

# Versatile Data Kit



# Preview

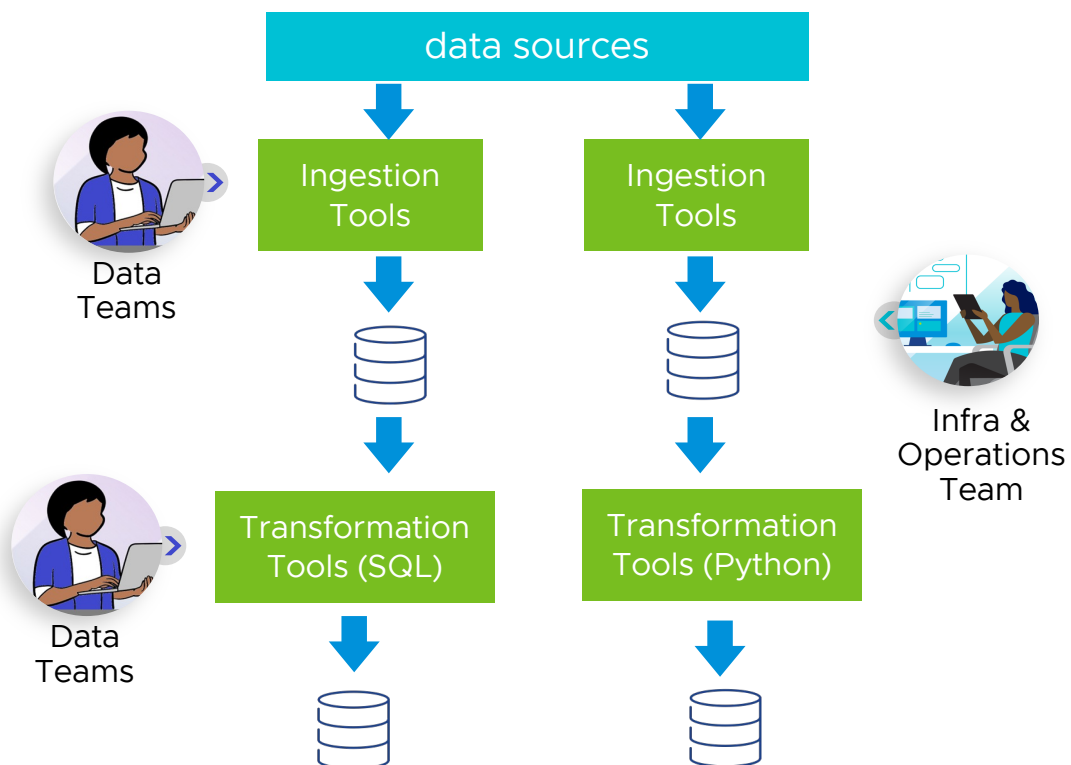
Enable everyone to focus on work that require their core skills



*Github repo*

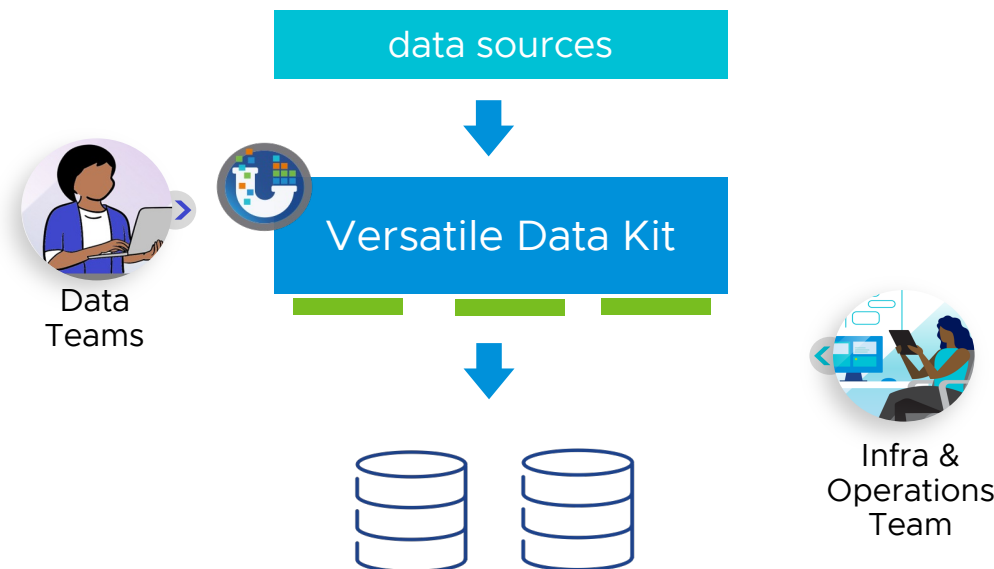
## Without Versatile Data Kit

- Fragmented Infrastructure
- Organization Silos
- Infra/Ops & Data Team tension



## With Versatile Data Kit

- Easier maintenance using unified managed interface
- Self-service, fully automated data teams
- Improved collaboration



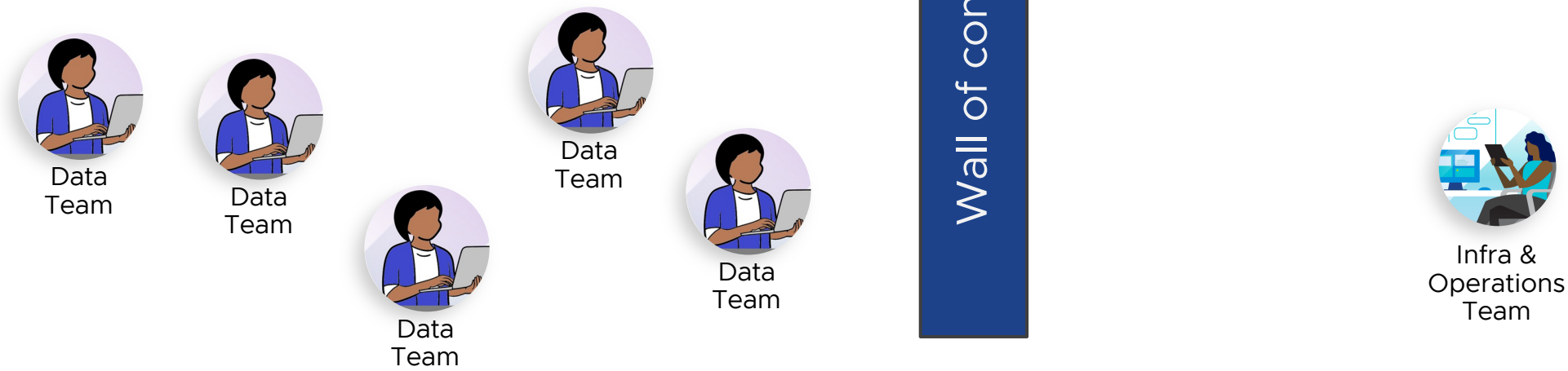
# What are the DataOps challenges?

