

Triton Server accelerates distribution of models based on Dragonfly

Installation

Prerequisites

Dragonfly Kubernetes Cluster Setup

[Prepare Kubernetes Cluster](#)

[Kind loads dragonfly image](#)

[Create dragonfly cluster based on helm charts](#)

[Expose the Proxy service port](#)

Install Dragonfly Repository Agent

[Set Dragonfly Repository Agent configuration](#)

[Set Model Repository configuration](#)

Triton Server integrates Dragonfly Repository Agent plugin

[Install Triton Server with Docker](#)

[Verify](#)

Performance testing

相关链接

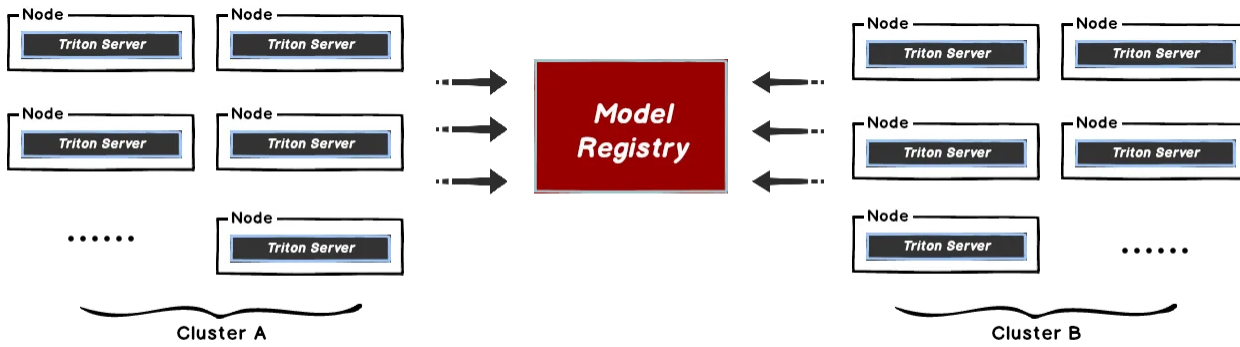
[Dragonfly 社区](#)

[NVIDIA Triton Inference Server](#)

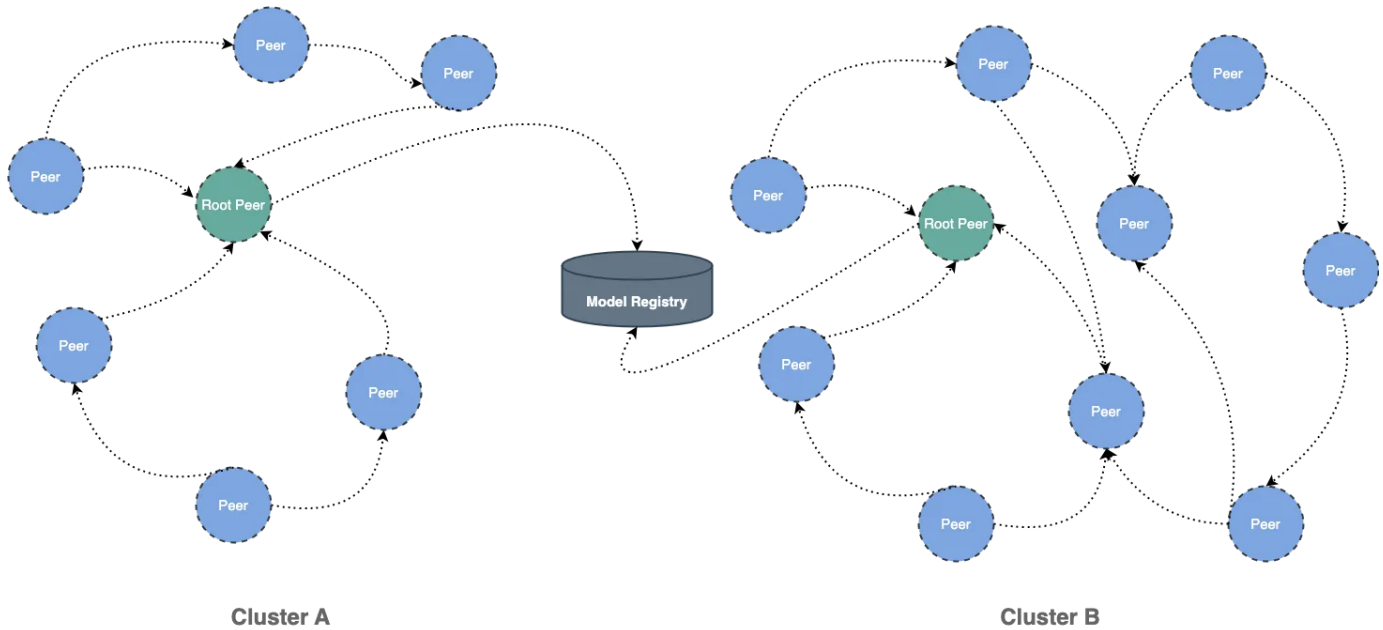
[二维码](#)

