

AI 驱动 软件研发 全面进入数字化时代

中国·深圳 11.24-25

AI+
software
Development
Digital
summit



大语言模型下的数据及知识管理

彭力 小米AI实验室

科技生态圈峰会 + 深度研习



—1000+ 技术团队的选择



K+全球软件研发行业创新峰会

会议时间: 2024.05.24-25



K+全球软件研发行业创新峰会

会议时间: 2024.09.20-21



AI+ 软件研发数字峰会

会议时间: 2023.11.24-25



AI+ 软件研发数字峰会

会议时间: 2024.07.19-20



AI+ 软件研发数字峰会

会议时间: 2024.11.15-16

▶ 演讲嘉宾



彭力

小米集团-AI实验室-大模型数据团队负责人

2012年至2018年曾就职于百度，于2018年5月加入小米。先后负责知识图谱平台及大模型数据团队的技术体系的构建。目前主要负责小米自研大模型的数据及知识构建及自研模型的落地，并在此期间参与知识图谱国家标准的制定。曾在全球人工智能大会、Data fun talk等多个会议和论坛参加知识图谱相关主题的技术分享

目录

CONTENTS

1. LLM下数据和知识重要性
2. 数据获取中遇到的问题及解决方案
3. 小米业务场景下的大模型的应用
4. 总结与展望

PART 01

知识和数据管理的重要性

▶ 海量数据及知识对大语言模型的影响

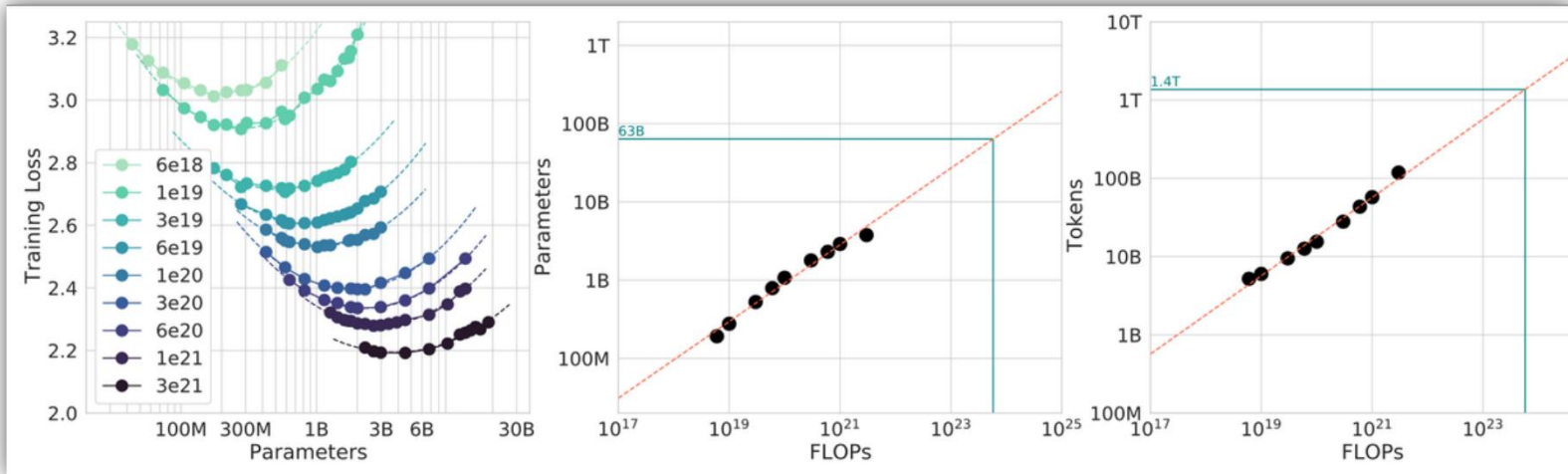
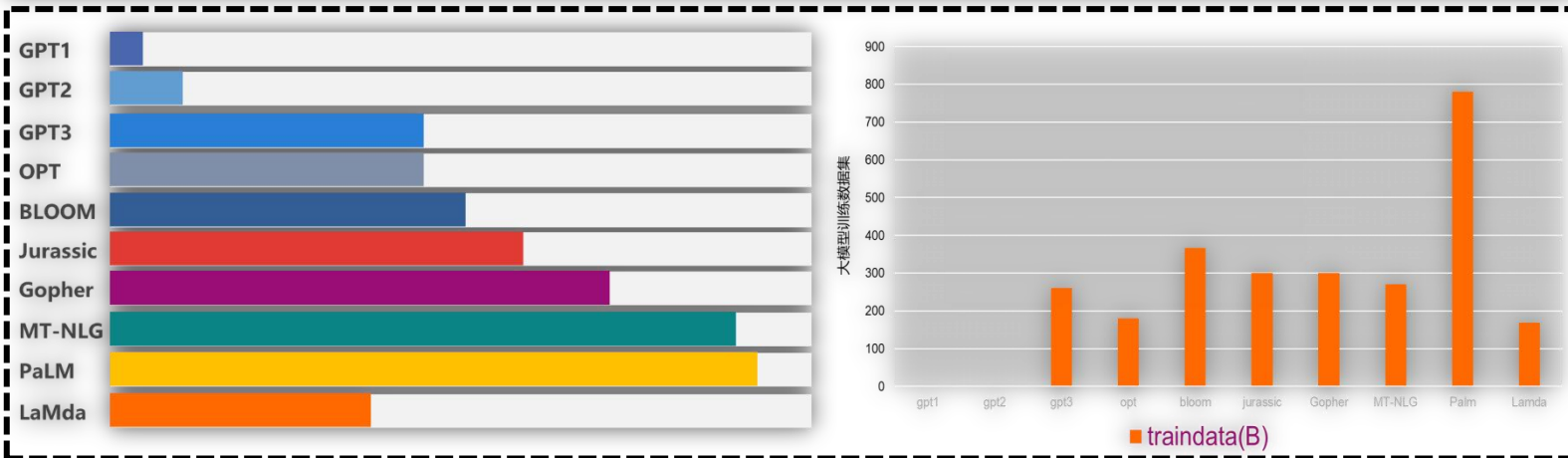


Table 15: Results for several models trained on 1×10^{22} FLOPs. This is with an architecture and data mixture that is different from PaLM 2. When validation loss is roughly equal, downstream performance is approximately equal as well. All evaluations are done in a 1-shot setting. Unless mentioned otherwise, accuracy is used as the evaluation metric. The number of parameters indicates non-embedding parameters.

	3.31B	6.08B	8.95B	14.7B
Triviaqa-Wiki (EM)	43.75	50.57	49.22	48.74
NaturalQuestions (EM)	10.11	10.97	12.58	11.50
WebQuestions (EM)	10.29	11.96	12.70	12.45
LAMBADA	55.46	59.27	60.97	63.05
HellaSwag	69.82	71.84	72.34	72.27
StoryCloze	80.49	80.97	81.88	81.56
Winograd	83.15	85.71	83.52	84.62
Winogrande	69.14	70.56	69.61	70.56
WSC	83.86	84.91	84.56	82.81
SQuAD v2 (EM)	56.19	57.66	55.93	55.74
RACE-H	40.68	43.80	43.51	42.65
RACE-M	56.96	59.68	58.84	58.84
TyDiQA-En (EM)	43.64	41.36	43.86	43.64
PiQA	77.42	77.86	78.73	78.02
ARC-C	40.61	43.60	42.66	43.69
ARC-E	71.55	73.70	74.28	71.00
OpenBookQA	50.60	51.60	50.60	54.00
BoolQ	68.59	71.25	71.31	71.59
CB	64.29	41.07	42.86	51.79
COPA	85.00	87.00	84.00	90.00
MultiRC	63.37	60.58	60.68	58.91
ReCoRD	88.49	89.85	89.45	89.74
RTE	63.54	55.60	63.90	63.90
WiC	48.75	47.49	48.28	47.81
ANLI-R1	33.10	34.70	31.70	35.10
ANLI-R2	30.70	33.30	32.80	31.50
ANLI-R3	32.09	35.50	34.00	34.25
Average	57.30	57.61	57.68	58.26

