

Atlanta, Georgia, April 19, 2010
in conjunction with IPDPS 2010



UNIBUS: ASPECTS OF HETEROGENEITY AND FAULT TOLERANCE IN CLOUD COMPUTING

Magdalena Slawinska
Jaroslaw Slawinski
Vaidy Sunderam

{magg, jaross, vss}@mathcs.emory.edu



EMORY
UNIVERSITY

Emory University, Dept. of Mathematics and Computer Science
Atlanta, GA, USA

Creating a problem

2

1. What do I want?

- Execute an MPI application

2. What do I need?

- Target resource: MPI cluster
- FT services: Checkpoint, Heartbeat

3. What do I have?

- Access to the Rackspace cloud



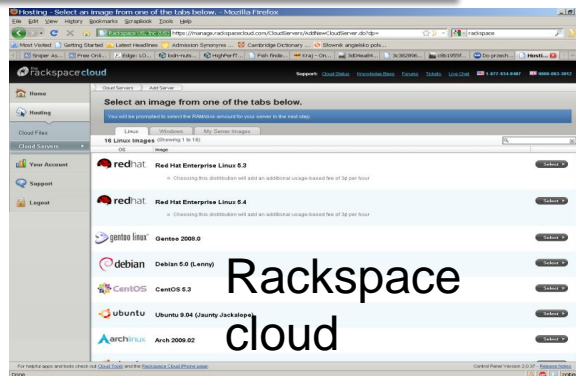
User

4. Why might I want FT on cloud?

- To reduce costs (money, time, energy, ...)
- Reliability
- ...

5. What is the overhead introduced by FT?

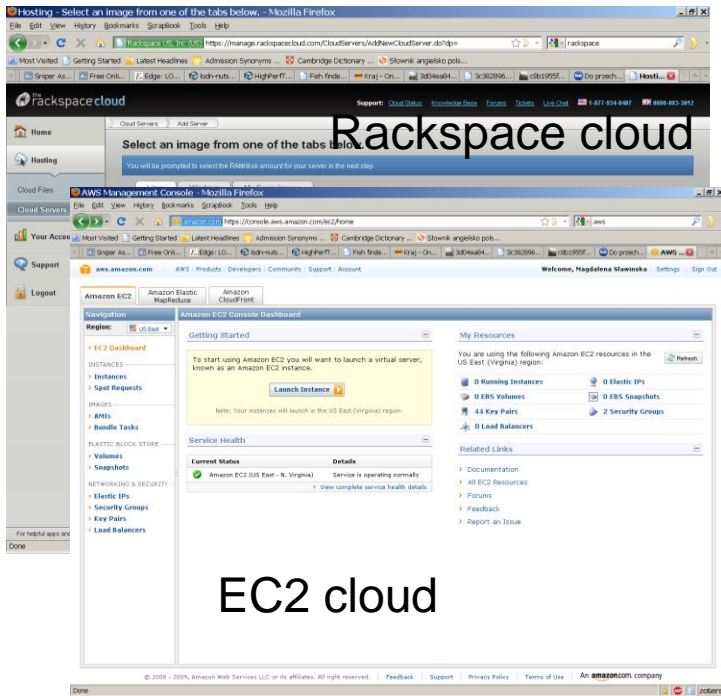
6. Can I do that? How?



Problem

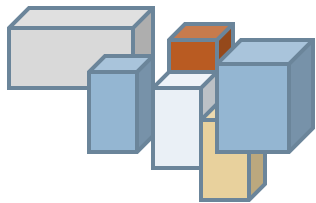
3

Available resource



Rackspace cloud

EC2 cloud



Target resource

Workstations



User

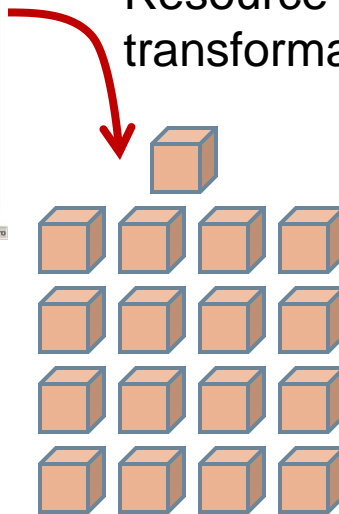
User's requirements

- Execute MPI software
- Target resource: MPI cluster
- Target platform: FT-flavor

User's resources

- Rackspace cloud (credentials)

Resource transformation



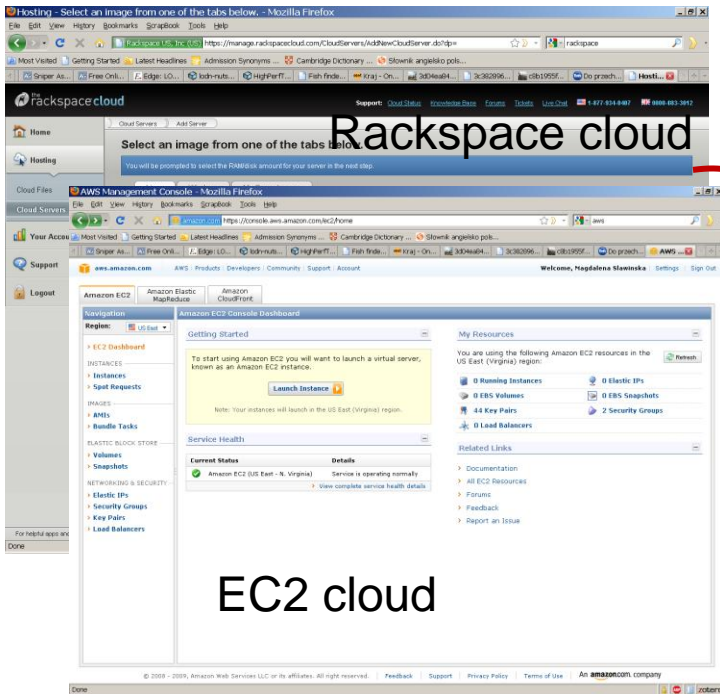
Manually:

- interaction with web page
- prepare the image: install required software and dependencies
- instantiate servers
- configure passwordless authentication
-
- 1 man-hour for 16+1 nodes

Unibus: a resource orchestrator

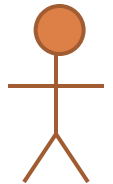
4

Available resource



Rackspace cloud

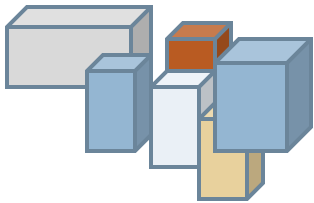
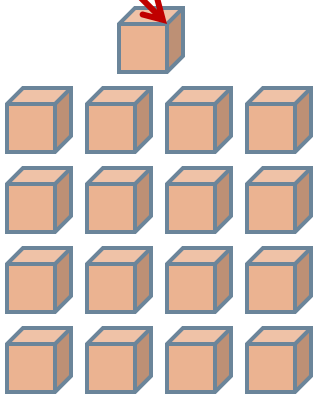
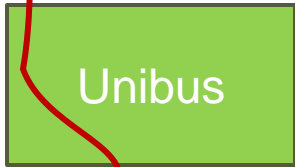
EC2 cloud



User

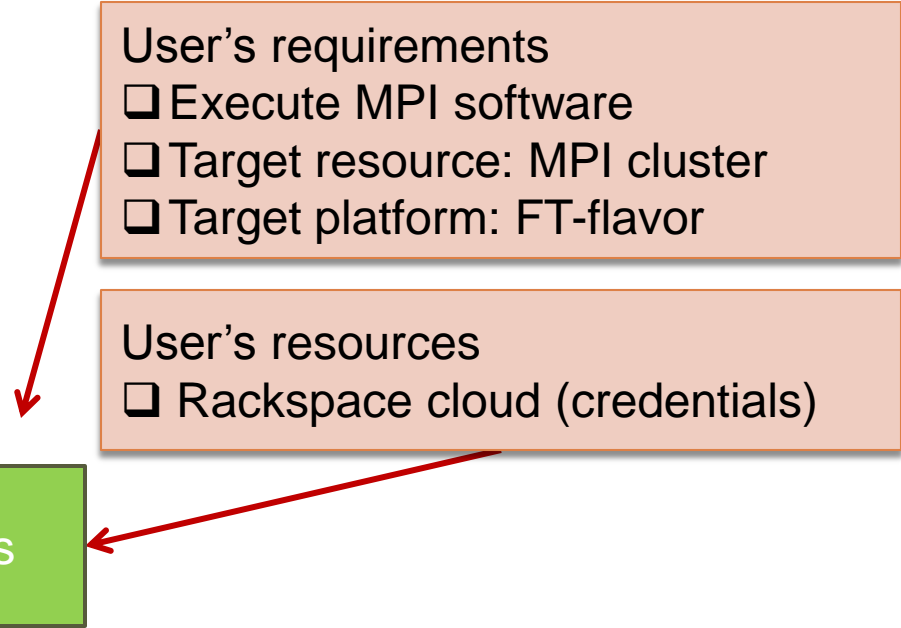
- User's requirements
- Execute MPI software
 - Target resource: MPI cluster
 - Target platform: FT-flavor

