



先进编译实验室
Advanced Compiler

编译论坛

动态shape深度学习算 子自动调优DietCode

嘉宾：王一帆



先进编译实验室
Advanced Compiler





DIETCODE: AUTOMATIC OPTIMIZATION FOR DYNAMIC TENSOR PROGRAMS

Bojian Zheng^{*123} Ziheng Jiang^{*4} Cody Hao Yu² Haichen Shen⁵ Josh Fromm⁶
Yizhi Liu² Yida Wang² Luis Ceze⁷⁶ Tianqi Chen⁸⁶ Gennady Pekhimenko¹²³

ABSTRACT

Achieving high performance for compute-intensive operators in machine learning (ML) workloads is a crucial but challenging task. Many ML and system practitioners rely on vendor libraries or auto-schedulers to do the job. While the former requires large engineering efforts, the latter only supports static-shape workloads in existing works. It is difficult, if not impractical, to apply existing auto-schedulers directly to dynamic-shape workloads, as this leads to extremely long auto-scheduling time.

We observe that the key challenge faced by existing auto-schedulers when handling a dynamic-shape workload is that they cannot construct a unified search space for all the possible shapes of the workload, because their search space is shape-dependent. To address this, we propose *DietCode*, a new auto-scheduler framework that efficiently supports dynamic-shape workloads by constructing a *shape-generic* search space and cost model. Under this construction, all shapes *jointly* search within the same space and update the same cost model when auto-scheduling, which is therefore more efficient compared with existing auto-schedulers.

We evaluate *DietCode* using state-of-the-art machine learning workloads on a modern GPU. Our evaluation shows that *DietCode* has the following key strengths when auto-scheduling an entire model end-to-end: (1) reduces the auto-scheduling time by up to $5.88\times$ less than the state-of-the-art auto-scheduler on the uniformly sampled dynamic shapes ($94.1\times$ estimated if all the possible shapes are included), (2) improves performance by up to 69.5% better than the auto-scheduler and 18.6% better than the vendor library. All these advantages make *DietCode* an efficient and practical solution for dynamic-shape workloads.

1 INTRODUCTION

Deep neural networks (DNNs) form an important class of ML algorithms (He et al., 2016; Vaswani et al., 2017; Amodi et al., 2016; Devlin et al., 2019). They are made of tensor operators which are often executed for tens of

ML frameworks rely on heavily optimized, hand-crafted vendor libraries (e.g., oneDNN (oneAPI, 2021) on Intel CPUs; cuDNN (Chetlur et al., 2014) and cuBLAS (NVIDIA, 2021) on NVIDIA GPUs) to provide highly optimized operators (Abadi et al., 2016; Paszke et al., 2019; Chen et al., 2015). Despite delivering high performance, the develop-



