



Identity and Access Management (IAM) and Quotas

Omar Ismail

Solutions Developer, Google
Cloud





Agenda

Course Intro

Beam and Dataflow Refresher

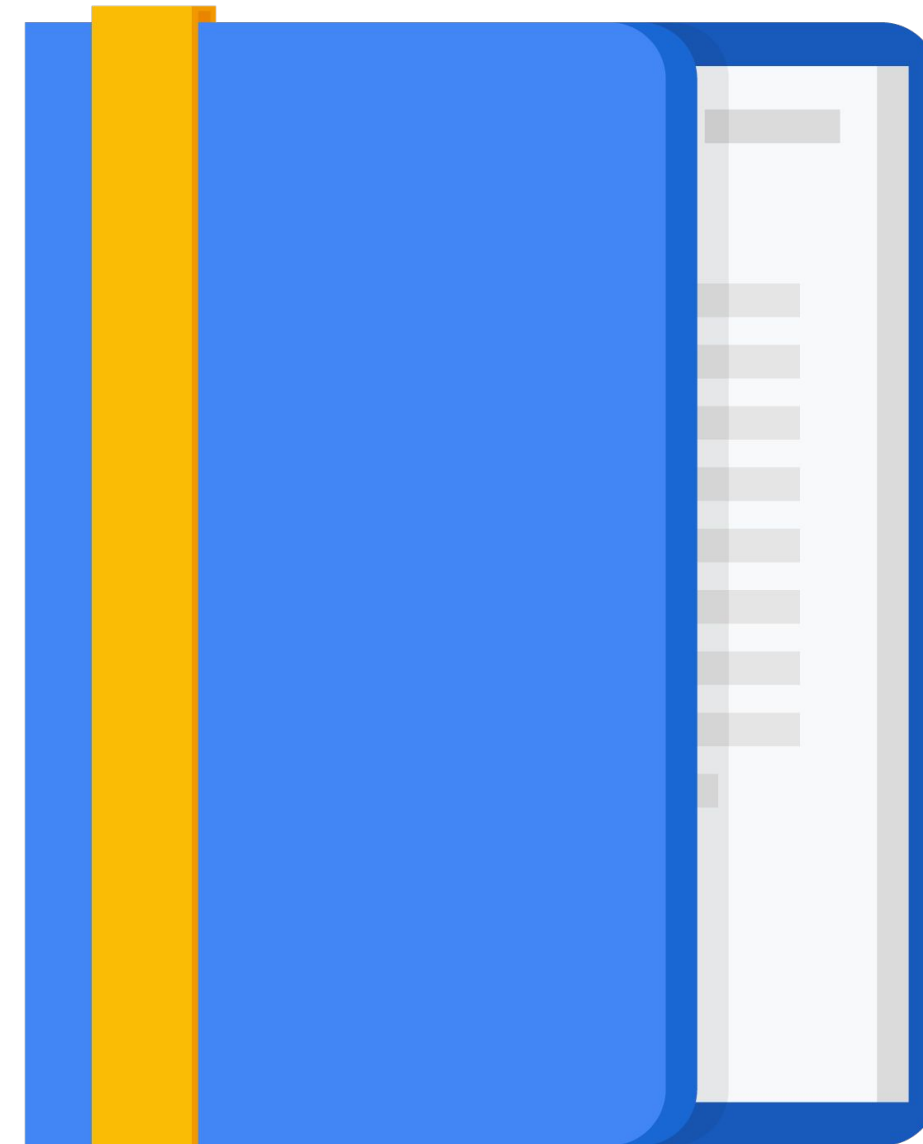
Beam Portability

Separating Compute and Storage

IAM, Quotas, and Permissions

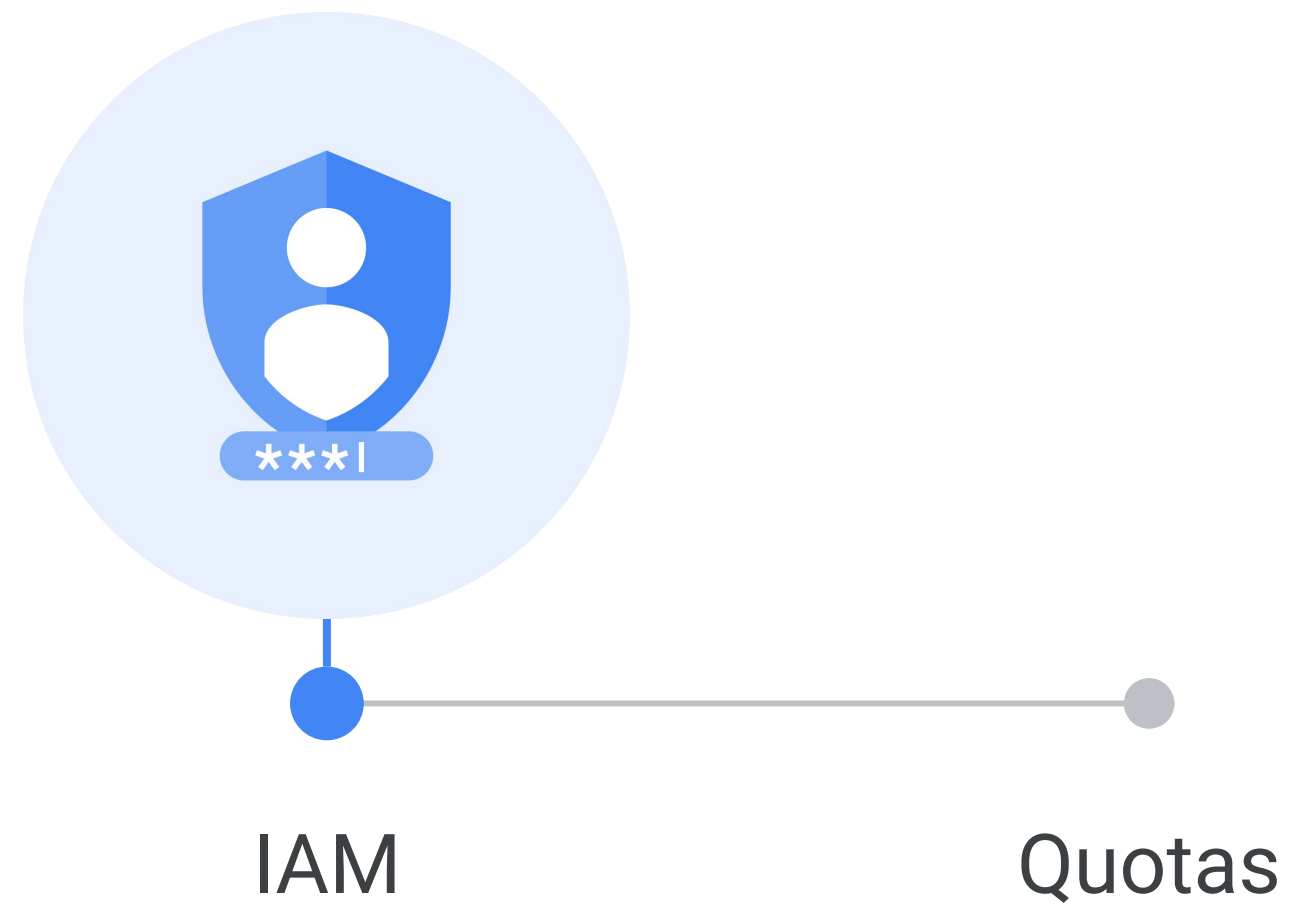
Security

Summary

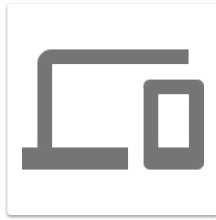


IAM, quotas, and permissions

Agenda

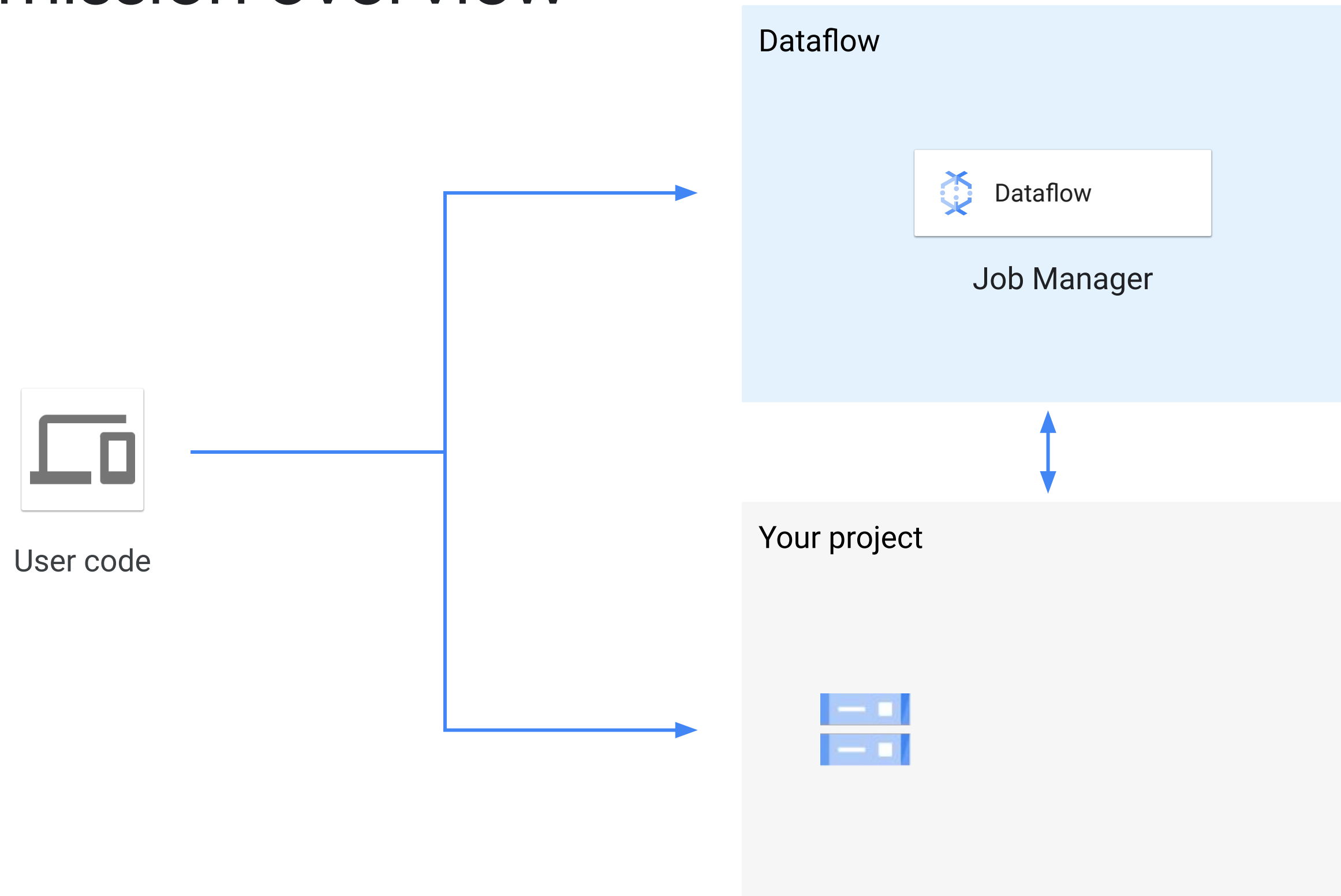


Job submission overview

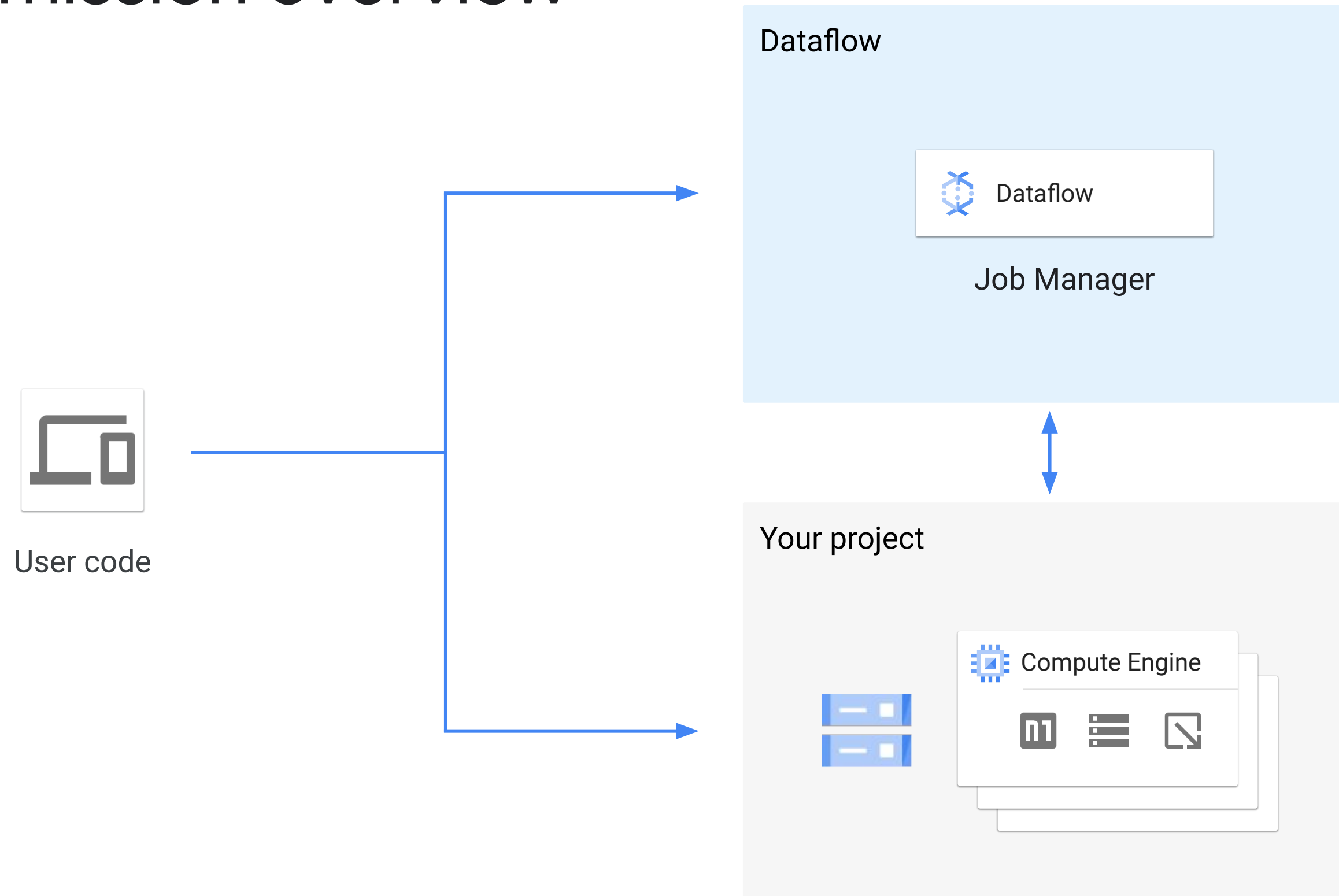


User code

Job submission overview



Job submission overview



Three credentials

- 1 User roles
- 2 Dataflow service account
- 3 Controller service account



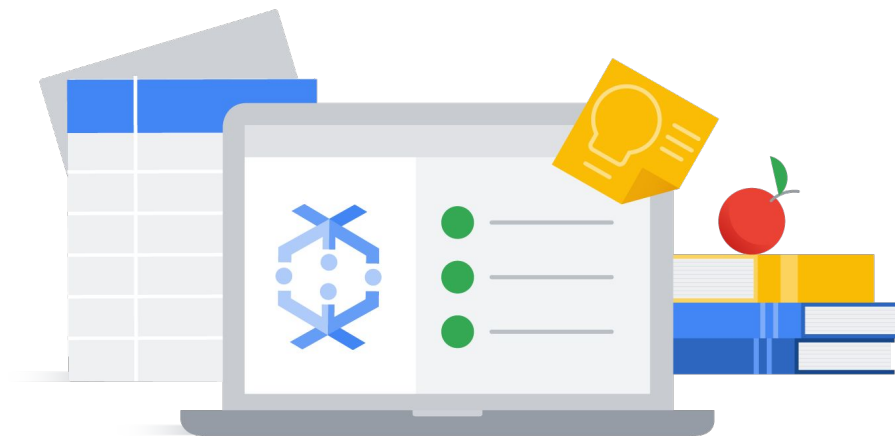
