

4.8 performance & scalability plans

برنامه‌های عملکرد و مقیاس‌پذیری - ProDecks

نسخه: 1.0

تاریخ: 1403/11/16

برنامه‌ریزی شده - MVP: وضعیت

فصل ۱: مرور کلی

هدف این سند: ۱.۱.

در مواجهه با رشد کاربران ProDecks تعیین استراتژی‌ها و برنامه‌های لازم برای تضمین عملکرد بهینه و مقیاس‌پذیری سیستم و داده‌ها.

چالش‌های پیش‌رو: ۱.۲.

- رشد تصاعدی تعداد کاربران و داده‌ها
- افزایش انتظارات کاربران از زمان پاسخ‌گویی
- (High Availability) نیاز به دسترسی‌پذیری بالا
- حفظ هزینه‌های زیرساخت در سطح منطقی

اهداف کلیدی: ۱.۳.

- برای ۹۵٪ درخواست‌ها زمان پاسخ API: ۲۰۰ms
- زمان لود صفحات: زیر ۳ ثانیه
- (Uptime) دسترسی‌پذیری: ۹۹.۹٪
- پشتیبانی از ۱۰,۰۰۰ کاربر همزمان
- مقیاس‌پذیری خطی هزینه با رشد کاربران

فصل ۲: معیارهای عملکرد (Performance Metrics)

۲.۱. (Application Metrics) معیارهای کاربردی

- API Response Time:
 - P50: < ۱۰۰ms
 - P95: < ۲۰۰ms
 - P99: < ۵۰۰ms

- Page Load Time:
 - First Contentful Paint: < ۱.۵s

۳.۲. بهینه‌سازی بکاند

۱. Caching:

- Redis برای داده‌های پر تکرار
- Memcached برای session storage
- CDN قابل cache برای API responses

۲. Query Optimization:

- ایندکس‌گذاری بهینه
- کوئری‌های بهینه شده
- استفاده از pagination

۳. Connection Pooling:

- مدیریت کارآمد اتصالات دیتابیس
- connection leaks جلوگیری از
- بر اساس بار pool تنظیم اندازه

۳.۳. بهینه‌سازی پایگاه داده

۱. Indexing Strategy:

- ایندکس‌های ترکیبی برای کوئری‌های پر تکرار
- حذف ایندکس‌های غیر ضروری
- استفاده از covering indexes

۲. Partitioning:

- بر اساس تاریخ ایجاد users جدول
- بر اساس space_id جدول
- horizontal partitioning پس از رسیدن به حد آستانه

۳. Query Optimization:

- برای تحلیل کوئری‌ها EXPLAIN استفاده از
- جلوگیری از N+1 queries
- برای گزارش‌ها materialized views استفاده از

فصل ۴: استراتژی‌های مقیاس‌پذیری

۴.۱. (Vertical Scaling) مقیاس‌پذیری عمودی

- فاز ۱ (MVP): سرورهای متوسط (4 vCPU, 8GB RAM)
- فاز ۲ (رشد اولیه): سرورهای بزرگ (8 vCPU, 16GB RAM)
- فاز ۳ (مقیاس بزرگ): سرورهای اختصاصی (16 vCPU, 32GB RAM)

۱. Web Tier:

- Load Balancer (NGINX/HAProxy)
- گروه سرورهای وب
- Stateless طراحی برای ease of scaling

۲. Database Tier:

- Primary-Secondary Replication
- Read Replicas برای بار خواندنی
- Sharding در صورت نیاز

۳. Cache Tier:

- Redis Cluster
- Memcached با چندین نود

۴.۳. الگوهای معماری مقياس پذیر

۱. Microservices (آینده):

- تفکیک سرویس‌های اصلی
- API Gateway
- Service Discovery

۲. Event-Driven Architecture:

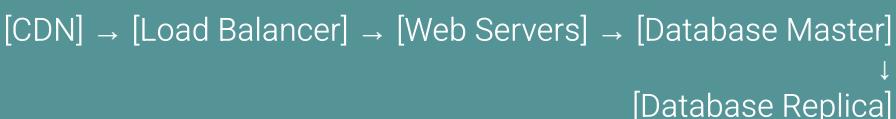
- Message Queue (RabbitMQ/Kafka)
- Event Sourcing برای موجودیت‌ها
- CQRS برای جداسازی read/write

