



Wilson Medeiros dos Santos Neto

Avaliação Experimental de Replicação em Banco de Dados para Recuperação de Desastres

Recife

2020

Wilson Medeiros dos Santos Neto

Avaliação Experimental de Replicação em Banco de Dados para Recuperação de Desastres

Artigo apresentado ao Curso de Bacharelado em Ciências da Computação da Universidade Federal Rural de Pernambuco, como requisito parcial para obtenção do título de Bacharel em Ciências da Computação.

Universidade Federal Rural de Pernambuco – UFRPE

Departamento de Computação

Curso de Bacharelado em Ciências da Computação

Orientador: Ermeson Carneiro de Andrade

Coorientador: Júlio Mendonça

Recife

2020

Dados Internacionais de Catalogação na Publicação
Universidade Federal Rural de Pernambuco
Sistema Integrado de Bibliotecas
Gerada automaticamente, mediante os dados fornecidos pelo(a) autor(a)

W751a dos Santos Neto, Wilson
Avaliação Experimental de Replicação em Banco de Dados para Recuperação de Desastres / Wilson dos Santos Neto. - 2020.
26 f. : il.

Orientador: Ermeson Carneiro de Andrade.
Coorientador: Julio Mendonca.
Inclui referências.

Trabalho de Conclusão de Curso (Graduação) - Universidade Federal Rural de Pernambuco,
Bacharelado em Ciência da Computação, Recife, 2020.

1. Replicação de Dados. 2. Recuperação de Desastres. 3. Banco de Dados Relacionais. I. Andrade, Ermeson Carneiro de, orient. II. Mendonca, Julio, coorient. III. Título



**MINISTÉRIO DA EDUCAÇÃO
UNIVERSIDADE FEDERAL RURAL DE PERNAMBUCO (UFRPE)
BACHARELADO EM CIÊNCIA DA COMPUTAÇÃO**

<http://www.bcc.ufrpe.br>

FICHA DE APROVAÇÃO DO TRABALHO DE CONCLUSÃO DE CURSO

Trabalho defendido por Wilson Medeiros às 16 horas do dia 18 de dezembro de 2020, no link meet.google.com/pyw-ajqu-agp, como requisito para conclusão do curso de Bacharelado em Ciência da Computação da Universidade Federal Rural de Pernambuco, intitulado **Avaliação Experimental de Replicação em Banco de Dados para Recuperação de Desastres**, orientado por Ermeson Andrade e aprovado pela seguinte banca examinadora:

Ermeson Andrade
DC/UFRPE

Gabriel Alves
DEINFO/UFRPE

Resumo

Os sistemas de TI são essenciais para as operações de qualquer negócio moderno. Tais sistemas precisam suportar as operações de suas empresas correspondentes sob quaisquer condições. Estratégias de Recuperação de Desastres (RD) têm sido implementadas para auxiliar as organizações a mitigar falhas inesperadas e reduzir gastos desnecessários. No entanto, no melhor do nosso conhecimento, nenhum trabalho analisa experimentalmente a replicação de dados na camada de banco de dados (BD) com foco em estratégias de RD. Desta forma, este trabalho avalia a replicação em BDs relacionais como uma forma de implementar uma solução de RD. Para isso, nós utilizamos um ambiente de testes real em nuvem pública para executar experimentos extensivos visando a implementação da replicação fornecida pelo MySQL, considerando vários cenários no contexto de RD. Nossos resultados mostram como o tempo de resposta, o *Recovery Point Objective* (RPO) e o *Recovery Time Objective* (RTO) variam de acordo com o tamanho dos dados replicados, a configuração da replicação (ex.: assíncrona ou semissíncrona) e a configuração das réplicas. Este trabalho pode auxiliar os coordenadores de RD ou indivíduos a decidir qual configuração de replicação de banco de dados para recuperação de desastres é melhor para seu ambiente de trabalho.

Palavras-chave: Replicação de Dados, Recuperação de Desastres, Banco de Dados Relacionais.

Abstract

IT systems are essential for the operations of any modern business. Such systems must support operations of their corresponding company under any conditions. Disaster Recovery (DR) strategies have been implemented to help organizations mitigate unexpected failures and reduce unnecessary expenses. However, to the best of our knowledge, no other work experimentally analyzes data replication at the database layer with a focus on DR strategies. Therefore, this work evaluates a relational database replication as a mean of implementing a DR solution. We use a real testbed in a public cloud environment to perform extensive experiments aimed at implementing the replication provided by MySQL, considering various scenarios in the context of DR. Our results show how response time, Recovery Point Objective (RPO) and Recovery Time Objective (RTO) vary according to the size of the replicated data, the synchronization type (ex.: asynchronous or semisynchronous) and the configuration of the slave servers. This work can assist DR coordinators or individuals to decide which database replication configuration for disaster recovery is best for their work environment.

Keywords: Data Replication, Disaster Recovery, Relational Databases.

Lista de ilustrações

Figura 1 – Tempo de resposta médio apresentado para diferentes tamanhos de <i>blob</i> e configurações.	18
Figura 2 – Tempos médios para o Δ relay.	19
Figura 3 – Tempo médio para aplicar uma atualização na fila.	20

Lista de tabelas

Tabela 1 – Resultados dos testes de hipótese não paramétricos	18
---	----

Lista de abreviaturas e siglas

RPO	Recovery Point Objective
RTO	Recovery Time Objective
SGBD	Sistema Gerenciador de Banco de Dados
SGBDR	Sistema Gerenciador de Banco de Dados Relacional
TIC	Tecnologia da Informação e Comunicação

Sumário

Lista de ilustrações	4
1 INTRODUÇÃO	8
2 FUNDAMENTOS	10
2.1 Recuperação de Desastres	10
2.2 Replicação de Dados no <i>MySQL</i>	10
3 TRABALHOS RELACIONADOS	13
4 ARQUITETURA EXPERIMENTAL	14
4.1 Ambiente de testes	14
4.2 Execução dos experimentos e Medições	14
4.3 Cenários Analisados	15
5 RESULTADOS EXPERIMENTAIS E DISCUSSÃO	17
5.1 Tempo de Resposta	17
5.2 <i>Recovery Point Objective e Recovery Time Objective</i>	19
5.3 Limitações	21
6 CONCLUSÃO E TRABALHOS FUTUROS	23
REFERÊNCIAS	24

