# MolReactGen: Generating Molecules and Reaction Templates with a Transformer Decoder Model

**Master Thesis Defensio**

Author
**Stephan Holzgruber**
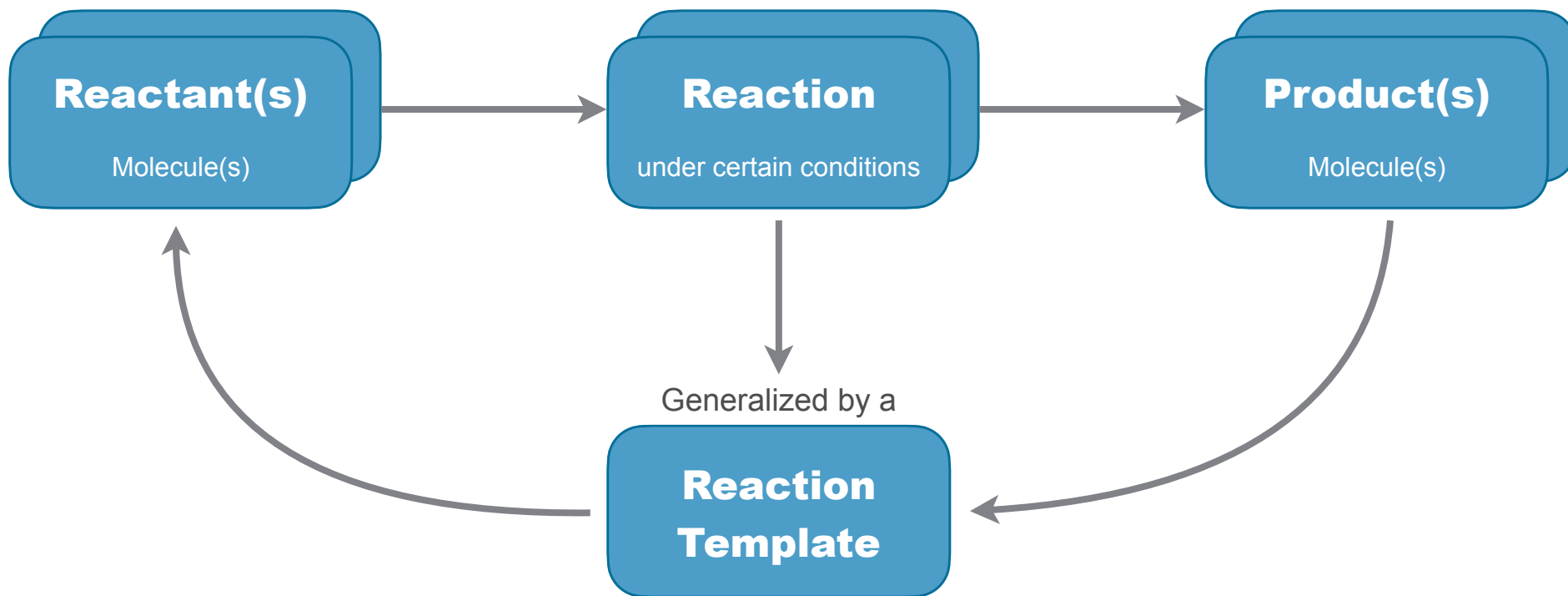
Supervisor
**Günter Klambauer**

Co-Supervisor
**Philipp Seidl**

May 29th, 2024

E-Mail          stephan.holzgruber@gmail.com
Master Thesis   JKU ePUB
Github          hogru/MolReactGen
Hugging Face    hogru/MolReactGen-GuacaMol-Molecules
                hogru/MolReactGen-USPTO50K-Reaction-Templates

