

分享主题

# 大模型时代的智能运维 (AIOps)

裴丹 清华大学



# 人工智能产业链联盟

星主： AI产业链盟主

 知识星球

微信扫描预览星球详情





面临哪些技术挑战?



与以往的AIOps小模型是什么关系?



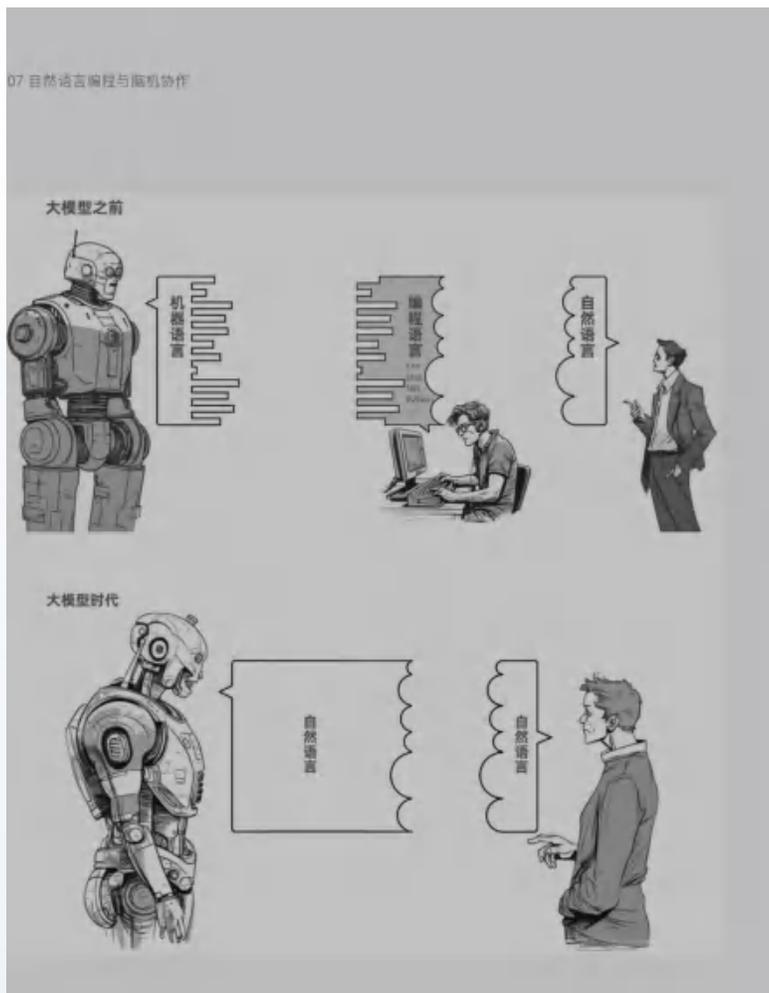
如何选择通识大模型底座?



近、中、长期有哪些应用?

问题

# 在大模型时代，AIOps可以“说人话”了





## GatesNotes

LOG IN

SIGN UP



### THE FUTURE OF AGENTS

## AI is about to completely change how you use computers

And upend the software industry.