- User's resources
- Rackspace cloud (credentials)



Target resource

Workstations



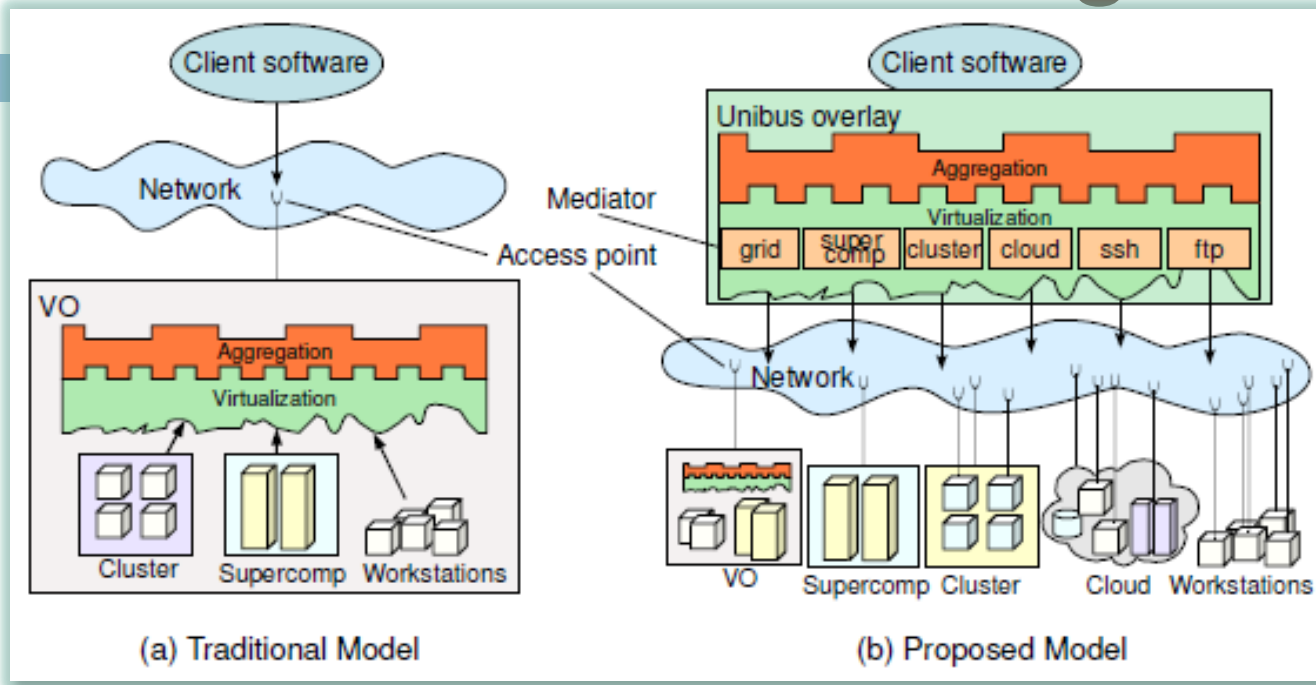
Outline

5

- Unibus – an infrastructure framework that allows to orchestrate resources
 - ▣ Resource access virtualization
 - ▣ Resource provisioning
- Unibus – FT MPI platform on demand
 - ▣ Automatic assembly of an FT MPI-enabled platform
 - ▣ Execution of an MPI application on the Unibus-created FT MPI-enabled platform
 - ▣ Discussion of the FT overhead

Unibus resource sharing model

6

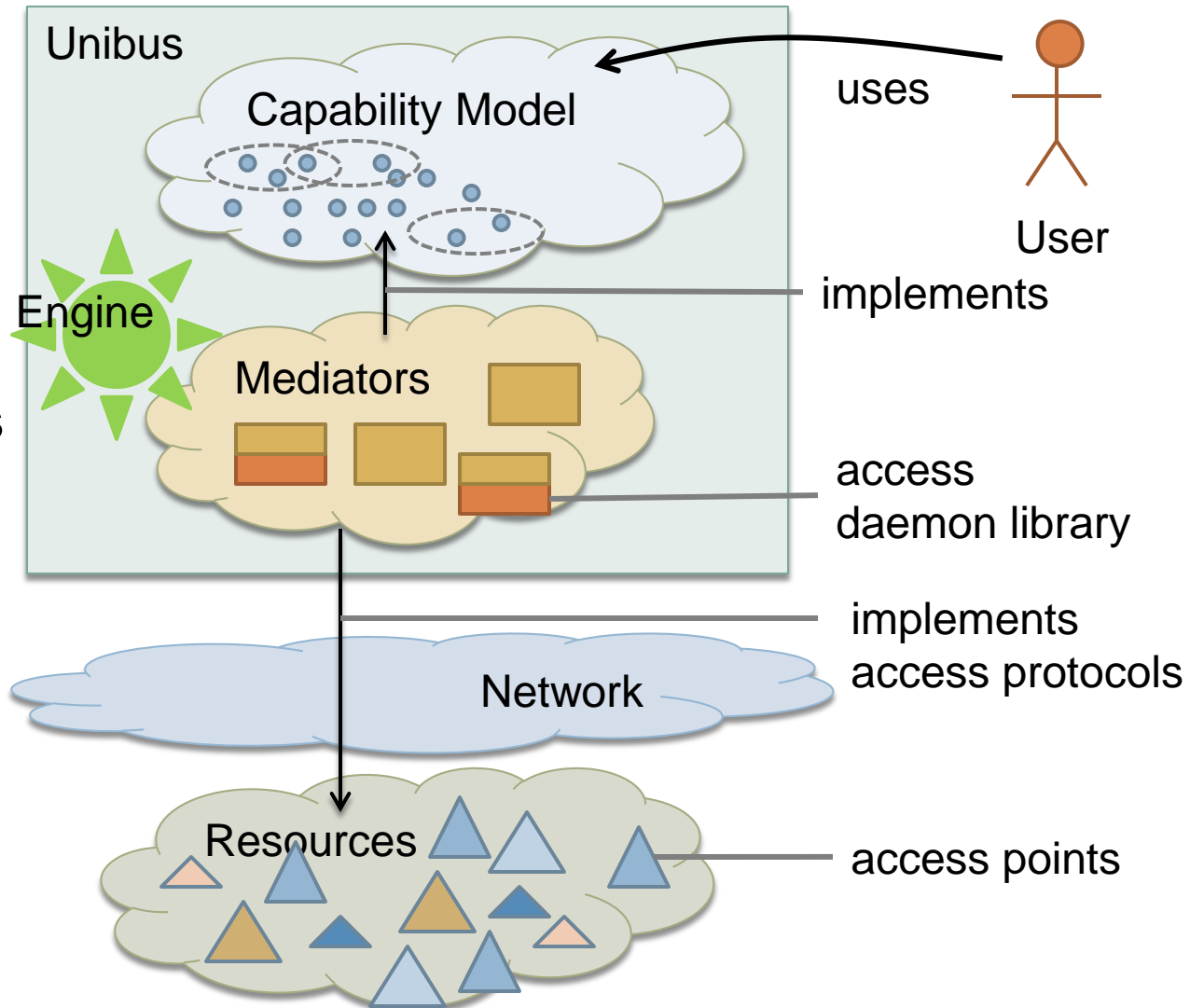


	Traditional Model	Proposed Model
Resource exposition	Virtual Organization (VO)	Resource provider
Resource usage	Determined by VO	Determined by a particular resource provider
Resource virtualization and aggregation	Resource providers belonging to VO	Software at the client side

