

# Geekbang>

极客邦科技

全球领先的技术人学习和交流平台

扫我，码上开启新世界



# Geekbang>

InfoQ | EGO NETWORKS | StuQ

## InfoQ

专注中高端技术人员  
的社区媒体

## EGO NETWORKS

EXTRA GEEKS' ORGANIZATION  
高端技术人员  
学习型社交网络

## StuQ

实践驱动的IT职业  
学习和服务平台

促进软件开发领域知识与创新的传播

InfoQ<sup>new</sup>

**QCon**  
全球软件开发大会

**[上海]** 2015年10月15-17日

**ArchSummit**  
全球架构师峰会

**[北京]** 2015年12月18日-19日



关注InfoQ官方微信  
及时获取ArchSummit演讲视频信息

# 大型商业银行故障处理实践 应急标准化方法论

张春林

2015年7月18日

- ∞ **黑天鹅事件** —— **哲学思辨**
- ∞ **应急处理的思考点** —— **理性分析**
- ∞ **应急标准化方法论** —— **成果分享**

- 黑天鹅事件，总是在我们自以为是的时刻发生。



宁夏银行数据库故障遭银保监会通报 业务中断达57小时

2016年06月06日 来源：和讯银行

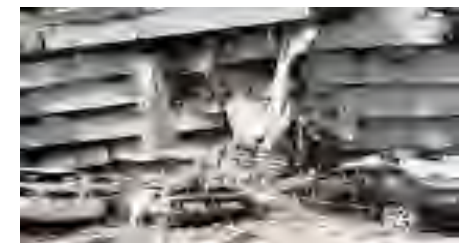


工商银行系统瘫痪 不敢以“道歉”结束

工商银行系统瘫痪 不敢以“道歉”结束

工商银行系统瘫痪 不敢以“道歉”结束

携程宕机的损失为平均每小时  
106.48 万美元



支付宝大面积瘫痪 灾备能力引发争议

支付宝大面积瘫痪 灾备能力引发争议

- 从天灾到人祸
- 从基础设施损坏到人为操作风险

- 无论机构规模大小
- 无论IT成熟与否
- “互联网+”背后的危机

- 金融体系的脆弱性理论告诉我们，金融危机的发生是由金融体系的脆弱性内生决定的。同理，信息科技风险也是由信息系统的脆弱性内生决定的。
- 硬盘的坏盘率
- 软件代码的BUG率
- 人员更替
- 以已知对抗未知

