Fail Fast, Learn Faster SRE (실패에서 배워나가는 SRE)

김재헌, 유호균 Service System Reliability

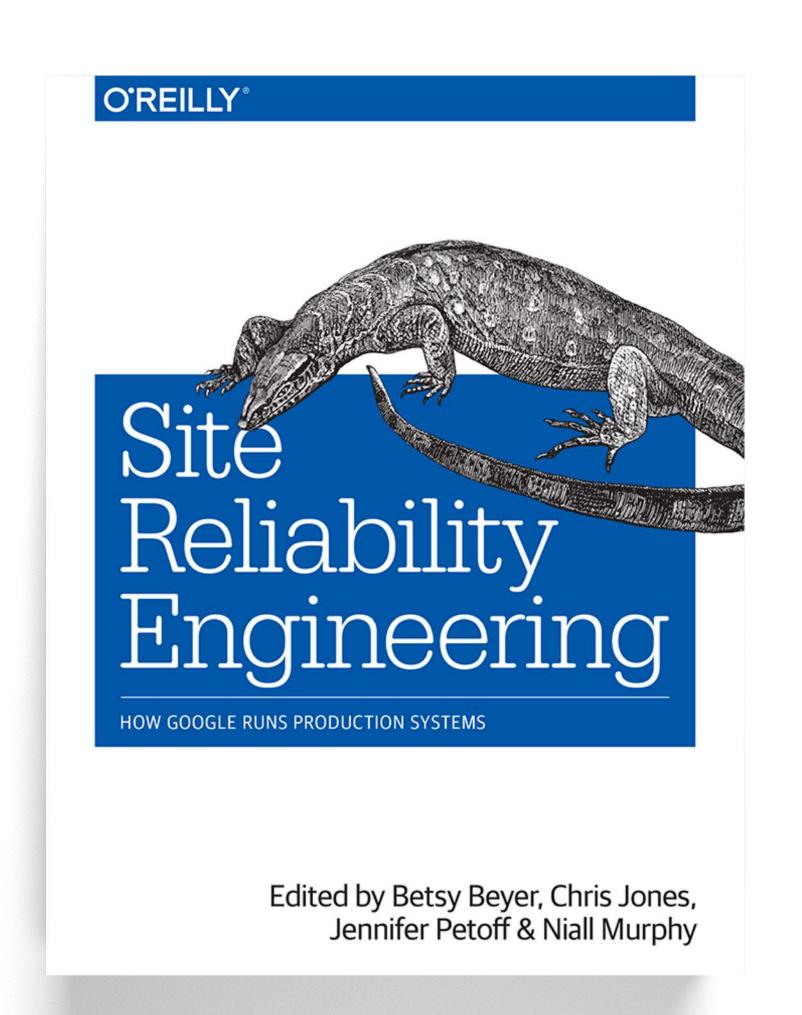
NAVER

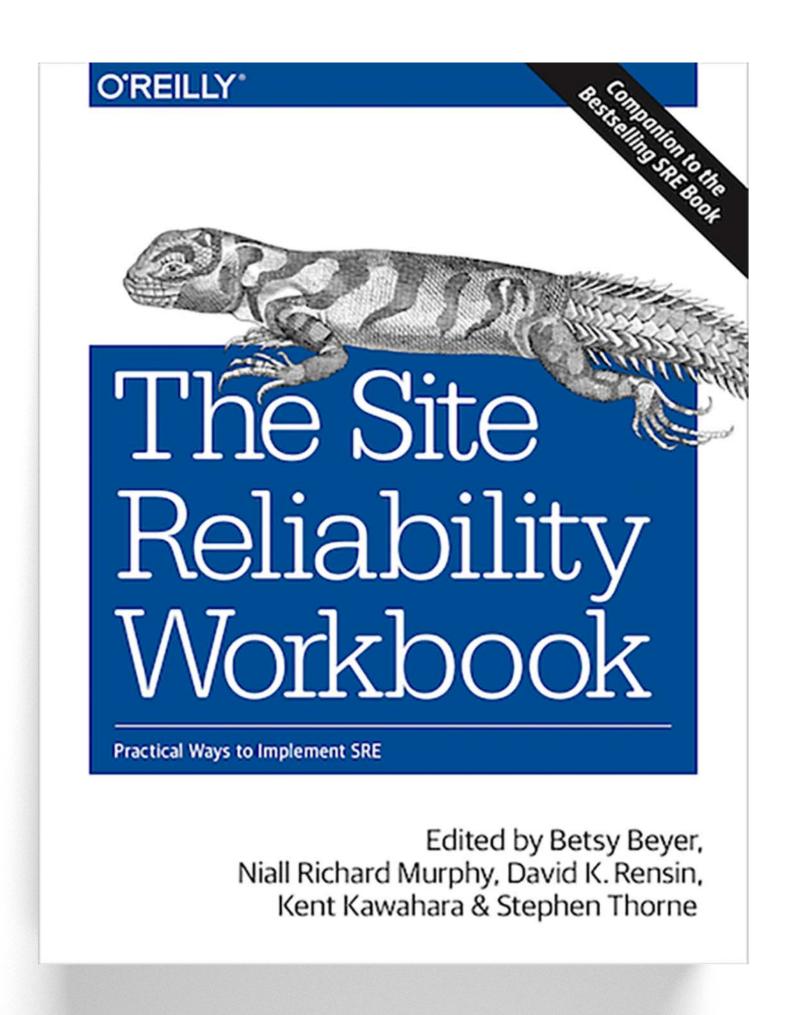
CONTENTS

- 1. 왜 SRE를 도입해야 하는가?
- 2. 우리들의 SRE
- 3. Metric + Meta Data = Insight
- 4. Fail Fast, Learn Faster
- 5. Lessons Learned

SRE: Site Reliability Engineering 사이트 신뢰성 엔지니어링

글로벌 스케일, 또는 대규모의 인터넷 서비스를 제공하면서 어떻게 하면 시스템의 신뢰성을 보장할 수 있을지 고민하는 기술 분야이자 방법론, 문화



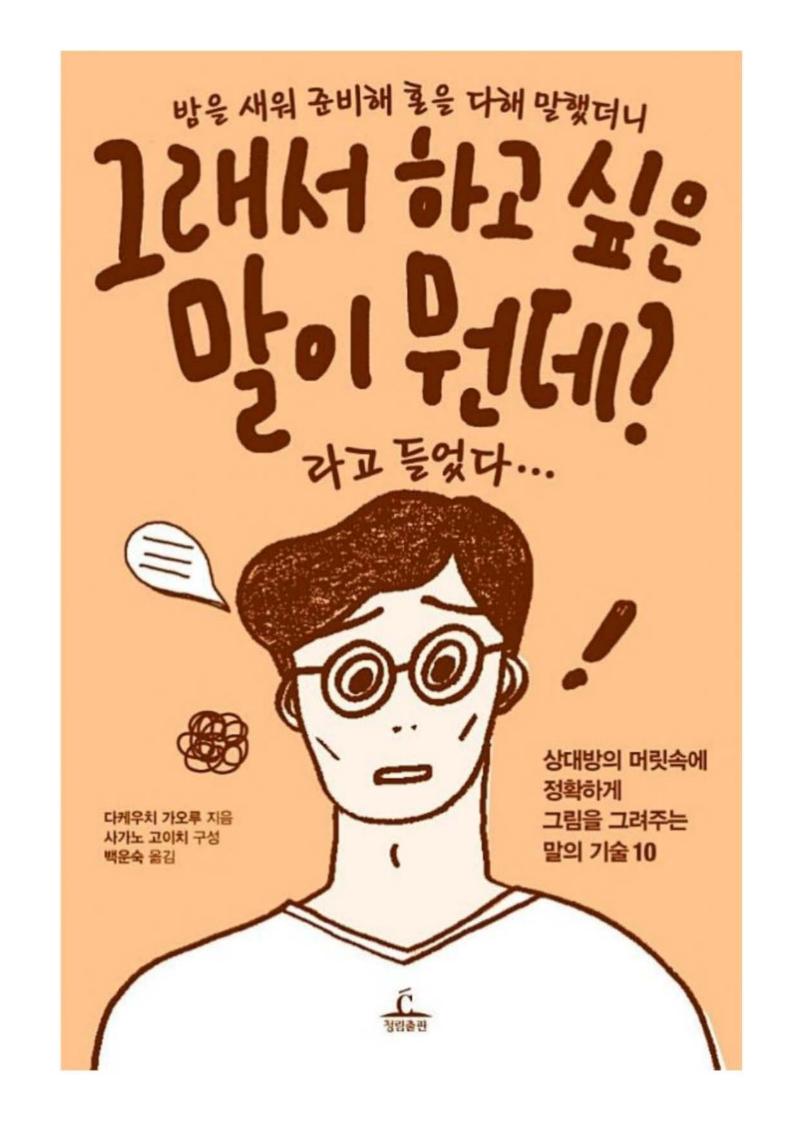


https://landing.google.com/sre/books

SRE: Site Reliability Engineering 사이트 신뢰성 엔지니어링

SLI, SLO, SLA, Error Budget, DevOps, Post-Mortem, Blameless, Toil management, Silo, Availability, Emergency Response, Capacity Planning, Change Management, On-Call, Monitoring

네... 좋은 말씀 잘 들었습니다...



SRE: Site Reliability Engineering 사이트 신뢰성 엔지니어링

글로벌 스케일, 또는 대규모의 인터넷 서비스를 제공하면서 어떻게 하면 시스템의 신뢰성을 보장할 수 있을지 고민하는 기술 분야이자 방법론, 문화

네이버 검색 SRE팀의 정의

신뢰성: 네이버 검색 서비스가 정상적으로 제공되는 것

방법론: Availability, Capacity Planning, Monitoring, Contingency Plan

SRE: Site Reliability Engineering 사이트 신뢰성 엔지니어링

글로벌 스케일, 또는 대규모의 근록 네 나 사고하면서 어떻게 하면 시스템의 신뢰성을 보장할 꼭 있을 지하는 요? 분야이자 방법론, 문화

네이버 검색 SRE팀의 정의

신뢰성: 네이버 검색 서비스가 정상적으로 제공되는 것

방법론: Availability, Capacity Planning, Monitoring, Contingency Plan,

SRE: Site Reliability Engineering
사이트 신뢰성 엔지니어링

글로벌 스케일/네이버 검색은 365일 / 24시간당게 하면 시스템의 신뢰성을 무충단제공되어야 한다^{다자 방법론, 문화}

네이버 검색 SRE팀의 정의

신뢰성: 네이버 검색 서비스가 정상적으로 제공되는 것

방법론: Availability, Capacity Planning, Monitoring, Contingency Plan,

만약...

1. 사용자가 검색을 했는데..



2. 검색 불가



검색은절대서비스제공이

중단되면 안됩니다.

온국민이 사용하는 공공재니까요!

SRE: Site Reliability Engineering 사이트 신뢰성 엔지니어링 2.

글로벌 스케일, 도선비수, 규모가 컨죄는 라쿠어떻게 하면 시스템의 신뢰성엔처 비어는 '늘어나지 않는다 방법론, 문화

네이버 검색 SRE팀의 정의

신뢰성: 네이버 검색 서비스가 정상적으로 제공되는 것

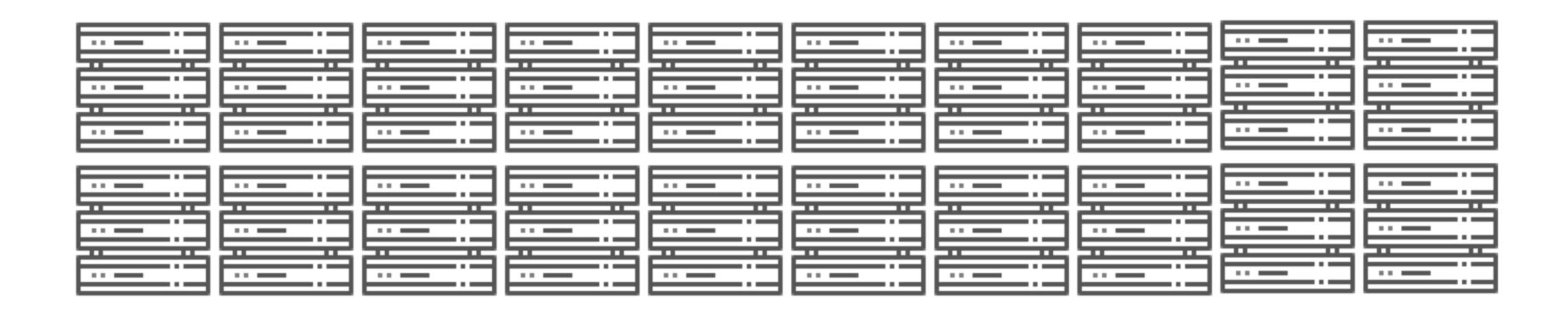
방법론: Availability, Capacity Planning, Monitoring, Contingency Engineering,

몇대의 서버는 1명이 관리 가능



몇대의 서버는 1명이 관리 가능

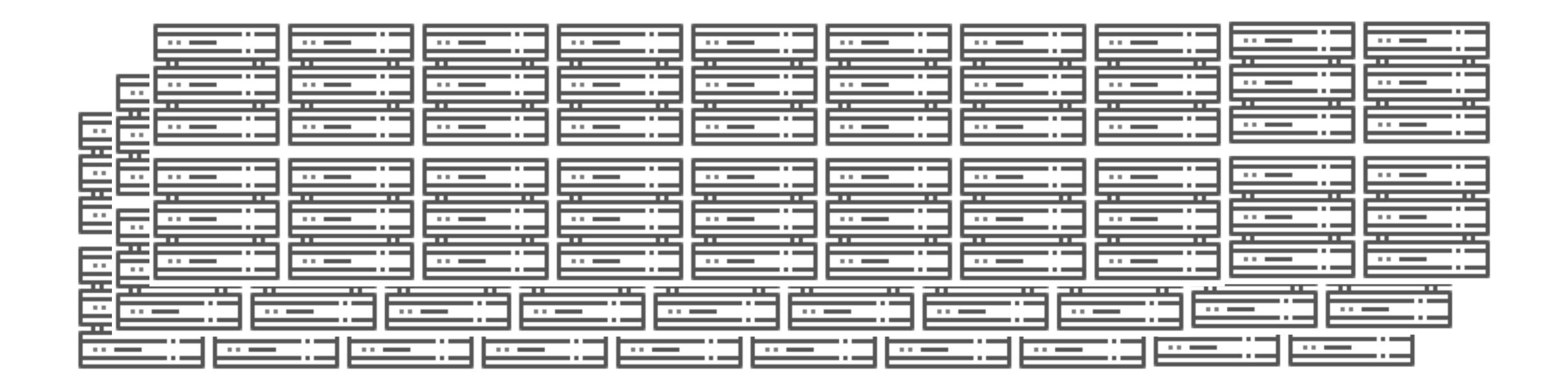
수십~ 수백대의 서버는 몇명이 관리 가능



몇대의 서버는 1명이 관리 가능

수십 ~ 수백대의 서버는 몇명이 관리 가능

그렇다면, 수천 ~ 수만대의 서버는?



서버는 무한히 늘어나도, 사람은 무한히 늘어날 수 없다.

SRE: Site Reliability Engineering 사이트 신뢰성 엔지니어링 3.

