite3

##### Import sources stage #####
E1 = prune.env_add(type='EC2',
image='ami-b06a98d8')
D1, D2 = prune.file_add( `nouns.txt',
`verbs.txt' )
```

# Sample code: Merge sort

---

```
##### Sort stage #####
```

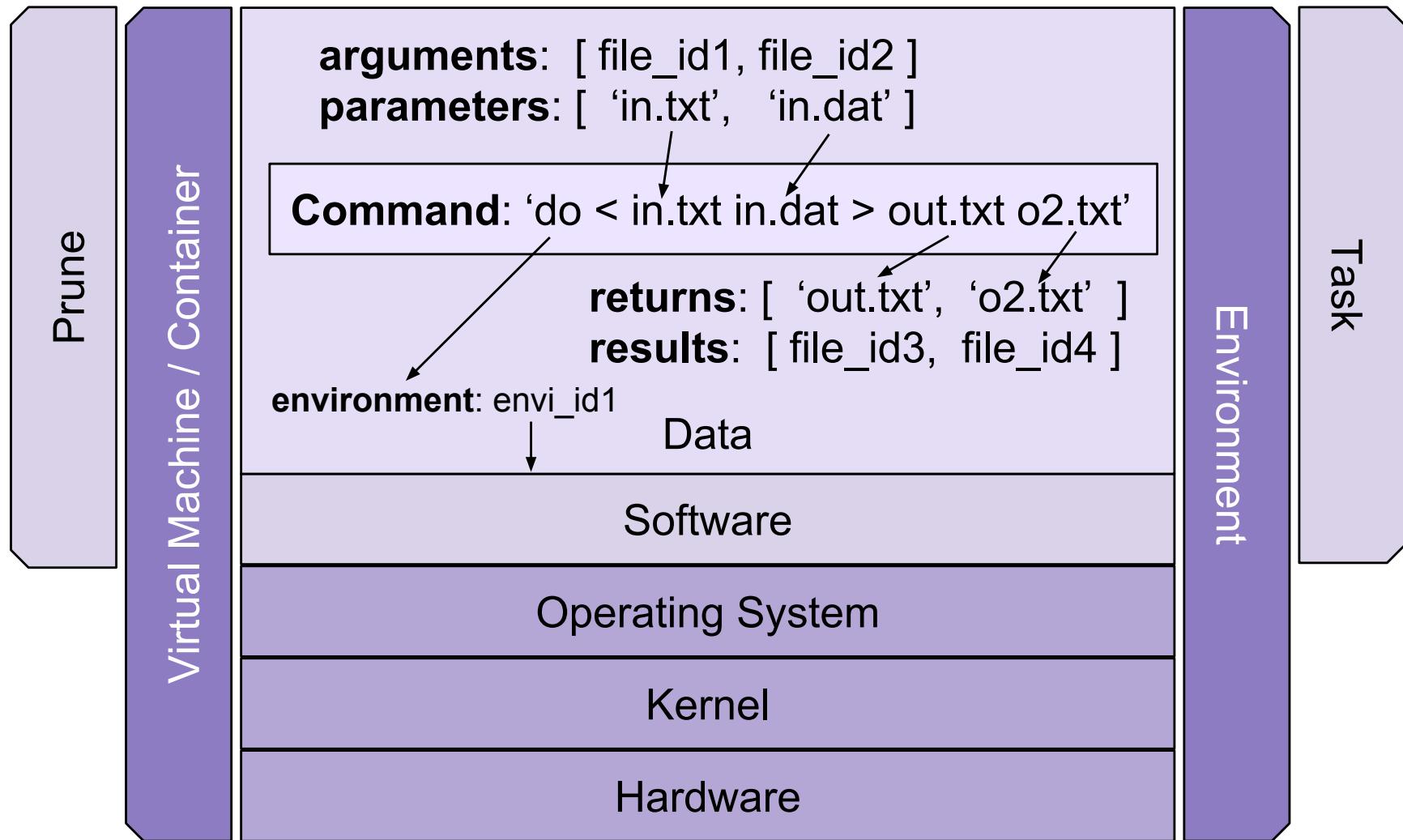
```
D3, = prune.task_add( returns=['output.txt'],  
    env=E1, cmd='sort input.txt > output.txt',  
    args=[D1], params=['input.txt'] )
```

```
D4, = prune.task_add( returns=['output.txt'],  
    env=E1, cmd='sort input.txt > output.txt',  
    args=[D2], params=['input.txt'] )
```

```
##### Merge stage #####
```

```
D5, = prune.task_add(  
    returns=['merged_out.txt'], env=E1,  
    cmd='sort -m input*.txt > merged_out.txt',  
    args=[D3,D4], params=['input1.txt', 'input2.txt'] )
```

---



# Sample code: Merge sort

---

```
##### Sort stage #####
D3, = prune.task_add( returns=['output.txt'],
    env=E1, cmd='sort input.txt > output.txt',
    args=[D1], params=['input.txt'] )
D4, = prune.task_add( returns=['output.txt'],
    env=E1, cmd='sort input.txt > output.txt',
    args=[D2], params=['input.txt'] )

##### Merge stage #####
D5, = prune.task_add(
    returns=['merged_out.txt'], env=E1,
    cmd='sort -m input*.txt > merged_out.txt',
    args=[D3,D4], params=['input1.txt', 'input2.txt'] )
```

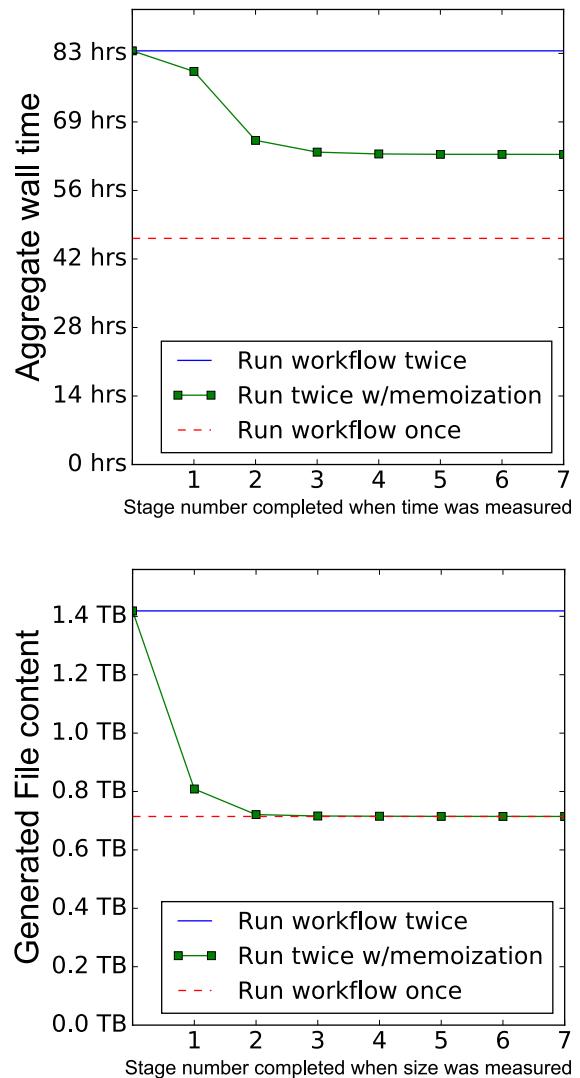
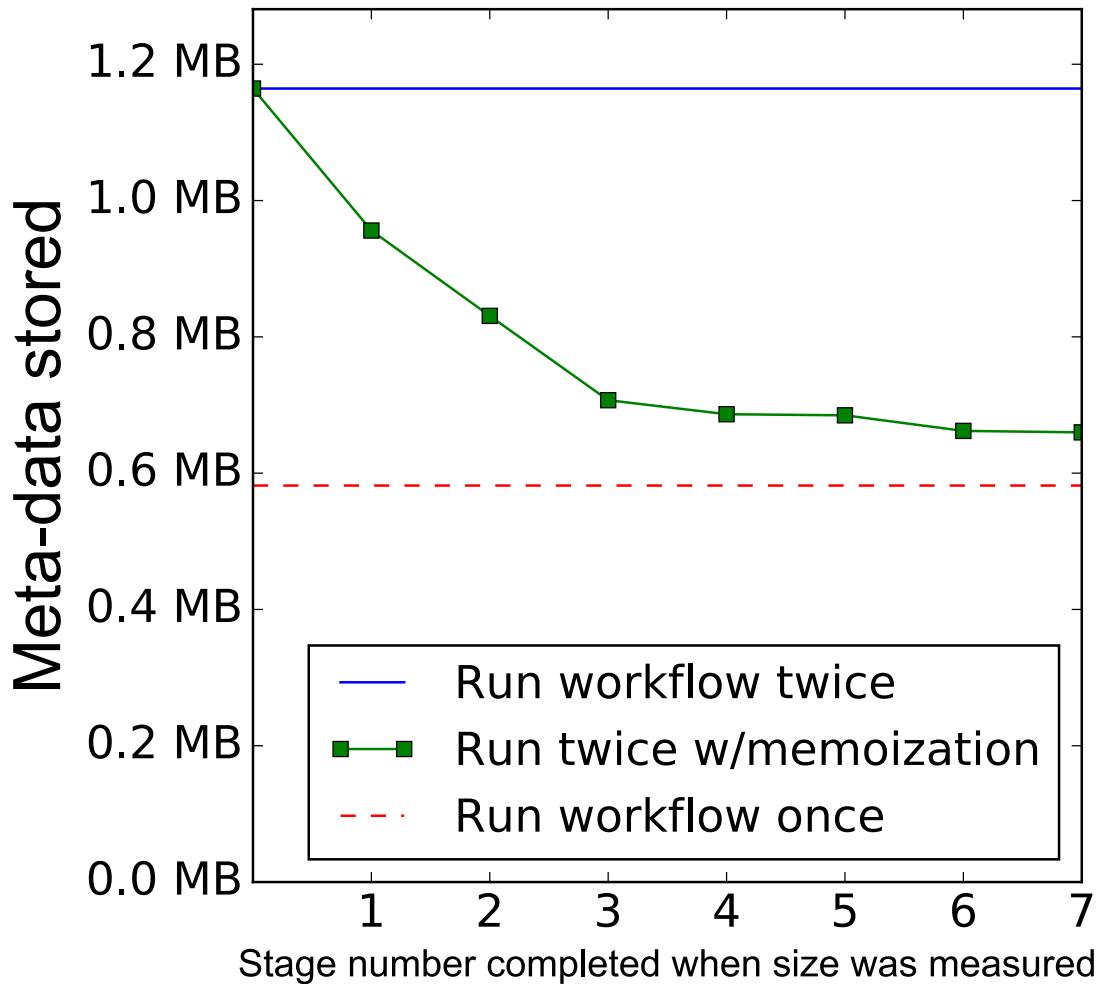
# Sample code: Merge sort

---

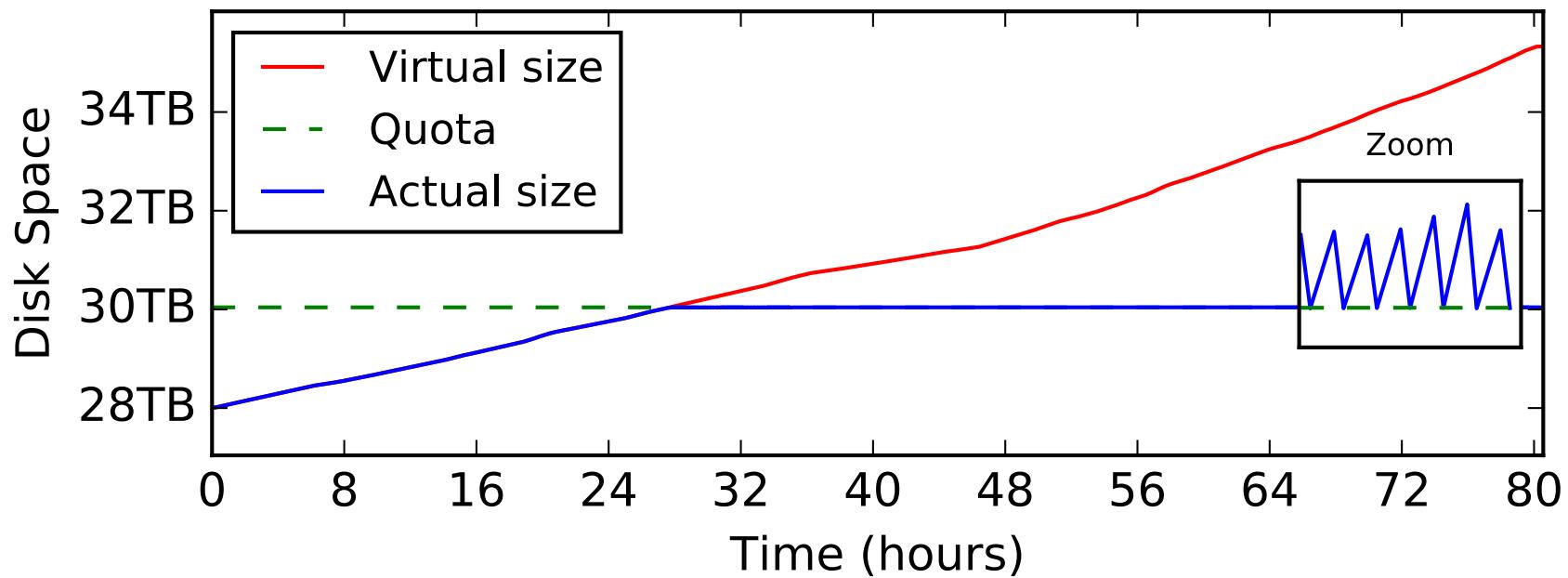
```
##### Execute the workflow #####
prune.execute( worker_type='local', cores=8 )

#####
# Export #####
prune.export( D5, `merged.txt` ) # Final data
prune.export( D5, `wf.prune`, lineage=2 )
```

# Derivation History = Cachable Results

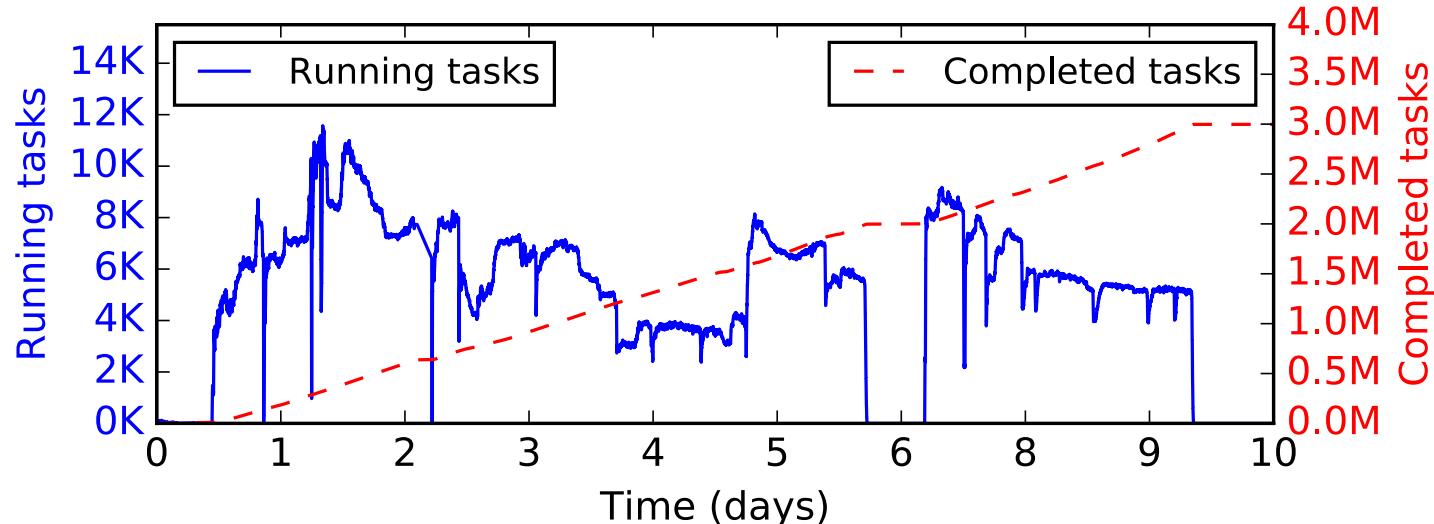


# Quotas



# Scalability

- ~12,000 parallel cores
- ~3 million tasks
- Wall clock overhead
  - ~1% above native





# Thank You!

- Sample workflows
- <http://ccl.cse.nd.edu/software/prune/prune.html>
  - Merge sort
  - Pairwise comparisons (US Censuses)
  - High-energy Physics

