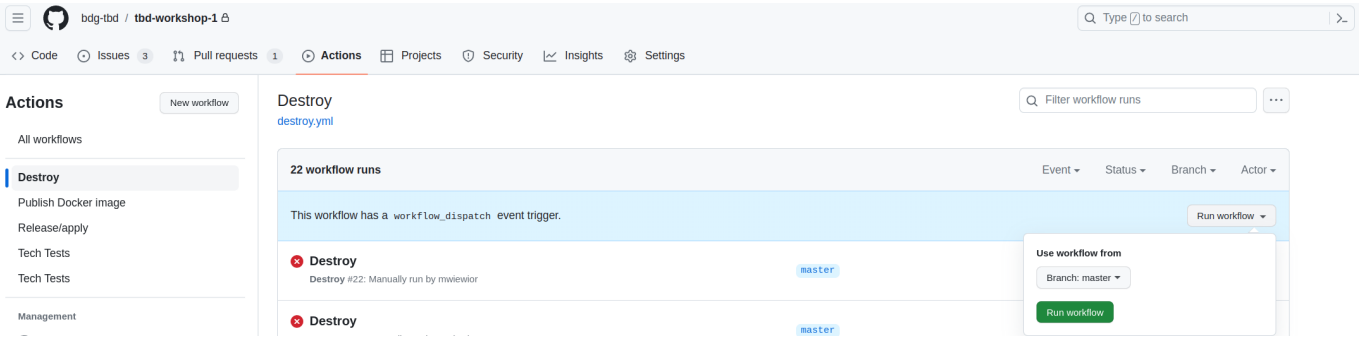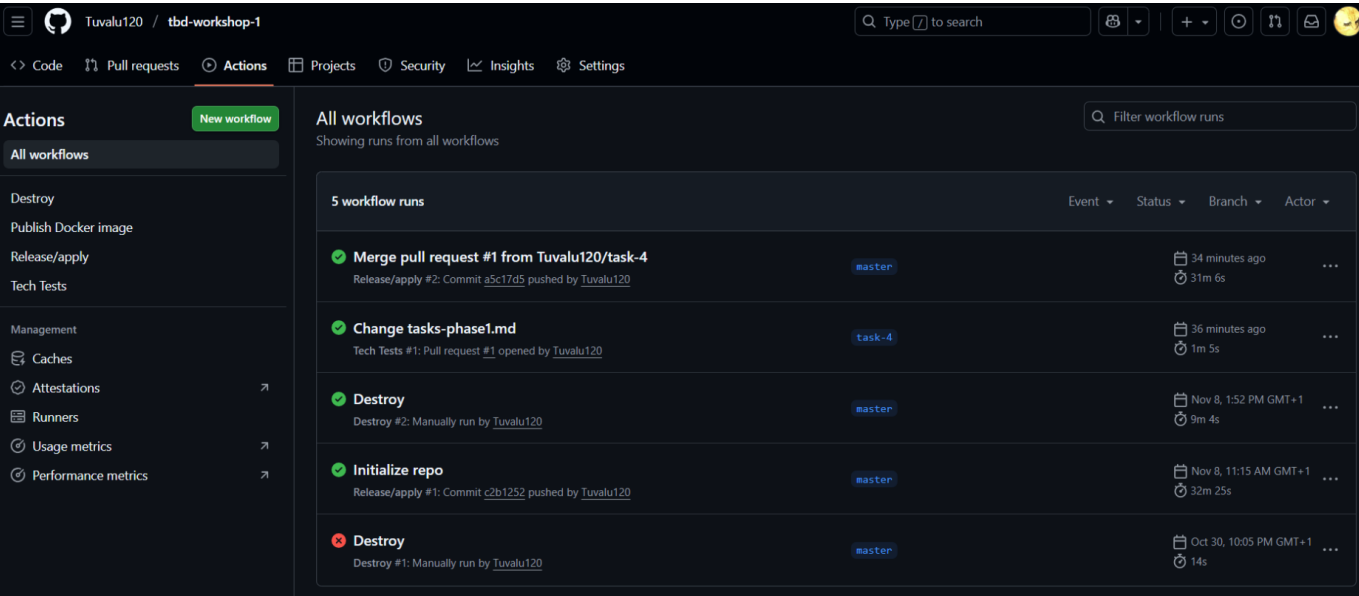IMPORTANT 🛑 Please remember to destroy all the resources after each work session. You can recreate infrastructure by creating new PR and merging it to master.



