Google Cloud
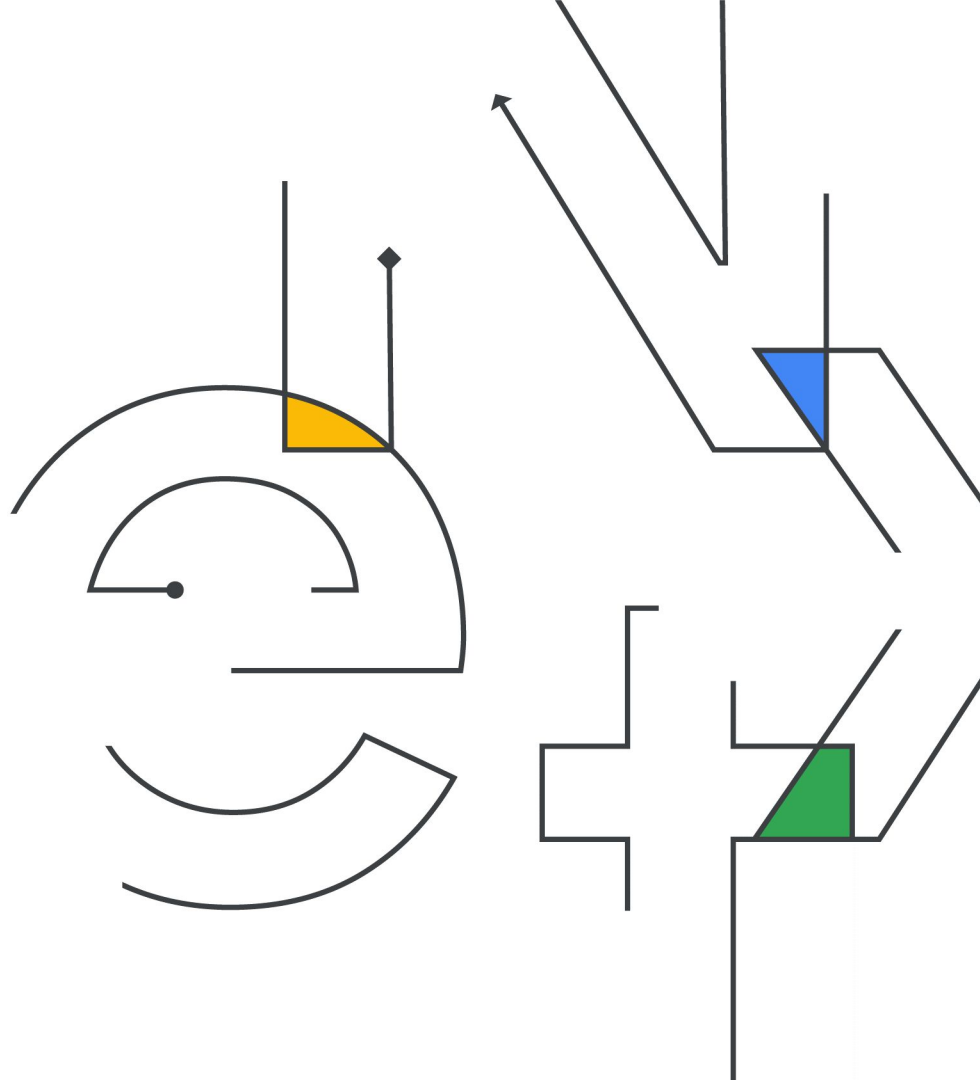
Next '22

# Better Hardware Provisioning for ML Experiments on GCP

Provision hardware for ML reliably.

Oct/
11–13

# Sayak Paul

## ML Engineer at Carted

October 13, 2022

"

