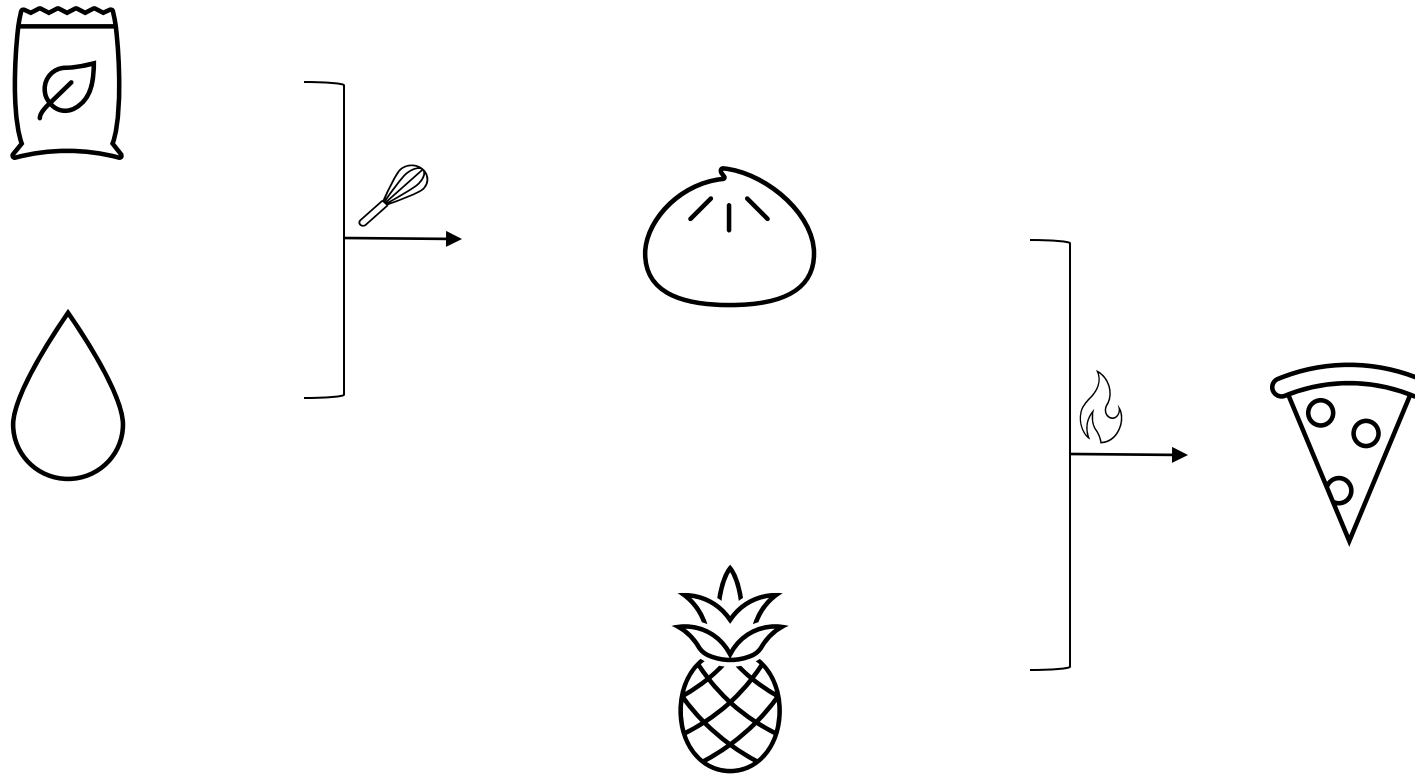


Pushing the boundaries of chemical synthesis with RetroChimera

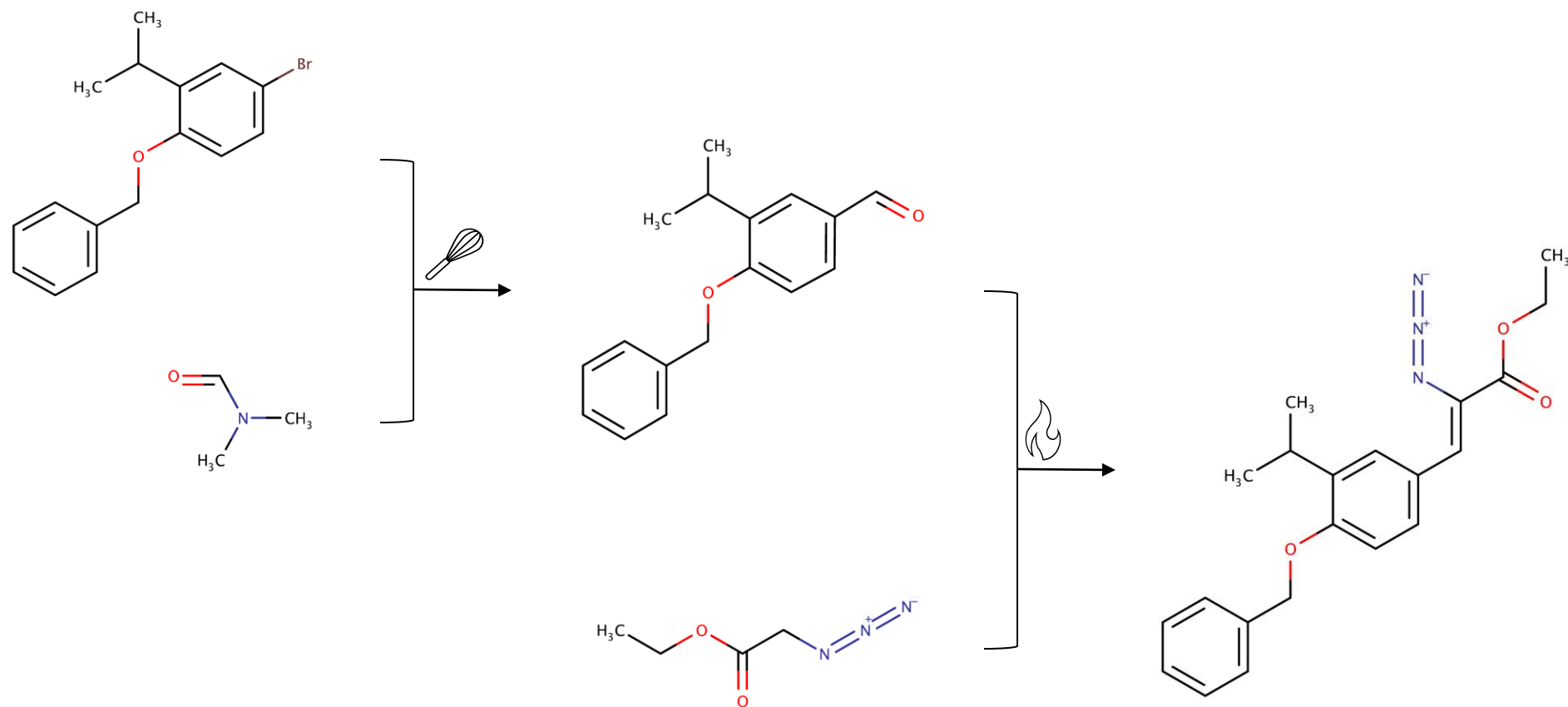
Krzysztof Maziarz
Principal Researcher @ Microsoft Research AI for Science



