# Kubeflow *as-a-service* on HPC clusters – first experience

Mohsin Ahmed Shaikh[1], Nasr Hassanien[2], Islam Elmas[2], Saber Feki[1]

1: KAUST Supercomputing Lab, 2: Brightskies Inc

Correspondance: mohsin.shaikh@kaust.edu.sa

# KAUST Supercomputing Core Lab

- Shaheen 2 (Flagship)
  - Cray XC40: 6174 nodes (Intel Haswell 32 cores 128GB)
- Ibex cluster:
  - Heterogeneous CPU (Intel/AMD): > 22K cores (250GB-3TB memory)
  - Heterogeneous GPU (NVIDIA): > 600 GPUs
- Shaheen 3 (Flagship refresh – Coming Soon)
  - HPE Cray EX4000: CPUs + GPUs

**SLURM Everywhere**

# User Personas – by jobtypes

- SLURM batch jobs
  - Large scale HPC (CFD, CompChem, Earth Sci, GeoSci, Bioscience, Datascience)
  - Singleton, arrays, dependency, burst buffer
  - Client/Server (Dask, Ray)
- SLURM Interactive jobs
  - Datascience on GPUs – model development
  - Jupyter, VS Code

# Motivation

- Ease-of-use for Interactive computing

- Enable resilient + reproducible workflows

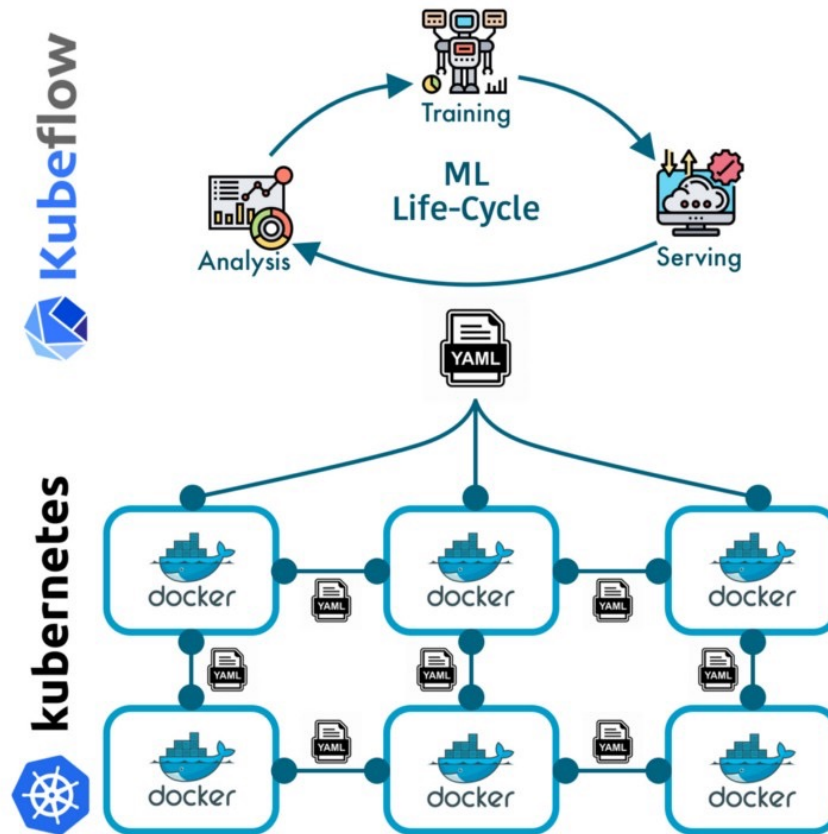- Cross-platform portability -- workstation/cloud/HPC cluster

