

Abstract

Heterogeneity continues to increase in all kinds of computing applications, with the rise of accelerators such as GPUs, FPGAs, APUs, and other co-processors.

Programming models, such as CUDA, OpenACC and OpenCL are designed to offload computeintensive workloads to co-processors efficiently.

Motivation and Goal

Drawbacks of current popular Kernel Level Pipelining

- End user must manually partition the task to multiple sub kernel chunks and then launch by multiple GPU streams.
- Splitting to multiple chunks may cause extra function call overhead.
- Parameters (#chunks ,#streams, etc.) must iii. be well tuned to provide optimal performance

We proposed a new block-level pipelining extension for OpenMP that:

- ✓ Handle data transfer and computation inside one kernel using different streaming multiprocessors.
- ✓ Trigger the computation using atomic operations as long as the data is ready thus pipeline the data transfer and computation.

Environment Setup

Machine1:

CPU: IBM Power9 processors GPU: NVIDIA TESLA V100 with NVLink2 Machine2: CPU: Intel Xeon Gold 6136 GPU: NVIDIA TESLA V100 with PCI-E **Benchmarks**: 1. 2D Convolution 2. Generic Matrix-Matrix Multiplication

CUDA Version: 9.0.176

 $\stackrel{\circ}{>}$ 50000





This work was performed under the auspices of the U.S. Department of Energy by Lawrence Livermore National Laboratory under Contract DE-AC52-07NA27344.

Performance Evaluation of the NVIDIA Tesla V100: Block Level Xuewen Cui, Thomas R. W. Scogland, Bronis R. de Supinski and Wu-chun Feng





		1.69	1.69	1.7	1.67	1.5	0.81	0.41	0.21	
		1.69	1.69	1.72	1.68	1.24	0.63	0.33	0.21	Speedup 2.00
		1.72	1.73	1.74	1.73	1.25	0.82	0.42	0.16	1.68
		1.64	1.73	1.76	1.72	1.14	0.82	0.39	0.17	1.50
	1.54	1.69	1.73	1.77	1.72	1.57	0.84	0.42	0.17	1.00
	1.56	1.68	1.74	1.74	1.74	1.52	0.85	0.43	0.17	- 0.50
31	1.54	1.63	1.65	1.62	1.46	1.21	0.7	0.44	0.23	- 0.00
1	1	0.96	0.95	0.89	0.81	0.71	0.45	0.36	0.17	
2	\triangleright	ଚ	~6	32	GA	20	150	512	024	Block Level
	#chunks									Pipeline
lized Speedup of Async version on V100 PCI-F										
meed speedup of Asyne version on vitoor ci-L										
		1.57	1.6	1.62	1.6	1.53	1.37	1.2	1.08	
		1.55	1.6	1.62	1.59	1.51	1.36	1.2	1.05	Speedup 2.00
		1.56	1.61	1.62	1.6	1.5	1.36	1.18	1.05	1.65
		1.57	1.6	1.62	1.6	1.51	1.36	1.16	1	
	1.46	1.56	1.61	1.62	1.6	1.48	1.31	1.13	0.91	1.00
	1.46	1.56	1.61	1.62	1.59	1.47	1.33	1.12	0.84	- 0.50
27	1.47	1.56	1.6	1.6	1.49	1.32	1.08	0.89	0.66	0.00
	0.99	0.99	0.98	0.96	0.95	0.89	0.81	0.71	0.58	
l	\triangleright	୫	~6	Br	6A	~28	250	52	024	Block Level
			#	chunk	S		•			Pipeline
Aatrix-Matrix Multiplication										
			••••				•			

with traditional kernel level pipeline on V100 GPUs.