# 动态控制流编译优化
# 论文分享：Cocktailer

COCKTAILER: Analyzing and Optimizing Dynamic Control Flow in Deep Learning

Chen Zhang[†£◇*]     Lingxiao Ma[◇]     Jilong Xue[◇]     Yining Shi[‡◇*]     Ziming Miao[◇]

Fan Yang[◇]     Jidong Zhai[†]     Zhi Yang[‡]     Mao Yang[◇]

[†]Tsinghua University     [‡]Peking University     [◇]Microsoft Research

嘉宾：王磊

# Cocktailer : Analyzing and Optimizing Dynamic Control Flow in Deep Learning　　　(OSDI '23)

模型表示

Tile-based IR

硬件抽象

Rammer (OSDI`20)

Welder (OSDI`23)

Roller (OSDI`22)

Grinder (OSDI`23)

\* 示意源于微软亚洲研究院文章
《微软亚洲研究院推出AI编译器界"工业量产级"部件》
https://www.msra.cn/zh-cn/news/features/ai-compilor

■ 深度学习模型：全连接 → 复杂控制逻辑结构：数据流 + 控制流

◆ 动态计算　　◆ 条件计算　　◆ 高效计算

■ 控制流分类

◆ 循环 Loop　　◆ 分支 Branch　　◆ 递归 Recursion

Cocktailer : Analyzing and Optimizing
Dynamic Control Flow in Deep Learning

先进编译实验室
Advanced Compiler

**4**

- 当前框架对控制流支持的实现

  - 控制流的表示：表示为特殊算子 or 使用编程语句

  - 计算的调度：数据流 → GPU等加速设备　控制流 → CPU

- 存在的问题

  - 带来频繁的同步和通信开销

  - 难以进行跨控制流边界优化

  - 忽略数据流的并行执行机会

在加速器侧协同调度数据流和控制流？ → 数据流并行性 ≠ 控制流并行性

控制流并行性 ≠ 数据流并行性



◆ 数据流并行中每个算子多层次并行，且共享相同的控制逻辑

◆ GPU等加速设备支持控制流指令

◆ 程序表示：以更细的粒度表示控制流，统一数据流和控制流表示

◆ 协同调度：正确映射到并行处理单元执行

在加速器侧协同调度数据流和控制流

uTask > uOperator > uProgram

● uTask

```
interface uTask {
 void compute();
 void get_input_data();
};
```

● Nested uTask

```
interface NestedUTask: uTask {
 void compute();
 void get_input_data();

 vector<uTask> body_uTasks;
};
```

● uOperator

```
interface uOperator {
 void compute(uTask_id);
 size_t get_uTask_num();

 set<uTask> uTasks;
};
```

● uProgram

```
interface uProgram {
   void compute(uTask_id);
   size_t get_uTask_num();
   set<uTask> uTasks;
   size_t unit_level;
};
```



accelerator's abstraction

MatMul   Add   Relu

- 调度接口

```
Function ScheduleProgram(g, D, unit_level):
    // g is a graph of operators in uTask representation or a control
      flow operation that the body has been scheduled
    resource = GetResource(D, unit_level); // calculate resource
    cfg = SetResource(D, unit_level, resource);
    if g ∈ D.ControlFlow then
        return ScheduleControlFlow(g, D, unit_level, cfg);
    else
        for op ∈ g.TopoSort() do
            ScheduleOperator(op, D, unit_level, cfg);
        return cfg.uProg;
```

