

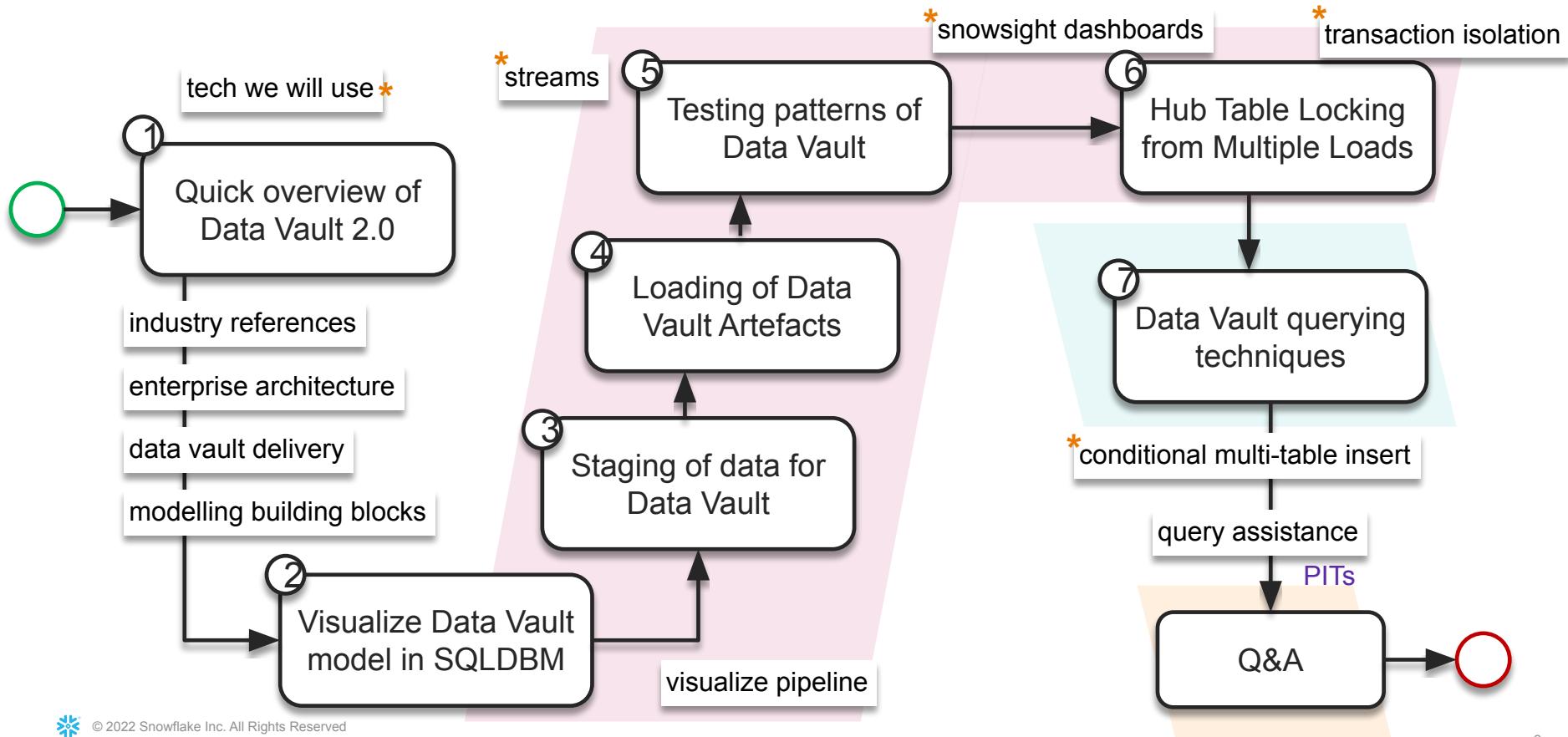


PRAGMATIC GUIDE TO BUILDING A DATA VAULT ON SNOWFLAKE

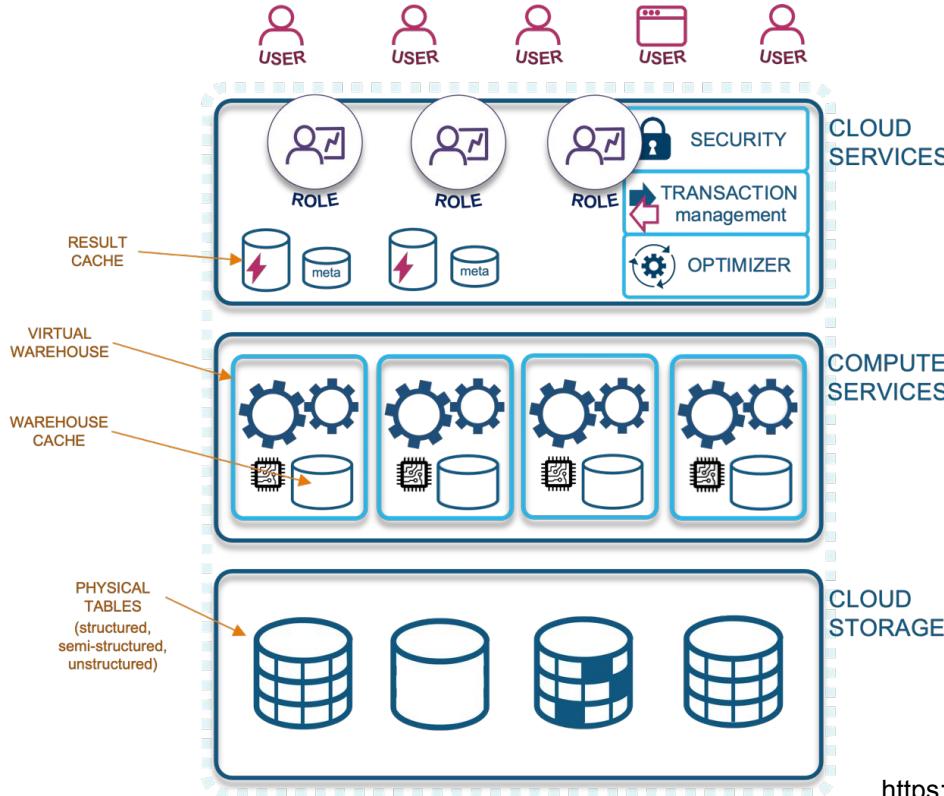
Overview, References & What Works

... explanation and a follow along demo

ZERO TO (Data Vault on) SNOWFLAKE



SNOWFLAKE CACHE



Services Layer

- **Metadata Cache**
 - Object Definitions, Statistics
- **Results Cache – same ROLE can access**
 - Exact Results from Exact Queries
 - Lasts for 24 Hours, 7-day max
 - Underlying data cannot have changed
 - Context functions case expiries

Compute Layer – autoscale, autosuspend

- Virtual Machines
- SSD cache
- Suspended flushes Cache
- Can have partial cache

Storage Services – centralized (“Remote”)

- Database, schema, tables, views, etc.
- S3 / Blob / Cloud Storage

https://docs.snowflake.com/en/sql-reference/functions/result_scan.html

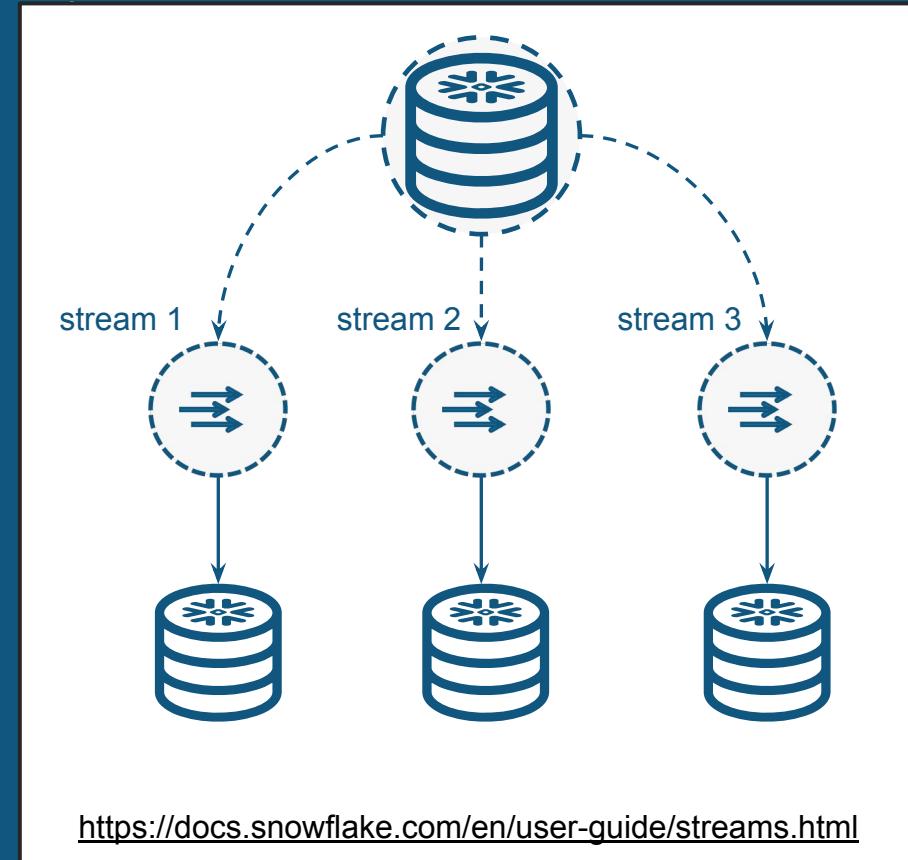
<https://community.snowflake.com/s/article/Caching-in-Snowflake-Data-Warehouse>



STREAMS (CDC)

- Change Data Capture (CDC) for Snowflake tables
- Provides a mechanism for detecting and recording data **CHANGES** to tables over time – this is the **OFFSET**.
- Provides a timestamp, operation and column values
- Useful for efficient, **INCREMENTAL** data processing within ELT
- Stream Types: Standard or **APPEND ONLY**
- **NO LIMIT** on Streams per object
- DML **consumes** streams, **progresses** OFFSET

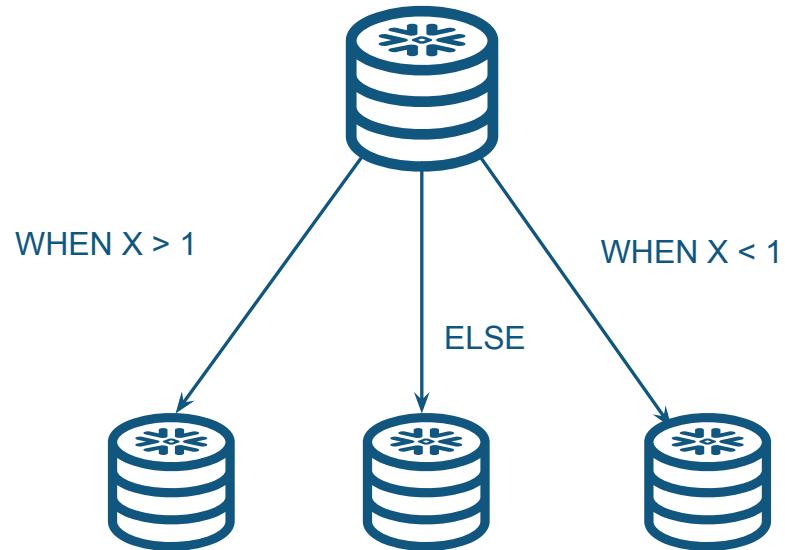
create or replace stream reconcile_hub_account
on table hub_account append_only = true



insert all when (aof_week_lastday=1) ...

