

基于文档的对话研究

孙润鑫 马龙轩 张伟男 刘 挺
(哈尔滨工业大学社会计算与信息检索研究中心 哈尔滨 150001)
(rxsun@ir.hit.edu.cn)

Research on Document Grounded Conversations

Sun Runxin, Ma Longxuan, Zhang Weinan, and Liu Ting
(Research Center for Social Computing and Information Retrieval, Harbin Institute of Technology, Harbin 150001)

Abstract Document grounded conversations is an emerging hot task in the field of dialogue system. Different from previous tasks, it needs to consider both the utterances and the given document. However, previous work focused on the relationship between the two, but ignored the utterances' difference in the effect of response generation. To solve this problem, a new dialectical approach to the dialogue history, which means the utterances before the last one, is proposed in this paper. At the encoding step, it divides the modeling of the semantic information into two parts: using history and ignoring history, and then uses the comparative integration method to summarize the branch results. In this way, when the dialogue history is not related to the current utterance, it can avoid being introduced as noise which will damage the performance of the model. Besides, it also strengthens the guiding role of the current utterance in the information filtering process. Experimental results show that compared with the existing baselines, this model can generate responses that are more in line with the current context and more informative, indicating that it can better understand dialogue information and conduct knowledge filtering. And through the ablation study, the effectiveness of each module in the modeling process is also verified.

Key words document grounded conversations; reply generation; information filtering; Transformer model; attention mechanism

摘 要 基于文档的对话是目前对话领域一个新兴的热点任务.与以往的任务不同,其需要将对话信息和文档信息综合进行考虑.然而,先前的工作着重考虑二者之间的关系,却忽略了对话信息中的句子对回复生成的作用具有差异性.针对这一问题,提出了一种新的辩证看待对话历史的方法,在编码阶段讨论利用历史和忽略历史 2 种情况进行语义信息的建模,并采用辩证整合的方式进行分支信息的汇总.由此避免了在历史信息与当前对话不相关时,其作为噪声被引入进而损害模型性能,同时也强化了当前对话对信息筛选的指导作用.实验结果表明,该模型与现有基线模型相比,能够生成更为符合当前语境且信息量更加丰富的回复,从而说明其能够更好地理解对话信息并进行知识筛选.并且通过进行消融实验,也验证了各模块在建模过程中的有效性.

关键词 基于文档的对话;回复生成;信息筛选;Transformer 模型;注意力机制

中图法分类号 TP391

收稿日期:2020-08-18;修回日期:2020-12-24
基金项目:国家自然科学基金项目(62076081,61772153,61936010);2030 新一代人工智能重大项目(2020AAA0108605)
This work was supported by the National Natural Science Foundation of China (62076081, 61772153, 61936010) and 2030 Major Project of New Generation Artificial Intelligence of China (2020AAA0108605).
通信作者:张伟男(wnzhang@ir.hit.edu.cn)

近年来,随着大数据的发展以及深度学习技术的兴起,对话系统这一研究领域已经取得了长足的进步.由于 Seq2Seq(sequence-to-sequence)模型^[1]的出现,许多模型采用了端到端的方式并基于大规模的人类对话语料进行训练,已经收到了一定的效果.然而,这些模型的问题在于其总是倾向于生成较为通用的回复^[2],这也是目前该领域面临的一大挑战.

事实上,通过引入外部知识使得模型生成信息更加丰富的回复,正逐渐成为这一问题可能且可行的方法之一.并且到目前为止,已经出现了一些富有建设性的成果.Ghazvininejad 等人^[3]是对一组候选事实进行检索,然后将其结果与对话信息相结合.Zhang 等人^[4]借助说话人的个人资料,使得模型可以生成一致性更好且更吸引人的对话.而 Mem2Seq^[5]则通过引入外部知识库的方式,以提升回复内容的丰富度.但是,这些工作大多是将现有的事实类知识整合到对话系统中,而此类知识的构建过程费时费力,且具有较大的局限性,从而在一定程度上制约了大规模数据的使用,以及现有系统的应用和落地.

基于文档的对话是一项在讨论特定文档的内容时,生成较自然对话回复的任务.与以往基于知识的对话不同,该任务采用文档作为知识源,因此能够使用较为广泛的知识,且无需对其预先进行构建.并且其与基于文档的问答也有所不同,相较于后者此时不仅需要考虑对话与文档之间的关系,同时还应考虑对话内部的联系以对其充分理解.针对这一任务,我们应主要从 3 个方面进行考虑:1)理解对话中各句之间的关系;2)理解对话与文档之间的关系;3)根据对话上下文,从文档中筛选有关信息以生成相应回复.

在本文中,我们提出了一种新颖的基于 Transformer^[6]的模型,用于更好地构建对话上下文,并利用其进行信息筛选从而生成更有意义的回复.主要思想是将信息筛选分为利用历史和忽略历史 2 种情况进行讨论:当忽略历史时,只需以当前对话为依据,对文档信息进行相应筛选;而当考虑历史时,应先利用其对上下文进行重新建模,然后以所得结果为依据,进行文档信息的筛选.注意到对话历史事实上并不总是与语境相关,因此在完成上述过程后我们需要对其与当前对话的相关程度进行度量,并以此为依据进行语义向量的整合,以避免引入不相关的信息作为噪声.实验结果表明,与现有的基线模型相比,我们的模型能生成语义更加连贯、知识更为丰富的回复.

本文的主要贡献可总结为 3 点:

1) 提出了一种辩证看待对话历史的方法,用于文档信息的筛选以及上下文信息的构建.在编码过程中,以分支的方式同时进行利用历史和忽略历史的信息筛选,并根据当前语境判断各分支的重要程度.

2) 提出了一种新的信息汇总的方式,通过判断当前对话与历史之间的相关程度,来对各分支所得到的信息进行整合,以得到在解码过程中所使用的语义向量.

3) 在 CMU-DoG 数据集^[7]上进行了模型的验证,并取得了当前已知的最好效果.同时借助消融实验,也显示出模型各个部分在语义向量建模过程中的有效性.

1 相关工作

目前,基于文档的对话已被证明在解决通用回复、改善回复质量方面能够取得不错的效果^[7-9].就现有工作而言,我们可以将其大致分为 2 类:基于抽取的方法和基于生成的方法.

