# 基于 FlagRelease 实现模型到不同芯片上的推理部署

# 支持FlagRelease实现模型多芯片自动发版

ModelScope/Huggingface/...

支持单一芯片的模型　　　　　　　　　　　　　　支持多款芯片的模型、docker镜像等

## ① 拉取
- 下载所选模型（含对应配置）
- 选择另一款芯片作为目标芯片

## ② 依赖分析
- 下载基础 docker 镜像
- 安装依赖包
- 安装 FlagOS 栈

## ③ 模型迁移
- 在目标硬件上实现 CICD
- 对模型进行量化（可选）
- 通过自动调优找到最佳策略
- 自动为模型提供服务

## ④ 模型评估
- 自动评估模型性能
- 观察系统监控
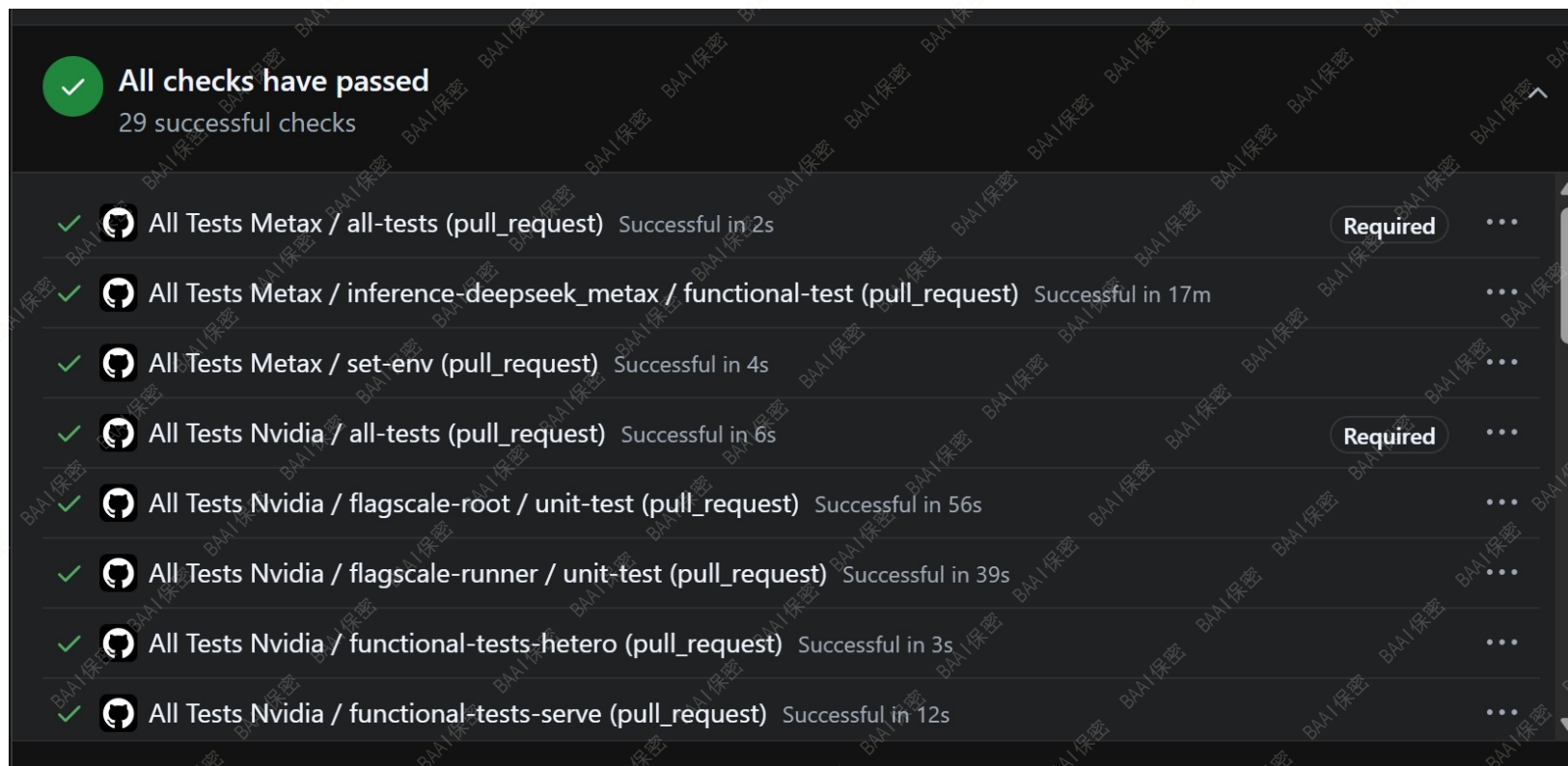- 记录评估结果和性能数据

## ⑤ 发布
- 发布生成的文档
- 发布评估结果
- 发布代码和 docker 镜像
- 发布量化后的模型（可选）

由 FlagScale 提供支持

芯片A集群

芯片B集群

芯片C集群

# 自动化CI/CD持续升级完善，更好支持FlagRelease

- 提交PR时自动触发

- 本地可手动出发

- 集成高性能算子库FlagGems测试

- 未来将继承高性能统一通信库FlagCX测试

- 为FlagRelease新模型的发布提供支撑

- 当前已上线沐曦CICD，更多国产芯片CI/CD未来持续上线中

- 下载镜像

docker pull flagrelease-registry.cn-beijing.cr.aliyuncs.com/flagrelease/flagrelease:nv_vllm085_gemscbbf39

- 启动容器

docker run --rm --init --detach \ --net=host --uts=host --ipc=host \ --security-opt=seccomp=unconfined \ --privileged=true \ --ulimit stack=67108864 \ --ulimit memlock=-1 \ --ulimit nofile=1048576:1048576 \ --shm-size=32G \ -v /nfs:/nfs \ --gpus all \ --name flagos \ flagrelease-registry.cn-beijing.cr.aliyuncs.com/flagrelease/flagrelease:nv_vllm085_gemscbbf39 \ sleep infinity docker exec -it flagos bash

- 安装modelscope并下载模型

pip install modelscope
modelscope download --model Qwen/Qwen3-4B --local_dir /nfs/Qwen3-4B

- 启动推理服务

flagscale serve qwen3

*https://www.modelscope.cn/models/FlagRelease/Qwen3-4B-FlagOS-Nvidia*

- 下载镜像

docker pull flagrelease-registry.cn-beijing.cr.aliyuncs.com/flagrelease/flagrelease:flagrelease_metax_qwen3

- 启动容器

docker run -it --device=/dev/dri --device=/dev/mxcd --group-add video --name test_flagrelease --device=/dev/mem --network=host --security-opt seccomp=unconfined --security-opt apparmor=unconfined --shm-size=100gb --ulimit memlock=-1 -v /usr/local/:/usr/local/ -v /nfs:/nfs flagrelease-registry.cn-beijing.cr.aliyuncs.com/flagrelease/flagrelease:flagrelease_metax_qwen3 /bin/bash

- 安装modelscope并下载模型

pip install modelscope
modelscope download --model Qwen/Qwen3-4B --local_dir /nfs/Qwen3-4B

- 启动推理服务

flagscale serve qwen3

https://www.modelscope.cn/models/FlagRelease/Qwen3-4B-FlagOS-Metax

谢谢聆听，
欢迎合作与建议。