

在地端 K8s 的排障故事分享

玉山銀行 | 核心組副主任工程師 李啓維

我是誰



玉山銀行智能金融處核心組

- 照料所有 MLaaS k8s
- 部署維運的疑難雜症
- 展覽咖啡長板

演講經歷

- 2023 k8s summit / CPU throttling
- 2024 hello devSecOps / 導入 OPA



✉ | sean22492249@gmail.com

🌐 | <https://kiwi-walk.com>

M | <https://sean22492249.medium.com/>

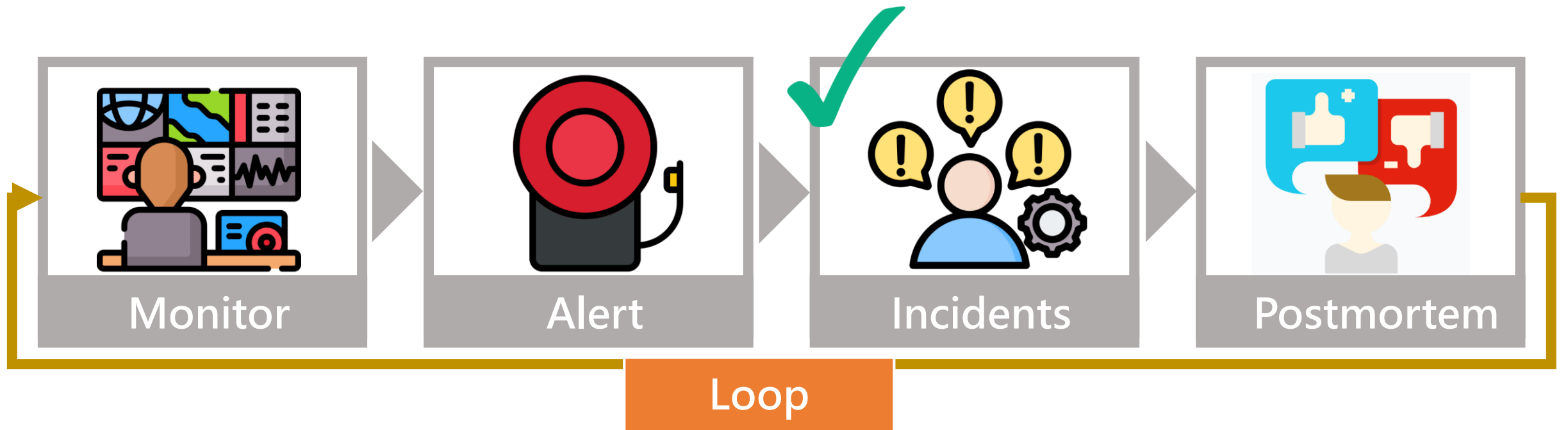


今天要來跟大家分享
「當初笑不出來的故事」

要用地端的 k8s 嗎







不探討

- monitor 如何規劃
- alert/incidents 如何偵測
- 不會過於深入技術

要分享的

- 出事時，我自己的經驗整理
- 配上血的故事
- 與許多的梗圖

A man in a white shirt is working on a yellow car seat in a workshop. The seat is mounted on a metal frame. The background is dark with some tools and equipment visible.

出事了阿伯

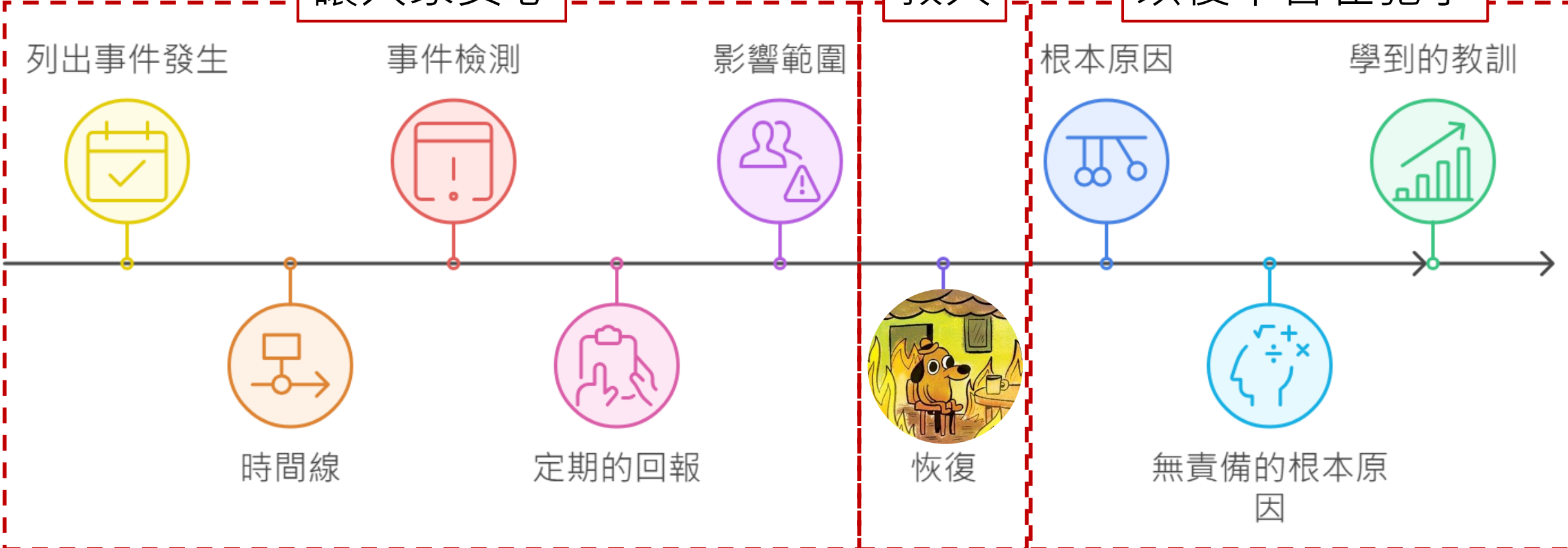
參戰!

