# FUNDAMENTALS OF DATA ENGINEERING.
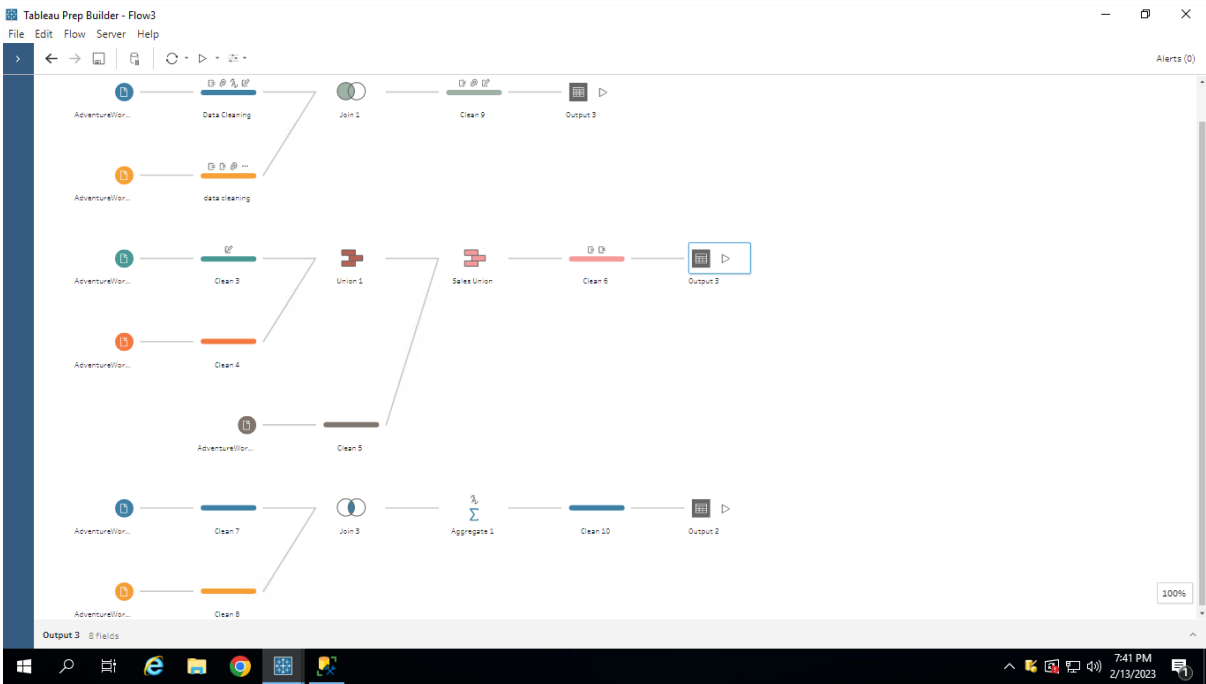
MAHIM DHUNGEL

# Overview: ETL Processes in Tableau

# Phase 1

1. Connect to Customer file



2. Add a "Clean Step

3. Rename Field LastNa to LastName



4. Group Education and Replace by College Degree

5. Change marital status M and S to "Married" and "Single"

6. Clean Prefix MrR to MR



7. Remove Numbers from FirstName

8. Remove punctuations from Occupation

Abc

**Occupation** 5

Clerical
Management
Manual
Professional
Skilled Manual

9. Split email address and remove the first letters prior to @

Abc

**EmailAddress - Split 1** 18K

aaron10
aaron11
aaron12
aaron13
aaron14
aaron15
aaron16
aaron17
aaron18
aaron19
aaron20
aaron21

Abc

**EmailAddress - Split 2** 1

adventure-works.com

# Phase 2

10. Add Adventure Works Customer New file



11. Replace Social Media Account Nulls with NoSocialMedia

## 12. Split Social Media Fields



## 13. Load Data toDimCustomer Tables

## Phase 3

- ·       Add AdventureWorks_Sales_2015
- ·       Add AdventureWorks_Sales_2016
- ·       Add AdventureWorks_Sales_2017