1 Introdução

Desastres muitas vezes podem não ser previstos ou evitados. Qualquer organização está propensa à ocorrência de desastres que podem resultar em danos catastróficos (REESE, 2009; ANDRADE et al., 2017). Em empresas que necessitam de sistemas de Tecnologia da Informação e Comunicação (TICs), além do possível prejuízo físico, a interrupção das operações pode resultar em perdas financeiras consideráveis. Uma pesquisa realizada pela Zetta (ZETTA, 2016) aponta que 67% das empresas pesquisadas poderiam ter um prejuízo acima de US\$ 20.000,00 por dia de indisponibilidade. Um outro estudo realizado pela Unitrends aponta que um grande número de empresas ainda sofrem taxas de perda de dados ou tempo de interrupção acima do aceitável (UNITRENDS, 2019).

Estratégias de Recuperação de Desastres (RD) vêm sendo adotadas para mitigar a perda de dados e garantir a continuidade das operações das empresas (ANDRADE et al., 2017). Existem diversas técnicas que podem ser aplicadas como estratégia de RD (ex.: *backup* ou migração de máquinas virtuais). Em especial, a replicação como estratégia de RD tem sido usada por várias empresas, uma vez que a mesma visa garantir redundância de dados em sistemas de TICs (MENDONÇA et al., 2019). Essa estratégia possibilita, por exemplo, manter um ou mais servidores secundários (Réplicas ou *Slaves*) atualizados em relação a um servidor primário (*Master*). Assim, se o servidor primário falhar, uma réplica poderá assumir o seu papel.

Apesar dos bancos de dados *NoSQL* terem surgido como alternativa aos relacionais, oferecendo melhor desempenho e escalabilidade, os Sistemas Gerenciadores de Banco de Dados Relacionais (SGBDR) ainda são os mais utilizados mundialmente. Sendo o *MySQL* o SGBDR de código aberto mais utilizado (DB engines, 2020; SHAY, 2018). Nesse sentido, alguns trabalhos têm sido desenvolvidos para avaliar ou comparar diferentes Sistemas Gerenciadores de Banco de Dados (SGBDs) (JOGI; SINHA, 2016; SANTANA; ARMENDÁRIZ-IÑIGO; MUÑOZ-ESCOÍ, 2016; WANG et al., 2014). No entanto, no melhor do nosso conhecimento, nenhum dos trabalhos disponíveis na literatura avalia a replicação no contexto de RD.

Este trabalho visa avaliar experimentalmente diferentes cenários utilizando a replicação Primário-Secundário, presente no *MySQL*, para fins de RD. Algumas métricas cruciais para a RD são estimadas, tais como: *Recovery Time Objective* (RTO), *Recovery Point Objective* (RPO), além de tempo de resposta. Nós realizamos testes de carga em um ambiente de testes real para coletar essas métricas, considerando uma gama de diferentes cenários. Desta forma, os resultados focados em RD que

apresentam o real comportamento de um sistema que utiliza replicação de dados em BD relacionais são apresentados. Também analisamos o impacto obtido no tempo de resposta ao utilizar diferentes configurações de replicação. Por fim, verificamos que, com as cargas consideradas, não houve sobrecarga relevante ao utilizar uma ou duas réplicas, o que possibilita a adição de redundâncias para o BD primário com pouco ou nenhum impacto no desempenho.

O restante deste trabalho está organizado da seguinte maneira: o Capítulo 2 apresenta os conceitos básicos de RD e descreve os diferentes tipos de replicação no *MySQL*. O Capítulo 3 apresenta alguns trabalhos relacionados. O Capítulo 4 detalha como os experimentos e medições foram realizados. O Capítulo 5 discute os resultados obtidos. O Capítulo 6 conclui o trabalho e cita possíveis trabalhos futuros.

2 Fundamentos

Este capítulo apresenta os principais conceitos para um melhor entendimento do trabalho.

2.1 Recuperação de Desastres

Qualquer infraestrutura computacional ou sistemas de TIC está vulnerável a um conjunto de interrupções. Algumas dessas vulnerabilidades podem ser eliminadas, ou pelo menos minimizadas, através do uso de estratégias de RD (MENDONÇA et al., 2019). Essas estratégias proveem soluções apropriadas para evitar perda de dados e/ou diminuir o tempo para a recuperação dos serviços depois de uma interrupção (BAUER; ADAMS; EUSTACE, 2011). Assim, podemos definir a RD como a prática de tornar sistemas capazes de sobreviver a falhas inesperadas ou extraordinárias (REESE, 2009).

Na análise de soluções de RD, duas métricas são primordiais, o RPO e o RTO. O RPO aponta a quantidade máxima de dados que pode ser perdida desde a realização do último backup até a ocorrência de uma falha, enquanto o RTO representa o tempo máximo necessário para a recuperação do serviço logo após uma interrupção inesperada (REESE, 2009). Assim, conhecer os valores dessas métricas pode auxiliar na escolha das melhores estratégias de RD para indivíduos ou empresas.

2.2 Replicação de Dados no MySQL

Este capítulo fornece um resumo sobre replicação de dados e detalha os diferentes tipos presentes no MySQL 8 (ORACLE, 2020a). A replicação de dados é usada para manter a sincronização entre diferentes dispositivos (nós) e é realizada quando ocorrem mudanças nas bases de dados. Os nós podem estar relacionados de diferentes formas:

- **Primário-Secundário:** é a relação mais comumente adotada. Um dos nós é considerado primário, enquanto os demais, nós secundários ou réplicas. O BD primário é o responsável por receber atualizações, aplicá-las e propagá-las a suas réplicas. Essa é a implementação padrão do MySQL. Como ela não utiliza protocolos para tratar possíveis conflitos quando mais de um nó recebe atualizações, somente o BD primário pode recebê-las.

