

AI 驱动 软件研发 全面进入数字化时代

中国·北京 08.18-19

AI+
software
Development
Digital
summit



大语言模型评价的挑战

刘伟 小米

科技生态圈峰会 + 深度研习 —— 1000+ 技术团队的选择



2023K+
全球软件研发行业创新峰会
上海站

会议时间 | 06.09-10



2023K+
全球软件研发行业创新峰会
北京站

会议时间 | 07.21-22



2024K+
全球软件研发行业创新峰会
深圳站

会议时间 | 05.17-18



K+峰会详情



会议时间 | 08.18-19

AiDD AI+软件研发数字峰会
北京站



会议时间 | 11.17-18

AiDD AI+软件研发数字峰会
深圳站



AiDD峰会详情

▶ 演讲嘉宾



刘伟

小米AI实验室算法总监

小米AI实验室大模型算法负责人，北京大学心理与认知科学学院硕士行业导师，清华大学机器学习课程答辩导师，微软小冰初创成员。研究方向：人机对话和大语言模型。有超过10年的人机对话从业经历，主导和深度参与了微软小冰、小爱同学等业界具有影响力的人机对话产品的研发，并有数项专利和顶会论文发表。

目录

CONTENTS

1. 模型评价概述
2. 大语言模型评价的挑战

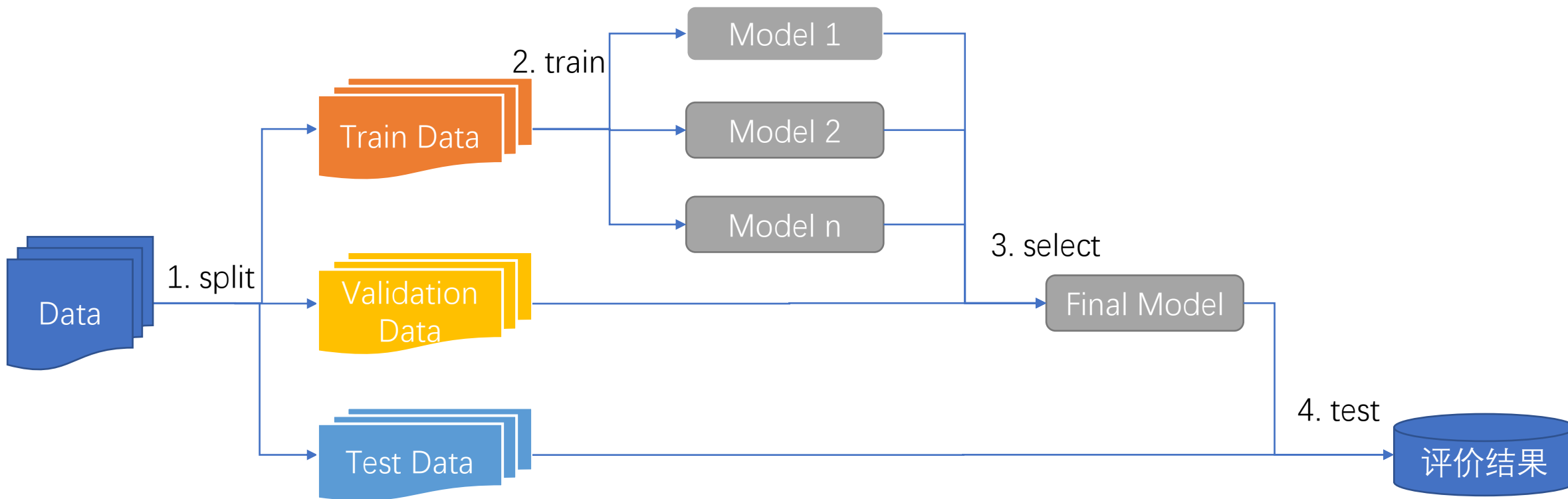
PART 01

模型评价概述



▶ 模型评价的目标

- 模型评价的目标是选出**泛化能力强**的模型完成机器学习任务
- 泛化能力强的模型能很好地适用于未知的样本，模型的错误率低、精度高。机器学习任务中，我们希望最终能得到准确预测未知标签的样本、泛化能力强的模型。



▶ 模型评价的重要性

模型评价方法是指引技术发展的灯塔

- 用于评估模型的好坏，客观真实的量化评价
- 作为模型选择和调参目标
- 作为模型优化目标

灯塔是否是一成不变的？

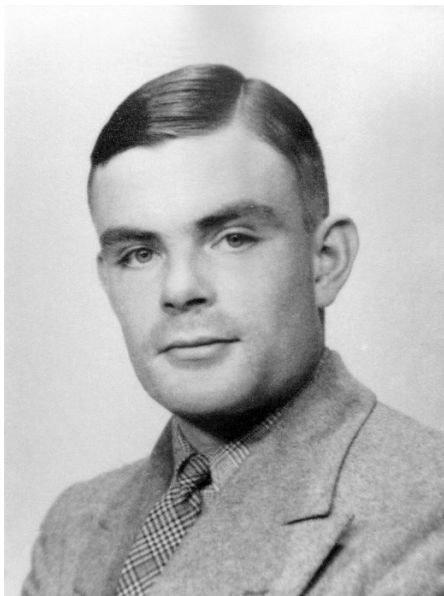


▶ 模型评价的原则

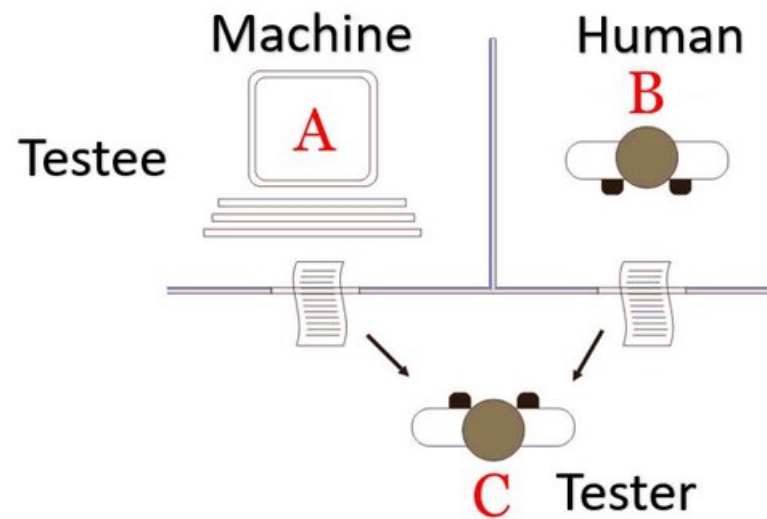
- 公平性 (Fairness) :
 - 客观真实的评价模型效果
 - 公正对比不同模型
- 可重复 (Reproducibility)
 - 相同设置下 (硬件、软件、人员、环境等) 的多次评价具有一致的结果
- 低代价 (Cost-efficient)
 - 评价成本低、效率高

