

# 基于分子图相似度的医药文献推荐方法

冯贤兵, 陶 涛, 吕肖庆

北京大学王选计算机研究所, 北京

收稿日期: 2022年11月20日; 录用日期: 2022年12月19日; 发布日期: 2022年12月26日

## 摘 要

当今生物医学等领域的文献快速增长, 一方面促进了科研交流, 但同时也为研究人员带来了巨大的阅读压力, 尽管业界已出现了一些论文搜索和推荐的方法, 但其大都只依据论文的元信息和文本信息, 而对文章内容, 尤其是插图等非文字对象尚未充分挖掘并利用, 因此现有系统在给读者的推荐结果中, 还存在着大量重复、泛化等低效情况。为此, 我们探索并建立了一个基于论文内容的文档级推荐系统, 具体包括: 文档解析、文本对象理解、内容相似性度量、多级索引机制、以及优化推荐结果等主要环节。其中, 针对生物医学类科技文献中特有的分子式图片, 我们提出了一种图相似度的度量方法, 即半分支编辑距离(Half-branch GED, 简称HB-GED)算法, 同时针对分子图表示和文档之间关系表示也提出了图卷积模型。在真实数据集上的实验结果表明, 本文提出的论文推荐方法, 可有效筛选出更符合查询者意图的候选论文。

## 关键词

图相似度, 分子图, 论文推荐, 图编辑距离, 二部图

# A Recommender via Similarity of Molecule Graphs for Medical Literature

Xianbing Feng, Tao Tao, Xiaoqing Lyu

Wangxuan Institute of Computer Technology, Peking University, Beijing

Received: Nov. 20<sup>th</sup>, 2022; accepted: Dec. 19<sup>th</sup>, 2022; published: Dec. 26<sup>th</sup>, 2022

## Abstract

Nowadays, the consistent growth of scientific and technical literature leads to formidable pressure on medical researchers. Researchers turn to the search engine and paper recommender systems and still have to spend more time keeping up with the trends and directions in their field. However, most existing recommender approaches mainly depend on text-based information and ignore

non-text objects, such as informative figures. To this end, we establish a document-to-document recommender system for medical literature. Specifically, we proposed a deep-learning-based segmentation method for extracting molecular graphs, a Half-branch GED algorithm for evaluating the similarity of molecules, and a bipartite-graph-based algorithm for paper similarity, respectively. Experimental results on real-world datasets demonstrate the effectiveness of the proposed recommender system.

## Keywords

Graph Similarity, Molecule Graph, Paper Recommender, Graphic Edit Distance, Bipartite Graph

Copyright © 2022 by author(s) and Hans Publishers Inc.

This work is licensed under the Creative Commons Attribution International License (CC BY 4.0).

<http://creativecommons.org/licenses/by/4.0/>



Open Access

## 1. 引言

当今生物学等领域的文献快速增长,这一方面积极地促进了各领域的科研交流,但同时,大量激增的文献阅读量也为研究人员带来了巨大的查阅负担,尤其是如何获得真正有帮助的最新文献,正在成为很多一线科研人员的阻碍。尽管业界已出现了一些论文搜索和推荐的方法,即现有研究中不乏可用的推荐方法,但这些方法大都只依据论文的元信息和文本信息,例如搜索相似的作者、引用关系,或者通过提取文章关键词或摘录文章内容来度量文章的相似性,进而推测文章间的关系,它们很少利用文章的非文字内容,尤其是插图等复杂对象,尚未充分挖掘并利用,因此,现有系统在给读者的推荐结果中,还存在着大量重复、泛化等低效情况,以基于关键词的方法为例,一个或多个关键词往往不能涵盖一篇文章的完整意义,这样就会导致推荐系统输出很多无关的论文。

为此,有必要深入挖掘并全面利用文献中的非文本信息,将更全面的整篇文章信息用于比较,以此可给推荐系统带来更多的准确推荐结果。在生物学领域,有一类重要的非文本的内容,是该类文章中的分子式,它们包含着非常重要的信息,值得进一步的提取和利用。但此类探索面临着诸多技术挑战,例如:1) 如何从文章中提取各式各样的图表信息,尤其是精细的分子图信息;2) 如何有效地计算分子图之间的相似度;3) 如何根据内容相关性建立查询文章和备选文章的关系网络,并实现有效的连接预测等。

对此,本项研究基于分子式检索技术的积累,通过检测包含有相同或相似的分子式的文献,实现更为有效的医药文献推荐方法,主要创新体现在:1) 提出了一套从文档到文档的推荐系统,简化了用户交互操作,但增加了输入信息量;2) 根据分子式的特点,对图表示和相似度计算改进了现有的图神经网络模型和图编辑距离算法,以折叠策略更有效地聚合了分子的图形特征,以半分支结构提高了分子相似度计算的准确度和效率;3) 对于文档与分子式之间存在的多对多的包含关系,建立了更全面的关系描述模型,可以更加准确地推测论文之间的相关程度。

## 2. 相关工作

过去数年中我们可以看到推荐系统的广泛运用,尤其是许多领域中(社交网络,电子商务等)与图神经网络的(GNN)的结合[1]。然而文献查询相比于商品查询目的有所不同。比如,一个买家在买了一样商品后,他/她可能不会再期待推荐更多同一类型的商品。而科研人员阅读一篇文章后,他/她可能会更加关注与这篇文献相关的,即,研究相同问题的,同类型文章。因此,有效的推荐系统更注重对文件内容相关

