

Tuning of BLAS level 1 and 2

Hans Henrik Brandenburg Sørensen
Section for Scientific Computing
DTU Informatics

BLAS (Basic Linear Algebra Subprograms)

- Basic routines for numerical applications.
- Legacy software package 1979-2002 (Netlib.org).
 - C. L. Lawson, R. J. Hanson, D. Kincaid, and F. T. Krogh, *Basic Linear Algebra Subprograms for FORTRAN usage*, 1979
 - J. J. Dongarra, J. Du Croz, S. Hammarling, and R. J. Hanson, *An extended set of FORTRAN Basic Linear Algebra Subprograms*, 1988
 - J. J. Dongarra, J. Du Croz, I. S. Duff, and S. Hammarling, *A set of Level 3 Basic Linear Algebra Subprograms*, 1990
 - L. S. Blackford, J. Demmel, J. Dongarra, I. Duff, S. Hammarling, G. Henry, M. Heroux, L. Kaufman, A. Lumsdaine, A. Petitet, R. Pozo, K. Remington, R. C. Whaley, *An Updated Set of Basic Linear Algebra Subprograms (BLAS)*, 2002

BLAS (Basic Linear Algebra Subprograms)



- Level 1 BLAS: (xAXPY, xAMAX, xDOT, xNRM2, etc.)
 - Vector of length N ($4 \times N$ bytes)
 - $4 \times N$ bytes : $O(N)$ flops
 - **Memory bound**

BLAS (Basic Linear Algebra Subprograms)

- Level 1 BLAS: (xAXPY, xAMAX, xDOT, xNRM2, etc.)
 - Vector of length N ($4 \times N$ bytes)
 - $4 \times N$ bytes : $O(N)$ flops
 - **Memory bound**
- Level 2 BLAS: (xGEMV, xSYMV, xHEMV, xTRSV, etc.)
 - Matrix of size $N \times N$ ($4 \times N \times N$ bytes) + Vector ($4 \times N$ bytes)
 - $4 \times N^2 + 4 \times N$ bytes : $O(N^2)$ flops
 - **Memory bound**

BLAS (Basic Linear Algebra Subprograms)

- Level 1 BLAS: (xAXPY, xAMAX, xDOT, xNRM2, etc.)
 - Vector of length N ($4 \times N$ bytes)
 - $4 \times N$ bytes : $O(N)$ flops
 - **Memory bound**
- Level 2 BLAS: (xGEMV, xSYMV, xHEMV, xTRSV, etc.)
 - Matrix of size $N \times N$ ($4 \times N \times N$ bytes) + Vector ($4 \times N$ bytes)
 - $4 \times N^2 + 4 \times N$ bytes : $O(N^2)$ flops
 - **Memory bound**
- Level 3 BLAS: (xGEMM, xSYMM, xHEMM, xTRSM, etc.)
 - 1 or 2 matrices of size $N \times N$ ($4 \times N \times N$ bytes)
 - $(2 \times) 4 \times N^2$ bytes : $O(N^3)$ flops
 - **Compute bound** for large N

BLAS (Basic Linear Algebra Subprograms)

- Level 1 BLAS: (xAXPY, xAMAX, xDOT, xNRM2, etc.)
 - Vector of length N ($4 \times N$ bytes)
 - $4 \times N$ bytes : $O(N)$ flops
 - Memory bound
- Level 2 BLAS: (xGEMV, xSYMV, xHEMV, xTRSV, etc.)
 - Matrix of size $N \times N$ ($4 \times N \times N$ bytes) + Vector ($4 \times N$ bytes)
 - $4 \times N^2 + 4 \times N$ bytes : $O(N^2)$ flops
 - Memory bound
- Level 3 BLAS: (xGEMM, xSYMM, xHEMM, xTRSM, etc.)
 - 1 or 2 m
 - $(2 \times) 4 \times N$
 - Compute bound for large N

Whenever possible, use Level 3
BLAS in your (GPU) applications!

BLAS (Basic Linear Algebra Subprograms)

- Level 1 BLAS: (xAXPY, xAMAX, xDOT, xNRM2, etc.)