1. Authors:

   14

   https://github.com/Tuvalu120/tbd-workshop-1

2. Follow all steps in README.md.

3. From avaialble Github Actions select and run destroy on main branch.

4. Create new git branch and:
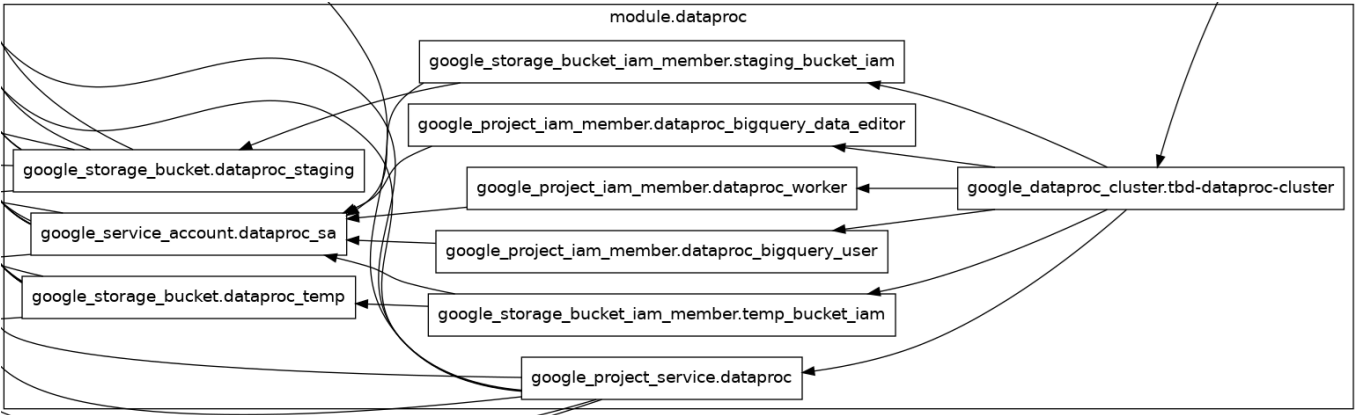
   1. Modify tasks-phase1.md file.

   2. Create PR from this branch to **YOUR** master and merge it to make new release.



5. Analyze terraform code. Play with terraform plan, terraform graph to investigate different modules.

   Moduł dataproc służy do stworzenia i konfiguracji klastra Apache Spark w wykorzystywanej w tym ćwiczeniu usłudze Google Cloud Dataproc. W architekturze naszego projektu pełni rolę środowiska obliczeniowego dla zadań przetwarzania danych oraz środowiska interaktywnego. Moduł ten jest ściśle powiązany z modułem sieciowym vpc (depends_on = [module.vpc]), czyli wymaga on uprzednio utworzonej sieci vpc. Zadaniami modułu jest zapewnienie, że sieć vpc jest gotowa, stworzenie klastra

w określonym projekcie (project_name = var.project_name) i regionie (var.region), umieszczenie węzłów klastra w odpowiedniej podsieci vpc (subnet = module.vpc.subnets[local.notebook_subnet_id].id). Dodatkowo instaluje on oprogramowanie oparte na obrazie Dataproc w wersji 2.2 z Ubuntu 22 (image_version = "2.2.69-ubuntu22") i używa maszyn wirtualnych typu e2-standard-2 (machine_type = "e2-standard-2").



6. Reach YARN UI

```
gcloud compute ssh --zone "europe-west1-b" "tbd-cluster-m" --tunnel-
through-iap --project "tbd-2025z-14" -- -L 8088:localhost:8088
```



7. Draw an architecture diagram (e.g. in draw.io) that includes:
   1. Description of the components of service accounts
   2. List of buckets for disposal

8. Create a new PR and add costs by entering the expected consumption into Infracost For all the resources of type: `google_artifact_registry`, `google_storage_bucket`, `google_service_networking_connection` create a sample usage profiles and add it to the Infracost task in CI/CD pipeline. Usage file example

```
google_artifact_registry_repository:
  storage_gb: 150
  monthly_egress_data_transfer_gb:
    europe_west1: 50

google_storage_bucket:
  storage_gb: 192
  monthly_class_a_operations: 100000
  monthly_class_b_operations: 500000
  monthly_data_retrieval_gb: 250
  monthly_egress_data_transfer_gb:
    same_continent: 100
    worldwide: 50

google_service_networking_connection:
  monthly_egress_data_transfer_gb:
    same_region: 10
    europe: 5
```

github-actions (bot) commented 33 minutes ago • edited ▾

💰 **Infracost report**

**Monthly estimate increased by $92** 📈

| Changed project | Baseline cost | Usage cost* | Total change | New monthly cost |
|---|---|---|---|---|
| Tuvalu120/tbd-workshop-1 | +$0 | +$65 | +$65 | $65 |
| Tuvalu120/tbd-workshop-1/bootstrap | +$0 | +$13 | +$13 | $13 |
| Tuvalu120/tbd-workshop-1/mlops | +$0.30 | +$14 | +$14 (+53%) | $42 |

*Usage costs were estimated using infracost-usage.yml, see docs for other options.

