

Rhizoma: A Runtime for Self-deploying, Self-managing Overlays

Qin Yin*, Adrian Schüpbach*, Justin Cappos+, Andrew Baumann*,
Timothy Roscoe*

* Systems Group, Department of Computer Science, ETH Zurich

+ Dept. of Computer Science and Engineering, University of Washington



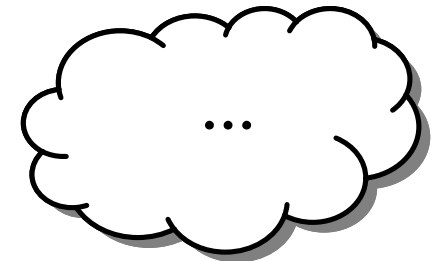
Cloud computing infrastructures are emerging



CPU	1	4	8
Memory	1.7G	7.5G	15G
Disk	160G	850G	1690G
Unit	1	2	2
Price	\$0.125	\$0.50	\$1.00



Storage	\$0.150 per GB
Transfer	\$0.100 per GB
Requests	\$0.01 per 10,000 GET



Deploying an online service



Challenge 1:
Selecting appropriate providers and resources

Service deployed with one provider



Service gains popularity

**Challenge 2:
Adapting to changes in application load**



Provider outage

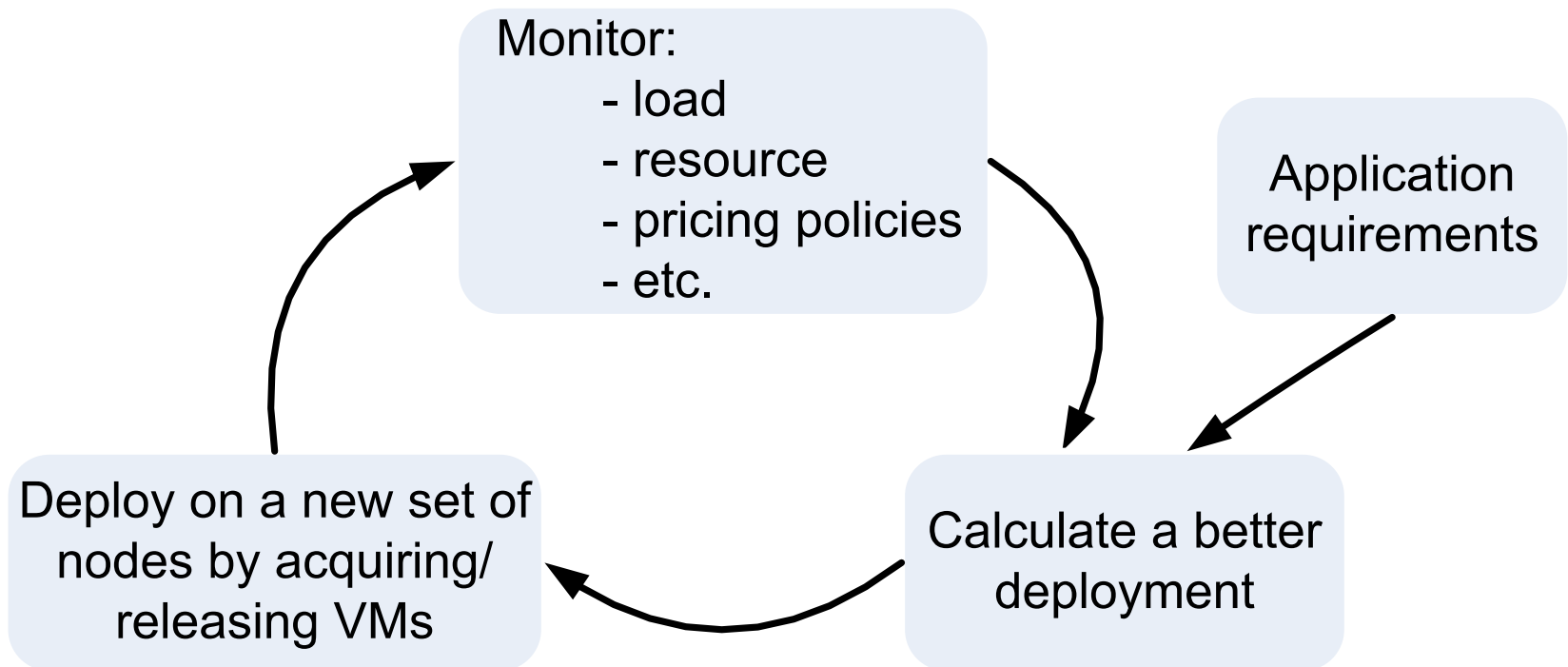


Problem summary

- Challenges
 - Changing resource availability
 - Changing application requirements
 - Changing pricing policies
- Current solutions
 - Human-in-the-loop
 - Separate management console (e.g. Rightscale)

Alternative approach: bundle management

- Management logic integrated with application instances
 - Instances use P2P-style self-organization
- Application “deploys itself”



Related work

- We build on existing work
 - Early PARC “worm” programs (1982)
 - Intelligent and mobile agents (1990s)
 - Autonomic computing (2000-)
- Apply the ideas to cloud environments...
 - Span multiple providers and configurations
 - Application controls resource allocation

Example application requirements (PsEPR)

- Planetary Scale Event Propagation and Router (PsEPR)
 - A unified reporting and event service for PlanetLab

P.Brett et al., “A Shared Global Event Propagation System to Enable Next Generation Distributed Services”, WORLDS’04.
- Requirements
 - 6-10 well-connected nodes
 - Lightly-loaded
 - Geographically distributed

 - Maximize CPU cycles
 - Minimize network diameter
 - Reduce budget and migration cost

Example application requirements (PsEPR)

- 6-10 well-connected nodes
- Lightly-loaded
- Geographically distributed

- Maximize CPU cycles
- Minimize network diameter
- Reduce budget and migration cost

Example application requirements (PsEPR)

- **Constraints**
 - 6-10 well-connected nodes
 - Lightly-loaded
 - Geographically distributed

 - Maximize CPU cycles
 - Minimize network diameter
 - Reduce budget and migration cost

Example application requirements (PsEPR)

- **Constraints**
 - 6-10 well-connected nodes
 - Lightly-loaded
 - Geographically distributed
- **Objective**
 - Maximize CPU cycles
 - Minimize network diameter
 - Reduce budget and migration cost

Example application requirements (PsEPR)

- **Constraints**
 - 6-10 well-connected nodes
 - Lightly-loaded
 - Geographically distributed
- **Objective**
 - Maximize CPU cycles
 - Minimize network diameter
 - Reduce budget and migration cost
- **Constraint Optimization Problem**

Idea: Apply constraint logic programming (CLP)

- Constraint Programming
 - Constraints
 - Objective function = Utility function - Cost function

Idea: Apply constraint logic programming (CLP)

- Constraint Programming
 - Constraints
 - Objective function = Utility function - Cost function

Utility function: value of a given configuration

Idea: Apply constraint logic programming (CLP)

- Constraint Programming
 - Constraints
 - Objective function = Utility function - Cost function

Cost function: cost of moving to a new configuration

Idea: Apply constraint logic programming (CLP)

- Constraint Programming
 - Constraints
 - Objective function = Utility function - Cost function

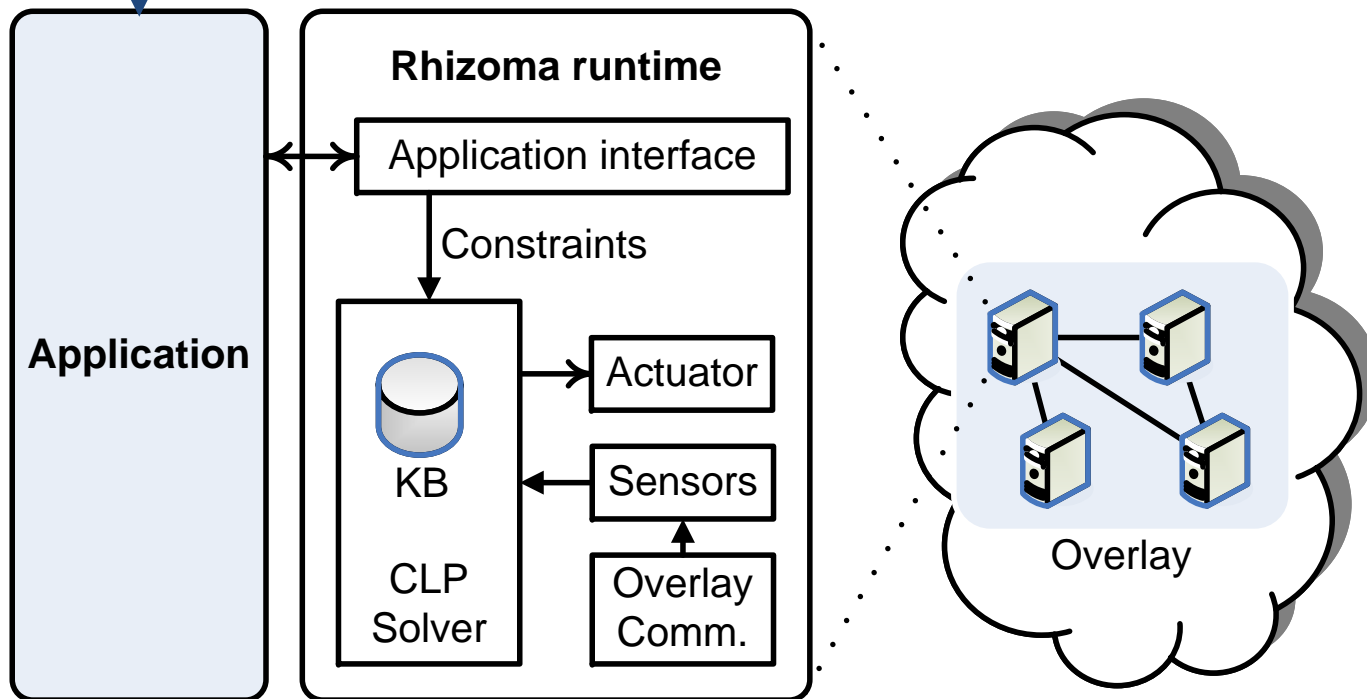
- Logic Programming
 - Heterogeneous resource information
 - Inference rules

