



ТИНЬКОФФ

Как построить data lineage.

Обзор решений и опыт нашей команды

План доклада

01

Что такое lineage? Каким он бывает?
Для чего можно использовать?

02

Подходы к формированию data
lineage. Pros&cons.

03

Обзор решений для работы с data
lineage.

04

Для чего мы используем lineage в
TEDI? Цели и кейсы

05

Наша реализация: что мы пробовали,
на чем остановились, на какие грабли
наступили, trade offs.

06

Выводы

Пирамида потребностей в **данных**



Базовые потребности (collection)

Доступность данных, актуальность

Удобство (description)

Data Quality, Data Catalog, [Lineage](#)

Высшие (growth)

Сложный анализ, предсказания

Определение lineage

Очень зависит от
потребностей, но попробуем



Определение #1

это набор сложных отношений между датсетами и трансформациями в пайплайне.



Определение #2

процесс *понимания*, документирования и визуализации того, как данные проходят через экосистему



Определение #3

информация, которая описывает движение данных от источника их происхождения по точкам обработки и применения

Области применения

Описать все возможные варианты невозможно.

Из нашего опыта:

- отслеживание зависимостей
- системный анализ для доработок etl
- поиск узких мест (оптимизация процессов)
- change management
- анализ первопричин: аномалии, сбои
- информационная безопасность
- комплаенс: персональные данные, чувствительные данные
- etc

Области применения, обобщение



Поиск источников

отслеживание зависимостей, системный анализ, поиск, узких мест (оптимизация процессов), change management



Анализ влияния

change management, анализ первопричин: аномалии, сбои, информационная безопасность, комплаенс



Подходы к построению lineage

Вручную

Суть подхода в явном указании зависимостей.

Кто использует: SSIS, Pentaho DI

Pros:

- Простота в реализации
- Можно связывать black boxes
- Поддержка множеством etl инструментов

Cons:

- Сложно разрабатывать и поддерживать
- Column-column lineage обычно не доступен

По мета-данным

Необходимо заполнять мета-данные трансформаций

Кто использует: Informatica, SAS

Pros:

- Автоматизированный подход
- Возможен column-column lineage
- Граф выполнения исходя из графа зависимостей

Cons:

- Мета-данные могут не соответствовать действительности
- Сложность в описании неявных зависимостей

Автоматически

Мета-данные извлекаются из кода (настроек)

Кто использует: dbt (table-table only), atlan, TEDI

Pros:

- Не требует дополнительных действий
- Всегда актуальное состояние
- Возможен column-column lineage
- Граф выполнения исходя из графа зависимостей