基于抽取的方法最初在机器阅读理解(machine reading comprehension, MRC)中被提出^[10],其是在一篇给定的文档中抽取特定的片段,以回答所给问题的任务.Vinyals 等人^[11]提出了一种名为指针网络(pointer networks, Ptr-Nets)的结构,其使用注意力机制^[12]从输入序列中选取元素作为输出.Seo 等人^[13]提出了 BiDAF 模型,其使用双向注意力流以获得问题感知的上下文表征.Wang 等人^[14]提出了 R-Net 模型,其使用一种自匹配注意力机制来对较长段落的表示进行完善.然而,这种方法实际上并不适合基于文档的对话.其原因主要有 2 点:1)在该任务中,事实上并没有一个明确的问题,这需要在充分理解对话内容的前提下,自行决定回复句的主题;2)由于方法本身的限制,只从原文中进行抽取使得对话的流利度难以得到保证^[15].

至于基于生成的方法,目前绝大多数模型还是基于经典的 Seq2Seq 结构^[1].Lian 等人^[16]利用后验知识来进一步指导知识的筛选过程.Liu 等人^[17]引入知识图谱作为又一知识源,从而与非结构化的文档相结合起到互为增强的效果.Li 等人^[18]提出了一种增量式 Transformer 编码器来对多轮对话及其相关文档进行联合建模,同时应用 2 阶段的推敲解码器以进一步增强对话的连贯性和信息的丰富性.

但是,之前的工作并没有深入研究对话历史与当前对话之间的关系.考虑真实的对话场景,由于对话过程中往往会出现主题的转换,且由于开放域对话的随意性,后续的主题并不一定总是与历史相关.在这种情况下,历史便成为了一种噪声,从而对知识的挑选以及后续对话的生成起到了负面作用.但当主题较为明晰、对话内容较为集中时,历史信息可以帮助我们更好地理解当前语境.由此便可生成内容更加丰富、语义更为连贯的回复.出于以上 2 点考虑,本文提出了一种辩证看待对话历史的方法,通过判断其与当前对话之间的相关程度,以更好地利用对话历史并避免引入不必要的噪声.

2 基于比较整合的 Transformer 模型

2.1 问题定义

形式上,我们对基于文档的对话这一任务做如下定义.令 $\mathbf{U}=(\mathbf{u}^{(1)}, \mathbf{u}^{(2)}, \cdots, \mathbf{u}^{(K)})$ 为由 K 句话组成的整段对话,而 $\mathbf{u}^{(k)}=(u_1^{(k)}, u_2^{(k)}, \cdots, u_l^{(k)})$ 则表示包含 l 个单词的第 k 句话,其中 $u_i^{(k)}$ 即为第 k 句话中的第 i 个单词.同样地,对每段对话 \mathbf{U} 而言, $\mathbf{d}^{(k+1)}=(d_1^{(k+1)}, d_2^{(k+1)}, \cdots, d_j^{(k+1)})$ 表示与第 $k+1$ 句话也即回复相关、包含 j 个单词的文档.从而基于文档的对话则可被定义为在给定相关文档 $\mathbf{d}^{(k+1)}$ 以及先前 k 句话 $\mathbf{U}^{(\leq k)}$ 的情况下,生成回复 $\mathbf{u}^{(k+1)}$ 的概率:

$$P(\mathbf{u}^{(k+1)} | \mathbf{U}^{(\leq k)}, \mathbf{d}^{(k+1)}; \boldsymbol{\theta}) =$$

$$\prod_{i=1}^l P(u_i^{(k+1)} | \mathbf{U}^{(\leq k)}, \mathbf{d}^{(k+1)}, \mathbf{u}_{<i}^{(k+1)}; \boldsymbol{\theta}), \quad (1)$$

其中 $\mathbf{u}_{<i}^{(k+1)}=(u_1^{(k+1)}, u_2^{(k+1)}, \cdots, u_{i-1}^{(k+1)})$.

2.2 模型架构

模型基于经典的 Transformer 结构,各组成模块如图 1 所示.可以看出,其主要由 3 部分构成:

1) 自注意编码器(self-attentive encoder, SA). SA 是经典 Transformer 模型^[6]中所使用的编码器,其负责将对话信息和文档信息分别预先进行编码.

2) 筛选式 Transformer 编码器(selective transformer encoder, STE).STE 是一种用于进行信息筛选的 Transformer 编码器,它通过注意力机制^[6]对编码后的对话历史和文档信息分别进行筛选,以得到更为符合当前语境的向量表示,用于后续语义信息的构建.

3) 推敲解码器(deliberation decoder, DD).DD 是一个 2 阶段的 Transformer 解码器,用于生成语义更为连贯的回复^[18].在第 1 阶段,其将经 SA 编码后的当前对话 $\mathbf{u}^{(k)}$ 和整合后的语义向量作为输入,借助对话上下文来生成回复.而第 2 阶段则是将经 SA 编码后的前一阶段输出和相关文档 $\mathbf{d}^{(k+1)}$ 作为输入,使用文档知识进一步完善回复.

接下来我们对模型中的各个环节进行详细介绍.

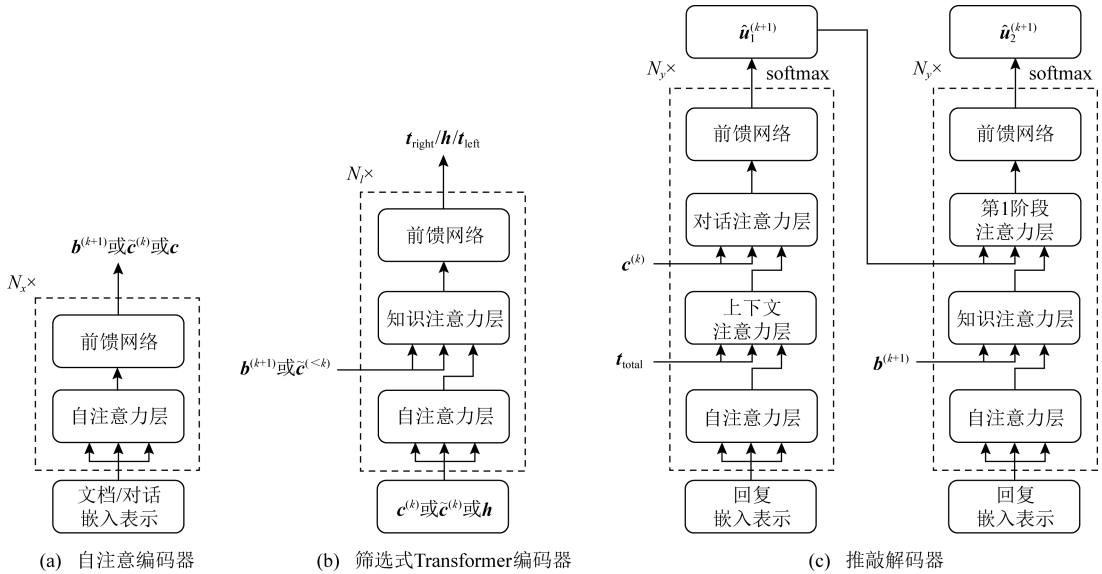


Fig. 1 Composition modules of our model

图 1 本文模型的各构成模块

2.2.1 输入信息的编码

由于文档信息相对较长,因此如何解决长距离

依赖问题,并得到语义更加丰富的单词表示便成为了首要问题.这里我们基于多头自注意力机制(multi-

head self-attention)^[6],使用自注意编码器来计算相关文档信息 $\mathbf{d}^{(k+1)}$ 的表示.编码器的输入 $\tilde{\mathbf{d}}^{(k+1)}$ 是一个融合了位置编码(positional encoding)^[6]的单词向量序列:

$$\tilde{\mathbf{d}}^{(k+1)} = (\tilde{\mathbf{d}}_1^{(k+1)}, \tilde{\mathbf{d}}_2^{(k+1)}, \dots, \tilde{\mathbf{d}}_f^{(k+1)}), \quad (2)$$

$$\tilde{\mathbf{d}}_j^{(k+1)} = \mathbf{e}_{d_j} + PE(j), \quad (3)$$

其中 \mathbf{e}_{d_j} 是 $\mathbf{d}_j^{(k+1)}$ 的词向量表示,而 $PE(\cdot)$ 表示位置编码函数.

如图 1(a)所示,自注意编码器由 N_x 个相同的层堆叠而成,每层都有 2 个子层.第 1 个子层是多头自注意力子层,而第 2 个子层则是简单的、位置完全连接的前馈网络(feed-forward network, FFN).同时每个子层也都采用了残差连接和标准化机制^[6],因此子层的输出实际为 $LayerNorm(x + Sublayer(x))$,其中 $Sublayer(x)$ 是由子层本身所实现的功能,为了表述简洁这里故将其省略.下面给出该过程的公式表示:

$$\mathbf{A}^{(1)} = MultiHead(\tilde{\mathbf{d}}^{(k+1)}, \tilde{\mathbf{d}}^{(k+1)}, \tilde{\mathbf{d}}^{(k+1)}), \quad (4)$$

$$\mathbf{D}^{(1)} = FFN(\mathbf{A}^{(1)}), \quad (5)$$

$$FFN(\mathbf{x}) = \max(\mathbf{0}, \mathbf{x}\mathbf{W}_1 + \mathbf{b}_1)\mathbf{W}_2 + \mathbf{b}_2, \quad (6)$$

其中, $\mathbf{A}^{(1)}$ 是经第 1 个自注意力子层计算出的隐层状态,而 $\mathbf{D}^{(1)}$ 则是第 1 层最终得到的文档表示.对每层而言,重复这一过程,从而我们有:

$$\mathbf{b}^{(k+1)} = \mathbf{D}^{(N_x)} = SA_{doc}(\tilde{\mathbf{d}}^{(k+1)}), \quad (7)$$

其中, $\mathbf{b}^{(k+1)}$ 是文档信息 $\mathbf{d}^{(k+1)}$ 的最终表示,而 $SA_{doc}(\cdot)$

表示上述编码的完整过程.

同样地,我们用 $\tilde{\mathbf{u}}^{(k)} = (\tilde{\mathbf{u}}_1^{(k)}, \tilde{\mathbf{u}}_2^{(k)}, \dots, \tilde{\mathbf{u}}_l^{(k)})$ 来表示融入了位置信息的每句对话 $\mathbf{u}^{(k)}$.然后重复上述过程对其进行编码,所得结果记为 $\mathbf{c}^{(k)}$.不过本次编码过程中所采用的 SA_u 与 SA_d 结构相同,但参数不同.

2.2.2 文档信息的筛选

为了借助注意力机制进行高效地信息筛选,我们使用了一种结构简单的筛选式 Transformer 编码器.由图 1(b)可知,其与 2.2.1 节中叙述的结构最主要的区别就是增加了一个用于进行信息筛选的注意力子层.在这一子层,输入不再满足 $\mathbf{Q} = \mathbf{K} = \mathbf{V}$ 这一特征,而是令 $\mathbf{K} = \mathbf{V}$ 且 $\mathbf{Q} \neq \mathbf{K}$,通过利用投影后 \mathbf{Q} 与 \mathbf{K} 之间的相似关系,来从 \mathbf{V} (也即是 \mathbf{K}) 中进行信息的抽取.

出于辩证看待对话历史的考虑,在编码阶段,我们分利用历史和忽略历史 2 种情况讨论文档信息的筛选,整体结构如图 2 所示.先从忽略历史开始,容易想到,只需将编码后的当前对话 $\mathbf{c}^{(k)}$ 作为下层输入,而将编码后的文档表示 $\mathbf{b}^{(k+1)}$ 作为上层输入即可,公式表述为

$$\mathbf{t}_{right} = STE_{doc}(\mathbf{c}^{(k)}, \mathbf{b}^{(k+1)}, \mathbf{b}^{(k+1)}), \quad (8)$$

其中 STE_{doc} 是用于进行文档信息筛选的编码器.然后将 \mathbf{t}_{right} ,也即忽略历史进行信息筛选的输出作为右支最终的语义向量.

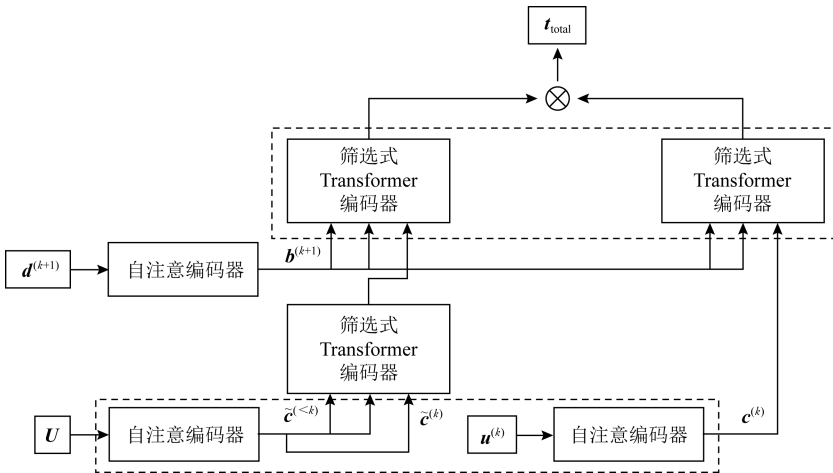


Fig. 2 Overall architecture of the encoding part of our model

图 2 本文模型编码部分的整体架构

接着考虑对话历史.类比上述过程,首先我们还是以编码后的当前对话 $\mathbf{c}^{(k)}$ 为依据,从编码后的对话历史 $\mathbf{c}^{(<k)}$ 中进行信息的筛选,得到符合当前语境的历史表示 \mathbf{h} ,然后再以其为依据,从编码后的文档