By Bill Gates | November 09, 2023 · 12 minute read



## 运维智能体

助理

教练

顾问

参谋

专家

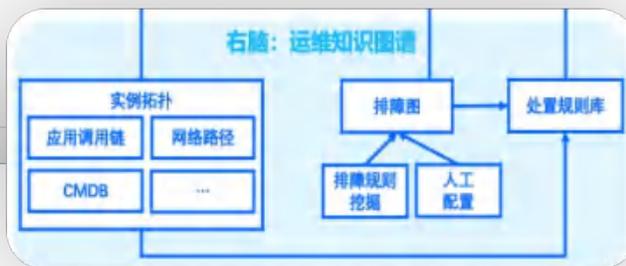
# AIOps中的智能体

## 大语言模型的模型栈

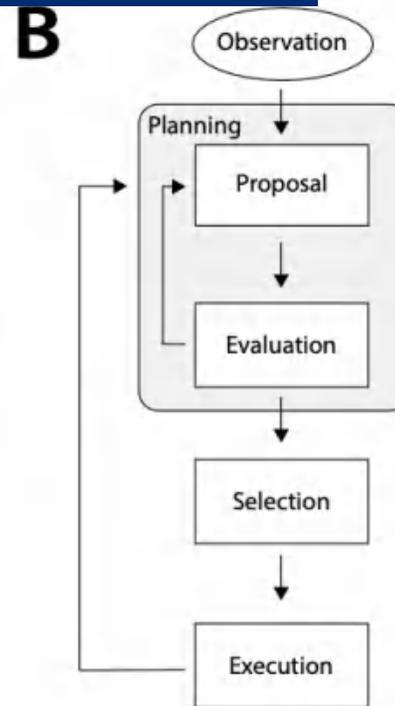
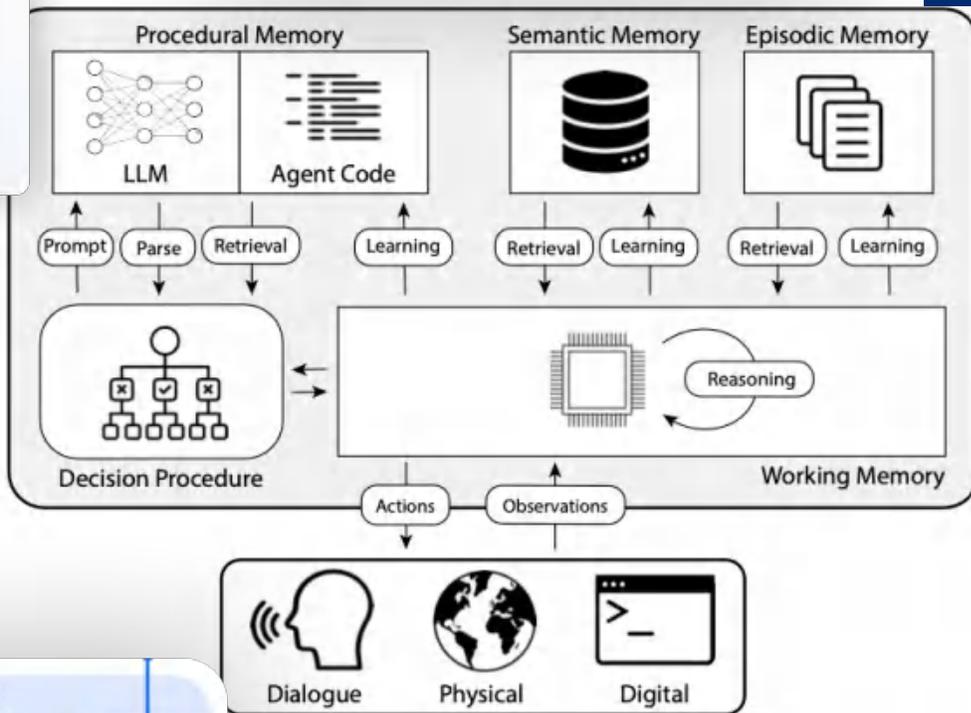
- L3 **私有部署运维大语言模型**  
基于私有运维数据、提示工程、外挂知识库检索
- L2 **运维大语言模型**  
基于公网运维资料、知识库、进行预训练、微调、提示工程
- L1 **松耦合的通识大语言模型底座**