Infra/Ops Team and Data Teams tension and conflict

*Inefficient Operations    Stalled development.*

Domain knowledge  
Implement business logic  
Optimizes for agility and speed

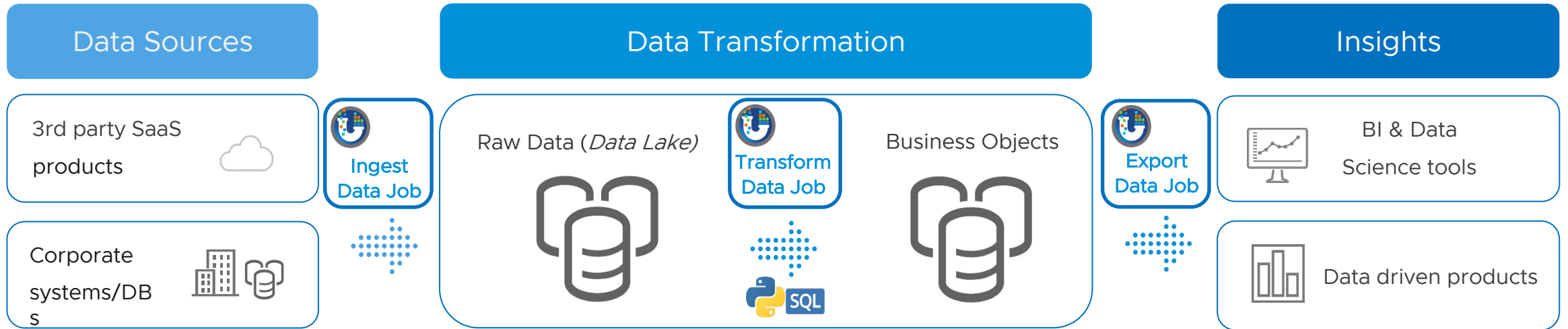
DevOps & Infrastructure knowledge  
Maintain infrastructure  
Optimizes reliability, availability and security



*Blurred lines of responsibility*

# Versatile Data Kit

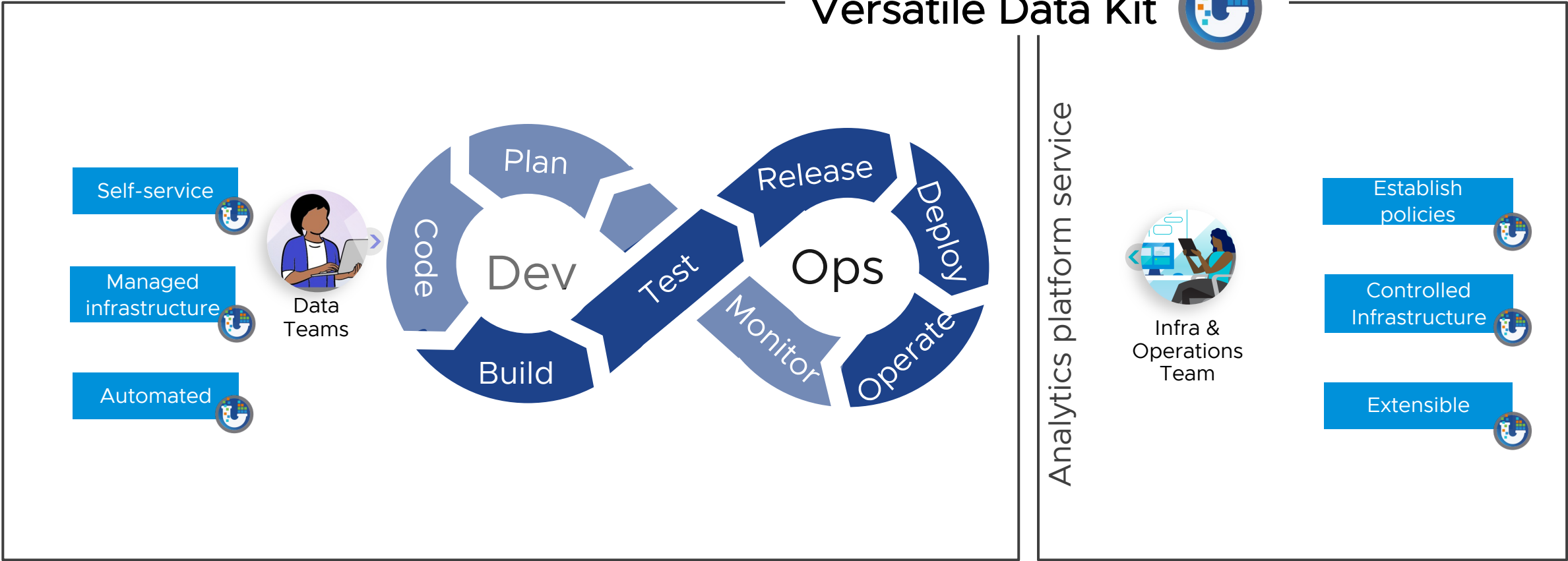
Data lifecycle (BI Journey) and where VDK fits in