User roles



Dataflow Viewer

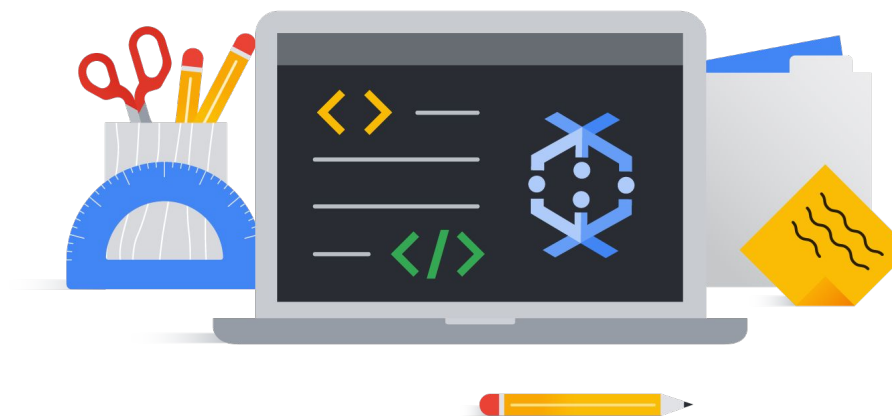
Provides read-only access to all Dataflow-related resources.

User roles



Dataflow Viewer

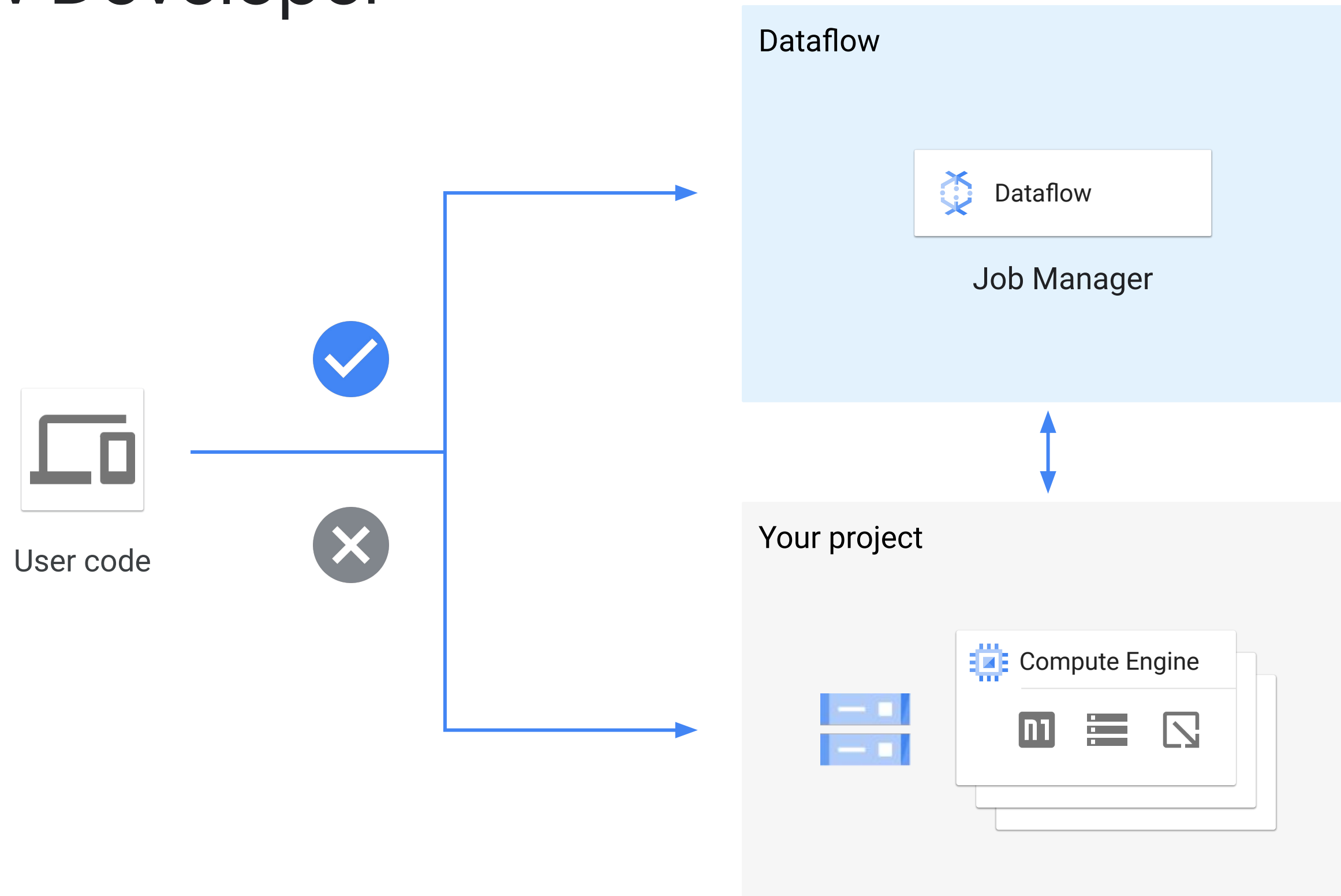
Provides read-only access to all Dataflow-related resources.



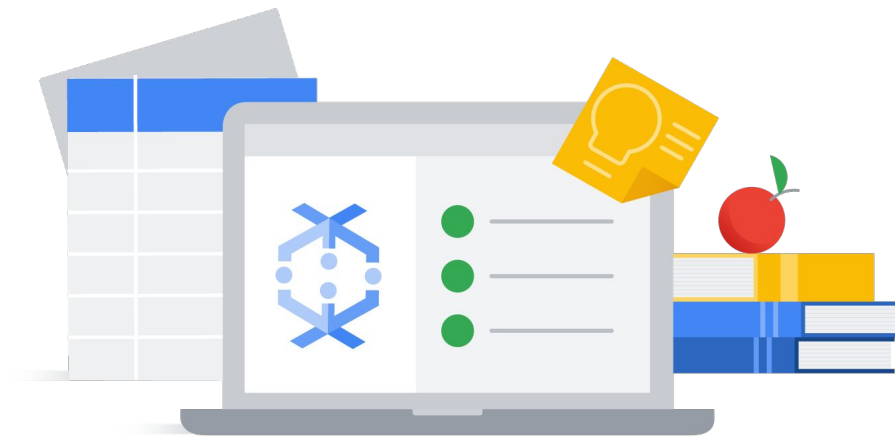
Dataflow Developer

Provides access to view, update, and cancel Dataflow jobs.

Dataflow Developer

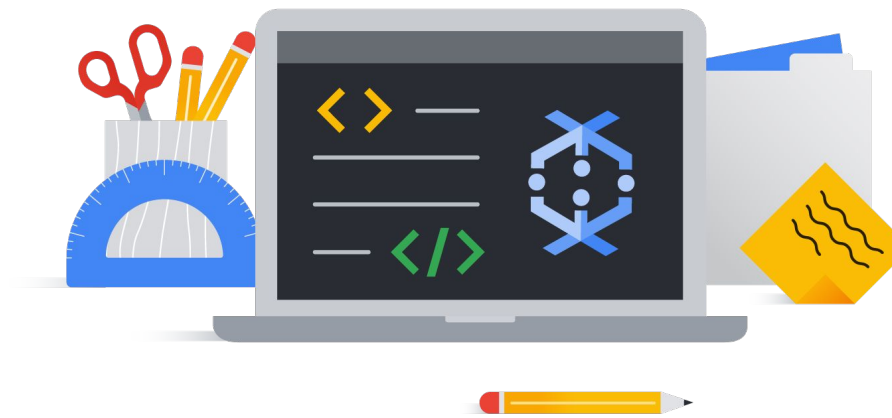


User roles



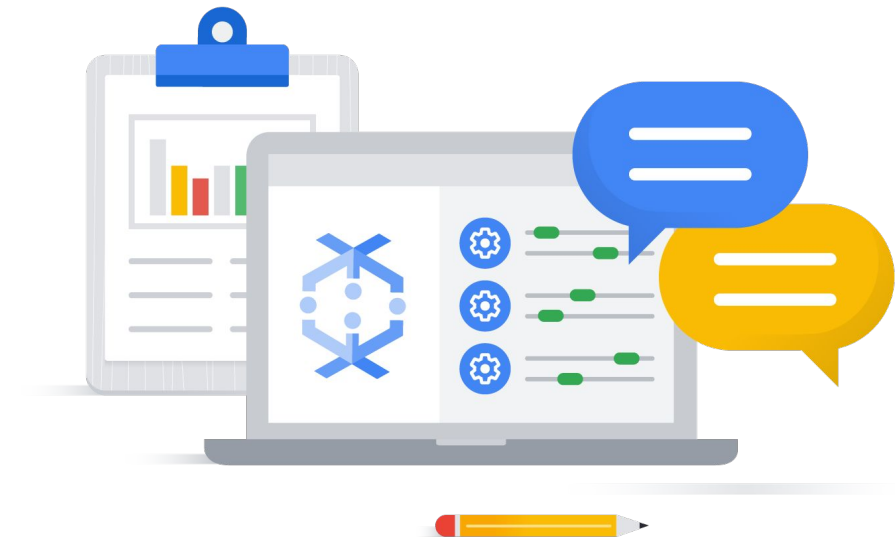
Dataflow Viewer

Provides read-only access to all Dataflow-related resources.



Dataflow Developer

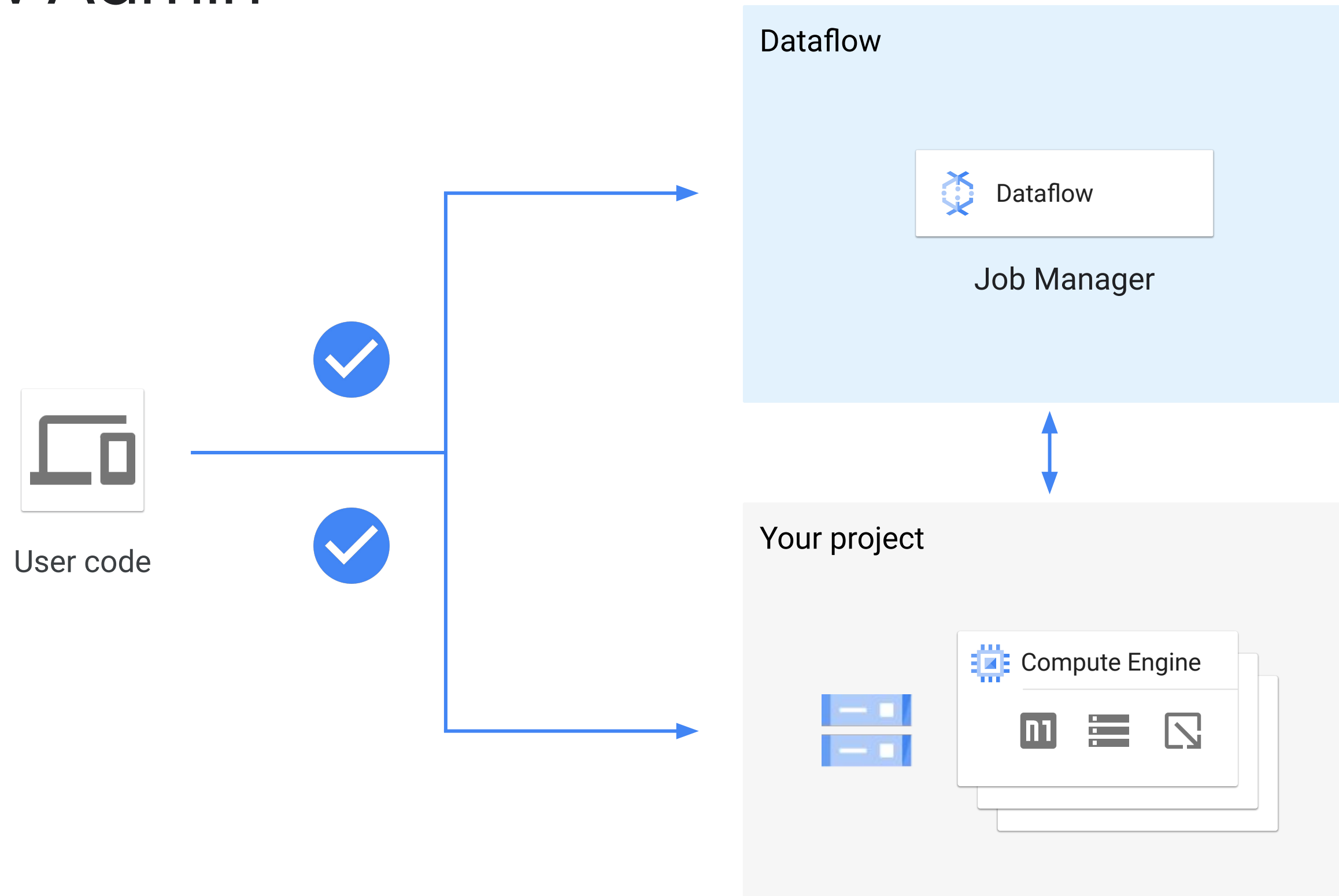
Provides access to view, update, and cancel Dataflow jobs.