MULTI-TABLE INSERT

- Single source query, can include joins
- **UNCONDITIONAL** – insert to all identified targets
- CONDITIONAL **WHEN** clause can be
- FIRST - executes for first TRUE condition
- ALL – executes for all TRUE conditions
- OVERWRITE – reloads target



<https://docs.snowflake.com/en/sql-reference/sql/insert-multi-table.html>



© 2022 Snowflake Inc. All Rights Reserved

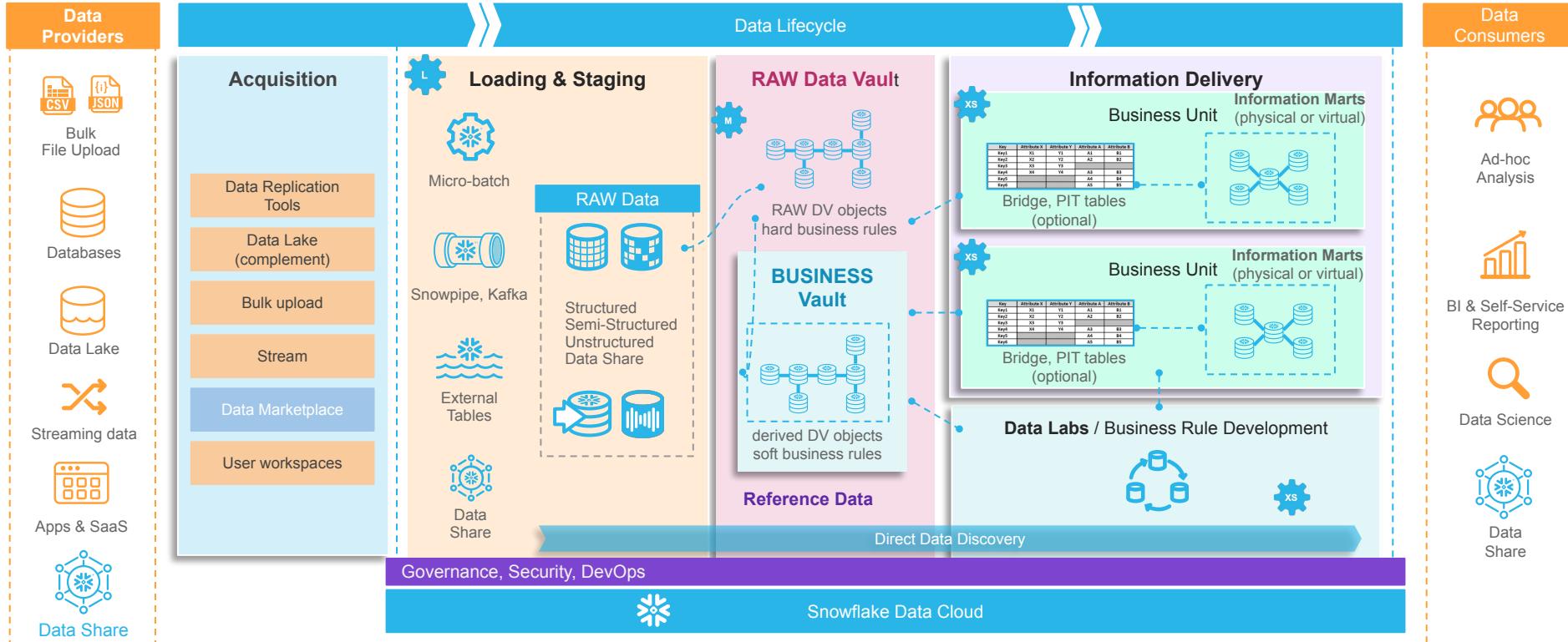
DATA VAULT 2.1

AUDIT, AGILE, AUTOMATED



MULTI-TIER (+CHARGEBACK) DATA VAULT ARCHITECTURE

Customer managed
Snowflake managed



“

*“A system of Business Intelligence containing the necessary components needed to accomplish **enterprise vision** in **Data Warehousing and Information Delivery**”*

Dan Linstedt, Data Vault inventor

*“A data warehouse is a subject-oriented, integrated (**by business key**), **time-variant** and **non-volatile** collection of data in support of management’s decision-making process, and/or in support of **auditability** as a system-of-record.”*

Bill Inmon, father of the data warehouse at World Wide Data Vault Consortium (wwdvc) 2019

“Data Vault reaches a tipping point – Data Vault modelling techniques are going to reach a tipping point in 2020 where a plurality of projects that involve building or refactoring the “hub” layer of a 3-tier data warehouse architecture will employ this modelling technique”.



ENTERPRISE ARCHITECTURE

Business process as defined as a behaviour element that groups behaviour based on an ordering of activities

A **business object** is a passive element that has relevance from a business perspective

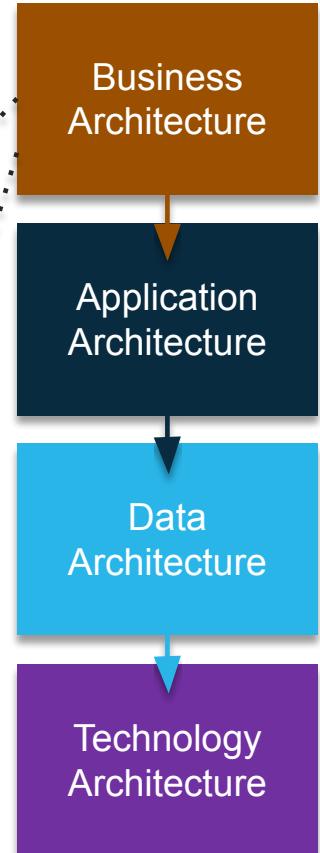
- ArchiMate 2.1

A **Business Capability** defines what the business does...

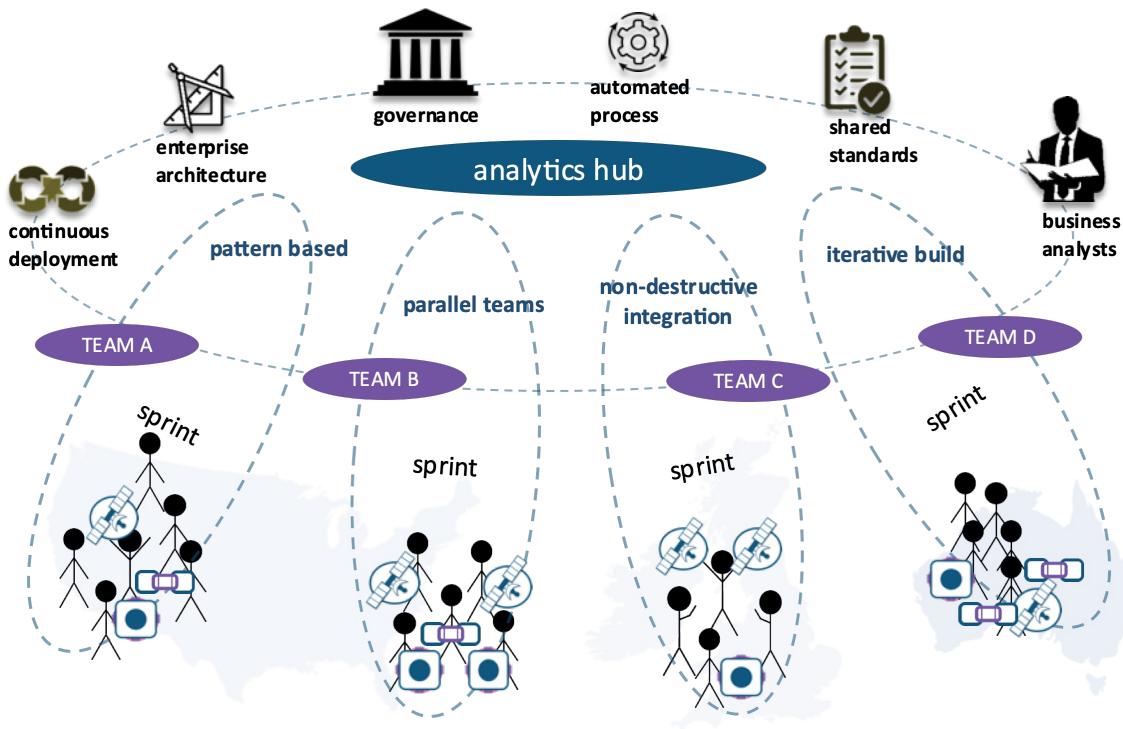
Capabilities are based on **Business Objects**

Business Objects are tangible things commonly recognized by the business, ex. agreement, customer, account, insurance policy, claim, asset, agent, plan, message, research and human resource

- *The Business Architecture Body of Knowledge (BIZBOK)*

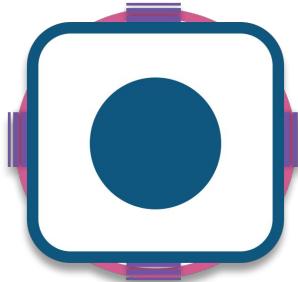


DATA VAULT DELIVERY



- **Shared approach & standards**
- **Central governance**
- **Independent agile teams**
- Data model that relates to **business architecture**
- Common Hubs but **autonomous** business cases
- **Scalable, flexible, rapid analytics**

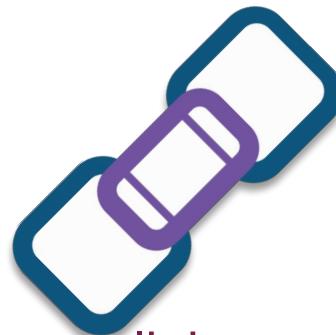
DATA VAULT MODELLING BUILDING BLOCKS



hubs

contains a unique list of **business objects** that represents a domain or concept within the enterprise, everything else connected to the hub gives us more context

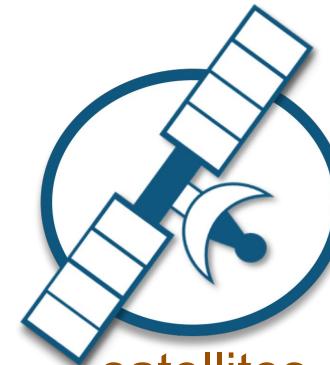
immutable business key



links

represents the relationship between two or more **business objects**, representing a part of one or many business processes or even value streams

unit of work

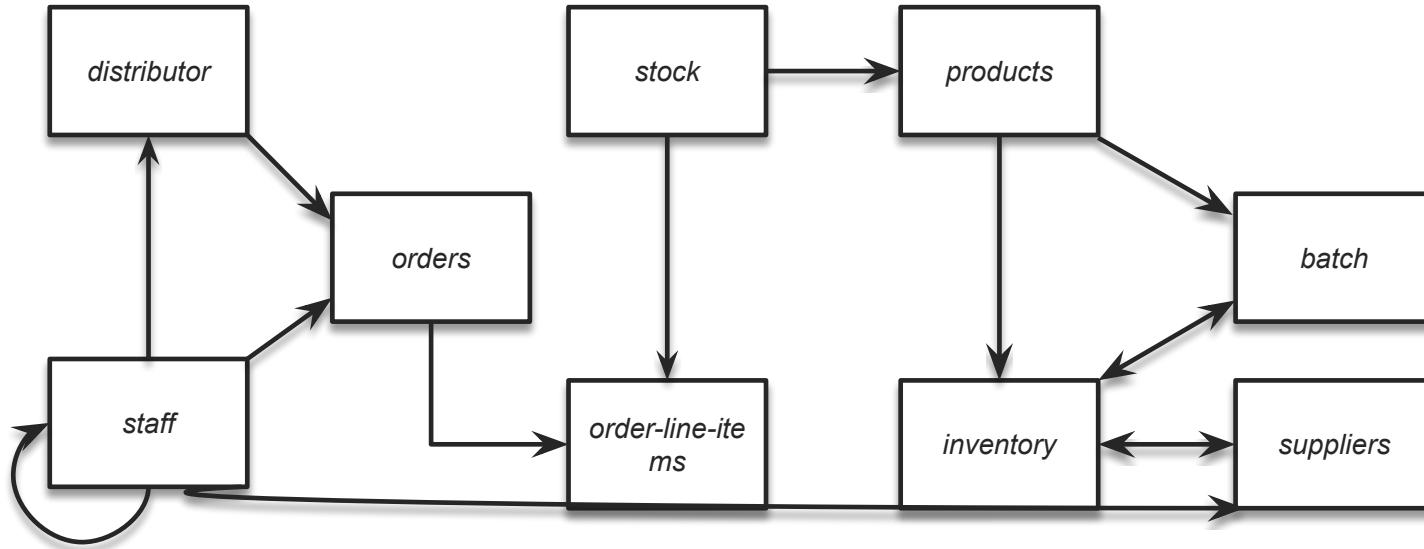


satellites

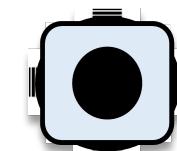
descriptive change-tracking content either describing the **business object** (hub-satellite) or the unit of work (link-satellite)

