

Uma abordagem para a seleção automática de soluções de recuperação de desastres baseada em modelos estocásticos

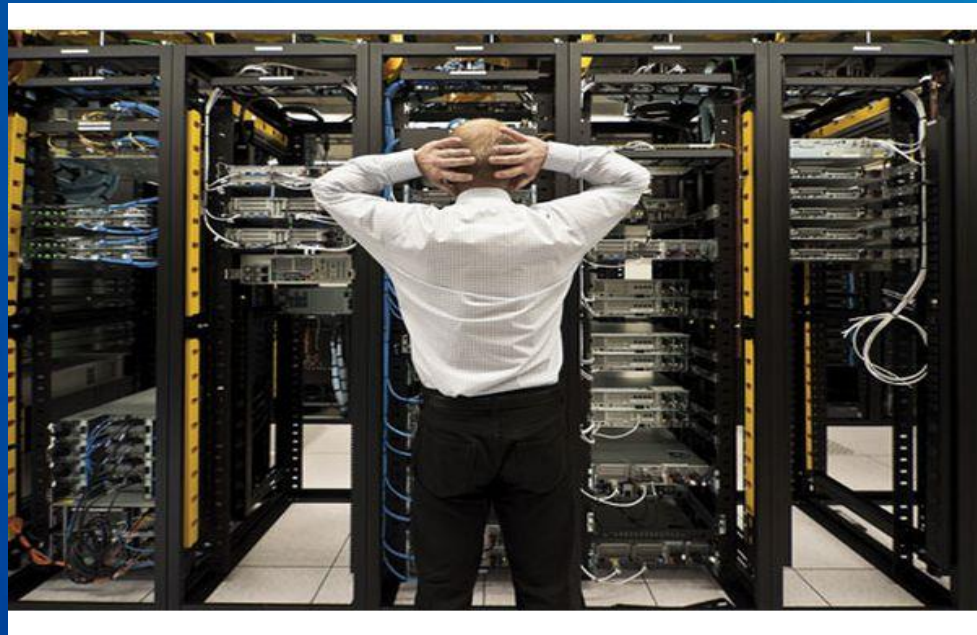
Júlio Mendonça
jrmn@cin.ufpe.br

Orientador: Prof. Dr. Ricardo Massa
Co-orientador: Prof. Dr. Ermeson Andrade

Agenda

- Motivação
- Fundamentação
- Status do Trabalho
- Próximos passos
- Referências

Motivação





Menu



Contact Sales

Products ▾

Solutions

Pricing

Getting Started

More ▾

English ▾

My Account ▾

Create an AWS Account

Summary of the Amazon S3 Service Disruption in the Northern Virginia (US-EAST-1) Region

We'd like to give you some additional information about the service disruption that occurred in the Northern Virginia (US-EAST-1) Region on the morning of February 28th, 2017. The Amazon Simple Storage Service (S3) team was debugging an issue causing the S3 billing system to progress more slowly than expected. At 9:37AM PST, an authorized S3 team member using an established playbook executed a command which was intended to remove a small number of servers for one of the S3 subsystems that is used by the S3 billing process. Unfortunately, one of the inputs to the command was entered incorrectly and a larger set of servers was removed than intended. The servers that were inadvertently removed supported two other S3 subsystems. One of these subsystems, the index subsystem, manages the metadata and location information of all S3 objects in the region. This subsystem is necessary to serve all GET, LIST, PUT, and DELETE requests. The second subsystem, the placement subsystem, manages allocation of new storage and requires the index subsystem to be functioning properly to correctly operate. The placement subsystem is used during PUT requests to allocate storage for new objects. Removing a significant portion of the capacity caused each of these systems to require a full restart. While these subsystems were being restarted, S3 was unable to service requests. Other AWS services in the US-EAST-1 Region that rely on S3 for storage, including the S3 console, Amazon Elastic Compute Cloud (EC2) new instance launches, Amazon Elastic Block Store (EBS) volumes (when data was needed from a S3 snapshot), and AWS Lambda were also impacted while the S3 APIs were unavailable.

S3 subsystems are designed to support the removal or failure of significant capacity with little or no customer impact. We build our systems with the assumption that things will occasionally fail, and we rely on the ability to remove and replace capacity as one of our core operational processes. While this is an operation that we have relied on to maintain our systems since the launch of S3, we have not completely restarted the index subsystem or the placement subsystem in our larger regions for many years. S3 has experienced massive growth over the last several years and the process of restarting these services and running the necessary safety checks to validate the integrity of the metadata took longer than expected. The index subsystem was the first of the two affected subsystems that needed to be restarted. By 12:26PM PST, the index subsystem had activated enough capacity to begin servicing S3 GET, LIST, and DELETE requests. By 1:18PM PST, the index subsystem was fully recovered and GET, LIST, and DELETE APIs were functioning normally. The S3 PUT API also required the placement subsystem. The placement subsystem began recovery when the index subsystem was functional and finished recovery at 1:54PM PST. At this point, S3 was operating normally. Other AWS services that were impacted by this event began recovering. Some of these services had accumulated a backlog of work during the S3 disruption and required additional time to fully recover.

