

A hybrid piece-wise slowdown model for concurrent kernel execution on GPU

B. López-Albelda, F. M. Castro, J. M. González-Linares and N. Guil





Motivation



Kernel execution on a GPU may not scale well due to architecture constraints and saturation of computational resources like pipeline stalls, L1-cache trashing, global memory bandwidth, etc.



CTA: Cooperative Threads Array



Waste of computational power



Kernels like Reduction do not scale well





GPU resources allocation to kernels: SMT



Spatial Multitask (**SMT**): CTAs are assigned to a subset of SMs (CTAs of one kernel in light blue, CTAs of other kernels in light orange)





Concurrent Kernel Execution (CKE)



A kernel can be co-executed with another kernel with complimentary requirements to improve performance



Slowdown Model



Best results obtained by co-executing compute bound (CB) kernels together with memory bound (MB) kernels

A slowdown model can be used to obtain the best mapping:

Predict Normalized Progress, NP, for each SM allocation

$$NP = \frac{IPC^{shared}}{IPC^{alone}}$$

Use some criteria like fairness, STP or ANTT to choose the best mapping

$$fairness = \min_{i,j} \frac{NP_i}{NP_j}$$



HSM: A Hybrid Slowdown Model for Multitasking GPUs, Zhao, Jahre and Eeckhout, ASPLOS'20

- NP of compute-bound kernels is predicted as $NP = \frac{R_{cb}^{shared}}{R_{cb}^{alone}}$ where R_{cb}^{shared} is the total number of SMs allocated and R_{cb}^{alone} is the total number of SMs in the GPU.
- NP of memory-bound kernels is predicted using the effective bandwidth utilization, which can be obtained from the Row Buffer Hit rate (RBH)





A hybrid piece-wise slowdown model for concurrent kernel execution on GPU

ΔD

GPU resources allocation to kernels: SMK



Simultaneous Multikernel (**SMK**): CTAs are assigned to all SMs (CTAs of one kernel in light blue, CTAs of other kernels in light orange)





CKE using SMK



Performance under SMK allocation can improve significatively





Objectives



Compare	Build	Implement
Compare SMK with respect to SMT in terms of performance using	Build a slowdown model for SMK	Implement a fairness- based CTA allocation policy
 Average Normalized Turn- around Time (ANTT) System Throughput (STP) 		





- GPGPU-sim v4.0 modified to support SMT and SMK
- Simulated system: Volta Titan V, with 80 SMs and HBM memory with 3 stacks and 24 channels
- Benchmarks: 14 kernels from CUDA SDK, Rodinia, Parboil and Chai
 - 8 MB and 6 CB
- Concurrent kernel execution test: both kernels are launched at the same time, and stopped when one of them has no new CTAs left to allocate
- Performance metrics:

JNIVERSIDAD

)F MÁLAGA

$$NP = \frac{IPC^{shared}}{IPC^{alone}} \quad STP = NP_i + NP_j \quad ANTT = \frac{1}{2} \left(\frac{1}{NP_i} + \frac{1}{NP_j} \right) \quad fairness = \min_{i,j} \frac{NP_i}{NP_j}$$

SMK vs SMT







Using HSM with SMK



NP prediction for CB kernels under SMK allocation is not very good





HSM slowdown prediction with SMK



1. Allocate half SMs use to each kernel





HSM slowdown prediction with SMK



- 1. Allocate half SMs use to each kernel. Predict NP of both kernels
- 2. Give (or take) SMs use to the CB kernel





A hybrid piece-wise slowdown model



- 1. Start at one extrema: detect whether kernel is CB or MB
- 2. Give (or take) SMs use to the CB kernel





A hybrid piece-wise slowdown model



- 2. Give (or take) SMs use to the CB kernel. Adjust NP with the real value
- 3. Give (or take) SMs use to the CB kernel.





A hybrid piece-wise slowdown model for concurrent kernel execution on GPU

A hybrid piece-wise slowdown model



- 3. Give (or take) SMs use to the CB kernel. Adjust NP with the real value
- 4. Give (or take) SMs use to the CB kernel...





A hybrid piece-wise slowdown model for concurrent kernel execution on GPU

HSM-Fair







A hybrid piece-wise slowdown model for concurrent kernel execution on GPU

20

HPSM-Fair







A hybrid piece-wise slowdown model for concurrent kernel execution on GPU

21

Prediction error







Performance evaluation







Conclusions and future work



We have compared two CTA distribution strategies: SMT and SMK

We have built a Hybrid Piece-wise Slowdown Model that predicts accurately the normalized progress of co-executing kernels

We have implemented a fairness-based CTA allocation policy

In the future, we plan to develop a more sophisticated model for kernel performance to reduce errors and number of approximations





THANK YOU FOR YOUR ATTENTION!

ANY QUESTIONS?