descriptive content

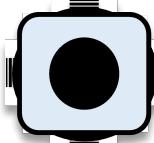
BUSINESS CAPABILITY EXAMPLE



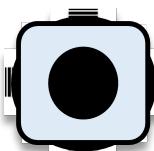
MODELLED BUSINESS OBJECTS



hub_distributor



hub_order



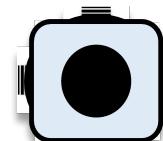
hub_stock



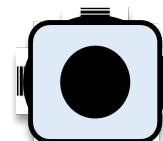
hub_batch



hub_staff

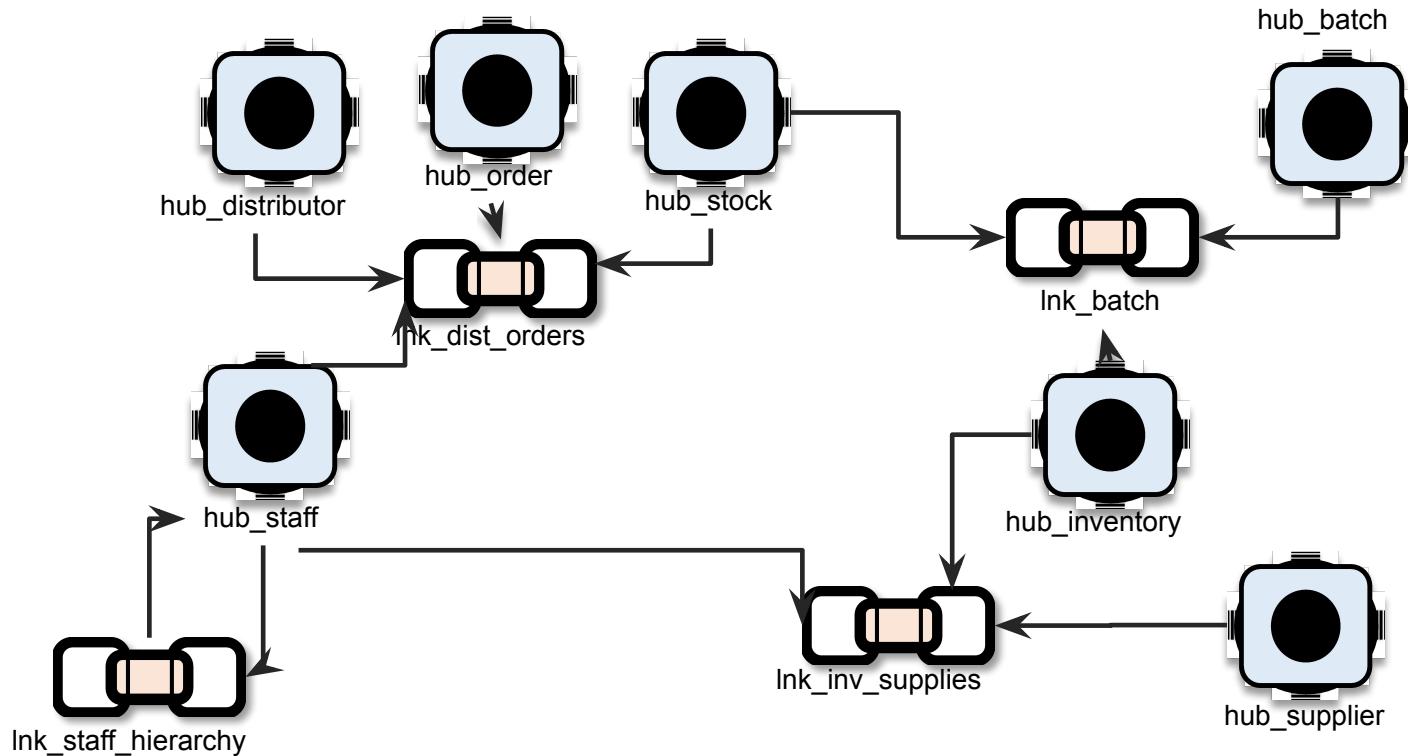


hub_inventory

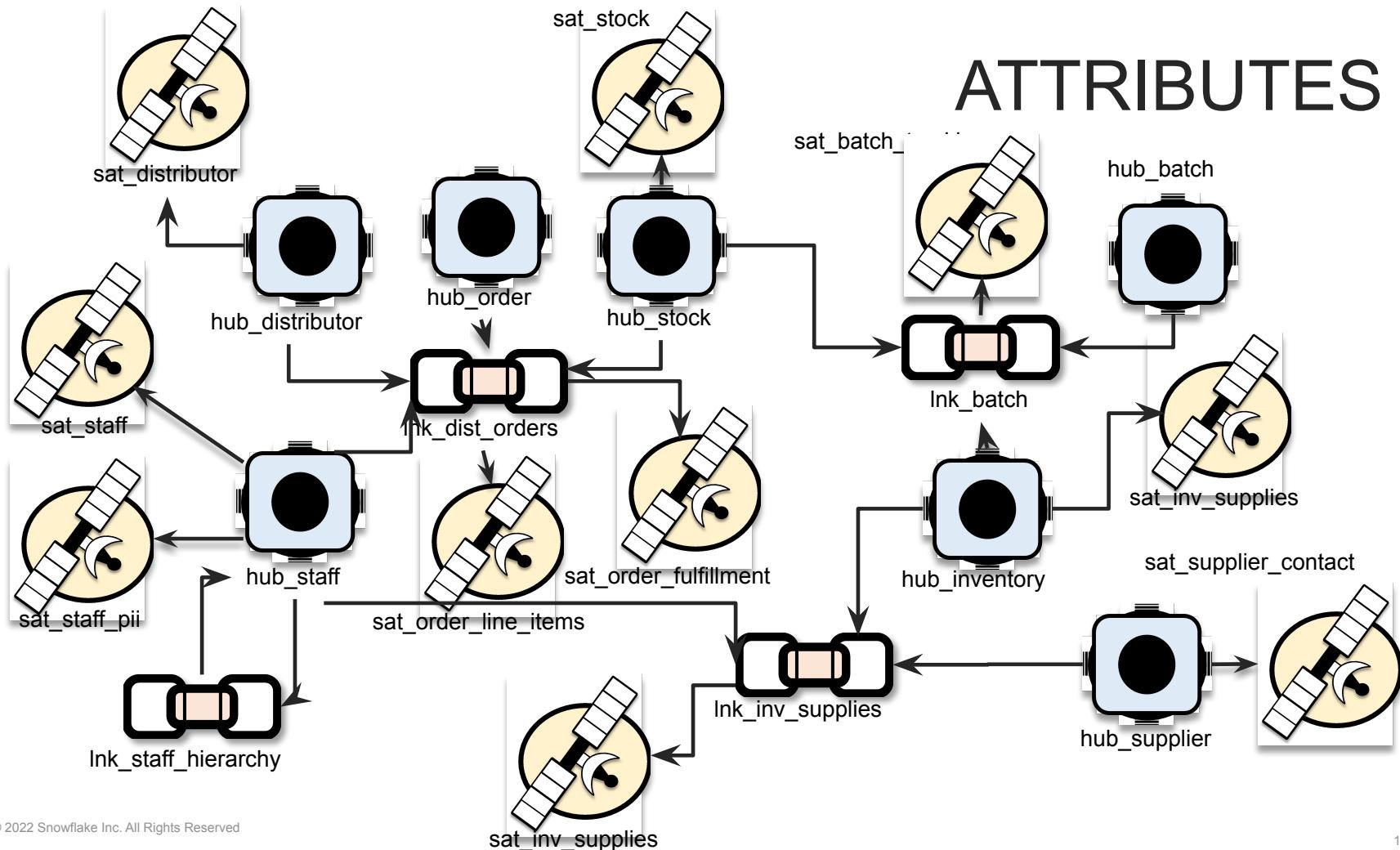


hub_supplier

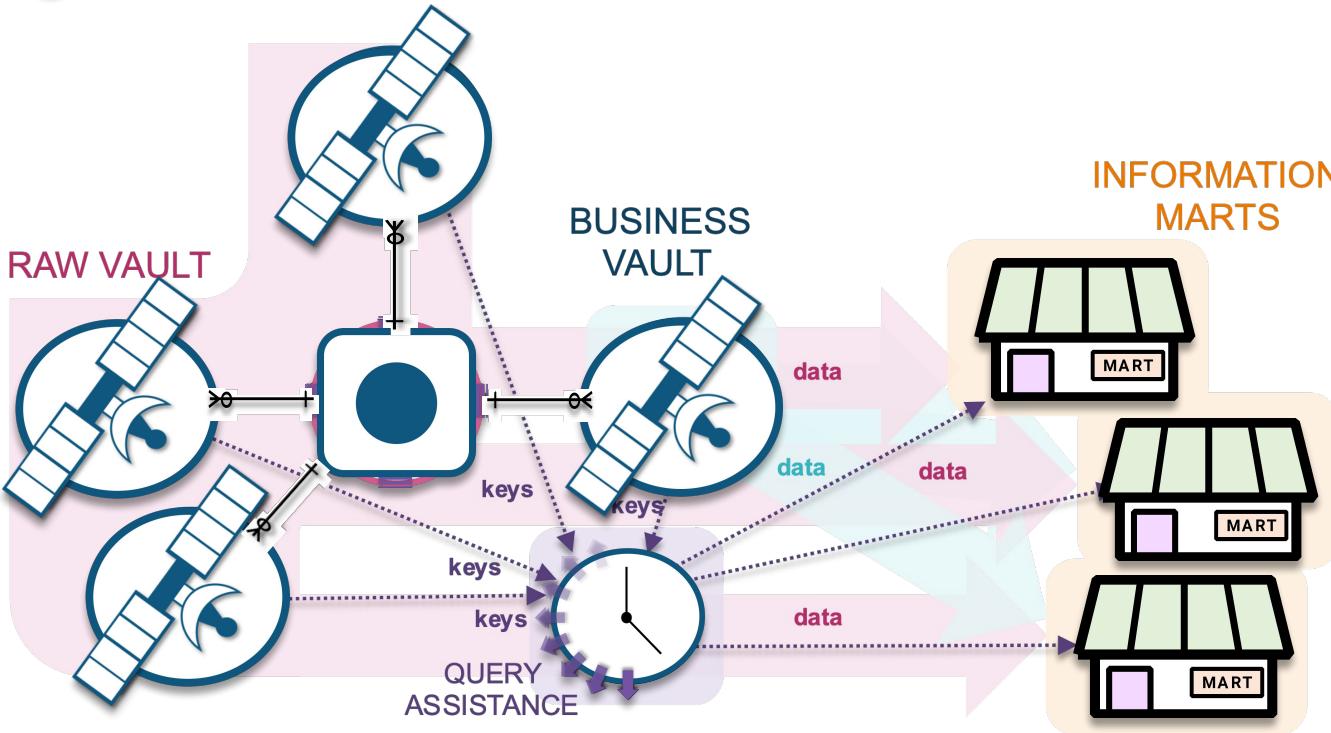
UNIT OF WORK / BUSINESS PROCESS



ATTRIBUTES



② DATA VAULT 2.0 MODEL



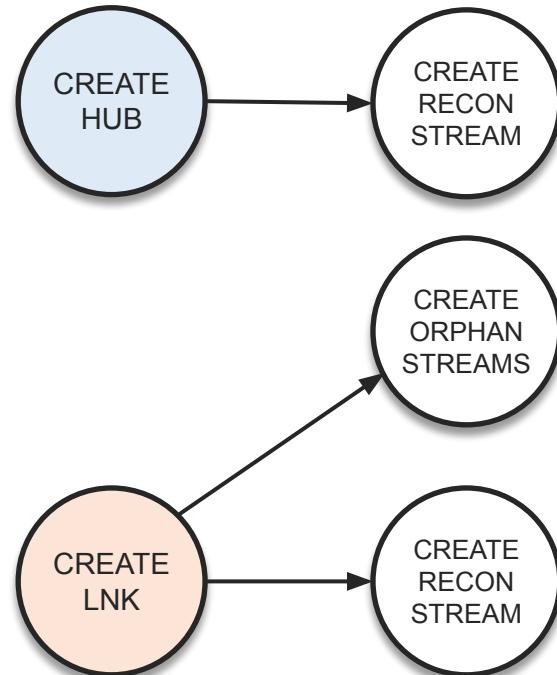
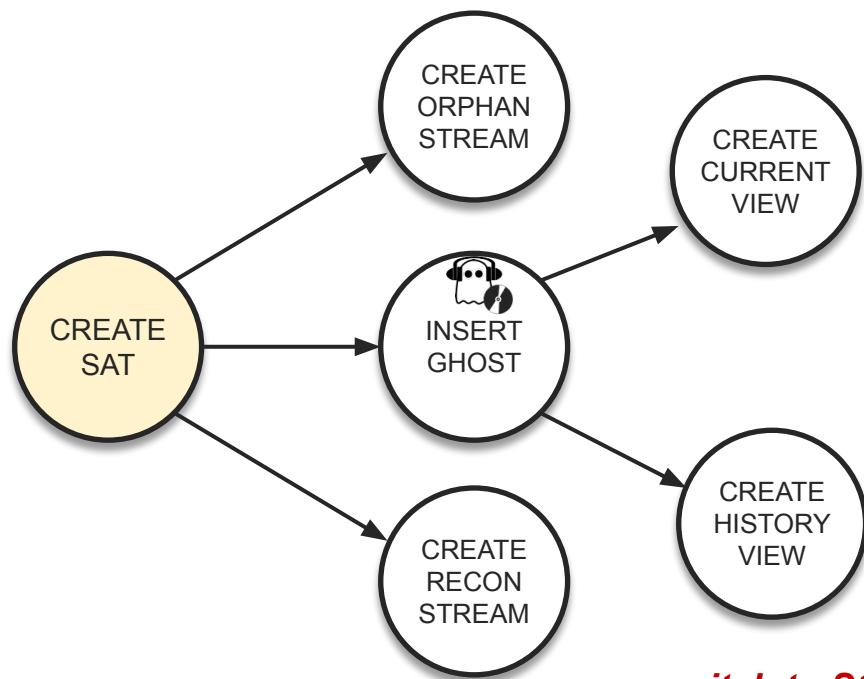
-switch to SQLDBM

- Raw Vault modelled as *hubs, links & satellites*
- Business Vault **sparingly** modelled
- Query Assistance provided by *PITs & Bridges*
- Information Marts delivered as *VIEWS*



INITIALISE DATA VAULT

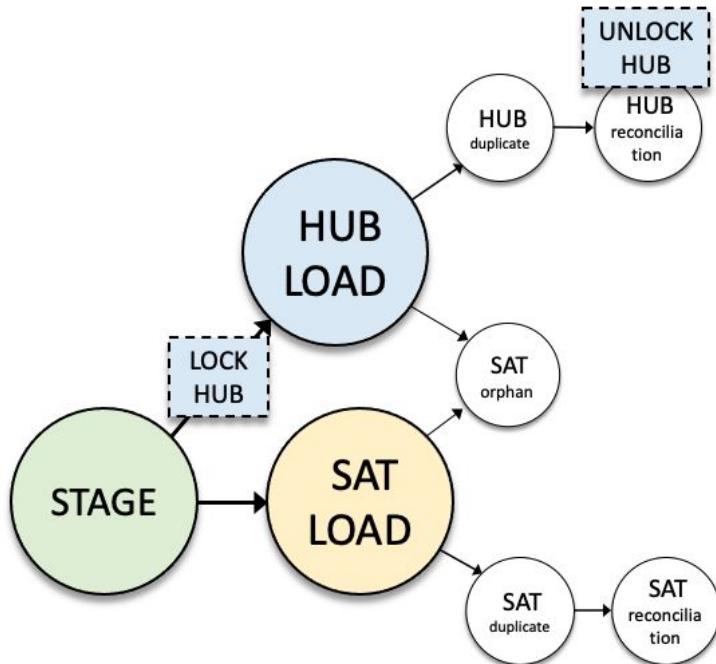
run once (still a template)



