



2024 AI+研发数字峰会

AI+ Development Digital summit

AI驱动研发变革 促进企业降本增效

北京站 08/16-17

基于大模型的根因分析实战

文吉 畅捷通信息技术股份有限公司

科技生态圈峰会 + 深度研习



—1000+ 技术团队的选择



上海站

K+全球软件研发行业创新峰会

时间: 2024.06.21-22



敦煌站

K+思考周®研习社

时间: 2024.10.17-19



香港站

K+思考周®研习社

时间: 2024.11.10-12



K+峰会详情



上海站

Ai+研发数字峰会

时间: 2024.05.17-18



北京站

Ai+研发数字峰会

时间: 2024.08.16-17



深圳站

Ai+研发数字峰会

时间: 2024.11.08-09



AiDD峰会详情



2024 AI+研发数字峰会

AI+ Development Digital summit

深圳站 11/08-09

AI 驱动研发变革 促进企业降本增效

2024深圳站-议题设置

AI+产品线	LLM驱动产品创新	LLM驱动需求与业务分析	AI驱动设计与用户体验
AI+开发线	AI 原生应用开发框架与技术	AI Agents在研发落地实践	LLM驱动编程与单测
AI+测试线	LLM驱动测试分析与设计	基于LLM生成测试脚本与数据	LLM和AI应用的评测
AI+工程线	AI+DevOps 与工具 (LLM 时代的平台工程)	大模型对齐与安全	端侧大模型与云端协同
AI+领域线	领域大模型 SFT 与优化	知识增强与数据智能	大厂专场

扫描右侧二维码
查看更多会议详情



早鸟票限时抢购中 (截止到9月30日)

¥ 3680

早鸟票

¥ 2800

学生票



文吉

十年以上SRE实战经验，特别是对ToB场景有丰富实战经验

用友集团 P9高级专家

多次对外分享，融合大模型能力升级智能运维

荣获了信通院颁发的“稳定性优秀案例”

目录

CONTENTS

1. 背景
2. 问题/痛点
3. 解决思路/整体方案
4. 具体实现/技术实践
5. 总结与展望

PART 01

背景

▶▶ 畅捷通是做什么的?

畅捷通信息技术股份有限公司是用友旗下成员企业，成立于2010年3月，于2014年在港交所上市，是中国领先的小微企业财税及业务云服务提供商。

业务架构复杂

公有云模式
接入阿里云、华为云、腾讯云等各类型公云，实现标准化运维
深入研究各云组件的特性、配置、使用，实现云原生

C端用户量 + B端客户体量

混合云模式
服务上云、转型公有云运维

要保障每个用户的体验

业务迭代速度快

纯IDC模式
平台层运维：
glusterfs、cloudfoundry



▶ 畅捷通运维转型之路——目标0-2-5-10

业务从自建机房逐步转向全面采用公有云容器化架构，为业务发展提供了更强大的基础，但同时也带来了运维复杂性的指数级增长。



PART 02

问题/痛点

▶ 从一次飞机撞鸟说起

2023年11月1日，旭日8409飞机起飞离地时，发动机遭遇鸟击。
情况万分危急，关系到机上183人的生命安全。

客机起飞时突遭鸟击，机组人员28分钟安全返航获嘉奖



新京报

2023-11-03 20:46 发布于北京 新京报官方账号

+ 关注

新京报讯（记者 吴采倩）“稍等一下，旭日8409，刚才撞鸟了。”11月1日，“航班遭鸟击机组带173名旅客安全返航”一事登上微博热搜，网友表示，驾驶舱内的对话“满满安全感”。

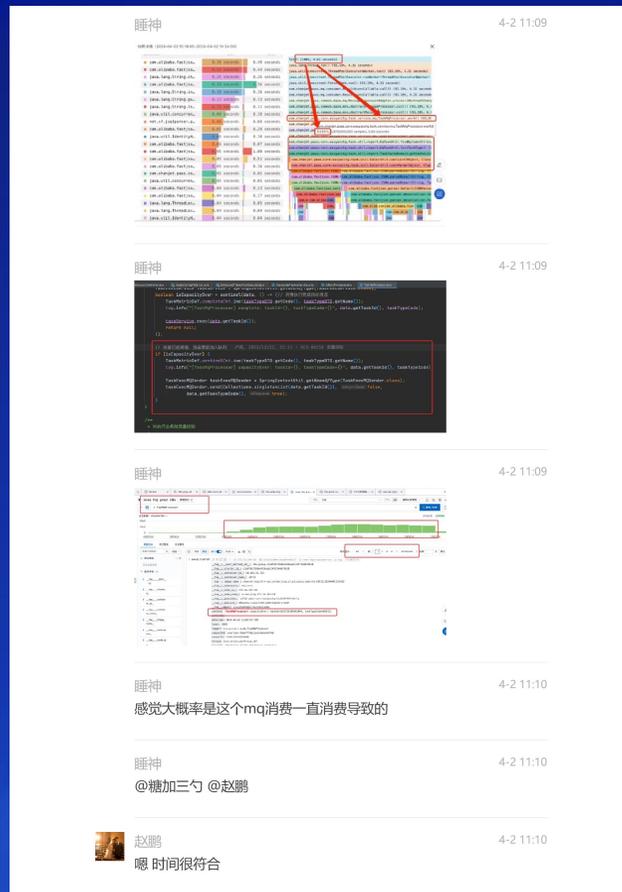
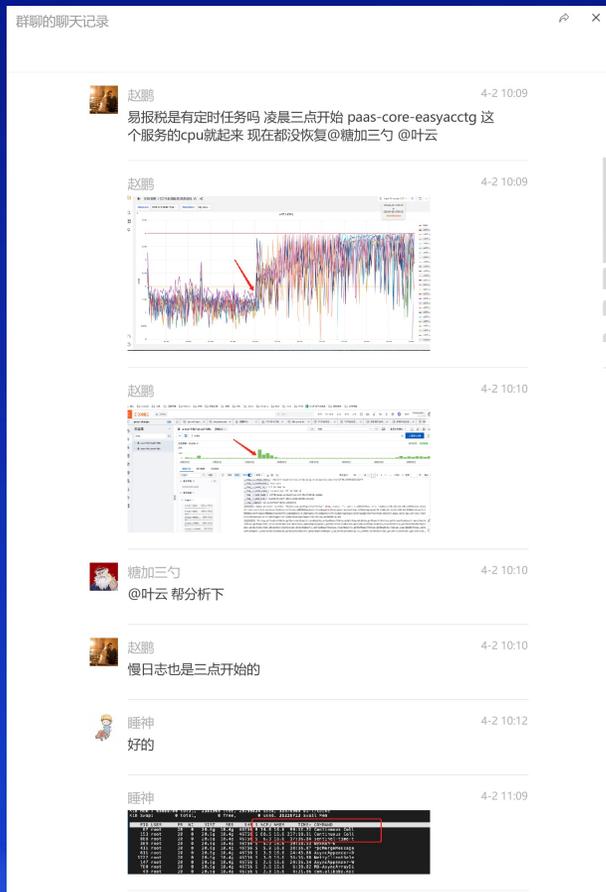
这段对话发生在8月26日，长安航空9H8409机组人员在执飞广西梧州至海南三亚航班时，飞机起飞离地时遭遇鸟击。从鸟击发生到安全返航，机组人员用28分钟化险为夷，保障了机上173名旅客和10名机组人员的生命安全。



▶ 畅捷通运维面临什么样的压力?



▶ 发生故障时难以定位



定位一个问题，需要：

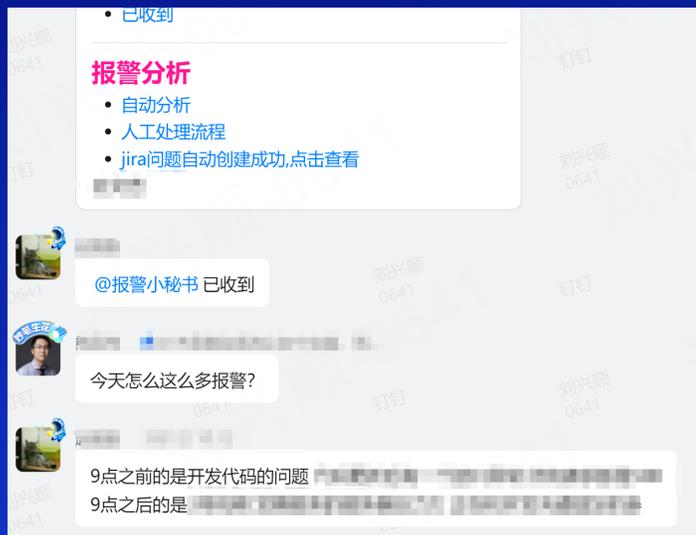
- 打开3-5个看板
- 执行2-4次分析脚本

90%的问题此时就能找到原因，耗时10分钟。

但另10%的问题，才会产生大的故障，且往往难以定位原因



无法快速判断爆炸半径



- 怎么判断报警严重性?
- 报警爆炸半径多大?
- 是否正在处理? 谁在处理?
- 恢复了吗?



