

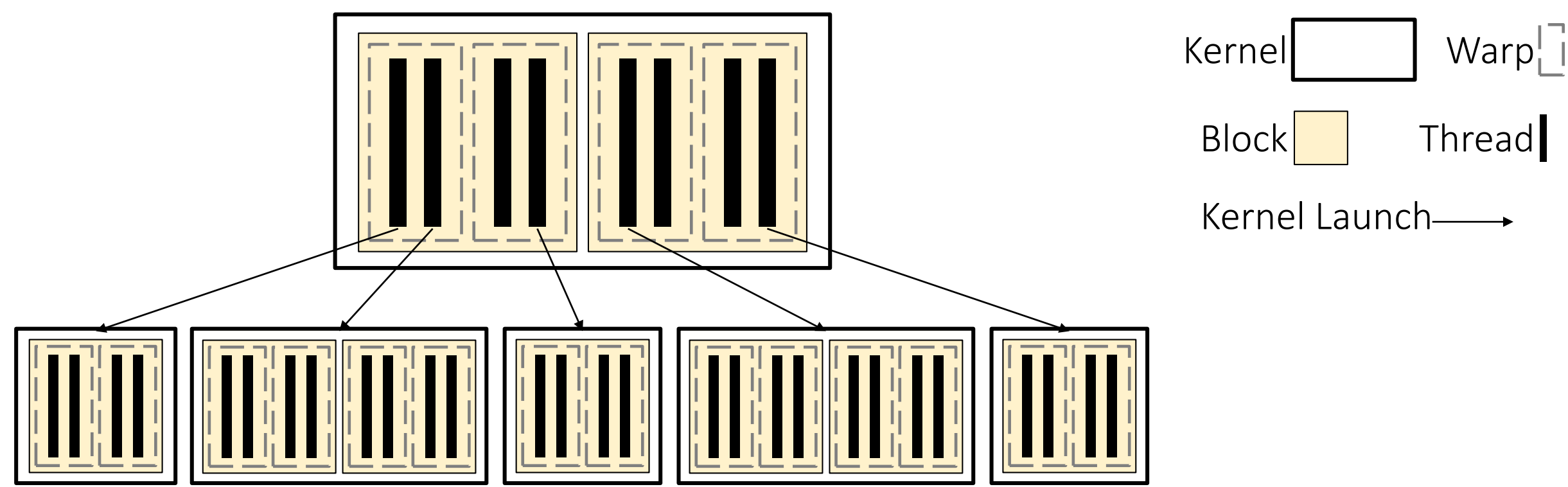
KLAP: Kernel Launch Aggregation and Promotion for Optimizing Dynamic Parallelism

Izzat El Hajj¹ (Illinois), Juan Gómez-Luna² (Córdoba), Cheng Li (Illinois), Li-Wen Chang (Illinois), Dejan Milojcic (HPE), Wen-mei Hwu (Illinois)
¹elhajj2@illinois.edu, ²el1goluj@uco.es

