

DataCenter optimization for Cloud Computing



Benjamín Barán
National University of Asuncion (UNA)
bbaran@pol.una.py
Paraguay



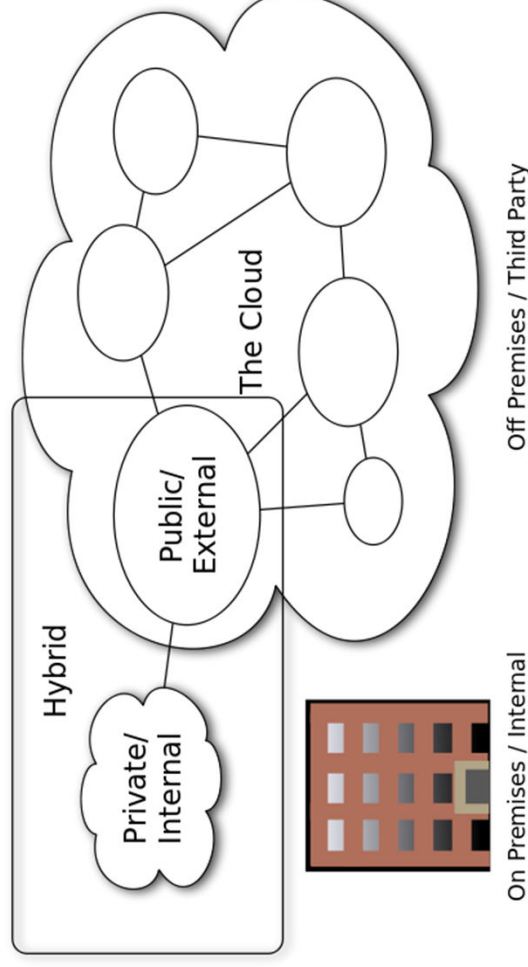
Content

- ▶ Cloud Computing
- ▶ Commercial Offerings
- ▶ Basic Problem Formulation
- ▶ Open Research
- ▶ Conclusions



Cloud Computing

Cloud computing is an *Internet - based* computing in which large groups of remote servers are networked to allow sharing of data processing tasks, centralized data storage and on-line access to computer services or resources.



Cloud Computing Types

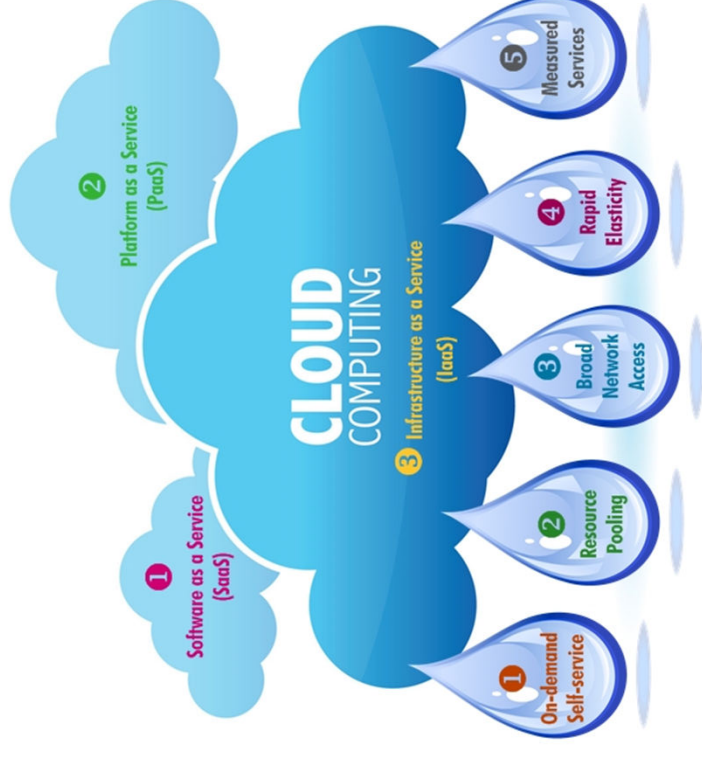
- 1- Public Cloud
- 2- Private Cloud
- 3- Hybrid Cloud

[http://en.wikipedia.org/wiki/Cloud_computing]

NIST definition of Cloud Computing

Cloud computing is a model for enabling ubiquitous, convenient, on-demand network access to a shared pool of configurable computing resources (as, networks, servers, storage, applications, and services) that can be rapidly provisioned and released with minimal management effort or service provider interaction.

This *cloud model* is composed of 5 essential characteristics, and 3 service models.



National Institute of Standards
and Technology (NIST)

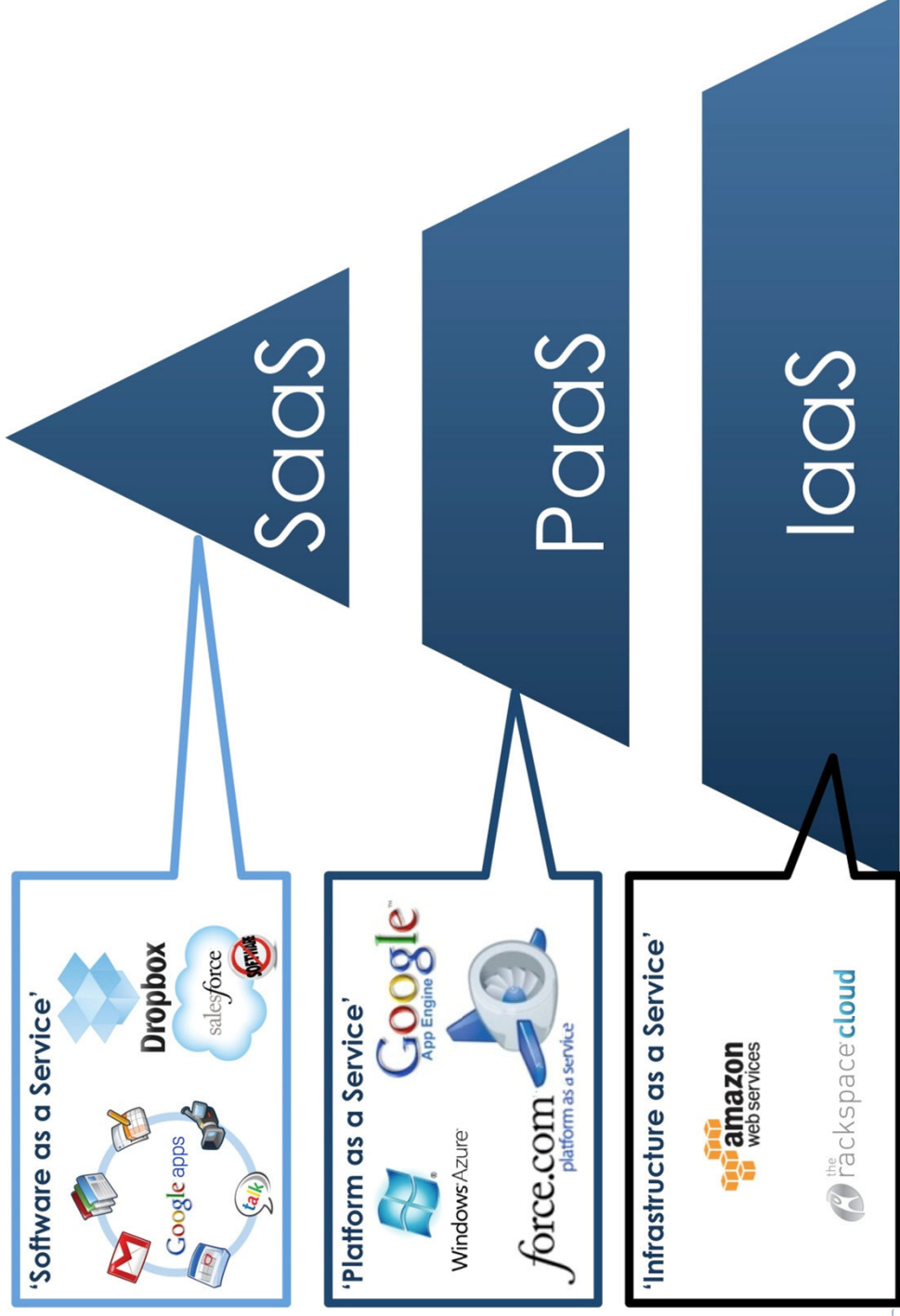
Cloud Computing

- ▶ The very definition of *cloud computing* still remains controversial.
- ▶ There are alternative definition as the following one: *Cloud Computing is the dynamic provisioning of IT capabilities (hardware, software, or services) from third parties over a network.*
- ▶ Cloud computing is a *computing model*, not a technology. In this model of computing, all elements (processing, storage, etc.) **related to DataCenters** are made available to end users via the *Internet*.
- ▶ **Virtualization** - as well as the cloud computing model within which it often runs - answers much of **DataCenters** needs.

