

# Proyecto: Carga y Análisis Concurrente de Datos Masivos de Viajes en Taxi de NYC

## Objetivos del Proyecto:

El objetivo principal de este proyecto es que los estudiantes demuestren su capacidad para integrar conocimientos de diseño de bases de datos relacionales y programación concurrente en Java para resolver un problema de ingeniería común: la ingesta y el análisis eficiente de grandes volúmenes de datos de transporte urbano.

Los objetivos específicos incluyen:

- Diseño de Bases de Datos:** Aplicar los principios de normalización (hasta la Tercera Forma Normal - 3FN) para diseñar un esquema de base de datos robusto y eficiente en MySQL, basado en los datos de viajes en taxi.
- Manejo de Datos Masivos:** Desarrollar una solución en Java para procesar un archivo CSV de grandes dimensiones (más de 3 millones de filas) de los registros de viajes en taxi de NYC.
- Optimización de E/S:** Implementar técnicas de programación concurrente para optimizar el tiempo de carga de los datos del CSV a la base de datos.
- Consultas Avanzadas:** Diseñar e implementar métodos concurrentes para realizar consultas complejas y analíticas sobre la base de datos.

## Fases del Proyecto y Requisitos Técnicos

### Paso 1: Adquisición y Preparación del Conjunto de Datos (CSV)

**Requisito:** El proyecto debe comenzar con el archivo CSV `yellow_tripdata_2025-02.csv` de los **NYC Yellow Taxi Trip Records** (registros de viajes en taxi amarillo de NYC, febrero de 2025), conteniendo **más de tres millones de filas**. Los datos están disponibles en el sitio web de la [Comisión de Taxis y Limusinas de NYC \(TLC\)](#).

**Estructura de Columnas Clave del CSV (Ejemplo):**

| Columna                            | Tipo de Dato          | Descripción  | Normalización  |
|------------------------------------|-----------------------|--|--|
| <code>VendorID</code>              | <code>INT</code>      | Identificador del proveedor de la TLC.                   | Redundante (Depende de la descripción del proveedor)     |
| <code>tpep_pickup_datetime</code>  | <code>DATETIME</code> | Fecha y hora de recogida.                                | Atributo de la entidad <code>Trip</code>                 |
| <code>tpep_dropoff_datetime</code> | <code>DATETIME</code> | Fecha y hora de destino.                                 | Atributo de la entidad <code>Trip</code>                 |
| <code>passenger_count</code>       | <code>INT</code>      | Número de pasajeros.                                     | Atributo de la entidad <code>Trip</code>                 |
| <code>trip_distance</code>         | <code>FLOAT</code>    | Distancia del viaje en millas.                           | Atributo de la entidad <code>Trip</code>                 |
| <code>RateCodeID</code>            | <code>INT</code>      | Código de tarifa (e.g., 1=Estándar, 2=JFK).              | Redundante (Depende de la descripción de la tarifa)      |
| <code>PULocationID</code>          | <code>INT</code>      | ID de la Zona de Recogida (Pick-up Location ID).         | Redundante (Depende del nombre de la zona y el distrito) |
| <code>DOLocationID</code>          | <code>INT</code>      | ID de la Zona de Destino (Drop-off Location ID).         | Redundante (Depende del nombre de la zona y el distrito) |
| <code>payment_type</code>          | <code>INT</code>      | Código del método de pago (e.g., 1=Tarjeta, 2=Efectivo). | Redundante (Depende de la descripción del pago)          |
| <code>fare_amount</code>           | <code>FLOAT</code>    | Monto de la tarifa.                                      | Atributo de la entidad <code>Trip</code>                 |
| <code>tip_amount</code>            | <code>FLOAT</code>    | Monto de la propina.                                     | Atributo de la entidad <code>Trip</code>                 |
| <code>total_amount</code>          | <code>FLOAT</code>    | Monto total cobrado.                                     | Atributo de la entidad <code>Trip</code>                 |

**Nota para el Estudiante:** Se adjunta el archivo de **TLC Taxi Zone Lookup Table** (tabla de búsqueda de zonas de taxi) `taxi_zone_lookup.csv` para obtener los nombres de las zonas (`Zone`) y los distritos (`Borough`) a partir de los `LocationID`. También el diccionario de datos `data_dictionary_trip_records_yellow.pdf` que describe la estructura del CSV.

### Paso 2: Diseño y Normalización de la Base de Datos (MySQL - 3FN)

**Requisito:** Diseñar y crear la base de datos en MySQL. El esquema debe estar normalizado a la **Tercera Forma Normal (3FN)** para eliminar dependencias transitivas y redundancia, asegurando la integridad de los datos.

**Proceso de Normalización Sugerido:**

- Identificación de Redundancias:** Las descripciones de `VendorID`, `RateCodeID`, `payment_type`, y la información de la zona (`Zone`, `Borough`) que depende de `PULocationID` / `DOLocationID` son dependencias funcionales que violan 2FN y 3FN.
- Creación de Tablas Normalizadas.**
- Implementación en MySQL:** Crear las sentencias `CREATE TABLE` con claves primarias, foráneas e índices adecuados.

