



Review

Retrosynthesis prediction using grammar-based neural machine translation: An information-theoretic approach

Vipul Mann, Venkat Venkatasubramanian*

Department of Chemical Engineering, Columbia University, New York, NY, 10027, USA

ARTICLE INFO

Article history:

Received 13 April 2021

Revised 19 August 2021

Accepted 7 September 2021

Available online 11 September 2021

Keywords:

Machine learning

Retrosynthetic analysis

Artificial intelligence

Synthesis Planning

Reaction prediction

ABSTRACT

Retrosynthetic prediction is one of the main challenges in chemical synthesis because it requires a search over the space of plausible chemical reactions that often results in complex, multi-step, branched synthesis trees for even moderately complex organic reactions. Here, we propose an approach that performs single-step retrosynthesis prediction using SMILES grammar-based representations in a neural machine translation framework. Information-theoretic analyses of such grammar-representations reveal that they are superior to SMILES representations and are better-suited for machine learning tasks due to their underlying redundancy and high information capacity. We report the top-1 prediction accuracy of 43.8% (syntactic validity 95.6%) and maximal fragment (MaxFrag) accuracy of 50.4%. Comparing our model's performance with previous work that used character-based SMILES representations demonstrate significant reduction in grammatically invalid predictions and improved prediction accuracy. Fewer invalid predictions for both known and unknown reaction class scenarios demonstrate the model's ability to learn the underlying SMILES grammar efficiently.

© 2021 Elsevier Ltd. All rights reserved.

1. Introduction

One of the important challenges in computational chemistry is the retrosynthetic analysis of desired molecules that satisfy property constraints, subject to the commercial availability of the precursors and the feasibility of the chemical reactions required for their synthesis. The immense interest in this problem over the recent years could be attributed to its practical applications across areas such as drug discovery, synthesis of novel organic compounds, and improvements in the reactions pathways from a commercial, social, or economic viability standpoint. The industrial applications of retrosynthetic analysis include automobiles, petrochemicals, specialty chemicals, and polymer science, with a great potential to revolutionize the entire industry if the right compound could be synthesized.

Retrosynthetic analysis often involves evaluating many potential candidate reaction pathways and molecules at multiple stages of the reaction, resulting in complex retrosynthesis trees that need to be searched and parsed efficiently. Computational approaches could significantly aid the chemist in solving different aspects of the retrosynthesis problem, such as the graph-theoretic search methodologies for efficient tree traversal to identify fea-

sible reaction pathways, dictionary-based methods to evaluate a large search-space of precursors, and chemistry-driven heuristics to eliminate practically infeasible routes. Multi-step retrosynthesis is usually formulated as a fundamentally different problem compared to the single-step retrosynthesis that we study in our work, and often involves either using efficient search techniques combined with one-step forward synthesis models, or using a sequence of single-step reverse transformations informed by chemistry-based heuristics.

One of the first attempts that leveraged computational tools and formalized the retrosynthesis problem was LHASA proposed in Pensak and Corey. This framework used logic and chemistry rules in the form of heuristics and transformations along with a chemical programming language to solve the retrosynthesis problem. Several subsequent approaches were proposed that utilized rule-based expert-systems (Salatin and Jorgensen, 1980; Corey et al., 1985; Jorgensen et al., 1990; Satoh and Funatsu, 1995; 1999; Chen and Baldi, 2009; Law et al., 2009; Gothard et al., 2012; Szymkuć et al., 2016; Segler and Waller, 2017a), with a few of them combining network theory to discover that chemistry networks follow scale-free properties, a ubiquitous class of networks reported to be optimal in several other areas (Mann et al., 2021; Biłozor et al., 2018; Zhang et al., 2020; Xu et al., 2020). However, such approaches were hard to scale beyond interesting prototypes as they required great human effort and expertise to develop (Venkatasubramanian, 2019).

* Corresponding author.

E-mail addresses: vm2583@columbia.edu (V. Mann), venkat@columbia.edu (V. Venkatasubramanian).

However, in recent years, the massive surge in computational capabilities combined with significant advances in machine learning have resulted in a renewed attack on this problem. This includes approaches that combine neural network models with known chemistry knowledge encoded in the form of reaction templates – e.g., Segler and Waller (2017b) leveraged neural networks for selecting the reactivity centers and most suitable transformations; Wei et al. (2016) predicted reaction types and used Smiles Molecular Arbitrary Target Specification (SMARTS) templates for predicting the likely products given a set of reactants and reagents, and (Coley et al., 2017) proposed selecting the suitable edit-based transformations in a reaction using reaction templates. Such methods, however, again address only certain limitations of the rules-based systems and the inherent limitation of the lack of their ability to suggest novel chemical reactions and a bias towards the common reaction types still exist.

This is overcome in purely data-driven approaches that use sophisticated machine learning architectures to learn the complex non-linear dynamics of a chemical reaction – both in the forward and the backward directions – primarily by modeling the chemical representations. This includes the neural sequence-to-sequence (or seq2seq) models introduced for the forward reaction prediction in Nam and Kim (2016) and the retrosynthetic prediction in Liu et al. (2017) that formulated the reaction prediction task as a sequence modeling problem. Other recent efforts for the retrosynthesis task include a seq2seq approach combined with a Monte Carlo tree search (Lin et al., 2020) and various transformer model-based approaches (Zheng et al., 2019; Mann and Venkatasubramanian, 2021; 2020; Karpov et al., 2019; Duan et al., 2020; Schwaller et al., 2020; Tetko et al., 2020).

Even though the prediction accuracy has significantly improved due to the increased complexity of model architectures, prior chemistry knowledge in such frameworks is still missing. The incorporation of this knowledge should, in principle, improve the model performance on out-of-sample examples. All previous works in this area use SMILES representations of molecules, treating them as merely character-based strings, except for the recent work by Ucak et al. (2021) that used substructure-based representations but suffered from lower prediction accuracy. In our earlier work on the forward prediction problem (Mann and Venkatasubramanian, 2021; 2020), we demonstrated that incorporating chemical and structural information about molecules ensures that the model learns the underlying chemical transformations with significantly fewer training parameters. As an extension of that work, we propose here a framework for solving the retrosynthesis problem using the rich, SMILES grammar-based representation of molecules and highlight the inherent benefits of such representations – both from an information-theoretic and model performance standpoint.

The rest of the paper is organized as follows: In Section 2, we formally define the retrosynthesis prediction problem as a sequence modeling task in the machine translation framework and present an overview of the methods underlying our work, such as the SMILES grammar, the transformer architecture and the beam search decoding procedure in Section 3. In Section 4, we present an information-theoretic analysis of the proposed grammar-representations and contrast them with the other representations (SMILES and molecular formula) to highlight the differences and quantify the advantages of using the underlying chemical structural information. The standard reaction dataset and the model training aspects of our work are presented in Section 5. The evaluation metrics used for assessing our model's performance, the results on the USPTO 50K reactions dataset, comparison with other works, and the limitations and future work in this direction are presented in Section 6. Finally, the concluding remarks summarizing the major contributions of this work appear in Section 7.

2. Problem formulation and objectives

We formulate the retrosynthesis prediction problem as a sequence modeling task and use a machine translation framework for predicting the precursors for a given target molecule. The objective is to translate a set of input tokens corresponding to the product molecule to an output sequence of tokens corresponding to the precursor molecules. The input sequence may be optionally prepended with an identifier that indicates the reaction class. To allow the model to differentiate between the different precursors (reactants), a separate identifier token is used to indicate the end of the representation of a given precursor and the start of another. This framework is depicted in Fig. 1.

In this framework, the participating product and reactants in a given reaction are represented using their corresponding grammar-based representation described in detail in Section 3.1. The representation starts with the token '1' and ends with the token '80' for all the molecules, the token '81' separates multiple reactants, and the token '82' signifies the end of all the precursor representations. The other identifiers (or tokens) correspond to the sequence of production rules required to obtain the given SMILES string, using the grammar productions described in Table A.12 in the Appendix. The sequence modeling task is performed using a transformer model, a state-of-the-art architecture for sequence modeling (Vaswani et al., 2017).

3. Methods

In this section, we describe the methods involving our approach, namely the SMILES grammar-based representations used for encoding molecules, the transformer architecture used for the sequence modeling task, and the beam search decoding procedure used for generating a set of most likely target sequences for a given input sequence.

3.1. SMILES grammar

One of the first works that attempted to formalize natural language through context-free grammars (CFGs) was proposed by Chomsky (1956) that was based on the idea that a group of words could be thought of as belonging to a constituent unit and that different constituent units could be grouped, hierarchically, to convey a given meaning. Formally, a context-free grammar could be thought of as a set of production rules that define the transformation of a set of non-terminal symbols to terminal symbols that correspond to strings with meaning in the natural language. In addition, there is a designated start symbol that indicates the start of a sentence. Therefore, a CFG consists of the following elements: S, a designated start symbol; Σ , the set of terminal symbols; N, the set of non-terminal symbols; and R, the set of production rules of the form $A \rightarrow \beta$ where $A \in N$ is non-terminal and $\beta \in \Sigma$ is a terminal symbol.

