

**КРАТКОЕ ОПИСАНИЕ
ПРОГРАММНОГО ОБЕСПЕЧЕНИЯ
BI.QUBE
METASTAGING**

ОГЛАВЛЕНИЕ

ВВЕДЕНИЕ.....	3
ГЛОССАРИЙ.....	3
1. ЦЕЛИ И НАЗНАЧЕНИЕ METASTAGING	4
2. АРХИТЕКТУРА METASTAGING	5
2.1. Описание компонентов системы	5
2.1.1. Утилита шифрования	6
2.1.2. Генерация JSON для поля Credentials в таблице stg.source	8
2.1.3. Экстрактор	9
2.1.4. Формирование слоя хранения данных в Greenplum.....	9
2.2. Структура проекта MetaStaging в файловой системе.....	12
2.3. Структура Базы данных.....	13
2.3.1. Настраиваемые таблицы MetaStaging.....	14
2.3.2. Таблицы логов MetaStaging	18
2.3.3. Таблицы дескрипторы	22
2.3.4. Хранимые процедуры и функции MetaStaging.....	23

ВВЕДЕНИЕ

Компонент MetaStaging позволяет консолидировать в стейджинговом слое хранилища данные из гетерогенных источников с поддержанием целостности и унифицированности метаданных, также уменьшает нагрузку на операционные базы при выполнении запросов, а кроме того, обеспечивает надежное подключение различных БД из разнородных источников для помещения данных в единый слой стейджинга (staging area) с поддержанием целостности метаданных в системе-назначения.

В документе приведено описание компонента и принципы работы с ним. Рассмотрены примеры загрузки данных с помощью компонента из разных источников.

Изучение данного документа позволит понять принцип работы компонента.

ГЛОССАРИЙ

1.	MetaStaging BI.Qube -	Инструмент, предназначенный для транспортировки данных.
2.	Хранимая процедура	Объект базы данных, представляющий собой набор SQL-инструкций, который компилируется один раз и хранится на сервере
3.	Представление	Виртуальная таблица, содержимое которой определяется запросом
4.	Бизнес-представление	Представление, в котором собраны Hub, Satellite и Link для сущности
5.	Материализация	Процесс сохранения результата запроса бизнес-представления в таблицу для ускорения выборки.
6.	Инкрементальная загрузка (загрузка с параметрами)	Регулярная загрузка данных в Greenplum. Извлекаются актуальные данные с даты последней загрузки. В таблице stg.session базы данных settings.db можно отследить историю всех загрузок.
7.	Полная загрузка (снэпшоты)	Загрузка данных в Greenplum без параметризации. Применяется, когда необходима полная перезагрузка всех данных в таблице на источнике (например, при отсутствии столбца, подходящего для секционирования).
8.	Полная загрузка с сохранением истории	Загрузка данных в Greenplum без параметризации. Но представления на Greenplum перенацеливаются на новые Parquet-файлы, а старые не удаляются из S3.
9.	Профиль	Добавляются в таблице stg.profile
10.	Экстрактор	Компонент системы для извлечения данных из источников в S3. Исполняемый файл находится в директории LoadingToS3. Вызывается в Airflow в соответствии с командами в настройке БД (settings.db).

11.	External Table (ET)	Вид таблицы в Greenplum, обеспечивающий доступ к внешним источникам данных, как к объекту самой БД Greenplum. В системе используется для получения доступа к файлам в S3. Используется фреймворк PXF.
12.	Сервисные процедуры	Процедуры, вызываемые автоматически в процессе работы компонента.

1. ЦЕЛИ И НАЗНАЧЕНИЕ METASTAGING

Цель MetaStaging – обеспечить транспортировку данных из систем источников в файловое S3-совместное хранилище данных (HDFS, ObjectStorage) с автоматической генерацией в СУБД Greenplum объектов типа «представление» на каждый полученный файл хранилищем.

Компонент MetaStaging, предназначен для передачи данных из различных источников, как правило, из учетных систем в целевое корпоративное хранилище данных (КХД) с поддержкой целостности метаданных систем-источников, при формировании промежуточного физического слоя хранения учитываются особенности целевой платформы.

Компонент MetaStaging входит в состав системы BI.Qube и может эксплуатироваться как отдельный компонент, так и в составе системы, так и под управлением компонента MetaOrchestrator, в такой конфигурации использование компонента является наиболее эффективной.

Компонент MetaStaging для развертывания, функционирования и настройки использует различные программные инструменты и фреймворки. Обязательным условием является наличие у них открытого исходного кода.

Поддерживаемые операционные системы: Linux (различные дистрибутивы, такие как Ubuntu, Mint, РЕД ОС), другие Unix-подобные системы, а также есть возможность развернуть компонент под Windows.

Настроечные данные компонента могут храниться посредством СУБД: PostgreSQL (9.0 и позднее) / Postgres Pro (10.22 и позднее) / Arenadata Postgres (ADPG) (14.2.1) / Greenplum на выбор заказчика.

Для тестирования корректности загрузки данных в S3 хранилище с помощью файлов «.parquet» использовались инструменты s3-browser и parquet-viewer.

Инструменты разработки DBeaver, Visual Studio Code

Среды выполнения Python, «.Net Core».

В качестве библиотек для взаимодействия с системами источниками и назначениями, а также для обеспечения интеграции данных используются:

- AWSSDK.S3 - Amazon Simple Storage Service для Amazon S3 (nuget.org)
- CommandLineParser - Terse syntax C# command line parser for .NET.
- ExcelDataReader - Lightweight and fast library written in C# for reading Microsoft Excel files

- Google.Cloud.BigQuery.V2 - Recommended Google client library to access the BigQuery API.
- Microsoft.Data.SqlClient - Provides the data provider for SQL Server.
- MySql.Data
- Newtonsoft.Json - Json.NET high-performance JSON framework for .NET
- Npgsql - is the open source .NET data provider for PostgreSQL.
- ParquetSharp -.NET library for reading and writing Parquet files.
- YandexDisk.Client - .NET library wrapper of Yandex Desktop RestAPI.
- Встроенные модули из стандартной библиотеки Python.
- Psycopg2 – PostgreSQL database adapter for the Python programming language.

В связи с высокой сложностью развертывания компонента в среде целевой СУБД установку компонента осуществляет вендор.

2. АРХИТЕКТУРА METASTAGING

2.1. Описание компонентов системы

Принцип работы MetaStaging сводится к взаимодействию программных блоков, которые отображены на рисунке ниже.