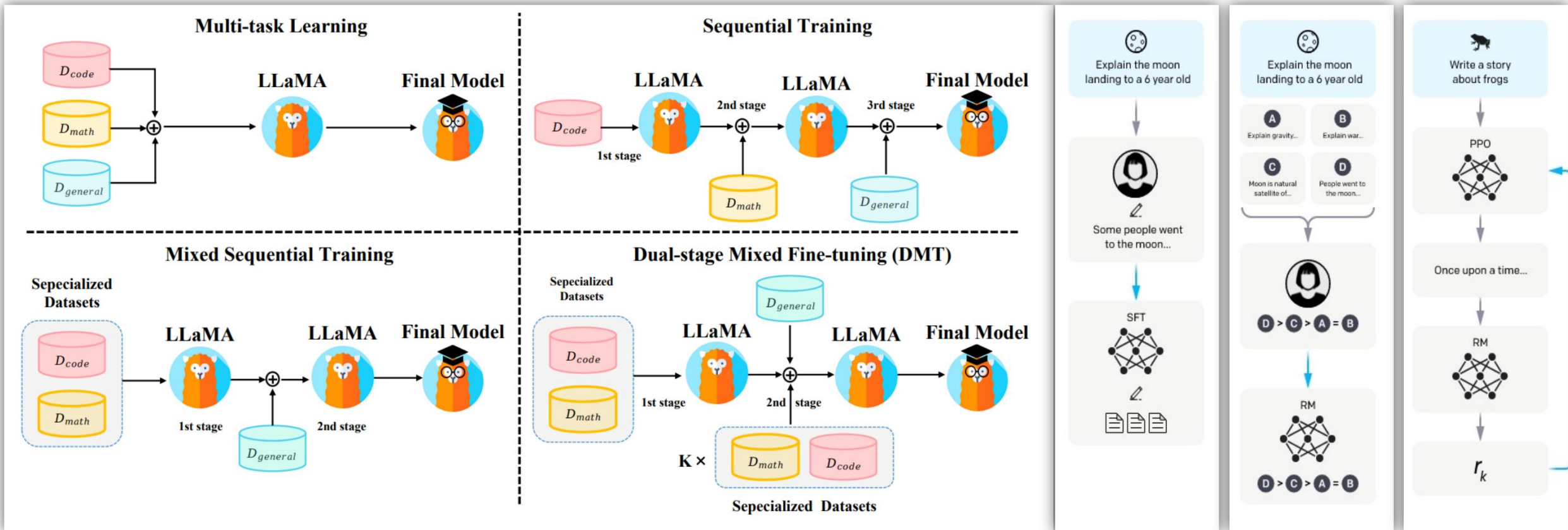


模型参数越大需更越多的知识

大语言模型基座训练需要更多更丰富的知识

▶ 海量数据及知识对大语言模型的影响

除基座训练以外全阶段的性能依赖高质量的知识



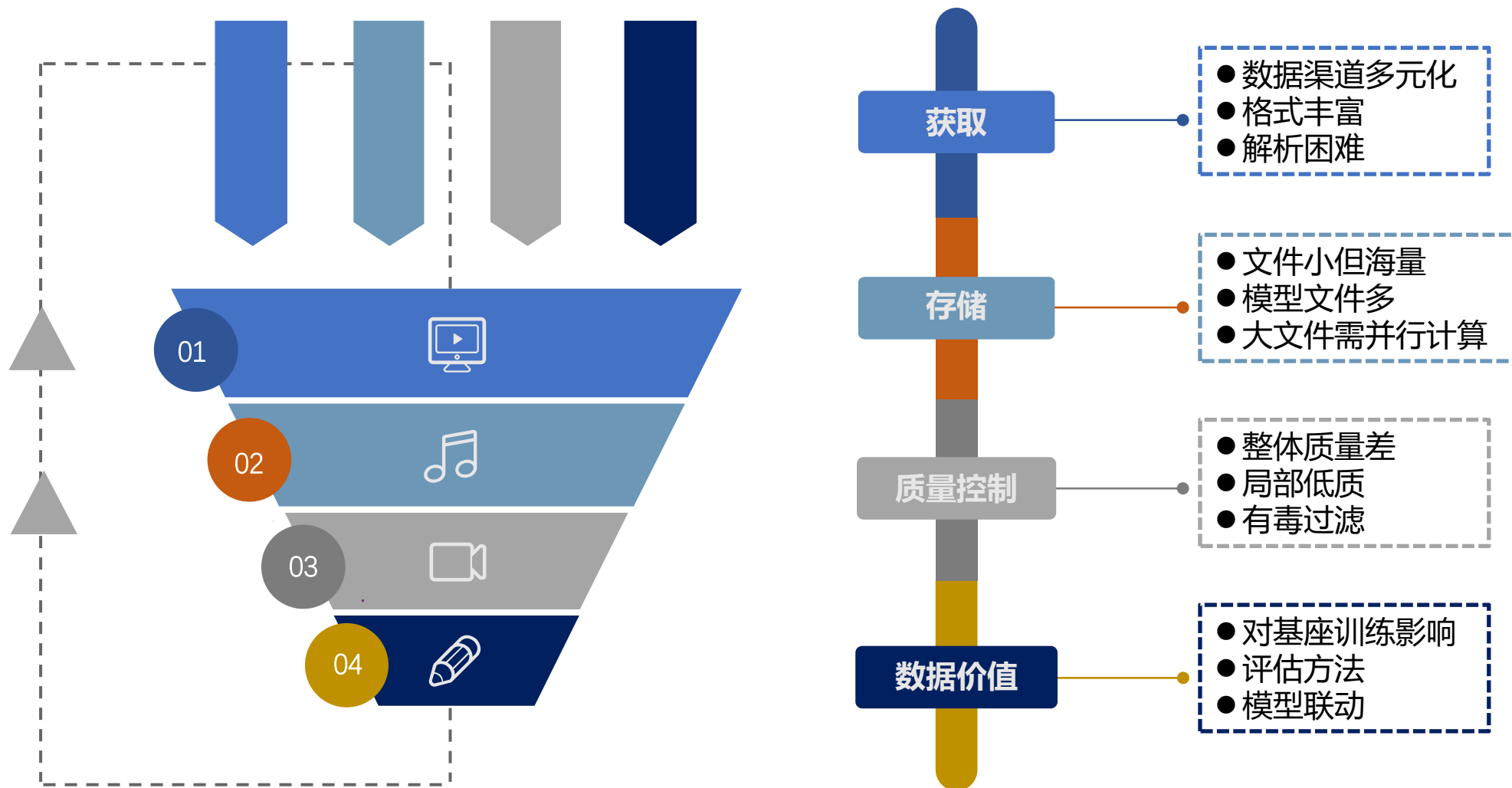
How Abilities in Large Language Models are Affected by Supervised Fine-tuning Data Composition

Introducing ChatGPT

PART 02

大模型下数据加工及管理的解决方案

▶ 大语言模型中数据及知识管理问题



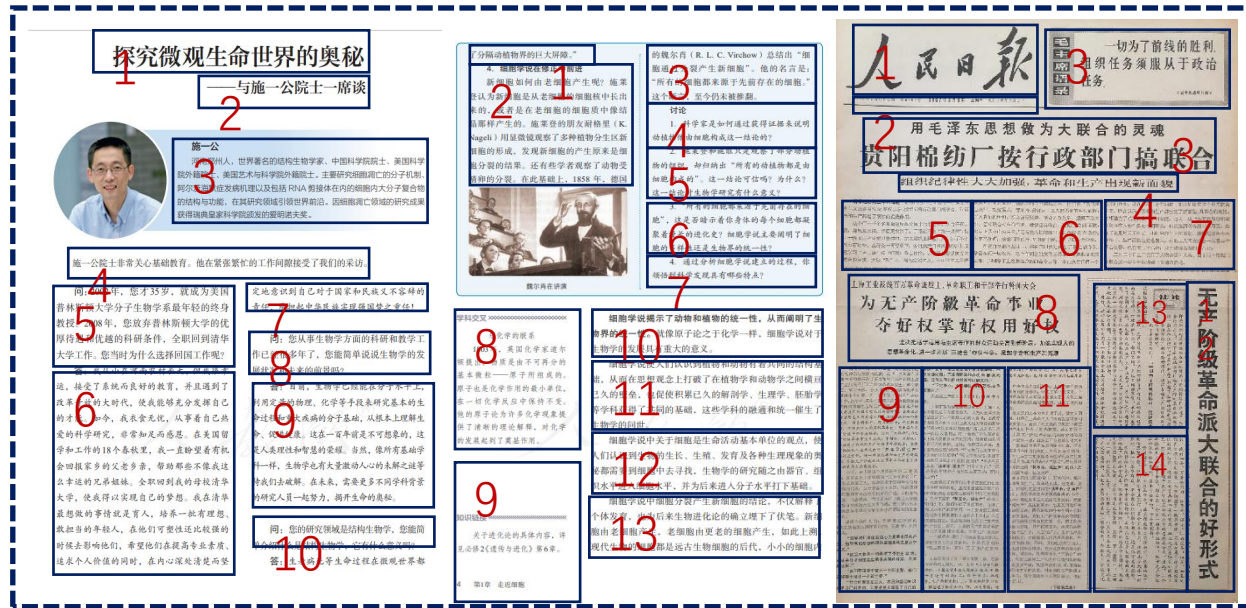
小米大语言模型中数据及知识挖掘

基于文档分析的不同格式的数据和知识提取依然存在挑战

布局识别

阅读顺序

格式解析
表格/公式/化学式/.....



● 布局复杂且包含丰富的元素

● 多栏之间转换影响阅读的顺序

● 多样式的公式及表格等需统一

AI驱动软件研发全面进入数字化时代

AI+ 软件研发数字峰会
AI+ software Development Digital summit

小米大语言模型中数据及知识挖掘

计算得到主要影响因素，包括驱动能力偏低、活门输出特性分散度大及驱动负载偏大，进一步研究其产生原因，认为是由发动机个体差异或热状态差异导致的，因此在故障现象上表现出一定的不确定性。进一步对上述3方面主要因素进行叠加分析，复现了故障现象。

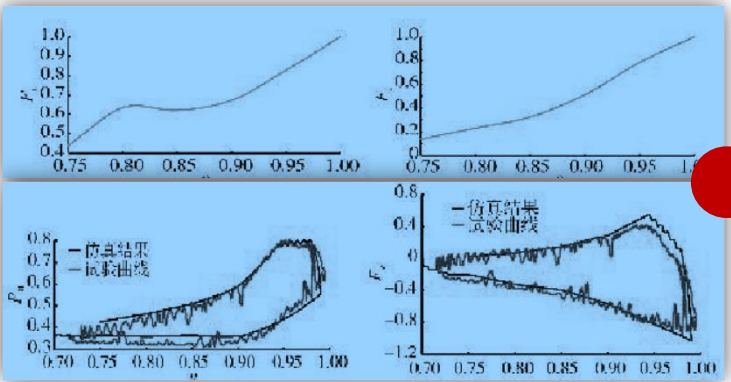


图8 作动筒B腔压力仿真 图9 驱动力仿真及试验结果(无量纲)

3.2 故障模式分析及验证

分析活塞工作状态,可得液压驱动力与腔压及活塞结构尺寸关系

式中: F_a, F_b 分别为转速增加,降低过器油压驱动力; P_a, P_b 分别为作动筒2腔腔压; A_a, A_b 分别为对应侧活塞作用面积; A_s 为活塞杆截面积

其中腔压可参考实际测量数据,活塞结构尺寸按实际情况给出,由于驱动负载 F_a, F_b 与 F_q, F_p 分别为作用力与反作用力,因此,联立式(1)-(4)可初步估算出气动力及摩擦阻力曲线,如图6所示。

(a) 气动力 (b) 摩擦阻力
图6 气动力及摩擦阻力曲线(无量纲)

对搭建的系统仿真模型的有效性进行了验证,将仿真计算结果与整机试验结果进行对比,可以看到腔压模拟(如图7.8所示)和驱动力模拟(如图9所示)计算结果与试验实测数据有较好的一致性,说明该模型可以较好地反映VSV系统各模块间的协调作用,且计算结果误差较小,可以作为进一步故障因素分析工作的仿真计算基础。

图7 作动筒A腔压力仿真及试验结果(无量纲)

3. 仿真计算及结果分析

以搭建的VSV闭环调节系统为研究平台,运用FMEA方法,注入可能的影响因素进行仿真计算,分析可能的故障因素,确定问题故障模式。并通过上述分析,对VSV闭环调节系统的设计流程进行优化改进。

