



# Dynamic Deployment and Scalability for the Cloud

Jerome Bernard  
Director, EMEA Operations  
Elastic Grid, LLC.

## Speaker's qualifications

- Jerome Bernard is a committer on Rio, Typica, JiBX and co-founder of Elastic Grid, LLC.
- Jerome Bernard speaks frequently on Cloud Computing
  - Recently: QCon, Devoxx, JavaZone, JavaOne, and the Open World Forum
- Jerome Bernard is working with many clients using EC2, from TV channels to specialized media processing companies.



## Agenda

Introduction to Cloud Computing  
Introduction to Amazon EC2  
Introduction to Elastic Grid

# Introduction to Cloud Computing

- Why Cloud Computing?
  - ▶ Next logical step after virtualization
    - Better usage of your IT infrastructure
    - Cost Savings
  - ▶ Can your traditional hosting scale to thousands of machines in a week?
  - ▶ Can you afford spending huge amounts buying hardware if you only need it for a week?

Virtualization is used for consolidation.

Cloud Computing allow you to rent resources when they are needed.

## Animoto Use Case

Company (US startup) creating cool videos based on a bunch of uploaded pictures. Really CPU intensive. Went from dozen of servers up to 3500 servers in a few days when their application was released on Facebook. But went down to a few hundreds after another week.

How would you cope with that situation in a few days? Would you be able to raise money from VCs, buy the hardware, have the dealer send you the machine, install them and put them in a datacenter in just a few days?

What would do a week after with all the servers you don't need anymore?

# Introduction to Cloud Computing

- Which Cloud Computing flavor?
  - ▶ Software as a Service (SaaS)
  - ▶ Platform as a Service (PaaS)
  - ▶ Infrastructure as a Service (IaaS)

IaaS: you rent some infrastructure -> some servers

PaaS: you rent access to a platform hosting your applications.

- References

- ▶ SaaS: Salesforce, Facebook, LinkedIn
- ▶ PaaS: Salesforce (EC2), Google App Engine, Microsoft Azure
- ▶ IaaS: Amazon EC2, GoGrid, Flexiscale

# Introduction to Cloud Computing

- Google App Engine
  - ▶ Make use of BigTable and Memcache
  - ▶ Integrate with Google Accounts
  - ▶ But in Python only...
- Microsoft Azure
  - ▶ Mostly for Windows and .Net solutions
  - ▶ Pricing model yet unclear

## PaaS vs IaaS

- PaaS Pros
  - ▶ Usually easier to use than IaaS
  - ▶ Integrate with specific environments (Google, Microsoft Live, Salesforce, etc.)
- PaaS Cons
  - ▶ Less/No control over the Infrastructure
  - ▶ Languages/Services chosen by the provider
  - ▶ Vendor Lock-in

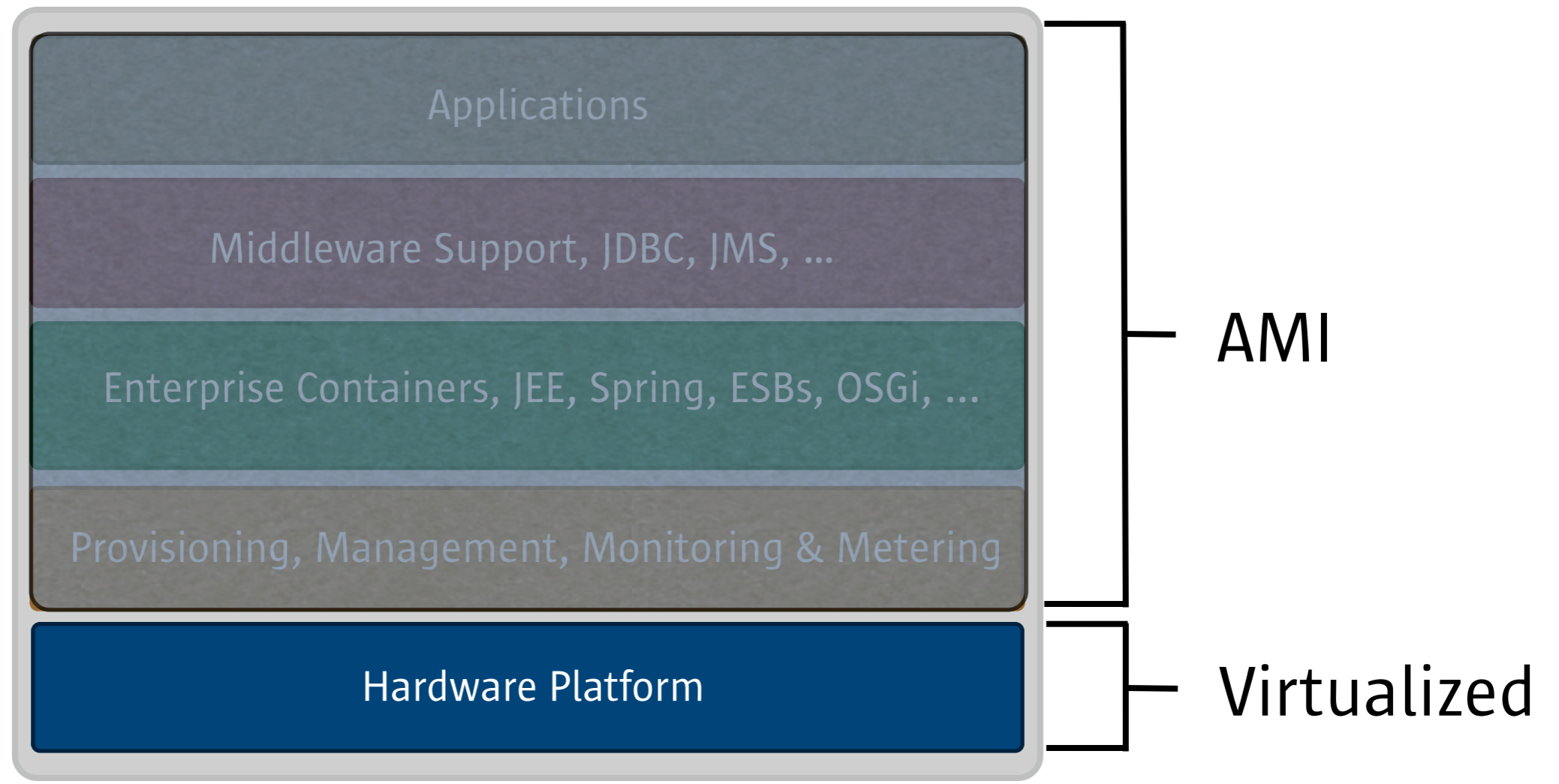
# Introduction to Amazon EC2

- Amazon EC2 is Infrastructure as a Service (IaaS)
  - ▶ Rent a server on a per hour base (from \$.10 to \$.80)
  - ▶ Many Operating Systems (Linux, Solaris, Windows)
- EC2 Amazon Machine Image (AMI)
  - ▶ Operating and system stack
  - ▶ Deployed to Amazon S3 (cheap storage)
- EC2 instances
  - ▶ Virtual machines that run AMI



