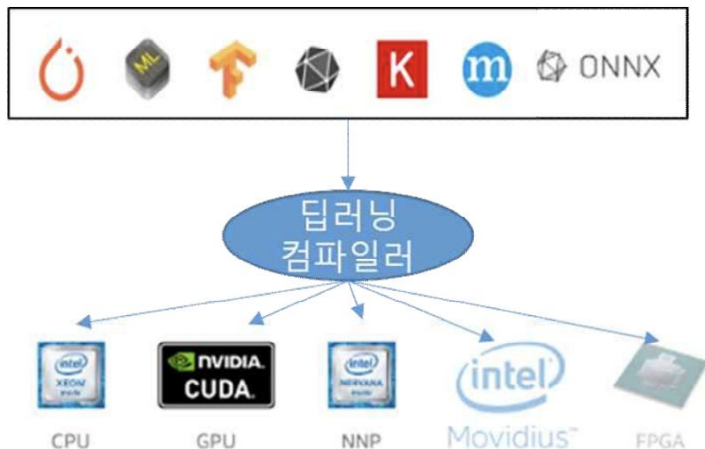
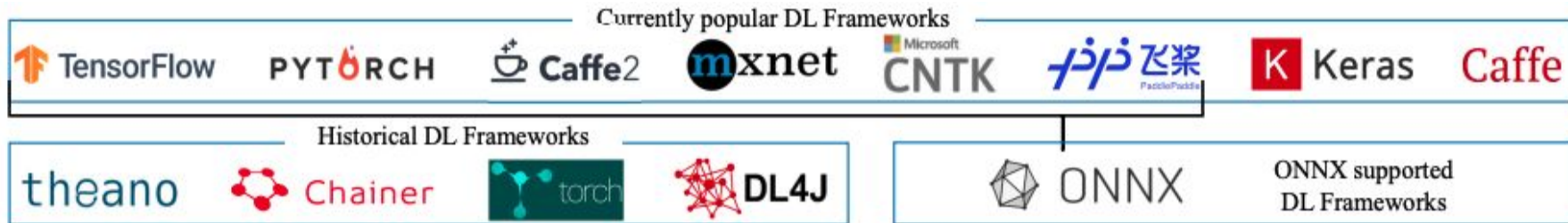


The Deep Learning Compiler

Tae Young Lee

DL Framework landscape



- 딥러닝 컴파일러는, 다양한 딥러닝 플랫폼에서 학습된 딥러닝 모델을 입력으로 받아,
- 특정 하드웨어에서 동작 가능한 머신 코드(또는, 백엔드 코드)를 자동으로 생성합니다.
- 최근 제안된 XLA, TVM, Glow 와 같은 딥러닝 컴파일러들은,
- TensorFlow, Pytorch, MxNet, ONNX 등의 프레임워크로 작성된 모델을 입력으로 하여,
- CPU 및 GPU 용 백엔드 코드를 생성하고 있습니다.
 - CPU 용 백엔드 코드 : llvm
 - GPU 용 백엔드 코드 : CUDA, OpenCL, Metal 등



- **계산 그래프 생성기**

- 계산 그래프 생성기는, 다양한 딥러닝 플랫폼에서 생성된 딥러닝 모델을 로딩하여, 그래프 구조로 재구성한 뒤,
- 모델에서 명시한 연산자들을, 프리미티브 연산자로 변환하여, 계산 그래프를 생성합니다.
- 이후, 불필요한 노드 제거, 연산자 결합(Operator Fusion), 양자화(Quantization), 상수 폴딩(constant-folding)과 같은, 그래프 최적화를 기법들을 적용하여,
- 최종적으로 최적화 된 계산 그래프를 생성합니다.

