



2023 CCF国际AIOps挑战赛决赛  
暨“大模型时代的AIOps”研讨会

# 基于大模型和多AGENT协同的 自主决策、自动修复运维方案

轻舟已过万重山

陈东辉、陈润、陈子扬、郭广宇、胡彬、吕颖、谢志鹏、袁俊、张楠、张曦

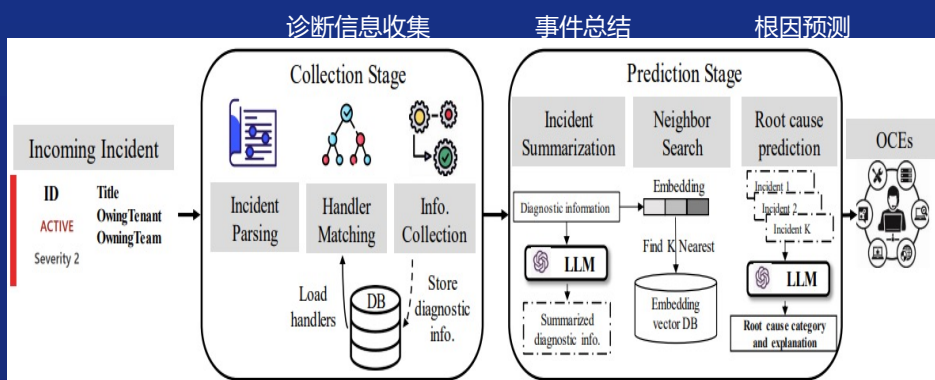
华为技术有限公司

主办单位：中国计算机学会（CCF）、清华大学、中国建设银行股份有限公司、南开大学

承办单位：中国计算机学会互联网专委会、清华大学计算机科学与技术系、中国建设银行股份有限公司运营数据中心、南开大学软件学院、北京必示科技有限公司

赞助单位：华为技术有限公司、国网宁夏电力有限公司电力科学研究院、软通动力信息技术（集团）股份有限公司

## 云厂商利用大模型对运维事故进行根因定位并给出故障缓解措施建议



7成以上运维人员对LLM的分析结果满意(>3分)<sup>[2]</sup>

Score	# of incident owners	In percent (%) of total
5	2	7.41
4	9	33.33
3	8	29.63
2	6	22.22
1	2	7.41