<https://aws.amazon.com/pt/message/41926/>



2/11 RCA - Storage – Service Management Operations

Summary of impact: Between 11:40 and 16:48 UTC on 02 Nov 2017, a subset of customers may have experienced issues with Service Management functions (Create, Update, Delete, GetAccountProperties etc.) for their Azure Storage resources. Storage customers may have been unable to provision new Storage resources or perform service management operations on existing resources. Other services with dependencies on Storage may have also experienced impact such as Virtual Machines, Cloud Services, Event Hubs, Backup, Azure Site Recovery, Azure Search and VSTS Load Testing. The impact for this issue was limited to Service Management functions. Service Availability for existing resources would not have been affected. Engineers received alerts and investigated the issue. The issue was understood due to high disk space utilization triggered by staging of OS updates under unexpected circumstance, which resulted in impacting the checkpoint processes of Storage Resource Provider (SRP) services. The incident was mitigated by removing unexpected additional staged images to allow the checkpoint processes of SRP to be succeeded in time.

Root cause and mitigation: As part of the scheduled monthly Guest OS update for the Azure Infrastructure, OS images were staged to each scale unit. During the October OS update cycle, Engineers detected the initial scheduled OS build had a .NET application compatibility issue after the build was staged to a subset of scale units across production. The staging of this build was subsequently stopped but the image was not removed from the staging list. This resulted in an additional size of images being staged to each node in the scale unit, resulting in a low disk space situation on some nodes. Storage Resource Provider (SRP) service handles storage account management operations for all regions (create/update/delete/list). SRP is a Paxos based service and uses disk for state checkpointing. Due to an increase in disk space utilization triggered by staging of OS updates above, the checkpointing process failed, which is critical for service operation. To mitigate the issue, engineers freed the required disk space, allowing state checkpointing to succeed and resume normal processing for all service management operation request. Updates for this incident were regularly communicated via <https://status.azure.com>. Due to a process error, a subset of customers relying on service health alerts from Azure Monitor or using the Azure Service Health experience in the Azure Management portal would not have been notified. The Azure service communications team sincerely apologizes for any inconvenience this may have caused. Customers are encouraged to enable Azure Monitor service health alerts for improved incident notifications. Details are available at the following link: <http://aka.ms/azurealerts>.

<https://azure.microsoft.com/en-gb/status/history/>

Por que adotar uma solução de recuperação de desastres?

The Cost of Downtime



More than half of companies (54%) report that they have experienced a downtime event that lasted **more than 8 hours** in the past five years.



In the event of a site outage, **67%** estimate that their business would lose **\$20K+** for every day of downtime.



Would lose up to \$20k



Would lose between \$20k - \$100k



Would lose between \$100k - \$500k

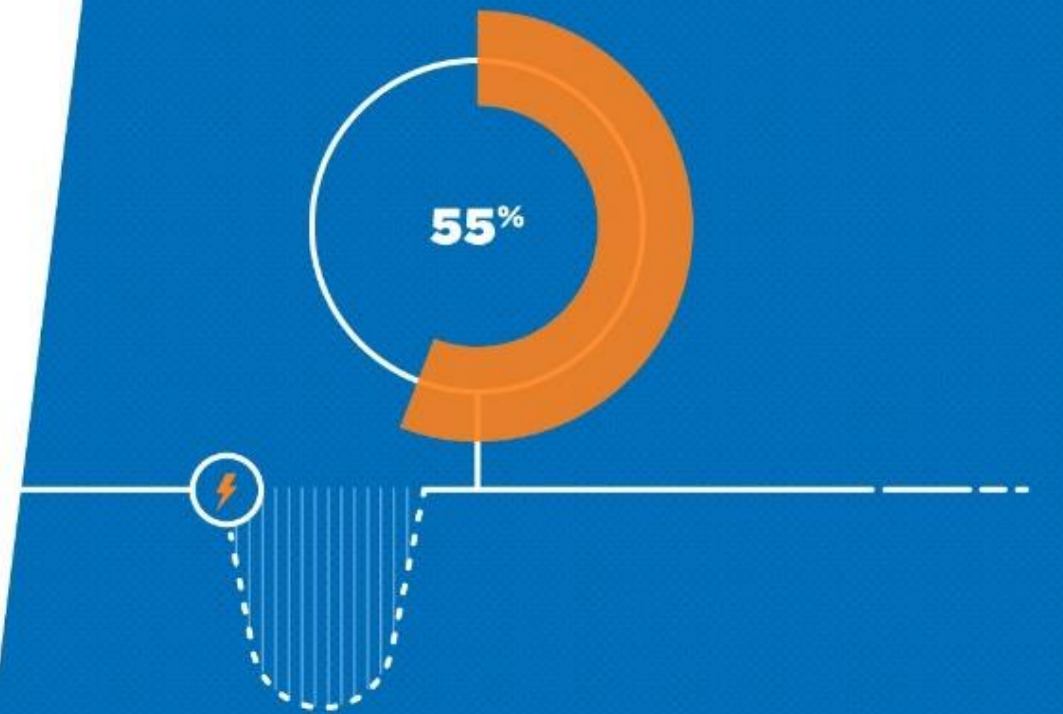


Would lose more than \$500k

Downtime Events Highlight Weaknesses in DR Plans

55% of those who have experienced a **downtime event** report **making changes to their DR plans** after the event occurred. It's even greater (**65%**) for those experiencing **outages greater than 8 hours**.

