

深度学习在反应和逆合成预测中的统一观点：现状与未来挑战

Ziqiao Meng^{1*}, Peilin Zhao^{2†}, Yang Yu² and Irwin King^{1†}

¹香港中文大学

²腾讯人工智能实验室

{zqmeng, king}@cse.cuhk.edu.hk, masonzhao@tencent.com

摘要

反应和逆合成预测是计算化学的基本任务，最近引起了机器学习界和药物发现界的关注。为了解决这些问题，人们提出了各种深度学习方法，其中一些已经取得了初步成功。在本调查中，我们对这些方法进行了比较。

我们对基于深度学习的高级反应和逆合成预判模型进行了深入研究。我们总结了设计机制、**SOTA的不足**

然后，我们讨论了当前解决方案的局限性以及问题本身所面临的挑战。然后，我们讨论了当前解决方案的局限性和问题本身的挑战。最后，我们提出了促进未来研究的可行方向。据我们所知，本文是第一份全面系统的研究报告，旨在提供对反应和逆合成预测的统一认识。

1 引言

药物研发对人类的医疗保健至关重要，但这一过程却是出了名的劳动密集型和低成本。正如埃鲁姆定律（Eroom's law）[Scannell *et al.*, 2012]所指出的，随着时间的推移，新药探索的速度越来越慢，成本越来越高。因此，利用机器学习技术来加速药物发现过程是自然而然的，也是意义重大的。近年来，随着深度学习的兴起，使用深度学习方法来促进药物发现的不同阶段已成为一种普遍现象。在这些阶段中，反应预测和逆合成预测是可以从深度学习工具中受益的两个基本步骤。

在实际生产环境中，化学家通常致力于设计合成路线，通过一系列化学反应获得目标分子。**一种常见的策略是将目标分子分解成更容易合成的简单前体结构，这一过程被称为逆合成分析。**自动规划

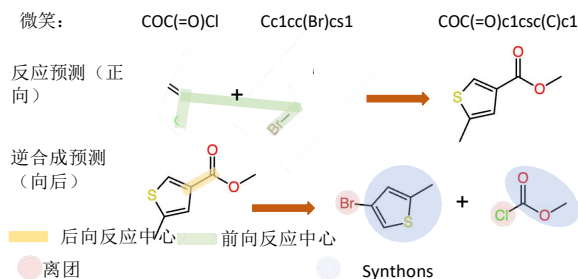


图 1: 该图说明了反应和逆合成预测的问题表述

使用深度学习的逆向合成过程对于发现和优化合成路线至关重要。**逆合成规划涉及两个子任务：多步骤逆合成规划和单步骤逆合成预测。**在本研究中，我们将重点讨论后者，因为**多步规划通常是作为搜索问题来处理的**，这与反应预测和单步逆合成预测所提出的**条件结构预测**问题有着本质区别。在本研究的其余部分，我们将单步逆合成预测简称为逆合成预测。**反应预测是有机合成分析的另一项关键任务。**一个强大的反应预测模型可以帮助我们深入了解生化反应的内在机理，并生成虚拟反应以扩展用于逆合成规划的数据库。总之，**反应预测和逆合成预测是相互关联、相互促进的。**关于反应和逆合成预测的研究有很多。Engkvist [Engkvist *et al.*, 2018]概述了从量子计算到深度学习工具的各种反应预测计算方法。然而，该调查缺乏每种方法的细节，对基于深度学习的方法的覆盖也很有限。另一方面，[Dong 等人, 2021]综述了一些基于深度学习的逆合成规划方法和数据集，但它落后于最先进的逆合成预测研究。值得注意的是，DualTF[Sun 等人, 2021]从基于能量的模型角度为逆合成预测提供了一个统一的框架，但它并没有提供对反应和逆合成的统一理解。

*这项工作是在孟子乔在腾讯人工智能实验室实习时完成的。

†通讯作者：赵培林和欧文-金

也没有具体讨论每种方法的局限性和主要挑战。总体而言，现有文献对以下问题缺乏全面统一的认识

先进的反应和逆合成预测模型。

与之前提到的调查和作品相比，我们的

调查首次统一了反应和逆合成预测的表述。针对这两个问题，我们从不同角度系统地讨论了每种方法的优缺点。此外，我们的研究还提出了新的挑战 and 局限性，而这些挑战和局限性在之前的研究中并未明确提出。最后，基于目前的现状，我们列出了未来进一步改进的几个方向，并进行了详细分析。

2 前期准备和问题提出

反应预测和逆合成预测是双重的

的任务。它们也分别被称为前向反应预测和后向反应预测。对

在深度学习中，这两项任务都被表述为条件生成任务。在本节中，我们首先介绍这两个任务的问题表述。然后，我们将介绍理解不同方法所需的基本背景知识。

分子式。化学分子 M 可以用两种主要数据格式表示：**SMILES 字符串**和**分子图**。(1) 对于 SMILES 格式，一个分子结构 M 被描述为一连串的字符串。

$M := m_1 m_2 \dots m_L$ ，其中 L 表示字符串的总长度。序列表示二维分子结构的生成树，每个字符

i 表示结构元素，如原子 e_i -

(2) 分子也可以抽象为无向图。(2) 分子也可以抽象为一个无向图

$G = \{V, E\}$ ，其中 $V = \{v_1, \dots, v_n\}$ 表示 n

原子， $E = \{e_1, \dots, e_m\}$ 表示 m 条边的集合。每个节点与特征向量 $h_i \in \mathbb{R}^d$ 相关联，其中包含

原子信息，如芳香性和电荷。然后，我们就有了一个特征矩阵 $H \in \mathbb{R}^{n \times d}$ ，其中包含所有--

原子信息。邻接矩阵 $A \in \mathbb{R}^{n \times n \times c}$ 描述了 M 的拓扑结构，其中 A_{ijk} 表示原子 i 和原子 j 之间是否存在 k 类型的化学键。上述两种格式可以轻松表示多个分子。对于 SMILES 格式，多个 SMILES 字符串可以用句号 "." 连接成一个 SMILES 序列。对于分子图，一组分子被视为一个单独的断开图，每个分子都是一个独立的连通组件。