- **Group Replication:** Essa relação de replicação visa obter consistência entre os dados dos nós. Nessa replicação, um protocolo deve ser adotado para garantir que todos os nós recebam as mensagens na ordem correta. O *atomic multicast* é um protocolo que pode ser adotado para garantir que todas as mensagens enviadas para um conjunto de nós sejam entregues a todos ou a nenhum deles. *One-phase commit*, *two-phase commit*, e *three-phase commit* também são apresentados como protocolos de confirmação distribuídos. No *one-phase commit*, utilizado na replicação Primário-Secundário, o BD primário propaga as atualizações mas não há uma fase para confirmar se as réplicas realizaram as mesmas com sucesso. Os protocolos de *two-phase commit* e *three-phase commit* adicionam mais uma etapa, permitindo que os nós decidam sobre uma transação. Assim, os nós tentam alcançar um consenso sobre confirmar ou descartar a transação recebida.

Neste trabalho, abordamos a replicação Primário-Secundário. Ela é realizada em três principais etapas: **(1)** ao receber uma atualização, o BD primário salva a transação em seu *binary log* e a envia a suas réplicas; **(2)** as réplicas, ao receber a transação, escrevem-na em seu *relay log*¹; e **(3)** uma ou mais *applier threads* aplicam a transação na réplica, salvando-a em seus respectivos *binary logs*, persistindo os dados (realizando *commit* nas transações) em seguida.

A replicação Primário-Secundário pode ser configurada para funcionar de forma **assíncrona** ou **semissíncrona**. Na replicação **assíncrona** o BD primário persiste uma transação sem nenhum tipo de sincronização com as réplicas, enquanto na **semissíncrona** existe uma sincronização na etapa **(2)** explicada anteriormente. Nessa sincronização, o BD primário espera por um aviso de recebimento, ou *acknowledgement* (*ack*), antes de dar uma resposta ao usuário e persistir a transação. Uma réplica comunica um *ack* ao salvar todos os metadados da transação em seu *relay log* (dados não persistidos ainda). Por esse motivo essa configuração é denominada semissíncrona. Além disso, o BD primário espera pelos *acks* por um tempo pré-definido, que caso seja alcançado, a replicação é reconfigurada automaticamente para a assíncrona. Em um sistema real, para evitar perda de dados, deve-se desabilitar essa reconfiguração automática da replicação semissíncrona para a assíncrona. Essa desabilitação, porém, degradaria a disponibilidade do sistema, uma vez que não é possível realizar operações de escrita enquanto não houver réplicas suficientes. A quantidade de *acks* necessários para que o BD primário persista uma transação também pode ser definida. Portanto, seria necessário utilizar mais réplicas e um número de *acks* menor que o total delas para mitigar esse problema. A decisão de utilizar ou não uma reconfiguração automática fica a critério de quem implementa a estratégia de RD. É importante perceber

¹ O *relay log* funciona como uma fila de transações a serem persistidas na réplica.

que o processo de sincronização tem impacto direto no tempo de resposta, visto que o BD primário aguardará a sincronização para dar resposta ao usuário. Mas o impacto não é tão grande quanto ao utilizar *Group Replication*.

3 Trabalhos Relacionados

Neste capítulo são apresentados alguns trabalhos que têm sido desenvolvidos para análise de BDs. Os autores (WANG et al., 2014), (JOGI; SINHA, 2016) e (SANTANA; ARMENDÁRIZ-IÑIGO; MUÑOZ-ESCOÍ, 2016) realizaram experimentos comparando diferentes implementações de BDs. (WANG et al., 2014) analisou a replicação utilizando dois BDs *NoSQL*: HBase e Cassandra. A abordagem utilizada focou em analisar os *tradeoffs* entre consistência e desempenho (tempo de resposta e vazão) entre diferentes configurações, além de escalabilidade. Em (JOGI; SINHA, 2016), além dos BDs *NoSQL* (HBase e Cassandra), Jogi e Sinha também utilizaram o *MySQL* (BD relacional) para comparar o desempenho desses BDs em relação as operações de escrita com grandes quantidades de dados. Já (SANTANA; ARMENDÁRIZ-IÑIGO; MUÑOZ-ESCOÍ, 2016) apresentou um estudo de replicação de BDs para ambientes de computação em nuvem. Através de experimentos, os autores avaliaram diferentes técnicas de replicação, focando principalmente em métricas como tempo de resposta e taxa de falha das transações.

Modelos formais também têm sido empregados para representar e avaliar sistemas com replicações de dados. (RODRIGUES et al., 2019) propôs o uso de Redes de Petri Generalizadas (GSPN) (MARSAN et al., 1991) para modelar um sistema que utiliza replicação. O trabalho utilizou o BD *NoSQL MongoDB* para avaliar métricas como vazão e disponibilidade. No entanto, no modelo adotado os BDs secundários são utilizadas apenas para aumentar a disponibilidade do sistema. No nosso trabalho anterior (MENDONÇA et al., 2019), avaliamos a replicação de dados em BDs relacionais através de modelos formais. No entanto, experimentos foram conduzidos apenas com o objetivo de parametrizar os modelos desenvolvidos.

Os trabalhos acima citados avaliam o uso de replicação de dados em BDs sem considerar características de RD, como as métricas de RTO e RPO. Diferentemente desses trabalhos, neste artigo avaliamos a replicação de dados em um BD relacional através de testes experimentais considerando diferentes cenários. No nosso estudo, consideramos a análise de RTO e RPO, verificando o comportamento das diferentes configurações da replicação de dados no *MySQL 8*, além de avaliar conjuntamente o desempenho delas através da métrica de tempo de resposta.

4 Arquitetura Experimental

Este capítulo detalha a metodologia e o ambiente utilizada na realização dos experimentos.

4.1 Ambiente de testes

Este capítulo descreve como o ambiente de teste foi configurado para realização dos experimentos. O ambiente de testes foi configurado em uma nuvem pública, o *Google Cloud* (Google, 2020). Assim, todo o ambiente foi virtualizado. No total, quatro Máquinas Virtuais (VMs) foram utilizadas: uma delas responsável por gerar a carga dos testes, representando requisições de usuários, foi configurada com o *Apache JMeter* (HALILI, 2008); Outra, realizando o papel de servidor de BD primário; e as outras representando servidores de BD secundários. Todas as VMs configuradas como servidor de BD utilizaram *MySQL 8* como SGDB.

