

自然语言语义表示的范畴论建模:系统综述与组合机制分析

李奕丹, 崔建英, 熊明辉

引用本文

李奕丹, 崔建英, 熊明辉. 自然语言语义表示的范畴论建模:系统综述与组合机制分析[J]. 计算机科学, 2026, 53(4): 337-346.

LI Yidan, CUI Jianying, XIONG Minghui. [Category-Theoretic Semantic Representation: Systematic Review and Compositional Mechanism Analysis](#) [J]. Computer Science, 2026, 53(4): 337-346.

相似文章推荐 (请使用火狐或 IE 浏览器查看文章)

Similar articles recommended (Please use Firefox or IE to view the article)

[基于混合量子经典长-短距离特征扩展网络的图像分类](#)

Image Classification Based on Hybrid Quantum-Classical Long-Short Range Feature Extension Network

计算机科学, 2026, 53(4): 277-283. <https://doi.org/10.11896/jsjcx.250600108>

[NISQ量子线路高频-密集量子门集策略优化算法](#)

High Frequency-Dense Quantum Gate Set Optimization Algorithm for Quantum Circuit in NISQ Era

计算机科学, 2026, 53(4): 112-120. <https://doi.org/10.11896/jsjcx.241200213>

[基于变量子量的离散对数求解算法](#)

Variational Quantum Algorithm for Solving Discrete Logarithms

计算机科学, 2026, 53(1): 353-362. <https://doi.org/10.11896/jsjcx.241100181>

[分区稀疏攻击:一种更高效的黑盒稀疏对抗攻击](#)

Section Sparse Attack:A More Powerful Sparse Attack Method

计算机科学, 2026, 53(1): 323-330. <https://doi.org/10.11896/jsjcx.241200002>

[基于邻域匹配概率与类型商图的实体对齐解释方法](#)

Explanation Method for Entity Alignment Based on Neighborhood Matching Probability and Type Quotient Graph

计算机科学, 2025, 52(12): 260-270. <https://doi.org/10.11896/jsjcx.241100081>

自然语言语义表示的范畴论建模：系统综述与组合机制分析

李奕丹¹ 崔建英¹ 熊明辉²

1 中山大学哲学系逻辑与认知研究所 广州 510275

2 浙江大学数字法治实验室 杭州 310008

(liyid66@mail2.sysu.edu.cn)

摘要 语义表示是自然语言处理(NLP)的核心挑战。当前的语义表示范式可归纳为两类：以逻辑形式为核心的符号主义方法以及基于分布式表示的联结主义方法。尽管后者在工程应用中取得了显著成效,但在刻画语言的组合结构、支持结构化推理以及实现可解释与可泛化的语义建模方面,逐渐暴露出被称为“组合性危机”的理论局限。现有方法中,基于范畴论的组合分布语义模型凭借其严谨的代数结构和类型驱动的组合范式,为统一符号的句法结构与分布的语义内容提供了一条极具潜力的数学路径。对此,从范畴论的数学视角出发,以“范畴(理论框架)-组合(核心机制)-量子(计算范式)”为主线,对基于范畴论的自然语言语义表示范式及其演进脉络进行系统梳理与评述。不同于按模型或任务划分的既有综述,聚焦语义组合机制本身,首先基于组合视角对句子语义表示模型进行归类与比较,剖析分布式语义方法在组合建模中的代表性进阶及其内在局限,进而梳理其向组合分布语义发展的内在逻辑和研究趋势。在此基础上,重点阐述以字符串图为演算工具的范畴组合语义框架,并结合典型模型(如 DisCoCat 与 DisCoCirc)说明这类框架的形式化特征及其在量子计算语境下的扩展方向,为理解和评估符号主义方法、联结主义方法与量子计算方法在自然语言处理中的融合路径提供统一的理论视角。

关键词 范畴论;字符串图;组合语义;量子计算;可解释性

中图分类号 B819

Category-Theoretic Semantic Representation: Systematic Review and Compositional Mechanism Analysis

LI Yidan¹, CUI Jianying¹ and XIONG Minghui²

1 Institute of Logic and Cognition, Sun Yat-sen University, Guangzhou 510275, China

2 Digital Rule of Law Laboratory, Zhejiang University, Hangzhou 310008, China

Abstract Semantic representation is a central challenge in natural language processing (NLP). Existing approaches can be broadly categorized into two paradigms: symbolic and connectionist methods. Although the latter have achieved remarkable practical success, they suffer from theoretical limitations—commonly referred to as the “compositionality crisis” in compositional modeling and semantic interpretability. In existing methods, categorical compositional distributional semantics provides a principled mathematical framework for unifying symbolic syntactic structure with distributed semantics via type-driven composition. From a categorical perspective, this paper surveys category-theoretic approaches to semantic representation along the conceptual line of “category theory-composition-quantum computation”. Unlike surveys organized by models or tasks, it focuses on semantic composition mechanisms, comparing sentence-level models from a compositional viewpoint, analyzing the limitations of distributed approaches, and outlining the theoretical shift toward compositional distributional semantics. Building on this, string diagram-based frameworks such as DisCoCat and DisCoCirc are presented, clarifying their formal properties and quantum extensions, offering a unified view of symbolic, connectionist, and quantum semantics.

Keywords Category theory, String diagrams, Compositional semantics, Quantum computing, Interpretability

1 引言

自人工智能(AI)诞生以来,自然语言处理(NLP)领域的发展经历了从侧重于逻辑和结构化知识的符号主义到依赖大

规模数据和统计模式的联结主义范式转换。进入 20 世纪 90 年代,以统计和数据驱动为核心的分布式语义建模逐渐成为主流范式。特别地,随着神经网络技术的发展,分布式表示逐渐从基于计数的方法演进为基于预测的词嵌入模型,例如

到稿日期:2025-10-28 返修日期:2026-01-26

基金项目:国家社会科学基金重大项目(19ZDA042)

This work was supported by the Major Program of National Social Science Foundation of China(19ZDA042).

通信作者:崔建英(cuijiany@mail.sysu.edu.cn)

Word2Vec^[1]和 GloVe^[2]以及近年来被广泛应用的上下文预训练语言模型 ELMo^[3]和 BERT^[4]等。这类方法在词相似度、语义相关性等词级语义任务中取得了显著成效。以 Chat-GPT 为代表的大规模预训练模型,通过在模型规模和训练数据上的持续扩展,进一步推动了联结主义方法在 NLP 领域占据主导地位。

分布式表示模型在 NLP 任务中表现出色,但其依赖隐式统计学习的特性,使其内部推理不透明,可解释性与结构化语义建模能力不足。为突破此瓶颈,研究者将组合性原则引入分布语义框架,以刻画复杂表达式的意义。早期的向量加法、逐点乘积等方法,因运算的交换性而限制了建模精度^[5-6];近年来,范畴论(Category Theory)作为连接“语言理论推理”与“数值计算”的数学工具,凭借其统一符号结构与分布式语义的优势而受到关注^[7-8]。基于范畴论的语言模型能够将句法到语义的推导过程形式化为张量网络上的代数运算,并借助字符串图(String Diagrams)实现语义组合的可视化与可计算化。更重要的是,以范畴论为核心的组合语义建模路径,既为联结主义模型的理论重构提供了全新视角,又为构建兼具可组合性、可解释性与量子优势的语义模型奠定了理论基础,在

