

 Open access • Posted Content • DOI:10.33774/CHEMRXIV-2021-ZV6F1-V2

Perplexity-based molecule ranking and bias estimation of chemical language models

— [Source link](#) 

Michael Moret, Francesca Grisoni, Paul Katzberger, Gisbert Schneider





Institutions: ETH Zurich, Eindhoven University of Technology

Published on: 27 Oct 2021 - ChemRxiv

Topics: Perplexity and Ranking

Related papers:

- [Better Language Models with Model Merging](#)
- [Active relevance feedback for difficult queries](#)
- [Feature ranking for multi-label classification using Markov Networks](#)
- [A Structured Prediction Approach for Label Ranking](#)
- [Boosted Generative Models](#)

Share this paper:    

View more about this paper here: <https://typeset.io/papers/perplexity-based-molecule-ranking-and-bias-estimation-of-3n4ggahrka>

Perplexity-based molecule ranking and bias estimation of chemical language models

Michael Moret,^{†,1} Francesca Grisoni,^{†,2*} Paul Katzberger,¹ Gisbert Schneider^{1,3*}

¹ETH Zurich, Department of Chemistry and Applied Biosciences, RETHINK, Vladimir-Prelog-Weg 4, 8093 Zurich, Switzerland;

²Eindhoven University of Technology, Institute for Complex Molecular Systems, Department of Biomedical Engineering, Groene Loper 7, 5612AZ Eindhoven, Netherlands;

³ETH Singapore SEC Ltd, 1 CREATE Way, #06-01 CREATE Tower, Singapore 138602, Singapore;

*Correspondence to Gisbert Schneider (gisbert@ethz.ch) and Francesca Grisoni (f.grisoni@tue.nl).

[†] M.M. and F.G. contributed equally to this work.

Abstract

Chemical language models (CLMs) can be employed to design molecules with desired properties. CLMs generate new chemical structures in the form of textual representations, such as the simplified molecular input line entry systems (SMILES) strings, in a rule-free manner. However, the quality of these de novo generated molecules is difficult to assess a priori. In this study, we apply the perplexity metric to determine the degree to which the molecules generated by a CLM match the desired design objectives. This model-intrinsic score allows identifying and ranking the most promising molecular designs based on the probabilities learned by the CLM. Using perplexity to compare “greedy” (beam search) with “explorative” (multinomial sampling) methods for SMILES generation, certain advantages of multinomial sampling become apparent. Additionally, perplexity scoring is performed to identify undesired model biases introduced during model training and allows the development of a new ranking system to remove those undesired biases.

Introduction

Generative deep learning has become a promising method for chemistry and drug discovery^{1-17,19,21}. Generative models learn the pattern distribution of the input data and generate new data instances based on learned probabilities¹². Among the proposed generative frameworks that have been applied to de novo molecular design^{1-4,6-17,19,21}, chemical language models (CLMs) have gained particular attention because of their ability to generate focused virtual chemical libraries

and bioactive compounds^{18,19,21,22}. CLMs are trained on string representations of molecules, *e.g.*, simplified molecular input line entry systems (SMILES) strings (Fig. 1a)²⁰. CLMs for molecular design iteratively predict the next SMILES character using all the preceding portions of the SMILES string (Fig. 1b). In this process, CLMs learn the conditional probability of sampling any SMILES character based on the preceding characters. After CLM training, the model can be used for molecular construction. CLMs have been demonstrated to both learn the SMILES syntax and capture semantic features of the training molecules, such as physicochemical properties^{19,21–23}, bioactivity, and chemical synthesizability^{1,2,21}. A desirable feature of CLMs is their ability to function in low-data regimes^{22,24}, *i.e.*, with limited training data (typically in the range of 5 to 40 molecules)^{1,2,21}. One of the most widely employed approaches for low-data model training is transfer learning^{19,25}. This method leverages previously acquired information on a related task for which more data are available ("pretraining") before training the CLM on a more specific limited dataset ("fine-tuning")²⁶.

Many prospective *de novo* design studies based on CLMs used weighted random sampling ("multinomial sampling", often in the form of temperature sampling) for molecule generation^{1,19,22,27}. This method samples the most likely SMILES string characters more frequently than the unlikely characters. This feature enables (i) extensive virtual molecule libraries to be generated and (ii) a certain chemical space to be investigated, owing to "fuzzy" (probability-weighted) random sampling. However, such a sampling strategy can result in molecules that do not possess the physicochemical and biological properties of the training data. Furthermore, because the number of molecules that can potentially be sampled from CLMs significantly exceeds synthetic capacities, and a natural ranking of the generated SMILES does not exist, an additional procedure is required for molecule prioritization, *e.g.*, that based on similarity assessment or activity prediction^{1,28}. We recently introduced the beam search algorithm as an alternative to multinomial sampling. During beam search, the most likely SMILES strings are generated based on the respective character probabilities, thereby alleviating the strict requirement for additional molecule prioritization²⁹. The beam search performs chemical space "exploitation" as the algorithm searches for the most probable SMILES strings in a greedy manner. This method generates only a few candidate molecules at the expense of chemical space exploration and design diversity.