# Introduction to Amazon EC2

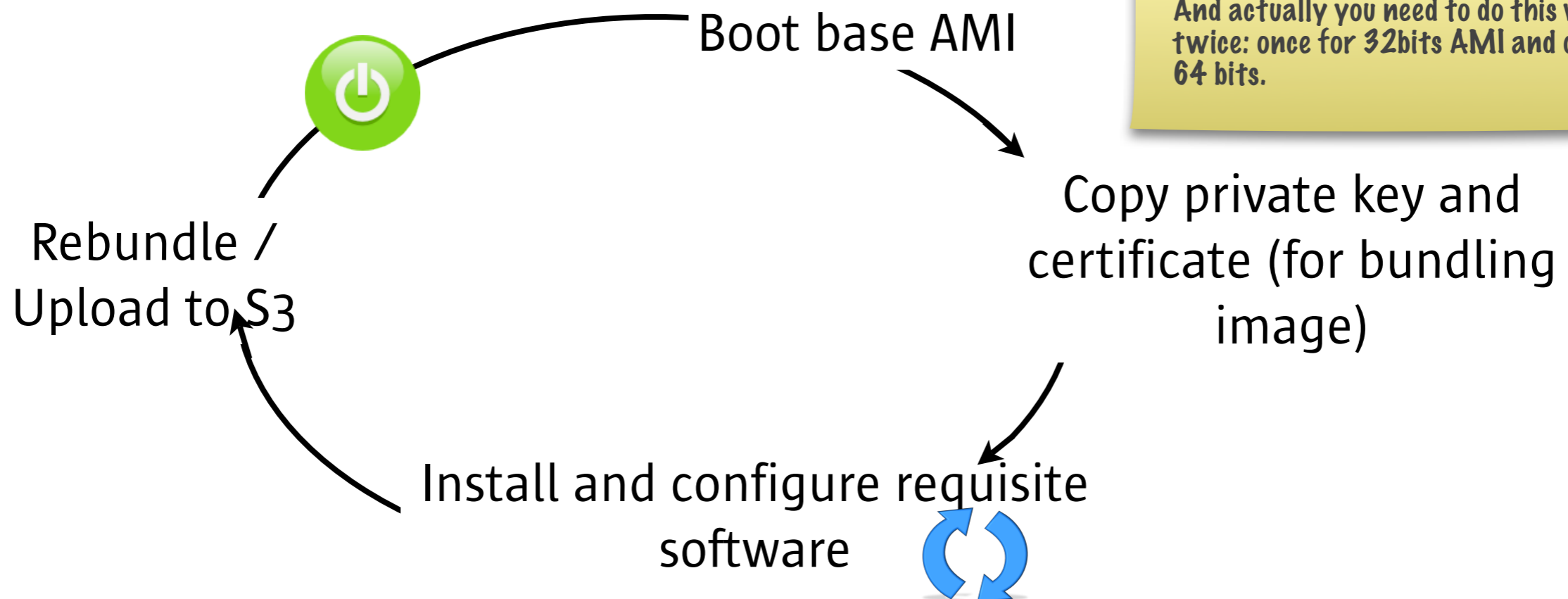
- ▶ Typical Architecture Taxonomy



# Amazon EC2 pitfalls

- EC2 AMI Challenges

- ▶ The EC2 AMI is a boot image, requires substantial system administrator knowledge
- ▶ As application code changes, AMIs typically need to change / be re-bundled



## Amazon EC2 pitfalls (continued)

- Infrastructure challenges
  - ▶ Networking: no multicast but this is what most Java framework uses for clustering (JGroups, Shoal, etc.)
  - ▶ Backup: the local filesystem has no durability guarantee
  - ▶ Significant boot latencies of EC2 instances (can be several minutes)
  - ▶ Failures: you have to design your application to be resilient to EC2 instance failures. Anyway you should always do so :-)

# Amazon EC2 Advice

- Some AWS Advice

- ▶ I/O: prefer an Elastic Block Storage (EBS) volume to a local filesystem
- ▶ Snapshot EBS volumes periodically (incremental backup) but export to S3 for complete backups
- ▶ Choose the right instance type
  - Don't use Small for production!
  - Don't choose based on disk space (think EBS)
  - Choose based on available memory and CPU virtual cores

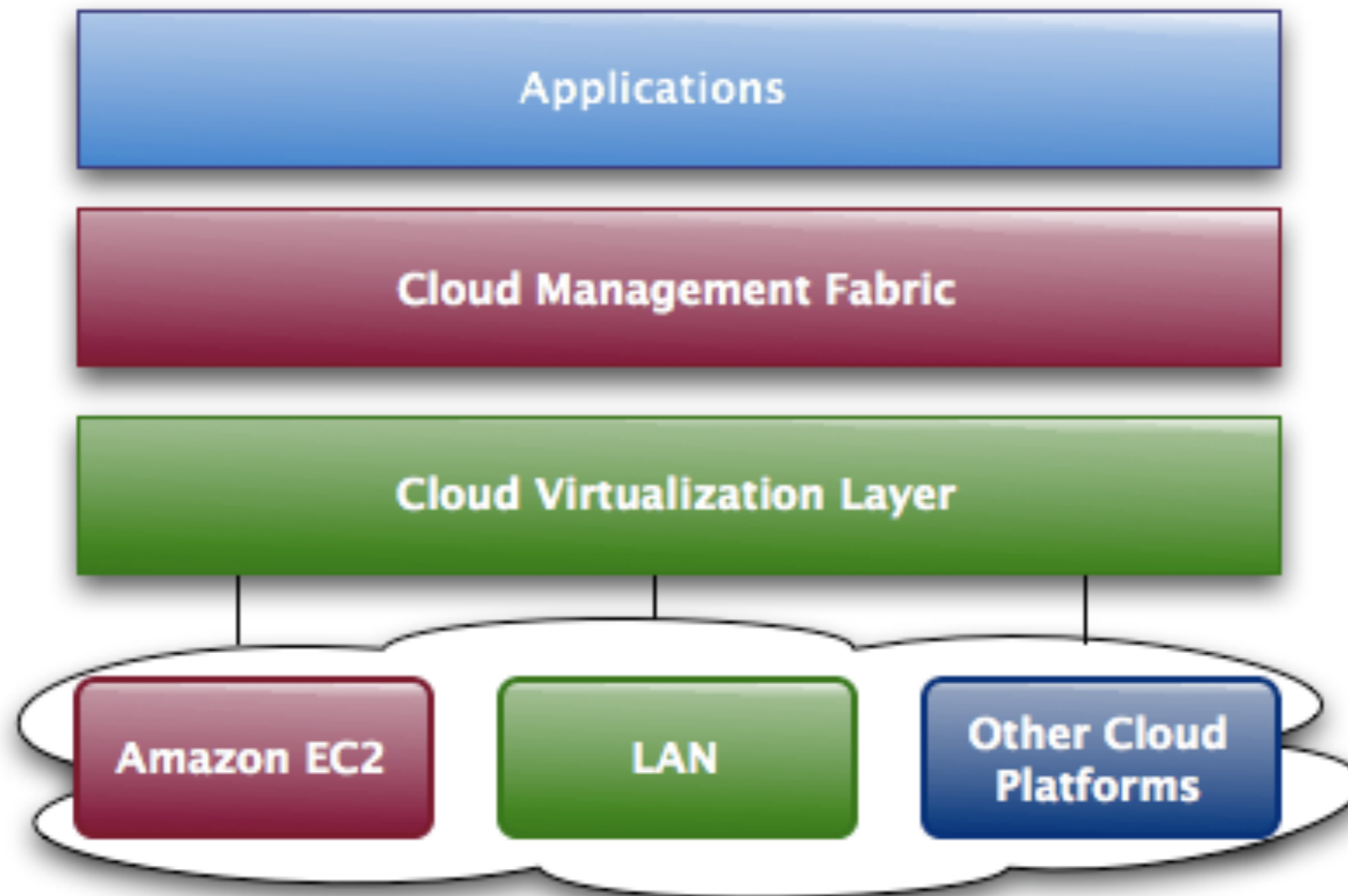
I/O are way better and you benefit from durability, snapshot supports, etc.

High CPU Medium is the best tradeoff usually unless you need a lot of memory.