基础研究与实际应用层面均具有重要的学术价值与广阔的发展前景。

本文以“范畴-组合-量子”框架为主线,通过梳理基于范畴论的语义建模方法在经典计算与量子计算语境中的发展脉络及其核心架构,揭示出以范畴论为核心的组合语义建模路线是融合符号主义、联结主义与量子计算的语义建模新范式,并通过实例验证这类语义建模进阶可为解决 NLP 中计算复杂度、可解释性弱与泛化能力不足等问题提供新的理论依据与实现路径。

为清晰地阐释该研究范式的内在结构演进,本文依时间顺序对该进阶发展的理论脉络进行梳理,并结合“基于范畴论的组合语义建模及其量子实现路径”的整体分析视角(对应图 1 所示的层次化框架)展开论述:第 2 和第 3 章围绕符号层进行论述,其中第 2 章主要介绍组合视角下,句子意义建模方法及其存在的问题,第 3 章概述范畴论作为语言计算工具的理论演进;第 4 章介绍了语义组合层工作,梳理了字符串图表示的意义组合的机制、瓶颈与研究前沿;第 5 章介绍了量子线路的扩展,从组合层映射到量子计算层的数据表示的质量直接关系任务应用层。

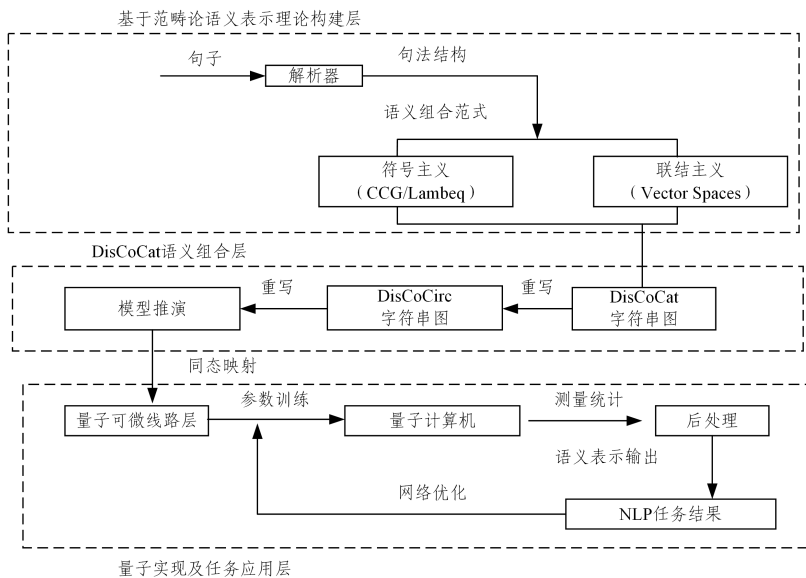


图 1 从 DisCoCat 到量子 NLP 实现的完整框架

Fig. 1 General framework of DisCoCat-based quantum NLP implementation pipeline

2 句子表示模型的演进与分类

词汇和句子的分布式表示一直是语言理解与人工智能交叉领域研究的核心议题。以 Word2Vec、GloVe 以及 BERT 等为代表的词嵌入方法,通过从大规模语料中学习连续向量表示,显著推动了分布式语义建模的发展。随着词嵌入技术的日趋成熟,研究重心逐渐由词级表示扩展至句子层面的语义表征,后者需涵盖更高层次的全局语义与句法结构特征。这一转向在语义表示理论与建模方法方面引发了广泛关注,并在近些年于理论与应用层面取得了丰富进展。

文献[9]指出,句子嵌入的核心问题在于如何通过适当的组合机制,将词层级的语义信息整合为句层级的意义表示。现有研究根据组合策略大致可分为 6 类:词集模型、语用连贯

方法、语义比较模型、混合多任务技术、特殊结构模型和预训练上下文模型。这一分类框架揭示了不同研究对“单词与句子之间的联系”以及“句子在人类语言理解中所扮演角色”的不同理解。为更全面地展现句子表征方法在语义组合机制与模型结构方面的发展脉络,对该分类体系进行了扩展(见表 1)。

如表 1 所列,现有句子组合模型可划分为两类:依赖符号化的逻辑推理与侧重于统计分布式的表示学习。前者通过显式的句法或语义结构对句子进行建模,强调符号之间的逻辑组合与规则约束,语义生成过程依赖预定义的推理规则。后者则通过在连续向量空间中以稠密分布形式表示词汇与句子,实现从离散符号到连续表征的转换,其语义信息主要源于对上下文关联的统计建模与参数优化,语法关系被“隐含”在

模型数十亿个权重参数中。通过在大规模语料库上进行自监督训练,这类分布式模型能够从词序列共现概率中学习语言规律,生成基于概率分布的语义表示,并已成功应用于多项 NLP 任务中。

然而,分布式模型的优势并非源于逼近“绝对理性”,而是其生成机制本质上是一种启发式搜索,能在高维空间中生成形式上“令人满意”的高概率样本。面对任务指令时,模型并非真正“理解”语义,而是通过统计关联生成在概率上最可能的近似解。这种任务驱动的优化目标,导致其生成的语义表

示在稳定性与一致性上存在不足。其深层原因在于,词与句子的语义形成机制存在本质差异:词嵌入通常通过上下文共现刻画局部语义,而句子语义不仅受外部语境影响,更受到内部句法结构的强约束,即句子的意义既是语境分布的产物,也是句法组合的结果。相比之下,符号主义模型具有内在驱动、演绎式的求解特征,能够显式追踪语义组合的逻辑路径。因此,如何在结构化与分布特征之间建立协调机制,有效融合分布式模型的统计学习能力与符号主义模型的结构化推理特性,成为当前研究的关键问题。

表 1 句子嵌入方法及其分类总结

Table 1 Summary of the sentence-embedding methods and their categories

3 类组合视角	5 类组合模型	核心机制	代表性模型
词汇驱动	词集模型	忽略词序,将句子视为词向量集合,通过加权求和、连接或者复杂神经网络的非线性映射组合词向量生成句子表示	WR ^[10] , p -mean ^[11] , PV ^[12] , SDAE ^[13] , MeanMax ^[14] , Ditto ^[15]
	语用连贯方法	通过建模相邻句子之间的语用连贯关系学习句子嵌入。侧重于学习语篇结构和上下文关系,依赖大规模语料以刻画上下文连续性	Skip-Thoughts ^[16] , Quick-Thoughts ^[17]
句子驱动	语义比较模型	通过理解和判断句子对之间的逻辑关系(如蕴涵、矛盾)学习句子表示,强调语义对比与推理能力	DisSent ^[18] , Discovery ^[19] , InferLite ^[20] , Sentence-BERT ^[21]
	混合多任务技术	综合多目标函数(如连贯性、语义比较等)共同训练,学习到的句子表示同时刻画多种语义和语用信息	DiscSent ^[22] , MILA/MSR ^[23] , USE ^[24] , Adapters ^[25]
结构建模驱动模型	显式结构模型	利用句子的内部结构(如句法解析树),显式建模句法进行句子建模,自底向上组合词义。	RNTN ^[26] , Tree-LSTM ^[27] , Hierarchical CNN ^[28] , DisCoCat ^[29]
	上下文预训练模型	通过大规模自监督预训练,学习上下文依存的词表示,动态建模语境信息,并通过池化或微调获得句子级表征	ELMo ^[3] , BERT ^[4] , ALBERT ^[30] , GPT ^[31]

