
Probabilistic hypergraph grammars for efficient molecular optimization

Egor Kraev

Mosaic Smart Data Ltd
egor.kraev@gmail.com

Mark Harley

Mosaic Smart Data Ltd
mharley.code@gmail.com

Abstract

We present a novel approach to the issue of molecular optimization. Our approach uses a hypergraph replacement grammar inferred from the ChEMBL database, with grammar construction optimized for molecular structure creation. While constructing the grammar, we also count the frequencies of particular rules being used to expand particular nonterminals in other rules, and use these as priors for the policy model.

Simulating random molecules from the resulting grammar, we show, using GuacaMol distribution benchmarks, that conditional priors result in a molecular distribution closer to the training set than using equal rule probabilities or unconditional priors (total rule frequencies).

We treat the optimization as a reinforcement learning problem, using a batch-advantage modification of the policy gradient algorithm - using individual rewards minus the batch average reward to weight the log probability loss. The reinforcement learning agent is tasked with building molecules using this grammar, with the goal of maximizing benchmark scores available from the literature. To do so, the agent has policies both to choose the next node in the graph to expand and to select the next grammar rule to apply. The policies are implemented using the Transformer architecture with the partially expanded graph as the input.

We show that using the empirical priors as the starting point for a policy eliminates the need for pre-training, and allows us to reach optima faster. We achieve competitive performance on common benchmarks from the literature, such as penalized logP and QED, with only hundreds of batches (without pre-training) on a budget GPU instance.

1 Introduction

A major problem facing the exploration of novel molecules for the purposes of drug design is the vast array of potentially useful compounds – estimated to be in the range of 10^{24} and 10^{60} possible drug-like structures [1, 2]. While it is of course necessary to experimentally determine the usefulness, and safety, of candidate drugs in the laboratory, de novo drug design is an approach to finding candidate molecules through either exhaustive search, or through various generative and machine learning models. This approach takes the form of an optimization procedure over given target scoring functions, giving pre-screened, promising molecules and thereby reducing drug discovery costs.

Deep learning has now been extensively investigated for encoding and generating molecular graphs [3–11], and remains an area of active research. Typically, the approach taken has been to generate a linear molecular representation, such as the SMILES format [12], with an encoder-decoder network architecture similar to that used in machine translation [10].

This route is, however, not optimal for this problem domain. Unlike written text, a molecule’s structure is non-linear – including both cycles and branches. The model is therefore forced not only

to learn to optimize molecules on the given benchmark, but also to learn to generate SMILES strings corresponding to chemically valid molecules. This task is non-trivial and robs capacity of the model from the true task at hand.

A recent development which partially remedied the issue was presented by Kusner, et. al. [9], in which the authors deduce a context-free grammar (CFG) for SMILES strings. This grammar guarantees that only valid SMILES strings will be produced, however not all valid SMILES strings are chemically possible molecules and so some model capacity must still be spent on learning the subset of chemically valid SMILES.

Moving away from linear representations, Kajino [13] proposed the use of a grammar defined on a hypergraph representation of molecular structure. The molecular hypergraph grammar (MHG), a special case of an hyperedge replacement grammar (HRG) [14], uses rules which can be pictured as splitting the molecular graph at non-terminal bonds, and replacing them with another subgraph, thereby constructing any desired molecular structure while guaranteeing that only chemically valid molecules are produced. In particular, this approach does not have issues of invalid atom valences or loss of stereochemistry that other approaches suffer.

Rather than reducing all cliques as in [13], we terminate early allowing cycles of length five and greater to remain. This permits an equally expressive grammar but allows non-trivial cycle structure to be expressed with far fewer rules, resulting in a grammar we call the hypergraph-optimized grammar (HOG). This allows the model to produce novelties without straying too far from reasonable drug-like molecules but, comes at the cost of introducing more rules. We optimize the model’s usage of the HOG by injecting conditional priors for rule selection. After parsing, we count all occurrences of rules conditional on the parent rule – counting all occurrences of (parent, child) pairs – and use these as priors to the rule selection with the effect of grounding the molecular structure in the region of those seen in the training set, but allowing the model to explore further substructure. We infer this grammar from the GuacaMol [15] training set of 1,273,104 SMILES strings.

The second innovation we present, first adapted to this problem in [16], is the use of the Transformer architecture [17] in place of the more typically applied RNNs [18, 19] or, recently, a graph convolutional network [20].¹ Unlike an RNN or CGN, the Transformer’s information distance between any two inputs is always one, giving the network full information about the sequence so far when selecting the next graph node and rule to expand. This will clearly have an impact on the memory performance of the algorithm, which grows as the square of the sequence length, and so it is beneficial that the HOG optimized MHG provides a concise representation of a given molecule.

