# NETAPP DATAOPS TOOL KIT

**for data management**

**Dr. Didier Gava**

EMEA & LATAM Senior Solution architect

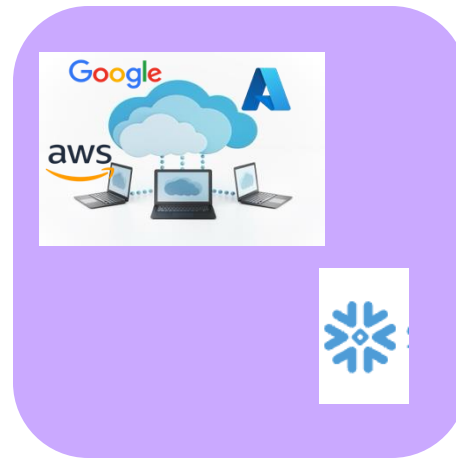September 2024

# "Data Scientist"/"Data Engineer" Journey
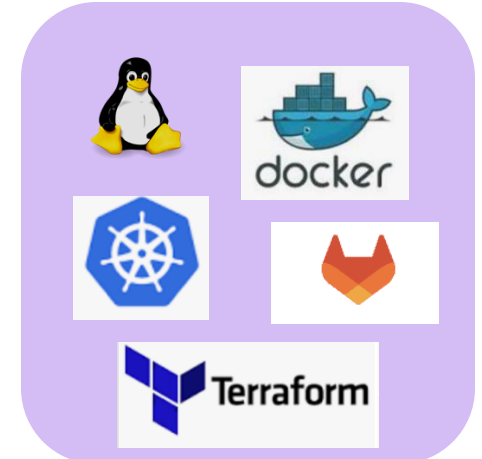


Base competencies

Data pipeline creation
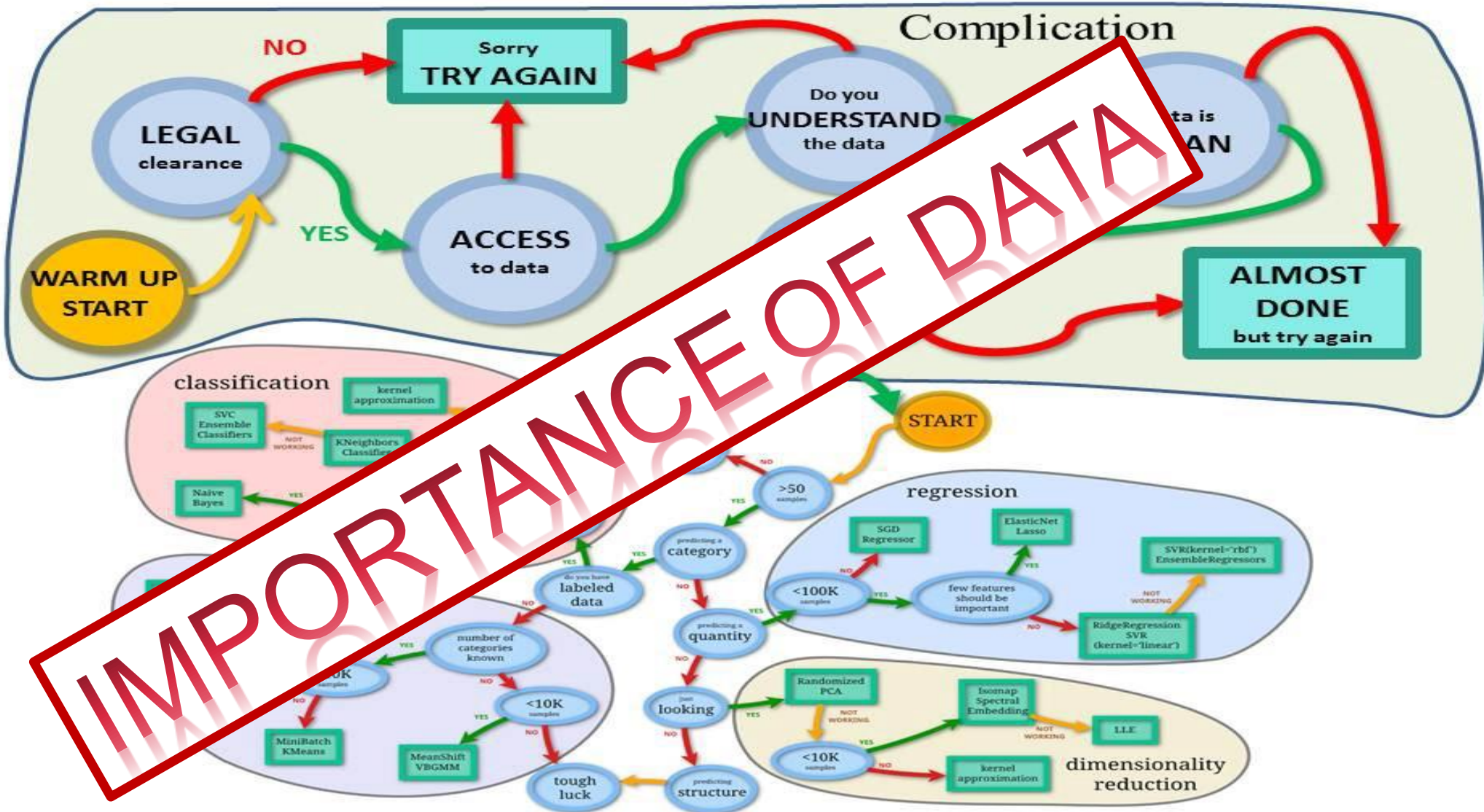
Big Data Platforms

