

Fujitsu Enterprise Postgres 18
for x86

Knowledge Data Management User's Guide

Windows/Linux

Preface

Purpose of this document

This document describes the knowledge data management features of Fujitsu Enterprise Postgres.

Intended readers

This document is intended for those who use the knowledge data management feature.

To read this document, you need to have the following knowledge:

- Fujitsu Enterprise Postgres
- PostgreSQL

Structure of this document

The structure and content of this manual is shown below.

[Chapter 1 Knowledge Data Management Feature](#)

This chapter describes the overview of the knowledge data management feature and examples of how knowledge data can be used.

[Chapter 2 Vector Data Management Feature](#)

This chapter describes the vector data management feature.



[Chapter 3 Semantic Text Search and Automatic Vectorization Feature](#)

This chapter describes the semantic text search and the automatic vectorization feature for semantic search.



[Chapter 4 Graph Management Feature](#)

This chapter describes the graph management feature.



[Chapter 5 Model Management in the Database](#)

This chapter describes the model management in the database.

Export restrictions

Exportation/release of this document may require necessary procedures in accordance with the regulations of your resident country and/or US export control laws.

Issue date and version

Edition 1.0: December 2025

Copyright

Copyright 2025 Fujitsu Limited

Contents

Chapter 1 Knowledge Data Management Feature.....	1
1.1 Overview of the Knowledge Data Management Feature.....	1
1.2 Examples of Using Knowledge Data.....	3
1.2.1 Example of Searching Text Data Based on Semantic Similarity.....	3
1.2.2 Example of Searching a Graph Based on Relationships.....	3
Chapter 2 Vector Data Management Feature.....	5
2.1 Setting Up the Vector Data Management Feature.....	5
2.2 Storing and Searching Vector Data.....	5
2.3 Protecting Vector Data.....	5
2.4 Performance Tuning for Similar Search of Vector Data.....	5
2.5 Using Vector Data by Applications.....	6
2.6 Quantitative Limits.....	6
Chapter 3 Semantic Text Search and Automatic Vectorization Feature.....	7
3.1 Overview of the Semantic Text Search and Automatic Vectorization Feature.....	7
3.2 Installation of Semantic Text Search and Automatic Vectorization.....	8
3.2.1 Operating Environment.....	8
3.2.2 Setup.....	8
3.2.2.1 Setting Up pgai.....	8
3.2.2.2 Setting Up pgx_vectorizer.....	8
3.2.2.3 Migration from Fujitsu Enterprise Postgres 17 SP1 and 17 SP2.....	9
3.2.3 Removing.....	10
3.2.4 Stopping the vectorize scheduler.....	10
3.2.5 Credentials Protection.....	10
3.2.5.1 Encrypting Credentials.....	10
3.2.5.2 Restrict Access to Credentials.....	11
3.3 Preparation for Semantic Text Search.....	11
3.3.1 Configuring Embedded Providers (for Workers).....	11
3.3.2 Configuring Embedded Providers (for Semantic Text Search).....	11
3.3.3 Definition of Vectorization.....	11
3.3.4 Granting Privilege to Execute Functions.....	14
3.4 Storing Vector Data for Semantic Text Search.....	14
3.5 Protecting Vector Data for Semantic Text Search.....	14
3.5.1 Encrypting Vector Data for Semantic Text Search.....	14
3.5.2 Restricting Access to Vector Data for Semantic Text Search.....	15
3.5.3 Recording Access to Vector Data for Semantic Text Search.....	17
3.6 Monitoring Vectorization Processing for Semantic Text Search.....	17
3.6.1 Checking the Vectorization Queue.....	17
3.6.2 Checking the Status of Vectorization Processing.....	18
3.6.3 Checking the Scheduler for Vectorization Processing.....	18
3.7 Temporarily Disabling Vectorization Processing for Semantic Text Search.....	19
3.8 Semantic Text Search.....	19
3.9 Changing the Vector Representation Used in Semantic Text Search.....	19
3.10 Performance Tuning of Semantic Text Search.....	20
3.11 Hybrid Search.....	20
3.12 Improving the Accuracy of Hybrid Search.....	21
3.12.1 Overview of Hybrid Search Tuning.....	21
3.12.2 Recording and Deleting Traces of Hybrid Search.....	22
3.12.3 Calculation of Search Accuracy Using External Tools.....	23
3.12.4 Calculation of Search Accuracy in the Database.....	23
3.12.4.1 Example of Calculating Search Accuracy in a Database.....	24
3.12.5 Tuning of Hybrid Search.....	25
3.13 Reference.....	26
3.13.1 Vectorization Functions.....	26
3.13.1.1 Defining Vectorization.....	26

3.13.1.2 Vectorization Schedule.....	26
3.13.1.3 Vectorizer Management Functions.....	26
3.13.2 Embedded Provider Management Functions.....	28
3.13.3 Search Functions.....	29
3.13.3.1 Semantic Text Search.....	29
3.13.3.2 Hybrid Search.....	29
3.13.3.3 Details of the pgx_hybrid_search Function.....	31
3.13.4 Tables/Views Created by Semantic Text Search and Automatic Vectorization Feature.....	34
3.13.5 Parameters.....	37
3.13.6 Evaluation of Knowledge Data Search.....	37
3.13.6.1 Concept of Evaluation for Knowledge Data Search.....	37
3.13.6.2 Evaluation Value per Record and Search Accuracy.....	38
3.13.6.3 Tuning the Combination Method of Hybrid Search.....	39
Chapter 4 Graph Management Feature.....	40
4.1 Overview of Graph Management Feature.....	40
4.2 Installation of Graph Management Feature.....	40
4.2.1 Setting Up the Graph Management Feature.....	40
4.2.2 Removing the Graph Management Feature.....	40
4.3 Creating a Graph.....	40
4.4 Storing Graph Data.....	41
4.5 Protecting Graph Data.....	41
4.5.1 Encrypting the Graph.....	41
4.5.2 Restricting Access to Graph.....	42
4.5.3 Recording Access to Graph.....	43
4.6 Searching Graph.....	43
4.7 Adding Labels to Graph.....	44
4.8 Performance Tuning of Graph Search.....	44
4.9 Using Graph Data in Applications.....	45
4.10 Visualizing Graph Data.....	45
4.11 Internal Structure of Graph Data.....	45
4.12 Quantitative Limits.....	45
4.13 Reference.....	45
Chapter 5 Model Management in the Database.....	46
5.1 Overview.....	46
5.2 Introduction/Setup.....	46
5.2.1 Setting Up Inference Server.....	46
5.2.2 Setting Up pgx_inference.....	47
5.2.3 Removing.....	48
5.2.4 Resource Control.....	49
5.2.5 Setting Privilege.....	49
5.2.6 Model Preparation.....	52
5.2.6.1 Database-side Model Specifications.....	52
5.2.7 Model Import.....	52
5.2.8 Model Load/Unload.....	53
5.2.9 Definition of Vectorization.....	53
5.2.10 Deletion of Imported Model.....	53
5.3 Usage Example.....	53
5.3.1 Model Import.....	53
5.3.2 Model Load/Unload.....	54
5.3.3 Definition of Vectorization.....	55
5.3.4 Deletion of Imported Model.....	55
5.4 Operation.....	55
5.4.1 Backup/Restore.....	55
5.4.2 Streaming Replication/Database Multiplexing.....	56
5.4.3 Security.....	57
5.4.4 Monitoring/Tuning.....	57

5.4.4.1 Monitoring Target.....	57
5.4.4.2 Parameter Tuning.....	57
5.5 Reference.....	58
5.5.1 Functions.....	58
5.5.2 Tables/Views.....	59
5.5.3 Parameters.....	60

Chapter 1 Knowledge Data Management Feature

This section describes an overview of the knowledge data management feature and examples of using knowledge data.

1.1 Overview of the Knowledge Data Management Feature

The Knowledge Data Management feature allows you to search based on semantic relationships using vectors and graphs, and manage those data.

When you use a large language model (LLM) in a retrieval-augmented generation (RAG) approach, you need external knowledge data to give to LLM. Fujitsu Enterprise Postgres can manage knowledge data in vector and graph formats, eliminating the need for a dedicated database for each format. You can also use Fujitsu Enterprise Postgres database multiplexing, access control, data encryption, and other features to securely manage your knowledge data.

In addition to existing searches, the Knowledge Data Management feature enables you to search knowledge data in three ways:

- Similar search of vector data

Performs a similarity search between the specified vector data and the vector data stored in Fujitsu Enterprise Postgres.



- Searching text data based on semantic similarity

Searches for semantic similarity between the specified text and the text stored in Fujitsu Enterprise Postgres using embedding in a semantic vector.

When you store text data for semantic search in Fujitsu Enterprise Postgres, a vector (semantic vector) that captures the semantic similarity of the text data is automatically generated and stored in the database. When searching, the specified text data is automatically converted to vector data, and semantic search is performed by the vector similarity search.

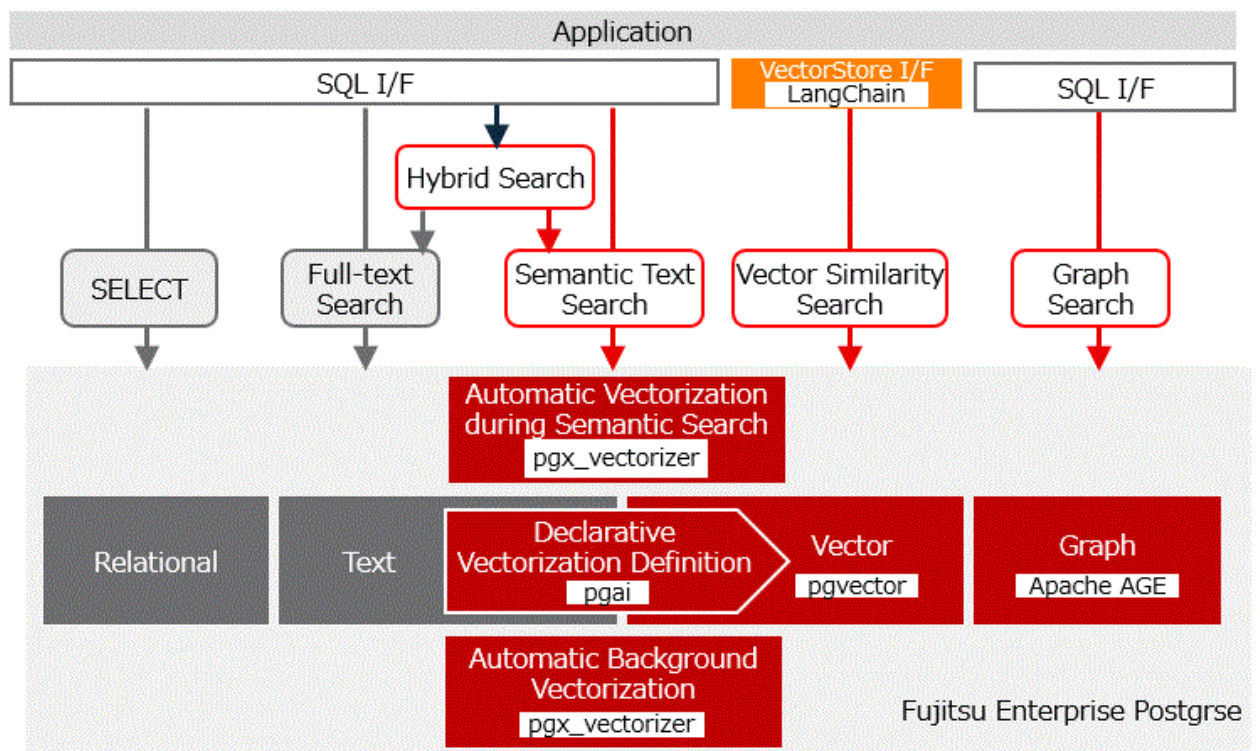
A hybrid search that combines text semantic similarity search and full-text search can also be utilized.



- Searching graphs based on relationships

A graph is a data structure that represents entities and their relationships as nodes and the edges that connect them. You can search graphs based on properties and relationships.

Knowledge data can be accessed by executing SQL queries from the application to Fujitsu Enterprise Postgres. You can also use LangChain, an AI application development framework, to access knowledge data in Fujitsu Enterprise Postgres.



What the knowledge data management feature can do

The Knowledge Data Management feature enables you to:

- Vector data storage and retrieval using the vector data management feature
 - Storage of float32, float16, bit vectors, and sparse vectors
 - Vector neighbor search by cosine or euclidean distance
 - HNSW and IVFFlat vector index
 - DiskANN-based vector index suitable for large data sets
- Semantic text search and automatic vectorization
 - Automated background vectorization of newly added text data
 - Semantic search with automatic query vectorization, consistent with stored vector representation
 - Use of external vector embedding models such as Ollama and OpenAI
 - Hybrid search combining semantic text search and full-text search
- Graph data storage and retrieval using the graph data management feature
 - Storing graphs
 - Exploring and updating graphs with openCypher
- Model management in the database
 - Management of models in ONNX format
 - Vectorization using ONNX format models in conjunction with Triton Inference Server
- Application development support
 - LangChain linkage

For access from Python applications using LangChain, refer to the technical documentation available at:

<https://www.postgresql.fastware.com/resource-center>

1.2 Examples of Using Knowledge Data

Here are some examples of using knowledge data.

1.2.1 Example of Searching Text Data Based on Semantic Similarity

Here is an example of a typical RAG-based application that utilizes vector embedding of text data. Instead of performing vectorization on the application side, we use the semantic search feature for text data provided by Fujitsu Enterprise Postgres.

1. Creating text data

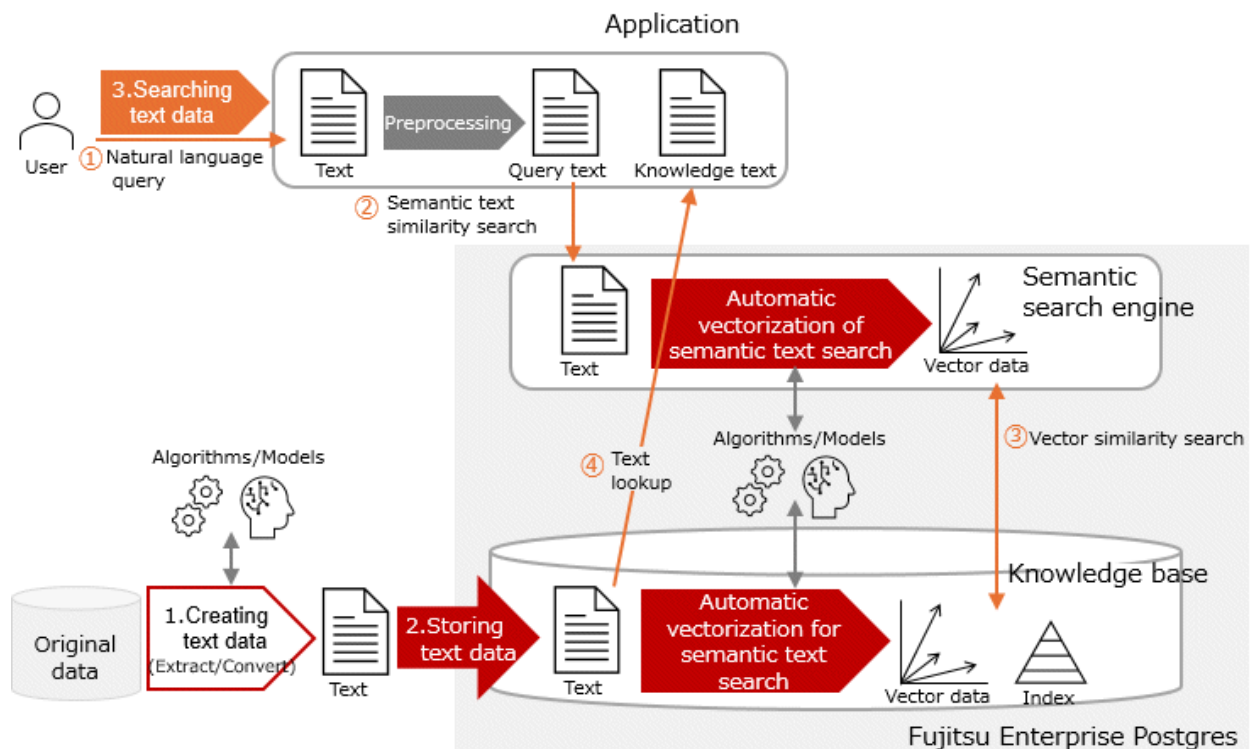
Use different algorithms and models to extract text from different data in your organization.

2. Storing text data in Fujitsu Enterprise Postgres

Store the text data in Fujitsu Enterprise Postgres. To use the semantic search function in text, declaratively define the vector representation to be used for the stored textual knowledge data. This automatically creates the vector data necessary for semantic search.

3. Searching text data

The text semantic search feature returns text that is semantically similar to the text you are searching for. No conversion to vector representation is required on the application side.



1.2.2 Example of Searching a Graph Based on Relationships

Here is an example of a graph as knowledge data is queried in natural language, and the obtained knowledge is converted into text for use.

1. Creating graph data

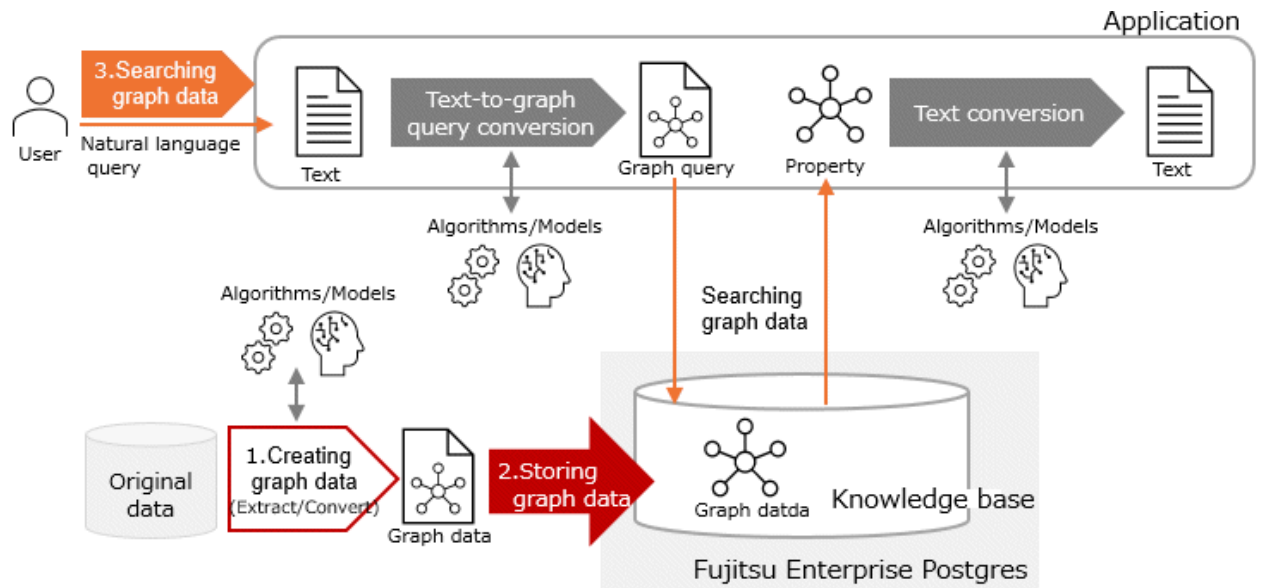
Extract graph data from a variety of data in your organization, such as using a relation extraction model.