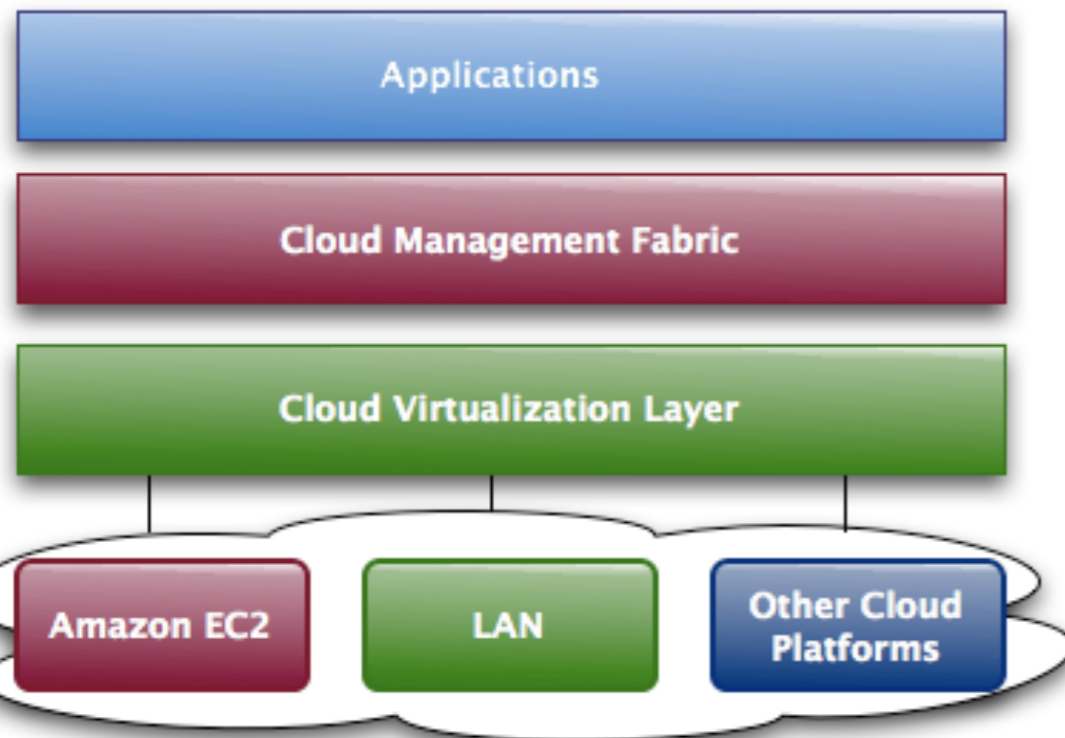
# Introduction to Elastic Grid

- Elastic Grid (abbreviated as EG)
  - ▶ Project initiated in early '08
  - ▶ AGPLv3 license
  - ▶ Part of the OW2 community
- Elastic Grid, LLC. founded in May '08
  - ▶ Team:
    - Dennis Reedy: Director US Operations
    - Jerome Bernard: we already went through this :-)

# Introduction to Elastic Grid



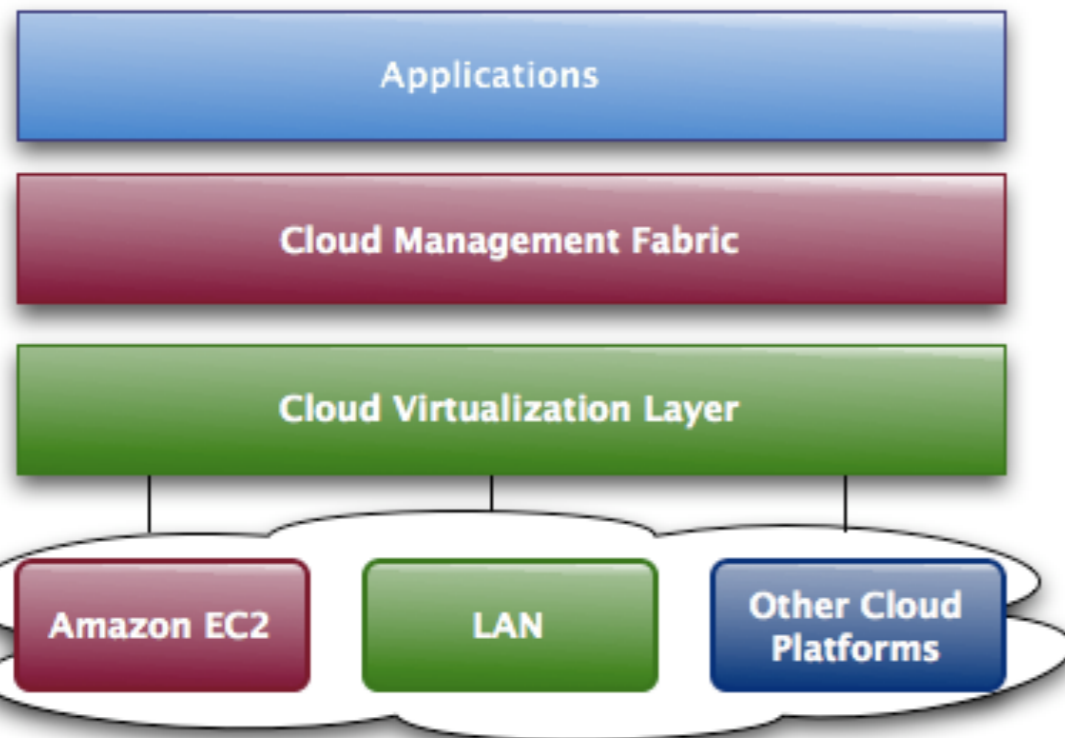
# Introduction to Elastic Grid



## Cloud Management Fabric

- Provides an adaptive capability to dynamically instantiate, monitor & manage application components
- The deployment provides context on service requirements, dependencies, associations and operational parameters
- Provisioning services additionally provides pluggable download distribution and resource

# Introduction to Elastic Grid



## Cloud Virtualization Layer

- Abstracts specific Cloud Computing provider technology
- Allows portability across specific implementations
- You can deploy on:
  - Private Cloud
  - Amazon EC2
  - More to come soon...



## Sound Familiar?

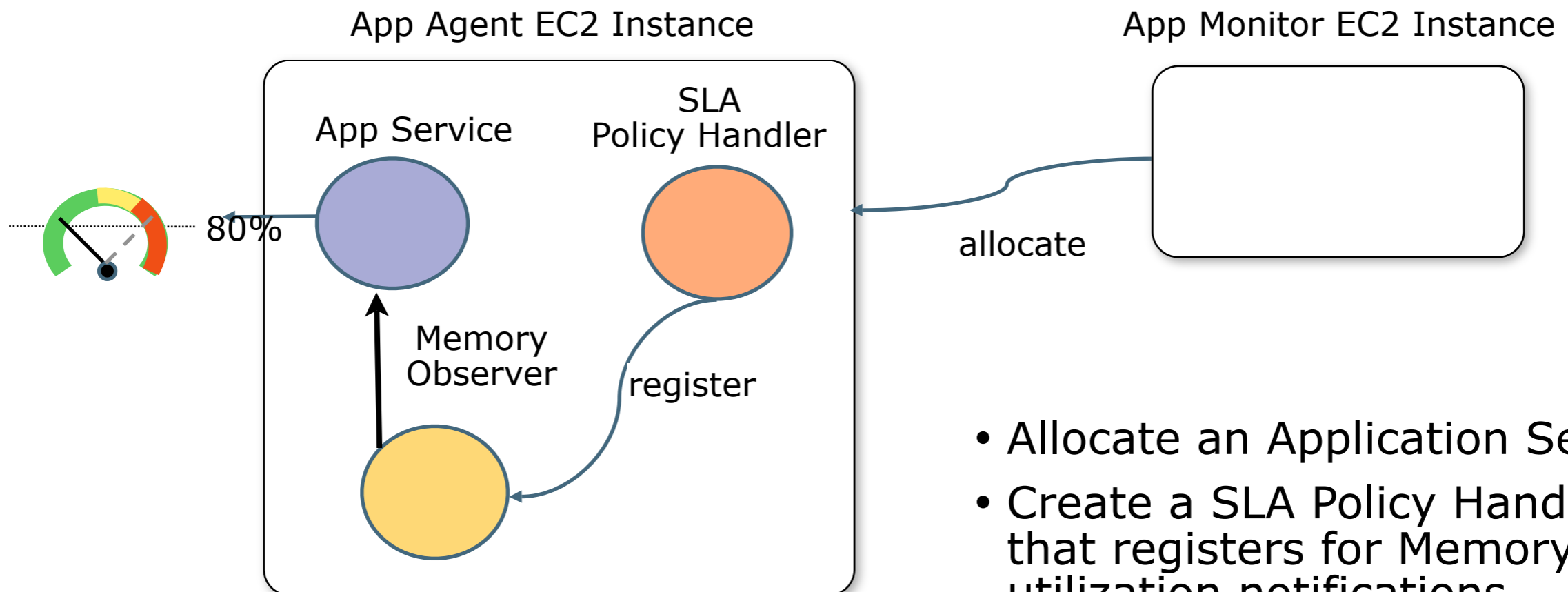
- The network is reliable
- Latency is zero
- Bandwidth is infinite
- The network is secure
- Topology doesn't change
- There is one competent administrator
- Transport costs are zero
- The network is homogenous