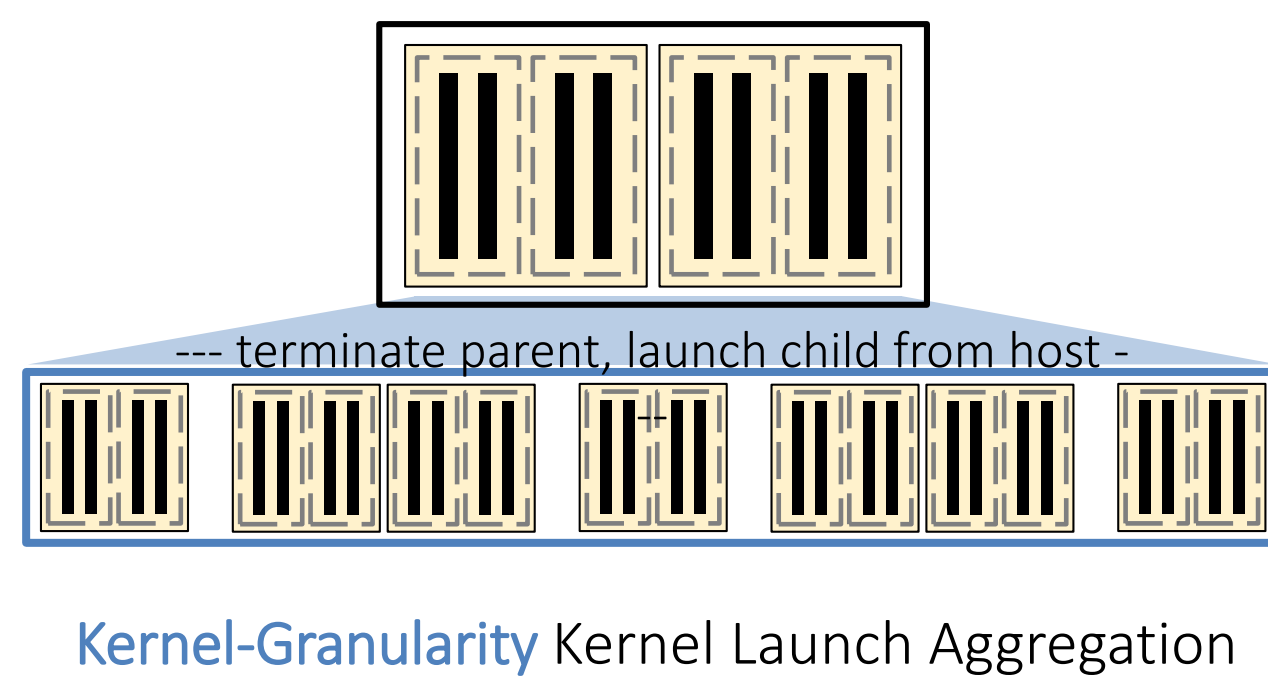
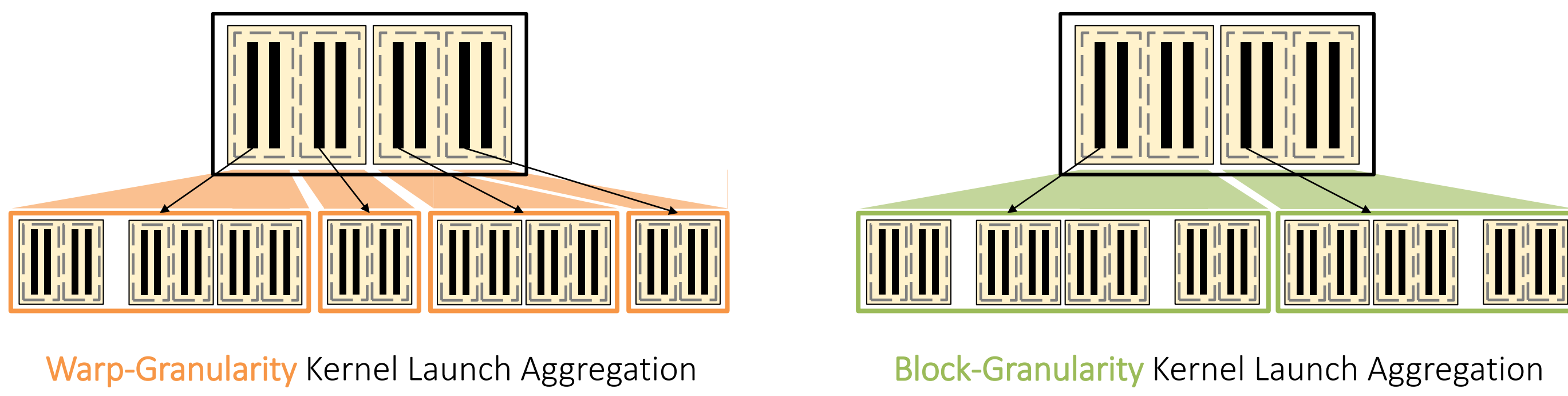
Kernel Launch Aggregation

Problem: Launching too many fine-grain kernels

- Large kernel count incurs too much launch overhead
- Fine granularity of kernels underutilizes GPU resources



Solution: One thread launches a kernel on behalf of a group



Code Transformation:

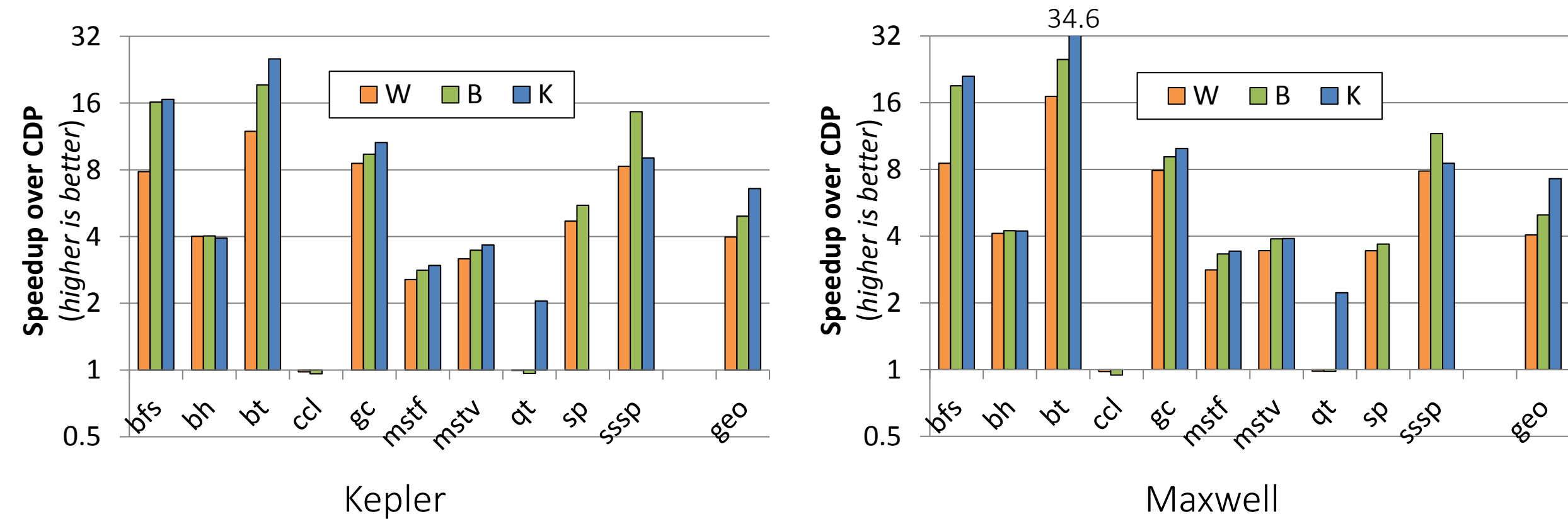
```
kernel<<<gD, bD>>>(arg1, arg2, arg3)    _global__ void kernel(params) {
                                        kernel body
}
                                        }

allocate arrays for args, gD, and bD    _global__ void kernel_agg(param arrays, gD array, bD array) {
store args in arg arrays                calculate index of parent thread
store gD in gD array, and bD in bD array load params from param arrays
new gD = sum of gD array across warp/block load actual gridDim/blockDim from gD/bD arrays
new bD = max of bD array across warp/block calculate actual blockIdx
if(threadIdx == launcher thread in warp/block) {
kernel_agg<<<new gD, new bD>>>         kernel body (with kernel launches transformed and with
                                        (arg arrays, gD array, bD array)         using actual gridDim/blockDim/blockIdx)
}
                                        }
```