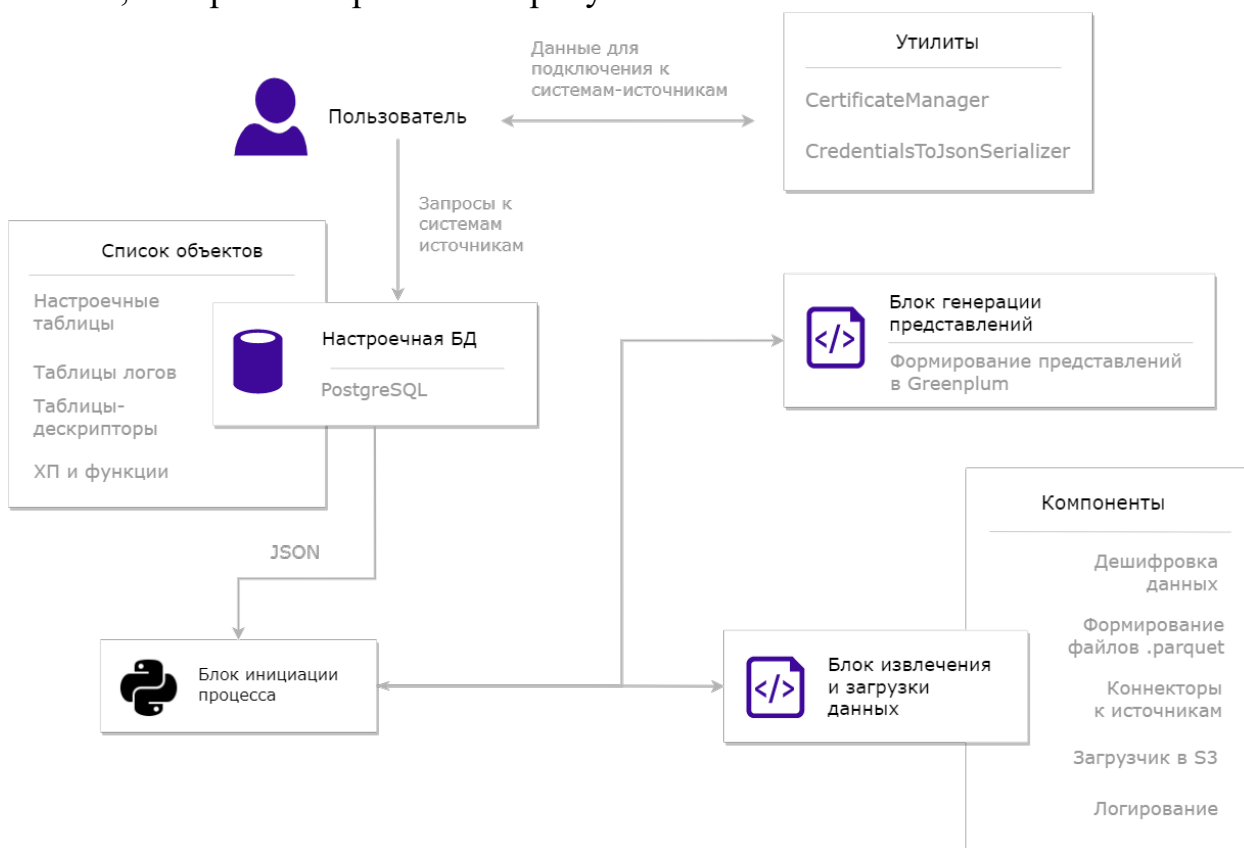


Рисунок 1. Компоненты системы MetaStaging

Краткое описание и назначение основных блоков компонента MetaStaging:

- Блок инициации процесса. Представляет собой Python3-скрипт и отвечает за запуск и координацию остальных блоков для интеграции данных.
- Блок извлечения и загрузки данных (Экстрактор). Представляет собой сборку «.Net Core». Загрузка может осуществляться в S3-совместимое хранилище в файлы «.Parquet».
- Блок генерации представлений. Отвечает за генерацию External tables и представлений в Greenplum, поддерживающих метаданные источников.
- Настроечная БД. Хранит информацию, необходимую для загрузки данных. Также служит интерфейсом для взаимодействия пользователя с MetaStaging. (см. п. Структура Базы данных).
- Утилиты. Предназначены для упрощения процесса заполнения настроечных таблиц
 - CertificateManager (Утилита шифрования) (см. п. Утилита шифрования).
 - CredentialsToJsonSerializer (Генератор Json для credentials источника) (см. п. Генерация JSON для поля Credentials в таблице stg.source).

2.1.1. Утилита шифрования

Шифрование ключей, паролей, строк подключения производится в ручном режиме с помощью программы CertificateManager. Программа расположена в директории /home/itpro_admin/CertificateManager/. Публичный и приватные ключи находятся в директории /home/itpro_admin/keytabs/

```
itpro_mgmt@prd-bietl-01:~/CertificateManager$ ./CertificateManager --help
CertificateManager 1.0.0
Copyright (C) 2022 CertificateManager

encrypt      Зашифровать текст
decrypt      Расшифровка текста
generate     Генерация двух ключей
help         Display more information on a specific command.
version      Display version information.
```

Рисунок 2. Пример применения утилиты

Пример команды для шифрования текста:

```
~/CertificateManager$ ./CertificateManager encrypt --public-key
../keytabs/public.crt --text "text to encrypt"
```

```
itpro_mgmt@prd-bietl-01:~/CertificateManager$ ./CertificateManager encrypt --public-key ../keytabs/public.crt --text "text to encrypt"
Зашифрованный текст: 'fePw4my/Ru8Iec0yda4SsAGU18DSgsBk9obxtVb0ASKHRAQP+oRyYAC4qf3Up7LiXbzhIEvOuZw0d0gmn/z+NLuNK1lP4mm2w1
9at+HugxJ7AZkGM0LDyaJ9NB05DbaC6UnH621zrpvyTGAYRR00Q+zf0ygtPzEngKp15CvzT2b1H0ju5VCYdwFzhNVKCPngzy/wqt1NBpC6320h7UhrReXN
bdUuDjyPmVJnpGwqx/k7T0oXSAPE618VKAwaQK0lPmn5lvr7X9M+MPFqGSpP9e73I62EZBu0kXaPv+Vzt0j+i4dqbaNKJZCK7xlmF4xIFSxe0d0FbNyeNwx1
1HkiLF7q7p0Xy97o5lXgQJw3WKhSzwIY+xUNtOrStmg7sKptEfuXi+11HV9j+eg5Wn6j9crNANkPGH+9UgX1gfZnez1b9mxRVRyZeG/QOZ405gmy2AAK1h2+
5Q6ycwtGv3kKVarcU6U5wfDZFyYnsRLQ00ecYRR2BQn+1tIBtCdc4pgZQfwfW/iW/namLq5y07NykFgu0fdqDGOSEHXfRcfHc9BPg8R0EBLc5+CnthW4wA90
6ESKGWCdq6qizbyejAw26o8WF/hwDJS2MyfICNwHihlRl3GFpEYOBHkRiz9LREd1SKbFk0Z5HNSgNgXaWR++xiVNMCC2FCA1lgnIgxZ58='
```

