# 蚂蚁代码大模型的评测实践

申敏　蚂蚁集团

# 科技生态圈峰会 + 深度研习

## ——1000＋技术团队的共同选择

**KEYLINK ing**

### K⁺峰会

**K⁺峰会 深圳站**
K⁺ 全球软件研发行业创新峰会
会议时间：2024.05.24-25

**K⁺峰会 上海站**
K⁺ 全球软件研发行业创新峰会
会议时间：2024.09.20-21

### AiDD峰会

**AiDD峰会 深圳站**
AI⁺ 软件研发数字峰会
会议时间：2023.11.24-25

**AiDD峰会 北京站**
AI⁺ 软件研发数字峰会
会议时间：2024.07.19-20

**AiDD峰会 深圳站**
AI⁺ 软件研发数字峰会
会议时间：2024.11.15-16

# ▶ 演讲嘉宾

**申敏**

蚂蚁集团-测试开发专家

蚂蚁集团测试开发专家，研究方向：大模型在代码领域的评测技术。
长期投入蚂蚁支付、账务、计收费等业务领域质量保障工作，熟悉企业级
编码风格及要求，当前，负责蚂蚁百灵大模型CodeFuse系列的代码能力
评测。

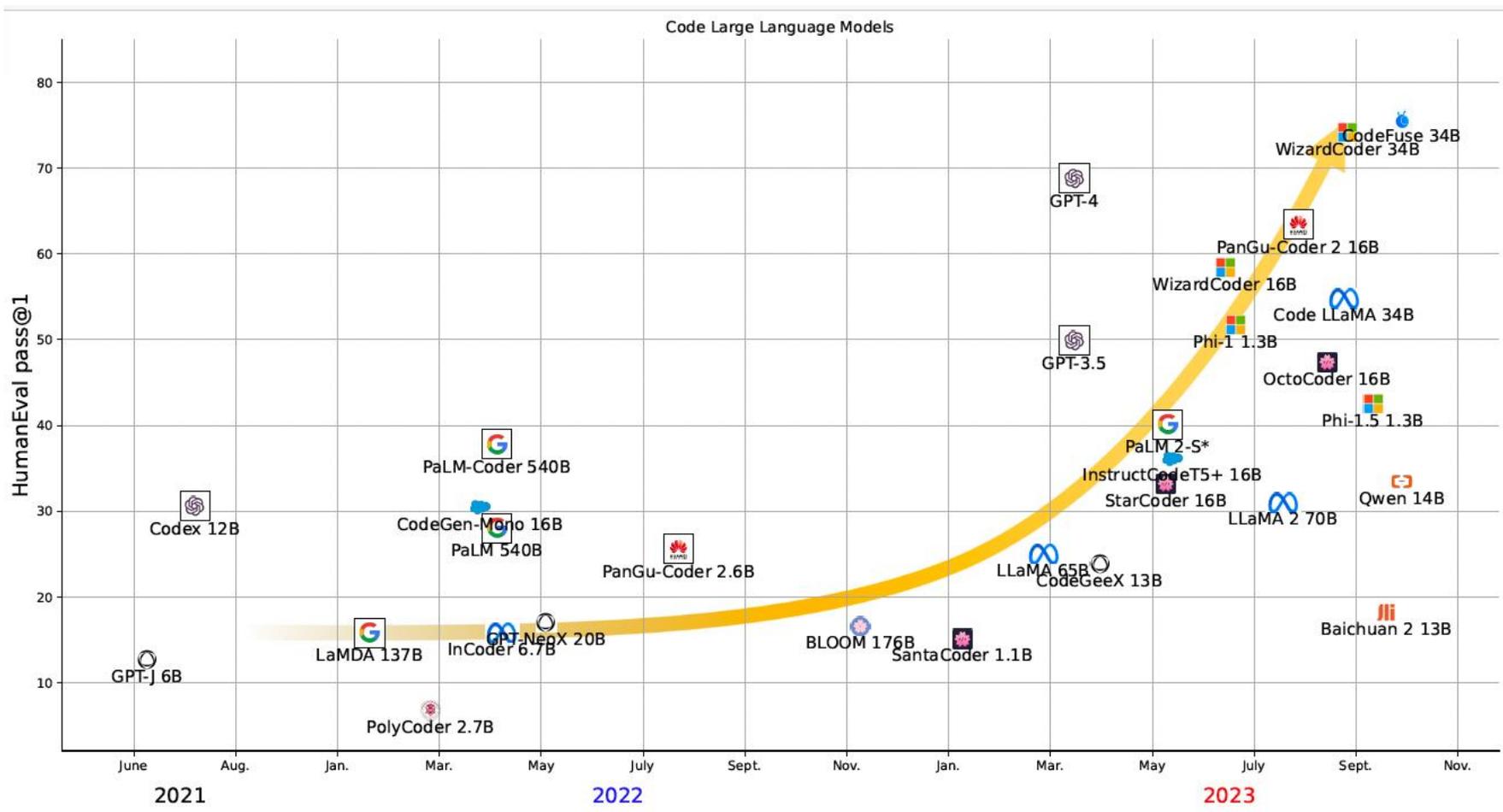# 目 录
## CONTENTS

PART 01
前言

# 前言：模型发展与模型评估

2023年大模型呈爆发式增长，截至2023年7月，中国累计有130个大模型问世，国外大模型138个，其中，美国大模型114个。*[赛迪顾问《IT2023秋》]* 模型发布必然离不开模型评估；AIGC编程是模型落地最为广泛的场景之一，充分衡量方能更好的运用or选用。



Code Large Language Models