合计  
70+%

## 针对运维领域海量知识快速获取、辅助诊断和故障分析能力

- 将LLM较为广泛的知识储备（横向能力）与运维领域专业知识（运维垂域）相结合
- 运维工具繁多，利用LLM+多智能体协同使能运维自动驾驶

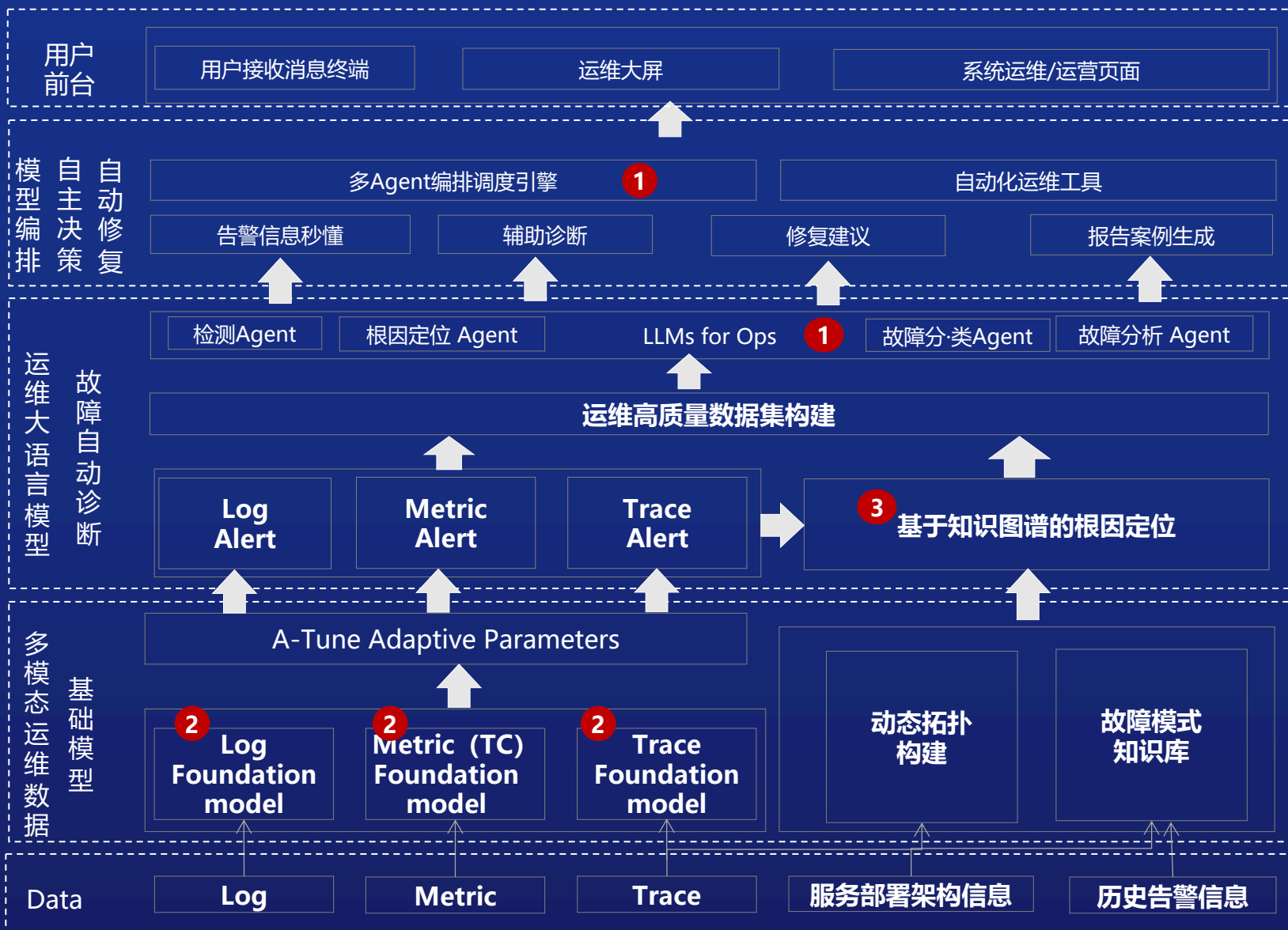
## 针对多模态数据进行快速高效准确的异常检测能力

- 多损益、多类型数据：Log、Metric、业务黄金指标、Trace
- 数据异常不等于故障发生，识别可能造成运维故障的数据异常波动

## 针对多源复杂部署的运维数据进行快速根因定位能力

- 告警繁多、区分故障传播节点和故障根因节点
- 历史排障经验积累数字化，形成知识库

# 基于大模型和多Agent协同自主决策、自动修复运维方案



## 关键技术

### 1 基于多Agent协同的编排调度

- ① 梳理异常检测告警信息和运维故障模式，形成运维领域高质量数据集，使用finetuning+外挂向量数据库的方案，使得LLM具备运维领域故障分析定位能力；
- ② 使用主管LLM对异常问题与子领域Agent进行桥接，多个子领域Agent协同工作，提升运维效率

### 2 更全能的多模态数据异常检测基础模型

- ① 针对成百上千条metric序列采取一定维度聚合的方法，化繁为简，一旦检测出异常，再细化分析关键指标的关键表现；
- ② 针对半结构化和非结构化日志采取针对性的解析方式：前者注重模板提取，后使用uADR/sADR进行异常检测，后者注重语义理解，使用BigLog预训练模型结合Deep SVDD+SAD进行异常检测
- ③ 针对结构化明显的trace数据，抽取分离调用关系并定义节点关键数据，通过转化时间序列有效识别异常节点，并通过拓扑关系分析帮助推测可能根因节点及故障传播方向

### 3 基于知识图谱的根因定位

- ① 利用DBSCAN从时间维度聚类，形成异常事件
- ② 利用Trace数据实时生成拓扑结构图
- ③ 根据故障模式知识库分析故障传播链路



## Metric异常检测

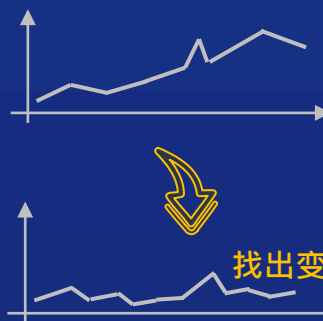
### 1 按采样频率分组

- 采样频率1
- 采样频率2
- 采样频率3
- ...



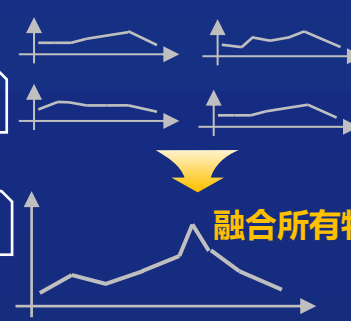
### 2 提取指标变化特征

- 原始指标曲线
- 指标标准化
- 窗口差分取绝对值
- 指标变化特征曲线



### 3 同类型指标时序融合

- 指标变化特征曲线
- 融合曲线
- 融合所有特征点



### 4 异常检测与时序聚类

- 融合曲线
- 多个异常检测器
- 针对单个融合指标的DBSCAN时序聚类
- 异常时间戳整体DBSCAN时序聚类
- 系统异常发现!



LLM主管  
Agent决策



### 5 检测Agent

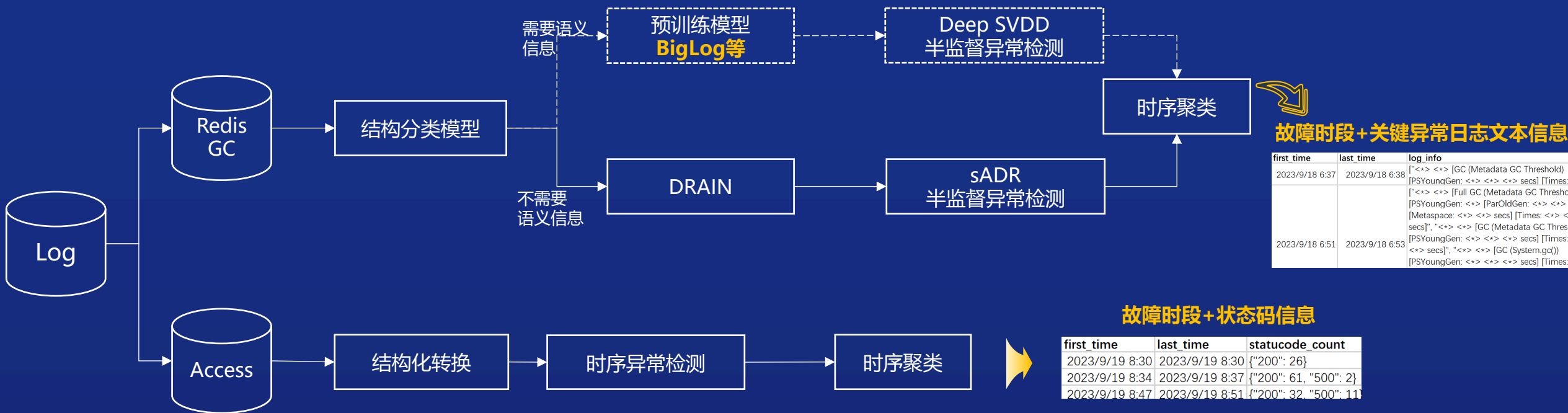
- 根因单元的异常对应类型  
全部指标变化特征曲线
- 多个异常检测器
- 异常指标



