# Predicting Reaction Success Using A Transformer Model Pretrained on Reaction SMILES Data

November 2023

Eric Gilbert, PhD

# Outline

- Motivation and challenges.

- Workflow for creating a fine-tuned BERT model for yield prediction.

- Benchmarking

- Multi-modal learning: experimental text and reaction SMILES.

- Use cases

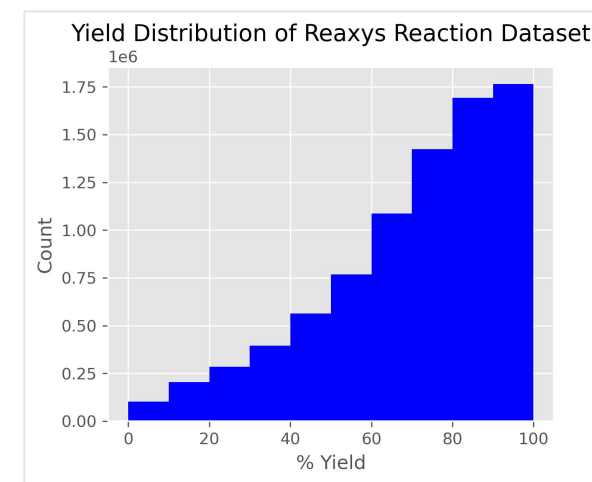# Predicting Reaction Yield

## Motivation

- ~20% of reactions fail or yield too low[1]

- wasteful- materials, human resources, opportunity cost, time

## Challenges

- Literature and patent data biased toward higher yielding reactions.

- Models need to learn from failed reactions.

## Strategy

- Pretrain a model from scratch using Reaxys reaction data.

- Fine-tune model on ELN data for predicting reaction success.
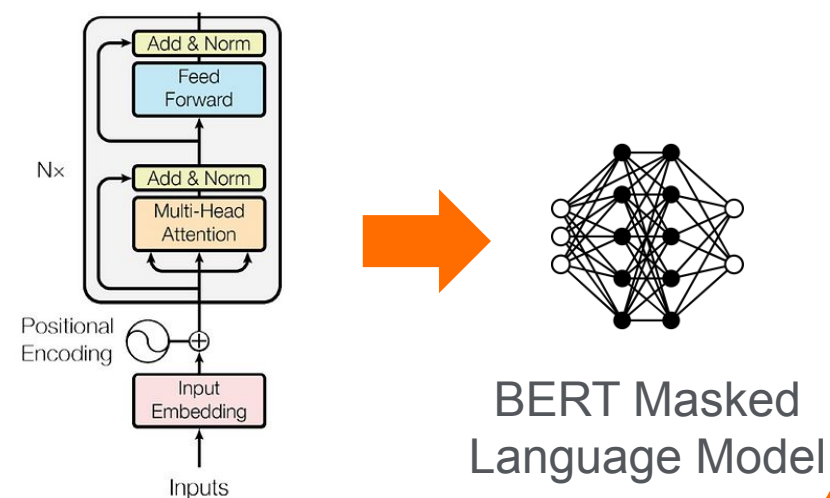
- Binary classification task- yield >5% or <5%



Yield Distribution of Reaxys Reaction Dataset

[1]Neves, P., McClure, K., Verhoeven, J. *et al.* Global reactivity models are impactful in industrial synthesis applications. *J Cheminform* **15**, 20 (2023). https://doi.org/10.1186/s13321-023-00685-0

# Workflow



**pre-training**

preprocessed reaction SMILES from Reaxys

BERT Masked Language Model

- pretraining with ~9M reactions, 4 GPUs, 30 epochs, ~10 days

**fine tuning**

ELN training data*

**inference**

test set

+ pretrained BERT MLM → fine-tuned BERT MLM → model performance metrics

*ELN data, model fine tuning, and inference on secure Amazon Workspaces hosted by Janssen.

ELSEVIER

4

# Suzuki Reaction Benchmark



5760 reactions[1]

11 ligands
6 boronic acids
4 aryl halides
7 bases
4 solvents

| Base model | $R^2$ |
|---|---|
| MLM-only pretrain | 0.80 ± 0.01 |
| Dual pretrain | 0.82 ± 0.01 |
| [2]Schwaller rxnfp | 0.79 ± 0.01 |



Eval R2 for 10 Splits
95% CI



Yield Distribution of Suzuki HTE Dataset

□ Dual pretraining leads to statistically significant improvement in model performance.