# DevOps for Data as a Service

Adopt and adapt DevOps to deliver value from data efficiently

## Versatile Data Kit

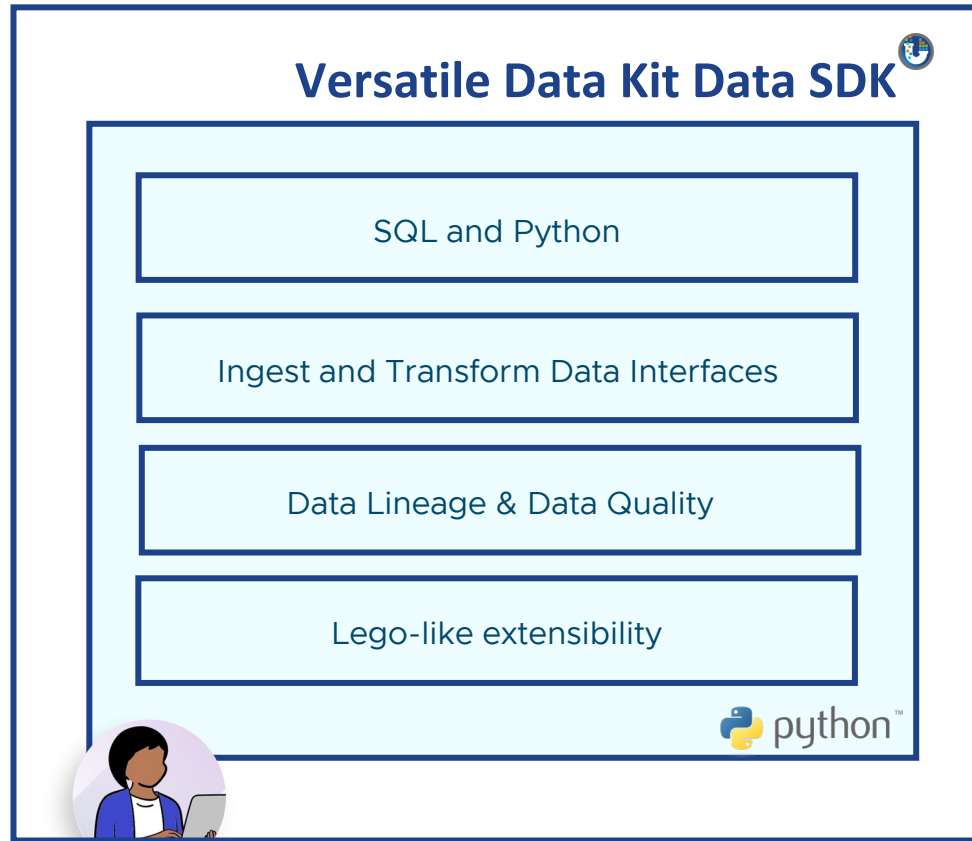




# Components

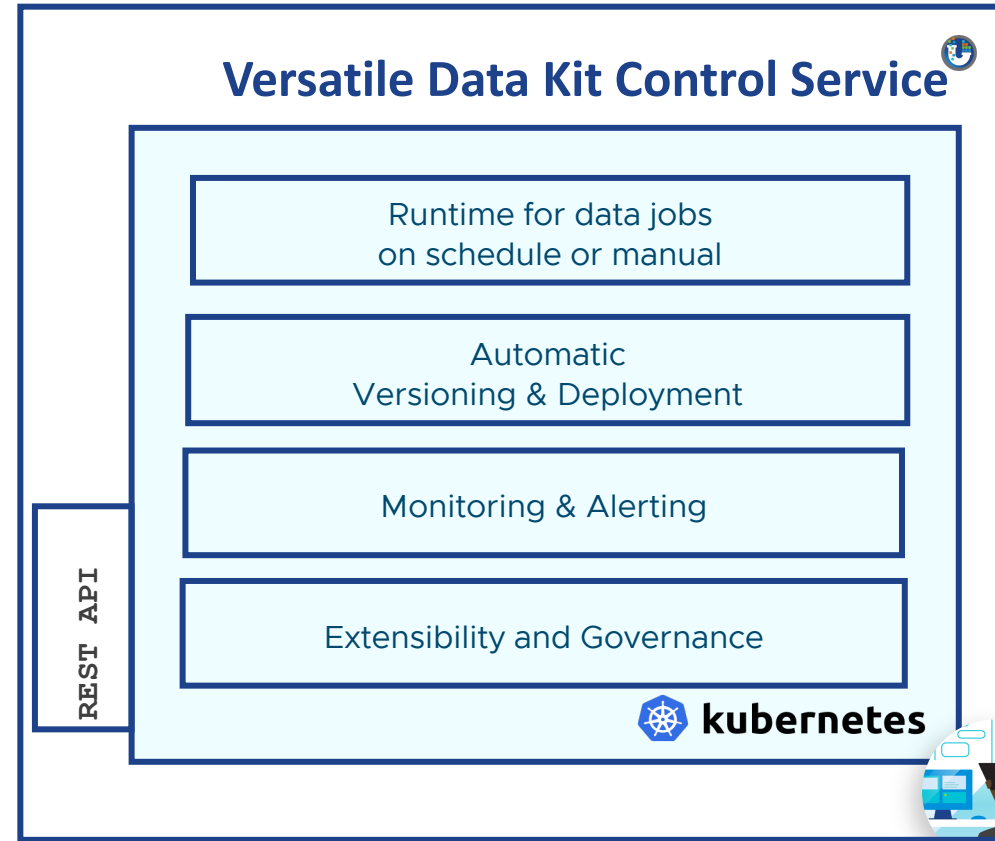
What does it take to run Versatile Data Kit and start deploying data jobs?

