CVPR2020 Video Modeling Tutorial

Deployment

Zhi Zhang, AWS, June 2020

©2020, Amazon Web Services, Inc. or its Affiliates. All rights reserved. Amazon Trademark







Outline

- Why we take deployment seriously
- Why deployment of deep video models is hard
- How we apply the best practices (notebook walkthrough)
- Things to expect in the near future
- Take-aways







Why deployment

In order to start using a deep learning model for practical workload in commercial products or production service, it needs to be efficiently deployed into production.

For example, you may train the model on cloud services but deploy the model in a self-driving vehicle with completely different hardware/software stack.













Why deep learning video models are so hard to deploy

- Diverging programming languages in development and deployment
- The gaps of compute power
- Portability and dependency(video)







Our approaches











Jupyter notebook walkthrough









Compact example for Jetson on device walkthrough

ſ	79%]	Building CXX object decord/CMakeFiles/decord.dir/src
ſ	82%]	Building CXX object decord/CMakeFiles/decord.dir/src
ſ	86%]	Linking CXX shared library libdecord.so
ſ	86%]	Built target decord
Scanning dependencies of target video_classification		
ſ	89%]	Building CXX object CMakeFiles/video_classification.
ſ	93%]	Building CXX object CMakeFiles/video_classification.
ſ	96%]	Building CXX object CMakeFiles/video_classification.
[100%] Linking CXX executable video_classification		
[100%] Built target video_classification		
xavier@xavier-a3:~/cvpr20-tutorial/cvpr2020-videomodeling-d		



video/ffmpeg/filter_graph.cc.o video/ffmpeg/threaded_decoder.cc.o

dir/tvm_runtime_pack.cc.o ir/src/classification.cpp.o lir/src/video.cpp.o

leployment/tvm_deploy/build\$





More to expect...

- TVM deployment with 3d Conv networks
- INT8 quantization for GluonCV video models
- GluonCV lite video model zoo: a curated list of models specially optimized for individual devices

Take-aways

- Thinking twice about deployment before you develop the models. Redesign model or rewrite the code for deployment can be costly. We use TVM to generated device specific library which is light weight and efficient.
- We now have `decord` library to deal with video inputs as straightforward as images.
- In case there's CPU/GPU available at the same time, make sure we figure out which one to target. For example, use GPUs instead of CPUs on Jetson devices
- The materials for this tutorial are all open sourced and available at: <u>https://</u> github.com/zhreshold/cvpr2020-videomodeling-deployment







©2020, Amazon Web Services, Inc. or its Affiliates. All rights reserved. Amazon Trademark



Thank you!



