

Blk-Tune

Blocking Parameter Auto-Tuning for Flash-based Out-of-Core Stencil Computations

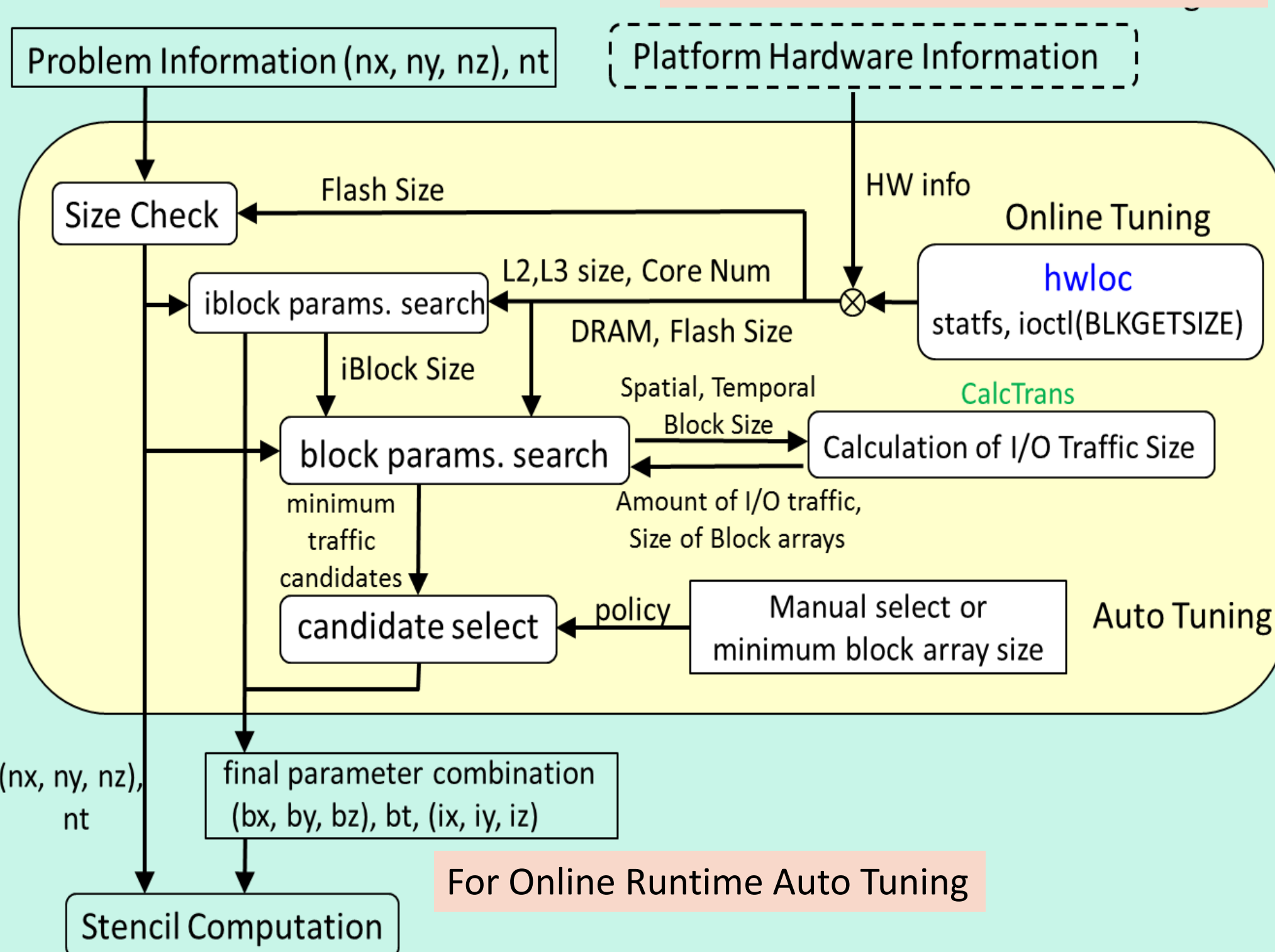
Hiroko Midorikawa JST CREST & Seikei University, Tokyo Japan
midori@st.seikei.ac.jp <http://www.ci.seikei.ac.jp/midori/paper>

