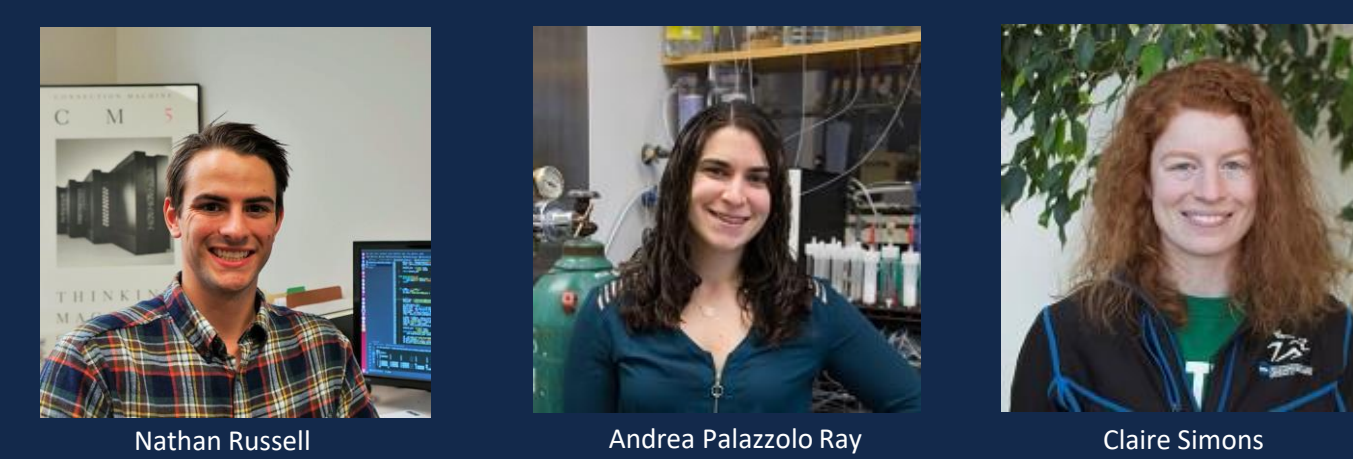


Synthetic Planning and Library Optimization for Automated Cross Coupling Synthesis Platforms

Nathan Russell^{1,†}, Andrea Palazzolo Ray^{2,†}, Claire Simons², Martin D. Burke², Jian Peng¹
Department of Computer Science¹ & Department of Chemistry², University of Illinois at Urbana-Champaign. † Equal Contribution.
Corresponding author: ntrusse2@illinois.edu



Introduction

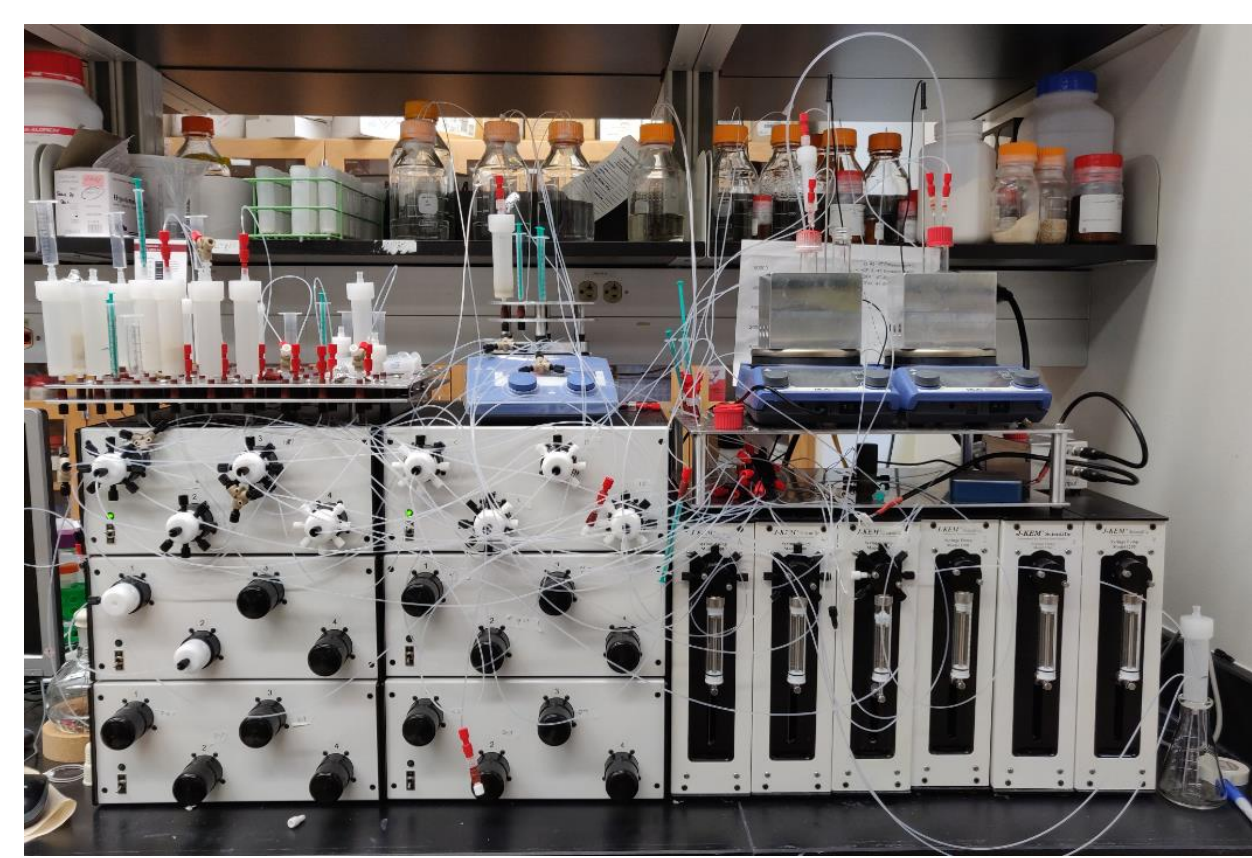
The Goal: Automatable Synthesis of Most Natural Product (NP) Chemical Space via Iterative Cross Coupling (ICC)