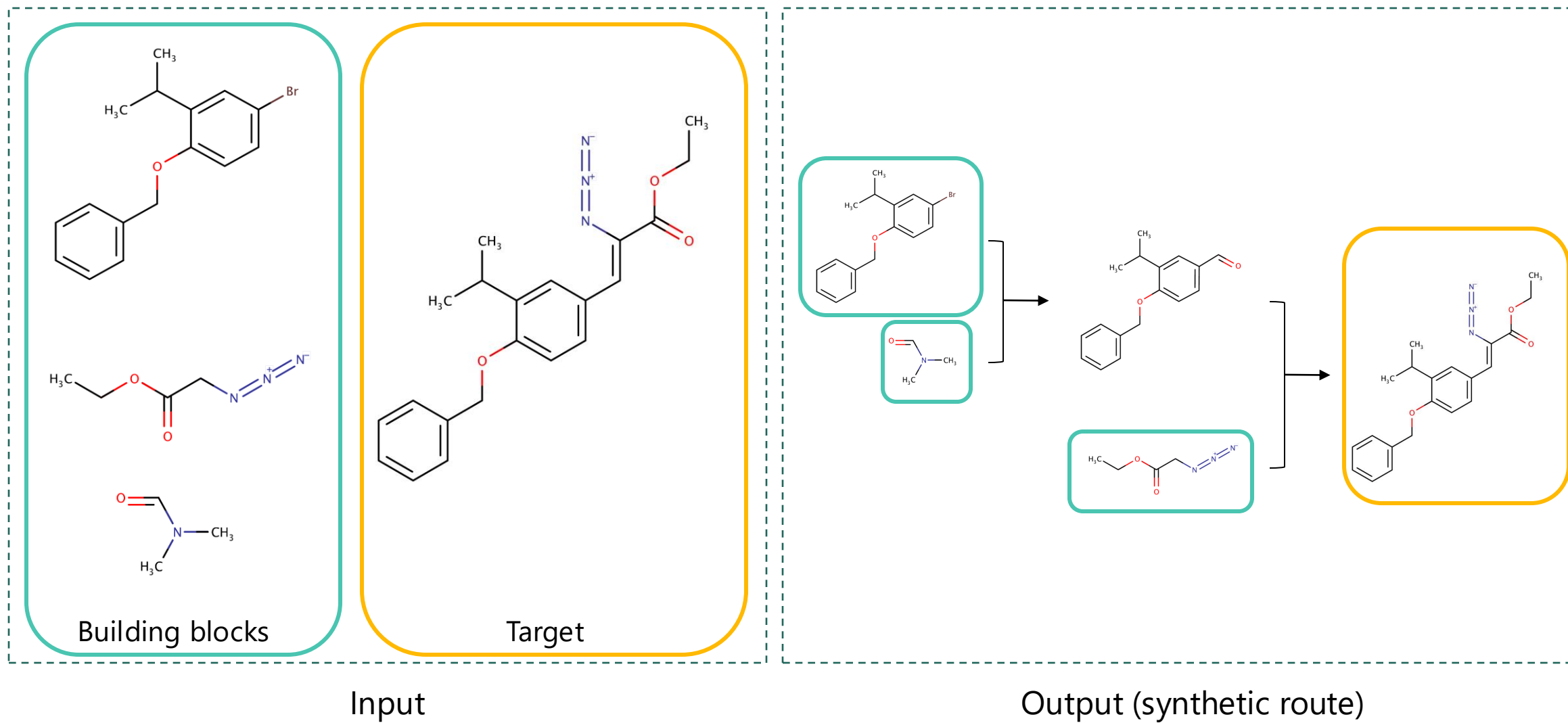
Synthesis: cooking with molecules



Synthesis: cooking with molecules



Synthesis planning

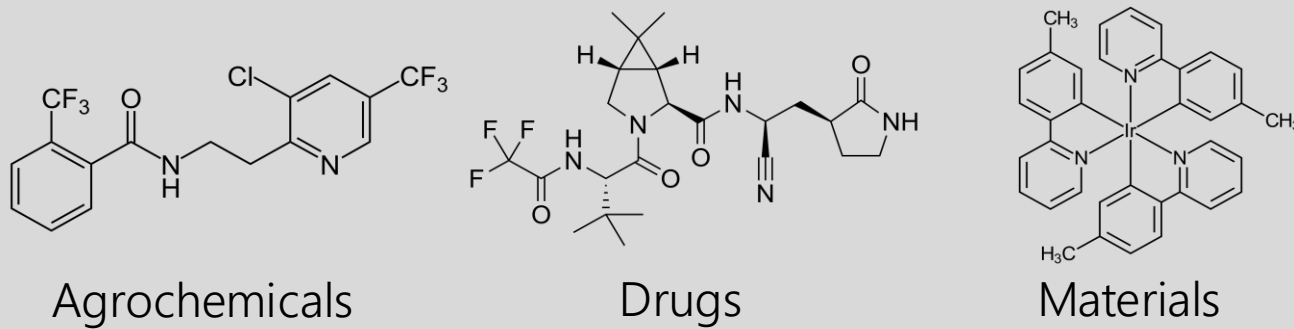


Demo time

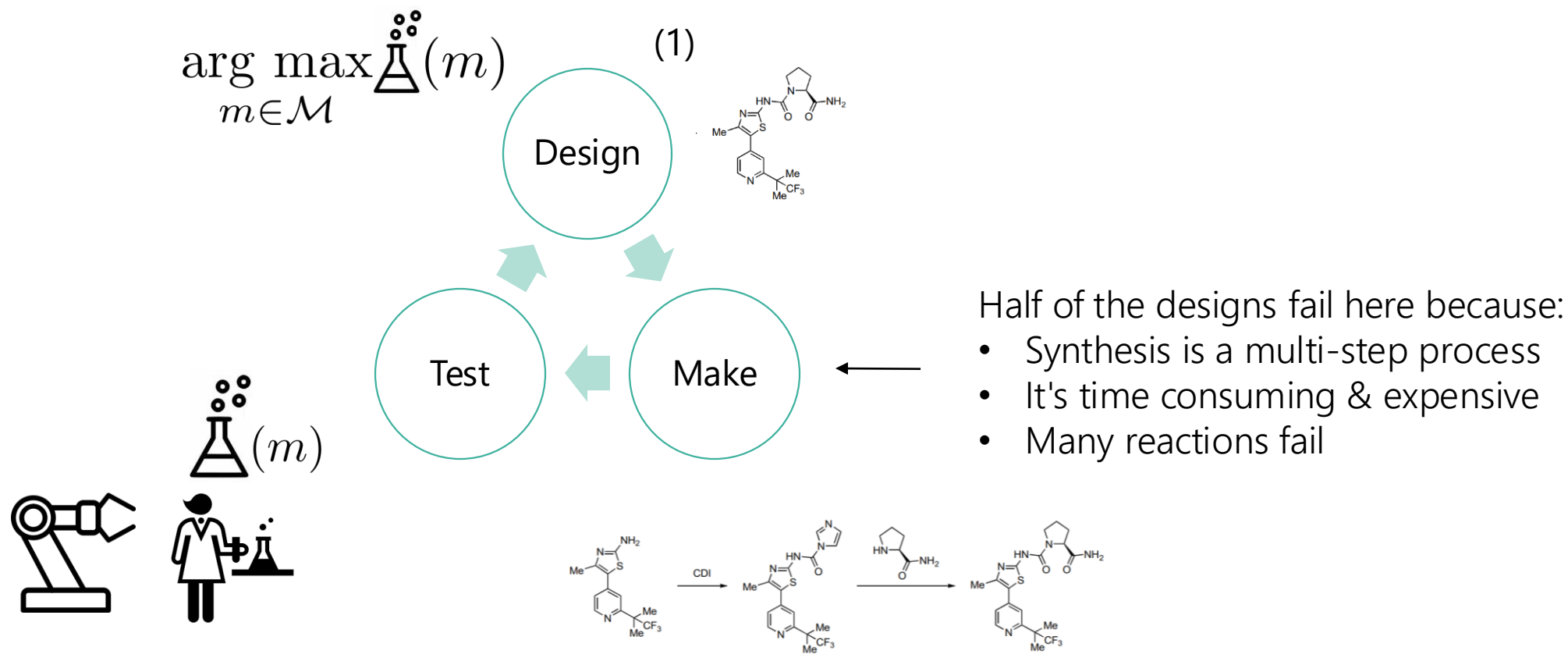


Why is synthesis planning important?

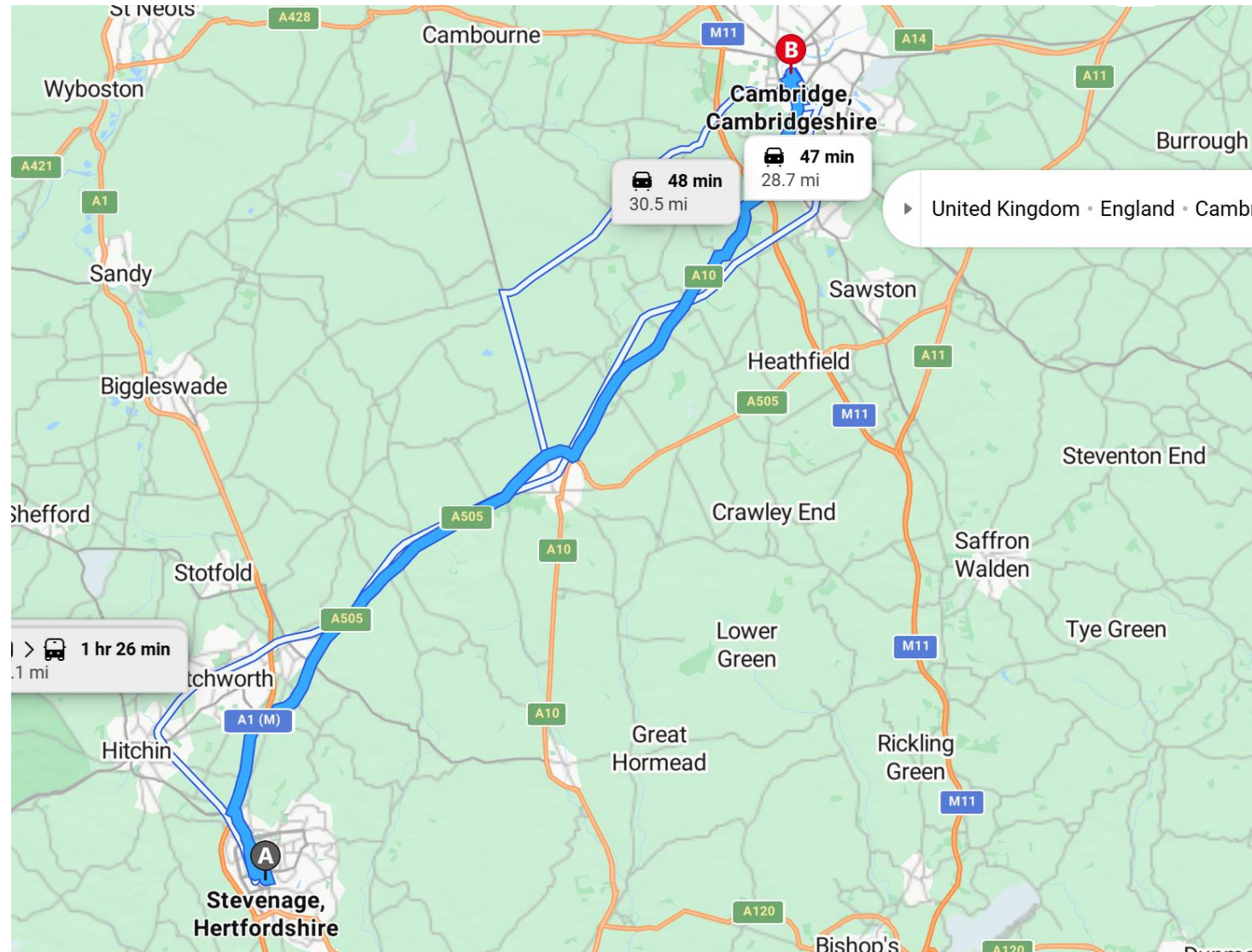




Synthesis is a bottleneck in drug discovery



GPS for molecules?



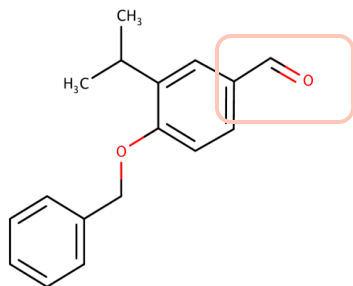
What kind of modelling is
RetroChimera based on?



How to model chemical reactions

Molecular Graph

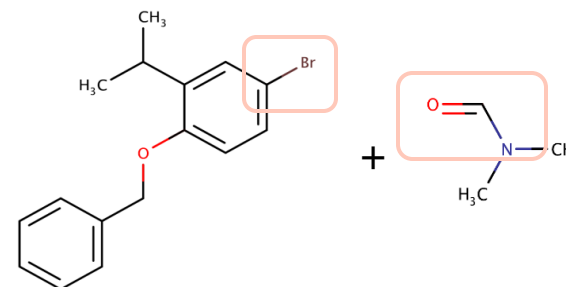
Input (Product)



CC(C)c1cc(C=O)ccc1OCc1ccccc1

SMILES
(Graph as Token Sequence)

Output (Reactants)

