

# 基于多源数据融合与大模型推理的 微服务根因定位

hwlyyzc战队 国泰海通证券股份有限公司

主办单位：中国计算机学会（CCF）

承办单位：中国计算机学会互联网专委会、中国科学院计算机网络信息中心、中国移动研究院、清华大学

协办单位：华为2012实验室、阿里云、中兴通讯、中国移动九天团队、南开大学、西安电子科技大学、清华大学计算机科学与技术系、神州灵云

# 目录 CONTENTS

第一章节 团队介绍

第二章节 赛题与现状分析

第三章节 方案设计

第四章节 效果检验

第五章节 总结与展望

# 第一节 团队介绍





**国泰海通证券**  
GUOTAI HAITONG SECURITIES

## 金融科技条线

系统运行部

SRE



**彭涛**

运维分析



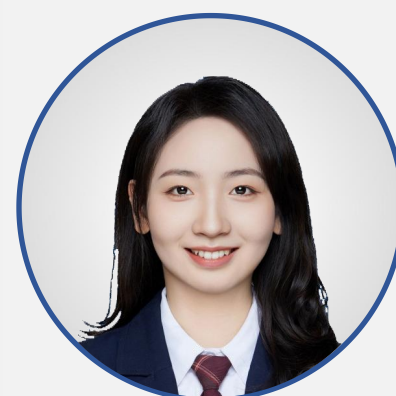
**郭建华**

运维分析



**房千淩**

运维质量



**黎治娴**

运维质量

研究运用先进运维技术栈，持续推进运维效能提升工具建设，不断提升运维服务质量，保障互联网应用及相关业务系统的高效稳定运行

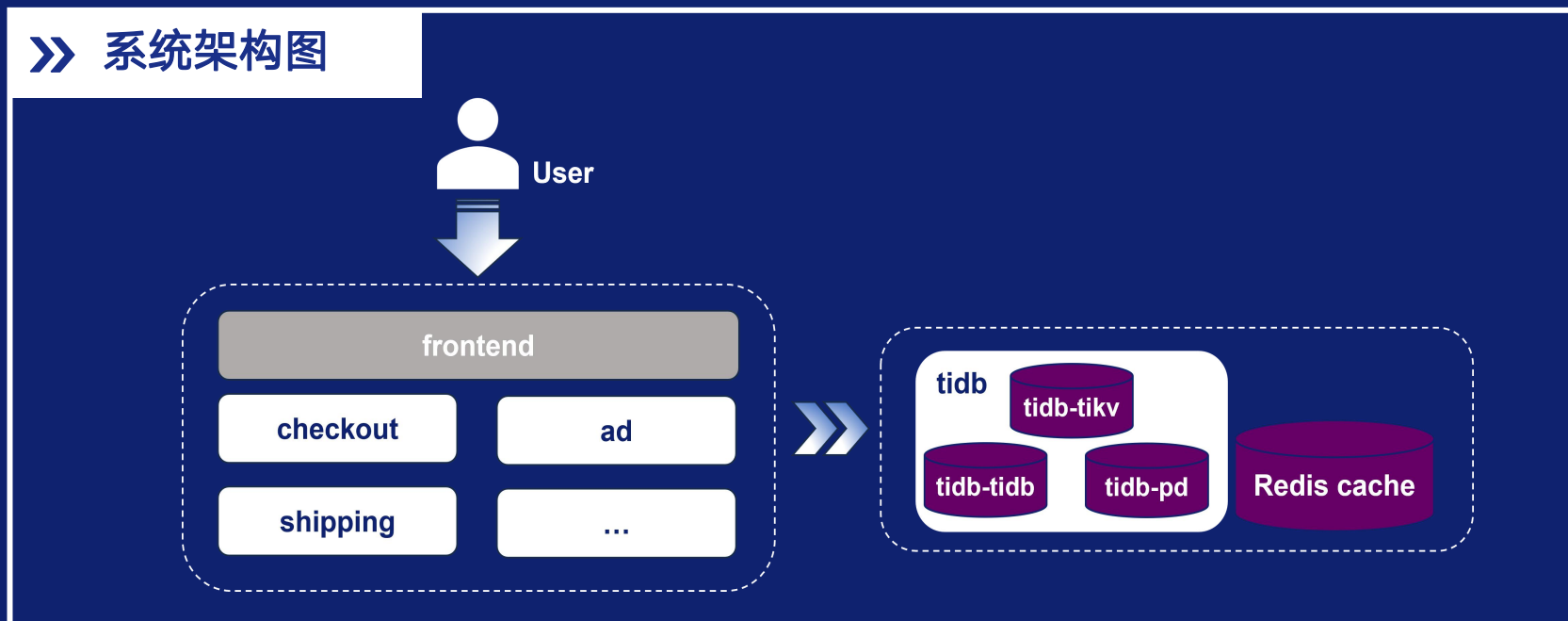
## 第二章节

# 赛题与现状分析

# 赛题及数据说明

- ◆ 赛题环境：基于 Google 开源的微服务示例系统 HipsterShop，每个微服务对应一个容器级别的监控系统。
- ◆ 系统架构：采用动态部署架构，共包含 10 个核心微服务和 8 台虚拟机，每个微服务各部署 3 个 Pod，这些 Pod 会被动态调度分布到 8 台虚拟机上；此外，系统中的 TiDB 组件也部署在虚拟机上，每个服务部署 1 个 Pod。
- ◆ 故障注入方式：Service 级、Pod 级和 Node 级三种层级的故障注入，Service 级故障表示该服务下的所有 3 个 Pod 同时注入故障；Pod 级故障仅对某一个具体的 Pod 注入故障。

## » 系统架构图



## » 多源数据



## 痛点

## 对策

## 方法设计

指标漏检严重  
噪声较多

多算法融合  
提升准确率



集成 Isolation Forest、HBOS、IQR  
结合全局+局部检测

调用链数据量大  
冗余信息多

异常筛选  
语义去重



基于耗时分布、Pod 分布筛选  
结合 Z-score 与 Levenshtein 压缩

日志关键字匹配粗糙  
缺乏模式识别

错误定位 + 模  