[<http://www.computerworld.com/article/2527305/cloud-computing/cloud-computing-definitions-and-solutions.html>]

NIST Service Models

National Institute of Standards and Technology



Everything/Anything as a Service - XaaS

- **BPaaS** - Business Process as a Service
- **CaaS** - Communication as a Service
- **DaaS** - Data as a Service
- **IaaS** - Infrastructure as a Service
- **ITaaS** - IT (*Information Technology*) as a Service
- **PaaS** - Platform as a Service
- **RaaS** – Resources as a Service
- **SaaS** - Software as a Service
- **SECaaS** - SECurity as a Service

Infrastructure as a Service - IaaS

Infrastructure as a Service – IaaS, provides grids or clusters or virtualized servers, networks, storage and systems software designed to augment or replace the functions of an entire DataCenter.

The highest-profile example is *Amazon's Elastic Compute Cloud [EC2]* and *Simple Storage Service [S3]*, but other traditional IT vendors are also offering services.

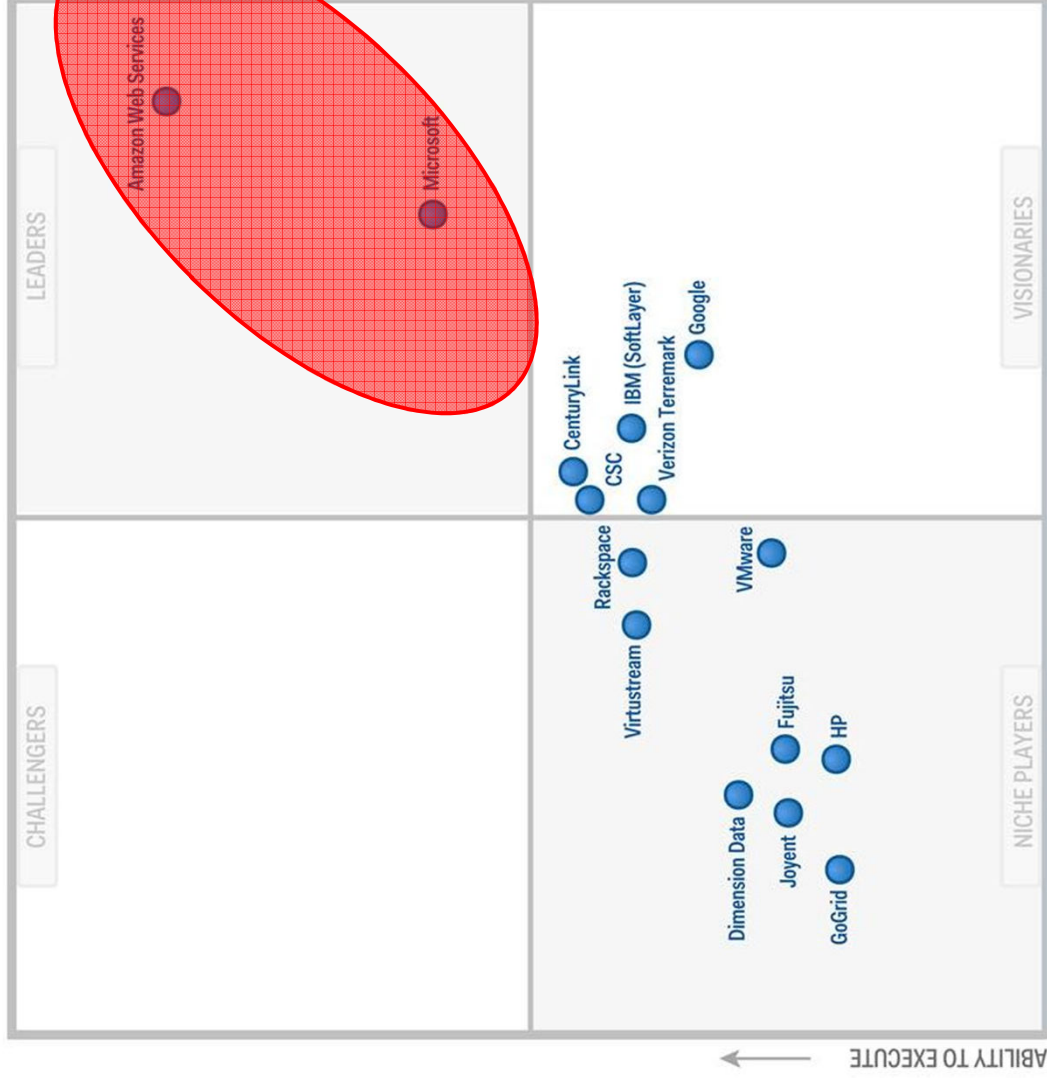


Microsoft

GOGRID



IaaS Gartner Magic Quadrant



Microsoft

[http://aws.amazon.com/resources/gartner-mq-2014-learn-more/?sc_icountry=en&sc_ichannel=ha&sc_idetail=ha_en_42&sc_icontent=ha_en_d_ed_42_1&sc_iplace=ha_en_ed&sc_icampaign=ha_en_Gartner&trk=/]

As of May 2014

COMPLETENESS OF VISION →
← ABILITY TO EXECUTE

Source: Gartner (May 2014)

AWS – Amazon Web Services



Amazon EC2 »

Servicio web que ofrece capacidad informática escalable en la nube.



Amazon S3 »

Almacenamiento de datos de alta escalabilidad y fiabilidad, y de baja latencia.



Amazon RDS »

Bases de datos MySQL, Oracle y SQL Server gestionadas.



Amazon SQS »

Cola escalable para almacenar los mensajes que van de un ordenador a otro.



Amazon CloudWatch »

Supervisión de recursos y aplicaciones en la nube de AWS.



AWS Data Pipeline »

Orquestación de flujos de trabajo basados en los datos.



Amazon DynamoDB »

Servicio de bases de datos NoSQL totalmente gestionado con una escalabilidad óptima.



Amazon CloudFront »

Un servicio web de entrega de contenido.



Amazon EBS »

Volumenes de almacenamiento predecibles y de alta disponibilidad y fiabilidad.



Amazon ELB »

Servicio web que ofrece escalabilidad y alta disponibilidad.



Amazon ElastiCache »

Caché ampliable gestionada.



Amazon SWF »

Servicio de flujos de trabajo para desarrollar aplicaciones escalables y adaptables.



Amazon SES »

Servicio de envío de correo electrónico rentable en la nube.



Amazon SNS »

Servicio web para establecer, ejecutar y enviar notificaciones desde la nube.



Amazon Elastic Transcoder »

Convierta sus archivos multimedia fácilmente, con un bajo coste y a escala.