*“Essentially everyone, when they first build a distributed application, makes the following assumptions. All prove to be false in the long-run and all cause big trouble and painful learning experiences”.*

Peter Deutsch - “Deutsch’s 8 Fallacies of Networking”

# Elastic Grid Scalability on EC2

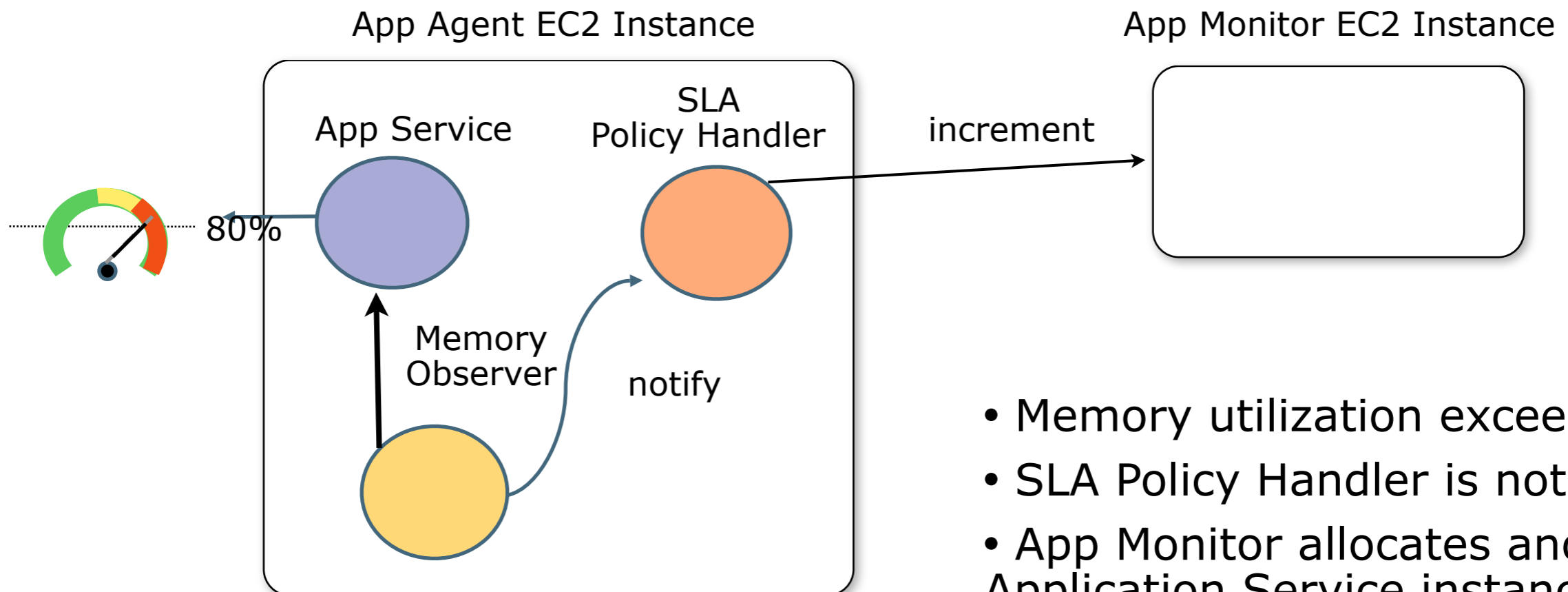
- Across existing EC2 instances



- Allocate an Application Service
- Create a SLA Policy Handler that registers for Memory utilization notifications
- SLA has upper limit set to 80%

# Elastic Grid Scalability on EC2

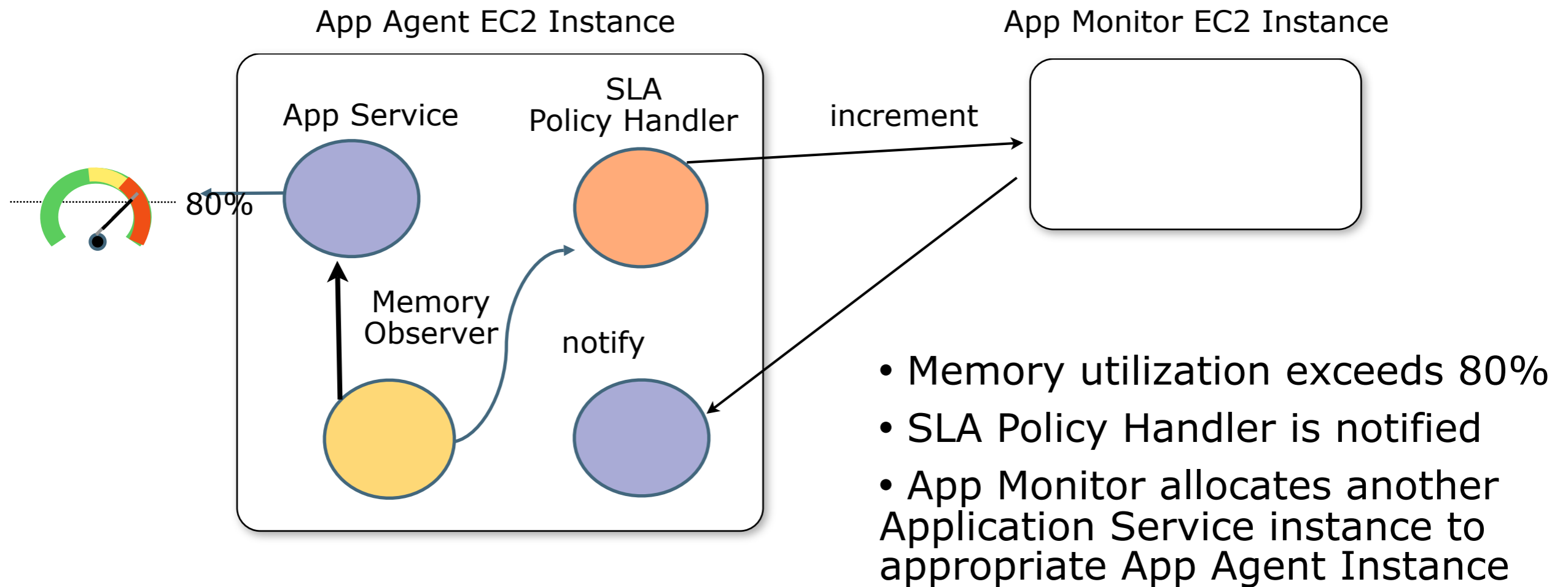
- Across existing EC2 instances



- Memory utilization exceeds 80%
- SLA Policy Handler is notified
- App Monitor allocates another Application Service instance to appropriate App Agent Instance