물로발스케일, 지진과 같은 <mark>재난상황</mark>에서는 명기하면 시스템의 기존의 대응방안이 동작하지 않는다. 문화

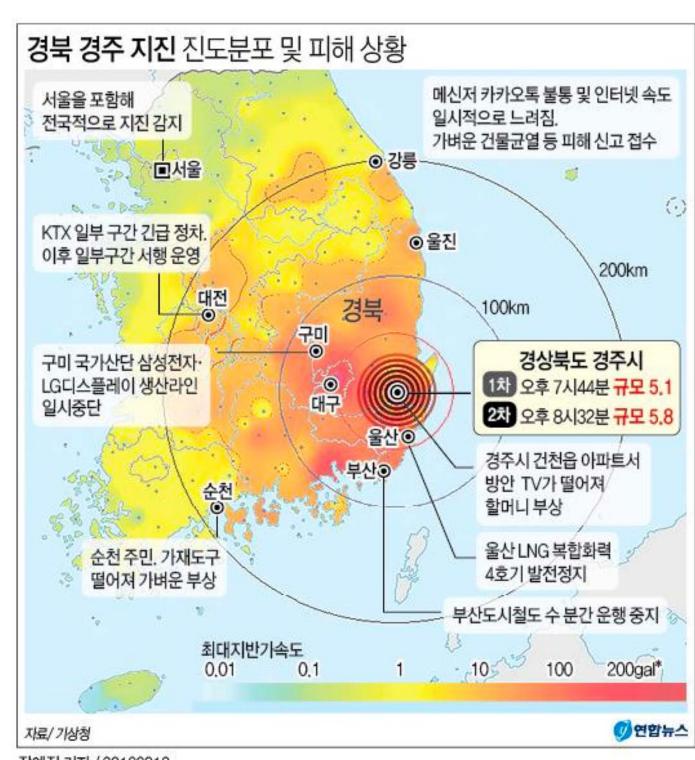
네이버 검색 SRE팀의 정의

신뢰성: 네이버 검색 서비스가 정상적으로 제공되는 것

방법론: Availability, Capacity Planning, Monitoring, Contingency Engineering,

미중 무역협상		Q
연예인 사망	——	Q
광화문 시위	——	Q
19호 태풍 하기비스		Q

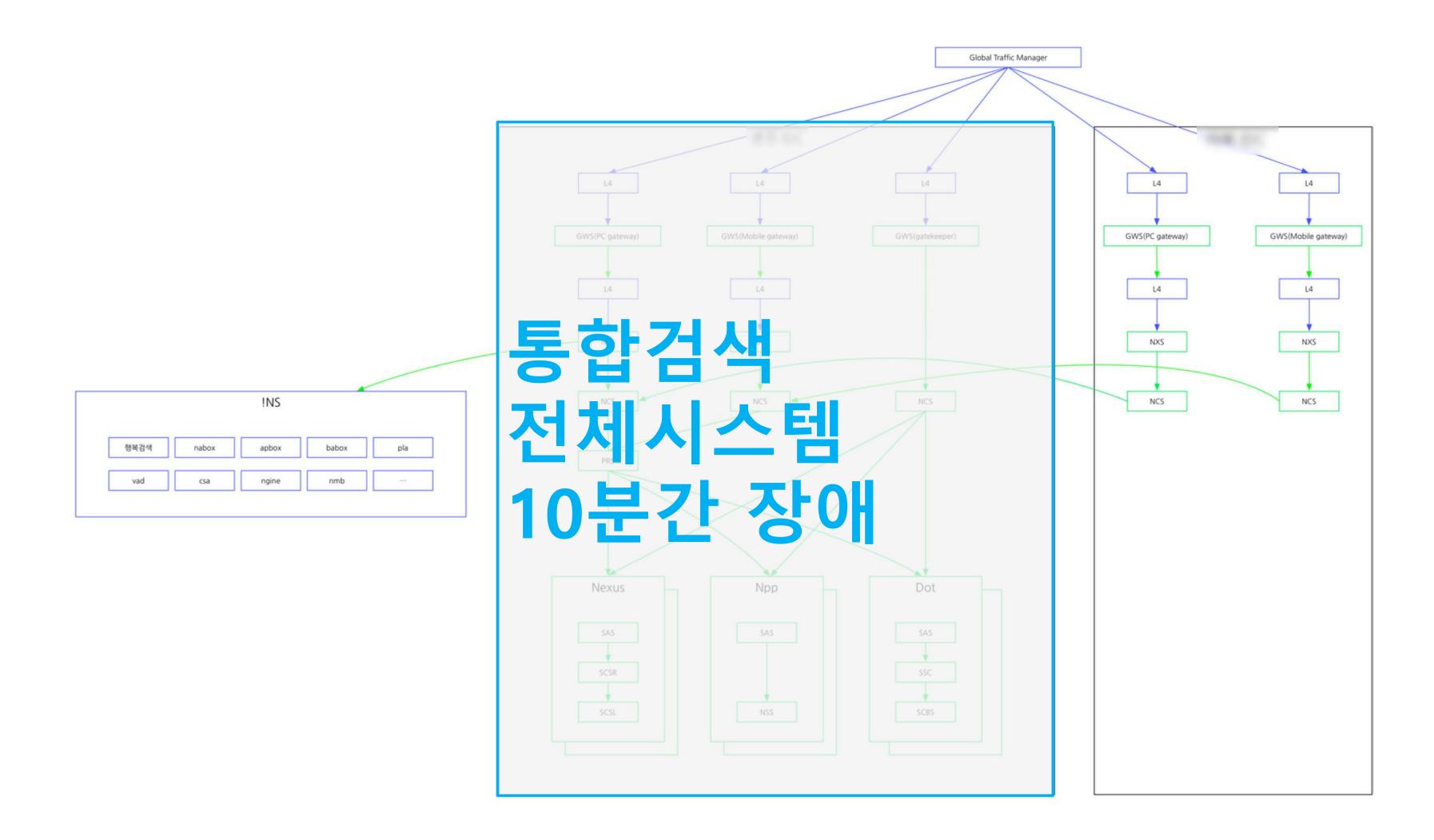
대한민국의 모든 이슈가 트래픽에 반영

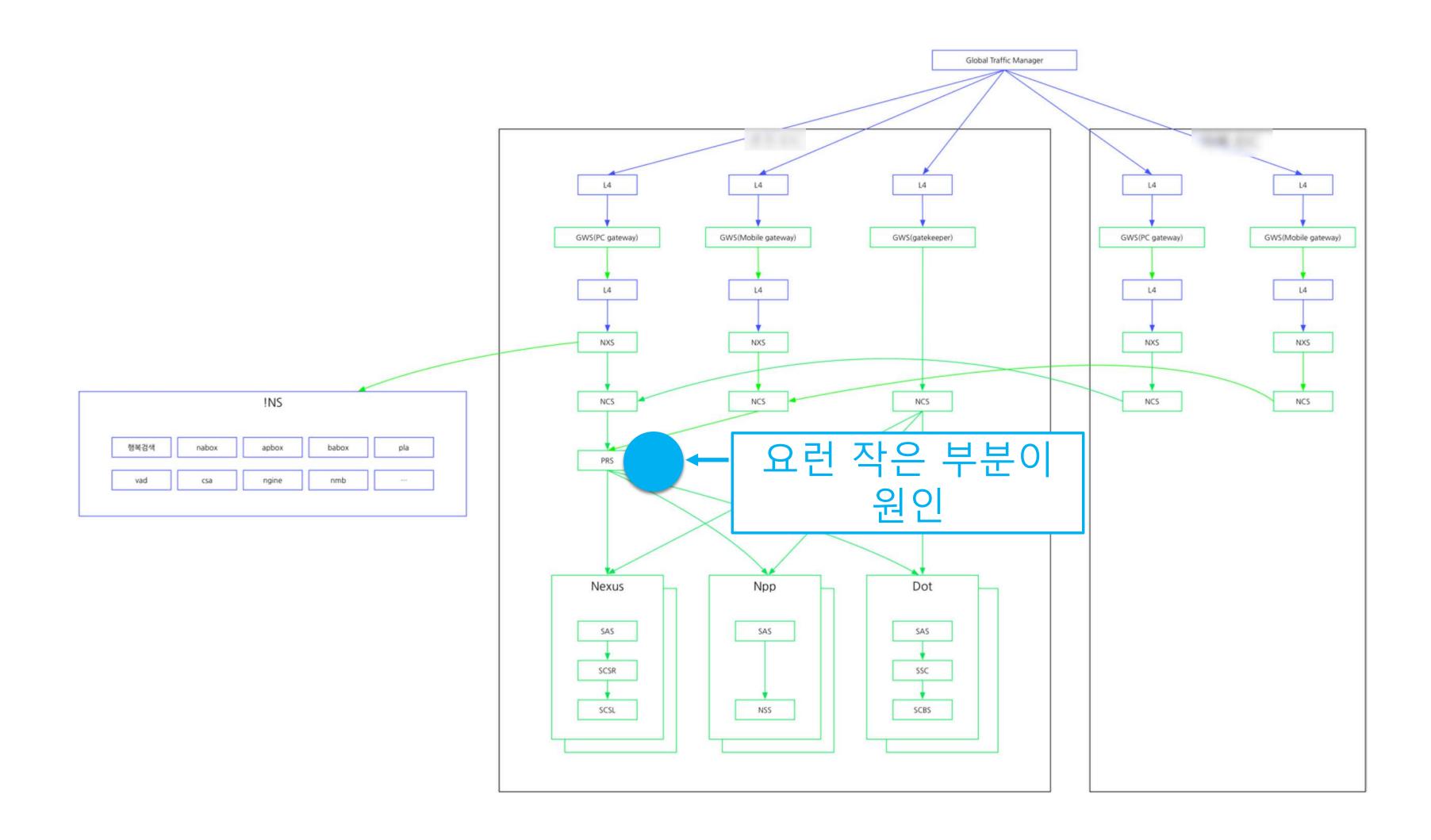


장예진 기자 / 20160912 트위터 @yonhap_graphics, 페이스북 tuney.kr/Le/N1

2016년 9월 경주 지진

지진	-	Q
경주 지진	—	Q
재난문자	—	Q





장애 탐지 : 2분

증상 완화 : 20분

원인 파악 : 40분

장애 영향도 파악 : 1시간

장애 복구 완료 : 수시간

포스트모텀 작성 : 48시간

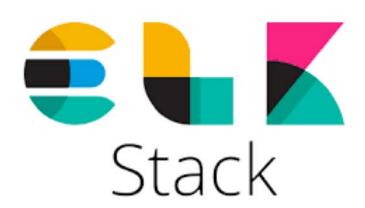
기존의 DevOps 방법론으로는 규모도, 재난도 감당할 수 없는 현실

SRE의 필요성을 절실히 느끼고 2016년 말 도입을 결정







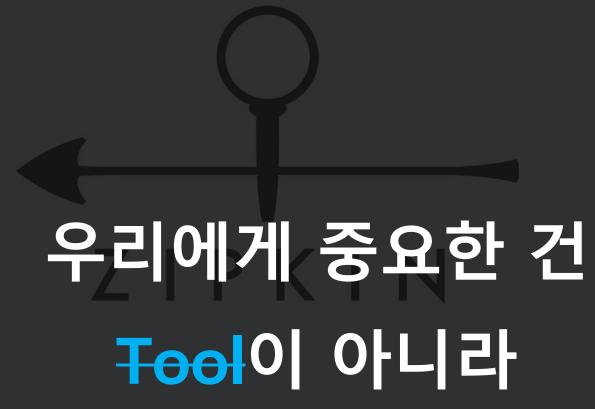


















문제점

- 1. 비상상황 대응 체계의 부재
- 2. 장애 발생 Detection의 어려움
- 3. Post-Mortem 문화 부재
- 4. 전체 시스템 현황 확인 불가

문제점

1. 비상상황 대응 체계의 부재

- 비상 상황(Contingency)이 무엇인지 정의가 없음
- 비상 상황에 어떤 행동을 해야 하는지 대응 체계가 없음

검색 contingency plan

개요

- 2가지 관점의 contingency plan 운영 (이하 CP)
 - o 경쟁사 서비스 불능에 대비한 CP_1
 - 국가적 이슈 발생에 대비한 CP_2

CP_1

- 정의
 - o daum, google 등의 경쟁사 서비스 마비로 naver 검색 유입이 증가
 - 검색 순유입량(unique query) 폭증
 - ㅇ 장시간 지속 가능
 - ㅇ 검색시스템 전체적인 부하 증가로 이어짐
- 대비
 - o N배 가용량 확보를 위한 scale out 정책 수립
 - o layer 간 cache server 도입
 - o IDC 이중화
 - 통합검색 일간 peak traffic 기준 배의 가용률 준비
 - 모든 경쟁사 outage에 대응 가능한 수준
- CP_1 발동
 - o peak 기준 배 이상의 traffic 유입시 CP_1 발동
 - ㅇ 성능 병목이 예상되는 검색서비스 판단하여 통검 결과에서 제외 준비
 - ㅇ 단계적 결과 제외로 검색시스템 신뢰성 유지

CP_2

- 정의
 - 자연 재해를 비롯한 국가적 이슈 발생으로 naver 검색 유입이 증가
 - 동일한 검색어 유입량(repetitive query) 폭증
 - o 이슈 발생 시점의 burst한 traffic 유입
 - 검색시스템 Front-end layer에 부하 집중
- 대비
 - 9/12 경주지진 계기로 CP_2 대응 시작
 - 서버의 역할(role) 별로 peak 기준 최대가용률 파악
 - 병목 지점의 scale out 정책 수립
 - 통합검색 일간 peak traffic 기준 배의 가용률 준비
 - 경주지진 발생 시점의 traffic 유입에 대응 가능
- CP_2 발동
 - 。 국가 차원의 이슈 발생시 CP_2 발동
 - 성능 병목이 예상되는 role의 scale out 수행하여 peak 기준 배 가용률 확보
 - CP_2 발동 중 CP_1 상황이 동시 발생할 경우 통검 결과 선별적 노출

DEVIEW 2019

부자

경쟁사 서비스 장애로 인한

트래픽 유입 집중상황 CP_1

해야 하는지 대응 체계가 없음

검색 contingency plan

개요

- 2가지 관점의 contingency plan 운영 (이하 CP)
 - 경쟁사 서비스 불능에 대비한 CP_1
 - 국가적 이슈 발생에 대비한 CP_2

CP_1

- 정의
 - o daum, google 등의 경쟁사 서비스 마비로 naver 검색 유입이 증가
 - o 검색 순유입량(unique query) 폭증
 - ㅇ 장시간 지속 가능
 - ㅇ 검색시스템 전체적인 부하 증가로 이어짐
- 대비
 - o N배 가용량 확보를 위한 scale out 정책 수립
 - o layer 간 cache server 도입
 - o IDC 이중화
 - 통합검색 일간 peak traffic 기준 배의 가용률 준비
 - 모든 경쟁사 outage에 대응 가능한 수준
- CP_1 발동
 - o peak 기준 배 이상의 traffic 유입시 CP_1 발동
 - ㅇ 성능 병목이 예상되는 검색서비스 판단하여 통검 결과에서 제외 준비
 - ㅇ 단계적 결과 제외로 검색시스템 신뢰성 유지

CP_2

- 정의
 - 자연 재해를 비롯한 국가적 이슈 발생으로 naver 검색 유입이 증가
 - 동일한 검색어 유입량(repetitive query) 폭증
 - 이슈 발생 시점의 burst한 traffic 유입
 - o 검색시스템 Front-end layer에 부하 집중
- 대비
 - 9/12 경주지진 계기로 CP_2 대응 시작
 - o 서버의 역할(role) 별로 peak 기준 최대가용률 파악
 - 병목 지점의 scale out 정책 수립
 - 통합검색 일간 peak traffic 기준 배의 가용률 준비
 - 경주지진 발생 시점의 traffic 유입에 대응 가능
- CP_2 발동
 - o 국가 차원의 이슈 발생시 CP_2 발동
 - 성능 병목이 예상되는 role의 scale out 수행하여 peak 기준 배 가용률 확보
 - CP_2 발동 중 CP_1 상황이 동시 발생할 경우 통검 결과 선별적 노출

DEVIEW 2019

부자

국가적 관심사로 인한

폭발적인 트래픽 급증 CP_2

해야 하는지 대응 체계가 없음

문제점

2. 장애 발생 Detection의 어려움

- 장애를 detection하는 기준이 필요함
- 아주 쉬운 기준으로!

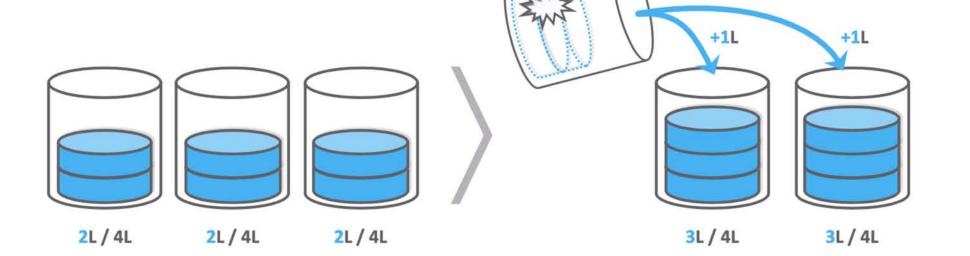
므데저

가용량 지표 개발

새로운 방식

부하증가배수

한 친구가 죽으면 나머지 친구들은 몇 배를 받나?

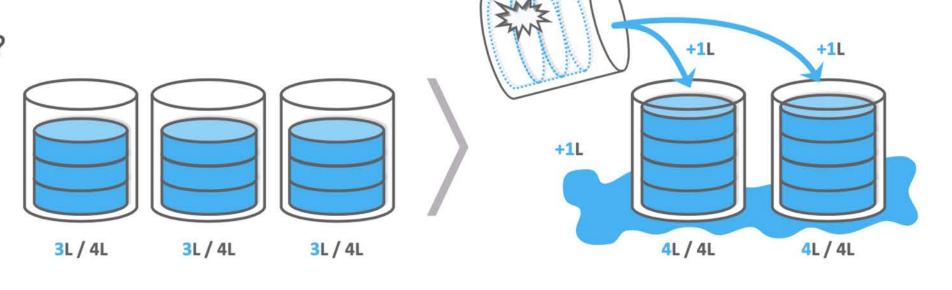


최대가용배수

한 명이 현재 몇 배까지 받을 수 있나?

"임계 상황" 판단

부하증가배수 > 최대가용배수



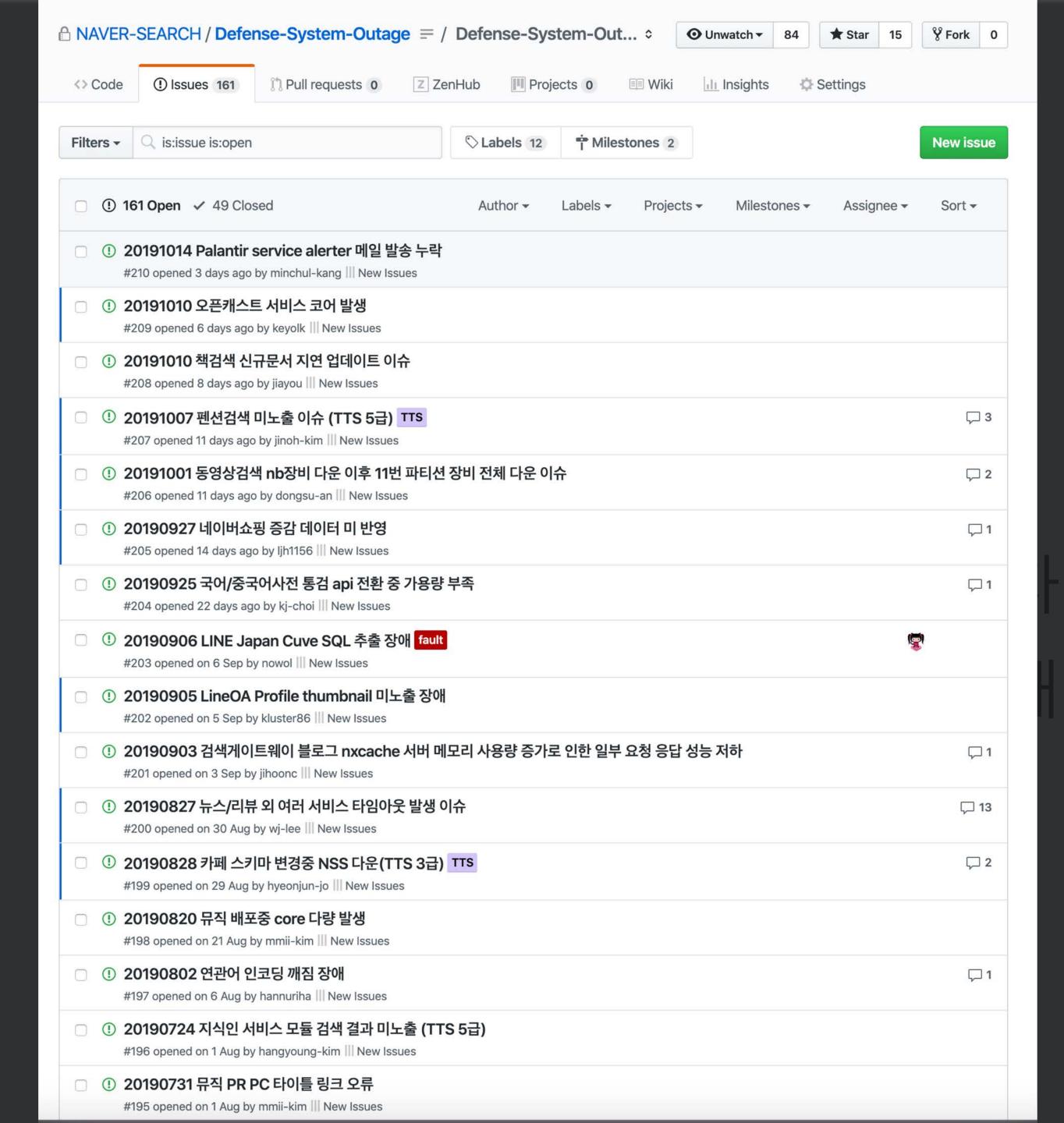
아주 쉬운 공식

가용량 지표를 정의, 특허 출원 가용량을 기반으로 한 Alerting

문제점

3. Post-Mortem 문화 부재

- 장애 복구 후 사후 처리 Report가 누락
- 비슷한 장애의 재발 시 대응 대책이 미비



DEVIEW 2019

2016년 말
Post-Mortem 작성 포맷 공유
처음 2-3 개월은 SRE가 작성
사후 분석을 습관화하는 문화를 전파 현재는 자율적 && 상시 운영

문제점

4. 전체 시스템 현황 확인 불가

- 다양한 서비스의 연계 장애 및 전체 현황 파악 불가
- Detection을 통해서 대응 필요시 Alerting이 필요함

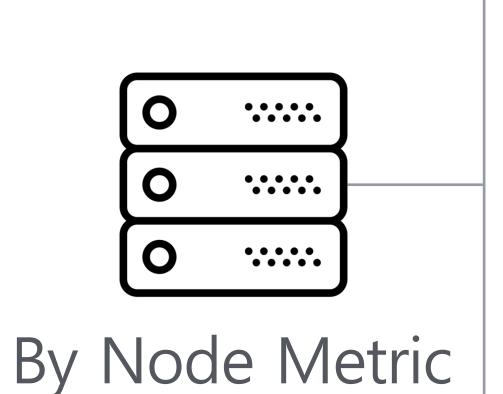
전체 현황을 파악하는

- 시스템은 어찌 만들어야 하지?