Dataflow Admin

Provides access for creating and managing Dataflow jobs.

Dataflow Admin



Question

A Dataflow job (2021-01-31_14_30_00-9098096469011826084) must be canceled. You are assigned the Dataflow Developer role. Can you run the following command:

```
$ gcloud dataflow jobs cancel  
2021-01-31_14_30_00-9098096469011826084 --region=$REGION
```

Three credentials

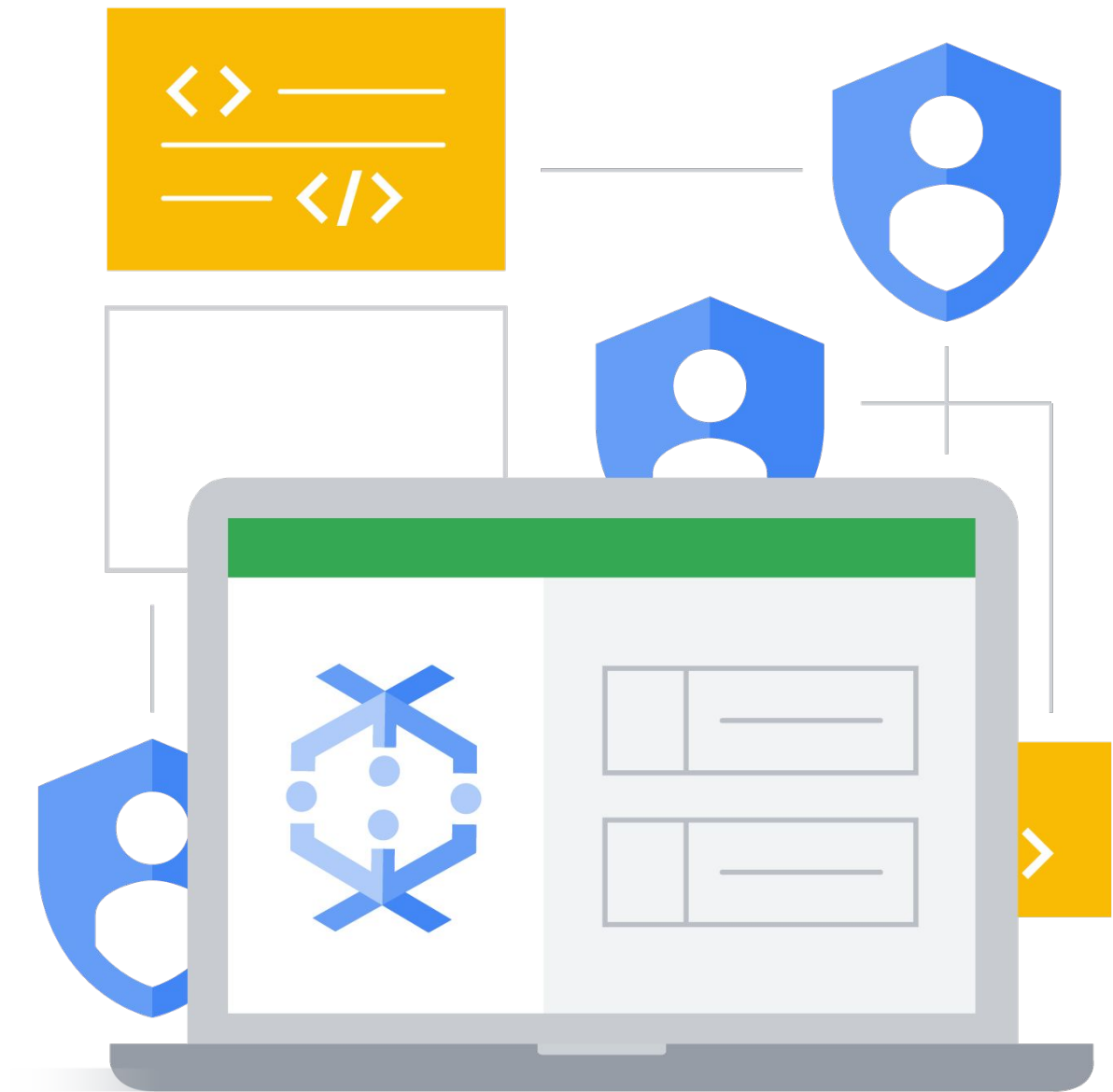
- 1 User roles
- 2 Dataflow service account
- 3 Controller service account



Dataflow service account

The Orchestrator

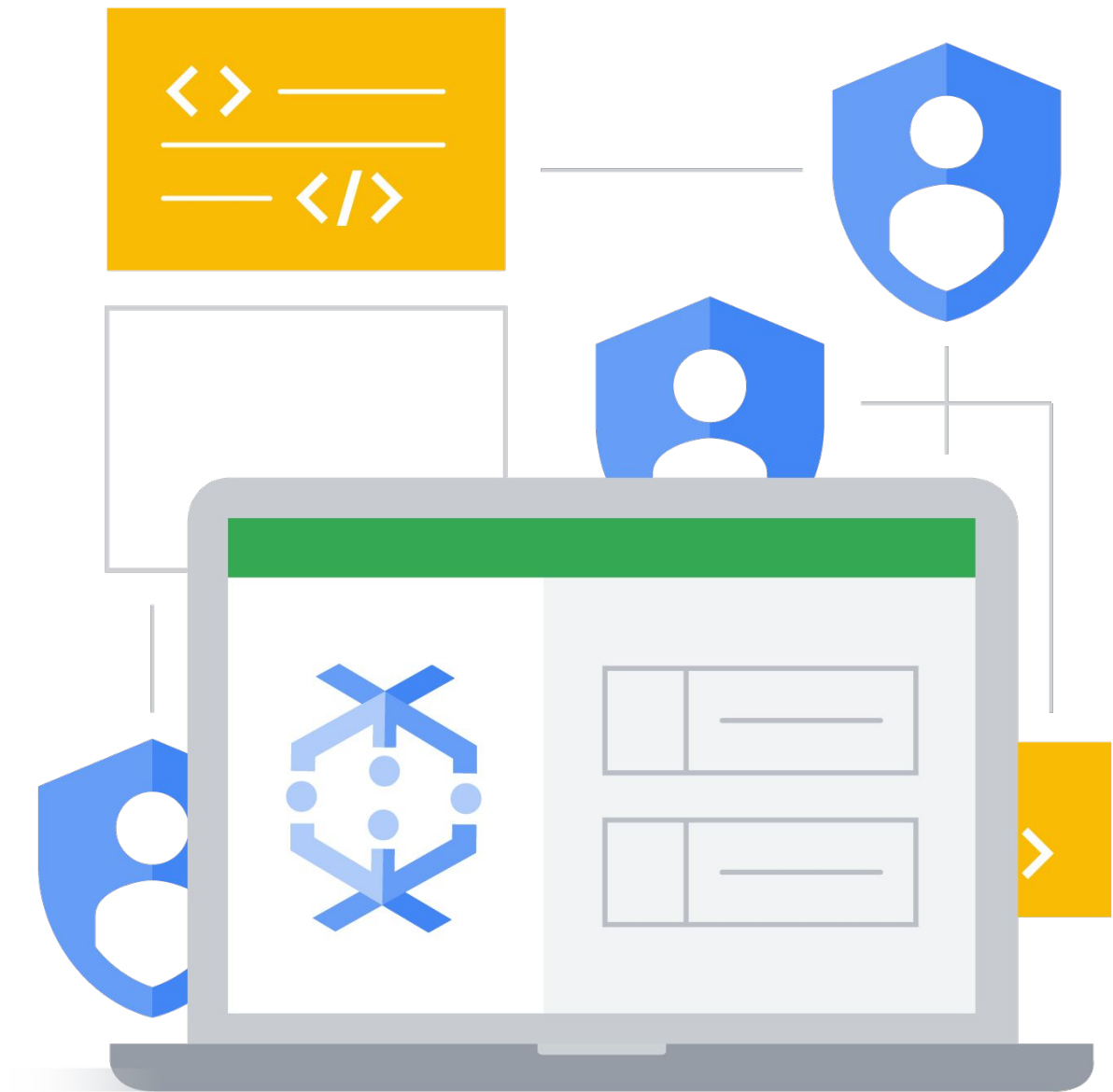
- Interacts between your project and Dataflow



Dataflow service account

The Orchestrator

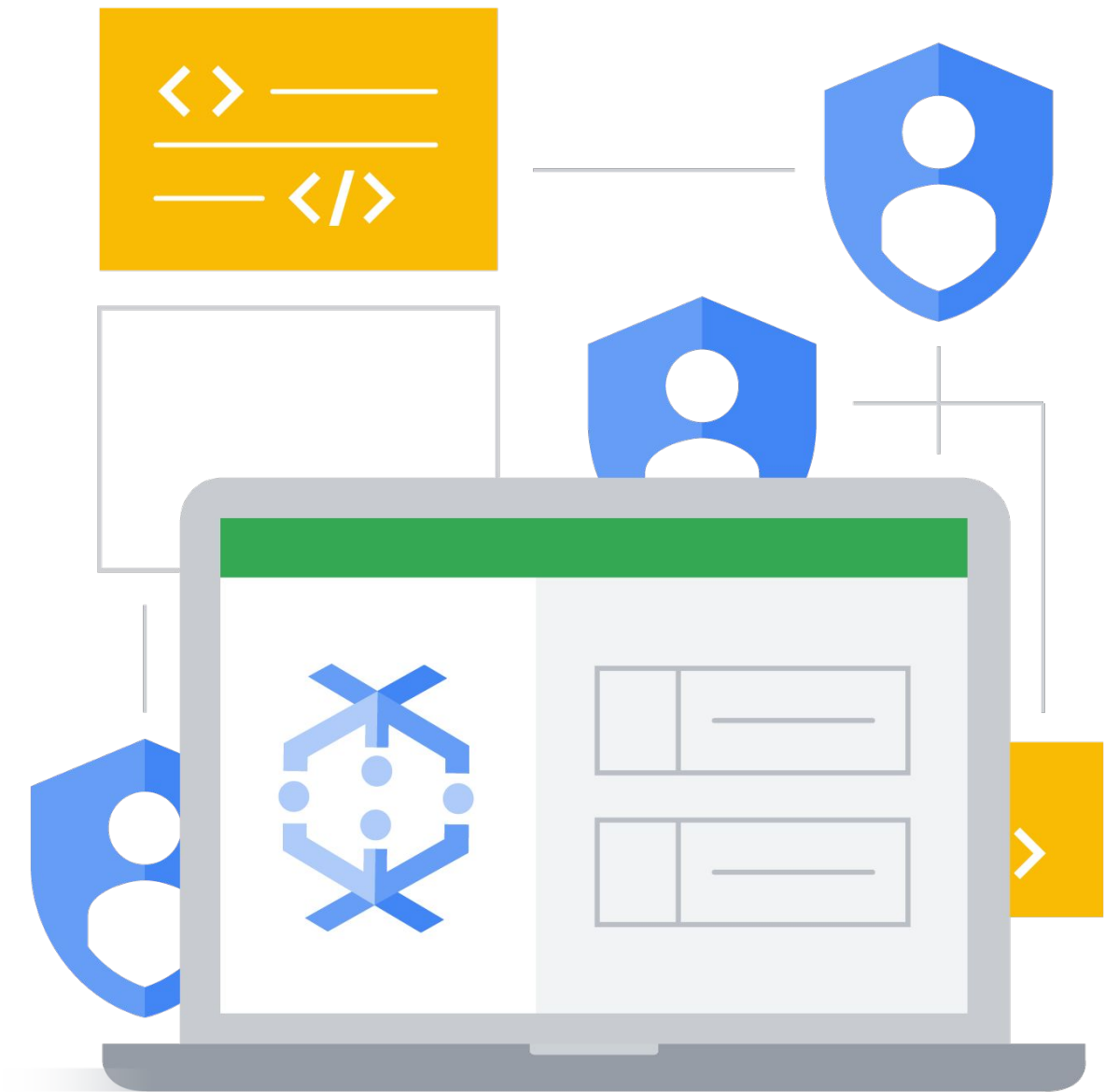
- Interacts between your project and Dataflow
- Used for worker creation and monitoring



