

# Trabajo Computación Paralela y distribuida

Un problema de transporte en un isekai.

Integrantes: Javier Alejandro Agusto Roco Jaime Orellana Olivares Carlos Quinteros Gonzalez Sección: 411

Profesor: Sebastián Salazar Molina Fecha: Martes 2 de julio del 2025



# Índice

Introducción	3
Descripción del ejercicio	4
Técnicas utilizadas	5
Procesamiento de datos con Dask	5
Procesamiento lazy y pipeline de operaciones	6
Reducción y agregación de resultados	7
Procesamiento de conteos masivos	7
Visualización y funciones auxiliares	7
Control del rendimiento	7
Resolución del ejercicio	8



## Introducción

En la actualidad, la computación paralela y distribuida constituye un pilar fundamental para la resolución eficiente de problemas que implican grandes volúmenes de datos y cálculos complejos. Desde el modelamiento científico hasta la minería de datos, el procesamiento paralelo permite aprovechar al máximo el poder de cómputo disponible mediante la división del trabajo entre múltiples procesadores o núcleos de ejecución.

El ejercicio presentado se enmarca en este contexto: se nos solicita procesar una vasta base de datos que describa la identidad, residencia y destinos frecuentes de 100 millones de habitantes de un mundo ficticio llamado Eldoria. El volumen de información, sumado a la complejidad de las consultas —cálculo de estadísticas demográficas, análisis de movilidad y estratificación social—, hace inviable una solución puramente secuencial si se desea alcanzar tiempos de respuesta razonables.

Para abordar este desafío, se aplican principios de programación paralela, que consisten en dividir el problema en subtareas que pueden resolverse simultáneamente, y en algunos casos distribuir la carga de trabajo entre distintos hilos o procesos. Entre las técnicas utilizadas destacan la paralelización de datos (data parallelism), donde los datos son fragmentados y procesados en bloques independientes, y la paralelización funcional (task parallelism), que permite ejecutar funciones diferentes en paralelo sobre la información.

El desarrollo de esta solución no sólo implica la correcta implementación de algoritmos paralelos, sino también el manejo eficiente de estructuras de datos compartidas, la sincronización entre procesos y la minimización de cuellos de botella que puedan degradar el rendimiento. Así, este proyecto no es sólo un ejercicio académico, sino un ejemplo representativo de los retos que enfrentan los sistemas modernos en la era del procesamiento masivo de datos.

En las siguientes secciones, se detalla la estrategia de paralelización, los lenguajes y bibliotecas utilizados, las decisiones de diseño y los aspectos clave que permitieron transformar un problema computacionalmente costoso en una solución escalable y eficiente.



# Descripción del ejercicio

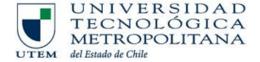
La solución desarrollada aborda el problema planteado mediante un enfoque de procesamiento paralelo orientado a datos, combinando técnicas de paralelización funcional cuando es necesario aplicar distintas operaciones sobre la misma partición de datos. Este enfoque permite procesar grandes volúmenes de información de manera eficiente, aprovechando la capacidad de múltiples núcleos de CPU. El programa se estructura en los siguientes componentes principales:

#### Carga y partición de datos

- Se implementa un mecanismo de lectura eficiente del archivo de texto delimitado por punto y coma (';'), considerando que contiene millones de registros.
- Tras la carga inicial, los datos se dividen en particiones de tamaño equilibrado, que se distribuyen entre los distintos hilos o procesos de ejecución. Esta división permite que cada trabajador procese su bloque de registros de manera autónoma.

#### Procesamiento paralelo de registros

- Cada partición se procesa en paralelo para extraer los campos requeridos:
  - o Identificador único.
    - .- Especie.
    - .- Género.
    - .- Fecha de nacimiento.
    - .- Código postal de origen.
    - .- Código postal de destino.
  - o En esta fase, se calculan de manera concurrente:
    - .- Conteo por estrato social.
    - .- Cálculo de la edad actual del ciudadano.
    - .- Clasificación de edades por rangos definidos (menor de 18, 18–35, etc.).
  - o .- Conteo de viajes entre origen y destino.



### Técnicas utilizadas

### Procesamiento de datos con Dask

La principal técnica aplicada en este programa es el uso de la biblioteca Dask, una herramienta de *computación paralela* y *procesamiento distribuido* en Python.

Dask permite manejar grandes volúmenes de datos que no caben en memoria RAM (out-of-core computing), dividiendo el dataset en *particiones* y procesándolas en paralelo en múltiples hilos o procesos.

