# Dokumentasi Data Pipelining Azure

*Dataset : Adventure Work*

## A. Configure Arch ADF, DBricks, Azure

### 1. Create and Config Architecture

#### a. Membuat Resource Group

Grup sumber daya dibuat sebagai wadah untuk semua komponen Azure yang akan digunakan.

**Create a resource group** ...

Basics     Tags     **Review + create**

⟳ Automation Link

**Basics**

| | |
|---|---|
| Subscription | Azure for Students |
| Resource group name | medalion-spark-dbt-rg |
| Region | Indonesia Central |

**Tags**

None

---

**Resource groups** ...
Default Directory (adiiknn9009gmail.onmicrosoft.com)

\+ Create   ⚙ Manage view ∨   ⟳ Refresh   ↓ Export to CSV   ⫘ Open query   |   ⦵ Assign tags      ☰ Group by none ∨

ⓘ You are viewing a new version of Browse experience. Some features may be missing. Click here to access the old experience.

▽ Filter for any field...    Subscription equals **all**    Location equals **all** ✕   \+ Add filter

| ☐ | Name ↑ | | Subscription | Location |
|---|---|---|---|---|
| ☐ | 〔⊛〕 medalion-spark-dbt-rg | ··· | Azure for Students | Indonesia Central |

#### b. Membuat Azure Data Lake Storage Gen2/ADLS Gen2

Dibuat sebagai penyimpanan utama untuk data.

**Create a storage account** ...

⟳ View automation template

**Basics**

| | |
|---|---|
| Subscription | Azure for Students |
| Resource group | medalion-spark-dbt-rg |
| Location | Indonesia Central |
| Storage account name | medalions |
| Primary service | |
| Performance | Standard |
| Replication | Locally-redundant storage (LRS) |

**Advanced**

| | |
|---|---|
| Enable hierarchical namespace | Enabled |
| Enable SFTP | Disabled |
| Enable network file system v3 | Disabled |
| Allow cross-tenant replication | Disabled |
| Access tier | Hot |
| Enable large file shares | Enabled |

**Security**

| | |
|---|---|
| Secure transfer | Enabled |

Your deployment is complete

Deployment name: medalions_1750125627293
Subscription: Azure for Students
Resource group: medalion-spark-dbt-rg

Start time: 6/17/2025, 9:01:36 AM
Correlation ID: 0171dae9-2e4c-41ea-b064-37fb22f0d690

**Deployment details**

| Resource | Type | Status | Operation details |
|---|---|---|---|
| medalions/default | Microsoft.Storage/storageAcco... | OK | Operation details |
| medalions/default | Microsoft.Storage/storageAcco... | OK | Operation details |
| medalions | Microsoft.Storage/storageAcco... | OK | Operation details |

Next steps

c. Membuat Container for Bronze, Silver, Gold level

Container untuk Bronze, Silver, Gold level: Tiga container dibuat di ADLS Gen2 untuk menyimpan data pada berbagai tahap pemrosesan:
   - Bronze: Data mentah dalam format aslinya
   - Silver: Data yang telah dibersihkan dan diubah
   - Gold: Data yang telah dimodelkan untuk analisis



d. Membuat Data Factory

Dibuat sebagai orkestrator utama pipeline.



Data factory
medalion-adf-1

Location
southeastasia

Subscription
11645350-20e4-4203-8ad9-40c75427d44e

Resource group
medalion-spark-dbt-rg

Provisioning state
Succeeded

Resource id
/subscriptions/11645350-20e4-4203-8ad9-40c75427d44e/resourceGroups/medalion-spark-dbt-rg/providers/Microsoft.DataFactory/factories/medalion-adf-1

Managed Identity Object ID
743a0d0f-d37f-42e0-9ff4-7cc6e0239622

Managed Identity Tenant
ce94ab30-05ac-4f14-9336-f865a3152f54

Your deployment is complete

Deployment name : Microsoft.DataFactory-20250617090636
Subscription : Azure for Students
Resource group : medalion-spark-dbt-rg

Start time : 6/17/2025, 9:09:17 AM
Correlation ID : a7404e4f-77a5-40af-87d6-de19dc2fba70

Deployment details

| Resource | Type | Status | Operation details |
|---|---|---|---|
| medalion-adf-1 | Data factory (V2) | OK | Operation details |

e. Membuat SQL Database dan Server

Digunakan untuk menyimpan data sumber.