性的深入挖掘和利用。

文本检索的文章主要基于 TF-IDF 技术来排序。之后的提升也是通过增加交互信息来改进[2]。近年来基于机器学习的方法主要关注两个方面: 1) 利用神经网络来计分或重新排序, 如 Co-PACRR [3], KNRM [4]; 2) 利用特征学习在向量空间来实现准确有效的查询[5] [6] [7] [8]。为了给图检索构建一个索引, Qian 等[9]利用文本单元作为节点, 关系作为边来构造图。现在效果好的检索系统主要关注如何利用机器学习方法和 NLP 的技术在索引、检索、排序等方面进行增强[10]。

和图相似度相关的重要工作有: Zheng 等[11]通过对边的池化(Hyperedge Pooling)和子图匹配来将正常的图转化成超图来评估相似度; Coupette 等[12]定义了一个标准的、极小的描述距离来衡量图之间的相似度; Raveaux 等[13]证明了在特定情况下图匹配工作可以用一个基于重配错误(Reformulated Error-based)模型来等价转化成一个 GED 问题; Riba 等[14]通过一个消息传递网络来构建图距离; Ling 等[15]通过学习图的特征与度量函数共同计算图相似度。

### 3. 文档推荐模型

要实现文档级的推荐, 需要建立一个基于论文内容理解的推荐系统, 至少包括: 文档解析、文本对象理解、内容相似性度量、索引机制、以及优化推荐结果等重要环节。本文针对生物医学类科技文献中特有的分子式图片提出了一个基于分子式相似度的论文推荐模型, 其流程及主要环节如图 1 所示。

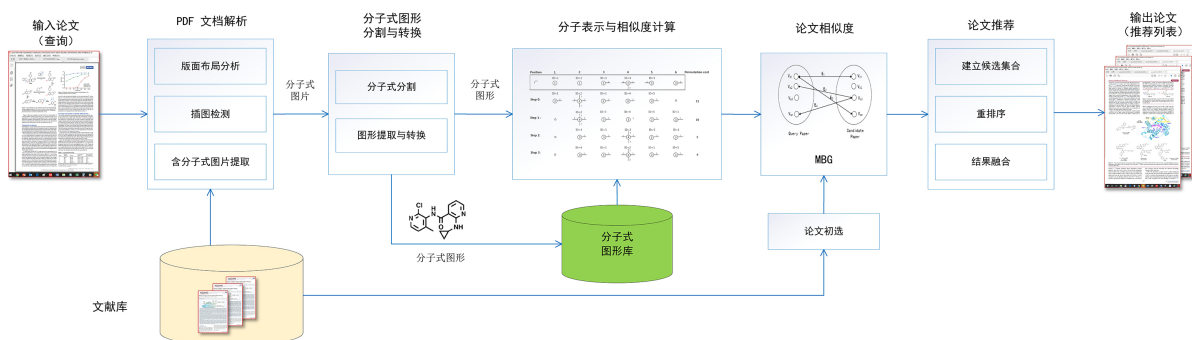


Figure 1. Medical literature recommender model based on molecule graphs similarity

图 1. 基于分子式相似度的医药论文推荐模型

该模型旨在通过分子式的图相似度来提升论文推荐系统的准确性和效率, 其工作流程可大致分为以下五个主要的步骤:

- 文档解析(PDF Analysis), 首先采用 PDF 分析器模块提取可用于推荐的内容信息, 其中包括, 通过深度学习算法定位、检测并提取其中的分子式的原始图片;
- 分子式分割与转换(Molecule Graph Extraction), 从分子式原始图片(大多为更复杂的化学反应式)中分割出单个分子式, 并将其转换为图的表示形式;
- 图的表征(Molecule Graph Embedding), 亦称为图表示, 利用新兴的深度图卷积神经网络方法计算图的深层结构信息;
- 分子图相似度度量(Molecule Similarity), 本文采用所提出的 HB-GED 算法来度量单个分子式之间的相似度;
- 文档相似度度量(Paper Similarity), 一篇文章往往含有多个分子图, 为此本文采用了分子二分图(MBG)来描述两篇文章中分子之间的相似关系, 进而计算文章之间的相似度。

通过以上关键步骤, 最后的推荐模块就可以根据文档相似度来筛选论文, 结合现有推荐方法就可以

提供更为精准的论文推荐列表。以下各小节主要针对模型中涉及的分子式分割、分子图表示、分子相似度量度和文章相似度度量几项关键技术进行介绍。

### 3.1. 分子式分割

由于分子式并不是科技文档中唯一的复杂对象，因此分子式的检测还要针对其与数学公式、插图和表格等做进一步的区分。同时，分子式很少单独出现，而是在化学反应式、化合物转换图、以及合成原理图等场景中出现，这无疑进一步加大了单个分子分割和提取的难度。现有的自然场景对象检测与分割方法可提供很好的借鉴，但文档页面图像不同于自然图像，它们没有丰富的纹理、颜色、频率等特征，因此常规的图像分割算法用于分子式分割后效果并不理想。

针对分子图所占版面区域的不规则性，需要更加细致地研究分子式区域的版面特征，有效挖掘并综合利用图形的上下文信息。据此，本文提出了一种分子式图片分割模型，其整体框架如图 2 所示，主要分为三个模块：