定义 1 (反应预测) 给定一组 N 个反应物

分子 $\{M^R\}_{i=1}^N$ ，目标是预测 M

可能的产品分子 $\{M^P\}_{i=1}^M$ 。

定义 2 (逆合成预测) 给定一组 M

产品分子 $\{M^P\}_{i=1}^M$ ，目标是预测一组

N 个反应物分子 $\{M^R\}_{i=1}^N$ ，可导致 $\{M^P\}_{i=1}^M$ 。

请注意，在实际应用中， $M = 1$ ，因为在公共基准数据集中只记录了主要产物。不幸的是，这个问题使得逆合成预测比反应预测更加困难，因为 $M < N$ 需要

重新合成来连接新出现的原子，从而导致更大的组合搜索空间。一般来说，再合成预测的目的是建立条件概率模型。

分布 $P(M^R | \{M^P\}_{i=1}^M)$ 而逆合成预测旨在模拟分布 $P(\{M^R\}_{i=1}^N | M^P)$ 。

反应中心和反应模板。在反应预处理中，反应中心 C 是原子对 $C = \{v_i, v_j\}$ 的子集。

$\{(v_i, v_j)\} \subseteq V \times V$ ，当化学反应发生时，这些键的类型会发生改变。在逆合成预测中，反应中心 C 被定义为现有键 $C = \{e_i\} \subseteq E$ 的子集，可以通过修改这些键来获得更简单的结构。反应模板库 T 是一组从大型化学反应数据库中提取的反应子图规则。反应模板 TE T 是从相应的化学反应数据库中提取的子图式。

具体来说， $T := t^R + t^R + \dots + t^R \rightarrow t_1^P, t_2^P, \dots, t_N^P$ ，其中 t^R 表示第 i 个反应物内部的子图模式、

和 t^P 表示乘积内部的子图模式。

原子映射 反应预测和逆合成预测都遵循原子映射原则。该原则规定，反应物/生成物中的每个原子在生成物/反应物中都有一个对应原子。这种基本的一对一映射关系从物理上限制了反应空间，并决定了化学反应主要是键的断裂和键的形成。

Synthon 和 Leaving Group。目标分子可以是分解为一组合成子 $S = \{G^S\}_{i=1}^N$ ，它们是

更简单的前体子结构，只需很少的额外键连接就能构成目标分子。请注意， S 所覆盖的原子集 V 与目标分子完全相同。

分子。离去基团 $L = \{G^L\}_{i=1}^N$ 是一个原子团

或原始反应物中的亚结构，反应发生后不会出现在目标分子中。简而言之，合子 S 和离去基团 L 可以形成反应物 $\{GR\}_{i=1}^N$ 。

评估指标。反应和逆合成前采用最高精确度来评估模型的性能。

准确度。Top- k 准确率是指

的前 k 个预测集合中的真实乘积。
分子。只要前 k 个预测产品中包括地面实况的主要产品，就可以算作正确预测。通常， k 的范围为 $\{1, 2, 3, 5, 10\}$ 。

3 反应和逆合成预测的深度学习

在本节中，我们将讨论不同的方法，并将其分为四类，即基于模板的方法、基于数据的方法、基于数据的方法和基于数据的方法。

在此基础上，我们对基于序列和图形的自回归模型、基于序列的自回归模型和基于图形的自回归模型进行了分析。

图模型、基于图形的两阶段模型和基于图形的非

自回归模型。具体而言，我们将介绍其学习机制的设计决策、弱点和优势。

3.1 基于模板的方法

基于模板 (TB) 的方法主要是利用反应模板库来推断可能的反应中心。模板法是模仿人类专家进行化学推理的方法。假设我们有一个反应模板库 $T =$

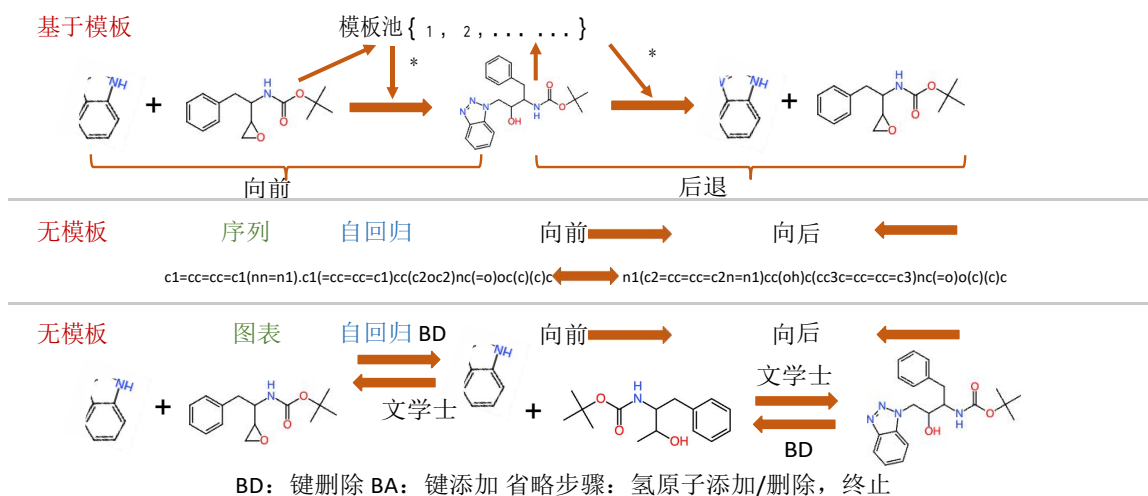


图 2: 该图展示了基于模板的模型、基于序列的自回归模型和基于图形的自回归模型。

$\{T_1, T_2, \dots\}$, 则 TB 方法的推理过程旨在选出最佳模板, 如下所示:

$$T^* = \arg \max_{T \in T} P(T = T^* | \{M_i\}_{i=1}^{RN}),$$

$$M^{P^*} = \arg \max_{M_i^P \in MP} P(M^P = M_i^P | T^*), \quad (1)$$

其中 M^P 是候选产品分子库。公式 (1) 描述了选择最佳模板 T^* 并根据给定的 T^* 生成可能的产品 M^{P^*} 的步骤。最关键的部分是将分子与模板进行匹配。基于某种相似性度量。具体来说, $P(T^* | \{M_i\}_{i=1}^{RN})$ 和 $P(M_i^P | T^*)$ 按以下方法进行评估:

$$p(t | \{m_i\}_{i=1}^{RN}) = \sigma_{T_j \in T} \frac{\exp(\text{sim}(T_i, \{M_i\}_{i=1}^{RN}))}{\exp(\text{sim}(T_j, \{M_i\}_{i=1}^{RN}))},$$

$$P(M_i^P | T^*) = \frac{\exp(\text{sim}(T^*, M_i^P))}{\sum_{M_j^P \in MP} \exp(\text{sim}(T^*, M_j^P))} \quad (2)$$