기존에도,

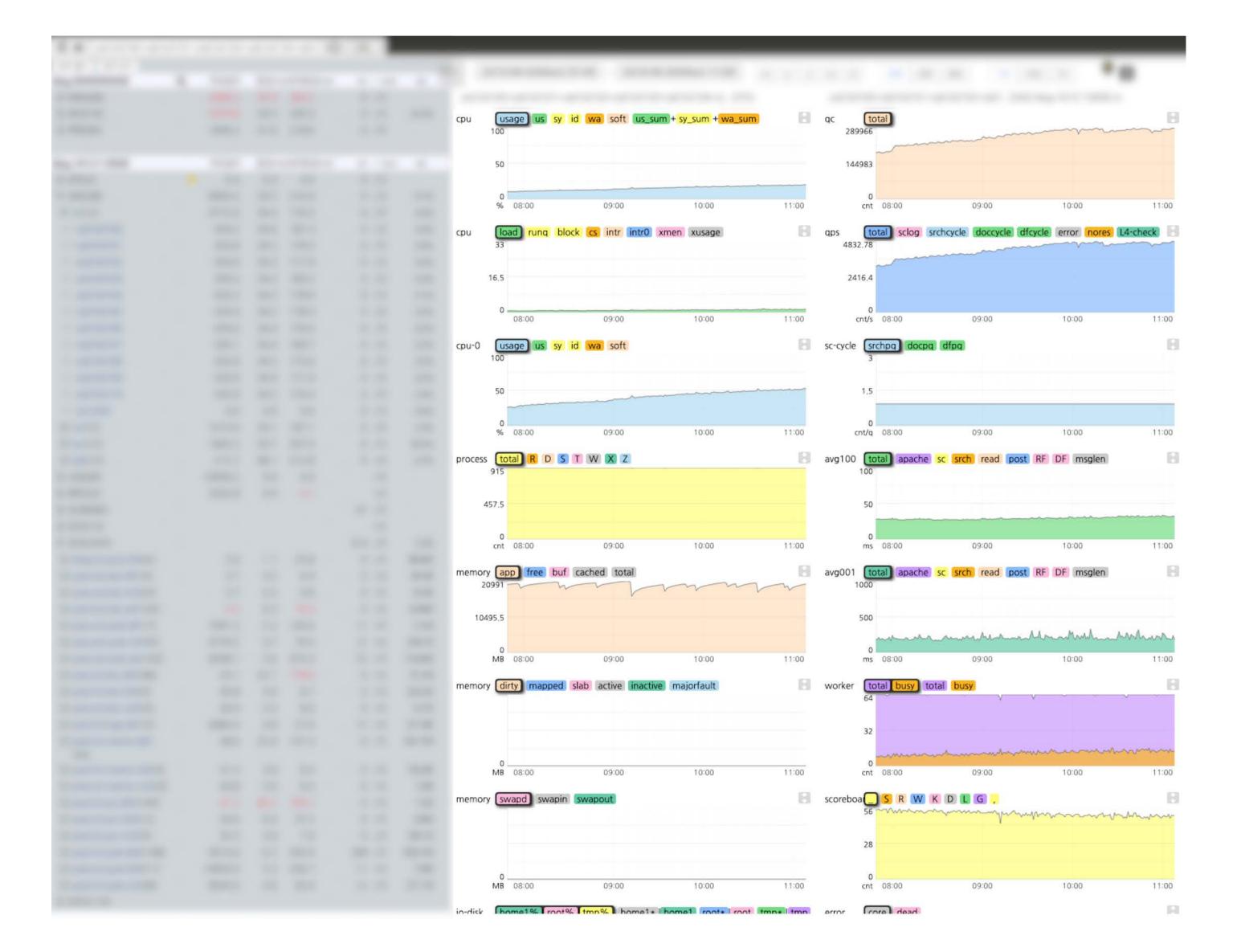
Monitoring은 하고 있었다.

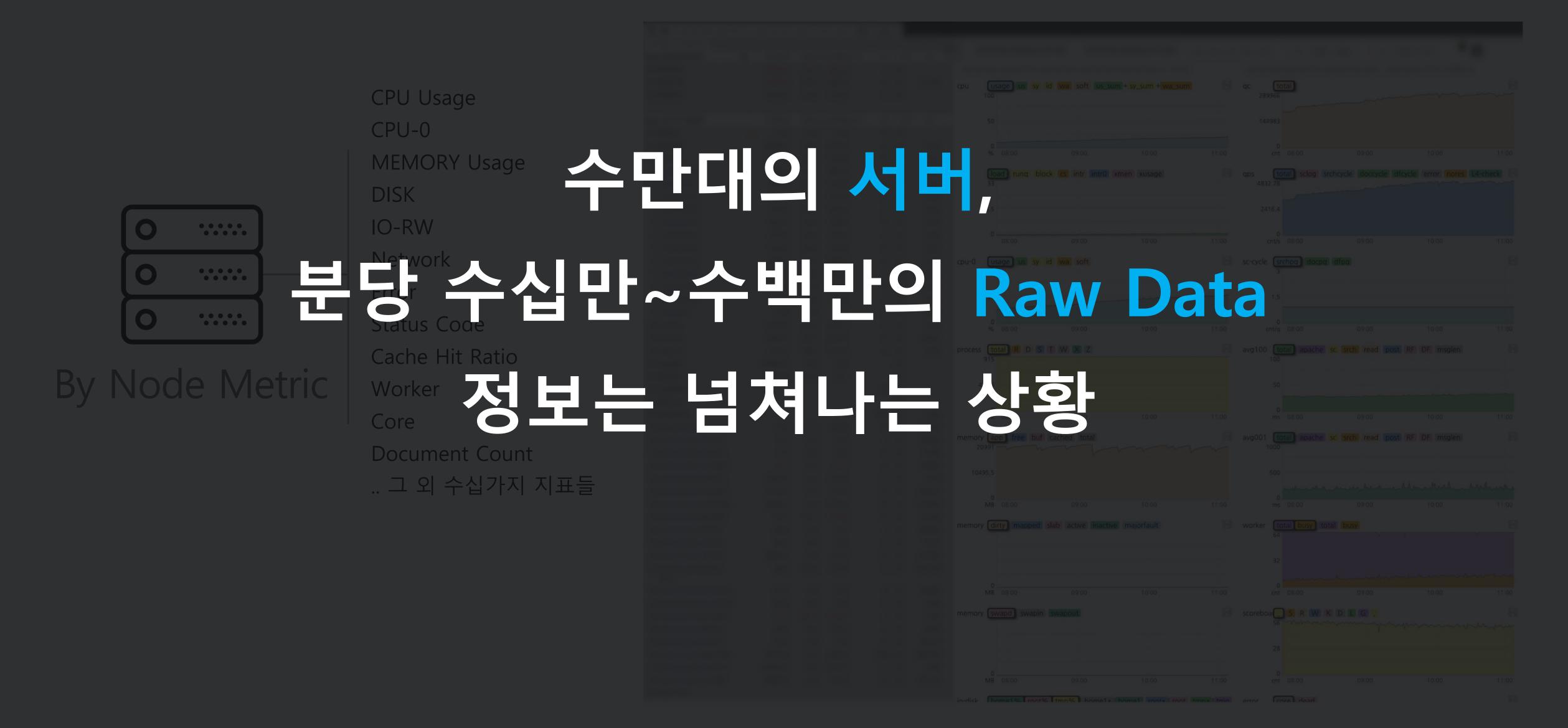
Metric도 잘 쌓고 있었다.

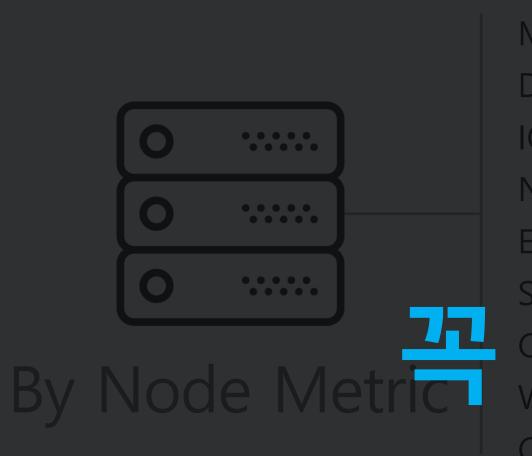


CPU Usage
CPU-0
MEMORY Usage
DISK
IO-RW
Network
Error
Status Code
Cache Hit Ratio
Worker
Core
Document Count

.. 그 외 수십가지 지표들







But,
MEMORY Usage
DISK
PHI 처나는 정보를 이용하여
Error
Status Code
Unsight를 제공하려면?

Core

Document Count

.. 그 외 수십가지 지표들



1. 서비스만에 정보 자동 Aggregation

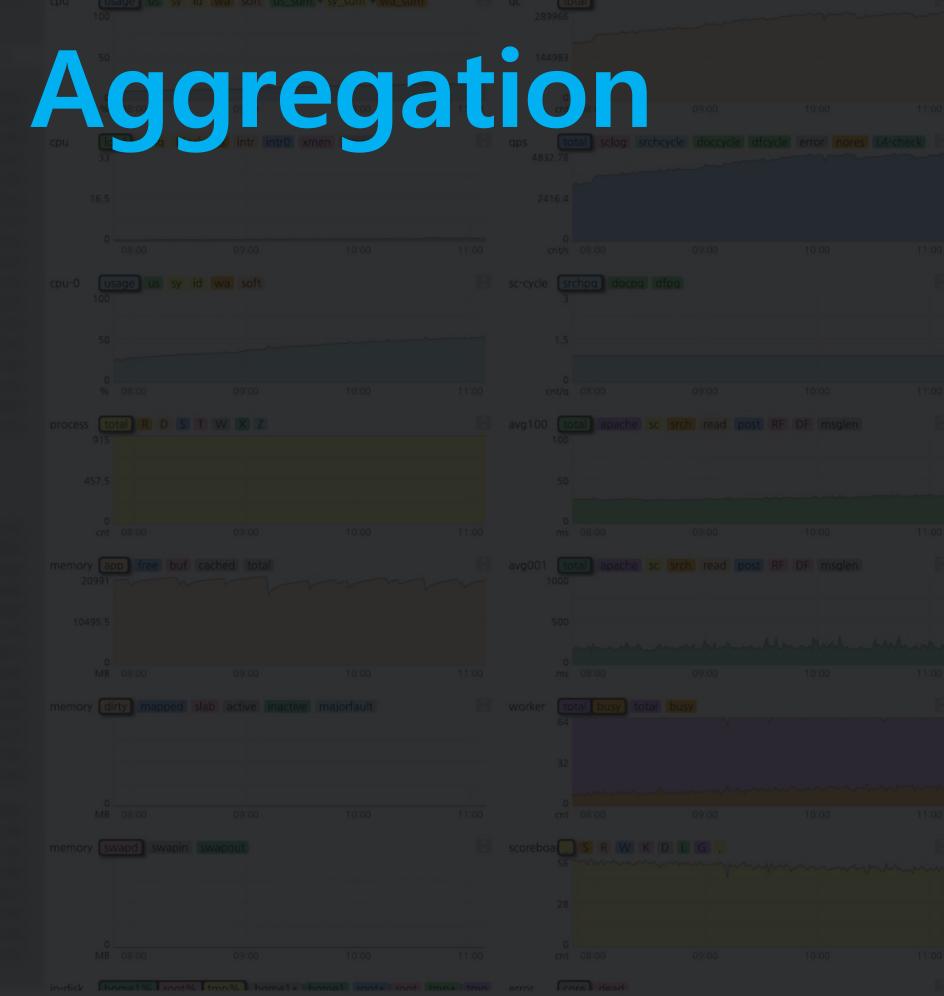
응 2. 담당자 정보의 정확도

By Node Metric

Status Code
Cache Hit Ratio
Worker
Core

Document Count

.. 그 외 수십가지 지표들



(1). SSUID (서비스 유니크 아이디) 발급

Needs: 복잡한 구성의 서비스를 묶어서 가시화 하고싶다.

Solution : 서비스에 유니크한 아이디를 부여, 모든 계층을 묶음으로 가시화

지역 검색

- Cache Layer : local, local_openapi

- View Layer : locallink, localnx

- Search Layer: m_local, nx_local, openapi_local, svc_local

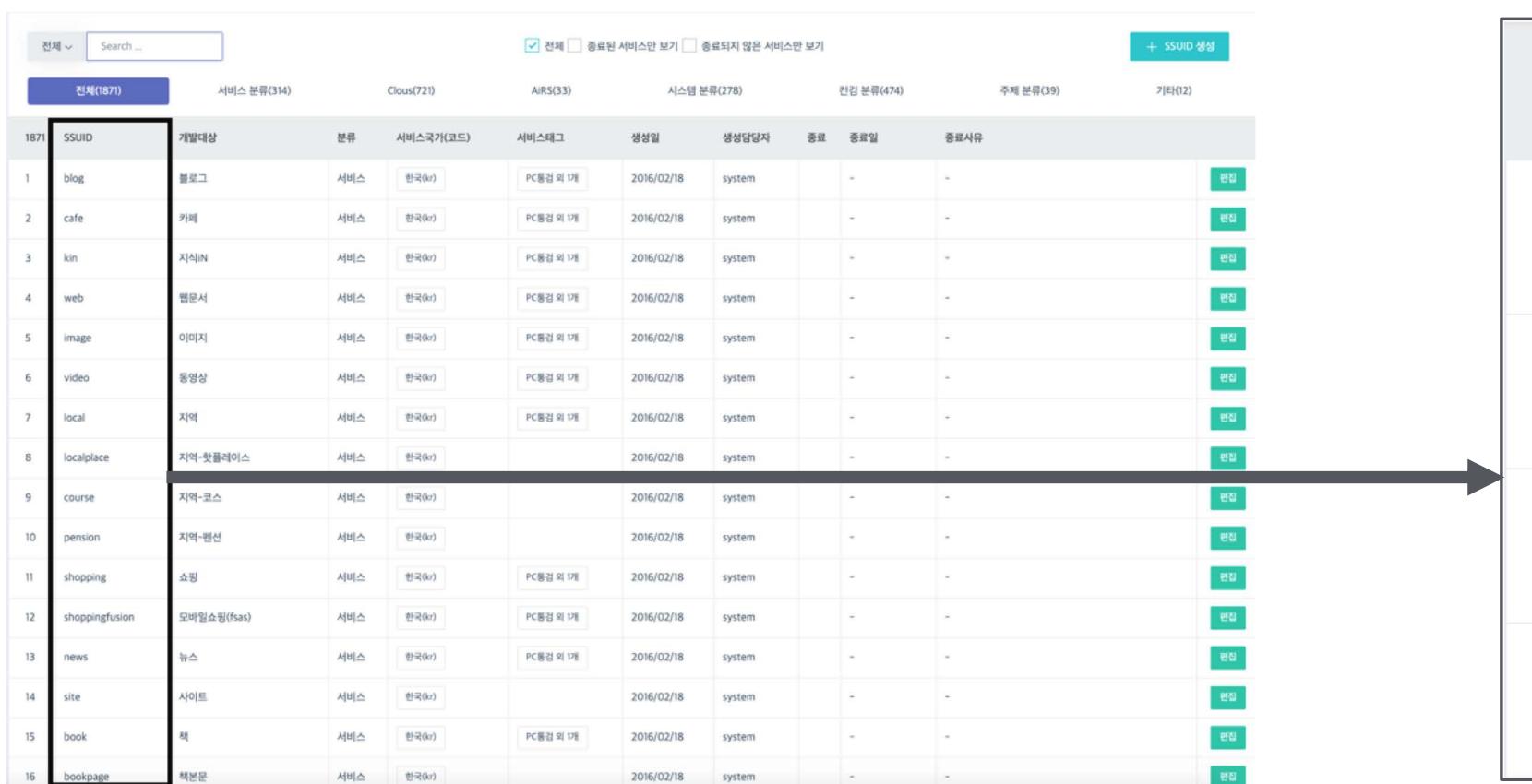
- Etc : local_xxx

[Aggregation]

SSUID: local

(서비스 유니크 아이디)

(1). SSUID (서비스 유니크 아이디) 발급

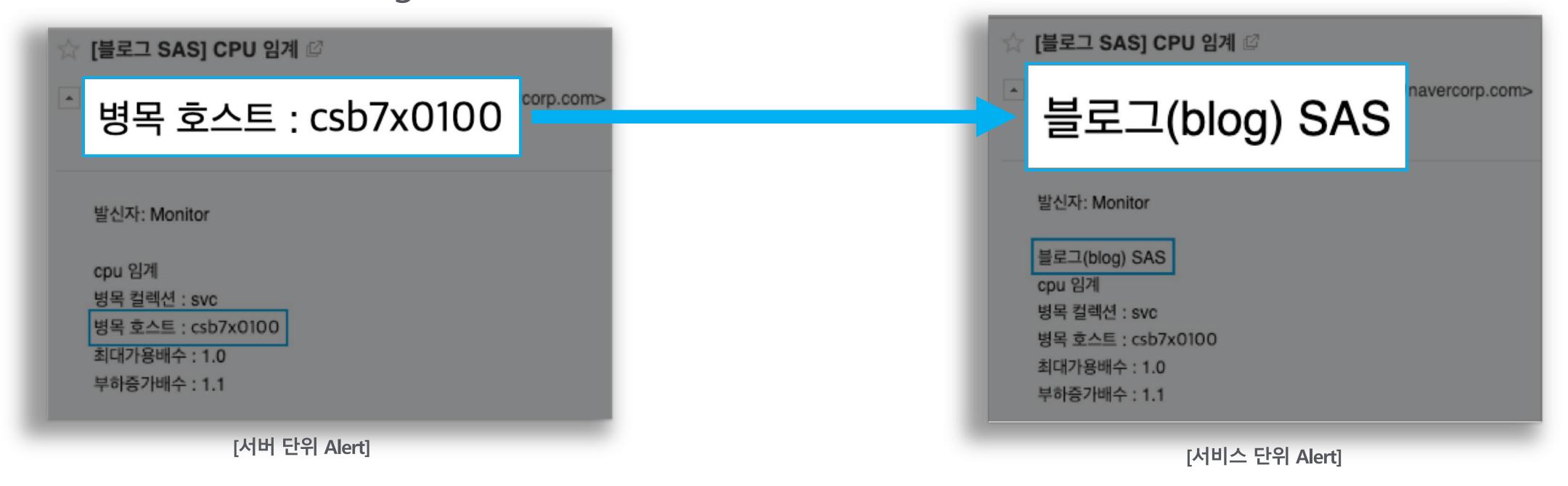


SSUID

1900	SSUID	개발대상
1	blog	블로그
2	cafe	카페
3	kin	지식iN
4	web	웹문서

(1). SSUID (서비스 유니크 아이디) 발급

서비스 유니크 아이디를 정리하니, [서버] 단위의 Alerting이 아니라 [서비스] 단위의 Alerting이 가능!