**JKU**
JOHANNES KEPLER
UNIVERSITÄT LINZ

# Application Area: Chemical Reactions

# Context / Motivation

| | |
|---:|:---|
| What's it all about? | Design new product molecules |
| What for? | New drugs to address gaps in disease treatment, ... |
| Why care? | Multi-drug resistant bacteria, neuro-degenerative diseases, ... |
| The issue? | Huge search space (~$10^{33}$[1]), cost, time, skills, equipment, ... |
| How? | Generative model (LLM) ➔ Similar, but different molecules |
| How is that useful? | Rapid, cheaper creation and screening of candidate drugs |
| Fine, and then? | How to synthesize those molecules? |
| Any ideas? | Start again, with generating reaction templates |

[1] P. G. Polishchuk et al., "Estimation of the size of drug-like chemical space based on GDB-17 data," J Comput Aided Mol Des, vol. 27, no. 8, pp. 675–679, Aug. 2013, doi: 10.1007/s10822-013-9672-4

**JOHANNES KEPLER UNIVERSITY LINZ**

# Research Questions

- **GuacaMol**[2] considered a reference paper/model for de novo molecule generation

- Research Questions
  - What is the **performance of a transformer decoder** architecture compared to GuacaMol?
  - What is the effect of different **tokenization approaches**?
  - Can we use a model pre-trained on natural language as a basis for **fine-tuning a "molecule language" model**?
  - Can the transformer decoder model also be used to **generate reaction templates**?

[2] N. Brown et al., "GuacaMol: Benchmarking Models for de Novo Molecular Design," *J. Chem. Inf. Model.*, vol. 59, no. 3, pp. 1096–1108, Mar. 2019, doi: 10.1021/acs.jcim.8b00839.

JOHANNES KEPLER
UNIVERSITY LINZ

# Data — Intuition

| Target | Source / Preprocessing | Format | Count | ø Length (train set) | Example |
|---|---|---|---|---|---|
| **Molecules** | ChEMBL[3] / GuacaMol | SMILES | 1.6 M | 48 | `O=C(O)C1CCC(OCC2CC(F)CN2C(=O)Cc2ccc(NC(=O)N3CCc4ccccc43)c(Cl)c2)CC1` |
| **Reaction Templates** | USPTO-50K[4] / MHNReact[5] | SMARTS | 12 K | 161 | `[#7;a:4]:[c:3]:[c;H0;D3;+0:1](:[#7;a:2])-[n;H0;D3;+0:9]1:[#7;a:5]:[c:6]:[#7;a:7]:[c:8]:1>>` `Cl-[c;H0;D3;+0:1](:[#7;a:2]):[c:3]:[#7;a:4].` `[#7;a:5]1:[c:6]:[#7;a:7]:[c:8]:[nH;D2;+0:9]:1` |

