# 一 书生·浦语大模型全链路开源体系

第一章精彩呈现了书生·浦语大模型的全链路开源体系，并探讨了通用人工智能的最新发展趋势，尤其是行业从专用模型逐渐过渡到通用大模型的重要转变。深入介绍了书生模型在近几个月中的关键升级，特别是其在8K语境支持、多模态输入及不同规模模型的支持方面的进步，这些提升极大增强了该模型在语言建模、对话交互以及智能体框架的处理能力。此外，还讨论了英特尔M2对于模型性能提升的重要影响，使其更加擅长处理复杂场景。

视频还详尽展示了书生·浦语大模型的不同尺寸和类型，旨在适应多样化的使用需求。通过详细的时间线，能够清楚地把握到每个关键环节的详细讲解，包括模型对长上下文的理解、智能体框架以及模型轻量化和多模态智能体工具箱的部署。同时，还特别强调了循环评测策略和全面的客观评测如何突出模型在语言知识推理、数学代码处理等领域的先进性。

最终，全面总结了书生·浦语大模型的开源体系，涵盖了从数据集的准备、模型的预训练、微调、部署，到评测和实际应用等各个环节，彰显了其在推动语言模型技术的进步及实际应用的全面性和深度。

# 大模型成为发展通用人工智能的重要途径

**专用模型：**
针对特定任务，一个模型解决一个问题

**通用大模型：**
一个模型应对多种
任务、多种模态

# 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 | ReWoo | AutoGPT |
|---|---|---|
| 输入 | 输入 | 输入 |
| 选择工具 | 计划拆分 | 选择工具 |
| 执行工具 | DAG | 人工干预 |
| | | 执行工具 |
| 结束条件 | 计划执行 | 结束条件 |
| 结束 | 结束 | 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 |

Since the introduction of ChatGPT and GPT-4 (OpenAI, 2023), Large Language Models (LLMs) have surged in popularity across the academic and industrial spheres. Models trained on billions of tokens have demonstrated profound empathy and problem-solving capabilities, leading to widespread speculation that the era of Artificial General Intelligence (AGI) may soon be upon us. Despite this enthusiasm, the path to developing models with capabilities comparable to those of ChatGPT or GPT-4 remains elusive. The open-source community has been working diligently to bridge the gap between proprietary LLMs and their open-source counterparts. In the past year, several notable open-source LLMs, such as LLaMA (Touvron et al., 2023a;b), Qwen (Bai et al., 2023a), Mistral (Jiang et al., 2023), and Deepseek (Bi et al., 2024), have made significant strides. In this paper, we introduce InternLM2, a new Large Language Model that outperforms the previously mentioned models.

The development of Large Language Models (LLMs) encompasses several main phases: pre-training, Supervised Fine-Tuning (SFT), and Reinforcement Learning from Human Feedback (RLHF) (Ouyang et al., 2022). Pre-training is chiefly based on leveraging a vast corpus of natural text, amassing trillions of tokens. This phase is aimed at equipping LLMs with a broad repository of knowledge and fundamental skills. The quality of data is considered the most crucial factor during pre-training. However, technical reports on LLMs (Touvron et al., 2023a;b; Bai et al., 2023a; Bi et al., 2024) in the past have seldom addressed the processing of pre-training data. InternLM2 extensively details how it prepares text, code, and long-context data for pre-training.

How to effectively extend the context length of LLMs is currently a hot research topic, since many downstream applications, such as Retrieval-Augmented Generation (RAG) (Gao et al., 2023) and agents (Xi et al., 2023), rely on long contexts. InternLM2 first employs Group Query Attention (GQA) to enable a smaller memory footprint when inferring long sequences. In the pre-training phase, we initially train InternLM2 with 4k context texts, then transit the training corpus to high-quality 32k texts for further training. Upon completion, through positional encoding extrapolation (LocalLLaMA, 2023), InternLM2 achieves commendable performance in the "Needle-in-a-Haystack" test within 200k contexts.

Following long-context pre-training, we utilize supervised fine-tuning (SFT) and reinforcement learning from human feedback (RLHF) to ensure the model adheres well to human instructions and aligns with human values. Notably, we also construct corresponding 32k data during these processes to further improve the long-context processing capability of

