

** ShahnamehMap** — سند ۴,۸: برنامه عملکرد و مقیاس‌پذیری

نسخه: ** ۱,۰ ***

تاریخ: ** ۱۴۰۳ *** ۱۰/۹/۲۰۲۳

** (CTO) مدیر فنی تهیه‌کننده: ***

وضعیت: ** فعال ***

** ۱. (Service Level Objectives -
SLOs)** اهداف سطح سرویس -

و ** (Reliability) های ما تعهد داخلی برای قابلیت اطمینان SLO
ها مبنای SLO سیستم است. این ** (Performance) کارایی تصمیم‌گیری‌های مهندسی و تجاری هستند.

| ** هدف | ** (SLI) سرویس/کندکتور *** | ** شاخص سطح سرویس | ** توضیح | ** (SLO) سطح سرویس | :--- | :--- | :--- | :--- | :--- |

| **API Gateway** درخواست‌های موفق | REST همه سرویس‌های (HTTP 2xx/3xx) نسبت به کل درخواست‌ها. | ۹۹,۵٪ در هر دوره ۲۸ ساعت (~۳,۵ روزه). | (در ماه Downtime خطای مجاز).

اتصالات | **Game Engine** سرویس Real-time) اتصالات | WebSocket پایدار (بدون قطعی غیرعمدی) نسبت به کل اتصالات. | ۹۹٪ در هر دوره ۲۸ روزه. | تحمل ۱٪ قطعی برای (سختگیرانه).

برای (P95) های اصلی | زمان پاسخ API (Latency) تأخیر *ms* ۲۰۰ برای endpoint های حیاتی • `POST /api/characters` • `GET /api/campaigns` | P95. | ۹۵٪ درخواست‌ها در این زمان پاسخ می‌گیرند.

زمان بین | **Game Action Latency** (Game Action Latency) تأخیر عملیات بازی | توسط همه State توسط کاربر و دریافت به روزرسانی Action ارسال یک برای تجربه بازی | ۵۰۰ ms زیر برای P99. | (P99) بازیکنان اتاق روان حیاتی است

درصد زمان قابل | **Overall Uptime** در دسترس بودن کلی | دسترسی بودن سرویس از دید کاربر نهایی. | ۹۹٪ در هر دوره ۳ ماهه. | معيار سطح بالای سلامت کسب و کار

*، بودجه خطای ما در **SLO 99.5% با (Error Budget):** بودجه خطای ماه ۰,۵٪ (۴۳,۲ دقیقه) است. اگر در هفته اول ماه ۰,۴٪ خطای مصرف کنیم، *تیم موظف است توسعه فیچرهای جدید را متوقف کند* و تمام تمرکز را بر بهبود پایداری بگذارد.

۲. Capacity Planning) برنامه اندازه‌گیری و ظرفیت‌سنجی

**۱۰. Key Performance Metrics) معیارهای کلیدی عملکرد

:مانیتور می‌کنیم **Real-time** ما این متریک‌ها را به صورت

| متریک | هدف | ابزار مانیتورینگ | هشدار | Alert) |

| :--- | :--- | :--- | :--- |

| * API Gateway | - | (RPS)** نرخ درخواست |

| **Prometheus** + **Grafana** | < ۵۰٪ | افزایش ناگهانی |

| **CPU Utilization** (Game Engine) | > میانگین
برای ۸۵٪ | Node Exporter + K8s Metrics | > ۸۰٪ |
| بیش از ۵ دقیقه

| **Memory Utilization** (Game Engine) | < ۷۰٪ |
Node Exporter + K8s Metrics | > ۸۵٪. |

| **WebSocket** فعال اتصالات | - | **Redis INFO**
نزدیک شدن به حد نظری هر نمونه | command** + Prometheus | (~۷۰۰). |

| **PostgreSQL Connections** | < ۸۰٪ از
'max_connections' | **pg_stat_activity** + Prometheus
| > ۹۰٪. |

| **Redis Latency** (P99) | < ۵ms | **Redis**
Monitoring** | > ۱۰ms. |

و نقاط گلوگاه شناخته **Capacity Model** (مدل ظرفیت .۲،۲) ***
***شده

و ***کاربر فعال ماهانه** (CCU) ما بر اساس ***کاربر فعال همزمان
برنامه‌ریزی می‌کنیم (MAU)**.