其中, $\text{sim}(-, -)$ 表示相似性度量函数。它可以是分子嵌入与深度神经网络获得的模板嵌入之间的一个简单内积, 也可以是一个更复杂的相似性函数, 由图匹配算法。

图匹配算法。TB 模型生成具有相同的

学习机制, 但方向相反:

$$T^* = \arg \max_{T \in T} P(T = T^* | M_i^P),$$

$$P(M^R | T^*) = \prod_{i=1}^N P(M_i^R | T^*), \quad (3)$$

$$\hat{M}_i^R = \arg \max_{M_j^R \in MR} P(M_i^R = M_j^R | T^*)$$

其中 M_i^R 表示候选反应物分子库, $M_i^R := m_1^R, m_2^R, \dots, m_N^R$ 。第二个和第三个等式在公式 (3) 中表示每个预测的反应物分子 M_i^R 与得分最高的子图模板 T^R 匹配。

优点(1) TB 方法是可靠的, 因为它们使用的是提取的人类知识, 而人类知识总能为预测提供很好的解释。(2) 训练和输入肺结核治疗方法的发散过程相对简单, 因此易于领域专家操作。

缺点(1) TB 方法的性能高度依赖于模板数据库的规模。因此, 模板数据库必须经常更新, 这显然非常昂贵。(2) TB 方法对域外未见反应的普适性较差。(3) 模板是

提取局部子图规则, 而忽略全局信息。

因此, TB 方法无法捕捉到全局信息的相互作用。因此, TB 方法无法捕捉到全局信息的交互作用, 很容易产生错误信息。基于本地规则的预测。TB 方法的说明

如图 2 所示。

3.2 基于序列和基于图形

自回归模型

基于序列的自回归 (SAR) 模型在前向和后向预测中被广泛采用。它认为这两个问题都是神经机器翻译问题。对于反应预测, 输入是反应物的 SMILES 字符串 $M^R := m_1^R m_2^R \dots m_N^R$, 长度为 L_1 , 输出是 SMILES。

长度为 L 的产品串 $M^P := m_1^P m_2^P \dots m_{L_2}^P$ 。

输入源和输出目标是翻转的, 以便进行回溯预测。具体来说, SAR 模型正在估计以下条件概率分布:

$$P(M^P) = \prod_{i=1}^N P(m_i^P | m_{<i}^P, M^R),$$

$$P(M^R) = \prod_{i=1}^N P(m_i^R | m_{<i}^R, M^P) \quad (4)$$

其中, $P(m_i^P | m_{<i}^P, M^R)$ 和 $P(m_i^R | m_{<i}^R, M^P)$ 是近似值。Vaswani 等人, 2017 年。每个生成步骤对标记空间进行贪婪搜索, 选出最佳标记。要生成前 k 个候选者, 我们只需

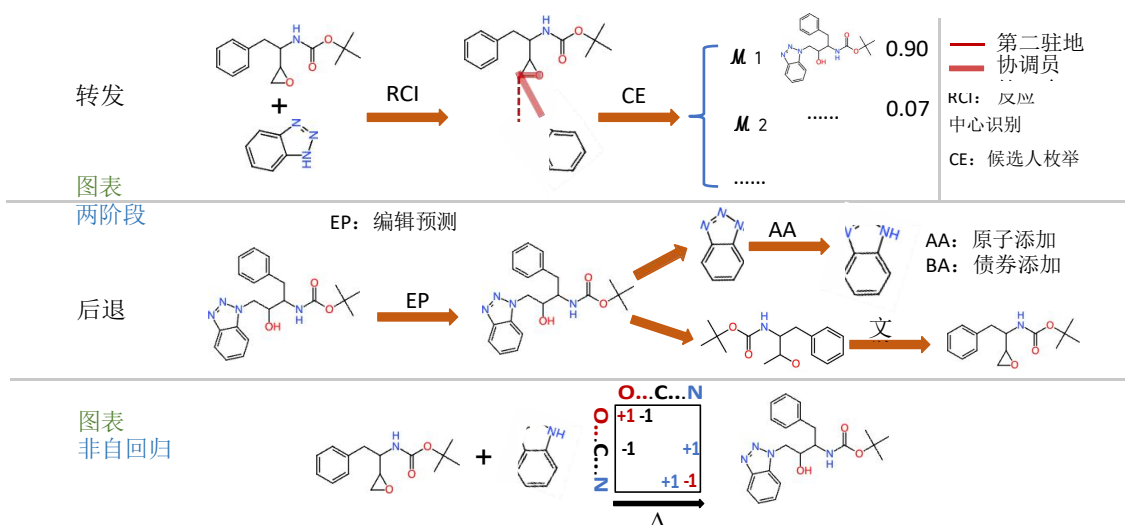


图 3: 该图展示了基于图形的两阶段模型和基于图形的非自回归模型。

需要对贪婪结果进行波束搜索。基于图的自回归 (GAR) 模型具有类似的学习机制，但生成序列定义不同。GAR 模型首先定义一个动作空间 π ，其中包括多个编辑动作，如原子添加/删除、键添加/删除和终止。这些动作的序列会将反应物/生成物转化为相应的产品/反应物。因此，他们估计概率分布：

$$P(G^P) = \prod_{i=1}^{L-1} P(\pi_i | \pi_{<i}), \quad (5)$$

$$P(G^R) = \prod_{i=1}^{L-1} P(\pi_i | \pi_{<i}, G^{P,<i}),$$

其中， $G^{R,<i}$ 和 $G^{P,<i}$ 表示已编辑反应物的状态和前 $i-1$ 步的产品。取样绘制过程与合成孔径雷达模型的绘制过程完全相同。

每个步骤 i 都会选择最优行动 π^* ，并将最优行动 π 应用于每个步骤。

将反应物 $G^{R,i-1}$ 和生成物 $G^{P,i-1}$ 分别转化为下一状态 $G^{R,i}$ 和 $G^{P,i}$ 。当预测所有生成步骤都结束时，最佳操作是“终止”。请注意，公式 (4) 和公式 (5) 都可以通过马尔可夫假设进行简化。

