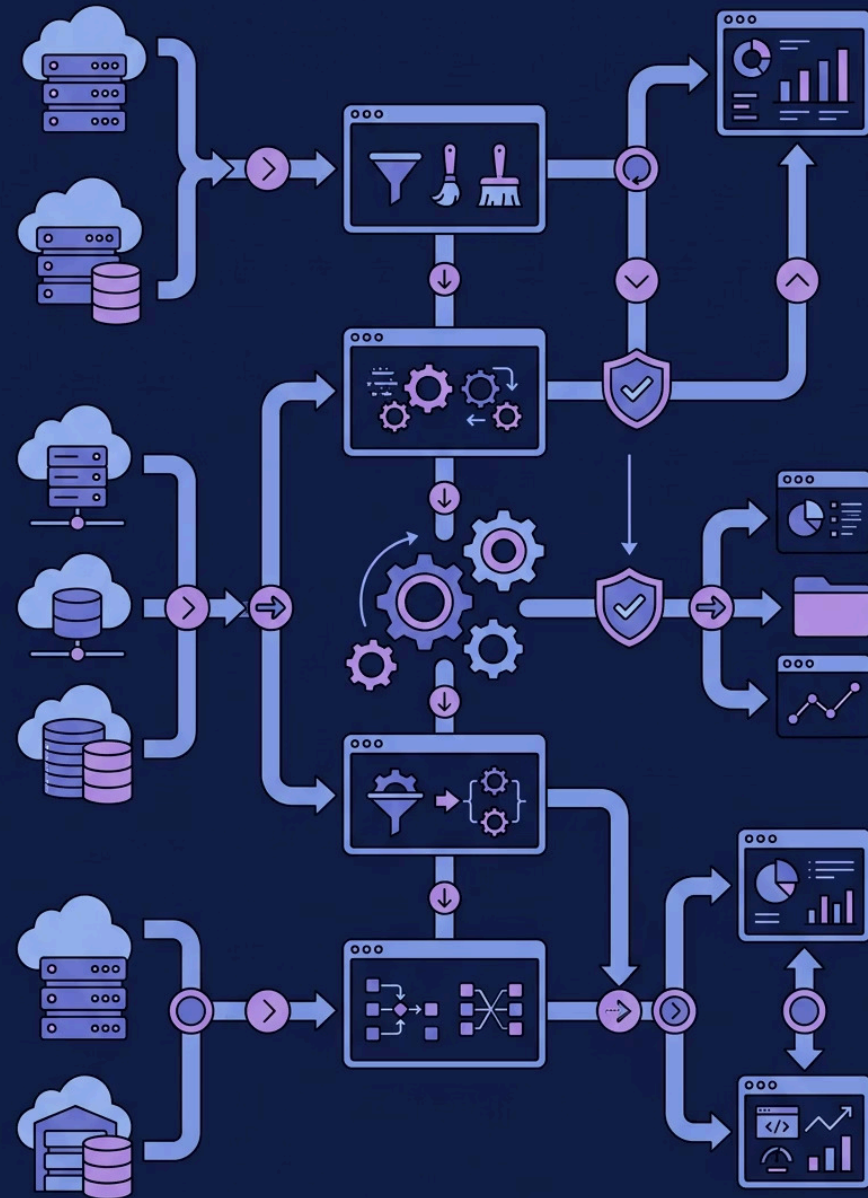




# CI/CD en datos: de DevOps a DataOps

Una aproximación técnica a la automatización de procesos de integración y despliegue continuo en datos

