# 论文分享：面向特定领域加速器的计算图高效调度

## Effectively Scheduling Computational Graphs of Deep Neural Networks toward Their Domain-Specific Accelerators

Jie Zhao
*Information Engineering University*

Siyuan Feng
*Shanghai Jiao Tong Univerisity*

Xiaoqiang Dan, Fei Liu, Chengke Wang, Sheng Yuan, Wenyuan Lv, Qikai Xie
*Stream Computing Inc.*

先进编译实验室
Advanced Compiler

先进编译实验室
Advanced Compiler

先进编译实验室
Advanced Compiler

嘉宾：

先进编译实验室
Advanced Compiler

# Effectively Scheduling Computational Graphs of Deep Neural Networks toward Their Domain-Specific Accelerators (OSDI '23)

- Introduction and Background

- Core Idea and Overview

- Scheduling Sub-graph Instances

- Kernel Generation for Sub-graph Instances

- Experimental Results

- Conclusion

■ 摩尔定律的放缓 → DNN DSA 的发展

■ DNN DSA 抽象

- ◆ d：Cluster
- ◆ c：Core
- ◆ u：CU（Compute Unit）
- ◆ LB
- ◆ GB
- ◆ DDR



先进编译实验室
Advanced Compiler

- **DNN Models**
  - ◆ Stage
  - ◆ Block
  - ◆ Layer

- **问题**
  - ◆ Kernel间频繁的数据移动
  - ◆ 错失跨层指令调度的机会
  - ◆ 没有充分利用本地内存

先进编译实验室
Advanced Compiler

- 问题

  ◆ Kernel间频繁的数据移动

  ◆ 错失跨层指令调度的机会

  ◆ 没有充分利用本地内存

- 解决方案

  ◆ 最大限度地使输入张量驻留于本地内存

  ◆ 子图划分后进行调度

  ◆ 跨Core利用并行性并使本地内存饱和

- 详细方案—GraphTurbo

  ◆ 调度子图实例

  ◆ Kernel生成



(a) Under-utilized.    (b) Saturated.

tensor used to store a batch of images    LB of Fig.1

images processed by the stage with the same color

■ 收集划分信息

$$SplitInfo = (split_d, n_d, f_d, d)$$   Core数量 LB容量等硬件信息   输出张量 → 中间变量 → 输入张量

■ 对子图进行分组

拓扑排序

三种分组模式：Straight、Diamond、Branch



(a) Straight.   (b) Diamond.   (c) Branch.



先进编译实验室
Advanced Compiler

- 收集划分信息
- 对子图进行分组
- 子图实例排序

消除实例间冗余依赖

确定调度顺序



- 推断Core绑定和Buffer范围
- 连接实例的输出

■ 层内循环融合

■ 跨层/块Buffer缝合



先进编译实验室
Advanced Compiler

■ 层内循环融合    ■ 跨层/块Buffer缝合

■ 内存分配和重用



(a) Execute $ct_1$. (b) Execute $ct_4$. (c) Execute $ct_5$. (d) Execute $ct_7$.

■ 跨层指令调度



先进编译实验室
Advanced Compiler

- TVM

  在子图内进行算子融合，生成的Kernel通过DDR交换数据

- Astitch

  没有对子图实例进行排序，也没有考虑计算密集型算子

- GraphTurbo

  TVM 11.15×　　　AStitch  6.16×

  供应商手工实现　1.04×

| label | model | batch size | batches per cluster | | throughput unit | TVM's result |
|---|---|---|---|---|---|---|
| | | | TVM | GraphTurbo | | |
| Ⓐ | ResNet-50 | 64 | 2 | 16 | images/s | 1064 |
| Ⓑ | BERT-128 | 32 | 4 | 8 | sentences/s | 512 |
| Ⓒ | BERT-256 | 16 | 2 | 4 | sentences/s | 412 |
| Ⓓ | BERT-384 | 8 | 1 | 2 | sentences/s | 36 |
| Ⓔ | BERT-512 | 8 | 1 | 2 | sentences/s | 324 |
| Ⓕ | DLRM | 1024 | 64 | 256 | queries/s | 131000 |
| Ⓖ | MobileNet-v2 | 128 | 2 | 32 | images/s | 1416 |
| Ⓗ | Vision_Transformer | 32 | 4 | 8 | images/s | 40 |
| Ⓘ | DenseNet | 32 | 4 | 8 | images/s | 456 |
| Ⓙ | Conformer | 12 | 1 | 3 | sentences/s | 184 |