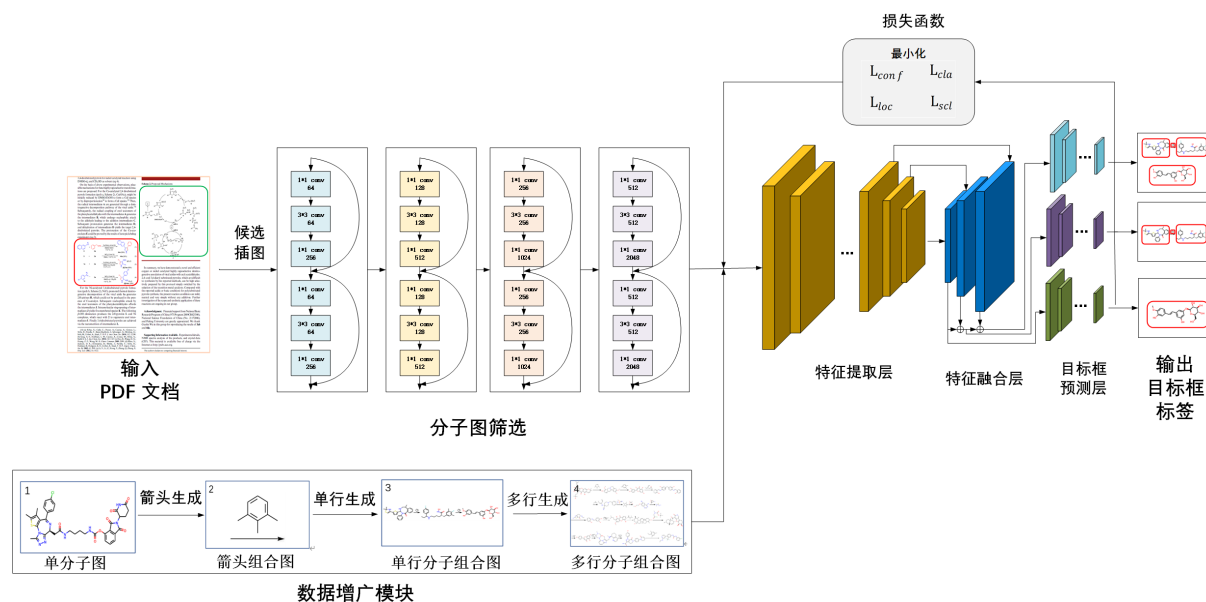


Figure 2. The model of molecule segmentation

图 2. 分子式分割模型

- 分子式插图筛选：从含有公式、表格、统计图表等其他类型的候选插图集合中，采用一个粗分类的网络，如 ResNet-50 模型，将含有分子式或化学反应方程式的图片筛选出来。
- 数据扩增：为了获得足够多分子式标注插图用于分割模型的学习，本文模拟分子方程式的生成方法，分三步从单个分子式的图片生成分子方程式的图片，第一步是生成多样的反应式中常见的箭头形式，第二步是将箭头和单个分子结合生成单行方程式，第三步是排列多个生成的单行方程式，构成复杂的多行方程式情况，在后两步中，为了进一步模拟真实插图，还需加入分子名称等文字框。
- 分割预测：基于 Yolo-v5 和 RetinaNet 等模型，设计了类似的预测架构，包括针对分子式的特征提取网络、特征融合网络、分类网络和位置回归网络，其中的特征提取使用 Darknet-53 网络，在特征融合使用 New CSP-PAN 网络，而其位置回归网络和分类网络则对应着 Faster-RCNN 模型中的全连接层。

在本模型中，我们的损失函数设计除了沿用 Yolo-v5 等模型中原有的置信度损失、分类损失、定位损失之外，针对现有损失仍不足以描述预测边框与实际边框的覆盖精细度的问题，新增设了外形约束损