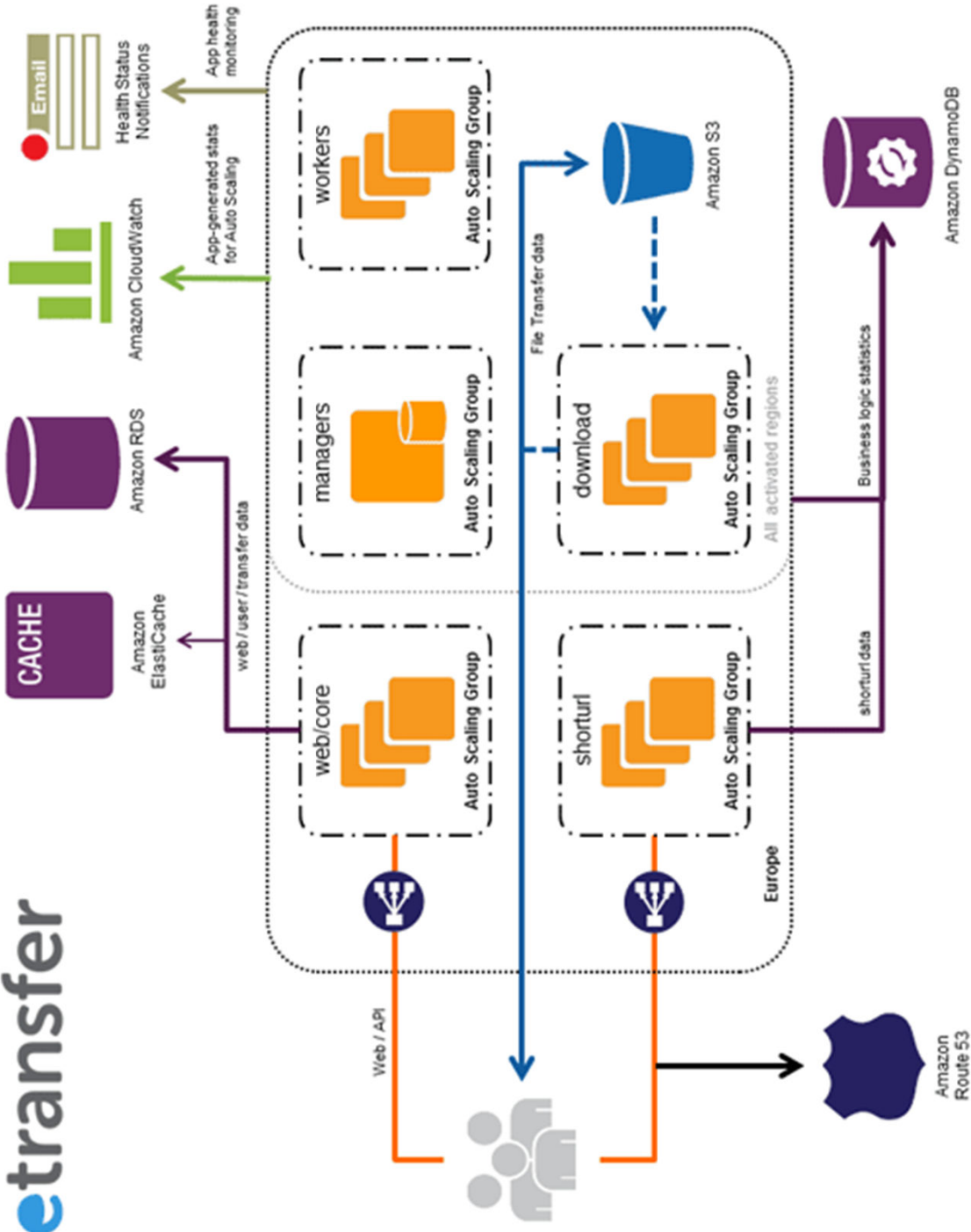
AWS Marketplace »

Software de socio preconfigurado para ejecutarse en AWS.



Case Study: using AWS

wetransfer



Companies using *Public Cloud Computing*



Instagram



Dropbox

[http://spectrum.ieee.org/computing/networks/escape-from-the-data-center-the-promise-of-peerto-peer-cloud-computing/?utm_source=techalert&utm_medium=email&utm_campaign=092514]

Cost Models

- ▶ **Static:** fixed prices (resource prices rarely change in time, as traditional Amazon EC2)
- ▶ **Dynamic Prices.** Resource prices fluctuates on demand on a day or weekly basis (e.g., weekend prices are different).
- ▶ **Spot Prices.** It is based on user's bids.
If user bid met or exceed the current spot price, he gains access to requested resources (as new Amazon EC2).





1 year Prices Example

INSTANCE	CPU	ECU	RAM [GiB]	Storage [GB]	Price per hour
t2.micro	1	Variable	1	EBS	\$0.013
t2.small	1	Variable	2	EBS	\$0.026
t2.medium	2	Variable	4	EBS	\$0.052
m3.medium	1	3	3.75	1 x 4 SSD	\$0.070
m3.large	2	6.5	7.5	1 x 32 SSD	\$0.140
m3.xlarge	4	13	15	2 x 40 SSD	\$0.280
m3.2xlarge	8	26	30	2 x 80 SSD	\$0.560

ECU ... EC2 Computing Unit (e.g. 1 ECU = 1.0-1.2 GHz 2007 Xeon)

EBS ... Elastic Block Storage (\$0.10 per GB-month)

SSD ... Solid State Drive, internal storage



Prices Example

aws marketplace

AMI: Amazon Machine Images



SUSE Linux Enterprise Server 11 (64-bit)

Version 3P3 | Sold by Amazon Web Services

\$0.03 to \$5.67/hr incl EC2 charges + other AWS usage fees

Amazon EC2 running SUSE Linux Enterprise Server is a proven platform for development, test, and production workloads. With more than 6,000 certified applications from over ...

Linux/Unix, SUSE Enterprise Server 11 SP3 | 64-bit Amazon Machine Image (AMI)



Microsoft Windows Server 2012 RTM

Version 2014.04.09 | Sold by Amazon Web Services

\$0.018 to \$9.348/hr incl EC2 charges + other AWS usage fees

Amazon EC2 running Microsoft Windows Server is a fast and dependable environment for deploying applications using the Microsoft Web Platform. Amazon EC2 enables you to run ...

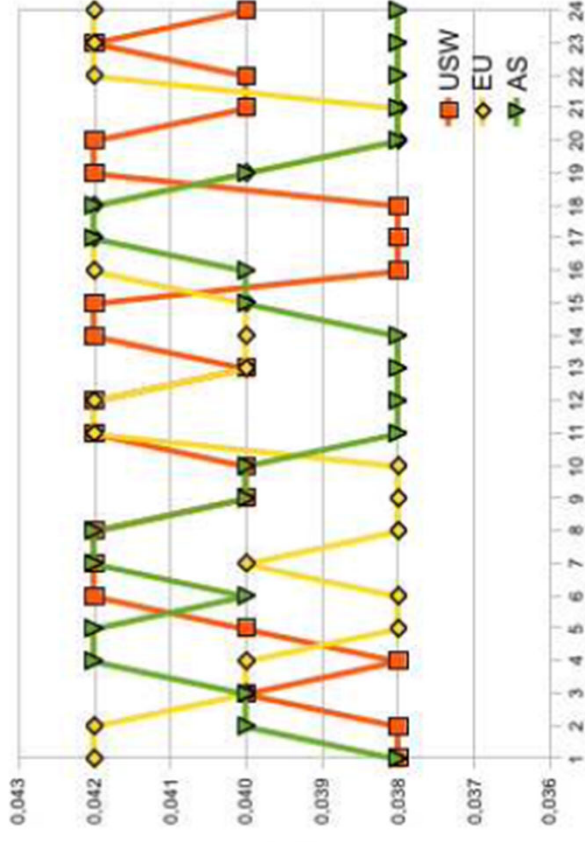
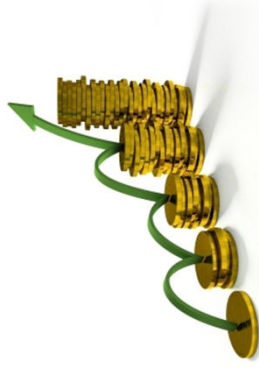
Windows, Windows Server 2012 6.2 | 64-bit Amazon Machine Image (AMI)

https://aws.amazon.com/marketplace/search/results/ref=mkt_ste_free_tier_ec2?page=1&restriction=%28or+asin%3A%27B00AA27RK4%27+asin%3A%27B00A6KUVBW%27+asin%3A%27B007ORS58%27+asin%3A%27B00AAEFK8%27+asin%3A%27B007O0H35O%27+asin%3A%27B00635Y2IW%27+asin%3A%27B007Z5YWX4%27%29