Herein, we used perplexity to assess the "goodness" of the designs generated by CLMs via multinomial sampling aiming to (i) preserve the advantage of intrinsic molecule ranking, as it can be achieved via beam search²⁹, (ii) benefit from the chemical space exploration provided by multinomial sampling^{22,29}, and (iii) provide model-based insights into the most promising molecules for follow-up analysis²⁹. Moreover, the perplexity metric identified undesired effects of transfer learning using CLMs, thereby qualifying as a criterion for detecting undesired model biases.

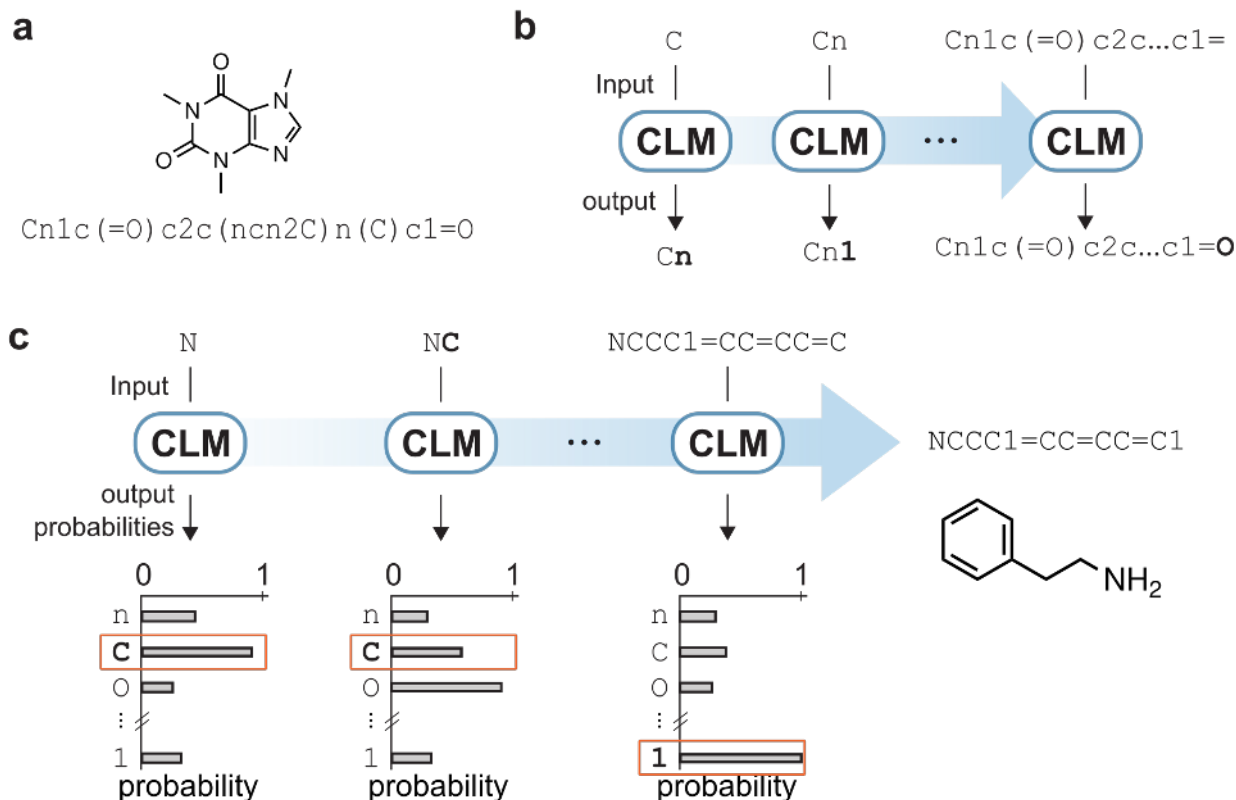


Fig. 1 | Principles of chemical language models (CLMs). **a**, Example of molecular structure (Kekulé structure) and corresponding SMILES string. **b**, CLMs are trained to iteratively predict next SMILES character based on preceding ones. **c**, Multinomial sampling used to generate new SMILES strings from trained CLMs, where SMILES characters are sampled with a weighted random sampling of probability distributions learned by CLM (i.e., multinomial sampling).