[1] Zhao J, Feng S, Dan X, et al. Effectively Scheduling Computational Graphs of Deep Neural Networks toward Their {Domain-Specific} Accelerators[C]//17th USENIX Symposium on Operating Systems Design and Implementation (OSDI 23). 2023: 719-737.
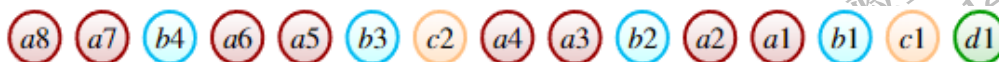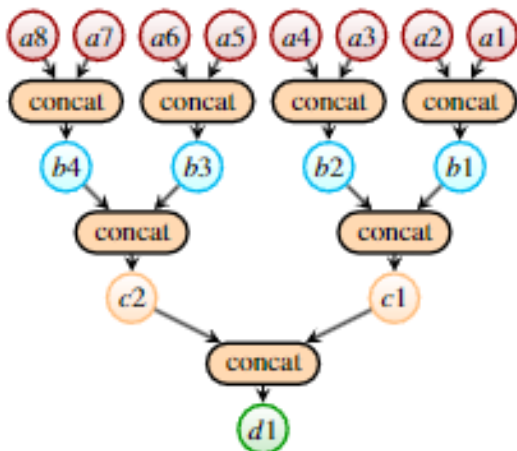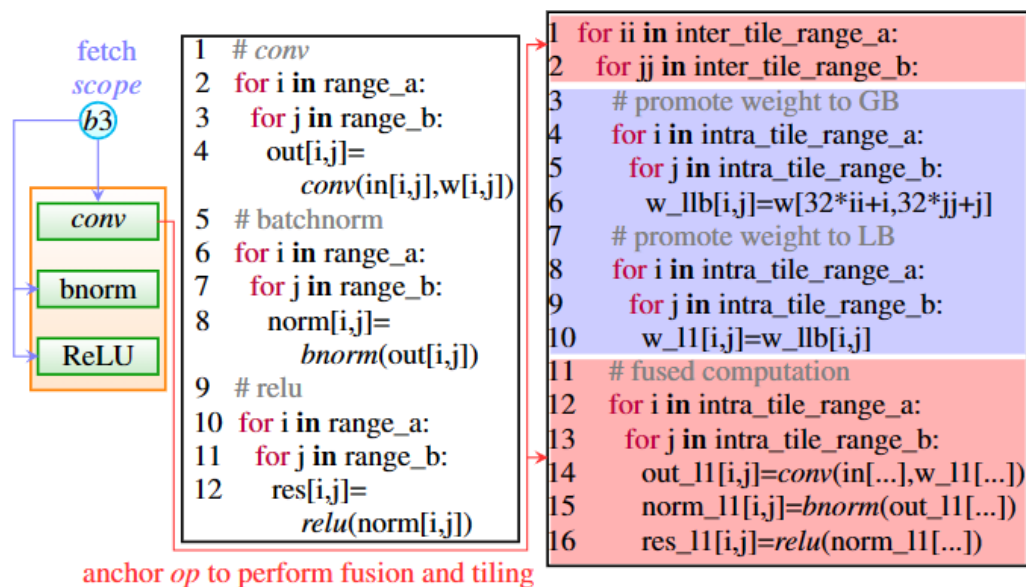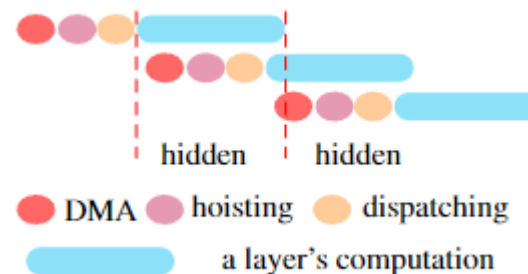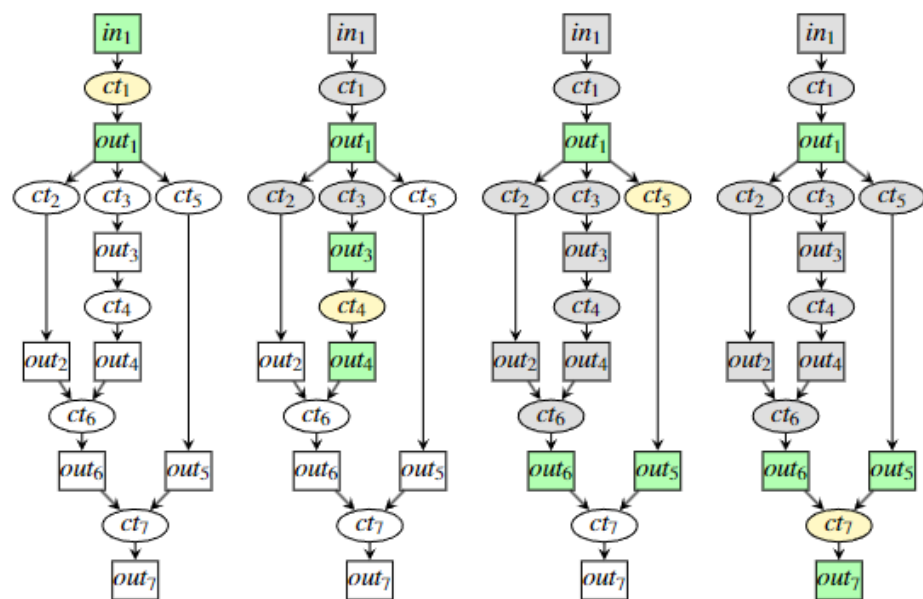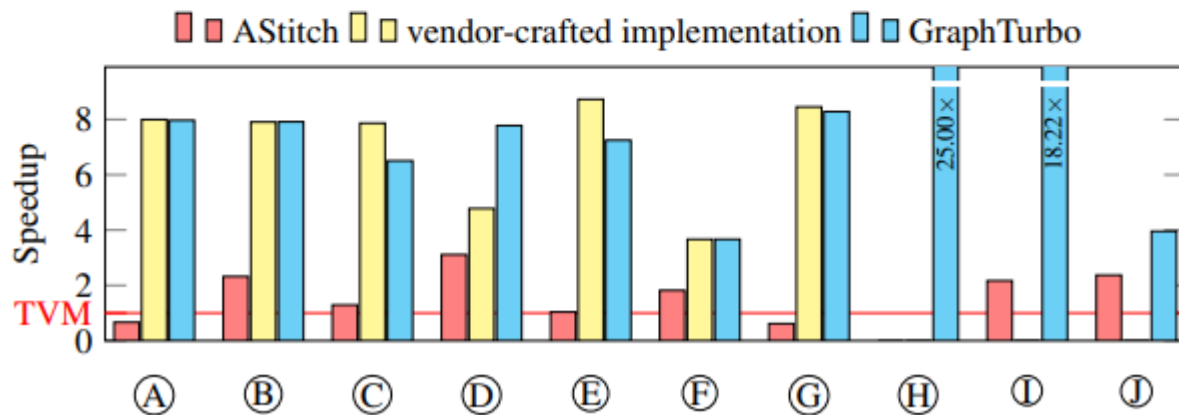
[2]