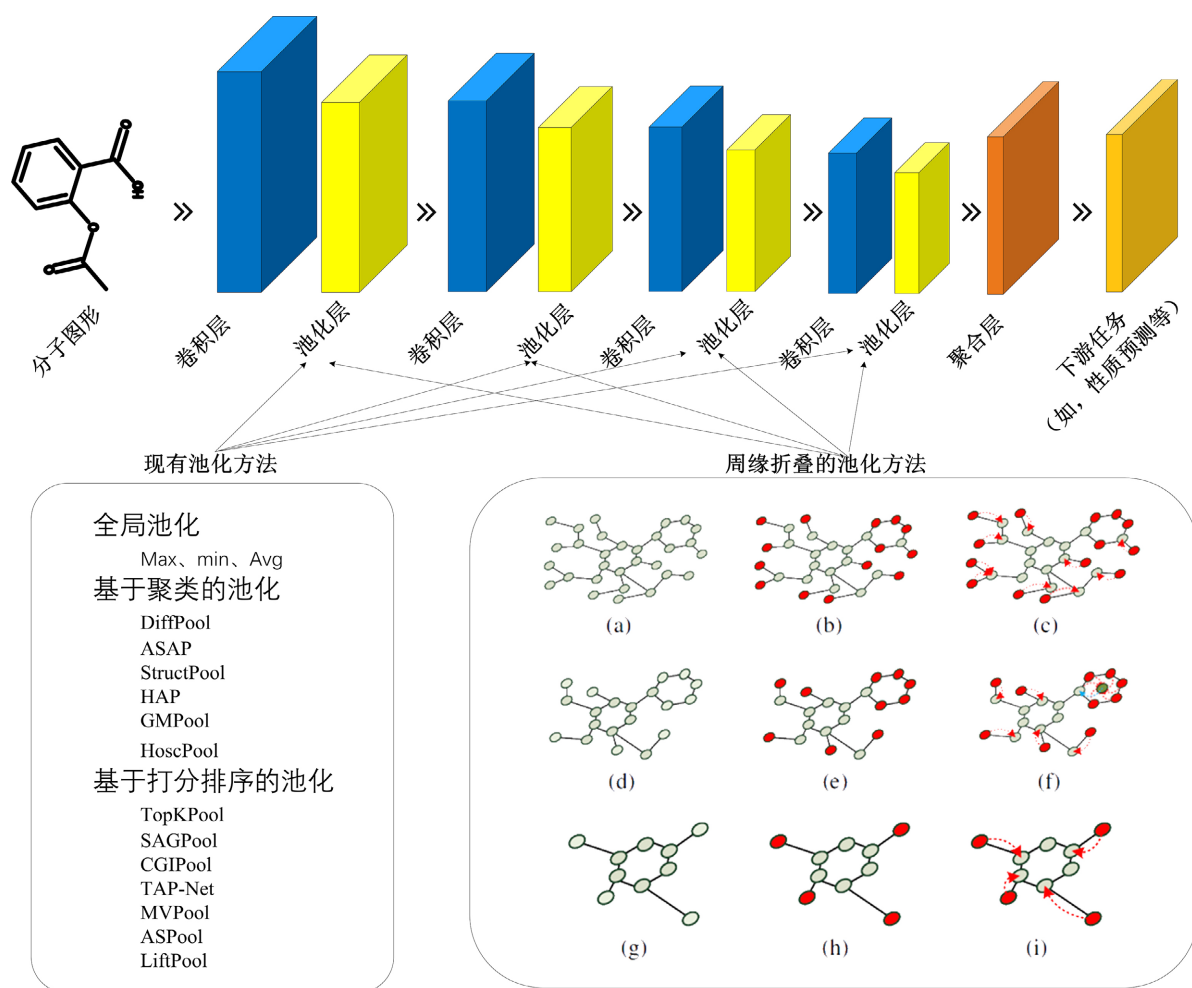
失函数  $L_{SCL}$ ，公式如下：

$$L_{SCL} = 1 - \frac{Inter(p, g)}{w(g) * h(g)} \quad (1)$$

其中， $Inter(p, g)$ 表示预测目标框和实际框的交集， $w$ 和 $h$ 表示实际目标框的宽和高。

### 3.2. 分子式图表示的深度图神经网络模型

以图形方式表示分子可以更深入地探索分子式的结构特征和拓扑特征，本文采用图神经网络(Graph Neural Network, GNN)作为基本的图表示方法。但此类方法中的池化(Pooling)操作对于卷积层之间信息聚合，不能有效地适应分子图的结构特点，因为灵活度很高的图结构在不同层之间信息继承机制上与较为规则的图像有本质不同。为此，本文提出了一种周缘折叠的图归约方法来替代现有的层间信息聚合方法，如图3所示，其出发点是：为真正发挥不同卷积层对图的抽象能力，引入了沿大图周边逐级向内进行动态折叠的图归约算法，逐级提高各个卷积层对全图的概括能力。



**Figure 3.** Margin-fold-based pooling is different from the existing pooling operations in graph neural networks for molecule expression

**图 3.** 在用于分子表示的图神经网络模型中，不同于现有池化方法的基于周缘折叠池化方法

具体地，首先将原始图作为最底层，检测其首批外边缘元素(包括节点、边和子结构)，作为首次折叠



的操作对象, 并找到其对应的内向邻居; 折叠过程亦可视为低层到高层的抽象过程, 即将边缘元素及属性信息融合到其内向邻居中之后, 丢弃边缘元素, 并形成新的卷积层; 重复以上步骤, 即重新检测新的边缘元素, 完成下一次折叠式的迭代。每次折叠可将图大幅度缩小, 形成抽象度更高的卷积层, 同时又通过信息融合有效地传递并保留了周边节点的局部信息。折叠的作用等同于传统卷积神经网络中池化作用, 同时又更加贴近分子图的拓扑结构, 可以自适应地提取和融合图的信息, 能更加准确地汇集不同图层之间的信息。此外, 在折叠过程中, 算法自动检测并识别分子中的具有化学语义的子结构, 如苯环等, 并将其直接塌缩为超级节点, 从而进一步提高了分子图的池化效率, 而且可获得更为简洁, 但又不破坏、不丢失化学语义信息的分子表示。

### 3.3. 分子相似度度量

在图相似性度量方面, 图形编辑距离(GED)是目前普遍采用的算法, 但在计算分子图相似度时, 笔者感受到其有两方面不足: 1) 由于大多数分子所含的原子数量较多, 导致了图形规模增大, 常规的 GED 算法处理能力非常有限; 2) 分子图中存在着大量特有的固定子结构, 而且频繁出现, 如苯环等, 而常规 GED 算法始终以最基本的节点和边作为计算对象, 不能高效地处理或利用这些子结构。

