Algonquin College

CST2112

Assignment 5

# National Address Register Dataset Normalization and ERD Design

**GROUP - 9**

**Shara Khandakar**

**Anam Vakil**

**Maksuda E Elahi**

**Murk Asad**

*Submitted to: Andreas Gausrab*

# Table of Contents

# 1. Introduction

The goal of this assignment was to transform the raw GA5 National Address Register (NAR) dataset into a properly normalized relational database. The raw data was provided as multiple CSV files containing complex address components such as civic information, street names, municipality descriptors, province codes, legal land descriptions (DLS), and mailing address fields.

To prepare the dataset for reliable querying and long-term storage, the files were ingested into a staging table (GA5_Address) and then transformed into a fully normalized Third Normal Form (3NF) schema. The resulting model separates distinct real-world entities (street, municipality, province, geography, and mailing address) and links them through foreign key relationships to a centralized GA5_AddressNormalized table that represents each physical address.

# 2. Identifying Entities and Relationships

An initial review of the raw NAR data revealed that many attributes were highly repetitive across rows:

- Street names and street types appeared thousands of times.
- Municipality names and subdivision codes were repeated for all addresses within the same region.
- Province codes were duplicated across the entire dataset.
- Geography fields (BG_DLS_*) described legal land descriptions shared by multiple addresses.
- Mailing address attributes repeated where postal routing was consistent.

From these patterns, the following entities were identified:

1. GA5_Street – Official street name, type, and direction
2. GA5_Municipality – Census Subdivision (CSD) names and municipal type codes
3. GA5_Province – Province code and mailing abbreviation
4. GA5_Geography – Legal land description fields and coordinate values
5. GA5_MailingAddress – Mailing street, direction, municipality name, and postal code

6. GA5_AddressNormalized – Central entity representing one physical address, with foreign keys to the above tables

The GA5_AddressNormalized table serves as the hub of the design. It contains address-specific attributes such as civic number, unit label, and unique identifiers, while referencing the supporting entity tables through foreign keys. This ensures consistency, avoids repetition, and supports efficient querying of address information.

# 3. Normalization Process

## 3.1 First Normal Form (1NF)

The original sourced NAR data was loaded into a staging table GA5_Address. In this structure, all data values were atomic (no repeating groups or multi-valued cells), which satisfies the definition of 1NF. However, this table still contained significant redundancy: street names, municipalities, provinces, and other attributes were repeated for every address row. This meant that updating a single attribute (such as a street name or municipality code) would require changes in multiple rows, increasing the risk of data inconsistencies.



**Figure1: 1NF GA5_Address (all attributes, still redundant)**

## 3.2 Second Normal Form (2NF)

To remove partial dependencies, logical groupings of attributes were identified that did not depend on the full address key. For example, street details, municipality details, and province information were repeated across multiple addresses. A 2NF representation was created by projecting these attributes into a grouped structure using a SELECT DISTINCT view, separating them from dwelling-specific elements such as civic numbers and unit labels. This eliminated partial dependencies by ensuring that attributes like street name or municipality type did not depend on the entire address record..



**Figure 2: 2NF – Street–Municipality–Province grouped (partial dependencies removed from full address)**

## 3.3 Third Normal Form (3NF)

The final design decomposed the address data into separate tables, each representing a single entity: GA5_Street, GA5_Municipality, GA5_Province, GA5_Geography, GA5_MailingAddress, and GA5_AddressNormalized. Each table contains only attributes that depend on its primary key. All transitive dependencies were removed — for example, province abbreviations are stored in GA5_Province and not in the address table. The central GA5_AddressNormalized table contains address-specific attributes and foreign keys linking to the normalized entity tables. This structure eliminates redundancy, maintains data consistency, and enables lossless reconstruction of the original address via joins

**Figure 3: Final Tables converted into 3NF**

## Summary:

The final 3NF schema ensures that each real-world concept is stored once, all non-key attributes depend only on the key of their table, and the original address can be reconstructed without data loss using foreign key relationships.

The series of figures demonstrate the transition from a single denormalized table to a fully normalized 3NF schema where each real-world entity is stored once, removing redundancy and preserving data integrity.

# 4. Implementation Summary

Three T-SQL scripts were created to automate the database pipeline.