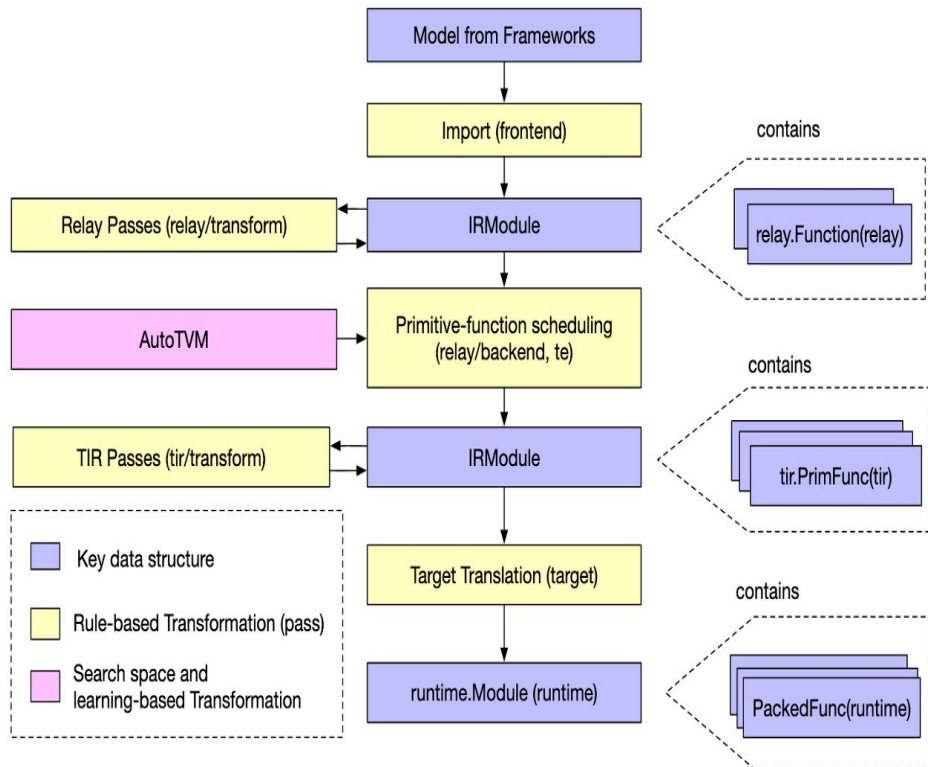
- **IR 생성기**

- 계산 그래프 생성기로부터 생성된 그래프는, 하드웨어에 독립적인 정보만 표현 할 수 있습니다.
 - ex) 연산자 간 입출력 관계, 연산자의 입출력 데이터 값, etc
- 그러나 컴파일러는, 딥러닝 모델의 실행 성능을 높이기 위해,
- 모델이 동작할 하드웨어의 아키텍처를 고려한, ‘하드웨어 의존적 최적화’를 필수적으로 수행해야 합니다.
- 따라서, 계산 그래프를 ‘하드웨어 의존적 최적화’에 적용 가능한 형태인,
- 중간 표현 (**IR: Intermediate Representation**)으로 변환하는 과정이 필요합니다.
 - 참고로, IR 은 입출력 변수(또는 Tensor)들에 대한 선언과,
 - 변수들을 사용하는 연산식에 대한 정보를 포함하고 있습니다.
- 컴파일러는, 생성된 **IR** 을 바탕으로 메모리 할당, 병렬화, 연산식 실행 순서 결정 등의 ‘하드웨어 의존적 최적화’ 작업을 수행한 뒤, 코드를 생성합니다.
- 요약하면, **IR** 생성기는 컴파일러가 하드웨어 의존적 최적화를 용이하게 수행하도록, 계산 그래프를 **IR** 로 변환하는 역할을 합니다.

- **백엔드 코드 생성기**

- 백엔드 코드 생성기는, **IR** 을 기반으로, 딥러닝 모델 워크로드가 배치될 타겟 하드웨어 (**CPU, GPU, TPU** 등)의 아키텍처에 최적화 된 백엔드 코드를 생성합니다.
- 즉, 타겟 하드웨어 별로, 병렬로 처리 가능한 데이터의 크기, 캐시 크기, 효율적으로 수행 가능한 딥러닝 오퍼레이터의 타입, 실행 순서 등을 참조하여,
- 스케줄 정보를 생성한 뒤, 해당 하드웨어에 최적화 된 백엔드 코드를 생성합니다.

TVM 실행 흐름



- **DL framework**로부터 생성한 모델을 **import** 한 뒤, **relay** 를 통해, 해당 모델을 **TVM** 의 **IRModule** 로 변환합니다.
 - IRModule에는 모델을 나타낼 수 있는 다양한 함수들이 포함되어 있습니다.
- **IRModule** 을 점점 **Lowering** 하는 **Transformation** 과정을 거칩니다.
 - IRModule 이 원하는 하드웨어 (메모리 구조, 데이터 타입 등)에 알맞게 최적화될 수 있도록 **Tensorize** 시킵니다.
 - 이를 위해 **Quantization**, **Graph-Packing** 등, 다양한 작업들이 포함되며,
 - 이때, **AutoTVM** 이 사용됩니다.
- 타깃 하드웨어에서 잘 작동하는 코드를 생성하기 위해, **Translation** 과정을 거칩니다.
 - 예를 들어, 타깃 하드웨어가 **GPU** 라면 **[CUDA/OpenCL/Metal 등]**, CPU라면 **[llvm]** 코드를 생성합니다.
 - 이는, **TVM** 의 **runtime.Module** 에 encapsulate 되어있고,
 - 이를 통해, 타깃 하드웨어에 쉽게 **export / load / execute** 시킬 수 있습니다.
- **runtime.Module** 를 로드하여, 타깃 하드웨어에서 최적화 된 모델을 동작시킵니다.

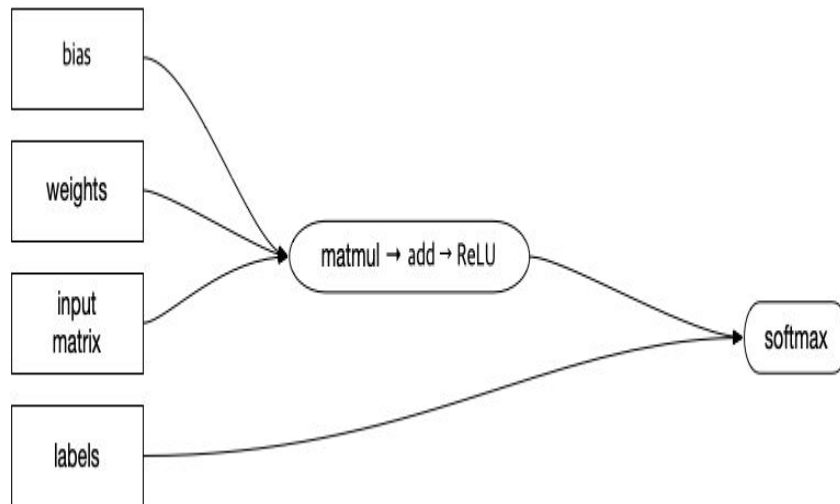
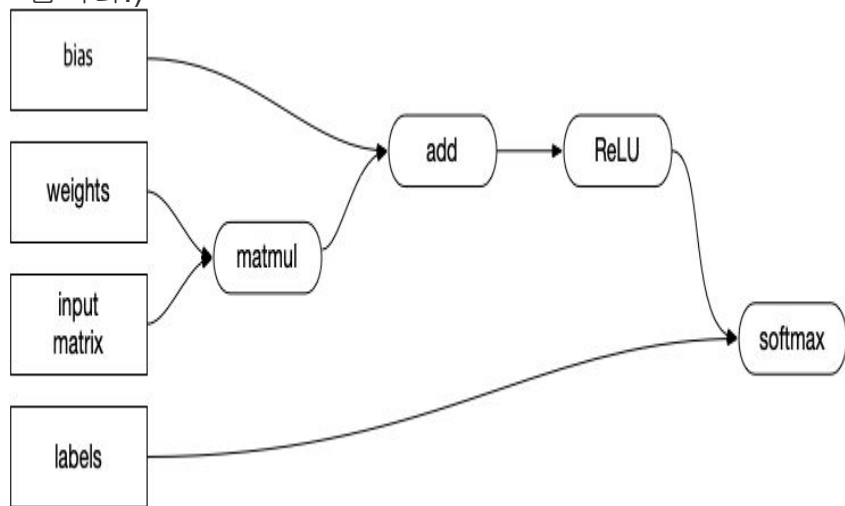
TFLite, TensorRT, TVM의 차이

