

Google SRE分享： 17LIVE的SRE太極拳

林毅民 (Sammy Lin),
17LIVE Engineering Director



何宗憲 (Shawn Ho),
GCP AppMod Specialist



```
apiVersion: v1
kind: Way-to-Succeed
metadata:
  name: SRE
data:
  culture: "Software oriented operations"
  monitoring:
    slis: |
      metrics=availability,latency,throughput
    slos: |
      service1=99.9%,95%in 28 days
  automation:
    - incident response
    - change management
    - monitoring and alerting
    - capacity planning
```



只重其意，不重其招
忘掉所有招式，你就練成了太極拳

Topics

- Google 的SRE 心法 (14分鐘)
- 面向Google SRE的服務設計(10分鐘)
- 17LIVE 的SRE實踐 (25分鐘)
- Take Away

Reliability Principles

SRE
的心法

1. Reliability is defined by the user
2. Sufficient reliability

3. Redundancy
4. Horizontal scalability

5. Overload tolerance

6. Rollback capability

7. Traffic spike prevention

17LIVE
SRE的實踐

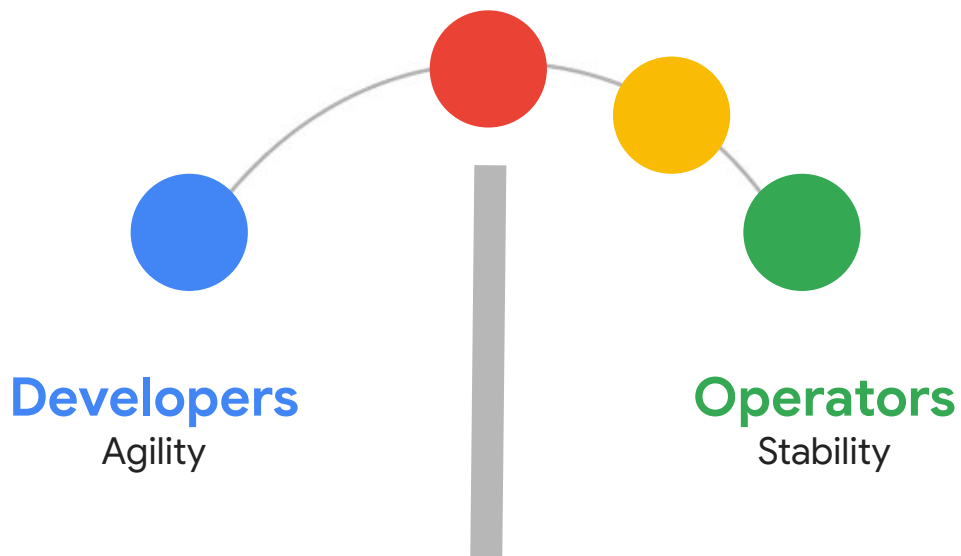
8. Failure recovery testing
9. Failure detection
10. Incremental change
11. Coordinated emergency response
12. Observability
13. Emergency response documentation and automation
14. Capacity Management
15. Toil Reduction



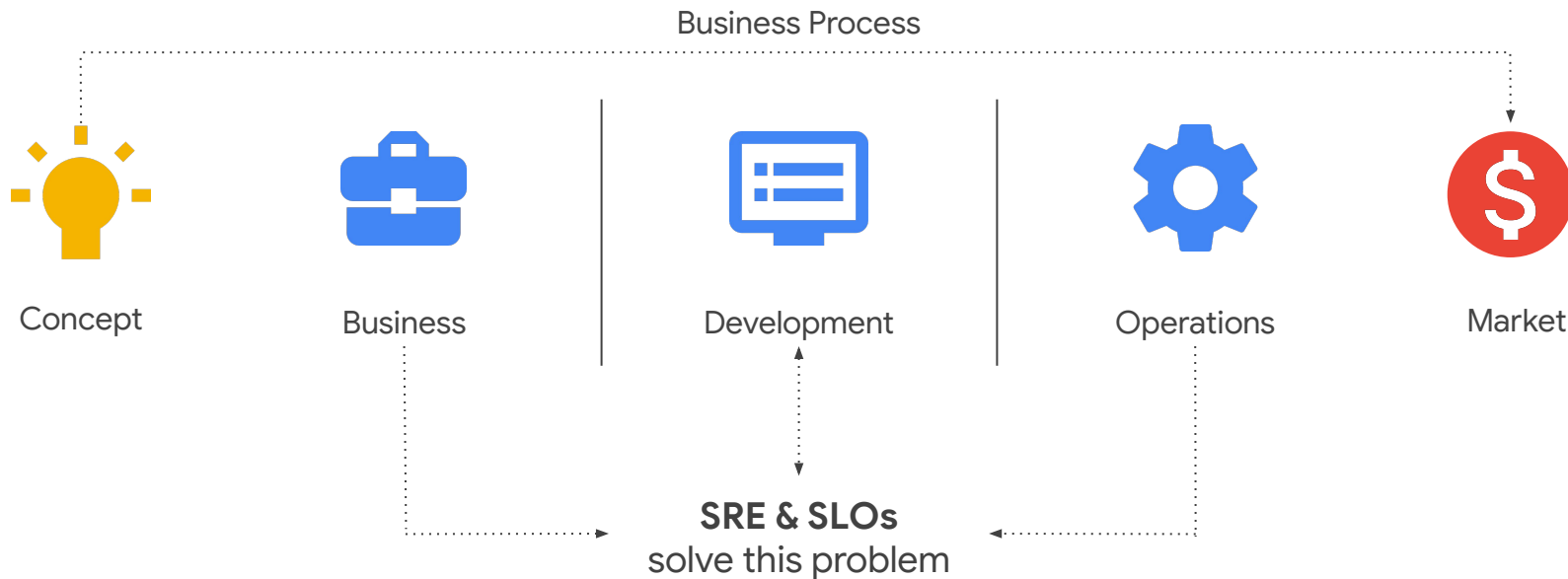
Google Cloud

面向SWE的
服務設計建
議

Incentives between groups aren't aligned



Reduce Product Lifecycle Friction



Google 的 SRE 心法



Principles

1. Reliability is defined by the user: For user facing workloads, measure the user experience, e.g query success ratio, as opposed to just server metrics such as CPU usage.

"Outage"?



Core Principle:

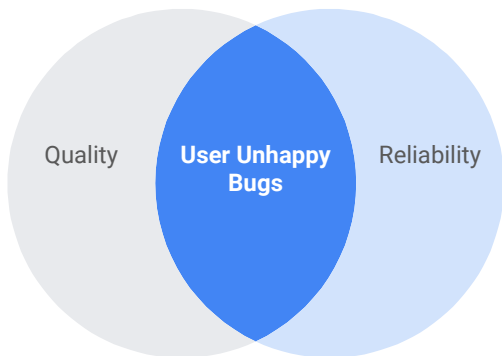
Reliability is Defined by the User

An "outage" is when the system is unusable, making the user unhappy.

Implication for measuring reliability:

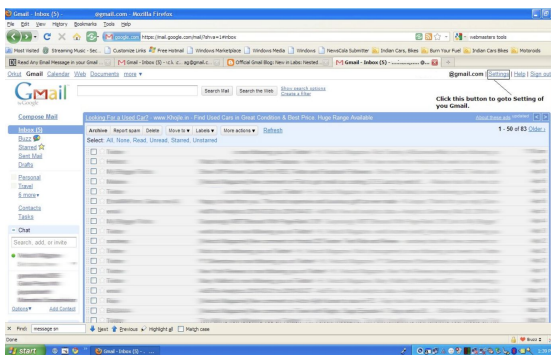
- Measure the user experience, e.g. request success ratio, as opposed to server-side metrics such as CPU usage.
- Measure as close to the user as possible.