(2). DRI (서비스 담당자 정보)

Needs: 1. 경보를 서비스 담당자에게 직접 보내야 함

2. 그런데 담당자 정보는 언제나 outdated

Solution: 경보에 특화된 Meta Data API를 만들고, 강제성 부여

System&Solution

(2). DRI (서비스 담당자 정보)

local									
비교에 Auto 표시	시가 있는 DRI는 정비정보에 의해 지	F동 생성된 DN 입니다. 장비정보에 따라 삭제	1될 수 있습니다. DHI가 삭제되길 원하지 않으시다면 Auto 를 한번 눌러주세요						
20	SSUID	개발대상	담당자	조직					
1	local	지역	조지선 ①		20	SSUID	개발대상	담당자	조직
2	local	지역	구동현 ①						
3	local	지역	최지훈 ①	AIRSPACE	1	local	지역	조지선 ①	
4	local	지역	강영길 ①	Local Search					
5	local	지역	정승재 ①	POI Refinemen			-1-1		
6	local	지역	최지훈 ①	AIRSPACE	2	local	지역	구동현 ①	
7	local	지역	이해웅 ①	POI Refinement					
8	local	지역	조명재 ①		3	local	지역	최지훈 ①	AiRSPACE
9	local	지역	이상서 ① 점비담당	System&Solution					
10	local	지역	박연주 ①	·		운영			
11	local	지역	유준현 ①	Local Search		모델링		DRI	
12	local	지역	조한성 ①	System&Solution		PRS			
13	local	지역	이해웅 ① 장비담당	POI Refinement		DCB			

운영

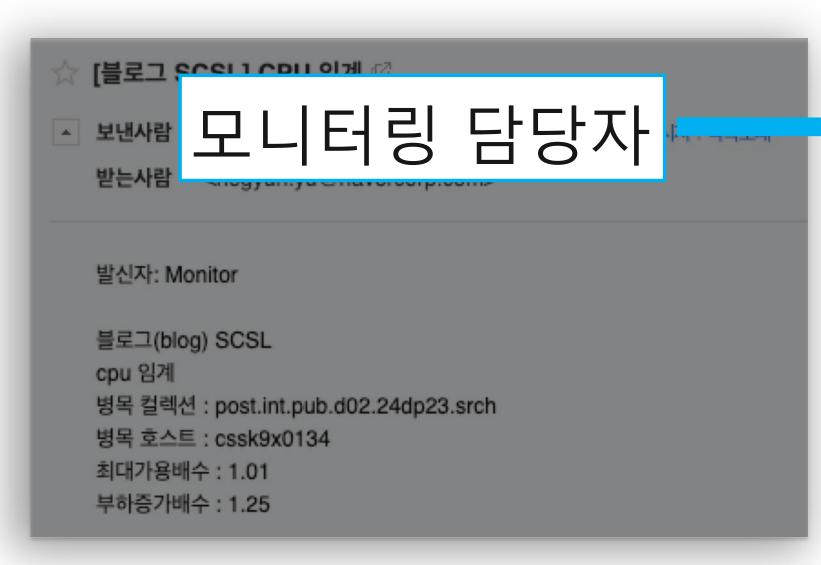
SAS

(2). DRI (서비스 담당자 정보)

경보를 위한 담당자 정보를 정리하니,

[모니터링 담당자] 가 아니라

[서비스 담당자] 가 경보를 수신



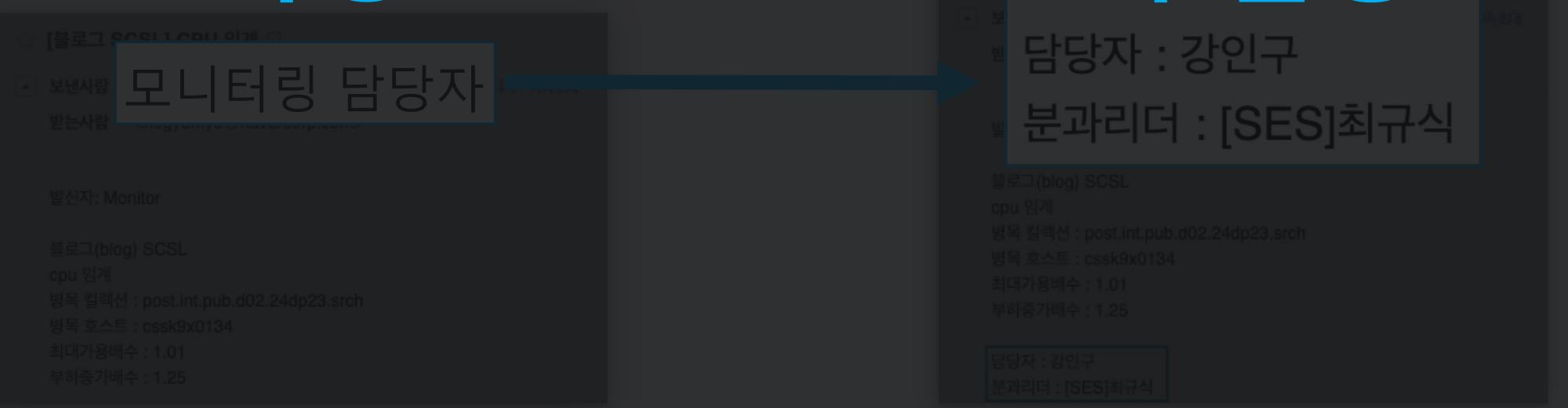
[모니터링 담당자 Alert]

(2). DRI (서비스 담당자 정보)

경보를 위한 담당자 정보를 정리하니,

[모니터링 담당자] 가 아니라

버스타양한 Meta-Data의 활용

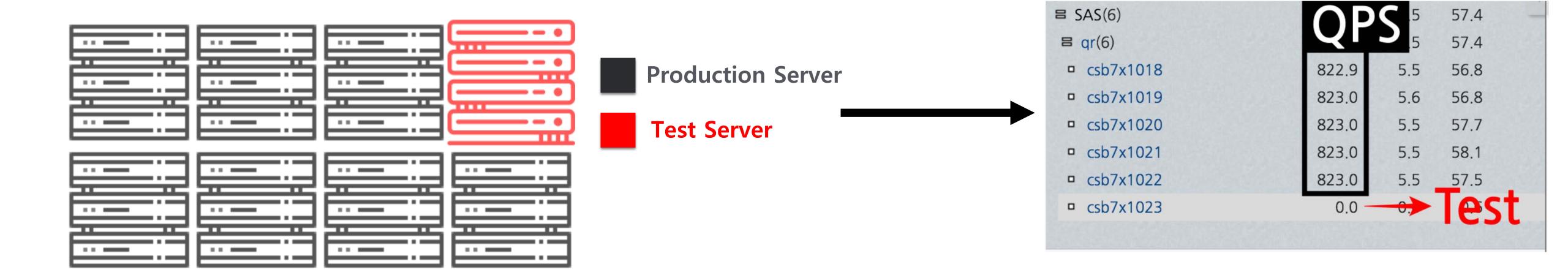


[서비스 담당자 Alert

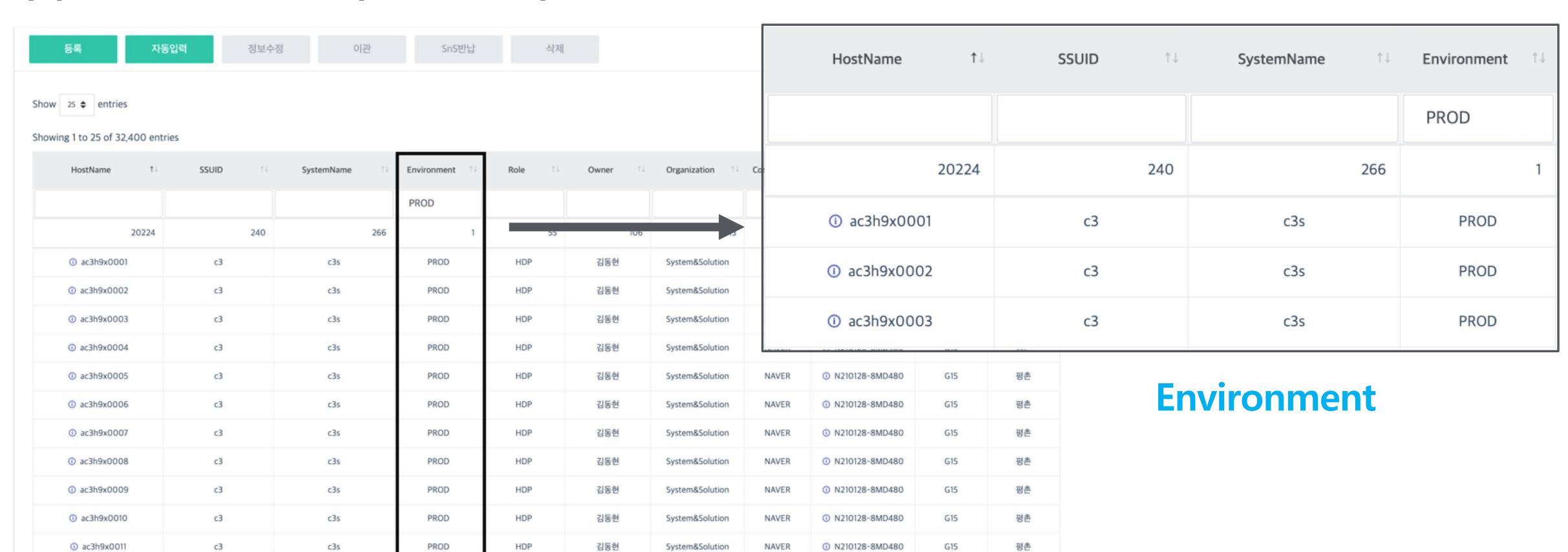
(3). Environment (장비 환경)

Problem: 수만대의 서버 중 테스트 서버가 섞여 있으면 합산 지표에 왜곡이 발생

Solution: dev장비, stage장비, test장비, prod장비 분류



(3). Environment (장비 환경)



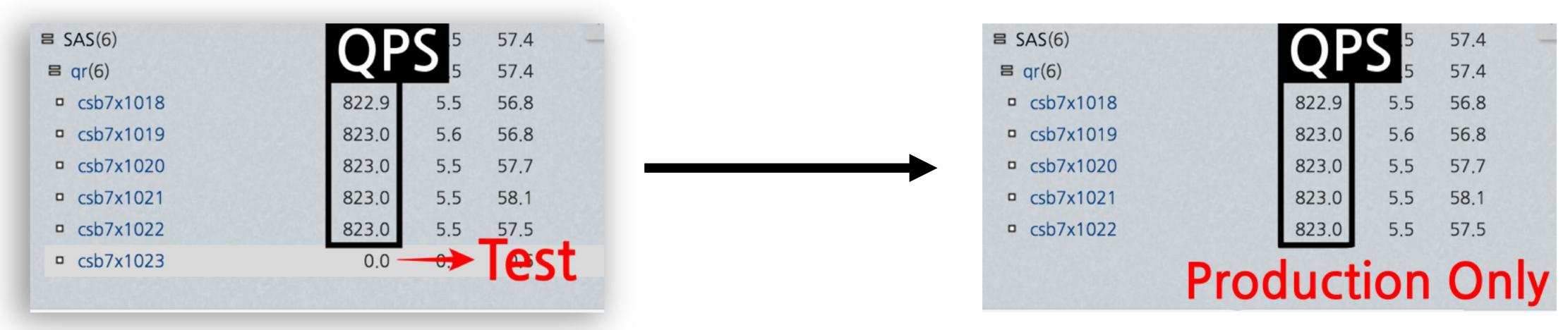
(3). Environment (장비 환경)

장비 Environment를 정리하니

[Test를 포함한 합산 지표]

[서비스 전체 합산 지표] 가 아니라

[Production 서비스 합산 지표] 를 알 수 있게 되었음



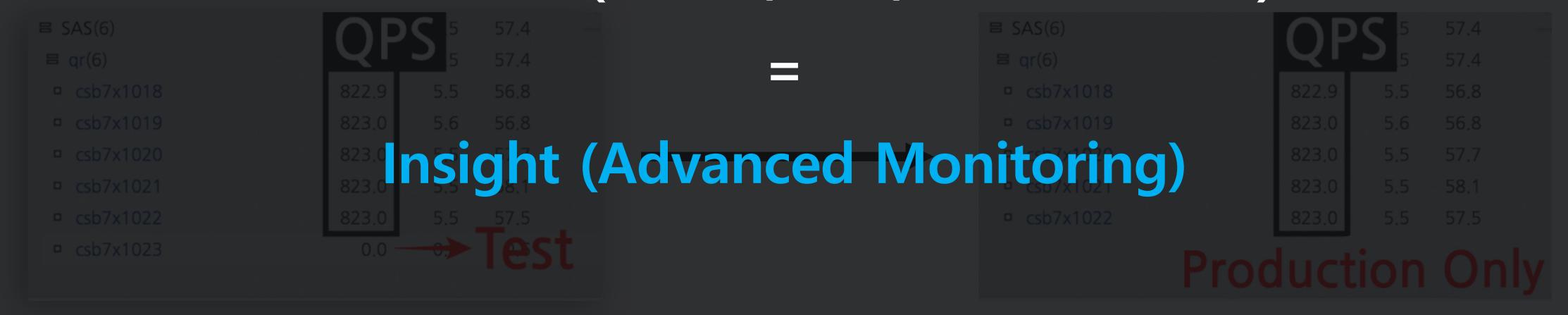
[Production 서비스 합산 지표]

(3). Environment (장비 환경)

장비 Environment를 정리하**Metric** (Raw Data)

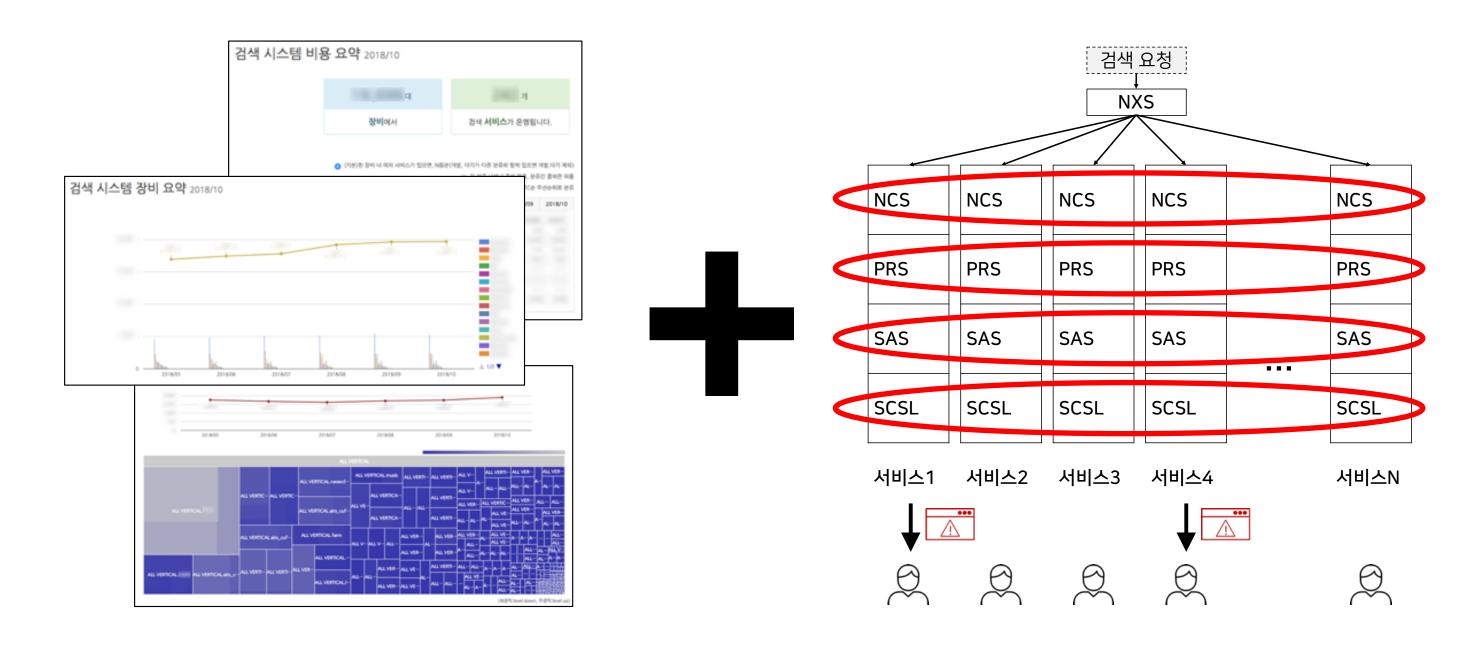
[서비스 전체 합산 지표] 가 아니라

Meta Data (SSUID, DRI, Environment)



ITest를 포함한 합산 지표]