# Kubeflow

# Kubeflow

Kubeflow is a framework for developing data science models and workflows

Depends on Kubernetes and abstract its use through a GUI

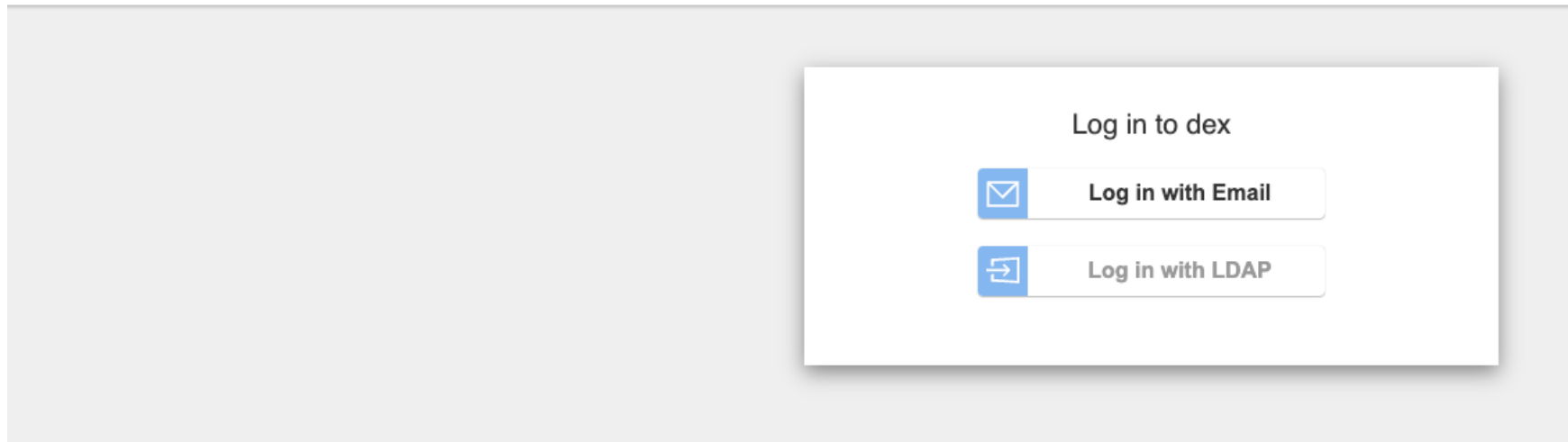Kubeflow *as a service* provides:

- Interactive notebooks
- Kubeflow pipelines – create reproducible workflows
- Katib – hyperparameter optimization experiments
- Kserve – model serving platform

# Functionality tested so far …

- Interactive computing
- Distributed training

# Kubeflow Authentication

# Main dashboard

# Persistent Volume -- attachable

# Creating Notebook

# Creating Notebook

# Creating Notebook

# Creating Notebook

# Tensorboard service

# Distributed training
## Trianing Opreator -- PytorchJob

```yaml
                                                containers:
                                                  - name: pytorch
                                                    image: docker.io/mshaikh/kubeflow-demo:kf-dist-torch-vanilla
                                                    imagePullPolicy: IfNotPresent
                                                    command: [ "torchrun",
                                                              "--nnodes", "1",
                                                              "--nproc_per_node","4",
                                                              "--node_rank","0",
                                                              "/workspace/mnist/src/ddp.py",
apiVersion: "kubeflow.org/v1"                              "--batch-size","32",
kind: "PyTorchJob"                                         "--num-worker","8",
metadata:                                                  "--epochs","4",
   name: "DDP"                                             "--lr","0.001" ]
   namespace: training01                            env:
spec:                                                 - name: 'NCCL_DEBUG'
   pytorchReplicaSpecs:                                 value: 'INFO'
      Worker:                                         - name: 'DATA_DIR'
         replicas: 1                                    value: '/data/tiny-imagenet-200'
       restartPolicy: Never                           - name: 'OMP_NUM_THREADS'
       template:                                        value: '1'
         metadata:                                   resources:
           annotations:                                limits:
             sidecar.istio.io/inject: "false"            cpu: 8
         spec:                                           memory: '200Gi'
            affinity:                                    nvidia.com/gpu: 4
               nodeAffinity:                        volumeMounts:
               requiredDuringSchedulingIgnoredDuringExecution:
                nodeSelectorTerms:                     - name: dshm
                  - matchExpressions:                    mountPath: /dev/shm
                   - key: nvidia.com/gpu.product  volumes:
                     operator: In                    - emptyDir:
                     values:                             medium: Memory
                       - "Tesla-V100-SXM2-32GB"       name: dshm
```

```
+ NVIDIA-SMI 535.86.10              Driver Version: 535.86.10    CUDA Version: 12.2    |
|-----------------------------------------+----------------------+----------------------+
| GPU  Name                 Persistence-M | Bus-Id        Disp.A | Volatile Uncorr. ECC |
| Fan  Temp   Perf          Pwr:Usage/Cap |         Memory-Usage | GPU-Util  Compute M. |
|                                         |                      |               MIG M. |
|=========================================+======================+======================|
|   0  Tesla V100-SXM2-32GB          Off | 00000000:61:00.0 Off |                    0 |
| N/A   62C    P0             190W / 300W |   5071MiB / 32768MiB |    93%       Default |
|                                         |                      |                  N/A |
+-----------------------------------------+----------------------+----------------------+
|   1  Tesla V100-SXM2-32GB          Off | 00000000:62:00.0 Off |                    0 |
| N/A   61C    P0             226W / 300W |   5031MiB / 32768MiB |    94%       Default |
|                                         |                      |                  N/A |
+-----------------------------------------+----------------------+----------------------+
|   2  Tesla V100-SXM2-32GB          Off | 00000000:89:00.0 Off |                    0 |
| N/A   61C    P0             246W / 300W |   5031MiB / 32768MiB |    94%       Default |
|                                         |                      |                  N/A |
+-----------------------------------------+----------------------+----------------------+
|   3  Tesla V100-SXM2-32GB          Off | 00000000:8A:00.0 Off |                    0 |
| N/A   66C    P0             279W / 300W |   5039MiB / 32768MiB |    95%       Default |
|                                         |                      |                  N/A |
+-----------------------------------------+----------------------+----------------------+

+-----------------------------------------------------------------------------------------+
| Processes:                                                                              |
|  GPU   GI   CI        PID   Type   Process name                             GPU Memory |
|        ID   ID                                                              Usage      |
|=========================================================================================|
|    0   N/A  N/A    2745778      C   /opt/conda/bin/python                     5066MiB |
|    1   N/A  N/A    2745779      C   /opt/conda/bin/python                     5026MiB |
|    2   N/A  N/A    2745780      C   /opt/conda/bin/python                     5026MiB |
|    3   N/A  N/A    2745781      C   /opt/conda/bin/python                     5034MiB |
```

