

GETL: An Extract-Transform-Load Framework Across Graph Models in Graph Warehouse

Feng Yao, Xiaokang Yang, Shufeng Gong, Qian Tao, Yanfeng Zhang,
Wenyuan Yu, and Ge Yu, *Senior Member, IEEE*

Abstract—Various graph models have emerged to meet diverse application needs, each with unique characteristics and specialties. Managing and analyzing graph data inevitably requires interactions across different models to serve upstream business requirements. Therefore, an Extract-Transform-Load (ETL) tool designed to bridge different graph models is desired. In this paper, we propose GETL, a generalized graph ETL framework capable of automatically identifying graph model schemas and performing seamless data conversion among RDF, RDF-star, labeled property graph, and the relational model. This is attributed to GETL’s unified graph representation model, constructed as nested $\langle \text{label}, \text{entity} \rangle$ pairs, offering powerful capabilities in graph representation and model compatibility. Additionally, we develop a unified programming interface to support complex graph transformation tasks. It is built upon the Gremlin syntax and provides strong expressive capabilities. Finally, our evaluation demonstrates that GETL outperforms state-of-the-art solutions in terms of model conversion efficiency and data manipulation language (DML) intelligibility.

Index Terms—graph ETL framework, unified graph representation, programming interface.

I. INTRODUCTION

GRAPH-structured data has now drawn widespread attention in various fields. Several models exist, including the Resource Description Framework (RDF), Labeled Property Graph (LPG), and the Relational Model, each with unique characteristics suited to different graph task scenarios. For example, RDF is well-suited for semantic linking and reasoning [1], as well as data integration [2]. LPG is ideal for representing complex networks, excelling in graph traversal and analysis [3]. The relational model and its database systems is optimal for applications that require strict transaction control and well-structured data [4]. However, for recently emerging graph warehouses [5], [6], supporting diverse graph applications requires fusing and converting graph data across different models as a prerequisite for meaningful analysis.

Feng Yao, Xiaokang Yang, Shufeng Gong, and Ge Yu are with the School of Computer Science and Engineering, Northeastern University, Shenyang 110819, China (e-mail: yaofeng@stumail.neu.edu.cn; yangxk@stumail.neu.edu.cn; gongsf@mail.neu.edu.cn; yuge@mail.neu.edu.cn).

Qian Tao and Wenyuan Yu are with Alibaba Group, Hangzhou 311121, China (e-mail: qian.tao@alibaba-inc.com; wenyuan.ywy@alibaba-inc.com).

Yanfeng Zhang is with the School of Computer Science and Engineering, Northeastern University, Shenyang 110819, China, also with the National Frontiers Science Center for Industrial Intelligence and Systems Optimization, Northeastern University, Shenyang 110819, China, and also with the Key Laboratory of Data Analytics and Optimization for Smart Industry (Northeastern University), Ministry of Education, Shenyang 110819, China (e-mail: zhangyf@mail.neu.edu.cn).

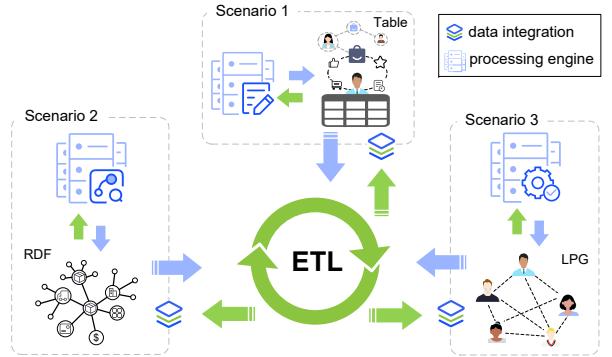


Fig. 1. A case of fusion and conversion involving various types of graph data across different business scenarios.

Example 1: Figure 1 illustrates how different graph models coexist and interact in large enterprises such as Alibaba. (1) In high-concurrency scenarios, e.g., Taobao [7] generates 4 billion behavioral records per day [8], relational tables are used to maintain the user–product relationship network with purchase and order attributes (Scenario 1). (2) In logistics and warehousing, RDF maintains a knowledge graph over products, inventory, and suppliers to integrate data from heterogeneous channels (Scenario 2). (3) Built on a large user base, e.g., 400 million daily active users on Taobao [9], LPG represents social networks in the ecosystem (Scenario 3). Each scenario is typically equipped with dedicated processing engines to perform specialized tasks, such as RDF engines for semantic queries and LPG engines for graph analytics [10]. To support cross-business analysis, enterprises need to fuse and convert data across these models. For example, integrating user purchase information from the relational tables in Scenario 1 into social LPG of Scenario 3 enriches user preferences and supports community discovery. Moreover, the same dataset may be modeled differently for different needs [11]. For instance, converting part of the RDF graph in Scenario 2 into LPG and reusing the engine in Scenario 3 enables detailed supply-chain flow analysis [12]. Such model fusion and conversion are essential to break data silos and promote data sharing and collaboration. □

From the above example, a natural requirement arises: graph data needs to be interconnected and converted between different models across systems. Additionally, the graph often requires a series of transformation operations (e.g., pruning of graph topology) before being loaded into the corresponding processing engines to meet the needs of the analysis task. Therefore, it is essential to develop a graph Extract-Transform-Load (ETL) tool that extracts graph data from source systems, seamlessly transforms it, and loads it into target systems, while

concealing the complexity of the underlying models. This is beneficial for practitioners, as they are inclined to focus on the graph area they specialize in without having to manually program cross-model conversion routines.

Currently, there are some efforts on graph model conversion and interoperability. One line of research [4], [13]–[19] focuses on constructing unidirectional or bidirectional mappings between any two of LPG, RDF, and the relational model. However, integrating all pairwise mappings in graph ETL requires incorporating a large number of conversion rules and mapping logic, which would be both complex and rigid. Direct mapping also hinders transformation operations in the ETL process. Another interesting proposal is to employ an intermediary to reconcile the disparities among models.

Unfortunately, the current mainstream graph models are not well-suited to serve as this intermediary. Each model’s design possesses unique characteristics while also revealing limitations in other aspects. Specifically, RDF represents data in the form of (*subject*, *predicate*, *object*) triples, while supporting global identifiers (IRIs) and a schema-less design [20], making it well-suited for efficient data integration. However, RDF lacks direct counterparts for higher-order relationships [21], *e.g.*, edge properties in LPG. Additionally, it does not support *multigraphs*, which allow multiple edges of the same type between two vertices [10]. Conversely, LPG is characterized as directed multigraphs, adopting the graph terminology to explicitly define labels and properties for vertices and edges. This makes LPG ideal for representing complex network structures and performing graph traversal. However, the flexibility of definitions on schema [19], [22] makes data integration challenging. The relational model can be used to represent graph data, where one table corresponds to a type of vertices in LPG, and edges are possible integrity constraints [23]. The structured organizational form makes various relationships transparent and normalized. While the relational model is valuable for well-structured data, its predefined schema limits the ability to accommodate dynamic changes in graph structures. In a nutshell, as an intermediary, the ability to accommodate various model representations and efficient data integration are goals pursued by the graph ETL. Clearly, none of them possess a fully compatible feature set.

On the other hand, Multilayer Graph [24] and Statement Graph [25] serve as intermediaries, exploring the representation of RDF (including RDF-star [26]) and LPG in *unified data model* manners. Multilayer Graph captures higher-order relationships by adding edge IDs to directed labeled edges for multilayer nesting. Statement Graph associates (*subject*, *predicate*, *object*) triples by setting internal nodes with statement identifiers and describes complex relationships through the connections between internal nodes. Both are designed with richer statements to support data representations from two graph models, but this in turn may require increased effort on data integration. Moreover, they cannot handle some of the more complex representational features specific to graph ETL, *e.g.*, paths as vertices, which enables users to manipulate graph data from different hierarchical views. The analysis above prompts us to design a unified data model to serve as an intermediary for graph ETL, which needs to accommodate

representations of three mainstream graph models and the specific complex representational features required by graph ETL, while also providing robust data integration capabilities.

Another key requirement for building a graph ETL tool is the design of a *domain-specific* programming interface to perform a wide range of complex transformations during the ETL process. This differs significantly from existing graph query languages, such as Gremlin [27] and Cypher [28], which are proficient in basic graph pattern matching and can serve as core components in graph ETL pipelines, but have limitations in addressing ETL-specific needs. Firstly, these languages struggle to express *special mapping operations* in graph ETL, such as mapping subgraphs to vertices and condensing edges between subgraphs into an edge. Secondly, they appear less elegant in handling *user-defined recursive traversal* operations that involve complex computational logic. For example, constructing weighted shortest paths in pathfinding operations requires complex programming or rigid predefined functions (*e.g.*, *Graph Data Science* library [29] in Neo4j and *VertexPrograms* [30] in TinkerPop).

In this paper, we develop GETL, a generalized graph ETL framework. Firstly, we propose a unified graph representation model. The model uses label-entity pairs to define all elements (*e.g.*, vertex, edge, property, and path) and employs nested binary relations to represent connections between entities at different levels, enhancing the graph representation capabilities. Furthermore, the overall consistent structural form emphasizes its data integration capabilities. Additionally, we define concise mapping rules that automatically identify the schemas of different models and enable seamless data conversion. Secondly, GETL features a three-layer architecture to accommodate different types of graph transformation operations, providing users with flexible and extensive operational space to customize the graph according to the requirements of upstream applications. Finally, we design a programming interface for GETL to facilitate user interactions with graph data within the three-layer architecture throughout the graph ETL process. It is built on Gremlin syntax to support complex mapping operations. Furthermore, by integrating Datalog’s recursive querying capabilities with Gremlin’s traversal query strengths, we introduce a hybrid Data Manipulation Language to optimize the needs of recursive aggregation operations.

Contributions. Overall, we make the following contributions:

- We develop GETL, a generalized graph ETL framework. To the best of our knowledge, it is the first graph ETL tool specifically designed to accommodate various graph models, serving graph warehouses and supporting broad types of upstream applications (§III).
- We propose a unified graph representation model as an intermediary for graph ETL, featuring robust representational power and strong data integration capabilities. (§IV-A).
- We design concise and intuitive mapping rules to enable the automatic identification of diverse graph model schemas and seamless data conversion. (§IV-B).
- We design unified programming interfaces tailored for the graph ETL process to meet users’ needs for extensive graph transformation tasks (§V).

II. PROBLEM FORMULATION AND MODEL ANALYSIS

In this section, we present the definition of the problem and analyze the characteristics of mainstream models for graphs.

A. Definition of Graph ETL Process

The ETL process involves extracting and transforming data from sources and then loading it into a target system. We define the *graph ETL process* as follows:

Extract. A set of graph sources is $S_G = \cup G_i$, where each G_i is a graph from a different system with a data model m_i (e.g., LPG, RDF, and the relational model). Given a unified data model m_* and a graph form d , the extract function $\epsilon : G_i, m_i \xrightarrow{m_*} d$ converts the graph G_i into a consolidated data form d relying on m_* .