# ▶ 前言：**CodeFuse 让研发变的更简单**

CodeFuse 是一款为国内开发者提供智能研发服务的产品，该产品是基于蚂蚁集团自研的基础大模型进行微调的**代码大模型**。

CodeFuse 具备**代码补全、添加注释、解释代码、生成单测，以及代码优化等功能**，以帮助开发者更快、更轻松地编写代码。



目前支持10款IDE

40+编程语言

包括python/java/js等

官网:https://codefuse.alipay.com/welcome/product

# ▶ 前言：CodeFuse模型

CodeFuse 模型：旨在支持整个软件开发生命周期的大型代码语言模型（Code LLMs），涵盖设计、需求、编码、测试、部署、运维等关键阶段。



**业界开源评测集上的roadmap**

Table 1: CODEFUSE project roadmap.

| Release date | Model | HUMANEVAL Pass@1 |
|---|---|---|
| Mar 2023 | CODEFUSE-1.3B-2K Seq-Length | 11.58% |
| Apr 2023 | CODEFUSE-6.5B-4K Seq-Length | 20.46% |
| May 2023 | CODEFUSE-13B-Base-4K Seq-Length | 32.93% |
| Jun 2023 | CODEFUSE-13B (opened in Sep) | 37.10% |
| Sep 2023 | CODEFUSE-CODELLAMA-34B (opened) | 74.40% |