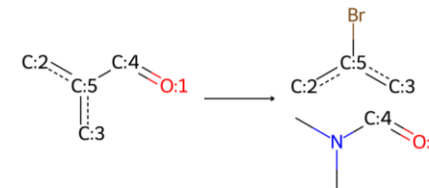


CN(C)C=O.CC(C)c1cc(Br)ccc1OCc1ccccc1

Predict Reactants *De Novo*

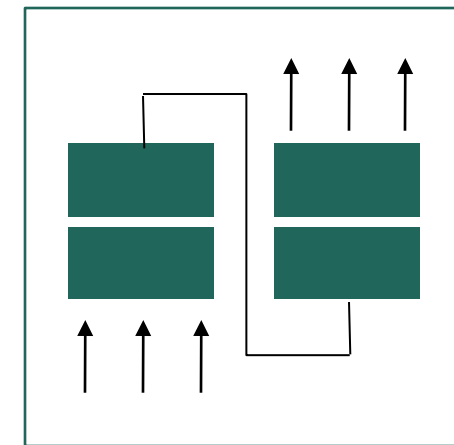
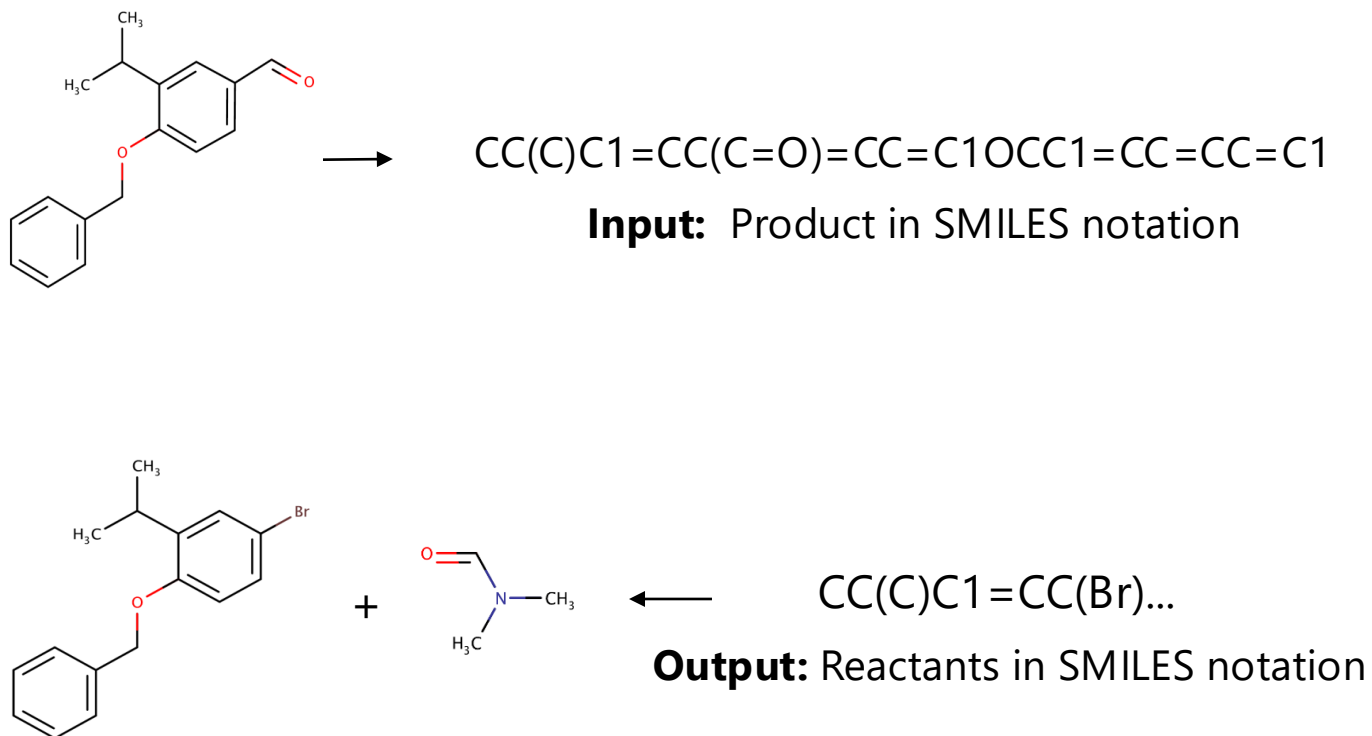
CN(C)C=O.CC(C)c1cc(Br)ccc1OCc1ccccc1

Predict Molecular *Edits*



Applying Edit Template to product yields reactants

Approach #1: end-to-end Transformer



Encoder-Decoder Transformer



Approach #1: end-to-end Transformer

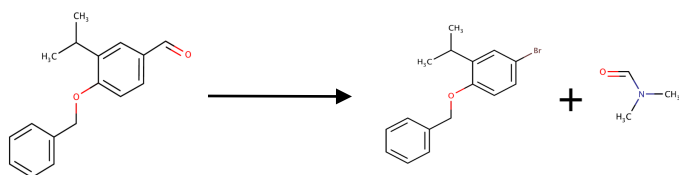
1. R-SMILES alignment and augmentation

2. Model architecture

Product

Reactants

Molecular Graph:



1. Canonical SMILES:

CC(C)c1cc(C=O)ccc1OCc1ccccc1 → CC(C)c1cc(Br)ccc1OCc1ccccc1.CN(C)C=O

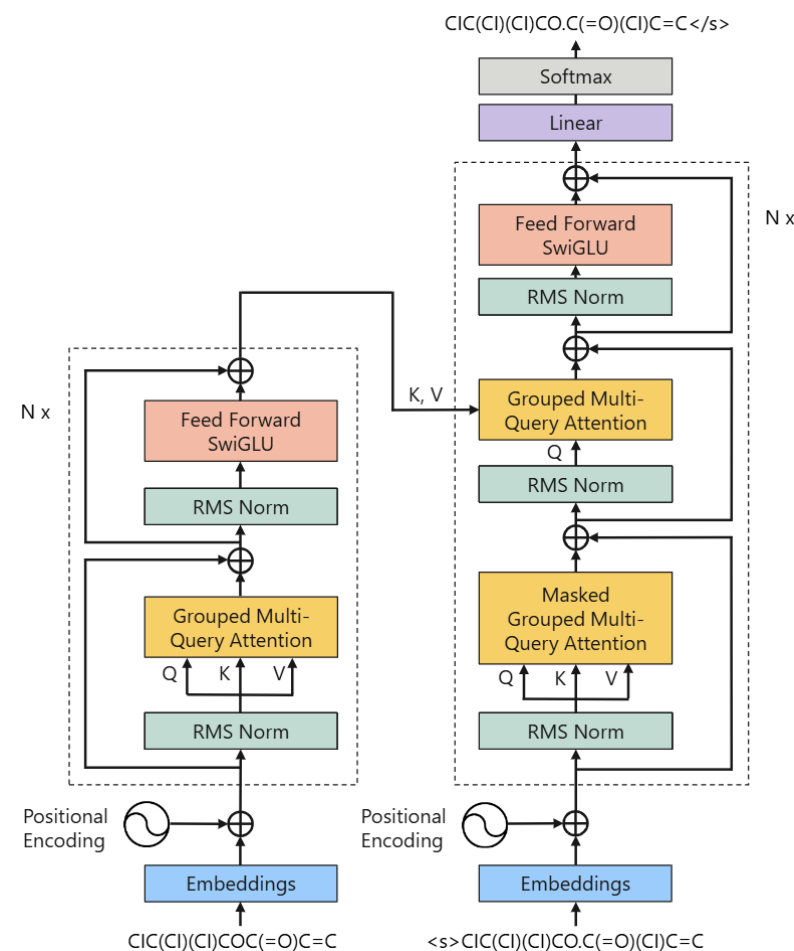
2. Randomized SMILES:

CC(C)c1cc(C=O)ccc1OCc1ccccc1 → CN(C)C=O.c1(C(C)C)c(OCc2ccccc2)ccc(Br)c1
c1ccccc1COc1c(C(C)C)cc(C=O)cc1 → CN(C)C=O.c1(COC2ccc(cc2C(C)C)Br)ccccc1
c1cccc(COc2c(C(C)C)cc(C=O)cc2)c1 → C(C)(C)c1cc(ccc1OCc1ccccc1)Br.CN(C=O)C

3. Root-aligned SMILES (plus N times augmentation):

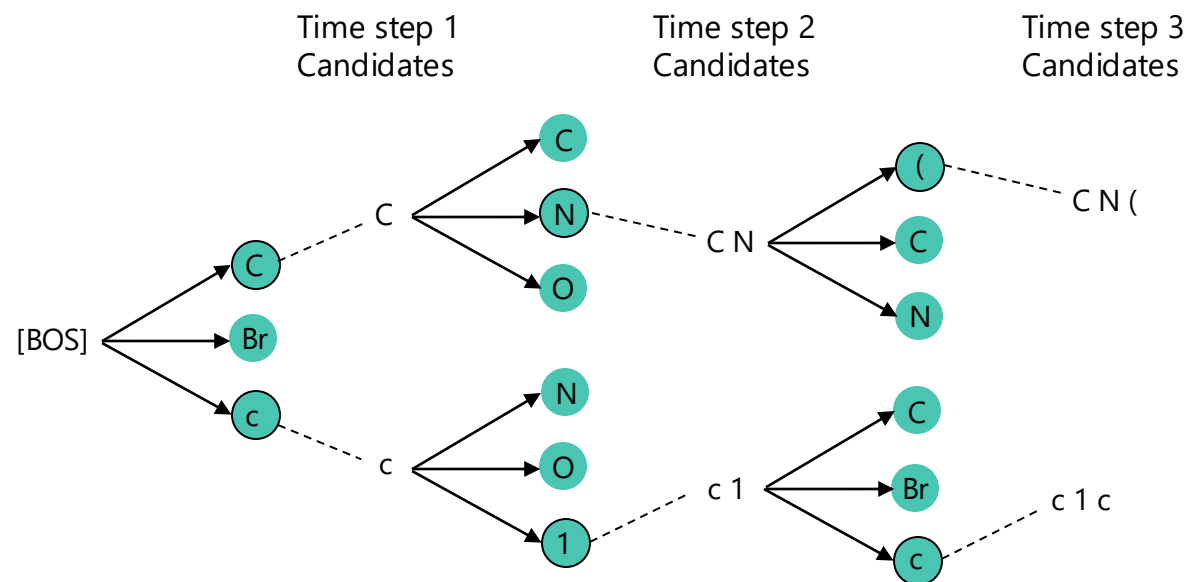
CC(C)c1cc(C=O)ccc1OCc1ccccc1 → CC(C)c1cc(Br)ccc1OCc1ccccc1.C(=O)N(C)C
c1ccccc1COc1c(C(C)C)cc(C=O)cc1 → c1ccccc1COc1c(C(C)C)cc(Br)cc1.C(=O)N(C)C
c1cccc(COc2c(C(C)C)cc(C=O)cc2)c1 → c1cccc(COc2c(C(C)C)cc(Br)cc2)c1.C(=O)N(C)C

Leverage atom mapping information to reduce edit distance between two sides



Approach #1: end-to-end Transformer

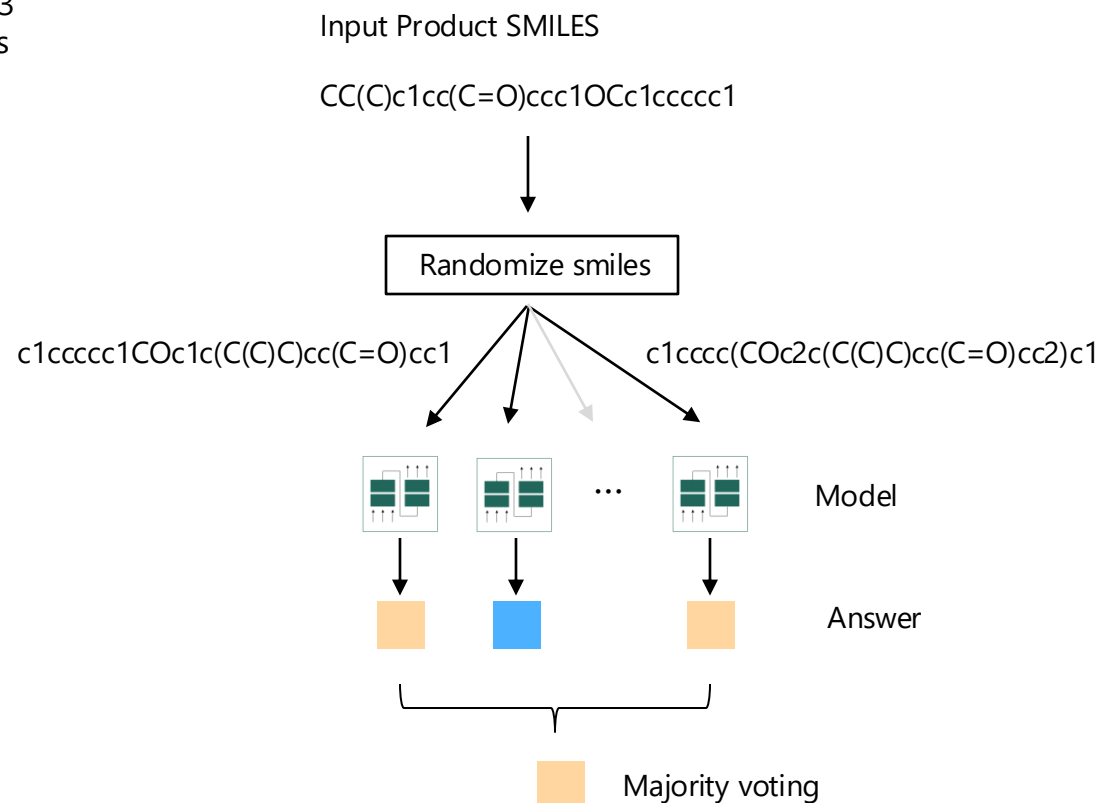
3. Optimized decoding for retrosynthesis prediction