Dataflow service account

The Orchestrator

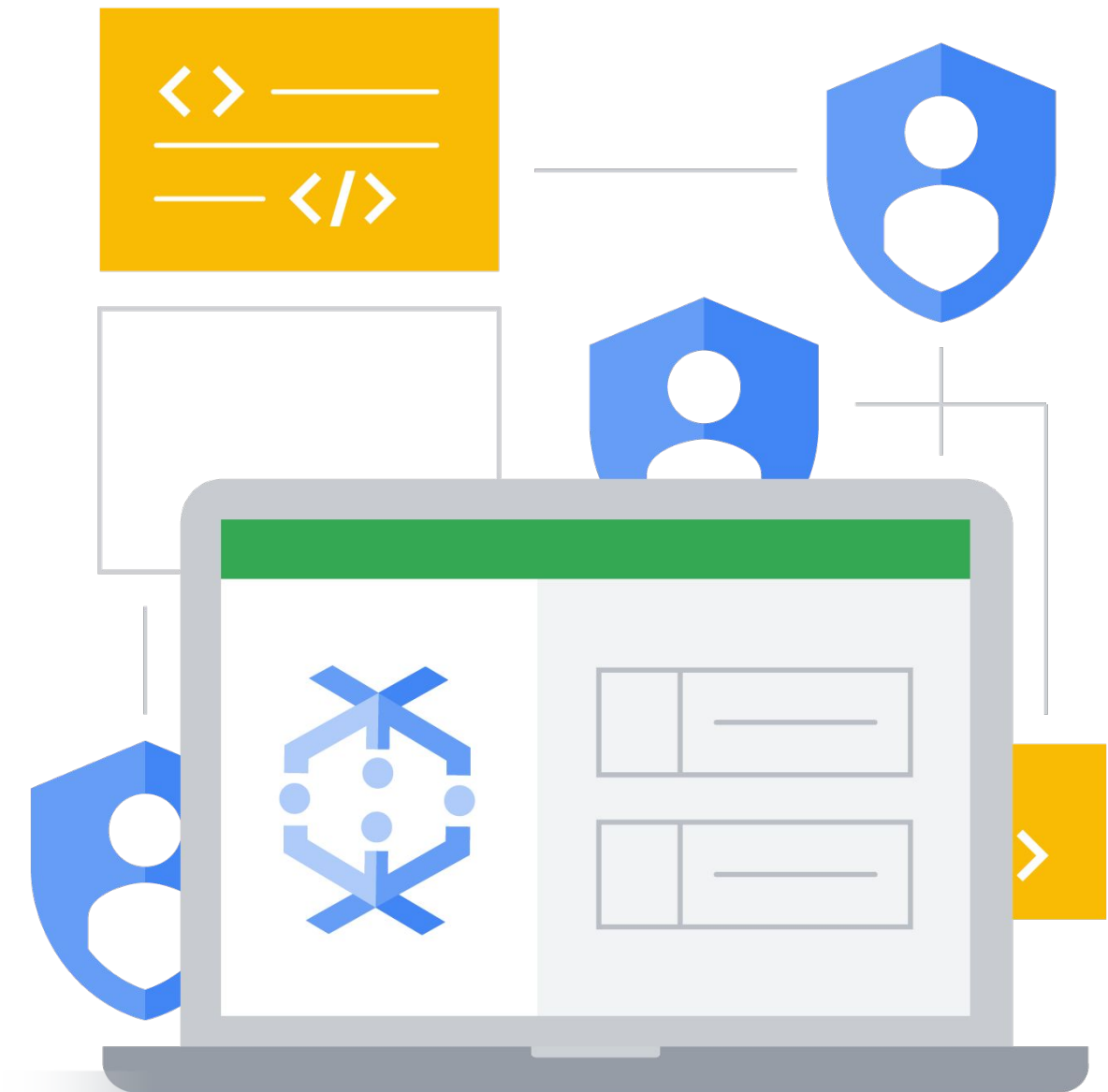
- Interacts between your project and Dataflow
- Used for worker creation and monitoring
- `service-<project_number>@dataflow-service-producer-prod.iam.gserviceaccount.com`



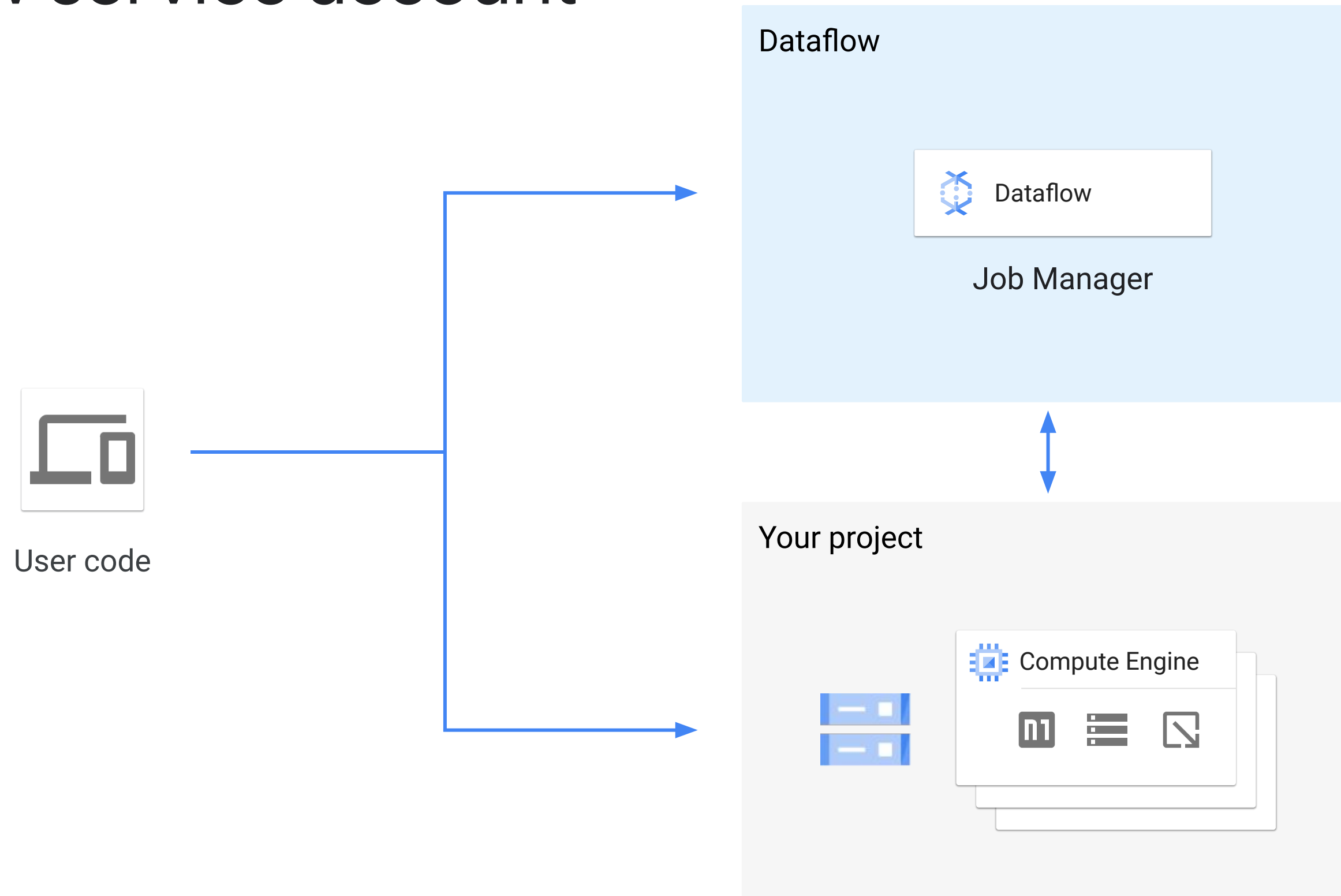
Dataflow service account

The Orchestrator

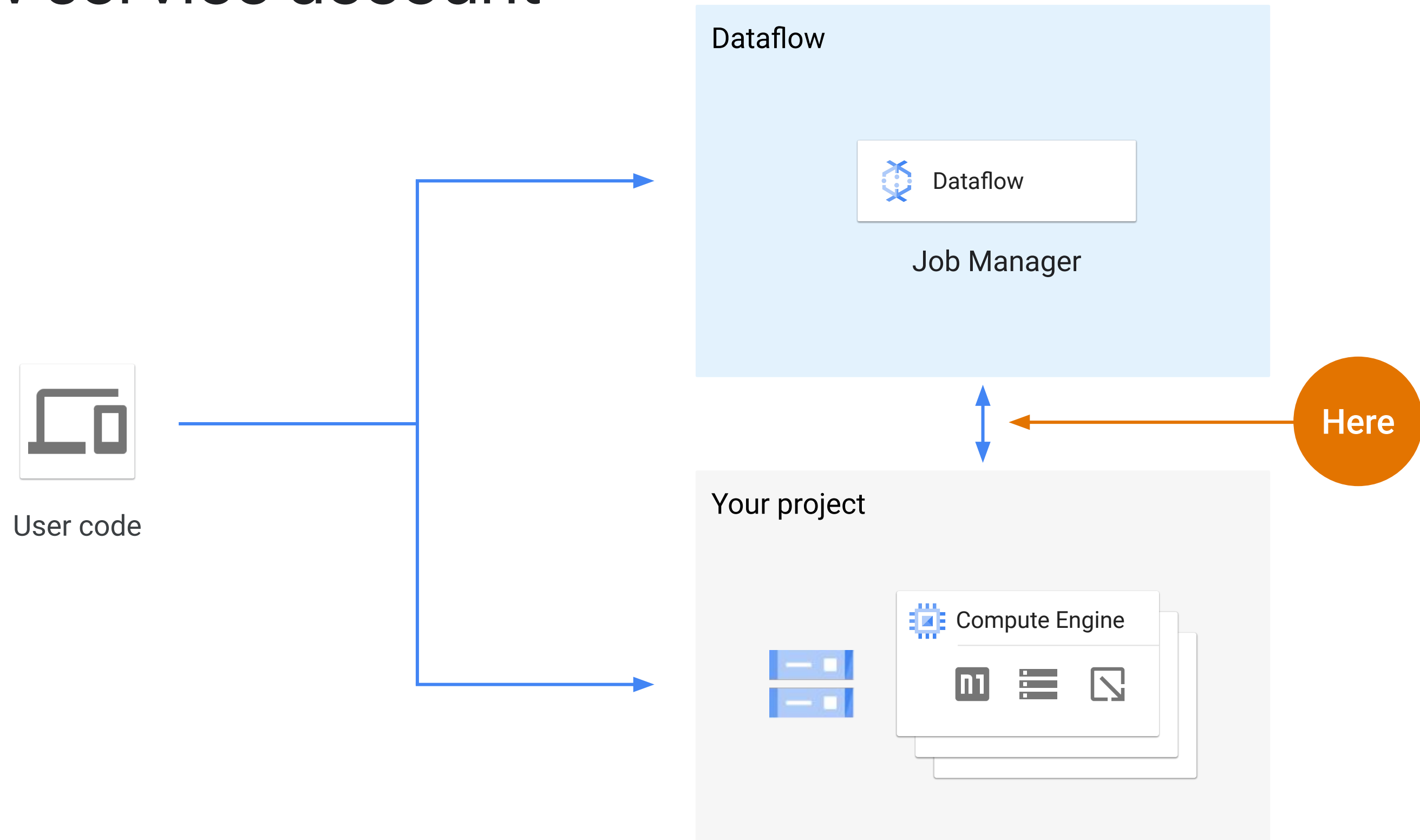
- Interacts between your project and Dataflow
- Used for worker creation and monitoring
- `service-<project_number>@dataflow-service-producer-prod.iam.gserviceaccount.com`
- Assigned the Dataflow Service Agent role



Dataflow service account



Dataflow service account



Three credentials

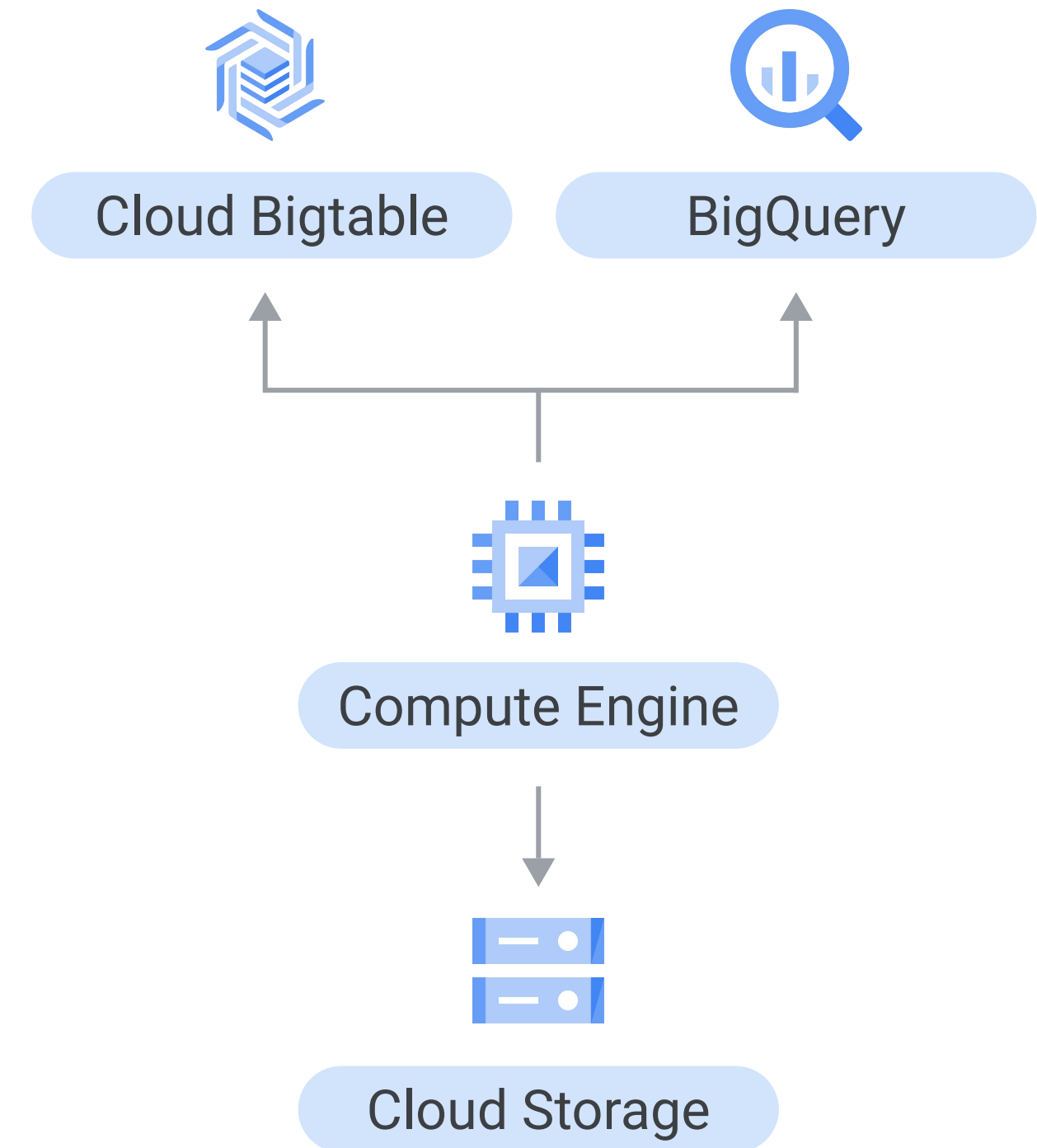
- 1 User roles
- 2 Dataflow service account
- 3 Controller service account