CodeFuse系列论文：

https://arxiv.org/abs/2311.02303
https://arxiv.org/abs/2310.06266

# ▶ 前言：大模型时代如何评估代码大模型

通用

代码
（垂类）

伴随蚂蚁代码大模型的投产，我们发现代码领域打榜与实际投产存在一定的差异，基于此，我们探索并构建了适合企业项目的代码大模型的评测范式。

多维： 代码能力、 基础能力、 安全能力等多维度

多样： 多编码语言、 编码规范、跨项目编码等任务多样
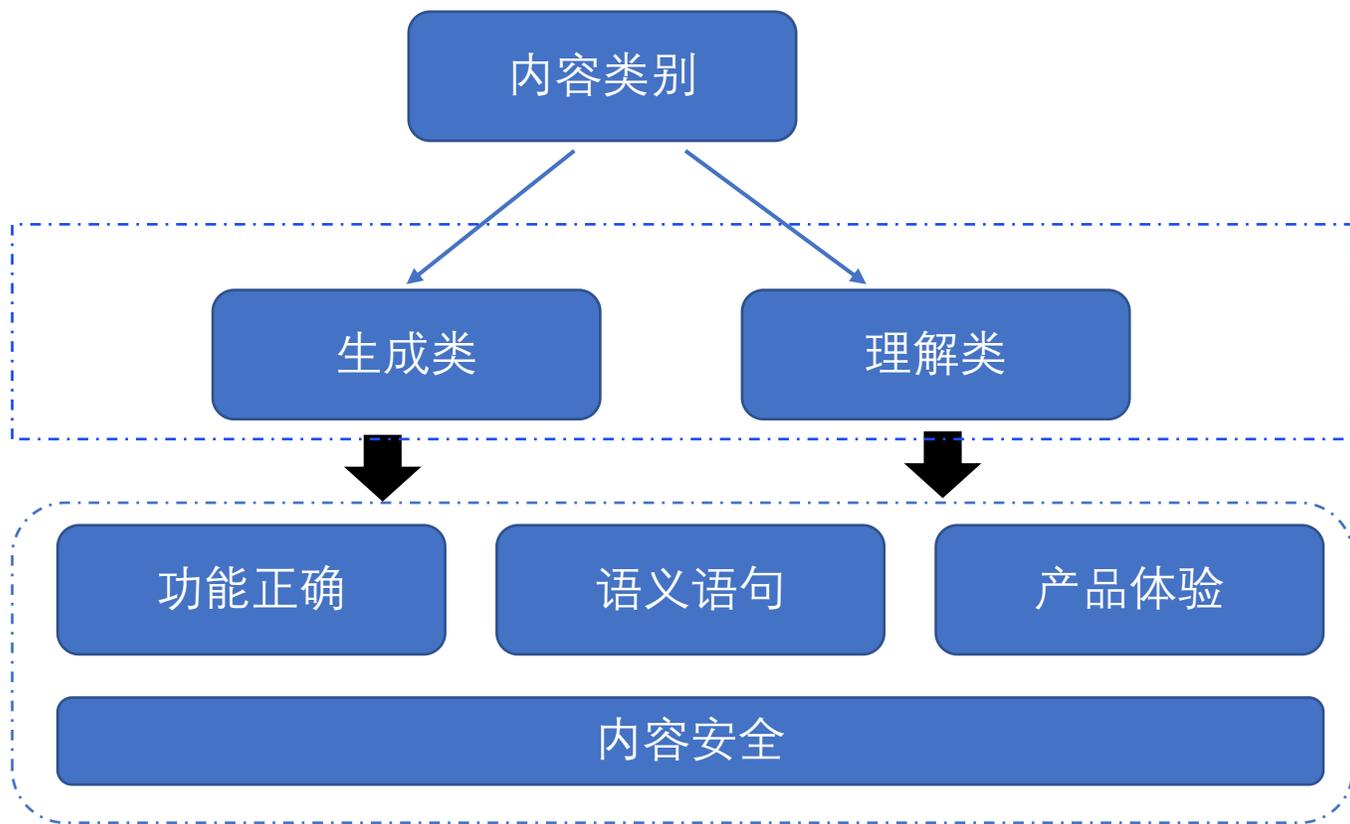
可信： 样本纯净、 标准开放、 客观中立

# PART 02
# 代码大模型的评估内容

# 代码大模型的评估内容

通过评估内容客观衡量代码大模型的各项能力，发现其长处和短处，反馈并协助模型做针对性的改进。

安全

确定

生成

理解

无害

使用

基础能力

编码领域

认知

...