## f. Membuat Key Vaults

Digunakan untuk menyimpan rahasia dan kredensial secara aman.





2. Linked on Data Factory

Data Factory dikonfigurasi untuk terhubung ke berbagai layanan:

a. Link SQL Database & Server

Link SQL Database & Server: Koneksi ke database sumber.

**New linked service**

Azure SQL Database  Learn more

**Name** *

linkMedalionDB

**Description**

**Connect via integration runtime** *

AutoResolveIntegrationRuntime

**Version**

○ 2.0 (Recommended)  ○ 1.0

⟳ Import from connection string

**Account selection method**

○ From Azure subscription  ○ Enter manually

**Azure subscription**

Azure for Students (11645350-20e4-4203-8ad9-40c75427d44e)

**Server name** *

medalion-db-server

**Database name** *

medalion-db-dev

b. Link Storage Account
Menambahkan linked ke storage account utama

**Edit linked service**

Azure Data Lake Storage Gen2  Learn more

**Name** *

linkDataLakeStoraqeBronze

**Description**

**Connect via integration runtime** *

AutoResolveIntegrationRuntime

**Authentication type**

Account key

**Account selection method**

○ From Azure subscription  ○ Enter manually

**URL** *

https://medalions.dfs.core.windows.net/

[Storage account key]  [Azure Key Vault]

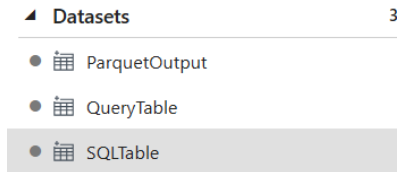**Storage account key** *

••••••••••

**Test connection**

○ To linked service  ○ To file path

**Annotations**

3. Pipeline Activity/Process on Data Factory
Membuat dataset antara lain :

a. Parquetoutput : digunakan untuk menerima data hasil dari for each copy data yang akan distore kan pada storage account

b. Query Table : memuata data set yang akan digunakan berdasarkan query yang digunakan disini menggunakan untuk memanggil semua table schema SalesST dari table_type

c. SQL Table : memuat data set yang akan digunakan sebagai copy dari for each yang akan menerima iterasi sebanyak hasil dari LOOKUP query sebelumnya

d. Create Dataset

◢ Datasets

● ▦ ParquetOutput

● ▦ QueryTable

● ▦ SQLTable

e. Create LOOKUP Activity on Pipeline
   - Menggunakan QueryTable sebagai sumber
   - Query: SELECT * FROM [predation-db-dev].information_schema.tables
     WHERE table_schema = "SelectT" AND table_name = "Base Table"
   - Digunakan untuk mendapatkan daftar tabel yang akan diproses

| Source dataset * | ▦ QueryTable |
|---|---|

✎ Open  + New  👓 Preview data  Learn more ⧉

| First row only | ☐ |
|---|---|
| Use query | ○ Table  ● Query  ○ Stored procedure |
| Query * | SELECT * FROM [medalion-db-dev].information_schema.tables WHERE table_schema = 'SalesLT' AND ▲ ▼  ✎ Edit |
| Query timeout (minutes) ⓘ | 120 |
| Isolation level ⓘ | Select... |
| Partition option ⓘ | ● None  ○ Physical partitions of table ⓘ  ○ Dynami |

ⓘ Please preview data to validate the partition settings.

f. Create FOREACH Activity on Pipeline
   - Memproses setiap table yang ditemukan oleh LOOKUP
   - Menjalankan iterasi secara sequential