纳西姆·尼古拉斯·塔勒布在

《黑天鹅》中给出的建议：

- 1、不要预测。
- 2、谨慎预防。
- 3、保持充足冗余。

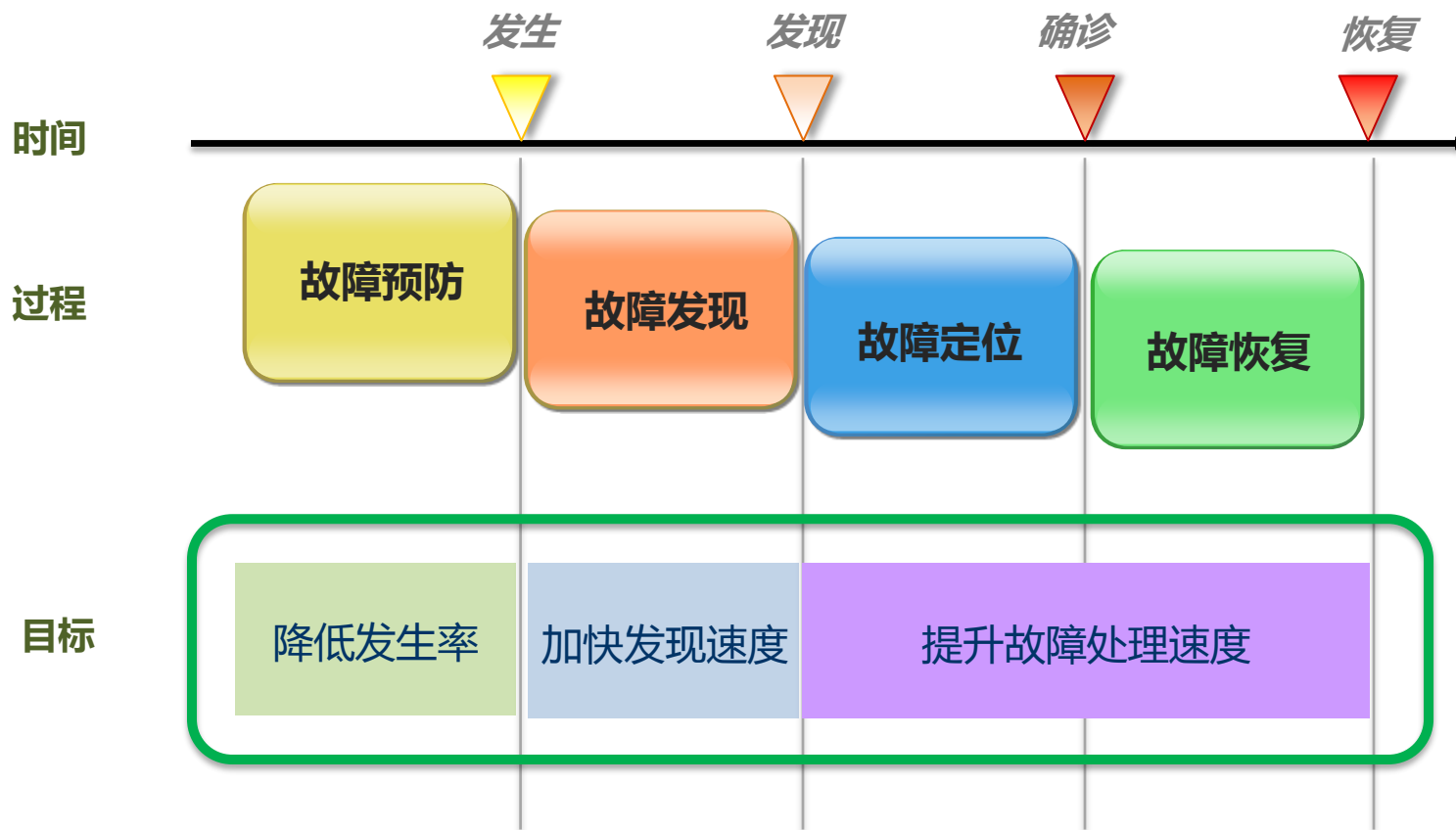


- ∞ 黑天鹅事件 —— 哲学思辩
- ∞ 应急处理的思考点 —— 理性分析
- ∞ 应急标准化方法论 —— 成果分享

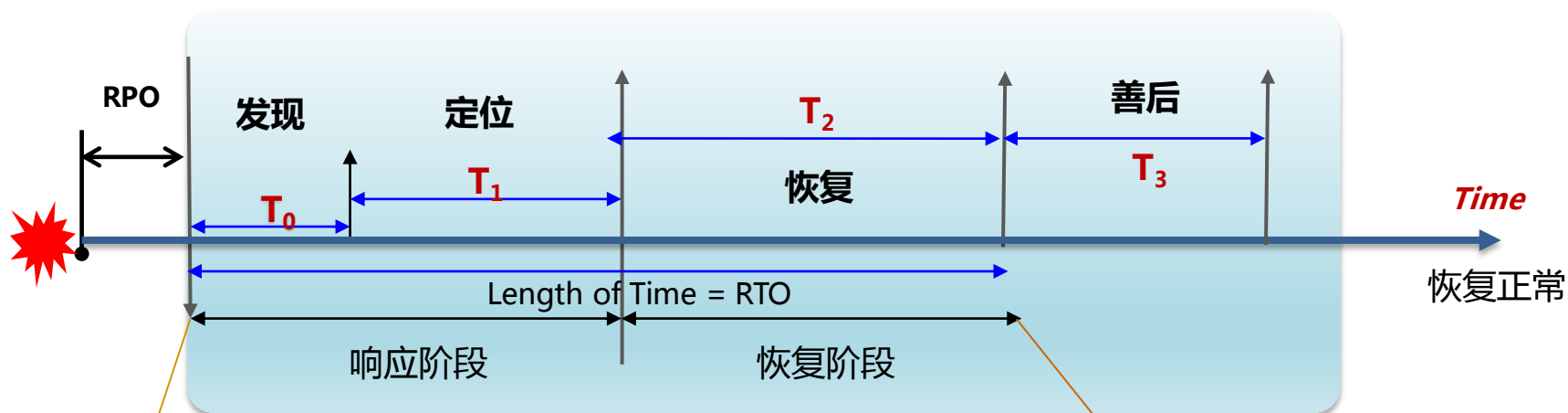


### 故障处理的思考点

- ❖ 故障恢复时长（RTO）是衡量一个业务连续性的关键指标，它的长短决定了业务影响。
- ❖ 从故障发生开始到故障恢复截止，整个过程包括发现、定位、恢复三阶段。



## 应急处置过程



### 现状：

$T_0 = 10$ 分钟  
 $T_1 = 30$ 分钟  
 $T_2 = 20$ 分钟

### 目标：

$T_0 = 0$ 分钟  
 $T_1 = 0$ 分钟  
 $T_2 \leq 5$ 分钟

$$RTO = T_0 + T_1 + T_2$$

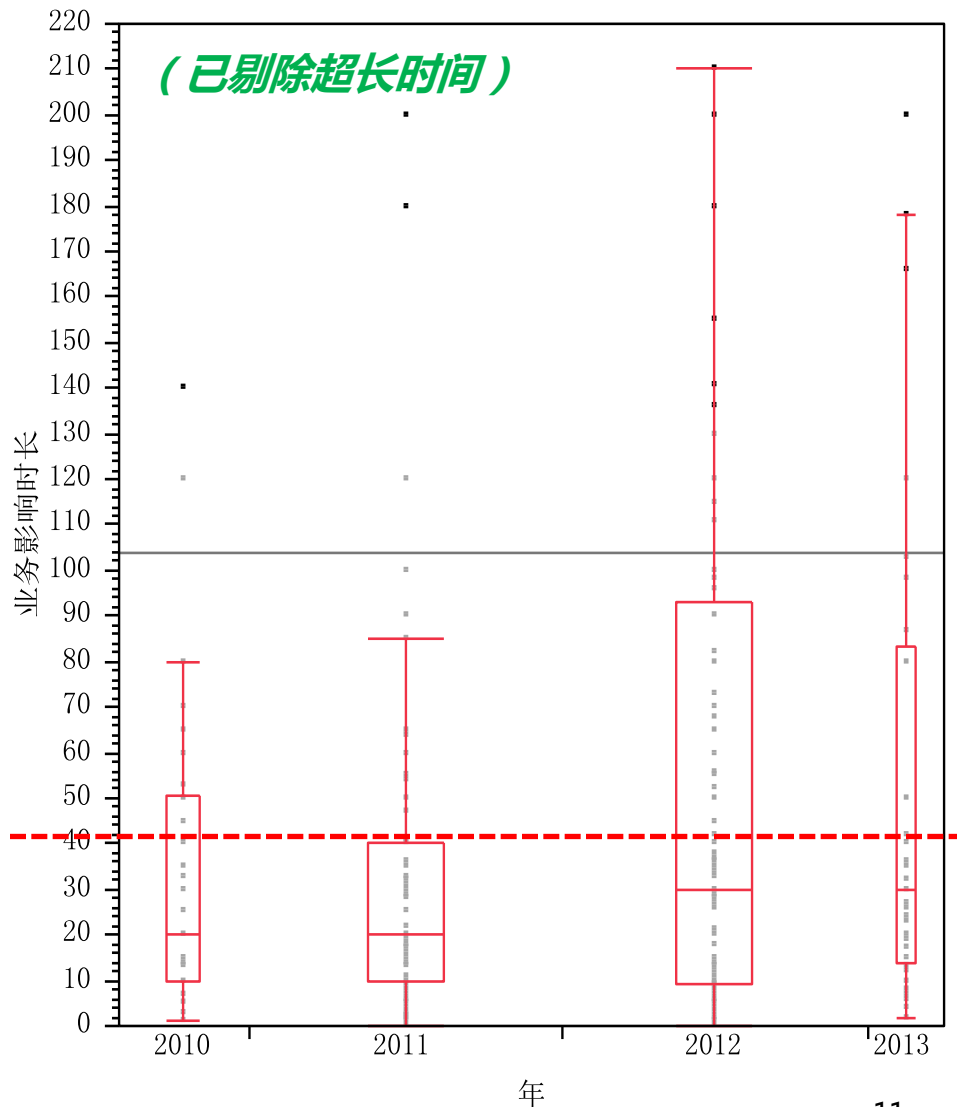
## 分析方法：

精益六西格玛

## 现状概览：

- 1、很多故障的处理时长超过30分钟；
- 2、很多故障无法通过切换来解决，平均每月2起。

## 大量的事件处理时长超过30分钟



## VOC (客户之声)

1. 出现故障后不知道如何处理，找不到适当的应急预案
2. 不知道找谁处理
3. 现场人员混乱

## VOP (流程之声)

1. 缺少通用故障处理流程，过于依赖A角；
2. 演练不到位，处理不够熟练
3. 应急环境（ECC、工具）缺少规范化管理

## VOB (管理之声)

1. 应急通讯手段落后
2. 无统一指挥和标准化流程
3. 30分钟以上故障必须报银监会

- 搜集用户的声音、流程的问题、管理者的声音
- VOC-Voice Of Customer ; VOP-Voice Of Procces ; VOB-Voice Of Business