** فرضیات مدل

* اتصال ۱ CCU هر کاربر فعال همزمان + در دقیقه REST درخواست ~۵.

* CCU یعنی ۱۰،۰۰۰ (پیک) CCU ≈ نسبت MAU ≈ ۲۰۰،۰۰۰ MAU).

| گلوگاه شناخته شده *** | * طرفیت فعلی (Resource) | * منبع | راهکار Scale-out | * حد آستانه برای CCU برای ۱۰۰۰ | * مقیاس پذیری

| :--- | :--- | :--- | :--- | :--- |

| * Game Engine (Node.js) | ** CPU
هر اتاق. | ۱ نمونه = State و حافظه برای Node.js | * تک هسته
~۷۰۰ اتصال پایدار. | ۶۰۰ اتصال در هر نمونه. | * مقیاس افقی:
در RoomID API اضافه کردن نمونه های جدید. * روتینگ مبتنی بر
Gateway. |

| * Redis (State + Pub/Sub) | * پهنه ای باند شبکه و | * تأخیر در صورت فشار بالا. | نسخه فعلی: ۱ نمونه
> | ** CPU > | * . تأخیر در صورت فشار بالا. | نسخه فعلی: ۱ نمونه
۷۰٪ Latency > ۵ms*. | ۱۰٪ Redis ارتقا به کلاستر

(Redis Cluster).
۲. به RoomID شارдинگ اتاق‌ها** بر اساس** چندین نمونه Redis. |

| **PostgreSQL** داده کاربر/کاراکتر) | **تعداد** Connection
| . همزمان** و **پرفورمنس کوئری‌های گزارش‌گیری.** | ۱ نمونه
Connection Pool > ۸۰٪ ۲۰۰ < ms**. |
۱. **Connection Pooling** با PgBouncer.
| . برای کوئری‌های گزارش (Read Replica)** خواندنی/نوشتني**
| با) شارдинگ** طولی‌مدت (Citus). |
| **ترافیک** (Network Bandwidth)** | **Real-
time** (WebSocket). | ~۱۰۰ استفاده < مگابیت بر ثانیه. |
| . ارتقای پلن شبکه با ارائه‌دهنده کلاد ۷۰٪. |
| . برای API Gateway (Kong)** | * مقدار حافظه (RAM)**
| . کش و تعداد اتصالات همزمان. | ۲ نمونه پشت Load Balancer. |
| . **CPU > ۷۵٪**. | * نمونه‌های Kong. |

۳. (Load Testing Plan)** برنامه تست بار.

** هدف: ** شناسایی محدودیت‌های سیستم قبل از مواجهه واقعی.

۳.۱. ابزارها و محیط

* JS به دلیل سادگی، اسکریپتنویسی با) k6 ** ابزار * و Grafana). یکپارچگی با

* Production (اما است که دقیقاً مشابه Staging ** محیط با نمونه‌های کوچک‌تر)

* Production (anonimized) ** داده تست: از دیتابیس برای تست واقعی‌تر استفاده می‌شود

۳.۲. سناریوهای تست بار

حفظ پیک **، (Ramp-up) ** هر تست ۳ مرحله دارد: ** رمپ آپ (Sustain) **، ** رمپ داون (Ramp-down) **.