General  **Settings**  Activities (1)  User properties

| Sequential | ☐ |
|---|---|
| Batch count ⓘ | |
| Items | @activity('Fetch All Tables').output.va... |

g. Create Copy Data Activity on FOREACH Activity
   Di dalam FOREACH, menyalin data dari sumber ke tujuan, dengan tujuan hasil akan
   disimpan pada container bronze yang digunakan untuk mengmount data raw

h. Debug Run Pipeline

Melakukan debug untuk menmountkan data set yang dilakukan pada pipeline ke dalam container bronze pada storage account dengan tipe data .parquet yang akan digunakan nantinya untuk dilakukan cleaning dan monitoring

5. Create Data Bricks



Your deployment is complete

Deployment name : medalion-spark-dbt-rg_medalion-spark-databricks   Start time    : 6/17/2025, 1:42:20 PM
Subscription        : Azure for Students                                                    Correlation ID : eb683b5d-cd05-4b40-b9b7-f61b3d5b6f74
Resource group    : medalion-spark-dbt-rg

∨ Deployment details

| | Resource | Type | Status | Operation details |
|---|---|---|---|---|
| ✔ | medalion-spark-databricks | Azure Databricks Service | OK | Operation details |

∨ Next steps

Go to resource

6. Menghubungkan databricks dengan storage account untuk masing masing container : bronze, silver, gold :

Menghubungkan Databricks ke ADLS Gen2 untuk akses data di container bronze/silver/gold.

```
dbutils.fs.mount(
    source="wasbs://bronze@medalions.blob.core.windows.net",
    mount_point="/mnt/bronze",

extra_configs={"fs.azure.account.key.medalions.blob.core.windo
ws.net":
dbutils.secrets.get('DataBricksScope','storageAccountKey')}
)

dbutils.fs.mount(
    source="wasbs://silver@medalions.blob.core.windows.net",
    mount_point="/mnt/silver",

extra_configs={"fs.azure.account.key.medalions.blob.core.windo
ws.net":
dbutils.secrets.get('DataBricksScope','storageAccountKey')}
)

dbutils.fs.mount(
    source="wasbs://gold@medalions.blob.core.windows.net",
    mount_point="/mnt/gold",

extra_configs={"fs.azure.account.key.medalions.blob.core.windo
ws.net":
dbutils.secrets.get('DataBricksScope','storageAccountKey')}
)
```
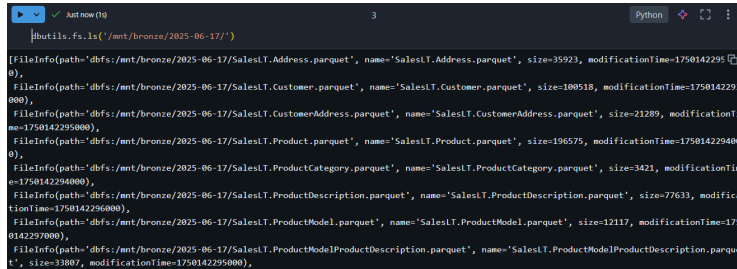
- source: Path ke container Azure Blob Storage (wasbs://[container]@[storage-account].blob.core.windows.net).
- mount_point: Direktori virtual di Databricks (/mnt/bronze).

- extra_configs: Mengambil kunci akses storage dari Azure Key Vault (DataBricksScope).

```
dbutils.fs.ls('/mnt/bronze/2025-06-17/')
```

Output Success :



Penjelasan : Container bronze/silver/gold bisa diakses via /mnt/[level] di Databricks.

```
BaseNotebook

folder_name = dbutils.widgets.get('folder_name')
table_schema = dbutils.widgets.get('table_schema')
table_name = dbutils.widgets.get('table_name')

# Buat database jika belum ada
spark.sql(f"CREATE DATABASE IF NOT EXISTS `{table_schema}`")

# Buat tabel jika belum ada
spark.sql(f"""
        CREATE TABLE IF NOT EXISTS
`{table_schema}`.`{table_name}`
        USING PARQUET
        LOCATION
'/mnt/bronze/{folder_name}/{table_schema}.{table_name}.parquet
'
        """)
```

**Penjelasan :**
Membuat tabel di Databricks dari file Parquet di ADLS.
- Parameter Di-pass dari ADF (e.g., folder_name=2025-06-17).
- LOCATION: Path ke file Parquet di ADLS.
- USING PARQUET: Format file sumber.