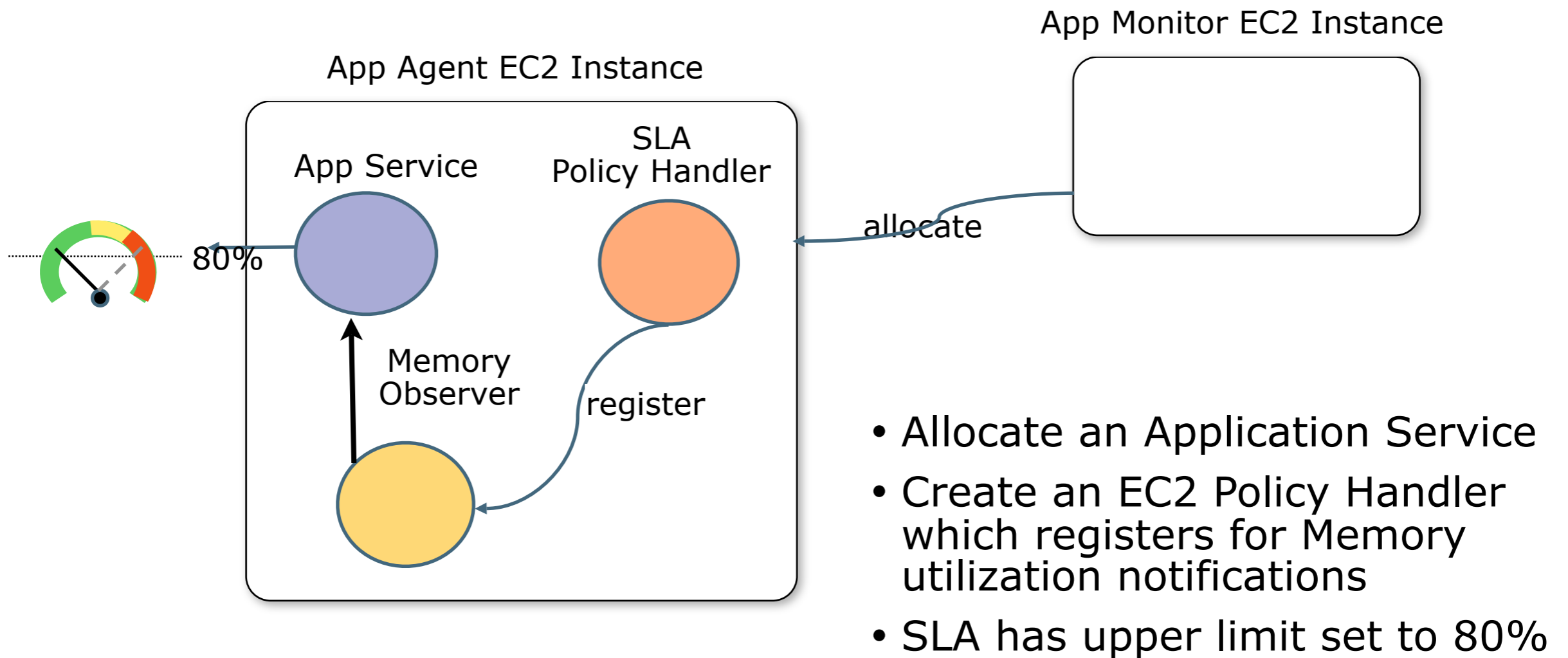
# Elastic Grid Scalability on EC2

- Across existing EC2 instances



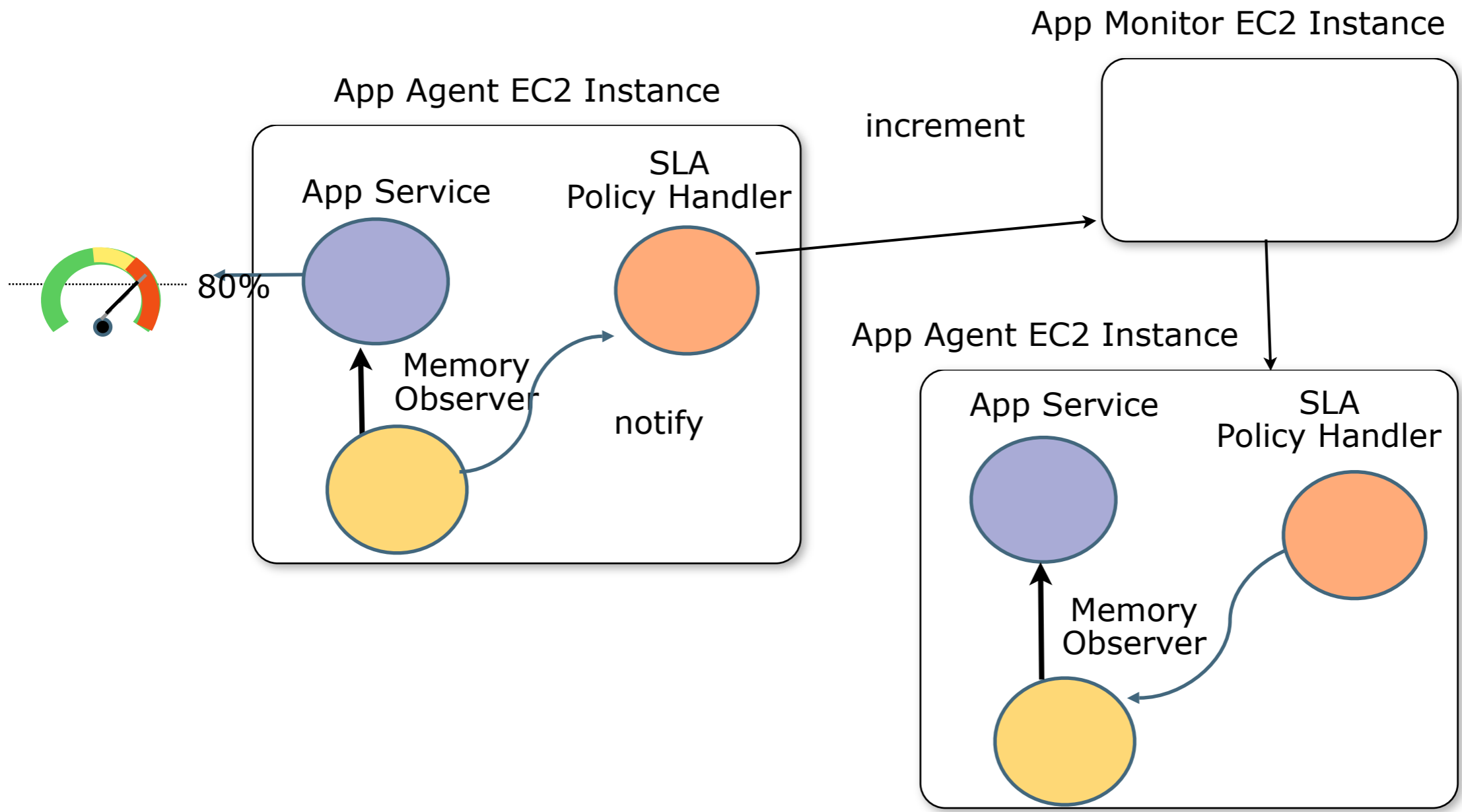
# Elastic Grid Scalability on EC2

- New EC2 instance