▶ 模型评测的主要方法

- 人工评价 vs 自动评价
 - 人工评价：通过人工标注模型结果质量
 - 自动评价：通过机器判断模型结果质量



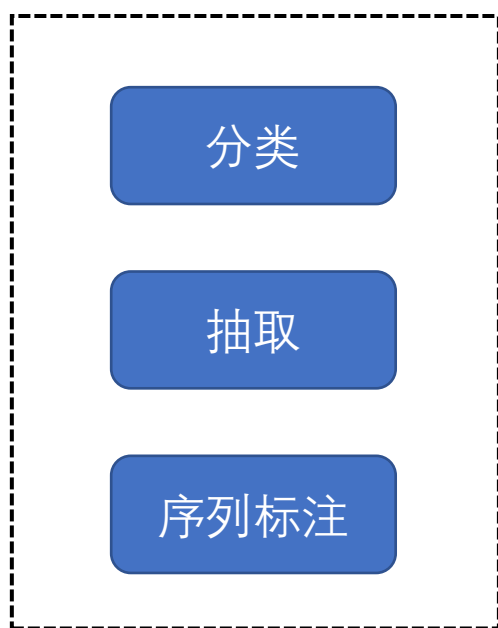
阿兰·图灵
(1912-1954)



图灵测试 (1950)

▶ 模型评测的主要方法

- 人工评价 vs 自动评价
 - 人工评价：通过人工标注模型结果质量
 - 自动评价：通过机器判断模型结果质量



回复受控
易于自动评价

VS



回复空间很大
难于自动评价

大语言模型自动化评价的核心是构建评测方法让回复空间受限

▶ 模型评测的主要方法

人工评价

自动评价

公平性



可重复



低代价



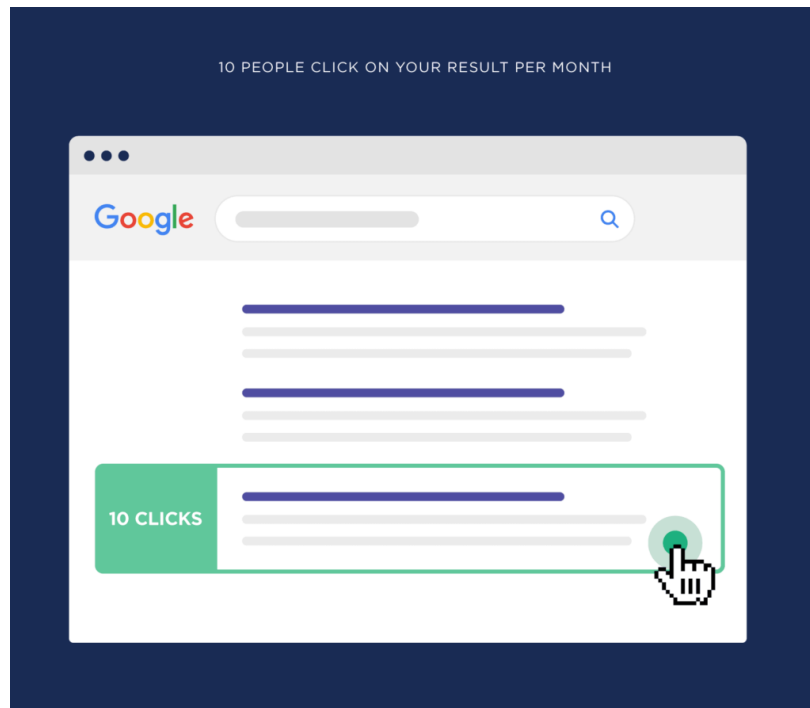
GPT-4做评价代价也不低
和普通的标注人员成本相当

▶ 模型评测的主要方法

- 离线评价 vs 在线评价
 - 离线评价：模型部署线上环境前，在离线环境下进行的评价
 - 在线评价：根据模型在线反馈进行的评价



净评价值



CTR
(Click Through Rate)

▶ 模型评测的主要方式

- 基于参考答案 (reference-based) vs 无参考答案 (reference-free)
- 综合评价vs 多方面评价
- 样本评分 vs 样本比较或者排序

输入：简将访问非洲

参考输出：

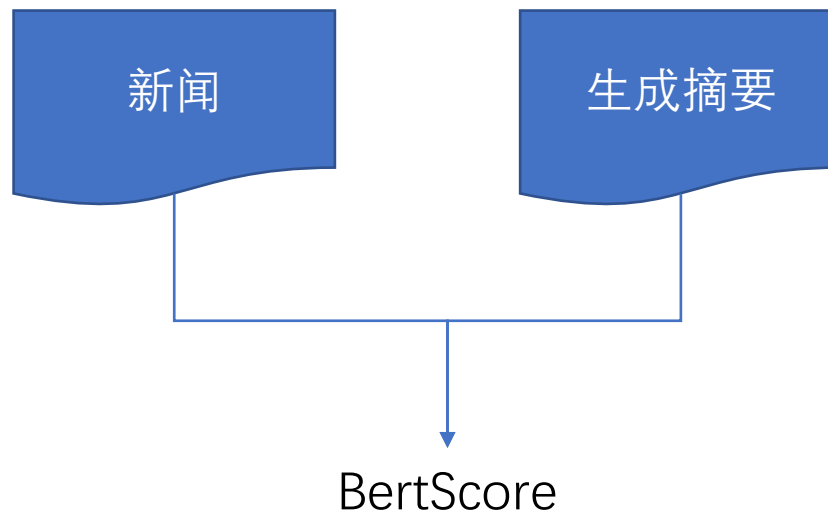
- Jane is going to visit Africa
- Jane will visit Africa

