



我们在生长……

永辉超市-故障演练实战与落地

2021年11月26日

目录

故障演练落地



大厂重大故障回顾



故障演练发展里程碑



平台落地实战



业务落地执行



Q&A环节



Amazon Web Services » Service Health Dashboard

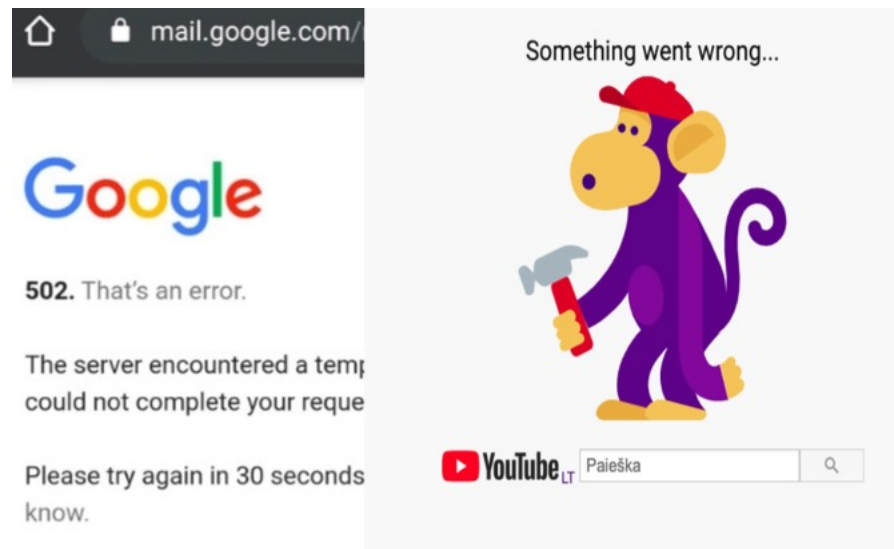
Increased Error Rates

Latest Update (6:23 PM PST): We'd like to provide an update on the issue affecting the Kinesis Data Streams API, and other dependent services, within the US-EAST-1 Region. We have now fully mitigated the impact to the subsystem within Kinesis that is responsible for the processing of incoming requests and are no longer seeing increased error rates or latencies. However, we are not yet taking the full traffic load and are working to relax request throttles on the service. Over the next few hours we expect to relax these throttles to previous levels. We expect customers to begin seeing recovery as these throttles are relaxed over this timeframe.

CloudWatch metrics remain delayed in the US-EAST-1 Region. Once we have restored the throttles for Kinesis to previous levels, we will be restoring CloudWatch metrics functionality. We expect to see recovery of CloudWatch metrics at that stage for new incoming metrics, but the backlog of metrics may take additional time to populate.

We will continue to keep you updated on our progress.

Currently Impacted Services: ACM, Amplify Console, AppStream2, AppSync, Athena, Batch, CodeArtifact, CodeGuru Profiler, CodeGuru Reviewer, CloudFormation, CloudMap, CloudTrail, Connect, Comprehend, DynamoDB, Elastic Beanstalk, EventBridge, IoT Services, Lambda, LEX, Macie, Managed Blockchain, Marketplace, MediaLive, MediaConvert, Personalize, RDS Performance Insights, Rekognition, SageMaker, and Workspaces.



2020.11.25 AWS北弗吉尼亚服务大面积故障，多个云产品服务受到影响，问题定位4个小时，完全修复花费15小时。

2020.12.15 谷歌出现今年第三次大规模宕机，宕机大约45分钟，波及20亿用户，预计损失170万美元广告收入。

- 2017年10月30日，阿里云部分服务器有30多分钟的时间无法正常访问，起因是电力故障。
- 2019年3月23日，腾讯多个产品出现大规模宕机，上海南汇网络光纤因施工被意外挖断。

要具备“用好云”能力，同时也要具备“面向云容灾”的能力

Netflix 宕机事件故事



- 2008年自建机房意外发生故障，让Netflix业务停了整整3天时间，Netflix考虑将所有数据中心搬到AWS云端；
- 2012年12月24日，AWS弹性负载均衡服务中断导致Netflix网站服务受到中断影响。
- 2015年9月20日，AWS数据中心平台上的20多种服务开始出现故障后，互联网上的一些知名网站和应用系统随之间歇性地无法使用，而Netflix网站可以很快恢复服务，问题来了他们问什么可以快速恢复故障，Netflix 这些年做了什么？

故障演练领域重要发展时间线

2010



Netflix Eng Tools团队开发出了Chaos Monkey。当时Netflix从物理基础设施迁移到AWS上，为了保证AWS实例的故障不会给Netflix的用户体验造成影响，他们开发了这个工具，用来测试系统稳定性。

2012



Netflix在 GitHub上开源了Chaos Monkey，并声称他们“已经找到了应对主要非预期故障的解决方案。通过经常性地制造故障，我们的服务因此变得更有弹性。”

2014



Chaos Monkey的升级版FIT诞生，实现微服务级别的故障注入。

2016



混沌工程(Chaos Engineering)的概念通过《Principles Of Chaos》这本书被提出。面向失败设计也逐渐成为构建一个高可用架构产品的基本要求，混沌工程也在更多公司进行实践。

2018



混沌工程（Chaos Engineering）成为CNCF的一个新的技术领域。

2019



阿里巴巴把集团内部故障演练工具进行整合开源出ChaosBlade项目，目前社区生态比较活跃，覆盖故障场景比较全面。

2020

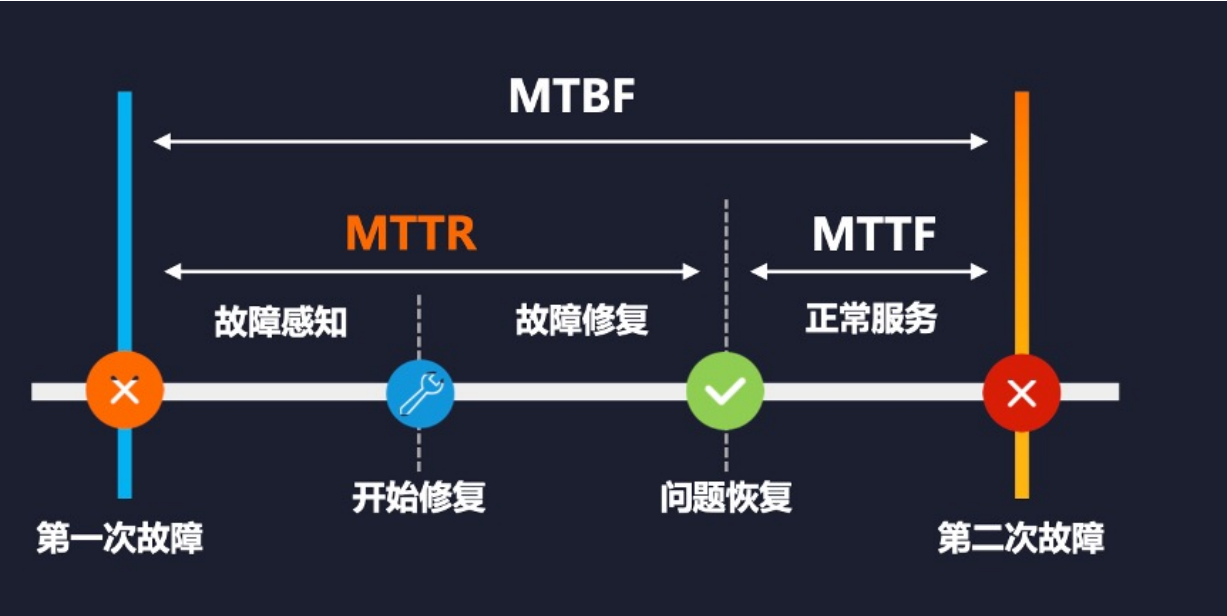


PingCap引入混沌工程项目用于测试TiDB在各种极端场景下数据产品的可用性和健壮性，并推出Chaos Mesh开源项目和进行社区运维，这个项目已经成功进入 CNCF 沙箱托管项目