-switch to Step_01_Setup_Model

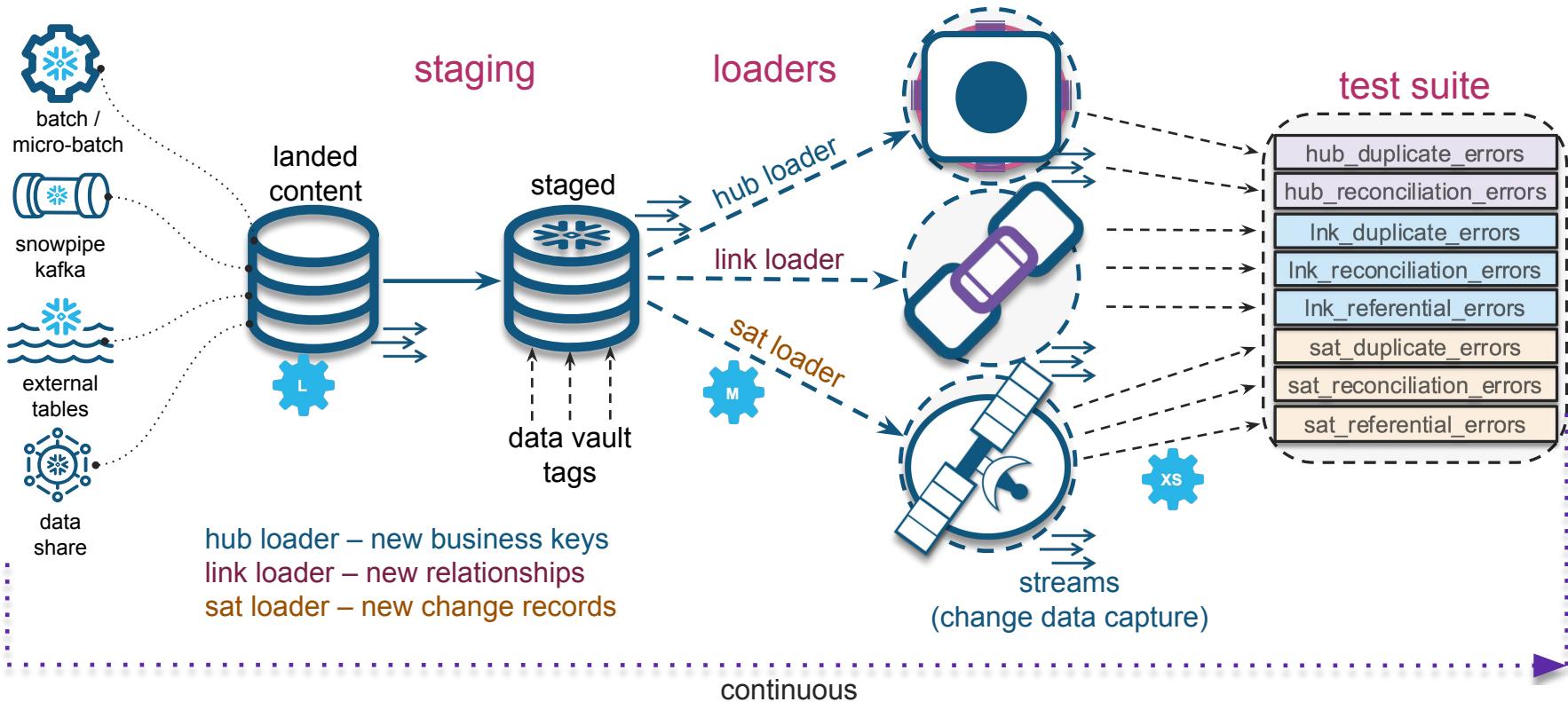


SAMPLE ORCHESTRATION



1. Landed content is STAGED with Data Vault tagging
2. The following will happen in parallel:
 - a) Lock is attained for Target HUB Table
 - b) HUB is loaded
 - c) HUB Tests
 - HUB duplicate check
 - HUB reconciliation check
 - d) Unlock HUB
2. SAT is loaded
 - b) SAT Tests
 - If HUB LOADED, SAT-orphan
 - SAT duplicate check

TEMPLATED INGESTION PIPELINES



③

STAGED DATA VAULT TAGS

Tag	Description	Data Type	Hub	Link	Sat
Business Key (not a tag)	<i>Immutable</i> key representing the business object	text	✓		
Surrogate Hash Key	<i>Durable</i> key optimized for join efficiency	binary/text	✓	✓	✓
Load Date	Timestamp of data <i>loaded</i> to vault	timestamp	✓	✓	✓
Record Source	<i>Lineage</i> describing business process source	text	✓	✓	✓
Applied Date	When the business events applies, extract date	timestamp	✓	✓	✓
Business Key Collision Code	Used to ensure business keys are <i>uniquely</i> represented in a hub	text	✓		
Surrogate Sequence Key	<i>Incremental</i> key, for optimizing joins	numeric			✓

there are more you can add to increase lineage, ex. jira-id, job-id etc

-switch to Step_02_populate_for_one_day hashing at least MD5, SHA1 is better



TEST HARNESS



⑤

DATA VAULT TEST FRAMEWORK

template based

hub_duplicate_errors

hub_reconciliation_errors

lnk_duplicate_errors

lnk_reconciliation_errors

lnk_referential_errors

sat_duplicate_errors

sat_reconciliation_errors

sat_referential_errors

Cat	Test	Hub	Link	Sat
1	Have we loaded duplicates into the target table?	✓	✓	✓
2	Does the target reconcile with the source?	✓	✓	✓
3	Have I maintained referential integrity?		✓	✓



TEST: CHECKING FOR DUPLICATES

hub business keys, link relationships, satellite current details

COLUMN	DESCRIPTION	EXAMPLE
TABLENAME	Hub / Link / Satellite name	HUB_ACCOUNT
SOURCE_TABLENAME	Lineage	CARD_MASTERFILE
LOADDATE	Same as the load date used in staging	\${LOADDATE}
RUNDATE	When test was executed	current_time()
HUB / LINK KEY ERROR COUNT	Hubs – duplicate hash key, duplicate business key collision code + business key Links – duplicate link-hash key, duplicate of all hub-hash keys Satellites – duplicate parent hash-key + current load date	err > 0

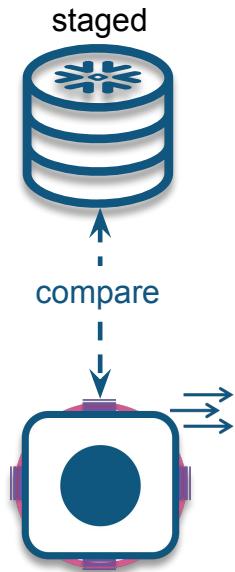
highlighted portion is the same for each test table

TEST: RECONCILIATION

staged source to data vault artefact

COLUMN	DESCRIPTION	NOTES
NEW KEYS / DATA	Hubs – new hash keys, business keys Links – new hash keys, relationships Satellites – new hash key + load date	Uses STREAMS
STAGED KEYS / DATA	Hubs – staged hash keys, business keys Links – staged hash keys, relationships Satellites – staged hash key + load dates	Uses METADATA CACHE
*DISTINCT KEYS / DATA	Hubs – distinct staged hash keys, business keys Links – distinct staged hash keys, relationships Satellites – distinct staged hash key + load dates	Must query all MICRO-PARTITIONS
TOTAL KEYS / DATA	Hubs – hub record total count after load Links – link record count after load Satellites – satellite record count after load	Uses METADATA CACHE
RECON ERROR COUNT	Hubs – missing hash keys, business keys Links – missing hub & link hash keys Satellites – missing hash key + load dates	err > 0

*optional test

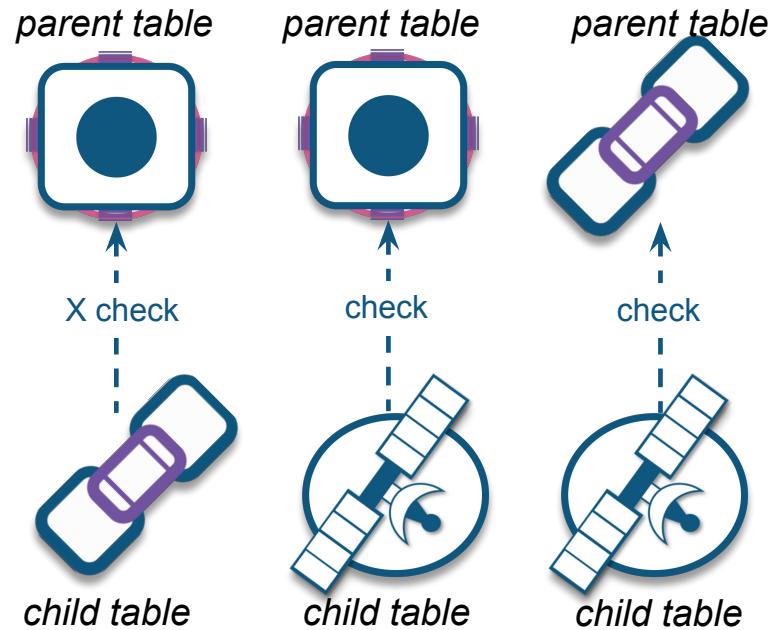


TEST: REFERENTIAL INTEGRITY

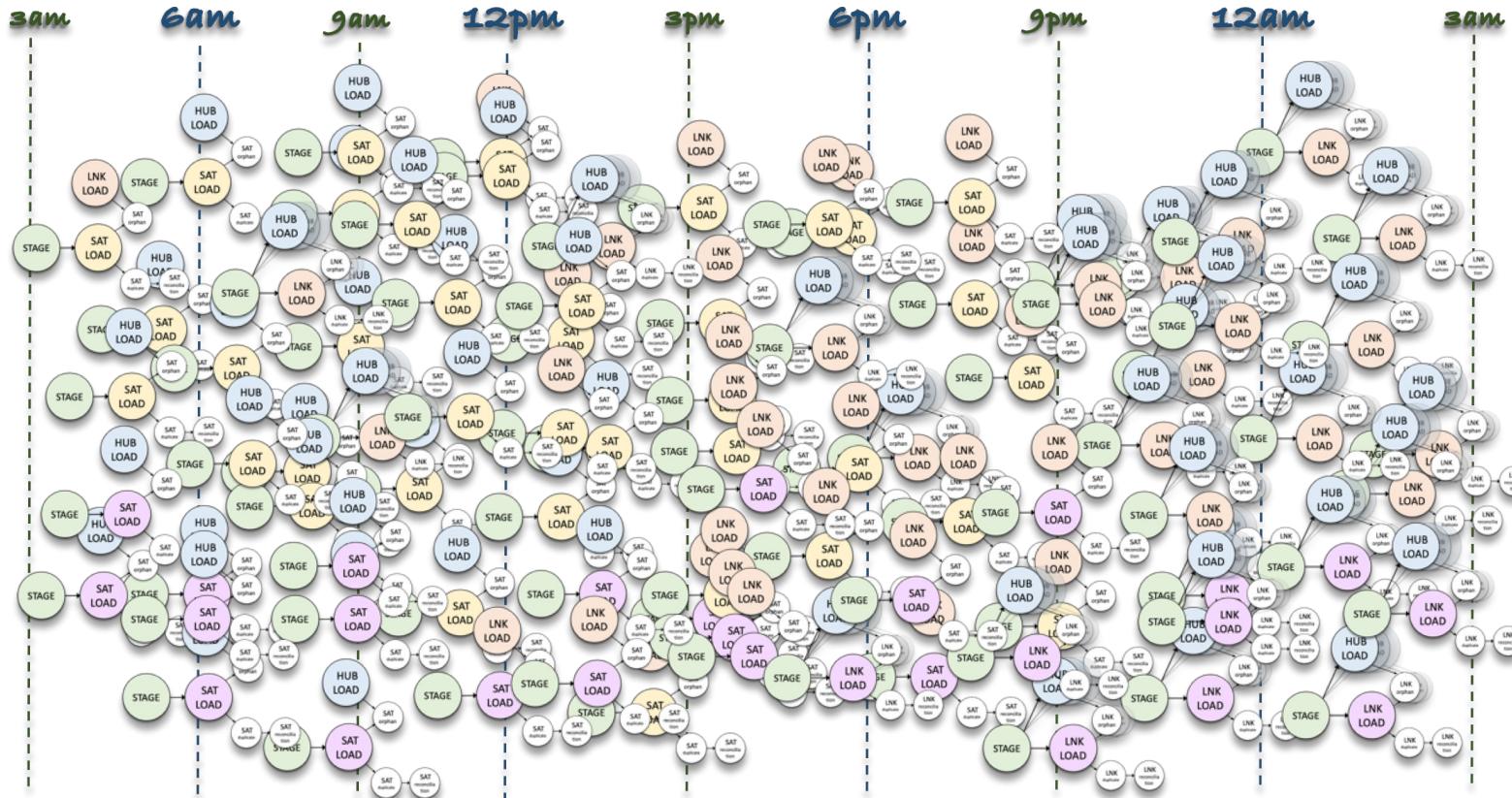
satellites to parent hub or link, link to parent hub(s)

COLUMN	DESCRIPTION
MISSING PARENT KEY ERROR COUNT	Links – Link hub hash keys not found in adjacent hub table(s) HUB-Satellites - Satellite hub hash key not found in Parent Hub table LINK-Satellites – Satellite link hash key not found in Parent Link table