Output : Tabel Databricks (e.g., saleslt.address) yang linked ke file Parquet di ADLS.

7. Menambahkan Notebook / Data Bricks pada pipeline FOREACH
   a. Konfigurasi Data Bricks

Penjelasna : menambahkan proses atau service pada pipeline didalam foreach untuk digunakan setelah data berhail di mount pada pipeline ketika copy maka selanjutna akan di mountkan ke dalam bronze container menggunakan databricks yang dihubngkan ke dalam network untuk dilakukan proses mountingnya untuk datanya agar dapat domount pada workload id databrick

b. Setting Params & Notebook Path



**Penjelasan :** Menambahkan konfigurasi parameter yang akan diload pada basenotebook atau db notebook yang dijalankan dimana akan melakukan check juga untuk data yang tidak tertera untuk table_schema mapun table_name nya.

c. Testing



**Penjelasan :**

Hasil dari data yang di mountkan pada copy yang digunakan pada notebook untuk diload pada workload database pada databrick azure berhasil dibuat sama dengan isi dataset pada bronze level

2. DBT (Data Build Tool) Transformation

    a. DBT Configuration

        i. Installation library :
- Dbt-databricks = 1.10.1
- databricks cli = 1.10.3

    b. Initialization dbt



    c. Creating file



| bronze.yml |
| --- |
| ```
version: 2

sources:
  - name: saleslt
    schema: saleslt
    description: This is the adventureworks database loaded
into bronze
    tables:
      - name: address
      - name: customer
      - name: customeraddress
``` |

```
        - name: product
        - name: productcategory
        - name: productdescription
        - name: productmodel
        - name: salesorderdetail
        - name: salesorderheader
```

**Penjelasan :**

Deklarasi Sumber Data: Daftar tabel mentah yang akan digunakan di seluruh proyek DBT. Mengintegrasikan DBT dengan data yang sudah ada di load pada bronze layer, dan membuat metadata yang bisa dilacak DBT docs

d. SQL Silver Level

Membuat snapshot historis data dari sumber ke silver layer.

i. Creating address.sql

address.sql

```
{% snapshot address_snapshot %}

{{
    config(
      file_format = "delta",
      location_root = "/mnt/silver/address",
      target_schema='snapshots',
      invalidate_hard_deletes=True,
      unique_key='AddressID',
      strategy='check',
      check_cols='all'
    )
}}

with source_data as (
    select
        AddressID,
        AddressLine1,
        AddressLine2,
        City,
        StateProvince,
        CountryRegion,
        PostalCode
    from {{ source('saleslt', 'address') }}
)
select *
from source_data

{% endsnapshot %}
```

**Penjelasan :**

- File_format : Delta Lake (format yang dipilih untuk silver).
- Unique_key : Kolom PK untuk tracking perubahan
- strategy='check': Bandingkan nilai kolom (check_cols='all') untuk deteksi perubahan.
- location_root: Penyimpanan fisik file Delta.

**Output**

Tabel snapshot (e.g., snapshots.address_snapshot) dengan kolom tambahan:
- dbt_valid_from: Timestamp awal validitas record.
- dbt_valid_to: Timestamp akhir validitas (NULL untuk data terbaru).

ii. Creating customer.sql

customer.sql

```
{% snapshot customer_snapshot %}

{{
    config(
      file_format = "delta",
      location_root = "/mnt/silver/customer",

      target_schema='snapshots',
      invalidate_hard_deletes=True,
      unique_key='CustomerId',
      strategy='check',
      check_cols='all'
    )
}}

with source_data as (
    select
        CustomerId,
        NameStyle,
        Title,
        FirstName,
        MiddleName,
        LastName,
        Suffix,
        CompanyName,
        SalesPerson,
        EmailAddress,
        Phone,
        PasswordHash,
        PasswordSalt
    from {{ source('saleslt', 'customer') }}
)
select *
from source_data
```