Handling heterogeneity in Unibus

7

- Resources exposed in an arbitrary manner as *access points*
- Capability Model to abstract operations available on provider's resources
- Mediators to implement the specifics of access points
- Knowledge engine to infer relevant facts

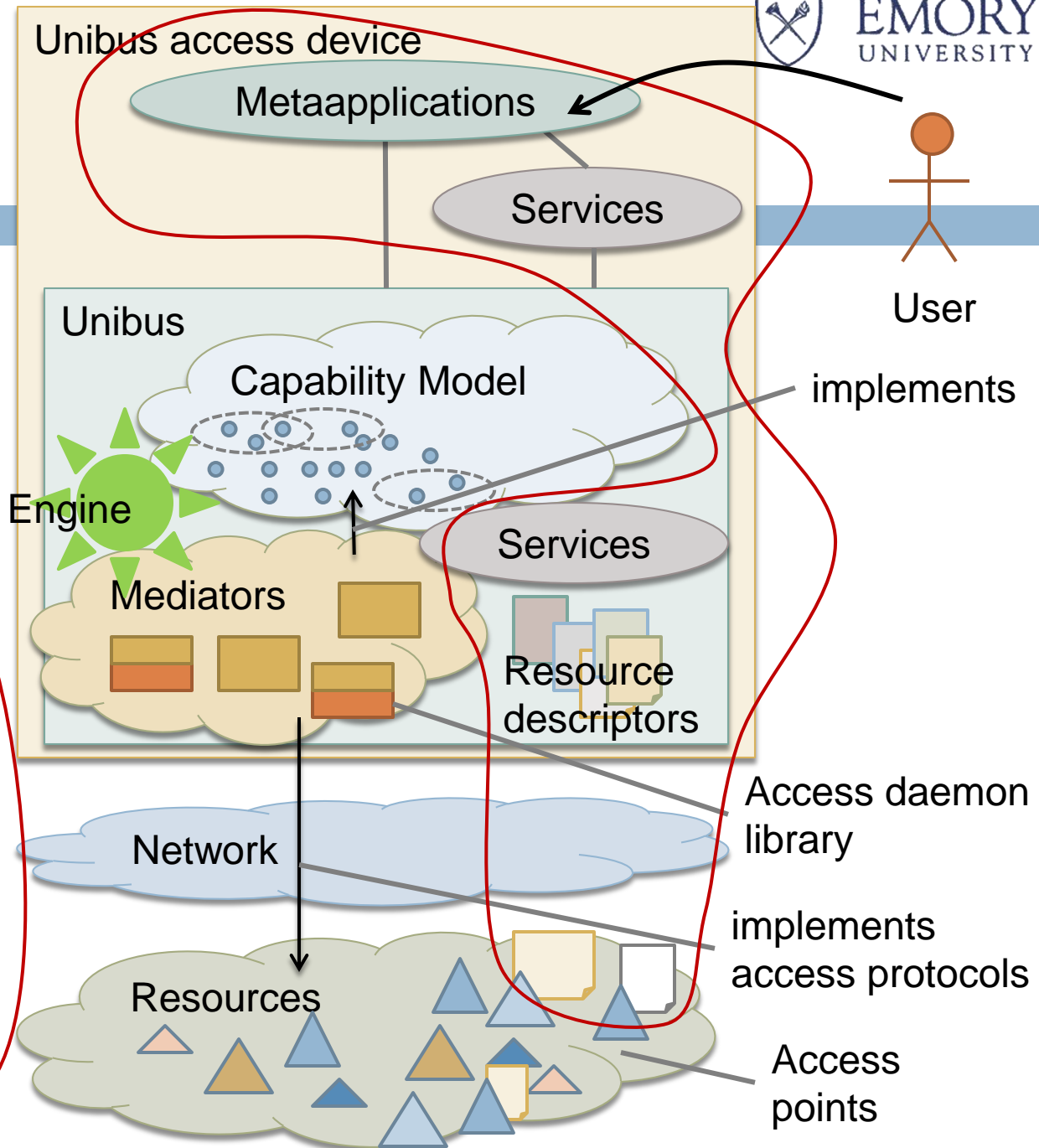




Complicating a big picture ...

8

- Resources exposed in an arbitrary manner as *access points*
- Capability Model to abstract operations on resources
- Mediators to implement the specifics of access points
- Knowledge engine to infer relevant fact
- Resource descriptors to describe resources *semantically* (OWL-DL)
- Services (standard and third parties), e.g., heartbeat, checkpoint, resource discovery, etc.
- Metaapplications to orchestrate execution of applications on relevant resources

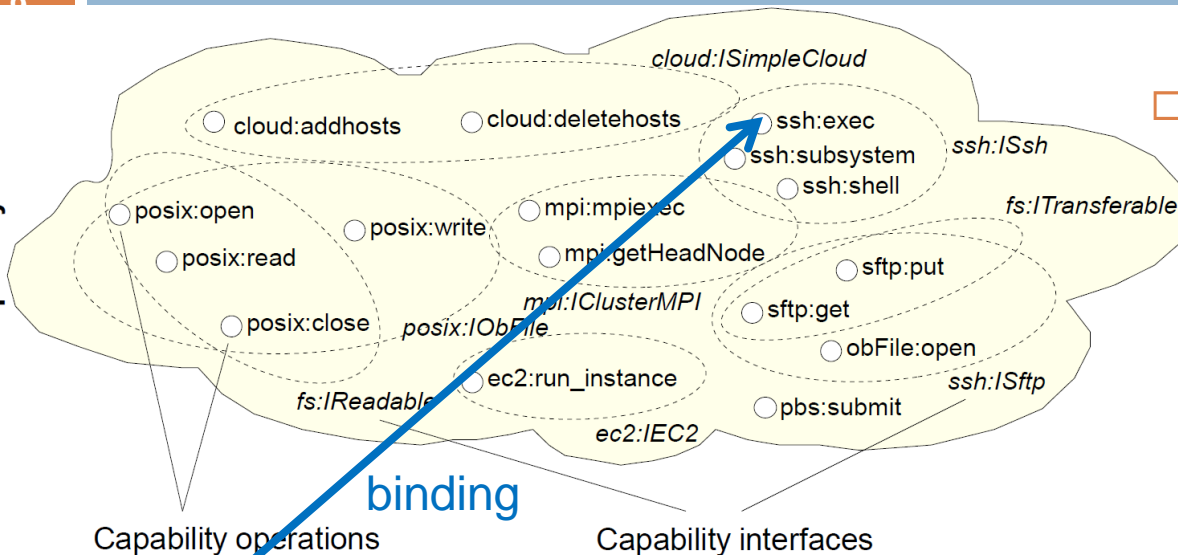




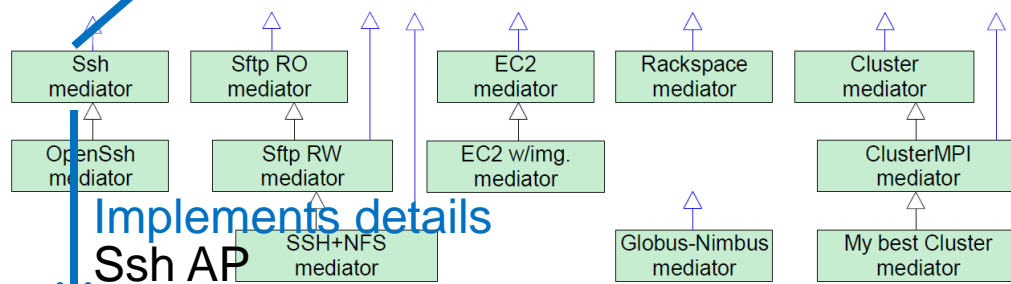
Virtualizing access to resources

Capability Model and mediators

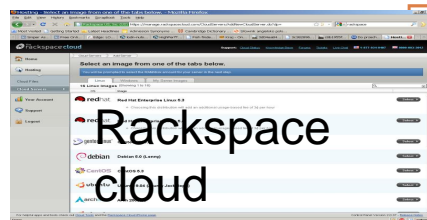
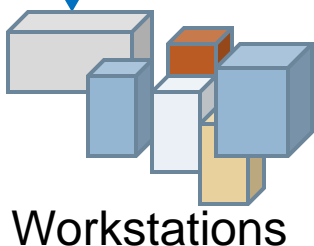
Capability model



