

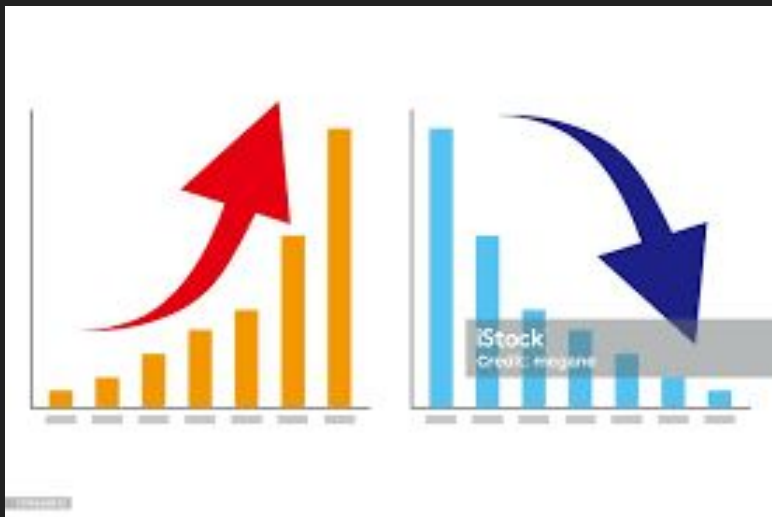
[1장] 사용자 수에 따른 규모 확장성

가상 면접 사례로 배우는 대규모 시스템 설계 기초

이민석 / unchaptered

확장성(scalability)이란 무엇인가?

클라우드에서 **확장성**은 증가하는 트래픽(Traffic)에 맞춰서 시스템의 규모가 확장되거나 축소될 수 있는 특징



그렇다면 가용성(Availability)는 무엇인가?

시스템이 필요한 시간에 **항상 사용 가능한 상태**를 유지하는 능력

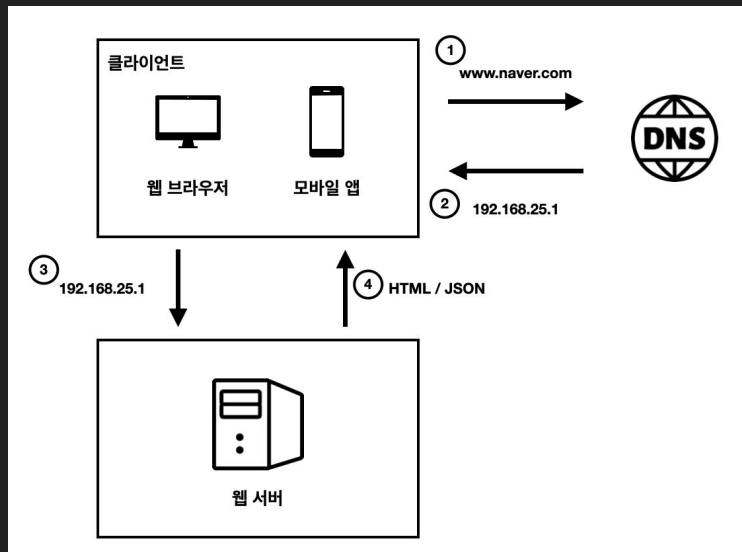
즉, 어떤 상황에서도 서비스가 다운되지 않는다면 가용성이 매우 높다고 볼 수 있습니다.

대부분의 클라우드 아키텍처는

가용성 / 확장성을 비롯하여 다양한 요소를 신경써서 설계를 진행함
따라서 이런 관점에 따라서 1장을 읽어보았습니다.