$$Y = F(X)$$

Y : 故障恢复时长

定义：故障恢复时长RTO，  
故障开始时间至故障恢复时间。  
测量：提取2010年至今的所有故障影响时长，  
取平均值。

X1 : 工具保障操作一致化

定义：**具备电子化工具**，包括监控工具、  
恢复工具和通讯工具  
测量：通过多个项目同步建设工具并联动，实现  
操作一致化。

X2 : 流程的范围与**熟练度**

定义：**具备流程**，包括通用应急指挥流程  
和通用故障处理流程  
测量：通过各种不定时演练抽查**提高熟练度**

## 快速诊断

故障定位时间

- 应急预案
- 监控平台
- 应急平台

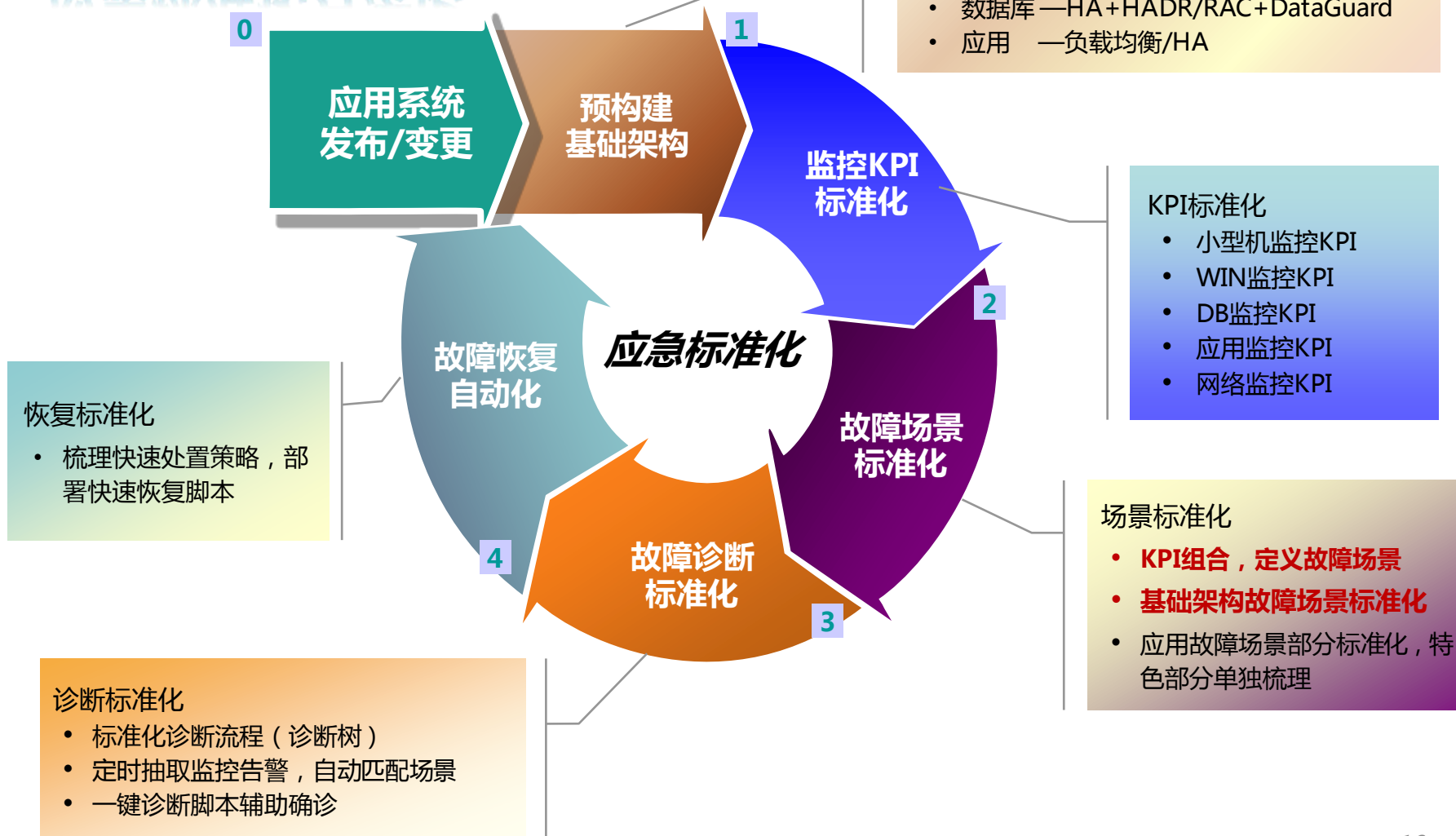