- 01. 背景与动机
- 02. 主体思路
- 03. 形状通用搜索空间
- 04. 成本模型
- 05. 自动调度器
- 06. 性能

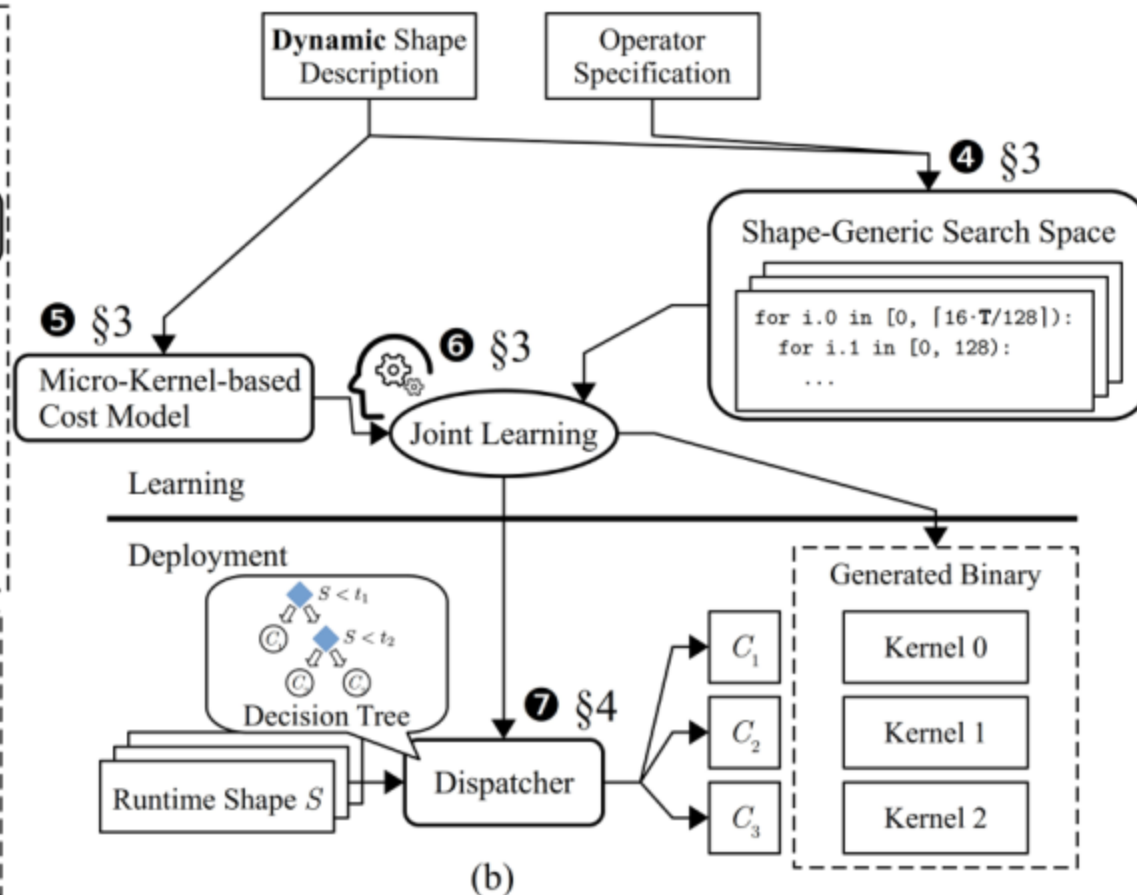
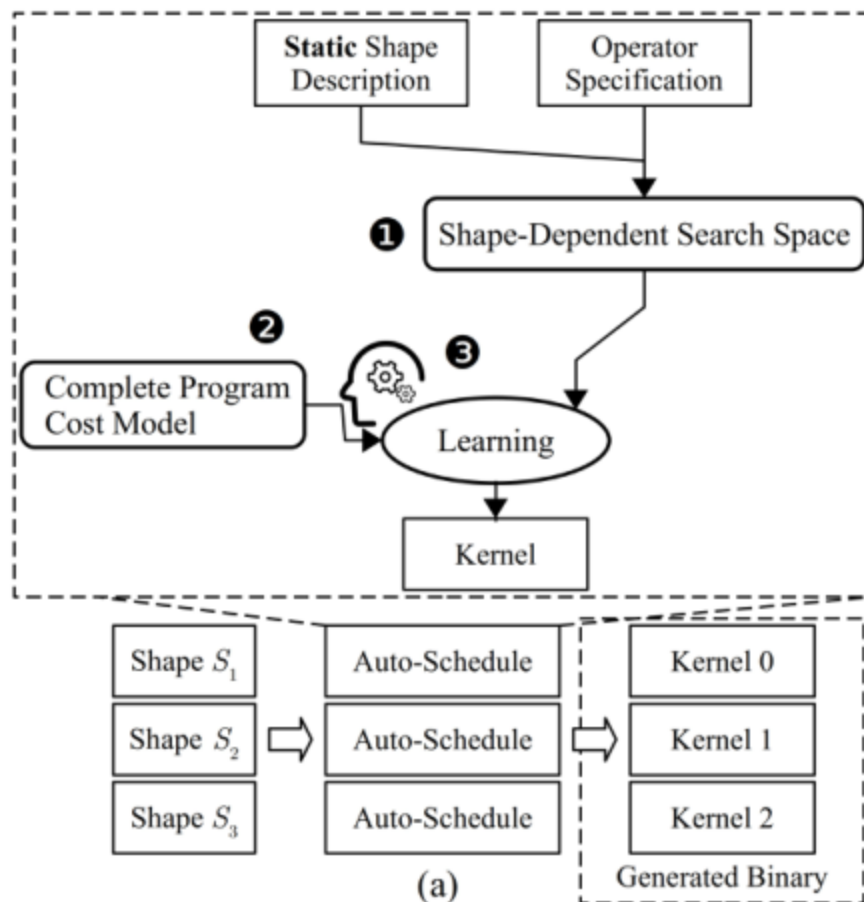