Furthermore, in place of the encoder-decoder architecture used in previous work [6–10], we treat this problem with a reinforcement learning approach (see [21] **TODOCITE [other RL papers]** for other RL based approaches) with policies selecting directly the next grammar rule and onto which hypergraph node it should be applied. This means that we do not have to learn a latent space representation of the molecules. We train these policies using a batch-advantage modification of the policy gradient algorithm.

Using this approach, we obtain state-of-the-art performance on both common benchmarks from the literature, such as penalized logP and QED, and on more advanced GuacaMol v2 goal-oriented benchmarks. Coupled with a Transformer based discriminator ², the model achieves competitive results on the GuacaMol distribution benchmarks with stable training over a range of hyperparameter values. This is accomplished with only hundreds of steps (without pre-training) on a budget GPU instance.

2 Hypergraph grammar

In order to address the issue highlighted in Sec. 1 relating to SMILES string grammars, we choose to use a *molecular hypergraph grammar* (MHG) as derived by Kajino in [13]. This works by first representing the molecular graph atoms as hyperedges and bonds as hypernodes of a molecular hypergraph. This will produce hypergraphs which are (i) 2-regular and (ii) have constrained cardinality

¹You, et. al. [20] also adopt a reinforcement learning approach with policy gradient optimization, but do not use an explicitly grammar based construction mechanism

²see [11, 22, 23] for applications of discriminator based method, though with a sequential SMILES molecular representation

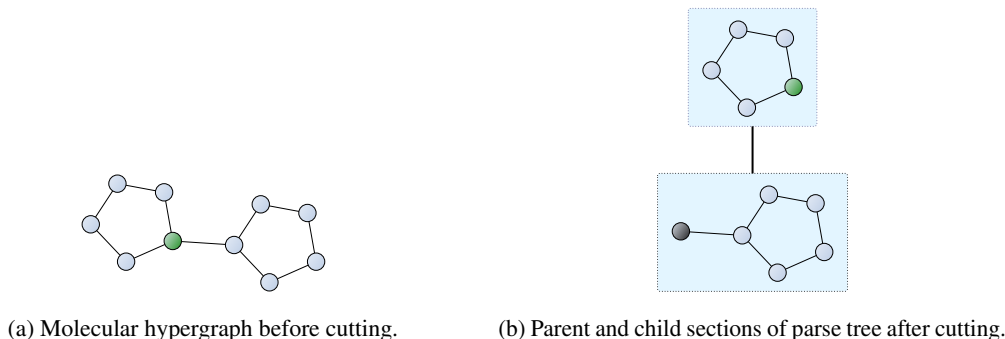


Figure 1: Molecular hypergraph parse tree creation step. Black node indicates the parent node that is merged into the green parent node when the child rule is applied.

of its hyperedges. (i) ensures that the hypergraph can be decoded back to a molecular graph and (ii) preserves the valence of each atom.

The MHG is defined over these hypergraphs as a hyperedge replacement grammar (HRG) [14], a context free grammar generating a hypergraph by replacement of hyperedges with other hypergraphs. This approach has a number of desirable properties, such as preserving the number of hypernodes belonging to each hyperedge, thereby preserving an atom’s valence – satisfying (i) above. Furthermore, stereochemistry can be encoded directly into the hyperedge replacement rules. MHG is thus guaranteed to produce only chemically valid molecules, allowing our model below to focus on optimizing the target benchmarks, without wasting network capacity on learning to generate valid outputs.

2.1 Parsing algorithm

In order to build the grammar which will be used by the RL agent outlined in 3, we use the approach of [13] to construct a parse tree for each molecule in a training set. By doing so for each molecule, we can identify the set of unique hyperedge replacement rules defining the grammar. Given an input molecular graph, the algorithm to deduce the MHG rules is as follows.

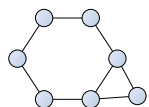
1. Find a node, n , at which the graph can be subdivided by cutting a single connected hyperedge
2. Split the graph at n producing a now reduced parent graph and a new child graph containing a new non-terminal node on the cut hyperedge, which we call the *parent node*
3. Iterate steps above until the parent graph can no longer be subdivided in this manner
4. Apply all of the above recursively to the new child graphs

After this, we have a parse tree where the original graph can be reconstructed by replacing the correct node in the parent graphs with children at their corresponding parent nodes. A tree constructed in this manner will contain only graphs containing a single hyperedge or graphs composed of cycles.