优点(1) 自回归模型不需要原子映射信息。(2) 自回归模型具有很强的自然采样过程，如光束搜索。(3) 序列基于模型的建模可以直接利用一些成熟的技术来自自然语言处理。

缺点(1) 自回归模型只能逐步生成预测，效率很低。(2)

自回归建模而无需元明主生成

分子。然而，分子生成顺序并不明确。(3) 基于序列的建模需要数据增强技术来提高性能。图 2 是 SAR 模型和 GAR 模型的示意图。

3.3 基于图形的两阶段模型

反应对于基于图的两阶段反应预测模型，他们将反应预测分为两个阶段，即反应中心识别阶段和候选物排序阶段。对于反应中心识别，其目的是选择具有高反应性得分的原子对：

$$C^* = \text{top-}k(\{s(v_i, v_j | \{G_i^{RN}\}_{i=1}^N)\}_{ij}), \quad (6)$$

其中， $s(-)$ 表示输出分数范围在 (0, 1)。根据已确定的反应中心 C^* ，将在 $P = \{G^P, G^P, \dots\}$ 中枚举出一组可能态产物 $G_2^P = \{G^P, G^P, \dots\}$ 。

通过手写规则或反应模板，以组合方式进行。第二阶段是学习如何对 G^P 中的可生成产物进行排序。为了对生成的产品进行排序，每对 $(\{G_i^{KN}\}_{i=1}^N, G_i^P)$ 的评估得分应如下：

$$s(\{G_i^{KN}\}_{i=1}^N, G_i^P) = \sigma(f(\{G_i^{KN}\}_{i=1}^N, G_i^P)), \quad (7)$$

其中， $f(-, -)$ 可以是复杂神经网络， $\sigma(-)$ 记为 sigmoid 函数。

回溯合成。对于基于图的两阶段逆合成预测模型，他们将逆合成分为两个阶段，即编辑预测阶段和合成阶段。在编辑预测阶段，他们选择一些预测反应性得分较高的现有边缘进行分解：

$$C^* = \text{top-}k(\{s(e_i^P)\}_{i=1}^m), \quad (8)$$

在获得预测的编辑中心 C^* 之后，我们通过打破预测的编辑中心来分解目标分子，即将产生一组合成子 $S = \{G^S\}$ ， $(N = k + 1)$ 。 $s_N(N = k + 1)$ 。

这种方法的基本假设是，合成子的数量与反应物的数量相同。在合成子完成阶段，它采用以下分布模型：

$$p(g_i^R | g_i^S) = p(g_i^L | g_i^S). \quad (9)$$

上式 (9) 表明, 同义词补全等同于将离群组 G^L 附加到相应的同义词 G^S 上。离群组附加可以等同于一个自回归条件生成问题或一个带有预定义二分法的分类问题。在预测的离去基团 G^L 通过预测的离去基团 G^L , 每个预测的 G_i 可以通过将 G_i 附加到相应的同义词 G^S 来恢复。

优点(1) 基于图的两阶段模型的推理过程与化学家的推理过程非常相似。(2) 基于图的两阶段模型将一项艰巨的任务分成两项更容易处理的简单任务。

缺点(1) 一些两阶段模型, 如 WLDN [Jin 等人, 2017] 需要昂贵的手工组合枚举。(2) 两阶段模型的整体性能受到每个阶段瓶颈的限制。上述方法的示意图如图 3 所示。

3.4 基于图形的非自回归模型

只有反应预测在 NERF [Bi 等人, 2021 年] 实现了基于图的非自回归模型, 目前还没有用于逆合成的非自回归模型。NERF 将问题重新表述为电子再分布模型。以往的方法采用三维邻接矩阵 $A \in \mathbb{R}^{n \times n \times c}$, 而 NERF 则采用二维邻接矩阵, 将单击类型编码转换为标量值, 从而使 $A \in \mathbb{R}^{n \times n}$ 。形式上, A_{ij} 是一个范围为 $[0, 3]$ 的标量值, 代表原子 i 和原子 j 之间的共享电子数 (键数)。请注意, 芳香键表示为 $A_{ij} = 1$, 原子 i, j 将被标记为芳香原子。NERF 的主要思想是通过结合自注意映射来预测 ΔA 。NERF 采用条件变异自动编码器 [Sohn 等人, 2015] (CVAE) 架构, 通过引入潜变量 z 来近似 $P(G^P | G^R)$ 。CVAE 不是直接最大化对数似然 $\log P(G^P | G^R)$, 而是最大化其证据下限 (ELBO):

$$\log P(G^P | G^R) \geq \mathbb{E}_{P(z|G, G)} [\log P(G^P | G^R, z)] \quad (10)$$

$$- \text{KL}(q(z|G^P, G^R) || P(z|G^R)),$$

其中 $q(z|G^P, G^R)$ 是反应编码器, 以反应 (G^R, G^P) 为输入, 以低维表示 h^z 为输出; $P(G^P | G^R, z)$ 是产品解码器, 以反应物 G^R 和潜在嵌入 h^z 为输入; $P(z|G^R)$ 表示潜在变量 z 的先验分布。KL 项是最小化 $q(z|G^P, G^R)$ 与 $P(z|G^R)$ 之间的差距。 $q(z|G^P, G^R)$ 的骨干网络结构是 GNN 和变压器的组合。在此架构下, G^R 和 G^P 将分别投射到嵌入 h^R 的反应物和嵌入 h^P 的产品。交叉注意层用于将 h^R 映射到潜在 h^z , h^P 作为训练期间的教师强迫。条件潜在嵌入的推导结果是 $h^{\wedge z} = h^R + h^z$ 。然后在 $h^{\wedge z}$ 上应用自我注意机制, 推导出两个电子再分布矩阵 W^+ 和 W^- , 分别用于键增加和键减少。那么 $\Delta A^{\wedge} = W^+ - W^-$ 和 $A^{\wedge P} = A^R + \Delta A^{\wedge}$ 。

优点(1) 非自回归模型可以进行准线采样, 采样速度比自回归模型快得多。(2) 非自回归模型的准确率已达到一流的 top-1, 这表明非自回归解码器在反应建模方面非常强大。(3) 非自回归模型不需要预先确定生成顺序。

