

When yield prediction does not yield prediction: an overview of the current challenges

Varvara Voinarovska, PhD student, AstraZeneca MAI group

19.09.2023



Plan:

Introduction

Data side

SotA side

Benchmarking

Outlook



Why do we care: Importance of Reaction Yield Prediction

 In CASP: We could use yield prediction to filter out unsuccessful reactions in Computer-Aided Synthesis Planning (CASP). This could significantly cut costs and boost sustainability.

- In HTE: In High-Throughput Experimentation (HTE) accurate yield prediction can help us optimize synthesis processes, saving time and resources.





Factors Influencing Yield of a Chemical Reaction

- Low Reactivity
- Side Reactions
- Reactant/Reagent/Catalyst Deactivation
- Thermodynamic and Kinetic Factors
- Contaminants
- Sensitivity to Environment
- Product Degradation/Reactivity
- Product Isolation



Data generation and sources



ORD



Flow Chemistry

High-Throughput Experimentation



Electronic Lab Notebooks



Reaction mining



What's the problem with the data storage?



What is the SotA?

For *low-scale* data the Active Learning approach is popular. The design of fingerprints, including advanced ones is popular.

For *large-scale* data the Deep Learning (Transformers, Graph Neural Networks) are current SoTA.





SotA overview

Data encoding:

- DFT
- Structural fingerprints (ECFP, DRFP)
- Learned representations
- Graph-based Encodings
- One-hot

Methods:

- XGBoost, Random Forest and other classics
- Bayesian modeling
- Active learning approach

Natural Language Processing:

- Yield-BERT
- Multimodal Transformer
- Augmented Transformer

Graph-based DL:

- MPNN with self-attention
- Uncertainty-aware MPNN
- YieldGNN

Schwaller, P.; Vaucher, A. C.; Laino, T.; Reymond, J.-L. Prediction of chemical reaction yields using deep learning. Machine Learning: Science and Technology 2021, 2, 015016 Saebi, M.; Nan, B.; Herr, J. E.; Wahlers, J.; Guo, Z.; Zura´nski, A. M.; Kogej, T.; Norrby, P-O.; Doyle, A. G.; Chawla, N. V.; Wiest, O. On the use of real-world datasets for reaction yield prediction. Chemical Science 2023, 14, 4997–5005

Baraka, S.; Kerdawy, A. M. E. Multimodal Transformer-based Model for Buchwald-Hartwig and Suzuki-Miyaura Reaction Yield Prediction. 2022; https://anxiv.org/abs/2204.14062 Kwon, Y.; Lee, D.; Choi, Y-S.; Kang, S. Uncertainty-aware prediction of chemical reaction yields with graph neural networks. Journal of Cheminformatics 2022, 14 Yarish, D.; Garkot, S.; Grvaorenko, O. O.; Radchenko, D. S.; Morzo, Y. S.; Gurbvch, O. Advancing molecular graphs with descriptors for the prediction of chemical reaction yields.

Journal of Computational Chemistry 2022, 44, 76-92

Neves, P.; McClure, K.; Verhoeven, J.; Dyubankova, N.; Nugmanov, R.; Gedich, A.; Menon, S.;zhicai Shi,; Wegner, J. Global Reactivity Models are Impactful in Industrial Synthesis Applications. 2022



"Solved" Case

Buchwald-Hartwig HTE dataset:

15 aryl halides, 23 additives, 4 palladium catalysts, and 3 bases

Dense and consistent

Toy dataset for yield prediction







t-SNE of BH HTE dataset demonstrates that there are distinctive clusters from which one could see the areas with lower or higher yield.





Gradient Boost Regression for BH HTE dataset in different featurizations





"Unsolved" Case

Buchwald-Hartwig Amination reaction

- Reaxys 7K entries
- AZ ELN 750 500 entries
- Ahneman's HTE Buchwald-Hartwig (BH HTE) - 4K entries
- USPTO 6K entries









t-SNE without conditions



t-SNE with conditions



USPTO ID01456115



BH HTE example











What could we do about that?

Data Standardization

Deep Integration of Chemistry

Yield Variability: a need for classification models with multiple bins to address data complexity.

Future Trajectory: The future of yield prediction involves *enhanced datasets*, uncertainty-based predictions, and the development of *reaction-specific descriptors*.



Yield prediction is influenced by inherent data noise. An analysis revealed standard deviations of around 16% in general datasets.



Acknowledgments

Supervisors:





Dr. Igor Tetko

Dr. Dr. Dmytro Samuel Dudenko Genheden





Dr. Mikhail Kabeshov





Advanced machine learning for Innovative Drug Discovery (AIDD) Horizon 2020 Marie Skłodowska-Curie Innovative Training Network - European Industrial Doctorate



The AIDD Team

Early-Stage Researchers (ESRs)



Peter Hartog



Emma Svensson





Peraire



Varvara Voinarovska



Julian

Cremer



Son Hà



Hassen

ESR15



Ana Sanchez





Yasmine Nahal



Rosa Friesacher



Vincenzo

Pamacci

ESR12



Mikhail Andronov



Allesio Fallani





Mathias

Hilfiker



Mariia Radaeva