板解析结合



错误模式识别 + Drain3 模板解析  
相似度匹配算法压缩数据量

大模型幻觉风险高  
可解释性不足

高质量 Prompt  
约束机制



融合拓扑与多源证据  
采用“下游优先+权重分配”打分体系



# 第三章节

# 方案设计



# 总体设计



The system experienced an anomaly from 2025-04-30T19:10:15Z to 2025-04-30T19:40:15Z. Please infer the possible cause.

理解问题  
提取异常时间段



LLM

- start\_time: 2025-04-30T19:10:15Z
- end\_time: 2025-04-30T19:40:15Z

构造 prompt

调用拓扑关系

组件有效命名

加权打分策略

Json格式化输出

大模型  
根因推理



根因定位输出结果

- ✓ 根因组件component
- ✓ 异常原因reason
- ✓ 推理步骤reasoning\_trace

单指标检测

IF+HBOS+IQR  
全局检测

候选异常点

时序突变  
局部检测

调用链检测

异常调用链  
定位

计算Z-score值

相似度  
计算去重

日志检测

日志格式化

模式识别



- ◆ 创新点1：引入基于耗时分布、错误标签与Pod分布比例的多维筛选，并结合Z-score突增检测与Levenshtein距离语义去重，实现了高信息密度的调用链压缩。

## » 异常调用链定位

### ERROR调用链异常

- ✓ (key: value)→(error: True)
- ✓ (key: value)→(http.status\_code: >400)

Span耗时异常： > 95分位数

Pod分布比例异常：依据显著性统计分析，位置位于分布三倍标准差外的记录

异常调用链  
记录集合



■ **难点1：**某些固定存在错误的调用链可能干扰后续大模型根因推理判断。

✓ **应对：**对于错误跨度较大的trace，采用基于Z-score的递归方法识别突发式错误行为。

■ **难点2：**调用链数据量大、存在重复数据，易超出大模型token限制。

✓ **应对：**引入Levenshtein距离算法，对超出大模型token上限的tags.message字段进行语义去重，相似度大于95%的message被合并，避免语义冗余导致的大模型token浪费。