Author: Yufei Chen, Miao Hao, Min Huang

This document will help you experience how to use dragonfly with [TritonServe](#). During the downloading of models, the file size is large and there are many services downloading the files at the same time. The bandwidth of the storage will reach the limit and the download will be slow.



Dragonfly can be used to eliminate the bandwidth limit of the storage through P2P technology, thereby accelerating file downloading.



Installation

By integrating Dragonfly Repository Agent into Triton, download traffic through Dragonfly to pull models stored in S3, OSS, GCS, and ABS, and register models in Triton. The Dragonfly Repository Agent is in the dragonfly-repository-agent repository.

Prerequisites

Name	Version	Document
Kubernetes cluster	1.20+	kubernetes.io
Helm	3.8.0+	helm.sh

Triton Server	23.08-py3	Triton Server
---------------	-----------	---------------

Notice: [Kind](#) is recommended if no kubernetes cluster is available for testing.

Dragonfly Kubernetes Cluster Setup

For detailed installation documentation, please refer to [quick-start-kubernetes](#).

Prepare Kubernetes Cluster

Create kind multi-node cluster configuration file `kind-config.yaml`, configuration content is as follows:

```
▼ kind-config.yaml YAML |  
1 kind: Cluster  
2 apiVersion: kind.x-k8s.io/v1alpha4  
3 nodes:  
4   - role: control-plane  
5   - role: worker  
6   - role: worker
```

Create a kind multi-node cluster using the configuration file:

```
▼ Shell |  
1 kind create cluster --config kind-config.yaml
```

Switch the context of kubectl to kind cluster:

```
▼ Shell |  
1 kubectl config use-context kind-kind
```

Kind loads dragonfly image

Pull dragonfly latest images:

```
▼ Shell |
1  docker pull dragonflyoss/scheduler:latest
2  docker pull dragonflyoss/manager:latest
3  docker pull dragonflyoss/dfdaemon:latest
```

Kind cluster loads dragonfly latest images:

```
▼ Shell |
1  kind load docker-image dragonflyoss/scheduler:latest
2  kind load docker-image dragonflyoss/manager:latest
3  kind load docker-image dragonflyoss/dfdaemon:latest
```

Create dragonfly cluster based on helm charts

Create helm charts configuration file `charts-config.yaml` and set `dfdaemon.config.agents.regx` to match the download path of the object storage. Example: add `regx:.*/models.*` to match download request from object storage bucket `models`. Configuration content is as follows:

```
1 scheduler:
2   image: dragonflyoss/scheduler
3   tag: latest
4   replicas: 1
5   metrics:
6     enable: true
7   config:
8     verbose: true
9     pprofPort: 18066
10
11 seedPeer:
12   image: dragonflyoss/dfdaemon
13   tag: latest
14   replicas: 1
15   metrics:
16     enable: true
17   config:
18     verbose: true
19     pprofPort: 18066
20
21 dfdaemon:
22   image: dragonflyoss/dfdaemon
23   tag: latest
24   metrics:
25     enable: true
26   config:
27     verbose: true
28     pprofPort: 18066
29   proxy:
30     defaultFilter: 'Expires&Signature&ns'
31     security:
32       insecure: true
33       cacert: ''
34       cert: ''
35       key: ''
36     tcpListen:
37       namespace: ''
38       port: 65001
39     registryMirror:
40       url: https://index.docker.io
41       insecure: true
42       certs: []
43       direct: false
44   proxies:
45     - regx: blobs/sha256.*
```

```
46         # Proxy all http download requests of model bucket path.
47         - regx: .*models.*
48
49 manager:
50   image: dragonflyoss/manager
51   tag: latest
52   replicas: 1
53   metrics:
54     enable: true
55   config:
56     verbose: true
57     pprofPort: 18066
58
59 jaeger:
60   enable: true
```