1) Perera et al., Science 359, 429–434 (2018)
2) Philippe Schwaller *et al* 2021 *Mach. Learn.: Sci. Technol.* **2** 015016

# Comparison of Embedding Projections for Pretrained Models

MLM only pretrain

Dual pretrain



Legend:
- addition to c=c
- hydrolysis of esters
- addition to c=o
- aliphatic n-acylation of amines
- n-alkylation
- addition to alkynes
- reduction of nitro
- palladium catalysed c-c coupling

dual pretrained model shows improved clustering of reaction classes

# Use Cases – Data Quality

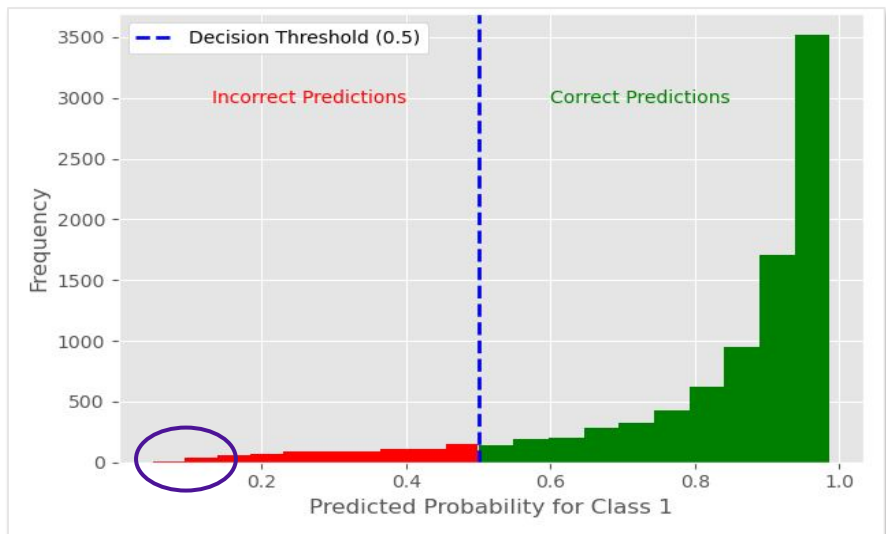☐ Data quality insights from test set inference.

- can be instructive to look at what the model got 'most wrong'

- helpful if slice data by reaction class when interrogating



- missing reagents?

- missing catalysts?

- erroneous reactants?

- etc.

☐ Data quality issues may not be obvious when looking at model metrics by class.

- models are surprisingly resilient at learning from messy data.

☐ Use insights to inform anomaly detection in training data embedding clusters.

# Use Cases - Synthesis

- Focus high throughput experimentation (HTE) efforts.

  - can create a combinatorial combination of potential reagents

  - rank order probabilities of reaction success

- Aid medicinal chemists on more focused synthetic queries:

  - which solvent is predicted to be best for this transformation?

  - rank order potential targets based on predicted probability of success

- Incorporate into multi-modal model.

# Multi-modal Deep Learning

- Contrastive learning on procedure text and reaction SMILES

  - 'foundation model'