Abstract

Blk-Tune is a runtime-based blocking parameter auto-tuning system that enables the use of flash memory as an extension of main memory. It obtains platform hardware information in runtime by using Portable Hardware Locality (hwloc) [3], which is widely portable to various kinds of OS and CPU architectures. By calculating the total amount of IO traffic to a flash device, it selects an optimal combination of spatial and temporal blocking sizes to suit the capacity of each memory layer (flash, DRAM, level-3 cache, and level-2 cache). A selected parameter combination minimizes the amount of data transferred between the flash device and DRAM, which is the dominant factor affecting the performance of out-of-core algorithms using flash. **Blk-Tune** allows users to easily achieve maximum performance of large-scale stencil computations for particular platforms and application settings.

Overview of Blk-Tune

Auto Tuning System Diagram



Blk-Tune: Auto-tuning system for out-of-core stencil computations using flash, DRAM, and cache

RunTime Parameter Tuning as a frontend

- User command parameters**
domain size(nx,ny,nz), Time step(nt), Flash device path
`./stencil7p -n 4094 4096 2048 -t 1000 -d /dev/sdc`
---- autotune start ----
(bx,by,bz), bt = (4094,1024,512), 125
(ix,iy,iz) = (4094,1,40) : 20961280 B (19.99 MiB)
---- autotune end ----
- Get hardware information**
device capacity, device block size, DRAM size, L3 cache size, # of Cores, # of CPU sockets, Flash info.
- Calculate optimal blocking sizes by adjusting to underlying hardware**
 - Efficient spatial block size and shape, temporal block size (optimal bx,by,bz, bt for Blocks on DRAM)
 - Efficient spatial inner block shape and size (optimal ix,iy,iz for inner blocks on L2 and L3 cache)
- Numa-aware computing**
 - Blocks are divided into sub-block areas according to # of CPU sockets
 - memory-bind & cpu-bind between each sub-block to each local socket
 - locality-aware sub-block computation on each CPU socket using local cores

Search Problem Definition

- Search of a Globally Optimal set to Minimize the IO traffic between Flash and DRAM
- Inputs to the auto-tuning system**
- Hardware information (from hwloc)**
Capacity of each memory layer
Sc2: level-2 cache, Sc3: level-3 cache, Sm: Main memory (DRAM) capacity,
Sf: Flash memory capacity, Sb: Flash device block size
CPU configurations
Nc: # of cores in a system, Ns: # of sockets
 - Problem information (User input)**
Domain data size & Number of iterations for stencil computations
nx, ny, nz: 3D array, nt: (time steps)

