

ValorizeAI: Documentação e Validação de uma Arquitetura Serverless Elasticamente Gerenciada

Title: ValorizeAI: Documenting and Validating a Managed Serverless Architecture

Felipe Tomkiel Malacarne, Prof. Me. Marcos André Lucas

¹ Universidade Regional Integrada do Alto Uruguai e das Missões
Departamento de Engenharias e Ciência da Computação
Caixa Postal 743 – 99.709-910 – Erechim – RS – Brasil

101090@uricer.edu.br, mlucas@uricer.edu.br

Abstract. *Digital financial workloads demand elastic architectures that remain observable under sustained stress. This paper documents and validates ValorizeAI, a serverless/containerized platform built on Cloud Run, Cloud SQL, Redis, Cloud Tasks, and WebSockets. Infrastructure-as-code modules provisioned replicable testbeds, while k6 executed two load scenarios (read-only and read/write) against SLOs of $P95 \leq 300$ ms and error rate $< 0.5\%$. The read scenario sustained 1,000 virtual users (mean 470 req/s, peak ≈ 970 req/s) with $P95 = 158$ ms and zero failures; the mixed scenario held 226 req/s at 650 VUs (violating the SLO at ≈ 539 VUs) with $P95 = 658$ ms and $p99 = 2.67$ s once the 10-instance Cloud Run quota was exhausted. The asynchronous pipeline processed 51.58k Cloud Tasks in 10 minutes (86 tasks/s) without loss or duplication. These results show that the architecture comfortably serves read-heavy traffic and that further horizontal scaling or write-path optimizations are required before supporting higher write concurrency.*

Keywords *Serverless Computing, Cloud Run, Financial Platforms, Performance Testing, Observability.*

Resumo. *Plataformas financeiras modernas exigem arquiteturas elásticas capazes de lidar com picos de carga transacional sem abrir mão de consistência forte nem observabilidade. Este trabalho documenta e valida o ValorizeAI, aplicação real de categorização automática construída sobre Google Cloud Run, Cloud SQL, Redis (Memorystore), Cloud Tasks e WebSockets. A pesquisa adota abordagem aplicada inspirada em SRE: os SLOs (latência $P95 \leq 300$ ms, erro $< 0,5\%$, disponibilidade $\geq 99,5\%$) são definidos antecipadamente, a infraestrutura é provisionada como código e os testes k6 cobrem cenários de leitura intensiva, mix leitura/escrita e o pipeline assíncrono. O cenário de leitura sustentou 1,000 VUs (média de 470 req/s, pico de ≈ 970 req/s) com $P95 = 158$ ms; o cenário misto manteve 226 req/s em 650 VUs e violou o SLO em ≈ 539 VUs, elevando o $P95$ para 658 ms e o $p99$ para 2,67 s ao saturar a cota de 10 instâncias do Cloud Run. O pipeline assíncrono processou 51,58 mil tarefas do Cloud Tasks em 10 minutos (86 tarefas/s) sem perda ou duplicação. Conclui-se que a arquitetura atende confortavelmente workloads intensivos em leitura e que novas instâncias HTTP ou otimizações no caminho de escrita são necessárias para acomodar escritas altamente concorrentes.*

Palavras-Chave *Computação Serverless, Cloud Run, Plataformas Financeiras, Testes de Carga, Observabilidade.*

1. Introdução

Aplicações digitais modernas, que abrangem desde plataformas de e-commerce e serviços financeiros até mídias sociais e sistemas de colaboração em tempo real, enfrentam um desafio operacional comum: a gestão de cargas de trabalho (workloads) voláteis e imprevisíveis [Google Cloud 2024]. Picos de tráfego, eventos contínuos de ingestão de dados e a necessidade de múltiplas integrações com sistemas externos demandam uma infraestrutura que reaja dinamicamente, muito além da capacidade de provisionamento manual ou pipelines de integração monolíticos e rigidamente acoplados.

A resposta da indústria a esse desafio é a *elasticidade na nuvem* (cloud elasticity), definida como a capacidade de uma infraestrutura de computação em nuvem alocar e desalocar recursos computacionais de forma automática e autônoma, com base na flutuação da demanda em tempo real [Google Cloud 2024]. Diferente da escalabilidade tradicional, que muitas vezes envolve intervenção humana para provisionar novas instâncias, a elasticidade é projetada para lidar com picos abruptos e vales de tráfego, garantindo tanto o desempenho quanto a eficiência de custos [Google Cloud 2024]. Os benefícios de negócios são diretos: (1) *Eficiência de Custo*, ao adotar um modelo *pay-as-you-go* que evita o superprovisionamento dispendioso de recursos ociosos; e (2) *Alta Disponibilidade e Desempenho*, ao assegurar que a aplicação permaneça responsiva e confiável, preservando a experiência do usuário mesmo sob demanda extrema [Google Cloud 2024].

Contudo, a adoção de arquiteturas que possibilitam essa elasticidade — notadamente microsserviços, contêineres e paradigmas *serverless* [Basteri 2023] — introduz um nível exponencial de complexidade. Sistemas distribuídos são intrinsecamente mais difíceis de depurar e monitorar do que monólitos. Em resposta, a *observabilidade* (observability) emergiu como uma prioridade estratégica [Basteri 2023]. A observabilidade transcende o monitoramento tradicional (que rastreia falhas *conhecidas*) ao fornecer *insights* sobre o estado interno do sistema a partir de seus outputs (logs, métricas e traces), permitindo a depuração de falhas *desconhecidas* [Basteri 2023].

Elasticidade e observabilidade não são, portanto, características independentes, mas um ciclo de *feedback* simbiótico. A elasticidade automática gera uma complexidade que só pode ser gerenciada pela observabilidade. Por sua vez, a observabilidade fornece os dados (na forma de Indicadores de Nível de Serviço, ou SLIs) que alimentam e validam o próprio mecanismo de elasticidade, garantindo que o escalonamento automático atenda aos objetivos de negócio (SLOs) sem incorrer em custos desnecessários [Basteri 2023]. A falha de processos manuais é uma consequência direta da alta velocidade desse ciclo, que opera em milissegundos.

O ValorizeAI, objeto deste trabalho, nasce como um estudo de caso completo para investigar essa simbiose. Trata-se de uma aplicação web modular que centraliza fluxos complexos de ingestão de dados, processamento síncrono e assíncrono, e entrega de notificações e painéis em tempo real, utilizando uma *stack* de tecnologias modernas.

1.1. Justificativa e Problema de Pesquisa

Workloads transacionais que concentram ingestão massiva de dados, gerenciamento de estados compartilhados e interfaces colaborativas — como os simulados pelo ValorizeAI — impõem requisitos rigorosos de arquitetura. Tais sistemas exigem consistência forte nos dados, rastreabilidade completa para fins de auditoria e, fundamentalmente, respostas de baixa latência (respostas instantâneas), mesmo quando o tráfego varia abruptamente.

Para atender a esses requisitos, arquiteturas modernas combinam múltiplos padrões especializados:

- **CDNs e Balanceamento Global:** Para reduzir a latência de entrega de *assets* estáticos, distribuindo o conteúdo para pontos de presença (PoPs) próximos ao usuário [Barri 2025].
- **Filas Assíncronas (EDA):** Para desacoplar tarefas pesadas (ex: processamento de relatórios, envio de e-mails) da resposta síncrona, garantindo resiliência e escalabilidade através de uma arquitetura orientada a eventos [Confluent 2024].
- **Cache Distribuído (Redis):** Para armazenar dados frequentemente acessados em memória, reduzindo drasticamente a latência de leitura e a carga sobre o banco de dados principal [Yadav 2019].
- **WebSockets (Tempo Real):** Para comunicação bidirecional persistente, essencial para painéis e notificações em tempo real sem a sobrecarga do HTTP *polling* [Fernando e Engel 2025].

Embora existam tutoriais e artigos pontuais sobre cada uma dessas tecnologias, a justificativa deste trabalho reside na lacuna da literatura acadêmica e técnica [Christidis et al. 2022, Abad et al. 2021]. É raro encontrar material que conecte, de ponta a ponta (end-to-end), a implementação de uma arquitetura híbrida (CaaS + Filas + WebSockets) aos resultados práticos e reproduzíveis de testes de desempenho. A literatura existente tende a ser fragmentada, focando em comparações de ferramentas de IaC [Pessa 2023] ou na validação de microsserviços específicos [Hebbbar 2025], mas raramente no sistema holístico.

A contribuição deste TCC é, portanto, metodológica e empírica. Ao construir o ValorizeAI e registrar rigorosamente os experimentos de carga e processamento assíncrono, este trabalho evidencia como as decisões arquiteturais se traduzem em métricas quantificáveis (ex: latência P95, *throughput* de tarefas, taxa de erro). Este trabalho busca responder à seguinte questão: *Como uma arquitetura híbrida e elástica, composta por Cloud Run, Redis, Cloud SQL, Cloud Tasks e WebSockets dedicados, se comporta sob estresse de carga frente a SLOs definidos e como validar esse comportamento de modo reproduzível?*

Documentar esse caminho, apoiado por Infraestrutura como Código (IaC) para reprodutibilidade e monitoramento de custo, fornece uma referência concreta para equipes técnicas que precisam justificar e implementar arquiteturas orientadas a eventos com elasticidade horizontal automática.