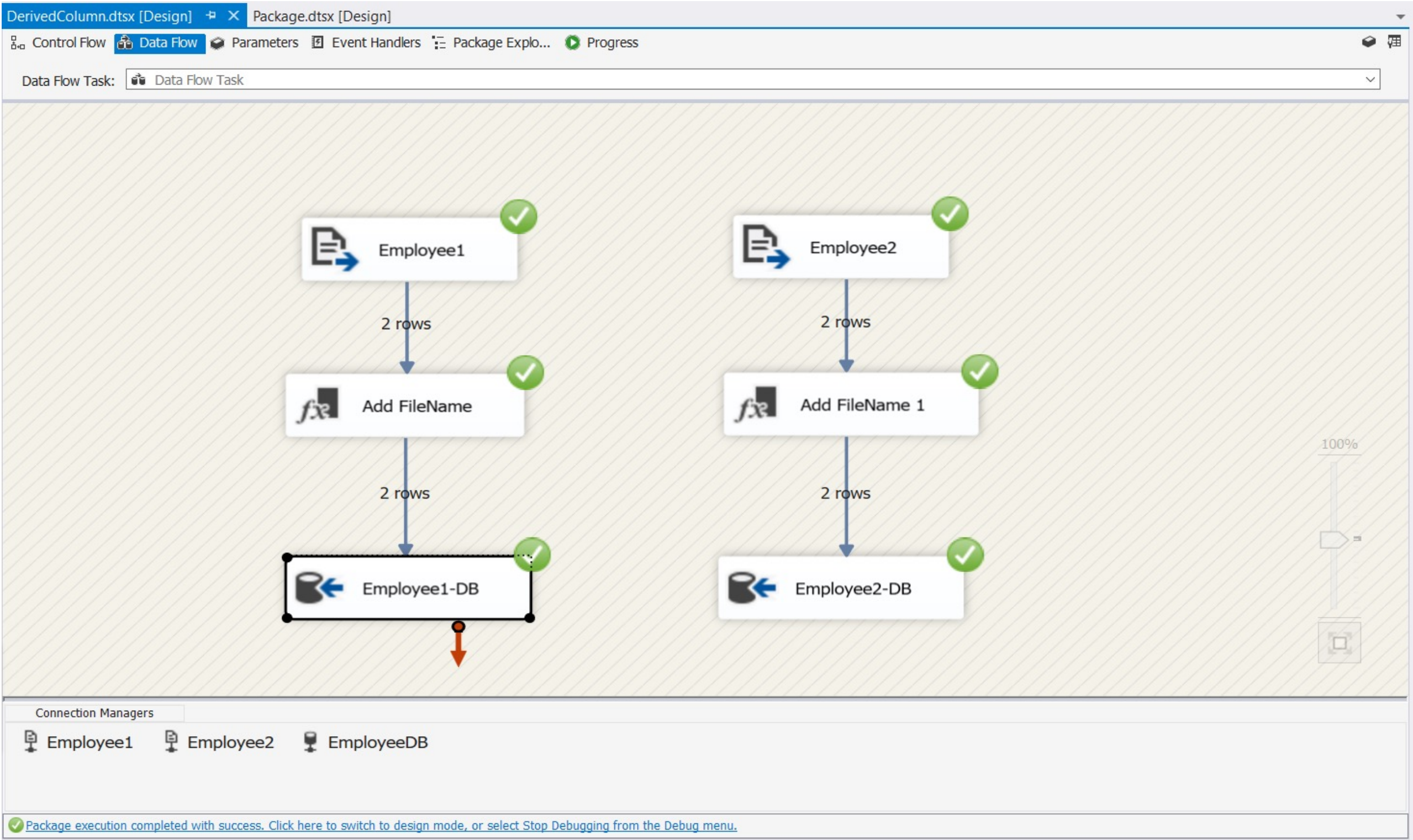
Cons:

- Ограниченное количество инструментов
- Большие различия у синтаксисов



Реализации

SSIS



Informatica

Informatica PowerCenter Designer - [Mapping Designer - Informatica_Tutorial - [Repsvc_Edureka]]

Repository Edit View Tools Layout Versioning Mappings Transformation Window Help

Informatica_Tutorial - [Repsvc_Edureka] 100%

m_Most_Selling_Category

Mapping Designer

Business Components

- psvc_Edureka
 - Edureka_Transformations
 - Informatica_PowerCenter_Tutorial
 - Informatica_Training
 - Informatica_Training_17_4
 - Informatica_Training_New
 - Informatica_Transformation
 - Informatica_Transformation_Edureka
 - Informatica_Tutorial
 - Business Components
 - Sources
 - FlatFile
 - Customer
 - Product
 - ORACLE_XE_HR_SRC
 - TRANSACTIONS
 - Targets
 - TGT_LOW_Customers
 - TGT_Marketing_Regions
 - TGT_Most_Loyal_Customers
 - Cubes
 - Dimensions
 - Transformations
 - Maplets
 - Mappings
 - User-Defined Functions
 - Neel_1

