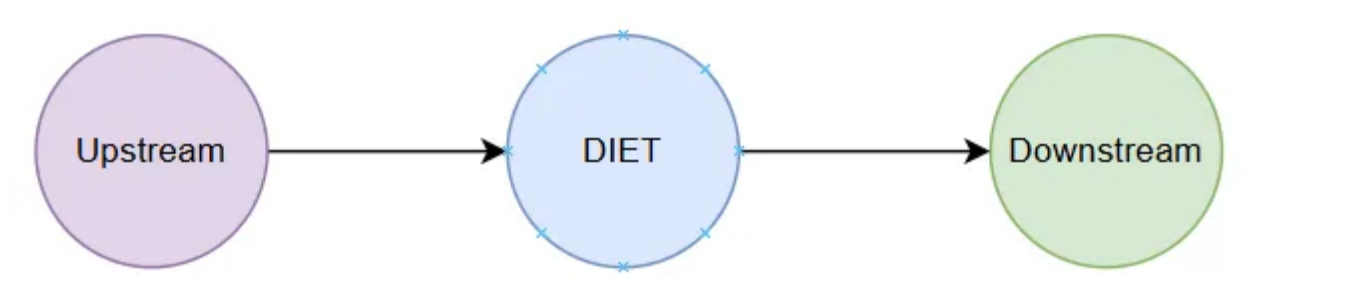


DIET (Data Ingestion and Extraction Tool) Guide

1. 概述

DIET (Data Ingestion and Extraction Tool) 是一个元数据驱动的ETL框架，旨在为部门内部各团队提供统一的数据接入和提取解决方案。它支持多种数据源和目标，包括文件系统 (SFTP/HDFS)、API(RESTful)、数据库(JDBC)等，并通过配置化的方式实现灵活的数据处理流程。



2. 核心概念

2.1 主要组件

组件	描述
Feed	数据流的基本单位，代表一个完整的数据源或目标
Step	处理流程中的单个步骤，可组合成Flow
Flow	由多个Step组成的有序执行序列
Metadata	驱动框架运行的配置信息，存储在数据库中

2.2 关键实体关系

Plain Text

1 Feed (1) → (*) Step

2 Step (*) → (1) Flow

3 Feed (*) ↔ (*) Configuration

3. 元数据配置详解

3.1 Feed基础配置 (Feed_Basic_Info)

字段名	类型	必填	描述
FEED_KEY	VARCHAR(50)	Y	唯一标识符
FEED_NAME	VARCHAR(100)	Y	父文件名
FEED_TYPE	ENUM	Y	类型: ZIP/GZ/CSV/EXCEL
SOURCE_TYPE	ENUM	Y	源类型: SFTP/API/JDBC/HDFS
SOURCE_APP_NAME	VARCHAR(20)	Y	上游应用
SOURCE_CONTACT	VARCHAR(200)	Y	上游联系方式
DESTINATION_APP_NAME	VARCHAR(20)	Y	下游应用
DESTINATION_CONTACT	VARCHAR(200)	Y	下游联系方式
DESTINATION_TYPE	ENUM	Y	目标类型: DB/HDFS/API
ACTIVE_FLAG	BOOLEAN	Y	是否激活
SCHEDULE_CRON	VARCHAR(50)	N	调度表达式
HOLIDAY_PROFILE	VARCHAR(50)	N	节假日配置ID

3.3 Feed列信息表 (Column_Info)

字段名	类型	必填	描述
SUB_FEED_ID	INT	Y	外键，关联 Sub_feed_info

Column_ID	INT	Y	唯一标识符
COLUMN_NAME	VARCHAR(100)	Y	列名
COLUMN_ORDER	INT	Y	列顺序号
COLUMN_TYPE	VARCHARE(50)	Y	列字段类型： int/varchar/date/datetime/numeric

3.4 子Feed配置表 (Sub_Feed_Info)

字段名	类型	必填	描述
SUB_FEED_ID	INT	Y	唯一标识符
FEED_KEY	VARCHAR(50)	Y	外键，关联 Feed_Basic_Info
SUB_FEED_NAME	VARCHAR(100)	Y	子文件名
CONTROL_FILE_NAME	VARCHAR(100)	Y	控制文件名称
FEED_TYPE	ENUM	Y	文件类型： dat/txt/csv/json/xls/parquet

3.5 步骤配置 (FEED_STEPS)

字段名	类型	描述	示例
STEP_ID	VARCHAR(50)	步骤唯一ID	VALI
STEP_NAME	VARCHAR(100)	步骤名称	"文件
STEP_TYPE	ENUM	步骤类型	VALI
STEP_ORDER	INT	执行顺序	2
CONFIG_REF	VARCHAR(50)	关联的配置ID	FILE

RETRY_COUNT	INT	重试次数	3
FAILURE_ACTION	ENUM	失败行为: STOP/SKIP/ALERT	STO

3.6 源配置示例 (SFTP_SOURCE_CONFIG)

SQL

```
1 CREATE TABLE SFTP_CONFIG (  
2     CONFIG_ID VARCHAR(50) PRIMARY KEY,  
3     HOST VARCHAR(100) NOT NULL,  
4     PORT INT DEFAULT 22,  
5     USERNAME VARCHAR(50),  
6     PASSWORD_REF VARCHAR(100), -- 引用密钥管理服务  
7     DIRECTORY_PATH VARCHAR(255),  
8     FILE_PATTERN VARCHAR(100),  
9     ENCRYPTION_TYPE ENUM('NONE', 'PGP', 'AES'),  
10    COMPRESSION_TYPE ENUM('NONE', 'ZIP', 'GZIP')  
11 );
```

3.7 验证规则配置 (VALIDATION_RULES)