*Automate and abstract the Data Journey*



Data Teams

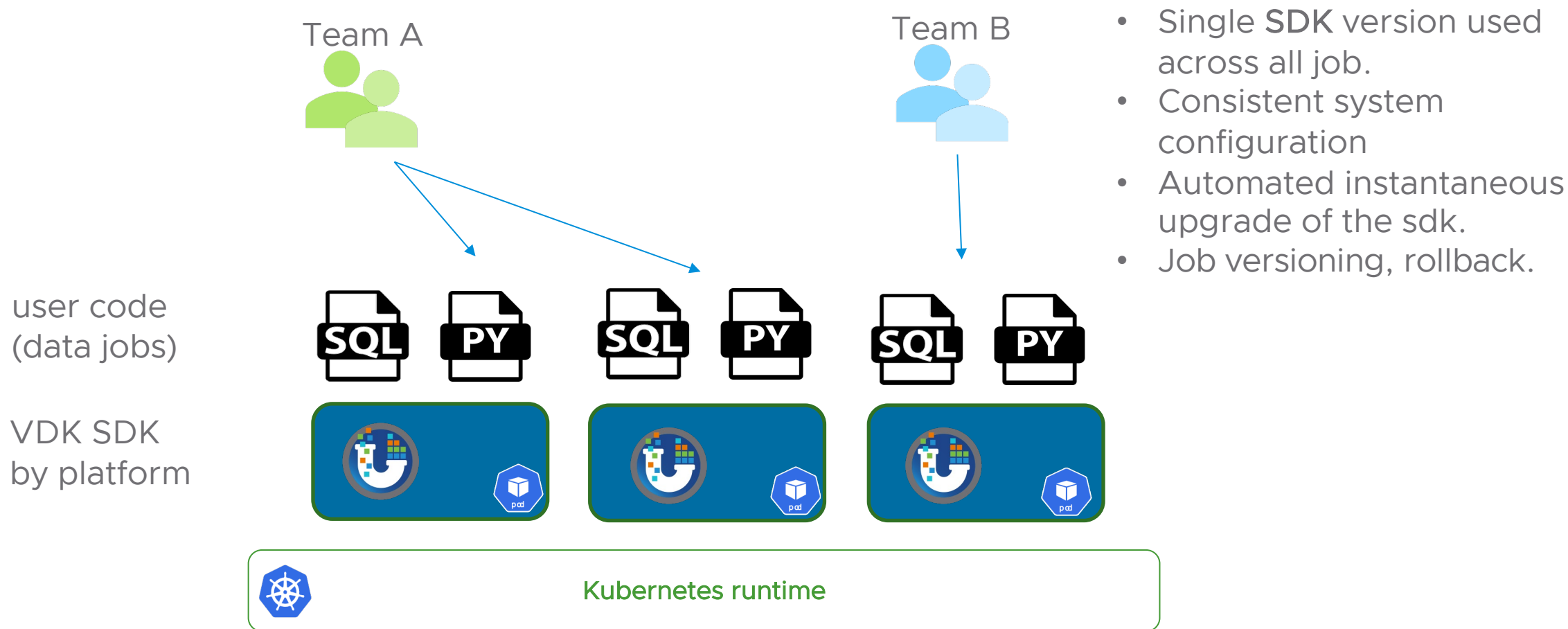
*Automate and abstract the DevOps Cycle*