| * هدف *** | * پروفایل بار *** | * معیارهای موفقیت *** |

| :--- | :--- | :--- | :--- |

| **LT-۱: API های CRUD ثبت نام و ساخت کاراکتر انبوه** | تست: در ۵ دقیقه. حفظ ۵ دقیقه. • نرخ خطای VU → ۵۰۰ VU → دیتابیس. | HTTP < ۱%.
• P95 Latency < ۳۰۰ ms.
• PostgreSQL CPU < ۸۰%. |

| **LT-۲: Real-time** شبیه سازی کاربران بازی | پیچیده: ایجاد ۲۰۰ "اتاق بازی" و Game Engine و Redis. | **LT-۲: Real-time** شبیه سازی کاربران بازی | پیچیده: ایجاد ۲۰۰ "اتاق بازی" و Game Engine و Redis. هر بازیکن هر ۱۰ ثانیه CCU مجموعاً ۸۰۰ (مجازی، هر کدام ۴ "بازیکن" هر بازیکن هر ۱۰ ثانیه). • نرخ قطعی Action می فرستد. | نرخ قطعی WebSocket < ۰,۵%.
• Game Action Latency (P99) < ۷۰۰ ms.
• Redis Latency < ۱۰ ms.
• Game Engine CPU < ۷۰%. |

| **LT-۳: API** پیک ترافیک (برگشت از کمپین مارکتینگ) | تست: در ۲ VU و مقیاس خودکار. | افزایش سریع از ۱۰۰ به ۱۵۰۰ API Gateway | . دقیقه. حفظ ۳ دقیقه | • API Gateway Latency (P95) < ۲۵۰ ms.
• Auto-scaling ۳ دقیقه فعال شود | سرویس ها ظرف ۳ دقیقه استقامت: |

| **LT-۴: Endurance Test** | استقامت: memory leak Degradation یا در طول زمان. | بار ثابت معادل ۵۰٪ از ظرفیت هیچ گونه افزایش تدریجی پیش بینی شده پیک برای ۲ ساعت. | • هیچ گونه افزایش تدریجی Memory Utilization یا Latency ۲۰٪ بیش از |

فرکانس اجرا .۳.

- * ***پس از هر تغییر عمدی معماري:*** (اجباری)
- * ***ماهانه:*** (برای اطمینان از حفظ عملکرد)
- * ***قبل از رویدادهای بزرگ بازاریابی:*** (اجباری)

به تفکیک لایه (Scalability Plan) برنامه مقیاس پذیری .۴.

- ### ** مقیاس پذیری برنامه .۱. (Application Scaling)**
- * *** استاتلس (Stateless - APIها، Gateway):*** **Auto-scaling افقی*** بر اساس CPU/Memory Utilization.
- * *** به مدت ۱ دقیقه `scale up` اگر CPU > ۷۰٪ K8s HPA:**
- * *** به مدت ۵ دقیقه `scale down` اگر CPU < ۳۰٪.
- * *** استیت فول (Stateful - Game Engine):*** **Auto-scaling افقی با هوشمندی
- * *** اتاق های خاصی را در حافظه دارد State چالش:*** هر نمونه،

* * راهکار: استفاده از **شاخص** زمانی که میانگین اتصالات هر نمونه از `connections_per_pod` است. ۶۰۰ گذر کند، نمونه جدید ایجاد می‌شود. **اتاق‌های جدید** به نمونه‌های State شود، اتاق‌هایش بر اساس Fail کمبارتر هدایت می‌شوند. اگر نمونه‌ای .. به نمونه‌های دیگر مهاجرت می‌کنند Redis ذخیره شده در

۴.۲. مقیاس‌پذیری داده (Data Scaling)

* **PostgreSQL:** مسیر تکامل:

۱. **Connection Pooling** با PgBouncer.
۲. **Read Replica** برای افزودن MAU > ۵۰,۰۰۰ در. کوئری‌های گزارش و تحلیل.
۳. **Sharding** (Sharding) برای شارдинگ MAU > ۵۰۰,۰۰۰ در. با استفاده از `user_id` اساس Citus**.

* **Redis:** مسیر تکامل:

۱. **یک نمونه قوی** (هم‌اکنون).
۲. **Redis Cluster** برای راهاندازی CCU > ۵,۰۰۰ در. توزیع داده و بار.

* **ClickHouse:** به طور ذاتی برای بارهای تحلیلی سنگین مقیاس‌پذیر افقی است.