研究背景

深度学习模型已经广泛地应用于日常的生产与生活。算子是模型的最小执行单元，其性能极大地影响着模型整体的执行效率。高性能的算子需要针对目标硬件的架构特点进行优化（称为调度），以充分利用目标硬件的计算能力。而软件对于硬件的适配需要大量软硬件跨度的专家来解决问题。他们会根据硬件的架构、指令等对算子进行适配。



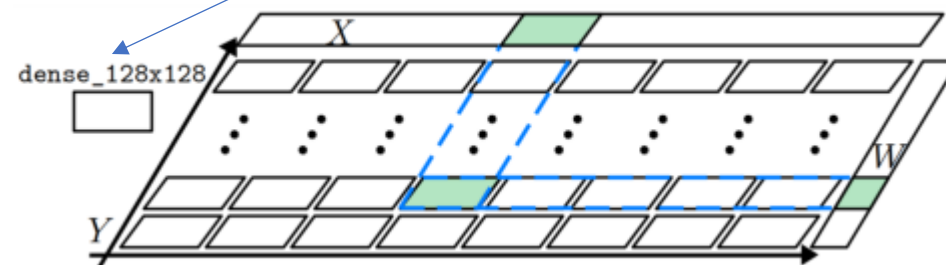


```
for (i = 0; i < T; ++i)  for (io = 0; io <  $\lceil T/t \rceil$ ; ++io)
  A[i] = ...             for (ii = 0; ii < t; ++ii)
                          A[io*t+ii] = ...
```

(a) (b)

观察到，不同形状搜索空间实际上可以重叠，因此可以潜在地形成形状通用搜索空间，而不是使用现有的自动调度器逐个调整每个可能的形状。

$$Y = XW^T, X : [128, 768], W : [128, 768]$$



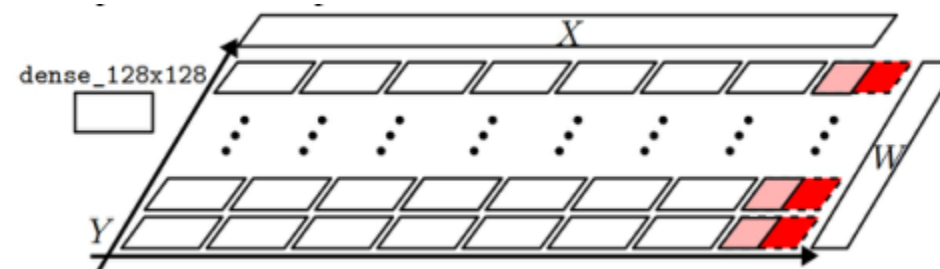
(a) $X : [1024, 768], W : [2304, 768]$
 $Y = XW^T$
(b) $X : [16 \times T, 768], W : [2304, 768],$
 $T \in [1, 128]$




```
for (i = 0; i < T; ++i)  A[i] = ...
for (io = 0; io <  $\lceil T/t \rceil$ ; ++io)
  for (ii = 0; ii < t; ++ii)
    A[io*t+ii] = ...
```

(a) (b)

```
for i.0 in [0, T):
  for i.1 in [0, t):
    if i.0*t+i.1 < T:
      X_local = X[...]
    if i.0*t+i.1 < T:
      Y_local = ...
    if i.0*t+i.1 < T:
      Y[...] = Y_local
```

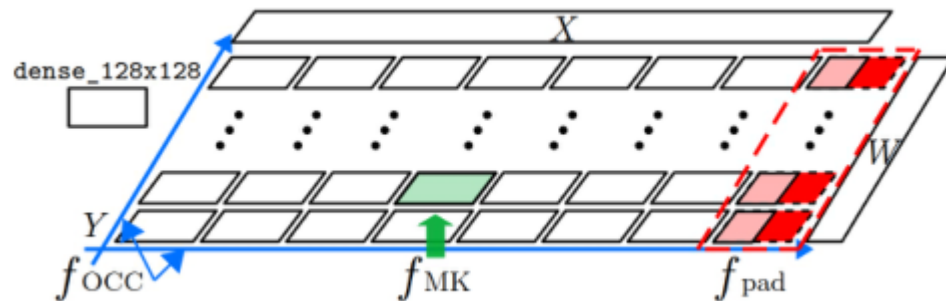


然而， workload实例通常不完全适合微内核，如图所示，最后一列的微内核没有完全具体化。在这些情况下，需要在微内核内部注入边界检查，以确保程序不会在无效数据值上运行，但它们也会给程序带来很大的性能下降这是因为这些检查会在生成的程序中引入额外的分支和计算指令。