Transformation Details

JNRTRANS Joiner

Name	Source Type	Datatype
INVOICE_NO	Detail	nstring
STOCK_CODE	Detail	nstring
DESCRIPTION	Detail	nstring
QUANTITY	Detail	nstring
INVOICE_DATE	Detail	nstring
UNIT_PRICE	Detail	nstring
CUSTOMER_ID	Detail	nstring
PRODUCT_ID	Detail	decimal
Product_Id1	Master	decimal
Product_Name	Master	string
Price	Master	decimal
Seller	Master	string

AGGTRANS Aggregator

Name	Expression
Product_Id1	Product_Id1
Department	Department
Count	Count(Depart

SRTTRANS Sorter

Name	Key
Product_Id1	
Department	
Count	Yes

Parsing expression Count(Department)

Field:Department found

Save Fetch Log Generate Validate Debugger Session Log Notifications

SAS Data Integration Studio 4.3.4 - prod

Edit View Check Outs Actions Debug Tools Window Help

Workspace Server

olders Inventory Transformations

My Folder
BK1_OUT
CC
Documents
DQ Monitor
DSO
DWH
DWH Front Tables
Models
Obsolete
ONLINE
OptmZ
Products
Shared Data
System
TOI
Users

Search

DDS LOAD TQM EVALUATION (Read-Only)