cnid	ip	device	进程名	first time	last time	cluster	first timestamp	last timestamp
Message_02	system.disk.io, system.dfs.used, system.fs...	huy	2.00045	2022-09-18 06:20:00	2022-09-18 06:40:00	1	1.05215e+00	1.05207e+00
Message_03	system.net.bytes_sent, system.net.packets_rcv...	enb	1.471271	2022-09-18 06:20:00	2022-09-18 06:40:00	1	1.05215e+00	1.05207e+00
Message_03	system.disk.io, system.dfs.used, system.fs...	huy	2.00045	2022-09-18 06:20:00	2022-09-18 06:40:00	1	1.05215e+00	1.05207e+00
Message_04	system.net.bytes_sent, system.net.packets_out...	enb	3.00000	2022-09-17 20:20:00	2022-09-17 20:40:00	4	1.04940e+00	1.04938e+00
Message_05	system.net.bytes_sent, system.net.packets_out...	enb	3.00005	2022-09-17 20:20:00	2022-09-17 20:40:00	4	1.04940e+00	1.04938e+00



# 多模态异常检测-Log基础模型



- ✓ **日志类型全覆盖:** 自动识别日志结构, 对结构化、半结构化和非结构化日志自适应匹配时序异常检测、基于日志解析提取模板的半监督异常检测和基于预训练语言模型BigLog理解语义的半监督检测算法
- ✓ **高效率在线检测:** 可对流式日志数据实时检测并实时更新模型

1

## 拓扑图构建

TraceId	cmdb_id	ParentSpan	...
e4670ad5630a1096874ffff6efaa6c28	Weblogic_10	Weblogic_27	...
76b0f51edd1753898209a29920a75062	Weblogic_11	Weblogic_19	...
76b0f51edd1753898209a29920a75062	Weblogic_16	-	...



2

## 时间序列构建

各节点真正耗时

=

当前节点耗时

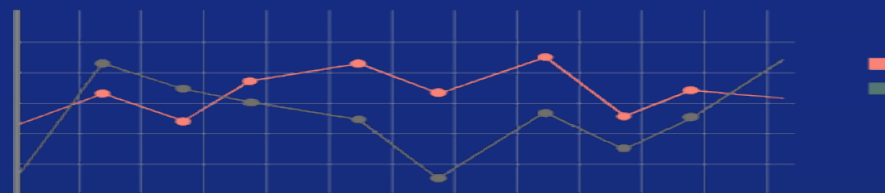
-

子节点总耗时

对每个二元组  
(cmdb\_id, kind)