# 日志异常检测模块

◆ 创新点2：将错误模式与Drain3模板解析相结合，既能捕获显性错误，又可发现模式级异常。



监控日志记录



错误模式匹配



Drain3日志解析



异常日志压缩

识别历史故障日志中的典型错误模式，形成错误汇总集合；在异常时段匹配日志内容，以快速锁定显著错误日志条目



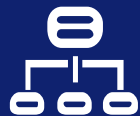
日志清洗

移除时间戳、IP地址等变量信息，并将数字、哈希值等替换为<\*>占位符



异常检测

在指定时间窗口内统计各模板出现频率，检测频率**突增、突降或新出现**的模板，以识别异常模式



解析树构建

按照日志消息的词数进行第一层分支，再按特定位置单词/标记进行中间层分支，最终在叶子节点**聚合日志模板**

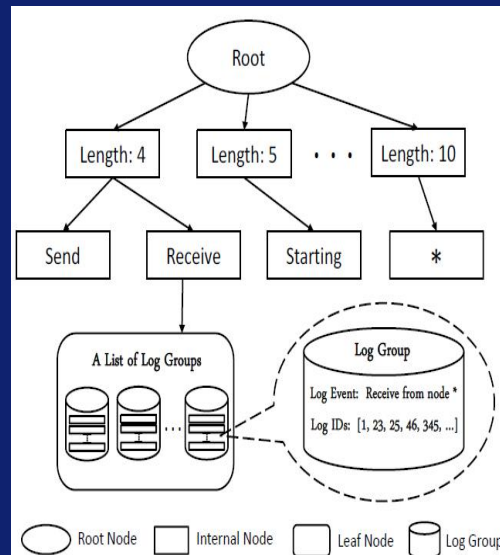


Fig. 2: Structure of Parse Tree in Drain (depth = 3)



- ◆ 创新点3：集成Isolation Forest、HBOS与IQR三类互补型算法，通过全局统计与局部时序双重检测，提高了对高维特征、分布稳定指标及突变模式的综合识别能力。

## » 全局检测模块

全局检测模块采用Isolation Forest、HBOS与IQR算法，通过加权方式综合生成异常分数，有效识别全局范围内的异常指标。

- 异常分数计算公式：

$$anomaly\_score_i = \frac{if\_scores_i - (0.1 \times hbos\_scores_i)}{2}$$

- 异常判定逻辑：

$$is\_anomaly_i = \begin{cases} false, & if (if\_scores_i > 0) \vee (anomaly\_score_i > 0) \\ (anomaly\_score_i < -0.5) \vee iqr\_anomalies_i & otherwise \end{cases}$$

## » 局部检测模块

局部检测模块专注于**时间序列突变模式**的识别，通过对象和指标的聚合汇总，形成最终的异常检测结果，为根因定位提供关键信息。该模块聚焦于指标在异常时间段内是否存在明显的**结构性异常趋势**，如：

- ✓ 突增后迅速回落
- ✓ 突降并维持低位
- ✓ 突增并持续维持高位
- ✓ 剧烈震荡（突增突降）
- ✓ 突降后迅速回升

当检测时间段中异常点比例大于阈值且  
存在局部时序异常时，判定为异常



- ◆ **创新点4：**在根因推理模块中，方案将系统架构、有效组件清单与多源异常数据融合为高质量Prompt，配合明确的输出约束与量化评分规则，实现了可解释、可复现的根因定位。



## » 设计优势

- **多源融合、规则引导：**以统一Prompt注入系统结构、服务调用拓扑关系、故障注入边界与评分准则，促使LLM“有据可依”。
- **强格式输出：**严格的JSON约束，确保结果可解析、可比对、可评测。
- **组件筛选防幻觉：**结合合法组件清单与“最下游优先”等打分机制，降低无关/虚构组件的选择概率。
- **时序因果：**以时间先后组织reasoning\_trace，强调“先症状后传播”，减少将次生故障当作根因的风险。