表示 $\mathbf{b}^{(k+1)}$ 中进行信息的筛选,得到用于回复生成的语义表示.需要注意的是,这里的 $\mathbf{c}^{(k)}$ 与 $\mathbf{c}^{(<k)}$ 并不相同.出于借助历史以更好地理解对话信息的考虑,此处我们将当前对话 $\mathbf{u}^{(k)}$ 与历史信息 $\mathbf{u}^{(<k)}$ 进行

拼接,并重新送入自注意编码器中:

$$\mathbf{c}^{(<k)}, \tilde{\mathbf{c}}^{(k)} = \text{SA}_{\text{utt}}([\mathbf{u}^{(<k)}; \mathbf{u}^{(k)}]). \quad (9)$$

从而考虑对话历史的筛选过程可以表示为

$$\mathbf{h} = \text{STE}_{\text{utt}}(\mathbf{c}^{(k)}, \tilde{\mathbf{c}}^{(<k)}, \mathbf{c}^{(<k)}), \quad (10)$$

$$\mathbf{t}_{\text{left}} = \text{STE}_{\text{doc}}(\mathbf{h}, \mathbf{b}^{(k+1)}, \mathbf{b}^{(k+1)}), \quad (11)$$

其中 STE_{utt} 是用于进行对话信息筛选的编码器.同样地,我们将最终得到的结果作为左支也即借助历史的语义向量.

2.2.3 语义向量的整合

在得到了左右支的语义向量之后,一个关键的步骤便是如何将二者进行整合,以得到最终的语义表示.Sankar 等人^[19]指出,当前对话 $\mathbf{u}^{(k)}$ 在回复生成的过程中起着决定性的作用.因此这里我们需要保留右支得到的信息,并根据对话历史与当前对话的相关程度,确定最终表示中是否包含左支信息及其所占的比重.

我们可以通过注意力机制^[12]并结合最大池化(max pooling)的方式,来计算当前对话中的各词与对话历史之间的相关程度^[20]:

$$s_{ij} = \mathbf{v}^T \tanh(\mathbf{W}_1 \mathbf{c}_i^{(k)} + \mathbf{W}_2 \tilde{\mathbf{c}}_j^{(<k)} + \mathbf{b}), \quad (12)$$

$$\alpha = \text{sigmoid}(\max_{\text{col}}(\mathbf{S})). \quad (13)$$

基于 Zheng 等人^[21]所给出的形式进行左右支信息的整合:

$$\mathbf{t}_{\text{total}} = \alpha \odot \mathbf{t}_{\text{left}} + (1 - \alpha) \odot \mathbf{t}_{\text{right}} + \mathbf{t}_{\text{right}}, \quad (14)$$

此时 $\mathbf{t}_{\text{total}}$ 即为最终用于进行回复生成的语义表示.

2.2.4 解码与回复生成

受现实世界中人类认知过程的启发,Li 等人^[18]设计了一种 2 阶段的推敲解码器,以提高知识的相关性和上下文的连贯性,这里我们同样使用这一结构.在第 1 阶段,其将编码后的当前对话 $\mathbf{c}^{(k)}$ 和整合后的语义向量 $\mathbf{t}_{\text{total}}$ 作为输入,以学习生成上下文连贯的回复.而第 2 阶段则是将编码后的前一阶段输出和相关文档 $\mathbf{b}^{(k+1)}$ 作为输入,试图将文档知识进一步注入到生成回复中.

如图 1(c)所示,推敲解码器由第 1 和第 2 阶段解码器构成.2 个解码器具有相同的架构,但子层的输入不同.同样地,其也由 N_y 个相同的层堆叠而成,每层都有 4 个子层.这里仅以第 1 阶段为例,对相应流程进行详细说明.首先第 1 个子层同样是多头自注意力子层:

$$\mathbf{H}_1^{(n)} = \text{MultiHead}(\mathbf{R}_1^{(n-1)}, \mathbf{R}_1^{(n-1)}, \mathbf{R}_1^{(n-1)}), \quad (15)$$

其中 $n=1, 2, \dots, N_y$, $\mathbf{R}_1^{(n-1)}$ 是前一层的输出,并且有 $\mathbf{R}_1^{(0)} = \mathbf{u}_{<_i}^{(k+1)}$,也即融合了位置编码的已生成部分

的向量表示.第 2 个子层是多头上下文注意力层:

$$\mathbf{B}_1^{(n)} = \text{MultiHead}(\mathbf{H}_1^{(n)}, \mathbf{t}_{\text{total}}, \mathbf{t}_{\text{total}}), \quad (16)$$

其中 $\mathbf{t}_{\text{total}}$ 即为式(14)所得的整合后的语义表示.第 3 个子层是多头对话注意力层:

$$\mathbf{C}_1^{(n)} = \text{MultiHead}(\mathbf{B}_1^{(n)}, \mathbf{c}^{(k)}, \mathbf{c}^{(k)}), \quad (17)$$

其中 $\mathbf{c}^{(k)}$ 即为编码后的当前对话.值得一提的是,式(16)(17)分别将文档信息和对话信息融入到生成回复中,以求在生成信息更加丰富的回复时,不失其连贯性.第 4 个子层是全连接的前馈网络:

$$\mathbf{R}_1^{(n)} = \text{FFN}(\mathbf{C}_1^{(n)}). \quad (18)$$

在 N_y 层之后,我们使用 softmax 函数获得第 1 阶段解码出的单词的概率分布:

$$P(\hat{\mathbf{u}}_1^{(k+1)}) = \text{softmax}(\mathbf{R}_1^{(N_y)}), \quad (19)$$

其中 $\hat{\mathbf{u}}_1^{(k+1)}$ 即为第 1 阶段解码器所解码出的回复.