Timestamp: RealDuration



3

## 异常检测算法

### 时间序列异常检测器

节点耗时期序



异常记录

基于知识消除误报

异常事件

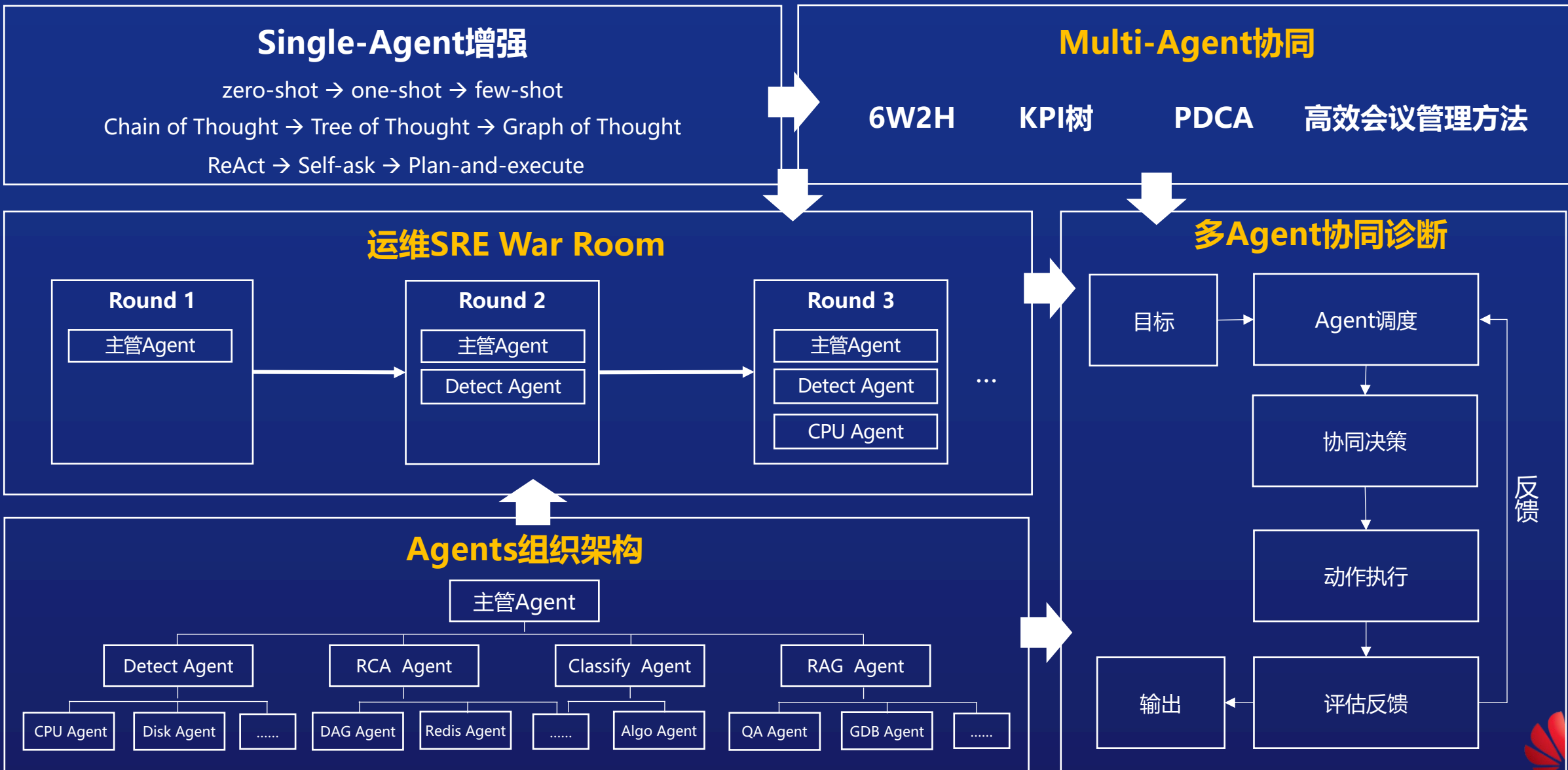


过滤后的异常记录

时序聚类

- ✓ **高效**: 将Trace数据转化为时间序列, 形成一套时序异常检测算法+不同参数的轻量化方案, 秒级异常检测
- ✓ **实时**: 通过实时拓扑关系分析帮助推测可能根因节点及故障传播方向

# 多Agent协同框架：基于企业组织管理方法的高效协同诊断



# 多Agent协同故障诊断

运维知识库 (企业内源)

LLM主管 Agent 决策

检测Agent

知识图谱 Trace拓扑结构

根因定位Agent

系统异常发现!

异常通知:  
2023/9/20 6:05检测到节点Weblogic\_16 eth0=\、eth0= vda 出现异常  
Weblogic\_17 eth0=\、eth0= vda出现异常

主管Agent:节点Weblogic\_16Weblogic\_17是web节点,需要关注!  
请检测Agent执行详细检测,请根因定位Agent根据详细检测结果查找根因  
请故障分析Agent给出故障分析报告返回给我

检测Agent执行磁盘模块检测...有异常,异常程度高,已获得高优先级检测结论,无需进一步检测CPU内存指标,输出详细检测结果  
节点Weblogic\_16 eth0=\ ['system.cpu.i\_dle']指标异常、异常程度0.849、异常开始时间2023/9/20 6:05  
节点Weblogic\_16 eth0=vda ['system.io.avg\_q\_sz', 'system.io.avg\_rq\_sz', 'system.io.await', 'system.io.r\_await', 'system.io.r\_s', 'system.io.util', 'system.io.w\_await', 'system.io.w\_s', 'system.io.wkb\_s']指标异常、异常程度1.049、异常开始时间2023/9/20 6:05  
节点Weblogic\_17 eth0=\ ['system.cpu.i\_dle', 'system.cpu.pct\_usage', 'system.load.1', 'system.load.15', 'system.load.5']

根因定位Agent执行根因定位分析中...  
根因为Weblogic\_16 eth0=vda

故障分类Agent:已知根因节点异常信息为节点Weblogic\_16 eth0=\ ['system.cpu.i\_dle']指标异常、异常程度0.849、异常开始时间2023/9/20 6:05  
节点Weblogic\_16 eth0=vda ['system.io.avg\_q\_sz', 'system.io.avg\_rq\_sz', 'system.io.await', 'system.io.r\_await', 'system.io.r\_s', 'system.io.util', 'system.io.w\_await', 'system.io.w\_s', 'system.io.wkb\_s']  
分析故障类别  
输出:故障类别为 磁盘故障

分析

故障原因分析: 根据以上分析,故障原因如下:  
1. 磁盘IO性能下降: 磁盘IO性能下降会导致数据库性能下降,进而影响业务系统的正常运行。  
2. 磁盘空间不足: 磁盘空间不足会导致数据库无法写入数据,进而影响业务的正常运行。  
3. 磁盘坏道: 磁盘坏道会导致数据损坏,进而影响业务的正常运行。  
4. 磁盘老化: 磁盘老化会导致性能下降,进而影响业务的正常运行。  
5. 磁盘配置不当: 磁盘配置不当会导致性能下降,进而影响业务的正常运行。  
6. 磁盘碎片化: 磁盘碎片化会导致性能下降,进而影响业务的正常运行。  
7. 磁盘温度过高: 磁盘温度过高会导致性能下降,进而影响业务的正常运行。  
8. 磁盘供电不稳定: 磁盘供电不稳定会导致性能下降,进而影响业务的正常运行。  
9. 磁盘固件问题: 磁盘固件问题会导致性能下降,进而影响业务的正常运行。  
10. 磁盘驱动程序问题: 磁盘驱动程序问题会导致性能下降,进而影响业务的正常运行。

修复建议

- 检查磁盘空间使用情况,清理无用数据,释放磁盘空间。
- 检查磁盘坏道情况,使用坏道修复工具进行修复。
- 检查磁盘老化情况,及时更换新磁盘。
- 检查磁盘配置是否正确,调整磁盘参数。
- 检查磁盘碎片化情况,定期执行碎片整理。
- 检查磁盘温度是否正常,确保散热良好。
- 检查磁盘供电是否稳定,更换电源适配器。
- 检查磁盘固件是否为最新版本,及时更新。
- 检查磁盘驱动程序是否为最新版本,及时更新。