A similar grammar for the SMILES representation of molecules also exists (Weininger, 1988) where the individual tokens in the SMILES string represent the terminal symbols that could be obtained through the sequential application of a set of production rules on the non-terminal symbols. Consider, for example, a subset of the official SMILES grammar presented in Table 1. The equivalent symbols similar to CFG for this grammar are:

- S: SMILES
- Σ : { (,), =, c, C, 0, 1, 2 }
- N: { SMILES, CHAIN, BRANCHED_ATOM, BOND, ATOM, RINGBOND, BB, RB, BRANCH, AROMATIC_ORGANIC, ALIPHATIC_ORGANIC, DIGIT }
- R: productions (rules) 1 through 20 in Table 1

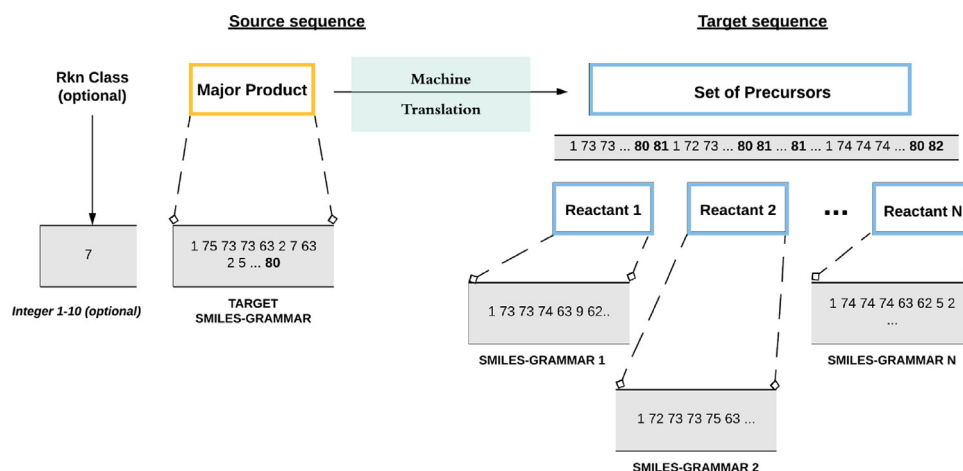


Fig. 1. The single-step retrosynthesis prediction problem formulation using machine translation. The reaction class information is optional.

Table 1
Reduced SMILES grammar.

S.No	Production rules
1	SMILES \rightarrow CHAIN
2	CHAIN \rightarrow CHAIN BRANCHED_ATOM
3	CHAIN \rightarrow CHAIN BOND BRANCHED_ATOM
4	CHAIN \rightarrow BRANCHED_ATOM
5	BRANCHED_ATOM \rightarrow ATOM RINGBOND
6	BRANCHED_ATOM \rightarrow ATOM
7	BRANCHED_ATOM \rightarrow ATOM BB
8	BRANCHED_ATOM \rightarrow ATOM RB
9	BB \rightarrow BRANCH
10	RB \rightarrow RINGBOND
11	BRANCH \rightarrow (CHAIN)
12	RINGBOND \rightarrow DIGIT
13	BOND \rightarrow =
14	ATOM \rightarrow AROMATIC_ORGANIC
15	ATOM \rightarrow ALIPHATIC_ORGANIC
16	AROMATIC_ORGANIC \rightarrow c
17	ALIPHATIC_ORGANIC \rightarrow C
18	ALIPHATIC_ORGANIC \rightarrow O
19	DIGIT \rightarrow 1
20	DIGIT \rightarrow 2

We leverage such underlying grammar to assign structure to a given SMILES string and derive from such structures the grammar-based representations. Consider benzene, with the SMILES string representation given by C1=CC=CC=C1. This representation could be obtained by applying the set of production rules in Table 1 sequentially with the corresponding parse-tree shown in Fig. 2. The grammar-representation that we work with, originally proposed in our earlier work (Mann and Venkatasubramanian, 2021), is obtained by extracting production rules from the parse-tree by parsing it in a bottom-up-left-corner strategy, i.e., starting at the top and going down the left-most branch, then coming back up to parse the immediate right branch, and so on until the entire tree is parsed. The grammar representation thus obtained corresponding to the parse-tree for benzene is given in the figure caption.

Clearly, as compared to a purely character-based SMILES string representation consisting merely of the tokens 'C', '1', '=', 'C', 'C', '=', 'C', 'C', '=', 'C', '1', without any additional information conveying the relationships between the tokens, the grammar-based representations are significantly richer, incorporate chemical and structural information, and contain hierarchical information about the underlying chemistry. This is leveraged by the model architecture for modeling the underlying SMILES grammar. We have shown that these representations are more efficient in modeling the underlying chemistry and eliminate overparameterization in

complex machine learning architectures (Mann and Venkatasubramanian, 2021). We present an information-theoretic analysis of the grammar representations and the text-based representations in Section 4 to establish the fundamental superiority of the grammar representations compared to other text-based representations such as SMILES.

3.2. Sequence-to-sequence models

We model the reaction prediction problem as a sequence modeling task that involves mapping the input sequence to a sequence of tokens corresponding to the output sequence. This framework has been used in recent years and has shown a significant promise in reaction modeling. We use the state-of-the-art model in this area, known as the transformer framework, proposed in Vaswani et al. (2017).

The transformer framework, shown in Fig. 3, consists of an encoder-decoder architecture where the encoder maps the input sequence to a latent space, and the decoder decodes from the latent space in an autoregressive manner, one element at a time, to give rise to the output sequence. The positional encodings in a transformer encode the position of a given word (or token) in the sequence to a high dimensional vector space, getting rid of recurrent or convolution operations that significantly improved the computational complexity of training the model architecture. These mappings are characterized by sines and cosines of different frequencies, given by

$$\vec{p}_{pos,i} = \begin{cases} \sin(pos/10000^{2k/d}), & \text{if } i = 2k \\ \cos(pos/10000^{2k/d}), & \text{if } i = 2k + 1 \end{cases} \quad (1)$$

An attention mechanism lets the transformer model relationships between groups of words in an input sequence at different stages of the network. The attention-mechanism used in Vaswani et al. (2017) is the 'Scaled-Dot Product Attention', characterized by a set of queries, keys, and values vectors. The query and key vectors are of dimensions d_k , and the value vector is of dimension d_v which are used to represent a given word and the corresponding key-value pairs for computing the attention function. The query, key, and value vectors are obtained from the output of the preceding layers in the network. The attention-score is computed as softmax function applied over the dot-products of the queries and key vectors, scaled down by a factor of $\sqrt{d_k}$, given by

$$Attention(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V \quad (2)$$

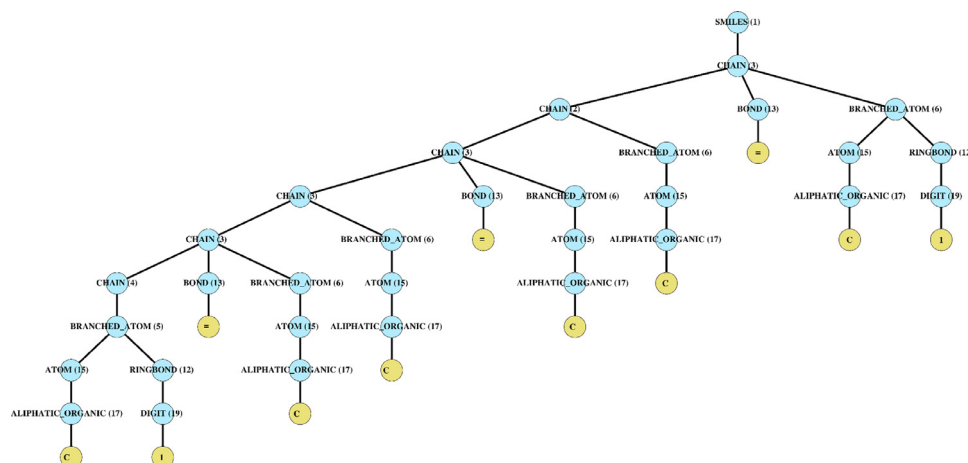


Fig. 2. The parse-tree obtained for benzene with SMILES string representation as C1=CC=CC=C1. The production rules from Table 1 applied at each stage are indicated next to the non-terminal symbols. Parsing this tree in a bottom-up-left-corner strategy gives rise to the grammar-representation given by: 1,3,2,3,3,3,4,5,15,17,12,19,13,6,15,17,6,15,17,13,6,15,17,6,15,17,13,6,15,17,13,6,15,17,12,19.

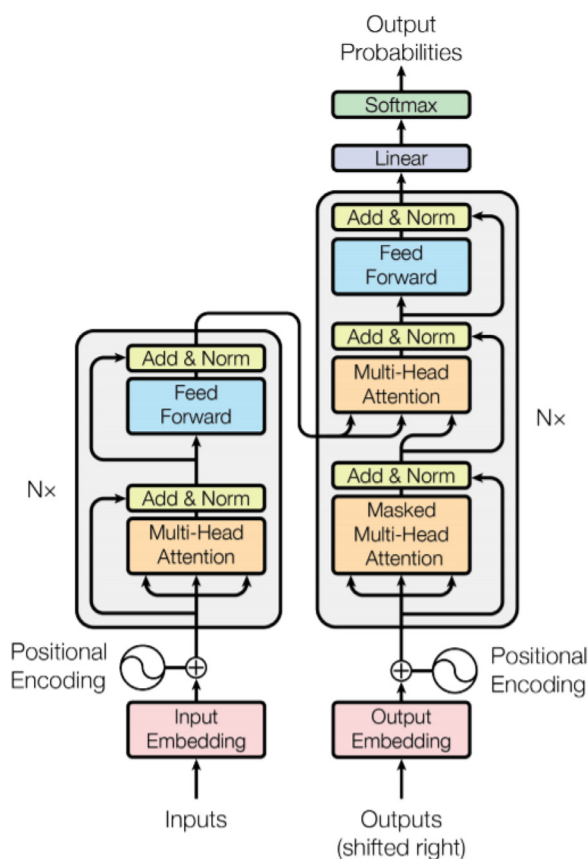


Fig. 3. The encoder-decoder model architecture of a transformer as depicted in Vaswani et al. (2017).

where Q , K , and V are the matrices of query, key, and values vectors, respectively. The attention score computed above determines the importance of different parts of an input sequence in the current context. In order to allow the model to jointly factor in information from different representation subspaces at different positions, multi-headed attention is computed, which involves computing multiple attention scores in parallel, which are then concatenated and projected using a linear transformation to compute

the multi-head attention scores as,

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \text{head}_2, \dots, \text{head}_h)W^O \quad (3)$$

where $\text{head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$, and $W_i^Q \in \mathbb{R}^{d_{\text{pos}} \times d_k}$, $W_i^K \in \mathbb{R}^{d_{\text{pos}} \times d_k}$, and $W_i^V \in \mathbb{R}^{d_{\text{pos}} \times d_v}$ are the projection matrices for Q , K , and V , respectively. The reader is referred to Vaswani et al. (2017) for further details on the transformer model architecture.