## 快速恢复

故障修复时间

- 故障自愈进程
- 一键恢复工具
- 业务补账工具
- 流量清洗服务

- ∞ 黑天鹅事件 —— 哲学思辩
- ∞ 应急处理的思考点 —— 理性分析
- ∞ 应急标准化方法论 —— 成果分享

## 应急标准化方法论

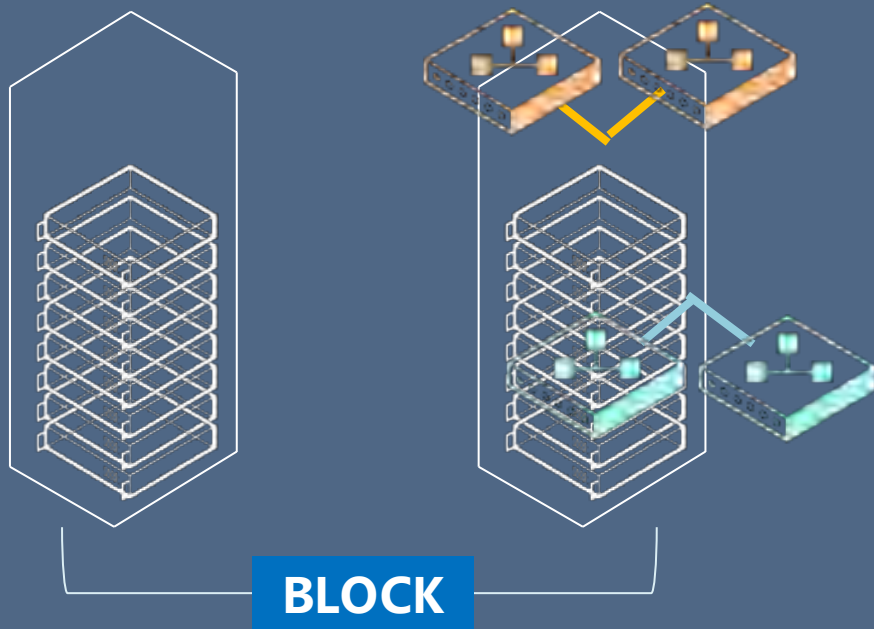




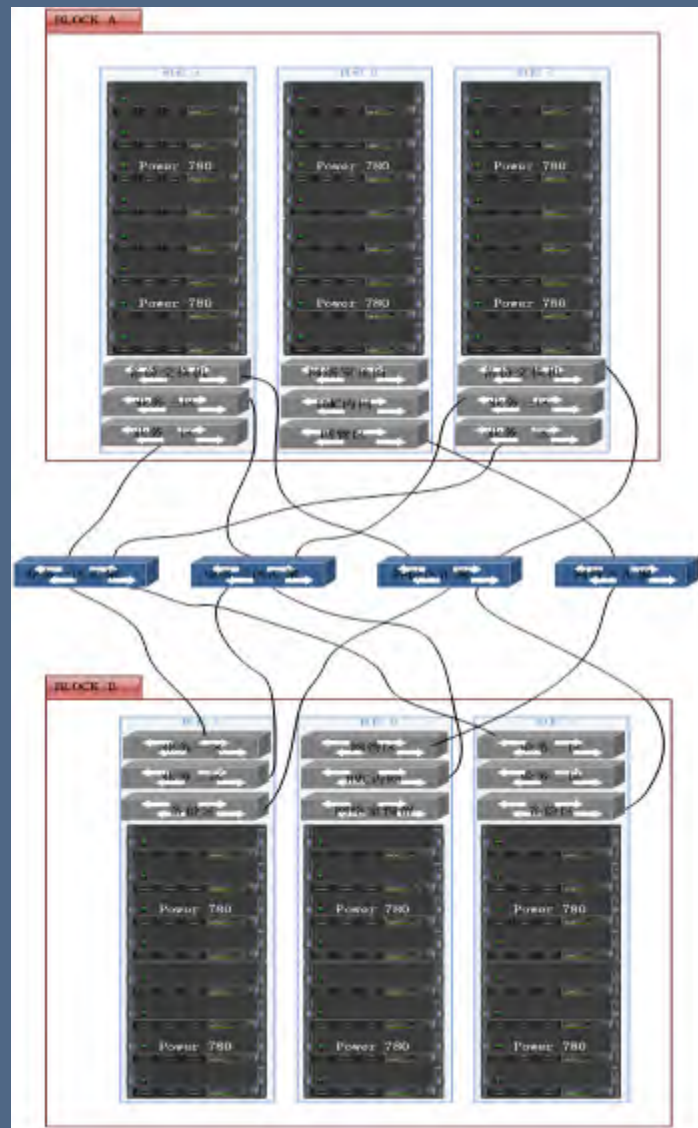
# 架构标准化



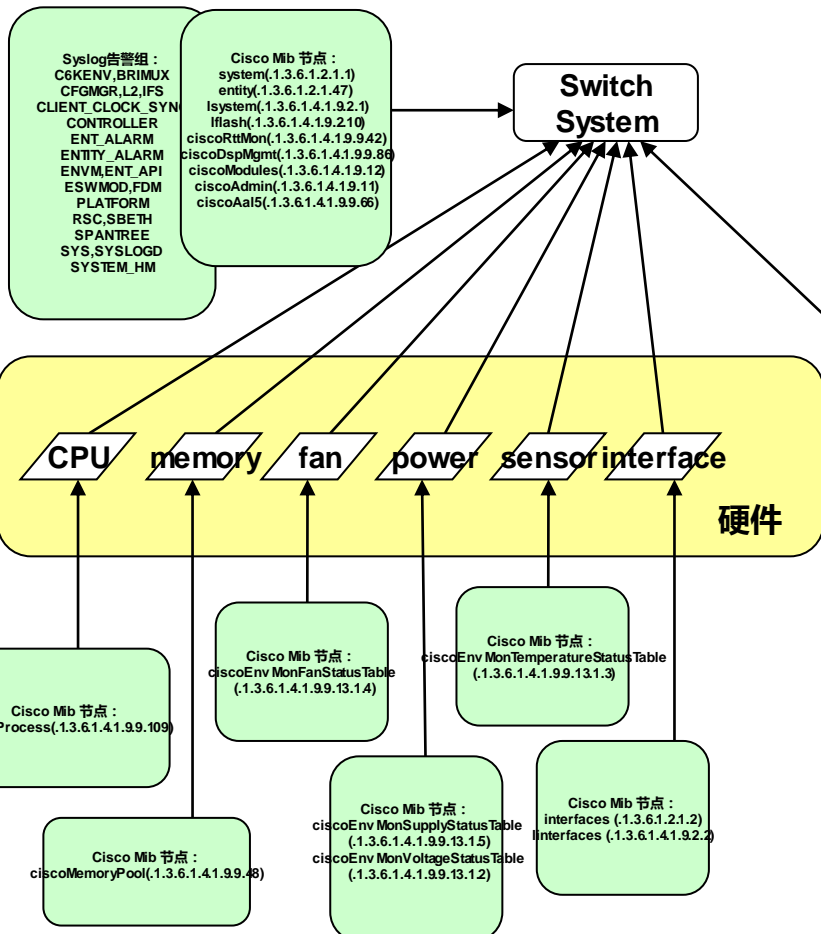
- 标准化硬件配置和部署方式
- 网络双路双区接入
- 标准化同时考虑高可用性
- 组建服务器集群