۴.۴. Database Scaling Plan

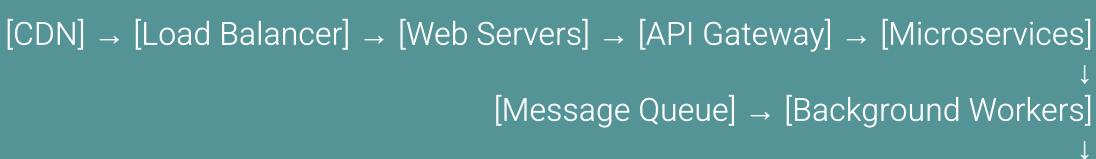
- تا: ۱۰,... کاربر: Single Master + ۲ Read Replicas
- تا: ۵۰,... کاربر: Master-Master + ۴ Read Replicas
- بیش از ۵۰,... کاربر: Sharding بر اساس space_id

فصل ۵: معماری مقياس پذیر

۵.۱. (MVP) معماری فعلی



۵.۲. معماری فاز ۲ (۱۰,... کاربر)



۷.۱. Load Testing:

- شبیه‌سازی کاربران واقعی
 - افزایش تدریجی بار
 - شناسایی bottlenecks

۷.۲. Stress Testing:

- بار بیش از حد ظرفیت
- شناسایی نقاط شکست
- تعیین حد نهایی سیستم

۷.۳. Endurance Testing:

- تست طولانی مدت (۷۲-۲۴ ساعت)
- memory leaks شناسایی
- stability بررسی

۷.۴. Spike Testing:

- افزایش ناگهانی بار
- شبیه‌سازی رویدادهای خاص
- resilience تست

۷.۵. ابزارهای تست

- Load Testing: Apache JMeter, k6
- Monitoring: Prometheus, Grafana
 - APM: New Relic, Datadog
- Database: pt-query-digest, Percona Toolkit

۷.۶. سناریوهای تست

- سناریو ۱: ایجاد کارت
- کاربران همزمان: ۱۰۰ → ۵۰۰ → ۱,۰۰۰
 - فرکانس: ۱۰ کارت/دقیقه به ازای هر کاربر
 - معیارها: زمان پاسخ, throughput, error rate

۷.۷. Drag & Drop: سناریو ۲

- کاربران همزمان: ۵۰ → ۲۰۰ → ۵۰۰
- عملیات: جابجایی ۵ کارت/دقیقه به ازای هر کاربر
- معیارها: real-time performance, consistency

۷.۸. گزارش‌گیری: سناریو ۳

- کاربران همزمان: ۱۰ → ۵۰ → ۱۰۰
- زیاد joins گزارش‌ها: پیچیده با
- معیارها: زمان اجرای کوئری, مصرف منابع

۸.۱. معیارهای مانیتورینگ

- Application Level:
 - Response times
 - Error rates
 - Throughput
 - Business transactions
- Infrastructure Level:
 - CPU/Memory/Disk usage
 - Network I/O
 - Service availability
- Database Level:
 - Query performance
 - Connection counts
 - Replication lag
 - Lock contention

۸.۲. تنظیمات هشدار

- Critical (PagerDuty/SMS):
 - Uptime < ۹۹%
 - Error rate > ۵%
 - Database down
- Warning (Email/Slack):
 - Response time P95 > ۵۰..ms
 - CPU usage > ۸۰% به مدت ۵ دقیقه
 - Disk usage > ۸۵%
- Informational (Dashboard):
 - Performance degradation
 - Unusual patterns
 - Capacity thresholds

۸.۳. dashboards

- وضعيت لحظه‌ای Real-time Dashboard:
- معیارهای کسب‌وکار Business Dashboard:
- روند عملکرد Performance Dashboard:
- مصرف منابع و Capacity Dashboard: trends

۹.۱. قوانین Auto-scaling

- origin برای کاهش بار استفاده از
- caching aggressive برای کاهش بار دیتابیس
- object storage ارزان برای داده‌های قدیمی استفاده از
- compression داده‌ها در انتقال و ذخیره‌سازی

۱۰.۳. مدیریت چرخه حیات داده

- ارزان‌تر storage انتقال داده‌های قدیمی به
- archiving پروژه‌های کامل شده
- حذف خودکار داده‌های موقت
- retention policies بر اساس ارزش داده