Outputs from the auto-tuning system

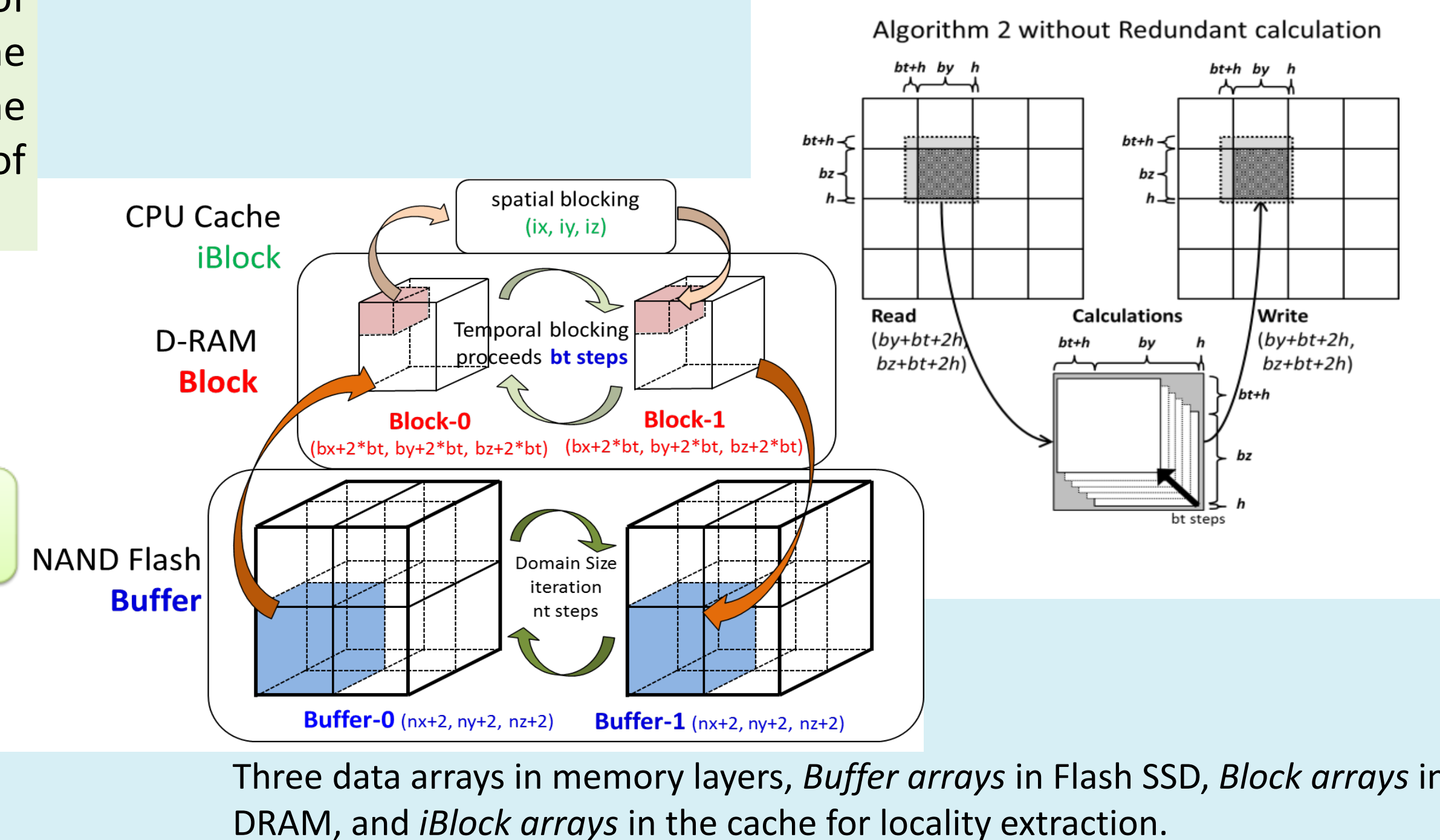
- An optimal combination of temporal and spatial block size**
bx, by, bz: spatial blocking size (Block Array size), bt: temporal size
It minimize IO traffic between DRAM and flash memory the volume of the Block array can be accommodated in DRAM capacity.
- Semi-optimal internal spatial block size combination**
ix, iy, iz: spatial blocking size (iBlock array)

Premises

$Sf > nx * ny * nz * k * Esz$
(k is constant, 2 for double buffer arrays. Esz is the data element size in byte)
nx is required to be a multiple of Sb. (tune for asynchronous IO)

Out-of-core Stencil Algorithm using Flash SSDs

Locality-aware algorithm:
temporal and spatial blocking for memory layers



Three data arrays in memory layers, Buffer arrays in Flash SSD, Block arrays in DRAM, and iBlock arrays in the cache for locality extraction.

Search procedures

Step1 find optimal iBlock size (ix, iy, iz)

- Objective Function**
 $\min (| k * iBsz(ix, iy, iz) * Esz - Sc2 * Nc |)$ (1a)
where $iBsz(ix, iy, iz) = ix * iy * iz$.
- Constraint conditions**
 - iz must be a multiple of Nc.
 - ix has to equal nx.