[3] D. Mendez et al., "ChEMBL: towards direct deposition of bioassay data," Nucleic Acids Res, vol. 47, no. D1, pp. D930–D940, Jan. 2019, doi: 10.1093/nar/gky1075
[4] D. M. Lowe, "Extraction of chemical structures and reactions from the literature," Thesis, University of Cambridge, 2012. doi: 10.17863/CAM.16293
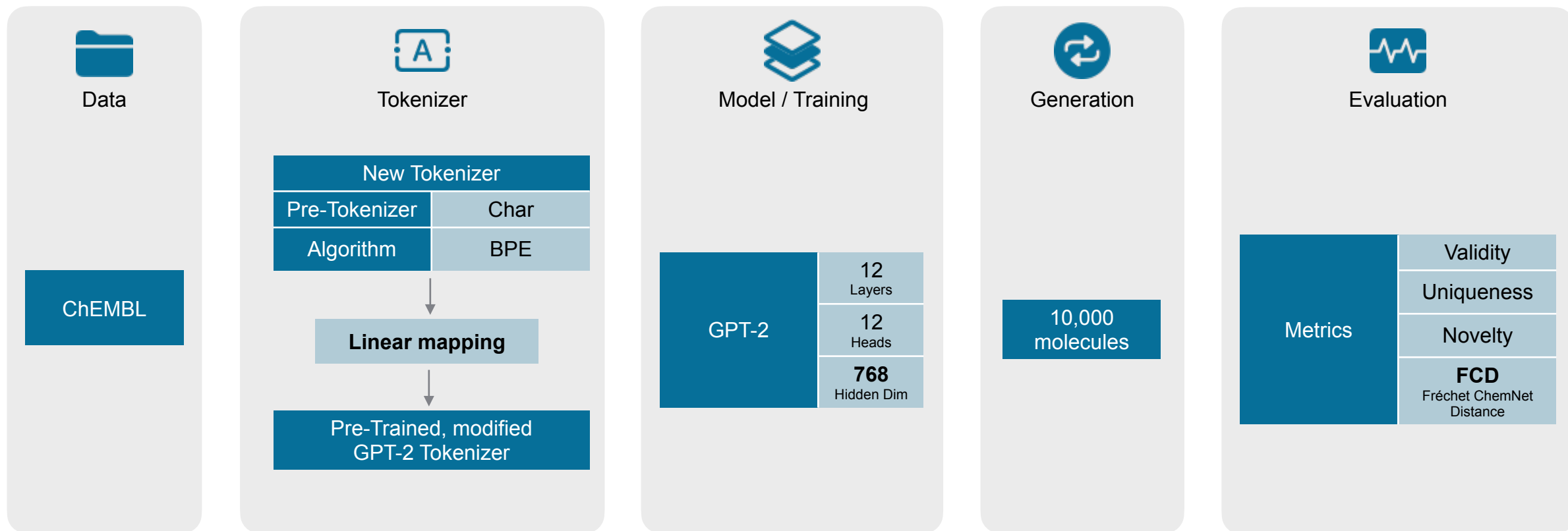[5] P. Seidl et al., "Modern Hopfield Networks for Few- and Zero-Shot Reaction Template Prediction," arXiv:2104.03279 [cs, q-bio, stat], Jun. 2021, Accessed: Nov. 02, 2021. [Online]. Available: http://arxiv.org/abs/2104.03279

JOHANNES KEPLER
UNIVERSITY LINZ

# Pipeline 1/3 — <u>Molecules From Scratch</u>



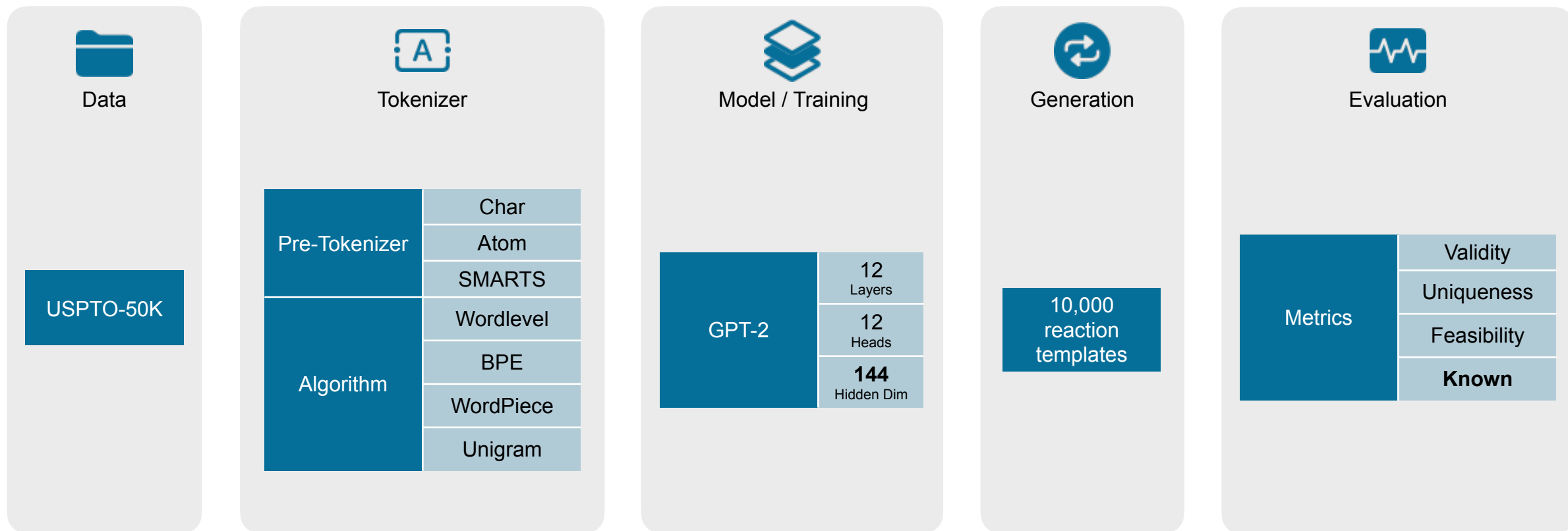| Data | Tokenizer | | Model / Training | | Generation | Evaluation | |
|------|-----------|--|------------------|--|------------|------------|--|
| ChEMBL | Pre-Tokenizer | Char | GPT-2 | 12 Layers | 10,000 molecules | Metrics | Validity |
| | | Atom | | 12 Heads | | | Uniqueness |
| | | SMARTS | | **144** Hidden Dim | | | Novelty |
| | Algorithm | Wordlevel | | | | | **FCD**[6] Fréchet ChemNet Distance |
| | | BPE | | | | | |
| | | WordPiece | | | | | |
| | | Unigram | | | | | |