2.3 优化目标

按照 Li 等人^[18]所给出的方法,这里我们摒弃了原始推敲解码器所采用的较为复杂的训练算法,其是一种基于蒙特卡洛思想的联合学习框架^[22].而是出于连贯性的考虑,选择对解码器的 2 个阶段分别计算损失,并采用加和的方式同时进行训练^[23],其公式为

$$L_{\text{mle}} = L_{\text{mle1}} + L_{\text{mle2}}, \quad (20)$$

$$L_{\text{mle1}} = - \sum_{n=1}^N \sum_{i=1}^I \ln P(\hat{\mathbf{u}}_{1(i)}^{(k+1)}), \quad (21)$$

$$L_{\text{mle2}} = - \sum_{n=1}^N \sum_{i=1}^I \ln P(\hat{\mathbf{u}}_{2(i)}^{(k+1)}), \quad (22)$$

其中 N 是参与训练的样本总数,而 I 是参考回复的长度.

3 实验设置

3.1 数据集

我们使用基于文档的对话数据集(CMU-DoG)^[7]来进行模型的评估.该数据集共包含基于 120 篇文档的 4 112 段对话,平均每段对话的轮数为 21.43.文档内容取自维基百科中有关热门电影的文章,而对话内容不仅可以基于文档,也可以是随意闲聊的.需要注意的是,这里我们将同一个人的连续对话视为一整句,并忽略每段对话起始部分的问候语.每个样例包含一段与回复相关的、平均词数为 200 的文档,最近 3 句作为上下文的对话和 1 句参考回复.数据集的统计信息详见表 1,更多细节请查阅文献^[7],在此不再赘述.

Table 1 Statistics of CMU-DoG Dataset
表 1 CMU-DoG 数据集的统计信息

属性	训练集	验证集	测试集
对话总数	3 373	229	619
语句总数	107 792	7 030	19 375
对话平均语句数	31.96±14.37	30.70±13.55	31.30±13.00
语句平均长度	12.57±36.83	13.87±11.98	12.43±11.68
样例数目	66 332	3 269	10 502

3.2 基线模型

我们将所提出的模型与以下 2 类基线模型进行比较:基于 RNN 的模型和基于 Transformer 的模型.

1) 基于 RNN 的模型

① S2S(Seq2Seq).S2S 是一种经典的编码器-解码器模型^[1].根据 Zhou 等人^[7]给出的方法,这里我们将编码得到的文档知识,与上一时刻生成的词拼接起来作为解码器的输入,以向生成回复中添加文档信息.

② S2SIF(S2S with input feeding).S2SIF 是一种 Seq2Seq 模型的变体^[24].其是利用注意力机制来对编码结果进行动态汇总,并将其输入到下一时刻的解码过程.这里我们对编码得到的文档表示进行这一操作.

2) 基于 Transformer 的模型

① T-MemNet (transformer memory network).T-MemNet 是一种将 Transformer 与 Memory Network 相结合^[9]的模型.首先利用注意力机制对知识进行预筛选,然后将其与对话信息拼接后送入编码器中进行二次编码.在这里由于相关知识已经给定,故省去筛选步骤.

② ITE-DD.ITE-DD 是一种使用增量式 Transformer 编码器和 2 阶段推敲解码器的模型^[18].增量式编码器使用多头注意力机制,将文档信息和对话上下文合并到每句对话的编码过程中.而推敲解码器采用 2 阶段的方式,分别利用对话上下文生成回复主干以及将文档知识进一步注入回复之中.

3.3 实验设置

我们使用 OpenNMT-py^[25]作为代码框架.对于所有模型,隐层规模均设置为 300.对基于 RNN 的模型,按照文献^[7]中的描述,我们使用 2 层的双向 LSTM 作为编码器,并使用 1 层单向的 LSTM 作为解码器.对基于 Transformer 的模型,编码器和解码器的层数均设置为 3.多头注意力层中的注意力头数为 6,过滤器(filter)的规模为 1 200.我们使用预训练

的词向量,即在 Google News 语料库上借助 Word2Vec 训练得到的 300 维词向量^[26]对嵌入矩阵进行初始化.并且这一表示会被文档、对话和回复所共享.我们使用 Adam 优化器^[27]来对模型进行优化.在解码阶段,集束规模(beam size)被设置为 5.

3.4 评价指标

1) 自动评价.我们使用 BLEU, ROUGE-1, ROUGE-2,ROUGE-L 作为生成回复的自动评价指标.Zheng 等人^[21]指出,由于文档中真正有意义的词往往出现频率较小,较低的困惑度(perplexity, PPL)反而会使得模型倾向于生成通用回复,因此这里我们没有使用这一指标.BLEU 主要度量生成回复和参考回复之间的 n 元组(n -gram)重叠,而 ROUGE 则从召回角度来评测回复的语义完整度.

2) 人工评价.我们随机选取了 100 个测试样例并邀请了 4 位评委,在给定之前 3 句对话和文档信息的条件下,从 3 个方面对回复质量进行评判:①连贯性(coherence),回复是否符合当前语境,语义逻辑是否连贯;②信息性(informativeness),回复是否包含有用信息,是否与所给文档相关而非通用回复;③自然性(naturalness),回复语义是否通顺,是否符合语法规定和人类表述习惯.所有指标均被划分为 3 个等级:0 代表较差;1 代表中等;2 代表较好.

4 结果分析

4.1 实验结果

表 2 列出了所有模型在第 3 节设置下的自动评价和人工评价结果.

在自动评价方面,我们的模型(CITE-DD)在所有的指标上都优于基线模型.其中在 BLEU 指标上模型取得了约 30%(即 $(1.31-1.01)/1.01\div 30\%$)的性能提升,这也在一定程度上体现了其结果与参考回复具有更高的重合度,从而表明其能够生成更有意义的回复.不仅如此,模型在 ROUGE 指标上也有不小的进步,由此便可以看出其结果对参考回复具有更高的覆盖度,也即更为全面地表达了参考回复的内容.综合来看,模型相较于基线模型,能够生成质量更高的回复,其结构的有效性也得到了较为显著的体现.

而在人工评价方面,我们的模型同样表现不错.值得一提的是,模型在信息性这一指标上相较于基线模型有着较为显著的提升.从而可以反映出其能够更好地筛选和利用文档信息,而这也是基于文档

的对话任务中的一大关键.另外对于连贯性和自然性而言,模型也显现出了一定的效果.不过在连贯性上,模型与 ITE-DD 相比区别不大,考虑主要是因为

推敲解码器已经大幅提升了生成回复的连贯性,若要在此基础上取得新的进展,需要对其有针对性的改进,但本文的重点并不在于此.

