

机器阅读理解研究综述

张天成¹⁾ 王雅婷¹⁾ 李 凡¹⁾ 孙相会¹⁾ 于明鹤²⁾ 于 戈¹⁾

¹⁾(东北大学计算机科学与工程学院 沈阳 110819)

²⁾(东北大学软件学院 沈阳 110819)

摘 要 机器阅读理解(Machine Reading Comprehension, MRC)作为自然语言处理领域的核心任务之一,旨在赋予机器理解文本并回答相关问题的能力。本文综述了 MRC 领域的最新研究进展,系统地梳理了自 2015 年以来的主要研究工作。文章首先概述了 MRC 任务的不同形式,包括填空式、多项选择式、抽取式和自由答案式,并对其评估方法进行了详细分析。随后,本文深入分析了机器阅读理解模型架构的发展历程,从基础的通用结构和依赖注意力机制的模型出发,进一步探讨了预训练语言模型(Pre-trained Language Models, PLMs)及大语言模型(Large Language Models, LLMs)技术的应用,并对融合推理结构的模型进行了详细阐述。此外,文章还讨论了 MRC 面临的挑战,如无答案问题、多答案问题、对话型、多轮交互型、跨语言与跨模态型阅读理解,以及零样本和少样本问题,并分析了相应的解决方案;最后讨论其应用和发展趋势。本文的主要贡献包括:(1)对 MRC 不同形式及其相关数据集的系统综述;(2)对 MRC 架构的深入探讨;(3)分析 MRC 面临的挑战。本文旨在为 MRC 领域的研究者提供全面的参考,以促进具有更强理解能力和泛化能力的模型设计与发展,推动相关领域的深入探索。

关键词 机器阅读理解;自然语言处理;预训练语言模型;注意力机制;推理结构

中图法分类号 TP18 **DOI 号** 10.11897/SP.J.1016.2026.00979

A Survey on Machine Reading Comprehension

ZHANG Tian-Cheng¹⁾ WANG Ya-Ting¹⁾ LI Fan¹⁾ SUN Xiang-Hui¹⁾ YU Ming-He²⁾ YU Ge¹⁾

¹⁾(Department of Computer Science and Engineering College, Northeastern University, Shenyang 110819)

²⁾(Department of Software College, Northeastern University, Shenyang 110819)

Abstract Machine Reading Comprehension (MRC) is one of the core tasks in Natural Language Processing (NLP), which aims at equipping machines with human-level text comprehension and question-answering abilities. This paper provides a comprehensive review of MRC research since 2015 and systematically traces its rapid evolution. In terms of model architecture, attention-based neural models have formed a universal framework; the word embedding layer maps text into dense vectors, the encoding layer models contextual dependencies, the interaction layer calculates question-text semantic matching, and the output layer generates predicted answers. This paradigm significantly enhances model performance, propelling MRC into the era dominated by deep learning. Traditional MRC models are primarily distinguished by the different attention mechanisms employed in their interaction layers, we introduce MRC models based on attention mechanisms, including interactive attention, self matching attention, and sparse attention. The emer-

收稿日期:2025-06-01;在线发布日期:2026-01-07。本课题得到国家自然科学基金面上项目(No. 62272093)、国家自然科学基金重点项目(No. 62137001)、国家自然科学基金国际(地区)合作与交流项目(No. 62461146205)资助。张天成,博士,教授,中国计算机学会(CCF)高级会员,主要研究领域为人工智能、时空数据管理、智慧教育。E-mail: tczhang@mail.neu.edu.cn。王雅婷(通信作者),硕士,中国计算机学会(CCF)会员,主要研究领域为知识追踪、自然语言处理、可解释人工智能。E-mail: 2301964@stu.neu.edu.cn。李 凡,博士,中国计算机学会(CCF)会员,主要研究领域为机器学习、深度学习、知识追踪、可解释人工智能。孙相会,硕士,中国计算机学会(CCF)会员,主要研究领域为智慧教育、可解释人工智能。于明鹤,博士,副教授,中国计算机学会(CCF)会员,主要研究领域为数据库、信息检索、智能教育。于 戈,博士,教授,中国计算机学会(CCF)会士,主要研究领域为分布式与并行数据库、OLAP 与数据仓库、数据集成、图数据管理。

gence and widespread adoption of Pre-trained Language Models (PLMs) have revolutionized the MRC field. By designing and optimizing pre-training tasks on large-scale unlabeled text data, PLMs have gained rich linguistic knowledge and semantic representations. Building on this foundation, Large Language Models (LLMs) have further advanced in expanded parameter scales and enhanced few-shot/zero-shot learning capabilities, enabling them to achieve state-of-the-art performance in MRC tasks. Additionally, we provide a detailed elaboration on models integrated with reasoning structures. These models aim to enhance the logical reasoning and deductive capabilities of MRC systems, enabling them to tackle complex questions that require deep semantic analysis and multi-step inference beyond simple information extraction. The current research mainly focuses on the following types of reasoning: multi-hop reasoning, logical reasoning, numerical reasoning, and commonsense reasoning.

Furthermore, we discuss the key task challenges currently facing the MRC field, along with the corresponding state-of-the-art solutions. These challenges cover six parallel types: Unanswerable MRC, where the model must identify and flag the absence of valid answers when the text lacks sufficient information; Multi-answer MRC, where the model is required to accurately extract all relevant answer fragments for a single question with multiple correct answers; Dialogue-based MRC, where dynamic interactions between the model and users demand contextual coherence across multiple Q&A rounds and continuous understanding updates based on prior interactions; Multi-round Interactive MRC, where a single question is broken down into multiple sub-questions, with the model leveraging multi-task learning to deeply mine semantic information; Cross-lingual/Cross-modal MRC, where the model needs to achieve semantic understanding and alignment across different languages or multiple modalities (e. g. , text, image); and Zero-shot/Few-shot MRC, where the model is required to complete comprehension tasks with extremely limited or no labeled training data. Our contributions are as follows: (1) We systematically organize the diverse question formats in the Machine Reading Comprehension domain and their corresponding datasets. (2) We conduct an in-depth analysis of the technological evolution within the MRC field. (3) We identify the prevailing hotspots and persistent challenges in MRC, and further present a comprehensive review of the existing solutions to these issues.

In summary, our work provides a systematic reference for MRC research, facilitating the design and development of models with stronger comprehension and generalization capabilities and further promoting in-depth exploration in related fields.

Keywords machine reading comprehension; natural language processing; pre-trained language models; attention mechanism; reasoning structure

1 引 言

自然语言处理是人工智能的重要分支,是实现人工智能的核心技术,主要研究如何处理、分析以及应用自然语言。教会机器阅读文本并且理解人类语言是自然语言处理领域的重要任务,机器阅读理解(Machine Reading Comprehension, MRC)的目标就是利用自然语言处理技术使得计算机能够像人一样阅读并且理解文章,它有着很多应用场景,如搜索引

擎中的智能问答,电商领域的智能客服以及对话系统等。早期研究可追溯至 20 世纪 70 年代,如 Lehnert 提出的 QUALM^[1] 系统,但受限于手工编码规则,泛化能力较弱。1999 年,首个自动阅读理解系统 DeepRead^[2] 采用词袋模型和手工规则,准确率仅达 40%,仍难以推广。传统机器学习方法依赖浅层特征提取,导致 MRC 发展缓慢。

2015 年,随着深度学习的兴起,DeepMind 的 Hermann 等人^[3] 提出基于神经网络和注意力机制的 Attentive Reader 与 Impatient Reader 模型,并

在 CNN&Daily Mail 数据集上取得突破性进展,标志着神经机器阅读理解(NMRC)的崛起。此后,大规模数据集(如 SQuAD^[4]、TriviaQA^[5]、HotpotQA^[6])的构建进一步推动了 MRC 研究,并在形式上逐渐细化为填空式、多项选择式、抽取式和自由答案式四类。在模型架构方面,基于注意力机制的神经模型形成通用框架:词嵌入层将文本映射为稠密向量,编码层建模上下文依赖,交互层计算“问题-文本”语义匹配,最终输出预测答案。这一范式显著提升了性能,使 MRC 进入深度学习主导时代。由于传统机器阅读理解模型的显著特点主要体现在其交互层所采用的注意力机制的差异性,因此本文介绍了基于注意力机制的 MRC 模型,包括交互注意力、自匹配注意力和稀疏注意力。预训练语言模型(Pre-trained Language Models, PLMs)的兴起,特别是以此为基础的大模型(Large Language Models, LLMs)技术的发展进一步推动 MRC 在多个基准(如 SQuAD)上超越人类表现。尽管 PLMs 在抽取式问答上表现优异,但多数模型只能做简单的模式匹配,MRC 的深层推理能力仍然不足。当前研究主要聚焦以下推理类型:多跳推理、逻辑推理、数值推理和常识推理。此外,针对 MRC 面临的任务挑战,如无答案问题型、多答案问题型、对话型、多轮交互型、跨语言与跨模态型的阅读理解,以及零样本和少样本问题,许多模型提出了很好的解决方案,本文对此进一步探讨。本文的目的是对从 2015 年以来机器阅读理解(MRC)领域的研究任务、相关数据集以及模型做综述。本文的主要贡献包括:

(1)系统综述 MRC 不同的形式分类及其评估方法,涵盖填空式、选择式、抽取式和生成式,并分析代表性数据集;

(2)深入探讨 MRC 架构演进,从早期注意力机制到 LLMs 再到聚焦推理结构的技术进展;

(3)探讨了 MRC 面临的任务挑战及与此相关的技术发展,讨论了 MRC 的应用前景,并展望未来方向。

本文旨在为 MRC 研究提供系统参考,以促进具有更强理解能力和泛化能力的模型设计与发展,推动相关领域的深入探索。文章整体安排如下:第 2 节概述 MRC 不同形式及其评估方法,分析相关的数据集;第 3 节介绍 MRC 模型架构发展,包括通用结构、基于注意力的 MRC 模型、大模型驱动的 MRC 模型和基于推理结构的 MRC 模型;第 4 节分析 MRC 任务目前的主要挑战及技术;第 5 节讨论 MRC 的应

用以及未来的发展趋势;第 6 节对全文做总结。

2 机器阅读理解概述及形式

机器阅读理解(MRC)任务是为了使得计算机具有对自然语言文本理解的能力,像人类一样阅读并且理解一篇文章。MRC 可以用一个三元组(D, Q, A)来描述,其中 D 代表文章(Document)^①, Q 表示问题(Question), A 表示答案(Answer),即给定一篇文章 D 和一些与文章 D 相关的问题 Q ,要求模型通过阅读 D 之后给出 Q 的正确答案 A ,建模给定 D 和 Q 的条件下预测 A 的概率: $P(A|D, Q)$ 。

2.1 任务形式

基于答案形式的不同,阅读理解任务可主要分为四类:填空式、多项选择式、抽取式和自由答案式^[7]。下面对这四种类型任务分别进行叙述并介绍相关的数据集。

(1)填空式。填空式阅读理解是指给定一篇文章 D 和一个与文章相关的问题 Q , Q 通过删除掉句子中某一个单词构成,要求模型根据 D 能够正确地填写出 Q 中缺失的单词 a ,且 $a \in D$ 。相关的数据集如 CNN&Daily Mail^[3], CBT^[8] 和 CLOTH^[9]。

(2)多项选择式。多项选择式任务要求从多个候选答案中选出所有正确答案,其任务是在给定文章 D 和问题 Q 的条件下,从候选答案集合 $A = \{A_1, A_2, \dots, A_n\}$ 中选出正确答案。该过程可通过建模概率: $P(A_i|D, Q)$ 来实现,其中 $A_i \in A$ 。相关数据集如 MCTest^[10] 和 RACE^[11]。

(3)抽取式。抽取式阅读理解任务是填空式任务的扩展,要求从原文中提取不定长度的连续文本作为答案,而不仅限于单个词。具体来说,给定文章 D 和问题 Q ,问题的答案由 D 中的一段连续的单词构成。可以表示为 $P(A|D, Q)$,其中 $A = \{t_i, t_{i+1}, \dots, t_{i+k}\} (1 \leq i \leq i+k \leq n)$, n 代表 D 中单词的个数, k 代表答案的长度。相关数据集如 SQuAD^[4], NewsQA^[12] 和 TriviaQA^[5]。

(4)自由答案式。抽取式任务受限于从原文中提取文本,难以满足更实际的问题需求,自由答案式阅读理解任务则更为灵活,允许答案超越原文词句,贴近人类的概括与自然表达方式。可以表示为 $P(A|D, Q)$,其中 $A \subseteq D$ 或 $A \not\subseteq D$ 。相关数据集如 MS MARCO^[13], DuReader^[14] 和 NarrativeQA^[15]。

① 本文中文章(Document)、段落(Passage)和文本均是同样的概念。

2.2 评估方法

不同的 MRC 任务采用不同的评估指标来衡量性能。填空式任务与多项选择式任务均属于客观题型,可以通过准确率 Acc(Accuracy)衡量模型的性能。例如对于测试集中的所有问题 $Q = \{Q_1, Q_2, \dots, Q_m\}$, 其中 m 代表问题的个数。若模型预测的 m 个答案中有 n 个是正确的,则模型的准确率为 n/m 。

抽取式任务属于半客观题型,通常采用精确匹配 EM(Exact Match)和 F1 分数来评估模型。EM 评估指标可以看作是准确率的扩展,就抽取式任务来讲,EM 要求预测出来的所有单词要和标准答案的所有单词要完全一致,EM 值才为 1,否则为 0。F1 值的计算方式是一种模糊匹配,它是精确率和召回率之间的调和平均数:

$$F_1 = \frac{2 \times P \times R}{P + R} \quad (1)$$

P 代表精确率(Precision),是指模型预测答案中的单词有多大比例是标准答案中的单词。 R 代表召回率(Recall),是指标准答案中的单词有多大比例在预测答案中出现。

对于自由答案式任务,由于其答案形式不固定,一般采用单词水平的匹配率作为评分标准,常用标准为 ROUGE-L^[16] 和 BLEU^[17]。ROUGE-L 用来计算标准答案和预测答案的最长公共子序列(Longest Common Subsequence, LCS),计算公式如下:

$$\begin{aligned} R_{LCS} &= \frac{LCS(\mathbf{X}, \mathbf{Y})}{m}, \\ P_{LCS} &= \frac{LCS(\mathbf{X}, \mathbf{Y})}{n}, \\ F_{LCS} &= \frac{(1 + \beta)^2 R_{LCS} P_{LCS}}{R_{LCS} + \beta^2 P_{LCS}} \end{aligned} \quad (2)$$

\mathbf{X} 代表长度为 m 的预测答案, \mathbf{Y} 代表长度为 n 的真实答案, $LCS(\mathbf{X}, \mathbf{Y})$ 代表 \mathbf{X} 和 \mathbf{Y} 的最长公共子序列, P_{LCS} 和 R_{LCS} 分别表示精确率和召回率, β 参数用于调控精确率与召回率的重要性水平。BLEU 指标最初被提出用于翻译结果的性能评估工作,拓展至 MRC 任务场景时,用于衡量模型所输出的预测答案与实际标注的真实答案二者之间的相似程度,计算公式如下:

$$BLEU = BP \cdot \exp\left(\sum_{n=1}^N w_n \log P_n(\mathbf{X}, \mathbf{Y})\right) \quad (3)$$

其中,

$$P_n(\mathbf{X}, \mathbf{Y}) = \frac{\sum_i \sum_k \min(h_k(\mathbf{x}_i), \max(h_k(\mathbf{y}_i)))}{\sum_i \sum_k h_k(\mathbf{x}_i)} \quad (4)$$

$h_k(\mathbf{x}_i)$ 计算候选答案中 k 元词组的个数,类似的, $h_k(\mathbf{y}_i)$ 计算标准答案中 k 元词组的个数。由于候选答案中,句子越短的答案 $P_n(\mathbf{X}, \mathbf{Y})$ 分数越高,因此添加一个惩罚因子 BP :

$$BP = \begin{cases} 1, & l_x > l_y \\ e^{-\frac{l_y}{l_x}}, & l_x \leq l_y \end{cases} \quad (5)$$

表 1 列举了本文介绍的所有数据集及相应的评估方法。

3 机器阅读理解模型架构

3.1 通用结构

机器阅读理解文本内容大致需要以下流程:(1)将段落和问题这种文本形式的无结构数据表示为计算机可以处理的形式;(2)使段落或问题中某个单词能够关注其上下文信息,更好地表示段落或问题;(3)根据问题检索出段落中与问题最相关的部分;(4)从检索出来的文章片段中归纳得到答案。从整个流程可以看出,每一个步骤都有明确的目的,对应着神经网络中的某一层。因此用于 MRC 任务的深度学习模型的整体框架主要包括如下几层:词嵌入层、编码层、交互层、答案输出层,如图 1(a)所示。

词嵌入层对应于步骤(1),作用是将段落和问题嵌入到低维的向量空间中,用每一个向量表示一个单词。经典的方法包括字符嵌入^[39]、词嵌入^[40]、上下文嵌入(Word2Vec^[41]、GloVe^[42]和 ELMo^[43])、特征级别的嵌入^[44]及预训练编码器。

编码层对应于步骤(2),作用是编码段落和问题中单词的语义信息,使得每一个单词可以关注到它的上下文,常用的特征提取器有基于循环神经网络(RNNs)的变体(如 LSTM^[45]、GRU^[46])和 seq2seq 结构的模型 Transformer^[47]。

交互层对应于步骤(3),作用是将段落的语义信息与问题的语义信息融合,让模型学习到段落中与问题最相关的部分。最常用的方法就是注意力机制(Attention),注意力机制可以被视为是一个查询向量(query)和一组键值对向量(key-value pairs)的映射过程。整个过程首先是利用函数 f 衡量 query 和 key 之间的相似度,生成一个权重分数向量,然后将权重分数向量归一化后对 value 加权求和,得到的结果就是 query 对 key-value pairs 的注意力。具体计算公式如下:

表 1 机器阅读理解数据集

数据集	发布时间	文章来源/数量	文章类型	问题来源/数量	答案类型	评估指标
CNN&Daily Mail ^[5]	2015	新闻/3×10 ⁵	单段落型	人工合成/1.4×10 ⁶	填空式	Acc
CBT ^[8]	2015	儿童读物/108	单段落型	人工合成/6.8×10 ⁵	填空式	Acc
CLOTH ^[9]	2016	英语考试/1×10 ⁵	单段落型	英语考试/1×10 ⁵	填空式	Acc
MCTest ^[10]	2013	儿童读物/500	单段落型	众包/2×10 ³	多项选择式	Acc
RACE ^[11]	2018	英语考试/5×10 ⁴	单段落型	英语考试/8.7×10 ⁵	多项选择式	Acc
SQuAD ^[4]	2016	维基百科/536	单段落型	众包/1×10 ⁵	抽取式	EM/F1
NewsQA ^[12]	2017	新闻/1×10 ⁴	单段落型	众包/1×10 ⁵	抽取式	EM/F1
SQuAD 2.0 ^[18]	2018	维基百科/536	单段落型	众包/1.5×10 ⁵	抽取式	EM/F1
TriviaQA ^[5]	2017	网页搜索/6.6×10 ⁵	多段落型	搜索日志/4×10 ⁴	抽取式	EM/F1
HotpotQA ^[6]	2018	维基百科	多段落型	众包/1.13×10 ⁵	抽取式	EM/F1
WIKIHOP ^[19]	2018	维基百科/5.1×10 ⁴	多段落型	众包/5.1×10 ⁴	抽取式	Acc
DROP ^[20]	2019	维基百科/7×10 ³	单段落型	众包/9.7×10 ⁴	自由答案式	EM/F1
QUOREP ^[21]	2019	维基百科/4.7×10 ³	单段落型	众包/2.4×10 ⁴	抽取式	EM/F1
MS MARCO ^[13]	2016	搜索引擎/2×10 ⁸	多段落型	搜索日志/1×10 ⁶	自由答案式	ROUGE-L/BLEU
DuReader ^[14]	2018	搜索引擎/1×10 ⁶	多段落型	搜索日志/2×10 ⁵	自由答案式	ROUGE-L/BLEU
NarrativeQA ^[15]	2017	小说和电影剧本 1.5×10 ³	多段落型	众包/4.6×10 ⁴	自由答案式	ROUGE-L/BLEU
LogiQA ^[22]	2020	公务员考试/8.7×10 ³	单段落型	公务员考试/8.7×10 ³	多项选择式	Acc
ReClor ^[23]	2020	研究生入学考试/6.1×10 ³	单段落型	研究生入学考试/6.1×10 ³	多项选择式	Acc
RACENum ^[24]	2020	英语考试/1.2×10 ³	单段落型	英语考试/1.2×10 ³	自由答案式	Acc
CommonsenseQA ^[25]	2019	ConceptNet 图知识库/3.2×10 ⁷	单段落型	众包/1.2×10 ⁴	多项选择式	Acc
Cosmos QA ^[26]	2019	个人叙事博客/2.2×10 ⁴	单段落型	众包/3.6×10 ⁴	多项选择式	Acc
ART ^[27]	2019	ROCStories 故事集/2×10 ⁴	单段落型	众包/2×10 ⁴	多项选择式	Acc
MA-MRC ^[28]	2023	DBpedia+ 维基百科/1.3×10 ⁵	单段落型	人工合成/1.3×10 ⁵	抽取式,多答案类型	Acc
UnAnswGen ^[29]	2024	维基百科/8.6×10 ³	单段落型	人工合成/1.2×10 ⁵	抽取式,无答案类型	EM/F1
CMRC ^[30]	2018	维基百科/2×10 ⁴	单段落型	众包/2×10 ⁴	抽取式	EM/F1
BiPaR ^[31]	2019	中英平行小说段落/3.7×10 ³	单段落型	人工合成/3.7×10 ³	抽取式	EM/F1
GCRC ^[32]	2022	中文高考/5×10 ³	多段落型	中文高考/8.7×10 ³	多项选择式	Acc
X-STA ^[33]	2023	多语言维基百科/1×10 ⁴	单段落型	人工合成/1.5×10 ³	抽取式	EM/F1
ArQuAD ^[34]	2024	阿拉伯语维基百科/4×10 ³	单段落型	专家标注/1.6×10 ⁴	抽取式	EM/F1
Tamil-SQuAD ^[35]	2025	泰米尔语维基百科/5.3×10 ⁴	单段落型	人工合成/2×10 ⁴	抽取式	EM/F1
VEGA ^[36]	2024	arXiv 论文(文本+图表)/5×10 ⁴	多段落型	人工合成/2×10 ⁵	多项选择式	Acc
Illusory VQA ^[37]	2025	合成图像/1×10 ⁴	单段落型	人工合成/1.2×10 ⁴	多项选择式	EM/F1
M3-Bench ^[38]	2025	机器人视角视频+网络视频/1×10 ⁴	多段落型	人工合成/4.5×10 ³	自由答案式	Acc

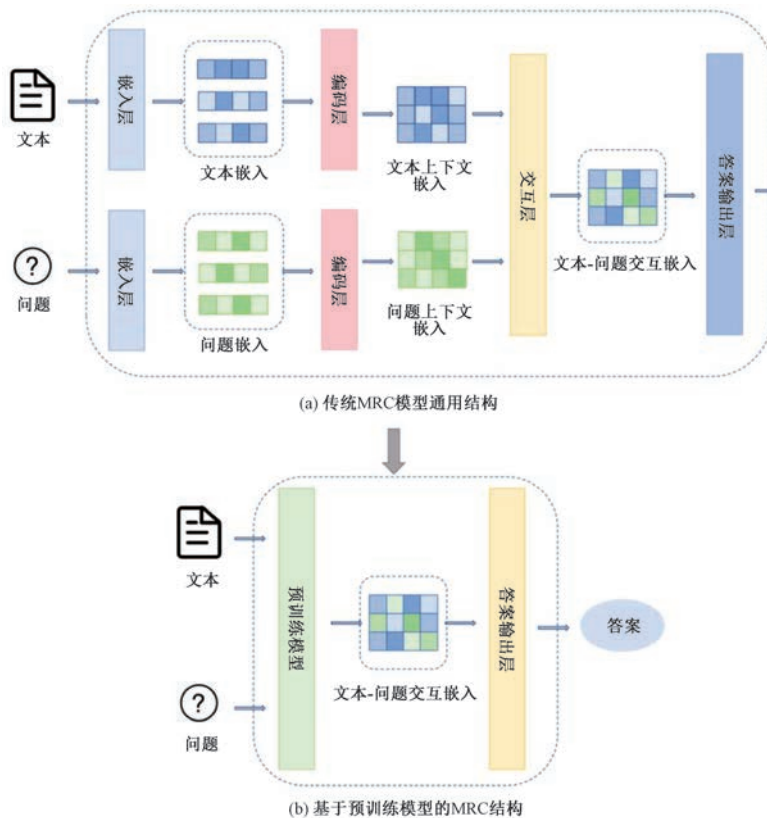


图 1 MRC 模型框架

$$\alpha_i = \text{softmax}(f(\mathbf{Q}, \mathbf{K}_i))$$

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \sum_{i=1}^n \alpha_i \mathbf{V}_i \quad (6)$$

其中, \mathbf{Q} 表示 query 的向量表示, $(\mathbf{K}_i, \mathbf{V}_i)$ 代表 key-value pairs 向量 (K, V) 表示的第 i 个值, 函数 f 常采用的计算方式有内积、二次型函数、前馈神经网络、双维度转换函数, 分别见如下公式:

$$f(\mathbf{p}_i, \mathbf{Q}) = \mathbf{p}_i^T \mathbf{Q} \quad \text{内积} \quad (7)$$

$$f(\mathbf{p}_i, \mathbf{Q}) = \mathbf{p}_i^T \mathbf{W} \mathbf{Q} \quad \text{二次型函数} \quad (8)$$

$$f(\mathbf{p}_i, \mathbf{Q}) = \mathbf{v}^T \tanh(\mathbf{W} \mathbf{p}_i + \mathbf{U} \mathbf{Q}) \quad \text{前馈神经网络} \quad (9)$$

$$f(\mathbf{p}_i, \mathbf{Q}) = \mathbf{p}_i^T \mathbf{W}^T \mathbf{U} \mathbf{Q} \quad \text{双维度转换函数} \quad (10)$$

答案输出层对应于步骤(4), 作用是从段落中查找出问题的答案。MRC 任务按照答案形式的不同大致分成四类, 因此这一层的设计需要考虑到答案形式。由于多项选择式任务的做法可以归结为填空式任务(将每一个选项看作是填空位置的候选项), 而填空式任务又是抽取式任务的特例, 这里主要介绍抽取式与自由答案式任务的输出层设计。对于抽取式阅读理解任务, Wang 等人^[48]提出了两种基于指针网络的输出模型, 第一种是序列式模型, 利用指针网络^[49](Pointer Network)以一种序列式的形式预测答案的每一个位置, 处理过程类似于 seq2seq 模型的解码过程。第二种是边界式模型, 利用指针网络仅仅预测答案的起始位置和终止位置。所预测答案的概率是预测这两个位置概率的乘积。抽取式模型的损失函数可以写为:

$$L(\theta) = -\frac{1}{N} \sum_{i=1}^N \log \mathbf{P}_{y_i^s}^S + \log \mathbf{P}_{y_i^e}^E \quad (11)$$

其中, θ 为模型参数, N 代表样本数目, y_i^s 表示第 i 个样本中标准答案的起始位置在文章中的位置, y_i^e 表示第 i 个样本中标准答案的终止位置在文章中的位置。对于自由答案式阅读理解任务, 典型的处理生成任务的架构有 seq2seq 模型, 将段落看作是 encoder 端的输入, decoder 端根据词汇表中的单词生成答案; 以及指针生成网络模型(Pointer-Generator Network, PGNet)^[50], 结合了 seq2seq 的生成机制以及指针网络的拷贝机制, 使得模型既能从词典中生成单词又能在原文中拷贝单词。

基于四层通用结构, 机器阅读理解模型的架构演进经历了显著的发展阶段, 从最初的基于注意力机制的模型, 逐步发展到以预训练语言模型为基础的架构, 最终聚焦于融入推理结构的 MRC 模型。其模型架构演进如图 2, 下面的小节由此展开论述。



图 2 MRC 模型架构发展

3.2 基于注意力的 MRC 模型

在神经阅读理解的四层结构中, 词嵌入层和编码层并不是 MRC 模型独有的, 其他 NLP 任务同样包含这两个层次。真正体现各个机器阅读理解模型特色的是交互层和答案输出层, 尤其是交互层中注意力机制的设计。由于当前大多数模型在交互层中使用的注意力机制较为复杂, 本节将根据注意力计算的方向和结构对各个模型进行分类。

早期研究聚焦于问题与上下文的交互注意力, 包括单向和双向注意力。单向注意力多为计算问题到段落(Q2C)注意力, 目的是突出段落中与问题最相似的部分, 如 Hermann 等人^[3]的 Attentive Reader 和 Impatient Reader 均是计算问题到文章(Q2C)的注意力, 且注意力的运算方式采用前馈神经网络(公式(9))。在此基础上, 利用双线性项(公式(10))取代原有的前馈神经网络计算方式, Kadlec 等人^[51]利用内积运算(公式(7))作为注意力计算方式。单向注意力所能交互的信息有限, 而双向注意力则可以融合 Q2C 和 C2Q 的双向交互, 达到两个方向的互补, 提供更加全面的交互信息, 典型的模型如 DCN^[52], BiDAF^[39]。

交互注意力过度依赖先验信息, 自匹配注意力机制可使文章中每个单词关注其余单词, 加深对文章理解。基于这一问题, 很多模型在交互注意力的基础上添加自匹配注意力机制, 如 R-Net^[53], 在 Match-LSTM^[48]的 C2Q 注意力的基础上添加一层自注意力。

预训练时代推动注意力机制高效化, 为使模型能够更准确地捕捉单词之间的复杂依赖关系并增强位置感知能力, DeBERTa^[54]提出解耦注意力, 使模型更灵活捕捉单词关系。针对长序列处理, Child 等^[55]提出分块注意力, 将输入序列分块稀疏计算。Longformer^[56]通过局部窗口与全局锚点策略处理长文本, 包括扩张滑动窗口注意力机制和全局注意力机制, Fin-EMRC 模型^[57]应用其处理金融文档中的长文本依赖问题。神经稀疏注意力 NSA^[58]进一步提出动态分层稀疏策略, 通过三个并行注意力分支处理输入序列, 既保证模型对全局上下文的感知, 又兼顾局部信息的精确性。

表 2 对比了本节介绍的经典的 MRC 模型在注意力机制设计上的差异。其中 Q2C 代表问题到段

落注意力, C2Q 代表段落到问题注意力, Bidirectional 代表双向注意力 (C2Q + Q2C), self-attention

代表对段落做自匹配注意力运算, one-hop 代表单跳结构, multi-hop 代表多跳结构。

表 2 基于注意力机制的模型对比

模型	词向量表示形式	建模方法	注意力方向	推理模式
Attentive Reader ^[53]	单词+上下文表示	GRU+Att	Q2C	one-hop
Impatient Reader ^[53]	单词+上下文表示	GRU+Att	Q2C	multi-hop
Stanford Reader ^[59]	单词+上下文表示	GRU+Att	Q2C	one-hop
AS Reader ^[51]	单词+上下文表示	GRU+Att	Q2C	one-hop
Match-LSTM ^[48]	单词+上下文表示	LSTM+Att	C2Q	multi-hop
GA Reader ^[60]	单词+上下文表示	GRU+Att	C2Q	multi-hop
BiDAF ^[39]	单词+上下文表示	GRU+Att	Bidirectional	one-hop
DCN ^[52]	单词+上下文表示	LSTM+Att	Bidirectional	one-hop
IA Reader ^[61]	单词+上下文表示	GRU+Att	Bidirectional	multi-hop
R-Net ^[53]	单词+字符+上下文表示	GRU+Att	C2Q+Self-Att	multi-hop
DIM Reader ^[62]	单词+字符+上下文表示	LSTM+Att	Bidirectional	multi-hop
RMR ^[63]	单词+字符+上下文+特征表示	LSTM+Att	Q2C+Self-Att	multi-hop

3.3 基于预训练语言模型的 MRC 模型

近年来, NLP 领域受到了广泛的关注, 以 Transformer 架构为基石的预训练模型大幅提升了语义理解能力, 进一步推进 MRC 发展。预训练方法源自迁移学习的概念: 首先在其他相关任务上预训练模型, 使得模型学习到一些知识, 然后在目标任务上做进一步优化, 实现模型所学知识的迁移^[64]。对于 NLP 领域来讲, 预训练过程就是在大量的文本数据上学习到通用的语言表示。针对机器阅读理解任务构建模型时, 仅需设计适配特定任务的输出模块, 并将其与预训练模型进行拼接便能实现理想的任务性能。传统 MRC 结构与基于预训练模型的 MRC 结构对比如图 1(b) 所示。从模型结构的角度看, 预训练模型相当于将传统模型通用结构的编码层和交互层融合在一起, 在编码的同时进行段落与问题的交互。

3.3.1 Transformer

鉴于目前几乎所有的预训练模型都采用 Transformer^[47] 结构或者其变体作为模型的特征提取器, 因此本节首先介绍 Transformer 结构。Transformer 是由 Vaswani 等人提出的一种用于机器翻译的序列到序列 (seq2seq) 结构。Encoder 端由六个相同的层堆叠而成, 每一层有两个子层, 第一个子层采用多头 (multi-head) 自注意力机制, 第二个子层采用前馈神经网络 (Feed-Forward Network, FFN) 构成。之所以用自注意力机制是因为它既可以通过结合位置嵌入捕获句子中每一个单词的全局依赖关系而不受距离影响, 又可以并行计算。对比公式 (6), 自注意力机制下 $Q=K=V$, Transformer 中采用的计算方式如下:

$$Attention(Q, K, V) = softmax\left(\frac{QK^T}{\sqrt{d_k}}\right)V \quad (12)$$

其中 d_k 代表张量维度。此外 Transformer 采用多头自注意力机制, 将 Q, K, V 三个张量线性映射为多份, 每一份之间做注意力的运算最后拼接。

$$head_i = Attention(QW_i^Q, KW_i^K, VW_i^V)$$

$$Multi-head(Q, K, V) =$$

$$concat(head_1, head_2, \dots, head_h)W^O \quad (13)$$

其中, h 代表头的数目, 是一个超参数, W_i^Q, W_i^K, W_i^V, W^O 都是训练参数。采用 multi-head 的目的是让模型联合关注序列中不同位置单词的不同表示子空间的信息, 可以类比于卷积神经网络中利用多个卷积核做特征提取, 目的同样是使得不同的卷积核关注不同特征。此外每一个子层都利用层归一化 (Layer Normalization^[65]) 和残差连接 (Residual Connection^[66]) 机制。

3.3.2 预训练语言模型

基于 Transformer^[47] 架构的预训练语言模型 (Pre-trained Language Models, PLMs) 经历了显著的迭代和发展。编码器-解码器架构中, 编码器将输入源文本编码为中间表示, 解码器将其解码为目标文本或其他输出, 如 T5^[67] 在翻译和摘要任务中表现优异。纯解码器架构包括因果解码器 (如 GPT 系列^[68-70]) 和前缀解码器 (也称为非因果解码器, 如 Chat-GLM^[71]), 二者分别在文本生成和补全任务中表现突出。

位置编码方法包括绝对位置编码 (公式 (14))、相对位置编码 (公式 (15))、旋转位置编码^[72] (RoPE, 公式 (16)), 被 Llama2/3^[73-74] 采用, 以及 ALiBi 位置编码^[75] (公式 (17)) 通过线性偏置项实

现位置感知。

$$\mathbf{x}_i = \mathbf{x}_i + \mathbf{p}_i \quad (14)$$

$$\mathbf{A}_{ij} = \mathbf{W}_q \mathbf{x}_i \mathbf{x}_j^T \mathbf{W}_k^T + \mathbf{r}_{i-j} \quad (15)$$

$$\mathbf{A}_{ij} = \mathbf{W}_q \mathbf{x}_i \mathbf{R}_{\theta, i-j} \mathbf{x}_j^T \mathbf{W}_k^T \quad (16)$$

$$\mathbf{A}_{ij} = \mathbf{W}_q \mathbf{x}_i \mathbf{x}_j^T \mathbf{W}_k^T - m(i-j) \quad (17)$$

其中, \mathbf{p}_i 为 i 处的位置嵌入向量, $\mathbf{R}_{\theta, i}$ 表示旋转角为 $i\theta$ 的旋转矩阵。

激活函数如 ReLU、GeLU^[76] (公式(18)) 和 SwiGLU^[77] (公式(19)) 用于增强模型非线性表达能力, 其中 SwiGLU 通过门控机制动态调节信息流, 在深层网络中可能表现更佳。

$$\text{GeLU}(\mathbf{x}) = 0.5 \otimes \left[1 + \text{erf} \left(\frac{\mathbf{x}}{\sqrt{2}} \right) \right],$$

$$\text{erf}(\mathbf{x}) = \frac{2}{\sqrt{\pi}} \int_0^{\mathbf{x}} e^{-t^2} dt \quad (18)$$

$$\text{Swish}(\mathbf{x}) = \mathbf{x} \otimes \text{sigmoid}(\mathbf{x}),$$

$$\text{SwiGLU}(\mathbf{x}_1, \mathbf{x}_2) = \text{Swish}(\mathbf{x}_1) \otimes \mathbf{x}_2 \quad (19)$$

其中, \otimes 表示逐元素乘法。

归一化方法包括 LayerNorm^[65] (公式(20)) 通过特征维度归一化稳定训练过程、RMSNorm^[78] (公式(21)) 去除均值计算, 较 LayerNorm 降低 20% 计算开销, 应用于 LLaMA 系列^[73-74, 79] 中, 以及 DeepNorm^[80] (公式(22)), 用于保证千层模型的训练稳定性, 应用于 GLM-130B^[81]。归一化方法在残差连接之前或之后均可执行。

$$\text{LayerNorm}(\mathbf{x}) = \frac{\mathbf{x} - \mu}{\sigma} \gamma + \beta,$$

$$\mu = \frac{1}{d} \sum_{i=1}^d \mathbf{x}_i,$$

$$\sigma = \sqrt{\frac{1}{d} \sum_{i=1}^d (\mathbf{x}_i - \mu)^2} \quad (20)$$

$$\text{RMSNorm}(\mathbf{x}) = \frac{\mathbf{x}}{\text{RMS}(\mathbf{x})},$$

$$\text{RMS}(\mathbf{x}) = \sqrt{\frac{1}{d} \sum_{i=1}^d \mathbf{x}_i^2} \quad (21)$$

$$\text{DeepNorm}(\mathbf{x}) = \text{LayerNorm}(a\mathbf{x} + \text{SubLayer}(\mathbf{x})) \quad (22)$$