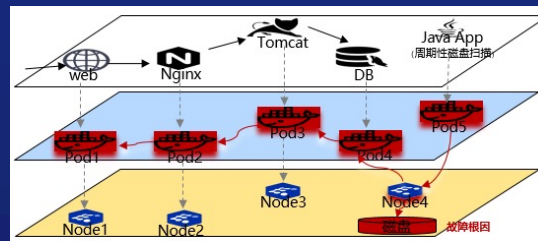
诊断建议

为了进一步确认故障原因,建议通过以下命令进行诊断:

```
1. iostat -x 1 1  
2. df -h  
3. dd if=/dev/zero of=/tmp/testfile bs=1M count=100  
4. mdadm --detail /dev/md0
```

总结

综上所述,故障原因可能是磁盘IO性能下降、磁盘空间不足、磁盘坏道、磁盘老化、磁盘配置不当、磁盘碎片化、磁盘温度过高、磁盘供电不稳定、磁盘固件问题、磁盘驱动程序问题。



提供故障爆炸半径,为恢复策略提供依据

调用运维工具

执行修复建议或进一步诊断

故障分析Agent

故障分类Agent

运维知识库 (企业内源)





## 创新

· 多Agent协同  
完成复杂运维任务

- 基于企业组织管理方法的多Agent协同框架，复杂运维任务处理更高效
- 多Agent协同完成运维主流程：异常检测->根因定位->故障分类->故障分析->修复建议

## 通用

· 解决运维领域  
共性问题

- 多模态异常检测基础模型，包含Trace、Metric、Log数据全方面处理能力，开箱即用
- 框架与算法不依赖具体特定应用场景，结合大模型实现较强的泛化能力

## 实用

· 故障恢复优先  
各模块松耦合

- 故障诊断报告体现可解释的故障爆炸半径，为实际生产运维故障快速恢复提供有力依据
- 各模块松耦合可插拔，可以全面应用于各类场景故障快恢需求，已在公司内部多场景落地

1. Zhang S, Pan Z, Liu H, et al. Efficient and Robust Trace Anomaly Detection for Large-Scale Microservice Systems. ISSRE, 2023.
2. Li D, Zhang S, Sun Y, et al. An Empirical Analysis of Anomaly Detection Methods for Multivariate Time Series. ISSRE, 2023.
3. Wang Z, Liu Z, Zhang Y, et al. RCAgent: Cloud Root Cause Analysis by Autonomous Agents with Tool-Augmented Large Language Models. arXiv, 2023.
4. Jin P, Zhang S, Ma M, et al. Assess and Summarize: Improve Outage Understanding with Large Language Models. ESEC/FSE, 2023.
5. Chen Y, Xie H, Ma M, et al. Empowering Practical Root Cause Analysis by Large Language Models for Cloud Incidents. arXiv, 2023.
6. Zhou X, Li G, Sun Z, et al. D-Bot: Database Diagnosis System using Large Language Models. arXiv, 2023.
7. Zhou X, Li G, Liu Z. Llm as dba. arXiv, 2023.
8. Wen Q, Gao J, Song X, et al. RobustSTL: A robust seasonal-trend decomposition algorithm for long time series. AAAI, 2019.
9. Liu Y, Tao S, Meng W, et al. LogPrompt: Prompt Engineering Towards Zero-Shot and Interpretable Log Analysis. arXiv, 2023.
10. Tao S, Liu Y, Meng W, et al. Biglog: Unsupervised large-scale pre-training for a unified log representation. IWQoS, 2023.
11. Ma L, Yang W, Xu B, et al. KnowLog: Knowledge Enhanced Pre-trained Language Model for Log Understanding. ICSE, 2023.
12. Zhong Z, Fan Q, Zhang J, et al. A Survey of Time Series Anomaly Detection Methods in the AIOps Domain. arXiv, 2023.
13. Wu H, Hu T, Liu Y, et al. Timesnet: Temporal 2d-variation modeling for general time series analysis. ICLR, 2023.
14. Yu G, Chen P, Li P, et al. Logreducer: Identify and reduce log hotspots in kernel on the fly. ICSE, 2023.