Рисунок 3. Пример применения утилиты

Пример команды для расшифровки текста:

```
~/CertificateManager$ ./CertificateManager decrypt --private-key
../keytabs/private.key --text "text"
```

```
itpro_mgmt@prd-bietl-01:~/CertificateManager$ ./CertificateManager decrypt --private-key ../keytabs/private.key --text "
fePw4my/Ru8Iec0yda4SsAGU18DSgsBk9obxtVb0ASKHRAQP+oRyYAC4qf3Up7LiXbzhIEvOuZw0d0gmn/z+NLuNK1lP4mm2w19at+HugxJ7AZkGM0LDyaJ9
NB05DbaC6UnH621zrpvyTGAYRR00Q+zf0ygtPzEngKp15CvzT2b1H0ju5VCYdwFzhNVKCPngzy/wqt1NBpC6320h7UhrReXNbdUuDjyPmVJnpGwqx/k7T0
oXSAPE618VKAwaQK0lPmn5lvr7X9M+MPFqGSpP9e73I62EZBu0kXaPv+Vzt0j+i4dqbaNKJZCK7xlmF4xIFSxe0d0FbNyeNwx1HkiLF7q7p0Xy97o5lXgQJ
w3WKhSzwIY+xUNtOrStmg7sKptEfuXi+11HV9j+eg5Wn6j9crNANkPGH+9UgX1gfZnez1b9mxRVRyZeG/QOZ405gmy2AAK1h2+5Q6ycwtGv3kKVarcU6U5wf
DZFyYnsRLQ00ecYRR2BQn+1tIBtCdc4pgZQfwfW/iW/namLq5y07NykFgu0fdqDGOSEHXfRcfHc9BPg8R0EBLc5+CnthW4wA906ESKGWCdq6qizbyejAw26o
8WF/hwDJS2MyfICNwHihlRl3GFpEYOBHkRiz9LREd1SKbFk0Z5HNSgNgXaWR++xiVNMCC2FCA1lgnIgxZ58="
Расшифрованный текст: 'text to encrypt'
```

Рисунок 4. Пример применения утилиты

Данные, которые будут зашифрованы утилитой, помещаются в настроенную БД. В экстракторе они дешифруются и используются для подключения к внешним системам.

Шифруются следующая информация:

- Поля *key* и *secret* в таблице *stg.subscription* для подключения к S3-хранилищу.
- Отдельные элементы JSON в поле *credentials* таблицы *stg.source*.
 - Для **реляционных источников** (PostgreSQL, MySQL, SQL Server) атрибутами являются *ConnectionString* и *ConnectionStringSecure*. В 1 передается нешифрованная часть строки подключения, во 2 – зашифрованная.

Например, если необходимо зашифровать только пароль, выполняется следующая команда:

```
~/CertificateManager$ ./CertificateManager encrypt --public-key
../keytabs/public.crt --text "Password=iii435iaf2Ma"
```

Результат передается в элемент ***ConnectionStringSecure*** для поля *credentials*.

```
{
  "ConnectionString": "Server=192.168.72.109;Database=Tests;Use
r
  Id=itpro_admin;TrustServerCertificate=true;",
  "ConnectionStringSecure": "<Результат выполнения утилиты>"
}
```

- Для **других источников** обязательно прописывать флаг, указывающий зашифрован ли атрибут.

Пример заполнения поля *credentials* для RestAPI (можно зашифровать логин и пароль отдельно):

```
{
  "User": "user@itprocomp.ru",
```

```

    "UserEncryption":false,
    "Password":"<Результат выполнения утилиты>",
    "PasswordEncryption":true,
    "AuthType":1
} "

```

2.1.2. Генерация JSON для поля Credentials в таблице stg.source

Для подключения к источникам и указания учетных данных (логинов, паролей, строк подключения) необходимо заполнить поле Credentials таблицы source.

Для упрощения процесса заполнения данного поля необходим воспользоваться программой CredentialsToJsonSerializer, предварительно задав параметры, специфичные для источника. Программа расположена в директории /home/itpro admin /CredentialsToJsonSerializer/.

```

itpro_mgmt@prd-bietl-01:~/CredentialsToJsonSerializer$ ./CredentialsToJsonSerializer --help
CredentialsToJsonSerializer 1.0.0
Copyright (C) 2022 CredentialsToJsonSerializer

sqlserver      SqlServer
postgresql     PostgreSql
mysql          MySql
onedriveexcel  OneDriveExcel
restapi        RestApi
yandexdiskexcel YandexDiskExcel
bigquery       BigQuery
help           Display more information on a specific command.
version        Display version information.

itpro_mgmt@prd-bietl-01:~/CredentialsToJsonSerializer$

```

Рисунок 5. Пример использования программы

Данная программа преобразует текст в формат, который необходим для успешной загрузки из источника. При помощи ключа --help можно получить дополнительную информацию по любому источнику:

```

itpro_mgmt@prd-bietl-01:~/CredentialsToJsonSerializer$ ./CredentialsToJsonSerializer restapi --help
CredentialsToJsonSerializer 1.0.0
Copyright (C) 2022 CredentialsToJsonSerializer

--user          Required. Имя пользователя для доступа к ресурсу
--user-enc      (Default: false) Зашифровано ли имя пользователя
--password      Required. Пароль для доступа к ресурсу
--password-enc  (Default: false) Зашифрован ли пароль
--auth          Required. Тип аутентификации Rest
--help          Display this help screen.
--version       Display version information.

```

Рисунок 6. Пример использования программы

Пример команды для источника restapi:

```

itpro_mgmt@prd-bietl-01:~/CredentialsToJsonSerializer$ ./CredentialsToJsonSerializer restapi --user abc --password abcd --password-enc --auth 1
JSON:
{"User":"abc","UserEncryption":false,"Password":"abcd","PasswordEncryption":true,"AuthType":1}
itpro_mgmt@prd-bietl-01:~/CredentialsToJsonSerializer$

```

Рисунок 7. Пример использования программы

```

~/CredentialsToJsonSerializer$ ./CredentialsToJsonSerializer restapi --
user abc --password abcd --password-enc --auth 1

```


2.1.3. Экстрактор

При вызове блока инициации процесса (python-скрипт, описанный в начале главы) из БД settings автоматически подтягиваются все включенные запросы для всех включенных источников. На основе этого списка генерируются вызовы исполняемого файла экстрактора *GetFromSourceToParquetConsole*.

Путь к файлу – */home/itpro_admint/LoadingToS3*.

Таким образом, для пользователя нет необходимости взаимодействовать с этим компонентом, вся нагрузка лежит на блоке инициации процесса или на специализированном оркестраторе (например, MetaOrchestrator, входит в состав системы BI.Qube). В некоторых случаях может потребоваться запустить вручную данный файл, например, чтобы протестировать запрос отдельно. Аналогично предыдущим компонентам работает `--help` в командной строке.

Перечислим параметры, которые необходимо передать скрипту для взаимодействия с источниками и назначениями:

- Source – необходим для указания модулю типа источника, из которого извлекаются данные. *Поле name в таблице stg.source_type..*
- BatchSize – количество записей при пакетной загрузке данных. *Поле batch_size в таблице stg.command.*
- Command – запрос на получение данных к источнику. Конструкция запроса предполагает SQL-подобную инструкцию Select с возможностью определения полей и фильтров. *Поле command в таблице stg.command.*
- Credentials – реквизиты для подключения к источнику в формате JSON. Опционально реквизиты можно хранить в зашифрованном виде. *Поле credentials в таблице stg.source.*
- FileName – путь к файлу Parquet в S3-хранилище для записи данных из источников. *Поле sink_filename в таблице stg.command.*
- Key, Secret, Region, BucketName, Address – прочие параметры, специфичные для загрузки в S3 object storage. *Таблицы subscription и bucket.*
- Metadata-json-path – путь к файлам JSON с версиями метаданных запросов.

2.1.4. Формирование слоя хранения данных в Greenplum

Для этой задачи разработан блок генерации внешних таблиц и представлений. Внешние таблицы (external table, ET) позволяют обращаться к файлам формата Parquet как к объектам БД. Представления (view) позволяют поддержать оригинальные наименования и типы данных источников.

