

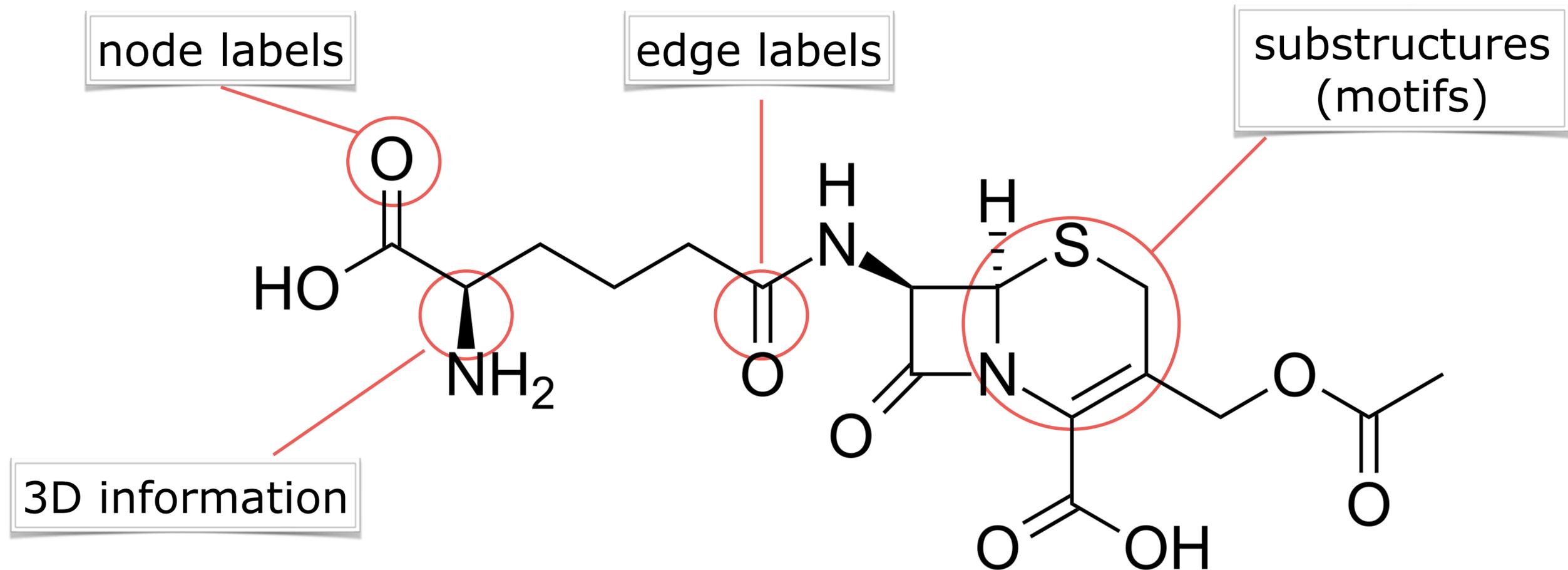
Representation and Generation of Molecular Graphs

Wengong Jin
MIT CSAIL

in collaboration with
Tommi Jaakkola, Regina Barzilay, Kevin Yang, Kyle Swanson

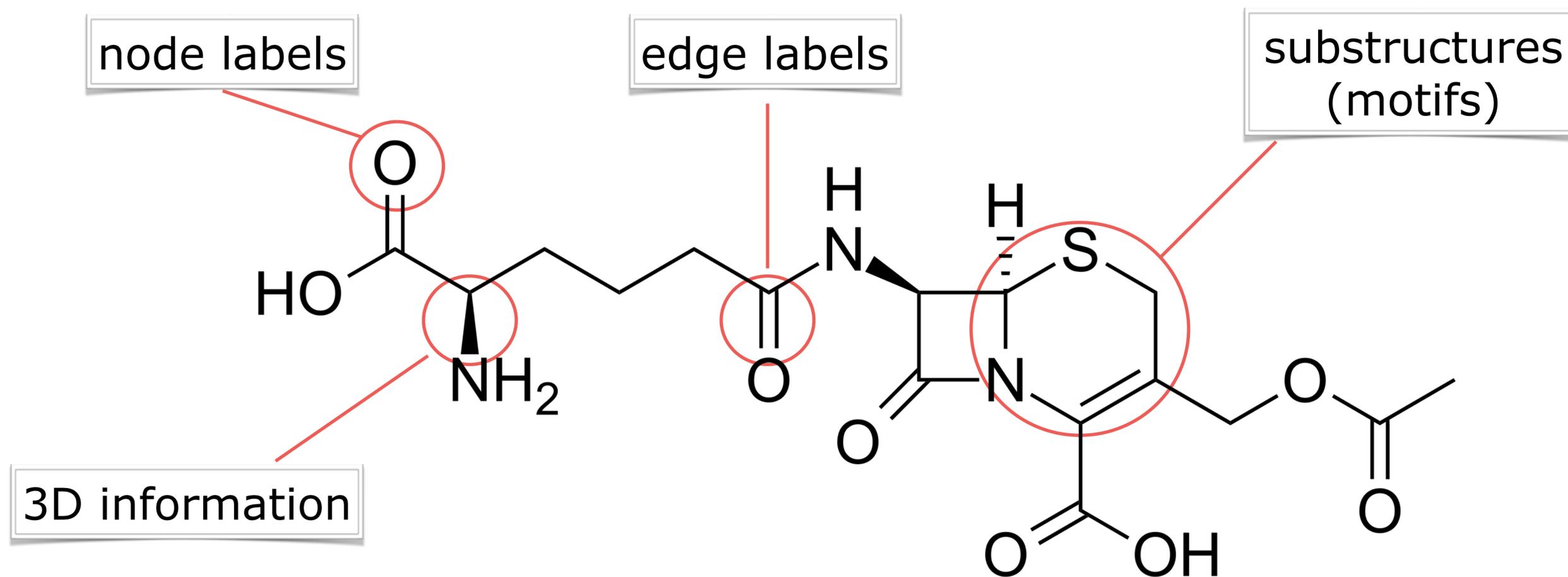
Why are molecules interesting for ML?

- ▶ E.g., antibiotic (cephalosporin)



Why are molecules interesting for ML?

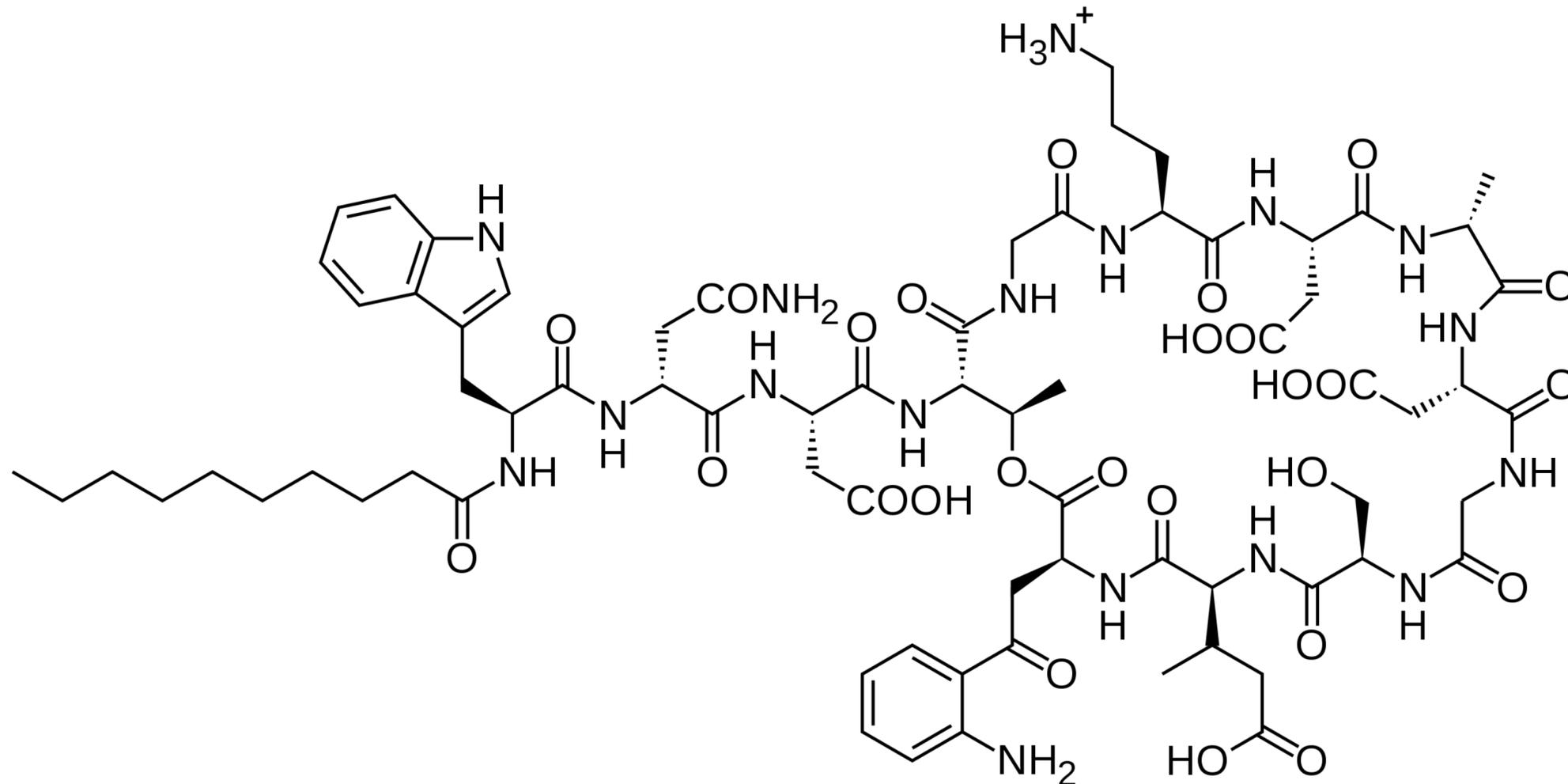
- ▶ E.g., antibiotic (cephalosporin)



Together give rise to various chemical properties
(e.g., solubility, toxicity, ...)

Why are molecules interesting for ML?

- ▶ Properties may depend on intricate structures;
- ▶ The key challenges are to automatically predict chemical properties and to generate molecules with desirable characteristics



(Daptomycin antibiotic)

Interesting ML Problems

- ▶ **Deeper into known chemistry**

- extract chemical knowledge from journals, notebooks (NLP)

- ▶ **Deeper into drug design**

- molecular property prediction (graph representation)
 - (multi-criteria) lead optimization (graph generation)

- ▶ **Deeper into reactions**

- forward reaction prediction (structured prediction)
 - forward reaction optimization (combinatorial optimization)

- ▶ **Deeper into synthesis**

- retrosynthesis planning (reinforcement learning)

Interesting ML Problems

- ▶ **Deeper into known chemistry**

- extract chemical knowledge from journals, notebooks (NLP)

- ▶ **Deeper into drug design**

- molecular property prediction (graph representation)
 - (multi-criteria) lead optimization (graph generation)

- ▶ **Deeper into reactions**

- forward reaction prediction (structured prediction)
 - forward reaction optimization (combinatorial optimization)

- ▶ **Deeper into synthesis**

- retrosynthesis planning (reinforcement learning)

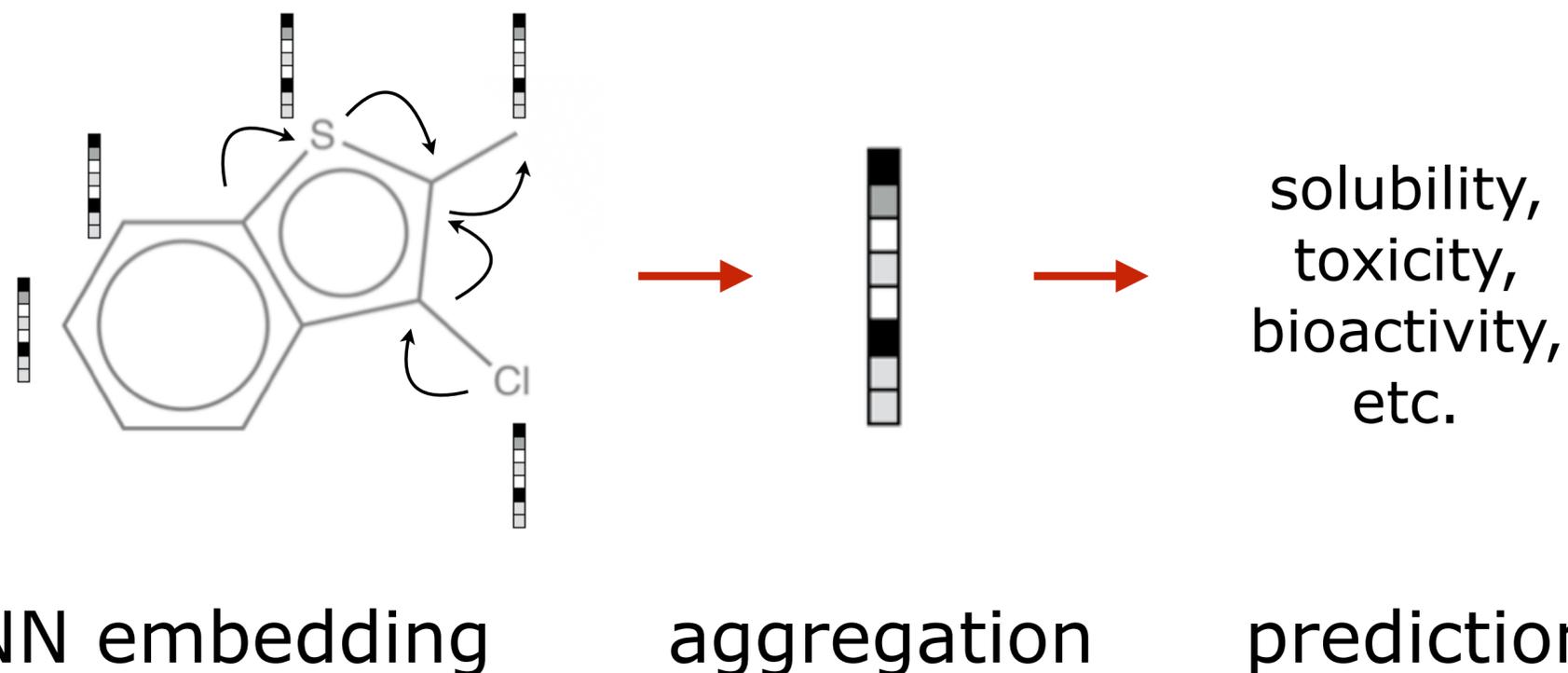
Automating Drug design

▸ Key challenges:

1. **representation and prediction:** learn to predict molecular properties
2. **generation and optimization:** realize target molecules with better properties programmatically
3. **understanding:** uncover principles (or diagnose errors) underlying complex predictions

GNNs for property prediction?

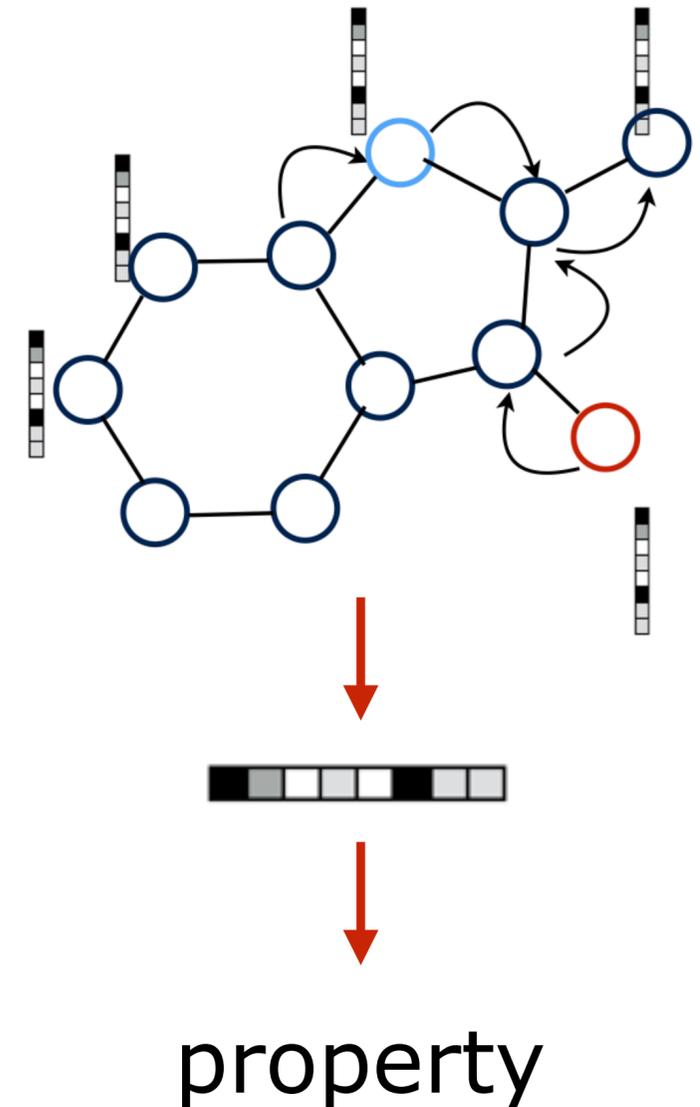
- ▶ Are GNN models operating on molecular graphs sufficiently expressive for predicting molecular properties (in the presence of “property cliffs”)?



- ▶ A number of recent results pertaining to the power of GNNs (e.g., Xu et al. 2018, Sato et al. 2019, Maron et al., 2019, ...);

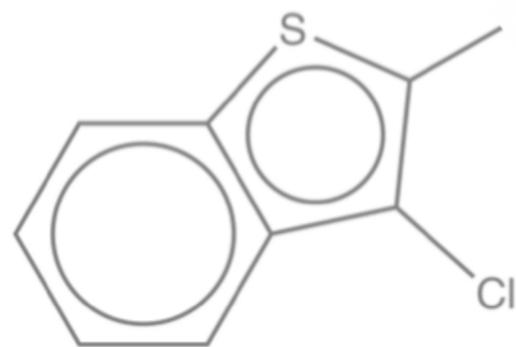
Are basic GNNs sufficiently expressive?

- ▶ **Theorem [Garg et al., 2019]:** GNNs with permutation invariant readout functions cannot “decide”
 - girth (length of the shortest cycle)
 - circumference (length of the longest cycle)
 - diameter, radius
 - presence of conjoint cycle
 - total number of cycles
 - presence of c -clique
 - etc. (?)
- ▶ (most results also apply to MPNNs)



Beyond simple GNNs: sub-structures

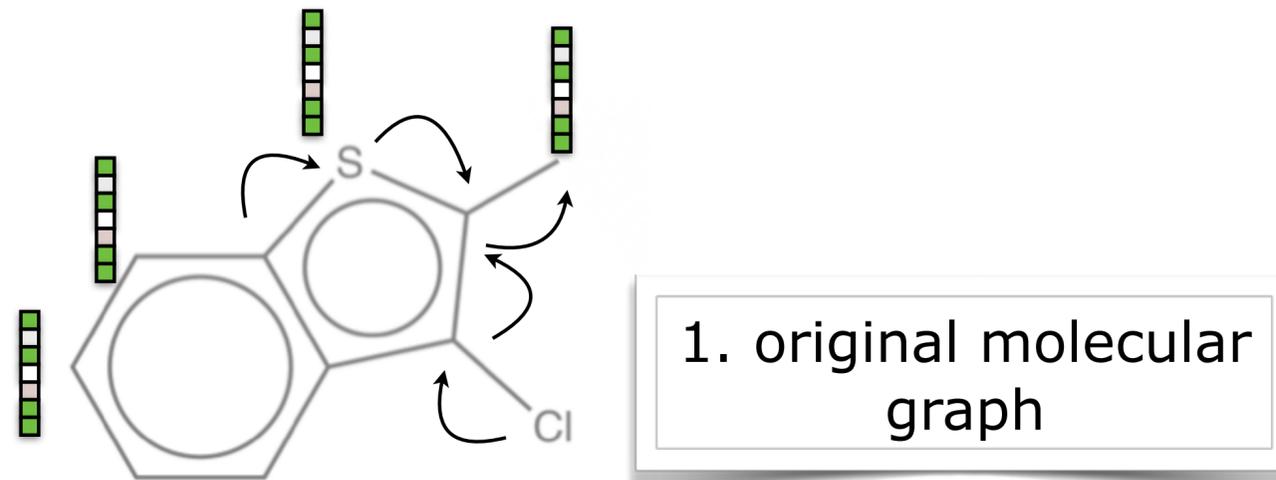
- ▶ Learning to view molecules at multiple levels [Jin et al., 2019]



1. original molecular graph

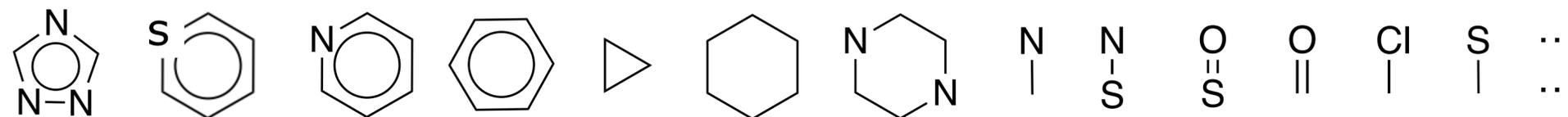
Beyond simple GNNs: sub-structures

- ▶ Learning to view molecules at multiple levels

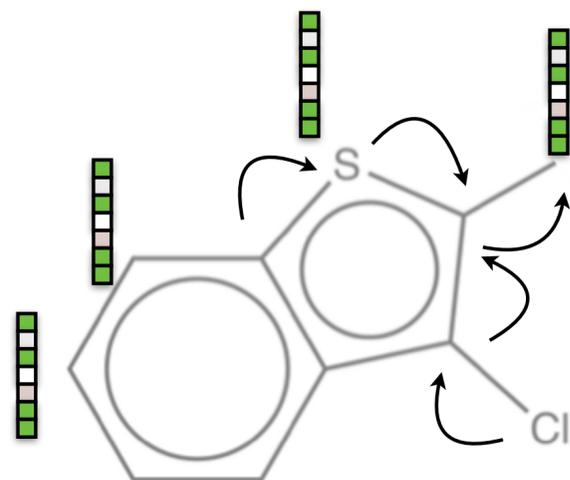


Beyond simple GNNs: sub-structures

- ▶ Learning to view molecules at multiple levels



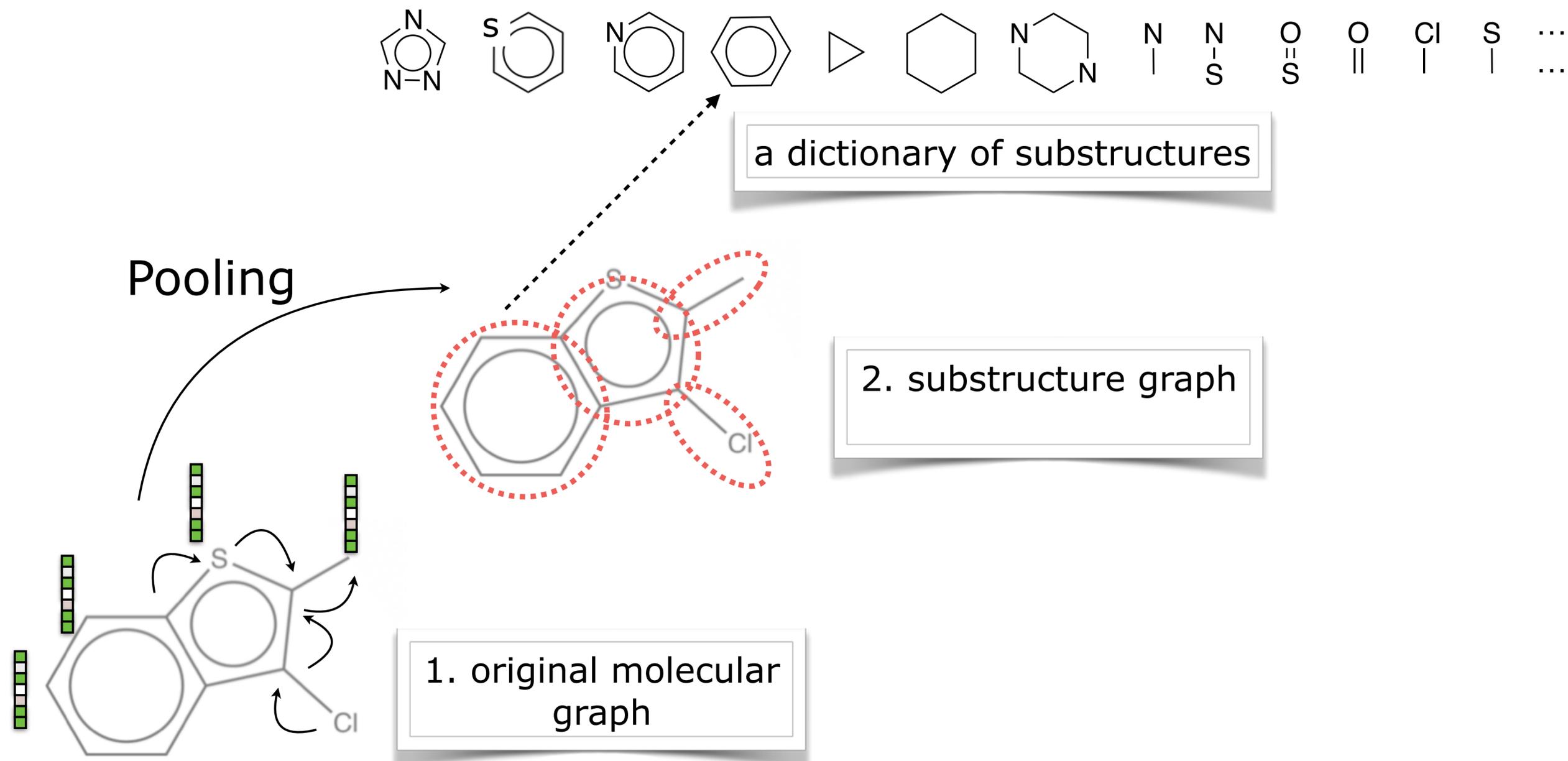
a dictionary of substructures



1. original molecular graph

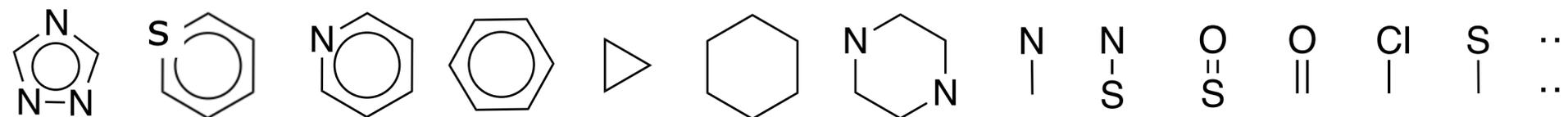
Beyond simple GNNs: sub-structures

- ▶ Learning to view molecules at multiple levels

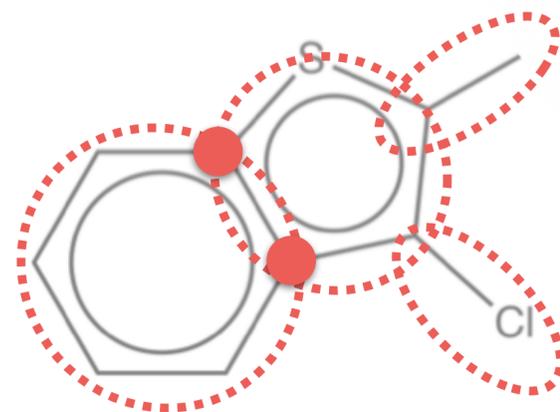


Beyond simple GNNs: sub-structures

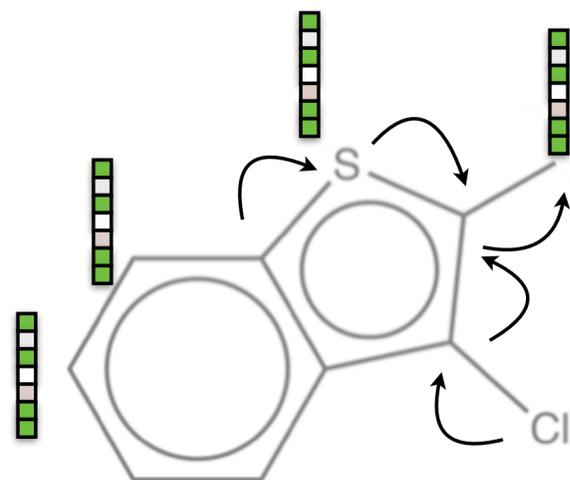
- ▶ Learning to view molecules at multiple levels



a dictionary of substructures



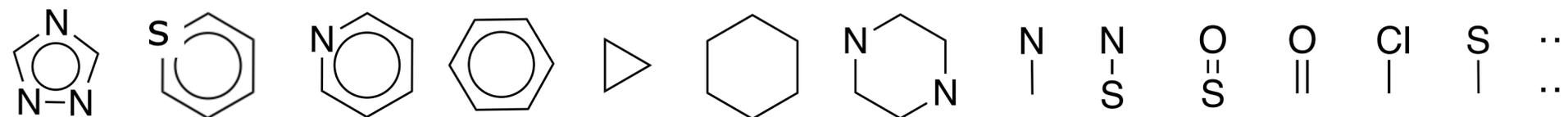
2. substructure graph with attachments



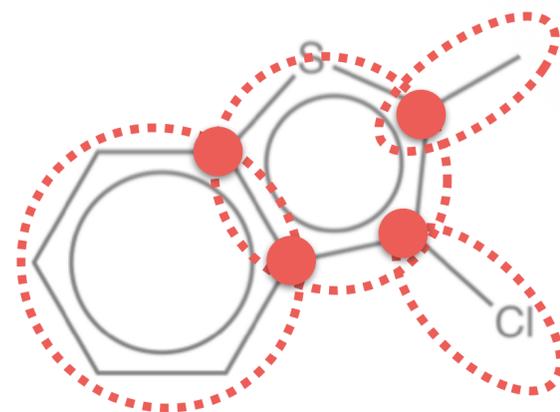
1. original molecular graph

Beyond simple GNNs: sub-structures

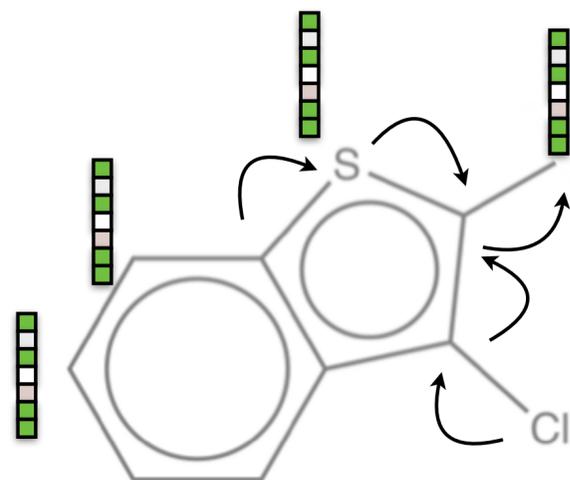
- ▶ Learning to view molecules at multiple levels



a dictionary of substructures



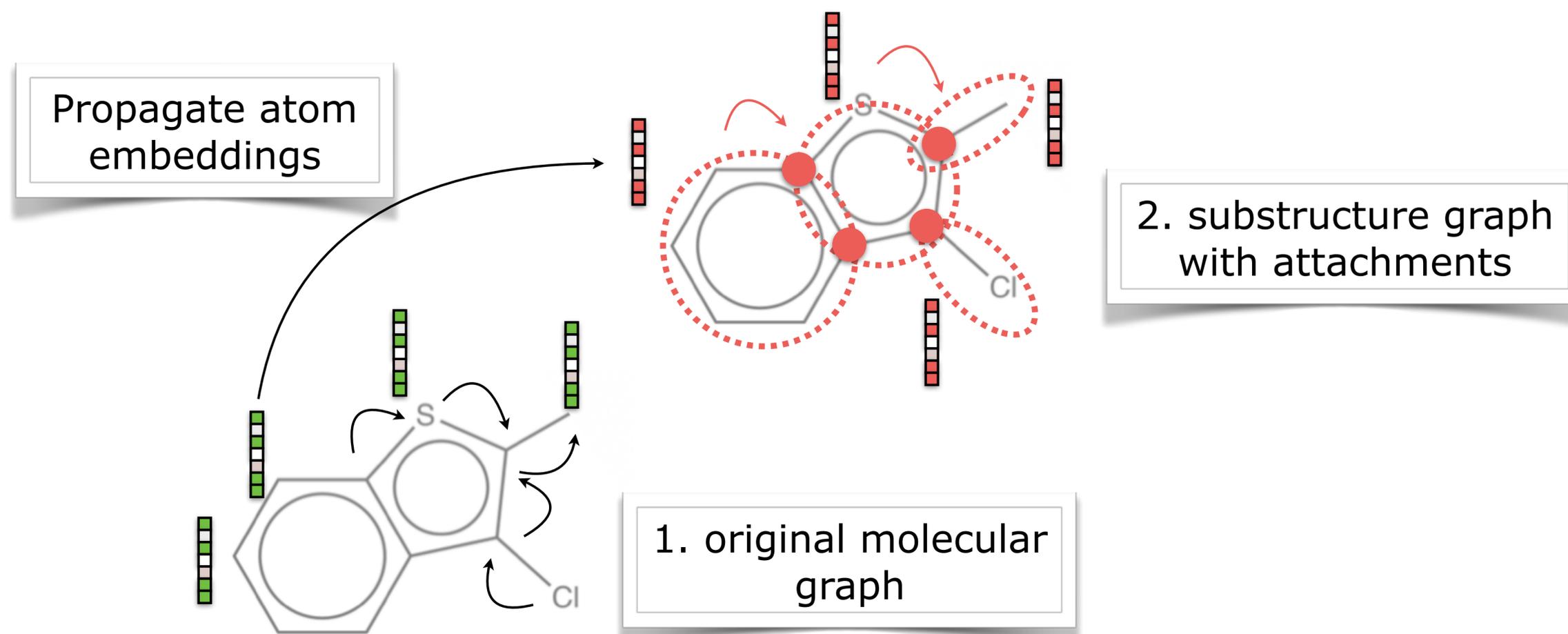
2. substructure graph with attachments



1. original molecular graph

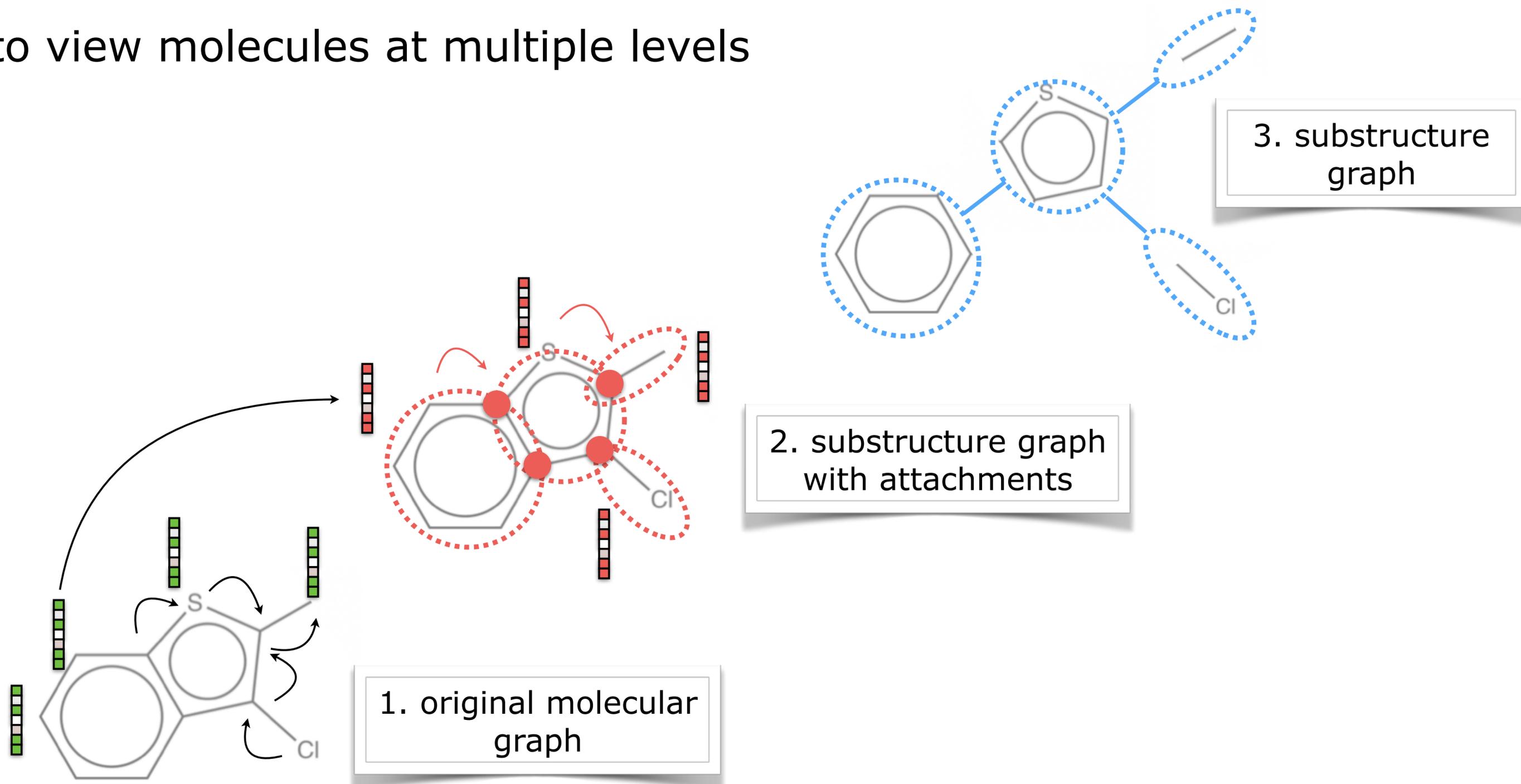
Beyond simple GNNs: sub-structures

- ▶ Learning to view molecules at multiple levels



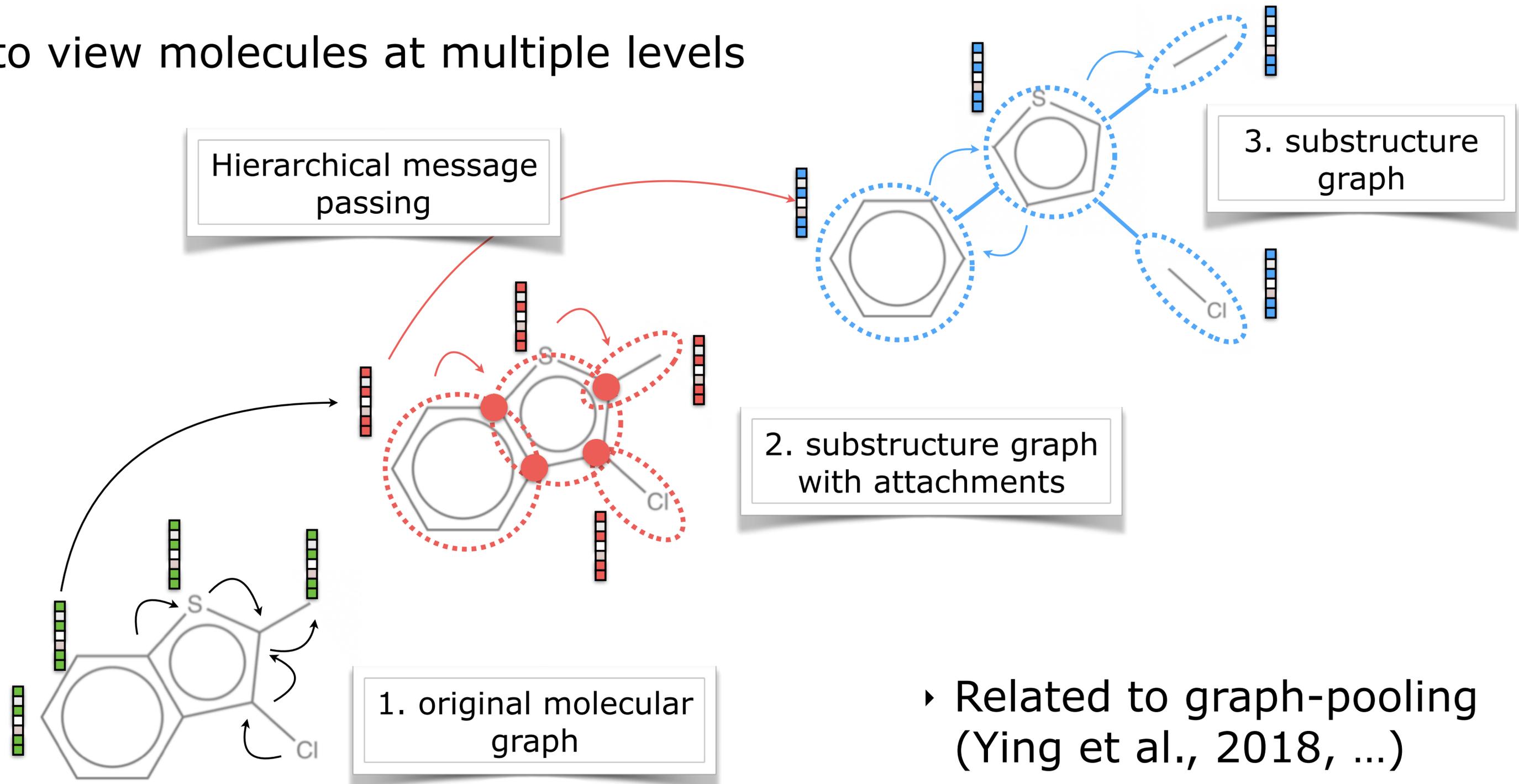
Beyond simple GNNs: sub-structures

- ▶ Learning to view molecules at multiple levels



Multi-resolution representations

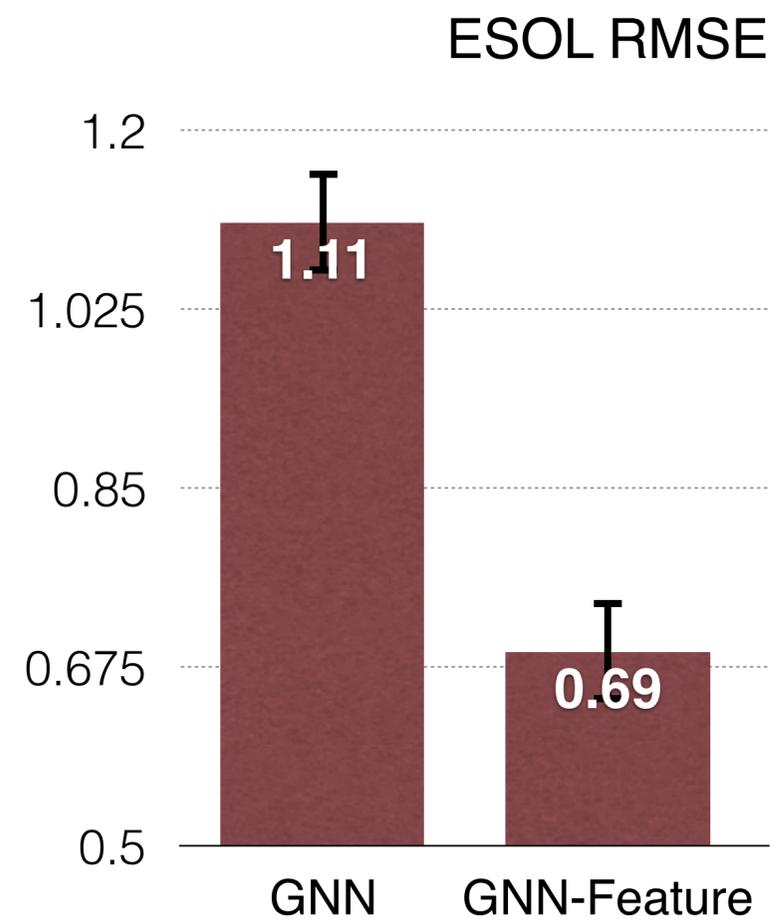
- ▶ Learning to view molecules at multiple levels



- ▶ Related to graph-pooling (Ying et al., 2018, ...)

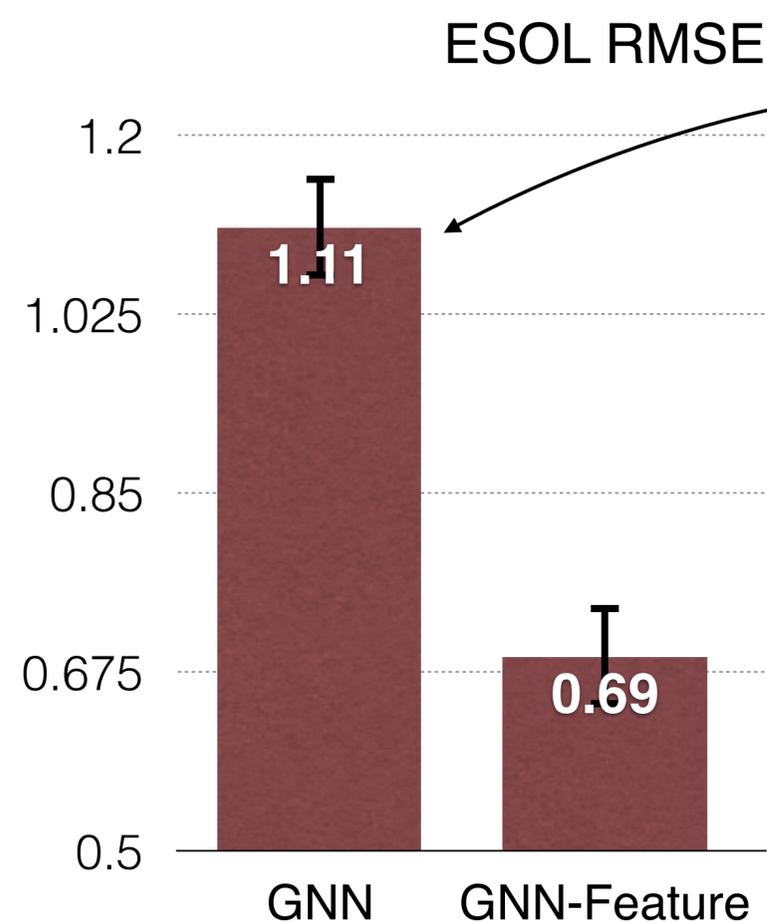
Experiments on solubility

- ▶ ESOL dataset (averaged over 5 folds)



Experiments on solubility

- ▶ ESOL dataset (averaged over 5 folds)

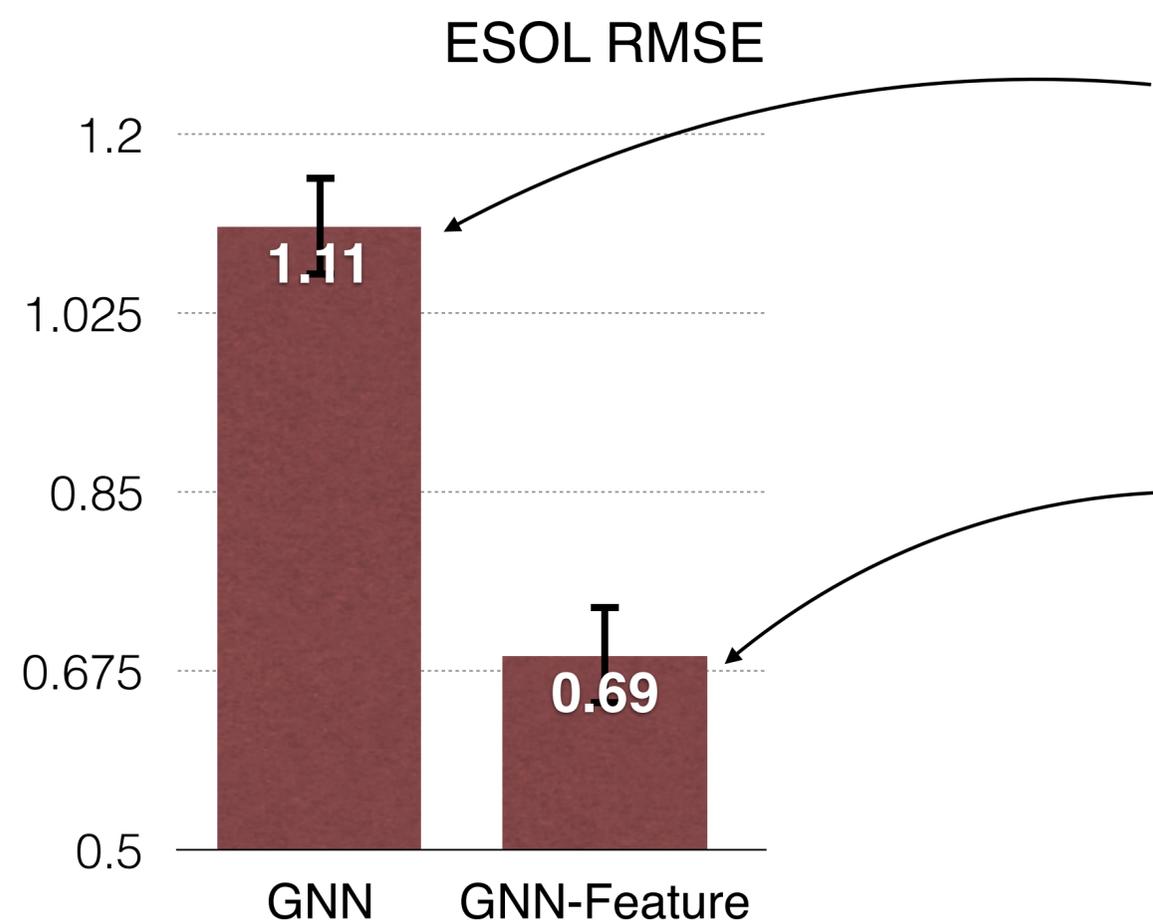


Raw GNN

- ▶ atom feature: only atom type label

Experiments on solubility

- ▶ ESOL dataset (averaged over 5 folds)



Raw GNN

- ▶ atom feature: only atom type label

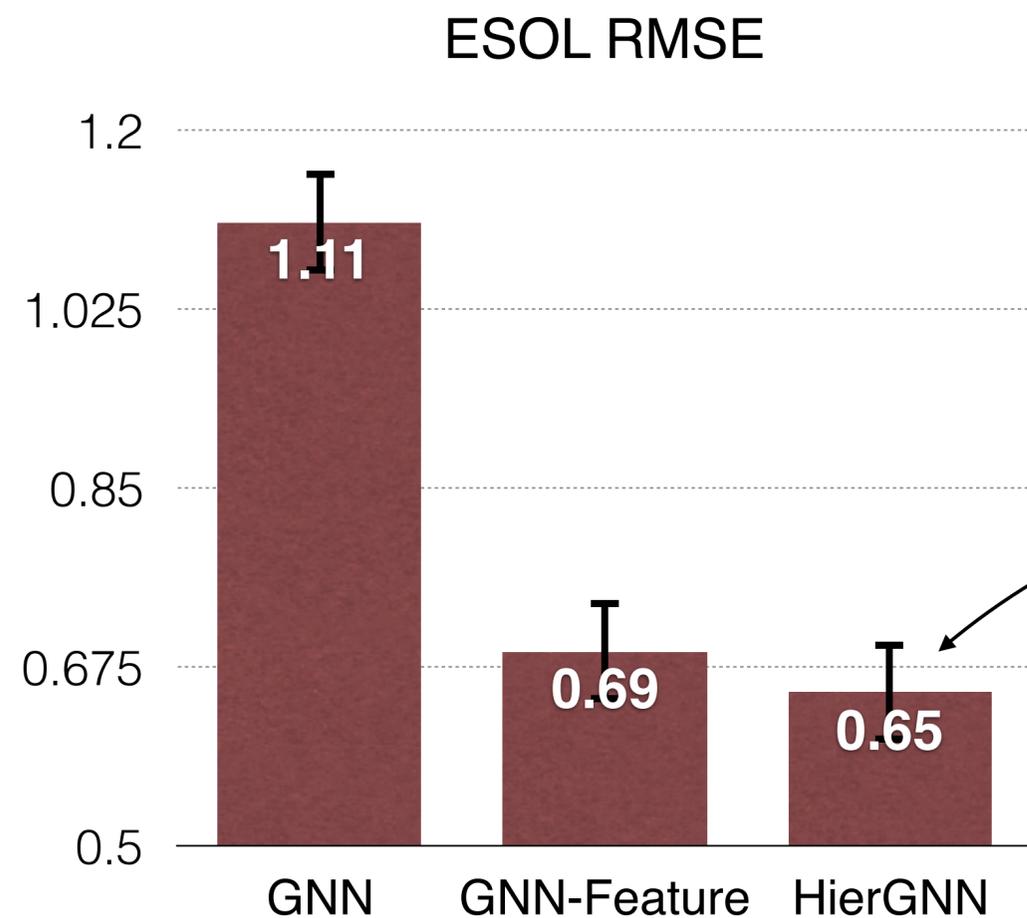
GNN with **features**

- ▶ atom type label
- ▶ degree
- ▶ valence
- ▶ **whether an atom is in a cycle**
- ▶ **whether an atom is in an aromatic ring**
- ▶

**Cycle
information**

Experiments on solubility

- ▶ ESOL dataset (averaged over 5 folds)

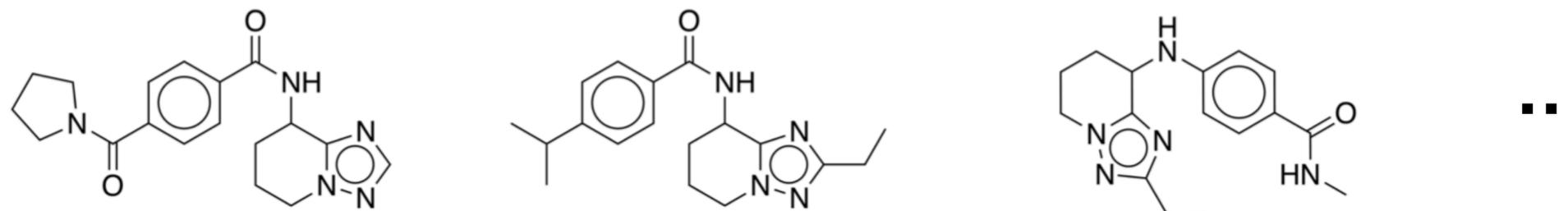


Hierarchical GNN

- ▶ Atom features: still just atom type
- ▶ But has extra substructure information built into the architecture

New Antibiotic Discovery

- ▶ If we can accurately predict molecular properties, we can screen (select and repurpose) molecules from a large candidate set



- ▶ Antibiotic Discovery [Stokes et al., 2019]
 - Trained a model to predict the inhibition against E. Coli (some bacteria...)
 - Data: ~2000 measured compounds from Broad Institute at MIT
 - Screened in total ~100 million compounds
 - Biologists tested 15 molecules (top prediction, structurally diverse) in the lab
 - 7 of them are validated to be inhibitive in-vitro
 - 1 of them demonstrate strong inhibition against other bacteria (e.g., *A. baumannii*)
 - All of them are new antibiotics distinct from existing ones!

Automating Drug design

▸ Key challenges:

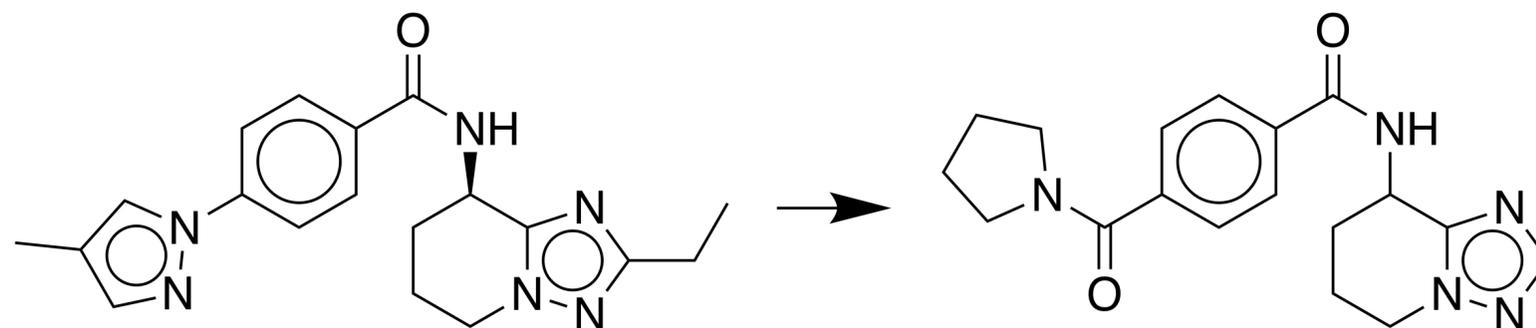
1. **representation and prediction:** learn to predict molecular properties
2. **generation and optimization:** realize target molecules with better properties programmatically
3. **understanding:** uncover principles (or diagnose errors) underlying complex predictions

De novo molecule optimization

- ▶ **Goal:** We aim to programmatically turn precursor molecules into molecules that satisfy given design specifications

De novo molecule optimization

- ▶ **Goal:** We aim to programmatically turn precursor molecules into molecules that satisfy given design specifications



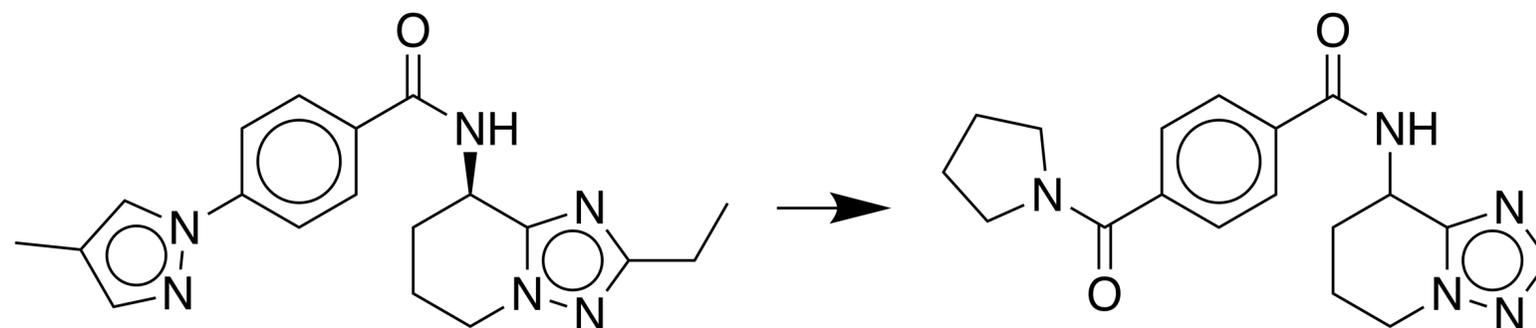
Source Molecule (QED=0.784)

QED=0.924

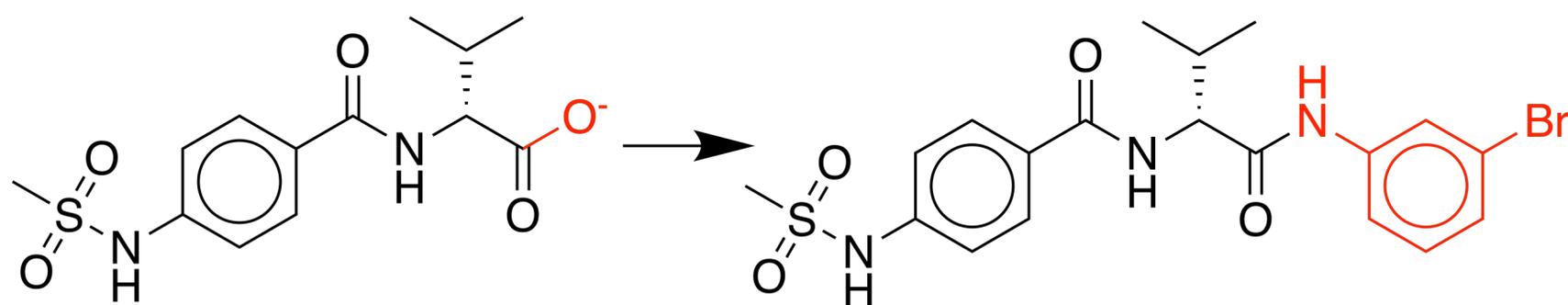
- ▶ Similar but ...
- ▶ Better drug-likeness

De novo molecule optimization

- ▶ **Goal:** We aim to programmatically turn precursor molecules into molecules that satisfy given design specifications



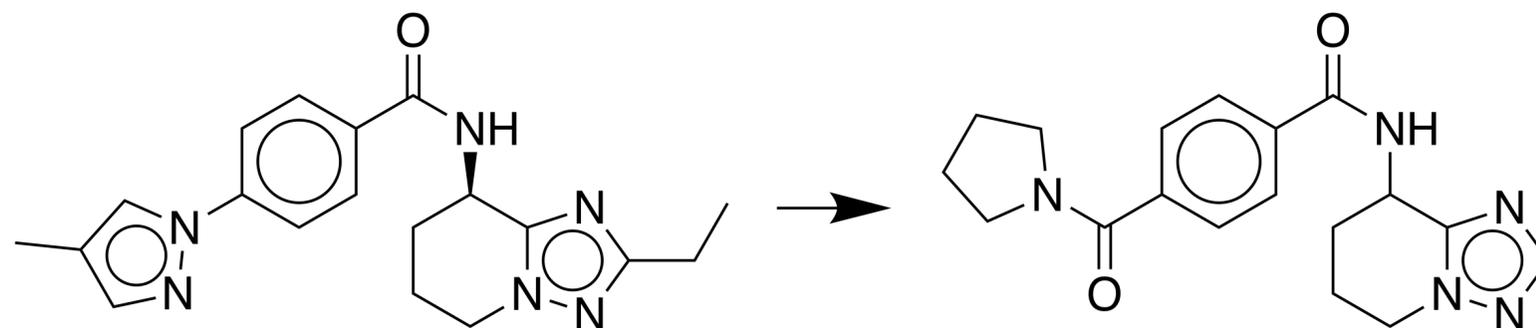
- ▶ Similar but ...
- ▶ Better drug-likeness



- ▶ Similar but ...
- ▶ Better solubility

De novo molecule optimization

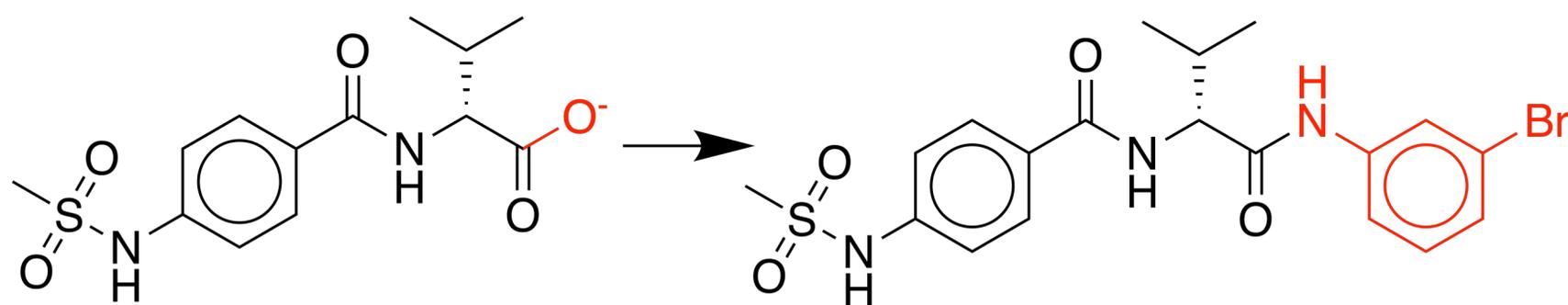
- ▶ **Goal:** We aim to programmatically turn precursor molecules into molecules that satisfy given design specifications



Source Molecule (QED=0.784)

QED=0.924

- ▶ Similar but ...
- ▶ Better drug-likeness

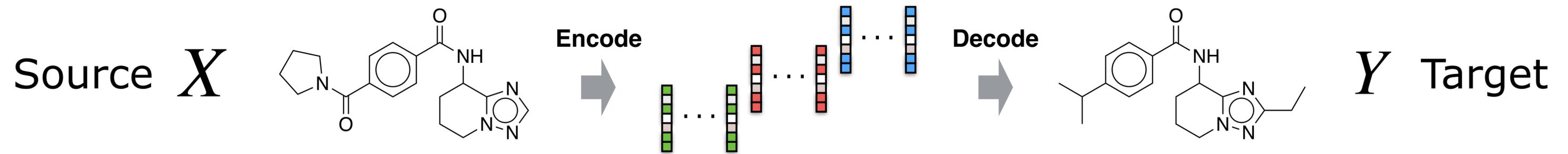


- ▶ Similar but ...
- ▶ Better solubility

- ▶ Need to learn a molecule-to-molecule mapping (i.e., graph-to-graph)

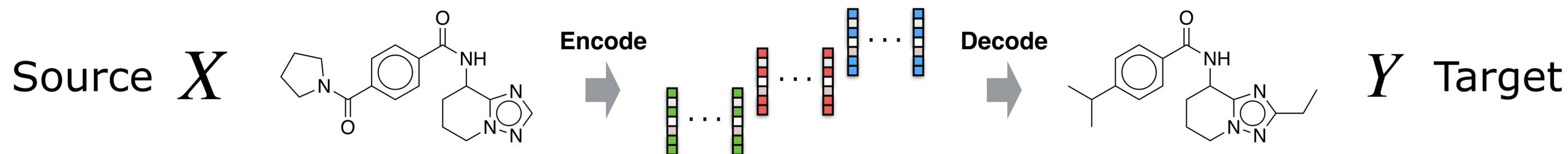
Molecule optimization as Graph Translation

- ▶ **Goal:** We aim to programmatically turn precursor molecules into molecules that satisfy given design specifications

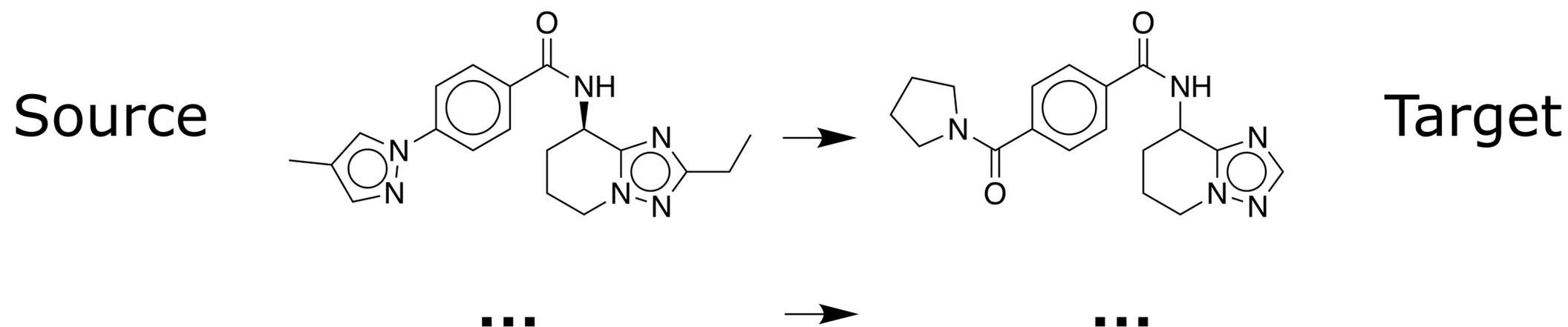


Molecule optimization as Graph Translation

- ▶ **Goal:** We aim to programmatically turn precursor molecules into molecules that satisfy given design specifications

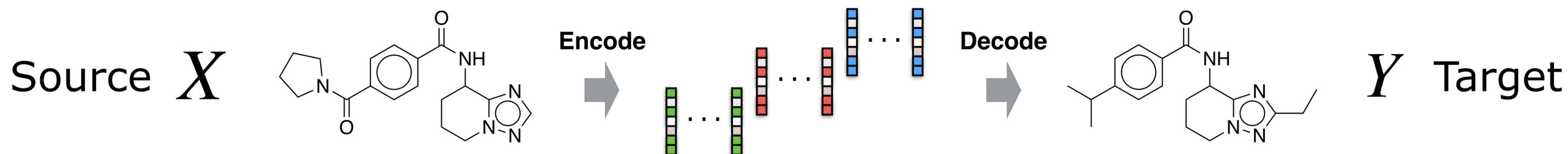


- ▶ The training set consists of (source, target) molecular pairs, e.g.,

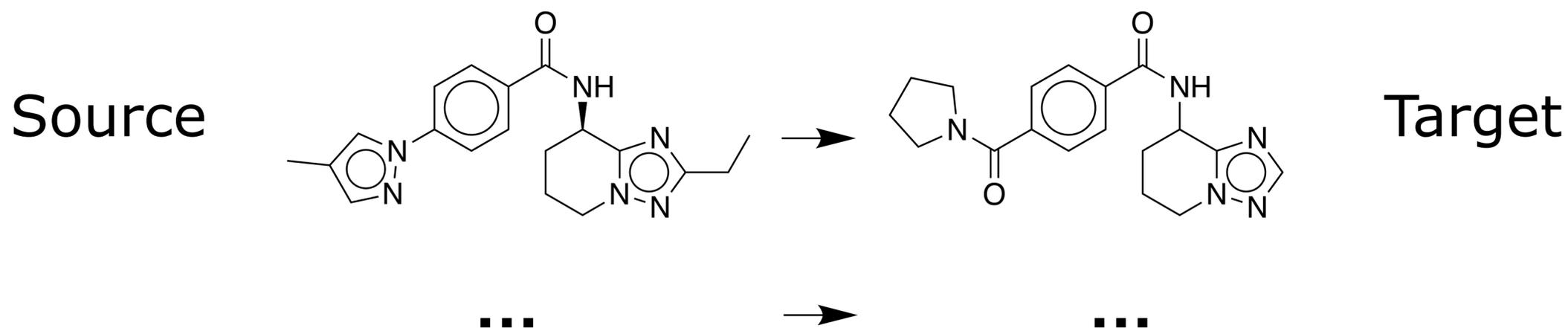


Molecule optimization as Graph Translation

- ▶ **Goal:** We aim to programmatically turn precursor molecules into molecules that satisfy given design specifications



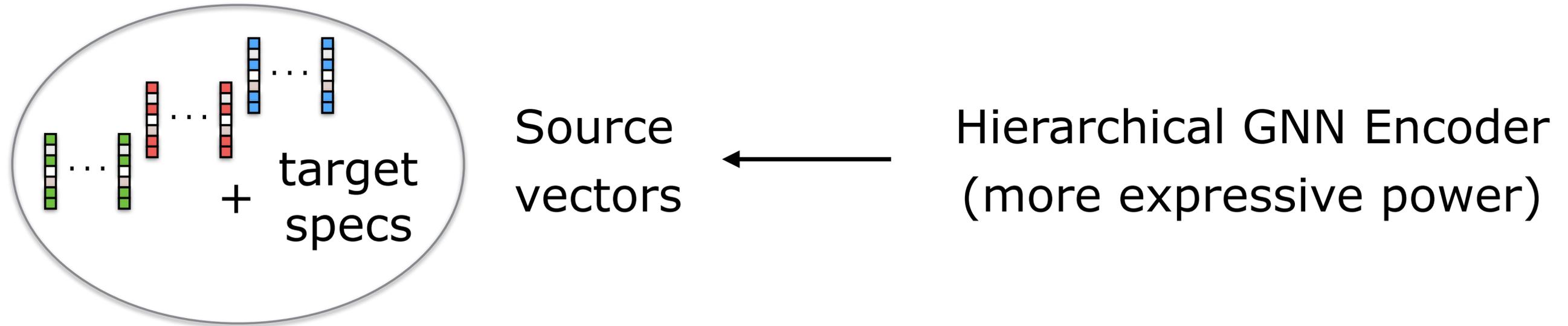
- ▶ The training set consists of (source, target) molecular pairs, e.g.,



- ▶ **Key challenges:** graph generation, diversity, multi-criteria optimization

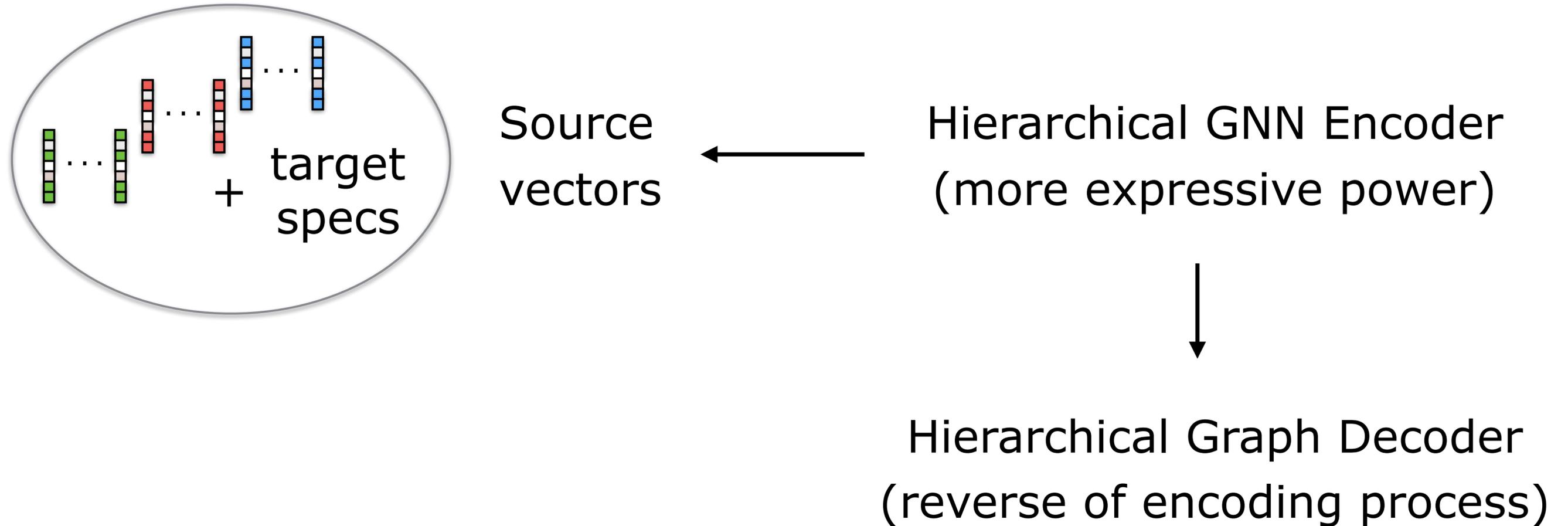
Graph-to-Graph Translation (Decoder)

- ▶ Modifying a pre-cursor to meet target specifications



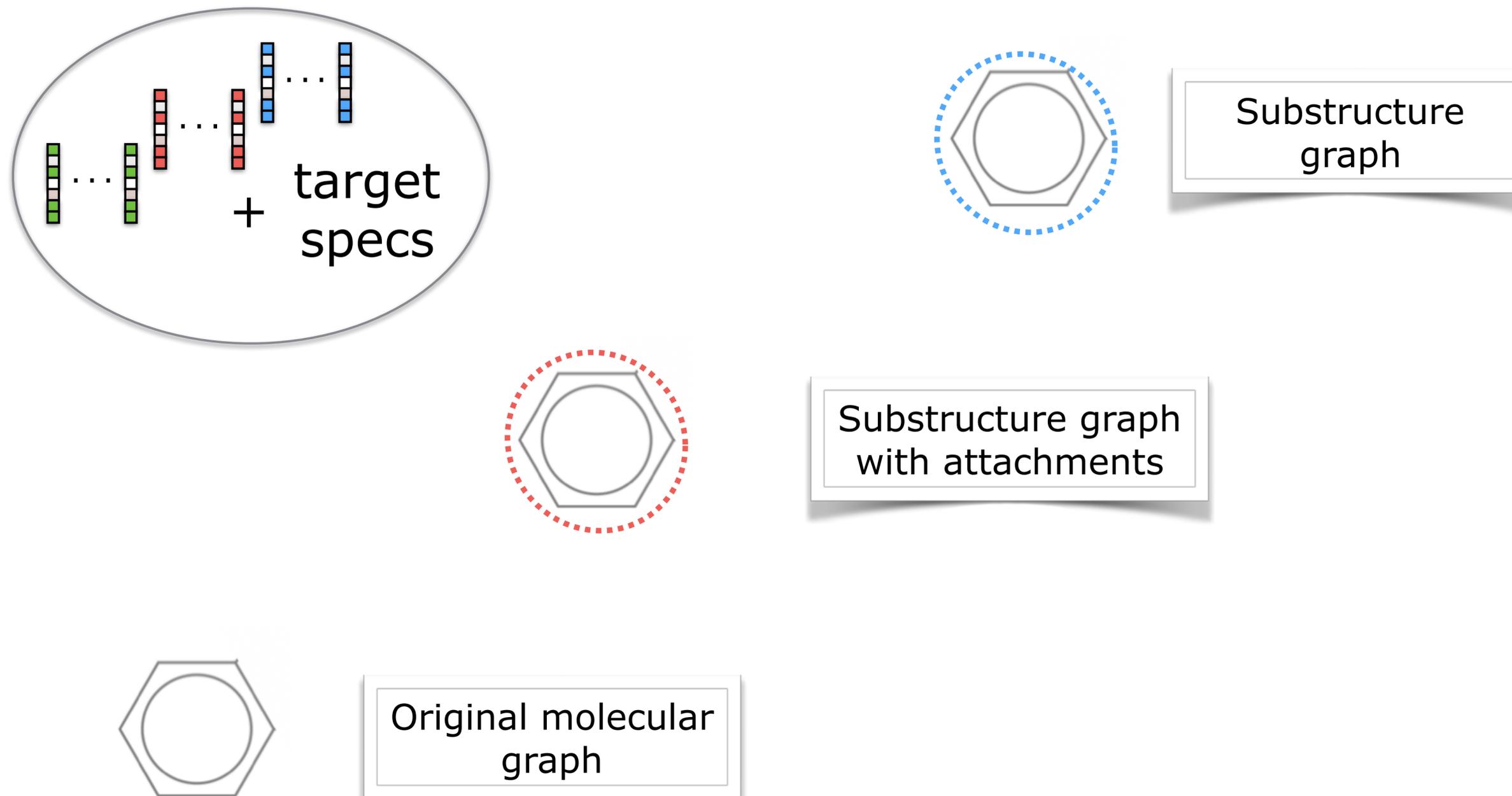
Graph-to-Graph Translation (Decoder)

- ▶ Modifying a pre-cursor to meet target specifications



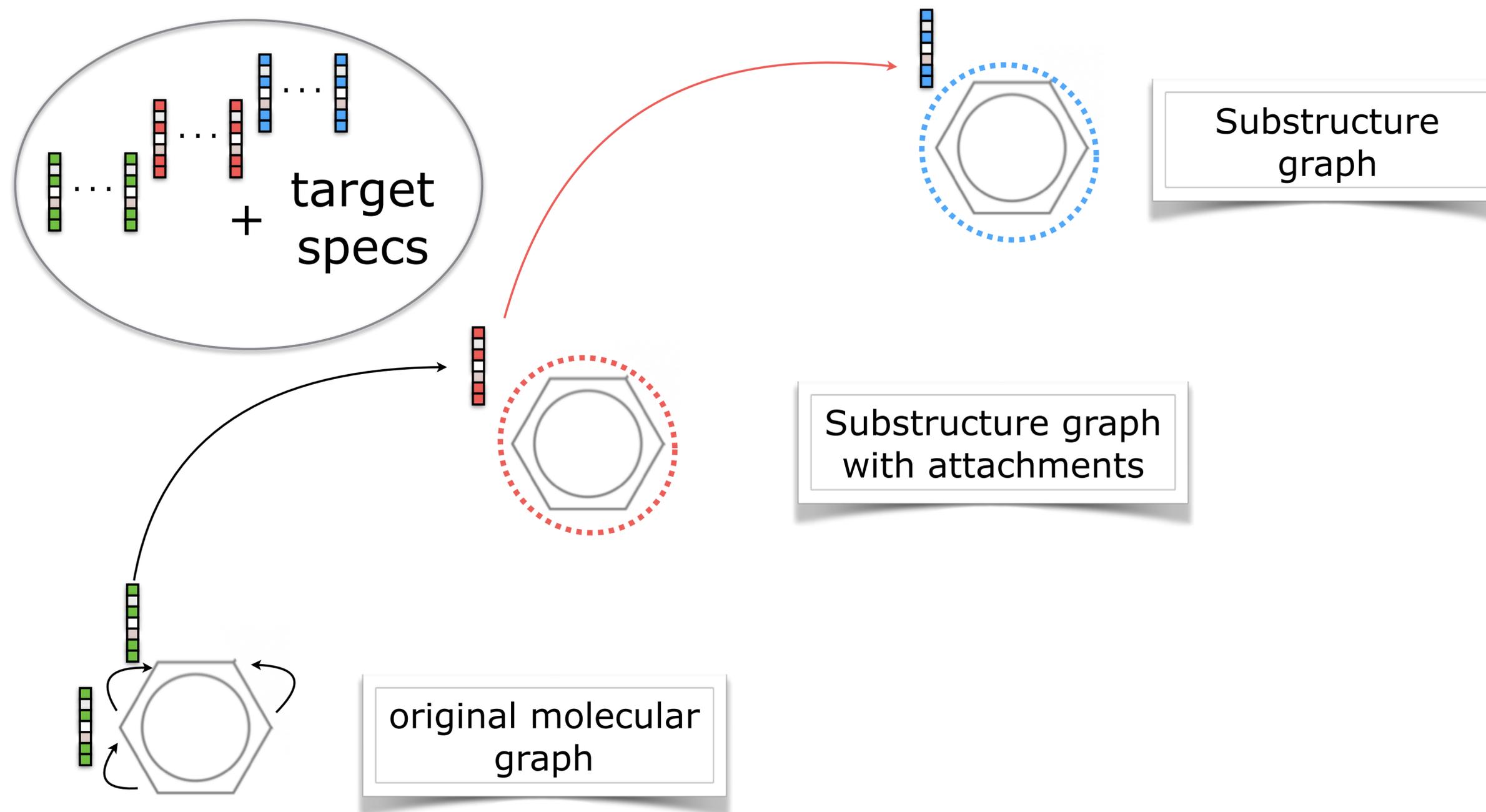
Graph-to-Graph Translation (Decoder)

- ▶ Modifying a pre-cursor to meet target specifications



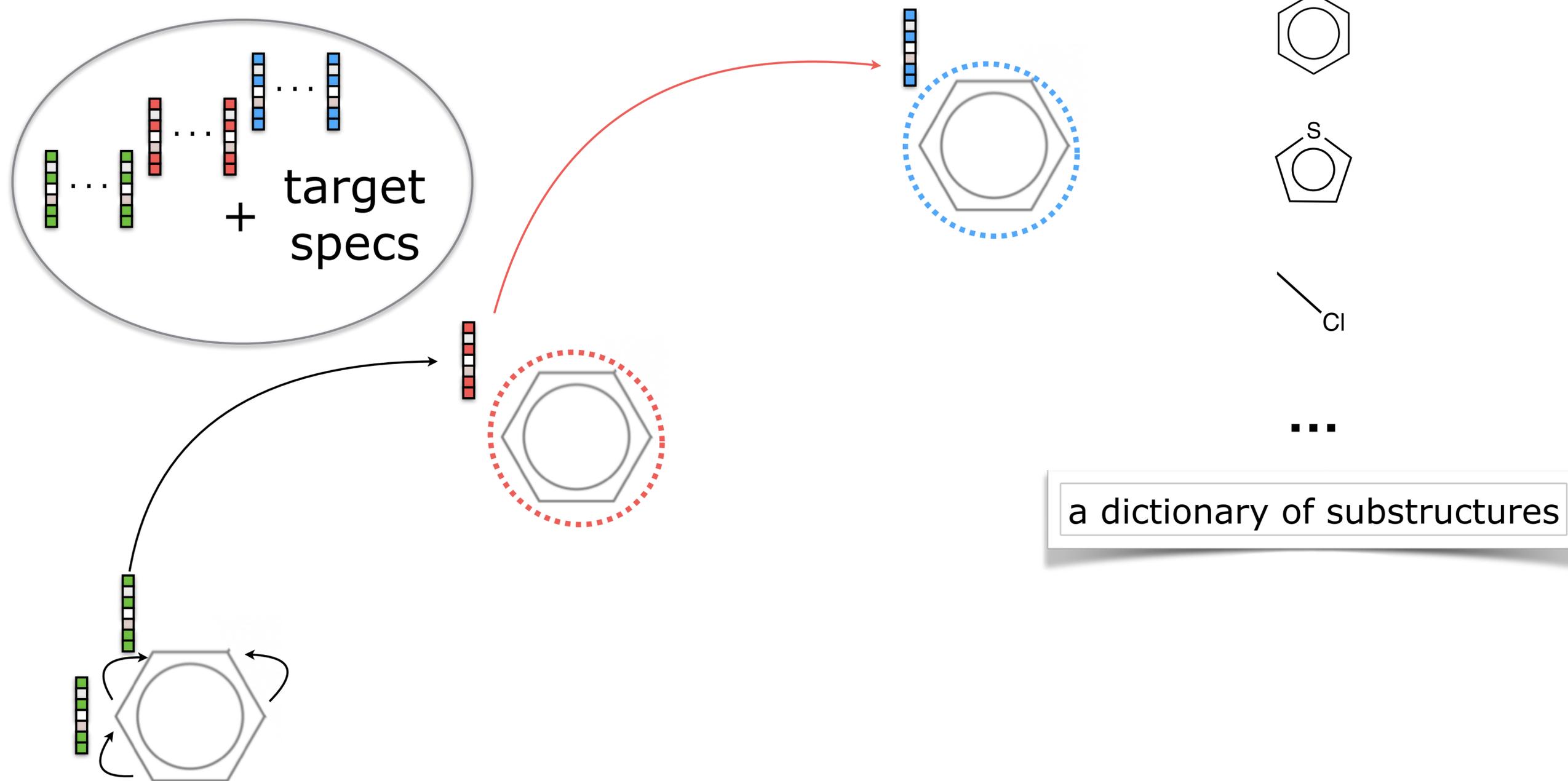
Graph-to-Graph Translation (Decoder)

- ▶ Modifying a pre-cursor to meet target specifications



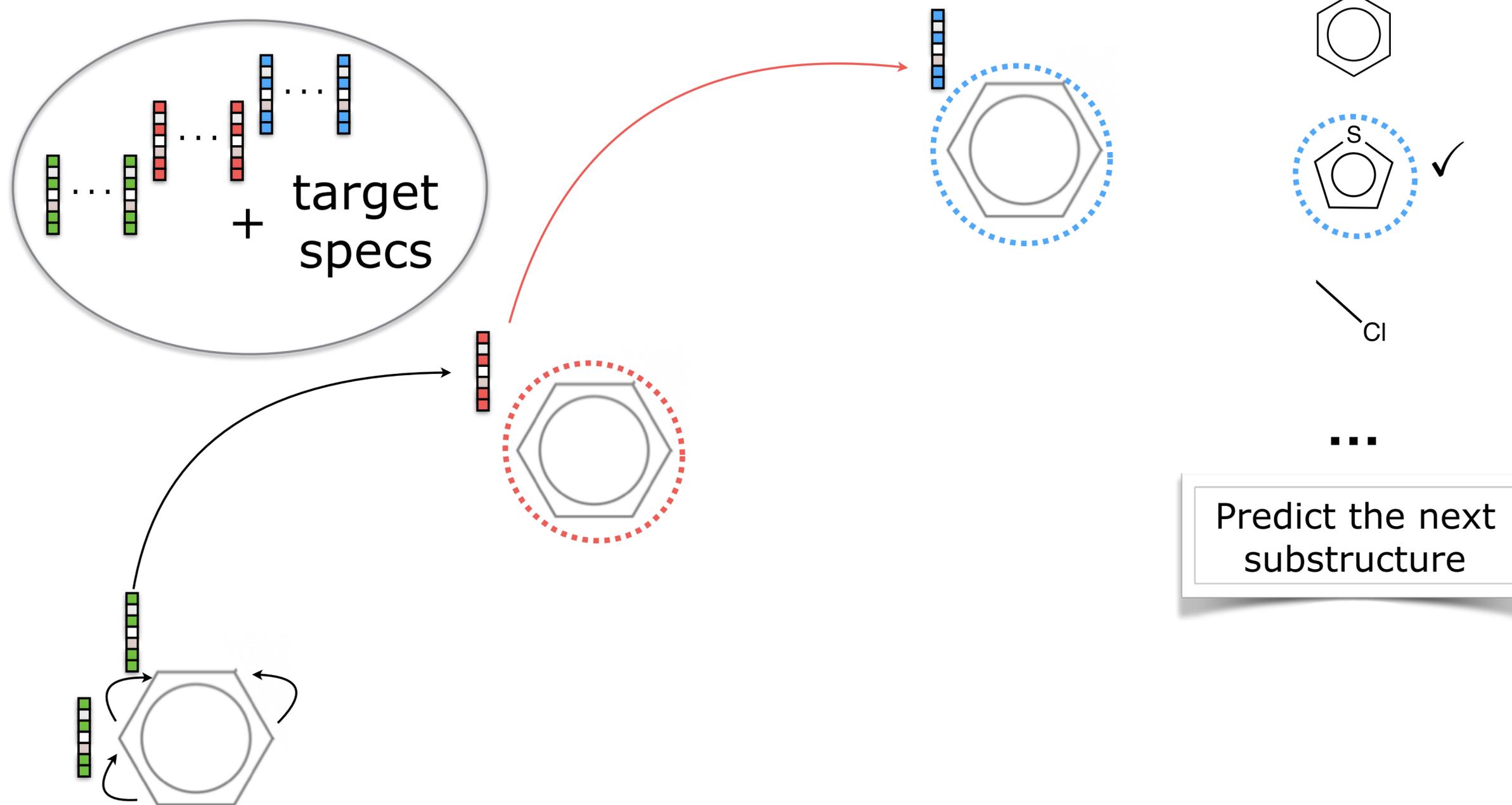
Graph-to-Graph Translation (Decoder)

- ▶ Modifying a pre-cursor to meet target specifications



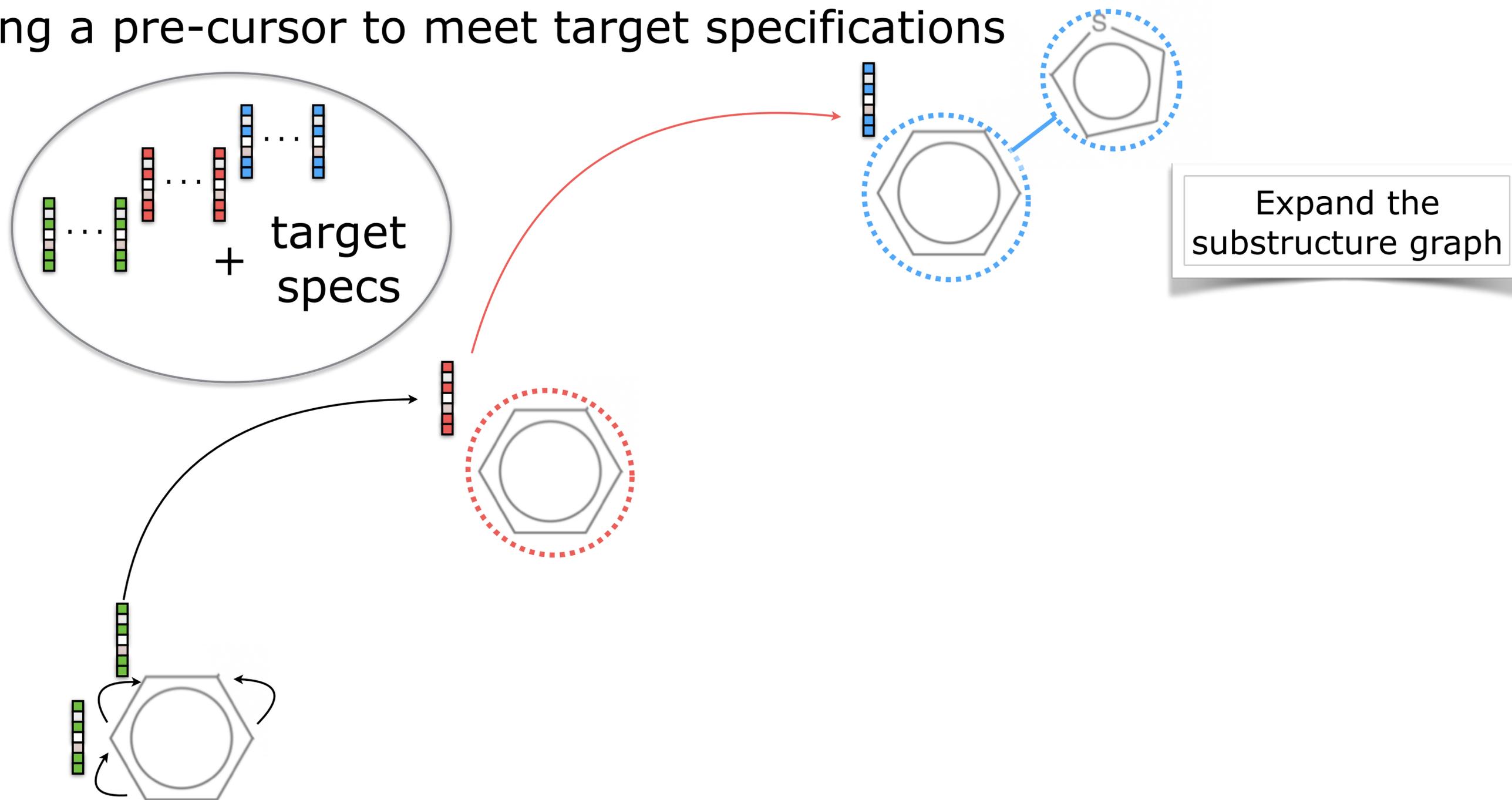
Graph-to-Graph Translation (Decoder)

- ▶ Modifying a pre-cursor to meet target specifications



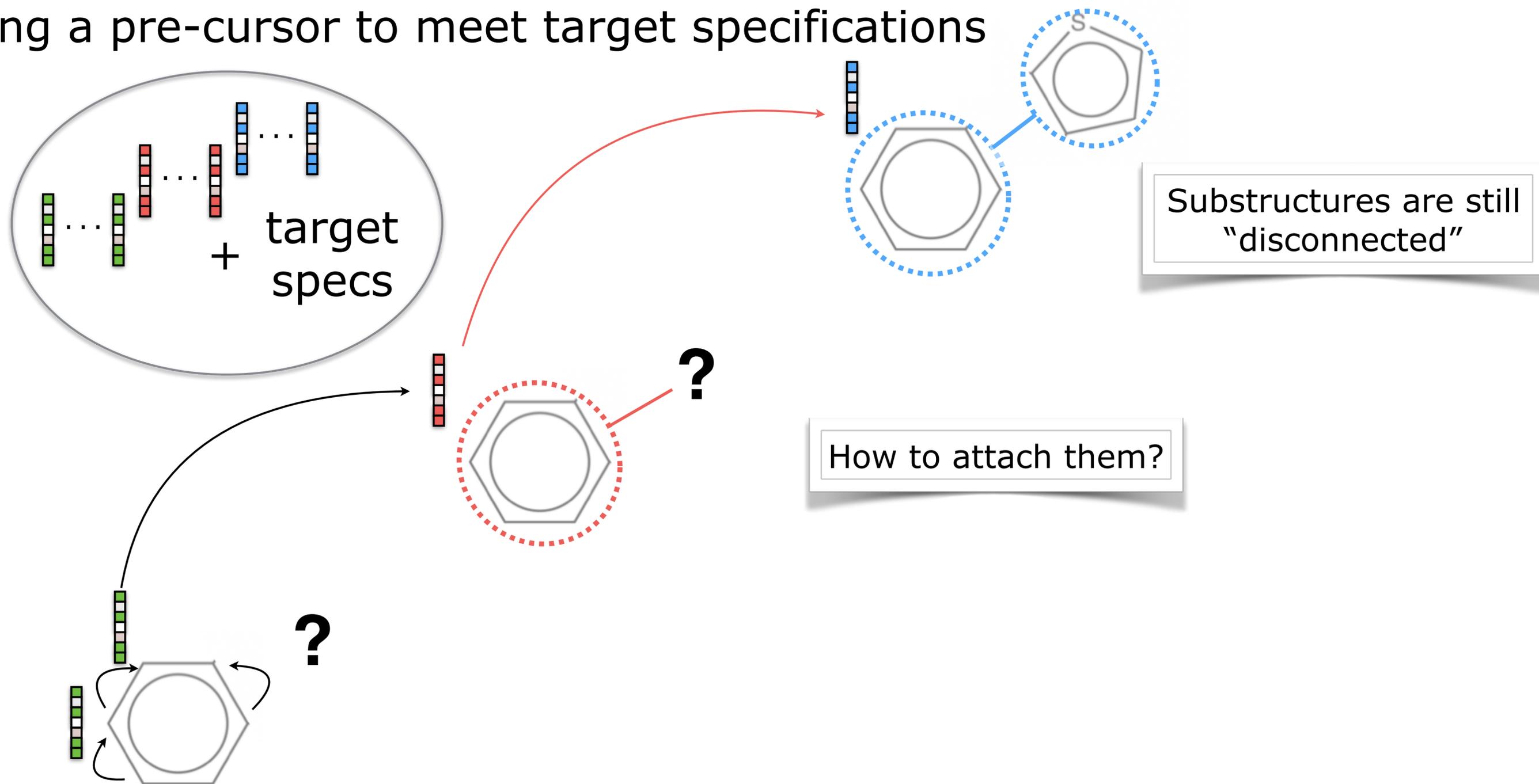
Graph-to-Graph Translation (Decoder)

- ▶ Modifying a pre-cursor to meet target specifications



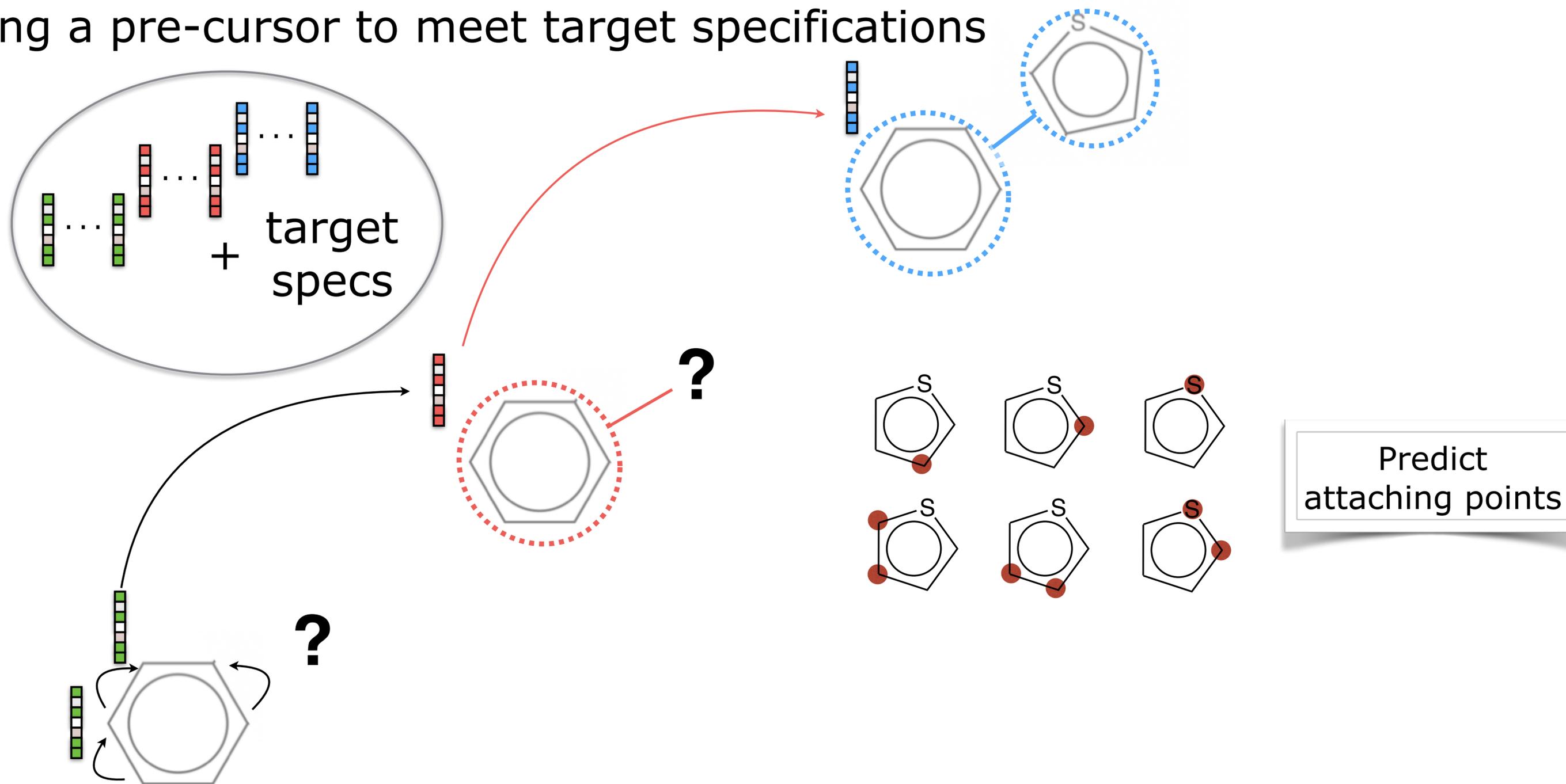
Graph-to-Graph Translation (Decoder)

- ▶ Modifying a pre-cursor to meet target specifications



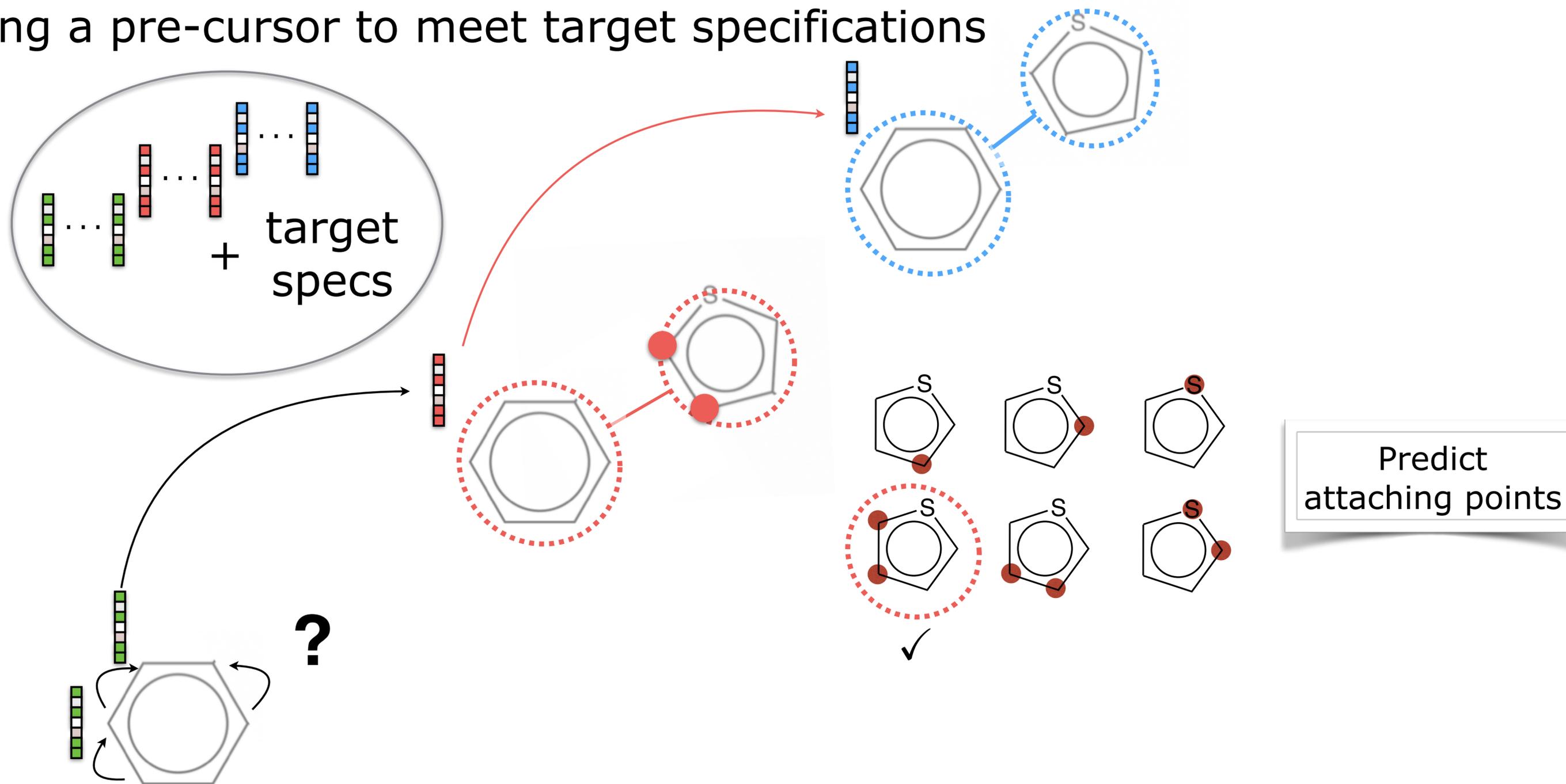
Graph-to-Graph Translation (Decoder)

- ▶ Modifying a pre-cursor to meet target specifications



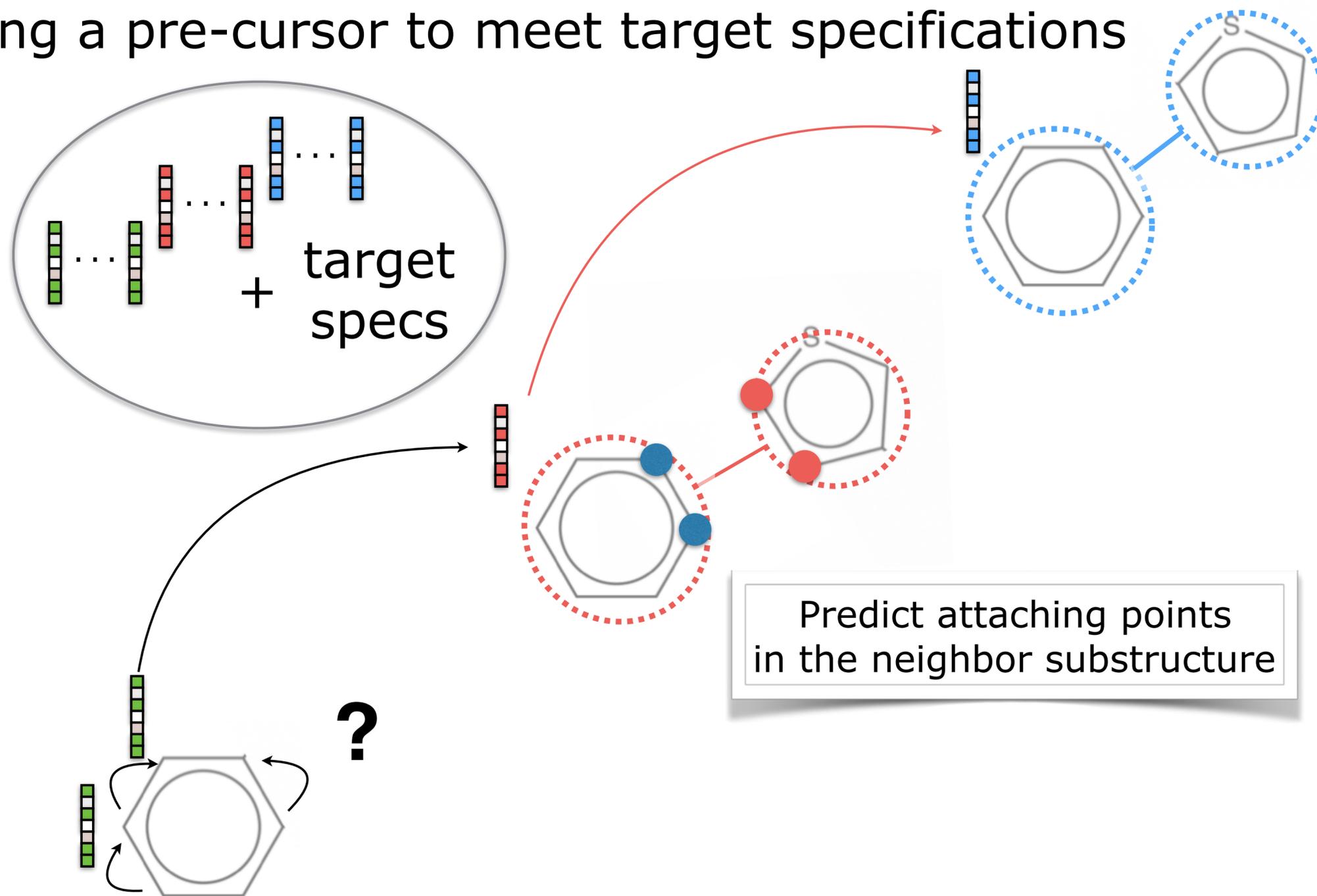
Graph-to-Graph Translation (Decoder)

- ▶ Modifying a pre-cursor to meet target specifications



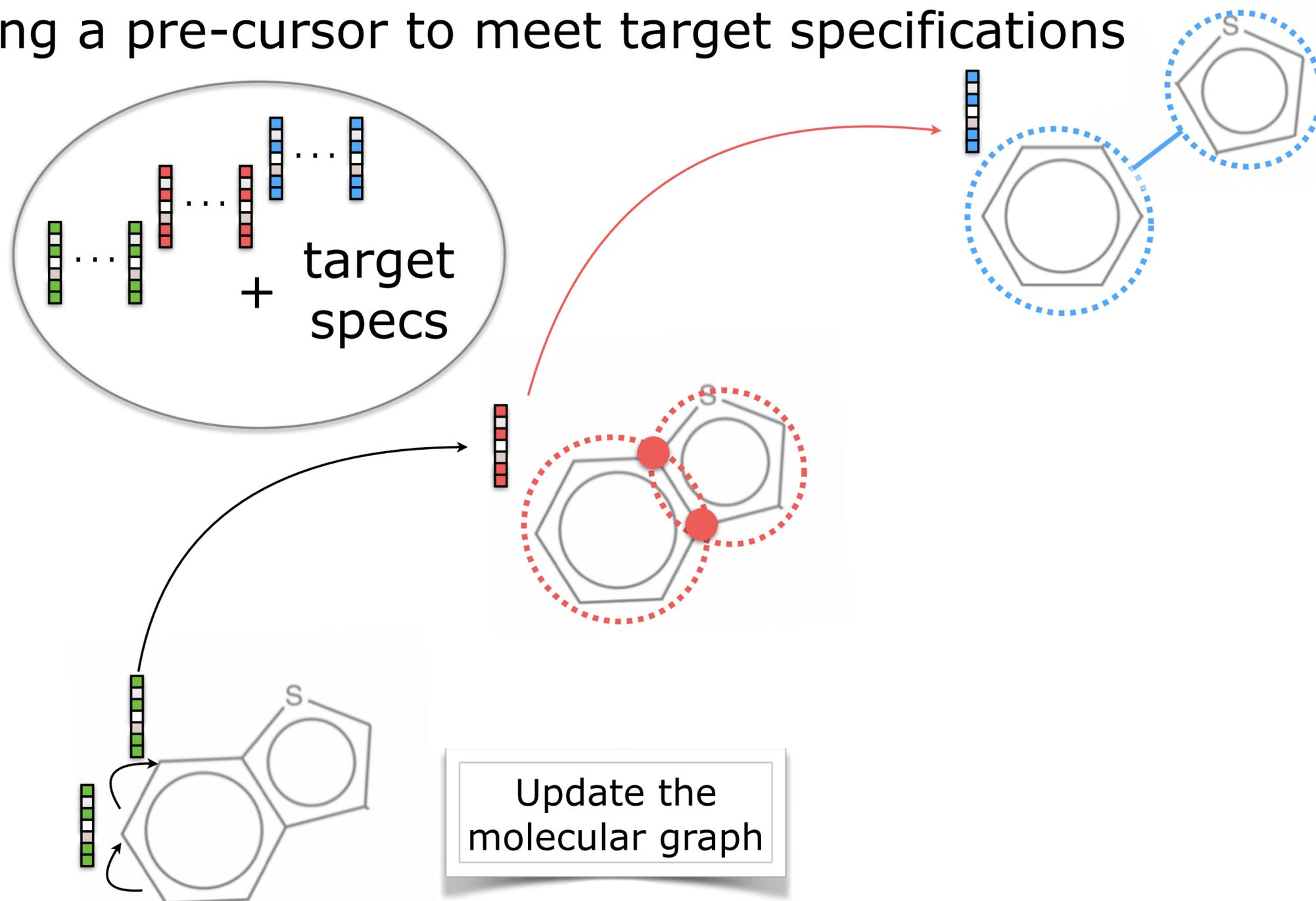
Graph-to-Graph Translation (Decoder)

- ▶ Modifying a pre-cursor to meet target specifications



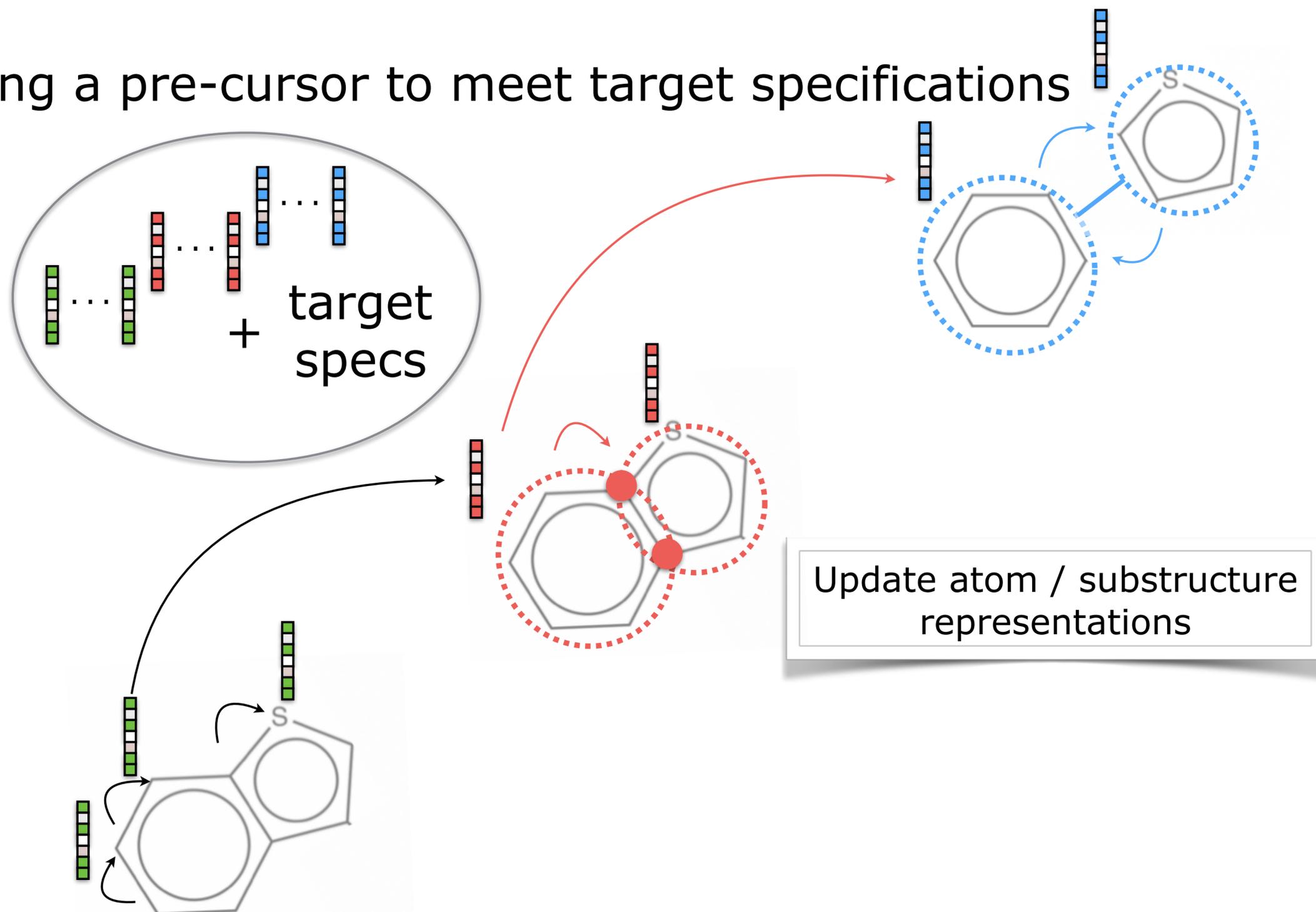
Graph-to-Graph Translation (Decoder)

- ▶ Modifying a pre-cursor to meet target specifications



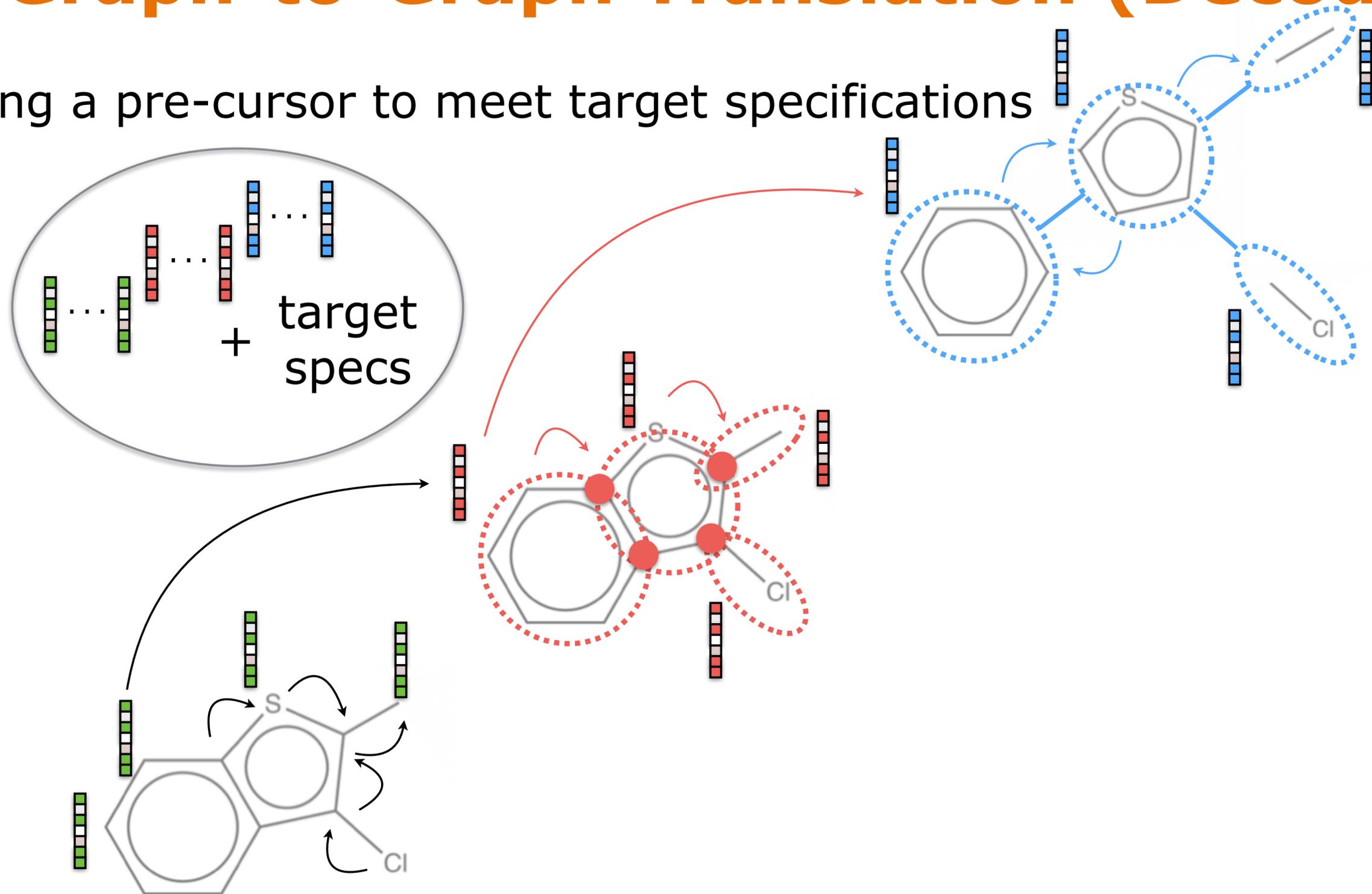
Graph-to-Graph Translation (Decoder)

- ▶ Modifying a pre-cursor to meet target specifications



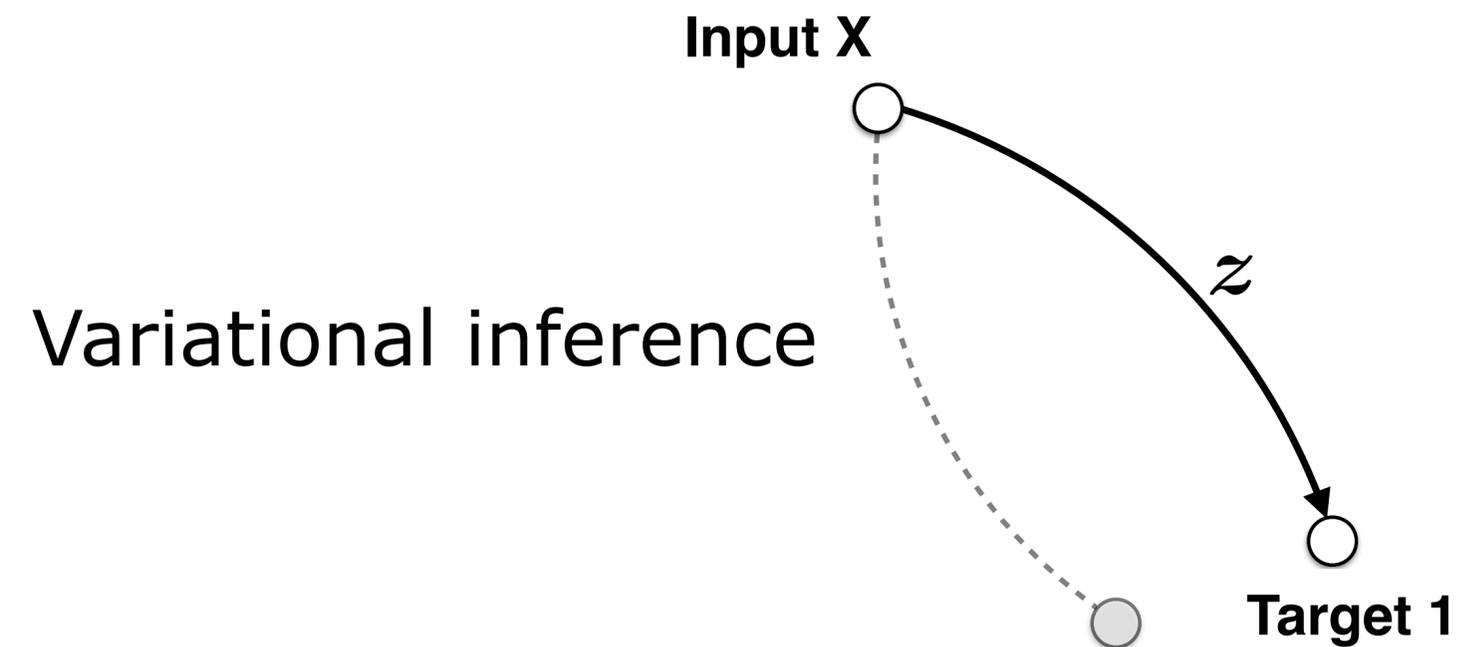
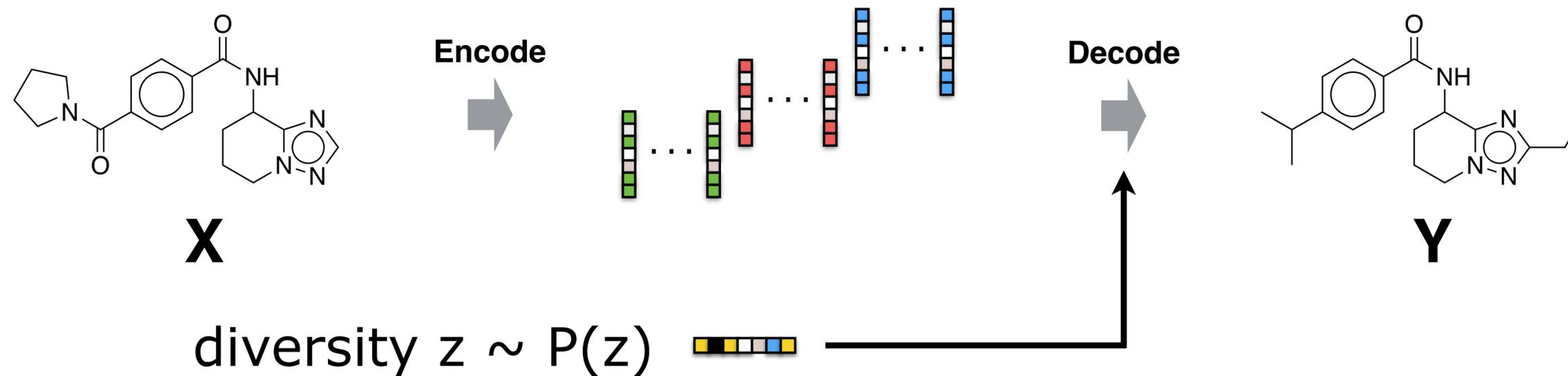
Graph-to-Graph Translation (Decoder)

- ▶ Modifying a pre-cursor to meet target specifications



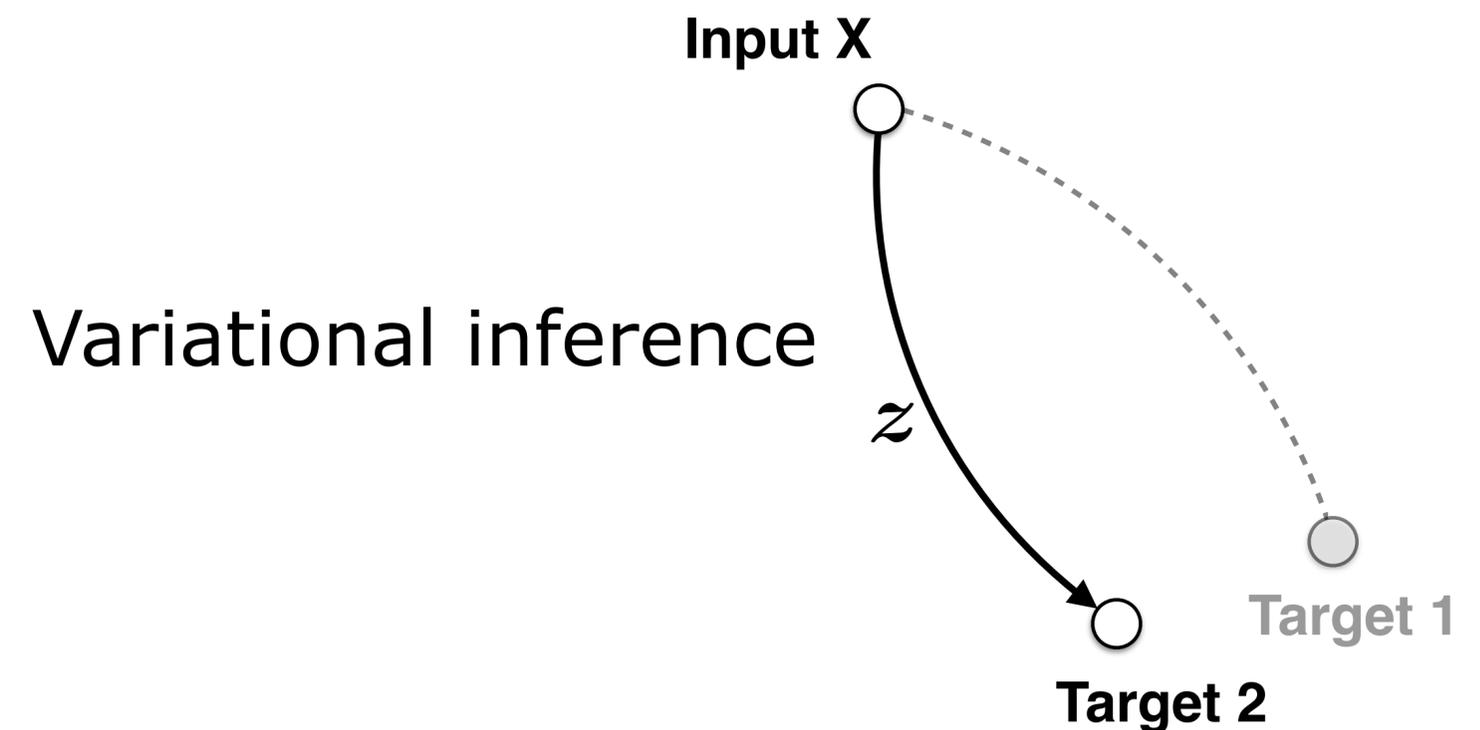
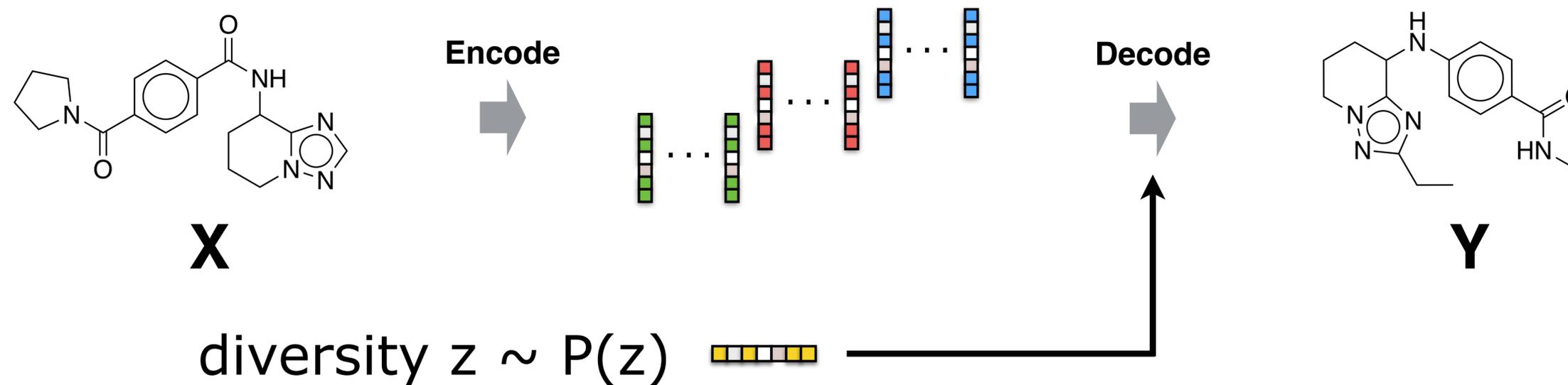
De novo molecule optimization: diversity

- ▶ **Goal:** We aim to programmatically turn precursor molecules into versions that satisfy given design specifications



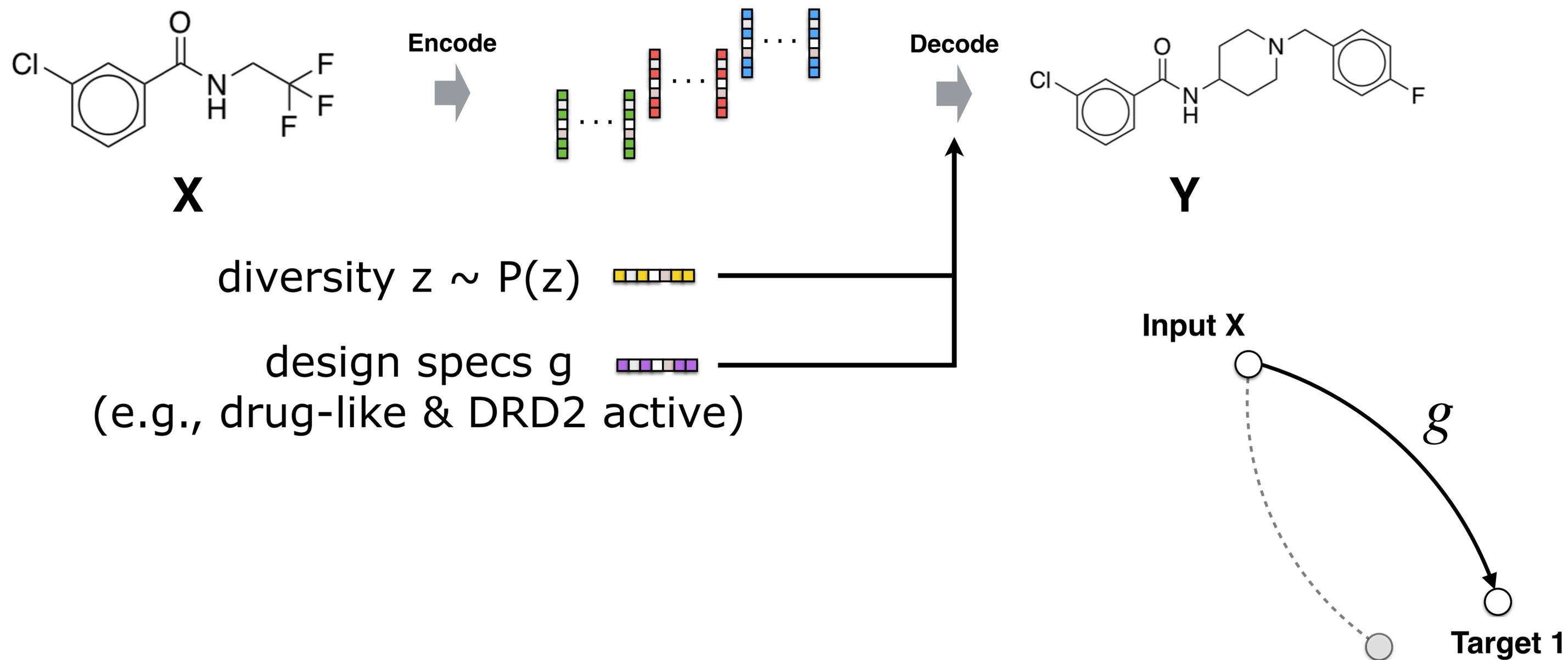
De novo molecule optimization: diversity

- ▶ **Goal:** We aim to programmatically turn precursor molecules into versions that satisfy given design specifications



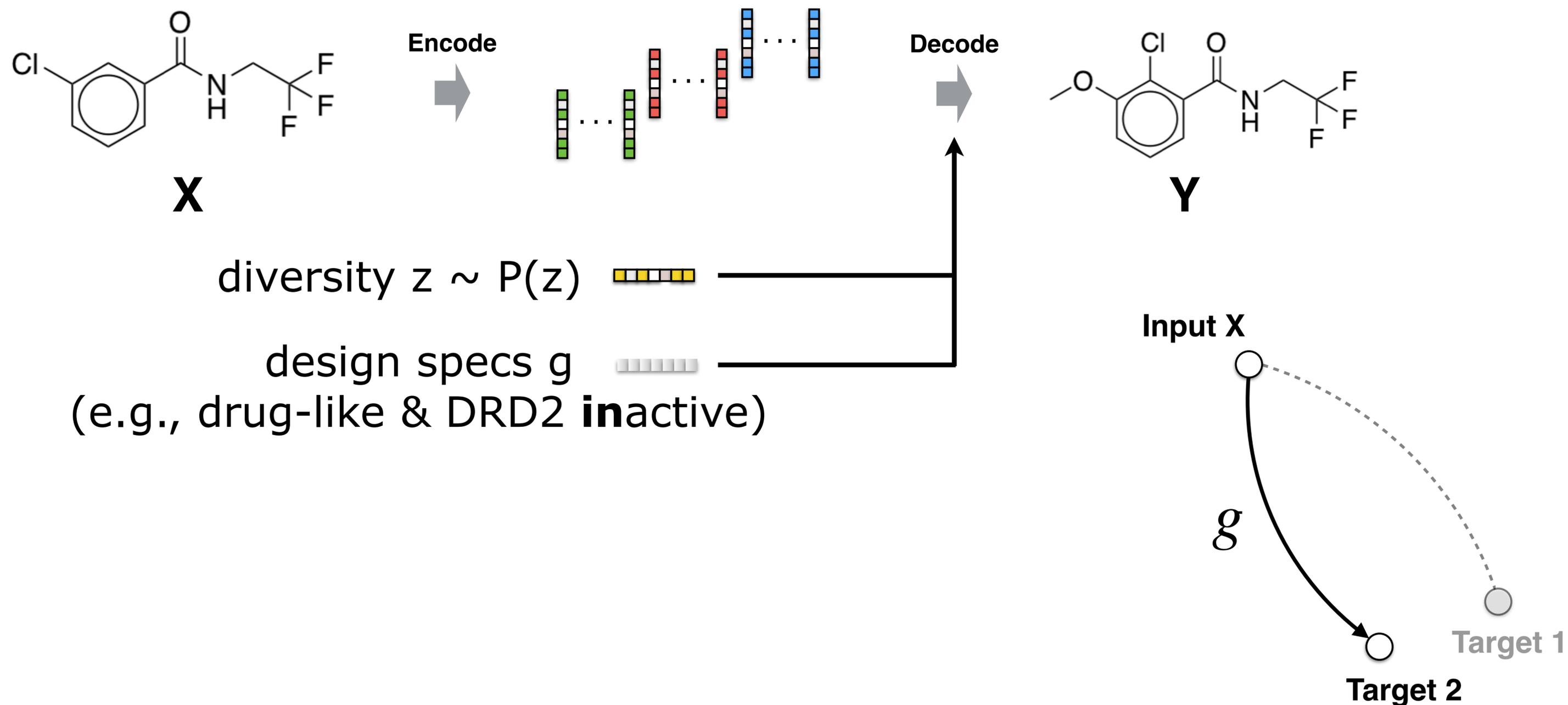
De novo molecule optimization: specs

- ▶ **Goal:** We aim to programmatically turn precursor molecules into versions that satisfy given design specifications



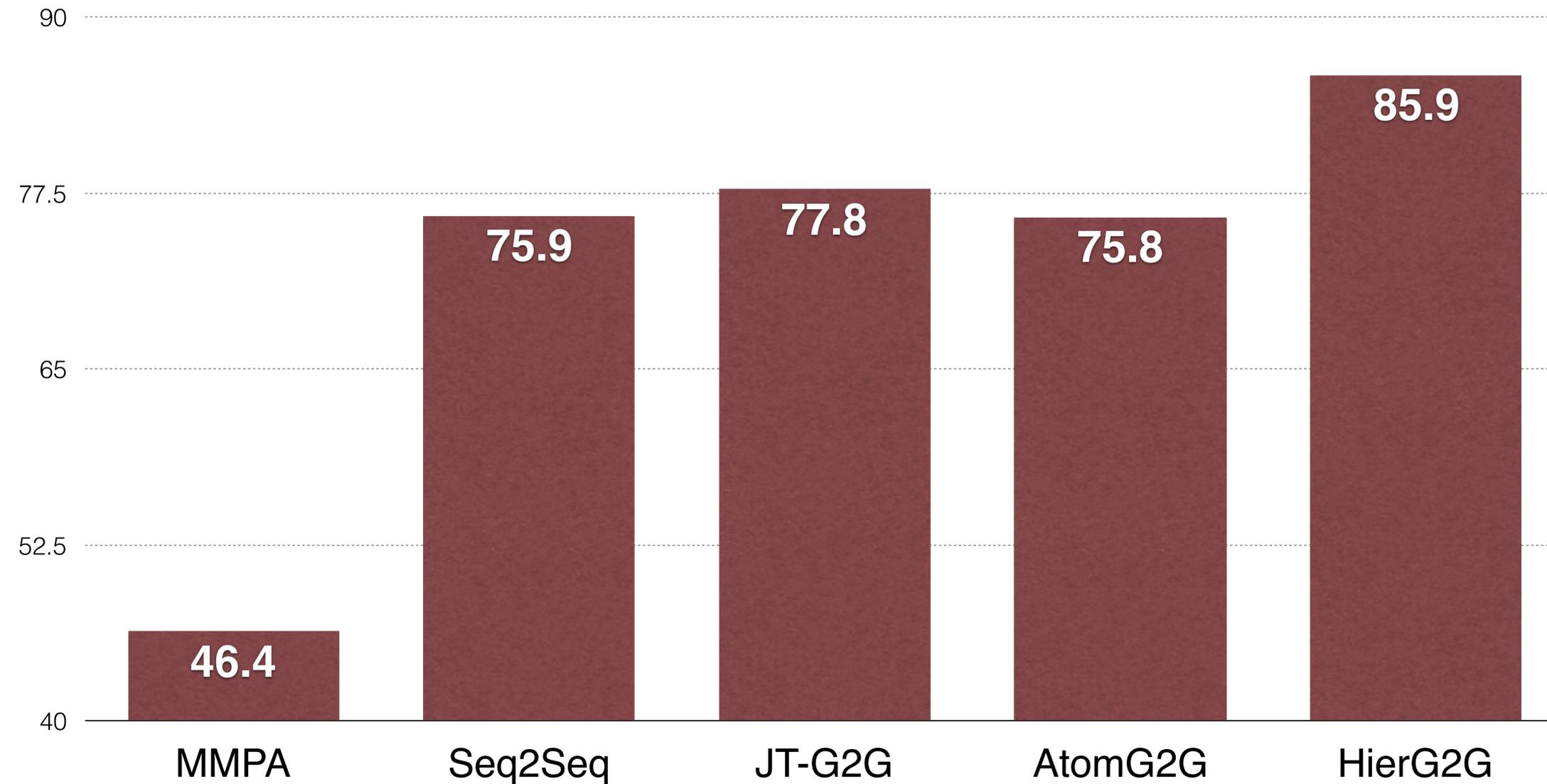
De novo molecule optimization: specs

- ▶ **Goal:** We aim to programmatically turn precursor molecules into versions that satisfy given design specifications



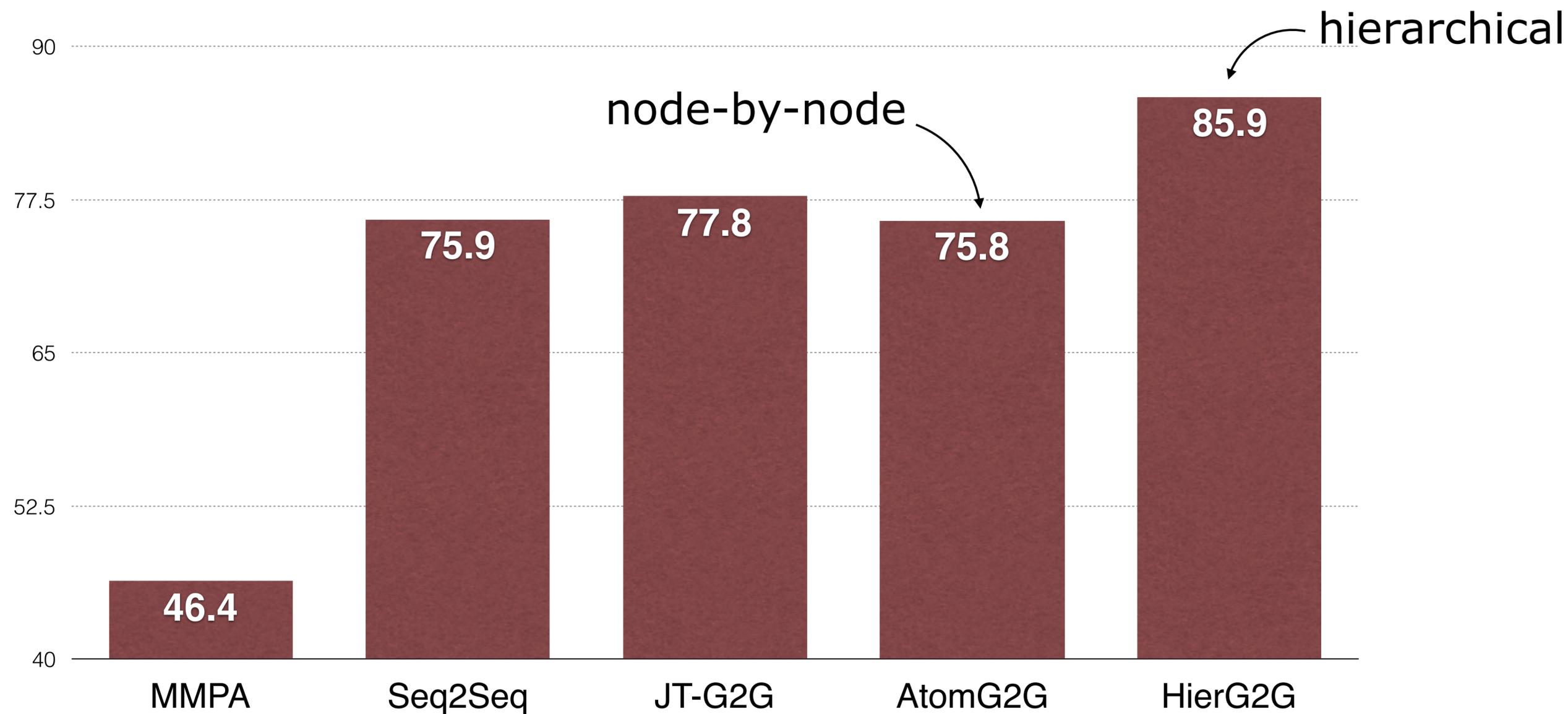
Example results (DRD2)

- ▶ Single property optimization: DRD2 success % (from inactive to active)



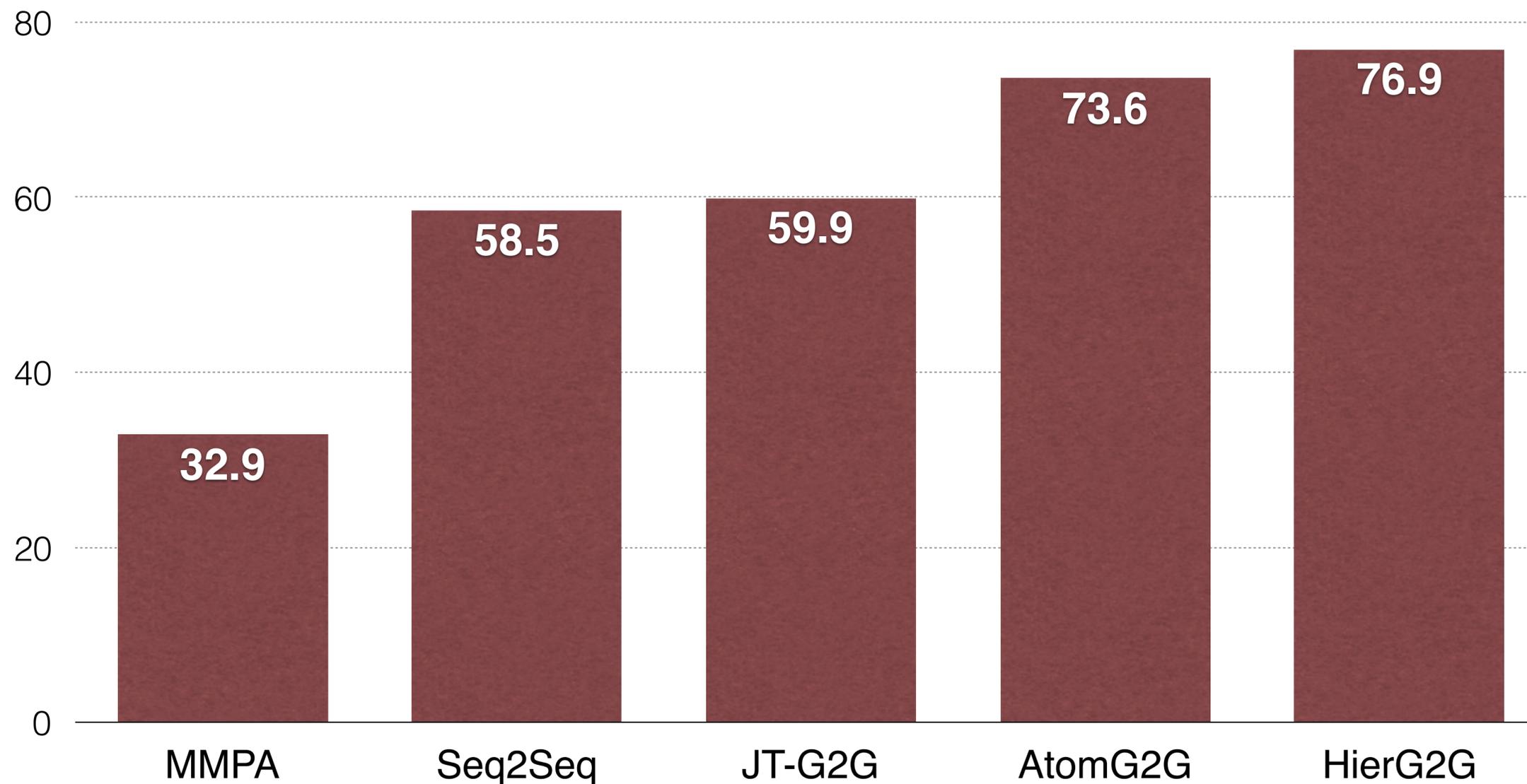
Example results (DRD2)

- Single property optimization: DRD2 success % (from inactive to active)



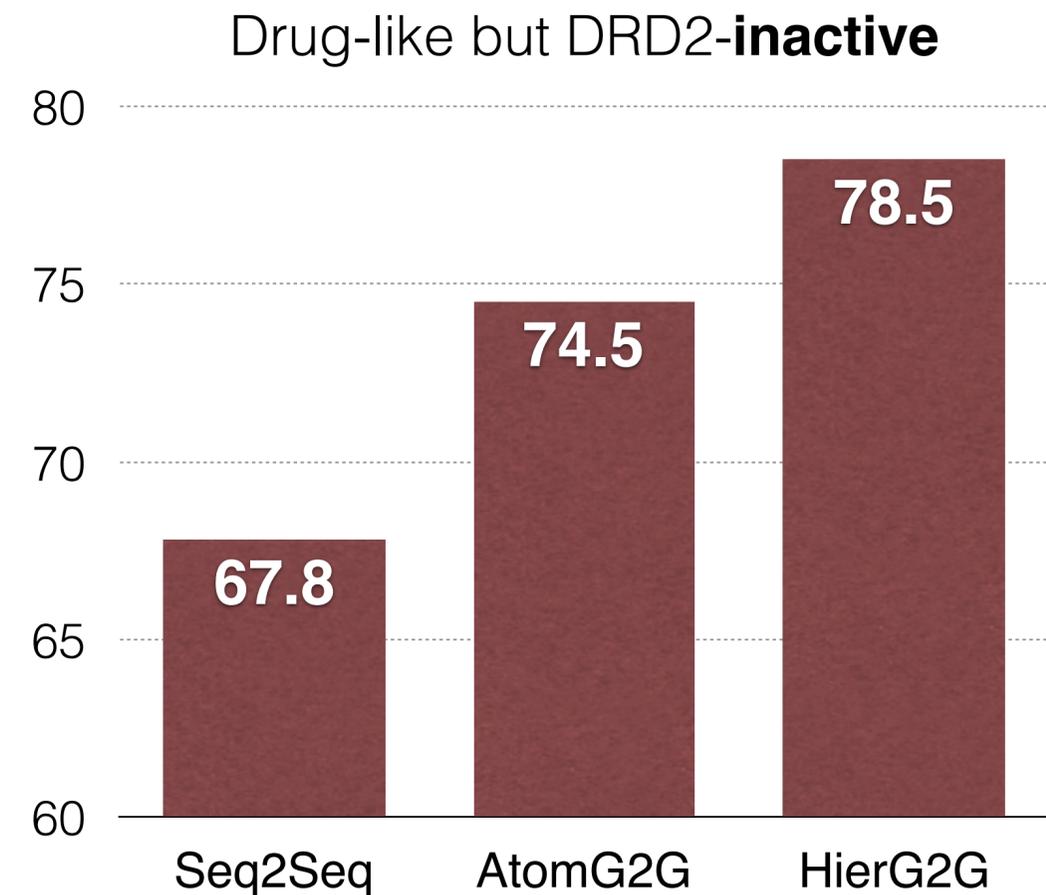
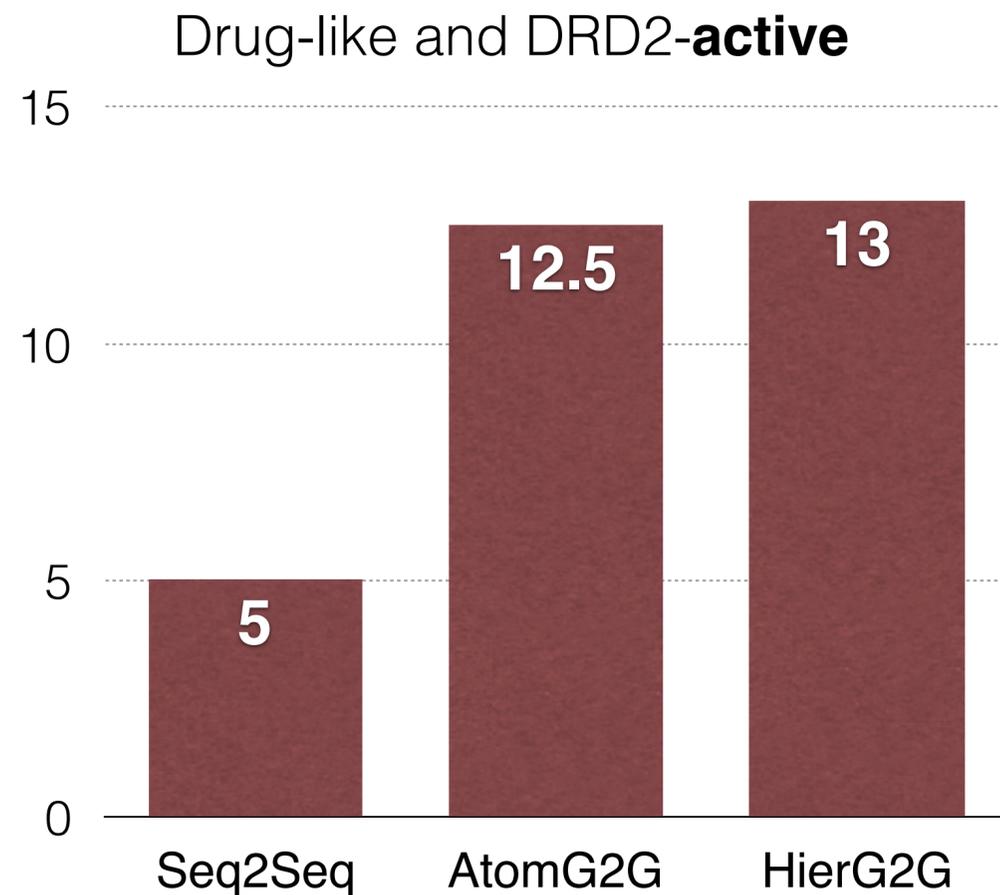
Example results (drug-likeness)

- ▶ Single property optimization: drug-likeness (QED) success % (QED > 0.9)



Example results (multiple design specs)

- Multi-criteria success % (design specs driven generation)



- Challenge:** only 1.6% training pairs are both drug-like and DRD2-active

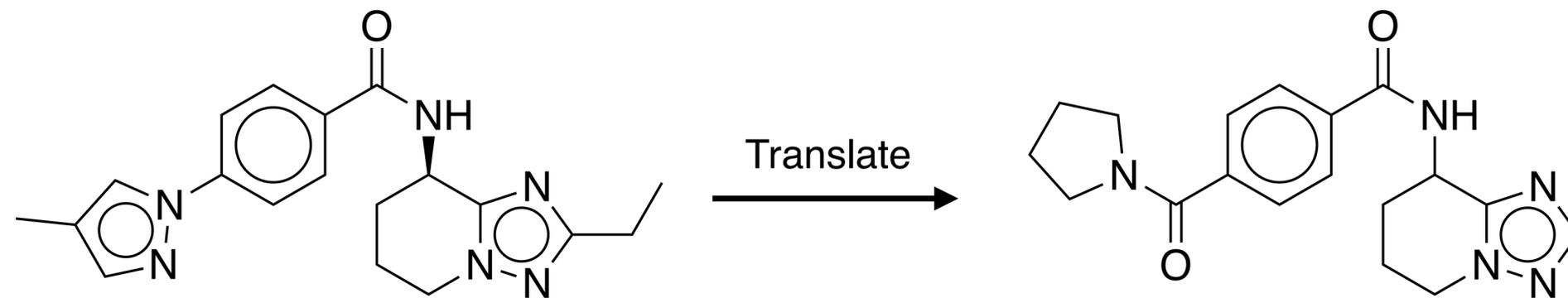
Disentangling what's important

- ▶ Models are complicated, important to assess how individual parts contribute to performance

Method	QED	DRD2
HierG2G	76.9%	85.9%
• atom-based decoder	76.1%	75.0%
• two-layer encoder	75.8%	83.5%
• one-layer encoder	67.8%	74.1%
• GRU MPN	72.6%	83.7%

Still many ways to improve

- ▶ Generating complex objects (e.g., molecules) is hard
- ▶ Assessing the quality of the object (property prediction) is substantially easier



hard to
realize

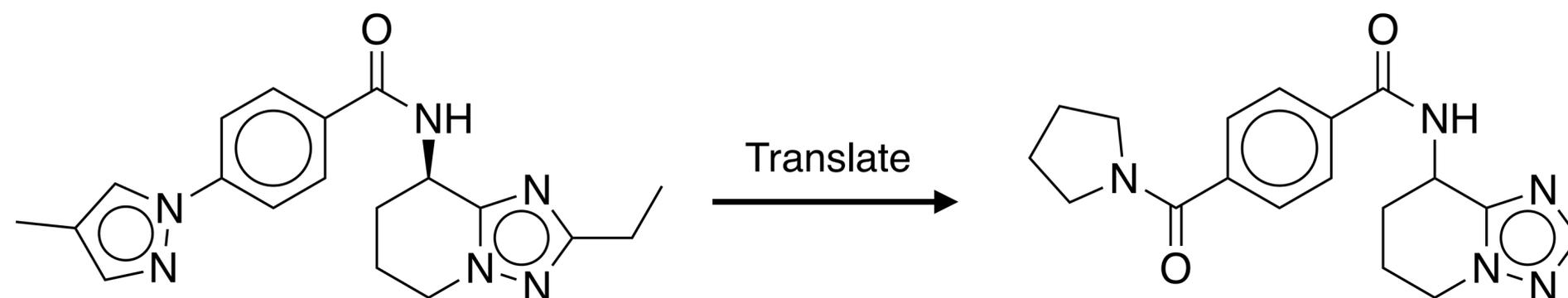
Constraints:

- Molecular Similarity $\text{sim}(X, Y) \geq 0.4$
- Drug-likeness $\text{QED}(Y) \geq 0.9$

easier to
check/predict

Property-guided generation

- ▶ Generating complex objects (e.g., molecules) is hard
- ▶ Assessing the quality of the object (property prediction) is substantially easier



Constraints:

- Molecular Similarity $\text{sim}(X, Y) \geq 0.4$
- Drug-likeness $\text{QED}(Y) \geq 0.9$

hard to
realize

easier to
check/predict

- ▶ **Target augmentation:** we can use property predictor to generate additional (self-supervised) data for the generative model

Target augmentation = stochastic EM

- ▶ **Objective:** maximize the log-probability that generated candidates satisfy the properties of interest (structure is now a latent variable)

$$\sum_{X \in \text{source set}} \log \left[\sum_Y P(\text{target specs} | Y) P(Y | X; \theta) \right]$$

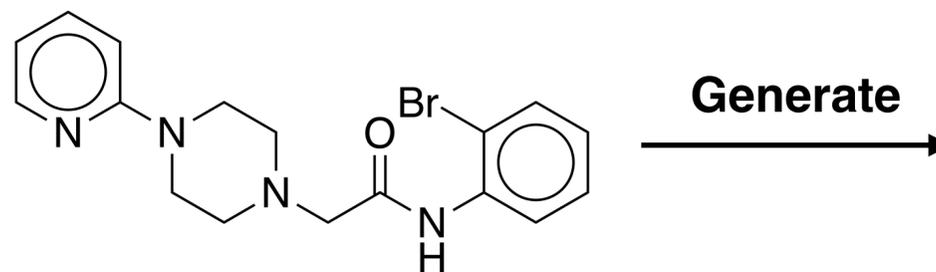
- ▶ **E-step:** generate candidates from the current model; filter/reweight by property predictor (\sim posterior samples)
- ▶ **M-step:** maximize the log-probability of new (weighted) targets

Target augmentation = stochastic EM

- ▶ **Objective:** maximize the log-probability that generated candidates satisfy the properties of interest (structure is now a latent variable)

$$\sum_{X \in \text{source set}} \log \left[\sum_Y P(\text{target specs} | Y) P(Y | X; \theta) \right]$$

- ▶ **E-step:** generate candidates from the current model; filter/reweight by property predictor (\sim posterior samples)
- ▶ **M-step:** maximize the log-probability of new (weighted) targets

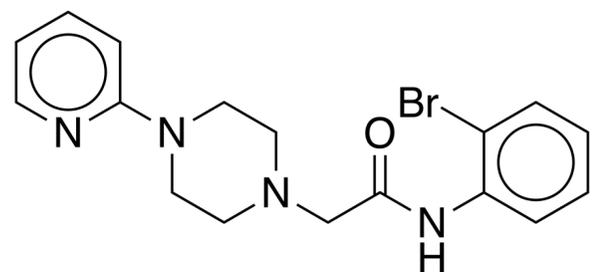


Target augmentation = stochastic EM

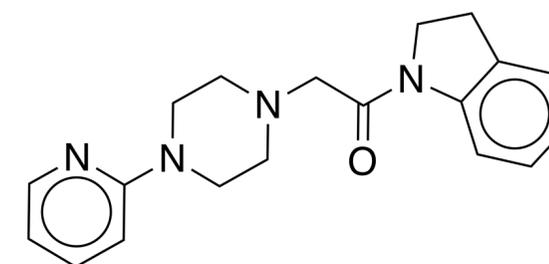
- ▶ **Objective:** maximize the log-probability that generated candidates satisfy the properties of interest (structure is now a latent variable)

$$\sum_{X \in \text{source set}} \log \left[\sum_Y P(\text{target specs} | Y) P(Y | X; \theta) \right]$$

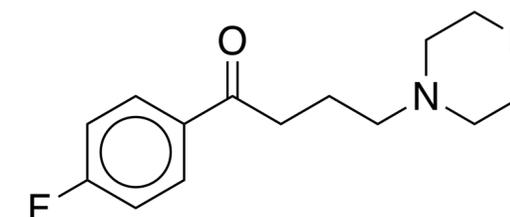
- ▶ **E-step:** generate candidates from the current model; filter/reweight by property predictor (\sim posterior samples)
- ▶ **M-step:** maximize the log-probability of new (weighted) targets



Generate



DRD2=0.933

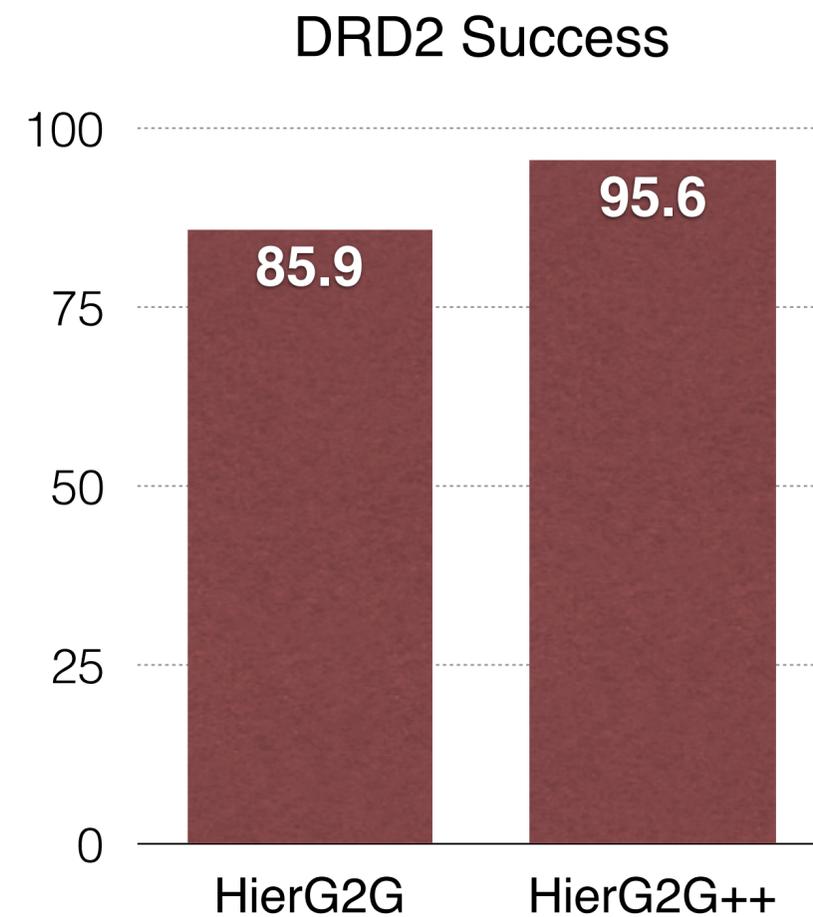
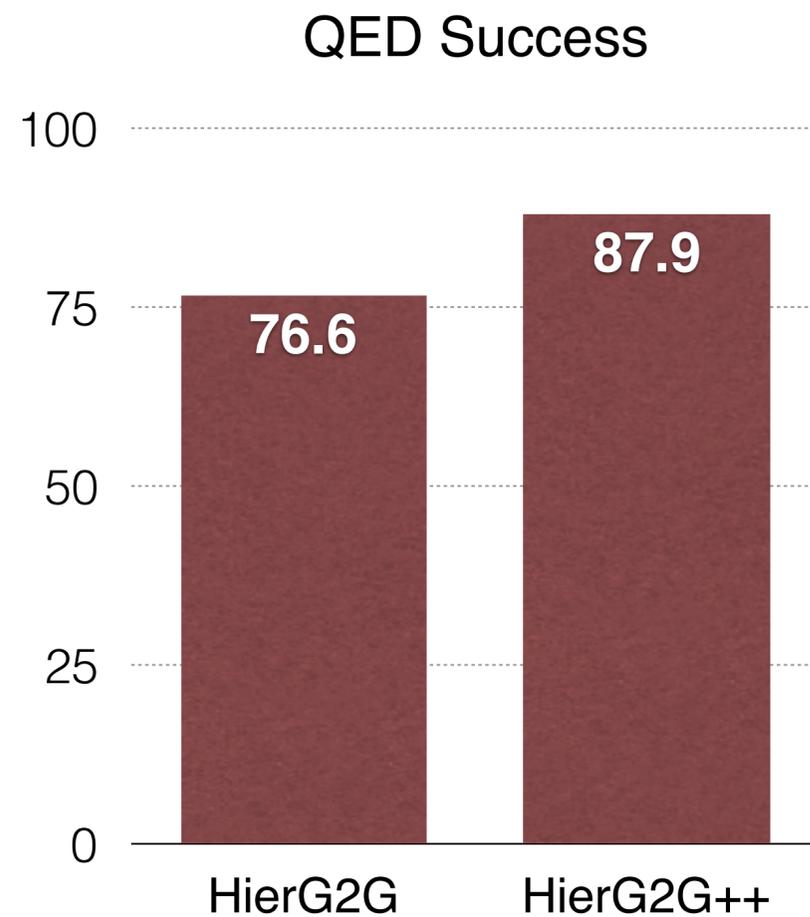


DRD2=0.008



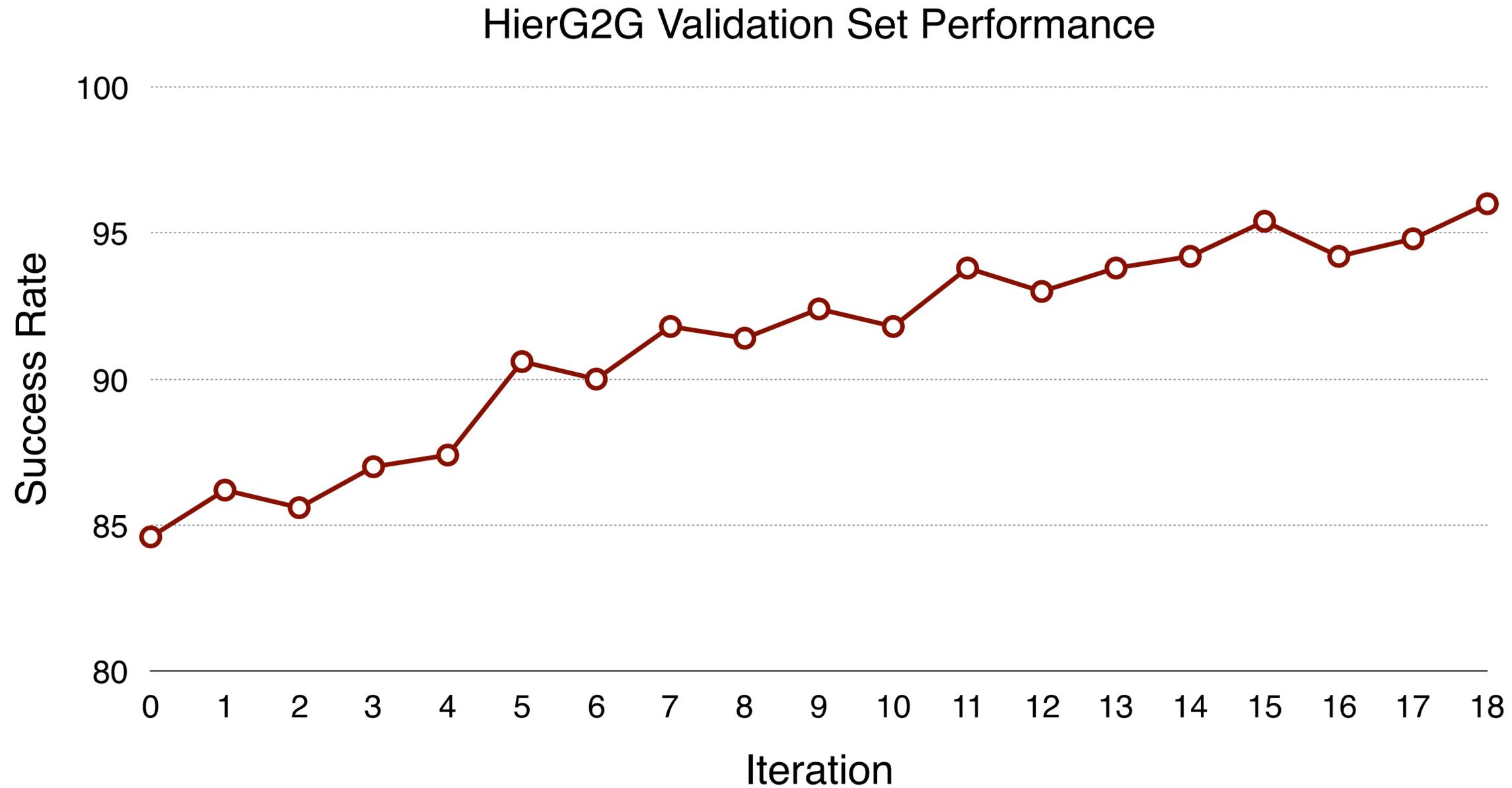
Example results: gains

- ▶ Substantial gains in translation/optimization success %



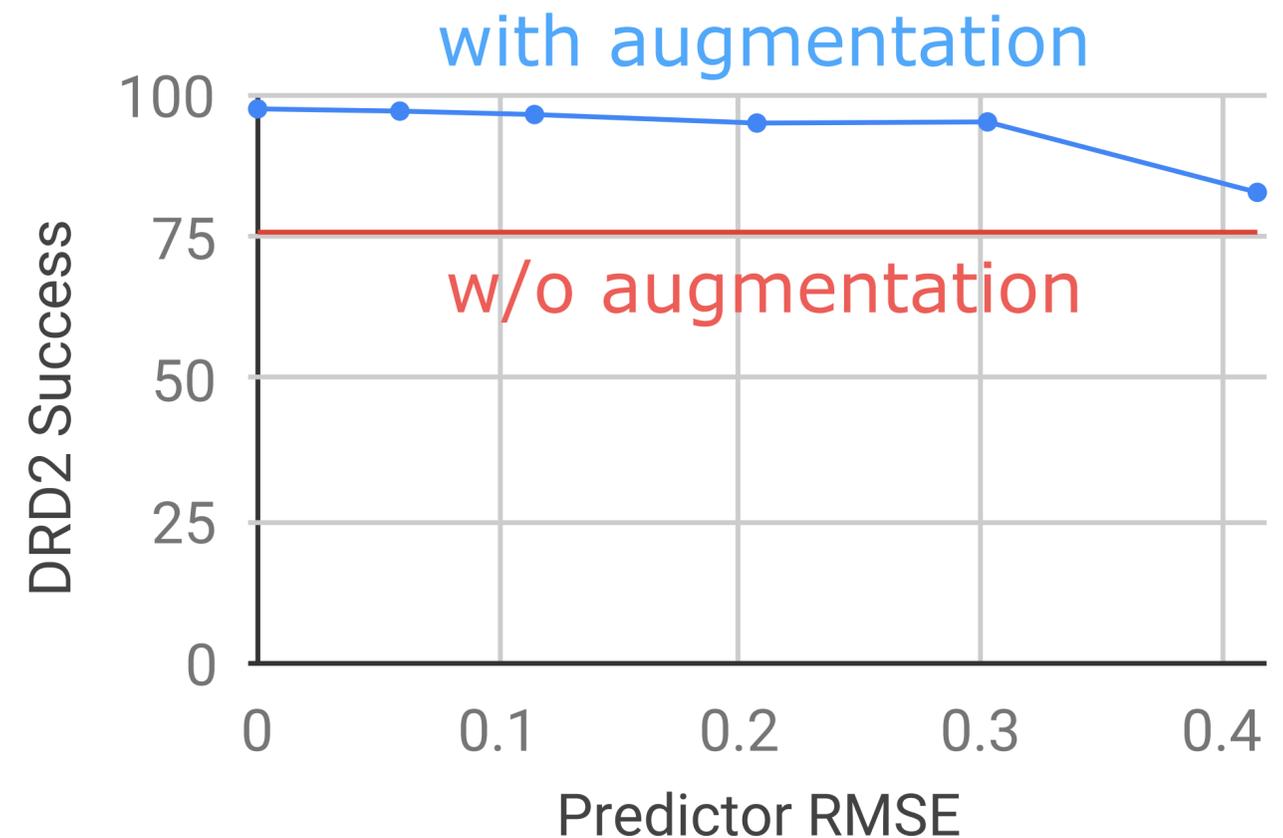
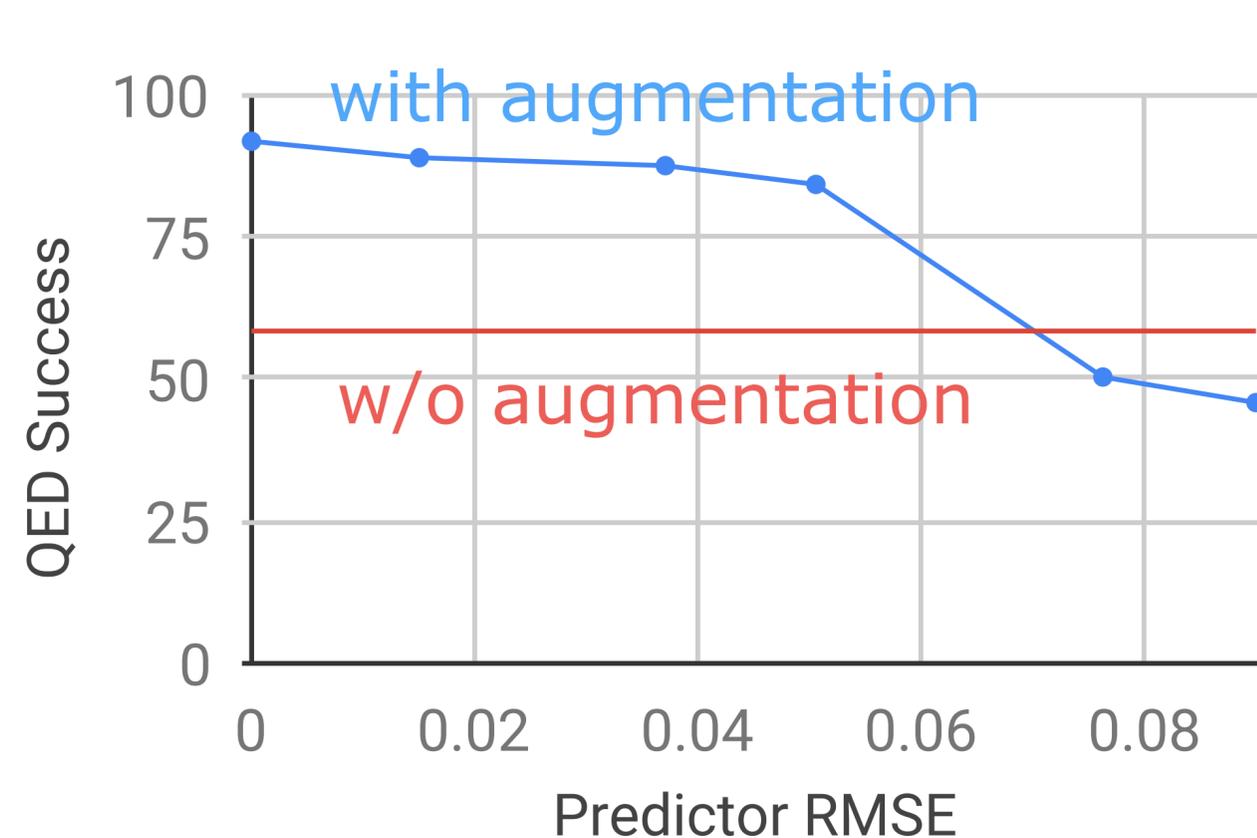
Example results: gains

- ▶ Consistently improving ...



Example results: robustness

- ▶ The gains are robust against errors in the property predictor



- ▶ Note: curves are for a weaker seq2seq model; baseline performance is much lower, but final performance with augmentation comparable to hierG2G.

Summary

- ▶ Molecules as structured objects provide a rich domain for developing ML tools; key underlying challenges shared with other areas involving generation/manipulation of diverse objects
- ▶ ML molecular design methods are rapidly becoming viable tools for drug discovery
- ▶ Several key challenges remain, however:
 - effective multi-criteria optimization
 - incorporating 3D features, physical constraints
 - generalizing to new, unexplored chemical spaces (domain transfer)
 - explainability, etc.