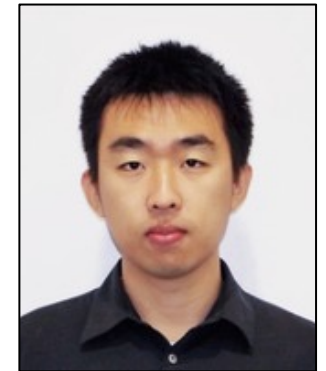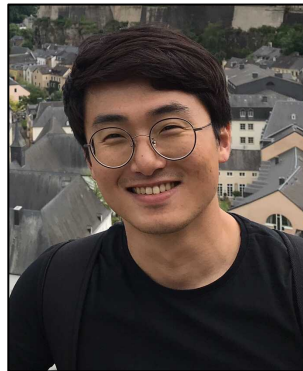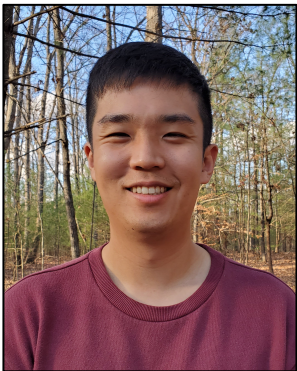# Collage: Automated Integration of Deep Learning Backends

Byungsoo Jeon*[1], Sunghyun Park*[2], Peiyuan Liao[1,4], Sheng Xu[3], Tianqi Chen[1,2], Zhihao Jia[1]

[1]*Carnegie Mellon University,* [2]*OctoML,* [3]*Amazon Web Services,* [4]*Praxis Pioneering*
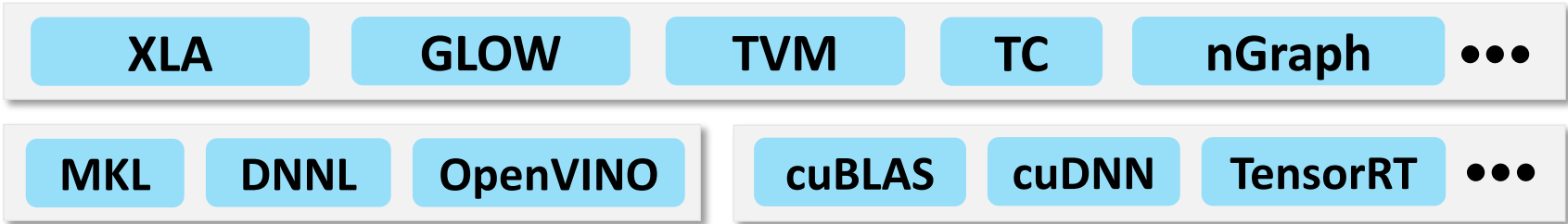
# Deep Learning (DL) Backend