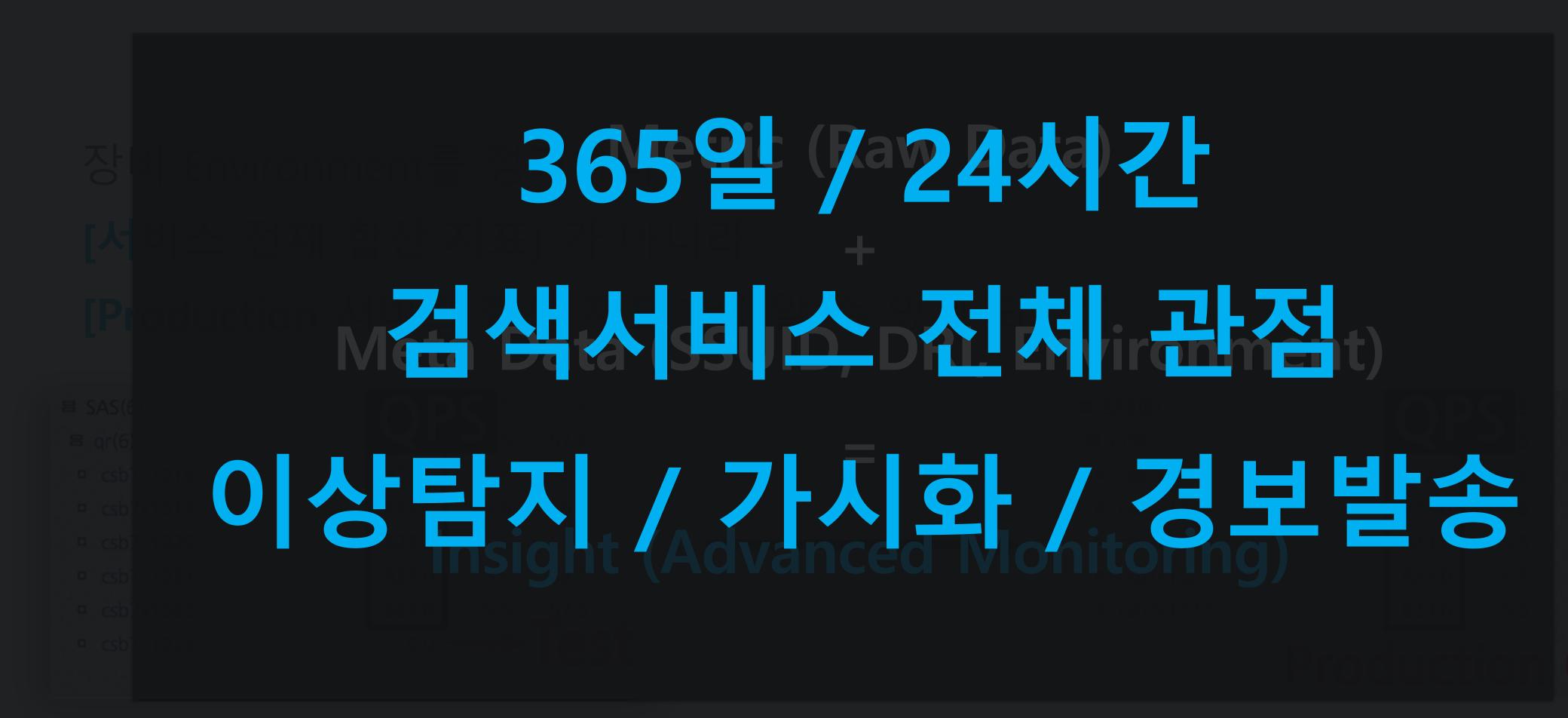
[Production 서비스 합산 지표



Meta Data Tagging, Mapping

서비스/계층 단위 Aggregation

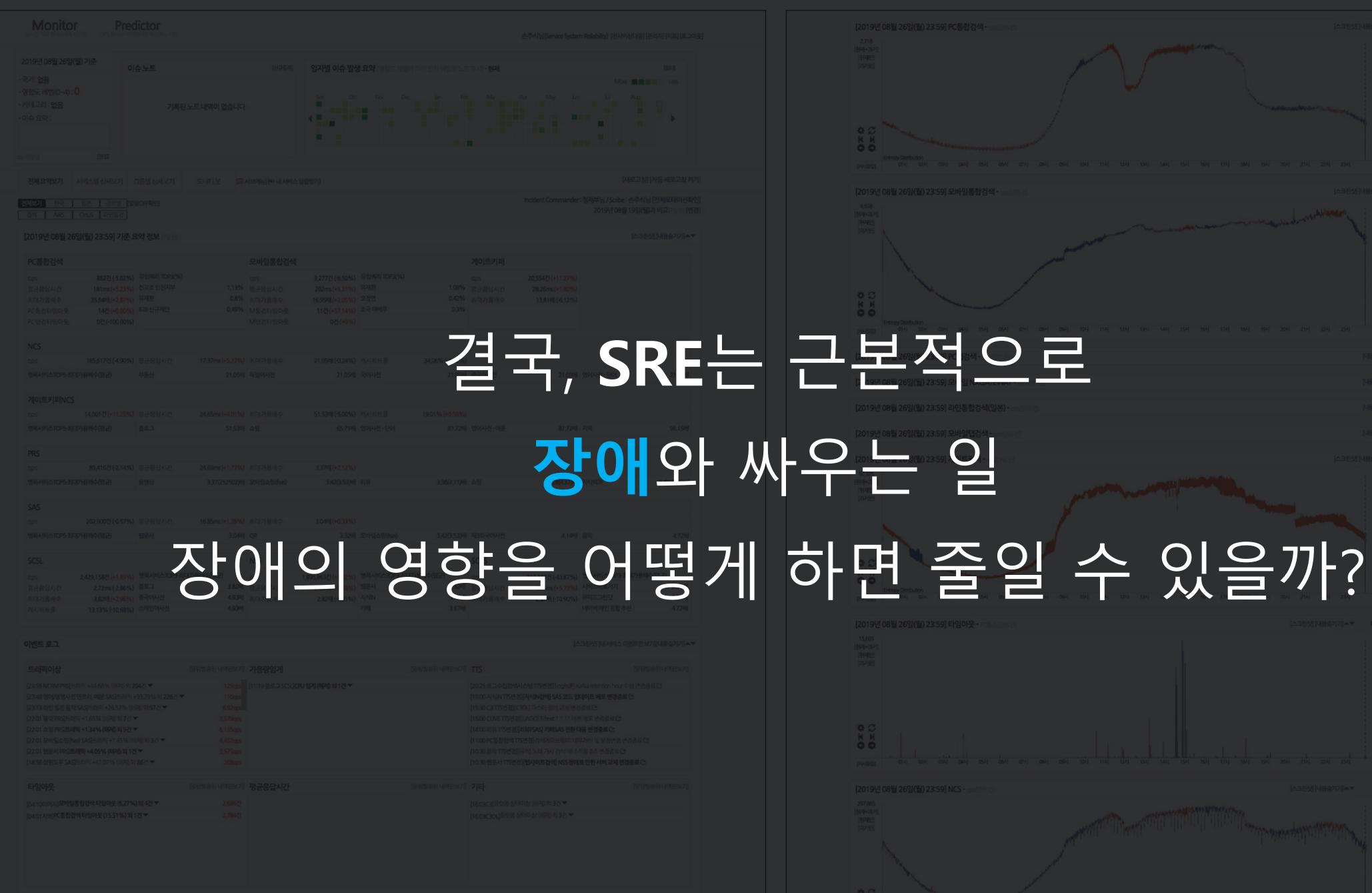


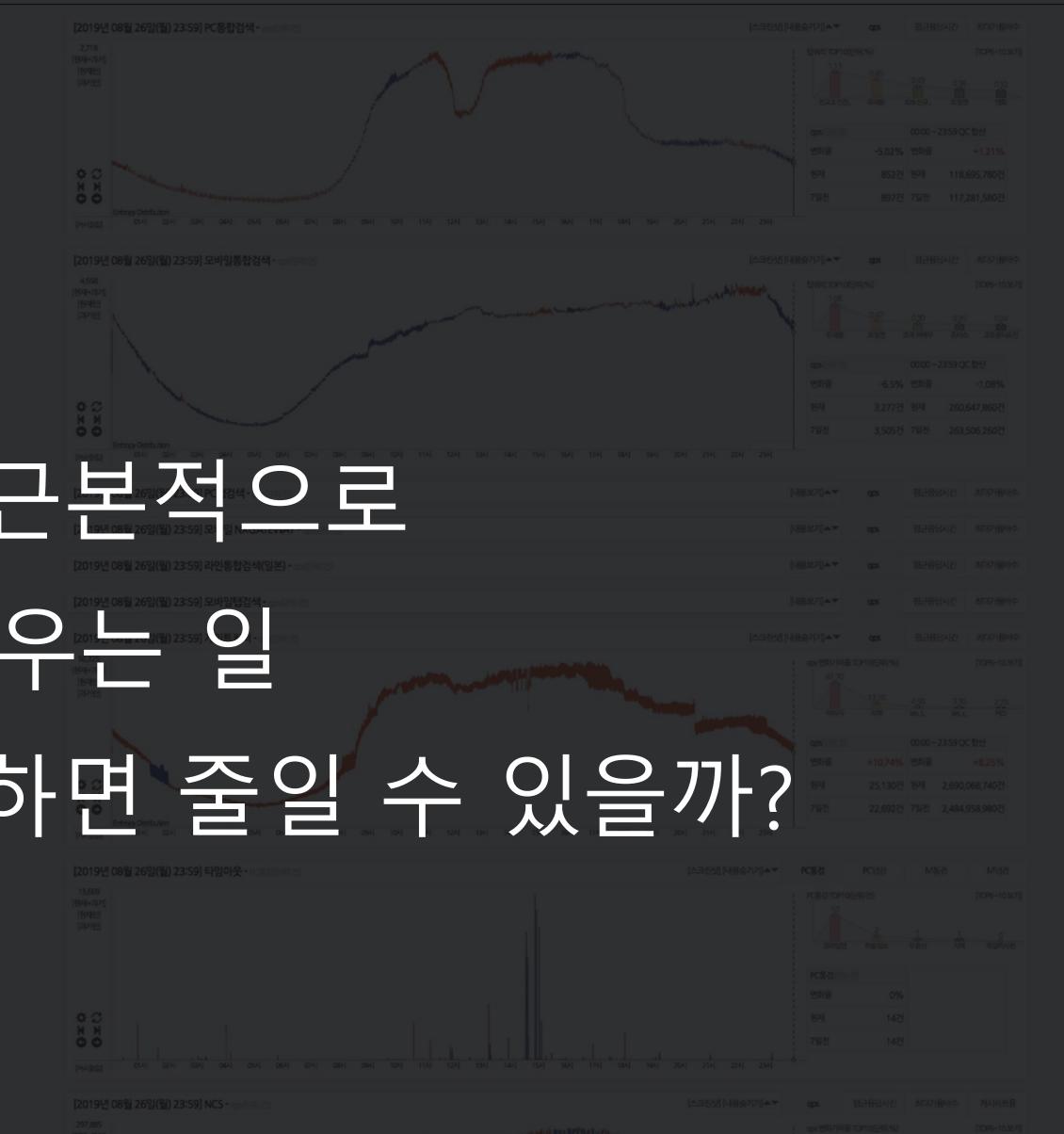


[Production 서비스 합산 지표

사보은 모두 나 Meta Data (SSU관찰I, Environment)

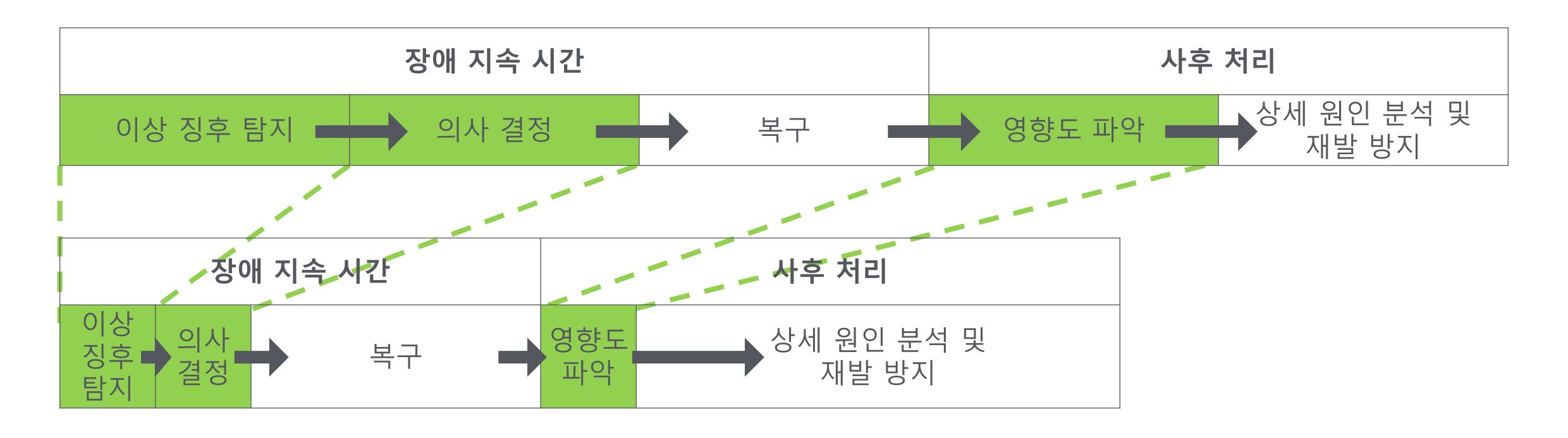
> 실패를통한 개선 Insignt Advanced To The Part of the Part of





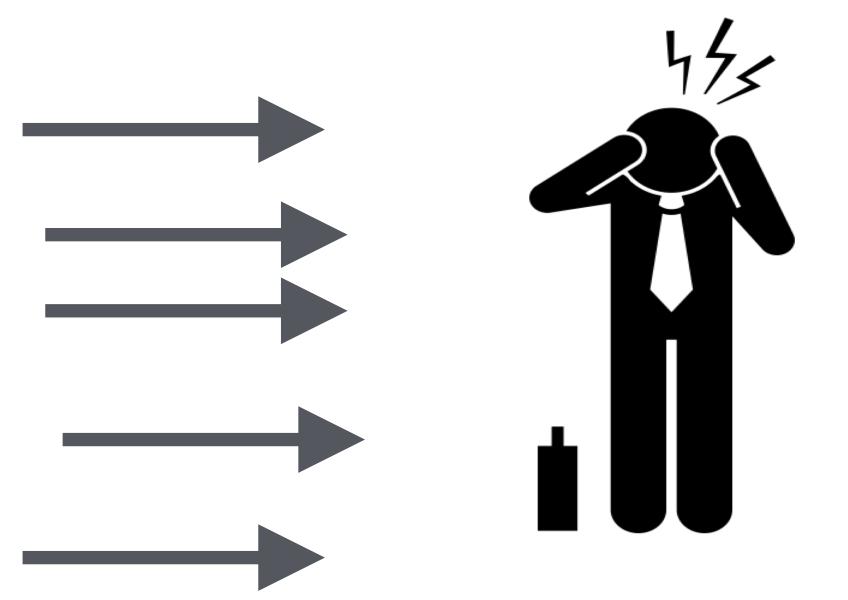


[장애 복구, 원인 분석, 재발 방지 대책]에 소요되는 시간은 줄이기 어려움



그렇지만 [이상 징후 탐지, 의사 결정, 영향도 파악] 시간은 줄일 수 있음

[검색서버팜(태국어) SAS] 가용성 임계상황 해제 [2]	2017-12-06 11:08
[검색서버팜(번체) SAS] 가용성 임계상황 해제 [2]	2017-12-06 11:08
[검색서버팝(영어) SAS] 가용성 임계상황 해제 [2]	2017-12-06 11:08
[검색서버팜(인도네시아어) SAS] 가용성 임계상황 해제 [2]	2017-12-06 11:08
[검색서버팜(간체) SAS] 가용성 임계상황 해제 [2	2017-12-06 11:08
[검색서버팜(일본어) SAS] 가용성 임계상황 해제 ②	2017-12-06 11:08
[검색서버팜(태국어) SAS] 임계상황 1.3<=1.3 farm c3a057 [2]	2017-12-06 11:07
[검색서버팜(한국어) SAS] 임계상황 1.3<=1.3 farm c3a057 [2]	2017-12-06 11:07
[검색서버팜(번체) SAS] 임계상황 1.3<=1.3 farm c3a057 [2]	2017-12-06 11:07
[검색서버팜(간체) SAS] 임계상황 1.3<=1.3 farm c3a057 [2	2017-12-06 11:07
[검색서버팜(영어) SAS] 임계상황 1.3<=1.3 farm c3a057 [2017-12-06 11:07
[검색서버팜(인도네시아어) SAS] 임계상황 1.3<=1.3 farm c3a057 [2017-12-06 11:07
[검색서버팜(일본어) SAS] 임계상황 1.3<=1.3 farm c3a057 [2	2017-12-06 11:07
[검색서버팜(태국어) SAS] 가용성 임계상황 해제 🖸	2017-12-06 11:04
[검색서버팜(일본어) SAS] 가용성 임계상황 해제 🖸	2017-12-06 11:04
[검색서버팜(번체) SAS] 가용성 임계상황 해제 🖸	2017-12-06 11:04
[검색서버팜(영어) SAS] 가용성 임계상황 해제 🖸	2017-12-06 11:04
[검색서버팜(인도네시아어) SAS] 가용성 임계상황 해제 🖸	2017-12-06 11:04
[검색서버팜(한국어) SAS] 가용성 임계상황 해제 [2]	2017-12-06 11:04
[검색서버팜(간체) SAS] 가용성 임계상황 해제 🖸	2017-12-06 11:04
[검색서버팜(인도네시아어) SAS] 임계상황 1.2<=1.3 farm c3a057 [2017-12-06 11:03
[검색서버팜(일본어) SAS] 임계상황 1.2<=1.3 farm c3a057 [2]	2017-12-06 11:03
[검색서버팜(번체) SAS] 임계상황 1.2<=1.3 farm c3a057 [2]	2017-12-06 11:03
[검색서버팜(간체) SAS] 임계상황 1.2<=1.3 farm c3a057 [2	2017-12-06 11:03
[검색서버팜(한국어) SAS] 임계상황 1.2<=1.3 farm c3a057 🖸	2017-12-06 11:03
[검색서버팜(태국어) SAS] 임계상황 1.2<=1.3 farm c3a057 [2]	2017-12-06 11:03
[검색서버팜(영어) SAS] 임계상황 1.2<=1.3 farm c3a057 [2	2017-12-06 11:03
[검색서버팜(영어) SAS] 가용성 임계상황 해제 🖸	2017-12-06 11:00
[검색서버팜(간체) SAS] 가용성 임계상황 해제 🖸	2017-12-06 11:00



[이상 징후 탐지 시간]을 줄이기 위해 [경보 빈도]를 늘리고

하루에 200통 이상 경보를 받기도 하면서

경보 피로가 발생

경보를줄이기 위한 개발이 아니라,

원칙이물요

하루에 200통 이상 경보를 받기도 하면서

원칙 1 [간단함] 오캄의 면도날 (Occam's Razor)

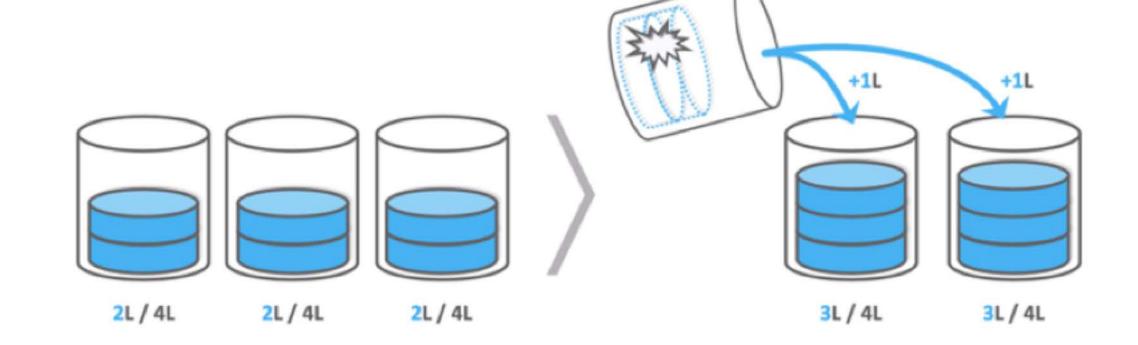
- 가능하면 간단한 추론을 하라.
- 조건이 같다면 간단한 것이 더 강력하다.