- **TFLite**
 - TensorFlow(Keras)의 모델만 최적화가 가능합니다.
 - 여러 하드웨어에서 추론이 가능합니다.
 - 타겟 하드웨어에 따른 추가적인 최적화는 없습니다.
- **TensorRT**
 - 다양한 프레임워크의 모델 최적화가 가능합니다.
 - **NVIDIA GPU** 만 추론 하드웨어로 사용 가능합니다.
 - 추론에 사용하려는 **NVIDIA GPU** 종류에 따라 추가적으로 최적화할 수 있습니다.
- **TVM**
 - 다양한 프레임워크의 모델 최적화가 가능합니다.
 - 여러 하드웨어에서 추론이 가능합니다.
 - 추론에 사용하려는 하드웨어 종류에 따라 추가적으로 최적화할 수 있습니다.

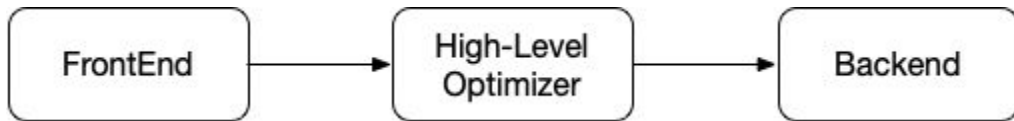
XLA 란?

우선 XLA (Accelerated Linear Algebra)는 Tensorflow의 서브 프로젝트로 그래프 연산의 최적화 / 바이너리 사이즈의 최소화 등을 목적으로 하는 컴파일러입니다.

결과를 먼저 말씀드리자면, XLA를 활용하는 것을 통해 그래프 연산 과정에 필요한 임시 버퍼도 덜 사용할 수 있고, 그래프 연산 속도도 단축시킬 수 있습니다. (구글 측 자료에 의하면 Nvidia GPU를 사용할 때 50%까지도 성능 향상이 있었다고 합니다.)

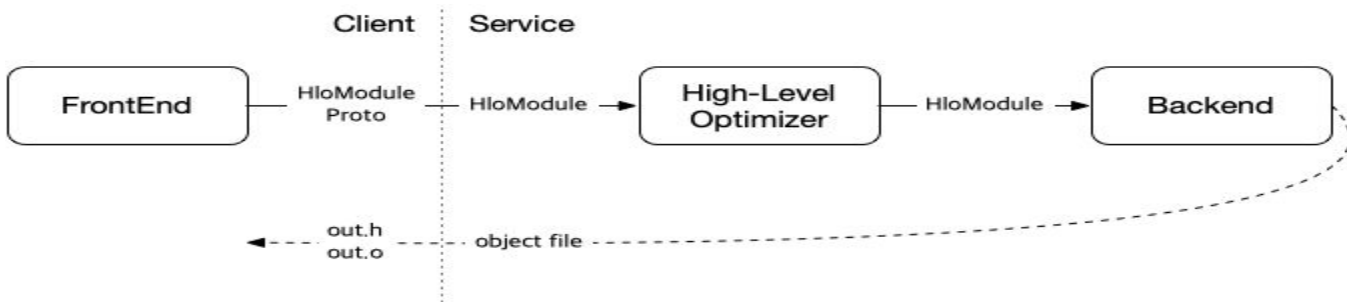


XLA



XLA도 일종의 컴파일러이기 때문에 일반적인 컴파일러 들이 동작하는 것처럼 **FrontEnd (Parser) / MiddleEnd (High-level Optimizer) / Backend (Target Code Generator)**의 구조를 가지고 있으며, 주어진 타겟 플랫폼에 맞춰 그래프를 최적화 한 후 실행가능한 형태의 무언가를 생성하는 역할을 수행

Backend에서는 LLVM을 사용하고 있으며, LLVM과 마찬가지로 **JiT compilation**과 **AoT compilation** 모두를 지원하는게 특징입니다. 두가지 중 어떤 방식을 사용하느냐에 따라 실제 사용 시나리오가 달라지게 되며, **AoT compilation**을 사용하는 경우 아래와 같이 컴파일의 결과로 **object file** 및 이를 사용하기 위한 **helper class**에 대한 **header file**이 생성됩니다.