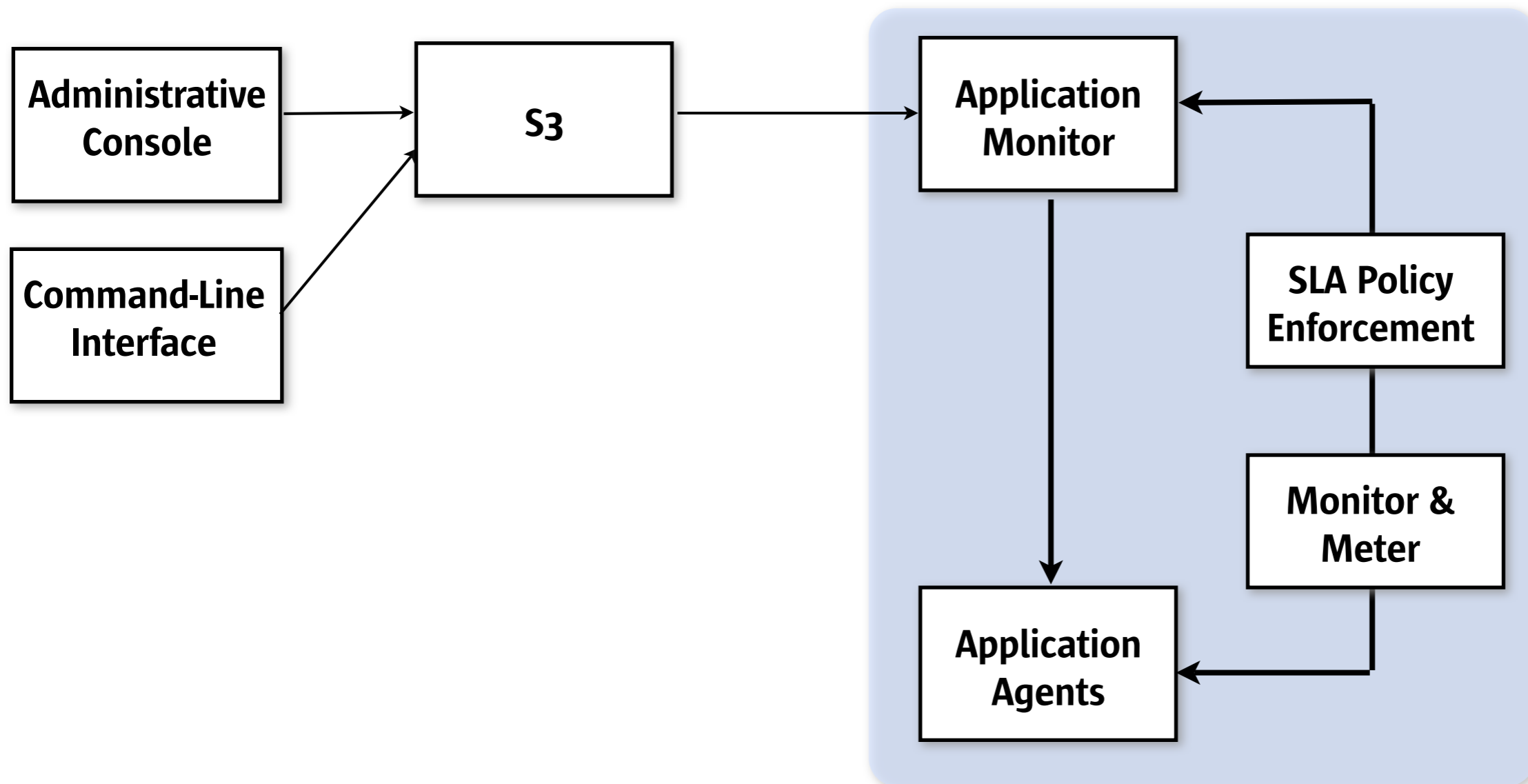


# Elastic Grid Scalability on EC2

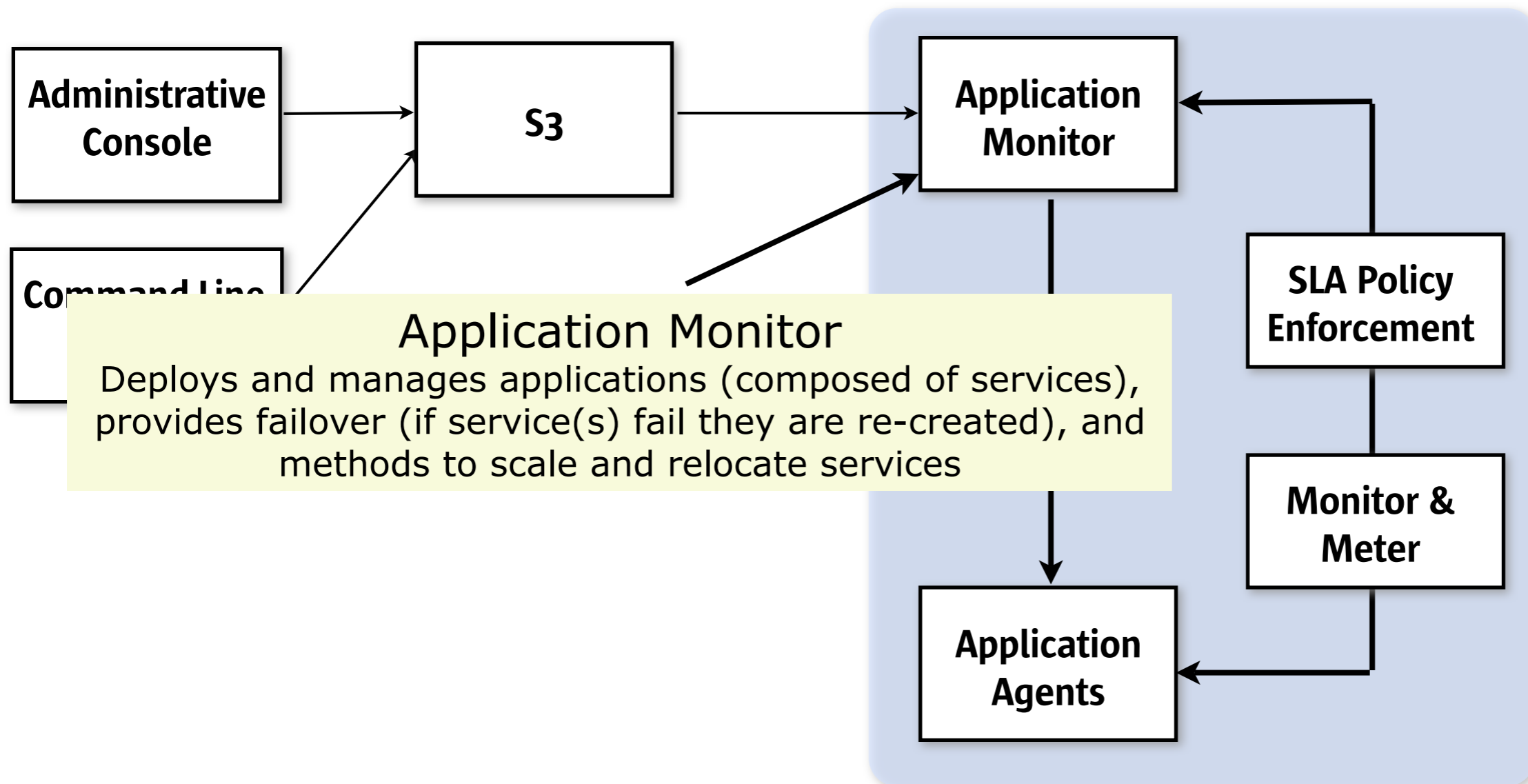
- New EC2 instance



# Elastic Grid Architecture

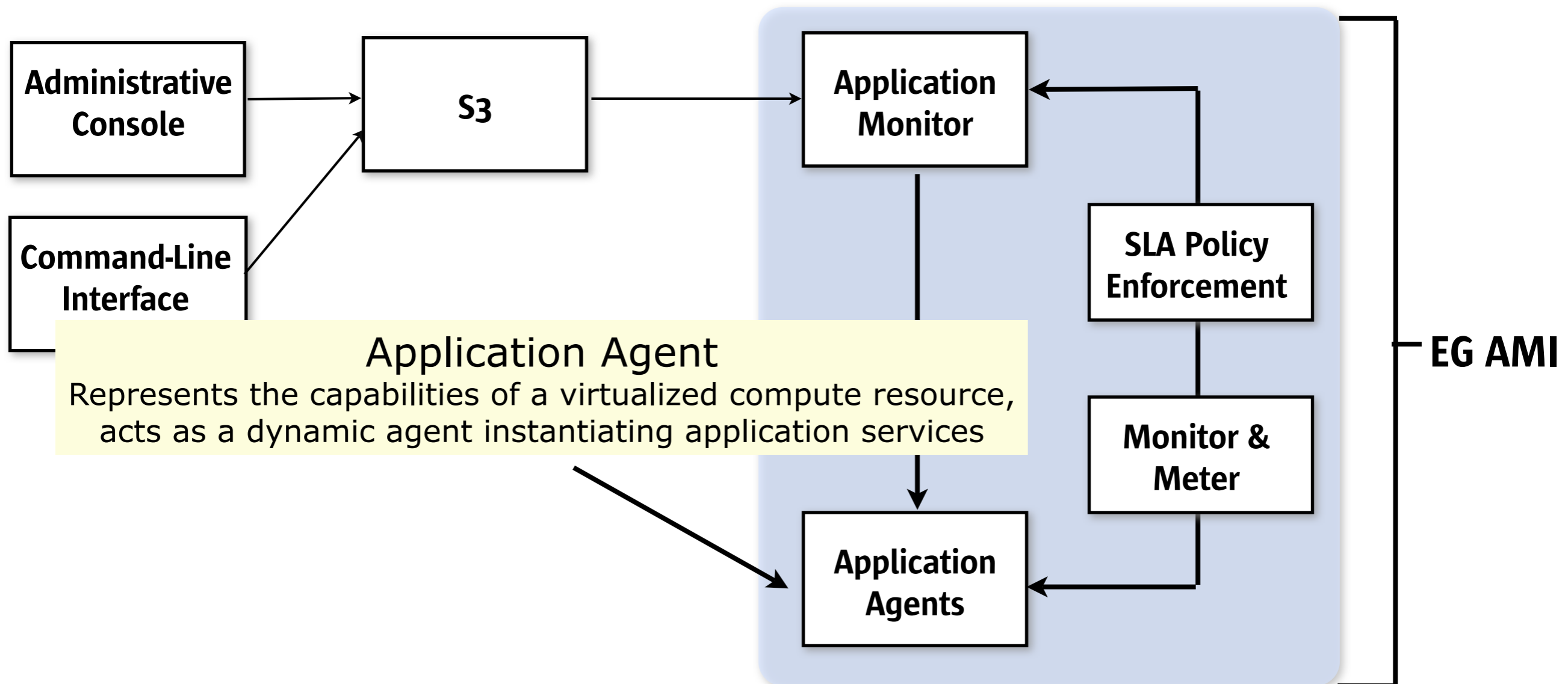


# Elastic Grid Architecture