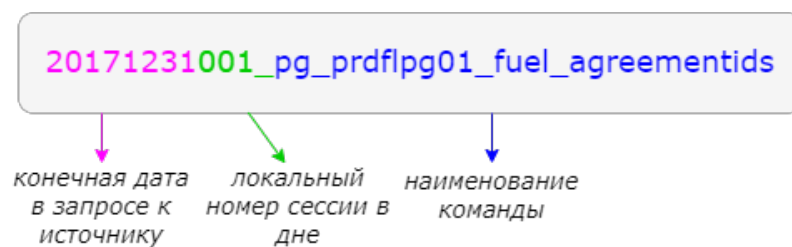
Генератор вызывается на последнем этапе пайплайна (блок инициации процесса). Сборка расположена в директории /home/itpro_admin/DbEntitiesGenerator/.

В ходе работы генератора на каждый файл в S3-совместимом хранилище создаются внешние таблицы в схеме “back”. Если таблицы уже существуют, то пересоздания не происходит.

Загрузка этих таблиц может быть 3 типов (*подробнее в главе 5*):

1. инкрементальная,
2. снимки (полная загрузка),
3. снимки (полная загрузка) с сохранением истории.

Каждый запрос из инкрементальной загрузки представлен в Greenplum перечнем ЕТ (1 загрузка в parquet = 1 ЕТ). Наименование для ЕТ состоит из следующих частей:



Поверх ЕТ формируются представления 2 типов:

1. представление, включающее ЕТ, удовлетворяющие последней версии метаданных (к названию добавляется постфикс, указывающий на номер версии “v003”);
2. общее представление, включает в себя первую версию данных. Далее дополняется вручную пользователем.

```
> bq_marketing_datalake_analytics_255573231_events
> bq_marketing_datalake_analytics_255573231_events_v001
> bq_marketing_datalake_analytics_255573231_events_v002
> bq_marketing_datalake_analytics_255573231_events_v003
```

Аналогичный сценарий с версиями используется для снимков. Файлы Parquet либо перезаписывается, либо обновляется с сохранением истории в S3

Пример SQL-запроса для генерации внешней таблицы:

```
CREATE EXTERNAL TABLE
monopolysun.public.20221104002_MS_DocumentsInsurance (
  id text,
  Sum numeric
)
LOCATION (
  'pxf://monopoly-sun-
temp/incremental/*/*/*DocumentsInsurance_*.parquet?PROFILE=s3:parquet&a
ccesskey=YCAJEcQEeSmEXJ2wUXJK_NyO&secretkey=YCO9wLy5_08C9mA3CgcV4kXnQ
rJCddd6ZLklp4&endpoint=storage.yandexcloud.net&SERVER=storage'
) ON ALL
FORMAT 'CUSTOM' (FORMATTER='pxfwritable_import')
ENCODING 'UTF8';
```

Пример SQL-запроса для генерации представлений:

```
CREATE OR REPLACE VIEW public.pg_prdflpg01_fuel_fuelsupplyschemes_v001
```

```

AS SELECT q.id,
          q.updatedat,
          q.filledtype,
          q.servicemethod
FROM
(
  SELECT
    "20171231001_pg_prdflpg01_fuel_fuelsupplieschemes"."Id"::uuid    AS
    id,
    to_timestamp(("20171231001_pg_prdflpg01_fuel_fuelsupplieschemes"."U
    pdatedAt"/ 1000000)::double precision)::timestamp without time zone
    AS updatedat,

    "20171231001_pg_prdflpg01_fuel_fuelsupplieschemes"."FilledType" AS
    filledtype,

    "20171231001_pg_prdflpg01_fuel_fuelsupplieschemes"."ServiceMethod"
    AS servicemethod

  FROM back."20171231001_pg_prdflpg01_fuel_fuelsupplieschemes"
) q
WHERE
  q.updatedat >= '1900-01-01 00:00:00'::timestamp without time zone
AND
  q.updatedat < '2017-12-31 21:00:00'::timestamp without time zone;

```

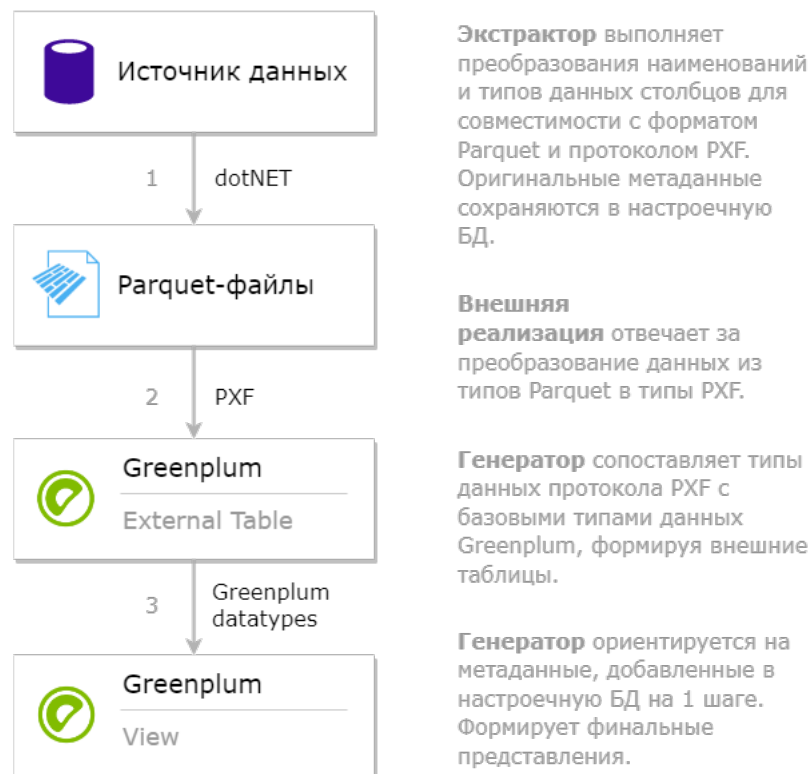


Рисунок 8. Алгоритм формирования слоя в Greenplum

Условие “WHERE” в данном случае помогает оптимизатору Greenplum не сканировать все ЕТ, когда нужно взять из одной конкретной ЕТ. Чтобы данное условие было добавлено необходимо заполнить partition_column.

На рисунке ниже представлен полный путь перекладки данных из источников в назначение.

Внешние таблицы ссылаются на файлы с помощью протокола PXF, рекомендуемого для чтения из S3-хранилища в документации Yandex-Cloud. Список поддерживаемых типов данных PXF и сопоставление с Greenplum можно посмотреть здесь:

[Reading and Writing HDFS Parquet Data | Pivotal Greenplum Docs](#)

2.2. Структура проекта MetaStaging в файловой системе

Путь к файлам MetaStaging на Windows и на Linux отличается.

Рисунок 9. Путь в файловой системе Windows

Файлы MetaStaging на Linux находятся в папке home пользователя itpro_admin.