训练范式方面, BERT^[82] 提出在预训练方式上采用的降噪自编码方式, 随机掩盖掉一些单词, 在输出层获得掩盖位置的概率分布, 让模型根据掩盖位置的上下文预测这个单词, 这种机制也叫掩码语言模型 (Masked Language Model, MLM) 或双向语言模型。MLM 的目标函数为

$$L(\Theta) = \sum_{i=1}^N (\log P(t_k | t_1, t_2, \dots, t_{k-1}, t_{k+1}, \dots, t_N)) \quad (23)$$

除了 MLM 任务外, 还利用下一句预测 (Next Sentence Prediction, NSP) 任务使得模型在诸如文本蕴含、问答这类需要判断两个句子关系的下游任务表现更好。BERT 的预训练流程属于多任务学习, 通过 MLM 和 NSP 两个任务从不同角度提高预训练模型的语义表达能力, 其中 MLM 任务用来学习句子中词与词之间的语义关联, 而 NSP 任务用来学习两个句子之间的逻辑关系。BERT 在 SQuAD^[4] 数据集上的效果超过了人类的水平, 在其他的 NLP 任务上也都有提升。

BERT 开启了 NLP 领域预训练模型的时代, 此后很多更加强大的预训练模型相继提出。这些模型大部分都是基于 BERT 的改进模型, 主要是针对 BERT 预训练阶段的两个任务 MLM 和 NSP 做改进, 如 RoBERTa^[83] 模型, 使用动态掩码替换 BERT 的静态掩码, 同时去除 NSP 任务, 在 160GB 数据与 500K 训练步长下, 发现单一 MLM 任务性能超越 BERT 的复合任务, 表明任务复杂度并非数据效率的必要条件; 由于 MLM 任务本质是双向语言模型, 需同时依赖上下文预测掩码 token, 这种设计在理解类任务, 如分类、问答中表现优异, 但在生成类任务中, 需要模型按从左到右的单向逻辑输出序列, BERT 的双向建模无法直接适配这种时序依赖, 因此, UNILM^[84] 扩展了 BERT 预训练的任务, 同步训练单向、双向及 Seq2Seq 语言建模范式, 并借助掩码机制针对性应对各类语言模型的约束挑战; ALBERT^[85] 模型则利用句子顺序预测 (Sentence Order Prediction, SOP) 任务改进了 BERT 的 NSP 任务, 以排序任务替代二分类, 使模型对句子间细粒度的语篇线索进行显式建模; ELECTRA^[86] 引入替换令牌检测任务 (RTD), 采用生成器 + 判别器的方式, 生成器首先生成新的 token, 判别器进一步区分生成的 token 和原始 token, 这种预训练方式在计算上更加高效, 可以涉及所有输入标记, 而不仅仅是被掩盖的那一小部分。除了在预训练阶段改进训练任务外, MT-DNN^[87] 模型在微调阶段引入了多任务学习机制, 使用多个任务微调模型参数使得模型具有更好的泛化性。

3.3.3 大模型技术进展

随着大模型 LLMs (Large Language Models) 技术的不断发展, 其应用向多语言、长文本处理、低

资源部署场景渗透,传统注意力机制在长上下文建模、推理效率、多场景适配性上的局限逐渐凸显,业界不断迭代出了多种维度的改进技术以提升模型性能与适用范围,本小节梳理大模型发展以来的技术改进。

一些工作尝试在注意力机制上做出适配性改进,多语言生成式大语言模型 LLaMA2^[73] 在 LLaMA^[79] 的基础上,将上下文长度扩展至 4096,并引入了分组查询注意力 GQA^[88] (Grouped Query Attention) 机制,有效降低了推理过程中的内存需求,提升了推理速度。LLaMA3^[74] 则进一步将 GQA 应用于小型模型,并采用更高效的分词器 TikToken,扩大了词汇表的数量,同时将上下文长度翻倍,并大幅增加了训练数据量。为进一步容纳更多示例,GPT-4 Turbo^[89] 将上下文窗口大幅扩展至 128k tokens,实现了超长序列的完整建模。

在提升模型输出准确性方面,示例学习方法展现出显著的效果与潜力。上下文学习 ICL^[90] 使用提示词方法,通过示例集 $D_k = \{f(x_1, y_1), f(x_2, y_2), \dots, f(x_k, y_k)\}$ 和任务描述 I 引导模型输出,其中 $f(x, y)$ 是由“输入-输出对”组成的样例到自然语言的映射函数。之后,为更好地应对复杂推理任务,Liu 等人^[91] 提出逻辑思维链指令微调方法 (Logical Chain-of-Thought Instruction Tuning, LogiCoT),针对不同的逻辑推理任务,设计多种类型的指令,旨在通过指令驱动的数据增强方法,提高模型在逻辑推理任务上的性能^[92]。在此基础上,为了确保这些推理不仅合理而且正确并与人类价值观对齐,基于人类反馈的强化学习 RLHF^[93] (Reinforcement Learning from Human Feedback) 进一步引入过程奖励模型,要求模型在生成最终答案前先产生详细推理轨迹,并通过在线强化学习(如 GRPO^[94] 算法)和偏好奖励模型对推理链中的每个逻辑步骤进行精细优化,显著提升了模型在数学推理和代码生成等复杂任务中的表现。

然而,大模型仍存在固有的知识局限性及高昂的迭代成本等局限,为解决这一问题,检索增强生成技术 RAG^[95] (Retrieval-Augmented Generation) 引入向量化检索机制,通过将外部知识库实时检索的信息作为上下文,使模型能动态获取外部知识库的最新信息,形成“检索-增强-生成”的闭环处理流程。RAG 聚焦外部知识补充,Function Call^[70] 则聚焦工具能力扩展,为更好地引入外部的函数及功能,函数调用(function call)技术通过将自然语言指令映射

为结构化 API 请求,使模型可以根据用户意图,自主决定调用外部工具(如计算器、API、数据库)来执行具体任务,实现了从认知理解到环境交互的关键跨越。

同时,为防止大量计算资源浪费,实现更加高效的推理,混合专家 MoE^[96] (Mixture-of-Experts) 模型使用门控网络来决定每个数据应该被哪个模型训练,从而减轻不同类型样本之间的干扰,避免了传统大模型全量计算的资源浪费。典型的混合专家指令微调的模型 Mistral 8x7B^[96] 采用 MoE 架构,通过八位“专家”和七十亿参数,结合稀疏门控机制,动态激活不同专家模块处理多样化指令,能够将数据高效地分配给各自擅长处理特定任务的神经网络部分。

3.3.4 迁移模型到 MRC 任务

预训练模型凭借其强大的文本表征能力,在机器阅读理解任务中表现出广泛的适用性,只需结合具体任务特点设计针对性的微调网络结构即可。迁移预训练模型的方式比较灵活,可以将预训练模型迁移在传统模型通用结构中除了答案预测层外的任意一层,而 3.1.4 节所介绍的两种输出层都可以直接拼接在预训练模型上。表 3 详细对比了本文介绍的所有预训练模型,其中 SBO 指跨度边界预测任务 (Span-Boundary Objective),EMD 指增强掩码解码器 (Enhanced Mask Decoder)。

3.4 基于推理结构的 MRC 模型

训练机器理解人类语言是人工智能的长期目标,具有广泛的应用场景,如问答和对话系统。尽管模型已在 CNN&Daily Mail、SQuAD、RACE 等数据集上取得了显著成果,进一步增强机器的复杂推理能力依然是值得深入探索的关键研究方向。根据现有研究^[59,104-105],推理能力可以大致分为多跳推理、逻辑推理、数值推理和常识推理^[106]。

3.4.1 多跳推理

根据段落与问题之间的交互,单跳结构是指段落与问题之间的交互仅仅计算一次。一种方法是将问题整体压缩为一个向量与段落计算一次注意力,如 Attentive Reader^[3]、AS Reader^[51] 等,另一种方法是问题与段落的整体表示采用并行化的计算方式,如 DCN^[52]、BiDAF^[39] 等。但单跳结构不能实现多步推理的效果,多跳推理要求模型能够跨越多个句子或段落,整合分散的信息来回答问题。这种推理类型模拟了人类在阅读复杂文本时的思维过程,需要模型具备良好的信息整合和逻辑推理能力。

表 3 预训练模型技术细节对比

模型名称	位置编码	自注意力机制	激活函数	归一化层	预训练范式	特点
BERT ^[82]	绝对位置编码	标准多头注意力	GeLU	LayerNorm	MLM+NSP	利用掩码语言模型 (MLM) 和下一句预测 (NSP) 共同作为预训练任务
RoBERTa ^[83]	绝对位置编码	标准多头注意力	GeLU	LayerNorm	MLM	采用动态掩码机制, 去除 NSP 任务
UNLM ^[84]	绝对位置编码	标准多头注意力	GeLU	LayerNorm	多任务联合训练	同时训练多种语言模型, 采用掩码机制解决不同语言模型的约束问题
ALBERT ^[85]	绝对位置编码	标准多头注意力	GeLU	LayerNorm	MLM+SOP	用句子顺序预测任务 (SOP) 取代下一个句子预测 (NSP)
ELECTRA ^[86]	绝对位置编码	标准多头注意力	ReLU	LayerNorm	RTD	生成器-判别器对抗训练
MT-DNN ^[87]	绝对位置编码	标准多头注意力	GeLU	LayerNorm	多任务联合训练	预训练过程与 BERT 一致, 微调阶段采用多任务学习
Longformer ^[56]	相对位置编码	稀疏注意力	GeLU	LayerNorm	字符级语言建模	适用于长文档的高效预训练模型, 采用稀疏注意力机制
DeBERTa ^[59]	绝对位置编码	解耦自注意力	GeLU	LayerNorm	MLM+EMD	采用解耦注意力机制和改进解码器结构提升性能
SpanBERT ^[97]	相对位置编码	标准多头注意力	GeLU	LayerNorm	MLM+SBO	基于跨度掩蔽策略预训练, 适合需要精确跨度预测的场景
T5 ^[67]	相对位置编码	标准多头注意力	ReLU	LayerNorm	文本到文本统一框架	任务通用框架, 所有任务均转换为输入-输出文本对
Gopher ^[98]	相对位置编码	标准多头注意力	GeLU	RMSNorm	自回归语言建模	DeepMind 专攻科学理解的 280B 参数模型
GLM-130B ^[81]	RoPE	混合注意力+掩码	GeLU	DeepNorm	自回归填空	高效预训练语言模型, 采用混合专家架构提升效率
LLaMA ^[79]	RoPE	标准多头注意力	SwiGLU	RMSNorm	多任务联合训练	高性能预训练模型, 具有多种参数规模, 适用于多种 NLP 任务
LLaMA2 ^[73]	RoPE	分组查询注意力	SwiGLU	RMSNorm	多任务联合训练	Meta 开源模型, 70B 参数接近 GPT-3.5
LLaMA3 ^[74]	RoPE	分组查询注意力	SwiGLU	RMSNorm	多模态混合训练	支持跨模态推理
GPT-4 ^[88]	RoPE	混合专家模型 MoE	GeLU	LayerNorm	多模态混合训练	16 专家架构实现多模态通用智能
PaLM ^[99]	RoPE	标准多头注意力	SwiGLU	LayerNorm	多任务联合训练	谷歌支持 100+ 语言的通用模型, 医疗领域超越人类专家
OpenAI o1 ^[100]	RoPE	混合专家模型 MoE	GeLU/GLU	LayerNorm/ RMSNorm	COT+RLHF+ 过程监管	通过强化学习训练, 在回答前进行链式思维推理, 具备更强复杂推理与安全策略遵循能力
DeepSeek-R1 ^[101]	RoPE	混合专家模型 MoE	SwiGLU	LayerNorm/ RMSNorm	长思维链监督微调 (Long-COT SFT)+ RLHF	首次尝试不使用 SFT, 仅使用强化学习 (RL) 提升模型性能的推理能力
Kim-K1.5 ^[102]	RoPE	混合专家模型 MoE	SwiGLU	LayerNorm/ RMSNorm	长思维链 SFT+COT+采样	128 KB 长上下文强化学习
Gemini 1.5 ^[103]	RoPE	稀疏混合专家模型 MoE	SwiGLU	LayerNorm/ RMSNorm	微调+ SFT+RLHF	多模态长上下文模型

图 3 是一个多跳推理过程的示意图, 数据来源于 Hotpot QA^[6] 数据集。

实现多步推理通常有以下几种方式: 第一种是基于注意力迭代的方法。可以利用 RNNs 这种基于上一时刻隐藏状态更新下一时刻隐藏状态的循环特性来达到多步推理, 如 Match-LSTM^[48], 类似的

模型如 R-Net^[53]、IA Reader^[61] 等。还可以基于之前时间步的问题感知段落表示, 计算下一时间步的段落和问题交互, 以序列式方式计算注意力^[3], 这种方式类似于人在阅读过程中不断地在问题和文章之间做关注。或者通过堆叠多个计算注意力的层数达到多步推理的目的, 如 Gated Attention Reader (GA

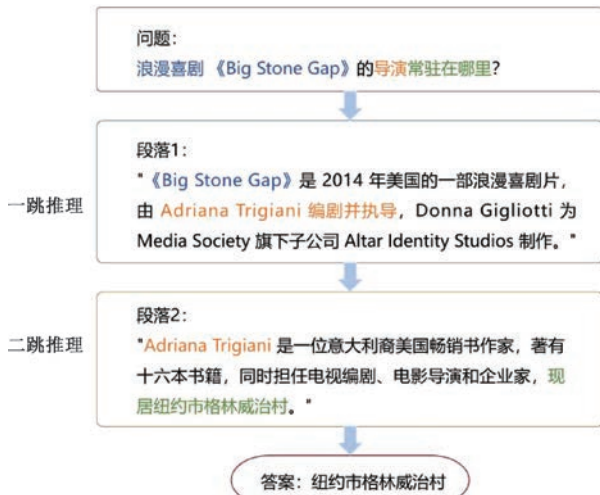


图3 多跳推理示意例图

Reader)^[60]模型,在每一步的推理过程中得到问题感知的段落表示输入下一层,这种处理过程类比于带着问题反复地阅读文章,每一次都加深对文章的语义理解。

在此基础上为进一步利用多跳信息,RMR(Reinforced Mnemonic Reader)^[63]提出一种重关注机制(reattention mechanism),通过直接利用之前层计算的注意力信息来微调当前层注意力分布的计算。基于双向注意力机制的模型在多跳过程中交替关注查询编码和文档编码,以发现文档、缺失查询词和查询之间的推理联系,如MHPGM^[107]和DIM Reader^[62]。模型S2G^[108]引入了句子感知自注意力(SaSA, Sentence-aware Self-Attention)机制,通过在句子开始前注入特殊标记,使模型能够明确地聚合句子内的所有词嵌入。同时,提出了证据引导注意力(EGA, Evidence Guided Attention)机制,迫使模型更多地关注从先前步骤中提取的证据。DECOMPRC^[109]提出基于问题分解的方法,将多跳问题分解为多个单跳子问题,通过答案选择策略,将多个子问题的答案整合起来,得到最终答案。与前面在推理过程中使用固定跳数或迭代的方法不同,ReasonNets^[110]、GAI^[111]引入了终止状态来放松对推理深度的限制。

第二种是基于图神经网络的方法,使用图神经网络通过消息传递机制传播多跳信息,节点的状态会根据其邻居节点的信息进行逐步更新和传播,从而捕捉到不同事实之间的多跳关系。为强化多跳任务中不同粒度信息的关联与推理能力,HDE^[112]构建了一个包含文档、候选答案、实体三种类型节点的异构文档实体 HDE (Heterogeneous Document-

Entity)图,使用图卷积网络(GCN)对HDE图进行编码,通过消息传递算法在图中传播信息,以更新节点表示。MHQA-GRN^[113]模型则聚焦实体关联的精细化建模,对实体引入同类型边、窗口类型边和共指类型边等多种边类型,构建出更复杂丰富的图结构,打破了传统仅依赖共指信息的局限。动态语义图AMR-SG^[114]通过构建抽象语义表示AMR(Abstract Meaning Representation)图,将问题、答案和相关事实的语义信息转化为图结构。然后,利用图卷积网络(GCN)对AMR图进行编码和推理,通过消息传递机制在图中传播信息,实现多跳信息的整合和推理。但GCN通过统一的矩阵运算处理所有边,难以差异化利用异构边的语义信息,针对这一局限,GAT可结合注意力机制为不同邻居分配不同权重,能更精准捕捉节点间重要关联,多跳编码融合网络MulQG^[115]使用图卷积网络(GCN)对实体感知答案编码器生成的动态实体图进行编码,通过双向注意力机制更新多跳答案编码。Zheng等^[116]利用图注意力网络对文档进行不同层次(词、句、段、文档)的表示学习,以捕捉文档的层次结构。多跳推理模型比较如表4所示。

表4 多跳推理模型比较

模型	知识表示形式	建模方法	多跳推理方法
S2G ^[108]	实体	Att	选择引导策略,由粗到细逐步提取
ReasonNets ^[110]	实体	Att	引入终止状态,放松对推理深度的限制
MHPGM ^[107]	事实单元	Att	通过多跳自注意力机制逐步推理
HDE ^[112]	实体	图模型	构建异构文档-实体(HDE)图
MHQA-GRN ^[113]	实体	图模型	基于实体提及和三种类型边构建图
MulQG ^[115]	实体	图模型	通过GCN在动态实体图上传播信息
DECOMPRC ^[109]	语篇单元	图模型	将多跳问题分解为多个子问题,合并答案
Zheng等 ^[116]	实体、短语、语篇单元	图模型	选择性门控注意力机制+ConceptNet常识知识
GAI ^[111]	事实单元	图模型	动态终止的多跳推理机制
AMR-SG ^[114]	事实单元	图模型	通过GCN在动态语义图(AMR)上传播信息

3.4.2 逻辑推理

逻辑推理是指模型能够根据给定的逻辑规则和事实,进行演绎、归纳或溯因推理,以得出正确的结论。这类推理要求模型具备严谨的逻辑思维能力,

能够理解文本中的因果关系、条件关系等,相关数据集如 LogiQA^[22] 和 ReClor^[23],如图 4 是一个基于条件关系的逻辑推理过程的示意图,数据来源于 LogiQA^[22] 数据集。