2023 CCF国际AIOps挑战赛决赛暨“大模型时代的AIOps”研讨会

# THANKS

## 针对运维领域海量知识快速获取、辅助诊断和故障分析能力

- 将LLM较为广泛的知识储备（横向能力）与运维领域专业知识（运维垂域）相结合
- 运维工具繁多，利用LLM+多智能体协同使能运维自动驾驶

## 针对多模态数据进行快速高效准确的异常检测能力

- 多损益、多类型数据：Log、Metric、业务黄金指标、Trace
- 数据异常不等于故障发生，识别可能造成运维故障的数据异常波动

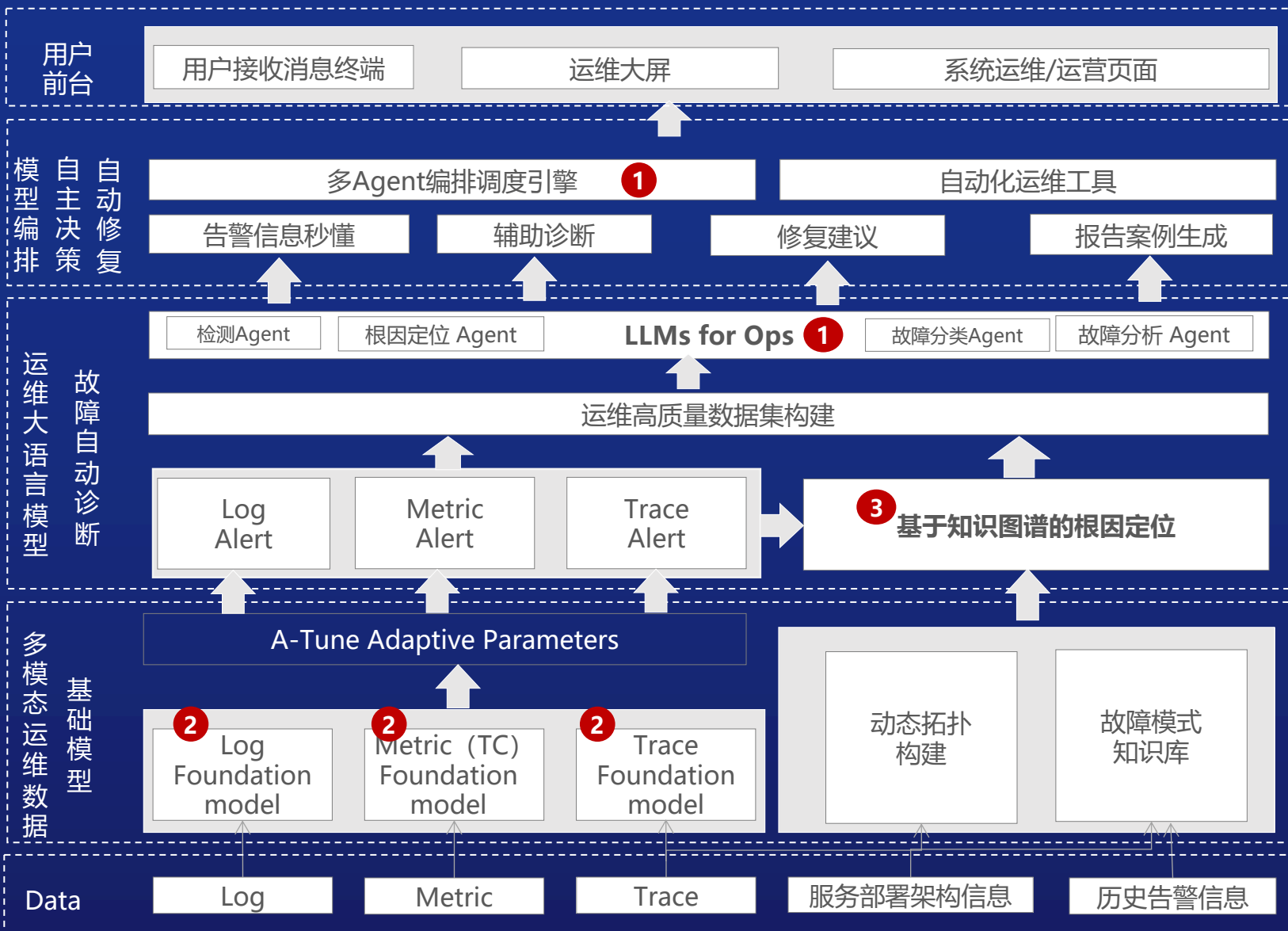
## 针对多源复杂部署的运维数据进行快速根因定位能力

- 告警繁多、区分故障传播节点和故障根因节点
- 历史排障经验积累数字化，形成知识库

	运维能力	是否主要	权重
1	多模态异常检测 (log、metric、trace)	是	0.5
2	多Agent协同故障诊断框架	是	0.5
3	根因定位能力	否	0.1
4	基于LLM的故障分类	否	0.1
5	基于LLM的故障分析报告与修复建议	是	0.1



# 基于大模型和多Agent协同自主决策、自动修复运维方案



## 关键技术

### ① 具备运维领域知识经验的LLM

- ①梳理异常检测告警信息和运维故障模式，形成运维领域高质量数据集，使用finetuning+外挂向量数据库的方案，使得LLM具备运维领域故障分析定位能力；
- ②使用主管LLM对异常问题与子领域Agent进行桥接，多个子领域Agent协同工作，提升运维效率

### ② 更全能的多模态数据异常检测基础模型

- ①针对成百上千条metric序列采取一定维度聚合的方法，化繁为简，一旦检测出异常，再细化分析关键指标的关键表现；
- ②针对半结构化和非结构化日志采取针对性的解析方式：前者注重模板提取，后使用uADR/sADR进行异常检测，后者注重语义理解，使用BigLog预训练模型结合Deep SVDD+SAD进行异常检测
- ③针对结构化明显的trace数据，抽取分离调用关系并定义节点关键数据，通过转化时间序列有效识别异常节点，并通过拓扑关系分析帮助推测可能根因节点及故障传播方向