- ✓ 预构建、可伸缩的基础架构单元
- ✓ 一组计算、存储和网络资源



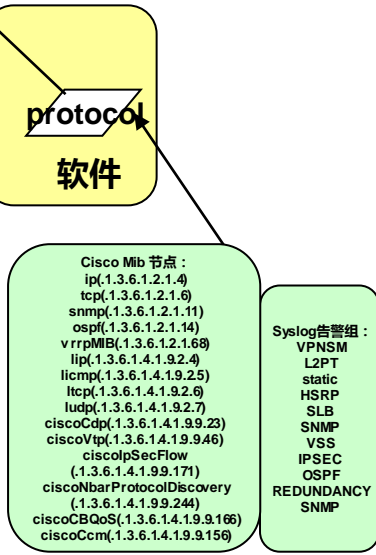
# 监控指标标准化



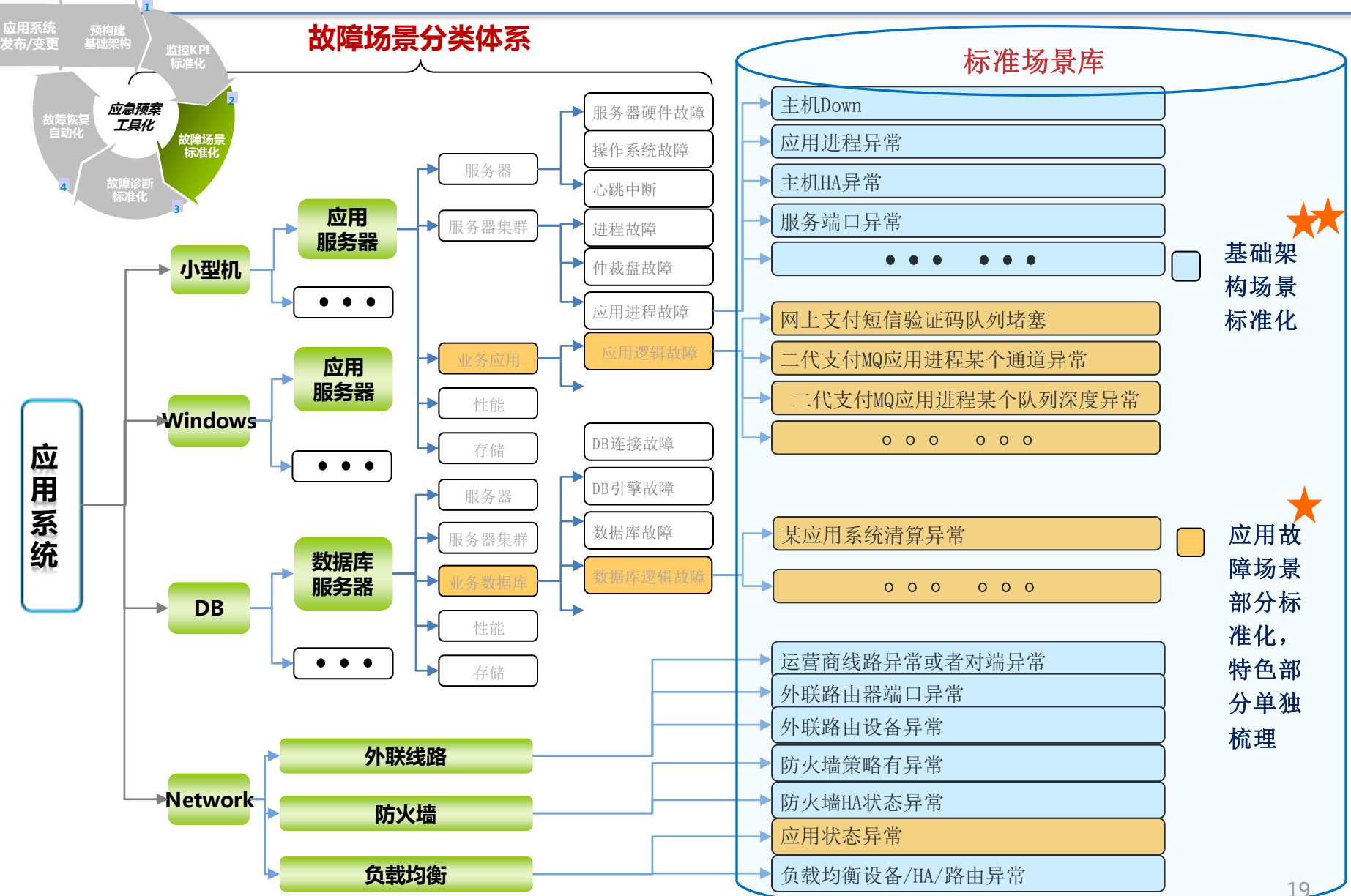
管理对象ID (001-99)	管理对象 (001-99)	管理对象属性 (001-999)	事件代码 (0001-9999)	事件名称	报警数据源	事件描述	报警事件	报警级别	
Network (11)	Switch_Cisco (02)	System (001)	FM-11-02-001-0001	系统可用报警	RD-11-02-001-0001	系统不可用	系统不可用	5分钟	
			FM-11-02-001-0002	系统运行时间过长	RD-11-02-001-0002	达到设备运行时间	1天		
			FM-11-02-001-0003	系统重启	RD-11-02-001-0003	系统重启	5分钟		
			FM-11-02-001-0004	系统老化	RD-11-02-001-0004	系统老化	5分钟		
			FM-11-02-001-0005	系统配置更改	RD-11-02-001-0005	系统配置更改	5分钟		
			FM-11-02-001-0006	系统PT Root登陆	RD-11-02-001-0007	系统PT Root登陆	5分钟		
			FM-11-02-001-0007	系统崩溃	RD-11-02-001-0008	系统崩溃	5分钟		
		FM-11-02-002-0001	CPU使用率报警	RD-11-02-002-0001	CPU使用率报警	5分钟			
		FM-11-02-002-0002	CPU可用报警	RD-11-02-002-0002	CPU可用报警	5分钟			
		FM-11-02-003-0001	内存使用率报警	RD-11-02-003-0001	内存使用率报警	5分钟			
FM-11-02-003-0002	内存不可用	RD-11-02-003-0002	内存不可用	5分钟					
Interface (004)	Interface (004)	FM-11-02-004-0001	端口不通	RD-11-02-004-0001	端口不通	5分钟			
		FM-11-02-004-0002	端口入流量报警	RD-11-02-004-0002	端口入流量报警	>=10%	5分钟		
		FM-11-02-004-0003	端口出流量报警	RD-11-02-004-0003	端口出流量报警	>=10%	5分钟		
		FM-11-02-004-0004	端口入CRC错包报警	RD-11-02-004-0004	端口入CRC错包报警	>=10%	5分钟		
		FM-11-02-004-0005	端口出CRC错包报警	RD-11-02-004-0005	端口出CRC错包报警	>=10%	5分钟		
Protocol (005)	Protocol (005)	FM-11-02-005-0001	协议报警	RD-11-02-005-0001	协议报警	5分钟			
		FM-11-02-006-0001	配置不可用	RD-11-02-006-0001	配置不可用	5分钟			
		FM-11-02-007-0001	配置不可用	RD-11-02-007-0001	配置不可用	5分钟			
Sensor (006)	Sensor (006)	FM-11-02-008-0001	传感器不可用	RD-11-02-008-0001	传感器不可用	5分钟			