Controller service account

The Worker

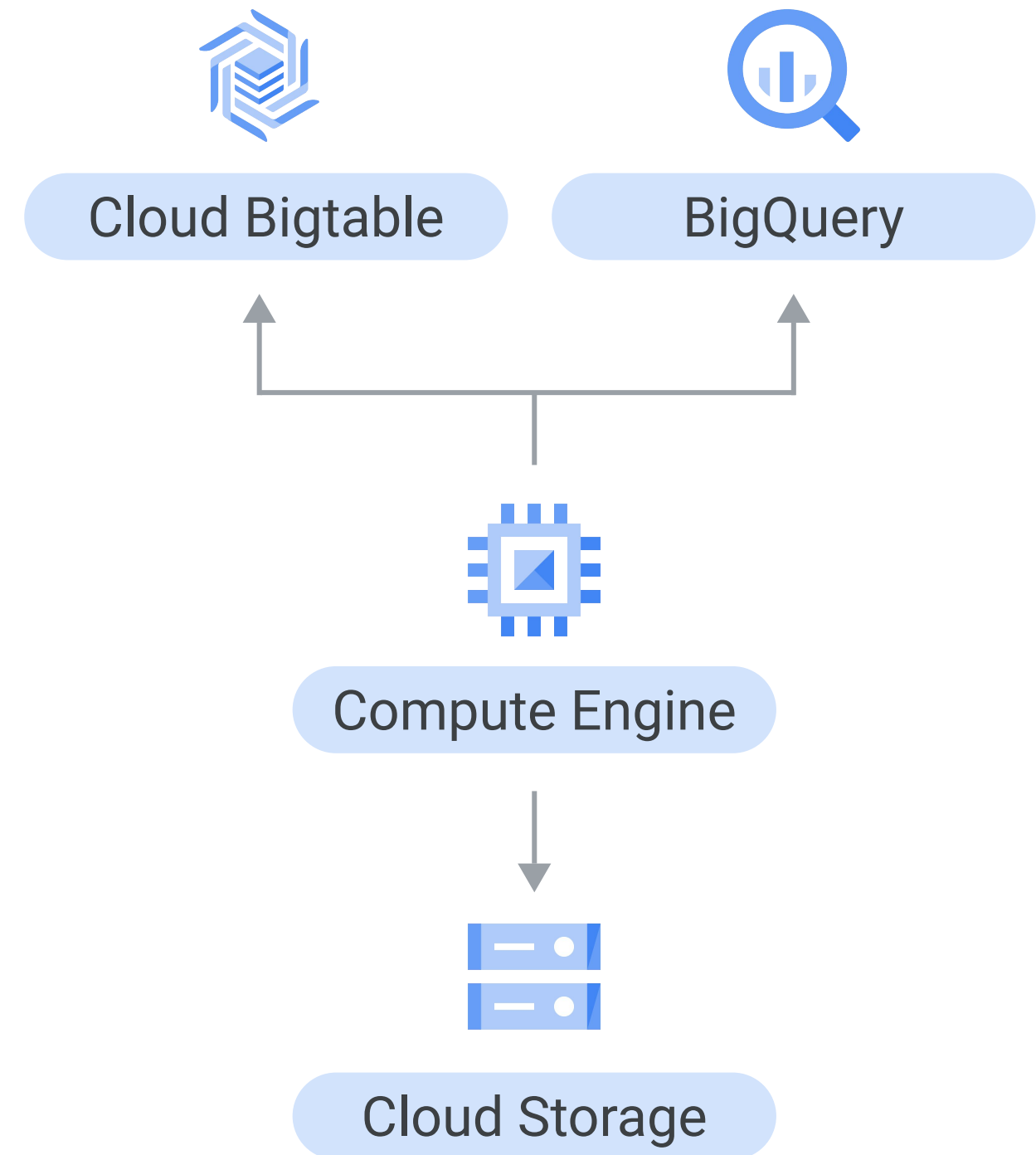
- Used by the workers to access resources needed by the pipeline



Controller service account

The Worker

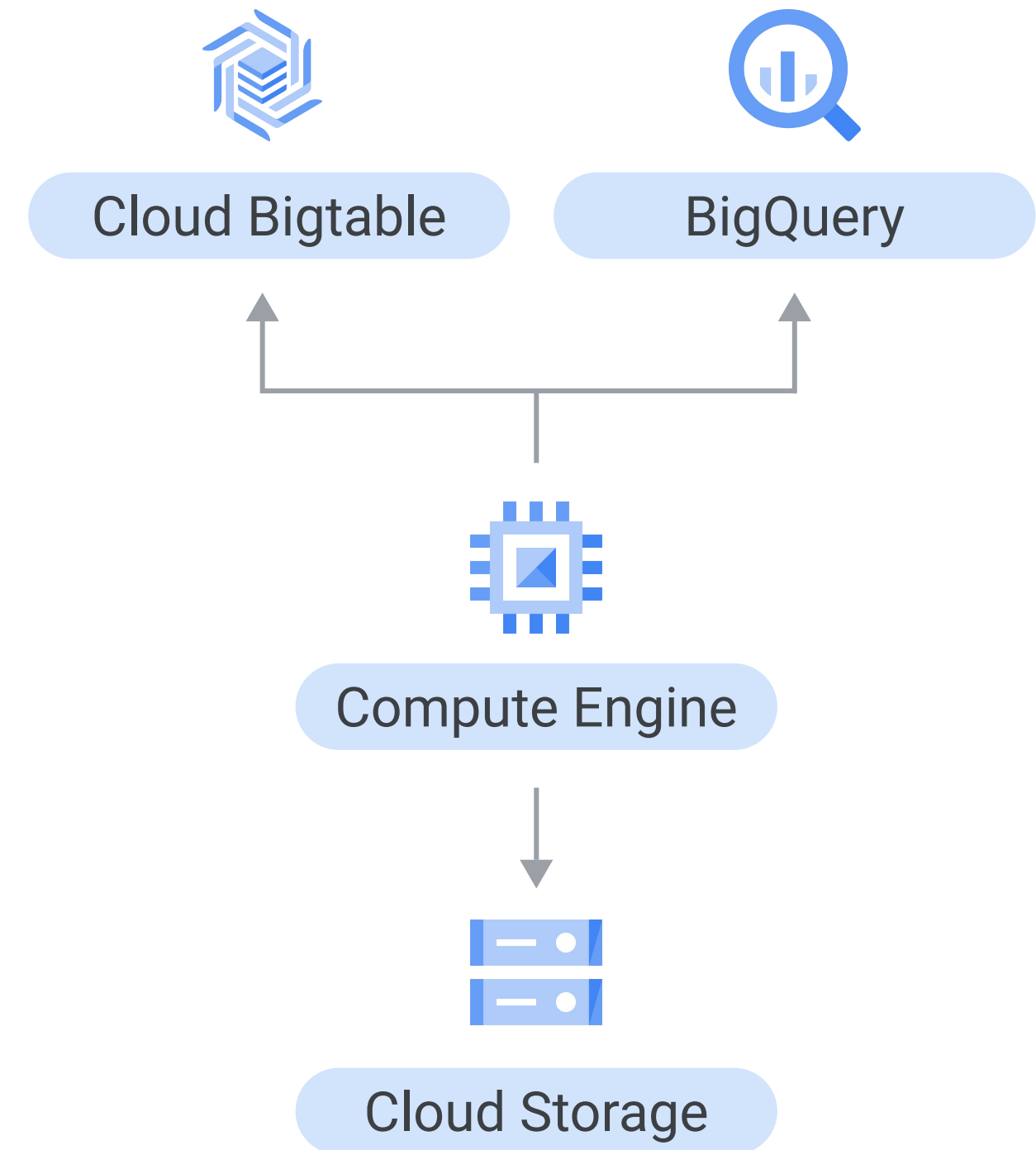
- Used by the workers to access resources needed by the pipeline
- <project-number>-compute@developer.gserviceaccount.com



Controller service account

The Worker

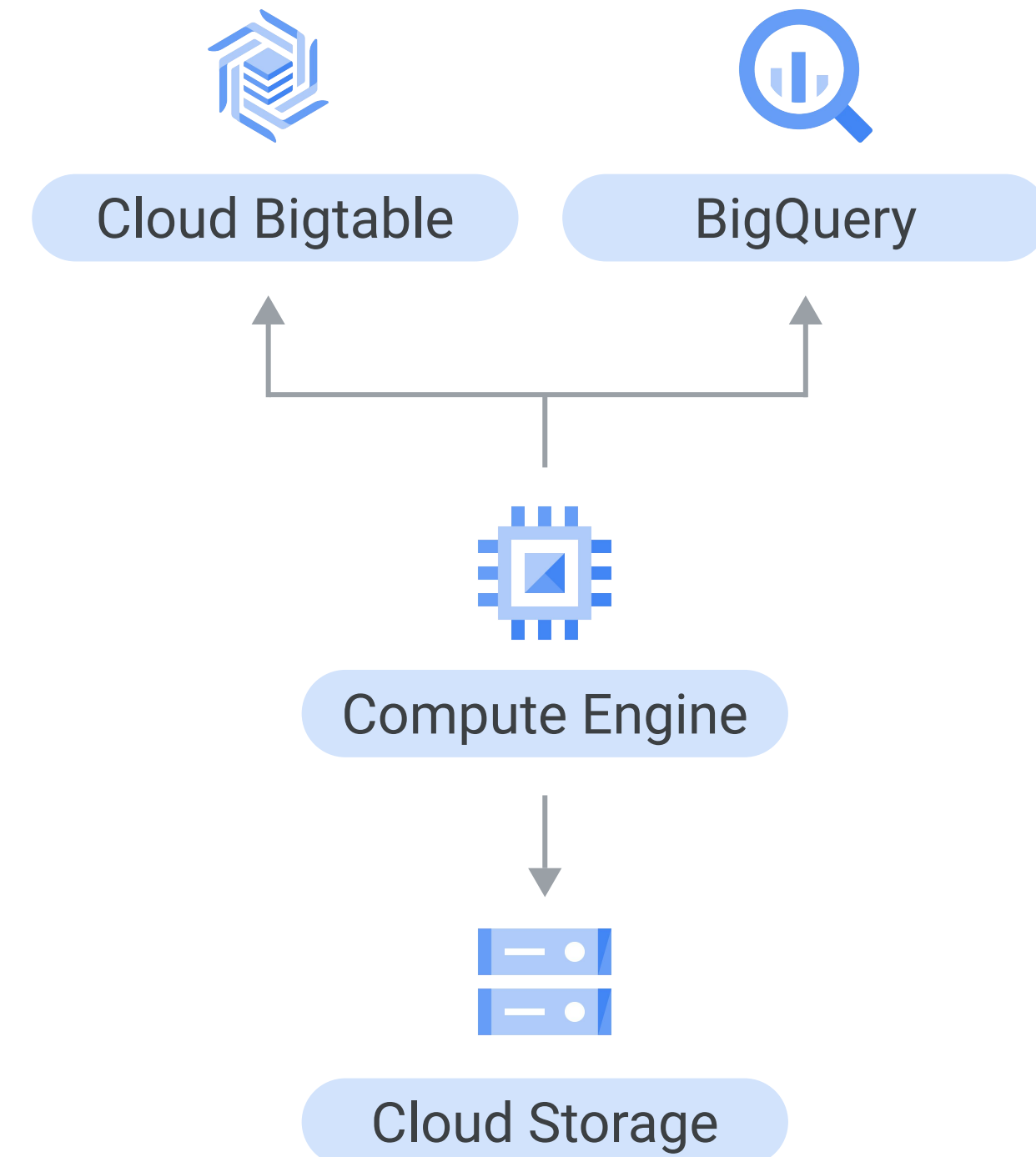
- Used by the workers to access resources needed by the pipeline
- <project-number>-compute@developer.gserviceaccount.com
- Flag to override default:
 - Python: `--service_account_email`
 - Java: `--serviceAccount`