۱۱. Disaster Recovery و Business Continuity

۱۱.۱. Backup استراتژی

- هر ساعت Backup + روزانه incremental: داده‌ها
- خودکار با Backup versioning: پیکربندی
- تست بازیابی: ماهانه

۱۱.۲. Recovery استراتژی

- ساعت ۴ Recovery Time Objective (RTO):
- ساعت ۱ Recovery Point Objective (RPO):
- اصلی region در صورت قطعی failover: قابل انتقال

۱۱.۳. High Availability طراحی

- برای سرویس‌های حیاتی Multi-AZ deployment
- Load balancing across zones
- Automated failover mechanisms

۱۲. نقشه راه مقیاس‌پذیری

۱۲.۱. FaaS (ماه‌های ۱-۳): پایه‌های مقیاس‌پذیری

- لایه caching پیاده‌سازی
- بهینه‌سازی کوئری‌های اصلی
- پایه monitoring تنظیم
- اوپریت load تست

۱۲.۲. FaaS (ماه‌های ۴-۶): مقیاس‌پذیری اوپریت

- read replicas پیاده‌سازی
- web tier برای auto-scaling
- statelessness بھبود معماری برای
- performance تست‌های پیشرفته

فاز ۳ (ماه‌های ۷-۹): مقیاس‌پذیری پیشرفته ۱۲.۳.

- microservices تفکیک به
- message queue پیاده‌سازی
- database sharding طراحی
- global CDN راهاندازی

فاز ۴ (ماه‌های ۱۰-۱۲): مقیاس‌پذیری enterprise ۱۲.۴.

- multi-region deployment
- advanced caching strategies
- predictive auto-scaling
- cost optimization پیشرفته

فصل ۱۳: تیم و مسئولیت‌ها

نقش‌های کلیدی ۱۳.۱.

- Performance Engineer: مسئول بهینه‌سازی
- DevOps Engineer: scaling مسئول زیرساخت و
- DBA: مسئول مقیاس‌پذیری دیتابیس
- SRE: reliability و monitoring مسئول

فرآیندها ۱۳.۲.

- Weekly Performance Reviews
- Capacity Planning جلسات ماهانه
- Load Testing قبل از هر release اصلی
- Performance Budgeting برای feature های جدید

آموزش و توسعه ۱۳.۳.

- best practices آموزش مقیاس‌پذیری
- knowledge sharing sessions
- communities مشارکت در
- new technologies تحقیق در مورد تکنولوژی‌های جدید

فصل ۱۴: ریسک‌ها و mitigations

ریسک‌های فنی ۱۴.۱.

- ریسک: Database becoming bottleneck
- Mitigation: Read replicas, caching, query optimization

- ریسک: Single point of failure
- Mitigation: Multi-AZ, failover automation

- ریسک: Cost overrun
- Mitigation: Cost monitoring, right-sizing, reserved instances

ریسک‌های عملیاتی ۱۴.۲.

- ریسک: Lack of monitoring
- Mitigation: Comprehensive monitoring, alerting

- ریسک: Insufficient testing
- Mitigation: Regular load testing, chaos engineering

- ریسک: Knowledge silos
- Mitigation: Documentation, cross-training

ریسک‌های کسب‌وکار ۱۴.۳.

- ریسک: Performance affecting user retention
- Mitigation: Performance SLAs, proactive monitoring

- ریسک: Inability to handle growth spikes
- Mitigation: Auto-scaling, capacity buffer

- ریسک: Competitive disadvantage
- Mitigation: Continuous performance improvement

فصل ۱۵: نتیجه‌گیری

یک رویکرد جامع و مرحله‌ای برای تضمین ProDecks برنامه عملکرد و مقیاس‌پذیری تجربه کاربری بهینه در تمام مراحل رشد کسب‌وکار ارائه می‌دهد. با ترکیب، بهینه‌سازی‌های فنی، استراتژی‌های مقیاس‌پذیری هوشمند و فرآیندهای عملیاتی دقیق ما آمادگی لازم برای رشد تصاعدی را خواهیم داشت.

- تمرکز ما نه تنها روی مقیاس‌پذیری فنی، بلکه روی مقیاس‌پذیری اقتصادی نیز هست. حفظ کیفیت سرویس در حالی که هزینه‌ها به صورت خطی با رشد افزایش می‌یابند.

ضمیمه‌ها