1.2. Objetivo Geral

Demonstrar, por meio de documentação técnica e experimentos de desempenho, que a arquitetura do ValorizeAI — composta por balanceador com CDN, contêineres escalados horizontalmente, processamento assíncrono em filas, servidor de WebSockets, *buckets* para artefatos e cache em Redis — sustenta os SLOs definidos para um produto transacional completo, mantendo todo o ciclo (modelagem, desenvolvimento, infraestrutura, observabilidade e testes) versionado no repositório.

1.3. Objetivos Específicos

1. **Mapear a arquitetura end-to-end**, destacando o papel do balanceador/CDN, das instâncias de contêineres, do servidor de WebSockets, das filas assíncronas, dos *buckets* de armazenamento e do Redis para garantir consistência e baixa latência.
2. **Documentar o desenvolvimento** do *backend* Laravel, do *frontend* React e dos fluxos síncronos/assíncronos, com foco nos módulos críticos (ingestão de dados, automações, notificações e painéis em tempo real).

3. **Planejar e executar os testes de carga** (k6, cenários de leitura e de leitura/escrita) e o teste de processamento assíncrono para validar horizontalmente a arquitetura frente aos SLOs e identificar gargalos.
4. **Interpretar os resultados e propor otimizações**, relacionando desempenho, elasticidade e custo (ex.: ajustes de limites de instância, estratégias de cache) e apontando como essas evidências fundamentam decisões para *workloads* transacionais de alta criticidade.

1.4. Contribuições Tangíveis

1. **Arquitetura documentada:** descreve-se o desenho completo do ValorizeAI — balanceador/CDN, Cloud Run, Redis, Cloud SQL, Cloud Tasks e Reverb — para que outros pesquisadores possam replicar o raciocínio arquitetural em aplicações transacionais.
2. **Infraestrutura reprodutível:** disponibiliza-se o conjunto de automações (Terraform, Docker, Makefile) que reconstrói fielmente o ambiente avaliado, mantendo a linha do tempo de mudanças de infraestrutura junto ao código-fonte.
3. **Experimentos de desempenho:** registra-se integralmente os cenários k6 (leituras e mix), com ramp-ups, parâmetros e métricas coletadas (latências, throughput, erros), assegurando transparência na validação dos SLOs.
4. **Teste assíncrono:** documenta-se o ensaio do pipeline baseado em Cloud Tasks e workers HTTP, explicitando como o sistema responde ao acúmulo e drenagem de tarefas.

Em conjunto, essas contribuições formam um pacote de replicação completo — código, infraestrutura e experimentos — para avaliações futuras de *workloads* transacionais em plataformas CaaS.

1.5. Estrutura do Trabalho

A organização segue esta lógica: a **Seção 2** revisa a literatura sobre paradigmas de execução em nuvem, padrões arquiteturais e metodologias de validação, situando a lacuna investigada. A **Seção 3** consolida os conceitos teóricos e técnicos que sustentam o design do sistema e os experimentos. A **Seção 4** descreve a abordagem experimental, do planejamento dos SLOs ao uso de k6 e Cloud Monitoring para coleta de evidências. As seções seguintes detalham a implementação da arquitetura ValorizeAI, discutem os resultados obtidos e encerram com conclusões e direções futuras.

2. Trabalhos Relacionados

A arquitetura proposta pelo ValorizeAI situa-se na interseção de três domínios de pesquisa em engenharia de software e sistemas distribuídos: (1) os paradigmas de execução em nuvem, (2) os padrões de design para resiliência e desempenho, e (3) as metodologias de validação empírica. Revisa-se o estado da arte em cada um desses eixos para posicionar a contribuição deste trabalho e fundamentar a lacuna de pesquisa identificada na Introdução.

2.1. Paradigmas de Execução: Serverless (FaaS) vs. Contêineres Gerenciados (CaaS)

A primeira decisão de design em arquiteturas elásticas modernas é a escolha do paradigma de computação. A literatura recente concentra-se no debate entre *Functions-as-a-Service* (FaaS) e *Containers-as-a-Service* (CaaS).

O paradigma FaaS, popularizado por serviços como AWS Lambda e Google Cloud Functions, abstrai completamente o gerenciamento de servidores, oferecendo um modelo de faturamento por execução e escalabilidade instantânea (incluindo *scale-to-zero*)

[Sonawane et al. 2024]. Esse modelo é ideal para *workloads* reativos, *stateless* e de curta duração. No entanto, a literatura aponta desafios significativos: (1) a latência de inicialização (cold start), que pode impactar o desempenho de aplicações sensíveis à latência [Sonawane et al. 2024]; (2) a complexidade de monitoramento e observabilidade em um ambiente altamente efêmero [Sonawane et al. 2024]; e (3) a inadequação para processos *stateful* ou de longa duração, como conexões de banco de dados persistentes ou servidores WebSocket [Datadog 2024].

Em contrapartida, o CaaS, exemplificado por plataformas como Google Cloud Run e AWS Fargate, emerge como um meio-termo estratégico [Lloyd et al. 2018]. O CaaS combina a elasticidade e o modelo de *scale-to-zero* do FaaS com a portabilidade, consistência e controle de ambiente fornecidos pelos contêineres (ex: Docker) [Datadog 2024].

A arquitetura do ValorizeAI requer explicitamente um servidor de WebSockets dedicado (Reverb) para comunicação em tempo real — um processo *stateful* e de longa duração. A revisão da literatura [Datadog 2024, Sonawane et al. 2024] demonstra que o FaaS é um paradigma inadequado para esse requisito. O CaaS (especificamente o Cloud Run) foi, portanto, escolhido por ser o paradigma que permite a execução de processos persistentes (o contêiner do Reverb) enquanto ainda fornece a elasticidade horizontal automática e a abstração de infraestrutura desejadas para os serviços web *stateless*.

2.2. Padrões Arquiteturais para Desempenho e Resiliência

Para atender aos requisitos de um sistema transacional em tempo real, o ValorizeAI combina padrões de comunicação síncronos e assíncronos.

2.2.1. Comunicação Assíncrona e Arquiteturas Orientadas a Eventos (EDA)

O processamento assíncrono por meio de filas, um pilar central do ValorizeAI, é a implementação prática de uma Arquitetura Orientada a Eventos (EDA). EDAs são definidas como sistemas que promovem o desacoplamento (loose coupling), a escalabilidade e a resiliência [Confluent 2024]. Ao utilizar um *message broker* (como RabbitMQ, Kafka ou serviços gerenciados como Google Cloud Tasks), os serviços "produtores" podem enfileirar tarefas (eventos) sem esperar que os "consumidores" as processem [Confluent 2024]. Isso permite que o sistema absorva picos de escrita e mantenha a responsividade da interface do usuário, além de garantir a entrega de tarefas mesmo que os serviços consumidores falhem temporariamente.

A pesquisa acadêmica neste domínio frequentemente se concentra em análises de desempenho comparativas dos *brokers* de mensagens. Por exemplo, estudos comparam o desempenho de RabbitMQ, Apache Kafka e Apache Pulsar em cenários de IoT, medindo *throughput* e latência sob diferentes tamanhos de mensagem [Thepphakan 2025]. Esses estudos validam que a escolha da tecnologia de fila deve estar alinhada com os requisitos específicos do *workload* (ex: baixa latência para mensagens pequenas vs. alto *throughput* para *streams* de dados).

2.2.2. Comunicação em Tempo Real (WebSockets) e Cache Distribuído (Redis)

Para os requisitos de "painéis em tempo real" e "notificações instantâneas", a arquitetura do ValorizeAI utiliza Laravel Reverb e Redis. A literatura não trata desses componentes isoladamente, mas sim como um padrão arquitetural combinado para escalar aplicações em tempo real.

O Redis é amplamente citado por seu papel como um *cache* distribuído em memória, fornecendo acesso a dados de baixa latência e alto *throughput*, o que reduz a carga sobre bancos de dados relacionais [Yadav 2019]. Servidores WebSocket (como o Reverb [Laravel Holdings Inc. 2025]) fornecem o canal de comunicação bidirecional persistente necessário para que o servidor envie dados aos clientes sem que eles precisem solicitá-los (push) [Twine 2022].

O desafio de escalar WebSockets reside na sua natureza *stateful* (o servidor deve manter o registro de cada conexão ativa). Em um ambiente de CaaS elástico como o Cloud Run, onde múltiplas instâncias *stateless* são criadas e destruídas dinamicamente, uma conexão WebSocket estabelecida com a "Instância A" não pode ser acessada pela "Instância B". A literatura e a documentação técnica [Laravel Holdings Inc. 2025] resolvem isso usando o mecanismo de Publicação/Subscrição (Pub/Sub) do Redis como um *backplane* de mensagens. Quando a Instância B precisa enviar uma mensagem para um usuário conectado à Instância A, ela publica a mensagem no canal Redis. A Instância A, que está inscrita (subscribed) nesse canal, recebe a mensagem e a retransmite ao cliente correto através de sua conexão WebSocket local.