Changed DR Strategy



<https://www.zetta.net/resource/state-disaster-recovery-2016>

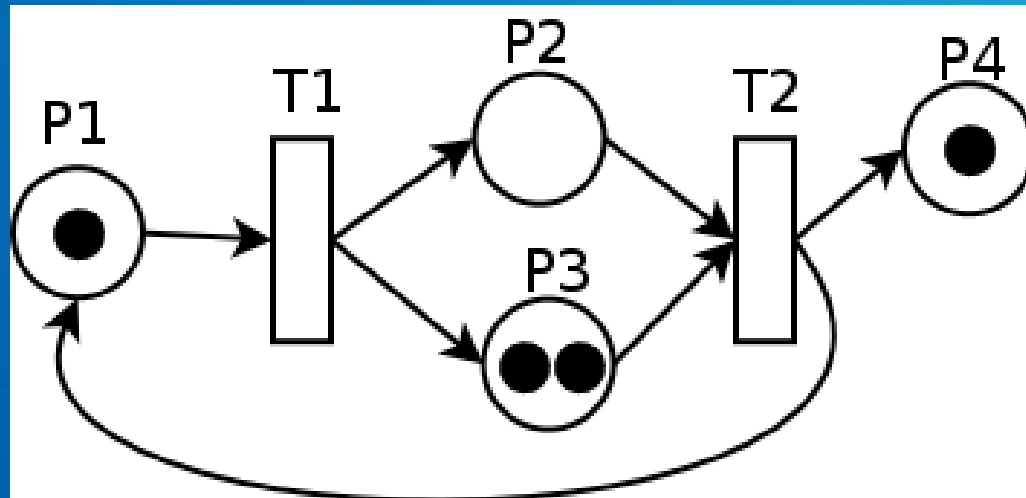
Por que não adotar simplesmente um DRaaS ou BaaS?

- Existem várias estratégias para recuperação de desastres:
 - Clouds;
 - Mecanismos de Backup;
 - Replicações;
 - Deduplicações, entre vários outros.
- Essas estratégias ainda possuem uma infinidade de combinações possíveis;
- A adoção dessas medidas depende da necessidade e condições de uma empresa.

“Criar uma **abordagem** que permita a **avaliação e escolha da melhor solução** de recuperação de desastres de forma **automática, eficiente e confiável.**”

Fundamentação

- Stochastic Petri Nets

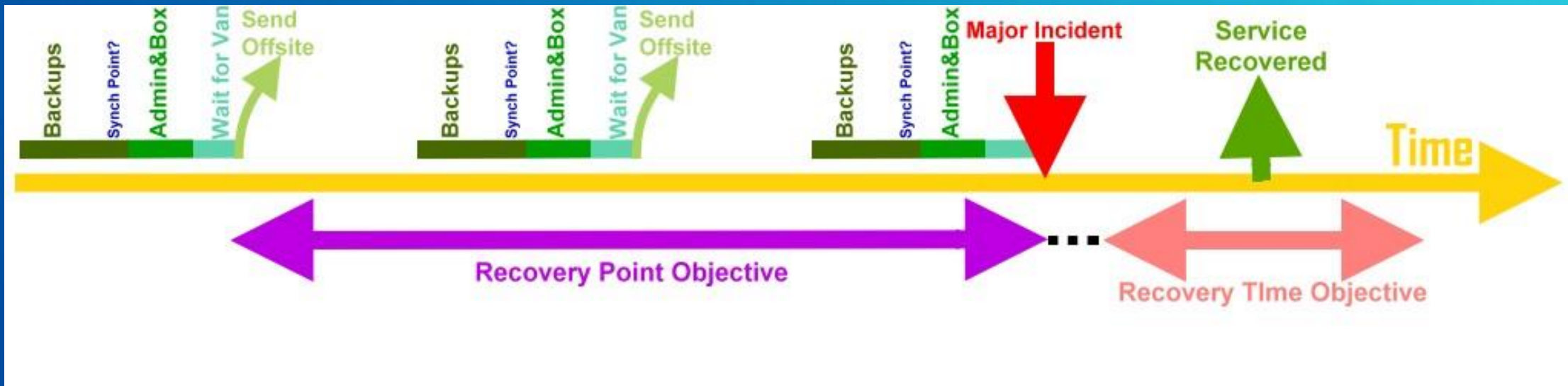


● Recuperação de Desastres

"Recuperação de desastres é a prática de fazer um sistema ser capaz de **sobreviver a falhas inesperadas ou extraordinárias.**" (Reese, 2009)

HIGH
AVAILABILITY





RTO – Tempo máximo esperado para retomar a operacionalização dos serviços;

RPO – Máximo de perda de dados aceitáveis; Quantidade de dados mínimos necessários para voltar as operações.

