

The logo of Shanghai Jiao Tong University (SJTU) is a circular emblem. It features a red border with the university's name in Chinese characters "上海交通大学" at the top and "SHANGHAITECH UNIVERSITY" at the bottom. In the center, there is a stylized red pagoda-like structure with the year "2013" at its base. The Chinese characters "立志" (Lìzhì) and "报国" (Bàoguó) are on the left, and "成才" (Chéngtái) and "裕民" (Yùmín) are on the right, arranged around the central structure.

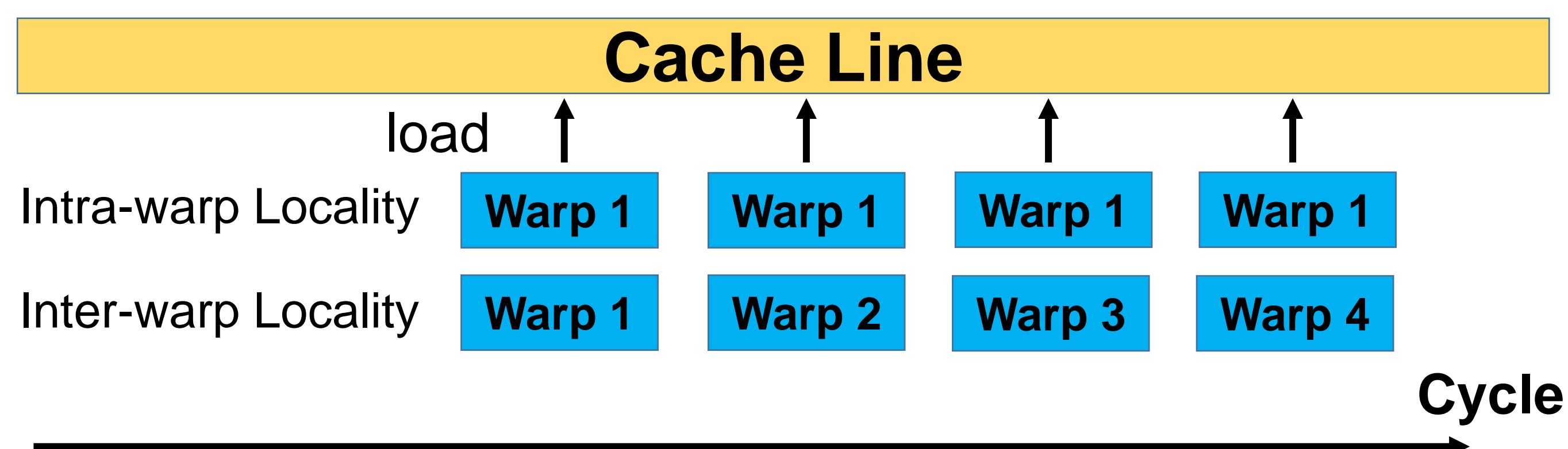
¹School of Information Science and Technology, ShanghaiTech University, Shanghai, China

{huwm1, louxin}@shanghaitech.edu.cn, {peterzhou, rainiequan, rogerwang}@glenfly.com