Step2 find the optimal Block size (bx, by, bz) and bt

- Objective function**
 $\min(Total_A) = \min \left(\sum_i A_i \right) = \min \left(\sum_{X,Y,Z} B(X,Y,Z,i) \right)$ (2a)
 $\sum_i A_i = \sum_{X,Y,Z} (B(X,Y,Z,1) * T_f + B(X,Y,Z,T_c) * (T_c - T_f))$ (2d)

where $1 \leq i \leq T_c$, $1 \leq X \leq ny$ and X is a multiple of bx ,
 $1 \leq Y \leq ny$ and Y is a multiple of by , $1 \leq Z \leq nz$ and Z is a multiple of bz

Constraint conditions

- bz must be multiple of iz
- by must be multiple of iy
- bx must equal nx, bx = nx
- bt must satisfy the following
 $2 \leq bt \leq nt$, $bt \leq \frac{bx}{2}$, $bt \leq \frac{by}{2}$, and $bt \leq \frac{bz}{2}$
bt = 2: temporal blocking for minimum size
bt must satisfy the following
Total size of block arrays must not exceed DRAM capacity (Sm)
 $k * Bsz_max * Esz < Sm - Km$
Bsz_max: The largest Block size, k: constant, Km: constant memory capacity reserved for OS
 $Bsz_max = \max_{X,Y,Z} (\max (Bsz_w(X,Y,Z,1), Bsz_r(X,Y,Z,1)))$

Amount of IO traffic for a Block: B(X, Y, Z, i)

- Typical block array size for algorithm-1 with redundant calculations

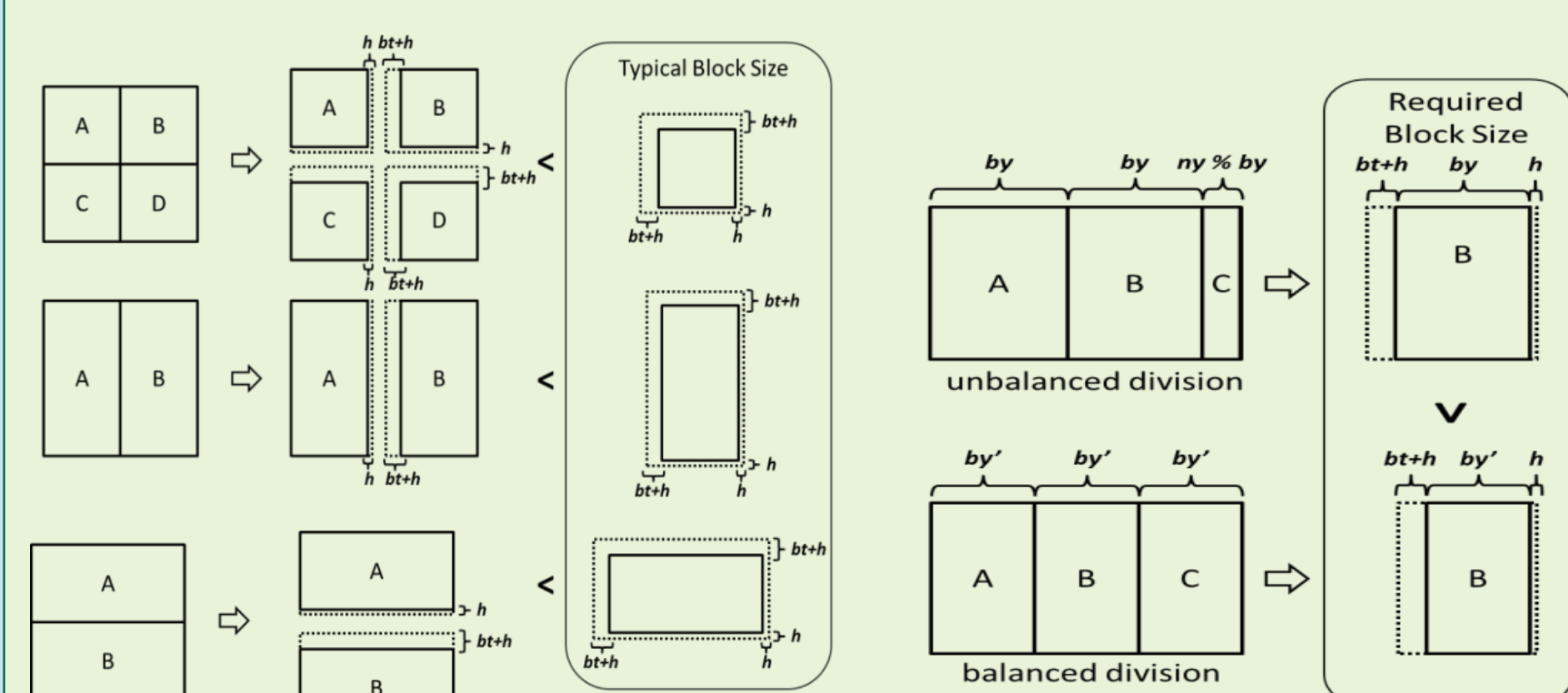
$$Bsz1_r = (bx + 2 * h + 2 * bt) * (by + 2 * h + 2 * bt) * (bz + 2 * h + 2 * bt)$$