能力发现 ——→ 衡量评估 ←—→ 应用与改进

# 代码大模型的评估内容



内容类别
├── 生成类
└── 理解类

功能正确　　语义语句　　产品体验

内容安全



```
代码补全    添加注释    解释代码

1  package com.alipay.codefuse;
2
   采纳 Tab      下一个 Option ]    上一个 Option [
3  public class I
4      /*判断字符串是否是英文*/
5      public static boolean isEnglish(String str) {
6          char[] chars = str.toCharArray();
7          for (char c : chars) {
8              if (!(c >= 'a' && c <= 'z') && !(c >= 'A' && c <= 'Z')) {
9                  return false;
10             }
11         }
12         return true;
13     }
14  }
15
```

```
代码补全    添加注释    解释代码

1  package com.alipay.codefuse;
2  public class IsEnglish{
3      // 判断字符串是否全部由英文字母组成的方法
4      public static boolean isEnglish(String str) {
5          // 将字符串转化为字符数组
6          char[] chars = str.toCharArray();
7          // 遍历字符数组
8          for (char c : chars) {
9              // 如果某个字符不是英文字符，则说明不是全部由英文字母组成的字符串
10             if (!Character.isLetter(c)) {
11                 return false;
12             }
13         }
14         // 如果遍历完整个字符数组都是英文字母，则说明是全部由英文字母组成的字符串
15         return true;
16     }
17  }
18
```

PART 03
代码大模型的评估方法

# ▶ 代码大模型的评估方法

按照是否由人来评估模型性能，评测方法分为主观和客观

**客观评测**

基于评测基准对生成内容进行各维度量化评估

**主观评测**

组织多位专业人员通过人模交互观察模型表现
并根据基础标准、专家知识和经验综合评估

# 代码大模型的评估方法

按照prompts设置方法评测又可分为：零样本（zero-shot）、小样本（few-shot）、零样本思维链（zero-shot-cot）、小样本思维链（few-shot-cot）

代码生成能力目前大部分采用的策略是：零样本（zero-shot）

# PART 04
# 代码大模型评估基准

# 评估基准

模型评估基准是优化模型，了解差距，衡量不同架构模型的同类场景性能的最有效的工具。

| 代码能力评测评测基准（截止10月份-不完全统计） | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| 名称 | 时间 | 语言 | 评估 | 创建方式 | 来源 | 粒度 | 任务数 | 用例数 | 输入方式 |
| Concode | 2018 | Java | BLEU | Automated | GitHub | Function-level | 2000 | - | NL |
| CoNaLA | 2018 | Python | BLEU | Automated | Stack Overflow | Statement-level | 500 | - | NL |
| APPS | 2021 | Python | TestCases | Automated | Contest Sites | Competitive | 5000 | 13.2 | NL + Example Inputs/Outputs |
| HumanEval | 2021 | Python | TestCases | Manual | - | Function-level | 164 | 7.7 | NL + Function Signature + Example Inputs/Outputs |
| MBPP | 2021 | Python | TestCases | Manual | - | Function-level | 974 | 3 | NL |
| math-qa | 2021 | Python | Acc | Manual | Math Study Sites | Statement-level | 2985 | - | NL |
| Multi-HumanEval | 2022 | 多语言 | TestCases | Manual | - | Function-level | 164 | 7.7 | NL + Function Signature + Example Inputs/Outputs |
| MBXP | 2022 | 多语言 | TestCases | Manual | - | Function-level | 974 | 3 | NL |
| multi-math-qa | 2022 | 多语言 | Acc | Manual | Math Study Sites | Statement-level | 2985 | - | NL |
| CodeContests | 2022 | Python、C++ | TestCases | Automated | Contest Sites | Competitive | 165 | 203.7 | NL + Example Inputs/Outputs |
| DS-1000 | 2022 | Python | TestCases+Surface Form Constraints | Automated | Stack Overflow | Statement-level | 1000 | 1.6 | NL + Function Signature + Example Inputs/Outputs |
| HumanEval+ | 2023 | Python | TestCases | Manual | - | Function-level | 164 | 7.7 | NL + Function Signature |
| CoderEval | 2023 | Python、Java | Compilation | Automated | Github | Function-level | 230 | - | |
| ClassEval | 2023 | Python | TestCases | Manual | - | Class-level | 100 | 30 | Class Skeleton |
| CodeFuseEval | 2023 | 多语言 | pass@k/ES/BLEU... | Manual | - | Muti-Tasks | 6000+ | - | NL/Code/Example |
| CCEval | 2023 | 多语言 | EM/ES/F1 | Manual | - | Class-level | 100 | 30 | Class Skeleton |