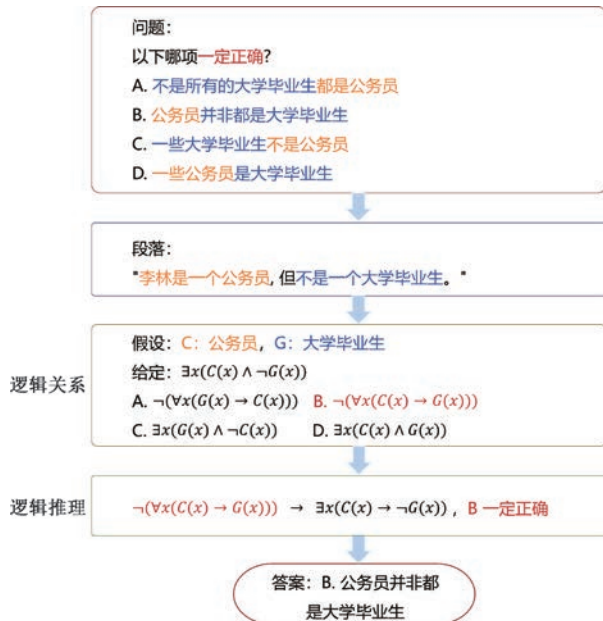


图 4 逻辑推理示意例图

根据模型架构的不同,主要分为以下四类:基于符号神经网络的方法、基于图神经网络的方法、基于预训练的方法和基于大模型微调的方法^[92]。基于符号神经网络的方法会先设定逻辑规则和拓展规则,再借助神经网络执行推理流程以获取对应结论。逻辑驱动推理器 LReasoner^[117] 使用逻辑推理的文本扩充框架,通过提取逻辑表达式,根据预定义的逻辑等价法则进行符号化推断并扩展文本以匹配答案。AdaLoGN^[118] 模型利用 Graphene^[119] 获取基本语篇单元(Elementary Discourse Units, EDUs^[120]),得到蕴含、否定等六种逻辑关系构建文本逻辑图,通过推理规则扩展逻辑图,并结合图神经网络消息传递过程,实现神经推理和符号推理的互相增强。

基于图神经网络的方法主要由图构建和图推理两部分组成。基于语篇结构的图网络 DAGN^[121] 参考 PDTB 2.0^[122] 区分语篇关系,构建语篇图并利用多跳注意力机制捕捉节点间长距离依赖。整体图网络 HGN(Holistic Graph Network)^[123] 模型构建包含元素语篇单元(EDU)和关键短语(Key Phrase, KPH)的图结构,通过分层交互机制,通过融合类型层面注意力与节点层面注意力,实现对图网络表示的优化更新;具体而言,类型层面注意力作用于不同类型的相邻节点,节点层面注意力则适用于同类型

相邻节点间。Logiformer^[124] 从文本中提取逻辑单元,构建逻辑图和语法图,分别建模文本中的因果关系及共现关系,利用图变换器结构捕捉逻辑单元间的长距离依赖关系。事实驱动的逻辑推理模型 Focal Reasoner^[106] 通过 Spacy^[125] 技术从文本中提取主谓宾语构成事实单元(Fact Unit),通过 levi 图^[126] 形式进一步构建超图,使用关系图卷积网络对超图前向传播,使用注意力机制融合选项与问题的交互信息。受 TransE^[127] 启发,尾论元的嵌入应该接近于头论元的嵌入加上关系相关向量在隐藏表示空间中的表示,即 $v_{\text{subject}} + v_{\text{predicate}} \rightarrow v_{\text{object}}$,应用逻辑事实正则化建模事实单元之间的显式关系,从而增强逻辑关系,模型联合训练选项与问题的交互损失和事实单元间的损失。

基于预训练的方法,通过使模型从大规模文本中学习逻辑模式或设计针对逻辑推理的额外预训练任务增强逻辑推理能力。MERIT^[128] 模型提出一种元路径引导的逻辑推理对比学习模型,将大规模无标注文本构建成实体级别的知识图,其训练样本采用元路径策略和反事实数据增强构造,之后采用对比学习机制与掩码语言建模辅助训练目标。APOLLO^[129] 模型通过修改掩码语言建模损失和句子分类损失提升语言建模的逻辑推理效能。Pi 等人^[130] 提出的对抗性预训练学习逻辑推理模型 LogiGAN 依托逻辑指示词,在大规模文本语料库中自动识别逻辑推理现象,进而驱动语言模型完成被掩码逻辑语句的预测任务。IDOL^[131] 则通过指示词进行逻辑预训练,结合经典的 MLM 任务和逻辑类别预测任务,使模型学习分析文本中的逻辑结构和推理过程。

近年来,大模型的崛起为 MRC 中的逻辑推理问题提供了强大工具。逻辑思维链指令微调方法 LogiCoT^[91] 通过在由强大模型(如 GPT-4)生成的高质量思维链(CoT)推理数据上进行监督微调(SFT),来增强 LLMs 的推理能力^[132-135]。然而,监督微调方法仅关注正向推理路径(即那些能得出正确答案的推理路径),却在很大程度上忽视了负向推理路径。因此,在此基础上,结合强化学习机制,尤其是基于群体相对策略优化 GRPO^[94] (Group Relative Policy Optimization)的方法,模型不仅能在每一步获得细粒度的过程奖励反馈,还能通过试错机制在无需人类标注的前提下,自主进化其推理策略,典型的应用如 DeepSeek-R1^[101] 和 Kimi-K1.5^[102]。逻辑推理模型比较如表 5 所示。

表 5 逻辑推理模型比较

模型	知识表示形式	建模方法	逻辑推理方法
LReasoner ^[117]	实体	Att	提取逻辑表达式,使用神经网络推理
AdaLoGN ^[118]	语篇单元	图模型	通过推理规则扩展逻辑图
DAGN ^[121]	语篇单元	图模型	构建注意力图+多轮推理
HGN ^[123]	语篇单元、短语	图模型	类型级别注意力+节点表示注意力更新图
Logiformer ^[124]	逻辑单元	图模型	分别构建逻辑图和分类图
Focal Reasoner ^[106]	事实单元	图模型	提取事实单元构建图
MERIt ^[128]	实体	预训练模型	MLM+元路径引导的逻辑推理对比学习
APOLLO ^[129]	实体	预训练模型	MLM+句子分类损失
LogiGAN ^[130]	实体	预训练模型	通过逻辑指示词预测被遮盖的逻辑语句
IDOL ^[131]	实体	预训练模型	MLM+逻辑类别预测任务
LogiCoT ^[82]	实体	预训练模型+微调	逻辑思维链指令微调

3.4.3 数值推理

数值推理要求模型对文本中的数值信息进行理解和计算,以得出正确的答案。这类推理在处理包含数据、统计信息或数学问题的文本时尤为重要。图 5 是一个数值推理过程的示意图,数据来源于 DROP^[20]数据集。Dua 等人^[20]构建了 DROP 数据集,专门用于评估模型在数值推理方面的表现,例如对数值的准确提取、单位的正确理解,以及复杂的数学运算。同样用于数值推理的还有 Chen 等构建的 RACENum^[24]数据集。其处理方式分为以下几种:

(1)基于符号神经网络的方法。通过定义数值操作规则和扩展规则,使用神经网络推理得出答案。NeRd^[136]模型使用一个阅读器编码段落和问题,再通过执行器生成可执行的程序,这些程序可以直接在段落上执行以产生答案。模型通过定义域特定语言 DSL(Domain Specific Language),将数值推理操作(如加减法、计数、排序)与文本中的信息结合起来。例如,通过 PASSAGE_SPAN 和 QUESTION_SPAN 操作符来选择段落或问题中的对应的文本,然后应用 COUNT、SUM、DIFF 等操作符来进行数值推理。OPERA^[137]模型定义了一组操作(如 ADDITION、DIFF、MAX、MIN、ARGMAX 等),这些

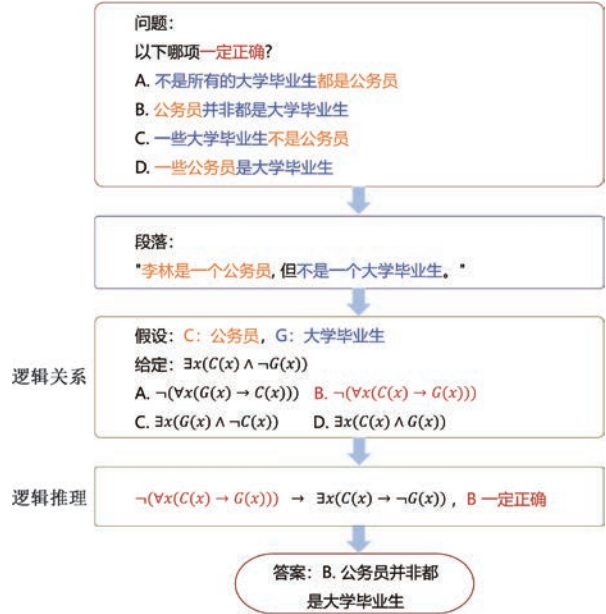


图 5 数值推理示意图

操作覆盖了数据集中所有问题的推理需求。模型通过操作选择器自动选择与问题相关的操作,然后通过操作执行器在给定的上下文中执行这些操作。操作选择器使用双线性函数计算每个操作与问题的相似度,操作执行器则基于多头交叉注意力机制实现,对选定的操作进行执行。模型对比示意图如图 6。

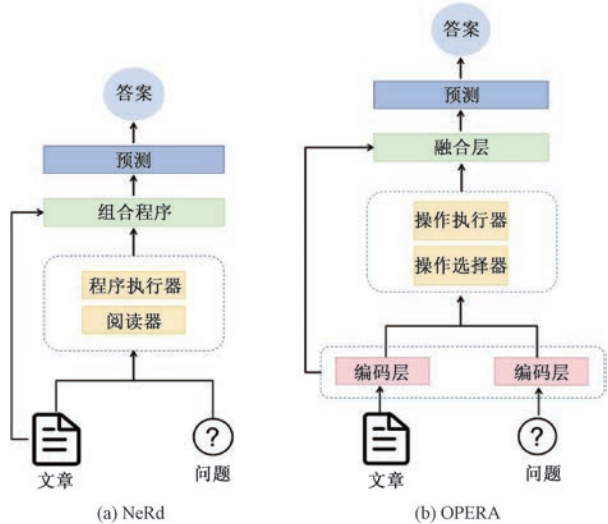


图 6 基于符号神经网络的数值推理模型

(2)基于图神经网络的方法。构建包含不同类型的数值和相关实体的异构图,编码它们之间的比较信息(如大于、小于或等于)和数值与实体之间的共现关系,通过图神经网络进行信息传播,实现数值推理。NumNet^[138]模型采用问题导向的图神经网络在推理过程中利用问题嵌入来引导图上的信息传播,动态地选择与问题相关的数值和实体,进行多步

推理,从而更好地处理数值推理问题。为了考虑与数字相关的实体信息,QDGAT^[24]模型构建了一个包含数字和额外实体的异构图,然后利用问题引导的图注意力网络来增强其数值推理能力。尽管这些基于GNN的模型表现出色,但它们的GNN模块将构建的图视为同质图,仅聚合直接邻居节点的信息,限制了对丰富异质邻域信息和节点类型的显式捕捉。为了明确整合丰富的异构邻域信息,HNNT-GAT^[139]模型采用带有重启的随机游走为每个节点采样异构邻居,并按不同节点类型分组。为了整合节点类型信息,在异构图的节点初始化阶段,进一步明确地将原始节点向量与节点类型嵌入连接起来。使用BiLSTM和注意力机制更新每个节点的向量,利用其丰富的异构邻域信息,并利用问题引导的图注意力网络进行数值推理。

(3)基于预训练的方法。通过改进预训练任务,显式增强语言模型对数字的敏感性和运算能力。NAQANet^[20]模型在预训练模型(如BERT)的基础上增加数值推理模块,模型首先使用BERT等预训练模型对段落和问题进行编码,然后通过特定的输出层来预测数值答案、计数或加减法等操作。例如,对于数值答案,模型会预测是否是计数或算术表达式,并进一步预测具体的数值。GenBERT^[140]模型首先在大规模的数值数据(如简单的数值计算表达式)和文本数据(如生成的需要数值推理的问题-段落对)上进行预训练。然后,在微调阶段,模型可以直接应用于数值推理任务,如DROP数据集。NC-BERT^[141]引入了数值上下文注意力掩码(NC-Mask)来引导模型利用数值相关的上下文知识,以减轻语言模型在数值推理过程中对参数化知识的过度依赖,并通过数值上下文知识增强模型的数值推理能力。POET^[142]通过让语言模型模仿程序执行器的行为,来增强其推理能力。具体来说,POET在预训练阶段使用程序及其执行结果作为训练数据,让模型学习程序执行的推理知识。POET提出了三种不同的实例化方法,包括POET-Math、POET-Logic和POET-SQL。这些方法分别针对不同的推理任务进行预训练,使模型能够学习到数值推理、逻辑推理和多跳推理等技能,通过不同的程序执行器来增强模型的多种推理能力。

(4)基于微调的方法则是利用预训练阶段学习的通用知识,针对MRC具体任务进行微调。如MTMSN^[143]基于预训练的BERT模型进行微调,在预训练的BERT模型基础上,添加特定的任务

层,包括多类型答案预测器和多跨度提取相关层。多类型答案预测器支持多种数值答案类型,如跨度、计数、否定和算术表达式等。对于算术表达式,模型为段落中的每个数字分配一个符号形成可执行的算术表达式,包括加、减或零。此外,还提出了算术表达式重排机制,通过束搜索解码多个表达式候选,并根据上下文信息重新排序,以进一步确认预测结果。通过在DROP数据集上进行监督训练,优化整个模型的参数,使其适应数值推理的MRC任务。

随着大模型技术发展,采用Function Call调用外部计算器的工具增强式推理范式已成为重要技术方向,其核心是通过结构化指令调用外部工具弥补模型原生计算能力的不足,目前已涌现出一系列具有代表性的模型。主流模型中,最早由OpenAI在GPT-4/GPT-3.5中推出,支持Function Call接口调用外部计算器,能自动解析复杂数学表达式。目前阿里Qwen^[144]、深度求索DeepSeek^[66]等国内模型,均已支持Function Call,可通过调用数学计算函数、实时数据查询函数等,适配中文场景下的数值推理及实时任务需求。数值推理模型比较如表6所示。

表6 数值推理模型比较

模型	知识表示形式	建模方法	数值推理方法
NeRd ^[136]	实体、数值	符号神经网络	使用领域特定语言DSL定义操作
OPERA ^[137]	实体、数值	符号神经网络	设计符号操作作为神经模块
NumNet ^[138]	实体、数值	图模型	构建数值感知图神经网络
QDGAT ^[24]	实体、数值	图模型	利用问题表示指导图上的信息传播
HNNT-GAT ^[139]	实体、数值	图模型	利用问题引导的图注意力网络进行数值推理
NAQANet ^[20]	实体、数值	预训练模型	输出层预测答案类型
GenBERT ^[140]	实体、数值	预训练模型	自动数据生成+多任务训练
NC-BERT ^[141]	实体、数值	预训练模型	引入数值上下文注意力掩码 NC-Mask 预训练语言模型
POET ^[142]	实体、数值	预训练模型	通过程序及其执行结果预训练语言模型
MTMSN ^[143]	实体、数值	预训练模型+微调	预训练模型+特定任务层微调

3.4.4 常识推理

常识推理是指模型能够理解并运用人类日常生活中的基本知识和常识来回答问题,作为机器理解人类语言的关键挑战之一,这类推理要求模型不仅具备语言理解能力,还能模拟人类的常识性思

维^[145]。图 7 是一个常识推理过程的示意图,数据来源于 CommonsenseQA^[25] 数据集。CommonsenseQA^[25]、Cosmos QA^[26] 和 ART^[27] 等数据集推动了常识推理在机器阅读理解中的研究。此外,为更好地评估模型的常识推理能力,oLMpics^[146] 提出了八种推理任务,包括比较、连接和组合等,引入了一种新的评估协议,通过零样本学习和学习曲线分析,来理解模型在不同任务上的表现。外部知识库包括 ConceptNet^[147]、WordNet^[148] 和 NELL^[149] 等,均以(主体,关系,客体)三元组的形式存储,每个三元组表示两个实体之间的特定关系,如图 8 为 Concept-Net 知识图谱示意图。

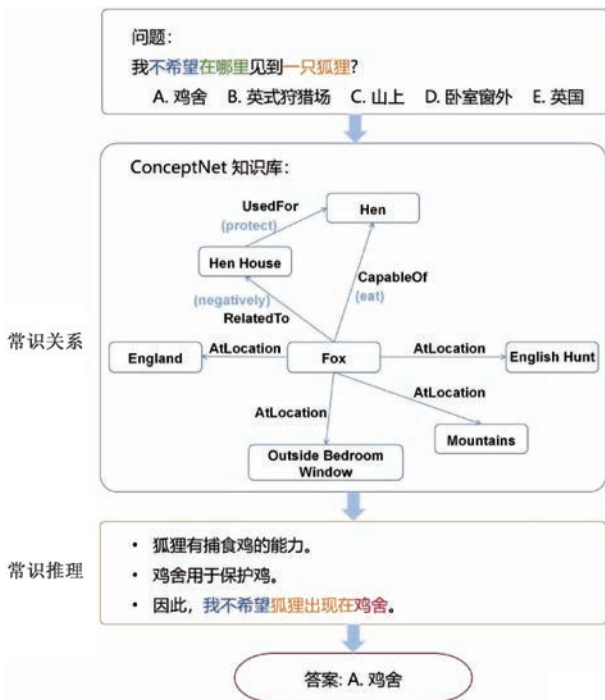


图 7 常识推理示意图

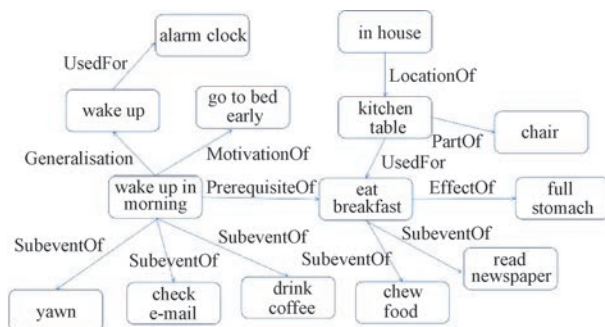


图 8 ConceptNet 知识图谱示意图

如何将外部知识与上下文信息进行融合是常识推理的重要问题。Chen 等人^[150] 通过结合 ConceptNet^[147] 常识知识图谱中的常识知识,增强模型对文本中实体关系的理解和推理能力,结合词性、命

名实体、关系等多方面特征对文本进行编码和交互建模,提出一种基于卷积神经网络和门控机制的模型,通过输入层、门控卷积层和输出层的结构,利用门控机制控制信息流动,最后通过双线性交互函数连接各层表达以产生最终答案。RekNet^[151] 从外部知识图谱(如 ConceptNet^[147])中检索与文本关键信息和候选答案相关的知识四元组(主体,关系,客体,置信度),使用图卷积网络(GCN)和长短期记忆网络(LSTM)对检索到的知识四元组进行编码,通过注意力机制将知识四元组与上下文信息进行融合。KagNet^[152] 通过在外部知识图谱中找到问题和答案概念之间的路径形成模式图,即外部知识图谱的相关子图,将外部知识图谱的问题和答案对从语义空间映射到知识表示的符号空间。KCF-NET^[153] 采用基于多头注意力机制的融合结构平衡知识和上下文的权重,从而在阅读理解过程中更有效地进行常识推理。