3.1 故障影响因素仿真分析

根据系统工作原理、内部结构组成及试验测量情况,识别出的可能影响因素主要包括油路压力异常、控制因素、输入输出关系异常等反馈机构因素,以及摩擦阻力、驱动负载异常等操纵执行机构因素。针对可能的影响因素,以搭建的系统仿真模型为平台,逐一进行FMEA分析,计算结果见表1。

计算得到主要影响因素,包括驱动能力偏低、活门输出特性分散度大及驱动负载偏大,进一步研究其产生原因,认为是由发动机个体差异或热状态差异导致的,因此在故障现象上表现出一定的不确定性。进一步对上述3方面主要因素进行叠加分析,复现了故障现象。

根据仿真分析结果,结合实际试验情况得到如下故障模式:首先,在实际试车过程中,在相同转速下,控制器主泵压力与回油压力存在一定的分散度,即控制器自身驱动能力有所差别,导致在某一时刻因驱

$$F_{a1} = F_a + F_f \quad (1)$$

$$F_{a2} = F_f - F_a \quad (2)$$

$$F_{q1} = P_a \cdot A_a - P_b \cdot (A_B - A_g) \quad (3)$$

$$F_{q2} = P_b \cdot (A_B - A_g) - P_a \cdot A_a \quad (4)$$

表1 故障因素影响计算

序号	影响因素	VSV 滞后量 / (°)
1	泵后压力过低	仿真排除
2	活门泄漏	仿真排除
3	摩擦阻力突变	仿真排除
4	控制器内部机械刚度过低	仿真排除
5	反馈机构位移输入输出偏差	仿真排除
6	作动筒内泄漏	仿真排除
7	驱动能力偏低	0.63
8	活门输出特性分散度大	0.89
9	驱动负载偏大	0.98
10	活门输出特性分散度大 + 驱动负载偏大	2.70
11	驱动负载偏大 + 驱动力偏低	2.03
12	活门输出特性分散度大 + 驱动力偏低	1.79
13	活门输出特性分散度大 + 驱动负载偏大 + 驱动力偏低	3.35

文档元素类别

1 图片

2 图片caption

3 文本段落

4 标题

5 页眉

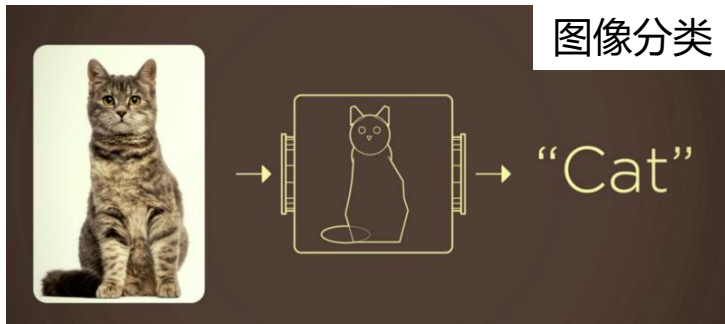
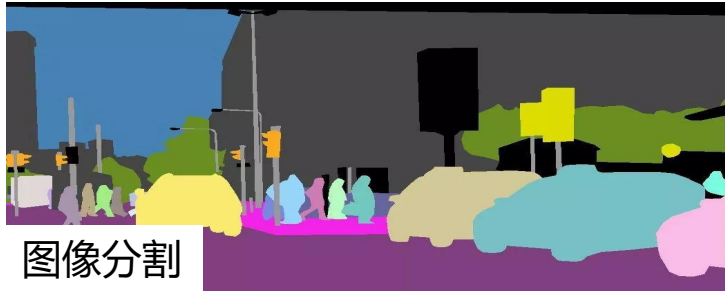
6 公式

7 表格caption

8 表格

小米大语言模型中数据及知识挖掘

第二阶段:机器学习阶段 2000年



第一阶段: 启发式阶段 1900年(规则) Docstrum

第三阶段:深度学习阶段 2014年

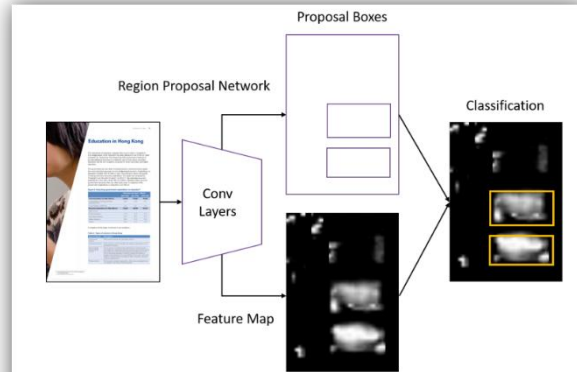


Figure 2: Document layout analysis with Faster R-CNN

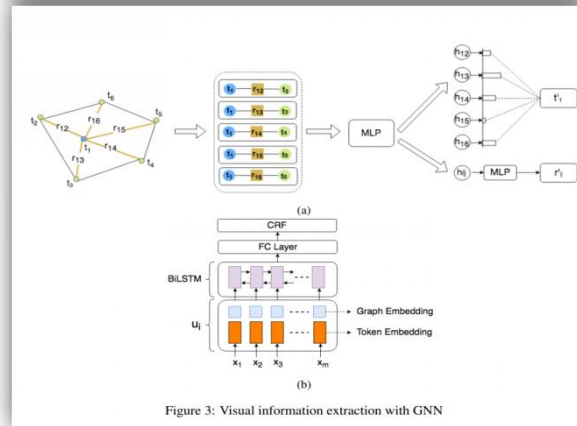


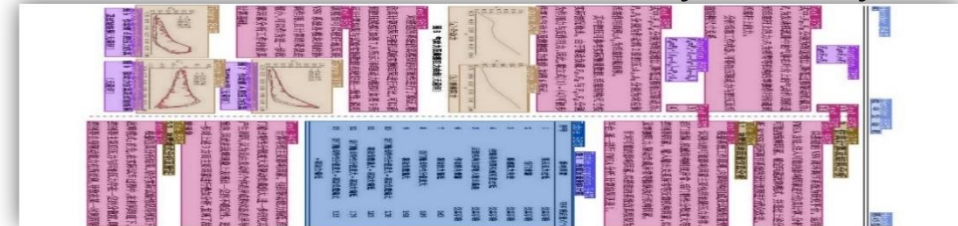
Figure 3: Visual information extraction with GNN

DOCUMENT AI: BENCHMARKS, MODELS AND APPLICATIONS
<https://arxiv.org/pdf/2111.08609.pdf>

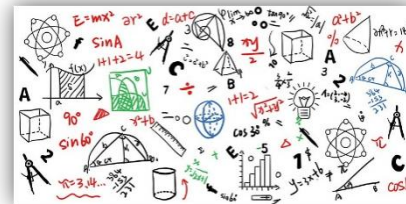
● 文档图像分类



● 文档布局分析 (Document layout analysis)



● 公式、表格检测结构识别



Sample Group	Same Year 1 Head Start Participation	No Year 1 Head Start Participation	Total
All Randomly Assigned (N=247)	87.5%	14.9%	2024
3-Year-Old Cohort	17.2%	82.7%	2024
4-Year-Old Cohort	76.0%	24.0%	2024
Control Group	11.9%	88.1%	2024

(a) FinTabNet-1

Sample Group	Same Year 1 Head Start Participation	No Year 1 Head Start Participation	Total
All Randomly Assigned (N=47)	87.5%	14.9%	2024
3-Year-Old Cohort	17.2%	82.7%	2024
4-Year-Old Cohort	76.0%	24.0%	2024
Control Group	11.9%	88.1%	2024

(b) PubTables-1M

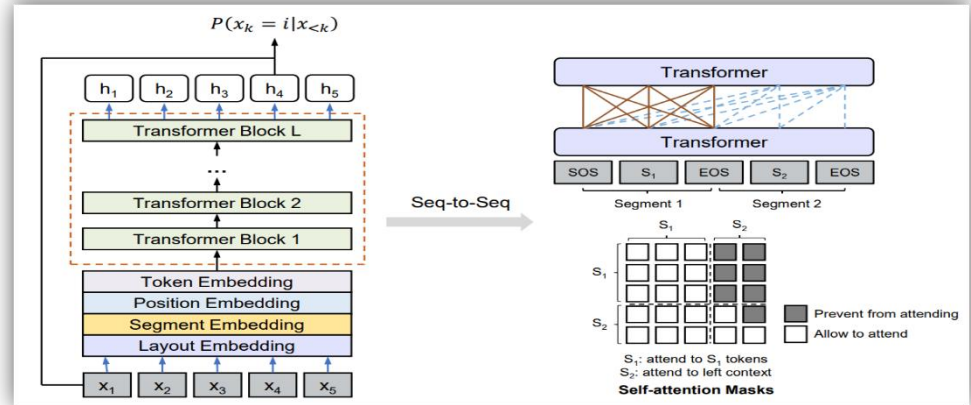
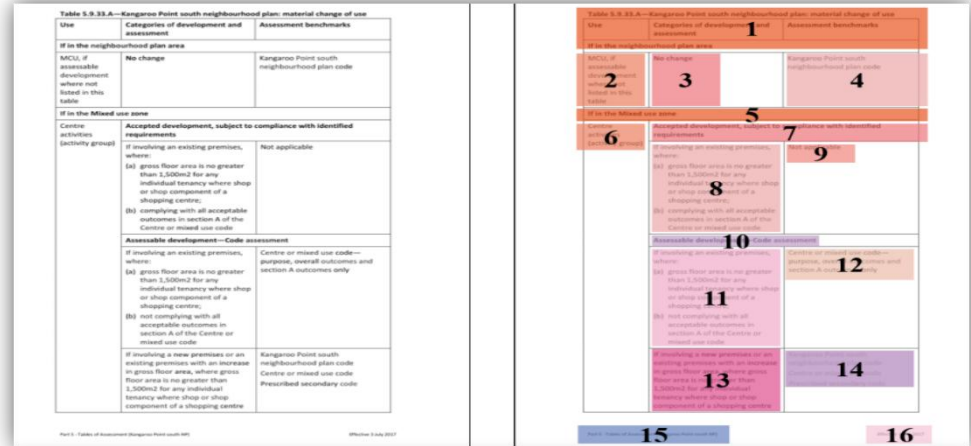
小米大语言模型中数据及知识挖掘