```
{% endsnapshot %}
```

iii.  Creating customeraddress.sql

customeraddress.sql

```
{% snapshot customeraddress_snapshot %}

{{
    config(
      file_format = "delta",
      location_root = "/mnt/silver/customeraddress",

      target_schema='snapshots',
      invalidate_hard_deletes=True,
      unique_key="CustomerId||'-'||AddressId",
      strategy='check',
      check_cols='all'
    )
}}

with source_data as (
    select
        CustomerId,
        AddressId,
        AddressType
    from {{ source('saleslt', 'customeraddress') }}
)
select *
from source_data

{% endsnapshot %}
```

iv.  Creating product.sql

product.sql

```
{% snapshot product_snapshot %}

{{
    config(
      file_format = "delta",
      location_root = "/mnt/silver/product",

      target_schema='snapshots',
      invalidate_hard_deletes=True,
      unique_key='ProductID',
      strategy='check',
      check_cols='all'
    )
```

```
}}

with product_snapshot as (
    SELECT
        ProductID,
        Name,
        ProductNumber,
        Color,
        StandardCost,
        ListPrice,
        Size,
        Weight,
        ProductCategoryID,
        ProductModelID,
        SellStartDate,
        SellEndDate,
        DiscontinuedDate,
        ThumbNailPhoto,
        ThumbnailPhotoFileName
    FROM {{ source('saleslt', 'product') }}
)

select * from product_snapshot

{% endsnapshot %}
```

v.  Creating productmodel.sql

productmodel.sql

```
{% snapshot productmodel_snapshot %}

{{
    config(
      file_format = "delta",
      location_root = "/mnt/silver/productmodel",
      target_schema='snapshots',
      invalidate_hard_deletes=True,
      unique_key='ProductModelID',
      strategy='check',
      check_cols='all'
    )
}}

with product_snapshot as (
    SELECT
        ProductModelID,
        Name,
        CatalogDescription
    FROM {{ source('saleslt', 'productmodel') }}
```

```
)

select * from product_snapshot

{% endsnapshot %}
```

vi.  Creating salesorderdetail.sql

| salesorderdetail.sql |
| --- |

```
{% snapshot salesorderdetail_snapshot %}

{{
    config(
      file_format = "delta",
      location_root = "/mnt/silver/salesorderdetail",
      target_schema='snapshots',
      invalidate_hard_deletes=True,
      unique_key='SalesOrderDetailID',
      strategy='check',
      check_cols='all'
    )
}}

with salesorderdetail_snapshot as (
    SELECT
        SalesOrderID,
        SalesOrderDetailID,
        OrderQty,
        ProductID,
        UnitPrice,
        UnitPriceDiscount,
        LineTotal
    FROM {{ source('saleslt', 'salesorderdetail') }}
)

select * from salesorderdetail_snapshot

{% endsnapshot %}
```

vii.  Creating salesorderheader.sql

| salesorderheader.sql |
| --- |

```
{% snapshot salesorderheader_snapshot %}

{{
    config(
      file_format = "delta",
      location_root = "/mnt/silver/salesorderheader",
```

```
        target_schema='snapshots',
        invalidate_hard_deletes=True,
        unique_key='SalesOrderID',
        strategy='check',
        check_cols='all'
    )
}}

with salesorderheader_snapshot as (
    SELECT
        SalesOrderID,
        RevisionNumber,
        OrderDate,
        DueDate,
        ShipDate,
        Status,
        OnlineOrderFlag,
        SalesOrderNumber,
        PurchaseOrderNumber,
        AccountNumber,
        CustomerID,
        ShipToAddressID,
        BillToAddressID,
        ShipMethod,
        CreditCardApprovalCode,
        SubTotal,
        TaxAmt,
        Freight,
        TotalDue,
        Comment
    FROM {{ source('saleslt', 'salesorderheader') }}
)