## Script 1 – Table Creation

Created the staging table and all normalized tables, including primary keys, identity columns, and foreign key constraints.

```
--     Raw National Address Register data (all VARCHAR for safety)
------------------------------------------------------------

CREATE TABLE dbo.GA5_Address (
        LOC_GUID                VARCHAR(50),
        ADDR_GUID               VARCHAR(50),
        APT_NO_LABEL            VARCHAR(50),
        CIVIC_NO                VARCHAR(50),
        CIVIC_NO_SUFFIX         VARCHAR(50),
        OFFICIAL_STREET_NAME    VARCHAR(255),
        OFFICIAL_STREET_TYPE    VARCHAR(100),
        OFFICIAL_STREET_DIR     VARCHAR(50),
        PROV_CODE               VARCHAR(20),
        CSD_ENG_NAME            VARCHAR(255),
        CSD_FRE_NAME            VARCHAR(255),
        CSD_TYPE_ENG_CODE       VARCHAR(100),
        CSD_TYPE_FRE_CODE       VARCHAR(100),
        MAIL_STREET_NAME        VARCHAR(255),
        MAIL_STREET_TYPE        VARCHAR(100),
        MAIL_STEET_DIR          VARCHAR(50),
        MAIL_MUN_NAME           VARCHAR(255),
        MAIL_PROV_ABVN          VARCHAR(50),
        MAIL_POSTAL_CODE        VARCHAR(255),
        BG_DLS_LSD              VARCHAR(100),
        BG_DLS_QTR              VARCHAR(100),
        BG_DLS_SCTN             VARCHAR(100),
        BG_DLS_TWNSHP           VARCHAR(100),
        BG_DLS_RNG              VARCHAR(100),
```

117 %

```
------------------------------------------------------------
-- 2. NORMALIZED ENTITY TABLES (3NF)
------------------------------------------------------------

-- 2.1 Street: official street information
CREATE TABLE dbo.GA5_Street (
        StreetID INT IDENTITY(1,1) PRIMARY KEY,
        OFFICIAL_STREET_NAME  VARCHAR(255) NOT NULL,
        OFFICIAL_STREET_TYPE  VARCHAR(100) NULL,
        OFFICIAL_STREET_DIR   VARCHAR(50)  NULL
);
GO


-- 2.2 Municipality: CSD-level information
CREATE TABLE dbo.GA5_Municipality (
        MunicipalityID    INT IDENTITY(1,1) PRIMARY KEY,
        CSD_ENG_NAME        VARCHAR(255) NOT NULL,
        CSD_FRE_NAME        VARCHAR(255) NULL,
        CSD_TYPE_ENG_CODE VARCHAR(100) NULL,
        CSD_TYPE_FRE_CODE VARCHAR(100) NULL
);
GO


-- 2.3 Province: province-level information
CREATE TABLE dbo.GA5_Province (
        ProvinceID        INT IDENTITY(1,1) PRIMARY KEY,
        PROV_CODE         VARCHAR(20)  NOT NULL,
        MAIL_PROV_ABVN VARCHAR(50)  NULL
```

117 %

# Script 2 – Data Loading

Used BULK INSERT to load every NAR CSV file into the staging table (GA5_Address). This step ensured that the raw files were ingested efficiently and consistently without manual intervention.

```sql
                     -- Address_12.csv
BULK INSERT dbo.GA5_Address
FROM 'c:\cst2112_data\ga5_NAR\Address_12.csv'
WITH (
    FIRSTROW = 2,
    FIELDTERMINATOR = ',',
    ROWTERMINATOR   = '\n',
    TABLOCK
);
GO

                     -- Address_13.csv
BULK INSERT dbo.GA5_Address
FROM 'c:\cst2112_data\ga5_NAR\Address_13.csv'
WITH (
    FIRSTROW = 2,
    FIELDTERMINATOR = ',',
    ROWTERMINATOR   = '\n',
    TABLOCK
);
GO


--------------------------------------------------------------
-- Address_24 parts
--------------------------------------------------------------

BULK INSERT dbo.GA5_Address
```

## Script 3 – Data Transformation

- Populated lookup tables using INSERT...SELECT DISTINCT.