Mediator capability implementations



- **Capability Model**
 - ▣ Provides virtually homogenized access to heterogeneous resources
 - ▣ Specifies *abstract operations*, grouped in *interfaces*
 - ▣ Interface hierarchy not appropriate (e.g. fs:ITransferable and ssh:ISftp)

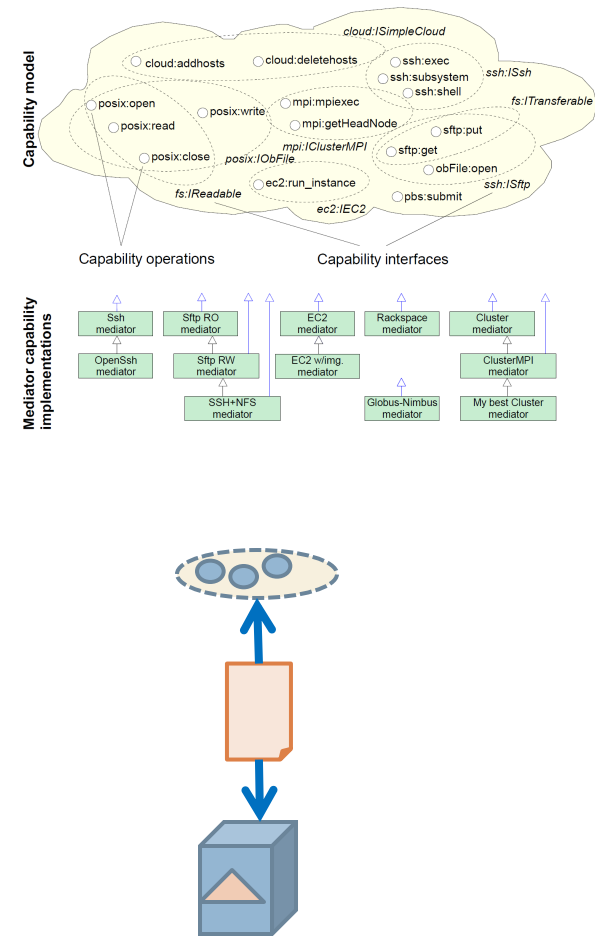
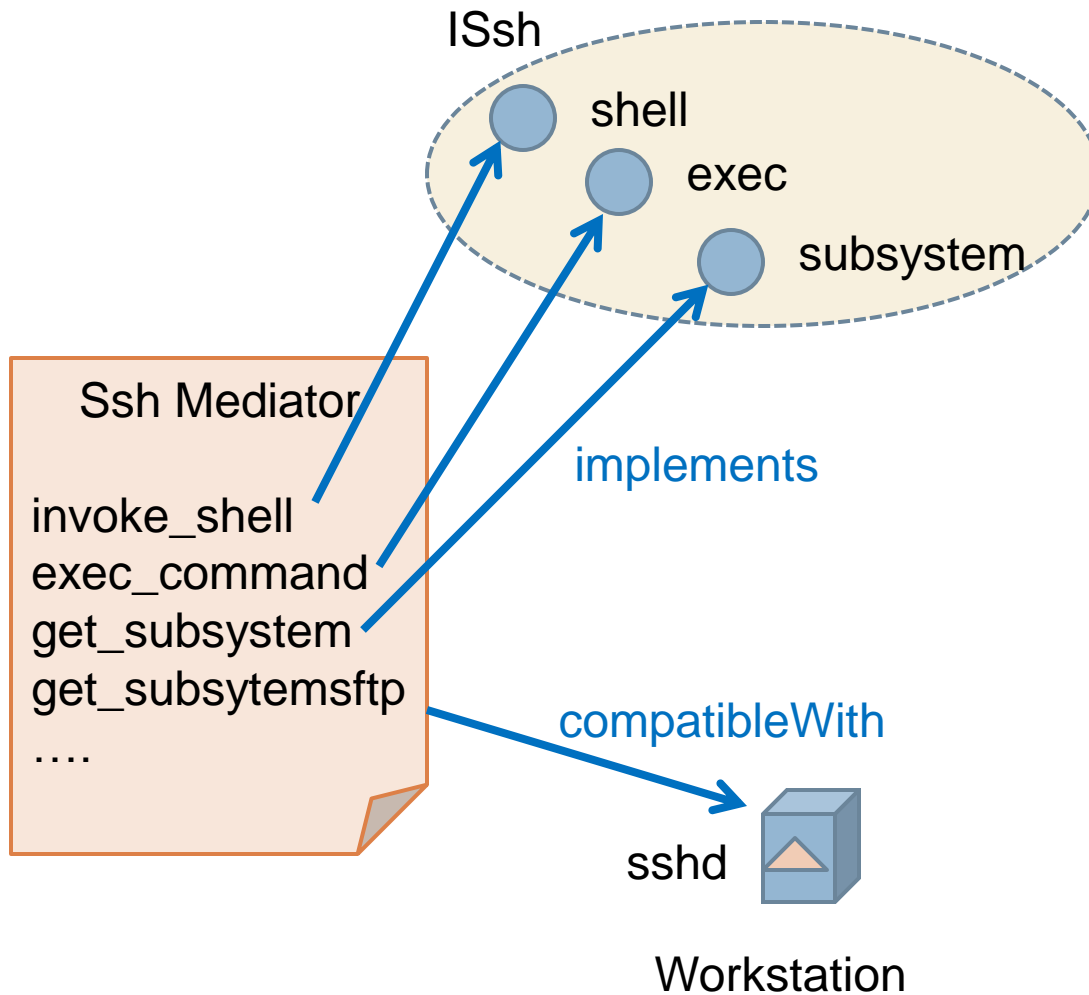


Mediators

- ▣ Implement resource access point protocols

Virtualizing access to resources

10



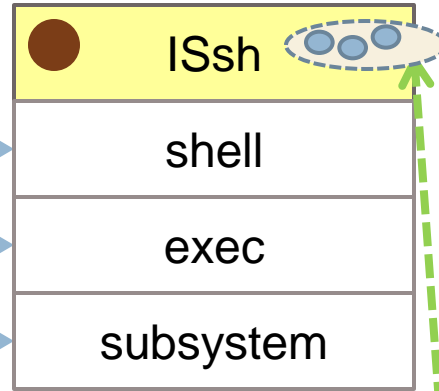
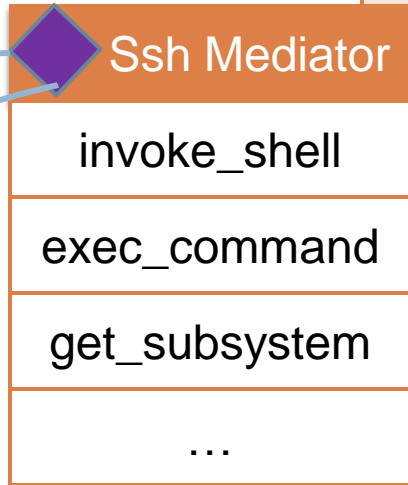
Knowledge engine

11



Mediator's Developer

Knowledge Set

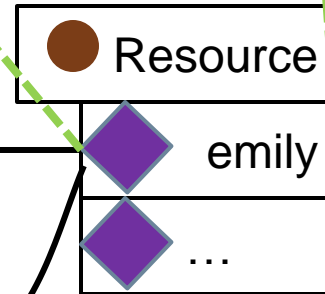
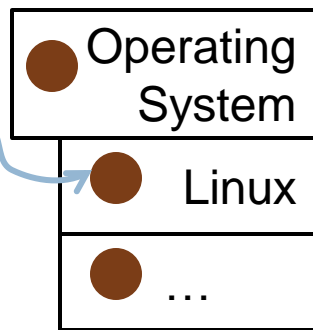
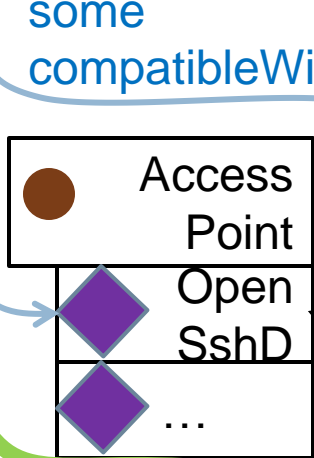


implements
implements
implements

compatibleWith

hasOperation

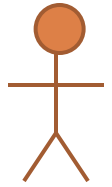
some compatibleWith



some hasOS

hasAccessPoint

Request interface ISsh

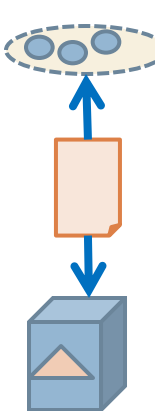


User

Resource: emily



Knowledge Engine (inferring)