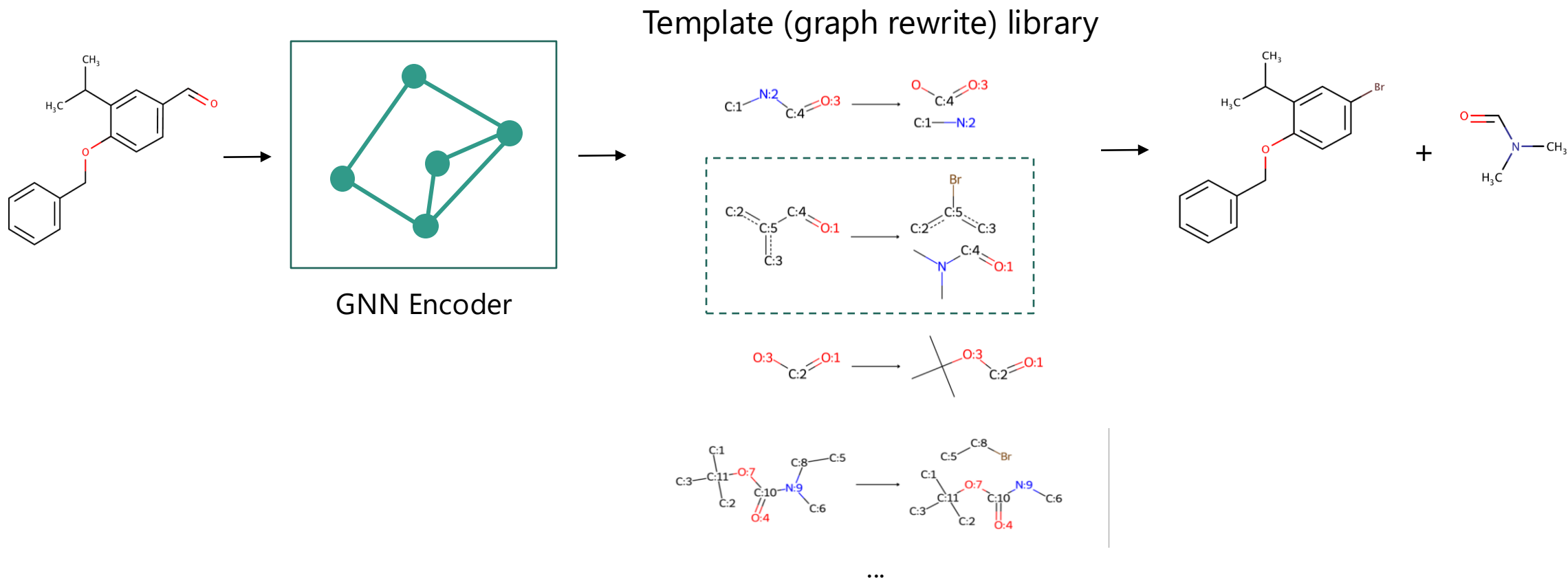
Beam Search Strategy 1 (utilized by OpenNMT): Maintain a pool to save ended sequences during search, search are done until two conditions are met: 1) the pool size = beam size; 2) the top-rated sequence in the beam has lower probability than all candidates in the pool

Beam Search Strategy 2: Keep ended sequences within the beam itself, search are done when each sequence in the beam finishes with the EOS token

4. Test-time augmentation

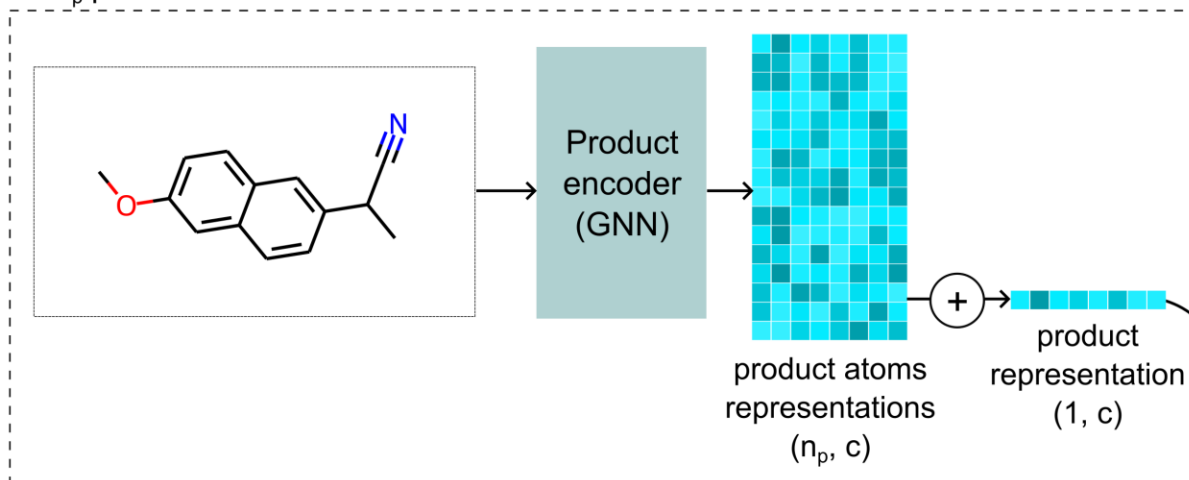


Approach #2: template-based GNN

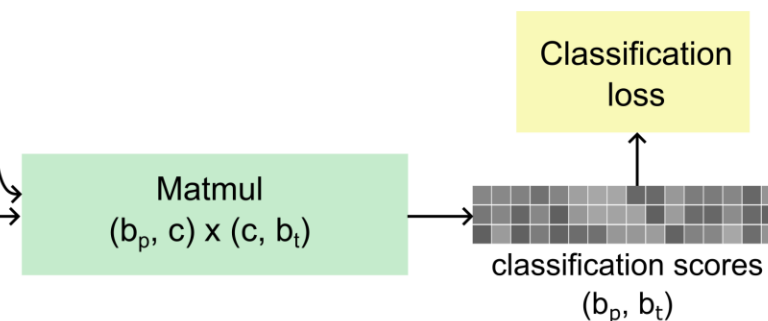
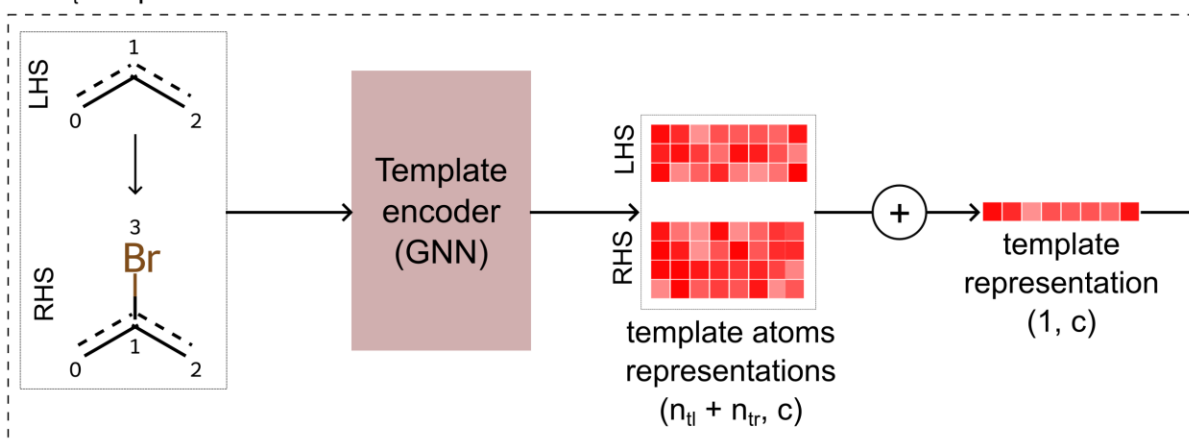