Estudos de desempenho de bibliotecas WebSocket, como o de Fernando e Engel [Fernando e Engel 2025], validam a importância de escolhas de implementação leves, focando em métricas-chave como *throughput* (mensagens/segundo) e latência de *Round Trip Time* (RTT) para garantir o desempenho em tempo real.

2.3. Metodologias de Validação Empírica

Provar que uma arquitetura complexa atende aos seus requisitos de desempenho exige uma metodologia de validação rigorosa e, idealmente, reproduzível.

2.3.1. Infraestrutura como Código (IaC) para Reprodutibilidade

A Infraestrutura como Código (IaC) é uma prática de DevOps onde a infraestrutura de TI (redes, máquinas virtuais, balanceadores de carga) é provisionada e gerenciada usando arquivos de definição legíveis por máquina (ex: Terraform, AWS CDK), em vez de configuração manual [Pessa 2023]. No contexto da pesquisa acadêmica e de engenharia, o principal benefício do IaC é a *reprodutibilidade*. Ao versionar a configuração da infraestrutura juntamente com o código da aplicação, o IaC garante que o ambiente de teste possa ser recriado de forma consistente, eliminando a "deriva de configuração" e tornando os resultados dos testes de desempenho verificáveis [Guerriero et al. 2019].

A literatura sobre IaC, como o estudo de Pessa [Pessa 2023], muitas vezes foca na comparação das próprias ferramentas de IaC (ex: AWS CDK vs. Terraform) em termos de desempenho de provisionamento e experiência do desenvolvedor, em vez de usar o IaC como um *meio* para validar o desempenho da *aplicação* que ele provisiona.

2.3.2. Validação de SLOs com Testes de Carga (k6)

A metodologia do ValorizeAI baseia-se nos princípios de Engenharia de Confiabilidade de Sites (SRE), onde o sucesso é medido pelo cumprimento dos Objetivos de Nível de Serviço (SLOs) [McCoy e Forsgren 2020]. A validação de SLOs requer testes empíricos sob carga.

A literatura acadêmica recente começa a adotar ferramentas de teste de carga modernas, como o k6 [Cervone 2024], para essa finalidade. Um exemplo notável é o trabalho de Hebbar [Hebbar 2025] sobre APIs reativas para o setor financeiro. Hebbar utiliza o k6 para criar perfis de carga (ex: rajada, estado estacionário) e simular tráfego contra um microserviço Spring WebFlux [Hebbar 2025]. A contribuição desse estudo é a validação de que a arquitetura consegue aplicar priorização de tráfego (tiering de SLA) em tempo real, medindo métricas de latência, taxa de descarte e saturação [Hebbar 2025]. Este estudo serve como um "espelho" metodológico, validando a abordagem do ValorizeAI (uso de k6 para medir métricas de latência P95 contra SLOs definidos) como academicamente rigorosa e alinhada com o estado da arte da pesquisa em desempenho de sistemas.

2.4. Síntese da Revisão e Identificação da Lacuna

A revisão da literatura revela que a pesquisa é frequentemente especializada e fragmentada. Encontramos estudos que comparam FaaS vs. CaaS [Lloyd et al. 2018], analisam o desempenho de *brokers* de EDA [Thepphakan 2025], comparam ferramentas de IaC [Pessa 2023], ou validam um microserviço específico usando k6 e SLOs [Hebbar 2025].

A lacuna na literatura, identificada em trabalhos como [Christidis et al. 2022] e [Abad et al. 2021], é a ausência de estudos de caso *end-to-end* que integrem *todos* esses componentes. Falta um trabalho que documente e valide empiricamente uma arquitetura holística e híbrida (CaaS + EDA + WebSockets + Cache) que seja:

1. Provisionada de forma reproduzível (via IaC).
2. Validada rigorosamente contra SLOs de latência e *throughput* (via k6).
3. Analisada em seus múltiplos componentes (fluxos síncronos e assíncronos).

A Tabela 1 visualiza essa lacuna. Enquanto trabalhos anteriores focam em colunas específicas, este TCC (ValorizeAI) é o único que propõe uma validação integrada de todos os eixos: Paradigma, Padrões Híbridos e Metodologia de Validação Completa. Este trabalho preenche, assim, a lacuna ao fornecer um "plano" de arquitetura e validação, completo e empiricamente verificado, para aplicações transacionais modernas em tempo real.

Tabela 1. Quadro Comparativo de Estudos sobre Desempenho de Arquiteturas em Nuvem (2018-2025)

Estudo (Autor)	Paradigma	Padrões Analisados	Metodologia de Validação
Lloyd et al. [Lloyd et al. 2018]	FaaS vs. CaaS	Fatores de desempenho em microserviços simples.	Benchmarking de desempenho (latência, custo). Não foca em IaC ou SLOs formais.
Pessa [Pessa 2023]	N/A (Foco na ferramenta)	Provisionamento de infraestrutura (FaaS, CaaS).	Comparação de ferramentas de IaC (CDK vs. Terraform). Não valida desempenho da aplicação.
Hebbar [Hebbar 2025]	Microserviço (Monolítico)	API Reativa (Spring WebFlux) com priorização.	SLOs e k6. Não foca em IaC, EDA ou WebSockets.

Tabela 1 – Continuação

Estudo (Autor)	Paradigma	Padrões Analisados	Metodologia de Validação
Thepphakan [Thepphakan 2025]	N/A (Foco no broker)	EDA (RabbitMQ vs. Pulsar).	Benchmarking de desempenho (latência, <i>throughput</i>). Não é um sistema E2E.
Abad et al. [Abad et al. 2021]	FaaS / Serverless	Revisão de aplicações serverless.	Análise de literatura. Aponta a lacuna em estudos E2E e workflows complexos.
ValorizeAI (este TCC)	CaaS híbrido com Cloud Run e Cloud Tasks	Arquitetura E2E combinando EDA (Pub/Sub), WebSockets e cache Redis com automações transacionais.	Validação reprodutível via Terraform + k6 + observabilidade (SLOs, logs e filas).

Os fatores comparados deixam claro que a literatura cobre fragmentos isolados do problema. Esses referenciais serão consolidados, a seguir, em um vocabulário comum (Clean Architecture, DDD, SRE, EDA) que sustenta a proposta do ValorizeAI e prepara o terreno para a metodologia e a implementação.

3. Fundamentação Teórica

Segue uma síntese do vocabulário e das bases conceituais utilizados no design, implementação e validação do ValorizeAI, cobrindo princípios de design de software, arquitetura dos componentes e fundamentos de engenharia de confiabilidade.

3.1. Princípios de Design de Software

O ValorizeAI adota uma abordagem de "arquitetura limpa", segregando responsabilidades com base em princípios estabelecidos de design de software.

3.1.1. Clean Architecture

Formalizada por Robert C. Martin, a *Clean Architecture* (Arquitetura Limpa) é um modelo arquitetural que advoga pela separação de interesses [Martin 2017]. Seu objetivo é criar sistemas que sejam: (1) Independentes de frameworks; (2) Testáveis; (3) Independentes da interface do usuário (UI); e (4) Independentes do banco de dados [Martin 2017].

O pilar central dessa arquitetura é a *Regra da Dependência* (The Dependency Rule). Esta regra estipula que as dependências do código-fonte devem apontar exclusivamente "para dentro— de camadas de baixo nível (detalhes voláteis, como frameworks e bancos de dados) para camadas de alto nível (políticas de negócio estáveis e abstrações) [Martin 2017]. No ValorizeAI, isso se manifesta na separação das regras de negócio (localizadas em *Actions* ou *Queries*) da lógica do framework (Controladores Laravel) ou da persistência (Modelos Eloquent).

3.1.2. Domain-Driven Design (DDD)

O *Domain-Driven Design* (DDD), introduzido por Eric Evans, é uma abordagem para o desenvolvimento de software que se concentra em modelar o software para corresponder a um domínio de negócio complexo [Evans 2003]. O DDD é essencial para gerenciar a complexidade em sistemas como o ValorizeAI. Os conceitos-chave utilizados neste trabalho incluem:

- **Linguagem Ubíqua (Ubiquitous Language):** Um vocabulário compartilhado e rigoroso, desenvolvido em colaboração entre os desenvolvedores e os especialistas do domínio (usuários). Essa linguagem é usada em todas as comunicações e reflete-se diretamente no código (nomes de classes, métodos e variáveis) [Evans 2003].
- **Contexto Delimitado (Bounded Context):** A fronteira explícita dentro da qual um modelo de domínio e sua Linguagem Ubíqua são aplicáveis e consistentes [Evans 2003].
- **Agregado (Aggregate):** Um cluster de objetos de domínio (Entidades e Objetos de Valor) que é tratado como uma única unidade para fins de consistência de dados. Um Agregado possui uma raiz (a *Aggregate Root*), que é o único ponto de entrada para modificações dentro do Agregado, garantindo que todas as regras de negócio (invariantes) sejam aplicadas [Evans 2003].