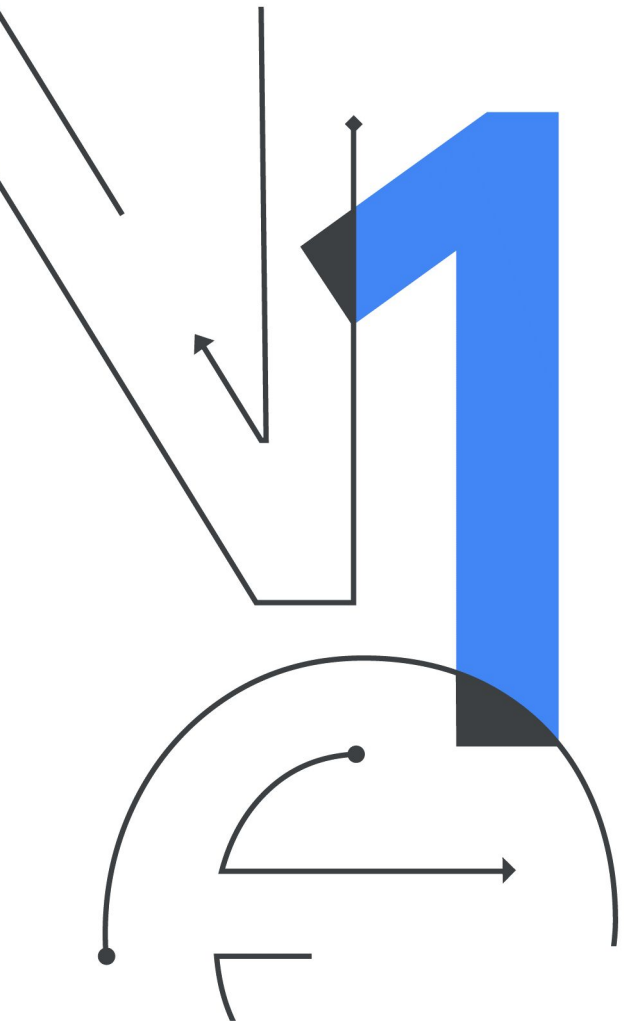I'm presenting work done by my colleagues (Nilabhra, Shane, Sam) and myself.

# Contents

Google Cloud
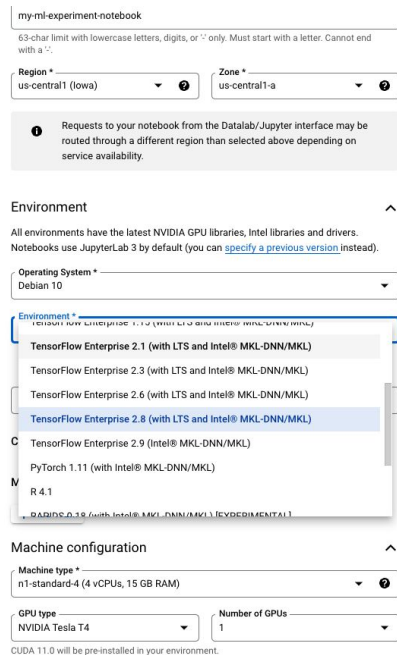
# The standard workflow

# Components of ML experiments

We need:

- Pre-configured environment (with libraries and CUDA stuff sorted).
- Notebook instance for better interaction with the modules.
- Access to a GPU or GPUs.
- Optionally a central storage location for storing experiment artifacts.

Google Cloud

# Vertex AI and GCS to the rescue

Vertex AI Workbench allows to do this with a few clicks!

# Vertex AI and GCS to the rescue

Same applies to <u>GCS</u> too!





Google Cloud

# Problems with the typical workflow

# Handling hardware components

- **Manual efforts**
  - Redundancy from repeating the instructions every time a new instance needs to be created.
  - What if a step is missed out?
- **Easy to forget about cleaning up the resources after experimentation**
  - Multiple moving pieces (notebook instance, GCS bucket, etc.)

Google Cloud

"Provisioning infrastructure through point-and-click GUI's or custom scripts is slow, error-prone, inefficient, and doesn't scale."

Terraform

# Further extensions

- Supporting custom dependencies
- Supporting Python virtual environments

Google Cloud

# Observations

- Bad:
  - Multiple interconnected steps that can be easily missed out and wreak havoc
- Good:
  - Single commands to provision and deprovision the resources

# Enter
# Terraform

With HashiCorp Terraform, provisioning and security can be automated based on infrastructure and policy as code. Infrastructure and policies are codified, shared, managed, and executed within a workflow that is consistent across all infrastructure.