Results and Discussion

Perplexity-based scoring of molecular *de novo* designs

The perplexity metric was used as a scoring function to reflect the probability of sampling a SMILES string as a function of its characters (Eq. 1)³⁰. Perplexity is typically used to assess the performance of language models in natural language processing³⁰⁻³². For a SMILES string of length N , the perplexity score can be computed by considering the CLM probability of any i -th character (p_i), as follows:

$$\text{perplexity} = 2^{-\frac{1}{N} \sum_{i=1}^N \log(p_i)} \quad (1)$$

Perplexity allows to quantify the CLM confidence that a specific SMILES string could have belonged to the training data. Because the training objectives of the CLM are implicitly encoded in the fine-tuning data, the perplexity score allows one to assess whether the generated SMILES strings match the objectives. A SMILES string composed of probable characters (high p_i values)

exhibits low perplexity, whereas a string containing many unlikely characters (low p_i values) exhibits high perplexity. Hence, low perplexity scores are desirable.

Perplexity as an indicator for comparing molecular sampling strategies

In a previous study, new bioactive compounds were successfully identified via beam search sampling²⁹, which is a heuristic greedy algorithm. Its search "breadth" is controlled by a width parameter (k), which represents the number of the most probable SMILES strings that the model considers during string extension. In this study, beam search was used as a reference method for comparison with multinomial sampling.

We investigated the difference in perplexity scores between molecules generated using a CLM via either beam search or multinomial sampling. In this regard, a CLM was pretrained with approximately 1.6 million molecules from ChEMBL (version 28)³³. Ten randomly selected targets were used for fine-tuning (Table 1). For each macromolecular target, ligands that possessed a pChEMBL value larger than 6 were selected, where pChEMBL is defined as $-\log_{10}(\text{molar IC}_{50}, \text{XC}_{50}, \text{EC}_{50}, \text{AC}_{50}, \text{K}_i, \text{K}_d, \text{ or potency})$. To emulate different low-data regimes typical of drug discovery, we prepared fine-tuning sets of different sizes that contained 5, 10, 20, or 40 randomly selected ligands for each target. For each of the 10 targets and each of the four fine-tuning sets, 10 or 50 SMILES strings were generated via beam search sampling ($k = 10$ or $k = 50$). A total of 1,000 SMILES strings were sampled every second epoch during 100 CLM fine-tuning epochs via multinomial sampling.

Table 1 | Macromolecular targets selected for CLM fine-tuning. ChEMBL target identifier, generic target name, and protein classification specified based on respective ChEMBL Target Report Card.

CHEMBL ID	Target	Protein classification
CHEMBL1836	Prostanoid EP4 receptor	G protein-coupled receptor
CHEMBL1945	Melatonin receptor 1A	G protein-coupled receptor
CHEMBL1983	Serotonin 1D (5-HT _{1D}) receptor	Family A G protein-coupled receptor
CHEMBL202	Dihydrofolate reductase	Oxidoreductase
CHEMBL3522	Cytochrome P450 17A1	Cytochrome P450
CHEMBL4029	Interleukin-8 receptor A	Family A G protein-coupled receptor
CHEMBL5073	CaM kinase I delta	Kinase
CHEMBL5137	Metabotropic glutamate receptor 2	G protein-coupled receptor
CHEMBL5408	Serine/threonine-protein kinase TBK1	Kinase
CHEMBL5608	NT-3 growth factor receptor	Kinase

In CLM fine-tuning using the smallest fine-tuning sets (five molecules), multinomial sampling consistently outperformed beam search sampling in terms of the perplexity score, as molecules with the best score (lowest perplexity) were obtained (Fig. 2a). Increasing the beam search width from $k = 10$ to $k = 50$ did not markedly improve the ability of this method in identifying molecules with higher perplexity scores (Fig. 2). These observations were confirmed for the larger fine-tuning sets (Supplementary Figs. S1–S3). A potential explanation for this observation is the “greedy” nature of the beam search³⁴, which explores only a limited number of possibilities for next-character addition. By contrast, the “fuzzy” nature of multinomial sampling allows the generation of a greater number of molecules, and hence a broader exploration of the chemical space of interest.

The 50 top-scoring molecules generated via multinomial sampling not only indicated lower median perplexity values but also spanned a narrower range of values (Fig. 2a). This suggests that multinomial sampling yields more high-scoring designs (low perplexity) for follow-up synthesis and biological testing than the beam search algorithm (Fig. 2a and Supplementary Fig. S1 to S3).

When filtering out designs with a fragment similarity (Tanimoto index on Morgan fingerprints³⁵) greater than 50% of the respective fine-tuning molecules, the difference between multinomial sampling and beam searching was less pronounced (Fig. 2b, Supplementary Figs. S3–S6). Multinomial sampling identified molecules with lower perplexity scores than the beam search in 72% of the cases involving the smallest fine-tuning sets (Supplementary Table S1). The deterioration in performance for highly diverse molecules was less pronounced with the larger fine-tuning sets (Supplementary Table S1).