▶ 畅捷通运维面临什么样的压力？

客户不能等

线上无法复现



压力山大

无迹可寻



PART 03

解决思路/整体方案

10.1

737 快速检查单

空速不可靠

状况： 怀疑空速或马赫数指示不可靠。（可能表明空速不可靠的项目在“其它信息”中列出）。

目的： 如可能，识别出可靠的空速指示，或使用“性能-飞行中”章节的“空速指示不可靠飞行”图表继续飞行。

- 1 自动驾驶（如接通） 脱开
 - 2 自动油门（如接通） 脱开
 - 3 飞行指引电门（两个） 关
 - 4 按以下数据调置起落架收上的俯仰姿态和推力：
 襟翼放出 10°和 80%N1
 襟翼收上 4°和 75%N1
-

1. 吸收了所有故障排查经验
2. 紧急时刻不需要思考
3. 谁都可以执行，无门槛
4. 资料集中，查阅方便



- 运维团队积累的专家经验很难编码到算法模型中。通常，这些经验会被简化为阈值或复杂的规则，不仅难以维护，也难以传承。

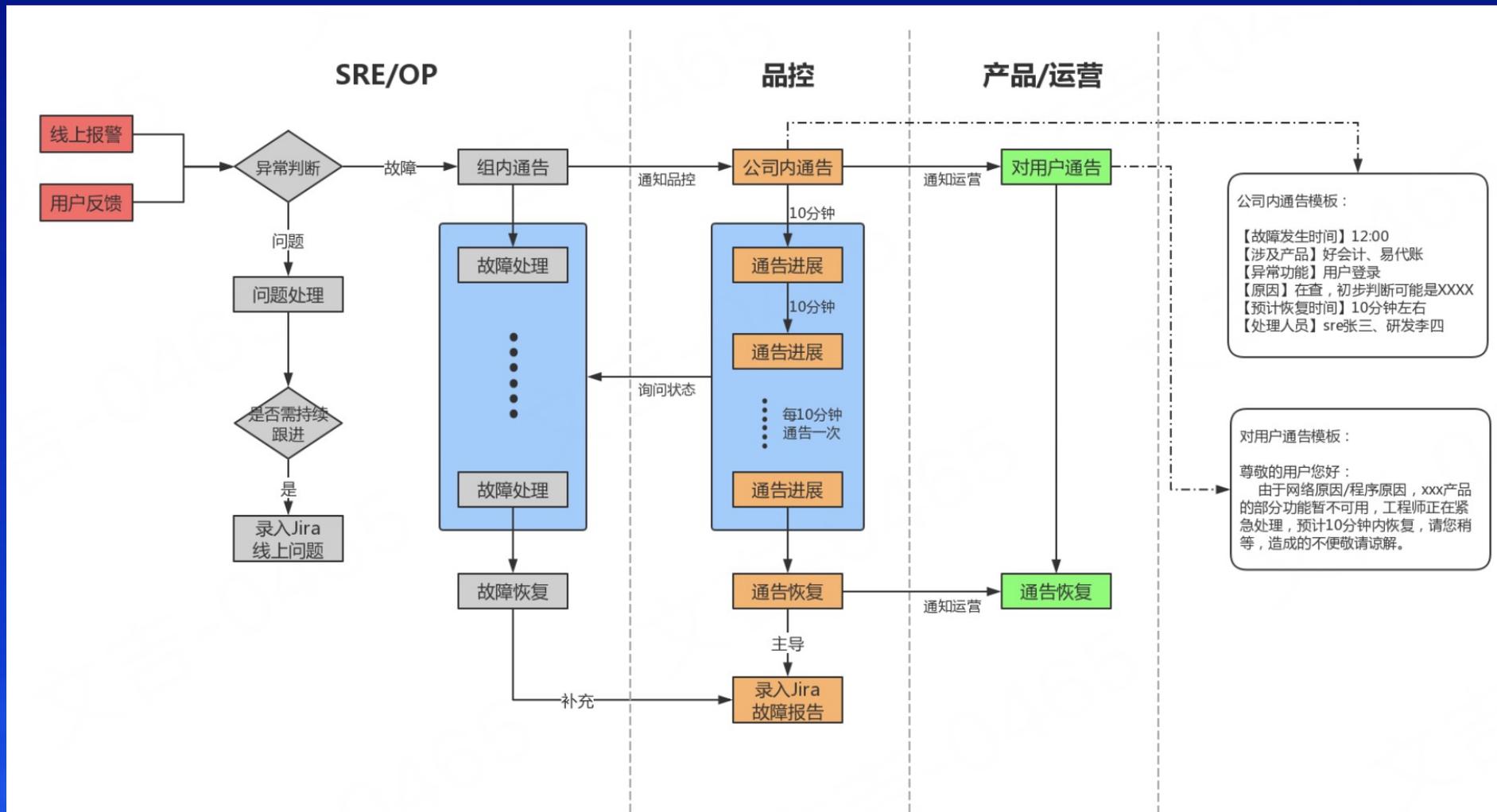
如何打造运维检查单？



- 接入和维护成本高，需要业务和算法团队深入理解业务逻辑和算法模型。
- 未遇到过的故障很难被解决，因为它们超出了模型的训练范围。
- 方案需要用户理解模型并精确地传递参数

