# Extending the Unified Mathematical Framework to Support Graph Data Structures

**Jaepil Jeong**
Cognica, Inc.
*Email: jaepil@cognica.io*
Date: December 25, 2024

*"Language is a labyrinth of paths. You approach from one side and know your way about; you approach the same place from another side and no longer know your way about."*
— Ludwig Wittgenstein

## Abstract

This paper presents a formal extension of the unified query algebra framework to incorporate graph data structures and operations. Building upon the rigorous mathematical foundation established by Jeong (2023), we introduce type-theoretic definitions for graph elements, define traversal and pattern matching operators, establish algebraic properties of graph operations, and extend the category-theoretic foundations to encompass graph structures. We prove that our extension preserves the completeness of the existing Boolean algebra while enabling seamless integration of graph operations with relational, textual, and vector operations. Additionally, we analyze the computational complexity and optimization strategies for cross-paradigm queries involving graphs. This extension provides a comprehensive theoretical foundation for next-generation data systems that must operate across relational, textual, vector, and graph paradigms.

## 1. Graph-Extended Type-Theoretic Foundation

### 1.1 Graph Type System

We begin by extending the existing type system to incorporate graph-specific types.

**Definition 1.1.1** (Graph Universe). Let $\mathcal{G}$ be the universe of all possible graphs in our data system.

**Definition 1.1.2** (Extended Type Hierarchy). We extend the existing type hierarchy with:

- $\mathcal{V}_G \subset \mathcal{G}$: The set of all vertices (nodes)
- $\mathcal{E}_G \subset \mathcal{G}$: The set of all edges
- $\mathcal{P}_G$: The set of all vertex and edge properties
- $\mathcal{L}_G$: The set of all edge labels

**Definition 1.1.3** (Graph Structure). A graph $G \in \mathcal{G}$ is defined as a tuple:

$$G = (V, E, \lambda_V, \lambda_E) \tag{1}$$

where:

- $V \subset \mathcal{V}_G$ is a set of vertices
- $E \subset \mathcal{E}_G$ is a set of edges, where $E \subseteq V \times V \times \mathcal{L}_G$
- $\lambda_V : V \times \mathcal{F} \to \mathcal{P}_G$ is a vertex property mapping function
- $\lambda_E : E \times \mathcal{F} \to \mathcal{P}_G$ is an edge property mapping function

**Definition 1.1.4** (Graph Posting List). A graph posting list $L_G \in \mathcal{L}_G$ is defined as an ordered sequence:

$$L_G = [(id_1, G_1), (id_2, G_2), \ldots, (id_n, G_n)] \tag{2}$$

where:

- $id_i \in \mathbb{N}$ is a unique identifier
- $G_i \in \mathcal{G}$ is a graph or subgraph

**Definition 1.1.5** (Graph-Document Bijection). We define bijective mappings between graphs and documents:

$$\phi_{G \to D} : \mathcal{G} \to 2^{\mathcal{D}} \tag{3}$$

$$\phi_{D \to G} : 2^{\mathcal{D}} \to \mathcal{G} \tag{4}$$

These mappings establish an isomorphism between graph structures and document collections, allowing us to leverage the existing framework.

**Theorem 1.1.6** (Graph-Posting List Isomorphism). The graph posting lists $\mathcal{L}_G$ and standard posting lists $\mathcal{L}$ form isomorphic algebraic structures under their respective operations.

*Proof.* We can establish an isomorphism $\Phi : \mathcal{L}_G \to \mathcal{L}$ such that:

$$\Phi(L_G) = PL(\bigcup_{i=1}^{n} \phi_{G \to D}(G_i)) \tag{5}$$

For any graph posting lists $L_G^1, L_G^2 \in \mathcal{L}_G$, we can verify that:

$$\Phi(L_G^1 \cup_G L_G^2) = \Phi(L_G^1) \cup \Phi(L_G^2) \tag{6}$$

$$\Phi(L_G^1 \cap_G L_G^2) = \Phi(L_G^1) \cap \Phi(L_G^2) \tag{7}$$

where $\cup_G$ and $\cap_G$ are the graph-specific union and intersection operations.

# 2. Graph Algebra of Posting Lists

## 2.1 Core Graph Algebraic Operations

We define fundamental operations on graph posting lists that integrate with the existing Boolean algebra.