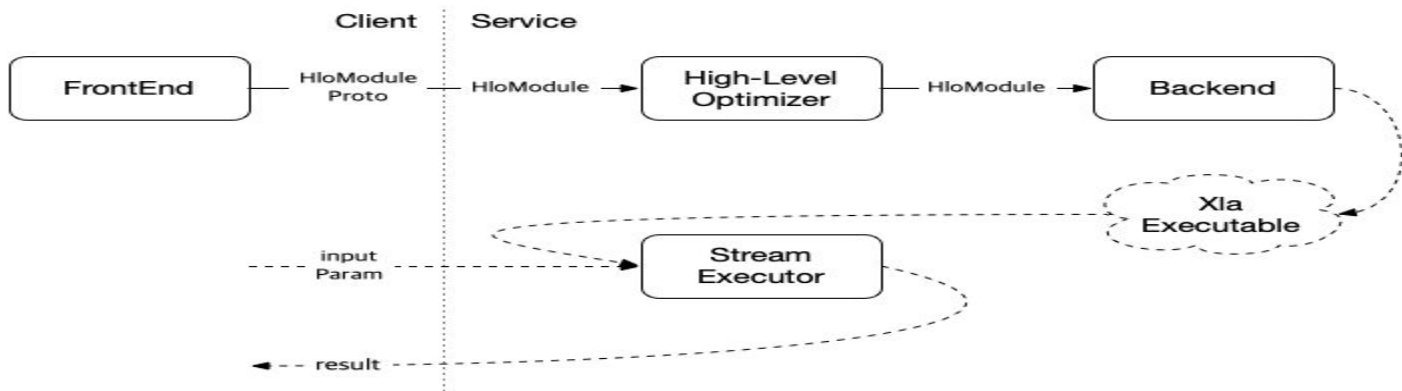


AOT

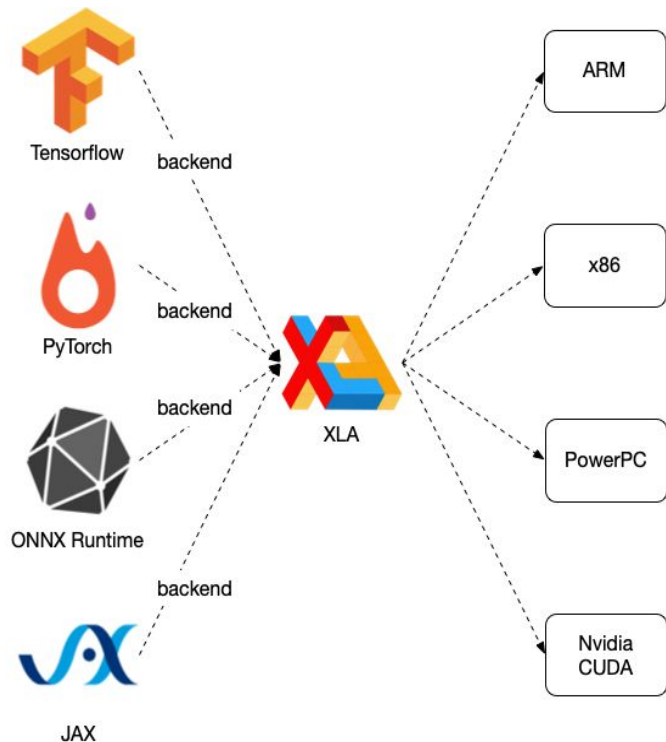
`tfcompile`은 Ahead-of-Time(AOT)이 TensorFlow 그래프를 실행 가능한 코드로 컴파일하는 독립 실행형 도구입니다. 이를 통해 총 바이너리 크기를 줄이고 일부 런타임 오버헤드를 방지할 수 있습니다. `tfcompile`은 일반적으로, 추론 그래프를 모바일 기기용 실행 코드로 컴파일하는 데 사용됩니다.

실행 바이너리를 만들기 위해 XLA 라이브러리나 Tensorflow 라이브러리 등을 포함할 필요가 없기 때문에 최종 타겟 바이너리 사이즈를 굉장히 작게 만들 수 있으나 기본 CPU backend의 경우 convolution, matmul 등의 연산에 대해 XLA runtime function을 호출하는 코드를 생성하도록 되어 있기 때문에 일부 XLA runtime function들과 링크되어야 할 수도 있습니다.

JiT compilation을 사용하는 경우는 `compile`의 결과로 `XlaExecutable` 데이터가 생성되며, `StreamExecutor`를 통해 input parameter를 feeding하고 그래프를 실행시킬 수 있는 형태입니다.



JitCompilation



JitCompilation을 사용하는 경우 런타임에 그래프를 입력으로 받아 실행할 수 있다는 장점이 있으나 애플리케이션이 **XLA** 런타임을 포함해야하기 때문에 **eigen**, **protobuf**, **LLVM** 등등 의존성들이 줄줄이 엮인다는 단점이 있습니다. (타겟 플랫폼에 해당 라이브러리들을 줄줄이 포팅해야할 수도 있습니다.)

애플리케이션을 배포하기 전 단계에 타겟 플랫폼이나 네트워크 형태에 대해 미리 알고 있는 경우에는 **AoT compilation**을 활용하는 방안을 고려해볼 수 있으며, 애플리케이션은 네트워크를 모르는 상태로 배포되어야 하는 환경이라면 **JiT compilation**을 사용해야합니다.

Tensorflow 외에도 다양한 머신러닝 프레임워크의 백엔드로 **XLA**를 활용할 수 있으며, **PyTorch**를 위해선 **PyTorch XLA**, **ONNX (Open Neural Network Exchange)**를 위해선 **onnx xla** 프로젝트가 개발되고 있습니다.

Compiler frontend

Transformation
Symbolic representation

- relay (TVM)
- bridge (nGraph)
- ATen (TC)
- direct translation
- ...

High-level IR / Graph IR
(Device independent)

- Representation**
DAG-based
Let-binding-based
Tensor computation
Lambda / Einstein
- Implementation**
Data representation
Operators supported

Computation graph

Computation graph Optimizations

- Algebraic simplification
- Operator fusion
- Operation sinking
- CSE
- DCE
- Static memory planning
- Layout transformation
- ...

Methods

- Pattern matcher
- Graph rewriting

Debug tools
(IR dumping)

- Text form
- DAG form
- ...

Optimized computation graph

Hardware specific Optimizations

- Intrinsic mapping
- Memory allocation
- Memory latency hiding
- Loop oriented opt
- Parallelization
- ...

Auto Scheduling
(e.g., polyhedral)

- Auto-tuning**
- Parameterization
 - Cost model
 - Parameter searching

Manual Scheduling
(e.g., Halide)

- Using kernel libraries**
- Intel DNNL
 - NV cuDNN / TensorRT
 - AMD MIOpen
 - Other customized libs

Low-level IR / Operator IR
(Device specific)

- Halide based
- Polyhedral model based
- Other unique IRs ...

Compilation Scheme

- Just-In-Time
- Ahead-Of-Time

- Code generation**
- LLVM
 - CUDA
 - OpenCL
 - OpenGL
 - ...