3.1.3. Padrões de Comunicação e Segregação

Para implementar a Regra da Dependência e gerenciar o fluxo de dados, o ValorizeAI utiliza padrões de segregação e transferência de dados.

- **DTO (Data Transfer Object):** Conforme popularizado por Martin Fowler, um DTO é um objeto simples, sem comportamento, cujo único propósito é transferir dados entre subsistemas ou camadas [Fowler 2002]. Em arquiteturas distribuídas ou em camadas, os DTOs são usados para agregar múltiplas chamadas em uma única, reduzindo a latência da rede e desacoplando os modelos internos (domínio) dos modelos de visualização (API/UI).
- **CQRS (Command Query Responsibility Segregation):** Um padrão, descrito por Martin Fowler [Fowler 2011] e Greg Young, que propõe a segregação dos modelos de dados e da lógica de aplicação em duas categorias: *Commands* (operações que alteram o estado, ou seja, escritas) e *Queries* (operações que leem o estado). O ValorizeAI adota esse princípio através da separação explícita de *Actions* (Commands) e *Queries* (Queries), permitindo otimizações distintas para os caminhos de escrita e leitura.

Nas rotas de leitura (consultas de transações, contas e categorias), os controladores invocam exclusivamente *Queries* especializadas que exploram o cache Redis e projetam os DTOs esperados pelo front-end. Já as operações de escrita (POST `/api/transactions`, importações em lote, automações) são canalizadas por *Actions* que encapsulam validações, regras de negócio e persistência, utilizando comandos explícitos que podem ser monitorados ou reprocessados. Essa separação facilita ajustes independentes: otimizações de leitura não afetam o caminho de escrita e vice-versa, e cada fluxo abriga métricas distintas no Cloud Monitoring.

3.2. Arquitetura e Componentes da Aplicação

A infraestrutura do ValorizeAI é composta por serviços gerenciados na nuvem, escolhidos por suas características de elasticidade e desempenho.

3.2.1. Google Cloud Run e Cloud Tasks

O Google Cloud Run é uma plataforma de computação CaaS (Container-as-a-Service) totalmente gerenciada. Ele permite a execução de contêineres *stateless* que escalam horizontalmente de forma automática, com a capacidade de escalar até zero instâncias quando não há tráfego, eliminando custos ociosos [Google Cloud 2024]. O serviço foi escolhido por combinar a elasticidade típica de funções serverless com a flexibilidade dos contêineres, executando tanto os serviços web *stateless* do ValorizeAI quanto o servidor *stateful* de WebSockets.

O Google Cloud Tasks é o serviço de enfileiramento de tarefas gerenciado. Ele é usado para implementar o processamento assíncrono (EDA), permitindo que a aplicação principal (síncrona) enfileire tarefas de longa duração (ex: processamento de lotes) para execução em *workers* separados, garantindo resiliência e baixa latência na resposta ao usuário.

3.2.2. Laravel Reverb (WebSockets)

O Laravel Reverb é o servidor WebSocket oficial de primeira-parte para aplicações Laravel, projetado para comunicação em tempo real de alto desempenho [Laravel Holdings Inc. 2025]. Ele utiliza o protocolo Pusher, integrando-se nativamente ao sistema de *broadcasting* do Laravel para facilitar o envio de notificações *push* aos clientes conectados.

A característica arquitetural mais importante do Reverb para este TCC é seu suporte à escalabilidade horizontal. Para operar em um ambiente elástico como o Cloud Run (com múltiplas instâncias de servidor), o Reverb utiliza um *backplane* de mensagens, que no caso do ValorizeAI é implementado com o Redis (detalhado na subseção sobre cache) [Laravel Holdings Inc. 2025].

3.2.3. Redis (Remote Dictionary Server)

O Redis (Remote Dictionary Server) é um armazenamento de estrutura de dados em memória, de código aberto, usado como banco de dados, *cache* e *message broker* [Kleppmann 2017]. No contexto da arquitetura ValorizeAI, o Redis desempenha dois papéis críticos e distintos, ambos fundamentais para o desempenho do sistema:

1. **Cache de Baixa Latência:** O Redis é usado como um *cache* para dados frequentemente acessados (ex: painéis, dados de sessão). Sua operação em memória permite latências de leitura e escrita na ordem de submilissegundos, reduzindo drasticamente a carga sobre o banco de dados PostgreSQL e melhorando a responsividade das *Queries* [Yadav 2019].
2. **Backplane Pub/Sub:** O Redis fornece um mecanismo de Publicação/Subscrição (Pub/Sub) de alto desempenho. Este mecanismo é utilizado como o *backplane* do Laravel Reverb. Quando uma instância do servidor (Instância A) precisa notificar um usuário que está conectado via WebSocket a outra instância (Instância B), a Instância A publica a mensagem em um canal Redis. Todas as outras instâncias, incluindo a Instância B, estão inscritas nesse canal, recebem a mensagem e a retransmitem aos seus clientes WebSocket conectados localmente.

3.3. Engenharia de Confiabilidade de Sites (SRE)

A metodologia de validação deste trabalho é baseada nos princípios de Engenharia de Confiabilidade de Sites (SRE), popularizados pelo Google [Beyer et al. 2016]. O SRE trata as

operações de infraestrutura como um problema de engenharia de software, utilizando métricas rigorosas para equilibrar a inovação (velocidade de desenvolvimento) com a confiabilidade do serviço.

3.3.1. SLIs, SLOs e Orçamentos de Erro

Os conceitos centrais do SRE utilizados para a validação do ValorizeAI são:

- **SLI (Service Level Indicator):** Um indicador de nível de serviço é uma medida quantitativa de um aspecto da qualidade do serviço fornecido [McCoy e Forsgren 2020]. Os SLIs são métricas diretas do desempenho do sistema, como latência de requisição, taxa de erro ou *throughput* do sistema [Beyer et al. 2016].
- **SLO (Service Level Objective):** Um objetivo de nível de serviço é um valor-alvo ou um intervalo de valores para um SLI, medido ao longo de um período [McCoy e Forsgren 2020]. Um SLO é a definição formal de "quão bom" o serviço precisa ser. Por exemplo, "95% das requisições de leitura (SLI: latência de leitura) devem ser concluídas em menos de 250ms (SLO) nos últimos 28 dias".
- **Orçamento de Erro (Error Budget):** O orçamento de erro é o complemento do SLO (ou seja, $100\% - SLO\%$) [Beyer et al. 2016]. Ele representa a quantidade de falhas "permitidas" (ex: requisições lentas ou com erro) durante o período. O orçamento de erro é uma ferramenta de gerenciamento: enquanto houver orçamento, a equipe de desenvolvimento tem "permissão" para lançar novas funcionalidades (que inerentemente trazem risco); se o orçamento se esgotar, o foco da equipe deve mudar para a melhoria da confiabilidade [Beyer et al. 2016].

Esses fundamentos orientam a abordagem metodológica detalhada na Seção 4, que explica como o planejamento dos SLOs, a infraestrutura como código e os experimentos com k6 e Cloud Tasks foram conduzidos para gerar as evidências analisadas posteriormente.

4. Metodologia

O estudo foi conduzido de ponta a ponta, do planejamento dos objetivos de nível de serviço (SLOs) à coleta e interpretação dos experimentos. O enfoque é aplicado e experimental: toda a instrumentação foi construída diretamente no repositório ValorizeAI, o que permite a reprodução dos resultados.

4.1. Tipo de Pesquisa e Estratégia Geral

O trabalho caracteriza-se como uma **pesquisa aplicada** conduzida como **estudo de caso** de um sistema real em produção. A estratégia seguiu quatro fases iterativas. No **planejamento**, foram definidos os SLOs (latência P95 de 300 ms, erro $<0,5\%$, disponibilidade $\geq 99,5\%$), mapeadas as cotas vigentes do Cloud Run (10 instâncias de 1 vCPU / 1 GiB, totalizando 10 vCPU) e estimado como essa limitação poderia afetar o *throughput* máximo almejado, o que orientou as cargas aplicadas. Em seguida veio a **preparação do ambiente**: módulos Terraform provisionaram rede, bancos e serviços gerenciados; Docker Compose reproduziu localmente PostgreSQL, Redis e a stack de observabilidade; o Makefile encapsulou tarefas de lint, testes e execução dos cenários. Na etapa de **execução controlada**, os cenários k6 de leitura e leitura/escrita foram disparados contra a API em Cloud Run enquanto o pipeline assíncrono recebia um lote adicional de tarefas no Cloud Tasks, exercitando os workers HTTP. Por fim, na **coleta e análise**, as métricas agregadas (latência, *throughput*, taxa de erro) foram extraídas dos CSVs e painéis do Cloud Monitoring, e as

observações qualitativas sobre o teste de filas foram registradas juntamente com o tempo total de drenagem, subsidiando os capítulos de implementação e resultados.

4.2. Arquitetura do Ambiente Experimental

A Figura 1 sintetiza os componentes usados nos experimentos. O tráfego HTTP/HTTPS entra por um **Cloud Load Balancer** com **Cloud CDN**, que reduz a latência de *assets* estáticos e protege o backend com inspeção WAF. Esse tráfego é encaminhado para dois serviços Cloud Run:

- **API Laravel**: processa requisições REST, expõe endpoints usados pelos testes k6 e orquestra o pipeline assíncrono.
- **Laravel Reverb**: mantém conexões WebSocket persistentes para eventos em tempo real; é tratado como serviço independente para permitir escalonamento específico.