Transform. By integrating over $\cup \epsilon_{m_*}(G_i, m_i)$, a unified graph representation, D_G , is constructed. Graph transformations based on D_G are then executed according to the requirements of upstream applications, categorized into the following types.

(1) *Graph pattern matching*, which essentially matches a set of *path bindings* in a graph [31]. A path $\rho = (v_0, e_0, v_1, e_1, \dots, e_{n-1}, v_n)$ in D_G represents a sequence of vertex v and edge e identifiers. The sequence is captured through pattern expressions, where each variable in the expression maps to a graph element (vertex, edge).

(2) *Special mapping*, which involves the abstraction of the graph at different layers. Specifically, the transformation aggregates different types of elements (vertices, edges, and *multidimensional properties*) to create an aggregate network [32]. For example, it can address questions such as, “*What is the network structure of a social graph when grouped by user gender?*” It makes sense for the user to customize graphs according to business needs. Additionally, consider a mapping function $f : G' \rightarrow v_{G'}$, which maps subgraph G' , constructed in different dimensional spaces on D_G , to a new vertex $v_{G'}$ in a new graph layer. Correspondingly, the edge set between subgraphs (e.g., $E(G'_1, G'_2)$) is condensed into a new edge, i.e., $f : E(G'_1, G'_2) \rightarrow e(v_{G'_1}, v_{G'_2})$. The f is sufficiently flexible to accommodate various mapping possibilities, such as mapping from subgraph to property and from path to edge.

(3) *User-defined recursive graph traversal*, which requires tailored recursive/iterative configurations and aggregation operations based on the recursive/iterative nature of graph constructs. Standard algorithms often fail to fully satisfy the needs of practical applications, necessitating a degree of customization by the user [33]. Examples include path planning, heuristic search algorithms, personalized graph ranking, weighted random walk, and GNN stochastic sampling.

It is worth noting that the three types of transformation operations mentioned above can function independently or be utilized in conjunction to achieve the desired outcomes.

Load. Generate graph data in a specified data model and load it into upstream graph management or analysis systems.

B. Different Data Models for Graphs

In search of characteristics of data models to lay the foundation for GETL compatibility, we introduce LPG, RDF, RDF-star, and the relational model.

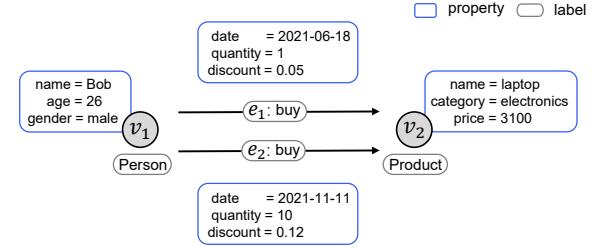


Fig. 2. An example of an LPG for online shopping.

Labeled Property Graph (LPG). We follow the definition from [10] and have included adaptations for data type descriptions. Assume that four infinite sets: \mathcal{L} (labels), \mathcal{K} (property names), \mathcal{V} (property values), and \mathcal{T} (data types).

Definition 1 (LPG). *An LPG is a tuple $(V, E, \pi, \lambda, \sigma)$ where:*

- (1) V is a finite set of vertices;
- (2) E is a finite set of edges such that $V \cap E = \emptyset$;
- (3) $\pi : E \rightarrow (V \times V)$ is a total function that associates each edge in E with a pair of vertices in V ;
- (4) $\lambda : (V \cup E) \rightarrow \mathcal{L}$ is total function that associates labels from \mathcal{L} with vertices/edges;
- (5) $\sigma : (V \cup E) \times \mathcal{K} \rightarrow (\mathcal{T}, \mathcal{V})$ is a partial function that associates property pairs (name, value) with vertices/edges.

Example 2: Figure 2 shows an example of an LPG for online shopping, featuring two types of vertices: v_1 , labeled as “Person”, and v_2 , labeled as “Product”. v_1 refers to a user named “Bob”, and v_2 is a “laptop” with a price of “3100”, both possessing their own properties. Edges e_1 and e_2 denote that “Bob” buys the “laptop” twice, with each edge’s properties capturing the respective order details. An example of mapping partial elements of LPG is as follows:

$V = v_1, v_2, \quad \lambda(v_1) = \text{Person}, \quad \lambda(v_2) = \text{Product},$
$E = e_1, e_2, \quad \lambda(e_1) = \text{buy}, \quad \pi(e_1) = (v_1, v_2),$
$\sigma(v_1, \text{name}) = (\text{String}, \text{Bob}), \quad \sigma(v_2, \text{category}) = (\text{String}, \text{electronics}),$
$\sigma(e_1, \text{date}) = (\text{DATETIME}, 2021 - 06 - 18).$

Notably, the definition of LPG schemas is relatively flexible and typically vendor-specific [10], which can lead to inconsistencies in integrating LPG data from different sources. □

RDF & RDF-star. The RDF model consists of *subject-predicate-object* triples. Assume that three pairwise disjoint sets: \mathcal{I} (IRIs), \mathcal{N} (literals), and \mathcal{B} (blank nodes).

Definition 2 (RDF). *An RDF triple is a tuple $t = (s, k, o)$ where subject $s \in (\mathcal{I} \cup \mathcal{B})$, predicate $k \in \mathcal{I}$, and object $o \in (\mathcal{I} \cup \mathcal{B} \cup \mathcal{N})$. An RDF graph is a finite subset of RDF triples.*

Example 3: The following provides an RDF example related to Example 2, serialized in the RDF Turtle syntax [34]. For brevity, we omit the namespace declarations.

$_v1 \text{ rdf:type } \text{ex:Person} ;$
$_v1 \text{ ex:name } \text{“Bob”} ;$
$_v1 \text{ ex:age } 26^{\text{xsd:integer}} .$
$_v2 \text{ rdf:type } \text{ex:Product} ;$
$_v2 \text{ ex:name } \text{“laptop”} ;$
$_v2 \text{ ex:category } \text{“electronics”} .$
$_v1 \text{ ex:buys } _v2 .$

Here, the blank nodes $_v_1$ and $_v_2$ represent the vertices v_1 and v_2 . The triple $\langle _v_1 \text{ ex:buys } _v_2 \rangle$ corresponds to an edge in the LPG, and XML Schema datatypes (*e.g.*, xsd:integer) are used to define the data types. \square

However, RDF does not provide built-in support for relationships involving more than two entities, *e.g.*, edge properties. RDF-star extends RDF with the *quoted triple*.

Definition 3 (RDF-star). An RDF-star triple is a tuple $t^* = (s^*, k, o^*)$ where subject $s^* \in \mathcal{I} \cup \mathcal{B} \cup \{t_1^*\}$, predicate $k \in \mathcal{I}$, and object $o^* \in \mathcal{I} \cup \mathcal{B} \cup \mathcal{N} \cup \{t_2^*\}$, with t_1^* and t_2^* being given RDF-star triples not equal to t^* . An RDF-star graph is a finite subset of RDF-star triples.

Here t_1^* and t_2^* are quoted triples, which are unique. By definition, an RDF graph is also an RDF-star graph.

Example 4: Building on Example 3, RDF-star represents edge properties with quoted triple as follows, serialized in Turtle-star format [35]:

```
<< \_v1 ex:buys \_v2 >> ex:date "2021-06-18"^^xsd:date;
<< \_v1 ex:buys \_v2 >> ex:quantity 1^^xsd:integer;
<< \_v1 ex:buys \_v2 >> ex:discount 0.05^^xsd:decimal.
```

Here, $\ll _v_1 \text{ ex:buys } _v_2 \gg$ is a quoted triple acting as the subject, with order information serving as the predicate and object to form RDF-star triples. \square

Regarding alignment with Example 2, RDF-star cannot effectively distinguish between edges e_1 and e_2 and their distinct properties. This is because they form the same triple, $\ll _v_1 \text{ ex:buys } _v_2 \gg$, but the *occurrences* they represent are different. For instance, when the properties of e_1 and e_2 are added to the triple, it is not clear which date corresponds to which discount.

Relational Model. The relational model organizes data into a set of tables, each representing a relation, which in turn is viewed as a collection of tuples.

Definition 4 (Relational Model). Let $R^* = (R_1, R_2, \dots, R_n)$ is a relational schema consisting of relations $R_i(A_1, A_2, \dots, A_m)$, where A_1, A_2, \dots, A_m are attributes. Each instance of relation R_i is an ordered n -tuple (a_1, a_2, \dots, a_m) , where each a_j is an element of attribute A_j .

Person				Product			
Sid	Name	Age	Gender	Did	Name	Category	Price
S001	Bob	26	male	D001	laptop	electronics	3100
...
Order							
Rid	Sid	Did	Date	Quantity	Discount		
R001	S001	D001	2021-06-18	1	0.05		
R002	S001	D001	2021-11-11	10	0.12		
...		

Fig. 3. The relational model corresponding to Example 2.

Example 5: Figure 3 illustrates the relational model representation corresponding to Example 2. There are three relational tables: Person, Product, and Order, each column representing an attribute, *e.g.*, “Name” in the Person table, with “Bob” being an element. Additionally, “Sid”, “Did”, and “Rid” are primary keys as unique identifiers. Instances of the Person and

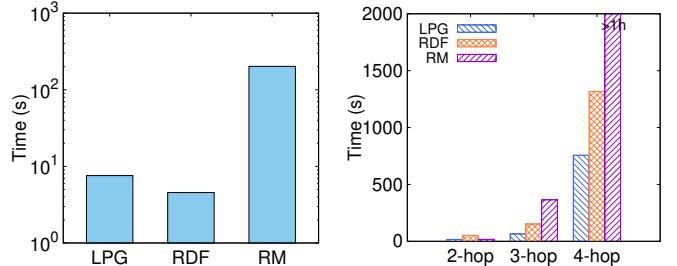


Fig. 4. Runtime of data integration. Fig. 5. Runtime of k -hop query.

Product tables correspond to vertices v_1 and v_2 in Example 2, respectively. Instances of the Order table correspond to edges labeled as “buy”, which are associated with the Person and Product tables through the foreign keys “Sid” and “Did”. \square

C. Model Design Philosophy and Insight

We discuss the uniqueness of these models from a design-philosophy perspective, and empirically compare LPG (TinkerPop [30]), RDF (RDF4J [36]), and RM (MySQL [37]) on a test graph with 20,000 vertices and 1,000,000 edges in terms of data integration and k -hop neighbor queries. These experiments provide insights into the design of the underlying unified graph model and reveal the need for flexible transformation operations in the graph ETL process.

