# EP31: Super High-performance NoSQL and MQ

**ALEX XU**
**NOV 5, 2022**

♡ 131          ⊙ 4          ↻                                   **Share**     •••

This week's system design refresher:

- ByteByteGo talent collective

- Data Platform

- What is an API gateway

- Super high-performance NoSQL and MQ

- How to scale from 0 to millions of users - spooky edition

## ByteByteGo talent collective

This week has been rough. Lots of layoffs and hiring freezes everywhere. My heart goes out to everyone who is going through this right now.

We are working on compiling a list of interview resources and will share out soon.

We also have ByteByteGo's talent collective here (dozens of companies are hiring) and hope this can be helpful.

# Data Pipeline by Semantix

Very nice illustration of the Data Pipeline by Semantix. It may provide some insights into understanding data pipelines.

The data platform ingests, processes, analyzes and presents data generated by different data sources. A data platform manages all aspects of the data puzzle.

Modern data platforms offer a number of benefits, including centralized access to data across an organization, which eliminates silos and provides actionable insights.

Thanks for reading ByteByteGo Newsletter!
Subscribe for free to receive new posts and
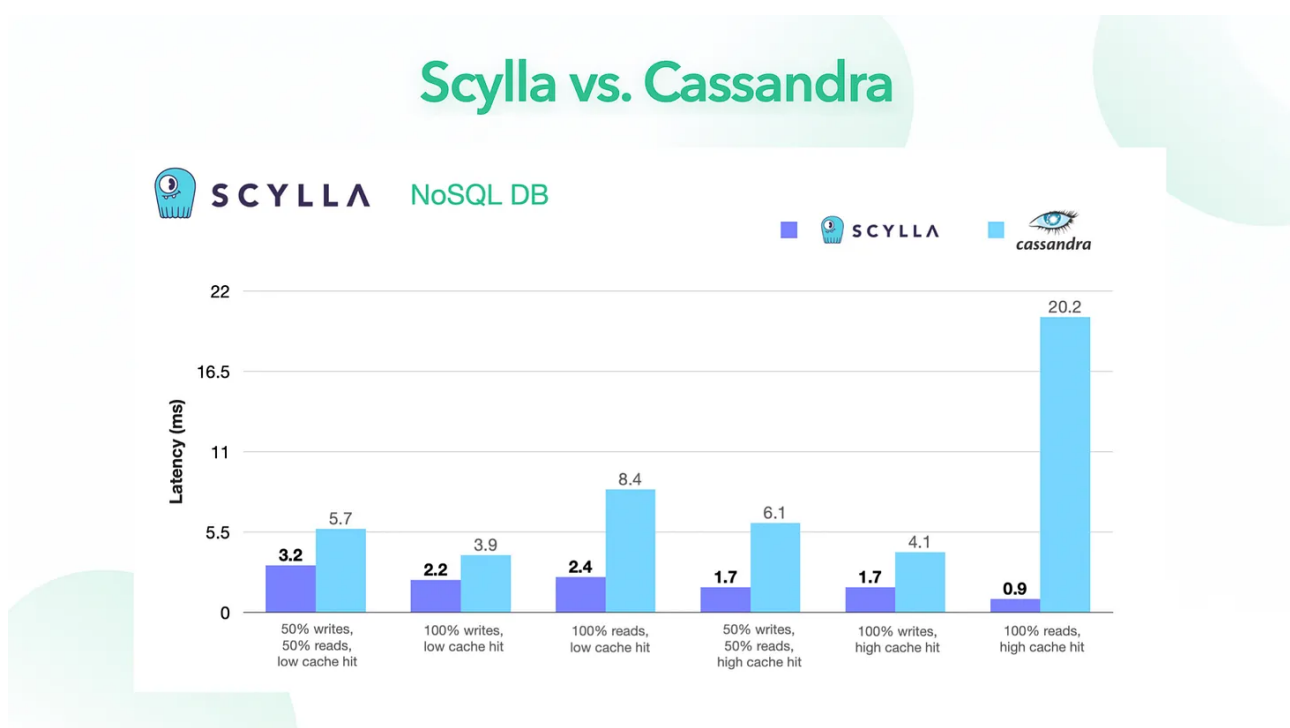support my work.
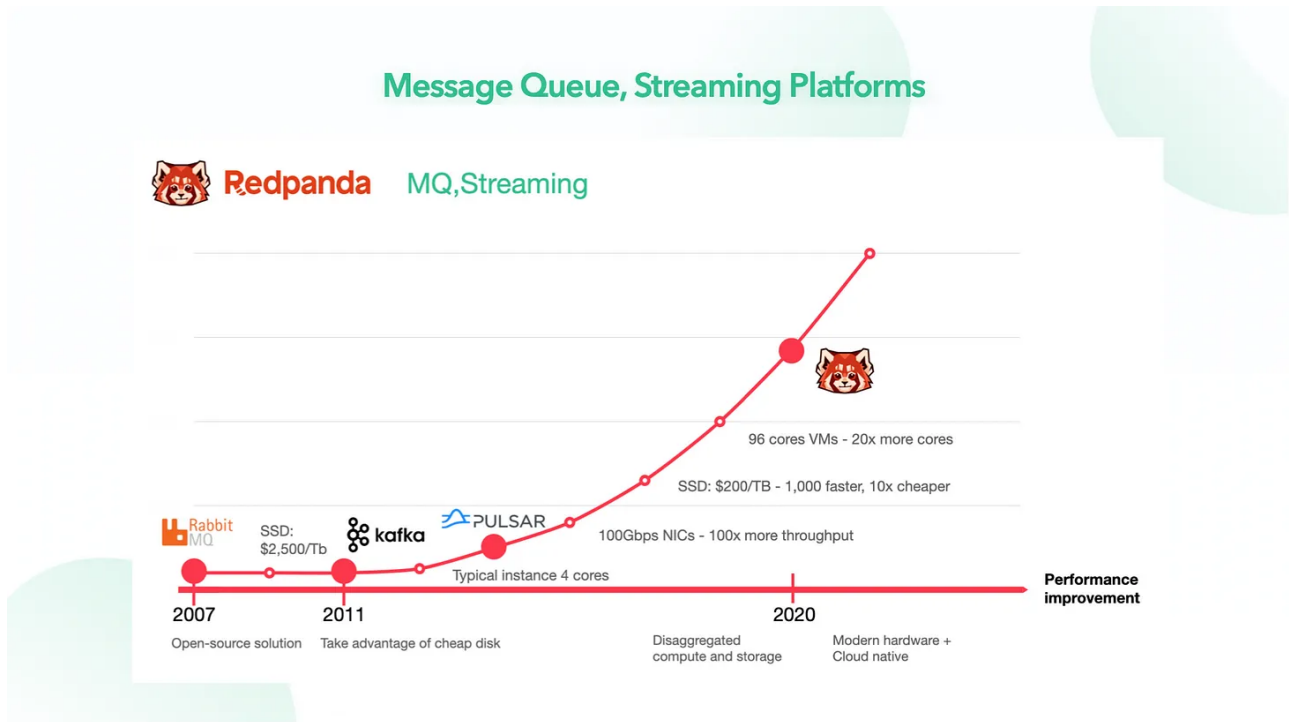
# What does API gateway do?