AKA *Orphan-checks*



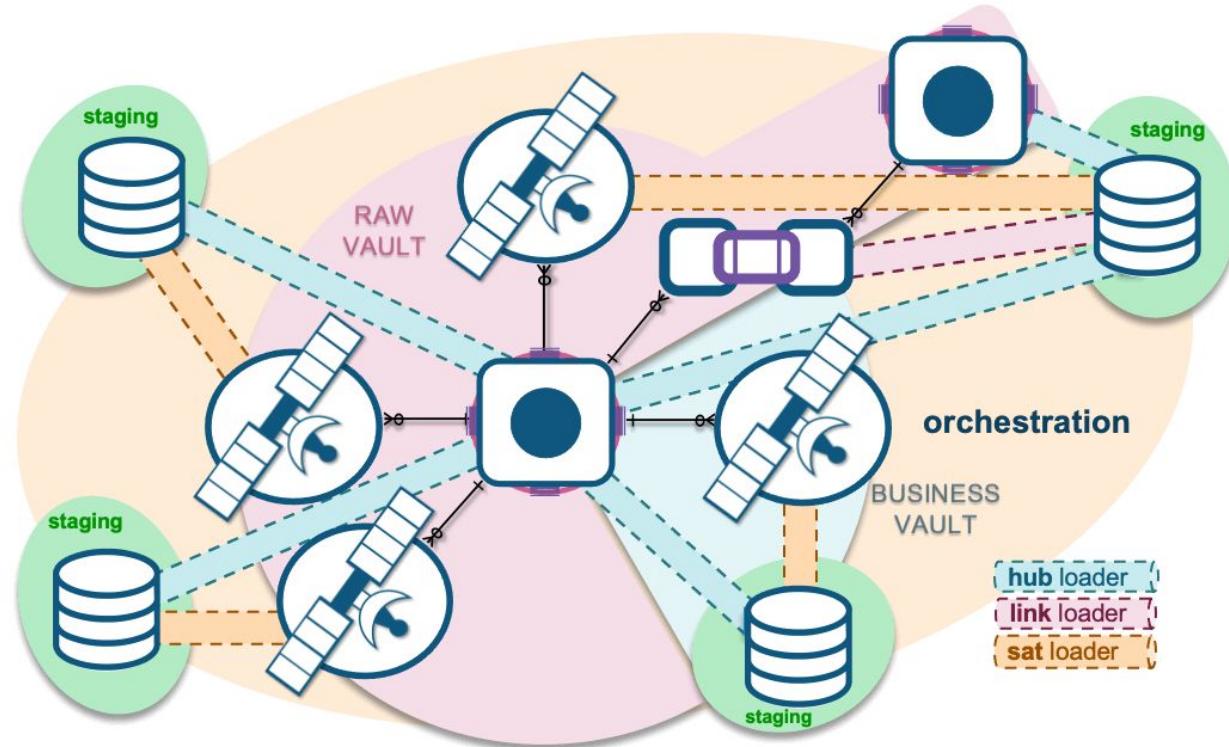
LOAD ANYTIME, ALL DAY



HUB LOCKING



⑥ DATA VAULT PARALLEL LOADING

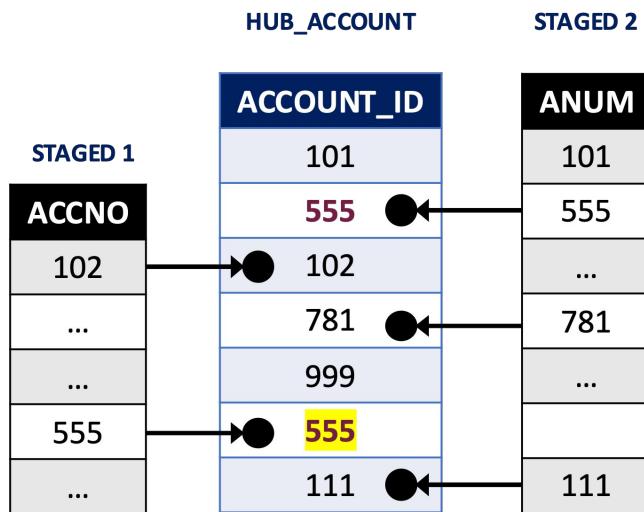


-switch to Parallel Scripts

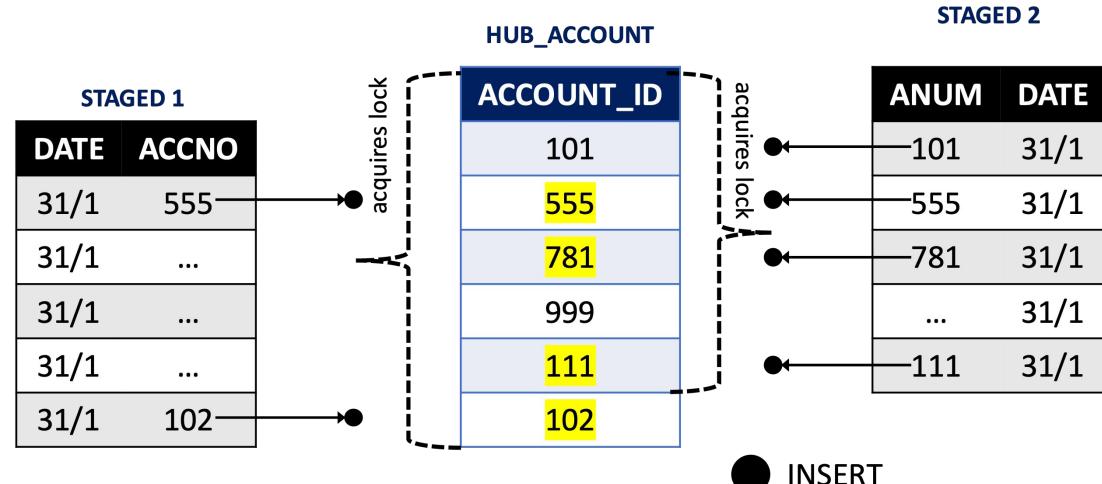


TRANSACTION ISOLATION

INSERT



MERGE



MERGE acquires locks



TRANSACTION ISOLATION

INSERT

HUB_ACCOUNT		STAGED 2	
STAGED 1	ACCOUNT_ID	ANUM	
ACCNO	101	101	
102	555	555	●
...	102	...	
...	781	781	●
555	999	...	
...	555	...	
555	111	111	●
...	111	111	●

MERGE (*optional update*)

HUB_ACCOUNT		STAGED 2	
STAGED 1	ACCOUNT_ID	LAST_SEEN	
DATE	ACCNO	DATE	
31/1	555	31/1	○
31/1	...	31/1	●
31/1	781	31/1	●
31/1	...	15/1	
31/1	111	31/1	●
31/1	102	31/1	●

acquires lock

acquires lock

○ UPDATE
● INSERT

applicable to HUBs and LINKs, SATELLITEs remain INSERT-ONLY



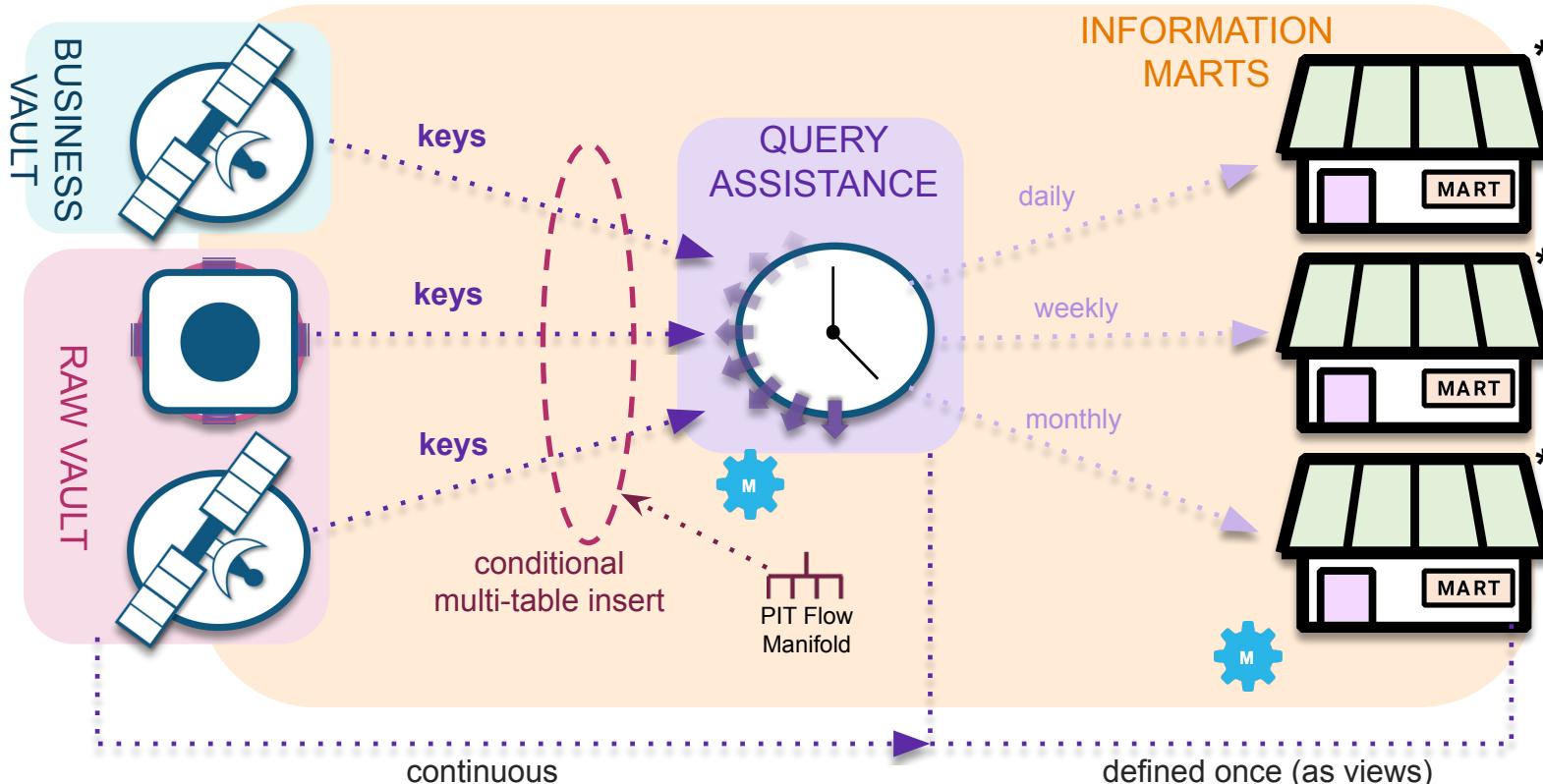
QUERYING DATA VAULT

Point in Time tables



⑦

TEMPLATED CONSUMPTION PIPELINES



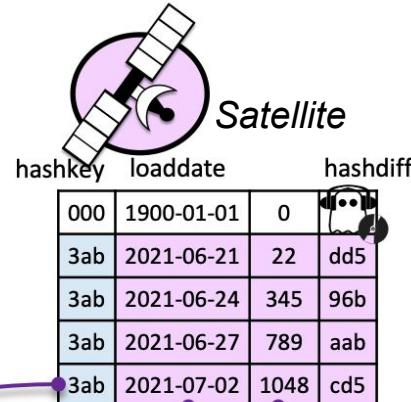
*first query will cache the result; Information Mart is in memory as result_scan
-switch to Step_04_Querying_Information_Marts



PIT & Sequence Number PIT

Traditional PIT

businesskey	hashkey loaddate	hashkey loaddate
snapshotdate	hashkey loaddate	hashkey loaddate
101	2021-06-21	000 1900-01-01
101	2021-06-22	3ab 2021-06-21
101	2021-06-23	3ab 2021-06-22
101	2021-06-24	3ab 2021-06-23
101	2021-06-25	3ab 2021-06-24
101	2021-06-26	3ab 2021-06-25
101	2021-06-27	3ab 2021-06-26
101	2021-06-28	3ab 2021-06-27
101	2021-06-29	3ab 2021-06-27
101	2021-06-30	3ab 2021-06-27
101	2021-07-01	3ab 2021-06-30
101	2021-07-02	3ab 2021-07-02
101	2021-07-03	3ab 2021-07-02
101	2021-07-04	3ab 2021-07-03



Sequence PIT

businesskey	snapshotdate	hashkey loaddate	sequence ids
101	2021-06-21	0	0
101	2021-06-22	22	1
101	2021-06-23	350	2
101	2021-06-24	345	3
101	2021-06-25	1111	4
101	2021-06-26	2048	5
101	2021-06-27	789	6
101	2021-06-28	5021	7
101	2021-06-29	1048	8
101	2021-06-30	7200	9
101	2021-07-01		10
101	2021-07-02		11
101	2021-07-03		12
101	2021-07-04		13

Enable EQUI-JOIN, like dimensions around a fact!