3.3. Beam search

In order to generate the output sequences in a transformer framework, a decoding procedure is used that decodes from the latent space in an autoregressive manner, with the current prediction as input while decoding the next token. Therefore, the decoding procedure could either use a greedy strategy that involves selecting the token with the maximum likelihood at each stage for decoding the next token, generating a single most-likely sequence in the end; or on the other hand, it could employ a beam search procedure that decodes a set of top- B tokens at each stage based on their likelihood and return them as the model output. We follow the latter approach for decoding. This allows us to evaluate our model's performance more extensively and compare it with the top- K accuracy reported in other similar works in this area. A schematic of the beam-search decoding procedure used in our work is shown in Fig. 4.

4. Information-theoretic analysis of chemical representations

Before discussing the model training aspects, we demonstrate the richness of the proposed grammar-based representations using an information-theoretic framework. We compare the information capacity, information gain, and redundancy characterizing the various symbols-based chemical representations, namely, molecular formula, SMILES, and grammar representations. We first provide a brief overview of the relevant information-theoretic concepts and their intuition in the next section, followed by their application to chemical representations and quantify the superiority of grammar-based representations from an information-theoretic standpoint.

4.1. Shannon entropy and information content

The development and formalization of information theory, mainly by Claude Shannon in Shannon (1948), offered a mathe-

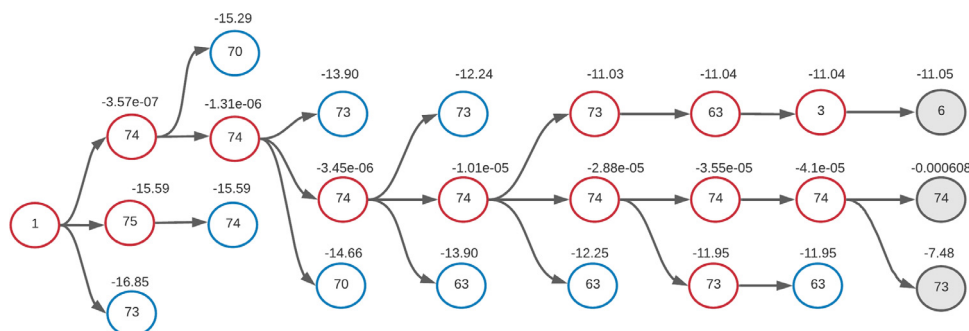


Fig. 4. A partially completed beam search output for a beam width of 3 for a reference input. At each stage, the most likely grammar-rules are predicted, that are used as input to decode the next most likely set of grammar-rules, and so on until the entire sequence of grammar-rules corresponding to a given SMILES string is reconstructed. The log-likelihood values are indicated above each node in the schematic.

mathematical definition of the *amount of information* communicated between any two components or channels of a given system. The primary motivation was the fundamental problem of decoding a source message passing through a noisy channel, either exactly or approximately, at any other point in the communication system. However, the applications and adaptations of it are not limited to communication systems alone but have had far-reaching consequences across most fields of science and engineering.

The Shannon entropy for a given probability distribution $p(x)$ of a random variable x is defined as,

$$H = - \sum_{i=1}^M p(x_i) \log_2 p(x_i) \quad (4)$$

where $p(x)$ is the probability mass function of x with M possible values. This is equivalent to the expected value of the Shannon information or self-information of a variable and is measured in units of *bits per symbol*. There is a direct correspondence between the amount of information in a message and the degree of uncertainty that is associated with it. That is, if a system can exist in one of a very large number of possible states, then there is a great amount of uncertainty associated with its state as opposed to another system that can exist only in a handful states. Therefore, the amount of information required is more for the former than the latter.

Consider the two extremes of zero-information content and maximum information content. The Shannon entropy in Eq. 4 attains a value of zero when the probability $p(x_i)$ of a x_i attaining a given value is 1 meaning that the outcome or the value that x_i could take is known with complete certainty, and hence, there is no information content (or gain) associated with knowing its value explicitly. On the other hand, when x_i could take any of the possible values with equal probability, i.e., $p(x_i) = 1/M$ where M is the total number of possible values that the symbols in the source message could take, the information content is maximized and is equal to $\log_2 M$. This implies that in such a scenario, specifying the value of a given bit in the sequence would result in the maximum information gain when compared to any other scenario.

The generalization of Eq. 4 when several random variables X_1, X_2, \dots, X_n are present is given by the joint Shannon entropy as,

$$H(X_1, X_2, \dots, X_n) = - \sum_{x_1, x_2, \dots, x_n} p(x_1, x_2, \dots, x_n) \log_2 p(x_1, x_2, \dots, x_n) \quad (5)$$

The joint entropy in Eq. 5 could be interpreted as an information measure corresponding to multiple random variables presented simultaneously. Similarly, the conditional entropy that quantifies the information content of a given random variable X_1 conditioned on

a set of other random variable X_2, X_3, \dots, X_n , is given as

$$H(X_1 | X_2, X_3, \dots, X_n) = - \sum_{x_1, x_2, \dots, x_n} p(x_1, x_2, x_3, \dots, x_n) \log_2 p(x_1 | x_2, x_3, \dots, x_n) \quad (6)$$

The conditional entropy could be used to measure the information gain when partial information or context of other random variables is known. Equipped with information theory concepts, we now apply these information measures to chemical systems and molecules.

4.2. Information theory and chemical representations

Studies in chemical information theory (Bonchev and Trinajstić, 1982) have demonstrated the promise of entropic perspective in chemistry (Chandler, 2017; Graham, 2002; Nalewajski and Parr, 2001; 2000). We analyze various chemical representations, namely, the SMILES representations, molecular formulas, and our proposed SMILES grammar-based representations from the perspective of Shannon entropy. We quantify the superiority of certain representations when compared to the others and highlight their inherent benefits when used in machine learning algorithms.

In our framework, we consider the individual tokens in various representations as random variables that contain *bits of information* required to reconstruct a given molecule. The representations are therefore a sequence of random variables, X_1, X_2, \dots, X_n , where n is the length of the representation for a given molecule and X_i could take any of the M possible tokens defined in the vocabulary of the representation. For instance, consider the earlier example of benzene from Section 3.1. The corresponding random variables for each of the three representations is given by,

- Molecular formula (C_6H_6): $X_i^{Mo} \in \{'C', 'H'\}$, where $M = 3$, $n = 4$
- SMILES ($C1 = CC = CC = C1$): $X_i^S \in \{'C', '1', '='\}$, where $M = 3$, $n = 11$
- Grammar¹ (1, 3, 2, 3, ..., 12, 19): $X_i^G \in \{1, 2, 3, 4, 5, 6, 12, 13, 15, 17, 19\}$, where $M = 11$, $n = 32$

Defining the random variables and computing their probability distributions over all the molecules in the dataset, we compute the corresponding information measures using Shannon entropy in Eq. 4. Since our objective is to quantify the information capacity for an entire representation instead of certain specific molecules, this distribution is computed over all the possible lengths of representations, n , in the dataset. Similarly, the conditional information measure in Eq. 6 could be computed using the conditional

¹ Using the representative grammar in 1.

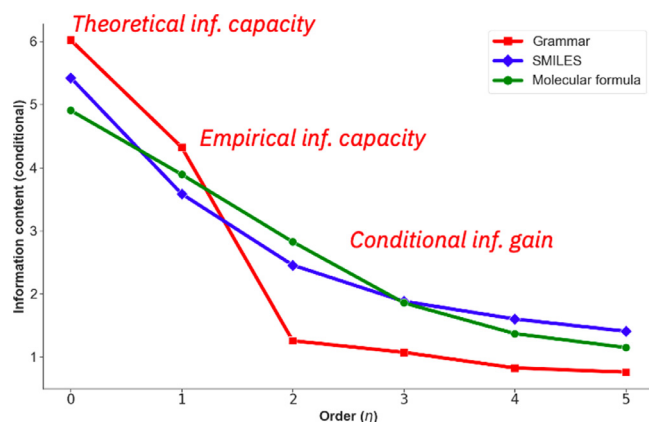


Fig. 5. Information content vs order of conditioning (η) for the three representations.

Table 2

Information content (i_η) for various orders of conditioning (η) for the three representations.

	SMILES	Grammar	Molecular Formula
i_0	5.426	6.022	4.906
i_1	3.583	4.322	3.891
i_2	2.453	1.254	2.823
i_3	1.879	1.070	1.855
i_4	1.599	0.822	1.367
i_5	1.404	0.756	1.146

distribution of random variables based on the co-occurrence matrices (up to a given order) of the random variables in the database. The order indicates the number of random variables under consideration, with $\eta - 1$ conditioned random variables for an order of η . An order $\eta = 1$ corresponds to Shannon entropy (Eq. 4), order $\eta = 0$ corresponds to Shannon entropy when the random variables follow a uniform distribution, and orders $\eta > 1$ correspond to conditional entropy with conditioning on $\eta - 1$ random variables (Eq. 6).

The Shannon entropy and conditional information measures, presented in Fig. 5 and Table 2, are computed (at various orders) using Eqs. 4 and 6, respectively. The USPTO 50K test set is used to estimate the required (conditional) probability distributions for the three representations (SMILES, grammar, and molecular formula) based on the co-occurrence matrices, conditioned on a given number of tokens according to the order of conditioning. The probability distributions for the random variables are computed using the three representations for all the molecules in the test-set of the USPTO 50K reaction dataset to limit computational requirements, especially for calculating the conditional distributions. We evaluate the maximum conditional distribution up to an order of $\eta = 5$. The molecular formulas are extracted from the SMILES representations of molecules using the 'rcdk' library in R.