3 范畴论的计算语言学转向

范畴论为语言建模中的根本挑战,即如何为“意义”建立可计算、可组合的形式化描述,提供了数学工具。本章通过梳理范畴论从哲学渊源、语法演算、语义表示到图形化编译的演进路径,阐明范畴论如何逐步发展为衔接句法与语义的统一建模框架。

范畴论的思想起源可追溯至亚里士多德对事物分类的传统,并在康德《纯粹理性批判》中被赋予认知论意义;范畴作为组织经验的先验形式,为“类型先于认识”的原则奠定了基础。此思想随后在类型论中得到形式化表达^[32],并直接催生了范畴语法(Categorical Grammar, CG)的提出^[33-34]。该理论奠定了“组合性”作为语言建模的核心原则,也标志着范畴从哲学概念向形式语法工具转变。

文献[35]在数学框架下定义了“范畴”与“函子”的概念,将函子作为范畴间的结构保持映射,使得句法推导可转化为路径组合操作? 范畴论成为一门关于“关系结构”的科学,推动了范畴论从具体对象描述转向结构关系抽象表征的研究范式。进一步,文献[36]将逻辑系统本身视为范畴,模型则定义为从逻辑范畴到其他范畴的结构保持函子,提出函子语义学,确立了“逻辑可范畴化”的语义研究范式,使范畴论成为数理逻辑与语义分析的重要桥梁,为形式语义学奠定了统一的数

学基础。而文献[37]通过兰贝克演算(Lambek Calculus)将语法推导与证明树结构建立同构关系,并基于柯里-霍华德同构(Curry Howard Isomorphism),将 AB-演算系统扩展为幺半范畴(Monoidal Category),使函子类型具有左右伴随性质。语法树编译为 λ -项,实现“逻辑-计算-范畴”的三位一体;逻辑公式和类型系统被视为对象,证明与程序和范畴中的态射对应。基于此,文献[38]证明词向量的组合语义可通过将预群语法(Pregroup Grammar)映射到一个幺半范畴的实例,实现符号语义与分布式语义的范畴化融合,完成语法层的范畴化,即从“AB-演算—Lambek 演算—Pregroup 语法”的过渡。

在文献[39]发展的字符串图语言中,句法树可被编译为二维图形,语法组合过程直接对应张量网络缩并。这种图形编译系统完成了从句法到语义平行推演的端到端显式编译,并通过 Lambeq^[40]和 DisCoPy^[41]等开源工具包实现自动化,支持在量子硬件上运行范畴化的 NLP 任务,标志着范畴论可从“元语言”编译为“工程语言”。至此,范畴论凭借其高度抽象与结构保持的特性,形成了一套能够贯通语法、语义与计算的形式化工具体系。其中,作为范畴论的表达工具的字符串图,也是跨领域的计算语义载体。

表 2 明确展示了自然语言、逻辑推理与计算过程之间的深层同构关系,使得语言的组合性、推理性与可计算性得以在同一数学框架下表达¹⁾。

¹⁾ 关于该同构关系的严格范畴论证明,可参阅文献[37]的原始论述。

表 2 代数、证明论、范畴语法及程序语言之间的同构关系

Table 2 Isomorphisms among algebra, proof theory, categorical grammar, and programming languages

范畴论	逻辑学	语言学	Python 实现
双闭范畴	证明系统	范畴语法	DisCoPy
对象	公式	类型	biclosed. Ty
生成元	公理	词库/词典	biclosed. Box
态射	证明树	规约	biclosed. Diagram

4 基于范畴论的组合语义模型及其实现

第 3 章揭示了范畴论如何以统一的“结构-意义-计算”框架关联句法形式与语义推理。本章将进一步转向具体模型的构建。首先,引入字符串图作为范畴论的图形演算工具,使语法组合可在语义空间中以张量网络的形式执行。进而,引入组合语义模型的公理化定义:通过结构保持的函子将句法范畴映射至语义范畴,保证组合性的计算一致性。作为组合模型的实例,本章剖析两个代表性模型:专注于短语与句子层级组合语义的 DisCoCat 模型,以及将其组合性原则扩展至语篇与上下文建模的 DisCoCirc 模型。

4.1 字符串图的形式化语义与范畴结构

在语言的形式语义学领域,研究者通常采用范畴论的数学原理及其配套的字符串图语言,对组合性原则进行形式化刻画。其中,字符串图具有直观且严格的图形约定:线表示对象;方框表示态射;顺序复合(\circ)通过垂直连接表示,对应常规函数复合;并行复合(\otimes)通过平行连接表示,对应张量积操作;单位对象 I 表示空输入与空输出结构,图形读写顺序遵循自左至右的约定。

在幺半范畴的语义定义下,字符串图为模型架构提供几何语法,包括恒等态射 (Identity)、复制 (Copy)、删除 (Delete) 与交换 (Swap) 等的过程图形表示,如图 2 所示,“复制”与“删除”等基本操作分别表示分支与终止结构,对应信息流的指派与消失。

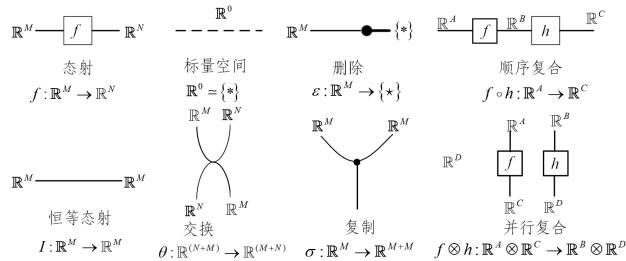


图 2 字符串图的图示语法

Fig. 2 String-diagrammatic syntax

作为字符串图语义模型的核心构件,恒等、复制、删除与交换 4 类普适操作适用于任意语义空间,构成了一套兼具直观性与形式严谨性的图形化语言。这不仅统一了语言结构的语法-语义接口,也因其与量子态操作(如张量积、交换、丢弃、复制的受限版本)具有天然结构同构性,为构建面向量子实现的语义模型奠定了基础。

自然语言组合语义形式化建模的核心挑战在于建立一种既能清晰分离抽象语法结构与具体语义解释,又能保持结构统一性的映射机制。范畴论所提供的组合模型恰好满足这一

需求,其为解决语义组合性问题提供了通用的数学框架。文献[38]证明,词向量的压缩过程(SVD)可以形式化为 Lambek 范畴的一个“紧闭范畴 (Compact Closed Category, CCC)”实例。在此框架中,句法范畴被映射到有限维 Hilbert 空间 (FVect),斜杠类型对应张量积的对偶操作,语法树平铺为张量缩并。

在字符串图中,对象被解释为特定的张量空间。例如,态射 $f: \mathbb{R}^M \rightarrow \mathbb{R}^N$ 对应于标量域上 N 维与 M 维向量空间之间的映射,其中 \mathbb{R}^M 表示所有 M 维实向量构成的空间。通过笛卡尔积自然扩展(笛卡尔积空间 $\mathbb{R}^{A \times C}$ 表示所有 $A \times C$ 实矩阵构成的空间),在线条上引入刻度以标注索引。该表示法同样可以表示高阶张量空间,例如 $\mathbb{R}^{A \times B \times C}$ 中的张量可分解为由带标注线条表示的张量积序列。如图 3 所示,等式与其对应的字符串图在形式上是严格等价的。因此,作用于张量空间中的多线性函数可以等价表示为字符串图。