What is API Gateway?

# Super high-performance NoSQL and MQ

Is it possible to achieve at least a 10x performance boost compared to the original Kafka and Cassandra? How to achieve that? What are the trade-offs?

There is an exciting class of storage software like **Redpanda** and **ScyllaDB** that boasts at least an order of magnitude improvement in performance.



Redpanda and ScyllaDB are used as examples in the diagram below. Redpanda can be compared to Kafka, while ScyllaDB is like NoSQL Cassandra.

- **No JVM, No GC**
  Kafka and Cassandra are written in JVM-compatible languages and usually suffer from high tail latency, where the average latency performs good but 99% latency is not so good due to GC (Garbage Collection).
  Redpanda and ScyllaDB are rewritten from scratch using C++ and leverages some new frameworks (For example, SeaStar). They are hard to code but can achieve much higher performance (see the diagram below for detailed performance metrics).

- **Share-nothing Architecture**
  Every request is pinned to a CPU core. There is no memory contention between cores. This is also friendly to NUMA (Non-Uniform Memory Access) architecture, so that thread can access the memory closer to the CPU core.

- **Zero-copy Networking**
  Using the SeaStar framework, both products can access network devices directly in user mode, and the kernel is not involved. Zero-copy, zero-lock, and zero-
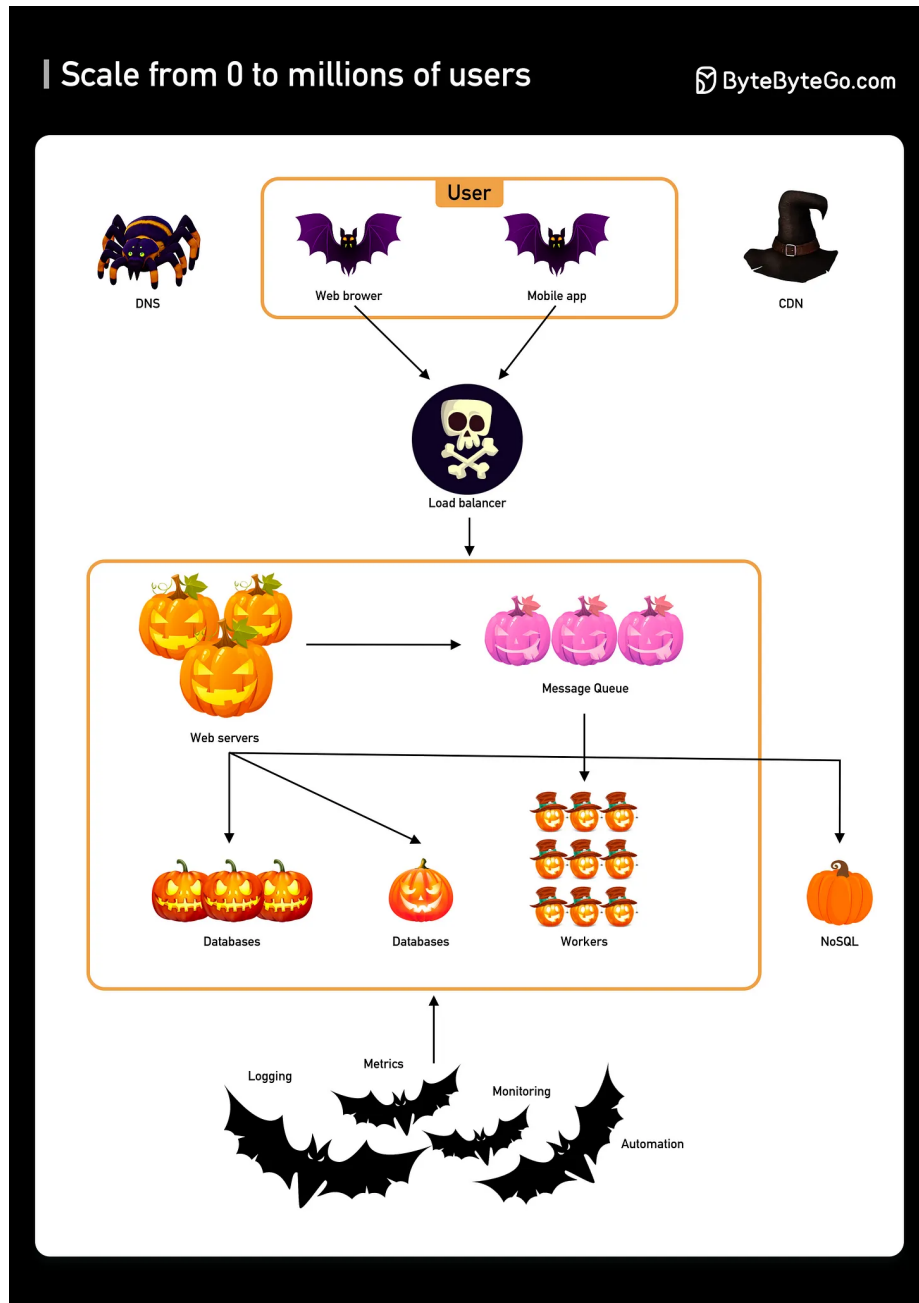
context-switch.

It's been a decade since Apache Kafka, and Apache Cassandra revolutionized how the software industry handled huge amounts of data.

Since then, the server CPU core count has grown 10x. Memory has grown from 64GB to half a TB. NVMe SSD drives are about 100 times faster than spinning disks from a decade ago. Network bandwidth at 25Gbps is commonplace.

A new class of software has come into the market to capitalize on this trend. We wrote this post to raise awareness about this trend.