PsEPR requirements in CLP

CONSTRAINTS	UTILITY FUNCTION
<pre>node_constraint(Host) :- comonnode{hostname: Host}, alive(Host), light_loaded(Host), is_avail(Host), get_node_attr(Host, cpuspeed, Cpuspeed), get_node_attr(Host, freecpu, Freecpu), Cpuspeed*Freecpu/100 > 1.5.</pre>	<pre>util_function(NodeList, Util, Params) :- % Compute utility values for different node attributes <i>fiveminloadUtil</i>(LoadMin, LoadMax, LoadWeight), assemble_values(NodeList, fiveminload, LoadList), util_value("<", LoadList, LoadMin, LoadMax, LoadUtil), % Omitting utility values for liveslices and freecpu, the same as above</pre>
<pre>group_constraint(NodeList) :- assemble_values(NodeList, location, Locs), length(NodeList, Len), Max is ((Len-1)//4)+1, (for(I, 1, 4), param(Locs, Max) do count_element(I, Locs, Num), Num =<= Max).</pre>	<pre>% Utility of max distance from fixed nodes to the overlay findall(P, fixednode(P), Fixed), <i>minlatency</i>(MinLat), <i>maxneighUtil</i>(_, NeighMax, NeighWeight), get_nearest_neighbor_list(Fixed, NodeList, NeighList), max(NeighList, MaxDist), util_value("<", [MaxDist], MinLat, NeighMax, NeighUtil),</pre>
<pre>path_constraint(LenList, Max) :- max(LenList, Max), <i>diameterUtil</i>(_, DiameterMax, _), Max < DiameterMax.</pre>	<pre>% Utility of overlay network diameter <i>diameterUtil</i>(_, DiamMax, DiamWeight), util_value("<", Params, MinLat, DiamMax, DiamUtil),</pre>
COST FUNCTION	
<pre>migration_cost(Actions, MigrateCost) :- count_element(add, Actions, AddLen), count_element(remove, Actions, RmvLen), <i>addCostParam</i>(AddParam), <i>removeCostParam</i>(RmvParam), MigrateCost is AddParam*AddLen + RmvParam*RmvLen.</pre>	<pre>% Weighted average of the utilities above weighted_avg([LoadUtil, SliceUtil, CpuUtil, NeighUtil, DiamUtil], [LoadWeight, SliceWeight, CpuWeight, NeighWeight, DiamWeight], Util). Definition of util_value: util_value_{<} = avg_i ((x_{max} - bound(x_i))/(x_{max} - x_{min})) util_value_{>} = avg_i ((bound(x_i) - x_{min})/(x_{max} - x_{min})) bound(x) = max(min(x, x_{max}), x_{min})</pre>
CONFIGURATIONS	
<pre><i>fiveminloadUtil</i>(0, 10, 2). <i>liveslicesUtil</i>(0, 10, 1). <i>freecpuUtil</i>(1, 4, 3). <i>maxneighUtil</i>(0, 500, 2). <i>diameterUtil</i>(0, 1000, 2). <i>addCostParam</i>(0.012). <i>removeCostParam</i>(0).</pre>	