기본적인 단일 서버

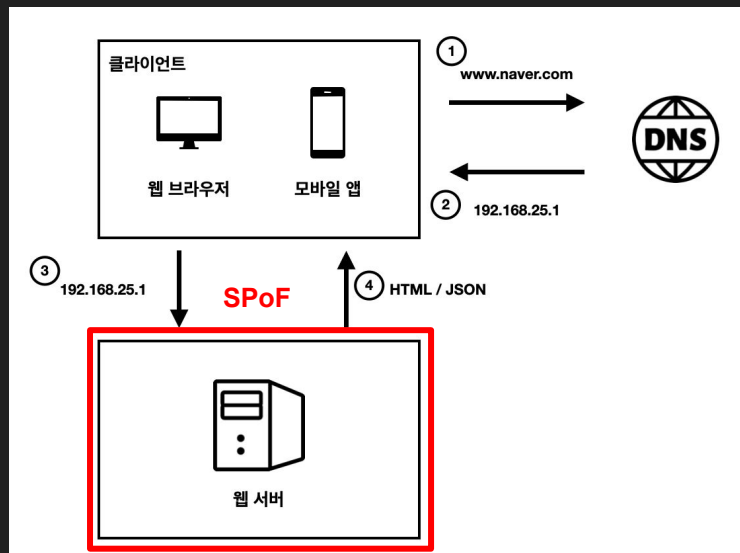
- DNS 서버
- 페이지(Frontend)
- 서버 (Backend) + 데이터베이스 (Database)



<https://ssdragon.tistory.com/152>

단일 서버의 SPoF 문제점

SPoF(Single Point of Failure)는
한 기능을 담당하는 시스템의 일부가 죽으면
시스템 전체가 마비되는 **지점(부분)**을 의미

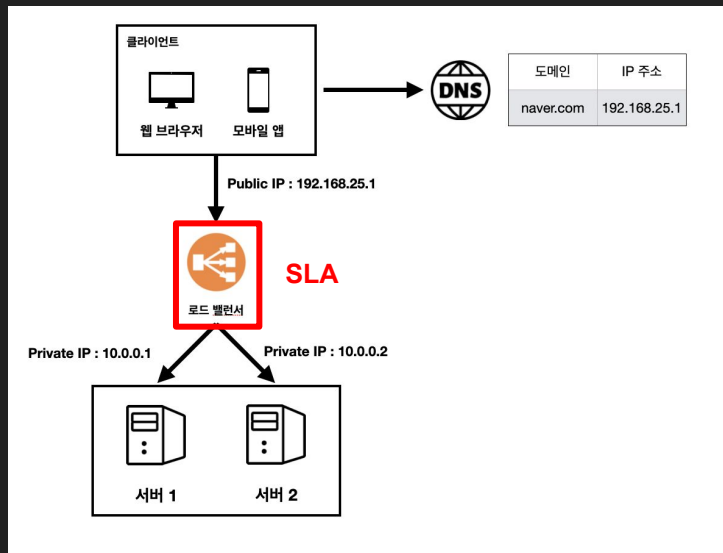


<https://ssdragon.tistory.com/152>

핵심은 SPoF를 줄여서 고가용성의 서비스를 구축

서버가 N개로 늘어나면서 서버 가용성 개선

하지만 로드 밸런서가 새로운 SPoF로 보임



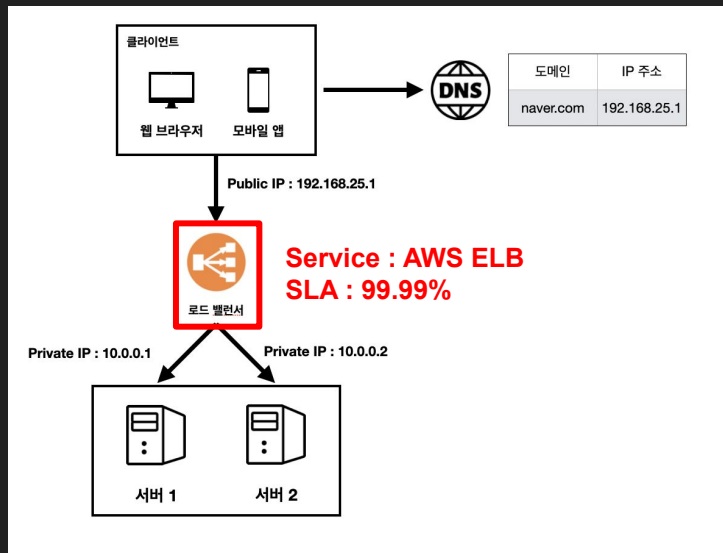
<https://ssdragon.tistory.com/152>

모든 서비스의 SPoF를 줄이기 어렵다.