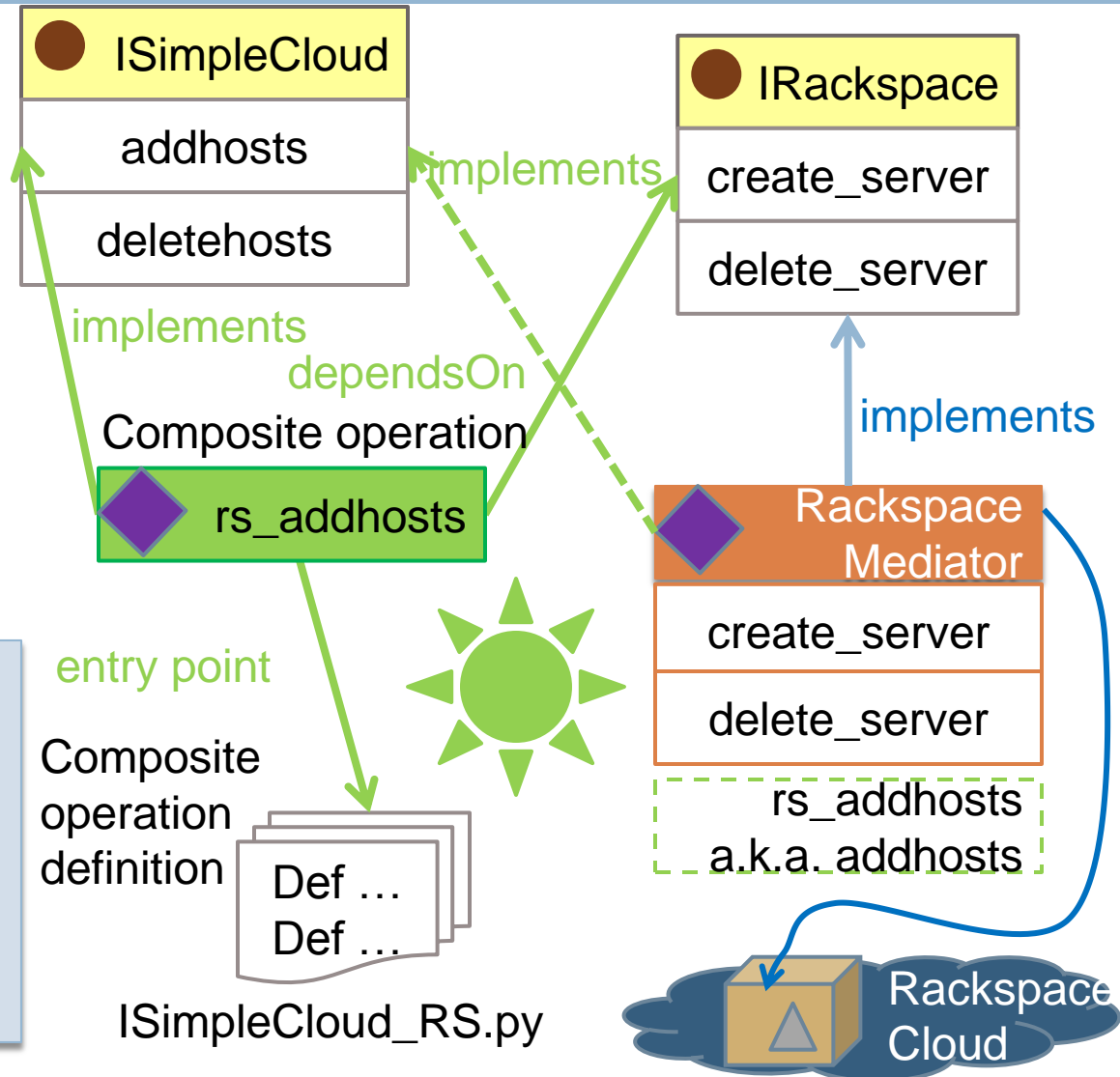
Composite operations

12

- Rs_addhosts dependsOn create_server
- Create_server is implemented by RS Mediator
- Rs_addhosts implements addhosts
- So RS mediator implements addhosts

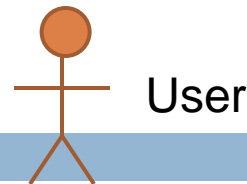
Composite operations

- Dynamically expand mediator's operations
- May result in classification of mediators and compatible resources to new interfaces