Utilizamos dois servidores como réplicas para verificar se existe diferença significativa ao utilizar réplicas em locais geograficamente diferentes (mais ou menos distantes do nó primário). Além disso, essa configuração também foi utilizada para verificar se há impacto no processo de replicação ao utilizar uma ou duas réplicas. Dessa forma, as localizações das VMs foram definidas da seguinte maneira: a VM configurada com o *JMeter*, responsável por representar o envio de requisições de usuários, ficou na região *west US-region*; o servidor de BD primário, na região *center US-region*; uma das réplicas, na região *east US-region* enquanto a outra réplica, na região *Europe-west4-a*. Todas as VMs utilizaram como sistema operacional o *Ubuntu 18.04 LTS*, tinham 2 v-CPU's, RAM de 7,5GB e SSD de 80GB.

4.2 Execução dos experimentos e Medições

Este capítulo detalha como os experimentos foram conduzidos e quais dados foram coletados. Utilizando o ambiente de testes explicado anteriormente, utilizamos a ferramenta *Apache JMeter* (HALILI, 2008) para coletar dados no ambiente configurado.

Como a replicação ocorre quando há alteração nos dados, a carga de trabalho que utilizamos contém apenas atualizações, desconsiderando a leitura de dados. Assim, a carga de trabalho foi caracterizada da seguinte forma: Usuários realizam operações SQL de *insert* e *update* durante um período de tempo T . Neste período de tempo, várias operações são realizadas pelos mesmos usuários. Na primeira operação de um

usuário, ele realiza um *insert*, nas demais operações um *update* no registro já gerado. Essas operações de *insert* e *update* inserem e modificam um campo do tipo *blob* (dado binário de tamanho variável) (ORACLE, 2020b) em uma única tabela no BD primário. Essa tabela possui como campos, apenas uma chave primária (do tipo inteiro) e um campo do tipo *blob*.

Dessa forma, coletamos os seguintes dados:

- **Tempo de resposta:** É o tempo decorrido para dar uma resposta ao usuário. Esse valor é disponibilizado pelo *JMeter*. Essa métrica é relacionada ao desempenho do ambiente. O objetivo é verificar como ela é afetada ao ativar as diferentes configurações de replicação.
- **Master commit e Slave commit:** Representam os tempos nos quais uma transação foi salva no *binary log* do BD primário e do BD secundário, respectivamente. Esses valores podem ser coletados através da realização de um *parsing* nos *binary logs* do BD primário e secundário.
- **Relay start e Relay end:** Representam os tempos nos quais a transação começou e finalizou de ser escrita no *relay log* da réplica, respectivamente. Esses valores são disponibilizados no *MySQL 8* através do *performance_schema*. Porém, como somente os valores da última transação podem ser consultados, coletamos as amostras através de uma *thread* no *JMeter* que realiza a consulta em intervalos de 1 segundo.

4.3 Cenários Analisados

Neste capítulo são descritos os diferentes cenários analisados nos experimentos. Foram adotados os seguintes parâmetros para a composição dos cenários:

- **Tipo da sincronização:** assíncrona ou semissíncrona.
- **Configuração das réplicas:** Se o tipo da sincronização for assíncrona:
 - 1 réplica perto do BD primário (*east US-region*);
 - 1 réplica longe do BD primário (*Europe-west4-a*);
 - 2 réplicas;

Se o tipo da sincronização for semissíncrona:

- 1 réplica perto do BD primário (*east US-region*);
- 1 réplica longe do BD primário (*Europe-west4-a*);

- 2 réplicas e 1 *ack*;
- 2 réplicas e 2 *acks*.

- **Carga dos usuários:**

- 0,1 requisições por segundo (0,1 req/s);
- 2 requisições por segundo (2 req/s).

Para a carga de 0,1 req/s, 1 requisição a cada 10 segundos é enviada por 1 usuário e tem duração de 1000 segundos. Já para a carga de 2 req/s, as requisições são enviadas por 50 usuários, com intervalos gerados aleatoriamente pelo *JMeter*, de forma a alcançar a vazão desejada de 2 req/s e tem duração de 240 segundos.

- **Tamanho do *blob* enviado:** 500 B, 25 KB, 50 KB, 75 KB, 100 KB.

Cada cenário analisado, é formado pela combinação desses parâmetros. Isto resulta em 30 cenários para a replicação assíncrona e 40 cenários para a replicação semi-síncrona, totalizando 70 cenários analisados. Além disso, levando em consideração que a latência pode sofrer variações no ambiente utilizado, cada cenário foi executado duas vezes num intervalo de aproximadamente 13 horas. Deste modo, são geradas duas amostras para cada cenário.

5 Resultados Experimentais e Discussão

Este capítulo discute os resultados obtidos através dos experimentos realizados no ambiente adotado onde, mais especificamente, são apresentando os resultados referente às métricas de tempo de resposta, RPO e RTO.

5.1 Tempo de Resposta

A Figura 1 apresenta os resultados obtidos para o tempo de resposta ao utilizar 1 réplica perto ou longe em relação ao BD primário. Como pode ser visto na figura, há quatro gráficos que representam combinações de dois parâmetros: **(1)** configuração da replicação (assíncrona ou semissíncrona) e **(2)** carga dos usuários (0,1 req/seg ou 2 req/s). Note que para cada cenário são apresentados os resultados de tempo de resposta (eixo Y) utilizando tamanhos distintos para o arquivo que é enviado na requisição (*blob* – apresentado no eixo X). Os pontos considerados representam os resultados obtidos na primeira execução dos experimentos.