- 调度约束
  - ◆ 保证同步点uTask 间依赖关系
  - ◆ 控制流 unit_level ≤ 数据流 unit_level

- 调度优化
  - ◆ 函数内联  ◆ 循环展开  ◆ 递归展开

- 调度策略：自底向上遍历

```
Function Schedule(g, D, unit_level = NULL):
    ulevel=unit_level, ulevel_max=D.unit_levels.size()-1, uProgs=[];
    if g ∈ D.Operators and ulevel is NULL then
        ulevel = 0;
    if ulevel is NULL then
        for op ∈ g.TopoSort() do
            g_op, ulevel_op = Schedule(op, D, NULL);
            ulevel = max(ulevel, ulevel_op);
            uProgram_p = uProgs[-1];
            ulevel_m = max(ulevel_op, uProgram_p.ulevel);
            if ulevel_m < ulevel_max then
                g_merge = uProgram_p.g + g_op;
                g_merge, ulevel_merge = Schedule(g_merge, D, ulevel_m);
                if ulevel_merge < ulevel_max then
                    uProgs[-1] = g_merge.uProgs[0];
                    ulevel = max(ulevel, ulevel_merge);
                    continue;
            uProgs.append(g_op.uProgs);
    else
        uProgs = g.uProgs
    if ulevel < ulevel_max then
        for ulevel_cur ∈ range(ulevel, ulevel_max) do
            uProgram_cur = ScheduleProgram(g, D, ulevel_cur);
            if uProgram_cur is not NULL then
                if Latency(uProgram_cur) < Latency(uProgs) then
                    uProgs = [uProgram_cur];
                    ulevel = ulevel_cur;
    g.uProgs = uProgs;
    return g, ulevel;
```
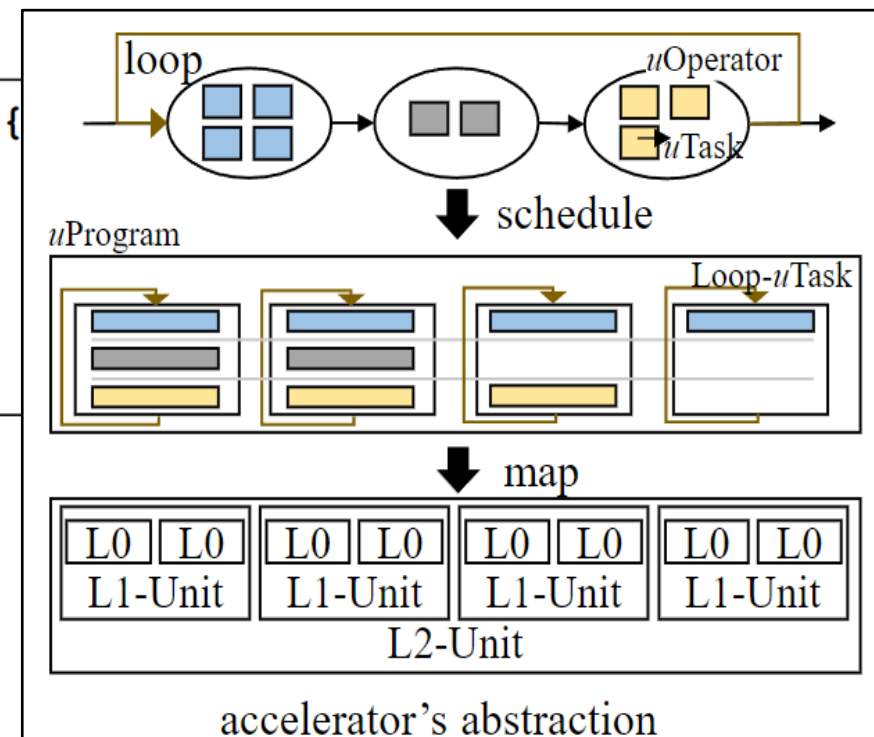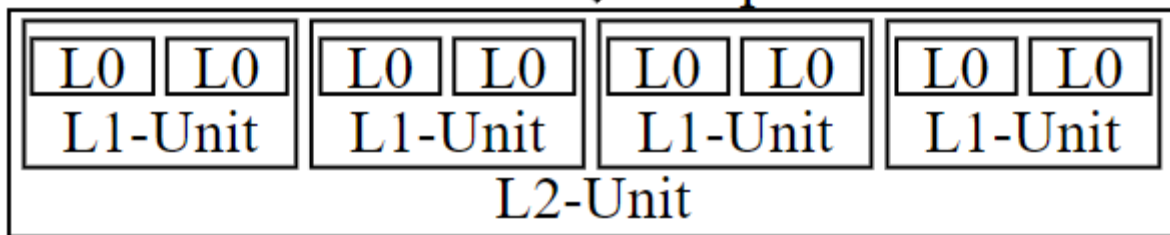
- 对接 框架/编译器