The RDF focuses on resource identification and format standardization to ensure generality, simplicity, and scalability. Each RDF triple is an independent unit of information, so integrating new data only requires appending triples rather than performing complex schema matching, which is consistent with our results in Figure 4, where RDF achieves lower data integration time than LPG and RM. However, RDF struggles to represent higher-order relationships, as described in Example 4, and its triple-based representation leads to more joins during traversal. Conversely, LPG emphasizes the richness of topology and efficient traversal: as shown in Figure 5, LPG has the lowest latency for 2–4 hop neighbor queries, outperforming both RDF and RM. Additionally, the relational model employs unique identifiers (keys) to ensure entity integrity. However, in RM, edges must be materialized by constructing foreign-key relationships during integration, and additional table joins are required during query processing, making it more time-consuming for both data integration and complex traversal.

Insight. In the graph ETL process, the underlying unified data model m_* (refer to Section II-A) must be capable of integrating data from different models, aligning with RDF’s design philosophy for scalability. It employs a standardized format that does not interfere with existing data and structures. However, this alone is not sufficient; m_* must also cover the representation capabilities of all models, such as representing higher-order relationships in LPG. From this perspective, the design of unique identifiers (keys) in the relational model ensures data integrity, which is crucial for distinguishing and referencing relationships across different orders. On the other hand, LPG’s efficient graph traversal capabilities effectively facilitate graph transformation operations during the Transform stage. Additionally, its flexible schema definition offers users a richer space for graph semantics and topological manipulation on integrated graphs.

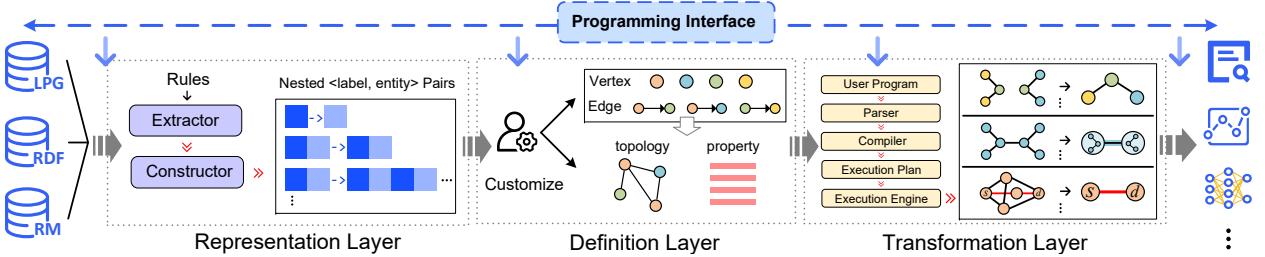


Fig. 6. The architecture overview of GETL.

III. SYSTEM OVERVIEW

Architecture. We build GETL, a generalized graph ETL framework. Figure 6 shows its system architecture, which consists of the following three layers:

Representation Layer. The representation layer provides a unified description of graphs from different sources across LPG, RDF, RDF-star, and the relational model (RM). In this layer, we design a unified graph representation model, constructed as nested $\langle \text{label}, \text{entity} \rangle$ pair, that acts as an intermediary for converting different data models, accommodating various graph representations. It employs a consistent representation form for different entities *without imposing a fixed schema*, facilitating data integration.

Definition Layer. The underlying model in the representation layer employs a schema-less design to enhance compatibility and integration. However, this flexibility conflicts with the structured and clearly defined syntax required in programming interfaces used for manipulating data. To bridge the gap between the two, we constructed the definition layer. This layer enables users to customize graph schemas using high-level language in the programming interface. Users can flexibly define different types of vertices, edges, and properties based on the underlying model, constructing a well-defined graph. It is worth noting that: 1) The graph defined based on the underlying model is not unique; it is customized by users, who can adapt the schema according to business needs. 2) The choice of property graph schema as the defining layer is due to its flexibility and superior graph traversal capabilities, which facilitate the execution of transformation operations. The design of this layer fulfills the programming interface's requirements for clear structuring, while simultaneously ensuring meaningful data semantics for the integrated graph.

Transformation Layer. Based on the user-customized graph in the definition layer, the transformation layer provides support for three categories of transformations. It features a built-in processing engine to execute transformation operations.

Workflow. We describe the workflow of GETL from the perspective of the graph ETL process.

Extract. Graphs from various sources based on LPG, RDF, and RM are automatically extracted into the representation layer. GETL establishes a set of rules, guiding the extractor in recognizing and extracting source graphs. Subsequently, the constructor unifies these into nested $\langle \text{label}, \text{entity} \rangle$ pairs.

Transform. After the extraction is complete, users can manipulate graph transformations through the programming interface.

First, customize the schema using the *Data Definition Language* (DDL) statements through the programming interface. For transformations, one approach involves executing simple transformations directly during the graph schema definition phase, such as applying constraints and conducting filtering operations. Alternatively, the defined graph schema can enter the transformation layer as a base graph layer. The difference between the two approaches is that the former is more efficient as it does not require the construction of multi-layered indexes, but it may struggle to support some complex operations. In contrast, the latter approach enables richer and more complex transformations. GETL empowers users by offering them the autonomy to make informed decisions tailored to their specific needs. At the transformation layer, users write transformation scripts through the programming interface. These scripts direct the transformations on the user-customized graph. Upon receiving these scripts, the transformation processing engine handles a series of processes, such as parsing and compiling.

Load. Finally, GETL converts the transformed graph into the graph form under the user-specified data model and loads it into the upstream system. Additionally, the output can serve as an intermediate graph and reintegrate into GETL.

IV. INTERMEDIARY DESIGN IN GETL

This section presents the design of the unified graph representation model (UGR) within GETL's representation layer.

A. Unified Graph Representation Model (UGR)

Design Idea. We adopt a schema-less design that treats all components of the graph equally as entities and describes the associative relationships expressed by the entities in this perspective. Thus, each entity is an independent unit and forms a consistent expression, which helps avoid complex matching or modifications when integrating new data, thereby facilitating integration. Additionally, we use nested structures to represent higher-order relationships and extend set representations to support specific features.

We first introduce a *label-entity pair* structure, which serves as the basic building block of the model. Let \mathcal{L} be a set of labels and \mathcal{E} be a set of entities representing concepts, things, or objects in the physical world within a given domain \mathcal{D} .

Definition 5 (*Label-Entity Pair*). A *label-entity pair* is an ordered pair (lb, et) where $lb \in \mathcal{L}$ and $et \in \mathcal{E}$. The pair (lb, et) forms an association in which lb categorizes, describes, or tags the entity et within the context of \mathcal{D} .

A label lb is a descriptor that can either be any element from \mathcal{L} or be *null*. The flexibility within \mathcal{L} allows for dynamic categorization of entities, accommodating various levels of granularity. An entity et can be identified in various formats, such as numbers, strings, tensors, and IRIs. \mathcal{D} refers to the definition of labels and entity relationships within a particular domain. It encompasses the rules, context, and knowledge system that influence how labels classify or describe entities. Furthermore, label-entity pairs support the description of a class of entities, specifically represented as $(lb, *)$, which denotes the set of all entities under the label lb .

To represent higher-order relationships among components in the graph, we introduce a nested binary relation based on label-entity pairs, termed *Nested <label, entity> Pair*.

Definition 6 (Nested <label, entity> Pair). *Let \mathcal{P} be a set in which each element $p \in \mathcal{P}$ is an ordered pair (p_1, p_2) . A nested binary relation \mathcal{R} on \mathcal{P} is defined recursively as follows:*

- (1) *A basic pair p in \mathcal{R} is a label-entity pair (lb, et) ;*
- (2) *A nested relation in \mathcal{R} can be an ordered pair $p = (p_1, p_2)$, where p_1 and p_2 are either basic pairs or themselves nested relations.*

Intuitively, everything (e.g., vertices, edges, properties, and even paths or an entire graph) can be identified as entities with descriptive labels. The content of entities can be detailed using nested binary relations. The nesting structure still follows the semantics of the label-entity pair form, even at varying levels of nesting. This approach effectively unifies representations of different graph components in a concise way.

Example 6: We continue referencing the example in Figure 2, illustrating its representation through nested label-entity pairs. A partial representation of the elements is provided below:

The basic pair p in \mathcal{R} :

$\text{Person} \rightarrow ps_id1$, $\text{Product} \rightarrow pd_id1$, $\text{buy} \rightarrow by_id1$, $\text{buy} \rightarrow by_id2$,
 $\text{String} \rightarrow \text{Bob}$, $\text{Integer} \rightarrow 26$, $\text{DATETIME} \rightarrow 2021 - 06 - 18$.

The nested relation pair p in \mathcal{R} :

#1 description of user and product information (vertex related)
 $(\text{name} : nm_id1) \rightarrow ((\text{Person} : ps_id1) \rightarrow (\text{String} : \text{Bob}))$,
 $(\text{age} : ag_id1) \rightarrow ((\text{Person} : ps_id1) \rightarrow (\text{Integer} : 26))$,
 $(\text{category} : cg_id1) \rightarrow ((\text{Product} : pd_id1) \rightarrow (\text{String} : \text{electronics}))$,
#2 description of purchase relationships and orders (edge related)
 $(\text{buy} : by_id1) \rightarrow ((\text{Person} : ps_id1) \rightarrow (\text{Product} : pd_id1))$,
 $(\text{date} : dt_id1) \rightarrow ((\text{buy} : by_id1) \rightarrow (\text{DATETIME} : 2021 - 06 - 18))$,
 $(\text{buy} : by_id2) \rightarrow ((\text{Person} : ps_id1) \rightarrow (\text{Product} : pd_id1))$,

Firstly, the model defines all components as entities, distinguishing them using unique identifiers, such as “ $\text{Person} \rightarrow ps_id1$ ”. In terms of expression, there is no difference among these components; all can be described using label-entity pairs. For constants, the data type and constant value also form the label-entity pair. Secondly, binary relations between entities are likewise regarded as entities and are identified using a labeled identifier, forming a nested relation pair, such as vertex properties, edges, and properties on edges. Formally, there is no structural difference between vertex-related entities and edge-related entities, despite their differing semantics and relational compositions. The design philosophy behind the

unified representation form for different entities is intended to facilitate effective data exchange and integration. \square

We further provide set operations on nested <label, entity> pairs to support the representational features of different levels of abstraction specific to graph ETL.

Definition 7 (Set Operations). *Given a nested binary relation \mathcal{R} and an ordered pair $p = (p_1, p_2)$ in \mathcal{P} , p_1 and p_2 can be expressed using set operations based on \mathcal{R} , including union (\cup), intersection (\cap), and difference (\setminus).*

Set operations that combine or modify relationships in the nested structure can support specific representational features of graph ETL. An example (*path as vertex*) is provided below:

The set operation on pair p in \mathcal{R} :

$(\text{specific_promotion} : sp_id1) \rightarrow ((\text{start_date} : sd_id1) \rightarrow (\text{DATETIME} : 2021 - 06 - 01)) \cup ((\text{duration} : dr_id1) \rightarrow (\text{Integer} : 18)) \cup ((\text{end_date} : ed_id1) \rightarrow (\text{DATETIME} : 2021 - 06 - 18))$.