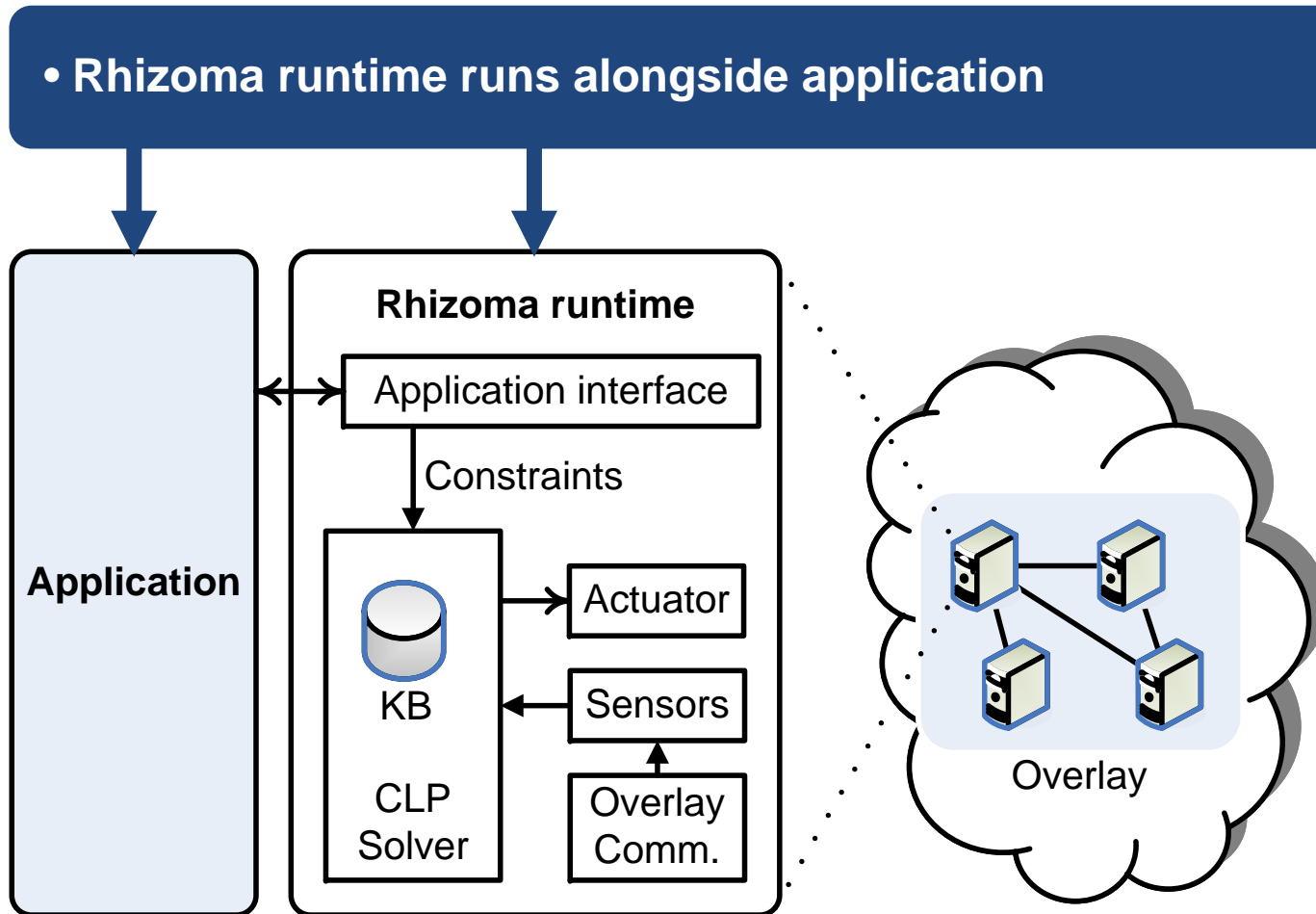
Our system: Rhizoma

- Application is a self-organizing overlay



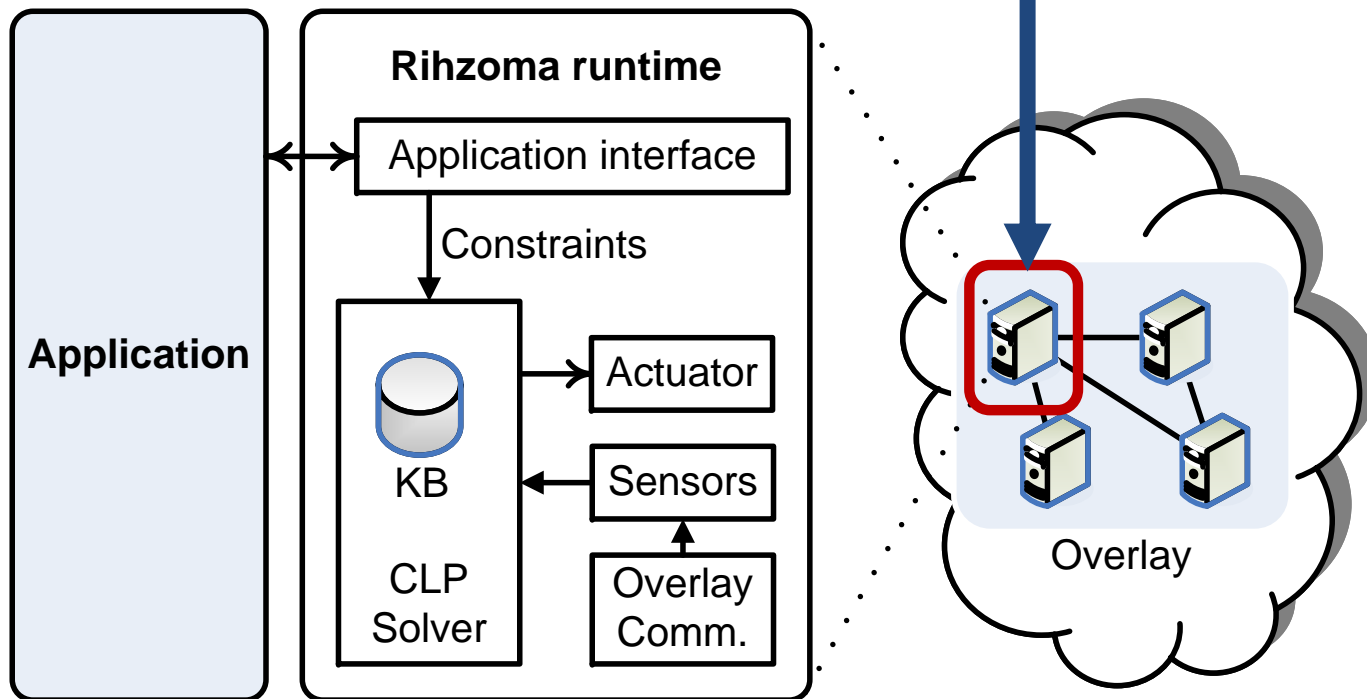
Our system: Rhizoma

- Rhizoma runtime runs alongside application



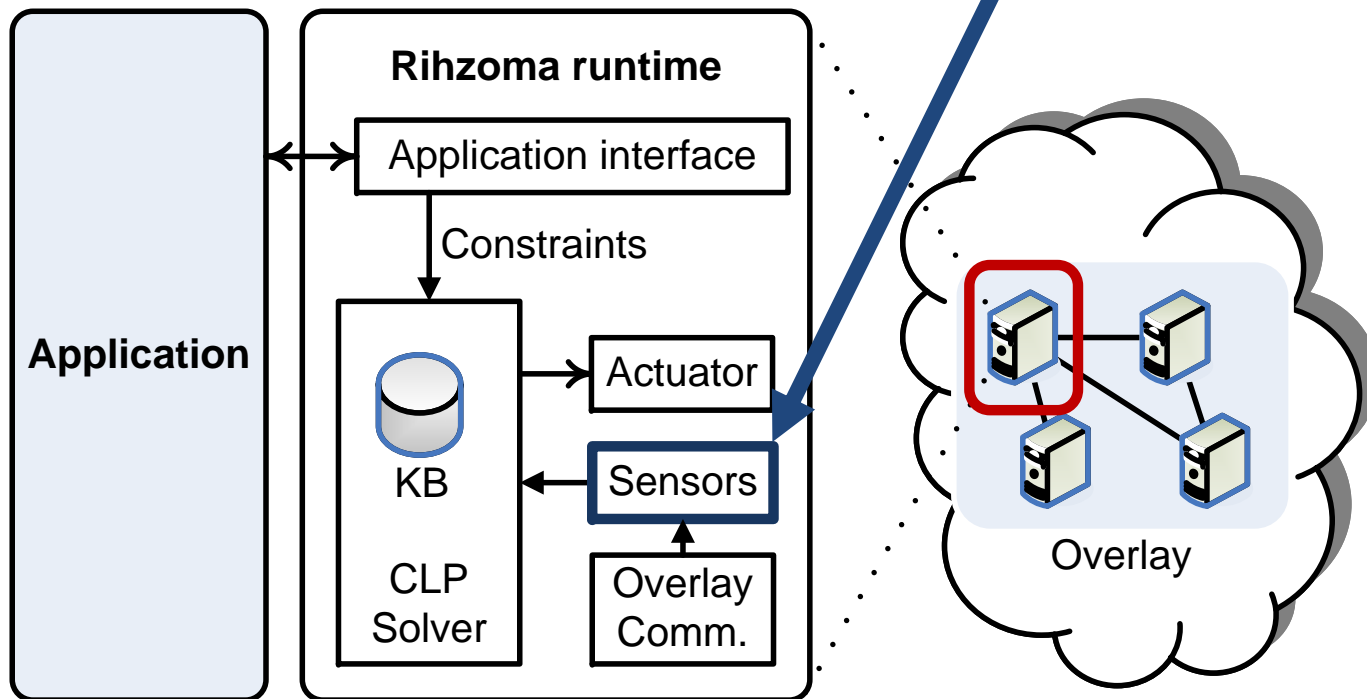
Our system: Rhizoma

- One node is elected as *coordinator* and runs CLP solver



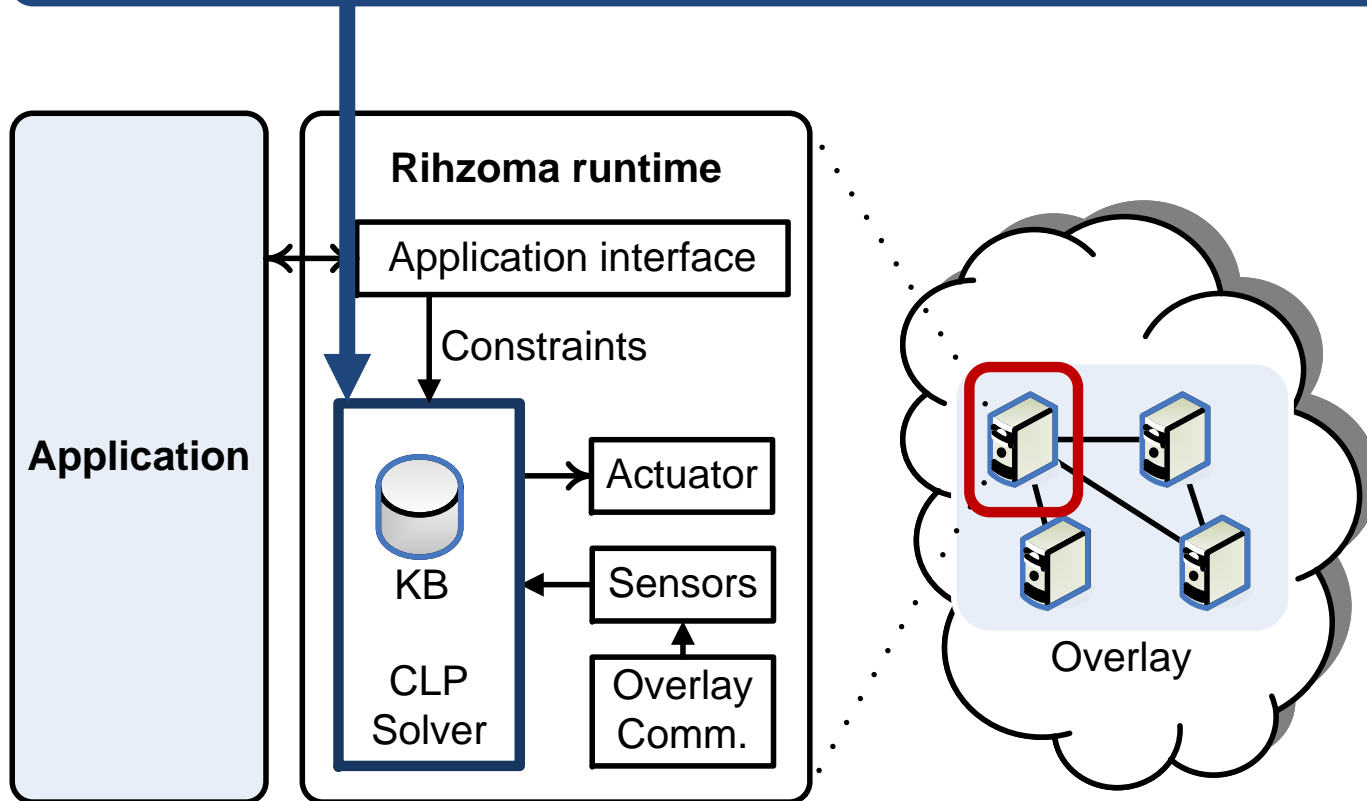
Our system: Rhizoma

- Sensors gather info from overlay and external service