Spot Price example

INSTANCE	LINUX	WINDOWS
m1.small	\$0.0071 per Hour	\$0.0171 per Hour
m1.medium	\$0.0081 per Hour	\$0.0331 per Hour
m1.large	\$0.0161 per Hour	\$0.0661 per Hour
m1.xlarge	\$0.0352 per Hour	\$0.1321 per Hour

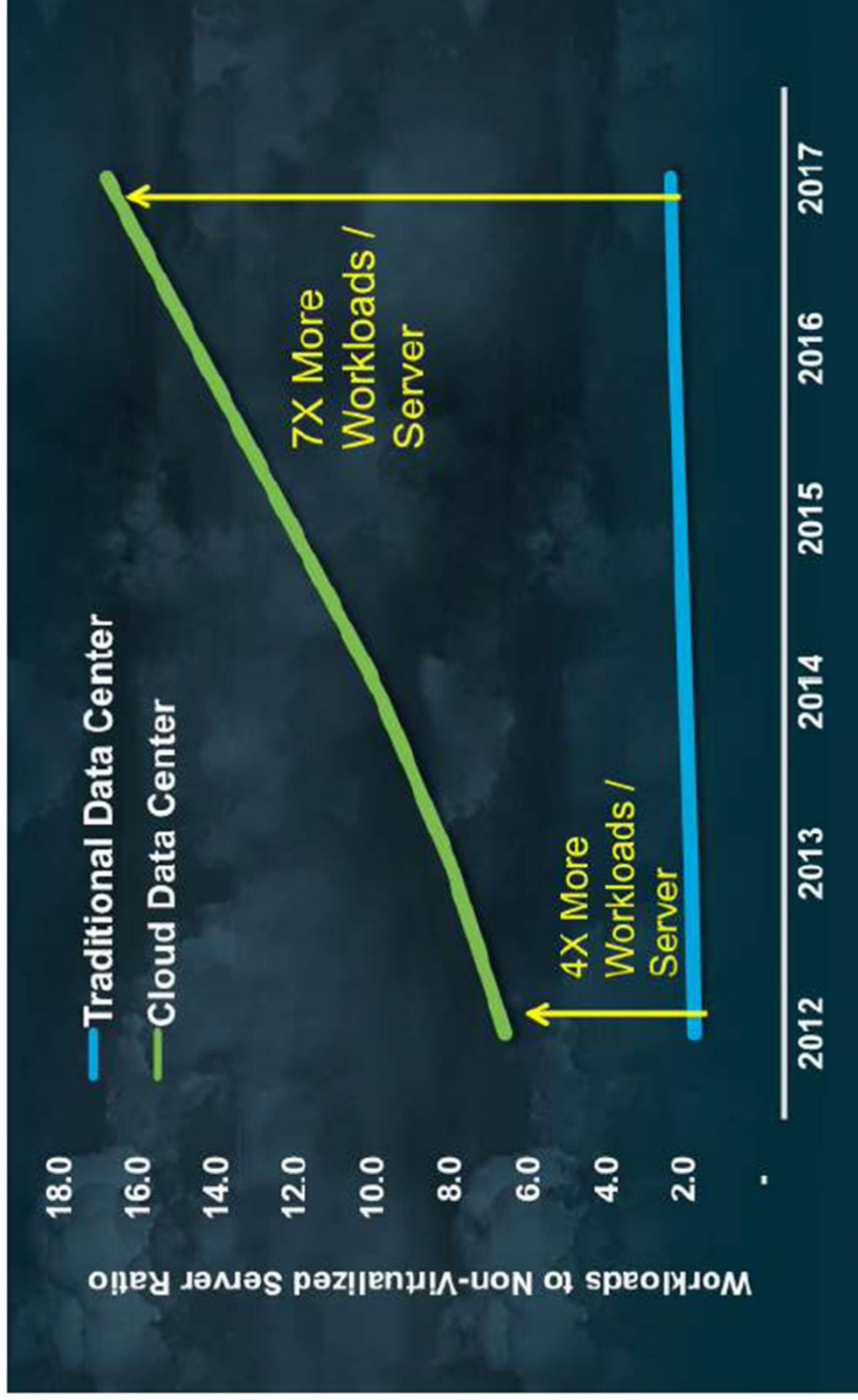


See **TUTORIALS** at:

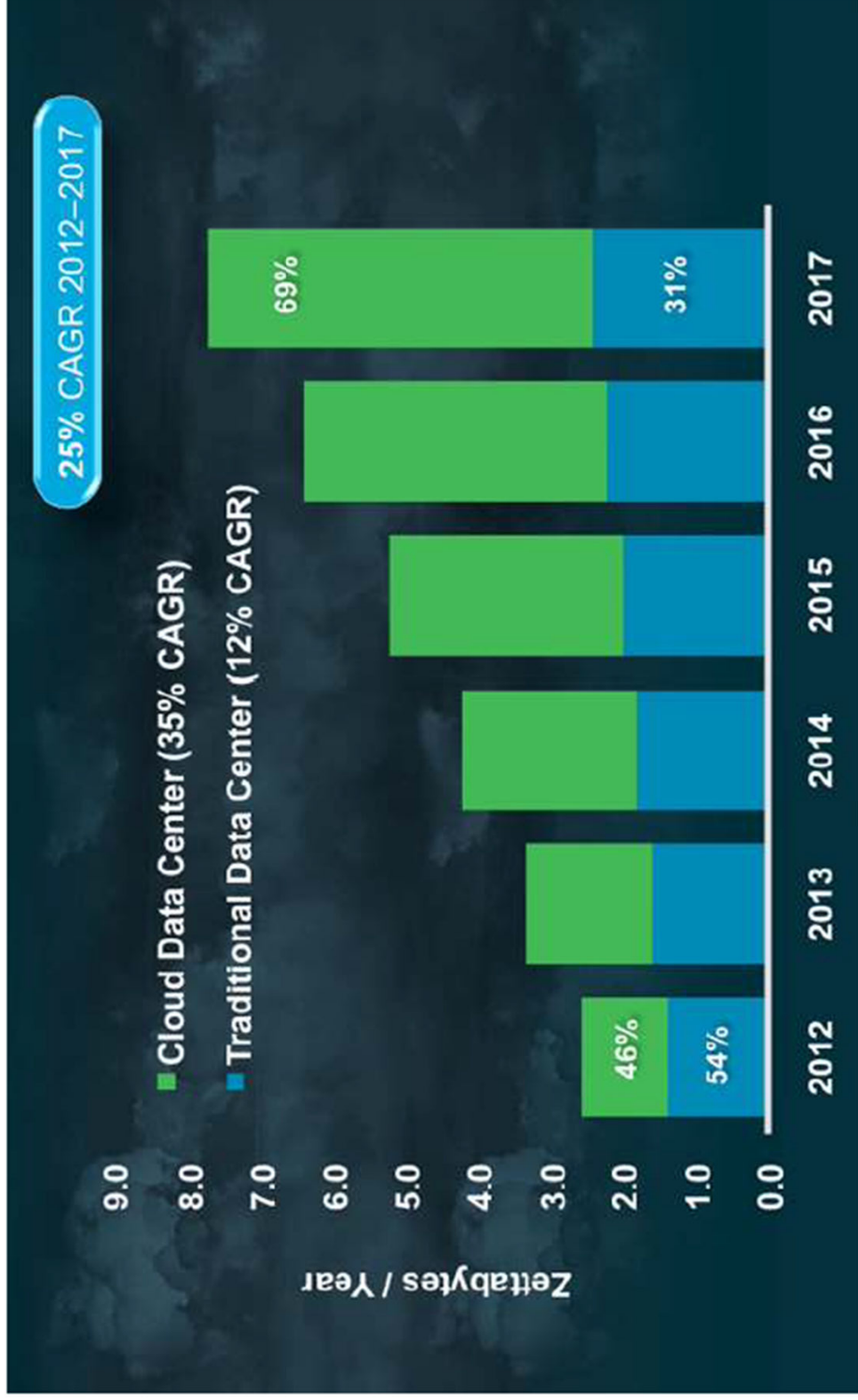
[<http://aws.amazon.com/ec2/purchasing-options/spot-instances/>]