select * from salesorderheader_snapshot

{% endsnapshot %}
```

e. GOLDEN SQL level (marts)
   Transformasi data silver ke dimensional model (star schema).

   i.  Customer
       Creating dim_customer.sql

       | dim_customer.sql |
       | --- |
       | ```
{{
    config(
        materialized = "table",
``` |

```
        file_format = "delta",
        location_root = "/mnt/gold/customers"
    )
}}

with address_snapshot as (
    select
        AddressID,
        AddressLine1,
        AddressLine2,
        City,
        StateProvince,
        CountryRegion,
        PostalCode
    from {{ ref('address_snapshot') }} where dbt_valid_to
is null
)

, customeraddress_snapshot as (
    select
        CustomerId,
        AddressId,
        AddressType
    from {{ref('customeraddress_snapshot')}} where
dbt_valid_to is null
)

, customer_snapshot as (
    select
        CustomerId,
        concat(ifnull(FirstName,' '),'
',ifnull(MiddleName,' '),' ',ifnull(LastName,' ')) as
FullName
    from {{ref('customer_snapshot')}} where dbt_valid_to
is null
)

, transformed as (
    select
    row_number() over (order by
customer_snapshot.customerid) as customer_sk, --
auto-incremental surrogate key
    customer_snapshot.CustomerId,
    customer_snapshot.fullname,
    customeraddress_snapshot.AddressID,
    customeraddress_snapshot.AddressType,
    address_snapshot.AddressLine1,
    address_snapshot.City,
    address_snapshot.StateProvince,
    address_snapshot.CountryRegion,
    address_snapshot.PostalCode
```

```
    from customer_snapshot
    inner join customeraddress_snapshot on
customer_snapshot.CustomerId =
customeraddress_snapshot.CustomerId
    inner join address_snapshot on
customeraddress_snapshot.AddressID =
address_snapshot.AddressID
)
select *
from transformed
```

**Penjelasan :**
- materialized = "table": Buat tabel fisik (bukan view).
- location_root: Path penyimpanan di gold layer.
- ref(): Referensi ke model DBT lain (e.g., address_snapshot).
- row_number() as customer_sk: Generate surrogate key.

Output:
- Tabel dimensi (dim_customers) di /mnt/gold/customers dengan schema:

Creating dim_customer.yml

| dim_customer.yml |
|---|

```
version: 2

models:
  - name: dim_customers
    columns:
      - name: customer_sk
        description: The surrogate key of the customer
        tests:
          - unique
          - not_null
      - name: customerid
        description: The natural key of the customer
        tests:
          - not_null
      - name: fullname
        description: The customer name. Adopted as
customer_fullname when person name is not null.
      - name: AddressId
        tests:
          - not_null
      - name: AddressType
      - name: AddressLine1
      - name: City
      - name: StateProvince
      - name: CountryRegion
```

```
                    -   name: PostalCode
```

ii.    Product
       Creating dim_product.sql

```
dim_product.sql

{{
    config(
        materialized = "table",
        file_format = "delta",
        location_root = "/mnt/gold/products"
    )
}}

with product_snapshot as (
    select
        productId,
        name,
        standardCost,
        listPrice,
        size,
        weight,
        productcategoryid,
        productmodelid,
        sellstartdate,
        sellenddate,
        discontinueddate
    from {{ ref("product_snapshot") }}
    where dbt_valid_to is null
),

product_model_snapshot as (
    select
        productmodelid,
        name,
        CatalogDescription,
        row_number() over (order by name) as model_id
    from {{ ref("productmodel_snapshot") }}
    where dbt_valid_to is null
),


transformed as (
    select
        row_number() over (order by p.productId) as
product_sk,
        p.name as product_name,
        p.standardCost,
        p.listPrice,
```

```
        p.size,
        p.weight,
        pm.name as model,
        pm.CatalogDescription as description,
        p.sellstartdate,
        p.sellenddate,
        p.discontinueddate
    from product_snapshot p
    left join product_model_snapshot pm on
p.productmodelid = pm.productmodelid
)

select * from transformed
```