A partir dos resultados obtidos, podemos fazer algumas observações com relação a replicação assíncrona e semissíncrona. Na replicação assíncrona o tempo de resposta apresentou diferença apenas ao variar o tamanho da requisição que era enviada (*blob*). Isso se deve ao fato do tempo de resposta nesta configuração de replicação depender apenas do quão rápido o servidor primário pode processar uma requisição feita pelo usuário. Por outro lado, na replicação semissíncrona, além do tamanho da requisição enviada (*blob*), houve diferenças significativas com cargas dos usuários e localizações distintas das réplicas. Como esperado, esta configuração de replicação apresentou um tempo de resposta maior que a replicação assíncrona, tendo essa diferença atenuada ao aumentar o tamanho da requisição enviada e, da mesma maneira, ao utilizar a réplica mais distante. Além disso, a carga de 0,1 req/s apresentou tempo de resposta maior em relação a carga de 2 req/s. Isso pode ser explicado se considerarmos que as atualizações podem ser enviadas às réplicas em “lotes”, diminuindo o tempo de envio de algumas atualizações às réplicas quando se utiliza uma carga requisições maior.

Para comparar os outros cenários nós utilizamos o Teste U de Mann-Whitney (HOLLANDER; WOLFE, 1999), que é um teste não paramétrico para comparar duas amostras independentes. Para realizar tal teste, foram juntadas as amostras obtidas nas duas execuções do experimento e também as amostras com tamanhos de *blob* diferentes, a fim de obter uma comparação mais generalizada. O parâmetro comparado foi a **configuração das réplicas** (ver Seção 4.3). A Tabela 1 exibe os resultados do

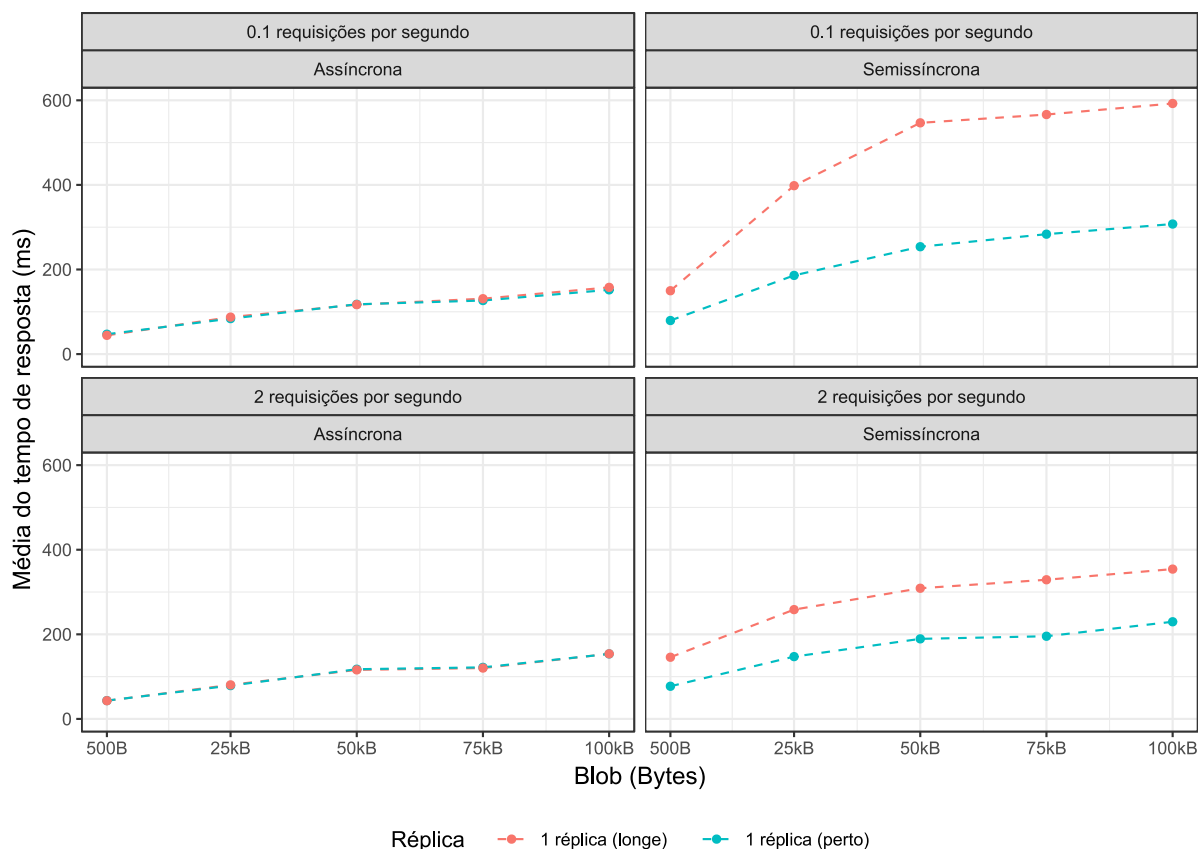


Figura 1 – Tempo de resposta médio apresentado para diferentes tamanhos de *blob* e configurações.

Tabela 1 – Resultados dos testes de hipótese não paramétricos

Comparações entre cenários			
Rep. & workload	Configurações	<i>p</i> -value (latência)	<i>p</i> -value (delta relay)
Assínc. & 0,1 req/s	1 réplica perto x 2 réplicas	0,7786966	0,1240706
	1 réplica longe x 2 réplicas	0,6967137	0,6118413
Assínc. & 2 req/s	1 réplica perto x 2 réplicas	0,2665331	0,1020656
	1 réplica longe x 2 réplicas	0,1237713	0,1413215
Semissínc. & 0,1 req/s	1 réplica perto x 2 réplicas (1 ack)	0,8865125	0,3204590
	1 réplica longe x 2 réplicas (2 acks)	0,8489201	0,9123771
Semissínc. & 2 req/s	1 réplica perto x 2 réplicas (1 acks)	0,6444780	0,1296032
	1 réplica longe x 2 réplicas (2 acks)	0,1002598	0,4275970

valor *p* para a comparação da latência em cenários distintos. Com um intervalo de confiança de 95%, os resultados do teste indicam que todas essas comparações não apresentam diferenças significativas, pois todos os valores de *p* são maiores que 0,05 (5%). Isso mostra que, com as cargas consideradas, utilizar 2 réplicas ao invés de 1 não impacta significativamente no tempo de resposta do sistema.