$$Bsz1_w = (bx) * (by) * (bz)$$

- Typical block array size for algorithm-2 without redundant calculations

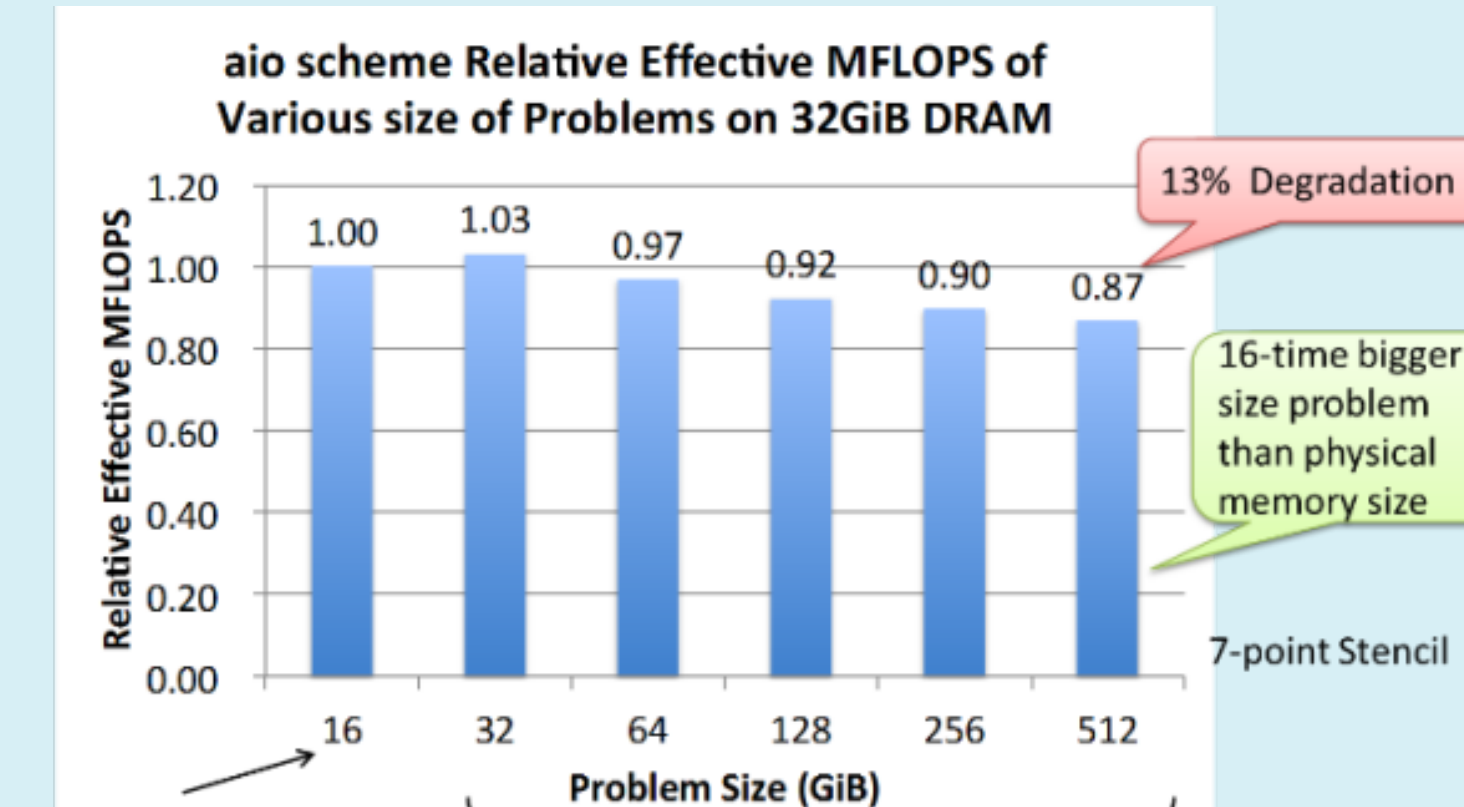
$$Bsz2_r = Bsz2_w = (bx + 2 * h + bt) * (by + 2 * h + bt) * (bz + 2 * h + bt)$$