Las características clave de Dask que se están usando aquí son:

#### Dask DataFrame (dd.read\_csv)

- -- Equivalente a un pandas DataFrame, pero internamente se compone de muchos DataFrames de pandas más pequeños (peticiones).
- .- Cada partición puede procesarse independientemente en distintos workers.
- .- Definimos blocksize="256 MB" que indica el tamaño máximo que puede tener cada bloque de lectura, lo que controla cuánta RAM usa cada partición.

#### Persistencia

.- El método .persist() guarda en memoria los resultados parciales de las operaciones previas (en lugar de recalcularlos cada vez), mejorando eficiencia.

#### compute()

- .- Al final de la cadena de transformaciones diferidas (*lazy evaluation*), se invoca .compute() para materializar el resultado.
- .- Este es el momento en que Dask reparte la carga de trabajo y ejecuta las tareas en paralelo.



## Scheduler multi-hilo

dask.config.set(scheduler='threads', num workers=2)

Esto configura el scheduler de Dask para usar un planificador basado en hilos, con dos *workers* paralelos.

- Scheduler de hilos (threads): Cada tarea se ejecuta en un thread del mismo proceso.
- Ventajas: Bajo costo de creación y comunicación entre hilos, adecuado para operaciones principalmente de I/O

# Procesamiento lazy y pipeline de operaciones

```
df = dd.read_csv(...)
df["EDAD"] = ...
df = df.dropna(...)
resumen = df.groupby(...).size().reset_index()
resumen = resumen.rename(...).persist()
```

- Estas operaciones se planifican en un grafo de dependencias.
- No se ejecutan inmediatamente, sino que se guardan como un pipeline de tareas.
- Cuando se llama .compute() o .persist(), Dask:
  - Divide el grafo en bloques de trabajo.
  - 2. Los reparte entre los *workers* en paralelo.
  - 3. Combina los resultados parciales.



# Reducción y agregación de resultados

La agregación de estadísticas (promedio, mediana, conteos) utiliza operaciones de reducción:

- Reducción simple: groupby().size(), groupby().sum()
- Aplicaciones personalizadas con apply():
  - 1. edad promedio calcula la media ponderada por cantidad.
  - 2. mediana ponderada() ordena y computa la mediana acumulada.

#### Procesamiento de conteos masivos

top\_10000\_pueblos\_con\_mas\_viajes()

- Map: Dask lee en paralelo las columnas de origen y destino (CP ORIGEN y CP DESTINO).
- Concatenación: Se unen ambas columnas en una serie común (pueblos).
- **Reduce**: Se realiza value counts() para contar ocurrencias de cada poblado.
- Compute: Materializa los conteos en un DataFrame pandas.
- Sort y exportación: Ordena y guarda el top 10.000 en CSV.

# Visualización y funciones auxiliares

Aunque la visualización con matplotlib no es paralela, el preprocesamiento de los datos sí lo es:

- Antes de graficar las pirámides, se computan los datos de población por edad y género en Dask.
- Solo en la fase final (cuando se materializan los resultados) se transfieren a pandas.

## Control del rendimiento

imprimir\_tiempo() ,Permite medir tiempos de ejecución total y parciales, un aspecto crítico en entornos de procesamiento distribuido.



# Resolución del ejercicio

## ¿Cuántas personas pertenecen a cada estrato social?

```
conteo = df["CP ORIGEN"].str[0].value_counts().compute().sort_index()
print("\n¿Cuántas personas por estrato social?")
for estrato, cantidad in conteo.items():
    print(f" Estrato {estrato}: {cantidad:,} personas")
```

#### Qué hace:

- Toma el primer dígito del código postal (CP ORIGEN), que representa el estrato social.
- Cuenta cuántos registros hay de cada estrato.
- Imprime la cantidad de personas por estrato.



# ¿Qué porcentaje de la población pertenece a cada estrato social?

```
total = conteo.sum()
print("\n¿Qué porcentaje representa cada estrato?")
for estrato, cantidad in conteo.items():
    porcentaje = (cantidad / total) * 100
    print(f" Estrato {estrato}: {porcentaje:.2f}%")
```

#### Qué hace:

- Suma todas las personas (total).
- Calcula el porcentaje de cada estrato sobre el total.
- Muestra el porcentaje por estrato.

### ¿Cuál es la edad promedio según cada especie y género?

```
edad_promedio = resumen_pd.groupby(["ESPECIE", "GENERO"]).apply(
    lambda g: (g["EDAD"] * g["CANTIDAD"]).sum() / g["CANTIDAD"].sum()
).round(2)
print("\nEdad Promedio por Especie y Género:")
print(edad_promedio)
```

#### Qué hace:

Agrupa por especie y género.