It follows from our discussion in the earlier section that the maximum information (corresponding to i_0) is achieved when the random variables follow a uniform distribution and all the bits have the same probability ($1/M$) of taking a given value. Thus, i_0 is independent of the dataset under consideration and is purely a property of the representation that is indicative of its information storing capacity. Based on Fig. 5, the grammar-representations have much higher information capacity, followed by the SMILES representation and then the molecular formulas, highlighting the theoretically high information capacity of grammar representations.

When the order of analysis is increased to 1, the information capacity decreases for all the representations, indicating that the underlying probability distributions are far from uniform, with cer-

tain values more likely than others. This is expected since in any chemical representation, the identifiers for atoms such as C and H are significantly more likely to occur when compared to others such as F or B. It could be inferred through the probability versus identifier index plot depicted in Fig. 6 that the SMILES and molecular formula representations are much more skewed, with a majority of the identifiers occurring much more frequently than the others. On the other hand, the grammar-based representations' identifiers exhibit a much smoother and slower decay, indicating more evenly distributed probabilities for the identifiers. This validates the richness of grammar-representations due to the incorporation of structural-hierarchy, an argument that we made qualitatively in our earlier work (Mann and Venkatasubramanian, 2021).

As the order of conditioning while computing the information measure is increased to $\eta = 2$, a drastic decrease in the information content is observed for grammar-representations, and the conditional information content remains significantly lesser than the other representations even for higher values of η . This could be attributed to the in-built redundancy in the grammar-representations incorporated by means of a hierarchical sequence of production rules encoded in a molecule's representation. This transforms into lower values of conditional probabilities when an identifier's context in terms of the preceding tokens is known. Qualitatively, this means that when the context of a token is provided, the uncertainty associated with the possible values it could take is much lesser than its equivalent in the SMILES representation and molecular formula-based representations.

It is interesting to also note from Fig. 5 that the conditional information content plots intercept twice for the SMILES and molecular formula representations, which could possibly be explained as follows – the first intercept is due to the relative differences in the theoretical information content (i_0) and the actual information content (i_1) computed using the conditional probabilities from the database, indicating that the conditional probabilities at order 1 are much more skewed for the SMILES representation (translating to lower entropy) since it has more tokens that are repeated compared to molecular formulas; and the second intercept at order 3 could be due to the trade-off between the number of tokens and redundancy in the representations where the higher number of tokens for SMILES begin to contribute more towards the conditional entropy (uncertainty) even after a reduction in entropy due to the partial knowledge of the context (preceding tokens). In contrast, the grammar-based representation consistently has the highest theoretical (i_0) and actual (i_1) information content, and the lowest conditional entropy (highest redundancy) beyond order 2. This clearly demonstrates the ability of grammar-based representations in overcoming the associated trade-off between higher number of tokens and redundancy that the SMILES and molecular formula-based representations suffer from.

In summary, the underlying redundancy in grammar-representations, indicated by i_η with $\eta \geq 2$, could be leveraged by machine learning algorithms that model the long and short-range dependencies between tokens in a given sequence, such as the class of sequence-to-sequence models used in our work. In addition, the higher information-storage capacity of these representations, as indicated by i_0 and i_1 , implies that they are much richer when compared to the other representations and therefore contain additional bits of information that is lacking in the other representations and could be crucial for the adequate differentiation between molecules in the latent space. There are other representations such as International Chemical Identifier (InChI) that are used to represent molecules and performing a similar information-theoretic analysis on such representations would be part of our future work in this direction.

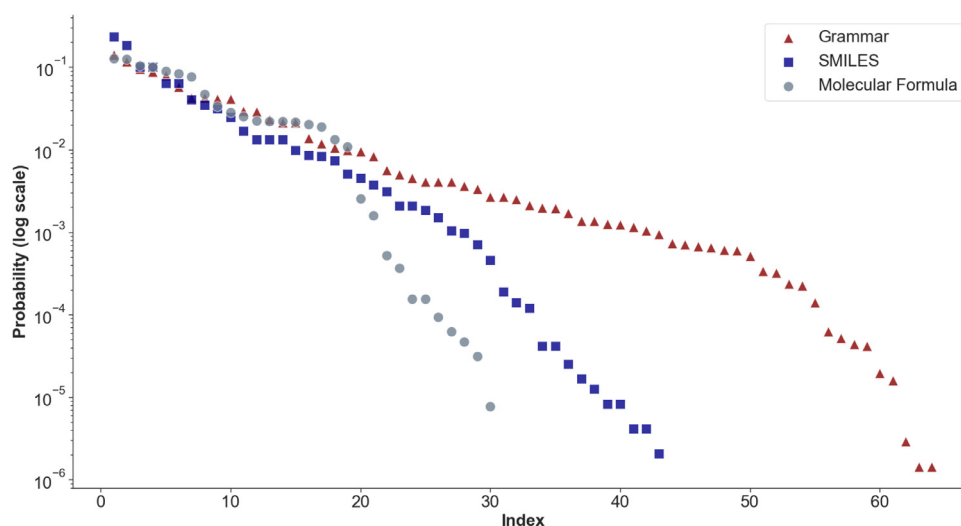


Fig. 6. Probability of occurrence of a given token versus the sorted index.

5. Data and model training

We demonstrate our model's performance using a standard retrosynthesis prediction dataset which is a filtered dataset derived from the text extraction work done on US Patents and Trademark Office's (USPTO) database (Lowe, 2012) and further classified into ten different reaction classes (Schneider et al., 2016). The filtered dataset contains only the reactants and products, with the reagent information removed and the SMILES strings canonicalized. Further, similar to Liu et al. (2017), the multiple product reactions are split into multiple reactions so that each reaction contains only a single major product. This dataset is referred to as the USPTO 50K dataset in the literature.

In order to use our approach, we encode the SMILES strings corresponding to all the molecules in the database into their equivalent grammar representations as described in Section 3.1. This implies that since we are working with a subset of the official OpenSMILES grammar, certain molecules that are not in grammar are skipped and therefore are not included in the model training stage. Table 3 summarizes the reaction database with the number of reactions that are in grammar along with the train, validation, and test-set splits. Table 4 summarizes the distribution of the various reactions across the 10 reaction classes.

Since the retrosynthesis prediction task involves predicting a set of precursors that could be used for obtaining a given product molecule, we define identifiers that distinguish the various reactant molecules (grammar-representation) from each other and also indicate the end of the set. These two additional tokens convey to the model the separation between various precursors' representations and also the end of the entire set of precursors. The reaction class identifiers are appended at the start of the source (product) molecule's representation while evaluating the model performance under known reaction type scenarios. This additional step is skipped when the model performance is evaluated for the unknown reaction type scenario. A schematic for this is shown in Fig. 7.

Table 3

An overview of the retrosynthesis dataset used in our work.

Dataset	train	valid	test	total
USPTO 50K				
with (sanitized) molecules	40,029	5004	5004	50,037
in grammar	38,995	4861	4861	48,717

We train the transformer model for this task using a cross-entropy-based loss function that minimizes the sequence-to-sequence translation error. The model was trained using the Adam optimizer (Kingma and Ba, 2014) with beta $\beta_1 = 0.9$, $\beta_2 = 0.98$, and $\epsilon = 10^{-9}$, and a cyclic learning rate schedule that is characterized by a fixed number of warmup steps given by

$$lr = d_{model}^{-0.5} \cdot \min(step_num^{-0.5}, step_num * warmup_steps^{-0.5}) \quad (7)$$

where d_{model} is the embedding dimension (positional). At the training stage, to avoid overfitting, a dropout layer is used for both the feed-forward networks and the attention-mechanism, for the encoder and the decoder. A masking approach similar to Kusner et al. (2017) is used for generating the output SMILES strings from the decoded grammar-representation. A loss function based on sparse categorical cross-entropy between the predicted and actual target sequences is minimized. The possible and the best hyperparameters identified for the model are given in Table 5. The lengths of the input and output representations to the model are fixed at 301 and 900, respectively.

Both the models were trained using TensorFlow 2.1 and python 3.7 for 12 cycles (~ 700 epochs). For generating the parse-trees and extracting grammar-based features, we used the Natural Language ToolKit (NLTK) 3.4.5 library. The molecular datasets were processed using the 2019 release of the RDKit library.

6. Results and discussion

In this section, we define the performance metrics, evaluate the model's performance on the test-set of the USPTO 50K dataset, and benchmark the performance of our approach against other similar works in this area, highlighting the advantages and limitations of this framework.

6.1. Evaluation metrics

We evaluate our model's performance using the following metrics – accuracy, which captures the ability to perfectly predict all the precursor molecules; fractional accuracy, which indicates the fraction of accurately predicted precursors from the set of molecules in the ground truth; and syntactic validity, meaning the percentage of grammatically valid predictions. In addition, we also compute the accuracy of prediction of the Maximal Fragment or MaxFrag Tetko et al. (2020) indicating the prediction accuracy of the longest reactant involved and report the average BLEU (bilingual evaluation understudy) Papineni et al. (2002) and similarity

Table 4
Distribution of reactions across different reaction classes that are in-grammar.

Reaction class	Reaction name	train	valid	test	total
1	Heteroatom alkylation and arylation	11,886	1,476	1478	14,840
2	Acylation and related processes	9358	1,165	1169	11,698
3	C – C bond formation	4324	544	539	5407
4	Heterocycle formation	710	89	90	889
5	Protections	513	64	62	639
6	Deprotections	6357	796	789	7942
7	Reductions	3607	448	452	4,507
8	Oxidations	629	80	79	788
9	Functional group interconversion (FGI)	1434	176	180	1790
10	Functional group addition (FGA)	177	23	23	223



Fig. 7. The retrosynthesis reaction encoding strategy used in the machine translation framework. The identifier '80' indicates the end of a given molecule's grammar-representation, '81' indicates the separation between two precursor molecules, and '82' indicates the end of the entire set of precursor molecules. The additional token indicating the reaction type is optional and we report the model performance under both the scenarios with known and unknown reaction classes.

Table 5
Possible and best hyperparameter values for the transformer model architecture described in Fig. 3.