Approach #2: template-based GNN

for b_p products in batch

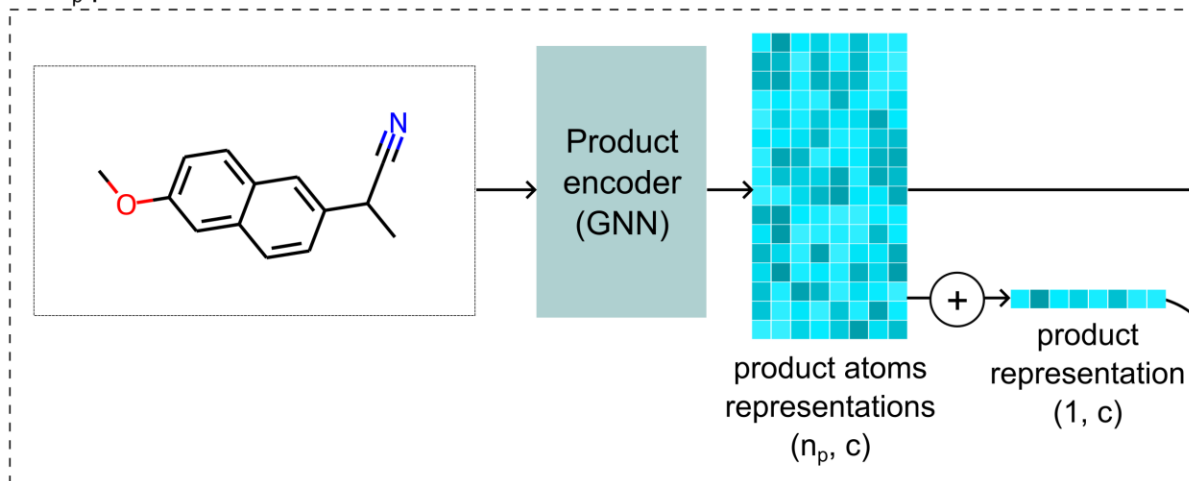


for b_t templates in batch

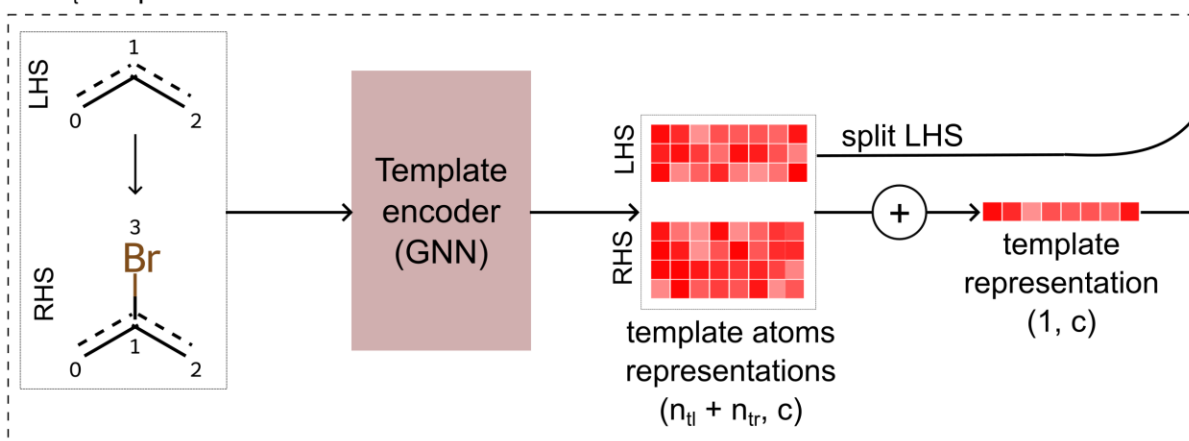


Approach #2: template-based GNN

for b_p products in batch



for b_t templates in batch



Matmul (b_p times)
 $(n_p, c) \times (c, \sum n_{tl})$



localization scores
 $b_p \times (n_p, \sum n_{tl})$

Localization
loss

Matmul
 $(b_p, c) \times (c, b_t)$

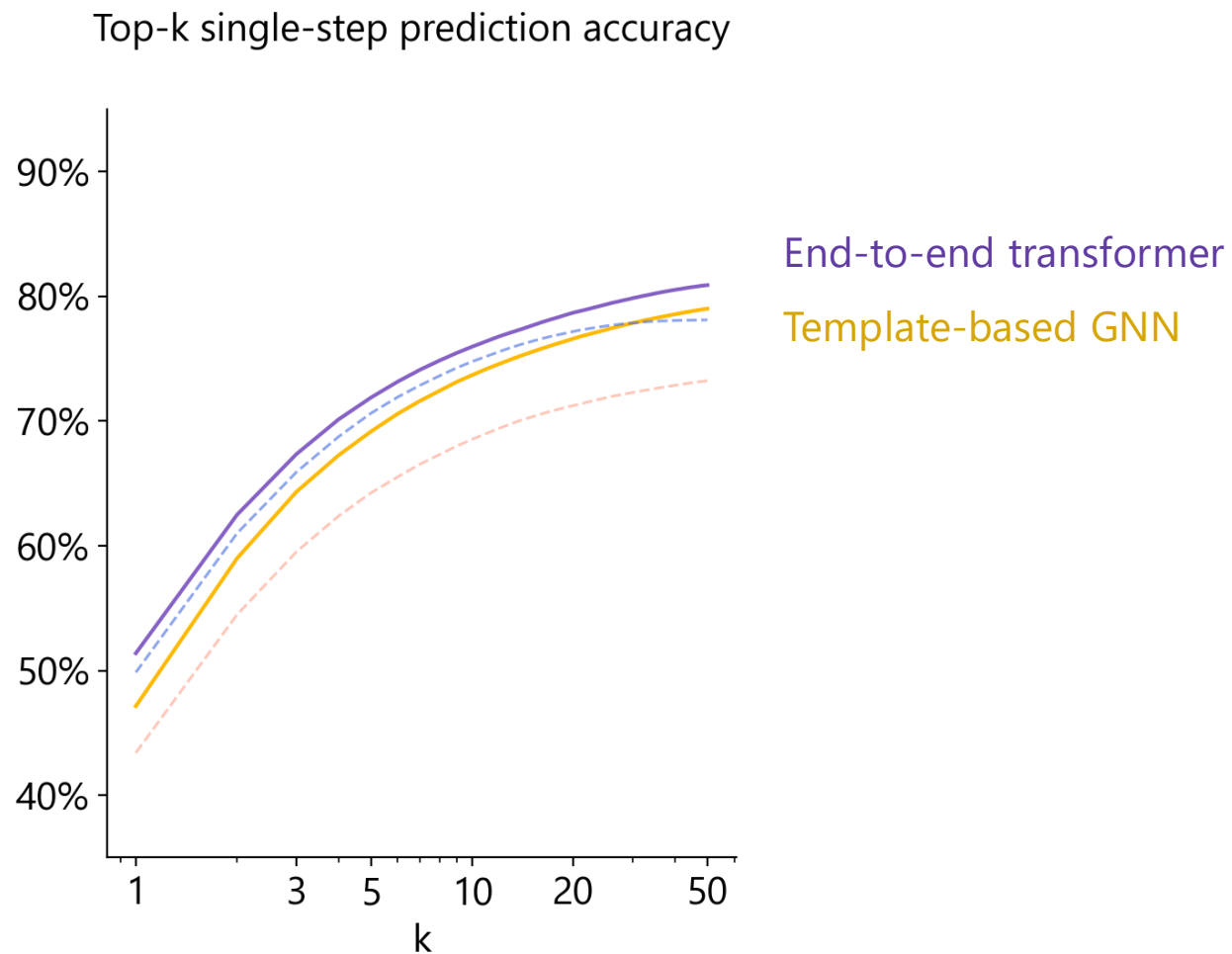


classification scores
 (b_p, b_t)