Compiler backend

Target platforms

CPU
(X86, ARM, RISC-V)

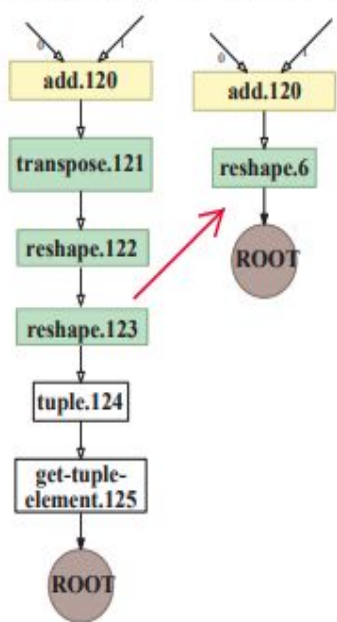
GPU
(NVIDIA, AMD)

ASIC
(TPU, Inferentia, NNP, ...)

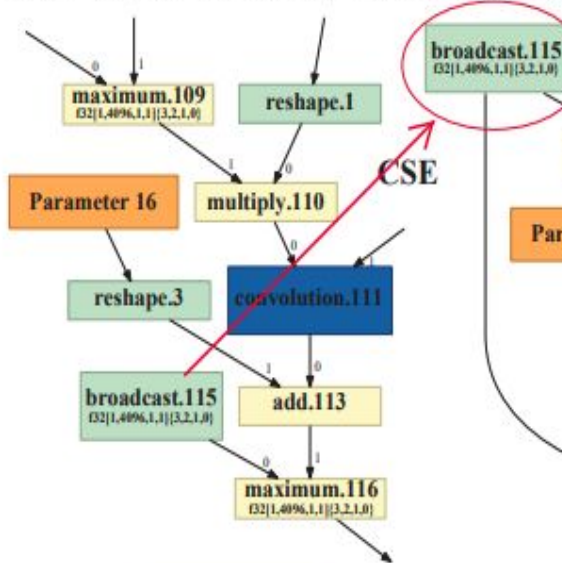
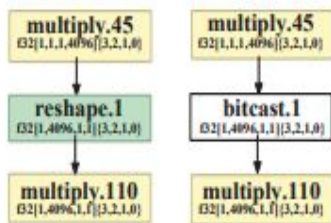
DSP

...

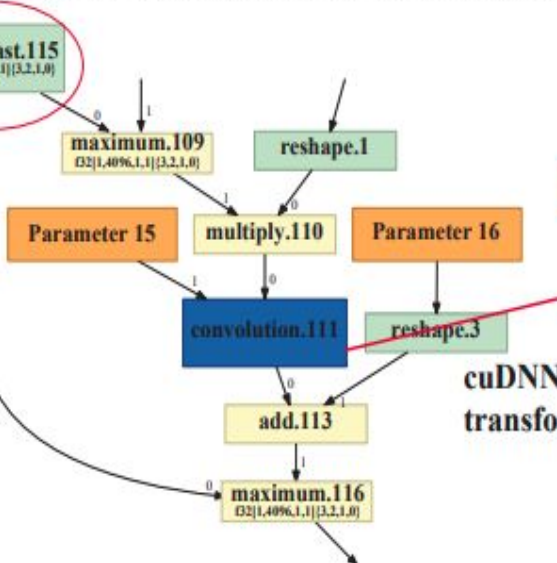
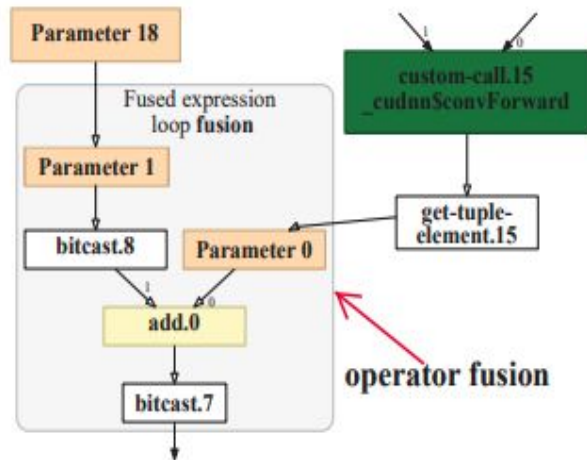
More and more Accelerators