Ambos os serviços acessam o **Memorystore for Redis**, usado simultaneamente como cache de leitura (padrão *cache-aside*) e como *backplane* Pub/Sub do Reverb. O armazenamento transacional permanece no **Cloud SQL for PostgreSQL**, que atende às operações de leitura e escrita executadas durante os testes. Para workloads assíncronos, a API publica tarefas em **Cloud Tasks**, que aciona workers HTTP também hospedados no Cloud Run. Diferentemente de filas baseadas em *polling*, o Cloud Tasks opera no modo *push*: quando novas tarefas são registradas, o serviço dispara requisições para os endpoints configurados e o Cloud Run instancia containers apenas enquanto houver demanda. Isso elimina processos *long living* dedicados a escutar filas e garante elasticidade automática durante rajadas. O **Cloud Storage**, embora componha o ambiente produtivo (armazenamento de extratos e relatórios), não foi exercitado nos testes de carga — a instrumentação concentrou-se nos caminhos transacionais da API e no pipeline assíncrono.

Para refletir o ambiente em produção, o Cloud SQL foi configurado como instância personalizada *db-perf-optimized-N-2* (2 vCPU e 16 GiB de RAM), garantindo buffers suficientes para os planos complexos das consultas analíticas sem alterar o perfil de custos do sistema.

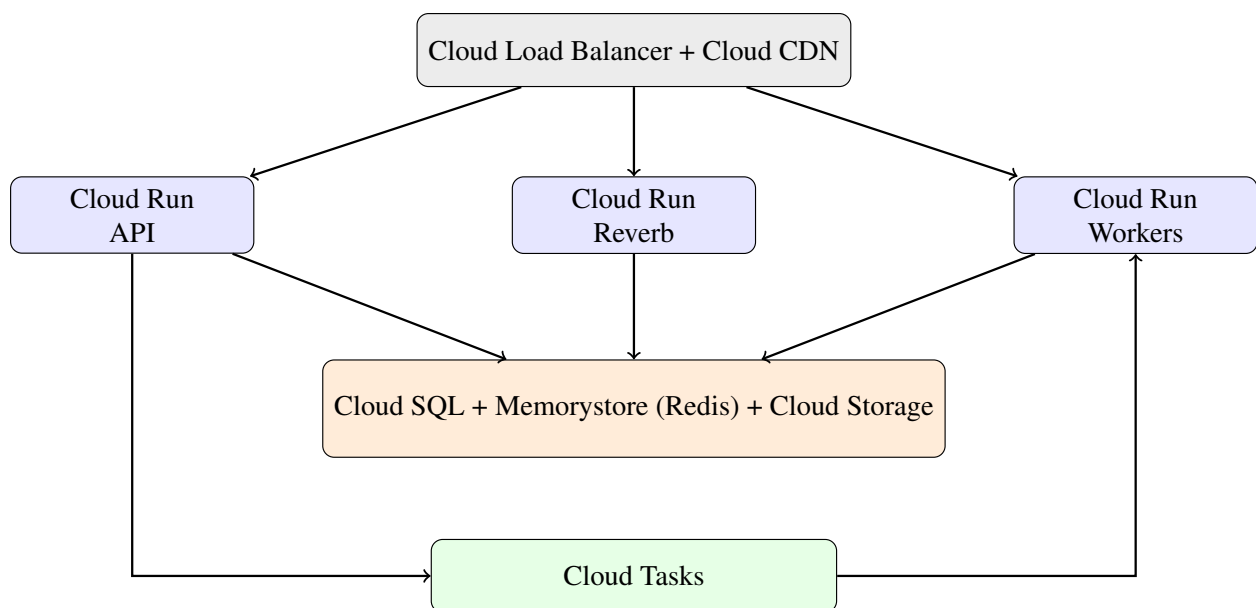


Figura 1. Arquitetura utilizada nos experimentos.

4.3. Desenvolvimento da Aplicação e Modelo de Dados

O desenvolvimento do ValorizeAI avançou em paralelo à infraestrutura. O *backend* Laravel concentra o domínio financeiro (contas, orçamentos, categorizações) e expõe APIs REST e eventos

WebSocket consumidos pelo painel React. Essa base inclui ingestão automática de extratos (OFX/CSV), reconciliação assistida por filas, categorização por embeddings e automações que propagam atualizações de saldo e metas em tempo real. Os workers HTTP são os responsáveis por processar os pipelines mais caros (importação de transações, geração de relatórios, envio de notificações) e publicam eventos para o Reverb, fechando o ciclo com a interface.

O modelo relacional, sintetizado na Figura 2, reflete os invariantes do domínio multi-inquilino. Usuários controlam múltiplas contas vinculadas a bancos e moedas distintas, cada conta agrega transações com divisão por categorias, e os orçamentos mensais impõem limites e alocações que retroalimentam o painel de metas. Camadas de auditoria (logs de importação, feedback de categorização e embeddings vetoriais em pgvector) conectam a coleta de dados com os recursos de IA usados para sugerir classificações.

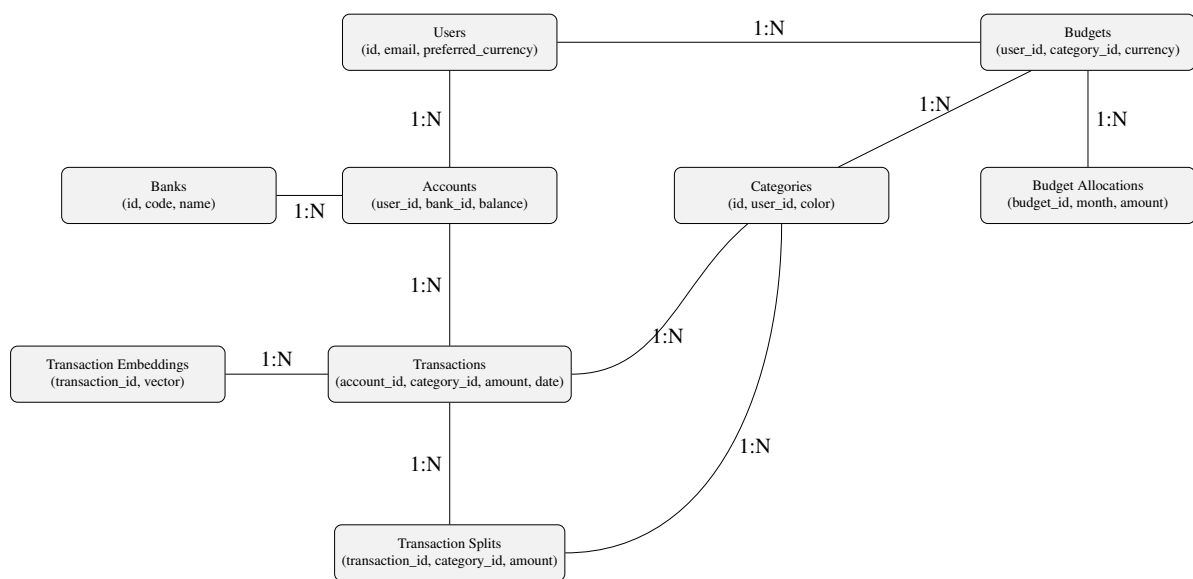


Figura 2. Modelo lógico central derivado do esquema relacional do ValorizeAI.

4.4. Ferramentas e Processo de Preparação

Do ponto de vista de engenharia, três pilares garantiram a reprodutibilidade: **(i) Infraestrutura como Código:** os módulos Terraform descrevem VPC, balanceadores, Cloud Run, Cloud SQL, Redis e Cloud Tasks. Cada mudança passa por *plan/apply* versionado, evitando deriva de ambiente. **(ii) Ambientes determinísticos:** o Makefile e os manifestos Docker recompõem o stack local (PostgreSQL, Redis e PHP 8.4) idêntico ao ambiente de teste antes de qualquer execução k6. **(iii) Observabilidade:** o Cloud Monitoring consolidou métricas de latência, uso de CPU/memória e backlog das filas. A adoção de OpenTelemetry e Loki/Tempo permanece como trabalho futuro.

A camada de infraestrutura como código segue um arranjo modular: cada elemento da arquitetura (Cloud SQL, Memorystore, Cloud Run, Cloud Tasks e o balanceador com CDN) é definido como módulo independente e parametrizado com as mesmas variáveis de projeto, região e rótulos. Esses módulos expõem *outputs* que encadeiam os serviços — por exemplo, o endereço privado do banco é injetado na definição do serviço Cloud Run, e o host/porta do Redis alimenta tanto a API quanto o servidor de WebSockets. A malha também contempla recursos de suporte, como a VPC dedicada ao tráfego serverless, a conexão ao Service Networking para disponibilizar IPs privados e os segredos no Secret Manager (credenciais de banco, tokens de API, chaves de serviço). Dependências explícitas garantem que os segredos e o emparelhamento de rede estejam

prontos antes do provimento dos contêineres, evitando condições de corrida. Como resultado, o ambiente completo pode ser criado, atualizado ou destruído com um único *apply*, assegurando rastreabilidade das mudanças de infraestrutura no mesmo repositório do código-fonte.