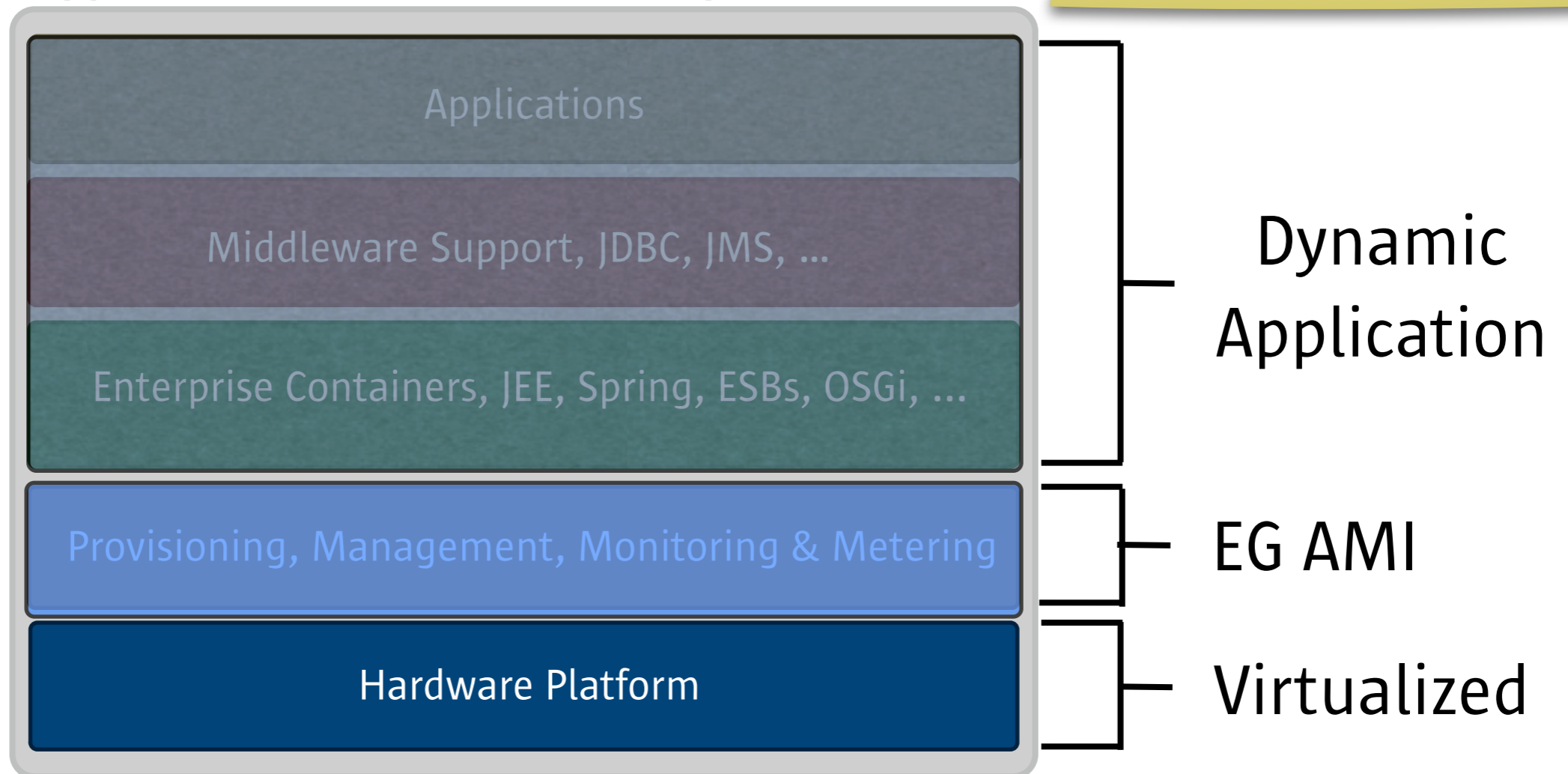




# Elastic Grid Architecture



▶ Typical Architecture Taxonomy



You won't have to create AMIs anymore and of course, you won't have to update an AMI when your application code changes!

