

Classes 1

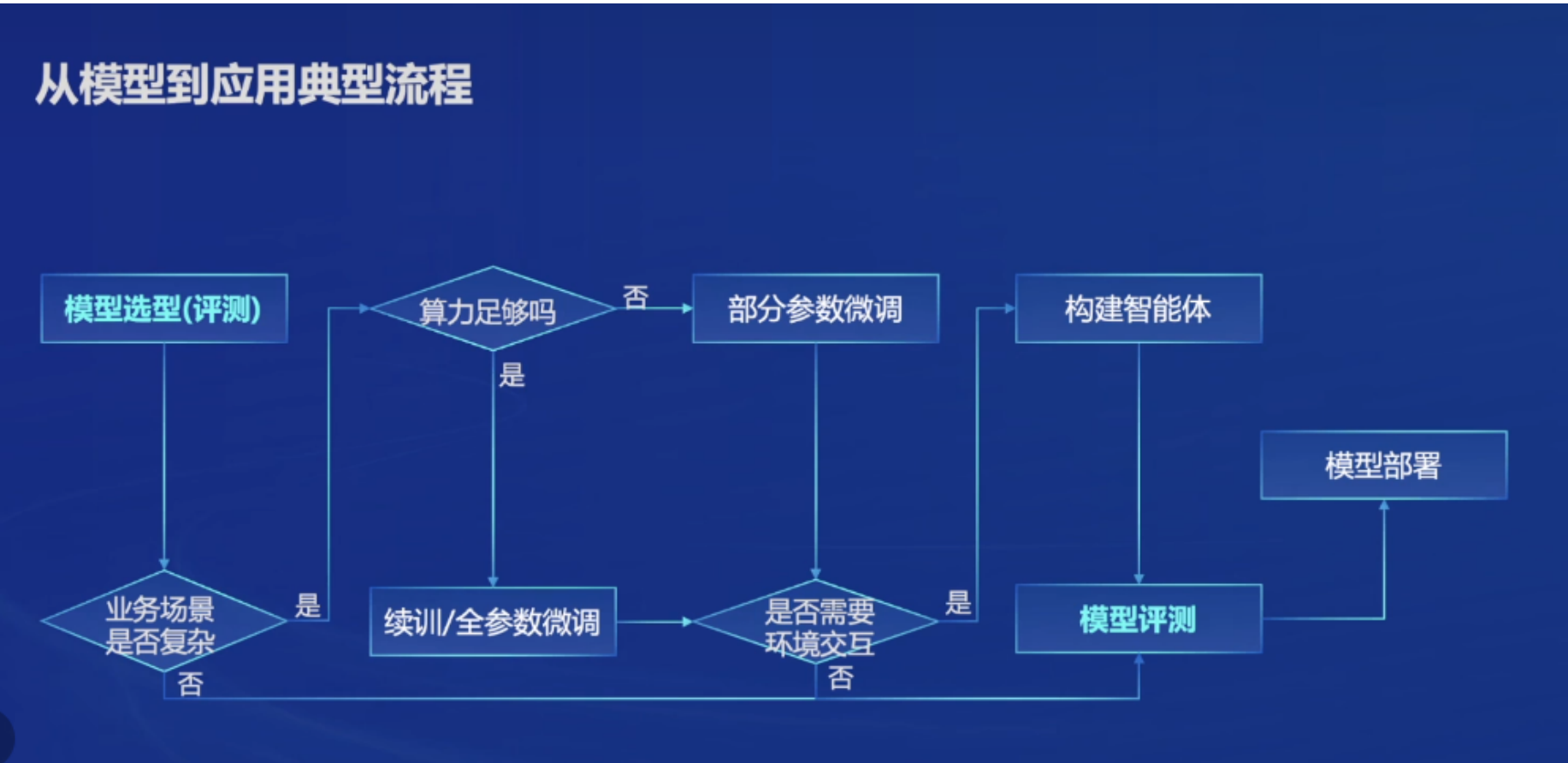
InternLM2的技术报告不仅展示了一个高性能的开源模型，也为如何准备和训练大型语言模型提供了实用的指导。作为一个对自然语言处理领域感兴趣的学习者，我对InternLM2在长文本处理和多语言能力方面的表现感到印象深刻。此外，COOLRLHF策略的引入为解决人类偏好对齐问题提供了新的思路。我期待在实际项目中尝试应用InternLM2，并探索其在不同场景下的潜力