## Backend

- a software library or a runtime framework that takes DL workloads as inputs and generates an optimized low-level target code
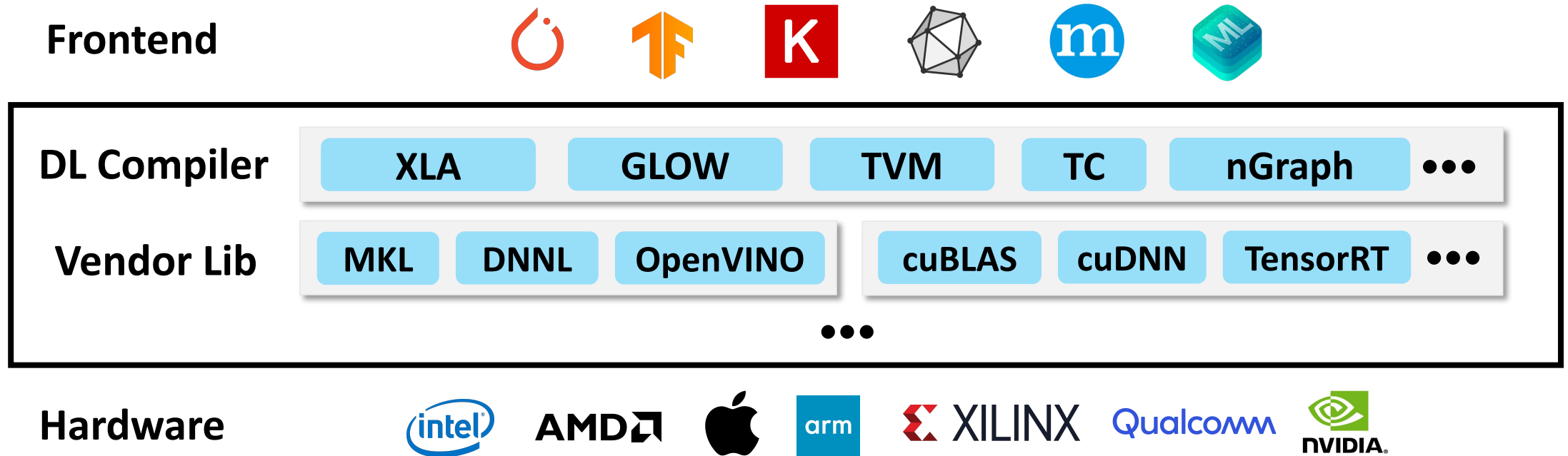
**Frontend**

**Backend**

| XLA | GLOW | TVM | TC | nGraph | ••• |

| MKL | DNNL | OpenVINO | cuBLAS | cuDNN | TensorRT | ••• |

•••

**Hardware**

# Deep Learning (DL) Backend

## Backend

o a software library or a runtime framework that takes DL workloads as inputs and generates an optimized low-level target code

**Frontend**

**DL Compiler**

| XLA | GLOW | TVM | TC | nGraph | • • • |

| MKL | DNNL | OpenVINO | cuBLAS | cuDNN | TensorRT | • • • |

• • •

**Hardware**

# Deep Learning (DL) Backend

## Backend

o a software library or a runtime framework that takes DL workloads as inputs and generates an optimized low-level target code

**Frontend**

| XLA | GLOW | TVM | TC | nGraph | ••• |

**Vendor Lib**

| MKL | DNNL | OpenVINO | | cuBLAS | cuDNN | TensorRT | ••• |

•••

**Hardware**

# Observation: Diversified DL Backends

DL backends are highly diversified and evolving fast

o Each backend has its own coverage (e.g., HW, DL operators) and strength

**Frontend**



**DL Compiler** | XLA | GLOW | TVM | TC | nGraph ●●●

**Vendor Lib** | MKL | DNNL | OpenVINO | cuBLAS | cuDNN | TensorRT ●●●

●●●

**Hardware**

# Problem: Backend Integration

Backend Integration = Backend Register + Backend Placement

# Problem: Backend Integration

Backend Integration = **Backend Register** + Backend Placement

# Problem: Backend Integration

Backend Integration = Backend Register + **Backend Placement**

# Existing Approach: Manual Backend Integration

**DL Workloads**



**Rule-based Heuristics**

**Existing DL Frameworks**

```
If op == "conv2d":
    If cudnn_enabled:
        lower_to_cudnn_kernel
    Else if …
Else if op == "batch_matmul":
    If cublas_enabled:
        lower_to_cublas_kernel
    ….
```