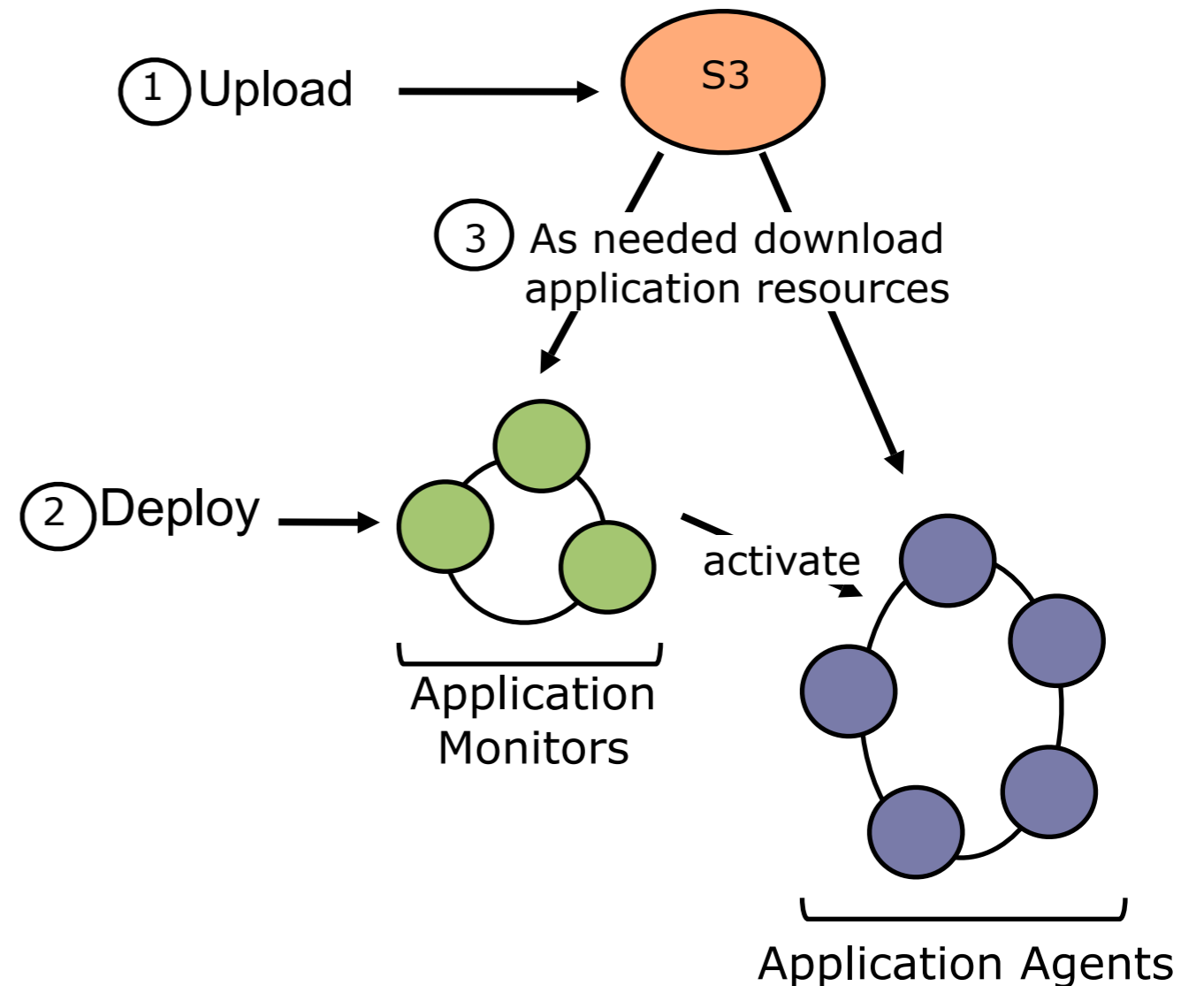
# Deployments with Elastic Grid

- EG AMIs are pre-set, no need to (re-)bundle
- As application code changes, upload to S3 and deploy
- Focuses on developer productivity



# Elastic Grid Deployments

- Deploy application code to S3
- Run the deploy command
- All code is dynamically served and instantiated
- Application is monitored and managed across EC2 instances





## Sample Deployments

# Sample: Calculator

- Calculator
  - Made of 4 core services (add, subtract, divide, multiply)
  - Consumed by a Calculator, with the help of dynamic associations

```
deployment(name: 'Calculator') {
  groups 'rio'

  resources id: 'impl.jars', 'calculator/lib/calculator.jar'
  resources id: 'client.jars', 'calculator/lib/calculator-dl.jar'

  service(name: 'Calculator') {
    interfaces {
      classes 'calculator.Calculator'
      resources ref: 'client.jars'
    }
    implementation(class: 'calculator.service.CalculatorImpl') {
      resources ref: 'impl.jars'
    }
    associations {
      association name: 'Add', type: 'requires', property: 'add'
      association name: 'Subtract', type: 'requires', property: 'subtract'
      association name: 'Multiply', type: 'requires', property: 'multiply'
      association name: 'Divide', type: 'requires', property: 'divide'
    }
    maintain 1
  }
}

['Add', 'Subtract', 'Multiply', 'Divide'].each { s ->
  service(name: s) {
    interfaces {
      classes "calculator.$s"
      resources ref: 'client.jars'
    }
    implementation(class: "calculator.service.${s}Impl") {
      resources ref: 'impl.jars'
    }
    maintain 1
  }
}
```

## Sample: Spring dm Server

- Indicate where Spring dm Server distribution is located
- Indicate where are located the sample applications

```
software(name: 'Spring DM Server', version: '1.0.1', removeOnDestroy: true) {
  install source: 'http://elastic-grid-examples.s3.amazonaws.com/spring-dm/springsource-dm-server-1.0.1.RELEASE.zip',
    target: 'springsource-dm-server',
    unarchive: true
}

// deploy Spring Travel Sample
data source: 'http://elastic-grid-examples.s3.amazonaws.com/spring-dm/spring-travel-1.2.0.zip', unarchive: true,
  target: 'springsource-dm-server/springsource-dm-server-1.0.1.RELEASE'

// deploy Form Tags Sample
data source: 'http://elastic-grid-examples.s3.amazonaws.com/spring-dm/formtags-1.4.0.zip', unarchive: true,
  target: 'springsource-dm-server/springsource-dm-server-1.0.1.RELEASE'
```

## Sample: Spring dm Server

- Hook into JMX Gauges in order to get the metrics
- Indicate how you'd like to automatically scale

```
// monitor number of threads and scale Spring dm Server instances
sla(id: 'thread-count', low: 80, high: 200) {
  policy type: 'scaling', max: 3
  monitor name: 'Thread Count',
           objectName: ManagementFactory.THREAD_MXBEAN_NAME,
           attribute: 'ThreadCount', period: 5000
}

// monitor maximum time spent on a HTTP call
sla(id: 'http-max-time', high: 5000) {
  policy type: 'notify'
  monitor name: 'HTTP Max Time',
           objectName: 'springsource.server.catalina:type=GlobalRequestProcessor,name=http-8080',
           attribute: 'maxTime', period: 1000
}

execute command: 'bin/startup.sh'
maintain 1 // we want at least 1 instance to be running in our cluster
maxPerMachine 1 // we don't want to host more than 1 instance per machine
```



# Summary

- So what does Elastic Grid do for the app?
  - ▶ Ease development and deployment of Java applications using Amazon services
  - ▶ Provides automated management, fault detection and scalability for the application
  - ▶ To be available soon: Cloud Bursting!
- Why you should use Elastic Grid?
  - ▶ Avoid Cloud Computing platforms pitfalls
  - ▶ Focus on development, not infrastructure





Thanks for your attention!

Register for the private Beta on  
Elastic Grid Website

Elastic Grid Website: <http://www.elastic-grid.com>

Elastic Grid Blog: <http://blog.elastic-grid.com>

Elastic Grid Wiki: <http://wiki.elastic-grid.com>