# Qualitative and quantitative study of modern GPU synchronization approaches

Ilya Pershin, Vadim Levchenko and Anastasia Perepelkina Keldysh Institute of Applied Mathematics, Moscow, Russia

September 27, 2021



## Parallelism levels

#### <u>One GPU</u>



#### **CPU Cluster**





 $\frac{10^2 \mbox{ SMs}}{10^4 \mbox{ FLOP} \mbox{ / clock}} =$ 

 $10^{2\div5}$  Nodes  $\times$  32 Cores  $\times$  32 FLOP / Core =  $10^{5\div8}$  FLOP / clock

#### High parallelism hardware issues

#### Memory wall Why: Computing >> Data transfer Model: Roofline



Load balancing Why: Computing heterogeneity, Sync latency Model: ?? (Spoiler: New Tau-model) Task:

- Create high-performance 3D stencil (LBM) code
- 2 Focus on SMs' synchronization approaches
- 3 Build a quantitative model

# Lattice Boltzmann method (LBM)

- In each cell, fluid flow  $\rightarrow \{f_i\}$  with constant  $ec{c}_i$ ,  $i=1\ldots Q$
- Two-stage stencil scheme
  - Collision:  $f_i^*(\vec{x}, t) = \Omega(f_1(\vec{x}, t), \dots, f_Q(\vec{x}, t))$ Streaming:  $f_i(\vec{x} + \vec{c}_i, t + 1) = f_i^*(\vec{x}, t)$ Data transfer
- Macroscopic parameters:  $\rho = \sum_{i} f_{i}, \ \rho \vec{u} = \sum_{i} f_{i} \vec{c}_{i}$ , etc.
- Specifically, LBM D3Q19 & BGK  $\Omega$  is used





## Recursive Domain Decomposition

# RDD boosts operational intensity of a stencil code

- SM' register file localization (256 KB per SM)
- $\blacksquare$  3D mesh  $\rightarrow$  recursive rectangular blocks
- Decomposition level = Parallelism level
- Blocks' data exchange:
  - $\blacksquare$  Threads in a warp  $\rightarrow$  Warp shuffle
  - $\blacksquare Warps \rightarrow Shared memory$
  - $\blacksquare \underline{SMs} \rightarrow L2 \ \underline{Cache}$
- Time loop inside the kernel
- Pershin I. et al. GPU Implementation of a Stencil Code with More Than 90% of the Peak Theoretical Performance // RuSCDays2019





# SMs synchronization approaches

#### Asynchronous

#### Synchronous

kernel <<<blocks, threads>>>(...) cudaLaunchCooperativeKernel(kernel, blocks, threads ,...)

- Well-known standard CUDA approach
- Automatic sync at every stencil step
- Global memory is intensely used
- Not suitable for RDD implementation

- Cooperative Groups

  Official API (CUDA)
  - $\geq$  6.1)
- All-SM barrier sync
- Generic instructions for any task
- RDD is fast enough

#### Semaphores

- Manual implementation
- SM-pairwise sync
- Task- and algorithm-specific code
- RDD is even faster

# How to compare sync methods? Tau-model

#### Model parameters:

- N update iter-s
- each iter = M computing op-s and K data sync op-s
- elapsed T(M, N, K) [sec]

Model derivatives:

iteration time 
$$\tau(M, K) = \lim_{\substack{N \to \infty \\ M \to \infty}} \frac{T(N, M, K)}{NM}$$
computation time  $\tau_I = \lim_{\substack{M \to \infty \\ K = const}} \tau(M, K)$ 
synchronization latency  $\lambda_s = \lim_{\substack{K \to \infty \\ M = const}} [\tau(M, K + 1) - \tau(M, K)]$ 
arithmetic intensity  $\iota = \frac{M}{K}$ 
performance  $\pi = \frac{1}{\tau}$ 

# Tau-model testing and fitting

- Model quantifies
   latency λ<sub>s</sub> for CG and
   Sem syncs
- $\tau(\iota)$  fits well to  $\tau = \tau_l + \frac{\lambda_s}{\iota}$ • Empirically, •  $\tau_s = \frac{\lambda_s}{\iota}$ Little's law
  - $\tau = \tau_s + \tau_l$ Additive law •  $\frac{1}{\pi} = \frac{1}{\pi_s} + \frac{1}{\pi_l}$ Harmonic law



# Tau-model testing and fitting

- Model quantifies
   latency λ<sub>s</sub> for CG and
   Sem syncs
- $\tau(\iota)$  fits well to  $\tau = \tau_l + \frac{\lambda_s}{\iota}$ • Empirically, •  $\tau_s = \frac{\lambda_s}{\iota}$ Little's law •  $\tau = \tau_s + \tau_l$ Additive law
  - Additive law  $\frac{1}{\pi} = \frac{1}{\pi_s} + \frac{1}{\pi_l}$ Harmonic law



# Conclusions

**I** RDD for D3Q19 LBM is developed with two sync options available

2 New *tau-model* is developed, linking *overall performance* with *comp performance*, *sync latency*, and *arithmetic intensity* 

$$rac{1}{\pi} = rac{1}{\pi_I} + rac{\iota}{\lambda_s}$$

3 Found out

- RTX 20: Semaphores are 2.0x faster than CG barriers
- RTX 30: Semaphores are 1.2x faster than CG barriers

Thank you for your attention!

Contact: pershin2010@gmail.com