Analytics / Visualisation

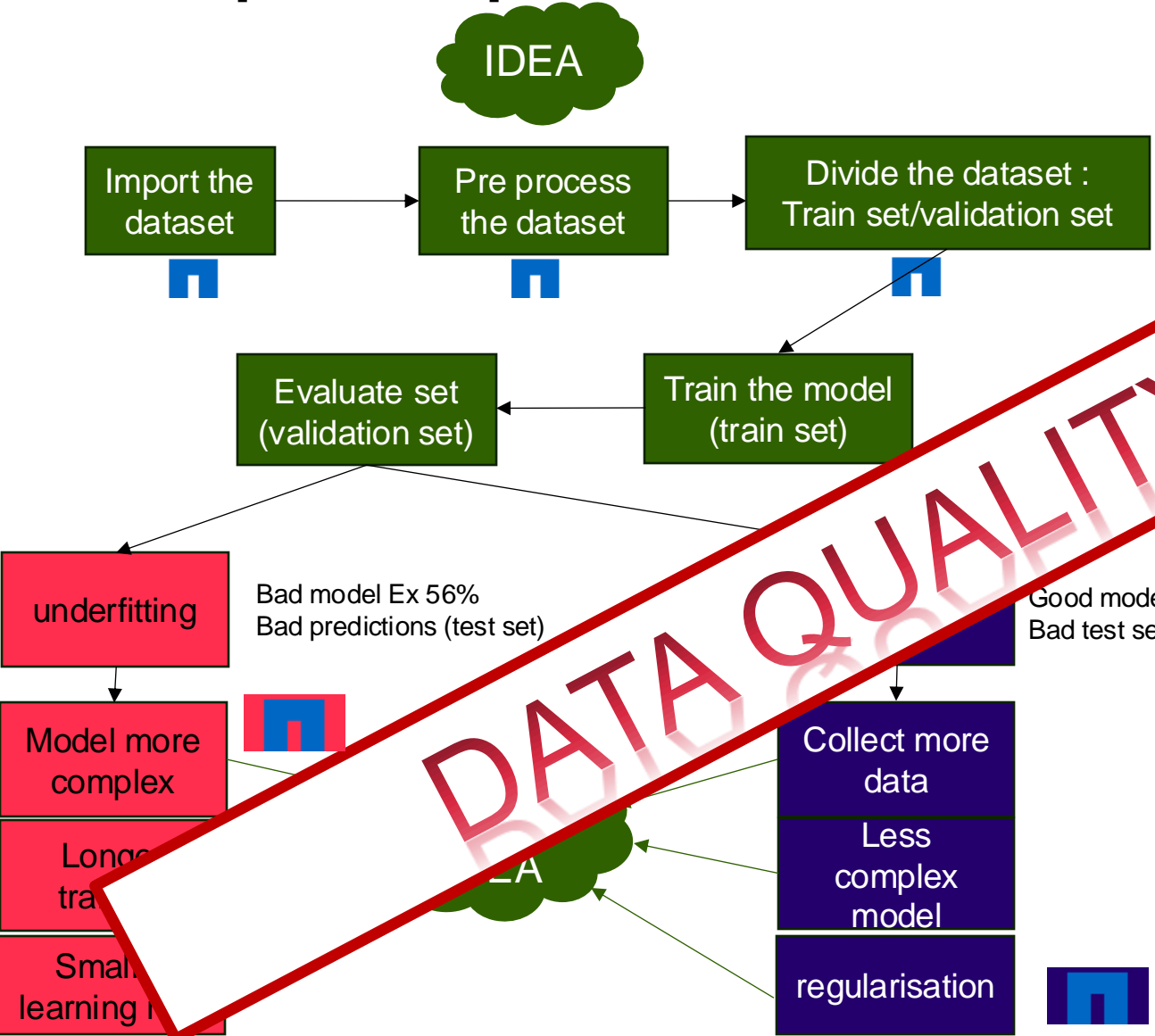Devops Infrastructure

# Data Scientist Data

# Development steps in ML
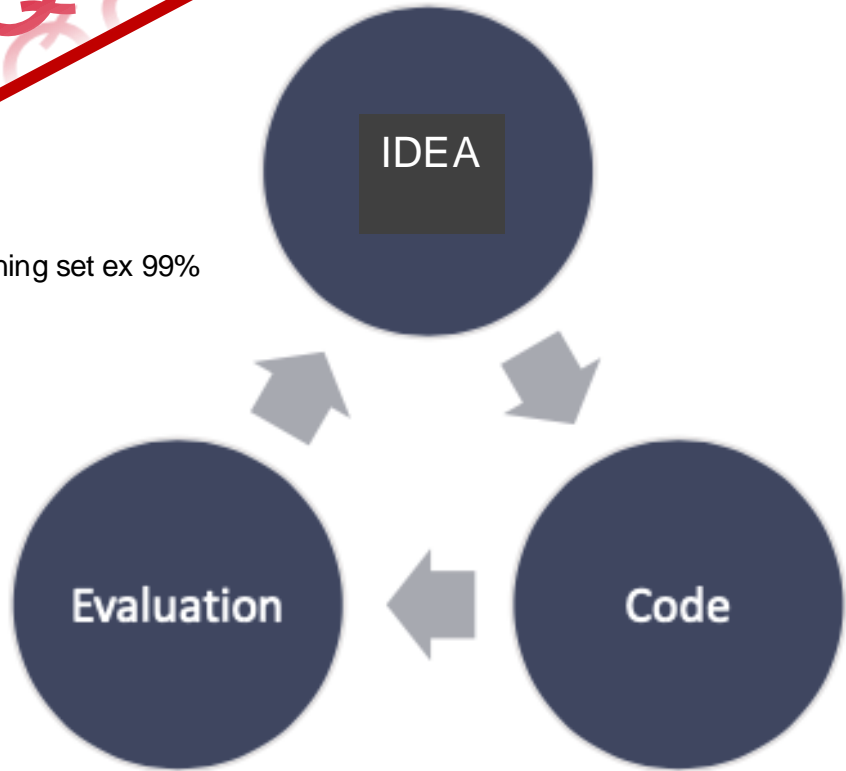
IDEA

Import the dataset → Pre process the dataset → Divide the dataset : Train set/validation set

Evaluate set (validation set) ← Train the model (train set)

underfitting

Bad model Ex 56%
Bad predictions (test set)

Model more complex

Longer tra...