采纳原则：  有效未被污染的      多样多维：多语言，多任务，多维衡量

# 评估基准的演进

MBPP (pass@k)

HumanEval(pass@k) — 单语言 - Python

APPS (pass@any)　　Codex-12B/CODE-T/GPT-4等（50+）

| 2018 | 2021 | 2022 | 2023 |

CoNaLA (BLEU) 单语言- python

TranX/ Reranker / PanGu-Coder-FT-I
（10+）

Concode(B LEU)　单语言 -Java

CodeT5/Redcoder-ext （2+）

10+语言
Python/Java/
Go/Ruby…

Multi-HumanEval (pass@k)

CodeContests (Test Set))

DS1000 (TestCases-Score))
......

HumanEval-X (pass@k)

CoderEval(pass@k, acc@k)

ClassEval (pass@k,)

CodeFuseEval (pass@k,ES,BLUE...)

CCEval (ES/EM?FQ,)

单语言 静态指标 单轮 ⟶ 多语言 动态指标 单轮 ⟶ 多语言 动静指标 单轮

# ▶ 评估基准（自建）

自建是开源的延伸:

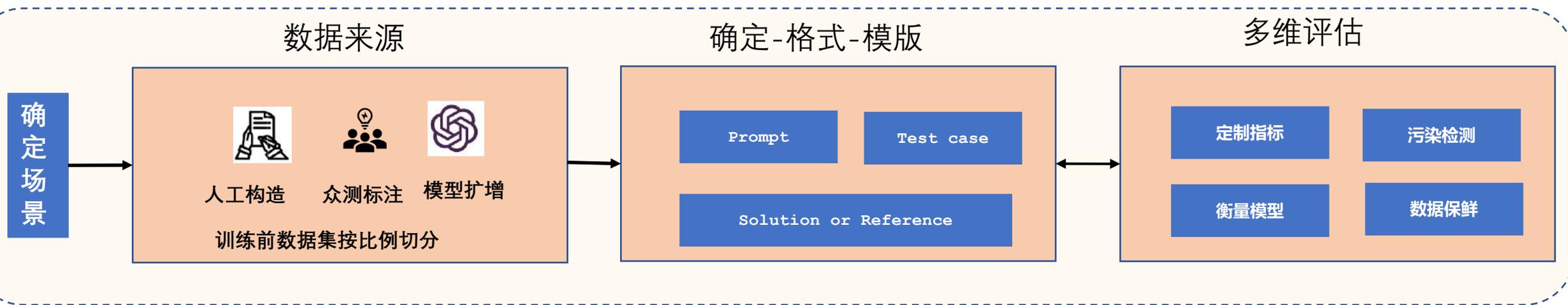1. 更贴合企业特定场景如衡量模型对企业项目代码的生成和理解能力

2. 补充开源评测集，全面有梯度的衡量模型性能

# ▶ 评测样本prompt&格式

```
from typing import List
```

代码补全-prompt

```
def has_close_elements(numbers: List[float], threshold: float) -> bool:
    """ Check if in given list of numbers, are any two numbers closer to each othe
    given threshold.
    >>> has_close_elements([1.0, 2.0, 3.0], 0.5)
    False
    >>> has_close_elements([1.0, 2.8, 3.0, 4.0, 5.0, 2.0], 0.3)
    True
    """
```

```
{
    "task_id": "Python/0",
    "prompt": "from typing import List\n\n\ndef has_close_elements(numbers
    "canonical_solution": "    for idx, elem in enumerate(numbers):\n
    "test": "\n\nMETADATA = {\n    'author': 'jt',\n    'dataset': 'test'\
    "text": "    Check if in given list of numbers, are any two numbers cl
    "prompt_text": "Check if in given list of numbers, are any two numbers
    "prompt_explain": "Check if in given list of numbers, are any two numb
    "func_title": "def has_close_elements(numbers: List[float], threshold:
    "prompt_text_chinese": "检查在给定的数字列表中，是否有任何两个数字比给定的阈值
}
```

```
Write a python function to remove first and last occurrence of a given character
from the string.
```