TABLE CLUSTERING

Tablename	Record count	Column count	Size
hub_account	150,000	9	4MB
sat_card_masterfile	17,253,768	613	3GB
sat_card_balancecategories	17,259,059	419	2GB
sat_card_transaction_header	17,257,668	12	875.3MB
sat_bv_account_card_summary	150,001	10	4.4MB

Column(s)	Average overlaps	Average depth
dv_hashkey_hub_account	1452.0	1453.0
dv_hashkey_hub_account, dv_loaddate	1452.0	1453.0
dv_loaddate	0	1.0
dv_sid	6.9897	7.9897

Tablename	Record count	Column count	Size
pit_cardaccount_daily	27,297,309	11	2.7GB
pit_cardaccount_weekly	3,900,000	11	400MB
pit_cardaccount_monthly	900,000	11	92.3MB
pit_cardaccount_current	150,000	11	15.4MB
snopit_cardaccount_daily	27,297,309	7	929MB
snopit_cardaccount_weekly	3,900,000	7	132.3MB
snopit_cardaccount_monthly	900,000	7	30.5MB
snopit_cardaccount_current	150,000	7	5.1MB

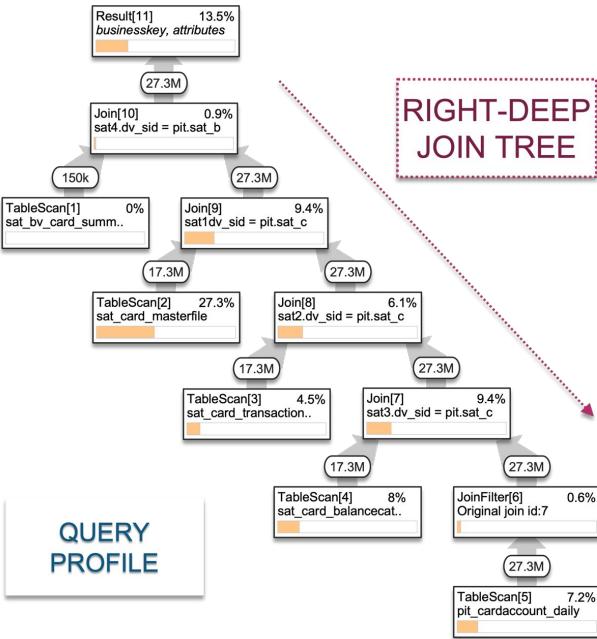
```
create table sat_card_masterfile
(dv_tenantid varchar(20) not null
, dv_hashkey_hub_account binary(20) not null
, dv_loaddate datetime not null
, dv_applieddate datetime not null
, dv_recsource varchar(100) not null
, dv_hashdiff binary(20) not null
, dv_sid int autoincrement(0, 1)
, card_type varchar(1)
, card_balance decimal
, card_status varchar(1)
, credit_limit decimal
, ...)
```

-switch to Step_04_Querying_Information_Marts



STAR JOIN QUERY

Query Planner – for PIT
(acts like a factless fact)

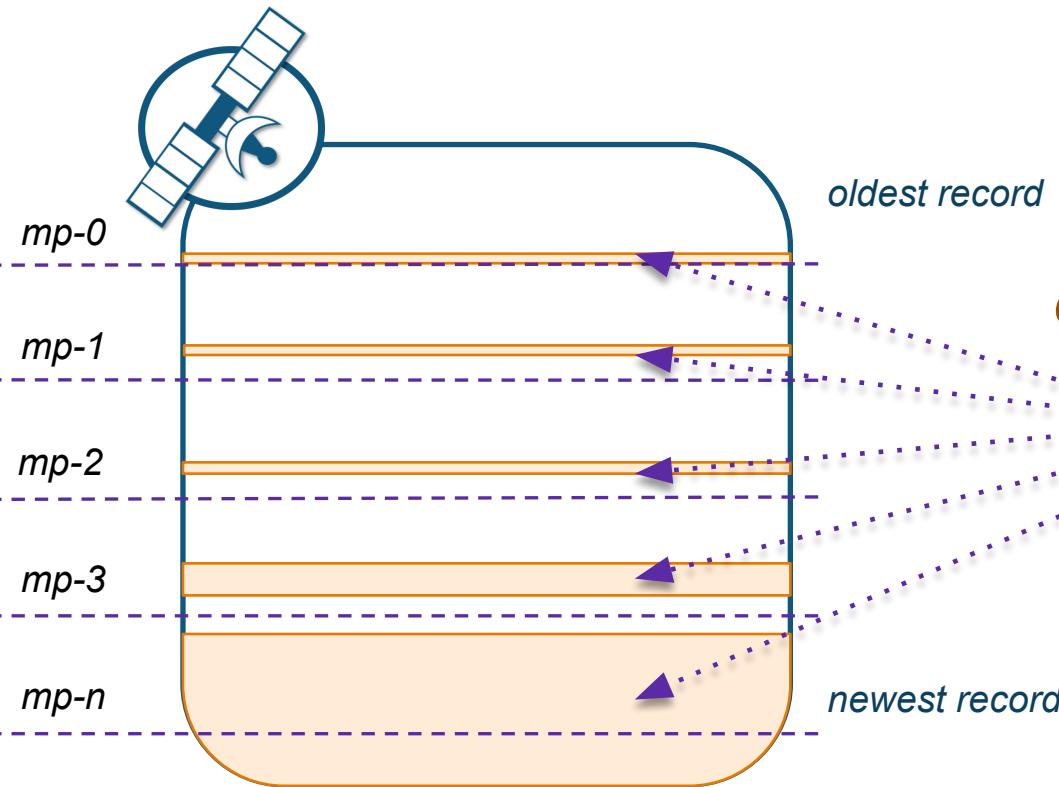


1) No PIT, 2) PIT, 3) Sequence-PIT

More on clustering depth: bit.ly/2SzCAH3

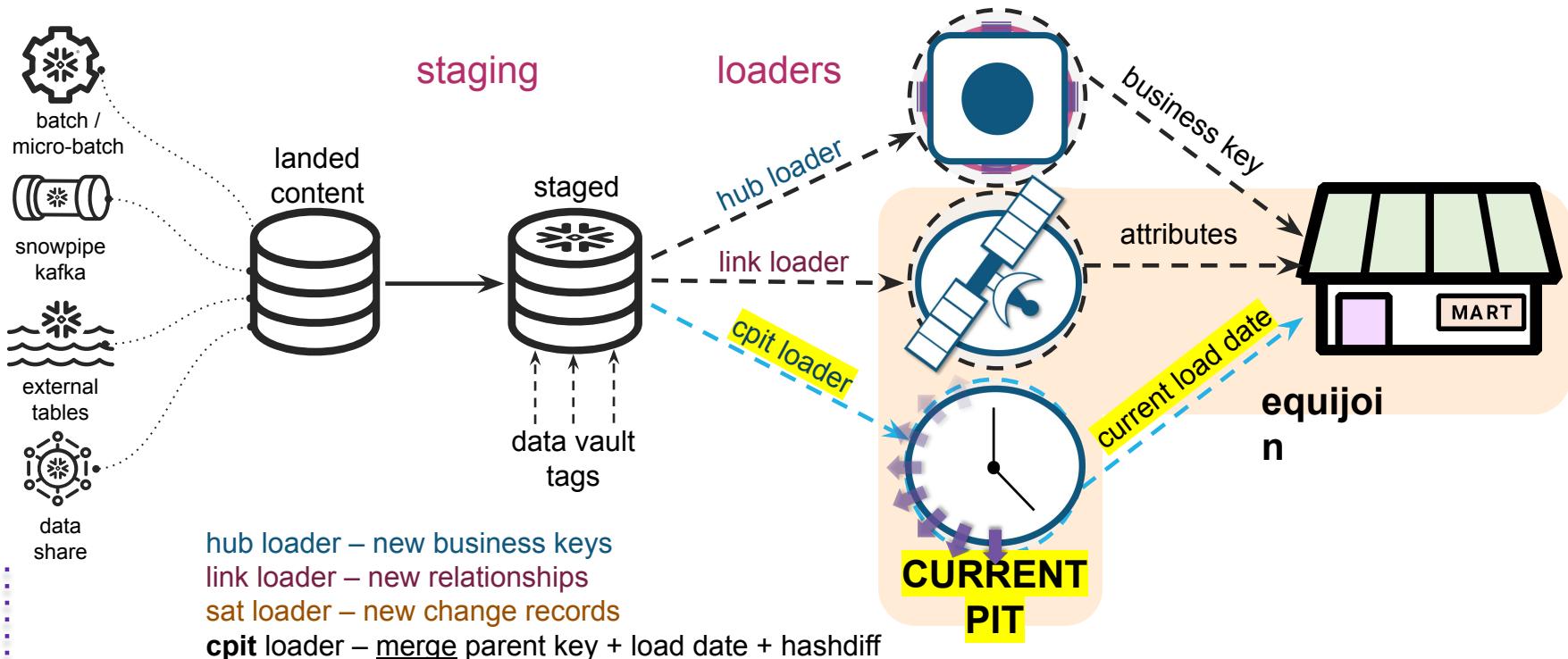


QUERYING BIG SATELLITES

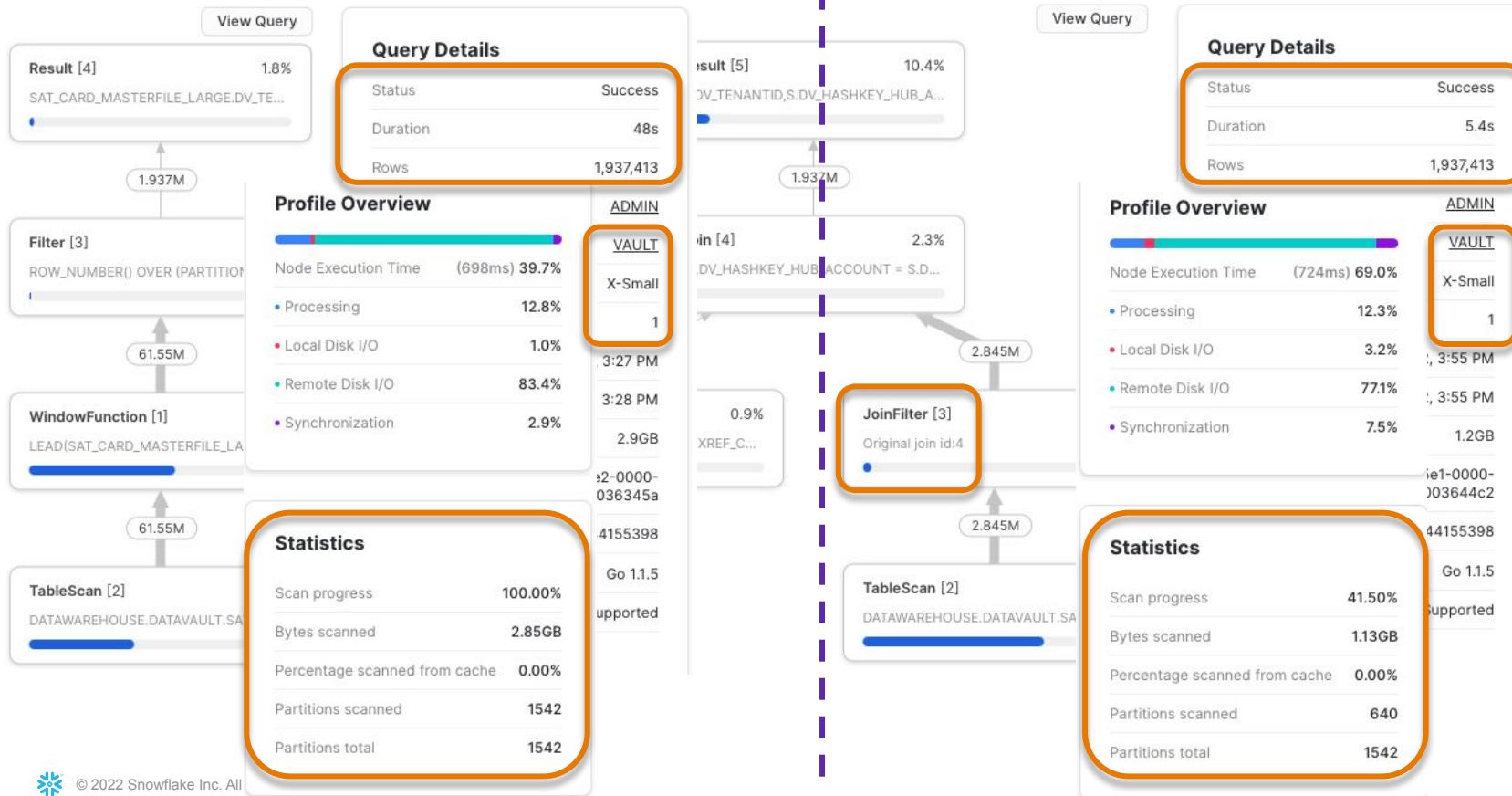


current record (high-date)
for *different parent keys*
may be *scattered* across a
satellite table in *different*
micro-partitions