- Medidas preventivas
- Medidas de detecção
- Medidas corretivas

Nossa pesquisa





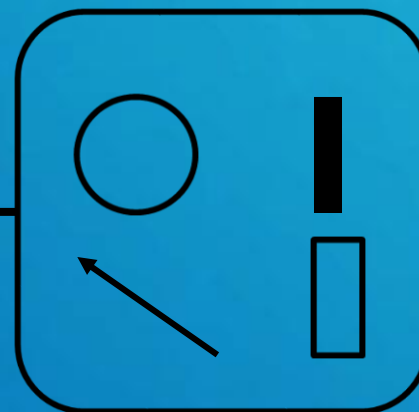
Interface gráfica



Modelos de alto nível



Mercury scripts



SPNs



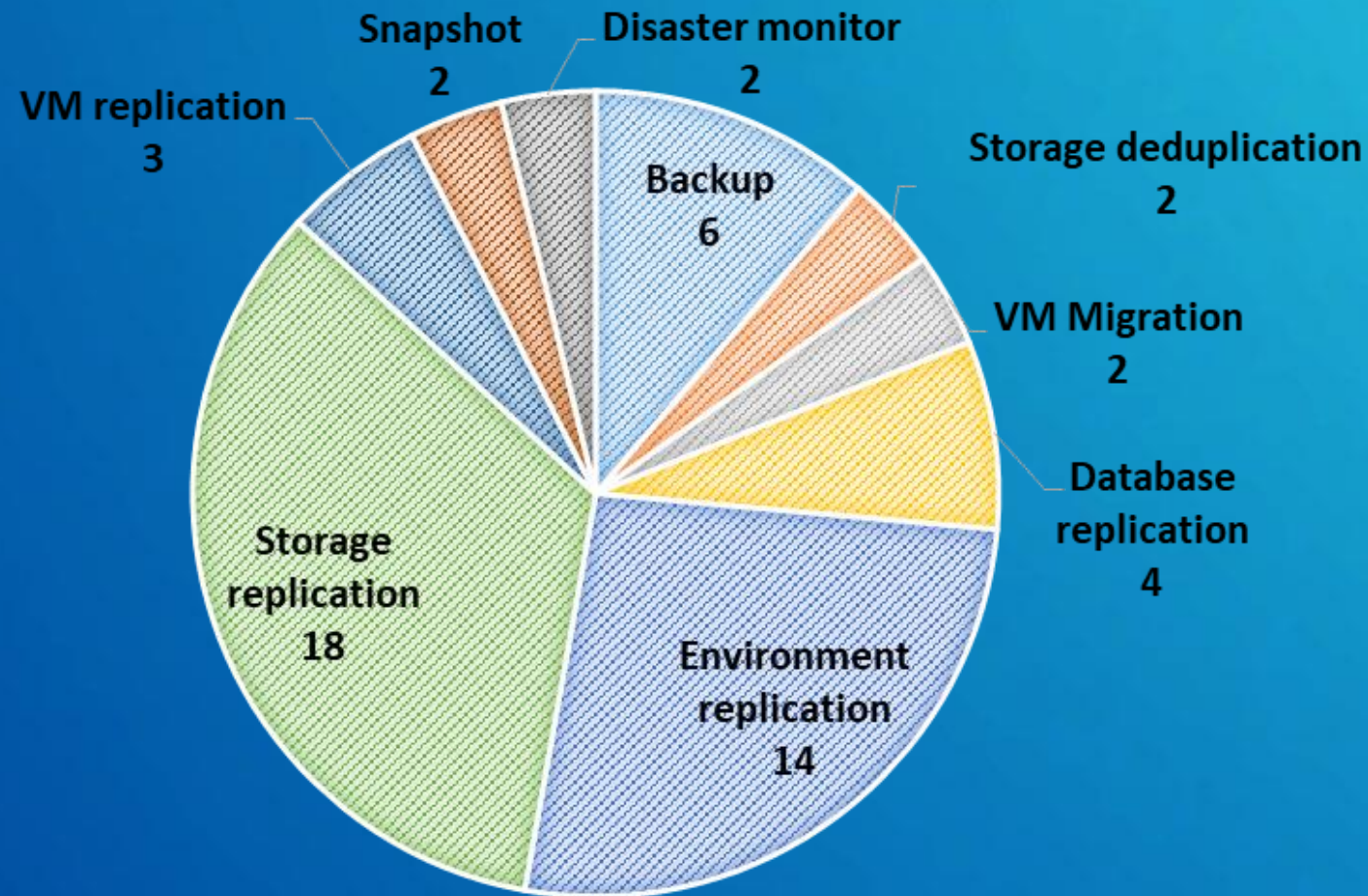
Métricas



