# A Model to Search for Synthesizable Molecules

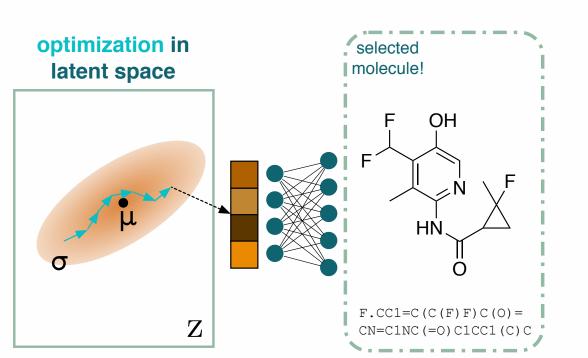
John Bradshaw<sup>1,2</sup>, Brooks Paige<sup>1,4</sup>, Matt J. Kusner<sup>3,4</sup>, Marwin H. S. Segler<sup>5,6</sup>, José Miguel Hernández-Lobato<sup>1,4,7</sup>

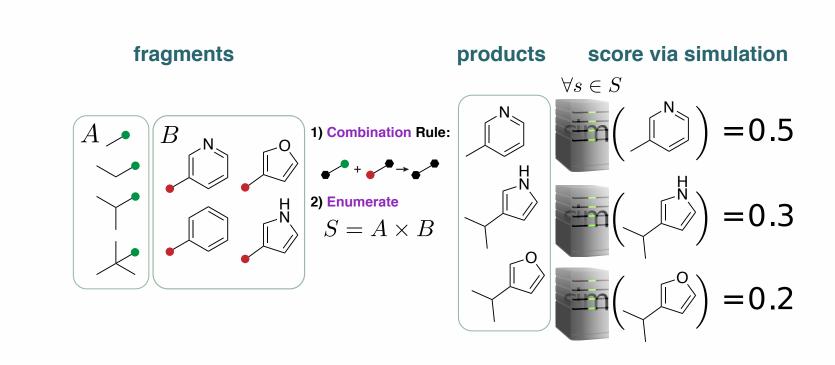
<sup>1</sup>University of Cambridge, <sup>2</sup>MPI for Intelligent Systems, <sup>3</sup>University College London, <sup>4</sup>The Alan Turing Institute, <sup>5</sup>BenevolentAI, <sup>6</sup>Westfälische Wilhelms-Universität Münster, <sup>7</sup>Microsoft Research Cambridge.

Aim: to design a generative model that enables the searching for useful molecules (eg for drugs) over its continuous latent space, whilst having a decoder which produces both stable chemical products and their synthetic routes.

#### We don't just want to know what molecule to make...

Virtual screening can be used to narrow down the number of candidates when searching for desirable drug molecules. However, this often consists of expensive combination and enumeration steps.



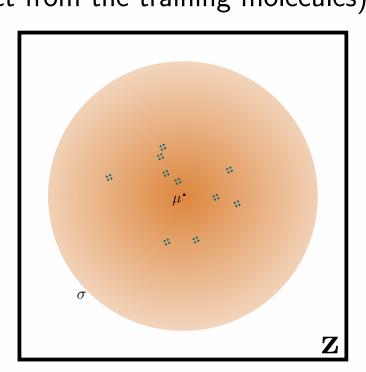


Recent machine learning approaches [1–3] have proposed learning an autoencoder to encode and decode molecules to a continuous latent space. The costly enumeration above can be instead replaced by search (eg local optimization).

But outstanding questions for these ML methods remain: (1) (how) are the proposed molecules synthesizable, and (2) how can we better imbue our models with inductive biases that result in semantically valid (chemically stable) molecules?