常見的事件處理 SOP

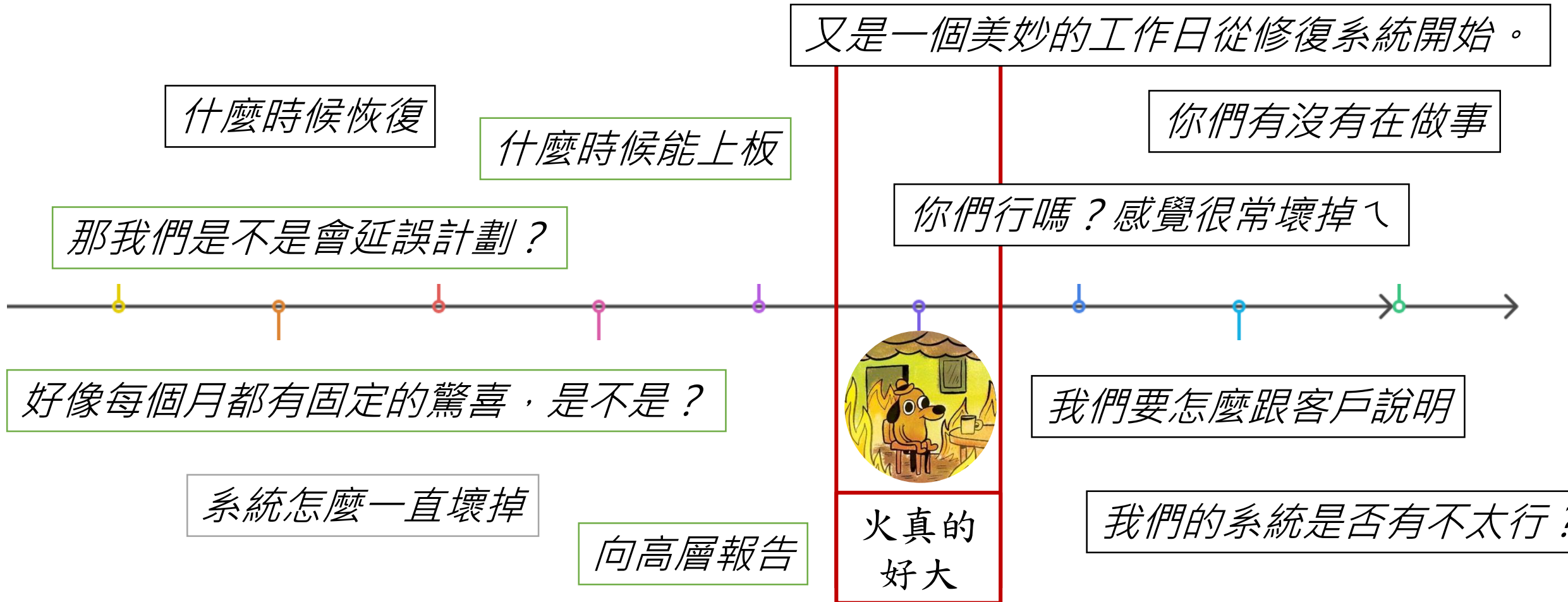
讓大家安心

救火

以後不會在犯了



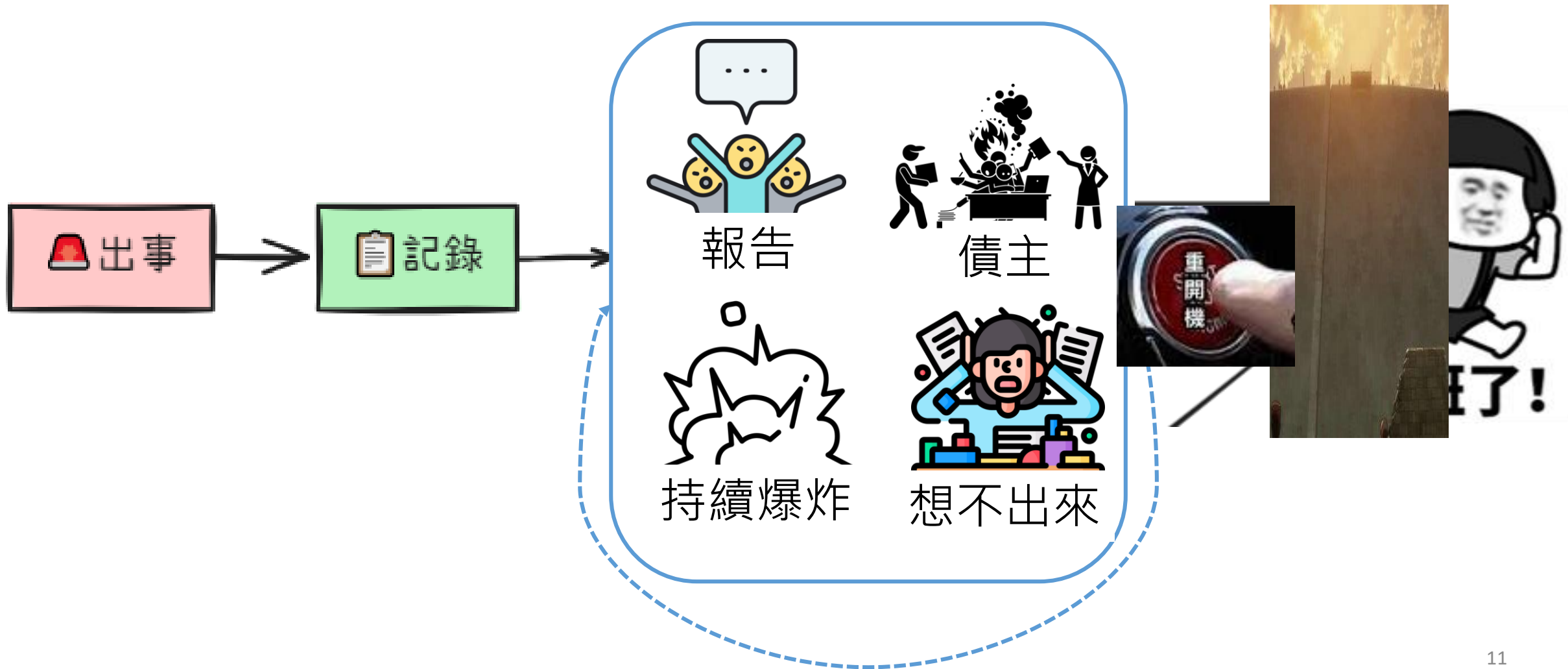
事件處理時會遇到的鼓勵



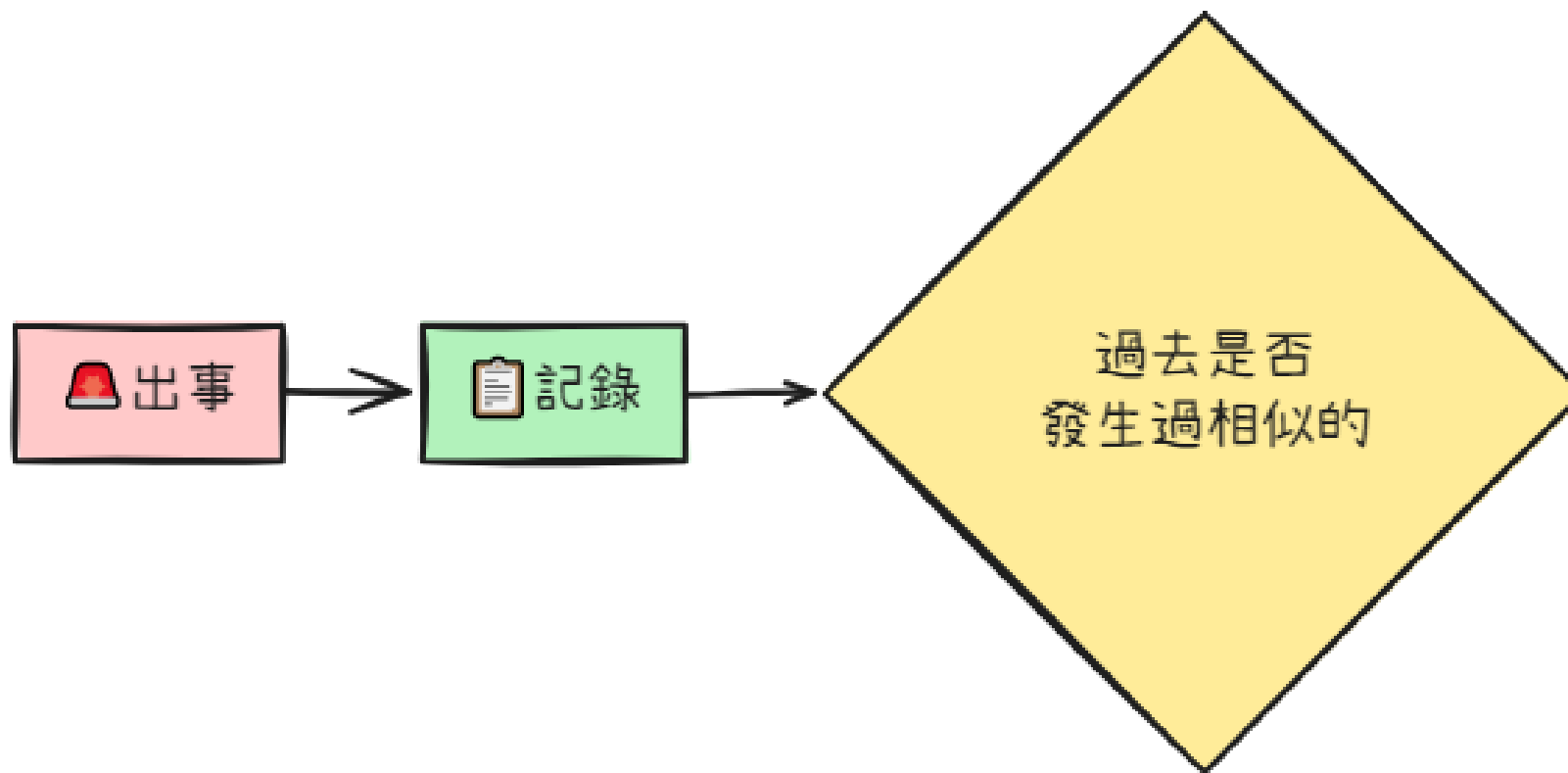
身為處理人的內心



實際上我們在(想)做的



第一步



優先確認資訊來源 & 記錄編號

- who, when, what



事件 → 記錄



所有人 上週六patch目前看還是有更新到網路相關套件，所以大家記得再確認一下有登記patch的系統有無異常，特別是服務如果是用container相關方式啟動的，有機會受到影響

感謝大家配合，都已填寫完畢，但請系統負責人注意，有填 自動重開機者，需參與月變更驗證服務工作，謝謝

慘 3 1



每個人 上週六patch目前看還是有更新到網路相關套件，所以...

目前有發現什麼狀況嗎？



目前有發現什麼狀況嗎？

昨天早上有發現到 連不上，初步看一下log看起來都是harbor相關元件在互相連線時發生問題，同步去查看syslog發現出問題時間大概就跟patch執行時間差不多，後來把服務全部重啟之後就恢復正常了

尋找過去的參考

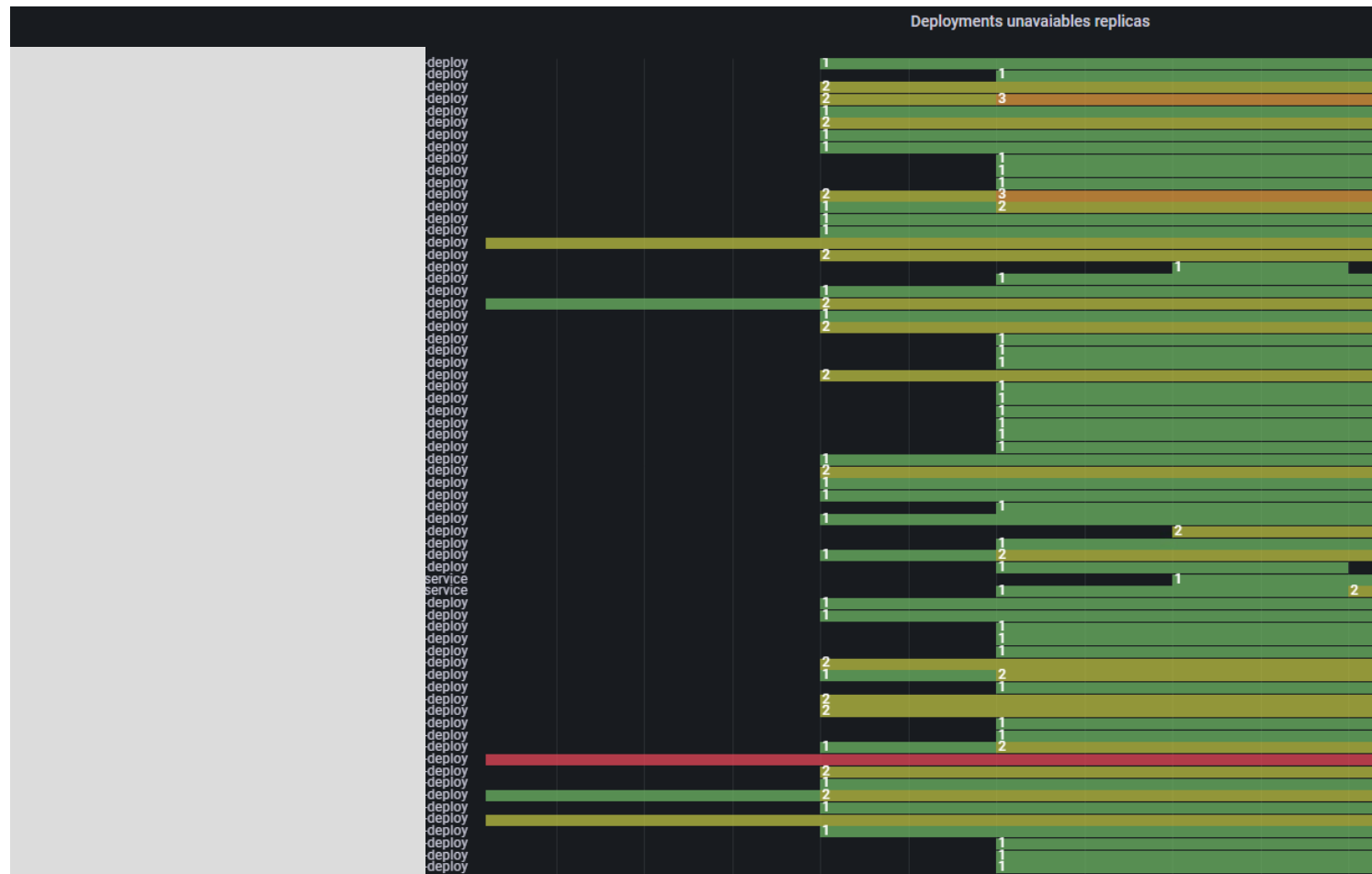
事件關鍵字：

- CVE Patch
- container



歷史痛苦：

- 2023.12 發生過
- 同樣是 container 問題



回到那時候，來自 CVE 的更新

主旨: [REDACTED] 通報單-(中)-GNU C函式庫存在緩衝區溢位漏洞(CVE-2023-4911)

各位好：

發佈 [REDACTED] 資安通報，請各位於回覆期限內回報處理方式，謝謝。

一、通報主旨：GNU C函式庫存在緩衝區溢位漏洞(CVE-2023-4911)

二、通報描述：

1. 風險等級： [REDACTED]
2. 通報回覆期限： [REDACTED] (註：若無法完成影響範圍之評估，請務必於期
3. 通報描述：
 - 擁有存取權限的攻擊者成功利用時可取得Root權限。
 - 細節請參考附件。

三、建議措施：請安排更新作業。

四、 [REDACTED]

NetworkManager 有重啟

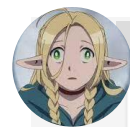


當時的出事



deployment 出現 unavailable

- 40+ | 10分鐘後自動恢復
- 8+ | 無法恢復



大家好，Linux 年度 auto patch，原預估 MLaaS 服務不受影響，但稍早發現 ...
今天有詢問 團隊，昨天MLaaS pod無預警重啟的影響
影響時間：18:45 ~ 18:50
影響筆數：19筆
影響筆數：424筆
pod重啟之後服務皆恢復正常

各種事件群組開始瘋狂回報



大家好，Linux 年度 auto patch，原預估 MLaaS 服務不受影響，但稍早發現 MLaaS pod 有發生無預警重啟現象，上圖是 pod 無預警重啟期間同時有發生 2 pods (含) 以上失效，可能影響服務的清單

下午 08:03



多數影響時間在 18:49 ~ 18:5

下午 08:04



以上專案需要請各自負責同仁，明日上班後確認一下影響服務狀況

下午 08:12



部分節點發生了問題

```

[2023/12/06 09:04:07] [error] [filter:kubernetes:kubernetes.0] kubelet upstream connection error
[2023/12/06 09:04:17] [error] [net] connection #71 timeout after 10 seconds to: kubernetes.default.svc:443
[2023/12/06 09:04:17] [error] [net] socket #71 could not connect to kubernetes.default.svc:443
[2023/12/06 09:04:17] [error] [filter:kubernetes:kubernetes.0] kubelet upstream connection error
[2023/12/06 09:04:27] [error] [net] connection #71 timeout after 10 seconds to: kubernetes.default.svc:443
[2023/12/06 09:04:27] [error] [net] socket #71 could not connect to kubernetes.default.svc:443
[2023/12/06 09:04:27] [error] [filter:kubernetes:kubernetes.0] kubelet upstream connection error
[2023/12/06 09:04:27] [error] [upstream] connection #79 to 443 timed out after 10
[2023/12/06 09:04:27] [error] [upstream] connection #80 to 443 timed out after 10
[2023/12/06 09:04:27] [error] [upstream] connection #82 to 443 timed out after 10
[2023/12/06 09:04:27] [error] [upstream] connection #84 to 443 timed out after 10
[2023/12/06 09:04:27] [error] [upstream] connection #85 to 443 timed out after 10
[2023/12/06 09:04:27] [error] [upstream] connection #86 to 443 timed out after 10
[2023/12/06 09:04:27] [error] [upstream] connection #70 to 443 timed out after 10
[2023/12/06 09:04:27] [error] [upstream] connection #81 to 443 timed out after 10
[2023/12/06 09:04:27] [info] [input:tail:tail.0] inode=840 link_rotated: /var/log/containers/kube-state-me
7430bac8d28ba2c582139073b7d2d835d40ea221980914569.log
[2023/12/06 09:04:27] [info] [input:tail:tail.0] inode=8409238 handle rotation(): /var/log/containers/kube
c776be0547430c828a2c582139073b7d2d835d40ea221980914569.log => /var/log/pods/kube-system_kube-state-metr
metrics/416.log (deleted)
[2023/12/06 09:04:27] [info] [input:tail:tail.0] inode=87 link_rotated: /var/log/containers/namespace-pro
51cf72fad4fcd3302228c0857d1baaedf2ae6d098e6859749c139a5ee092af0.log
[2023/12/06 09:04:27] [info] [input:tail:tail.0] inode=8409240 handle rotation(): /var/log/containers/name
isioner-761cf72fad4fcd3302228c0857d1baaedf2ae6d098e6859749c139a5ee092af0.log => /var/log/pods/default_names

level=info msg="regenerating all endpoints" reason="one or more identities created or deleted" subsys=endpoint-manager
level=info msg="regenerating all endpoints" reason="one or more identities created or deleted" subsys=endpoint-manager
level=info msg="regenerating all endpoints" reason="one or more identities created or deleted" subsys=endpoint-manager
level=info msg="regenerating all endpoints" reason="one or more identities created or deleted" subsys=endpoint-manager
level=info msg="regenerating all endpoints" reason="one or more identities created or deleted" subsys=endpoint-manager
level=info msg="regenerating all endpoints" reason="one or more identities created or deleted" subsys=endpoint-manager
level=warning msg="Network status error received, restarting client connections" error="an error on the server side:
thook/start-kube-apiserver-admission-initializer ok\n[+]poststarthook/generic-apiserver-start-informers ok\n[+]poststarthook/priority-and-fairn
ok/priority-and-fairness-filter ok\n[+]poststarthook/start-apiserver-extensions-informers ok\n[+]poststarthook/start-apiextensions-controllers ok\n[+]poststar
[+]poststarthook/bootstrap-controller ok\n[+]poststarthook/rbac/bootstrap-roles ok\n[+]poststarthook/scheduling/bootstrap-system-priority-classes ok\n[+]
fairness-config-producer ok\n[+]poststarthook/start-cluster-authentication-info-controller ok\n[+]poststarthook/agggregator-reload-proxy-client-cert ok\n[+]
agggregator-informers ok\n[+]poststarthook/apiservice-registration-controller ok\n[+]poststarthook/apiservice-status-available-controller ok\n[+]poststar
tration ok\n[+]autoregister-completion ok\n[+]poststarthook/apiservice-openapi-controller ok\n[+]poststarthook/apiservice-openapi3-controller ok\n[+]neal
level=info msg="regenerating all endpoints" reason="one or more identities created or deleted" subsys=endpoint-manager
level=info msg="regenerating all endpoints" reason="one or more identities created or deleted" subsys=endpoint-manager
level=info msg="regenerating all endpoints" reason="one or more identities created or deleted" subsys=endpoint-manager
level=warning msg="Received delete event for key which re-appeared within delay time window" key=def
ndow=30s
level=warning msg="Received delete event for key which re-appeared within delay time window" key=def
ndow=30s
level=warning msg="Received delete event for key which re-appeared within delay time window" key=def
ndow=30s
storeName=store-cilium/state/nodes/
storeName=store-cilium/state/nodes/
storeName=store-cilium/state/nodes/

```

Cilium log

fluent-bit

```

[root@mlaas20webm01p ~]# kubectl -n argocd get po -owide
NAME                                READY   STATUS             RESTARTS   AGE   IP              NODE                                NOMINATED NODE   READINESS GATES
mlaas-argocd-application-controller-0 0/1     CrashLoopBackOff   9 (3m32s ago)  17m   [redacted]      [redacted]                        <none>            <none>
mlaas-argocd-applicationset-controller-7b66475c46-xm9tm 0/1     CrashLoopBackOff   7 (2m57s ago)  17m   [redacted]      [redacted]                        <none>            <none>
mlaas-argocd-notifications-controller-86c89d4d6-bnvk8 1/1     Running            0          17m   [redacted]      [redacted]                        <none>            <none>
mlaas-argocd-redis-bd44f8f56-7mxvf 2/2     Running            0          17m   [redacted]      [redacted]                        <none>            <none>
mlaas-argocd-repo-server-f44f5598c-lmbmd 1/1     Running            0          17m   [redacted]      [redacted]                        <none>            <none>
mlaas-argocd-server-75cc89f64d-bt7h 0/1     CrashLoopBackOff   21 (3m14s ago)  54m   [redacted]      [redacted]                        <none>            <none>

[root@mlaas20webm01p ~]#

```

service recover when pods moved to other nodes

```

NAME                                READY   STATUS             RESTARTS   AGE   IP              NODE                                NOMINATED NODE   READINESS GATES
mlaas-argocd-application-controller-0 1/1     Running            0          113s   [redacted]      [redacted]                        <none>            <none>
mlaas-argocd-applicationset-controller-7b66475c46-tq867 1/1     Running            0          34s    [redacted]      [redacted]                        <none>            <none>
mlaas-argocd-notifications-controller-86c89d4d6-bnvk8 1/1     Running            0          20m    [redacted]      [redacted]                        <none>            <none>
mlaas-argocd-redis-bd44f8f56-7mxvf 2/2     Running            0          20m    [redacted]      [redacted]                        <none>            <none>
mlaas-argocd-repo-server-f44f5598c-lmbmd 1/1     Running            0          20m    [redacted]      [redacted]                        <none>            <none>
mlaas-argocd-server-75cc89f64d-bt7h 1/1     Running            0          34s    [redacted]      [redacted]                        <none>            <none>

[root@mlaas20webm01p ~]#

```

Cilium 只會在一開始修改 iptables

```
ff28c83 cilium / pkg / datapath / iptables / iptables.go
Code Blame 1858 lines (1648 loc) · 58.8 KB
276 type params struct {
288 func newIptablesManager(p params) *Manager {
289     iptMgr := &Manager{
290         logger:    p.Logger,
291         modulesMgr: p.ModulesMgr,
292         cfg:        p.Cfg,
293         sharedCfg:  p.SharedCfg,
294         haveIp6tables: true,
295     }
296
297     p.Lifecycle.Append(iptMgr)
298
299     return iptMgr
300 }
301
302 // Start initializes the iptables manager and checks for iptables kernel modules availability.
303 func (m *Manager) Start(ctx hive.HookContext) error {
304     if os.Getenv("CILIUM_PREPEND_IPTABLES_CHAIN") != "" {
305         m.logger.Warning("CILIUM_PREPEND_IPTABLES_CHAIN env var has been deprecated. Please use
306             "env var or '--prepend-iptables-chains' command line flag instead")
```



• To my surprise, cilium doesn't periodically synchronize those rules like kube-proxy. If you somehow remove a rule in its custom chain, you have to add it back manually or restart cilium-agent. Is this a bug or feature ?

規則：
只會在一開始 onStart, onStop

導致：
流量卡住出不去

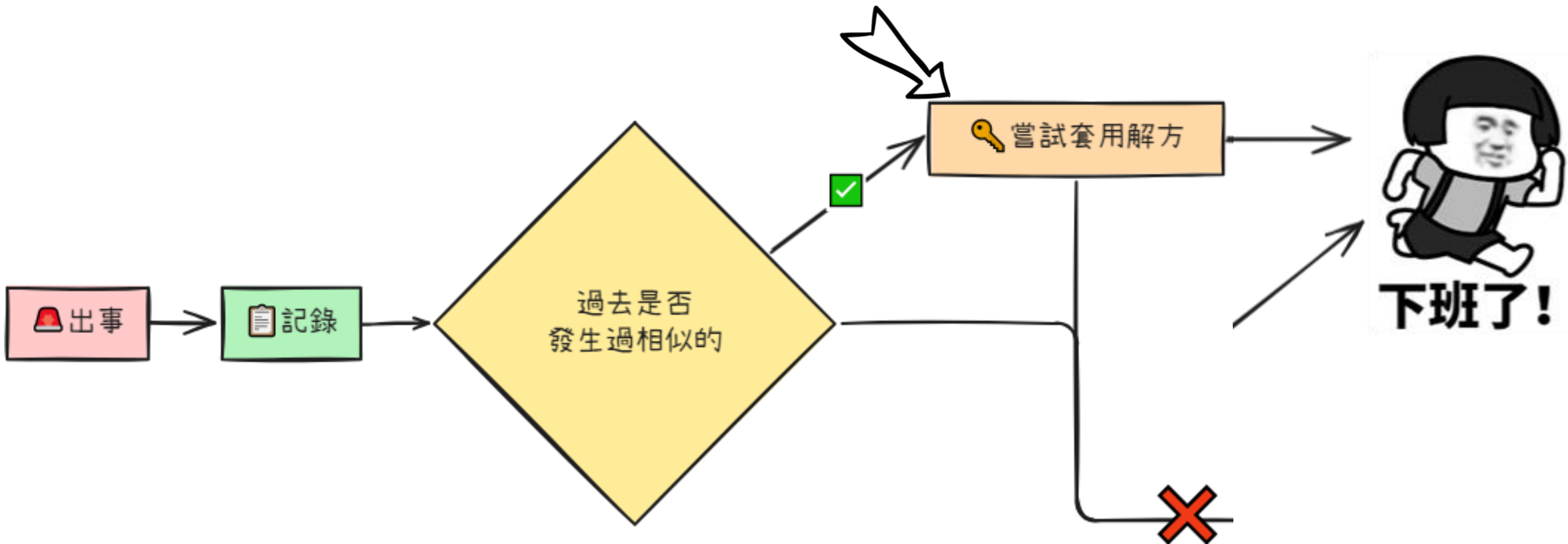


其他苦主

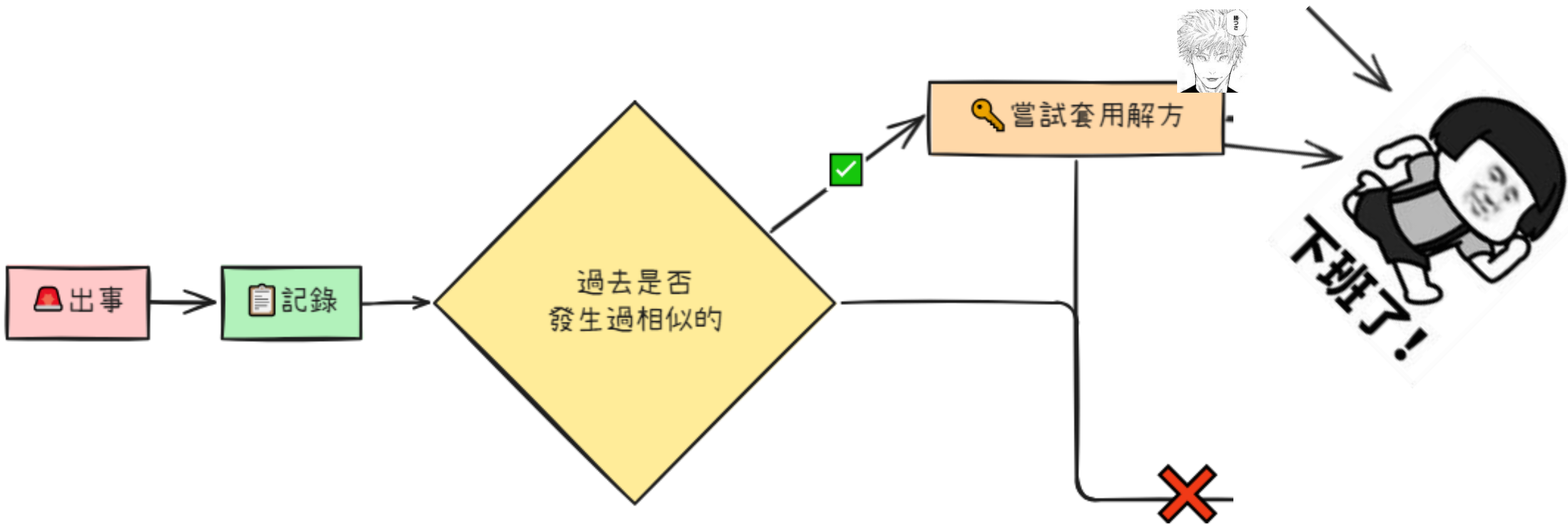
重啟 remote docker daemon 後已修復

結合前面事件的參考

- 知道「重啟與 container network 有關的即可」



記錄的重要性在於：不要考驗自己的記憶力



只是如果

- 解方不管用的話呢？
- 或是沒有可以參考的呢？

另一個事件：response failed



嘗試了重啟 cilium

事件影響範圍：

apiserver, kube-scheduler, cilium-operator, cilium-agent 發生服務重啟

幾乎每個 node 的 cilium pod 都有重啟，其中 更陷入無限的重啟，而 則發生無法新增 pod 的情形
導致當天 airflow 的 dag 無法順利執行。

當時擔心損害擴散，所以將 w04p cordon，拒絕新的 pod 部署

先前另一件事情 202 服務不少 deployment 下陷

在處理完成觀察後，未將 uncordon

因此在 prod cluster 有 的 node 無法部署下，使得 airflow 的 dag，因 node 限制 110 個 pods 運行，kube

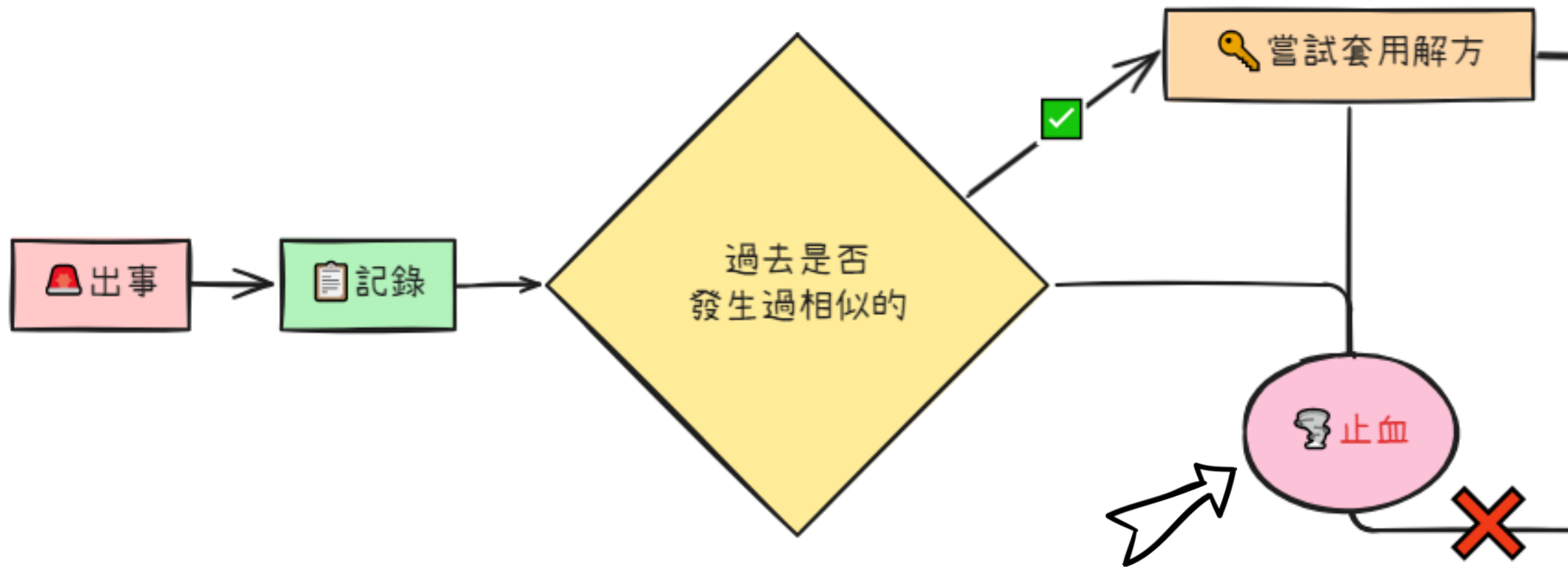
後來有將 uncordon，一切恢復正常

跟隨先前的解法

- 重啟 cilium pod
- 一時無法解除

止血

- 判斷與 node 相關
- 先 cordon 節點



止血

【對外】

- 發公告
- 提出預計修復時間
- 定時回報

各位主管同仁大家好：

[緊急公告] kubernetes 同步作業異常處理。

維護範圍：esunisgood.com 相關的服務

維護起始時間：2024/10/24(四) 18:45

預計結束時間：2024/10/24(四) 20:45

維護說明：

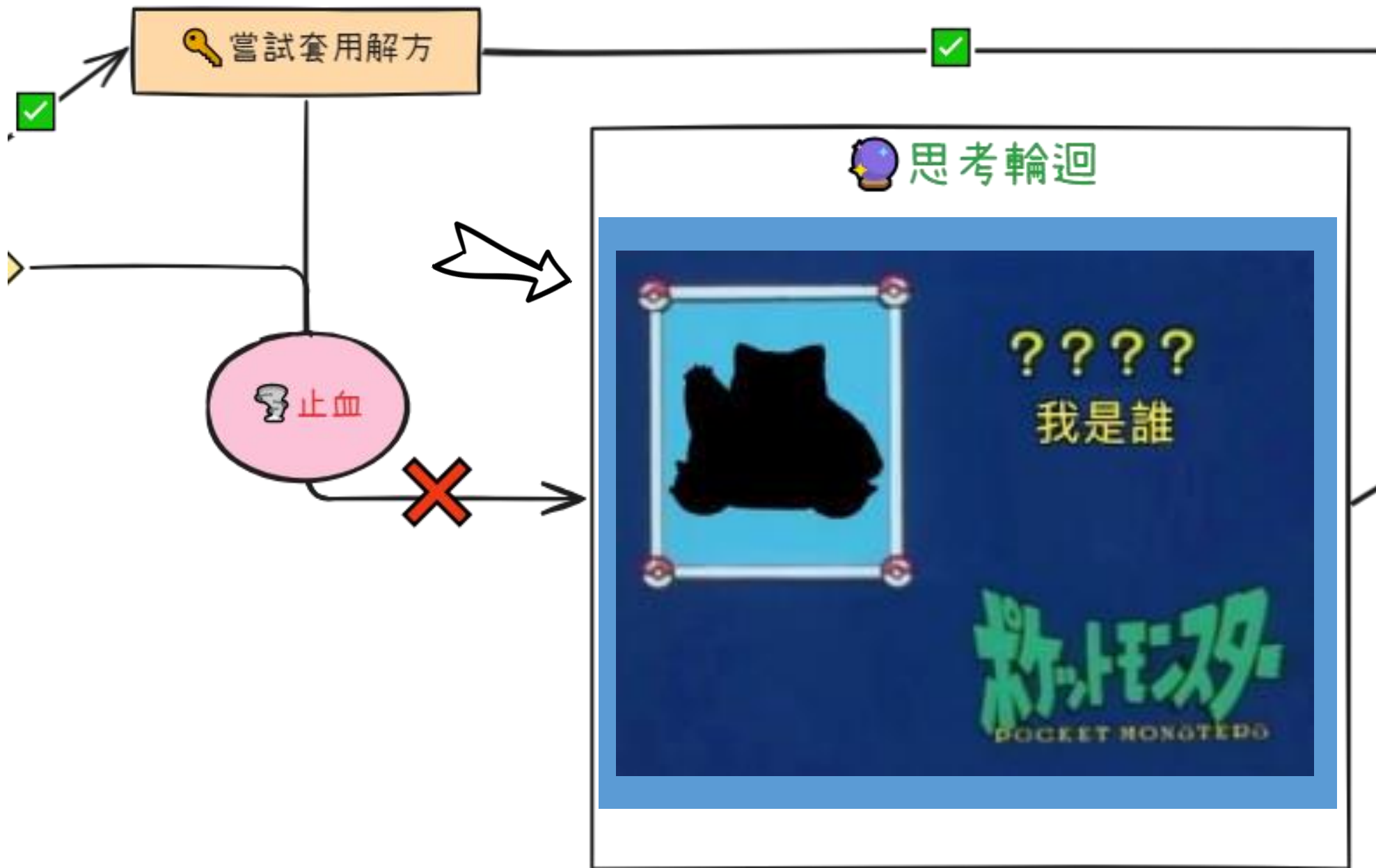
1. 目前 k8s hello world 這台因同步異常，需暫停與 k8s hi 同步後恢復服務
2. 作業期間若有使用 k8s hello world 的請暫時改用 esunistmp.com

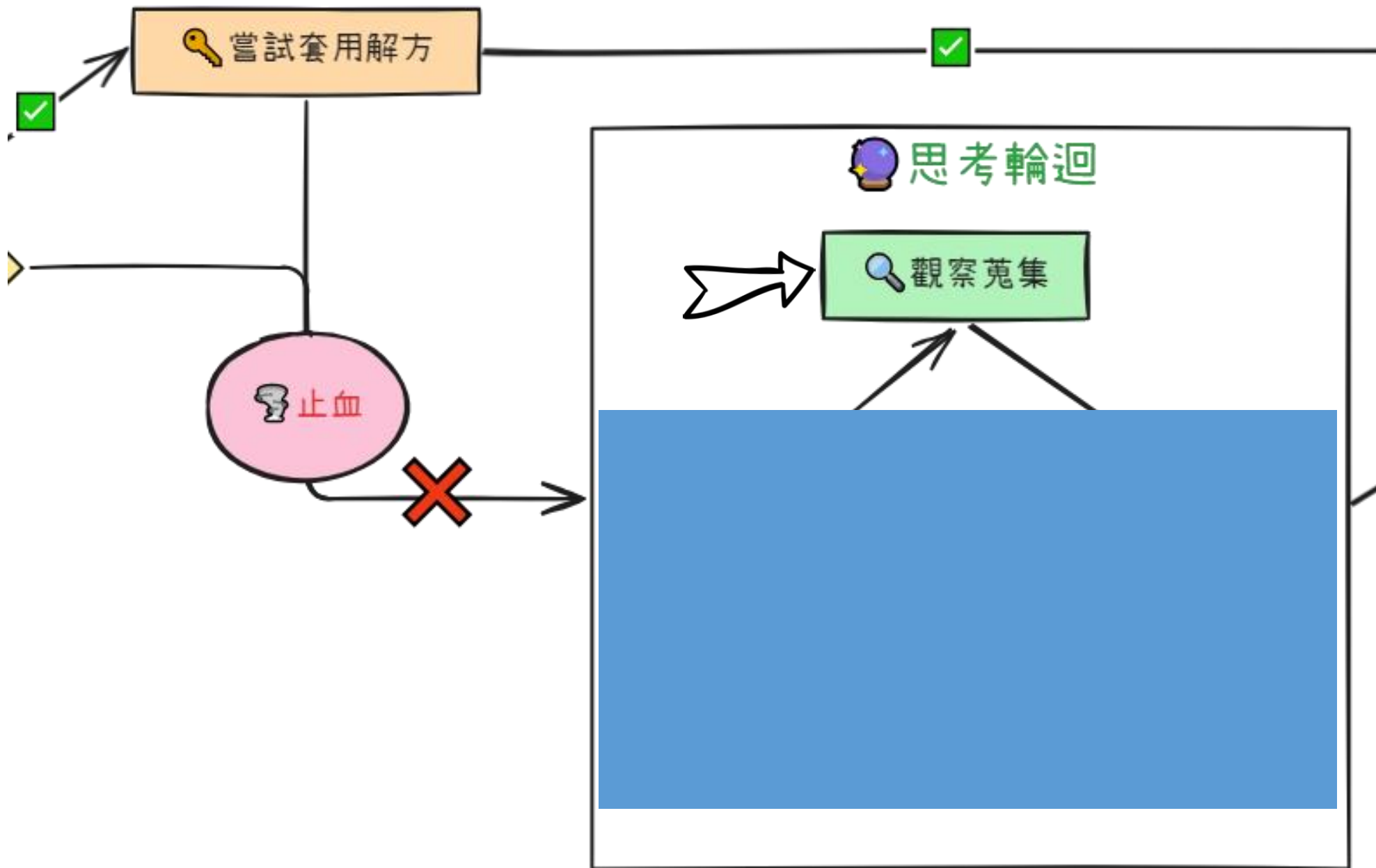
維護負責人：kiwi

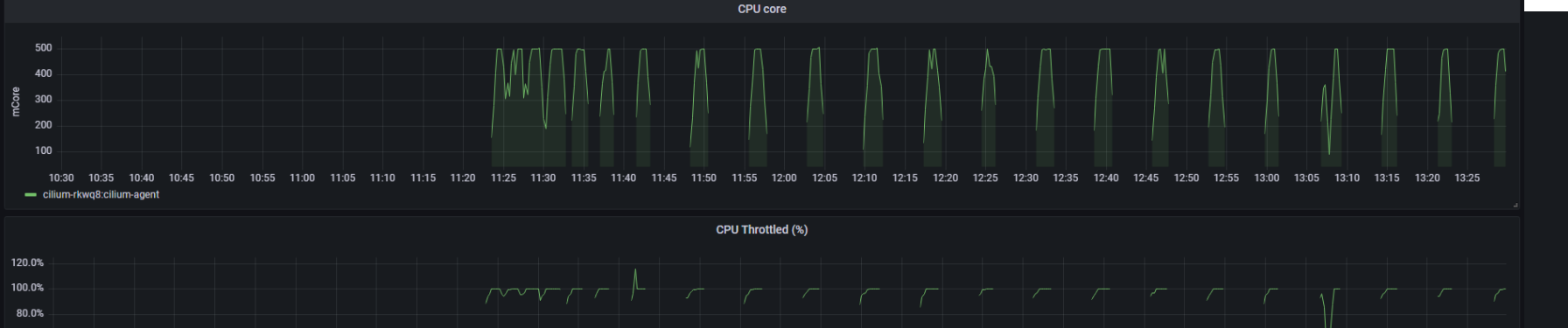
謝謝

【對內】

- 看要停止哪些服務
 - pod 驅離
 - node cordon
 - load balancer 直接導轉
 - 啟動 HA 機制
- 相信自己的判斷
 - 錯了趕快修正







node 使用量

```

level=error msg="Command execution failed" cmd="[tc filter replace dev lxc7bd94e69e05d ingress prio 1 handle 1 bpf da obj 1100_next/bpf_lxc.o sec from-container]" error="signal: killed" subsystem=datapath-load
level=error msg="Command execution failed" cmd="[tc filter replace dev lxc3b47842030f ingress prio 1 handle 1 bpf da obj 3376_next/bpf_lxc.o sec from-container]" error="signal: killed" subsystem=datapath-load
level=error msg="Command execution failed" cmd="[tc filter replace dev lxc8c5595eafad ingress prio 1 handle 1 bpf da obj 1448_next/bpf_lxc.o sec from-container]" error="signal: killed" subsystem=datapath-load
level=error msg="Command execution failed" cmd="[tc filter replace dev lxc8e6a056eaf3f ingress prio 1 handle 1 bpf da obj 2340_next/bpf_lxc.o sec from-container]" error="signal: killed" subsystem=datapath-load
level=error msg="Command execution failed" cmd="[tc filter replace dev lxc477270ade973 ingress prio 1 handle 1 bpf da obj 1965_next/bpf_lxc.o sec from-container]" error="signal: killed" subsystem=datapath-load
level=error msg="Command execution failed" cmd="[tc filter replace dev lxc4c733564269 ingress prio 1 handle 1 bpf da obj 2730_next/bpf_lxc.o sec from-container]" error="signal: killed" subsystem=datapath-load
level=error msg="Command execution failed" cmd="[tc filter replace dev lxc9492526cc87 ingress prio 1 handle 1 bpf da obj 3885_next/bpf_lxc.o sec from-container]" error="signal: killed" subsystem=datapath-load
level=error msg="Command execution failed" cmd="[tc filter replace dev lxcfaaf789c6fae ingress prio 1 handle 1 bpf da obj 214_next/bpf_lxc.o sec from-container]" error="signal: killed" subsystem=datapath-load
level=error msg="Command execution failed" cmd="[tc filter replace dev lxc0115cb12ae45 ingress prio 1 handle 1 bpf da obj 360_next/bpf_lxc.o sec from-container]" error="signal: killed" subsystem=datapath-load
level=error msg="Command execution failed" cmd="[tc filter replace dev lxc112b87efe47 ingress prio 1 handle 1 bpf da obj 875_next/bpf_lxc.o sec from-container]" error="signal: killed" subsystem=datapath-load
level=error msg="Command execution failed" cmd="[tc filter replace dev lxcd24f88351c7 ingress prio 1 handle 1 bpf da obj 3834_next/bpf_lxc.o sec from-container]" error="signal: killed" subsystem=datapath-load

```

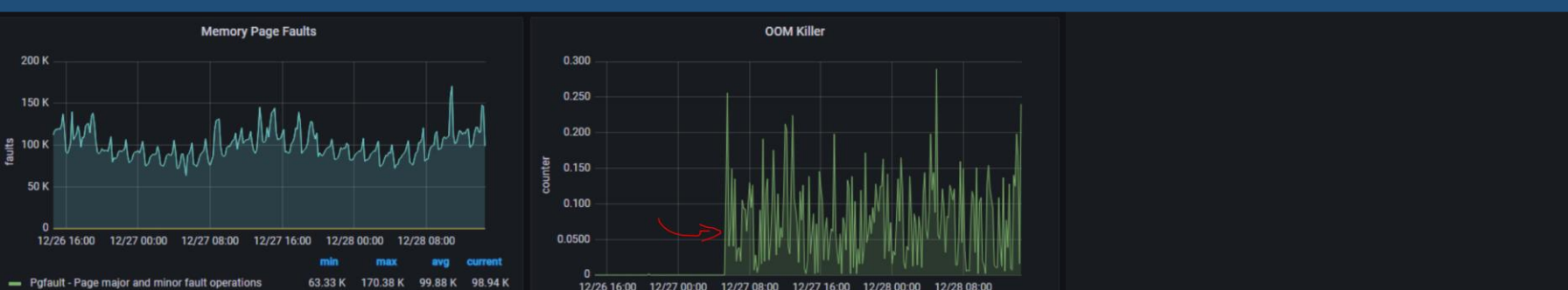
各個地方的 log (假設 monitor infra 沒建好)

```

128771: lxc8506afda8a6@if128770: <BROADCAST,MULTICAST,UP,LOWER_UP> mtu 1500 qdisc noqueue state UP mode DEFAULT group default qlen 1000
128773: lxc7e870a7c084@if128772: <BROADCAST,MULTICAST,UP,LOWER_UP> mtu 1500 qdisc noqueue state UP mode DEFAULT group default qlen 1000
128775: lxc1e17e513782c@if128774: <BROADCAST,MULTICAST,UP,LOWER_UP> mtu 1500 qdisc noqueue state UP mode DEFAULT group default qlen 1000
128777: lxc86b0d58e2908@if128776: <BROADCAST,MULTICAST,UP,LOWER_UP> mtu 1500 qdisc noqueue state UP mode DEFAULT group default qlen 1000
99087: lxc6e04432c37ef@if99086: <BROADCAST,MULTICAST,UP,LOWER_UP> mtu 1500 qdisc noqueue state UP mode DEFAULT group default qlen 1000
120081: lxc4c635f66abe6@if120080: <BROADCAST,MULTICAST,UP,LOWER_UP> mtu 1500 qdisc noqueue state UP mode DEFAULT group default qlen 1000
99089: lxc15339e9e44b8@if99088: <BROADCAST,MULTICAST,UP,LOWER_UP> mtu 1500 qdisc noqueue state UP mode DEFAULT group default qlen 1000
17: lxc3183896e4c34@if16: <BROADCAST,MULTICAST,UP,LOWER_UP> mtu 1500 qdisc noqueue state UP mode DEFAULT group default qlen 1000
120083: lxc29620ebf6c1@if120082: <BROADCAST,MULTICAST,UP,LOWER_UP> mtu 1500 qdisc noqueue state UP mode DEFAULT group default qlen 1000
117523: lxc3a7b34ab498e@if117522: <BROADCAST,MULTICAST,UP,LOWER_UP> mtu 1500 qdisc noqueue state UP mode DEFAULT group default qlen 1000
99091: lxc8f8e0bb396e@if99090: <BROADCAST,MULTICAST,UP,LOWER_UP> mtu 1500 qdisc noqueue state UP mode DEFAULT group default qlen 1000
99093: lxcd71eae72366a@if99092: <BROADCAST,MULTICAST,UP,LOWER_UP> mtu 1500 qdisc noqueue state UP mode DEFAULT group default qlen 1000
117527: lxc3a64084f63f@if117526: <BROADCAST,MULTICAST,UP,LOWER_UP> mtu 1500 qdisc noqueue state UP mode DEFAULT group default qlen 1000
99095: lxc2e280670a16a@if99094: <BROADCAST,MULTICAST,UP,LOWER_UP> mtu 1500 qdisc noqueue state UP mode DEFAULT group default qlen 1000

```

Node 的 state (iptabels,iplink,netfilter,route...)



系統元件的 metrics (cilium,apiserver,etcd,cadvisor...)

找到是關於 cilium state 的儲存機制

```
[root@m1aas20itcm0lp state]# kubectl -n kube-system get po cilium-l57xx -oyaml | grep hostPath -A 10
- hostPath:
  path: /var/run/cilium
  type: DirectoryOrCreate
  name: cilium-run
- hostPath:
  path: /sys/fs/bpf
  type: DirectoryOrCreate
  name: bpf-maps
- hostPath:
  path: /proc
  type: Directory
  name: hostproc
- hostPath:
  path: /run/cilium/cgroupv2
  type: DirectoryOrCreate
  name: cilium-cgroup
- hostPath:
  path: /opt/cni/bin
  type: DirectoryOrCreate
  name: cni-path
- hostPath:
  path: /etc/cni/net.d
  type: DirectoryOrCreate
  name: etc-cni-netd
- hostPath:
  path: /lib/modules
  type: ""
  name: lib-modules
- hostPath:
  path: /run/xtables.lock
  type: FileOrCreate
  name: xtables-lock
```

hostPath 儲存

state 狀態

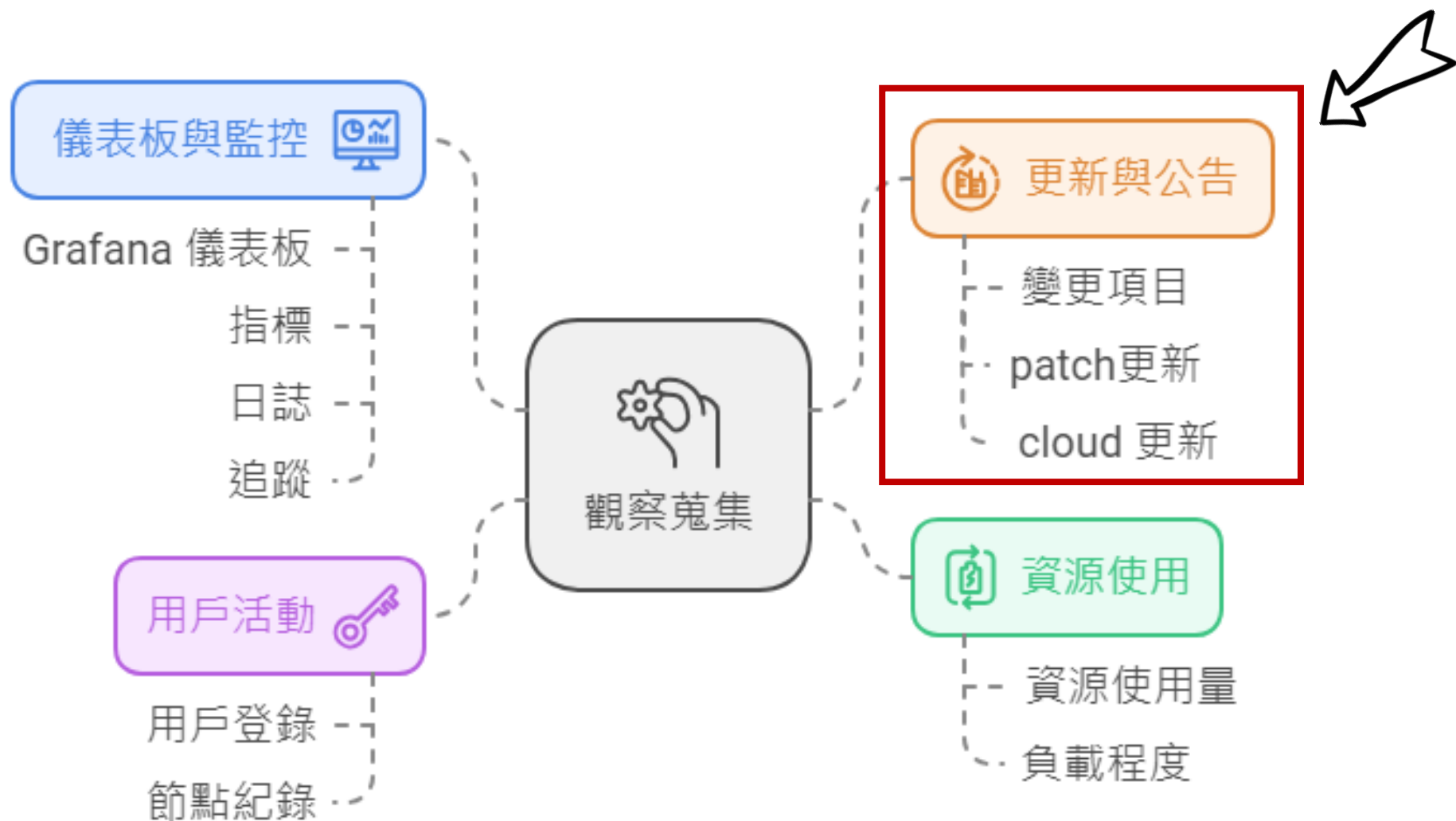
-> 重啟 pod 沒有用

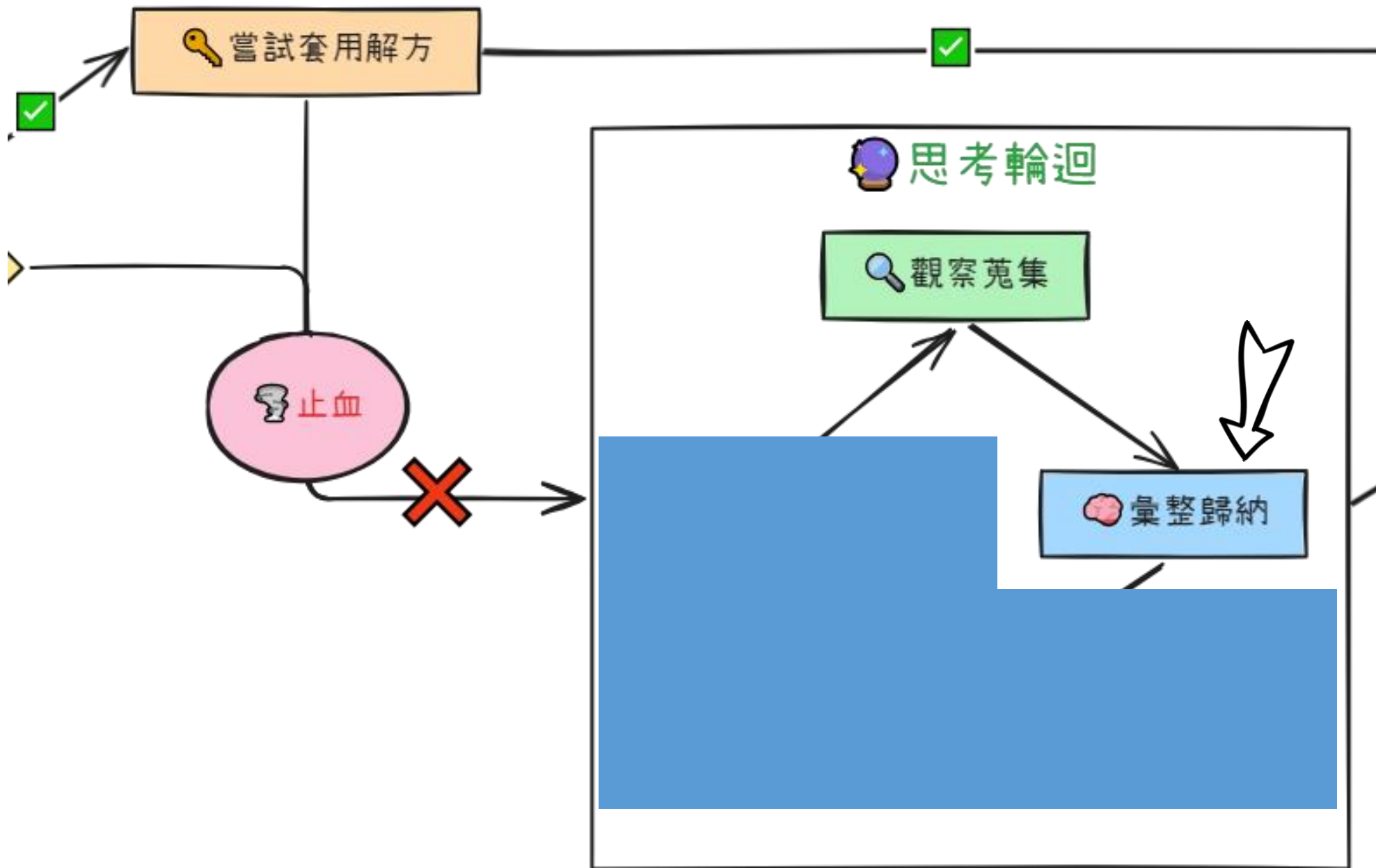


解法

1. 在 pod 重啟前
清掉 state
2. 重開機

觀察蒐集

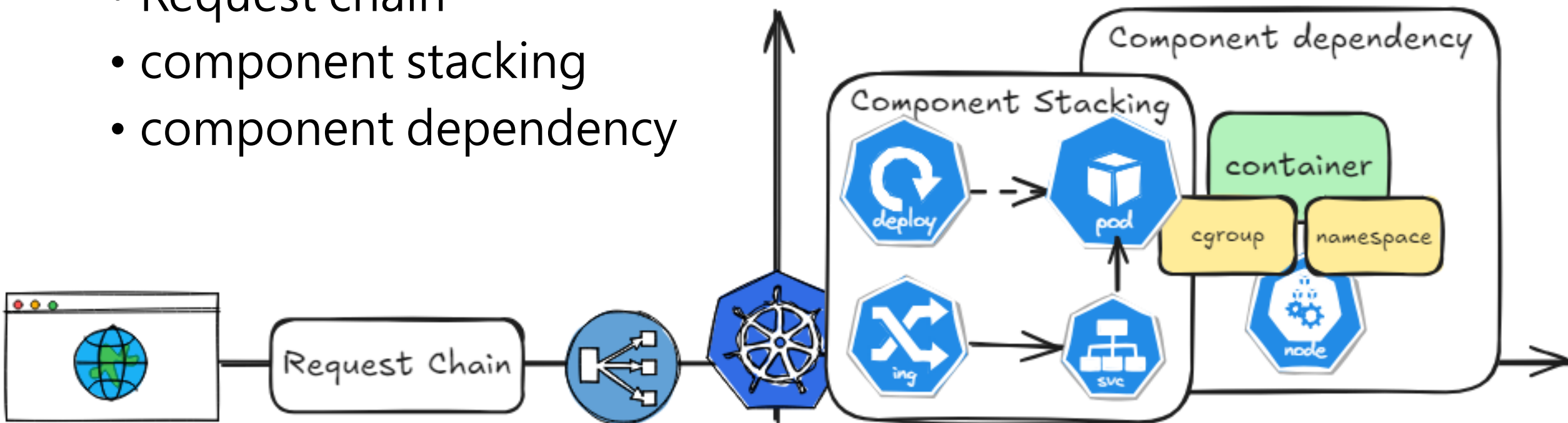




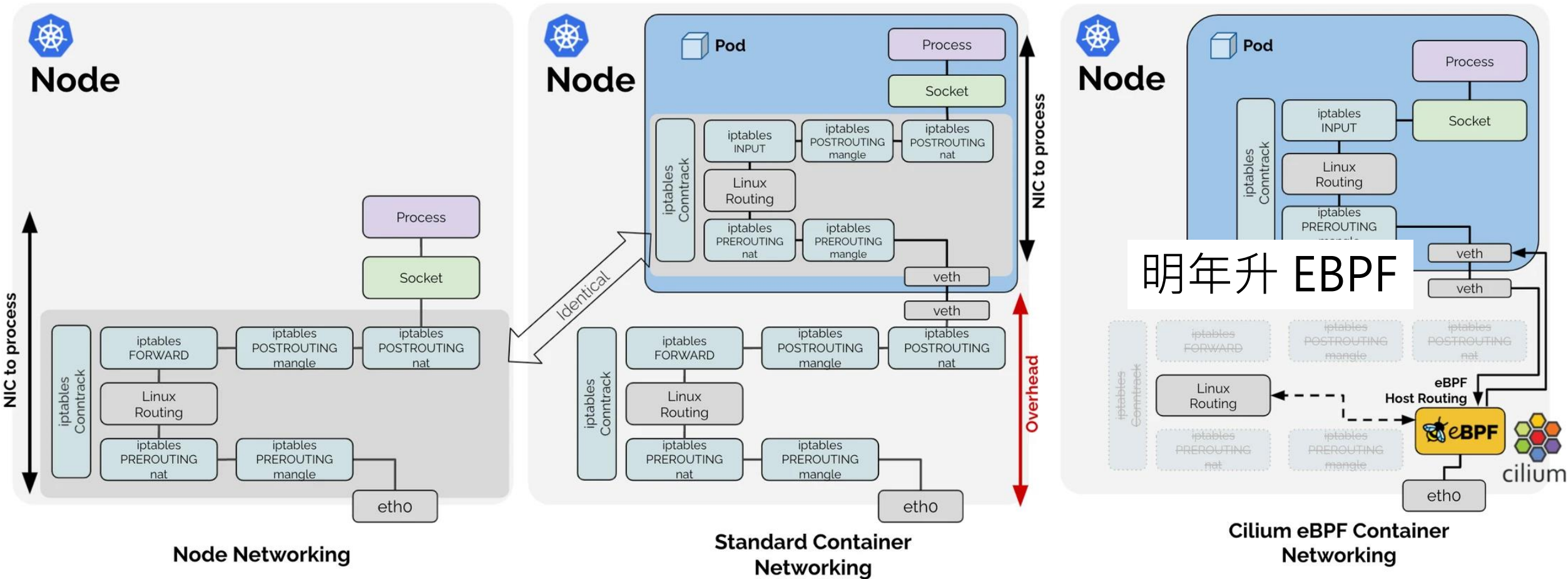
彙整歸納

釐清災難範圍



- Request chain
- component stacking
- component dependency

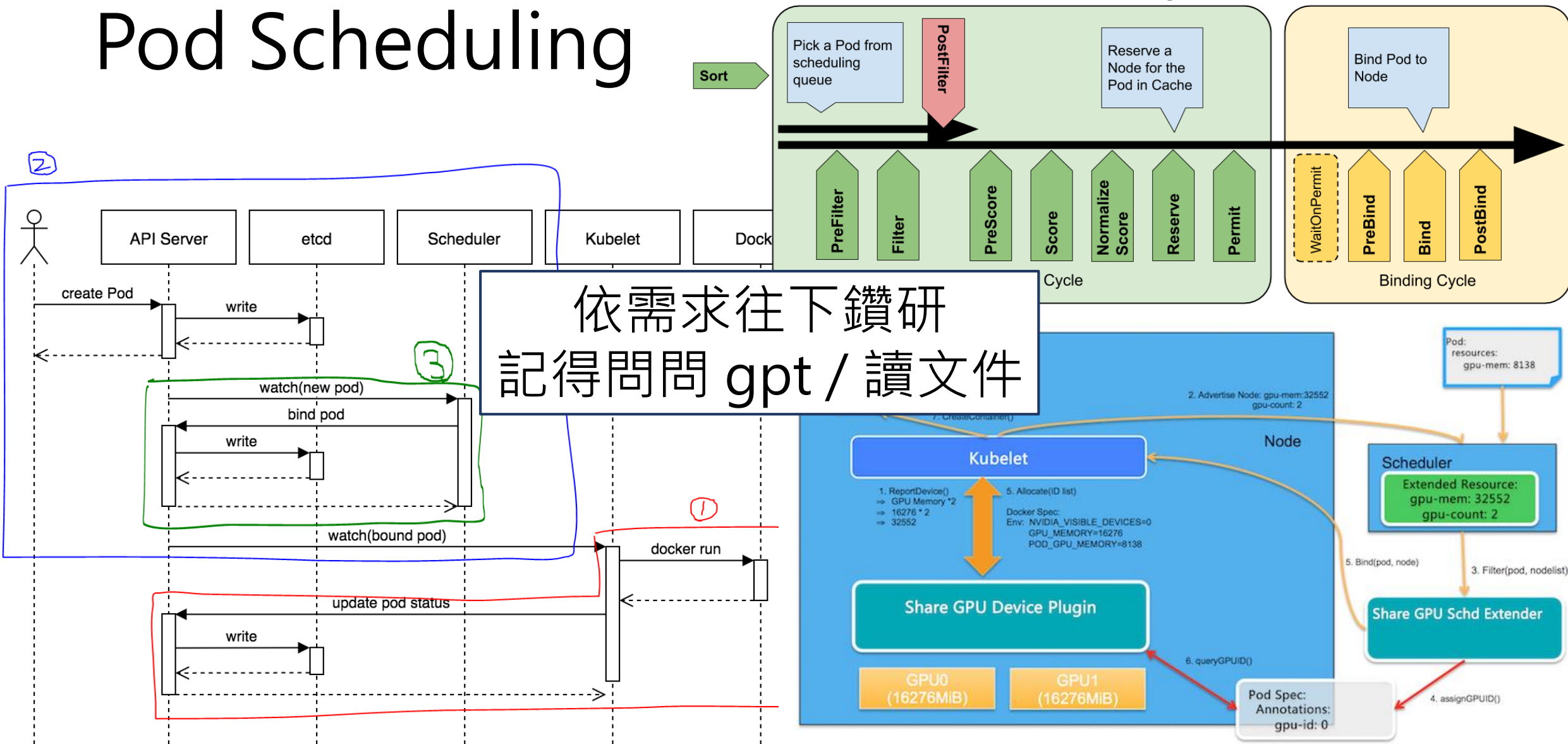


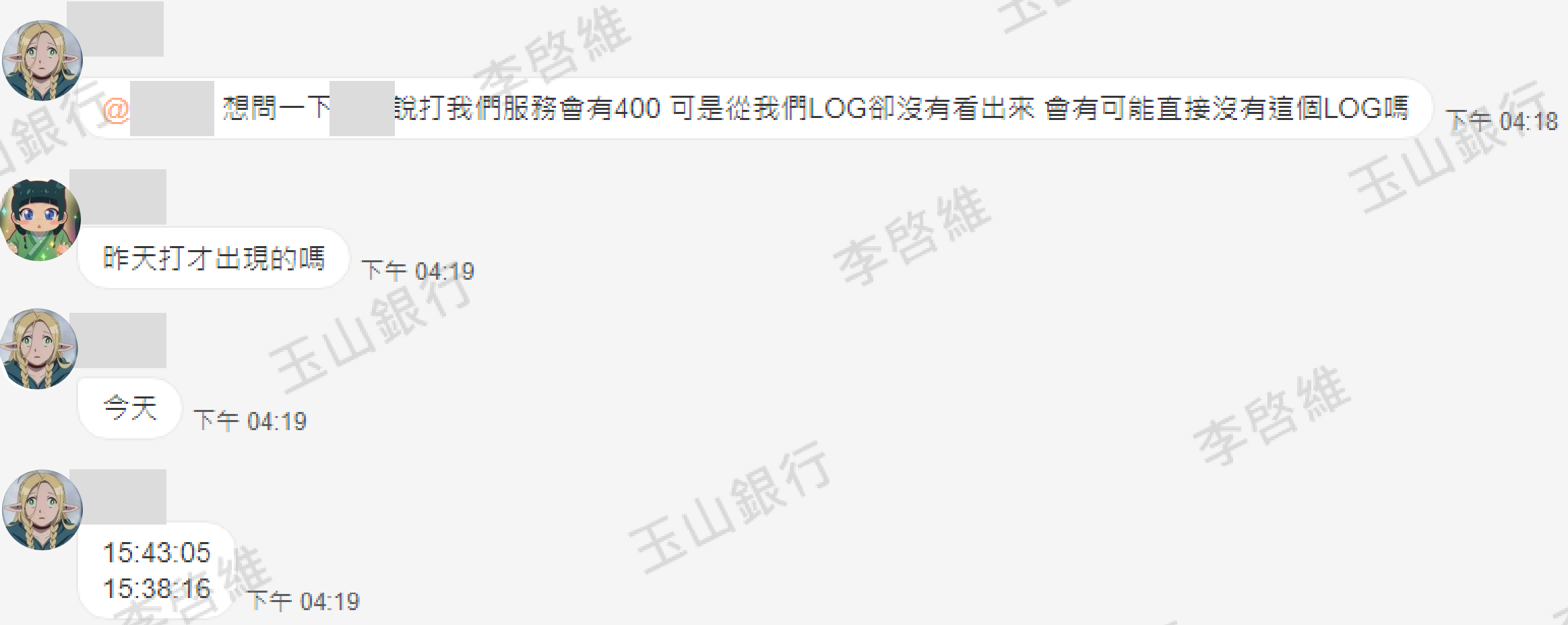
Cilium CNI 解釋



Pod Scheduling

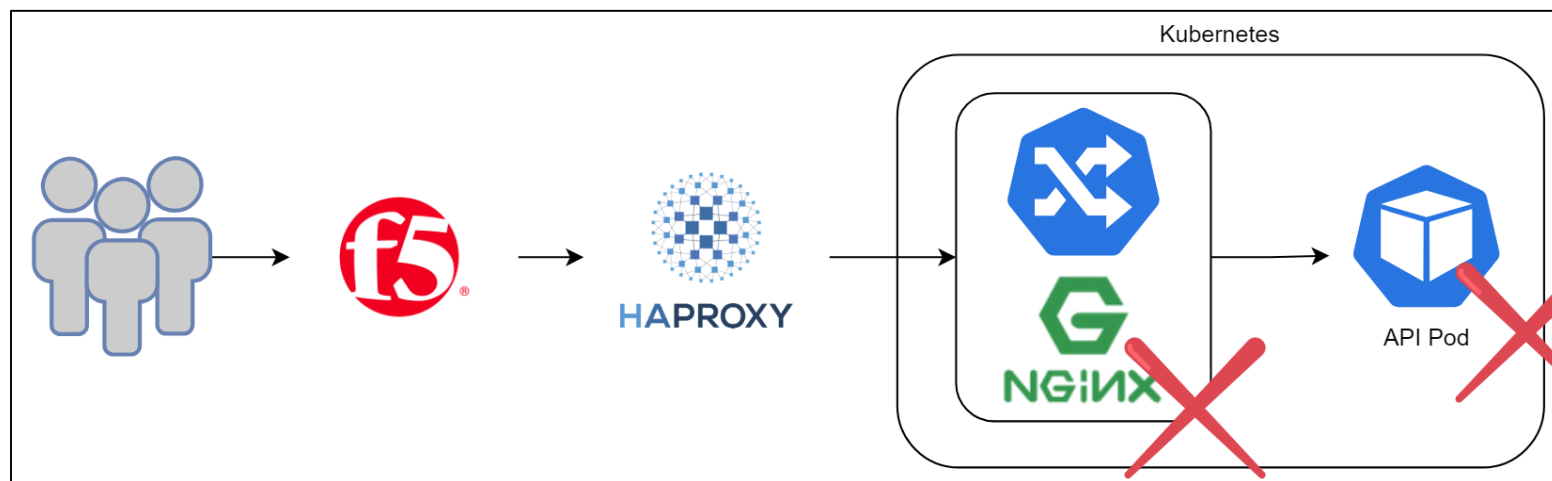
 Extensible API
 Internal API





專案團隊說：「response 400」

Request Chain 解決

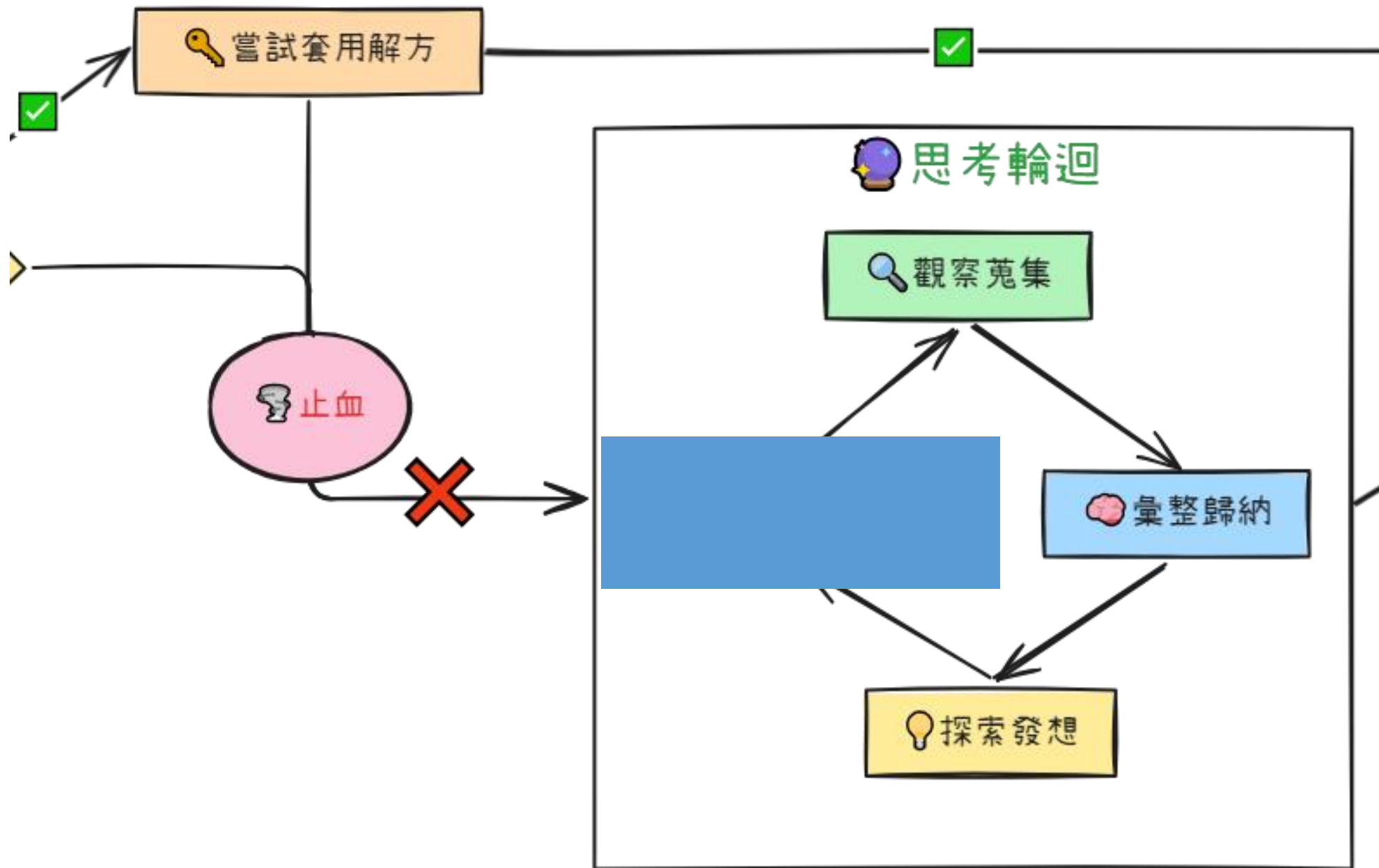


```
11 localhost haproxy | mlaas mlaas/<NOSRV> -1/-1/-1/-1/0 400 187 - - PR-- 42/42/0/0/0 0/0 {} " <BADREQ> "
14 localhost haproxy | mlaas mlaas/<NOSRV> -1/-1/-1/-1/0 400 187 - - PR-- 37/36/0/0/0 0/0 {} " <BADREQ> "
16 localhost haproxy | mlaas mlaas/<NOSRV> -1/-1/-1/-1/0 400 187 - - PR-- 24/23/0/0/0 0/0 {} " <BADREQ> "
17 localhost haproxy | mlaas mlaas/<NOSRV> -1/-1/-1/-1/0 400 187 - - PR-- 19/18/0/0/0 0/0 {} " <BADREQ> "
18 localhost haproxy | mlaas mlaas/<NOSRV> -1/-1/-1/-1/0 400 187 - - PR-- 17/16/0/0/0 0/0 {} " <BADREQ> "
19 localhost haproxy | mlaas mlaas/<NOSRV> -1/-1/-1/-1/0 400 187 - - PR-- 17/16/0/0/0 0/0 {} " <BADREQ> "
19 localhost haproxy | mlaas mlaas/<NOSRV> -1/-1/-1/-1/0 400 187 - - PR-- 13/12/0/0/0 0/0 {} " <BADREQ> "
20 localhost haproxy | mlaas mlaas/<NOSRV> -1/-1/-1/-1/0 400 187 - - PR-- 17/16/0/0/0 0/0 {} " <BADREQ> "
```

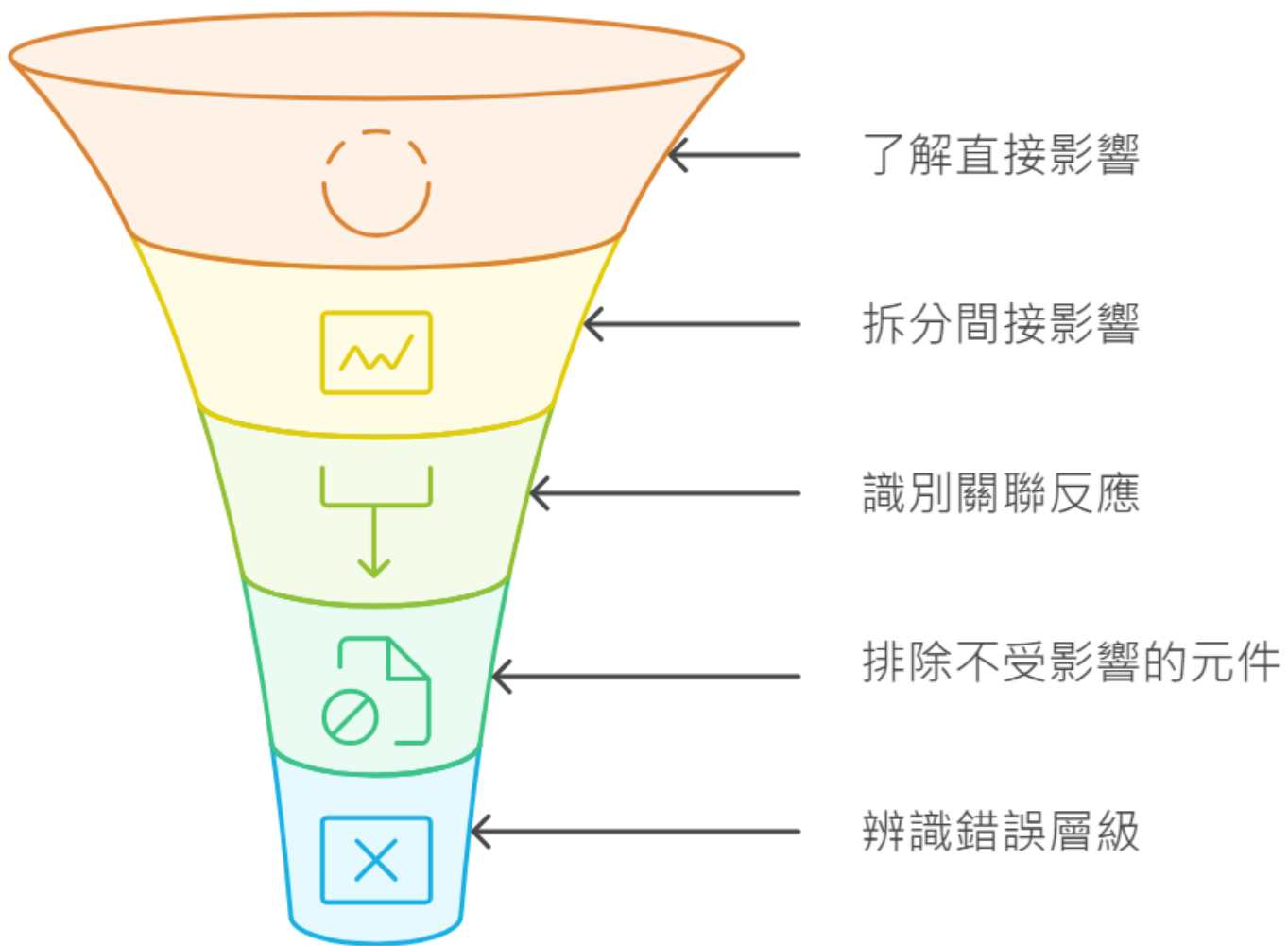


客戶端呼叫錯誤 😊

前回應 400 bad request , MLaaS 有收到 HTTP header 但沒收到 HTTP content body
- 回應 400 bad request , MLaaS 只有 TCP 連線成功 , 沒有收到 HTTP header & content body



探索發想



問問自己

錯誤是如何擴散呢？

是否集中在某個服務/任務呼叫？
會不會是來源端就出錯了？

ex. cluster, node, pod, container

自己弄得錯誤

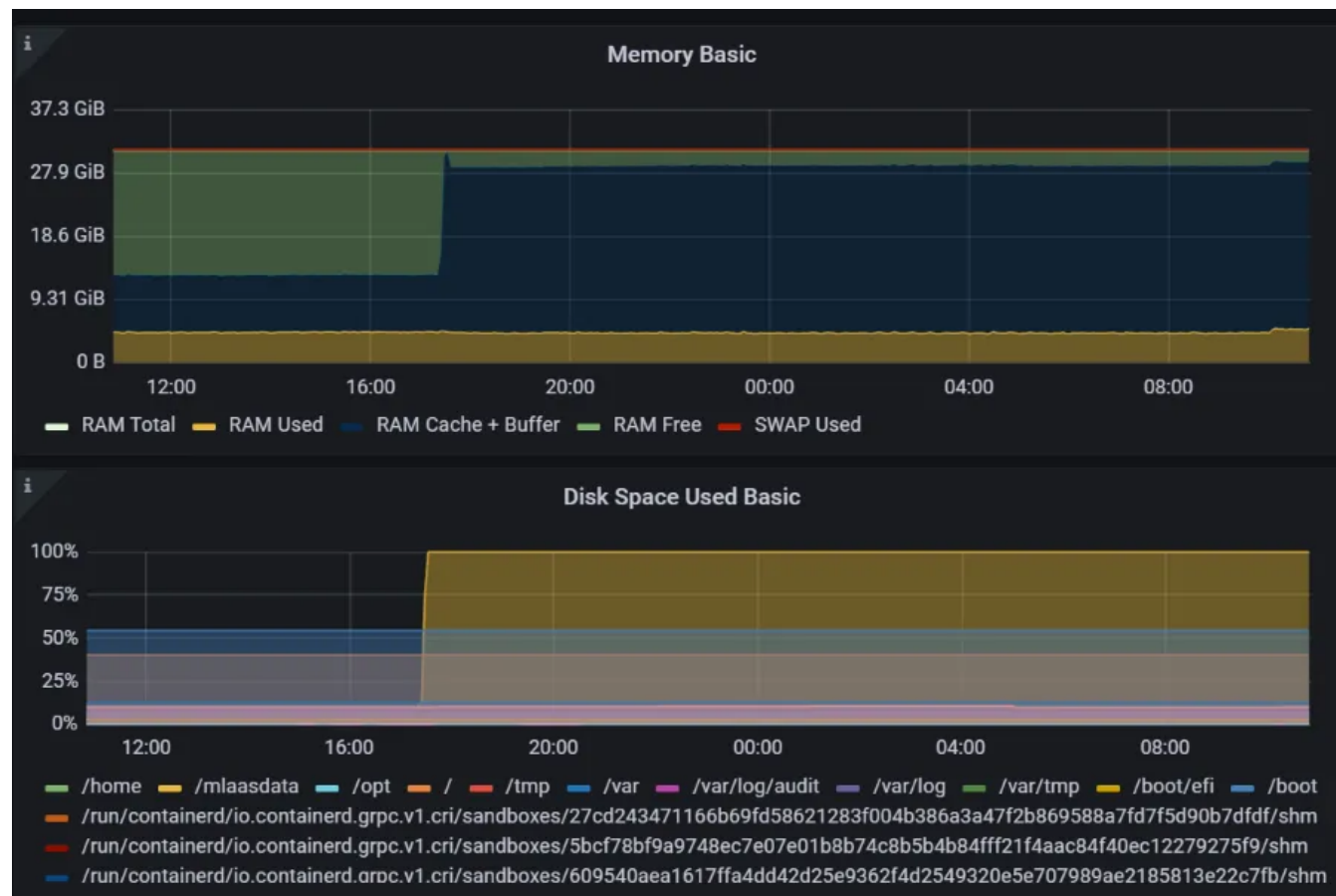
Actions

```
:22:34| nerdctl pull [redacted]  
:31:41| nerdctl run --entrypoint bash -it [redacted]  
:31:58| du -h  
:32:01| df -h  
:32:14| nerdctl run --entrypoint bash -it [redacted]  
:32:27| nerdctl images  
:32:34| nerdctl image prune  
:32:38| nerdctl image prune --all  
:33:59| kubectl get nodes -owide
```

緊急處置後

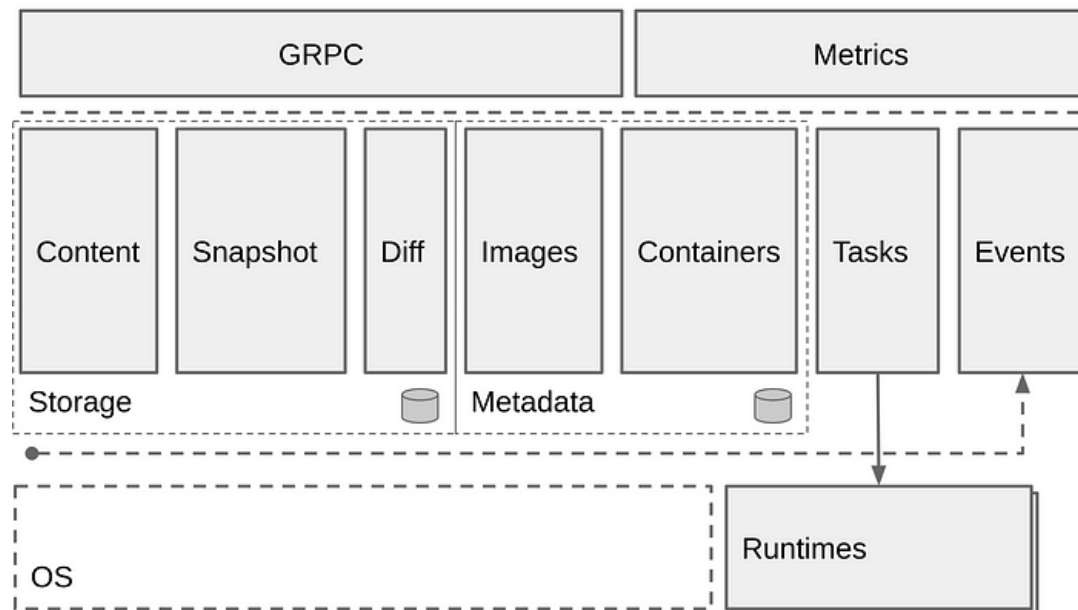
- ctr images 沒看到
- ctr container 沒看到
resources 仍然沒降下來

Result



彙整歸納後，嘗試解看看

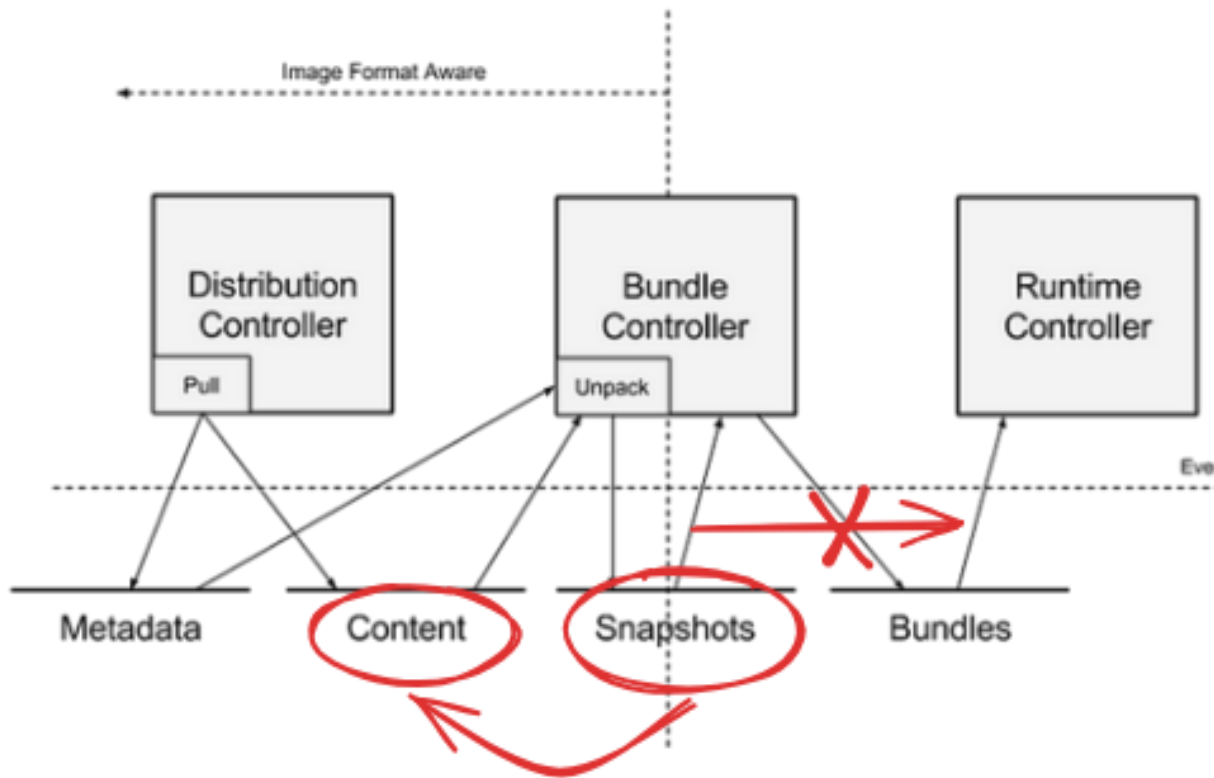
```
[root@████████ containerd]# du -h -d 1 ██████████ containerd/
3.7G ██████████ containerd/
1.1G ██████████ containerd/io.containerd.content.v1.content
72K ██████████ containerd/io.containerd.snapshotter.v1.overlayfs
4.0K ██████████ containerd/io.containerd.runtime.v2.task
4.0K ██████████ containerd/tmpmounts
8.0K ██████████ containerd/io.containerd.runtime.v1.linux
268K ██████████ containerd/io.containerd.grpc.v1.introspection
91G ██████████ containerd/io.containerd.grpc.v1.cri
2.7M ██████████ containerd/io.containerd.snapshotter.v1.native
4.0K ██████████ containerd/io.containerd.metadata.v1.bolt
96G ██████████ containerd/io.containerd.snapshotter.v1.btrfs
██████████ containerd/
```



```
sha256:308ef665d277cad0b42fb039b1a4d1c702228646a126ddb4da3062a2ddb7c9cb
sha256:55608c1965eb0f277f8572e0cc953d1f8a34e80965711f5fcd888deeb067ff50
sha256:717f9a2779ca8929b3100ecaa9506a1d23c5d15c4b9f58d57b3c1e665b3f4ffc
├─ extract-536596800 -fPin sha256:3450ad4047717bef9f23b8a873b31cbeedcebe91c39b4fea143d3c72bd8ed908
├─ extract-748197178 -ByGC sha256:3450ad4047717bef9f23b8a873b31cbeedcebe91c39b4fea143d3c72bd8ed908
└─ sha256:3450ad4047717bef9f23b8a873b31cbeedcebe91c39b4fea143d3c72bd8ed908
sha256:5b1fa8e3e100361047c8bcd5553ab6329b9c713c1d4eb87a646760329cea5b3a
├─ sha256:f18795e2ae0e0fba2c4fad9a76f7d2432e3c4159dbf3722bf7194daf19b55b80
sha256:87b6a930c8d02125eedb40c3bc2801f3cfd99fcc8c05134a33b0cc74a2fa6b45
└─ sha256:8cc32bcda2b1d188e9e3a931e4d39ac3bde8e2337d8cc7b700af954f7b63e71e
[root@████████ ~]# ctr -n k8s.io snapshot --snapshotter rm sha256:3450ad4047717bef9f23b8a873b31cbeedcebe91c39b4fea143d3c72bd8ed908
No help topic for 'sha256:3450ad4047717bef9f23b8a873b31cbeedcebe91c39b4fea143d3c72bd8ed908'
[root@████████ ~]# ctr -n k8s.io snapshot --snapshotter rm sha256:717f9a2779ca8929b3100ecaa9506a1d23c5d15c4b9f58d57b3c1e665b3f4ffc
No help topic for 'sha256:717f9a2779ca8929b3100ecaa9506a1d23c5d15c4b9f58d57b3c1e665b3f4ffc'
```