- Joined the staging data to the lookup entities using exact attribute matching.

- Filled the GA5_AddressNormalized table with foreign key references to each entity.

```sql
------------------------------------------------------------
-- 1. Load GA5_Street (official street info)
------------------------------------------------------------
INSERT INTO dbo.GA5_Street (
    OFFICIAL_STREET_NAME,
    OFFICIAL_STREET_TYPE,
    OFFICIAL_STREET_DIR
)
SELECT DISTINCT
    OFFICIAL_STREET_NAME,
    OFFICIAL_STREET_TYPE,
    OFFICIAL_STREET_DIR
FROM dbo.GA5_Address
WHERE OFFICIAL_STREET_NAME IS NOT NULL;
GO


------------------------------------------------------------
-- 2. Load GA5_Municipality (CSD info)
------------------------------------------------------------
INSERT INTO dbo.GA5_Municipality (
    CSD_ENG_NAME,
    CSD_FRE_NAME,
    CSD_TYPE_ENG_CODE,
    CSD_TYPE_FRE_CODE
)
SELECT DISTINCT
    CSD_ENG_NAME,
    CSD_FRE_NAME,
```

117 %

```sql
        GA5_Address.BU_N_CIVIC_ADD,
        GA5_Address.BU_USE
    FROM dbo.GA5_Address
    INNER JOIN dbo.GA5_Street
        ON  GA5_Address.OFFICIAL_STREET_NAME = GA5_Street.OFFICIAL_STREET_NAME
        AND ISNULL(GA5_Address.OFFICIAL_STREET_TYPE,'') = ISNULL(GA5_Street.OFFICIAL_STREET_TYPE,'')
        AND ISNULL(GA5_Address.OFFICIAL_STREET_DIR,'')  = ISNULL(GA5_Street.OFFICIAL_STREET_DIR,'')
    INNER JOIN dbo.GA5_Municipality
        ON  GA5_Address.CSD_ENG_NAME = GA5_Municipality.CSD_ENG_NAME
        AND ISNULL(GA5_Address.CSD_FRE_NAME,'')      = ISNULL(GA5_Municipality.CSD_FRE_NAME,'')
        AND ISNULL(GA5_Address.CSD_TYPE_ENG_CODE,'') = ISNULL(GA5_Municipality.CSD_TYPE_ENG_CODE,'')
        AND ISNULL(GA5_Address.CSD_TYPE_FRE_CODE,'') = ISNULL(GA5_Municipality.CSD_TYPE_FRE_CODE,'')
    INNER JOIN dbo.GA5_Province
        ON  GA5_Address.PROV_CODE = GA5_Province.PROV_CODE
        AND ISNULL(GA5_Address.MAIL_PROV_ABVN,'') = ISNULL(GA5_Province.MAIL_PROV_ABVN,'');
    GO


    ------------------------------------------------------------
    USE [25f_cst2112_group_09];
    GO


    -- 1NF: all attributes in a single wide table (no repeating groups)
    SELECT TOP (20) *
    FROM dbo.GA5_Address;
```

117 %

# 5. Challenges and Lessons Learned

## Challenge 1 – Extremely Large Data Volume

The dataset contained millions of records. Lookup tables such as GA5_Street and GA5_Geography had tens of thousands to millions of distinct values. Efficient joins and appropriate indexing were required to avoid performance bottlenecks.

## Challenge 2 – Matching Lookup Values

Some fields contained NULL values or minor spelling inconsistencies. These were addressed using:

- ISNULL() comparisons during joins

- LEFT JOIN for optional relationships instead of INNER JOIN

## Challenge 3 – Normalization vs. Analytical Schemas

Initially, there was uncertainty between designing a 3NF relational model versus a dimensional star schema. After clarifying this doubt during the lecture, the focus was therefore on eliminating redundancy and ensuring functional dependency, not on analytical reporting measures.