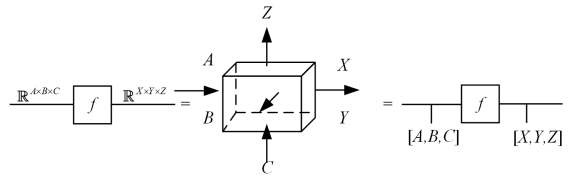


图 3 SMC 中张量算子的字符串表示

Fig. 3 String diagrams representation of tensor operators in SMCs

定义 1 [36] 一个组合模型 M 是一个三元组 $M = (G, C, [-])$ 。其中,符号集 G 定义了模型的抽象语法,包含了对象集、态射集以及等式集,通过它自由生成范畴 $S = Free(G)$,即模型的所有合法组合结构; C 表示语义范畴;函子 $[-]: S \rightarrow C$ 。

定义 1 中,结构范畴 S 是由类型集 G 自由生成的范畴,对象为 G 中的变量列表,其态射正是由 G 中的生成元通过顺序和并行复合构造出的所有可能的字符串图。语义范畴 C 提供对模型对象与过程进行具体实例化的语义空间。例如,在向量空间语义模型中, C 可以是欧几里德空间及其连续函数构成的 CMC。 $[-]$ 是一个从结构范畴 S 到语义范畴 C 的结构保持映射,将 S 中的抽象变量(线)映射为 C 中的具体对象(如向量空间),并将生成元(框)映射为 C 中的具体态射(如 FVect 或者量子线路)。二者之间通过一个保持结构的函子相连,从而形成函子语义学:句法的组合结构被函子映射为语义的组合结构。

4.2 基于范畴论的组合模型实例

本节以自然语言处理中的 DisCoCirc 与 DisCoCat 模型为例,阐释“组合模型”的核心概念。其结构(由对象与生框表示)定义了模型的组合逻辑,语义范畴的选择(如向量空间范畴 FVec)决定了字符串图中组件的具体语义解释(如线性变换)。

4.2.1 DisCoCat: 从句法范畴到向量空间的映射

在 DisCoCat 模型中,通过引入语法,给每个词汇指派特定的类型(例如名词或者动词),并通过语法规则进行组合,获得短语或者句子的语义表示。Lambek 的“预群”概念为该模型提供了一个可计算自然语言的句法结构的范畴化语法形式。预群是一个严格的幺半范畴,又是一个偏序集,即对于对

象 x, y 至多具有唯一的态射 $x \rightarrow y$ 。预群语法在代数结构上是偏序幺半群,其中每一个元素 x 都有一个左/右伴随,用于刻画词项的组合归约过程。

定义 2^[42] 一个预群语法是一个具有幺半结构 $(1, \cdot)$ 的偏序集 (G, \leq) ,使得每一个元素 x 都有一个伴随函子,即 x^l 和 x^r 满足:

$$x^l \cdot x \leq 1 \leq x \cdot x^l, x \cdot x^r \leq 1 \leq x^r \cdot x \quad (1)$$

上述不等式展示了语法的运算规则。

在语义表示中,一个预群语法是由每一个带有语法类型的词汇 $w:t$ 组成的列表 Σ ,其中类型 $t \in G(B)$ 。原子类型集 $B = \{n, s\}$,其中 n 和 s 分别代表名词和句子。那么生成的自由预群 $G(B)$ 中的类型,例如 $n, s, n^r \cdot s \cdot n^l$ 等,就可以作为它的元素。当类型为 t_1, t_2, \dots, t_n 的词汇序列的类型 $t = t_1, \dots, t_n$ 满足 $t \leq s$,则通过语法推导生成的句子就是一个合法语句。例如,句子“小明喜欢小红”的预群推导:

$$n \cdot n^r \cdot s \cdot n^l \cdot n \leq 1 \cdot s \cdot 1 \leq s \quad (2)$$

其推导图如图 4 所示。

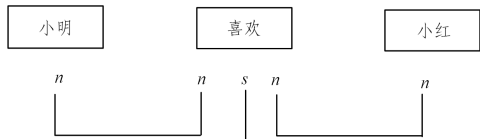


图 4 句子“小明喜欢小红”的预群语法推导图

Fig. 4 Pregroup grammar diagram of the sentence “Xiaoming likes Xiaohong”

为将预群语法所给出的句法推导赋予可计算的语义解释,需要引入具体的组合语义模型,将预群范畴中的句法结构通过函子映射到相应的语义范畴中。在 DisCoCat 模型中,常取有限维向量空间范畴 FdVect 作为语义范畴,此时句法推导可被函子化地表示为字符串图,从而得到图 5 所示的语义组合结构。

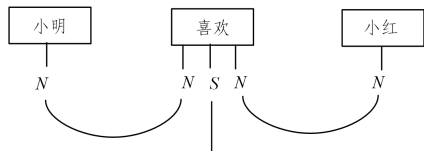


图 5 句子“小明喜欢小红”的字符串图

Fig. 5 String diagram of the sentence of “Xiaoming likes Xiaohong”

这个字符串图由方框中的词汇与预群类型推导决定词汇所处的向量空间。语法范畴 S 中,对于词汇列表 $(w:t) \in \Sigma$,函子都会指派一个语义 $\llbracket w \rrbracket$,作为语义范畴 C 的一个态射,即从单位对象 I 到类型 t 的语义 $\llbracket t \rrbracket$ 态射。其中, N, S 是向量空间,“小明”和“小红”两个名词由 N 中的向量表示,及物动词“喜欢”的语法类型 $n^r \cdot s \cdot n^l$,它的语义 $\llbracket vt \rrbracket$ 就是 $(I \rightarrow \llbracket n^r \cdot s \cdot n^l \rrbracket) \in C$ 。由于函子能够保持幺半结构, $\llbracket n^r \cdot s \cdot n^l \rrbracket \simeq \llbracket n \rrbracket^r \otimes \llbracket s \rrbracket \otimes \llbracket n \rrbracket^l$,因此,一个及物动词对应 FVec 中的三阶张量,在字符串图中对应于 $N \otimes S \otimes N$ 中的三阶张量。基于以上两个推导图结构上的相似性以及预群语法推导对张量及其收缩运算的形状的决定程度,在 DisCoCat 模型中使用预群语法的优势就变得十分清晰了。

在语义范畴 C 中,用杯子来定义张量网络中的张量收

缩,以确保范畴中的态射可以通过张量收缩实现语义组合。

杯子通过空间 $V \otimes V$ 中的状态:

$$\sum_{i=1}^n |i\rangle \otimes |i\rangle \quad (3)$$

帽子则通过映射:

$$|i\rangle \otimes |j\rangle \mapsto \begin{cases} 1, & i=j \\ 0, & \text{otherwise} \end{cases} \quad (4)$$

在该语义范畴中,每个词被映射为一个张量,如句子“小明喜欢小红”的语义可表示为:

$$| \text{小明喜欢小红} \rangle_k = \sum_{i,j=1}^n | \text{小明} \rangle_i | \text{喜欢} \rangle_{i,k,j} | \text{小红} \rangle_j \quad (5)$$

其中,张量收缩反映了语法依存关系引导的语义计算。

基于字符串图的自然语言意义组合模型的理论基础可追溯至文献[29]提出的 DisCoCat 模型的数学基础。其核心思想是将句子的语法推导过程精确地转化为向量空间中的张量代数运算。范畴语法与双闭范畴之间的联系为使用范畴论重新构建了许多语法的语言模型——常规语法、上下文无关语法、预群语法等^[43],然后将其融入 DisCoCat 更广泛的组合分布模型。文献[44]详细阐述了从预群语法到更具覆盖能力的 CCG 的过渡,推动了广覆盖组合 NLP 工具 lambeq 的诞生。