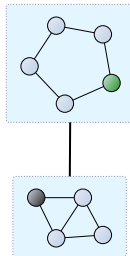
2.1.1 Extraction of Hypergraph Cliques

Though the above procedure will produce a MHG able to represent any molecule, it is not the most efficient representation given the combinatorial complexity in the construction of molecular graph cycles leading to an unnecessarily large set of rules. We now continue to reduce the cycle containing graphs by indentifying any cliques of the remaining graph of length greater than or equal to three, but which do not contain the parent node. This is achieved by replacing the entire clique with a new non-terminal graph node in the parent. The child consists of the clique and all edges exiting the clique connected to the new non-terminal parent node.

After removing all such cliques, we continue to remove 2-cliques from any remaining cycles of length greater than five. This greatly reduces the complexity of remaining cycles, and so the number of resultant rules, but discourages the production of triangles and boxes which are undesirable in output molecules from the model.



(a) Molecular hypergraph before removal of the 3-clique.



(b) Parent and child sections of parse tree after removing the 3-clique.

Figure 2: Molecular hypergraph clique extraction step. Black node indicates the parent node that is merged into the green parent node when the child rule is applied.

2.2 Using the grammar to create new molecules

After following the procedures outlined above, we can collect all unique graphs, taking into account also the position of the parent node. This set of unique graphs are the rules of the HOG. In the sense of a CFG, a rule maps a non-terminal node to a hypergraph containing further terminal and non-terminal nodes. To generate a new molecule we first select a rule which by definition will contain one or more non-terminal nodes. Next, for each non-terminal node in this first graph, we select a rule whose parent node will expand the non-terminal in the parent rule. We recursively apply the same procedure now over the selected child rules until all non-terminals are replaced with terminals. This is guaranteed to produce a valid molecular hypergraph.

2.2.1 Conditional and unconditional rule frequencies

To assist the model in selecting these rules, we collect the unconditional and conditional rule frequencies as they appear in the inference data set after having inferred the final set of grammar rules. Unconditional frequencies are simply the total count of appearances of each unique rule in the inference set parse trees, whereas conditional frequencies count the occurrence of each rule conditional on the given parent it will expand in the parse tree. These are used as priors for rule selection by the model before learning.

2.2.2 Ensuring expansions terminate

The final challenge we address is making sure the rule expansion terminates before the maximum allowed number of expansion rules has been reached by the agent. To do this, we define the concept of *terminal distance* of a hypergraph node, defined as the length of the shortest sequence of rules needed to expand said node into a sub-hypergraph consisting only of terminal nodes.

We calculate this distance for each hypergraph in our set of grammar rules by means of the following algorithm:

1. Define a set R of all rule hypergraphs observed so far, and seed it with the root token molecule.
2. Initialize all parent node terminal distances to ∞ .
3. Iterate over all parent nodes n_p of rules r in R , defining R_{n_p} as the set of rules with equivalent parent node n_p .
 - (a) For each n_p , the terminal distance is defined as one plus the minimum of terminal distances for all possible child rules of this graph. Defining $C(r)$ as the set of child nodes of the rule r ,

$$TD_n(n_p) = 1 + \min_{r \in R_{n_p}} \left(\sum_{c \in C(r)} TD_r(c) \right), \quad (1)$$

where TD_n, TD_r are the node and rule terminal distances respectively. If the child rule, c , is terminal it is assigned terminal distance $TD_r(c) = 0$.

We repeat step 3. until convergence, that is until the computed terminal distances do not change between iterations. This is most efficiently implemented as a dynamic programming problem since there are many overlapping sub-problems.

Finally, we define the terminal distance of a rule hypergraph as the sum of terminal distances of the child nodes,

$$\text{TD}_r(r) = \sum_{c \in C(r)} \text{TD}_r(c). \quad (2)$$

We use the terminal distance concept to make sure the rule expansion terminates before the maximum number of steps, in the following manner: at each step, we consider the number s of steps left and the terminal distance, td , of the sequence generated so far.

At each step, we make sure $\text{td} \leq s$, using induction. First, we choose the maximum rule sequence length to be larger than the terminal distance of the root graph. Second, at each rule selection step we consider all child rules, r , who can expand the next nonterminal in the graph, and only allow those where $\Delta\text{TD}_r(r) + \text{td} \leq s - 1$. This must be a nonempty set because $\Delta\text{TD}_r(r') = -1$ for at least one applicable rule r' . Thus, by the time we run out of steps, that is, $s = 0$, we know $\text{td} = 0$, that is our token sequence consists only of nonterminals.