Up Run Stop

Diagram Code Log Output

Details

Mappings Status Warnings and Errors Statistics Control Flow

Source table: All Tables

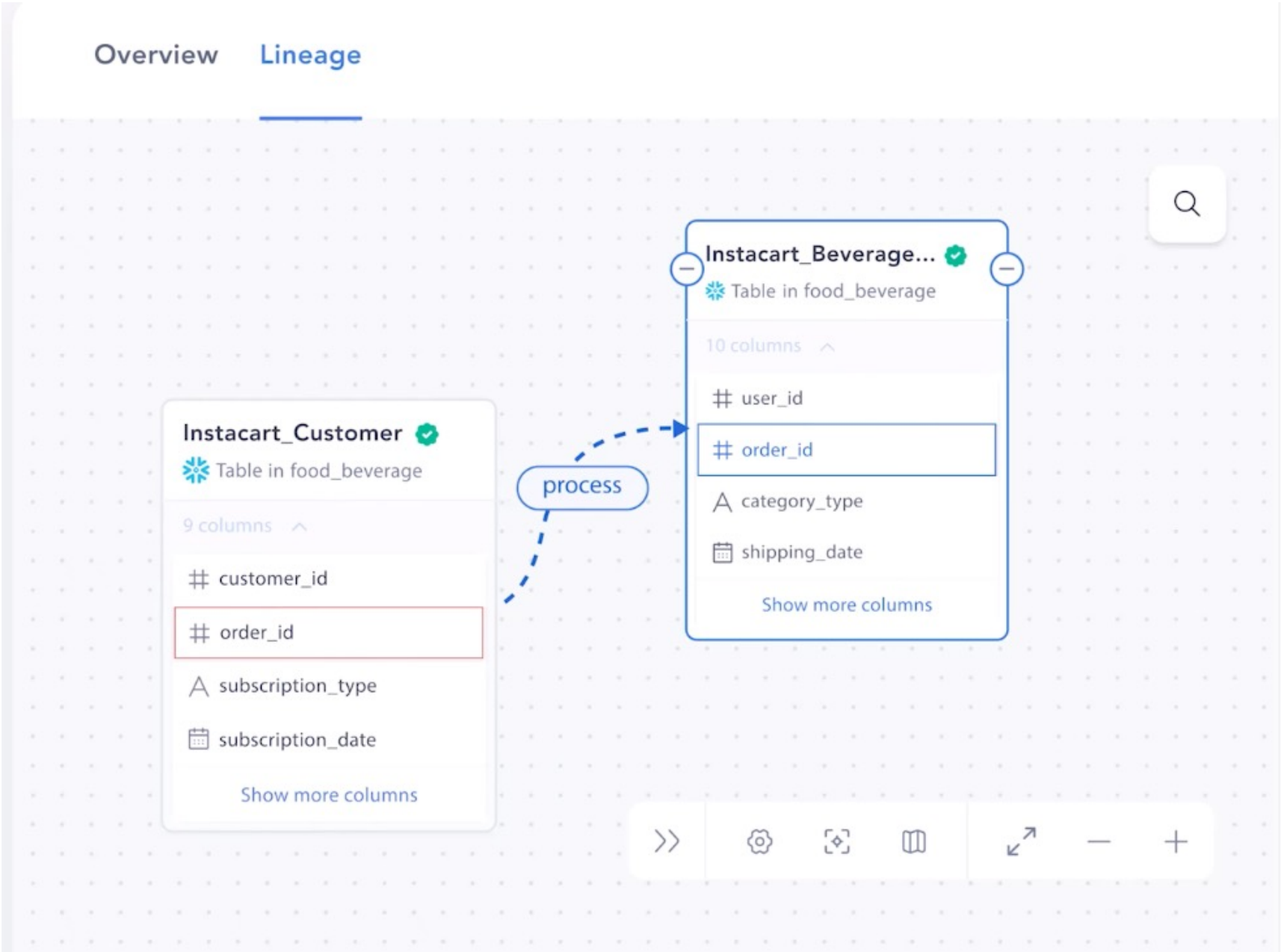
#	Column	Column Description	Type	Length	Informat	For
35	call_rk		Numeric	820.	20.	
36	web_call_rk		Character	21	\$21.	\$21.
37	tqm_evaluation_form_dk	evaluation_form_te...	Numeric	89.	9.	
38	tqm_evaluation_type_dk	type	Numeric	84.	4.	
39	score	score	Numeric	830.6	30.6	
40	calibration_flg	has_calibration_ev...	Numeric	84.	4.	
41	eduid	eduid	Character	32	\$32.	\$32.
42	create_dtm	created_at	Numeric	8	DATETIME25.6	DATETI
43	update_dtm	modified_at	Numeric	8	DATETIME25.6	DATETI
44	deleted_flg	is_deleted	Numeric	84.	4.	
45	deleted_dtm		Numeric	8	DATETIME25.6	DATETI
46	rating_participation_flg		Numeric	84.	4.	
47	tqm_evaluation_period_dk	calculation_period_id	Numeric	89.	9.	
48	total_penalty_amt	penalty_score	Numeric	812.3	12.3	
49	tqm_evaluation_pensly_dk		Character	32	\$32.	\$32.
50	valid_from_dtm		Numeric	8	DATETIME25.6	DATETI
51	valid_to_dtm		Numeric	8	DATETIME25.6	DATETI
52	processed_dtm		Numeric	8	DATETIME25.6	DATETI
53	tqm_evaluation_action_type...		Numeric	84.	4.	
54	action_rk	ext_task_id	Numeric	820.	20.	
55	tqm_evaluation_rk	evaluation_id	Numeric	820.	20.	

Target table: Join (d_deleted_dtm)

#	Column	Column Description	E
7	end_dtm		
8	call_segment_direction_cd	call_direction	
9	call_segment_duration_sec	billing_duration	
10	call_rk		
11	web_call_rk		
12	tqm_evaluation_form_dk	evaluation_form_te...	
13	tqm_evaluation_type_dk	type	
14	score	score	
15	calibration_flg	has_calibration_ev...	
16	eduid	eduid	
17	tqm_evaluation_penalty_dk	penalty_id	
18	create_dtm	created_at	
19	update_dtm	modified_at	
20	deleted_flg	is_deleted	
21	deleted_dtm		case when TQM_EVAL
22	valid_from_dtm		
23	rating_participation_flg		
24	tqm_evaluation_period_dk	calculation_period_id	
25	total_penalty_amt	penalty_score	
26	tqm_evaluation_action_type_dk		
27	action_rk	ext_task_id	

