

NebulaGraph Database Manual

v3.6.0

Table of contents

1. Welcome to NebulaGraph 3.6.0 Documentation	6
1.1 Getting started	6
1.2 Release notes	6
1.3 Other Sources	6
1.4 Symbols used in this manual	6
1.5 Modify errors	7
2. Introduction	8
2.1 What is NebulaGraph	8
2.2 Data modeling	12
2.3 Path types	14
2.4 VID	16
2.5 NebulaGraph architecture	18
3. Quick start	37
3.1 Quickly deploy NebulaGraph using Docker	38
3.2 Deploy NebulaGraph on-premise	0
3.3 nGQL cheatsheet	0
4. nGQL guide	0
4.1 nGQL overview	0
4.2 Data types	0
4.3 Operators	0
4.4 Functions and expressions	0
4.5 General queries statements	0
4.6 Clauses and options	0
4.7 Variables and composite queries	0
4.8 Space statements	0
4.9 Tag statements	0
4.10 Edge type statements	0
4.11 Vertex statements	0
4.12 Edge statements	0
4.13 Native index statements	0
4.14 Full-text index statements	0
4.15 Query tuning and terminating statements	0
4.16 Job manager and the JOB statements	0
5. Deploy and install	0
5.1 Prepare resources for compiling, installing, and running NebulaGraph	0

5.2	Compile and install	0
5.3	Local single-node installation	0
5.4	Deploy a NebulaGraph cluster with RPM/DEB package on multiple servers	0
5.5	Deploy NebulaGraph with Docker Compose	0
5.6	Install NebulaGraph with ecosystem tools	0
5.7	Manage NebulaGraph Service	0
5.8	Connect to NebulaGraph	0
5.9	Manage Storage hosts	0
5.10	Upgrade NebulaGraph to 3.6.0	0
5.11	Uninstall NebulaGraph	0
6.	Configure and log	0
6.1	Configurations	0
6.2	Log management	0
7.	Monitor	0
7.1	Query NebulaGraph metrics	0
7.2	RocksDB statistics	0
8.	Data security	0
8.1	Authentication and authorization	0
8.2	SSL encryption	0
9.	Backup and restore	0
9.1	NebulaGraph BR Community	0
9.2	Backup and restore data with snapshots	0
10.	Synchronize and migrate	0
10.1	BALANCE syntax	0
11.	Import and export	0
11.1	Import tools	0
11.2	NebulaGraph Importer	0
11.3	NebulaGraph Exchange	0
12.	Connectors	0
12.1	NebulaGraph Spark Connector	0
12.2	NebulaGraph Flink Connector	0
13.	Best practices	0
13.1	Compaction	0
13.2	Storage load balance	0
13.3	Graph data modeling suggestions	0
13.4	System design suggestions	0
13.5	Execution plan	0
13.6	Processing super vertices	0

13.7	Enable AutoFDO for NebulaGraph	0
13.8	Best practices	0
14.	Clients	0
14.1	Clients overview	0
14.2	NebulaGraph Console	0
14.3	NebulaGraph CPP	0
14.4	NebulaGraph Java	0
14.5	NebulaGraph Python	0
14.6	NebulaGraph Go	0
15.	Studio	0
15.1	About NebulaGraph Studio	0
15.2	Deploy and connect	0
15.3	Quick start	0
15.4	Troubleshooting	0
16.	Dashboard (Community)	0
16.1	What is NebulaGraph Dashboard Community Edition	0
16.2	Deploy Dashboard Community Edition	0
16.3	Connect Dashboard	0
16.4	Dashboard	0
16.5	Metrics	0
17.	NebulaGraph Operator	0
17.1	What is NebulaGraph Operator	0
17.2	Getting started	0
17.3	NebulaGraph Operator management	0
17.4	Cluster administration	0
17.5	FAQ	0
18.	Graph computing	0
18.1	NebulaGraph Algorithm	0
19.	NebulaGraph Bench	0
19.1	Scenario	0
19.2	Release note	0
19.3	Test process	0
20.	FAQ	0
20.1	About manual updates	0
20.2	About legacy version compatibility	0
20.3	About execution errors	0
20.4	About design and functions	0
20.5	About operation and maintenance	0

20.6 About connections	0
21. Appendix	0
21.1 Release Note	0
21.2 Ecosystem tools overview	0
21.3 Port guide for company products	0
21.4 How to Contribute	0
21.5 History timeline for NebulaGraph	0
21.6 Error code	0

1. Welcome to NebulaGraph 3.6.0 Documentation

Note

This manual is revised on 2024-3-6, with GitHub commit [574ca11e9b](#).

NebulaGraph is a distributed, scalable, and lightning-fast graph database. It is the optimal solution in the world capable of hosting graphs with dozens of billions of vertices (nodes) and trillions of edges (relationships) with millisecond latency.

1.1 Getting started

- [Quick start](#)
- [Preparations before deployment](#)
- [nGQL cheatsheet](#)
- [FAQ](#)
- [Ecosystem Tools](#)
- [Live Demo](#)

1.2 Release notes

- [NebulaGraph Community Edition 3.6.0](#)
- [NebulaGraph Dashboard Community](#)
- [NebulaGraph Studio](#)

1.3 Other Sources

- [To cite NebulaGraph](#)
- [Forum](#)
- [NebulaGraph Homepage](#)
- [Blogs](#)
- [Videos](#)
- [Chinese Docs](#)

1.4 Symbols used in this manual

Note

Additional information or operation-related notes.

Caution

May have adverse effects, such as causing performance degradation or triggering known minor problems.

**Warning**

May lead to serious issues, such as data loss or system crash.

**Danger**

May lead to extremely serious issues, such as system damage or information leakage.

**Compatibility**

The compatibility notes between nGQL and openCypher, or between the current version of nGQL and its prior ones.

**Enterpriseonly**

Differences between the NebulaGraph Community and Enterprise editions.

1.5 Modify errors

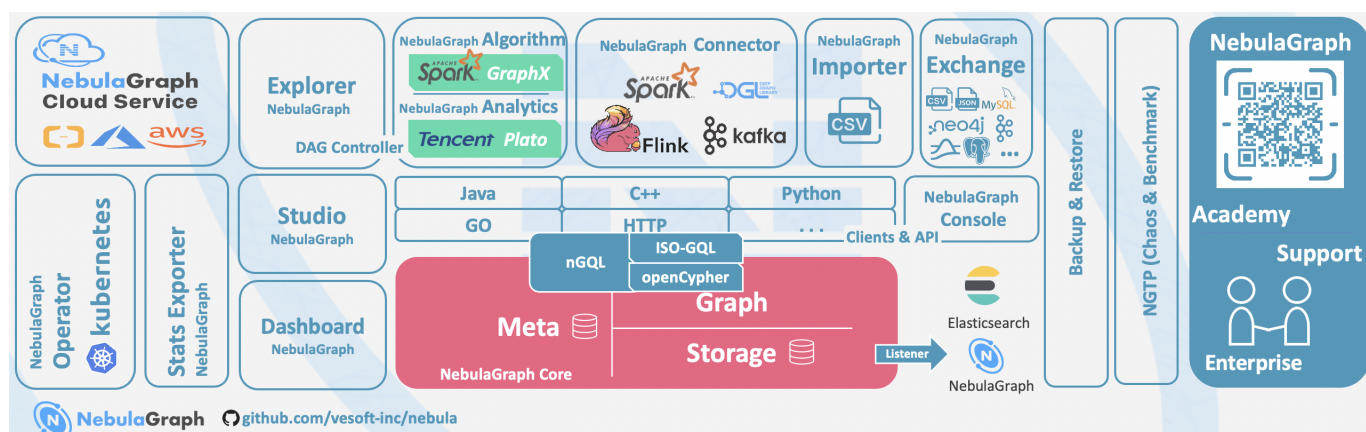
This NebulaGraph manual is written in the Markdown language. Users can click the pencil sign on the upper right side of each document title and modify errors.