SPoF를 피할 수 없는 경우,
CSP 업체는 일부 가용성을 보장

AWS의 로드 밸런서의 경우,
AWS ELB(ALB, L7)가 있으며,
99.99%의 SLA를 보장한다.

- 하루 8.6 초 / 주 1분 / 월 4분 21초 / 연 52분



<https://ssdragon.tistory.com/152>

<https://aws.amazon.com/ko/elasticloadbalancing/>

https://d1.awsstatic.com/legal/AmazonElasticLoadBalancing/Amazon_Elastic_Load_Balancing_Service_Level_Agreement_2022-07-25_KO-KR.pdf

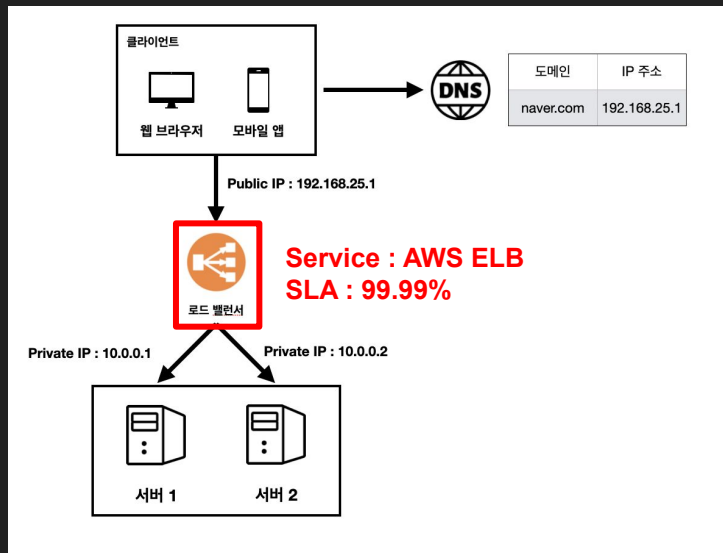
<https://uptime.is/>

그리고 모든 SPoF를 줄이는 것이 옳지 않을 수 있다.

SPoF를 피하기 위해서 다중화를 하는 경우,
다중화 갯수에 만큼 비용 증가가 발생한다.

또한 인프라 복잡성 증가에 따라,
개발 및 운영 복잡성이 증가한다.

따라서 SPoF 지점에 대해서 인지하고
비용/가용성/확장성에 **타협**이 필요할 수 있다.



<https://ssdragon.tistory.com/152>

<https://aws.amazon.com/ko/elasticloadbalancing/>

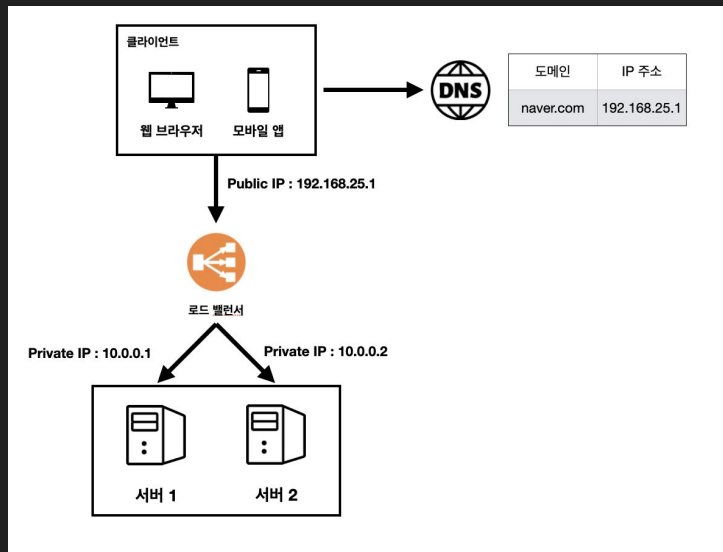
https://d1.awsstatic.com/legal/AmazonElasticLoadBalancing/Amazon_Elastic_Load_Balancing_Service_Level_Agreement_2022-07-25_KO-KR.pdf

<https://uptime.is/>

서버의 수평적 확장

기존에는 단일 서버였으나
현재는 서버1과 서버2로 확장되었습니다.