### ③ 基于知识图谱的根因定位

- ①利用DBSCAN从时间维度聚类，形成异常事件
- ②利用Trace数据实时生成拓扑结构图
- ③根据故障模式知识库分析故障传播链路





## 创新

· 多Agent协同  
完成复杂运维任务

- 基于LLM实现故障分类、故障分析报告及修复建议
- 多Agent协同完成运维主流程：异常检测->根因定位->故障分类->故障分析->修复建议

## 通用

· 解决运维领域  
共性问题

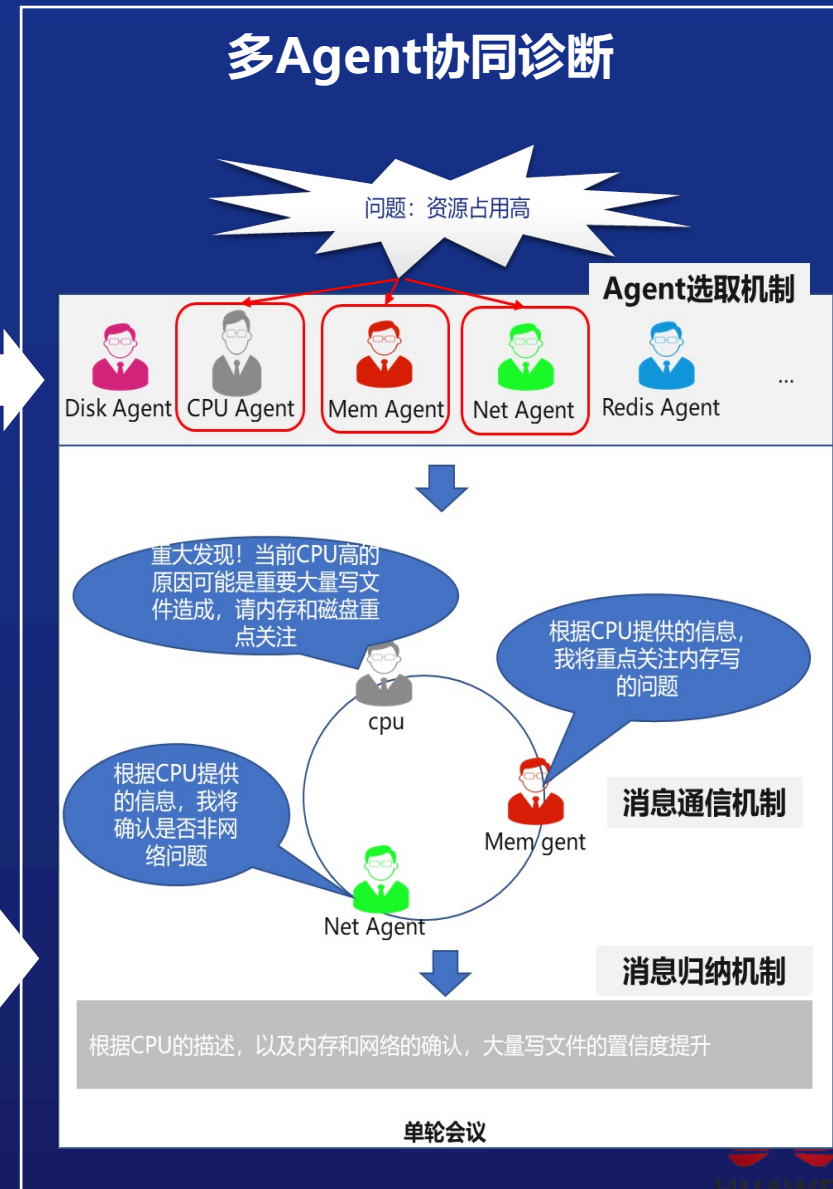
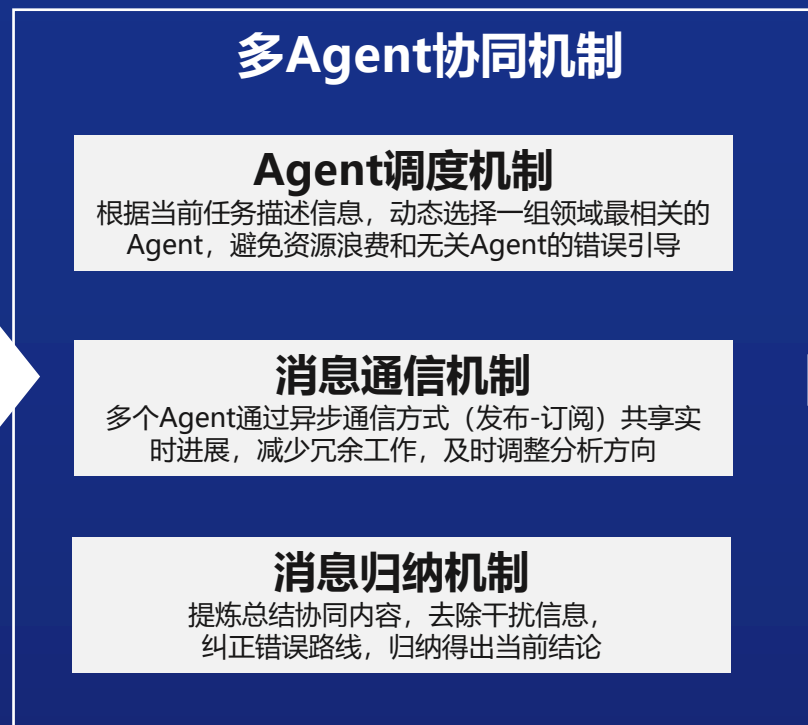
- 多模态异常检测基础模型，近开箱即用，少量参数即可完成参数调优
- 基于公域运维语料、知识库，完成运维领域的结构化知识问答

## 实用

· 各模块松耦合

- 松耦合的LLM底座，低成本替换、可插拔
- 各模块可插拔，既可以解决运维故障感知问题也可以全面应用于故障定位、排障、修复

# 多Agent协同框架：基于专业组织管理方法的高效协同诊断



# 多Agent协同故障诊断

运维知识库 (企业内源)

Trace模块

Web日志模块

知识图谱  
Trace拓扑结构

网络模块

Redis日志模块

xx模块

CPU模块

LLM主管  
Agent**决策**

检测Agent

根因定位Agent

系统异常  
发现!