**Dedup data**. We then apply a composite safety strategy to filter the data, resulting in **Safe data**. We have adopted different quality filtering strategies for data from various sources, ultimately obtaining **High-quality pre-training data**.
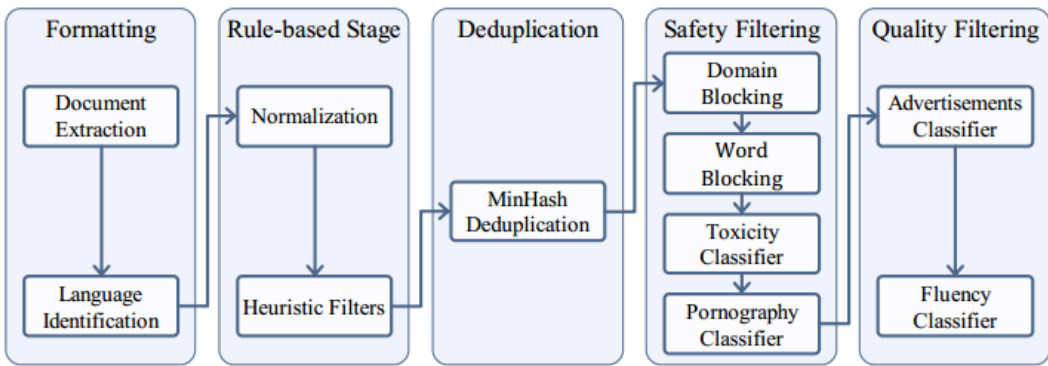


Figure 3: Data Process Pipeline

**Data Formatting**    We will detail the data processing pipeline using web page data as an example. Our web page data mainly comes from Common Crawl[1]. Firstly, we need to decompress the original Warc format files and use Trafilatura (Barbaresi, 2021) for HTML parsing and main text extraction. Then, we use the pycld2[2] library for language detection and classification of the main text. Finally, we assign a unique identifier to the data and store it in jsonl (JSON lines) format and obtained **Format data**.

**Rule-based Stage**    Web page data randomly extracted from the internet often contains a large amount of low-quality data, such as parsing errors, formatting errors, and non-natural language text. A common practice is to design rule-based regularization and filtering methods to modify and filter the data, as seen in Gopher (Rae et al., 2021), C4 (Dodge et al., 2021), and RefinedWeb (Penedo et al., 2023). Based on our observations of the data, we have designed a series of heuristic filtering rules that focus on anomalies in separation and line breaks, frequency of abnormal characters, and distribution of punctuation marks. By applying these filters, we obtained **Clean data**.

**Prompt in CIBench**

**System Prompt:**
You are an assistant who can utilize external tools.
IPythonInterpreter: It can run Python code in a manner as jupyter notebook. The code must be a valid code that contains only python method.
To use a tool, please response with the following format:

**Thought:** Think what you need to solve, do you need to use tools?

**Action:** The tool name, should be one of IPythonInterpreter.

**Action Input:** The input to the tool that you want to use.

The tool will give you response after your response using the following format:

**Response:** the results after call the tool.

Therefore DO NOT generate tool response by yourself.

Also please follow the guidelines:
1. Always use code interpreter to solve the problem.
2. The generated codes should always in a markdown code block format.
3. The generated codes will be executed in an ipython manner and the results will be cached.
4. Your responded code should always be simple and only solves the problem in current step.

For example:

File url: xxxx
### Step 1. Load the dataset from the url into a pandas DataFrame named df.

**Thought:** We should use pandas to solve this step.
**Action:** IPythonInterpreter
**Action Input:**

```
import pandas as pd \\
url = "xxxx" \\
data = pd.read\_csv(url) \\
```

**Response:** The code is succeed without any outputs.
Let us begin from here!

**User Prompt:**

{Question}. Please use {modules} modules.

Figure 19: Prompt used in CIBench

**Prompt in CIBench**

**System Prompt:**
You are an assistant who can utilize external tools.
{tool_description}
To use a tool, please use the following format:

```
\{thought\}Think what you need to solve, do you need to use
    tools? \\
\{action\}the tool name, should be one of [\{action\_names\}]
     \\
\{action\_input\}the input to the action \\
```

The response after utilizing tools should using the following format:

```
\{response\}the results after call the tool. \\
```

If you already know the answer, or you do not need to use tools,
please using the following format to reply:

```
\{thought\}the thought process to get the final answer \\
\{finish\}final answer \\
```

Begin!

**Few-shot Prompt:**

**HUMAN:** Find the coefficient of $x^3$ when $3(x^2 - x^3 + x) + 3(x + 2x^3 - 3x^2 + 3x^5 + x^3) - 5(1 + x - 4x^3 - x^2)$ is simplifie.
**BOT:**
**Tool:**PythonInterpreter
**Tool Input:**

```
from sympy import symbols, simplify

def solution(): \\
\quad x = symbols('x') \\
\quad expr = $3*(x**2 - x**3 + x) + 3*(x + 2*x**3 - 3*x**2 +
    3*x**5 + x**3) - 5*(1 + x - 4*x**3 - x**2)$ \\
\quad simplified\_expr = simplify(expr) \\
\quad x3\_coefficient = simplified\_expr.as\_coefficients\
    _dict()[x**3] \\
\quad result = x3\_coefficient \\
\quad return result \\
```

**SYSTEM:** Response:26

**BOT:** FinalAnswer: The final answer is 26. I hope it is correct.
...

Figure 20: Prompt used in MATH