이를 수평적 확장이라고 부릅니다.



<https://ssdragon.tistory.com/152>

<https://aws.amazon.com/ko/elasticloadbalancing/>

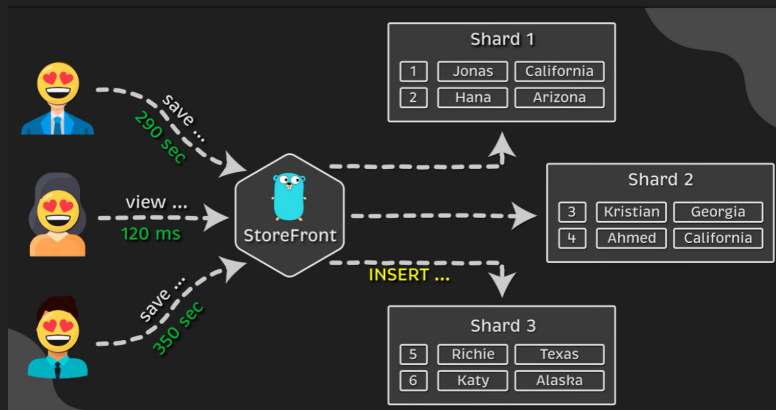
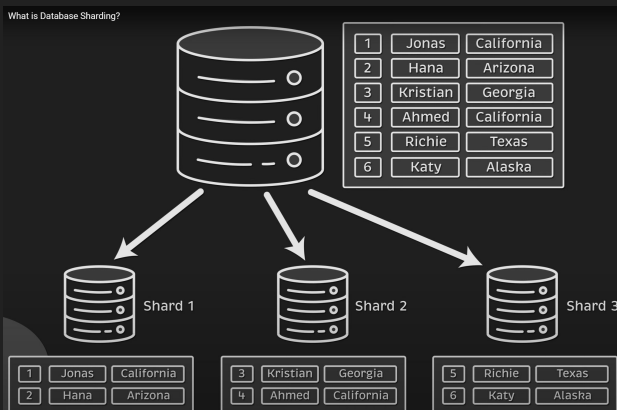
https://d1.awsstatic.com/legal/AmazonElasticLoadBalancing/Amazon_Elastic_Load_Balancing_Service_Level_Agreement_2022-07-25_KO-KR.pdf

<https://uptime.is/>

다중 데이터베이스

수직적 확장 : 단일 데이터 베이스는 치명적인 SPoF로 작용할 수 있다.

수평적 확장 : 단일 데이터 베이스를 여러 개의 샤드(Shard)로 분할하는 기술
모든 샤드는 같은 스키마를 쓰지만 데이터에는 중복이 없음
하지만 샤딩은 **복잡성을 증가시킴**



🚩 샤드 ? 샤딩 ?

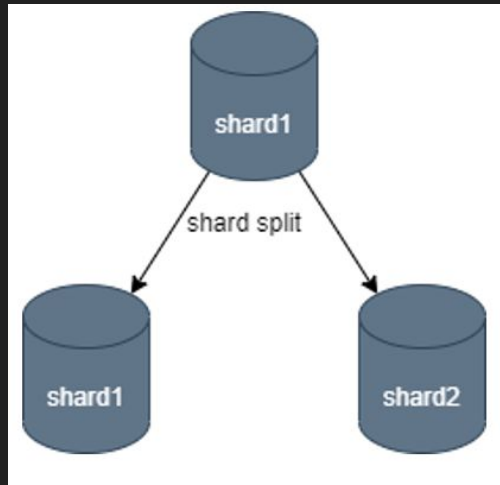
파티션키(**Partition Key**)에 따라서 샤딩 진행

샤드 소진(**Shard Exhaustion**)이 있으면, 재샤딩 진행

유명인사(**Celebrity**) 이슈가 발생할 수 있음

What is the Database Sharding?

