
ANSWERING COUNTERFACTUAL QUERIES ON GRAPH DATABASES

Anonymous authors

Paper under double-blind review

ABSTRACT

Counterfactual analysis on graph data is central to causal reasoning and interpretability, yet existing graph-based methods rely on ad hoc perturbations and remain tied to model behavior rather than underlying data. To address this challenge, we introduce **Counterfactual Graph Database (CF-GDB) queries**, the first query-based framework for counterfactual reasoning on graphs that grounds counterfactuals in verifiable database instances. Our approach abstracts graphs into semantically meaningful concepts and compares them using a hypergraph-based distance that integrates local structure with global semantics. To ensure efficiency and scalability, we propose two complementary indices: the *Concept Distribution Index (CDI)*, a histogram that provides certified lower bounds, and the *Concept Semantic Index (CSI)*, a continuous embedding that provides upper bounds. These indices yield provably tight sandwich guarantees and enable efficient candidate pruning while preserving the fidelity of counterfactual retrieval. Using 8 read data sets across 4 domains, CF-GDB improves accuracy by over 20% and achieves up to 20× faster performance, demonstrating both fidelity and scalability.

1 INTRODUCTION

Counterfactual analysis (Rubin, 1974) has found broad applications in machine learning (Wachter et al., 2017) and data science (Karimi et al., 2021). Unlike factual explanations that justify observed outcomes, counterfactuals identify minimal yet semantically valid perturbations sufficient to alter a model’s prediction (Wachter et al., 2017; Mothilal et al., 2020). Such perturbations correspond to actionable alternatives, such as treatment adjustments in healthcare (Zhang et al., 2025), portfolio rebalancing in finance (Zhang et al., 2025), or precedent revisions in law (Zhang et al., 2023), providing evidence that is both interpretable and verifiable. Moreover, counterfactuals improve generalization (Tan et al., 2022), enhance robustness (Bajaj et al., 2021), and are indispensable in high-stakes domains where explanations must translate into practical guidance.

Recent work has explored counterfactual analysis in graph learning (Ma et al., 2022; Verma et al., 2024b; He et al., 2024; Fournier & Medya, 2025; Giorgi et al., 2025) by perturbing nodes, edges, or attributes to generate alternative graphs. However, these methods primarily focus on *model behavior* rather than *data-grounded evidence*, and face three key limitations: (i) **model dependence** on gradients or weights (Bajaj et al., 2021; Armgaan et al., 2024; Lucic et al., 2022), which confines them to white-box settings and renders counterfactuals unstable in practice (McCoy et al., 2022; Shu et al., 2024); (ii) **domain validity**, as arbitrary edits often violate structural constraints, leading to, for example, invalid molecules (Zhang et al., 2025) or unrealistic transactions (Gan et al., 2021); and (iii) **limited transferability**, since instance-specific edits lack grounding in database schemas, making them difficult to reuse across queries or datasets (Giorgi et al., 2025).

In this paper, as our first contribution, we propose the **Counterfactual Graph Database (CF-GDB)**, a novel framework that reframes counterfactual reasoning as a query problem over graph databases. Unlike prior approaches that generate perturbed graphs to flip model predictions (Ying et al., 2019; Lucic et al., 2022), CF-GDB retrieves dataset-grounded,

054 domain-valid counterfactuals, ensuring explanations are anchored in verifiable instances—an
055 essential requirement in high-stakes domains such as law (Zhang et al., 2023), finance (Gan
056 et al., 2021), and healthcare (McCoy et al., 2022). This task is challenging because coun-
057 terfactuals involve coordinated changes to node features, edges, and topology, where even
058 small edits can fundamentally alter semantics (e.g., substituting a carbon atom with ni-
059 trogen transforms a benzene ring into a pyridine ring, changing its toxicity (Zhang et al.,
060 2025)). Traditional similarity measures, such as graph edit distance (Gao et al., 2010) or
061 whole-graph embeddings (Bai et al., 2019; Zheng et al., 2014), are either too fine-grained
062 or too coarse, and are further constrained by Weisfeiler–Leman expressivity (Xu et al.,
063 2019; Morris et al., 2023), offering limited guidance for localizing meaningful counterfactual
064 differences (Armgaan et al., 2024; Giorgi et al., 2025).

065 To implement CF-GDB efficiently, as our second contribution, we introduce **Concept-**
066 **Based Counterfactual Graph Query (C²GQ)**, which reframes counterfactual reasoning
067 as a query problem in a shared **graph concept space**. Graph concepts serve as semantic
068 prototypes that cluster structurally and semantically similar subgraphs across a database
069 of graphs (e.g., rings in molecules, communities in networks), enabling application-aware
070 counterfactual reasoning beyond naive edge and vertex edits. Each graph is represented as a
071 distribution over concepts, and differences are measured by a **hypergraph-based concept**
072 **distance** grounded in unbalanced optimal transport (Vayer et al., 2020). This distance
073 jointly captures fine-grained local changes (e.g., removing a molecular ring) and global
074 distributional shifts (e.g., reconfiguring community structures), ensuring interpretability
075 and domain validity across multiple scales.

076 To achieve scalable counterfactual queries, as the third contribution, we introduce two in-
077 dices inspired by graph homomorphism theory (Dell et al., 2017): the **Concept Distribu-**
078 **tion Index (CDI)**, a histogram that aggregates concept counts to provide certified lower
079 bounds, and the **Concept Semantic Index (CSI)**, a continuous embedding that encodes
080 concept semantics to yield upper bounds. Beyond efficiency, our framework offers guaran-
081 tees absent in prior work, including query success rates, provably correct pairwise rankings,
082 and safe pruning rules that never discard valid counterfactuals. As a whole, these results
083 establish the first theoretical foundation for counterfactual queries on graphs, in contrast
084 to previous heuristic perturbation methods. Experiments using 8 real data sets across 4
085 domains demonstrate that C²GQ improves accuracy by over 20% with 20× faster queries.

086 2 RELATED WORK

088 **Counterfactual Analysis.** Counterfactual analysis explains model predictions by iden-
089 tifying minimal input changes that alter outcomes (Wachter et al., 2017), and has been
090 widely applied in high-stakes domains (Zhang et al., 2025; 2023). While effective for tabular
091 and textual data (Pawelczyk et al., 2020; Yang et al., 2021), graphs pose unique chal-
092 lenges because local edits can propagate globally (Ying et al., 2019). Early graph-based
093 approaches (Ying et al., 2019; Cai et al., 2025; Lucic et al., 2022; Tan et al., 2022) primarily
094 used mask-based deletions to highlight influential substructures (Prado-Romero et al., 2024),
095 but these methods are closer to model explanations than true counterfactuals (Wachter
096 et al., 2017). Subsequent work has explored generative models (Ma et al., 2022), heuristic
097 edits (Abrate & Bonchi, 2021; Bajaj et al., 2021), and bi-level optimization (Giorgi et al.,
098 2025), followed more recently by motif-based (He et al., 2024), contrastive (Fournier &
099 Medya, 2025), and RL-driven strategies (Verma et al., 2024b). However, these methods
100 largely focus on single-graph perturbations to flip model predictions, and thus suffer from
101 strong model dependence, weak domain validity, and limited transferability. In contrast, we
102 are the first to formalize query-based counterfactual analysis on graphs, enabling retrieval
103 of counterfactuals with semantically meaningful concept differences from databases. Our
104 framework ensures domain validity by leveraging graphs in the input database, introduces
105 concept-based indices for scalability and fidelity, and provides certified bounds on query
106 correctness, bridging local explanations with database-grounded reasoning.

107 **Graph Database.** Graph queries (Zhao & Han, 2010) provide principled tools for re-
retrieval and analysis, for subgraph matching (Sun & Luo, 2020), reachability (Castanón
et al., 2015), shortest-path search (Sommer, 2014), and motif discovery (Lin et al., 2016),

with applications in many domains, such as social networks (Angles et al., 2013), recommendation (Ren et al., 2017), and knowledge graphs (Hogan et al., 2021). To scale these tasks, graph databases such as Neo4j (Robinson et al., 2015), Titan (DB-Engines, 2020), and Geabase (Fu et al., 2019) integrate native storage, indexing, and query languages. Graph learning (Kipf & Welling, 2017) further advanced this line of work by providing neural representations, enabling tighter integration with databases for downstream tasks (Walke et al., 2024; Zhou et al., 2023) through platforms such as GraphScope (Fan et al., 2021), AliGraph (Yang, 2019), DGL (Zheng et al., 2020), and neural graph databases (Besta et al., 2022). However, most graph learning databases prioritize scalability for model development or deployment. In contrast, our **Counterfactual Graph Database (CF-GDB)** shifts the focus from prediction accuracy to interpretability by querying semantic counterfactuals that flip outcomes. Instead of embedding graphs solely for prediction, CF-GDB grounds counterfactual queries in real database instances and reframes counterfactual analysis as a retrieval problem, bridging scalable database support with explainable reasoning.

3 PROBLEM FORMULATION

Consider a graph $G = (V_G, E_G, \mathbf{X}_G)$, where V_G is the set of nodes, $E_G \subseteq V_G \times V_G$ the edges, and $\mathbf{X}_G \in \mathbb{R}^{|V_G| \times d}$ the node-feature matrix, with each row $\mathbf{x}_v \in \mathbb{R}^d$ representing the feature vector of node $v \in V_G$. The parameter d denotes the feature dimensionality. Each graph is associated with a class label $y \in \mathcal{Y}$, where \mathcal{Y} is a finite label set. A **graph database** is defined as $\mathcal{D} = (G_i, y_i) \mid 1 \leq i \leq n, y_i \in \mathcal{Y}$, a collection of labeled graphs. A complete notation table is provided in Table 3 of Appendix A.

Problem. Given a graph database \mathcal{D} and a query $Q = (G, y)$, where G is a graph and $y \in \mathcal{Y}$ its label, the *counterfactual* of Q with respect to \mathcal{D} is defined as $(\tilde{G}^*, \tilde{y}^*) = \arg \min_{(\tilde{G}, \tilde{y}) \in \mathcal{D}, \tilde{y} \neq y} \Delta(G, \tilde{G})$, where $\Delta(G, \tilde{G}) \geq 0$ denotes a distance measure between graphs G and \tilde{G} . Note that the query graph Q does not have to belong to \mathcal{D} . A **counterfactual graph query (CGQ)** is the task of retrieving such a counterfactual from the database.

Defining an effective distance measure is critical for counterfactual graph queries. A natural choice is edit-based metrics, but these are computationally expensive, as they typically require both graph alignment (e.g., subgraph isomorphism or homomorphism, which are #P-hard (Dell et al., 2017)) and subsequent distance computation (e.g., graph edit distance, NP-hard (Gao et al., 2010)). More importantly, naive node-by-node alignment is semantically uninformative, as it fails to capture structural patterns shared across graphs (Armgaan et al., 2024). For example, in molecular graphs, replacing a carbon atom with nitrogen is not merely a single-atom substitution (Ma et al., 2022); it transforms a benzene ring into a pyridine ring, thereby altering key molecular properties such as toxicity or reactivity. The true unit of change is thus a higher-level structural element, such as a *functional group*, which naturally serves as a graph concept. Similarly, in social network classification, removing a single friendship edge rarely affects the predicted label. What matters instead is the presence or absence of larger structures, such as *communities* (Bajaj et al., 2021), which serve as graph concepts that distinguish between “fragmented” and “cohesive” networks.

These limitations highlight the need for a higher-level abstraction that enables semantically meaningful comparison and retrieval. We introduce **graph concepts** (or simply concepts), defined as probabilistic clusters of nodes and relations forming coherent semantic units (e.g., communities in social networks or functional groups in molecules). Graph concepts act as reusable *principal components* across graphs and serve as atomic units for counterfactual editing and retrieval. Identifying such concepts for counterfactual queries is nontrivial: they must be learned directly from the database and capture semantic differences between graphs that drive label changes. This requires both a principled method for concept learning and a mechanism to leverage them in queries. Unlike supervised approaches such as Concept Bottleneck Models (Koh et al., 2020; Barbiero et al., 2024), which rely on human-defined labels as auxiliary targets, our approach (C²GQ) automatically learns and extracts concepts for counterfactual queries without external supervision.

4 CONCEPT-BASED COUNTERFACTUAL GRAPH QUERY

To overcome the limitations of existing approaches—either computationally prohibitive and overly localized (edit-based) or semantically coarse and unable to capture fine-grained counterfactual differences (embedding-based)—we propose **Concept-Based Counterfactual Graph Query (C²GQ)**. Unlike prior methods tied to individual instances, C²GQ mines graph concepts from multiple graphs in a database, uncovering recurring structures (e.g., functional groups, communities, motifs) that serve as prototypes for capturing differences that drive prediction flips. By introducing a **hypergraph-based concept distance**, C²GQ enables multi-scale, domain-valid comparisons, yielding a more expressive metric for counterfactual queries. Building on this foundation, the **Concept Distribution Index (CDI)** and **Concept Semantic Index (CSI)** provide certified candidate-quality bounds by approximating $\Delta(G, \tilde{G})$ with lower and upper guarantees. These indices ensure that index-based search delivers model- and task-independent validity, certifies retrieval rankings, and safely prunes irrelevant candidates (Section 5). Taken as a whole, these guarantees establish C²GQ as a new paradigm for scalable and interpretable graph counterfactual queries in databases. Detailed pseudocode is provided in Appendix A.4.

4.1 CONCEPT EXTRACTION VIA STRUCTURAL EMBEDDINGS

A central challenge in counterfactual graph queries is to find a representation that captures meaningful structures while remaining scalable. Prior approaches often fall into two extremes: node-level representations, which are overly fine-grained and sensitive to small perturbations (Gao et al., 2010; Ying et al., 2019; Ma et al., 2022), and graph-level embeddings, which scale well but collapse local variations and overlook higher-order substructures (Bai et al., 2019; Wang et al., 2024; Li et al., 2025; Abrate & Bonchi, 2021). To bridge this gap, we introduce **graph concepts** as intermediate representations. Concepts serve as recurring prototypes (e.g., functional groups, communities) that act as stable yet sensitive anchors across instances, enabling large-scale database operations.

To effectively extract concepts, we employ Graph Neural Networks (GNNs) (Kipf & Welling, 2017; Veličković et al., 2018) to obtain node embeddings \mathbf{h}_v , and denote their collection by $\mathcal{H} = \{\mathbf{h}_v \mid v \in V_G, G \in \mathcal{D}\}$. Importantly, GNN-derived embeddings are structure-aware and semantically aligned across graphs (Armgaan et al., 2024), and thus capture higher-order patterns beyond raw features, crucial for reliable and transferable concept extraction (see Appendix A.1). We then cluster \mathcal{H} into K prototypes via K -means:

$$\min_{\{a(v)\}, \{\mathbf{c}_k\}_{k=1}^K} \sum_{\mathbf{h}_v \in \mathcal{H}} \|\mathbf{h}_v - \mathbf{c}_{a(v)}\|_2^2, \quad (1)$$

where $a : \mathcal{H} \rightarrow [K]$ assigns each embedding to a cluster index, and the corresponding centroid is $\mathbf{c}_{a(\mathbf{h}_v)} \in \mathbf{C}$. Each centroid \mathbf{c}_k serves as a *concept prototype*, representing a recurring structural pattern across the database. The collection of prototypes \mathbf{C} forms a global semantic dictionary that (i) amortizes computation by shifting clustering offline, and (ii) provides a shared basis that renders distance measures both efficient and semantically meaningful for counterfactual graph queries.

4.2 HYPERGRAPH-BASED CONCEPT DISTANCE

With concepts as the representation basis, the next challenge is to define a principled distance for comparing and mining multiple database graphs in concept space. Existing measures fall short: edit-based metrics (Gao et al., 2010; Chang et al., 2017) are computationally prohibitive, while embedding similarities (Bai et al., 2019; Zheng et al., 2014) obscure fine-grained counterfactual edits. In contrast, Optimal Transport (OT) (Petric Maretić et al., 2019) offers a more general framework, as it preserves structural correspondences, accommodates imbalance, and naturally models both discrete edits (via mass movement) and semantic shifts (via distances in embedding space). Yet node-level OT (Yu et al., 2025) remains insufficient for counterfactual queries, as it overlooks higher-order dependencies and cannot capture concept-level transformations such as ring substitutions in molecules or community rewiring in social networks. To address these limitations, we introduce the

hypergraph-based concept distance, which extends unbalanced OT (Vayer et al., 2020) to jointly model discrepancies at both node and concept levels across multiple graphs—a capability not achieved by prior approaches (Lucic et al., 2022; Abrate & Bonchi, 2021; Bajaj et al., 2021)—thereby providing a more expressive metric for counterfactual reasoning.