[6] K. Preuer et al., "Fréchet ChemNet Distance: A Metric for Generative Models for Molecules in Drug Discovery," J. Chem. Inf. Model., vol. 58, no. 9, pp. 1736–1741, Sep. 2018, doi: 10.1021/acs.jcim.8b00234

JOHANNES KEPLER
UNIVERSITY LINZ

# Pipeline 2/3 — <u>Molecules from Pre-Trained Model</u>

**Data**

ChEMBL

**Tokenizer**

| New Tokenizer | |
|---|---|
| Pre-Tokenizer | Char |
| Algorithm | BPE |

↓

**Linear mapping**

↓

Pre-Trained, modified GPT-2 Tokenizer

**Model / Training**

| GPT-2 | 12 Layers |
|---|---|
| | 12 Heads |
| | **768** Hidden Dim |

**Generation**

10,000 molecules

**Evaluation**

| Metrics | Validity |
|---|---|
| | Uniqueness |
| | Novelty |
| | **FCD** Fréchet ChemNet Distance |

# Pipeline 3/3 — <u>Reaction Templates From Scratch</u>

**Data**

USPTO-50K

**Tokenizer**

| Pre-Tokenizer | Char |
| | Atom |
| | SMARTS |
| Algorithm | Wordlevel |
| | BPE |
| | WordPiece |
| | Unigram |

**Model / Training**

| GPT-2 | 12 Layers |
| | 12 Heads |
| | **144** Hidden Dim |

**Generation**

10,000 reaction templates

**Evaluation**

| Metrics | Validity |
| | Uniqueness |
| | Feasibility |
| | **Known** |

JⱯU JOHANNES KEPLER
UNIVERSITY LINZ

# Selected Results
# Generation of 10,000 <u>Molecules</u>

| Model | Tokenizer<br>Pre-Tokenizer \| Algorithm \| Vocab Size | Metrics | | | |
|---|---|---|---|---|---|
| | | **Validity** ↗ | **Uniqueness** ↗ | **Novelty** ↗ | **FCD** ↘ |
| GuacaMol | | 0.959 | **1.000** | **0.994** | 0.455 |
| MolReactGen *from scratch* | Char \| Wordpiece \| 176 | **0.976** ± 0.001 | 0.999 ± 0.000 | 0.935 ± 0.002 | **0.219** ± 0.005 |
| | | 0.976 ± 0.001 | **0.999** ± 0.000 | **0.935** ± 0.002 | 0.219 ± 0.005 |
| MolReactGen *fine-tuned* | Char \| BPE \| 50,527 | **0.990** ± 0.001 | 0.999 ± 0.000 | 0.797 ± 0.004 | **0.209** ± 0.006 |