- https://youtu.be/XP98YCr-iXQ?si=u_XTF7SKiaUcsTJO



유명인사(Celebrity) 이슈

한 샤드에 유명인사가 포함되어 있으면,

한 샤드에 **read** 요청이 집중적으로 몰리게 됩니다.

이 현상을 유명인사 이슈라고 부릅니다.

Alibaba Cloud의 핫스팟 키 검색 및 공통 솔루션

- https://www.alibabacloud.com/blog/redis-hotspot-key-discovery-and-common-solutions_594446

🚩 파티션 키의 종류에 따라

Range-based Sharding

- 'A' to 'I' → Shard 1
- 'J' to 'S' → Shard 2
- 'T' to 'Z' → Shard 3

Hashed Sharding

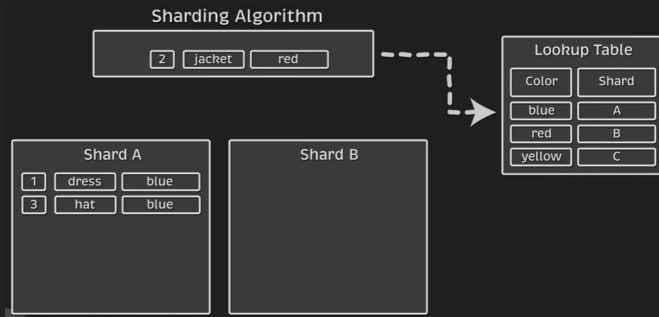
Directory Sharding

Geo Sharding

hash_function(, ,) = 2



Hashed Sharding

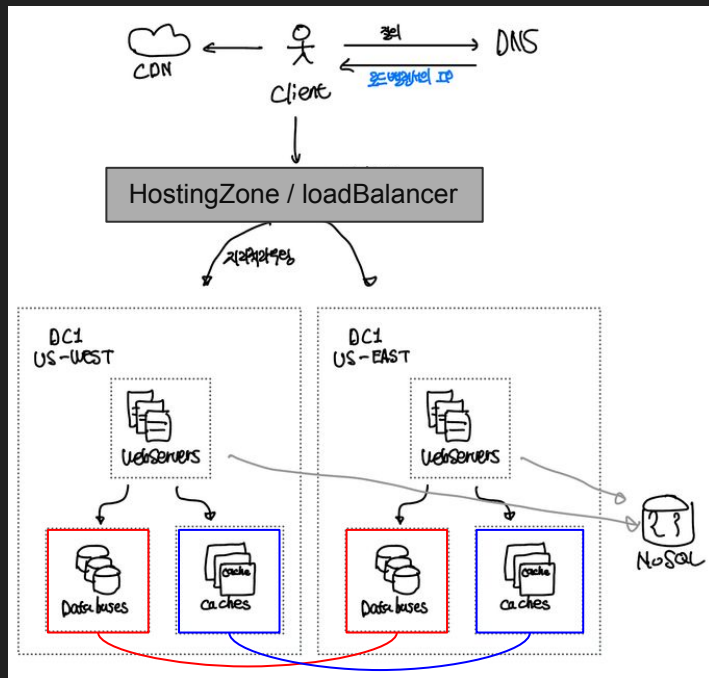


Directory Sharding

데이터 센터(지리적 라우팅)