Create a dragonfly cluster using the configuration file:

```
1 $ helm repo add dragonfly https://dragonflyoss.github.io/helm-charts/
2 $ helm install --wait --create-namespace --namespace dragonfly-system drag
  onfly dragonfly/dragonfly -f charts-config.yaml
3 LAST DEPLOYED: Wed Nov 29 21:23:48 2023
4 NAMESPACE: dragonfly-system
5 STATUS: deployed
6 REVISION: 1
7 TEST SUITE: None
8 NOTES:
9 1. Get the scheduler address by running these commands:
10 export SCHEDULER_POD_NAME=$(kubectl get pods --namespace dragonfly-syste
  m -l "app=dragonfly,release=dragonfly,component=scheduler" -o jsonpath={.i
  tems[0].metadata.name})
11 export SCHEDULER_CONTAINER_PORT=$(kubectl get pod --namespace dragonfly-s
  ystem $SCHEDULER_POD_NAME -o jsonpath="{.spec.containers[0].ports[0].cont
  ainerPort}")
12 kubectl --namespace dragonfly-system port-forward $SCHEDULER_POD_NAME 80
  02:$SCHEDULER_CONTAINER_PORT
13 echo "Visit http://127.0.0.1:8002 to use your scheduler"
14
15 2. Get the dfdaemon port by running these commands:
16 export DFDAEMON_POD_NAME=$(kubectl get pods --namespace dragonfly-syste
  m -l "app=dragonfly,release=dragonfly,component=dfdaemon" -o jsonpath={.it
  ems[0].metadata.name})
17 export DFDAEMON_CONTAINER_PORT=$(kubectl get pod --namespace dragonfly-s
  ystem $DFDAEMON_POD_NAME -o jsonpath="{.spec.containers[0].ports[0].contai
  nerPort}")
18 You can use $DFDAEMON_CONTAINER_PORT as a proxy port in Node.
19
20 3. Configure runtime to use dragonfly:
21 https://d7y.io/docs/getting-started/quick-start/kubernetes/
22
23
24 4. Get Jaeger query URL by running these commands:
25 export JAEGER_QUERY_PORT=$(kubectl --namespace dragonfly-system get serv
  ices dragonfly-jaeger-query -o jsonpath="{.spec.ports[0].port}")
26 kubectl --namespace dragonfly-system port-forward service/dragonfly-jaeg
  er-query 16686:$JAEGER_QUERY_PORT
27 echo "Visit http://127.0.0.1:16686/search?limit=20&lookback=1h&maxDurati
  on&minDuration&service=dragonfly to query download events"
```

Check that dragonfly is deployed successfully:

```
Shell |
1 $ kubectl get pods -n dragonfly-system
2 NAME                                READY   STATUS    RESTARTS   AGE
3 dragonfly-dfdaemon-8qcpd            1/1    Running   4 (118s ago) 2m45s
4 dragonfly-dfdaemon-qhkn8            1/1    Running   4 (108s ago) 2m45s
5 dragonfly-jaeger-6c44dc44b9-dfjfv   1/1    Running   0           2m45s
6 dragonfly-manager-549cd546b9-ps5tf  1/1    Running   0           2m45s
7 dragonfly-mysql-0                  1/1    Running   0           2m45s
8 dragonfly-redis-master-0            1/1    Running   0           2m45s
9 dragonfly-redis-replicas-0          1/1    Running   0           2m45s
10 dragonfly-redis-replicas-1          1/1    Running   0           2m7s
11 dragonfly-redis-replicas-2          1/1    Running   0           101s
12 dragonfly-scheduler-0              1/1    Running   0           2m45s
13 dragonfly-seed-peer-0              1/1    Running   1 (52s ago) 2m45s
```

Expose the Proxy service port

Create the `dfstore.yaml` configuration file to expose the port on which the Dragonfly Peer's HTTP proxy listens. The default port is `65001` and set `targetPort` to `65001`.


```
dfstore.yaml | YAML |
1  kind: Service
2  apiVersion: v1
3  metadata:
4    name: dfstore
5  spec:
6    selector:
7      app: dragonfly
8      component: dfdaemon
9      release: dragonfly
10
11   ports:
12     - protocol: TCP
13       port: 65001
14       targetPort: 65001
15
16   type: NodePort
```

Create service:

```
Shell |
1  kubectl --namespace dragonfly-system apply -f dfstore.yaml
```