#### MOLECULE CHEF generates a wide range of stable molecules

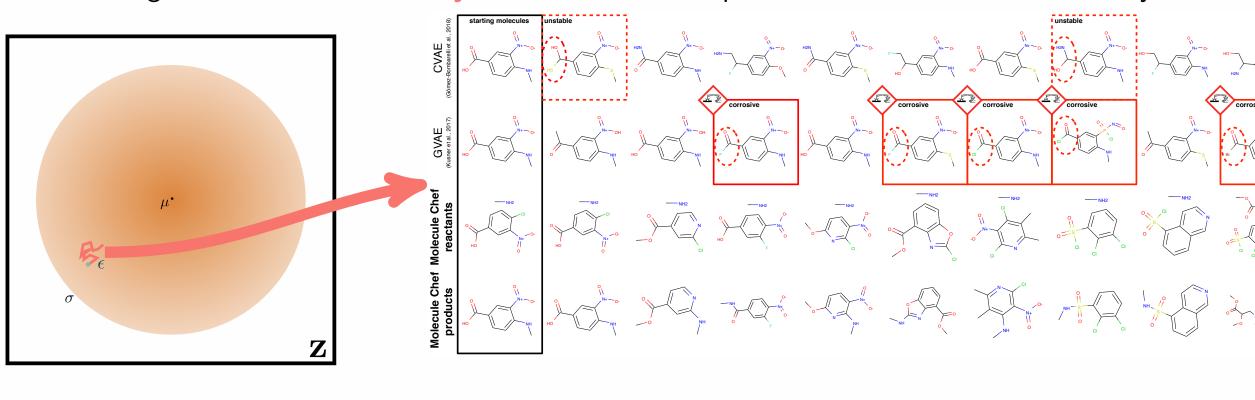
Following previous work we sample 20000 molecules from the prior. We assess these molecules for validity (whether they can be parsed by chemoinformatics software), and conditioned on that, uniqueness and novelty (whether they are distinct from the training molecules).



Model Name	Validity	Uniqueness	Novelty	Quality	FCD
Molecule Chef + MT	99.05	95.95	89.11	95.30	0.73
AAE	85.86	98.54	93.37	94.89	1.12
CGVAE	100.00	93.51	95.88	44.45	11.73
CVAE	12.02	56.28	85.65	52.86	37.65
GVAE	12.91	70.06	87.88	46.87	29.32
LSTM	91.18	93.42	74.03	100.12	0.43
	51.10	JJ.42	1 1.00	100.12	υ. τ.

Tab. 1: Table showing the validity, uniqueness and novelty (all as %, with uniqueness and novelty conditioned on validity ) of the products/or molecules generated from decoding from 20k random samples from the prior  $p(\mathbf{z})$ . MT stands for the Molecular Transformer [4].

We also qualitatively evaluate the semantic validity of our molecules (eg are they stable and non-toxic). To do this we start from a training molecule and randomly walk in the latent space to decode to molecules nearby.



Jorge Aguilera-Iparraguirre, Timothy D. Hirzel, Ryan P. Adams, and Alán Aspuru-Guzik Automatic Chemical Design Using a Data-Driven

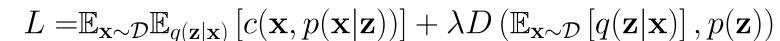
### ...we also want to know how to make it! Therefore, our model, MOLECULE CHEF, generates reactant bags!

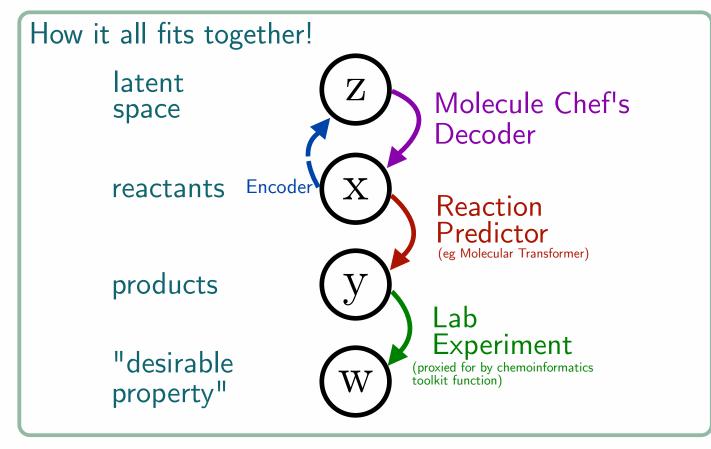
Our model, Molecule Chef, decodes using a two step process.

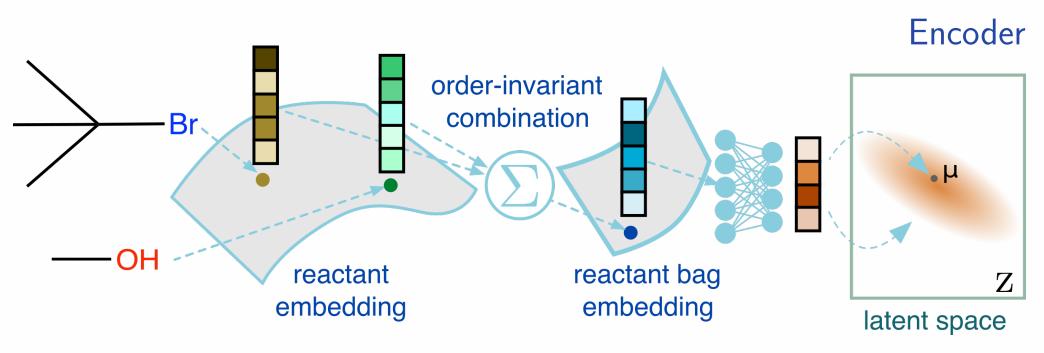
- 1. The decoder first maps from latent space to a reactant bag.
- 2. This is then fed through a reaction predictor model (we use the Molecular Transformer [4]) to predict a final product.