일정한 조건에 맞는 곳에 라우팅

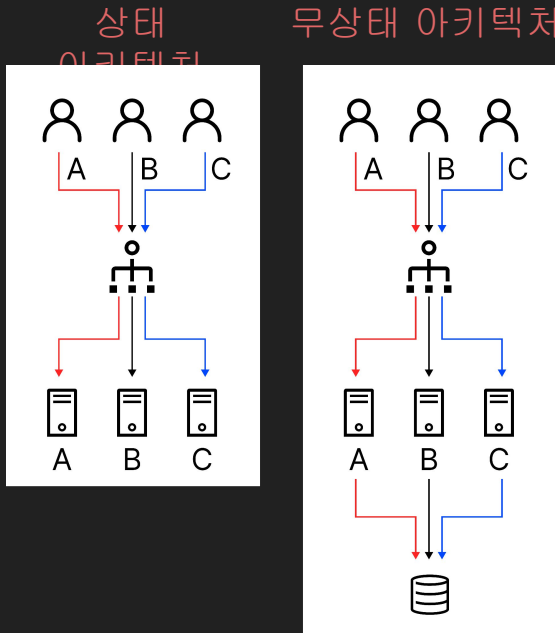
- 트래픽 우회
- 데이터 동기화
- 테스트와 배포



상태(stateful)와 무상태(stateless)

JWT가 아닌 Session 방식의 서비스에서

별도의 Session Storage의 유무에 따라서 상태 아키텍처 / 무상태 아키텍처 구분



캐시(Cache)

주로 많이 접근을 하는 데이터들에 캐시(Cache)를 도입하거나...

[NHN FORWARD 2021] Redis 야무지게 사용하기

조회수 3.1만회 • 2년 전



세션 설명 레디스, 얼마나 알고 사용하고 계신가요? 레디스를 사용하는 사람들이라면 이것만큼은 꼭 알아야 할 내용들을 공유하고자 ...

자막

<https://youtu.be/92NizoBL4uA?si=rCs6mboYZIfNIGB7>

콘텐츠 전송 네트워크(CDN)

많이 제공 되는 파일(프론트엔드 페이지)들은 CDN을 사용한다던가

What is Cloudflare CDN?

조회수 1.1만회 • 1년 전

 Work Smart

what is Cloudflare CDN? How does it work? And why use it? Learn more: →
[https://www.cloudflare.com/what-is-cloudflare/ ...](https://www.cloudflare.com/what-is-cloudflare/)

자막

https://youtu.be/ZBR2Ub325Uc?si=8OMUz6MB3zecwA_G

Edge?

전세계 주요 지점에 거점을 만들고 한 번 조회한 데이터는 거점에 저장됨



🚩 Edge Computing, Edge Function

전세계 주요 지점에 서버가 있으면 Edge Computing 함수가 있으면 Edge Function



메세지 큐

독립적인 2개의 서비스를 하나의 서비스로 **유연하게 연결**하기 위해

What is a Message Queue?

조회수 7.1만회 · 2년 전



IBM Technology ✓

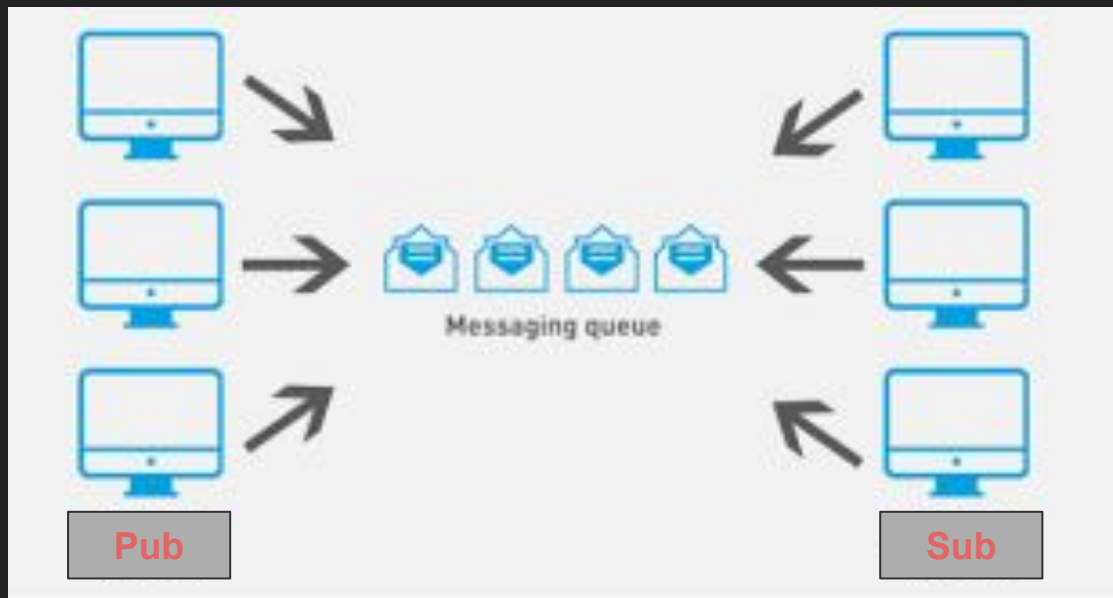
A message queue is a component of messaging middleware solutions that enables independent applications and services to ...