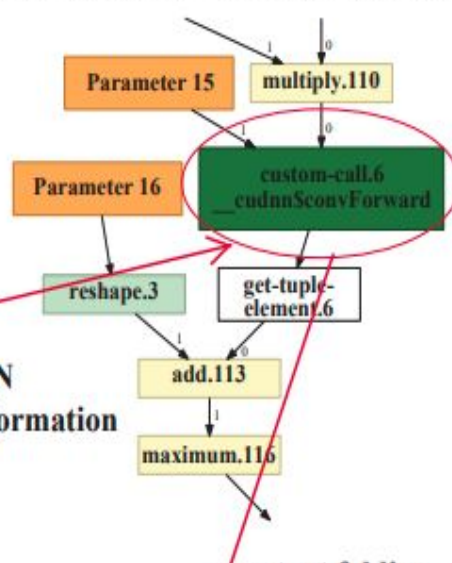
algebraic simplification



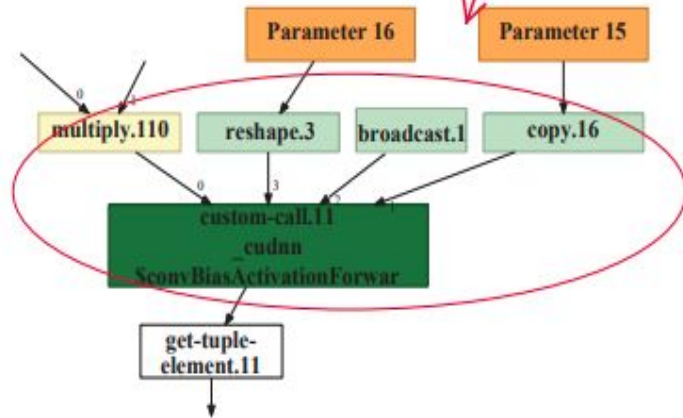
operator fusion



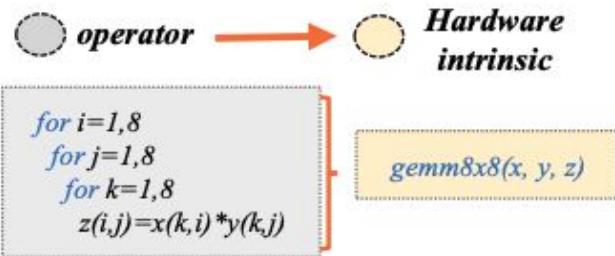
cuDNN transformation



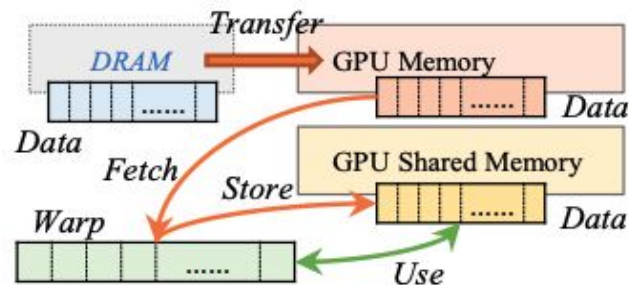
constant folding



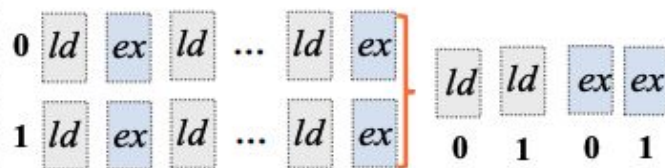
Hardware Intrinsic Mapping



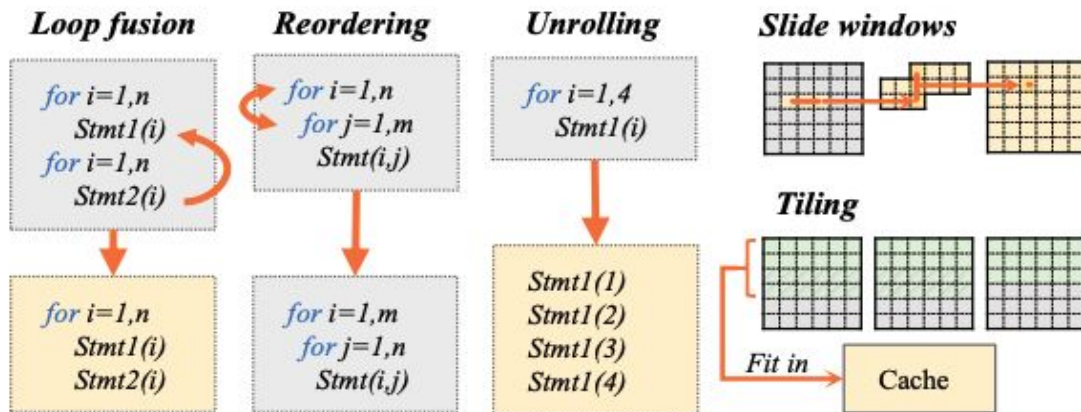
Memory Allocation & Fetching



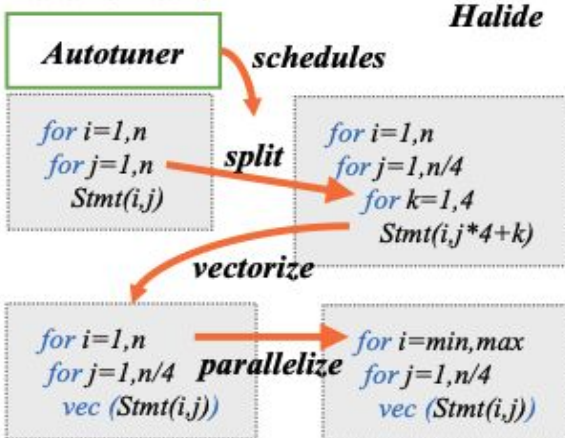
Memory Latency Hiding



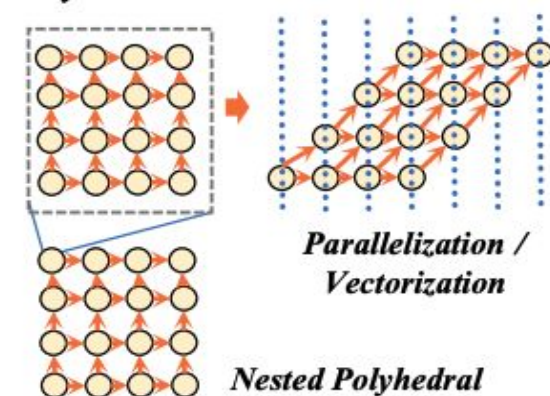
Loop Oriented Optimization Techniques



Parallelization



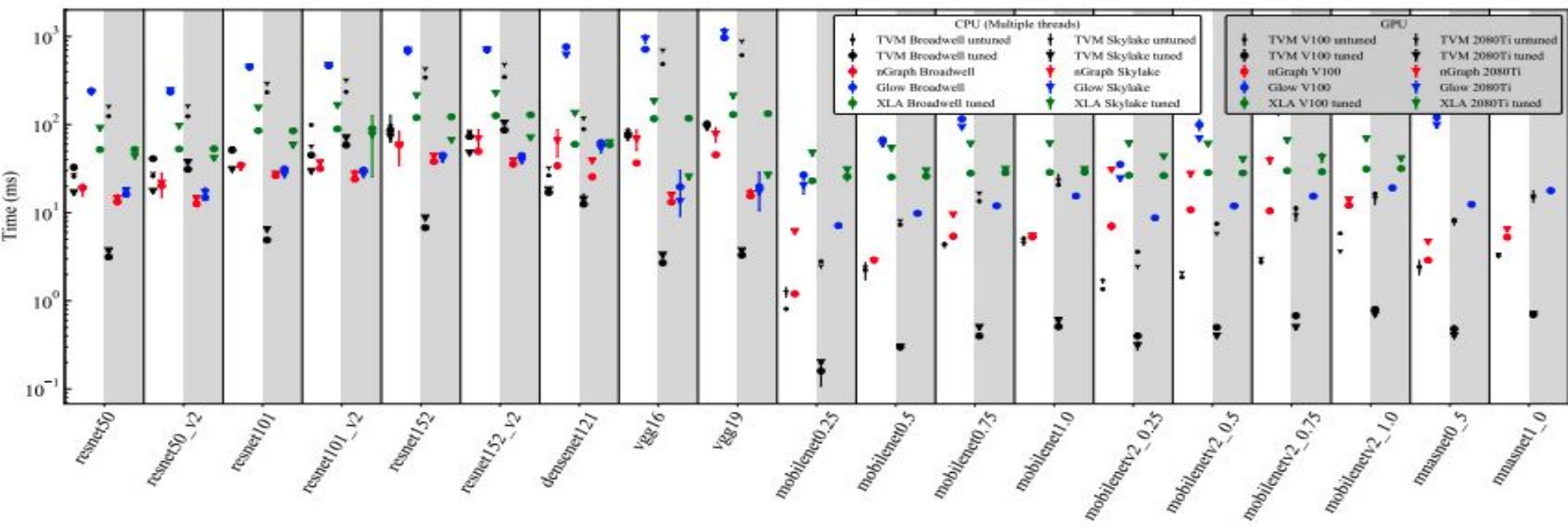
Polyhedral



		TVM	nGraph	TC	Glow	XLA
Frontend	Developer	Apache	Intel	Facebook	Facebook	Google
	Programming	Python/C++ Lambda expression	Python/C++ Tensor expression	Python/C++ Einstein notation	Python/C++ Layer programming	Python/C++ Tensorflow interface
	ONNX support	✓ tvm.relay.frontend .from_onnx (built-in)	✓ Use ngraph-onnx (Python package)	×	✓ ONNXModelLoader (built-in)	✓ Use tensorflow-onnx (Python package)
	Framework support	tvm.relay.frontend .from_* (built-in) tensorflow/tflite/keras pytorch/caffe2 mxnet/coreml/darknet	tensorflow paddlepaddle (Use *-bridge, act as the backend)	(Define and optimize a TC kernel, which is finally called by other frameworks.) pytorch/other DLPack supported frameworks	pytorch/caffe2 tensorflowlite (Use built-in ONNXIFI interface)	Use tensorflow interface