Software Quality vs Reliability



Quality: Does every feature in the product work or not?

Reliability: Do all Critical User Journeys work for the user or not?



Temporary error () (500)

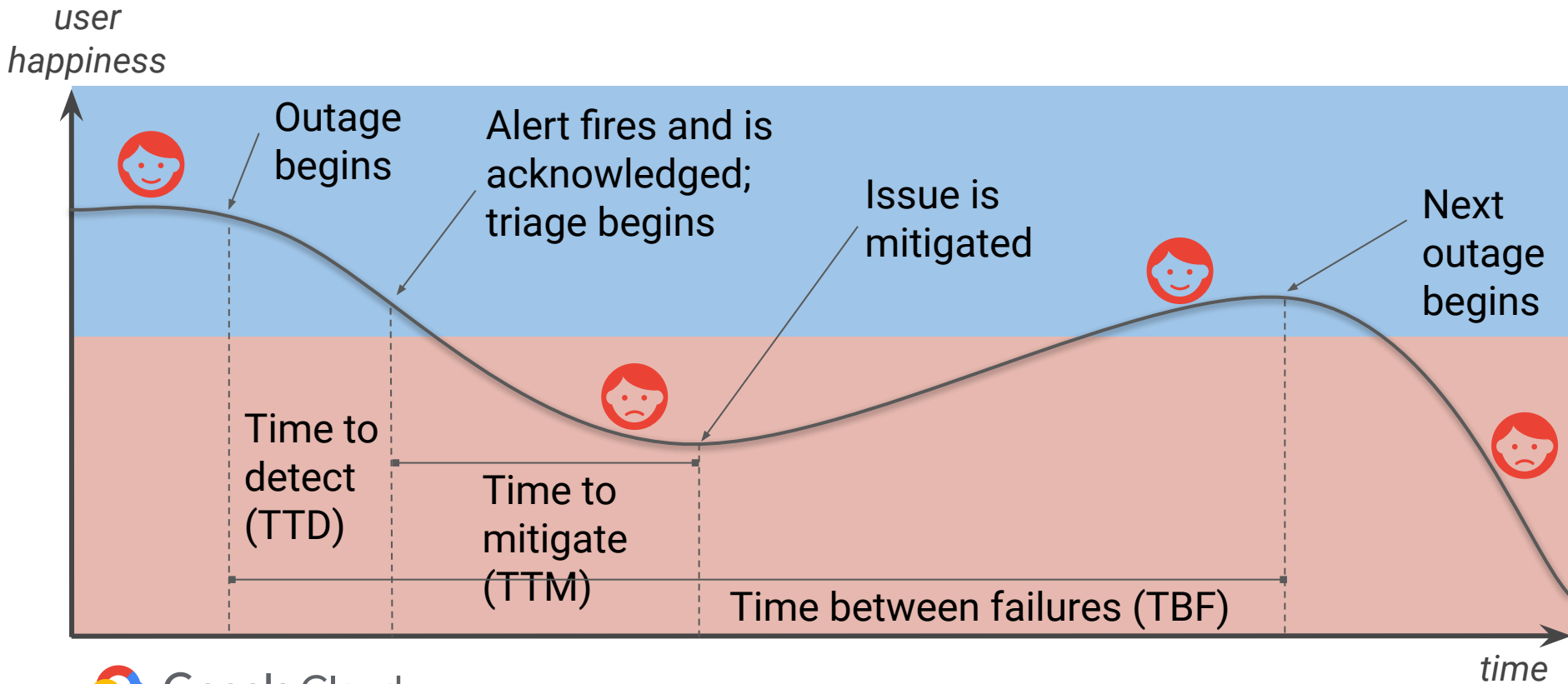
We're sorry, but your account is temporarily unavailable. We apologise for the inconvenience and suggest

If the issue persists, please visit the [Help Centre](#).

[Try Again](#) [Sign Out](#)

[Show Detailed Technical Info](#)

Outage Timelines and User Happiness



Principles

2. Sufficient reliability: Systems should be just reliable enough so that users are happy. Higher reliability comes at a steep cost.

“

100% is the wrong reliability target for basically everything.”

Benjamin Treynor Sloss, Vice President of 24x7 Engineering, Google



Glossary of Terms

SLI

service level
indicator: a
well-defined
measure of
success

SLO

service level
objective: a
top-line target
for fraction of
successful
interactions

SLA

service level
agreement:
business
consequences

Error Budget

proportion of “**affordable**”
unreliability;
one minus the SLO

CUJ

critical **user journey:**
specific steps that a user
takes to accomplish a goal



[SLO-Generator](#)

Watch 10

Fork 39

Star 247

```
apiVersion: sre.google.com/v2
kind: ServiceLevelObjective
metadata:
  name: prom-metrics-availability
  labels:
    service_name: prom
    feature_name: metrics
    slo_name: availability
spec:
  description: 99.9% of Prometheus requests return a good HTTP code
  backend: prometheus
  method: query_sli
  exporters:
  - prometheus
  service_level_indicator:
    expression: >
      sum(rate(prometheus_http_requests_total{handler="/metrics", code=~"2.."}[window]))
      /
      sum(rate(prometheus_http_requests_total{handler="/metrics"}[window]))
  goal: 0.999
```

Math for "Make This More Reliable"

$$E \propto \frac{TTD + TTM}{TBF} \times Impact$$

E is the rate at which error budget is being consumed by an outage

Increase **Reliability** \Leftrightarrow Reduce **Error Budget Burn Rate**

1

Increase Time
Between Failures
(**Fewer** Failures)

2

Reduce Time To Detect
(**Shorter** Failures)

3

Reduce Time To
Mitigate
(**Shorter** Failures)

4

Reduce Impact
(**Smaller** Failures)

Google SRE 架構設計原則





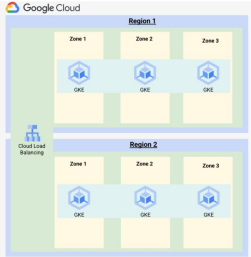
“

SRE is what happens when you ask a software engineer to design an operations team.”

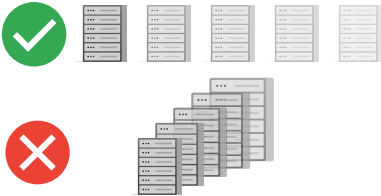
Benjamin Treynor Sloss, Vice President of 24x7 Engineering, Google

Designing for Reliability (Architecture Check)