Division scheme and boundary effects in block size

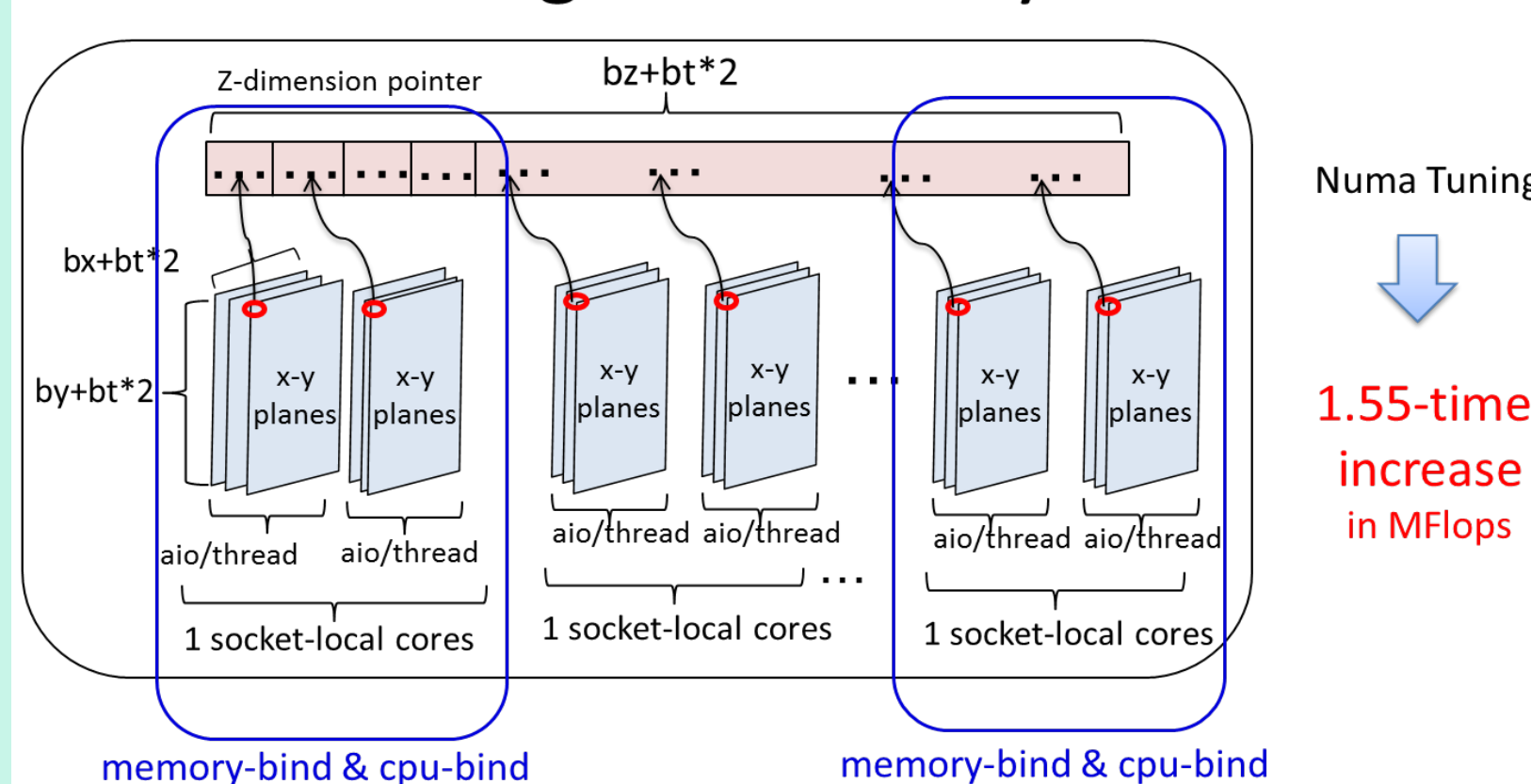


REFERENCES

- H. Midorikawa, et al., "An Evaluation of the Potential of Flash SSD as Large and Slow Memory for Stencil Computations," Proc. IEEE Int. Conf. on High Performance Computing and Simulation IEEE-HPCS2014, pp. 268-277
- H. Midorikawa et al., "Locality-Aware Stencil Computations using Flash SSDs as Main Memory Extension," Proc. IEEE/ACM Int. Symp. on Cluster, Cloud and Grid Computing CCGrid2015, pp. 1163-1168
- Portable Hardware Locality (hwloc), <https://www.open-mpi.org/projects/hwloc/>
- H. Midorikawa: "Blk-Tune: Blocking Parameter Auto-Tuning to Minimize Input-Output Traffic for Flash-based Out-of-Core Stencil Computations", IEEE Int. Parallel and Distributed Process. Symp (IPDPS) Workshop IWAPT , May 2016



Tuning for Numa systems



Blk-Tune searches an optimal spatial & temporal blocking parameter set for particular platforms and application settings