By using stable reactant building blocks, we hope that our model proposes more semantically valid molecules, ie molecules that are non toxic or not about to break down.

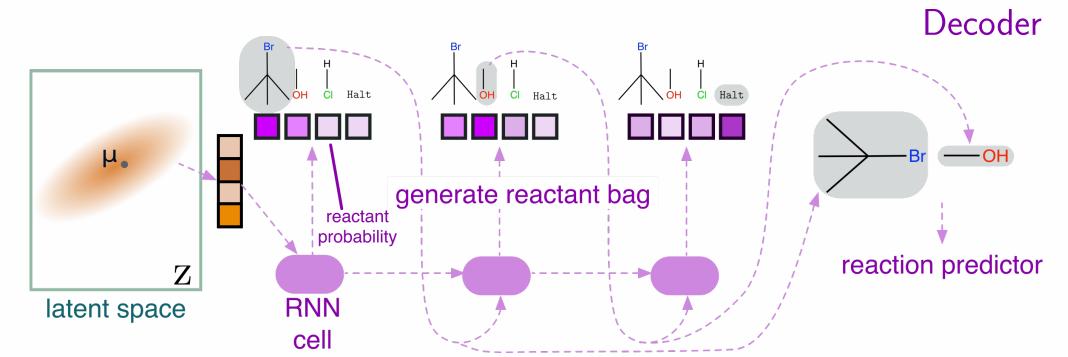
We train the model (on a dataset derived from USPTO) using the WAE objective [7], which involves







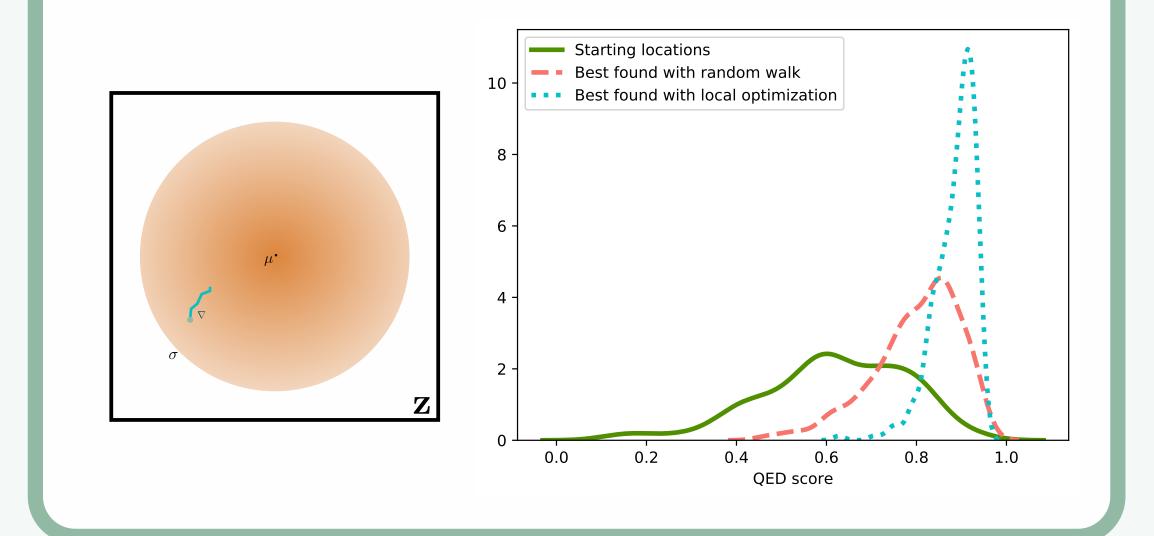
The encoder embeds each reactant to a vector using a graph neural network. These embeddings are summed to produce a reactant bag embedding that is invariant to the order of the reactants in the bag.



The decoder uses a RNN to sequentially output reactants belonging to the reactant bag. Reactants are selected from a fixed set of 3180 easily obtainable and common molecules.

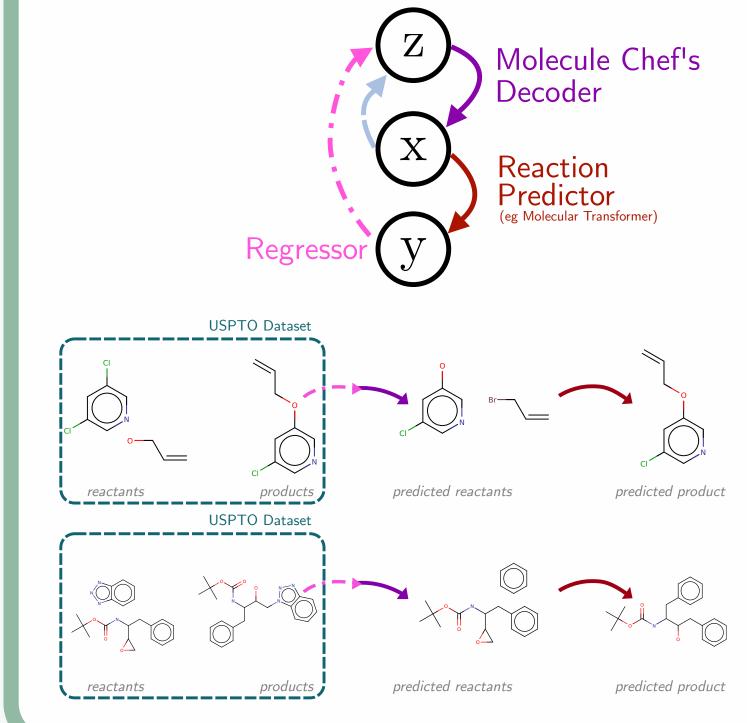
## MOLECULE CHEF's latent space can be used for optimization

We train a regression network from the latent space to a property of interest, the QED (quantitative estimate of drug-likeness). This regression network can be used for local optimization and we compare the QED of the best molecules found this way to those found from a random walk.



## We can also use MOLECULE CHEF for retrosynthesis!

We train a regression network from the products to the latent space, based on graph neural networks. This allows us to do retrosynthesis, ie predicting what reactants created a product.



Can we find molecules that are easier to make but with similar properties? Even if the reconstructed product is not correct we are interested in whether it has similar properties, as our model may be useful in suggesting products with similar properties but that are easier to make. We assess correlations on two subsets of the test set below: left shows results for subset of reactions for which the reactants are all in Molecule Chef's vocabulary, right corresponds to subset of reactions with at least one reactant

## 0.0 0.2 0.4 0.6 0.8 1.0 0.0 0.2 0.4 0.6 0.8 1.0 product's QED product's QED

#### References

Continuous Representation of Molecules ACS Cent. Sci, 2018 [2] Kusner, Matt J., Brooks Paige, and José Miguel Hernández-Lobato. Grammar Variational Autoencoder. ICML, 2017 [1] Gómez-Bombarelli Rafael, Jennifer N. Wei, David Duvenaud, José Miguel Hernández-Lobato, Benjamín Sánchez-Lengeling, Dennis Sheberla, [3] Jin, Wengong, Regina Barzilay, and Tommi Jaakkola. Junction Tree Variational Autoencoder for Molecular Graph Generation. ICML, 2018.

[4] Schwaller, Philippe, Teodoro Laino, Théophile Gaudin, Peter Bolgar, Costas Bekas, and Alpha A. Lee. Molecular Transformer for Chemical [6] Pyzer-Knapp, Edward O., Changwon Suh, Rafael Gómez-Bombarelli, Jorge Aguilera-Iparraguirre, and Alán Aspuru-Guzik. What Is High-Reaction Prediction and Uncertainty Estimation. arXiv preprint arXiv:1811.02633, 2018. Throughput Virtual Screening? A Perspective from Organic Materials Discovery. Annual Review of Materials Research, 2015. [5] Segler, Marwin HS, Thierry Kogej, Christian Tyrchan, and Mark P. Waller. Generating Focused Molecule Libraries for Drug Discovery with [7] Tolstikhin, Ilya, Olivier Bousquet, Sylvain Gelly, and Bernhard Schölkopf. Wasserstein Auto-Encoders. ICLR, 2018. Recurrent Neural Networks. ACS Cen. Sci., 2017. [8] Lowe, Daniel M. Extraction of chemical structures and reactions from the literature PhD Thesis (University of Cambridge), 2012