Smal... learning r...

Good model on the training set ex 99%
Bad test set predictions

Collect more data

Less complex model

regularisation

- Developing a machine lea... done the first time.

- We often sta... a simple **and** quick m... en we analyze if we ha... as and we try a new **idea** ...blems encountered, etc.

DATA QUALITY/QUANTITY

IDEA

Evaluation ← Code
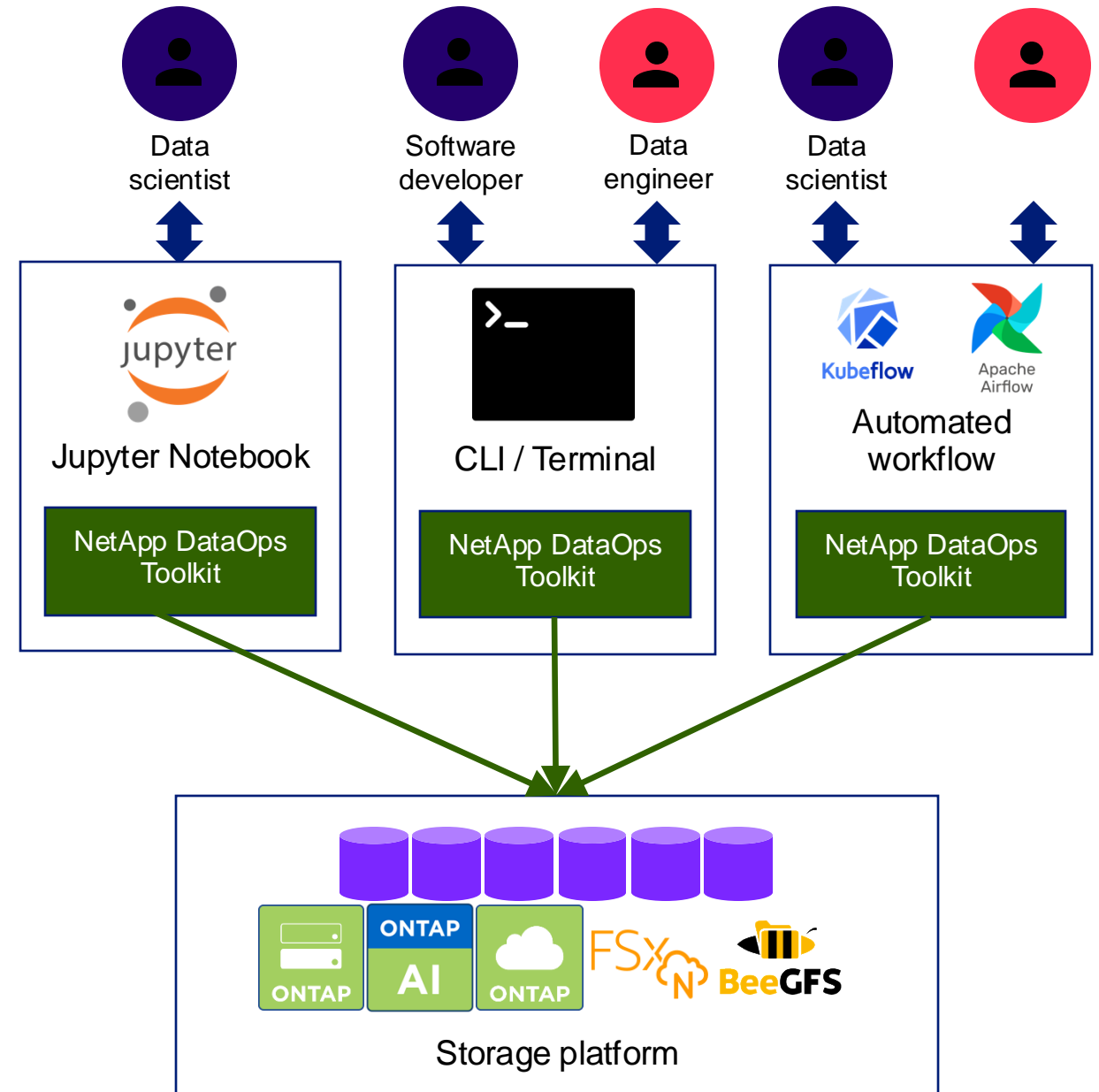
# Training accurate models