Hyperparameter	Possible values	Final model
Embedding dimensions	64, 128, 256	256
Attention heads	4, 8, 16	8
Feedforward network units	512, 1024, 2048	512
Number of layers	4, 6	4
Dropout	0.1, 0.2	0.1
Warmup steps	4k, 8k, 12k	8k

Table 6
Accuracy, fractional accuracy, and syntactic validity on the test set with known reaction class.

Performance measure	top-1	top-3	top-5	top-10
Accuracy	43.8	57.2	61.4	66.6
Fractional accuracy	53.8	65.4	69.2	73.7
Syntactic validity	95.6	92.8	91.6	90.4

scores for this maximal fragment. The BLEU score is a standard metric used for evaluation of the quality of machine-translated texts against a reference translation, and the similarity scores² are computed using the similarities between the string substructures of the predictions and the ground truth. These metrics are reported for three example predictions in Fig. 8.

6.2. Results on USPTO 50K dataset

The performance evaluation measures computed on the test set of the USPTO 50K dataset are presented in Tables 6 and 8 for the known reaction class scenario and in Tables 7 and 9 for the scenario when reaction classes are not known. We observe from Table 6 that though the top-10 accuracy is 66.6%, the fractional accuracy at 73.7% is much higher and indicates that a major fraction of the ground truth reactants is accurately predicted across reactions. The syntactic validity is as high as 95.6% for the top-1 predictions and 90.4% for the top-10 predictions. The decreasing trend in syntactic validity is expected since as the number of predictions

² Computed using the SequenceMatcher routine in python that matches the longest contiguous matching sub-sequence that does not contain any unwanted (or junk) elements.

Table 7
Accuracy, fractional accuracy, and syntactic validity on the test set with unknown reaction class.

Performance measure	top-1	top-3	top-5	top-10
Accuracy	32.1	44.3	48.9	54.0
Fractional accuracy	39.6	51.5	56.2	61.8
Syntactic validity	94.9	92.6	91.6	90.3

Table 8
MaxFrag accuracy and the corresponding BLEU and similarity scores on the test set with known reaction class.

Performance measure	top-1	top-3	top-5	top-10
MaxFrag accuracy	50.4	62.1	65.7	70.2
BLEU score	74.8	83.4	85.2	87.4
Similarity score	80.0	87.2	88.6	90.2

Table 9
MaxFrag accuracy and the corresponding BLEU and similarity scores on the test set with unknown reaction class.

Performance measure	top-1	top-3	top-5	top-10
MaxFrag accuracy	38.1	49.1	53.2	58.4
BLEU score	67.5	76.0	78.3	81.1
Similarity score	75.7	82.2	83.8	85.7

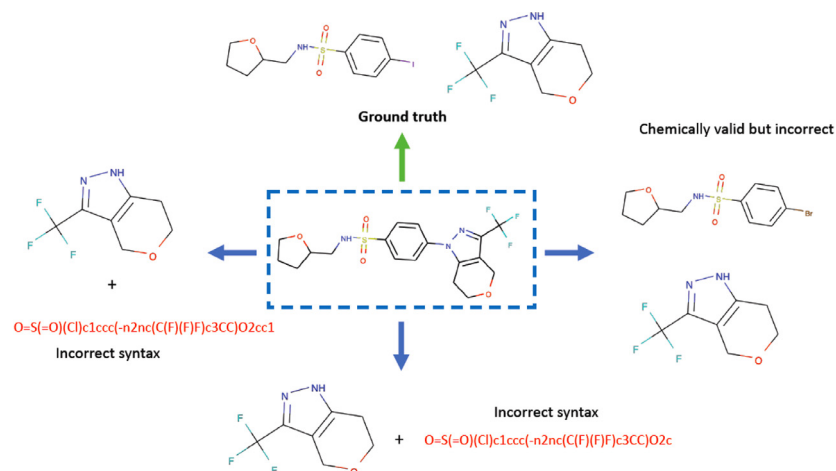
increases, the invalid predictions go up because of the model's susceptibility to decode grammatically invalid strings.

The similarity scores in Table 8 indicate that the MaxFrag precursor is predicted with a top-10 accuracy of over 70% and a similarity score of over 90%, highlighting the model's ability to correctly identify the characteristics of the most critical molecule (in classical retrosynthesis) with a fairly high degree of accuracy. The corresponding BLEU scores also indicate the good quality of translation that is achieved for the MaxFrag molecule.

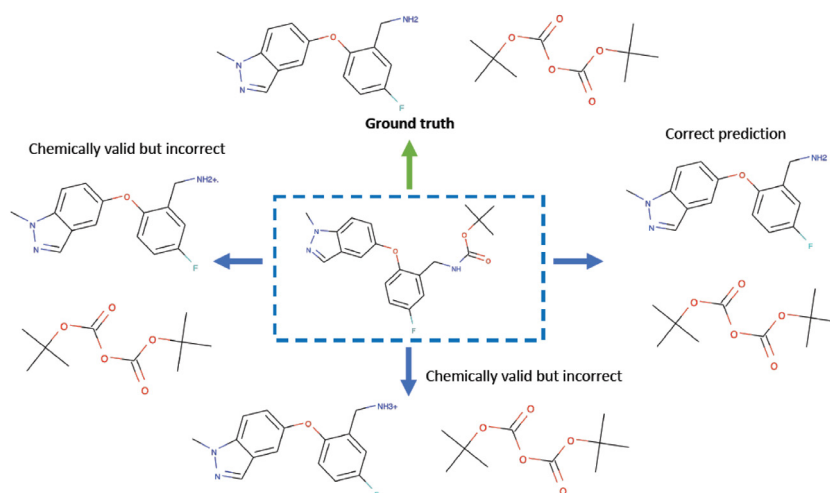
Few of the example top-3 predictions along with the prediction inaccuracies and performance metrics are presented in Fig. 8.

6.3. Performance across reaction classes

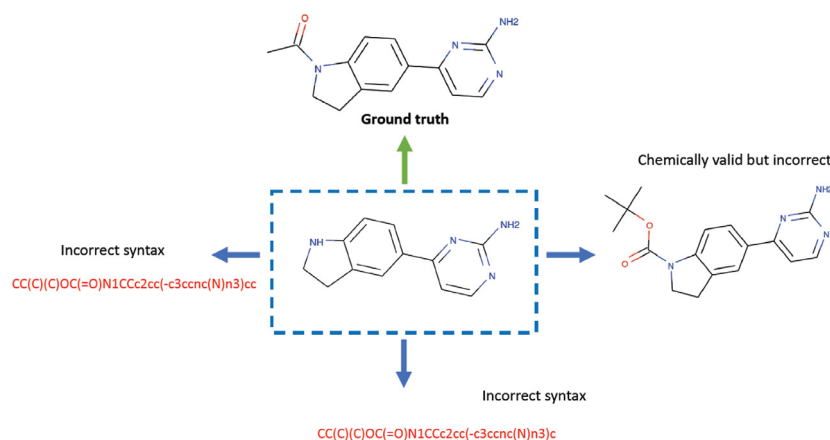
In order to further understand the performance of our model across reaction classes, we increase the granularity of the anal-



(a) Example from reaction class 1; accuracy: 0.0, fractional accuracy: 0.5; syntactic validity: 0.67, MaxFrag accuracy: 0.0, MaxFrag similarity: 0.56, MaxFrag BLEU: 0.36



(b) Example from reaction class 5; accuracy: 1.0, fractional accuracy: 1.0; syntactic validity: 1.0, MaxFrag accuracy: 1.0, MaxFrag similarity: 1.0, MaxFrag BLEU: 1.0



(c) Example from reaction class 6; accuracy: 0.0, fractional accuracy: 0.0; syntactic validity: 0.33, MaxFrag accuracy: 0.0, MaxFrag similarity: 0.89, MaxFrag BLEU: 0.79

Fig. 8. Example top-3 predictions made by our model and their corresponding evaluation metrics indicated in the figure captions.

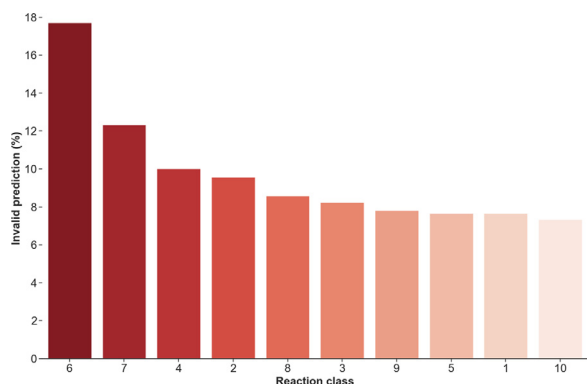


Fig. 9. Invalid percentages for Top-10 predictions with known reaction class.

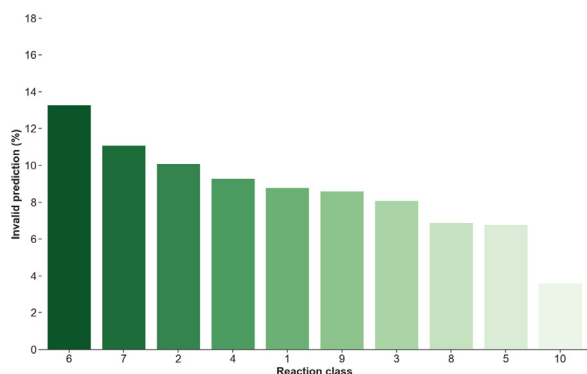


Fig. 10. Invalid percentages for Top-10 predictions without reaction class.

ysis and compute the five metrics— accuracy, fractional accuracy, MaxFrag accuracy, similarity score, and syntactic validity across the 10 reaction classes. The detailed measures of these metrics are summarized in Tables B.13, B.14, B.15, and B.16 for the known reaction class scenario, and in Tables B.17, B.18, B.19, and B.20 for the unknown reaction class scenario in Appendix 2. The fraction of invalid predictions across the various reaction types for top-10 analysis are presented in Figs. 9 and 10.