Redundancy



Horizontal AutoScale



Tolerate Overload



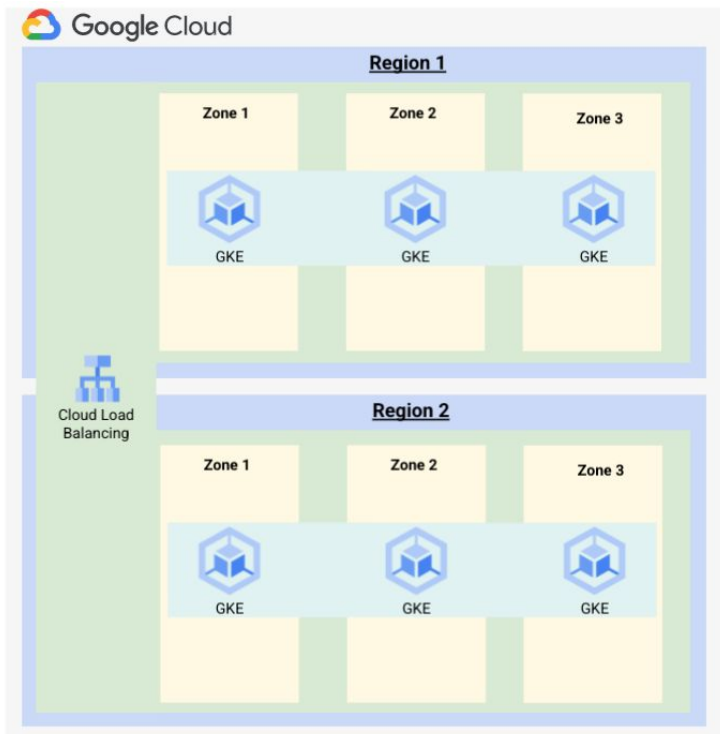
Rollback Support



Prevent Traffic Spikes



Design for Redundancy



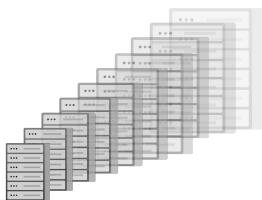
Reliable systems

- Have no single points of failure
- Have resources replicated across multiple failure domains, with automatic **failover**

Google's most critical products are deployed with **N+2 redundancy**.

- Top two zones can be down without the service falling below minimum capacity to handle peak load
- One zone hosting a GKE regional cluster may fail, but the workloads keep running a long as the other zones are fine.

Enable Horizontal Scalability



Horizontal scalability: Every component can handle more traffic/data by adding more resources.

- *Vertical Scaling* = Buy a **bigger, faster** machine / disk / database / cluster / zone
- *Horizontal Scaling* = Add **more** machines / disks / databases / clusters / zones

Sharding: Partition compute effort between tasks. Partition data with each compute shard handling a separate data partition.

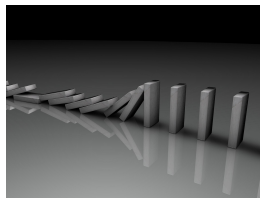
Cloud Spanner and sharding

Great apps run on great databases.

The only [enterprise-grade](#), [globally distributed](#), and [strongly consistent](#) database service built for the cloud specifically to combine the benefits of relational database structure with non-relational horizontal scale.



Tolerate Overload



Design systems to degrade gracefully under load

- Return slower or lower quality results

Where possible, systems should automatically scale up horizontally under load

Each replica must handle overload independently, with request throttling and circuit breakers.

Worst case scenario for systems where replicas crash when overloaded is **cascading failure**

[Managing Load](#) chapter in SRE Workbook

Support Rapid Rollback



Anything an operator can do to a service to change it must have a well-defined, well-tested method to undo.

Design, implement, and regularly test the rollback procedure for every operation.

Prevent Traffic Spikes



Requests must not be synchronized across clients.

Introduce exponential backoff with randomized delay in client error handling code.

Key definition:

- **Jitter:** Time interval added to or subtracted from the initially computed retry time.
- Each client generates a different randomly chosen jitter value in order to break synchronization with other clients and smooth out traffic.

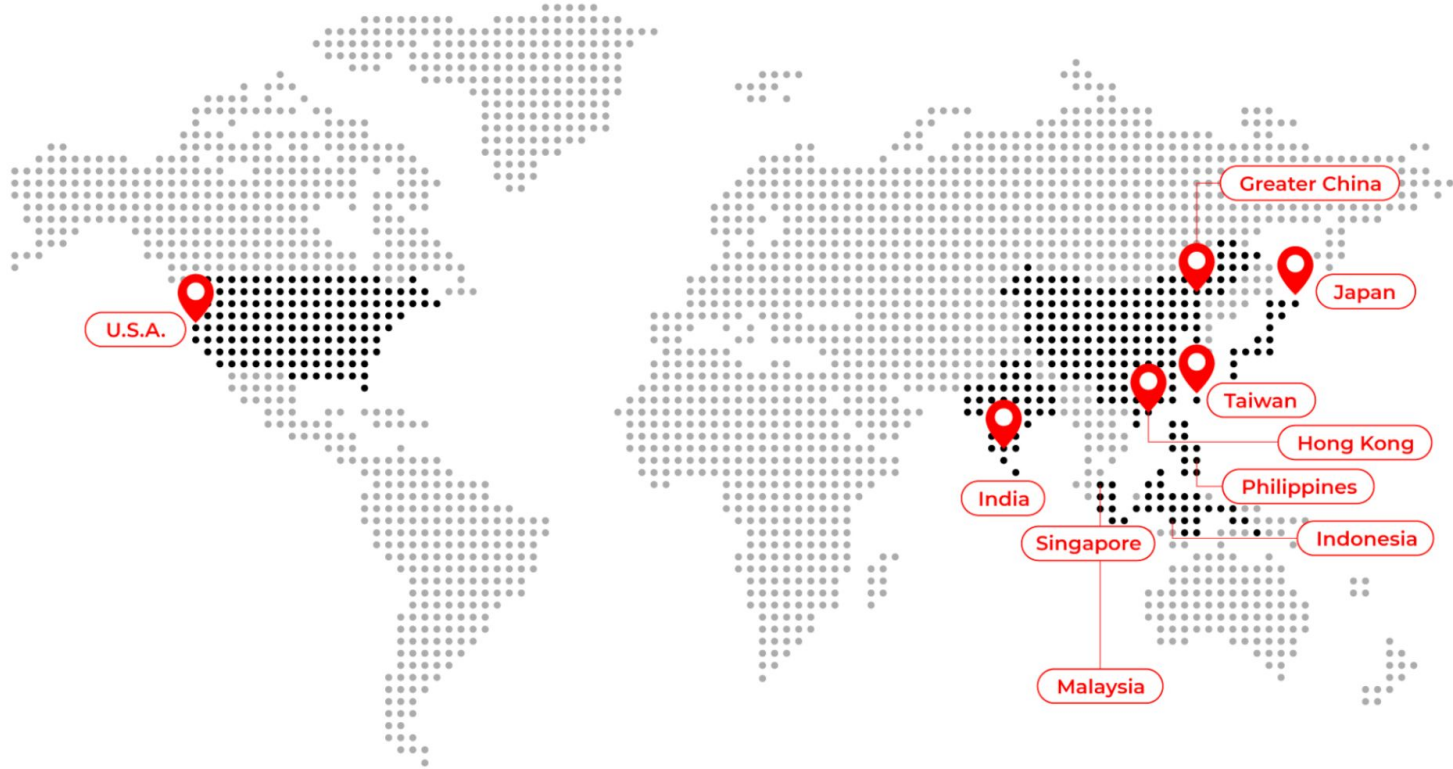
17LIVE的SRE實踐



林毅民 (Sammy Lin),
17LIVE Engineering Director



17LIVE IT Coverage



17LIVE SRE



2015

- 17Media founded
- Built on AWS



2016/8

- First DevOps Engineer joined



2017/1

- Migrated from beanstalk to ECS



2018/5

- Migrated to GCP



2018/6

- Migrated from node.js to golang



2019/8

- Building GRE

17LIVERs 每天最關心的事情: SRE's CUJ

Received Gift

Day Week Month









TOP 1

珈菲 🍷 Garfield 🍷 225 🍷 LIVE

791,027

Follow

- | | | | | |
|---|--|-------------------------|---------|--------|
| 2 |  | 濃濃 🍷 Serena 🍷 64 🍷 LIVE | 700,037 | Follow |
| 3 |  | 陳波波 🍷 Chan 🍷 150 🍷 LIVE | 444,826 | Follow |
| 4 |  | 孫卉彤 Candy 🍷 159 🍷 LIVE | 429,528 | Follow |
| 5 |  | Akemi 🍷 花花 🍷 153 🍷 LIVE | 389,365 | Follow |
| 6 |  | 靚 🍷 Qian 🍷 73 | 388,988 | Follow |
| 7 |  | 雯雯 🍷 Wai 🍷 101 🍷 LIVE | 345,529 | Follow |