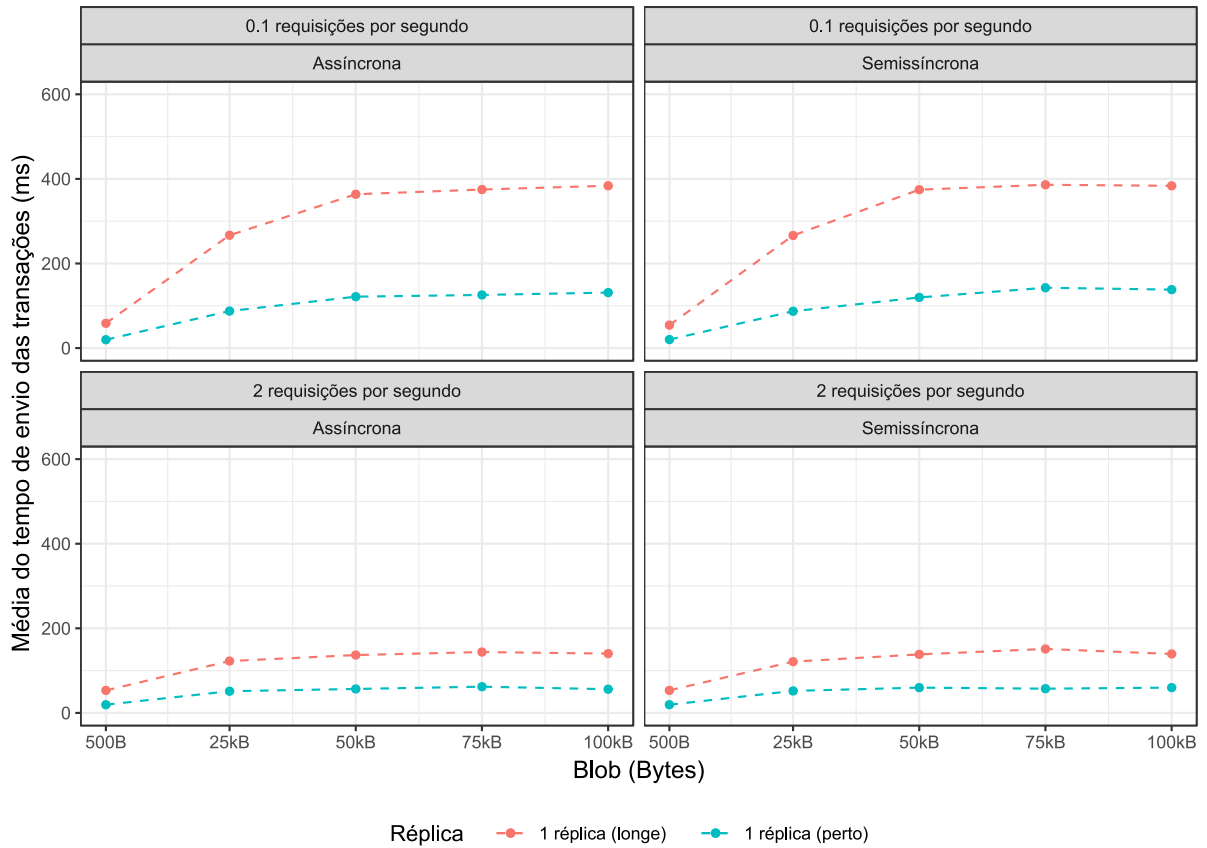


Figura 2 – Tempos médios para o Δ relay.

5.2 Recovery Point Objective e Recovery Time Objective

Para estimar o RPO, primeiro é preciso identificar em quais cenários algum dado pode ser perdido caso o BD primário falhe. Em um BD relacional, podemos considerar que uma transação foi perdida se o BD primário a persistiu (realizou *commit*) mas falhou antes de enviá-la completamente a ao menos uma réplica. Assim, para estimar tal valor, podemos calcular a diferença entre o tempo *Master commit* (T_{Master_commit}) e o tempo *Relay end* ($T_{Slave_relay_end}$) da última transação persistida no BD primário antes da falha. Neste trabalho, denominaremos essa diferença de *Delta Relay* (Δ relay). A Equação 5.1 detalha esse cálculo.

$$RPO = \Delta relay = T_{Slave_relay_end} - T_{Master_commit} \quad (5.1)$$

Para a replicação assíncrona, consideramos que os valores de Δ relay coletados são os possíveis intervalos do RPO, pois o BD primário poderia vir a falhar em qualquer momento durante a execução. Para a replicação semissíncrona é diferente, pois ela oculta a transação e não dá resposta ao usuário até que um número desejado de réplicas enviem *acks*. Conseqüentemente, o RPO é zero enquanto a replicação semissíncrona for usada.