	Training support	×	✓ Only on NNP-T processor	✓ (Support auto differentiation)	✓ (Limited support)	✓ Use tensorflow interface
	Quantization support	✓ int8/fp16	✓ int8 (include training)	×	✓ int8	✓ int8/int16 (Use tensorflow interface)
IR	High-/low-level IR	Relay/Halide	nGraph IR/None	TC IR/Polyhedral	Its own high-/low-level IR	HLO (Both high- and low- level)
	Dynamic shape	✓ (Any)	✓ (PartialShape)	×	×	✓ (None)
Optimization	Frontend opt	Hardware independent optimizations (refer to Section 4.3) Hardware specific optimizations (refer to Section 4.4) And hybrid optimizations				
	Backend opt					
	Autotuning	✓ (To select the best schedule parameters)	×	✓ (To reduce JIT overhead)	×	✓ (On default convolution and gemm)
	Kernel libraries	✓ mkl/cudnn/cublas	✓ eigen/mkldnn/cudnn/ Others	×	×	✓ eigen/mkl/ cudnn/tensorrt
Backend	Compilation methods	JIT AOT (experimental)	JIT	JIT	JIT AOT (Use built-in executable bundles)	JIT AOT (Generate executable libraries)
	Supported devices	CPU/GPU/ARM FPGA/Customized (Use VTA)	CPU/Intel GPU/NNP GPU/Customized (Use OpenCL support in PlaidML)	Nvidia GPU	CPU/GPU Customized (Official docs)	CPU/GPU/TPU Customized (Official docs)

Table 2. The hardware configuration.

	CPU	GPU
Platform a	Broadwell E5-2680v4 *2 (28 physical cores, 2.4GHz)	Tesla V100 32GB (15.7TFlops, FP32)
Platform b	Skylake Silver 4110 *2 (16 physical cores, 2.1GHz)	Turing RTX2080Ti 11GB (13.4TFlops, FP32)