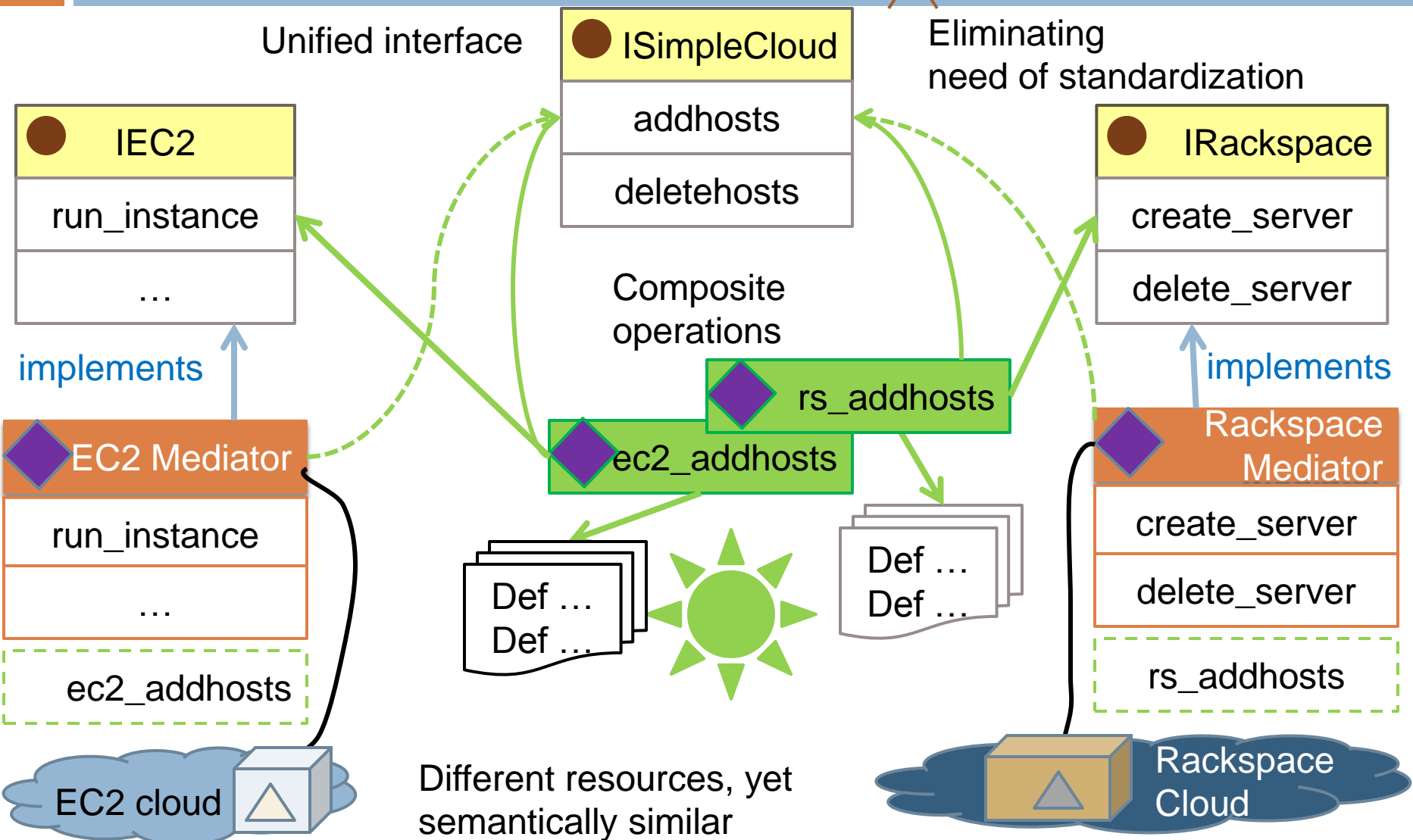


Resource access unification via composite operations



Unified interface

Eliminating need of standardization



Resource provisioning

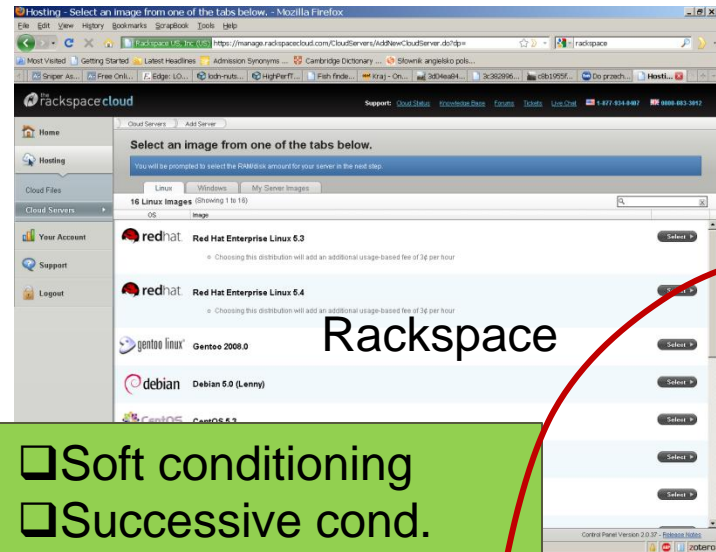
Homogenizing resource heterogeneity

14

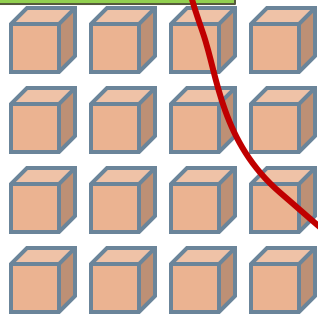
- Conditioning increases resource specialization levels
 - **Soft conditioning**
 - changes resource **software capabilities**
 - e.g., installing MPI enables execution of MPI apps
 - **Successive conditioning**
 - enhances resource capabilities in terms of available **access points** (may use soft conditioning)
 - e.g., deploying Globus Toolkit makes the resource accessible via Grid protocols

Transforming Rackspace to FT-enabled MPI platform

15



- Soft conditioning
- Successive cond.
- Composite ops
- ...



FT MPI cluster

Unibus

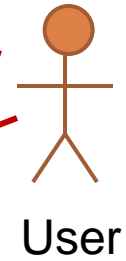
User's credentials

Rackspace descriptor

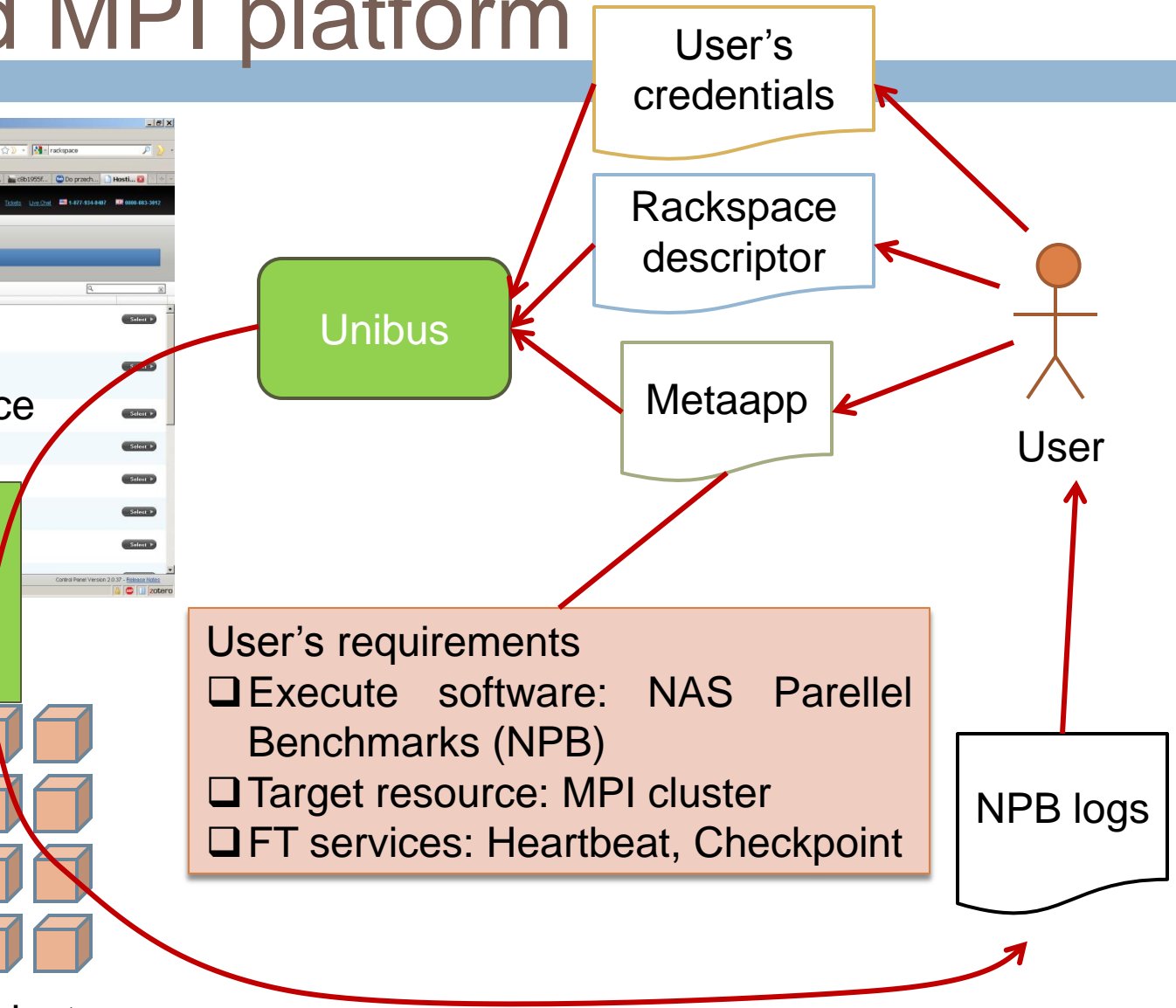
Metaapp

- User's requirements
- Execute software: NAS Parellel Benchmarks (NPB)
 - Target resource: MPI cluster
 - FT services: Heartbeat, Checkpoint