Forward request to Dragonfly Peer's HTTP proxy:

```
Shell |
1  kubectl --namespace dragonfly-system port-forward service/dfstore 65001:65001
```

Install Dragonfly Repository Agent

Set Dragonfly Repository Agent configuration

Create the `dragonfly_config.json` configuration file, the configuration is as follows:

```

dragonfly_config.json Shell |
1 {
2   "proxy": "http://127.0.0.1:65001",
3   "header": {
4     },
5   "filter": [
6     "X-Amz-Algorithm",
7     "X-Amz-Credential&X-Amz-Date",
8     "X-Amz-Expires",
9     "X-Amz-SignedHeaders",
10    "X-Amz-Signature"
11  ]
12 }

```

- proxy: The address of Dragonfly Peer's HTTP Proxy.
- header: Adds a request header to the request.
- filter: Used to generate unique tasks and filter unnecessary query parameters in the URL.

In the filter of the configuration, set different values when using different object storage:

type	value
OSS	["Expires","Signature","ns"]
S3	["X-Amz-Algorithm", "X-Amz-Credential", "X-Amz-Date", "X-Amz-Expires", "X-Amz-SignedHeaders", "X-Amz-Signature"]
OBS	["X-Amz-Algorithm", "X-Amz-Credential", "X-Amz-Date", "X-Obs-Date", "X-Amz-Expires", "X-Amz-SignedHeaders", "X-Amz-Signature"]

Set Model Repository configuration

Create `cloud_credential.json` cloud storage credential, the configuration is as follows:

```
cloud_credential.json Shell |
1 {
2   "gs": {
3     "": "PATH_TO_GOOGLE_APPLICATION_CREDENTIALS",
4     "gs://gcs-bucket-002": "PATH_TO_GOOGLE_APPLICATION_CREDENTIALS_2"
5   },
6   "s3": {
7     "": {
8       "secret_key": "AWS_SECRET_ACCESS_KEY",
9       "key_id": "AWS_ACCESS_KEY_ID",
10      "region": "AWS_DEFAULT_REGION",
11      "session_token": "",
12      "profile": ""
13    },
14    "s3://s3-bucket-002": {
15      "secret_key": "AWS_SECRET_ACCESS_KEY_2",
16      "key_id": "AWS_ACCESS_KEY_ID_2",
17      "region": "AWS_DEFAULT_REGION_2",
18      "session_token": "AWS_SESSION_TOKEN_2",
19      "profile": "AWS_PROFILE_2"
20    }
21  },
22  "as": {
23    "": {
24      "account_str": "AZURE_STORAGE_ACCOUNT",
25      "account_key": "AZURE_STORAGE_KEY"
26    },
27    "as://Account-002/Container": {
28      "account_str": "",
29      "account_key": ""
30    }
31  }
32 }
```

In order to pull the model through Dragonfly, the model configuration file needs to be added following code in `config.pbtxt` file:

```
1  model_repository_agents
2  {
3  ▾  agents [
4  ▾    {
5      name: "dragonfly",
6    }
7  ]
8  }
```

The [densenet_onnx example](#) contains modified configuration and model file. Modified `config.pbtxt` such as:

```
1  name: "densenet_onnx"
2  platform: "onnxruntime_onnx"
3  max_batch_size : 0
4  ▾ input [
5  ▾  {
6      name: "data_0"
7      data_type: TYPE_FP32
8      format: FORMAT_NCHW
9  ▾  dims: [ 3, 224, 224 ]
10 ▾  reshape { shape: [ 1, 3, 224, 224 ] }
11  }
12 ]
13 ▾ output [
14 ▾  {
15      name: "fc6_1"
16      data_type: TYPE_FP32
17 ▾  dims: [ 1000 ]
18 ▾  reshape { shape: [ 1, 1000, 1, 1 ] }
19      label_filename: "densenet_labels.txt"
20  }
21 ]
22 model_repository_agents
23 {
24 ▾  agents [
25 ▾    {
26      name: "dragonfly",
27    }
28  ]
29 }
```

Triton Server integrates Dragonfly Repository Agent plugin

Install Triton Server with Docker

Pull `dragonflyoss/dragonfly-repository-agent` image which is integrated Dragonfly Repository Agent plugin in Triton Server, refer to [Dockerfile](#).

```
▼ Shell |
1  docker pull dragonflyoss/dragonfly-repository-agent:latest
```

Run the container and mount the configuration directory:

```
▼ Shell |
1  docker run --network host --rm \
2  -v ${path-to-config-dir}:/home/triton/ \
3  dragonflyoss/dragonfly-repository-agent:latest tritonserver \
4  --model-repository=${model-repository-path}
```

- `path-to-config-dir` : The files path of `dragonfly_config.json` & `cloud_credential.json`.
- `model-repository-path` : The path of remote model repository.