The results of this study corroborate the potential of multinomial sampling, not only for chemical space exploration to obtain chemically diverse molecular designs, but also for generating high-scoring compounds that are sufficiently diverse from the fine-tuning compounds. Multinomial sampling allows a user-defined number of SMILES strings (potentially exceeding 1000 molecules) to be generated via a CLM. It is reasonable to assume that better scoring molecules (i.e., with lower perplexity) can be created by generating more molecular designs via multinomial sampling.

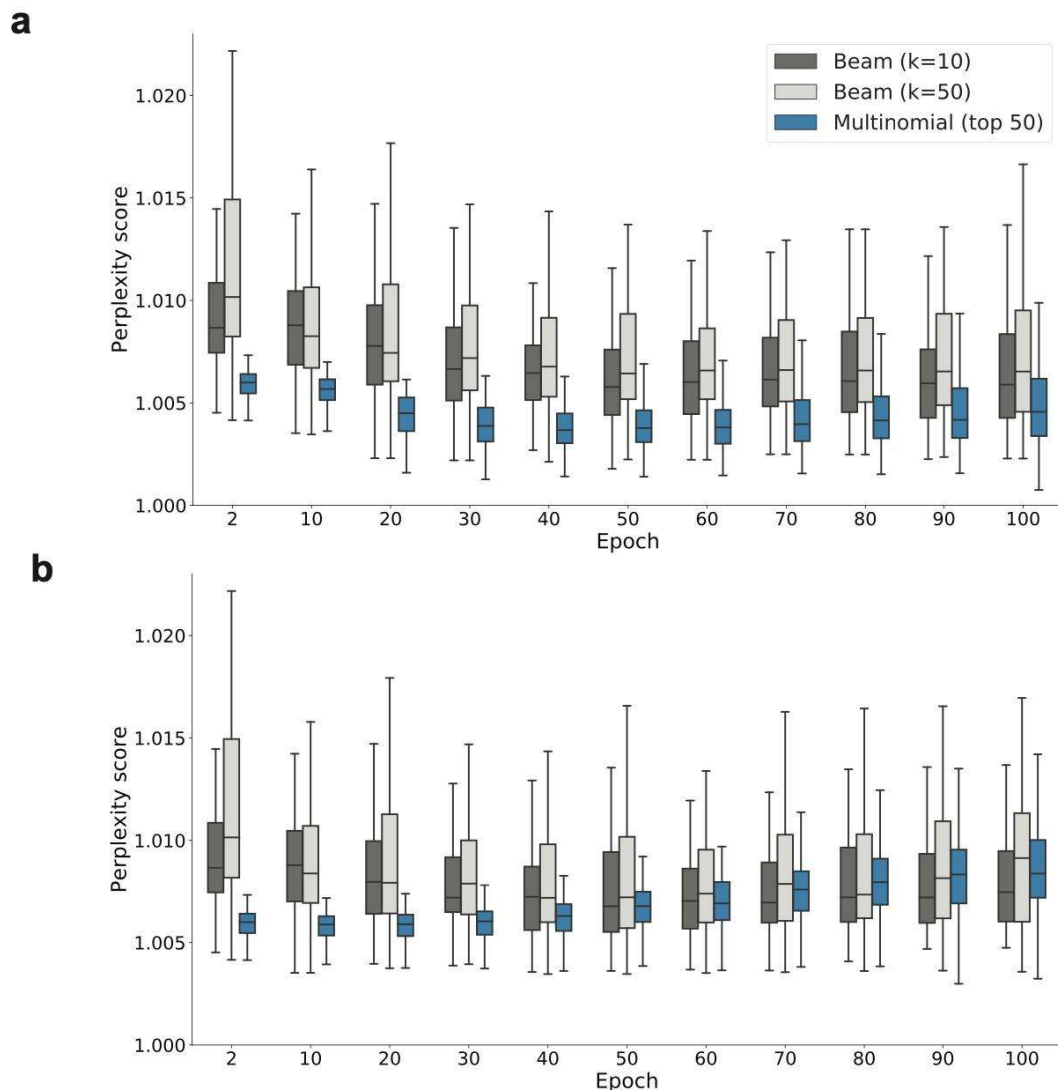


Fig. 2 | Variation in perplexity during fine-tuning. **a**, Distribution of top-scoring compounds for each method over 100 fine-tuning epochs (only every 10 epochs shown in graph for clarity). **b**, Distribution of top-scoring compounds by considering only molecules with similarity below 50% (Tanimoto index computed on Morgan fingerprints) to closest molecule in their respective fine-tuning set. Median and lower to upper quartile values reported using boxplots for 10 different protein-specific fine-tuning sets, which contain five molecules each. Boxplots for fine-tuning sets of size 10, 20 and 40 molecules provided as supporting information.

