

# Efficient GPU Hardware Transactional Memory through Early Conflict Resolution

Sui Chen and Lu Peng

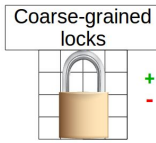
Division of Electrical and Computer Engineering, Louisiana State University

Accepted & presented at the 22th IEEE Symposium on High Performance Computer Architecture (HPCA), Barcelona, Spain

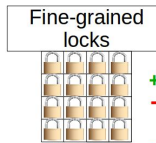
## Motivation and Goal

Goal: Conflict  $\Rightarrow$  Performance Energy

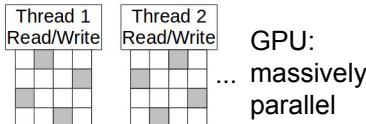
Parallel programs with critical sections



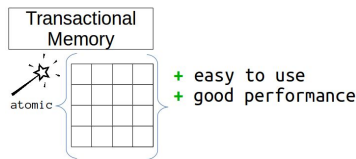
- + easy to use
- slow (serializes everything)



- + fast
- hard to program  
Example: RBTree
- error-prone



GPU: massively parallel



TM: Get the best of both coarse- and fine- locks

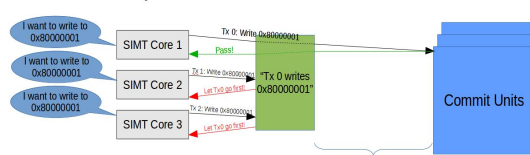
## Opportunity for Speedup

### Early abort

- SIMT Cores **can not see each other** directly



- What if they can?

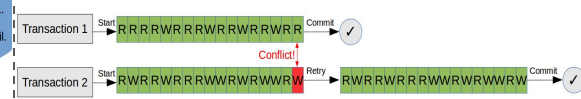


- Detect **Spatial 2-3** conflicts in **SIMT Cores**
- Reduces contention in **Commit Units**

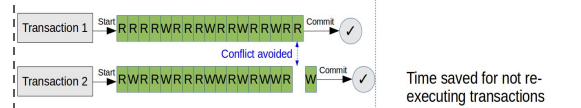
### Pause-and-Go

- Aborting and retrying may be expensive

- One small conflict wastes all transactional work



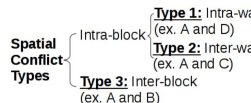
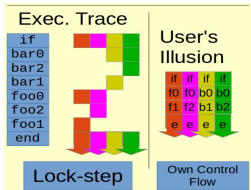
- If conflict can be avoided ...



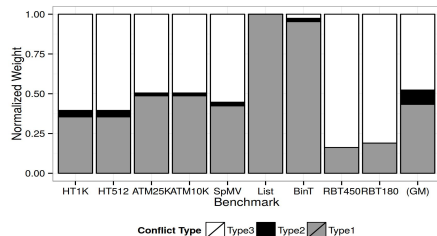
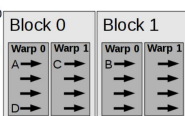
More **Temporal 2** conflicts resolved in **SIMT Cores**

Time saved for not re-executing transactions

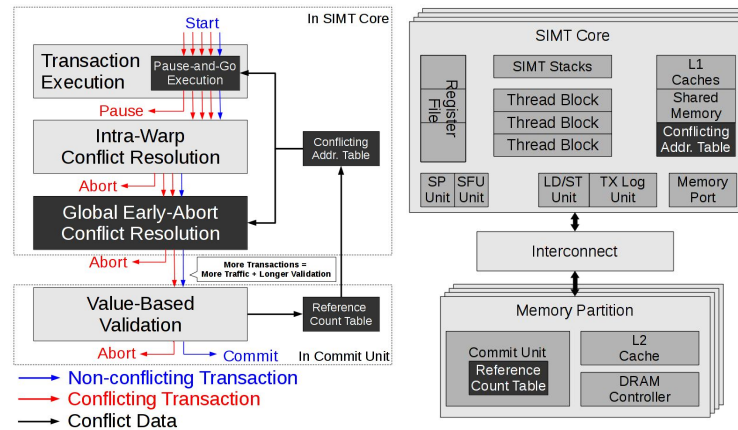
## GPU-Specific Architecture & Conflict Types



SIMT execution gives 3 spatial types of conflicts

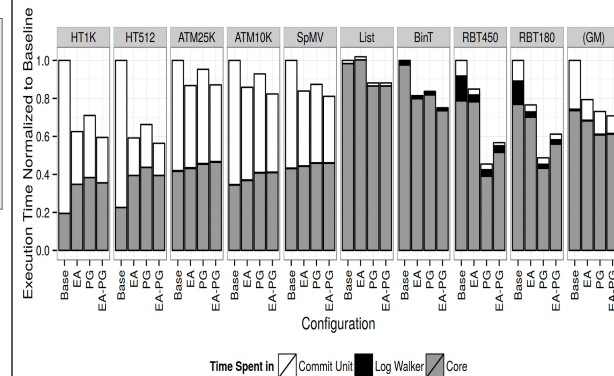


## Interaction with Hardware



- One set of hardware change facilitates both approaches
- Modified transaction execution flow

## Results



- 1.41x speedup
- 0.8x energy consumption
- Table size chosen using sensitivity study
- 5 workloads are CU-heavy; the rest are SIMT Core-heavy

Spatial types 1 and 3 frequent in workloads