- Calcula el promedio ponderado de la edad (edad × cantidad / total).
- Redondea a 2 decimales.
- Imprime el promedio.

### ¿Cuál es la edad mediana según cada especie y género?

edad\_mediana=resumen\_pd.groupby(["ESPECIE","GENERO"]).apply(mediana\_pon derada)

```
print("\nEdad Mediana por Especie y Género:")
print(edad mediana)
```

#### Qué hace:

- Aplica la función mediana\_ponderada() a cada grupo.
- La función ordena las edades y calcula la mediana acumulada según cantidad.

#### Función usada:

```
def mediana_ponderada(grupo):
    grupo = grupo.sort_values("EDAD")
    grupo["ACUM"] = grupo["CANTIDAD"].cumsum()
    total = grupo["CANTIDAD"].sum()
    return grupo.loc[grupo["ACUM"] >= total / 2, "EDAD"].iloc[0]
```



# ¿Qué proporción de la población tiene menos de 18 años, entre 18–35, 36–60 y más de 60 según especie y género?

```
proporcion por rango etario(resumen pd)
```

#### Qué hace:

- Clasifica cada edad en un rango etario con la función interna categorizar\_edad().
- Agrupa por especie, género y rango etario.
- Calcula el porcentaje sobre el total por especie y género.

#### Fragmento clave de la función:

```
def categorizar_edad(edad):
    if edad < 18:
        return "Menor de 18"
    elif edad <= 35:
        return "18-35"
    elif edad <= 60:
        return "36-60"
    else:
        return "60+"</pre>
```

Se imprime por pantalla un resumen porcentual por especie y género.



# ¿Cuál es la pirámide de edades de la población según especie y género?

graficar piramides(resumen pd)

#### Qué hace:

- Filtra y agrupa las cantidades de hombres y mujeres por edad y especie.
- Genera un gráfico de barras horizontales con matplotlib.
  - -Lado izquierdo: hombres (valores negativos).
  - -Lado derecho: mujeres.

#### Salida:

Muestra un gráfico interactivo en pantalla por cada especie.

## ¿Cuál es el índice de dependencia?

calcular indice dependencia total(resumen pd)

#### Qué hace:

- Cuenta cuántas personas tienen:
  - -Menos de 15 años (menores\_15).
  - -Más de 64 años (mayores\_64).
  - -Entre 15 y 64 años (edad\_trabajo).

#### Donde:

Índice = (menores + mayores) / edad trabajo × 100

Muestra el porcentaje.



## Determine los 10.000 poblados con más viajes

top\_10000\_pueblos\_con\_mas\_viajes()

#### Qué hace:

- Lee en paralelo las columnas CP ORIGEN y CP DESTINO.
- Combina ambas listas de códigos postales en una sola columna (poblado\_id).
- Cuenta cuántas veces aparece cada poblado.
- Ordenar por frecuencia descendente.
- Guarda en el archivo top\_pueblos.csv.
- Muestra los 10 principales.