由于预训练模型参数量极大,注入新知识时重新训练参数会带来损耗和灾难性遗忘,一些模型尝试将外部知识直接引入预训练模型,而无需从头开始预训练。LIBERT 模型^[154] 通过在预训练过程中融入词汇级别的语义相似性来增强模型的语义理解能力。通过多任务学习框架,结合了 BERT 的语言模型目标和一个新的任务,即词汇语义分类,以提升模型在语义相似性任务上的性能。E-BERT 模型^[155] 则通过将 Wikipedia2Vec^[156] 实体向量与 BERT 的词向量空间对齐,然后将这些对齐后的实体向量作为输入,无需额外的预训练,就能有效地将事实性知识注入 BERT 模型中,从而在常识推理的 MRC 任务中取得更好的效果。常识推理模型比较如表 7 所示。

表 7 常识推理模型比较

模型	知识表示形式	建模方法	常识推理方法
Chen ^[150]	事实单元	CNN+GRU	引入外部常识增强文本表示
RekNet ^[151]	事实单元	图模型	引入外部常识四元组增强文本表示
KagNet ^[152]	事实单元	图模型	构建模式图表示常识知识
KCF-NET ^[153]	事实单元	预训练模型	采用基于多头注意力机制的融合结构平衡知识和上下文的权重
LIBERT ^[154]	事实单元	预训练模型	在预训练过程中融入词汇语义相似性知识
E-BERT ^[155]	事实单元	预训练模型	在预训练的语言模型中注入事实知识

由于大模型知识限于训练数据,而主流模型训练集多为公开网络数据,缺乏实时、非公开或私域信息,且存在幻觉问题,因此,引入向量化检索机制以降低训练成本的检索增强生成技术(RAG),已成为突破这些局限的主流趋势。当前,GPT-4 Turbo^[89]、Gemini 1.5^[103]等前沿大模型,均通过动态检索机制实现了能力升级。面对复杂专业问题时,它们能从学术数据库、专业知识库等多源渠道检索信息,并融合自身预训练知识生成答案。在传统 RAG 基础上,GraphRAG^[157]进一步实现技术迭代。它通过构建知识图谱来增强检索过程,利用图谱的结构化信息快速定位相关知识节点,最终在答案的全面性与多样性上,显著优于传统 RAG 方法。

4 机器阅读理解任务挑战及技术

在机器阅读理解的的实际应用中,真实场景往往面临更复杂的任务挑战。传统抽取式数据集(如 SQuAD 等)仅要求从单一段落提取答案,难以满足现实需求。为考察模型深层阅读理解与推理能力,本章将重点探讨六个机器阅读理解领域面临的主要任务挑战,并分析相应的技术解决方案,包括无答案问题的阅读理解、多答案问题的阅读理解、多段落问题的阅读理解、对话型阅读理解、多轮交互型阅读理解以及零样本和少样本任务上的阅读理解。

4.1 无答案问题的阅读理解

早期 MRC 数据集默认每个问题都能在文本中找到答案,但实际一段文本所包含的知识是有限的,因此存在两类无答案问题:(1)与文本内容无关的问题;(2)与文本内容类似但含义不同的问题。Moradisani 等人^[29]构建 UnAnswGen 数据集,通过语言学变换生成候选无答案问题并筛选。对于这类任务,模型需分为答案抽取模块和判别无答案问题模块,前者预测答案位置,第 3 节所介绍的 MRC 模型大部分都可以作为答案抽取模块,后者判断问题是否可回答。

Clark 等人^[158]尝试在原有的答案抽取模块的基础上额外添加一个专门用来预测无答案情况的网络层,此时损失函数定义如下:

$$L_{joint} = -\log\left(\frac{(1-\delta)e^z + \delta e^{\alpha_a \beta_b}}{e^z + \sum_{i=1}^{l_p} \sum_{j=1}^{l_p} e^{\alpha_i \beta_j}}\right) \quad (24)$$

其中, z 表示模型预测该问题是不可回答问题的分数,如果问题是可以回答的那么 $\delta=1$,反之 $\delta=0$ 。 α

和 β 分别表示输出层预测的文章中每一个单词作为答案起始位置和终止位置的概率, a 和 b 分别代表标准答案在文章中的起始位置和终止位置。预测的答案跨度分数 α_a 、 β_b 和判断无答案问题的分数 z 是共同归一化的。

Hu 等人^[159]认为两个分数共同归一化会出现冲突,如果模型过分信任预测的答案跨度分数,那么就会在预测无答案问题时产生较低分数。此外之前的模型并没有验证答案抽取模块预测的答案跨度的合理性。为了解决以上问题,他们提出 Read + Verify 架构。其中 Read 模块包括答案抽取模块和判别无答案问题模块,Verify 模块用来进一步验证答案抽取模块预测的答案跨度所在的句子(原文中称为 answer sentence)是否就是标准答案所在的句子。为了解决上面提到的冲突问题,在 Read 模块中额外增加了两个辅助损失函数:

$$L_{indep-span} = -\log\left(\frac{e^{\tilde{\alpha}_a \tilde{\beta}_b}}{\sum_{i=1}^{l_p} \sum_{j=1}^{l_p} e^{\tilde{\alpha}_i \tilde{\beta}_j}}\right) \quad (25)$$

$$L_{indep-unknown} = -(1-\delta)\log\sigma(\mathbf{z}) - \sigma\log(1-\delta(\mathbf{z})) \quad (26)$$

其中, $L_{indep-span}$ 代表答案抽取模块的损失函数,而此时的答案抽取模块是独立的预测答案片段而不考虑问题是否可以回答, $\tilde{\alpha}_a$ 和 $\tilde{\beta}_b$ 表示的就是答案抽取模块所预测出来的答案跨度。 $L_{indep-unknown}$ 代表判断问题无答案的损失函数,同样它是独立于答案抽取模块的。 σ 代表 sigmoid 函数。最后整个 Read 模块的损失函数定义为:

$$L_{Read} = L_{joint} + \gamma L_{indep-span} + \lambda L_{indep-unknown} \quad (27)$$

γ 和 λ 是两个超参数。实验表明去掉 $L_{indep-unknown}$ 后模型在判断无答案问题上的准确率显著下降,证明了上述提出的冲突确实存在。对于 Verify 模块,他们采用三种结构。第一种将预测出来的答案片段连同问题以及 answer sentence 连接成一个句子送入预训练模型 GPT^[68]中预测无答案的概率。第二种采用交互式结构,通过注意力机制计算它们之间的关联。第三种结构是前两个结构的结合,具体的就是将前两个结构的输出张量拼接,实验证明这种混合结构使得模型效果更好。

然而,上面这些方法并没有明确指出在不可回答的情况下,问题与候选答案之间不匹配的具体位置,因此容易选择一个看似合理但实际上错误的答案。NeurQuRI^[160]受使用成分词作为检查表中条件的想法的启发^[161],通过提取问题中的条件列表

并检查候选答案是否满足这些条件来改进不匹配问题。将问题拆解为需满足的条件集合(如“时间”和“地点”),通过循环单元累积注意力计算候选答案与每个条件的匹配程度,生成条件满足度向量,并与问题自身作为伪答案的满足度对比,模型利用交叉熵损失和条件满足度损失联合训练。满足度分数损失:

$$L_a = -\gamma \log(\min(\mathbf{a}_T^{q \rightarrow q})) - \phi_d \log(1 - \min(\mathbf{a}_T^{x \rightarrow q})) - (1 - \phi_d) \log(\min(\mathbf{a}_T^{x \rightarrow q})) \quad (28)$$

该损失旨在使候选答案在不可回答情况下至少有一个条件未满足,即在“候选答案到问题”满足度分数向量 $\mathbf{a}_T^{x \rightarrow q}$ 中至少有一个低值,但在其他情况下所有值都高。此外,所有“问题到问题”满足度分数向量 $\mathbf{a}_T^{q \rightarrow q}$ 中的分数都被强制为高值,因为问题本身应该满足所有条件,将此添加到损失中,权重为 γ 。通过这种方式,NeurQuRI 模型能够更准确地识别出那些看似合理但实际上不满足问题所有条件的候选答案,从而提高对无答案问题的识别能力。关于处理无答案问题阅读理解任务的其他相关模型可以参考文献[162-165]。

4.2 多答案问题的阅读理解

在多选答案阅读理解任务中,一个问题 Q 会对应着多个相关的段落 (D_1, D_2, \dots, D_K) 。模型的目标是从这 K 个段落中寻找所有可能的最佳答案 A_1, A_2, \dots, A_M ,并为每个答案计算其如下条件概率:

$$P(A_m | D_1, D_2, \dots, D_K, Q) \quad (29)$$

其中, M 是可能答案的数量。与单答案阅读理解任务不同,多选答案任务允许多个正确答案存在。模型需要能够识别并区分这些答案的相关性和正确性,同时处理段落之间的信息冗余和互补性,典型数据集如 MA-MRC^[28]。

传统的指针网络使用 softmax 函数来预测序列的开始和结束位置,这种方法在多分类任务中效果很好,但只能提取单个答案。Wang 等^[166]通过改进的指针网络(pointer network)来处理多答案和零答案的情况。使用两个相同的二元分类器来预测答案的开始和结束位置,而不是使用 softmax 进行多分类。每个标记被赋予一个二元标签(0 或 1),1 表示该标记是答案的开始或结束位置,0 表示不是。对于每个编码器输出的 token h_j ,计算其作为开始位置或结束位置的得分,然后使用 sigmoid 函数将得分转换为概率值:

$$p_{start,j} = \sigma(\alpha_{start}^T \tanh(W_{start} h_j) + \beta_{start}) \quad (30)$$

$$p_{end,j} = \sigma(\alpha_{end}^T \tanh(W_{end} h_j) + \beta_{end}) \quad (31)$$

其中, α_{start} 、 α_{end} 和 W_{start} 、 W_{end} 是学习矩阵, β_{start} 和 β_{end} 是偏置项,用于计算开始和结束位置的得分。为每个 token 分配一个二进制标签 0/1,基于邻近原则,标记为 1 的两个最近的标记形成答案,由此模型就能够更灵活地处理多答案和零答案的情况。为了处理多个答案的情况,Liu 等^[167]在指针网络中采用了最可能对匹配策略做改进,对开始和结束位置进行排序,选择概率最高的开始和结束位置对,确保答案片段不重叠。

Zhao 等^[168]通过 BIOES 标签(BIOES boundary labels)体系,将实体边界识别转化为序列标注问题,能够更细粒度地捕捉答案存在性。对于每个输入的隐藏状态 h_j 预测 BIOES 边界标签,可以将候选 BIOES 标签的概率计算为

$$p(\hat{l}) = \text{softmax}(W h_j + b) \quad (32)$$

其中, W 、 b 是学到的参数, \hat{l} 表示预测的边界标签,由此即可通过从标签序列中识别边界来提取候选实体。

4.3 多段落问题的阅读理解

单段落阅读理解仅要求模型从单一文本段落中寻找答案,对深层理解能力的考查较为有限。而真实场景下的阅读理解往往涉及多段落信息处理,需要从多篇文档中检索相关内容,通过跨段落信息比对与整合得出最准确的答案。多段落阅读理解任务中一个问题 Q 会对应着多个相关的段落 D_1, D_2, \dots, D_K ,模型需要从这 K 个段落中寻找最佳答案 A ,通过建模如下概率来实现:

$$P(A | D_1, D_2, \dots, D_K, Q) \quad (33)$$

多段落阅读理解也可以认为是开放领域(open-domain)问答的一种形式。Open-domain 问答目的是从广泛的领域资源(如维基百科,网页搜索等)寻找问题的答案而不仅仅在某段文本中,这更贴近于真实场景但同时具有相当大的难度,其核心挑战在于处理跨段落信息冗余、长距离推理逻辑链构建以及文本-表格等多模态数据融合。现有的长文本机器阅读理解模型可归纳为长文本模型、滑动窗口模型和粗到精模型三种方法^[169]。

长文本模型 Longformer^[56]和 Big Bird^[170]等提出采用稀疏自注意力机制替代传统全连接自注意力,通过选择性地关注部分 token 对(如通过局部窗口、固定步长位置或重要性得分聚焦部分 token),显著提高了长文本 MRC 的性能,但依赖先验知识

手动设计注意力规则的方式难以捕捉文本中动态变化的语义依赖关系。

为解决这一问题,滑动窗口模型将文档分割成小块,分别预测局部答案,再比较得分选择最高者^[5,12,171]。每个单独的块被放入一个读者模型中以预测局部答案,由来自多个组块的局部答案进行集成得到全局答案。Tan 等人^[172]提出 S-Net 模型,先通过片段抽取模块提取出一段文本作为答案的预测依据,然后利用生成模块生成答案。其中片段抽取模块采用多任务学习策略,答案生成模块采用 seq2seq 模型,其中 encoder 端的输入是问题和段落的向量表示,同时将片段抽取模块的输出作为额外的特征拼接到段落中。实验证明 S-Net 在 MS MARCO 数据集上的效果显著优于 R-Net^[53],ReasonNet^[110]这些单独做片段抽取任务的模型。然而,这种方法在模型需要综合文档多段落内容进行推理时存在问题。

另一研究方向是使用分层网络,即粗到细的方法将阅读字段从块级扩展到文档级,包含两个阶段,

例如块阅读器和文档阅读器,块阅读器首先通过选择相关句子^[173-174]或聚合局部答案输出创建文档的精简版本^[175];然后将精简版本输入文档阅读器以找到全局答案。Chen 等人^[44]提出利用检索+阅读(Retrieve+Read)的模式处理 open-domain 问答,先利用检索模块(Document Retriever)从维基百科中获取 5 个与问题最相关的段落,然后利用阅读器(Document Reader)预测出答案所在的位置。Retrospective Reader^[165]则设计两阶段验证框架,首阶段生成候选答案后通过路径评分修正答案置信度。

尽管这些方法可以扩展阅读范围,但仍存在两个缺点:(1)将整个过程分为两个连续步骤会导致错误累积;(2)使用少量文本(块级)而不是完整文本(文档级)会导致上下文信息不足。针对这些问题,RAiO^[169]采用多层上下文架构,将全局信息(文档级)丰富到每个考虑的块(块级)中,通过考虑文档的完整上下文来提取全局答案,从而在不重新训练输入文本的情况下实现全局答案的提取。上面三种模型的对比如图 9 所示。

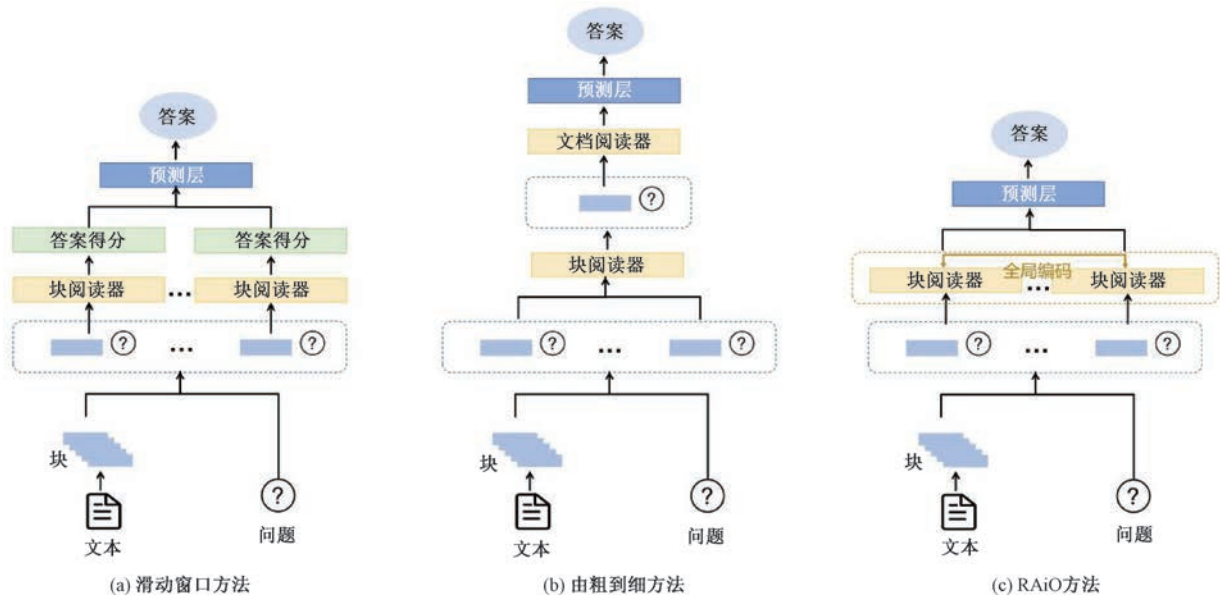


图 9 多段落阅读理解模型对比

由于不同的段落都有可能包含与标准答案类似的句子,但是有些答案并不是正确的。基于这个问题,Wang 等人^[176]提出一种模型使得来自不同段落的候选答案在基于它们所在的上下文内容里相互验证对方的正确性。具体的就是将每一篇段落中预测出来的答案与其他段落预测的答案做交互验证。这样做的原因是经观察发现,相比于错误的答案正确答案中的单词往往会在多个段落中重复出现,因此通过交互验证可以凸显出正确答案。

4.4 对话型阅读理解

无论是单段落阅读理解还是多段落阅读理解任务,它们都属于单次对话问答,即问答的形式仅有一次,后续的问题与此前的问题和答案没有关联,每一个问题皆为相互独立的。而在实际场景里,人们往往借助多次对话的形式展开交流,每一次的问题和答案都建立在之前问答情况的基础之上。因此,针对对话型阅读理解问题,在回应当前次的问题时,不仅要参考文章内容,还需考虑前几次的问题和答案。

具体可以表示为:给定 $\mathbf{D}, \mathbf{Q}_i, \mathbf{Q}_{i-1}, \dots, \mathbf{Q}_{i-k}$ 以及 $\mathbf{A}_{i-1}, \dots, \mathbf{A}_{i-k}$ 要求模型给出 \mathbf{A}_i 。其中 $\mathbf{Q}_i, \mathbf{A}_i$ 表示第 i 次的问题和答案, \mathbf{D} 表示文章, $\mathbf{Q}_{i-1}, \dots, \mathbf{Q}_{i-k}$ 和 $\mathbf{A}_{i-1}, \dots, \mathbf{A}_{i-k}$ 分别表示前 k 次的问题和答案。即通过如下建模概率实现:

$$P(\mathbf{A}_i | \mathbf{D}, \mathbf{Q}_i, \mathbf{Q}_{i-1}, \dots, \mathbf{Q}_{i-k}, \mathbf{A}_{i-1}, \dots, \mathbf{A}_{i-k}) \quad (34)$$

目前典型的对话型阅读理解数据集有 CoQA^[177]、QuAC^[178] 和 DoQA^[179]。不同之处在于 CoQA 数据集的答案形式较为简单,类似于 SQuAD^[4],但是包含有是非问题(答案是 yes/no)以及无答案问题。对于这类任务模型必须解决指代消解以及如何利用对话历史信息的问题。