The above trend indicates that except for reaction class 6 (deprotections) and the surprisingly more accurate predictions on reaction class 10 (functional group addition) when the reaction class is unknown, all the reaction types result in nearly the same percentage of invalid predictions. A likely possibility for this observation could be the model learning the underlying grammar, irrespective of the number of samples in each class or the chemical transformations occurring across the different reaction types. This behavior is not trivial since the corresponding top-10 prediction accuracy in Tables B.16 and B.20 do not follow the same trend across reaction classes. Moreover, the percentage of invalid predictions shows only minor variations across the two scenarios with known and unknown reaction classes. This observation again highlights the ability of our proposed SMILES grammar-based representations to force the model to learn the underlying grammar and consequently generate grammatically correct predictions, irrespective of the other factors. The high percentage error in deprotection reactions could be attributed to several factors that could be specific to the reaction class and could be analyzed through chemistry-driven heuristics that we envision as a hybrid explanation-generation system as a future extension of this work.

6.4. Comparison with other works

Here, we compare the performance of our model against other similar works in this area. One of the first benchmarks in retrosynthesis prediction using seq2seq models on SMILES string representations is by Liu et al. (2017). Their framework is similar to ours in that there are no post-processing of predictions, data augmentation strategies, and model performance-boosting methods used for further improving the model performance – techniques that usually result in improved accuracy custom-fit to a given setting. Our objective is to propose an alternative formulation that is fundamentally different from the other approaches in that it ensures incorporation of chemistry knowledge, forcing the model to learn the underlying SMILES grammar and minimize invalid predictions.

Table 10 compares the prediction accuracy against those reported in Liu et al. We observe that our model improves the prediction accuracy by a margin of $\sim 5\%$ across all the top-N measures and reduces the percentage of invalid predictions by 53% – 64% when the reaction class is known. We attribute the higher accuracy and the reduced invalid predictions to the grammar-representations that incorporate structural information about the molecules and are characterized by much higher redundancies when compared to SMILES strings as demonstrated using the information-theoretic analysis in Section 4. Fig. 11 demonstrates our model's ability to outperform the top-10 accuracy reported in Liu et al. across reaction classes, often by a significant margin. For completeness, we also compare our models' performance with that reported in Liu et al. when the reaction class is unknown. Since they did not evaluate the model under this setting, we use the implementation of Lin et al. (2020) that evaluated the performance of this model with unknown reaction class. This is of interest for retrosynthetic planning under certain scenarios where no chemistry information about the target molecule is known apriori. A comparison of the accuracy reported under this scenario is presented in Table 11. We report the detailed class-wise results for both the models (with and without reaction class information) in Appendix 2.

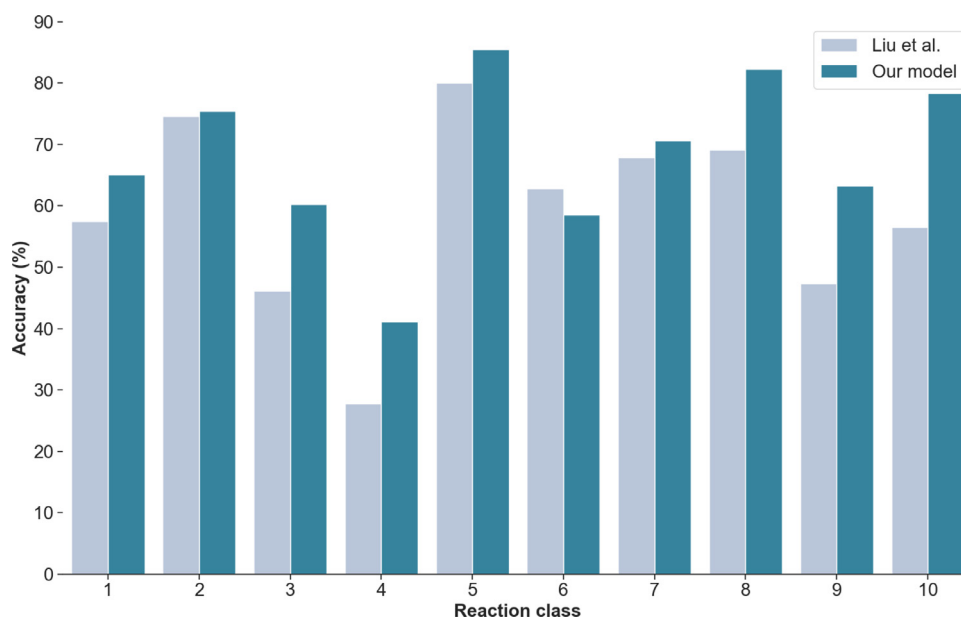
As mentioned earlier, it is possible to achieve even higher prediction accuracy through additional performance boosting techniques as demonstrated in the following studies. Zheng et al. (2020) used an additional transformer model that takes as input the output of another transformer model to correct the invalid predictions. Tetko et al. (2020) proposed data augmentation strategies that significantly increased the size of the dataset used for building a transformer model for retrosynthesis. Karpov et al. (2019) used model ensembling, snapshot learning methods, and increasing beam search temperature to improve the model performance. Lin et al. (2020) used averaging of model weights and combination with Monte Carlo Tree Search (MCTS) strategies for proposing retrosynthesis routes. The accuracy for such augmented (template-based and template-free) models vary significantly and the top-1 accuracy could be as high as 65%. However, we emphasize here that the objective of this work is not to pursue the state-of-the-art but to highlight the benefits of incorporating prior chemistry knowledge into such black-box models. We have shown that incorporation of this knowledge translates to higher accuracy and fewer invalid predictions when compared to purely black-box models using the same framework. Applying such data augmentation and transfer-learning strategies to boost the model performance further would be the future extension of our work which we conjecture would further improve the model accuracy significantly.

Finally, we would like to highlight here that building models that leverage as much prior chemistry knowledge as possible would be more reliable, acceptable, and explainable as com-

Table 10

Comparison with other similar works involving purely seq2seq models and USPTO 50K dataset with known reaction classes.

Model	Top-N measure (with reaction class)			
	accuracy (%) invalid (%)			
	1	3	5	10
Liu Seq2Seq (Liu et al., 2017)	37.4 12.2	52.4 15.3	57.0 18.4	61.7 22.0
Our work	43.8 4.4	57.2 7.2	61.4 8.4	66.6 9.6

**Fig. 11.** Comparison of top-10 accuracies across different reaction classes.**Table 11**

Comparison with other similar works involving purely seq2seq models and USPTO 50K dataset with unknown reaction classes.

Model	Top-N measure (without reaction class)			
	accuracy (%) invalid (%)			
	1	3	5	10
Liu Seq2Seq ³	28.3 -	42.8 -	47.3 -	52.8 -
Our work	32.1 5.1	44.3 7.4	48.9 8.4	54.0 9.7

³ As implemented in Lin et al. (2020); the invalid fractions were not reported for this model.

pared to purely data-driven, black-box models that completely disregard known underlying chemistry. Such prior chemistry knowledge could be in the form of information about molecules (grammar, molecular-graph, or structure-based representations), possible reactions (reaction class information, molecular descriptors), model architecture and workflow (that mimic expert chemists), and other similar approaches utilizing deeper integration of first-principles with machine learning-based models.

7. Conclusions

Retrosynthetic analysis is a challenging problem since it involves predicting the precursors with limited information, searching a combinatorially large number of possible synthesis pathways, and approximating an often complex multi-step analysis as a single-step prediction problem. Naturally, incorporating additional information about the reaction or the molecules involved would

be of considerable use given the complexity of the task and the limited information often present for making the predictions. Towards that goal, we have proposed grammar-based representations of molecules that incorporate chemical and structural information extracted from their SMILES string representations. We have shown in our earlier work (Mann and Venkatasubramanian, 2021) that such representations successfully overcome over-parameterization in models for the forward reaction prediction. Here, we have quantified the superiority of SMILES grammar-based representations compared to the character-based SMILES representations from an information-theoretic standpoint. We have shown that such representations have higher information capacity captured by the Shannon entropy computed for molecules in the USPTO 50K dataset. Moreover, the conditional entropy measures highlighted the higher redundancy built-in to these representations, making them better-suited for machine learning architectures.

The performance of our model reinforced the above observations. We report the top-1 prediction accuracy of 43.8% and syntactic validity of 95.6% as opposed to 37.4% and 87.8%, respectively, reported in Liu et al. We have shown that not only does our model outperform the aggregate statistics reported in Liu et al., the performance of our model across the various reaction classes is much better. An interesting observation is that owing to the grammar representations, our model results in nearly the same percentage of invalid predictions across reaction classes – independent of reaction type, the class-wise number of reactions in the training set, and the known or unknown reaction class scenarios. Moreover, the MaxFrag similarity, which could be as high as 90%, indicates that the model predicts the major precursors required for synthesis fairly accurately. The future extension of our work would

involve solving the multi-step retrosynthesis problem and incorporating additional contextual information about the reactions into the same framework.

Declaration of Competing Interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

CRediT authorship contribution statement

Vipul Mann: Conceptualization, Methodology, Writing – original draft. **Venkat Venkatasubramanian:** Conceptualization, Funding acquisition, Formal analysis, Writing – original draft.

Acknowledgments

This work was supported in part by the Center for the Management of Systemic Risk (CMSR) at Columbia University.

As the corresponding author, I take this opportunity to express my sincere appreciation and gratitude to Professor George Stephanopoulos for his unwavering support and encouragement throughout my career. As an outsider in the PSE community, lack-

ing the well-known PSE pedigrees of my distinguished colleagues, I recall finding the start of my academic start particularly challenging. Furthermore, I was among the very small minority of researchers exploring the use of artificial intelligence (AI) in chemical engineering in the 1980s, an area that was generally considered a lost cause in those early years. Despite these twin challenges, if I seem to have survived the trials and tribulations, it is due to the invaluable mentoring assistance from Professor George Stephanopoulos. It seems apropos to acknowledge and record my debt to him here in this special issue in his honor.