- Vector of length N ($4 \times N$ bytes)

- $4 \times N$ bytes

- Memory bound

Why consider Level 1 and Level 2:

- Matvecs, orthogonalizations, norms, triangular solves, LAPACK building blocks, e.g., factorizations.

- Level 2 BLAS: (xGEMV, xGBMV, etc.)

- Matrix

- $4 \times N^2$ bytes

- Memory bound

- Still very little attention from GPU community.

- $4 \times N$ bytes

- Level 3 BLAS: (xGEMM, xSYMM, xHEMM, xTRSM, etc.)

- 1 or 2 matrices of size $N \times N$ ($4 \times N \times N$ bytes)

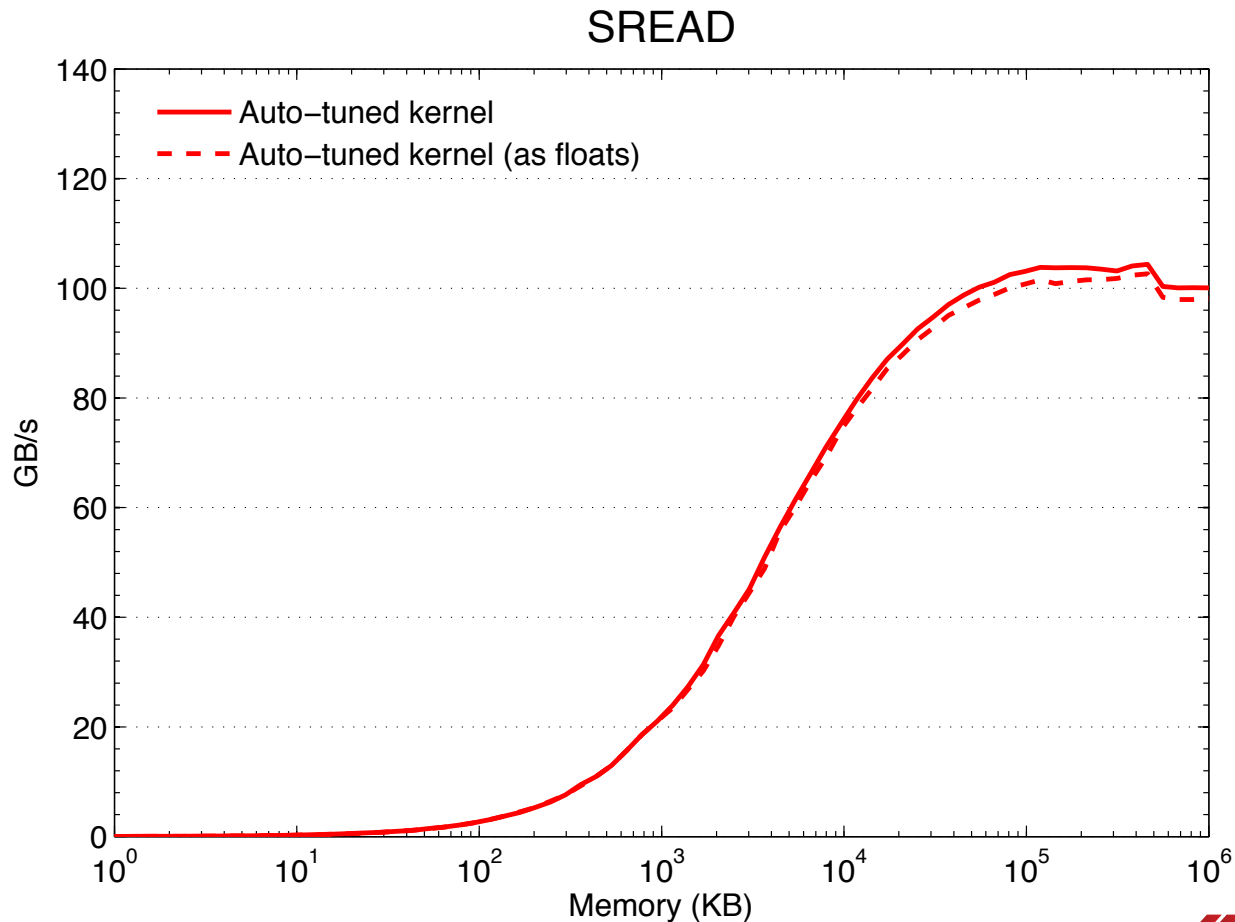
- $(2 \times) 4 \times N^2$ bytes : $O(N^3)$ flops