원칙 1 [간단함]

가용량 규칙

부하증가배수

한 서버가 죽으면 나머지 서버들은 몇 배를 받나?



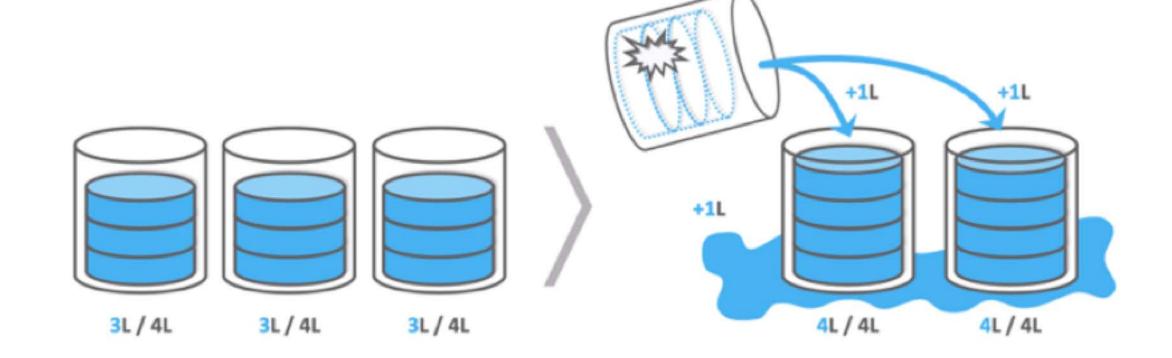
최대가용배수

한 서버가 현재 몇 배까지 받을 수 있나?

So Simple!

"임계 상황 판단 공식"

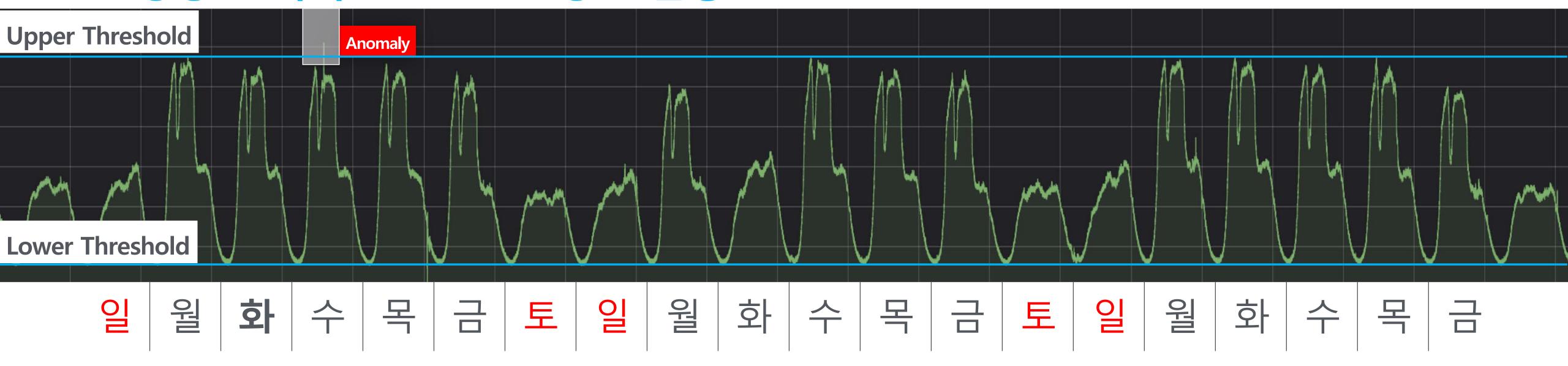
부하증가배수 > 최대가용배수



원칙 1 [간단함] 트래픽 경보

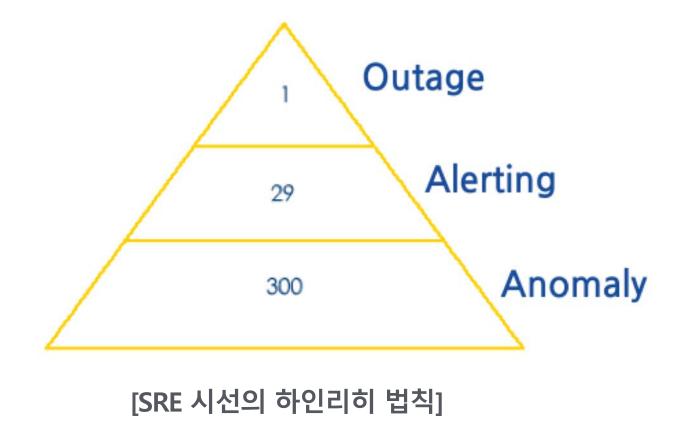
- 사람들의 생체 주기를 따라가는 트래픽 (PC 통합검색)

"정상 트래픽"에서 벗어나면 "경보 발송"



원칙 2 [1:29:300 법칙] 하인리히 법칙 (Heinrich's Law)

- 어떤 **대형 사고**가 일어나기 전에는 수십 차례의 경미한 사고, 수백 번의 징후가 반드시 일어난다.

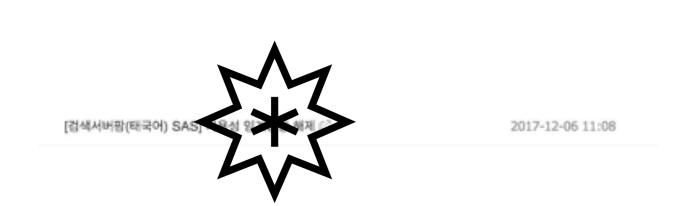


DEVIEW 2019

4. Fail Fast, Learn Faster

원칙 2 [1:29:300 법칙]

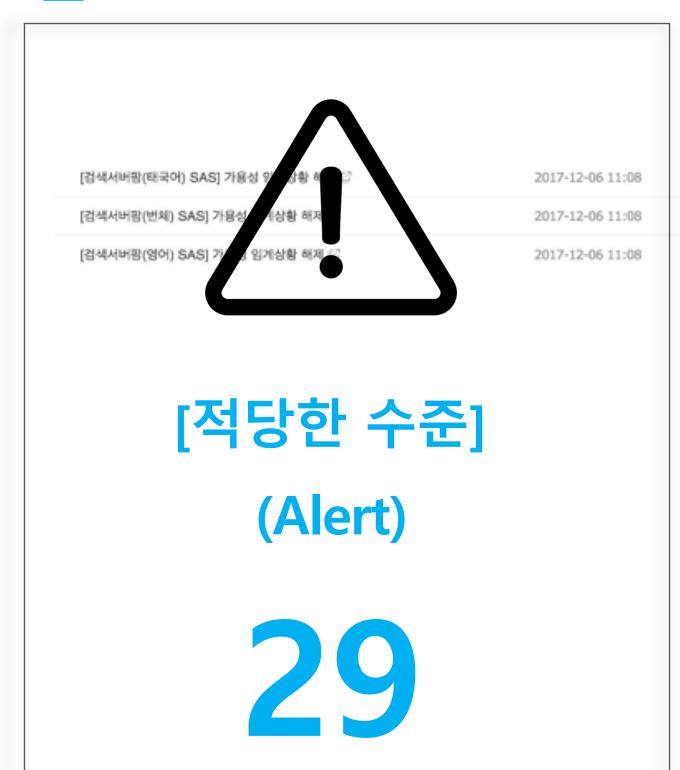
적당한 수준의 경보 빈도



[적음]

(Outage)

1





[많음]

(Anomaly)

300

원칙 2 [1:29:300 법칙]

적당한 수준의 경보 빈도

- Pending:

즉각적으로 대응해야 하는 경우가 아니라면 경보를 잠시 보류

[검색서버팜(태국어) SAS] 가용성 임계상황 웩제 🗗	2017-12-06 11:08	
[검색서버팜(번째) SAS] 가용성 임계상황 해제 ©	2017-12-06 11:08	
[검색서버팜(영어) SAS] 가용성 임계상황 해제 😅	2017-12-06 11:08	
[검색사바람(민도네시아어) SAS] 가용성 임계상황 해제 ②	2017-12-06 11:08	
[검색서버팜(간체) SAS] 가용성 임계상황 해제 🕖	2017-12-06 11:08	
[검색서버팜(일본어) SAS] 가용성 임계상황 백제 ②	2017-12-06 11:08 - [06:21 라인	- [06:21 라인뉴스-카테고리 추천 SAS]CPU 임계 [singleNode] (1분 판
[검색서버랑(태국어) SAS] 임계상황 1.3<=1.3 farm c3a057 🗇	2017-12-06 11:07	
[검색서버림(한국어) SAS] 임개상황 1.3<=1.3 farm c3a057 🗇	2017-12-06 11:07	
[검색서버림(번째) SAS] 임계상황 1.3-c=1.3 farm c3a057 🗇	2017-12-06 11:07	[Pending]
[검색서버림(간체) SAS] 임계상황 1.3-c=1.3 farm c3a057 🗭	2017-12-06 11:07	
[검색서버용(영어) SAS] 임계상황 1.3-c=1.3 farm c3a057 ©	2017-12-06 11:07	
[검색서버림(인도네시아이) SAS] 임기상황 1.3-c=1.3 farm c3a057 🗊	2017-12-06 11:07	
[검색서버림(일본어) SAS] 임계상황 1.3<=1.3 farm c3a057 ©	2017-12-06 11:07	
[검색서버림(태국어) SAS] 가용성 임계상황 레제 🗇	2017-12-06 11:04	
[검색서버림(일본어) SAS] 가용성 임계상황 배제 ②	2017-12-06 11:04	

[Threshold 넘으면 바로 경보]

[검색서버림(번째) SAS] 가용성 임계상황 해제 😅

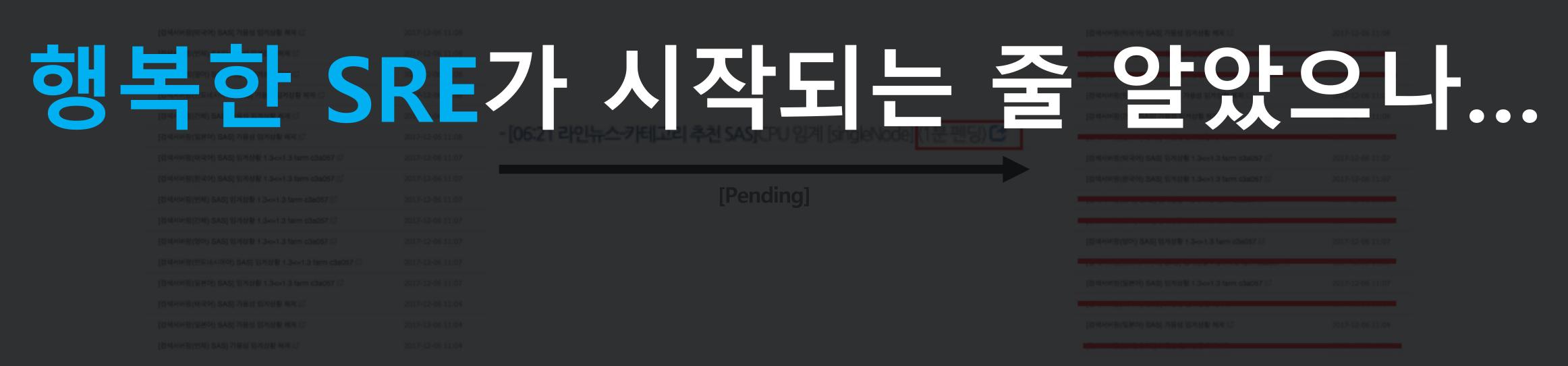


[Pending으로 경보피로 완화]

원칙 2 [1:29:300 법칙]

적당한 수준의 경보 빈도

Pending: 경보자상당수 줄어들고



[Threshold 넘으면 바로 경보]

[Pending으로 경보피로 완화

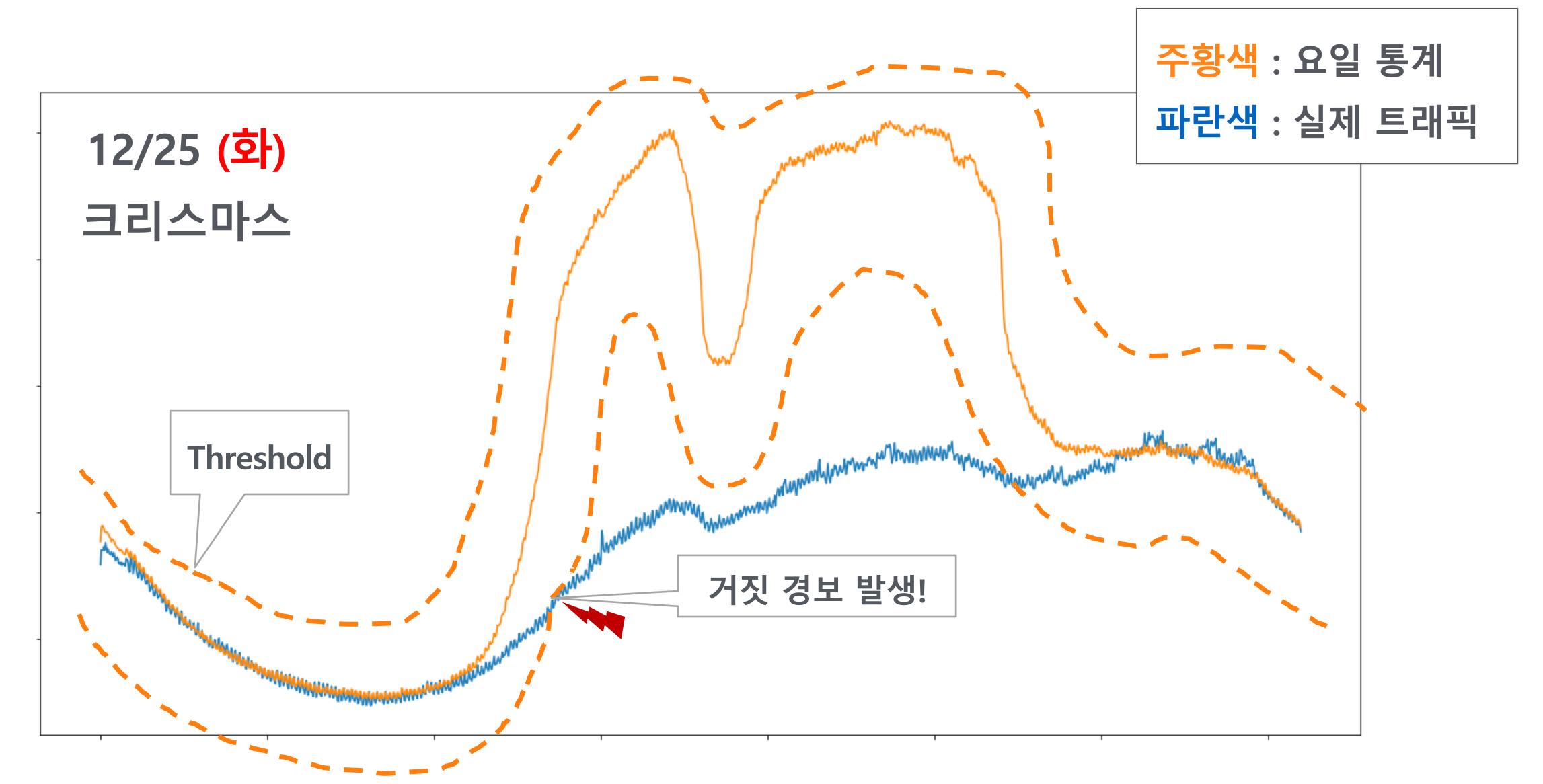
평일 공휴일

그런데... 평화로운 **크리스마스**에 경보가 대량 발생 SRE가 긴급 대응 시작



평일 공휴일 (ex. 크리스마스 화요일) 대량의 거짓 경보 발생

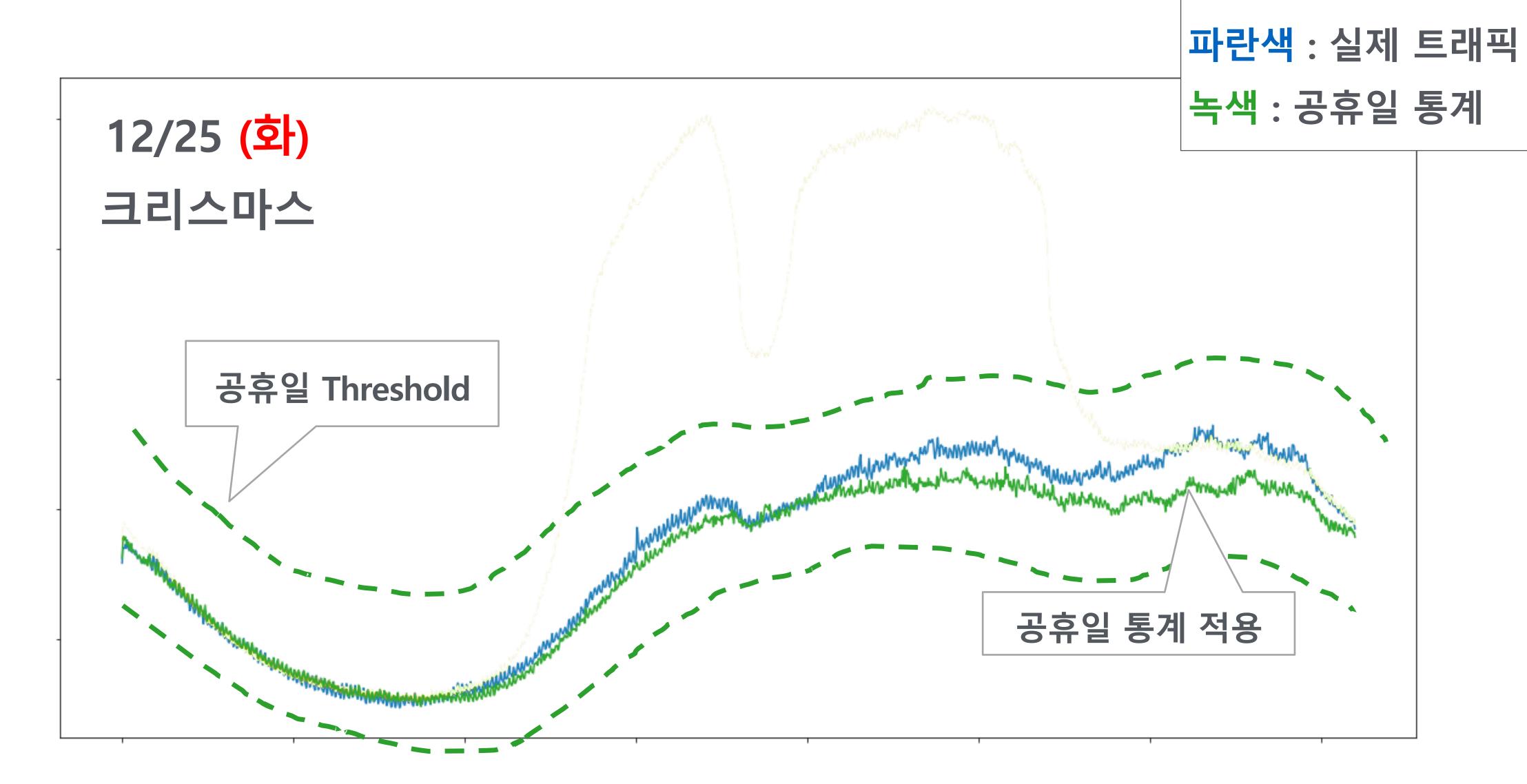




12/25(화) [평일 and 공휴일]은 기존 규칙이 아니라

따로 통계 & 학습을 통해 적용

거짓 경보 발생!