缺点(1) 非自回归模型的不确定性建模非常棘手。top-k 取样过程不像自回归模型中的波束搜索那么自然。(2) 非自回归模型依赖于原子映射信息, 而这也需要额外的配准算法。图 3 展示了 NERF 学习机制。

4 局限与挑战

在本节中, 我们将讨论当前解决方案中存在的一些重要限制和挑战。

4.1 辅助产品

反应。美国专利商标局的公开基准数据集中缺少副产品, 导致监督信号不完整。特别是在非自回归模型中, 副产品的缺失会使电子再分布矩阵违反守恒规则, 从而导致反应空间不完全受限。如何完成和推断这种缺失的形成是反应预测的一个基本挑战。

逆合成。逆合成直接使用 USPTO 中的反应数据, 因此所有单步分析只包含一个单一结果。结果侧的侧产物缺失可能不会影响逆合成预处理的一般过程。不过, 侧产物仍能提供有关离去基团的重要信息。

4.2 数据集的局限性

反应。USPTO-479K 数据集有两个主要问题。首先, 反应类型非常不平衡。表示为 Bi [Bi 等人, 2021 年] 中, 线性拓扑结构的反应是主要的反应类型, 而环状拓扑结构的反应很少。

数据集。如何从罕见的反应中学习可迁移的知识是反应建模的重要一课。其次, 在实际应用中, 同一组反应物在不同的物理条件下会产生不同的产物, 这就是反应预测中的多模态。多模态可以为条件生成模式提供丰富的信息, 从而生成有效和多样的候选产物。然而, USPTO-MIT 中的大多数反应都是一对一映射, 这意味着同一组反应物只能生成唯一的主要产物。

回溯合成。USPTO-50K 数据集有两大局限性。首先, USPTO-50K 的规模不够大, 只包含 50K 个逆反应。考虑到最近的许多方法在 top-k 精确度上只有微小的数值差异, 目前的小规模数据集不足以测试模型的能力。其次, 目前的数据集会使编辑预测和离组分析产生偏差。大多数逆反应只有一个单一的编辑

	模板使用	方法	端到端	一代人	图表/序列	Top-1	前五名
反应预测（正向）	基于模板	NN 反应 [Wei 等人, 2016]	×	NA	序列	NR	NR
		NeuralSym [Segler and Waller, 2017]	×	NA	序列图	NR	NR
		Symbolic [Qian <i>et al.</i>]	×	NA		90.4	95.0
			×	NA		90.8	96.3
	模板 - 免费	WLDN [Jin <i>et al.</i>]	×	NA	图表	79.6	89.2
		全球贸易点网络 [Do <i>et al.</i>]	✓	AA	图表 图表	83.2	86.5
		MEGAN [Sacha 等人, 2020]	✓	AA	序列 序列	89.3	95.6
		ELECTRO [Bradshaw 等人, 2019]	✓	AA	序列 序列	77.8	94.7
		Motif-Reaction [Zhao 等人, 2022]	✓	AA	图表 图表	91.0	95.7
		MT-base [Schwaller 等人, 2019]	✓	AA	图表	88.8	94.4
		MT [Schwaller 等人, 2019]	✓	AA		90.4	95.3
		Graph2SMILES [Tu 和 Coley, 2021]	✓	A		90.3	94.8
		Chemformer [Irwin 等人, 2022]	×	NA		91.3	93.7
		Stranformer [Lee 等人, 2019]	✓	NA		NR	NR
		ReactionT5 [Lu 和 Zhang, 2022]	✓			88.9	95.2
		Aug-Transformer [Tetko 等人, 2020]	✓			90.6	96.1
NERF [Bi 等人, 2021]				90.7	93.7		
反应汇[Meng 等人, 2023]				91.3	94.0		
逆合成预测（向后）	基于模板	RetroSim [Coley <i>et al.</i>]	×	NA	序列	52.9/37.3	81.2/63.3
		RetroComposer [Yan 等人, 2022]	×	A	图表 图表	65.9/54.5	89.5/83.2
		NeuralSym [Segler和Waller, 2017]	×	NA	序列	55.3/44.4	81.4/72.4
		GLN [Dai 等人, 2020]	×	NA		64.2/52.5	85.2/75.6
		LocalRetro [Chen and Jung, 2021]	×	A		63.9/53.4	92.4/85.9
		DualTB [Sun <i>et al.</i>]	×	NA		67.7/55.2	88.9/80.5
	模板 - 免费	MEGAN [Sacha <i>et al.</i>]	✓	A	图表	60.7/48.1	87.5/78.4
		AutoSynRoute [Lin 等人, 2020]	✓	aaa	序列图 序	54.6/43.1	80.2/71.8
		SCROP [Zheng 等人, 2020]	✓	ana	列图 序列	59.0/43.7	78.1/65.2
		低压变压器 [Chen 等人, 2019]	✓	aaa	图 序列图	NR/40.5	NR/72.8
		DualTF [Sun 等人, 2021]	✓	aaa	序列图 序	65.7/53.6	84.7/74.6
		GET [Mao 等人, 2021]	✓	aaa	列图 序列	57.4/44.9	74.8/62.4
		Retroprime [Wang 等人, 2021]	✓		图 序列图	64.8/51.4	85.0/74.0
		GTA [Seo 等人, 2021]	✓		序列图 序	NR/51.1	NR/74.8
		G2Gs [Shi 等人, 2020]	×		列图 序列	61.0/48.9	86.0/72.5
		RetroXPert [Yan 等人, 2020]	✓		图 序列图	62.1/50.4	78.5/62.3
		Retroformer [Wan 等人, 2022]	✓		序列图 序	64.0/53.2	86.7/76.6
		Tied-transformer [Kim 等人, 2021]	✓		列图 序列	NR/47.1	NR/73.1
		RetroLSTM [Liu 等人, 2017]	✓		图 序列图	NR/37.4	NR/57.0
		G2GT [Lin 等人, 2022]	✓		序列图 序	NR/54.1	NR/74.5
		Stranformer [Lee 等人, 2019]	×		列图 序列	NR/43.8	NR/NR
		GraphRetro [Somnath 等人, 2021]	✓		图 序列图	63.9/53.7	85.2/72.2
		MARS [Liu 等人, 2022]			序列图 序	66.2/54.6	90.2/83.3
					列图 序列		
					图 序列图		
					序列图 序		
					列图 序列		
					图 序列图		
					序列图 序		
					列图 序列		
					图 序列图		
					序列图 序		
			列图 序列				
			图 序列图				
			序列图 序				
			列图 序列				
			图 序列图				
			序列图 序				
			列图 序列				
			图 序列图				
			序列图 序				
			列图 序列				
			图 序列图				
			序列图 序				
			列图 序列				
			图 序列图				
			序列图 序				
			列图 序列				
			图 序列图				
			序列图 序				
			列图 序列				
			图 序列图				
			序列图 序				
			列图 序列				
			图 序列图				
			序列图 序				
			列图 序列				
			图 序列图				
			序列图 序				
			列图 序列				
			图 序列图				
			序列图 序				
			列图 序列				
			图 序列图				
			序列图 序				
			列图 序列				
			图 序列图				
			序列图 序				
			列图 序列				
			图 序列图				
			序列图 序				
			列图 序列				
			图 序列图				
			序列图 序				
			列图 序列				
			图 序列图				
			序列图 序				
			列图 序列				
			图 序列图				
			序列图 序				
			列图 序列				
			图 序列图				
			序列图 序				
			列图 序列				
			图 序列图				
			序列图 序				
			列图 序列				
			图 序列图				
			序列图 序				
			列图 序列				
			图 序列图				
			序列图 序				
			列图 序列				
			图 序列图				
			序列图 序				
			列图 序列				
			图 序列图				
			序列图 序				
			列图 序列				
			图 序列图				
			序列图 序				
			列图 序列				
			图 序列图				
			序列图 序				
			列图 序列				
			图 序列图				
			序列图 序				
			列图 序列				
			图 序列图				
			序列图 序				
			列图 序列				
			图 序列图				
			序列图 序				
			列图 序列				
			图 序列图				
			序列图 序				
			列图 序列				
			图 序列图				
			序列图 序				
			列图 序列				
			图 序列图				
			序列图 序				
			列图 序列				
			图 序列图				
			序列图 序				
			列图 序列				
			图 序列图				
			序列图 序				
			列图 序列				
			图 序列图				
			序列图 序				
			列图 序列				
			图 序列图				
			序列图 序				
			列图 序列				
			图 序列图				
</							