Cloud Computing Trend

Figure 6. Increasing Cloud Virtualization



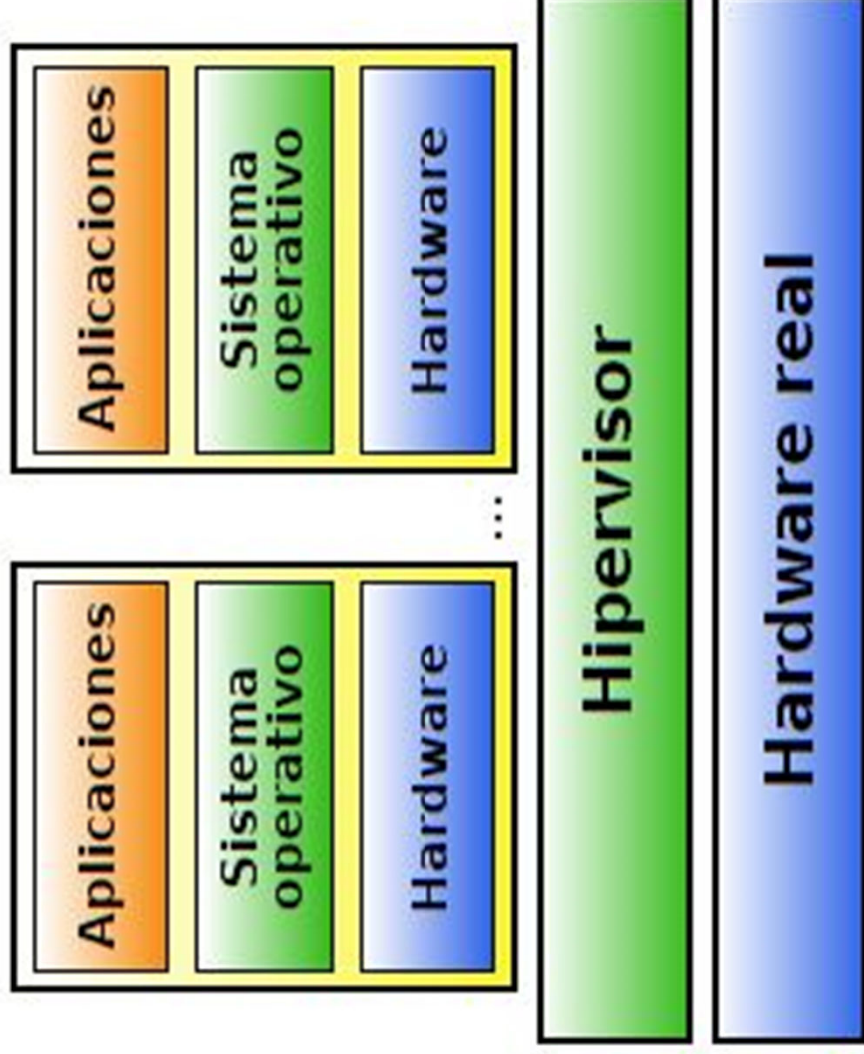
Cloud Computing Trend



1 ZB = 10²¹ B

CAGR ... Compound Annual Growth Rate

Virtualization

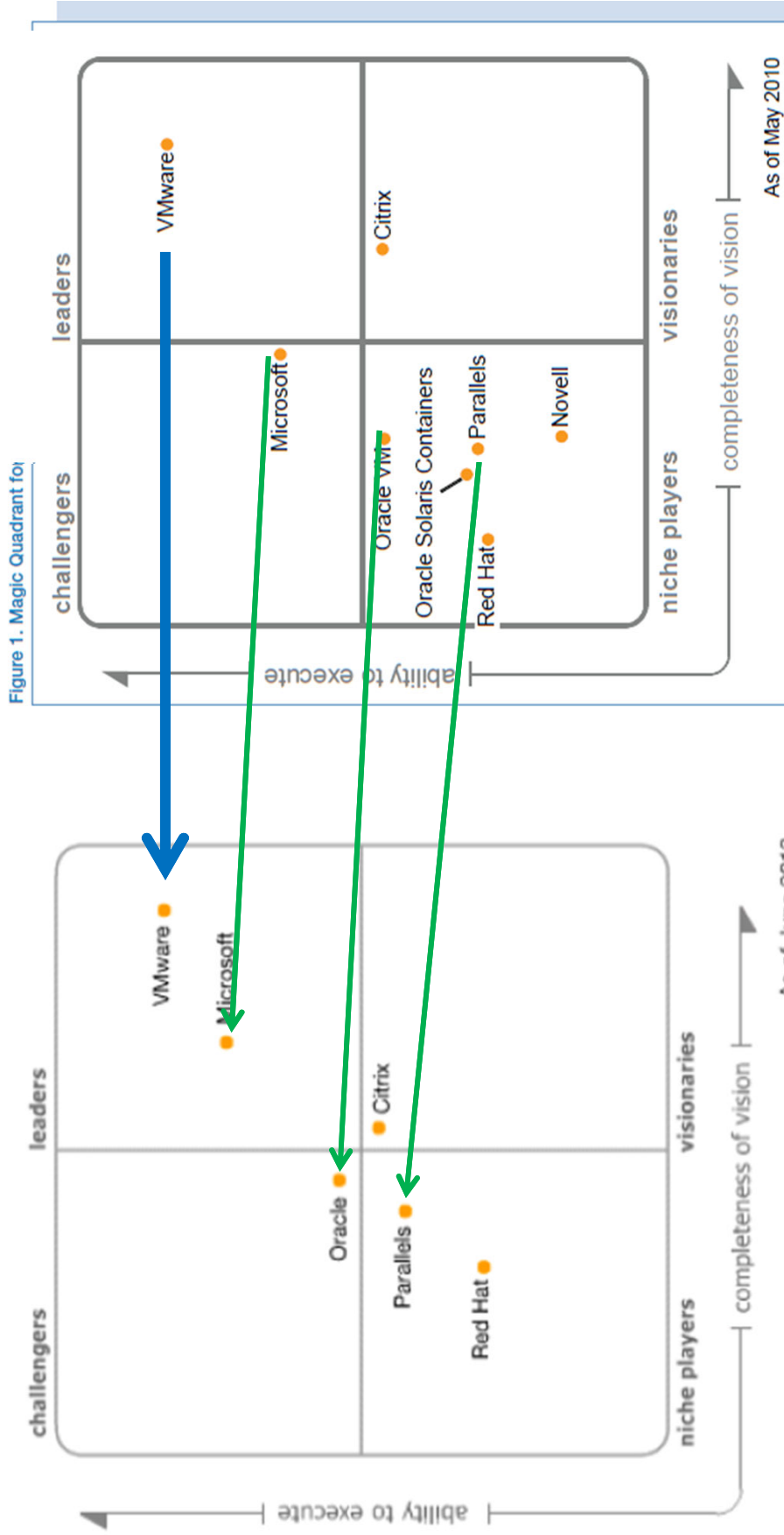


Virtualization



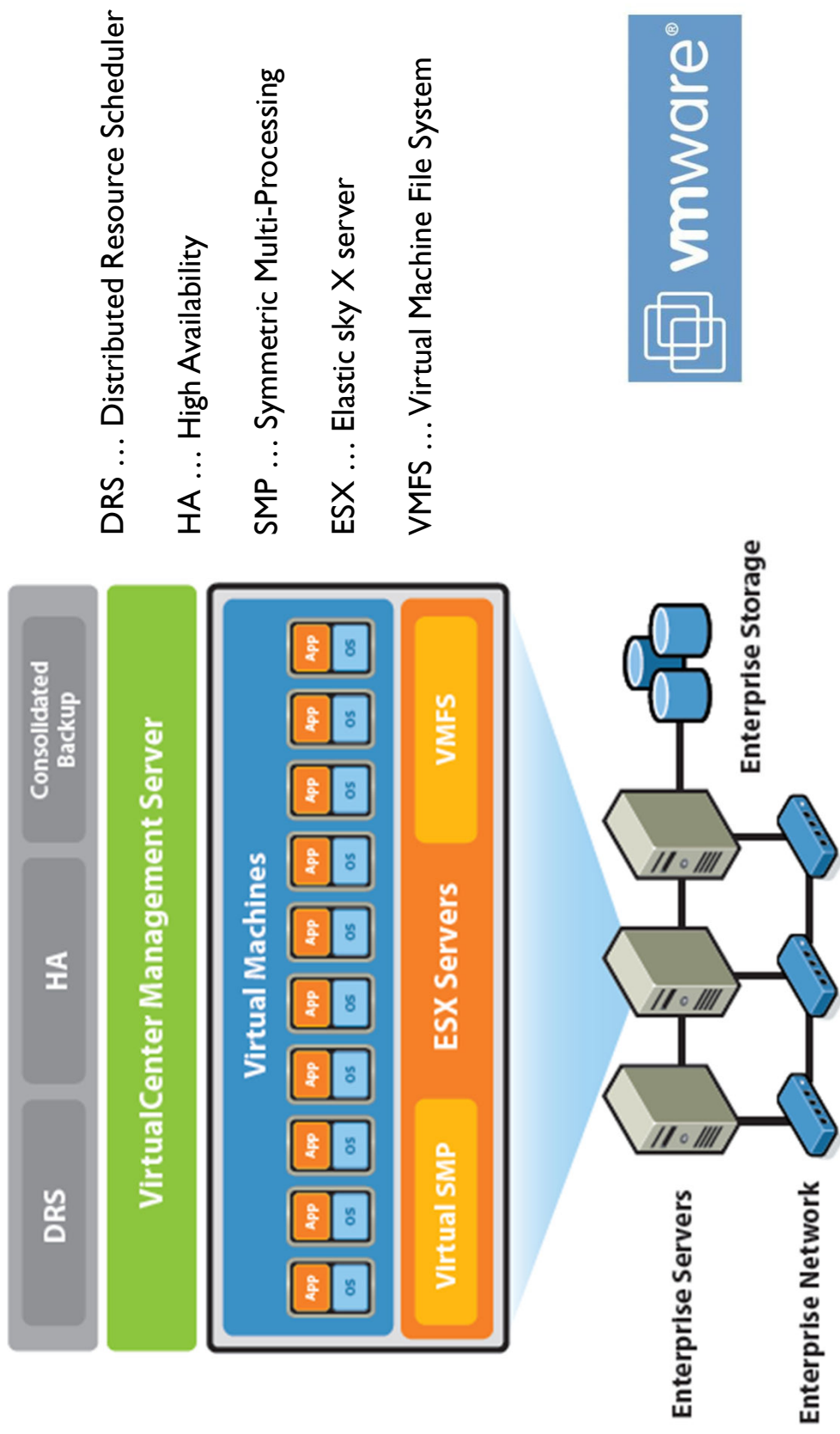
Magic Quadrant