Infra & Operations Team

# VDK SDK and VDK Runtime (Control Service)

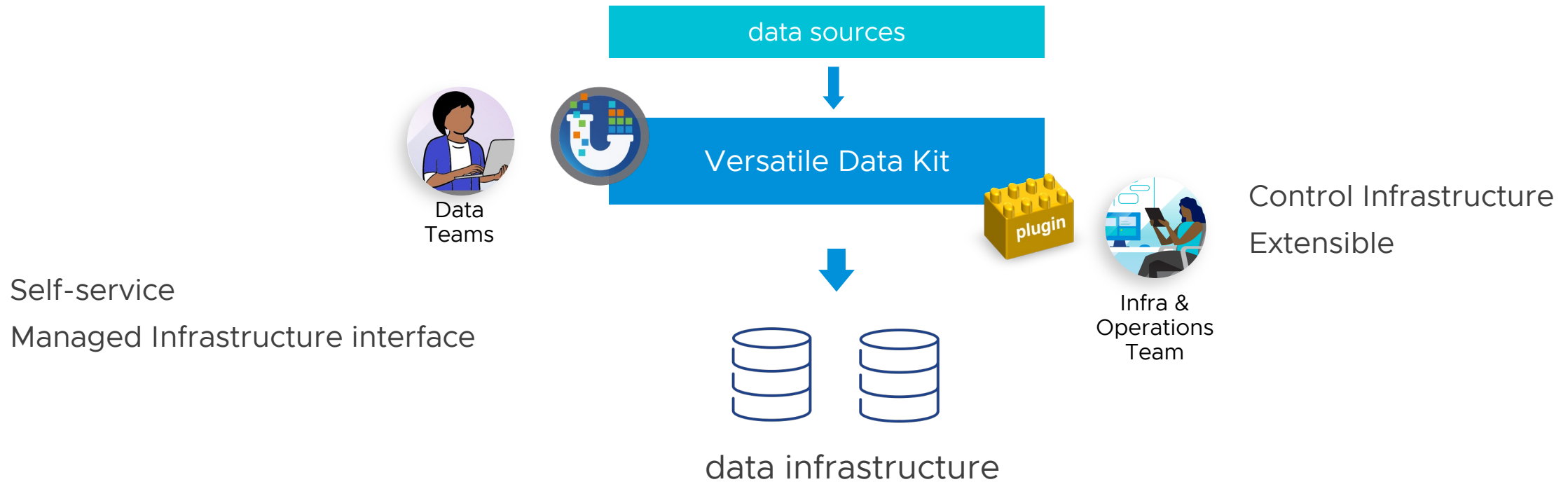
Data library and Job runner



# Automate and Abstract the Data journey

Simplify and hide complexity of data infrastructure for Data teams

Give power to establish best Infrastructure practices



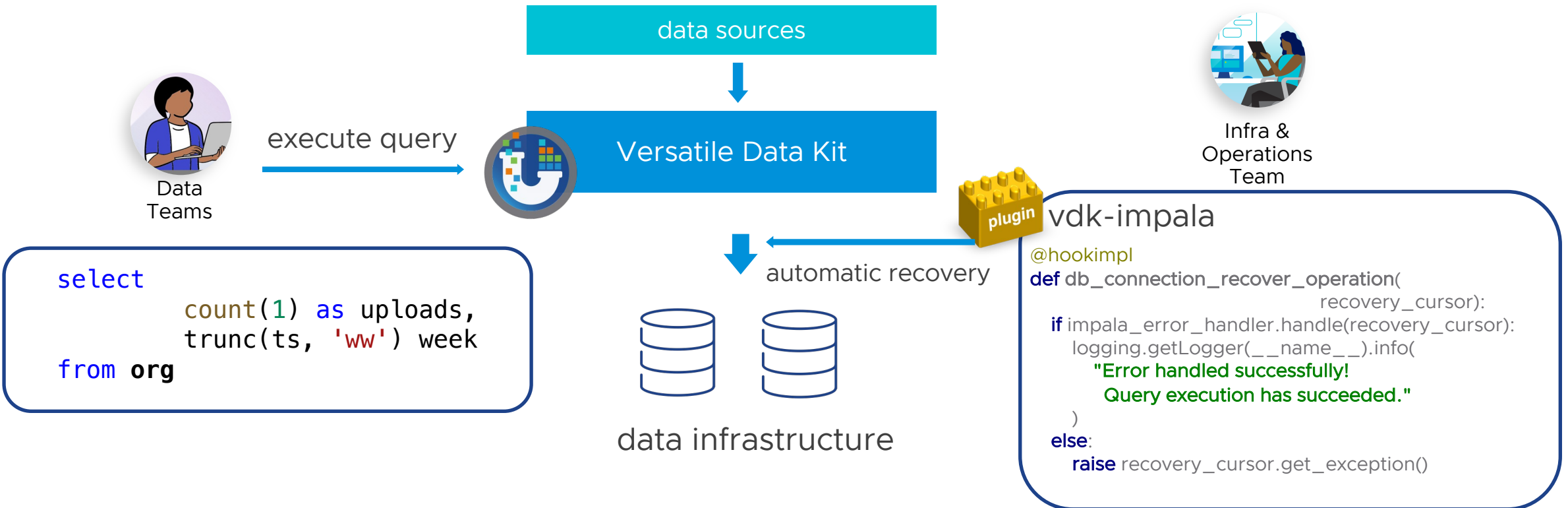


# Quick example: vdk-impala, vdk-trino

Simplify and hide complexity of data infrastructure for Data teams

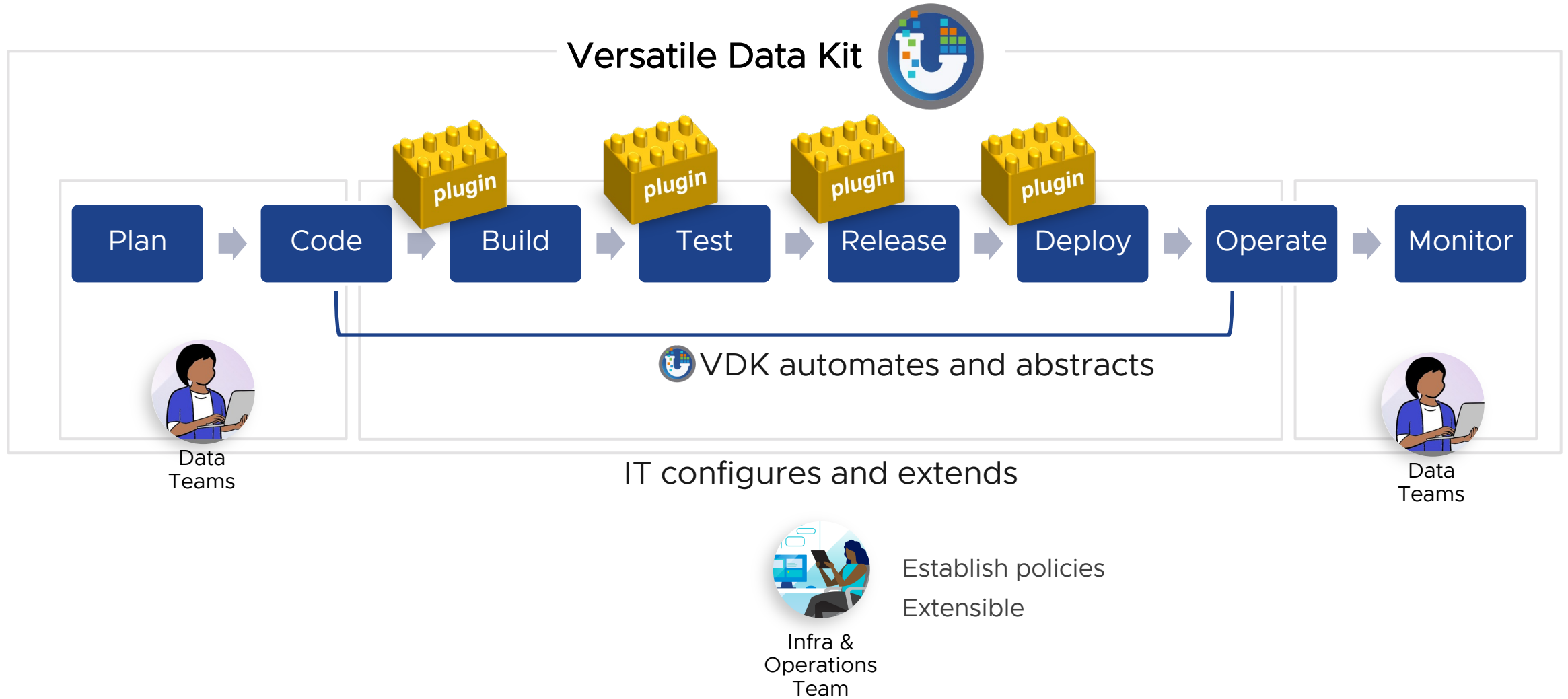
Give power to establish best Infrastructure practices

<https://github.com/vmware/versatile-data-kit/tree/main/projects/vdk-plugins/vdk-impala>



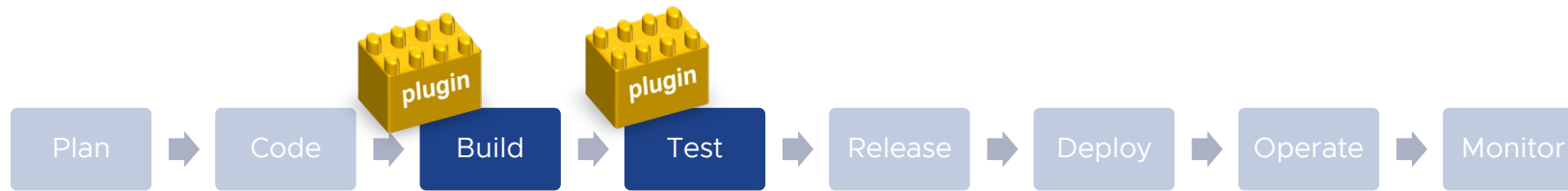
# Automate and Abstract the Development Process

Give power to Operators to establish best dev practices; Ease data job development



## Quick example: DevOps Plugin

Establish standard system tests and security hardening



`helm install --set job-builder=my-job-builder-image`

```
2  >> FROM versatiledatakit/job-builder
3      Xxx    -/tree/main
4      # Run system test before accepting the new job code
5      RUN pytest system_test.py || die 'Failed system test'
6
7      # Remove execution privileges from files during container build
8      RUN chmod -R -x $job_name/
9
```