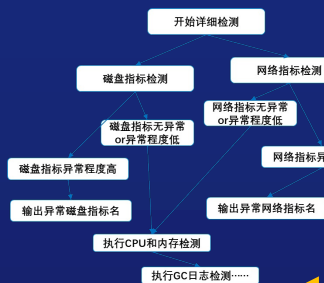
异常通知:  
2023/9/20 6:05检测到节点Weblogic\_16 eth0=\、eth0= vda 出现异常  
Weblogic\_17 eth0=\、eth0= vda出现异常

主管Agent:节点Weblogic\_16Weblogic\_17是web节点, 需要关注!  
请检测Agent执行详细检测, 请根因定位Agent根据详细检测结果查找根因  
请故障分析Agent给出故障分析报告返回给我

检测Agent执行磁盘模块检测...有异常, 异常程度高, 已获得高优先级检测结论, 无需进一步检测CPU内存指标, 输出详细检测结果  
节点Weblogic\_16 eth0=\ ['system.cpu.i\_dle']指标异常、异常程度0.849、异常开始时间2023/9/20 6:05  
节点Weblogic\_16 eth0=vda ['system.io.avg\_q\_sz', 'system.io.avg\_rq\_sz', 'system.io.await', 'system.io.r\_await', 'system.io.r\_s', 'system.io.util', 'system.io.w\_await', 'system.io.w\_s', 'system.io.wkb\_s']指标异常、异常程度1.049、异常开始时间2023/9/20 6:05  
节点Weblogic\_17 eth0=\ ['system.cpu.i\_dle', 'system.cpu.pct\_usage', 'system.load.1', 'system.load.15', 'system.load.5']

根因定位Agent执行根因定位分析中...  
根因为Weblogic\_16 eth0=vda

故障分类Agent:已知根因节点异常信息为节点Weblogic\_16 eth0=\ ['system.cpu.i\_dle']指标异常、异常程度0.849、异常开始时间2023/9/20 6:05  
节点Weblogic\_16 eth0=vda ['system.io.avg\_q\_sz', 'system.io.avg\_rq\_sz', 'system.io.await', 'system.io.r\_await', 'system.io.r\_s', 'system.io.util', 'system.io.w\_await', 'system.io.w\_s', 'system.io.wkb\_s']  
分析故障类别  
输出: 故障类别为 磁盘故障



```
分析
根据提供的异常指标信息, 我们可以从以下几个方面进行分析:
1. 系统CPU使用率异常: 这通常意味着系统负载过高, 可能是由于应用程序的某个进程占用了大量的CPU资源。
2. 磁盘I/O等待时间异常: 这通常与磁盘性能有关, 可能是由于磁盘空间不足、磁盘碎片化或磁盘老化等原因。
3. 网络I/O等待时间异常: 这通常与网络性能有关, 可能是由于网络拥塞、网络配置错误或网络设备故障等原因。
4. 系统负载异常: 这通常与系统资源分配有关, 可能是由于系统配置不当、应用程序bug或系统资源不足等原因。
5. 系统内存使用率异常: 这通常与内存管理有关, 可能是由于内存泄漏、内存分配不当或内存碎片化等原因。
6. 系统垃圾回收(GC)异常: 这通常与垃圾回收算法有关, 可能是由于垃圾回收配置不当、垃圾回收器bug或垃圾回收器资源不足等原因。
7. 系统日志异常: 这通常与系统日志记录有关, 可能是由于日志配置不当、日志记录器bug或日志记录器资源不足等原因。
8. 系统配置异常: 这通常与系统配置有关, 可能是由于系统配置不当、系统配置错误或系统配置冲突等原因。
9. 系统安全异常: 这通常与系统安全有关, 可能是由于系统安全配置不当、系统安全漏洞或系统安全攻击等原因。
10. 系统其他异常: 这通常与系统其他方面有关, 可能是由于系统其他方面的配置不当、系统其他方面的bug或系统其他方面的资源不足等原因。
综合以上几个方面的分析, 我们可以初步判断: 可能存在某个进程占用了大量的CPU资源, 导致了CPU使用率异常, 同时该进程可能还在进行大量的磁盘I/O操作, 导致了磁盘I/O等待时间异常。
修建议议
1. 检查系统的负载情况, 确认是否存在异常的进程。
2. 检查系统的磁盘空间使用情况, 确认是否存在磁盘空间不足的情况。
3. 检查系统的网络配置情况, 确认是否存在网络配置错误的情况。
4. 检查系统的内存使用情况, 确认是否存在内存使用率异常的情况。
5. 检查系统的垃圾回收配置情况, 确认是否存在垃圾回收配置不当的情况。
6. 检查系统的日志配置情况, 确认是否存在日志配置不当的情况。
7. 检查系统的配置情况, 确认是否存在系统配置不当的情况。
8. 检查系统的安全情况, 确认是否存在系统安全漏洞或系统安全攻击的情况。
9. 检查系统的其他方面情况, 确认是否存在系统其他方面的配置不当、系统其他方面的bug或系统其他方面的资源不足等情况。
诊断建议
为了进一步确认问题的原因, 我们可以尝试以下操作:
1. 使用top命令查看系统的负载情况。
2. 使用df命令查看系统的磁盘空间使用情况。
3. 使用netstat命令查看系统的网络配置情况。
4. 使用free命令查看系统的内存使用情况。
5. 使用jstat命令查看系统的垃圾回收配置情况。
6. 使用tail命令查看系统的日志配置情况。
7. 检查系统的配置文件。
8. 检查系统的安全日志。
9. 检查系统的其他方面情况。
总结
根据以上分析, 我们初步判断: 可能存在某个进程占用了大量的CPU资源, 导致了CPU使用率异常, 同时该进程可能还在进行大量的磁盘I/O操作, 导致了磁盘I/O等待时间异常。
```

调用运维工具

执行修复建议或进一步诊断

故障分析Agent

故障分类Agent

运维知识库 (企业内源)