**Rapid experimentation is necessary**

# NetApp DataOps Toolkit
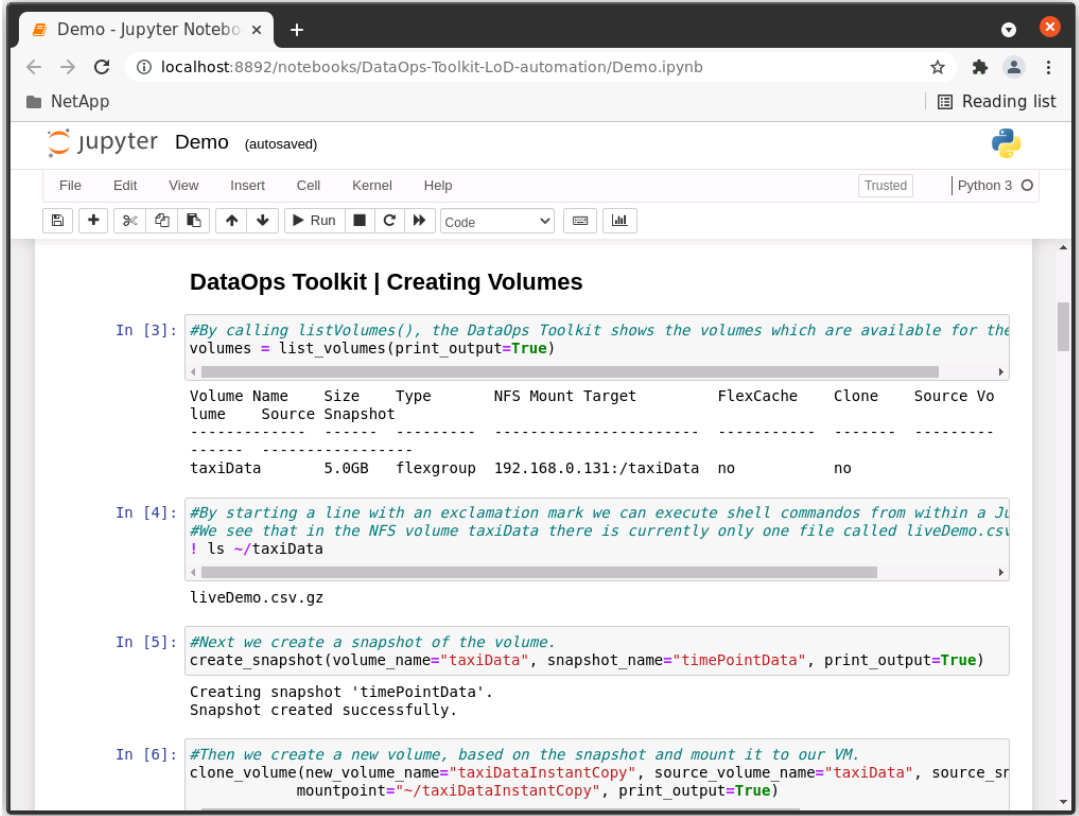
## Simplifying AI data management

- NetApp's industry-leading, multitenant data management capabilities
  - Traditional toolkit: Supports ONTAP (AFF, FAS, FSx, Cloud, Select)
  - Kubernetes toolkit: Supports ONTAP (AFF, FAS, FSx, Cloud, Select), Azure NetApp Files (ANF), Cloud Volumes Service (CVS), and BeeGFS (limited)