# Data Journey : Ingestion

## Data Sources

3rd party SaaS  
products



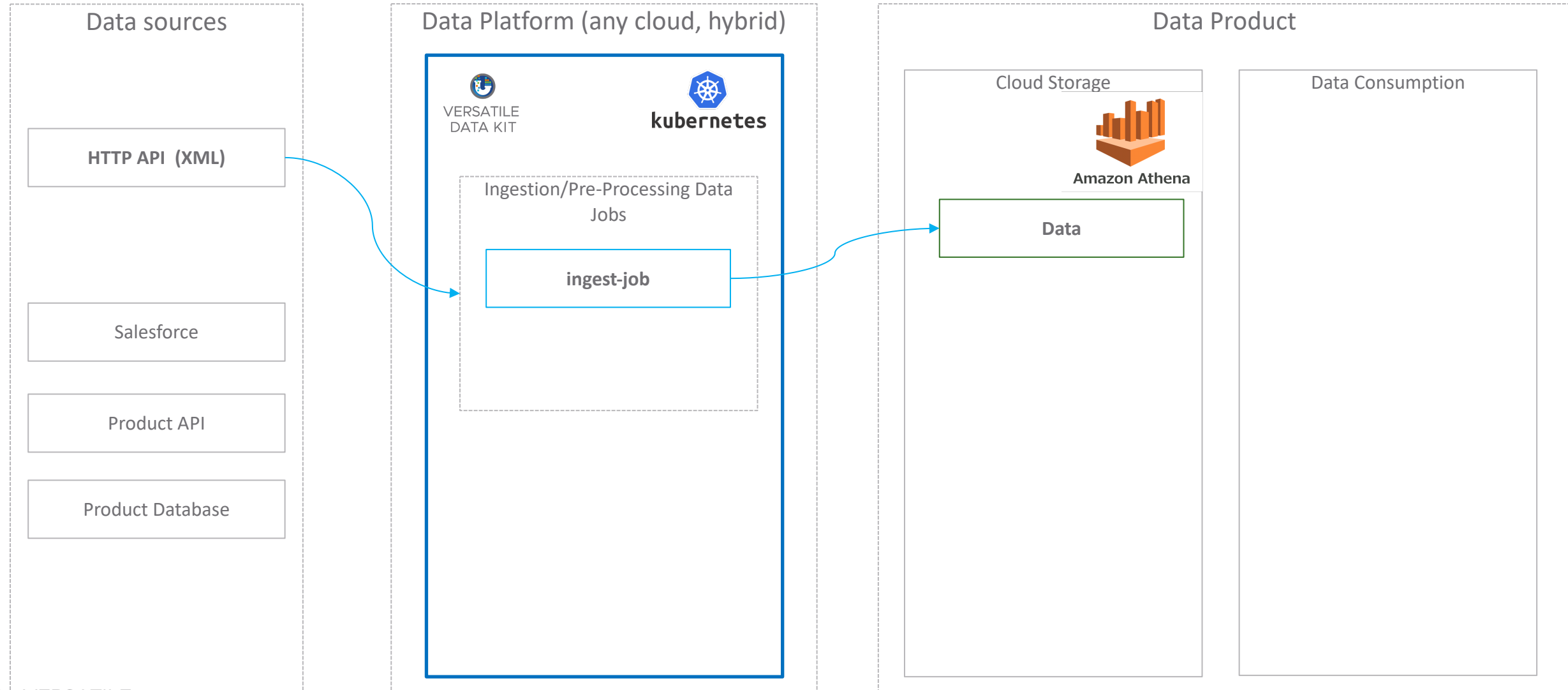
Ingest  
Data Job

Corporate  
systems/DB  
s



Data Teams

# Ingestion



# Ingestion Job

vdk run ingest-job

User code

```
job_input.send_object_for_ingestion(  
    payload=xmltodict.parse(response),  
    destination_table="rates"  
)
```

response =

```
<Rates>  
  <Rate>  
    <No>C/NBP/2</No>  
    <Date>2020-10-05</Date>  
    <Bid>4.347</Bid>  
    <Ask>4.431</Ask>  
  </Rate>  
  ....  
</Rates>
```



Versatile Data Kit



No	Date	Bid	Ask
C/NBP/2	5/10/20	4.347	4.431
C/NBP/1	5/10/20	4.234	4.532



trino





# Data Journey : Transformation



Data Teams

## Data Transformation

Raw Data (*Data Lake*)