Send Gift

Day Week Month










TOP 1

一個兩個三個四個五個 163

772,704

Follow

- | | | | | |
|---|---|----------------------|---------|--------|
| 2 |  | 一個兩個三個四個五個 166 | 700,030 | Follow |
| 3 |  | 莊孝偉 🍷 Zhang 🍷 143 | 388,374 | Follow |
| 4 |  | racemangt3 🍷 220 | 345,190 | Follow |
| 5 |  | jack 🍷 痲蛤蟆 🍷 51 | 299,518 | Follow |
| 6 |  | 辣個 🍷 嗆辣的 Ganny 🍷 125 | 285,841 | Follow |
| 7 |  | 心屬芋芋 🍷 皓皓 🍷 122 | 200,002 | Follow |



成為咒術師 領域展開 咒術之戰 Event Description


04/11 (Mon) 18:00:00-04/15 (Fri) 23:59:59
Time Remaining: 0 Day 07:39:41

活動攻略


- LIVER 累積 50,000 領域展開獲得個人頁背景
- LIVER 排名前 8 名獲得獎金

Search Following

Leaderboard



Top 1
舒服哥哥 🍷 Kimochi
709,310



04.09 04.16

戀愛腮紅 Event Description

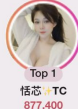


04/09 (Sat) 18:00:00-04/16 (Sat) 23:59:59
Time Remaining: 1 Day 00:42:12

活動攻略

LIVER 前 15 名可獲得妝漂亮獎金。

Search Following

Leaderboard

Top 1
恬恬 🍷 TC
877,400

Top 2
你的小可愛 🍷 ...
870,000

Top 3
宥宥 🍷
754,200

17LIVERs 每天最關心的事情: ~~SRE's~~ CUJ

Received Gift

Day Week Month



TOP 1

珈菲 🍷 Garfield 🍷 225 🍷 LIVE

791,027

Follow

- 2 濃濃 🍷 Serena 🍷 64 🍷 LIVE
444,037 Follow
- 3 陳波波 🍷 150 🍷 LIVE
444,826 Follow
- 4 孫卉彤 Candy 🍷 159 🍷 LIVE
429,528 Follow
- 5 Akemi 🍷 花花 🍷 153 🍷 LIVE
389,365 Follow
- 6 靚 🍷 73
388,988 Follow
- 7 雯雯 🍷 101 🍷 LIVE
345,529 Follow

Send Gift

Day Week Month




TOP 1

一個兩個三個四個五個 163

772,704

Follow

- 2 一個兩個三個四個五個 166
700,030 Follow
- 3 莊孝偉 143
388,374 Follow
- 4 racemang13 220
345,190 Follow
- 5 jack 癩蛤蟆 51
299,518 Follow
- 6 辣個 🍷 嗆辣の 🍷 anny 🍷 125
285,841 Follow
- 7 心屬芋芋 🍷 皓皓 🍷 122
200,002 Follow



成為咒術師 領域展開 咒術之戰 Event Description


04/11 (Mon) 18:00:00-04/15 (Fri) 23:59:59
Time Remaining: 0 Day 07:39:41

活動攻略


- 1. LIVER 累積 50,000 領域展開獲得個人頁背景
- 2. LIVER 排名前 8 名獲得獎金

Search Following

Leaderboard



Top 1
舒服哥哥 🍷 Kimochi
709,310



04.09 04.16

戀愛腮紅 Event Description

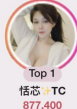


04/09 (Sat) 18:00:00-04/16 (Sat) 23:59:59
Time Remaining: 1 Day 00:42:12

活動攻略

LIVER 前 15 名可獲得妝漂亮獎金。

Search Following

Leaderboard

Top 1
恬恬 🍷 TC
877,400

Top 2
你的小可愛 🍷 ...
870,000

Top 3
宥宥 🍷
754,200

Product







Reliability is defined by the user

17LIVERs 每天最關心的事情: ~~SRE's~~ CUJ

Received Gift

Day Week Month






TOP 1
 珈菲 🍷 Garfield 🍷 225 🍷 LIVE
 791,027
 Follow

- 2  濃濃 🍷 Serena 🍷 64 🍷 LIVE
 700,037 Follow
- 3  陳波波 🍷 Chan 🍷 150 🍷 LIVE
 444,826 Follow
- 4  孫丹彤 Candy 🍷 Candy 🍷 159 🍷 LIVE
 429,528 Follow
- 5  Akemi 🍷 花花 🍷 153 🍷 LIVE
 389,365 Follow
- 6  靚 🍷 Qian 🍷 73
 388,988 Follow
- 7  雯雯 🍷 Wai 🍷 101 🍷 LIVE
 345,529 Follow

Send Gift

Day Week Month


TOP 1
 一個兩個三個四個五個 🍷 163
 772,704
 Follow

- 2  一個兩個三個四個五個 🍷 166
 700,030 Follow
- 3  莊孝偉 🍷 Zhang 🍷 143
 388,374 Follow
- 4  racemangt3 🍷 225
 345,190 Follow
- 5  jack 癩蛤蟆 🍷 51
 299,518 Follow
- 6  辣個 🍷 嗆辣的 🍷 Ganny 🍷 125
 285,841 Follow
- 7  心屬芋芋 🍷 皓皓 🍷 122
 200,002 Follow


4.11 - 4.21
咒術之戰 領域展開
 男 LIVER 專屬活動
 成為咒術師 領域展開 咒術之戰 Event Description
 04/11 (Mon) 18:00:00-04/15 (Fri) 23:59:59
 Time Remaining: 0 Day 07:39:41

活動攻略

- LIVER 累積 50,000 領域展開獲得個人頁背景
- LIVER 排名前 8 名獲得獎金

Search Following

Leaderboard


Top 1
 舒服哥哥 🍷 Kimochi
 709,310

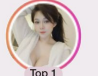
17 妝漂亮
 戀愛腮紅
 04.09 04.16
 戀愛腮紅 Event Description
 04/09 (Sat) 18:00:00-04/16 (Sat) 23:59:59
 Time Remaining: 1 Day 00:42:12


活動攻略


LIVER 前 15 名可獲得妝漂亮獎金。

Search Following

Leaderboard


Top 1
 恬恬 🍷 TC
 877,400


Top 2
 你的小可愛 🍷...
 870,000

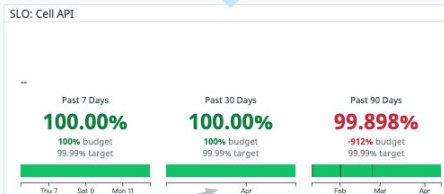

Top 3
 宥宥 🍷
 754,200



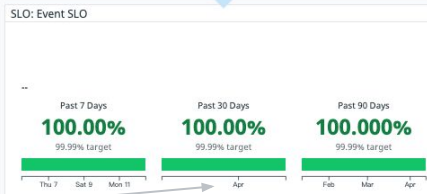
來源: <https://www.facebook.com/groups/2616981278627207/permalink/3208110619514267/>

以CUJ為主的加權警示(Failure Detection)

系統API 加權分數(SLO)



活動/榜單系統 加權分數(SLO)



付費系統 加權分數(SLO)



