

# Mathematical Formulation of Intermediate Generation

Let:

- $\mathcal{T}_V$ : a finite set of *element types* (e.g., types of vertices or nodes)
- $\mathcal{T}_E$ : a finite set of *edge types* (e.g., bond types or connection types)

Define a graph  $G = (V, E, \tau_V, \tau_E)$  where:

- $V$  is a finite set of *vertices* (elements)
- $E \subseteq \{\{u, v\} \mid u, v \in V, u \neq v\}$  is a set of *undirected edges* (no loops or multi-edges)
- $\tau_V : V \rightarrow \mathcal{T}_V$  assigns a type to each vertex
- $\tau_E : E \rightarrow \mathcal{T}_E$  assigns a type to each edge

## Equivalence of Graphs

Two graphs  $G_1 = (V_1, E_1, \tau_{V_1}, \tau_{E_1})$  and  $G_2 = (V_2, E_2, \tau_{V_2}, \tau_{E_2})$  are *isomorphic* (or equivalent) if there exists a bijection

$$\phi : V_1 \rightarrow V_2$$

such that:

$$\begin{aligned} \tau_{V_1}(v) &= \tau_{V_2}(\phi(v)) \quad \text{for all } v \in V_1 \\ \{u, v\} \in E_1 &\iff \{\phi(u), \phi(v)\} \in E_2 \\ \tau_{E_1}(\{u, v\}) &= \tau_{E_2}(\{\phi(u), \phi(v)\}) \end{aligned}$$

This ensures both structure and labeling are preserved under isomorphism.

## Subgraph Definition

A graph  $H = (V_H, E_H, \tau_{V_H}, \tau_{E_H})$  is a *subgraph* of  $G$  (denoted  $H \subseteq G$ ) if:

- $V_H \subseteq V$
- $E_H \subseteq \{\{u, v\} \in E \mid u, v \in V_H\}$
- $\tau_{V_H} = \tau_V|_{V_H}$ , i.e., the restriction of  $\tau_V$  to  $V_H$
- $\tau_{E_H} = \tau_E|_{E_H}$

Trivially,  $G \subseteq G$ .

## Question: How to Get All Non-Trivial Unique Subgraphs of $G$ ?

Let  $G = (V, E, \tau_V, \tau_E)$ . To enumerate all non-trivial (at least one vertex) unique subgraphs up to isomorphism, with the additional constraint that any subgraph of size at least two must be *connected* as a single component (i.e., the subgraph must not be separable into disconnected parts):

### Step-by-Step Procedure

1. **Enumerate All Vertex Subsets:**

For every non-empty subset  $V' \subseteq V$ , proceed to edge enumeration.

2. **Enumerate All Edge Subsets:**

For each vertex subset  $V'$ , consider the induced edge set:

$$E' = \{\{u, v\} \in E \mid u, v \in V'\}$$

3. **Filter by Connectivity:**

Accept the subgraph  $H = (V', E')$  only if:

- If  $|V'| = 1$ , the single vertex must have degree at least 1 in the original graph  $G$
- If  $|V'| \geq 2$ , the induced subgraph  $H$  must be *connected*, meaning there exists a path between any pair of nodes in  $H$

4. **Construct Labeled Subgraphs:**

Each valid  $(V', E')$  pair defines a candidate subgraph  $H$  with inherited labeling functions  $\tau_V|_{V'}$  and  $\tau_E|_{E'}$ .

5. **Canonical Labeling (Isomorphism Reduction):**

Use a graph isomorphism algorithm that respects vertex and edge types to convert each subgraph into a canonical form.

6. **Filter Unique Subgraphs:**

Maintain a set of canonical representations to eliminate isomorphic duplicates.

### Output

The set of all non-trivial (i.e.,  $|V'| \geq 1$ ) and connected (if  $|V'| \geq 2$ ) unique subgraphs of  $G$ , up to isomorphism.

## Problem Motivation and Hashing Approach

Directly comparing each candidate subgraph via isomorphism testing is computationally expensive, especially as the number of candidate subgraphs grows exponentially with graph size.