为有效利用对话历史信息回答当前问题,Reddy 等人^[177]将前几次问题与答案结合到段落中,拼接后输入 DrQA^[44] 模型预测答案。Choi 等人^[178]利用 BiDAF++ 模型,在文章中设置标记向量标记历史答案单词,并添加问题轮次作为额外特征。Mark 等人^[180]借助 BiDAF++ 模型的上下文编码能力,对段落和问题深度编码,捕捉词间关系及交互信息,通过多层神经网络结构提炼整合信息。尽管 BiDAF++ 性能优于 DrQA,但其建模与交互能力不如 BERT 等预训练模型。Ohsugi 等人^[181]将对话历史中的问题和答案作为独立输入,使用 BERT 分别对段落、当前问题、历史问题和历史答案进行编码,拼接后作为预测当前答案的依据。为进一步避免对历史信息的复杂建模,Qu 等人^[182]引入了一种历史答案嵌入 HAE^[183] 方法,通过为对话历史中的答案添加特殊的嵌入信息以表示该词是否属于历史答案,如位置编码、词编码、段落编码,将会话历史无缝集成到基于 BERT 的会话问答(ConvQA)模型中,使模型能够简单但高效地理解当前问题与历史答案之间的关系。

上述的方法仅简单添加之前轮次的问题和答案,而忽略了在回答之前轮问题时模型对整篇文章的推理过程状态。FlowQA^[184] 模型采用 Integration-Flow(IF) 机制,将每次问答的语义理解状态传递到下一步,利用双向循环神经网络编码文章,单向循环神经网络编码对话历史,能集成更深层次的对话历史状态。但 FlowQA 的 IF 机制存在不足:段落词和问题轮次处理方向交织,需跨轮次细化推理结果,且对段落文本推理时未深入探索词间丰富语义关系。GraphFlow^[185] 模型通过构建动态感知问题和对话历史的上下文图,捕捉词间丰富语义关系,

每轮对话动态构建的上下文图可捕捉问题与对话历史的关联。

与上面的抽取式模型不同,生成式模型对答案生成模块要求较高。因此如何提高模型的答案生成能力是值得进一步研究的方向。预训练模型 UNILM^[83] 改进了 BERT 的训练任务,增加了自回归语言模型以及 seq2seq 语言模型,使得其在生成式任务上的效果很好,在 CoQA 数据集上远超过 Reddy 等^[175] 提出的生成式基准模型,迁移预训练模型在对话型阅读理解任务上的相关工作还包括^[186] 等。

4.5 多轮交互型阅读理解

为针对文本或段落中的信息进行深度挖掘,多轮交互阅读理解模型通常通过设计多个角度的问题,并结合多任务学习的方式深度挖掘语义信息。Chen 等^[187]认为构建双向多轮交互的 MRC 框架能够更好地捕捉方面情感三元组中的信息。通过采用双向结构,从方面(aspects)到意见表达(opinion expressions)(A→O)和意见表达到方面(O→A)两个方向进行提取,使模型能更全面地捕捉信息,减少遗漏,提升多轮推理效率和结果的准确性。Zhou 等^[188]认为引入语法关系来引导多轮 MRC 模型能更准确地识别方面和意见术语之间的依赖关系。通过语法引导网络(SG-NET)将语法关系融入自注意力层,增强输入表示,在多轮查询(非限制性查询、限制性提取查询和情感分类查询)中依次抽取头实体,关系和尾实体,更准确地捕捉词与词之间的依赖关系,减少干扰信息的影响,提高方面情感三元组提取的效率和准确性。

为进一步缓解多轮交互式抽取产生的噪声累积问题,Zhao 等^[168]提出通过加权投票策略整合多轮答案,首先基于开发集上各问题的 F1 分数为其分配权重,分数越高权重越大;然后通过加权投票机制融合不同角度提取的答案,实现多源信息的有效整合。针对多轮交互式抽取中的长距离依赖衰减问题 Chai 等^[56]利用 Longformer 模型的稀疏注意力机制处理金融文本中的长文本依赖问题,通过关系选择多轮查询模板,减少无用推理,降低计算复杂度,提高多轮推理效率。

4.6 跨语言与跨模态阅读理解

跨语言阅读理解要求模型在不同自然语言间迁移知识并回答关于文本的问题,而跨模态阅读理解则进一步将语言理解扩展至图像、语音等多源信号,实现异构信息的一致语义解析与问答。

其中跨语言阅读理解的发展高度依赖针对性数据集构建与跨语言适配技术的双重支撑, BiPaR^[31]作为首个中英双语平行的小说体裁 MRC 数据集, 包括 14668 个双语平行(段落-问题-答案)三元组, 其问题需模型具备指代消解、多句推理等复杂能力, 且支持单语、多语及跨语言三类 MRC 任务。CMRC^[30]系列数据集则聚焦中文研究的基础支撑, 通过与 SQuAD 2.0 的横向对比, 凸显中文句法结构与语义表达的独特性, 为跨语言场景下中文与其他语言的理解差异分析提供数据支持。GCRC^[32]数据集基于高考语文的现代文、古文混合试题, 对答案做可解释标注, 补充其题型分布、难度及评测指标。另有小样本语言的 MRC 数据集, 如阿拉伯语和泰米尔语维基百科数据集^[34-35]可见表 1。

技术层面上, 为解决跨语言场景中答案跨度位置偏移导致的知识迁移难题, X-STA^[33]遵循“共享、教学、对齐”三大原则: 通过梯度解耦知识共享阻断目标语言对源语言表征的梯度干扰并引入可训练修正项, 借助注意力教师引导校准实现语义对齐, 结合多粒度语义对齐强化跨语言表征一致性, 最终为低资源语言 MRC 提供高效技术路径。随着大语言模型技术发展, 尽管检索增强生成 RAG^[95]通过检索外部知识提升信息准确性, 长上下文语言模型(LCLMs, Long-Context Language Models)^[103]擅长处理长文本输入, 但两者在低资源语言的专业领域的连贯推理与可靠性方面仍存在局限。为解决这一问题, 采用强化学习的方法^[189]在语义文本相似性(STS, Semantic Textual Similarity)任务上微调 BERT, 并将调整后的模型作为机器翻译的 REINFORCE 奖励, 结合专业领域策略优化, 强化推理连贯性与结果一致性。在此基础上, Pawitsapak 等^[190]在泰语法律问答(QA)等专业领域通过群体相对策略优化(GRPO)对齐大语言模型的方法, 通过动态调整模型参数使输出更贴合泰语法律术语的表述规范与逻辑体系。

跨模态阅读理解从多种不同模态, 如文本、图像、声音等多种数据中提取和融合信息, 要求计算机在理解图像等内容的同时, 准确解析自然语言问题, 并给出恰当的回答。例如, 图文多模态 MRC 数据集 VEGA^[36]基于 arXiv 论文构建, 数据包含交互性的图像和文本, 考察模型对交错图文理解的能力。Illusory VQA^[37]聚焦于评估多模态模型对图像的认识与解读能力。视频多模态 MRC 数据集 M3-Bench^[38]含机器人视角的 M3-Bench-robot 和网络

来源的 M3-Bench-web 两部分, 以长视频问答为核心评估多模态模型的长期记忆与推理能力。以上数据集见表 1。

技术层面上, Dosovitskiy 等^[191]提出的 ViT(Vision Transformer)模型将 Transformer 架构引入计算机视觉领域, 通过将图像分割并进行序列化处理, 实现了与文本处理类似的统一建模方式, 为多模态融合提供了架构上的基础。此后, 基于对比学习的多模态模型 CLIP^[192](Contrastive Language-Image Pre-training)通过大规模的“图像-文本对”对数据集进行预训练, 学习图像和文本之间的匹配关系, 将图像和文本编码到同一向量空间中, 通过对比学习使得相似的图像和文本在空间中距离更近, 从而实现跨模态的语义理解和检索。LLaVA^[193]引入指令微调范式, 通过构建多模态指令数据集对模型进行微调。后续研究 InstructBLIP^[194]、MiniGPT-4^[195]等持续跟进, 通过对指令设计逻辑、数据集构建策略及微调训练流程的优化与拓展, 进一步丰富并完善了多模态指令微调的技术体系。除了图像数据, VideoBERT^[196]聚焦于视频和文本的统一建模问题, 将视频量化为离散的视觉词汇, 实现对视频内容的有效建模。OpenAI 发布的语音识别模型 Whisper^[197]采用编码器-解码器架构处理多种语音任务, 具备强大的语音识别和翻译能力。

4.7 零样本与少样本阅读理解

在实际应用中, 获取大规模标注数据往往面临困难, 特别是在新领域或新任务出现时, 机器阅读理解模型如何在零样本或少样本条件下快速适应和准确预测是一项重要挑战。近年来, 针对这一挑战提出了多种有效的方法, 包括以下几种。

预训练与微调策略是解决该问题的关键方法之一。通过在大规模数据集上进行预训练, 模型能够学习到丰富的知识和语义信息, 之后在少量目标数据上进行微调, 以适应特定的 MRC 任务。例如 Chen 等^[198]针对跨语言机器阅读理解的少样本问题, 通过预训练+微调方式以适应特定语言的 MRC 任务。Yu 等^[167]采用预训练和微调的训练策略, 将跨域槽填充任务抽象为机器阅读理解任务。Zhao 等^[199]认为现有的 MRC 模型通常依赖于大规模预训练语言模型, 这些模型在资源受限的环境中部署困难, 而现有的知识蒸馏方法在逻辑推理能力的迁移上存在不足。由此提出多教师多阶段知识蒸馏方法 MTMs, 基于图神经网络(Graph Neural Networks)通过逻辑教师(logical teacher)和语义教师

(semantic teacher)的指导,帮助学生模型学习更丰富的表示,学生模型基于图变换器(Graph Transformer)网络。模型引入多阶段对比学习方法,逐步对齐学生模型与教师模型的表示。

对比学习数据增强方法通过引入对比学习机制或数据增强技术,提升模型在少样本场景下的学习能力和泛化性能。Chen 等^[198]通过 hard-learning 算法最大化 $top-k$ 个预测结果中包含正确答案的可能性,其中 k 在本节中指模型对给定问题所生成的候选答案列表的长度,同时采用答案感知对比学习机制,利用高置信度的预测作为负样本,进一步提升模型性能。MRC-PASCL^[200]在微调阶段,引入答案跨度对比学习,选择围绕正确答案的跨度作为负样本,通过多任务学习增强模型对细粒度答案的预测能力。

但上述两种方法均不能很好地结合跨领域和跨语言的专业知识,在一些专业性的领域难以做出更专业化的推理。知识增强与结构化信息利用方法通过引入外部知识或结构化信息,增强模型对任务的理解和推理能力。例如,Fin-EMRC^[57]在金融领域的实体关系抽取任务中,引入金融知识图谱进行知识增强,利用结构化知识提高模型对金融文本的理解和实体关系抽取能力。此外,Xu 等^[201]通过引入外部知识构建跨句依赖图(ISDG)等结构化表示,利用通用依赖关系作为不同语言间的锚点,增强模型在跨语言场景下的理解能力。

5 讨论

5.1 机器阅读理解应用

本章主要讨论机器阅读理解模型的应用以及目前面临的主要问题和未来的发展趋势。

(1)智能客服。客服机器人是一种基于自然语言处理技术的模拟人类进行对话的服务程序,通过文字或语音与客户进行多轮交流。将用户的查询需求看作问题,通过阅读产品说明文档,利用机器阅读理解模型回答用户对产品的相关问题。

(2)辅助决策。将机器阅读理解模型引入专业性较强的领域,如医疗病例报告、法律裁判文书等,可以帮助用户更好地决策。以医疗为例,Suster 等人^[202]发布了关于医疗的阅读理解数据集 CliCR,包含大量的病例报告,同时设计了包括确诊疾病、治疗用药等问题,这些与医生的日常工作息息相关。模型如果对数据集中的问题与文本有较深的理解就可

以用于临床辅助医生诊断。

(3)智能问答。搜索引擎的目标是根据用户的查询返回相关度较高的网页,机器阅读理解在搜索引擎中的一个重要应用就是智能问答。将网页内容视为文章,从网页文本中抽取或者生成最相关的答案,避免用户主动从网页中寻找答案。

5.2 未来的发展趋势

基于预训练语言模型的能力深化 预训练语言模型作为机器阅读理解的核心基础,其能力的深化将对 MRC 模型的性能产生深远影响。未来的研究将聚焦于进一步提升 PLMs 的语义理解能力和逻辑推理能力,以更好地适应复杂的阅读理解任务。例如,通过引入更复杂的预训练目标,如因果推理、反事实推理等,增强模型对文本中隐含逻辑关系的理解。同时,结合知识图谱等外部知识资源,进一步丰富模型的语义表示,使其能够更好地处理知识密集型任务。此外,针对特定领域的预训练将使模型在专业领域表现出更高的准确性和适应性。

(1)多模态与跨模态融合。目前机器阅读理解主要集中于非结构化的文本领域,未来可扩展模型对非结构化数据(如图表、视频)的理解能力,如 CLIP^[203]风格的模型可对齐文本与医学影像特征,支持“影像报告+病史文本”联合推理。将文本、图像、表格等数据映射为统一的知识表示,如在金融领域,将财务报表中的数字与新闻文本中的事件关联,辅助投资决策问答。结合增强现实(AR)技术,实现“文本+视觉”实时交互,如维修指导场景中,用户通过文字提问并拍摄设备故障部位,模型生成图文结合的维修步骤。

(2)低资源环境的能力提升。主流 MRC 模型依赖大规模标注数据,在专业领域(如法律、医疗)和小语种场景中面临标注成本高、数据稀疏的瓶颈,未来可通过更优的开发少样本学习方法及领域自适应的预训练目标,设计轻量化模型蒸馏框架等实现跨任务迁移学习架构,提升低资源环境下的自适应学习能力。

(3)多任务协同与生态化应用。未来 MRC 将深度融入 NLP 技术生态,实现任务协同与场景扩展,支持阅读理解、文本摘要、对话生成等多任务联合训练,共享语义表示。针对医疗、法律、教育等场景开发专用模型。可设计人类反馈闭环系统,允许用户修正模型错误并迭代优化。例如,律师在使用法律问答系统时,可标注错误答案并提供修正依据,同时模型动态更新知识库。

6 总 结

本文主要回顾了机器阅读理解近年来的研究进展,对比了各个不同的阅读理解任务以及介绍了相应的数据集和评估指标,梳理出机器阅读理解发展过程中的范式转换,早期以通用框架与基于注意力机制的编码器为核心的机器阅读理解模型旨在解决问题与文本的关联问题,预训练语言模型(PLMs)使模型摆脱对特定任务数据的依赖,具备通用语言理解能力,成为MRC性能提升的关键节点,特别是大模型(LLMs)的一些关键技术的发展,极大程度降低了MRC在各类复杂场景下的应用门槛;近年来,融合显式推理结构的架构成为研究方向,旨在突破表面文本匹配的局限,提升模型在逻辑分析、因果推断等深层理解任务中的表现,以响应MRC向“类人理解”发展的需求。针对一些复杂形式的任务,我们将其视为机器阅读理解任务目前面临的挑战并且介绍了对应的解决方案。

当前,MRC的应用场景已从传统问答系统延伸至智能客服、辅助决策、智能问答等领域,今后的发展中,研究亟需聚焦更复杂推理能力构建、跨模态阅读理解融合、低资源场景性能优化等方向,不断缩小机器理解与人类理解的差异。本文通过梳理MRC任务与数据集的系统分类、剖析架构发展里程碑、总结开放挑战与应对策略,为领域研究者提供了全面的参考框架,也为推动MRC向更强理解性与泛化性的发展奠定了基础。

参 考 文 献

- [1] Lehnert W G. The process of question answering. Yale University, USA, 1977
- [2] Hirschman L, Light M, Breck E, et al. Deep Read: A reading comprehension system//Proceedings of the 37th Annual Meeting of the Association for Computational Linguistics. College Park, USA, 1999: 325-332
- [3] Hermann K M, Kočiský, T., Grefenstette, E., et al. Teaching machines to read and comprehend//Proceedings of the 29th International Conference on Neural Information Processing Systems-Volume 1. Montreal, Canada, 2015: 1693-1701
- [4] Rajpurkar P, Zhang J, Lopyrev K, et al. SQuAD: 100,000+ Questions for machine comprehension of text//Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing (EMNLP 2016). Austin, USA, 2016: 2383-2392
- [5] Joshi M, Choi E, Weld D S, et al. TriviaQA: A large scale distantly supervised challenge dataset for reading comprehension//Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (ACL 2017). Vancouver, Canada, 2017: 1601-1611
- [6] Yang Z, Qi P, Zhang S, et al. HotpotQA: A dataset for diverse, explainable multi-hop question answering//Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing (EMNLP 2018). Brussels, Belgium, 2018: 2369-2380
- [7] Chen D. Neural Reading Comprehension and Beyond. Stanford University, Stanford, USA, 2018
- [8] Hill F, Bordes A, Chopra S, Weston J. The goldilocks principle: Reading children's books with explicit memory representations//4th International Conference on Learning Representations (ICLR 2016). San Juan, Puerto Rico, 2016:1-13
- [9] Xie Q, Lai G, Dai Z, Hovy E H. Large-scale cloze test dataset designed by teachers. CoRR, 2017, abs/1711.03225
- [10] Richardson M, Burges C J C, Renshaw E. MCTest: A challenge dataset for the open-domain machine comprehension of text//Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing (EMNLP 2013). Seattle, USA, 2013: 193-203
- [11] Lai G, Xie Q, Liu H, et al. RACE: Large-scale reading comprehension dataset from examinations//Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing (EMNLP 2017). Copenhagen, Denmark, 2017: 785-794
- [12] Trischler A, Wang T, Yuan X, et al. NewsQA: A machine comprehension dataset//Proceedings of the 2nd Workshop on Representation Learning for NLP (Rep4NLP@ACL 2017). Vancouver, Canada, 2017: 191-200
- [13] Nguyen T, Rosenberg M, Song X, et al. MS MARCO: A human generated machine reading comprehension dataset//Proceedings of the Workshop on Cognitive Computation: Integrating neural and symbolic approaches 2016 (co-located with NIPS 2016). Barcelona, Spain, 2016: 1-11
- [14] He W, Liu K, Liu J, et al. DuReader: A chinese machine reading comprehension dataset from real-world applications//Proceedings of the Workshop on Machine Reading for Question Answering@ACL 2018. Melbourne, Australia, 2018: 37-46
- [15] Kociský T, Schwarz J, Blunsom P, et al. The narrativeQA reading comprehension challenge. Transactions of the Association for Computational Linguistics, 2018, 6: 317-328
- [16] Lin C-Y. ROUGE: A package for automatic evaluation of summaries//Text Summarization Branches Out. Barcelona, Spain, 2004: 74-81
- [17] Papineni K, Roukos S, Ward T, Zhu W-J. Bleu: A method for automatic evaluation of machine translation//Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics (ACL 2002). Philadelphia, USA, 2002:

- 311-318
- [18] Rajpurkar P, Jia R, Liang P. Know what you don't know: Unanswerable questions for SQuAD//Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (ACL 2018). Melbourne, Australia, 2018; 784-789
- [19] Welbl J, Stenetorp P, Riedel S. Constructing datasets for multi-hop reading comprehension across documents. Transactions of the Association for Computational Linguistics, 2018, 6: 287-302
- [20] Dua D, Wang Y, Dasigi P, et al. DROP: A reading comprehension benchmark requiring discrete reasoning over paragraphs//Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics (NAACL-HLT 2019). Minneapolis, USA, 2019; 2368-2378
- [21] Dasigi P, Liu N F, Marasovic A, et al. Quoref: A reading comprehension dataset with questions requiring coreferential reasoning//Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP 2019). Hong Kong, China, 2019; 5924-5931
- [22] Liu J, Cui L, Liu H, et al. LogiQA: A challenge dataset for machine reading comprehension with logical reasoning//Proceedings of the 29th International Joint Conference on Artificial Intelligence (IJCAI-20), Yokohama, Japan, 2020; 3953-3957
- [23] Yu W, Jiang Z, Dong Y, et al. ReClor: A reading comprehension dataset requiring logical reasoning//Proceedings of the 8th International Conference on Learning Representations (ICLR 2020). Addis Ababa, Ethiopia, 2020; 7577-7597
- [24] Chen K, Xu W, Cheng X, et al. Question directed graph attention network for numerical reasoning over text//Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP). Online, 2020; 6759-6768
- [25] Talmor A, Herzig J, Lourie N, et al. CommonsenseQA: A question answering challenge targeting commonsense knowledge//Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL-HLT 2019), Minneapolis, USA, 2019; 4149-4158
- [26] Huang L, Bras R L, Bhagavatula C, et al. Cosmos QA: Machine reading comprehension with contextual commonsense reasoning//Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP 2019). Hong Kong, China, 2019; 2391-2401
- [27] Bhagavatula C, Bras R L, Malaviya C, et al. Abductive commonsense reasoning//Proceedings of the 8th International Conference on Learning Representations (ICLR 2020). Addis Ababa, Ethiopia, 2020; 2921-2938
- [28] Yue Z, Liu J, Zhang C, et al. MA-MRC: A multi-answer machine reading comprehension dataset//Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval. Taipei, China, 2023; 2144-2148
- [29] Moradisan H, Zarrinkalam F, Serbanescu J, Noorian, Z. UnAnswGen: A systematic approach for generating unanswerable questions in machine reading comprehension//Proceedings of the 2024 Annual International ACM SIGIR Conference on Research and Development in Information Retrieval in the Asia Pacific Region (SIGIR-AP 2024). Tokyo, Japan, 2024; 280-286
- [30] Cui Y, Liu T, Che W, et al. A span-extraction dataset for Chinese machine reading comprehension//Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP 2019), Hong Kong, China, 2019; 5882-5888
- [31] Jing Y, Xiong D, Yan Z. BiPaR: A bilingual parallel dataset for multilingual and cross-lingual reading comprehension on novels//Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing. Hong Kong, China, 2019; 2452-2462
- [32] Tan H, Wang X, Ji Y, et al. GCRC: A new challenging MRC dataset from Gaokao Chinese for explainable evaluation//Proceedings of the Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021. Online, 2021; 1319-1330
- [33] Cao T, Wang C, Tan C, et al. Sharing, teaching and aligning: Knowledgeable transfer learning for cross-lingual machine reading comprehension. 2023. arXiv:2311.06758
- [34] Obeidat R, Al-Harbi M, Al-Ayyoub M, et al. ArQuAD: An expert-annotated arabic machine reading comprehension dataset. Cognitive Computation, 2024, 16(3): 984-1003
- [35] Sinthusha A V A, Charles E Y A, Weerasinghe R. Machine reading comprehension for the tamil language with translated SQuAD. IEEE Access, 2025, 13: 13312-13328
- [36] Zhou C, Zhang M, Chen P, et al. VEGA: Learning interleaved image-text comprehension in vision-vision-language large models. Technical Report; arXiv:2406.10228, 2024
- [37] Rostamkhani M, Ansari B, Sabzevari H, et al. Illusory VQA: Benchmarking and enhancing multimodal models on visual illusions//Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops. Nashville, USA, 2025; 2995-3004
- [38] Long L, He Y, Ye W, et al. Seeing, listening, remembering, and reasoning: A multimodal agent with long-term memory. Computing Research Repository, 2025, abs/2508.09736
- [39] Seo M J, Kembhavi A, Farhadi A, et al. Bidirectional attention flow for machine comprehension//Proceedings of the 5th International Conference on Learning Representations (ICLR

- 2017). Toulon, France, 2017:2649-2661
- [40] Rumelhart D E, Hinton G E, Williams R J. Learning representations by back-propagating errors//Neurocomputing: Foundations of Research. Cambridge, USA; MIT Press, 1988:696-699
- [41] Mikolov T, Sutskever I, Chen K, et al. Distributed representations of words and phrases and their compositionality//Advances in Neural Information Processing Systems 26 (NIPS 2013). Lake Tahoe, USA, 2013; 3111-3119
- [42] Pennington J, Socher R, Manning C D. GloVe: Global vectors for word representation//Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP 2014). Doha, Qatar, 2014: 1532-1543
- [43] Peters M E, Neumann M, Iyyer M, et al. Deep contextualized word representations//Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics (NAACL-HLT 2018). New Orleans, USA, 2018: 2227-2237
- [44] Chen D, Fisch A, Weston J, et al. Reading wikipedia to answer open-domain questions//Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (ACL 2017). Vancouver, Canada, 2017: 1870-1879
- [45] Hochreiter S, Schmidhuber J. Long Short-Term Memory. Neural Computation, 1997, 9(8): 1735-1780
- [46] Cho K, Merriënboer B van, Gülçehre Ç, et al. Learning phrase representations using RNN encoder-decoder for statistical machine translation//Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP). Doha, Qatar, 2014: 1724-1734
- [47] Vaswani A, Shazeer N, Parmar N, et al. Attention is all you need//Advances in Neural Information Processing Systems 30 (NIPS 2017). Long Beach, USA, 2017: 5998-6008
- [48] Wang S, Jiang J. Machine comprehension using match-LSTM and answer pointer//Proceedings of the 5th International Conference on Learning Representations (ICLR 2017). Toulon, France, 2017:1223-1233
- [49] Vinyals O, Fortunato M, Jaitly N. Pointer networks//Advances in Neural Information Processing Systems 28 (NIPS 2015). Montreal, Quebec, Canada, 2015: 2692-2700
- [50] See A, Liu P J, Manning C D. Get to the point: Summarization with pointer-generator networks//Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (ACL 2017). Vancouver, Canada, 2017: 1073-1083
- [51] Kadlec R, Schmid M, Bajgar O, et al. Text understanding with the attention sum reader network//Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (ACL 2016). Berlin, Germany, 2016: 951-961
- [52] Xiong C, Zhong V, Socher R. Dynamic coattention networks for question answering//Proceedings of the 5th International Conference on Learning Representations (ICLR 2017). Toulon, France, 2017:2396-2409
- [53] Wang W, Yang N, Wei F, et al. Gated self-matching networks for reading comprehension and question answering//Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (ACL 2017). Vancouver, Canada, 2017: 189-198
- [54] He P, Liu X, Gao J, et al. DeBERTa: Decoding-enhanced BERT with disentangled attention//9th International Conference on Learning Representations (ICLR 2021). Virtual Event, Austria, 2021
- [55] Child R, Gray S, Radford A, et al. Generating long sequences with sparse transformers. arXiv preprint arXiv:1904.10509, 2019
- [56] Beltagy I, Peters M E, Cohan A. Longformer: The long-document transformer. arXiv preprint arXiv:2004.05150, 2020
- [57] Chai Y, Chen M, Wu H, et al. Fin-EMRC: An efficient machine reading comprehension framework for financial entity-relation extraction. IEEE Access, 2023, 11: 82685-82695
- [58] Yuan J, Gao H, Dai D, et al. Sparse attention in its native form: Hardware-aligned and natively trainable sparse attention. arXiv preprint arXiv:2502.11089, 2025
- [59] Chen D, Bolton J, Manning C D. A comprehensive analysis of the CNN/Daily mail reading comprehension task//Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (ACL 2016). Berlin, Germany, 2016: 2356-2365
- [60] Dhingra B, Liu H, Yang Z, et al. Gated-attention readers for text comprehension//Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (ACL 2017). Vancouver, Canada, 2017: 1832-1846
- [61] Sordani A, Bachman P, Bengio Y. Iterative alternating neural attention for machine reading. arXiv preprint arXiv:1606.02245, 2016
- [62] Liu Z, Huang D, Huang K, Zhang J. DIM Reader: Dual interaction model for machine comprehension//Chinese Computational Linguistics and Natural Language Processing Based on Naturally Annotated Big Data (CCL 2017). Lecture Notes in Computer Science, Vol. 10565. Springer, Cham, 2017: 387-397
- [63] Hu M, Peng Y, Huang Z, et al. Reinforced mnemonic reader for machine reading comprehension//Proceedings of the Twenty-Seventh International Joint Conference on Artificial Intelligence (IJCAI 2018). Stockholm, Sweden, 2018: 4099-4106
- [64] Chen Qiu-Yi. Research on question answering algorithm based on multi-granularity and validator scoring mechanism. Chongqing University of Technology, Chongqing, 2023 (in Chinese)
(陈秋怡. 基于多粒度和验证器打分机制的问答算法研究. 重庆理工大学, 重庆, 2023)
- [65] Ba L J, Kiros J R, Hinton G E. Layer normalization. arXiv preprint arXiv:1607.06450, 2016
- [66] He K, Zhang X, Ren S, Sun J. Deep residual learning for

- image recognition//2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR 2016). Las Vegas, USA, 2016: 770-778
- [67] Raffel C, Shazeer N, Roberts A, et al. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research*, 2020, 21: 140:1-140:67
- [68] Radford A, Narasimhan K, Salimans T, et al. Improving language understanding by generative pre-training. San Francisco, USA, 2018
- [69] Radford A, Wu J, Child R, et al. Language models are unsupervised multitask learners. *OpenAI Blog*, 2019, 1(8): 9
- [70] Brown T B, Mann B, Ryder N, et al. Language models are few-shot learners//*Advances in Neural Information Processing Systems 33 (NeurIPS 2020)*, Vancouver, Canada, 2020: 159-183
- [71] Zeng A, Xu B, Wang B, et al. ChatGLM: A family of large language models from GLM-130B to GLM-4 all tools. *arXiv preprint arXiv:2406.12793*, 2024
- [72] Su J, Ahmed M H M, Lu Y, et al. RoFormer: Transformer with rotary position embedding enhanced. *Neurocomputing*, 2024, 568: 127063
- [73] Touvron H, Martin L, Stone K, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023
- [74] Dubey A, Jauhri A, Pandey A, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024
- [75] Press O, Smith N A, Lewis M. Train short, test long: Attention with linear biases enables input length extrapolation//*Proceedings of the 10th International Conference on Learning Representations (ICLR 2022)*. Virtual, 2022
- [76] Hendrycks D, Gimpel K. Gaussian error linear units (GELUs). *arXiv preprint arXiv:1606.08415*, 2023
- [77] Shazeer N. GLU variants improve transformer. *arXiv preprint arXiv:2002.05202*, 2020
- [78] Zhang B, Sennrich R. Root mean square layer normalization//*Advances in Neural Information Processing Systems 32 (NeurIPS 2019)*. Vancouver, Canada, 2019: 12360-12371
- [79] Touvron H, Lavril T, Izacard G, et al. LLaMA: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023
- [80] Wang H, Ma S, Dong L, et al. DeepNet: Scaling transformers to 1,000 layers. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024, 46(10): 6761-6774
- [81] Aohan Z, Xiao L, Zhengxiao D, et al. GLM-130B: An open bilingual pre-trained model. *arXiv preprint arXiv:2210.02414*, 2023
- [82] Devlin J, Chang M-W, Lee K, Toutanova K. BERT: Pre-training of deep bidirectional transformers for language understanding//*Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019*. Minneapolis, USA, 2019: 4171-4186
- [83] Liu Y, Ott M, Goyal N, et al. RoBERTa: A robustly optimized BERT pretraining approach. *arXiv preprint arXiv:1907.11692*, 2019
- [84] Dong L, Yang N, Wang W, et al. Pre-training of unified language model for natural language understanding and generation//*Proceedings of the 33rd Conference on Neural Information Processing Systems (NeurIPS 2019)*. Vancouver, Canada, 2019: 1170-1182
- [85] Lan Z, Chen M, Goodman S, et al. ALBERT: A lite BERT for self-supervised learning of language representations//*8th International Conference on Learning Representations (ICLR 2020)*. Addis Ababa, Ethiopia, 2020: 11574-11587
- [86] Clark K, Luong M-T, Le Q V, Manning C D. ELECTRA: Pre-training text encoders as discriminators rather than generators//*8th International Conference on Learning Representations (ICLR 2020)*. Addis Ababa, Ethiopia, 2020: 3295-3312
- [87] Liu X, He P, Chen W, Gao J. Multi-task deep neural networks for natural language understanding//*Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics (ACL 2019)*. Florence, Italy, 2019: 4487-4496
- [88] Hudson D A, Manning C D. GQA: A novel dataset for real-world visual reasoning and compositional question answering. *Technical Report*; *arXiv:1902.09506*, 2019
- [89] OpenAI, Achiam J, Adler S, et al. GPT-4 Technical report. *Technical Report*; *arXiv:2303.08774*, OpenAI, 2024
- [90] Xie S M, Raghunathan A, Liang P, et al. An interpretation of in-context learning as implicit Bayesian inference//*The 10th International Conference on Learning Representations, ICLR 2022*. Virtual, 2022: 1-15
- [91] Liu H, Teng Z, Cui L, et al. LogiCoT: Logical chain-of-thought instruction tuning//*Findings of the Association for Computational Linguistics: EMNLP 2023*. Singapore, 2023: 2908-2921
- [92] Li Qing, Li Yan-Ling, Dong Jie, et al. Research survey on machine reading comprehension based on logical reasoning. *Journal of Computer Science and Exploration*, 2024, 18(8): 1998-2013 (in Chinese)
(李晴, 李艳玲, 董杰等. 基于逻辑推理的机器阅读理解综述. *计算机科学与探索*, 2024, 18(8): 1998-2013)
- [93] Stiennon N, Ouyang L, Wu J, et al. Learning to summarize from human feedback//*Proceedings of the 34th International Conference on Neural Information Processing Systems*. Vancouver, Canada, 2020: 253-266
- [94] Shao Z, Wang P, Zhu Q, et al. DeepSeekMath: Advancing the boundaries of mathematical reasoning in open language models. *Technical Report*; *arXiv:2402.03300*, 2024
- [95] Lewis P, Perez E, Piktus A, et al. Retrieval-augmented generation for knowledge-intensive NLP tasks//*Advances in Neural Information Processing Systems*. 2020, 33: 9459-9474

- [96] Jiang A Q, Sablayrolles A, Roux A, et al. Mixtral of experts. CoRR, 2024, abs/2401.04088
- [97] Joshi M, Chen D, Liu Y, et al. SpanBERT: Improving pre-training by representing and predicting spans. Transactions of the Association for Computational Linguistics, 2020, 8: 64-77
- [98] Rae J W, Borgeaud S, Cai T, et al. Methods, analysis & insights from training Gopher: Scaling language models. CoRR, 2021, abs/2112.11446
- [99] Chowdhery A, Narang S, Devlin J, et al. PaLM: Scaling language modeling with pathways. Journal of Machine Learning Research, 2023, 24(240): 1-113
- [100] OpenAI, Jaech A, Kalai A, et al. OpenAI O1 system card. 2024
- [101] Guo D, Yang D, Zhang H, et al. DeepSeek-R1: Enhancing reasoning ability in LLMs through reinforcement learning. arXiv preprint arXiv:2501.12948, 2025. 1, 3
- [102] Kimi Team, Du A, Gao B, et al. Kimi K1.5: Scaling reinforcement learning with LLMs. arXiv preprint arXiv:2501.12599, 2025. 1, 3
- [103] Reid M, Savinov N, Teplyashin D, et al. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. 2024. arXiv:2403.05530
- [104] Kaushik D, Lipton Z C. How much reading does reading comprehension require? A critical investigation of popular benchmarks//Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing. Brussels, Belgium, 2018; 5010-5015
- [105] Zhou M, Duan N, Liu S, Shum H Y. Progress in neural NLP: Modeling, learning, and reasoning. Engineering, 2020, 6(3): 275-290
- [106] Ouyang S, Zhang Z, Zhao H. Logical reasoning driven by facts for machine reading comprehension//Proceedings of the 38th AAAI Conference on Artificial Intelligence (AAAI 2024). Vancouver, Canada, 2024; 18851-18859
- [107] Bauer L, Wang Y, Bansal M. Commonsense for generative multi-hop question answering tasks//Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing. Brussels, Belgium, 2018; 4220-4230
- [108] Wu B, Zhang Z, Zhao H. Graph-free multi-hop reading comprehension: A select-to-guide strategy. CoRR, 2021, abs/2107.11823
- [109] Min S, Zhong V, Zettlemoyer L, Hajishirzi, H. Multi-hop reading comprehension through question decomposition and rescoring//Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics (ACL 2019). Florence, Italy, 2019; 6097-6109
- [110] Shen Y, Huang P-S, Gao J, Chen W. ReasoNet: Learning to cease reading in machine comprehension//Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining. Halifax, Canada, 2017; 1047-1055
- [111] Song J, Tang S, Qian T, et al. Reading document and answering question via global attentional inference//Advances in Multimedia Information Processing PCM2018. Cham, 2018; 335-345
- [112] Tu M, Wang G, Huang J, et al. Multi-hop reading comprehension across multiple documents through reasoning over heterogeneous graphs//Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics. Florence, Italy, 2019; 2704-2713
- [113] Song L, Wang Z, Yu M, et al. Exploring graph-structured passage representation for multi-hop reading comprehension with graph neural networks. CoRR, 2018, abs/1809.02040
- [114] Xu W, Zhang H, Cai D, Lam W. Dynamic semantic graph construction and reasoning for explainable multi-hop science question answering//Findings of the Association for Computational Linguistics: ACL/IJCNLP 2021. Online, 2021; 1044-1056
- [115] Su D, Xu Y, Dai W, et al. Multi-hop question generation with graph convolutional network//Findings of the Association for Computational Linguistics; EMNLP 2020. Online, 2020; 4636-4647
- [116] Zheng B, Wen H, Liang Y, et al. Graph attention networks for document modeling in multi-grained machine reading comprehension//Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics (ACL 2020). Online, 2020; 6708-6718
- [117] Wang S, Zhong W, Tang D, et al. Logic-driven context extension and data augmentation for logical reasoning of text//Findings of the Association for Computational Linguistics: ACL 2022. Dublin, Ireland, 2022; 1619-1629
- [118] Li X, Cheng G, Chen Z, et al. AdaLoGN: Adaptive logic graph network for reasoning-based machine reading comprehension//Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers). Dublin, Ireland, 2022; 7147-7161
- [119] Cetto M, Niklaus C, Freitas A, et al. Graphene: Semantically-linked propositions in open information extraction//Proceedings of the 27th International Conference on Computational Linguistics (COLING 2018). Santa Fe, USA, 2018; 2300-2311
- [120] Hou S, Zhang S, Fei C. Rhetorical structure theory: A thorough review of theory, parsing approaches and applications. Expert Systems with Applications, 2020, 157: 113421
- [121] Huang Y, Fang M, Cao Y, et al. DAGN: Discourse-aware graph network for logical reasoning//Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL-HLT 2021). Online, 2021; 5848-5855
- [122] Habernal I, Wachsmuth H, Gurevych I, et al. The argument reasoning comprehension task: Identification and re-