监控KPI标准库

监控对象组件模型

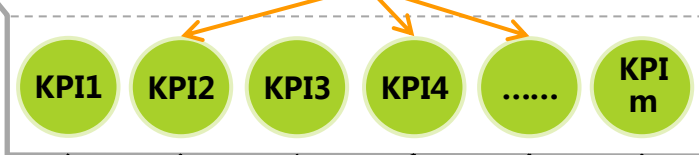


# 故障场景标准化——标准故障场景库





监控平台



应急平台

确诊

一键诊断脚本  
(小型机)

一键诊断脚本  
(WIN)

一键诊断脚本  
(DB)

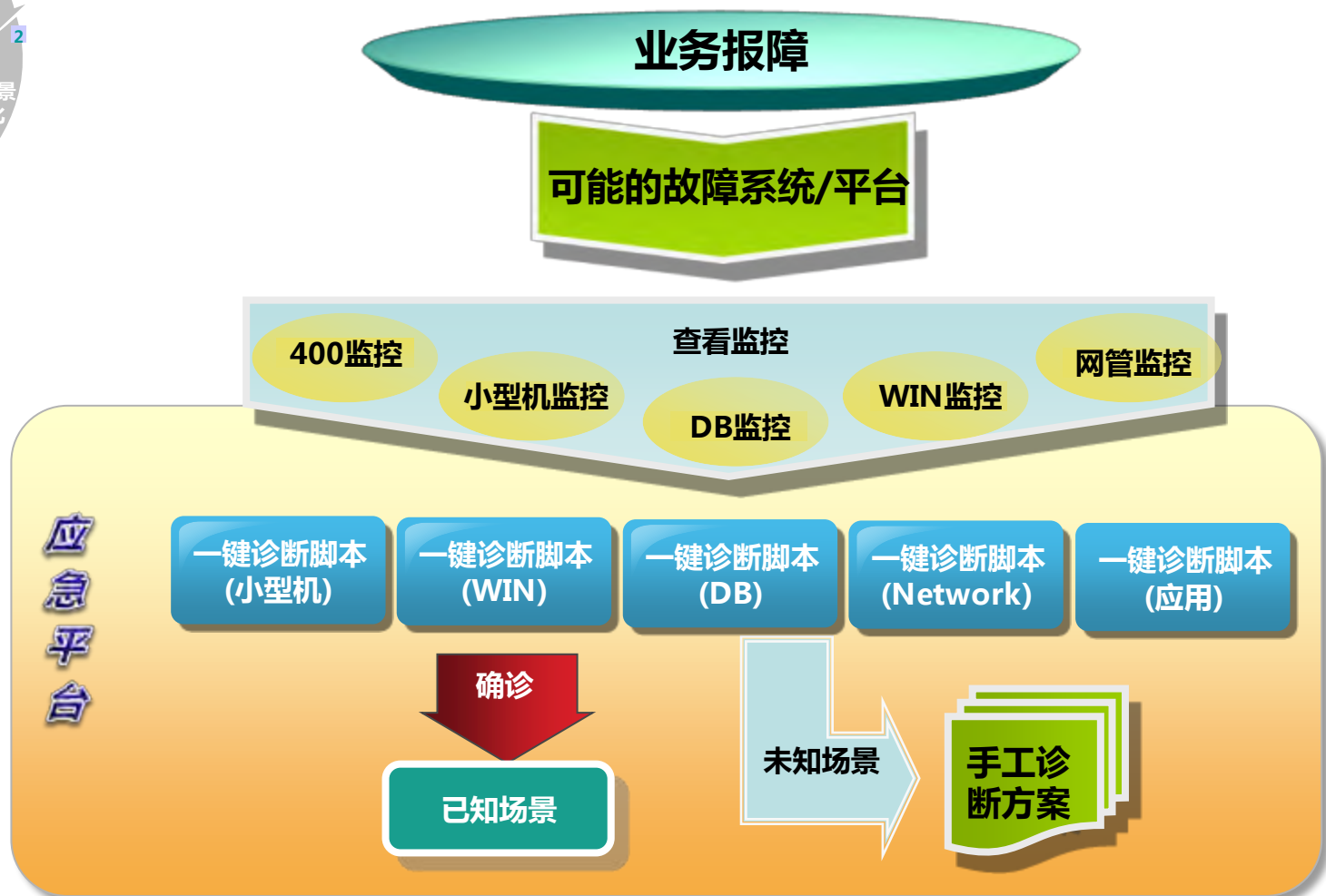
一键诊断脚本  
(Network)

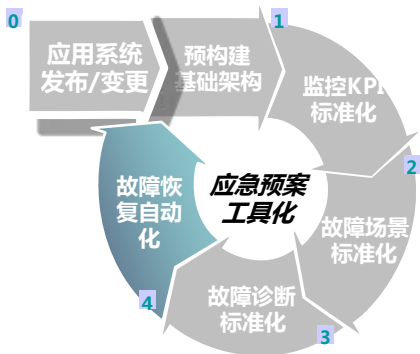
一键诊断脚本  
(应用)