自然语言到代码

```
#include <bits/stdc++.h>
using namespace std;
```

代码翻译：【c++ 】-> target language

```
/**
 * Write a function to check if the given tuple list has all k elements.
 * > checkKElements(vector<vector<int>>{{4, 4}, {4, 4, 4}, {4, 4}, {4, 4, 4, 4},
{4}}, 4)
 * true
 * > checkKElements(vector<vector<int>>{{7, 7, 7}, {7, 7}}, 7)
 * true
 * > checkKElements(vector<vector<int>>{{9, 9}, {9, 9, 9, 9}}, 7)
 * false
 */
bool checkKElements(vector<vector<int>> testList, int k) {
    for (vector<int> i: testList)
        for (int j: i)
            if (j != k)
                return false;
    return true;
}
```

# 多任务评估基准



第一阶段
（10-11月）

CodeFuseEval

生成代码
- 代码补全
- 自然语言到代码
- 代码翻译
- 代码优化
- 用例生成
- 缺陷修复

代码应用
- 数据科学
- 文件处理
- 数学运算
- 工具调用
- 存储操作
- API接入

代码安全
- 代码漏洞
- 内容安全
- 业务合规

领域知识
- 计算机知识
- 其他行业知识

理解代码
- 代码意图
- 生成代码注释
- 解释代码
- 缺陷检测
- 代码优化建议
- 白盒测试建议

第二阶段(12月)



代码补全-多编程语言
代码补全-中文
自然语言到代码
代码补全-英文
测试生成
代码翻译
代码生成-数据科学

- CodeFuse-CodeLlama-34B
- GPT-3.5
- GPT-4

开源地址：https://github.com/codefuse-ai/codefuse-evaluation

# ▶ 多类型评测指标

生成结果要求：有用 真实 无害

## 评估指标



- manual_eval
- pass@K
- code_match
- comment_match
- bleu
- rouge
- sql_va
- accuracy
- F1
- sql_ex
- bleurt
- clean_proportion

| 指标 | 类型 | 描述 |
|------|------|------|
| manual_eval | 自定义指标 | 人工评测，该指标将模型推理结果，通过人工打分来判断评测结果，详见主观评分表 |
| comment_match | 自定义指标 | comment_fix场景测试 |
| code_match | 自定义指标 | bug_fix场景测试 |
| clean_proportion | 自定义指标 | 评测集纯净度，检测评测数据是否被污染 |
| markdown_pct | 自定义指标 | markdown_caculate |
| bleurt | 基础指标 | BLEURT（Bilingual Evaluation Understudy for Machine Translation）不仅关注词汇层面的匹配，而且能够捕捉生成文本与参考文本之间的语义信息 |
| bleu | 基础指标 | 语句相似性评估指标，通过比较源语言句子和参考语言句子之间的 n-gram 匹配来计算。 |
| rouge | 基础指标 | rouge 是一个衡量自动文本摘要和句子摘要性能的测试指标 |
| sql_ex | 基础指标 | sql 执行正确性指标 |
| sql_va | 基础指标 | sql 逻辑语法正确性指标 |
| accuracy | 基础指标 | 精确度（Accuracy）：是所有分类正确的样本数目占样本总数的比例 |
| f1 | 基础指标 | 精确率和召回率的调和平均数 |
| EvalGPT_GoodCaseRate | 自定义指标 | 大模型评测结果后输出当前数据集的点赞率 |
| pass@K | 基础指标 | 用于评估code任务中代码功能的正确性。 |