* **۴.۳. CDN** مقیاس‌پذیری شبکه و

- * **CDN:** استفاده از asset های برای ارائه Cloudflare است. با رشد ترافیک، پلن‌های بالاتر خریداری (CSS، SLS تصاویر،) استاتیک می‌شود.
 - * پنهانی باند: ارتقای پلن شبکه با ارائه دهنده کلاد*

* **۵. (Cost Projection)** برآورد هزینه رشد.

ما هزینه‌ها را بر اساس مدل مصرف منابع در پلتفرم کلاد داخلی (مثالاً دیجی‌کالا کلاد) برآورد می‌کنیم.

| **سناپریو کاربری** | **منابع مورد نیاز** | **هزینه ماهانه** | **تخمینی** | **توضیح**

| :--- | :--- | :--- | :--- |

| * **نمونه ۴ Kپیک): ۱ MAU: ۲۰K, CCU (سال اول) *
Game Engine (۲۲ هسته، ۲ GB)
• نمونه ۲ API/Safer (۲ گیگابایت RAM)
• Redis (۲GB)
• PostgreSQL (۴GB RAM)
• Redis (۲GB)
• ۲۰۰ میلیون تومان ~۸-۱۲ GB پهنای باند: | * **MVP در حال رشد |

| * **نمونه ۱۰ Kپیک): ۱۰ MAU: ۲۰۰K, CCU (سال دوم) *
Game Engine
• نمونه ۶ API/Safer
• PostgreSQL + Read Replica (۸GB)
• Redis Cluster (۳ Node)
• ۲ میلیون تومان ~۴۰-۶۰ TB پهنای باند: | * **مرحله ۲ رشد سریع |

| * **نمونه ۲۰ Kپیک): ۲۰ MAU: ۱M, CCU (سال سوم بین‌المللی) *
Game Engine
• نمونه ۷۵ Kپیک): ۵۰ API/Safer
• PostgreSQL Sharded (Citus)
• Redis Cluster
• ۱۰ بزرگ TB+ | * **۳۵۰-۲۰۰ میلیون تومان * هزینه کلاد بین‌المللی برای دیاسپورا. | پلتفرم پایدار و سودآور |

* نکته: این هزینه‌ها فقط زیرساخت است. هزینه‌های نیروی انسانی، مارکتینگ و سایر عملیات جداگانه است.

** جمع‌بندی: آیا سیستم از پا درمی‌آید؟ ۶.

بر اساس این برنامه، پاسخ **خیر** است — اگر و تنها اگر مانیتورینگ فعال** ما هشدارها را به موقع شناسایی کند**. ۱. تست‌های بار منظم** گلوگاه‌های جدید را قبل از بروز مشکل آشکار کنند.

۲. **Auto-scaling Rules** به درستی پیکربندی شده باشند

۴. 转折 (Inflection Point) وجود داشته باشد. برای ارتقای منابع در نقطه** ظرفیت مالی**.

** بزرگ‌ترین ریسک‌های مقیاس‌پذیری

* Real-time در مقیاس دهها هزار کاربر همزمان** است. راهکار شارдинگ Redis بازی Stateful:** گلوگاه** مدیریت ما اصلی مشکل است. Game Engine و معماری بر این چالش متمرکز است.

* می‌تواند هزینه زیرساخت را به طور خطی CCU هزینه:*** رشد نمایی ** باشد که گونه‌ای باشد که (Freemium) افزایش دهد. *** مدل درآمدی ما LTV کاربر پولی، هزینه زیرساخت کاربران رایگان را پوشش دهد **.

این سند نقشه راهی است که تضمین می‌کند *** مهندسی همگام با رشد کسب و کار حرکت می‌کند**. ما نه تنها می‌دانیم سیستم امروز چگونه کار می‌کند، بلکه می‌دانیم فردا با دوباره شدن کاربران، چه اقداماتی باید انجام دهیم و این اقدامات چقدر هزینه خواهد داشت. این *** قابلیت پیش‌بینی ***، کلید رشد کنترل شده و بدون شوک است.