Last update: November 29, 2023

2. Introduction

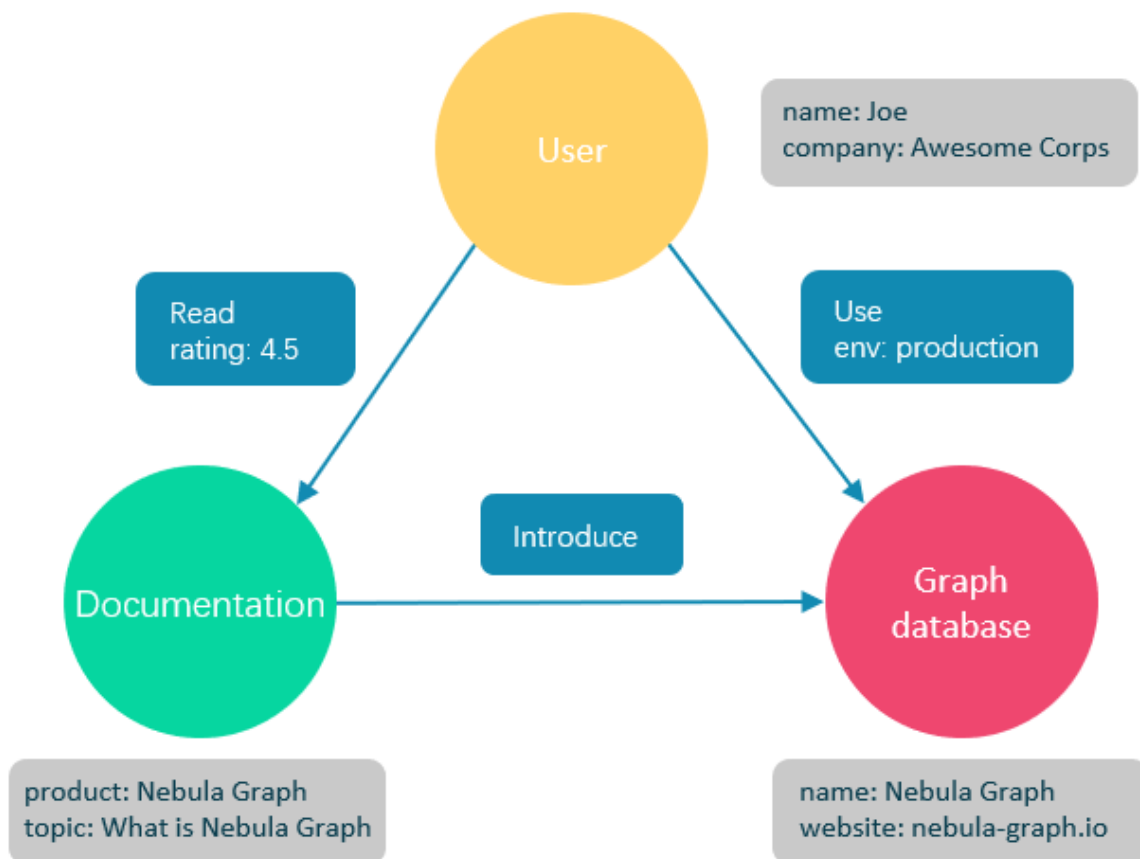
2.1 What is NebulaGraph

NebulaGraph is an open-source, distributed, easily scalable, and native graph database. It is capable of hosting graphs with hundreds of billions of vertices and trillions of edges, and serving queries with millisecond-latency.



2.1.1 What is a graph database

A graph database, such as NebulaGraph, is a database that specializes in storing vast graph networks and retrieving information from them. It efficiently stores data as vertices (nodes) and edges (relationships) in labeled property graphs. Properties can be attached to both vertices and edges. Each vertex can have one or multiple tags (labels).



Graph databases are well suited for storing most kinds of data models abstracted from reality. Things are connected in almost all fields in the world. Modeling systems like relational databases extract the relationships between entities and squeeze them into table columns alone, with their types and properties stored in other columns or even other tables. This makes data management time-consuming and cost-ineffective.

NebulaGraph, as a typical native graph database, allows you to store the rich relationships as edges with edge types and properties directly attached to them.

2.1.2 Advantages of NebulaGraph

Open source

NebulaGraph is open under the Apache 2.0 License. More and more people such as database developers, data scientists, security experts, and algorithm engineers are participating in the designing and development of NebulaGraph. To join the opening of source code and ideas, surf the [NebulaGraph GitHub page](#).

Outstanding performance

Written in C++ and born for graphs, NebulaGraph handles graph queries in milliseconds. Among most databases, NebulaGraph shows superior performance in providing graph data services. The larger the data size, the greater the superiority of NebulaGraph. For more information, see [NebulaGraph benchmarking](#).

High scalability

NebulaGraph is designed in a shared-nothing architecture and supports scaling in and out without interrupting the database service.

Developer friendly

NebulaGraph supports clients in popular programming languages like Java, Python, C++, and Go, and more are under development. For more information, see NebulaGraph [clients](#).

Reliable access control

NebulaGraph supports strict role-based access control and external authentication servers such as LDAP (Lightweight Directory Access Protocol) servers to enhance data security. For more information, see [Authentication and authorization](#).

Diversified ecosystem

More and more native tools of NebulaGraph have been released, such as [NebulaGraph Studio](#), [NebulaGraph Console](#), and [NebulaGraph Exchange](#). For more ecosystem tools, see [Ecosystem tools overview](#).

Besides, NebulaGraph has the ability to be integrated with many cutting-edge technologies, such as Spark, Flink, and HBase, for the purpose of mutual strengthening in a world of increasing challenges and chances.

OpenCypher-compatible query language

The native NebulaGraph Query Language, also known as nGQL, is a declarative, openCypher-compatible textual query language. It is easy to understand and easy to use. For more information, see [nGQL guide](#).

Future-oriented hardware with balanced reading and writing

Solid-state drives have extremely high performance and [they are getting cheaper](#). NebulaGraph is a product based on SSD. Compared with products based on HDD and large memory, it is more suitable for future hardware trends and easier to achieve balanced reading and writing.

Easy data modeling and high flexibility

You can easily model the connected data into NebulaGraph for your business without forcing them into a structure such as a relational table, and properties can be added, updated, and deleted freely. For more information, see [Data modeling](#).

High popularity

NebulaGraph is being used by tech leaders such as Tencent, Vivo, Meituan, and JD Digits. For more information, visit the [NebulaGraph official website](#).

2.1.3 Use cases

NebulaGraph can be used to support various graph-based scenarios. To spare the time spent on pushing the kinds of data mentioned in this section into relational databases and on bothering with join queries, use NebulaGraph.

Fraud detection

Financial institutions have to traverse countless transactions to piece together potential crimes and understand how combinations of transactions and devices might be related to a single fraud scheme. This kind of scenario can be modeled in graphs, and with the help of NebulaGraph, fraud rings and other sophisticated scams can be easily detected.

Real-time recommendation

NebulaGraph offers the ability to instantly process the real-time information produced by a visitor and make accurate recommendations on articles, videos, products, and services.

Intelligent question-answer system

Natural languages can be transformed into knowledge graphs and stored in NebulaGraph. A question organized in a natural language can be resolved by a semantic parser in an intelligent question-answer system and re-organized. Then, possible answers to the question can be retrieved from the knowledge graph and provided to the one who asked the question.

Social networking

Information on people and their relationships is typical graph data. NebulaGraph can easily handle the social networking information of billions of people and trillions of relationships, and provide lightning-fast queries for friend recommendations and job promotions in the case of massive concurrency.

2.1.4 Related links

- [Official website](#)
- [Docs](#)
- [Blogs](#)
- [Forum](#)
- [GitHub](#)

Last update: October 25, 2023

2.2 Data modeling

A data model is a model that organizes data and specifies how they are related to one another. This topic describes the Nebula Graph data model and provides suggestions for data modeling with NebulaGraph.

2.2.1 Data structures

NebulaGraph data model uses six data structures to store data. They are graph spaces, vertices, edges, tags, edge types and properties.

- **Graph spaces:** Graph spaces are used to isolate data from different teams or programs. Data stored in different graph spaces are securely isolated. Storage replications, privileges, and partitions can be assigned.
- **Vertices:** Vertices are used to store entities.
- In NebulaGraph, vertices are identified with vertex identifiers (i.e. `VID`). The `VID` must be unique in the same graph space. `VID` should be `int64`, or `fixed_string(N)`.
- A vertex has zero to multiple tags.



Compatibility

In NebulaGraph 2.x a vertex must have at least one tag. And in NebulaGraph 3.6.0, a tag is not required for a vertex.

- **Edges:** Edges are used to connect vertices. An edge is a connection or behavior between two vertices.
- There can be multiple edges between two vertices.
- Edges are directed. `->` identifies the directions of edges. Edges can be traversed in either direction.
- An edge is identified uniquely with `<a source vertex, an edge type, a rank value, and a destination vertex>`. Edges have no EID.
- An edge must have one and only one edge type.
- The rank value is an immutable user-assigned 64-bit signed integer. It identifies the edges with the same edge type between two vertices. Edges are sorted by their rank values. The edge with the greatest rank value is listed first. The default rank value is zero.
- **Tags:** Tags are used to categorize vertices. Vertices that have the same tag share the same definition of properties.
- **Edge types:** Edge types are used to categorize edges. Edges that have the same edge type share the same definition of properties.
- **Properties:** Properties are key-value pairs. Both vertices and edges are containers for properties.



Note

Tags and Edge types are similar to "vertex tables" and "edge tables" in the relational databases.

2.2.2 Directed property graph

NebulaGraph stores data in directed property graphs. A directed property graph has a set of vertices connected by directed edges. Both vertices and edges can have properties. A directed property graph is represented as:

$$G = \langle V, E, P_V, P_E \rangle$$

- V is a set of vertices.
- E is a set of directed edges.
- P_V is the property of vertices.
- P_E is the property of edges.

The following table is an example of the structure of the basketball player dataset. We have two types of vertices, that is **player** and **team**, and two types of edges, that is **serve** and **follow**.