▶ 可落地的协同处理流程

建立故障处理流程；
高效协同多个组织；



建立业务高峰期预防应急机制。

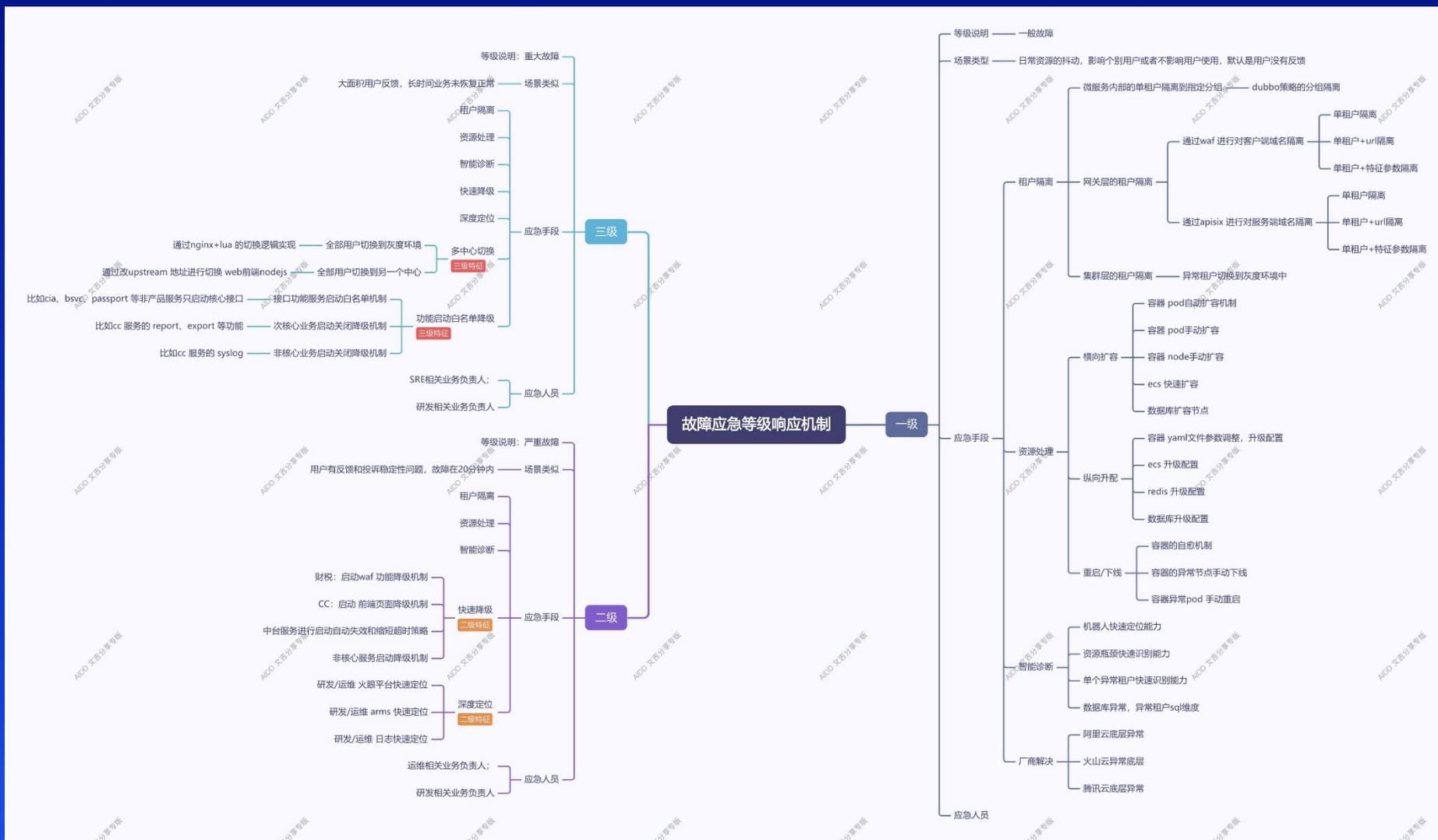
4.6 财税业务--大税期预防和应急机制

由 文吉创建, 最后修改于六月 30, 2024

- 预防篇：
 - 一级响应机制-处置方案：（高等风险处理机制）
 - 运维团队：
 - 1、遵循2-5-10原则，告警2分钟内须响应，5分钟内须定位到问题，10分钟内须对问题进行处理
 - 2、线上报警异常后，及时定位到问题，如问题不能立刻修复，须打开紧急预案，保障用户核心功能可用
 - 3、线上发生问题，及时打开快速定位面板，最短时间内定位问题
 - 4、如果是单个节点异常，应及时从集群摘除节点，防止影响扩大
 - 5、需要根据经验对问题严重性进行评估，如果达到一级响应级别，要及时通知研发，拉好问题处理专项群进行问题跟进与处理
 - 研发测试团队：
 - 1、ui自动化报警须及时跟进问题，进行解决，保证报警的准确性
 - 2、ui自动化大规模报警，须及时发出预警，拉好问题处理相关负责人的群进行专项问题跟进与处理
 - 3、研发同学在进入专项问题跟进群后，须集中精力处理线上问题，直到问题解决
 - 服务团队：
 - 1、共性问题及时反馈，及时发现和处理共性问题
 - 2、大规模共性异常问题走紧急流程，要求开发运维及时跟进处理
 - 二级响应机制-处置方案：（中等风险处理机制）
 - 运维团队：
 - 1、值班人员负责对告警面板进行巡检，不可存留告警信息，及时反馈和解决。前期是grafana面板，后期为监控中心的驾驶舱
 - 2、智能巡检提供的异常巡检及时验证。（24小时内的全部异常都要解决掉）
 - 3、容量评估预测风险：
 - 4、历史问题进行分析和提前预警：
 - 5、线上须研发跟进问题要及时创建jira给对应研发负责人进行跟进
 - 研发团队：
 - 1、ui自动化要求必须全部有效和正常运行。（24小时内的全部异常都要解决掉）
 - 2、研发团队及时分析异常日志信息。
 - 服务团队：