We extend the example of the online shopping graph to describe promotion information with a two-hop path, which represents an 18-day promotion from the “ $\text{start_date} : sd_id1$ ” to the “ $\text{end_date} : ed_id1$ ”. This path entity can be treated as a vertex for establishing further associations, such as the entity is linked as a vertex to a promotional product category (“ $\text{category} : cg_id1$ ”) or a discount (“ $\text{discount} : dc_id1$ ”). Additionally, we can use $(lb, *)$ to represent a class of entities and use nested relations on it to represent graph entities, e.g., $(\text{graph} : pg_id1) \rightarrow ((\text{Person} : *) \rightarrow (\text{Product} : *))$, represents a user-product graph entity. Clearly, higher-order associations can be established on the abstraction layer of graph entities.

B. Automated Graph Model Conversion

GETL automatically identifies various graph model schemas and converts them using the concise and natural rules prescribed by UGR. Here, we provide a detailed discussion of the schemas and mapping rules for different models.

Labeled Property Graph. According to Definition 1, for the LPG tuple $(V, E, \pi, \lambda, \sigma)$, there exists the following mapping Φ_{LPG} between UGR and LPG:

- (R1) $\lambda(V \cup E) \rightarrow \text{pair } (lb, et)$, where $lb \in \mathcal{L}$ and $et \in (V \cup E)$.
- (R2) $\pi(E) \rightarrow \text{pair } p = (p_1, p_2)$, where p_1 is a label-entity pair (lb, et) , with $p_1 \in (\mathcal{L}, E)$, and p_2 is a nested label-entity pair (p_1, p_2) , with $p_1, p_2 \in (\mathcal{L}, V)$.
- (R3) $\sigma((V \cup E), \mathcal{K}) \rightarrow \text{pair } p = (p_1, p_2)$, where p_1 is a label-entity pair (lb, et) , with $p_1 \in (\mathcal{K}, IDs)$, and p_2 is a nested label-entity pair (p_1, p_2) , with $p_1 \in (\mathcal{L}, (V \cup E))$ and $p_2 \in (\mathcal{T}, V)$.

As demonstrated in Example 6, all entities are identified by labeled unique identifiers, forming label-entity pairs. It is worth noting that these unique identifiers can distinguish *multiple edges* of the same type between vertices. The mapping rule R1 specifies the labeled vertex and edge entities themselves. In R2, p_1 specifies the edge entity, while p_2 represents the binary relation between the two vertex entities that constitute the edge. R3 is a mapping of properties, where p_1 designates the

property entity, and p_2 establishes the binary relation between the vertex or edge entity and its corresponding property value.

RDF & RDF-star. Based on Definitions 2 and 3, the mapping Φ_{RDF} between UGR and RDF & RDF-star is as follows:

- (R1) Given a tuple $t = (s, k, o)$ with $k = \text{rdf : type}$, then $t \rightarrow$ pair $p^s = (o, s)$, where p^s is a label-entity pair on entity s , and $s, o \in (\mathcal{I} \cup \mathcal{B})$.
- (R2) Given a tuple $t = (s, k, o)$ with $k \in \{\mathcal{I} \setminus \{\text{rdf : type}\}\}$, then $t \rightarrow$ pair $p = (p_1, p_2)$, where p_1 is a label-entity pair (lb, et) , with $p_1 = (k, id)$, and p_2 is a nested label-entity pair (p'_1, p'_2) , with $p'_1 \in p^s \cup (null, s)$ and $p'_2 \in p^o \cup (null, o)$.
- (R3) Given a tuple $t^* = (s^*, k, o^*)$, where $s^* \in \mathcal{I} \cup \mathcal{B} \cup \{t_1^*\}$, $k \in \mathcal{I}$, and $o^* \in \mathcal{I} \cup \mathcal{B} \cup \mathcal{N} \cup \{t_2^*\}$, with t as specified in cases (R1) and (R2). There exists $\Phi'_{RDF} : t_1^*, t_2^* \rightarrow p'$, where $p' = p^s$ in case (R1) or $p' = p_1$ in case (R2). Then, $t^* \rightarrow$ pair p^* , where p^* is a nested binary relation pair that embeds p' into the corresponding position within the mapping results of case (R1) or case (R2).

The RDF syntax does not explicitly specify labels, but objects linked by the predicate “rdf : type” can semantically be interpreted as labels. Correspondingly, for R1, triples with the predicate “rdf : type”, whose subject s serves as the entity identifier, and the object o serves as the label, form label-entity pairs. Here, p^s is a label-entity pair denoting a labeled entity with the entity identifier positioned at the subject s . When the predicate k is not “rdf : type”, the triple in R2 is interpreted as a binary relation between the subject s and the object o concerning the predicate k . Therefore, p_1 specifies the entity k and assigns a unique identifier $id \in IDs$, while p_2 represents the binary relation between the entities s and o . In R3, RDF-star allows triples to be contained within the subject s and object o . Consequently, according to RDF mapping rules R1 and R2, the quoted triple t is mapped first, followed by the mapping of the outer triple. Specifically, the quoted triple t is initially mapped to p^s in R1 or to p_1 in R2. Subsequently, t is replaced at the corresponding positions in the mapping functions of R1 or R2. For example, in an RDF-star triple $t^* = ((s_1, k_1, o_1), k, o)$, the triple (s_1, k_1, o_1) is identified by R1 as $p^{s_1} = (o_1, s_1)$, and then t^* is mapped to $p^* = (k, id) \rightarrow ((o_1, s_1) \rightarrow p^o)$ after R2 quoting p^{s_1} as subject. Note that R3 describes only first-order RDF-star nesting; multi-layer nesting can be handled by a layer-by-layer mapping decomposition.

Relational Model. Following Definition 4, the mapping Φ_{RM} between UGR and RM can be expressed as follows:

- (R1) Given a relational table $R_i(A_1, A_2, \dots, A_m)$, where attributes A_j and A_{j+1} are set as foreign keys, then $R_i(A_j, A_{j+1}) \rightarrow$ set of pairs $P = \{p | p = (p_1, p_2)\}$, where p_1 is a label-entity pair, $p_1 \in (R_i, PK)$, and p_2 is a nested label-entity pair (p'_1, p'_2) , with $p'_1 \in (A_j, a_j)$ and $p'_2 \in (A_{j+1}, a_{j+1})$.
- (R2) Given a relational table $R_i(A_1, A_2, \dots, A_m)$ excluding the attributes A_j and A_{j+1} , then, $R_i(A_1, A_2, \dots, A_m) \rightarrow$ set of pairs $P = \{p | p = (p_1, p_2)\}$, where p_1 is a label-entity

TABLE I
SUMMARY OF VARIOUS MODEL FEATURES

	LPG	RDF	RDF-star	RM	MG	SG	UGR
Vertex Label	✓			✓		✓	✓
Edge Label	✓	✓	✓	✓	✓	✓	✓
Vertex Property	✓	✓	✓	✓	✓	✓	✓
Edge Property	✓		✓	✓	✓	✓	✓
Multiple Edges	✓			✓	✓	✓	✓
Edge as Vertex				✓	✓	✓	✓
Path as Vertex							✓
Path as Edge							✓
Graph as Vertex							✓

pair, $p_1 \in (A, IDs)$, and p_2 is a nested label-entity pair (p'_1, p'_2) , with $p'_1 \in (R_i, PK)$ and $p'_2 \in (\mathcal{T}, a)$.

The mapping in rule R1 specifies the binary relation between entities, denoted by foreign keys (*i.e.*, p_2), and uniquely identifies the relation using the table name R_i and primary keys PK (*i.e.*, p_1). In R2, non-foreign key attributes are specified as entities by combining attribute names A and unique identifiers IDs in p_1 . p_2 establishes a binary relation between instances identified by the primary key PK of R_i and attribute values a with data types \mathcal{T} . Rule R1 captures relational tables whose schema already encodes a binary relation via foreign keys. For purely tabular RM input, the data are first extracted into UGR using R2, and their graph structure is then explicitly defined either via the GETL API `addRMSchema(RMSchemarmSchema)` (as shown in Table II) or via DDL in the definition layer.

Note that while the unified graph representation model can accommodate various representation features from different graph data models, it cannot ensure lossless conversion between the input and output models. This limitation arises from the proprietary characteristics of these models. For example, although the unified graph representation model can represent graph data with multiple edges, such structures cannot be losslessly encoded back into RDF.

C. Comparison of Representational Features between Models

UGR provides support for a rich set of graph representation features. In Table I, we compare the representational capabilities of various models [24], including LPG, RDF, RDF-star, relational model (RM), Multilayer Graphs (MG) [24], Statement Graphs (SG) [25], and UGR used in GETL. Additionally, we provide the following introduction to the features specifically required for graph ETL.

- *Edge/Path/Graph as Vertex*: Edges, specific paths, and entire graphs can be referenced as vertices.
- *Path as Edge*: Specific paths can be referenced as edges.

We emphasize that these higher-order features are essential for graph ETL, as they are built from different abstraction perspectives and enable users to customize graphs at multiple layers without reconstructing complex structures from scratch. For example, *Edge as Vertex* allows edges to connect to other vertices in the graph [24], while *Path as Edge* supports assigning properties to specific paths. In UGR, Definition 7 is an intrinsic part of the model semantics: it introduces set operations over nested structures and treats higher-order constructs (paths/graphs) as first-class set-valued entities $(lb, *)$, so that

their compositional semantics are defined algebraically and higher-level entities (*e.g.*, paths or graphs) can be constructed from their constituent relations. This makes abstractions such as *Path/Graph as Vertex* native constructions in UGR, whereas an enhanced SG or MG would typically require additional reification or extra layers to emulate similar abstractions.

V. PROGRAMMING INTERFACE FOR GETL

This section presents a high-level overview of the programming interface of GETL, which extends the basic components of Gremlin [27]. The interface supports graph schema definition and transformation manipulation.