Element	Name	Property name (Data type)	Description
Tag	player	name (string) age (int)	Represents players in the team.
Tag	team	name (string)	Represents the teams.
Edge type	serve	start_year (int) end_year (int)	Represents actions taken by players in the team. An action links a player with a team, and the direction is from a player to a team.
Edge type	follow	degree (int)	Represents actions taken by players in the team. An action links a player with another player, and the direction is from one player to the other player.

Note

NebulaGraph supports only directed edges.

Compatibility

NebulaGraph 3.6.0 allows dangling edges. Therefore, when adding or deleting, you need to ensure the corresponding source vertex and destination vertex of an edge exist. For details, see [INSERT VERTEX](#), [DELETE VERTEX](#), [INSERT EDGE](#), and [DELETE EDGE](#).

The MERGE statement in openCypher is not supported.

Last update: November 3, 2023

2.3 Path types

In graph theory, a path in a graph is a finite or infinite sequence of edges which joins a sequence of vertices. Paths are fundamental concepts of graph theory.

Paths can be categorized into 3 types: `walk`, `trail`, and `path`. For more information, see [Wikipedia](#).

The following figure is an example for a brief introduction.



2.3.1 Walk

A `walk` is a finite or infinite sequence of edges. Both vertices and edges can be repeatedly visited in graph traversal.

In the above figure C, D, and E form a cycle. So, this figure contains infinite paths, such as `A->B->C->D->E`, `A->B->C->D->E->C`, and `A->B->C->D->E->C->D`.

Note

`GO` statements use `walk`.

2.3.2 Trail

A `trail` is a finite sequence of edges. Only vertices can be repeatedly visited in graph traversal. The Seven Bridges of Königsberg is a typical `trail`.

In the above figure, edges cannot be repeatedly visited. So, this figure contains finite paths. The longest path in this figure consists of 5 edges: `A->B->C->D->E->C`.

Note

`MATCH`, `FIND PATH`, and `GET SUBGRAPH` statements use `trail`.

There are two special cases of trail, `cycle` and `circuit`. The following figure is an example for a brief introduction.



- cycle

A **cycle** refers to a closed **trail**. Only the terminal vertices can be repeatedly visited. The longest path in this figure consists of 3 edges: $A \rightarrow B \rightarrow C \rightarrow A$ or $C \rightarrow D \rightarrow E \rightarrow C$.

- circuit

A **circuit** refers to a closed **trail**. Edges cannot be repeatedly visited in graph traversal. Apart from the terminal vertices, other vertices can also be repeatedly visited. The longest path in this figure: $A \rightarrow B \rightarrow C \rightarrow D \rightarrow E \rightarrow C \rightarrow A$.

2.3.3 Path

A **path** is a finite sequence of edges. Neither vertices nor edges can be repeatedly visited in graph traversal.

So, the above figure contains finite paths. The longest path in this figure consists of 4 edges: $A \rightarrow B \rightarrow C \rightarrow D \rightarrow E$.

Last update: October 25, 2023

2.4 VID

In a graph space, a vertex is uniquely identified by its ID, which is called a VID or a Vertex ID.

2.4.1 Features

- The data types of VIDs are restricted to `FIXED_STRING(<N>)` or `INT64`. One graph space can only select one VID type.
- A VID in a graph space is unique. It functions just as a primary key in a relational database. VIDs in different graph spaces are independent.
- The VID generation method must be set by users, because NebulaGraph does not provide auto increasing ID, or UUID.
- Vertices with the same VID will be identified as the same one. For example:
- A VID is the unique identifier of an entity, like a person's ID card number. A tag means the type of an entity, such as driver, and boss. Different tags define two groups of different properties, such as driving license number, driving age, order amount, order taking alt, and job number, payroll, debt ceiling, business phone number.
- When two `INSERT` statements (neither uses a parameter of `IF NOT EXISTS`) with the same VID and tag are operated at the same time, the latter `INSERT` will overwrite the former.
- When two `INSERT` statements with the same VID but different tags, like `TAG A` and `TAG B`, are operated at the same time, the operation of `Tag A` will not affect `Tag B`.
- VIDs will usually be indexed and stored into memory (in the way of LSM-tree). Thus, direct access to VIDs enjoys peak performance.

2.4.2 VID Operation

- NebulaGraph 1.x only supports `INT64` while NebulaGraph 2.x supports `INT64` and `FIXED_STRING(<N>)`. In `CREATE SPACE`, VID types can be set via `vid_type`.
- `id()` function can be used to specify or locate a VID.
- `LOOKUP` or `MATCH` statements can be used to find a VID via property index.
- Direct access to vertices statements via VIDs enjoys peak performance, such as `DELETE xxx WHERE id(xxx) == "player100"` or `GO FROM "player100"`. Finding VIDs via properties and then operating the graph will cause poor performance, such as `LOOKUP | GO FROM $-.ids`, which will run both `LOOKUP` and `|` one more time.

2.4.3 VID Generation

VIDs can be generated via applications. Here are some tips:

- (Optimal) Directly take a unique primary key or property as a VID. Property access depends on the VID.
- Generate a VID via a unique combination of properties. Property access depends on property index.
- Generate a VID via algorithms like snowflake. Property access depends on property index.
- If short primary keys greatly outnumber long primary keys, do not enlarge the `N` of `FIXED_STRING(<N>)` too much. Otherwise, it will occupy a lot of memory and hard disks, and slow down performance. Generate VIDs via BASE64, MD5, hash by encoding and splicing.
- If you generate int64 VID via hash, the probability of collision is about 1/10 when there are 1 billion vertices. The number of edges has no concern with the probability of collision.

2.4.4 Define and modify a VID and its data type

The data type of a VID must be defined when you [create the graph space](#). Once defined, it cannot be modified.

A VID is set when you [insert a vertex](#) and cannot be modified.

2.4.5 Query `start vid` and global scan

In most cases, the execution plan of query statements in NebulaGraph (`MATCH`, `GO`, and `LOOKUP`) must query the `start vid` in a certain way.

There are only two ways to locate `start vid`:

1. For example, `GO FROM "player100" OVER` explicitly indicates in the statement that `start vid` is "player100".
2. For example, `LOOKUP ON player WHERE player.name == "Tony Parker"` or `MATCH (v:player {name:"Tony Parker"})` locates `start vid` by the index of the property `player.name`.