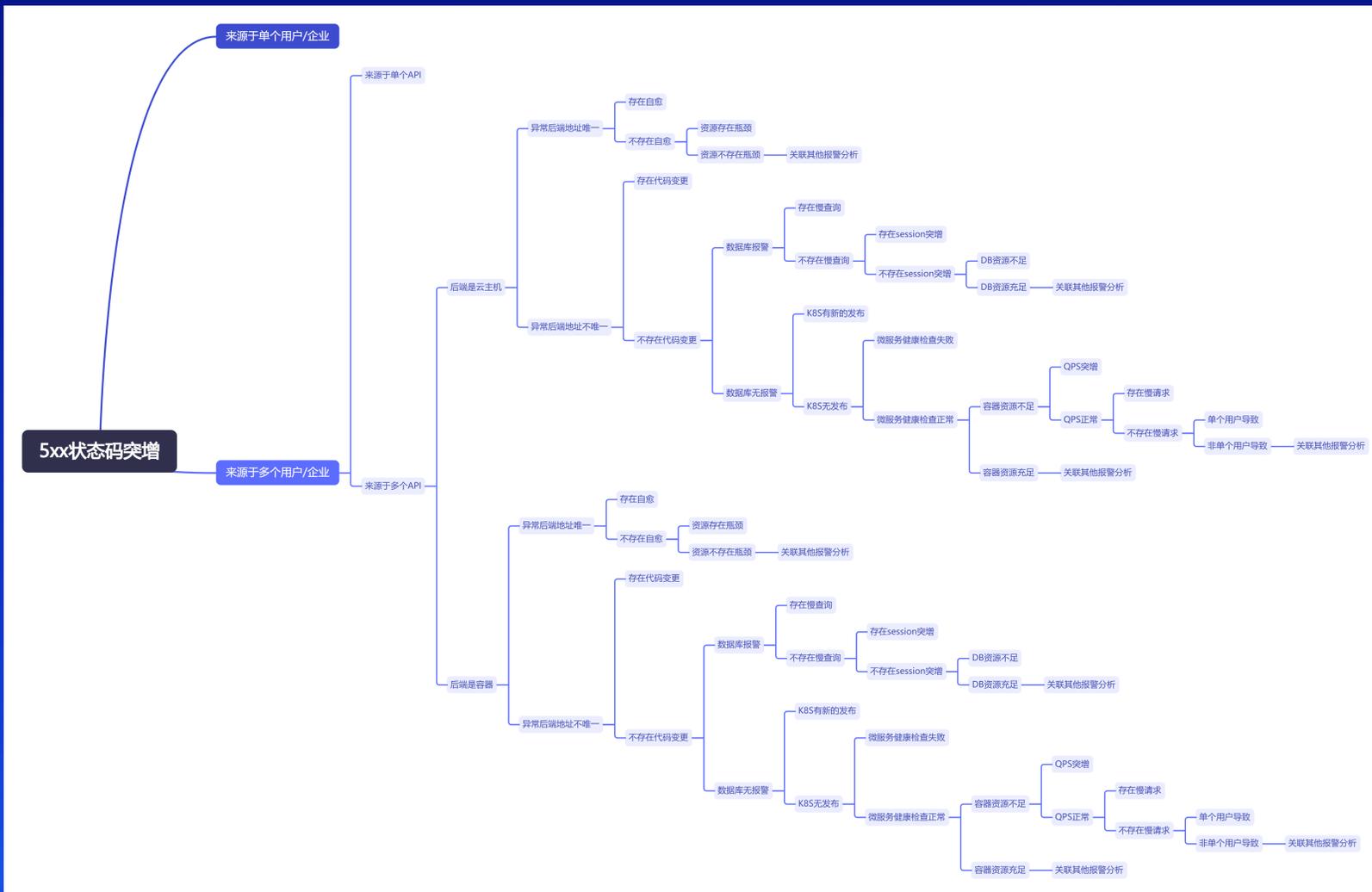
▶ 应急止损方法论——应急止损

建立应急止损操作流程和工具。



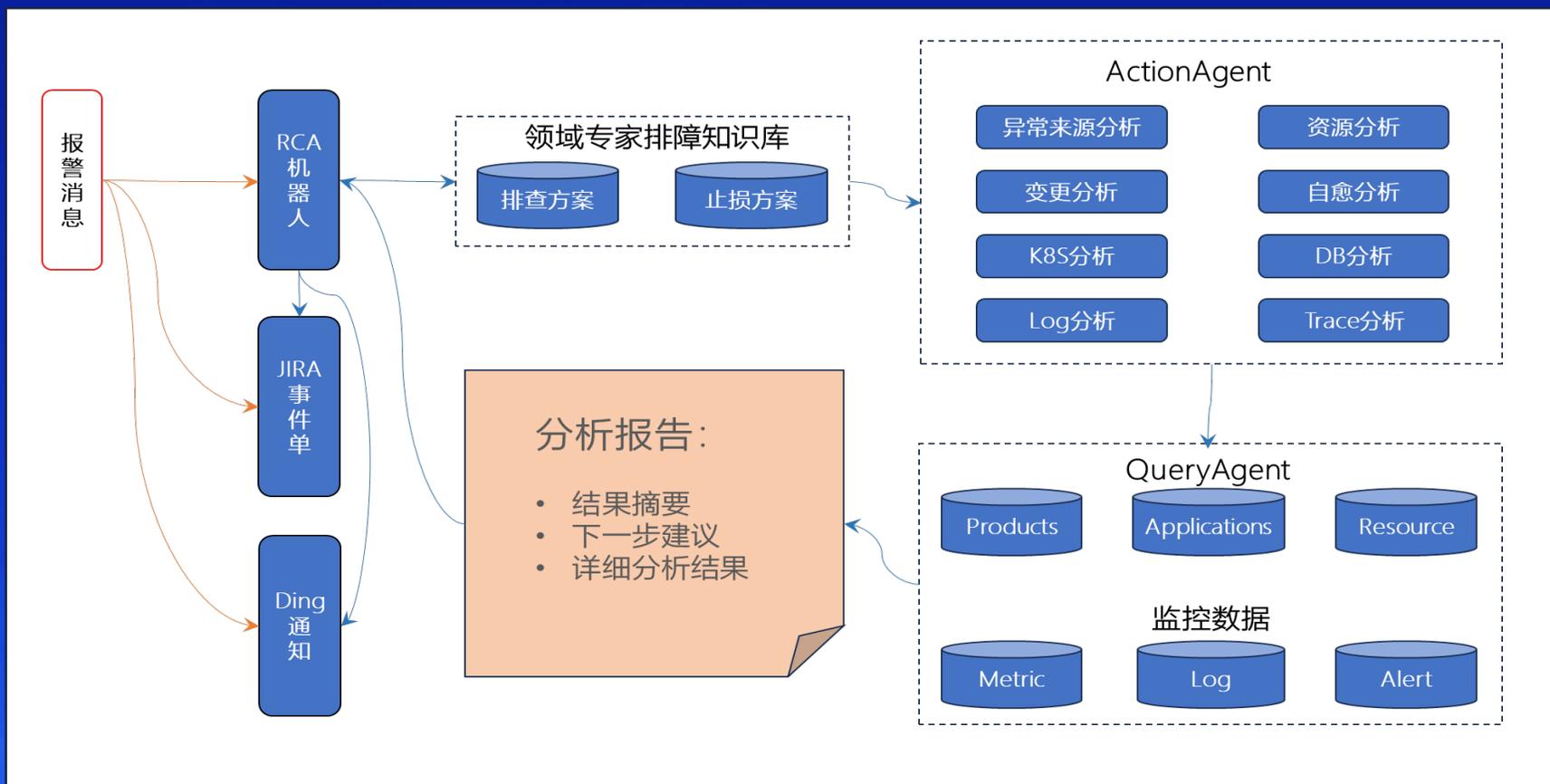
应急止损方法论——排障树

建立故障排查的专家经验排障树。



▶ 基于大语言模型的根因诊断 (RCA) Agent 框架

我们定义了一些工具和插件，这些工具和插件是用于出现故障时进行检测。除了工具和插件，我们还设计了工作流编排，可以自动化的故障处理流程。此外我们构建了一个知识库，它包含了历史故障数据、专家经验和故障处理策略，这些都是进行有效根因分析的关键资源。



将传统的针对多模态运维数据的异常检测方法变成工具（Agent），用户仅需维护指标项即可。

比如我们定义变更查询工具，该工具可以用于确定问题是否由线上变更导致。这样的工具有很多，一般都是基于运维专家日常的排障经验，可以是一个简单的脚本，也可以是一个API，或者是一个命令，这些工具可以完成故障排查过程中某一个环节的任务。

服务器资源瓶颈分析

The screenshot shows a REST client interface with a POST request to `h[redacted]/server_resource/`. The request body is a JSON object with the following content:

```
1 {
2   "ip": "192.168.1.100",
3   "alert_time": "2024-07-29 14:49:21"
4 }
```

The response is also in JSON format, showing a successful status (200 OK) and the following data:

```
1 {
2   "code": 1,
3   "reason": "服务器内存自愈. 命令: win_app_clean_mem 发起时间: 2024-07-29 06:44:02",
4   "data": {}
5 }
```

域名错误量upstream分布分析

The screenshot shows a REST client interface with a POST request to `h[redacted]ain/upstream/`. The request body is a JSON object with the following content:

```
1 {
2   "msg": "产品302状态码报警,当前302状态码条数821.0",
3   "domain": "h[redacted].com",
4   "alert_time": "2024-07-29 14:49:21"
5 }
```

The response is a detailed JSON object showing a status of 200 OK and the following data:

```
1 {
2   "code": 1,
3   "reason": "单upstream导致:192.168.1.100, 占比: 100.00%, 绝对值: 869.0",
4   "data": {
5     "upstream_count": 91,
6     "upstreams": {
7       "192.168.1.100:80": "6015",
8       "192.168.1.100:80": "99",
9       "192.168.1.100:80": "97",
10      "192.168.1.100:80": "8054",
11      "192.168.1.100:80": "1114",
12      "192.168.1.100:80": "1851",
13      "192.168.1.100:80": "1231"
14     }
15   }
16 }
```

构建工作流，我们在prompt和文档中预先设置了不同报警的分析流程，即应该先后检查哪些数据，从而得出结论。

这个工作流类似飞机检查单（SOP），不同的现象对应不同的检查项，类似一个树状结构，最终一定会递归找到一个叶子节点然后返回。比如当某个域名出现5xx状态码报警，我们需要先判断这些状态码是否来源于同一个用户的请求，再判断这些请求是否都打到了同一个upstream节点，后端承载流量的微服务、容器和node是否存在问题，最后再检查是否是第三方依赖存在问题等。



这是一种妥协方案，我们可以选择对通用大语言模型进行训练，它能够根据用户的SOP文档直接生成工作流，但是大模型训练的成本是非常高的，一方面是资源成本，另一方面是对大模型人才需求的成本。

PART 04

具体实现/技术实践

▶▶ 数据治理——CMDDB建设

将资产标签化，将标签目录化，得到完整的产品六级目录，既有业务信息，又有资产实例的关联关系，每种资源都拥有自己的身份证号：六级目录。这是AIOps落地的基石。

The screenshot displays the CMDDB interface with three main components:

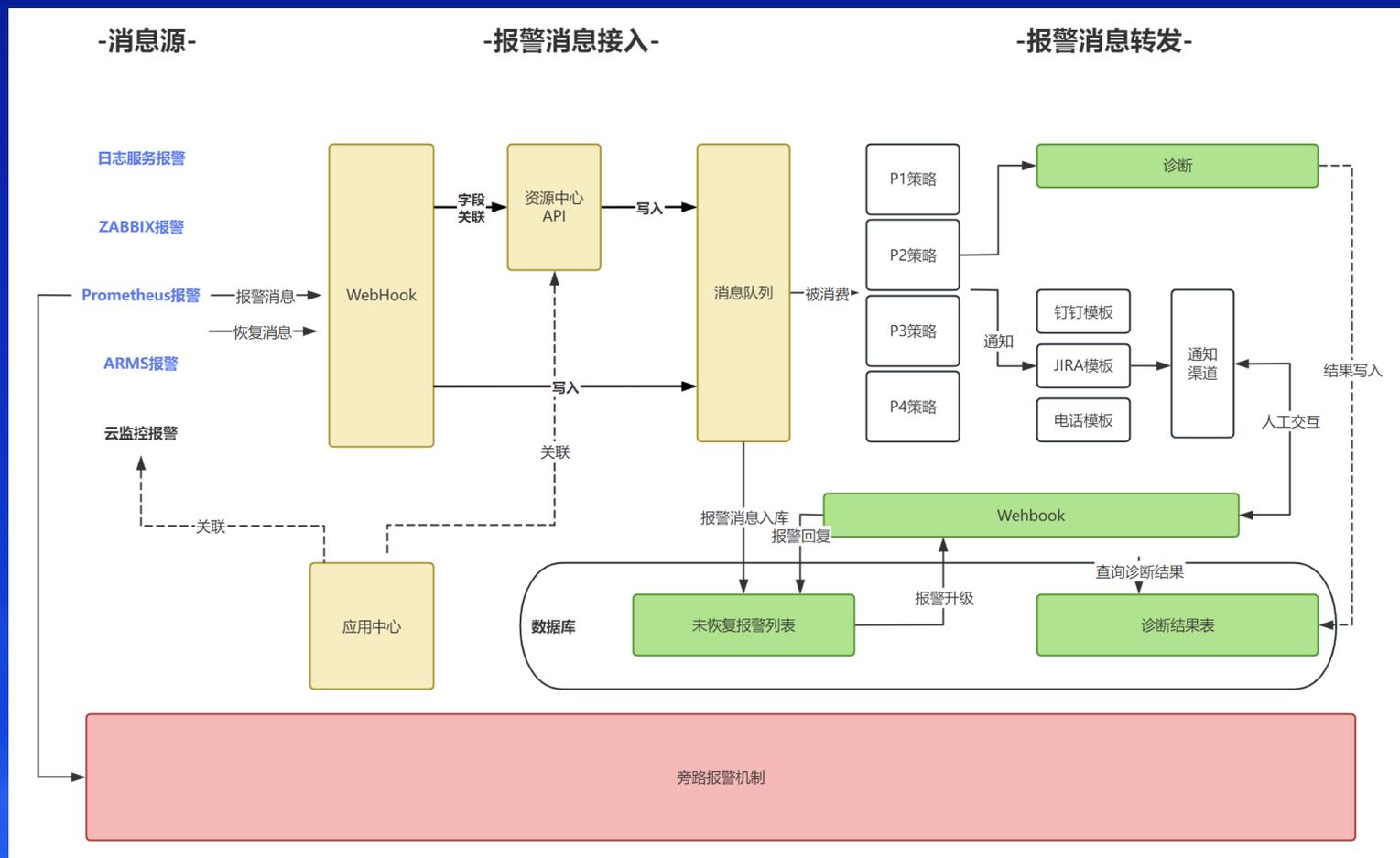
- Left Panel (Directory Structure):** A tree view under '应用中心' (Application Center) showing a hierarchy: '8125' > 'T+cloud' > '研发测试中心' > '测试环境' > '普及版' > '普及版集群' > 'C8125-h2t_pop_t7'.
- Middle Panel (Bar Chart):** Titled '目录下资源统计' (Resource Statistics under Directory). The chart shows the count of resources for different categories: '主机' (Host) has 7, '域名' (Domain) has 1, 'RDS' has 3, and 'Redis' has 3.
- Right Panel (Resource Table):** A table titled '主机' (Host) listing resource instances. The table has columns for IP, 实例ID/名称 (Instance ID/Name), 状态 (Status), 系统/配置 (System/Config), 账号平台 (Account Platform), 可用区 (Availability Zone), 付费方式 (Billing Method), 创建时间 (Creation Time), 标签 (Tags), 使用人 (User), and 操作 (Action).