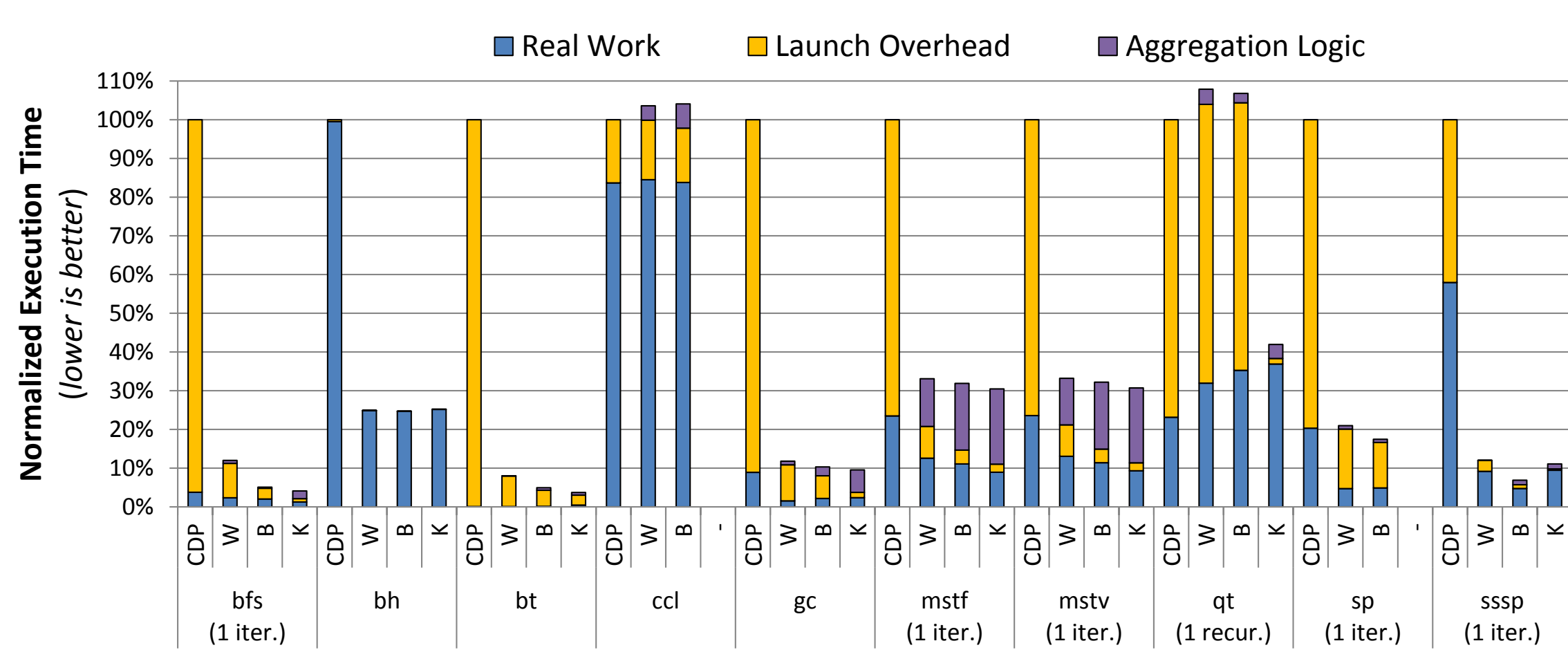
Original Kernel Call Original Kernel

Transformed Kernel Call (block-granularity aggregation example) Transformed Kernel (block-granularity aggregation example)

Results:



Increasing aggregation granularity improves performance (geomean speedup of 6.58x for K-aggregation on Kepler)