针对上述问题, 我们提出了一种半分支编辑距离算法(Half-branch GED, 简称 HB-GED)用于分子图的相似度度量。具体地, HB-GED 将编辑距离的计算对象扩大到“半分支”, 即将当前节点和它所连带的边看作一个基本计算单元, 因此放大了比较计算的对象, 相应地就减少了全图范围内编辑距离的计算操作。同时, HB-GED 将计算 GED 的问题转化为找到半分支图项(Graph Item)最佳排列的问题, 以求解最优排列的策略替代传统的遍历式求解策略, 即, 我们设计了一种在置换空间中的迭代局部搜索算法, 在搜索过程中, 我们只考虑在一个局部的半分支图项并导出一个次优距离, 可以用更少的迭代步骤接近最优或局部最优结果。此外, HB-GED 还采用了剪枝策略, 基于下限过滤掉不太可能的半分支图项排列, 进一步缩小搜索范围, 提高了算法的效率。

### 3.4. 基于二分图的文章相似度

为了评估论文的相似度, 我们设计了一种基于分子二分图(MBG)的文章相似度度量算法。在 MBG 中, 如图 4 所示, 左侧节点集  $V_q$  表示查询论文中的分子, 右侧节点集  $V_c$  表示一张候选论文中的分子。边集合  $E$  表示两个分子之间可能存在的相似性。

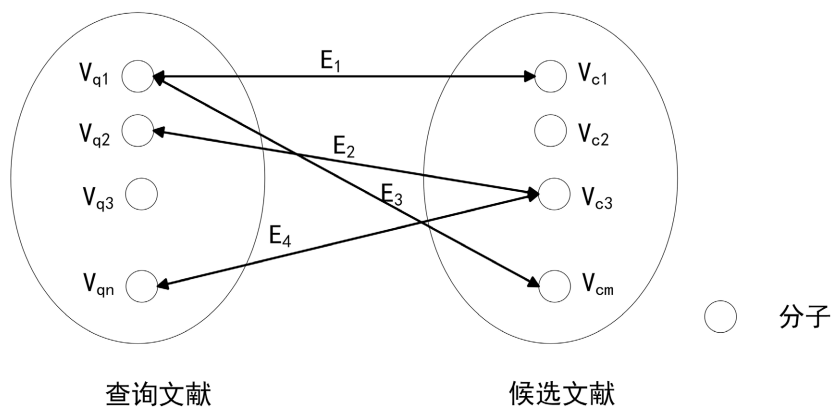


Figure 4. MBG for calculating the paper similarity  
图 4. 用于计算文章相似度的分子二部图(MBG)

本文根据三个因素(平均权重、全局相似度和自适应权重)计算文章相似度, 步骤如下:

- 平均权重，利用边数量将分子间的相似值标准化得到  $Mean(E)$ ；
  - 全局相似度，作为文章相似度的最重要部分，通过二分图之间的边数  $|E|$  和最多可能边数  $|V_q| \times |V_r|$  的比例来计算；
  - 自适应权重，为了强调查询文章的影响而引入，其值为与其它文章有相似度边的分子式数量  $N_{Degree>0}(V_q)$  和文章总分子式数量  $|V_q|$  的比例。
- 最终，可用公式

$$S = Mean(E) \times \frac{|E|}{|V_q| \times |V_r|} \times \frac{N_{Degree>0}(V_q)}{|V_q|} \quad (2)$$

计算两篇含有一个或多个相似分子式的文章的相似度。

## 4. 实验

为评估本文提出技术的有效性，我们对其分别进行了对比实验，特别是量化评估了分割算法、分子相似度算法与现有算法相比的改进效果。

### 4.1. 分子式分割对比试验

在数据集建设方面，本文首先从万方数据的论文中直接提取的 561 张分子图片，并据此采用数据增广与合成技术，扩充至 2500 张图片。为评价分子分割模型的性能，我们采用了平均准确率(mAP)和召回率两个指标，该方法与常见图像分割模型的实验对比情况如表 1 所示。

**Table 1.** Comparison of detection results of segmentation model on molecular picture dataset

**表 1.** 分割模型在分子式图片数据集上的检测结果对比

模型	mAP (0.5~0.95)	mAP@0.5	mAP@0.75	平均召回率
Faster-RCNN [16]	0.446	0.684	0.500	0.545
Mask-RCNN [17]	0.471	0.713	0.525	0.614
RetinaNet [18]	0.499	0.776	0.563	0.630
FCOS [19]	0.431	0.724	0.457	0.499
本文方法	0.533	0.784	0.581	0.695

对比实验的结果显示，本文方法较其他方法的准确率有所提升。同时，图 5 显示了本文所提出的分割模型在真实 PDF 医药文献中检测并分割分子式图片的效果，即，在较为复杂的版面中，该模型仍可获得较为理想的分子式分割边界。

### 4.2. 分子相似度实验

在分子式检索的实验中，我们从万方数据官方语料库中，选取了 187 篇医学论文建立了一个资源文献数据集。借助于分子检测和图形提取模块，我们得到 2829 张分子图和 1325 张分子图，在真实数据的环境中验证了推荐系统的有效性。为了量化检验 HB-GED 的有效性，我们与现有的相似度算法进行了比较，在评价标准方面，在 AIDS 数据集和 LINUX 数据集上[15]采用了平方误差(MSE)。

在 WIKI 数据集上的检索实验中，我们采用了公制平均精度(mAP)。

表 2 和表 3 的实验结果表明，HB-GED 在准确性和有效性方面优于现有方法。

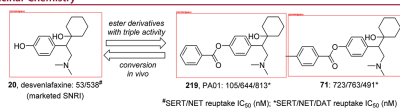


Figure 39. Conversion of the SNRI 20 into the corresponding benzoate ester derivatives to incorporate DAT affinity (Luye Pharma).

concentration (5-HT, 180%; NE, 400%; DA, 270%) in mPFC (microdialysis).