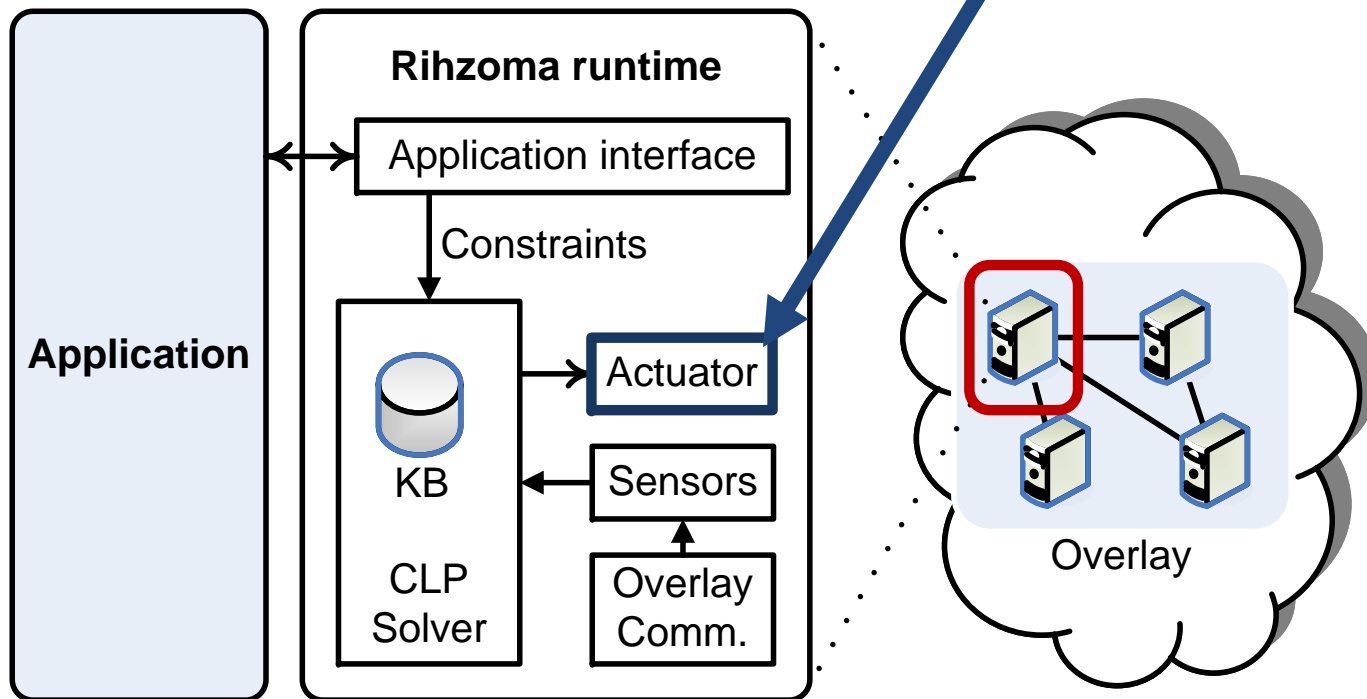
Our system: Rhizoma

- CLP solver calculates better deployment



Our system: Rhizoma

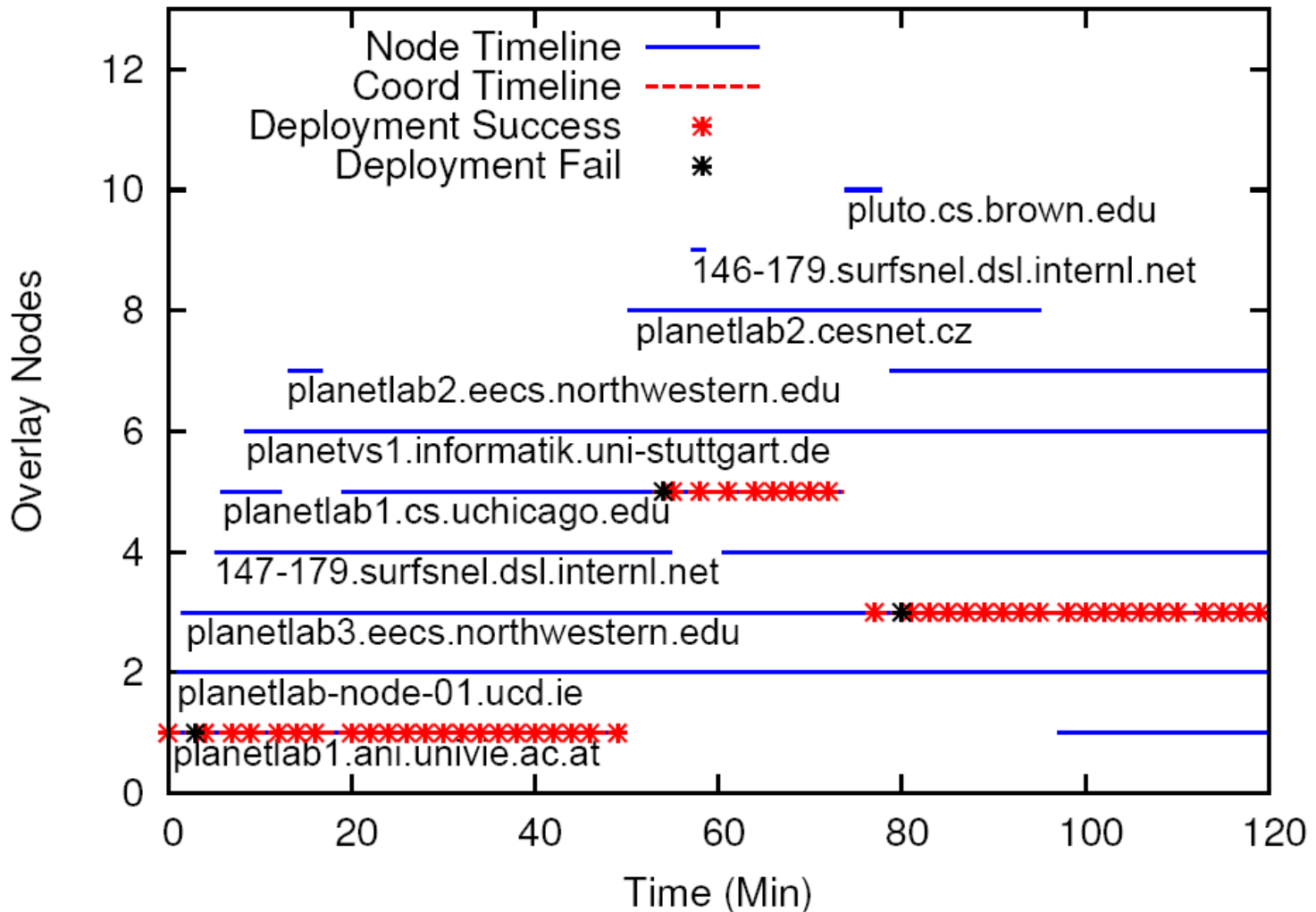
- Actuator acquires/releases VMs, starts/stops application



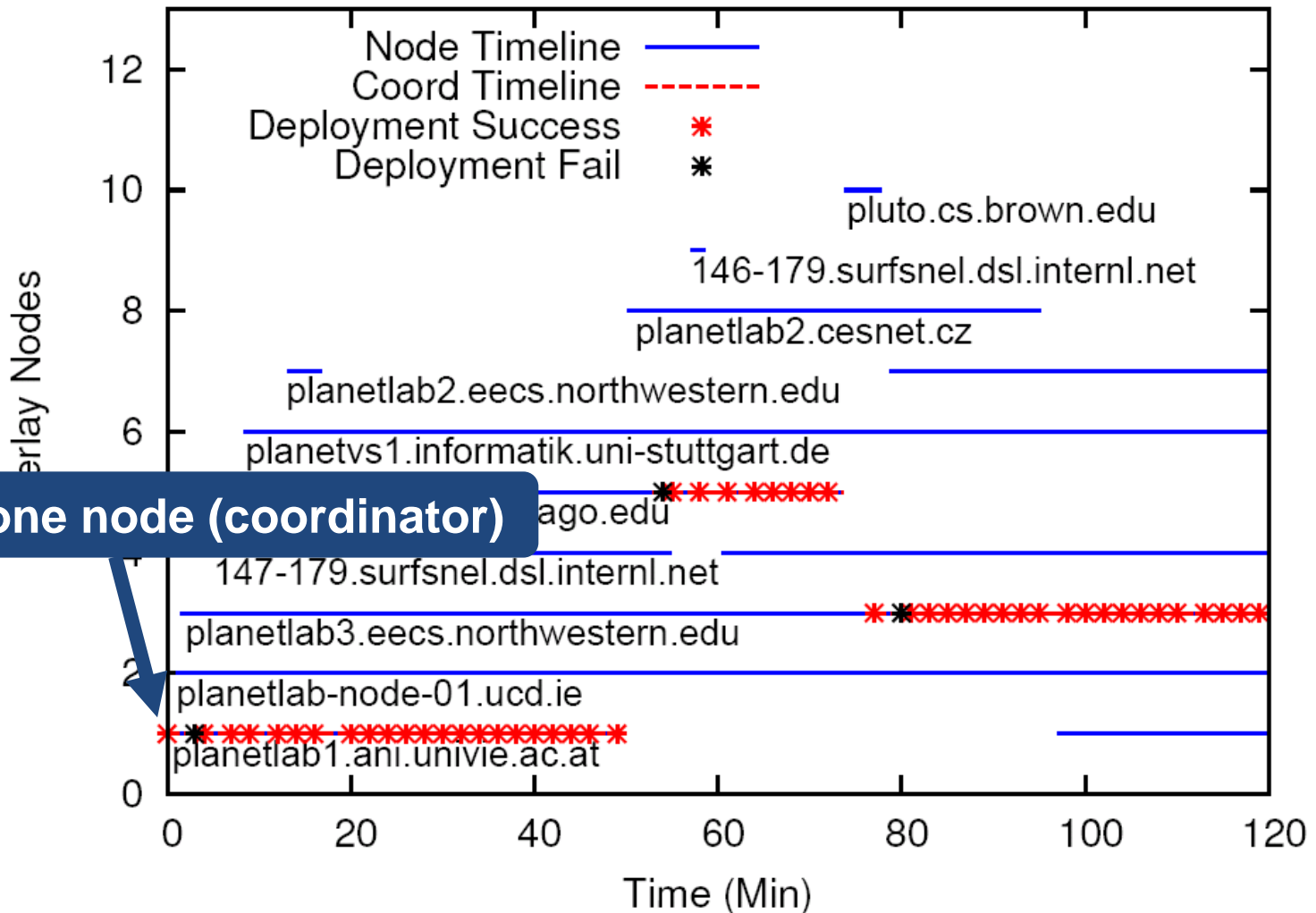
Evaluation

- Rhizoma runtime
- PsEPR requirements
- 8 hour trace running on PlanetLab
- Overlay behavior
- Utility value
- Weighted attributes