14. Clean OrderQuantity from "Quantite de Ventes" to "OrderQuantity"

15. Union all three Files



16. Add a calculation to extract "Year" from "Order Date"

17. Remove Table Names



18. Load Data into FactSales table

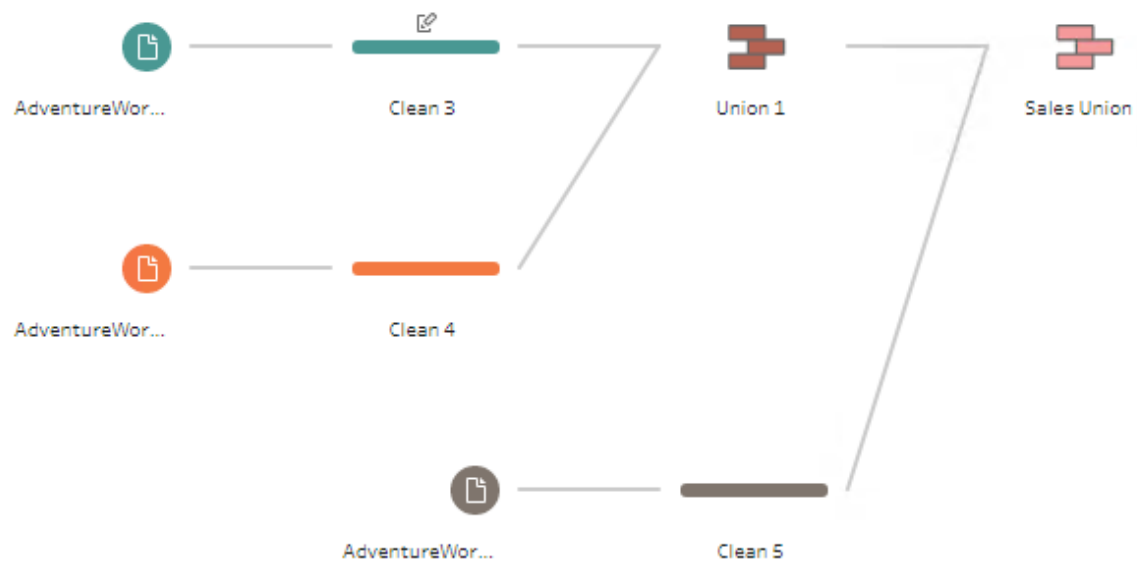| | OrderDate | StockDate | OrderNumber | ProductKey | CustomerKey | TerritoryKey | OrderLineItem | OrderQuantity |
|---|---|---|---|---|---|---|---|---|
| 1 | 2017 | 2003-12-13 | SO61285 | 529 | 23791 | 1 | 2 | 2 |
| 2 | 2017 | 2003-09-24 | SO61285 | 214 | 23791 | 1 | 3 | 1 |
| 3 | 2017 | 2003-09-04 | SO61285 | 540 | 23791 | 1 | 1 | 1 |
| 4 | 2017 | 2003-09-28 | SO61301 | 529 | 16747 | 1 | 2 | 2 |
| 5 | 2017 | 2003-10-21 | SO61301 | 377 | 16747 | 1 | 1 | 1 |
| 6 | 2017 | 2003-10-23 | SO61301 | 540 | 16747 | 1 | 3 | 1 |
| 7 | 2017 | 2003-09-04 | SO61269 | 215 | 11792 | 4 | 1 | 1 |
| 8 | 2017 | 2003-10-21 | SO61269 | 229 | 11792 | 4 | 2 | 1 |
| 9 | 2017 | 2003-10-24 | SO61286 | 528 | 11530 | 6 | 2 | 2 |
| 10 | 2017 | 2003-09-27 | SO61286 | 536 | 11530 | 6 | 1 | 2 |
| 11 | 2017 | 2003-10-23 | SO61298 | 530 | 18155 | 10 | 1 | 2 |

Query executed successfully.                                    RAS-RDSH-04\SQLEXPRESS (15....  | CLARKU

Ln 1            Col 1            Ch 1            INS

## Phase 4

19. Add Returns table and Add Product Table

AdventureWor...        Clean 8

AdventureWor...        Clean 9

20. Inner Join on Product key

AdventureWor...        Clean 8        Join 3

AdventureWor...        Clean 9

21. Add aggregation

AdventureWor...        Clean 8        Join 3        Aggregate 1

AdventureWor...        Clean 9

22. Sum of Prices for Product Returns



23. Create a new table with the following columns:
· Product Key
· Territory Key

24. Create a new table in the database called DMReturns

| | TerritoryKey | ProductKey | ProductCost |
|---|---|---|---|
| 1 | 7 | 377 | 2641.3676 |
| 2 | 8 | 350 | 1898.0944 |
| 3 | 8 | 581 | 2165.02 |
| 4 | 9 | 592 | 308.2179 |
| 5 | 4 | 575 | 2963.8758 |
| 6 | 4 | 482 | 10.0869 |
| 7 | 10 | 576 | 1481.9379 |
| 8 | 1 | 475 | 26.1763 |
| 9 | 9 | 215 | 72.1668 |
| 10 | 4 | 214 | 209.3808 |
| 11 | 10 | 563 | 1481.9379 |

Query executed successfully.