Assessing pretraining bias based on perplexity

CLM pretraining might impose a greater effect on model performance than CLM fine-tuning, as model pretraining is typically performed with data that are at least two orders of magnitude higher in amount than fine-tuning^{1,2,22,28,29}. If a molecule is generated by a CLM due to pretraining only, it will not necessarily match the design objectives (as represented by the fine-tuning data). We analyzed the degree to which new molecules were generated due to the sole effect of pretraining, i.e., we verified whether “pretraining bias” occurred. In principle, for a CLM, perplexity can be used

to score any molecule, including those that are not generated by the model. This can be achieved by computing the conditional probabilities of each SMILES character using the CLM. Therefore, the perplexity score was employed to differentiate between the information learned by the CLM during pretraining and during fine-tuning to score the molecules generated at a specified fine-tuning epoch. First, for each fine-tuning epoch, molecules were scored and ranked by the perplexity of the model used to generate them. Subsequently, each de novo design was scored and ranked based on the perplexity of the CLM after pretraining (i.e., prior to any fine-tuning). We hypothesized that a suitable ranking method should favor molecules that were generated based on information learned by the model during fine-tuning (capturing the final objectives of the experiment), and downrank the molecules generated based solely on pretraining (capturing “generic” information).

To seize this concept quantitatively, we subtracted the rank yielded by the pretrained model ($rank_{pt}$) from that of the fine-tuned CLM ($rank_{ft}$) for each molecule, and defined this difference as the “Delta” score (Eq. 2), as follows:

$$Delta = rank_{ft} - rank_{pt}. \quad (2)$$

Molecules with a positive delta score were considered more likely to be output by the fine-tuned CLM. A negative Delta score suggests that the fine-tuning procedure renders a certain molecule less likely to be output than after pretraining, which thus does not satisfy the design objectives.

The Delta score was computed for all sampled molecules during fine-tuning experiments (Fig. 3b). The percentage of molecules with a negative Delta exceeded 40% for the first 20 fine-tuning epochs and remained above 10% until the end of fine-tuning (Fig. 3c). This suggests that 10%–40% of the molecules were generated based on “generic” pretraining, instead of “task-focused” fine-tuning. As such, this highlights the practicability of the proposed Delta score as an indicator to (i) detect potential pretraining bias, (ii) identify the best-suited epoch for a productive sampling of molecules that fulfill the study goals, and (iii) select the most promising de novo designs.

To expand the analysis, we focused on the 50 top scoring molecules generated via multinomial sampling. We discovered that among them, only up to 3% of the molecules received a negative Delta score (Supplementary Fig. S7). This shows that using perplexity alone reduces the pretraining bias. However, the pretraining bias was not completely removed, which highlights the benefits of using both the perplexity and Delta score for molecule prioritization prior to synthesis and biological testing.

In summary, these results suggest that the potential of generative CLMs in medicinal chemistry can be expanded by employing the SMILES perplexity for molecule prioritization and for detecting potential pretraining bias.