2.3 Grammar conciseness and expressiveness

When constructing a grammar, there is a trade-off between the number of rules taken to express a given molecule and the total number of unique rules in the grammar. This tension is caused by the complexity of the individual rules. If we allow far more complex substructures to remain in the final rule-set then molecules containing such substructure are able to be expressed succinctly. However, given that these substructures can be arbitrarily complex we can swiftly see an explosion of the number of rules remaining in the grammar.

3 Model choice

3.0.1 Graph embedding

In a model that generates a discrete object via a series of decisions, we need to choose how to represent the intermediate state to feed into making the next choice. In models that produce sequences, such as SMILES strings, the natural choice is a sequence of one-hot embedded string characters. In models that construct the molecular graph directly, we have more options: we could take as our representation the sequence of the grammar rules chosen so far, like in [9, 16], or a representation of the intermediate graph state, as in [20].

We choose a simple encoding of three concatenated parts, with a sequence of length equaling the number of nodes in the graph. That embedding is composed of the concatenation of the connectivity matrix of the graph (with values 0, 1, 2, etc. denoting no connection, single bond, double bond, etc.), the identity matrix (so that the vector knows which node it refers to, as the ordering of the sequence itself is arbitrary), and a one-hot encoding of any data on each node, such as atomic number and chirality. This is then linearly transformed to the dimension of the downstream model.

3.0.2 Network architecture

We choose the Transformer architecture [17] rather than the more often used RNNs, because the nodes in the graph are not naturally ordered, and as we don’t add the sinusoidal position encoding from the original Transformer model to the inputs, the outputs are invariant to the ordering of the input sequence. The Transformer was used for graph inputs by [24], **TODO(WERE YOU THINKING OF THIS? They use a transformer based model for graph-based NLP parsing)** but never to our knowledge in the context of molecular optimization. To assess the usefulness of this approach, we provide ablation tests using an RNN architecture. We use model dimension 512, 6 heads, head size 64, and 5 layers.

3.0.3 Policy network and Discriminator

We implement both a policy model and for some runs a discriminator model using the above architecture. The policy model has two heads, one that supplies the logits for choosing the next nonterminal node to expand, and one that supplies the logits for the next rule to use on that node.

The discriminator network has a single head that tries to calculate the probability of whether a given molecule comes from the training dataset or has been created by the model.

3.1 Training

Most of the literature on molecular optimization either trains a VAE and then optimizes over the latent space, or uses a reinforcement learning approach; most of the reinforcement learning approaches don’t exchange any information between the results of a simulation batch, with the exception of [25, 26] and [16], who simulate the whole model and then take the k results with the best reward, and rewards the log-likelihood of those decision paths; which can be regarded as a basic kind of Monte Carlo RL. That is motivated by the fact that in contrast to many other RL use cases, we don’t know the reward until the end of the episode, and so it’s worth trying to reward the more successful paths as a whole after the fact.

3.1.1 Non-originality penalty

To avoid the optimization getting stuck in a local minimum, we penalize repeated occurrences of a molecule. Denoting the batch size used for generator training as b , we keep buffers of b , $10 \cdot b$, and $100 \cdot b$ last unique molecules, and decrease the rewards of new molecules if these are found in the training dataset or one of these buffers, with the exact formulation found in Appendix ??

3.1.2 Batch Advantage

We extend the approach of [25] by introducing the concept of batch advantage, defined as the reward for a particular molecule in the batch minus the average reward for the batch. This can be regarded as a special case of the advantage concept [27] **TODO(CHECK THIS REFERENCE - WASN’T CLEAR WHERE ADVANTAGE COMES FROM. Probably also want to add more recent developments like the async learning paper)** being applied to the molecular optimization, regarding the whole sequence of decisions going into a particular molecule as one ‘action’. The batch reward average can then be regarded as an estimate of the state value of the initial state, and the rewards for the batch members as their exact action values. To keep the learning rate uniform, we normalize these advantage values so that the sum of their absolute values equals 1, and use them as weights for the log likelihood of the respective molecule.

Note that batch advantage is especially effective in encouraging exploration when used together with non-originality penalty. When these two are combined, if a batch returns many molecules with very similar rewards (before applying the non-originality penalty), some of them new, and some repeated, then after applying the penalty and batch advantage, the repeated molecules’ log likelihood will actually have a negative weight, so we will discourage the model from reproducing these and encourage it to learn new ones.