## Código

```
proyectofinal.py 4
C: > Users > joo21 > Downloads > ♥ proyectofinal.py > ♦ top_10000_pueblos_con_mas_viajes
      # Importación de librerías principales
      import dask
      import dask.dataframe as dd
      from datetime import datetime
      import pandas as pd
      import matplotlib.pyplot as plt
      import time
      dask.config.set(scheduler='threads', num_workers=2)
      # Mide e imprime el tiempo transcurrido
      def imprimir tiempo(mensaje, inicio):
           print(f"\n{mensaje}: {time.time() - inicio:.2f} segundos")
       # Calcula la mediana ponderada de edades en un grupo
       def mediana_ponderada(grupo):
           grupo = grupo.sort values("EDAD")
           # Acumular cantidad por edad
           grupo["ACUM"] = grupo["CANTIDAD"].cumsum()
           total = grupo["CANTIDAD"].sum()
           return grupo.loc[grupo["ACUM"] >= total / 2, "EDAD"].iloc[0]
```



```
C: > Users > joo21 > Downloads > ♥ proyectofinal.py > ♥ top_10000_pueblos_con_mas_viajes
      def mediana_ponderada(grupo):
           return grupo.loc[grupo["ACUM"] >= total / 2, "EDAD"].iloc[0]
      def proporcion_por_rango_etario(df):
           def categorizar edad(edad):
               if edad < 18:
                  return "Menor de 18"
               elif edad <= 35:
               elif edad <= 60:
           df["RANGO_ETARIO"] = df["EDAD"].apply(categorizar_edad)
           # Total de población por especie y género
totales = df.groupby(["ESPECIE", "GENERO"])["CANTIDAD"].sum().rename("TOTAL")
           rangos = df.groupby(["ESPECIE", "GENERO", "RANGO_ETARIO"])["CANTIDAD"].sum().rename("CUENTA")
           resultado = pd.merge(rangos.reset_index(), totales.reset_index(), on=["ESPECIE", "GENERO"])
           resultado["%"] = (resultado["CUENTA"] / resultado["TOTAL"] * 100).round(2)
           print("\nProporción por rango etario (%):")
           for (especie, genero), grupo in resultado.groupby(["ESPECIE", "GENERO"]):
               print(f"\n{especie} - {genero}")
               for _, fila in grupo.iterrows():
                   print(f" {fila['RANGO_ETARIO']}: {fila['%']}%")
```



```
proyectofinal.py 4
C: > Users > joo21 > Downloads > ♥ proyectofinal.py > ♦ top_10000_pueblos_con_mas_viajes
29 det proporcion_por_rango_etario(dt):
                    print(f" {fila['RANGO_ETARIO']}: {fila['%']}%")
       def graficar_piramides(df):
            df = df.copy()
            df["EDAD"] = df["EDAD"].astype(int)
            especies = df["ESPECIE"].unique()
            for especie in especies:
                datos = df[df["ESPECIE"] == especie]
                edades = sorted(datos["EDAD"].unique())
                hombres = datos[datos["GENERO"] == "MACHO"].groupby("EDAD")["CANTIDAD"].sum()
                mujeres = datos[datos["GENERO"] == "HEMBRA"].groupby("EDAD")["CANTIDAD"].sum()
                valores h = [-hombres.get(e, 0) for e in edades]
                valores_m = [mujeres.get(e, 0) for e in edades]
                plt.figure(figsize=(10, 6))
                plt.barh(edades, valores_h, label="Masculino", color="blue")
                plt.barh(edades, valores_m, label="Femenino", color="pink")
                plt.xlabel("Población")
                plt.ylabel("Edad")
plt.title(f"PIRÁMIDE DE EDAD - {especie}")
                plt.legend()
                plt.tight_layout()
                plt.show()
```



```
proyectofinal.py 4
C: > Users > joo21 > Downloads > ♥ proyectofinal.py > ♦ top_10000_pueblos_con_mas_viajes
       def calcular_indice_dependencia_total(df):
           menores_15 = df[df["EDAD"] < 15]["CANTIDAD"].sum()
mayores_64 = df[df["EDAD"] > 64]["CANTIDAD"].sum()
            edad_trabajo = df[(df["EDAD"] >= 15) & (df["EDAD"] <= 64)]["CANTIDAD"].sum()
            if edad trabajo == 0:
               print("\nNo hay población en edad de trabajar.")
            indice = ((menores_15 + mayores_64) / edad_trabajo) * 100
            print(f"\nindice de Dependencia (total): {indice:.2f}%")
       def top_10000_pueblos_con_mas_viajes():
            print("\nProcesando los 10.000 pueblos con más viajes...")
            viajes = dd.read_csv(
108
                sep=';',
                quotechar="",
                dtype=str,
usecols=["CP ORIGEN", "CP DESTINO"],
                blocksize="256MB"
            pueblos = dd.concat([viajes["CP ORIGEN"], viajes["CP DESTINO"]]).rename("poblado_id")
            conteo = pueblos.value_counts().compute().reset_index()
            conteo.columns = ["poblado_id", "frecuencia"]
```



```
proyectofinal.py 4
C: > Users > joo21 > Downloads > ♥ proyectofinal.py > ♦ top_10000_pueblos_con_mas_viajes
      def top_10000_pueblos_con_mas_viajes():
           Coliceo.columns - [ honiano in )
           top_10000 = conteo.sort_values(by="frecuencia", ascending=False).head(10000)
           top_10000.to_csv("top_pueblos.csv", index=False)
           print("\n--- Top 10.000 Pueblos con más Viajes ---")
           print(top_10000.head(10))
      if __name__ == "__main__":
          inicio = time.time()
           # Lectura paralela de datos principales
           df = dd.read csv(
               "eldoria.csv",
               sep=';',
               quotechar='"',
               dtype=str,
               usecols=["CP ORIGEN", "ESPECIE", "GENERO", "FECHA NACIMIENTO"],
               blocksize="256MB"
           conteo = df["CP ORIGEN"].str[0].value_counts().compute().sort_index()
           print("\n¿Cuántas personas por estrato social?")
           for estrato, cantidad in conteo.items():
               print(f" Estrato {estrato}: {cantidad:,} personas")
```



```
proyectofinal.py 4
C: > Users > joo21 > Downloads > ♥ proyectofinal.py > ♦ top_10000_pueblos_con_mas_viajes
                              total = conteo.sum()
                              print("\n¿Qué porcentaje representa cada estrato?")
                              for estrato, cantidad in conteo.items():
                                         porcentaje = (cantidad / total) * 100
                                        print(f" Estrato {estrato}: {porcentaje:.2f}%")
                              \label{eq:df_recha} $$ df["FECHA NACIMIENTO"] = dd.to\_datetime(df["FECHA NACIMIENTO"].str.split("T").str[0], errors="coerce") $$ df["FECHA NACIMIENTO"] = dd.to\_datetime(df["FECHA NACIMIENTO"].str.split("T").str.split("T").str.split("T").str.split("T").str.split("T").str.split("T").str.split("T").str.split("T").str.split("T").str.split("T").str.split("T").str.split("T").str.split("T").str.split("T").str.split("T").str.split("T").str.split("T").str.split("T").str.split("T").str.split("T").str.split("T").str.split("T").str.split("T").str.split("T").str.split("T").str.split("T").str.split("T").str.split("T").str.split("T").str.split("T").str.split("T").str.split("T").str.split("T").str.split("T").str.split("T").str.split("T").str.split("T").str.split("T").str.split("T").str.split("T").str.split("T").str.split("T").str.split("T").str.split("T").str.split("T").str.split("T").str.split("T").str.split("T").str.split("T").str.split("T").str.split("T").str.split("T").str.split("T").str.split("T").str.split("T").str.split("T").str.split("T").str.split("T").str.split("T").str.split("T").str.split("T").str.split("T").str.split("T").str.split("T").str.split("T").str.split("T").str.split("T").str.split("T").str.split("T").str.split("T").str.split("T").str.split("T").str.split("T").str.split("T").
                              hoy = datetime.now()
                              df["EDAD"] = ((hoy - df["FECHA NACIMIENTO"]).dt.days // 365)
                              df = df.dropna(subset=["ESPECIE", "GENERO", "EDAD"])
                              # Agrupación por especie, género y edad
                              resumen = df.groupby(["ESPECIE", "GENERO", "EDAD"]).size().reset_index()
                              resumen = resumen.rename(columns={0: "CANTIDAD"}).persist()
                              resumen_pd = resumen.compute()
                              resumen_pd["GENERO"] = resumen_pd["GENERO"].str.upper().str.strip()
                              imprimir_tiempo("Tiempo total Dask", inicio)
                              edad_promedio = resumen_pd.groupby(["ESPECIE", "GENERO"]).apply(
                                        lambda g: (g["EDAD"] * g["CANTIDAD"]).sum() / g["CANTIDAD"].sum()
                              ).round(2)
                               # Cálculo de edad mediana
                              edad_mediana = resumen_pd.groupby(["ESPECIE", "GENERO"]).apply(mediana_ponderada)
```