Table 2 Experimental Results of All Models on CMU-DoG Dataset

表 2 所有模型在 CMU-DoG 数据集上的实验结果

模型	BLEU	ROUGE-1	ROUGE-2	ROUGE-L	连贯性	信息性	自然性
S2S	0.24	10.26	1.37	7.42	0.74	0.29	1.25
S2SIF	0.18	10.32	1.28	7.37	0.78	0.15	1.29
T-MemNet	1.01	13.97	2.09	10.05	0.87	0.48	1.26
ITE-DD	1.01	14.14	2.37	10.48	1.13	0.69	1.29
CITE-DD(本文)	1.31	15.74	2.67	11.59	1.20	0.96	1.43

注:黑体值为最优值.

4.2 消融实验

为了验证模型各部分的有效性,我们还设计并进行了消融实验,结果如表 3 所示:

Table 3 Experimental Results of Ablation Study

表 3 消融实验的结果

模型	BLEU	ROUGE-1	ROUGE-2	ROUGE-L
-left branch	0.88	13.29	2.09	9.75
-right branch	1.27	15.30	2.54	11.37
-integration	1.14	15.28	2.57	11.31
CITE-DD(本文)	1.31	15.74	2.67	11.59

注:黑体值为最优值.

在这里,-left branch 表示将编码部分的左支(式(10)(11)),也即利用历史的部分去掉,此时模型便退化为基于文档的单轮对话.显然,完全忽略对话历史会对回复的生成有着极大的影响,这也较为明显地在表 3 指标中得以体现.考虑主要有 2 个原因:1)其包含了丰富的语义信息,利用其进行文档信息的筛选可以得到更多、更准确的结果;2)其亦可对更好地理解当前对话起到一定的帮助作用.当然上面 2 点都建立在历史信息与当前对话具有较强相关性的基础之上.而-right branch 则表示将编码部分的右支,也即忽略历史的部分去掉,此时模型便仅能从利用历史的角度进行文档信息的筛选以及语义信息的建模.可以看出,回复生成的质量也受到了部分影响,这也在一定程度上说明了当历史信息与当前对话不相关时,忽略历史的方法使得模型避免了不必要噪声的引入,从而可以更为准确地理解当前对话和构建对话上下文.

另外,-integration 表示将语义向量的整合部分

去掉,这里采用直接加和的方式进行替换,此时模型便失去了辩证看待对话历史的能力.注意在这种情况下,模型的性能反倒比只用左支时要差,由此便进一步印证了本文所提方法的有效性与合理性.

4.3 权重可视化

如图 3 所示,这里我们将注意力矩阵进行可视化处理,并以此说明整合方式的合理性.在该例中,历史信息为“i did! she really delivered a knockout in mean girls. what was your favorite scene in mean girls? i personally like the scene where cady met the plastics. <SEP> i love the revenge plot the best i think”,而当前对话为“oh yeah! the plan of revenge against regina? that was awesome! did you know that mean girls was partially based on a book?”.

根据图 3 中的结果,可以看出模型认为当前对话标出的部分与历史信息较为相关,而事实也同样如此.显然,2 段文字具有高度的重合性,这时左支得到的信息才可视为右支对应位置的补充而非噪声.由此便可说明模型最终的效果与我们的预期相吻合,从而进一步印证了本文方法的有效性.

4.4 样例分析

我们从测试集中选取了一个例子,以较为直观地说明所提模型的有效性,内容详如表 4 所示.可以看出,对话是一段典型的基于电影知识的闲聊:其没有明确的问题引出下句,当前对话与历史之间的关系也并不密切.在这种情况下,历史信息如果不加筛选反而会成为一种噪声,从而在一定程度上分散了回复句主题的确定.此时,已有的基线方法就不能很好地理解对话信息并从文档中提取出相关知识,从

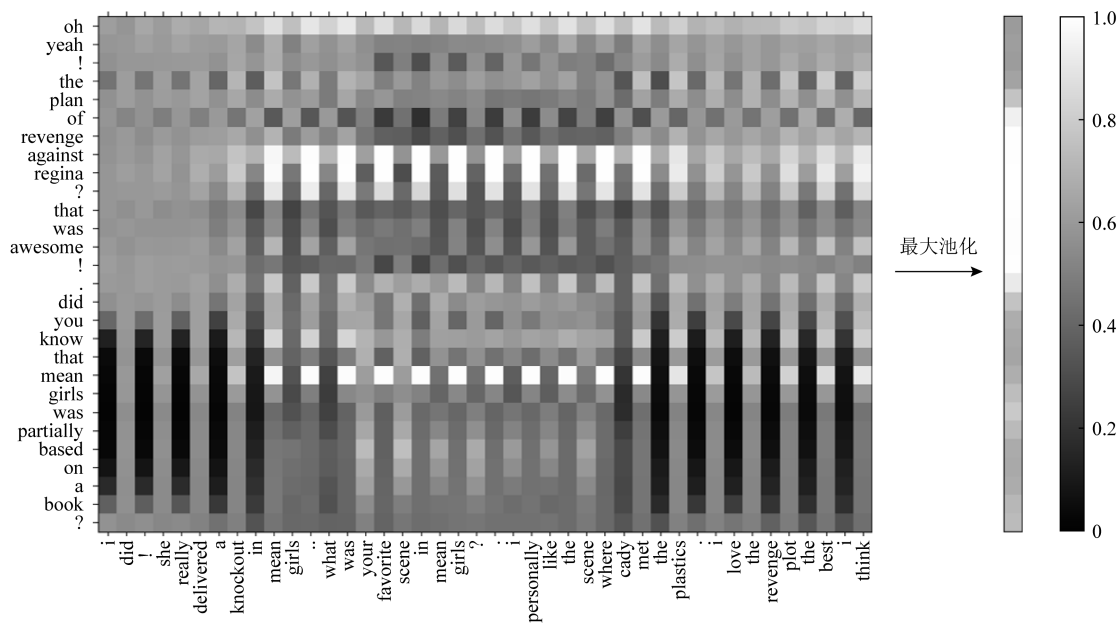


Fig. 3 Visualization results of the attention matrix

图3 注意力矩阵的可视化结果

而倾向于生成通用回复.然而,我们的模型却可以很好地做到这一点,并利用文档信息生成了正确且自然的回复,使得当前对话可以顺利地延续下去.由此便可体现出辩证看待对话历史的重要性.