A. Graph Schema Definition

The definition layer allows users to customize graph schemas and perform simple data filtering. To facilitate this, we design a *Data Definition Language* (DDL) as part of GETL’s programming interface, utilizing a syntax similar to Gremlin. We illustrate the DDL syntax through Figure 7, which defines the underlying data model described in Example 6 as a shopping graph. First, register a graph named “shopping_graph” using `g.register`. The `Config()` specifies the nested label-entity pairs in the model associated with this graph, such as “{ ‘Person’: * }”, representing all entities associated with the “Person” label. Second, define vertices with `putV()` and indicate the labels of entities set as vertices (*e.g.*, “Person”). Define edges with `putE()`, specifying the edge label (*e.g.*, “buy”), direction `Direction.OUT`, and endpoints. Follow these with `property()` to assign properties and `label()` to reset labels if necessary. Third, DDL supports pruning the newly registered graph using `where()`, filtering structures that meet specific conditions, such as vertices where “age” is greater than 18. Finally, `setIdTransform()` finalizes the user-registered graph by constructing a mapping from UGR entity identifiers to internal graph-element IDs, which are then used for traversal and indexing in subsequent transformations.

```
# define a shopping graph configured by three types of label-entity sets
shopping_graph = g.register.Config({‘Person’: *}, {‘Product’: *}, {‘buy’: *})
    # define vertex with the ‘Person’ label and add properties
    .putV(‘Person’)
        .property(‘name’).property(‘age’).property(‘gender’)
        .where(values(‘age’).is(gt(18)))
    # define vertex with the ‘Product’ label and add properties
    .putV(‘Product’)
        .property(‘name’).property(‘category’).property(‘price’)
    # define edge labeled ‘buy’ from ‘Person’ to ‘Product’
    .putE(‘buy’, Direction.OUT, ‘Person’, ‘Product’)
        .property(‘date’).property(‘quantity’).property(‘discount’)
        .where(values(‘date’).is(lt(‘2021-06-20’)))
    .setIdTransform()
```

Fig. 7. Example DDL for defining a shopping graph.

DDL also provides graph schema definitions at various abstraction levels, through different model representational features, such as *Edge/Path/Graph as Vertex* and *Path as Edge*.

B. Graph Transformation Manipulation

GETL defines a set of syntactic rules to describe three types of transformation operations. We demonstrate the programming syntax features and commands for each type.

Basic Syntax. We utilize Gremlin to perform basic transformation operations, *i.e.*, *graph pattern matching*. The Gremlin syntax supports high-level and declarative programming, enabling users to define traversal instructions easily. The traversal source, `g`, serves as an interface for accessing data in the graph, using the operators `V()` and `E()` to select vertices and edges. Moreover, Gremlin provides a comprehensive set of operators for filtering (`has`, `where`), projection (`select`, `by`), looping (`repeat`), aggregation (`group`), sorting (`order`, `limit`), and more. Further details are available in [38].

```
#Q1: query “Person” who purchased “ProductA”
g.V().has(‘Product’, ‘name’, ‘productA’).in(‘buy’).dedup()
#Q2: query “Person” who purchased both “productA” and “productB”
g.V().match(
    __.as(‘person’).out(‘buy’).has(‘Product’, ‘name’, ‘productA’),
    __.as(‘person’).out(‘buy’).has(‘Product’, ‘name’, ‘productB’))
    .select(‘person’).values(‘name’).dedup()
```

Fig. 8. Gremlin for graph pattern matching.

Figure 8 demonstrates using Gremlin for graph pattern matching. Query Q1 searches for individuals who purchased “productA”. Here, `has()` filters and locates vertices labeled “productA”. The `in()` traverses incoming edges labeled “buy” to find vertices connected to “ProductA” vertices. The `dedup()` operator removes duplicates. Query Q2 uses `match()` for declarative pattern matching, containing two clauses that both start with the alias “person” and confirm whether the person purchased specific products. Finally, `select()` extracts the “name” from matched “Person”.

Layer Building. GETL introduces the `layer()` operator to construct new layers that support *special mapping* and rewrites existing Gremlin steps to support its execution semantics. We illustrate this enhancement with the following two use cases.

Figure 9 shows an example of extracting subgraphs from the “Person” graph based on “age” ranges and mapping these to “Cohort” vertices. The process includes several steps: Initially, `layer()` creates a new graph layer, grouping “Person” vertices by “age” intervals of 10 to form vertices labeled “Cohort”. The `groupBy()` enables grouping operations using embedded expressions. Subsequently, vertex properties are defined to specify the maximum and minimum ages within the “Cohort”. The `unfold()` is then used to expand the “Cohort” vertices and extract their “age” properties. Edges are constructed between “Cohort” vertices, incorporating the original edges between “Person” vertices, with an edge property labeled as “weight” that counts all original edges. Finally, the `finish()` completes all the layer-building operations.

```
# define the graph’s abstraction at different layers
cohort_graph = layerGraphBuilder(person_graph)
    .layer(g.V().hasLabel(‘Person’))
        .label(‘Cohort’).as(‘cohort’)
    # group ‘Person’ by ‘age’ range interval of 10
    .groupBy(‘age’, v -> Math.floor(v.get().value(‘age’) / 10) * 10)
    # add maximum and minimum age as vertex property
    .property(‘maxAge’, select(‘cohort’).unfold().values(‘age’).max())
    .property(‘minAge’, select(‘cohort’).unfold().values(‘age’).min())
    .addE(‘connection’, Direction.OUT, ‘cohort’)
    # add count of original edges as edge property
    .property(‘weight’, select(‘cohort’).bothE().unfold().count())
    .setIdTransform().finish()
```

Fig. 9. Define the graph’s abstraction at different layers.

```

path_graph = layerGraphBuilder(person_graph)
    .layer(g.V().hasLabel('Person'))
    .label('Person')
    # map 3-hop paths to edges
    .addE(g.V().as('start')
        .repeat(out().simplePath().times(3).as('end'))
        .path(), 'connection', Direction.OUT, 'start', 'end')
    .setIdTransform().finish()

```

Fig. 10. Map three-hop paths to edges.

Figure 10 provides another example, where three-hop paths are mapped directly as edges. The `addE()` accepts a `path()` parameter, utilizing the “start” and “end” vertices of the path as the two endpoints of the new edge.

User-defined Recursive Graph Traversal. User-defined recursive graph traversals are an essential tool for complex transformation operations, yet current approaches lack adequate support. Figure 11 illustrates a standard method, `ShortestPathVertexProgram`, provided by TinkerPop [38] for solving the single-source shortest path problem (SSSP), which includes several configurable parameters. However, this method is overly rigid, thereby restricting users’ operational flexibility. In addition, as depicted in Figure 12, implementing SSSP (*i.e.*, vertex-centric iterative algorithms) using Gremlin scripts requires maintaining additional vertex states, along with complex iteration control and arithmetic expressions, significantly increasing the program’s complexity.

```

# create a shortest path vertex program and configure it
spvp = ShortestPathVertexProgram.build()
    .source(sourceId).distanceProperty('weight').create()
result = graph.compute().program(svpv).submit().get()
# output the shortest path from the results memory
result.memory().get(ShortestPathVertexProgram.SHORTEST_PATHS)

```

Fig. 11. `ShortestPathVertexProgram` for SSSP.

```

# define the source vertex and initialize distances
def sourceId = 0
g.V().property('distance',
    __.choose(_.id().is(sourceId), 0, Double.MAX_VALUE)
    .iterate())
def vertexCount = g.V().count().next()
# |V| - 1 iterations to find the shortest path
g.V().repeat(
    __.as('v').outE().as('e').inV().as('w')
    .select('e', 'w').by('weight').as('b').by('distance').as('c')
    .by(select('v')).values('distance').as('a')
    .where(__.math('a + b').is(lt('c')))
    .select('w').property('distance', __.math('a + b')))
    .times(vertexCount - 1).iterate()
)

```

Fig. 12. An implementation of SSSP via a Gremlin script.

To address the above issues, we develop a hybrid Data Manipulation Language (DML) that combines Datalog and Gremlin. We choose the Datalog + Gremlin combination based on the following observations. First, Datalog and Gremlin have been widely used in knowledge-oriented applications [39] and graph applications [33] domains, making them easy for new users to learn. Second, Gremlin can easily express complex navigational graph queries, excelling in path filtering and control, but handling recursive/iterative logic is challenging, which is where Datalog excels. Datalog uses concise rule definitions for recursion and possesses logical reasoning capabilities. Moreover, its rule definition approach is well-suited for managing intermediate states and arithmetic expressions in iterative computations. Thus, combining Gremlin and Datalog leverages the strengths of both. The syntax is as follows:

```

FROM GRAPH example_graph
...
R(x1, ..., xn, [Aggregate(z)]) :- B1(x1, ..., xn, z);
...
:- Bm(x1, ..., xn, z).
-----
where  $\forall 1 \leq i \leq m B_i(x_1, \dots, x_n, z)$  : R(x1, ..., xn, y),
<Gremlin-expression>, P1(x1, ..., xn, c), ...
Pl(x1, ..., xn, c), z = Function(y, c),
[Condition(z)];

```

The program begins with the `FROM GRAPH` clause, specifying the input graph view. The core of the program consists of a finite set of *rules*. Within these rules, the predicate `R` on the left side of “`:`” serves as the rule head. `Aggregate` specifies an aggregation function applied to the arguments `z`. The right side of “`:`” contains multiple rule bodies `B`, separated by semicolons, where each `B` can include multiple predicates separated by commas. In the rule body `B`, the predicate `R` with the same name as the head predicate forms a recursive rule. The `Gremlin-expression` guides recursive traversal on the graph view, while the predicate `P` provides parameters for the computation of `z`-Function. In addition, some programs include rules with conditions for terminating recursion, *e.g.*, `[sum[Δrank] < 0.001]` in delta-based PageRank [40].

Figure 13 presents an example program using hybrid DML for SSSP, consisting of two rules. The first rule identifies the source vertex and sets the initial distance. The second rule recursively calculates the distance `dis_y` from the source to `y`, based on the path length `dis_x` from vertex `x` and the edge weight between `x` and `y`. Gremlin expressions guide edge filtering and path matching throughout this process. The aggregation function `min` determines the current shortest path to `y`. Compared to the previously mentioned Gremlin implementation of SSSP, the hybrid DML requires only two SSSP rules but enhances the flexibility of algorithm design.

```

FROM GRAPH road_graph
SSSP(x, dis) :- x = 0, dis = 0.
SSSP(y, min[dis_y]) :- SSSP(x, dis_x),
g.V(x).outE().has('weight').as('e')
.inV().hasLabel('location').as('y'),
dis_y = dis_x + e.weight.

```

Fig. 13. Example SSSP via hybrid DML.

```

# weighted random walk with a maximum step count of i + 1
FROM GRAPH social_graph
degree(x, count[y]) :- g.V(x).bothE().as('y').
WRW(0, x) :- x = 0.
WRW(i + 1, random[y, w_y]) :- WRW(i, x),
g.V(x).outE().has('weight').as('e')
.inV().hasLabel('Person').as('y'),
degree(y, d),
w_y = 0.35 * e.weight + 0.65 * d.

```

Fig. 14. Example weighted random walk via hybrid DML.

Another example, shown in Figure 14, involves a user-defined weighted random walk algorithm. This algorithm adjusts the transition probabilities of the random walk based on the user-defined importance of each vertex, aiming to identify interest groups within social networks. Specifically, the algorithm incorporates both vertex degree and edge weight to determine the importance `w_y` of vertex `y`. Then, `w_y` is

TABLE II
GETL PROGRAMMING INTERFACES.