Rhizoma hosting PsEPR on PlanetLab



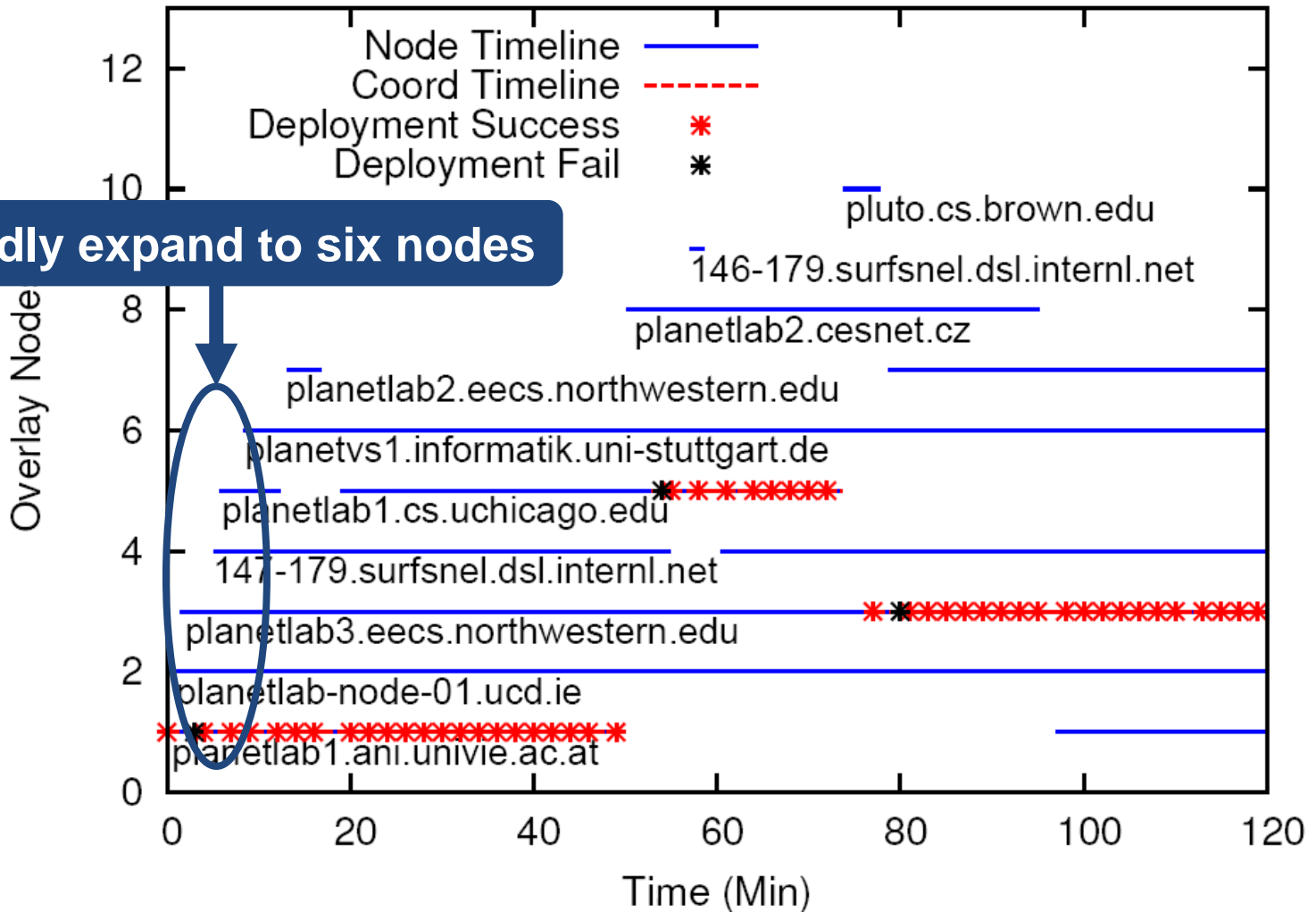
Rhizoma hosting PsEPR on PlanetLab



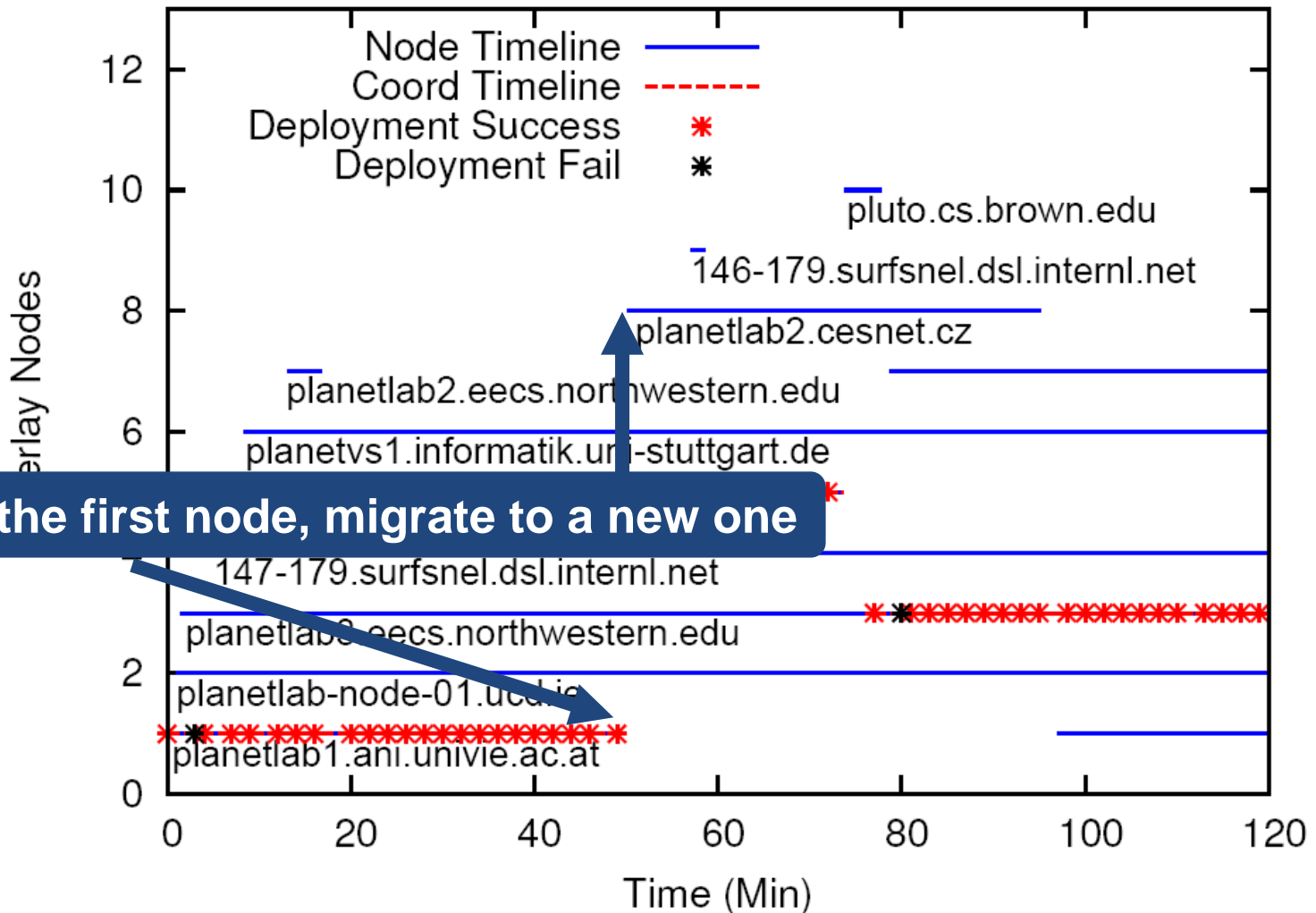
Start on one node (coordinator)