자막

<https://youtu.be/xErwDaOc-Gs?si=ttm2YzzNCfZ9Y2sx>

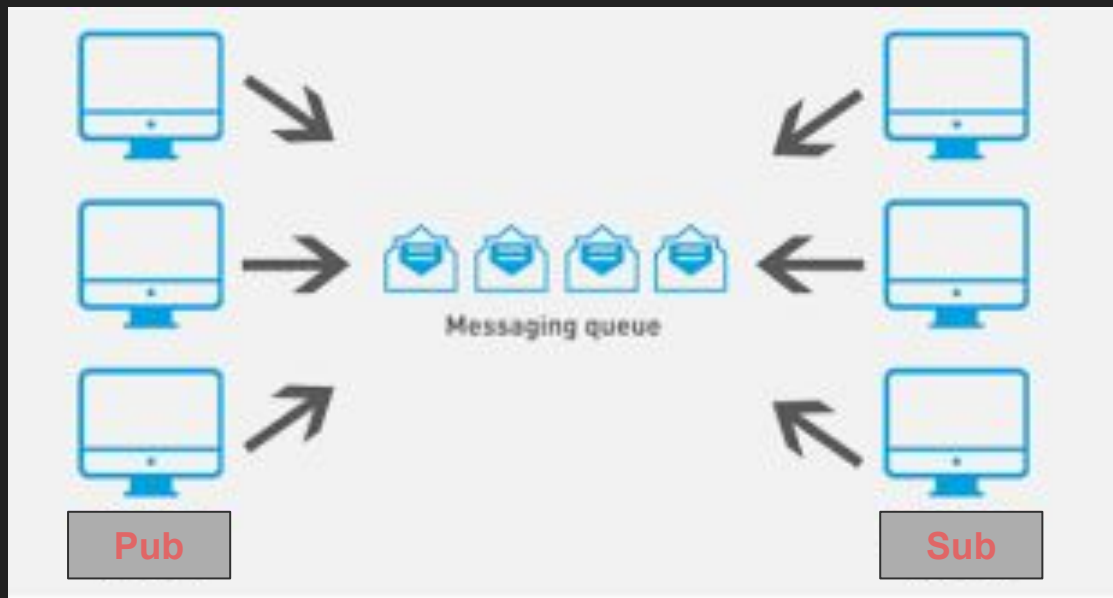
🚩 Producer&Consumer or Publisher&Subscriber

메세지(작업)을 만드는 사람을 Pub으로 이를 꺼내서 사용하는 사람을 Sub으로



🚩 Decouple?

개별 모듈(pub/sub) 간 연결고리를 끊고 의존성을 없앤 것을 **Decouple**이라 한다.



로그, 매트릭 그리고 자동화

어플리케이션(L7) 계층에서 기록한 로그

다양한 어플리케이션(L7)과 시스템 구성요소의 지표(CPU/MEM/Networking I/O)를 부르는 매트릭

시스템을 더욱 안정적이고 견고하게 만들어주는 자동화 등

백만 사용자, 그 이상

이 작업들을 반복하면서 더욱 **고가용성** 및 **고확장성**에 맞게 시스템을 고도화해 나가면서 **xx Million** 단위의 트래픽을 버틸 수 있는 서비스로 성장시킬 수 있다.

단, 이 과정에서의 **비용(노동, 유지/보수, 사용량)**을 반드시 고려해야한다.

감사합니다