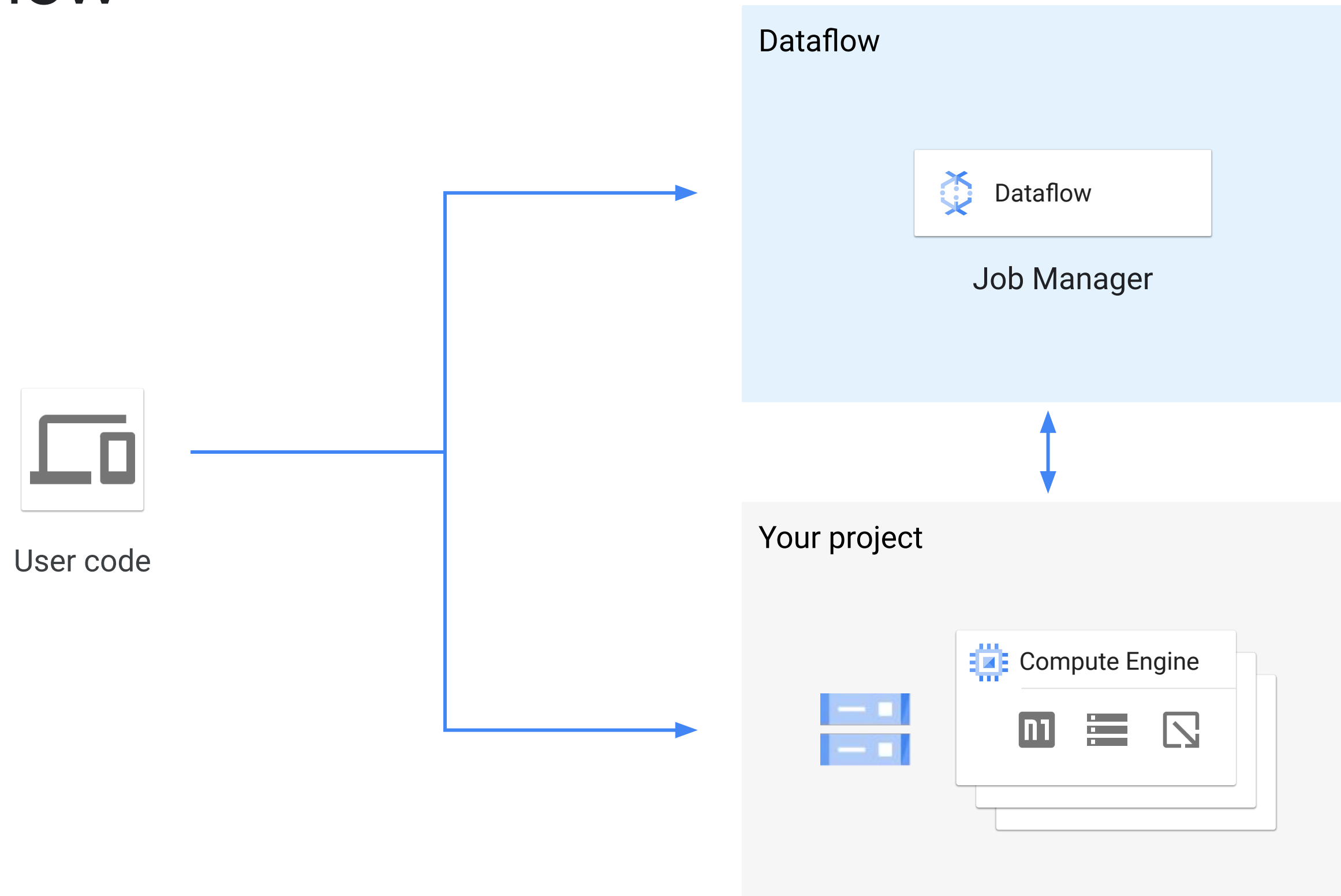
Controller service account

The Worker

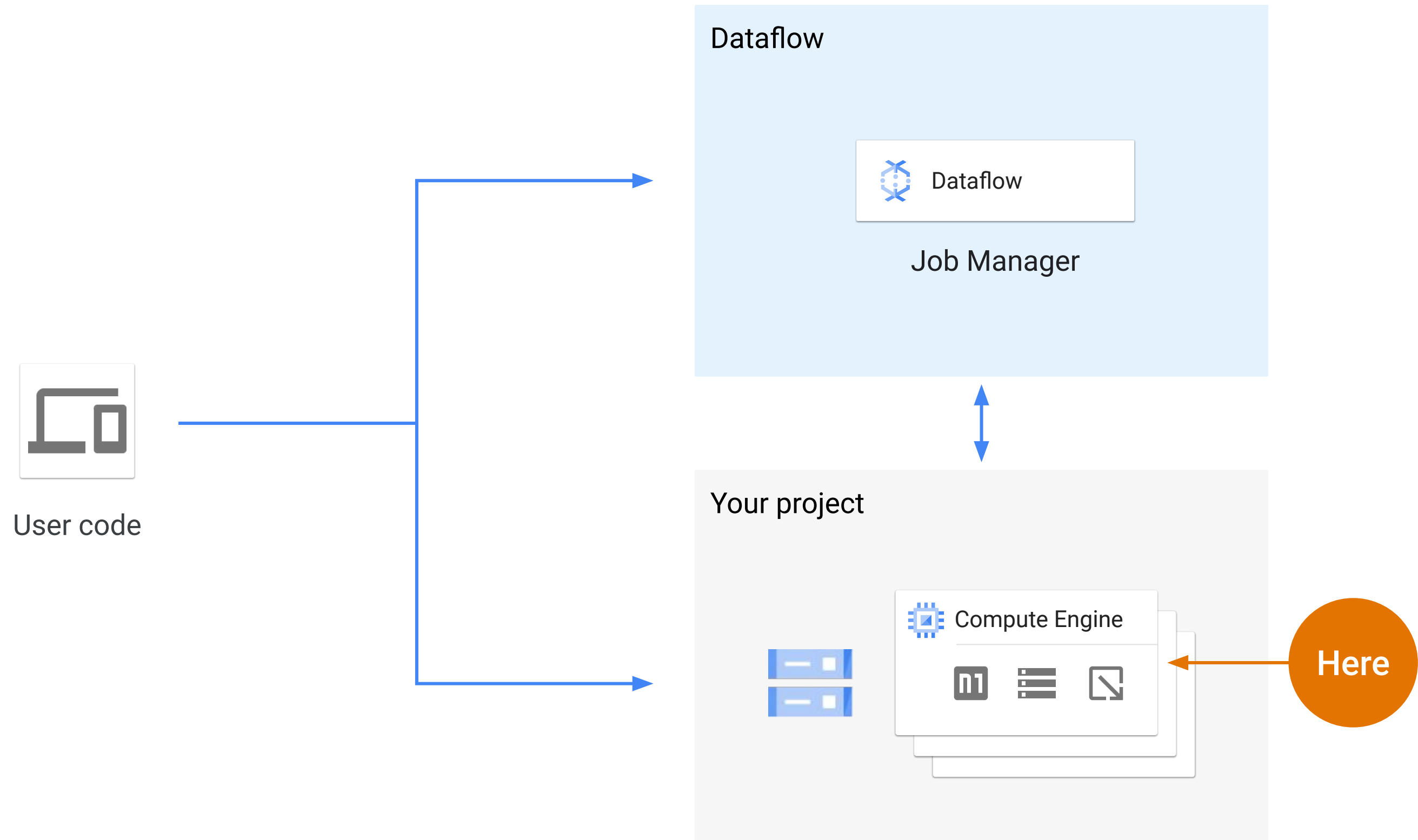
- Used by the workers to access resources needed by the pipeline
- <project-number>-compute@developer.gserviceaccount.com
- Flag to override default:
 - Python: `--service_account_email`
 - Java: `--serviceAccount`
- At a minimum, your custom service account must have the Dataflow Worker role.



Let's review

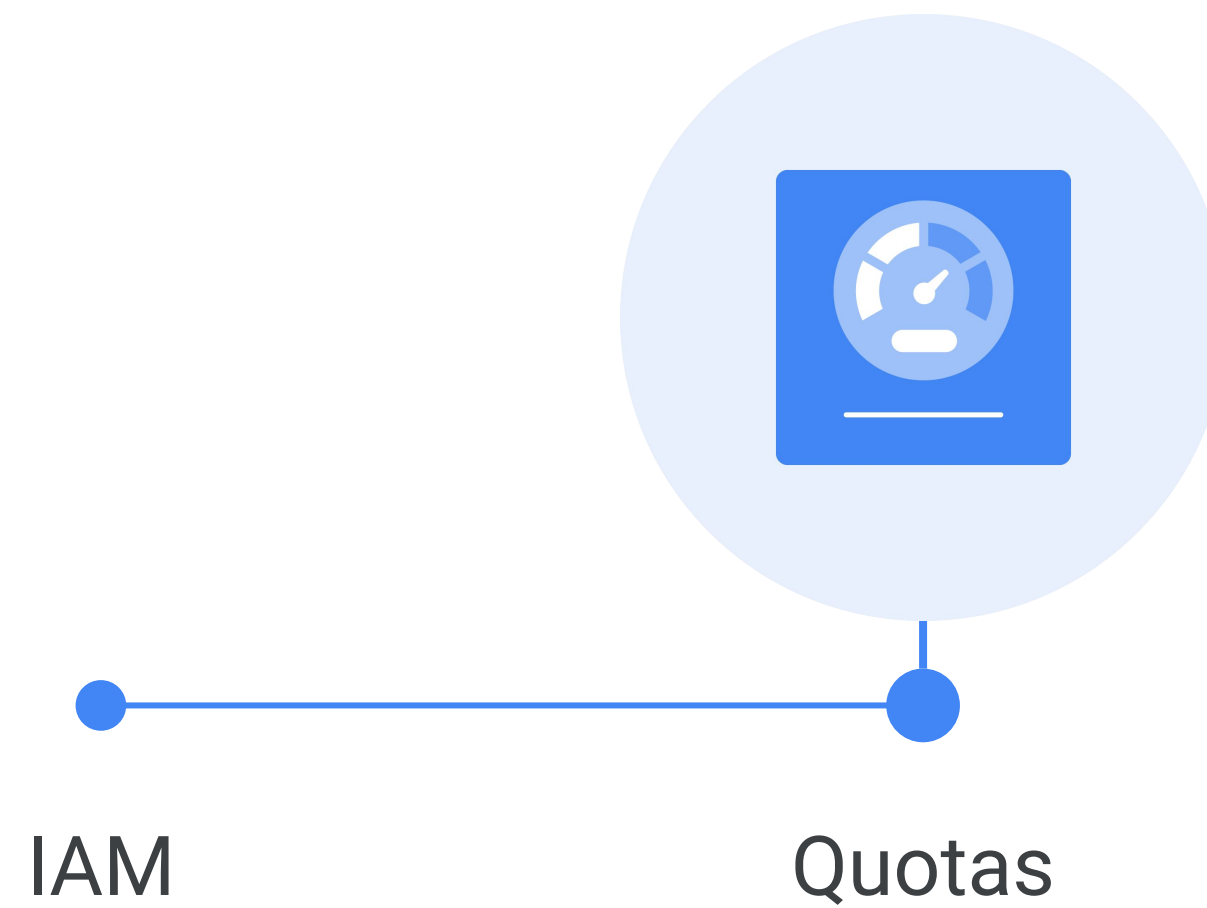


Answer



IAM, quotas, and permissions

Agenda



Quotas

CPUs

CPU quota based on number of cores

←

Quota metric details

EDIT QUOTAS

Service: Compute Engine API

Name: CPUs

Quota Metric: compute.googleapis.com/cpus

Limit Name: CPUS-per-project-zone/CPUS-per-project-region

Filter table

<input type="checkbox"/>	Location	Current Usage	7 Day Peak Usage	↑	Limit	
<input type="checkbox"/>	northamerica-northeast1	<div><div></div></div> 219	<div><div></div></div> 222		600	View hierarchy
<input type="checkbox"/>	us-central1	<div><div></div></div> 12	<div><div></div></div> 14		600	View hierarchy
<input type="checkbox"/>	us-east1	<div><div></div></div> 6	<div><div></div></div> 7		600	View hierarchy

Quotas

IPs

In-use IP addresses must be sufficient to accommodate the desired number of instances

←

Quota metric details

EDIT QUOTAS

Service: Compute Engine API

Name: In-use IP addresses

Quota Metric: compute.googleapis.com/regional_in_use_addresses

Limit Name: IN-USE-ADDRESSES-per-project-region

Filter table

	Location	Current Usage	7 Day Peak Usage	↑	Limit	
	northamerica-northeast1	<div><div></div></div> 207	<div><div></div></div> 210		575	View hierarchy
	us-central1	<div><div></div></div> 3	<div><div></div></div> 4		575	View hierarchy
	us-east1	<div><div></div></div> 3	<div><div></div></div> 3		575	View hierarchy

Quotas

Persistent Disks

Choose either HDD or SSD

	<div><div>←</div><div>Quota metric details</div><div> EDIT QUOTAS</div></div>																												
	Service: Compute Engine API																												
	Name: Persistent Disk Standard (GB)																												
	Quota Metric: compute.googleapis.com/disks_total_storage																												
	Limit Name: DISKS-TOTAL-GB-per-project-zone/DISK S-TOTAL-GB-per-project-region																												
	<div><div>Filter table</div><table><thead><tr><th><input type="checkbox"/></th><th>Location</th><th>Current Usage ↑</th><th>7 Day Peak Usage</th><th>Limit</th><th></th></tr></thead><tbody><tr><td><input type="checkbox"/></td><td>northamerica-northeast1</td><td><div><div></div></div>20,498 GB (20.498 TB)</td><td><div><div></div></div>22,498 GB (22.498 TB)</td><td>102,400 GB (102.4 TB)</td><td>View hierarchy</td></tr><tr><td><input type="checkbox"/></td><td>us-central1</td><td><div><div></div></div>1,500 GB (1.5 TB)</td><td><div><div></div></div>1,700 GB (1.7 TB)</td><td>102,400 GB (102.4 TB)</td><td>View hierarchy</td></tr><tr><td><input type="checkbox"/></td><td>us-east1</td><td><div><div></div></div>604 GB</td><td><div><div></div></div>604 GB</td><td>102,400 GB (102.4 TB)</td><td>View hierarchy</td></tr></tbody></table></div>	<input type="checkbox"/>	Location	Current Usage ↑	7 Day Peak Usage	Limit		<input type="checkbox"/>	northamerica-northeast1	<div><div></div></div> 20,498 GB (20.498 TB)	<div><div></div></div> 22,498 GB (22.498 TB)	102,400 GB (102.4 TB)	View hierarchy	<input type="checkbox"/>	us-central1	<div><div></div></div> 1,500 GB (1.5 TB)	<div><div></div></div> 1,700 GB (1.7 TB)	102,400 GB (102.4 TB)	View hierarchy	<input type="checkbox"/>	us-east1	<div><div></div></div> 604 GB	<div><div></div></div> 604 GB	102,400 GB (102.4 TB)	View hierarchy				
<input type="checkbox"/>	Location	Current Usage ↑	7 Day Peak Usage	Limit																									
<input type="checkbox"/>	northamerica-northeast1	<div><div></div></div> 20,498 GB (20.498 TB)	<div><div></div></div> 22,498 GB (22.498 TB)	102,400 GB (102.4 TB)	View hierarchy																								
<input type="checkbox"/>	us-central1	<div><div></div></div> 1,500 GB (1.5 TB)	<div><div></div></div> 1,700 GB (1.7 TB)	102,400 GB (102.4 TB)	View hierarchy																								
<input type="checkbox"/>	us-east1	<div><div></div></div> 604 GB	<div><div></div></div> 604 GB	102,400 GB (102.4 TB)	View hierarchy																								

Quotas

Persistent Disks

Python: Use the `--worker_disk_type` flag

```
$ python3 -m apache_beam.examples.wordcount \
--input gs://dataflow-samples/shakespeare/kinglear.txt \
--output gs://$BUCKET/results/outputs --runner DataflowRunner \
--project $PROJECT --temp_location gs://$BUCKET/tmp/ --region $REGION \
--worker_disk_type compute.googleapis.com/projects/$PROJECT/zones/$ZONE/diskTypes/pd-ssd
```