加權指數
+7 days

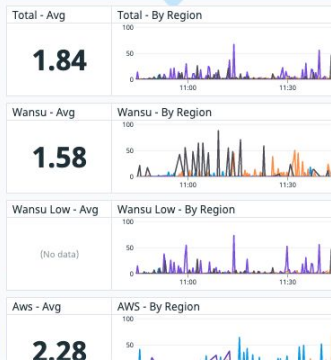
Video Latency (sec)

(TW: Blue, JP: Red, US: Purple, HK: Gra)



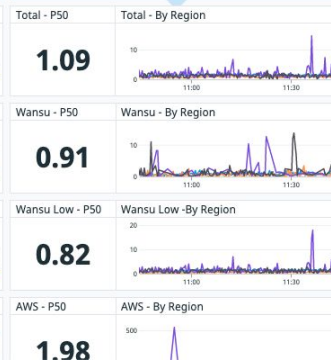
Buffering Ratio

(TW: Blue, JP: Red, US: Purple, HK: Gra)

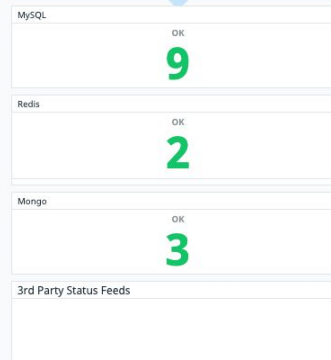


First Render Time(Sec)

(TW: Blue, JP: Red, US: Purple, HK: Gra)



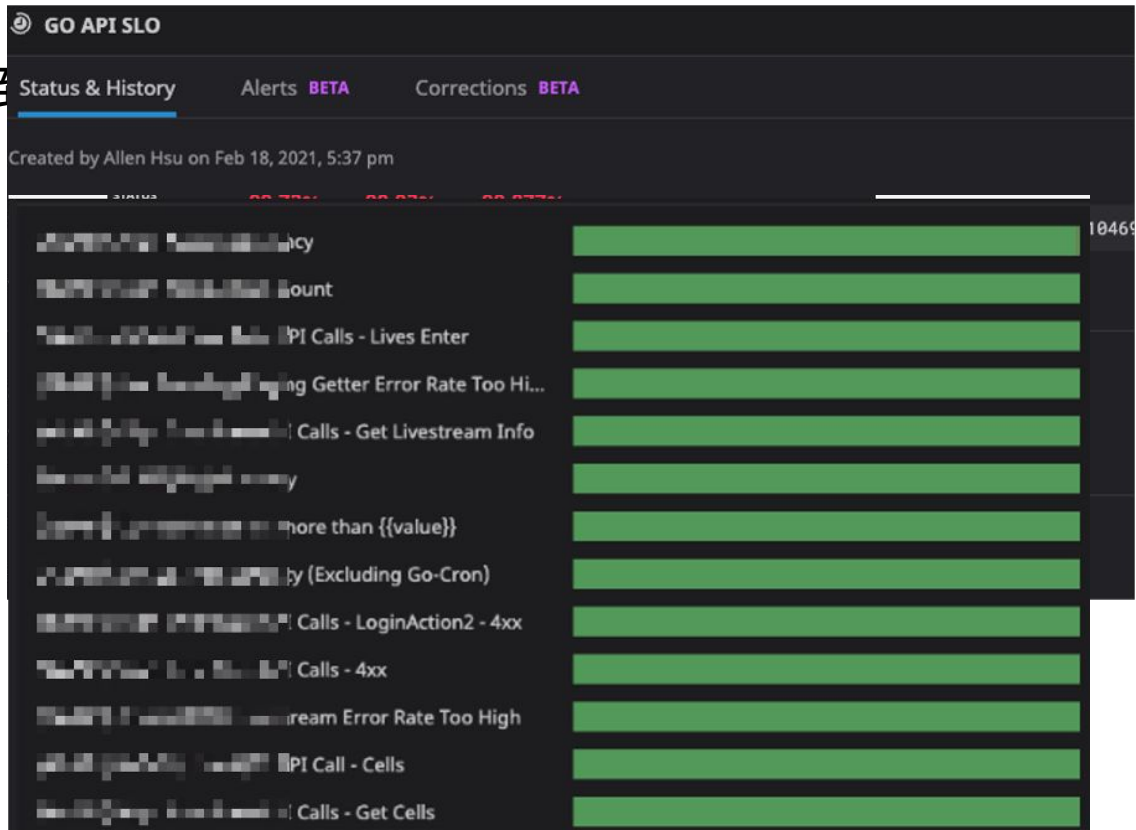
Infra Stability



CUJ下的關鍵監控: Observability Dashboard

當CUJ的警示發生後，快速切換到第二層支援服務的狀態監測

- API Latency
- API 錯誤率
- 關鍵子系統的效能
- 工作負荷量(QPS)



Capacity Planning: 預估上限並提前採取措施

NEWS

「SUPERSONIC on 17LIVE」 配信アーティスト発表！

The poster features a stylized landscape with a blue sky, white clouds, and green foliage on the left and right sides. The event title 'SUPERSONIC' is prominently displayed in the center. Below the title, the dates '9.18 sat' and '9.19 sun' are listed. The lineup of artists is presented in two columns, with names in white text on colored rectangular backgrounds. The 17LIVE logo is at the bottom.

9.18 sat	9.19 sun
SKY-HI	R3HAB
どんぐりず	DIGITALISM
AURORA	NiziU
石野 卓球	きゃりーぱみゅぱみゅ
BE:FIRST	STEVE AOKI
ALAN WALKER	special closing B2B set STEVE AOKI/ZEDD/ALAN WALKER
CLEAN BANDIT (DJ SET)	
ZEDD	

17LIVE

- 與PM與行銷人員確認預估人數

行銷：伺服器有多少就開多少，資料庫有多大就開多大！？ 反正錢不是問題？

Capacity Planning: 預估上限並提前採取措施

NEWS

「SUPERSONIC on 17LIVE」 配信アーティスト発表！

The poster features a stylized landscape with a sunset sky, clouds, and a body of water. The event title 'SUPERSONIC' is prominently displayed in the center. Below it, the dates '9.18 sat' and '9.19 sun' are listed. The lineup of artists is presented in two columns, each with a colored bar representing the artist's name. The 17LIVE logo is at the bottom.

9.18 sat	9.19 sun
SKY-HI	R3HAB
どんぐりず	DIGITALISM
AURORA	NiziU
石野 卓球	きゃりーぱみゅぱみゅ
BE:FIRST	STEVE AOKI
ALAN WALKER	special closing B2B set STEVE AOKI/ZEDD/ALAN WALKER
CLEAN BANDIT (DJ SET)	
ZEDD	