Each graph G is mapped to a hypergraph $H_G = (V_G, F_G)$, where each factor $f \in F_G$ groups nodes with the same concept assignment $a(v)$. This reflects the tendency of nearby nodes to share concepts, as well as the alignment of actors exhibiting similar structural patterns in the graph. Factor embeddings are then computed as $\mathbf{g}_f^G = \psi(\mathbf{h}_v^G : v \in f)$, where ψ is a permutation-invariant operator such as mean, sum, or attention (Xu et al., 2019; Lee et al., 2019). In this way, factor representations capture higher-order dependencies—such as communities in social networks or functional groups in molecules—that cannot be expressed solely through pairwise edges (see Appendix A.2 for details).

For two graphs G and \tilde{G} , we define a block-diagonal transport plan

$$T = \begin{bmatrix} T^{V_G V_{\tilde{G}}} & 0 \\ 0 & T^{F_G F_{\tilde{G}}} \end{bmatrix}, \quad T^{V_G V_{\tilde{G}}}, T^{F_G F_{\tilde{G}}} \geq 0, \quad (2)$$

where each row of $T^{V_G V_{\tilde{G}}}$ (resp. $T^{F_G F_{\tilde{G}}}$) corresponds to a node $v \in V_G$ and each column to a node $\tilde{v} \in V_{\tilde{G}}$, with $T_{v, \tilde{v}}^{V_G V_{\tilde{G}}}$ denoting the transported mass aligning v with \tilde{v} . Standard partial matching constraints (Chapel et al., 2020) enforce valid alignments, disentangling node- and concept-level distances, and yielding interpretable edits such as removals or substitutions.

By solving the transportation plan, the hypergraph-based distance is formally defined as

$$\Delta(G, \tilde{G}) = \min_{T^{V_G V_{\tilde{G}}}, T^{F_G F_{\tilde{G}}} \geq 0} \lambda_{\text{sub}}^V \sum_{v, v'} T_{vv'}^{V_G V_{\tilde{G}}} \delta_V(\mathbf{h}_v^G, \mathbf{h}_{v'}^{\tilde{G}}) + \lambda_{\text{del}}^V (|V_G| - \sum_{v, v'} T_{vv'}^{V_G V_{\tilde{G}}}) + \lambda_{\text{ins}}^V (|V_{\tilde{G}}| - \sum_{v, v'} T_{vv'}^{V_G V_{\tilde{G}}}) \\ + \lambda_{\text{sub}}^F \sum_{f, f'} T_{ff'}^{F_G F_{\tilde{G}}} \delta_F(\mathbf{g}_f^G, \mathbf{g}_{f'}^{\tilde{G}}) + \lambda_{\text{del}}^F (|F_G| - \sum_{f, f'} T_{ff'}^{F_G F_{\tilde{G}}}) + \lambda_{\text{ins}}^F (|F_{\tilde{G}}| - \sum_{f, f'} T_{ff'}^{F_G F_{\tilde{G}}}). \quad (3)$$

Here, λ_{sub}^* , λ_{del}^* , $\lambda_{\text{ins}}^* \geq 0$ denote the penalties for substitution, deletion, and insertion at both node and factor levels, while δ_V and δ_F are squared ℓ_2 costs (Bai et al., 2019; Dixit et al., 2022). These weights can be estimated from data by applying controlled perturbations and fitting regression models that map histogram shifts to unit costs.

Notably, the hypergraph-based distance provides flexible control over the granularity of concept comparison: setting $F = \emptyset$ recovers the classical graph edit distance (Gao et al., 2010); choosing $F = V$ with a single factor yields whole-graph embedding similarity (Bai et al., 2019); and defining local neighborhoods as factors enables structure-aware alignment through graph coarsening (Ying et al., 2018; Lee et al., 2019). By aligning concepts at appropriate granularities, the distance faithfully captures the minimal edits that drive label flips. Moreover, domain-specific choices of F enforce validity under real-world constraints, such as valence preservation in molecules (Gilmer et al., 2017), regulatory plausibility in finance (Rao et al., 2021), and spatio-temporal smoothness in traffic (Li et al., 2018).

4.3 GRAPH-LEVEL INDICES FROM CONCEPTS

Although the hypergraph-based concept distance is expressive, its direct application at scale is computationally prohibitive. To address this, we introduce two lightweight graph-level indices: (i) the **Concept Distribution Index (CDI)**, a histogram σ_G with $\sigma_G(k) = \sum_{v \in V_G} \mathbf{1}(a(v) = k)$, which provides a bag-of-concepts view and reflects discrete edits such as insertions or deletions; and (ii) the **Concept Semantic Index (CSI)**, a dense embedding $\mathbf{z}_G = \sum_{k=1}^K \frac{\sigma_G(k)}{|V_G|} \mathbf{c}_k$ that aggregates across multiple graphs to capture concept-level semantics. Theoretically, CDI tracks homomorphism counts to yield a certified lower bound, while CSI encodes these counts in continuous space to provide an upper bound. Together, they sandwich the hypergraph-based concept distance $\Delta(G, \tilde{G})$ (Theorems 1–6), offering the first certified guarantees for scalable and faithful counterfactual queries—unlike prior single-graph methods that rely on heuristic perturbations without guarantees (Ying et al., 2019; Abrate & Bonchi, 2021; Ma et al., 2022) (detailed Appendix A.3).

Finally, we introduce a three-step method for efficient counterfactual graph query answering:

-
1. **Label hashing.** Partition the database into buckets by label, and restrict the search to buckets with labels different from the query label y , i.e., $\mathcal{B}_{\neg y}$.
 2. **Dual pruning.** Within each bucket, first prune candidates using the ℓ_1 histogram distance $\delta_{\text{freq}}(G, \tilde{G}) = \|\sigma_G - \sigma_{\tilde{G}}\|_1$ (retaining the αM closest), and then refine with the ℓ_2 semantic distance $\delta_{\text{sem}}(G, \tilde{G}) = \|\mathbf{z}_G - \mathbf{z}_{\tilde{G}}\|_2$ to select the top- M candidates.
 3. **Fine re-ranking.** Apply the hypergraph-based distance to the pruned set $\mathcal{D}_{\text{pruned}}$, yielding interpretable edits $\Delta(G, \tilde{G})$ from the optimal transport plan T^* .

To further accelerate counterfactual queries, steps (1) and (2) are precomputed offline: label partitioning is stored in a hash map, and both indices are organized in a KD-tree. At query time, only step (3) is executed on a small portion of the database, with the prebuilt index enabling rapid retrieval of the relevant candidate subset.

Complexity. For a query graph G and a candidate \tilde{G} , the per-query time complexity is

$$O\left(K [\log |\mathcal{B}_{\neg y}| + \alpha M \log(\alpha M)] + d \alpha M + MI|V||\tilde{V}|\right). \quad (4)$$

The first term arises from CDI, which uses KD-tree search over K prototypes to select αM candidates from a bucket of size $|\mathcal{B}_{\neg y}|$. The second term corresponds to CSI, computing pairwise distances in d -dimensional embeddings with partial sorting. The final term accounts for the Sinkhorn step in OT, requiring I iterations over $|V|$ and $|\tilde{V}|$ nodes. Theorem 3 shows that enlarging the candidate set with both indices monotonically increases the probability of finding the optimal counterfactual. In experiments, we set $\alpha = 3$ and $M = 10$, retaining 30 candidates from CDI and 10 from CSI, which is sufficient to identify counterfactuals.

5 THEORETICAL ANALYSIS

Unlike edit-based distances that are computationally intractable or embedding-based distances that collapse local variations, our hypergraph-based concept distance explicitly captures counterfactual edits at both node and concept levels. This makes it particularly suitable for counterfactual reasoning. In this section, we show that this distance can be efficiently approximated by two complementary indices with certified guarantees. CDI captures bounded-depth homomorphism statistics and provides reliable lower bounds for safe pruning, while CSI reconstructs these features in a prototype space and yields computable upper bounds for similarity search. Together, they establish a sandwich bound on the true distance, yielding guarantees on query success rate, ranking correctness, and database-scale efficiency (see also Theorems 5–6 in Appendix B, which provide detailed homomorphism-based analysis of why C²GQ can find frequent subgraphs as concepts).

Lower bounds via CDI. We first show that differences in CDI yield certified lower bounds on the hypergraph-based concept distance.

Theorem 1. *Let $\lambda = \min\left\{\lambda_{\text{del}}^V, \lambda_{\text{ins}}^V, \frac{\lambda_{\text{sub}}^V}{2}\right\} + \alpha \cdot \min\left\{\lambda_{\text{del}}^F, \lambda_{\text{ins}}^F, \frac{\lambda_{\text{sub}}^F}{2}\right\}$. Then,*

$$\Delta(G, \tilde{G}) \geq \lambda \|\sigma_G - \sigma_{\tilde{G}}\|_1.$$

Here $\alpha = 1/(r\rho_{\max})$, where r is the maximum factor arity (i.e., nodes per factor) and ρ_{\max} is the maximum factor-load of any node (i.e., factors per node).

Theorem 1 formalizes the intuition that every edit necessarily modifies at least one WL-type count. Intuitively, α is a conservative scaling ensuring that factor-level edits are fully reflected by concept changes in the original graph G . Hence, the ℓ_1 gap between CDIs provides a certified lower bound on edit cost, ensuring that pruning by CDI never discards valid counterfactuals.

Upper bounds via CSI. We now introduce CSI as a complementary continuous perspective that provides upper bounds on the hypergraph-based distance.

Theorem 2. *Assume each factor feature is given by a permutation-invariant, L_ψ -Lipschitz aggregator $\mathbf{g}_f^G = \psi(\{\mathbf{h}_v^G : v \in f\})$, and let r be the maximum factor arity, ρ_{\max} the maximum*

factor-load of any node, and $\mathbf{C} \in \mathbb{R}^{d \times K}$ the prototype matrix with pseudoinverse \mathbf{C}^+ . Then,

$$\begin{aligned} \Delta(G, \tilde{G}) \leq & \left(\frac{\lambda_{\text{sub}}^V}{2} |V_G| + \lambda_{\text{sub}}^F \rho_{\text{max}} r L_\psi \right) \|\mathbf{C}^+\|_{1 \leftarrow 2} \|\mathbf{z}_G - \mathbf{z}_{\tilde{G}}\|_2 \\ & + \max\{\lambda_{\text{del}}^V, \lambda_{\text{ins}}^V\} \left| |V_G| - |V_{\tilde{G}}| \right| + \max\{\lambda_{\text{del}}^F, \lambda_{\text{ins}}^F\} \left| |F_G| - |F_{\tilde{G}}| \right|. \end{aligned}$$

$\|\mathbf{C}^+\|_{1 \leftarrow 2} := \sup_{x \neq 0} \|\mathbf{C}^+ x\|_1 / \|x\|_2$ converts the CSI’s ℓ_2 gap into an ℓ_1 mass movement.

Theorem 2 shows that CSI enables counterfactual analysis by quantifying smooth trajectories between a query graph and its perturbations, while also supporting concept mining by clustering structurally similar subgraphs. The CSI gap provides a computable upper bound: differences in CSI encode the amount of “semantic mass” to be shifted, certifying feasibility without solving the optimal transport problem exactly.

Remark. By combining the Theorems 1 and 2, we obtain a sandwich bound (Corollary 1 in Appendix B.4) that certifies counterfactual queries with three guarantees: (i) query success rates improve monotonically, since enlarging the CSI candidate set cannot reduce recall (Theorem 3 in Appendix B.5); (ii) pairwise rankings are provably correct, as CDI and CSI jointly certify margins between candidates (Theorem 4 in Appendix B.6); and (iii) pruning is safe, since no valid counterfactual lies outside the retained set (Corollary 2 in Appendix B.7).

6 EXPERIMENTS

We evaluate C²GQ on eight datasets spanning molecular, social, biological, and traffic domains. Experimental results demonstrate that C²GQ consistently outperforms strong baselines: it improves query accuracy by over 20%, achieves up to a 20× speed-up, and identifies counterfactuals that are both interpretable and domain-valid. Additional details on the experimental setup, dataset statistics, ablation studies, scalability analysis, and sensitivity tests are provided in Appendix C.

6.1 SETUP

Datasets. We evaluate on eight datasets spanning four domains. In the molecular domain, we use *NC11* (Wale et al., 2008), *Mutag* (Maron & Ames, 1983), and *AIDS* (Ivanov et al., 2019) for binary classification of chemical properties such as anticancer activity, mutagenicity, and HIV activity. In the social domain, we use *Reddit-Binary* (Yanardag & Vishwanathan, 2015), which consists of graphs of online discussion threads. In the biological domain, we evaluate on *PROTEINS* and *ENZYMES* (Borgwardt et al., 2005), covering enzyme classification and protein function prediction. In the traffic domain, we construct subgraph datasets from *METR-LA* and *PEMS-BAY* (Li et al., 2018) for predicting local congestion patterns. Since no benchmark provides ground-truth counterfactuals (Ying et al., 2019; Lucic et al., 2022; Zhang et al., 2023), we follow Giorgi et al. (2025) to synthesize them. Perturbations are designed as minimal yet valid edits to edges or node features that flip predictions while respecting domain constraints. We frame evaluation as a query task: given a graph G , the system must rank and retrieve its counterfactual \tilde{G}^* .

Evaluation Protocol. We compare C²GQ against several graph query baselines. *GED* applies graph edit distance with bipartite and anchor-aware bounds (Riesen & Bunke, 2009; Chang et al., 2017), producing accurate alignments but at high computational cost. Embedding-based methods include *GCN* (Kipf & Welling, 2017), *GIN* (Xu et al., 2019), *DiffPool* (Ying et al., 2018), *SAGPool* (Lee et al., 2019), *Graphormer* (Ying et al., 2021), and *SimGNN* (Bai et al., 2019). For scalability, we also evaluate a *+LSH* variant of each method, which accelerates search by 5–10× with only modest accuracy loss. We assess query quality using *Recall@k* and *Mean Reciprocal Rank (MRR)* (Wu et al., 2022), which capture coverage (whether a ground-truth counterfactual appears), position (average rank), and consistency (ranking quality). To evaluate the correctness of predicted concept-edit sets, we report their alignment with ground-truth transformations as *counterfactuals accuracy*, measured by *Precision (P)*, *Recall (R)*, and *F1-score (F1)* (Giorgi et al., 2025). Efficiency is measured by query latency T , defined as the average runtime (in seconds) per 100 queries.

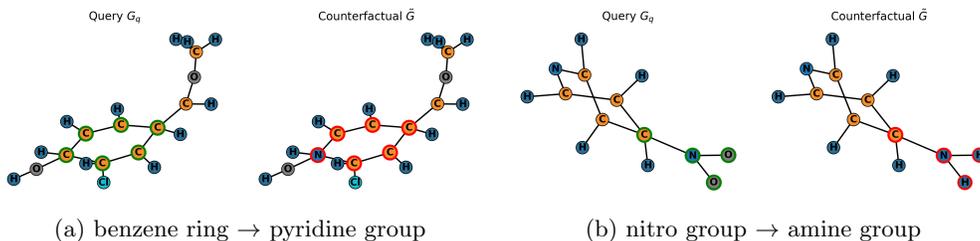


Figure 1: Examples of counterfactual graph queries on the Mutag dataset.

Method	Molecular									Social		
	NCI			Mutag			AIDS			Reddit-Binary		
	R@1	MRR	T	R@1	MRR	T	R@1	MRR	T	R@1	MRR	T
GED	0.460	0.579	410	0.449	0.566	22	0.473	0.590	36	0.424	0.561	1020
GCN	0.504	0.612	75	0.489	0.600	4	0.516	0.625	8	0.465	0.591	65
GIN	0.542	0.649	68	0.527	0.635	4	0.558	0.658	8	0.488	0.604	58
DiffPool	0.572	0.653	80	0.553	0.648	3	0.590	0.673	10	0.510	0.624	50
SAGPool	0.565	0.648	70	0.560	0.654	2	0.582	0.670	11	0.517	0.629	46
Graphormer	0.556	0.651	74	0.540	0.641	3	0.573	0.664	9	0.500	0.613	55
SimGNN	0.584	0.662	72	0.568	0.666	2	0.601	0.688	9	0.524	0.635	42
SimGNN+LSH	0.557	0.659	15	0.544	0.647	1	0.574	0.668	2	0.497	0.611	13
C ² GQ	0.704	0.790	29	0.683	0.773	3	0.724	0.808	5	0.628	0.728	42
C ² GQ (w/o index)	0.710	0.795	520	0.690	0.780	30	0.730	0.813	55	0.634	0.734	1323

Method	Biological						Traffic					
	PROTEINS			ENZYMES			METR-LA			PEMS-BAY		
	R@1	MRR	T	R@1	MRR	T	R@1	MRR	T	R@1	MRR	T
GED	0.440	0.576	138	0.409	0.547	95	0.405	0.540	470	0.413	0.555	722
GCN	0.475	0.600	24	0.445	0.573	16	0.427	0.561	60	0.435	0.569	85
GIN	0.503	0.616	21	0.470	0.590	13	0.446	0.581	54	0.454	0.588	78
DiffPool	0.520	0.627	22	0.480	0.592	14	0.452	0.584	47	0.468	0.597	70
SAGPool	0.512	0.622	20	0.488	0.598	11	0.460	0.589	45	0.462	0.594	67
Graphormer	0.511	0.620	23	0.478	0.589	12	0.452	0.583	55	0.460	0.592	74
SimGNN	0.529	0.632	19	0.493	0.605	12	0.463	0.596	44	0.471	0.602	64
SimGNN+LSH	0.500	0.610	6	0.470	0.585	4	0.438	0.573	15	0.445	0.580	20
C ² GQ	0.632	0.763	11	0.623	0.748	7	0.621	0.716	34	0.635	0.724	48
C ² GQ (w/o index)	0.648	0.779	200	0.640	0.765	140	0.640	0.733	652	0.672	0.758	927

Table 1: Query-level retrieval performance grouped by domain.