2. Storing graph data in Fujitsu Enterprise Postgres

Store the created graph data in Fujitsu Enterprise Postgres. Nodes and edges can be added with Cypher queries using SQL functions.

3. Searching graph data

To search a graph, you specify a Cypher query to the SQL function. Converts the retrieved properties into a text representation for use.



Chapter 2 Vector Data Management Feature

This chapter describes the vector data management feature, which provides the storage and search of vector data.



pgvectorscale, described below, is only available on Linux.

2.1 Setting Up the Vector Data Management Feature

Vector data management features are provided by the OSS's pgvector and pgvectorscale.

Refer to "pgvector" in the Installation and Setup Guide for Server to set up pgvector.

If you use the StreamingDiskANN index, refer to "pgvectorscale" in the Installation and Setup Guide for Server to set up pgvectorscale.

In addition, to use the StreamingDiskANN index, a CPU that supports the AVX2 and FMA instruction sets is required. If you use a CPU that does not support AVX2 and FMA, an error will occur when creating a StreamingDiskANN index for a table that stores vector data, or when inserting data into a table with a StreamingDiskANN index created. Refer to the documentation for each CPU to see if the CPU supports AVX2 and FMA.

2.2 Storing and Searching Vector Data

The vector data management feature provides a new vector data type for storing vector data. Create a table with columns of vector data type and store vector data.

The similarity search of vector data is performed by calculating the distance between two vector data using the distance operator added by the vector data management feature, and by the closeness and ordering of the calculated distance.

Example) Vector similarity search

```
SELECT * FROM items ORDER BY embedding <-> '[3,1,2]' < 5 LIMIT 5;
```

Creating a vector index on a column of the vector data type causes the vector index to be used when searching for similarities using the distance operator.



Point

Similarity searches using vector index are approximate similarity searches.



See

Refer to the pgvector documentation for information about vector data types, vector operations, distance types between vectors, and HNSW and IVFFlat vector indexes.

Refer to the pgvectorscale documentation for information about StreamingDiskANN indexes.

2.3 Protecting Vector Data

Vector data is stored as column values in tables, and can be protected by encryption, backup, multiplexing, replication, access control, and audit settings that specify the table, tablespace, database, or instance in which the vector data is stored.

2.4 Performance Tuning for Similar Search of Vector Data

You can verify that indexes are being used for similar searches of vector data by checking the access plan in the EXPLAIN statement.

Example) Access Plan for Vector Similar Search

Here is an example using pgvectorscale:

```
EXPLAIN SELECT * FROM items ORDER BY embedding <=> '[3,1,2]' LIMIT 5;
               QUERY PLAN
-----
Limit  (cost=24.75..24.96 rows=5 width=33)
  ->   Index Scan using idx_diskann on items (cost=24.75..445.75 rows=10000 width=33)
        Order By: (embedding <=> '[3,1,2]':vector)
(3 rows)
```



Each vector index has parameters for tuning. See the documentation for pgvector and pgvectorscale for details.

2.5 Using Vector Data by Applications

If your application wants to work directly with vector data types, introduce a driver for each language in your application. If no driver is used, the vector data type is returned to the application as a string or array type.



For information on pgvector drivers for various languages, refer to below.

<https://github.com/pgvector/pgvector?tab=readme-ov-file#languages>

2.6 Quantitative Limits

Refer to the pgvector documentation for the quantitative limits of the data types provided by pgvector.

Chapter 3 Semantic Text Search and Automatic Vectorization Feature

This chapter describes the semantic text search and the automatic vectorization feature for semantic text search.

3.1 Overview of the Semantic Text Search and Automatic Vectorization Feature

Semantic text search

The semantic text search is a feature that searches for highly relevant texts based on semantic similarity of the text. This is achieved by utilizing vector representations that maintain the semantic similarity of text data. It is possible to use hybrid search that combines semantic text search using vector representation and full-text search based on string matching.

Automatic vectorization feature

The automatic vectorization feature automatically generates and stores vector representations corresponding to the text data inserted into the table. Use an external service called an embedded provider to generate vector representations of text.

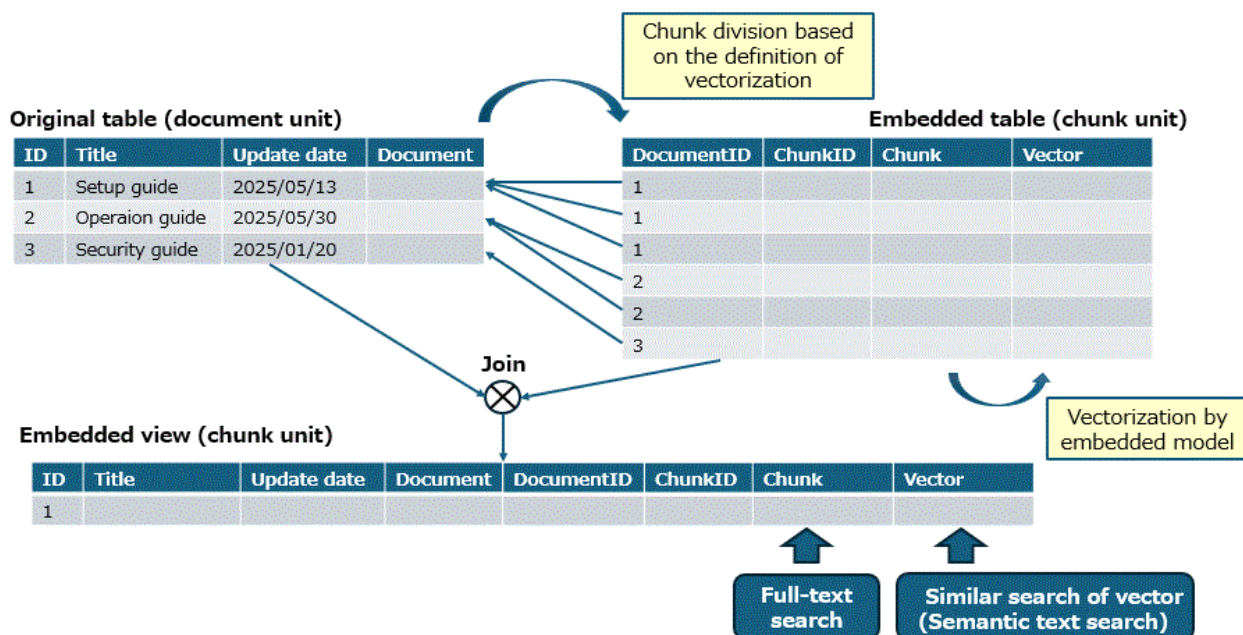
The automatic vectorization feature is executed in the background by defining vectorization for the table to be searched. The definition of vectorization includes the embedded model that determines the vector representation used for semantic text search, the method of text segmentation, and the definition of the index.

When a vectorization is defined for table, a corresponding "embedded table" is created. This embedded table stores text chunks, which are units of the original table's text data divided into chunks, and the corresponding vector data. At the same time as defining the vectorization, you can define a vector index for the vector data and a full-text search index for the text chunks. An "embedded view" that combines the original table and the embedded table is also created.

Semantic text search and embedded view

Semantic text search and hybrid search are executed using embedded view as the search target.

When semantic text search and hybrid search are performed for embedded view, the text given as a query is internally converted into a vector representation, and a similarity search is conducted between the vector representation of the stored knowledge data.



3.2 Installation of Semantic Text Search and Automatic Vectorization

3.2.1 Operating Environment

For details about the packages required for this feature to work, refer to the following sections in the Installation and Setup Guide for Server.

- Required Packages
- Related Software

3.2.2 Setup

Support for semantic text search and automatic vectorization is provided as an extension called `pgx_vectorizer`. `pgx_vectorizer` uses `pgai`, which relies on `pgvector` and `plpython3u`.

Refer to "pgvector" in the Installation and Setup Guide for Server to set up `pgvector`.

If you want to manage ONNX format models in a database and use them for vectorization, refer to "[Chapter 5 Model Management in the Database](#)".

If you are already using the `pgx_vectorizer` extension when migrating from Fujitsu Enterprise Postgres 17 SP1 and 17 SP2, refer to "[3.2.2.3 Migration from Fujitsu Enterprise Postgres 17 SP1 and 17 SP2](#)" to upgrade the `pgx_vectorizer` extension and the `pgai` extension.

3.2.2.1 Setting Up `pgai`

As a superuser, execute the following command to set up `pgai`. "<x>" indicates the product version.

```
$ su -  
# cp -r /opt/fsepv<x>server64/OSS/pgai-extension/* /opt/fsepv<x>server64/
```

Configure the following before starting the instance.

Set the <Python package installation destination> to the installation destination listed in the "Related Software" of the "Installation and Setup Guide for Server".

```
PYTHONPATH=/opt/fsepv<x>server64/psycopy/python3.11/site-packages/:<Python package  
installation destination>:$PYTHONPATH
```

3.2.2.2 Setting Up `pgx_vectorizer`

Setting parameters in the `postgresql.conf` file

Set the following parameters.

- Add `pgx_vectorizer` to the `shared_preload_libraries` parameter.
- Specify the maximum parallelism for vectorization processing in the `pgx_vectorizer.max_vectorize_worker` parameter.
- Add the following value to the value of the `max_worker_processes` parameter:
 - number of databases to enable `pgx_vectorizer` functionality + `pgx_vectorizer.max_vectorize_worker` + 2
- If you have changed the installation destination of the Fujitsu Enterprise Postgres Server feature to a location other than the standard installation destination, specify the following for the `pgx_vectorizer.pgai_worker_path` parameter.

```
<Fujitsu Enterprise Postgres server feature Installation Directory>/OSS/pgai-  
worker/bin/pgai
```

Enabling the `pgx_vectorizer` extension

Execute `CREATE EXTENSION` for the database that will use this feature.

Adding the CASCADE option will also enable the dependent pgai, pgvector, and plpython3u at the same time.

After CREATE EXTENSION, execute the start_vectorize_scheduler function to start the vectorize scheduler. The vectorize scheduler is a feature that schedules workers to perform vectorization.

Example) Connecting to the database "rag_database" using the psql command

```
rag_database=# CREATE EXTENSION IF NOT EXISTS pgx_vectorizer CASCADE;  
CREATE EXTENSION  
rag_database=# SELECT pgx_vectorizer.start_vectorize_scheduler(); -- Starting the  
vectorize scheduler
```

After enabling the extended features, update the postgresql.conf parameter settings using commands such as pg_ctl reload.

Create a database user and set up a connection

Create a database user that the automatic vectorization feature will use when converting vectors in the background, and register it as the user that will convert vectors. Specify a password to use password authentication.

```
CREATE ROLE <worker_user> PASSWORD `<worker password>` ... LOGIN;  
SELECT pgx_vectorizer.set_worker_setting('user', 'VECTORIZE_USER', '<worker_user>');
```

This database user will connect to Fujitsu Enterprise Postgres as an application, so modify the pg_hba.conf file for client authentication. Set it so that the database user created above can connect to the database that uses the pgx_vectorizer function from localhost using password authentication.

host	<ai-database>	<worker_user>	127.0.0.1/32	scram-sha-256
host	<ai-database>	<worker_user>	:::1/128	scram-sha-256

The background vectorization process runs with the privileges of the OS user that starts Fujitsu Enterprise Postgres. The password required for the above database user to connect to Fujitsu Enterprise Postgres is referenced from the password file of the OS user that starts Fujitsu Enterprise Postgres. Specify the information required to connect to the vectorization process. The password file used is in the default location. For details about the password file, refer to "The Password File" in the PostgreSQL Documentation.

3.2.2.3 Migration from Fujitsu Enterprise Postgres 17 SP1 and 17 SP2

If you have set up the pgx_vectorizer extension with Fujitsu Enterprise Postgres 17 SP1 and 17 SP2, you need to upgrade the pgx_vectorizer extension. Follow the steps below to perform the upgrade.

1. If the database server is running, stop the database server.

```
$ pg_ctl stop -D <data_directory>
```

2. Remove pgx_vectorizer from shared_preload_libraries parameter in postgresql.conf.
3. Upgrade the product.
4. Start the database server.

```
$ pg_ctl start -D <data_directory>
```

5. Update pgvector, pgai, and pgx_vectorizer.

```
ALTER EXTENSION vector UPDATE;  
ALTER EXTENSION ai UPDATE;  
ALTER EXTENSION pgx_vectorizer UPDATE;
```

6. Add pgx_vectorizer to shared_preload_libraries parameter in postgresql.conf.
7. Restart the database server.

```
$ pg_ctl restart -D <data_directory>
```


3.2.3 Removing

1. Connect to the database that is using this feature and execute DROP EXTENSION.

```
rag_database=# DROP EXTENSION pgx_vectorizer;  
DROP EXTENSION
```

2. As a superuser, execute the following command to remove pgai. "<x>" indicates the product version.

```
$ su -  
# rm /opt/fsepv<x>server64/<Files copied during setup>
```



Information

The files copied during setup can be checked below.

```
# find /opt/fsepv<x>server64/OSS/pgai-extension
```

3.2.4 Stopping the vectorize scheduler

To stop vectorize scheduler, refer to pg_stat_activity to check the pid of vectorize scheduler connected to the target database, and then use the pg_ctl kill TERM <pid> command or the SQL function pg_terminate_backend() to stop it. Because vectorize scheduler is continuously connected to the database, if you want to delete or change a database that has CREATE EXTENSIONed pgx_vectorizer, you must stop vectorize scheduler before performing the operation.

You can check the pid of vectorize scheduler connected to the current database by executing the following SQL. The backend_type of vectorize scheduler includes 'vectorize scheduler'.

```
SELECT pid FROM pg_stat_activity WHERE datid = (SELECT oid FROM pg_database WHERE datname =  
current_database()) AND backend_type LIKE '%vectorize scheduler%';
```

Below is an example of using pg_terminate_backend().

```
rag_database=>SELECT pg_terminate_backend(pid) FROM pg_stat_activity WHERE datid = (SELECT  
oid FROM pg_database WHERE datname = current_database()) AND backend_type LIKE '%vectorize  
scheduler%';
```

```
pg_terminate_backend  
-----  
t  
(1 row)
```

3.2.5 Credentials Protection

3.2.5.1 Encrypting Credentials

The pgx_vectorizer.worker_setting_table, which is created when pgx_vectorizer is set up, should be protected because it stores the API key that the worker performing vectorization uses to access the embedded provider that offers the vector embedding feature.

To encrypt this table, either encrypt the entire database to which the pgx_vectorizer feature is being added in advance, or move this table to an encrypted tablespace after adding the extension feature.

Example) When encrypting the entire database for which this feature is enabled

```
postgres=# CREATE DATABASE rag_database TABLESPACE = encrypted_tablespace;  
  
rag_database=> CREATE EXTENSION pgx_vectorizer;
```

Example) Moving to encrypted table space

```
rag_database=> ALTER TABLE pgx_vectorizer.worker_setting_table SET TABLESPACE
encrypted_tablespace;
```

3.2.5.2 Restrict Access to Credentials

The `pgx_vectorizer.worker_setting_table`, which stores the credential information, is configured so that only the user who executed the `CREATE EXTENSION` of the `pgx_vectorizer` feature can access it. If you want to allow access by other users, use the `GRANT` statement to grant access privilege to those users.

3.3 Preparation for Semantic Text Search

Semantic text search is achieved by similarity search between vectors (semantic vectors) that store the semantic similarity of text. Therefore, in order to use the semantic text search feature, you need to define vectorization (what kind of vector representation to use) for the text data to be searched and create vector data.

3.3.1 Configuring Embedded Providers (for Workers)

Before defining a vectorization, you must configure the worker process that will perform the vectorization to access the embedded provider.

Example) Settings when using Ollama as an embedded provider

```
rag_database=> SELECT pgx_vectorizer.set_worker_setting('ollama', 'OLLAMA_BASE_URL',
'http://your.ollama.server:11434');
```

Example) Settings when using OpenAI as an embedded provider

```
rag_database=> SELECT pgx_vectorizer.set_worker_setting('openai', 'OPENAI_API_KEY', 'your
api key');
```



Information

Automatic vectorization is performed by worker processes that are regularly initiated. Additionally, the worker processes are initiated by Fujitsu Enterprise Postgres.



Point

In this configuration, credentials for the worker processes to access the embedded provider are stored in an area accessible by Fujitsu Enterprise Postgres. Because the configuration information for the embedded provider, including the API key, is stored in the table as plain text, place the table in an encrypted tablespace and protect it with the Fujitsu Enterprise Postgres transparent data encryption feature. For instructions, refer to "[3.2.5 Credentials Protection](#)".

3.3.2 Configuring Embedded Providers (for Semantic Text Search)

To access embedded providers for semantic text search, use the `pgai` settings. Refer to the `pgai` documentation for details.

3.3.3 Definition of Vectorization

The definition of vectorization is done with the `pgx_vectorizer.pgx_create_vectorizer` function provided by `pgx_vectorizer`. In the vectorization definition, you can specify information about the table that contains the text data to be vectorized, the embedding model and dimensions directly related to vector representation, chunking processing, as well as specify the timing of vectorization as a schedule. If you want to immediately schedule automatic vectorization in the background with Fujitsu Enterprise Postgres, specify the `pgx_vectorizer.schedule_vectorizer` function. Additionally, if you want to specify the vectorization later or manually control the timing of vectorization, specify the `ai.scheduling_none` function.

When using hybrid search, you can simultaneously define a GIN or GiST index for full-text search used in hybrid search.

When using an embedded model managed in the database, specify `pgx_vectorizer.pgx_embedding_onnx` as the argument embedding for the `pgx_vectorizer.pgx_create_vectorizer` function. For more details, refer to "[5.2.9 Definition of Vectorization](#)".

When defining vectorization, an embedded view combining a table containing text data and a table containing vector data is defined. The embedded view is a view that adds the following columns to the original table, and for each piece of text data in the original table, it holds records divided into chunks.

Column	Type	Description
embedding_uuid	uuid	Text chunk ID
chunk_seq	integer	Sequential number of chunks in a text data
chunk	text	Text chunk
embedding	vector	Vector representation of the chunk

Information

In the following cases, refer to "[Definition of full-text search index](#)" to define the index.

- When using an index other than GIN or GiST in hybrid search
- When you have already performed only vectorization and want to add an index for full-text search later

Fujitsu Enterprise Postgres provides `pg_trgm` and `pg_bigm` as full-text search features using GIN and GiST indexes. When using these as full-text search features, set up `pg_trgm` and `pg_bigm` in advance. Refer to the following for setup instructions.

- `pg_trgm`
"Additional Supplied Modules and Extensions" in the PostgreSQL Documentation
- `pg_bigm`
"pg_bigm" in the Installation and Setup Guide for Server

Example) Definition of vectorization and full-text search index

Below is an example of simultaneously defining a vectorization and a GIN index for full-text search on the knowledge data table `sample_table` in text format.

```
rag_database=> SELECT pgx_vectorizer.pgx_create_vectorizer(
    'sample_table'::regclass,
    destination => 'sample_embeddings',
    embedding => ai.embedding_ollama('all-minilm', 384),
    chunking => ai.chunking_recursive_character_text_splitter('contents'),
    processing => ai.processing_default(batch_size => 200, concurrency => 1),
    scheduling => pgx_vectorizer.schedule_vectorizer(interval '1 hour'),
    indexing => ai.indexing_hnsw(min_rows => 50000, opclass => 'vector_cosine_ops'),
    fulltext_indexing => pgx_vectorizer.pgx_fulltext_indexing_gin(opclass =>
'gin_trgm_ops')
);
pgx_create_vectorizer
-----
1 - The ID of the created vectorizer
(1 row)
```

The definition of the full-text search index created by the `pgx_create_vectorizer` function can be checked in the `pgx_vectorizer.pgx_fulltext_index` table.