Last update: October 25, 2023

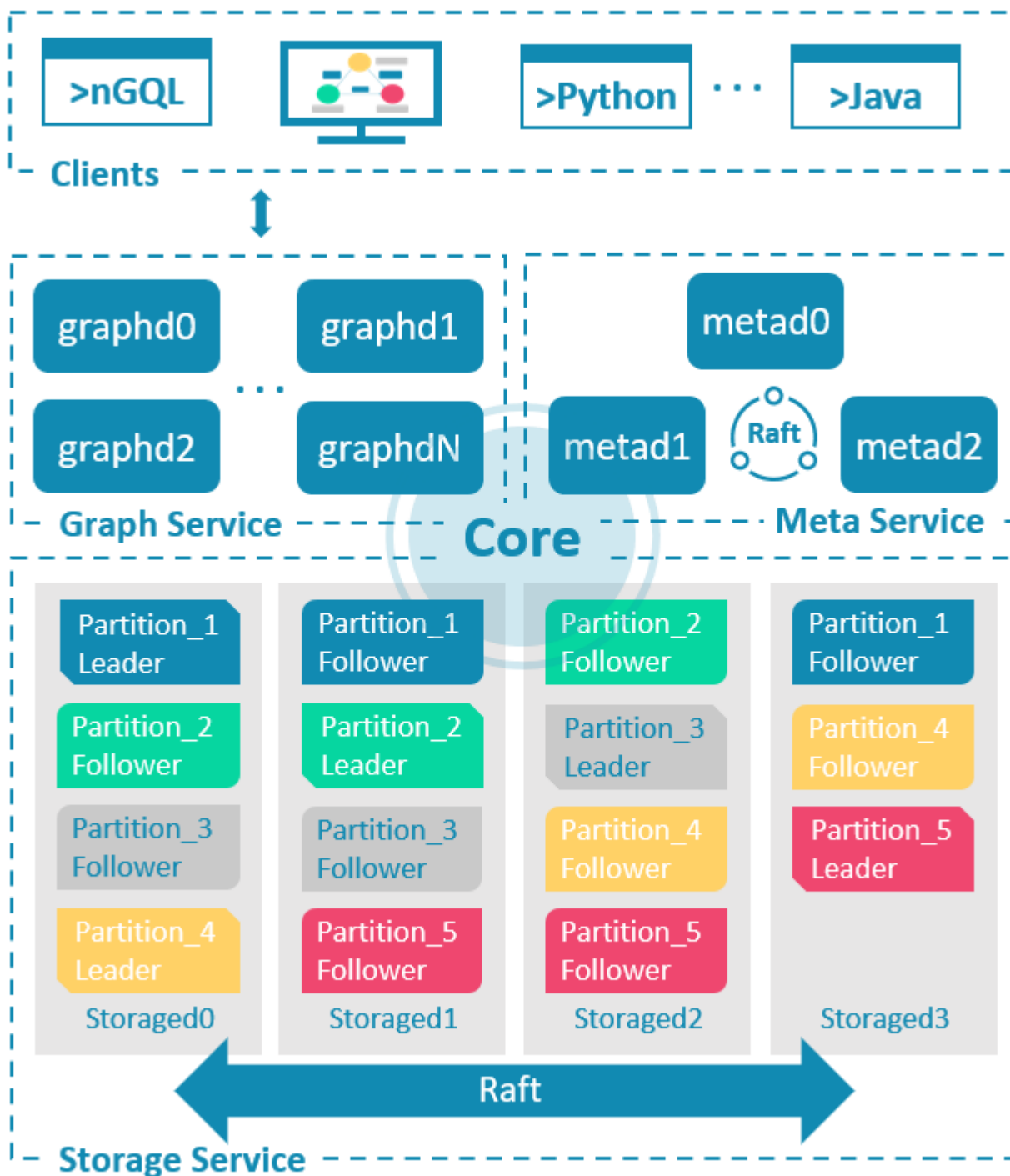
2.5 NebulaGraph architecture

2.5.1 Architecture overview

NebulaGraph consists of three services: the Graph Service, the Storage Service, and the Meta Service. It applies the separation of storage and computing architecture.

Each service has its executable binaries and processes launched from the binaries. Users can deploy a NebulaGraph cluster on a single machine or multiple machines using these binaries.

The following figure shows the architecture of a typical NebulaGraph cluster.



The Meta Service

The Meta Service in the NebulaGraph architecture is run by the `nebula-metad` processes. It is responsible for metadata management, such as schema operations, cluster administration, and user privilege management.

For details on the Meta Service, see [Meta Service](#).

The Graph Service and the Storage Service

NebulaGraph applies the separation of storage and computing architecture. The Graph Service is responsible for querying. The Storage Service is responsible for storage. They are run by different processes, i.e., nebula-graphd and nebula-storaged. The benefits of the separation of storage and computing architecture are as follows:

- Great scalability

The separated structure makes both the Graph Service and the Storage Service flexible and easy to scale in or out.

- High availability

If part of the Graph Service fails, the data stored by the Storage Service suffers no loss. And if the rest part of the Graph Service is still able to serve the clients, service recovery can be performed quickly, even unfelt by the users.

- Cost-effective

The separation of storage and computing architecture provides a higher resource utilization rate, and it enables clients to manage the cost flexibly according to business demands.

- Open to more possibilities

With the ability to run separately, the Graph Service may work with multiple types of storage engines, and the Storage Service may also serve more types of computing engines.

For details on the Graph Service and the Storage Service, see [Graph Service](#) and [Storage Service](#).

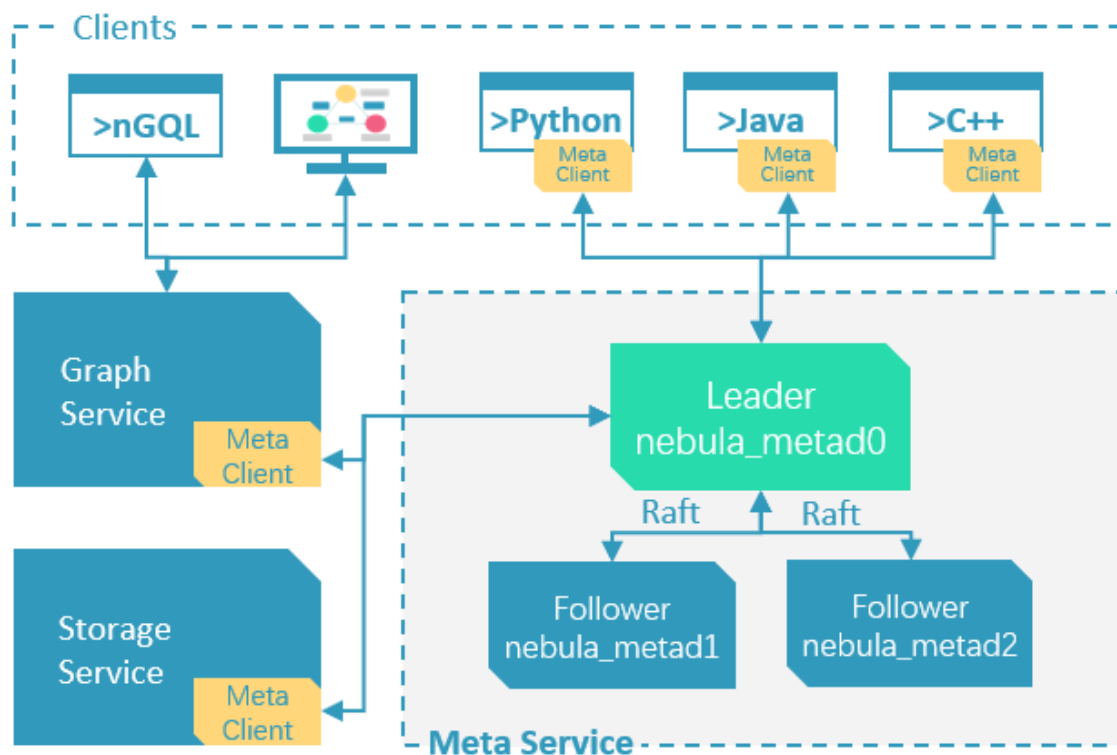
Last update: October 25, 2023

2.5.2 Meta Service

This topic introduces the architecture and functions of the Meta Service.

The architecture of the Meta Service

The architecture of the Meta Service is as follows:



The Meta Service is run by nebula-metad processes. Users can deploy nebula-metad processes according to the scenario:

- In a test environment, users can deploy one or three nebula-metad processes on different machines or a single machine.
- In a production environment, we recommend that users deploy three nebula-metad processes on different machines for high availability.

All the nebula-metad processes form a Raft-based cluster, with one process as the leader and the others as the followers.

The leader is elected by the majorities and only the leader can provide service to the clients or other components of NebulaGraph. The followers will be run in a standby way and each has a data replication of the leader. Once the leader fails, one of the followers will be elected as the new leader.

Note

The data of the leader and the followers will keep consistent through Raft. Thus the breakdown and election of the leader will not cause data inconsistency. For more information on Raft, see [Storage service architecture](#).

Functions of the Meta Service

MANAGES USER ACCOUNTS

The Meta Service stores the information of user accounts and the privileges granted to the accounts. When the clients send queries to the Meta Service through an account, the Meta Service checks the account information and whether the account has the right privileges to execute the queries or not.

For more information on NebulaGraph access control, see [Authentication](#).

MANAGES PARTITIONS

The Meta Service stores and manages the locations of the storage partitions and helps balance the partitions.

MANAGES GRAPH SPACES

NebulaGraph supports multiple graph spaces. Data stored in different graph spaces are securely isolated. The Meta Service stores the metadata of all graph spaces and tracks the changes of them, such as adding or dropping a graph space.

MANAGES SCHEMA INFORMATION

NebulaGraph is a strong-typed graph database. Its schema contains tags (i.e., the vertex types), edge types, tag properties, and edge type properties.

The Meta Service stores the schema information. Besides, it performs the addition, modification, and deletion of the schema, and logs the versions of them.

For more information on NebulaGraph schema, see [Data model](#).

MANAGES TTL INFORMATION

The Meta Service stores the definition of TTL (Time to Live) options which are used to control data expiration. The Storage Service takes care of the expiring and evicting processes. For more information, see [TTL](#).

MANAGES JOBS

The Job Management module in the Meta Service is responsible for the creation, queuing, querying, and deletion of jobs.

Last update: November 3, 2023

2.5.3 Graph Service

The Graph Service is used to process the query. It has four submodules: Parser, Validator, Planner, and Executor. This topic will describe the Graph Service accordingly.

The architecture of the Graph Service



After a query is sent to the Graph Service, it will be processed by the following four submodules:

1. **Parser:** Performs lexical analysis and syntax analysis.
2. **Validator:** Validates the statements.
3. **Planner:** Generates and optimizes the execution plans.
4. **Executor:** Executes the plans with operators.

Parser

After receiving a request, the statements will be parsed by Parser composed of Flex (lexical analysis tool) and Bison (syntax analysis tool), and its corresponding AST will be generated. Statements will be directly intercepted in this stage because of their invalid syntax.

For example, the structure of the AST of `GO FROM "Tim" OVER Like WHERE properties(edge).likeness > 8.0 YIELD dst(edge)` is shown in the following figure.



Validator

Validator performs a series of validations on the AST. It mainly works on these tasks:

- Validating metadata

Validator will validate whether the metadata is correct or not.

When parsing the `OVER`, `WHERE`, and `YIELD` clauses, Validator looks up the Schema and verifies whether the edge type and tag data exist or not. For an `INSERT` statement, Validator verifies whether the types of the inserted data are the same as the ones defined in the Schema.