- Simple, easy-to-use interfaces; designed for data scientists and data engineers
  - CLI utility
  - Importable library of Python functions

- Provides access to advanced features that would normally require help from storage admin

- Key capabilities
  - Rapidly provision a new data volume
  - Near-instantaneously clone a data volume
  - Snapshot a data volume for traceability/versioning
  - Trigger data sync

Data scientist — Jupyter Notebook — NetApp DataOps Toolkit
Software developer / Data engineer — CLI / Terminal — NetApp DataOps Toolkit
Data scientist — Automated workflow (Kubeflow, Apache Airflow) — NetApp DataOps Toolkit
Storage platform (ONTAP, ONTAP AI, Cloud ONTAP, FSx N, BeeGFS)

# NetApp DataOps Toolkit (Traditional)

**Simplify access to NetApp solutions from Data Science environments**



Jupyter Notebook:



Key Functions:


Cloning of volumes


Deleting volumes
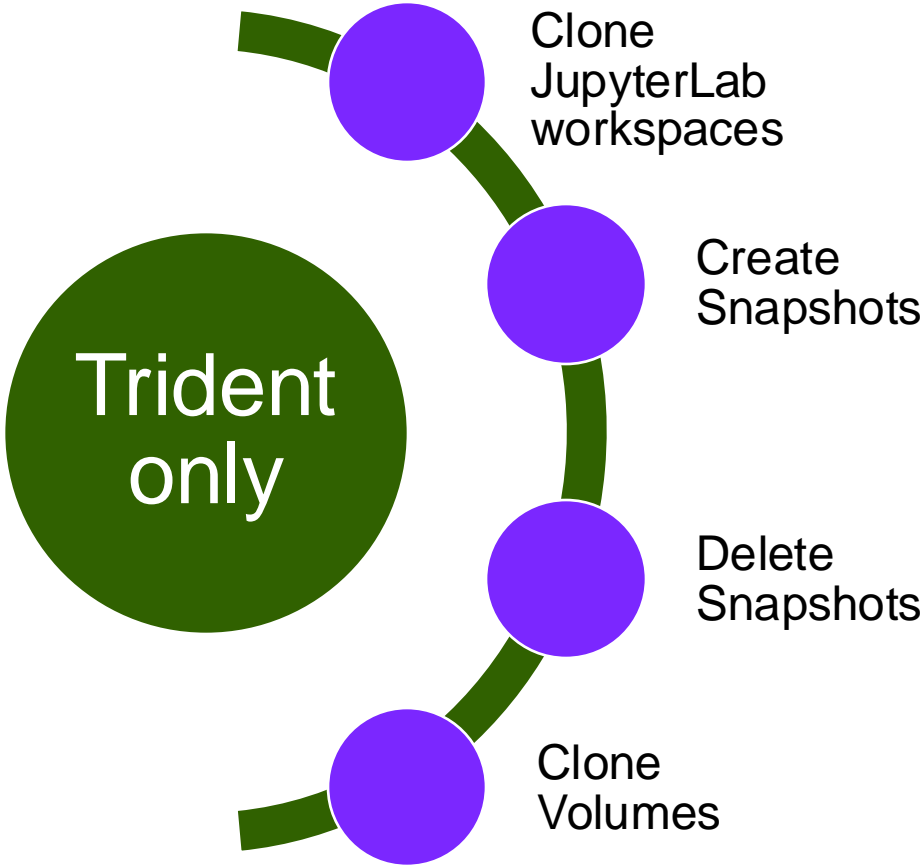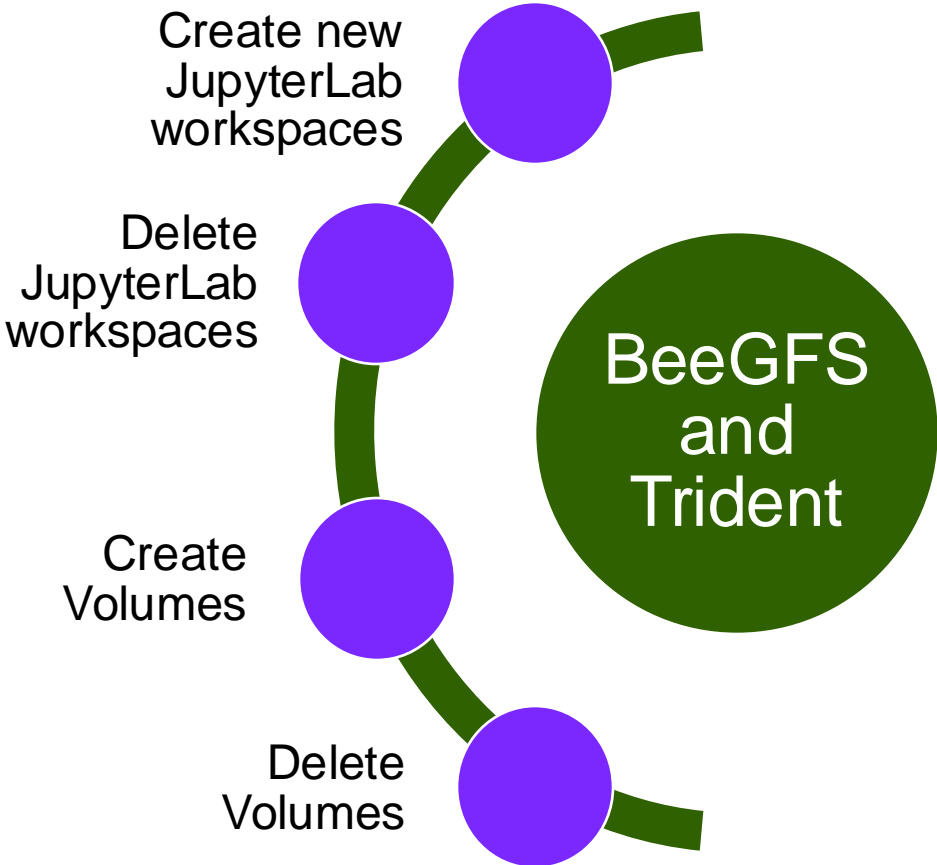