近年来出现的高阶 DisCoCat 研究,通过引入更丰富的范畴结构,扩展了传统模型对词汇与句法类型的刻画能力,使其能够更灵活地处理长距离依赖、语义歧义以及否定、量词等复杂语言现象,并提升了组合语义建模的表达能力^[42]。

4.2.2 DisCoCirc 模型

尽管早期的 lambeq 专注于句子层面的语义处理,但范畴论的组合原理使其能够扩展到语篇层面。DisCoCirc 模型正是基于此发展而来,该模型旨在处理动态话语表征。

文献[45]介绍了动态话语表征模型 DisCoCirc。当句子的预群语法树构建完成后,其语法结构将进一步被转换为一系列嵌套的图框,获得句子级的字符串图。在 DisCoCirc 框架中,文本结构通过由状态、方框和框架生成的字符串图来建模。其中,状态用于表示文本中出现的名词,状态中所编码的信息沿着图中的连线进行传递,并通过方框进行更新。这些方框对应于文本中的单词,用于更新状态或促使不同状态之间发生交互。例如,及物动词描述两个名词之间的作用关系,形容词则用于对单个名词进行语义限定。图 6 展示了句子级的 DisCoCirc 字符串图。

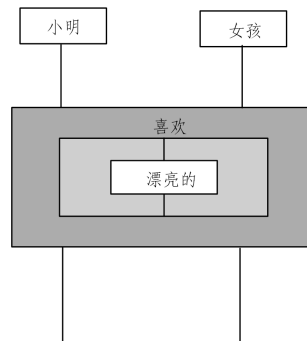


图 6 句子“小明喜欢漂亮的小红”的 DisCoCirc 字符串图

Fig. 6 DisCoCirc diagram for the sentence “Xiaoming loves beautiful Xiaohong”

上述 DisCoCirc 字符串图引入了高阶态射¹⁾的概念。在字符串图中,这类高阶过程通常被表示为带有可插入的留有“空洞”的方框所示,作用于方框、子图或其他框架。如图 7 中的形容词“确实”方框,空洞部分可以嵌入其他语义成分图,从而对其语义进行动态变换。

这种基于高阶态射与可插入字符串图的动态组合机制,使模型能够以字符串图形式刻画复杂的句法嵌套与语义依赖关系,生成被称为文本线路的结构化语义表示。类似于自然语言中句子可以组合成篇章的方式,DisCoCirc 框架中的字符串图也可沿着同一名词对应的连线进行组合,形成更大规模的语义结构。在获得句子级线路后,将线路依次进行组合,得到整个文本对应的整体线路(见图 7)。由于这些线路是自上而下读取的,因此较早出现的句子位于上方。

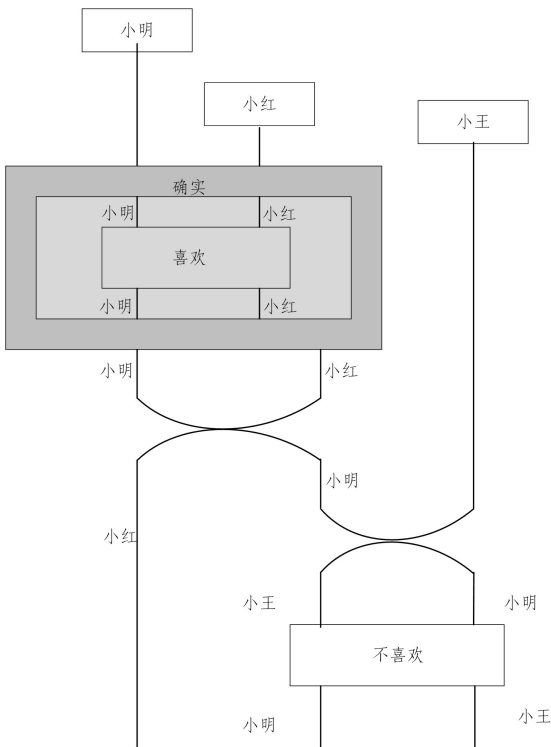


图 7 句子“小明确实喜欢小红,但小红不喜欢小明”的字符串图
Fig. 7 String diagram of the sentence of “Xiaoming really likes Xiaohong, but Xiaowang dislikes Xiaoming”

从范畴论与字符串图的角度,形式化刻画“文本线路”为一种可组合、可生成的结构化语义表示。文本线路是一个具有 n 个输入线和 n 个输出线的字符串图,这些输入和输出分别对应于文本中不同的语言成分。令 Σ 为一个带类型的词汇列表,即 $w:t. TextCirc_{\Sigma}(n)$ 表示所有可能的字符串图的集合。在该集合中,生成元经过有限多个步骤组合插到对应的“空洞”中,直到不留下可插入的方框,则生成了一个完整的文本线路图。

DisCoCirc 模型是对 DisCoCat 模型的扩展,旨在解决 3 个核心问题。

- 1) 如何组合句子的语义获得文本的意义?
- 2) 当获得新的知识时,单词意义如何更新?

3) 句子的意义空间是什么类型?

在 DisCoCat 中,每个词在词典中的意义都是固定的(在 DisCoCirc 模型中,每个名词指称对象指派独立的语义空间,在字符串图中表示为“并行”的线),高阶态射支持单词意义的动态更新。全局语义通过线路图拓扑结构(顺序复合或者嵌套)组合获得,单词意义的更新是通过高阶态射的可组合性实现的,即动态插入和状态传递。新知识插入到高阶方框后,语义获得更新,而每个态射框作为一个“可输入”的结构,接收句子作为输入/输出(I/O)过程相互连接。DisCoCirc 框架下的语篇级字符串图结构既可以接收前序名词状态,又可以输入到后续线路。近年来,Coecke 团队又补充了更多细节^[46]。

5 量子化扩展

范畴论结构与量子系统所依赖的线性代数形式之间具有天然的相容性,这为量子自然语言处理(QNLP)提供了稳固的理论基础。在该框架下,语言推导过程可被形式化为量子线路,词汇的向量表示对应量子态,而语法组合则通过张量积与线性变换加以实现。其中,伴随函子在形式上统一了语法中的“句法-类型消除”与量子线路中的“纠缠-消解”过程,使得语言的句法推导能够以结构保持的方式映射为量子计算中的图形操作,构成连接语言组合机制与量子计算模型的核心枢纽。

在 DisCoCirc 框架中,字符串图的每条线被指派为固定数量的量子比特(Quantum Qubit),名词被编码为量子态,而方框则对应作用于这些量子态上的幺正运算。将语义范畴从 $FVec$ 扩展至复希尔伯特空间范畴 $Philb$ 后,语义基底扩展至复数域,词项与句子的语义表示不再局限于实数值张量,而是可表征为复向量空间中的量子态,以模拟语义的纠缠与叠加等量子现象。仍然以“小明喜欢小红”这个语句为例,图 8 展示了语义函子 $\llbracket - \rrbracket$ 将 DisCoCirc 字符串图的原子组件转换为参数化量子线路元素的过程。

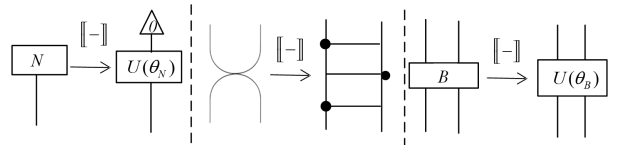


图 8 DisCoCirc 原子组件到参数化量子线路元素的语义函子 $\llbracket - \rrbracket$
Fig. 8 Semantic functor $\llbracket - \rrbracket$ from atomic components of DisCoCirc to parameterized quantum circuit elements