- Validating contextual reference

Validator will verify whether the cited variable exists or not, or whether the cited property is variable or not.

For composite statements, like `$var = GO FROM "Tim" OVER Like YIELD dst(edge) AS ID; GO FROM $var.ID OVER serve YIELD dst(edge)`, Validator verifies first to see if `var` is defined, and then to check if the `ID` property is attached to the `var` variable.

- Validating type inference

Validator infers what type the result of an expression is and verifies the type against the specified clause.

For example, the `WHERE` clause requires the result to be a `bool` value, a `NULL` value, or `empty`.

- Validating the information of `*`

Validator needs to verify all the Schema that involves `*` when verifying the clause if there is a `*` in the statement.

Take a statement like `GO FROM "Tim" OVER * YIELD dst(edge), properties(edge).likeness, dst(edge)` as an example. When verifying the `OVER` clause, Validator needs to verify all the edge types. If the edge type includes `like` and `serve`, the statement would be `GO FROM "Tim" OVER like,serve YIELD dst(edge), properties(edge).likeness, dst(edge)`.

- Validating input and output

Validator will check the consistency of the clauses before and after the `|`.

In the statement `GO FROM "Tim" OVER like YIELD dst(edge) AS ID | GO FROM $-.ID OVER serve YIELD dst(edge)`, Validator will verify whether `$-.ID` is defined in the clause before the `|`.

When the validation succeeds, an execution plan will be generated. Its data structure will be stored in the `src/planner` directory.

Planner

In the `nebula-graphd.conf` file, when `enable_optimizer` is set to be `false`, Planner will not optimize the execution plans generated by Validator. It will be executed by Executor directly.

In the `nebula-graphd.conf` file, when `enable_optimizer` is set to be `true`, Planner will optimize the execution plans generated by Validator. The structure is as follows.



- Before optimization

In the execution plan on the right side of the preceding figure, each node directly depends on other nodes. For example, the root node `Project` depends on the `Filter` node, the `Filter` node depends on the `GetNeighbor` node, and so on, up to the leaf node `Start`. Then the execution plan is (not truly) executed.

During this stage, every node has its input and output variables, which are stored in a hash table. The execution plan is not truly executed, so the value of each key in the associated hash table is empty (except for the `Start` node, where the input variables hold the starting data), and the hash table is defined in `src/context/ExecutionContext.cpp` under the `nebula-graph` repository.

For example, if the hash table is named as `ResultMap` when creating the `Filter` node, users can determine that the node takes data from `ResultMap["GN1"]`, then puts the result into `ResultMap["Filter2"]`, and so on. All these work as the input and output of each node.

- Process of optimization

The optimization rules that Planner has implemented so far are considered RBO (Rule-Based Optimization), namely the pre-defined optimization rules. The CBO (Cost-Based Optimization) feature is under development. The optimized code is in the `src/optimizer/` directory under the `nebula-graph` repository.

RBO is a “bottom-up” exploration process. For each rule, the root node of the execution plan (in this case, the `Project` node) is the entry point, and step by step along with the node dependencies, it reaches the node at the bottom to see if it matches the rule.

As shown in the preceding figure, when the `Filter` node is explored, it is found that its children node is `GetNeighbors`, which matches successfully with the pre-defined rules, so a transformation is initiated to integrate the `Filter` node into the `GetNeighbors` node, the `Filter` node is removed, and then the process continues to the next rule. Therefore, when the `GetNeighbor` operator calls interfaces of the Storage layer to get the neighboring edges of a vertex during the execution stage, the Storage layer will directly filter out the unqualified edges internally. Such optimization greatly reduces the amount of data transfer, which is commonly known as filter pushdown.

Executor

The Executor module consists of Scheduler and Executor. The Scheduler generates the corresponding execution operators against the execution plan, starting from the leaf nodes and ending at the root node. The structure is as follows.



Each node of the execution plan has one execution operator node, whose input and output have been determined in the execution plan. Each operator only needs to get the values for the input variables, compute them, and finally put the results into the corresponding output variables. Therefore, it is only necessary to execute step by step from `start`, and the result of the last operator is returned to the user as the final result.

Source code hierarchy

The source code hierarchy under the `nebula-graph` repository is as follows.

```

|--src
|  |--graph
|     |--context    //contexts for validation and execution
|     |--executor   //execution operators
|     |--gc         //garbage collector
|     |--optimizer  //optimization rules
|     |--planner    //structure of the execution plans
|     |--scheduler  //scheduler
|     |--service    //external service management
|     |--session    //session management
|     |--stats      //monitoring metrics
|     |--util       //basic components
|     |--validator  //validation of the statements
|     |--visitor    //visitor expression

```

Last update: October 25, 2023

2.5.4 Storage Service

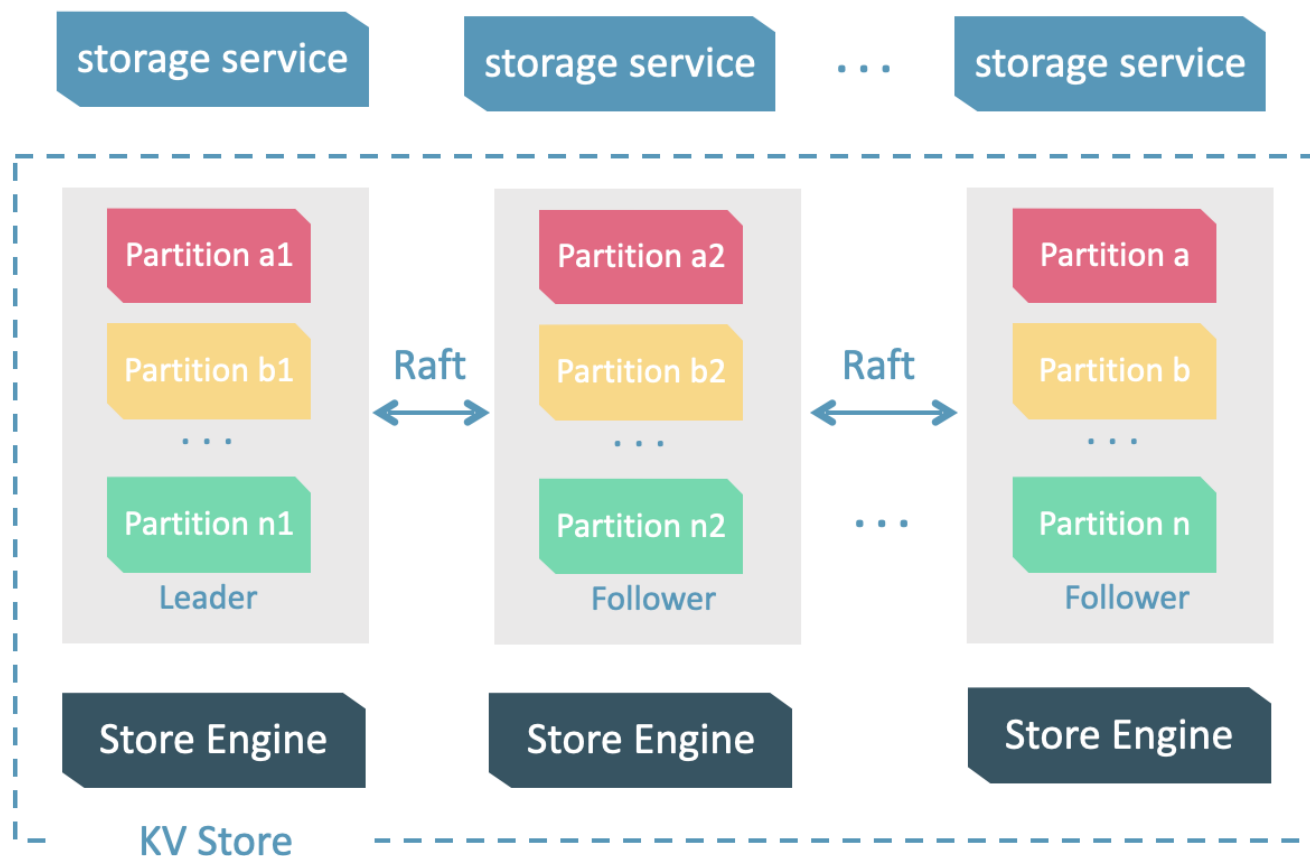
The persistent data of NebulaGraph have two parts. One is the [Meta Service](#) that stores the meta-related data.

The other is the Storage Service that stores the data, which is run by the nebula-storaged process. This topic will describe the architecture of the Storage Service.

Advantages

- High performance (Customized built-in KVStore)
- Great scalability (Shared-nothing architecture, not rely on NAS/SAN-like devices)
- Strong consistency (Raft)
- High availability (Raft)
- Supports synchronizing with the third party systems, such as [Elasticsearch](#).

The architecture of the Storage Service



The Storage Service is run by the nebula-storaged process. Users can deploy nebula-storaged processes on different occasions. For example, users can deploy 1 nebula-storaged process in a test environment and deploy 3 nebula-storaged processes in a production environment.

