

## 2) Anticipated Datasets & Proposed Domain Model

---

### 2.1 Design Principles

- **Keep grains clean:** **f\_order** is at an order-line grain, while **shipping** is at a customer grain. I avoid mixing these grains in the same table.
  - **Conform dimensions:** I reuse customer, country, and product dimensions across fact tables for consistent joins and filters.
  - **Debias joins:** Because **shipping** is at the customer level, I prevent data fan-out by collapsing it to a single row per customer whenever a report needs shipping status.
  - **Data quality built-in:** Data quality checks, such as name hygiene (digit/symbol mapping, stripping, and trimming), age guardrails, and referential integrity checks, are applied as part of the data load process.
- 

### 2.2 Core Warehouse Entities

#### Fact: f\_order (order-line grain)

- **Grain:** 1 row = 1 order line from the source data.
- **Columns:**
  - order\_id (PK, from source)
  - customer\_id (FK → d\_customer)
  - product\_id (FK → d\_product)
  - country\_id (FK → d\_country) ← Stamped from the customer data at load time
  - quantity (INT, default 1)
  - amount (DECIMAL(18,2), amount > 0)
- **Notes:** This table does not include shipping status, as it's a customer-grain attribute. This design keeps the fact table clean and prevents double-counting.

#### Dim: d\_customer (customer grain)

- **Columns:** customer\_id (PK), first\_name, last\_name, age, country\_id (FK → d\_country)
- **Data Quality Rules Applied at Load:**
  - Map digits/symbols to letters (e.g., 0→o, 1→i, 3→e, 4→a, 5→s, 7→t, @→a, !→i).
  - Strip disallowed punctuation, then trim and collapse spaces.
  - Enforce age to be BETWEEN 10 AND 100.

#### Dim: d\_country (country grain)

- **Columns:** country\_id (PK), country\_name (UNIQUE)
- **Population:** Populated with distinct countries from the customers source data.

### Dim: d\_product (product grain)

- **Columns:** product\_id (PK), product\_name (from Item in source data)
- **Future-proof Fields (nullable):** category, unit\_price

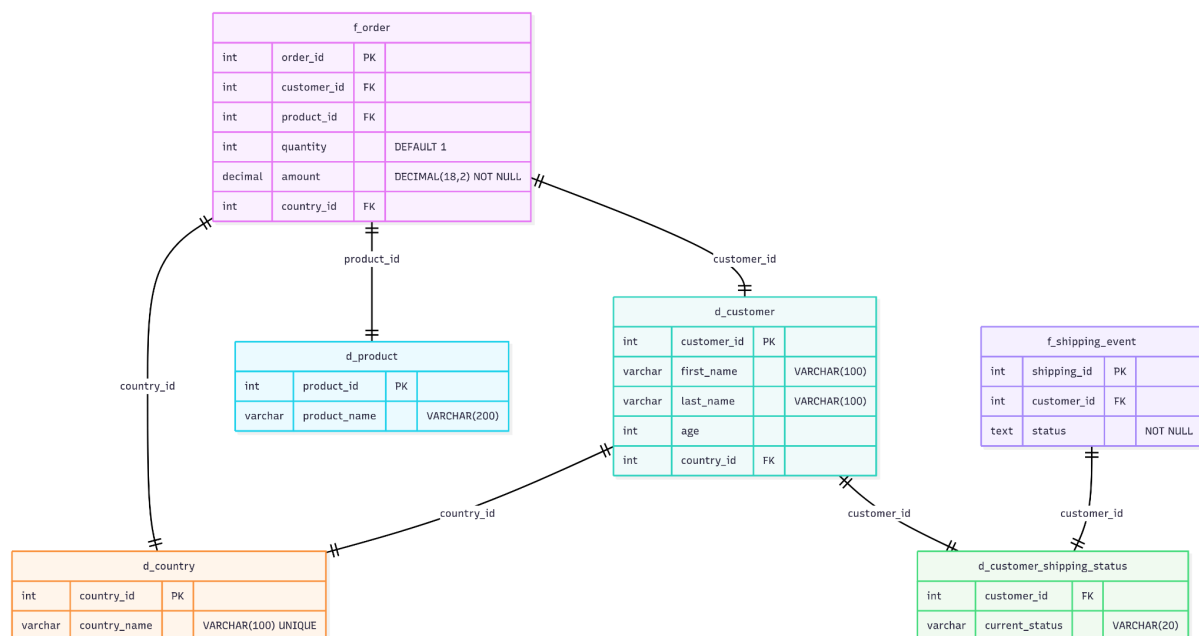
### Fact (event): f\_shipping\_event (event grain, optional now)

- **Purpose:** To preserve raw shipping rows for lineage and history (multiple rows per customer).
- **Columns:** shipping\_id (degenerate key), customer\_id (FK), status\_code (FK → d\_customer\_shipping\_status).
- **Use Cases:** QA, tracking trends of status changes, and calculating time-to-delivered once timestamps become available.