## How to scale from 0 to millions of users - spooky edition

Designing a system that supports millions of users is challenging, and it is a journey that requires continuous refinement and endless improvement. Let's take a quick look at what are some of the key components powering the system.

**Load balancer**

A load balancer evenly distributes incoming traffic among web servers that are defined in a load-balanced set.

**Web servers**

Web server returns HTML pages or JSON response for rendering.

**Databases**: vertical scaling and horizontal scaling

**Cache**

A cache is a temporary storage area that stores the result of expensive responses or

frequently accessed data in memory so that subsequent requests are served more quickly.

**CDN**

A CDN is a network of geographically dispersed servers used to deliver static content. CDN servers cache static content like images, videos, CSS, JavaScript files, etc.

**Message queue**

A message queue is a durable component, stored in memory, that supports asynchronous communication.

**Logging, metrics, automation**

When working with a small website that runs on a few servers, logging, metrics, and automation support are good practices but not a necessity. However, now that your site has grown to serve a large business, investing in those tools is essential.

Thanks for reading ByteByteGo Newsletter!
Subscribe for free to receive new posts and
support my work.

131 Likes

## 4 Comments

Write a comment...

**Eli** Writes The Machine Learning PhD   Nov 5, 2022

I believe one of the "Databases" pumpkins in the spooky figure is supposed to be "Cache"? Did I understand correctly?

♡ LIKE (3)      💬 REPLY      ⬆ SHARE                                              ⋯

1 reply by Alex Xu

**Hoang Cong Tuan Anh**   Nov 18, 2022

Hi. Can you share which tools did you use to create those diagrams in the post? They have really nice theme. Especially the one with the animation.

Thanks.

♡ LIKE      💬 REPLY      ⬆ SHARE                                                  ⋯

**2 more comments...**