DisCoCat 结合了语义和语法,在量子硬件上表示为不同的状态和测量,如图 9 所示。

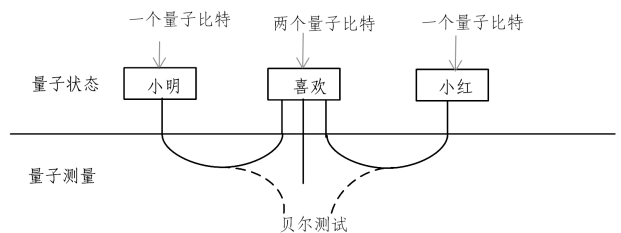


图 9 语法结构与量子测量、词义与量子态的对应关系
Fig. 9 Correspondence between syntactic structure and quantum measurement, and between lexical semantics and quantum states

¹⁾ 高阶态射是指从一组具有固定类型的输入态射产生新态射的函数。

在上述映射关系中,语法结构被解释为量子测量,词汇意义被解释为量子态,从而实现了句子层面的语义组合。然而,自然语言的意义生成是层级化的,即它不仅依赖于句内结构,还包括句间与语篇层面的语义关联和层级嵌套。为此,进一步探讨如何在量子环境中表征语篇框架,并给出其对应的量子线路表示,文献[47]提出通过语义函子构建框架及其内部组件的线路表示。

由于框架的内部状态与其输出状态未必构成双射关系

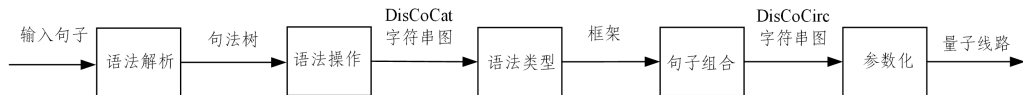


图 10 基于 DisCoCirc 模型的语篇级 QNLP 计算流程

Fig. 10 Discourse-level QNLP computational workflow based on the DisCoCirc model

1) 选定组合模型,通过调用一个统计解析器来获得句子的句法树。

2) 使用 CCG 替换预群语法及其实践操作,将句法树转换成字符串图。

3) 根据语法类型确定不同层级的框架,框架中每个嵌套的子线路均匹配独立的算子,以确定内部的组合机制。

4) 根据参数化酉算子进行句子组合,获得 DisCoCirc 字符串图。

5) DisCoCirc 字符串图可以转换成量子线路,这一转换基于特定的参数化量子线路选择。

上述步骤完成后,参数化量子线路的输出态即可作为训练数据。该线路经量子编译器翻译为具体的量子门操作序列,再提交至量子计算设备运行。

在量子计算框架下,上述过程体现了一种结构保持的同态关系:语言中的句法组合结构,经由范畴论函子映射,被系统地转化为字符串图中的态射组合,并进一步映射为量子线路中的并行与顺序操作。具体而言,每个词项被指派为量子态或参数化量子门,句法推导中的类型消除在量子层面对应于通过 Bell 测量实现的态消解与纠缠生成。该映射在保持组合顺序与并行结构不变的前提下,实现了从语法结构到量子线路的统一表示,从而保证了语言推导、语义组合与量子计算过程在结构层面的一致性。整个句子的语义可以表示为一个参数化量子线路:由预群语法规则映射为量子线路的拓扑结构;语义信息被编码在量子线路的可调参数之中;句子整体表示对应于线路末端的输出量子态,测量结果可被解释为句子的真值或语义向量。

这种“句子被编码为量子线路”的机制,使得语言的组合性在量子层面得以保持,并保留了 DisCoCat 模型中“句法控制语义组合”的核心思想。如图 11 所示,句中 3 个词语分别对应 3 个量子态,语法约简通过量子纠缠实现语义组合,最终在输出线(对应句法类型 s)上得到句子的整体语义表示。

图 11 展示了将自然语言句子编码为量子线路的完整流程。第一阶段(深灰色框)根据预群语法完成句法推导,将词汇的语法类型初始化为量子比特状态。第二阶段(浅灰色框)执行贝尔测量,以此实现语义组合并提取句子的最终语义信息。

(如在连词或反身动词等情形下,I/O 的语义关联通常并非是一一对应),因此需引入基于参数化酉算子(Unitary Operators)的框架表示方法,以在量子层面保留其嵌套结构与语义关系。具体而言,每个框架被指派到一对参数化酉算子,作用于其输入态,算子权重通过训练任务学习。框架内各嵌套子线路对应独立算子,使模型能够自主决定框架对其内部组件的作用方式。图 10 展示了句子处理的 5 个主要阶段。

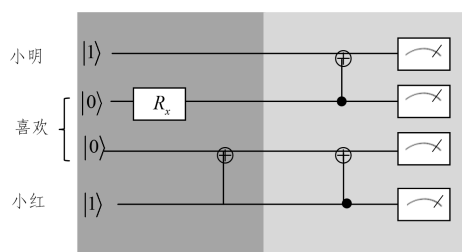


图 11 句子“小明喜欢小红”的量子线路编码

Fig. 11 Quantum circuit encoding the sentence “Xiaoming likes Xiaohong”

从字符串图到量子线路的函子化映射过程中,依据的形式化定义与编码规则可参考文献[48]提出的量子自然语言处理框架。该框架系统阐述了句法-语义结构向量子线路的函子化映射机制、参数化量子线路的构造方法,以及 Bell 测量在语义组合过程中的作用。

为更清晰地展示范畴 QNLP 模型研究的整体图景,并对不同方法的理论基础与应用定位进行对比,对相关工作进行了系统性汇总,如表 3 所列。

表 3 现有基于范畴 QNLP 方法的比较概览

Table 3 Comparative overview of existing QNLP approaches

类型	任务	主要工作
受量子启发的方法 (经典计算机硬件)	信息检索	Sordani et al. (2013) ^[49]
		Xie et al. (2015) ^[50]
		Li et al. (2018) ^[51]
		Jiang et al. (2020) ^[52]
		Pasin et al. (2024) ^[53]
问答	Li et al. (2019) ^[54]	
	Mansky et al. (2023) ^[55]	
	Duneau et al. (2024) ^[56]	
情感分析	Zhang et al. (2021) ^[57]	
	Ruskanda et al. (2023) ^[58]	
量子计算方法 (量子计算机硬件)	多模态信息分类	Qu et al. (2023) ^[59]
	问答	Meichanetzidis et al. (2020) ^[60]
	文本分类	Meichanetzidis et al. (2023) ^[61]
	机器翻译	Vicente(2021) ^[62]

将目前文献中的 QNLP 路线按照“理论-经典模拟-量子实现”3 类属性加以区分,并结合其对应的任务类型进行举例说明。特别地,那些已在真实量子硬件上完成实验的研究虽然数量有限,但可视为对理论模型的应用性验证,展示了量子

语言模型在特定 NLP 任务上可能具备的结构优势与可靠性。由于当前量子硬件仍受规模与噪声的限制,这些实验通常基于小或中规模的数据集进行,但已为后续量子语义模型的可扩展性研究奠定了基础。