## 右脑：运维知识图谱



历史工单、告警、  
操作记录、文档等



## 左脑：AIOps算法



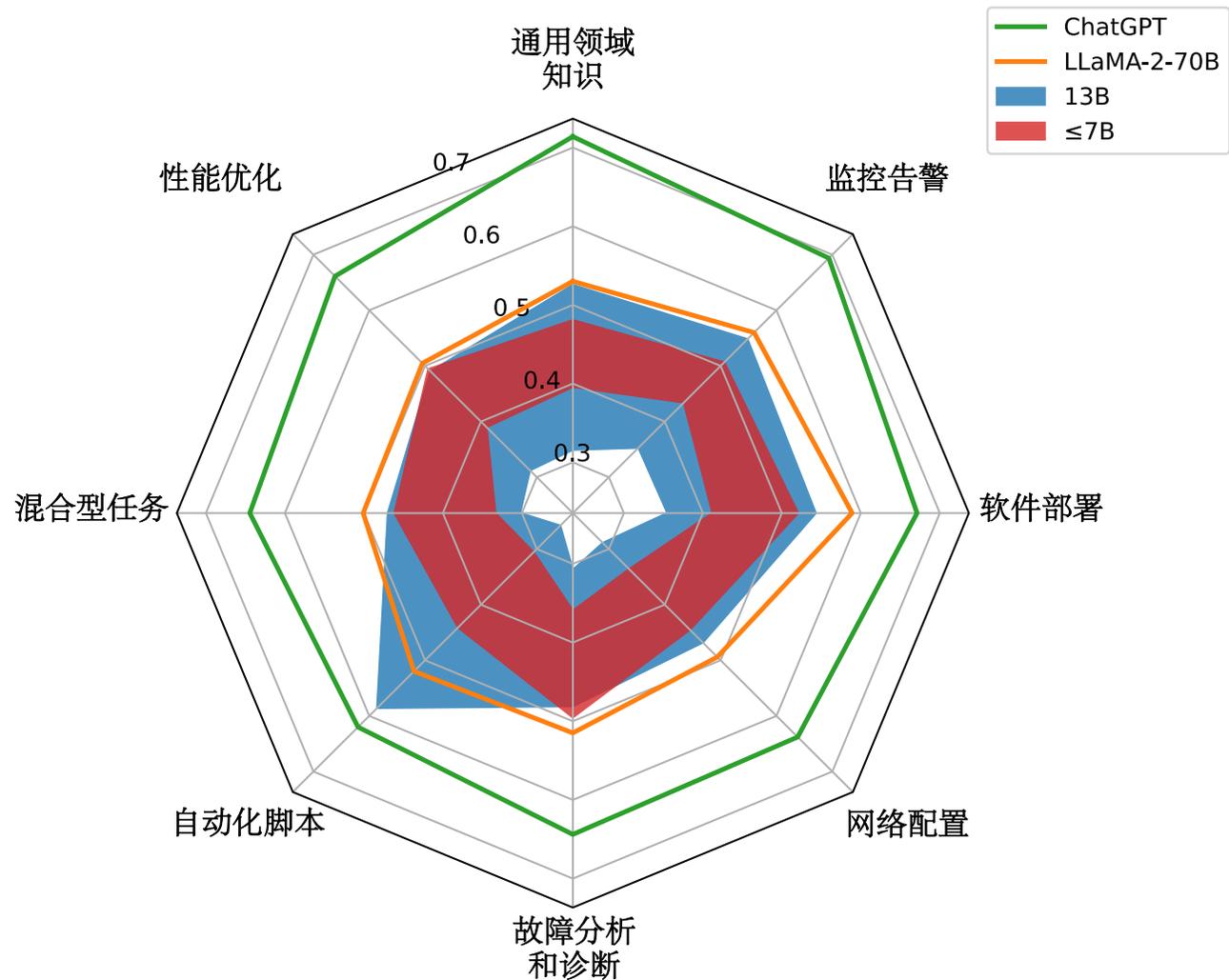
architectures for language agents (CoALA). **A**: CoALA defines a set of interacting components. The **decision procedure** executes the agent's source code. This source code consists of interactions with the LLM (prompt templates and parsers), internal memories (retrieval and learning), and the external environment (grounding). **B**: Temporally, the agent's decision procedure executes a **decision cycle** in a loop with the external environment. During each cycle, the agent uses **retrieval** and **reasoning** to plan by proposing and evaluating candidate **learning** or **grounding** actions. The best action is then selected and executed. An observation may be made, and the cycle begins again.

# 通识大模型在运维领域表现整体不如人意，而且参差不齐

欢迎为评测榜单贡献题目和模型：

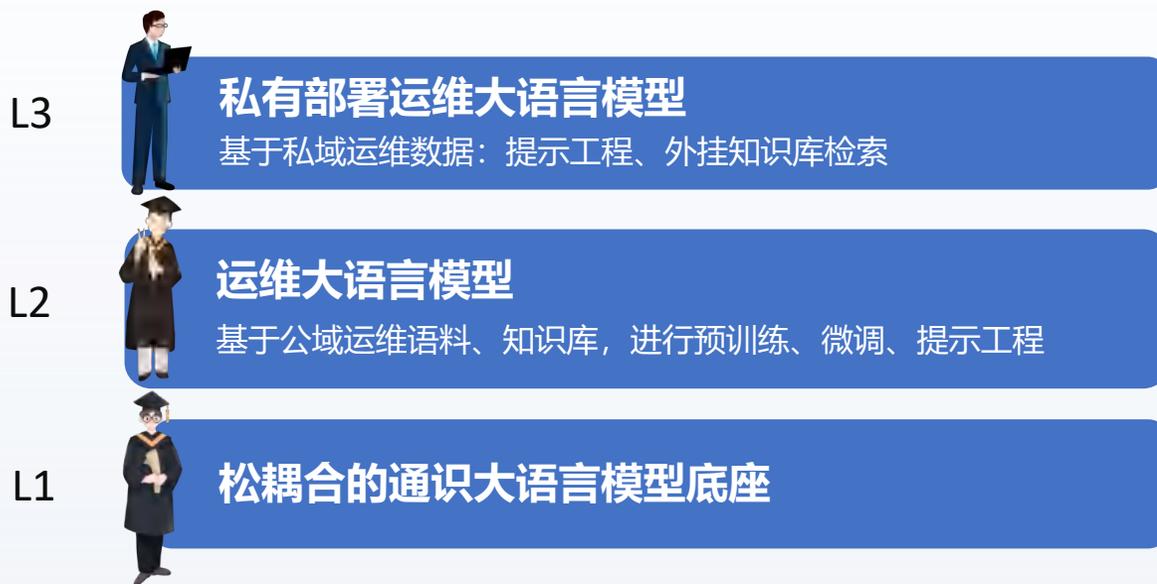
<https://opseval.cstcloud.cn/content/leaderboard>