Creating dim_customer.yml

dim_product.yml

```
version: 2

models:
  - name: dim_products
    columns:
      - name: product_sk
        description: The surrogate key of the product
        tests:
          - unique
          - not_null
      - name: product_name
        description: The name of the product
        tests:
          - not_null
      - name: standard_cost
        description: The standard cost of the product
      - name: list_price
        description: The list price of the product
      - name: size
        description: The size of the product
      - name: weight
        description: The weight of the product
      - name: category
        description: The category of the product
      - name: model
        description: The model of the product
      - name: description
        description: The description of the product
      - name: sellstartdate
        description: The date when the product is
available for sale
        tests:
```

```
          - not_null
        - name: sellenddate
          description: The date when the product is no
longer available for sale
        - name: discontinueddate
          description: The date when the product is
discontinued
```

iii.  Sales
Creating dim_customer.sql

sales.sql

```
{{
    config(
        materialized = "table",
        file_format = "delta",
        location_root = "/mnt/gold/sales"
    )
}}

with salesorderdetail_snapshot as (
    SELECT
        SalesOrderID,
        SalesOrderDetailID,
        OrderQty,
        ProductID,
        UnitPrice,
        UnitPriceDiscount,
        LineTotal
    FROM {{ ref("salesorderdetail_snapshot") }}
),

product_snapshot as (
    SELECT
        ProductID,
        Name,
        ProductNumber,
        Color,
        StandardCost,
        ListPrice,
        Size,
        Weight,
        SellStartDate,
        SellEndDate,
        DiscontinuedDate,
        ThumbNailPhoto,
        ThumbnailPhotoFileName
    FROM {{ source('saleslt', 'product') }}
),
```

```
saleorderheader_snapshot as (
    SELECT
        SalesOrderID,
        RevisionNumber,
        OrderDate,
        DueDate,
        ShipDate,
        Status,
        OnlineOrderFlag,
        SalesOrderNumber,
        PurchaseOrderNumber,
        AccountNumber,
        CustomerID,
        ShipToAddressID,
        BillToAddressID,
        ShipMethod,
        CreditCardApprovalCode,
        SubTotal,
        TaxAmt,
        Freight,
        TotalDue,
        Comment,
        row_number() over (partition by SalesOrderID
order by SalesOrderID) as row_num
    FROM {{ source('saleslt', 'salesorderheader') }}
),

transformed as (
    select
        sod.SalesOrderID,
        sod.SalesOrderDetailID,
        sod.OrderQty,
        sod.ProductID,
        sod.UnitPrice,
        sod.UnitPriceDiscount,
        sod.LineTotal,
        p.Name,
        p.ProductNumber,
        p.Color,
        p.StandardCost,
        p.ListPrice,
        p.Size,
        p.Weight,
        p.SellStartDate,
        p.SellEndDate,
        p.DiscontinuedDate,
        p.ThumbNailPhoto,
        p.ThumbnailPhotoFileName,
        soh.RevisionNumber,
        soh.OrderDate,
```

```
        soh.DueDate,
        soh.ShipDate,
        soh.Status,
        soh.OnlineOrderFlag,
        soh.SalesOrderNumber,
        soh.PurchaseOrderNumber,
        soh.AccountNumber,
        soh.CustomerID,
        soh.ShipToAddressID,
        soh.BillToAddressID,
        soh.ShipMethod,
        soh.CreditCardApprovalCode,
        soh.SubTotal,
        soh.TaxAmt,
        soh.Freight,
        soh.TotalDue,
        soh.Comment
    from salesorderdetail_snapshot sod
    left join product_snapshot p on sod.ProductID =
p.ProductID
    left join saleorderheader_snapshot soh on
sod.SalesOrderID = soh.SalesOrderID
)

select * from transformed
```