[Pending] + [평일 공휴일 통계, 예측]

(검색서버림(태국어) SAS] 가용성 임계상황 배제 ②	2017-12-06 11:08
[검색서버림(번째) SAS] 가용성 임계상황 해제 ©	2017-12-06 11:08
[검색시버팜(영어) SAS] 가용성 임계상황 해제 😅	2017-12-06 11:08
[검색서버팜(인도네시아어) SAS] 가용성 임계상황 해제 😅	2017-12-06 11:08
[검색시버팜(간체) SAS] 가용성 임계상황 하제 ②	2017-12-06 11:08
[검색서버팜(일본어) SAS] 가용성 임계상황 레제 ②	2017-12-06 11:08
[검색시버팜(태국어) SAS] 임계상황 1.3<=1.3 farm c3a057 🖸	2017-12-06 11:07
[검색서버팜(한국어) SAS] 임개상황 1.3-c=1.3 farm c3a057 🗭	2017-12-06 11:07
[검색서버팜(번째) SAS] 임계상황 1.3-c=1.3 farm c3a057 🕖	2017-12-06 11:07
[검색서버팜(간체) SAS] 임계상황 1.3-c=1.3 farm c3a057 🕖	2017-12-06 11:07
[검색서버팜(영어) SAS] 임계상황 1.3-c=1.3 farm c3a057 🕖	2017-12-06 11:07
[검색서버림(인도네시아어) SAS] 임계상황 1.3<=1.3 farm c3a067 🖸	2017-12-06 11:07
[검색서버팜(일본어) SAS] 임계상황 1.3<=1.3 farm c3a057 ©	2017-12-06 11:07
[검색서버팜(태국어) SAS] 가용성 임계상황 해제 🗗	2017-12-06 11:04
[검색서버림(일본어) SAS] 가용성 임계상황 해제 🗗	2017-12-06 11:04
[검색서버핌(번째) SAS] 가용성 임계상황 해제 😅	2017-12-06 11:04

[Pendng] + [평일 공휴일 통계, 학습 값] 적용 후

(거짓 경보 90% 감소)



[Pending, 평일 공휴일 통계 적용]

[기존 경보]

Pending + 평일 공휴일 통계, 예측

- 1. 트래픽을 자세히 관찰하고
- 2. 다양한 경보들을 수신하고 통계에 작가적용 후
- 3. 고통스러워질 때쯤 개선해보니 감소)

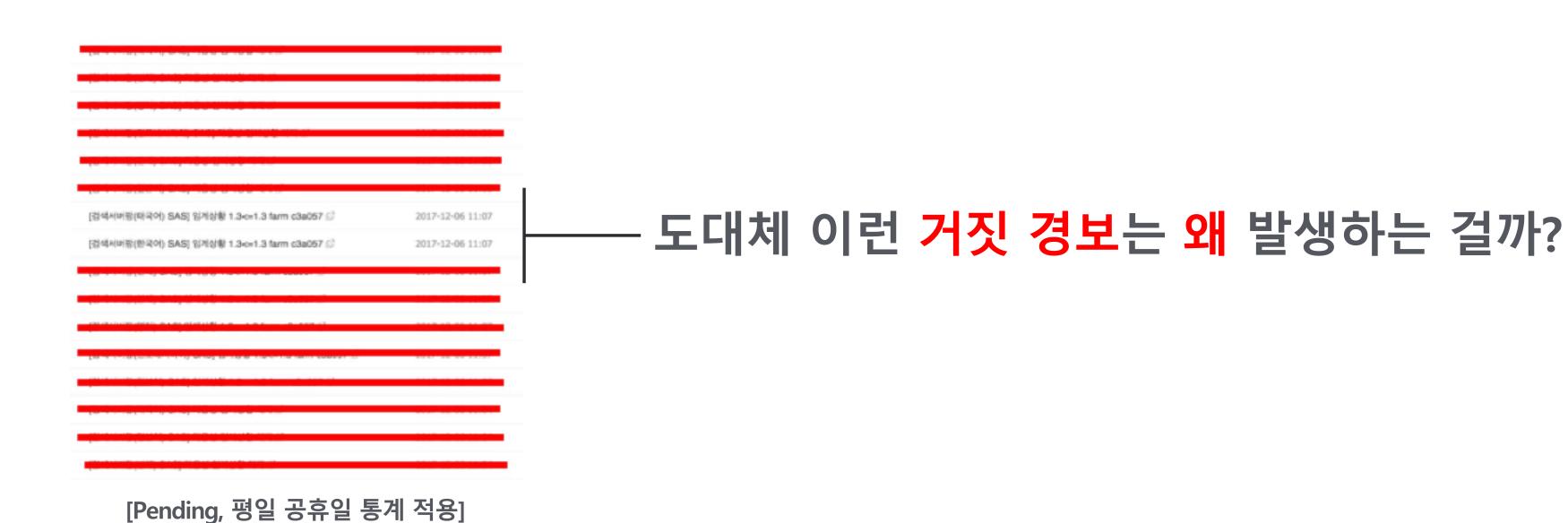
자신감이 생겨서 조금 더 Advanced한 SRE에 도전!

[기존 알람]

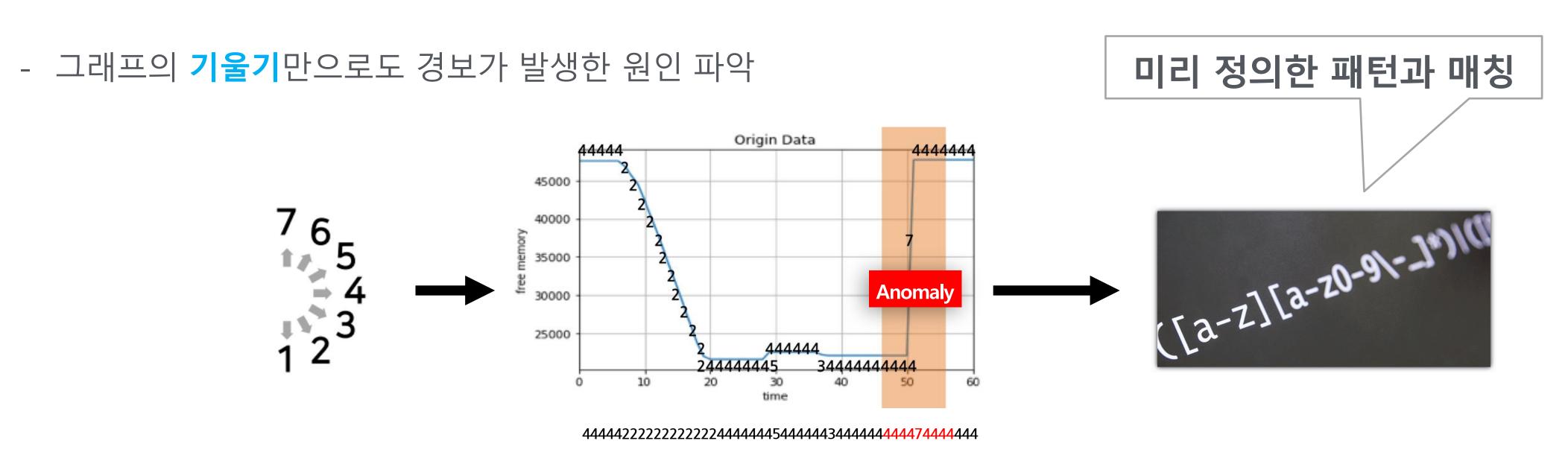
[Pending, 평일 공휴일 통계 적용]

Advanced 1. 트래픽 패턴 분석

- 현재도 발생하는 거짓 경보는 어떤 이유로 발생하는 것일까?



Advanced 1. 트래픽 패턴 분석



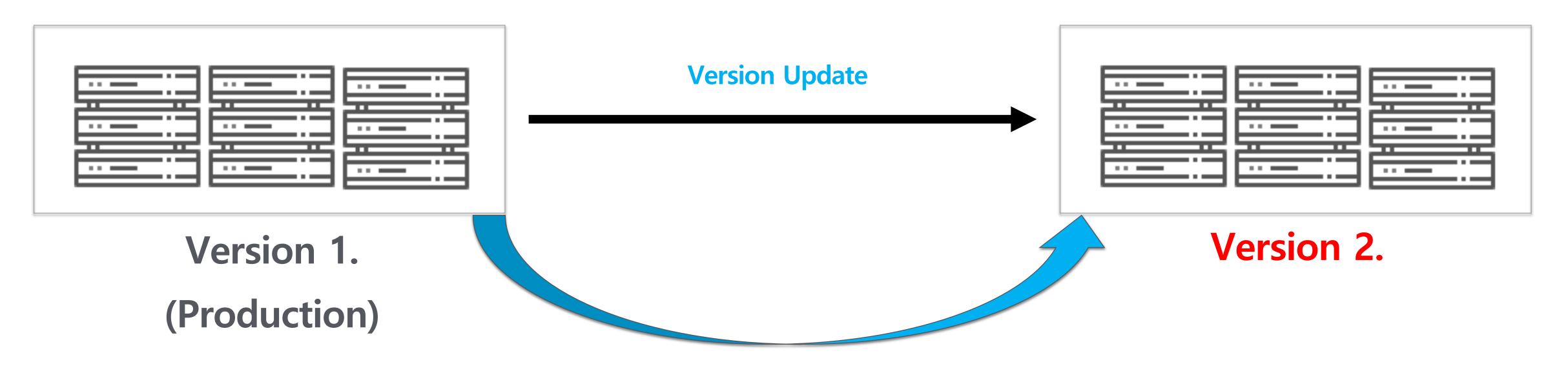
[기울기]

[그래프 인코딩]

[정규 표현식]

Advanced 1. 트래픽 패턴 분석

- 패턴 예제: 트래픽(Traffic)이 다른 곳으로 이전(Migration) 되었는지 판단

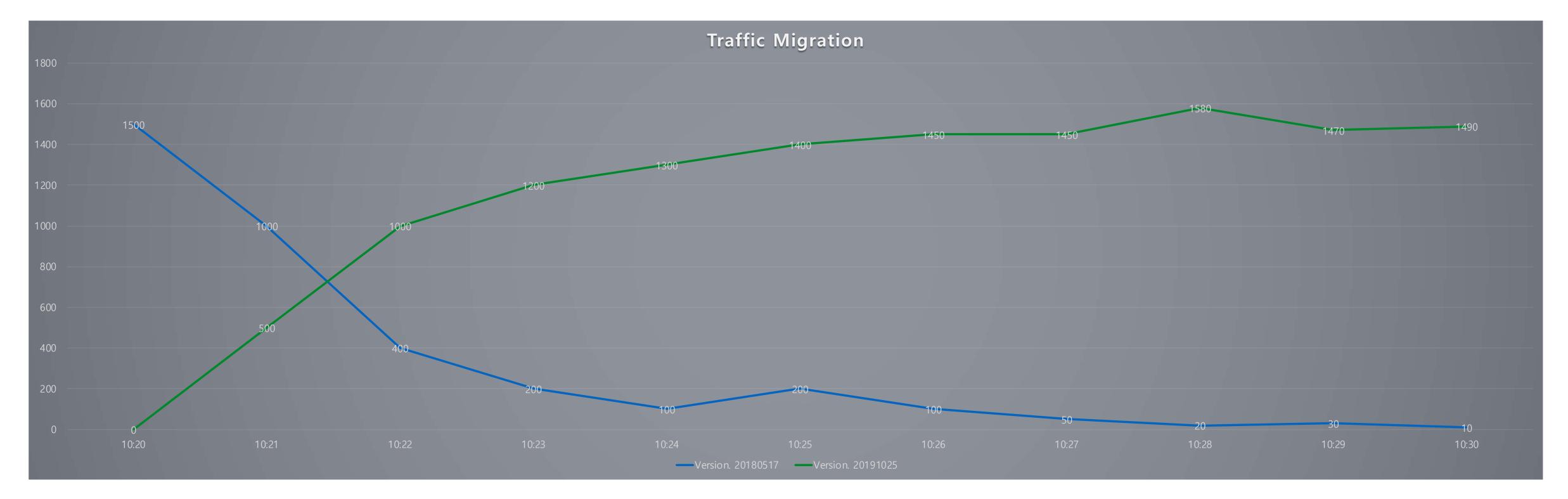


Traffic이 점진적으로 Migration

Advanced 1. 트래픽 패턴 분석

- Traffic Migration: 트래픽(Traffic)이 다른 곳으로 이전(Migration) 되었는지 판단



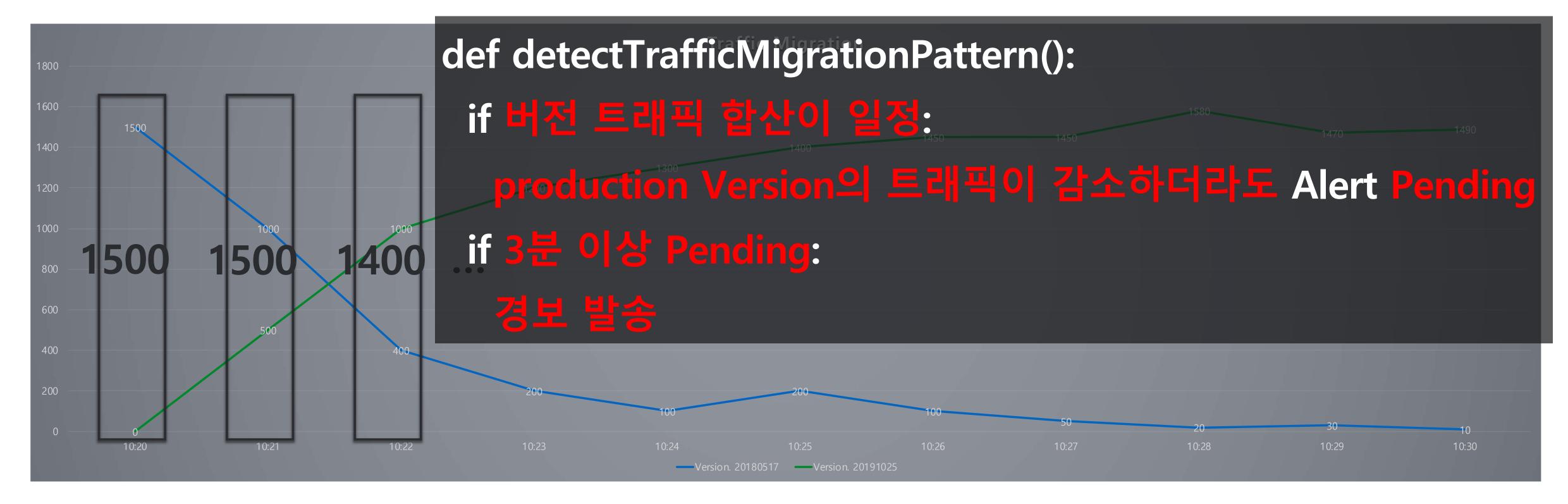


Advanced 1. 트래픽 패턴 분석

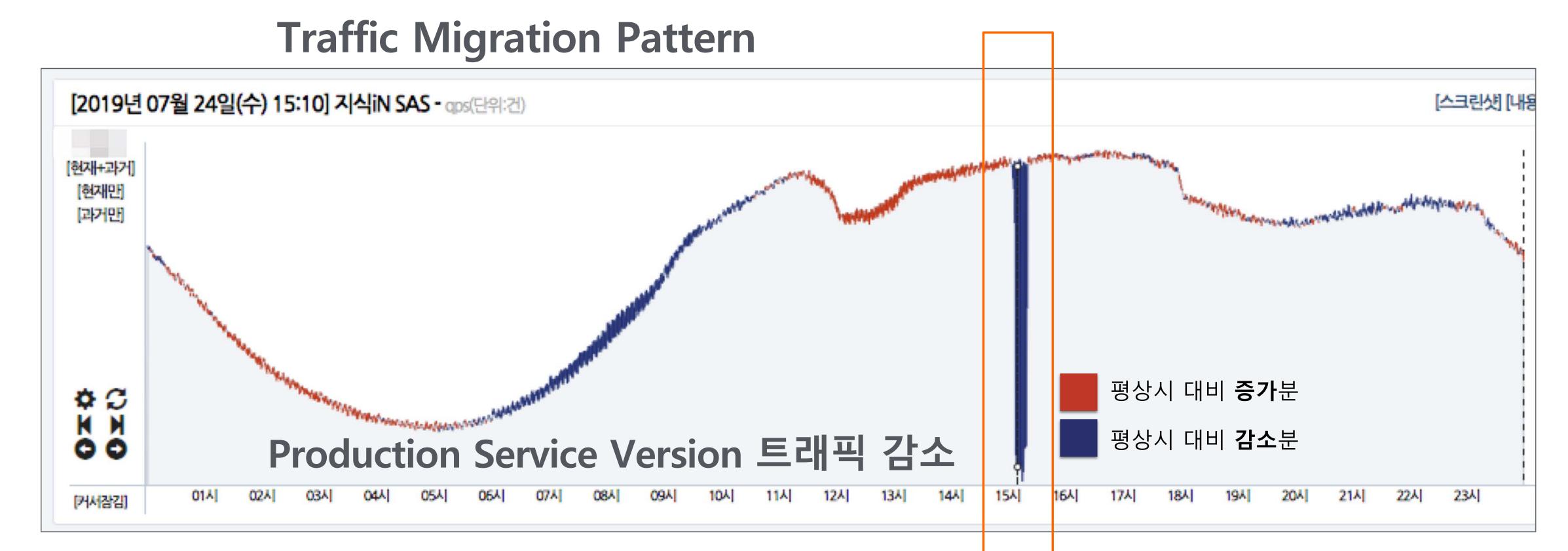
- Traffic Migration: 트래픽(Traffic)이 다른 곳으로 이전(Migration) 되었는지 판단

