

中文大模型基准测评报告

— 2024年度中文大模型阶段性进展评估

2024.2.27

目录

01 国内大模型关键进展

- 2024年大模型关键进展与中文大模型全景图
- 国内外大模型发展趋势

02 测评体系

- 中文大模型基准SuperCLUE介绍
- 测评体系、层次、方法及示例说明

03 大模型综合测评结果

- SuperCLUE模型象限
- 国内外大模型总体表现及竞争格局
- 大模型对战胜率、成熟度指数及开源格局

04 SuperCLUE2.0专项与行业

- 各行业及专项测评基准
- 未来两个月基准发布计划

05 四大维度测评分析及示例介绍

- 四大维度测评结果及示例

06 优秀模型案例介绍

- 优秀模型案例介绍

第1部分

2024年2月大模型关键进展



1.2023-2024大模型关键进展 2.中文大模型全景图 3.国内外大模型发展趋势

2023- 2024大模型关键进展

◆自2022年11月30日ChatGPT发布以来，AI大模型在全球范围内掀起了有史以来规模最大的人工智能浪潮。国内学术和产业界在过去一年也有了实质性的突破。大致可以分为三个阶段，即准备期（ChatGPT发布后国内产学研迅速形成大模型共识）、成长期（国内大模型数量和质量开始逐渐增长）、爆发期（各行各业开源闭源大模型层出不穷，形成百模大战的竞争态势）。