- 自动诊断
- 一键诊断
- 通用流程
- 专家诊断

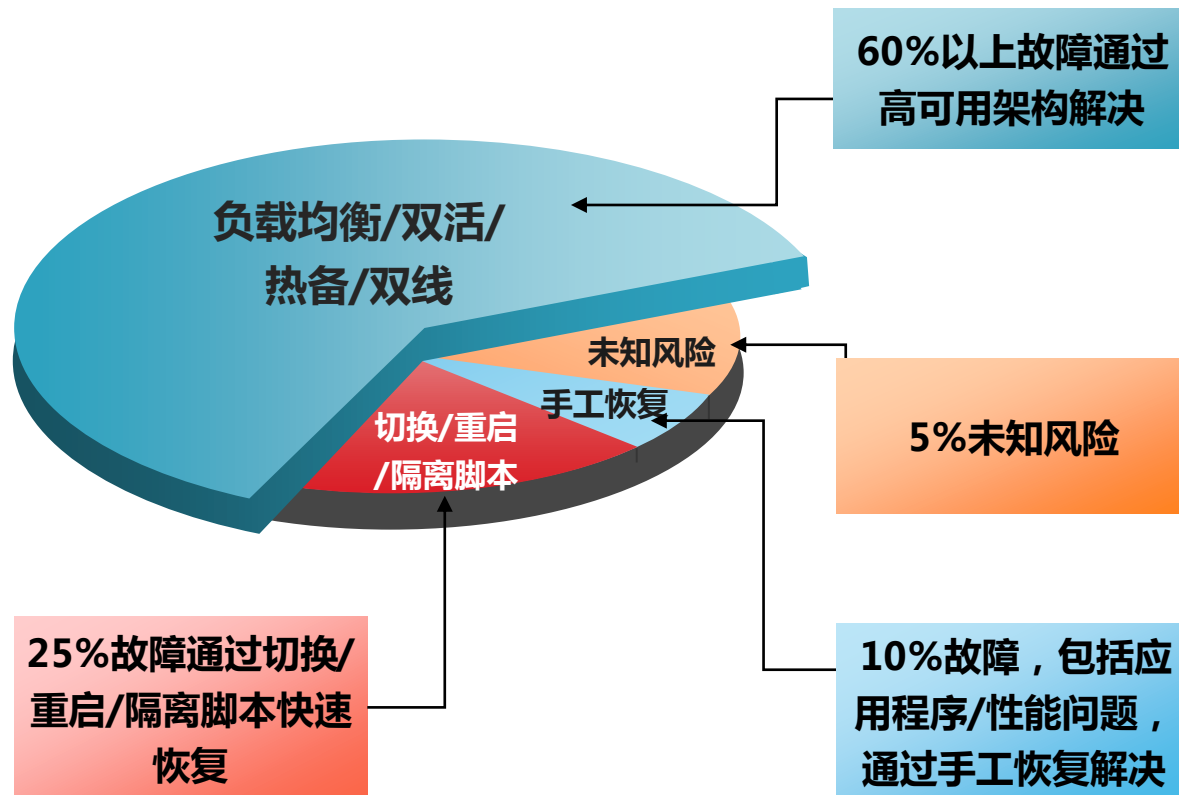
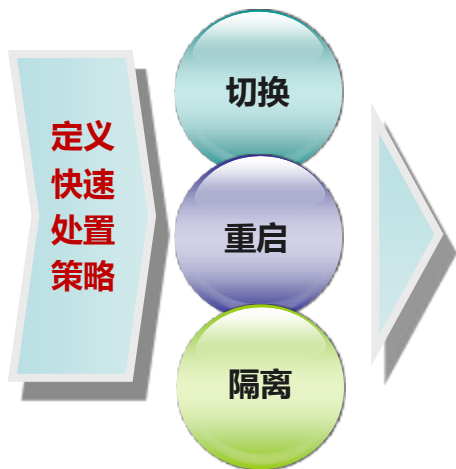