The correct output is as follows:

```

1  =====
2  == Triton Inference Server ==
3  =====
4  successfully loaded 'densenet_onnx'
5  I1130 09:43:22.595672 1 server.cc:604]
6  +-----+
7  | Repository Agent | Path
8  |                 |
9  | dragonfly       | /opt/tritonserver/repoagents/dragonfly/libtritonrepa
10 | gent_dragonfly.so |
11 +-----+
12 I1130 09:43:22.596011 1 server.cc:631]
13 +-----+
14 | Backend      | Path
15 |             | Config
16 |             |
17 | pytorch      | /opt/tritonserver/backends/pytorch/libtriton_pytorch.so
18 |             | {}
19 |             |
20 | onnxruntime  | /opt/tritonserver/backends/onnxruntime/libtriton_onnxrunti
21 | me.so | {"cmdline":{"auto-complete-config":"true","backend-directory":"/op
22 | t/tritonserver/backends","min-compute-capability":"6.000000","default-max-
    batch-size":"4"}} |

```

```

23 +-----+-----+-----+
24 | densenet_onnx | 1 | READY |
25 +-----+-----+-----+
26
27 I1130 09:43:22.598318 1 metrics.cc:703] Collecting CPU metrics
28 I1130 09:43:22.599373 1 tritonserver.cc:2435]
29 +-----+-----+-----+
30 | Option | Value
31 +-----+-----+-----+
32 | server_id | triton
33 | server_version | 2.37.0
34 | server_extensions | classification sequence model_reposit
    ory model_repository(unload_dependents) schedule_policy model_configuratio
    n system_shared_memory cuda_shared_memory binary_tensor_data parameters st
35 | model_repository_path[0] | s3://192.168.36.128:9000/models
36 | model_control_mode | MODE_NONE
37 | strict_model_config | 0
38 | rate_limit | OFF
39 | pinned_memory_pool_byte_size | 268435456

```

```

40 | min_supported_compute_capability | 6.0
41 | strict_readiness | 1
42 | exit_timeout | 30
43 | cache_enabled | 0
44 +-----+
45 +-----+
46 I1130 09:43:22.610334 1 grpc_server.cc:2451] Started GRPCInferenceService
47 at 0.0.0.0:8001
48 I1130 09:43:22.612623 1 http_server.cc:3558] Started HTTPService at 0.0.0.
49 0:8000
50 I1130 09:43:22.695843 1 http_server.cc:187] Started Metrics Service at 0.
0.0.0:8002

```

Execute the following command to check the Dragonfly logs:

```

Shell |
1 kubectl exec -it -n dragonfly-system dragonfly-dfdaemon-<id> -- tail -f /va
r/log/dragonfly/daemon/core.log

```

Check downloaded successfully through Dragonfly:


```

1 {
2   "level":"info","ts":"2024-02-02 05:28:02.631",
3   "caller":"peer/peertask_conductor.go:1349",
4   "msg":"peer task done, cost: 352ms",
5   "peer":"10.244.2.3-1-4398a429-d780-423a-a630-57d765f1ccfc",
6   "task":"974aaf56d4877cc65888a4736340fb1d8fecc93eadf7507f531f9fae650f1b4d",
7   "component":"PeerTask",
8   "trace":"4cca9ce80dbf5a445d321cec593aee65"
9 }

```

Verify

Call inference API:

```

1 docker run -it --rm --net=host nvcr.io/nvidia/tritonserver:23.08-py3-sdk /w
  orkpace/install/bin/image_client -m densenet_onnx -c 3 -s INCEPTION /works
  pace/images/mug.jpg