SuperCLUE：AI大模型2023年关键进展

（关键进展）



时间

2024年值得关注的中文大模型全景图

通用大模型

闭源



开源



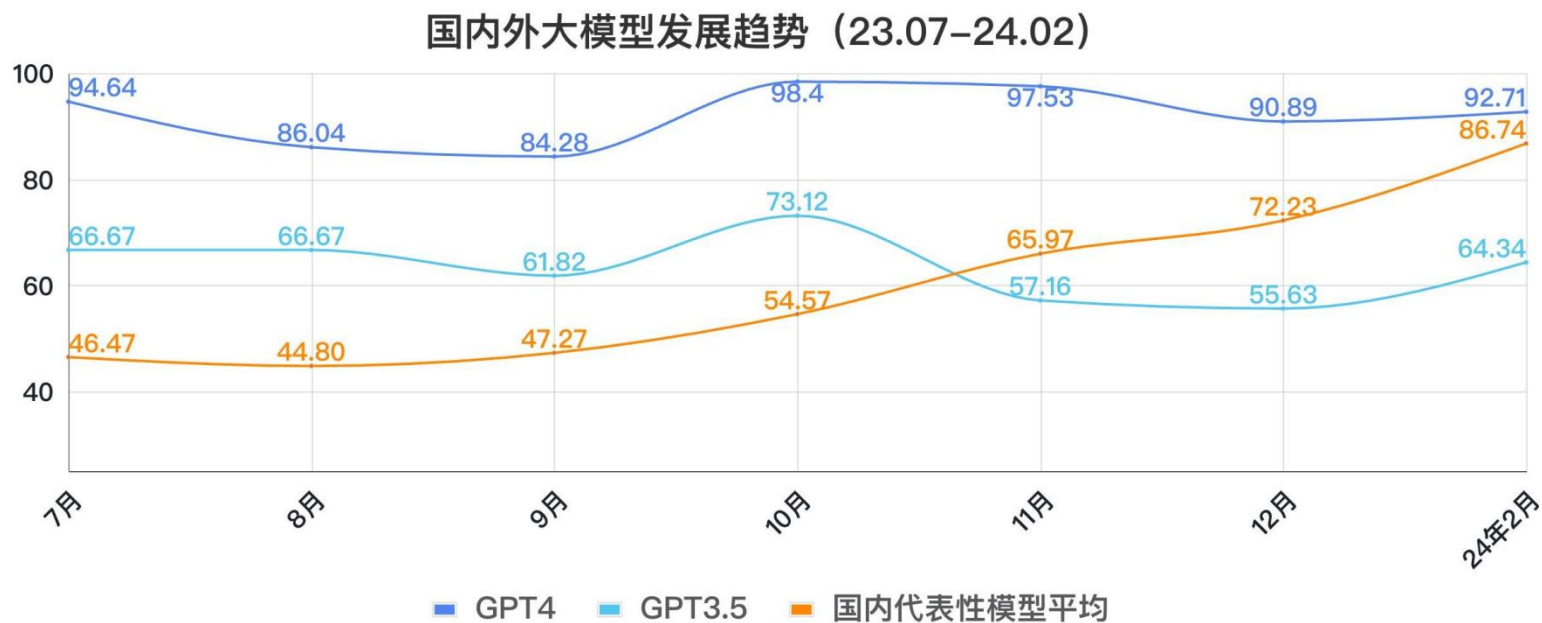
行业大模型

部分领域



国内外大模型发展趋势

过去八个月国内外代表性模型的发展趋势



部分国内代表性模型SuperCLUE基准得分(23年7月-24年2月)

模型	7月	8月	9月	10月	11月	12月	24年2月
文心一言	50.48	54.18	53.72	61.81	73.62	75	87.75
通义千问	-	41.73	33.78	43.36	61.01	71.78	85.70
ChatGLM	42.46	38.49	54.31	58.53	63.27	69.91	86.77

趋势说明

过去1年，国内领军大模型企业实现了大模型代际追赶的奇迹，从7月份与GPT3.5的20分差距，每个月都有稳定且巨大的提升，到24年2月份测评时已经完成总分上对GPT3.5的超越。

我们可以看到GPT 3.5和GPT 4在中文上的表现情况基本一致，在11月份测评结果中显示，在中文能力都有一定的下滑，而国内头部模型则展现了继续稳健提升的能力。在24年2月份的测评结果中可以看到，国内第一梯队模型与GPT4的差距在持续缩小。

说明：

趋势展示，选取了23年7月-24年2月Super CLUE- OPEN测评分数。国内代表性模型，选取了文心一言、通义千问、Chat GLM。原因是综合考虑了过去半年Super CLUE测评结果、长期稳定迭代及对国内大模型生态的贡献；GPT4成绩，由GPT4-API(7-9月)与GPT4-Turbo(10-2月)组成，用以表现国外最好模型发展。

第2部分

测评体系



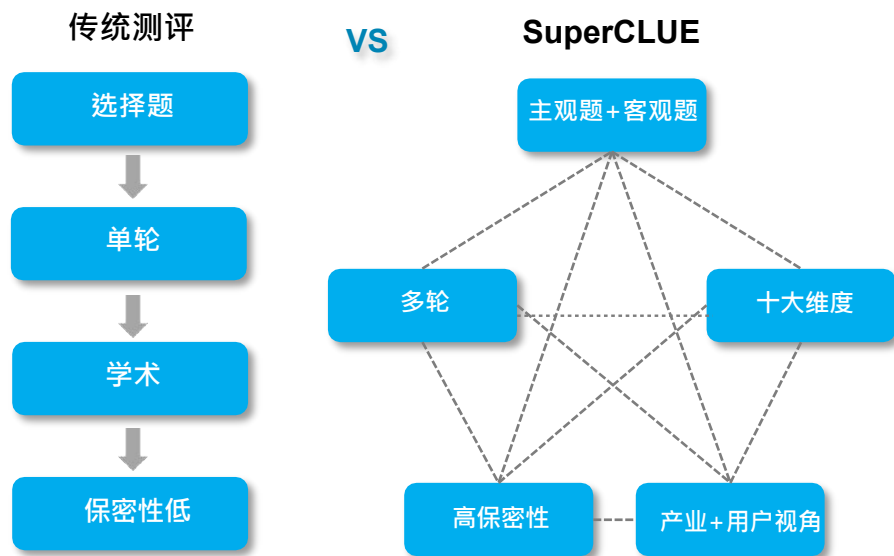
1.SuperCLUE介绍 2.测评体系 3.测评方法及示例

SuperCLUE介绍

中文语言理解测评基准 **CLUE** (**The Chinese Language Understanding Evaluation**) 是致力于科学、客观、中立的语言模型评测基准，发起于2019年。陆续推出CLUE、FewCLUE、KgCLUE、DataCLUE等广为引用的测评基准。

Super CLUE是大模型时代CLUE基准的发展和延续。聚焦于通用大模型的综合测评。传统语言模型测评往往局限于学术范围单轮选择题，Super CLUE根据多年的测评经验，基于通用大模型在学术、产业与用户侧的广泛应用，构建了多层次、多维度的综合性测评基准。

传统测评与SuperCLUE的区别



SuperCLUE 三大特征

第三方测评，不与模型厂商竞争

01

SuperCLUE始终秉持中立、客观的第三方测评理念，不会预设立场或偏向特定的模型方。同时，SuperCLUE采用自动化方式的客观评估，大大降低评测过程中的人为评估的不确定性。

测评方式与真实用户体验目标一致

02

不同于传统测评通过选择题形式的测评，SuperCLUE目标是与真实用户体验目标保持一致，所以纳入了开放主观问题的测评。通过多维度多视角多层次的测评体系以及对话的形式，真实模拟大模型的应用场景，真实有效的考察模型生成能力。

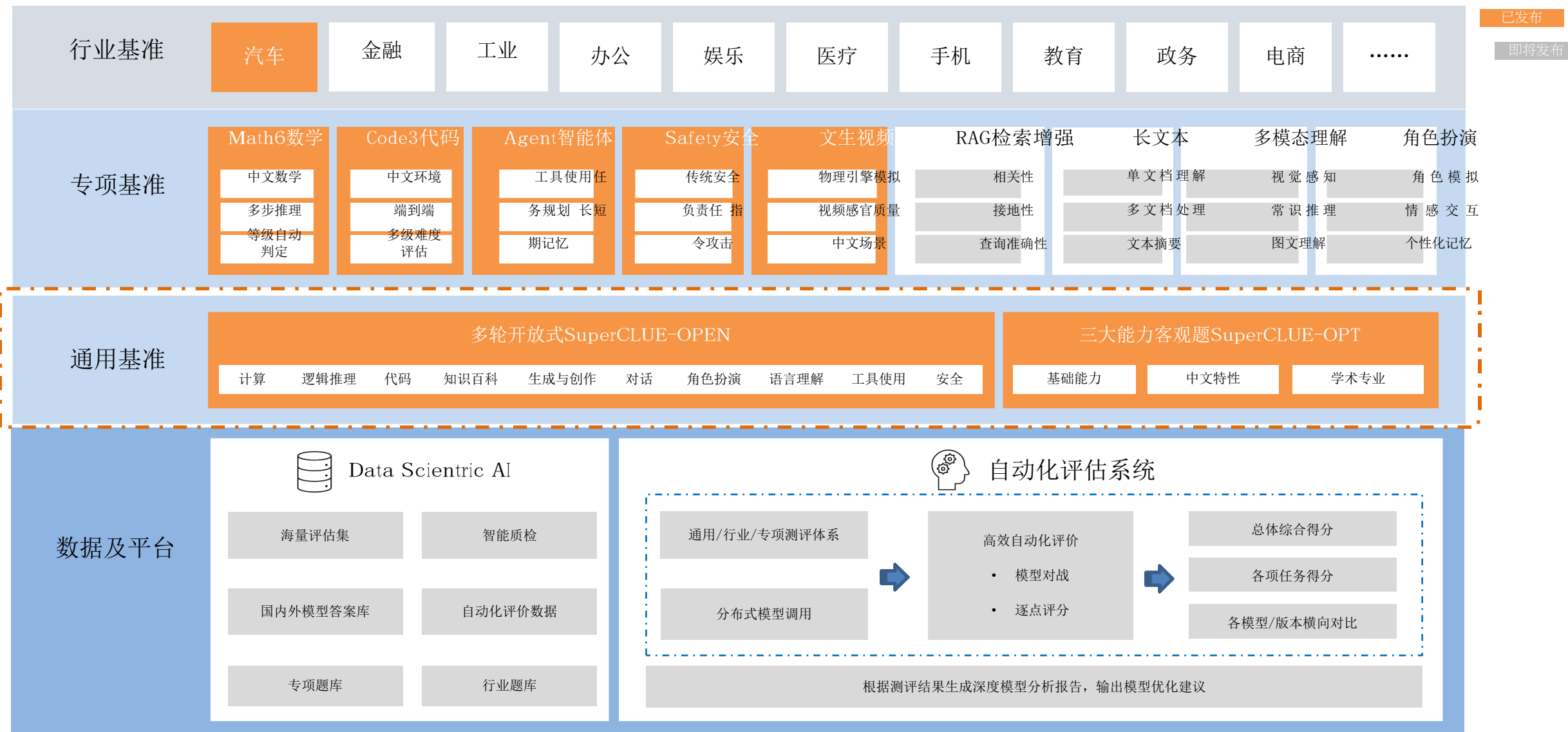
同时，SuperCLUE通过构建多轮对话场景，更深层次考察大模型在真实多轮对话场景的应用效果。对大模型的上下文、记忆、对话能力全方位评测。

不限于学术领域的测评，更为了服务产业界

03

不同于传统学术领域的评测，SuperCLUE从通用基准维度的选择、安全和智能体专项测评的设计，到行业大模型测评基准的推出，所有评测的目的都是为产业和应用服务。真实反应通用大模型与产业应用之间的差距，引导大模型提升技术落地效果，在通用能力的基础上更好的进行垂直领域的应用。

测评体系

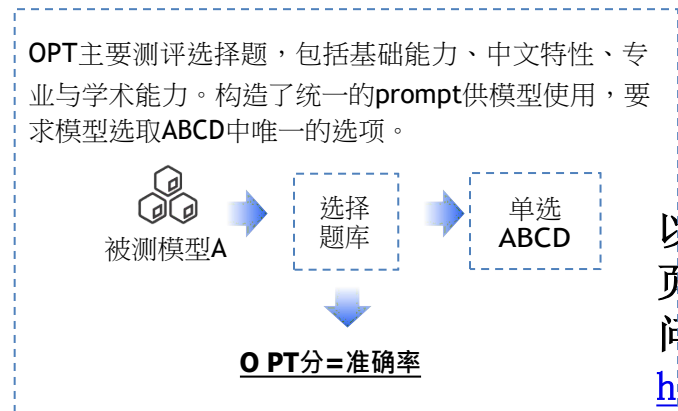
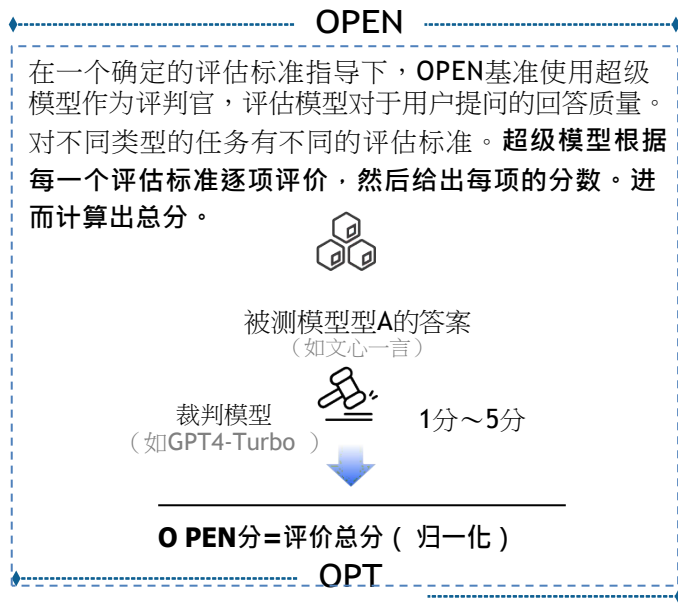


已发布
即将发布

即将发布

测评方法

为更真实反应大模型能力，本次测评采用多维度、多视角的综合性测评方案，由多轮开放问题Super CLUE- OPEN和三大能力客观题Super CLUE- OPT两部分测评结果组成。测评集共4572题，其中1504道多轮简答题（OPEN），3068道客观选择题（OPT），以下为评测集与测评方法简述。



$$\text{SuperCLUE总分} = 8 * \text{OPEN分} + 0.2 * \text{OPT分}$$

注：多轮简答题OPEN更能反应模型真实能力，故权重设置提高。

以上内容仅为本文档的试下载部分，为可阅读页数的一半内容。如要下载或阅读全文，请访问：

<https://d.book118.com/538134106062006061>