Classification
loss

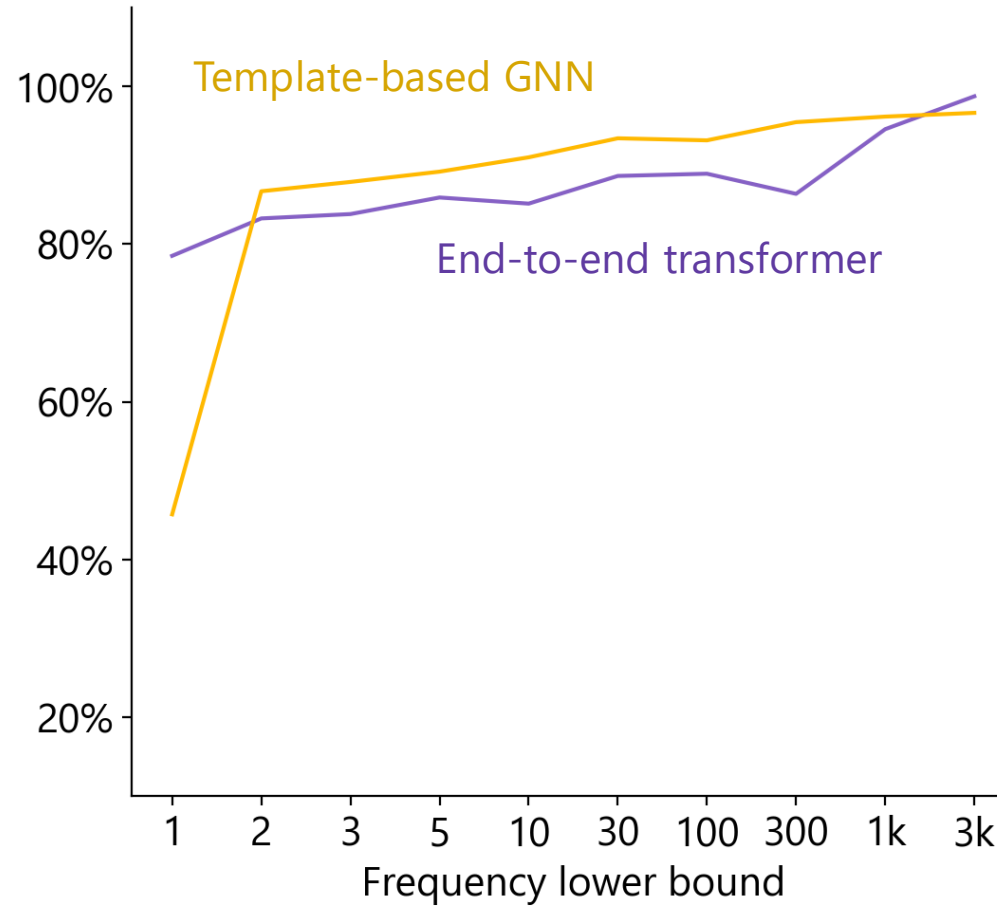
Comparison

aka.ms/RetroChimeraNeurIPS



Comparison: trade-offs

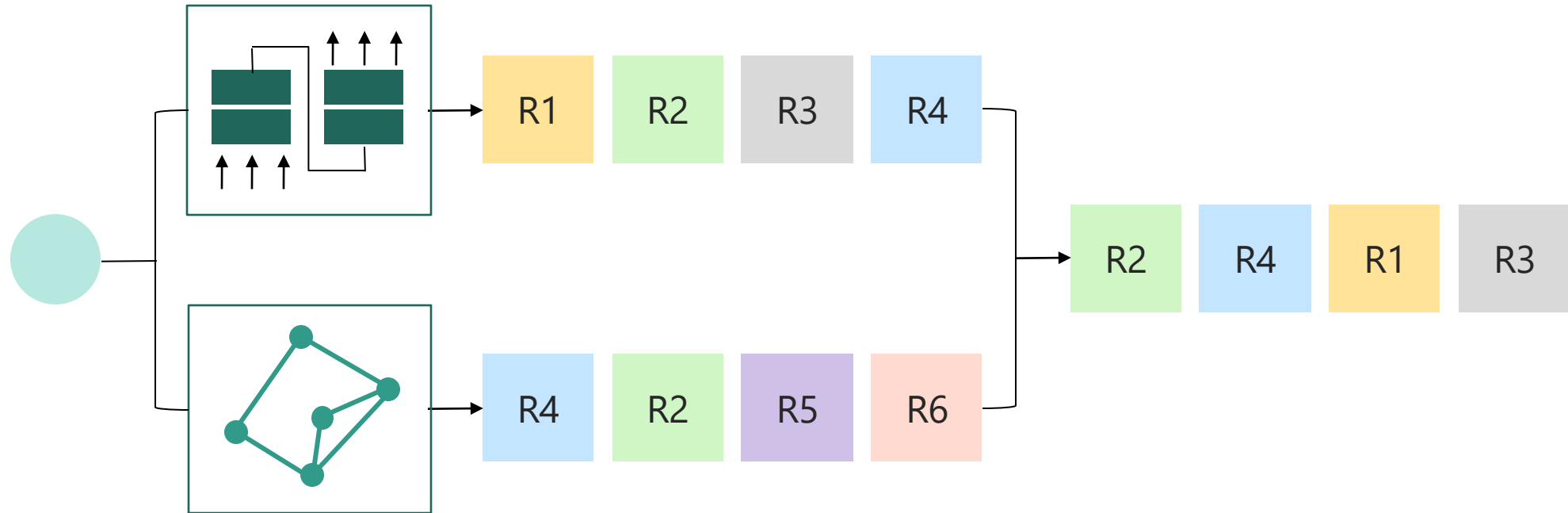
Top-50 accuracy vs template support



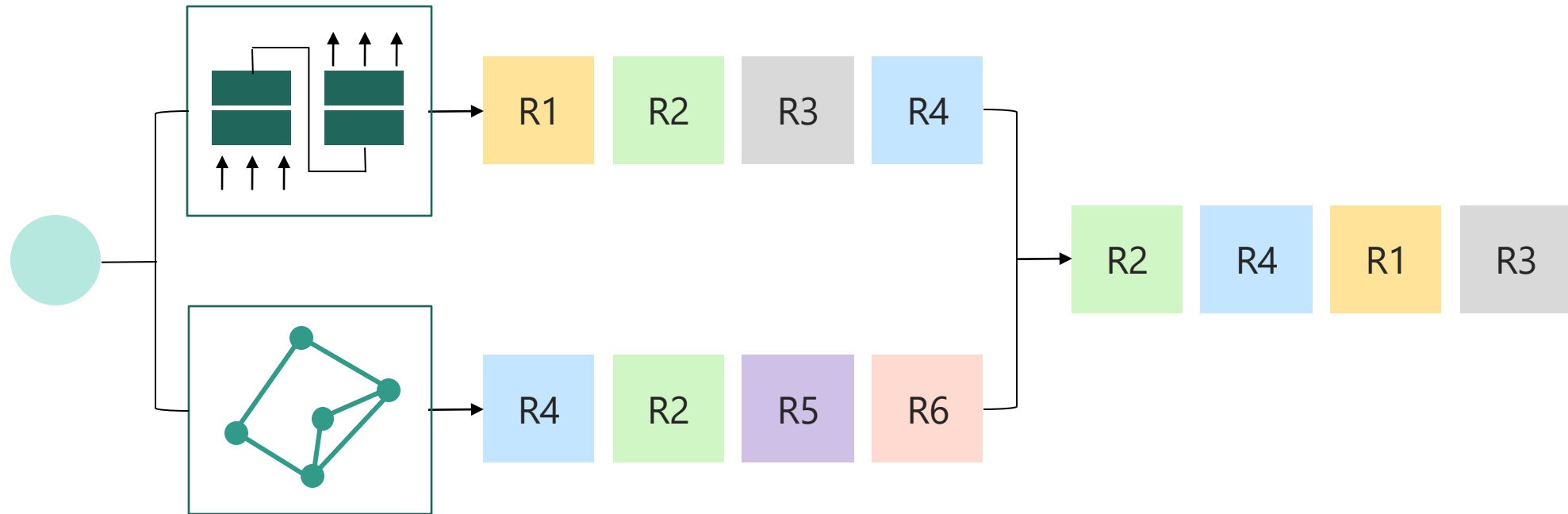
How does RetroChimera fuse results from its submodels?



Learned ensembling



Learned ensembling

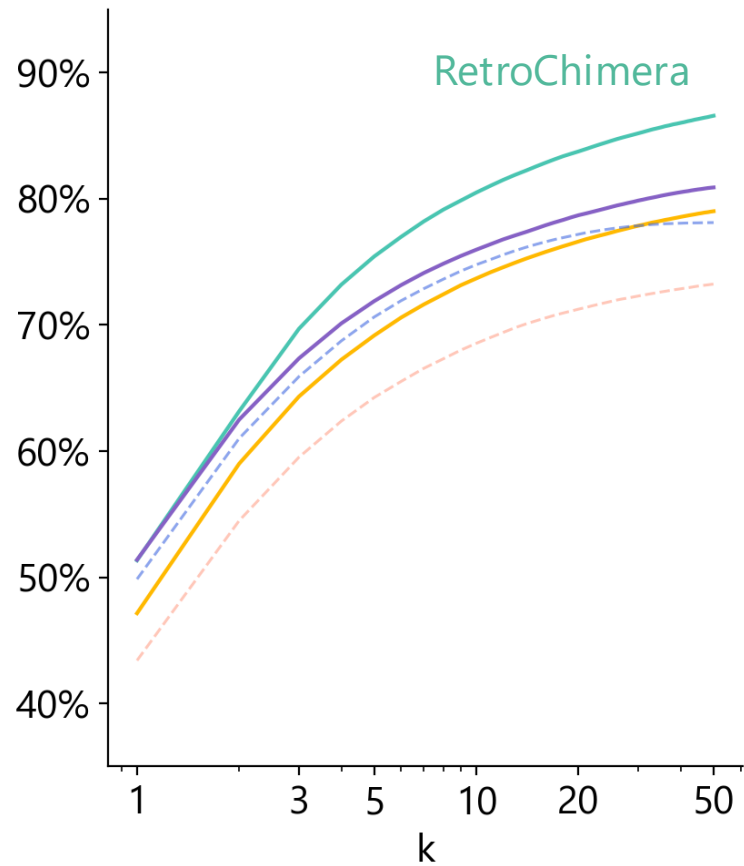


$$\text{score}(r) = \sum_{i=1}^m \sum_{k=1}^{k_{\max}} \mathbb{1}[r = r_{i,k}] \cdot \theta_{i,k}$$

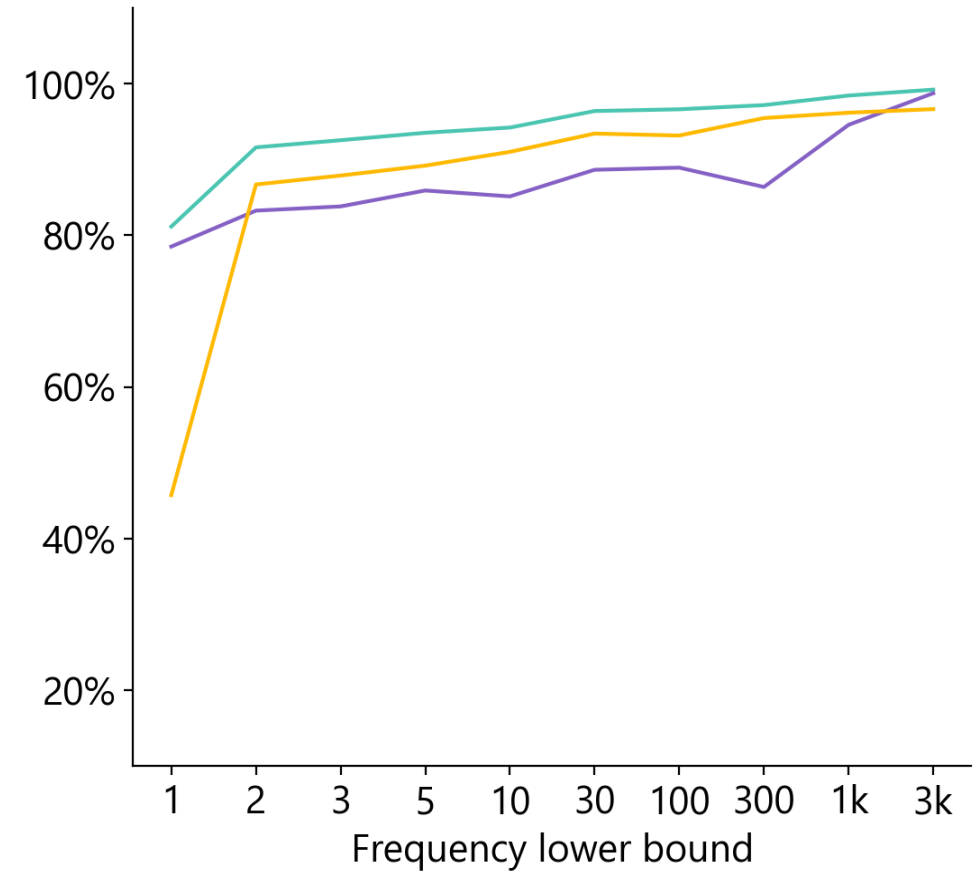
$$\mathcal{L}_{\text{rank}}(r^+, r^-) = \sigma \left(\frac{\text{score}(r^-) - \text{score}(r^+) + \epsilon}{T} \right)$$

RetroChimera's performance

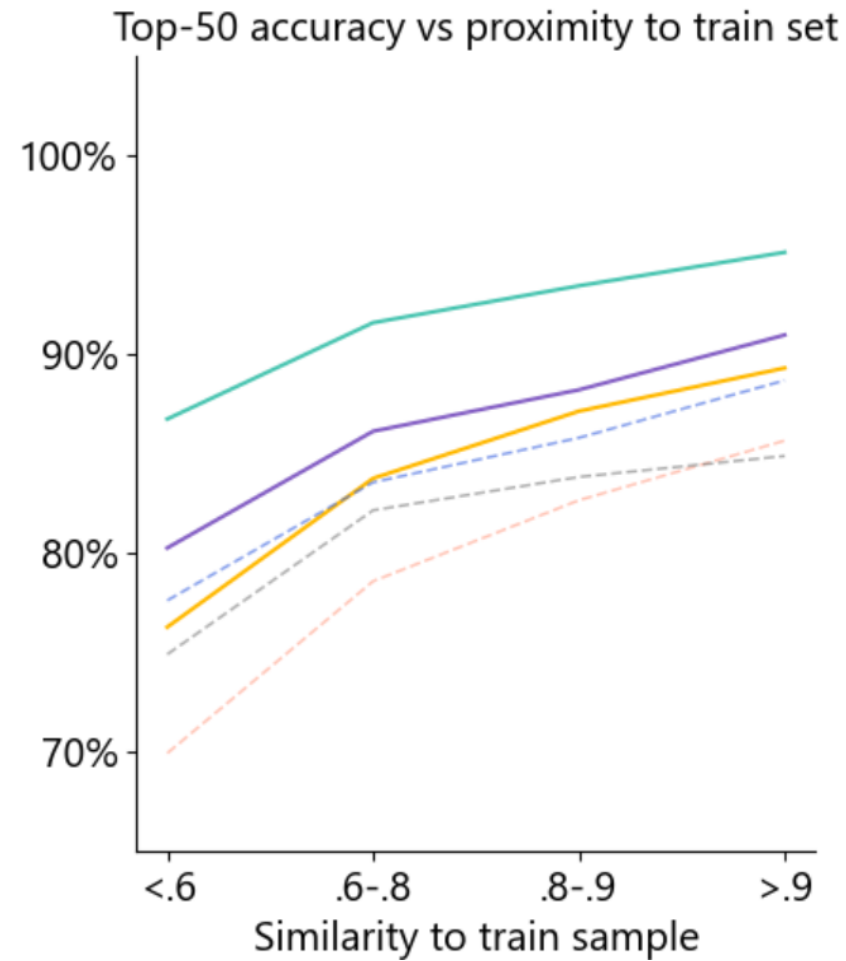
Top-k single-step prediction accuracy



Top-50 accuracy vs template support



Performance is robust OOD

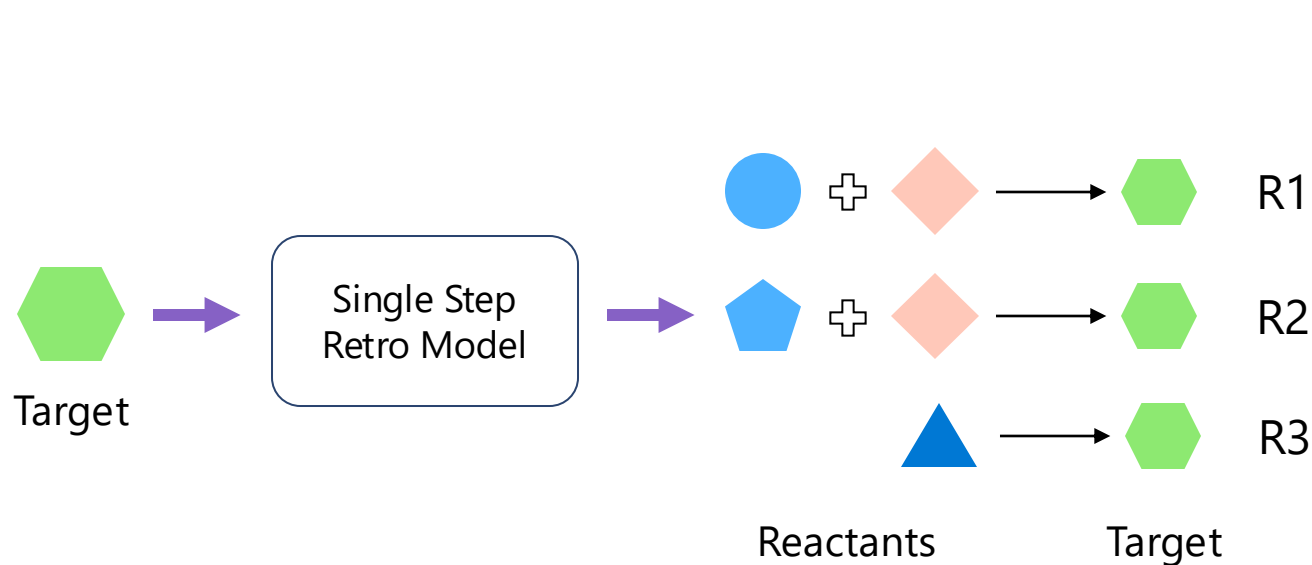


How to combine individual
reactions into full plans?