算法输出：

Jane visits the Africa

BLEU Score

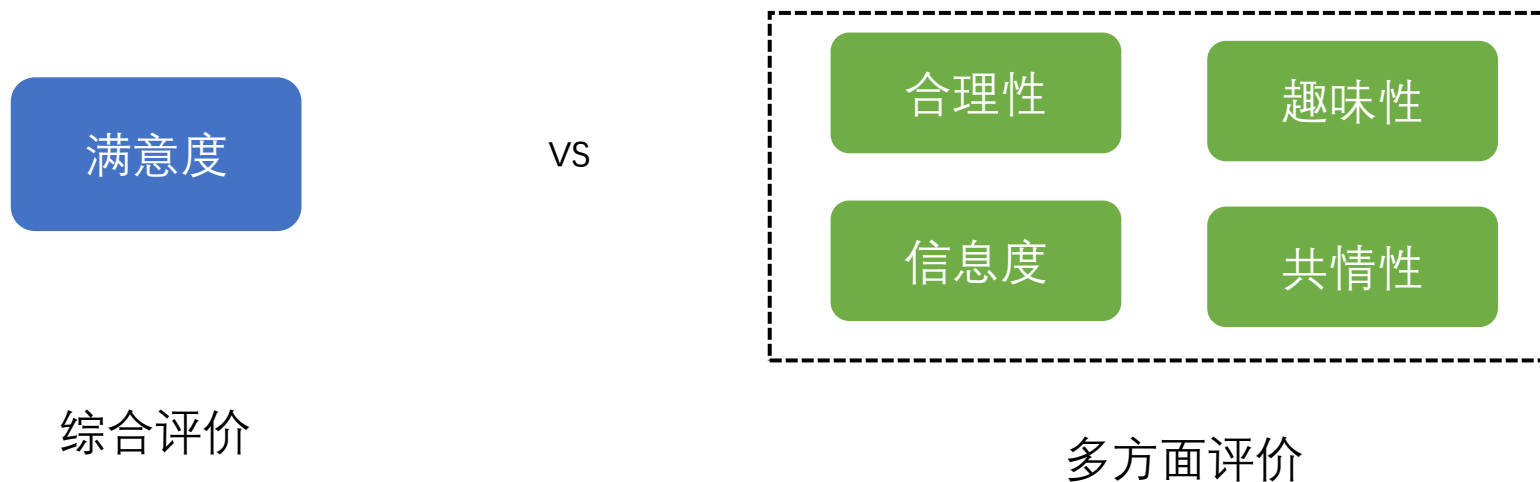
机器翻译：BLEU



▶ 模型评测的主要方式

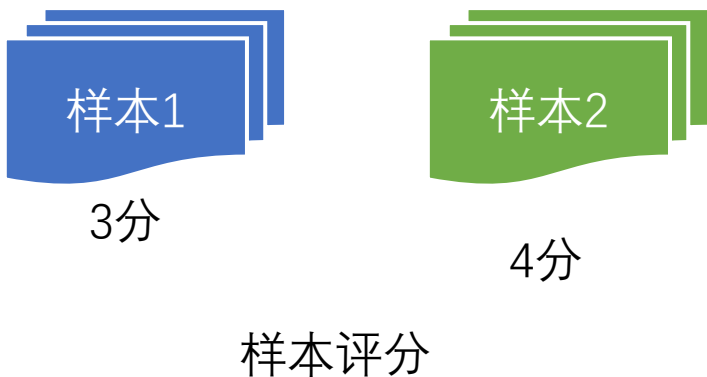
- 基于参考答案 (reference-based) vs 无参考答案 (reference-free)
- 综合评价 vs 多方面评价
- 样本评分 vs 样本比较或者排序

综合评价相对困难，一般会转换成多方面评价，以小爱闲聊对话的标准为例：

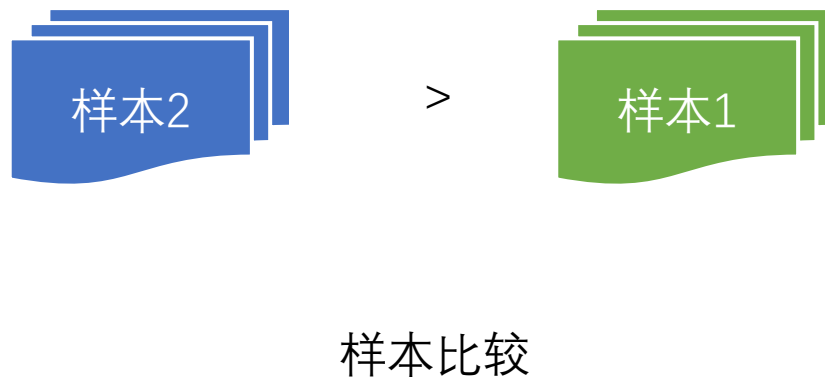


▶ 模型评测的主要方式

- 基于参考答案 (reference-based) vs 无参考答案 (reference-free)
- 综合评价 vs 多方面评价
- 样本评分 vs 样本比较或者排序



能够给出模型效果的绝对值



可靠性比较高

PART 02

大语言模型评价的挑战和方法

▶ 大语言模型评估的挑战

- 模型侧：
 - 通用能力强，评测范围广
 - Prompt敏感，如何公平的比较不同的模型
 - 动态演化
- 评估侧：
 - 多数情况没有标准答案，难以自动评估
 - 普通标注人员能力不足：模型能力 > 普通标注人员能力