4.5. Planejamento dos SLOs e Desenho dos Cenários

Com base nas premissas de negócio e na literatura de SRE [McCoy e Forsgren 2020, Beyer et al. 2016], o sistema foi avaliado contra três metas: latência $P95 \leq 300$ ms, taxa de erro $< 0,5\%$ e disponibilidade mensal $\geq 99,5\%$. Essa abordagem segue a disciplina de SLOs/SLIs/Orçamentos de Erro descrita por Beyer et al. e McCoy & Forsgren, na qual o desempenho é continuamente confrontado com metas explícitas e o “espaço” para experimentação deriva do orçamento de erro restante. A cota vigente do Cloud Run (10 instâncias de 1 vCPU/1 GiB) limita o total de CPU disponível; no nosso cenário isso significou que os testes deveriam aumentar a carga até consumir essas 10 vCPU — ambos os cenários chegaram ao teto de 10 instâncias ativas, permitindo documentar o comportamento imediatamente antes do esgotamento.

Dois cenários foram modelados:

1. **Leitura intensiva:** 1.000 usuários virtuais consultando listas de transações por 17 minutos em seis estágios, exercitando cache Redis + réplica de leitura do PostgreSQL.
2. **Mistura leitura/escrita:** 650 usuários virtuais alternando consultas e criação de transações durante 21 minutos, forçando locks no banco e pressionando o pipeline de escrita.

Além desses ensaios HTTP, foi planejado um **teste de filas** no qual um volume elevado de tarefas artificiais percorre o fluxo Cloud Tasks → workers HTTP, permitindo observar o tempo de drenagem e a elasticidade dos consumidores assíncronos.

O roteiro de leitura replica o padrão de navegação predominante no produto: cada VU monta filtros distintos para o endpoint `GET /api/transactions`, alternando entre contas e categorias, e mantém um *think time* de 1 s para simular o tempo de leitura na interface. As etapas de carga (150→1000 VUs) formam uma curva contínua que aquece a cache, sustenta o platô de 600–900 VUs e aplica um pico curto de 1.000 VUs apenas para identificar o ponto de saturação com todas as instâncias do Cloud Run ativas.

O roteiro misto introduz um endpoint de provisionamento que entrega um token e dados exclusivos para cada VU, impedindo que múltiplos usuários escrevam na mesma conta. O fluxo executa leituras em 65% das iterações, `POST /api/transactions` em 20% e `GET /api/accounts` no restante, com cargas crescendo até 650 VUs. Esse percentual reproduz a relação leitura/escrita observada nos logs reais, enquanto o estágio final concentra o esforço em operações que seguram locks no PostgreSQL.

Para o ensaio assíncrono, um script independente publica 51,58 mil tarefas em lotes para o Cloud Tasks. Cada inserção ativa automaticamente os workers HTTP no Cloud Run, que escalam horizontalmente conforme o backlog. As métricas coletadas (vazão em tarefas/min e tempo de drenagem) servem para validar o comportamento da arquitetura orientada a eventos.

4.6. Execução dos Experimentos

Cada rodada segue os passos:

1. **Preparação dos dados:** seeds e factories povoam o PostgreSQL com contas, transações e orçamentos realistas; a instância Redis é pre-aquecida com métricas e dashboards frequentes.

2. **Disparo do cenário:** os perfis do k6 focados em leitura e no mix leitura/escrita são executados via Makefile, apontando para o domínio público do Cloud Load Balancer; estágios, VUs e SLIs monitorados seguem o planejamento experimental.
3. **Registro automático:** os resultados agregados são gravados em CSVs (latência, taxa de erro, uso de VUs) e correlacionados com as métricas de infraestrutura capturadas pelo Cloud Monitoring.
4. **Teste de filas:** um script HTTP produz um lote adicional de tarefas e a drenagem é acompanhada por meio das métricas do Cloud Tasks e dos logs dos workers.

4.7. Coleta e Integração das Evidências

As evidências produzidas sustentam as análises de arquitetura, implementação e resultados:

- **Planilhas de latência e throughput:** derivadas dos CSVs exportados pelo k6, utilizadas posteriormente para comparar as métricas observadas com os SLOs.
- **Series temporais de infraestrutura:** capturas dos dashboards do Cloud Monitoring registram uso de CPU das instâncias Cloud Run, saturação do Redis e backlog do Cloud Tasks durante cada rodada.
- **Relatos de execução:** cada rodada é registrada em um diário experimental (mantido no repositório, quando permitido pela instituição) com horários, parâmetros e observações qualitativas sobre o comportamento do sistema.

Essa metodologia garante rastreabilidade completa entre arquitetura, implementação e resultados, pois cada passo experimental está ancorado em artefatos versionados do projeto.

5. Resultados e Discussão

Esta seção consolida as métricas obtidas nos testes de carga (k6) e no ensaio assíncrono com Cloud Tasks. As medições seguem os SLIs definidos na metodologia: latência P95, taxa de erro e tempo de drenagem do pipeline assíncrono.

Tabela 2. Resumo dos cenários de carga e conformidade com os SLOs

Cenário	Latência P95	Throughput médio	Throughput pico	Taxa de erro
Leitura intensiva	158 ms	470 req/s	970 req/s	0,00%
Mistura leitura/escrita	658 ms	226 req/s	450 req/s	0,00%

Os valores de *throughput* foram obtidos diretamente dos dashboards do Cloud Run para cada estágio de VUs; por exemplo, o platô de 650 VUs no cenário misto correspondeu a aproximadamente 226 requisições por segundo, enquanto o patamar final de 900→1000 VUs no cenário de leitura atingiu pico de cerca de 970 req/s (a média global do teste ficou em 470 req/s). Assim, VUs e req/s descrevem o mesmo nível de pressão gerada nos serviços.

5.1. Cenário de Leitura Intensiva

Os 1.000 usuários virtuais do cenário de leitura mantiveram o backend em média de 470 requisições por segundo durante 17 minutos, com pico aproximado de 970 req/s no platô final (900→1000 VUs). A latência P95 permaneceu em 158 ms (Tabela 2), confortavelmente dentro do SLO de 300 ms mesmo nesse pico, e nenhuma requisição falhou. Isso confirma que o caminho otimizado para consultas — balanceador/CDN, API em Cloud Run, cache Redis e leituras do PostgreSQL — absorve picos de leitura sem degradação perceptível. As métricas do Cloud Monitoring confirmaram o escalonamento até as 10 instâncias de Cloud Run disponíveis, com CPU média de 72% e memória de 31%, mostrando que a cota foi integralmente utilizada sem comprometer a latência.

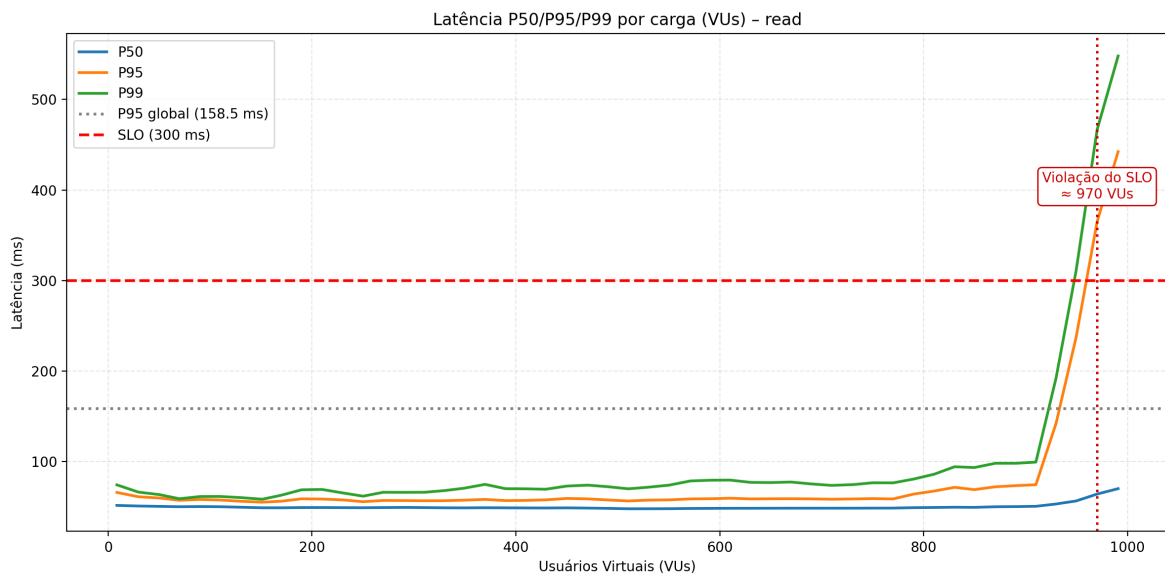


Figura 3. Latências P50/P95/P99 por carga (VUs) no cenário de leitura. A linha tracejada indica o SLO de 300 ms, não violado durante o teste.