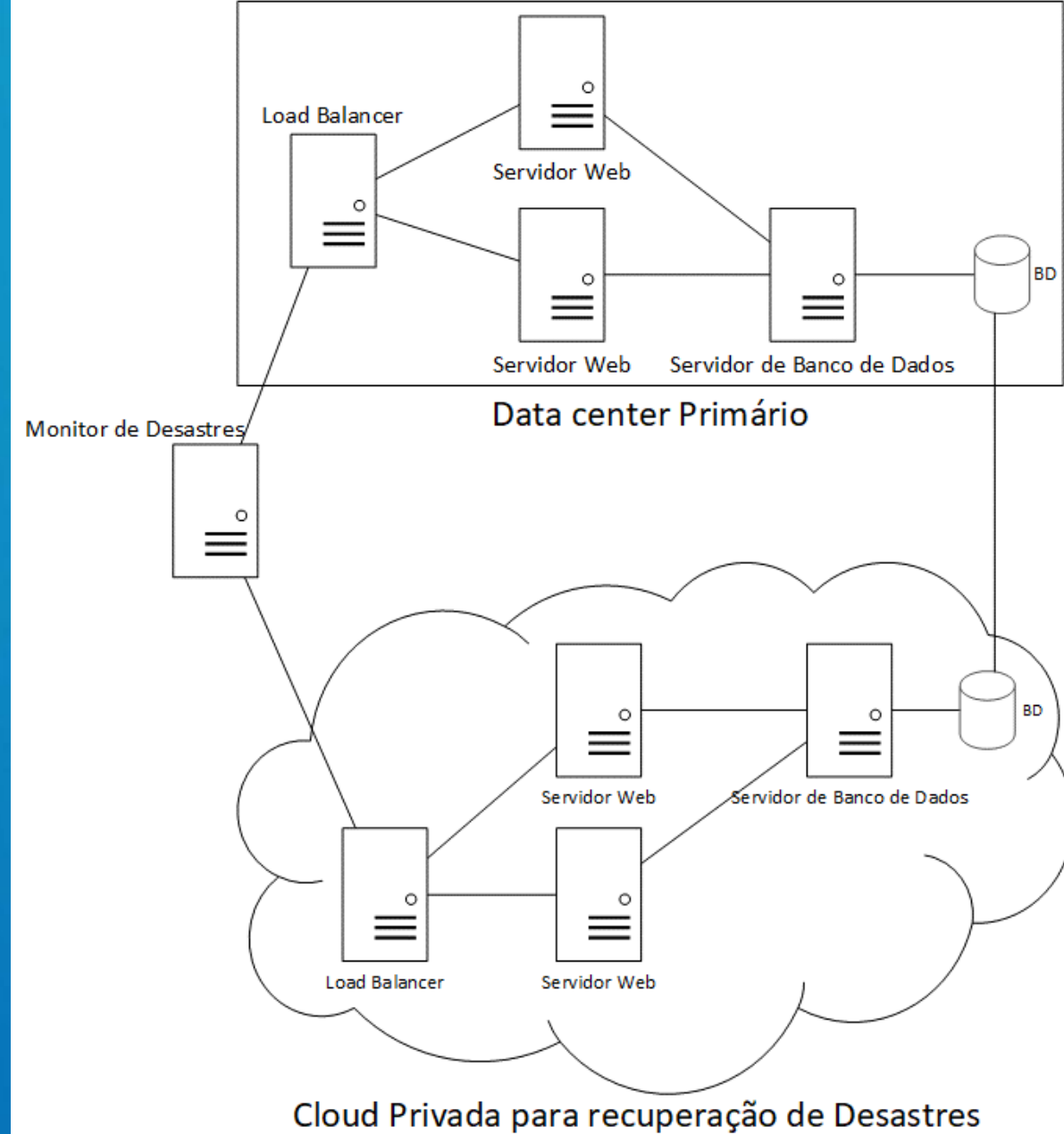
Ranking das soluções

Status do trabalho

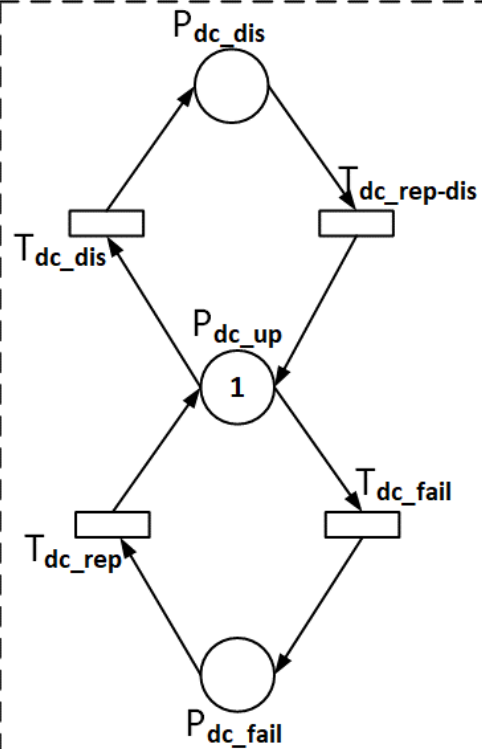
- Realizamos um estudo de mapeamento sistemático para identificar que tipos de estratégias que são mais utilizadas