Appendix A

The SMILES grammar used in this work is the same as that used in our previous work on the forward prediction problem (Mann and Venkatasubramanian, 2021). This grammar comprises 80 production rules with 24 non-terminal symbols specifying the different structural components of a SMILES string. All the production rules for the grammar used in our work are summarized in Table A.12. The first and the last production rules, SMILES \rightarrow CHAIN and NOTHING \rightarrow NONE, are additional rules signifying the start and end of a SMILES string, which are analogous to the <START> and <END> tokens in natural language processing marking the beginning and the end of sentences, respectively.

Table A.12

SMILES grammar used in Grammar Ontology-based Prediction of Reaction Outcomes (or GO-PRO) (Mann and Venkatasubramanian, 2021).

S.No	Production rules
1	SMILES \rightarrow CHAIN
2	ATOM \rightarrow BRACKET_ATOM ALIPHATIC_ORGANIC AROMATIC_ORGANIC
3	ALIPHATIC_ORGANIC \rightarrow B C N O S P F I Cl Br
4	AROMATIC_ORGANIC \rightarrow c n o s p
5	BRACKET_ATOM \rightarrow [BAI]
6	BAI \rightarrow ISOTOPE SYMBOL BAC SYMBOL BAC ISOTOPE SYMBOL SYMBOL
7	BAC \rightarrow CHIRAL BAH BAH CHIRAL
8	BAH \rightarrow HCOUNT BACH BACH HCOUNT
9	BACH \rightarrow CHARGECLASS CHARGE CLASS
10	SYMBOL \rightarrow ALIPHATIC_ORGANIC AROMATIC_ORGANIC ELEMENT_SYMBOLS
11	ISOTOPE \rightarrow DIGIT DIGIT DIGIT DIGIT DIGIT DIGIT
12	DIGIT \rightarrow 1 2 3 4 5 6 7 8
13	CHIRAL \rightarrow @ @@
14	HCOUNT \rightarrow H H DIGIT
15	CHARGE \rightarrow - - DIGIT - DIGIT DIGIT + + DIGIT + DIGIT DIGIT
16	BOND \rightarrow - = # / \ \
17	RINGBOND \rightarrow DIGIT BOND DIGIT
18	BRANCHED_ATOM \rightarrow ATOM ATOM RB ATOM RB BB
19	RB \rightarrow RB RINGBOND RINGBOND
20	BB \rightarrow BB BRANCH BRANCH
21	BRANCH \rightarrow (CHAIN) (BOND CHAIN)
22	CHAIN \rightarrow BRANCHED_ATOM CHAIN BRANCHED_ATOM CHAIN BOND BRANCHED_ATOM
23	CLASS \rightarrow DIGIT
24	ELEMENT_SYMBOLS \rightarrow H
25	NOTHING \rightarrow NONE

Appendix B

The detailed results capturing the model performance for the five metrics – accuracy, fractional accuracy, syntactic validity, maximal fragment (MaxFrag) accuracy and maximal fragment (MaxFrag) similarity are reported here. [Tables B.13, B.14, B.15, and B.16](#) present the results for the top-1, top-3, top-5, and top-10 predictions, respectively, when the reaction class is known. [Tables B.17, B.18, B.19, and B.20](#) present the results for the top-1, top-3, top-5, and top-10 predictions, respectively, when the reaction class is unknown.

Scenario 1: Reaction class known

Table B.13

The top-1 accuracy, fractional accuracy, MaxFrag accuracy and MaxFrag similarity scores across the reaction classes (in %).

Top-1 measure	Reaction class									
	1	2	3	4	5	6	7	8	9	10
Accuracy	40.9	52.2	37.7	26.7	66.1	35.4	50.4	69.6	38.3	56.5
Fractional accuracy	54.9	67.0	49.9	39.4	81.5	35.4	50.4	75.3	45.0	71.7
Syntactic validity	96.7	95.8	96.4	94.1	97.6	90.9	95.2	97.1	98.8	97.8
MaxFrag accuracy	50.3	61.3	44.9	37.8	83.9	35.4	50.4	77.2	44.4	60.9
MaxFrag similarity	82.0	86.2	78.6	71.2	91.8	68.0	79.5	89.0	78.4	82.9

Table B.14

The top-3 accuracy, fractional accuracy, MaxFrag accuracy and MaxFrag similarity scores across the reaction classes (in %).

Top-3 measure	Reaction class									
	1	2	3	4	5	6	7	8	9	10
Accuracy	54.3	66.6	51.4	31.1	80.6	49.9	61.3	77.2	51.7	73.9
Fractional accuracy	66.3	77.4	62.5	49.4	89.5	49.9	61.3	82.9	57.8	80.4
Syntactic validity	94.5	92.6	93.9	91.8	94.4	86.3	89.6	94.4	95.6	91.9
MaxFrag accuracy	61.0	72.1	58.3	46.7	90.3	49.9	61.3	84.8	58.9	73.9
MaxFrag similarity	87.0	91.5	84.3	80.3	98.4	81.0	89.2	96.0	88.0	89.4

Table B.15

The top-5 accuracy, fractional accuracy, MaxFrag accuracy and MaxFrag similarity scores across the reaction classes (in %).

Top-5 measure	Reaction class									
	1	2	3	4	5	6	7	8	9	10
Accuracy	59.3	70.9	54.5	36.7	83.9	53.9	64.4	79.7	56.7	73.9
Fractional accuracy	70.7	80.9	66.0	53.9	91.1	53.9	64.4	84.8	61.7	80.4
Syntactic validity	93.4	91.6	92.6	90.9	92.6	84.0	88.9	93.2	94.3	91.6
MaxFrag accuracy	65.2	75.7	61.0	50.0	93.5	53.9	64.4	87.3	62.2	73.9
MaxFrag similarity	88.5	92.2	85.6	81.2	98.5	83.3	90.3	99.2	88.2	89.2

Table B.16

The top-10 accuracy, fractional accuracy, MaxFrag accuracy and MaxFrag similarity scores across the reaction classes (in %).

Top-10 measure	Reaction class									
	1	2	3	4	5	6	7	8	9	10
Accuracy	65.1	75.4	60.3	41.1	85.5	58.6	70.6	82.3	63.3	78.3
Fractional accuracy	75.5	84.2	70.5	58.3	91.9	58.6	70.6	86.7	67.8	82.6
Syntactic validity	92.3	90.4	91.8	90.0	92.3	82.3	87.7	91.4	92.2	92.7
MaxFrag accuracy	69.8	79.4	65.7	53.3	93.5	58.6	70.6	88.6	67.8	78.3
MaxFrag similarity	90.1	93.5	87.8	82.9	98.6	85.0	92.0	98.9	90.5	91.8

Scenario 2: Reaction class unknown

Table B.17

The top-1 accuracy, fractional accuracy, MaxFrag accuracy and MaxFrag similarity scores across the reaction classes (in %).

Top-1 measure	Reaction class									
	1	2	3	4	5	6	7	8	9	10
Accuracy	32.2	39.9	25.8	10.0	35.5	26.7	35.8	38.0	16.1	60.9
Fractional accuracy	42.6	52.2	33.9	15.0	44.4	26.8	35.8	40.5	22.5	63.0
Syntactic validity	95.4	94.8	96.1	92.4	96.4	92.0	94.4	96.8	96.3	100.0
MaxFrag accuracy	40.7	48.9	33.0	11.1	50.0	26.9	35.8	41.8	21.7	60.9
MaxFrag similarity	76.7	80.6	72.9	57.3	79.2	70.7	74.3	83.0	72.4	81.6

Table B.18

The top-3 accuracy, fractional accuracy, MaxFrag accuracy and MaxFrag similarity scores across the reaction classes (in %).

Top-3 measure	Reaction class									
	1	2	3	4	5	6	7	8	9	10
Accuracy	43.5	54.8	37.1	15.6	43.5	39.9	48.0	48.1	23.9	69.6
Fractional accuracy	53.5	65.7	46.0	21.1	53.2	39.9	48.0	51.3	30.3	71.7
Syntactic validity	93.2	92.7	93.8	90.0	94.8	89.4	92.0	95.5	93.3	97.7
MaxFrag accuracy	50.3	60.7	44.2	17.8	58.1	39.9	48.0	51.9	28.9	69.6
MaxFrag similarity	82.5	86.7	78.3	65.1	85.1	79.2	83.1	87.2	78.2	89.0

Table B.19

The top-5 accuracy, fractional accuracy, MaxFrag accuracy and MaxFrag similarity scores across the reaction classes (in %).

Top-5 measure	Reaction class									
	1	2	3	4	5	6	7	8	9	10
Accuracy	48.5	59.8	41.7	17.8	46.8	44.2	51.5	48.1	30.0	69.6
Fractional accuracy	58.7	70.6	51.5	27.8	58.9	44.2	51.5	51.9	35.3	71.7
Syntactic validity	92.3	91.4	93.3	90.5	94.0	88.0	90.6	94.6	92.8	96.4
MaxFrag accuracy	54.6	64.7	48.8	24.4	64.5	44.2	51.5	51.9	32.8	69.6
MaxFrag similarity	83.8	88.2	80.1	68.0	89.1	81.5	84.8	87.8	78.8	88.8

Table B.20

The top-10 accuracy, fractional accuracy, MaxFrag accuracy and MaxFrag similarity scores across the reaction classes (in %).