**Entregable:** El script SQL (`schema.sql`) con la definición completa de la base de datos.

### Paso 3: Carga Concurrente de Datos del CSV a la BD (Java)

**Requisito:** Implementar una aplicación en Java que utilice programación concurrente para leer el archivo CSV y cargar los datos en las tablas normalizadas de MySQL.

**Conceptos Clave a Aplicar:**

- Patrón Productor-Consumidor:** Un hilo productor lee el CSV, múltiples hilos consumidores procesan las líneas, realizan la lógica de inserción y ejecutan **inserciones por lotes (Batch Inserts)** en JDBC.
- Medición de Rendimiento:** Medir y reportar el tiempo total de carga para diferentes configuraciones de hilos.

### Paso 4: Consultas Concurrentes y Simulación de Carga Distribuida (Java)

**Requisito:** Implementar métodos en Java que utilicen programación concurrente para ejecutar un conjunto de consultas analíticas complejas a la base de datos, simulando un entorno de alta demanda.

**Conceptos Clave a Aplicar:**

- Ejecución Asíncrona:** Utilizar un `ExecutorService` y tareas para gestionar la ejecución concurrente de las consultas.
- Pool de Conexiones:** Uso de un Pool de Conexiones JDBC para gestionar eficientemente las conexiones bajo alta concurrencia.

**Consultas Analíticas Complejas (Ejemplos):**

- Consulta 1: Análisis de Propina por Zona y Hora Pico:** Calcular el **monto promedio de propina** por **método de pago** y **Borough de recogida** para viajes realizados en un rango de horas específico (e.g., 7:00 AM - 9:00 AM).
- Consulta 2: Identificación de Rutas Rentables:** Identificar el **Top 5 de pares de zonas (Recogida-Destino)** que generan la **mayor tarifa promedio por milla** para viajes con más de 2 pasajeros.
- Consulta 3: Distribución de Viajes por Tipo de Tarifa:** Contar el número total de viajes por cada `RateCodeDescription` (e.g., JFK, Newark, Estándar) y calcular el porcentaje que representa del total de viajes.
- Consultas 4-8:** Deben ser generadas por los integrantes del equipo y debidamente explicadas en el informe de entrega.

**Entregable:** Código Java con el método principal que orquesta la carga de datos (Paso 3) y la ejecución de consultas concurrentes (Paso 4), incluyendo la medición de rendimiento.

## Estructura del Informe de Entrega

El informe debe ser un documento técnico y reflexivo que demuestre la comprensión de los conceptos teóricos y las decisiones prácticas tomadas durante la implementación. Debe tener un mínimo de 10 páginas y un máximo de 15.

### I. Introducción

- Objetivos del Proyecto:** Reafirmar los objetivos planteados.
- Contexto y Motivación:** Descripción del problema de la ingesta y análisis de datos de transporte urbano de NYC.
- Estructura del Informe:** Descripción concisa de las secciones siguientes.

### II. Diseño de la Base de Datos

- Análisis del CSV y Dependencias Funcionales:**
  - Descripción de la estructura del CSV de viajes en taxi.
  - Identificación de las dependencias funcionales (parciales y transitivas) que justifican la normalización (e.g., `LocationID` → `ZoneName`, `Borough`).
- Esquema Normalizado (3FN):**
  - Diagrama Entidad-Relación (DER) del esquema final (`Trips`, `Vendors`, `RateCodes`, `PaymentTypes`, `TaxiZones`).
  - Justificación detallada de cómo se alcanzó la 3FN para cada tabla.

### III. Implementación de la Carga Concurrente (Paso 3)

- Arquitectura Concurrente:**
  - Diagrama de flujo o de componentes que ilustre el patrón **Productor-Consumidor**.
  - Justificación de la elección del número de hilos consumidores.
- Manejo de la Persistencia:**
  - Explicación del uso de **Batch Inserts** en JDBC.
  - Descripción de la lógica para manejar la inserción en las tablas de dimensiones (uso de caché concurrente).
- Ánalisis de Rendimiento:**
  - Tabla comparativa del tiempo de carga total para diferentes configuraciones de hilos (e.g., 1, 4, 8).
  - Gráfico de barras o líneas que visualice la ganancia de rendimiento (Speedup) obtenida por la concurrencia.

### IV. Implementación de Consultas Concurrentes (Paso 4)

- Diseño de la Capa de Consultas:**
  - Explicación del uso del `ExecutorService` y el pool de conexiones.
- Detalle de las Consultas:**
  - Para cada consulta compleja implementada, incluir: la pregunta de negocio, la sentencia SQL utilizada y el tiempo de ejecución medido.

### V. Conclusiones y Trabajo Futuro

- Conclusiones:** Resumen de los logros del proyecto y las lecciones aprendidas.
- Recomendaciones:** Identificación de las debilidades o áreas de mejora.
- Trabajo Futuro:** Propuestas para extender el proyecto.

### VI. Declaración del uso de IA Generativa

- Declaramos haber hecho uso del sistema de IA Generativa (Nombre del sistema/herramienta IA y versión: ChatGPT, Gemini, Copilot...)** para:

- Ejemplos de instrucciones (prompts):**

### Adjuntar

- Código Fuente.
- Script SQL de Creación de la BD.