PART 05
代码大模型多任务评估

# ▶ 代码大模型多任务评估

多任务评估让我们全面了解模型的专业能力和基础性能、检验模型的泛化能力、并辅助模型优化。

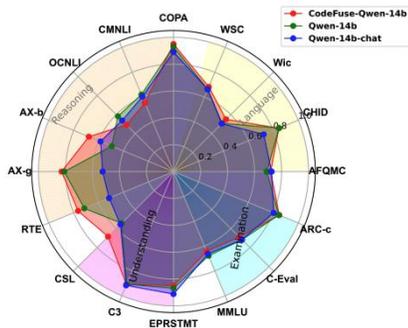| Model | Size | Python | Java | C++ | JavaScript | Golang | Avgerage |
|---|---|---|---|---|---|---|---|
| QWen-base | 14B | 32.3%* | 35.37% | 30.49% | 32.93% | 21.34% | 30.49% |
| CodeFuse-QWen-MFT | 14B | 48.78% | 41.46% | 38.41% | 46.34% | 26.83% | 40.36% |
| Llama-base | 65B | 23.7%* | 29.26% | 20.73% | 23.78% | 18.9% | 23.27% |
| CodeFuse-Llama-MFT | 65B | 34.76% | 37.2% | 29.88% | 32.93% | 23.78% | 31.71% |
| Llama2-base | 70B | 29.9%* | 39.02% | 31.10% | 35.98% | 23.78% | 31.96% |
| CodeFuse-Llama2-MFT | 70B | 40.85% | 35.98% | 32.32% | 38.41% | 27.44% | 35.00% |
| StarCoder-base | 15B | 33.6%* | 34.15% | 25.61% | 22.56% | 22.56% | 29.48% |
| CodeFuse-StarCoder-MFT | 15B | 54.9% | 47.56 | 46.34% | 48.17% | 37.20% | 46.83% |
| CodeGeex2-base | 6B | 35.9%* | 30.8%* | 29.3%* | 32.2%* | 22.5%* | 30.14% |
| CodeFuse-CodeGeex2-MFT | 6B | 45.12% | 45.73% | 37.2% | 37.2% | 28.05% | 38.66% |
| CodeLlama-Python-base | 13B | 43.3%* | 41.46% | 34.76% | 38.41% | 29.27% | 37.44% |
| CodeFuse-CodeLlama-Python-MFT | 13B | 60.37% | 57.32% | 46.34% | 54.27% | 45.12% | 52.68% |
| CodeLlama-34B-Python-base | 34B | 53.7%* | 45.73% | 42.68% | 45.73% | 31.71% | 43.91% |
| CodeFuse-CodeLLama-Python-MFT | 34B | **74.4%** | **61.6%** | **54.3%** | **61.0%** | **50.6%** | **60.38%** |

Table 5: Performance(pass@1) comparison of CODEFUSE with previous models on code translation using greedy decoding

