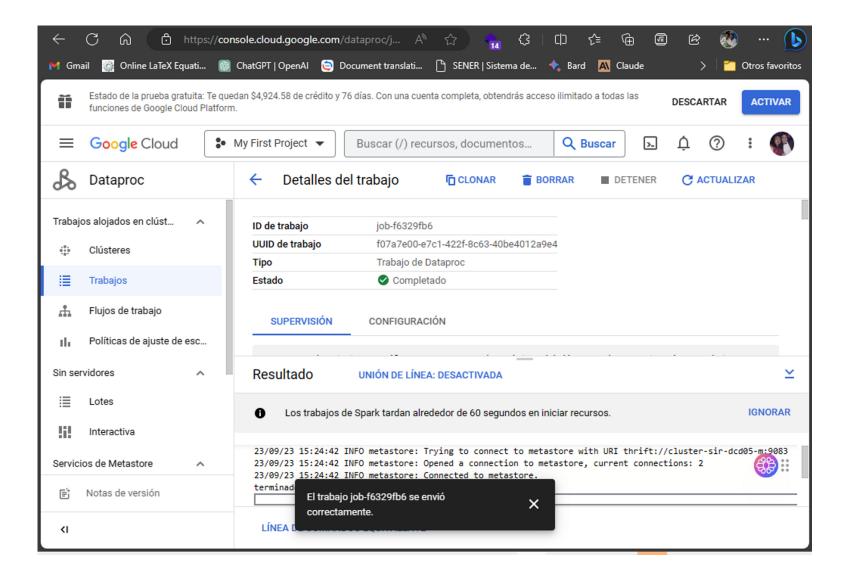**Diplomado en Ciencia de Datos UNAM
Módulo 13 Datos Masivos
Sesión 3 del 23 Septiembre de 2023**

**Sergio Ibarra Ramírez**
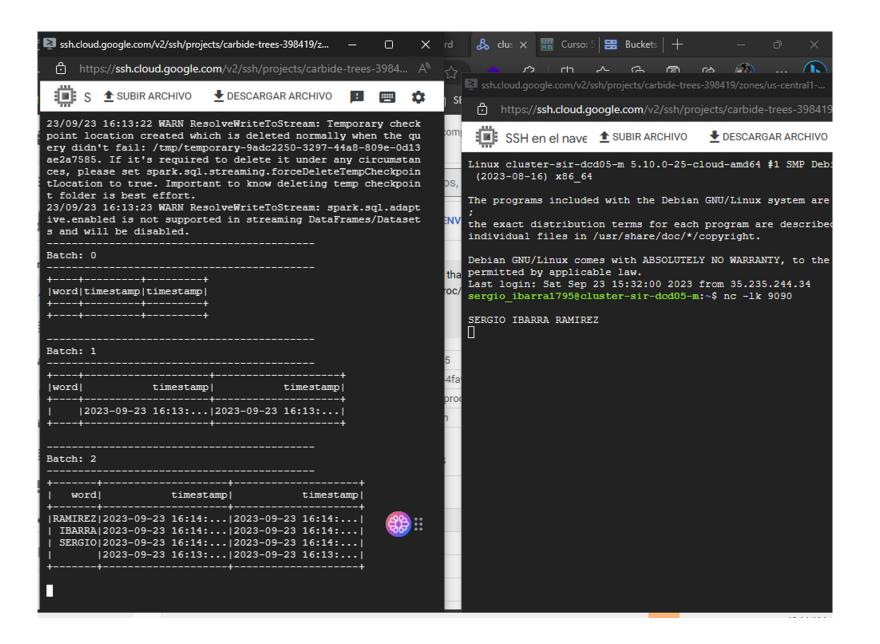
# Pantalla de Job de pysparkSQL ejecutado

# Pantalla 2

# Pantalla 3

# Pantalla 4