Implementation Details. Concept embeddings are obtained from a pretrained GCN (Lu et al., 2021) with hidden dimensionality 128, depth $L = 2$, and batch size 64, trained using Adam with a learning rate of 10^{-3} . Final-layer node embeddings $\{\mathbf{h}_v\}$ are clustered into $K = 64$ prototypes via K -means, yielding for each graph in \mathcal{D} its CDI σ_G and CSI \mathbf{z}_G . Penalty weights λ in the hypergraph-based distance are estimated from dataset statistics to balance edits across concepts. Retrieval is performed in three stages: (i) **Label hashing** to discard graphs with mismatched labels, (ii) **Dual indexing** on σ_G and \mathbf{z}_G for coarse pruning ($\alpha = 3$, $M = 10$), and (iii) **Fine re-ranking** of remaining candidates via the hypergraph-based concept distance with entropic Sinkhorn iterations ($\varepsilon = 0.01$, max 200).

6.2 EXPERIMENTAL RESULTS

Case Study. Unlike prior methods that reduce counterfactuals to isolated node- or edge-level perturbations, C²GQ identifies coherent, concept-level counterfactuals grounded in the database. As illustrated in Fig. 1, the highlighted regions (green in the query graph and red in the counterfactual) correspond to chemically meaningful substructures rather than arbitrary atom-level edits—for example, substituting a benzene ring with a pyridine ring or replacing a nitro group with an amine group, both of which induce plausible functional modifications (Reiser et al., 2022). These findings demonstrate that C²GQ uncovers interpretable, domain-valid counterfactual trajectories at the functional-group level, yielding explanations that are not only chemically meaningful but also informative of model behavior.

Validation Study. To evaluate the contribution of learned concepts, we conducted a control experiment by randomly permuting candidate graph concepts prior to retrieval. This procedure preserves the size of the concept space while destroying semantic alignment between prototypes and substructures. Relative to this randomized baseline, C²GQ achieved

Method	Molecular									Social		
	NCI			Mutag			AIDS			Reddit-Binary		
	R	P	F1	R	P	F1	R	P	F1	R	P	F1
GED	0.552	0.534	0.543	0.566	0.548	0.557	0.574	0.556	0.565	0.493	0.478	0.485
SimGNN	0.711	0.694	0.702	0.726	0.707	0.716	0.739	0.720	0.729	0.591	0.574	0.582
C ² GQ	0.805	0.839	0.822	0.807	0.781	0.794	0.835	0.806	0.820	0.671	0.656	0.663
C ² GQ (w/o index)	0.798	0.834	0.816	0.858	0.781	0.819	0.827	0.806	0.816	0.678	0.660	0.669
Method	Biological						Traffic					
	PROTEINS			ENZYMES			METR-LA			PEMS-BAY		
	R	P	F1	R	P	F1	R	P	F1	R	P	F1
GED	0.491	0.472	0.481	0.458	0.442	0.450	0.426	0.411	0.418	0.414	0.398	0.406
SimGNN	0.609	0.591	0.600	0.573	0.555	0.564	0.532	0.514	0.523	0.519	0.502	0.510
C ² GQ	0.697	0.672	0.684	0.655	0.637	0.646	0.619	0.601	0.610	0.598	0.589	0.593
C ² GQ (w/o index)	0.693	0.676	0.684	0.708	0.629	0.667	0.614	0.601	0.608	0.656	0.593	0.623

Table 2: Counterfactual set accuracy grouped by domain.

an average improvement of 27% in F1 across datasets. A paired t -test confirmed the statistical significance of this gain ($p \approx 0.029$), indicating that the improvements arise from structured concept modeling rather than chance. These results demonstrate that C²GQ uncovers domain-valid mechanisms, underscoring its value in high-stakes applications such as drug discovery and financial risk

Query Accuracy and Efficiency. Table 1 shows that C²GQ consistently outperforms the strongest baseline (SimGNN), with average gains of 12.4% in MRR and 13.5% in Recall@1. These improvements highlight the benefit of concept-level modeling in capturing counterfactual similarity, yielding robust generalization across molecular, social, biological, and traffic domains. The advantage is most pronounced on traffic datasets, where C²GQ achieves more than 20% higher accuracy. In contrast, GED computes edit distance directly on nodes and edges, requiring full OT alignment without concept abstraction, which makes it both slow and semantically limited. On Reddit-Binary, for example, GED takes over 1000 seconds yet still underperforms C²GQ in ranking. Embedding-based methods are faster than GED but, lacking concept-awareness, lag by more than 15% in accuracy. LSH-based variants further accelerate queries but incur substantial accuracy loss, underscoring the efficiency–accuracy trade-off. By contrast, our two concept-aware indices approximate the hypergraph-based distance with certified bounds (Theorems 3 and 4), achieving an average 20× speedup with less than a 1% accuracy drop. This demonstrates that indexing is not only practically effective but also theoretically grounded. Finally, scalability tests (Appendix C.3) show that C²GQ grows sublinearly with graph size $|G|$ and dataset size $|\mathcal{D}|$, whereas GED scales quadratically and linearly, respectively.

Counterfactuals Accuracy. Table 2 reports *counterfactuals accuracy*, which measures the alignment between predicted edits and ground-truth transformations. Overall, C²GQ achieves the highest F1, improving by an average of 44% over GED and 14% over SimGNN. These gains arise from jointly learning shared embeddings and constructing concept-prototypical sets, which provide a stable foundation for counterfactual reasoning. By leveraging the hypergraph-based concept distance, our method captures both fine-grained and coarse-grained structural variations. Moreover, the indices are theoretically guaranteed to preserve essential subgraph information via graph homomorphisms (Theorems 5 and 6), enabling scalability without sacrificing fidelity and interpretability.

7 CONCLUSIONS

We introduced the **Counterfactual Graph Database (CF-GDB)**, the first retrieval-based framework that grounds counterfactual reasoning in verifiable graph instances. At its core, **Concept-Based Counterfactual Graph Query (C²GQ)** abstracts graphs into semantic concepts and measures differences via a hypergraph-based distance, accelerated by two indices (CDI and CSI) with certified bounds. Our theoretical analysis established sandwich guarantees, while experiments across four domains demonstrated over 20% accuracy gains and up to 20× speedups against strong baselines. CF-GDB thus provides a principled foundation for scalable, interpretable, and domain-valid counterfactual analysis, paving the way for conditional queries, fairness auditing, and other high-stakes applications.

REPRODUCIBILITY STATEMENT

Implementation details such as concept extraction, hypergraph-based distance, indexing strategies (CDI/CSI), and retrieval pipeline pseudocode are presented in Appendix A. All theoretical contributions, including formal definitions of counterfactual queries, certified bounds, and proofs of correctness and efficiency (Theorems 1-6), are provided in detail in Appendix B to enable independent verification. For empirical evaluation, we describe datasets across four domains (molecular, social, biological, and traffic) and the associated tasks in Section 6, with additional dataset statistics, hyperparameters, and ablation studies reported in Appendix C. To further enhance transparency, we release our implementation as anonymous supplementary material at <https://anonymous.4open.science/r/CF-GDB-2DD6>, which includes scripts for preprocessing datasets, computing indices, and executing C²GQ queries. For theoretical claims, all assumptions are explicitly stated, and complete proofs are provided for guarantees on query success, ranking correctness, and safe pruning.

REFERENCES

- Carlo Abrate and Francesco Bonchi. Counterfactual graphs for explainable classification of brain networks. In *Proceedings of the 27th ACM SIGKDD conference on knowledge discovery & data mining*, pp. 2495–2504, 2021.
- Omid Amini, Fedor V Fomin, and Saket Saurabh. Counting subgraphs via homomorphisms. *SIAM Journal on Discrete Mathematics*, 26(2):695–717, 2012.
- Renzo Angles, Arnau Prat-Pérez, David Dominguez-Sal, and Josep-Lluís Larriba-Pey. Benchmarking database systems for social network applications. In *First International Workshop on Graph Data Management Experiences and Systems*, pp. 1–7, 2013.
- Burouj Armgaan, Manthan Dalmia, Sourav Medya, and Sayan Ranu. Graphtrail: Translating gnn predictions into human-interpretable logical rules. *Advances in Neural Information Processing Systems*, 37:123443–123470, 2024.
- Yunsheng Bai, Hao Ding, Song Bian, Ting Chen, Yizhou Sun, and Wei Wang. Simgnn: A neural network approach to fast graph similarity computation. In *Proceedings of the twelfth ACM international conference on web search and data mining*, pp. 384–392, 2019.
- Mohit Bajaj, Lingyang Chu, Zi Yu Xue, Jian Pei, Lanjun Wang, Peter Cho-Ho Lam, and Yong Zhang. Robust counterfactual explanations on graph neural networks. *Advances in neural information processing systems*, 34:5644–5655, 2021.
- Pietro Barbiero, Francesco Giannini, Gabriele Ciravegna, Michelangelo Diligenti, and Giuseppe Marra. Relational concept bottleneck models. *Advances in Neural Information Processing Systems*, 37:77663–77685, 2024.
- Maciej Besta, Patrick Iff, Florian Scheidl, Kazuki Osawa, Nikoli Dryden, Michal Podstawski, Tiancheng Chen, and Torsten Hoefer. Neural graph databases. In *Learning on Graphs Conference*, pp. 31–1. PMLR, 2022.
- Karsten M Borgwardt, Cheng Soon Ong, Stefan Schönauer, SVN Vishwanathan, Alexander J Smola, and Hans-Peter Kriegel. Protein function prediction via graph kernels. In *ISMB*, pp. 47–56. World Scientific, 2005.
- Ruichu Cai, Yuxuan Zhu, Xuexin Chen, Yuan Fang, Min Wu, Jie Qiao, and Zhifeng Hao. On the probability of necessity and sufficiency of explaining graph neural networks: A lower bound optimization approach. *Neural Networks*, 184:107065, 2025.
- Gregory Castanón, Yuting Chen, Ziming Zhang, and Venkatesh Saligrama. Efficient activity retrieval through semantic graph queries. In *Proceedings of the 23rd ACM international conference on Multimedia*, pp. 391–400, 2015.
- Lijun Chang, Xing Feng, Xuemin Lin, Lu Qin, and Wenjie Zhang. Efficient graph edit distance computation and verification via anchor-aware lower bound estimation. *arXiv preprint arXiv:1709.06810*, 2017.

540 Laetitia Chapel, Mokhtar Z Alaya, and Gilles Gasso. Partial gromov-wasserstein with ap-
541 plications on positive-unlabeled learning. *Advances in Neural Information Processing*
542 *Systems*, 2020.

543 Sajad Darabi, Piotr Bigaj, Dawid Majchrowski, Artur Kasymov, Pawel Morkisz, and Alex
544 Fit-Florea. A framework for large-scale synthetic graph dataset generation. *IEEE Trans-*
545 *actions on Neural Networks and Learning Systems*, 2025.

547 DB-Engines. Db-engines ranking of graph dbms. [https://db-engines.com/en/ranking/
548 graph+dbms](https://db-engines.com/en/ranking/graph+dbms), 2020.

549 Holger Dell, Martin Grohe, and Gaurav Rattan. Homomorphism counts and related param-
550 eters. *Journal of Computer and System Sciences*, 89:133–157, 2017.

552 Tanay Dixit, Bhargavi Paranjape, Hannaneh Hajishirzi, and Luke Zettlemoyer. Core: A
553 retrieve-then-edit framework for counterfactual data generation. In *Findings of the As-*
554 *sociation for Computational Linguistics: EMNLP 2022*, pp. 2964–2984, 2022.

555 Wenfei Fan, Tao He, Longbin Lai, Xue Li, Yong Li, Zhao Li, Zhengping Qian, Chao Tian,
556 Lei Wang, Jingbo Xu, et al. Graphscope: a unified engine for big graph processing.
557 *Proceedings of the VLDB Endowment*, 14(12):2879–2892, 2021.

559 Gregoire Fournier and Sourav Medya. Comrecgc: Global graph counterfactual explainer
560 through common recourse. *arXiv preprint arXiv:2505.07081*, 2025.

561 Zhisong Fu, Zhengwei Wu, Houyi Li, Yize Li, Min Wu, Xiaojie Chen, Xiaomeng Ye, Benquan
562 Yu, and Xi Hu. Geabase: A high-performance distributed graph database for industry-
563 scale applications. *International Journal of High Performance Computing and Network-*
564 *ing*, 15(1-2):12–21, 2019.

566 Jingwei Gan, Shinan Zhang, Chi Zhang, and Andy Li. Automated counterfactual generation
567 in financial model risk management. In *2021 IEEE International Conference on Big Data*
568 *(Big Data)*, pp. 4064–4068. IEEE, 2021.

569 Xinbo Gao, Bing Xiao, Dacheng Tao, and Xuelong Li. A survey of graph edit distance.
570 *Pattern Analysis and applications*, 13(1):113–129, 2010.

572 Justin Gilmer, Samuel S Schoenholz, Patrick F Riley, Oriol Vinyals, and George E Dahl.
573 Neural message passing for quantum chemistry. In *International conference on machine*
574 *learning*, pp. 1263–1272. Pmlr, 2017.

575 Flavio Giorgi, Fabrizio Silvestri, and Gabriele Tolomei. Combinex: A unified counterfactual
576 explainer for graph neural networks via node feature and structural perturbations. *arXiv*
577 *preprint arXiv:2502.10111*, 2025.

579 Yinhan He, Wendy Zheng, Yaochen Zhu, Jing Ma, Saumitra Mishra, Natraj Raman, Ning-
580 hao Liu, and Jundong Li. Global graph counterfactual explanation: A subgraph mapping
581 approach. *arXiv preprint arXiv:2410.19978*, 2024.

582 Aidan Hogan, Eva Blomqvist, Michael Cochez, Claudia d’Amato, Gerard De Melo, Claudio
583 Gutierrez, Sabrina Kirrane, José Emilio Labra Gayo, Roberto Navigli, Sebastian Neu-
584 maier, et al. Knowledge graphs. *ACM Computing Surveys (Csur)*, 54(4):1–37, 2021.

586 Sergei Ivanov, Sergei Sviridov, and Evgeny Burnaev. Understanding isomorphism bias in
587 graph data sets. *arXiv preprint arXiv:1910.12091*, 2019.

588 Amir-Hossein Karimi, Bernhard Schölkopf, and Isabel Valera. Algorithmic recourse: from
589 counterfactual explanations to interventions. In *Proceedings of the 2021 ACM conference*
590 *on fairness, accountability, and transparency*, pp. 353–362, 2021.

592 Thomas N. Kipf and Max Welling. Semi-supervised classification with graph convolutional
593 networks. In *International Conference on Learning Representations*, 2017. URL <https://openreview.net/forum?id=SJU4ayYgl>.

594 Pang Wei Koh, Thao Nguyen, Yew Siang Tang, Stephen Mussmann, Emma Pierson, Been
595 Kim, and Percy Liang. Concept bottleneck models. In *International conference on ma-*
596 *chine learning*, pp. 5338–5348. PMLR, 2020.

597 Junhyun Lee, Inyeop Lee, and Jaewoo Kang. Self-attention graph pooling. In *International*
598 *conference on machine learning*, pp. 3734–3743. pmlr, 2019.

600 Xunkai Li, Daohan Su, Sicheng Liu, Ru Zhang, Rong-Hua Li, and Guoren Wang. Two sides
601 of the same optimization coin: Model degradation and representation collapse in graph
602 foundation models. *arXiv preprint arXiv:2509.08401*, 2025.

604 Yaguang Li, Rose Yu, Cyrus Shahabi, and Yan Liu. Diffusion convolutional recurrent neural
605 network: Data-driven traffic forecasting. In *ICLR*, 2018.