Function Name	Description
refreshLPG() / refreshRM() / refreshRDF()	Convert UGR to LPG / RM / RDF.
loadLPG(LPGGraph lpgGraph) / loadRM(RMGraph rmGraph) / loadRDF(RDFGraph rdfGraph)	Load LPG / RM / RDF into UGR.
addRMSchema(RMSchema rmSchema)	Add the RM schema for mapping rules with UGR.
addLPGMappingConfig(LPGConfig lpgConfig)	Add user-customized graph configuration.
readRDFFile(RdfFormat rdfFormat, String filePath)	
queryFromMySQL(MysqlSessions sessions)	Read RDF files / MySQL to RMGraph / LPG
loadVertex(String fileName, String vertexLabel, String...properties)	vertex and edge (including Property) files.
loadEdge(String filePath, String edgeLabel, String from, String to, String...properties)	
Modified or Extended Gremlin Steps	Enhanced Gremlin Functions
register(), Config(), putV(), putE(), property(), label(), setIdTransform()	Define the graph schema.
layer(), groupBy(), unfold(), addE(), drop(), path(), groupCount(), count(), materialize(), finish()	Build a new graph layer.

used as a parameter to adjust the random selection process for determining the next vertex in the walk.

VI. SYSTEM IMPLEMENTATION

We propose a graph ETL framework, GETL. Here, we discuss the details of some components.

Interface. GETL’s data interface supports multiple input and output data formats for graph models, such as CSV, JSON, XML, and TTL. Additionally, GETL facilitates direct integration with data management and analysis systems, such as GraphScope [33], MySQL [37], Neo4j [29], Amazon Neptune [41], and TinkerPop [30]. GETL also provides programming interfaces for the entire ETL process, with the main components listed in Table II. These interfaces include system-level APIs and extended Gremlin steps based on Gremlin v3.4.9 [27]. Furthermore, GETL leverages existing entity alignment approaches [42], [43] to address data consistency.

Indexes. GETL does not allocate new storage for user-customized graph schemas but instead constructs indexes based on underlying nested pairs. The advantage lies in a more streamlined and efficient construction of graph schemas, which also reduces memory consumption. Additionally, since the definition layer supports only simple transformations, using indexes does not significantly impact transformation performance. In addition, GETL provides the materialize() interface, which allows intermediate graph layers to physically materialize the constructed temporary graphs.

Garbage Collection (GC). Constructing different graph layers typically incurs significant storage overhead. Therefore, as an in-memory framework, effective GC is crucial for GETL to prevent excessive memory consumption during continuous operations. GETL performs data cleanup in stages. After constructing a new graph layer, GETL traverses the previous layer to identify unnecessary data objects based on the connections between the current and previous layers. The GC thread then releases these objects. Furthermore, GETL also supports selective retention of data from any layer to serve data integration and reuse.

Hybrid DML Execution Engine. Although combining DataLog with Gremlin increases expressiveness for recursive ETL transformations, naive evaluation may generate large intermediate results. GETL mitigates this issue by compiling the

hybrid program into a frontier-based execution plan and executing recursion via semi-naive evaluation: in each iteration, the Gremlin expression is evaluated only from the newly derived tuples ΔR_t of the recursive predicate, which restricts traversal to the current frontier, enables early pruning via Gremlin filters and path constraints, and avoids full-graph scans. Moreover, the rule head supports aggregation (e.g., min for SSSP and sum for delta-based PageRank [40]) that can be incrementally maintained to deduplicate derivations and bound intermediate growth.

VII. EXPERIMENTAL EVALUATION

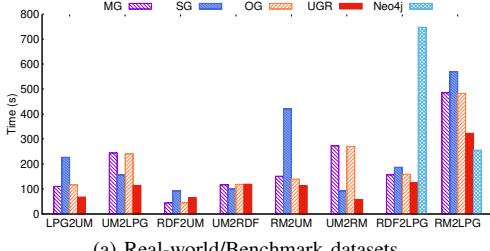
A. Experimental Setup

Environments. All the experiments in this section are conducted on an AliCloud ECS server (ecs.c7.16xlarge instance) equipped with 64 vCPU cores and 512GB of memory.

Datasets. We evaluate GETL on real-world and benchmark datasets covering RM, RDF, and LPG. For the relational model, we use MovieLens [44] with 18 columns and 41,741,454 rows. For RDF, we use Wikipedia [45] with 22,791,171 triples. For LPG, we use the LDBC SNB benchmark [46], which contains 3,966,203 vertices and 23,031,794 edges, along with 19,736,348 vertex properties and 6,854,988 edge properties. To validate data integration across models, we further derive three sub-datasets from LDBC: an LPG G_1 that stores edges with properties, an RDF graph G_2 that stores edges without properties and a subset of vertex properties, and a relational dataset G_3 that stores the remaining vertex properties; specifically, G_1 has 3,966,203 vertices, 19,740,638 edges, and 6,854,988 edge properties, G_2 has 14,883,937 triples, and G_3 has 23 columns and 3,948,592 rows.

We also construct three synthetic scalability datasets to evaluate the scalability of model conversion in GETL: an LPG dataset with 36,802,338 vertices and 231,371,359 edges (with 198,880,051 vertex properties and 77,985,155 edge properties), an RDF dataset with 100 million triples, and a relational dataset with 18 columns and 208,707,270 rows.

Competitors. In the conversion performance comparison of the unified model (UM), we use Multilayer Graphs (MG), Statement Graphs (SG), and OneGraph [47] (OG) as competitors. All three models support LPG, RDF, and RDF-star representations. To align the evaluation, we implement routines that convert between these models and RM. We choose GraphScope [33] and Ra-SQL [48] for comparison, both of which



(a) Real-world/Benchmark datasets.

Fig. 15. Runtime of model conversion.

have designed DMLs capable of handling recursive algorithms based on traversal queries. GraphScope extends the “Scatter-Gather” step on Gremlin to support vertex-centric computation logic. RaSQL supports recursive aggregation through an extension of the SQL standard.

B. Model Conversion Performance

We first evaluate the performance of MG, SG, OG, and UGR in bidirectional mapping between LPG, RDF (including RDF-star), and RM, and additionally report Neo4j as a specialized LPG baseline for RDF2LPG and RM2LPG. As shown in Figure 15, the mapping is divided into two directions: from graph models to UM, *e.g.*, LPG2UM, and from the unified model to specific graph models, *e.g.*, UM2LPG. Across both the real-world/benchmark and synthetic datasets, UGR consistently achieves the lowest runtime for LPG2UM and UM2LPG, as well as for RM2UM and UM2RM. This benefit mainly comes from the simplicity of UGR’s mapping rules and its direct bidirectional support for RM. In contrast, MG and OG are faster on RDF2UM and UM2RDF, since their relationship-identity mechanisms can be naturally aligned with RDF triples. We also evaluate conversions between LPG and RDF via the unified models, *i.e.*, RDF2LPG and RM2LPG. Across both datasets, UGR achieves the best performance among the unified models for RDF2LPG and RM2LPG. Neo4j, as a specialized LPG system, can directly convert from MySQL to Neo4j without going through an intermediate unified model. Therefore, in RM2LPG conversion, Neo4j is faster.

C. Graph ETL Execution Performance

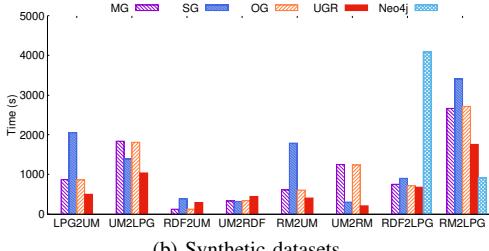
We introduce seven queries to describe the graph ETL process under diverse requirements and evaluate the overall performance of GETL in executing these queries. These queries involve input and output formats across different models and various types of transformation operations, encompassing the core functionalities of GETL and hold practical relevance.

Q1: Directly convert MovieLens from RM to an RDF graph.

Q2: In the LDBC graph, execute a custom graph schema operation to identify vertices representing individuals who posted on forums with more than 20 members from 2010-01 to 2010-05, and retrieve the received comments. Output the results in RM.

Q3: In Wikipedia, construct graph schema based on IRIs, filter vertices with a degree greater than 20 along with connected edges, and convert these results to LPG.

Q4: In the LDBC graph, identify the paths: personA → comment → post → personB, and map the path to edge labeled “isFanOf” from personA to personB. Then, output in RDF.



(b) Synthetic datasets.

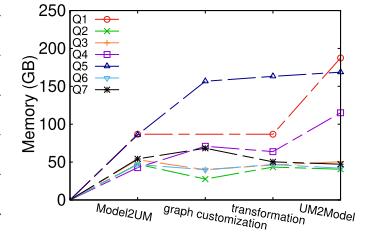


Fig. 16. Memory cost by stage.

Q5: In MovieLens, categorize movies into different subgraphs based on the “genome-scores.relevance” property and map each subgraph to a vertex. Construct edges between these subgraphs by identifying users who rate movies in both subgraphs with a score of 5, labeling these edges as “relevance” to reflect the relevance of user preferences. Additionally, count the occurrences of such rating patterns to use as an edge property, representing the weight of movie category relevance. Finally, export the graph layer as LPG.

Q6: In Wikipedia, similar to query Q3, extract vertices with a degree greater than 20, along with their connecting edges, to construct an undirected graph. Apply the Connected Component algorithm [49] to identify subgraphs. Finally, output these connected subgraphs as relational tables.

Q7: Integrating G_1 , G_2 , and G_3 from different models enriches the properties on vertices and edges. Subsequently, for community detection among user groups, conduct a biased two-hop random walk on the integrated graph to sample paths with edge labels “likes” or “hasCreator”: person → post → person. Calculate transition probabilities based on tag similarity (intersection of person and post tags) and the number of comments and likes on the post, which influence the selection of vertices at each hop. The sampled paths are then mapped to new edges, *i.e.*, person → person, with “common_interest” labels, and the results are converted into RM.

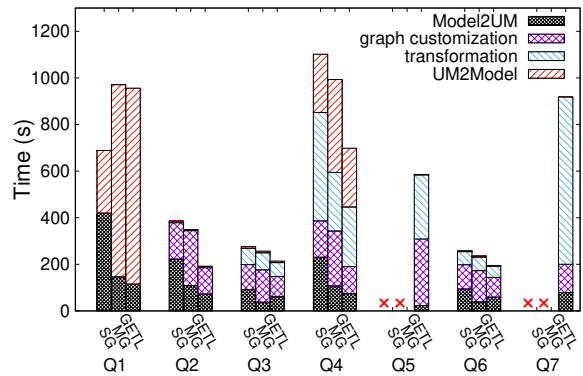


Fig. 17. Runtime breakdown.