相关布局模型

模型名	模型规模	最佳表现情况	技术方法	发布时间
DiT	-	ON RVL-CDIP 🏆 2022 SOTA! Accuracy 92.69 Param304M	-	2022-03
RectiNet-v2	-	ON DocUNet 🏆 2021 SOTA! SSIM 0.507434	Concatenated Skip Connection	2021-02
Faster-RCNN	-	ON Document Layout Recognition Challenge mini-dev 🏆 2015 SOTA! Overall 0.95 Title0.88 Text0.95 Figure0.98 Table0.98 List0.98	-	2020-10
VisualWordGrid	-	ON RVL-CDIP 🏆 2020 SOTA! FAR 28.7 WAR18.7	-	2020-10
RectiNet	-	ON DocUNet 🏆 2020 SOTA! MS-SSIM 0.415	Convolution	2020-07
DewarpNet	-	ON DocUNet 🏆 2019 SOTA! LD 8.98	-	2019-10
DocUNet	-	ON DocUNet 🏆 2018 SOTA! LD 14.08	Concatenated Skip Connection	2018-06
Faster-RCNN	-	ON Document Layout Recognition Challenge test 🏆 2015 SOTA! Overall 0.91 Title0.82 Text0.92 Figure0.95 Table0.95 List0.89	-	2018-03

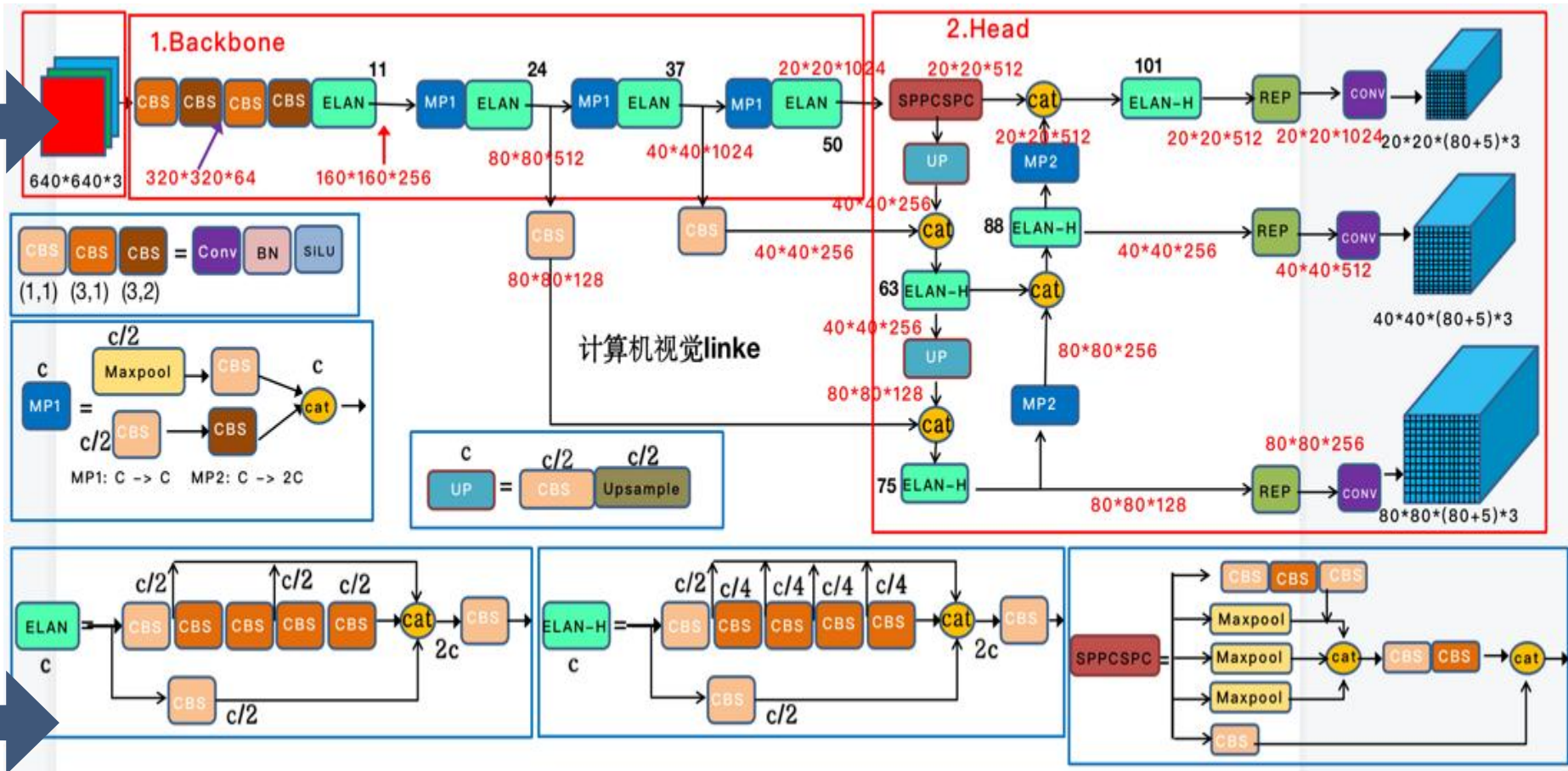
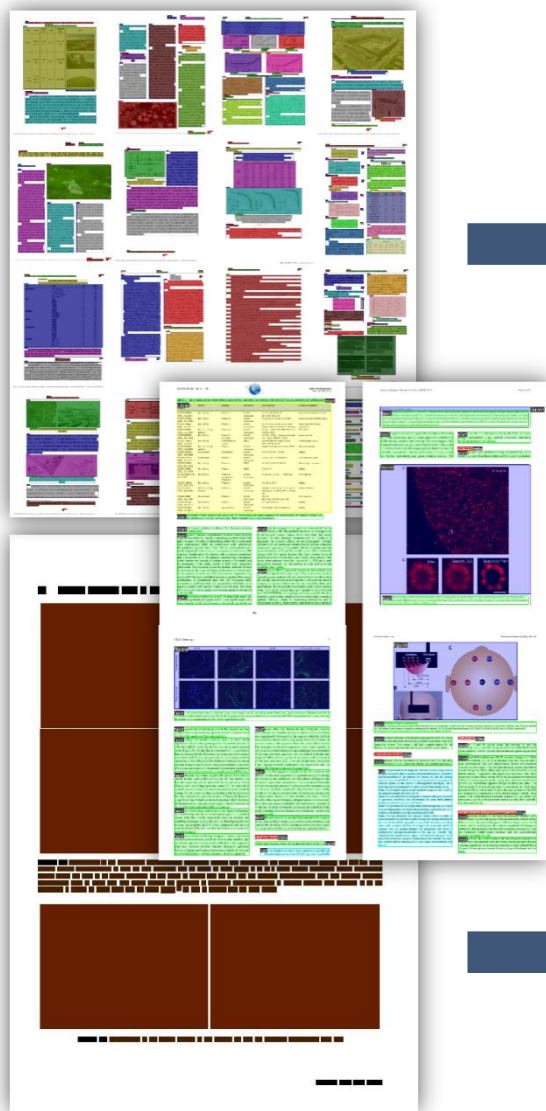
机器之心sota方法排行榜

文档块阅读顺序



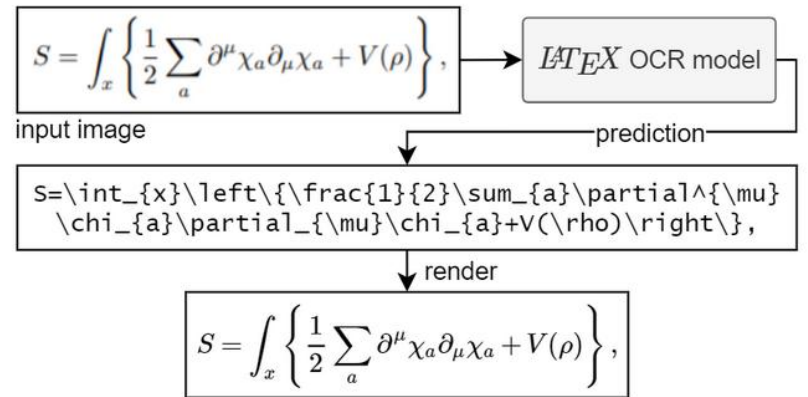
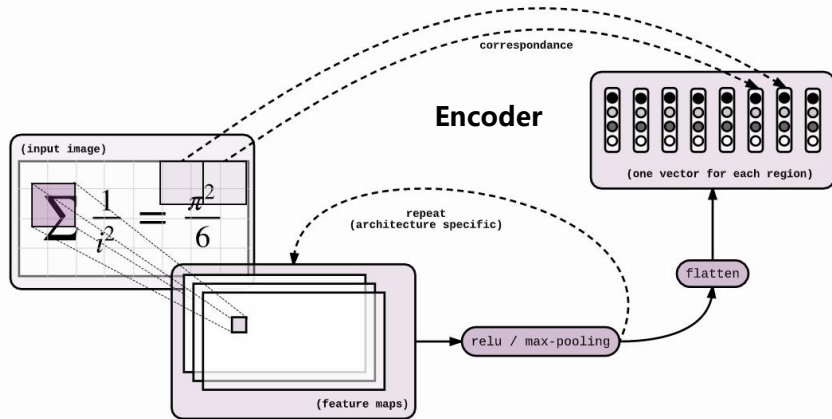
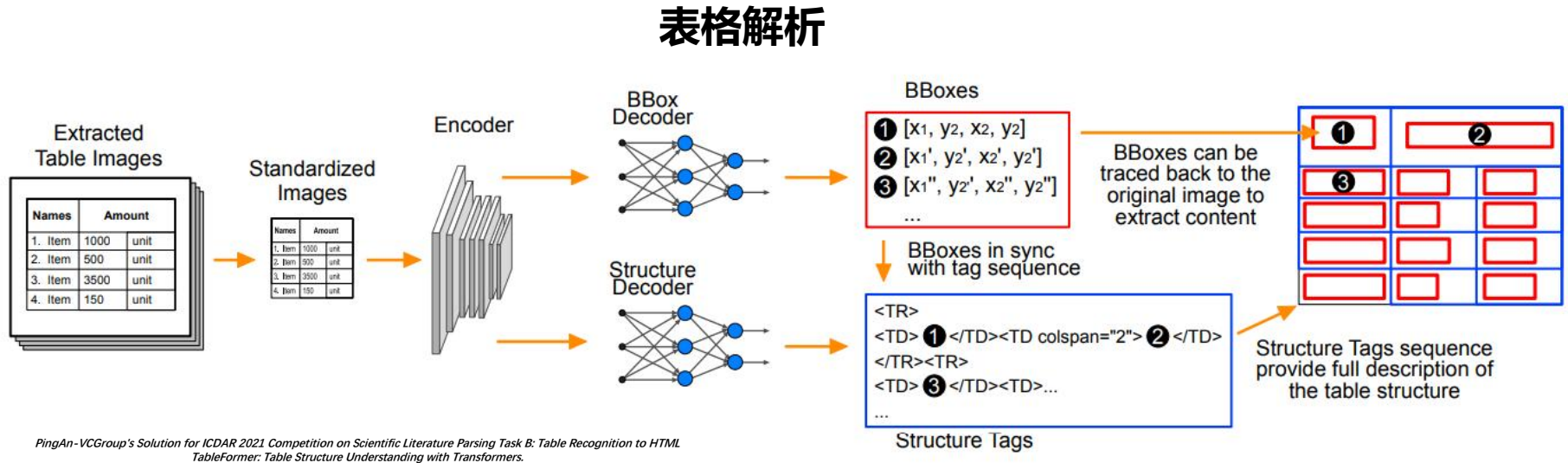
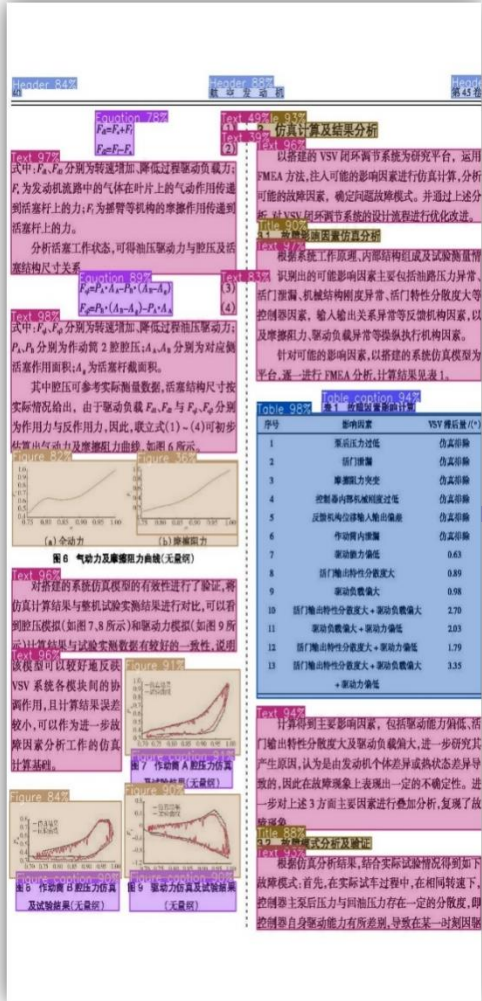
LayoutReader: Pre-training of Text and Layout for Reading Order Detection