**Execution**

o Heuristics are often sub-optimal and susceptible to be outdated

o Direct code modification to the DL framework is required

# Our Approach: Automated Backend Integration



o It eliminates manual efforts to design heuristics and change codes

o It provides fast and stable performance across different models and hardwares

# System Overview

# Overview

# Overview



**~ 70 LoC to integrate one backend**

# Overview

# Overview



**Built-in patterns support most of popular backends (e.g., cuDNN, cuBLAS, TensorRT, TVM, MKL, etc.)**

# Overview

# Overview

# Overview

# End-to-end Evaluation: NVIDIA V100, Intel Xeon



o Stable performance across different networks and hardwares

# Optimized Backend Placement



o *Collage* leverages unique strength of each backend

o Collage maps same type of operators to different backends based on the performance landscape
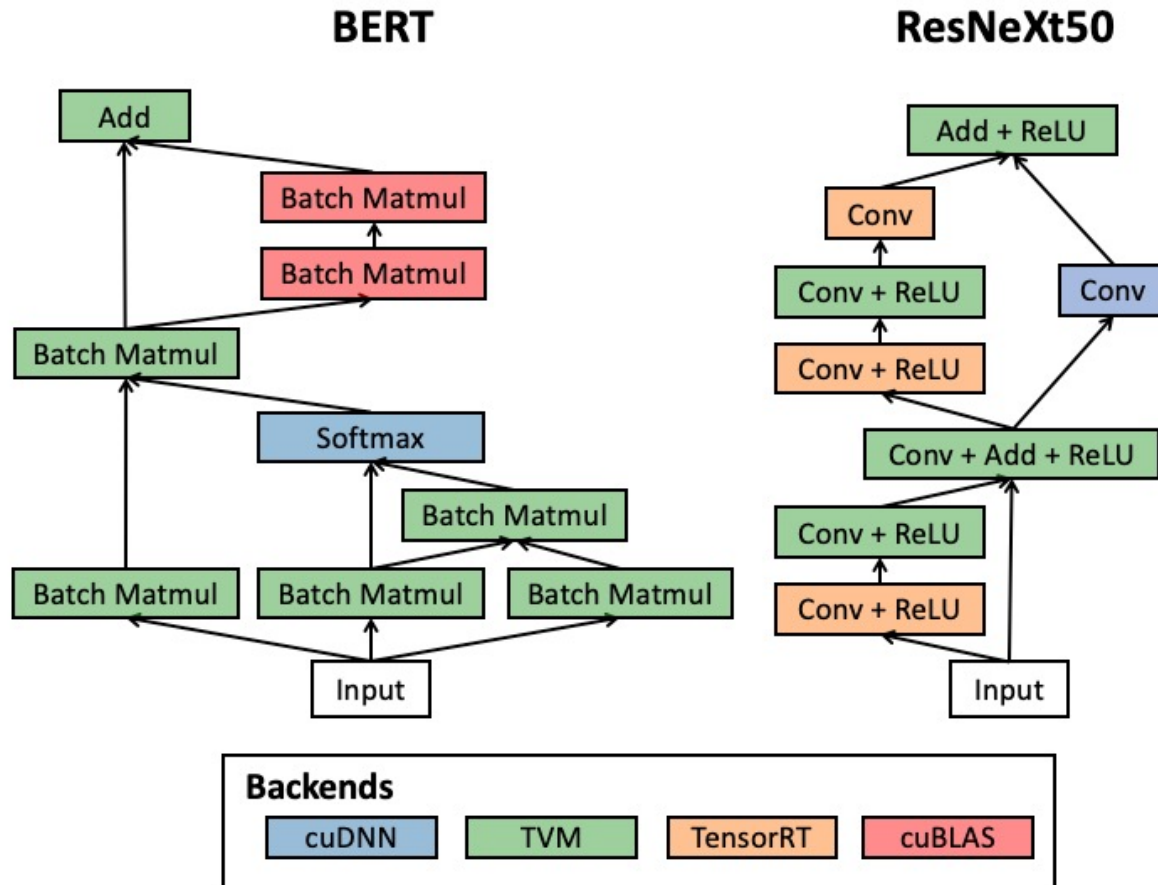
o *Collage* employs diverse operator fusion patterns

# References

Arxiv Paper: https://arxiv.org/abs/2111.00655

Code: https://github.com/cmu-catalyst/collage