基于范畴论的组合语义模型在量子自然语言处理中的实践性研究大致经历了两个阶段:2021—2024 年间,研究重心从单一的 Lambeq 实现逐渐转向与机器学习相结合的混合模型,即在 DisCoCat 模型中引入深度学习机制,以提升模型性能与适应性;2024 年之后,基于量子神经网络(Quantum Neural Network, QNN)的自然语言处理模型数量已为 Lambeq 的两倍多,显示出深度学习驱动的量子语言建模方法正逐渐成为主流。

结束语 本研究表明,基于范畴论的自然语言语义表示范式的核心是通过函子映射在语法范畴与语义向量空间之间建立严格、可计算的对应关系。这一范式提供了将组合性、分布性与可计算性统一起来的数学基础。首先,利用 DisCoCat 提供的图解理论,为 QNLP 提供数学基础。然后,使用该图形系统,自然语言可以被解释为量子过程。此外,字符串图可以翻译为量子线路。这一步对于该理论在量子硬件上的真正实现至关重要。最后,DisCoCat 框架中的图形语言为比较不同语言句子的语法结构提供了统一的形式化工具。

当前基于范畴论的语义表示的范畴研究在理论层面为统一句法结构与分布语义提供了系统而严格的数学基础,但在走向实际应用与理论深化时,仍面临以下核心挑战和研究突破。

1)模型表达能力方面,现有模型(如基于 Pregroup/CCG 的 DisCoCat)对句内组合语义的刻画已较为成熟,但对语篇级语义、长距离依赖、嵌套结构及隐喻等复杂现象的处理能力有限。未来需引入更丰富的范畴论构造(如高阶函子)来增强模型的表达力与推理能力。例如,文献[63]探索了在态射上附加概率信息以构建更精细的语义模型。

2)语法-语义接口方面,当前模型严重依赖预设的语法范畴(如 CCG, Pregroup),其覆盖范围与复杂语言现象的适配性不足。关键挑战在于设计可学习的函子映射,使模型能从数据中自动归纳语法类型与语义空间的对应关系,实现神经符号融合。例如,文献[64]探索了从 CCG 推导语篇线路的路径,但如何自动化并推广此过程仍是具有挑战的难题。

3)在计算实现层面,范畴模型(尤其是涉及高阶张量运算或量子线路模拟的模型)仍面临显著的可扩展性挑战。一方面,现有方法在参数规模、计算复杂度与训练效率方面难以适配大规模自然语言语料,亟需发展更高效的参数化策略、近似计算方法与轻量化组件设计,以提升模型在资源受限环境中的可部署性。另一方面,在量子计算语境下,范畴组合语义模型的实际落地仍受到近中期量子设备(NISQ)的严格限制。文献[65]表明,在将 DisCoCat 等模型映射为量子线路时,受限于量子比特数量与噪声水平,当前实验仅能处理极为简化的语法结构与有限的词汇规模,尚难扩展至真实自然语言场景。这一现状表明,范畴组合语义在量子计算中的潜在优势,仍有赖于量子硬件能力、算法设计与模型结构三方面的协同推进。

参考文献

- [1] MIKOLOV T, CHEN K, CORRADO G, et al. Efficient estimation of word representations in vector space [C] // Proceedings of the 1st International Conference on Learning Representations. 2013.
- [2] PENNINGTON J, SOCHER R, MANNING C D. GloVe: Global vectors for word representation [C] // Proceedings of the Conference on Empirical Methods in Natural Language Processing. 2014:1532-1543.
- [3] PETERS M E, NEUMANN M, IYYER M, et al. Deep contextualized word representations [C] // Proceedings of the Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies. 2018:2227-2237.
- [4] DEVLIN J, CHANG M W, LEE K, et al. BERT: Pre-training of deep bidirectional transformers for language understanding [C] // Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies. 2019:4171-4186.
- [5] MITCHELL J, LAPATA M. Vector-based models of semantic composition [C] // Proceedings of ACL-08: Human Language Technologies. 2008:236-244.
- [6] BLACOE W, LAPATA M. A comparison of vector-based representations for semantic composition [C] // Proceedings of the 2012 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning. 2012:546-556.
- [7] DE FELICE G. Categorical tools for natural language processing [J]. arXiv:2212.06636, 2022.
- [8] TULL S, LORENZ R, CLARK S, et al. Towards compositional interpretability for XAI [J]. arXiv:2406.17583, 2024.
- [9] LI R, ZHAO X, MO M F. A brief overview of universal sentence representation methods: A linguistic view [J]. ACM Computing Surveys, 2022, 55(3):1-42.
- [10] ARORA S, LIANG Y, MA T. A simple but tough-to-beat baseline for sentence embeddings [C] // Proceedings of the 5th International Conference on Learning Representations. 2017.
- [11] RÜCKLÉ A, EGER S, PEYRARD M, et al. Concatenated p-mean word embeddings as universal cross-lingual sentence representations [C] // Proceedings of the 2018 Conference. 2018.
- [12] LE Q V, MIKOLOV T. Distributed representations of sentences and documents [C] // Proceedings of the 31st International Conference on Machine Learning. 2014:1188-1196.
- [13] HILL F, CHO K, KORHONEN A. Learning distributed representations of sentences from unlabelled data [C] // Proceedings of the Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies. 2016:1367-1377.
- [14] ZHANG M, WU Y, LI W K, et al. Learning universal sentence representations with mean-max attention autoencoder [C] // Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing. 2018.
- [15] CHEN Q, WANG W, ZHANG Q L, et al. Ditto: A simple and