17LIVE

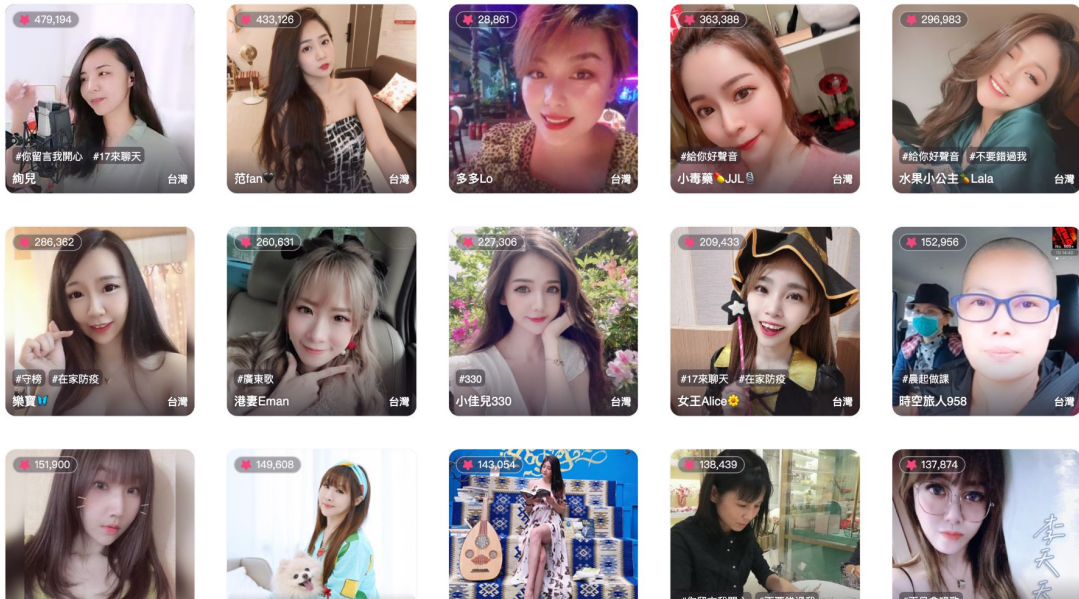
- 與PM與行銷人員確認預估人數
- PreProduction以1:1測試+Locust:
 - Auto scale 機制的 最大值是否充足?
 - 使用者體驗是否符合預期
 - Pre-scaling 機制的調教
 - GCP Quota 預先調整
 - 介接的第三方服務以及資料庫能否能乘載當時的量級
 - 效能瓶頸?
 - [Kubernetes 上限](#)
- 以測試結果推測資源增加的趨勢

自動降載機制：避免服務崩潰

關鍵直播不能停



女生



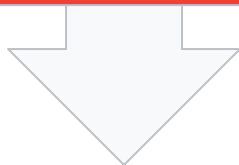
API呼叫頻率
(Rate Limit)

關閉次要功能
(Secondary Functions)

限制用戶數
(Crowd Capacity)

非關鍵服務關閉，擴
充關鍵服務容量

想盡辦法
Workaround



減少重複工作(Toil Reduction): 活動資源優化流程



alfred APP 6:30 PM

Scheduled scaling

Success

```
Patch 'k8sprod-golives-main' successfully, minReplicas: '5'->'85' output: horizontalpodautoscaler.autoscaling/k8sprod-golives-main patched
```

```
Patch 'k8sprod-gotrade-main' successfully, minReplicas: '4'->'12' output: horizontalpodautoscaler.autoscaling/k8sprod-gotrade-main patched
```

```
Patch 'k8sprod-gocells-main' successfully, minReplicas: '10'->'25' output: horizontalpodautoscaler.autoscaling/k8sprod-gocells-main patched
```

```
Patch 'k8sprod-gousersearch-main' successfully, minReplicas: '2'->'10' output: horizontalpodautoscaler.autoscaling/k8sprod-gousersearch-main patched
```

```
Patch 'k8sprod-revprox-jp-main' successfully, minReplicas: '18'->'36' output: horizontalpodautoscaler.autoscaling/k8sprod-revprox-jp-main patched
```

```
Patch 'k8sprod-goapi-main' successfully, minReplicas: '9'->'70' output: horizontalpodautoscaler.autoscaling/k8sprod-goapi-main patched
```



alfred APP 9:30 PM

Scheduled scaling

Success

```
Reset 'k8sprod-gotrade-main' successfully, minReplicas: '4' output: horizontalpodautoscaler.autoscaling/k8sprod-gotrade-main patched
```

Scheduled scaling

Success

```
Reset 'k8sprod-gousersearch-main' successfully, minReplicas: '2' output: horizontalpodautoscaler.autoscaling/k8sprod-gousersearch-main patched
```

Scheduled scaling

Success

```
Reset 'k8sprod-gocells-main' successfully, minReplicas: '10' output: horizontalpodautoscaler.autoscaling/k8sprod-gocells-main patched
```

減少重複工作(Toil Reduction): 快速簽核

大家如果還在用 UI 跟簽核流程的話

Add principals to " ABCDE "

Add principals and roles for " ABCDE " resource

Enter one or more principals below. Then select a role for these principals to grant them access to your resources. Multiple roles allowed. [Learn more](#)

New principals

abchakra@google.com

Role *
Kubernetes Engine Admin

Condition
[Add condition](#)

Full management of Kubernetes Clusters and their Kubernetes API objects.

+ ADD ANOTHER ROLE

SAVE CANCEL



跟17LIVE 一起把他自動化吧!

Supervisor Approval WORKFLOW Dec 13th, 2021 at 5:44 PM

A New Request from @Brent Chang (HQ SRE) is created.

This form has not yet been submitted. Please review the below information and click Review to submit the form.

- No: 1639388632654
- Applicant: @Brent Chang (HQ SRE),
- Request for: Kubernetes Engine Admin role for installing grafana agent in Prod,
- Form Url: https://docs.google.com/spreadsheets/d/1K9ycQC4edRHhs-A8_zSo2T6WmlfGdc95u01Z7vEhTV0/edit#gid=1812408247,
- Supervisor: @Sammy (HQ ENG)

Description

- Timestamp: Mon Dec 13 2021 17:43:48 GMT+0800 (Taipei Standard Time)
- Email Address: brentchang@17.media
- Slack Username: Brent Chang (HQ SRE)
- Department: HQ/Reliability
- Is it bot account?: No
- What item do you want to apply for access: Kubernetes Engine Admin role for installing grafana agent in Prod
- These permissions must be kept secret and for work usage only.: Agree

6 replies

Brent Chang (HQ SRE) 4 months ago
學昂有空再麻煩 Approve, 要裝 grafana tracing agent 在 Prod 上 🙏

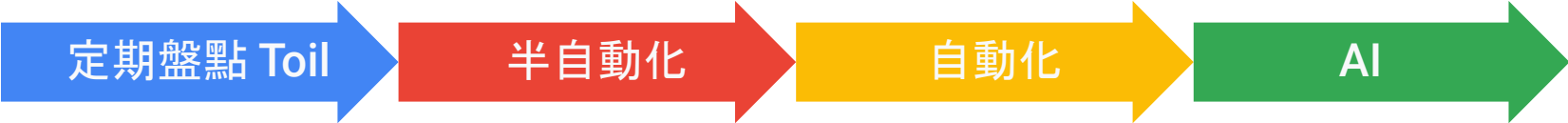
Supervisor Approval WORKFLOW 4 months ago

Confirmation submission from @Sammy (HQ ENG)

Approve / Reject
Approve

Comments
ok

減少重複工作(Toil Reduction):



打造Incremental Change: 資料庫Schema更新的竅門

rubenv/sql-migrate:
Golang DB Migrate Tool

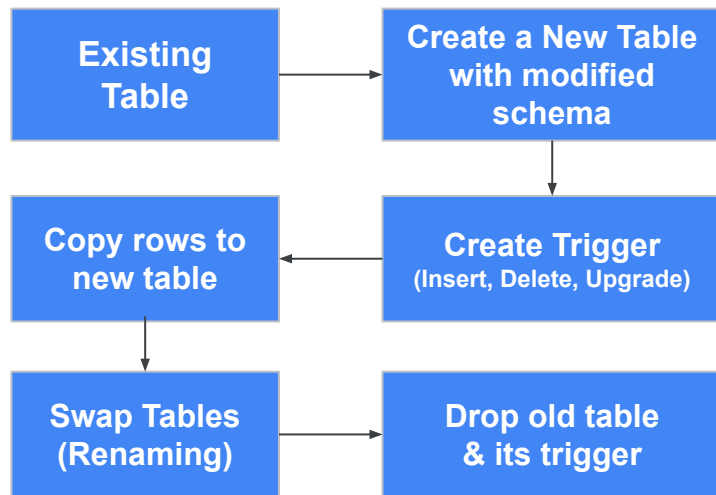


```
-- +migrate Up
CREATE TABLE IF NOT EXISTS `users` (
  `id` VARCHAR(28) NOT NULL,
  `email` VARCHAR(100) NOT NULL,
  `name` VARCHAR(20) NOT NULL,
  `age` int(10) UNSIGNED,
  `birthday` DATETIME,
  `created_at` DATETIME NOT NULL,
  `updated_at` DATETIME NOT NULL,
  PRIMARY KEY (`id`),
  CONSTRAINT email_unique UNIQUE(email)
)ENGINE = InnoDB DEFAULT CHARSET=utf8mb4;

-- +migrate Down
DROP TABLE `users`;
```