- construction of implicit warrants//Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics; Human Language Technologies (NAACL-HLT 2018). New Orleans, USA, 2018; 1930-1940
- [123] Chen J, Zhang Z, Zhao H. Modeling hierarchical reasoning chains by linking discourse units and key phrases for reading comprehension//Proceedings of the 29th International Conference on Computational Linguistics (COLING 2022). Gyeongju, Republic of Korea, 2022; 1467-1479
- [124] Xu F, Liu J, Lin Q, et al. Logiformer: A two-branch graph transformer network for interpretable logical reasoning//SIGIR '22: The 45th International ACM SIGIR Conference on Research and Development in Information Retrieval. Madrid, Spain, 2022; 1055-1065
- [125] Honnibal M. SpaCy 2: Natural language understanding with bloom embeddings, convolutional neural networks and incremental parsing. To appear, 2017
- [126] Levi F W. Finite geometrical systems: Six public lectures delivered in february, 1940, at the university of calcutta. Calcutta, India: The University of Calcutta, 1942
- [127] Bordes A, Usunier N, García-Durán A, et al. Translating embeddings for modeling multi-relational data//Advances in Neural Information Processing Systems 26 (NIPS 2013). Lake Tahoe, USA, 2013; 2787-2795
- [128] Jiao F, Guo Y, Song X, et al. MERIt: Meta-path guided contrastive learning for logical reasoning//Findings of the Association for Computational Linguistics; ACL 2022. Dublin, Ireland, 2022; 3496-3509
- [129] Sanyal S, Xu Y, Wang S, et al. APOLLO: A straightforward method for adaptive pretraining of language models in logical reasoning//Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1; Long Papers). Toronto, Canada, 2023; 6308-6321
- [130] Pi X, Zhong W, Gao Y, et al. LogiGAN: Adversarial pre-training for learning logical reasoning//Advances in Neural Information Processing Systems 35 (NeurIPS 2022). New Orleans, USA, 2022
- [131] Xu Z, Yang Z, Cui Y, et al. IDOL: Indicator-oriented logic pre-training for logical reasoning//Findings of the Association for Computational Linguistics; ACL 2023. Toronto, Canada, 2023; 8099-8111
- [132] Thawakar O, Dissanayake D, More K, et al. Llamav-01: Rethinking step-by-step visual reasoning in llms. arXiv preprint arXiv:2501.06186, 2025
- [133] Xu G, Jin P, Hao L, Song Y, Sun L, Yuan L. Llava-01: Let vision language models reason step-by-step. arXiv preprint arXiv:2411.10440, 2024
- [134] Yao H, Huang J, Wu W, et al. Mulberry: Enhancing mllm with o1-like reasoning and reflection through collective Monte Carlo tree search. arXiv preprint arXiv:2412.18319, 2024
- [135] Zhang R, Zhang B, Li Y, et al. Improve vision language model chain-of-thought reasoning. arXiv preprint arXiv: 2410.16198, 2024
- [136] Chen X, Liang C, Yu A W, et al. Neural symbolic reader: Scalable integration of distributed and symbolic representations for reading comprehension//8th International Conference on Learning Representations (ICLR 2020). Addis Ababa, Ethiopia, 2020;11377-11392
- [137] Zhou Y, Bao J, Duan C, et al. OPERA: Operation-pivoted discrete reasoning over text//Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics; Human Language Technologies (NAACL 2022). Seattle, USA, 2022; 1655-1666
- [138] Ran Q, Lin Y, Li P, et al. NumNet: Machine reading comprehension with numerical reasoning//Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP2019). Hong Kong, China, 2019; 2474-2484
- [139] Han S, Gao N, Guo X, et al. Aggregating heterogeneous neighbors and node types for numerical reasoning over text//ICCAI'22: 8th International Conference on Computing and Artificial Intelligence. Tianjin, China, 2022; 200-206
- [140] Geva M, Gupta A, Berant J. Injecting numerical reasoning skills into language models//Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics (ACL 2020). Online,2020; 946-958
- [141] Kim J, Kang J, Kim K-m, et al. Exploiting numerical-contextual knowledge to improve numerical reasoning in question answering//Findings of the Association for Computational Linguistics; NAACL 2022. Seattle, USA, 2022; 1811-1821
- [142] Pi X, Liu Q, Chen B, et al. Reasoning akin to program executors//Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing (EMNLP 2022). Abu Dhabi, United Arab Emirates, 2022; 761-779
- [143] Hu M, Peng Y, Huang Z, et al. A multi-type and multi-span network for reading comprehension requiring discrete reasoning//Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP 2019). Hong Kong, China, 2019; 1596-1606
- [144] Yang A, Li A, Yang B, et al. Qwen3: Technical report. arXiv preprint arXiv:2505.09388, 2025
- [145] Davis E, Marcus G. Commonsense reasoning and commonsense knowledge in artificial intelligence. Communications of the ACM, 2015, 58(9): 92-103
- [146] Talmor A, Elazar Y, Goldberg Y, et al. oLMPics-On what language model pre-training captures. Transactions of the Association for Computational Linguistics, 2020, 8: 743-758

- [147] Speer R, Chin J, Havasi C. ConceptNet 5.5: An open multilingual graph of general knowledge//Proceedings of the 31st AAAI Conference on Artificial Intelligence (AAAI 2017). San Francisco, USA, 2017; 4444-4451
- [148] Kilgarriff A, Fellbaum C. WordNet: An electronic lexical database. *Language*, 2000, 76(3): 706
- [149] Ji S, Pan S, Cambria E, et al. A survey on knowledge graphs: Representation, acquisition, and applications. *IEEE Transactions on Neural Networks and Learning Systems*, 2022, 33(2): 494-514
- [150] Chen W, Quan X, Chen C. Gated convolutional networks for commonsense machine comprehension//Neural Information Processing. Cham, 2018: 297-306
- [151] Zhao Y, Zhang Z, Zhao H. Reference knowledgeable network for machine reading comprehension. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 2022, 30:1461-1473
- [152] Lin B Y, Chen X, Chen J, et al. KagNet: Knowledge-aware graph networks for commonsense reasoning//Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP). Hong Kong, China, 2019: 2829-2839
- [153] Gong P, Liu J, Yang Y, et al. Towards knowledge enhanced language model for machine reading comprehension. *IEEE Access*, 2020, 8:224837-224851
- [154] Lauscher A, Vulic I, Ponti E M, et al. Specializing unsupervised pretraining models for word-level semantic similarity//Proceedings of the 28th International Conference on Computational Linguistics (COLING 2020). Barcelona, Spain, 2020: 1371-1383
- [155] Poerner N, Waltinger U, Schütze H. E-BERT: Entity embeddings for BERT that are efficient yet effective//Findings of the Association for Computational Linguistics; EMNLP 2020. Online, 2020: 803-818
- [156] Yamada I, Asai A, Sakuma J, et al. Wikipedia2Vec: An efficient toolkit for learning and visualizing the embeddings of words and entities from Wikipedia//Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations (EMNLP 2020-Demos). Online, 2020: 23-30
- [157] Edge D, Trinh H, Cheng N, et al. From local to global: a Graph RAG approach to query-focused summarization. Technical Report; arXiv:2404.16130, 2024
- [158] Clark C, Gardner M. Simple and effective multi-paragraph reading comprehension//Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (ACL 2018). Melbourne, Australia, 2018; 845-855
- [159] Hu M, Peng Y, Huang Z, Li D. Read + Verify: Machine reading comprehension with unanswerable questions//Proceedings of the 33rd AAAI Conference on Artificial Intelligence (AAAI 2019). Honolulu, Hawaii, USA, 2019: 6529-6537
- [160] Back S, Chinthakindi S C, Kedia A, et al. NeurQuRI: Neural question requirement inspector for answerability prediction in machine reading comprehension//8th International Conference on Learning Representations (ICLR 2020). Addis Ababa, Ethiopia, 2020:4932-4943
- [161] Kiddon C, Zettlemoyer L, Choi Y. Globally coherent text generation with neural checklist models//Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing (EMNLP 2016). Austin, USA, 2016: 329-339
- [162] Levy O, Seo M, Choi E, et al. Zero-shot relation extraction via reading comprehension//Proceedings of the 21st Conference on Computational Natural Language Learning (CoNLL 2017). Vancouver, Canada, 2017: 333-342
- [163] Sun F, Li L, Qiu X, et al. U-Net: Machine reading comprehension with unanswerable questions. arXiv preprint arXiv:1810.06638, 2018
- [164] Wu Z, Xu H. A multi-task learning machine reading comprehension model for noisy document (student abstract)//Proceedings of the 34th AAAI Conference on Artificial Intelligence (AAAI 2020). New York, USA, 2020: 13963-13964
- [165] Zhang Z, Yang J, Zhao H. Retrospective Reader for Machine Reading Comprehension//Proceedings of the 35th AAAI Conference on Artificial Intelligence (AAAI 2021). Virtual, 2021: 14506-14514
- [166] Wang X, Liu J, Wang J, et al. Document-level relation extraction based on machine reading comprehension and hybrid pointer-sequence labeling. *ACM Transactions on Asian Low-Resource Language Information Processing*, 2024, 23(7): 100:1-16
- [167] Liu J, Yu M, Chen Y, Xu J. Cross-domain slot filling as machine reading comprehension: A new perspective. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 2022, 30: 673-685
- [168] Zhao T, Yan Z, Cao Y, Li Z. A machine reading comprehension based framework for joint entity-relation extraction: Asking effective and diverse questions//Proceedings of the 29th International Joint Conference on Artificial Intelligence (IJCAI 2020). Virtual, 2020: 3948-3954
- [169] Phan T-A, Jung J J, Bui K-H N. Read-all-in-once (RAiO): Multi-layer contextual architecture for long-text machine reading comprehension. *IEEE Access*, 2023, 11: 77873-77879
- [170] Zaheer M, Guruganesh G, Dubey K A, et al. Big Bird: Transformers for longer sequences//Advances in Neural Information Processing Systems 33 (NeurIPS 2020), Vancouver, Canada, 2020: 1450-1464
- [171] Soleimani A, Monz C, Worring M. NLQuAD: A non-factoid long question answering data set//Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics (EACL 2021). Online, 2021:

- 1245-1255
- [172] Tan C, Wei F, Yang N, et al. S-Net: For machine reading comprehension, from answer extraction to answer generation. arXiv preprint arXiv:1706.04815, 2017
- [173] Choi E, Hewlett D, Uszkoreit J, et al. Long document question answering in a coarse-to-fine manner//Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (ACL 2017). Vancouver, Canada, 2017; 209-220
- [174] Ding M, Zhou C, Yang H, Tang J. CogLTX: Applying BERT to long texts//Advances in Neural Information Processing Systems 33 (NeurIPS 2020), Vancouver, Canada, 2020; 1073-1715
- [175] Zhao J, Bao J, Wang Y, et al. RoR: Read-over-read for long document machine reading comprehension//Findings of the Association for Computational Linguistics: EMNLP 2021. Virtual, 2021; 1862-1872
- [176] Wang Y, Liu K, Liu J, et al. Multi-passage machine reading comprehension with cross-passage answer verification//Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (ACL 2018). Melbourne, Australia, 2018; 1918-1927
- [177] Reddy S, Chen D, Manning C D. CoQA: A challenge for conversational question answering. Transactions of the Association for Computational Linguistics, 2019, 7: 249-266
- [178] Choi E, He H, Iyyer M, et al. QuAC: Question answering in context//Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing (EMNLP 2018). Brussels, Belgium, 2018; 2174-2184
- [179] Campos J A, Otegi A, Soroa A, et al. DoQA-Accessing domain-specific FAQs via conversational QA//Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics (ACL 2020). Online, 2020; 7302-7314
- [180] Yatskar M. A qualitative comparison of CoQA, SQuAD 2.0 and QuAC//Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics (NAACL-HLT 2019). Minneapolis, USA, 2019; 2318-2323
- [181] Ohsugi Y, Saito I, Nishida K, et al. A straightforward yet efficient approach to integrate multi-turn context with BERT for conversational machine comprehension. arXiv preprint arXiv:1905.12848, 2019
- [182] Qu C, Yang L, Qiu M, et al. BERT with history answer embedding for conversational question answering//Proceedings of the 42nd International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR2019). Paris, France, 2019; 1133-1136
- [183] Qu C, Yang L, Qiu M, et al. Attentive history selection for conversational question answering//Proceedings of the 28th ACM International Conference on Information and Knowledge Management (CIKM 2019). Beijing, China, 2019; 1391-1400
- [184] Huang H-Y, Choi E, Yih W-t. FlowQA: Grasping flow in history for conversational machine comprehension//Proceedings of the 7th International Conference on Learning Representations (ICLR 2019). New Orleans, USA, 2019; 2036-2050
- [185] Chen Y, Wu L, Zaki M J. GraphFlow: Utilizing conversation flow through graph neural networks for conversational machine comprehension//Proceedings of the 29th International Joint Conference on Artificial Intelligence (IJCAI 2020). Virtual, 2020; 1230-1236
- [186] Zhu C, Zeng M, Huang X. SDNet: A deep network with contextualized attention for conversational question answering. arXiv preprint arXiv:1812.03593, 2018
- [187] Chen S, Wang Y, Liu J, et al. Aspect sentiment triplet extraction based on bidirectional machine reading comprehension//Proceedings of the 35th AAAI Conference on Artificial Intelligence (AAAI 2021). Virtual, 2021; 12666-12674
- [188] Zhou Y, Huang W, Wang J, et al. An aspect sentiment triplet extraction method based on syntax-guided multi-turn machine reading comprehension//Proceedings of the 2024 16th International Conference on Machine Learning and Computing (ICMLC 2024). Shenzhen, China, 2024; 7-13
- [189] Yasui G, Tsuruoka Y, Nagata M. Using semantic similarity as reward for reinforcement learning in sentence generation//Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics: Student Research Workshop. Florence, Italy, 2019; 400-406
- [190] Akarajaradwong P, Chaksangchaichot C, Pothavorn P, et al. Can group relative policy optimization improve thai legal reasoning and question answering? arXiv preprint arXiv: 2507.09638, 2025
- [191] Dosovitskiy A, Beyer L, Kolesnikov A, et al. An image is worth 16x16 words; Transformers for image recognition at scale. 2020 <https://arxiv.org/abs/2010.11929>
- [192] Radford A, Kim J W, Hallacy C, et al. Learning transferable visual models from natural language supervision. 2021. <https://arxiv.org/abs/2103.00020>
- [193] Liu H T, Li C Y, Wu Q Y, et al. Visual instruction tuning//Proceedings of the 37th International Conference on Neural Information Processing Systems. New Orleans, USA, 2023; 34892-34916
- [194] Dai W, Li J N, Li D X, et al. InstructBLIP: Advancing towards general-purpose vision-language models through instruction tuning//Proceedings of the 37th International Conference on Neural Information Processing Systems. New Orleans, USA, 2023, 36: 49250-49267
- [195] Zhu D, Chen J, Shen X, et al. Minigpt-4: Enhancing vision-language understanding with advanced large language models//Proceedings of the International Conference on Representation Learning (ICLR). Arlington, USA, 2024; 40652-40668
- [196] Sun C, Myers A, Vondrick C, et al. VideoBERT: A joint

- model for video and language representation learning//Proceedings of the 2019 IEEE/CVF International Conference on Computer Vision. Seoul, Republic of Korea, 2019: 7464-7473
- [197] Radford A, Kim J W, Xu T, et al. Robust speech recognition via large-scale weak supervision//Proceedings of the 40th International Conference on Machine Learning. Honolulu, USA, 2023: 28492-28518
- [198] Chen N, Shou L, Gong M, et al. From good to best: Two-stage training for cross-lingual machine reading comprehension//Proceedings of the 36th AAAI Conference on Artificial Intelligence(AAAI 2022). Virtual, 2022: 10501-10508
- [199] Zhao Z, Xie Z, Zhou G, et al. MTMS: Multi-teacher multi-stage knowledge distillation for reasoning-based machine reading comprehension//Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR2024). Washington, USA, 2024: 11
- [200] Li R, Xiao Q, Yang J, et al. MRC-PASCL: A few-shot machine reading comprehension approach via post-training and answer span-oriented contrastive learning. IEEE/ACM Transactions on Audio, Speech, and Language Processing, 2024, 32: 4838-4849
- [201] Xu L, Zhang X, Zong B, et al. Zero-shot cross-lingual machine reading comprehension via inter-sentence dependency graph//Proceedings of the 36th AAAI Conference on Artificial Intelligence (AAAI 2022). Virtual, 2022: 11538-11546
- [202] Suster S, Daelemans W. CliCR: A dataset of clinical case reports for machine reading comprehension//Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL-HLT 2018). New Orleans, USA, 2018: 1551-1563
- [203] Radford A, Kim J W, Hallacy C, et al. Learning transferable visual models from natural language supervision//Proceedings of the 38th International Conference on Machine Learning (ICML 2021). Virtual, 2021: 8748-8763



ZHANG Tian-Cheng, Ph. D. , professor. His research interests include artificial intelligence, spatiotemporal data management, and smart education.

WANG Ya-Ting, M. S. Her research interests include knowledge tracing, natural language processing and explainable artificial intelligence.

LI Fan, Ph. D. His research interests include machine learning, deep learning, knowledge tracing, and explainable

artificial intelligence.

SUN Xiang-Hui, M. S. His research interests include smart education and explainable artificial intelligence.

YU Ming-He, Ph. D. , associate professor. Her research interests include databases, information retrieval, and intelligent education.

YU Ge, Ph. D. , professor. His research interests include distributed and parallel databases, OLAP and data warehousing, data integration, and graph data management.

Background

The early machine reading comprehension systems relied on manually coded rules and suffered from the limitation of poor generalization ability. Later, the emergence of attention-based neural network models changed this landscape, achieving a qualitative leap in performance on benchmark datasets such as CNN&Daily Mail. Subsequently, the advent of Pre-trained Language Models and the advancement of Large Language Models technology have enabled them to integrate richer semantic information when utilized as feature extractors. These models not only excel at capturing long-range features but also possess robust parallel processing capabilities. To further enhance the reasoning capabilities of MRC models, numerous existing works have conducted in-depth research and achieved breakthroughs in four key directions: multi-hop reasoning, logical reasoning, numerical reasoning, and commonsense reasoning. In addition, many excellent solutions have been proposed for the task challenges

faced by MRC, such as unanswerable MRC, multi-answer MRC, dialogue-based MRC, multi-turn interactive MRC, cross-lingual/cross-modal MRC, and zero-shot/few-shot MRC, which we further explore.

This study presents a comprehensive review of MRC research since 2015, dissecting task taxonomies, datasets, and model architectures. By addressing key research gaps in reasoning capabilities and domain generalization, this work aims to deepen the understanding of MRC and foster the development of more robust, generalizable models, thereby laying the foundation for next-generation intelligent MRC solutions.

We gratefully acknowledge support from the General Program of the National Natural Science Foundation of China (No. 62272093), the State Key Program of the National Natural Science Foundation of China (No. 62137001), and the International Cooperation and Exchange of the National Natural Science Foundation of China (No. 62461146205).