Business Objects

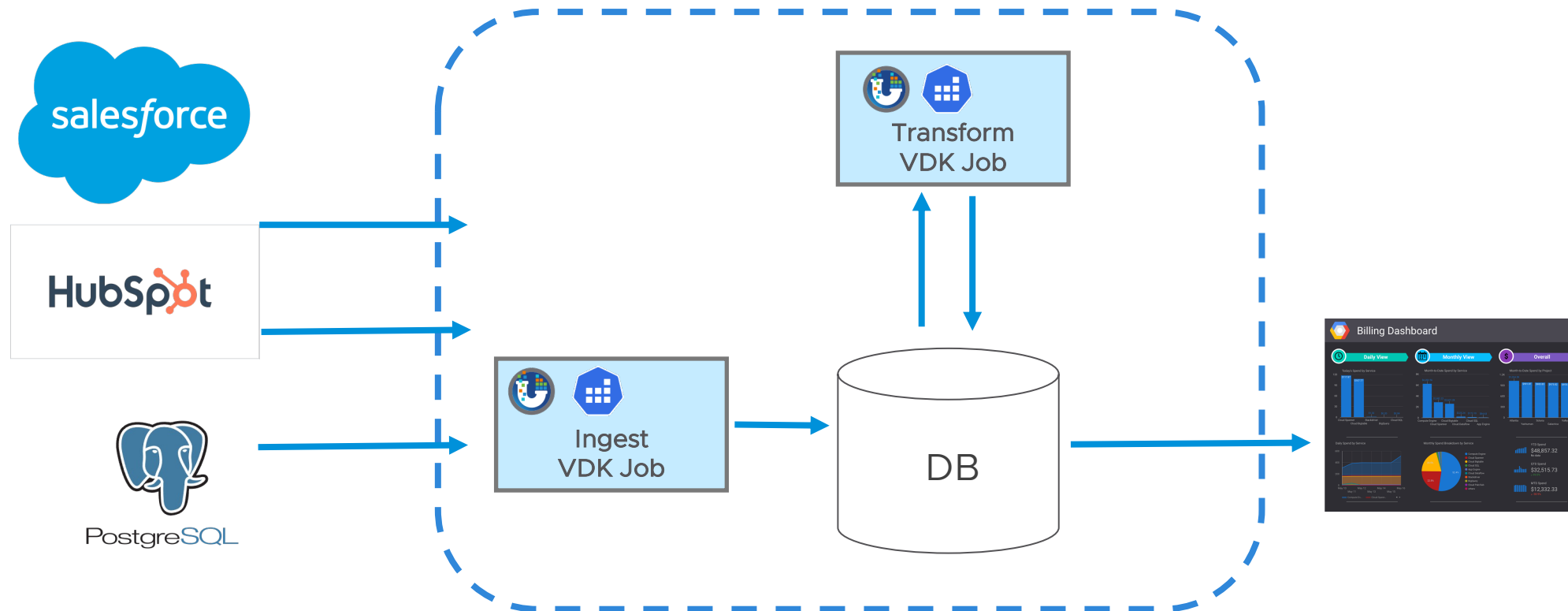


# Data Transformation

Pre-process/ Clean big data quickly, create business KPIs

Out of the box functionality for pre-processing data in Star schema data model (fact/dimensions)

Data engineers focus on data transformations





## Raw Data

Table:  
vmc\_console\_org

id	sla	type
ab4234	core	partner
fh3hf3f	internal	internal
df3f422	core	customer

Table:  
company\_name

org	company	rn
ab4234	Amazon	1
1235ff	Amazon	2
df3f422	MIT	1

Table:  
billing\_start\_date

org	date	rn
ab4234	11/10/20	1
df3f422	7/10/20	1
df3f422	5/10/20	1

```
1  SELECT
2      org.id AS org_id,
3      org.org_sla,
4      org.org_type,
5      COALESCE(CASE
6          WHEN ocn.comp_name_type = 'map_org_company_name'
7              THEN ocn.company_name
8          WHEN org.org_type IN ('INTERNAL_CORE', 'INTERNAL')
9              THEN 'VMware'
10         WHEN org.org_type = 'INTERNAL_AWS'
11             THEN 'Amazon AWS'
12         ELSE ocn.company_name
13         END, org.org_display_name) AS company_name,
14      bsd.billing_start_date
15  FROM vmc_console_org as org
16  LEFT JOIN company_name ocn ON ocn.dim_org_id = org.id
17      AND ocn.rn = 1
18  LEFT JOIN billing_start_date bsd ON bsd.org_id = org.id
19      AND bsd.rn = 1
```

transform and load into



## Business objects

Table:  
dimension vmc organization

id	type	company	created
ab4234	partner	Amazon	11/10/20
df3f422	customer	MIT	5/10/20
fh3hf3f	internal	VMware	7/10/20

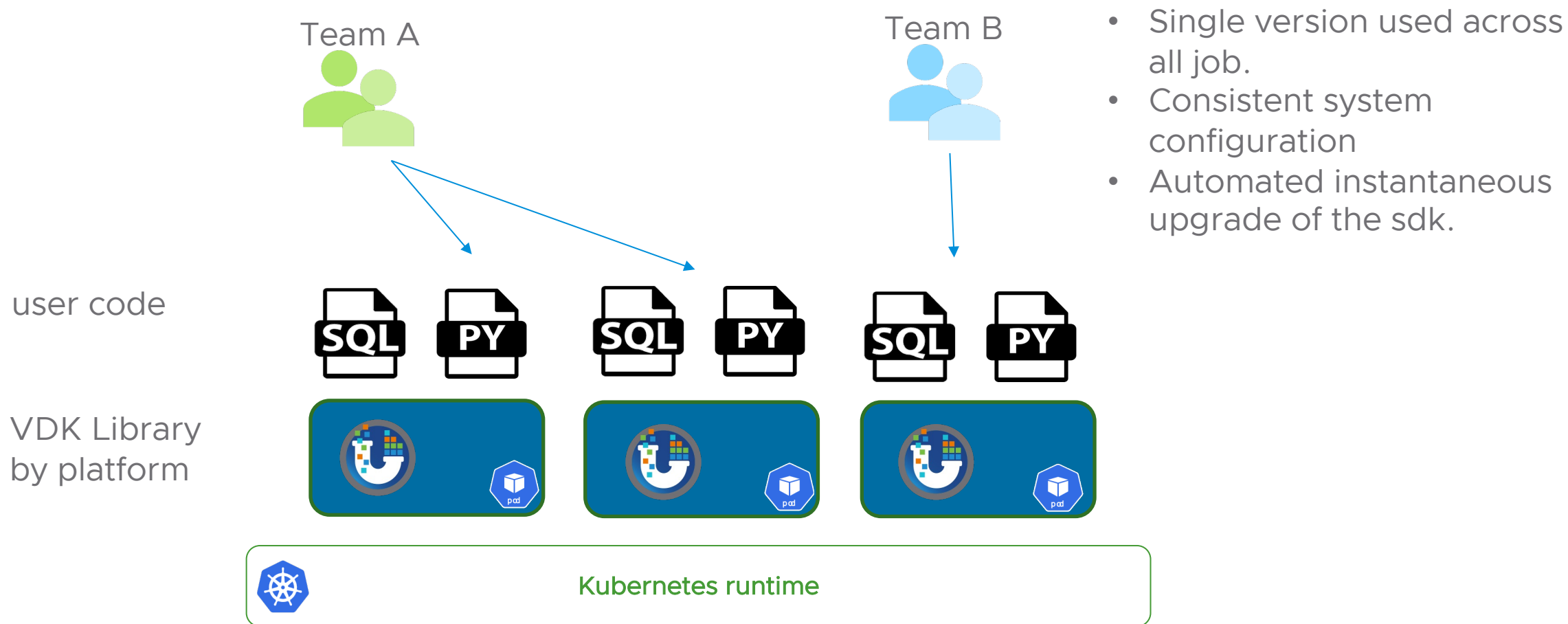
# Deployment



Data Teams

# VDK SDK and VDK Runtime (Control Service)

Data library and Job runner





# Operations



Data Teams



# Data Jobs statistics, status, notifications, alerts

## Data Jobs

Data Jobs help Data Engineers develop, deploy, run, and manage data processing workloads



Jobs with Failing executions - Last 24h	
pmitev-test	48
e2e-cp-stg-test-job	8
e2e-cy-dp-failing	2
e2e-cp-stg-failing-job	2