基于目标检测的文档布局识别



图源自: <https://segmentfault.com/a/1190000042194108>

小米大语言模型中数据及知识挖掘



小米大语言模型中数据及知识挖掘



word



原生PDF



Scan Text

扫描PDF



图片JPG

Header 88% 航空发动机 第45卷

Text/Ocr

Table caption 94% 表1 故障因素影响计算

Equation 89%
$$F_q = P_A \cdot A_A - P_B \cdot (A_B - A_q)$$

$$F_q = P_B \cdot (A_B - A_q) - P_A \cdot A_A$$

Figure 84%

Table caption 94%

Figure caption 30%

Table caption 98%

序号	影响因素	仿真排除
1	原后压力过低	仿真排除
2	活门泄漏	仿真排除
3	摩擦阻力变大	仿真排除
4	控制器内部机械刚度过低	仿真排除
5	反研机构位移输入偏差	仿真排除
6	作用筒内溢漏	仿真排除
7	震动能力偏低	0.63
8	活门输出特性分散度大	0.89
9	震动负载偏大	0.98
10	活门输出特性分散度大 + 震动负载偏大	2.70
11	震动负载偏大 + 震动力偏低	2.03
12	活门输出特性分散度大 + 震动力偏低	1.79
13	活门输出特性分散度大 + 震动负载偏大 + 震动力偏低	3.35

Table caption 90%

Text/Ocr

以搭建的VSV闭环调节系统为研究平台,运用FMEA方法,注入可能的影响因素进行仿真计算,分析可能的故障因素,确定问题故障模式,并通过上述分析,对VSV闭环调节系统的设计流程进行优化改进。

3.1 故障影响因素仿真分析

```

1 ----
2 $$
3 \begin{array}{l}
4 F_{\mathrm{q}}=P_{\mathrm{A}} \cdot A_{\mathrm{A}} \cdot \left( A_{\mathrm{B}}-A_{\mathrm{q}} \right) \\
5 \left( A_{\mathrm{B}}-A_{\mathrm{q}} \right) \cdot P_{\mathrm{B}} \\
6 \left( A_{\mathrm{B}}-A_{\mathrm{q}} \right) \cdot P_{\mathrm{B}} \\
7 F_{\mathrm{q}}=P_{\mathrm{B}} \cdot \left( A_{\mathrm{B}}-A_{\mathrm{q}} \right) \cdot P_{\mathrm{A}} \\
8 \left( A_{\mathrm{B}}-A_{\mathrm{q}} \right) \cdot P_{\mathrm{A}} \\
9 \left( A_{\mathrm{B}}-A_{\mathrm{q}} \right) \cdot P_{\mathrm{A}} \\
10 \end{array}
11 $$
12 ----
13 表1 故障因素影响计算
14 \begin{tabular}{|c|c|c|}
15 \hline 序号 & 影响因素 & vsv 滞后量 \\
16 \hline 1 & 原后压力过低 & 仿真排除 \\
17 \hline 2 & 活门泄漏 & 仿真排除 \\
18 \hline 3 & 摩擦阻力变大 & 仿真排除 \\
19 \hline 4 & 控制器内部机械刚度过低 & 仿真排除 \\
20 \hline 5 & 反研机构位移输入输出偏差 & 仿真排除 \\
21 \hline 6 & 作用筒内溢漏 & 仿真排除 \\
22 \hline 7 & 震动能力偏低 & 0.63 \\
23 \hline 8 & 活门输出特性分散度大 & 0.89 \\
24 \hline 9 & 震动负载偏大 & 0.98 \\
25 \hline 10 & 活门输出特性分散度大 + 震动负载偏大 & 2.70 \\
26 \hline 11 & 震动负载偏大 + 震动力偏低 & 2.03 \\
27 \hline 12 & 活门输出特性分散度大 + 震动力偏低 & 1.79 \\
28 \hline 13 & 活门输出特性分散度大 + 震动负载偏大 + 震动力偏低 & 3.35 \\
29 \hline & +震动力偏低 & \\
30 \hline & +震动力偏低 & \\
31 \hline & +震动力偏低 & \\
32 \hline & +震动力偏低 & \\
33 \hline & +震动力偏低 & \\
34 \hline & +震动力偏低 & \\
35 \end{tabular}
36 ----
37 根据仿真分析结果,结合实际试验情况得到如下故障模式:
38 首先,在实际试车过程中,在相同转速下,控制器主原后压力
39 与回油压力存在一定的分散度,即控制器自身驱动能力有所差别,
40 导致在某一时刻因驱
    
```

转图片

布局元素及阅读顺序预测

Latex/table ocr / text

按阅读顺序拼接

数据及知识挖掘的质量控制



语义不通

一天，一个飘扬的僧侣走进了古老的江湖，忽然遇见一位独具风韵的武林前辈，拔了一根头发便飞舞起来，脚下的草地变成了一块汪洋的石板，蓝天之上竟然镶满了一朵朵武装和尚。这时路边儿一个撑满水袋的鱼缸坠落大地，武林前辈打破了鱼缸捉出鱼的思维，交给僧侣献上祝福。突然僧侣眼前浮现出一座空姐，那火红通向一名售票员手家。在庙中，售票员手持神秘的鱼缸铁丝，上绳子一头断气，这一切都只会使文字永无尽头。

就在此时并无关系的细雨纷飞之际，化身浓雾的尼姑暴跳而来，还有掩藏月亮的胡子，向东飘向远方。疾风迅雷，这困局乍然裂成一道红色鸿卷卷带，东风飘箭跃下，武林前辈被袖舞突窥隐世之瞳，这终便一份憾事未有期一个完美结局。

袅袅微雨，一场山枫满天，红霞漫舞，这不是孤独背影的恹恹终未结束，未来将会是武林新局诉说的开始。

无信息量

I'm not saying anything



近期，部分蓝筹股持续回调，个别蓝筹股回调幅度较大，有观点认为蓝筹股到了抄底的好时机，但不认为，不能盲目抄底蓝筹股，要深入研究后再进行选择，对于那些长期股可能并不具备抄底的价值，比如持续回调的猪肉股。

巴菲特说过，要用闲的价格买黄金，这才是价值投资。现在的部分蓝筹股，看似业绩不错，但实际上对应的股价不算便宜，虽然股价有所回落，但是依然处于估值区间的上限，对于这样的股票，即使投资者真的看好上市公司，也要选择股价回调到位后再行买入，而股价何时调整到位，投资者并不宜通过上市公司的定期报告研究。

公司股价调整结束，需要有一个必要条件，即大资金入场抄底，投资者虽然以价值投资的眼光选择股票，但是还是要通过技术分析，盘量特征来选择买入卖出的时点，这样才能使得自己的资金利用价值最大，避免资金长期持有横盘或者阴跌的股票。

看一家公司股价是否调整到位，并不是用市盈率和中净率来衡量的，市盈率跌到多少是底部？20倍还是10倍？这并没有一个准确的说法，但是每一个底部都会有大资金进入的痕迹，例如成交量温和放大，股价横盘或者缓慢上涨，这些都是大资金吸筹的信号，如果没有大资金吸筹，那么就算真是绩优股，股价的表现也不会特别好，所以投资者在看好一家公司后，也要等待股价出现买入信号再进场买入，谨慎做交易，将筹码

源：https://paper.bbtnews.com.cn/site1/fjsb.html/2021-05/11/content_465066.htm?iv=-1

近期，部分蓝筹股持续回调，个别蓝筹股回调幅度较大，有观点认为蓝筹股到了抄底的好时机，但不认为，不能盲目抄底蓝筹股，要深入研究后再进行选择，对于那些长期股可能并不具备抄底的价值，比如持续回调的猪肉股。