Table 4 A Sample in the Test Set and Responses
Generated by Each Model

表4 测试集中的样例及各模型的生成回复

属性	内容
文档信息	genius, billionaire, and playboy tony stark, who has inherited the defense contractor stark industries from his father. ... he is captured and imprisoned in a cave by a terrorist group , the ten rings; ... but tony and yinsen know that raza will not keep his word.
对话信息	U1: yes i agree. i can't believe that rotten tomatoes gave a 7.7/10 i think it deserves a much higher score U2: yes me too right. it deserves a higher score. i love pepper potts acting U3: yes she was great too. i can't believe this movie came out in 2008 it was good enough even for todays age
S2S	it was a very good movie
S2SIF	i agree
T-MemNet	i think it was a great movie
ITE-DD	yeah, it was a great movie for sure
CITE-DD	i liked the scene where stark was captured by terrorists

注: 黑体文字为本文模型所参考的正文.

5 总结与展望

本文提出了一种新颖的基于 Transformer 的模型,以辩证的观点来看待对话历史,同时充分发挥当前对话在上下文构建和信息筛选中的指导作用,以应用到基于文档的对话任务.模型分为利用历史和忽略历史 2 种情况进行文档信息的筛选,然后根据当前对话与历史之间的相关程度,来决定最终是否包含利用历史的信息及其所占的比重.在公开数据集上的实验结果表明,与现有的基线模型相比,其可以更好地理解对话中的主题转移,从而生成在相关性和信息性上质量更高的回复.

至于后续的研究方向,考虑到目前并没有一种较好的方式,能够在长文档上对相关知识进行更为准确和精细的抽取.因此考虑文档本身的结构信息,同时引入粒度的概念可能是一个有效的改进思路.

参 考 文 献

[1] Sutskever I, Vinyals O, Le Q V. Sequence to sequence learning with neural networks [C] //Proc of the 28th Conf on Neural Information Processing Systems. Cambridge, MA: MIT Press, 2014: 3104-3112

[2] Gao Jianfeng, Galley M, Li Lihong. Neural approaches to conversational AI [C] //Proc of the 41st Int ACM SIGIR Conf on Research & Development in Information Retrieval. New York: ACM, 2018: 1371-1374

- [3] Ghazvininejad M, Brockett C, Chang Mingwei, et al. A knowledge-grounded neural conversation model [C] //Proc of the 32nd AAAI Conf on Artificial Intelligence. Menlo Park, CA: AAAI, 2018: 5110-5117
- [4] Zhang Saizheng, Dinan E, Urbanek J, et al. Personalizing dialogue agents: I have a dog, do you have pets too? [C] //Proc of the 56th Annual Meeting of the ACL. Stroudsburg, PA: ACL, 2018: 2204-2213
- [5] Madotto A, Wu Chien-Sheng, Fung P. Mem2Seq: Effectively incorporating knowledge bases into end-to-end task-oriented dialog systems [C] //Proc of the 56th Annual Meeting of the ACL. Stroudsburg, PA: ACL, 2018: 1468-1478
- [6] Vaswani A, Shazeer N, Parmar N, et al. Attention is all you need [C] //Proc of the 31st Conf on Neural Information Processing Systems. Cambridge, MA: MIT Press, 2017: 5998-6008
- [7] Zhou Kangyan, Prabhume S, Black A W. A dataset for document grounded conversations [C] //Proc of the 2018 Conf on Empirical Methods in Natural Language Processing. Stroudsburg, PA: ACL, 2018: 708-713
- [8] Moghe N, Arora S, Banerjee S, et al. Towards exploiting background knowledge for building conversation systems [C] //Proc of the 2018 Conf on Empirical Methods in Natural Language Processing. Stroudsburg, PA: ACL, 2018: 2322-2332
- [9] Dinan E, Roller S, Shuster K, et al. Wizard of Wikipedia: Knowledge-powered conversational agents [J]. arXiv preprint, arXiv:1811.01241, 2018
- [10] Rajpurkar P, Zhang Jian, Lopyrev K, et al. SQuAD: 100 000+ questions for machine comprehension of text [C] //Proc of the 2016 Conf on Empirical Methods in Natural Language Processing. Stroudsburg, PA: ACL, 2016: 2383-2392
- [11] Vinyals O, Fortunato M, Jaitly N. Pointer networks [C] //Proc of the 29th Conf on Neural Information Processing Systems. Cambridge, MA: MIT Press, 2015: 2692-2700
- [12] Bahdanau D, Cho K, Bengio Y. Neural machine translation by jointly learning to align and translate [J]. arXiv preprint, arXiv:1409.0473, 2014
- [13] Seo M, Kembhavi A, Farhadi A, et al. Bidirectional attention flow for machine comprehension [J]. arXiv preprint, arXiv:1611.01603, 2016
- [14] Wang Wenhui, Yang Nan, Wei Furu, et al. Gated self-matching networks for reading comprehension and question answering [C] //Proc of the 55th Annual Meeting of the ACL. Stroudsburg, PA: ACL, 2017: 189-198
- [15] Meng Chuan, Ren Pengjie, Chen Zhumin, et al. RefNet: A reference-aware network for background based conversation [J]. arXiv preprint, arXiv:1908.06449, 2019
- [16] Lian Rongzhong, Xie Min, Wang Fan, et al. Learning to select knowledge for response generation in dialog systems [C] //Proc of the 28th Int Joint Conf on Artificial Intelligence. San Francisco: Morgan Kaufmann, 2019: 5081-5087
- [17] Liu Zhibin, Niu Zhengyu, Wu Hua, et al. Knowledge aware conversation generation with reasoning on augmented graph [C] //Proc of the 2019 Conf on Empirical Methods in Natural Language Processing and the 9th Int Joint Conf on Natural Language Processing. Stroudsburg, PA: ACL, 2019: 1782-1792
- [18] Li Zekang, Niu Cheng, Meng Fandong, et al. Incremental transformer with deliberation decoder for document grounded conversations [C] //Proc of the 57th Annual Meeting of the ACL. Stroudsburg, PA: ACL, 2019: 12-21
- [19] Sankar C, Subramanian S, Pal C, et al. Doneural dialog systems use the conversation history effectively? An empirical study [C] //Proc of the 57th Annual Meeting of the ACL. Stroudsburg, PA: ACL, 2019: 32-37
- [20] Zhang Yangjun, Ren Pengjie, de Rijke M. Improving background based conversation with context-aware knowledge pre-selection [J]. arXiv preprint, arXiv:1906.06685, 2019
- [21] Zheng Yinhe, Zhang Rongsheng, Mao Xiaoxi, et al. A pre-training based personalized dialogue generation model with persona-sparse data [J]. arXiv preprint, arXiv:1911.04700, 2019
- [22] Xia Yingce, Tian Fei, Wu Lijun, et al. Deliberation networks: Sequence generation beyond one-pass decoding [C] //Proc of the 31st Conf on Neural Information Processing Systems. Cambridge, MA: MIT Press, 2017: 1784-1794
- [23] Xiong Hao, He Zhongjun, Wu Hua, et al. Modeling coherence for discourse neural machine translation [C] //Proc of the 33rd AAAI Conf on Artificial Intelligence. Menlo Park, CA: AAAI, 2019: 7338-7345
- [24] Luong Minh-Thang, Pham H, Manning C D. Effective approaches to attention-based neural machine translation [C] //Proc of the 2015 Conf on Empirical Methods in Natural Language Processing. Stroudsburg, PA: ACL, 2015: 1412-1421
- [25] Klein G, Kim Y, Deng Yuntian, et al. OpenNMT: Open-source toolkit for neural machine translation [C] //Proc of the 55th Annual Meeting of the ACL. Stroudsburg, PA: ACL, 2017: 67-72
- [26] Mikolov T, Sutskever I, Chen Kai, et al. Distributed representations of words and phrases and their compositionality [C] //Proc of the 27th Conf on Neural Information Processing Systems. Cambridge, MA: MIT Press, 2013: 3111-3119
- [27] Kingma D P, Ba J. Adam: A method for stochastic optimization [J]. arXiv preprint, arXiv:1412.6980, 2014