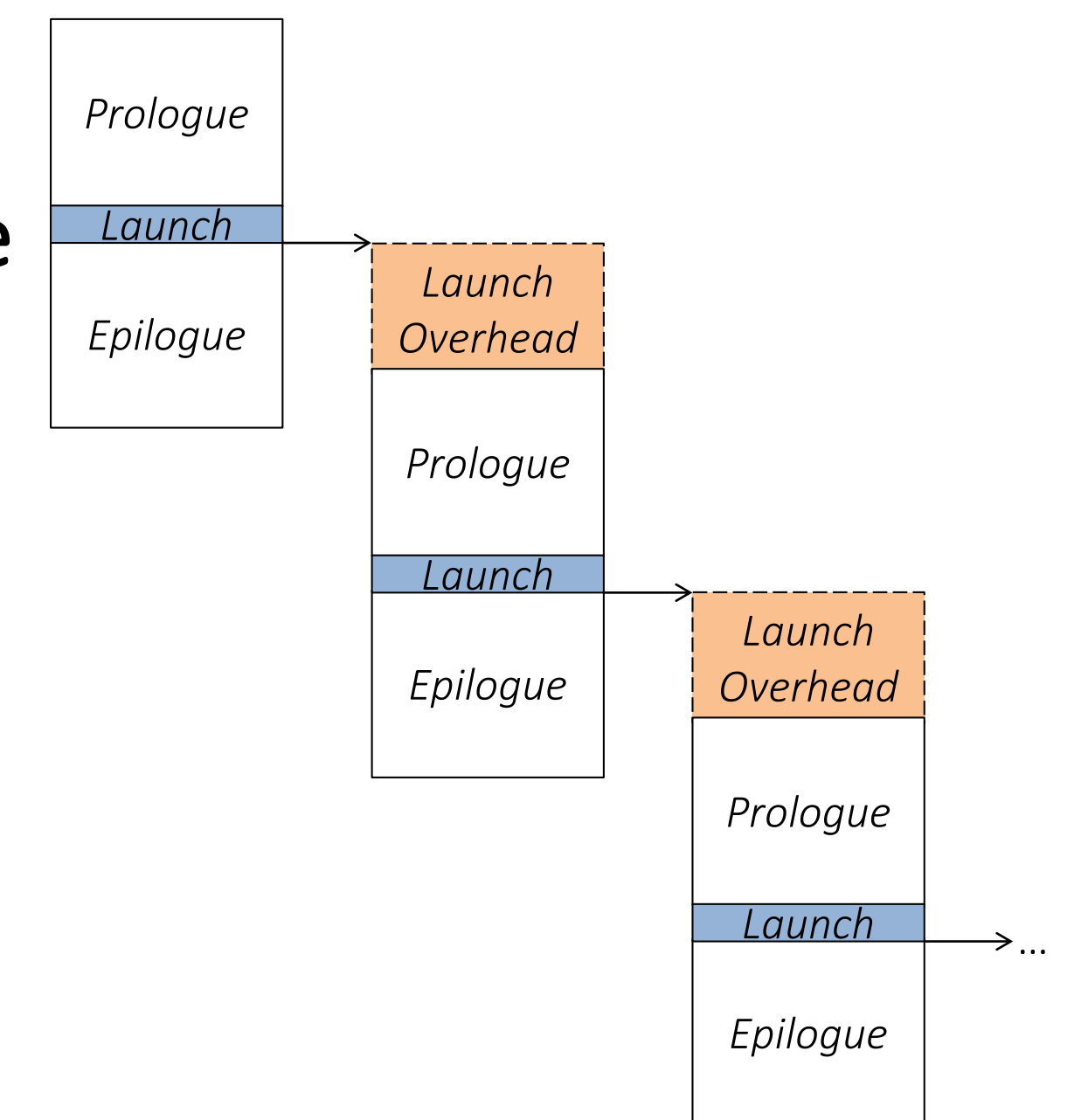


Performance improvement comes from reduced launch overhead and better resource utilization