Note

Regarding the table containing text subject to vectorization, the following cannot be changed:

- Changes to the schema containing the table
- Changes to the table name
- Changes to the primary key, addition or removal of columns constituting the primary key
- Name and data type of columns included in the primary key
- Name and data type of columns containing text data subject to vectorization

Changing these may cause automatic vectorization or semantic text search to fail. If you want to change them, or if you have already changed them, you need to redefine the vectorization. In this case, re-vectorization of all text in the target table is required.

Point

Users can implement their own full-text search functionality, other than `pg_trgm` or `pg_bigm` provided by Fujitsu Enterprise Postgres, by defining indexes supported by that full-text search functionality and specifying search conditions using its operators, allowing for hybrid search with custom full-text search capabilities.

Definition of full-text search index

If you use an index other than GIN or GiST index for full-text search, manually define the index for the embedding table created after defining the vectorization. The name of the embedding table is the name with `"_store"` added to the end of the embedding view name specified in the destination argument when defining the vectorization. The column name for text chunks in the embedding table is `chunk`.

In the following example, a full-text search index `my_iam`, other than GIN and GiST, is defined in the embedded table `sample_embeddings2_store`.

```
CREATE INDEX sample_fti ON sample_embeddings2_store (chunk) USING my_iam;
```

Deleting full-text search index

You can delete the full-text search index by removing the information of the full-text search index from the vectorization and then executing the `DROP INDEX` statement.

```
rag_database=> SELECT
pgx_delete_fulltext_index_config(pgx_vectorizer.get_vectorizer_id(view_name =>
'sample_embeddings'));
rag_database=> SELECT indexname FROM pg_indexes where tablename =
'sample_embeddings_store';
indexname
-----
sample_embeddings_store_chunk_idx - It is created for the column called "chunk".

rag_database=> DROP INDEX sample_embeddings_store_chunk_idx;
```

Changing the full-text search index

When changing the full-text search index defined by the `pgx_create_vectorizer` function, first delete the full-text search index and then define a new index.

```
--Deleting full-text search index
rag_database=> SELECT
pgx_delete_fulltext_index_config(pgx_vectorizer.get_vectorizer_id(view_name =>
'sample_embeddings'));
rag_database=> SELECT indexname FROM pg_indexes where tablename =
```

```
'sample_embeddings_store';
indexname
-----
sample_embeddings_store_chunk_idx - It is created for the column called "chunk".

rag_database=> DROP INDEX sample_embeddings_store_chunk_idx;

-- Adding full-text search index
rag_database=> CREATE INDEX ON sample_embeddings_store USING gin (chunk gin_bigm_ops);
```

3.3.4 Granting Privilege to Execute Functions

When executing the semantic text search and full-text search, you must grant the executing user function privileges that exist in the ai schema.

```
GRANT EXECUTE ON ALL FUNCTIONS IN SCHEMA ai TO <user>;
```

3.4 Storing Vector Data for Semantic Text Search

Vectorization for text semantic text search is performed automatically according to the vectorization definition, and is saved as a vector data type column value in an internally created table so that the same text does not need to be vectorized again.

Vector data is generated and stored in the background according to the schedule specified when defining the vectorization. If new text data is added to the table on which the vectorization is defined, the corresponding vector data is automatically added asynchronously.



Information

Until the vector data has been created, the text data cannot be used for semantic text search.

3.5 Protecting Vector Data for Semantic Text Search

Vector data for semantic text search is automatically protected by backup, multiplexing, or replication settings that specify the tablespace, database, or instance in which it is stored.

This section describes the points to keep in mind when using Fujitsu Enterprise Postgres features to protect vector data for semantic text search.

Vectorization data using the semantic vector embedding model is a conversion that corresponds the semantic similarity between the original data to the closeness of the distance between vectors. Because there is a risk that the meaning of the original data can be inferred from vector data, it should be protected at the same level as the original data.

Multiple database objects are created by the vectorization definition. These database objects are the targets of protection.

3.5.1 Encrypting Vector Data for Semantic Text Search

Several database objects are created for semantic text search, and if you want to encrypt these database objects together, place the entire database in an encryption tablespace.

Example) Encrypting the entire database

```
postgres=# CREATE DATABASE rag_database TABLESPACE = encrypted_tablespace;

rag_database=> CREATE EXTENSION pgx_vectorizer;
```

To encrypt database objects for semantic text search when a single tablespace cannot be used per database, temporarily change the default tablespace to an encrypted tablespace before defining vectorization.

Indexes for vector data are created when the index creation conditions are met. To encrypt an index, change the tablespace after index creation, or create the index manually without specifying an index in the create_vectorizer function. For

information about tables containing vector data created by this feature, refer to "3.13.4 Tables/Views Created by Semantic Text Search and Automatic Vectorization Feature".

3.5.2 Restricting Access to Vector Data for Semantic Text Search

Access to vector data is restricted by controlling access to the tables in which the vector data is stored. Access restrictions are set using the confidentiality management feature of Fujitsu Enterprise Postgres.

The following is an example of using the confidentiality management feature to grant reference rights to the table sample, which contains text data, and the table sample_embeddings_store, which contains vector data generated by this feature.

1. Refer to the "Confidentiality Management" in the Security Operations Guide to define confidentiality management role, confidentiality matrix, confidentiality level, and confidentiality group.
2. Grant confidentiality privilege for the table to the confidentiality group.

```
SELECT pgx_grant_confidential_privilege('rag_matrix', 'level1', 'group1', '{"schema": ["ALL"], "table": ["SELECT"]}');
```

3. Add tables containing text data and embedded tables as confidentiality object to the confidentiality level.

```
SELECT pgx_add_object_to_confidential_level ('rag_matrix', 'level1',
'[{
  "type": "table",
  "object": [
    {
      "schema": "public",
      "table": ["sample"]
    },
    {
      "schema": "public",
      "table": ["sample_embeddings_store"]
    }
  ]
}]');
```

4. Add the role to the confidentiality group.

```
SELECT pgx_add_role_to_confidential_group('rag_matrix', 'group1', '{"rag_user"}');
```



Information

A vector table contains foreign keys, vector data, and chunks, but it is not necessary to set access privileges on a column-by-column basis; setting access privileges on a table-by-table basis is sufficient.

The following is an example of setting row-level security for table sample, which contains text data, and table sample_embeddings_store, which contains vector data generated by this function. In this example, sample and sample_embeddings_store contain user names in a column called username, and sample_embeddings_store has the primary key (id) of sample as a foreign key (id). When using row-level security, grant the BYPASSRLS attribute to the user set in VECTORIZE_USER.

1. Refer to the "Confidentiality Management" in the Security Operations Guide to define confidentiality management role, confidentiality matrix, confidentiality level, and confidentiality group.
2. Grant confidentiality privilege for the rowset to the confidentiality group.

```
SELECT pgx_grant_confidential_privilege('rag_matrix', 'level1', 'group1', '{"table": ["SELECT"], "schema": ["ALL"], "rowset": ["SELECT"]}');
```


3. Enable row-level security for the table.

```
ALTER TABLE sample ENABLE ROW LEVEL SECURITY;
ALTER TABLE sample_embeddings_store ENABLE ROW LEVEL SECURITY;
```

4. Ensure that the tables on which the target view is based are checked against the privileges of the view's user.

```
ALTER VIEW sample_embeddings SET (security_invoker = true);
```

5. Add the ai schema, tables containing text data and embedded tables and embedded views as confidentiality objects to the confidentiality level.

```
SELECT pgx_add_object_to_confidential_level ('rag_matrix', 'level1',
'[{
  "type": "schema",
  "object": [
    { "schema": "ai" },
    { "schema": "pgx_vectorizer" }
  ]
},
{
  "type": "table",
  "object": [
    {
      "schema": "public",
      "table": ["sample", "sample_embeddings_store", "sample_embeddings"]
    },
    {
      "schema": "ai",
      "table": ["vectorizer"]
    }
  ]
}]');
```

6. Add the table containing text data and the rowset of the embedded table as confidentiality objects to the confidentiality level.

```
SELECT pgx_add_object_to_confidential_level ('rag_matrix', 'level1',
'[{
  "type": "rowset",
  "object": [
    {
      "schema": "public",
      "table": "sample",
      "rowset_name": "rowset1",
      "rowset_expression": [
        {
          "as": "permissive",
          "using": "username = current_user"
        }
      ]
    },
    {
      "schema": "public",
      "table": "sample_embeddings_store",
      "rowset_name": "rowset1",
      "rowset_expression": [
        {
          "as": "permissive",
          "using": "EXISTS (SELECT 1 FROM public.sample WHERE public.sample.id =
public.sample_embeddings_store.id AND public.sample.\"username\" = current_user)"
        }
      ]
    }
  ]
}]');
```

```
]
}]);
```

7. Add the role to the confidentiality group.

```
SELECT pgx_add_role_to_confidential_group('rag_matrix', 'group1', '["rag_user"]');
```

3.5.3 Recording Access to Vector Data for Semantic Text Search

The audit log feature of Fujitsu Enterprise Postgres is used to record access to vector data in the audit log. When specifying individual database objects to be audited, also specify the tables that store vector data as audit targets.

For information about tables containing vector data created by this feature, refer to ["3.13.4 Tables/Views Created by Semantic Text Search and Automatic Vectorization Feature"](#).

3.6 Monitoring Vectorization Processing for Semantic Text Search

If you use automatic background vectorization in your database, the corresponding vector data will be automatically changed when text data is added, updated, or deleted. However, because the generation of vector data for new data is performed asynchronously, newly added data will not be reflected immediately in semantic similarity search. Also, if you use an external service to create vector data, you will need to check for errors.

3.6.1 Checking the Vectorization Queue

You can check the number of texts waiting to be vectorized for each vectorization definition by referencing the `ai.vectorizer_status` view, and you can check the number for a specific vectorization definition by using the `ai.vectorizer_queue_pending` function.

Example) Check the `ai.vectorizer_status` view

```
rag_database=> SELECT * FROM ai.vectorizer_status;
-[ RECORD 1 ]-----
id           | 1
source_table | public.sample_table
target_table | public.sample_embeddings_store
view         | public.sample_embeddings
pending_items | 1000
disabled     | f
```

Example) Check with the `ai.vectorizer_queue_pending` function

```
SELECT ai.vectorizer_queue_pending( pgx_vectorizer.get_vectorizer_id(view_name =>
'sample_embeddings') );
-[ RECORD 1 ]-----+--
vectorizer_queue_pending | 1000
```

If this value remains at 0 or a low value, you can determine that the vectorization process is on time. If this value tends to increase beyond the execution interval specified in the schedule or data addition interval, refer to [3.6.2 Checking the Status of Vectorization Processing](#) to check whether an error has occurred in the vectorization process, and refer to [3.6.3 Checking the Scheduler for Vectorization Processing](#) to check whether the vectorize scheduler is running.

If no errors have occurred, the vectorization processing speed is likely slow compared to the data addition speed. If the load is temporarily high, start a temporary vectorization process with the `pgx_vectorizer.run_vectorize_worker` function. If not, change the schedule with the `pgx_vectorizer.alter_vectorizer_schedule` function, or change the parallelism or the upper limit of the amount of data to be processed in one startup with the `pgx_vectorizer.alter_vectorizer_processing` function.

Example) Changing the execution interval to 5 minutes

```
SELECT
pgx_vectorizer.alter_vectorizer_schedule(pgx_vectorizer.get_vectorizer_id(view_name =>
'sample_embeddings'), interval '5 m');
```

Example) When changing the parallelism to 2 and the amount of data to be processed at one time to 200

```
SELECT
pgx_vectorizer.alter_vectorizer_processing(pgx_vectorizer.get_vectorizer_id(view_name =>
'sample_embeddings'), batch_size => 200, concurrency => 2);
```

Use a monitoring tool to check the time trend of the number of texts waiting for vectorization. If there is an error in the parameter value set by `set_worker_setting`, the vectorization process will not be executed. A message will be output to the server log, so please check it together with the `ai.vectorizer_status` view.

3.6.2 Checking the Status of Vectorization Processing

By referring to the `ai.vectorizer_errors` view, you can check the details of the errors that occurred during the vectorization process, the time of occurrence, the number of occurrences, etc.

Example) Check the details of the most recent error

```
rag_database=> SELECT * FROM ai.vectorizer_errors ORDER BY recorded DESC LIMIT 50;
-[ RECORD 1 ]-----
id          | 1
message     | embedding provider failed
details     | {"provider": "ollama", "error_reason": "model \"all-minilm\" not found, try pulling it first"}
recorded    | 2025-02-03 06:47:35.958882+00
-[ RECORD 2 ]-----
id          | 1
message     | embedding provider failed
details     | {"provider": "ollama", "error_reason": "model \"all-minilm\" not found, try pulling it first"}
recorded    | 2025-02-03 06:47:41.250279+00
```

Example) Check the number of errors for a vectorization definition

```
rag_database=> SELECT COUNT(*) FROM ai.vectorizer_errors WHERE id =
pgx_vectorizer.get_vectorizer_id(view_name => 'sample_embeddings');
count
-----
      20
(1 row)
```

Check the details of the error and remove the cause. Some embedded providers have a maximum load per period. If an error occurs because the load exceeds these conditions, use the `pgx_vectorizer.alter_vectorizer_schedule` or `pgx_vectorizer.alter_vectorizer_processing` function to adjust the worker schedule or parallelism.

3.6.3 Checking the Scheduler for Vectorization Processing

You can confirm that the scheduler for vectorization is running in the `pg_stat_activity` view.

```
SELECT * FROM pg_stat_activity WHERE backend_type LIKE '%pgx_vectorizer%';
```

3.7 Temporarily Disabling Vectorization Processing for Semantic Text Search

The vectorization process of this feature is performed periodically according to the specified schedule. For example, if you want to temporarily stop vectorization processing to reduce the impact on other work, you can disable vectorization processing using the `ai.disable_vectorizer_schedule` function.

```
rag_database=> SELECT
ai.disable_vectorizer_schedule(pgx_vectorizer.get_vectorizer_id(view_name =>
'sample_embeddings')); -- Specify the ID of the vectorizer you want to disable
```

To enable it, run the `ai.enable_vectorizer_schedule` function.

```
rag_database=> SELECT
ai.enable_vectorizer_schedule(pgx_vectorizer.get_vectorizer_id(view_name =>
'sample_embeddings')); -- Specify the ID of the vectorizer you want to enable
```

3.8 Semantic Text Search

For text data with vectorization defined, you can use the semantic text search feature to search for semantically similar text. The `pgx_vectorizer.pgx_similarity_search` function is used for the semantic text search.

The values in the distance column represent the distance between the text specified as an argument and the chunk, and the smaller the distance value, the higher the similarity to the text.

Example) Semantic text search

```
SELECT * FROM pgx_vectorizer.pgx_similarity_search('sample_embeddings'::regclass, 'text
for search', 5,
'<=>');
embedding_uuid | chunk | distance
-----+-----+-----
89A-927B-4271-82A3-6A73E8962B1C | Action items assigned. | 0.1027381927364
8E73B5F2-461A-4622-89A9-C1D364F4E19B | Next steps discussed. | 0.2938471928374
5D39A6E1-B226-4197-9F03-A78B80A509C2 | Timeline adjusted. | 0.5183749128735
2F97C1E5-6E8A-400F-8692-177D77740B6A | Budget approved. | 0.7815239187521
A0EEBC99-9C0B-4EF8-BB6D-6BB9BD380A11 | Status update. | 0.9274618273645
(5 rows)
```



Information

The `pgx_vectorizer.pgx_similarity_search` function specifies the type of distance to be used in vector similarity searches performed for semantic text searches. If this specification differs from the distance specified when defining an index for vector data, the index will not be used during the search.

3.9 Changing the Vector Representation Used in Semantic Text Search

You can have multiple vector representations for the text you want to perform semantic text search on. A vector representation is associated one-to-one with a vectorization definition, and you can have a different vector representation by defining a new vectorization. In the vectorization definition, you can specify the embedding model to be used, etc. If you do not need the previous vector representation, delete the vectorization definition.

Example) Defining a new vectorizer and deleting the old one

```
rag_database=> SELECT ai.drop_vectorizer(pgx_vectorizer.get_vectorizer_id(view_name =>
'sample_embeddings'), drop_all => true);
```

```
rag_database=> SELECT pgx_vectorizer.pgx_create_vectorizer(...);
```

3.10 Performance Tuning of Semantic Text Search

A semantic text search internally performs a similarity search on vector data. If an index is not specified for the vector data, all vector data is scanned for each semantic text search. You can check whether an index is being used in a semantic text search by running `pgx_similarity_search_checking_index`.

Example) Check if an index is used in a semantic text search

```
rag_database=> SELECT * FROM
pgx_similarity_search_checking_index('sample_embeddings'::regclass, 'text for search', 5,
'<=>');
ERROR: opclass of index is vector_ip_ops, but the distance_operator is <+>
```

If search results are output when you execute `pgx_similarity_search_checking_index`, the index is being used. If the search ends with an error as shown above, the distance operator specified in the index definition may differ from the distance type specified in the semantic text search. If the distance type specified in the semantic text search is incorrect, correct the distance specified. If the distance operator specified in the index definition is incorrect, delete the index and re-create the correct index.

Indexes for vectors explicitly define for the table that stores vectors that are created internally when vectorization is defined.

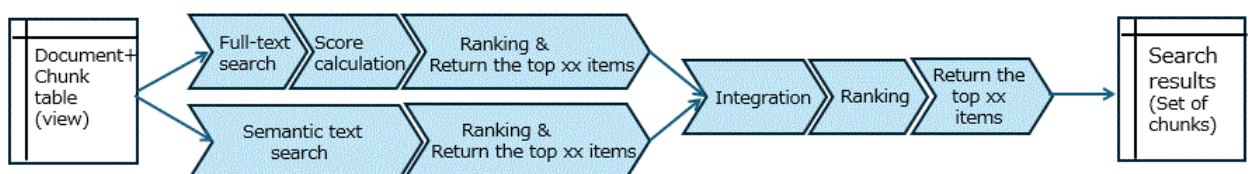
Example) Changing the vector table index

```
-- Check and delete the old index name
rag_database=> SELECT indexname FROM pg_indexes where tablename =
'sample_embeddings_store';
            indexname
-----
sample_embeddings_store_pkey
sample_embeddings_store_id_chunk_seq_key
sample_embeddings_store_embedding_idx - It is created for a column called embedding.
(3 rows)
rag_database=> DROP INDEX sample_embeddings_store_embedding_idx
-- Add the correct index
rag_database=> CREATE INDEX ON sample_embeddings_store USING hnsw (embedding
vector_l1_ops);
```

3.11 Hybrid Search

In hybrid search, results from both semantic text search and full-text search are combined, ranked based on scores considering both results, and the specified number of results are returned in order of highest score.

Hybrid search flow



Hybrid search uses the SQL function `pgx_vectorizer.pgx_hybrid_search` function. For `pgx_vectorizer.pgx_hybrid_search` function, refer to "[3.13.3.2 Hybrid Search](#)".

In the following example, a search is performed on the embedded view `sample_embeddings` created by the vectorization definition. The union of the top 10 results of full-text search and semantic text search of text is taken, and the top 20 results based on the final score are returned. In this example, partial match search using the LIKE operator is performed as a full-text search.