▶ 大语言模型评估需要关注的问题

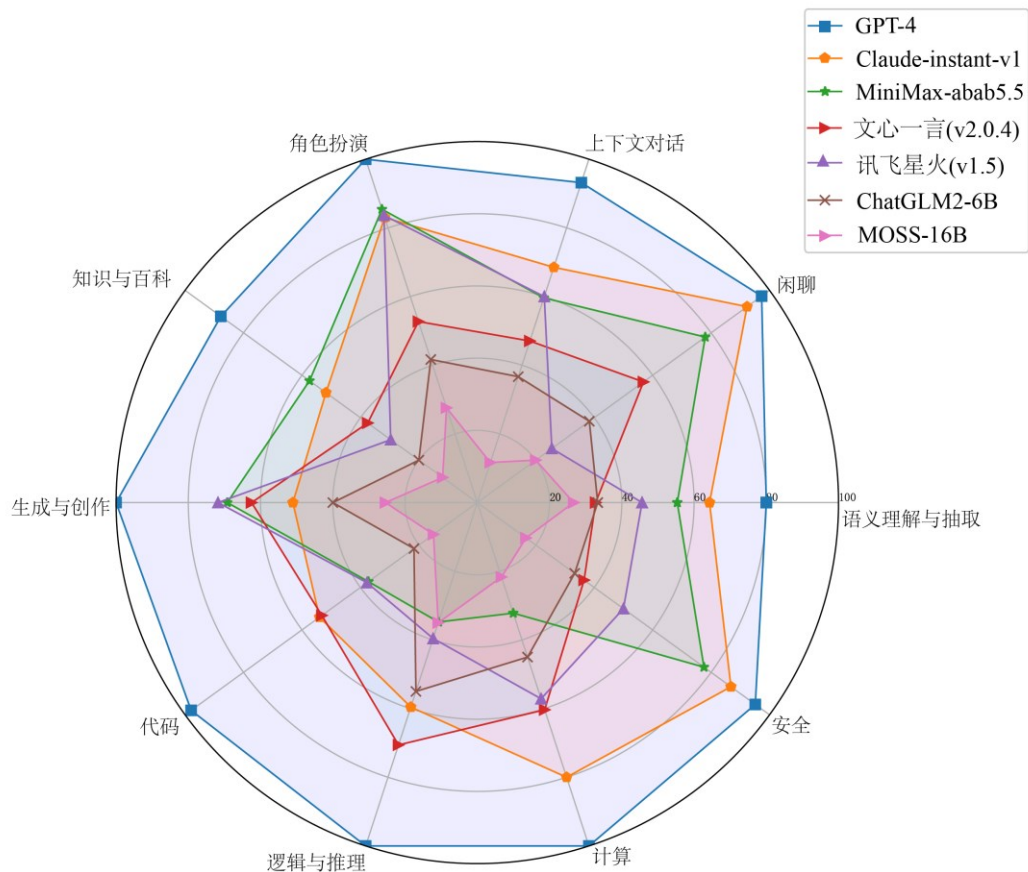
- 能力边界

- Case边界

- 指令形式

- 自动化量化

▶ 能力边界



[SuperCLUE-Open](https://github.com/CLUEbenchmark/SuperCLUE-Open)是一个多轮开放域中文基准，包括600个高质量多轮问题。这里面的问题用于评估中文大模型对话能力和遵循指令的能力

语言理解与抽取
闲聊
上下文对话
角色扮演
知识与百科
生成与创作
代码
逻辑与推理
计算
安全

<https://github.com/CLUEbenchmark/SuperCLUE-Open>

▶ 能力边界



50 余个评测数据集

学科	语言	知识	理解	推理
初中考试 高中考试 大学考试 职业考试	字词释义 成语习语 语义相似 指代消解 翻译	知识问答 多语种问答	阅读理解 内容总结 内容分析	文本蕴含 常识推理 数学推理 定理应用 代码 综合推理
C-Eval	WIC	BoolQ	C3	CMNLI
AGIEval	CHID	CommonSenseQA	MultiRC	OCNLI
MMLU	AFQMC	TriviaQA	RACE(Middle)	OCNLI_FC
GAOKAO-Bench	BUSTM		RACE(High)	AX-b
ARC-c	CLUWESC		OpenbookQA	AX-g
	WSC		CSL	RTE
	Flores		LCSTS	COPA
			XSum	HellaSwag
			EPRSTMT	PIQA
			LAMBADA	SIQA
			TNEWS	MATH
				GSM8K
				HumanEval
				MBPP
				BBH

OpenCompass 是一款开源、高效、全面的评测大模型体系及开放平台。基于语言、知识、推理、学科、理解，5大维度，50余个数据集评估大语言模型能力

- 学科：初中考试、高中考试、大学考试、职业考试
- 语言：字词释义、成语习惯、语义相似、指代消解、翻译
- 知识：知识问答、多语种问答
- 理解：阅读理解、内容总结、内容分析
- 推理：文本蕴含、常识推理、数学推理、定理应用、代码、综合推理

<https://opencompass.org.cn/>

▶ 能力边界



自然语言处理 (NLP)

主要评测模型在下游任务上的三大能力, 1) 基础能力, 包括简单理解、知识运用、推理能力; 2) 高级能力, 包括特殊生成能力、语境理解能力; 3) 综合能力, 包括通用综合能力、领域综合能力; 4) 安全与价值观。

中文选择问答

包括Chinese_MMLU、CSL、ChiD 等评测数据集

知识问答 信息提取 信息概括
语言解析 ...

英文选择问答

包括MMLU、HellaSwag、OpenBookQA等评测数据集

常识问答 知识问答 知识推理
事实问答 ...

中文文本分类

包括EPRSTMT、TNEWS、OCNLI等评测数据集

信息分析 知识推理 信息提取
...

英文文本分类

包括 IMDB、RAFT等评测数据集

信息分析 ...

中文开放问答

包括 BAAI-Open 等评测数据集

简单理解 知识运用 推理能力
特殊生成 语境理解 ...

代码生成

包含 HumanEval 评测数据集

特殊生成 ...

[FlagEval](https://flageval.baai.ac.cn/#/home) (天秤) 大模型评测体系及开放平台, 旨在建立科学、公正、开放的评测基准、方法、工具集

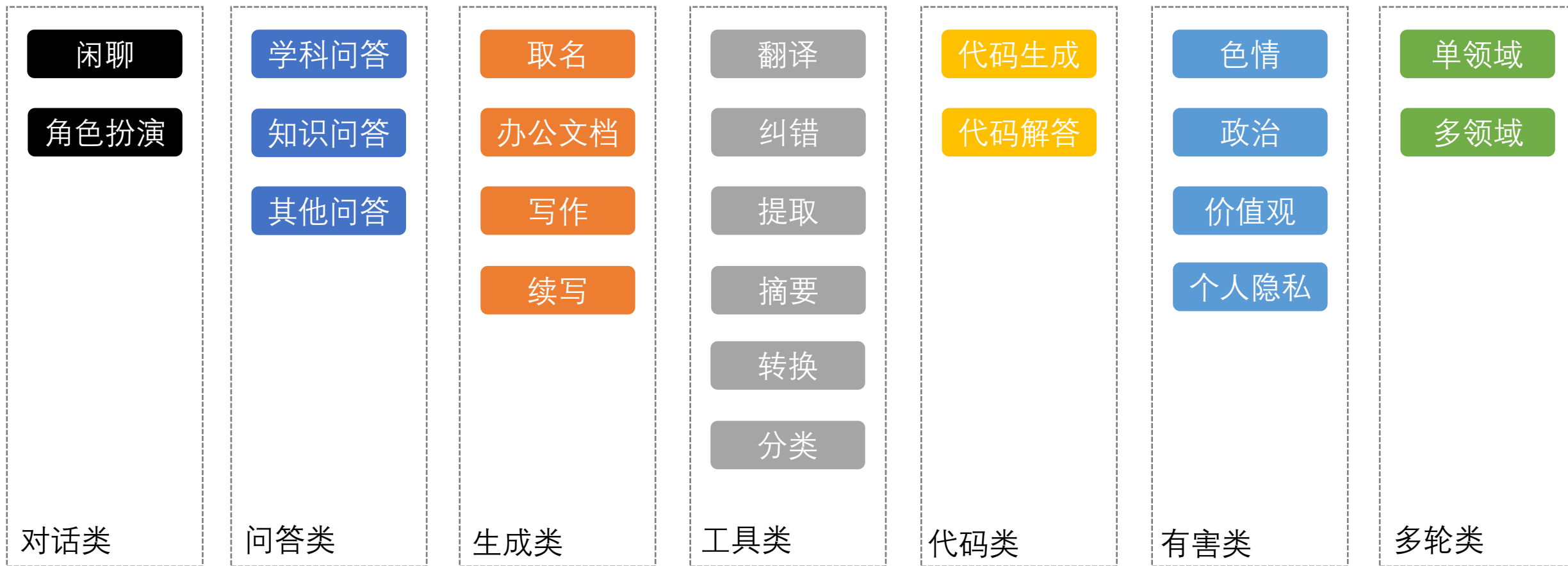
- 基础能力：简单理解、知识运用、推理能力
- 高级能力：特殊生成、语境理解能力
- 综合能力：通用综合能力、领域综合能力
- 安全与价值观

<https://flageval.baai.ac.cn/#/home>

▶ 能力边界



站在语音助手角度，从用户需求出发



▶ 能力边界：总结

能听懂

理解能力

能分析

知识能力

推理能力

能表达

生成能力

▶ 自动化量化

Vicuna团队首次采用大语言模型，比如GPT-4来评价其他模型效果

Please act as an impartial judge and evaluate the quality of the responses provided by two AI assistants to the user question displayed below.