3.1.3 Record rewards

As a final way to encourage exploration, we keep a buffer of the 10 best rewards observed in the current simulation. In the goal-optimization benchmarks, if a new reward is bigger than the smallest of these, we add the respective log likelihood to the objective function for that iteration, to especially encourage such molecules

3.1.4 Training step

A training step for the generator consists of simply computing the objective function as above and doing one step of the Adam optimizer with **SETTINGS**.

For the discriminator, the key thing was to balance a quick reaction to innovations from the generator with long-term memory. To achieve that, we use the above mentioned buffers of size b , $10 \cdot b$, and $100 \cdot b$ last unique molecules, let’s call them buf_1 , buf_{10} , and buf_{100} . Then each batch of b_g

molecules fed to the generator consists of $0.5b_g$ molecules from the training set, $0.25b_g$ molecules from buf_1 , $0.125b_g$ molecules from buf_{10} , and $0.125b_g$ molecules from buf_{100} .

4 Results

4.1 GuacaMol Benchmarks

The GuacaMol framework formulates three sets of benchmarks. The distributional benchmarks measure a model’s ability to sample from the ‘same’ distribution as that of the training set. The goal-oriented benchmarks ...

4.2 Ablation Studies

best-of-batch vs batch advantage rnn vs transformer

5 Discussion

We are the first to construct and use a probabilistic grammar for molecular optimization. That allows us to guarantee that all molecules generated by the model are chemically valid by construction, without the need for runtime checks.

When we decompose the training dataset using that grammar, we use frequencies of particular rule combinations to calculate conditional priors for a given rule being used to expand a given nonterminal in another rule. Sampling using these probabilities is a good first approximation of the training set distribution, as demonstrated by the ?? GuacaMol distributional benchmarks.

Having reasonable priors makes it much easier for a model to proceed to objective optimization. That allowed us to achieve competitive results without pre-training on all the ‘trivial’ goal-oriented GuacaMol benchmarks, in a very small number of steps; and ??? on their distributional benchmarks.

We are also the first to apply the Transformer architecture to graph embeddings for the purpose of molecular optimization, and ?? showed them to have superior performance to RNNs in that context. Comparison with the graph convolutional networks will be the subject of further work.

Further, we introduced the concept of batch advantage, and showed that it encouraged exploration and sped up convergence, especially when used together with a penalty for repetition.

?? is our best result on non-trivial benchmarks at least better than the max over dataset???

Acknowledgments

References

References

- [1] W. P. Walters, “Virtual chemical libraries,” *Journal of Medicinal Chemistry*, vol. 62, no. 3, pp. 1116–1124, 2019. PMID: 30148631.
- [2] L. Ruddigkeit, R. van Deursen, L. C. Blum, and J.-L. Reymond, “Enumeration of 166 billion organic small molecules in the chemical universe database gdb-17,” *Journal of Chemical Information and Modeling*, vol. 52, no. 11, pp. 2864–2875, 2012. PMID: 23088335.
- [3] D. K. Duvenaud, D. Maclaurin, J. Iparraguirre, R. Bombarell, T. Hirzel, A. Aspuru-Guzik, and R. P. Adams, “Convolutional networks on graphs for learning molecular fingerprints,” in *Advances in Neural Information Processing Systems 28* (C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama, and R. Garnett, eds.), pp. 2224–2232, Curran Associates, Inc., 2015.
- [4] S. Kearnes, K. McCloskey, M. Berndl, V. Pande, and P. Riley, “Molecular graph convolutions: moving beyond fingerprints,” *Journal of Computer-Aided Molecular Design*, vol. 30, pp. 595–608, Aug. 2016.
- [5] J. Gilmer, S. S. Schoenholz, P. F. Riley, O. Vinyals, and G. E. Dahl, “Neural message passing for quantum chemistry,” *CoRR*, vol. abs/1704.01212, 2017.
- [6] H. Dai, Y. Tian, B. Dai, S. Skiena, and L. Song, “Syntax-directed variational autoencoder for structured data,” *CoRR*, vol. abs/1802.08786, 2018.

