

End-to-End Deep Learning Compiler Stack

AWS AI

Presenter: Animesh Jain Amazon SageMaker Neo

Deep Learning is Pervasive



ARM – Unique Role in AWS Ecosystem





How to Accelerate Deep Learning?



Agenda

- Overview of Neo
- Relay Graph Optimizations
- TVM Tensor IR Optimizations
- Evaluation

Deep Learning Inference



amazon echo

aws

Models and hardware targets are far away!





TVM: end-to-end optimization stack







^{© 2019,} Amazon Web Services, Inc. or its Affiliates. All rights reserved.

Agenda

- Overview of TVM
- Relay Graph Optimizations
- TVM Tensor IR Optimizations
- Evaluation

Computation Graph Optimization

Represent high-level deep learning computations



Target Independent

- Constant propagation
- Dead code elimination
- Operation fusion

Target Dependent

- Data layout transform
- Graph partitioning
- Legalization



Operation Fusion Example





Target Dependent Layout Transformation





Agenda

- Overview of TVM
- Relay Graph Optimizations
- TVM Tensor IR Optimizations
- Evaluation

TVM Tensor IR

Compute definition

```
C = tvm.compute((m, n),
    lambda i, j: tvm.sum(A[i, k] * B[k, j], axis=k))
```

- TVM Schedule Developer-friendly loop transformations
 - Do not need hardware ISA knowledge to perform loop optimizations

```
s = tvm.create_schedule(C.op)
> xo, yo, xi, yi = s[C].tile(C.op.axis[0], C.op.axis[1], bn, bn)
> ko, ki = s[C].split(k, factor=4)
> s[C].reorder(ko, xi, ki, yi)
> s[C].unroll(ki)
> s[C].vectorize(yi)
```

> s[C].parallel(xo)

. . .

Large Search Space for Schedule





- Relatively low experiment cost
- Domain-specific problem structure
- Large quantity of similar tasks

NeurIPS'19 Learning to Optimize Tensor Programs

Agenda

- Overview of TVM
- Relay Graph Optimizations
- TVM Tensor IR Optimizations
- Evaluation

Experiment Setup

- Server EC2 A1 Instance
 - 16-core ARMv8

- Edge device Acer aiSage
 - Rockchip RK3399 SoC + ARM Mali GPU T-860



TVM – Evaluation on ARM A1 Server

Speedup of TVM execution normalized to Tensorflow-Eigen



TVM – Evaluation on Edge Device Acer aiSage



ICPP'19 A Unified Optimization Approach for CNN Model Inference on Integrated GPUs

aws

Effects of Tuning Convolution operators in TVM



ICPP'19 A Unified Optimization Approach for CNN Model Inference on Integrated GPUs



Takeaways

- Deep learning compilation is essential for portability and performance across a variety of targets.
- Optimizations are important at all levels graph- and tensor-level.
- Abstracting compute and HW-dependent schedule enables developers to write kernels without extensive knowledge of HW ISA.
- Open-source collaborations are the key to achieve the dream of running deep learning everywhere.



Linaro Community

- TVM enjoys LLVM ARM codegen support
 - Better support for Int8 instructions
 - Better support for different ARM variants
- Better schedules
 - Data layout optimizations are hardware dependent
 - TVM performance of ARM GPUs can be improved
- ACL support for TVM

Check it out!

Thank you!

https://github.com/neo-ai/

https://github.com/dmlc/tvm



Amazon SageMaker Neo