Chaos Engineering

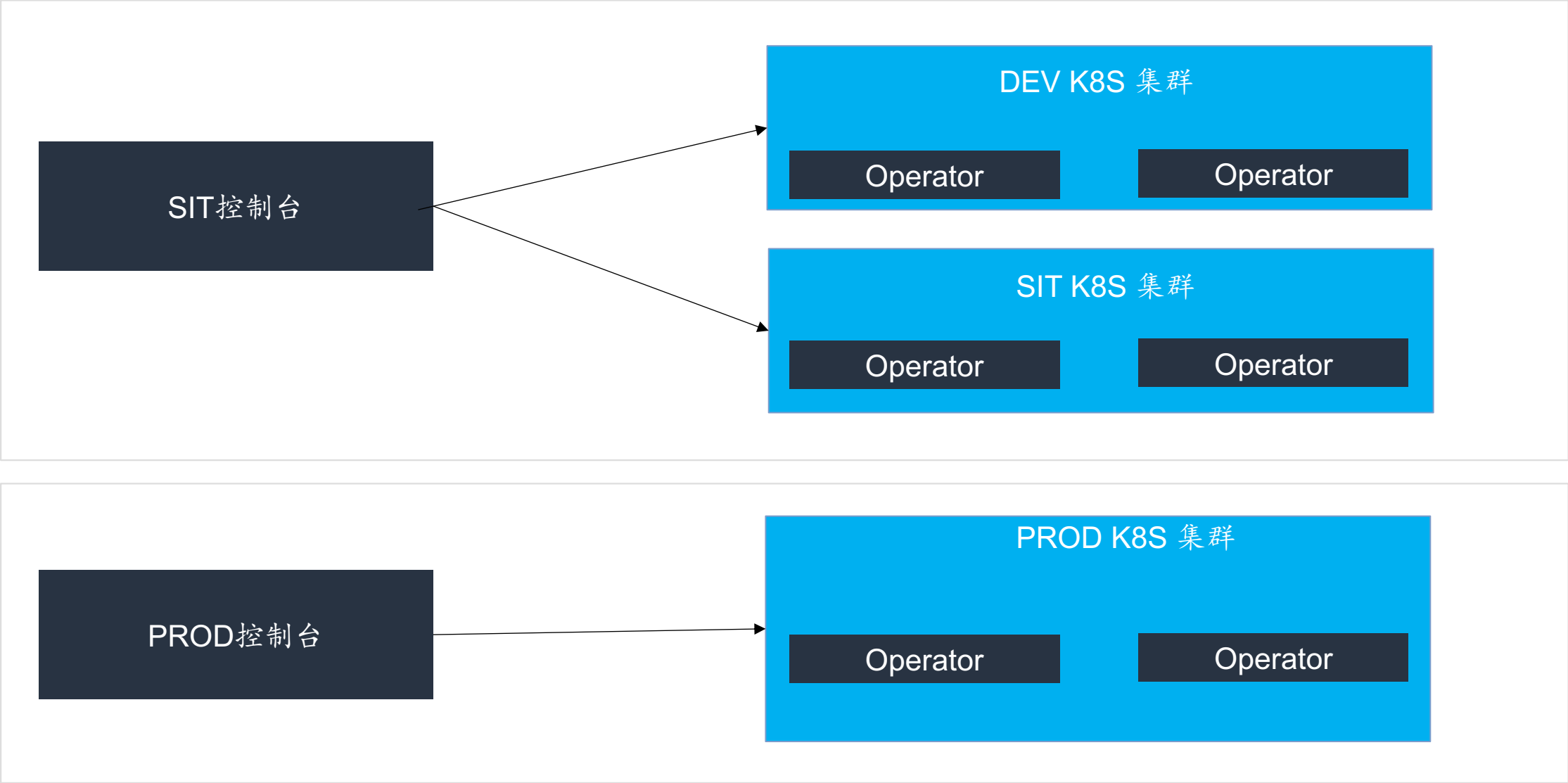


故障演练价值



平均恢复时间 MTTR	<ul style="list-style-type: none"> 发现缺失的系统监控、告警，缩短故障发现时间 提高技术人员应急效率，缩短故障解决时间
平均无故障时间 MTTF	<ul style="list-style-type: none"> 历史故障持续进行回放，避免相同的故障再次发生 架构薄弱点持续进行验证，避免架构稳定性衰弱
平均时效间隔 MTBF	<ul style="list-style-type: none"> 通过提供平均无故障时间和降低平均修复时间，可以有效提升平均时效间隔

L3	L2	L1	L0
1小时	2小时		
2小时			
0.5小时	1小时	2小时	
1小时	2小时		
2小时			
0.5小时	1小时	2小时	
1小时	2小时		
2小时			
5分钟	0.5小时	1小时	2小时
0.5小时	1小时	2小时	4小时
1小时	2小时	4小时	
5分钟	0.5小时		
0.5小时	1小时		
1小时	2小时		
0.5小时	1小时	2小时	
1小时	2小时	3小时	
2小时			
0.5小时	1小时	2小时	3小时
1小时	2小时	3小时	4小时



平台落地实战



版本组合搭配

Chaosblade v1.5.0

Chaos Box v0.4.2

优点:

- 版本发布周期长，稳定
- 集群接入比较简单





平台落地实战

版本组合搭配

Chaosblade v1.6.1

Chaos Box v1.0.1

优点：

- 商业化控制台体验
- 演练库经验能力沉淀
- 租户化环境数据隔离能力

应用高可用服务

概览

我的空间

演练场景

演练经验

应用管理

探针管理

数据管理

应用高可用服务 / 故障演练 / My spac

default

My spac

● 执行过的演练 426

⌚ 失败:1

⌚ 运行中:0

✅ 成功:425

● 未执行的演练 54

● 总演练数480

新建演练 >

全部停止

请选择状态

已选0个标签

请输入演练名称

只看定时演练

应用接入

演练名称	标签	场景	创建时间	定时任务	最近运行状态	最近运行时间	操作
fp-view-api 宕机演练0707		C	2022-07-07 21:00:53		成功	2022-07-07 21:00:57	演练 拷贝 删除
fp-view-api 宕机演练		C	2022-07-06 17:21:58		不符合预期	2022-07-08 15:35:58	演练 拷贝 删除
Examples-Cpu-Fullload		C C	2022-07-05 16:40:03		成功	2022-07-11 10:53:07	演练 拷贝 删除
fp-view-api 线程升高			2022-07-04 04:52:55				演练 拷贝 删除
fp-view-api 线程升高			2022-07-04 04:37:28				演练 拷贝 删除
fp-view-api dubbo线程池满			2022-07-04 04:27:06				演练 拷贝 删除

业务落地实战

实战案例：互联网前端业务线故障演练

故障演练组织核心阶段：

- 规划演练计划
- 制定演练目标清单
- 确认故障演练参与相关方角色
- 确认演练业务线配合时间、地点
- 组织演练执行
- 复盘演练结果

实战案例-制定演练目标清单

演练SOP模板地址: <https://docs.qq.com/sheet/DS1JPbIFibUNRWGV3?tab=BB08J2>

演练模板Sheet功能介绍如下:

演练模板填写说明: 这里面介绍故障演练具有哪些能力, 演练发起人可以参考这些能力指定演练目标; 也定义了演练清单里面每个字段的填写说明;

演练计划与任务事项: 这个sheet里面把演练计划和任务事项进行公开, 可以根据实际情况进行填写;

演练参与人员名单: 故障演练是一个团队参与的活动, 需要很多角色参与配合; 这个里面把每个人的工作安排给明确行下;

演练清单: 里面每一行介绍了本次演练一次动作编排, 和演练预期行为需要补充完善;

演练复盘: 演练结束后针对本次演练清单的结构进行复盘, 针对问题需要给出跟进人和修复时间;



实战案例-规划演练计划

标题	投入	十一月 17日	十一月 18日	十一月 19日	十一月 22日	十一月 23日	十一月 24日	十一月 25日	十一月 26日	十一月 29日	十一月 30日	十二月 1日	十二月 2日	十二月 3日
▼ 1) 方案设计	2 周 3 天	▼												
• 1.1) 方案设计	4 天				演练发起人; 业务研发									
• 1.2) 目标方案评审	3 天					业务研发; 业务负责人; 业务测试								
• 1.3) 演练计划评审	3 天					业务研发; 业务负责人; 业务测试								
• 1.4) 演练参与人员确认	3 天					业务研发; 业务负责人; 业务测试								
▼ 2) 方案准备	4 周 4 天				▼									
• 2.1) 资源信息梳理	1 周 3 天							业务研发; 业务测试; 演练发起人; 演练负责人						
• 2.2) 监控渠道梳理	1 周 3 天							业务研发; 业务测试; 演练发起人; 演练负责人						
• 2.3) 故障场景编排	1 周 3 天							业务研发; 业务测试; 演练发起人; 演练负责人						
▼ 3) 故障正式演练	2 周 4 天 6.75 小时						▼							
• 3.1) 相关人员准备事项确认	2 天 7.75 小时													
• 3.2) 故障发起	2 天 7.75 小时												演练发起人	
• 3.3) 故障业务测试	2 天 7.75 小时												业务测试	
• 3.4) 故障演练结论评估	2 天 7.75 小时												业务负责人	
• 3.5) 故障演练恢复	2 天 7.75 小时												演练发起人	
▼ 4) 故障演练复盘	1 周									▼				
• 4.1) 故障演练结论复盘	1 周												业务研发; 业务测试; 演练发起人; 业务负责人; 演练负责	

实战案例-组织演练执行



Q&A



融合共享 成于至善