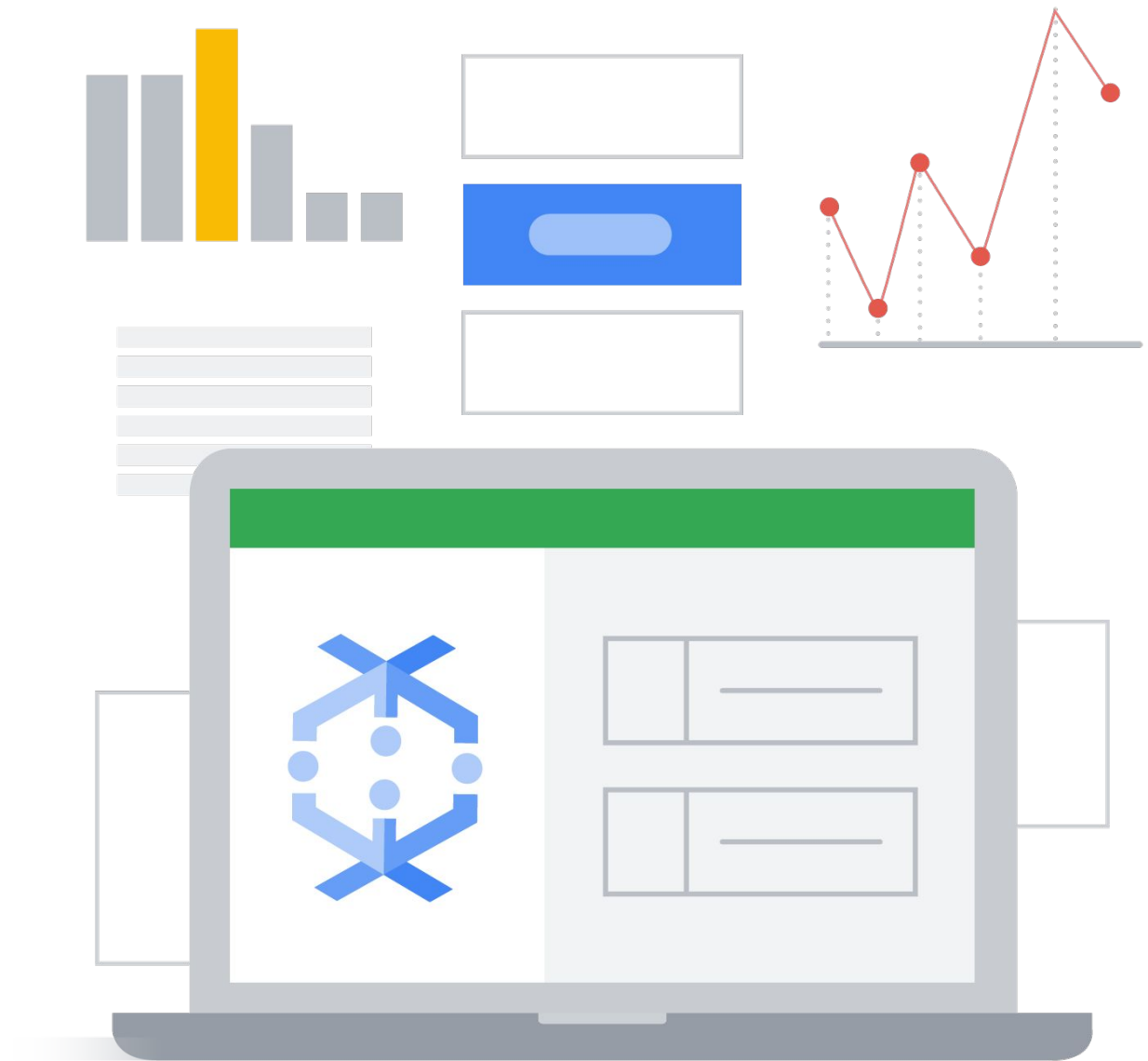
Java: Use the `--workerDiskType` flag

```
$ gradle clean execute -DmainClass=org.apache.beam.examples.WordCount -Dexec.args="\
--inputFile=gs://apache-beam-samples/shakespeare/kinglear.txt \
--output=gs://$BUCKET/results/outputs --runner=DataflowRunner \
--project=$PROJECT --tempLocation=gs://$BUCKET/tmp/ --region=$REGION \
--workerDiskType=compute.googleapis.com/projects/$PROJECT/zones/$ZONE/diskTypes/pd-ssd"
```

Quotas

Persistent Disks - Batch Pipeline

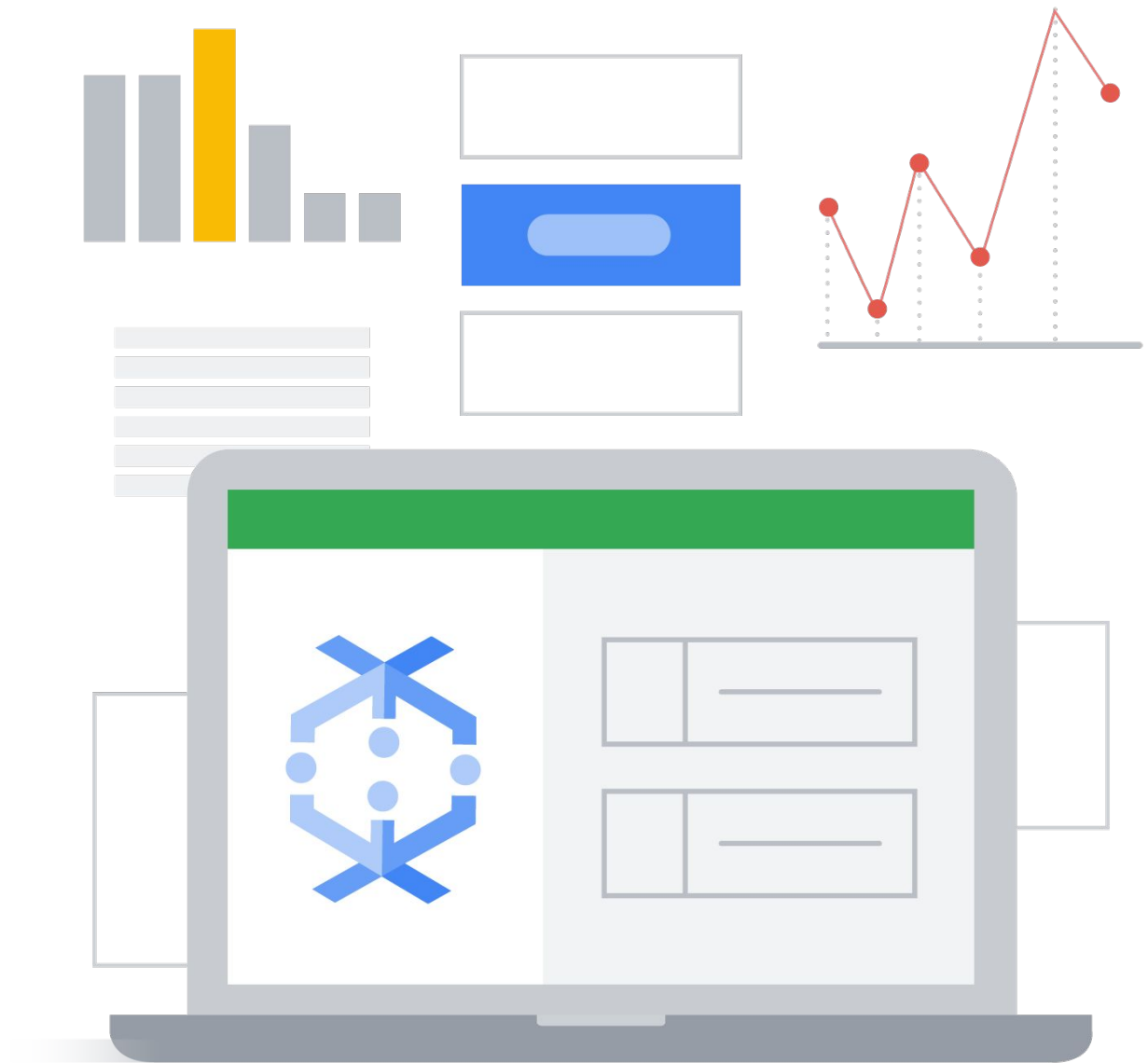
- VM to PD ratio is 1:1 for Batch



Quotas

Persistent Disks - Batch Pipeline

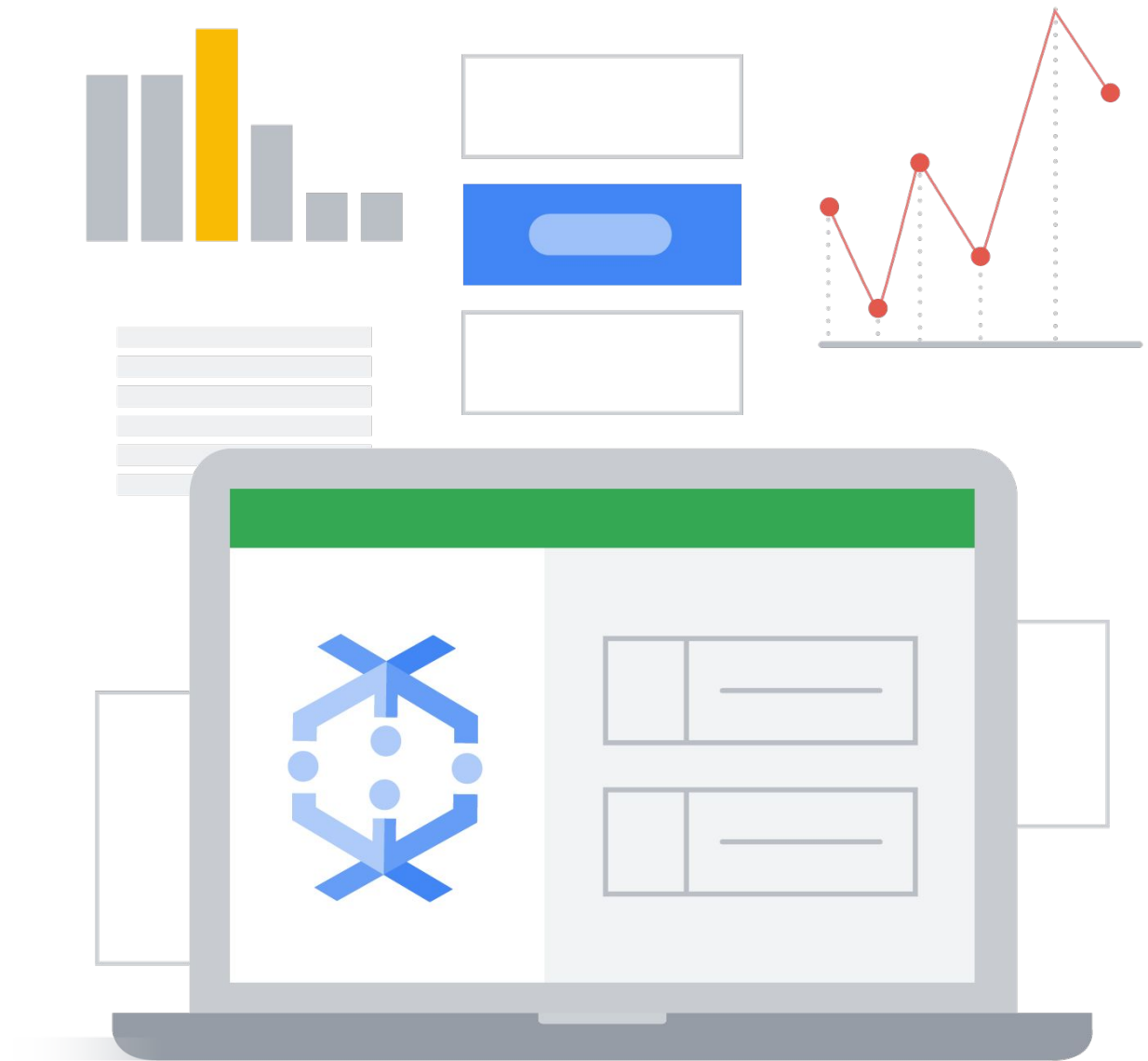
- VM to PD ratio is 1:1 for Batch
- Size if Shuffle on VM: 250 GB



Quotas

Persistent Disks - Batch Pipeline

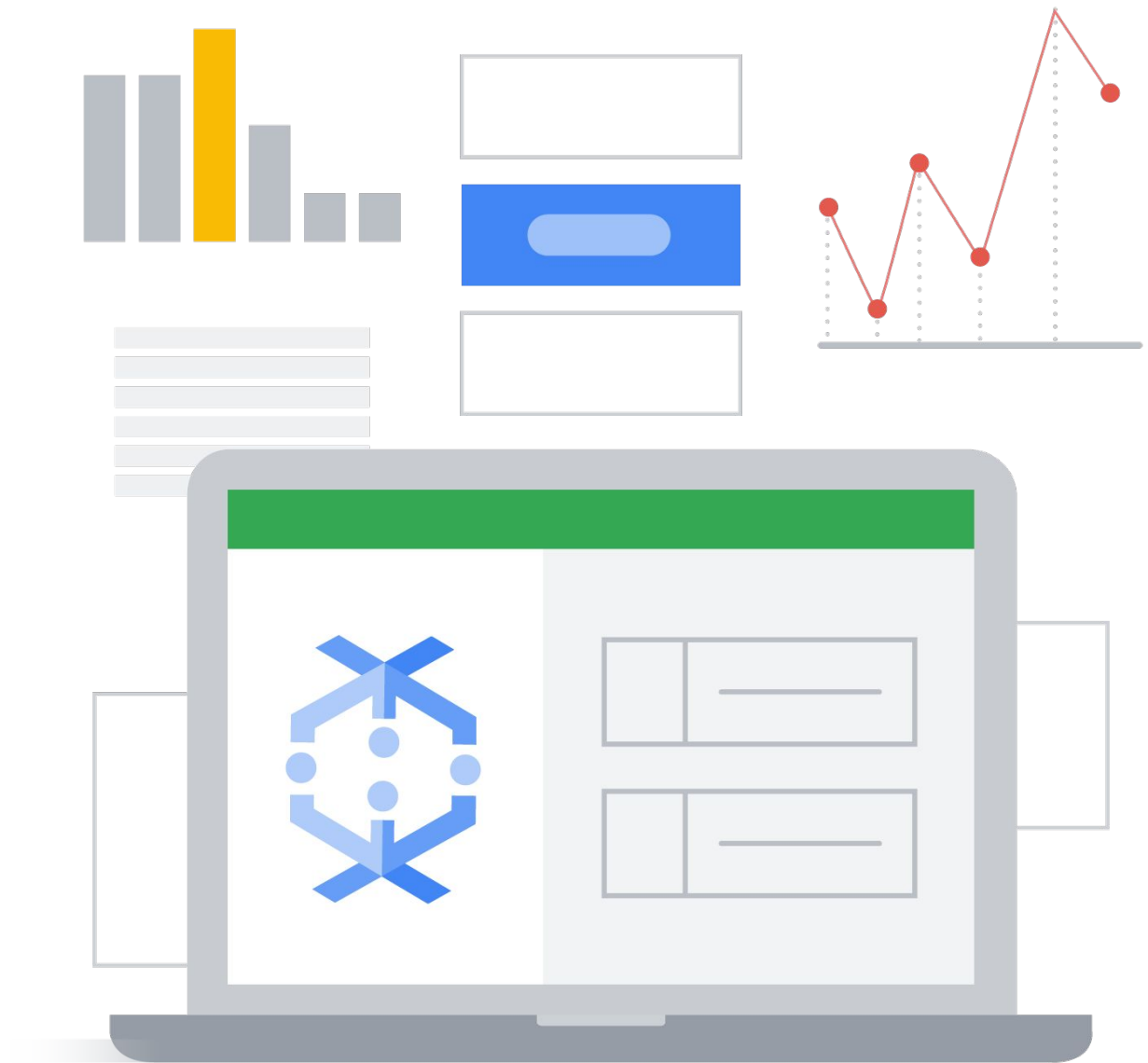
- VM to PD ratio is 1:1 for Batch
- Size if Shuffle on VM: 250 GB
- Size if Shuffle Service: 25 GB