Molecule generation results
Numbers represent the mean and standard deviation (superscript) across five training and generation runs
FCD metric not stated in GuacaMol paper, calculated here as $-5 \log \mathrm{FCD_{GuacaMol}}$

# Selected Results
# Generation of 10,000 <u>Reaction Templates</u>

| Model | Tokenizer<br><br>Pre-Tokenizer \| Algorithm \| Vocab Size | Metrics | | | |
|---|---|---|---|---|---|
| | | Validity ↗ | Uniqueness ↗ | Feasibility ↗ | Known ↗ |
| MolReactGen *from scratch* | SMARTS \| Wordlevel \| 947 | $0.804 \pm 0.022$ | $0.795 \pm 0.008$ | $0.110 \pm 0.004$ | $735.2 \pm 27.0$ |

Reaction template generation results
Numbers represent the mean and standard deviation (superscript) across five training and generation runs

**JⱮU** JOHANNES KEPLER
UNIVERSITY LINZ
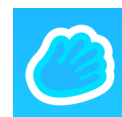
# Conclusion wrt Research Questions

- Transformer decoder (GPT-2) architecture exhibits a **better FCD** than the GuacaMol baseline

- **Tokenizers** exert "some" influence on FCD (max `0.022`, but statistically significant)

- **Fine-tuning** a pre-trained language model **works**, but at a (computational) cost

- Approach (with different optimal tokenizer) also **works with reaction templates**

**JꞰU** **JOHANNES KEPLER
UNIVERSITY LINZ**

# "Giving back to the community"

- `PySmilesUtils`: 2 accepted PRs ([#1](#), [#2](#))

- `rdkit`: [bug report](#) (fixed)

- `Hugging Face`
  - Implementation of P. Schwaller's [feature request](#)
  - 4 issue contributions ([#1](#), [#2](#), [#3](#), [#4](#))
  - Member of HF "Helping Hands"
  - **> 1K pre-trained model downloads** ([#1](#), [#2](#))

- [Initial port](#) of GuacaMol evaluation code to current packages

- `FCD`: 2 accepted PRs ([#1](#), [#2](#))

- [First citation](#) of master's thesis 😉

The efficacy of MHNreact has been assessed in various studies (Chen et al., 2023; Liu et al., 2022), and its integration and testing on additional datasets have yielded promising results (Torren-Peraire et al., 2023). Further benchmarking efforts, particularly under multi-step conditions as demonstrated in (Maziarz et al., 2023), would provide valuable insights into its performance and applicability across different scenarios. Incorporating generated reaction-templates (Holzgruber, 2024) into MHNreact is possible due to its zero shot capability, and could overcome the often discussed limitation of template-based methods, namely their upper accuracy bound due to unreachable reactants in the test-set.

JOHANNES KEPLER
UNIVERSITY LINZ

# Details — Links into the Master's Thesis

- Molecules — From Scratch
  - [Data and Tokenization](#)
  - [Hyper Params](#)
  - [Metrics](#)
  - [Tokenizer Selection Process](#)
  - [Tokenizer Results](#)
  - [Results](#)
- Molecules — Pre-trained
  - [Tokenizer Mapping](#)
  - [Hyper Params](#)
  - [Results](#)

- Reaction Templates
  - [Data and Tokenization](#)
  - [Hyper Params](#)
  - [Metrics](#)
  - [Reaction Template Feasibility](#)
  - [Tokenizer Results](#)
  - [Results](#)
- [RegEx Patterns](#)

**JⴑU** JOHANNES KEPLER
UNIVERSITY LINZ