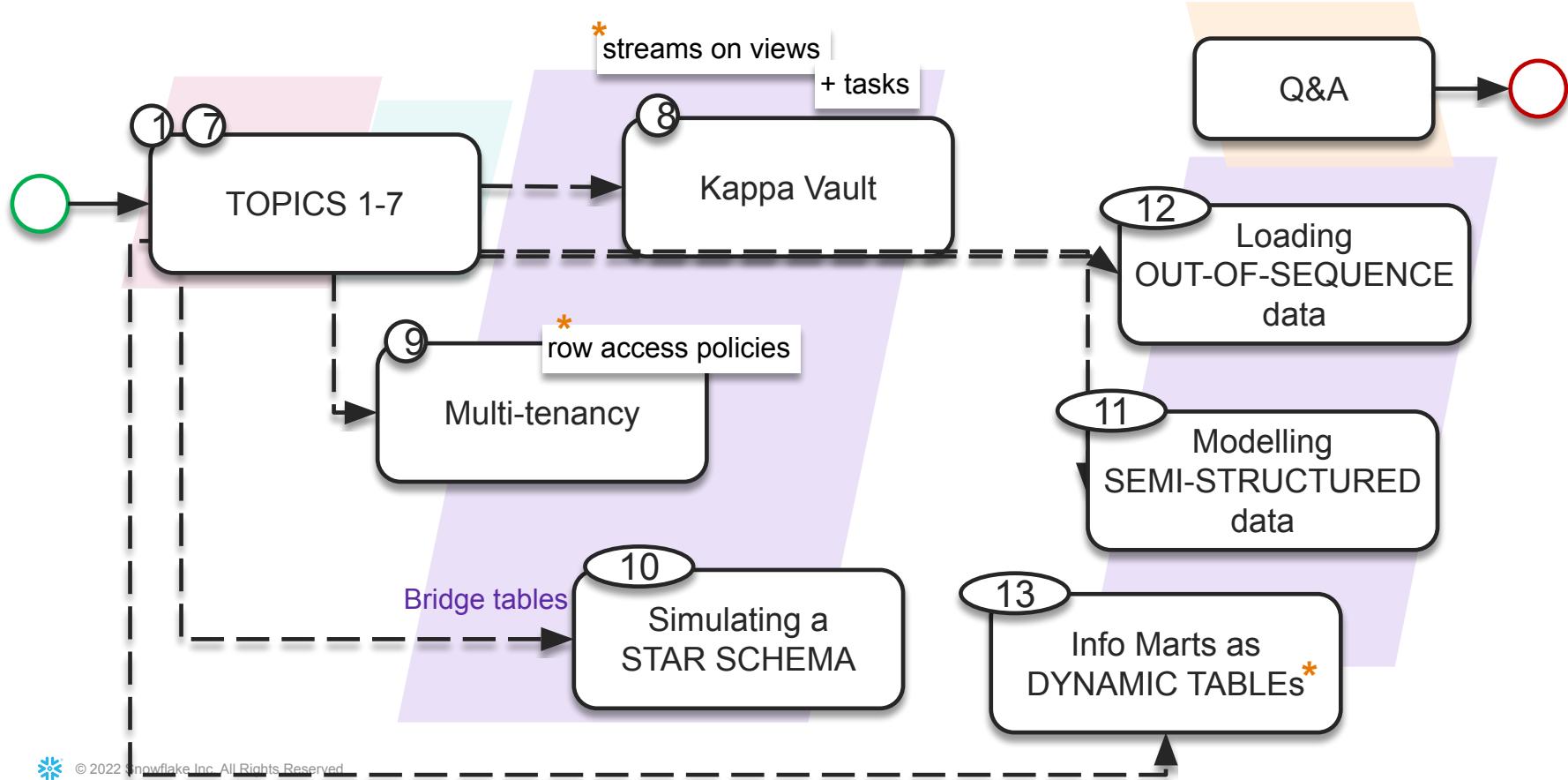
PARALLEL LOADING A PIT



DYNAMIC PRUNING



ADVANCED TOPICS & PATTERNS



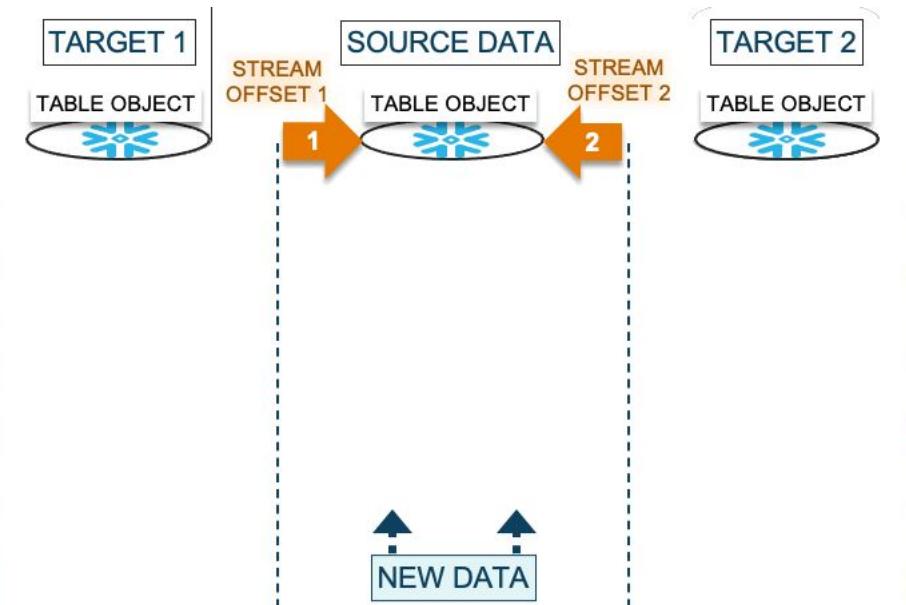
KAPPA VAULT



STREAMS & TASKS

How Streams Work

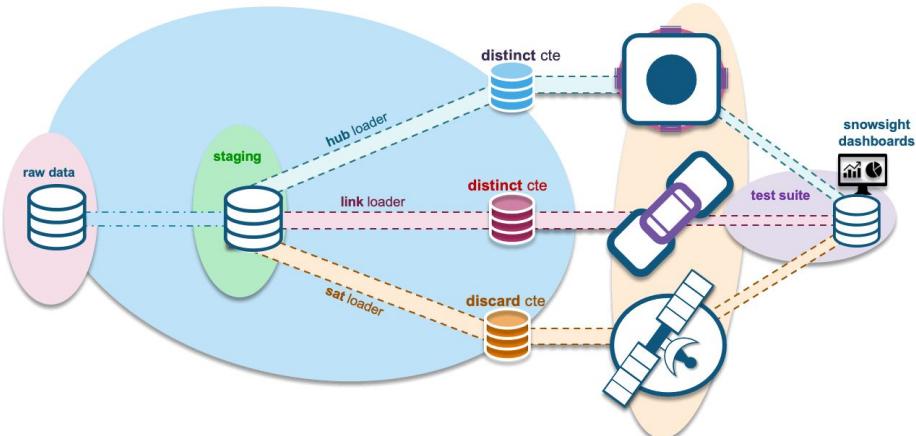
- Tracks **Change Data Capture (CDC)**
- You can have as many streams on an object as you like
- Streams on Tables, Shares, External Tables, Views
- Load & Test on the *same* data requires **Repeatable READ Isolation**



STREAMS ON VIEWS

Tasks firing Streams

⑧



-switch to Step_06_Set_and_Forget

a pattern to delta with *true change* in the staged content

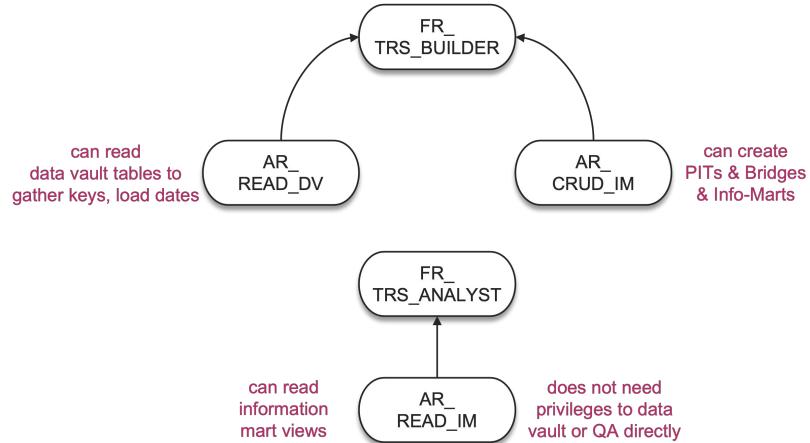
- 1 view with DV metadata columns
- 1:M streams on the view for each data pipeline
- 1x Stream: 1x Task for each table that needs loading and each test executed
- CTE discards non-true change in staging before standard load pattern
- Repeatable-read isolation if you include a test harness



MULTI-TENANCY



RBAC & ENTITLEMENTS



Tenant	Tenant Code	Role	Description
Enterprise	Default	ENGINEER	Custodians of the Vault
Treasury	TRS	TRS_BUILDER	This role can use data vault but only see data the role is entitled to. This role will build Query Assistance tables, i.e., PITs & Bridges
Treasury	TRS	TRS_ANALYST	Has no access to Data Vault directly, but using views defined using PITs & Bridges

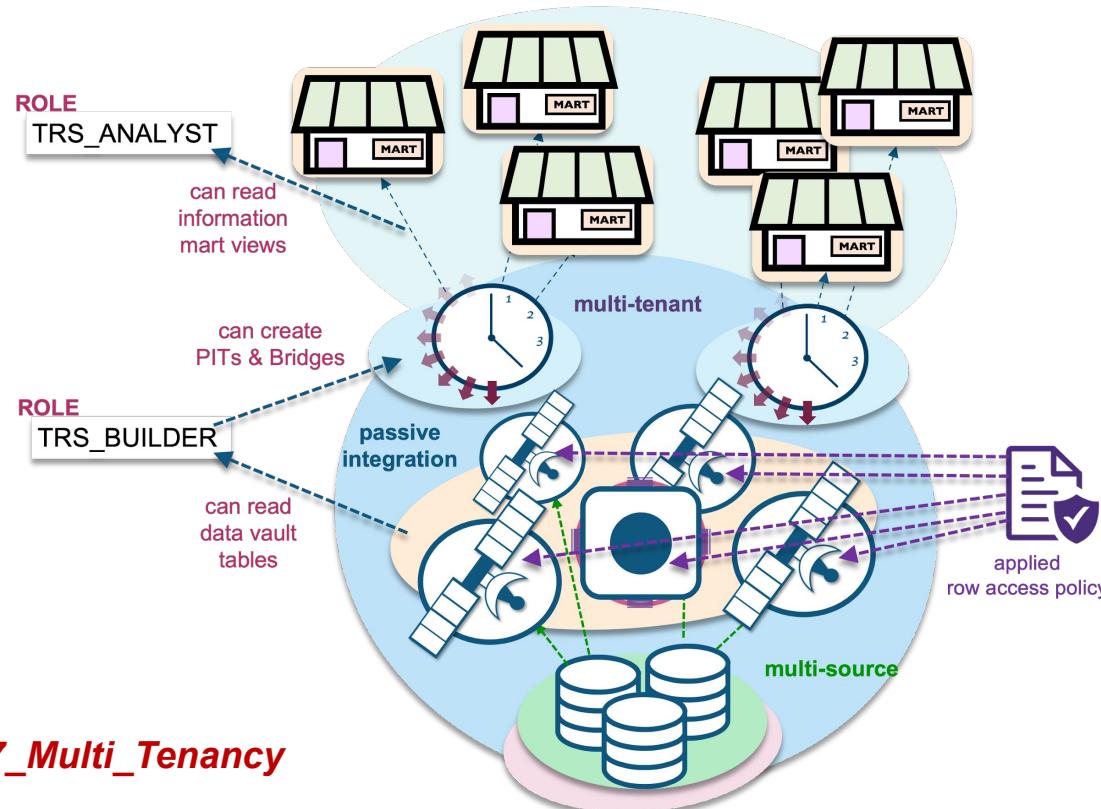
Combine DV multi-tenancy +
RBAC + Row Access Policies
(RAP)

- Every DV table has a tenant code
- Build look entitlements lookup
- Deploy RAP
- Use RAP + Role to Build PITs
- VIEW role does not need explicit access to DV
- Consider INVOKER_ROLE and Memoizable function