Creating dim_customer.yml

**sales.yml**

```
version: 2

models:
  - name: dim_sales
    description: This is the fact table for sales
    columns:
      - name: saleOrderID
        description: The surrogate key of the sale order
        tests:
          - unique
          - not_null
      - name: saleOrderDetailID
        description: The surrogate key of the sale order
detail
        tests:
          - unique
          - not_null
      - name: orderQty
        description: The quantity of the order
        tests:
```

```
          - not_null
      - name: productID
        description: The surrogate key of the product
        tests:
          - not_null
      - name: unitPrice
        description: The unit price of the product
        tests:
          - not_null
      - name: unitPriceDiscount
        description: The unit price discount of the
product
      - name: lineTotal
        description: The line total of the product
        tests:
          - not_null
      - name: name
        description: The name of the product
        tests:
          - not_null
      - name: productNumber
        description: The product number of the product
        tests:
          - not_null
      - name: color
        description: The color of the product
      - name: standardCost
        description: The standard cost of the product
        tests:
          - not_null
      - name: listPrice
        description: The list price of the product
        tests:
          - not_null
      - name: size
        description: The size of the product
      - name: weight
        description: The weight of the product
      - name: sellStartDate
        description: The date when the product is
available for sale
        tests:
          - not_null
      - name: sellEndDate
        description: The date when the product is no
longer available for sale
      - name: discontinuedDate
        description: The date when the product is
discontinued
      - name: thumbNailPhoto
        description: The thumbnail photo of the product
```

```yaml
      - name: thumbnailPhotoFileName
        description: The thumbnail photo file name of the
product
      - name: revisionNumber
        description: The revision number of the sale
order
      - name: orderDate
        description: The order date of the sale order
        tests:
          - not_null
      - name: dueDate
        description: The due date of the sale order
      - name: shipDate
        description: The ship date of the sale order
      - name: status
        description: The status of the sale order
      - name: onlineOrderFlag
        description: The online order flag of the sale
order
      - name: salesOrderNumber
        description: The sales order number of the sale
order
      - name: purchaseOrderNumber
        description: The purchase order number of the
sale order
      - name: accountNumber
        description: The account number of the sale order
      - name: customerID
        description: The surrogate key of the customer
        tests:
          - not_null
      - name: shipToAddressID
        description: The surrogate key of the ship to
address
      - name: billToAddressID
        description: The surrogate key of the bill to
address
      - name: shipMethod
        description: The ship method of the sale order
      - name: creditCardApprovalCode
        description: The credit card approval code of the
sale order
      - name: subTotal
        description: The sub total of the sale order
        tests:
          - not_null
      - name: taxAmt
        description: The tax amount of the sale order
        tests:
          - not_null
      - name: freight
```

```
          description: The freight of the sale order
          tests:
            - not_null
      - name: totalDue
          description: The total due of the sale order
          tests:
            - not_null
      - name: comment
          description: The comment of the sale order
```

9. Data at silver level result

# dbt

## Tables and Views

- hive_metastore
  - saleslt
    - address
    - customer
    - customeraddress
    - dim_customer
    - dim_product
    - product
    - productcategory
    - productdescription
    - productmodel
    - sales
    - salesorderdetail
    - salesorderheader
  - snapshots
    - address_snapshot
    - customer_snapshot
    - customeraddress_snapshot
    - product_snapshot
    - productmodel_snapshot
    - salesorderdetail_snapshot
    - salesorderheader_snapshot

## customer_snapshot  snapshot

Details    Description    Columns    Referenced By    Depends On    SQL

### Details

| TAGS | OWNER | TYPE | PACKAGE | LANGUAGE | RELATION |
|------|-------|------|---------|----------|----------|
| untagged | (Bearer token)adiiknn9009 | table | medalion_dbt_spark | sql | hive_metastore.sn... |

### Description

This snapshot is not currently documented

### Columns

| COLUMN | TYPE | DESCRIPTION | CONSTRA... |
|--------|------|-------------|------------|
| CustomerId | integer | | |
| NameStyle | boolean | | |
| Title | string | | |
| FirstName | string | | |

## Lineage Graph

- saleslt.customer
- customer_snapshot
- dim_customer

## 10. Data at gold level result