NPB logs



User



Rackspace Cloud to MPI cluster

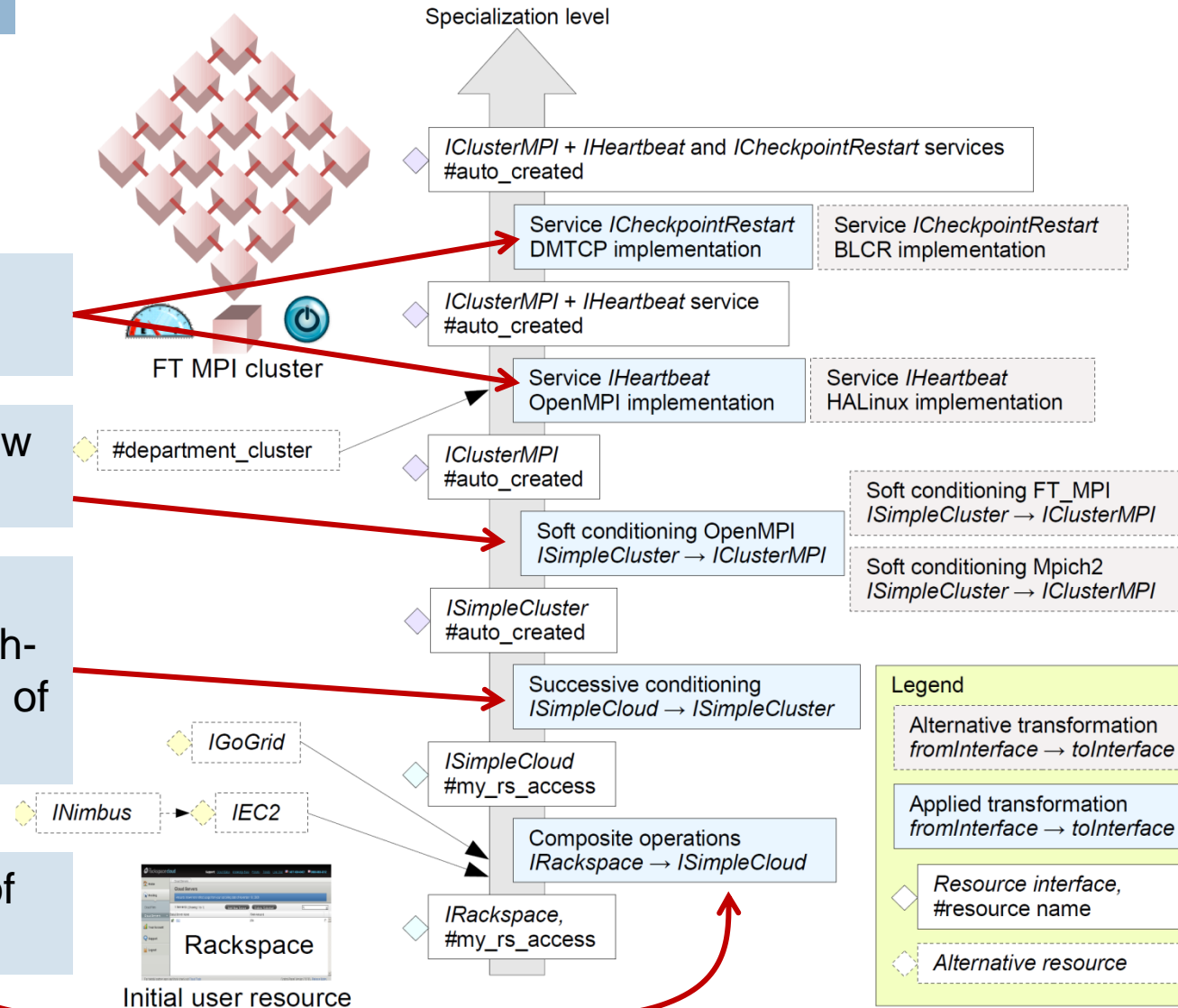
16

Installing other services (FT)

Deployment of MPI on new resources

Creating a new group of resources (Rackspace ssh-enabled servers) in terms of new access points

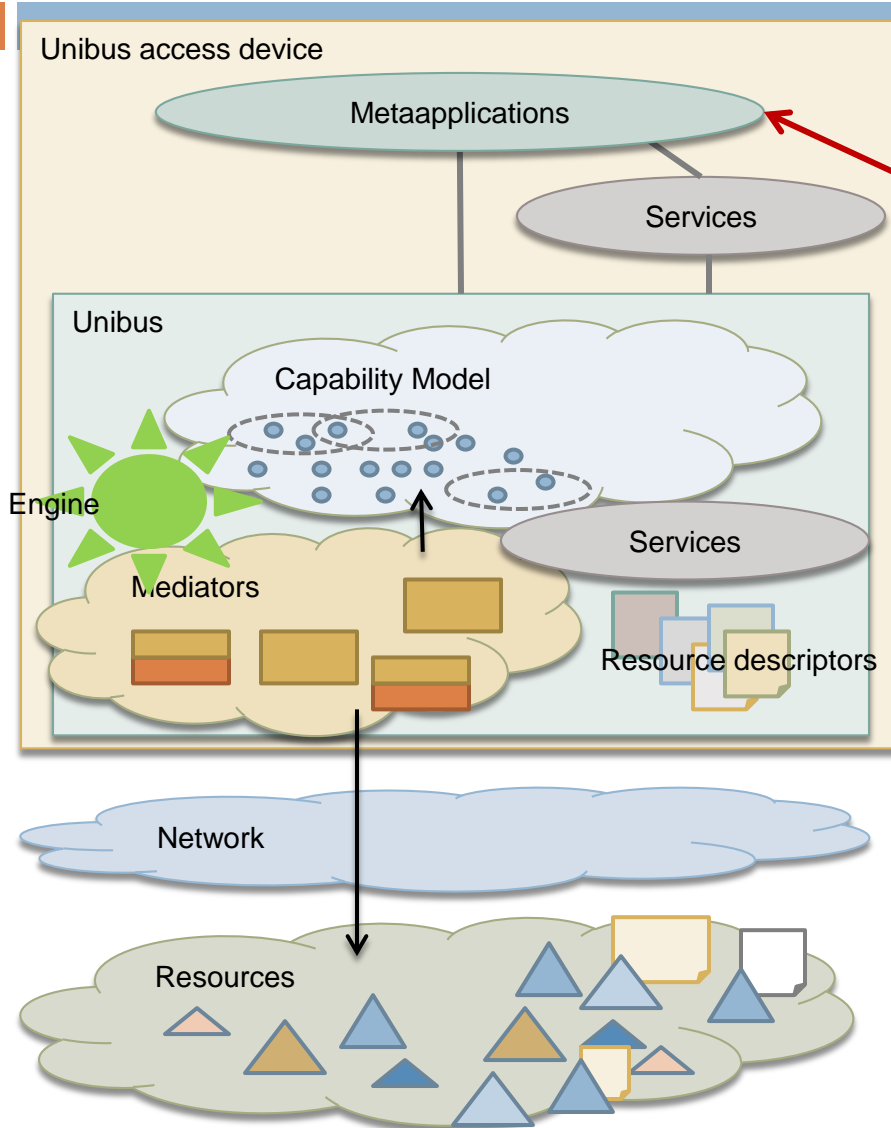
Obtaining a higher level of abstraction



Metaapplications

17

- User's requirements
- Execute software: NAS Parellel Benchmarks (NPB)
 - Target resource: MPI cluster
 - FT services: Heartbeat, Checkpoint



Metaapplication

18

Metaapplication

- Requests
 - IClusterMPI
 - FT services:
 - IHeartbeat
 - ICheckpointRestart
- Specifies available resources
- Performs benchmarks
- Transfers benchmarks execution logs to the head node
 - Requests ISftp

```

1  def create_cluster(resource):
2      proxy = resource.createProxy('IClusterMPI', services=('IHeartbeat',
3          'ICheckpointRestart'), cpus/4) #4 core processors
4      #install NPB3.3 and dependencies (gcc, gfortran, OpenMPI)
5      #checkpoint every 1 minutes, heartbeat every 10 seconds
6      proxy.services['IHeartbeat'].init(interval = 10, death_after = 2)
7      proxy.services['ICheckpointRestart'].init(interval = 1*60)
8
9      import_resources('file:resources.owl')
10     #the resource needs to be MPI cluster and have some features
11     res = get_resource('a cluster:IClusterMPI;\
12         uc:compatibleWithService uc:ICheckpointRestart, uc:IHeartbeat;\
13         cluster:cpuCores ?n . filter(?n >= %d)' % cpus)
14     proxy = create_cluster(res)
15
16     tests = BenchmarkTestsIterator()
17     for test in tests:
18         test_file = '..'.join(test, class_, cpus)
19         cmd = '..'.join(path, test_file, '>', log_path)
20         try:
21             proxy['#mpiexec'](cmd, cpus)
22         except proxy.services['IHeartbeat']['#Exception']():
23             # is any chkp files available (to restart)?
24             if proxy.services['ICheckpointRestart']['#get']():
25                 while True:
26                     new_proxy = create_cluster(res)
27                     ckpt = proxy.services['ICheckpointRestart']['#get']()
28                     new_proxy.services['ICheckpointRestart']['#set'](ckpt)
29                     #remove old, failed cluster
30                     proxy = new_proxy
31                     try:
32                         proxy.services['ICheckpointRestart']['#restart']()
33                     except proxy.services['IHeartbeat']['#Exception']():
34                         continue
35                     else:
36                         break
37                 #no ckpt available yet — just restart the same test
38             else:
39                 tests.repeat()
40
41     # test done — save the result log locally
42     with open('..'.join((test_file, time.strftime('%y%m%d%H%M%S'))), 'w') as f:
43         #get the ISftp interface to the head node
44         head = proxy['#get_head_node']().createProxy('ISftp')
45         f.write(head['#open'](log_path).read())

```

Rackspace testbed

19

- ❑ 16 working nodes (WN) + 1 head node (HN)
- ❑ Node: 4-core, 64-bit, AMD 2GHz
- ❑ Debian 5.0 (Lenny)
- ❑ OpenMPI v. 1.3.4 (GNU suite v. 4.3.2 (gcc, gfortran))
- ❑ NAS Parallel Benchmarks v.3.3, class B