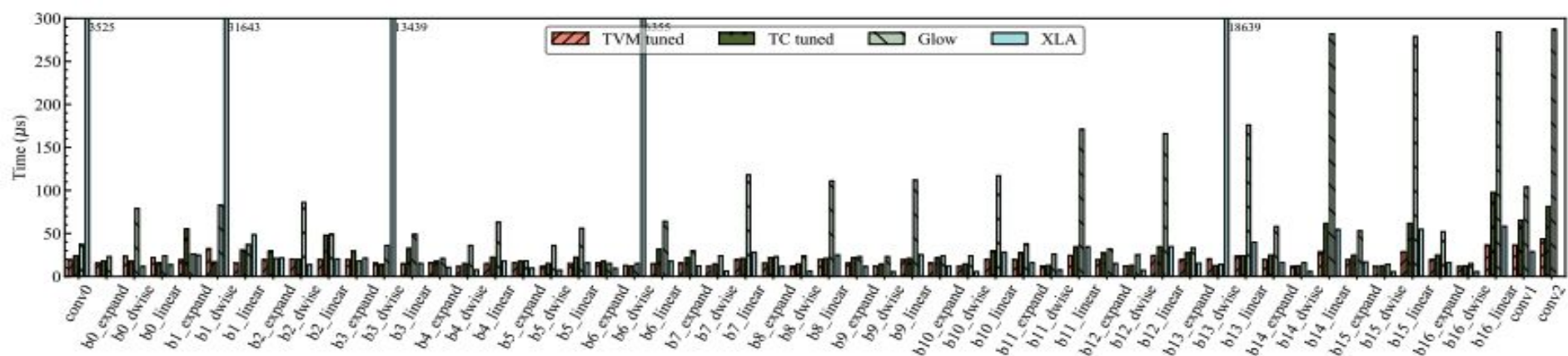
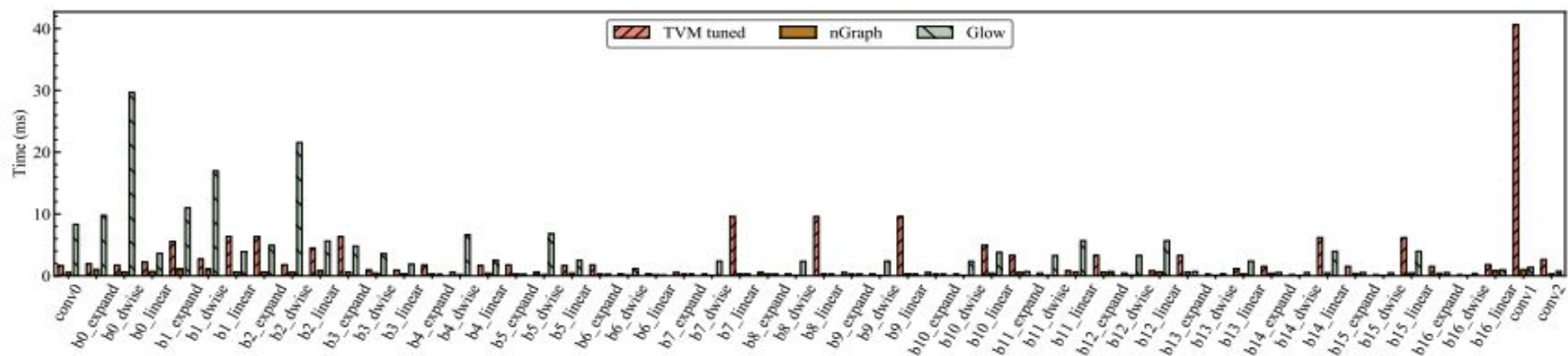


Fig. 6. The performance comparison of convolution layers in MobileNetV2_1.0 across TVM, TC, Glow and XLA on V100 GPU.



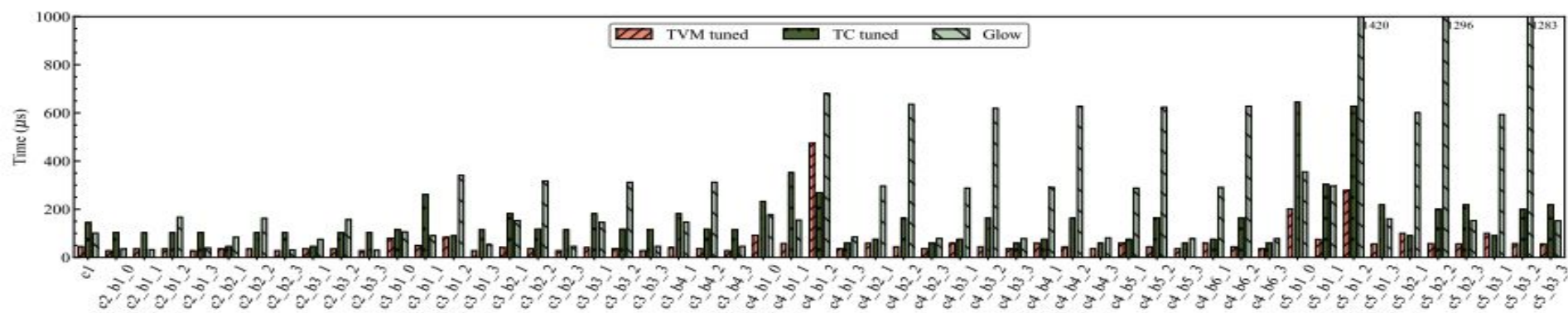


Fig. 8. The performance comparison of convolution layers in ResNet50 across TVM, TC and Glow on V100 GPU.

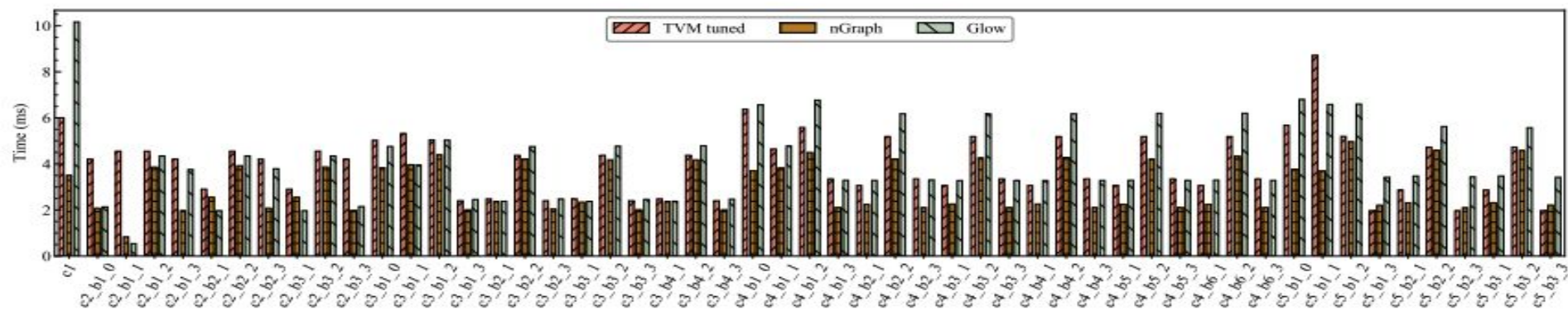


Fig. 9. The performance comparison of convolution layers in ResNet50 across TVM, nGraph and Glow on Broadwell CPU.

Table 3. The number of the clustered and non-clustered convolutions of XLA on V100 GPU and Broadwell CPU.

	MobileNetV2_1.0		ResNet50	
	Clustered	Non-clu-	Clustered	Non-clu-
V100	5	47	0	53
Broadwell	17	35	53	0