ROLE + RAP + DV MULTI-TENANCY

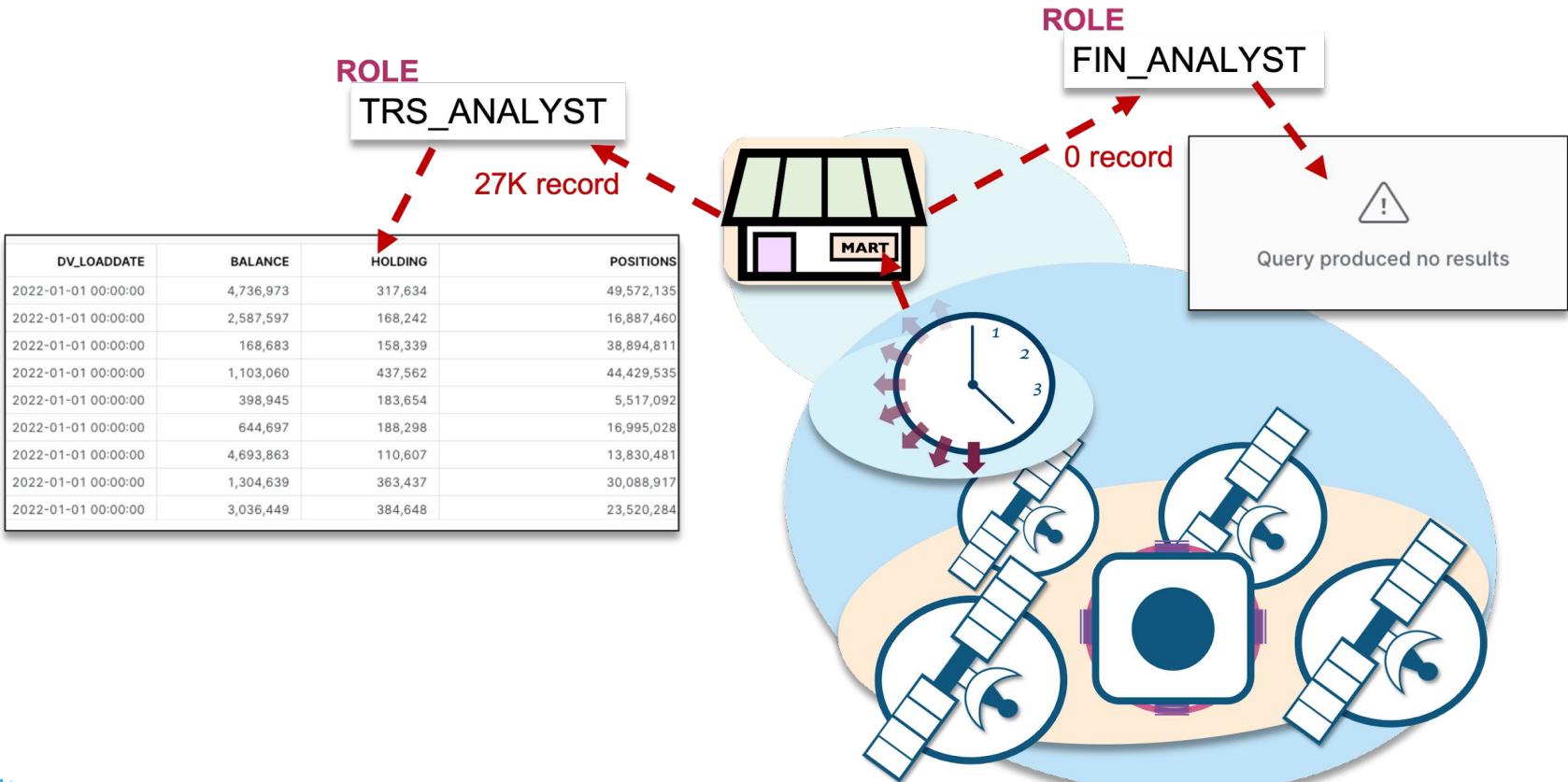
9



-switch to Step_07_Multi_Tenancy



ROLE with and without TRS ACCESS

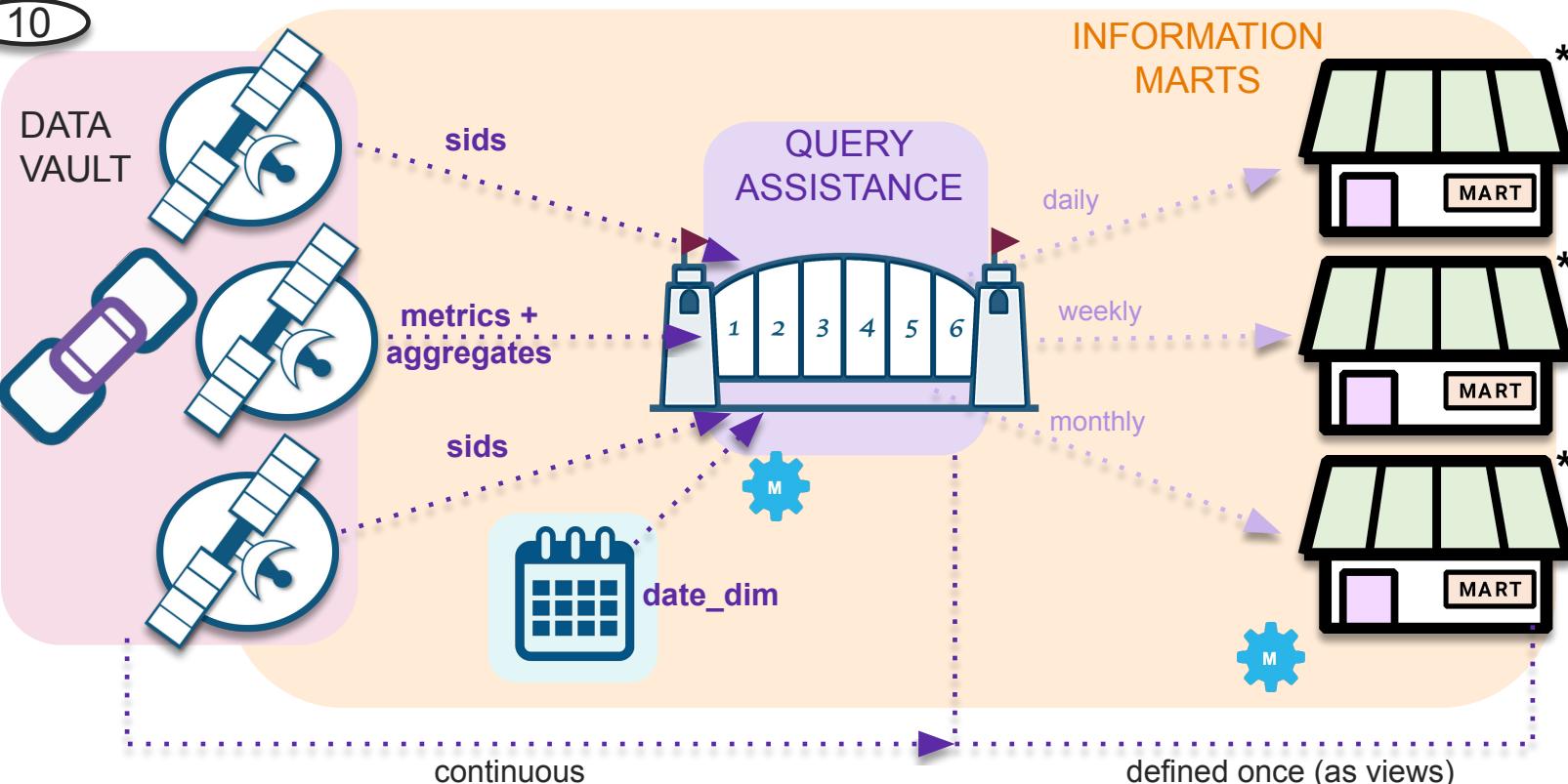


SIMULATING A STAR SCHEMA



PRE-AGGREGATE FACTS

10



*satellites act as dims, bridge tables act as facts with metrics & aggregations
-switch to Step_08_Bridge_and_Metrics*

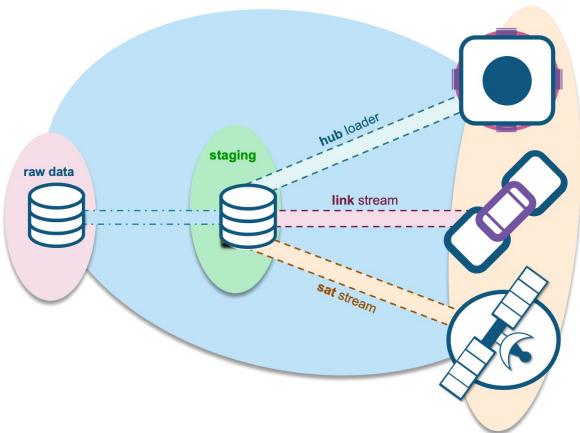


Modelling **SEMI-STRUCTURE** DATA



GUIDELINES / BEST PRACTICES

11



**-switch to
Step_09_SemiStructure_Query_and_Design**

1. Load semi-structure+structured together
2. Persist only key-pairs that you need
3. Consider natural key DV
4. Load DV first or after
5. Streaming use non-historized load
6. Extract bkeys & event timestamps
7. Consider sensitive key-pairs
8. Must schedule test harness



TIME CRIME

(loading OUT-OF-SEQUENCE data)



The Problem: Data Arrives Out of Sequence

STAGED

ID	City	HD	Date
4	Sydney	6D	1-October
4	Sydney	6D	3-October
4	Brisbane	6E	2-October

SAT

ID	City	HD	Date
4	Sydney	6D	1-October
4	Brisbane	6E	2-October

not a **TRUE** change

HUB

ID	Date
4	1-October

This Object's record arrived OUT-OF-SEQUENCE
It differs from previous record, **WE MUST INSERT**

Object's current STATE
is **INCORRECT**



The Solution: COPY correction w/ XTS

STAGED

ID	City	HD	Date
4	Sydney	6D	1-October
4	Sydney	6D	3-October
4	Brisbane	6E	2-October

SAT

ID	City	HD	Date
4	Sydney	6D	1-October
4	Brisbane	6E	2-October
4	Sydney	6D	3-October

enforces a **COPY** as previous record

XTS

HUB

ID	Date
4	1-October

ID	HD	Date
4	6D	1-October
4	6D	3-October
4	6E	2-October

STATE is now **CORRECT**

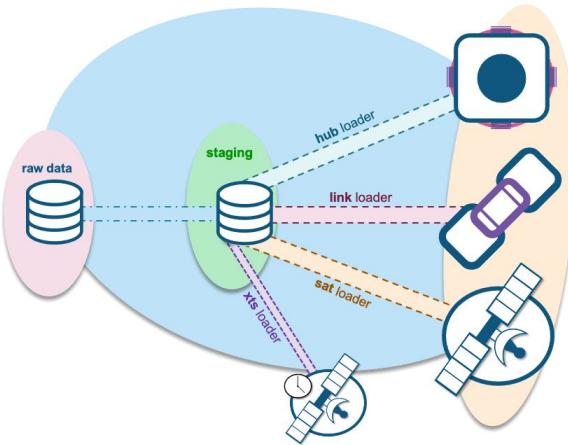
NEW RECORD
DIFFERS from
PREVIOUS record
and **NEXT** record

-switch to Step_11_Solving_Time_Crime_Part_1



GUIDELINES / BEST PRACTICES

12



1. One XTS per Hub or Link or one XTS per Satellite
2. Each XTS has the record_target column denoting which hashdiff it is recording
3. XTS can be populated before or after (dirty load) the current record being loaded, don't let XTS get too dirty

-switch to Step_11_Solving_Time_Crime_Part_1



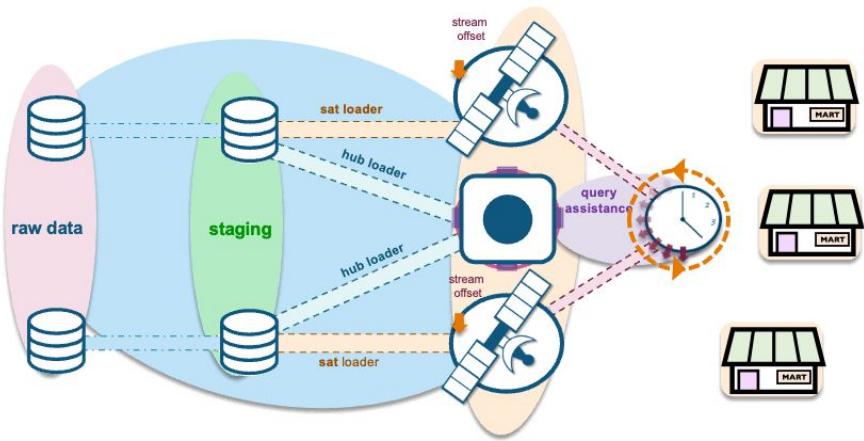
© 2022 Snowflake Inc. All Rights Reserved

Info Marts as DYNAMIC TABLES



GUIDELINES / BEST PRACTICES

13



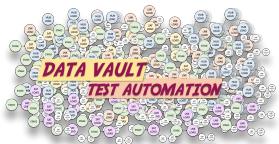
1. Should NOT be used to replace Hub, Link and Satellite Tables
2. Can be used as Materialized Views but beware of *LAG*..
3. Should be considered based on NON-HISTORISED Links and Satellites for *STREAMING* use cases
4. Can VIEWS be used instead?
5. Beware that Dynamic Tables are NOT FREE

-switch to Step_10_Using_Dynamic_Tables_for_BV

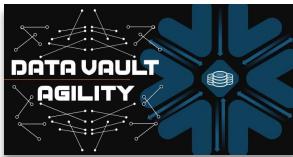


© 2022 Snowflake Inc. All Rights Reserved

DATA VAULT LITERATURE



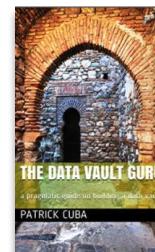
Data Vault Blogs



Data Vault on Snowflake Techniques (transactions & pruning)



Snowflake Data Vault Resource Kit (lots of Snowflake DV content)



The Data Vault Guru (sample code and modelling examples)

Recommended coaching and training, and

join a community of data vault modelers

<https://forum.ukdatavaultusergroup.co.uk/>

Data Vault Accelerator: <https://bit.ly/3yGTzas>
(take a copy as PDF)



BUILD OR BUY - PLAYERS

 WhereScape®

Internal Demo: <https://bit.ly/3s2Qejp>
Webpage: <https://www.wherescape.com/>

 VAULTSPEED

Internal Demo: <https://bit.ly/3FqYheM>
Webpage: <https://vaultspeed.com/>

 AutomateDV
formerly known as dbtvault

Internal Demo: <https://bit.ly/3xCstzL>
Webpage: <https://dbtvault.readthedocs.io/en/latest/>

 erwin®
by Quest
DataJoinery.

Internal Demo: <https://bit.ly/3iyRW95>
Webpage: <https://bit.ly/3qKIEfD>

Internal Demo: <https://bit.ly/2VIMibb>
Webpage: <https://datajoinery.io>

 Datavault
BUILDER

Internal Demo: <https://bit.ly/3LWSCSi>
Webpage: <https://datavault-builder.com/>

 Varigence

Internal Demo: <https://bit.ly/3DoJ93m>
Webpage: <https://www.varigence.com/>

 coalesce®

Internal Demo: <https://bit.ly/3n0oa1B>
Webpage: <https://coalesce.io/>

Head office in NZ
Thick Client

Head office in Belgium
Thin Client

Head office in UK
Thin Client, dbt package

Head office in US
Thick Client, no orchestration

Head office in South Africa
Thin Client

Head office in Germany
Thin Client

Head office in Australia
Thin Client

Head office in US
Thin Client



Q&A



THANK YOU



Any organization that designs a system will inevitably produce a design whose structure is a copy of the organization's communication structure.

-Melvin E. Conway