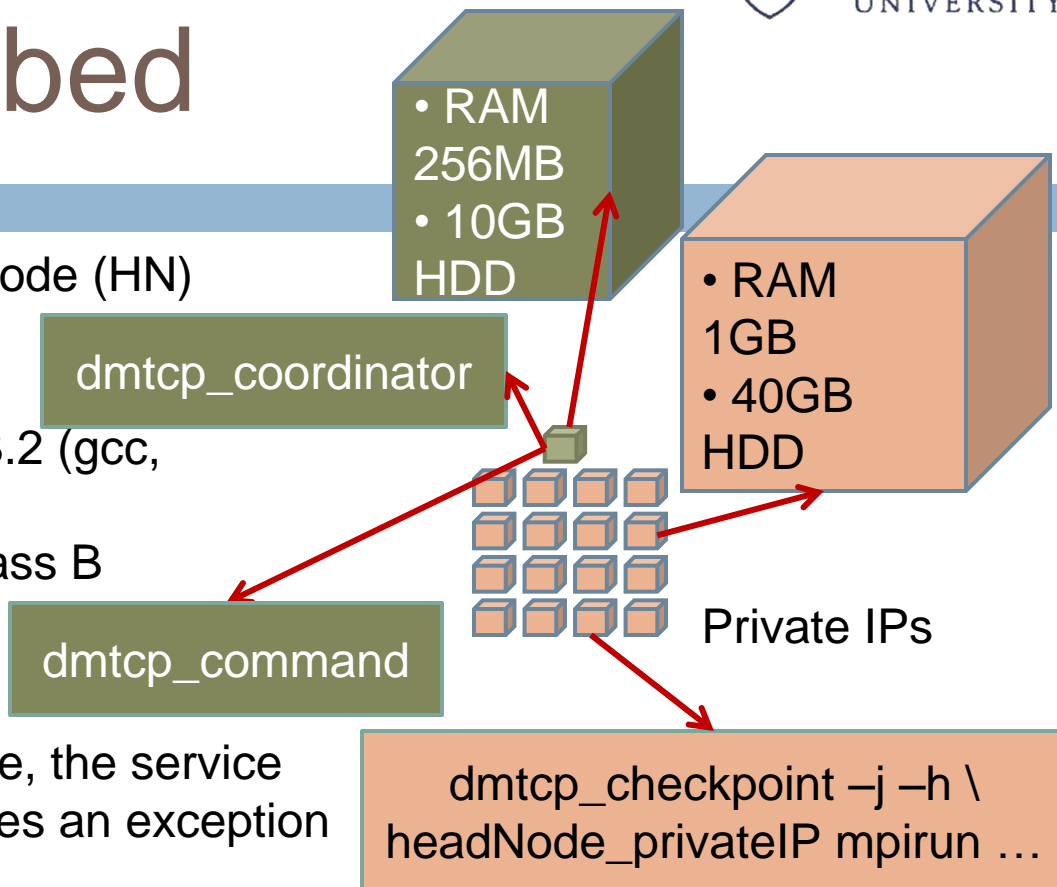
FT setup:

Heartbeat service:

- ❑ OpenMPI-based – in case of failure, the service determines failed node(s) and raises an exception

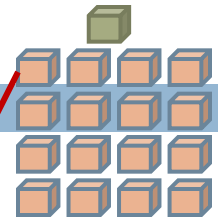
Checkpoint/restart service:

- ❑ DMTCP – Distributed MultiThreaded CheckPointing
 - ❑ user-level transparent checkpointing
- ❑ Executes `dmtcp_command` every 60 secs on HN to checkpoint 81 processes (64 MPI processes, 16+1 OpenMPI supervisor processes)
- ❑ Moves local checkpoint files from WN to HN (in parallel)
- ❑ Checkpoint time – 5 sec; moving checkpoints from WN -> HN less than 10 sec; compressed checkpoint size c.a.1GB



Results: NPB, class B, Rackspace, DMTCP, OpenMPI Heartbeat

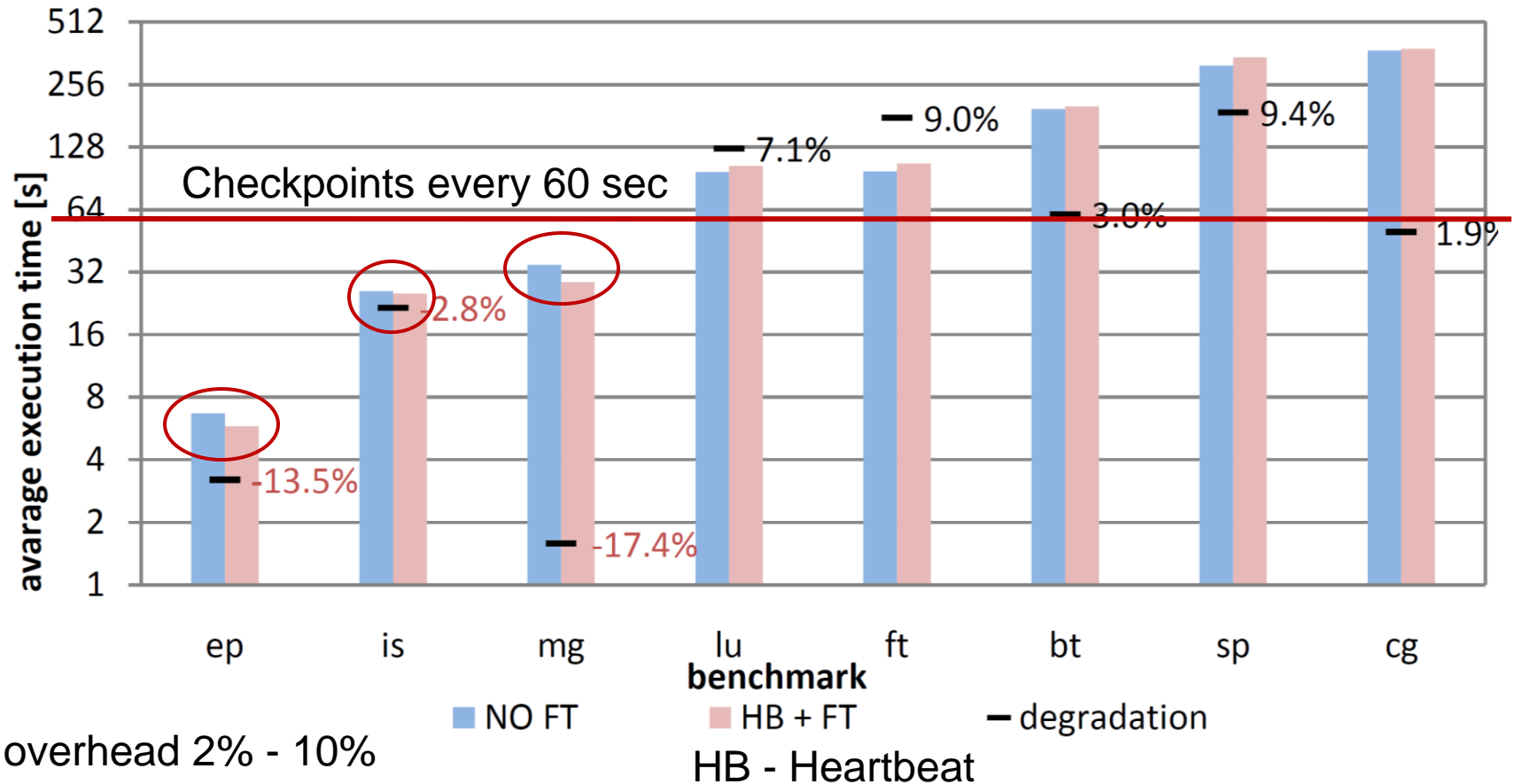
20



16 Worker Nodes (WN) + 1 Head Node

WN: 4-core, 64-bit, AMD Opteron 2GH, 1GB RAM, 40 GB HDD

Checkpoints every 60 sec, average of 8 series



- The Unibus infrastructure framework
 - ▣ Virtualization of access to various resources
 - ▣ Automatic resource provisioning
- Innovatively used to assemble an FT MPI execution platform on cloud resources
 - ▣ Reduces effort to bare minimum (servers instantiation, etc)
 - 15-20 min from 1 man-hour
 - ▣ Observed FT overhead 2%-10% (expected at least 8%)
- Future work
 - ▣ Migration and restart of MPI-based computations on two different clouds or a cloud and a local cluster
 - ▣ Work with an MPI application