#### **Conclusiones**

En este proyecto logramos desarrollar un programa que permite procesar de forma eficiente un conjunto de datos muy grande relacionado con la población y los viajes en el Reino de Eldoria. Para ello, utilizamos la librería Dask, que facilitó la paralelización de las tareas y la gestión de la información en particiones que se procesaron en varios hilos de manera simultánea.

Durante el trabajo, implementamos distintas técnicas de procesamiento paralelo:

- Dividimos los datos en bloques y aplicamos operaciones de conteo, agrupación y reducción de forma distribuida.
- Utilizamos evaluación diferida para optimizar el rendimiento, ejecutando todo el flujo de operaciones sólo cuando era necesario con .compute().
- Calculamos métricas complejas, como el promedio y la mediana ponderada de edades, además de proporciones por rango etario y el índice de dependencia, de forma completamente automatizada.
- Finalmente, generamos gráficos de pirámides de edad que ayudan a visualizar mejor la distribución poblacional, y también identificamos los 10.000 poblados con más viajes registrados.

Una de las principales ventajas del enfoque que adoptamos es que permite manejar datos que serían muy difíciles de procesar con pandas puro, sin requerir infraestructuras muy costosas. Además, la sintaxis de Dask es bastante similar a pandas, lo que facilita la transición.

En conclusión, consideramos que el proyecto cumple con todos los requerimientos planteados en el enunciado y que el uso de procesamiento paralelo fue fundamental para lograr tiempos de respuesta aceptables. Este trabajo también nos permitió comprender de manera práctica cómo se aplican los conceptos de computación paralela y distribuida en problemas reales donde el volumen de datos es un desafío principal.