Figure 1. Magic Quadrant for x86 Server Virtualization Infrastructure



Source: Gartner (June 2013)

Virtualization example: VMware

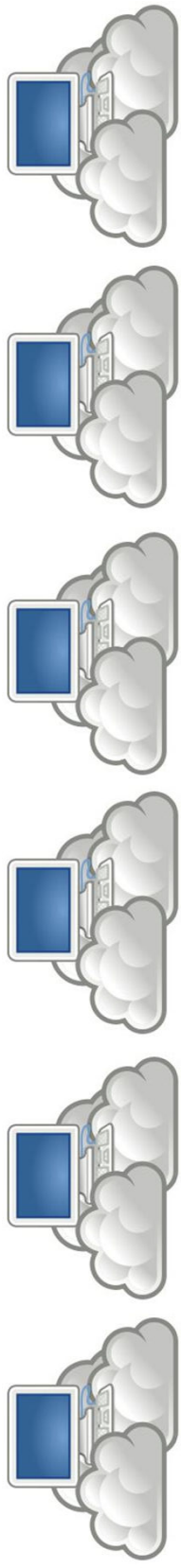


Basic Problem Formulation

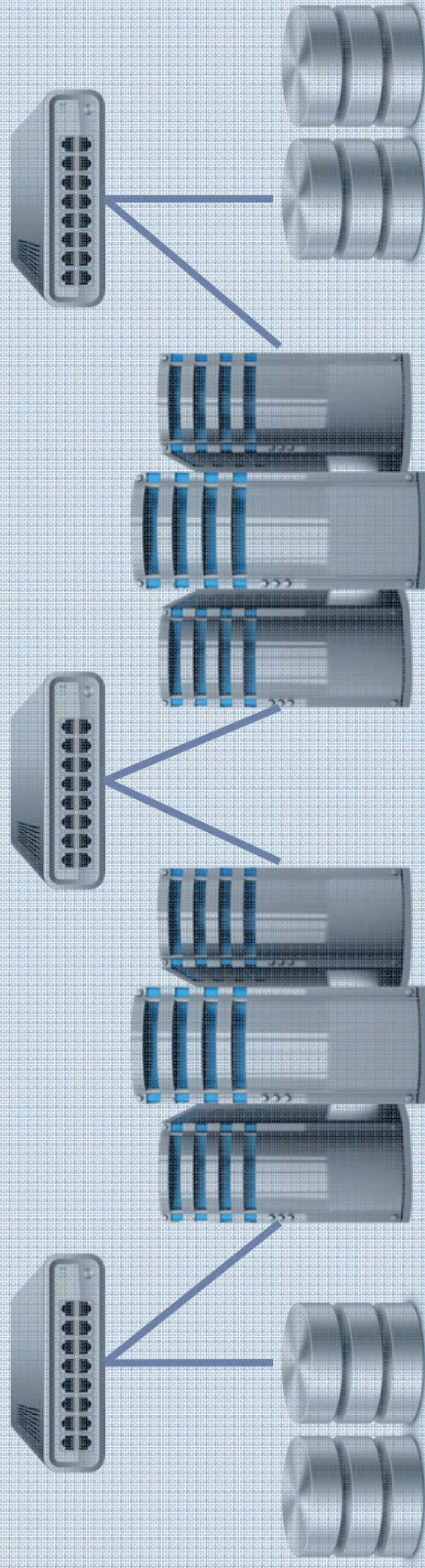
Virtual Machine Placement

Which virtual machines should be located at each physical machine?