这将降低预测产物的有效性。总之，不确定性估计是非自回归模型面临的一个重要挑战。正确估计不确定性的非自回归模型将为反应预测研究带来革命性的变化。

逆合成。非自回归逆合成模型非常棘手，因为在连接离去基团之前，必须预先确定反应物的原子总数。然而，反应物原子总数的分布受到基准数据集的影响。因此，为了以非自回归方式生成反应物，应为预测的离去基团保留足够的空白条目。不幸的是，如果三维加成矩阵和原子特征矩阵中的空白条目过多，这项任务就会变成一项非常艰巨的不平衡分类任务，因为只有少数空白条目会被预测的原子和化学键填满。

4.5 推广方面的挑战

在对反应和逆合成预测进行归纳时，有两个共同的挑战。**第一个挑战是分布外预测。**分布外预测主要评估在主要反应类型上训练的模型能否泛化到稀有反应类型。这种不平衡问题自然存在于化学反应中，因为我们可以预料到有些反应很容易发生，而有些反应却很少发生。因此，设计不同分布情况下稳定的反应和逆合成预测模型非常重要。**第二个挑战来自低层次的再现质量。**这两个问题都需要低层次的分子和反应表征学习技术。因此，推导出强大的分子和反应表征并将其推广到相关任务中具有重要意义。

5 结论和未来方向

从上述讨论中，我们知道当前的方法虽然取得了可喜的成果，但仍存在一些重大缺陷。因此，我们列出了未来进一步完善现有解决方案的几个方向。

5.1 三维分子信息

化学分子本质上是由一组具有笛卡尔坐标的三维点组成的点云。然而，三维分子信息尚未在文献中得到应用，而分子的序列和图形表示法却已得到广泛探索。三维位置矢量为每对原子提供了重要的互补距离信息。与二维分子图相比，三维欧几里得几何中的相对成对距离可能会有很大不同。例如，原子 v_1 和原子 v_2 在非欧几里得分子图中可能相距甚远，而在三维欧几里得空间中可能相距很近。这对反应中心排序特别有用。因此，有效地将三维分子信息纳入建模有助于更准确地预测反应和逆合成。

5.2 多样化的基准数据集和新的评估指标

对于反应预测而言，当前 USPTO-479K 数据集的规模已经足够大。但是，其模式和多样性

仍然不够。新的基准数据集应包括更多的反应类型和更复杂的反应。此外，还应包括不同的数据集拆分，如支架拆分和时间拆分，以进行交叉验证。对于逆合成预测，USPTO-50K 数据集的规模较小。新的基准数据集应该是至少包含 10 万个样本的大规模数据集。此外，未来的基准数据集应包含更多具有多重编辑的目标分子。此外，新的逆合成评估指标也是必要和迫切的。FusionRetro [Liu 等人, 2023] 尝试在多步骤规划的背景下评估单步骤逆合成模型。未来还可以设计更多样化的评估指标。

5.3 非自回归模型

正如上一节所述，非自回归模型的不确定性估计仍然不准确。因此，为非自回归模型探索更有效的不确定性建模方案至关重要。我们相信，如果非自回归模型的 top-k 精度能与自回归模型相媲美，那将是一个开创性的时刻。非自回归的逆合成预测在目前的研究中还没有涉及。由于数据集的离群规模较小，推理速度的重要性在目前的研究中被忽视。我们希望在分析复杂的合成路线时，推理效率能得到重视。

5.4 自我监督学习

正如之前讨论的局限性一样，反应和回溯论文预测受到数据集本身带来的一些问题的困扰。在无标注数据集上利用自监督学习 (SSL) 策略来克服这些局限是一个自然而然的探索方向。虽然针对一般分子表征学习已经提出了多种 SSL 策略，但目前很少有针对反应和逆合成预编译设计的特定 SSL 策略。PMSR [Jiang 等人, 2021 年] 探索了逆合成的 SSL 策略，但未能取得很好的效果。因此，针对这两个问题探索更强大的 SSL 策略是一个可行的方向。

