# Large-scale push architectures

How many concurrent connections can one server handle?

**SOLVED**

## C10K
problem

client

client

connection

server

client

client

client

client

## C10M
problem

**SOLVED**

handling concurrent connections ≠ handling concurrent requests

is about efficient scheduling of connections     is about speed of request processing
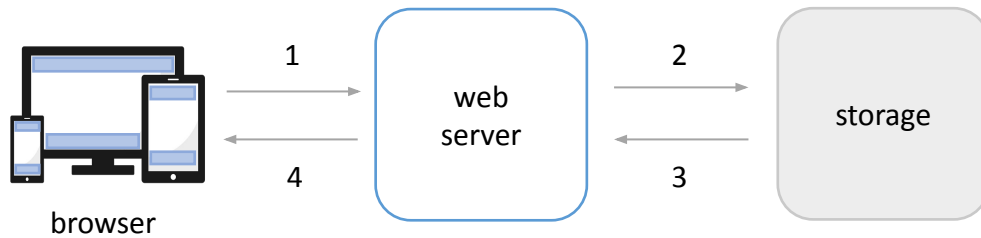
# Large-scale push architectures

## Mail.Ru
email service

### short polling

2017
50,000 requests per second
60% returns empty results

to reduce server load and
speed up mail delivery

### WebSocket

3,000,000 connections per server

browser — 1 → web server — 2 → storage
browser ← 4 — web server ← 3 — storage

storage — 1 → message queue — 2 → web server — 3 ↔ browser

# Large-scale push architectures

## Netflix
streaming service

e.g. Cassandra, DynamoDB, Redis



push registry

register user

look up server
(by user id)

Netflix
client app

API
gateway
(zuul)

message
processor

e.g. Kafka

message
queue

backend
service
A

backend
service
B

backend
service
C

- authentication
- logging
- request routing
- rate limiting
- push messaging
  (WebSocket, SSE)

- authorization
- monitoring
- load balancing
- load shedding

# Large-scale push architectures

| problem | solution |
| --- | --- |
| thread per connection model scales poorly for long lived connections | choose non-blocking IO servers for push architectures |
| server restarts | • migrate connections without reconnecting clients<br>• force clients to reconnect to a different server |
| server failures | prefer multiple small servers to one big server |
| older load balancer versions cut WebSocket connections | • use load balancers with native support for WebSockets<br>• balance the load at layer 4 (TCP) instead of layer 7 (HTTP) |