All the nebula-stored processes consist of a Raft-based cluster. There are three layers in the Storage Service:

- Storage interface

The top layer is the storage interface. It defines a set of APIs that are related to the graph concepts. These API requests will be translated into a set of KV operations targeting the corresponding [Partition](#). For example:

- `getNeighbors` : queries the in-edge or out-edge of a set of vertices, returns the edges and the corresponding properties, and supports conditional filtering.
- `insert vertex/edge` : inserts a vertex or edge and its properties.
- `getProps` : gets the properties of a vertex or an edge.

It is this layer that makes the Storage Service a real graph storage. Otherwise, it is just a KV storage.

- Consensus

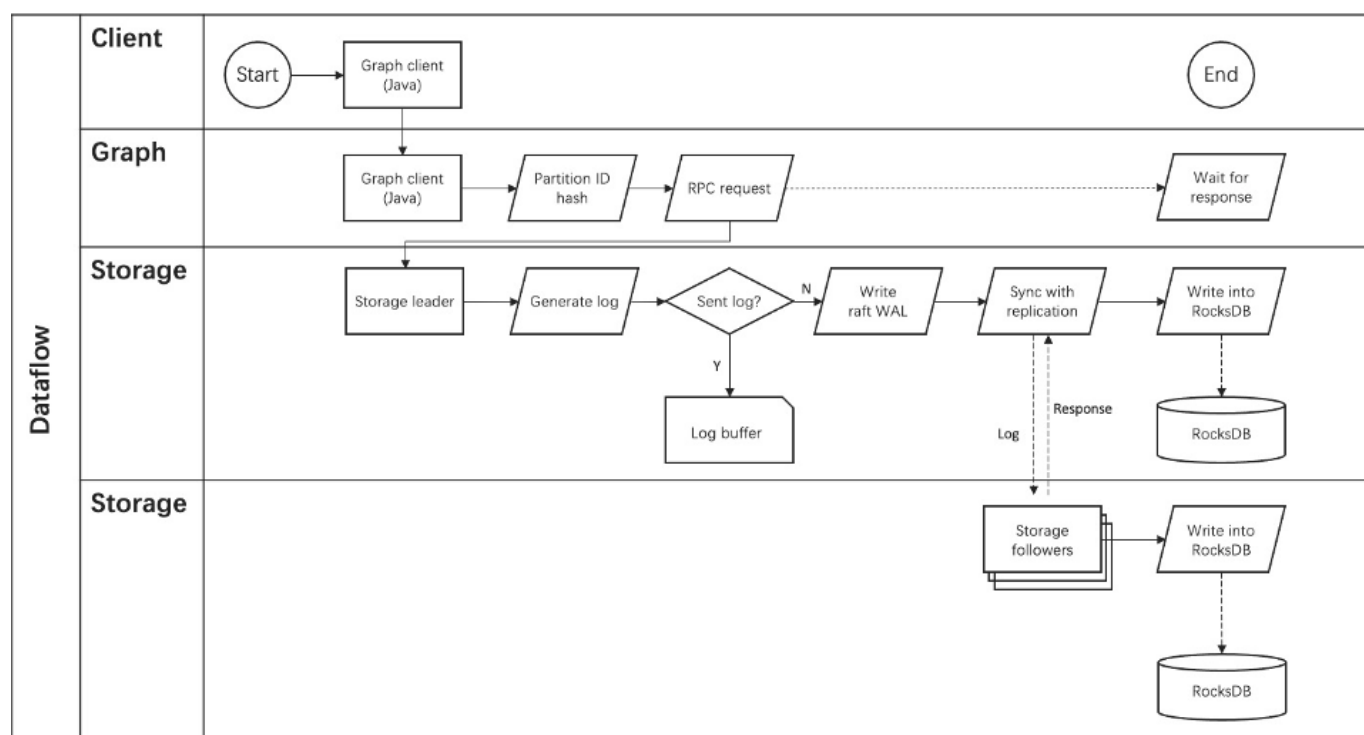
Below the storage interface is the consensus layer that implements [Multi Group Raft](#), which ensures the strong consistency and high availability of the Storage Service.

- Store engine

The bottom layer is the local storage engine library, providing operations like `get`, `put`, and `scan` on local disks. The related interfaces are stored in `KVStore.h` and `KVEngine.h` files. You can develop your own local store plugins based on your needs.

The following will describe some features of the Storage Service based on the above architecture.

Storage writing process



KVStore

NebulaGraph develops and customizes its built-in KVStore for the following reasons.

- It is a high-performance KVStore.
- It is provided as a (kv) library and can be easily developed for the filter pushdown purpose. As a strong-typed database, how to provide Schema during pushdown is the key to efficiency for NebulaGraph.
- It has strong data consistency.

Therefore, NebulaGraph develops its own KVStore with RocksDB as the local storage engine. The advantages are as follows.

- For multiple local hard disks, NebulaGraph can make full use of its concurrent capacities through deploying multiple data directories.
- The Meta Service manages all the Storage servers. All the partition distribution data and current machine status can be found in the meta service. Accordingly, users can execute a manual load balancing plan in meta service.



Note

NebulaGraph does not support auto load balancing because auto data transfer will affect online business.

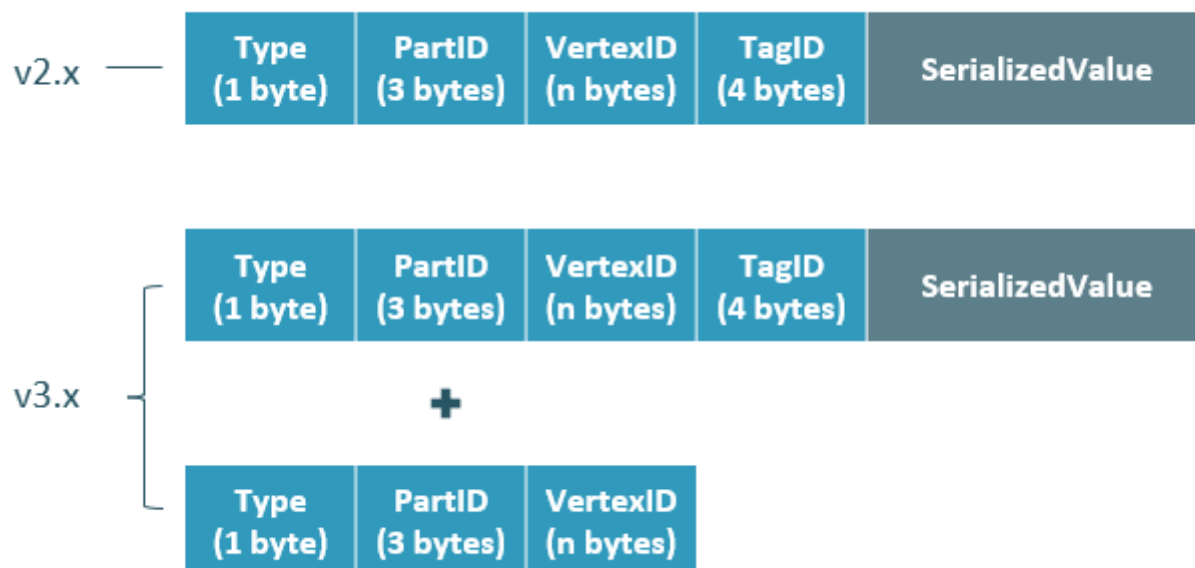
- NebulaGraph provides its own WAL mode so one can customize the WAL. Each partition owns its WAL.
- One NebulaGraph KVStore cluster supports multiple graph spaces, and each graph space has its own partition number and replica copies. Different graph spaces are isolated physically from each other in the same cluster.

Data storage structure

Graphs consist of vertices and edges. NebulaGraph uses key-value pairs to store vertices, edges, and their properties. Vertices and edges are stored in keys and their properties are stored in values. Such structure enables efficient property filtering.

- The storage structure of vertices

Different from NebulaGraph version 2.x, version 3.x added a new key for each vertex. Compared to the old key that still exists, the new key has no `TagID` field and no value. Vertices in NebulaGraph can now live without tags owing to the new key.



Field	Description
Type	One byte, used to indicate the key type.
PartID	Three bytes, used to indicate the sharding partition and to scan the partition data based on the prefix when re-balancing the partition.
VertexID	The vertex ID. For an integer VertexID, it occupies eight bytes. However, for a string VertexID, it is changed to <code>fixed_string</code> of a fixed length which needs to be specified by users when they create the space.
TagID	Four bytes, used to indicate the tags that vertex relate with.
SerializedValue	The serialized value of the key. It stores the property information of the vertex.