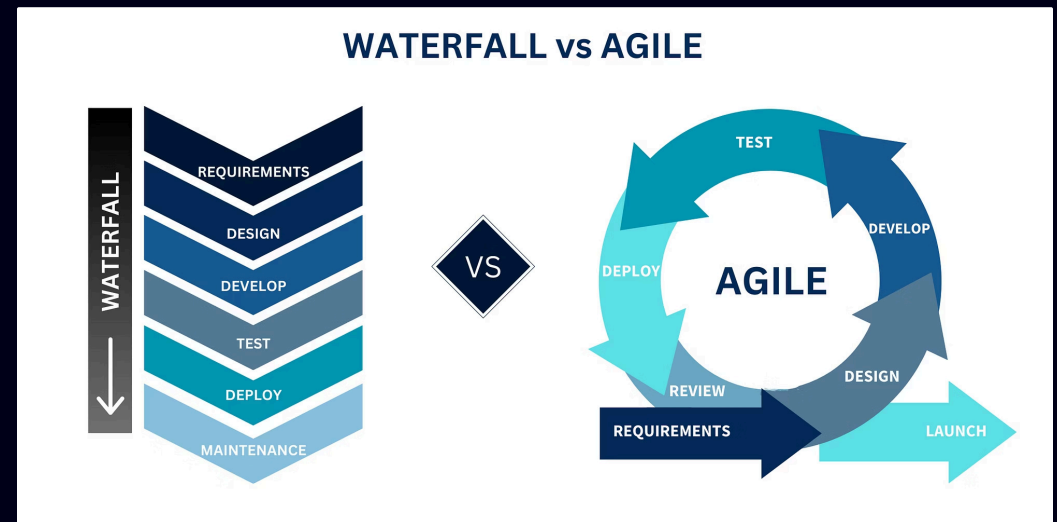


# Introducción y contexto histórico

El origen de CI/CD se remonta a principios de los 2000 como respuesta a los desafíos clásicos del desarrollo de software:

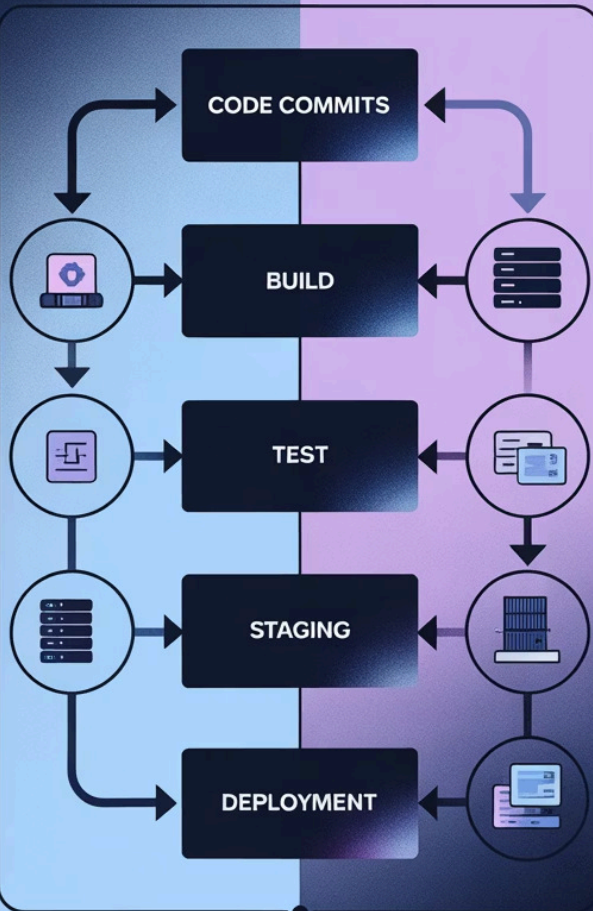
- Ciclos de desarrollo extensos (6-12 meses)
- Despliegues manuales propensos a errores
- Integraciones problemáticas ("infierno de integración")
- Retroalimentación tardía sobre problemas

Las metodologías ágiles como XP, Scrum y posteriormente DevOps impulsaron la adopción de prácticas CI/CD para reducir fricciones y acelerar entregas.



**Evolución:** De ciclos largos y entregas manuales a integración y despliegue continuos automatizados.

# Conceptos base de CI/CD



1

## Integración Continua (CI)

Proceso que integra cambios de código **frecuentemente** en un repositorio compartido, ejecutando pruebas automatizadas para validar la calidad.

Objetivo: detectar problemas tempranamente.

2

## Entrega Continua (Continuous Delivery)

Extiende CI para mantener la aplicación **siempre lista** para ser desplegada a producción, aunque requiere aprobación manual final.

3

## Despliegue Continuo (Continuous Deployment)

Lleva CD un paso más allá: cada cambio que pasa todas las pruebas se despliega **automáticamente** a producción sin intervención humana.

El pipeline típico de CI/CD en software combina estos elementos para crear un flujo de trabajo fluido desde el desarrollo hasta la implementación.

# El desafío del mundo de datos

La adopción de CI/CD en entornos de datos ha sido más lenta debido a varias complejidades:

- En software testearnos principalmente **código**; en datos, debemos validar tanto el **código como los datos**
- Los datos son dinámicos, cambiantes y mucho más voluminosos
- Los entornos de desarrollo y producción son difíciles de mantener sincronizados
- Las pruebas requieren datos representativos pero seguros (sin exponer información sensible)

**⚠ Problemas sin CI/CD en datos:** pipelines rotos, dashboards inconsistentes, pérdida de confianza en los datos y decisiones empresariales comprometidas.