Terraform

# Realize the same workflow but with code

```
Define variables (names,          Specify what should be              terraform init
GPU type, machine         →       provisioned (notebook,      →       terraform plan
type, etc.)                       bucket, etc.) with          terraform apply —auto-approve
                                  variables
```

Google Cloud

# Realize the same workflow but with code

```
Define variables (names, GPU type, machine type, etc.)
```
→
```
Specify what should be provisioned (notebook, bucket, etc.) with variables
```
→
```
terraform init
terraform plan
terraform apply —auto-approve
```

```
Optionally, supply a startup script
```

```
terraform destroy
```

Google Cloud

# Standard anatomy for Terraform code

- `variables.tf` **(define variables)**
- `main.tf` **(specify what should be provisioned with** `variables.tf`**)**
- **scripts/**
  - `auto_shutdown.sh` **(configure automatic shutdown)**
  - `dependencies.sh` **(configure additional dependencies)**

`.tf` is the extension for HCL (Hashicorp Configuration Language)

# The `variables.tf` file

- **GCP related information**
  - Project ID
  - Location
  - Region
- **Notebook related information**
  - Name
  - Machine type
  - GPU type and count
  - Image family
- **GCS related information**
  - Name
  - Region

# The `variables.tf` file

- GCP related information

```
variable "gcp" {
  type = object(
    {
      project_id = string
      location   = string
      region     = string
    }
  )
  default = {
    project_id = "my-gcp-project"
    location   = "us-central1-a"
    region     = "us-central1"
  }
}
```

# The `variables.tf` file

- Notebook related information

```
variable "notebook" {
  type = object(
    {
      notebook_name      = string
      machine_type       = string
      gpu_type           = string

      …
    }
  )
  default = {
    notebook_name       = "my-ml-nb"
    machine_type        = "n1-standard-8"
    gpu_type            = "NVIDIA_TESLA_T4"

    …
  }
}
```

# The `variables.tf` file

- **GCS related information**

```
variable "gcs" {
  type = object(
    {
      name         = string
      location     = string
      force_destroy = bool
    }
  )
  default = {
    name         = "my-ml-bucket"
    location     = "us-central1"
    force_destroy = false
  }
}
```

Google Cloud

# The `main.tf` file

- Provision notebook (Vertex AI Workbench) resource

```
resource "google_notebooks_instance" "notebook_instance" {
  project      = var.gcp.project_id
  name         = var.notebook.notebook_name
  machine_type = var.notebook.machine_type
  location     = var.gcp.location
  ...
}
```

# The `main.tf` file

- Provision notebook (Vertex AI Workbench) resource

```
resource "google_notebooks_instance" "notebook_instance" {
  ...
  network = ...
  subnet  = ...

  vm_image {
    project      = local.image_project
    image_family = var.notebook.image_family
  }

  accelerator_config {
    type       = var.notebook.gpu_type
    core_count = var.notebook.gpu_count
  }
  ...
}
```

Google Cloud

# The `main.tf` file

- Provision notebook (Vertex AI Workbench) resource

```
resource "google_notebooks_instance" "notebook_instance" {
  ...
  install_gpu_driver = var.notebook.install_gpu_driver
  boot_disk_size_gb  = ...
  data_disk_size_gb  = ...
}
```

# The `main.tf` file

- Provision GCS bucket resource

```
resource "google_storage_bucket" "default" {
  project                     = var.gcp.project_id
  name                        = var.gcs.name
  location                    = var.gcs.location
  storage_class               = "REGIONAL"
  force_destroy               = var.gcs.force_destroy
  uniform_bucket_level_access = true
}
```

# Optionals

- Pass additional scripts to
  - Install additional dependencies
  - Automatically shutdown the instance w.r.t low CPU utilization
- Pass a URI of a predefined container image (`container_image`)
- Store the Terraform states to a remote location for better concurrency and collaboration

Google Cloud

# State management: Shell-scripts vs. Terraform

- Terraform comes with state management
    - Keeps track of what's successfully provisioned or failed to provision
    - Control over the entire lifecycle of the infrastructure
    - Easy "planning" of the hardware to be provisioned for previewing changes
    - Enforces version control providing ease of collaboration

Google Cloud

[bit.ly/next-22-tfm](bit.ly/next-22-tfm)

# Demo

# Wrapping up

- Manual hardware provisioning is faulty, introducing redundancy and rough edges.
- Streamlining the provisioning workflow establishes DRY and helps teams to have reproducibility when provisioning infrastructure.
- With infrastructure as code, we get:
  - Version control
  - Maintainability
  - Reusability

Google Cloud

# Thank you

Sayak Paul (@RisingSayak)

bit.ly/next-22-tfm

Google Cloud

Next '22