- The storage structure of edges

Type (1 byte)	PartID (3 bytes)	VertexID (n bytes)	EdgeType (4 bytes)	Rank (8 bytes)	VertexID (n bytes)	Placeholder (1 byte)	SerializedValue
------------------	---------------------	-----------------------	-----------------------	-------------------	-----------------------	-------------------------	-----------------

Field	Description
Type	One byte, used to indicate the key type.
PartID	Three bytes, used to indicate the partition ID. This field can be used to scan the partition data based on the prefix when re-balancing the partition.
VertexID	Used to indicate vertex ID. The former VID refers to the source VID in the outgoing edge and the dest VID in the incoming edge, while the latter VID refers to the dest VID in the outgoing edge and the source VID in the incoming edge.
Edge Type	Four bytes, used to indicate the edge type. Greater than zero indicates out-edge, less than zero means in-edge.
Rank	Eight bytes, used to indicate multiple edges in one edge type. Users can set the field based on needs and store weight, such as transaction time and transaction number.
Placeholder	One byte. Reserved.
SerializedValue	The serialized value of the key. It stores the property information of the edge.

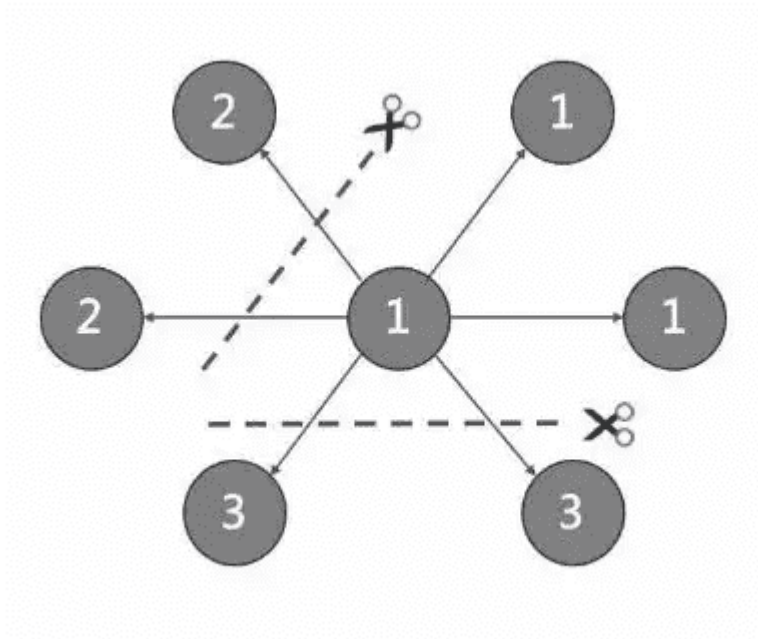
PROPERTY DESCRIPTIONS

NebulaGraph uses strong-typed Schema.

NebulaGraph will store the properties of vertex and edges in order after encoding them. Since the length of fixed-length properties is fixed, queries can be made in no time according to offset. Before decoding, NebulaGraph needs to get (and cache) the schema information in the Meta Service. In addition, when encoding properties, NebulaGraph will add the corresponding schema version to support online schema change.

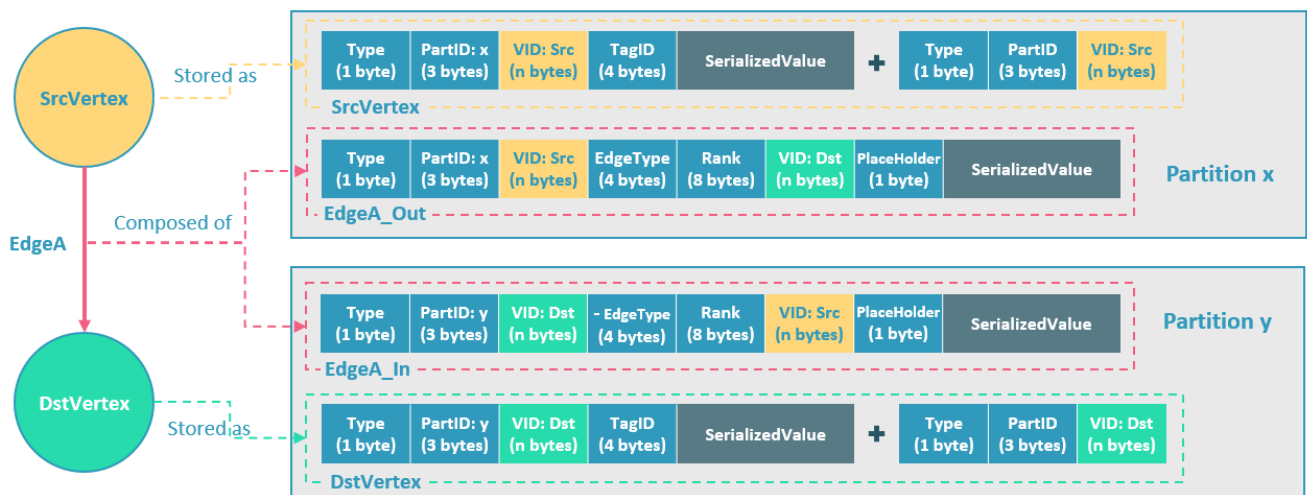
Data partitioning

Since in an ultra-large-scale relational network, vertices can be as many as tens to hundreds of billions, and edges are even more than trillions. Even if only vertices and edges are stored, the storage capacity of both exceeds that of ordinary servers. Therefore, NebulaGraph uses hash to shard the graph elements and store them in different partitions.



EDGE PARTITIONING AND STORAGE AMPLIFICATION

In NebulaGraph, an edge corresponds to two key-value pairs on the hard disk. When there are lots of edges and each has many properties, storage amplification will be obvious. The storage format of edges is shown in the figure below.



In this example, SrcVertex connects DstVertex via EdgeA, forming the path of (SrcVertex)-[EdgeA]->(DstVertex). SrcVertex, DstVertex, and EdgeA will all be stored in Partition x and Partition y as four key-value pairs in the storage layer. Details are as follows:

- The key value of SrcVertex is stored in Partition x. Key fields include Type, PartID(x), VID(Src), and TagID. SerializedValue, namely Value, refers to serialized vertex properties.
- The first key value of EdgeA, namely EdgeA_Out, is stored in the same partition as the SrcVertex. Key fields include Type, PartID(x), VID(Src), EdgeType(+ means out-edge), Rank(0), VID(Dst), and Placeholder. SerializedValue, namely Value, refers to serialized edge properties.
- The key value of DstVertex is stored in Partition y. Key fields include Type, PartID(y), VID(Dst), and TagID. SerializedValue, namely Value, refers to serialized vertex properties.
- The second key value of EdgeA, namely EdgeA_In, is stored in the same partition as the DstVertex. Key fields include Type, PartID(y), VID(Dst), EdgeType(- means in-edge), Rank(0), VID(Src), and Placeholder. SerializedValue, namely Value, refers to serialized edge properties, which is exactly the same as that in EdgeA_Out.

EdgeA_Out and EdgeA_In are stored in storage layer with opposite directions, constituting EdgeA logically. EdgeA_Out is used for traversal requests starting from SrcVertex, such as (a)-[]->(); EdgeA_In is used for traversal requests starting from DstVertex, such as ()-[]->(a).

Like EdgeA_Out and EdgeA_In, NebulaGraph redundantly stores the information of each edge, which doubles the actual capacities needed for edge storage. The key corresponding to the edge occupies a small hard disk space, but the space occupied by Value is proportional to the length and amount of the property value. Therefore, it will occupy a relatively large hard disk space if the property value of the edge is large or there are many edge property values.

PARTITION ALGORITHM

NebulaGraph uses a **static Hash** strategy to shard data through a modulo operation on vertex ID. All the out-keys, in-keys, and tag data will be placed in the same partition. In this way, query efficiency is increased dramatically.

Note

The number of partitions needs to be determined when users are creating a graph space since it cannot be changed afterward. Users are supposed to take into consideration the demands of future business when setting it.

When inserting into NebulaGraph, vertices and edges are distributed across different partitions. And the partitions are located on different machines. The number of partitions is set in the CREATE SPACE statement and cannot be changed afterward.

If certain vertices need to be placed on the same partition (i.e., on the same machine), see [Formula/code](#).

The following code will briefly describe the relationship between VID and partition.

```
// If VertexID occupies 8 bytes, it will be stored in int64 to be compatible with the version 1.0.
uint64_t vid = 0;
if (id.size() == 8) {
    memcpy(static_cast<void*>(&vid), id.data(), 8);
} else {
    MurmurHash2 hash;
    vid = hash(id.data());
}
PartitionID pId = vid % numParts + 1;
```

Roughly speaking, after hashing a fixed string to int64, (the hashing of int64 is the number itself), do modulo, and then plus one, namely:

```
pId = vid % numParts + 1;
```

Parameters and descriptions of the preceding formula are as follows:

Parameter	Description
%	The modulo operation.
numParts	The number of partitions for the graph space where the VID is located, namely the value of partition_num in the CREATE SPACE statement.
pId	The ID for the partition where the VID is located.

Suppose there are 100 partitions, the vertices with VID 1, 101, and 1001 will be stored on the same partition. But, the mapping between the partition ID and the machine address is random. Therefore, we cannot assume that any two partitions are located on the same machine.