## Challenge 4 – Foreign Key Load Order

Foreign key dependencies required careful load sequencing. When reloading data, tables were cleared and repopulated in the following order:

1. GA5_AddressNormalized

2. GA5_MailingAddress

3. GA5_Geography

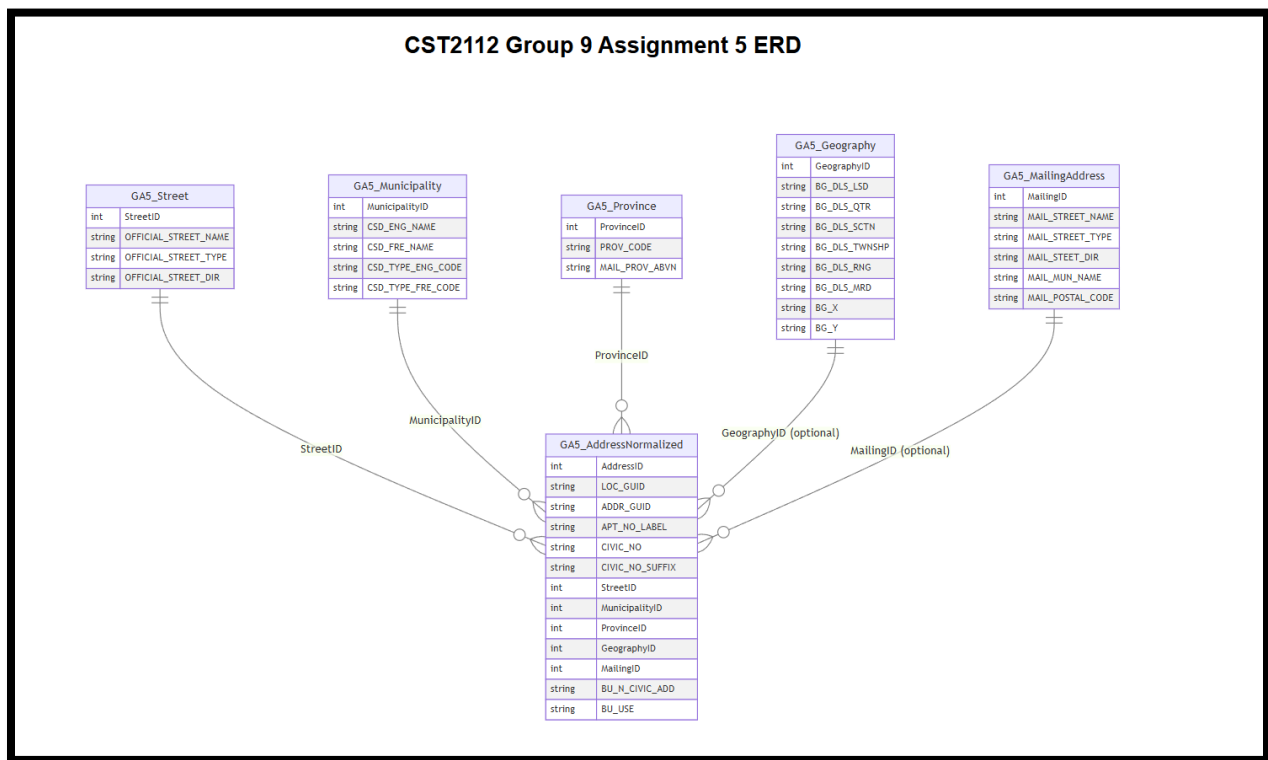4. GA5_Province

5. GA5_Municipality

6. GA5_Street

This ordering prevented constraint violations and was incorporated into Script 3.

# 6. Final ERD

The final ERD (included in the submission) demonstrates:

- Five lookup entities (Street, Municipality, Province, Geography, MailingAddress)

- One central address entity (GA5_AddressNormalized)

- One-to-many relationships from lookup entities to addresses

- Optional relationships for geography and mailing using nullable foreign keys

This accurately represents the final 3NF schema implemented in SQL Server.



# 7. Conclusion

The normalization process transformed a duplicated, semi-structured CSV dataset into a clean relational database that supports reliable querying and avoids update anomalies.

Key outcomes included:

- Identifying real-world entities and functional dependencies

- Separating repeating attributes into lookup tables

- Enforcing referential integrity through foreign keys

- Automating creation, loading, and population using T-SQL

- Producing a fully normalized 3NF structure consistent with address data management principles

This project strengthened my understanding of database normalization, staging pipelines, and T-SQL scripting, and provided practical experience working with large-scale government datasets.

.