现状：专用模型 → 通用模型

InternLM2 （base， IntenLM2， InternLM2）

- 7B
- 20B

语言建模 → 高质量语料 => 数据清洗要求（价值评估，数据富集，数据补齐）



符合价值观的语料

全链条开源开放体系 | 开放高质量语料数据

数据集获取： OpenDataLab <https://opendatalab.org.cn/>

书生·万卷 1.0

总数据量：2TB
发布日期：2023年8月14日
符合主流中国价值观的中文语料

数据构成

文本数据

5亿个文档
数据量超 1TB

图像-文本数据集

超2,200万个文件
数据量超140GB

视频数据

超2,200万个文件
数据量超 140GB

多模态融合

万卷包含文本、图像和视频等多模态数据，涵盖科技、文学、媒体、教育和法律等多个领域。该数据集对模型的知识内容、逻辑推理和泛化能力的提升有显著效果。

精细化处理

万卷经过语言筛选、文本提取、格式标准化、数据过滤和清洗（基于规则和模型）、多尺度去重和数据质量评估等精细数据处理环节，能够很好地适应后续模型训练的要求。

价值观对齐

在万卷的构建过程中，研究人员注重将数据内容与主流中国价值观进行对齐，并通过算法和人工评估的结合提高语料库的纯净度。

书生·万卷 CC

总数据量：400GB
发布日期：2024年3月6日
安全、信息密度更高的英文语料

三大优势

时间跨度长

横跨了 2013-2023 年互联网公开内容

来源丰富多样

从 90 个 dumps 的 1300 亿份原始数据中“萃取” 1.38% 内容

安全密度高

唯一在毒性、色情和个人隐私都进行了安全加固处理

高质量语料驱动效率性能双提升

万卷CC 作为 InternLM2 预训练语料，在不同数据规模上取得的任务性能分布显示，大幅提升模型训练效率

四重“萃取”高质量数据

通过原创技术，对原始数据进行多阶段处理，得到了高信息密度的万卷CC

数据质量高、模型更可靠

由万卷CC作为训练数据的模型在多项验证中取得了更优效果

Efficiency

Quality

Validation

Eval Datasets	Wanjuan-CC	Refinedweb
gpt-4o-mini	12.54	12.56
gpt-4o-mini-2	11.86	11.96
gpt-4o-mini-cn	8.68	8.81
tiny-stories	5.78	6.15

微调方法

全链条开源开放体系 | 微调

大语言模型的下游应用中，增量续训和有监督微调是经常会用到两种方式。

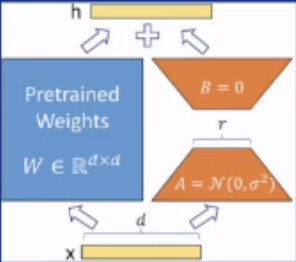
增量续训

使用场景：让基座模型学习到一些新知识，如某个垂类领域知识
训练数据：文章、书籍、代码等

有监督微调

使用场景：让模型学会理解各种指令进行对话，或者注入少量领域知识
训练数据：高质量的对话、问答数据

全量参数微调 部分参数微调



高效微调框架XTuner

全链条开源开放体系 | 微调

高效微调框架 XTuner



适配多种生态

- **多种微调算法**
多种微调策略与算法，覆盖各类 SFT 场景
- **适配多种开源生态**
支持加载 HuggingFace、ModelScope 模型或数据集
- **自动优化加速**
开发者无需关注复杂的显存优化与计算加速细节

适配多种硬件

- **训练方案覆盖 NVIDIA 20 系以上所有显卡**
- **最低只需 8GB 显存即可微调 7B 模型**

评测框架 → OpenCompass (CompassRank,CompassKit,CompassHub

全链条开源开放体系 | 评测



全链条开源开放体系 | 智能体

轻量级智能体框架 Lagent

支持多种类型的智能体能力

ReAct

输入

选择工具

执行工具

结束条件

结束

ReWoo

输入

计划拆分

DAG

计划执行

结束

AutoGPT

输入

选择工具

人工干预

执行工具

结束条件

结束

灵活支持多种大语言模型

 GPT-3.5/4

 InternLM

 Hugging Face Transformers

 Llama

简单易拓展，支持丰富的工具

AI 工具	能力拓展	Rapid API
文生图	搜索	出行 API
文生语音	计算器	财经 API
图片描述	代码解释器	体育资讯 API

32

工具箱

全链条开源开放体系 | 智能体

多模态智能体工具箱 AgentLego

- 丰富的工具集合，尤其是提供了大量视觉、多模态相关领域的前沿算法功能
- 支持多个主流智能体系统，如 LangChain，Transformers Agent，lagent 等
- 灵活的多模态工具调用接口，可以轻松支持各类输入输出格式的工具函数
- 一键式远程工具部署，轻松使用和调试大模型智能体



3

预训练过程

InternLM2的预训练过程包括了文本、代码和长文本数据的准备。模型从处理4k个token的数据开始，逐步扩展到32k个token，有效捕捉长期依赖性，并在“大海捞针”测试中表现优异。

模型结构

InternLM2采用了Transformer架构的改进版本，遵循LLaMA的结构设计原则。模型通过合并某些矩阵和重新配置矩阵布局来提高效率，支持多样化的张量并行变换。

数据准备

报告详细介绍了预训练数据的收集和处理流程，包括文本、代码和长文本数据。特别强调了高质量数据的重要性，并介绍了如何通过一系列过滤和格式化步骤来准备这些数据。

对齐策略

InternLM2采用了监督微调（SFT）和基于人类反馈的条件在线强化学习方法（COOLRLHF）。COOLRLHF通过条件奖励模型来解决偏好冲突，并减少奖励策略滥用问题。

InternLM2的成功开发为开源社区提供了一个强大的LLM工具，适用于多种自然语言处理任务。报告中提供的详细信息和评估结果为后续研究提供了宝贵的参考，有助于推动开源LLM技术的发展。