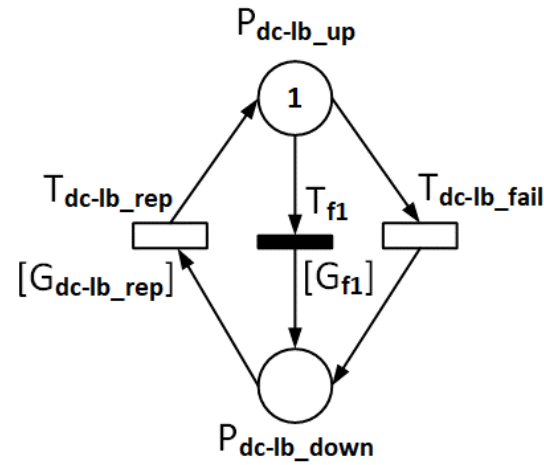
- Fizemos experimentos e modelagem de uma das infraestruturas;
- Avaliamos métricas de disponibilidade e downtime;



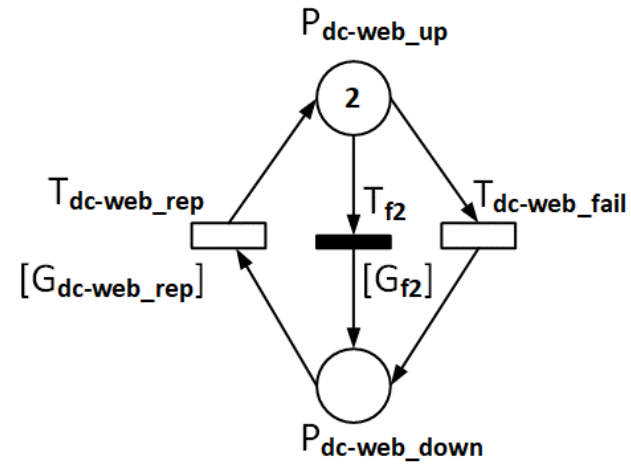
Primary Data center



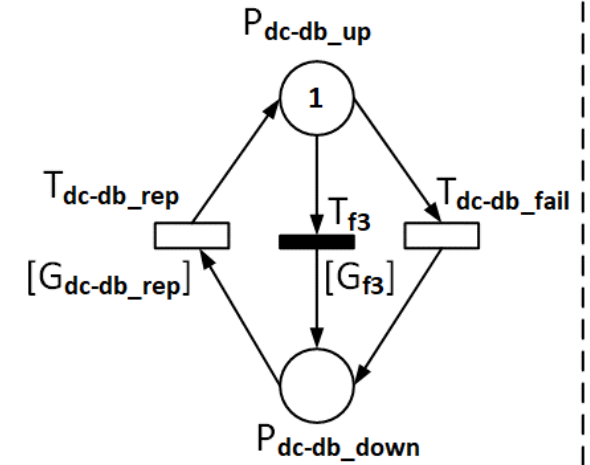
(a) SPN for the Primary data center



(b) SPN for the load balancer

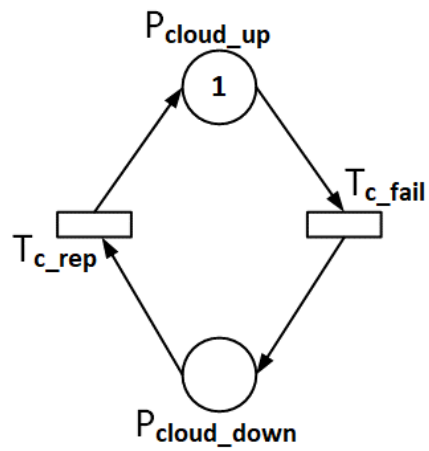


(c) SPN for the web servers

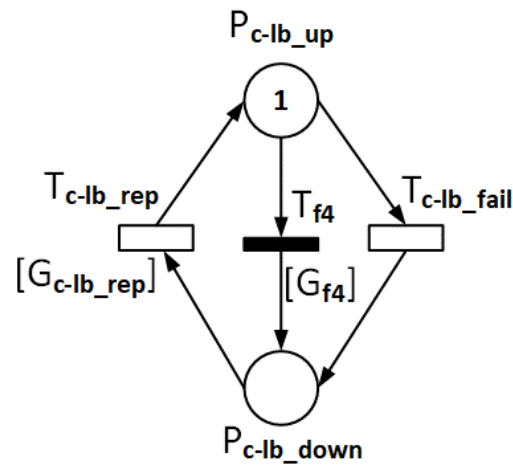


(d) SPN for the database server

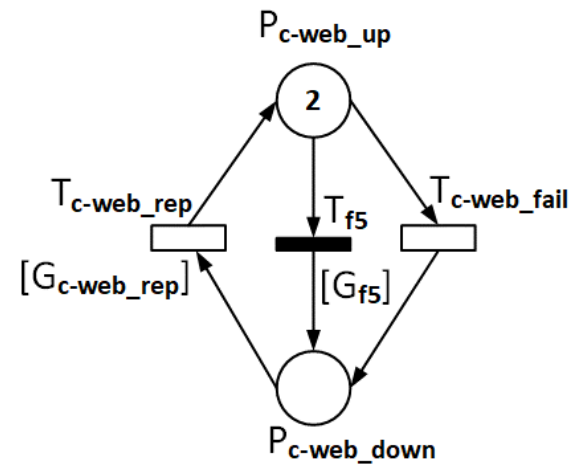
DR Cloud



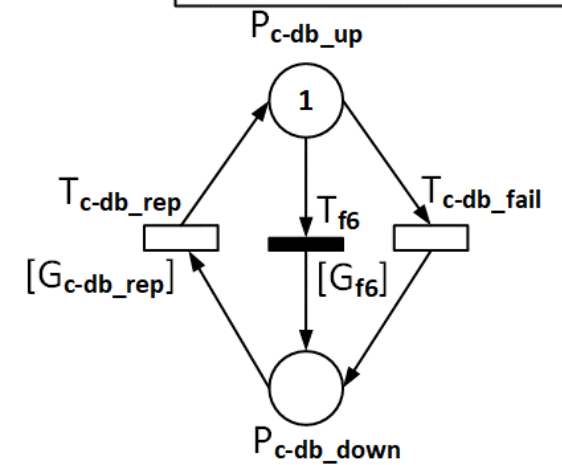
(e) SPN for the cloud server



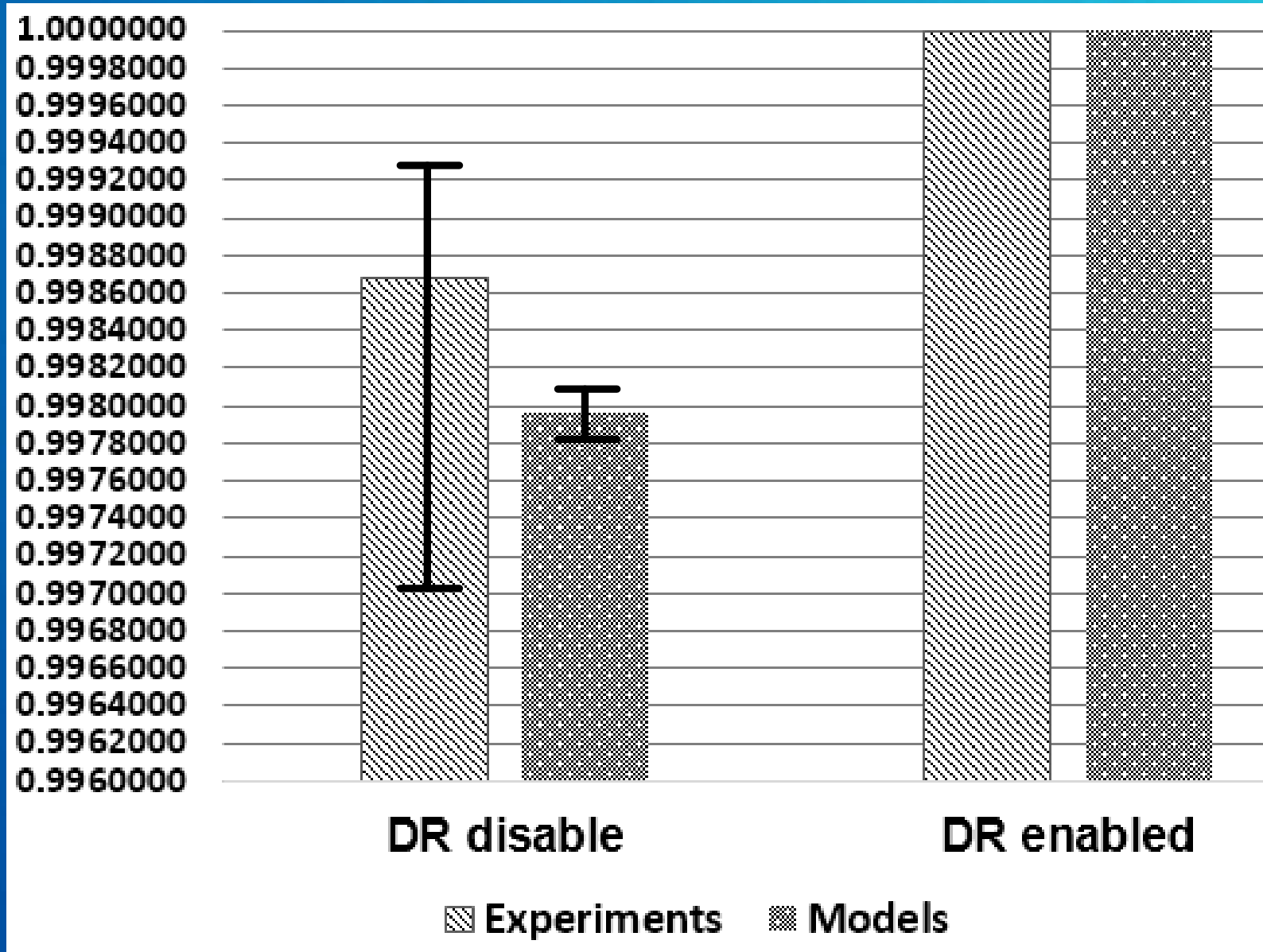
(f) SPN for the cloud load balancer



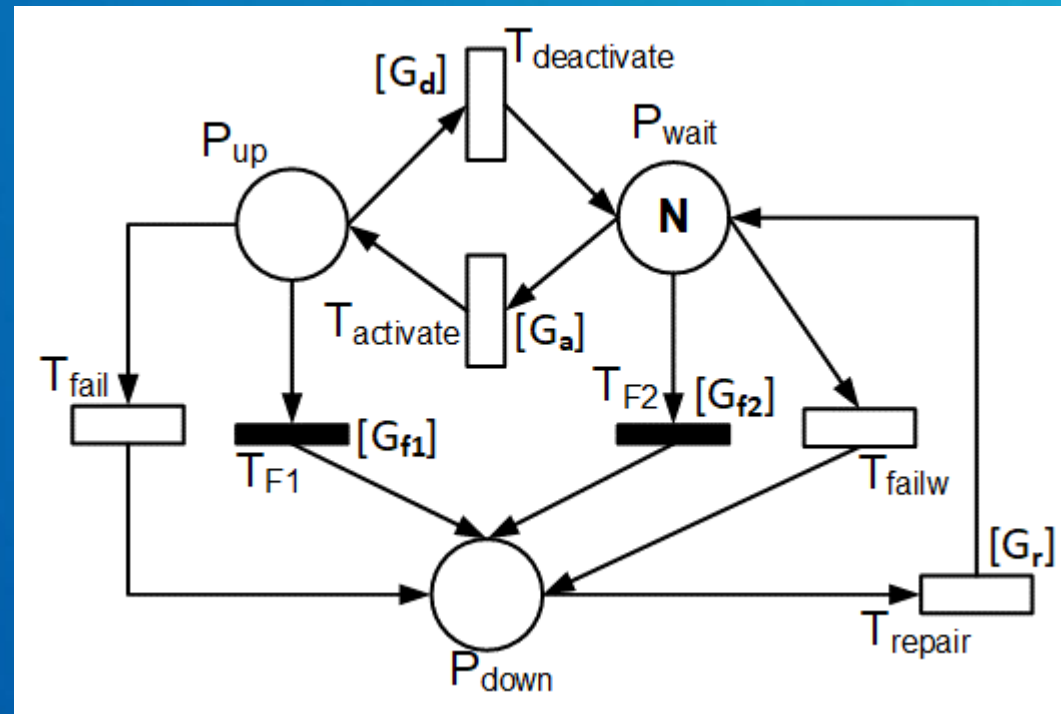
(g) SPN for the cloud web servers



(h) SPN for the cloud database server



- Hot-standby (active/standby) component



