

# Py-DataVault: Sistema de Backup Seguro con Dask

Alejandro Ríos Muñoz  
Lina Sofía Ballesteros Merchán  
Juan Esteban García Galvis  
David Grisales Posada

Universidad EAFIT  
Sistemas Operativos (ST0257)  
2024-1

24 de mayo de 2025

## Índice

<b>1. Resumen Ejecutivo</b>	<b>3</b>
<b>2. Arquitectura del Sistema</b>	<b>3</b>
2.1. Visión General . . . . .	3
2.2. Componentes Principales . . . . .	3
2.2.1. Módulo Core . . . . .	3
2.2.2. Módulo Storage . . . . .	3
<b>3. Implementación del Paralelismo con Dask</b>	<b>3</b>
3.1. Estrategia de Paralelización . . . . .	3
<b>4. Algoritmos de Compresión</b>	<b>4</b>
4.1. ZIP (DEFLATE) . . . . .	4
4.2. Alternativas Soportadas . . . . .	4
<b>5. Justificación Tecnológica</b>	<b>4</b>
5.1. Lenguaje: Python . . . . .	4
5.2. Bibliotecas Principales . . . . .	5
<b>6. Instrucciones de Uso</b>	<b>5</b>
6.1. Instalación . . . . .	5
6.2. Comandos Principales . . . . .	5
6.3. Verificación y Pruebas . . . . .	5
<b>7. Métricas de Rendimiento</b>	<b>5</b>
7.1. Compresión . . . . .	5
7.2. Transferencia . . . . .	5



# 1. Resumen Ejecutivo

Py-DataVault es un sistema de respaldo seguro que implementa compresión paralela, encriptación, y múltiples opciones de almacenamiento incluyendo fragmentación USB y respaldo en la nube. El sistema está diseñado para optimizar el rendimiento mediante el uso de procesamiento paralelo con Dask, mientras mantiene altos estándares de seguridad con encriptación AES-256.

## 2. Arquitectura del Sistema

### 2.1. Visión General

La arquitectura del sistema está organizada en módulos especializados que manejan diferentes aspectos del proceso de backup:

- **Core:** Módulos fundamentales para compresión, encriptación y restauración
- **Storage:** Gestión de almacenamiento local y en la nube
- **Interface:** CLI para interacción con el usuario

### 2.2. Componentes Principales

#### 2.2.1. Módulo Core

- `compressor.py`: Implementa algoritmos de compresión (ZIP, GZIP, BZIP2)
- `encryptor.py`: Maneja la encriptación AES-256 con PBKDF2
- `restorer.py`: Gestiona la restauración de backups
- `utils.py`: Utilidades generales y gestión de procesos

#### 2.2.2. Módulo Storage

- `cloud.py`: Integración con Google Drive
- `uploader.py`: Gestión de subidas
- `local.py`: Manejo de almacenamiento local
- `splitter.py`: Fragmentación y distribución USB

## 3. Implementación del Paralelismo con Dask

### 3.1. Estrategia de Paralelización

El sistema utiliza Dask para optimizar cuatro operaciones principales:

#### 1. Compresión Paralela

- Procesamiento simultáneo de múltiples archivos

- División de archivos grandes en chunks procesables
- Balanceo de carga automático

## **2. Encriptación Paralela**

- División de datos en bloques para encriptación simultánea
- Gestión de memoria optimizada
- Procesamiento paralelo de chunks

## **3. Transferencia de Datos**

- Operaciones I/O paralelas
- Buffering optimizado
- Gestión de concurrencia

## **4. Fragmentación USB**

- División paralela de archivos
- Escritura simultánea en múltiples dispositivos
- Verificación paralela de integridad

# **4. Algoritmos de Compresión**

## **4.1. ZIP (DEFLATE)**

- Algoritmo principal utilizado
- Basado en LZ77 y codificación Huffman
- Mejor balance entre velocidad y ratio de compresión
- Implementado mediante la biblioteca `zipfile` de Python

## **4.2. Alternativas Soportadas**

- **GZIP**: Para archivos individuales grandes
- **BZIP2**: Para máxima compresión cuando el tiempo no es crítico

# **5. Justificación Tecnológica**

## **5.1. Lenguaje: Python**

- Amplio soporte para procesamiento paralelo
- Excelente ecosistema de bibliotecas
- Facilidad de integración con servicios cloud
- Soporte multiplataforma robusto

## 5.2. Bibliotecas Principales

- **Dask:** Framework de computación paralela
- **PyDrive2:** API robusta para Google Drive
- **pycryptodome:** Implementación segura de AES
- **click:** CLI intuitiva y bien documentada

## 6. Instrucciones de Uso

### 6.1. Instalación

```
1 pip install -r requirements.txt
```

### 6.2. Comandos Principales

```
1 # Proceso completo de backup
2 python main.py full-backup-process \
3     --folders ./carpeta1,./carpeta2 \
4     --backup-name mi_backup \
5     --compression zip \
6     --encrypt \
7     --password mi_password \
8     --usb-paths usb1,usb2 \
9     --cloud
```

### 6.3. Verificación y Pruebas

```
1 python test_full_backup_process.py
```

## 7. Métricas de Rendimiento

### 7.1. Compresión

- Mejora de 2-4x en tiempo de procesamiento
- Uso eficiente de múltiples núcleos
- Escalabilidad lineal hasta 8 cores

### 7.2. Transferencia

- Optimización de operaciones I/O
- Reducción de tiempos de espera
- Mejor utilización del ancho de banda

## 8. Conclusiones

El sistema demuestra la efectividad de combinar procesamiento paralelo con Dask, algoritmos de compresión clásicos y prácticas modernas de seguridad. Los resultados muestran mejoras significativas en rendimiento mientras se mantienen altos estándares de seguridad y confiabilidad.