Example) Hybrid search

```
rag_database=> SELECT context_id, chunk, score FROM pgx_vectorizer.pgx_hybrid_search('
{
  "target_view": "sample_embeddings",
  "search_fusion": "UNION",
  "topN": 20,
  "semantic": {
    "search_text": "text for search",
    "num_result": 10,
    "score_weight": 10
  },
  "fulltext": {
    "search_condition": "chunk LIKE '%"text%' ",
    "score_expression": "CASE WHEN chunk LIKE '%"text%' THEN 1 ELSE 0 END"
  }
}':::JSONB
);
```

For the return results of the hybrid search, refer to "[Hybrid search results](#)".



Information

Combined search results are ordered using a method called weighted RRF (Reciprocal Rank Fusion). For more details, refer to "[Score calculation](#)".

3.12 Improving the Accuracy of Hybrid Search

Hybrid search is a search method that combines multiple search techniques to leverage their strengths and complement their weaknesses. To improve search accuracy by combining multiple search techniques, it is necessary to understand the search accuracy of each technique and choose a combination method accordingly. Below, we will explain how to calculate the search accuracy of individual search techniques (subqueries) that make up hybrid search.

3.12.1 Overview of Hybrid Search Tuning

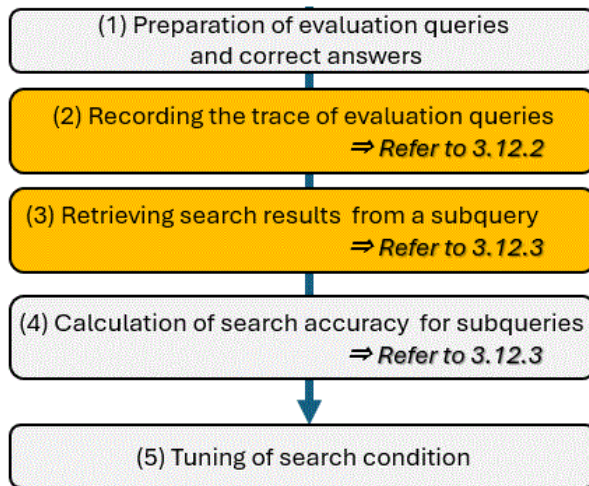
The accuracy of search processing can be calculated by preparing pairs of query conditions for evaluation and correct answers (for example, the search results that should be returned for that query) and comparing the results searched with those query conditions to the correct answers. The result obtained from hybrid search is only the final search result. Since we cannot know what the search results of individual subqueries were, their search accuracy cannot be calculated. Therefore, to calculate the search accuracy of subqueries, we use hybrid search trace information. Trace information is a record of the executed hybrid search, including the search conditions and search results of each subquery executed internally. The search accuracy can be calculated from the subquery search results included in the trace information. The accuracy of hybrid search using trace information can be calculated in the following two ways:

- a. Obtain the search results of each subquery recorded as traces and calculate using external evaluation tools.
- b. Calculate search accuracy within the database by inputting evaluation results for each text chunk that is a search result into the database.

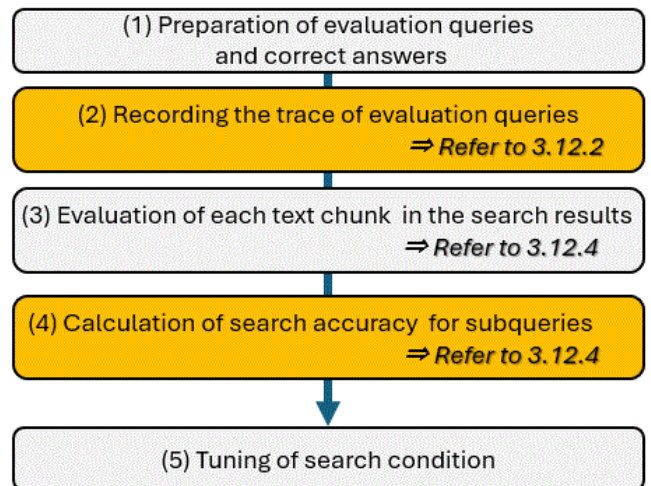
When evaluating the entire system, such as RAG-based applications, it is appropriate to use method a., which utilizes external evaluation tools. On the other hand, if you want to analyze the accuracy of individual subqueries or text chunks in detail without using tools, method b., which calculates accuracy within the database, is recommended.


Below are the steps for each method.

a. Calculating search accuracy using external tools



b. Calculating search accuracy into a database



 : Features provided by Fujitsu Enterprise Postgres

For basic ideas on the evaluation of knowledge data search, refer to "[3.13.6 Evaluation of Knowledge Data Search](#)".

3.12.2 Recording and Deleting Traces of Hybrid Search

To record trace for evaluation queries, a user with trace recording privileges enables the trace and executes a hybrid search using the evaluation queries. As a result, the trace information is recorded in the table. Trace recording can only be done on the primary server. Even if a hybrid search with trace enabled is performed on a hot standby server, the trace will not be recorded. The trace privileges are granted to the database user performing the search by the extension owner using the `pgx_grant_access_on_hybrid_search_trace` function. To enable the trace, the database user performing the hybrid search sets the `enable_hybrid_search_trace` parameter to on in their session. Regarding the `enable_hybrid_search_trace` parameter, refer to "[3.13.5 Parameters](#)".

Example) Privilege granted by the owner of the extension

```
rag_database=# SELECT pgx_vectorizer.pgx_grant_access_on_hybrid_search_trace('app_user');
```

Example) Setting session parameters by users performing searches and executing hybrid searches

```
rag_database=> SET pgx_vectorizer.enable_hybrid_search_trace = 'on';
rag_database=> SELECT * FROM pgx_vectorizer.pgx_hybrid_search(...);
```

When a trace is recorded, a query ID that identifies the search request required for subsequent operations is issued. The query ID is returned as a result of a hybrid search, and can also be confirmed by referring to the table where trace information is recorded (`pgx_trace_query` table). This table records the database user who executed the hybrid search, the time, and the application name, allowing you to identify the query ID using these details.

Example) The database user who executes the command refers to the query ID of the hybrid search executed within an hour

```
rag_database=> SELECT queryid FROM pgx_vectorizer.pgx_trace_query WHERE retrieval_time >=
(current_timestamp - interval '1 hour') ORDER BY retrieval_time;
queryid
-----
    10001
(1 row)
```

After completing the evaluation and tuning, trace information and evaluation values will be deleted. By deleting records with the target query ID from the `pgx_trace_query` table, all related information will be removed.

Example) Delete trace information by specifying the query ID

```
rag_database=> DELETE FROM pgx_vectorizer.pgx_trace_query WHERE queryid = 10001;
```

Example) The administrator deletes all trace information older than 30 days

```
rag_database=# DELETE FROM pgx_vectorizer.pgx_trace_query WHERE retrieval_time <=
(current_timestamp - interval '30 days') ;
```

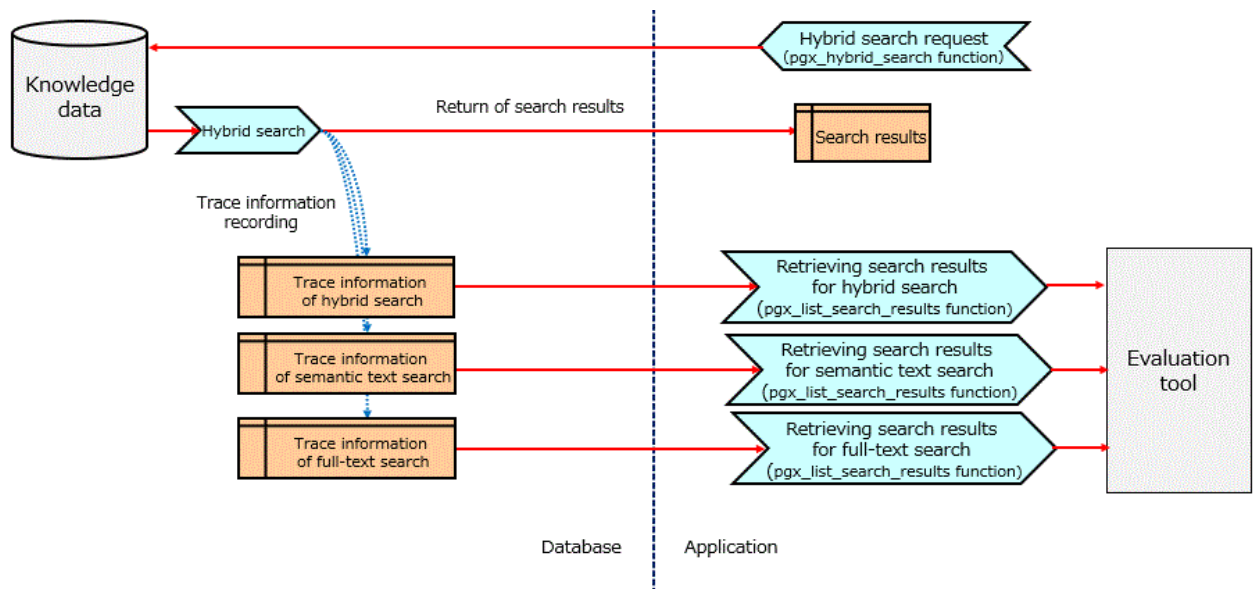


Information

The trace information includes the search conditions and the knowledge data itself as search results, and is recorded in a common table within the database. Privileges are set so that only the database users and administrators who executed the search process can refer to, update, and delete that trace information.

3.12.3 Calculation of Search Accuracy Using External Tools

When calculating the search accuracy of subqueries using external tools, the search results of each subquery are obtained from trace information.



To retrieve the search results of a subquery, execute the `pgx_list_search_results` function by specifying the query ID and the type of subquery. When executed, a set of text chunks is returned in descending order of score, in the same format as the hybrid search results.

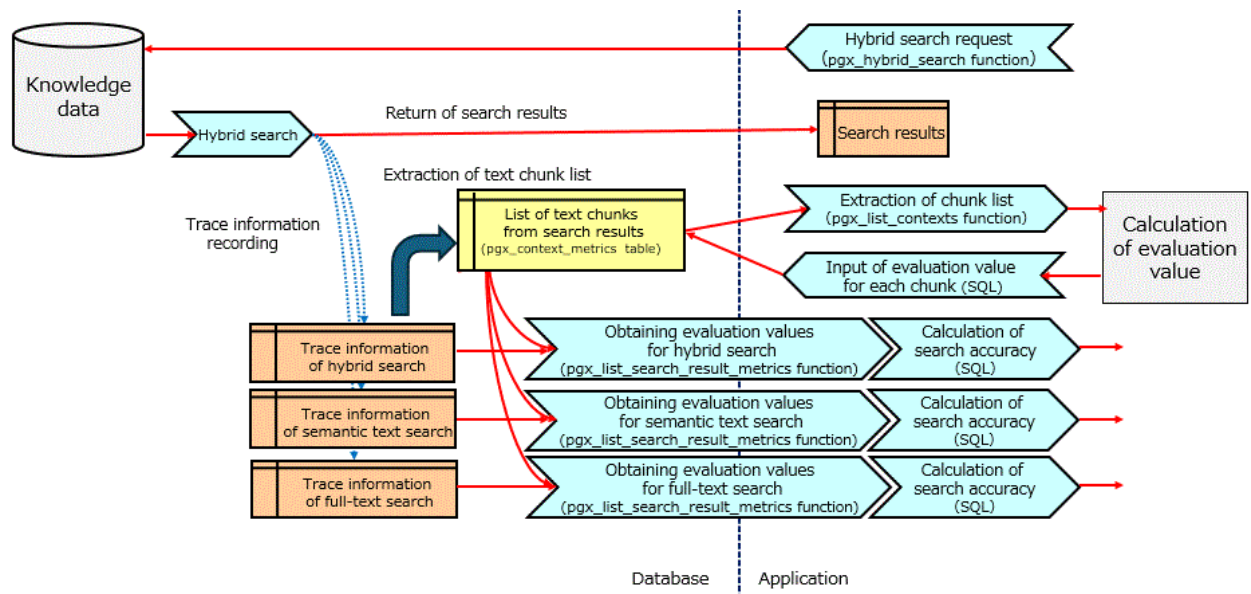
Example) Retrieve the search results of the full-text search subquery with query ID 10001

```
rag_database=> SELECT pgx_vectorizer.pgx_list_search_results(queryid => 10001,
subquery_type => 'fulltext');
```

As a result, by providing the search results of subqueries for evaluation queries and the correct answers corresponding to evaluation queries to the evaluation tool, the search accuracy of subqueries can be calculated.

3.12.4 Calculation of Search Accuracy in the Database

The calculation of search accuracy within a database consists of two major steps: "the calculation of evaluation values for each text chunk in the search results" and "the calculation of search accuracy based on those evaluation values".



Calculation of evaluation values for each text chunk in the search results

Use the `pgx_list_contexts` function to extract the text chunks returned by the hybrid search executed and any of its constituent subqueries.

Example) Extract the text chunk returned by the hybrid search with query ID 10001 and any of the subqueries

```
rag_database=> SELECT pgx_vectorizer.pgx_list_contexts(queryid => 10001);
```

Match the extracted text chunks with the correct answers for the evaluation queries, calculate the evaluation values for each, and input them into the evaluation value table (`pgx_context_metrics` table). To input the evaluation values, use the `UPDATE` statement on the `pgx_context_metrics` table. For an example of using the `UPDATE` statement, refer to "[3.12.4.1 Example of Calculating Search Accuracy in a Database](#)". The evaluation values can be stored in JSON format to allow calculation of search accuracy using any metric.

Calculation of search accuracy based on evaluation value

The text chunk of the search result from the subquery and its evaluation value can be referenced using the `pgx_list_search_result_metrics` function. This can be used to calculate search accuracy.

Example) Refer to the search results of full-text search and their evaluation values

```
rag_database=> SELECT pgx_vectorizer.pgx_list_search_result_metrics(queryid => 10001,
subquery_type => 'fulltext');
```

3.12.4.1 Example of Calculating Search Accuracy in a Database

Explain an example of calculating search accuracy.

First, store the evaluation value in the `pgx_context_metrics` table. Here, input the evaluation value into the `relevance_indicator` in JSON format.

```
{
  "relevance_indicator": 1
}
```

The evaluation value is set to 1 if the context IDs of the search results are 'AAAA', 'BBBB', or 'DDDD', indicating relevance to the search criteria, and 0 if they are not relevant. In this case, an `UPDATE` statement is executed on the `pgx_context_metrics` table using `queryid` as the key, and the evaluation values for all text chunks are input collectively.

Example) Input of evaluation value using the `UPDATE` statement

```
rag_database=> UPDATE pgx_vectorizer.pgx_context_metrics
SET context_metrics = jsonb_set(
    context_metrics,
    '{relevance_indicator}',
    CASE
        WHEN context_id IN ('AAAA', 'BBBB', 'DDDD') THEN '1'::jsonb ELSE '0'::jsonb
    END
)
WHERE queryid = 10001;
```

Next, refer to the search results of subqueries that include evaluation values and calculate the search accuracy for each subquery.

Search accuracy is calculated using a search accuracy metric called precision@k. precision@k is a metric that represents the precision when focusing on the top k search results, and can be calculated using the following formula.

$$\text{precision@}k = \frac{\text{The number of records related to the query among the top } k \text{ search results returned}}{k}$$

Below, the top 10 search precision (precision@10) for semantic text search is calculated using the relevance_indicator. The same can be calculated for full-text search.

```
rag_database=> SELECT
SUM( (context_metrics ->> 'relevance_indicator'):: integer )::float / COUNT(*) AS
precision_at_10
FROM(
    SELECT context_metrics FROM
    pgx_vectorizer.pgx_list_search_result_metrics(queryid => 10001, subquery_type =>
'semantic')
    ORDER BY score DESC LIMIT 10
);
```

3.12.5 Tuning of Hybrid Search

Tune hybrid search based on calculated search accuracy.

To reduce application modifications due to changes in hybrid search parameters, it is useful to create search functions for each application as user-defined functions.

Below is an example of defining a SQL function that only takes the string you want to search as input. By defining such a function, you can change the topN parameter and score_weight parameter on the database server side.

```
rag_database=> CREATE OR REPLACE FUNCTION mysearch(in text, in text, out uuid, out text,
out float8)
AS $$
SELECT context_id, chunk, score FROM pgx_vectorizer.pgx_hybrid_search(
jsonb_set(
    jsonb_set('{
        "target_view": "sample_embeddings",
        "search_fusion": "UNION",
        "topN": 20,
        "semantic": {
            "num_result": 10,
            "score_weight": 10
        },
        "fulltext": {
            "score_expression": "1"
        }
    }'::jsonb,
    '{semantic, search_text}', ('' || $1 || ''')::jsonb),
    '{fulltext, search_text}', ('' || $2 || ''')::jsonb)
```

```
)  
$$ LANGUAGE SQL
```

3.13 Reference

3.13.1 Vectorization Functions

3.13.1.1 Defining Vectorization

For the parameters to specify in the definition of vectorization, refer to the pgai documentation. If you want to perform vectorization within the database, you must specify `pgx_vectorizer.schedule_vectorizer` as the scheduling.

3.13.1.2 Vectorization Schedule

You can generate vector data from text data in the following ways.

- Manually start vectorization processing
- Perform vectorization processing periodically within the database

The timing of vectorization is determined by the schedule specified when defining vectorization. If `schedule_none` is specified, periodic vectorization will not be performed. To perform vectorization at a specified time, run the `run_vectorize_worker` function. If `schedule_vectorizer` is specified, periodic vectorization will be performed within the database.

Automatic vectorization can be disabled with the `ai.disable_vectorizer_schedule` function. It can also be re-enabled with the `ai.enable_vectorizer_schedule` function.

3.13.1.3 Vectorizer Management Functions

The following functions are provided to use for vectorization feature.