- ◆ **创新点4**：在根因推理模块中，方案将系统架构、有效组件清单与多源异常数据融合为高质量Prompt，配合明确的输出约束与量化评分规则，实现了可解释、可复现的根因定位。

## ● 系统描述信息

服务部署信息

调用拓扑关系

组件命名清单

## ● 多源数据检测

metrics检测异常

logs检测异常

traces检测异常

计算异常组件得分，  
推理根因组件



## ● 评分准则

- ✓ +1 — 多数据源多次出现；
- ✓ +2 — 异常span中的被调用方且状态码 $\geq 400$ 或timeout关键字；
- ✓ +1 — 异常关键词关联；
- ✓ +4 — 位于服务调用依赖关系中最下游；
- ✓ +10 — 重启restart关键字。

## ● 平分决策

**下游优先**：依赖图中更下游的组件优先

**时间优先**：异常时间更早的组件优先

异常类型优先级：重启 > 5xx / timeout > 异常关键词 > 频次。



# 第四章 效果检验





• 根因定位准确率

准确率：**57.75%**，部分误差源于故障层级判断有误或pod有误，另有部分误差来自调用链推理方向错误，导致跨模块定位。

• 原因描述一致性

约**57%**的cases在reason字段中能够体现与参考答案中Metrics、Logs、Traces关键词一致的异常特征，且与检测到的异常数据相符。

• 总体分析

优势：Traces检测召回率较高，能较好覆盖调用链相关异常信息。Logs检测精度中等且较为稳定。

不足：Metrics检测薄弱，存在明显漏检与噪声；推理链覆盖不完整，部分原因描述泛化或缺乏依据。

• 召回率与精确率

Traces 维度**表现最佳**，召回率较高，表明系统对调用链相关关键词（如suspected\_service / trace\_info等）捕捉较为充分。

Logs 维度**相对均衡**，说明系统对日志关键信息提取具有一定覆盖度和可靠性，但在异常日志模式的精细化匹配上仍有提升空间。

Metrics 维度**问题最突出**，反映出指标检测环节存在漏检与噪声并存的现象，可能是异常检测算法阈值问题或过拟合导致。

维度	召回率	精确率
Trace	62.86%	31.87%
Log	50.67%	36.50%
Metric	33.28%	10.82%



## 第五章节

# 总结与展望



## 模型创新点

构建了由**数据输入层、异常检测层、大模型推理层与结果输出层**组成的根因定位体系。其中异常检测层运用Isolation Forest、HBOS、IQR、Z-score、Levenshtein、Drain3等多种算法协助异常定位；大模型推理层将系统架构、有效组件清单、多源异常数据、输出约束与量化评分规则融合为高质量Prompt，实现了可解释、可复现的根因定位。

## 模型缺点

方案仍存在指标检测精度不足、推理链覆盖不够完整以及部分原因描述泛化或出现虚构的不足，产生**幻觉性诊断**，生成看似合理但实际错误的故障原因描述，这在一定程度上影响了结果的可靠性。

## 总结

## 展望

## 运维通用性

方案具有较强的通用性，单个错误实例的**平均检测时长较短（约15s）**，可应用于生产环境中的告警分析与故障溯源，缩短人工排障时间，提升运维效率。  
根因定位**总体准确率为57.75%**，多数案例能够正确识别故障组件与原因特征，尤其在调用链主导场景下效果较佳。

## 后续改进

方案在多维数据融合和可解释性推理方面表现出良好效果与落地潜力，后续通过优化指标异常模式识别、引入规则约束与白名单过滤以及加强推理链细粒度验证，**进一步提升检测精度与根因定位的稳定性**，从而为复杂微服务环境下的智能运维提供更可靠的技术支撑。



OpenAIOps AIOPS | 2025 CCF国际AIOps挑战赛  
2025 CCF International AIOps Challenge

# THANKS

主办单位：中国计算机学会（CCF）

承办单位：中国计算机学会互联网专委会、中国科学院计算机网络信息中心、中国移动研究院、清华大学

协办单位：华为2012实验室、阿里云、中兴通讯、中国移动九天团队、南开大学、西安电子科技大学、清华大学计算机科学与技术系、神州灵云