为了准确预测基于微内核的完整程序的性能，我们在成本模型中设计了三个项，分别说明：①微内核的性能；②硬件内核占用惩罚，其与执行该微内核以组成完整程序的次数相关；③填充惩罚。成本模型的数学表达式如下：

$$\text{Cost}_M(P) = f_{\text{MK}}(\text{FeatureExtractor}(M)) \textcircled{1} \cdot \underbrace{\textcircled{2} f_{\text{OCC}}(P/M) \cdot f_{\text{pad}}(P, M) \textcircled{3}}_{f_{\text{adapt}}(P, M)}$$



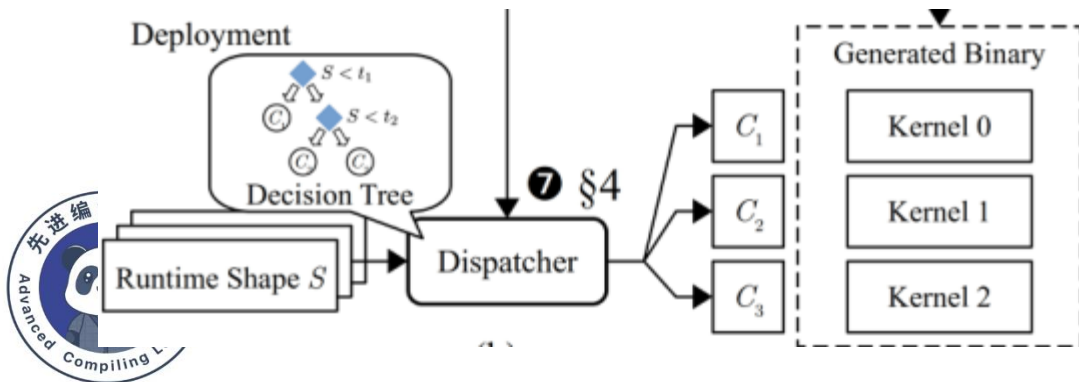
图说明了方程中三个项之间的相关性以及微内核如何形成完整程序，其中我们可以看到方程和完整程序组成之间的一一对应关系。这种对应关系使我们能够准确预测基于微内核的完整程序的性能，从而在联合学习过程中使用进化搜索有效地搜索高性能微内核。

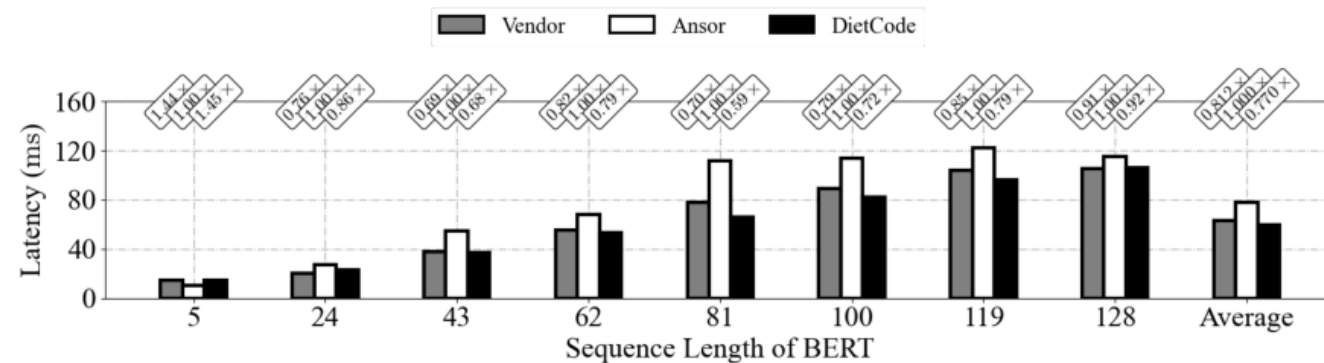


学习过程完成后，DietCode将生成一组微内核作为自动调度结果。为了将所有可能的形状分配给这些微内核，我们让每个形状 S 根据等式3中的成本公式投票给其最喜欢的微内核

$$\text{vote}(S) = \underset{M}{\operatorname{argmax}} (\text{Cost}_M(P(S, M)))$$

我们使用scikit学习框架训练决策树。决策树的输入是所有可能的形状，输出标签是它们选择的微内核。





评估表明，在DietCode上，可以将自动调度时间显著缩短5.88倍（如果包括所有可能的形状，则预测为94.1倍），同时在端到端的完整最先进DNN模型上，性能比最先进的自动调度程序高69.5%，比供应商库高18.6%。