GLENFLY

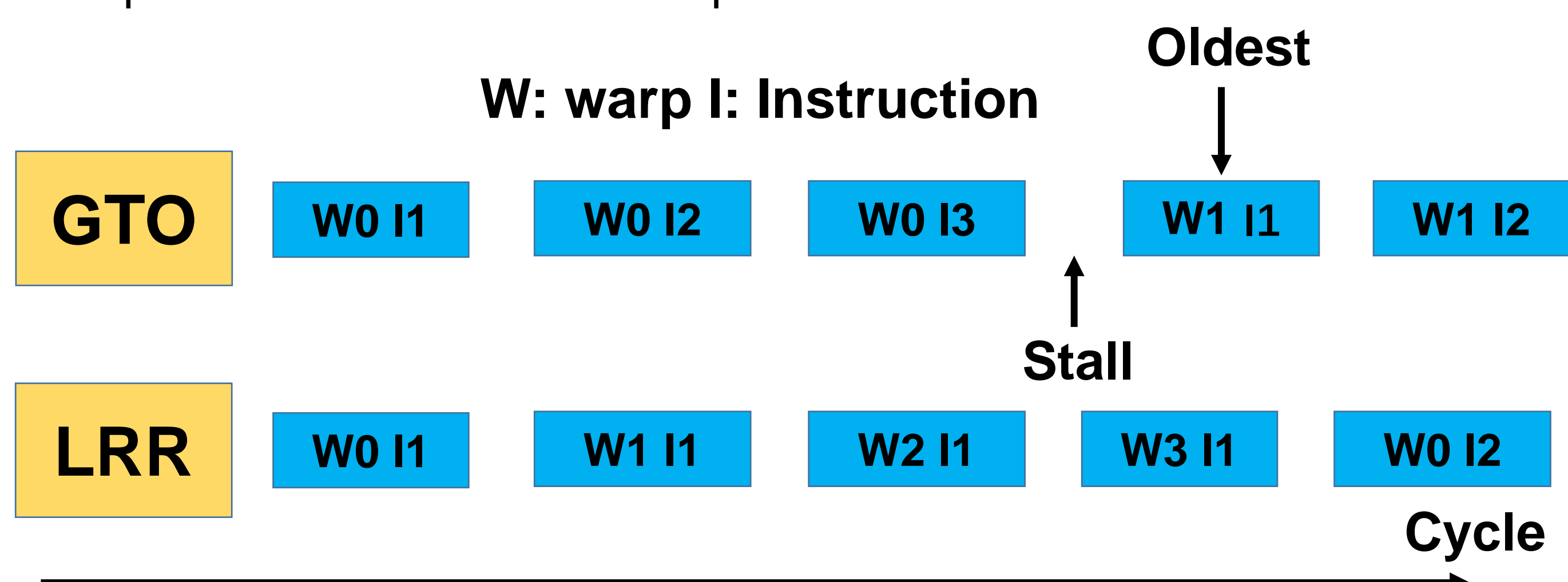
Method

- **Inter-warp locality:** a cache line is accessed by the adjacent warp
- **Intra-warp locality:** a cache line is accessed by the same warp