Rhizoma hosting PsEPR on PlanetLab

Rapidly expand to six nodes

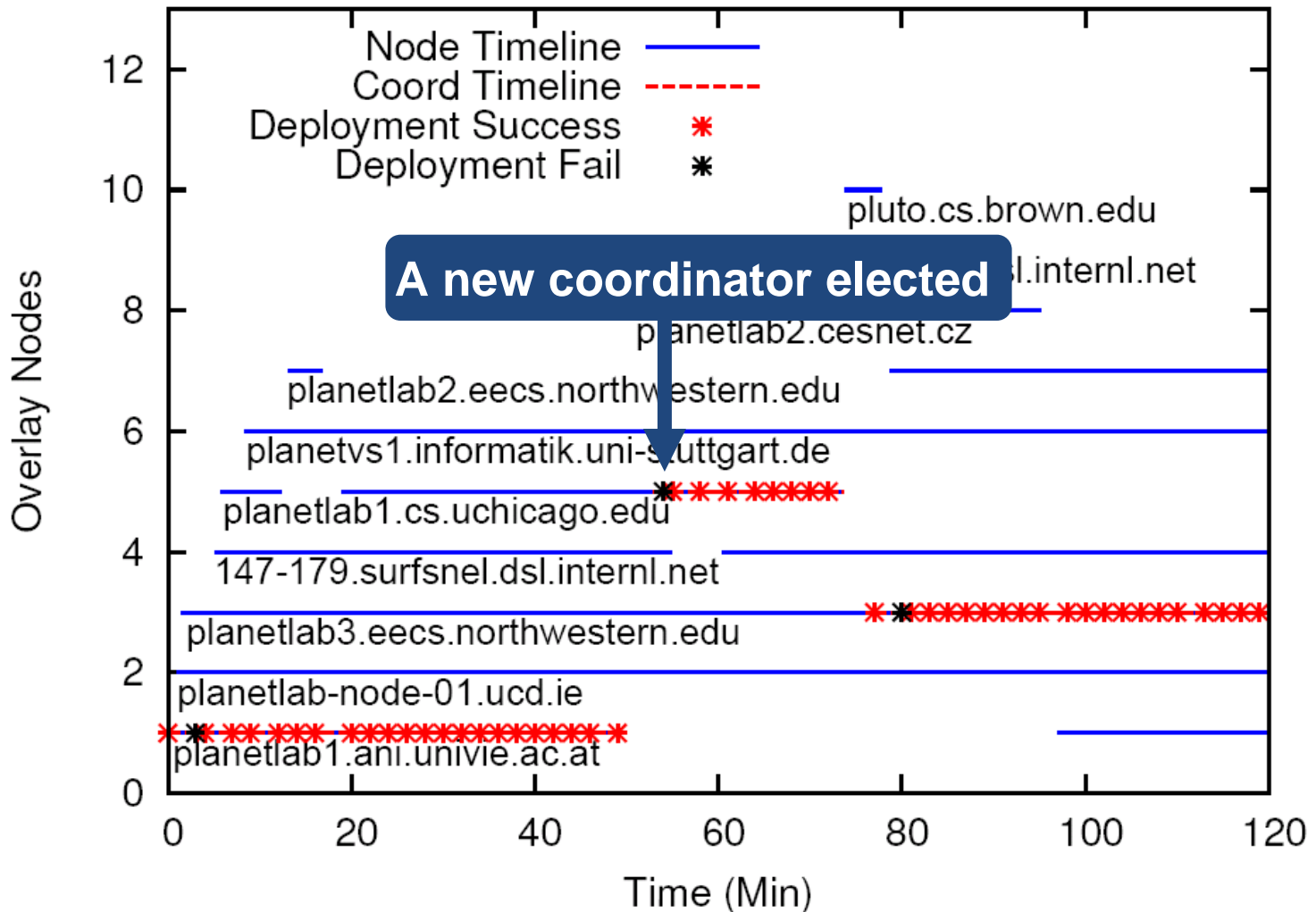


Rhizoma hosting PsEPR on PlanetLab

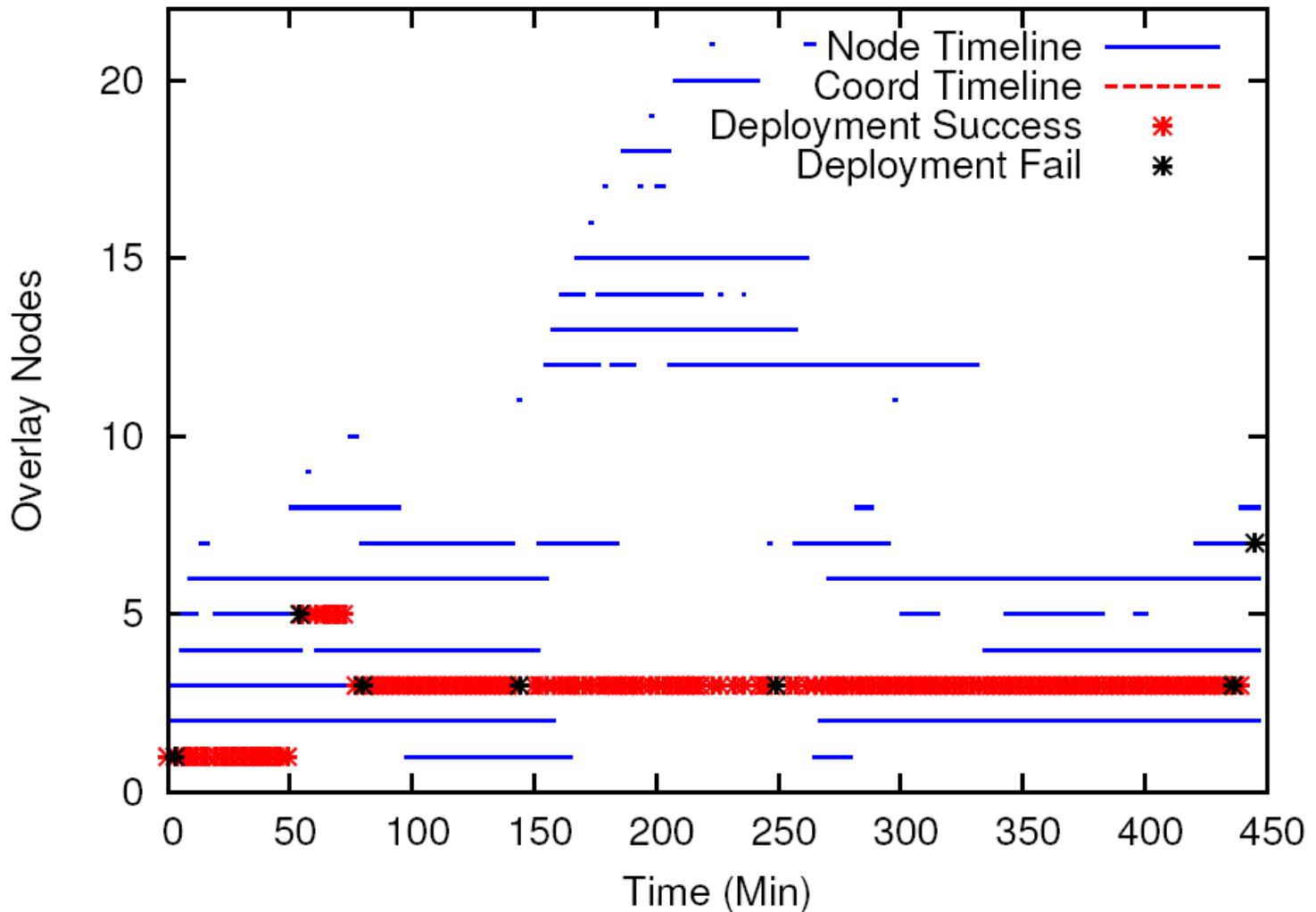


Vacate the first node, migrate to a new one

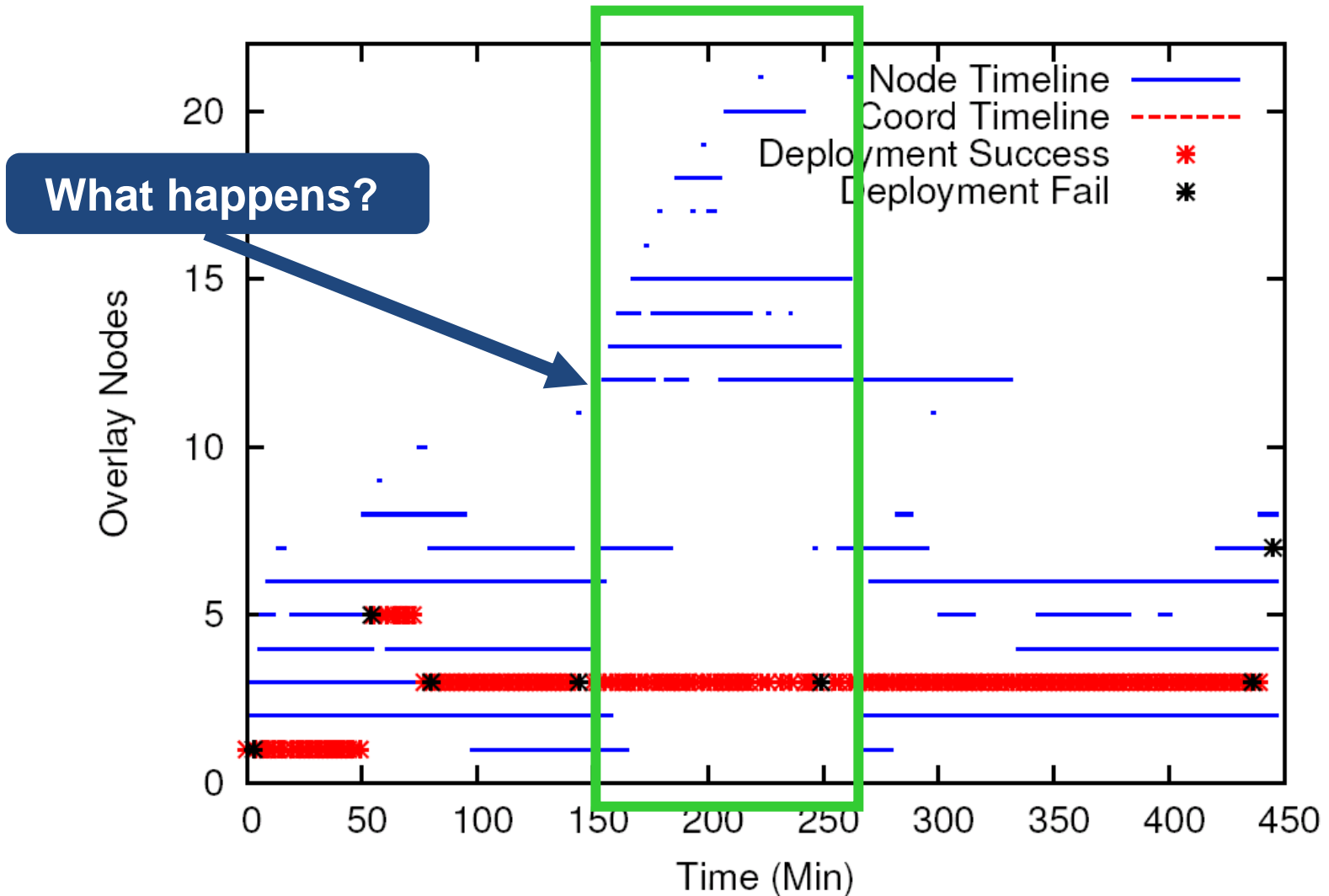
Rhizoma hosting PsEPR on PlanetLab



Rhizoma hosting PsEPR – longer trace