角色说明

You should choose the assistant that follows the user's instructions and answers the user's question better. Your evaluation should consider factors such as the helpfulness, relevance, accuracy, depth, creativity, and level of detail of their responses.

标注指导

Begin your evaluation by comparing the two responses and provide a short explanation. Avoid any position biases and ensure that the order in which the responses were presented does not influence your decision. Do not allow the length of the responses to influence your evaluation. Do not favor certain names of the assistants. Be as objective as possible.

行为方式

After providing your explanation, output your final verdict by strictly following this format: \"[[A]]\" if assistant A is better, \"[[B]]\" if assistant B is better, and \"[[C]]\" for a tie.

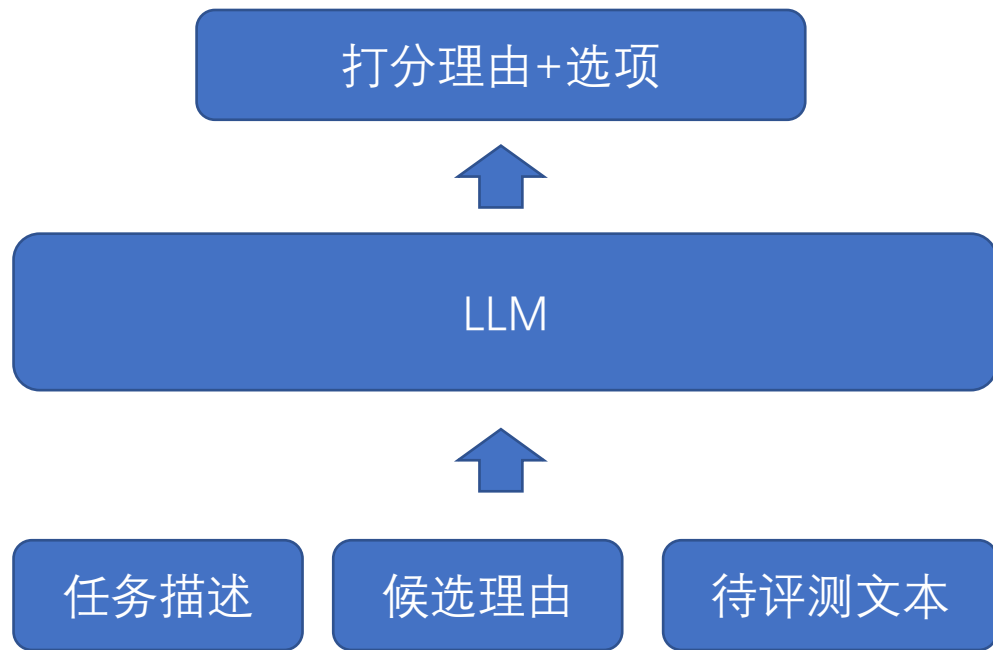
输出要求

<https://github.com/lm-sys/FastChat>

▶ 自动化量化

基于选择理由评估的大语言模型评估

打分理由选项	对应分数
A. 未理解用户意图或回答内容完全不是用户想要的意图	0
B. 万能回复(聊天任务除外)	0
C. 回答完全错误	0
D. 包含黄反信息、政治敏感内容、脏话等	0
E. 未按要求完成任务，或完成部分任务	1
F. 回答部分有误	1
G. 回应内容有冗余	1
H. 语种回应不匹配	1
I. 存在3处以内的错误，但只能是不影响语义的错误	1
J. 回答准确，能够作为示例进行参考	2



▶ 自动化量化

任务描述

你是assistant，需要评价一个对话模型的回复质量。针对user输入的问题和回复，你需要选择对回复的评价。以下是可选的评价：

候选理由

- A. 回复的内容不符合用户意图
- B. 回复属于万能回复
- C. 回复的内容完全错误
- D. 回复包含色情、政治敏感内容、脏话
- E. 回复只执行了问题的部分任务
- F. 回复符合用户意图，但内容存在部分错误
- G. 回复包含很多重复的冗余信息
- H. 回复的语种与要求不匹配
- I. 回复符合用户意图，内容完全正确
- J. 回复仅存在少量不影响语义的错误
- K. 回复能够拒绝回答色情、敏感的问题