We evaluate both the overall runtime and the per-stage runtime for executing seven queries in GETL. The total runtime is decomposed into four stages: extraction (*Model2UM*), graph schema customization (*graph customization*), transformation operations (*transformation*), and output model conversion (*UM2Model*). For MG and SG, we use TinkerPop to perform graph customization and transformation. Consequently, their breakdown includes not only the extra conversion overhead but also the cost of running ETL logic inside a general-purpose property-graph engine. Moreover, TinkerPop does not

support the layered graph operations required by Q5 and Q7, so these two queries are unsupported for MG and SG and are marked with red crosses. The results are shown in Figure 17. Across all supported queries, UGR consistently achieves the lowest overall runtime, outperforming MG and SG. In Q1, the mapping from UM to RDF dominates the runtime because it needs to generate a large number of RDF triples. For Q4 and Q5, the *transformation* stage involves searching for path and subgraph patterns and constructing new graph layers, which leads to higher transformation costs. The overhead in Q7 primarily arises from the serial integration of three independent datasets and the execution of complex multi-layer transformation operations. In addition, GETL outperforms the competitors in both graph schema customization and transformation, thanks to its tight integration of ETL functionalities and its efficient execution engine.

D. Memory Cost

We further evaluate the memory cost. GETL performs a phased garbage collection (GC) strategy to optimize memory usage. Figure VII reports the trend of memory cost changes at different stages of query execution. It can be observed that GETL maintains stable memory usage in most cases. Additionally, the method of building indexes based on underlying nested pairs ensures that memory growth is not significant for most queries during the graph customization phase. The sharp increase in memory during the UGR2Model stage of Q1 is due to the mapping generating a large number of RDF triple results. In Q5, we materialize graph schema to efficiently group subgraphs and condense edges. Moreover, these grouped vertices and edges need to be preserved to support the newly mapped layer. Thus leading to limited memory savings.

E. Intelligibility of DMLs

We finally analyze the intelligibility of programs for user-defined recursive queries. Intelligibility refers to the conciseness of DML expressions, which is crucial for understanding the program and enabling non-expert users to master it quickly. We use *Halstead Metrics* [50] for quantitative measurement, which analyzes program effort, denoted as P_E (where $P_E = P_D \cdot P_V$), by combining the operators and operands in the code. It involves two metrics: the program volume $P_V = (N_1 + N_2) \cdot \log_2(n_1 + n_2)$ and the program difficulty $P_D = (n_1/2) \cdot (N_2/n_2)$, where n_1 and n_2 are the numbers of distinct operators and operands, and N_1 and N_2 are the total numbers of operators and operands, respectively.

TABLE III
INTELLIGIBILITY METRICS (P_E) FOR VARIOUS DMLs

	SSSP	CC	PR	CP
Gremlin	6939.3	568.87	8401.77	114.12
GraphScope	3233	1512.67	5614.41	240.53
RaSQL	2712.26	1580.23	7227.4	1856.95
GETL	1437.39	317.84	3849.66	990.4

We compare the program effort of algorithms implemented using the Gremlin, GraphScope, RaSQL, and hybrid DML of GETL. The algorithms include Single-Source Shortest Path (SSSP) [51], Connected Component (CC) [49], PageRank

(PR) [52], and Count Paths (CP) [53]. Partial examples of the algorithms are available in their respective papers. The results are shown in Table III. It can be observed that in most cases, GETL requires the least program effort P_E . Gremlin and GraphScope excel at the CP algorithm because CP does not require complex recursive aggregation logic or temporary variables, relying primarily on traversal operations, which is a core strength of Gremlin-like syntax. RaSQL demonstrates strong intelligibility by extending SQL's declarative syntax, but it struggles to handle complex join operations succinctly.

VIII. RELATED WORKS

Conversion between data models. Numerous efforts have been dedicated to establishing direct mappings between the LPG, RDF, RDF-star, and the relational model [10], [14], [15], [17]–[19], [54]–[58]. Although these studies offer valuable insights into model-to-model conversions, they are not the ideal solution for GETL's unified approach to supporting interoperability among multiple models. On the other hand, OneGraph [47], Multilayer Graphs [24], and Statement Graphs [25] achieve interoperability between LPG, RDF, and RDF-star from the perspective of a unified data model. OneGraph associates each statement with a statement identifier (SID) and explores the technical challenges involved. Multilayer Graph builds relationships across multiple layers by adding edge IDs as connectors. Statement Graph constructs a directed acyclic graph centered on statements. While these models are intuitive and general, they lack direct support for ETL-specific higher-order features such as *Path as Vertex*. By contrast, UGR treats higher-order constructs as set-valued entities and provides algebraic set operations over nested pairs, making features like *Path as Vertex* native constructions in UGR, whereas SG- or MG-based designs typically need additional reification layers to emulate similar abstractions.

Programming Interfaces. Some efforts have been made for graph queries from two perspectives: (1) Describing graph analysis algorithms. StarPlat [59] enables code generation for multi-core, distributed, and many-core systems based on the same interface specification. Fregel [60] proposes a functional approach to vertex-centric graph processing, abstracting the computation of each vertex into a high-order function. GraphIt [61] facilitates the rapid implementation of algorithms with varying performance characteristics. (2) Developing graph query languages for graph databases [27], [28], [31], [62]–[67]. However, graph analytics programming interfaces fail to effectively support analysis on graph traversal queries, as they typically require predefined graphs as input. Graph database query languages, although effective for graph pattern matching, lack adequate support for special mappings and user-defined recursive queries. Conversely, GETL introduces a more targeted programming interface, specifically tailored for the graph ETL process.

ETL tools for graph data. Several ETL tools have been developed for graph data. Neo4j [29] is equipped with an ETL tool that enables the importation of data from relational tables to Neo4j. However, this tool supports only a single mapping

for LPG. It does not accommodate special mappings, such as Q4 and Q5. Table2Graph [68] facilitates the conversion of multiple relational tables into an LPG-based ETL process. Berro et al. [69] proposes an ETL chain for warehousing based on RDF. GraphBuilder [70] is designed to alleviate many complexities of graph construction, such as graph formation, tabulation, and partitioning. The motivation behind GETL fundamentally differs from these tools. GETL seeks to bridge the graph ETL process across diverse data models, thereby enhancing interactions between applications and systems.

IX. CONCLUSIONS

This work makes the first step toward a graph ETL framework specialized for multiple graph data models. The highlights of GETL are threefold: (1) a unified graph representation that is compatible with various graph models; (2) a layered architecture that offers considerable flexibility and enables the customization of graph schemas and graph transformations; and (3) a programming interface tailored to the graph ETL process, designed to support operations that are intrinsic to graph ETL workflows. In the future, we plan to further extend GETL by implementing additional adapters (*i.e.*, extractors/loaders) that define the corresponding mapping rules, so as to support richer data models such as hypergraphs and temporal graphs.