巴菲特说过，要用闲的价格买黄金，这才是价值投资。现在的部分蓝筹股，看似业绩不错，但实际上对应的股价不算便宜，虽然股价有所回落，但是依然处于估值区间的上限，对于这样的股票，即使投资者真的看好上市公司，也要选择股价回调到位后再行买入，而股价何时调整到位，投资者并不宜通过上市公司的定期报告研究。

公司股价调整结束，需要有一个必要条件，即大资金入场抄底，投资者虽然以价值投资的眼光选择股票，但是还是要通过技术分析，盘量特征来选择买入卖出的时点，这样才能使得自己的资金利用价值最大，避免资金长期持有横盘或者阴跌的股票。

看一家公司股价是否调整到位，并不是用市盈率和中净率来衡量的，市盈率跌到多少是底部？20倍还是10倍？这并没有一个准确的说法，但是每一个底部都会有大资金进入的痕迹，例如成交量温和放大，股价横盘或者缓慢上涨，这些都是大资金吸筹的信号

关注XX公众号，获取更多机会

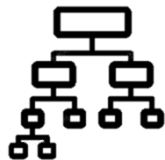
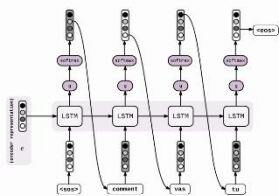
<image>图XX

查看全文

- ! 宗教信仰
 - ! 身份攻击
 - ! 侮辱
 - ! 威胁
 - ! 亵渎
 - ! 色情



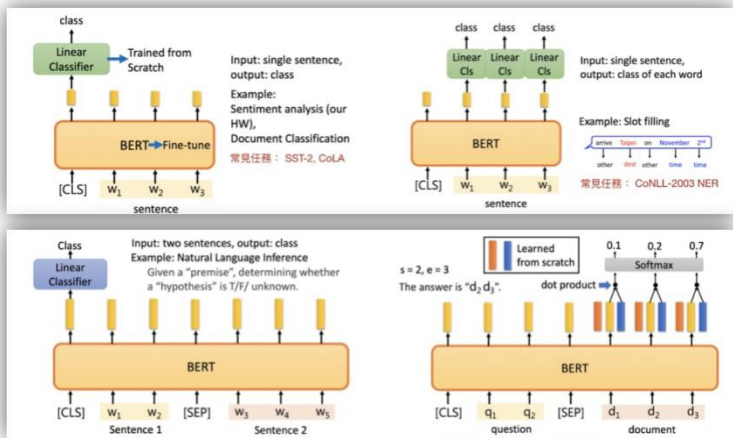
大模型与小模型联动提升模型迭代速度



MILM

小模型快速的质量发现及识别清洗

高质量的数据影响模型的质量



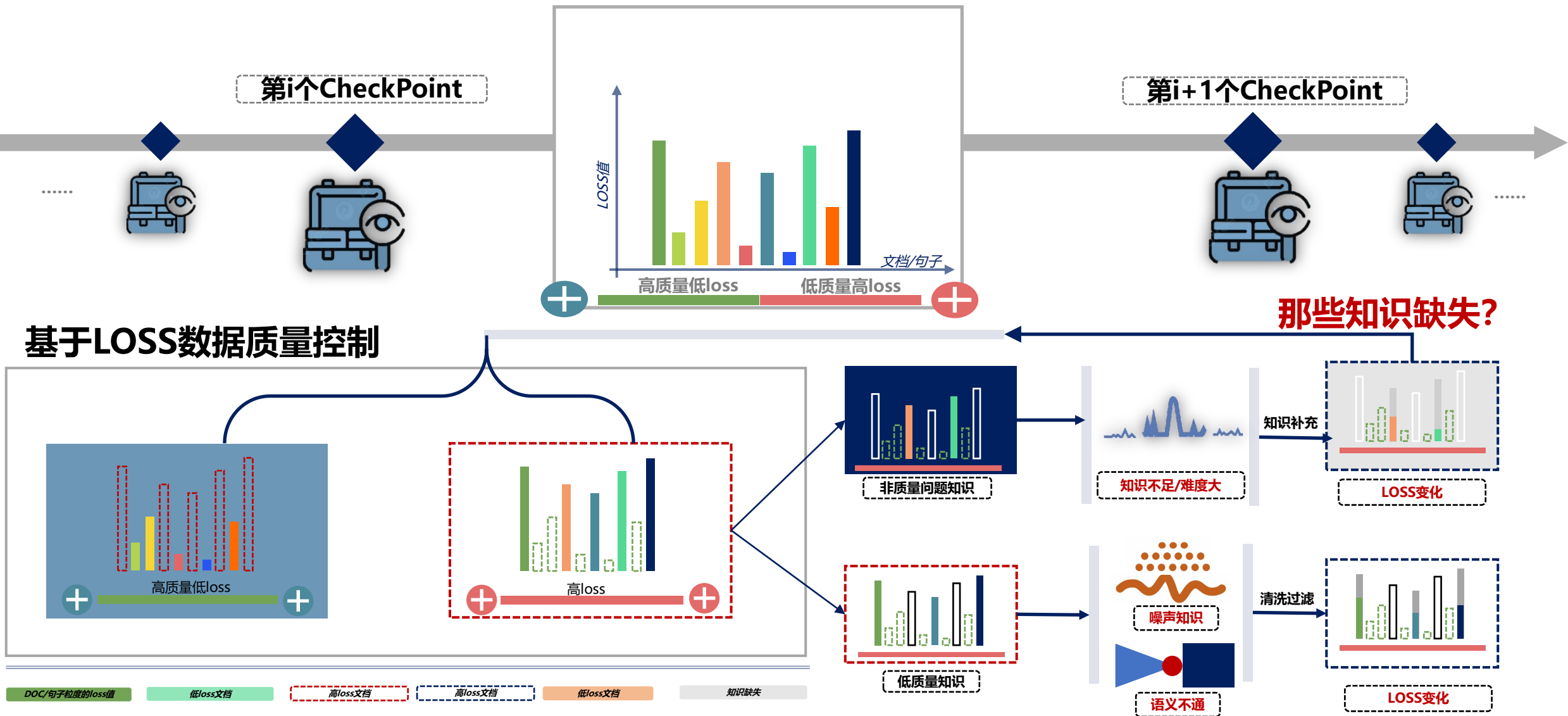
大模型能否发现质量问题?