5.2. Cenário Misto (Leitura/Escrita)

O cenário misto foi executado com provisionamento automático de usuários/token por VU, eliminando o gargalo artificial observado quando todos competiam pela mesma conta de teste. Com 650 usuários alternando entre 65% de leituras e 35% de escritas durante 21 minutos, o sistema sustentou 226 requisições por segundo (283 mil iterações) sem erros e preservou o SLO de 300 ms até aproximadamente 450 RPS (a violação é observada ao redor de 539 VUs). Assim que a carga ultrapassou 550 VUs, o Cloud Run atingiu novamente o teto de 10 instâncias e o P95 consolidado do ensaio ficou em 658 ms ($p99 = 2,67$ s). A pressão veio do aumento de CPU e latência nas operações de escrita, já que cada instância precisou executar validações, persistência e invalidação de cache antes de responder; não houve evidência de saturação no Cloud SQL, cujos tempos de consulta permaneceram estáveis durante todo o gradiente. Esse comportamento indica que, antes de pensar em particionamento de dados, é necessário permitir que mais instâncias HTTP sejam acionadas ou otimizar o custo das rotas de escrita para consumir menos CPU por requisição. Também ficou evidente o benefício do endpoint de provisionamento: cada VU pôde criar contas próprias e evitar competição por saldo, aproximando o teste de um uso real.

A Figura 5 evidencia o distanciamento entre os dois perfis: na leitura os percentis crescem suavemente até o pico, sinalizando que o gargalo ainda está restrito à cota de instâncias; já no misto, a curva dispara logo após 500 VUs porque cada POST exige mais trabalho por instância (validações, escrita e invalidação de cache), reduzindo o número de requisições atendidas por contêiner quando o limite de 10 instâncias é atingido. Esse comportamento é consistente com a literatura sobre sistemas transacionais centralizados [Kleppmann 2017], que destaca a dificuldade de um banco relacional único absorver altas taxas de escrita e sincronização sem ajustes adicionais. Assim, a arquitetura atual sustenta leituras em larga escala, mas depende de mais instâncias HTTP (ou otimizações focadas em escrita) antes de considerar estratégias mais complexas de particionamento ou réplicas.

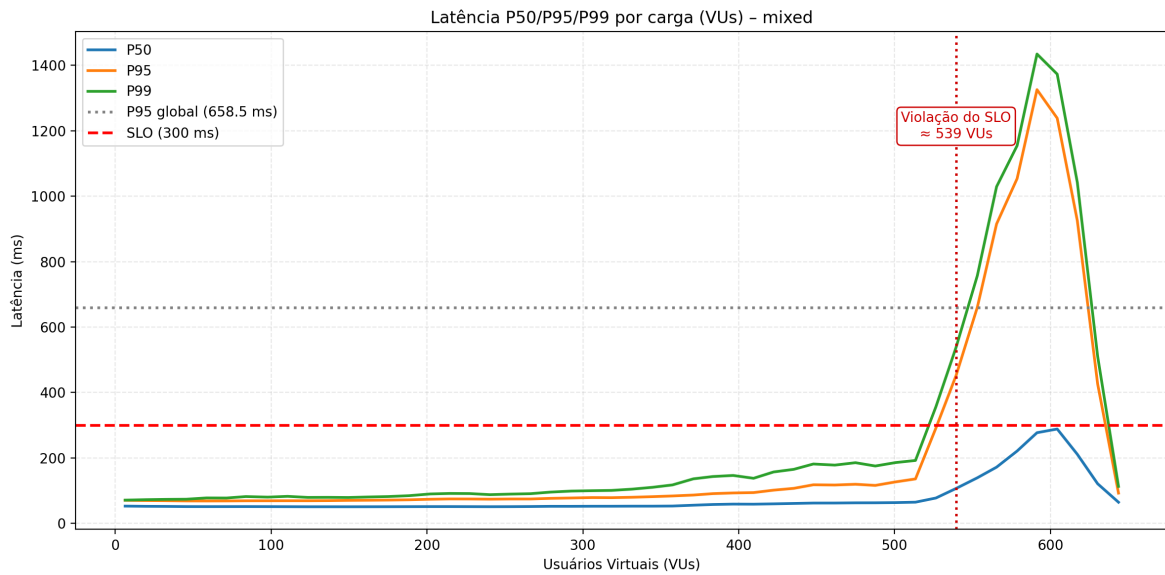


Figura 4. Latências P50/P95/P99 por carga (VUs) no cenário misto atualizado. A faixa tracejada marca o SLO de 300 ms; as etapas acima de 539 VUs apresentam violações persistentes desse SLO.

5.3. Processamento Assíncrono com Cloud Tasks

O ensaio de filas injetou 51,58 mil tarefas em um único lote e monitorou a drenagem entre 01:08 e 01:18, totalizando aproximadamente 10 min. Isso corresponde a um ritmo médio de 86 tarefas/segundo por worker HTTP em Cloud Run, comprovando que a combinação Cloud Tasks + workers escala horizontalmente sem perder entregas: cada aumento de backlog gerou novos *containers* apenas durante a janela de pico, sem repetição indevida ou perda de tarefas. A latência ponta-a-ponta permaneceu estável porque o Redis atuou somente como *backplane* para eventos, reduzindo contenção. Como Cloud Run cobra por instância ativa e Cloud Tasks por requisições enfileiradas, a elasticidade observada evita o superprovisionamento típico de VMs fixas: as 10 instâncias só permaneceram ativas enquanto havia tarefas, o que, mesmo sem detalhar faturamento, sugere controle de custo alinhado à carga real. resultado valida que o pipeline assíncrono consegue absorver rajadas significativas sem intervenção manual.

6. Conclusão

O ValorizeAI demonstrou que é possível documentar e validar, de ponta a ponta, uma arquitetura orientada a eventos construída exclusivamente com serviços gerenciados. O trabalho conecta decisões de design (Cloud Run, Redis, Cloud SQL, Cloud Tasks, Reverb) ao domínio transacional modelado no banco relacional, oferecendo um roteiro explícito para equipes que precisam de elasticidade sem perder consistência forte. A metodologia experimental, baseada em SRE, complementa essa documentação ao detalhar como planejar SLOs, materializar infraestrutura como código e acionar testes repetíveis, fortalecendo a reprodutibilidade típica de uma pesquisa aplicada.

Os cenários de carga comprovaram a eficácia das otimizações de leitura: com 1.000 usuários virtuais (≈ 470 req/s em média e pico de ≈ 970 req/s no platô final), o sistema manteve latência P95 em 158 ms e zero erros, satisfazendo os SLOs estabelecidos. O cenário misto, por sua vez, revelou o limiar atual da solução; o patamar de 650 VUs (≈ 226 req/s) levou o P95 a 658 ms e o p99 a 2,67 s porque o Cloud Run estacionou no limite de 10 instâncias e cada requisição de escrita

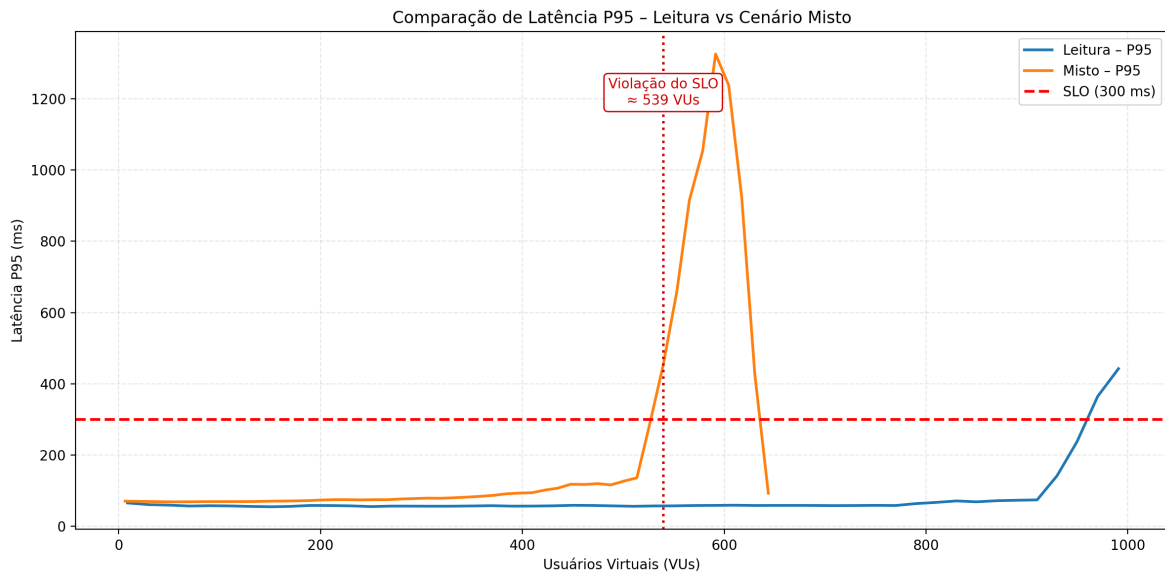


Figura 5. Comparação direta entre os P95 das cargas de leitura e mista. O cenário misto viola o SLO aos 539 VUs, enquanto a curva de leitura permanece abaixo de 200 ms mesmo com 1.000 VUs.

consome mais CPU devido aos passos adicionais (validações, persistência, invalidação de cache). O Cloud SQL manteve tempos estáveis, indicando que o gargalo imediato ainda é computacional na camada HTTP. No plano assíncrono, a drenagem de 51,58 mil tarefas em 10 min confirmou que Cloud Tasks e workers em Cloud Run sustentam bursts significativos sem intervenção humana, atendendo aos requisitos de elasticidade e rastreabilidade.