606 Wenqing Lin, Xiaokui Xiao, Xing Xie, and Xiao-Li Li. Network motif discovery: A gpu
607 approach. *IEEE transactions on knowledge and data engineering*, 29(3):513–528, 2016.

609 Yuanfu Lu, Xunqiang Jiang, Yuan Fang, and Chuan Shi. Learning to pre-train graph neural
610 networks. In *Proceedings of the AAAI conference on artificial intelligence*, volume 35, pp.
611 4276–4284, 2021.

613 Ana Lucic, Maartje A Ter Hoeve, Gabriele Tolomei, Maarten De Rijke, and Fabrizio Sil-
614 vestri. Cf-gnnexplainer: Counterfactual explanations for graph neural networks. In *In-*
615 *ternational Conference on Artificial Intelligence and Statistics*, pp. 4499–4511. PMLR,
616 2022.

617 Jing Ma, Ruocheng Guo, Saumitra Mishra, Aidong Zhang, and Jundong Li. Clear: Gener-
618 ative counterfactual explanations on graphs. *Advances in neural information processing*
619 *systems*, 35:25895–25907, 2022.

621 Dorothy M Maron and Bruce N Ames. Revised methods for the salmonella mutagenicity
622 test. *Mutation Research/Environmental Mutagenesis and Related Subjects*, 113(3-4):173–
623 215, 1983.

624 Liam G McCoy, Connor TA Brenna, Stacy S Chen, Karina Vold, and Sunit Das. Believing in
625 black boxes: machine learning for healthcare does not need explainability to be evidence-
626 based. *Journal of clinical epidemiology*, 142:252–257, 2022.

628 Christopher Morris, Yaron Lipman, Haggai Maron, Bastian Rieck, Nils M Kriege, Martin
629 Grohe, Matthias Fey, and Karsten Borgwardt. Weisfeiler and leman go machine learning:
630 The story so far. *Journal of Machine Learning Research*, 24(333):1–59, 2023.

631 Ramaravind K Mothilal, Amit Sharma, and Chenhao Tan. Explaining machine learning
632 classifiers through diverse counterfactual explanations. In *Proceedings of the 2020 confer-*
633 *ence on fairness, accountability, and transparency*, pp. 607–617, 2020.

635 Martin Pawelczyk, Klaus Broelemann, and Gjergji Kasneci. Learning model-agnostic coun-
636 terfactual explanations for tabular data. In *Proceedings of the web conference 2020*, pp.
637 3126–3132, 2020.

638 Hermina Petric Maretic, Mireille El Gheche, Giovanni Chierchia, and Pascal Frossard. Got:
639 an optimal transport framework for graph comparison. *Advances in Neural Information*
640 *Processing Systems*, 32, 2019.

642 Mario Alfonso Prado-Romero, Bardh Prenkaj, Giovanni Stilo, and Fosca Giannotti. A sur-
643 vey on graph counterfactual explanations: definitions, methods, evaluation, and research
644 challenges. *ACM Computing Surveys*, 56(7):1–37, 2024.

645 Susie Xi Rao, Shuai Zhang, Zhichao Han, Zitao Zhang, Wei Min, Zhiyao Chen, Yinan Shan,
646 Yang Zhao, and Ce Zhang. xfraud: Explainable fraud transaction detection. *Proceedings*
647 *of the VLDB Endowment*, 15(3):427–436, 2021.

648 Patrick Reiser, Marlen Neubert, André Eberhard, Luca Torresi, Chen Zhou, Chen Shao,
649 Houssam Metni, Clint van Hoesel, Henrik Schopmans, Timo Sommer, et al. Graph neural
650 networks for materials science and chemistry. *Communications Materials*, 3(1):93, 2022.
651

652 Yongli Ren, Martin Tomko, Flora Dilys Salim, Jeffrey Chan, Charles LA Clarke, and Mark
653 Sanderson. A location-query-browse graph for contextual recommendation. *IEEE Trans-*
654 *actions on Knowledge and Data Engineering*, 30(2):204–218, 2017.

655 Kaspar Riesen and Horst Bunke. Approximate graph edit distance computation by means
656 of bipartite graph matching. *Image and Vision computing*, 27(7):950–959, 2009.
657

658 Ian Robinson, Jim Webber, and Emil Eifrem. *Graph databases: new opportunities for*
659 *connected data.* ” O’Reilly Media, Inc.”, 2015.

660 Donald B Rubin. Estimating causal effects of treatments in randomized and nonrandomized
661 studies. *Journal of educational Psychology*, 66(5):688, 1974.
662

663 Mengying Shu, Zhuxuanzi Wang, Jiayu Liang, et al. Early warning indicators for financial
664 market anomalies: A multi-signal integration approach. *Journal of Advanced Computing*
665 *Systems*, 4(9):68–84, 2024.
666

667 Christian Sommer. Shortest-path queries in static networks. *ACM Computing Surveys*
668 *(CSUR)*, 46(4):1–31, 2014.

669 Shixuan Sun and Qiong Luo. In-memory subgraph matching: An in-depth study. In *Pro-*
670 *ceedings of the 2020 ACM SIGMOD International Conference on Management of Data*,
671 pp. 1083–1098, 2020.
672

673 Juntao Tan, Shijie Geng, Zuohui Fu, Yingqiang Ge, Shuyuan Xu, Yunqi Li, and Yongfeng
674 Zhang. Learning and evaluating graph neural network explanations based on counterfac-
675 tual and factual reasoning. In *Proceedings of the ACM web conference 2022*, pp. 1018–
676 1027, 2022.

677 Titouan Vayer, Laetitia Chapel, Rémi Flamary, Romain Tavenard, and Nicolas Courty.
678 Fused gromov-wasserstein distance for structured objects. *Algorithms*, 13(9):212, 2020.
679

680 Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and
681 Yoshua Bengio. Graph attention networks. In *International Conference on Learning*
682 *Representations*, 2018.

683 Sahil Verma, Varich Boonsanong, Minh Hoang, Keegan Hines, John Dickerson, and Chirag
684 Shah. Counterfactual explanations and algorithmic recourses for machine learning: A
685 review. *ACM Computing Surveys*, 56(12):1–42, 2024a.
686

687 Samidha Verma, Burouj Armgaan, Sourav Medya, and Sayan Ranu. Induce: Inductive
688 counterfactual explanations for graph neural networks. *Transactions on Machine Learning*
689 *Research*, 2024b.

690 Sandra Wachter, Brent Mittelstadt, and Chris Russell. Counterfactual explanations without
691 opening the black box: Automated decisions and the gdpr. *Harv. JL & Tech.*, 31:841,
692 2017.
693

694 Nikil Wale, Ian A Watson, and George Karypis. Comparison of descriptor spaces for chemical
695 compound retrieval and classification. In *ICDM*, pp. 678–689. IEEE, 2008.

696 Daniel Walke, Daniel Micheel, Kay Schallert, Thilo Muth, David Broneske, Gunter Saake,
697 and Robert Heyer. The importance of graph databases and graph learning for clinical
698 applications. *Database*, 2024:baad045, 2024.
699

700 Jingdong Wang and Shipeng Li. Query-driven iterated neighborhood graph search for large
701 scale indexing. In *Proceedings of the 20th ACM international conference on Multimedia*,
pp. 179–188, 2012.

702 Jingjing Wang, Hongjie Zhu, Haoran Xie, Fu Lee Wang, Xiaoliang Xu, and Yuxiang Wang.
703 Graph similarity computation via interpretable neural node alignment. *arXiv preprint*
704 *arXiv:2412.12185*, 2024.

705 Shiwen Wu, Fei Sun, Wentao Zhang, Xu Xie, and Bin Cui. Graph neural networks in
706 recommender systems: a survey. *ACM Computing Surveys*, 55(5):1–37, 2022.

707 Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural
708 networks? In *International Conference on Learning Representations*, 2019.

709 Pinar Yanardag and SVN Vishwanathan. Deep graph kernels. In *KDD*, pp. 1365–1374.
710 ACM, 2015.

711 Hongxia Yang. Aligraph: A comprehensive graph neural network platform. In *Proceedings of*
712 *the 25th ACM SIGKDD international conference on knowledge discovery & data mining*,
713 pp. 3165–3166, 2019.

714 Linyi Yang, Jiazheng Li, Padraig Cunningham, Yue Zhang, Barry Smyth, and Ruihai Dong.
715 Exploring the efficacy of automatically generated counterfactuals for sentiment analysis.
716 In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics*
717 *and the 11th International Joint Conference on Natural Language Processing (Volume*
718 *1: Long Papers)*, pp. 306–316. Association for Computational Linguistics (ACL), 2021.

719 Chengxuan Ying, Tianle Cai, Shengjie Luo, Shuxin Zheng, Guolin Ke, Di He, Yanming
720 Shen, and Tie-Yan Liu. Do transformers really perform badly for graph representation?
721 *Advances in neural information processing systems*, 34:28877–28888, 2021.

722 Zhitao Ying, Jiaxuan You, Christopher Morris, Xiang Ren, Will Hamilton, and Jure
723 Leskovec. Hierarchical graph representation learning with differentiable pooling. *Ad-*
724 *vances in neural information processing systems*, 31, 2018.

725 Zhitao Ying, Dylan Bourgeois, Jiaxuan You, Marinka Zitnik, and Jure Leskovec. Gnnex-
726 plainer: Generating explanations for graph neural networks. *Advances in neural informa-*
727 *tion processing systems*, 32, 2019.

728 Zhitao Ying, Dylan Bourgeois, Jiaxuan You, Marinka Zitnik, and Jure Leskovec. Gnnex-
729 plainer: Generating explanations for graph neural networks. *Advances in neural informa-*
730 *tion processing systems*, 32, 2019.

731 Qi Yu, Zhichen Zeng, Yuchen Yan, Lei Ying, R Srikant, and Hanghang Tong. Joint optimal
732 transport and embedding for network alignment. In *Proceedings of the ACM on Web*
733 *Conference 2025*, pp. 2064–2075, 2025.

734 Kun Zhang, Chong Chen, Yuanzhuo Wang, Qi Tian, and Long Bai. Cfgl-lcr: A counterfac-
735 tual graph learning framework for legal case retrieval. In *Proceedings of the 29th ACM*
736 *SIGKDD Conference on knowledge discovery and data mining*, pp. 3332–3341, 2023.

737 Xiuping Zhang, Qun Liu, and Rui Han. Mmgcf: generating counterfactual explanations for
738 molecular property prediction via motif rebuild. *Journal of Computer and Communica-*
739 *tions*, 13(1):152–168, 2025.

740 Peixiang Zhao and Jiawei Han. On graph query optimization in large networks. *Proceedings*
741 *of the VLDB Endowment*, 3(1-2):340–351, 2010.

742 Da Zheng, Xiang Song, Chao Ma, Zeyuan Tan, Zihao Ye, Jin Dong, Hao Xiong, Zheng
743 Zhang, and George Karypis. Dgl-ke: Training knowledge graph embeddings at scale. In
744 *Proceedings of the 43rd international ACM SIGIR conference on research and development*
745 *in information retrieval*, pp. 739–748, 2020.

746 Weiguo Zheng, Lei Zou, Xiang Lian, Dong Wang, and Dongyan Zhao. Efficient graph
747 similarity search over large graph databases. *IEEE Transactions on Knowledge and Data*
748 *Engineering*, 27(4):964–978, 2014.

749 Xuanhe Zhou, Guoliang Li, Jianhua Feng, Luyang Liu, and Wei Guo. Grep: A graph
750 learning based database partitioning system. *Proceedings of the ACM on Management of*
751 *Data*, 1(1):1–24, 2023.

752

753

754

755

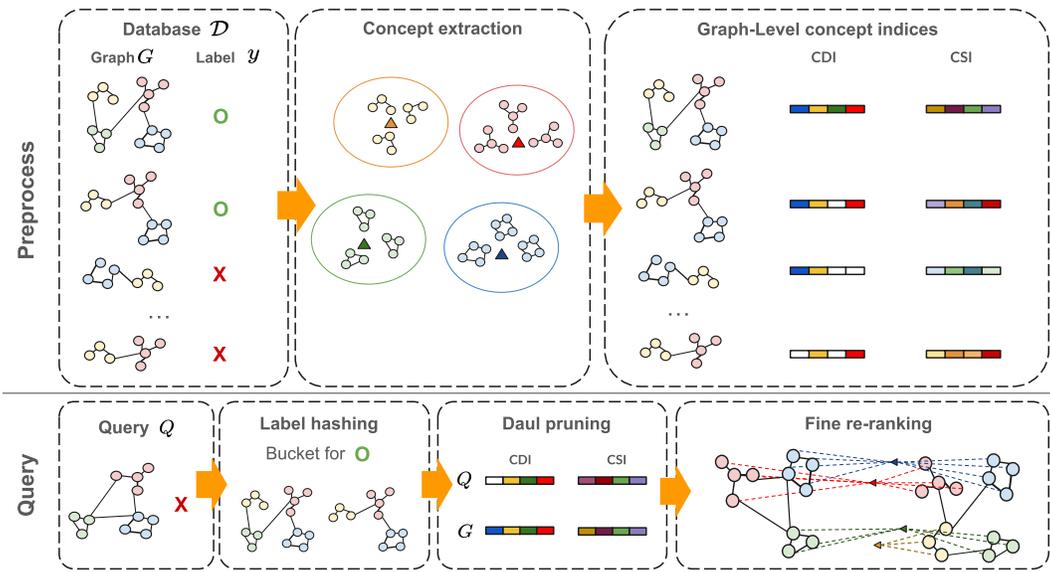


Figure 2: Overview of the C²GQ framework.

A DETAIL OF CONCEPT-BASED COUNTERFACTUAL GRAPH QUERY

We begin by formally defining the *concept assignment*, a mapping learned across multiple graphs in the database that projects nodes into a shared vocabulary of semantic prototypes. This induces concept groups globally, in contrast to prior methods that operate on single graphs (Ying et al., 2019; Verma et al., 2024a; Abrate & Bonchi, 2021; Lucic et al., 2022). Building on these assignments, we introduce a *hypergraph-based concept distance* that jointly aligns nodes and concept-level factors, capturing edits at multiple granularities while respecting domain-specific constraints (e.g., chemical valency or structural consistency). To ensure scalability, we further propose two lightweight *concept-based indices*—the Concept Distribution Index (CDI) and the Concept Semantic Index (CSI)—which approximate the hypergraph distance with certified lower and upper bounds. Together, these indices enable efficient yet faithful counterfactual search. The notations are summarized in Table 3.

As illustrated in Fig. 2, our C²GQ framework operates in two stages. The first is preprocessing, where we learn and extract graph concepts by applying K-means clustering ($K = 4$, shown as differently colored triangles) to group structurally similar subgraphs across multiple graphs in the database \mathcal{D} . Each cluster centroid, represented by a triangle, serves as a prototype concept. These concepts are then used to construct two indices per query, which are stored in a KD-tree for efficient retrieval.

In the query stage, we proceed in two steps. First, during label hashing, we filter the counterfactual bucket—i.e., graphs whose labels differ from that of the query graph. Second, we extract the indices of the query graph and compare them against candidates in the counterfactual bucket. Finally, we compute the hypergraph-based concept distance to rerank the candidates and identify the optimal counterfactual.

A.1 CONCEPT EXTRACTION VIA GNNs.