Under which criteria?



Virtual Infrastructure

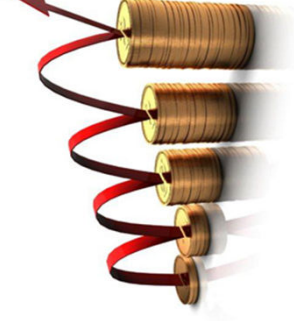
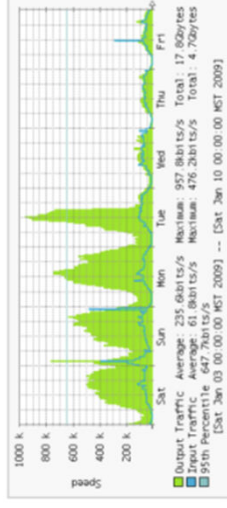


Objective Functions

- ▶ Main objective functions [3]

[F. López Pires, B. Barán, “Taxonomy of Optimal Virtual Machine Placement in Efficient Datacenters,” IEEE Aranducon’ 2012]

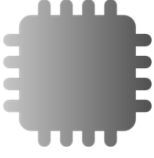



- (1) Energy Consumption Minimization
- (2) Economical Revenue Maximization
- (3) Network Traffic Minimization



- ▶ Mathematical formulation without SLA [4]

[F. López Pires, B. Barán, “Multi-Objective Virtual Machine Placement with Service Level Agreement,” 6th IEEE/ACM International Conference on Utility and Cloud Computing, UCC’2013. Dresden – Alemania]

Physical Resources Matrix

$$H = \begin{pmatrix} H_{cpu_1} & H_{ram_1} & H_{hdd_1} & \rho_{max_1} \\ H_{cpu_2} & H_{ram_2} & H_{hdd_2} & \rho_{max_2} \\ \vdots & \vdots & \vdots & \vdots \\ H_{cpu_n} & H_{ram_n} & H_{hdd_n} & \rho_{max_n} \end{pmatrix} \begin{matrix} i = 1 \\ i = 2 \\ \vdots \\ i = n \end{matrix}$$





where:

n : Number of physical machines

H_i : Virtual machine with identification i

H_{cpu_i} : Processing resource of the physical machine H_i in [MIPS]

H_{ram_i} : RAM memory resource of the physical machine H_i in [MB]

H_{hdd_i} : Storage resource of the physical machine H_i in [GB]

ρ_{max_i} : Maximum power consumption of the physical machine H_i in [W]

Virtual Requirement Matrix

$$V = \begin{cases} V_{cpu_1} & V_{ram_1} & V_{hdd_1} & r_1 & SLA_1 & j = 1 \\ V_{cpu_2} & V_{ram_2} & V_{hdd_2} & r_2 & SLA_2 & j = 2 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ V_{cpu_m} & V_{ram_m} & V_{hdd_m} & r_m & SLA_m & j = m \end{cases}$$

where:

m : Number of virtual machines

V_j : Virtual machine with identification j

V_{cpu_j} : Processing requirement of the virtual machine V_j in [MIPS]

V_{ram_j} : RAM memory requirement of the virtual machine V_j in [MB]

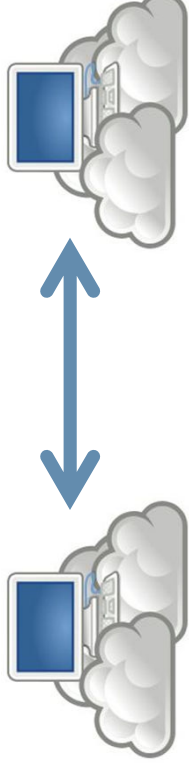
V_{hdd_j} : Storage requirement of the virtual machine V_j in [GB]

r_j : Economical revenue for placement of virtual machine V_j in [\$]

SLA_j : Service level agreement of virtual machine V_j

Network Traffic Matrix

$$T = \begin{matrix} & k = 1 & k = 2 & \dots & k = m \\ \begin{matrix} j = 1 \\ j = 2 \\ \vdots \\ j = m \end{matrix} & \begin{pmatrix} T_{11} & T_{12} & \dots & T_{1m} \\ T_{21} & T_{22} & \dots & T_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ T_{m1} & T_{m2} & \dots & T_{mm} \end{pmatrix} \end{matrix}$$



where:

m : Number of virtual machines

V_j : Virtual machine with identification j

V_k : Virtual machine with identification k

T_{jk} : Network Communication rate between V_j and V_k in [Kbps]

Basic Problem Formulation

$$H = \{4 \times n\}$$

$$V = \{5 \times m\}$$

$$T = \{m \times m\}$$

INPUT

$$P = \{m \times n\}$$

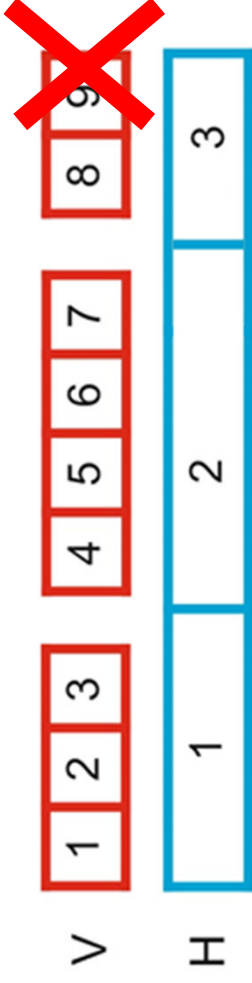
OUTPUT

$P = \{P_{ji}\}$ where $P_{ji} \in \{0, 1\}$

$(P_{ji} = 0)$ indicates that V_j IS NOT located in H_i

$(P_{ji} = 1)$ indicates that V_j IS located in H_i ($P_{ji} : V_j \rightarrow H_i$)

Placement Matrix



$$P = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 1 & 1 & 1 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

- $P_{11} : V_1 \rightarrow H_1$
- $P_{21} : V_2 \rightarrow H_1$
- $P_{31} : V_3 \rightarrow H_1$
- $P_{42} : V_4 \rightarrow H_2$
- $P_{52} : V_5 \rightarrow H_2$
- $P_{62} : V_6 \rightarrow H_2$
- $P_{72} : V_7 \rightarrow H_2$
- $P_{83} : V_8 \rightarrow H_3$
- V_9 is not assigned

Constraints

- ▶ Unique placement of virtual machines

$$\sum_{i=1}^n P_{ji} \leq 1 \quad \forall j \in \{1, 2, \dots, m\}$$

Constraint I

where:

n : Number of physical machines

P_{ji} : Binary variable equals 1 if the virtual machine V_j is located to run on the physical machine H_i ; 0 otherwise

m : Number of virtual machines

Constraints

- ▶ Service Level Agreement (SLA) provision

$$\sum_{i=1}^n P_{ji} = 1 \quad \forall j \text{ such that } SLA_j = 1$$

Constraint 2

where:

n : Number of physical machines

P_{ji} : Binary variable equals 1 if the virtual machine V_j is located to run on the physical machine H_i ; 0 otherwise