  ◆ 框架：PyTorch程序 → ONNX计算图 → Cocktailer → 自定义PyTorch算子

  ◆ 编译器：数据流算子Kernel → Rammer生成uTask组成的uOperator → Cocktailer
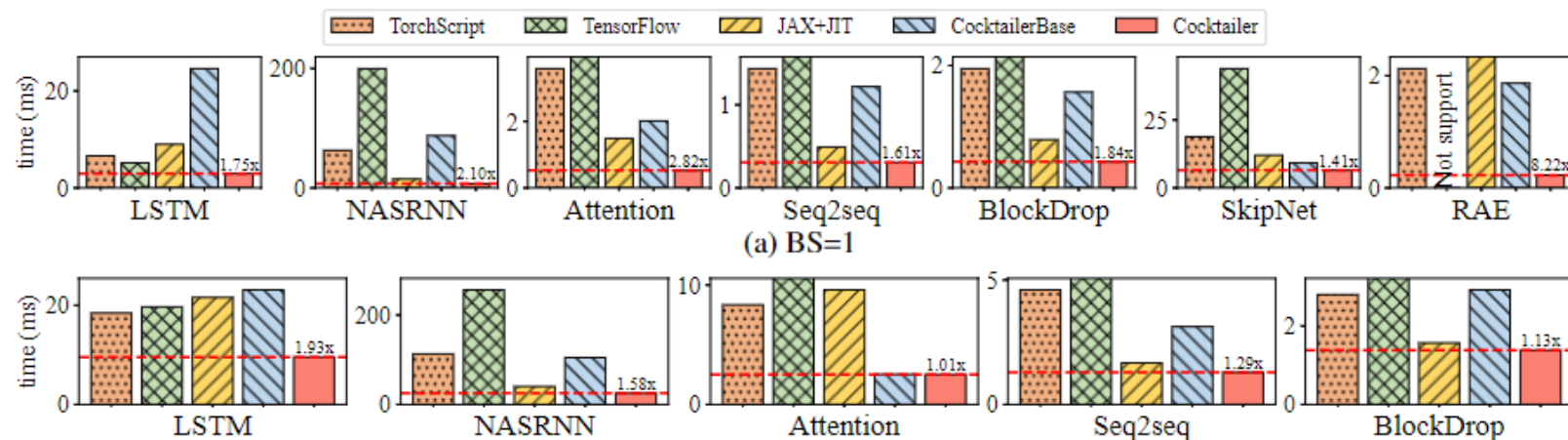
- 对接 加速器
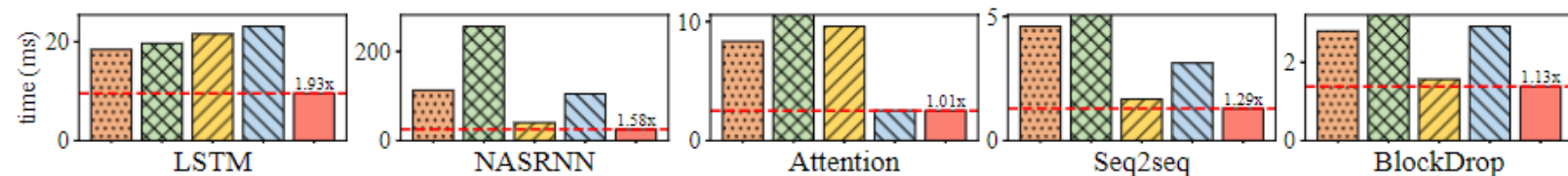
  ◆ NIVDIA GPU   ◆ AMD GPU   ◆ Graphcore IPU



- 其他实现细节

  ◆ Block alignment      ◆ Memory management      ◆ Simulation of GPU stack

  ◆ Register pressure      ◆ Branch reclustering

# 性能

- NIVDIA GPU



(a) BS=1

(b) BS=64. RAE and SkipNet cannot be batched for execution.

Figure 14: End-to-end DNN inference on NVIDIA V100 GPU

- AMD GPU

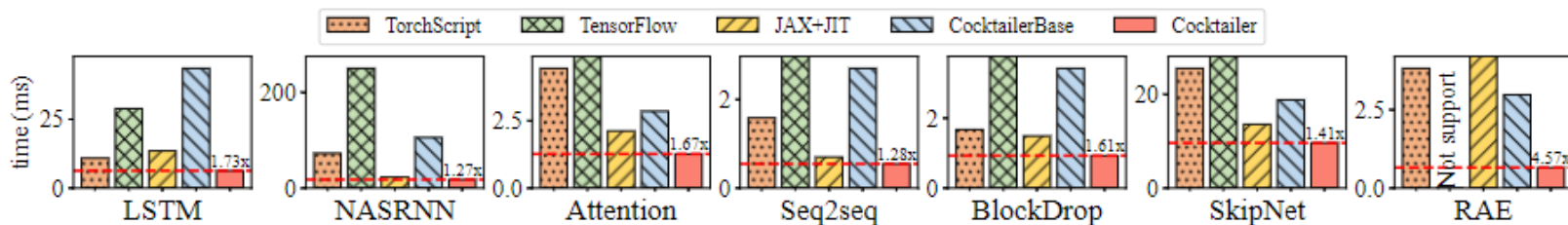

Figure 20: End-to-end DNN inference on AMD MI100 GPU with BS=1

[1] 微软亚洲研究院推出AI编译器界 "工业重金属四部曲"：https://www.msra.cn/zh-cn/news/features/ai-compilor

[2] Cocktailer: Analyzing and Optimizing Dynamic Control Flow in Deep Learning. Chen Zhang, Lingxiao Ma, Jilong Xue, Yining Shi, Ziming Miao, Fan Yang, Jidong Zhai, Zhi Yang, Mao Yang. OSDI 2023.

先进编译实验室
Advanced Compiler