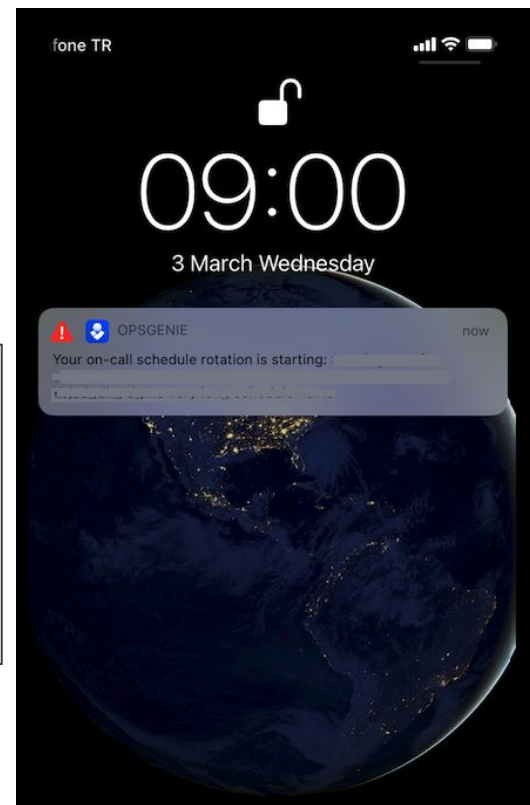
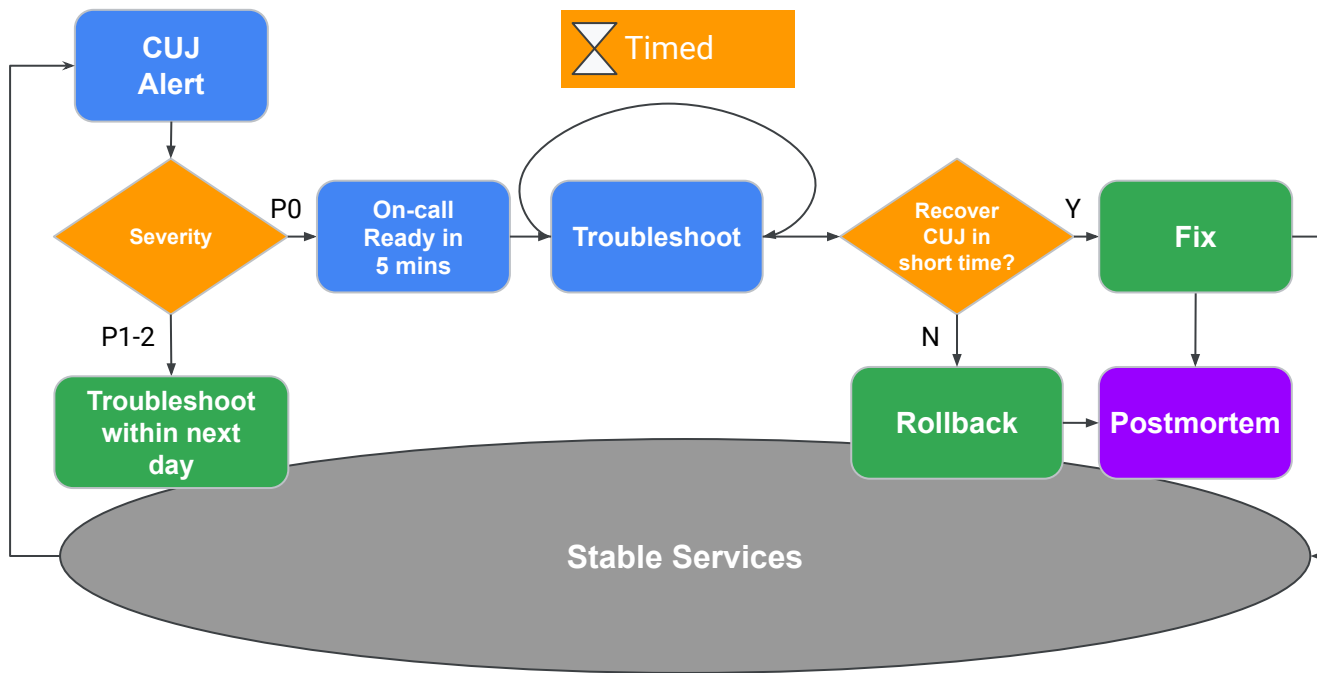


pt-online-schema-change:
Percona-toolkits, by Percona



PO處理自動化/手動流程:

(Emergency response documentation and automation)



緊急障礙處理流程: (Coordinated Emergency Response)



A: 什麼時候好?

B: 為什麼會壞掉?

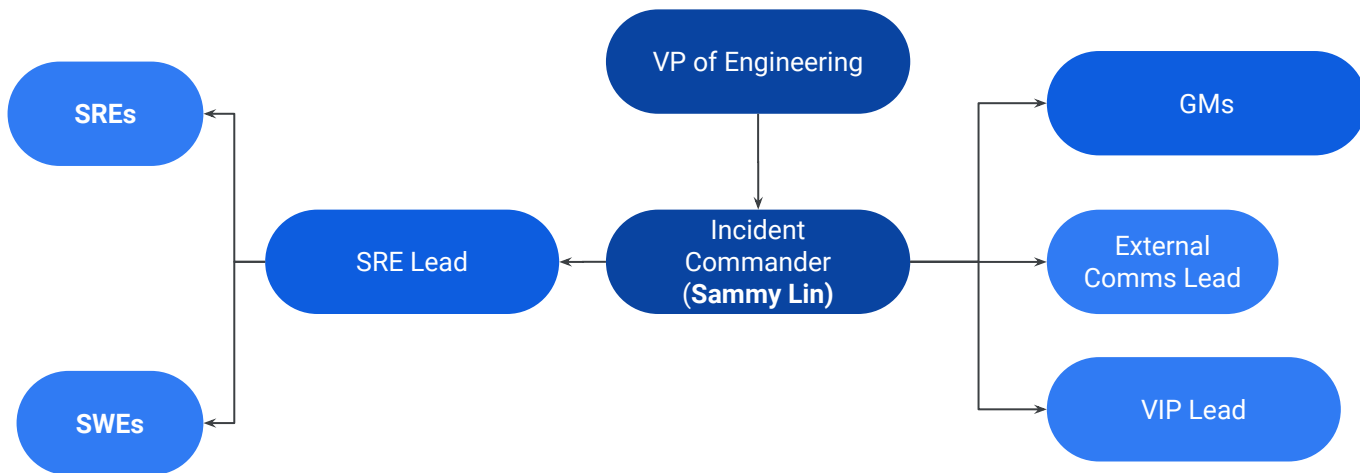
C: 客戶在抱怨了

D: 這不就是跟上次一樣, 就照上次的方法處理啊。

緊急障礙處理流程: (Coordinated Emergency Response)

首要工作:

- 讓後端SRE/SWE能專心查找問題
- 讓StakeHolder可以快速了解進度



PostMortem範本

Google SRE Book Template:

Shakespeare Sonnet++ Postmortem (incident #465)

Date: 2015-10-21

Authors: jennifer, martym, agoogler

Status: Complete, action items in progress

Summary: Shakespeare Search down for 66 minutes during period of very high interest in Shakespeare d sonnet.