**5.9. Luye Series.** The marketed SNRI desvenlafaxine (20) is known for relatively low permeability, presumably due to the presence of an ionizable phenol group. Ester derivatives with TRI activity were sought assuming that they would enhance efficacy by adding a DA component and by improving the CNS exposure through enhanced permeability. These compounds can also behave as prodrugs in vivo. Hou et al. have disclosed PA01 (219), the benzoate ester of the SNRI 20 (Figure 39).<sup>137</sup> Compound 219 showed a dose-dependent increase of mobility time in both the rat FST and mouse TST models with higher efficacy than 20 while showing no stimulatory effect on spontaneous locomotor activity. The anti-immobility effect of 219 was significantly reversed by pretreatment of the mice with 4-chlorophenylalanine, an inhibitor of 5-HT synthesis, the D<sub>1</sub> antagonist SCH23390 (220), and the D<sub>2</sub> antagonist sulpiride (221) (Figure 40). Compound 71, which is being developed

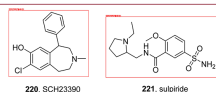


Figure 40. Structures of DA receptor antagonists.

for MDD, is a 4-methylbenzoate derivative of 20. Compound 71 was found to partially undergo cleavage in vivo to release 20 in rats.<sup>138</sup> Compared to 20, 71 elevated 5-HT, NE, and DA to higher levels in rat striatum following acute or chronic administration as an oral suspension in a microdialysis

experiment and exhibited increased reduction of immobility time in the rat FST model (po). Unlike 20, 71 did not induce a reduction in the levels of extracellular 5-HT in the early phase following iv administration of a solution of the compound. In a phase I clinical trial, 71 was found to show dose-proportional plasma exposure with a good safety profile and an absence of a food effect on the bioavailability.<sup>139</sup> A phase II clinical trial of this compound was initiated in China in September 2015.<sup>78</sup>

**5.10. D Series.** Dutta et al. reported the design of TRI pyran analogues based on their previous work of structurally constrained *cis*-3,6-disubstituted piperidine derivatives as selective DAT inhibitors (e.g., 222) that were explored for cocaine antagonism (Figure 41).<sup>139,140</sup> While the benchmark pyran compound 223 showed potential as a TRI, the introduction of a H-bond acceptor/donor substituent on the phenyl ring of the benzyl group, explored in the context of 224–228 (D-161, D-391, D-141, D-185, and D-411) which incorporate OH, OMe, or NH<sub>2</sub> resulted in an improvement in the affinity for SERT and NET, suggesting the importance of H-bonding interactions. Relocation of the 4-hydroxyl group in 224 to the 3-position as in 226 was less tolerated. Compound 225 significantly reduced immobility in the rat FST model at a dose of 10 mpk with efficacy comparable to 5. Compound 229 showed lower triple reuptake inhibitory potency than enantiomer 224. SAR exploration of ring substituents resulted in trisubstituted pyran analogues D-142 (230) and D-165 (231), in which both stereochemistry and regioselectivity played an important role.<sup>141</sup> The (–)-2S,4R absolute configuration and a *cis*-relationship between the benzhydryl moiety and the amino group were required for acceptable potency at SERT and NET. Compound 230 demonstrated high potency with a significant antidepressant

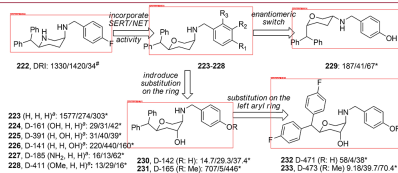


Figure 41. Evolution of DAT-selective piperidines into triple-acting pyran analogues with the introduction of SERT/NET affinity by isomeric heterocyclic replacement.

## Figure 5. Examples of molecule segmentation in medical literature

### 图 5. 医药文献中分子式分割效果示例

Table 2. Comparative experimental results of graph similarity algorithms on AIDS and LINUX datasets

表 2. 图形相似度算法在 AIDS 数据集和 LINUX 数据集上对比实验结果

	AIDS		LINUX	
方法	MSE	运行时间	MSE	运行时间
GED	0.395	0.001	0.060	0.000799
HB-GED	0.320	0.006	0.017	0.003478

Table 3. Comparative experimental results of retrieval on the WIKI dataset

表 3. 在 WIKI 数据集上检索的对比实验结果

方法	mAP@2	mAP@4	mAP@6	mAP@8	mAP@10
GED	0.960	0.878	0.778	0.680	0.613
HB-GED	0.973	0.906	0.906	0.726	0.695

## 5. 总结与展望

基于内容的论文精准推荐是大势所趋，其中，非文字类型的内容对象不仅是重要的推荐依据，其理解算法也是推荐系统中的技术难点。本文探索了依据分子式图形相似度的医药文献推荐途径，实验结果表明，该途径不仅可行，而且也是基于文本类型内容推荐算法的重要拓展。但距离医药领域的实际应用

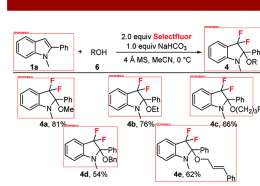
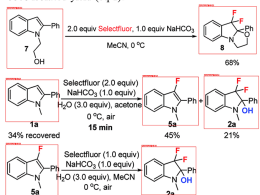