- efficient approach to improve sentence embeddings [C]// Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing. 2023.
- [16] KIROS R, ZHU Y, SALAKHUTDINOV R, et al. Skip-thought vectors [C]// Proceedings of the Advances in Neural Information Processing Systems. 2015.
- [17] LOGESWARAN L, LEE H. An efficient framework for learning sentence representations [C]// Proceedings of the 6th International Conference on Learning Representations. 2018.
- [18] NIE A, BENNETT E D, GOODMAN N D. DisSent: Sentence representation learning from explicit discourse relations [C]// Proceedings of the 2017 Conference. 2017.
- [19] SILEO D, VAN DE CRUYS T, PRADEL C, et al. Mining discourse markers for unsupervised sentence representation learning [C]// Proceedings of the Conference of the North American Chapter of the Association for Computational Linguistics; Human Language Technologies. 2019; 3477-3486.
- [20] KIROS J, CHAN W. InferLite: Simple universal sentence representations from natural language inference data [C]// Proceedings of the Conference on Empirical Methods in Natural Language Processing. 2018.
- [21] REIMERS N, GUREVYCH I. Sentence-BERT: Sentence embeddings using siamese BERT-networks [C]// Proceedings of the Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP). 2019.
- [22] JERNITE Y, BOWMAN S R, SONTAG D. Discourse-based objectives for fast unsupervised sentence representation learning [J]. arXiv:1705.00557, 2017.
- [23] SUBRAMANIAN S, TRISCHLER A, BENGIO Y, et al. Learning general purpose distributed sentence representations via large scale multi-task learning [C]// Proceedings of the 6th International Conference on Learning Representations. 2018.
- [24] CER D, YANG Y, KONG S Y, et al. Universal sentence encoder for English [C]// Proceedings of the Conference on Empirical Methods in Natural Language Processing. 2018; 169-174.
- [25] SCHOPF T, SCHNEIDER D, MATTHES F. Efficient domain adaptation of sentence embeddings using adapters [J]. arXiv: 2307.03104, 2023.
- [26] SOCHER R, PERELYGIN A, WU J, et al. Recursive deep models for semantic compositionality over a sentiment treebank [C]// Proceedings of the Conference on Empirical Methods in Natural Language Processing. 2013; 1631-1642.
- [27] TAI K S, SOCHER R, MANNING C D. Improved semantic representations from tree-structured long short-term memory networks [C]// Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing. 2015; 1556-1566.
- [28] GEHRING J, AULI M, GRANGIER D, et al. Convolutional sequence to sequence learning [C]// Proceedings of the 34th International Conference on Machine Learning. 2017; 1243-1252.
- [29] COECKE B, SADRZADEH M, CLARK S. Mathematical foundations for a compositional distributional model of meaning [J]. arXiv:1003.4394, 2010.
- [30] LAN Z, CHEN M, GOODMAN S, et al. ALBERT: A lite BERT for self-supervised learning of language representations [J]. arXiv:1909.11942, 2019.
- [31] OPENAI. GPT-5 system card [EB/OL]. <https://openai.com/index/gpt-5-system-card/>.
- [32] RUSSELL B. The principles of mathematics [M]. London: Routledge, 1903.
- [33] AJDUKIEWICZ K. Die syntaktische Konnexität [J]. *Studia Philosophica*, 1935, 1: 1-27.
- [34] BAR-HILLEL Y. Logical syntax and semantics [J]. *Language*, 1954, 30(2): 230.
- [35] EILENBERG S, MAC LANE S. General theory of natural equivalences [J]. *Transactions of the American Mathematical Society*, 1945, 58: 231-294.
- [36] LAWVERE F W. Functorial semantics of algebraic theories [J]. *Proceedings of the National Academy of Sciences of the United States of America*, 1963, 50(5): 869-872.
- [37] LAMBEK J. On the calculus of syntactic types [C]// *Structure of Language and Its Mathematical Aspects*. Providence: American Mathematical Society, 1961; 166-178.
- [38] COECKE B, GREFFENSTETTE E, SADRZADEH M. Lambek vs. Lambek: Functorial vector space semantics and string diagrams for Lambek calculus [J]. *Annals of Pure and Applied Logic*, 2013, 164(11): 1079-1100.
- [39] JOYAL A, STREET R. The geometry of tensor calculus II [R]. Unpublished manuscript, 1995.
- [40] KARTSAKLIS D, FAN I, YEUNG R, et al. Lambeq: An efficient high-level Python library for quantum NLP [J]. arXiv: 2110.04236, 2021.
- [41] DE FELICE G, DI LAVORE E, ROMÁN M, et al. Functorial language games for question answering [C]// Proceedings of the 3rd Annual International Applied Category Theory Conference (ACT 2020). 2020; 311-321.
- [42] TOUMI A, DE FELICE G. Higher-order DisCoCat (Peirce-Lambek-Montague semantics) [J]. arXiv:2311.17813, 2023.
- [43] ZENG W, COECKE B. Quantum algorithms for compositional natural language processing [J]. *Electronic Proceedings in Theoretical Computer Science*, 2016, 221: 67-75.
- [44] YEUNG R, KARTSAKLIS D. A CCG-based version of the DisCoCat framework [J]. arXiv:2105.07720, 2021.
- [45] COECKE B. The mathematics of text structure [C]// *The Interplay of Mathematics, Logic, and Linguistics*. 2021; 181-217.
- [46] WANG-MASCIANICA V, LIU J, COECKE B. Distilling text into circuits [J]. arXiv:2301.10595, 2023.
- [47] LAAKKONEN T, MEICHANETZIDIS K, COECKE B. Quantum algorithms for compositional text processing [J]. *Electronic Proceedings in Theoretical Computer Science*, 2024, 406: 162-196.
- [48] LE DU S, HERNÁNDEZ SANTANA S, SCARPA G. A gentle introduction to quantum natural language processing [J]. arXiv: 2202.11766, 2022.
- [49] SORDONI A, NIE J Y, BENGIO Y. Modeling term dependencies with quantum language models for IR [C]// Proceedings of

- the 36th International ACM SIGIR Conference on Research and Development in Information Retrieval. 2013:653-662.
- [50] XIE M, HOU Y, ZHANG P, et al. Modeling quantum entanglements in quantum language models [C] // Proceedings of the Twenty-Fourth International Joint Conference on Artificial Intelligence. 2015.
- [51] LI Q, MELUCCI M, TIWARI P. Quantum language model-based query expansion [C] // Proceedings of the 2018 ACM SIGIR International Conference on Theory of Information Retrieval. 2018:183-186.
- [52] JIANG Y, ZHANG P, GAO H, et al. A quantum interference inspired neural matching model for ad-hoc retrieval [C] // Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval. 2020:19-28.
- [53] PASIN A, CUNHA W, GONÇALVES M A, et al. A quantum annealing instance selection approach for efficient and effective transformer fine-tuning [C] // Proceedings of the 2024 ACM SIGIR International Conference on Theory of Information Retrieval. 2024:205-214.
- [54] LI Q, WANG B, MELUCCI M. CNM: An interpretable complex-valued network for matching [C] // Proceedings of the Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies. 2019:4139-4148.
- [55] MANSKY B, WÖRLE F, STEIN J K, et al. Adapting the DisCoCat-model for question answering in the Chinese language [C] // Proceedings of the IEEE International Conference on Quantum Computing and Engineering. 2023:591-600.
- [56] DUNEAU T, BRUHN S, MATOS G, et al. Scalable and interpretable quantum natural language processing: An implementation on trapped ions [J]. arXiv:2409.08777, 2024.
- [57] ZHANG P, ZHANG J, MA X, et al. TextTN: Probabilistic encoding of language on tensor network [C] // Proceedings of the International Conference on Learning Representations. 2021.
- [58] RUSKANDA F Z, ABIWARDANI M R, AL BARI M A, et al. Quantum representation for sentiment classification [C] // Proceedings of the IEEE International Conference on Quantum Computing and Engineering. 2022:67-78.
- [59] QU Z, MENG Y, MUHAMMAD G, et al. QMFND: A quantum multimodal fusion-based fake news detection model for social media [J]. Information Fusion, 2024, 104:102172.
- [60] MEICHANETZIDIS K, TOUMI A, DE FELICE G, et al. Grammar-aware question-answering on quantum computers [J]. arXiv:2012.03756, 2020.
- [61] MEICHANETZIDIS K, TOUMI A, DE FELICE G, et al. Grammar-aware sentence classification on quantum computers [J]. Quantum Machine Intelligence, 2023, 5(1):10.
- [62] NIETO V. Towards machine translation with quantum computers [D]. Stockholm: University of Stockholm, 2021.
- [63] BRADLEY T, TERILLA J, VLASSOPOULOS Y. An enriched category theory of language: from syntax to semantics [J]. La Matematica, 2022, 1(2):551-580.
- [64] LIU J, SHAIKH R A, RODATZ B, et al. A pipeline for discourse circuits from CCG [J]. arXiv:2311.17892, 2023.
- [65] LIU T, WEI Y, WANG J. Research on distributional compositional categorical model in both classical and quantum natural language processing [C] // Proceedings of the SNPD 2024. 2024:1-6.



LI Yidan, born in 1994, Ph.D candidate. Her main research interests include formal semantics, category theory, artificial intelligence logic and natural language processing.



CUI Jianying, born in 1975, Ph.D, associate professor, Ph.D supervisor. Her main research interests include formal argumentation theory, artificial intelligence logic and natural language processing.

(责任编辑:何杨)