Old Service Version
(Production Service Version)

New Service Version



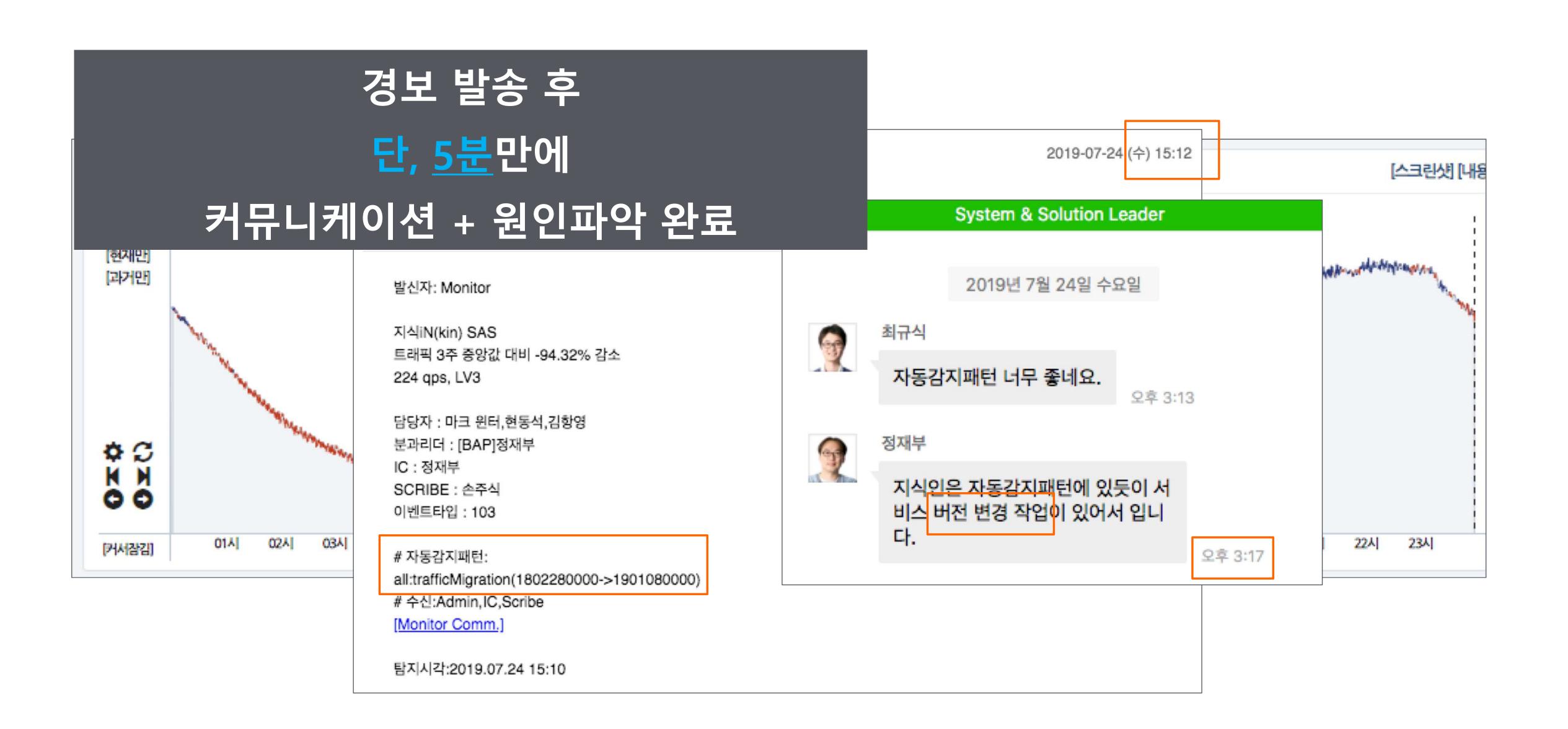
Advanced 1. 트래픽 패턴 분석 (사례)

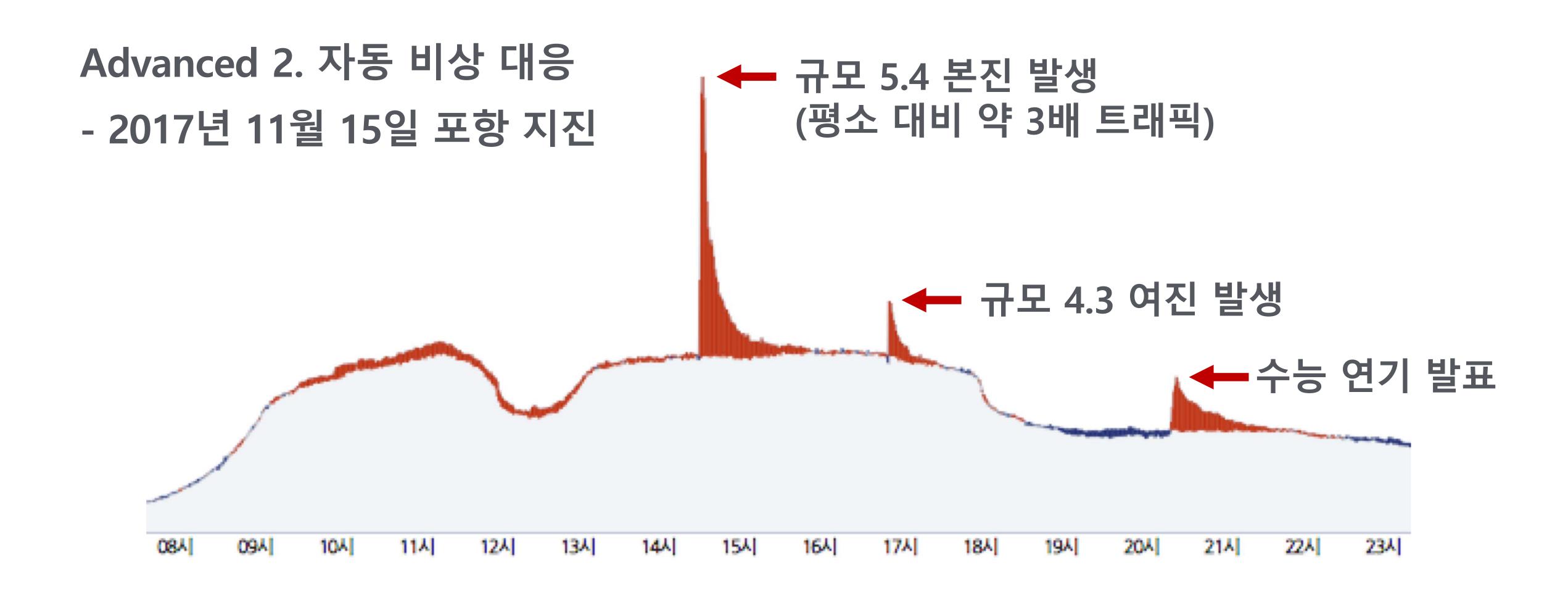


Advanced 1. 트래픽 패턴 분석 (사례)

Traffic Migration Pattern



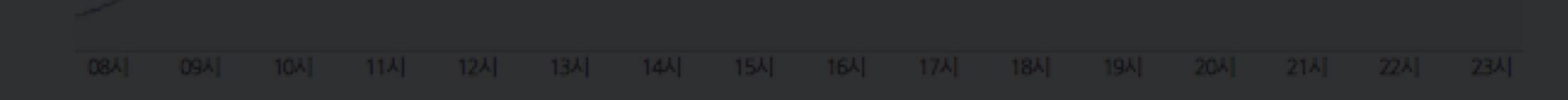




Advanced 2. 자동 비상 대응 - 2017년 11월 15일 포항 지진

→ 규모 5.4 본진 발생 (평소 대비 약 3배 트래픽)

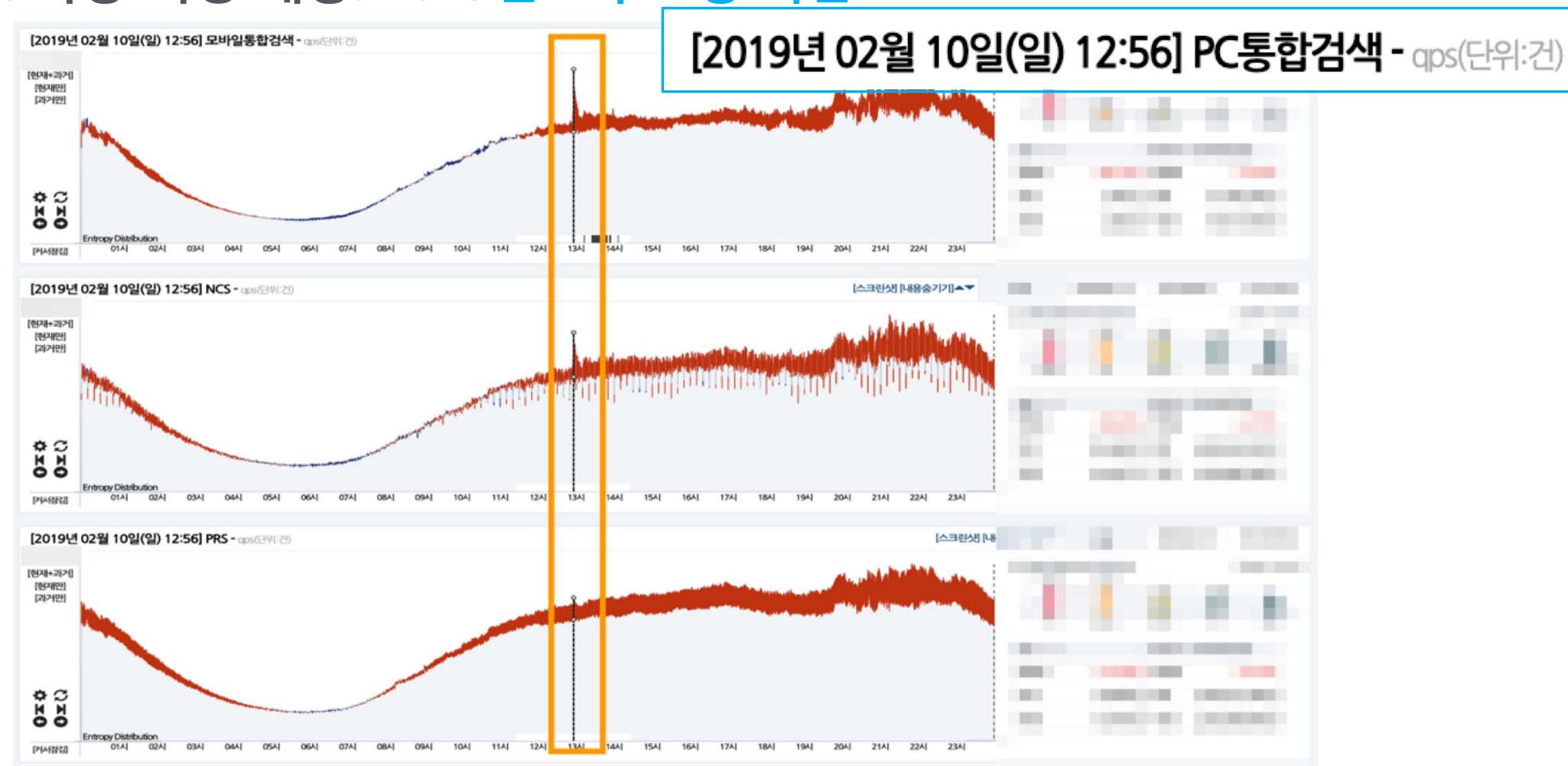
2017년 말 자동 비상 대응 시스템^{4,3} 여진 발생 개발 시작



Advanced 2. 자동 비상 대응 시스템 개발 완료 (2018년)



Advanced 2. 자동 비상 대응: 2019년 2차 포항 지진



Advanced 2. 자동 비상 대응: 2차 포항 지진

트래픽 급증하여 유입되기 전 자동 방어 장치 가동 때의 공적목적



Advanced Result. 빨라진 영향도 분석

장애 Detection

증상 완화

원인 파악

장애 영향도 파악

장애 복구 완료

포스트모텀 작성

[SRE Before]

2분

20분

40분

1시간

수시간

48시간

Mornior Predictor

Mornior Predi

전체 현황, 영향도 분석



[SRE After]

오캄의 면도날

하인리히 법칙

1:29:300

거짓 경보

평일 공휴일

트래픽 경보

최대가용배수

부하증가배수

가용량

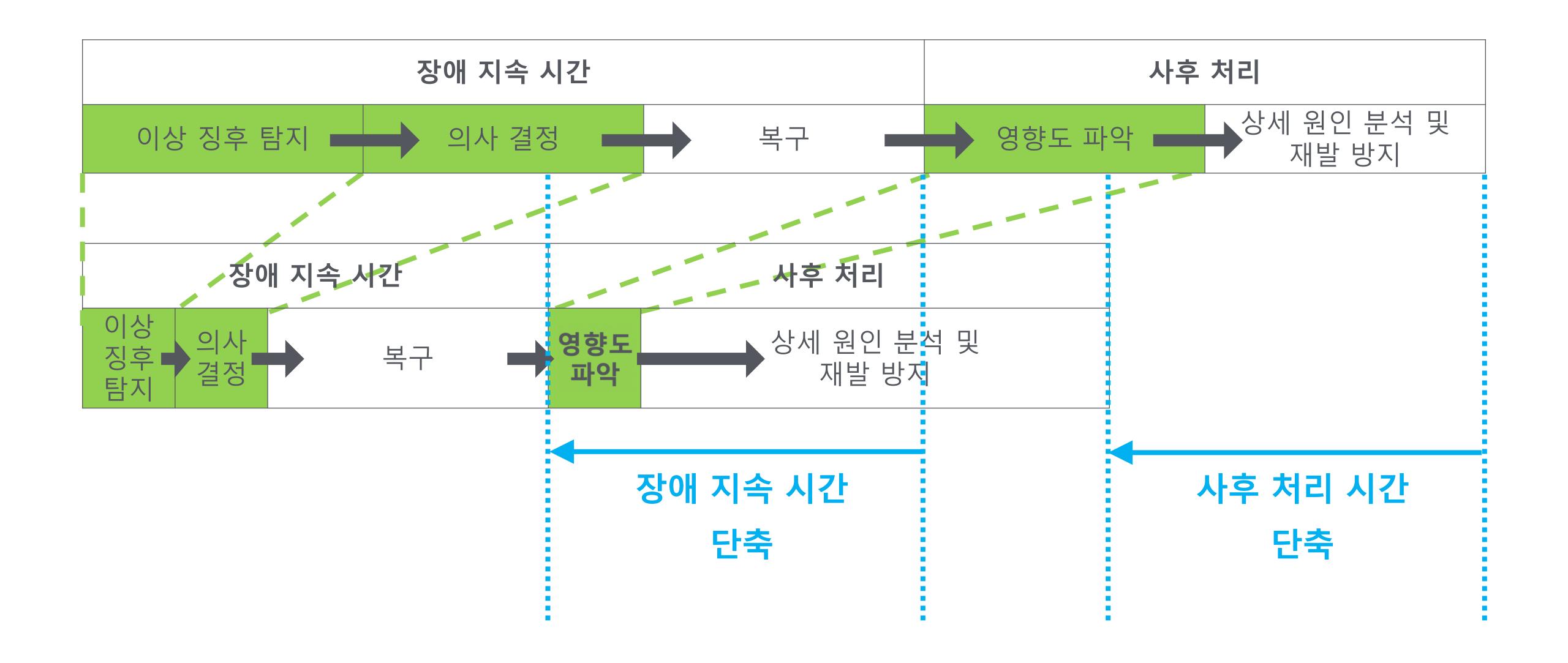
임계상황

Summary

트래픽 분석

영향도 분석

비상 대응 시스템



단순하고, 효과적인 방법으로

[가용량, Treshold기반의 경보]

단순하고, 효과적인 방법으로 빠르게 실패하고

[많은 거짓 경보, 평일 공휴일 에러, 지진]

단순하고, 효과적인 방법으로 빠르게 실패하고, 빠르게 배워서

[Pending + 평일 공휴일 통계 + 트래픽 패턴 분석 + 자동 비상 대응]

단순하고, 효과적인 방법으로

빠르게 실패하고, 빠르게 배워서 [개선]한다.

[업데이트 대비 장애 비율 감소]

많아진 변경 작업에도, 줄어든 장애 (19.09.03 기준)

	2017년	2018년	2019년
업데이트 (건)	1311	2110	1566
장애 (건)	57	49	18
변경 대비 장애 비율 (%)	4.35	2.32	1.15
증감 YoY (%)		-46.67	-50.43

(1) SRE = 재난엔지니어링

본질적인 SRE에 집중하다 보면, 결국 재난 엔지니어링으로 귀결된다.

[이상탐지(anomaly detection)의 자동화 + 경보압축(alert compression)] 추구

- = 시스템 문제점 해결 + 안정적인 시스템 구축
- = 재난 엔지니어링

(2) SRE를 하기 위해서 필요한 Attitude

- 이상적이고 멋들어진 도구가 아닌 실질적인 문제 해결에 주력하기 (Realist)
- 무식한 방법이더라도 힘겹더라도 포기하지 않고 끊임없이 개선하기 (Grit)
- 본인이 속한 도메인에 적절한 방법론을 적용하기 (Domain Specific)

(2) SRE 이런 Attitude를 장착하고

- 이상적인과 들어진것인것이 있는 질문에 함께 유력한지 Real 우리와 재미있는 일을 함께할
- 무식한방법이더라도함겹더라도포기하지않고끊임없이개선하기 (Gree 당신을 찾습니다.
- 본인이속한도메인에적절한방법론을적용하기 (Domain Spe(채용, 엄격, 근엄, 진지)

(2) SRE를 하기 위해서 필요한 Attitude

이렇게만 하면 아무도 연락 안 할 테니,

- 이상적이고 멋들어진yoohoogun114으로 [카톡] 주세요! 에 주력하기 (Realist) 그린팩토리에서 커피 사드립니다
- 무식한 방법이더라도[채용공과링크 Click] 입고 끊임없이 개선하기

(Greet)

- 본인이 속한 **도메인**에 적절한 방 (Domain Specific)

:smile:



Q & A

Thank You