JSON

```
1 {  
2     "ruleId": "DATE_VALIDATION_001",  
3     "ruleType": "DATE",  
4     "params": {  
5         "dateFormat": "yyyy-MM-dd",  
6         "position": "HEADER",  
7         "lineNumber": 1,  
8         "columnRange": "10-19",  
9         "expectedValue": "${COB_DATE}"  
10    },  
11    "errorAction": "REJECT"  
12 }
```

4. 典型流程配置

4.1 SFTP文件处理流程

mermaid

```
1 graph TD
2   A[Feed Monitor] --> B[File Unzip]
3   B --> C[File Validation]
4   C --> D[Data Loading]
5   D --> E[File Archival]
```

配置示例：

SQL

```
1 INSERT INTO FEED_STEPS VALUES
2 ('SFTP_SALES_001', 'MONITOR', 'MONITOR', 1, 'SFTP_MONITOR_CFG', 3, 'STOP'),
3 ('SFTP_SALES_001', 'UNZIP', 'PROCESS', 2, 'UNZIP_CFG', 2, 'SKIP'),
4 ('SFTP_SALES_001', 'VALIDATE', 'VALIDATION', 3, 'SALES_VALIDATION_CFG', 1,
```

4.2 API数据加载到大数据平台

mermaid

```
1 graph TD
2   A[Data Readiness Check] --> B[API Loading]
3   B --> C[HDFS Upload]
4   C --> D[EAP Transform]
5   D --> E[Data Archival]
```

配置示例：

JSON

```
1 {
2   "apiConfig": {
3     "endpoint": "https://api.sales.com/v1/data",
4     "authType": "OAUTH2",
5     "pagination": {
6       "type": "OFFSET",
7       "pageSize": 1000
8     },
9     "responseFormat": "JSON"
10  },
11  "hdfsConfig": {
12    "path": "/data/lake/sales/${COB_DATE}",
```

```
13     "format": "CSV",
14     "partitionBy": ["region"]
15 }
16 }
```

5. 使用指南

5.1 新Feed配置流程

1. 注册基础信息

SQL

```
1 INSERT INTO FEED_MASTER VALUES (
2   'API_MONTHLY_INVENTORY',
3   'Monthly Product Inventory',
4   'INGESTION',
5   'API',
6   'HDFS',
7   TRUE,
8   '0 0 5 1 * ?',
9   'GLOBAL_HOLIDAYS'
10 );
```

2. 定义处理步骤

SQL

```
1 INSERT INTO FEED_STEPS VALUES
2 ('API_MONTHLY_INVENTORY', 'READINESS_CHECK', 'VALIDATION', 1, 'INV_READINES
3 ('API_MONTHLY_INVENTORY', 'API_LOAD', 'PROCESS', 2, 'INV_API_CFG', 3, 'ALER
```

3. 配置API源

JSON

```
1 {
2   "configId": "INV_API_CFG",
3   "httpMethod": "GET",
4   "queryParams": {
5     "reportDate": "${COB_DATE}",
6     "detailLevel": "FULL"
7   },
8   "pagination": {
```

```
9     "type": "CURSOR",
10    "nextPageField": "pagination.next"
11  }
12 }
```

5.2 执行流程

通过Autosys调用：

Bash

```
1  #!/bin/bash
2  # diet_launcher.sh
3  FEED_KEY=$1
5  COB_DATE=$2
6  STEP_NAME=$3
8  java -jar diet-core.jar \
9    --feed=$FEED_KEY \
10   --date=$COB_DATE \
11   --step=$STEP_NAME \
12   --config=/path/to/config.properties
```

典型作业定义：

jil

```
1  insert_job: DAILY_SALES_LOAD
2  job_type: CMD
3  command: /opt/diet/bin/launcher.sh SFTP_DAILY_SALES ${COB_DATE} ALL
4  machine: diet-worker-01
5  date_conditions: 1
6  days_of_week: mo,tu,we,th,fr
```

6. 高级功能

6.1 数据质量检查

SQL

```
1  INSERT INTO DQ_RULES VALUES (
2    'SALES_AMOUNT_CHECK',
3    'COLUMN_RANGE',
4    'SALES_FACT.AMOUNT',
5    '{
```

```
6     "min": 0,
7     "max": 1000000,
8     "threshold": 99.5,
9     "action": "QUARANTINE"
10  }'
11 );
```

6.2 错误处理机制

错误处理配置示例：

YAML

```
1  errorHandling:
2    maxRetries: 3
3    retryInterval: 300000 # 5分钟
4    notifications:
5      - type: EMAIL
6        recipients: ["team@company.com"]
7      - type: SLACK
8        channel: "#data-alerts"
9    escalation:
10     afterAttempts: 3
11     contact: "oncall-engineer@company.com"
```

7. 监控与维护

7.1 关键监控指标

指标名称	类型	描述
diet.step.duration	Timer	步骤执行时间
diet.feed.success	Counter	成功运行次数
diet.validation.errors	Gauge	验证错误数
diet.records.processed	Counter	处理记录数

7.2 维护操作

元数据版本控制：

SQL

```
1 CREATE TABLE METADATA_VERSION (  
2   config_id VARCHAR(50),  
3   version INT,  
4   effective_date TIMESTAMP,  
5   config_content JSON,  
6   PRIMARY KEY (config_id, version)  
7 );
```

(注：此处应包含完整的ER图，展示所有配置表及其关系)

附录B：典型错误代码

代码	含义	建议操作
DIET-4001	源文件未找到	检查SFTP连接
DIET-5003	API认证失败	验证令牌/证书
DIET-3008	验证规则不匹配	检查数据内容和规则配置
DIET-2002	目标存储空间不足	清理目标系统或扩容