5.5 多任务学习

将这两项任务与其他相关任务相结合是一个很有前景的方向。DualTF [Sun 等人, 2021] 首次尝试提出了一个双模型，表明反应和回溯预测模型可以相互促进。FusionRetro [Liu 等, 2023] 和 GNN-Retro [Han 等, 2022] 分别利用反应上下文和分子上下文来改进逆合成分析。Lu 和 Zhang [Lu and Zhang, 2022] 的研究表明，多任务学习可以进一步提高反应预测的准确性。直观地说，反应分类、反应产率预测和许多其他任务都与反应和逆合成预测密切相关。如何综合利用这些任务对于进一步改进这两个问题至关重要。

致谢

本文所述工作得到了国家重点研发计划（编号：2018AAA0100204）和中国香港特别行政区研究资助局（中大14222922，研资局GRF，编号：2151185）的部分资助。

参考资料

- [Bi *et al.*, 2021] 毕杭瑞、王恒义、史晨思、Connor W. Coley、唐健和郭宏宇。用于反应预测的非自回归电子再分布建模。In *ICML 2021*, volume 139 of *Proceedings of Machine Learning Research*, pages 904-913. PMLR, 2021.
- [Bradshaw *et al.*, 2019] John Bradshaw, Matt J. Kusner, 布鲁克斯-佩奇、马文-H-S-塞格勒和何塞-米格尔 Hernánde z-Lobato. 电子路径的生成模型。In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net, 2019.
- [Chen and Jung, 2021] Shuan Chen and Yousung Jung. 利用局部反应和全局注意力的深度逆合成反应预测。 *JACS Au*, 1(10):1612-1620, 2021.
- [Chen and Jung, 2022] Shuan Chen and Yousung Jung. 基于通用模板的图神经网络用于有机反应性预测。 *Nature Machine Intelligence*, 4:1-9, 09 2022.
- [陈等人, 2019] Benson Chen、Tianxiao Shen、Tommi S. Jaakkola 和 Regina Barzilay。Jaakkola 和 Regina Barzilay. 学习为逆合成做出可生成的多样化预测。 *CoRR*, abs/1910.09688, 2019.
- [Coley 等人, 2017] Connor W. Coley、Luke Rogers、William H. Green 和 Klavs F. Jensen。基于分子相似性的计算机辅助逆合成。 *ACS Central Science*, 3(12):1237-1245, 2017.
- [戴等人, 2020] 戴汉军、李成涛、Connor W. Coley、戴波和宋乐。用条件图逻辑网络进行逆合成预测。 *CoRR*, abs/2001.01408, 2020.
- [Do 等人, 2019] Kien Do、Truyen Tran 和 Svetha Venkatesh。化学反应预测的图转换策略网络。In *KDD 2019*, pages 750-760. ACM, 2019.
- [Dong 等人, 2021 年] Jingxin Dong、Mingyi Zhao、Yuansheng Liu、Yansen Su 和 Xiangxiang Zeng。深度学习在逆合成规划中的应用：数据集、模型和工具。 *Briefings in Bioinformatics*, 23(1), 2021.
- [Engkvist *et al.*, 2018] Ola Engkvist、Per-Ola Norrby、Nidhal Selmi、Yu hong Lam、Zhengwei Peng、Edward C. Sherer、Willi Amberg、Thomas Erhard 和 Lynette A. Smyth。化学反应的计算预测：现状与展望。 *今日药物发现*, 23(6):1203-1218, 2018.
- [Han *et al.*, 2022] Peng Han, Peilin Zhao, Chan Lu, Junzhou Huang, Jiaxiang Wu, Shuo Shang, Bin Yao, and Xiangliang Zhang. Gnn-retro: 利用图神经网络的逆合成规划。 *美国人工智能学会会议论文集*, 36(4), 2022.
- [Irwin 等人, 2022] 罗斯-欧文、斯皮里宗-迪米特里阿迪斯、何继振和埃斯本-雅尼克-比耶鲁姆。Chemformer: a pre-trained transformer for computational chemistry. *机器学习：科学与技术*, 3(1):015022, 2022 年 1 月。
- [Jiang 等人, 2021 年] Yinjie Jiang、Ying Wei、Fei Wu、Zhengxiang Huang、Kun Kuang 和 Zhihua Wang。通过预训练学习逆合成的化学规则。在 *AAAI 2021* 中, 第 531-539 页。AAAI Press, 2021.
- [Jin *et al.*, 2017] Wengong Jin, Connor W. Coley, Regina Barzilay, and Tommi S. Jaakkola. Jaakkola. 预测有机再魏斯费勒-莱曼网络的行动结果。在 *NIPS 2017* 年, 第 2607-2616 页, 2017 年。
- [Eunji Kim、Dongseon Lee、Youngchun Kwon、Min Sik Park 和 Youn-Suk Choi。使用具有潜在变量的绑定双向变换器进行有效、可信和多样化的逆合成。 *化学信息与建模期刊*, 61(1):123-133, 2021.
- [David Kovacs、William McCorkindale 和 Alpha Lee。定量解释用于化学反应预测的机器学习模型并揭示偏差。 *自然通讯*, 12, 03 2021.
- [Lee *et al.*, 2019] Alpha Albert Lee、Qingyi Yang、Vishnu Sresht、Peter Bolgar、Xinjun Hou、Jacquelyn L. Klug-McLeod 和 Christopher R. Butler。分子转换器统一了制药化学领域的反应预测和逆合成。 *化学通讯*, 2019 年。
- [林等, 2020] 林康杰、徐友军、裴剑锋、赖璐华。使用无模板模型的自动逆合成路线规划。 *化学科学*, 11, 03 2020.
- [Lin *et al.*, 2022] Zaiyun Lin, Shiqiu Yin, Lei Shi, Wenbiao Zhou, and Yingsheng J. Zhang. G2gt: 使用图到图注意力神经网络和自我训练的逆合成预编译。 *ArXiv*, abs/2204.08608, 2022.
- [Liu *et al.*, 2017] Bowen Liu, Bharath Ramsundar, Prasad Kawthekar, Jade Shi, Joseph Gomes, Quang Luu Nguyen, Stephen Ho, Jack Sloane, Paul Wender, and Vijay Pande。利用神经序列模型预测逆合成反应。 *ACS Central Science*, 3(10):1103-1113, 2017.
- [刘等人, 2022 年] 刘嘉涵、闫超超、于洋、吕婵、黄俊洲、杨乐欧、赵培林。MARS: 基于图案的回溯论文预测自回归模型。 *CoRR*, abs/2209.13178, 2022.
- [刘松涛、涂正凯、徐敏凯、张作柏、林璐、应睿、唐健、赵培林、吴定浩。Fusionretro: 通过上下文反应进行分子表征融合以实现逆合成规划。In *ICML 2023, Proceedings of Machine Learning Research*. PMLR, 2023.