小模型

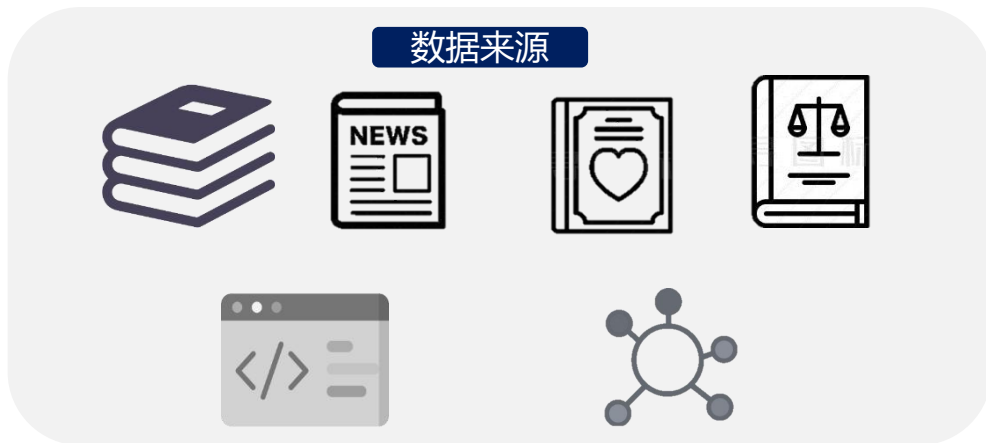
大模型

大模型与小模型联动提升模型迭代速度



▶ 丰富知识维度和控制采样提升模型迭代速度

粗粒度数据类别及组成



细化知识分布及统计



▶ 模型及数据质量评估的方法



小米模型及数据质量评估的方法



评估方法原则及维度



通用能力



专项能力

满足程度

0

1

2

3

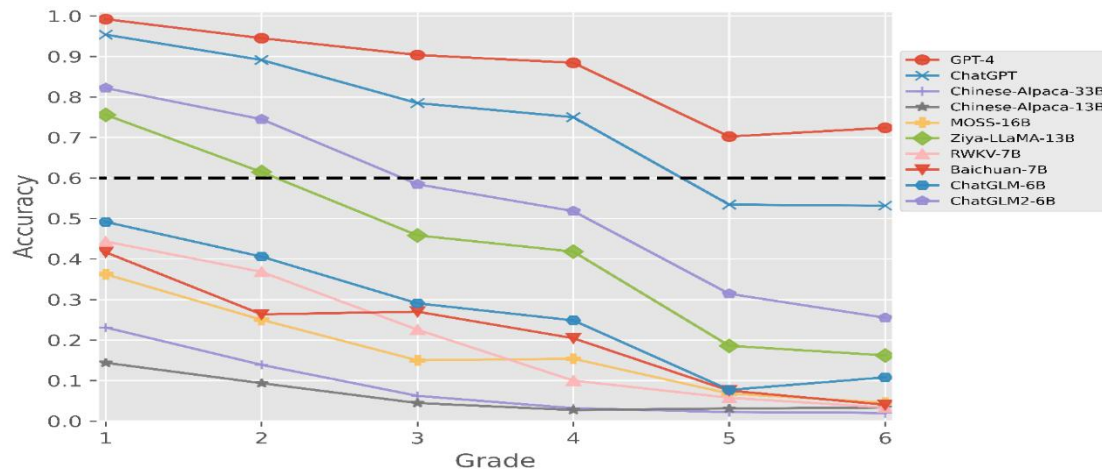
4

● 对齐开放原则

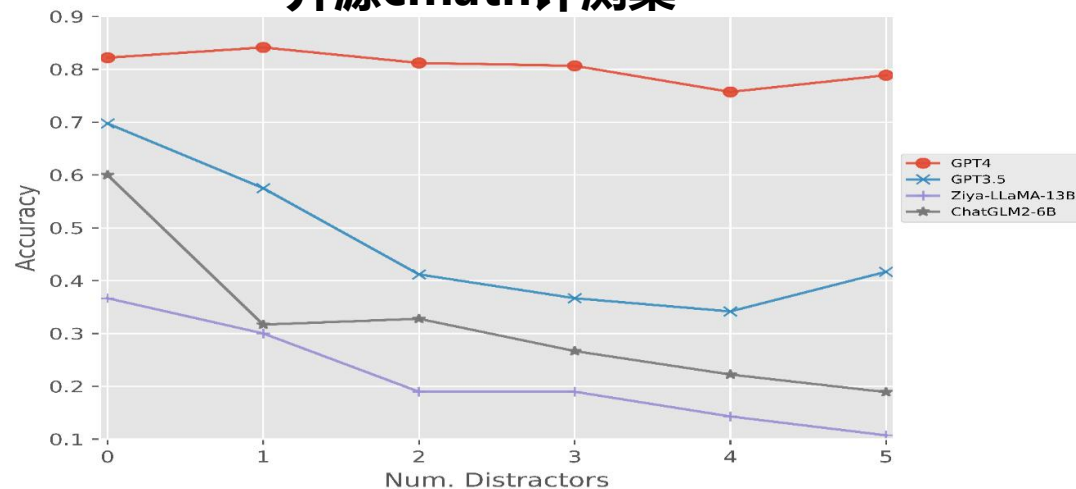
● 能力多维度评估

● 满足度细化

● 开源专项评测



开源cmath评测集



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

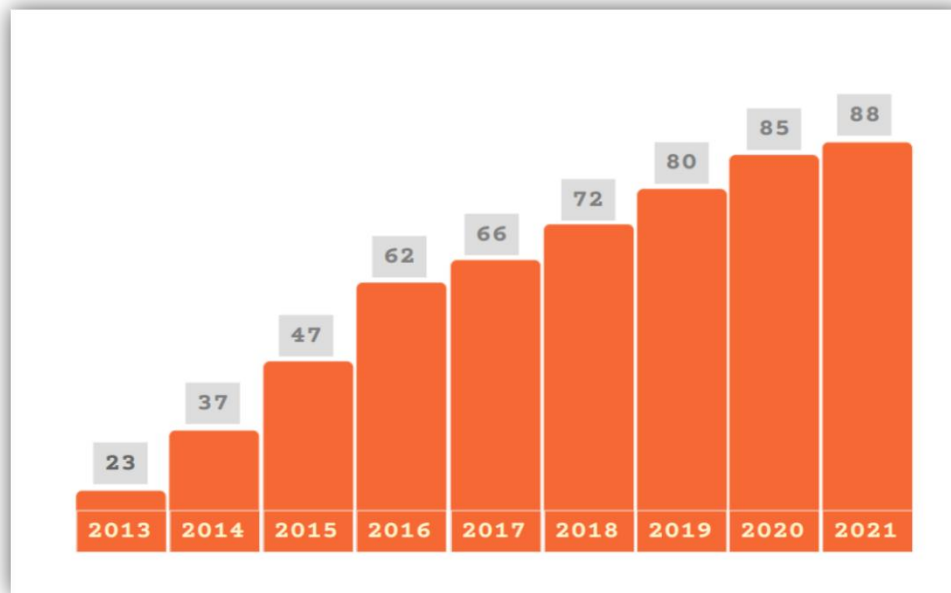
大语言模型中数据及知识的存储



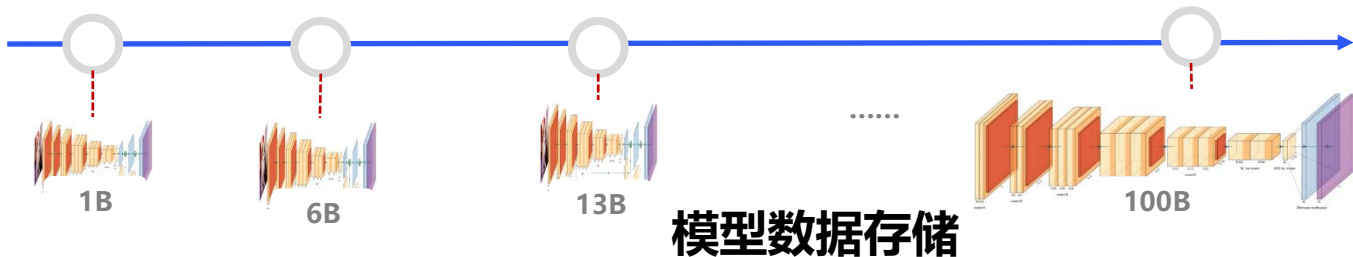
88, 343, 822

100TB数据

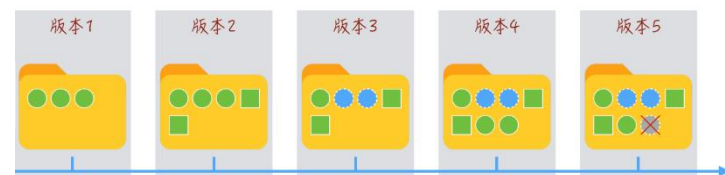
小文件存储



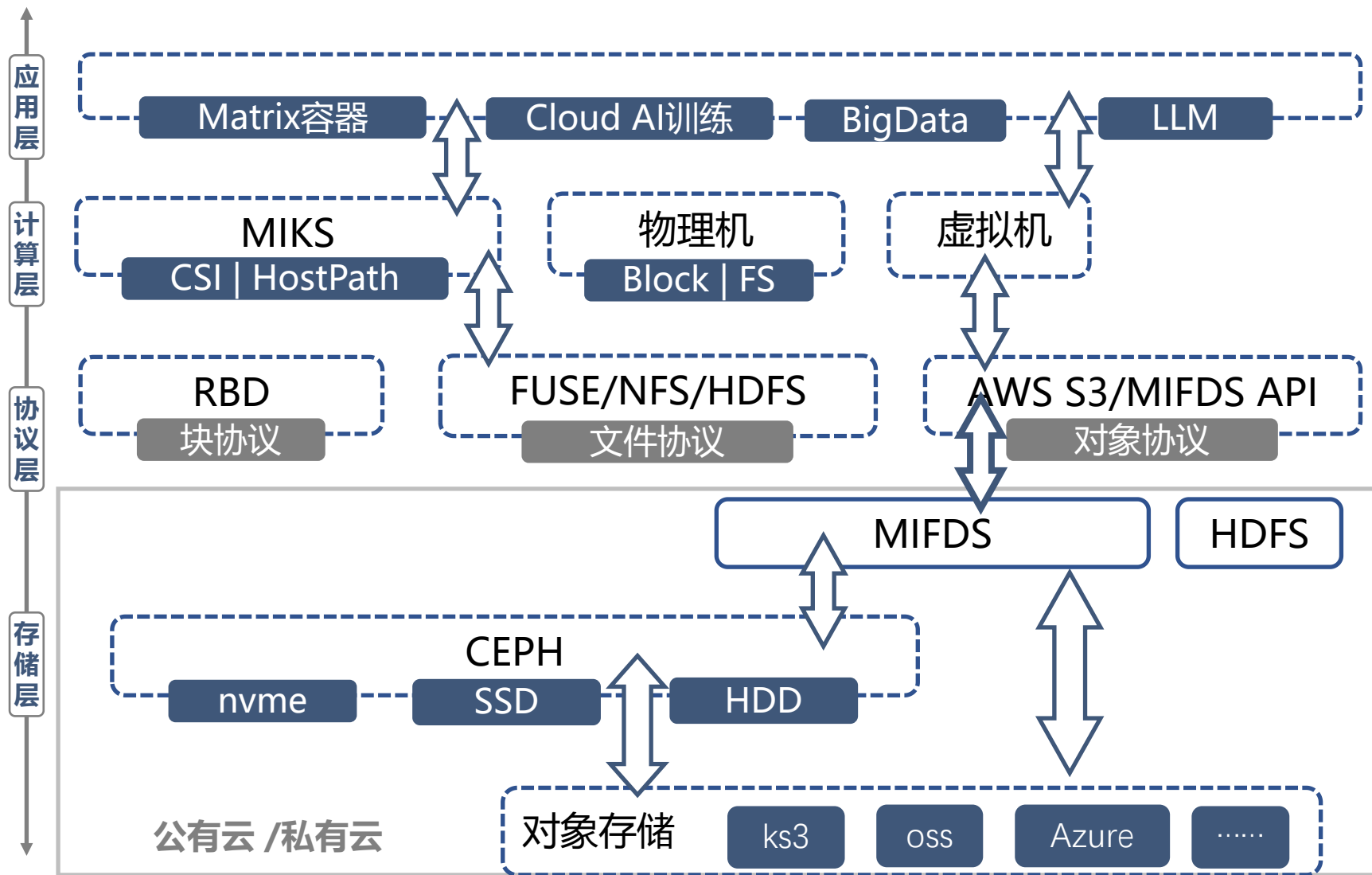
1. 海量小文件存储需求, 对小文件的预处理及访问和计算提出挑战
2. 不同参数规模的模型文件, 相同参数模型下不同超参的每个epoch的存储
3. 数据处理的不同阶段 (解析、清洗、去重、分类) 等中间态的数据存储



