GOPTX: FINE-GRAINED GPU KERNEL FUSION BY PTX-LEVEL INSTRUCTION FLOW WEAVING

Kan Wu, Zejia Lin, Mengyue Xi, Zhongchun Zheng, Wenxuan Pan, Xianwei Zhang#, Yutong Lu# Sun Yat-sen University, Guangzhou, China

{wukan3, linzj39, ximy, zhengzhch3, panwx5}@mail2.sysu.edu.cn, {zhangxw79, luyutong}@mail.sysu.edu.cn

Abstract

GPUs have been heavily utilized in diverse applications, and numerous approaches, including kernel fusion, have been proposed to boost GPU efficiency through concurrent kernel execution. However, these approaches generally overlook the opportunities to mitigate warp stalls and improve instruction level parallelism (ILP) in inter-kernel resource sharing. To address this issue, we introduce GoPTX, a novel design for kernel fusion that improves ILP through deliberate weaving instructions at the PTX level. GoPTX establishes a merged control flow graph (CFG) from original kernels, enabling to interleaving of instructions that were sequentially executed by default and minimizing pipeline stalls on data hazards. We further propose a latency-aware instruction weaving algorithm for more efficient instruction scheduling and an adaptive code slicing method to enlarge the scheduling space. Experimental evaluation demonstrates that GoPTX achieves an average speedup of 11.2% over the baseline concurrent execution, with a maximum improvement of 23%. The hardware resource utilization statistics show significant enhancements in eligible warps per cycle and resource use.



Interleaving independent instructions extends dependency chains, avoiding stalls and hiding latency by filling pipeline bubbles.

Weaved Kernel

1 mov %f1, 1

2 mov %s1, 1

3 mov %f2, %f1

4 mov %s2, %s1

5 mov %f3, %f2

Ouput PTX Overall workflow of GoPTX, mov %1f,1 with three key phases of Slicing, Weaving and Merging, and 17 mov %9s,9 additional phases of Prep-18 mov %9f,9 rocessing and Postprocessing. Postprocessing Selected First CFG Merged CFG without deadlock Deadlock solved 3 4 5 1: input: $F \neq \emptyset, S \neq \emptyset$ {the first and the second input CFG} Input Flow A 1 2 by inserting dummy sync0 empty >sync0> 2: **output**: M {the merged CFG} >empty dummy dummy sync 1 3: $M \leftarrow \emptyset$ Input Flow B 6 7 7 8 9 blocks before syn-Second CFG 4: M.add node (< 0, 0 >)5: $Stack \leftarrow \{ < 0, 0 > \}$ chronizations, del- I dummy>dummy>sync 1> 6: $Pushed \leftarrow \{ < 0, 0 > \}$ Output Flow 1 6 2 7 aying the barriers top while $Stack \neq \emptyset$ do



Results

	First Kernel	Ŷ	Second Kernel
1	globallaunch_bounds(3	2) 1	globallaunch_bounds(32)
2	<pre>void first(float x) {</pre>	2	<pre>void second(float y) {</pre>
3	x1 = 1; x2 = 2;	3	y1 = 1; y2 = 2;
4	if (x < 0) {	4	if (y < 0) {
5	x2 = 1; x1 = 2;	5	y2 = 1; y1 = 2;
6	}	6	}
7	}	7	}

Illustrative examples of first and second input kernel with coarse-grained and fine- grained (Desired) fusion.

Desire Kernel _global__ _launch_bounds__(32) void desire(float x, float y) { x1 = 1; y1 = 1; x2 = 2; y2 = 2;if (x < 0 && y < 0) { $x^2 = 1; y^2 = 1;$ x1 = 2; y1 = 2;} else if (x < 0) { x2 = 1; x1 = 2;else if (y < 0) { y2 = 1; y1 = 2;11 12 13 }