ic Properties

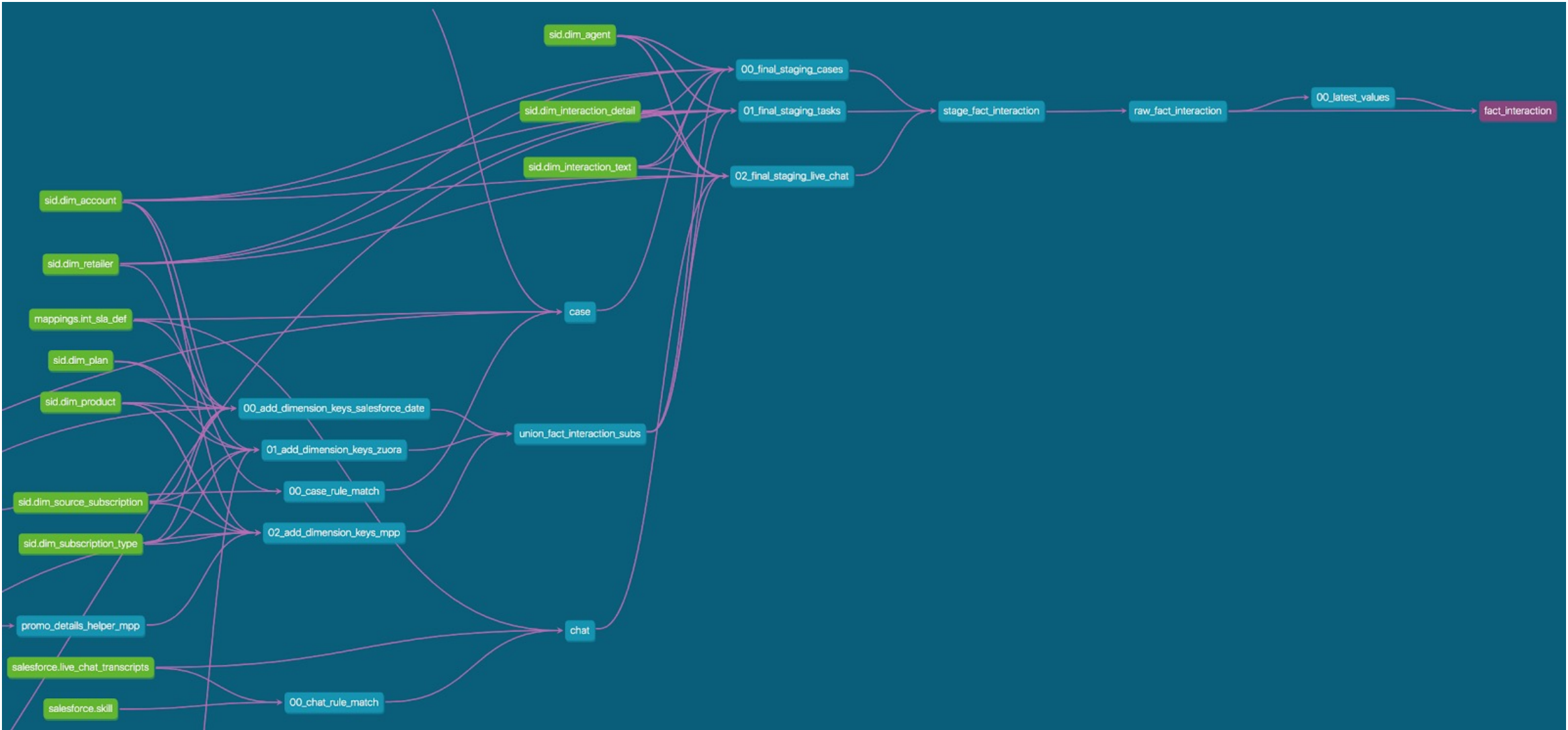
Name	Value
ie	[CL] deleted_dtm
cription	
le Type	Join
r Written	No
gnostic Mode	No
ckpoint	No
always when restarting	No
rol Order	15
adata ID	A5ESKX7U.BB0BONAT