В общем виде структура проекта выглядит следующим образом:

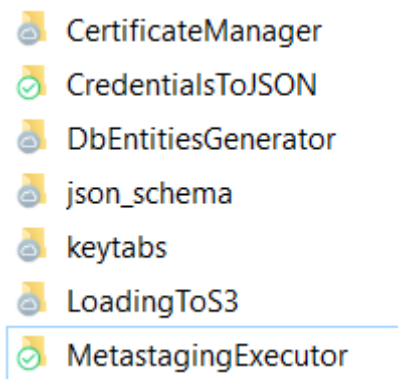


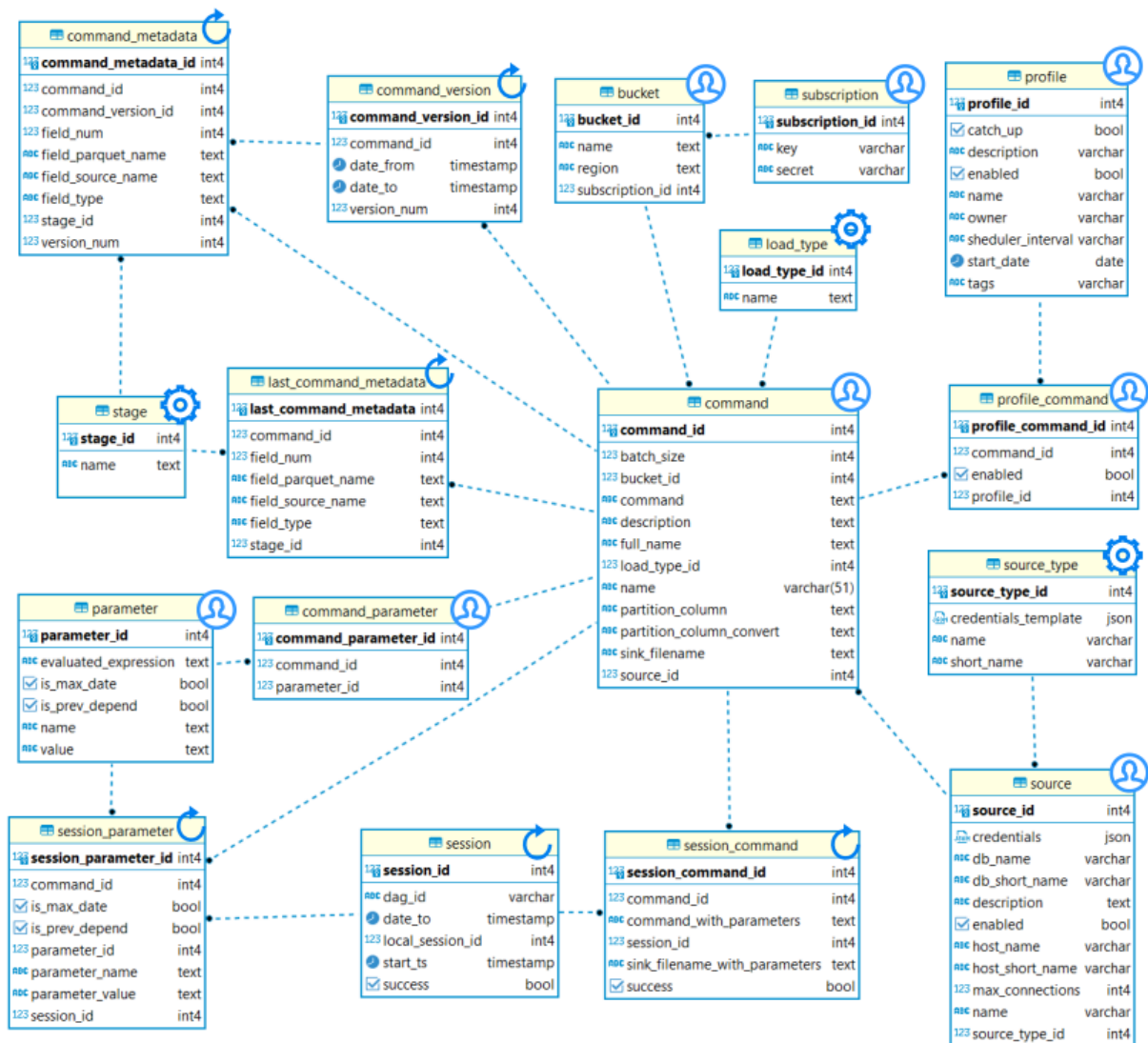
Рисунок 1. Файловая система компонента

- CertificateManager – это «Утилита шифрования», описанная в предыдущем разделе.
- CredentialsToJSON – содержит сборку утилиты «Генерация JSON для поля Credentials».
- DbEntitiesGenerator – это сборка для «Формирования слоя хранения данных на Greenplum».
- Json_schema – данный каталог содержит JSON-файлы с метаданными запросов из источника, которые выполнялись в рамках конкретных сессий. Эти данные записываются в БД в таблицу *stg.last_command_metadata*.
- Keytabs – данный каталог содержит ключи для шифрования и расшифрования конфиденциальной информации. Каталог используется экстрактором (в файле app.config проекта указан путь к каталогу) и утилитой CertificateManager.
- LoadingToS3 – содержит сборку с экстрактором.
- MetastagingExecutor – содержит python-скрипт инициации процесса.

2.3. Структура Базы данных

Таблицы MetaStaging делятся на три категории:

- *Настроечные таблицы* – данные вносит пользователь в соответствии с правилами заполнения.
- *Таблицы логов* (заполняются автоматически) – данные вносит система по итогам выполнения очередной сессии.
- *Таблицы дескрипторы* (заполнены разработчиком) – данные внесены предварительно в соответствии с реализованным функционалом.



Настроечные таблицы

Таблицы логов

Таблицы дескрипторы

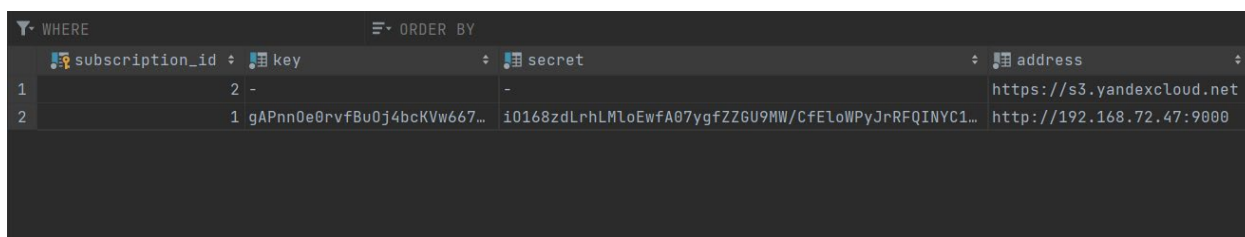
Рисунок 2. Структура таблиц MetaStaging

2.3.1. Настраиваемые таблицы MetaStaging

1. «**stg.subscription**» – список ключей необходимых для доступа к сервисам S3. Можно записывать как в зашифрованном, так и в незашифрованном виде.

Поля таблицы:

Имя столбца	Тип данных	Источник	Назначение
subscription_id	int	Автоинкрементен	Идентификатор ключа для доступа к сервисам
key	varchar	Вручную	Ключ
secret	varchar	Вручную	Секретный ключ
address	text	Вручную	Адрес хоста с установленным S3 MinIO



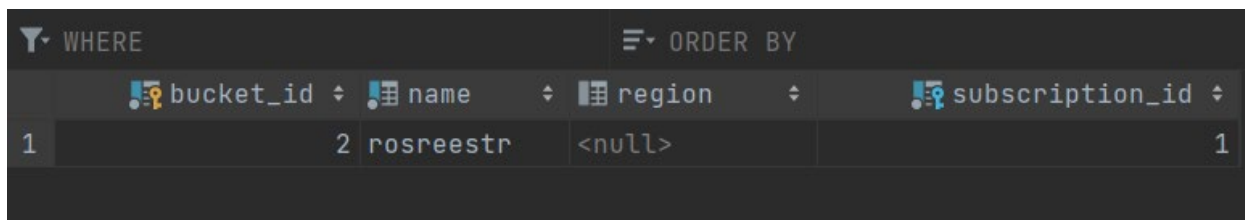
	subscription_id	key	secret	address
1	2	-	-	https://s3.yandexcloud.net
2	1	gAPnn0e0rvfBu0j4bcKVw667...	i0168zdLrhLMloEwfA07ygfZZGU9MW/CfEl0WPYJrRFQINYC1...	http://192.168.72.47:9000

Рисунок 3. Таблица «stg.subscription»

2. «**stg.bucket**» – список каталогов S3-хранилища, доступных для размещения файлов parquet.

Поля таблицы:

Имя столбца	Тип данных	Источник	Назначение
bucket_id	int	Автоинкремент	Идентификатор каталога
name	text	Вручную	Имя каталога
region	text	Вручную	Регион, в котором создаются бакеты (для MinIO не указывается)
subscription_id	int	Вручную	Идентификатор ключа для доступа к сервисам



	bucket_id	name	region	subscription_id
1	2	rosreestr	<null>	1

Рисунок 4. Таблица «stg.bucket»

3. «**stg.source**» – информация, необходимая системе для доступа к источникам данных.

Поля таблицы:

Имя столбца	Тип данных	Источник	Назначение
source_id	int	Автоинкремент	Идентификатор источника
name	varchar	Вручную	Название источника, не влияет на загрузку данных
description	text	Вручную	Описание источника, необязательное поле
credentials	json	Вручную	Учетные данные для подключения определенного источника в формате JSON
enabled	bool	Вручную	Метка о необходимости выполнения/игнорирования всех запросов к данному источнику
max_connections	int	Вручную	Максимальное количество одновременных запросов на выполнение, разрешенных для источника
source_type_id	int	Вручную	Ссылка на тип источника

source_id	name	description	source_type_id	credentials	enabled	max_connections
2	Arenadata-pg-kinopoisk		7	{\"ConnectionString\": \"...\"}	true	8
1	test_postgre	Тест загрузки данны...	7	{\"ConnectionString\": \"...\"}	true	8
3	MySQL-view	Компиляция отзывов...	4	{\"ConnectionString\": \"...\"}	true	8
4	SQL Server	<null>	6	{\"ConnectionString\": \"...\"}	true	8

Рисунок 5. Таблица «stg.source»

4. «stg.command» – список запросов к источникам.

Поля таблицы:

Имя столбца	Тип данных	Источник	Назначение
command_id	int	Автоинкремент	Идентификатор команды
batch_size	int	Вручную	Количество строк выгружаемых из источников в файл Parquet за одну итерацию при пакетной загрузке
command	text	Вручную	Текст запроса с возможностью параметризации. Используется Select-подобный синтаксис, при работе экстрактора запрос корректируется в соответствии с требованиями источника
description	text	Вручную	Описание команды
full_name	text	Вручную	Полное наименование команды
name	varchar	Вручную	Наименование команды, длина до 51 символа! Влияет на наименование файла Parquet
sink_filename	text	Вручную	Путь к Parquet-файлу в S3-совместимом хранилище
source_id	int	Вручную	Ссылка на источник
partition_column	text	Вручную	Поле в запросе, по которому выполняется секционирование на представлениях Greenplum. Если данное поле задано (например, <i>UpdatedAt</i>), в запросе можно писать так <code>< /*{partition_column}*/ >= /*{datefrom}*/ ></code>
partition_column_convert	text	Вручную	Поле содержит логику конвертации для значения в partition_column. Данная логика будет отражена в представлении на Greenplum. <i>Пример: cast(/*{partition_column}*/ as bigint)</i>

command_id	source_id	name	description	command	sink_filename	batch_size	bucket_id	load_type_id
1	2	1 Футболист...	Перечень футболи...	select id, n...	football	10000	2	2
2	3	2 Kinopoisk...	Русскоязычные От...	select * fro...	kinopoisk_rus_revi...	1000	2	2
3	4	3 Multilang...	Представление - ...	select * fro...	multilanguage_revi...	100000	2	2
4	5	4 Dialogs-e...	<null>	select * fro...	dialogs_expanded	10000	2	2

Рисунок 6. Таблица «stg.command»

5. «stg.profile» – список профилей загрузки. На основе данной таблицы (совместно с profile_command) запросы к источникам будет разделены на несколько DAG’ов в оркестраторе.

Поля таблицы:

Имя столбца	Тип данных	Источник	Назначение
profile_id	int	Автоинкремент	Идентификатор профиля
name	varchar	Вручную	Название профиля
description	varchar	Вручную	Описание профиля
enabled	bool	Вручную	Метка о необходимости выполнения/игнорирования данного профиля
owner	varchar	Вручную	Владелец графа выполнения (<i>используется только при наличии оркестратора</i>)
scheduler_interval	varchar	Вручную	Интервал запуска DAG’ов, может принимать в себя явное стоп-выражение, алиас (@daily, @monthly и тд) или спец-слово None, что приведет к отключению автоматического планирования и запуска задач по расписанию (<i>используется только при наличии оркестратора</i>)
start_date	date	Вручную	Дата начала профиля

profile_id	name	description	owner	start_date	enabled	scheduler_interval	tags	catch_up
1	test_rosreestr	<null>	Дьячков Ники...	2022-11-20	true	None	<null>	<null>

Рисунок 16. Таблица «stg.profile»

6. «profile_command» – соответствие выполняемых команд профилям загрузки. Одна команда может использоваться несколькими профилями, один профиль может выполнять несколько команд.

Поля таблицы:

Имя столбца	Тип данных	Источник	Назначение
profile_command_id	int	Автоинкремент	Идентификатор профиля, команды
profile_id	int	Вручную	Идентификатор профиля
command_id	int	Вручную	Идентификатор команды
enabled	bool	Вручную	Метка о необходимости выполнения/игнорирования данного запроса

	profile_command_id	profile_id	command_id	enabled
1	1	1	2	false
2	2	1	3	false
3	4	1	5	true
4	3	1	4	false

Рисунок 17. Таблица «profile_command»

7. **«parameter»** – список параметров для использования в запросах к источникам.

Поля таблицы:

Имя столбца	Тип данных	Источник	Назначение
parameter_id	int	Автоинкремент	Идентификатор параметра
name	text	Вручную	Наименование параметра, указывается в запросе
value	text	Вручную	Значение параметра, подставляется в запрос при работе экстрактора
evaluated_expression	text	Вручную	SQL-выражение, которое преобразует value к требуемому формату.
is_max_date	bool	Вручную	Указание экстрактору, ссылаться ли на value в запросе или не указывать в нем верхнюю границу (брать актуальные данные, начиная от конкретной даты)
is_prev_depend	text	Вручную	Указание экстрактору, ссылаться ли на value или использовать дату, зависящую от последней сессии

	parameter_id	name	value	is_prev_depend	is_max_date	evaluated_expression
--	--------------	------	-------	----------------	-------------	----------------------

Рисунок 18. Таблица «parameter»

8. **command_parameter** – ссылки на таблицы command и parameter.

Поля таблицы:

Имя столбца	Тип данных	Источник	Назначение
command_parameter_id	int	Автоинкремент	Идентификатор параметра, команды
command_id	int	Вручную	Идентификатор команды
parameter_id	int	Вручную	Идентификатор параметра

WHERE	ORDER BY
command_parameter_id	command_id
parameter_id	
1	3
2	4

Рисунок 19 Таблица «command_parameter»

2.3.2. Таблицы логов MetaStaging

Данные таблицы заполняются в процессе работы компонента

1. «session» – список запусков пайплайна.

Поля таблицы:

Имя столбца	Тип данных	Источник	Назначение
sessionID	int	Автоинкремент	Идентификатор сессии
start ts	timestamp	Автоматически	Дата и время запуска работы MetaStaging
date to	timestamp	Автоматически	Дата окончания сессии
success	bool	Автоматически	Идентификатор успешности загрузки в рамках данной сессии
dag_id	varchar	Автоматически	Наименования графа загрузки (<i>используется только при наличии оркестратора</i>)
local_session_id	int	Автоматически	Локальный идентификатор сессии относительно дня (когда запрос выполняется несколько раз за день)

WHERE	ORDER BY
session_id	start_ts
date_to	success
dag_id	local_session_id
1	2
2	3
3	4
4	5
5	6
6	7
7	8
8	9
9	10
10	11
11	12
12	13
13	14
14	15
15	16
16	17
17	18
18	19

Рисунок 7. Таблица «session»

2. «session_command» – список запросов для извлечения данных из источников в S3-хранилище. Запросы разделены на секции, в рамках которых выполнялись.

Поля таблицы:

Имя столбца	Тип данных	Источник	Назначение
session_command ID	int	Автоинкремент	Идентификатор строки
command_with_parameters	text	Автоматически	Окончательный текст запроса, отправляемого на источник данных
success	bool	Автоматически	Метка об успешном выполнении запроса
sink_filename_with_parameters	text	Автоматически	Полный путь к файлу parquet в рамках S3-совместимого хранилища.
session ID	int	Автоматически	Идентификатор сессии
command ID	int	Автоматически	Идентификатор команды

session_id	command_id	command_with_parameters	success	sink_filename_with_parameters
1	1	2 SELECT id, "Name", age, nationalit..	false	snapshot/football.parquet
2	2	2 SELECT id, "Name", age, nationalit..	true	snapshot/football.parquet
3	3	2 SELECT id, "Name", age, nationalit..	false	snapshot/football.parquet
4	4	2 SELECT id, "Name", age, nationalit..	false	snapshot/football.parquet
5	5	2 SELECT id, "Name", age, nationalit..	false	snapshot/football.parquet
6	6	2 SELECT id, "Name", age, nationalit..	false	snapshot/football.parquet
7	7	2 SELECT id, "Name", age, nationalit..	false	snapshot/football.parquet
8	8	2 SELECT id, "Name", age, nationalit..	false	snapshot/football.parquet
9	9	2 SELECT id, "Name", age, nationalit..	false	snapshot/football.parquet
10	10	2 SELECT id, "Name", age, nationalit..	false	snapshot/football.parquet
11	11	2 SELECT id, "Name", age, nationalit..	false	snapshot/football.parquet
12	12	2 select id, "Name", age, nationalit..	false	snapshot/football.parquet
13	13	2 select id, name, age, nationality,..	true	snapshot/football.parquet
14	14	2 select id, name, age, nationality,..	true	snapshot/football.parquet
15	15	2 select id, name, age, nationality,..	true	snapshot/football.parquet
16	16	2 select id, name, age, nationality,..	true	snapshot/football.parquet
17	17	2 select id, name, age, nationality,..	true	snapshot/football.parquet
18	18	2 select id, name, age, nationality,..	true	snapshot/football.parquet

Рисунок 8. Таблица «session_command»

3. «command_version» – список существующих версий для конкретных команд с указанием периода их действия.

Поля таблицы:

Имя столбца	Тип данных	Источник	Назначение
command_version ID	int	Автоинкремент	Идентификатор версии
date_from	timestep	Автоматически	Фактическая дата начала загрузки
date_to	timestep	Автоматически	Фактическая дата окончания загрузки
version num	int	Автоматически	Номер версии запроса к источнику
command_id	int	Автоматически	Ссылка на команду

command_version_id	date_from	date_to	version_num	command_id
1	2022-11-24 19:18:37.144075	9999-01-01 00:00:00.000000	1	2
2	2022-12-01 09:00:49.148210	9999-01-01 00:00:00.000000	1	3
3	2022-12-01 16:03:02.410113	9999-01-01 00:00:00.000000	1	4

Рисунок 9. Таблица «command_version»

4. «**command_metadata**» – список наименований и типов данных полей, включенных в запросы к источникам за всю историю выполнений этих запросов.

Поля таблицы:

Имя столбца	Тип данных	Источник	Назначение
command_metadata ID	int	Автоинкремент	Идентификатор метаданных команды
command_version ID	int	Автоматически	Идентификатор команды, версии
field_source_name	text	Автоматически	Наименование поля в источнике
field_type	text	Автоматически	Тип данных поля на различных этапах интеграции
field_parquet_name	text	Автоматически	Наименование поля, допустимое для записи в Parquet и чтения в Greenplum через External Table (замена пробелов и специальных символов)
version_num	int	Автоматически	Версия метаданных запроса. По умолчанию равна 1. При изменении метаданных значение увеличивается
field_num	int	Автоматически	Номер поля в запросе к источнику, используется для сохранения порядка полей в Greenplum относительно порядка в источнике

	command_metadata_id	command_version_id	field_source_name	field_parquet_name	field_type	stage_id
1	1	1	id	id	integer	
2		2	id	id	System.Int32	
3		3	id	id	Int(bitWidth=32, isSigned=true)	
4		4	id	id	Int32	
5		5	name	name	character varying	
6		6	name	name	System.String	
7		7	name	name	String	
8		8	name	name	ByteArray	
9		9	age	age	integer	
10		10	age	age	System.Int32	
11		11	age	age	Int(bitWidth=32, isSigned=true)	
12		12	age	age	Int32	
13		13	nationality	nationality	character varying	
14		14	nationality	nationality	System.String	
15		15	nationality	nationality	String	
16		16	nationality	nationality	ByteArray	
17		17	club	club	character varying	
18		18	club	club	System.String	

Рисунок 10. Таблица «command_metadata»

5. «**last_command_metadata**» – список наименований и типов данных полей при последнем выполнении конкретного запроса. Здесь всегда хранится актуальная версия метаданных. При формировании представлений на Greenplum данные из этой таблицы сравниваются с данными из таблицы *command_metadata* для проверки изменений (добавлено поле, изменен тип и тд).