To effectively extract concepts, we employ Graph Neural Networks (GNNs) (Kipf & Welling, 2017; Veličković et al., 2018) to compute structure-aware node embeddings. At layer ℓ , each node $v \in V_G$ updates its representation through a message-passing scheme:

$$\mathbf{m}_v^{(\ell)} = \text{AGGREGATE}^{(\ell)}(\{\mathbf{h}_u^{(\ell-1)} : u \in N(v)\}), \quad \mathbf{h}_v^{(\ell)} = \text{COMBINE}^{(\ell)}(\mathbf{h}_v^{(\ell-1)}, \mathbf{m}_v^{(\ell)}), \quad (5)$$

Symbol	Meaning	Type / Shape
<i>Graph</i>		
$G = (V_G, E_G, X_G)$	graph with nodes, edges, node features	$X_G \in \mathbb{R}^{ V_G \times d}$
$y \in \mathcal{Y}$	class label of G	finite label set
$\mathcal{D} = \{(G, y)\}$	graph database (graphs with labels)	set
$Q = (G, y)$	query (graph, label)	pair
\tilde{G}	candidate graph (counterfactual)	graph
<i>Hypergraph</i>		
$H_G = (V_G, F_G)$	hypergraph over concepts (factors)	nodes V_G , factors F_G
$f \in F_G$	factor (group of nodes sharing $a(\cdot)$)	set of nodes
$e(f)$	incident node set of factor f	subset of V_G
$\mathbf{g}_f = \psi(\{h_v : v \in f\})$	factor embedding via aggregator ψ	\mathbb{R}^d
ψ	permutation-invariant aggregator (mean/sum/attn)	$\mathbb{R}^{d \times e(f) } \rightarrow \mathbb{R}^d$
L_ψ	Lipschitz constant of ψ	scalar
r	max factor arity ($ e(f) \leq r$)	integer
$\rho(v)$	factor-load of node v	integer
ρ_{\max}	$\max_v \rho(v)$	integer
<i>Concepts embedding</i>		
h_v	node embedding of v	\mathbb{R}^d
d	embedding dimension	integer
$H = \{h_v\}$	set of node embeddings over \mathcal{D}	multiset in \mathbb{R}^d
K	# concept prototypes (vocabulary size)	integer
$\{c_k\}_{k=1}^K$	concept prototypes (cluster centroids)	$c_k \in \mathbb{R}^d$
$C = [c_1, \dots, c_K]$	prototype matrix	$\mathbb{R}^{d \times K}$
C^+	pseudoinverse of C	$\mathbb{R}^{K \times d}$
$a(v) \in [K]$	concept assignment of node v	index
<i>Hypergraph-based concept distance</i>		
$\Delta(G, \tilde{G})$	hypergraph-based concept distance	scalar ≥ 0
$T = \begin{bmatrix} T_{V_G V_{\tilde{G}}} & 0 \\ 0 & T_{F_G F_{\tilde{G}}} \end{bmatrix}$	transport plan (nodes & factors)	nonnegative matrices
δ_V, δ_F	node/factor matching costs (e.g., ℓ_2^2)	$\mathbb{R}_{\geq 0}$
$\lambda_V^{\text{sub}}, \lambda_V^{\text{del}}, \lambda_V^{\text{ins}}$	node substitution/deletion/insertion penalties	$\mathbb{R}_{\geq 0}$
$\lambda_F^{\text{sub}}, \lambda_F^{\text{del}}, \lambda_F^{\text{ins}}$	factor substitution/deletion/insertion penalties	$\mathbb{R}_{\geq 0}$
I	# Sinkhorn iterations (OT)	integer
ε	entropic regularization (OT)	scalar
<i>Concept indices</i>		
$\sigma_G(k) = \sum_{v \in V_G} \mathbf{1}[a(v) = k]$	Concept Distribution Index (CDI)	\mathbb{N}^K
$z_G = \sum_{k=1}^K \frac{\sigma_G(k)}{ V_G } c_k$	Concept Semantic Index (CSI)	\mathbb{R}^d
$\delta_{\text{freq}}(G, \tilde{G}) = \ \sigma_G - \sigma_{\tilde{G}}\ _1$	histogram (CDI) distance	scalar
$\delta_{\text{sem}}(G, \tilde{G}) = \ z_G - z_{\tilde{G}}\ _2$	semantic (CSI) distance	scalar
$L(G, \tilde{G})$	CDI-based certified lower bound	scalar
$U(G, \tilde{G})$	CSI-based certified upper bound	scalar
$\ C^+\ _{1 \leftarrow 2} = \sup_{x \neq 0} \frac{\ C^+ x\ _1}{\ x\ _2}$	operator norm for CSI bound	scalar
$\alpha = \frac{1}{r \rho_{\max}}$	scaling constant in CDI lower bound	scalar
<i>Querying</i>		
$\{\mathcal{B}_y\}_{y \in \mathcal{Y}}$	label-specific buckets	partition of \mathcal{D}
$\mathcal{B}_{-y} = \bigcup_{y' \in \mathcal{Y}, y' \neq y} \mathcal{B}_{y'}$	union of buckets with labels different from y	subset of \mathcal{D}
$N_{-y} = \mathcal{B}_{-y} $	size of \mathcal{B}_{-y}	integer
M	# retained candidates after pruning	integer
α (prune)	pruning ratio in dual indexing	scalar

Table 3: Notation Table.

where AGGREGATE summarizes neighbor information and COMBINE integrates it with the node’s own state. After L layers, we obtain the final node embeddings $\mathbf{h}_v^{(L)}$ (or simply \mathbf{h}_v), and denote the collection across the database as

$$\mathcal{H} = \{\mathbf{h}_v \mid v \in V_G, G \in \mathcal{D}\}. \quad (6)$$

We then cluster \mathcal{H} into K prototypes via K -means, obtaining a global semantic dictionary $\mathbf{C} = [\mathbf{c}_1, \dots, \mathbf{c}_K]$. This dictionary amortizes computation by shifting clustering offline and provides a shared semantic basis for efficient and meaningful distance computation in counterfactual graph queries.

Definition 1 (Concept Assignment). *Given prototypes \mathbf{C} , a concept assignment for a graph $G = (V_G, E_G)$ is a mapping*

$$a : V_G \rightarrow \{1, \dots, K\}, \quad a(v) = \arg \min_{k \in [K]} \|\mathbf{h}_v - \mathbf{c}_k\|_2^2, \quad (7)$$

which assigns each node $v \in V_G$ to its closest prototype.

The assignment a induces a partition $\{V_G^k\}_{k=1}^K$, grouping nodes that share the same concept. These groups form the building blocks for higher-level structures such as concept histograms (CDI), semantic embeddings (CSI), and concept-level hyperedges in the hypergraph representation.

A.2 HYPERGRAPH-BASED CONCEPT DISTANCE

Having defined concept assignments, we now show how they induce a hypergraph representation of a graph, which in turn enables a distance measure that jointly aligns node-level and concept-level structures.

Each graph G (and \tilde{G}) is mapped to a hypergraph $H_G = (V_G, F_G)$, where each factor $f \in F_G$ groups nodes sharing the same concept assignment $a(v)$. Factor embeddings are computed as

$$\mathbf{g}_f^G = \psi(\{\mathbf{h}_v^G : v \in f\}),$$

where ψ is a permutation-invariant operator such as mean, sum, or attention (Xu et al., 2019; Lee et al., 2019). In this way, factor representations capture higher-order dependencies—such as communities in social networks or functional groups in molecules—that cannot be represented by pairwise edges.

We define a block-diagonal transport plan

$$T = \begin{bmatrix} T^{V_G V_{\tilde{G}}} & 0 \\ 0 & T^{F_G F_{\tilde{G}}} \end{bmatrix}, \quad T^{V_G V_{\tilde{G}}}, T^{F_G F_{\tilde{G}}} \geq 0,$$

where $T_{v, \tilde{v}}^{V_G V_{\tilde{G}}}$ denotes the transported mass from node $v \in V_G$ to $\tilde{v} \in V_{\tilde{G}}$, and analogously for factors. Standard partial matching constraints (Chapel et al., 2020) enforce valid alignments, disentangling node- and concept-level distances and yielding interpretable edits such as removals or substitutions.

Definition 2 (Hypergraph-based Concept Distance). *Given two graphs G and \tilde{G} , the hypergraph-based concept distance is defined as*

$$\begin{aligned} \Delta(G, \tilde{G}) = & \min_{T^{V_G V_{\tilde{G}}}, T^{F_G F_{\tilde{G}}} \geq 0} \lambda_{\text{sub}}^V \sum_{v, v'} T_{vv'}^{V_G V_{\tilde{G}}} \delta_V(\mathbf{h}_v^G, \mathbf{h}_{v'}^{\tilde{G}}) + \lambda_{\text{del}}^V (|V_G| - \sum_{v, v'} T_{vv'}^{V_G V_{\tilde{G}}}) + \lambda_{\text{ins}}^V (|V_{\tilde{G}}| - \sum_{v, v'} T_{vv'}^{V_G V_{\tilde{G}}}) \\ & + \lambda_{\text{sub}}^F \sum_{f, f'} T_{ff'}^{F_G F_{\tilde{G}}} \delta_F(\mathbf{g}_f^G, \mathbf{g}_{f'}^{\tilde{G}}) + \lambda_{\text{del}}^F (|F_G| - \sum_{f, f'} T_{ff'}^{F_G F_{\tilde{G}}}) + \lambda_{\text{ins}}^F (|F_{\tilde{G}}| - \sum_{f, f'} T_{ff'}^{F_G F_{\tilde{G}}}). \end{aligned} \quad (8)$$

subject to the partial matching constraints

$$\sum_{v'} T_{vv'}^{V_G V_{\tilde{G}}} \leq 1, \quad \sum_v T_{vv'}^{V_G V_{\tilde{G}}} \leq 1, \quad \sum_{f'} T_{ff'}^{F_G F_{\tilde{G}}} \leq 1, \quad \sum_f T_{ff'}^{F_G F_{\tilde{G}}} \leq 1.$$

Here, $\lambda_{\text{sub}}^V, \lambda_{\text{del}}^V, \lambda_{\text{ins}}^V, \lambda_{\text{sub}}^F, \lambda_{\text{del}}^F, \lambda_{\text{ins}}^F \geq 0$ are penalties for substitution, deletion, and insertion at the node and factor levels, while δ_V and δ_F denote squared ℓ_2 costs.

These penalty weights are estimated from data by: (i) sampling graphs from the database, (ii) applying controlled micro-edits (e.g., node deletion, insertion with sampled attributes, or feature substitution via nearest neighbors), (iii) recording the induced changes in concept histograms, and (iv) fitting a linear or nonnegative least-squares regression that maps histogram shifts to unit edit costs. The fitted coefficients then provide calibrated estimates of each λ .

918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971

A.3 CONCEPT-BASED INDICES

While the hypergraph-based distance provides a flexible notion of similarity, computing it exactly can be expensive. To enable scalable search, we introduce lightweight indices that approximate this distance while preserving key guarantees.

Definition 3 (Concept-based Indices). *Given a graph G with concept assignment $a : V_G \rightarrow [K]$, we define:*

- **Concept Distribution Index (CDI):** $\sigma_G(k) = \sum_{v \in V_G} \mathbf{1}\{a(v) = k\}$, $k \in [K]$, which counts the number of nodes assigned to each concept.
- **Concept Semantic Index (CSI):** $\mathbf{z}_G = \sum_{k=1}^K \frac{\sigma_G(k)}{|V_G|} \mathbf{c}_k$, where $\mathbf{c}_k \in \mathbb{R}^d$ is the embedding of concept k .

CDI captures discrete structural edits such as insertions, deletions, or substitutions, whereas CSI encodes the overall semantic distribution of concepts and varies smoothly under small perturbations (e.g., minor attribute changes or connectivity edits). Together, CDI and CSI provide complementary coarse representations that enable efficient yet faithful approximation of the hypergraph-based distance.

A.4 PSEUDOCODE

Algorithm 1 C²GQ Preprocessing (Offline)

- 1: **Input:** Graph database $\mathcal{D} = \{(G, y)\}$; embedding dim d ; prototypes K
 - 2: **Output:** Prototypes $C = [c_1, \dots, c_K]$; per-graph indices $\{(\sigma_G, z_G)\}$
 - 3: Train node encoder (GNN) on \mathcal{D} to obtain embeddings $\{h_v\}$
 - 4: Collect all node embeddings $H \leftarrow \{h_v\}$
 - 5: Run K-means on H to obtain prototypes $\{c_k\}$ and assignments $a(v)$
 - 6: **for** each graph $G \in \mathcal{D}$ **do**
 - 7: $\sigma_G(k) \leftarrow \sum_{v \in V_G} \mathbf{1}[a(v) = k]$ {CDI (histogram)}
 - 8: $z_G \leftarrow \sum_{k=1}^K \frac{\sigma_G(k)}{|V_G|} c_k$ {CSI (semantic embedding)}
 - 9: **end for**
 - 10: Return $C, \{(\sigma_G, z_G)\}_{G \in \mathcal{D}}$
-

Algorithm 2 C²GQ Query (Online)

- 1: **Input:** Query $Q = (G, y)$; prototypes C ; indices $\{(\sigma_{G'}, z_{G'})\}$; buckets $\{B_{y'}\}$; shortlist size M
 - 2: **Output:** Top- M counterfactual candidates
 - 3: $\mathcal{C} \leftarrow \bigcup_{y' \neq y} B_{y'}$ {Opposite-label search space}
 - 4: Compute (σ_G, z_G) for query G
 - 5: **for** each $G' \in \mathcal{C}$ **do**
 - 6: $\delta_{\text{freq}}(G, G') \leftarrow \|\sigma_G - \sigma_{G'}\|_1$
 - 7: **end for**
 - 8: Keep $\alpha \cdot M$ smallest δ_{freq} candidates \mathcal{C}_1
 - 9: **for** each $G' \in \mathcal{C}_1$ **do**
 - 10: $\delta_{\text{sem}}(G, G') \leftarrow \|z_G - z_{G'}\|_2$
 - 11: **end for**
 - 12: Keep M smallest δ_{sem} candidates \mathcal{C}_2
 - 13: Optionally re-rank \mathcal{C}_2 with final distance $\Delta(G, G')$
 - 14: Return \mathcal{C}_2 sorted by $\Delta(G, G')$
-

972 B DETAILED PROOF

973 B.1 GRAPH HOMOMORPHISM

974 We then by introducing bounded-depth rooted homomorphisms, a notion fundamental to
 975 subgraph (concept) mining (Amini et al., 2012).

976 **Definition 4** (Rooted Homomorphism). *Let (R, root) be a rooted tree and let G be a graph.
 977 A rooted homomorphism from (R, root) to (G, v) is a mapping $\varphi : V_R \rightarrow V_G$ such that (i)
 978 $\varphi(\text{root}) = v$; (ii) for every $(u, u') \in E_R$, we have $(\varphi(u), \varphi(u')) \in E_G$; and (iii) node labels
 979 (if present) are preserved by φ . The rooted homomorphism count is*

$$980 \text{hom}_o((R, \text{root}), G) := \sum_{v \in V_G} \text{hom}_o((R, \text{root}), (G, v)), \quad (9)$$

981 *i.e., the total number of root-preserving homomorphisms of R into G .*

982 Let $\mathcal{R}_{\leq L}$ denote the set of rooted L -types (e.g., rooted L -neighborhood computation trees
 983 under 1-WL refinement). For a graph G , the Concept Distribution Index (CDI) is $\sigma_G \in \mathbb{N}^K$
 984 with $K := |\mathcal{R}_{\leq L}|$, where the k -th entry counts the number of nodes of type c_k .¹

985 Building on this perspective, we establish a linear correspondence between CDI and rooted
 986 homomorphism counts (Theorem 5). Having demonstrated the effectiveness of the discrete
 987 CDI representation, we next introduce the Concept Semantic Index (CSI) as a complemen-
 988 tary continuous view. In particular, we show that CSI recovers bounded-depth homomor-
 989 phism features, ensuring that the continuous representation retains the expressive power of
 990 the discrete statistics (Theorem 6).

991 B.2 PROOF OF THEOREM 1

992 **Theorem 1.** *Let $\lambda = \min\left\{\lambda_{\text{del}}^V, \lambda_{\text{ins}}^V, \frac{\lambda_{\text{sub}}^V}{2}\right\} + \alpha \cdot \min\left\{\lambda_{\text{del}}^F, \lambda_{\text{ins}}^F, \frac{\lambda_{\text{sub}}^F}{2}\right\}$. Then,*

$$993 \Delta(G, \tilde{G}) \geq \lambda \|\sigma_G - \sigma_{\tilde{G}}\|_1.$$

994 Here $\alpha = 1/(r\rho_{\max})$, where r is the maximum factor arity (i.e., nodes per factor) and ρ_{\max}
 995 is the maximum factor-load of any node (i.e., factors per node).

996 **Assumption 1** (Factor construction and loads). *Each factor $f \in F$ is determined, up to
 997 its “type”, by a permutation-invariant rule on the multiset of incident node concepts: there
 998 exists Ψ such that $\text{type}(f) = \Psi(\{a(u) : u \in f\})$. The factor feature is $\mathbf{g}_f = \psi(\{h_u : u \in f\})$,
 999 where ψ is permutation-invariant and L_ψ -Lipschitz. Let the factor arity be bounded by
 1000 $|f| \leq r$ and the node factor-load $\rho(v) := |\{f \in F : v \in f\}|$ be bounded with $\rho_{\max} :=$
 1001 $\max_v \rho(v) < \infty$.*

1002 **Assumption 2** (Node-local feasibility). *Feasible counterfactual edits are node-local: any
 1003 change to a factor must be induced by edits to its incident nodes (or their weights). Equiv-
 1004 alently, the feasible region forbids “pure factor-only” edits.*

1005 **Assumption 3** (Multiset sensitivity and nondegeneracy). *The map Ψ is strictly sensitive
 1006 to changes in the incident multiset of node concepts: if $\{a(u) : u \in f\}$ changes, then $\text{type}(f)$
 1007 changes. Let $\rho_{\min} := \min_v \rho(v)$, and assume the nondegeneracy condition*

$$1008 r \rho_{\max} \rho_{\min} \geq 2. \quad (10)$$

1009 For a graph G , write its CDI as $\sigma_G \in \mathbb{N}^K$ and define

$$1010 D := \|\sigma_G - \sigma_{\tilde{G}}\|_1. \quad (11)$$

1011 **Lemma 1** (Node edits versus CDI ℓ_1 gap). *Let N be the minimum number of node-level
 1012 edits (deletions, insertions, substitutions) to transform G into \tilde{G} . Then*

$$1013 N \geq \frac{D}{2}. \quad (12)$$

1014 ¹Here K is the dimensionality of the representation. In practice one may also set the number of
 1015 clusters in K -means to K , reflecting a partition of nodes into K concept categories.