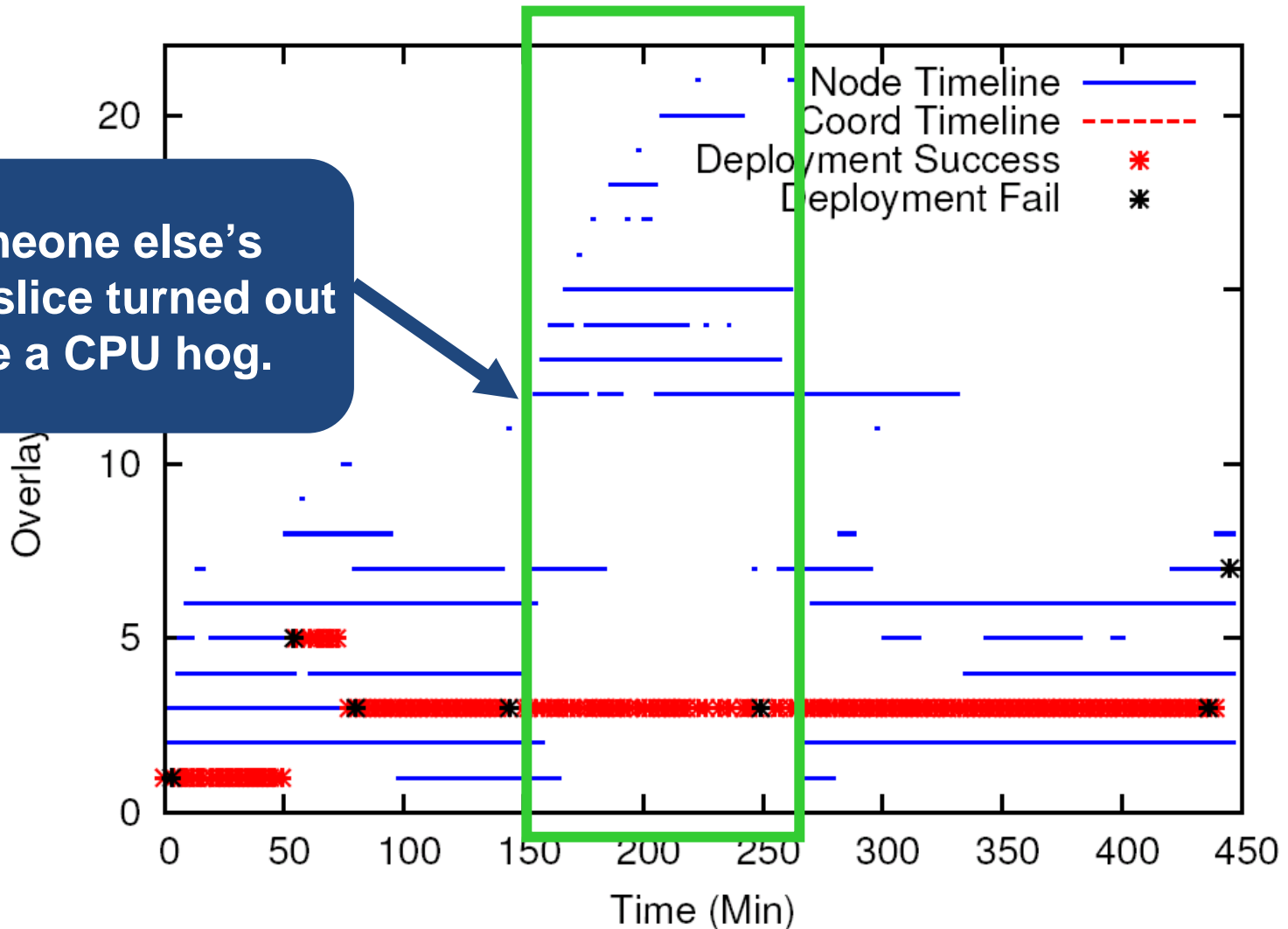


Rhizoma hosting PsEPR – longer trace



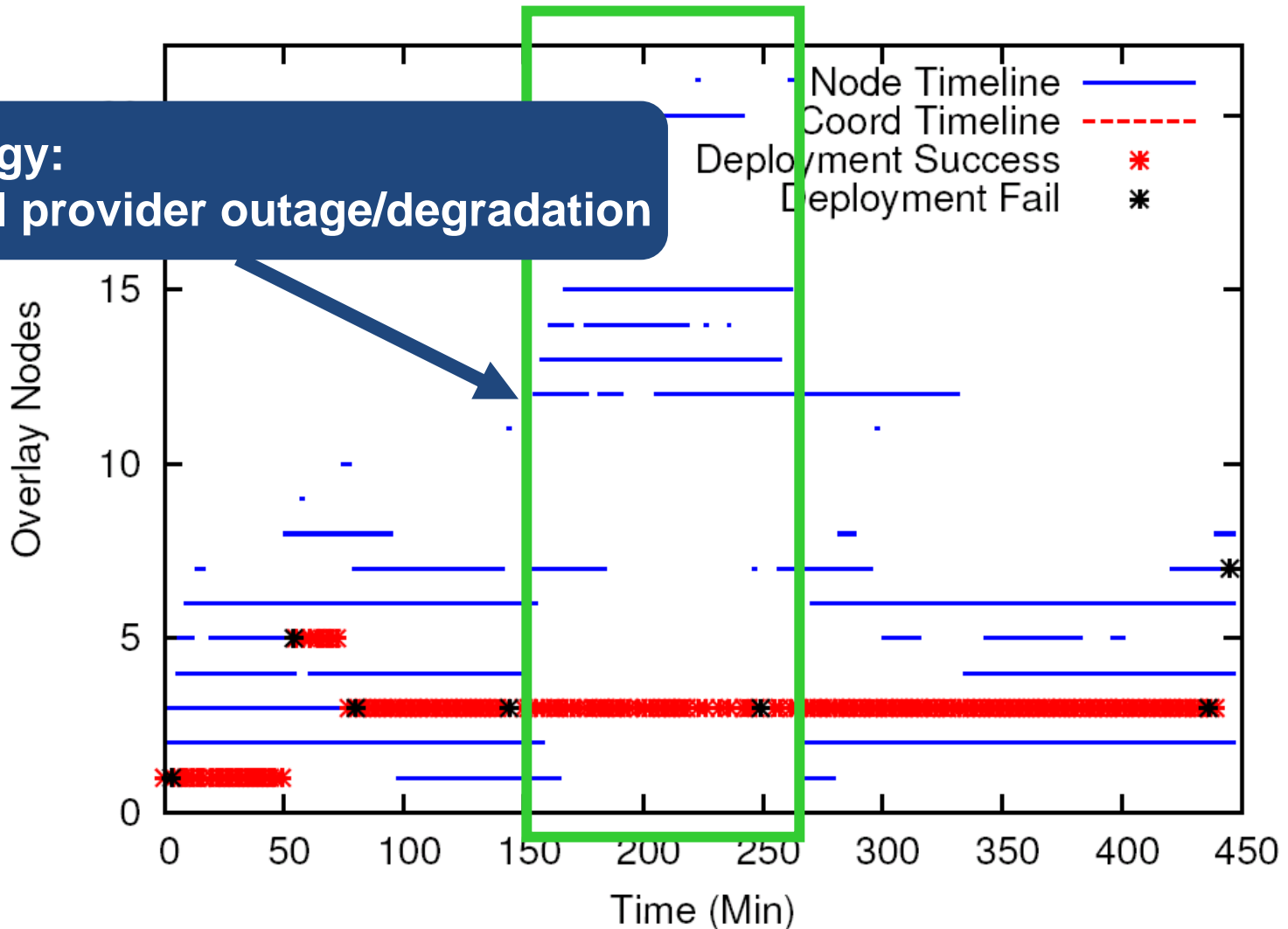
Rhizoma hosting PsEPR – longer trace

Someone else's
buggy slice turned out
to be a CPU hog.

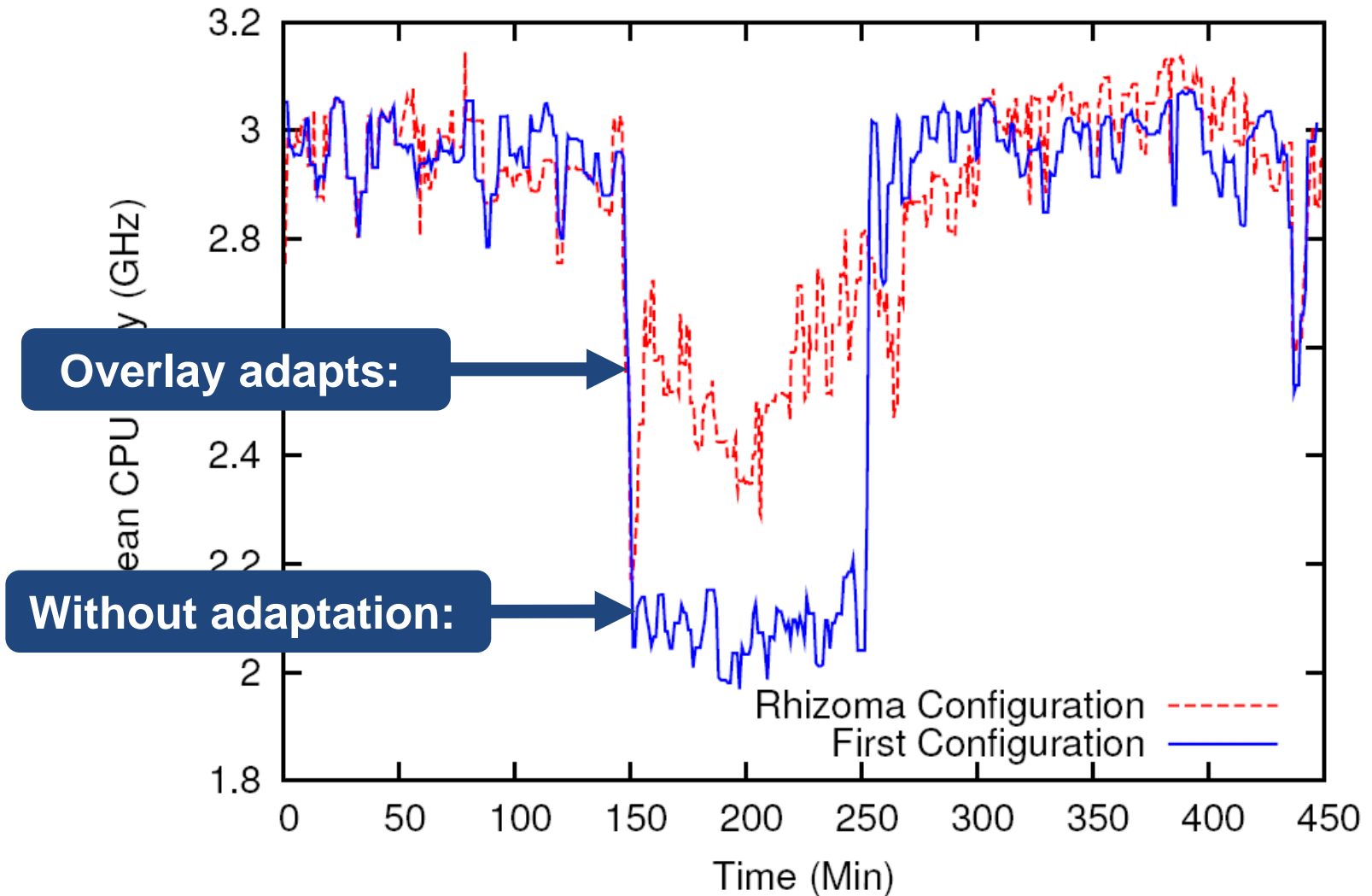


Rhizoma hosting PsEPR – longer trace

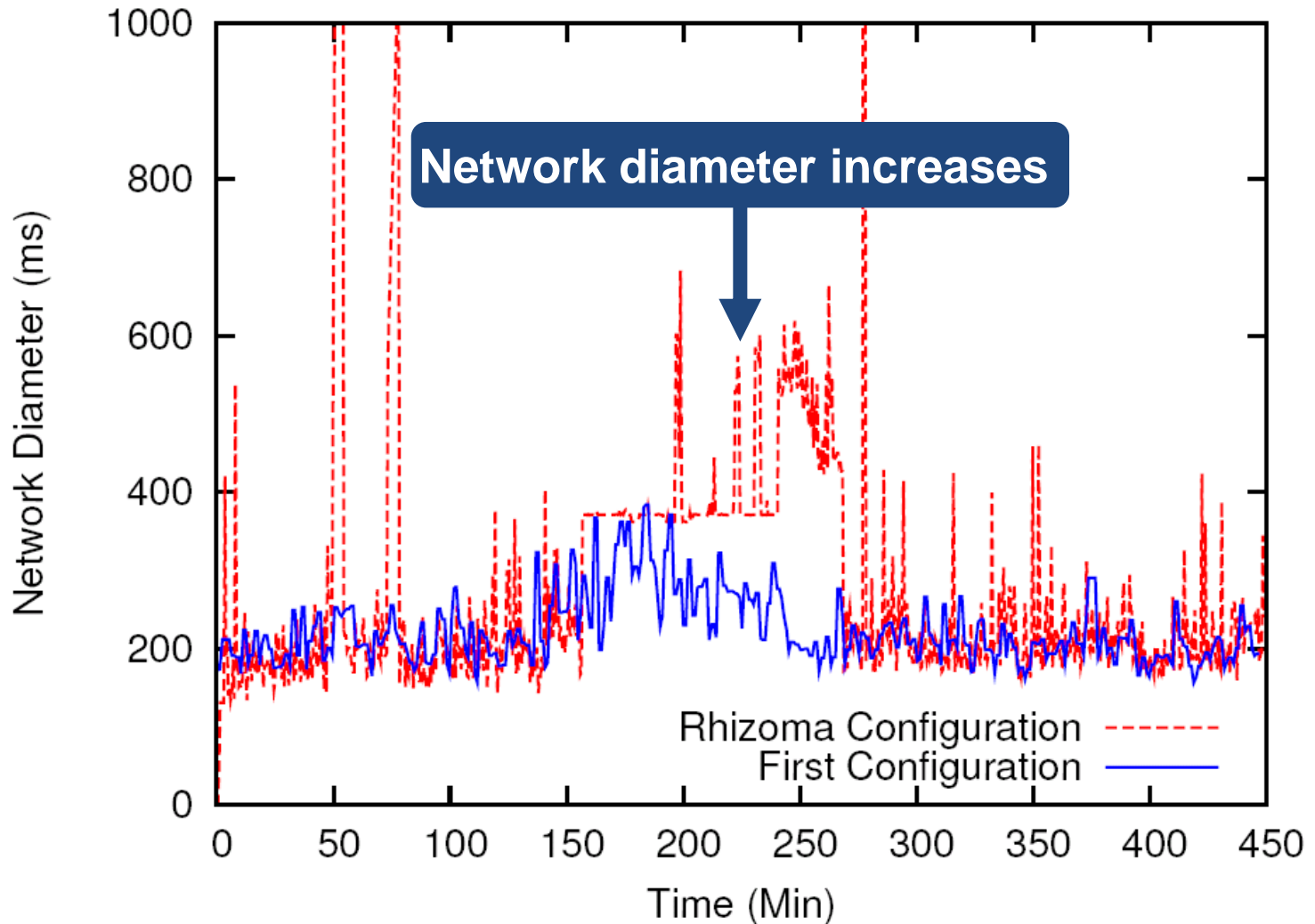
**Analogy:
Partial provider outage/degradation**