▼ Estimate details (includes details of unsupported resources)

Key: * usage cost, ~ changed, + added, - removed

9. Create a BigQuery dataset and an external table using SQL

```
CREATE SCHEMA IF NOT EXISTS dataset;

CREATE OR REPLACE EXTERNAL TABLE dataset.shakespeare
OPTIONS (

format = 'ORC',
uris = ['gs://tbd-2025z-14-data/data/shakespeare/.orc']);

SELECT FROM dataset.shakespeare
ORDER BY sum_word_count;
```

ORC to format plików używany do przechowywania danych w systemach Big Data. Jedną z jego głównych zalet jest to, iż pliki ORC przechowują informacje o swoim schemacie (nazwy kolumn i typy danych) wewnątrz struktury. Dzięki temu systemy takie jak BigQuery mogą odczytać plik i automatycznie odtworzyć strukturę danych, bez konieczności ręcznego definiowania schematu podczas tworzenia tabeli zewnętrznej.

10. Find and correct the error in spark-job.py

Problem był spowodowany tym, iż w pliku spark-job.py zawarty byłniewłaściwy Bucket. Informacje o złym Buckecie można było odczytać z logów nieudanego joba w konsoli Goggle Cloud i dataproc. Po zamianie nazwy Bucketa na:

```
DATA_BUCKET = "gs://tbd-2025z-14-data/data/shakespeare/"
```

job został wykonany pomyślnie.

```
The top words in shakespeare are
25/11/24 14:35:38 INFO GhfsGlobalStorageStatistics: periodic connector metrics: {action_http_delete_request=2, action_http_delete_request_duration=76,
[CONTEXT ratelimit_period="5 MINUTES" ]
+----+--------------+
|word|sum_word_count|
+----+--------------+
| the|         25568|
|   I|         21028|
| and|         19649|
|  to|         17361|
|  of|         16438|
|   a|         13409|
| you|         12527|
|  my|         11291|
|  in|         10589|
|  is|          8735|
|that|          8561|
| not|          8395|
|  me|          8030|
| And|          7780|
|with|          7224|
|  it|          7137|
| his|          6811|
|  be|          6724|
|your|          6244|
| for|          6154|
+----+--------------+
only showing top 20 rows
```

| | Job ID | Status | Region | Type | Cluster | Start time |
|---|---|---|---|---|---|---|
| ☐ | d53d2e1f-b099-41d4-a9dd-805aa2eb27c9 | ✅ Succeeded | europe-west1 | PySpark | tbd-cluster | 24 Nov 2025, |

11. Add support for preemptible/spot instances in a Dataproc cluster

```
[modules/dataproc/main.tf](modules/dataproc/main.tf)

~~~
preemptible_worker_config {
  num_instances = 2
  preemptibility = "SPOT"

  disk_config {
    boot_disk_type    = "pd-standard"
    boot_disk_size_gb = 100
  }
}
~~~
```

12. Triggered Terraform Destroy on Schedule or After PR Merge. Goal: make sure we never forget to clean up resources and burn money.

Add a new GitHub Actions workflow that:

1. runs terraform destroy -auto-approve
2. triggers automatically:    a) on a fixed schedule (e.g. every day at 20:00 UTC)    b) when a PR is merged to main containing [CLEANUP] tag in title

Steps:

1. Create file .github/workflows/auto-destroy.yml
2. Configure it to authenticate and destroy Terraform resources
3. Test the trigger (schedule or cleanup-tagged PR)

[.github/workflows/auto-destroy.yml](#)





Automatycznie czyszczenie środowiska o zadanej godzinie pomaga zminimalizować zużycie zasobów, kiedy nie są one używane (np. kiedy ktoś zapomniał o jego czyszczeniu po skończonej pracy), a czyszczenie środowiska z uzyciem taga [CLEANUP] znacząco ułatwia osobom zajmującym się projektem uporządkowanie zasobów po skończonej pracy nad projektem.