资源类型	数量
主机	7
域名	1
RDS	3
Redis	3

IP	实例ID/名称	状态	系统/配置	账号平台	可用区	付费方式	创建时间	标签	使用人	操作
...	...	运行中	4C/8G/G	阿里云	cn-beijing-k	按量付费	2024-01-07 16:28:00	详情
...	...	运行中	4C/8G/G	阿里云	cn-beijing-k	按量付费	2024-01-07 16:28:00	详情
...	...	运行中	4C/8G/G	阿里云	cn-beijing-k	按量付费	2024-01-07 16:28:00	详情

数据治理——监控统一

来源于不同监控工具的报警必须满足最小字段集合，这样以来所有的报警都能标准的关联到具体的业务、产品，从而关联出所有的资源、中间件等信息。同时我们也完成了CMDB的自动化维护，形成了包含业务、基础资源、人员、代码仓库、配置等关联关系的大型数据字典，本身也为webUI提供了许多API，这些API都将作为Agent被注册。



数据治理——监控统一

云擎平台 全局搜索菜单/资源 值班表

监控中心 **驾驶舱**

P1-Disaster 0个

P2-High 0个

P3-Info 171个

P4-Notclassified 218个

全量未恢复报警 P1 x P2 x 报警来源

报警来源	报警时间	持续时间	报警级别	报警内容	资源类型	报警实例	产品线	中心	环境	cid	操作	ACK
暂无数据												

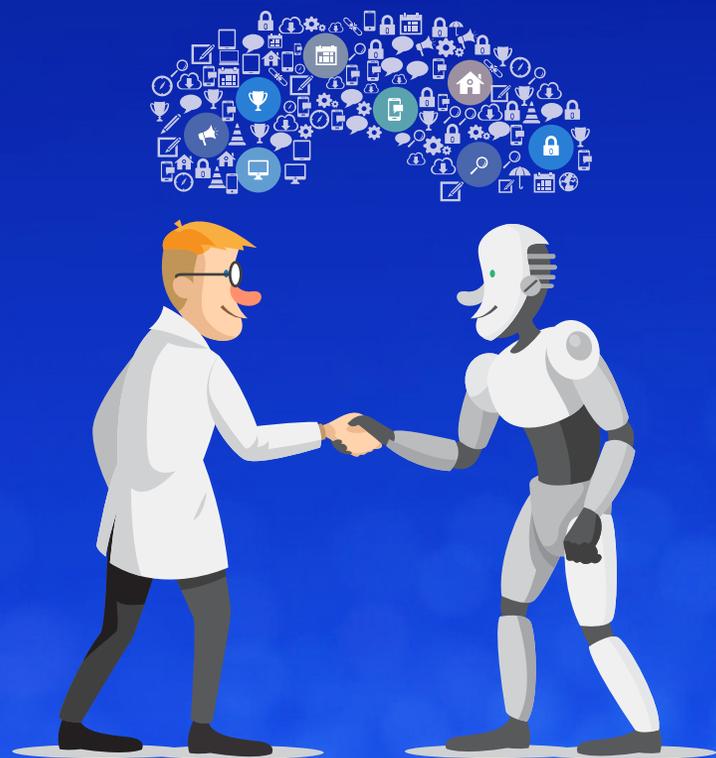
历史报警流水 近8小时 报警级别 报警来源

时间	报警级别	资源类型	报警实例	报警内容	产品线	中心	环境	CID	报警状态	持续时间	报警来源	手动关闭 / ACK
2024-08-16 16:46:24	P4	ecs	172.16.43.1...	win内存使用...	T+cloud	中心二	线上环境	C2053-devh...	firing	27	ZABBIX	No / No
2024-08-16 16:43:05	P4	pod	c0850-hsy-...	容器 CPU 资...	好生意-好业财	中心一	线上环境	C0850-publi...	firing	226	Prometheus	No / No

▶ SOP定义——专家经验的沉淀

针对每种现象，我们都梳理了运维专家的排障脑图，将故障排查过程固化下来。

根节点表示现象（报警），分支节点代表一个分析操作，每个分支节点会再分化出是和否两个分支，直到最外围的叶子节点，无法再进行下一步分析为止。外围节点会有两种状态：根因和非根因



► 工具构建之查询类Agent

查询类Agent融合了CMDB（产品、应用、资源的关联关系）、IT资产清单、CICD配置、config数据的查询。查询类Agent的还包含了历史故障单的查询，让AI具备寻找历史相似事件的能力。

Schema | [查看 OpenAPI-Swagger 规范](#) + 从 URL 中导入 例子

```
"servers": [
  {
    "url": "https://[redacted].chanjet.com"
  }
],
"paths": {
  "[redacted]_ost_info/": {
    "get": {
      "summary": "Get machine host info",
      "parameters": [
        {
          "name": "pageSize",
          "in": "query",
          "required": true,
          "schema": {
```

可用工具

名称	描述	方法	路径	操作
assetecsmachine_host_info_get	Get machine host info	get	[redacted]_ost_info/	测试
tagV2show_ctags_get	获取服务器的应用标签	get	[redacted]_show_ctags/	测试



工具构建之动作类Agent

动作类的Agent就是前文提到的，对于排障脑图中某个具体节点的对象的分析过程，我们可以非常原子化的进行这些Agent的定义，比如下面是我们定义的一些Agent

服务器资源瓶颈检查

```
1 POST http://.../api/monitor/server-resource-bottleneck/
2 {
3   "ip": "172.16.163.169",
4   "alert_time": "2024-07-17 13:00:33"
5 }
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
```

```
1 {
2   "code": 0,
3   "reason": "云主机资源正常",
4   "improvement": "云主机资源正常",
5   "affected_products": "",
6   "data": {}
7 }
```

异常访问来源检查

```
1 POST http://.../api/monitor/abnormal-access-source/
2 {
3   "msg": "产品响应时间(p90)波动,当前平均响应时间3.28s",
4   "instance": "devops.rd.chanjet.com",
5   "alert_time": "2024-07-16 17:28:42"
6 }
7
```

```
1 {
2   "code": 1,
3   "reason": "单接口导致:/prod-api/api/productEnv/deploy/save, 占比: 66.67%, 绝对值: 8.0",
4   "improvement": "请检查接口/prod-api/api/productEnv/deploy/save, 优化接口",
5   "affected_products": "其他",
6   "data": {}
7 }
```

异常访问upstream检查

```
1 POST http://.../api/monitor/abnormal-access-upstream/
2 {
3   "msg": "产品302状态码报警,当前302状态码数821.0",
4   "domain": "...",
5   "alert_time": "2024-07-29 14:49:21"
6 }
7
```

```
1 {
2   "code": 1,
3   "reason": "单Upstream导致:100.00%, 占比: 100.00%, 绝对值: 869.0",
4   "data": {
5     "upstream_count": 91,
6     "upstreams": {
7       "5932",
8       "6226",
9       "3146",
10      "7700",
11      "1555",
12      "8358",
13      "6873",
14      "6095",
15      "6901",
16      "7511"
17     }
18   }
19 }
```

The screenshot shows a GitLab repository interface. The main content is a file named 'prompt' (4.14 KB) with a 'Web IDE' button. The file content is a markdown document for an AI root cause analysis tool. The document includes a description of the tool's purpose, input parameters, analysis steps, and output requirements.

```
1 你是一个根因分析专家，专注于解析报警信息以确定问题根本原因。报警信息将以JSON格式提供，需仔细审查以下核
2  - `instance_type`：指示报警涉及的实例类型，涵盖domain、ecs、log、数据库、K8S组件等。
3  - `instance`：具体的报警实例详情。
4  - `product`、`center`、`env`、`cid`：这些字段共同标识特定的业务环境。
5  - `pd`：表示特殊产品线中的子产品名称；若pd为空，则无需将其纳入API查询参数
6
7  ### 针对提供的instance_type进行系统化分析流程，遵循以下步骤以精准定位问题根因，并按规范格式返回结果
8  1. **代码变更优先检查**： - 首先检查过去30分钟内是否有代码变更。若有变更，收集变更内容及变更人作为根
9  2. **instance_type特定分析**：
10 - **domain实例**：
11 - 查看异常来源分布，异常来源分布结果会有以下四类：
12 - 如果是单IP导致，直接返回作为根因。
13 - 或者是单IP+URL导致，直接返回作为根因。
14 - 如果是单接口导致，请勿将此作为根因返回，而是继续下面的步骤。
15 - 或者其他情况，请勿将此作为根因返回，而是继续下面的步骤。
16 - 若异常来源分布未明确根因，分析后端异常分布；若非单一upstream，则总结当前分析并通知未找到根因。
17 - 继续对单一upstream执行服务器资源分析，发现自愈行为或资源瓶颈即为根因；否则，告知未找到根因。
18 - **ecs实例**：
19 - 直接进行服务器资源分析，发现自愈或资源瓶颈作为根因；否则，告知未找到根因。
20 - **其他instance_type**：
21 - 根据实际情况推理最合适的分析路径。
22 3. **终极分析步骤**：若上述所有步骤未能确定根因，无论instance_type是什么，都进行关联报警分析。分析内
23 - 关联报警为空，此时告知用户该业务线10分钟内无其他报警。
24 - 关联报警不为空，此时告知用户10分钟内有多少条报警，并列举每条报警的关键信息（什么资源，发生了什么事）。
25
26 ### 分析任务要求 你的响应必须结构严谨，包含以下组成部分：
27 1. **result**（布尔值）：表明是否成功识别出根因。`True`表示已发现，`False`表示未发现。
28 2. **root_cause**：明确指出根因（如找到）。
29 3. **analysis_abstract**：中文总结分析过程，包括检查的要素及得出的结论，不要使用换行符等特殊字符。
30 4. **improvement**：提出针对发现的问题的具体改进建议（中文表述）。
31 5. **markdown_result**：整合上述所有信息为Markdown格式的字符串，遵循特定风格指南，确保颜色编码正
32 #### Markdown格式示例 ``## <font color=#0096fe >AI大模型智能诊断结果 </font> \n ### <font
```