Impact:¹⁶³ Estimated 1.21B queries lost, no revenue impact.

Root Causes:¹⁶⁴ Cascading failure due to combination of exceptionally high load and a resource leak when searches failed due to terms not being in the Shakespeare corpus. The newly discovered sonnet used a word that had never before appeared in one of Shakespeare's works, which happened to be the term users searched for. Under normal circumstances, the rate of task failures due to resource leaks is low enough to be unnoticed.

Trigger: Latent bug triggered by sudden increase in traffic.

Resolution: Directed traffic to sacrificial cluster and added 10x capacity to mitigate cascading failure. Updated index deployed, resolving interaction with latent bug. Maintaining extra capacity until surge in public interest in new sonnet passes. Resource leak identified and fix deployed.

Detection: Borgmon detected high level of HTTP 500s and paged on-call.

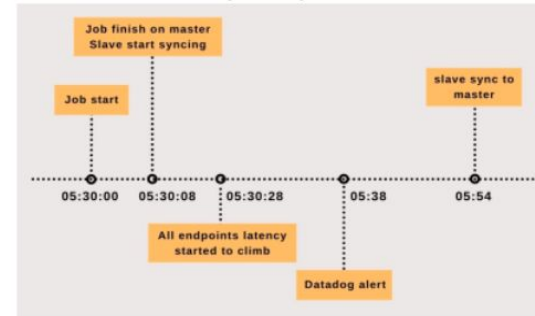
Action Items:¹⁶⁵

Action Item	Type	Owner	Bug
Update playbook with instructions for responding to cascading failure	mitigate	jennifer	n/a **DONE**
Use flux capacitor to balance load between clusters	prevent	martym	Bug 5554823 **TODO**
Schedule cascading failure test during next DiRT	process	docbrown	n/a **TODO**

17LIVE Postmortem Template:

RCA Report

1. Problem :
 - a. < Fill Customer-facing issues here >
2. Issue Level : < P0-P4 >
3. Impact : < Global or Critical or Regional >
4. Expectation of Revenue Loss
 - a. < Fill estimated loss here >
5. Details of the Root Cause
 - a. **Root Cause:**
 - b. **Details:**
< Details of the outage, including the timeline >



6. Recovery Time:
< Time to recovery >
7. Action Items and ETA:

Take Away

客戶體驗為SRE之尊



承認吧！世界是不完美的



Dev+Op共同的OKR

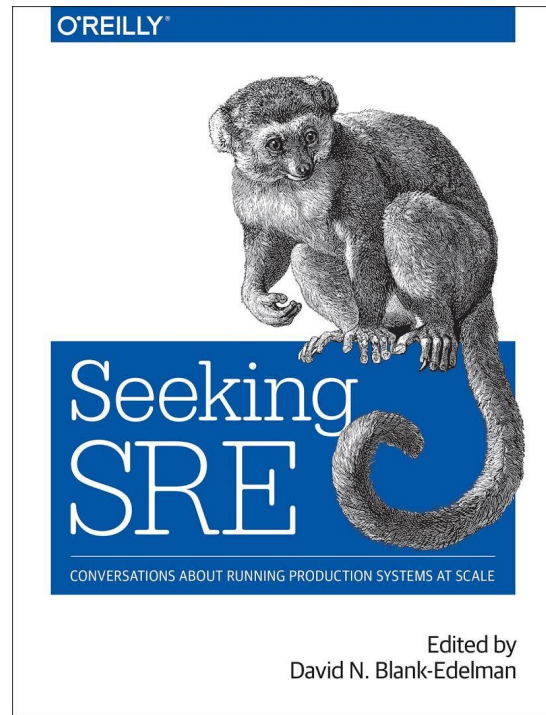
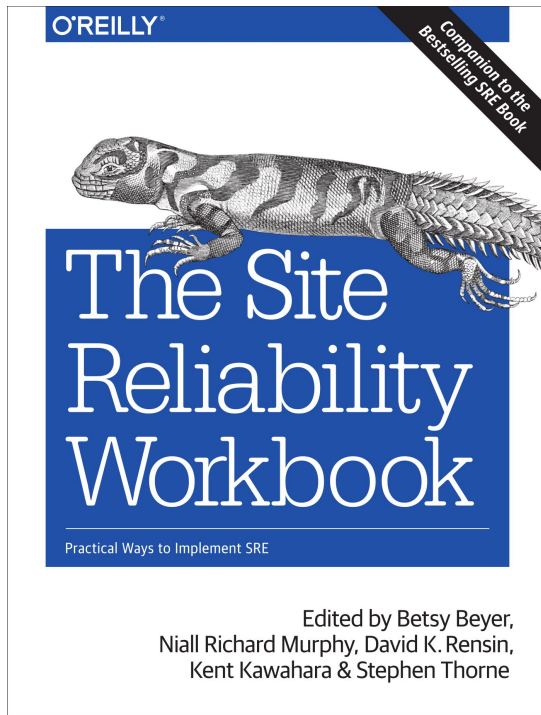
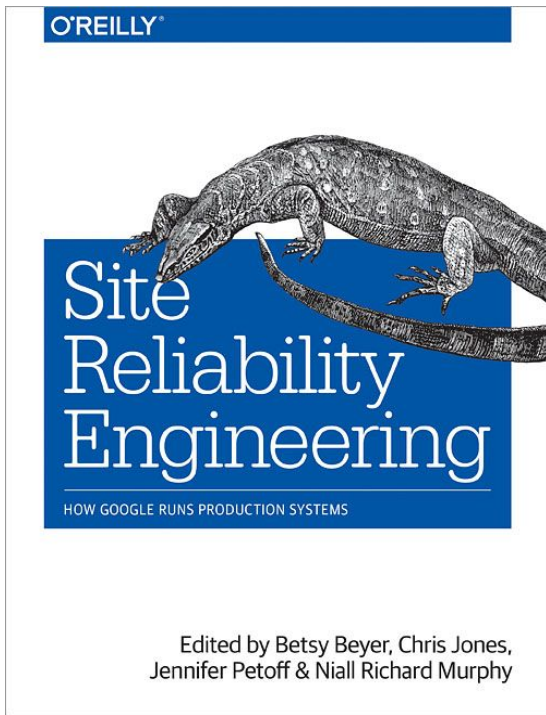


準備好Undo Button



不要累死你的SRE





Error Budgets Must Be Used!

Reliability level	Allowed unreliability window		
	per year	per quarter	per 30 days
90%	36.5 days	9 days	3 days
95%	18.25 days	4.5 days	1.5 days
99%	3.65 days	21.6 hours	7.2 hours
99.5%	1.83 days	10.8 hours	3.6 hours
99.9%	8.76 hours	2.16 hours	43.2 minutes
99.95%	4.38 hours	1.08 hours	21.6 minutes
99.99%	52.6 minutes	12.96 minutes	4.32 minutes
99.999%	5.26 minutes	1.30 minutes	25.9 seconds

0% error budget burn is the wrong target

Error Budget is the basis for innovation and agility

Unused error budget indicates over-investment in reliability and inadequate risk-taking