To address this, we make use of the *Weisfeiler-Lehman (WL) graph hash* as a surrogate for canonical labeling. This enables us to efficiently test whether two graphs are isomorphic up to node and edge types, without performing pairwise isomorphism tests.

### Why Full Isomorphism Checking is Expensive

Let  $\mathcal{G}_k$  denote the set of all size- $k$  subgraphs from the input graph  $G$ . To identify all unique subgraphs up to isomorphism, a naïve approach would require  $O(|\mathcal{G}_k|^2)$  pairwise isomorphism checks, where each check can be expensive depending on labeling.

### Weisfeiler-Lehman Graph Hashing

The Weisfeiler-Lehman method is a vertex refinement technique used to encode graph structure into compact hash values. It works as follows:

- Initially, each node is labeled with its original label (e.g., node type).
- At each iteration, the label of each node is updated by hashing its current label together with a multiset of its neighbors' labels (along with edge types, if specified).
- This process iteratively refines the labels in a way that captures structural information, producing a final string representation that is invariant under graph isomorphism.

This process yields a graph-invariant fingerprint: if two graphs have different WL hashes, they are guaranteed to be non-isomorphic. If they share the same hash, they are *likely* isomorphic (WL test is not complete for all graphs, but is highly effective in practice, especially for graphs with rich labelings).

## Implementation of Weisfeiler-Lehman Hashing in `networkx`

In practice, the Weisfeiler-Lehman (WL) graph hash can be computed efficiently using the `networkx` library, which provides a built-in function called `weisfeiler_lehman_graph_hash`.

## API Call and Usage

The function is accessed as follows:

```
from networkx.algorithms.graph_hashing import weisfeiler_lehman_graph_hash

hash_string = weisfeiler_lehman_graph_hash(
    G,
    node_attr="type",
    edge_attr="type"
)
```

### Parameters:

- **G**: the input graph (can be a subgraph)
- **node\_attr**: name of the node attribute to use in hashing (e.g., "type")
- **edge\_attr**: name of the edge attribute to use (e.g., "type")

**Output:** A canonical string that encodes the graph's topology and labeled structure. Two graphs with the same hash string are treated as equivalent (isomorphic with respect to labels).

## Preprocessing and Notes

- To ensure that node labels do not affect the hash (only structure and types), each subgraph is relabeled to use integer node indices before hashing:

```
H.relabel = networkx.convert_node_labels_to_integers(H)
```

- To maximize performance, type labels are stored as integers (or strings) rather than complex objects (e.g., Enums).
- The hash values are used to deduplicate subgraphs during enumeration:

```
if wl_hash not in seen_hashes:    ...
```

## Complexity Analysis

The problem of enumerating all non-trivial, fully connected, unique subgraphs of a labeled graph  $G = (V, E, \tau_V, \tau_E)$  involves several computational steps, each with its own complexity considerations.

## Subgraph Enumeration

The number of labeled subgraphs is exponential:

$$O(2^{|V|} \cdot 2^{|E|})$$

and each subgraph requires isomorphism checking, so efficient hashing or pruning is critical for tractability on large graphs.

## Overall Complexity of Weisfeiler-Lehman Hash Computation

For each valid subgraph  $H \subseteq G$ , a WL hash is computed. The runtime for WL hashing is:

$$O(k(|V_H| + |E_H|))$$

where  $k$  is the number of refinement iterations (typically a small constant), and  $V_H$ ,  $E_H$  are the vertex and edge sets of the subgraph. Since each subgraph is typically small, this step is efficient and can be considered nearly linear per subgraph.

After computing the hash string for each subgraph, we store it in a hash set to check for uniqueness. This provides:

$$O(1)$$

expected time per subgraph for lookup and insertion (amortized), assuming good hash distribution and a suitable hash function (e.g., Weisfeiler-Lehman).

Let  $N$  denote the number of candidate subgraphs generated and tested. The overall complexity is:

$$O(N \cdot (|V_H| + |E_H|)) \quad (\text{WL hash per subgraph})$$

where  $|V_H|$  and  $|E_H|$  are typically small constants for molecular graphs or motifs.