Function	Return type	Description
<code>pgx_create_vectorizer</code> ([Arguments that can be specified with <code>ai.create_vectorizer</code>], <code>fulltext_indexing => jsonb</code>)	integer	It is used to automatically define full-text search indexes for vectorization and chunks. The ID of the created vectorizer will be returned. In addition to the arguments that can be specified for <code>ai.create_vectorizer</code> , specify the definition of the full-text search index (<code>pgx_fulltext_indexing_gin()</code> or <code>pgx_fulltext_indexing_gist()</code>). <code>fulltext_indexing</code> is optional when only defining vectorization and when using full-text search indexes other than GiST or GIN. If omitted, a full-text search index will not be created.
<code>pgx_fulltext_indexing_gin</code> (<code>target_column text</code> , <code>fastupdate boolean</code> , <code>gin_pending_list_limit integer</code> , <code>opclass text</code> , <code>min_rows integer</code> , <code>create_when_queue_empty boolean</code>)	jsonb	It is used when creating a full-text search index as GIN. The value specified for <code>fulltext_indexing</code> in <code>pgx_create_vectorizer</code> will be returned. <code>target_column</code> is used when using expression indexes or covering

Function	Return type	Description
		<p>indexes.</p> <p>The column name for text chunks is chunk, and the default value for target_column is chunk.</p> <p>Specify the operator class to use in opclass.</p> <p>min_rows specifies the threshold for creating an index. The default is 100000. A full-text search index is created when the number of records in the embedded table exceeds this threshold.</p> <p>create_when_queue_empty is a parameter that specifies the timing for creating a full-text search index. The default is true, and the full-text search index is created after the vector conversion is completed.</p> <p>Other specifiable parameters refer to index storage parameters.</p>
pgx_fulltext_indexing_gist(target_column text, fillfactor integer, buffering text, opclass text, min_rows integer, create_when_queue_empty boolean)	jsonb	<p>It is used when creating a GIST as a full-text search index.</p> <p>The value specified for fulltext_indexing in pgx_create_vectorizer will be returned.</p> <p>Specify the operator class to use in opclass.</p> <p>The settings for target_column, min_rows, and create_when_queue_empty are the same as those of pgx_fulltext_indexing_gin.</p> <p>Other specifiable parameters refer to index storage parameters.</p>
pgx_delete_fulltext_index_config(vectorizer_id integer)	void	<p>It is used when deleting the vector transformation defined by pgx_create_vectorizer.</p> <p>Delete information related to full-text search of vectorization created by the pgx_create_vectorizer function.</p>
get_vectorizer_id(view_name pg_catalog.pg_regclass)	integer	Returns the ID of the vectorization definition that corresponds to the specified embedded view.
schedule_vectorizer(schedule_interval interval)	json	<p>Specify the interval for the vectorization process in schedule_interval.</p> <p>This function returns a JSON to be specified in the scheduling argument of the create_vectorizer function.</p> <p>If schedule_interval is not specified, 10 minutes will be specified.</p>

Function	Return type	Description
<code>alter_vectorizer_processing(vectorizer_id integer, batch_size integer, concurrency integer)</code>	void	Changes the amount of data to be converted at one time and the worker multiplicity for the vectorization definition with the id specified in <code>vectorizer_id</code> .
<code>alter_vectorizer_schedule(vectorizer_id integer, schedule_interval interval)</code>	void	Changes the interval for the vectorization process for the vectorization definition with the id specified in <code>vectorizer_id</code> . If <code>schedule_interval</code> is not specified, 5 minutes will be specified.
<code>run_vectorize_worker(vectorizer_id integer)</code>	integer	Immediately starts the vectorization process for the vectorization definition with the ID specified in <code>vectorizer_id</code> , and starts the vectorization process in the background. Returns the PID of the started process.
<code>start_vectorize_scheduler(void)</code>	void	This will start the vectorize scheduler that connects to the database where this SQL function was executed. If the vectorize scheduler is already running, an error will occur.
<code>pgx_embedding_onnx(model text, dimension integer)</code>	jsonb	Generate a configuration JSON object for use with the <code>pgx_create_vectorizer</code> function. It is used when executing automatic vectorization with a model imported into the database. Specify the model name and the number of dimensions of the vector to be generated. This cannot be omitted. If the specified model is not imported, a warning message will be output.



See

For information about the arguments that can be specified with `ai.create_vectorizer`, refer to the `pgai` documentation.

3.13.2 Embedded Provider Management Functions

The following functions are provided to set and reference parameters of the embedded provider.

Function	Return type	Description
<code>get_worker_setting(type text, param text)</code>	text	Specify a combination of type and parameter (param) to get the value set for that parameter. Only the user who executed <code>CREATE EXTENSION</code> can execute this command.

Function	Return type	Description
set_worker_setting(type text, param text, value text)	void	Specify the combination of type and parameter (param), and set the value (value) for that parameter. Only the user who executed CREATE EXTENSION can execute this command.

The possible embedded provider names and parameters are:

type	param	Description
openai	OPENAI_API_KEY	OpenAI API key value
voyage	VOYAGE_API_KEY	VoyageAI API key value
ollama	OLLAMA_BASE_URL	Ollama API base url
user	VECTORIZE_USER	Username to connect to the database the worker that performs the vectorization.

3.13.3 Search Functions

3.13.3.1 Semantic Text Search

The following functions are provided for semantic text search.

Function	Return type	Description
pgx_similarity_search(view pg_catalog.regclass, query text, num_result integer default 5, distance_operator text default '<=>', OUT embedding_uuid uuid, OUT chunk text, OUT distance float8);	SETOF record	Searches the embedding view specified in view to obtain text similar to the text specified in query. You can display results up to the number specified in num_result. You can specify the distance calculation method using distance_operator.
pgx_similarity_search_checking_index(view pg_catalog.regclass, query text, num_result integer default 5, distance_operator text default '<=>', OUT embedding_uuid uuid, OUT chunk text, OUT distance float8);	SETOF record	An error occurs if the index operator defined in the embedding column of the table that references the view does not match the operator specified in distance_operator. Other than the above, it is the same as the pgx_similarity_search function.

3.13.3.2 Hybrid Search

The following functions are provided for hybrid search.

Function	Return type	Description
<code>pgx_hybrid_search(query jsonb, OUT queryid bigint, OUT context_id uuid, OUT chunk text, OUT score float8)</code>	SETOF record	Search the specified embedded view and retrieve the relevant text. For the return value, refer to "3.13.3.3 Details of the pgx_hybrid_search Function" . For query details, refer to "Hybrid search results" .
<code>pgx_list_search_results(queryid bigint, subquery_type text)</code>	SETOF record	If trace information for the specified query ID is recorded, it returns the search results. If <code>subquery_type</code> is not specified, it returns the results of a hybrid search. If <code>subquery_type</code> is semantic or fulltext, it returns the results of semantic text search or full-text search conducted internally, respectively.
<code>pgx_list_search_result_metrics(queryid bigint, subquery_type text)</code>	SETOF record	Returns search results along with the evaluation value of the text chunk. Other specifications follow the <code>pgx_list_search_results</code> function.
<code>pgx_list_contexts(queryid bigint)</code>	SETOF record	If trace information for the specified query ID is recorded, it extracts and returns a list of text chunks returned as search results in either hybrid search processing or subqueries executed internally. The extracted list of text chunks is inserted into the evaluation value table.
<code>pgx_hybrid_search_trace_size()</code>	bigint	Return the size of the trace information and evaluation value table in bytes.
<code>pgx_grant_access_on_hybrid_search_trace(role name)</code>	None	Grant the necessary access rights to the role for recording and evaluating trace information.
<code>pgx_revoke_access_on_hybrid_search_trace(role name)</code>	None	Revoke the access rights necessary for recording and evaluating trace information assigned to the role.

The `pgx_list_search_results` function returns a set of records of the following type, in descending order of score.

Column	Type	Description
<code>queryid</code>	bigint	Query ID for hybrid search feature.
<code>context_id</code>	uuid	Identifier of the returned text chunk. Unique within the embedded view.
<code>chunk</code>	text	Returned text chunk.

Column	Type	Description
score	real	Hybrid search score. A higher score indicates a better match with the conditions.

The `pgx_list_search_results_metrics` function returns a set of records of the following type, in descending order of score.

Column	Type	Description
queryid	bigint	Query ID for hybrid search feature.
context_id	uuid	Identifier of the returned text chunk. Unique within the embedded view.
chunk	text	Returned text chunk.
score	real	Hybrid search score. A higher score indicates a better match with the conditions.
metrics	jsonb	Evaluation value for each input text chunk.

The `pgx_list_contexts` function returns a set of records of the following type, in descending order of score.

Column	Type	Description
queryid	bigint	Query ID for hybrid search feature.
context_id	uuid	Identifier of the returned text chunk. Unique within the embedded view.
chunk	text	Returned text chunk.

3.13.3.3 Details of the `pgx_hybrid_search` Function

query argument

Specify search conditions in JSON format for the query argument. Specify the following for each key.

Key	JSON format	Description
target_view	string	Embedded view to be searched. Interpreted as <code>pg_catalog.regclass</code> . Cannot be omitted.
search_fusion	string	How to combine semantic text search and full-text search. The default is "UNION". "UNION": Returns the union of each search result. "INTERSECT": Returns the intersection of each search result. "TEXT_ONLY": Returns the intersection with all results of the full-text search. "VECTOR_ONLY": Returns the intersection with all results of the semantic text search. "MINUS_TEXT": Returns the results of the semantic text search excluding the intersection. "MINUS_VECTOR": Returns the results of the full-text search excluding the intersection.
rrf_k	number	Weighted RRF's weight. When a floating-point type is specified, the

Key	JSON format	Deccription
		decimal part is truncated. The default is 60.
topN	number	Maximum number of results returned from hybrid search. When a floating-point type is specified, an error occurs. The default is 20.
semantic.search_text	string	Text to search with semantic text search. Cannot be omitted.
semantic.distance_operator	string	Distance operator in vector comparison. The default is "<=>".
semantic.num_result	number	Maximum number of items to be retrieved by semantic text search. When a floating-point type is specified, an error occurs. If omitted, it is calculated based on topN. If the combination method is "UNION", it is half of topN, otherwise it is the same as topN.
semantic.score_weight	number	Relative weighting (importance) of semantic text search. When a floating-point type is specified, the decimal part is truncated. The default is 1.
fulltext.search_condition	string	Full-text search condition (an expression that returns a boolean value specified in the WHERE clause). The name of the column targeted for full-text search in the embedded view is chunk. Refer to the " 3.3.3 Definition of Vectorization " for the definition of the embedded view. Cannot be omitted.
fulltext.num_result	number	Maximum number of items to be retrieved by full-text search. When a floating-point type is specified, an error occurs. If omitted, it is calculated in the same way as semantic.num_result.
fulltext.score_weight	number	Relative weight (importance) of full-text search. When a floating-point type is specified, the decimal part is truncated. The default is 1.
fulltext.score_expression	string	An expression that returns a floating-point score for full-text search. You can refer to tableoid and ctid. It cannot be omitted.



See

For the function used in calculating scores in full-text search, refer to "Full Text Search" in the PostgreSQL Documentation.

Information

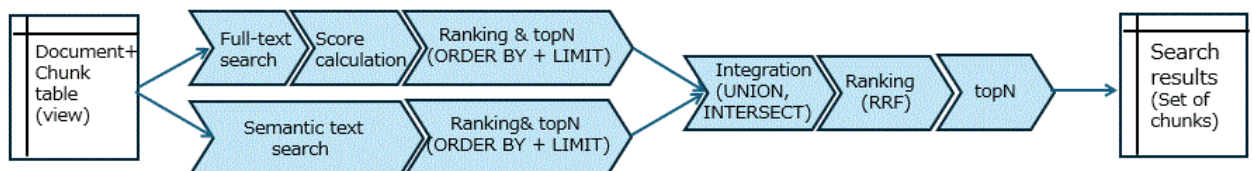
You can specify a matching operator for full-text search as a condition for full-text search. Specify the operator supported by the full-text search index defined in the table to be searched. The search text for full-text search must be specified in the pattern required by the matching operator.

Example) Query argument of the `pgx_hybrid_search` function

```
{
  "target_view": "sample_embeddings",
  "search_fusion": "UNION",
  "topN": 20,
  "semantic": {
    "search_text": "text for search",
    "distance_operator": "<=>",
    "num_result": 10,
    "score_weight": 1
  },
  "fulltext": {
    "search_condition": "chunk @@ websearch_to_tsquery('search text')",
    "num_result": 10,
    "score_weight": 1,
    "score_expression": "ts_rank_cd(to_tsvector(chunk), websearch_to_tsquery('search text'))"
  }
}
```

The search conditions specified by the query argument

The search conditions specified by the query argument correspond to the following flow. (This flow is for explaining the meaning of the search process and is not the same as the actual access plan.)



- Full-text search: Perform a full-text search based on the conditions specified in `fulltext.search_condition`.
- Full-text search score calculation: Perform score calculation using `fulltext.score_expression`. If not specified, full-text search will result in an error.
- Full-text search ranking & topN: Rank the results in order of highest score. Only the number of search results specified by `semantic.num_result` from the top will be adopted, and the rest will be discarded.
- Semantic text search: Perform semantic text search based on vector similarity search using `semantic.search_text` and `semantic.distance_operator`.
- Semantic text search ranking & topN: Rank results in order of smallest distance from vector similarity search results. Only the number of search results specified by `semantic.num_result` from the top will be adopted, and the rest will be discarded.
- Integration: Integrate the results of full-text search and semantic text search according to the method specified in `search_fusion`.
- Ranking: Calculate scores based on the ranks of full-text search and semantic text search and assign final ranks. Use `fulltext.score_weight`, `semantic.score_weight` and `rrf_k`.
- topN: Return only the specified number of results from the top based on the final rank.

Score calculation

The score of the integrated result of multiple ranked results is calculated using a method called weighted RRF (Reciprocal Rank Fusion). Unweighted RRF considers the reciprocal of the rank of each search method as the score and takes the sum of multiple search methods as the final score. The final rank is determined based on that score. Weighted RRF calculates the sum of multiple search methods by multiplying each score by the weight for each search method as a coefficient.

The method for calculating the score $WRRF(d)$ using weighted RRF for a certain search result d is shown below. i represents the search method, and the sum of scores for all search methods is obtained.

$$WRRF(d) = \sum_i w_i \times \frac{1}{k + rank_i(d)}$$

The value of k can be changed with the `rrf_k` parameter, and the default is 60. The smaller this parameter is, the greater the impact of differences in the original rank.

The rank ($rank_i(d)$) for search results not returned by one of the search methods is considered a sufficiently large value (considered to have no impact on the final score).

Hybrid search results

The `pgx_hybrid_search` function returns a set of records of the following type.

Column	Type	Description
queryid	bigint	Query ID of the hybrid search feature used to refer to trace information. A valid value will be returned if trace information is recorded. Otherwise, 0 will be returned.
context_id	uuid	Identifier of the returned text chunk. Unique within the embedded view.
chunk	text	Returned text chunk.
score	real	Hybrid search score. A higher score indicates a better match with the conditions.

3.13.4 Tables/Views Created by Semantic Text Search and Automatic Vectorization Feature

For information about the tables and views that `pgai` creates, refer to the `pgai` documentation.

`pgx_vectorizer` creates the following tables.

`pgx_vectorizer.worker_setting_table` table

Stores information about the parameters used by the embedded provider.

Column	Type	Constraint	Description
type	text	PRIMARY KEY	The type of parameter to set. The name of the embedded provider or user
parameter	text	PRIMARY KEY	Parameter name
value	text	NOT NULL	Value to set for the parameter

`pgx_vectorizer.pgx_fulltext_index` table

The `pgx_fulltext_index` table stores the definition information of the full-text search index created by the `pgx_create_vectorizer` function.

Column	Type	Reference destination	Description
vectorizerid	integer	ai.vectorizer.id	Vectorization identifier
config	jsonb		Definition information of full-text search index. Configurable values are the return values of the <code>pgx_fulltext_indexing_gin</code> function or the <code>pgx_fulltext_indexing_gist</code> function.

pgx_vectorizer.pgx_trace_query table

Provides information on hybrid search where trace information is recorded.

When a hybrid search is performed with tracing enabled, one record is inserted.

Deleting records from this table also deletes records from the `pgx_trace_subquery` table, `pgx_trace_results` table, and `pgx_context_metrics` table with the same queryid.

Column	Type	Description
queryid	bigint	Hybrid search ID
retrieval_time	timestamp with timezone	Time when hybrid search was conducted
retrieval_user	oid	Database user who performed the hybrid search
application_name	text	Application name that issued the hybrid search. The value of <code>application_name</code> specified
query	jsonb	Search criteria for hybrid search

pgx_vectorizer.pgx_trace_subquery table

Provides information about the subqueries for hybrid searches where trace information is recorded.

When a hybrid search is performed with tracing enabled, multiple records are inserted for each subquery.

Column	Type	Description
queryid	bigint	Hybrid search ID
subquery_seq	smallint	Serial number of subquery conducted within a hybrid search(A serial number closed in a certain hybrid search process)
subquery_type	text	Types of subqueries '-' : Overall hybrid search (main query) 'semantic' : Semantic text search 'fulltext' : Full-text search
subquery_text	text	Content of the executed subquery (SQL statement)

pgx_vectorizer.pgx_trace_results table

Provides information about the text chunks returned by each search process of hybrid search and subquery.

When performing a hybrid search with trace information enabled, multiple records corresponding to the text chunks of the search results of each subquery are inserted.

Column	Type	Description
queryid	bigint	Hybrid search ID
subquery_seq	smallint	Serial number of subquery conducted within a hybrid search
context_id	uuid	Returned text chunk ID
chunk	text	Returned text chunk
rank	integer	Rank of text chunks within the same subquery type
score	real	Score of text chunk. Hybrid search overall : Score calculated by RRF Semantic text search : Score based on distance between vectors in text meaning search Full-text search : Score calculated by score_expression in hybrid search processing

In the case of text meaning search, the score is calculated using the distance d between the vector representations of the text chunk and the query text, according to the following formula.

$$score = \frac{1}{1 + d}$$

pgx_vectorizer.pgx_context_metrics table

This is a list of text chunks returned by either hybrid search or subquery. When the `pgx_list_contexts` function is executed, records for the specified query ID are inserted. You can enter any evaluation value in the `context_metrics` column.

Column	Type	Description
queryid	bigint	Hybrid search ID
context_id	uuid	Returned text chunk ID
context_metrics	jsonb	You can input evaluation values for text chunks.

pgx_vectorizer.pgx_trace_contexts view

Provide information about the text chunks returned by either hybrid search or subqueries.

For hybrid searches conducted with trace information enabled, there are records for each text chunk of search results for each subquery.

Column	Type	Description
queryid	bigint	Hybrid search ID
retrieval_time	timestamp with timezone	Time when hybrid search was conducted
retrieval_user	oid	Database user who performed the hybrid search
application_name	text	Application name that issued the hybrid search. The value of <code>application_name</code> specified
query	jsonb	Search criteria for hybrid search
subquery_type	text	Types of subqueries '-' : Overall hybrid search (main query) 'semantic' : Semantic text search

Column	Type	Description
		'fulltext' : Full-text search
subquery_text	text	Content of the executed subquery (SQL statement)
context_id	uuid	Returned text chunk ID
chunk	text	Returned text chunk
rank	integer	Rank of text chunks within the same subquery type
score	real	Score of text chunk. Hybrid search overall : Score calculated by RRF Semantic text search : Score based on distance between vectors in text meaning search Full-text search : Score calculated by score_expression in hybrid search processing
context_metrics	jsonb	You can input evaluation values for text chunks.

3.13.5 Parameters

Describes the parameters to be set in the postgresql.conf file when using the semantic text search and automatic vectorization feature.

- pgx_vectorizer.max_vectorize_worker

Specify the maximum number of workers that can run simultaneously within an instance and perform vectorization. The number of workers that perform vectorization is determined by the number of vectorization definitions created. This parameter can only be set at server startup. The default is 1. Because workers act as background workers, add the value set for this parameter plus the number of databases with the pgx_vectorizer feature enabled plus 2 to the max_worker_processes parameter, which specifies the maximum number of background workers. If the value set for max_worker_processes is insufficient, the instance cannot start.

- pgx_vectorizer.pgai_worker_path

Specify the path to the program that performs vectorization processing. Specify <Fujitsu Enterprise Postgres server feature installation directory>/OSS/pgai-worker/bin/pgai. The default value is /opt/fsepv<x>server64/OSS/pgai-worker/bin/pgai (where "<x>" indicates the product version). You can reflect the changes by reloading the configuration file.

- pgx_vectorizer.enable_hybrid_search_trace(boolean)

Enable trace information for hybrid search. The default is disabled. This parameter can be specified with the SET statement by any user, but additional access rights are required to record trace information. Also, even if this parameter is enabled on a hot standby server, trace information will not be recorded.

3.13.6 Evaluation of Knowledge Data Search

3.13.6.1 Concept of Evaluation for Knowledge Data Search

A knowledge database exists to provide users with appropriate knowledge. It is necessary to evaluate whether it is providing users with appropriate knowledge. For the evaluation of knowledge data (offline evaluation), a set of pairs of queries and results related to those queries (search results that should be returned for those queries) is prepared, and search accuracy is calculated using it.