Quotas

Persistent Disks - Batch Pipeline

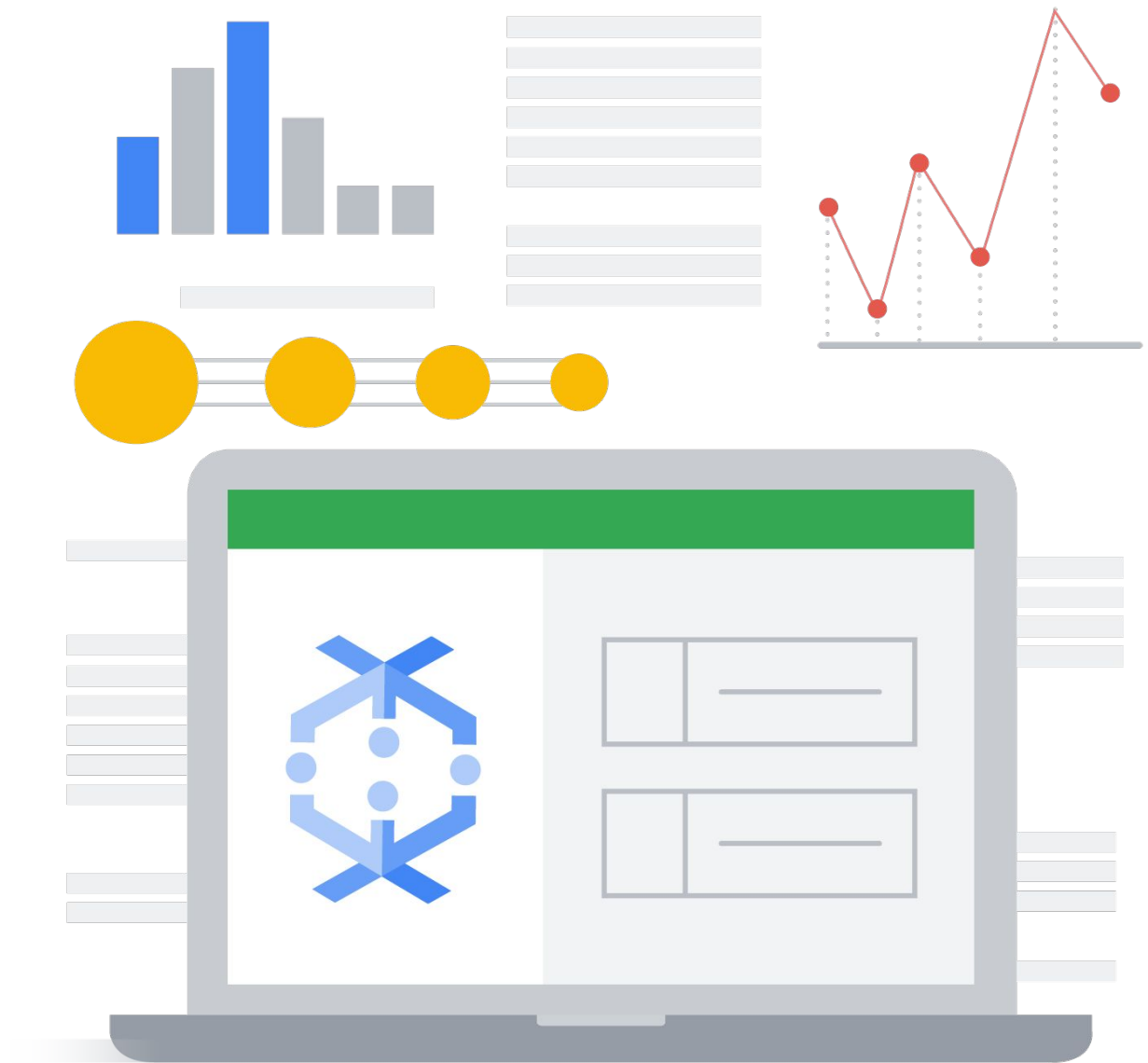
- VM to PD ratio is 1:1 for Batch
- Size if Shuffle on VM: 250 GB
- Size if Shuffle Service: 25 GB
- Flag to override default:
Python: `--disk_size_gb`
Java: `--diskSizeGb`



Quotas

Persistent Disks - Streaming Pipeline

- Fixed number of PDs



Quotas

Persistent Disks - Streaming Pipeline

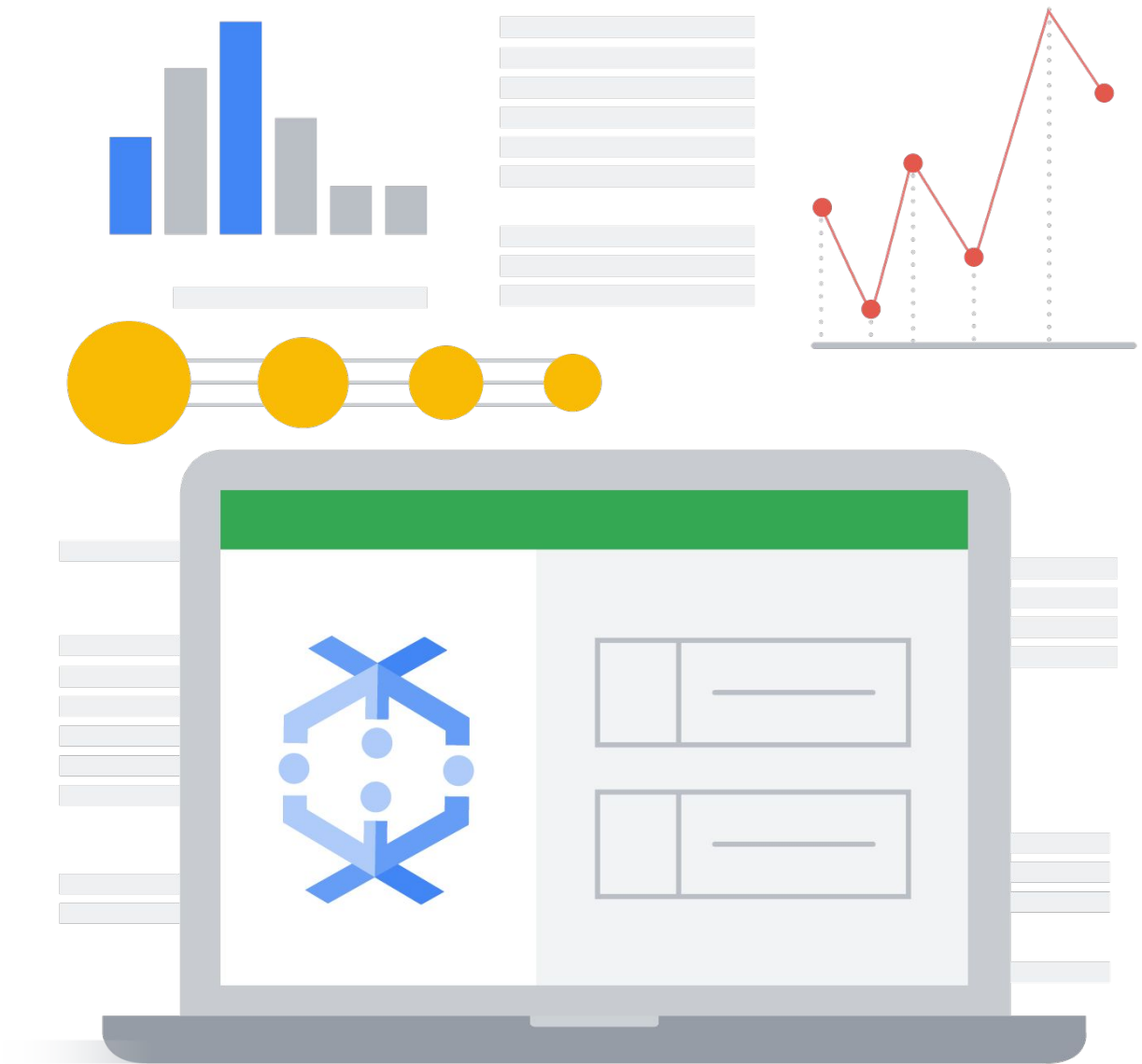
- Fixed number of PDs
- Default size if shuffle on VM: 400 GB



Quotas

Persistent Disks - Streaming Pipeline

- Fixed number of PDs
- Default size if shuffle on VM: 400 GB
- Default size if Streaming Engine: 30 GB



Quotas

Persistent Disks - Streaming Pipeline

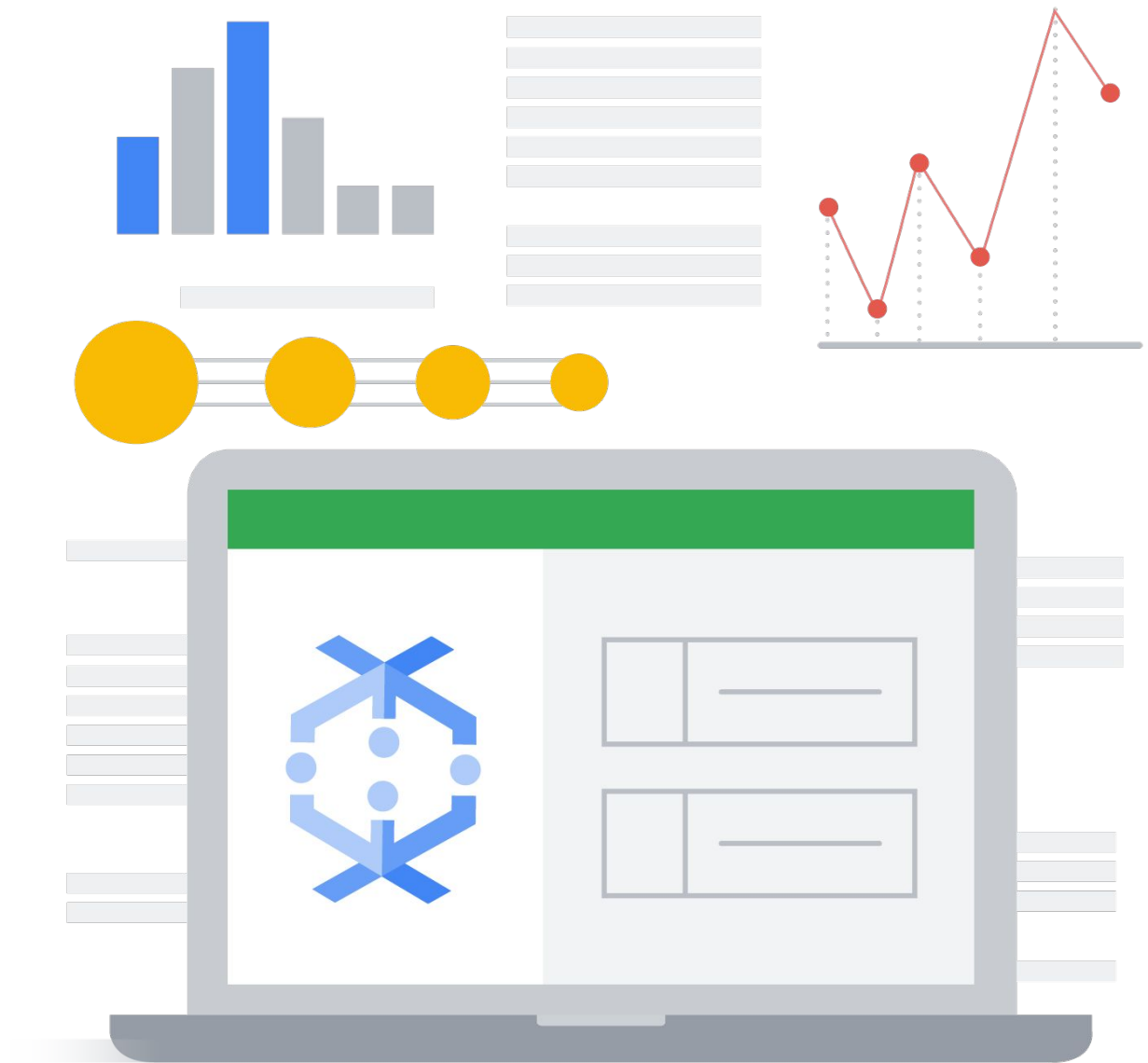
- Fixed number of PDs
- Default size if shuffle on VM: 400 GB
- Default size if Streaming Engine: 30 GB
- Flag to override default:
Python: **--disk_size_gb**
Java: **--diskSizeGb**



Quotas

Persistent Disks - Streaming Pipeline

- Amount of disk allocated == Maximum number of workers



Quotas

Persistent Disks - Streaming Pipeline

- Amount of disk allocated == Maximum number of workers
- Flag to set maximum number of workers:
Python: `--max_num_workers`
Java: `--maxNumWorkers`
- Flag required for streaming with shuffle on VMs
- Maximum number of workers that can be launched is 1000