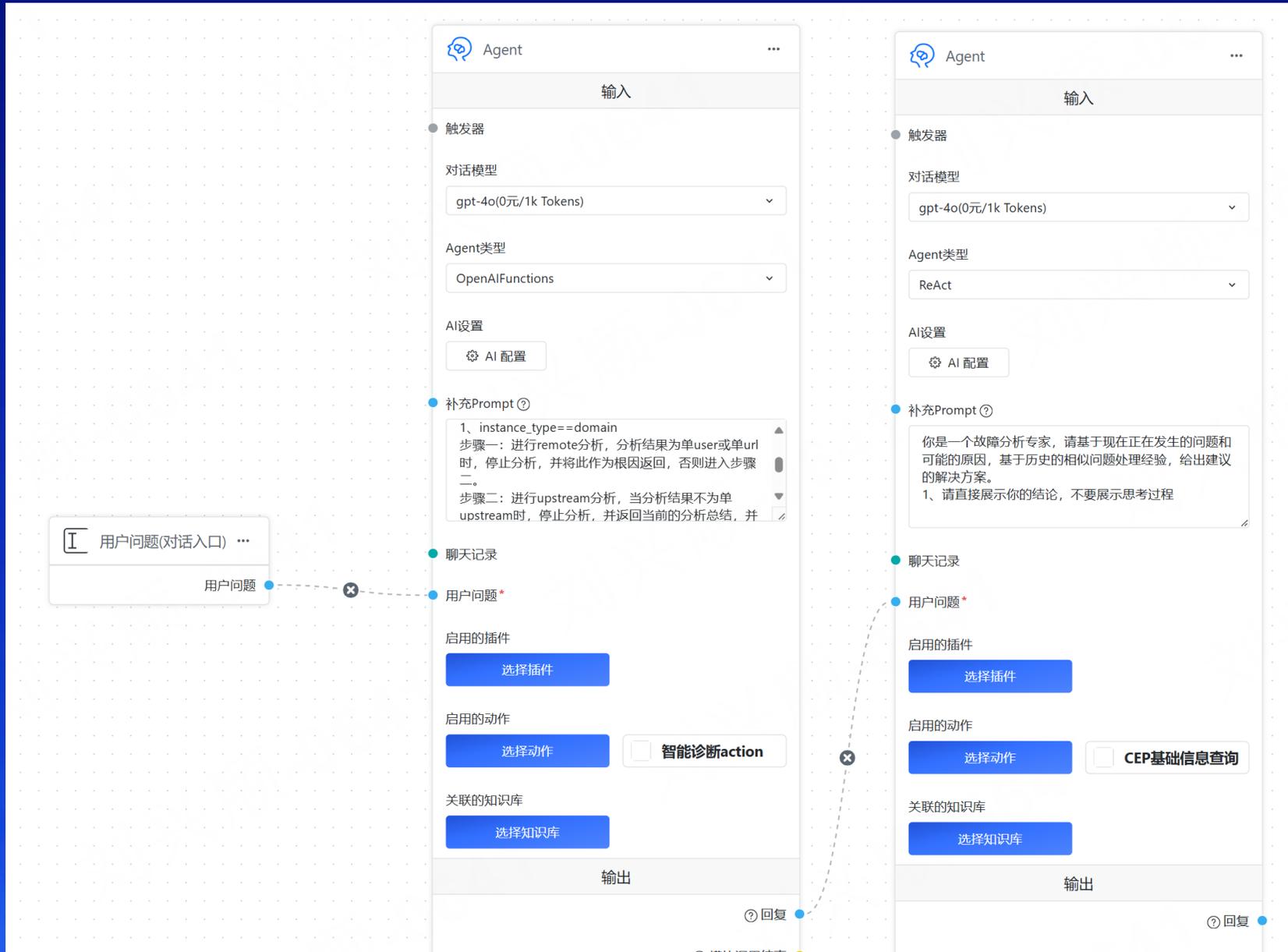
The screenshot shows the '编辑动作' (Edit Action) dialog box. It contains the following information:

- 名称** (Name): 智能诊断action
- Schema**: 查看 OpenAPI-Swagger 规范 (View OpenAPI-Swagger Specification)
- Schema Content**:

```
openapi: 3.1.0
info:
  title: Root_Cause_Analysis_API
  description: API for monitoring and analyzing root causes of alerts
  version: 1.0.0
servers:
  - url: https://...jet.com
    description: Production server
paths:
  /monitor/root_cause_analysis/ai/domain:
    post:
      operationId: remote_root_cause_analysis
      summary: Analyze the root cause of alerts pertaining to either a single user or domain,
        exclusively when the instance_type is designated as domain
      description: Analyze the root cause of alerts pertaining to either a single user or domain
```
- 可用工具** (Available Tools) table:

名称	描述	方法	路径
remote_root_cause_analysis	Analyze the root cause of alerts pertaining to either a single user or domain,	post	/monitor/root_cause_analysis/ai/domain
- 操作** (Actions): 移除 (Remove), 取消 (Cancel), 保存 (Save)

▶ 流程编排



▶▶ 效果升级

降低编码的复杂性和成本输出

云掌平台

全局搜索菜单

事件中心 策略管理 × 模板管理 × 钉钉模板-修... ×

策略配置

- 策略管理
- 策略组管理

配置中心

- 数据接入
- 模板管理**
- 报警屏蔽规则
- 默认接收规则

* 模板类型: 钉钉

* 模板标题: \${content}

模板备注: 模板备注信息

* 模板内容:

```
- **报警日期:** ${happen_time} \n\r - **产品:** ${product_app} \n\r  
- **应用:** ${product_app} \n\r - **报警内容:** <font  
color=#2F7FFE> ${content} </font> \n\r - **中心:** ${idc_name}  
${cloudsys} \n\r - **报警IP:** ${ip} \n\r ----- \n\r - **运维:**  
${operator} \n\r - **研发:** ${developer} \n\r - **事件:**  
${_ruleCaseName} \n\r - **应用ID:** ${cid} \n\r - **等级:**  
${level_name} \n\r ----- \n\r # <font color=#FF1493>报警回  
复</font> \n\r - [已收到](dtmd://dingtalkclient/sendMessage?
```