**Definition 2.1.1** (Graph Posting List Operations). We define:

- Union: $\cup_G : \mathcal{L}_G \times \mathcal{L}_G \to \mathcal{L}_G$ where $L_G^1 \cup_G L_G^2 = [(id, G)|G \in G_1 \cup G_2]$
- Intersection: $\cap_G : \mathcal{L}_G \times \mathcal{L}_G \to \mathcal{L}_G$ where $L_G^1 \cap_G L_G^2 = [(id, G)|G \in G_1 \cap G_2]$
- Difference: $\setminus_G : \mathcal{L}_G \times \mathcal{L}_G \to \mathcal{L}_G$ where $L_G^1 \setminus_G L_G^2 = [(id, G)|G \in G_1 \setminus G_2]$
- Complement: $\overline{\cdot}_G : \mathcal{L}_G \to \mathcal{L}_G$ where $\overline{L_G} = [(id, G)|G \in \mathcal{G} \setminus G_L]$

**Theorem 2.1.2** (Extended Boolean Algebra). The structure $(\mathcal{L}_G, \cup_G, \cap_G, \overline{\cdot}_G, \emptyset_G, \mathcal{G})$ forms a Boolean algebra satisfying commutativity, associativity, distributivity, identity, and complement properties.

*Proof.* The proof follows analogously to Theorem 2.1.2 in the original paper, leveraging the isomorphism established in Theorem 1.1.6.

## 2.2 Graph-Specific Operations

**Definition 2.2.1** (Graph Traversal Operator). For a vertex $v \in \mathcal{V}_G$, label $l \in \mathcal{L}_G$, and number of hops $k \in \mathbb{N}$, we define:

$$Traverse_{v,l,k} : \mathcal{L}_G \to \mathcal{L}_G \tag{8}$$

where $Traverse_{v,l,k}(L_G)$ returns a posting list containing all subgraphs that can be reached from vertex $v$ by traversing edges with label $l$ up to $k$ hops.

**Definition 2.2.2** (Pattern Matching Operator). For a pattern graph $P \in \mathcal{G}$, we define:

$$Match_P : \mathcal{L}_G \to \mathcal{L}_G \tag{9}$$

where $Match_P(L_G)$ returns a posting list containing all subgraphs that match pattern $P$ through subgraph isomorphism.

**Definition 2.2.3** (Vertex-Centric Aggregation). For a property $p \in \mathcal{P}_G$ and aggregation function $agg$, we define:

$$VertexAgg_{p,agg} : \mathcal{L}_G \to \mathcal{A} \tag{10}$$

where $VertexAgg_{p,agg}(L_G)$ applies the aggregation function $agg$ to property $p$ across all vertices in graphs contained in $L_G$.

# 3. Cross-Paradigm Integration

## 3.1 Graph-Relational Integration

**Definition 3.1.1** (Document-to-Graph Operator). We define an operator that constructs a graph from relational data:

$$ToGraph_{v_f,e_f} : \mathcal{L} \to \mathcal{L}_G \tag{11}$$

where $v_f$ is a field defining vertices and $e_f$ is a field defining edges.

**Definition 3.1.2** (Graph-to-Document Operator). We define an operator that flattens a graph into relational data:

$$FromGraph_{v_f,e_f} : \mathcal{L}_G \to \mathcal{L} \tag{12}$$

where $v_f$ is a field for storing vertex data and $e_f$ is a field for storing edge data.

## 3.2 Graph-Vector Integration

**Definition 3.2.1** (Vertex Embedding Function). We define:

$$vec_V : \mathcal{V}_G \times \mathcal{F} \to \mathcal{V} \tag{13}$$

such that $vec_V(v, f)$ maps a vertex $v$ and property $f$ to a vector representation.

**Definition 3.2.2** (Graph Embedding Function). We define:

$$vec_G : \mathcal{G} \to \mathcal{V} \tag{14}$$

such that $vec_G(G)$ maps an entire graph to a vector representation.

**Definition 3.2.3** (Vector-Enhanced Graph Pattern Match). For a pattern $P$, query vector $q$, and threshold $\theta$, we define:

$$VectorMatch_{P,q,\theta} : \mathcal{L}_G \to \mathcal{L}_G \tag{15}$$

where $VectorMatch_{P,q,\theta}(L_G)$ returns subgraphs that both match pattern $P$ and have a vector similarity with $q$ above threshold $\theta$.

## 3.3 Graph-Text Integration

**Definition 3.3.1** (Text-to-Graph Operator). For text documents $T \subset \mathcal{D}$, we define:

$$TextToGraph : \mathcal{L} \to \mathcal{L}_G \tag{16}$$

which constructs a graph where vertices are terms and edges represent co-occurrence.

**Definition 3.3.2** (Semantic Graph Search). For a query text $q \in \mathcal{T}^*$ and threshold $\theta$, we define:

$$SemanticGraphSearch_{q,\theta} : \mathcal{L}_G \to \mathcal{L}_G \tag{17}$$

which returns subgraphs semantically related to the query text with similarity above $\theta$.

# 4. Graph Pattern Join Operations

## 4.1 Graph Join Operators

**Definition 4.1.1** (Graph Join). For graph posting lists $L_G^1, L_G^2 \in \mathcal{L}_G$ and join vertices $v_1, v_2 \in \mathcal{V}_G$, we define:

$$GraphJoin_{v_1=v_2} : \mathcal{L}_G \times \mathcal{L}_G \to \mathcal{L}_G \tag{18}$$

where $GraphJoin_{v_1=v_2}(L_G^1, L_G^2)$ returns a posting list of graphs formed by joining graphs from $L_G^1$ and $L_G^2$ where vertices $v_1$ and $v_2$ match.

**Definition 4.1.2** (Cross-Paradigm Graph Join). For a graph posting list $L_G$ and document posting list $L$, with join fields $v_f \in \mathcal{F}$ and $d_f \in \mathcal{F}$, we define:

$$CrossJoin_{v_f=d_f} : \mathcal{L}_G \times \mathcal{L} \to \mathcal{L}_G \tag{19}$$

where $CrossJoin_{v_f=d_f}(L_G, L)$ joins graphs and documents on matching field values.

## 4.2 Graph Join Properties

**Theorem 4.2.1** (Graph Join Commutativity). The graph join operation is commutative up to isomorphism:

$$GraphJoin_{v_1=v_2}(L_G^1, L_G^2) \cong GraphJoin_{v_2=v_1}(L_G^2, L_G^1) \tag{20}$$

*Proof.* Similar to the proof in Theorem 4.3.1 of the original paper, we can establish a bijection between the result sets that preserves graph structure.

**Theorem 4.2.2** (Graph Join Associativity). The graph join operation is associative:

$$GraphJoin_{v_2=v_3}(GraphJoin_{v_1=v_2}(L_G^1, L_G^2), L_G^3) \cong GraphJoin_{v_1=v_2}(L_G^1, GraphJoin_{v_2=v_3}(L_G^2, L_G^3)) \tag{21}$$

*Proof.* The proof follows the structure of Theorem 4.3.2 in the original paper, extended to graph-specific join semantics.

# 5. Graph Path Expressions and Pattern Matching

## 5.1 Regular Path Queries

**Definition 5.1.1** (Regular Path Expression). A regular path expression $R$ over edge labels $\mathcal{L}_G$ is defined recursively:

- $l \in \mathcal{L}_G$ is a regular path expression
- If $R_1$ and $R_2$ are regular path expressions, then so are:
  - $R_1 \cdot R_2$ (concatenation)
  - $R_1|R_2$ (alternation)
  - $R_1^*$ (Kleene star)

**Definition 5.1.2** (Regular Path Query Operator). For a regular path expression $R$, we define:

$$RPQ_R : \mathcal{L}_G \to \mathcal{L}_G \tag{22}$$

where $RPQ_R(L_G)$ returns a posting list of all vertex pairs connected by paths matching $R$.

## 5.2 Graph Pattern Expressions

**Definition 5.2.1** (Graph Pattern Expression). A graph pattern expression $P$ is defined as a tuple:

$$P = (V_P, E_P, C_V, C_E) \tag{23}$$

where:

- $V_P$ is a set of vertex variables
- $E_P \subseteq V_P \times V_P \times \mathcal{L}_G$ is a set of edge patterns
- $C_V$ is a set of constraints on vertex properties
- $C_E$ is a set of constraints on edge properties

**Definition 5.2.2** (Graph Pattern Match Operator). For a graph pattern expression $P$, we define:

$$GMatch_P : \mathcal{L}_G \to \mathcal{L}_G \tag{24}$$

where $GMatch_P(L_G)$ returns a posting list of all subgraphs that match pattern $P$.

**Theorem 5.2.3** (Equivalence of Pattern Expressions). For graph pattern expressions $P_1$ and $P_2$, if they are equivalent under graph homomorphism, then:

$$GMatch_{P_1}(L_G) = GMatch_{P_2}(L_G) \tag{25}$$

for all $L_G \in \mathcal{L}_G$.

*Proof*. If $P_1$ and $P_2$ are equivalent under graph homomorphism, then there exists a bijective mapping between the variables in $V_{P_1}$ and $V_{P_2}$ that preserves the edge patterns and constraints. Therefore, any subgraph matching $P_1$ will also match $P_2$ and vice versa.

# 6. Query Optimization for Graph Operations

## 6.1 Graph-Specific Rewrite Rules

**Theorem 6.1.1** (Graph Pattern Pushdown). For a graph pattern $P$ and a filter condition $F$, we have:

$$Filter_F(GMatch_P(L_G)) \equiv GMatch_{P'}(L_G) \tag{26}$$

where $P'$ is $P$ with the filter condition incorporated.

*Proof*. The filter condition $F$ can be incorporated into the pattern constraints $C_V$ or $C_E$ to create an augmented pattern $P'$. Any subgraph matching $P'$ will necessarily match both $P$ and satisfy condition $F$.

**Theorem 6.1.2** (Join-Pattern Fusion). For graph patterns $P_1$ and $P_2$ with a common join vertex variable, we have:

$$GMatch_{P_1}(L_G) \bowtie GMatch_{P_2}(L_G) \equiv GMatch_{P_1 \sqcup P_2}(L_G) \tag{27}$$

where $P_1 \sqcup P_2$ is the pattern formed by merging $P_1$ and $P_2$ on their common variables.

*Proof*. Joining the results of two pattern matches on common variables produces the same result as matching the merged pattern directly. This follows from the definition of graph pattern matching and join semantics.

## 6.2 Cost Model for Graph Operations

**Definition 6.2.1** (Graph Operation Cost Functions). We define cost models for graph operations:

- For pattern matching: $Cost(GMatch_P, G) = O(|V_G|^{|V_P|})$ in the worst case
- For traversal: $Cost(Traverse_{v,l,k}, G) = O(\sum_{i=1}^{k} d^i)$ where $d$ is the average degree
- For regular path queries: $Cost(RPQ_R, G) = O(|V_G|^2 \cdot |R|)$ where $|R|$ is the size of the regular expression

## 6.3 Cardinality Estimation for Graph Queries

**Definition 6.3.1** (Graph Pattern Selectivity). For a graph pattern $P$ and graph $G$, the selectivity $sel(P, G)$ is defined as:

$$sel(P, G) = \frac{|GMatch_P(G)|}{|\text{possible subgraphs of size } |V_P| \text{ in } G|} \tag{28}$$

**Theorem 6.3.2** (Independence-Based Cardinality Estimation). Under the assumption of edge independence, the estimated cardinality of a pattern match is:

$$|GMatch_P(G)| \approx |V_G|^{|V_P|} \cdot \prod_{(u,v,l) \in E_P} \frac{|E_G(l)|}{|V_G|^2} \tag{29}$$

where $|E_G(l)|$ is the number of edges with label $l$ in $G$.

*Proof.* Assuming independence between edges, the probability of a random edge having label $l$ is $\frac{|E_G(l)|}{|V_G|^2}$. For a pattern with $|V_P|$ vertices, there are $|V_G|^{|V_P|}$ possible assignments of graph vertices to pattern vertices. Each edge constraint in the pattern reduces this count by multiplying by the edge probability.

# 7. Category-Theoretic Extensions for Graphs

## 7.1 Graph Category Integration

**Definition 7.1.1** (Graph Category). We define a category $\mathcal{C}_G$ where:

- Objects are graph posting lists $L_G \in \mathcal{L}_G$
- Morphisms are graph query operators $op_G : L_G^1 \to L_G^2$
- Composition is operator composition
- Identity morphisms are identity operators $I_{L_G} : L_G \to L_G$

**Definition 7.1.2** (Graph Functor). We define a functor $F_G : \mathcal{C}_G \to \mathcal{C}$ from the graph category to the original query category:

- For objects: $F_G(L_G) = \Phi(L_G)$ where $\Phi$ is the isomorphism from Theorem 1.1.6
- For morphisms: $F_G(op_G) = \Phi \circ op_G \circ \Phi^{-1}$

**Theorem 7.1.3** (Functor Properties). The functor $F_G$ preserves the algebraic structure:

$$F_G(L_G^1 \cup_G L_G^2) = F_G(L_G^1) \cup F_G(L_G^2) \tag{30}$$

$$F_G(L_G^1 \cap_G L_G^2) = F_G(L_G^1) \cap F_G(L_G^2) \tag{31}$$

*Proof.* These properties follow directly from the isomorphism established in Theorem 1.1.6.

## 7.2 Monad Extensions for Graph Operations

**Definition 7.2.1** (Graph Monad). We define a graph monad $(T_G, \eta_G, \mu_G)$:

- $T_G : \mathcal{C} \to \mathcal{C}_G$ maps query results to graph query results
- $\eta_G : 1_{\mathcal{C}} \to F_G \circ T_G$ is the unit natural transformation
- $\mu_G : T_G \circ F_G \circ T_G \to T_G$ is the multiplication natural transformation

**Theorem 7.2.2** (Adjunction Relationship). The functors $F_G : \mathcal{C}_G \to \mathcal{C}$ and $T_G : \mathcal{C} \to \mathcal{C}_G$ form an adjunction $T_G \dashv F_G$.

*Proof.* For any morphism $f : T_G(C) \to G$ in $\mathcal{C}_G$, there is a corresponding morphism $g : C \to F_G(G)$ in $\mathcal{C}$ defined by $g = F_G \circ f \circ \eta_G$. Conversely, for any morphism $g : C \to F_G(G)$ in $\mathcal{C}$, there is a corresponding morphism $f : T_G(C) \to G$ in $\mathcal{C}_G$ defined by $f = \mu_G \circ T_G(g)$.

# 8. Computational Complexity and Theoretical Limitations

## 8.1 Complexity Analysis of Graph Operations

**Theorem 8.1.1** (Pattern Matching Complexity). The problem of subgraph isomorphism (implemented by $GMatch_P$) is NP-complete in general. For a pattern $P$ and graph $G$:

$$\text{Time}(GMatch_P, G) = O(|V_G|^{|V_P|}) \tag{32}$$

*Proof.* Subgraph isomorphism is a well-known NP-complete problem. The brute-force approach requires checking all possible mappings from pattern vertices to graph vertices, yielding the stated complexity.

**Theorem 8.1.2** (Regular Path Query Complexity). For a regular path expression $R$ and graph $G$:

$$\text{Time}(RPQ_R, G) = O(|V_G|^2 \cdot |R|) \tag{33}$$

where $|R|$ is the size of the regular expression.

*Proof.* Using dynamic programming techniques based on context-free language recognition, regular path queries can be evaluated in the stated complexity.

## 8.2 Theoretical Limitations

**Theorem 8.2.1** (Expressivity Hierarchy). The query languages form a strict hierarchy of expressivity:

$$\text{Relational} \subsetneq \text{Relational} + \text{Graph} \subsetneq \text{Relational} + \text{Graph} + \text{RPQ} \tag{34}$$

*Proof.* We can construct queries in the extended language that cannot be expressed in the more restricted languages. For example, path reachability queries in graphs cannot be expressed in standard relational algebra with a fixed number of joins.

**Theorem 8.2.2** (Complexity Separation). There exist graph queries whose evaluation is inherently exponential in the worst case unless P = NP.

*Proof.* Since subgraph isomorphism is NP-complete, and our framework includes it as a primitive operation, there exist queries whose evaluation is inherently exponential in the worst case unless P = NP.

# 9. Practical Implementation Considerations

## 9.1 Graph Indexing Structures

To support efficient graph operations, we propose specialized indexing structures:

1. **Vertex-centric indexes**: Map vertex properties to vertex IDs

2. **Edge-centric indexes**: Map edge labels and properties to edge IDs

3. **Path indexes**: Precompute common path patterns for faster traversal

4. **Subgraph indexes**: Index frequent substructures for faster pattern matching

## 9.2 Distributed Graph Processing

The framework can be extended to distributed settings:

1. **Vertex-cut partitioning**: Partition graphs by cutting vertices to minimize cross-partition edges

2. **Edge-cut partitioning**: Partition graphs by cutting edges to balance vertex distribution

3. **Pregel-style computation**: Implement operations using a vertex-centric "think like a vertex" model

4. **Bulk synchronous processing**: Process graph operations in synchronized supersteps

## 9.3 Incremental Graph Maintenance

For dynamic graphs that change over time:

1. **Delta operations**: Define incremental versions of graph operators that only process changes

2. **Version management**: Maintain graph versions for temporal queries

3. **Consistency models**: Define consistency guarantees for concurrent graph modifications

# 10. Conclusion and Future Directions

This extension of the unified mathematical framework incorporates graph data structures and operations into the existing algebraic foundation. We have shown that:

1. Graph operations can be represented within the existing Boolean algebra structure

2. Cross-paradigm operations between graphs, relations, text, and vectors are well-defined

3. Query optimization techniques can be extended to incorporate graph-specific rules

4. The category-theoretic foundation elegantly accommodates graph structures

Future research directions include:

1. Extending the framework to temporal graphs and streaming graph operations

2. Incorporating graph neural network operations into the algebraic structure

3. Developing specialized algorithms for cross-paradigm optimization involving graphs

4. Exploring the theoretical limits of expressivity in the unified query language

5. Investigating probabilistic graph models within the framework

This extension provides a rigorous foundation for building next-generation data systems that seamlessly integrate graph data with relational, textual, and vector data, enabling complex analytics across previously siloed paradigms.