Model	Naive	Zero-shot		C	C	C	C	C
		SC	CoT					
GPT-4	/	/	/					
GPT-3.5-turbo	66.60	66.80	69.60					
LLaMA-2-70B	55.80	57.20	61.90					
LLaMA-2-13B	41.80	46.50	53.10					
LLaMA-2-7B	39.50	40.00	45.40					
Qwen-7B-Chat	45.90	46.00	47.30					
Baichuan2-13B-Chat	37.90	38.30	42.70					
InternLM-7B	38.70	38.70	43.90	43.90	43.20	43.20	51.40	51.40
Chinese-Alpaca-2-13B	37.70	37.70	49.70	49.70	48.60	48.60	50.50	50.50
ChatGLM2-6B	24.80	24.70	36.60	36.50	37.60	37.60	40.50	40.50
Chinese-LLaMA-2-13B	29.40	29.40	37.80	37.80	40.40	40.40	28.80	28.80
Baichuan-13B-Chat	18.30	20.40	28.60	37.00	24.10	26.70	18.20	17.80



第三梯队 7

## 大语言模型的模型栈





# 也要避免悲观：前述所有技术挑战都有技术思路可以解决

挑战	解决思路
避免幻觉	
可解释性强	检索增强 增大显式知识占比：思维链、思维树、思维图、知识图谱 “有据可依”的生成策略提供可解释性
严肃语料不足	训练策略：课程学习，由易到难
私有部署开销低	
私域数据数量、质量不足	模型分层 在公域做预训练、微调、提示工程 私有部署时避免预训练、微调 通过检索方式融合本地知识库 文档、提示作为便捷的知识工程手段 通过降低模型精度降低私有部署推理开销
底座模型选择	与开源大语言模型底座尽量解耦
大语言模型无法直接处理结构化、多模态、实时数据	多模态基础模型群； AIOps结构化大模型：左脑AIOps算法，右脑运维知识图谱
如何结合大量存量AIOps、自动化运维工具	基于智能体方式，采用基础模型编程框架，融合小模型运维工具（插件（Plug-in），工具（Tool））

## 运维大语言模型带来的直接增量价值

-  眼：监控
-  手：自动化运维
-  脑：AIOps
  - 左脑：算法
  - 右脑：运维知识图谱
-  心：专家人工审核、决策
-  耳：基于运维大语言模型的输入增强
-  嘴：基于运维大语言模型的输出增强



## 近中期应用的定位

助理

教练

顾问

参谋

## 中长期应用的定位

+内部专家

## 对结构化知识检索



问: X应用交易响应时间比平时增加1倍, 可能原因是什么?

答: 可能的原因包括:

1. 依赖Y应用处理时间过长
2. 中间件Weblogic问题, 请检查 ...
3. 网络丢包问题, 请检查相关服务间见的网络链路丢包率是否有异常
4. 第三方应用问题, 请检查...

问: Oracle活跃会话数的确飙高了, 请问可能的原因是什么?

答: Oracle实例B的活跃会话数的确飙高可能的原因包括:

1. log-file-syn告警
2. 实例B的主机异常
3. 有CPU占用过高的SQL

问: 中间件WebLogic的确异常, 请问下面怎么办?

答: 中间件Weblogic异常的原因包括:

1. WebLogic堆空间满, 请检查 ...
2. Oracle数据库“平均活跃会话数”指标飙高

问: 查到实例B上有个SQL C占用CPU 90%, 该如何处理?