Mean CPU (GHz)

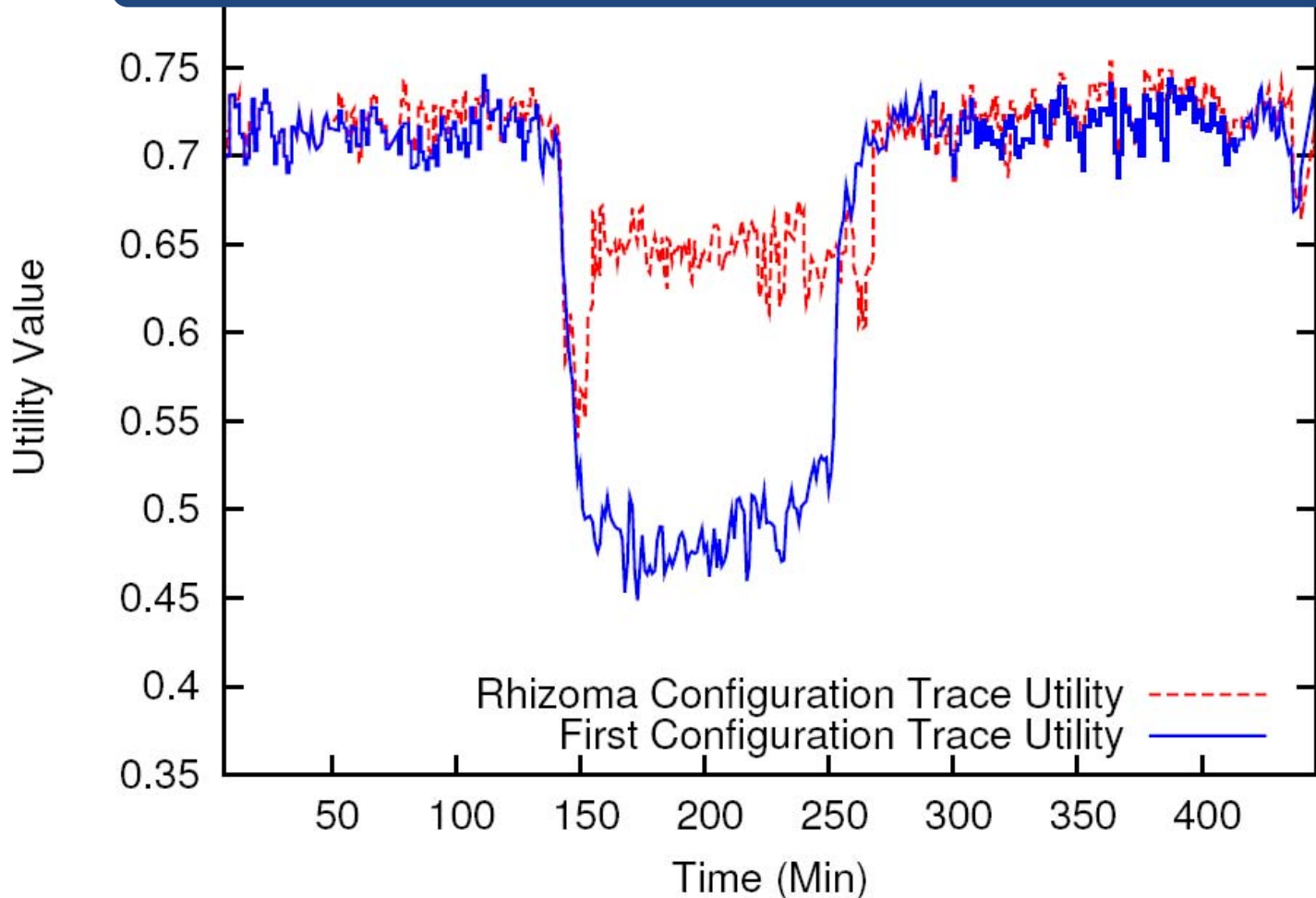


Network diameter (ms)

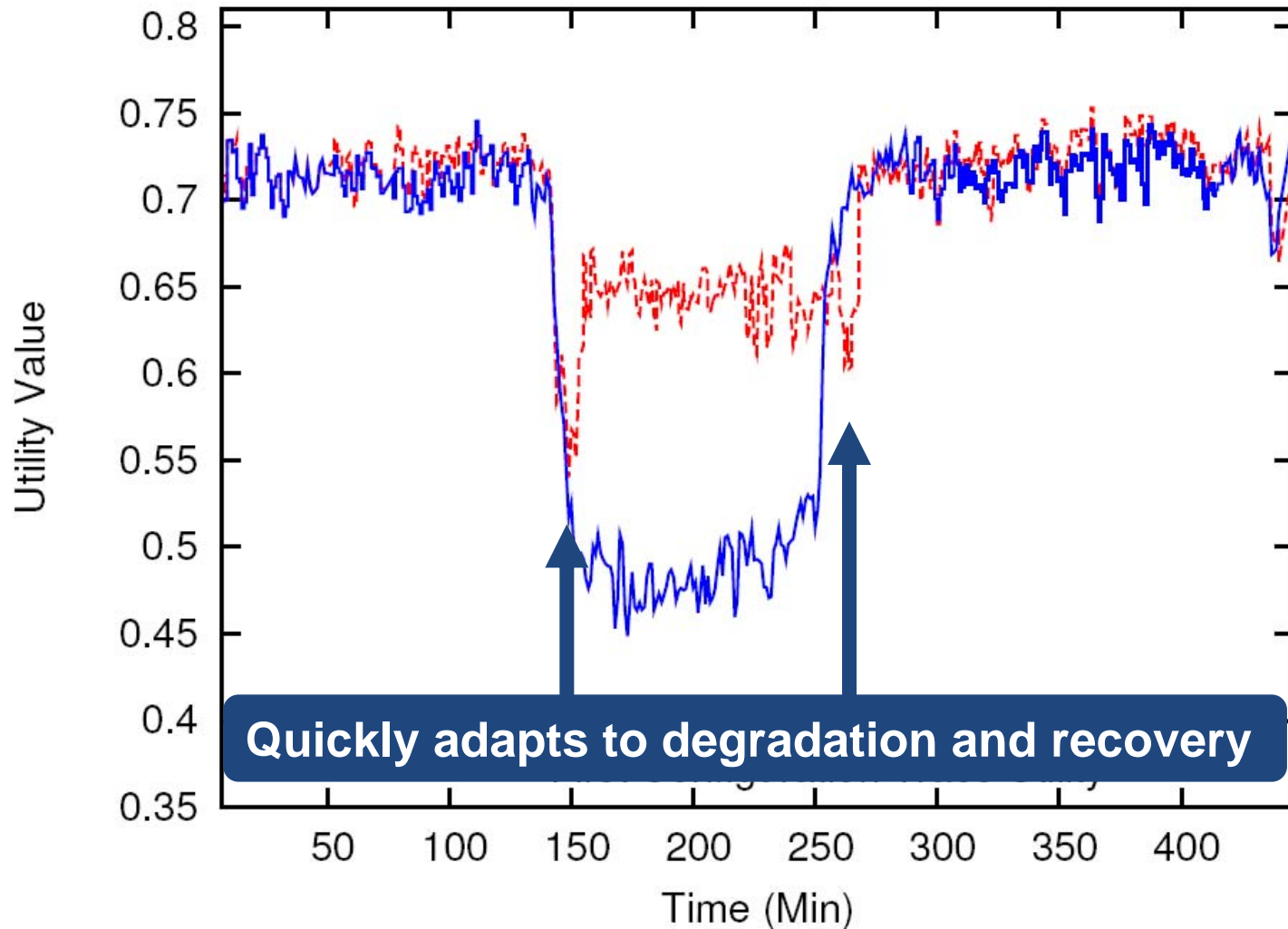


Utility values

Mean CPU is weighted more than network diameter



Utility values



Summary

- Key ideas:
 - **Runtime** for **self-deploying, self-managing** overlays
 - **CLP** to express desired application behavior

Summary

- Key ideas:
 - **Runtime** for **self-deploying, self-managing** overlays
 - **CLP** to express desired application behavior
- Appealing features:
 - **Expressiveness** of applications' resource requirements
 - **Optimization** in terms of performance and cost
 - **Responsiveness** to changes in resources, load, policies

Summary

- Key ideas:
 - **Runtime** for **self-deploying, self-managing** overlays
 - **CLP** to express desired application behavior
- Appealing features:
 - **Expressiveness** of applications' resource requirements
 - **Optimization** in terms of performance and cost
 - **Responsiveness** to changes in resources, load, policies
- Future work:
 - Platform integration (cluster, PlanetLab, Amazon EC2)
 - Feed application-level metrics to optimization process