SLA_j : Service Level Agreement $SLA_j = 1$ if V_j is critical, or 0 otherwise

Constraints

- ▶ Resource capacity of physical machines

$$\sum_{j=1}^m Vcpu_j \times P_{ji} \leq Hcpu_i$$

Constraint 3

$$\sum_{j=1}^m Vram_j \times P_{ji} \leq Hram_i$$

Constraint 4

$$\sum_{j=1}^m Vhdd_j \times P_{ji} \leq Hhdd_i$$

Constraint 5

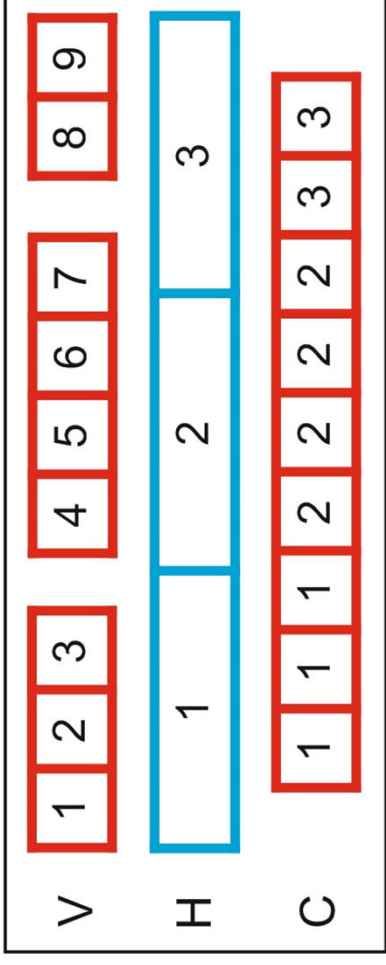
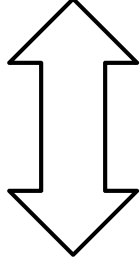
where:

- $Vcpu_j$: Processing requirement [MIPS] of virtual machine V_j
- $Vram_j$: RAM memory requirement [MB] of virtual machine V_j
- $Vhdd_j$: Storage requirement [GB] of virtual machine V_j

Multi-Objective Memetic Algorithm

- ▶ Chromosome representation

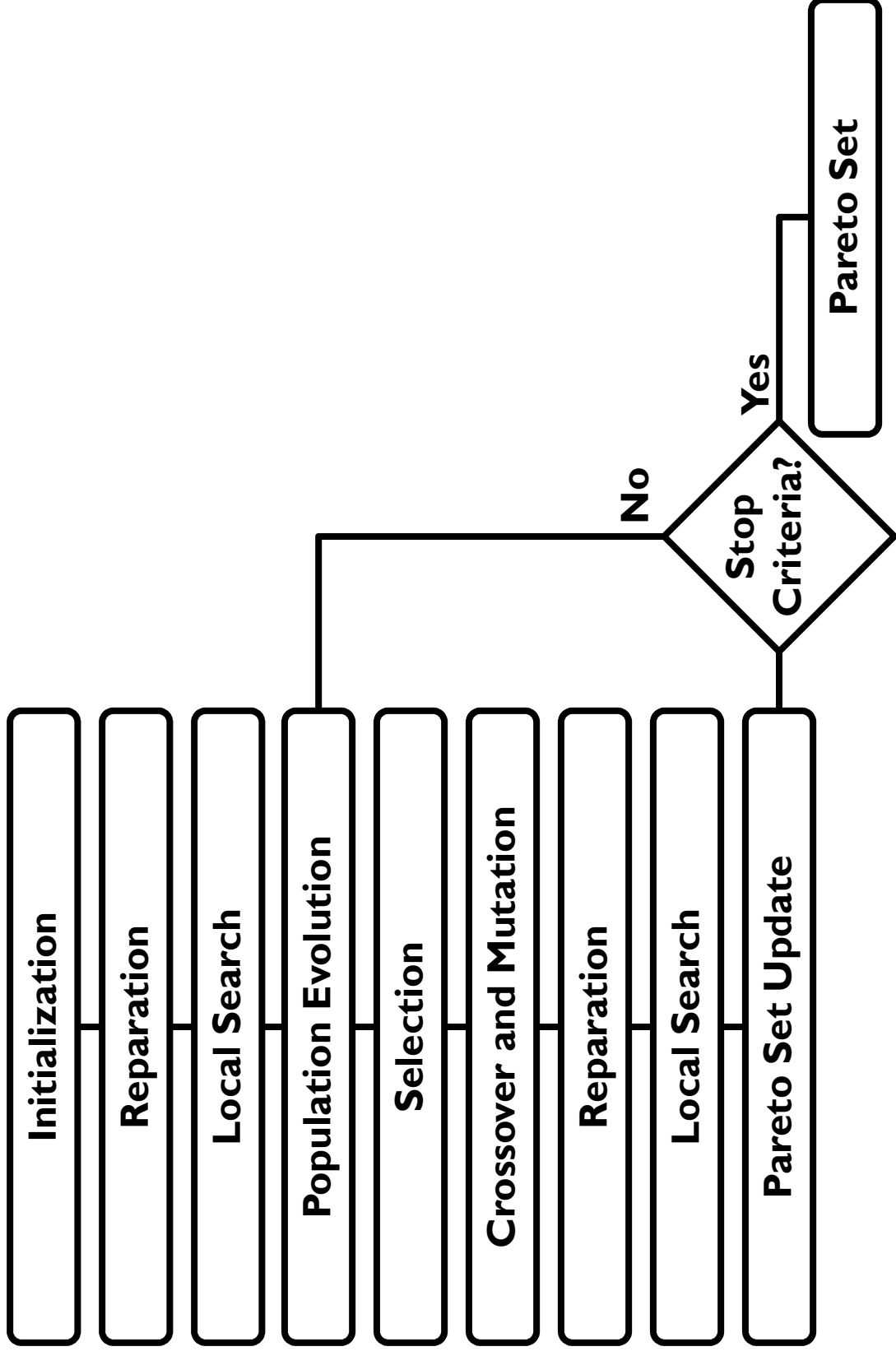
$$\text{Solution } P = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 1 \\ 0 & 0 & 1 & 0 & 1 \end{bmatrix}$$



Proposed Formulation

Proposed Chromosome Representation

Multi-Objective Memetic Algorithm



Experimental Results

- ▶ Testing Environment
 - ▶ Algorithms in ANSI C (GNU C)
 - ▶ GNU/Linux Ubuntu 11.10 Operating System
 - ▶ Intel Core i7 de 1.2 GHz Processor
 - ▶ 8 GB of RAM Memory

- ▶ Real Input Data



Parque Tecnológico
Itaipu



Experimental Results

▶ Experimental Test I:

Scenario	Number of Physical Machines	Number of Virtual Machines	Critical SLA Percentage	Number of Y_{known} Elements	Number of X_{known} Elements
10x20	10	20	50%	48	48

- ▶ Exhaustive search algorithm can not complete calculation in useful time.
- ▶ It is necessary to implement alternatives to exhaustive search.

Experimental Results

- ▶ Experimental Test 2:

Scenario	Number of Physical Machines	Number of Virtual Machines	Critical SLA Percentage
3x5	3	5	0, 10, 20, 30, 40, 50, 60, 70, 80, 90, 100%
4x10	4	10	0, 10, 20, 30, 40, 50, 60, 70, 80, 90, 100%

- ▶ Relation of variables:
 - ▶ Execution Time and Critical SLA Percentage
 - ▶ Number of Solutions and Critical SLA Percentage

Future Work

- ▶ **Alternative formulations for the problem:**
 - ▶ Considering more SLA levels and constrains (as geographical)
 - ▶ Considering more SLA metrics: *response time, jitter, etc.*
- ▶ **Formulation with other objective functions (more than 80 different objective functions were found in the specialized literature).**
- ▶ Testing other bio-inspired meta-heuristic, given the novelty of the proposed context.
- ▶ **Pure Dynamical Context and its uncertainty.**
- ▶ Use of a third-party Broker.
- ▶ **Consider Hybrid clouds.**
- ▶ Case studies and commercial applications.



Thanks!



Universidad Nacional de Asunción

Benjamín Barán
National University of Asuncion (UNA)
bbaran@pol.una.py
Paraguay

