

# مقیاس‌پذیری سرویس کوتاه‌کننده لینک

این سند به بررسی راهکارهای مقیاس‌پذیر کردن سرویس کوتاه‌کننده URL که با API و PostgreSQL FastAPI پیاده‌سازی شده است، می‌پردازد. هدف، حفظ عملکرد مناسب سرویس در شرایط بار بالا، جلوگیری از ایجاد نقاط تک خرابی (Single Point of Failure) و امکان گسترش افقی سیستم روی چندین سرور است.

## جداولی لاغ‌گیری سنگین از مسیر اصلی درخواست‌ها

هنگامی که حجم لاغ‌گیری برای هر درخواست افزایش یابد و نیاز باشد اطلاعات (مانند IP کاربر، زمان دسترسی و کد کوتاه) به صورت لحظه‌ای به یک سرویس جانبی یا دیتابیس جداگانه ارسال شود، اجرای این عملیات به صورت Synchronous می‌تواند زمان پاسخ‌دهی درخواست‌های اصلی (به‌ویژه ریدایرکت) را به شدت افزایش دهد.

برای جلوگیری از این مشکل، پیشنهاد می‌شود لاغ‌گیری را از مسیر بحرانی درخواست جدا کنیم. راهکار پیشنهادی استفاده از یک **Message Queue** مانند RabbitMQ یا Kafka است. در اطلاعات لاغ به صورت غیرهمزمان (Asynchronous) به یک تاپیک در کافکا ارسال می‌شود. این کار باعث می‌شود پاسخ به کاربر (HTTP 302 Status Code) بدون انتظار برای ذخیره لاغ برگردانده شود.

در سمت دیگر، Worker‌های Celery (با Redis به عنوان Broker) پیام‌ها را از صف دریافت کرده و پردازش می‌کنند. این Worker‌ها می‌توانند داده‌ها را تجمعی کرده و در یک دیتابیس تحلیلی جداگانه (مثلًا PostgreSQL شاردشده یا ClickHouse) ذخیره کنند. برای کاهش بار بیشتر، می‌توان لاغ‌ها را به صورت دسته‌ای (Batch) ارسال کرد؛ مثلًا هر ۰۰۰ رکورد یا هر ۵۰۰ میلی‌ثانیه یک بار.

در صورت بروز مشکل در صف (Back-Pressure)، با استفاده از Circuit Breaker می‌توان لاغ‌های غیرحیاتی را به طور موقت دور اندخت تا از تأثیر منفی بر سرویس اصلی جلوگیری شود. این رویکرد زمان پاسخ‌دهی عملیات اصلی را زیر ۱۰۰ میلی‌ثانیه نگه می‌دارد و با اصول Observability (مانند Tracing با OpenTelemetry) نیز سازگار است.

## آماده‌سازی سیستم برای استقرار چند-نمونه (Instance)

اگر قرار باشد سرویس روی چندین سرور به صورت همزمان اجرا شود، باید از حالت تک‌نمونه (Single-Instance) به معماری توزیع‌شده تغییر کنیم.

ابتدا یک **load balancer** (مانند NGINX) در مقابل نمونه‌های FastAPI قرار می‌گیرد تا ترافیک را بین آن‌ها تقسیم کند. در مرحله بعد، تمام وابستگی‌های دارای حالت (Stateful) باید **External** شوند:

**دیتابیس:** به جای PostgreSQL محلی، از یک کلاستر مدیریت شده با قابلیت Read Replication و Replication استفاده می شود. برای مقیاس پذیری بهتر نوشت، می توان داده ها را بر اساس Hash کوتاه شارد کرد.

**کش:** Redis به صورت یک سرویس خارجی (ترجیحاً Redis Cluster) برای کش کردن ریدایبرکت های پر تکرار و شمارش بازدیدها به کار گرفته می شود تا از مشکلات هم زمانی جلوگیری شود.

**صفها و لاغری:** Kafka یا RabbitMQ نیز به صورت کلاستر مستقر می شوند و Worker های Celery به طور مستقل مقیاس پذیر خواهند بود.

برای تولید کد کوتاه بدون تداخل (Collision) در محیط توزیع شده، می توان از یک **Key Generation** استفاده کرد یا کلیدها را به صورت پیش تولید در Redis ذخیره نمود.

### رسکهای مهم شامل موارد زیر است:

احتمال ایجاد کدهای تکراری → با استفاده از عملیات اتمیک در Redis (مانند INCR) یا قفل Zookeeper مدیریت می شود.

خرابی یک نمونه از دیتابیس یا Redis → با Auto-Failover و Replication برطرف می گردد.

افزایش تأخیر شبکه → با مانیتورینگ Tracing و در صورت نیاز استفاده از gRPC قابل کاهش است.

## مقابله با ترافیک سنگین ناشی از کمپین تبلیغاتی

در شرایطی که یک کمپین تبلیغاتی باعث ورود چندین هزار درخواست در ثانیه شود، باید ترکیبی از تصمیم های موجود و اقدامات جدید را به کار بگیریم تا سرویس از دسترس خارج نشود.

در طراحی فعلی از Connection Pooling در SQLAlchemy و Index روی short\_code و افزایش اتمیک تعداد بازدیدها استفاده شده است که پایه خوبی برای عملکرد فراهم می کند. برای تحمل بار بسیار بالا:

**لایه کش قوی:** با Redis برای کدهای پربازدید (Hot Keys) اضافه می شود. طبق اصل Pareto، معمولاً ۲۰٪ کدها ۸۰٪ ترافیک را تشکیل می دهند؛ کش کردن این موارد می تواند بار دیتابیس را به شدت کاهش دهد.

**Rate Limiting:** با استفاده از slowapi در Middleware اعمال می شود تا از سوءاستفاده جلوگیری کند و در صورت لزوم IP های کمپین را Whitelist کنیم.

**CDN:** (مانند Cloudflare) برای کش کردن پاسخ های 302 در لبه شبکه به کار گرفته می شود تا بخش زیادی از ترافیک از سرورهای اصلی عبور نکند.

**مانیتورینگ پیشرفته:** با Jaeger، Prometheus و Grafana با Tracing برای شناسایی سریع گلوگاهها و تنظیم آلام (مثلاً وقتی اتصالات دیتابیس بیش از ۸۰٪ طرفیت Pool شود).

**Graceful Degradation:** اعمال می شود: مثلاً در صورت فشار زیاد روی دیتابیس، در شرایط بحرانی،

آمار بازدید از کش قدیمی ارائه شود.

## منابع الهام‌بخش

- System Design : Scalable URL shortener service like TinyURL•
- a simple link shortener with FastAPI•