In the worst case (e.g., for dense graphs with many small motifs), the algorithm remains exponential in  $|V|$ , but the combination of pruning (via degree filtering) and hashing makes it tractable for moderate-sized sparse graphs.

## Subgraph Models: Induced vs. Edge-Deletion Connected Subgraphs

Subgraphs can be yielded with varying levels of structural flexibility. In this work, we distinguish between two valid but distinct method of getting subgraph, which differ in the set of allowed edges among selected vertices.

## 1. Induced Connected Subgraphs

Let  $G = (V, E, \tau_V, \tau_E)$  be the full input graph.

A subgraph  $H = (V_H, E_H, \tau_V|_{V_H}, \tau_E|_{E_H})$  is called an **induced connected subgraph** if:

- $V_H \subseteq V, V_H \neq \emptyset$
- $E_H = \{\{u, v\} \in E \mid u, v \in V_H\}$
- $H$  is connected

That is, for each subset of nodes  $V_H$ , the subgraph includes all edges among those nodes from the original graph. Only those node subsets where the induced graph is connected are retained.

**Consequence:** This formulation does not allow partial deletion of edges within a node subset. If two nodes in  $V_H$  are connected in  $G$ , they must be connected in  $H$ .

## 2. Edge-Deletion Connected Subgraphs

We now generalize to a more inclusive definition.

A subgraph  $H = (V_H, E_H, \tau_V|_{V_H}, \tau_E|_{E_H})$  is called an **edge-deletion connected subgraph** if:

- $V_H \subseteq V, V_H \neq \emptyset$
- $E_H \subseteq \{\{u, v\} \in E \mid u, v \in V_H\}$
- $H$  is connected

Here,  $H$  is a connected subgraph formed by selecting any non-empty subset of the edges among  $V_H$ . Unlike the induced case, not all internal edges are required to be included. This allows for a richer space of meaningful substructures, such as loops missing a bond or partial motifs.

## Mathematical Implication

Given an induced subgraph  $G_2 = (V, E, \tau_V, \tau_E)$ , let  $G_1 = (V, E' \subset E, \tau_V, \tau_E|_{E'})$  be formed by deleting one or more edges while maintaining connectivity.

Then  $G_1 \subseteq G_2$ , and

$$G_2 - G_1 = (\emptyset, \Delta E = E \setminus E', \emptyset, \Delta \tau_E)$$

Such a  $G_1$  is *excluded* from the induced subgraph model, but *included* in the edge-deletion connected model.

## Computational Consequence

- **Induced subgraphs** are simpler to enumerate: only vertex subsets are needed, and the edge set is fully determined by the node set.
- **Edge-deletion subgraphs** require enumeration over both node subsets and all non-empty, connected subsets of the induced edge set.
- Thus, the number of connected subgraphs grows much faster under the edge-deletion model.

## Implementation Differences

- Induced model: For each  $V_H$ , construct the induced subgraph using `G.subgraph(V_H)` and check connectivity.
- Edge-deletion model: For each  $V_H$ , enumerate all edge subsets  $E_H \subseteq E[V_H]$  such that  $(V_H, E_H)$  is connected.

To efficiently generate only connected subgraphs without checking all edge subsets, we implement a recursive DFS-style edge-growing strategy that builds connected subgraphs from individual edges, avoiding disconnected candidates entirely.

## Complexity Comparison

Let  $n = |V|$  and  $m = |E|$ .

- **Induced connected subgraphs:**

$$O(2^n)$$

since we only consider non-empty vertex subsets and check if the induced subgraph is connected.

- **Edge-deletion connected subgraphs:**

$$O(2^n \cdot 2^e) \quad (\text{in worst case, for each vertex subset, all edge subsets})$$

but pruned via connectivity. Using DFS-based generation limits this growth to only viable connected subgraphs.

- **Deduplication:** In both cases, isomorphism or hash-based filtering (e.g., Weisfeiler-Lehman graph hashing) is required to ensure uniqueness.

## Summary

- We use induced subgraphs to generate the list of complexes to significantly reduce the computational cost.
- The scaling is still exponential but now only to the number of vertices, not edges.