6.1. Limitações

Dois fatores limitam os resultados. Primeiro, o ambiente operou com as cotas padrão de uma conta recém-provisionada (10 instâncias de 1 vCPU), o que restringe a observação de comportamentos em estágios superiores de escala. Segundo, o banco de dados permaneceu concentrado em uma única instância Cloud SQL regional; embora consistente, essa configuração impõe contenções em workloads que concentram escritas nas mesmas contas ou categorias.

6.2. Trabalhos Futuros

As próximas etapas incluem ampliar gradualmente a cota de Cloud Run para liberar mais instâncias e repetir os testes sob 20 ou 40 vCPU, ao mesmo tempo em que se avaliam os próximos gargalos: réplicas de leitura do Cloud SQL para aliviar consultas analíticas, escalonamento vertical das instâncias primárias para absorver picos de escrita e mecanismos de *failover* multi-região para banco e Redis. Também se pretende complementar o pipeline assíncrono com um orquestrador baseado em eventos (ex.: Workflows) para medir a resiliência a falhas parciais, e estender o uso de embeddings e automações de categorização para outras áreas (como detecção de anomalias).

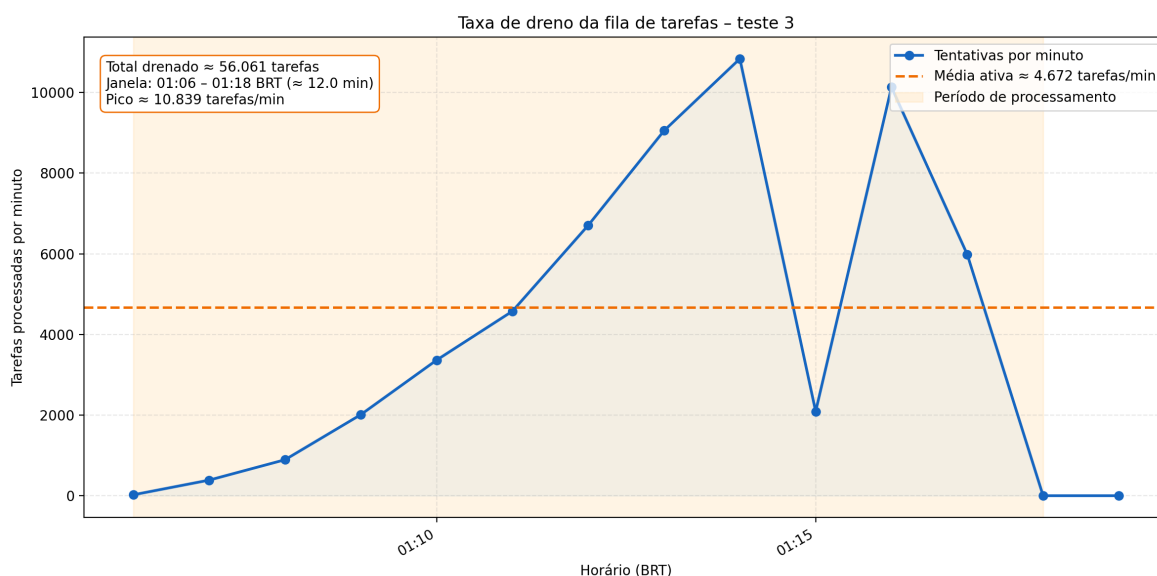


Figura 6. Taxa de processamento das filas (tarefas/min) e janela de drenagem entre 01:08 e 01:18 no ensaio com Cloud Tasks.

Referências

- Abad, C., Foster, I. T., Herbst, N., e Iosup, A. (2021). Serverless computing: One step forward, two steps back. *IEEE Computer*, 54(3):48–58. Refereciado por [1].
- Barri (2025). Mastering website scalability for large traffic: Preparing for surges. Technical report, Queue-Fair. <https://queue-fair.com/website-scalability-for-large-traffic>, Acessado em: 2025-10-14.
- Basteri, A. (2023). What makes observability a priority. Technical report, New Relic. <https://newrelic.com/resources/white-papers/observability-as-a-priority>, Acessado em: 2025-10-14.
- Beyer, B., Jones, C., Petoff, J., e Murphy, N. R. (2016). *Site Reliability Engineering: How Google Runs Production Systems*. O'Reilly Media, Inc., Sebastopol, CA, USA.
- Cervone, V. (2024). k6 - performance testing for developers. Technical report, Grafana Labs. <https://k6.io/>, Acessado em: 2025-10-14.
- Christidis, K. et al. (2022). Serverless cloud architectures for machine learning model deployment: A systematic review and case study. *World Journal of Advanced Engineering and Technology (WJAETS)*, 5(2). Disponível em: <https://wjaets.com/sites/default/files/WJAETS-2022-0025.pdf>.
- Confluent (2024). Event-driven architecture (eda): A complete introduction. Technical report, Confluent Inc. <https://www.confluent.io/learn/event-driven-architecture/>, Acessado em: 2025-10-14.
- Datadog (2024). Serverless vs. containers: What's the difference? Technical report, Datadog Research. <https://www.datadoghq.com/knowledge-center/serverless-architecture/serverless-vs-containers/>, Acessado em: 2025-10-14.

- Evans, E. (2003). *Domain-Driven Design: Tackling Complexity in the Heart of Software*. Addison-Wesley Professional, Boston, MA, USA.
- Fernando, L. e Engel, M. M. (2025). Comparative performance benchmarking of websocket libraries on node.js and golang. *Sinkron : Jurnal dan Penelitian Teknik Informatika*, 9(4).
- Fowler, M. (2002). *Patterns of Enterprise Application Architecture*. Addison-Wesley Professional, Boston, MA, USA.
- Fowler, M. (2011). Cqrs. Technical report, martinowler.com. <https://martinfowler.com/bliki/CQRS.html>, Acessado em: 2025-10-14.
- Google Cloud (2024). What is cloud elasticity? understanding elastic computing. Technical report, Google Cloud. <https://cloud.google.com/discover/what-is-cloud-elasticity>, Acessado em: 2025-10-14.
- Guerriero, A. et al. (2019). Adoption, support, and challenges of infrastructure-as-code: Insights from industry. In *Proceedings of the 2019 IEEE/ACM 41st International Conference on Software Engineering: Software Engineering in Practice (ICSE-SEIP)*.
- Hebbar, K. S. (2025). Priority-aware reactive apis: Leveraging spring webflux for sla-tiered traffic in financial services. *European Journal of Electrical Engineering and Computer Science*, 9(5).
- Kleppmann, M. (2017). *Designing Data-Intensive Applications*. O'Reilly Media, Inc.
- Laravel Holdings Inc. (2025). Laravel reverb official documentation. Technical report, Laravel Holdings Inc. <https://reverb.laravel.com>, Acessado em: 2025-10-14.
- Lloyd, W., Ramesh, S., Chinthalapati, S., Ly, L., e Pallickara, S. (2018). Serverless computing: An investigation of factors influencing microservice performance. In *2018 IEEE International Conference on Cloud Engineering (IC2E)*.
- Martin, R. C. (2017). *Clean Architecture: A Craftsman's Guide to Software Structure and Design*. Prentice Hall, Hoboken, NJ, USA.
- McCoy, J. e Forsgren, N. (2020). *SLO Adoption and Usage in Site Reliability Engineering*. O'Reilly Media, Inc., Sebastopol, CA, USA.
- Pessa, A. (2023). Comparative study of infrastructure as code tools for amazon web services. Master's thesis, Tampere University. Disponível em: <https://trepo.tuni.fi/bitstream/handle/10024/149567/PessaAntti.pdf>.
- Sonawane, S. et al. (2024). The role of serverless architecture in scalable and efficient web development. *International Journal of Scientific Research in Science and Technology*. Disponível em: https://www.researchgate.net/publication/389615854_The_Role_of_Serverless_Architecture_in_Scalable_and_Efficient_Web_Development.
- Thepphakan, A. (2025). Study of real-time data communication using pulsar and rabbitmq (case study of stock price in lao securities exchange). *International Journal of Advanced Research in Computer Science and Software Engineering*.
- Twine (2022). Twine case study: A scalable and open-source raas. Technical report, Twine Realtime Initiative. <https://twine-realtime.github.io/case-study>, Acessado em: 2025-10-14.

Yadav, P. S. (2019). Designing a high-performance real-time leaderboard system using redis: Scalability, efficiency, and fault tolerance. *Journal of Scientific and Engineering Research*, 6(3):313–320.