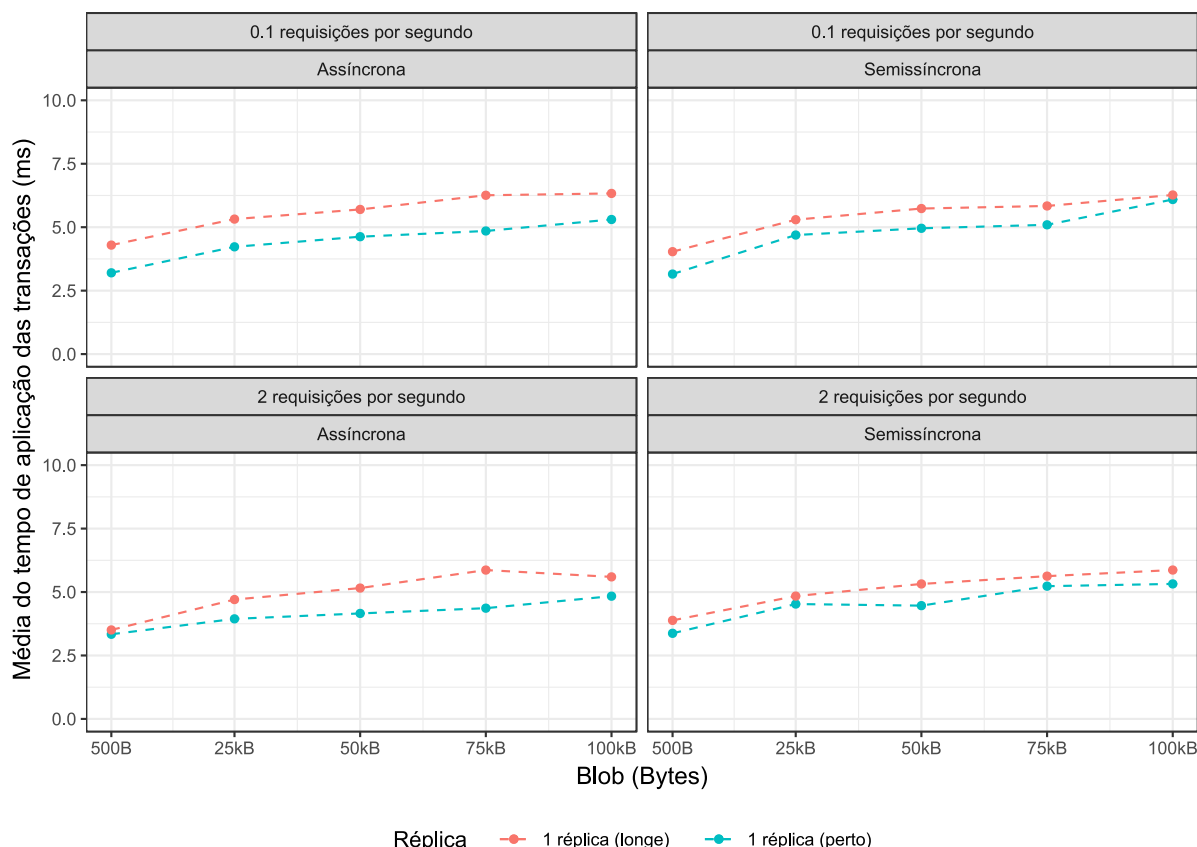


Figura 3 – Tempo médio para aplicar uma atualização na fila.

A Figura 2 apresenta os resultados obtidos para o Δ relay, gerados pelas transações realizadas ao utilizar 1 réplica perto ou longe em relação ao BD primário. O gráfico é estruturado de maneira similar ao do tempo de resposta (ver Figura 1). O envio das transações ocorre em ambas as configurações e, por isso, diferentemente do tempo de resposta, os resultados se comportam de maneira similar tanto na replicação assíncrona quanto na semissíncrona. Apesar disso, o Δ relay de cada configuração influencia diferentemente as métricas de interesse. Na replicação assíncrona, ele indica o intervalo de tempo em que alguma atualização pode ser perdida, ou seja, o RPO. Por outro lado, na replicação semissíncrona, ele indica a sobrecarga esperada no tempo de resposta. Devido à influência do Δ relay no tempo de resposta da replicação semissíncrona, as ponderações feitas sobre ela na Seção 5.1 também são válidas aqui. No entanto, a variação do tempo com o aumento do *blob* é mais sutil, pois o Δ relay depende da comunicação feita entre dois nós (*primário-secundário*) ao invés de três nós (*usuário-primário-secundário*), como ocorre para o caso do tempo de resposta com replicação semissíncrona. Por fim, também foram executados testes estatísticos e não houve diferença relevante nesse tempo ao utilizar 1 ou 2 réplicas (ver coluna *p-value (delta relay)*) na Tabela 1.

O RTO refere-se ao tempo que o sistema volta ao estado operacional após a

ocorrência de um desastre. Geralmente, um monitor de desastres deve detectar a ocorrência de um desastre ou falha do BD primário para que um processo de *failover*¹ se inicie. Ou seja, uma das réplicas deve assumir o papel de BD primário após um desastre. Para que isso ocorra, é necessário esperar que a réplica termine de aplicar todas as transações pendentes (enfileiradas em seu *relay log*). Pois caso novas transações sejam realizadas nele antes disso, seus dados assumirão um estado diferente daqueles no BD primário, além de que as aplicações das transações pendentes podem vir a falhar (pela inconsistência gerada). Isso se deve à falta de consenso entre os nós antes do *commit*. Assim, se após o processo de *failover* ainda houver transações na fila para a réplica aplicar, assumimos que o RTO pode ser calculado pelo tempo que ela levará para aplicar todas as transações já enfileiradas. Isso se dá subtraindo o tempo *Slave commit* (T_{Slave_commit}) do tempo *Relay end* ($T_{Slave_relay_end}$) da última transação enfileirada antes do BD primário falhar, como mostrado na Equação 5.2. Então o RTO seria o máximo entre o tempo de *failover* ($T_{failover}$) e o tempo de aplicação da última transação ($T_{Slave_apply_relay}$), como mostrado na (Equação 5.3).

$$T_{Slave_apply_relay} = T_{Slave_commit} - T_{Slave_relay_end} \quad (5.2)$$

$$RTO = \max(T_{Slave_apply_relay}, T_{failover}) \quad (5.3)$$

A Figura 3 mostra os resultados obtidos para o T_{apply_relay} , gerados pelas transações realizadas ao utilizar 1 réplica perto ou longe em relação ao BD primário. É importante ser destacado que esse tempo de processamento depende apenas de quão rápido a réplica consegue processar e persistir uma atualização no seu *binary log*, sem necessidade de comunicação com outros nós. Por este motivo, a quantidade de réplicas não influencia no resultado. Além disso, as diferenças entre os cenários exibidos no gráfico foram mais sutis. Dessa forma, o RTO com os parâmetros considerados apresenta diferenças mínimas entre todas as configurações.

5.3 Limitações

Os cenários considerados buscam representar situações reais de interação entre os usuários e o sistema implementado. Porém, nem sempre a carga de trabalho esperada de um sistema será similar às que aqui foram utilizadas. Ademais, há cenários que poderiam apresentar diferenças relevantes que não foram considerados neste trabalho, como utilizar mais de duas réplicas e realizar testes de estresse (sob cargas extremas) a fim de verificar as principais limitações da implementação. Por fim, neste

¹ Processo no qual um servidor secundário assume as operações no lugar de um servidor primário que entrou em estado de falha.

trabalho não realizamos as operações de recuperação de falhas no ambiente adotado, pois seu escopo foca na análise dos mecanismos de replicação.