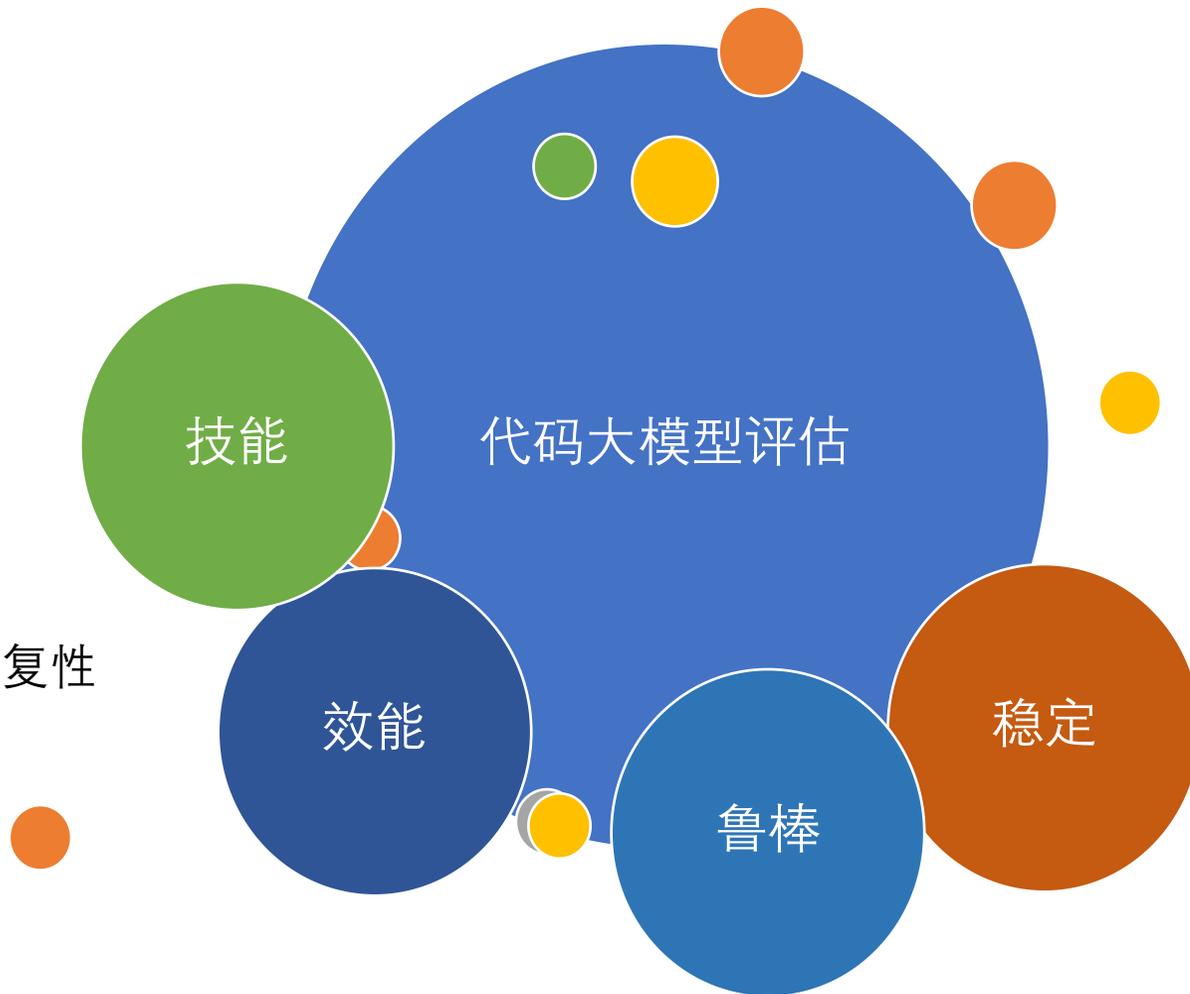| Models | Java to Py | C++ to Py | C++ to Java | Java to C++ | Py to Java | Py to C++ | Average |
|---|---|---|---|---|---|---|---|
| CODEFUSE-13B-Base | 53.66% | 55.49% | 41.46% | 37.80% | 48.10% | 50.00% | 47.75% |
| CODEFUSE-13B-SFT | 66.46% | 59.15% | 54.27% | 47.56% | **56.31%** | **55.40%** | **56.53%** |
| CODEGEN-multi-16B | 52.73% | 33.83% | 43.20% | 41.42% | 29.27% | 35.94% | 39.40% |
| CODEGEEX-13B | 43.41% | 27.18% | 22.56% | 39.33% | 25.84% | 26.54% | 30.81% |
| CODEGEEX-13B-FT | **75.03%** | **62.79%** | **71.68%** | **49.67%** | 41.98% | 34.16% | 55.89% |

代码能力



NLP能力

PART 06
展望

# ▶ 展望

评测对象（模型）：
1. 评测任务多样化，贴合企业级编码场景
2. 多维评估：技能、效能、服务（稳定/鲁棒）

评测技术（基准）：

1. 开源开放，快速迭代，匹配模型能力增长速度
2. 框架升级，模型评测，提高主观评测效能及可重复性

代码大模型评估

技能

效能

稳定

鲁棒