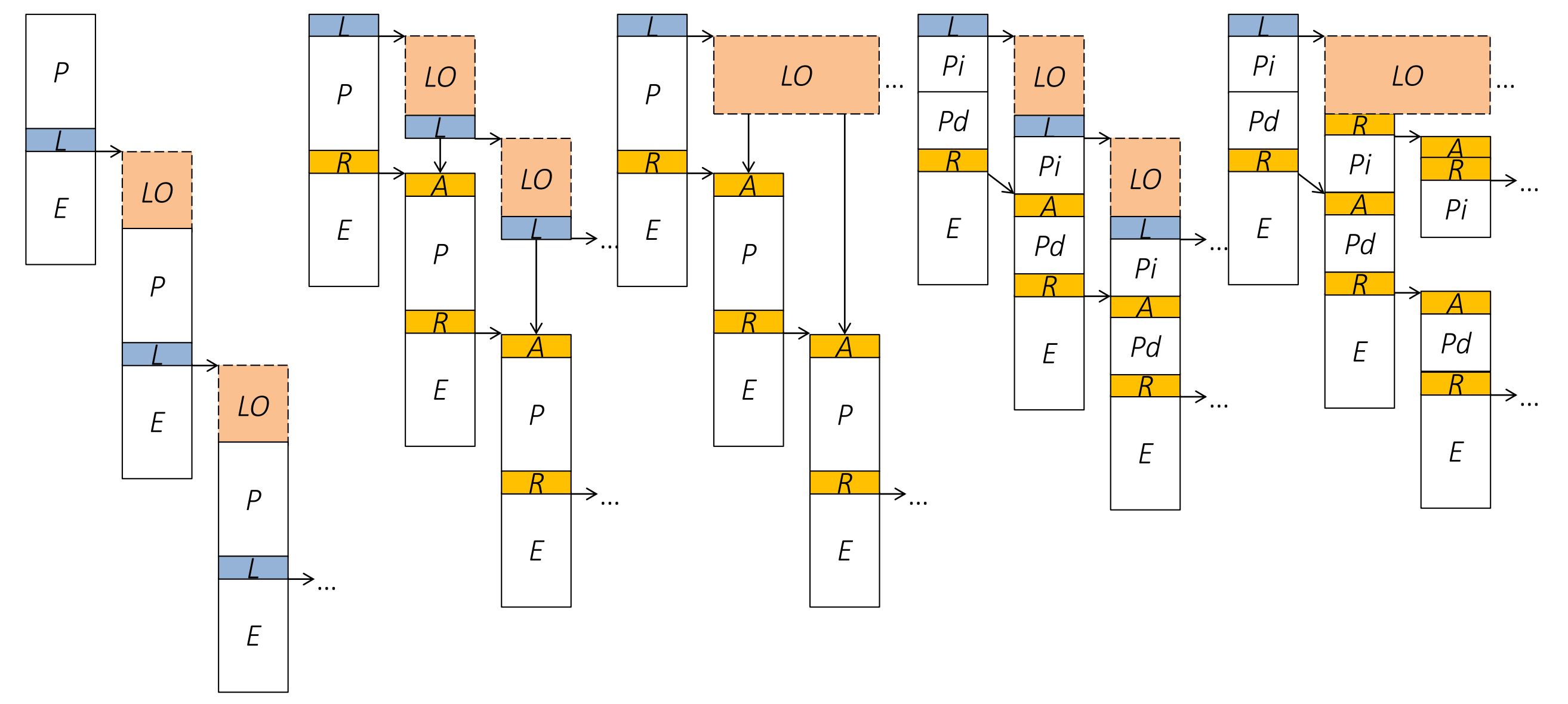
Kernel Launch Promotion

Problem: Deep call stacks

- Launch overhead dominates the critical path
- Call stack has limited depth



Solution: One thread launches a kernel on behalf of a group



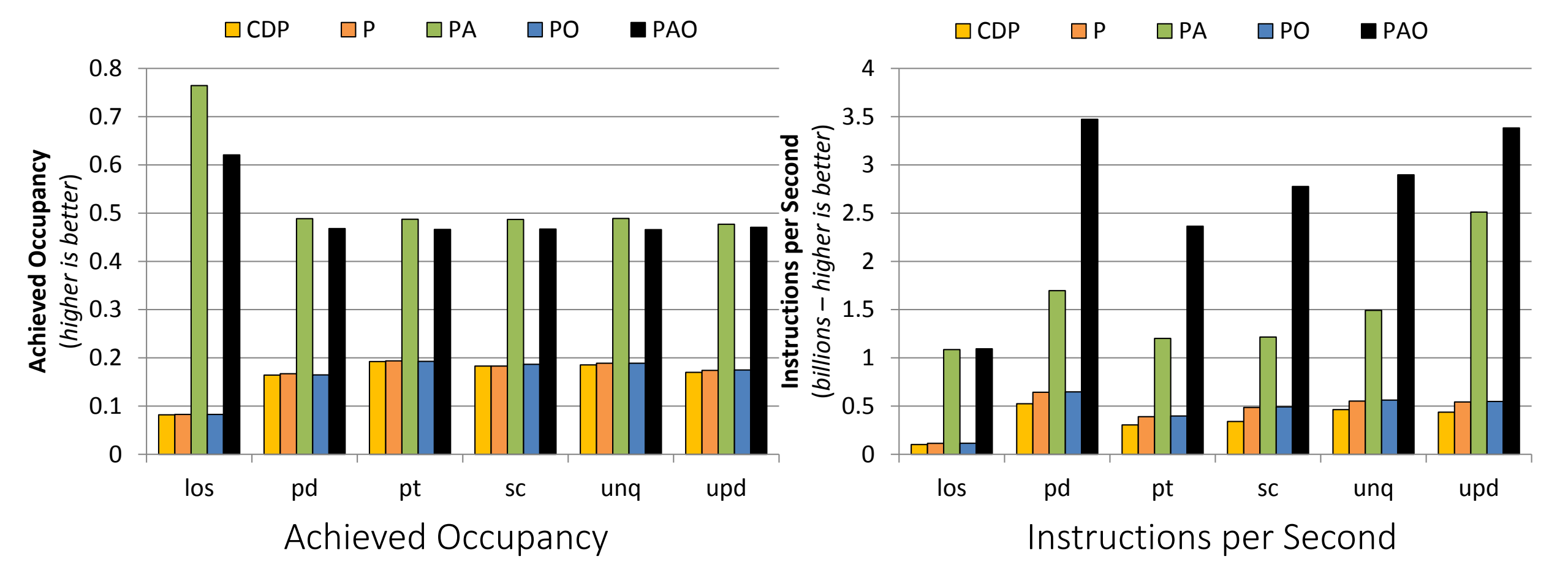
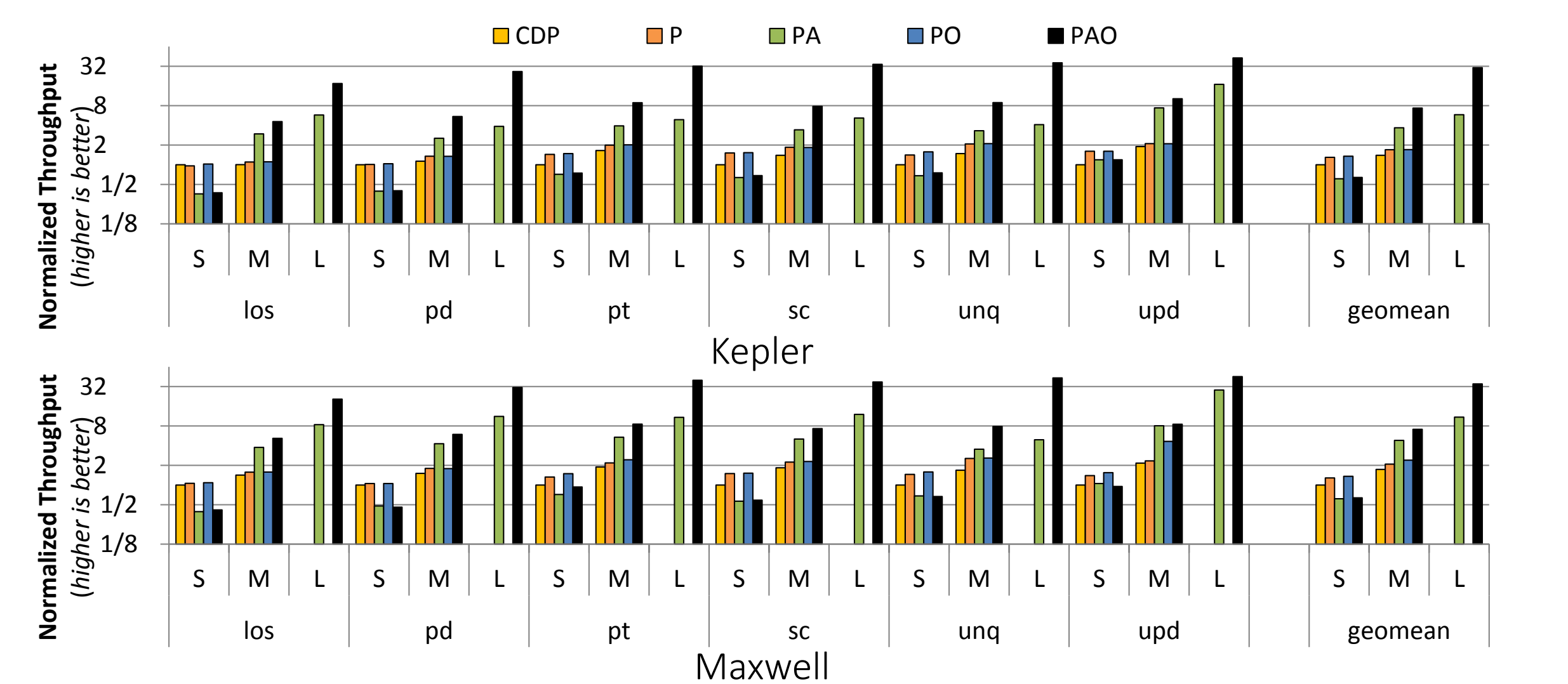
Code Transformation:

```
_global__ void kernel(params_avail, params_post) {
    prologue
    if(launcher thread) {
        kernel<<<1, nThreads>>>(args_avail, args_post)
    }
    epilogue
}
                                        }
                                        }

if(launcher thread) {
    allocate postponed arg buffers
    allocate child flag
    kernel_from_kernel<<<1, nThreads>>>
        (args_avail, postponed arg buffers, child flag)
}
wait to acquire flag
load params_post from postponed param buffers
prologue
if(launcher thread) {
    store args_post in postponed arg buffers
    memory fence
    set child flag to release child
}
epilogue
}
                                        }
                                        }
```

Original Kernel Transformed Kernel (called from host) Transformed Kernel (called from kernel)

Results:



Aggregation improves occupancy due to fewer coarser-grain kernels
 Overlap improves instructions per second due to more work parallel work available