REFERENCES

- [1] M. Mountantonakis and Y. Tzitzikas, “Large-scale semantic integration of linked data: A survey,” *CSUR*, vol. 52, no. 5, pp. 1–40, 2019.
- [2] M. Buron, F. Goasdoué, I. Manolescu, and M.-L. Mugnier, “Obi-wan: ontology-based rdf integration of heterogeneous data,” in *PVLDB*, 2020.
- [3] R. Angles, “The property graph database model,” in *AMW*, 2018.
- [4] H. R. Vyawahare, P. P. Karde, and V. M. Thakare, “A hybrid database approach using graph and relational database,” in *RICE*, 2018.
- [5] P. Yi, L. Liang, D. Zhang, Y. Chen, J. Zhu, X. Liu, K. Tang, J. Chen, H. Lin, L. Qiu *et al.*, “Kgfabric: A scalable knowledge graph warehouse for enterprise data interconnection,” in *PVLDB*, 2024.
- [6] X. Li, W. Zeng, Z. Wang, D. Zhu, J. Xu, W. Yu, and J. Zhou, “Graphar: An efficient storage scheme for graph data in data lakes,” *PVLDB*, vol. 18, no. 3, p. 530–543, 2024.
- [7] Taobao, “Taobao,” <https://www.taobao.com/>, 2025.
- [8] S. Yu, S. Gong, Q. Tao, S. Shen, Y. Zhang, W. Yu, P. Liu, Z. Zhang, H. Li, X. Luo *et al.*, “Lsmgraph: A high-performance dynamic graph storage system with multi-level csr,” *SIGMOD*, vol. 2, no. 6, pp. 1–28, 2024.
- [9] H. Li, Q. Tao, S. Yu, S. Gong, Y. Zhang, F. Yao, W. Yu, G. Yu, and J. Zhou, “Gastcoco: Graph storage and coroutine-based prefetch co-design for dynamic graph processing,” *PVLDB*, vol. 17, no. 13, p. 4827–4839, 2025.
- [10] S. Khayatbashi, S. Ferrada, and O. Hartig, “Converting property graphs to RDF: a preliminary study of the practical impact of different mappings,” in *SIGMOD*, 2022.
- [11] A. Group, “Graphscope,” <https://graphscope.io/journey/>, 2025.
- [12] A. Amiri, “Designing a distribution network in a supply chain system: Formulation and efficient solution procedure,” *Eur. J. Oper. Res.*, vol. 171, no. 2, pp. 567–576, 2006.
- [13] D. Tomaszuk, “RDF data in property graph model,” in *MTSR*, 2016.
- [14] D. Tomaszuk, R. Angles, and H. Thakkar, “PGO: describing property graphs in RDF,” *IEEE Access*, vol. 8, pp. 118 355–118 369, 2020.
- [15] V. M. de Sousa and L. M. del Val Cura, “Logical design of graph databases from an entity-relationship conceptual model,” in *iiWAS*, 2018.
- [16] E. Bytyçi, L. Ahmed, and G. Gashi, “RDF mapper: Easy conversion of relational databases to RDF,” in *WEBIST*, 2018.
- [17] G. Abuoda, D. Dell’Aglio, A. Keen, and K. Hose, “Transforming rdf-star to property graphs: A preliminary analysis of transformation approaches,” in *ISWC*, 2022.
- [18] R. Angles, H. Thakkar, and D. Tomaszuk, “Mapping rdf databases to property graph databases,” *IEEE Access*, vol. 8, pp. 86 091–86 110, 2020.
- [19] K. Rabbani, M. Lissandrini, A. Bonifati, and K. Hose, “Transforming rdf graphs to property graphs using standardized schemas,” *SIGMOD*, vol. 2, no. 6, pp. 1–25, 2024.
- [20] J. E. L. Gayo, E. Prud’Hommeaux, I. Boneva, and D. Kontokostas, *Validating RDF data*. Morgan & Claypool Publishers, 2017.
- [21] D. Wu, H.-T. Wang, and A. U. Tansel, “A survey for managing temporal data in rdf,” *Information Systems*, vol. 122, p. 102368, 2024.
- [22] R. Angles, A. Bonifati, S. Dumbrava, G. Fletcher, A. Green, J. Hidders, B. Li, L. Libkin, V. Marsault, W. Martens *et al.*, “Pg-schema: Schemas for property graphs,” *SIGMOD*, 2023.
- [23] D. W. Wardani and J. Kiing, “Semantic mapping relational to graph model,” in *IC3INA*, 2014.
- [24] R. Angles, A. Hogan, O. Lassila, C. Rojas, D. Schwabe, P. A. Szekely, and D. Vrgoc, “Multilayer graphs: a unified data model for graph databases,” in *GRADES & NDA*, 2022.
- [25] E. Gelling, G. Fletcher, and M. Schmidt, “Statement graphs: Unifying the graph data model landscape,” in *DASFAA*, 2024, pp. 364–376.
- [26] R.-D. C. Group, “Rdf-star and sparql-star,” <https://www.w3.org/2021/12/rdf-star.html>, 2021.
- [27] M. A. Rodriguez, “The gremlin graph traversal machine and language (invited talk),” in *DBPL*, 2015.
- [28] N. Francis, A. Green, P. Guagliardo, L. Libkin, T. Lindaaker, V. Marsault, S. Plantikow, M. Rydberg, P. Selmer, and A. Taylor, “Cypher: An evolving query language for property graphs,” in *SIGMOD*, 2018.
- [29] “Neo4j,” <https://neo4j.com/>, 2025.
- [30] “Tinkerpop,” <https://tinkerpop.apache.org/>, 2025.
- [31] A. Deutsch, N. Francis, A. Green, K. Hare, B. Li, L. Libkin, T. Lindaaker, V. Marsault, W. Martens, J. Michels *et al.*, “Graph pattern matching in gsql and sql/pgq,” in *SIGMOD*, 2022.
- [32] P. Zhao, X. Li, D. Xin, and J. Han, “Graph cube: on warehousing and olap multidimensional networks,” in *SIGMOD*, 2011.
- [33] W. Fan, T. He, L. Lai, X. Li, Y. Li, Z. Li, Z. Qian, C. Tian, L. Wang, J. Xu, Y. Yao, Q. Yin, W. Yu, K. Zeng, K. Zhao, J. Zhou, D. Zhu, and R. Zhu, “Graphscope: A unified engine for big graph processing,” in *PVLDB*, 2021.
- [34] D. Beckett, T. Berners-Lee, E. Prud’hommeaux, and G. Carothers, “Rdf 1.1 turtle,” <https://www.w3.org/TR/2014/REC-turtle-20140225/>, 2014.
- [35] O. Hartig, P.-A. Champin, G. Kellogg, and A. Seaborne, “Rdf-star and sparql-star,” <https://www.w3.org/2021/12/rdf-star.html>, 2021.
- [36] “Rdf4j,” <https://github.com/eclipse-rdf4j/rdf4j>, 2025.
- [37] “Mysql,” <https://www.mysql.com/>, 2025.
- [38] “Tinkerpop documentation,” <https://tinkerpop.apache.org/docs/3.7.2/reference/>, 2025.
- [39] A. Shkapsky, M. Yang, M. Interlandi, H. Chiu, T. Condie, and C. Zaniolo, “Big data analytics with datalog queries on spark,” in *SIGMOD*, 2016.
- [40] Y. Zhang, Q. Gao, L. Gao, and C. Wang, “Maiter: An asynchronous graph processing framework for delta-based accumulative iterative computation,” *TPDS*, vol. 25, no. 8, pp. 2091–2100, 2013.
- [41] “Amazon neptune,” <https://aws.amazon.com/cn/neptune/>, 2023.
- [42] C. Ge, X. Liu, L. Chen, B. Zheng, and Y. Gao, “Largeea: Aligning entities for large-scale knowledge graphs,” in *PVLDB*, 2022.
- [43] K. Xin, Z. Sun, W. Hua, W. Hu, J. Qu, and X. Zhou, “Large-scale entity alignment via knowledge graph merging, partitioning and embedding,” in *CIKM*, 2022.
- [44] “Movielens,” <https://grouplens.org/datasets/movielens/25m/>, 2019.
- [45] “Wikipedia,” https://databus.dbpedia.org/dbpedia/mappings/mappingbased-objects/2022.12.01/mappingbased-objects_lang=en.ttl.bz2, 2022.
- [46] “Ldbc social network benchmark graphs,” <https://repository.surfsara.nl/datasets/cwi/ldbc-snb-interactive-v1-dataset-v100>, 2022.
- [47] O. Lassila, M. Schmidt, O. Hartig, B. Bebee, D. Bechberger, W. Broekema, A. Khandelwal, K. Lawrence, C. López-Enriquez, R. Sharda, and B. B. Thompson, “The onegraph vision: Challenges of breaking the graph model lock-in¹,” *Semantic Web*, vol. 14, no. 1, pp. 125–134, 2023.
- [48] J. Gu, Y. H. Watanabe, W. A. Mazza, A. Shkapsky, M. Yang, L. Ding, and C. Zaniolo, “Rasql: Greater power and performance for big data analytics with recursive-aggregate-sql on spark,” in *SIGMOD*, 2019.
- [49] T. Hsu, V. Ramachandran, and N. Dean, “Parallel implementation of algorithms for finding connected components in graphs,” in *DIMACS*, 1994.
- [50] M. H. Halstead, *Elements of Software Science (Operating and programming systems series)*. Elsevier Science Inc., 1977.

- [51] V. T. Chakaravarthy, F. Checconi, P. Murali, F. Petrini, and Y. Sabharwal, “Scalable single source shortest path algorithms for massively parallel systems,” *TPDS*, vol. 28, no. 7, pp. 2031–2045, 2017.
- [52] M. Bianchini, M. Gori, and F. Scarselli, “Inside pagerank,” *ACM Trans. Internet Techn.*, vol. 5, no. 1, pp. 92–128, 2005.
- [53] E. T. Bax, “Algorithms to count paths and cycles,” *Inf. Process. Lett.*, vol. 52, no. 5, pp. 249–252, 1994.
- [54] G. Abuoda, D. Dell’Aglio, A. Keen, and K. Hose, “Transforming rdf-star to property graphs: A preliminary analysis of transformation approaches - extended version,” *CoRR*, vol. abs/2210.05781, 2022.
- [55] O. Hartig, “Reconciliation of rdf* and property graphs,” *CoRR*, vol. abs/1409.3288, 2014. [Online]. Available: <http://arxiv.org/abs/1409.3288>
- [56] R. Angles, H. Thakkar, and D. Tomaszuk, “RDF and property graphs interoperability: Status and issues,” in *AMW*, 2019.
- [57] G. Yuan, J. Lu, Z. Yan, and S. Wu, “A survey on mapping semi-structured data and graph data to relational data,” *ACM Comput. Surv.*, vol. 55, no. 10, pp. 218:1–218:38, 2023.
- [58] O. Hartig, “Foundations to query labeled property graphs using SPARQL,” in *SEMANTiCS*, 2019.
- [59] N. Behera, A. Kumar, E. R. T. S. Nitish, R. P. Muniasamy, and R. Nasre, “Starplat: A versatile DSL for graph analytics,” *CoRR*, vol. abs/2305.03317, 2023.
- [60] K. Emoto, K. Matsuzaki, Z. Hu, A. Morihata, and H. Iwasaki, “Think like a vertex, behave like a function! a functional DSL for vertex-centric big graph processing,” in *ICFP*, 2016.
- [61] Y. Zhang, M. Yang, R. Baghdadi, S. Kamil, J. Shun, and S. P. Amarasinghe, “Graphit: a high-performance graph DSL,” in *OOPSLA*, 2018.
- [62] W3C, “Sparql 1.1 query language,” <https://www.w3.org/TR/sparql11-query/>, 2013.
- [63] A. Deutsch, Y. Xu, M. Wu, and V. Lee, “Tigergraph: A native mpp graph database,” in *arXiv:1901.08248*, 2019.
- [64] O. Van Rest, S. Hong, J. Kim, X. Meng, and H. Chafi, “Pgql: a property graph query language,” in *GRADES*, 2016.
- [65] Z. Pan, T. Wu, Q. Zhao, Q. Zhou, Z. Peng, J. Li, Q. Zhang, G. Feng, and X. Zhu, “Geaflow: A graph extended and accelerated dataflow system,” *Proc. ACM Manag. Data*, vol. 1, no. 2, pp. 191:1–191:27, 2023.
- [66] R. Angles, M. Arenas, P. Barceló, P. A. Boncz, G. H. L. Fletcher, C. Gutierrez, T. Lindaaker, M. Paradies, S. Plantikow, J. F. Sequeda, O. van Rest, and H. Voigt, “G-CORE: A core for future graph query languages,” in *SIGMOD*, 2018.
- [67] “A new standard for a property graph database language,” <https://www.iso.org/standard/76120.html>, 2024.
- [68] S. Lee, B. H. Park, S. Lim, and M. Shankar, “Table2graph: A scalable graph construction from relational tables using map-reduce,” in *Big-DataService*, 2015.
- [69] A. Berro, I. Megdiche, and O. Teste, “Graph-based ETL processes for warehousing statistical open data,” in *ICEIS*, 2015.
- [70] N. Jain, G. Liao, and T. L. Willke, “Graphbuilder: scalable graph ETL framework,” in *GRADES*, 2013.



Feng Yao received the MS degree in computer science from Northeastern University, China, in 2021. He is currently working toward a PhD degree in computer science at Northeastern University, China. His research interests include distributed graph processing and data mining.



Xiaokang Yang is currently a Master’s student in Computer Science at Northeastern University, China. His research interests include distributed and parallel computation and data mining.



Shufeng Gong received the PhD degree in computer science from Northeastern University, China, in 2021. He is currently a lecturer with Northeastern University, China. His research interests include cloud computing, distributed graph processing, and vector databases.



Qian Tao received the PhD degree in computer science and technology from Beihang University, China in 2021. He is currently an engineer in Alibaba Group, China. His research interests include graph neural networks, graph computations, and large language models.



Yanfeng Zhang received the PhD degree in computer science from Northeastern University, China, in 2012. He is currently a professor with Northeastern University, China. His research consists of distributed systems and big data processing. He has published many papers in the above areas. His paper in SoCC 2011 was honored with “Paper of Distinction”.



Paper at SIGMOD 2017 and VLDB 2010, and the SIGMOD Research Highlight Award in 2018. Prior to Alibaba, Wenyuan was a founding member of 7Bridges Ltd and a Research Scientist at Facebook.



Ge Yu (Senior Member, IEEE) received the PhD degree in computer science from the Kyushu University of Japan, in 1996. He is now a professor with Northeastern University, China. His current research interests include distributed and parallel systems, cloud computing, big data management, and blockchain techniques and systems. He has published more than 200 papers in refereed journals and conferences. He is the CCF fellow and the ACM member.