- Estamos trabalhando nos modelos de outras infraestruturas;
- Iniciamos a montagem de uma segunda infraestrutura para avaliação;
- Começamos o estudo de viabilidade e concepção da ferramenta.

Próximos passos

- Modelos

- Criação de modelos de infraestruturas de RD;
- Validação de alguns modelos;
- Inclusão de métricas de RTO e custo;

- Infraestruturas

- Seleção de próximas infraestruturas a serem testadas;
- Montar infraestruturas;
- Avaliar infraestruturas;

- Ferramenta

- Finalizar estudo de viabilidade e concepção;
- Desenvolvimento;
- Testes.

Referências

- Reese, G. (2009). *Cloud Application Architectures: Building Applications and Infrastructure in the Cloud*. O'Reilly Media, Inc.
- Nguyen, T. A., Kim, D. S., and Park, J. S. (2016). Availability modeling and analysis of a data center for disaster tolerance. *Future Generation Computer Systems*, 56:27-50
- Zetta (2016). State of disaster recovery 2016. <http://www.zetta.net/resource/statedisaster-recovery-2016>.
- Andrade, E., Nogueira, B., Matos, R., Callou, G., and Maciel, P. (2017). Availability modeling and analysis of a disaster-recovery-as-a-service solution. *Computing*, pages 1-26.
- SILVA, B. (2016). *A Framework For Availability, Performance And Survivability Evaluation Of Disaster Tolerant Cloud Computing Systems*. PhD thesis, Federal University of Pernambuco.
- Keeton, K., Santos, C., Beyer, D., Chase, J., & Wilkes, J. (n.d.). Designing for disasters. Retrieved from <https://pdfs.semanticscholar.org/4718/60e9f10ef0003af7cf7a7c77fdce066a7b18.pdf>

Uma abordagem para a seleção automática de soluções de recuperação de desastres baseada em modelos estocásticos



Júlio Mendonça
jrmn@cin.ufpe.br

09/11/2017