Top-10 measure	Reaction class									
	1	2	3	4	5	6	7	8	9	10
Accuracy	53.2	64.4	49.0	20.0	58.1	50.1	55.5	57.0	34.4	69.6
Fractional accuracy	64.8	74.4	59.1	31.7	71.8	50.1	55.5	59.5	40.6	73.9
Syntactic validity	91.2	89.9	91.9	90.7	93.2	86.7	88.9	93.1	91.4	96.4
MaxFrag accuracy	59.8	68.6	55.8	30.0	79.0	50.1	55.5	60.8	37.8	69.6
MaxFrag similarity	85.4	89.4	82.3	71.9	92.9	83.9	87.3	88.7	80.1	88.2

Supplementary material

Supplementary material associated with this article can be found, in the online version, at [10.1016/j.compchemeng.2021.107533](https://doi.org/10.1016/j.compchemeng.2021.107533)

References

- Biłozor, A., Kowalczyk, A.M., Bajerowski, T., 2018. Theory of scale-free networks as a new tool in researching the structure and optimization of spatial planning. *J. Urban Plann. Dev.* 144 (2), 04018005.
- Bonchev, D., Trinajstić, N., 1982. Chemical information theory: structural aspects. *Int J Quantum Chem* 22 (S16), 463–480. doi:10.1002/qua.560220845.
- Chandler, J., 2017. An introduction to the foundations of chemical information theory. tarski-Lesniewski logical structures and the organization of natural sorts and kinds. *Information* 8 (1), 15. doi:10.3390/info8010015.
- Chen, J.H., Baldi, P., 2009. No electron left behind: a rule-based expert system to predict chemical reactions and reaction mechanisms. *J Chem Inf Model* 49 (9), 2034–2043.
- Chomsky, N., 1956. Three models for the description of language. *IRE Trans. Inf. Theory* 2 (3), 113–124. doi:10.1109/TIT.1956.1056813.
- Coley, C.W., Barzilay, R., Jaakkola, T.S., Green, W.H., Jensen, K.F., 2017. Prediction of organic reaction outcomes using machine learning. *ACS Cent Sci* 3 (5), 434–443.
- Corey, E., Long, A.K., Greene, T.W., Miller, J.W., 1985. Computer-assisted synthetic analysis. selection of protective groups for multistep organic syntheses. *J. Org. Chem.* 50 (11), 1920–1927.
- Duan, H., Wang, L., Zhang, C., Guo, L., Li, J., 2020. Retrosynthesis with attention-based nmt model and chemical analysis of “wrong” predictions. *RSC Adv* 10 (3), 1371–1378.
- Gothard, C.M., Soh, S., Gothard, N.A., Kowalczyk, B., Wei, Y., Baytekin, B., Grzybowski, B.A., 2012. Rewiring chemistry: algorithmic discovery and experimental validation of one-pot reactions in the network of organic chemistry. *Angew. Chem. Int. Ed.* 51 (32), 7922–7927.
- Graham, D.J., 2002. Information and organic molecules: structure considerations via integer statistics. *J Chem Inf Comput Sci* 42 (2), 215–221. doi:10.1021/ci0102923. PMID: 11911689.
- Jorgensen, W.L., Laird, E.R., Gushurst, A.J., Fleischer, J.M., Gothe, S.A., Helson, H.E., Paderes, G.D., Sinclair, S., 1990. Cameo: a program for the logical prediction of the products of organic reactions. *Pure Appl. Chem.* 62 (10), 1921–1932.
- Karpov, P., Godin, G., Tetko, I.V., 2019. A transformer model for retrosynthesis. In: *International Conference on Artificial Neural Networks*. Springer, pp. 817–830.
- Kingma, D. P., Ba, J., 2014. Adam: A Method for Stochastic Optimization.
- Kusner, M. J., Paige, B., Hernández-Lobato, J. M., 2017. Grammar variational autoencoder.
- Law, J., Zsoldos, Z., Simon, A., Reid, D., Liu, Y., Khew, S.Y., Johnson, A.P., Major, S., Wade, R.A., Ando, H.Y., 2009. Route designer: a retrosynthetic analysis tool utilizing automated retrosynthetic rule generation. *J Chem Inf Model* 49 (3), 593–602.
- Lin, K., Xu, Y., Pei, J., Lai, L., 2020. Automatic retrosynthetic route planning using template-free models. *Chem. Sci.* 11 (12), 3355–3364.
- Liu, B., Ramsundar, B., Kawthekar, P., Shi, J., Gomes, J., Luu Nguyen, Q., Ho, S., Sloane, J., Wender, P., Pande, V., 2017. Retrosynthetic reaction prediction using neural sequence-to-sequence models. *ACS Cent Sci* 3 (10), 1103–1113. doi:10.1021/acscentsci.7b00303. PMID: 29104927.
- Lowe, D.M., 2012. *Extraction of Chemical Structures and Reactions from the Literature*. University of Cambridge.
- Mann, V., Sivaram, A., Das, L., Venkatasubramanian, V., 2021. Robust and efficient swarm communication topologies for hostile environments. *Swarm Evol Comput* 62, 100848.
- Mann, V., Venkatasubramanian, V., 2020. A formal grammar-based machine learning approach for predicting reaction outcomes. 2020 Virtual AIChE Annual Meeting. AIChE.
- Mann, V., Venkatasubramanian, V., 2021. Predicting chemical reaction outcomes: a grammar ontology-based transformer framework. *AIChE J.* 67 (3), e17190. doi:10.1002/aic.17190.
- Nalewajski, R.F., Parr, R.G., 2000. Information theory, atoms in molecules, and molecular similarity. *Proceedings of the National Academy of Sciences* 97 (16), 8879–8882. doi:10.1073/pnas.97.16.8879.
- Nalewajski, R. F., Parr, R. G., 2001. Information Theory Thermodynamics of Molecules and Their Hirshfeld Fragments <https://pubs.acs.org/sharingguidelines>. doi:10.1021/jp004414q.
- Nam, J., Kim, J., 2016. Linking the neural machine translation and the prediction of organic chemistry reactions. *arXiv preprint arXiv:1612.09529*.
- Papineni, K., Roukos, S., Ward, T., Zhu, W.-J., 2002. Bleu: a method for automatic evaluation of machine translation. In: *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, pp. 311–318.
- Pensak, D. A., Corey, E. J., Lhasa-logic and Heuristics Applied to Synthetic Analysis. ACS Publications.
- Salatin, T.D., Jorgensen, W.L., 1980. Computer-assisted mechanistic evaluation of organic reactions. 1. overview. *J. Org. Chem.* 45 (11), 2043–2051.
- Satoh, H., Funatsu, K., 1995. Sophia, a knowledge base-guided reaction prediction system-utilization of a knowledge base derived from a reaction database. *J Chem Inf Comput Sci* 35 (1), 34–44.
- Satoh, K., Funatsu, K., 1999. A novel approach to retrosynthetic analysis using knowledge bases derived from reaction databases. *J Chem Inf Comput Sci* 39 (2), 316–325.
- Schneider, N., Stiefl, N., Landrum, G.A., 2016. What's what: the (nearly) definitive guide to reaction role assignment. *J Chem Inf Model* 56 (12), 2336–2346.
- Schwaller, P., Petraglia, R., Zullo, V., Nair, V.H., Haeuselmann, R.A., Pisoni, R., Bekas, C., Iuliano, A., Laino, T., 2020. Predicting retrosynthetic pathways using transformer-based models and a hyper-graph exploration strategy. *Chem. Sci.* 11 (12), 3316–3325.
- Segler, M.H., Waller, M.P., 2017. Modelling chemical reasoning to predict and invent reactions. *Chemistry—A European Journal* 23 (25), 6118–6128.
- Segler, M.H., Waller, M.P., 2017. Neural-symbolic machine learning for retrosynthesis and reaction prediction. *Chemistry—A European Journal* 23 (25), 5966–5971.
- Shannon, C.E., 1948. A mathematical theory of communication. *Bell System Technical Journal* 27 (3), 379–423. doi:10.1002/j.1538-7305.1948.tb01338.x.
- Szymkuć, S., Gajewska, E.P., Klucznik, T., Molga, K., Dittwald, P., Startek, M., Bajczyk, M., Grzybowski, B.A., 2016. Computer-assisted synthetic planning: the end of the beginning. *Angew. Chem. Int. Ed.* 55 (20), 5904–5937.
- Tetko, I.V., Karpov, P., Van Deursen, R., Godin, G., 2020. State-of-the-art augmented nlp transformer models for direct and single-step retrosynthesis. *Nat Commun* 11 (1), 1–11.
- Ucak, U.V., Kang, T., Ko, J., Lee, J., 2021. Substructure-based neural machine translation for retrosynthetic prediction. *J Cheminform* 13 (1), 4. doi:10.1186/s13321-020-00482-z.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L., Polosukhin, I., 2017. Attention is all you need. In: *Guyon, I., Luxburg, U.V., Bengio, S., Wallach, H., Fergus, R., Vishwanathan, S., Garnett, R. (Eds.), Advances in Neural Information Processing Systems 30*. Curran Associates, Inc., pp. 5998–6008. <http://papers.nips.cc/paper/7181-attention-is-all-you-need.pdf>.
- Venkatasubramanian, V., 2019. The promise of artificial intelligence in chemical engineering: is it here, finally? *AIChE J.* 65 (2), 466–478. doi:10.1002/aic.16489.
- Wei, J.N., Duvenaud, D., Aspuru-Guzik, A., 2016. Neural networks for the prediction of organic chemistry reactions. *ACS Cent Sci* 2 (10), 725–732.
- Weininger, D., 1988. Smiles, a chemical language and information system. 1. introduction to methodology and encoding rules. *J Chem Inf Comput Sci* 28 (1), 31–36. doi:10.1021/ci00057a005.
- Xu, S., Xia, Y., Ouyang, M., 2020. Effect of resource allocation to the recovery of scale-free networks during cascading failures. *Physica A* 540, 123157.
- Zhang, Y., Li, Y., Zhou, Y., Ma, J., 2020. Optimal link rewiring strategy for transport efficiency on scale-free networks with limited bandwidth. *International Journal of Modern Physics C* 31 (02), 2050033.
- Zheng, S., Rao, J., Zhang, Z., Xu, J., Yang, Y., 2019. Predicting retrosynthetic reactions using self-corrected transformer neural networks. *J Chem Inf Model* 60 (1), 47–55.
- Zheng, S., Rao, J., Zhang, Z., Xu, J., Yang, Y., 2020. Predicting retrosynthetic reactions using self-corrected transformer neural networks. *J Chem Inf Model* 60 (1), 47–55. doi:10.1021/acs.jcim.9b00949. PMID: 31825611.