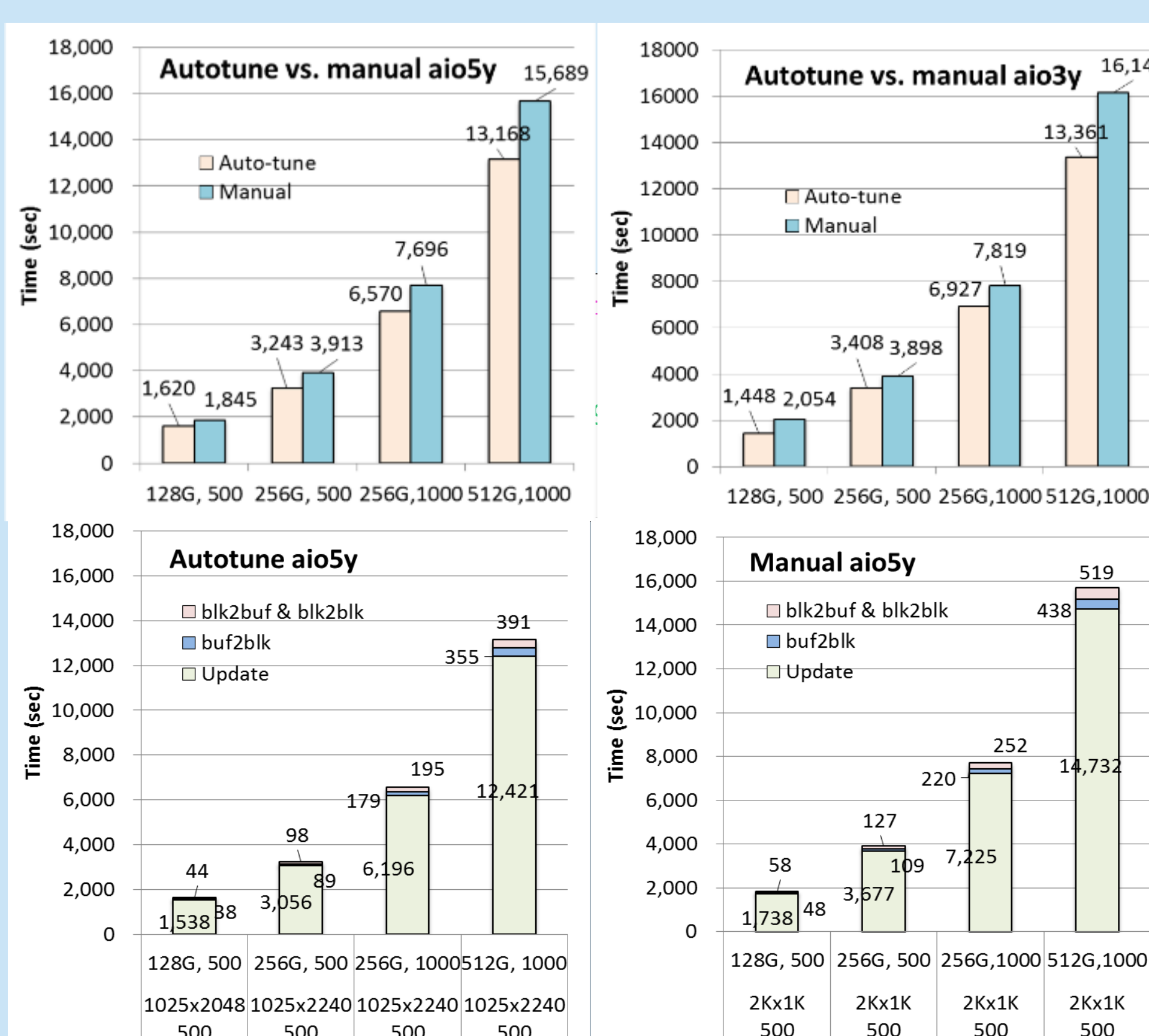
Table I. Parameter examples selected by Blk-Tune

Problem size	nx	ny	nz	nt		
128GiB-problem	2048	2048	2048	500		
Server (DRAM, cores)	bx, ix	by	bz	bt	iy	iz
crest6 (128GiB, 20)	2048	1025	2048	500	1	160
crest4 (64GiB, 16)	2048	1025	1280	500	1	128
crest0 (32GiB, 8)	2048	513	1152	250	1	64

Problem size	nx	ny	nz	nt		
256GiB-problem	2048	2048	4096	1000		
Server (DRAM, cores)	bx, ix	by	bz	bt	iy	iz
crest6 (128GiB, 20)	2048	1025	2240	500	1	160
crest4 (64GiB, 16)	2048	683	1536	334	1	128
crest0 (32GiB, 8)	2048	513	832	250	1	64

Table II. Platforms

server	L1 cache (KiB)	L2 cache (KiB)	L3 cache (MiB)	Phys Mem (GiB)	Flash Mem (TiB)	CPU Xeon E5, (GHz)	Total cores	socket	cores/socket
crest0	32	256	20	32	1.2	2650, (2)	8	1	8
crest4	32	256	20	64	0.785	2687W, (3.1)	16	2	8
crest6	32	256	25	128	1	2687W v3,(3.1)	20	2	10



Manual selection vs. Blk-Tune for various-size problems on crest6 server for two algorithms, aio3y and aio5y, temporal blocking without redundant calculations.