Figure 3. Difluorination of 1a using alcohols as nucleophiles.

tricyclic tetrahydrooxazolo[3,2-*a*]indole was obtained in 68% isolated yield (eq 1).<sup>15</sup>



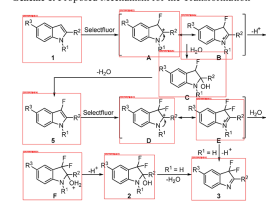
During the process of optimizing the conditions, when the reaction of 1a was stopped at 15 min in acetone as the solvent, the desired difluorinated product 2a (21%) and monofluorinated product 5a (45%) were obtained with 34% of 1a recovered (eq 2). Furthermore, 5a was then resubjected to the standard reaction conditions (only using 1.0 equiv of Selectfluor), which produced 2a in 80% yield (eq 3). These results indicate that the monofluorinated product serves as an intermediate in the transformation.

On the basis of the above results and information from the literature,<sup>16,18</sup> the mechanism of this transformation is proposed (Scheme 1). Initially, reaction of 1 with Selectfluor yields the unstable 3-fluoroindolinolium cation A or its resonance 3-fluoroindolenine cation B, and then the proton is extracted by the base (NaHCO<sub>3</sub> or DABCO from Selectfluor<sup>16</sup>) quickly to give the 3-fluoroindole 5. Alternatively, the unstable 3-fluoroindolinolium cation A or its resonance 3-fluoroindolenine cation B can also be attacked by H<sub>2</sub>O to furnish the

(15) During the revision of this manuscript, a similar intramolecular reaction has been reported; see: Lonzo, O.; Biondi, G.; del Campo, T. M.; Thompson, A. L.; Giuffridè, G. T.; Bettari, M.; Walker, M.; Borman, R.; Gouverneur, V. *Angew. Chem., Int. Ed.* DOI: 10.1002/anie.201101151.

(16) Radwan-Okocinska, K.; Palacios, F.; Kafarski, P. *J. Org. Chem.* 2011, 76, 1170.

### Scheme 1. Proposed Mechanism for the Transformation



3-fluoroindolinol-2-ol. C. Dehydration of C gives the 3-fluoroindole 5, which then undergoes the same process to produce the unstable 3,3-difluoroindolinolium cation D or 3,3-difluoroindolenine cation E. Finally, the carbon cation of D or E is attacked by H<sub>2</sub>O to produce 3,3-difluoroindolinol-2-ol 2. When the substituent group R<sup>2</sup> is an aryl group, which can stabilize the carbon cation, the substrates would generally achieve good yields as shown in Table 2. In contrast, when the R<sup>2</sup> was changed to an ester group, the carbon cation was not stable enough and the yield dropped (21, Table 2). We also tried the methyl group, although we could obtain the desired product, it was not stable enough at room temperature and decomposed quickly. In the cases of unprotected indole derivatives (R<sup>1</sup> = H) as the substrates, the direct deprotonation of intermediate D or E, or dehydration of 2 leads to the formation of 3. When alcohols are used as nucleophiles instead of H<sub>2</sub>O, the reactions proceed through a similar mechanism.

In summary, we have developed an efficient method for the synthesis of 3,3-difluoroindolinol-2-ol. In this method, the indole ring was difluorinated highly regioselectively at the C3 carbon site with equivalent (not excess) Selectfluor. Additionally, mild conditions and practical convenience would make it a valuable synthetic tool in organic chemistry. When alcohols were used as the nucleophiles instead of H<sub>2</sub>O, the reactions proceeded through a similar mechanism.

The study of the bioactivity of these novel compounds is ongoing in our laboratory.

**Acknowledgment.** Financial support from Peking University, the National Science Foundation of China (2087-2003), and the National Basic Research Program of China (973 Program 2009CB825300) is greatly appreciated. We also thank Guolin Wu in this group for reproducing the results of 2g and 4a.

**Supporting Information Available.** Experimental procedures, characterization data, and X-ray crystallographic data. This material is available free of charge via the Internet at <http://pubs.acs.org>.



要求而言, 该项研究还存在着图形提取完整度低、图形大小差异度大、以及分子图与其他类型插图不易区分等问题。在推荐算法方面, 也有一些新方向值得关注, 例如, 分子的表示具有多模态的特征, 即图像、图形和字符串等, 可借鉴计算机视觉领域中多模态或跨模态的融合方法, 进一步提升计算分子相似度的准确率。再如, 由于分子研究的门槛较高, 标注成本亦同步升高, 大规模的分子标注数据集的建设令人望而却步, 特别是分子相似度的人工标注数据集几乎还是空白, 因此, 在加大数据集建设力度的同时, 也可以考虑挖掘分子集合自身的内在特征, 更多地利用自监督或无监督的学习模型, 获得更准确的分子表征。