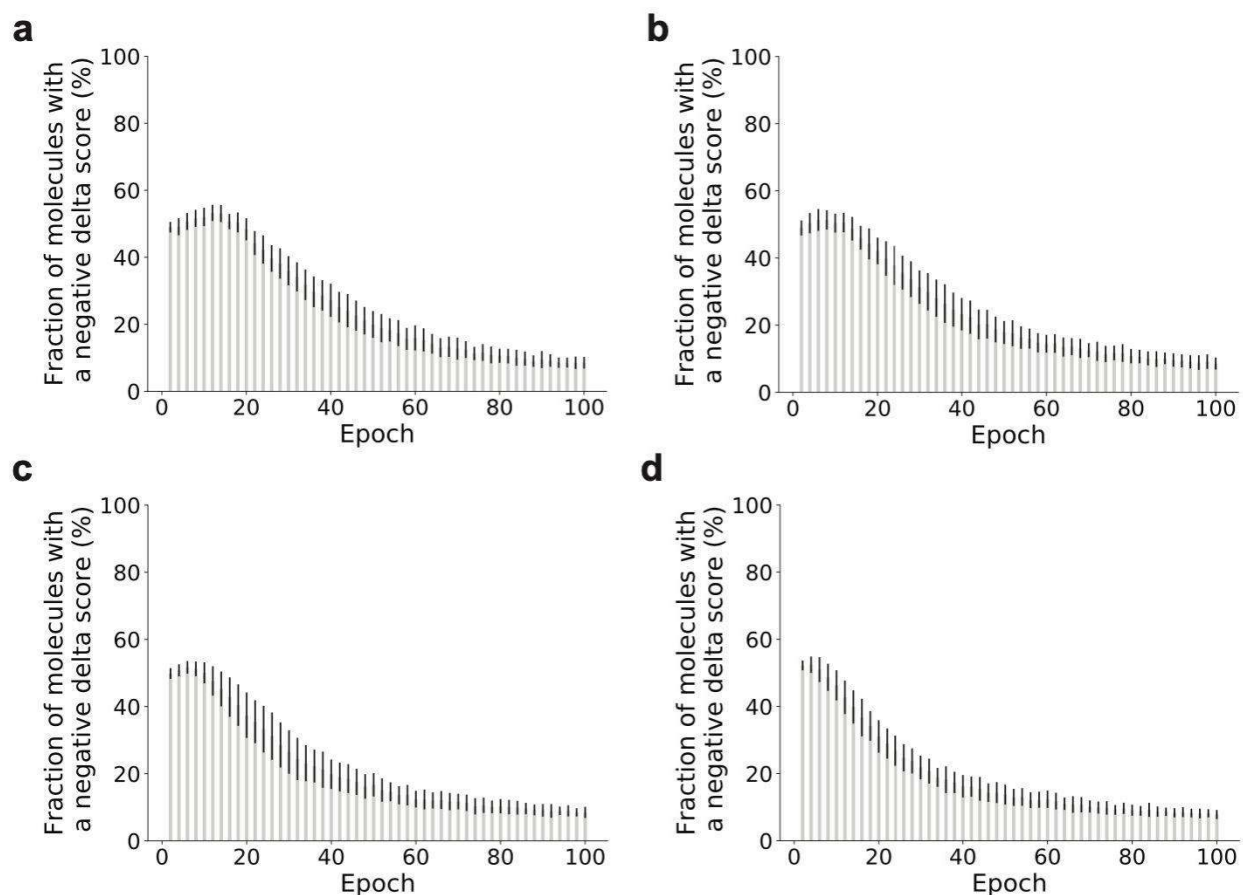


Fig. 3 | Delta score during fine-tuning experiments in low-data regime. Percentage of molecules with negative Delta score (1,000 sampled SMILES strings; mean \pm standard deviation reported for 10 different target proteins). Fine-tuning sets of **a**, 5, **b**, 10, **c**, 20, and **d**, 40 molecules.

Conclusions

This present study constitutes a step forward towards automated, self-supervised de novo design. By serving as a model-intrinsic score, perplexity enables the quality assessment of generated molecules. In particular, perplexity might be useful for identifying the most promising molecules, i.e., those that match the probability distribution of the training data as captured by the CLM. This approach enables the comparison of two different methods for SMILES sampling from a trained CLM. The results revealed certain advantages of multinomial sampling over the beam search method for molecule generation. Because perplexity can be used to score SMILES strings based on the information learned by a CLM, the pretraining bias can be identified based on the newly introduced Delta score. Perplexity combined with the Delta score can reveal the most promising molecules, in terms of the fine-tuning objectives, for synthesis and testing. These features can further accelerate drug discovery using CLMs. Future studies will focus on the combination of perplexity with the temperature parameter of multinomial sampling or SMILES augmentation^{27,36}. Furthermore, the combination of CLMs and perplexity scoring bears promise for screening large collections of commercially-available compounds to accelerate model validation³⁷. More

experiments should be performed to determine the effect of the new approach on molecular de novo designs involving CLMs.

Code and data availability

The computational framework presented herein, pretrained neural network weights, and data used for model training are available in a GitHub repository from URL: https://github.com/ETHmodlab/CLM_perplexity.

Methods

Data processing. Molecules were represented as canonical SMILES strings using the RDKit (2019.03.2). SMILES strings were standardized in Python (v 3.6.5) by removing salts and duplicates, and only SMILES strings with 20–90 characters were retained.

Pretraining set. The molecules were retrieved from ChEMBL28³⁸. After data processing, the pretraining dataset contained 1,683,181 molecules encoded as SMILES strings. This set was further segregated randomly into a training set (1,599,021 molecules) and a validation set (84,160 molecules).

Fine-tuning sets. Target selection was limited to molecules satisfying the following conditions (ChEMBL annotation): (i) Organism: *Homo sapiens*; (ii) protein classification (L1): enzymes, membrane receptors, transcription factors, and single proteins; (iii) number of compounds: 962–2057 molecules (range defined by ChEMBL); (iv) number of activities: at least 2,000 reported pChEMBL values. Ten target proteins were randomly selected from the list of targets. For each of the 10 selected target proteins, sets of 5, 10, 20, and 40 molecules with pChEMBL > 6 were compiled randomly.

CLM implementation and training. All computational experiments were implemented in Python (v3.6.5) using Keras (v2.2.0, <https://keras.io/>) with the Tensorflow GPU backend (v1.9.0, <https://www.tensorflow.org/>). CLMs were implemented using a recurrent neural network with long short-term memory cells (LSTM)³⁹. The network, which was composed of four layers comprising 5,820,515 parameters (layer 1: batch normalization; layer 2: LSTM with 1024 units; layer 3: LSTM with 256 units; layer 4: batch normalization), was trained with SMILES strings encoded as one-hot vectors. We used the Adam optimizer with a learning rate of 10^{-4} for the CLM training during 90 epochs⁴⁰, where one epoch was defined as one pass over all the training data. Fine-tuning was performed by further training the CLM on the fine-tuning set for 100 epochs.

Multinomial sampling. Multinomial sampling was performed based on the CLM output for each SMILES string character. In particular, the probability of each i -th character to be sampled (p_i) was computed using Eq. 3):

$$p_i = \frac{e^{z_i}}{\sum_j e^{z_j}}, \quad (3)$$

where z_i is the CLM output for the i -th character before the *softmax* function was applied.

Beam search sampling. We use the implementation provided in Ref. 29 (https://github.com/ETHmodlab/molecular_design_with_beam_search) using two different beam widths ($k = 10$ and $k = 50$).

Acknowledgments

This study was financially supported by the Swiss National Science Foundation (grant no. 205321_182176 to G.S.) and by the RETHINK initiative at ETH Zurich.

Author contributions

F.G., M.M., and G.S. conceived the study. M.M. and F.G. designed the methodology, and M.M. implemented the software. All authors analyzed the results and contributed to the writing of the manuscript.

Competing interests

G.S. declares a potential financial conflict of interest as he is a consultant in the pharmaceutical industry and the co-founder of inSili.com GmbH, Zurich, Switzerland. No other potential conflicts of interest are declared.

References

1. Merk, D., Friedrich, L., Grisoni, F. & Schneider, G. De novo design of bioactive small molecules by artificial intelligence. *Mol. Inf.* **37**, 1700153 (2018).
2. Grisoni, F. *et al.* Combining generative artificial intelligence and on-chip synthesis for de novo drug design. *Sci. Adv.* **7**, eabg3338 (2021).
3. Chen, H., Engkvist, O., Wang, Y., Olivecrona, M. & Blaschke, T. The rise of deep learning in drug discovery. *Drug Discov. Today* **23**, 1241–1250 (2018).
4. Tang, B., Ewalt, J. & Ng, H.-L. Generative AI models for drug discovery. in *Topics in Medicinal Chemistry* 1–23 (Springer Berlin Heidelberg, 2021).
5. Schmidhuber, J. Deep learning in neural networks: an overview. *Neural Netw.* **61**, 85–117 (2015).
6. Putin, E. *et al.* Reinforced adversarial neural computer for de novo molecular design. *J. Chem. Inf. Model.* **58**, 1194–1204 (2018).
7. Olivecrona, M., Blaschke, T., Engkvist, O. & Chen, H. Molecular de-novo design through deep reinforcement learning. *J. Cheminform.* **9**, 48 (2017).
8. Born, J. *et al.* PaccMannRL: De novo generation of hit-like anticancer molecules from transcriptomic data via reinforcement learning. *iScience* **24**, 102269 (2021).
9. De Cao, N. & Kipf, T. MolGAN: An implicit generative model for small molecular graphs. *Preprint at <http://arxiv.org/abs/1805.11973>* (2018).
10. Zhavoronkov, A. *et al.* Deep learning enables rapid identification of potent DDR1 kinase inhibitors. *Nat. Biotechnol.* **37**, 1038–1040 (2019).
11. Sanchez-Lengeling, B. & Aspuru-Guzik, A. Inverse molecular design using machine learning: Generative models for matter engineering. *Science* **361**, 360–365 (2018).
12. Salakhutdinov, R. Learning deep generative models. *Annu. Rev. Stat. Appl.* **2**, 361–385 (2015).
13. Gómez-Bombarelli, R. *et al.* Automatic chemical design using a data-driven continuous representation of molecules. *ACS Cent. Sci.* **4**, 268–276 (2018).
14. Jin, W., Barzilay, R. & Jaakkola, T. Junction tree variational autoencoder for molecular graph generation. *Preprint at <http://arxiv.org/abs/1802.04364>* (2018).
15. Guimaraes, G. L., Sanchez-Lengeling, B., Outeiral, C., Farias, P. L. C. & Aspuru-Guzik, A. Objective-reinforced generative adversarial networks (ORGAN) for sequence generation models. *Preprint at <http://arxiv.org/abs/1705.10843>* (2017).
16. Skalic, M., Jiménez Luna, J., Sabbadin, D. & De Fabritiis, G. Shape-based generative modeling for de novo drug design. *J. Chem. Inf. Model.* **59**, 1205–1214 (2019).
17. Kang, S. & Cho, K. Conditional molecular design with deep generative models. *J. Chem. Inf. Model.* **59**, 43–52 (2019).
18. Walters, W. P. Virtual chemical libraries. *J. Med. Chem.* **62**, 1116–1124 (2019).
19. Segler, M. H. S., Kogej, T., Tyrchan, C. & Waller, M. P. Generating focused molecule libraries for drug discovery with recurrent neural networks. *ACS Cent. Sci.* **4**, 120–131 (2018).
20. Weininger, D. SMILES, a chemical language and information system. 1. Introduction to methodology and encoding rules. *J. Chem. Inf. Comput. Sci.* **28**, 31–36 (1988).
21. Yuan, W. *et al.* Chemical space mimicry for drug discovery. *J. Chem. Inf. Model.* **57**, 875–882 (2017).
22. Moret, M., Friedrich, L., Grisoni, F., Merk, D. & Schneider, G. Generative molecular design in low data regimes. *Nat. Mach. Intell.* **2**, 171–180 (2020).
23. Skinnider, M. *et al.* A deep generative model enables automated structure elucidation of novel psychoactive substances. *Preprint at https://chemrxiv.org/articles/preprint/A_Deep_Generative_Model_Enables_Automated_Structure_Elu_cidation_of_Novel_Psychoactive_Substances/14644854/1* (2021).
24. Skinnider, M. A., Greg Stacey, R., Wishart, D. S. & Foster, L. J. Chemical language models enable navigation in sparsely populated chemical space. *Nat. Mach. Intell.* **3**, 759–770 (2021).
25. Peters, M., Ruder, S. & Smith, N. A. To tune or not to tune? Adapting pretrained representations to diverse tasks. *Preprint at <http://arxiv.org/abs/1903.05987>* (2019).
26. Yosinski, J., Clune, J., Bengio, Y. & Lipson, H. How transferable are features in deep neural networks? *Preprint at <http://arxiv.org/abs/1411.1792>* (2014).
27. Gupta, A. *et al.* Generative recurrent networks for de novo drug design. *Mol. Inf.* **37**, 1700111 (2018).
28. Merk, D., Grisoni, F., Friedrich, L. & Schneider, G. Tuning artificial intelligence on the de novo design