Most recent failed executions - UTC time		
pmitev-test-1662670800	Platform Error	Sep 8, 2022, 09:01:11 PM
pmitev-test-1662669000	Platform Error	Sep 8, 2022, 08:31:16 PM
pmitev-test-1662667200	Platform Error	Sep 8, 2022, 08:01:12 PM
pmitev-test-1662665400	Platform Error	Sep 8, 2022, 07:31:09 PM
e2e-cp-stg-test-job-1662664523	User Error	Sep 8, 2022, 07:16:25 PM

## Manage Data Jobs

REFRESH

ENABLE

DISABLE

EXECUTE NOW

All

Enabled

Disabled

Not Deployed

Search

	Job name	Deployment Status	Last Execution End (UTC)	Last Execution Duration	Last Execution Status	Success rate	Schedule (in UTC)	Logs	Details
<input type="radio"/>	tms-staging-vdk-1644321902	✓	Sep 8, 2022, 09:30 PM	12s	✓ Success	100.00%	Every minute	<a href="#">Logs</a>	<a href="#">Details</a>
<input type="radio"/>	pmitev-test	✓	Sep 8, 2022, 09:01 PM	57s	⚠ Platform Error	0.00% (337 failed)	Every 30 minutes	<a href="#">Logs</a>	<a href="#">Details</a>
<input type="radio"/>	e2e-cp-stg-test-job	✓	Sep 8, 2022, 07:16 PM	1m 1s	⚠ User Error	0.00% (115 failed)	At 11:11 PM, on day 5 of the month, and on Monday, only in August	<a href="#">Logs</a>	<a href="#">Details</a>





Infra &  
Operations  
Team

# Managed Infrastructure

## Improve data infrastructure security (demo)

Check out:



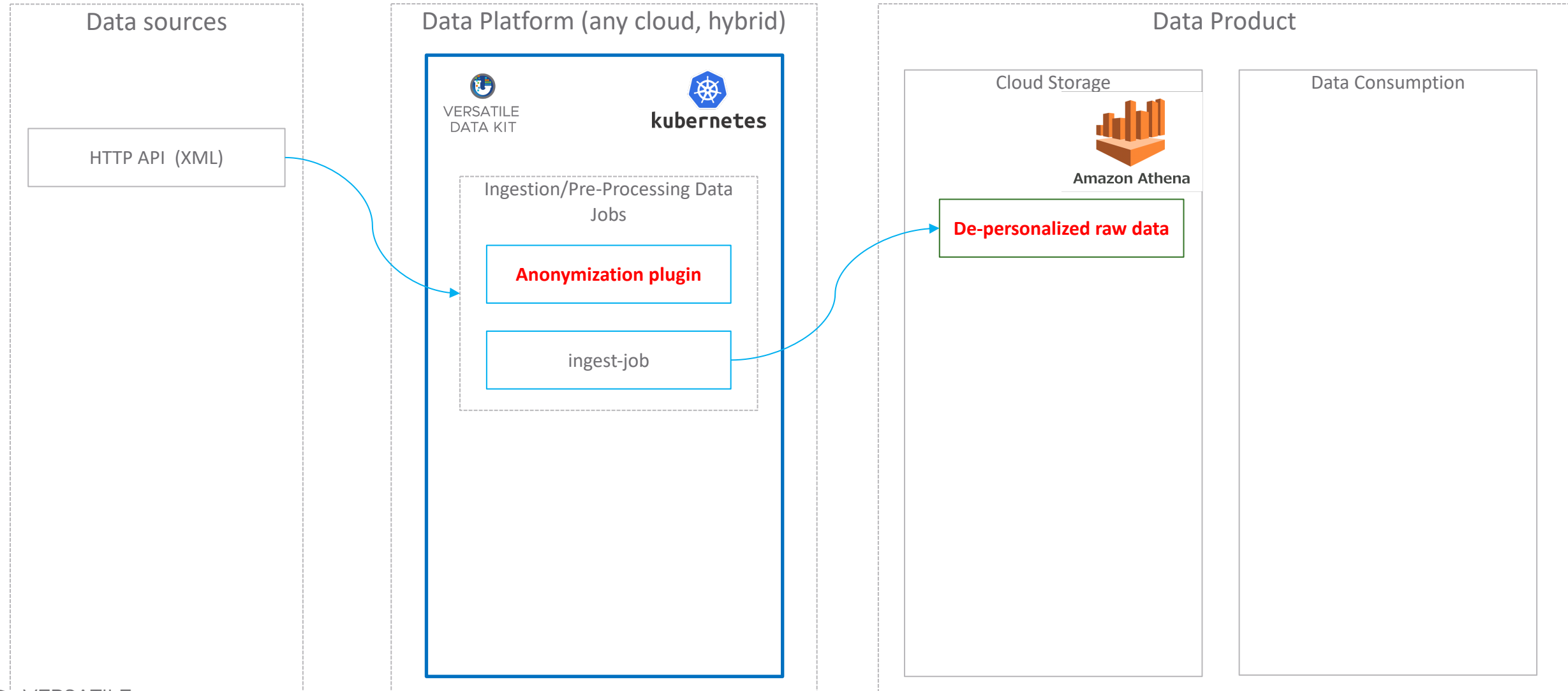
<https://github.com/vmware/versatile-data-kit/tree/main/examples/ingest-and-anonymize-plugin>



©2022 VMware, Inc.



# Data Platform – Ingest and Anonymization POC





# What are we going to do?



vdk run ingest-job

```
job_input.send_object_for_ingestion(  
    payload=response.json(),  
    destination_table="rest_table"  
)
```

data:

userId	id	title	completed
1	11	delectus	1
2	12	autem	0

intercepted



vdk-poc-anonymization (plugin)

---

userId	id	title	completed
1	11	325959af	1
2	12	f0a4cedd	0



## Improve data infrastructure stability (demo)



# What we did

vdk run sql-job    cursor.execute(...)    vdk query -q "..."

```
select
    count(1) as uploads,
    trunc(arrival_ts, 'ww') week
from org
```

↓ intercepted



vdk-query-validation (plugin)





# Deliver analytics platform for your business quickly



Infra &  
Operations  
Team



Check out:

<https://github.com/vmware/versatile-data-kit/wiki/Install-VDK-Control-Service-with-custom-SDK>