1026 *Proof.* A deletion (resp. insertion) changes one concept count by -1 (resp. $+1$), increasing
 1027 D by 1. A substitution decreases one count by 1 and increases another by 1, increasing D
 1028 by 2. Hence each node edit increases D by at least 1 (substitutions by 2), so $2N \geq D$. \square
 1029

1030 **Lemma 2** (Factor-to-node charging). *Under Assumptions 1–3, the total number of factor*
 1031 *changes between G and \tilde{G} satisfies*

$$1032 \quad \#(\text{factor changes}) \geq \frac{\rho_{\min}}{2} D. \quad (13)$$

1033
 1034 *Proof.* By Assumption 2, every factor change must be induced by edits on its incident nodes.
 1035 By Assumption 3, whenever a node’s concept changes, *each* incident factor must change its
 1036 type. Thus one node-concept change induces at least ρ_{\min} factor changes. Let N be the
 1037 minimal number of node edits; then factor changes $\geq \rho_{\min} N$. By Lemma 1, $N \geq D/2$, so
 1038 the claim follows. \square
 1039

1040 *Full proof of the lower bound.* Let

$$1041 \quad \kappa_V := \min\left\{\lambda_{\text{del}}^V, \lambda_{\text{ins}}^V, \frac{\lambda_{\text{sub}}^V}{2}\right\}, \quad \kappa_F := \min\left\{\lambda_{\text{del}}^F, \lambda_{\text{ins}}^F, \frac{\lambda_{\text{sub}}^F}{2}\right\}. \quad (14)$$

1042
 1043 **Node level.** Each unit increase of the CDI ℓ_1 gap incurs at least κ_V cost from node edits,
 1044 hence

$$1045 \quad \text{node-cost} \geq \kappa_V D. \quad (15)$$

1046
 1047 **Factor level.** Each factor change costs at least κ_F . By Lemma 2, factor changes \geq
 1048 $(\rho_{\min}/2) D$, hence

$$1049 \quad \text{factor-cost} \geq \kappa_F \cdot \frac{\rho_{\min}}{2} D. \quad (16)$$

1050 Summing (15)–(16) yields

$$1051 \quad \Delta(G, \tilde{G}) \geq \left(\kappa_V + \frac{\rho_{\min}}{2} \kappa_F\right) D. \quad (17)$$

1052
 1053 Now define $\alpha := 1/(r \rho_{\max})$. By (10), we have $\alpha \leq \rho_{\min}/2$, hence

$$1054 \quad \kappa_V + \alpha \kappa_F \leq \kappa_V + \frac{\rho_{\min}}{2} \kappa_F. \quad (18)$$

1055 Combining with (17) gives

$$1056 \quad \Delta(G, \tilde{G}) \geq \left(\kappa_V + \alpha \kappa_F\right) D = \left(\min\{\lambda_{\text{del}}^V, \lambda_{\text{ins}}^V, \frac{\lambda_{\text{sub}}^V}{2}\} + \alpha \min\{\lambda_{\text{del}}^F, \lambda_{\text{ins}}^F, \frac{\lambda_{\text{sub}}^F}{2}\}\right) \|\sigma_G - \sigma_{\tilde{G}}\|_1. \quad (19)$$

1057
 1058 The scaling $\alpha = 1/(r \rho_{\max})$ conservatively converts factor-level edits into CDI ℓ_1 mass while
 1059 avoiding over-counting across incident factors. Since at least one penalty is strictly positive,
 1060 the coefficient is $\lambda > 0$, the theorem follows. \square
 1061

1062 B.3 PROOF OF THEOREM 2

1063 **Theorem 2.** *Assume each factor feature is given by a permutation-invariant, L_ψ -Lipschitz*
 1064 *aggregator $\mathbf{g}_f^G = \psi(\{\mathbf{h}_v^G : v \in f\})$, and let r be the maximum factor arity, ρ_{\max} the maximum*
 1065 *factor-load of any node, and $\mathbf{C} \in \mathbb{R}^{d \times K}$ the prototype matrix with pseudoinverse \mathbf{C}^+ . Then,*

$$1066 \quad \Delta(G, \tilde{G}) \leq \left(\frac{\lambda_{\text{sub}}^V}{2} |V_G| + \lambda_{\text{sub}}^F \rho_{\max} r L_\psi\right) \|\mathbf{C}^+\|_{1 \leftarrow 2} \|\mathbf{z}_G - \mathbf{z}_{\tilde{G}}\|_2$$

$$1067 \quad + \max\{\lambda_{\text{del}}^V, \lambda_{\text{ins}}^V\} \left||V_G| - |V_{\tilde{G}}|\right| + \max\{\lambda_{\text{del}}^F, \lambda_{\text{ins}}^F\} \left||F_G| - |F_{\tilde{G}}|\right|.$$

1068 $\|\mathbf{C}^+\|_{1 \leftarrow 2} := \sup_{x \neq 0} \|\mathbf{C}^+ x\|_1 / \|x\|_2$ converts the CSI’s ℓ_2 gap into an ℓ_1 mass movement.

1069 *Proof.* Decompose the objective into node-/factor-substitution and size-gap parts:

$$1070 \quad \Delta(G, \tilde{G}) = \frac{\lambda_{\text{sub}}^V}{2} \sum_{v, v'} T_{vv'}^{V_G V_{\tilde{G}}} \|\mathbf{h}_v^G - \mathbf{h}_{v'}^{\tilde{G}}\|_2^2 + \frac{\lambda_{\text{sub}}^F}{2} \sum_{f, f'} T_{ff'}^{F_G F_{\tilde{G}}} \|\mathbf{g}_f^G - \mathbf{g}_{f'}^{\tilde{G}}\|_2^2 + (\text{del/ins terms}).$$

1080 **CSI gap \Rightarrow minimal concept mass movement.** Let $\mathbf{C} = [c_1, \dots, c_K] \in \mathbb{R}^{d \times K}$ be the
 1081 prototype matrix with pseudoinverse \mathbf{C}^+ , and $\mathbf{z}_G = \mathbf{C}p(G)$ for the normalized concept
 1082 histogram $p(G) \in \Delta^{K-1}$. Define

$$1083 m^* := \min_{y: \mathbf{C}y = \mathbf{z}_G - \mathbf{z}_{\tilde{G}}, \mathbf{1}^\top y = 0} \|y\|_1 \leq \|\mathbf{C}^+\|_{1 \leftarrow 2} \|\mathbf{z}_G - \mathbf{z}_{\tilde{G}}\|_2. \quad (1)$$

1084 Let $y^* = s - t$ with $s, t \geq 0$, $\mathbf{1}^\top s = \mathbf{1}^\top t = m^*$, and let $\Gamma = (\Gamma_{ij})$ be a concept-level coupling
 1085 with row/column sums s, t (so $\sum_{i,j} \Gamma_{ij} = m^*$).

1086 **Node substitution bound.** Couple node features to prototypes inside each concept pair
 1087 (*prototype coupling*):

$$1088 \sum_{v,v'} T_{vv'}^{V_G V_{\tilde{G}}} \|\mathbf{h}_v^G - \mathbf{h}_{v'}^{\tilde{G}}\|_2^2 \leq \sum_{i,j} \Gamma_{ij} \|\mathbf{c}_i - \mathbf{c}_j\|_2^2 \leq C_{\max} m^*,$$

1089 where $C_{\max} := \max_{i,j} \|\mathbf{c}_i - \mathbf{c}_j\|_2^2$. Using (1),

$$1090 \sum_{v,v'} T_{vv'}^{V_G V_{\tilde{G}}} \|\mathbf{h}_v^G - \mathbf{h}_{v'}^{\tilde{G}}\|_2^2 \leq C_{\max} \|\mathbf{C}^+\|_{1 \leftarrow 2} \|\mathbf{z}_G - \mathbf{z}_{\tilde{G}}\|_2. \quad (20)$$

1091 (If desired, absorb C_{\max} into constants or keep it explicit; below we absorb it.)

1092 **Factor substitution via incidence lift.** Assume $\mathbf{g}_f^G = \psi(\{\mathbf{h}_v^G : v \in f\})$ with ψ
 1093 permutation-invariant and L_ψ -Lipschitz; let $|f| \leq r$. Let ρ_{\max} be the maximum factor-
 1094 load of any node (factors per node). A node edit can affect at most ρ_{\max} factors, and each
 1095 factor aggregates at most r node edits; Lipschitz stability gives the lift

$$1096 \sum_{f,f'} T_{ff'}^{F_G F_{\tilde{G}}} \|\mathbf{g}_f^G - \mathbf{g}_{f'}^{\tilde{G}}\|_2^2 \leq 2 \rho_{\max} r L_\psi \sum_{v,v'} T_{vv'}^{V_G V_{\tilde{G}}} \|\mathbf{h}_v^G - \mathbf{h}_{v'}^{\tilde{G}}\|_2^2, \quad (21)$$

1097 where the factor 2 accounts for two-sided changes from insertions/deletions (if one uses a
 1098 one-sided lift, drop the 2 and keep a $\frac{\lambda_{\text{sub}}^F}{2}$ in the final coefficient).

1099 Combining (20)–(21) and absorbing C_{\max} into constants yields

$$1100 \sum_{f,f'} T_{ff'}^{F_G F_{\tilde{G}}} \|\mathbf{g}_f^G - \mathbf{g}_{f'}^{\tilde{G}}\|_2^2 \leq 2 \rho_{\max} r L_\psi \|\mathbf{C}^+\|_{1 \leftarrow 2} \|\mathbf{z}_G - \mathbf{z}_{\tilde{G}}\|_2. \quad (22)$$

1101 **Size gaps.** Unmatched nodes (resp. factors) cost at least $\max\{\lambda_{\text{del}}^V, \lambda_{\text{ins}}^V\}$ (resp.
 1102 $\max\{\lambda_{\text{del}}^F, \lambda_{\text{ins}}^F\}$) per unit, giving the two size-gap terms in the statement.

1103 **Aggregation.** Multiply (20) by $\frac{\lambda_{\text{sub}}^V}{2}$ and (22) by $\frac{\lambda_{\text{sub}}^F}{2}$; the factor 2 in (22) cancels the
 1104 $\frac{1}{2}$, producing the coefficient λ_{sub}^F as stated. Finally, (harmlessly) relax the node term by a
 1105 factor $|V_G| \geq 1$ to write the dependence as $\frac{\lambda_{\text{sub}}^V}{2} |V_G|$. This gives

$$1106 \Delta(G, \tilde{G}) \leq \left(\frac{\lambda_{\text{sub}}^V}{2} |V_G| + \lambda_{\text{sub}}^F \rho_{\max} r L_\psi \right) \|\mathbf{C}^+\|_{1 \leftarrow 2} \|\mathbf{z}_G - \mathbf{z}_{\tilde{G}}\|_2$$

$$1107 + \max\{\lambda_{\text{del}}^V, \lambda_{\text{ins}}^V\} \left| |V_G| - |V_{\tilde{G}}| \right| + \max\{\lambda_{\text{del}}^F, \lambda_{\text{ins}}^F\} \left| |F_G| - |F_{\tilde{G}}| \right|,$$

1108 which is the desired bound. The norm $\|\mathbf{C}^+\|_{1 \leftarrow 2} := \sup_{x \neq 0} \|\mathbf{C}^+ x\|_1 / \|x\|_2$ converts the CSI
 1109 ℓ_2 gap into an ℓ_1 concept-mass movement, as used in (1). The theorem follows. \square

1110 B.4 PROOF OF COROLLARY 1

1111 **Corollary 1.** For any candidate \tilde{G} , $\mathcal{C}^2 GQ$ provides a sandwich bound $L(G, \tilde{G}) \leq \Delta(G, \tilde{G}) \leq$
 1112 $U(G, \tilde{G})$, where L is the CDI-based lower bound and U is the CSI-based upper bound.

1113 *Proof.* By combining Theorems 1 and 2, we obtain the sandwich bound. \square

B.5 PROOF OF THEOREM 3

Theorem 3. *Let $B_{\neg y}$ be the opposite-label (counterfactual) bucket for G . For a distance threshold $\tau > 0$, define the CSI-feasible set $\mathcal{F}_\tau(G) := \{\tilde{G} \in B_{\neg y} : U(G, \tilde{G}) \leq \tau\}$. If candidates are drawn from a sampling distribution π on $B_{\neg y}$, then*

$$\Pr_{\tilde{G} \sim \pi} [\Delta(G, \tilde{G}) \leq \tau] \geq \sum_{\tilde{G} \in \mathcal{F}_\tau(G)} \pi(\tilde{G}).$$

Moreover, if M candidates are drawn i.i.d. from π , the success probability is lower bounded by $1 - (1 - \sum_{\tilde{G} \in \mathcal{F}_\tau(G)} \pi(\tilde{G}))^M$ (e.g., $1 - (1 - |\mathcal{F}_\tau(G)|/|B_{\neg y}|)^M$ for uniform distribution).

Proof. Let $\mathcal{F}_\tau(G) = \{\tilde{G} \in B_{\neg y} : U(G, \tilde{G}) \leq \tau\}$. By the sandwich inequality $L \leq \Delta \leq U$, we have $U(G, \tilde{G}) \leq \tau \Rightarrow \Delta(G, \tilde{G}) \leq \tau$, hence

$$\mathcal{F}_\tau(G) \subseteq \{\tilde{G} \in B_{\neg y} : \Delta(G, \tilde{G}) \leq \tau\}. \quad (23)$$

Taking probability under π gives

$$\Pr_{\tilde{G} \sim \pi} [\Delta(G, \tilde{G}) \leq \tau] \geq \Pr_{\tilde{G} \sim \pi} [\tilde{G} \in \mathcal{F}_\tau(G)] = \sum_{\tilde{G} \in \mathcal{F}_\tau(G)} \pi(\tilde{G}), \quad (24)$$

which proves the one-shot bound.

For M i.i.d. draws, let

$$p := \Pr_{\tilde{G} \sim \pi} [\tilde{G} \in \mathcal{F}_\tau(G)] = \sum_{\tilde{G} \in \mathcal{F}_\tau(G)} \pi(\tilde{G}). \quad (25)$$

The success event is that at least one of the M samples falls in $\mathcal{F}_\tau(G)$, hence

$$\Pr[\exists \tilde{G} \in \mathcal{F}_\tau(G)] = 1 - (1 - p)^M \geq 1 - e^{-Mp}. \quad (26)$$

Under uniform sampling on $B_{\neg y}$, $p = |\mathcal{F}_\tau(G)|/|B_{\neg y}|$, yielding

$$\Pr[\exists \tilde{G} \in \mathcal{F}_\tau(G)] \geq 1 - \left(1 - \frac{|\mathcal{F}_\tau(G)|}{|B_{\neg y}|}\right)^M. \quad (27)$$

(If sampling is without replacement, the exact success probability is $1 - \binom{|B_{\neg y}| - |\mathcal{F}_\tau(G)|}{M} / \binom{|B_{\neg y}|}{M}$, which is $\geq 1 - (1 - p)^M$.) The theorem follows. \square

B.6 PROOF OF THEOREM 4

Theorem 4. *If $\tilde{G}_b \prec^* \tilde{G}_a$, then $\Delta(G, \tilde{G}_b) < \Delta(G, \tilde{G}_a)$. Thus every certified comparison is strictly correct and induces an acyclic partial order.*

Proof. By definition, $\tilde{G}_b \prec^* \tilde{G}_a$ means $L(G, \tilde{G}_a) > U(G, \tilde{G}_b)$. Using the sandwich bound $L \leq \Delta \leq U$ gives

$$\Delta(G, \tilde{G}_a) \geq L(G, \tilde{G}_a) > U(G, \tilde{G}_b) \geq \Delta(G, \tilde{G}_b), \quad (28)$$

hence $\Delta(G, \tilde{G}_b) < \Delta(G, \tilde{G}_a)$. In particular, \prec^* is irreflexive (since $L(G, \tilde{G}) \leq U(G, \tilde{G})$) and asymmetric.

For acyclicity, suppose there is a directed cycle $\tilde{G}_1 \prec^* \tilde{G}_2 \prec^* \dots \prec^* \tilde{G}_k \prec^* \tilde{G}_1$. Applying the strict inequality along the cycle yields

$$\Delta(G, \tilde{G}_1) < \Delta(G, \tilde{G}_2) < \dots < \Delta(G, \tilde{G}_k) < \Delta(G, \tilde{G}_1), \quad (29)$$

a contradiction. Thus the digraph induced by \prec^* is acyclic. Consequently, the transitive closure of \prec^* is a strict partial order (it is irreflexive, asymmetric, and transitive), and every certified comparison is strictly correct. The theorem follows. \square

1188
1189
1190
1191
1192
1193
1194
1195
1196
1197
1198
1199
1200
1201
1202
1203
1204
1205
1206
1207
1208
1209
1210
1211
1212
1213
1214
1215
1216
1217
1218
1219
1220
1221
1222
1223
1224
1225
1226
1227
1228
1229
1230
1231
1232
1233
1234
1235
1236
1237
1238
1239
1240
1241

B.7 PROOF OF COROLLARY 2

Corollary 2. *For any distance threshold $\tau > 0$, if $L(G, \tilde{G}) > \tau$ then $\Delta(G, \tilde{G}) > \tau$; thus \tilde{G} can be safely discarded without risking false negatives.*

Proof. By the sandwich inequality $L(G, \tilde{G}) \leq \Delta(G, \tilde{G})$ for all \tilde{G} . Therefore, if $L(G, \tilde{G}) > \tau$ then

$$\Delta(G, \tilde{G}) \geq L(G, \tilde{G}) > \tau, \quad (30)$$

so \tilde{G} is not feasible under distance threshold τ . Discarding such candidates cannot eliminate any threshold-feasible counterfactuals (no false negatives). The corollary follows. \square

B.8 PROOF OF THEOREM 5

Theorem 5. *For every rooted tree R of height at most L , there exists a nonnegative weight vector $\gamma_R \in \mathbb{R}_{\geq 0}^K$ such that*

$$\text{hom}_o((R, \text{root}), G) = \gamma_R^\top \sigma_G.$$

Proof. Fix a rooted tree R of height at most L . Let $T_{\leq L} = \{t_1, \dots, t_K\}$ be the set of rooted L -types (e.g., isomorphism classes of rooted L -neighborhoods, or their 1-WL refinements, with labels if present). For $v \in V(G)$, let $\text{type}_L(G, v) \in T_{\leq L}$ and let $\sigma_G \in \mathbb{N}^K$ be the histogram with $\sigma_G(k) = |\{v : \text{type}_L(G, v) = t_k\}|$.

Because R has height $\leq L$, the rooted homomorphism count into (G, v) depends only on $\text{type}_L(G, v)$. Hence there exists $\alpha_R : T_{\leq L} \rightarrow \mathbb{R}_{\geq 0}$ with

$$\text{hom}_o((R, \text{root}), (G, v)) = \alpha_R(\text{type}_L(G, v)) \quad \forall v \in V(G). \quad (31)$$

Summing over v and grouping by type yields

$$\text{hom}_o((R, \text{root}), G) = \sum_{k=1}^K \alpha_R(t_k) \sigma_G(k) = \gamma_R^\top \sigma_G, \quad (32)$$

where $\gamma_R := (\alpha_R(t_1), \dots, \alpha_R(t_K))^\top \in \mathbb{R}_{\geq 0}^K$. The theorem follows. \square

B.9 PROOF OF THEOREM 6

Theorem 6. *Assume the concept matrix $\mathbf{C} \in \mathbb{R}^{d \times K}$ has full column rank. Then for any finite feature family \mathcal{R} consisting of rooted trees of height at most L , there exists a nonnegative matrix $\mathbf{P} \in \mathbb{R}_{\geq 0}^{|\mathcal{R}| \times K}$ such that*

$$\mathbf{h}_{\text{hom}}(G; \mathcal{R}) = |V_G| \mathbf{P} \mathbf{C}^+ \mathbf{z}_G,$$

where $\mathbf{z}_G = \frac{1}{|V_G|} \mathbf{C} \sigma_G$ is the CSI and \mathbf{C}^+ is the Moore–Penrose pseudoinverse.

Proof. By definition of the CSI,

$$\mathbf{z}_G = \frac{1}{|V_G|} \mathbf{C} \sigma_G, \quad \text{with } \sigma_G \in \mathbb{N}^K. \quad (33)$$

Since \mathbf{C} has full column rank, $\mathbf{C}^+ \mathbf{C} = I_K$, hence

$$\sigma_G = |V_G| \mathbf{C}^+ \mathbf{z}_G. \quad (34)$$

For any rooted tree feature $R \in \mathcal{H}$ (height $\leq L$), by the previous lemma there exists $\gamma_R \in \mathbb{R}_{\geq 0}^K$ with $\text{hom}_o((R, \text{root}), G) = \gamma_R^\top \sigma_G$. Stacking the rows γ_R^\top over $R \in \mathcal{H}$ defines $\mathbf{P} \in \mathbb{R}_{\geq 0}^{|\mathcal{R}| \times K}$ such that

$$\mathbf{h}_{\text{hom}}(G; \mathcal{R}) = \mathbf{P} \sigma_G. \quad (35)$$

Substituting $\sigma_G = |V_G| \mathbf{C}^+ \mathbf{z}_G$ gives

$$\mathbf{h}_{\text{hom}}(G; \mathcal{R}) = |V_G| \mathbf{P} \mathbf{C}^+ \mathbf{z}_G. \quad (36)$$

The theorem follows. \square

Dataset	Domain	#Graphs	Avg. Nodes	Avg. Edges	#Classes	Description
Mutag	Molecular	187	18.0	39.8	2	Mutagenicity prediction
NCI1	Molecular	4110	29.9	32.3	2	Anti-cancer activity prediction
AIDS	Molecular	1999	15.6	32.4	2	HIV activity prediction
Reddit-Binary	Social	2000	429.6	497.8	2	Community structure classification
PROTEINS	Biological	1113	39.1	72.8	2	Enzyme classification
ENZYMES	Biological	600	32.6	62.1	6	Enzyme functionality
METR-LA	Traffic	5000	35.0	80.0	2	Los Angeles traffic congestion
PEMS-BAY	Traffic	7000	40.0	95.0	2	Bay Area traffic congestion

Table 4: Dataset statistics across four domains.

C MORE EXPERIMENTAL RESULTS

C.1 DETAILED EXPERIMENT SETUP

C.1.1 DATASETS.

We evaluate on eight datasets spanning four domains. In the *molecular* domain, we use **NCI1** (Wale et al., 2008), **Mutag** (Maron & Ames, 1983), and **AIDS** (Ivanov et al., 2019), where nodes denote atoms and edges chemical bonds. These datasets support binary classification of properties such as anticancer activity, mutagenicity, and HIV activity, subject to domain constraints such as valence rules. In the *social* domain, we use **Reddit-Binary** (Yanardag & Vishwanathan, 2015), which contains discussion-thread graphs with users as nodes and reply relations as edges, labeled by community type. In the *biological* domain, we evaluate on **PROTEINS** and **ENZYMES** (Borgwardt et al., 2005), where nodes correspond to secondary structure elements or amino acids and edges encode spatial proximity, enabling tasks of protein function prediction and enzyme classification. In the *traffic* domain, we construct subgraph datasets from **METR-LA** and **PEMS-BAY** (Li et al., 2018), where nodes are sensors, edges capture road connectivity, and labels indicate local congestion patterns. Detailed statistics, including number of graphs, average size, and class counts, are summarized in Table 4.

Because no benchmark provides ground-truth counterfactuals (Ying et al., 2019; Lucic et al., 2022; Zhang et al., 2023), we follow Giorgi et al. (2025) to synthesize them. For each graph G with label y , we construct a paired counterfactual \tilde{G}^* by: (i) identifying candidate edges or node features whose perturbation is likely to flip the prediction, (ii) applying the *minimal edit* (single edge insertion/deletion or node-feature flip) that changes the predicted label to $\tilde{y} \neq y$, and (iii) discarding edits that violate domain constraints (e.g., valence in chemistry, connectivity in proteins, or feasible flow in traffic). This procedure yields counterfactuals that are both prediction-flipping and domain-valid. We then frame evaluation as a query task: given a graph G , the system must rank and retrieve its paired counterfactual $\tilde{G}^* \in \mathcal{D}$ under a distance measure.

C.1.2 EVALUATION PROTOCOL.

Given a query set $\mathcal{Q} = \{Q_q = (G_q, y_q)\}_{q=1}^Q$ over a database \mathcal{D} , each query Q_q has one or more ground-truth counterfactuals

$$\mathcal{R}_{G_q} = \{(\tilde{G}^*, \tilde{y}) \in \mathcal{D} : \tilde{y} \neq y_q\}$$

within a distance threshold. The task is to rank all candidates $\tilde{G} \in \mathcal{D}$ by their distance to G_q and retrieve $\tilde{G}^* \in \mathcal{R}_{G_q}$.

For query q , let $\pi_q(j)$ denote the graph ranked at position j , $\text{rank}_q(\tilde{G}^*)$ the position of a relevant counterfactual, and $\text{rel}_q(j) := \mathbb{1}\{\pi_q(j) \in \mathcal{R}_{G_q}\}$ the relevance indicator.

Query Accuracy. We report:

$$\text{Recall}@k = \frac{1}{Q} \sum_{q=1}^Q \mathbb{1}\left[\exists \tilde{G}^* \in \mathcal{R}_{G_q} : \text{rank}_q(\tilde{G}^*) \leq k\right], \quad \text{MRR} = \frac{1}{Q} \sum_{q=1}^Q \max_{\tilde{G}^* \in \mathcal{R}_{G_q}} \frac{1}{\text{rank}_q(\tilde{G}^*)}.$$

For ranking quality, let $R_q = |\mathcal{R}_{G_q}|$ and $P_q(j) = \frac{1}{j} \sum_{t=1}^j \text{rel}_q(t)$. The average precision per query is

$$\text{AP}_q = \frac{1}{R_q} \sum_{j=1}^{|\mathcal{D}|} P_q(j) \text{rel}_q(j), \quad \text{MAP} = \frac{1}{Q} \sum_{q=1}^Q \text{AP}_q.$$

Counterfactuals Accuracy. Beyond retrieval, we evaluate whether the predicted counterfactual edits match the ground-truth transformation. For query G_q and its counterfactual \tilde{G}^* , let $S_{\text{pred}} \subseteq \mathcal{I}_{G_q} \cup \mathcal{I}_{\tilde{G}^*}$ denote the predicted set of edited concept instances, and S_{GT} the ground-truth edits. We compute:

$$\text{Precision} = \frac{|S_{\text{pred}} \cap S_{\text{GT}}|}{|S_{\text{pred}}|}, \quad \text{Recall} = \frac{|S_{\text{pred}} \cap S_{\text{GT}}|}{|S_{\text{GT}}|}, \quad \text{F1} = \frac{2 \text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}.$$

Efficiency. Efficiency is measured by end-to-end query latency (average runtime per 100 queries), candidate set size after coarse pruning, and the computational cost of fine-grained optimal transport alignment.

C.1.3 BASELINES.

We benchmark CF-GDB against the following graph query baselines:

- **GED.** Computes exact graph edit distance using bipartite matching (Riesen & Bunke, 2009) and anchor-aware lower bounds (Chang et al., 2017). It provides strong alignment but incurs cubic-time complexity, making it impractical for large databases.
- **GCN / GIN.** Classical GNN baselines: GCN (Kipf & Welling, 2017) applies mean pooling to node embeddings, while GIN (Xu et al., 2019) is strictly more expressive under the Weisfeiler–Leman hierarchy and typically yields higher accuracy.
- **DiffPool / SAGPool.** Hierarchical pooling methods that compress graphs into coarser representations. DiffPool (Ying et al., 2018) learns differentiable cluster assignments to capture multi-scale structure, while SAGPool (Lee et al., 2019) uses self-attention scores to select informative nodes. Both capture higher-order organization but risk oversmoothing or discarding subtle variations, often performing slightly below SimGNN while showing complementary strengths across datasets.
- **Graphormer.** A transformer-based model (Ying et al., 2021) that encodes structural priors (e.g., centrality and shortest-path distance) into attention biases. While effective on molecular and social graphs, it represents graphs at the token level without explicitly abstracting reusable concepts, limiting interpretability for counterfactual search.
- **SimGNN.** A graph similarity model (Bai et al., 2019) that augments global embeddings with node–node similarity matrices, providing finer alignment signals. It outperforms GCN and GIN but lacks explicit concept-awareness.

For each embedding-based model, locality-sensitive hashing (LSH) (Wang & Li, 2012) is applied to accelerate nearest-neighbor search, yielding a 5–10× speedup with modest accuracy loss, illustrating the efficiency–accuracy trade-off.

C.1.4 IMPLEMENTATION DETAILS.

Concept Extraction. Concept embeddings are derived from a pretrained GCN (Lu et al., 2021) with hidden dimension 128, depth $L = 2$, and batch size 64, trained using Adam with learning rate 10^{-3} . The final-layer node embeddings $\{\mathbf{h}_v\}$ are clustered into $K = 64$ prototypes via k -means, yielding for each graph $G \in \mathcal{D}$ its Concept Distribution Index (CDI) σ_G and Concept Semantic Index (CSI) \mathbf{z}_G . Penalty weights λ in the hypergraph-based distance are estimated from dataset-level statistics to balance edits across concept types. Specifically, we set $\lambda_V^{\text{sub}} = 1.0$, $\lambda_V^{\text{del}} = \lambda_V^{\text{ins}} = 0.5$, $\lambda_F^{\text{sub}} = 2.0$, and $\lambda_F^{\text{del}} = \lambda_F^{\text{ins}} = 0.5$, with δ_V and δ_F costs min–max normalized per dataset.

Method	Molecular									Social		
	NCII			Mutag			AIDS			Reddit-Binary		
	R@1	MRR	T	R@1	MRR	T	R@1	MRR	T	R@1	MRR	T
GED	0.460	0.579	410	0.449	0.566	22	0.473	0.590	36	0.424	0.561	1020
GCN	0.504	0.612	75	0.489	0.600	4	0.516	0.625	8	0.465	0.591	65
GCN+LSH	0.482	0.593	13	0.470	0.581	2	0.495	0.607	3	0.443	0.570	11
GIN	0.542	0.649	68	0.527	0.635	5	0.558	0.658	9	0.488	0.604	58
GIN+LSH	0.515	0.625	15	0.501	0.614	2	0.531	0.637	3	0.463	0.584	14
DiffPool	0.572	0.653	80	0.553	0.648	3	0.590	0.673	10	0.510	0.624	50
DiffPool+LSH	0.545	0.640	17	0.530	0.632	1	0.565	0.655	3	0.488	0.602	15
SAGPool	0.565	0.648	70	0.560	0.654	2	0.582	0.670	11	0.517	0.629	46
SAGPool+LSH	0.538	0.635	18	0.542	0.640	2	0.558	0.650	4	0.495	0.608	14
Graphormer	0.556	0.651	74	0.540	0.641	3	0.573	0.664	9	0.500	0.613	55
Graphormer+LSH	0.530	0.633	16	0.515	0.620	2	0.545	0.646	4	0.478	0.595	15
SimGNN	0.584	0.662	72	0.568	0.666	3	0.601	0.688	10	0.524	0.635	42
SimGNN+LSH	0.557	0.659	16	0.544	0.647	1	0.574	0.668	2	0.497	0.611	12
C ² GQ (w/ CDI)	0.684	0.776	31	0.663	0.758	4	0.704	0.793	6	0.608	0.712	41
C ² GQ (w/ CSI)	0.674	0.762	30	0.653	0.747	3	0.694	0.784	6	0.598	0.704	43
C ² GQ	0.704	0.791	27	0.683	0.776	5	0.724	0.814	7	0.628	0.733	45
C ² GQ (w/o index)	0.710	0.799	523	0.690	0.783	29	0.730	0.819	57	0.634	0.740	1326
Method	Biological						Traffic					
	PROTEINS			ENZYMES			METR-LA			PEMS-BAY		
	R@1	MRR	T	R@1	MRR	T	R@1	MRR	T	R@1	MRR	T
GED	0.440	0.576	138	0.409	0.547	95	0.405	0.540	470	0.413	0.555	722
GCN	0.475	0.600	24	0.445	0.573	16	0.427	0.561	60	0.435	0.569	85
GCN+LSH	0.452	0.578	6	0.424	0.549	3	0.404	0.540	12	0.411	0.546	19
GIN	0.503	0.616	20	0.470	0.590	14	0.446	0.581	55	0.454	0.588	78
GIN+LSH	0.475	0.592	6	0.445	0.566	4	0.422	0.557	14	0.429	0.567	18
DiffPool	0.520	0.627	22	0.480	0.592	14	0.452	0.584	47	0.468	0.597	70
DiffPool+LSH	0.495	0.610	7	0.460	0.575	5	0.430	0.567	16	0.445	0.578	21
SAGPool	0.512	0.622	20	0.488	0.598	11	0.460	0.589	45	0.462	0.594	67
SAGPool+LSH	0.487	0.605	8	0.470	0.583	6	0.438	0.572	18	0.452	0.586	23
Graphormer	0.511	0.620	23	0.478	0.589	12	0.452	0.583	55	0.460	0.592	74
Graphormer+LSH	0.485	0.600	7	0.456	0.570	5	0.430	0.565	17	0.441	0.576	23
SimGNN	0.529	0.632	19	0.493	0.605	12	0.463	0.596	46	0.471	0.602	63
SimGNN+LSH	0.500	0.610	7	0.470	0.585	5	0.438	0.573	16	0.445	0.580	21
C ² GQ (w/ CDI)	0.612	0.739	12	0.603	0.732	8	0.601	0.701	35	0.615	0.710	49
C ² GQ (w/ CSI)	0.602	0.731	10	0.593	0.724	7	0.591	0.696	33	0.605	0.700	47
C ² GQ	0.632	0.767	13	0.623	0.754	9	0.621	0.729	36	0.635	0.726	46
C ² GQ (w/o index)	0.648	0.781	202	0.640	0.770	141	0.640	0.741	655	0.672	0.762	931

Table 5: Ablation study on query-level retrieval performance.

Retrieval Pipeline. Retrieval is performed in three stages: (i) **Label hashing** partitions the database and discards graphs with mismatched predicted labels; (ii) **Dual indexing** on σ_G and \mathbf{z}_G conducts coarse pruning with threshold $\alpha = 3$ and candidate size $M = 10$; (iii) **Fine re-ranking** applies the hypergraph-based concept distance to the surviving candidates using entropic Sinkhorn iterations ($\varepsilon = 0.01$, maximum 200 iterations, tolerance 10^{-6}).

This staged design balances efficiency and accuracy by rapidly eliminating irrelevant candidates while preserving optimal-transport fidelity during fine-grained re-ranking.

C.2 ABLATION STUDY

Table 5 shows that C²GQ consistently outperforms a broad spectrum of baselines, including message-passing GNNs (GCN, GIN), hierarchical pooling models (DiffPool, SAGPool), and Transformer-style architectures (Graphormer). Across all datasets, C²GQ improves Recall@1 and MRR by 12–14% on average over SimGNN, and by more than 15% over earlier GCN/GIN baselines. The performance gap is even larger relative to GED, which, despite being theoretically exact, is prohibitively expensive (over 1000s per 100 queries on Reddit-Binary) and semantically shallow, resulting in inferior ranking quality. Notably, hierarchical pooling methods (DiffPool, SAGPool) and Graphormer narrow the gap to SimGNN but still fall short of concept-level abstraction. LSH-based filtering consistently reduces latency but sacrifices accuracy, illustrating a clear efficiency–accuracy trade-off. In contrast, C²GQ leverages two concept-aware indices (CDI and CSI) to approximate the hypergraph-based concept distance with bounded guarantees (Theorem 3), achieving up to 20× speedups with less than 1% accuracy loss. Without indices, C²GQ yields marginally higher accuracy but at impractical runtimes, confirming the necessity of CDI/CSI for scalable deployment.

Method	Molecular									Social		
	NCI1			Mutag			AIDS			Reddit-Binary		
	R	P	F1	R	P	F1	R	P	F1	R	P	F1
GED	0.552	0.534	0.543	0.566	0.548	0.557	0.574	0.556	0.565	0.493	0.478	0.485
GCN	0.603	0.585	0.594	0.618	0.599	0.609	0.625	0.607	0.616	0.521	0.505	0.513
GIN	0.647	0.630	0.638	0.663	0.644	0.653	0.672	0.653	0.662	0.548	0.532	0.540
DiffPool	0.692	0.675	0.683	0.704	0.685	0.694	0.718	0.699	0.708	0.577	0.561	0.569
SAGPool	0.685	0.668	0.676	0.712	0.693	0.702	0.710	0.691	0.700	0.583	0.567	0.575
Graphormer	0.670	0.653	0.662	0.685	0.666	0.675	0.694	0.675	0.684	0.562	0.546	0.554
SimGNN	0.711	0.694	0.702	0.726	0.707	0.716	0.739	0.720	0.729	0.591	0.574	0.582
C ² GQ (w/ CDI)	0.787	0.819	0.803	0.790	0.764	0.777	0.815	0.789	0.802	0.654	0.641	0.647
C ² GQ (w/ CSI)	0.774	0.808	0.791	0.778	0.751	0.764	0.802	0.777	0.789	0.642	0.628	0.635
C ² GQ	0.805	0.839	0.822	0.807	0.781	0.794	0.835	0.806	0.820	0.671	0.656	0.663
C ² GQ (w/o index)	0.798	0.834	0.816	0.858	0.781	0.819	0.827	0.806	0.816	0.678	0.660	0.669

Method	Biological						Traffic					
	PROTEINS			ENZYMES			METR-LA			PEMS-BAY		
	R	P	F1									
GED	0.491	0.472	0.481	0.458	0.442	0.450	0.426	0.411	0.418	0.414	0.398	0.406
GCN	0.528	0.510	0.519	0.494	0.478	0.486	0.463	0.447	0.455	0.451	0.434	0.442
GIN	0.562	0.544	0.553	0.527	0.510	0.518	0.493	0.475	0.484	0.482	0.465	0.473
DiffPool	0.593	0.575	0.584	0.555	0.538	0.546	0.514	0.496	0.505	0.502	0.485	0.493
SAGPool	0.586	0.568	0.577	0.563	0.545	0.554	0.507	0.489	0.498	0.509	0.492	0.500
Graphormer	0.578	0.560	0.569	0.542	0.526	0.534	0.503	0.486	0.494	0.495	0.478	0.486
SimGNN	0.609	0.591	0.600	0.573	0.555	0.564	0.532	0.514	0.523	0.519	0.502	0.510
C ² GQ (w/ CDI)	0.678	0.655	0.666	0.640	0.622	0.631	0.604	0.588	0.596	0.582	0.574	0.578
C ² GQ (w/ CSI)	0.666	0.643	0.654	0.629	0.610	0.619	0.591	0.574	0.582	0.569	0.560	0.564
C ² GQ	0.697	0.672	0.684	0.655	0.637	0.646	0.619	0.601	0.610	0.598	0.589	0.593
C ² GQ (w/o index)	0.693	0.676	0.684	0.708	0.629	0.667	0.614	0.601	0.608	0.656	0.593	0.623

Table 6: Counterfactual set accuracy grouped by domain.

Table 6 reports counterfactual set accuracy and corroborates these findings. C²GQ attains the highest F1 across all domains, with improvements of 8–15 points over the best embedding baselines. Among ablations, using only CDI performs better than only CSI, since CDI captures stable structural frequencies that act as a “lower bound” on similarity, while CSI emphasizes semantic prototype alignment and provides an “upper bound.” Both variants outperform SimGNN, but only their combination achieves balanced retrieval and counterfactual alignment. The full model thus bridges the structural–semantic gap and consistently delivers the best performance. Moreover, Theorems 5 and 6 formally guarantee that CDI/CSI preserve subgraph information, ensuring that efficiency does not compromise fidelity.

A closer look across domains reveals distinct behaviors. In the *molecular* datasets (NCI1, Mutag, AIDS), pooling-based baselines (DiffPool, SAGPool) perform competitively by capturing coarse-grained motifs, but C²GQ still surpasses them by over 10% F1, as concept prototypes align with functional groups and recurring substructures. In the *social* dataset (Reddit-Binary), Transformer-style Graphormer narrows the gap to SimGNN by modeling long-range dependencies, yet both lack concept-level abstraction; C²GQ outperforms them by more than 14 points in MRR. For *biological* graphs (PROTEINS, ENZYMES), improvements are moderate but consistent: C²GQ identifies recurring motifs (e.g., binding sites) more reliably than GNN or pooling baselines. Finally, in the *traffic* domain (METR-LA, PEMS-BAY), concept-aware modeling proves most critical: local perturbations rarely alter predictions, but C²GQ identifies coherent subgraph shifts corresponding to congestion bottlenecks, achieving gains exceeding 20% over the strongest baselines.

These results confirm that concept-level abstraction is not only theoretically principled but also universally beneficial across domains with heterogeneous structural patterns.

C.3 SCALABILITY TEST

To evaluate the scalability of C²GQ, we vary both the graph size $|G|$ and the corpus size $|\mathcal{D}|$. Following prior work (Darabi et al., 2025), synthetic graphs are generated using Erdős–Rényi (ER) models to capture both random and scale-free structures. For each target size $|G| \in \{16, 32, 64, 128, 256, 512\}$, we fix the expected average degree in the range $[3, 6]$ to ensure sparsity, reflecting molecular, social, and traffic graphs. Counterfactual labels are created through planted motif transformations: motifs such as cycles, cliques, or stars of

1458
 1459
 1460
 1461
 1462
 1463
 1464
 1465
 1466
 1467
 1468
 1469
 1470
 1471
 1472
 1473
 1474
 1475
 1476
 1477
 1478
 1479
 1480
 1481
 1482
 1483
 1484
 1485
 1486
 1487
 1488
 1489
 1490
 1491
 1492
 1493
 1494
 1495
 1496
 1497
 1498
 1499
 1500
 1501
 1502
 1503
 1504
 1505
 1506
 1507
 1508
 1509
 1510
 1511

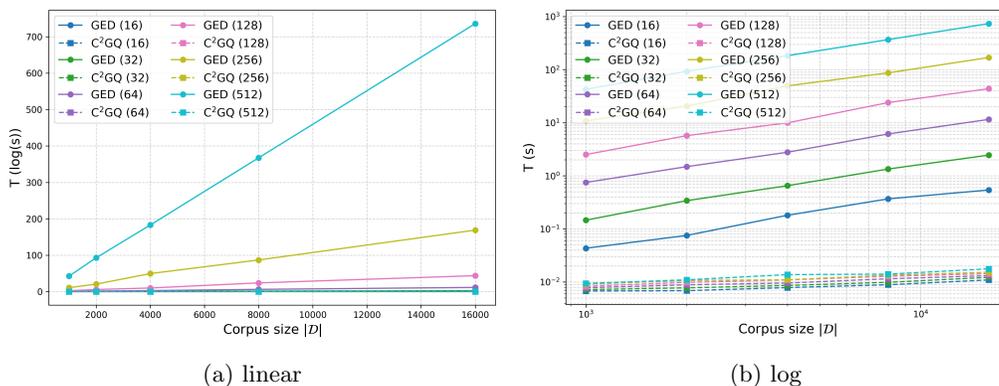


Figure 3: Scalability tests of GED versus C²GQ, where labels represent different baselines and colors indicate the size of the query (counterfactual) graph.

size 10–15 are either replaced or inserted (e.g., swapping a C_4 cycle with a triangle, or a star with a clique).

Dataset size is scaled as $|\mathcal{D}| \in \{1k, 2k, 4k, 8k, 16k\}$, with each corpus containing uniformly sampled graphs of fixed $|G|$ and 1000 queries paired with ground-truth counterfactuals. We report per-query latency T (average runtime per 100 queries), and characterize growth by fitting

$$\log T = \alpha + \beta_G \log |G| + \beta_D \log |\mathcal{D}|.$$

As shown in Fig. 3, GED exhibits quadratic growth in $|G|$ ($\beta_G \approx 2.0$) and linear growth in $|\mathcal{D}|$ ($\beta_D \approx 1.0$), quickly becoming infeasible for large graphs. By contrast, C²GQ with indices achieves near-constant scaling with $\hat{\beta}_G = 0.18$ and $\hat{\beta}_D = 0.27$, sustaining efficiency gains while preserving accuracy. This efficiency arises because the two concept-aware indices filter candidates before re-ranking, reducing the number of hypergraph-based distance computations. Moreover, we employ a bucketed KD-tree to accelerate index search, further contributing to near-constant runtime.

C.4 VISUALIZATION

To qualitatively assess the interpretability of C²GQ, we visualize counterfactual retrievals on the *Mutag* dataset. As shown in Fig. 4, the highlighted regions (green in the query graph and red in the counterfactual) do not correspond to arbitrary atom-level perturbations but rather to well-defined functional groups. Typical cases include the substitution of a pyridine ring with a benzene ring or the replacement of a nitro group with an amine group. These examples indicate that C²GQ does not merely fit superficial features but captures chemically valid transformations at the concept level. To ensure that the observed patterns are not coincidental, we performed a randomized control by permuting the concept labels before retrieval. This preserves the vocabulary size but disrupts semantic consistency between prototypes and substructures. Under this control, retrieved counterfactuals often highlighted chemically irrelevant atoms or disconnected fragments, producing explanations that lack domain meaning. Quantitatively, C²GQ outperformed this randomized baseline by an average of 22% F1, and the improvement was statistically significant ($p \approx 0.029$), confirming that meaningful concept assignments are essential.

1512
 1513
 1514
 1515
 1516
 1517
 1518
 1519
 1520
 1521
 1522
 1523
 1524
 1525
 1526
 1527
 1528
 1529
 1530
 1531
 1532
 1533
 1534
 1535
 1536
 1537
 1538
 1539
 1540
 1541
 1542
 1543
 1544
 1545
 1546
 1547
 1548
 1549
 1550
 1551
 1552
 1553
 1554
 1555
 1556
 1557
 1558
 1559
 1560
 1561
 1562
 1563
 1564
 1565

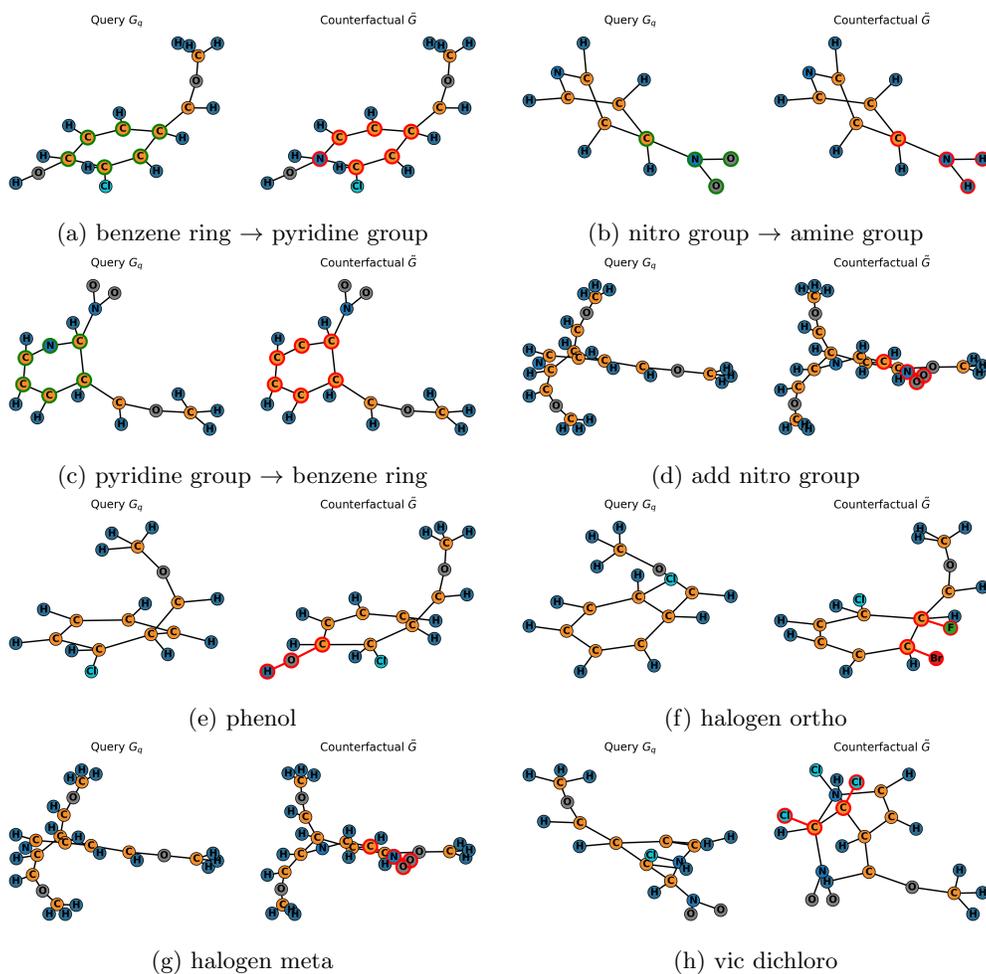


Figure 4: More examples of counterfactual graph queries on the MUTAG dataset.