Raft

RAFT IMPLEMENTATION

In a distributed system, one data usually has multiple replicas so that the system can still run normally even if a few copies fail. It requires certain technical means to ensure consistency between replicas.

Basic principle: Raft is designed to ensure consistency between replicas. Raft uses election between replicas, and the (candidate) replica that wins more than half of the votes will become the Leader, providing external services on behalf of all replicas. The rest Followers will play backups. When the Leader fails (due to communication failure, operation and maintenance commands, etc.), the rest Followers will conduct a new round of elections and vote for a new Leader. The Leader and Followers will detect each other's survival through heartbeats and write them to the hard disk in Raft-wal mode. Replicas that do not respond to more than multiple heartbeats will be considered faulty.

Note

Raft-wal needs to be written into the hard disk periodically. If hard disk bottlenecks to write, Raft will fail to send a heartbeat and conduct a new round of elections. If the hard disk IO is severely blocked, there will be no Leader for a long time.

Read and write: For every writing request of the clients, the Leader will initiate a Raft-wal and synchronize it with the Followers. Only after over half replicas have received the Raft-wal will it return to the clients successfully. For every reading request of the clients, it will get to the Leader directly, while Followers will not be involved.

Failure: Scenario 1: Take a (space) cluster of a single replica as an example. If the system has only one replica, the Leader will be itself. If failure happens, the system will be completely unavailable. Scenario 2: Take a (space) cluster of three replicas as an example. If the system has three replicas, one of them will be the Leader and the rest will be the Followers. If the Leader fails, the rest two can still vote for a new Leader (and a Follower), and the system is still available. But if any of the two Followers fails again, the system will be completely unavailable due to inadequate voters.

Note

Raft and HDFS have different modes of duplication. Raft is based on a quorum vote, so the number of replicas cannot be even.

MULTI GROUP RAFT

The Storage Service supports a distributed cluster architecture, so NebulaGraph implements Multi Group Raft according to Raft protocol. Each Raft group stores all the replicas of each partition. One replica is the leader, while others are followers. In this way, NebulaGraph achieves strong consistency and high availability. The functions of Raft are as follows.

NebulaGraph uses Multi Group Raft to improve performance when there are many partitions because Raft-wal cannot be NULL. When there are too many partitions, costs will increase, such as storing information in Raft group, WAL files, or batch operation in low load.

There are two key points to implement the Multi Raft Group:

- To share transport layer

Each Raft Group sends messages to its corresponding peers. So if the transport layer cannot be shared, the connection costs will be very high.

- To share thread pool

Raft Groups share the same thread pool to prevent starting too many threads and a high context switch cost.

BATCH

For each partition, it is necessary to do a batch to improve throughput when writing the WAL serially. As NebulaGraph uses WAL to implement some special functions, batches need to be grouped, which is a feature of NebulaGraph.

For example, lock-free CAS operations will execute after all the previous WALs are committed. So for a batch, if there are several WALs in CAS type, we need to divide this batch into several smaller groups and make sure they are committed serially.

TRANSFER LEADERSHIP

Transfer leadership is extremely important for balance. When moving a partition from one machine to another, NebulaGraph first checks if the source is a leader. If so, it should be moved to another peer. After data migration is completed, it is important to [balance leader distribution](#) again.

When a transfer leadership command is committed, the leader will abandon its leadership and the followers will start a leader election.

PEER CHANGES

To avoid split-brain, when members in a Raft Group change, an intermediate state is required. In such a state, the quorum of the old group and new group always have an overlap. Thus it prevents the old or new group from making decisions unilaterally. To make it even simpler, in his doctoral thesis Diego Ongaro suggests adding or removing a peer once to ensure the overlap between the quorum of the new group and the old group. NebulaGraph also uses this approach, except that the way to add or remove a member is different. For details, please refer to `addPeer/removePeer` in the `RaftPart` class.

Differences with HDFS

The Storage Service is a Raft-based distributed architecture, which has certain differences with that of HDFS. For example:

- The Storage Service ensures consistency through Raft. Usually, the number of its replicas is odd to elect a leader. However, `DataNode` used by HDFS ensures consistency through `NameNode`, which has no limit on the number of replicas.
- In the Storage Service, only the replicas of the leader can read and write, while in HDFS all the replicas can do so.
- In the Storage Service, the number of replicas needs to be determined when creating a space, since it cannot be changed afterward. But in HDFS, the number of replicas can be changed freely.
- The Storage Service can access the file system directly. While the applications of HDFS (such as HBase) have to access HDFS before the file system, which requires more RPC times.

In a word, the Storage Service is more lightweight with some functions simplified and its architecture is simpler than HDFS, which can effectively improve the read and write performance of a smaller block of data.

Last update: November 3, 2023

3. Quick start

3.1 Quickly deploy NebulaGraph using Docker

You can quickly get started with NebulaGraph by deploying NebulaGraph with Docker Desktop or Docker Compose.

Using Docker Desktop Using Docker Compose

NebulaGraph is available as a [Docker Extension](#) that you can easily install and run on your Docker Desktop. You can quickly deploy NebulaGraph using Docker Desktop with just one click.

1. Install Docker Desktop

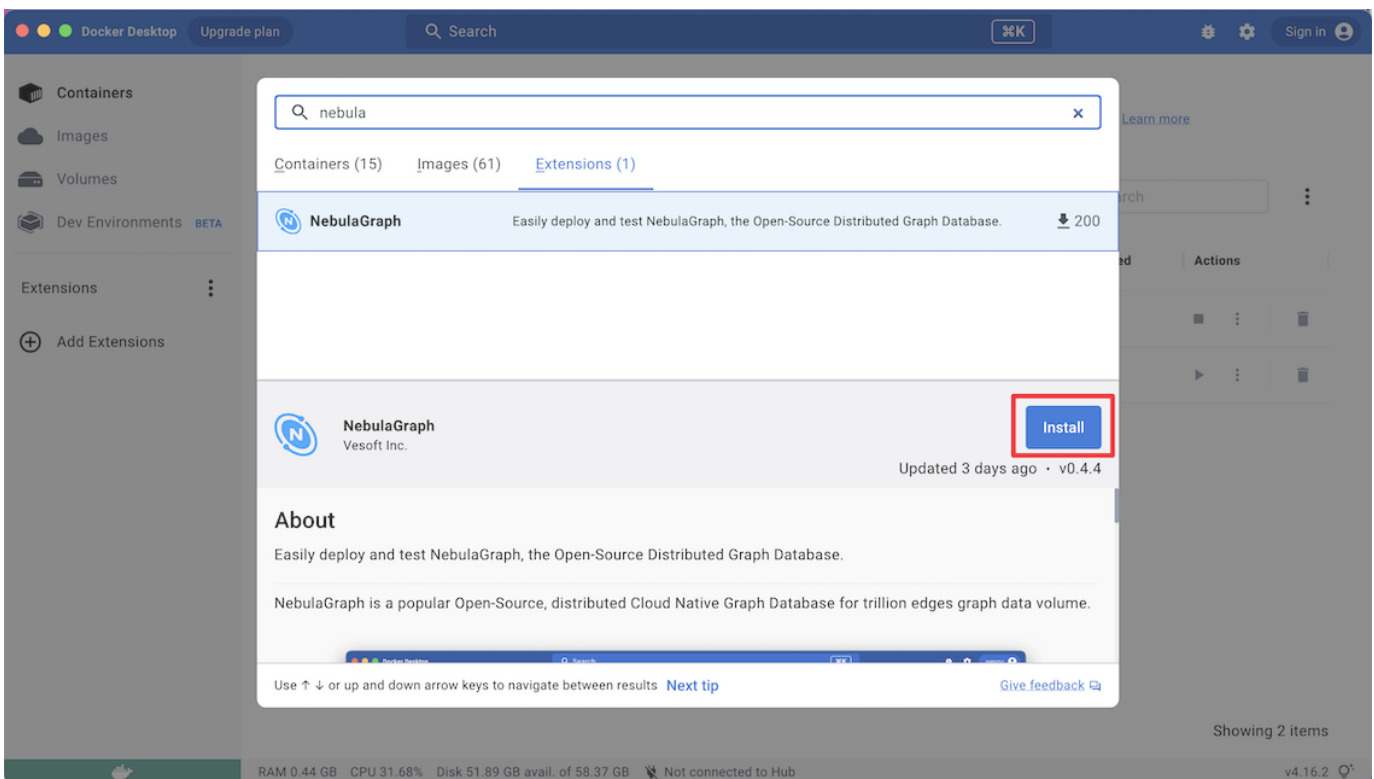
- [Install Docker Desktop on Mac](#)
- [Install Docker Desktop on Windows](#)

Caution

To install Docker Desktop, you need to install [WSL 2](#) first.

2. In the left sidebar of Docker Desktop, click **Extensions** or **Add Extensions**.

3. On the Extensions Marketplace, search for NebulaGraph and click **Install**.



Click **Update** to update NebulaGraph to the latest version when a new version is available.