- [Lu and Zhang, 2022] Jieyu Lu and Yingkai Zhang.多任务反应预测的联合深度学习模型与解释。 *J. Chem. Inf. Model.*, 62(6):1376-1387, 2022.
- [Mao 等人, 2021] 毛克龙、肖茜、徐廷扬、荣瑜、黄俊洲和赵培林。用于逆合成预测的分子图增强变换器。 *神经计算*, 457: 193-202, 2021。
- [孟子乔等, 2023] 孟子乔、赵培林、于洋和欧文-金。基于双随机图的非自回归反应预测。 In *IJCAI 2023*. *ijcai.org*, 2023.
- [钱伟等, 2020] Wesley Wei Qian、Nathan T. Russell、Claire L. W. Simons、Yunan Luo、Martin D. Burke 和 Jian Peng。集成深度神经网络和符号推理进行有机反应性预测。 *Chem- Rxiv*, 2020.
- [Mikolaj Sacha、Mikolaj Blaz、Piotr Byrski、Pawel Wlodarczyk-Pruszyński 和 Stanislaw Jas-trzebski。分子编辑图注意网络：将化学反应建模为图编辑序列。 *CoRR*, abs/2006.15426, 2020.
- [Scannell 等人, 2012 年] Jack Scannell、Alex Blanckley、Helen Boldon 和 Brian Warrington。诊断制药研发效率的下降。 *自然评论. 药物发现*, 11:191-200, 2012 年 3 月。
- [Philippe Schwaller, Teodoro Laino, Theophile Gaudin, Peter Bolgar, Christopher A. Hunter, Costas Bekas, and Alpha A. Lee.分子转换器：不确定性校准化学反应预处理模型。 *ACS Central Science*, 5(9), 2019.
- [Segler and Waller, 2017] Marwin H. S. Segler and Mark P. Waller.用于回溯和反应预测的神经符号机器学习。 *Chemistry - A European Journal*, 23(25):5966-5971, 2017.
- [Seung-Woo Seo、You Young Song、June Yong Yang、Seohui Bae、Hankook Lee、Jinwoo Shin、Sung Ju Hwang 和 Eunho Yang。GTA：用于逆合成的图截断注意力。 In *AAAI 2021*, pages 531-539. *AAAI Press*, 2021.
- [Shi 等人, 2020] Chence Shi、Minkai Xu、Hongyu Guo、Ming Zhang 和 Jian Tang。用于逆合成预测的图到图框架。 In *ICML 2020*, volume 119 of *Proceedings of Machine Learning Research*, pages 8818-8827. *PMLR*, 2020.
- [Sohn 等人, 2015] Kihyuk Sohn、Honglak Lee 和 Xinchun Yan。使用深度条件生成模型学习结构化输出表示。 In *Advances in Neural Information Processing Systems 28: Annual Conference on Neural Information Processing Systems 2015, December 7-12, 2015, Montreal, Quebec, Canada*, pages 3483-3491, 2015.
- [Somnath 等人, 2021 年] Vignesh Ram Somnath、Charlotte Bunne、Connor W. Coley、Andreas Krause 和 Regina Barzilay。学习逆合成预词典的图模型。 In *NeurIPS 2021*, 2021.
- [Sun et al., 2021] Ruoxi Sun, Hanjun Dai, Li Li, Steven Kearnes, and Bo Dai.通过基于能量的模型理解回溯论文。 In *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, pages 10186-10194, 2021.
- [Tetko 等人, 2020] Igor Tetko、Pavel Karpov、Ruud Deursen 和 Guillaume Godin。用于直接和单步逆合成的最新增强型 NLP 变换器模型。 *自然通讯*, 11, 11 2020。
- [Tu and Coley, 2021] Zhengkai Tu and Connor W. Coley。用于无模板逆合成和反应预测的置换不变图-序列模型。 *CoRR*, abs/2110.09681, 2021.
- [Vaswani et al., 2017] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin.注意力就是你所需要的一切。 In *NIPS 2017*, pages 5998-6008, 2017.
- [Wan 等人, 2022 年] Yue Wan、Chang-Yu Hsieh、Ben Liao 和 Shengyu Zhang.逆合成变压器：挑战端到端逆合成变换器的极限。 In *International Conference on Machine Learning, ICML 2022, 17-23 July 2022, Baltimore, Maryland, USA*, volume 162 of *Proceedings of Machine Learning Research*, pages 22475-22490. *PMLR*, 2022.
- [王晓瑞、李玉泉、邱杰忠、陈光勇、刘焕祥、廖本本、谢长瑜、姚小军。Retroprime：基于转换器的单步合成预测方法。 *化学工程学报*, 420:129845, 2021。
- [Wei 等人, 2016] Jennifer N. Wei、David Duvenaud 和 Alan Aspuru-Guzik。用于预测有机化学反应的神经网络。 *ACS Central Science*, 2(10):725-732, 2016.
- [Yan et al., 2020] Chaochao Yan, Qianggang Ding, Peilin Zhao, Shuangjia Zheng, Jinyu Yang, Yang Yu, and Junzhou Huang. Retroxpert：像化学家一样分解逆合成前词典。 In *NeurIPS 2020*, 2020.
- [Yan 等人, 2022 年] Chaochao Yan、Peilin Zhao、Chan Lu、Yang Yu 和 Junzhou Huang。RetroComposer：基于模板的逆合成预测模板。 *生物分子*, 12 (9) : 1325, 2022 年 7 月。
- [赵明等人, 2022 年] 赵明、方磊、谭立、Yves LePage 和楼建光。利用反应感知子结构进行逆合成和反应预测。 *CoRR*, abs/2204.05919, 2022.
- [Zheng 等人, 2020] Shuangjia Zheng、Jiahua Rao、Zhongyue Zhang、Jun Xu 和 Yuedong Yang。利用自校正变压器神经网络预判逆合成反应。 *化学信息与建模学报*, 60 (1) : 47-55, 2020.