Sun Runxin, born in 1998. Master candidate. His main research interests include man-machine dialogue and natural language processing.
孙润鑫, 1998 年生. 硕士研究生. 主要研究方向为人机对话和自然语言处理.



Ma Longxuan, born in 1983. PhD candidate. His main research interests include man-machine dialogue and natural language processing.
马龙轩, 1983 年生. 博士研究生. 主要研究方向为人机对话和自然语言处理.



Zhang Weinan, born in 1985. PhD, associate professor, master supervisor. Member of CCF. His main research interests include man-machine dialogue and natural language processing.
张伟男, 1985 年生. 博士, 副教授, 硕士生导师, CCF 会员. 主要研究方向为人机对话和自然语言处理.



Liu Ting, born in 1972. PhD, professor. Senior member of CCF. His main research interest is natural language processing.
刘挺, 1972 年生. 博士, 教授, CCF 高级会员. 主要研究方向为自然语言处理.

《计算机研究与发展》2019 年论文高被引 TOP10

排名	论文信息
1	施巍松, 张星洲, 王一帆, 张庆阳. 边缘计算: 现状与展望[J]. 计算机研究与发展, 2019, 56(1): 69-89 Shi Weisong, Zhang Xingzhou, Wang Yifan, Zhang Qingyang. Edge Computing: State-of-the-Art and Future Directions [J]. Journal of Computer Research and Development, 2019, 56(1): 69-89
2	黄继鹏, 史颖欢, 高阳. 面向小目标的多尺度 Faster-RCNN 检测算法[J]. 计算机研究与发展, 2019, 56(2): 319-327 Huang Jipeng, Shi Yinghuan, Gao Yang. Multi-Scale Faster-RCNN Algorithm for Small Object Detection [J]. Journal of Computer Research and Development, 2019, 56(2): 319-327
3	彭宇新, 基金玮, 黄鑫. 多媒体内容理解的研究现状与展望[J]. 计算机研究与发展, 2019, 56(1): 183-208 Peng Yuxin, Qi Jinwei, Huang Xin. Current Research Status and Prospects on Multimedia Content Understanding [J]. Journal of Computer Research and Development, 2019, 56(1): 183-208
4	郑庆华, 董博, 钱步月, 田锋, 魏笔凡, 张未展, 刘均. 智慧教育研究现状与发展趋势[J]. 计算机研究与发展, 2019, 56(1): 209-224 Zheng Qinghua, Dong Bo, Qian Buyue, Tian Feng, Wei Bifan, Zhang Weizhan, Liu Jun. The State of the Art and Future Tendency of Smart Education [J]. Journal of Computer Research and Development, 2019, 56(1): 209-224
5	夏清, 李帅, 郝爱民, 赵沁平. 基于深度学习的数字几何处理与分析技术研究进展[J]. 计算机研究与发展, 2019, 56(1): 155-182 Xia Qing, Li Shuai, Hao Aimin, Zhao Qiping. Deep Learning for Digital Geometry Processing and Analysis: A Review [J]. Journal of Computer Research and Development, 2019, 56(1): 155-182
6	纪守领, 李进锋, 杜天宇, 李博. 机器学习模型可解释性方法、应用与安全研究综述[J]. 计算机研究与发展, 2019, 56(10): 2071-2096 Ji Shouling, Li Jinfeng, Du Tianyu, Li Bo. Survey on Techniques, Applications and Security of Machine Learning Interpretability [J]. Journal of Computer Research and Development, 2019, 56(10): 2071-2096
7	任家东, 刘新倩, 王倩, 何海涛, 赵小林. 基于 KNN 离群点检测和随机森林的多层入侵检测方法[J]. 计算机研究与发展, 2019, 56(3): 566-575 Ren Jiadong, Liu Xinqian, Wang Qian, He Haitao, Zhao Xiaolin. An Multi-Level Intrusion Detection Method Based on KNN Outlier Detection and Random Forests [J]. Journal of Computer Research and Development, 2019, 56(3): 566-575
8	赵志远, 王建华, 徐开勇, 郭松辉. 面向云存储的支持完全外包属性基加密方案[J]. 计算机研究与发展, 2019, 56(2): 442-452 Zhao Zhiyuan, Wang Jianhua, Xu Kaiyong, Guo Songhui. Fully Outsourced Attribute-Based Encryption with Verifiability for Cloud Storage [J]. Journal of Computer Research and Development, 2019, 56(2): 442-452
9	陈游旻, 陆游游, 罗圣美, 舒继武. 基于 RDMA 的分布式存储系统研究综述[J]. 计算机研究与发展, 2019, 56(2): 227-239 Chen Youmin, Lu Youyou, Luo Shengmei, Shu Jiwu. Survey on RDMA-Based Distributed Storage Systems [J]. Journal of Computer Research and Development, 2019, 56(2): 227-239
10	赵洪科, 吴李康, 李微, 张兮, 刘淇, 陈恩红. 基于深度神经网络结构的互联网金融市场动态预测[J]. 计算机研究与发展, 2019, 56(8): 1621-1631 Zhao Hongke, Wu Likang, Li Zhi, Zhang Xi, Liu Qi, Chen Enhong. Predicting the Dynamics in Internet Finance Based on Deep Neural Network Structure [J]. Journal of Computer Research and Development, 2019, 56(8): 1621-1631