- [7] W. Jin, R. Barzilay, and T. S. Jaakkola, "Junction tree variational autoencoder for molecular graph generation," *CoRR*, vol. abs/1802.04364, 2018.
- [8] M. Simonovsky and N. Komodakis, "Graphvae: Towards generation of small graphs using variational autoencoders," *CoRR*, vol. abs/1802.03480, 2018.
- [9] M. J. Kusner, B. Paige, and J. M. Hernández-Lobato, "Grammar Variational Autoencoder," *ArXiv e-prints*, Mar. 2017.
- [10] R. Gómez-Bombarelli, D. K. Duvenaud, J. M. Hernández-Lobato, J. Aguilera-Iparraguirre, T. D. Hirzel, R. P. Adams, and A. Aspuru-Guzik, "Automatic chemical design using a data-driven continuous representation of molecules," *CoRR*, vol. abs/1610.02415, 2016.
- [11] G. Lima Guimaraes, B. Sanchez-Lengeling, C. Outeiral, P. L. Cunha Farias, and A. Aspuru-Guzik, "Objective-Reinforced Generative Adversarial Networks (ORGAN) for Sequence Generation Models," *arXiv e-prints*, p. arXiv:1705.10843, May 2017.
- [12] D. Weininger, "Smiles, a chemical language and information system. 1. introduction to methodology and encoding rules," *J. Chem. Inf. Comput. Sci.*, vol. 28, no. 1, pp. 31–36, 1988.
- [13] H. Kajino, "Molecular hypergraph grammar with its application to molecular optimization," *CoRR*, vol. abs/1809.02745, 2018.
- [14] F. Drewes, H.-J. Kreowski, and A. Habel, "Hyperedge replacement graph grammars," in *Handbook of Graph Grammars and Computing by Graph Transformation* (G. Rozenberg, ed.), pp. 95–162, River Edge, NJ, USA: World Scientific Publishing Co., Inc., 1997.
- [15] P. Pogány, N. Arad, S. Genway, and S. D. Pickett, "De novo molecule design by translating from reduced graphs to smiles," *Journal of Chemical Information and Modeling*, vol. 59, no. 3, pp. 1136–1146, 2019.
- [16] E. Kraev, "Grammars and reinforcement learning for molecule optimization," *CoRR*, vol. abs/1811.11222, 2018.
- [17] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention is all you need," *CoRR*, vol. abs/1706.03762, 2017.
- [18] X. Yang, J. Zhang, K. Yoshizoe, K. Terayama, and K. Tsuda, "Chemts: an efficient python library for de novo molecular generation," *Science and Technology of Advanced Materials*, vol. 18, no. 1, pp. 972–976, 2017. PMID: 29435094.
- [19] M. Olivecrona, T. Blaschke, O. Engkvist, and H. Chen, "Molecular de novo design through deep reinforcement learning," *Journal of Cheminformatics*, vol. 9, 04 2017.
- [20] J. You, B. Liu, R. Ying, V. Pande, and J. Leskovec, "Graph convolutional policy network for goal-directed molecular graph generation," in *Proceedings of the 32nd International Conference on Neural Information Processing Systems, NIPS'18, (USA)*, pp. 6412–6422, Curran Associates Inc., 2018.
- [21] M. Popova, O. Isayev, and A. Tropsha, "Deep reinforcement learning for de novo drug design," *Science Advances*, vol. 4, no. 7, 2018.
- [22] Sanchez-Lengeling, Benjamin, Outeiral, Carlos, Guimaraes, G. L., Aspuru-Guzik, and Alan, "Optimizing distributions over molecular space. an objective-reinforced generative adversarial network for inverse-design chemistry (organic)," Aug 2017.
- [23] E. Putin, A. Asadulaev, Y. Ivanenkov, V. Aladinskiy, B. Sanchez-Lengeling, A. Aspuru-Guzik, and A. Zhavoronkov, "Reinforced adversarial neural computer for de novo molecular design," *Journal of Chemical Information and Modeling*, vol. 58, no. 6, pp. 1194–1204, 2018. PMID: 29762023.
- [24] D. Kondratyuk, "75 languages, 1 model: Parsing universal dependencies universally," *CoRR*, vol. abs/1904.02099, 2019.
- [25] M. H. S. Segler, T. Kogej, C. Tyrchan, and M. P. Waller, "Generating focussed molecule libraries for drug discovery with recurrent neural networks," *CoRR*, vol. abs/1701.01329, 2017.
- [26] D. Neil, M. Segler, L. Guasch, M. Ahmed, D. Plumbley, M. Sellwood, and N. Brown, "Exploring deep recurrent models with reinforcement learning for molecule design," 2018.
- [27] L. C. Baird, "Reinforcement learning in continuous time: advantage updating," 1994.