# SLURM vs Kubeflow

- PyTorch DDP training
- Dataset: TinyImageNet200
- Model: ResNet50
- Nodes: 1
- GPUs per node: 4 V100 SXM2 (32GB)

- Training for 20 epochs
- Batch size: 256

| SLURM | | Kubeflow training opreator | |
|---|---|---|---|
| Time to solution(s) | Accuracy | Time to solution(s) | Accuracy |
| 1939.232 | 48% | 1881.21 | 48% |

# Cluster Infrastructure

# Cloud cluster



| Provision | Terraform |
|-----------|-----------|
| Manage | Ansible |
| OS | CentOS 8 |
| K8s Distro | k3s |
| Container Runtime | containerd |
| Storage class | NFS |

Active Directory

Ansible host

Control place

VCN

CPU worker nodes

GPU worker nodes

Gitlab CI/CD

KAUST
SUPERCOMPUTING
CORE LAB

# On-Prem cluster



| | |
|---|---|
| Provision | None |
| Manage | Ansible |
| OS | Rocky Linux 9.2 |
| K8s Distro | RKE2 |
| Container Runtime | containerd |
| Storage class | Longhorn on node-local storage |

Central AD

LDAP

Ansible host

Control place

Ethernet

CPU worker nodes

GPU worker nodes

Gitlab CI/CD

KAUST SUPERCOMPUTING CORE LAB

# Lessons Learnt

- Master in HA

- K3s vs RKE2 (mysql vs etcd)

- DiskPressure – `CrashLoopBackOff`

- Multi-node cluster is important for PoC

- NVIDIA GPU operator vs Device plugin

- Node feature discovery add-on

- Appropriate Container Network Interface (Ethernet + IB + …..)

# Lessons Learnt

- Identity management in a container

  - Per user image vs injecting credentials via Configmaps

- The `kubectl` CLI use in Jupyter session – restricted to namespace

- Customize base images form Kubeflow

  - https://www.kubeflow.org/docs/components/notebooks/container-images/

- The Kubeflow UI is useful, until there is a problem. `kubectl` access is needed to investigate what went wrong. That's a user support ticket for HPC centers.

- The `kubectl describe` output and `kubectl logs` are not always helpful

- For high throughput workflows (with SPMD), `restartPolicy=OnFailure` provides resiliency

# Thank you

Questions and contact:
[mohsin.shaikh@kaust.edu.sa](mailto:mohsin.shaikh@kaust.edu.sa)