### Dim (outrigger): d\_customer\_shipping\_status (customer grain, collapsed)

- **Grain:** 1 row per customer, representing their current/collapsed shipping status.
- **Columns:** customer\_id (PK), current\_status ∈ {Pending, Delivered, NoRecord}
- **Collapse Rule:** "Pending wins"; otherwise, Delivered; if no shipping rows exist, set to NoRecord.
- **Usage:** Used to join when a report needs the "current" shipping status (e.g., requirement A) without duplicating orders.

### ER Diagram:



## 2.3 Dataset Layer (BI-facing views)

I'm exposing thin, reusable views on top of the core model so the BI tool can plug in directly.

- **DS-1 — Pending spend by country (Req A):**
  - **Logic:** f\_order → d\_country, filtered by customer\_id with a Pending status via d\_customer\_shipping\_status.
  - **Fields:** country\_name, total\_amount\_spent.
- **DS-2 — Customer × Product summary (Req B):**
  - **Logic:** f\_order → d\_customer → d\_product.
  - **Metrics:** total\_transactions (count orders), total\_quantity, total\_amount.
  - **Dims Shown:** customer name, product name.
- **DS-3 — Top product by country (Req C):**
  - **Logic:** Aggregate SUM(quantity) by (country, product) and select the max per country.
  - **Fields:** country\_name, product\_name, max\_quantity.
- **DS-4 — Most purchased product by age category (Req D):**
  - **Logic:** Derive age\_category (<30 vs. >=30) from d\_customer, then aggregate SUM(quantity) by (age\_category, product) and select the max per band.
  - **Fields:** age\_category, product\_name, total\_quantity.
- **DS-5 — Country with minimum transactions and sales (Req E):**
  - **Logic:** f\_order → d\_country; group by country and order by transactions ASC, sales\_amount ASC, returning the top row.
  - **Fields:** country\_name, transactions, sales\_amount.

---

## 2.4 Data Preparation & Transformation Requirements (for ELT jobs)

### Sources (as landed):

- customers\_raw (from Customer.csv)
- orders\_raw (from Order.csv)
- shipping\_raw (from Shipping.json)

### Load/Transform Sequence:

1. **d\_country:** INSERT DISTINCT country\_name from customers\_raw to get country\_id.
2. **d\_product:** INSERT DISTINCT product\_name from orders\_raw.Item.
3. **d\_customer:**
  - Clean First and Last names (map, strip, trim/collapse).
  - Lookup country\_id from d\_country.
4. **f\_order:**
  - Join orders\_raw to the cleaned d\_customer to get customer\_id and country\_id.
  - Lookup product\_id from d\_product.
  - Set quantity = 1.
  - Hard check: amount > 0.
5. **d\_customer\_shipping\_status:**

- Collapse shipping\_raw per customer using the "Pending wins" rule. Outer join to all customers so those without shipping data get a NoRecord status.
- Keep f\_shipping\_event (optional) to persist source events for lineage.

#### Data Quality:

- **Mandatory keys:** Customer\_ID, Order\_ID, Shipping\_ID must be unique.
- **FK coverage:** All foreign keys must exist in their respective dimension tables.
- **Domains:** shipping.status is limited to {'Pending','Delivered'}; amount > 0; age is between 10-100.
- **Name Policy:** Implemented via a documented transformation.
- **Shipping Caveat:** The current feed is at a customer-grain. For future order-level shipping support, the feed would need to include Order\_ID and timestamps.

---

## 2.5 Acceptance Criteria (Model-Level)

- **Grain integrity:**
    - f\_order is strictly order-line; no customer attributes beyond keys.
    - d\_customer\_shipping\_status is strictly customer-grain (1 row per customer).
  - **Row parity & join safety:**
    - COUNT(f\_order) must equal COUNT(orders\_raw) (after filtering).
    - Joining f\_order to d\_customer\_shipping\_status must not change the row count (no fan-out).
    - d\_customer row count must equal the number of valid customers.
  - **Conformance:** All five reporting datasets (A–E) can be produced using only these entities/views without ad-hoc fixes.
  - **Data Quality Gates:** No NULL foreign keys; amount > 0; names are cleaned; domain values are locked.
-

## 2.6 How This Model Answers the Business Questions

Requirement	Dataset	Join Path	Notes
A. Pending spend by country	DS-1	f_order → d_country + d_customer_shipping_status	Status collapsed per customer; avoids duplicates.
B. Customer totals w/ product	DS-2	f_order → d_customer → d_product	Counts, quantities, spend per customer × product.
C. Max product per country	DS-3	f_order → d_country → d_product	Max SUM(quantity) per country; deterministic tie-break if needed.
D. Top product by age band	DS-4	f_order → d_customer(age) → d_product	age_category derived at query/view layer.
E. Min country by txns & sales	DS-5	f_order → d_country	Order by (transactions, sales_amount) ascending.