## 基金项目

本研究受到北京市自然科学基金——海淀原始创新联合基金资助项目(项目编号: L192024)的资助。

## 参考文献

- [1] Wu, S., Sun, F., Zhang, W. and Cui, B. (2022) Graph Neural Networks in Recommender Systems: A Survey. *ACM Computing Surveys*, **55**, Article No. 97. <https://doi.org/10.1145/3535101>
- [2] Li, H. (2014) Learning to Rank for Information Retrieval and Natural Language Processing. In: Hirst, G., Ed., *Synthesis Lectures on Human Language Technologies*, 2nd Edition, Springer, Berlin, 121 p. <https://doi.org/10.2200/S00607ED2V01Y201410HLT026>
- [3] Hui, K., Yates, A., Beberich, K. and Melo, G.D. (2018) Co-PACRR: A Context-Aware Neural IR Model for Ad-Hoc Retrieval. *Proceedings of the 11th ACM International Conference on Web Search and Data Mining (WSDM '18)*, Los Angeles, 5-9 February 2018, 279-287. <https://doi.org/10.1145/3159652.3159689>
- [4] Xiong, C., Dai, Z., Callan, J., Liu, Z. and Power, R. (2017) End-to-End Neural Ad-Hoc Ranking with Kernel Pooling. *Proceedings of the 40th International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR '17)*, Tokyo, 7-11 August 2017, 55-64. <https://doi.org/10.1145/3077136.3080809>
- [5] Dai, Z. (2020) Neural Matching and Importance Learning in Information Retrieval. Ph.D. Thesis, Tsinghua University, Beijing.
- [6] Gao, L., Dai, Z., Chen, T., Fan, Z., Durme, B.V. and Callan, J. (2021) Complementing Lexical Retrieval with Semantic Residual Embedding. In: Hiemstra, D., Moens, M.F., Mothe, J., Perego, R., Potthast, M. and Sebastiani, F., Eds., *Advances in Information Retrieval. Lecture Notes in Computer Science*, Vol. 12656, Springer, Cham, 146-160. [https://doi.org/10.1007/978-3-030-72113-8\\_10](https://doi.org/10.1007/978-3-030-72113-8_10)
- [7] Xiong, L., Xiong, C., Li, Y., Tang, K.F., Liu, J., Bennett, P., Ahmed, J. and Overwijk, A. (2021) Approximate Nearest Neighbor Negative Contrastive Learning for Dense Text Retrieval. *The 9th International Conference on Learning Representations (ICLR 2021)*, Virtual Event Austria, 3-7 May 2021, 16 p. <https://openreview.net/pdf?id=zeFrfgYzLn>
- [8] Lin, S.C., Yang, J.H. and Lin, J. (2020) Distilling Dense Representations for Ranking Using Tightly-Coupled Teachers. ArXiv Preprint arXiv: 2010.11386.
- [9] Qian, Y., Santus, E., Jin, Z., Guo, J. and Barzilay, R. (2018) GraphIE: A Graph-Based Framework for Information Extraction. ArXiv Preprint arXiv: 1810.13083.
- [10] Trabelsi, M., Chen, Z., Davison, B.D. and Heflin, J. (2021) Neural Ranking Models for Document Retrieval. *Information Retrieval Journal*, **24**, 400-444. <https://doi.org/10.1007/s10791-021-09398-0>
- [11] Zhang, Z., Bu, J., Ester, M., Li, Z., Yao, C., Yu, Z. and Wang, C. (2021) H2mn: Graph Similarity Learning with Hierarchical Hypergraph Matching Networks. *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*, Virtual Event Singapore, 14-18 August 2021, 2274-2284. <https://doi.org/10.1145/3447548.3467328>
- [12] Coupette, C. and Vreeken, J. (2021) Graph Similarity Description: How Are These Graphs Similar? *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining (KDD '21)*, Virtual Event Singapore, 14-18 August 2021, 185-195. <https://doi.org/10.1145/3447548.3467257>
- [13] Raveaux, R. (2021) On the Unification of the Graph Edit Distance and Graph Matching Problems. *Pattern Recognition Letters*, **145**, 240-246. <https://doi.org/10.1016/j.patrec.2021.02.014>
- [14] Riba, P., Fischer, A., Lladós, J. and Fornés, A. (2020) Learning Graph Edit Distance by Graph Neural Networks. *Pattern Recognition*, **120**, 108-132. <https://doi.org/10.1016/j.patcog.2021.108132>
- [15] Ling, X., Wu, L., Wang, S., Ma, T., Xu, F., Liu, A.X., Wu, C. and Ji, S. (2021) Multilevel Graph Matching Networks

- for Deep Graph Similarity Learning. *IEEE Transactions on Neural Networks and Learning Systems*.  
<https://doi.org/10.1109/TNNLS.2021.3102234>
- [16] Ren, S., He, K., Girshick, R. and Sun, J. (2017) Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, **39**, 1137-1149.  
<https://doi.org/10.1109/TPAMI.2016.2577031>
- [17] He, K., Gkioxari, G., Dollár, P. and Girshick, R. (2017) Mask R-CNN. *IEEE International Conference on Computer Vision (ICCV 2017)*, Venice, 22-29 October 2017, 2980-2988. <https://doi.org/10.1109/ICCV.2017.322>
- [18] Lin, T.-Y., Goyal, P., Girshick, R., He, K. and Dollár, P. (2017) Focal Loss for Dense Object Detection. *Proceedings of the IEEE International Conference on Computer Vision (ICCV 2017)*, Venice, 22-29 October 2017, 2999-3007.  
<https://doi.org/10.1109/ICCV.2017.324>
- [19] Tian, Z., Shen, C., Chen, H. and He, T. (2019) FCOS: Fully Convolutional One-Stage Object Detection. *IEEE/CVF International Conference on Computer Vision (ICCV 2019)*, Seoul, 27 October-2 November 2019, 9626-9635.  
<https://doi.org/10.1109/ICCV.2019.00972>