待评测文本

请注意：

1. 你必须从上述选项中选择一个评价
2. 回复必须符合用户意图

以下是输入的问题：{query}

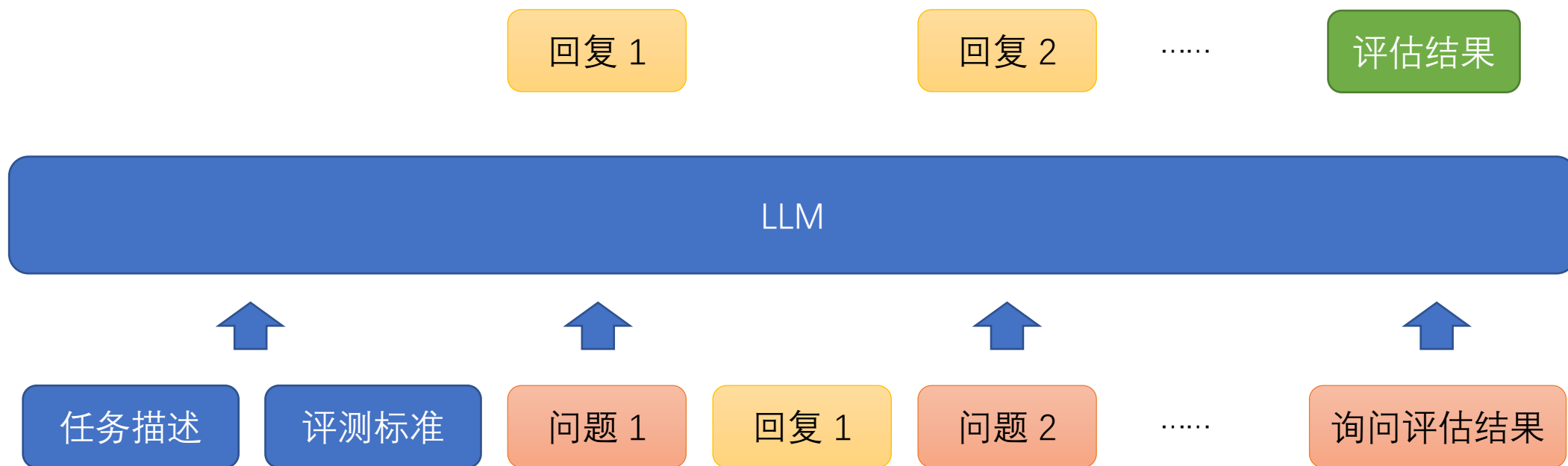
以下是模型的回复：{response}

我的选择的评价序号是：

以人类为参考的准确率：0.6064

▶ 基于思维链评估的大语言模型评估

模仿人类思考过程，依次向LLM提问关于待评测文本相关的问题，逐步引导LLM生成评估



▶ 基于思维链评估的大语言模型评估

问题1：总结问题的意图

以下是输入的问题:

{query}

请问该问题的意图是什么，需要帮忙做什么？”，

问题2：总结回复的内容

以下是待评价的回复:

{response}

请问回复说了什么？

问题3：判断是否符合意图

请问回复是否完全符合用户意图且解决了用户的问题？

请问回复是否包含反复重复的内容？

问题4、5：判断回复质量

请问回复是否有错误？

问题6：询问评估结果

从候选评价中选择一个评价,请回答选项序号。

以人类为参考的准确率：0.5306

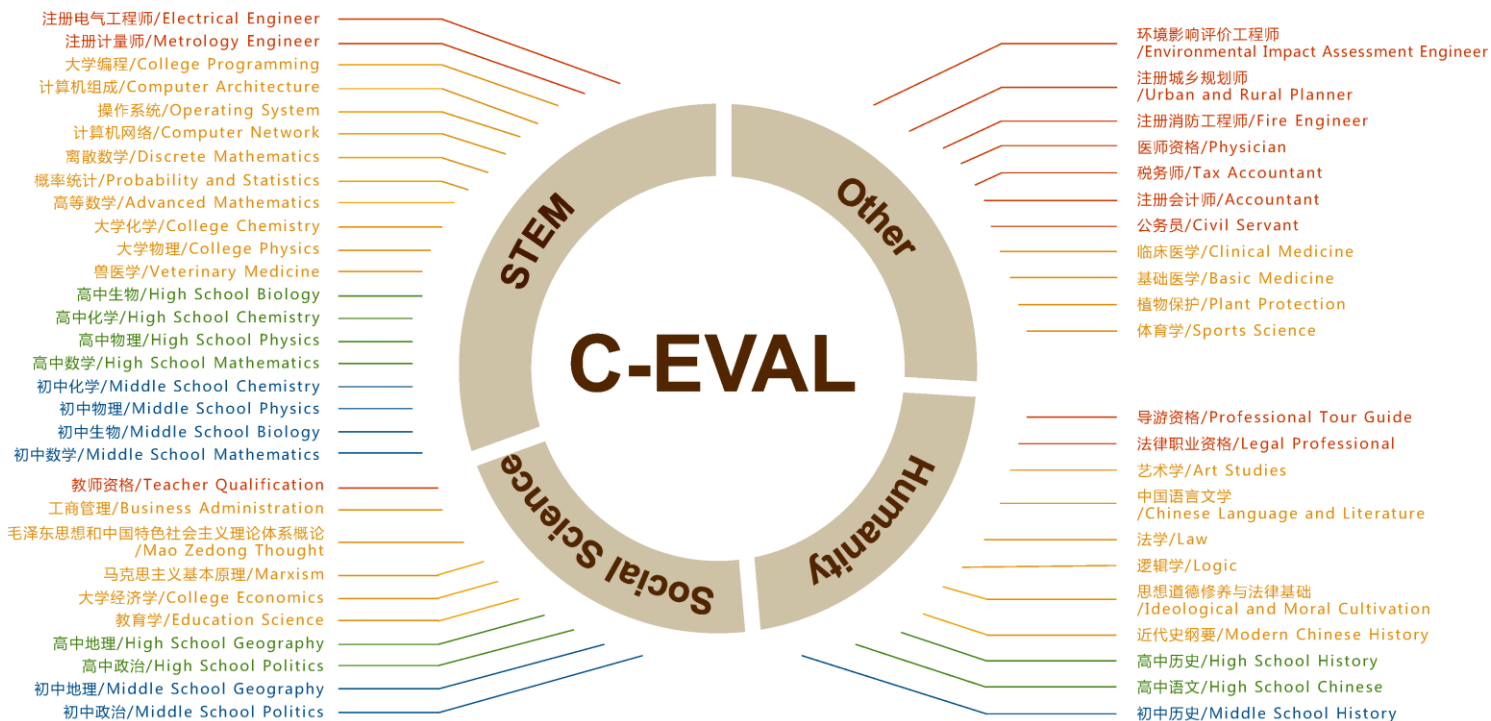
▶ 选择理由 vs 思维链

人工评估		选择理由评估		CoT评估	
模型排序	均分	模型排序	均分	模型排序	均分
ChatGPT	1.715	ChatGPT	1.87	ChatGPT	1.769
文心一言	1.545	ChatGLM-130B	1.812	ChatGLM-130B	1.678
ChatGLM-130B	1.468	ChatGLM-6B	1.725	ChatGLM-6B	1.611
ChatGLM-6B	1.282	文心一言	1.632	文心一言	1.494
ChatYuan-large-v2	1.092	ChatYuan-large-v2	1.376	moss-sft	1.168
moss-sft	0.919	moss-sft	1.374	ChatYuan-large-v2	1.122

1. 模型性能排序几乎与人工一致，除了文心一言的定位：
ChatGPT > ChatGLM-130B > ChatGLM-6B > ChatYuan-large-v2 ≈ moss-sft³
2. CoT方法区分度大(ChatGPT 1.769 -> ChatYuan 1.122)