```

Check the response successful:

```

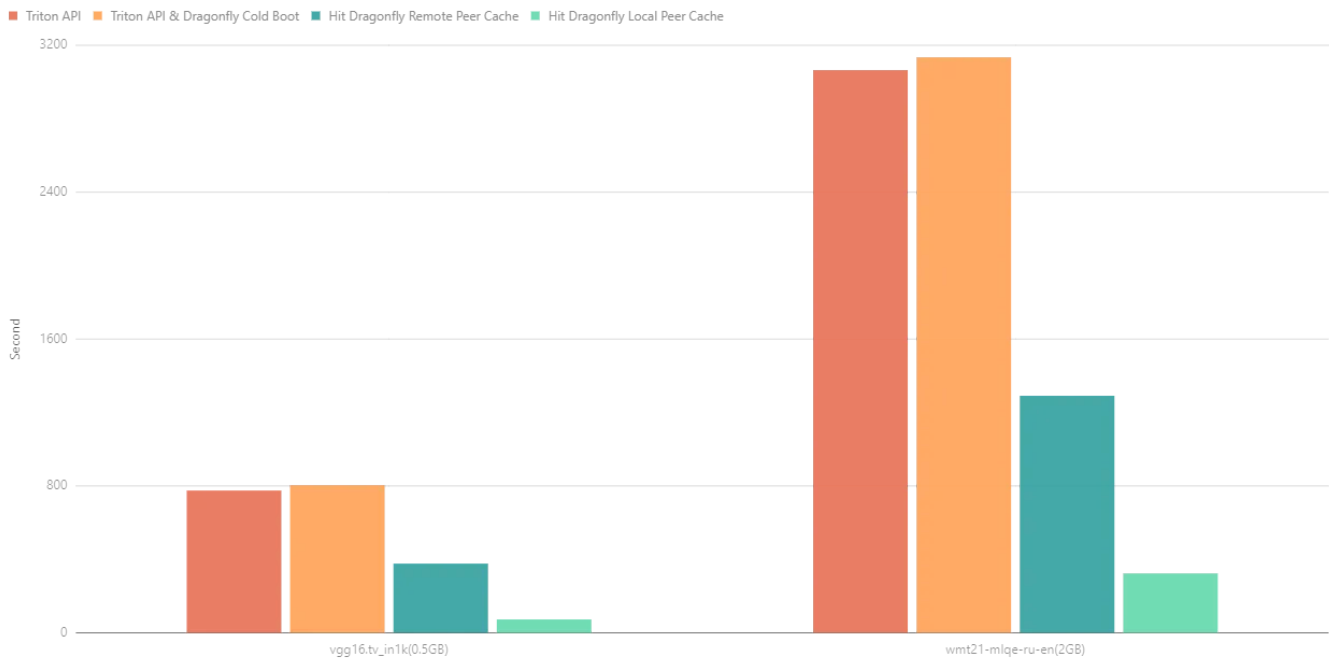
1 Request 01
2 Image '/workspace/images/mug.jpg':
3   15.349563 (504) = COFFEE MUG
4   13.227461 (968) = CUP
5   10.424893 (505) = COFFEEPOT

```

Performance testing

Test the performance of single-machine model download by Triton API after the integration of Dragonfly P2P. Due to the influence of the network environment of the machine itself, the actual download time is not important, but The proportion of download speed in different scenarios is more meaningful:

Time to download large



- Triton API: Use signed URL provided by Object Storage to download the model directly.
- Triton API & Dragonfly Cold Boot: Use `Triton Serve API` to download model via Dragonfly P2P network and no cache hits.
- Hit Remote Peer: Use `Triton Serve API` to download model via Dragonfly P2P network and hit the remote peer cache.
- Hit Local Peer: Use `Triton Serve API` to download model via Dragonfly P2P network and hit the local peer cache.

Test results show Triton and Dragonfly integration. It can effectively reduce the file download time. Note that this test was a single-machine test, which means that in the case of cache hits, the performance limitation is on the disk. If Dragonfly is deployed on multiple machines for P2P download, the models download speed will be faster.

相关链接

Dragonfly 社区

- Website: <https://d7y.io/>
- Github Repo: <https://github.com/dragonflyoss/Dragonfly2>
- Dragonfly Repository Agent Github Repo: <https://github.com/dragonflyoss/dragonfly-agent>

[repository-agent](#)

- Slack Channel: [#dragonfly](#) on [CNCF Slack](#)
- Discussion Group: dragonfly-discuss@googlegroups.com
- Twitter: [@dragonfly_oss](#)

NVIDIA Triton Inference Server

- Website: <https://developer.nvidia.com/triton-inference-server>
- Github Repo: <https://github.com/triton-inference-server/server>

二维码

Dragonfly Github 仓库:

