

流计算 Oceanus

ETL 开发指南



腾讯云

【 版权声明 】

©2013–2024 腾讯云版权所有

本文档（含所有文字、数据、图片等内容）完整的著作权归腾讯云计算（北京）有限责任公司单独所有，未经腾讯云事先明确书面许可，任何主体不得以任何形式复制、修改、使用、抄袭、传播本文档全部或部分内容。前述行为构成对腾讯云著作权的侵犯，腾讯云将依法采取措施追究法律责任。

【 商标声明 】

及其它腾讯云服务相关的商标均为腾讯云计算（北京）有限责任公司及其关联公司所有。本文档涉及的第三方主体的商标，依法由权利人所有。未经腾讯云及有关权利人书面许可，任何主体不得以任何方式对前述商标进行使用、复制、修改、传播、抄录等行为，否则将构成对腾讯云及有关权利人商标权的侵犯，腾讯云将依法采取措施追究法律责任。

【 服务声明 】

本文档意在向您介绍腾讯云全部或部分产品、服务的当时的相关概况，部分产品、服务的内容可能不时有所调整。您所购买的腾讯云产品、服务的种类、服务标准等应由您与腾讯云之间的商业合同约定，除非双方另有约定，否则，腾讯云对本文档内容不做任何明示或默示的承诺或保证。

【 联系我们 】

我们致力于为您提供个性化的售前购买咨询服务，及相应的技术售后服务，任何问题请联系 4009100100或 95716。

文档目录

ETL 开发指南

概述

ETL 作业术语表

上下游开发指南

数据源表 MySQL

数据目的表 MySQL

数据目的表 PostgreSQL

数据目的表 ClickHouse

数据目的表 Elasticsearch

ETL 开发指南

概述

最近更新时间：2023-06-21 15:21:59

ETL 是将业务系统的数据经过抽取、清洗转换之后加载到数据仓库的过程，目的是将企业中的分散、零乱、标准不统一的数据整合到一起，为企业的决策提供分析依据。ETL 管道任务可以从数据源表获取数据，进行一些转换操作或信息补充，再将结果加载到目的源表中。开发人员甚至不需要了解编程语言，只需要选择数据源表和目的表，根据业务逻辑完成字段映射的配置，即可启动 ETL 作业。

本章节主要提供在独享集群上开发 ETL 作业的指南。通过阅读本章节，您将学习到以下内容：

- [ETL 作业术语表](#)
- [数据源表 MySQL 开发](#)
- [数据目的表 MySQL 开发](#)
- [数据目的表 PostgreSQL 开发](#)
- [数据目的表 ClickHouse 开发](#)
- [数据目的表 Elasticsearch 开发](#)

ETL 作业术语表

最近更新时间：2023-06-21 15:22:00

ETL 作业常用术语如下：

术语	详细说明
流计算	流计算是面向流式数据的计算，它从一个或多个流式数据源读取持续不断产生的数据，经过引擎中多个算子的组合进行高效计算，再根据实际需要，将结果输出至下游的多种数据目的，例如消息队列、数据库、数据仓库、存储服务。
数据源表 (Source)	为流计算系统持续提供输入数据。
数据目的表 (Sink)	流计算系统处理结果输出的地方。
Schema	表示一个表的结构信息，例如各个列名、列类型等。对于 PostgreSQL 而言，Schema 是介于 Database 和 Table 之间的一个层级，可以理解成数据库内部的命名空间。
MySQL	一种常用数据库，在 ETL 作业中可用作数据源表与数据目的表。
PostgreSQL	类似 MySQL 的关系型数据库。
ClickHouse	ClickHouse 是一个用于联机分析（OLAP）的列式数据库管理系统（DBMS），在 ETL 作业中可用作数据目的表。
Elasticsearch	实时的搜索与数据分析引擎。
字段映射	字段映射实现了从数据源表中抽取数据，对数据进行计算、清洗，再把数据加载到目的表中。
常量字段	可以输入一个自定义常量字段到目的源表相应的字段中。
计算字段	可以对从数据源表抽取出来的字段数据进行 内置函数 数值转换或者计算。

上下游开发指南

数据源表 MySQL

最近更新时间：2023-06-21 15:22:00

介绍

MySQL 数据源表支持对 MySQL 数据库的全量和增量读取，并保证 Exactly Once 语义。MySQL 数据源表底层使用 Debezium 来做 CDC（Change Data Capture）。其工作机制如下：

1. 获取一个全局读锁，从而阻塞住其他数据库客户端的写操作。
2. 开启一个可重复读语义的事务，来保证后续在同一个事务内读操作都是在一个一致性快照中完成的。
3. 读取 Binlog 的当前位置。
4. 读取连接器中配置的数据库和表的模式（schema）信息。
5. 释放全局读锁，允许其他的数据库客户端对数据库进行写操作。
6. 扫描全表，当全表数据读取完后，会从第3步中得到的 Binlog 位置获取增量的变更记录。

Flink 作业运行期间会周期性执行快照，记录下 Binlog 位置，当作业崩溃恢复时，便会从之前记录的 Binlog 点继续处理，从而保证 Exactly Once 语义。

类型映射

MySQL 字段类型	Flink 字段类型
TINYINT	TINYINT
SMALLINT	SMALLINT
TINYINT UNSIGNED	
INT	INT
MEDIUMINT	
SMALLINT UNSIGNED	
BIGINT	BIGINT
INT UNSIGNED	
BIGINT UNSIGNED	DECIMAL(20, 0)
FLOAT	FLOAT
DOUBLE	DOUBLE

DOUBLE PRECISION	
NUMERIC(p, s)	DECIMAL(p, s)
DECIMAL(p, s)	
BOOLEAN	BOOLEAN
TINYINT(1)	
DATE	DATE
TIME [(p)]	TIME [(p)] [WITHOUT TIMEZONE]
DATETIME [(p)]	TIMESTAMP [(p)] [WITHOUT TIMEZONE]
TIMESTAMP [(p)]	TIMESTAMP [(p)]
TIMESTAMP [(p)] WITH LOCAL TIME ZONE	
CHAR(n)	STRING
VARCHAR(n)	
TEXT	
BINARY	BYTES
VARBINARY	
BLOB	

注意事项

用户权限

用于同步的源数据库的用户必须拥有以下权限 SHOW DATABASES、REPLICATION SLAVE、REPLICATION CLIENT、SELECT 和 RELOAD。

数据库参数设置

binlog_row_image 参数的参数运行值应当设置为 FULL。

WITH 参数

MySQL 数据源表基于数据库 MySQL CDC 开发，两者具有相同的 WITH 参数，具体参数配置方式可参见 [数据库 MySQL CDC](#)。

数据目的表 MySQL

最近更新时间：2023-06-21 15:22:00

介绍

MySQL 数据目的表支持将数据写入到 MySQL 数据库中。

注意事项

主键说明

由于 ETL 数据源表产生的数据都为 Upsert 数据，因此 MySQL 数据库的表必须正确定义主键。

WITH 参数

MySQL 数据目的表基于 [JDBC](#) 开发，可以使用其中用于目的表的相关配置项：

参数值	必填	默认值	描述
sink.buffer-flush.max-rows	否	100	批量输出时，缓存中最多缓存多少数据。如果设置为0，表示禁止输出缓存。
sink.buffer-flush.interval	否	1s	批量输出时，每批次最大的间隔（毫秒）。如果 <code>'sink.buffer-flush.max-rows'</code> 设为 <code>'0'</code> ，而这个选项不为零，则说明启用纯异步输出功能，即数据输出到算子、从算子最终写入数据库这两部分线程完全解耦。
sink.max-retries	否	3	数据库写入失败时，最多重试的次数。

数据目的表 PostgreSQL

最近更新时间：2023-06-21 15:22:01

介绍

PostgreSQL 数据目的表支持将数据写入到 PostgreSQL 数据库中。

注意事项

主键说明

- 由于 ETL 数据源表产生的数据都为 Upsert 数据，因此 PostgreSQL 数据目的表**必须**定义主键。
- 数据目的表定义的主键**必须**为物理表中定义的主键，否则任务启动后会出错。

WITH 参数

PostgreSQL 数据目的表基于 [JDBC](#) 开发，可以使用其中用于目的表的相关配置项：

参数值	必填	默认值	描述
sink.buffer-flush.max-rows	否	100	批量输出时，缓存中最多缓存多少数据。如果设置为0，表示禁止输出缓存。
sink.buffer-flush.interval	否	1s	批量输出时，每批次最大的间隔（毫秒）。如果 <code>'sink.buffer-flush.max-rows'</code> 设为 <code>'0'</code> ，而这个选项不为零，则说明启用纯异步输出功能，即数据输出到算子、从算子最终写入数据库这两部分线程完全解耦。
sink.max-retries	否	3	数据库写入失败时，最多重试的次数。

数据目的表 ClickHouse

最近更新时间：2023-06-21 15:22:01

介绍

ClickHouse 数据目的表支持将数据写入到 ClickHouse。

⚠ 注意

ClickHouse 数据目的表的表引擎必须使用 CollapsingMergeTree。

常见数据类型映射

关于 ClickHouse 支持的数据类型定义及其使用，可参考 [ClickHouse data-types](#)，这里列举了常用的数据类型，及其与 Flink 类型的对应关系。

Flink 数据类型	ClickHouse 对应数据类型
VARCHAR	String/FixedString(N)
STRING	String/FixedString(N)
BOOLEAN	没有单独类型存储，可以使用 UInt8 来存储布尔类型，将取值限制为0或1；或者使用字符串存储 true/false 来表示
DECIMAL	Decimal32(S)/Decimal64(S)/Decimal128(S)
TINYINT	Int8
SMALLINT	Int16
INTEGER	Int32
BIGINT	Int64
FLOAT	Float32
DOUBLE	Float64
DATE	Date
TIMESTAMP	DateTime
TIMESTAMP WITH LOCAL TIME ZONE	DateTime, 示例 DateTime64(3, 'Asia/Shanghai')

注意事项

主键说明

使用 ClickHouse 数据目的表时，需要按照建表语句正确的定义主键，否则有可能无法正确同步修改与删除操作。

折叠字段

ClickHouse 的 CollapsingMergeTree 引擎在合并算法中添加了折叠行的逻辑。折叠字段在使用 CollapsingMergeTree 引擎建表时所指定：`ENGINE = CollapsingMergeTree(Sign)`。对 ClickHouse 折叠详细说明可参考 [ClickHouse 官方文档](#)。

WITH 参数

ClickHouse 数据目的表基于数据仓库 ClickHouse 开发，两者具有相同的 WITH 参数，具体参数含义用法可参考 [数据仓库 ClickHouse](#)。

数据目的表 Elasticsearch

最近更新时间：2023-06-21 15:22:01

介绍

Elasticsearch 数据目的表支持将数据写入到 Elasticsearch 中。

⚠ 注意

Elasticsearch 数据目的表暂时只支持 Elasticsearch 6 或 Elasticsearch 7 版本。

常见数据类型映射

关于 Elasticsearch 支持的数据类型定义及其使用，可参考 [Elasticsearch data-types](#)，这里列举了常用的数据类型，及其与 Flink 类型的对应关系。

Flink 数据类型	Elasticsearch 对应数据类型
text	STRING
match_only_text	STRING
binary	STRING
keyword	STRING
wildcard	STRING
search_as_you_type	STRING
ip	STRING
short	SMALLINT
integer	INT
long	BIGINT
unsigned_long	BIGINT
float	FLOAT
half_float	FLOAT
double	DOUBLE

boolean	BOOLEAN
date	TIMESTAMP(3)
date_nanos	TIMESTAMP(6)

⚠ 注意

暂时不支持上述表格没有提到的类型。

注意事项

主键说明

Elasticsearch 必须设置主键，设置为主键的字段会被写入到 `_id` 字段中，相同 ID 的数据会进行覆盖。

版本差异

Elasticsearch 6 版本与 Elasticsearch 7 版本在配置上有一些不同，Elasticsearch 6 版本需要配置 `document-type` 而 Elasticsearch 7 版本不需要。

WITH 参数

Elasticsearch 数据目的表基于数据分析引擎 Elasticsearch 开发，两者具有相同的 WITH 参数，具体参数含义用法参见 [数据分析引擎 Elasticsearch](#)。