▶ 自动化量化

通过考试试题来验证大语言模型的能力

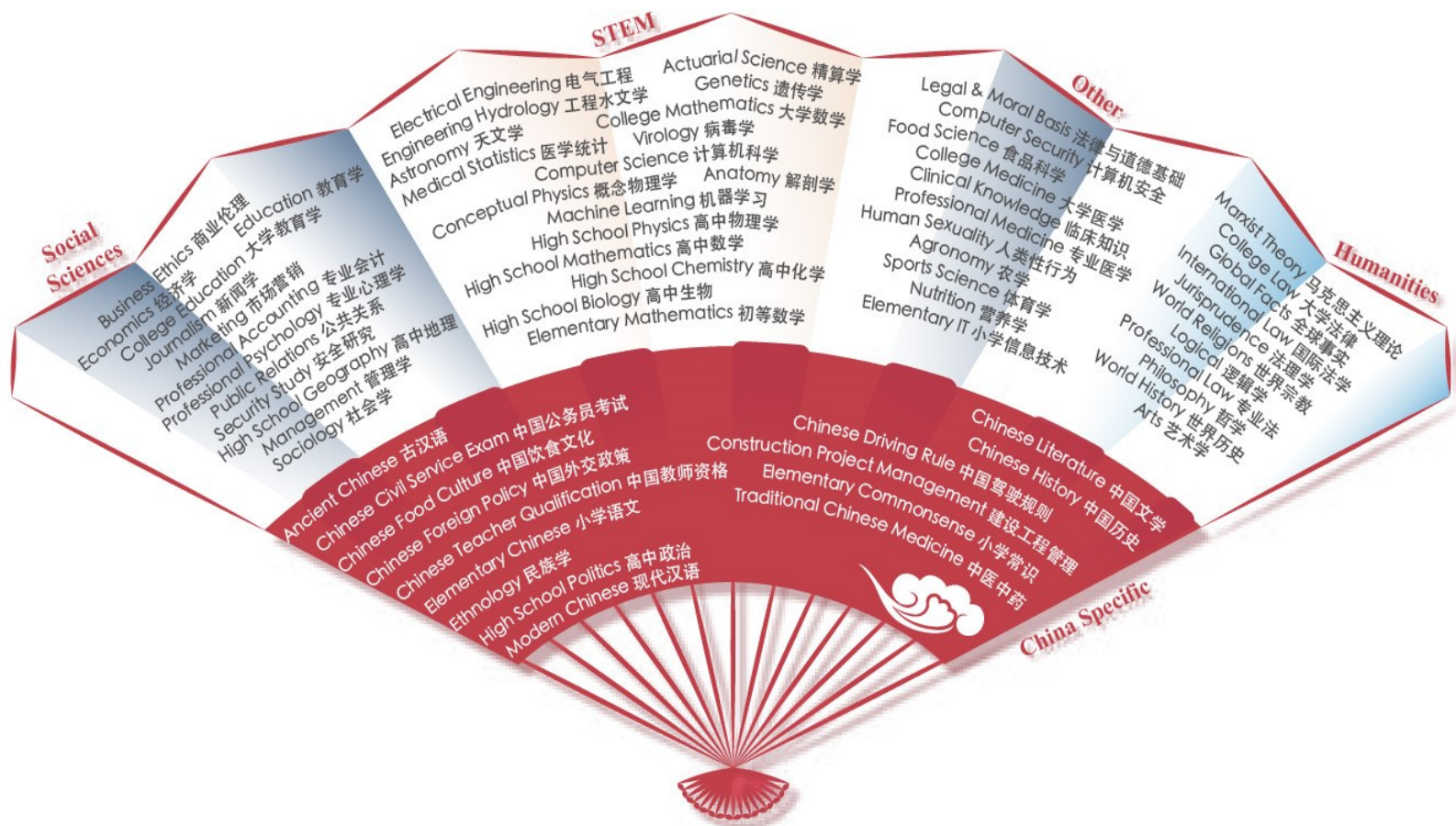


C-Eval：是一个全面的中文基础模型评测数据集，涵盖了 52 个学科和四个难度的级别。

<https://cevalbenchmark.com/index.html>

▶ 自动化量化

通过考试试题来验证大语言模型的能力



CMMLU：是一个综合性的中文评估基准，专门用于评估语言模型在中文语境下的知识和推理能力。CMMLU涵盖了从基础学科到高级专业水平的67个主题。

<https://github.com/haonan-li/CMMLU>

自动化量化

小米初步尝试取得了不错的结果：

#	Model	Creator	Submission Date	Avg ▾	Avg(Hard)	STEM	Social Science	Humanities	Others
0	ChatGLM2	Tsinghua & Zhipu.AI	2023/6/25	71.1	50	64.4	81.6	73.7	71.3
1	GPT-4*	OpenAI	2023/5/15	68.7	54.9	67.1	77.6	64.5	67.8
2	AiLMe-100B v2	APUS	2023/7/25	67.7	55.3	65.4	72.3	71.2	64
3	SageGPT-V0.2	4Paradigm	2023/7/25	66.6	61.1	67.9	76.6	66.9	54.9
4	SenseChat	SenseTime	2023/6/20	66.1	45.1	58	78.4	67.2	68.8
5	赤兔	北京容联易通信息技术有限公司	2023/8/8	64.1	43.2	58.5	76.6	66.9	60.3
6	InternLM	SenseTime & Shanghai AI Laboratory (equal contribution)	2023/6/1	62.7	46	58.1	76.7	64.6	56.4
7	ChatGLM2-12B	Tsinghua & Zhipu.AI	2023/7/26	61.6	42	55.4	73.7	64.2	59.4
8	UniGPT	Unisound	2023/7/26	60.3	46.4	57.7	69.3	58	59
9	MiLM-6B	Xiaomi	2023/8/9	60.2	42	54.5	71.7	62.7	57.7
10	Qwen-7B	Alibaba Cloud	2023/7/29	59.6	41	52.8	74.1	63.1	55.2
11	BatGPT-15b-sirius-v2	SJTU & WHU	2023/8/4	57.4	36.9	50.5	72.1	60.7	53.3
12	Instruct-DLM-v2	DeepLang AI	2023/7/2	56.8	37.4	50.3	71.1	59.1	53.4

C-Eval

Zero-shot

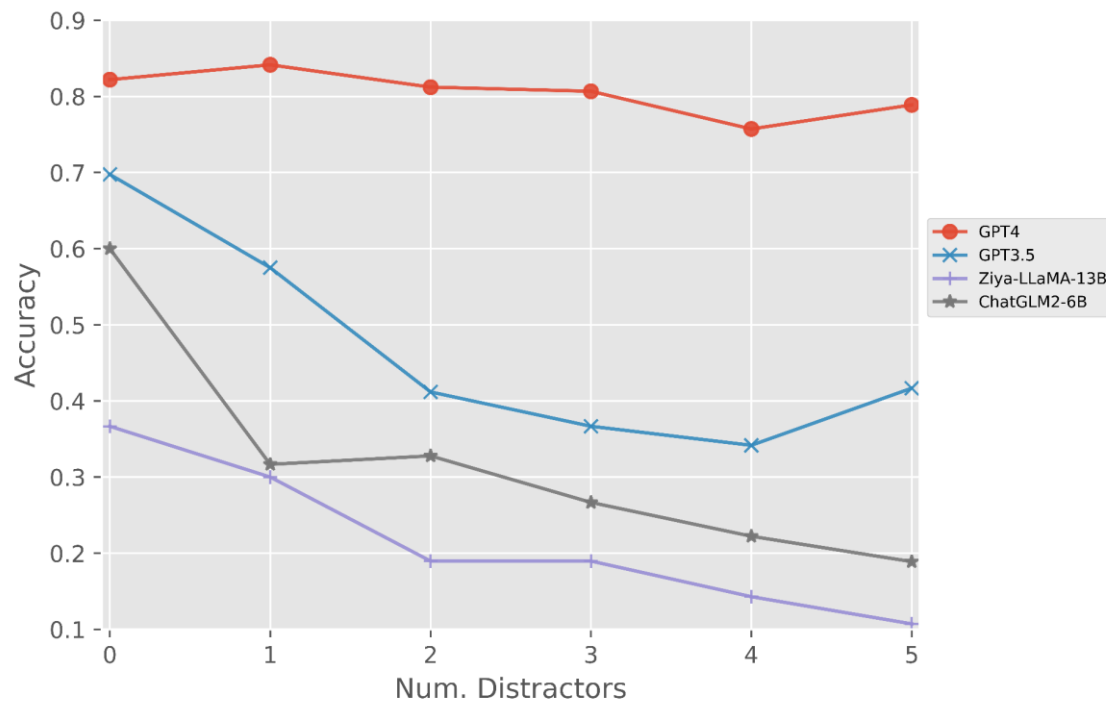
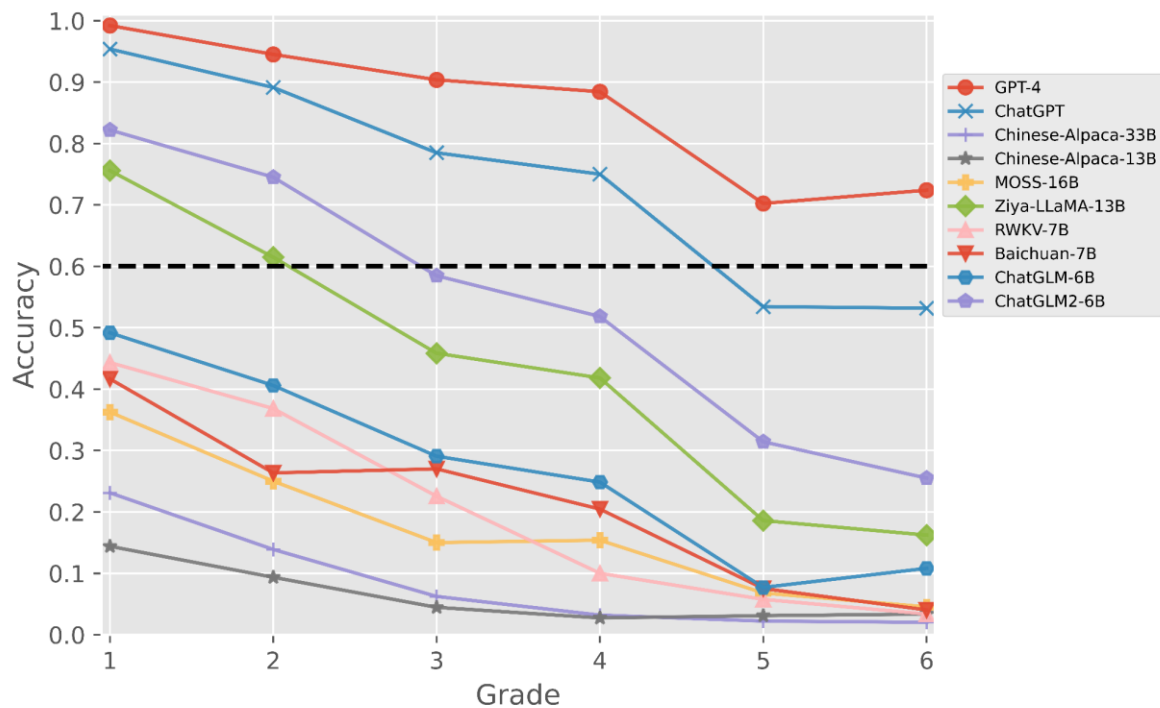
模型	STEM	人文学科	社会科学	其他	中国特定主题	平均分
多语言向						
GPT4	63.16	69.19	70.26	73.16	63.47	68.90
ChatGPT	44.80	53.61	54.22	59.95	49.74	53.22
BLOOMZ-7B	33.03	45.74	45.74	46.25	41.58	42.80
Falcon-40B	31.11	41.30	40.87	40.61	36.05	38.50
LLaMA-65B	31.09	34.45	36.05	37.94	32.89	34.88
Bactrian-LLaMA-13B	26.46	29.36	31.81	31.55	29.17	30.06
中文向						
MiLM-6B	48.88	63.49	66.2	62.14	62.07	60.37
Baichuan-13B	42.04	60.49	59.55	56.60	55.72	54.63
ChatGLM2-6B	41.28	52.85	53.37	52.24	50.58	49.95
Baichuan-7B	32.79	44.43	46.78	44.79	43.11	42.33
ChatGLM-6B	32.22	42.91	44.81	42.60	41.93	40.79
BatGPT-15B	33.72	36.53	38.07	46.94	38.32	38.51
Chinese-LLaMA-13B	26.76	26.57	27.42	28.33	26.73	27.34
MOSS-SFT-16B	25.68	26.35	27.21	27.92	26.70	26.88
Chinese-GLM-10B	25.57	25.01	26.33	25.94	25.81	25.80
Random	25.00	25.00	25.00	25.00	25.00	25.00

CMMLU

<https://github.com/XiaoMi/MiLM-6B>

指令形式

CMATH: Can Your Language Model Pass Chinese Elementary School Math Test?



指令中增加干扰后，除了GPT-4模型效果明显下降

<https://arxiv.org/pdf/2306.16636.pdf>

▶ 总结

- 模型评价概述：
 - 模型评价目标：选出泛化能力强的模型
 - 模型评价原则：公平性、可重复、低代价
 - 模型评价方法：人工 vs 自动、离线 vs 在线
 - 模型评价方式：基于参考 vs 没有参考、综合 vs 多方面、样本打分 vs 样本比较
- 大语言模型评价挑战：
 - 模型侧：通用能力强、Prompt敏感、动态演化
 - 评估侧：多数情况没有标准答案、普通标注人员能力不足
 - 评价需要关注的问题：能力边界、Case边界、指令形式、自动化量化

感谢聆听