- Text features are associated with reaction SMILES
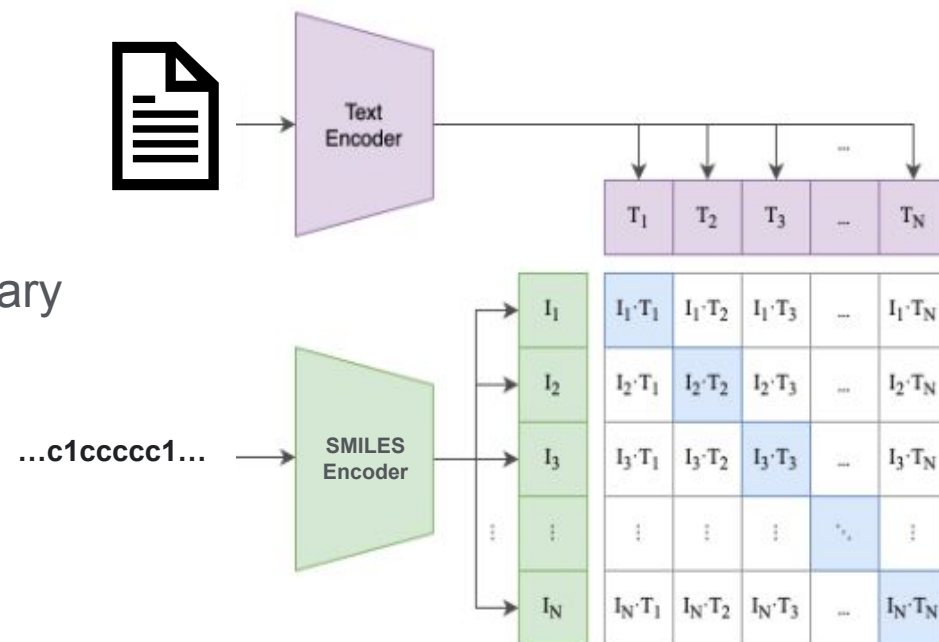
  - domain adaptation between literature and patent text necessary

- Enables new applications:

  - Zero / Few-shot learning

  - Cross-modality search

- Example of association of text and reaction SMILES:

  - low temperature in procedure correctly associated with solvent in SMILES:
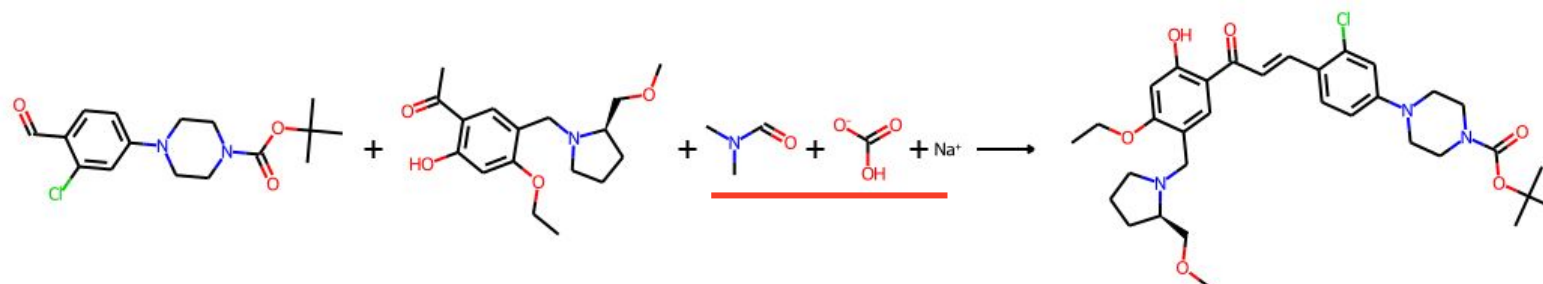
# Use Cases – Data Quality

- Example from Reaxys with data inconsistency:

**Procedure text:**

'Add 303mg B-122 to 4ml <span style="color:red">EtOH</span>, then add 221mg <span style="color:red">KOH</span>, 640mg B-10, stir at RT, and monitor LC-MS until there is no B-122 left.'

- Inconsistency identified using cosine similarity between vector embeddings of text and SMILES.
- Potential application:
  - Prevent errors in ELN data entry ☐ improve data quality ☐ reproducibility, better AI models.

# Use Cases – Search

- SMILES-to-Text search

    - can we suggest conditions / workup from existing procedures using reaction SMILES as input?

    - use similarity between vector embeddings of text and SMILES

    - Advantage: Suggest applicable existing procedures rather than black box prediction of conditions.

ELSEVIER

# Summary

- Pretrained a BERT Masked Language Model from scratch using Reaxys reaction SMILES.
    - investigated impact of adding reaction classification task to pretraining
- Fine tuned models on Janssen ELN data to predict reaction success (>5% yield).
- Benchmarked MLM-only & dual-pretrained models on Suzuki benchmark.
    - use cases- data quality, anomaly detection & synthesis.
- Multi-modal model- contrastive training with reaction SMILES & procedure text.
    - use cases- data quality / consistency, SMILES-to-Text search
- Demonstrated using proprietary ELN data hosted on server by pharma partner.
- Pretrained model can readily be used with other companies and their ELN data.

# Acknowledgements