SQL Query

```
SELECT order_id,
       office_sought,
       Format(Sum(total_$), 2)
FROM   combined_party_data
WHERE  election_date =
       "11/12/2019"
GROUP BY candidate,
         office_sought,
         election_year

13 lines hidden
office_sought, election_year
HAVING Sum(total_$) BETWEEN 3000000
AND 18000000
```



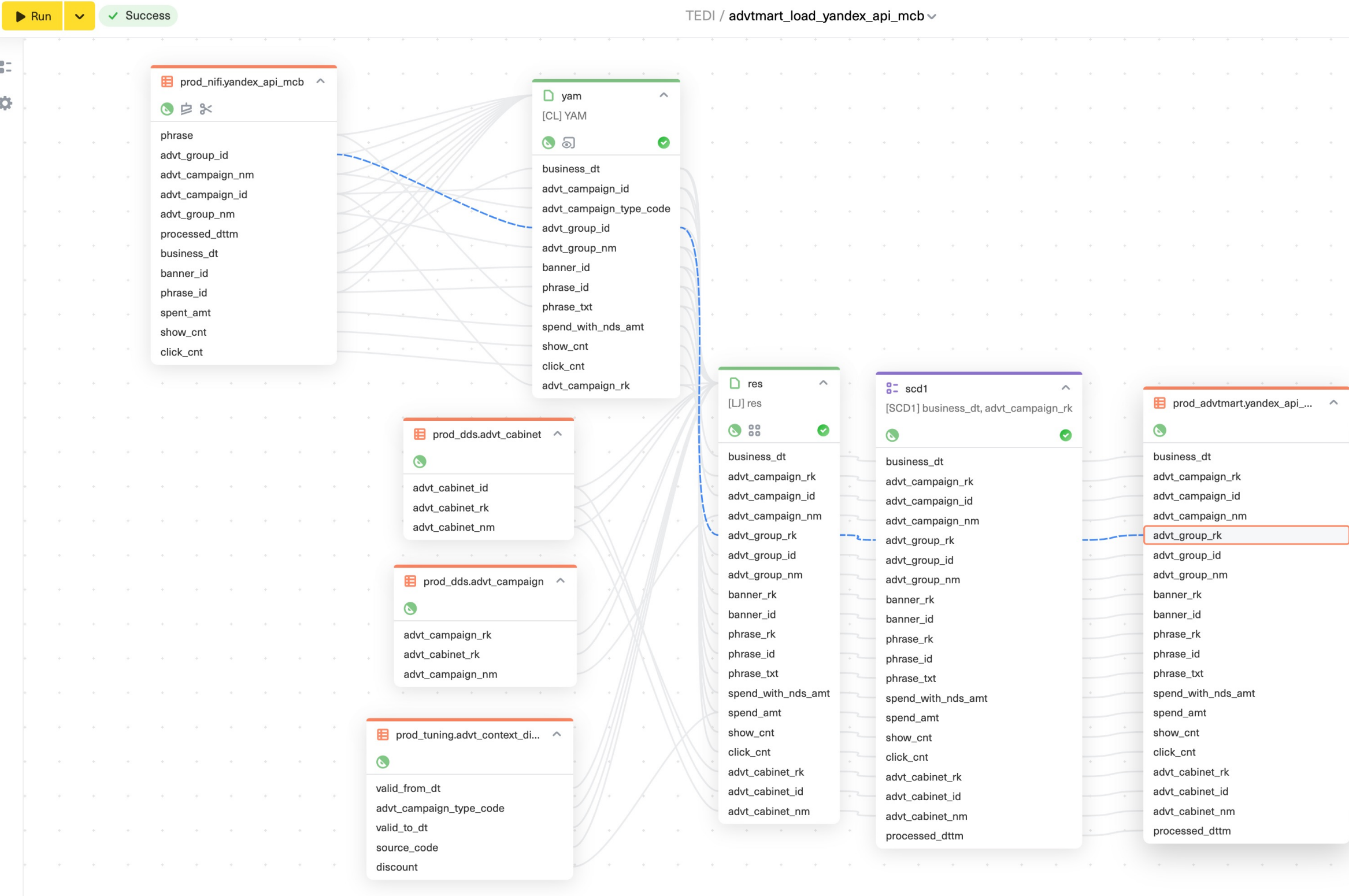
dbt: реализация

dbt обрабатывает шаблонный код для материализации запросов.

- Описание моделей в yml
- Шаблон на jinja2
- *ref* функция для нахождения зависимостей в шаблоне

Можно сломать lineage использовав функцию *ref*, например в комментарии к sql:

```
SELECT * FROM {{ ref('tbl0') }}  
-- INNER JOIN {{ ref('tbl1') }} USING (a)
```



TEDI: редактор

Основные абстракции редактора:

- Источники и таргеты (таблицы)
- Операторы (настройки):
 - SQL трансформация
 - Загрузчики: SCD1, SCD2
 - Генератор ключей
 - И еще 10 операторов
- «Временки» (объект в базе, созданный оператором)
- Связи (между объектами БД):
 - Источник-оператор
 - Оператор-оператор («временка»-«временка»)
 - Оператор-таргет

TEDI: оркестратор

Основные абстракции оркестратора:

- Операторы (исполняемая логика, Airflow Operator):
 - Создание таблиц и вьюх
 - Запуск библиотечного кода загрузчиков и трансформаций
- Связи: порядок исполнения

Порядок исполнения строится на связях между объектами

SQL оператор

Код трансформации

Любой* SELECT запрос, в том числе с CTE и рекурсивными CTE.

* С некоторыми ограничениями

des | + New node ▾

ow all nodes ▾

ij_portal_application

```
select
pa.hid,
pa1.hid as linked_appl_hid