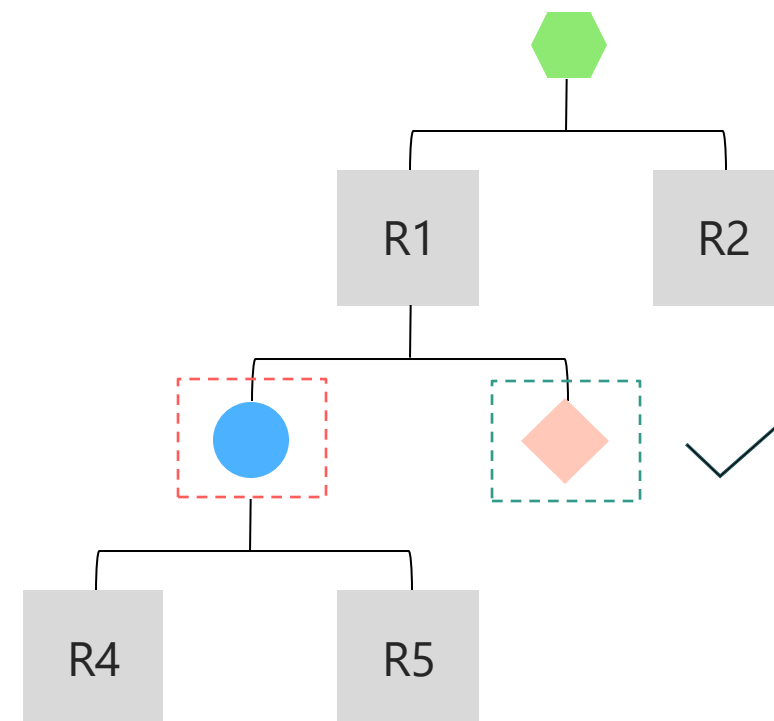


Synthesis prediction models and search

(1) Given target, predict different possible reactants

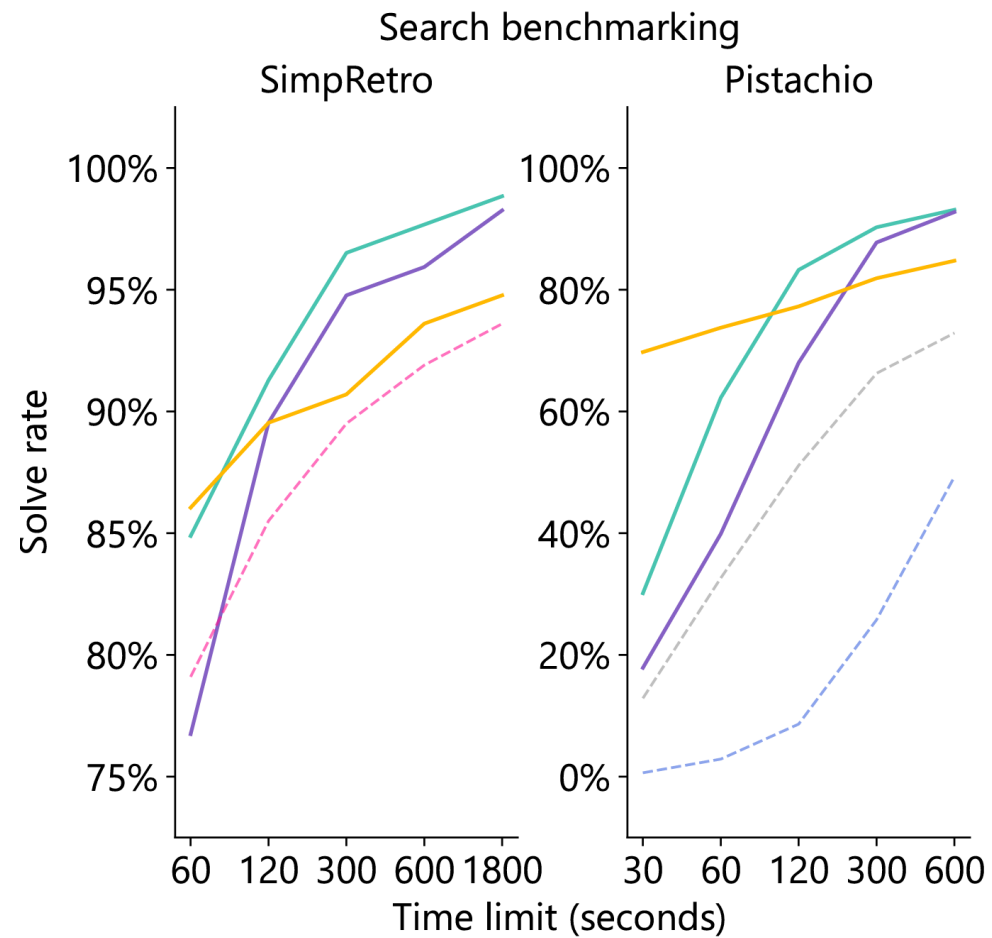


(2) Apply model recursively to obtain routes

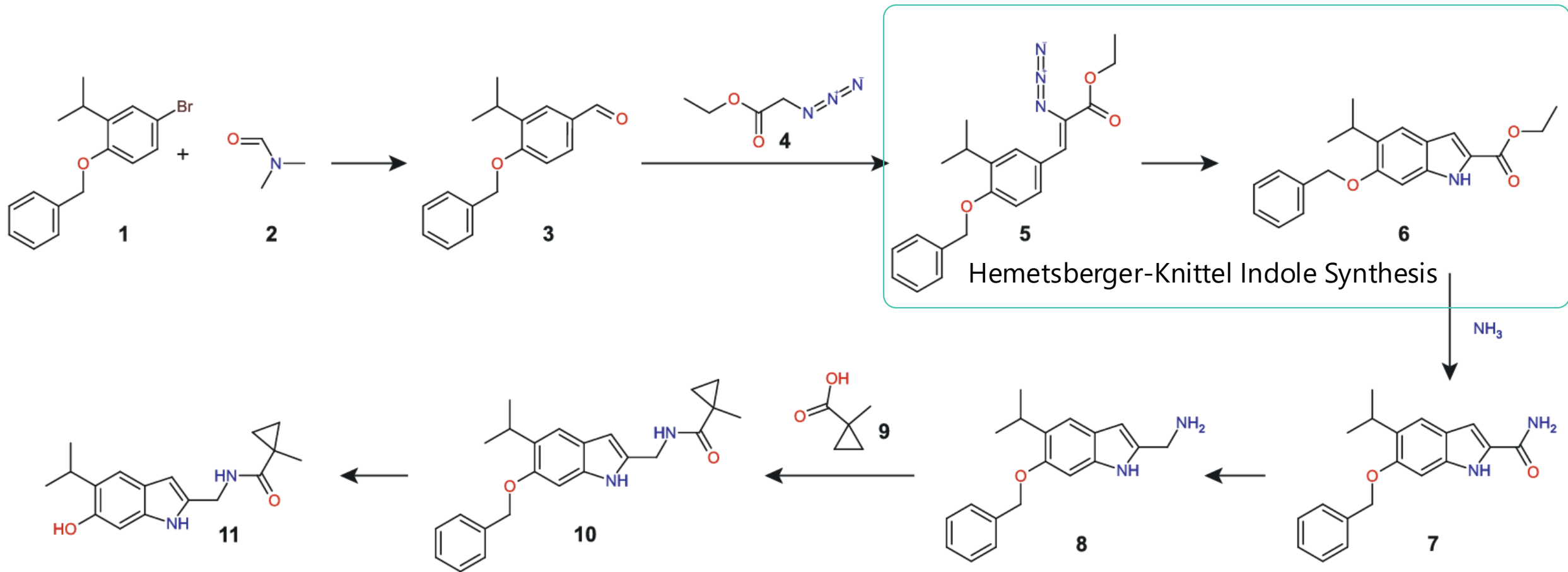


(2) Maziarz et al, "Re-evaluating Retrosynthesis Algorithms with Syntheseus", Faraday Discuss. 2024
Liu et al, "Retrosynthetic Planning with Dual Value Networks", ICML 2023

Search with RetroChimera



Example synthesis plan

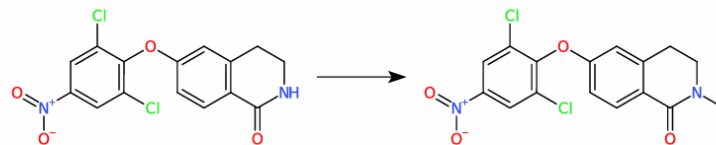
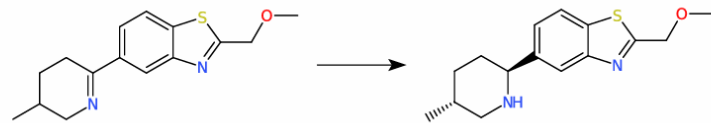
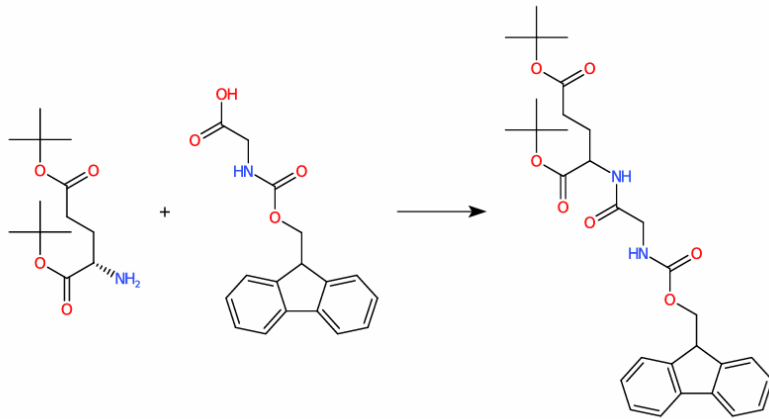


RetroChimera can find non-trivial routes

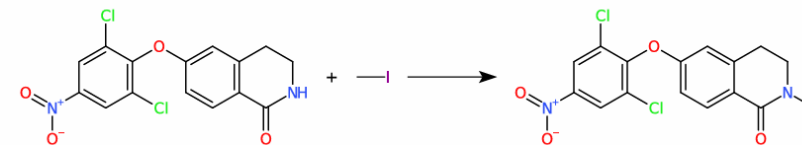
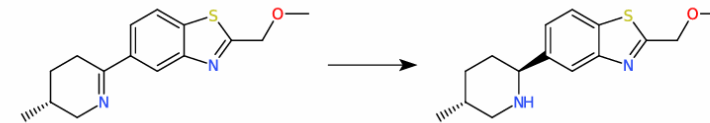
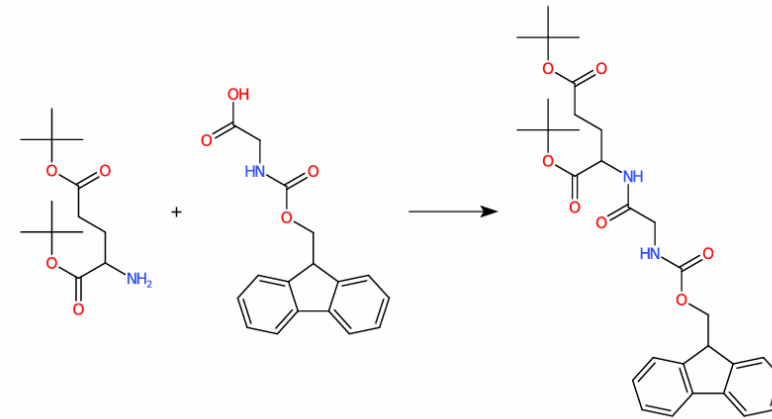
Exploring RetroChimera's predictions qualitatively



Denoising



original patent reaction

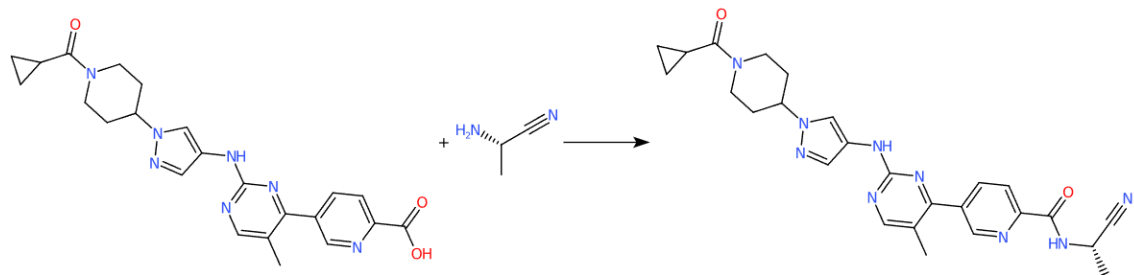


denoised by RetroChimera

Pairwise (blind) expert comparisons

Please pick your preferred reaction.

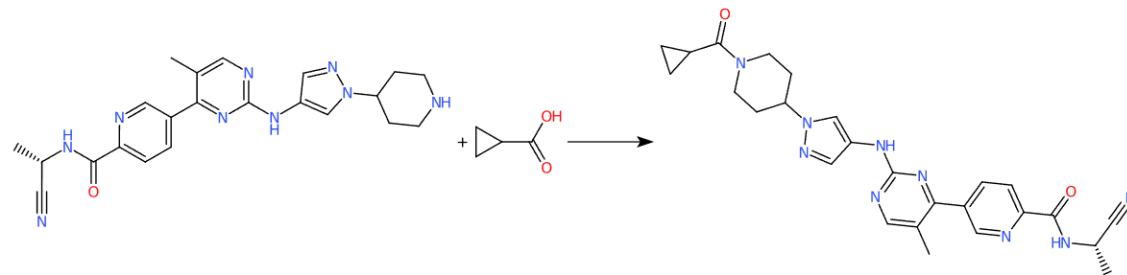
Model 1:



Model 1 better

Model 2 better

Model 2:



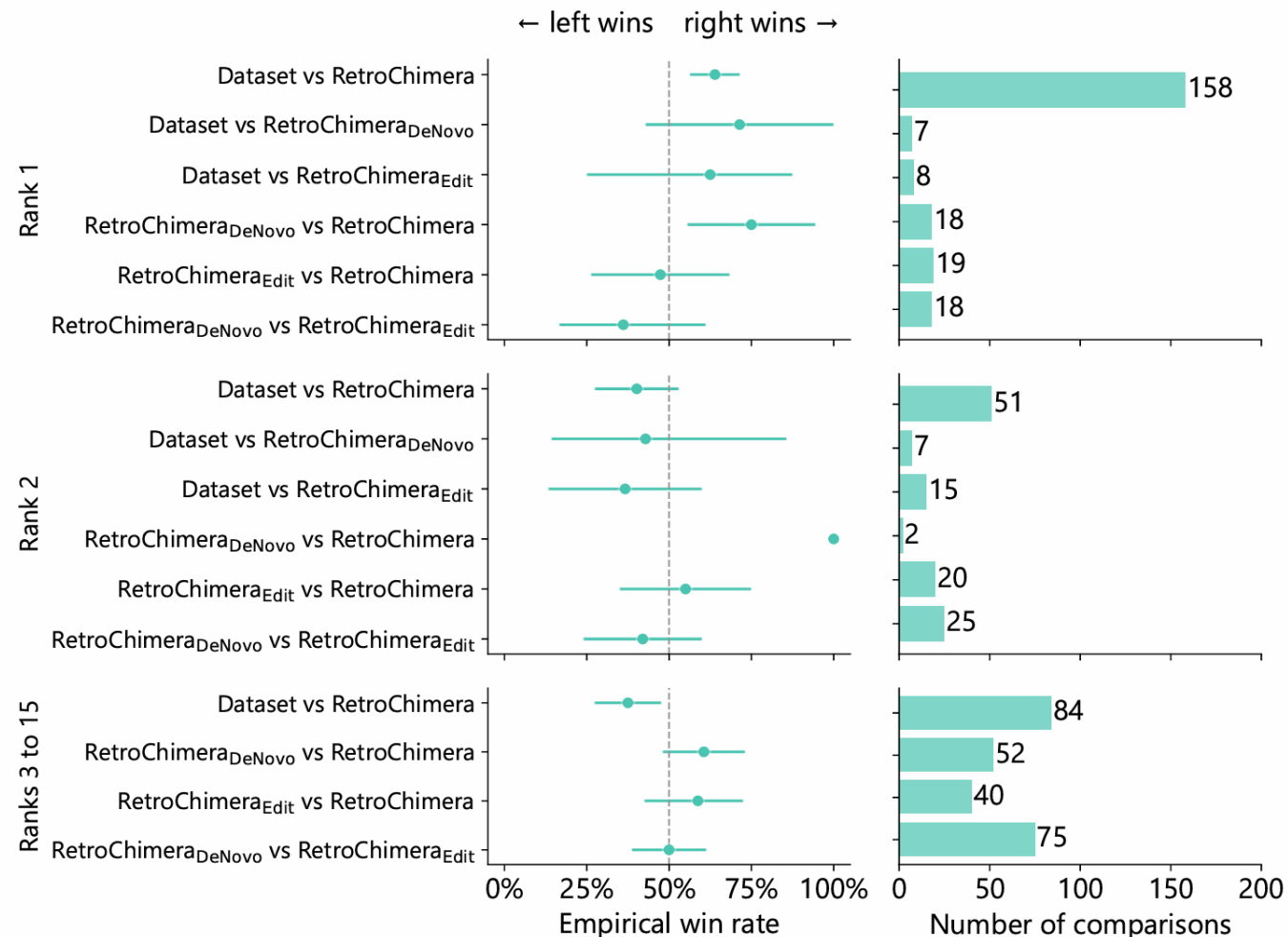
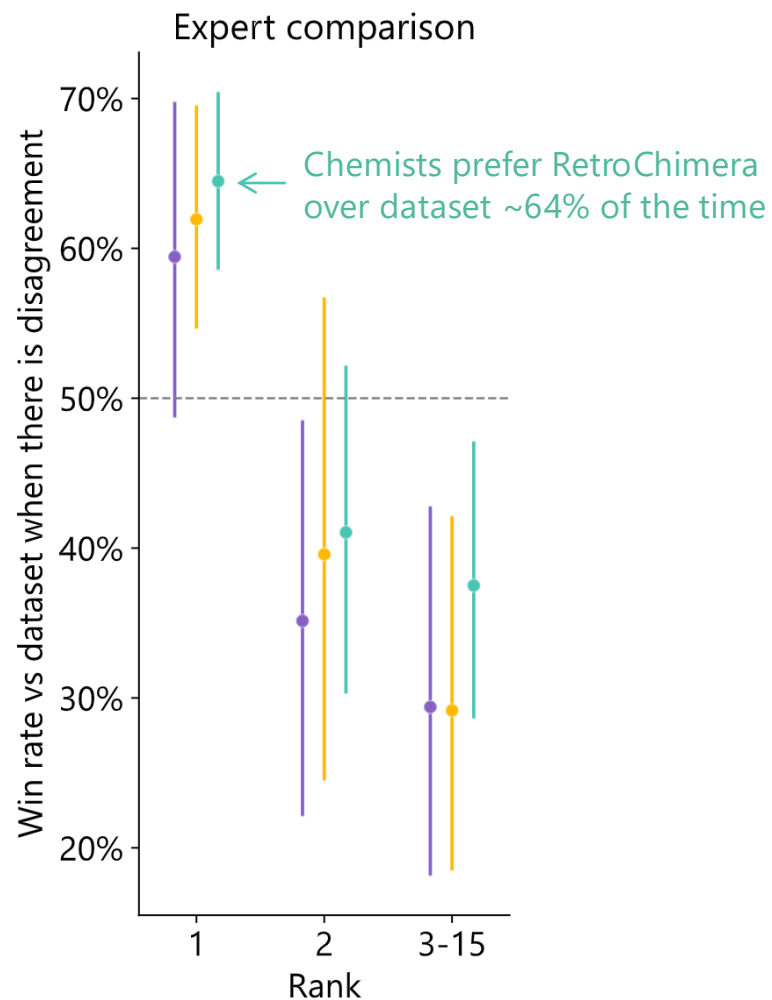
Both good

Both bad

Sample ID: ed2e9b78-bed5-43a5-b399-928248f08c2f

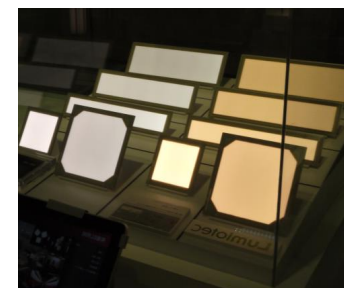
Experiment ID: 12

Preferences in pairwise expert comparisons



Questions

Our goal: Help chemists discover new essential molecules with predictive synthesis



aka.ms/RetroChimeraNeurIPS

aka.ms/RetroChimeraPaper

aka.ms/RetroChimeraCode (github.com/microsoft/retrochimera)

aka.ms/CondPredPaper (cond. pred. model QFANG, coming soon)

github.com/microsoft/syntheseus