模板验证

验证

我们目前已经实现了所有线上报警的自动分析，目前根因的召回率已经超过了50%，随着Agent和流程编排的完善，召回率还会逐渐提升。

对于成功召回根因的报警，机器人会自动关闭报警工单，同时支持钉群交互，形成闭环。

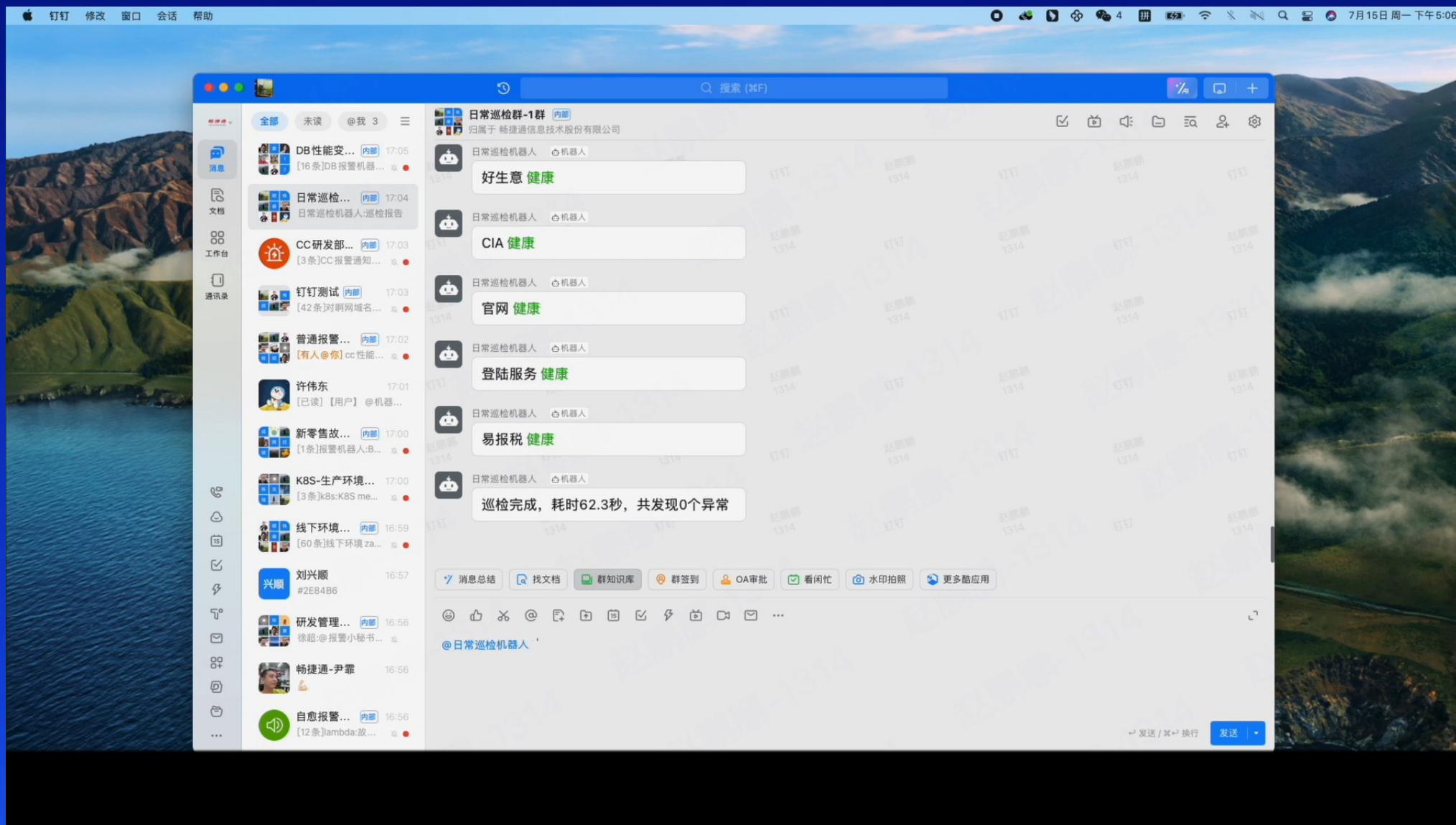


▶ RCA效果展示

The image displays a software interface for configuring a Root Cause Analysis (RCA) agent. It is divided into three main sections:

- Component List (Left):** A sidebar titled "高级编排组件列表" (Advanced编排Component List) containing various modules such as "输入引导 [5]", "语言模型 [3]", "知识增强 [3]", "基础组件 [4]", "智能服务 [3]", "其他功能 [3]", "系统插件 [0]", "用户插件 [0]", and "动作请求 [7]".
- Workflow Diagram (Center):** A visual flowchart showing a "用户问题(对话入口)" (User Question) component connected to an "Agent" component.
- Agent Configuration Panel (Right):** A detailed configuration window for the "Agent" component, organized into "输入" (Input) and "输出" (Output) sections.
 - Input Section:**
 - 触发器 (Trigger):** Includes "对话模型" (Dialog Model) set to "gpt-4o(0元/1k Tokens)" and "Agent类型" (Agent Type) set to "OpenAIFunctions".
 - AI设置 (AI Settings):** Includes an "AI配置" (AI Configuration) button.
 - 补充Prompt (Supplement Prompt):** A text area containing a system prompt: "你是一个根因分析专家，需要对报警进行分析。你只会接收到报警消息，报警消息是json格式，需要注意报警中携带的重要信息：1、instance_type: 代表这个报警的实例类型，有可能是domain (域名)、ecs (服务器)、log (错误日..."
 - 聊天记录 (Chat History):** Includes a "用户问题*" (User Question) field.
 - 启用的插件 (Enabled Plugins):** A "选择插件" (Select Plugin) button.
 - 启用的动作 (Enabled Actions):** A "选择动作" (Select Action) button and a checkbox for "智能诊断action" (Smart Diagnosis Action).
 - 关联的知识库 (Associated Knowledge Base):** A "选择知识库" (Select Knowledge Base) button.
 - Output Section:** Includes "回复" (Reply) and "模块调用结束" (Module Call End) options.
- Debug Preview (Far Right):** A "调试预览" (Debug Preview) panel with a "提示词预览" (Prompt Preview) section and a "用户输入" (User Input) section for testing the agent's response.

▶ 我们更进一步的尝试

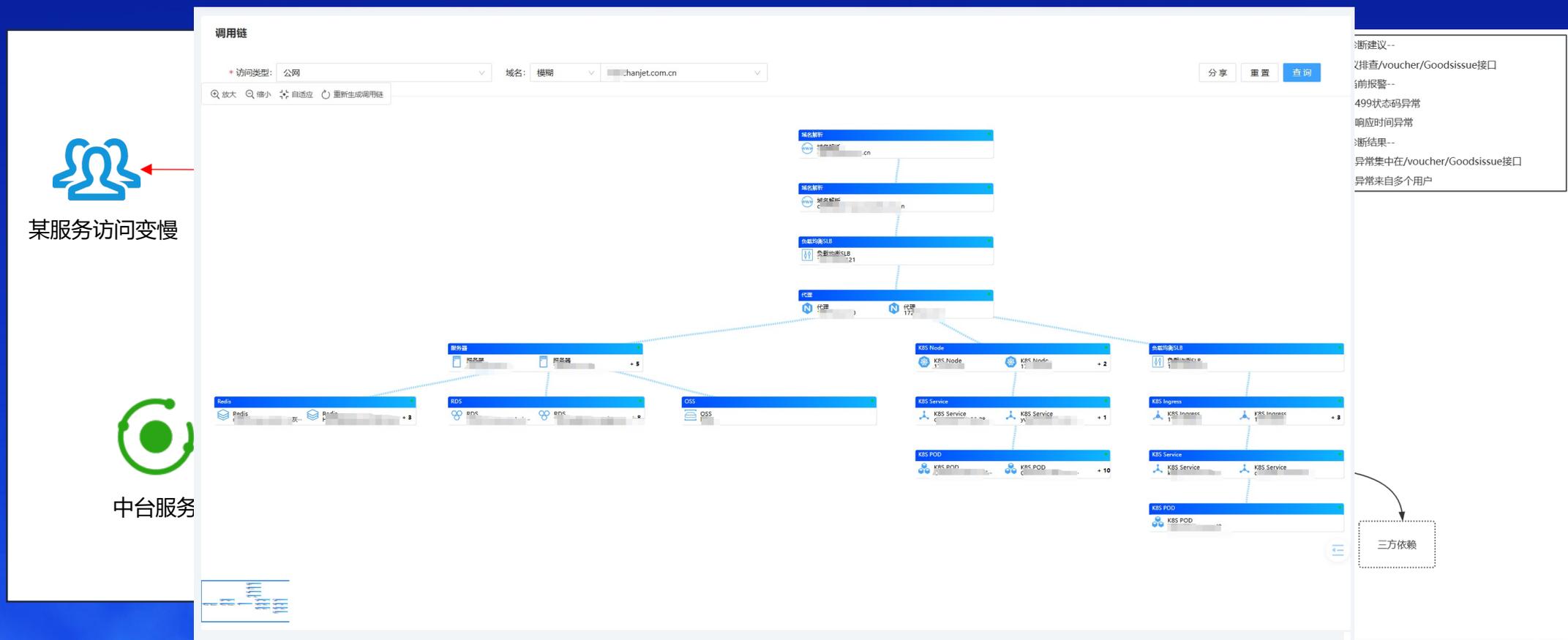


PART 05

总结与展望

▶ 方案总结——望、闻、问、切

本方案通过构建根因排查逻辑树、建立统一的报警字段集规范，建立多模态Agent集合，充分调度AI大模型文本推理的能力，对报警通知、报警事件单和根因分析过程进行了整合，实现了报警的自动化分析，整体耗时在1分钟以内，对于90%常见的报警都能分析出根因所在，即便是10%的不常见报警，也能完成分析过程，运维人员无需重复分析，为应急止损和故障定位争取了更多时间，保证了业务稳定性。