Our Contributions:

- Algorithms for ICC compatible synthesis planning of large target libraries and building block library selection
- Impact prioritized list of building blocks which enable efficient synthesis of most link-cyclic NP chemical space
- Characterization of cross coupling reactions and substrate scopes ranked by impact to synthesis of NP chemical space

The Approach:

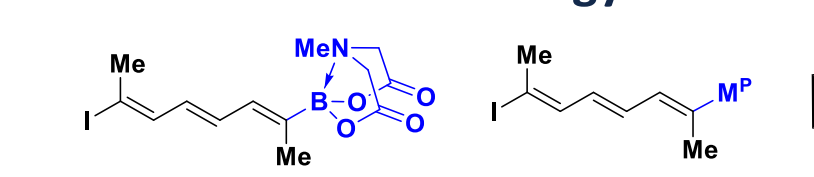
- Gather a representative set of 100k Link-cyclic Natural Products
- Given the constraints of iterative cross coupling, generate synthetic plans, allowing all valid blocks and couplings
- Find the smallest block library which enables coverage of the most NP chemical space using the fewest reaction steps
- Given block set, identify an "optimal" synthetic plan per target, and record the coupling reactions involved



Automated Cross Coupling Synthesis Platform

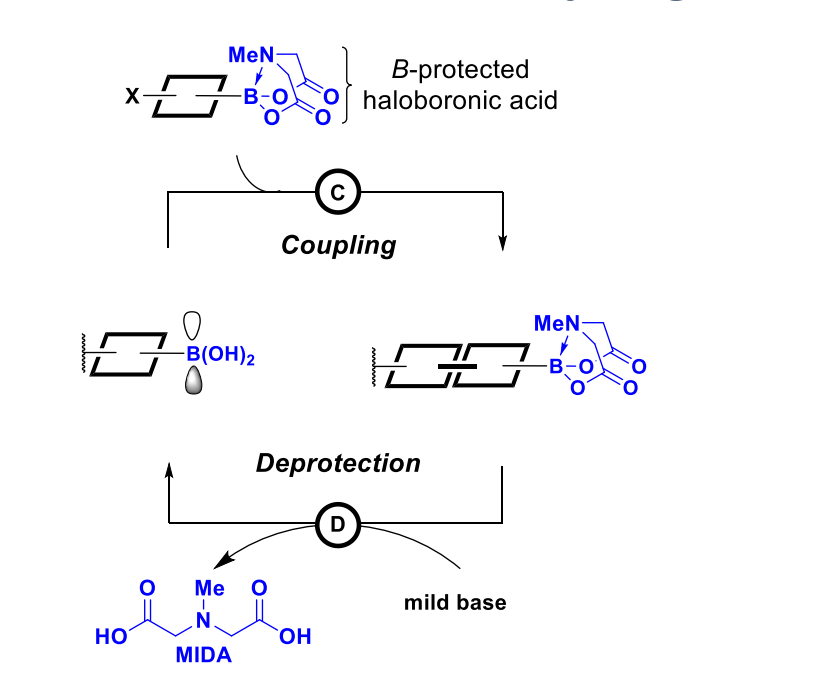
Synthesis Constraints

Notation & Terminology

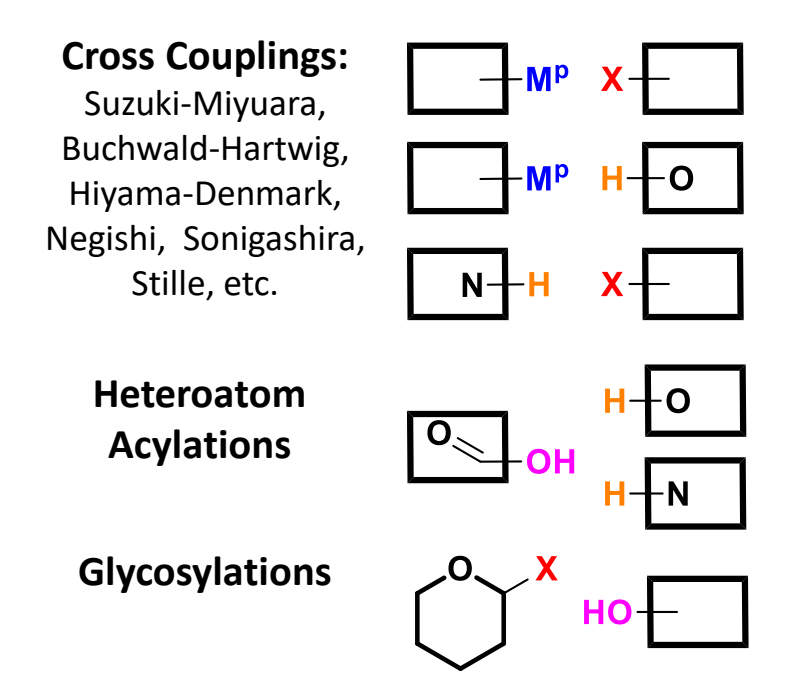


Termini: reactive group facilitates and is consumed during coupling, possibly protected
Block: molecular fragment with preinstalled termini

Iterative Cross coupling



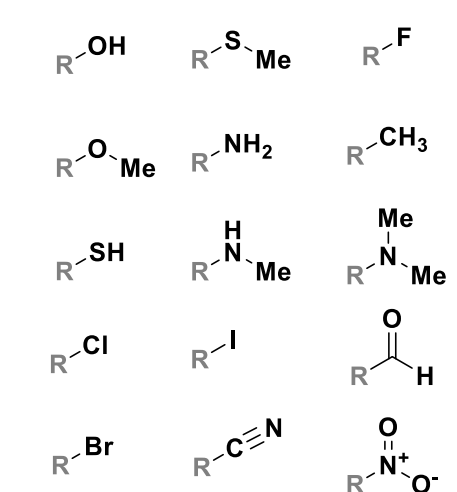
Coupling Reactions Allowed



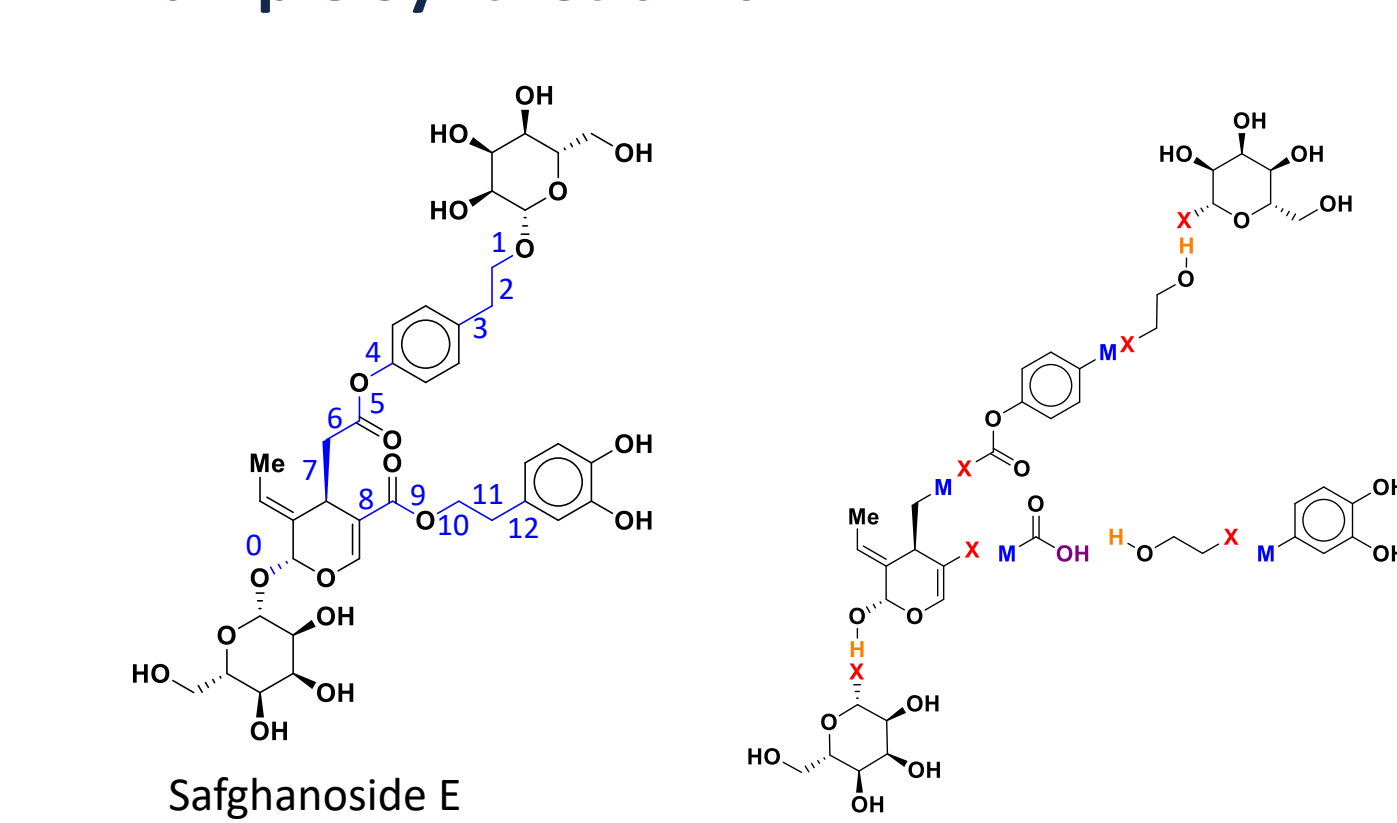
Primary Synthesis Constraints

- Only Rotatable Bonds can be broken
- Bonds to R Groups can NOT be broken
- Blocks must have at least 2 heavy atoms
- Blocks can have at most 3 Termini*
- An atom can have at most 2 termini
- Post Coupling Product, must have at least 1 of the following MP¹, MP², OH, H, or equivalent for Purification turnover unless terminal coupling in plan

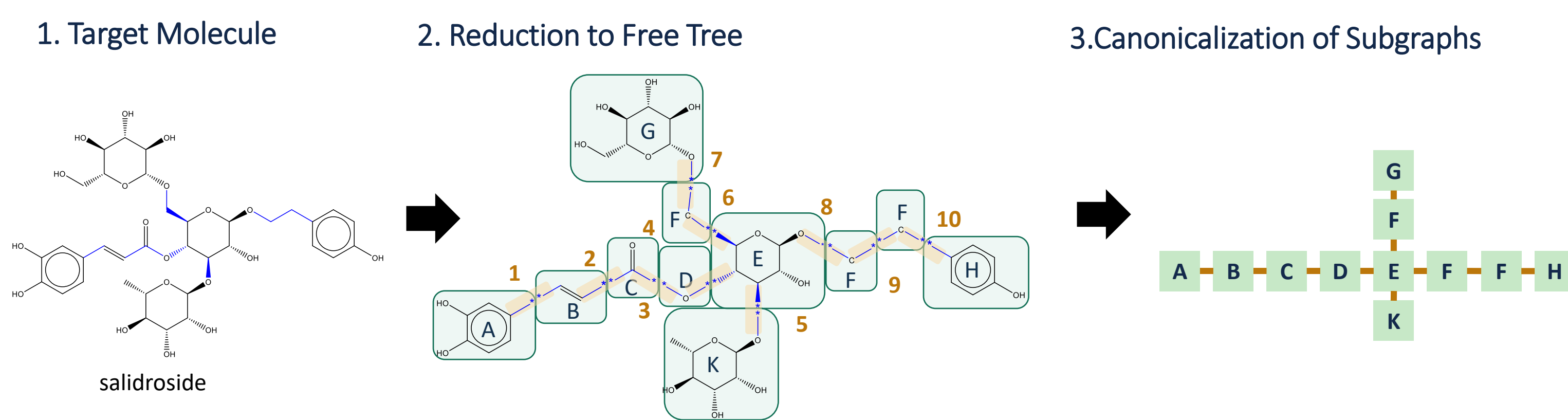
R Groups



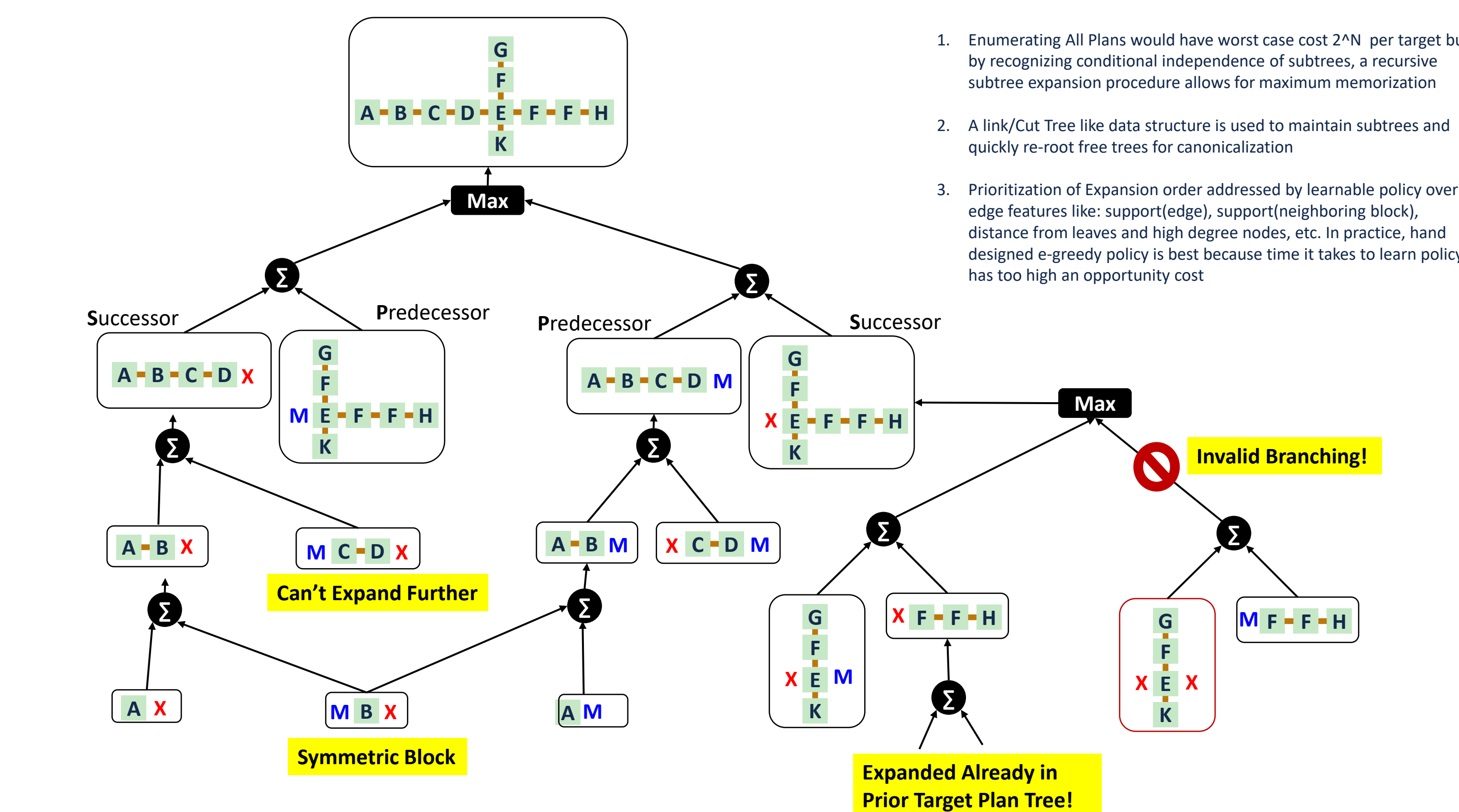
Example Synthetic Plan



Molecule Representation



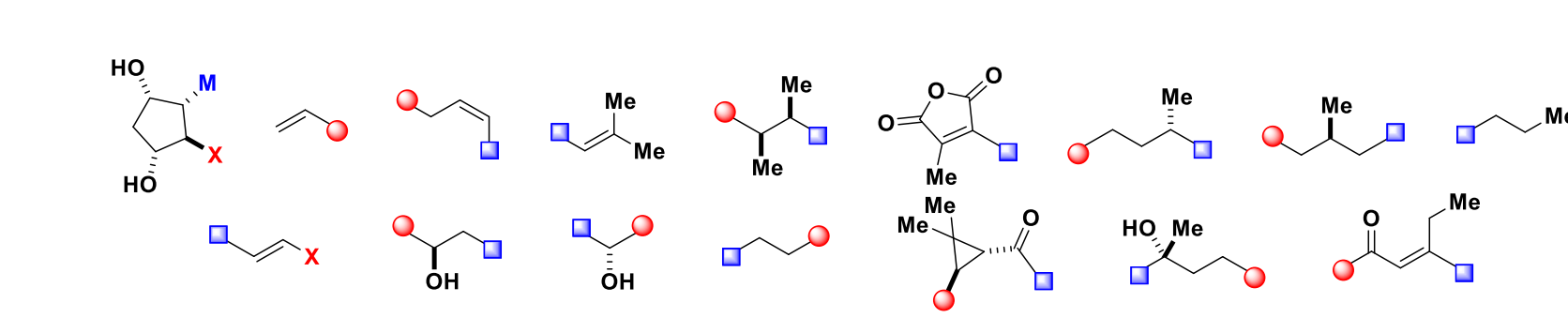
Synthesis Plan Tree Expansion



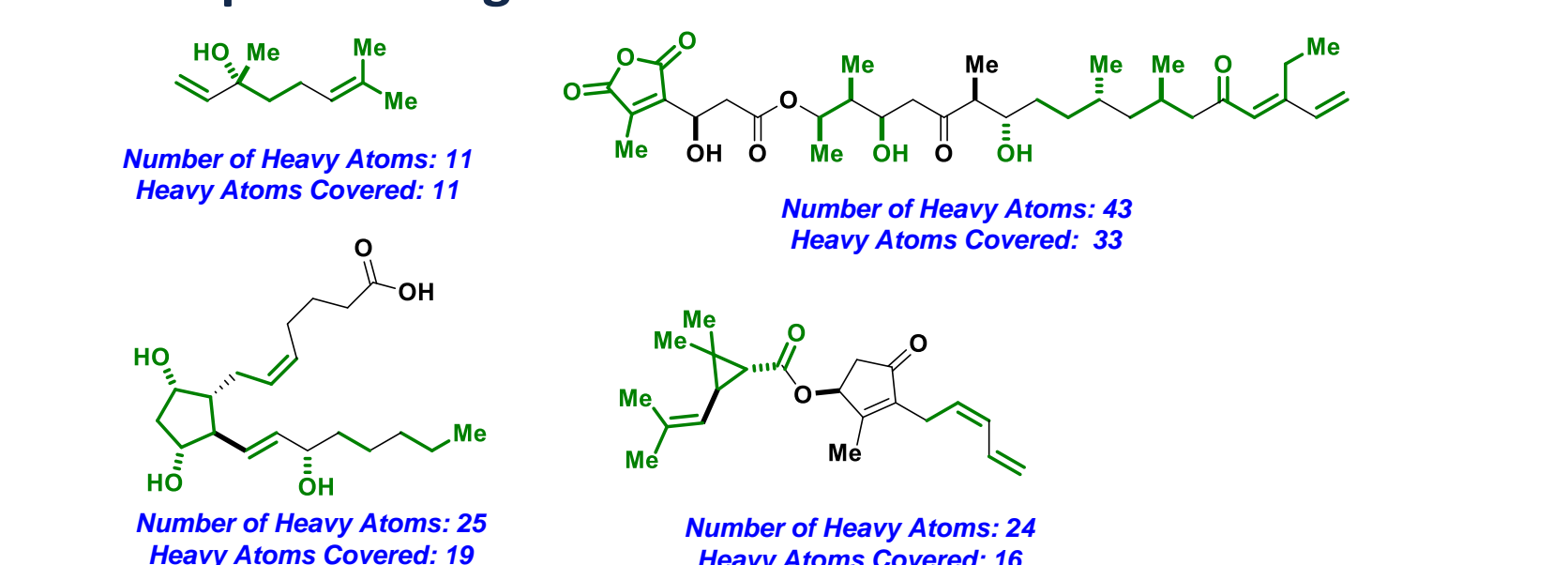
- Enumerating All Plans would have worst case cost 2^N per target but by recognizing conditional independence of subtrees, a recursive subtree expansion procedure allows for maximum memorization
- A link/Cut Tree like data structure is used to maintain subtrees and quickly re-root free trees for canonicalization
- Prioritization of Expansion order addressed by learnable policy over edge features like: support(edge), support(neighboring block), distance from leaves and high degree nodes, etc. In practice, hand designed e-greedy policy is best because time it takes to learn policy has too high an opportunity cost

Optimization / Block + Plan Selection

Candidate Block Set



Sample Coverage



Optimization Problem

Coverage **Library Size** **Tractability**

$$\max_x \sum_{i=1}^N z^i \quad \min_x \sum_{k=1}^K x^k \quad \max_x \sum_{i=1}^N \sum_{j \in SP^i} (y_j^i - 1) n c_j^i$$

s.t.

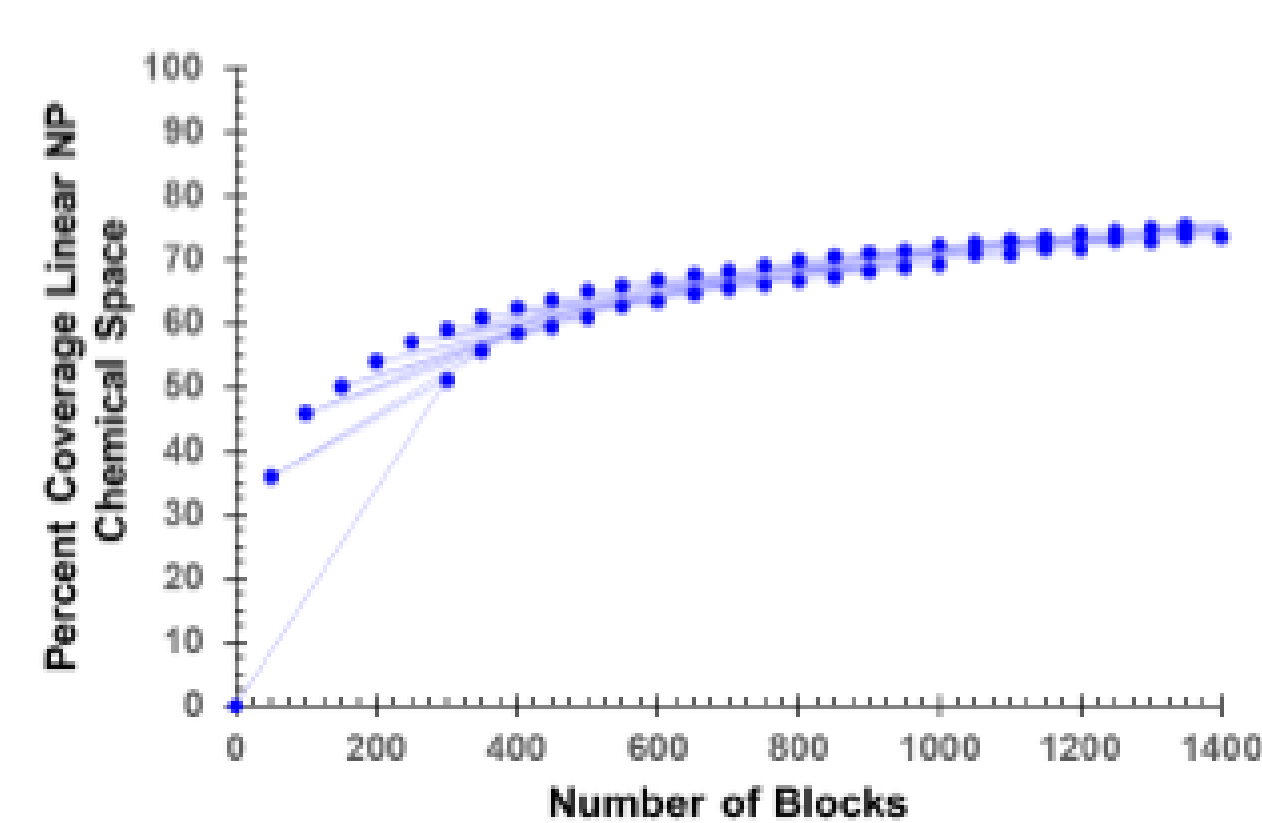
$$z^i \leq \sum_{j \in SP^i} \bar{x}_j^i + M y_j^i \quad \forall j \in SP^i \quad x^k \in \{0,1\} \quad \forall k$$
$$\sum_{j \in SP^i} y_j^i = |SP^i| - 1 \quad \forall i$$

$n c_j^i$ is the number of couplings for synthetic plan i, j
 z^i is the maximum coverage achievable for natural product i
 \bar{x}_j^i is the vector of block decision variables for synthetic plan i, j
 \bar{w}_j^i is the vector of heavy atom count constants for synthetic plan i, j
 M is a Constant $\geq \max(\bar{x}_j^i) - \min(\bar{w}_j^i) \quad \forall i, j$
 $N \geq 10^5$: # of Natural Products, indexed by i
 $K \geq 10^7$: # of Block Candidates, indexed by k
 $SP^i \approx 2^{n_{breakable}}$: Set of candidate synthetic plans, plans indexed by j .

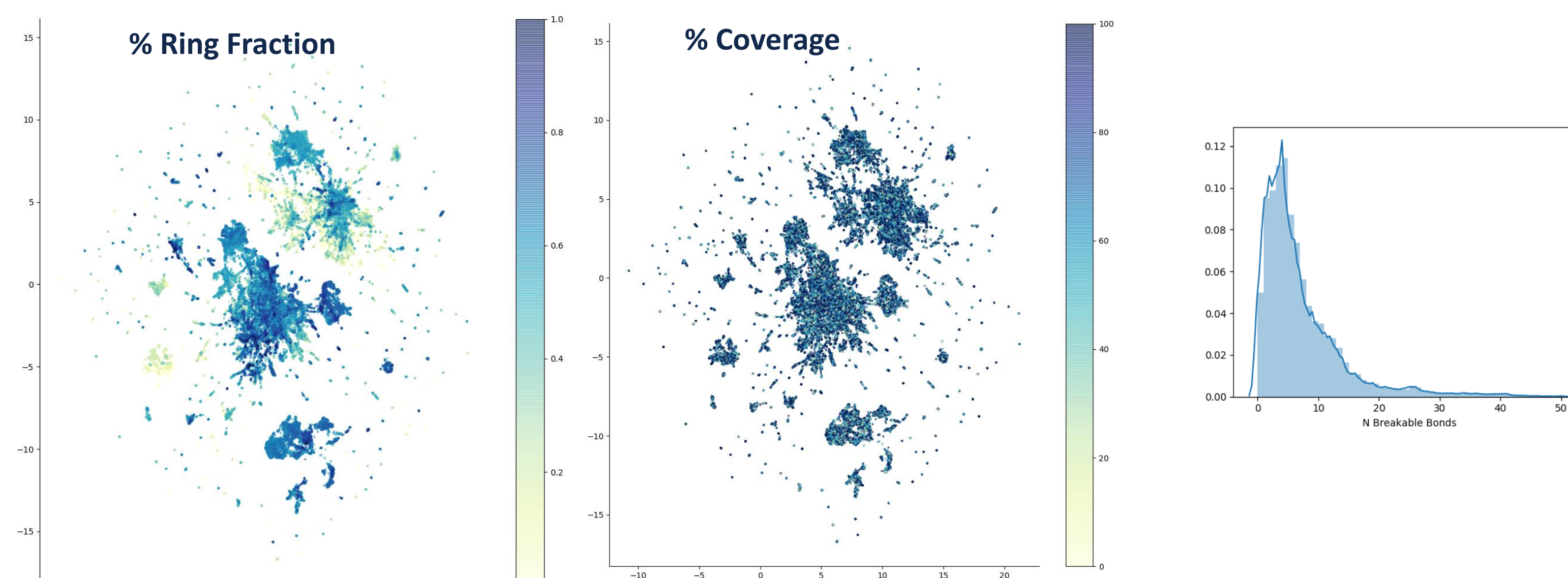
Algorithm: Bidirectional Lazy Greedy Subset Selection for Blocks

Sort blocks by descending marginal gain per block, maintain in a MaxHeap
While total_coverage < N%
For f Forward Steps:
max_gain = 0
for block in sorted_block_heap:
gain = conditional_marginal_gain(block, A)
max_gain = max(max_gain, gain)
if max_gain >= UB_gain:
Add max_block, break
Add max_block
For b Backward Steps:
add block with min loss upon removal from A

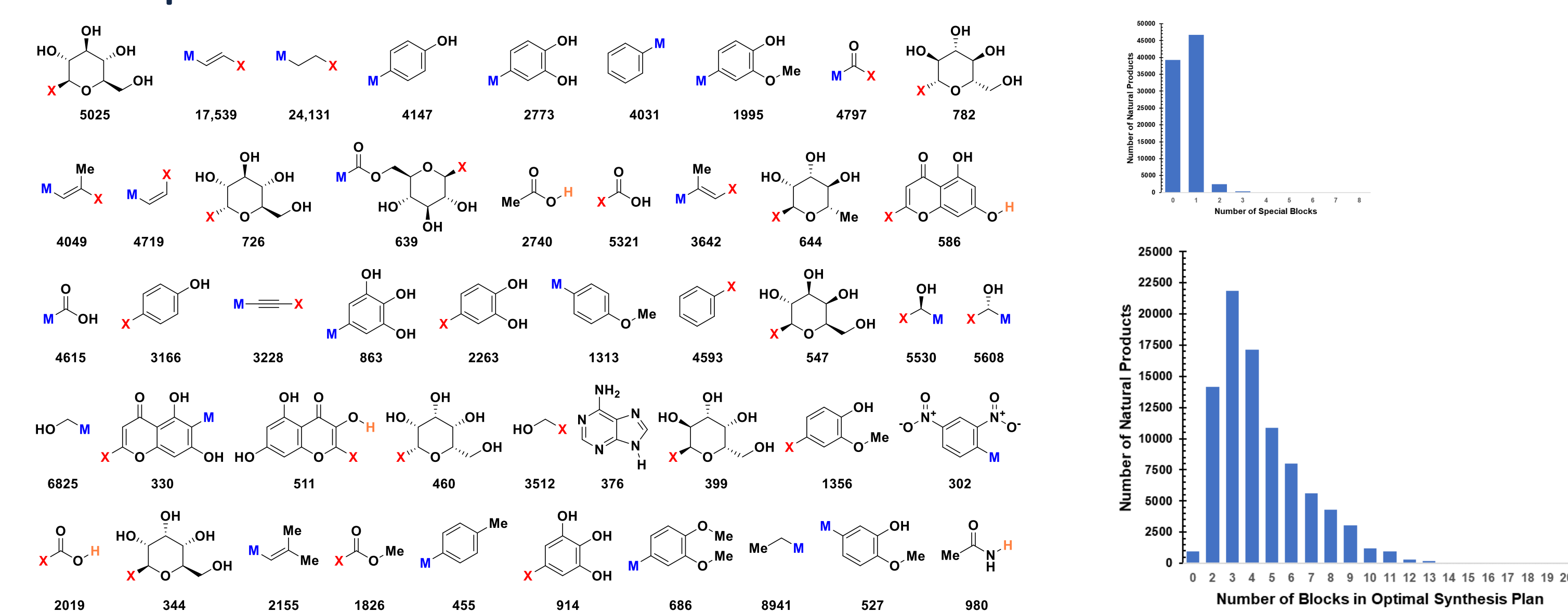
Optimization Curve



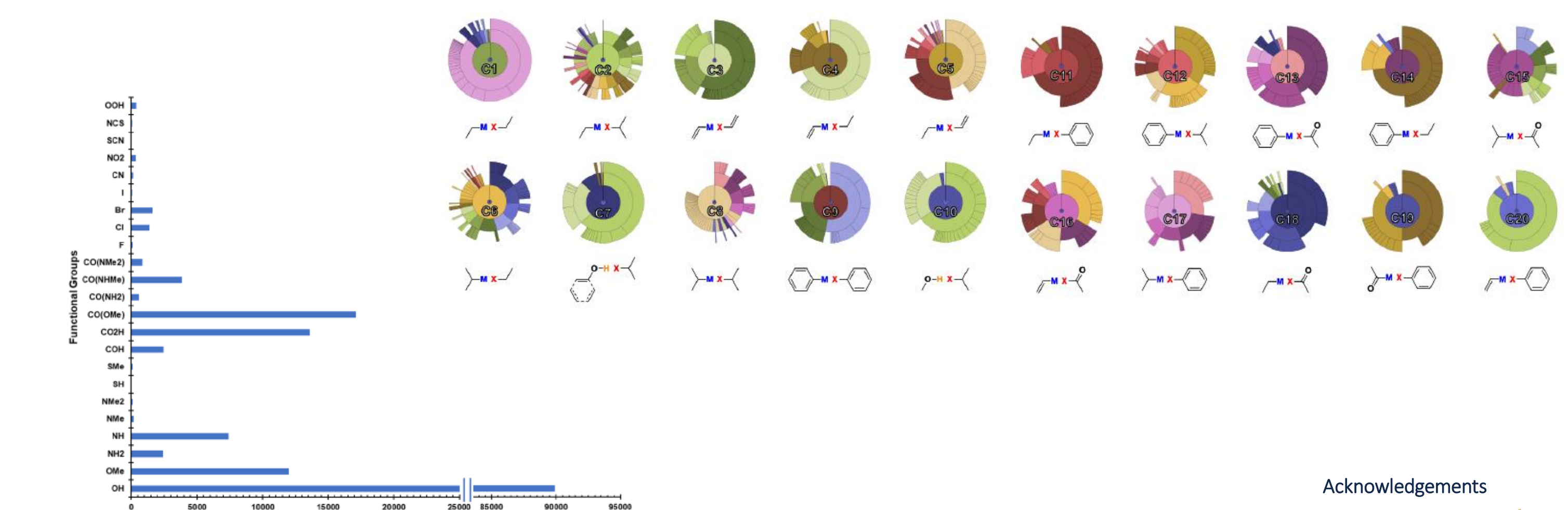
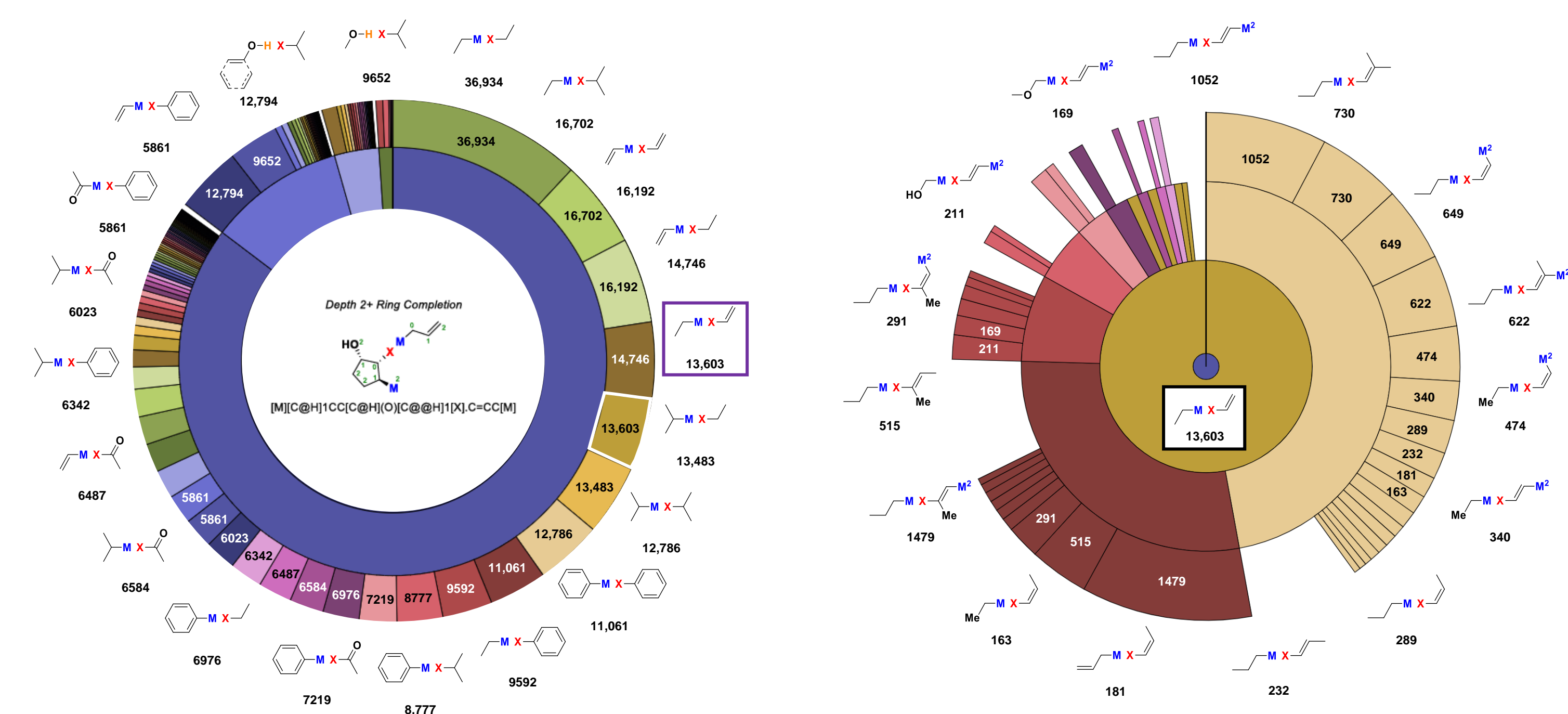
Natural Product Space



Example Blocks



Prioritized Cross Coupling Substrate Scopes



Acknowledgements