探索發想 - 從 mount 繞圈



卡了一段時間

- 在 snapshot, namespace
- 挖到 overlays, mount 的關聯

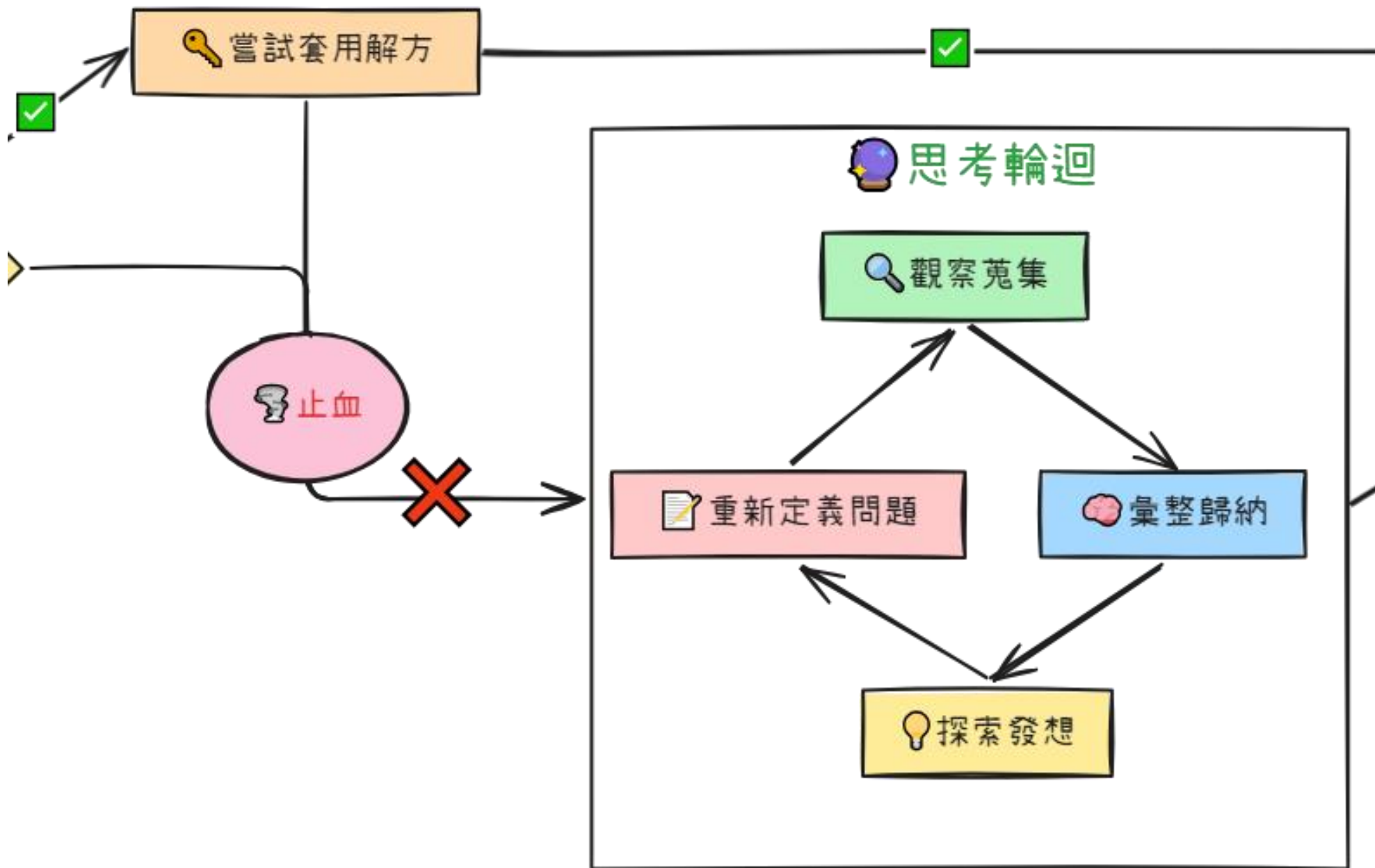
然後意外發現

- content 還沒檢查過
- 之前只有到 image


```
[root@██████████ sha256]# ctr -n k8s.io content ls | grep ██████████ | awk '{print $1}' | xargs ctr -n k8s.io content rm
sha256:093f9cb4ed5900be5b069aa140bee4a6533cc0734efe69b68011cb00bc533761
sha256:0c177adbc1a6746a3e0a7fe09576028f3e2522a58b4fd15305d2f1ab6411c519
sha256:1841a82df382ec6f438928bdb39f5e40378de67b1ac84e1424e4b65f818d4f21
sha256:382d92eb8ea3e46c02c6c3f64e34d0564a425211b14730e99b13c4c172465019
sha256:3f87ccbacbd7fa17cc187efd79dfc2759caafc1e4a7657714bc342add5f9d30
sha256:5cf4eab5d1dafcdd75784e4c9f4ab10d43bda0b3d4397d8a492a31e549b59bfd
sha256:64266c6e10350c49d372bb295f5845a9c2d04b2484841363ede54998d19fb193
sha256:68ef693ccc58a8aa9c387ac7658ae8f56d3f1db06f8e8fad3dca42903c9c4653
sha256:7e77df6a2c111332ec2ad5e83ae851e004544a3933cac31a2988578be4362586
sha256:806911a0f330b943ce4398bfc852cd3a93559e91f9a505cd6baa50242a6cfe33
sha256:852236cde8166ae74149b017340dfcf02d57f02b3e5adbbbbaa8383ff069f935e
sha256:924a30603b3c77cd2dc59cf9996b4f294d8ebb7f22b50411d018559132223b79
sha256:9295acad8e8d5fe32390043ab325f91d11a1ddc6c76a1ef5ab854095b45160d0
sha256:9e08999d446baf672263edd667fd0b50cb1c40303530deb803deb0671ceaaf4a
sha256:a2c53e669ed826a359f466e4de5a31619d5efb411f32b5b004bffe9fac305469
sha256:a542df56279ad54cbc6d09c07b0e0549f827d3e9e06f0563292b91a9f8defd39
sha256:a775e38b5f142ba50f73121074d9208f4a7fe5f7d6b35cdf7a1e3ad484664740
sha256:a7b60a932d48762a6a1f3ac8b82f0ca9cd3f9accec09b84f2fc9fb25a4babae79
sha256:a812797e651bb55279090fc711aee33a8ce66905018212c7cf4c0e531007af40
sha256:ab353a332cea8162589830e13948955024e24763883826db65ca28242c36e272
sha256:b07590bc38a13cfcacd7d3ae30925b9b0926959ecd268bcf548dbbef59a76ae9
sha256:ba7e7656f6366696323d20f4de10f052cd938b9038574a8f0c9ce8674fb01e7e
sha256:dc9affc4ae35e5b2340249841e0ebb584a905c4fedda1fd9e195a47f333470a1
sha256:dcc2a1c0b4a6d4453c477c65a05d2a40cb226a1c4a1e34fd4fd50a12eb1fa280
sha256:ee706bdc96476b63d0ff21bdbaf77a33184e33f474000cba5a2e
sha256:f775f244232d4c19f65c76d713ac7df9dc5f681e7eeff37fa274
sha256:fb668870d8a72b5d72a3b6d98ee626e00f9f7c29c6f4f7d3a636
```

DONE!!!

```
[root@██████████ snapshots]# du -h -d 0 ██████████ containerd/io.containerd.snapshotter.v1.native/snapshots/
2.9G ██████████ containerd/io.containerd.snapshotter.v1.native/snapshots/
[root@██████████ snapshots]# du -h -d 0 ██████████ containerd/io.containerd.content.v1.content/blobs/sha256/
384M ██████████ containerd/io.containerd.content.v1.content/blobs/sha256/
```



重新定義

- **保持冷靜 & 維持好心情**
 - 承認自己知識的不足，與他人分享可能會獲得靈感
- **重新聚焦**
 - 根據前面的資料，重新說明問題的範圍和特徵。
- **持續溝通回報**
 - 尋求外部協助，及系統上可能的取捨





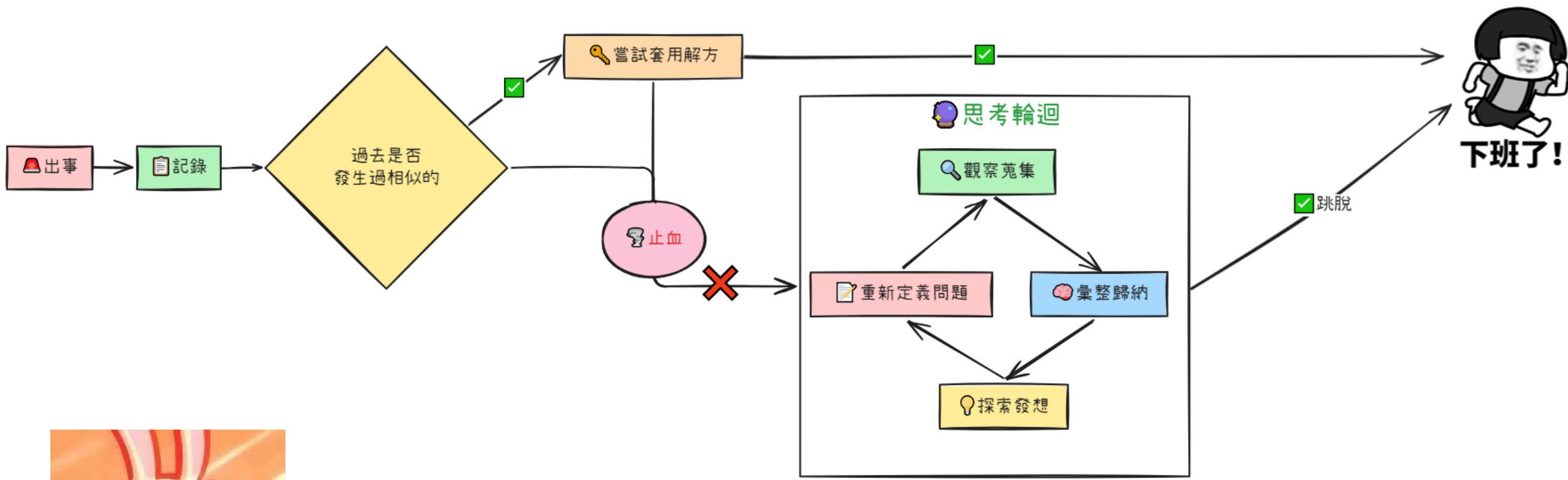
撞牆時，相信直覺

- **跳脫慣性思維**

- 每次都是網路問題啦
- 每次都是 DNS
- 每次都是我們系統問題 🤔

- **弄髒手**

- 直接挖底層 (tcp, packet 追查)
- 硬體更換 (曾經的麵包板)



頭一光即即

事件解決後呢？

- 系統性反思 / Deep Dive
 - 回顧解謎流程，拼湊框架
 - 梳理紀錄 & 錯誤經過
 - 整理成文件，並分享
- 持續
 - 監控系統的擴充與完善 (metrics, logs, traces)
 - 預防重於治療，alert 與 playbook 的定期盤點
 - 接觸更廣的技術維度



Takeaway

- 不要過度信任元件（包括開源項目）的完美運行
- 認識元件流程與工作原理
 - 使用 ChatGPT、Perplexity 等工具輔助擴展想法
 - 看文件，看之前的事件處理
- 保持冷靜與好心情

Thank you!

✉ | sean22492249@gmail.com

🌐 | <https://kiwi-walk.com>

M | <https://sean22492249.medium.com/>