6 Conclusão e Trabalhos Futuros

Neste trabalho avaliamos a implementação da replicação presente no *MySQL 8*. Através de um sistema implantado em uma nuvem pública e nós geograficamente distantes, testes experimentais foram realizados para este fim. Apresentamos o impacto gerado no desempenho ao utilizar replicação semissíncrona no lugar de assíncrona, considerando nós em localidades diferentes. Também estimamos o RPO esperado ao utilizar replicação assíncrona e verificamos que a utilização de duas réplicas resultou em sobrecarga irrelevante em relação a apenas um com as cargas consideradas. Finalmente, verificamos que o RTO não apresenta diferenças significativas entre os cenários considerados. Espera-se que os resultados apresentados e os pontos destacados auxiliem coordenadores de DR e indivíduos que considerem adotar essa técnica como solução de RD.

Como trabalho futuro, almejamos avaliar e comparar outras implementações presentes em BDs relacionais. Também pretendemos verificar o impacto no desempenho ao utilizar diferentes níveis de consistência entre os nós e propor possíveis alternativas a fim de minimizá-lo. Por fim, objetivamos realizar testes de estresse para melhor apontar as limitações de cada implementação estudada.

Considerações Finais

Este trabalho foi publicado nos Anais do XIX Workshop em Desempenho de Sistemas Computacionais e de Comunicação ([MEDEIROS et al., 2020](#)).

Referências

- ANDRADE, E. et al. Availability modeling and analysis of a disaster-recovery-as-a-service solution. *Computing*, Springer, v. 99, n. 10, p. 929–954, 2017. Citado na página 8.
- BAUER, E.; ADAMS, R.; EUSTACE, D. *Beyond Redundancy: How Geographic Redundancy Can Improve Service Availability and Reliability of Computer-Based Systems*. [S.l.]: Wiley, 2011. Citado na página 10.
- DB engines. *DB-Engines Ranking*. 2020. [Online]. <<https://bit.ly/2s90Xvl>>. Citado na página 8.
- Google. *Google Cloud*. 2020. [Online]. <<https://cloud.google.com>>. Citado na página 14.
- HALILI, E. H. *Apache JMeter: A practical beginner's guide to automated testing and performance measurement for your websites*. [S.l.]: Packt Publishing Ltd, 2008. Citado na página 14.
- HOLLANDER, M.; WOLFE, D. *Nonparametric Statistical Methods*. [S.l.]: Wiley, 1999. (Wiley Series in Probability and Statistics). ISBN 9780471190455. Citado na página 17.
- JOGI, V. D.; SINHA, A. Performance evaluation of MySQL, Cassandra and HBase for heavy write operation. In: *2016 3rd International Conference on Recent Advances in Information Technology (RAIT)*. [S.l.]: IEEE, 2016. p. 586–590. Citado 2 vezes nas páginas 8 e 13.
- MARSAN, M. et al. An introduction to generalized stochastic petri nets. *Microelectronics Reliability*, v. 31, n. 4, p. 699 – 725, 1991. ISSN 0026-2714. Disponível em: <<http://www.sciencedirect.com/science/article/pii/0026271491900105>>. Citado na página 13.
- MEDEIROS, W. et al. Avaliação experimental de replicação em banco de dados para recuperação de desastres. In: SBC. *Anais do XIX Workshop em Desempenho de Sistemas Computacionais e de Comunicação*. [S.l.], 2020. p. 121–132. Citado na página 23.
- MENDONÇA, J. et al. Evaluating database replication mechanisms for disaster recovery in cloud environments. In: IEEE. *2019 IEEE International Conference on Systems, Man and Cybernetics (SMC)*. [S.l.], 2019. p. 2358–2363. Citado na página 13.
- MENDONÇA, J. et al. Disaster recovery solutions for it systems: A systematic mapping study. *Journal of Systems and Software*, v. 149, p. 511 – 530, 2019. ISSN 0164-1212. Citado 2 vezes nas páginas 8 e 10.
- ORACLE. *MySQL 8.0 Reference Manual*. 2020. [Online]. <<https://dev.mysql.com/doc/refman/8.0/en/>>. Citado na página 10.

ORACLE. *The BLOB and TEXT Types*. 2020. [Online]. <<https://dev.mysql.com/doc/refman/8.0/en/blob.html>>. Citado na página 15.

REESE, G. *Cloud application architectures: building applications and infrastructure in the cloud*. [S.l.]: "O'Reilly Media, Inc.", 2009. Citado 2 vezes nas páginas 8 e 10.

RODRIGUES, M. et al. Evaluation of nosql dbms in private cloud environment: An approach based on stochastic modeling. In: IEEE. *2019 IEEE International Systems Conference (SysCon)*. [S.l.], 2019. p. 1–7. Citado na página 13.

SANTANA, M.; ARMENDÁRIZ-IÑIGO, J. E.; MUÑOZ-ESCOÍ, F. D. Evaluation of Database Replication Techniques for Cloud Systems. *Computing and Informatics*, v. 34, n. 5, p. 973–995, 2016. Citado 2 vezes nas páginas 8 e 13.

SHAY, T. *Most popular databases in 2018 according to StackOverflow survey*. 2018. [Online]. <<https://bit.ly/2DCwqhj>>. Citado na página 8.

UNITRENDS. *DATA PROTECTION, CLOUD, AND PROOF DRaaS DELIVERS – UNITRENDS 2019 SURVEY RESULTS*. 2019. Tech Report. [Online]. <<https://bit.ly/2XubXDo>>. Citado na página 8.

WANG, H. et al. Benchmarking replication and consistency strategies in cloud serving databases: Hbase and cassandra. In: SPRINGER. *Workshop on Big Data Benchmarks, Performance Optimization, and Emerging Hardware*. [S.l.], 2014. p. 71–82. Citado 2 vezes nas páginas 8 e 13.

ZETTA, I. *State of Disaster Recovery 2016*. 2016. [Online]. <<https://bit.ly/2H6TwhN>>. Citado na página 8.