Basic indicators in offline evaluation of search systems are recall and precision. A high recall state means that records related to the request are included in the search results without omission. A high precision state means that only records related to the request are included in the search results. Recall and precision can be calculated using the following formulas.

$$\text{recall} = \frac{\text{The number of records related to the query in the returned search results}}{\text{The number of records related to the inquiry}}$$

$$\text{precision} = \frac{\text{The number of records related to the query in the returned search results}}{\text{The number of returned records}}$$

Recall and precision are indicators that well represent whether necessary data is included or unnecessary data is not included. However, there are the following challenges in directly using recall and precision for evaluating knowledge data retrieval.

- Difficulty in defining query-related results

When the same knowledge is contained in multiple records, returning even one record as a search result satisfies the user's request. Therefore, it is impossible to define a single set of search results that should be returned for that query.

- High cost of preparing query-related results

Strictly defining query-related results requires determining whether each record in the database is related to the query or not.

To address the first issue, instead of defining the correct set of records to be returned, there is a method to define it by the content and meaning of the data. For example, there is a method to evaluate whether the sentences included in the records support the expected claim.

The preparation cost for the correct answer, which is the second issue, can also be reduced by the above approach. This is because the definition of the correct answer can be made from the content of the query or the content of the correct answer to the query, rather than from (all) records included in the database.

"Judging based on the content of the text" refers to things like the following.

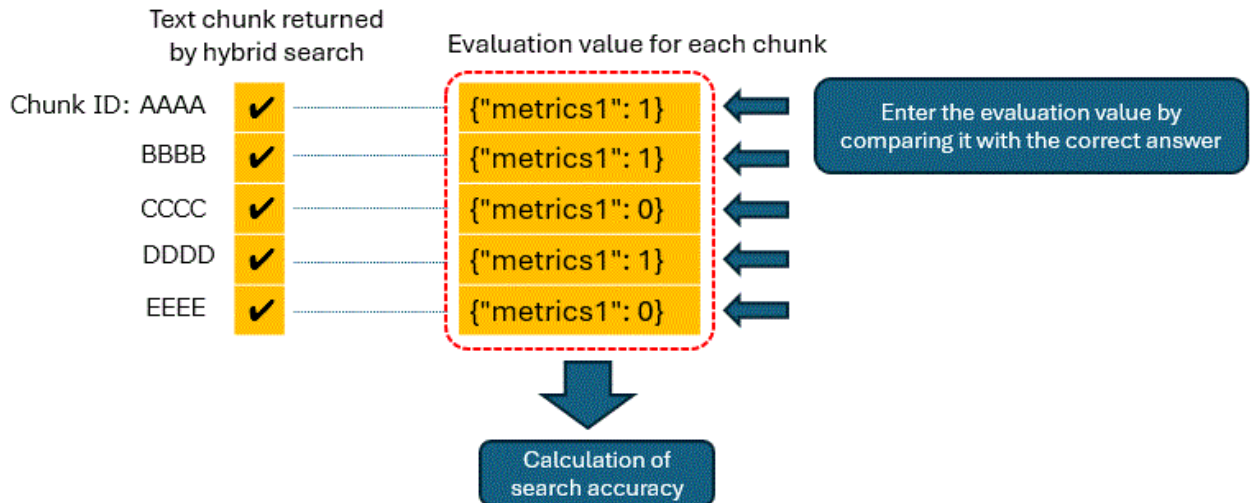
When the statement "Mount Fuji is the highest peak in Japan" is considered correct, if a record contains the sentence "Mount Fuji is 3776 meters high and is recorded as the tallest mountain in Japan", this sentence can be judged to support the correct fact.

In the method of judging based on the content and meaning of the text, it is necessary to extract the statement of fact from the query or correct content and determine whether these statements are included in the search results. While it is certainly possible for humans to perform this task, there is also a method to have language models do it. By doing so, the preparation cost for correctness can be further reduced.

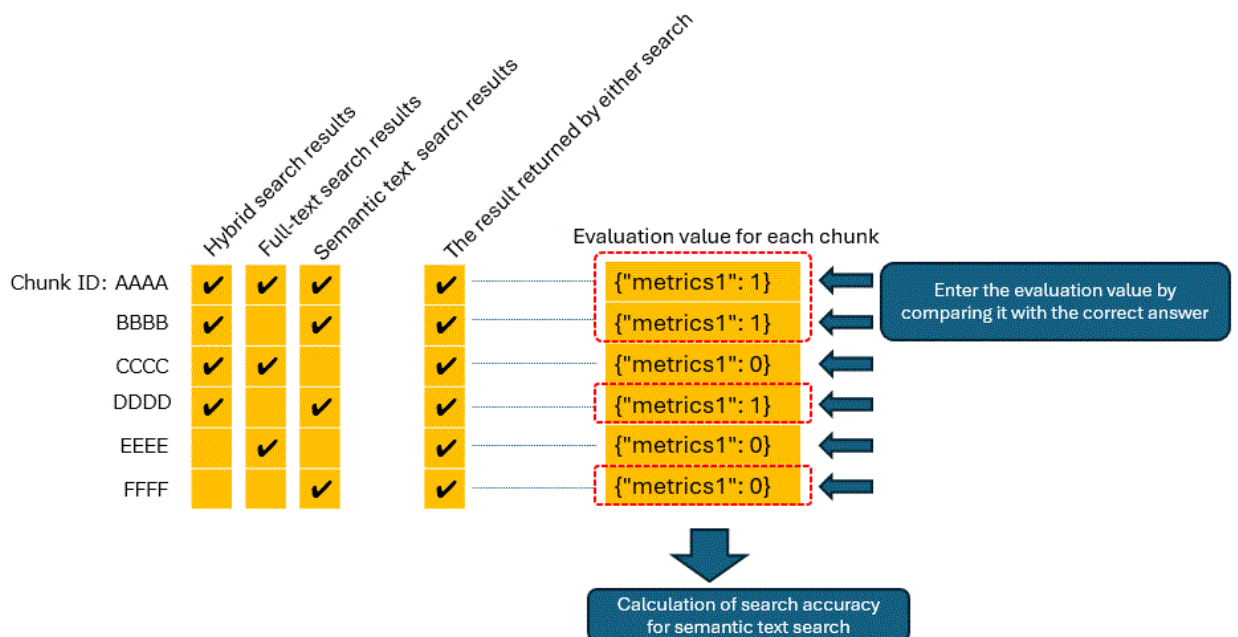
When adopting a method to evaluate search accuracy based on the content included in the records, it is important to note that the evaluated target is not only the implementation method of the search process. This method is a direct evaluation approach of the content obtained as search results. Therefore, not only the search method of the database but also what kind of knowledge is included in the database and how the user converts the content they want to know into queries for the database will affect the evaluation. When evaluating the accuracy of search methods with this approach, it is necessary to identify which part is causing the issue before starting to tune the search method.

3.13.6.2 Evaluation Value per Record and Search Accuracy

There are various metrics to represent search accuracy, but many calculations of search accuracy consist of two steps: obtaining evaluation values for individual records obtained as search results, and calculating search accuracy using the overall evaluation values of individual records (text chunks).



The search accuracy of subqueries can be calculated in the same way. The search accuracy can be evaluated in the same way for both full-text search subqueries and semantic text search. However, the text chunks returned by subqueries and the final hybrid search contain the same ones. Since the evaluation values for the same text chunks are common, these evaluation values can be determined and inputted together, and from there, only the evaluation values related to the text chunks for each subquery can be used to calculate the search accuracy of each subquery.



3.13.6.3 Tuning the Combination Method of Hybrid Search

Based on the search accuracy of each subquery, you can tune the combination method of subqueries. Adjustable parameters include, for example, the weighting between full-text search and semantic text search, and the number of retrievals each subquery returns. If the semantic text search has a higher precision compared to full-text search, it is possible to increase the weight of semantic text search used in the final ranking calculation. Additionally, by setting the number of retrievals (num_results) in subqueries to a larger value than usual to obtain trace information, and calculating search accuracy such as recall and precision for each number of retrievals, recall@k, and precision@k, you can determine the appropriate number of retrievals from subqueries. As a result of evaluating search accuracy, you may detect issues that cannot be improved by the combination method. In such cases, tuning of both semantic text search and full-text search is performed.

Chapter 4 Graph Management Feature

This chapter describes the graph management feature that provide graph storage and search.

Refer to the Apache AGE documentation for more details.

4.1 Overview of Graph Management Feature

The graph management feature enables you to store and search property graphs in Fujitsu Enterprise Postgres. A graph is a data structure that represents relationships between entities using nodes and edges that connect the nodes. A property graph is a type of graph in which nodes and edges can have information called properties. Graphs that represent relationships can be searched based on relationships and properties. Graph searches are performed using openCypher, a query language for graph databases.

4.2 Installation of Graph Management Feature

Graph management features are provided by the OSS's Apache AGE.

4.2.1 Setting Up the Graph Management Feature

Refer to "Apache AGE" in the Installation and Setup Guide for Server to set up Apache AGE.

4.2.2 Removing the Graph Management Feature

Refer to "Apache AGE" in the Installation and Setup Guide for Server to remove Apache AGE.

4.3 Creating a Graph

Graphs are stored in Fujitsu Enterprise Postgres as a single virtual data object called a graph. Internally, they are saved as multiple database objects under a schema that has one-to-one correspondence with the graph. Graphs are created and deleted using the `create_graph` and `drop_graph` functions, which are SQL functions provided by the graph management feature. If the second argument of the `drop_graph` function is set to true, the database object under the schema that corresponds to the graph will also be deleted.

Example) Creating a graph

```
SELECT create_graph('new_graph');
NOTICE: graph "new_graph" has been created
create_graph
-----
(1 row)
```

Example) Deleting a graph

```
SELECT drop_graph('new_graph', true);
NOTICE: drop cascades to 2 other objects
DETAIL: drop cascades to table new_graph._ag_label_vertex
drop cascades to table new_graph._ag_label_edge
NOTICE: graph "new_graph" has been dropped
drop_graph
-----
(1 row)
```



Information

When you create a graph, a schema with the same name as the graph is created. Specify a name that does not overlap with existing schemas. In particular, you cannot use reserved names beginning with `pg_`. Also, do not use names beginning with `pgx_`.

The schema corresponding to a graph will be deleted when the graph is deleted. Do not create objects under this schema.

You can view a list of graphs and their corresponding schemas (namespaces) stored in a database with the graph management feature enabled by using the following method.

Example) Graph list

```
SELECT * FROM ag_catalog.ag_graph ;
 graphid |      name      | namespace
-----+-----+-----
   81957 | new_graph      | new_graph
   82576 | new_graph2     | new_graph2
   82598 | sample         | sample
(3 rows)
```

4.4 Storing Graph Data

Adding nodes and edges to a graph is done through Cypher queries using cypher functions, which are SQL functions provided by the graph management feature.

Example) Adding nodes and edges

```
SELECT * FROM cypher('new_graph', $$
CREATE (:Person {name: 'Daedalus'})-[ :FATHER_OF ]->(:Person {name: 'Icarus'})
$$) AS (a agtype);
a
---
(0 rows)
```

In addition to the above, you can also load from a CSV file using the SQL functions `load_labels_from_file` and `load_edges_from_file`.

4.5 Protecting Graph Data

Graphs are automatically protected by backup, multiplexing, or replication settings that specify the tablespace, database, or instance in which the graph is saved.

This section describes the points to keep in mind when protecting graph data using Fujitsu Enterprise Postgres features.

Within Fujitsu Enterprise Postgres, graphs are stored as multiple objects under a schema that has one-to-one correspondence with the graph. To protect a graph, configure protection for the database objects that make up the graph.

4.5.1 Encrypting the Graph

A graph consists of multiple database objects, and even after the graph is created by the `create_graph` function, a new table will be created when a new label is added. If you want to encrypt multiple database objects that make up a graph, including those that will be created in the future, set the default tablespace of the entire database that stores the graph to an encrypted tablespace.

Example) When encrypting the entire database for which this feature is enabled

```
CREATE DATABASE rag_database TABLESPACE = encrypted_tablespace;
CREATE EXTENSION age;
```

To encrypt a graph when you cannot use a single tablespace per database, temporarily change the default tablespace to an encrypted tablespace before the operation that creates the graph.

Example) Specifying table space and creating a graph

```
SET default_tablespace = 'secure_tablespace';
SELECT create_graph('graph1');
```

4.5.2 Restricting Access to Graph

Access to a graph is restricted by access control for the database objects that make up the graph. Access restrictions are set using the confidentiality management feature of Fujitsu Enterprise Postgres. Since graphs have a one-to-one correspondence with schema objects, you can allow or deny access to a graph by specifying access rights for that schema using the confidentiality management feature as follows:

1. Refer to the "Confidentiality Management" in the Security Operations Guide to define confidentiality management role, confidentiality matrix, confidentiality level, and confidentiality group.
2. Grant confidentiality privilege on the schema to the confidentiality group.

```
SELECT pgx_grant_confidential_privilege('rag_matrix', 'level1', 'group1', '{"schema":
[ "USAGE" ]}');
```

3. Add the schema corresponding to the graph as a confidentiality object to the confidentiality level.

```
SELECT pgx_add_object_to_confidential_level ('rag_matrix', 'level1',
'[ {
  "type": "schema",
  "object": [
    {
      "schema": "new_graph"
    }
  ]
}]');
```

4. Add roles to the confidentiality group you created to set access rights to the graph.

```
SELECT pgx_add_role_to_confidential_group('rag_matrix', 'group1', '['rag_user']');
```

If you want to set fine-grained access privileges, such as allowing only searches of graphs but not updating them, you can use SQL statements to directly set access privileges for database objects such as the tables that make up the graph.

The privileges required to access a graph are as follows:

- Privileges required to create a graph
 - CREATE privilege for the database
- Privileges required to create and delete new nodes and edges
 - CREATE and USAGE privilege for the schema with the same name as the graph name
 - When adding a node, ownership of the `_ag_label_vertex` table under the schema with the same name as the graph name
 - When adding an edge, ownership of the `_ag_label_edge` table under the schema with the same name as the graph name
 - UPDATE privilege for `_label_id_seq` under the schema with the same name as the graph name
- Privileges required to create and delete nodes and edges that use existing labels
 - USAGE privilege for the schema with the same name as the graph name

- When adding a node, privileges under the schema with the same name as the graph name
 - INSERT privilege (to create), SELECT privilege, and UPDATE privilege (to delete) for the `_ag_label_vertex` table or a table with the same name as the label name to be added
 - USAGE privilege for the `_ag_label_vertex_id_seq` sequence or a sequence with the same name as the label name to be added
- Privileges when adding an edge
 - INSERT (create), SELECT, and UPDATE (delete) privileges for the `_ag_label_edge` table under the schema with the same name as the graph name or the table with the same name as the label name to be added.
 - USAGE privilege for the `_ag_label_edge_id_seq` sequence or the sequence with the same name as the label name to be added.
- Privileges required to search graphs.
 - USAGE privilege for the `ag_catalog` schema.
 - USAGE privilege for the schema with the same name as the graph name.
 - SELECT privilege for the table object in the schema with the same name as the graph name.



Information

Graph data structures do not have the concept of rows and columns, so PostgreSQL's row-level security and column-based access control features cannot be applied.

4.5.3 Recording Access to Graph

The audit log feature of Fujitsu Enterprise Postgres is used to record access to the graph in the audit log. When specifying database objects to be audited individually, specify each database object that constitutes the graph.

When adding new labels to a graph, new corresponding tables are created, so those tables must also be specified as targets for auditing.

4.6 Searching Graph

Graph searches are performed using Cypher queries. Cypher queries are specified as strings as arguments to the `cypher` function, which is an SQL function.

For more information about Cypher queries, refer to the Apache AGE documentation.

Example) Graph search

```
SELECT * FROM cypher('new_graph', $$
MATCH (:Person {name: 'Daedalus'})-[:FATHER_OF]->(person)
WITH person.name AS name ORDER BY person.name RETURN name
$$) AS (v agtype);
      v
-----
"Icarus"
(1 row)
```

Graph searches allow you to retrieve specific property values or sets of properties for nodes or edges that match a condition. Properties are returned in `agType` data type added by the graph management feature, which is compatible with JSON types. Graph searches are closed to cypher functions, but the retrieved values can also be combined with searches of other tables in the database.

Example) Combination

```
SELECT (x::json->'properties')->'name' FROM cypher('new_graph', $$
MATCH (x)
```



```

RETURN x $$) AS (x agtype);
?column?
-----
"Daedalus"
"Icarus"
(2 row)

```

4.7 Adding Labels to Graph

Labels in a graph represent the type of node or edge, and are an important element in expressing the schema of knowledge data stored as a graph. Adding a label is easier than adding a column to a table, but you should plan carefully and consider the impact on your application. By default, only the database user who created the graph can add new labels.

When adding nodes or edges with new labels to a graph, new tables corresponding to each label are created. To encrypt a graph, change the default tablespace to an encrypted tablespace before adding labels. The access rights for these newly created tables are inherited from the access rights for the tables for unlabeled nodes or unlabeled edges, so no additional steps are required if you want to reuse those access rights.

You can check the tablespaces in which tables corresponding to each label of the graph nodes and edges are located with the following SQL.

```

SELECT g2.graphname AS graphname, c.relname AS relname, c.reltablespace AS reltablespace
FROM pg_class c,
rag_databse-> (SELECT g.name AS graphname, l.relation AS relation FROM
ag_catalog.ag_graph AS g, ag_catalog.ag_label AS l WHERE g.graphid = l.graph) AS g2 WHERE
c.oid = g2.relation;

```

graphname	relname	reltablespace
new_graph	_ag_label_vertex	0
new_graph	_ag_label_edge	0
new_graph	x	0
new_graph	y	0
new_graph2	_ag_label_vertex	0
new_graph2	_ag_label_edge	0
sample	_ag_label_vertex	80722
sample	_ag_label_edge	80722

(8 rows)

4.8 Performance Tuning of Graph Search

Graph data is stored internally as multiple tables. Graph data searches are implemented as SELECT statements on those tables. The access plan can be checked with the EXPLAIN statement in the same way as for normal SQL statements.

Graph data consists of a table with a record for each node and a table with a record for each edge, and queries corresponding to graph data are implemented by joining and filtering those tables. Therefore, the approach to performance issues in graph searches is the same as for SELECT statements that include joins.

Example) Checking the access plan of a graph search

```

SELECT * FROM cypher('new_graph', $$
EXPLAIN (COSTS off) MATCH (:Person {name: 'Daedalus'})-[:FATHER_OF]->(person)
WITH person.name AS name ORDER BY person.name RETURN name
$$) AS (v agtype);
QUERY PLAN
-----
Sort
Sort Key: (agtype_access_operator(VARIADIC ARRAY[_agtype_build_vertex(person.id,
_label_name('25995'::oid, person.id), person.properties), '"name"'::agtype]))
-> Hash Join
Hash Cond: (person.id = _age_default_alias_1.end_id)
-> Append

```

```

-> Seq Scan on _ag_label_vertex person_1
-> Seq Scan on "Person" person_2
-> Hash
-> Hash Join
Hash Cond: (_age_default_alias_1.start_id = _age_default_alias_0.id)
-> Seq Scan on "FATHER_OF" _age_default_alias_1
-> Hash
-> Seq Scan on "Person" _age_default_alias_0
Filter: (properties @> '{"name": "Daedalus"}'::agtype)
(14 rows)

```

4.9 Using Graph Data in Applications

The cypher function returns a value of the agType data type. The agType type is a subset of the json type. If you want to handle this data type directly in your application, please install a driver for each language in your application. If you do not use a driver, the agType type will be returned to your application as a string or JSON type.



See

For driver installation, refer to below.

<https://github.com/apache/age/tree/master/drivers>

4.10 Visualizing Graph Data

You can use the OSS age-viewer as a tool to visualize graph data.



See

For age-viewer installation, refer to below.

<https://github.com/apache/age?tab=readme-ov-file#graph-visualization-tool-for-age>

4.11 Internal Structure of Graph Data

When you create a graph, a corresponding schema is created, and tables, indexes, and sequence objects are created under it. For more information, see the Apache AGE documentation.

4.12 Quantitative Limits

The total number of nodes and edges must be less than or equal to $2^{48} - 1$. This restriction applies only to one particular graph.

4.13 Reference

The following SQL functions are provided, refer to Apache AGE documentation for details:

Function	Description
create_graph	Creating a graph
drop_graph	Deleting a graph
cypher	Manipulating graphs with Cypher queries
load_labels_from_file and load_edges_from_file	Loading the graph data

Chapter 5 Model Management in the Database

Vectorization is performed using pre-trained models called embedding models. By managing models in a database, you can integrate model deployment (loading) with database operations and achieve fine-grained access control at the model level. When managing models in a database, you can use ONNX format models that integrate a series of processes for vectorization into a single file. Additionally, to perform vectorization (inference by the model), an inference server called Triton Inference Server is used.

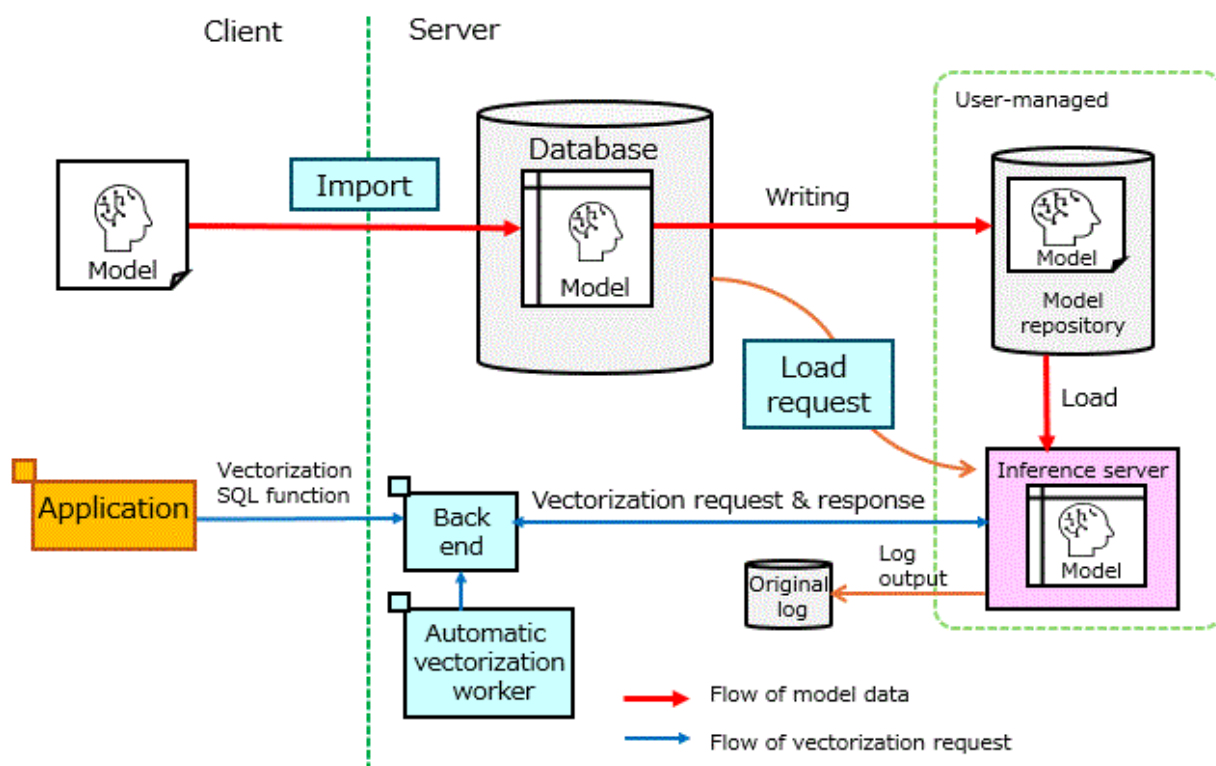
For the formats of models that can be handled by Fujitsu Enterprise Postgres, refer to the "[5.2.6.1 Database-side Model Specifications](#)". Also, for requirements regarding the inference server, refer to "[5.2.1 Setting Up Inference Server](#)".

5.1 Overview

Import the ONNX format text embedding model for vectorization into the database and operate it on an inference server running outside the database. This allows the model to be imported and made available for use through database operations alone, enabling vectorization and semantic text search with the model.

This feature provides an import command as a client feature to INSERT the model file into a table within the database. This operation is called import. After importing the model, a request is executed to make it ready for vectorization using the model. This is called a load request. When a load request is made, the model is loaded onto the inference server, making vectorization using the model possible. The inference server is the server that actually performs calculations based on the model.

Once the model is loaded, vectorization using the model becomes possible. The vectorization process is realized by requesting the inference server to transform text data into vector data and receiving the vector data. The backend process handles the request for vectorization to the inference server and the reception of results.



5.2 Introduction/Setup

5.2.1 Setting Up Inference Server

The Triton Inference Server to be used as the inference server for this feature must meet the following requirements.

- The version of Triton Inference Server must be 2.20.0 or later and meet the following requirements:
 - The following APIs are available.

For details about the API, refer to the official documentation of Triton Inference Server.

 - Model loading: RepositoryModelLoad
 - Model unloading: RepositoryModelUnload
 - Inference: ModelInfer/ModelStreamInfer
 - Model status check: RepositoryIndex
 - Model interface check: ModelConfig
 - The model repository has the structure described in the official documentation of Triton Inference Server.
 - The automatic completion feature of the configuration file is available.
- The model repository of the Triton Inference Server must be located in a place accessible as a local file by Fujitsu Enterprise Postgres.
- It must be accessible via gRPC API.
- ONNX Runtime must be enabled as the backend.
- Including onnxruntime-extensions built without Python, supporting custom operators for the tokenizer.
- The --model-control-mode option is specified as explicit.

Before using this feature, start the Triton Inference Server on the same machine as the database.

Fujitsu Enterprise Postgres provides a sample Dockerfile that meets the above requirements. To create a container image using this sample file, execute the following command. Set the appropriate label for the volumepath option (v option) so that the directory on the host side is visible from the container. "<x>" indicates the product version.

```
cp /opt/fsepv<x>server64/share/triton_dockerfile.sample ./triton_dockerfile
$ podman build -f ./triton_dockerfile -t triton_image
$ mkdir -p /path/to/model/repository
$ podman run -d --name triton_container -p8001:8001 -p8002:8002 -v /path/to/model/
repository:/models \
triton_image tritonserver --model-repository=/models --http-port=0 --grpc-port=8001 \
--metrics-port=8002 --backend-config=onnxruntime,device=cpu --model-control-mode=explicit \
--log-info=true --log-warning=true --log-error=true --log-verbose=0
$ podman container ls # Confirm that the container is running.
```

By using systemd, you can automatically start the created container. Below is an example of starting an inference server as a service using a sample file necessary for the automatic start of the container.

```
$ cp /opt/fsepv<x>server64/share/triton.container.sample \
~/.config/containers/systemd/triton.container
$ systemctl --user daemon-reload
$ systemctl --user start triton.service
```

To achieve model-level access control, configure mutual TLS authentication(mTLS authentication) on the inference server. For details, refer to "[5.4.3 Security](#)".

Enable metrics on the inference server to identify the cause when problems occur. Also, enable timestamp formatting for log output.

5.2.2 Setting Up pgx_inference

This feature is provided as an extension called pgx_inference. Additionally, this feature is used in conjunction with automatic vectorization and semantic text search by pgx_vectorizer. If you want to perform semantic text search and automatic

vectorization using a model imported into the database, set up `pgx_vectorizer` along with this feature. For instructions on setting up `pgx_vectorizer`, refer to ["3.2.2.2 Setting Up pgx_vectorizer"](#).

Setting parameters in the postgresql.conf file

Set the following parameters.

- `shared_preload_libraries`

Add `pgx_inference`. As a result, when the database instance is started, the load launcher is initiated internally.

- `max_worker_processes`

Add the number of databases to enable this feature +1. This is the number of worker processes that will be started with this feature.

- `pgx_inference.triton_model_repository_path`

Set the path of the model repository specified when setting up the inference server. You can only set the path on the machine where the database is running. It is recommended to prepare a dedicated directory for this function, as model files will be generated or deleted in the directory set in this parameter during loading.

- `pgx_inference.triton_grpc_port`

This is the port number used when sending requests to the inference server via gRPC. Set the gRPC port number specified during the setup of the inference server.

- `pgx_inference.triton_ort_extensions_library_filename`

Specify the absolute path of the shared library (onnxruntime-extensions) on the server where the Triton Inference Server is running.

If mTLS authentication is configured on the inference server, Set the following parameters.

- `pgx_inference.triton_use_ssl`

Specify whether to enable SSL communication with the Triton Inference Server. If enabled, set the following three parameters.

- `pgx_inference.triton_grpc_root_certificates`

Specify the path to the CA file for verifying the server certificate chain used in the gRPC connection.

- `pgx_inference.triton_grpc_certificate_chain`

Specify the path of the client certificate chain used for the gRPC connection.

- `pgx_inference.triton_grpc_private_key`

Specify the path of the client private key used for the gRPC connection.

Enabling `pgx_inference` extension

Execute `CREATE EXTENSION` for the database that will use this feature.

```
CREATE EXTENSION pgx_inference CASCADE;
```

Starting load launcher

Start the load launcher with the user who executed `CREATE EXTENSION`.

```
SELECT pgx_inference.pgx_launch_load_launcher();
```

5.2.3 Removing

Before deleting this feature, unload and delete all models used by this feature from the inference server.

Unload model

```
SELECT pgx_inference.pgx_unload_all_models();
```

The unload process is performed asynchronously. Refer to the `pgx_triton_model_status` view to confirm that the unload process is complete.

Delete model

After unloading is complete, delete all models.

```
SELECT pgx_inference.pgx_drop_all_models();
```

After unloading and deleting all models, connect to the database where this feature is used and execute `DROP EXTENSION`.

```
DROP EXTENSION pgx_inference;
```

When using semantic text search and automatic vectorization features with `pgx_vectorizer`, execute `DROP EXTENSION` for `pgx_vectorizer` with the `CASCADE` option specified.

```
DROP EXTENSION pgx_vectorizer CASCADE;
```

5.2.4 Resource Control

In this feature, the model is loaded into memory for calculations based on the model by the inference server. During actual calculations, CPU resources are used. In this feature, resources on the same machine are shared between the database and the inference server. Therefore, it is necessary to consider the impact of resource usage such as memory, CPU, and disk by the database and inference server on the main operations and appropriately control the resources. Additionally, for each resource, consider the following to ensure that the resources used by this feature are not insufficient.

- Memory
 - Total size of the model to be loaded
- CPU
 - Number of simultaneous vectorization processes
 - Parallelism of a single vectorization process
- Disk resources
(The directory specified in the `pgx_inference.triton_model_repository_path` parameter)
 - Total size of the model to be loaded

The memory used by this feature can be controlled with the following parameters.

- `pgx_inference.total_model_size_limit`

Set the upper limit for the total size of models that can be loaded into the inference server. Basically, use the default value of -1 (unlimited). Set a value other than unlimited when the memory size on the machine is small or there are multiple system administrators who can load models, and you want to limit the mass loading of models by those users.

Check and manage the official documentation of Triton Inference Server regarding the resources used by Triton Inference Server.

https://docs.nvidia.com/deeplearning/triton-inference-server/user-guide/docs/onnxruntime_backend/README.html

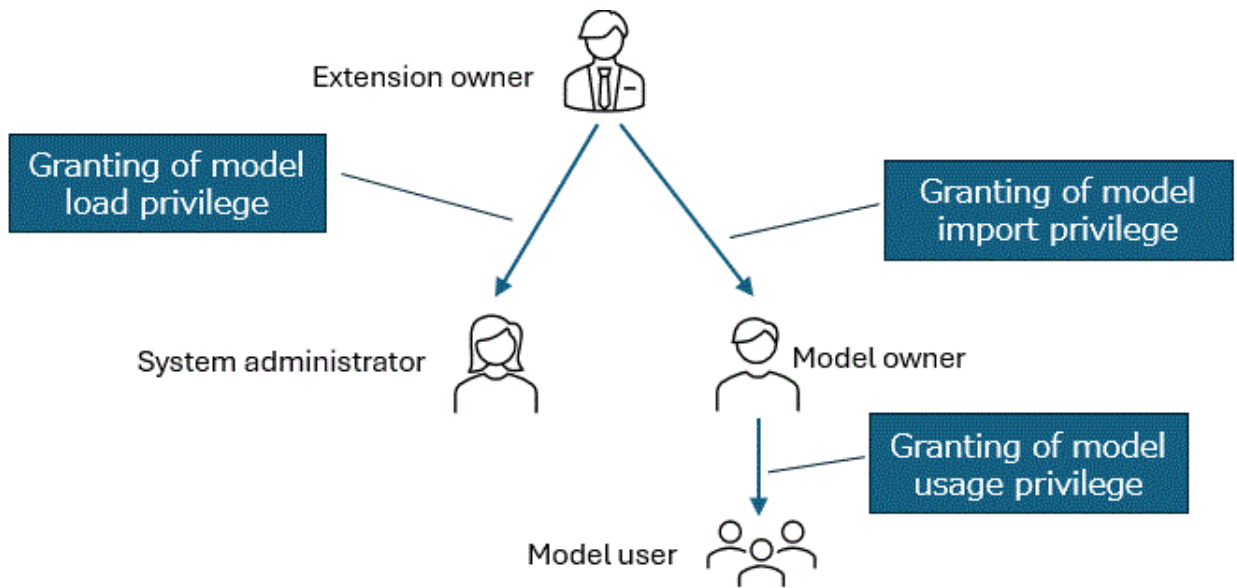
5.2.5 Setting Privilege

This feature is initially available only to the extension owner (the database user who enabled the extension) immediately after installation. For other database users to use this feature, privilege must be granted by the extension owner.

Users with model import privileges can import their models into the database. The user who imports the model becomes the owner of that model. Only users authorized by the model owner can use that model.

Loading a model requires separate model loading privileges from importing.

The extension owner grants model import and loading privileges to appropriate users. Meanwhile, the model owner grants model usage privileges to appropriate users. To utilize this extension feature, the necessary users are referred to as the extension owner, model owner, system administrator, and model user.



The roles of each user are as follows:

Extension owner: Grants appropriate privileges to each user.

Model owner: Manages each model and grants execution privileges for the model.

System administrator: Properly manages disk consumption and memory usage by the model.

Model user: Uses the model for inference (vectorization).

Additionally, each user is assigned "possible operations (execution of commands and functions)" and "privileges that can be granted or revoked to others". These are shown in the table below.

Operation/Privilege		Extension owner	Model owner	System administrator	Model user
Operation (Execution of commands and functions)	Import model	Possible	Possible	Impossible	Impossible
	Delete model	Possible	Only my model is possible	Possible	Impossible
	Update model	Possible	Only my model is possible	Impossible	Impossible
	Refer to model file	Possible	Only my model is possible	Impossible	Impossible
	Load/Unload model	Possible	Impossible	Possible	Impossible
	Reference to the total size of the loaded models	Possible	Impossible	Possible	Impossible
	Use model	Possible	Only my model is possible	Impossible	Only permitted models are possible
	Refer to metadata	Possible	Only my model is possible	Possible	Only permitted models are possible

Operation/Privilege		Extension owner	Model owner	System administrator	Model user
	Refer to model status	Possible	Only my model is possible	Possible	Only permitted models are possible
Granting/ revoking privilege to others	Import model	Possible	Impossible	Impossible	Impossible
	Delete model	Possible	Impossible	Impossible	Impossible
	Update model	Possible	Impossible	Impossible	Impossible
	Load/Unload model	Possible	Impossible	Impossible	Impossible
	Use model	Possible	Only my model is possible	Impossible	Impossible
	Execution privilege for referencing metadata (function)	Possible	Only my model is possible	Impossible	Impossible
	Execution privilege for referencing model status (function)	Possible	Only my model is possible	Impossible	Impossible

Granting and revoking the necessary privileges for model owners and system administrators is done using the `pgx_grant_model_role` function and the `pgx_revoke_model_role` function. These functions can only be executed by extension owners.

Example) Granting of privilege

```
SELECT pgx_inference.pgx_grant_model_role('target_user_name','model_owner') - Grant model owner privilege to the user named "target_user_name"
```

Example) Revoking of privilege

```
SELECT pgx_inference.pgx_revoke_model_role('target_user_name','model_owner') - Revoke model owner privilege to the user named "target_user_name"
```

Model users need EXECUTE privileges for the function corresponding to the model ID. The owner of the model should execute the GRANT and REVOKE statements to grant or revoke these privileges.

When a model is imported, it is stored in a table called `pgx_models`. At this time, a model ID is assigned to the imported model. A function is created for this model ID, which is called the function corresponding to the model ID. The function name corresponding to the model ID can be obtained using the `pgx_get_function_by_model_name` function.

```
SELECT pgx_inference.pgx_get_function_by_model_name('target_model');
pgx_get_function_by_model_name
-----
pgx_embedding_1
(1 row)

GRANT EXECUTE ON FUNCTION pgx_inference.pgx_embedding_1(text) TO target_user_name;
```

To use the functions provided by this feature, you need USAGE privilege on the `pgx_inference` schema. These privileges are granted to system administrators and model owners through the `pgx_grant_model_role` function. For model users, the extension owners needs to grant USAGE privilege. As shown in the table above, simply granting USAGE privilege does not allow users to reference or use the model, so it is possible to grant USAGE privilege in advance.

To achieve model-level access control, access control in Triton Inference Server is also necessary. For more details, refer to ["5.4.3 Security"](#).

5.2.6 Model Preparation

This feature allows you to import ONNX pipeline models. An ONNX pipeline model is an embedding model that takes text as input and generates embeddings as output. The core of a text embedding model is a neural network model expressed in ONNX format or similar (hereafter referred to as the core model). However, to perform text embedding (vectorization), not only the core model is needed, but also preprocessing such as tokenization and post-processing such as pooling. There are various methods for tokenization and pooling, but they are not independent of the core model and need to be used in combination as intended by the developer. Fujitsu Enterprise Postgres uses a model called the ONNX pipeline model, which consolidates the entire pipeline process for performing text embedding into a single ONNX format file. The specifications for ONNX pipeline models that can be imported with this feature are defined as ["5.2.6.1 Database-side Model Specifications"](#).

5.2.6.1 Database-side Model Specifications

You can use a model with the following specifications as an embedded model.

Item	Constraint	Remarks
Model format	ONNX format	
Model type (semantics)	Inference Model (Stateless Model)	It is not a learning model (not a stateful model).
Input tensor	1-dimensional tensor of STRING type. It may optionally have an additional dimension for batch processing. It should not require any other inputs (without default values).	
Output tensor	1-dimensional tensor of FLOAT16 type. It may optionally have an additional dimension for batch processing.	Because the vector type in Fujitsu Enterprise Postgres is a single-precision floating-point type.
Operator set and operator	Only operators included in the following domain/version operator sets should be used. ai.onnx: Version X to Y ai.onnx.ml: Version X to Y (Operators provided by onnxruntime-extensions) Ensure that the versions of ai.onnx and ai.onnx.ml are supported by the ONNX Runtime that the Triton server uses as its backend.	Even if the conditions are met, operators not provided by onnxruntime's CPUProvider cannot be used.
Model size	4TB or less	

5.2.7 Model Import

A user with privilege can import a model that meets the database-side model specifications into the database by executing the import subcommand of the `pgx_aimodel_tool` command. The user who imports the model becomes the owner of that model.

The imported model name and metadata can be checked by executing `pgx_get_model_metadata` function.

For the model name of the model to be imported with this function, use the following characters. Including characters other than the following may result in failure to load or call the API.

- Alphanumeric (A-Z, a-z, 0-9)
- Underscore (_)
- Hyphen (-)
- Dot (.)

However, the following usages are not allowed.

- Dots cannot be used at the beginning or end. It may be mistaken for a hidden/special name.
- Consecutive dots (..) cannot be used as they may cause path interpretation issues.
- The same model name cannot be used within an instance.



See

For details on the `pgx_aimodel_tool` command, refer to "`pgx_aimodel_tool`" in the Reference.

5.2.8 Model Load/Unload

Users who qualify as system administrators for this feature can request the loading of imported models into the database by executing the `pgx_load_model` function.

When a model load is requested, the load launcher running in the database loads the specified model file into the inference process. At this time, the model file is output to the directory set in `pgx_inference.triton_model_repository_path` parameter.

When the `pgx_unload_model` function is executed, it requests unloading for the model specified in the argument. As with loading, the requested model is unloaded from the inference server by the load launcher. Once the model is unloaded, vectorization and semantic text search using that model can no longer be performed.

To check whether the model is available/unavailable after executing the `pgx_load_model` function and the `pgx_unload_model` function, check the `pgx_triton_model_status` view. If the model is unavailable, the reason will be displayed in the reason column of the `pgx_triton_model_status` view.

Models requested to be loaded will automatically load even after instance restart. Also, load requests propagate to all standby servers and are loaded on each server. Unload requests are also synchronized to all standby servers.

5.2.9 Definition of Vectorization

To achieve semantic text search using the imported model with this feature, specify the imported model when defining the vectorization. You can use the imported model for vectorization by specifying `pgx_vectorizer.pg_embedding_onnx` in the embedding argument of `pgx_vectorizer.pg_create_vectorizer`. Specify the model name and the number of dimensions of the generated vector in `pgx_vectorizer.pg_embedding_onnx`.

If an unavailable model is specified when executing `pgx_create_vectorizer`, a warning message will be output. Follow the warning message to load the specified model.

5.2.10 Deletion of Imported Model

The imported model can be deleted using the `pgx_drop_model` function. The `pgx_drop_model` function specifies the model name. Model deletion can be performed by the model owners and system administrators.

If the model is in use by the automatic vectorization feature, the model cannot be deleted. The model cannot be deleted if it is loaded.

5.3 Usage Example

5.3.1 Model Import

Import the model into the database.

Import the model file on the client.

```
pgx_aimodel_tool import -f ./model.onnx -n sample-model-v1 -F embedding -m
'{"version":"v1"}' -d rag_db -h localhost -p 27500 -U model_owner
```

Check the imported model

```
SELECT * FROM pgx_inference.pgx_list_model_metadata();
id |      name      | function |      metadata      | model_file_oid |
file_size | owner |
-----+-----+-----+-----+-----+
+-----+-----+
1 | sample-model-v1 | embedding | {"metadata": "sample1"} | 16767 | 90000000 |
16329|
2 | sample-model-v2 | embedding | {"metadata": "sample2"} | 16782 | 60000000 |
16330|

(1 row)

SELECT * FROM pgx_inference.pgx_get_model_metadata('sample-model-v1');
id |      name      | function |      metadata      | model_file_oid |
file_size | owner |
-----+-----+-----+-----+-----+
+-----+-----+
1 | sample-model-v1 | embedding | {"metadata": "sample1"} | 16767 | 90000000 |
16329|

(1 row)
```

Example) When an existing model name is specified as an argument

```
pgx_aimodel_tool import -f ./model.onnx -n sample-model-v1 -F embedding -m
'{"version":"v1"}' -d rag_db -h localhost -p 27500 -U model_owner
pgx_aimodel_tool: error: query failed: ERROR: model "sample-model-v1" already exist
HINT: To overwrite the existing model, specify the "--update-on-conflict" option.
```

Example) When overwriting an imported model with the same name

```
pgx_aimodel_tool import -f ./model.onnx -n sample-model-v1 -F embedding -m
'{"version":"v1"}' --update-on-conflict -d rag_db -h localhost -p 27500 -U model_owner
pgx_aimodel_tool: model "sample-model-v1" import completed
```

5.3.2 Model Load/Unload

Load/Unload the imported model.

Example) Model load request and retrieval of list

```
SELECT pgx_inference.pgx_load_model('sample-model-v1');
pgx_load_model
-----

SELECT * FROM pgx_inference.pgx_list_model_load_requests();
      Name      | requests |      last_updated      |
-----+-----+-----+
sample-model_v1 | load     | Wed Aug 02 19:09:16.200357 2025 PDT
```

After loading the target model, if the total size of all loaded models exceeds the value set in `total_model_size_limit`, the loading will result in an error.

```
SELECT pgx_inference.pgx_load_model('model_name');
ERROR:  could not load model "good"
DETAIL:  toatl model size exeed 1000MB
```

Example) Model unload request

```
SELECT pgx_inference.pgx_unload_model('model_name');
pgx_unload_model
-----
```

5.3.3 Definition of Vectorization

To define vectorization using a model imported into the database, use the `pgx_vectorizer.pgx_embedding_onnx` function. In this example, vectorization is defined for the text format knowledge data table "sample_table".

```
SELECT pgx_vectorizer.pgx_create_vectorizer(
'sample_table'::regclass,
destination => 'sample_embeddings',
embedding => pgx_vectorizer.pgx_embedding_onnx('sample-model-v1',384),
chunking => ai.chunking_recursive_character_text_splitter('contents'),
processing => ai.processing_default(batch_size => 200, concurrency => 1),
scheduling => pgx_vectorizer.schedule_vectorizer(interval '1 hour'),
indexing => ai.indexing_hnsw(min_rows =>50000, opclass => 'vector_cosine_ops'),
);
```

5.3.4 Deletion of Imported Model

Delete the model imported into the database.

```
SELECT pgx_inference.pgx_drop_model('model_name');
pgx_drop_model
-----
```

5.4 Operaion

5.4.1 Backup/Restore

All data saved in the tables created by this feature will be subject to backup. For the tables created by this feature and the data saved in them, refer to "[5.5.2 Tables/Views](#)".

Files generated in the directory set in `pgx_inference.triton_model_repository_path` parameter do not need to be backed up. This is because these files are created based on the data stored in the table each time the model is loaded.

If the restore destination instance or database is using this feature, you need to unload and delete the models used in the restore destination instance or database from the inference server before performing the restore. Follow the procedures described in "[5.2.3 Removing](#)" to unload and delete all models.

When restoring backup data created by the `pg_dump` command or the `pg_dumpall` command, perform it on a database that is not using this feature. If the target database for restoration already has a model imported, restoration may not be possible. The procedures for restore varies depending on whether or not the semantic text search and the automatic vectorization feature are used.

When using semantic text search and the automatic vectorization feature

Backups are obtained using the `pg_dump` command. You cannot restore using a backup obtained with the `pg_dumpall` command.

When performing a restore, the parameters listed in "[5.2.2 Setting Up pgx_inference](#)" must be set on the instance where the restore is being performed. Additionally, restore separately for each schema where extensions are created.

1. Obtain a backup using the `pg_dump` command.
Since the model file is saved as a large object, explicitly specify the `-b(--large-objects)` option.

Example)

```
pg_dump -U postgres -d olddb -b -Fc -v -f ./backup_file.dmp
```

2. Restore only definitions with the `--schema-only` option of the `pg_restore` command.

Example)

```
pg_restore -U postgres -d newdb --schema-only ./backup_file.dmp
```

3. Restore data other than `pgai` and `pgx_vectorizer` using the `pg_restore` command.

Example)

```
pg_restore --exclude-schema=ai --exclude-schema=pgx_vectorizer --data-only -d newdb ./backup_file.dmp
```

4. Disable triggers to avoid unique constraint violations with the `pg_restore` command, and restore the data for `pgai` and `pgx_vectorizer`.

Example)

```
pg_restore --schema=ai --schema=pgx_vectorizer --data-only -d newdb --disable-triggers ./backup_file.dmp
```

When not using semantic text search and the automatic vectorization feature

Use the backup obtained with the `pg_dump` command to restore with the `pg_restore` command.

Since model files are saved as large objects, explicitly specify the `-b(--large-objects)` option of the `pg_dump` command. When performing a restore, the parameters listed in "[5.2.2 Setting Up pgx_inference](#)" must be set on the instance where the restore is being performed.

Example)

```
pg_dump -U postgres -d olddb -b -Fc -v -f ./backup_file.dmp
pg_restore -d newdb ./backup_file.dmp
```

Additionally, you can use the backup obtained with the `pg_dumpall` command to restore with the `psql` command.

Example) Backup the instance to be backed up using the `pg_dumpall` command

```
pg_dumpall -U postgres > pg_dumpall.sql
```

Example) Execute a restore on the instance to be restored

```
psql -U postgres -f pg_dumpall.sql
```

5.4.2 Streaming Replication/Database Multiplexing

When using this feature in a streaming replication or database multiplexing configuration, set up and start the Triton Inference Server on the standby server as well.

In a streaming replication or database multiplexing system, the following operations can only be performed on the primary server.

- Model import
- Model load request/unload request
- Deletion of imported model
- Granting access privilege to model/revoking access privilege

These operations performed on the primary server are all replicated to the standby server.

5.4.3 Security

The access privilege for the data stored in the table of this feature are automatically set according to the privileges described in ["5.2.5 Setting Privilege"](#).

You can use PostgreSQL's security features for the data stored in the table of this feature. If you use transparent data encryption, place the database that uses this function in an encrypted tablespace. This is because the model file is saved in a system catalog called `pg_largeobject`.

If you want to restrict users who can use the model on a per-model basis, configure the inference server to allow access only from Fujitsu Enterprise Postgres. Specifically, enable mTLS authentication on the inference server and allow only the client certificate used by Fujitsu Enterprise Postgres.



Point

.....

The model files created in the directory specified by the `pgx_inference.triton_model_repository_path` parameter during model loading are not protected by transparent data encryption. If you want to reduce the risk of model file leakage even slightly, specify a volatile temporary area such as `tmpfs` for `pgx_inference.triton_model_repository_path` parameter, and promptly delete unnecessary model files yourself after use. If high confidentiality is required, consider encrypting the entire file system or disk.

.....

5.4.4 Monitoring/Tuning

5.4.4.1 Monitoring Target

To monitor the state of the load launcher, refer to `pg_stat_activity`.

If you want to check the total size of model files loaded on the inference server, execute the `pgx_get_total_model_size` function. Additionally, the status of models loaded on the inference server can be checked in the `pgx_triton_model_status` view.

The logs output by the inference server and the PostgreSQL server logs are output as separate logs. Check the log file specified during the setup of the inference server. Refer to ["3.6 Monitoring Vectorization Processing for Semantic Text Search"](#) or ["3.10 Performance Tuning of Semantic Text Search"](#) to monitor vectorization and semantic text search by the imported model, and if there is a possibility of failure, check the server logs. Refer to the timestamps output in the logs to link the PostgreSQL server logs with the inference server logs.

In this feature, the model imported into the database is managed with the prefix `"pgx_"` added to the model name managed by Fujitsu Enterprise Postgres on the Triton Inference Server, and the suffix `"System Identifier (hereinafter referred to as SystemID)"` which is a system-specific value generated when creating the database cluster. For example, if imported with the name `"sample_model"` into a database cluster with SystemID `"7553888845569639366"`, it will be treated as `"pgx_sample_model_7553888845569639366"` on the Triton Inference Server. If an error message related to a specific model is output in the server log of Fujitsu Enterprise Postgres, check the log related to the model with the prefix `"pgx_"` added to the model name in the Triton Inference Server log.

The SystemID of the database cluster can be checked with the following SQL.

```
SELECT system_identifier FROM pg_control_system();
```

5.4.4.2 Parameter Tuning

The directory set in `pgx_inference.triton_model_repository_path` parameter will output large model files. To avoid the backup size from becoming bloated, set it outside the database cluster.

For GUC parameters related to memory, basically use the default values. If there are other tasks that you want to prioritize over this function, set each parameter to a value other than the default. If the machine's memory size is small or there are multiple system administrators with model loading privileges, loading the model may deplete memory and affect processes

other than this function. In this case, set a value other than unlimited for `pgx_inference.total_model_size_limit` parameter and impose an upper limit on the total size of models that can be loaded.



See

For more information on managing the resources used by the inference server, refer to the official documentation of Triton Inference Server.

https://docs.nvidia.com/deeplearning/triton-inference-server/user-guide/docs/onnxmlruntime_backend/README.html

5.5 Reference

5.5.1 Functions

Function	Return type	Description
<code>pgx_list_model_metadata()</code>	SETOF <code>pgx_models</code>	Get a list of imported models.
<code>pgx_get_model_metadata(model_name text)</code>	SETOF <code>pgx_models</code>	Get detailed information of the specified model.
<code>pgx_get_total_model_size()</code>	<code>bigint</code>	Get the total size of the currently loaded model files.
<code>pgx_load_model(model_name text)</code>	<code>void</code>	Execute the load request for the specified model. If the specified model is not imported, and after loading the model specified in the argument, the total size of all loaded models exceeds the value set by the <code>pgx_inference.total_model_size_limit</code> parameter, the execution will fail.
<code>pgx_unload_model(model_name text, force boolean DEFAULT FALSE)</code>	<code>void</code>	Executes the unload request for the specified model. If the specified model is not imported, and there is a vectorizer referencing the specified model, execution will fail. If the argument <code>force</code> is set to true, the unload request is executed forcibly.
<code>pgx_get_model_load_request(model_name text)</code>	SETOF record	Retrieve records from the <code>pgx_model_load_requests</code> table that correspond to the specified model. The model name is returned instead of the model ID.
<code>pgx_list_model_load_requests()</code>	SETOF record	Retrieve load and unload requests for all models. The model name will be returned instead of the model ID.
<code>pgx_drop_model(model_name text)</code>	<code>void</code>	Delete the imported model.
<code>pgx_grant_model_role(user_name text, role text)</code>	<code>void</code>	Grant the privileges held by the specified role to the user specified by <code>user_name</code> . This function can only be executed by the user who performed <code>CREATE EXTENSION</code> for this feature. The role can be specified as <code>model_owner</code> , meaning the model

Function	Return type	Description
		owner, or system_admin, meaning the system administrator.
pgx_revoke_model_role(user_name text, role text)	void	Revoke the privileges held by the specified role from the user specified by user_name. This function can only be executed by the user who created the extension for this feature. The role can be specified as model_owner, meaning the model owner, or system_admin, meaning the system administrator.
pgx_get_function_by_model_name(model_name text)	text	Get the name of the function created for the model specified in model_name.
pgx_onnx_embed(model text, text_to_infer text)	vector	The specified model is used to convert the text specified in text_to_infer into vectors. If the model specified in the model is not imported or is unavailable, the execution will fail.
pgx_launch_load_launcher()	void	Start the load launcher, which is a worker process that executes load processing.
pgx_unload_all_models()	void	Execute unload requests for all models. Even if there is a vectorizer referencing the model, a message will be output and the unload request will be executed. This function can only be executed by the extension owner.
pgx_drop_all_models()	void	All models will be deleted. Even if you are using models with the automatic vectorization feature, a message will be output and deletion will be performed. This function can only be executed by the extension owner. This function should only be used when returning this feature to an unused state.

5.5.2 Tables/Views

pgx_models table

It can only be referenced as the return value of the pgx_list_model_metadata function and the pgx_get_model_metadata function.

Column	Type	Constraint	Description
id	bigint	PRIMARY KEY	ID assigned to the model
name	text	UNIQUE NOT NULL	Name of the imported model

Column	Type	Constraint	Description
function	text	NOT NULL	The process executed by the imported model. Only embedding can be specified.
metadata	jsonb	NOT NULL	Metadata of the imported model
model_file_oid	oid	NOT NULL	The oid of the model file imported as a large object
file_size	bigint	NOT NULL	Size of the imported model file
owner	oid	NOT NULL	Model owner

pgx_model_load_requests table

It can only be referenced as the return value of the `pgx_list_model_load_requests` function and the `pgx_get_model_load_request` function.

Column	Type	Constraint	Description
id	bigint	PRIMARY KEY FOREIGN KEY (id) REFERENCES model_metadata ON DELETE CASCADE	Model ID
request	text		Request sent to the model
last_updated	timestamp with timezone	NOT NULL DEFAULT clock_timestamp()	The time the request was sent

pgx_triton_model_status view

Column	Type	Constraint	Description
id	bigint	PRIMARY KEY FOREIGN KEY (id) REFERENCES model_metadata ON DELETE CASCADE	Model ID
is_available	boolean	NOT NULL	Whether the model is available or not TRUE: Available FALSE: Not Available
reason	text		Reasons why the model is available or unavailable

5.5.3 Parameters

Parameters to set in the `postgresql.conf` file are shown.

Parameters related to memory

- `pgx_inference.total_model_size_limit` (integer)

Set the upper limit for the total size of model files that can be loaded. The default value is -1 (unlimited). If the unit is not specified, it is considered as MB.

Parameters related to the inference server

- `pgx_inference.triton_model_repository_path` (string)

Set the path of the model repository specified during the setup of the inference server used in conjunction with Fujitsu Enterprise Postgres. This cannot be omitted. If omitted, the database instance will fail to start. This directory must have permissions set so that the OS user starting Fujitsu Enterprise Postgres can read and write. If the privileges are inappropriate, a warning appears in the server log.

- `pgx_inference.triton_ort_extensions_library_filename` (string)

This corresponds to the `model_operations.op_library_filename` of the Triton Inference Server. Specify the absolute path of the shared library (`onnxruntime-extensions`) on the server where the Triton Inference Server is running. To reflect changes in the value, restart the database instance. It cannot be omitted. If omitted, the database instance will fail to start.

- `pgx_inference.triton_grpc_port` (integer)

This is the port number used when sending requests to the inference server. set the port number specified during the setup of the inference server. The default value is 8001.

- `pgx_inference.triton_use_ssl` (boolean)

Specify whether to enable SSL communication with the Triton Inference Server. The default value is false.

- `pgx_inference.triton_grpc_root_certificates` (string)

Specify the path to the CA file for verifying the server certificate chain used for gRPC connection. If the `pgx_inference.triton_use_ssl` parameter is true, it cannot be omitted. If omitted, the instance will fail to start.

- `pgx_inference.triton_grpc_certificate_chain` (string)

Specify the path of the client certificate chain used for gRPC connection. If the `pgx_inference.triton_use_ssl` parameter is true, it cannot be omitted. If omitted, the instance will fail to start.

- `pgx_inference.triton_grpc_private_key` (string)

Specify the path to the client private key used for gRPC connection. If the `pgx_inference.triton_use_ssl` parameter is true, it cannot be omitted. If omitted, the instance will fail to start.