► 大模型时代，做AI的主人

大模型技术诞生之后，已经颠覆了IT从业者的工作和思维习惯，大家的技术水平差距已经被大模型抹平了，而善于思考，能把问题想明白变这个事情，变得更加重要了。

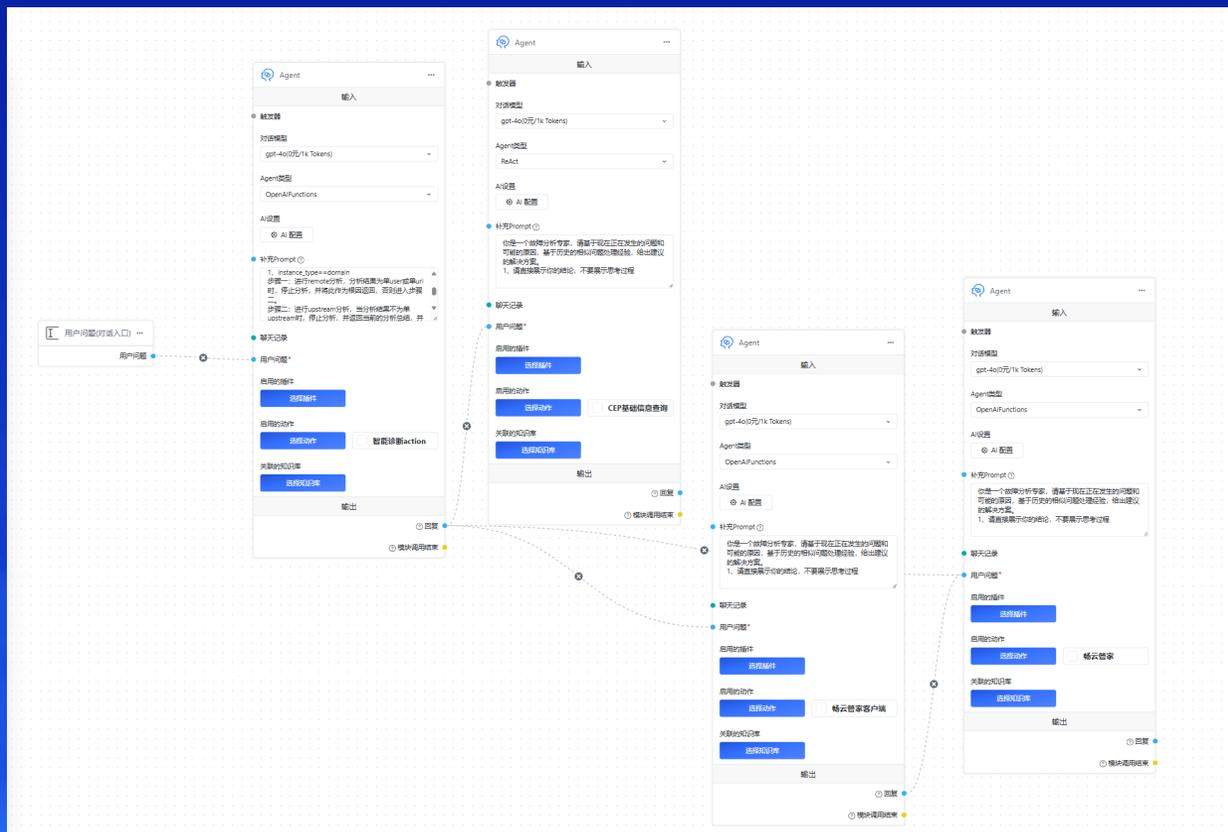
其实用大模型技术完成推理+检索实现RCA应用的过程，其实就是在prompt或者知识库中定义了各种if else的逻辑，理论上只要能说得清的逻辑，就可以通过传统编码的方式实现，我们为什么要承担AI大模型偶尔“一本正经胡说八道”的风险？

把问题想明白说清楚 > 专业技术强大

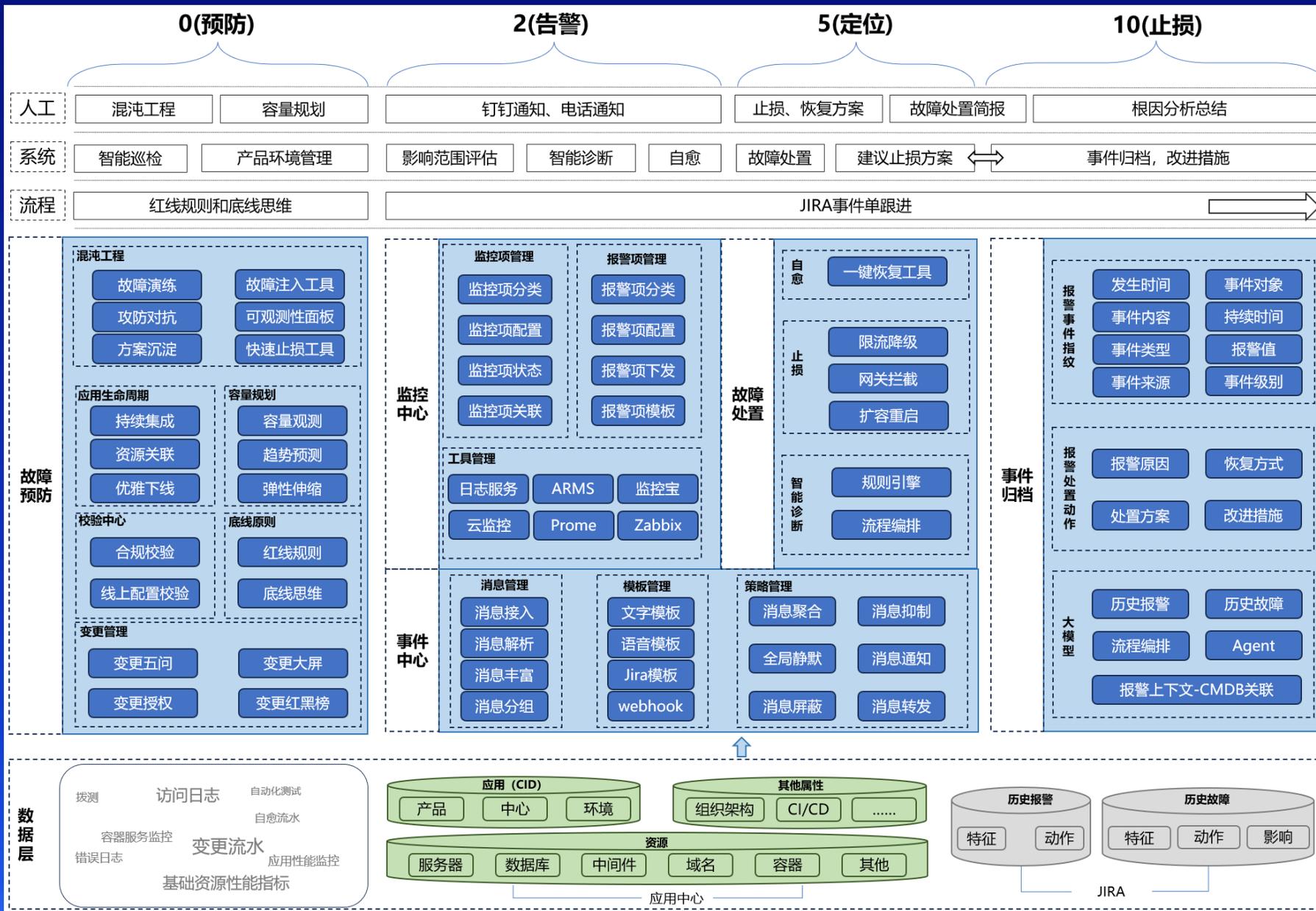


▶ 我们接下来会做的事情

- 更多 workflows: 让AI串联更多工作流程, 比如监控、巡检、故障止损、智能容量预防、智能风险识别等
- 工作流插件化: 让这些工作流变成插件, 从而可以在大模型应用中进行调用
- 大总管的模式: 面向对话框工作, 所有的交互不再需要设计webUI, 也不再需要设计问题, 简化开发的过程, 充分释放AI的能力



▶ 我们接下来会做的事情



► 未来的发展趋势

关键词：全文检索、逻辑推理、低代码

目标：

- 1、基于大模型的，减少知识获取难度
- 2、利用大模型擅于汇总、总结的能力



query



XXX当前有少量500状态码，都来自同一个用户的请求，结合历史判断目前正常.....

科技生态圈峰会 + 深度研习



—1000+ 技术团队的选择



上海站

K+ 全球软件研发行业创新峰会

时间: 2024.06.21-22



敦煌站

K+ 思考周®研习社

时间: 2024.10.17-19



香港站

K+ 思考周®研习社

时间: 2024.11.10-12



K+峰会详情



上海站

Ai+研发数字峰会

时间: 2024.05.17-18



北京站

Ai+研发数字峰会

时间: 2024.08.16-17



深圳站

Ai+研发数字峰会

时间: 2024.11.08-09



AiDD峰会详情



2024 AI+研发数字峰会

AI+ Development Digital summit

深圳站 11/08-09

AI 驱动研发变革 促进企业降本增效

2024深圳站-议题设置

AI+产品线	LLM驱动产品创新	LLM驱动需求与业务分析	AI驱动设计与用户体验
AI+开发线	AI 原生应用开发框架与技术	AI Agents在研发落地实践	LLM驱动编程与单测
AI+测试线	LLM驱动测试分析与设计	基于LLM生成测试脚本与数据	LLM和AI应用的评测
AI+工程线	AI+DevOps 与工具 (LLM 时代的平台工程)	大模型对齐与安全	端侧大模型与云端协同
AI+领域线	领域大模型 SFT 与优化	知识增强与数据智能	大厂专场

扫描右侧二维码
查看更多会议详情



早鸟票限时抢购中 (截止到9月30日)

¥3680

早鸟票

¥2800

学生票

THANKS