- of natural-product-inspired retinoid X receptor modulators. *Commun. Chem.* **1**, 68 (2018).
29. Moret, M., Helmstädter, M., Grisoni, F., Schneider, G. & Merk, D. Beam search for automated design and scoring of novel ROR ligands with machine intelligence. *Angew. Chem. Int. Ed.* **60**, 19477–19482 (2021).
 30. Manning, C. & Schütze, H. *Foundations of statistical natural language processing*. (MIT Press, 1999).
 31. Radford, A. *et al.* Language models are unsupervised multitask learners. *OpenAI blog* **1**, 9 (2019).
 32. Hu, J., Gauthier, J., Qian, P., Wilcox, E. & Levy, R. P. A systematic assessment of syntactic generalization in neural language models. *Preprint at <http://arxiv.org/abs/2005.03692>* (2020).
 33. Gaulton, A. *et al.* The ChEMBL database in 2017. *Nucleic Acids Res.* **45**, D945–D954 (2017).
 34. Wilt, C. M., Thayer, J. T. & Ruml, W. A comparison of greedy search algorithms. in *Third Annual Symposium on Combinatorial Search* (2010).
 35. Rogers, D. & Hahn, M. Extended-connectivity fingerprints. *J. Chem. Inf. Model.* **50**, 742–754 (2010).
 36. Arús-Pous, J. *et al.* Randomized SMILES strings improve the quality of molecular generative models. *J. Cheminform.* **11**, 71 (2019).
 37. Moret, M., Grisoni, F., Brunner, C. & Schneider, G. Leveraging molecular structure and bioactivity with chemical language models for drug design. *Preprint at <https://chemrxiv.org/engage/chemrxiv/article-details/615580ced1fc334326f9356e>* (2021).
 38. Mendez, D. *et al.* ChEMBL: towards direct deposition of bioassay data. *Nucleic Acids Res.* **47**, D930–D940 (2019).
 39. Hochreiter, S. & Schmidhuber, J. Long short-term memory. *Neural Comput.* **9**, 1735–1780 (1997).
 40. Kingma, D. P. & Ba, J. Adam: a method for stochastic optimization. *Preprint at <http://arxiv.org/abs/1412.6980>* (2014).