from test_marti.advt_portal_application as pa
inner join test_marti.advt_portal_application
on (pa.linked_appl_id::varchar(32)=pa1.porta
where pa.linked_appl_id is not null
and pa1.portal_application_id is not null
and pa.hid is not null
and pa1.hid is not null
```

SQL оператор

Настройки трансформации

Дополнительные параметры для создания объекта: тип объекта, распределение и тд

ij_portal_application

Options Script

Name (task_id)

ij_portal_application

Description

[IJ] linked_appl_id, portal_application_id

Object type

table

Distributed fields

hid X

WITH Clause

Analyze columns

hid X

Set optimizer



Lineage из sql кода

Статический анализ SQL

Статический анализ - анализ программного обеспечения, производимый (в отличие от динамического анализа) без реального выполнения исследуемых программ.

Что может sql парсер:

- Проверить синтаксис
- Построить AST

Примеры статических анализаторов SQL: Calcite, python-sqlparse, libpg_query

ANTLR

грамматики

ANTLR принимает в качестве входных данных грамматику, определяющую язык, и генерирует код парсера для целевого языка (Java, Python, C++).

Pros:

- Генерирует код в множество языков: Java, Python, C++, Go, Rust
- Доступно много SQL грамматик: postgresql, mysql, hive, trino

Cons:

- Производительность
- Большинство грамматик *неофициальные*

Парсинг запросов GP

Greenplum популярная MPP база, НО:

- Нет официальной ANTLR грамматики
- Внутри GP лексер и парсер реализованы на flex и bison
- Расширенный синтаксис относительно sql92
- Generic решения типа Calcite и python-sqlparse не справляются

Попытка #1

Пробовали ANTLR:

- Взяли грамматику для PG
- Допилили поддержку GP: distributed в ddl, etc
- Ломалось на сложных запросах
- Высокая сложность фиксов, в некоторых случаях казалось, что починить уже не получится
- Плохой перфоманс, жрало много ресурсов

Попытка #2

Пробовали libpg_query (pg_last):

- Нативный парсер грамматики **Postgresql**
- *C* код с биндингами под питон
- Прожевывает селекты любой сложности
- Очень производительно
- Blackbox

Trade offs

**В итоге мы остановились на варианте с pg_last.
Какие ограничения остались:**

- только DML запросы
- Поддержка синтаксиса pg/gr, другие движки не получится использовать


Ограничения связанные с статическим анализом:

- У нас нет state of the world, мы не знаем реальный ddl источника
 - SELECT * запрещен
 - Возможно обращение к несуществующей таблице
 - Возможно обращение к несуществующему полю
- Статически невозможно разрулить ambiguous column name (e.g. SELECT t0.col, t1.col FROM t0, t1)
 - Обращение к полю без указания имени таблицы
 - Обязательное указание алиасов в селекте

Пример оператора (SCD1)

Настройки трансформации

Входные параметры для исполнения оператора: входная/выходная таблица, ключи загрузки, алгоритм и тд

 Greenplum SCD1

scd1

OptionsScript

Name (task_id)

scd1

Description

[SCD1] application_rk, event_code,event_dttm,appsflyer_device_id, install_dttm

Input table

to_load

Target table

test_advtmart.apf_mobile_application_event

Loading method

delete-insert

KEY

application_rk Xevent_code Xevent_dttm Xappsflyer_device_id X

install_dttm X

Compare fields



Lineage из параметров оператора

Lineage операторов

TEDI поддерживает не только SQL преобразования.

А значит и на lineage они тоже влияют. Каким бы не был сложным оператор, он ссылается на сущности в бд или временную таблицу. Некоторые операторы образуют связи неявно.

В операторах работает валидация, НО только в случае с ссылками на «наблюдаемый» сущности. В случае с ссылками на объекты БД, подсказать или провалидировать мы не можем. У нас нет **state of the world**.

Lineage в TEDI

Цели и кейсы



Нужды системных аналитиков

Оценка влияния изменений, анализ аномалий, поиск узких мест



Нужды ETL разработчиков

Визуализация всех зависимых полей в процессе: явные зависимости и зависимости из фильтров и условий соединений; отображение различных трансформаций



Нужды ETL тестировщиков

Валидация процесса end-to-end: список источников, схема целевой таблицы.

Demo

test: smart_data

← → ↺ 🏠 ⚠ Не защищено | [tedi-ui-mr-tedi-1462.tedi-ui.v2.dev2.k8s.tcsbank.ru/dag/smart_data?taskId=tmp0](#) 🔖 ☆ 🛑 🌐 ⚙ ⌵ 🖨 📱 🕒 ⋮

[SURVIVAL | Faceb...](#) [SURVIVAL&CO.專...](#) [🚲 Календарь](#) [📰 62 полезных инс...](#) [🎬 The Secrets Of St...](#) [🐞 The Iterative UDP...](#) [🔧 Программирован...](#) [☀ Быстрое веб-при...](#) [🚲 Календарь любит...](#) [🌐 www.erlang-facto...](#) » | 📁 Другие закладки

▶ Run ▼

TEDI / smart_data ▼

Nodes + New node ✕

⚙ Show all nodes ▼

Выводы



Автоматический lineage возможен

За счет статического анализа можно строить полноценный lineage



Generic решение сложно реализовать

Синтаксисы движков очень различаются, только простые запросы



На рынке почти нет готовых решений

Atlan и dbt ближайшие соседи TEDI, но у каждого свой путь



Планируем добавить мету из БД

К статическому анализу добавим данные о типах полей, существовании таблиц и тд



ТИНЬКОФФ

Спасибо!

SmartData 2022