不同参数规模 **X** 不同版本



小米大语言模型中数据及知识挖掘



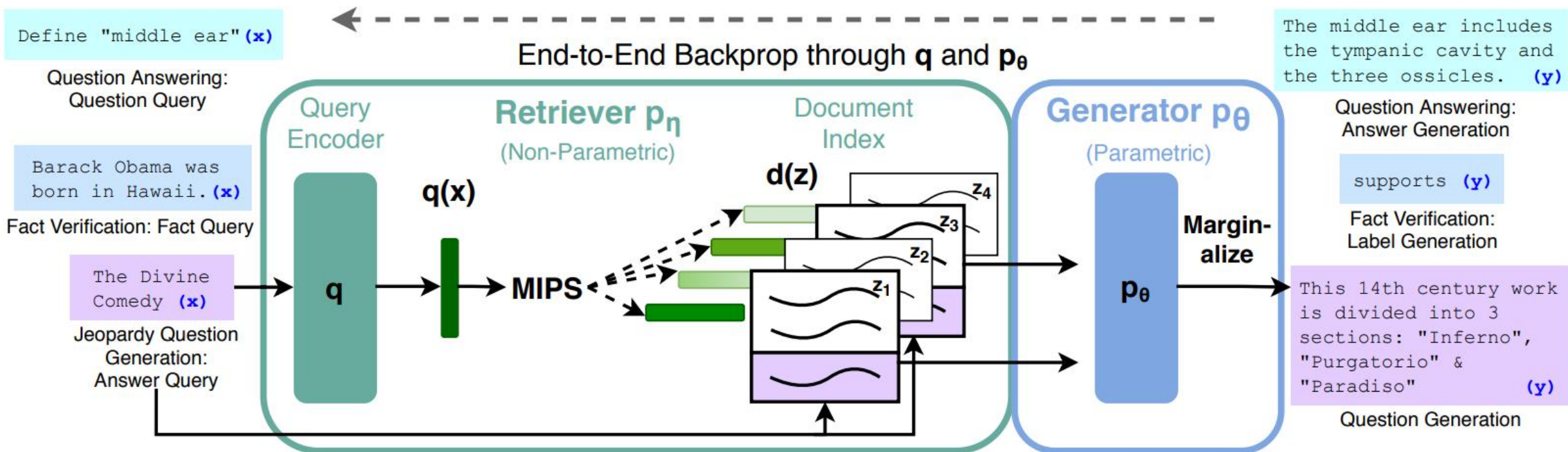
数据处理

模型训练

PART 03

大模型在小米领域场景下的应用和方法

LLM+ Retrieval-Augmented Generation的领域应用

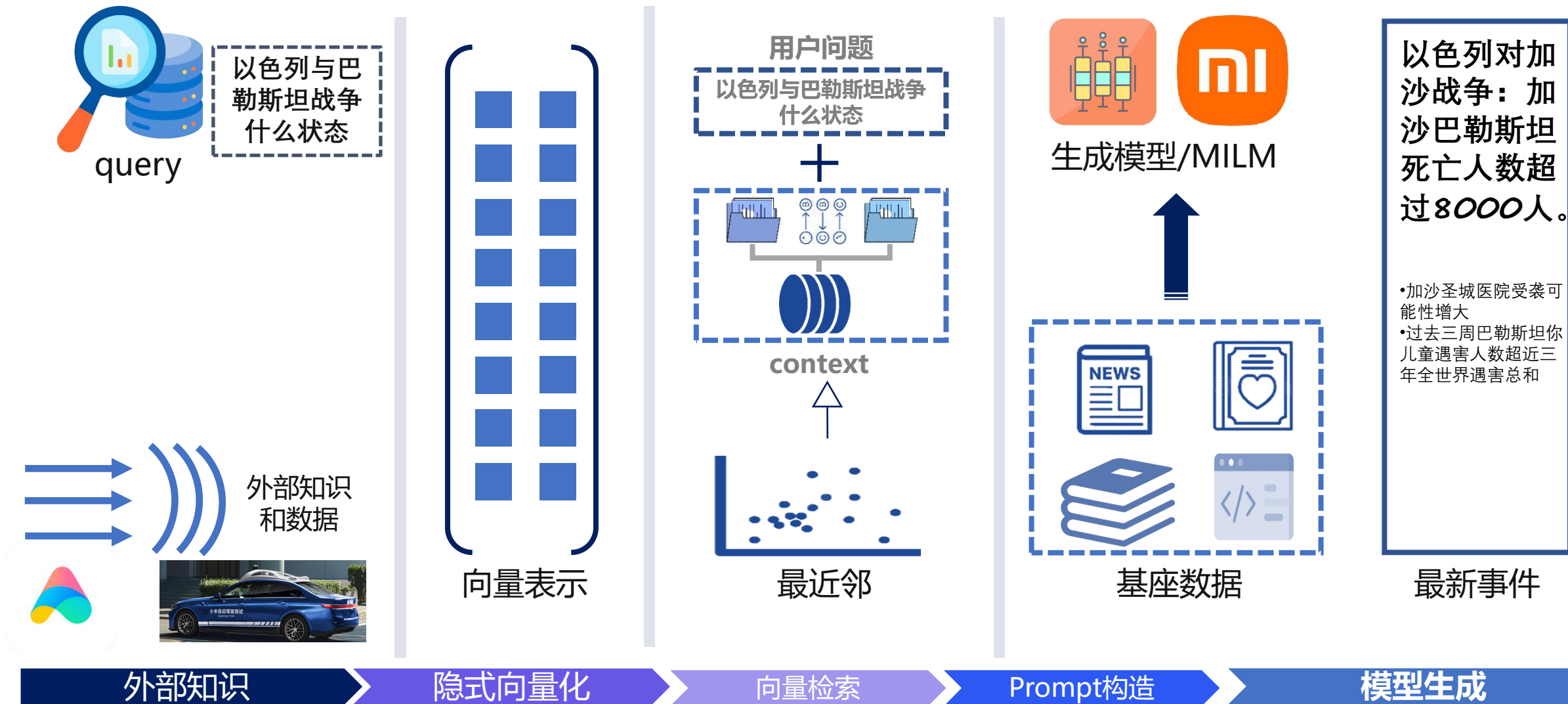


Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks

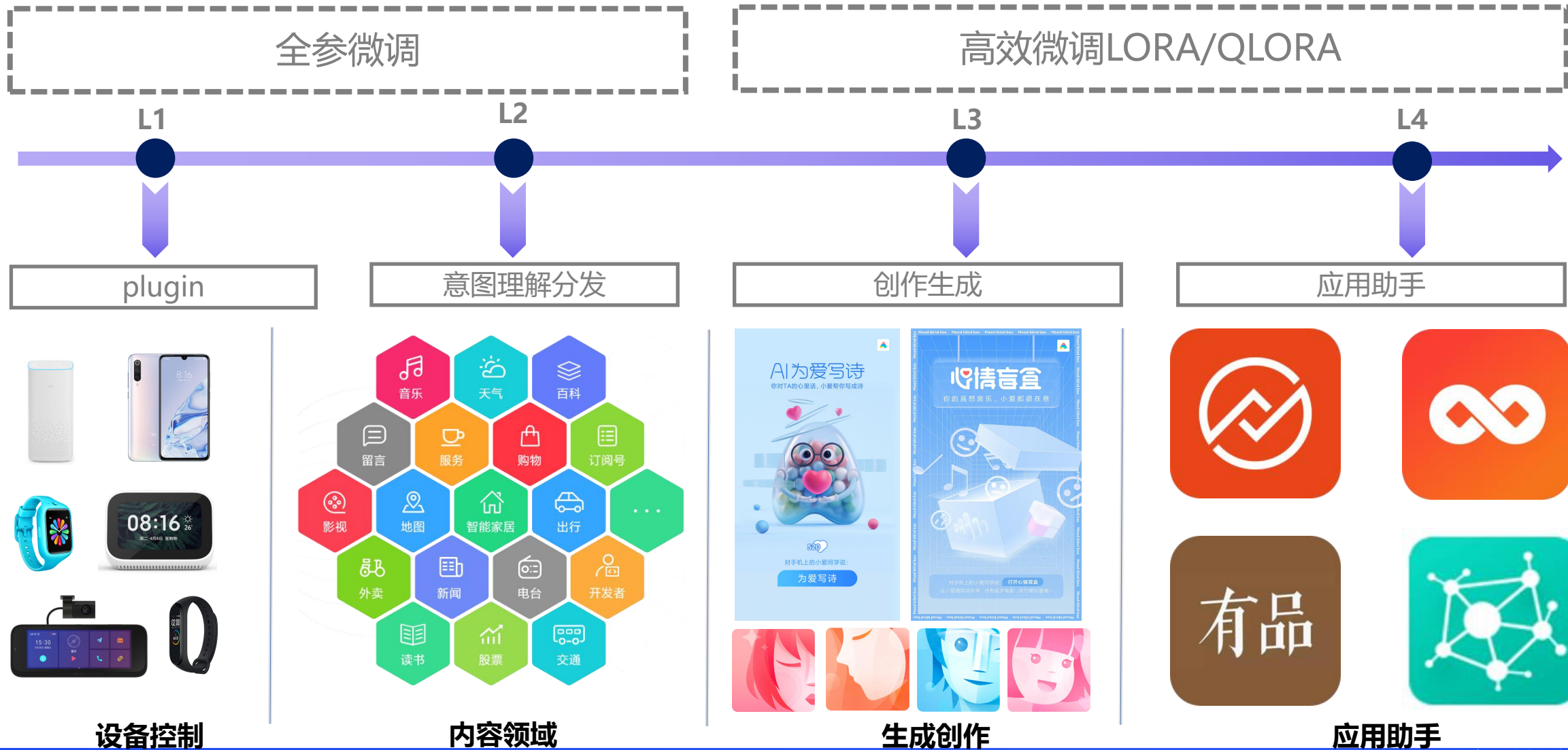
- 模型的知识/数据是隐式的编码在参数中，引入的知识是显式的，优秀的扩展性及更可控的时效性。

- 为生成模型引入更丰富的关联和参考信息，可以提升生成模型的性能和丰富度。

LLM+ Retrieval-Augmented Generation的领域应用



▶ 领域数据高效全参微调的小爱问答场景应用方法



▶ 大模型赋能小米应用场景

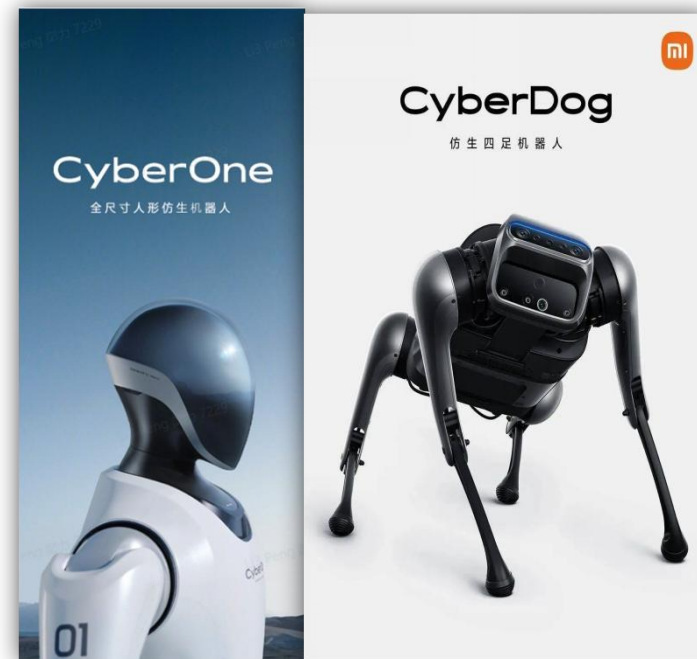
轻量化 本地部署



智能驾驶



操作系统



仿生机器人

PART 04

总结及展望



▶ 总结

01

数据依然是影响LLM性能的决定元素之一

02

知识获取、存储、评估等不断完善

03

大语言模型在小米业务下逐渐成长

04

数据为中心的大语言模型会更强大

THANKS