Mismo código, datos distintos → métricas y conclusiones diferentes

# Adaptación de CI/CD a datos (DataOps)

**DataOps** es la extensión natural de DevOps al mundo de datos, enfocada en:

- Automatizar el ciclo de vida completo de datos
- Mejorar la colaboración entre equipos
- Garantizar calidad y disponibilidad de datos
- Acelerar entrega de valor

Esta adaptación requiere cambios fundamentales en cómo abordamos:

- Testing específico para datos
- Versionado de entornos y datasets
- Validaciones continuas de calidad



El ciclo DataOps enfatiza la **calidad de los datos** como pilar fundamental del proceso.

# Tipos de pruebas en datos



## Unit Tests

Validan la lógica de transformaciones con datasets pequeños y controlados. Ejemplo: comprobar que una función de cálculo de descuento opere correctamente.



## Tests de Integridad

Verifican aspectos estructurales: claves primarias, relaciones, valores nulos, duplicados, tipos de datos, etc.



## Tests de calidad de datos

Aseguran que los datos cumplan reglas de dominio: outliers, rangos esperados, relaciones lógicas entre entidades, etc.

Testing de Software	Testing de Datos
Comportamiento predecible	Resultados dependen de los datos de entrada
Entradas controladas	Datos cambiantes y voluminosos
Resultados deterministas	Necesidad de validar patrones y distribuciones
Pruebas rápidas	Pruebas potencialmente costosas y lentas

La combinación de estos tipos de pruebas permite crear un enfoque integral de validación para pipelines de datos.

# Integración Continua en Ingeniería de Datos

La Integración Continua (CI) en el contexto de datos automatiza la verificación y validación temprana de los cambios de código y datos. Su objetivo es asegurar la calidad y fiabilidad de los pipelines de datos desde las primeras etapas del desarrollo, minimizando errores y acelerando la entrega.



## Linting y Formateo

Análisis estático del código (SQL, Python, etc.) para verificar estilo, sintaxis y adherencia a estándares. Ayuda a mantener la consistencia y legibilidad del código.



## Smoke Tests

Pruebas rápidas que validan la conectividad a las fuentes de datos y la ejecución básica del pipeline. Aseguran que los componentes principales están operativos.



## Pruebas Unitarias

Verifican la lógica de transformaciones y funciones individuales con datasets controlados y pequeños. Aseguran que cada parte del código de datos produce el resultado esperado.

# Pre-commits: calidad antes de la integración

Los **pre-commits** son ganchos locales que ejecutan validaciones automáticas y scripts de calidad antes de cada `git commit`. Actúan como una primera línea de defensa para garantizar que solo el código de alta calidad ingrese al repositorio.

## Beneficios clave



### Detección Temprana

Identifican problemas, errores de sintaxis o violaciones de estilo en tu máquina local, antes de enviar los cambios al repositorio central.



### Ciclos de Feedback Rápidos

Reducen drásticamente el tiempo de retroalimentación, evitando que la CI tenga que procesar cambios que ya están rotos o mal formateados.

## Ejemplos de pre-commits en proyectos de datos

- **Formateo automático:** Herramientas como `black` (Python) o `isort` para organizar imports.
- **Linting:** `flake8` para Python, `sqlfluff` para SQL (especialmente útil con dbt).
- **Validación de esquema:** Chequeos para asegurar que los modelos de dbt o esquemas de datos compilen correctamente.
- **Chequeos de seguridad:** `detect-secrets` para prevenir la subida accidental de credenciales o información sensible.



Desarrollador Local



Pre-commits



Pipeline CI



Despliegue CD

# Implementación de flujos CI/CD para datos

**Flujo de trabajo DataOps:** Desde el desarrollo local hasta la implementación en producción, con validaciones en cada etapa.

## Elementos esenciales:

1

Repositorio Git como fuente única de verdad

2

Ambientes separados (desarrollo, staging, producción)

3

Automatización de pruebas y despliegues

4

Controles de calidad obligatorios

## Beneficios clave:

✓ Mayor confianza en los datos

✓ Reducción de errores en producción

✓ Equipos más autónomos y veloces

✓ Mejora continua del pipeline de datos