答: 建议杀掉该SQL C, 具体指令是 ... 。建议联系Oracle数据库管理员执行该操作。



## 私有部署运维大语言模型

基于私域运维数据：提示工程、外挂知识库检索



## 运维大语言模型

基于公域运维语料、知识库,进行预训练、微调、提示工程



## 松耦合的通识大语言模型底座

## 检索插件

## 私有文档

生成技术：检索增强

生成策略：通过“有据可依”提供可解释性  
让模型能够正确的给自己生成的内容打上引用（reference），清晰地展示出每一段是来自于哪个文本



## 对存量脚本进行文字解读



对运维数据中的本体、实体、属性、字段、标签进行注释



# 近中期应用: Lang2Query, 为单个存量工具提供自然语言交互增强, 提供意图识别、总结等能力

基础: 数据标准化、工具接口标准化

Lang2API  
(自动生成API  
调用)

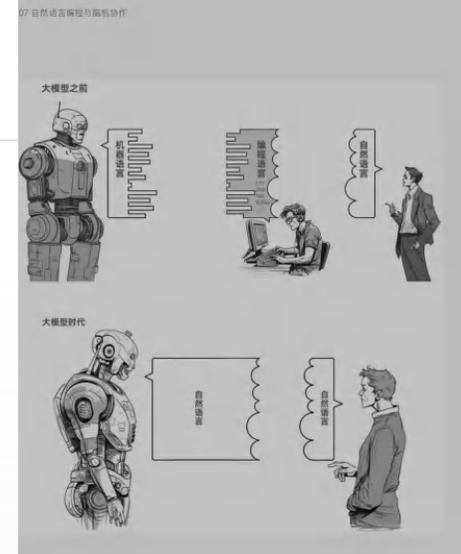
Lang2SQL  
(自动生成SQL)

Lang2GSQL  
(自动生成图  
SQL)

Lang2SPL  
(自动生成日志  
查询语句)

Lang2Scripts  
(自动生成脚本)

Lang2Config  
(自动生成配置)



## 举例：基于大语言模型的实时故障工单自动生成

基于大语言模型，以实时日志、调用链、指标等数据为输入，结合故障检测、定位、根因分析、影响分析等AIOps工具的输出，自动生成实时故障工单。



## 避免过于乐观：运维大模型仍面临不少挑战

杜绝幻觉、可解释性强、私有部署开销低、私有语料质量数据均不足、融合存量知识、工具、多模态数据、通识大语言模型底座不易选择

## 避免过于悲观：挑战都可解

模型分层：通识大语言模型、运维大语言模型、私有部署运维大语言模型

区分、整合非结构化大模型与结构化大模型

### 关键组件

运维大语言模型

多模态基础模型群

结构化大模型：左脑AIOps算法、右脑运维图谱

通过检索融合本地知识库

智能体&基础模型编程框架

### 运维大语言模型是核心基础

- 检索增强、有据可依
- 课程学习、由易到难
- 知识工程：从文档到知识图谱，增加显示知识
- 检索本地知识库、降低模型精度
- 与通识大语言模型底座尽量解耦

## 应用及路径建议：小步快跑、以用促建

应用的定位：从助手、教练、顾问、参谋到内部专家



近期



近中期



中长期

数字化运维助手，私有运维文档问答，运维脚本解读，运维数据注释

为单个运维工具提供自然语言交互增强

基于智能体，编排多个工具完成更复杂运维任务

## 谨慎乐观

大势所趋、前景可期、机遇与挑战并存、协同创新、以用促建

# AI人工智能产业链联盟

#每日为你摘取最重要的商业新闻#

更新 · 更快 · 更精彩



Zero

AI人工智能产业链联盟创始人  
河北清华发展研究院智能机器人中心运营经理



base:北京



主要业务:AI商业化答疑及课程应用场景探索, 各类AI产品学习手册, 答疑及课程



欢迎扫码交流

提供: 学习手册/工具/资源链接/商业化案例/  
行业报告/行业最新资讯及动态



人工智能产业链联盟创始人

邀请你加入星球, 一起学习

## 人工智能产业链联盟报 告库



星主: 人工智能产业链联盟创始人

每天仅需0.5元, 即可拥有以下福利!  
每周更新各类机构的最新研究成果。立志将人工智能产业链联盟打造成市面上最全的AI研究资料库, 覆盖券商、产业公司、研究院所等...

知识星球

微信扫码加入星球 ▶



The logo for ifenix features the word "ifenix" in a white, lowercase, sans-serif font. The dot on the "i" is replaced by a white arrowhead pointing to the right. The background is a solid blue color with a large, stylized white arrow shape on the right side, composed of multiple parallel lines that create a sense of depth and movement.

专 注 数 字 化