Creating Snapshots


Pull Content from S3


Push files to S3

# NetApp DataOps Toolkit (Kubernetes)

**Simplified data management in Kubernetes environments**

# Quick and easy installation and config

**Get started in seconds**

- NetApp DataOps Toolkit for Kubernetes – 1 step:

  ```
  1. pip install netapp-dataops-k8s
  ```

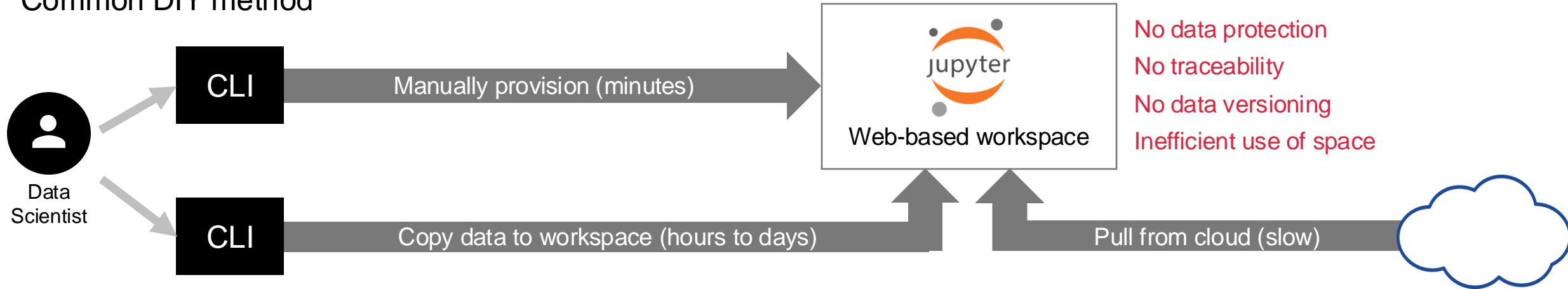- NetApp DataOps Toolkit for Traditional Environments – 2 steps

  ```
  1. pip install netapp-dataops-traditional
  2. netapp_dataops_cli.py config
  ```
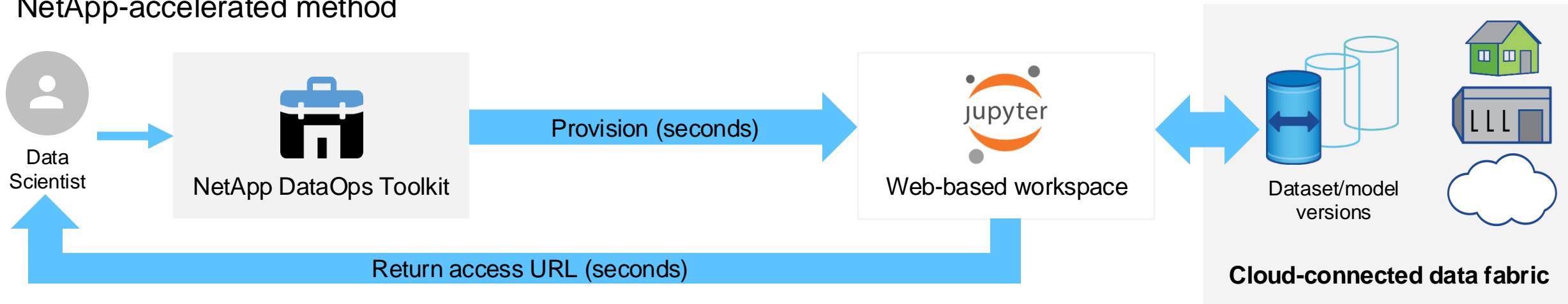
**~~Get up and running in under…~~ 2 minutes**

# Self-service data science workspace creation

## Common DIY method



Data Scientist

CLI → Manually provision (minutes) → Web-based workspace (jupyter)

CLI → Copy data to workspace (hours to days)

Pull from cloud (slow)

No data protection
No traceability
No data versioning
Inefficient use of space

## NetApp-accelerated method



Data Scientist → NetApp DataOps Toolkit → Provision (seconds) → Web-based workspace (jupyter) ↔ Dataset/model versions

Return access URL (seconds)

**Cloud-connected data fabric**

# Provision Volumes

# Dataset-to-model traceability

## Common DIY method



**Data Scientist**

🚫

**Workspace** — jupyter

CLI → Manual copy (hours/days) → **Dataset archive**

## NetApp-accelerated method

**Workflow driven:** part of automated workflow
NetApp DataOps Toolkit
*Simple Python and CLI interfaces*

**User driven / interactive:** user self-service
NetApp DataOps Toolkit
*Simple Python and CLI interfaces*

# Cold data tiering

Without NetApp



Cold data consumes
valuable storage space
in high-performance tier

High-performance tier

With NetApp FlexCache and FabricPool



Automatic cold data tiering

**FlexCache**

Automatic cold data tiering

**FabricPool**

High-performance tier

ONTAP    ONTAP AI    ONTAP    FSx N

StorageGRID®    S3

# Data movement and sync

## Common DIY method



## NetApp-accelerated method

**KEY TAKEAWAYS**

**01**

**Increase of Speed** by creating clones within seconds and mounting volumes within Jupyter notebook

**02**

**Access to Data** easily as they are "located" in servers

**03**

**Traceability of Experiments** by creating clones and/ or snapshots of working environments

THANK YOU