- 自动诊断
- 一键诊断
- 通用流程
- 专家诊断





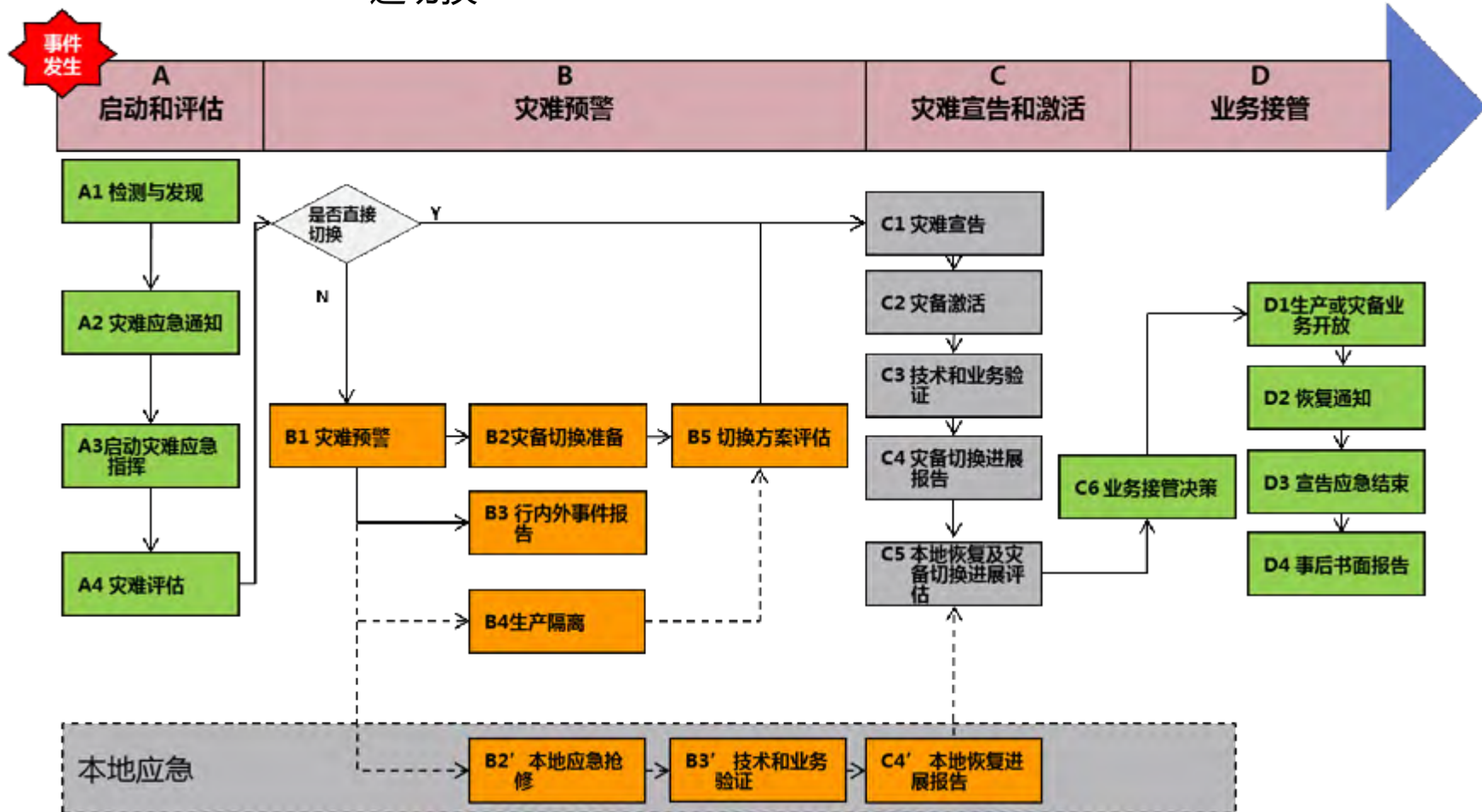
- 故障自愈
- 一键恢复
- 手工恢复



## 低频高损事件

**快速切换：**  
通过灾备导航工具，一键快速切换

**全程监控：**  
执行过程全景监控



## 高频低损事件

通过与监控系统建立接口，对该类事件发生信息进行  
**及时推送、记录和汇总。**

### 提醒：

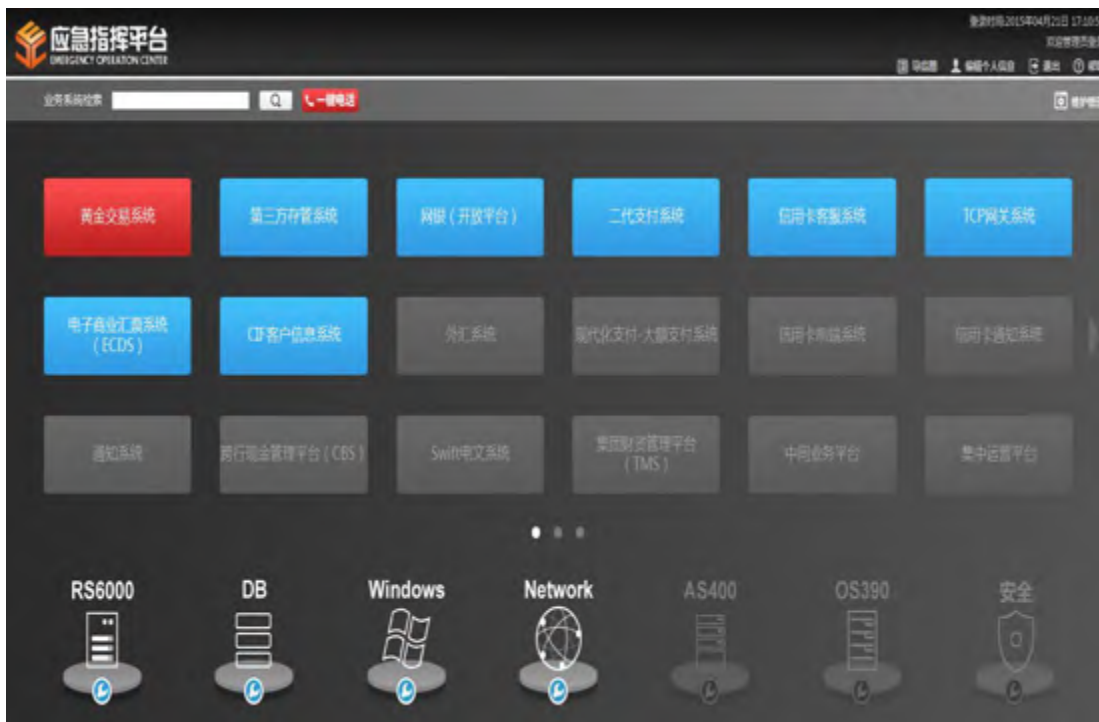
通过自动恢复可以解决的事件，在事件发现阶段、自动恢复阶段、恢复完成阶段，均有信息发送至值班经理和值班管理员公告板处进行提醒。

### 留存：

通过公告板及日志模块，可以查看事件处理的过程信息。

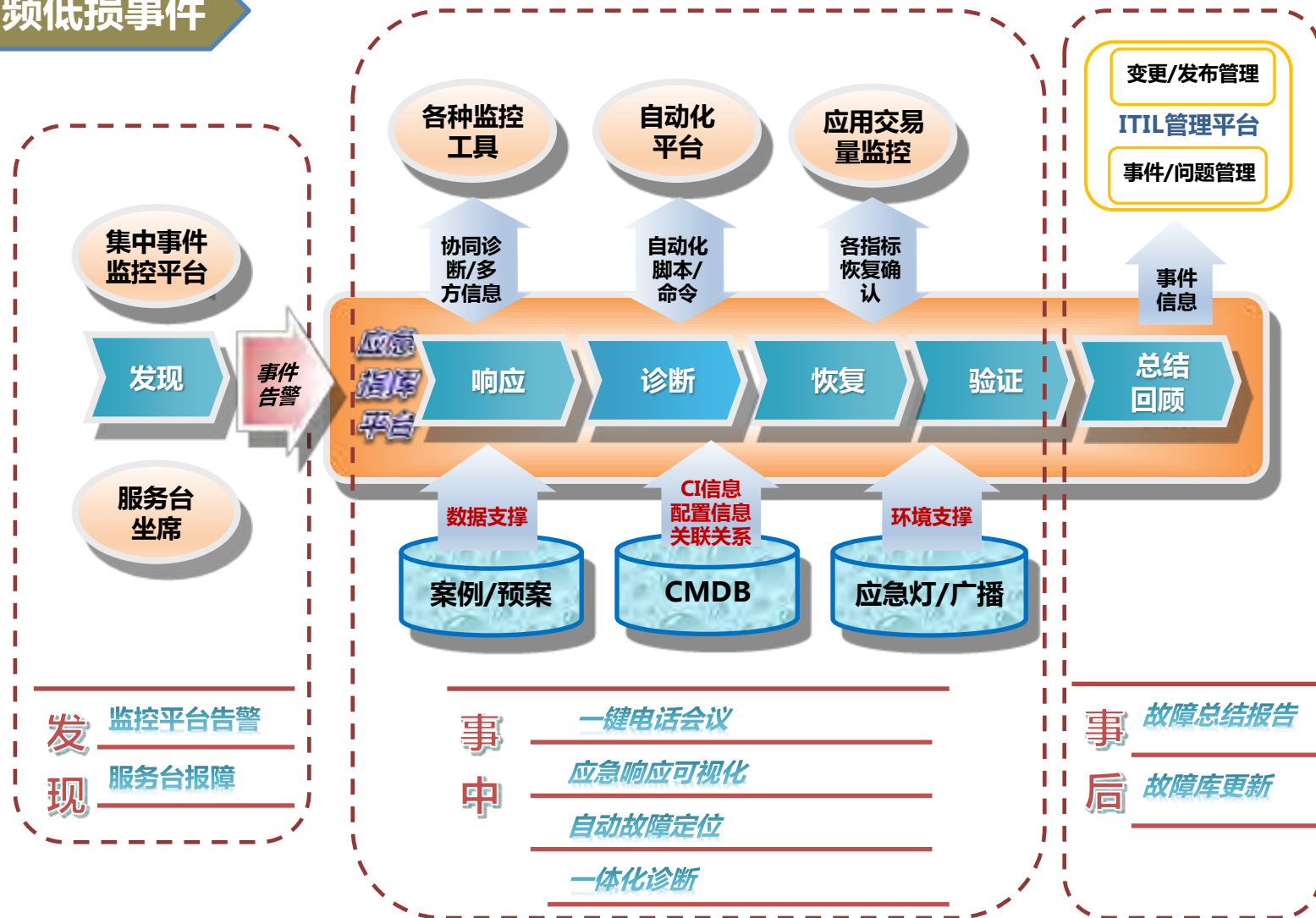
### 分析：

定期汇总生成报告，并进行统计分析。





## 高频低损事件





# ArchSummit 2015

**谢谢！**