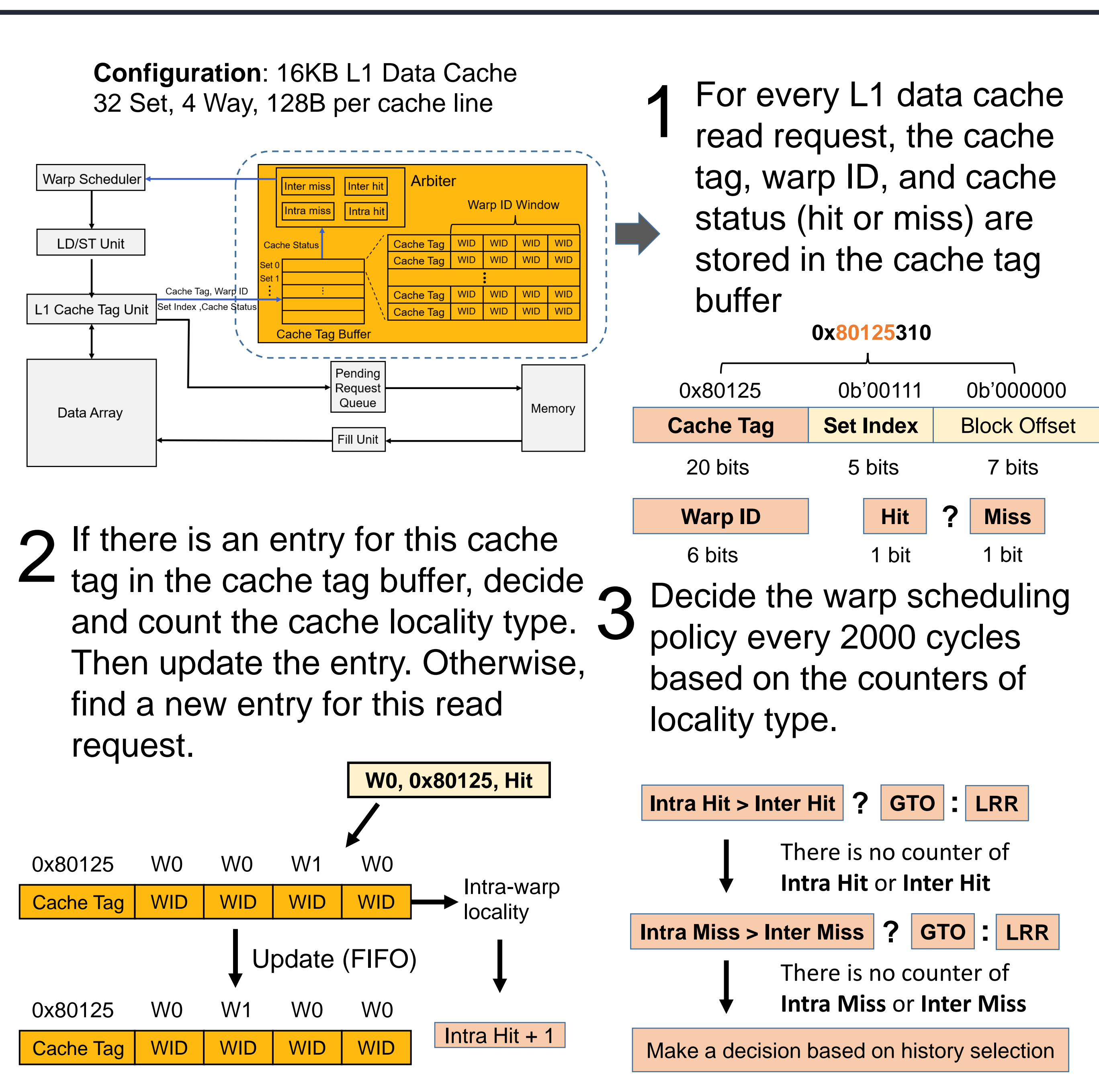
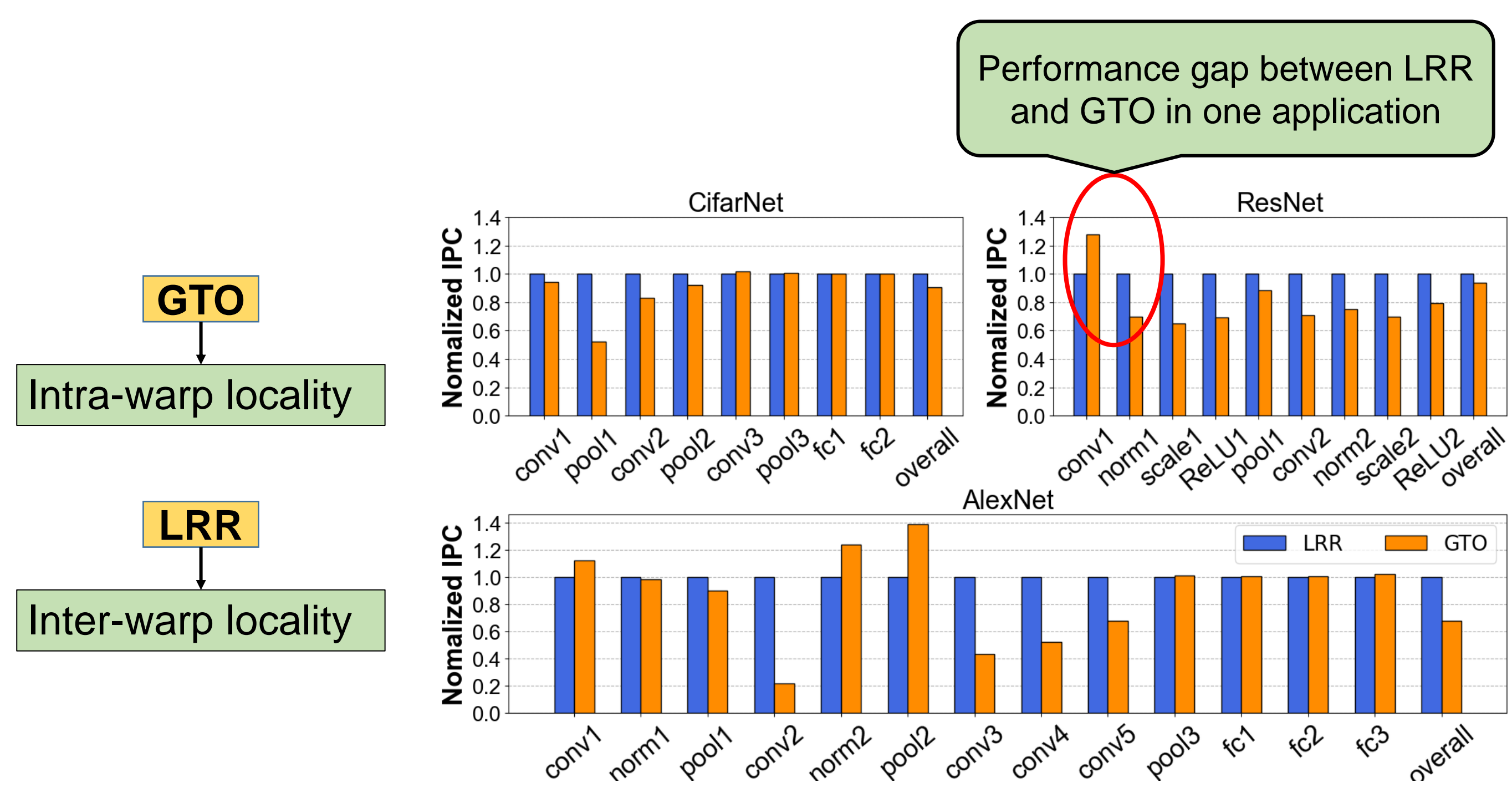
When one warp is stalled, the scheduler will issue another ready warp. Basically, there are two commonly used baseline scheduling policies in typical schedulers:

- **Loosely-Round-Robin (LRR):** equal priority has been given to each warp and switching warps in every cycle
- **Greedy-Then-Oldest (GTO):** prioritizing older warps over younger warps. It does not switch a warp until it is stalled.



Evaluation

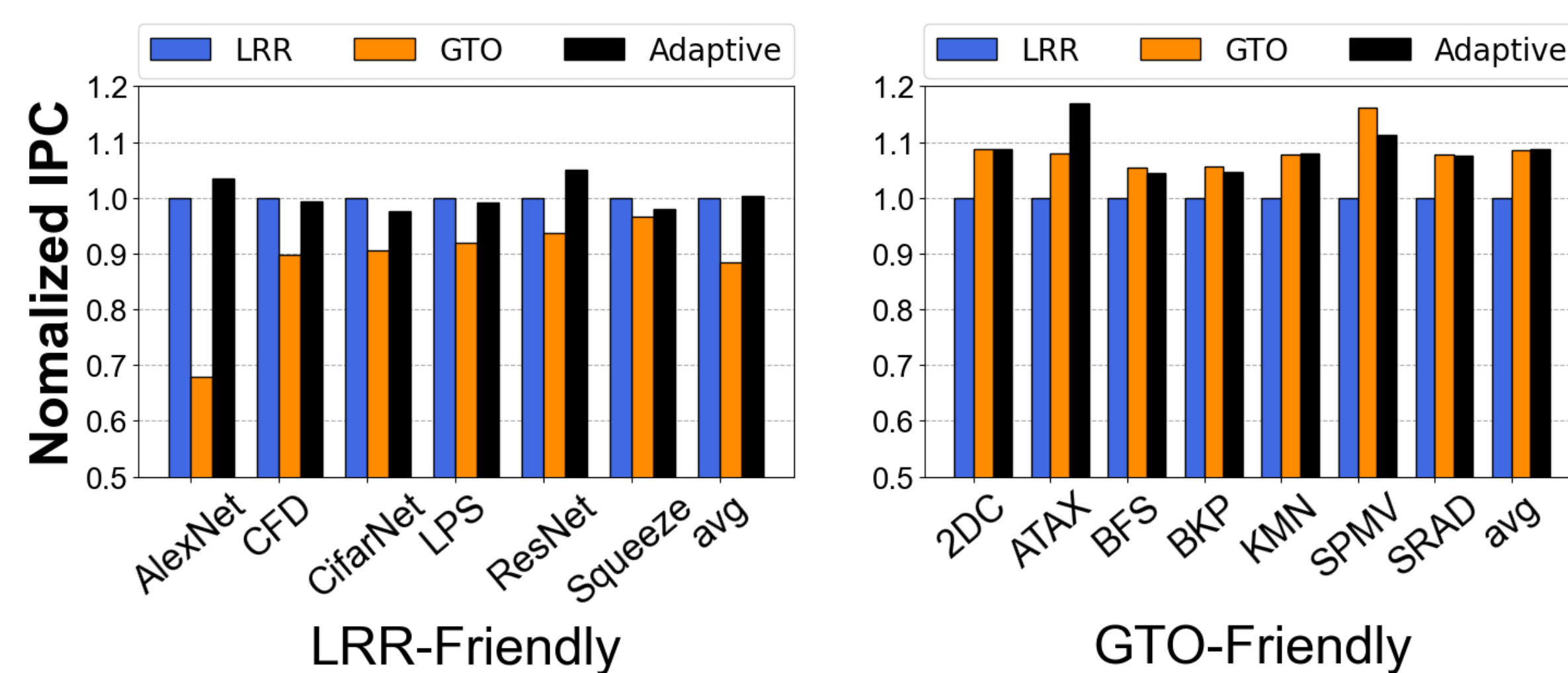
- LRR-based scheduler in general can better exploit inter-warp locality
- GTO-based scheduler would be more suitable for intra-warp locality
- Single warp scheduling policy cannot fit all workloads



Evaluation

For most applications, the proposed scheduling mechanism is able to select the appropriate scheduling policy and approaches the better performance of two baseline schedulers.

Parameter	Value
Number of SMs	15
Max number of Warps per SM	48
Max number of Blocks per SM	8
Number of Schedulers per SM	2
Number of Registers per SM	32768
L1 Data Cache	16KB per SM (32-sets/4-ways)
L1 Inst Cache	2KB per SM (4-sets/4-ways)
L2 Cache	768KB (64-sets/8-ways)
DRAM	974MHZ FR-FCFS scheduler



- we propose a scheduling mechanism that analyzes the cache locality type and adaptively selects the warp scheduler between GTO and LRR at runtime.
- Evaluation results show that the proposed scheduling mechanism can select the better scheduler between LRR and GTO in most cases.