Поля таблицы:

Имя столбца	Тип данных	Источник	Назначение
last_command_metadata_id	int	Автоинкремент	Идентификатор строки
field_source_name	text	Автоматически	Наименование поля в источнике
field_type	text	Автоматически	Тип данных поля на этапах интеграции
field_parquet_name	text	Автоматически	Наименование поля, допустимое для записи в Parquet и чтения в Greenplum через External Table (замена пробелов и специальных символов)
field_num	int	Автоматически	Номер поля в запросе к источнику, используется для сохранения порядка полей в Greenplum относительно порядка в источнике

	last_command_metadata	field_source_name	field_parquet_name	field_type	stage_id	command_id
1	509	sentence_uuid	sentence_uuid	VARCHAR	1	4
2	510	sentence_uuid	sentence_uuid	System.String	2	4
3	511	sentence_uuid	sentence_uuid	String	3	4
4	512	sentence_uuid	sentence_uuid	ByteArray	4	4
5	513	sentence_text	sentence_text	VARCHAR	1	4
6	514	sentence_text	sentence_text	System.String	2	4
7	515	sentence_text	sentence_text	String	3	4
8	516	sentence_text	sentence_text	ByteArray	4	4
9	517	language	language	VARCHAR	1	4
10	518	language	language	System.String	2	4
11	519	language	language	String	3	4
12	520	language	language	ByteArray	4	4
13	521	domain	domain	VARCHAR	1	4
14	522	domain	domain	System.String	2	4
15	523	domain	domain	String	3	4
16	524	domain	domain	ByteArray	4	4
17	525	labelled	labelled	VARCHAR	1	4
18	526	labelled	labelled	System.String	2	4

Рисунок 11. Таблица «last_command_metadata»

6. «session_parameter» – информация о параметрах запросов, выполняемых в рамках сессии.

Поля таблицы:

Имя столбца	Тип данных	Источник	Назначение
session_parameter_id	int	Автоинкремент	Идентификатор строки
session_id	int	Автоматически	Ссылка на сессию
Parameter_id	int	Автоматически	Ссылка на параметр
Parameter_value	text	Автоматически	Значение параметра
Is_prev_depend	bool	Автоматически	Поле, аналогичное полю из <i>stg.parameter</i>
Is_max_date	bool	Автоматически	Поле, аналогичное полю из <i>stg.parameter</i>
Command_id	int	Автоматически	Ссылка на команду
Parameter_name	text	Автоматически	Наименование параметра из таблицы <i>stg.parameter</i>

	session_parameter_id	session_id	parameter_id	parameter_value	is_prev_depend	is_max_date	command_id	parameter_name
1	1	32	1	2020	false	false		7 year_from
2	2	32	2	2022	false	false		7 year_to
3	3	33	1	2020	false	false		7 year_from
4	4	33	2	2022	false	false		7 year_to
5	5	34	1	2020	false	false		7 year_from
6	6	34	2	2022	false	false		7 year_to
7	7	35	1	'1900-01-01 00:00:00'	true	false		8 date_from
8	8	35	2	'2022-12-05 08:51:23'	false	true		8 date_to
9	9	36	1	'2022-12-05 08:51:23'	true	false		8 date_from
10	10	36	2	'2022-12-05 14:45:15'	false	true		8 date_to
11	11	37	1	'2022-12-05 08:51:23'	true	false		8 date_from
12	12	37	2	'2022-12-06 14:11:21'	false	true		8 date_to

Рисунок 12. Таблица «session_parameter»

2.3.3. Таблицы дескрипторы

Данные в этих таблицах описывают функциональные части компонента. Они заполняются автоматически при развертывании MetaStaging. И используются только в качестве ссылок при заполнении настроечных таблиц.

1. «stage» – список этапов интеграции данных.

Поля таблицы:

Имя столбца	Тип данных	Источник	Назначение
Stage_ID	int	Автоинкремент	Идентификатор строки
name	text	Автоматически	Наименование этапа интеграции <i>source</i> – метаданные системы-источника, <i>logical</i> – логический тип хранения в Parquet, <i>physical</i> – физический тип хранения Parquet, <i>dotnet</i> – тип данных «.Net Core» используемый в экстракторе для извлечения из источника

	stage_id	name
1	1	source
2	2	logical
3	3	physical
4	4	dotnet

Рисунок 26. Таблица «stage»

2. «load_type» – список режимов загрузки данных из источников.

Поля таблицы:

Имя столбца	Тип данных	Источник	Назначение
load_type_id	int	Автоинкремент	Идентификатор строки
name	text	Автоматически	Наименование типа загрузки

	load_type_id	name
1	1	Секционированная загрузка
2	2	Полная загрузка
3	3	Полная загрузка с сохранением истории

Рисунок 27. Таблица «load_type»

3. «source_type» – список поддерживаемых источников.

Поля таблицы:

Имя столбца	Тип данных	Источник	Назначение
source_type_id	int	Автоинкремент	Идентификатор строки
name	varchar	Автоматически	Наименование источника
credentials_template	json	Автоматически	Справочная информация – JSON шаблон заполнения поля credentials в <i>stg.source</i>
short name	varchar	Автоматически	Краткое наименование

	source_type_id	name	credentials_template	short_name
1	1	bigquery	{ "type": "<service_account>", "type_encryption": bq	
2	2	onedriveexcel	{ "Email": "<sample@gmail.com>", "EmailEncryption": od	
3	3	yandexdiskexcel	{ "Token": "<token> Get Yandex token URL: https://oau yd	
4	4	mysql	{"ConnectionString": "server=<server_name>;uid=<user_nam my	
5	5	restapi	{ "User": "<User>", "UserEncryption": true, ra	
6	6	sqlserver	{ "ConnectionString": "Data Source=<servername>;Tru ms	
7	7	postgresql	{ "ConnectionString": "Host=<localhost>;Port=<5432> pg	

Рисунок 28. Таблица «source_type»

2.3.4. Хранимые процедуры и функции MetaStaging

Представленные далее процедуры и функции являются сервисными, то есть, вызов пользователями запрещен.

stg.init_session – инициализирует новую сессию загрузки данных, и выполняет подстановку параметров в запросы к источникам, заполняя таблицы логов (*stg.session_command*).

Параметры:

- *profile_id* – идентификатор профиля загрузки.
- *start_ts_var* – дата запуска сессии.
- *dag_id* – название профиля (при отсутствии оркестратора, можно передавать любое текстовое значение).

Возвращает:

- JSON с данными, необходимыми для запуска экстрактора (см. раздел 4.1.3.), включая *id* новой сессии и список запросов к источникам (*id* команды, текст команды, путь к файлу в S3-хранилище).

stg.get_sample_commands

Параметры:

- *profile_id* – идентификатор профиля загрузки.

Возвращает:

- JSON с данными, необходимыми для создания каталога в файловой системе «*json_schema*» (п. Генерация JSON для поля *Credentials* в таблице *stg.source*) и для запуска экстрактора (п. Экстрактор).

stg.compare_metadata – вызывается для конкретного запроса к источнику с целью проверки изменения метаданных текущей таблицы с таблицей, загруженной в прошлой сессии. Заполняет таблицы *stg.command_version* и *stg.command_metadata*.

Параметры:

- *command_id* – идентификатор запроса к источнику.

Возвращает:

- Номер актуальной версии метаданных.

stg.mark_command_success – Помечает запросы, которые были выполнены корректно (включая создание представлений на Greenplum), как успешные: *поле success в stg.session_command = True*.

Параметры:

- *command_id*

stg.mark_session_success - Помечает сессии, которые были выполнены корректно как успешные: *поле success в stg.session = True*.

Параметры:

- *command_id*