- Compute bound for large N

Performance Prediction

C2050: Theoretical bandwidth 144 GB/s

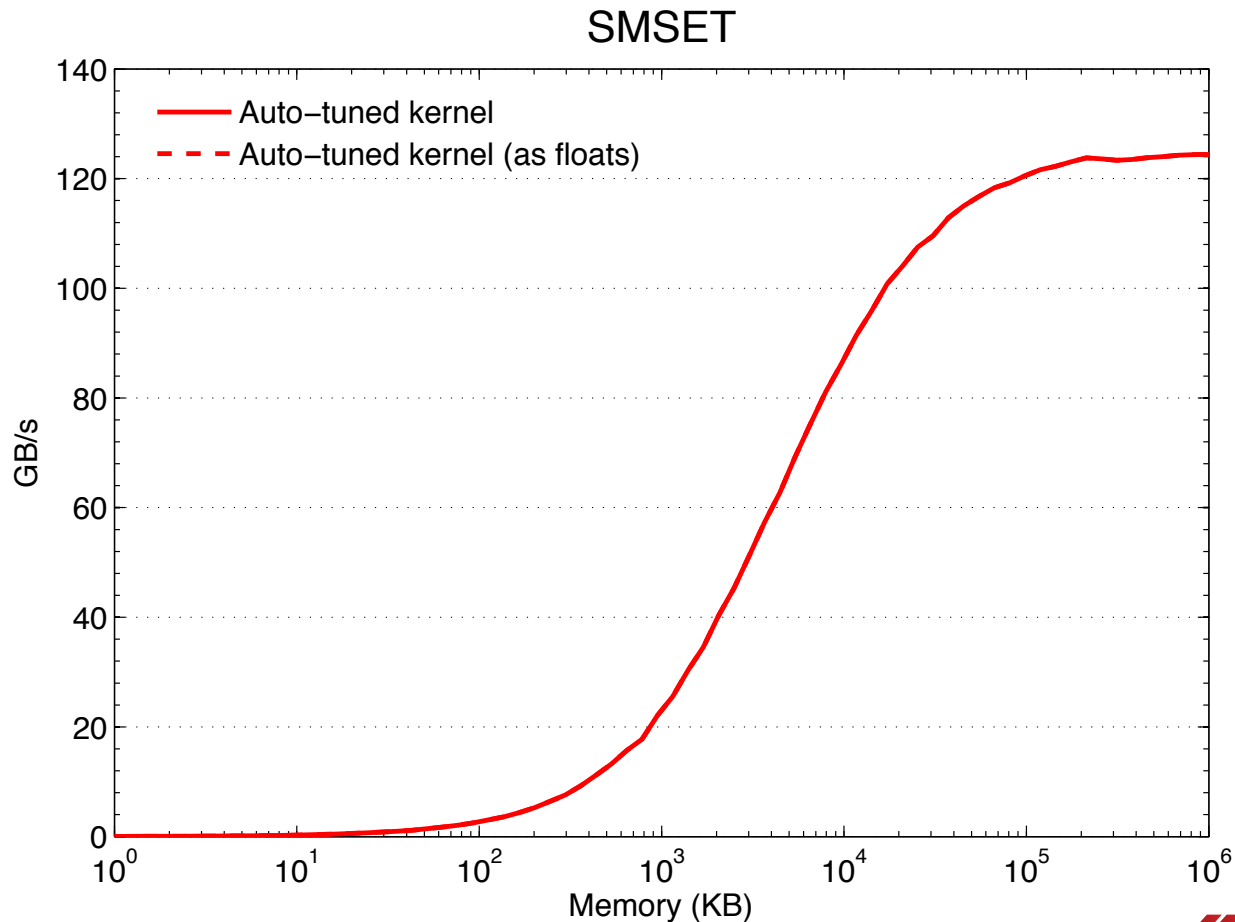
What is the effective bandwidth?

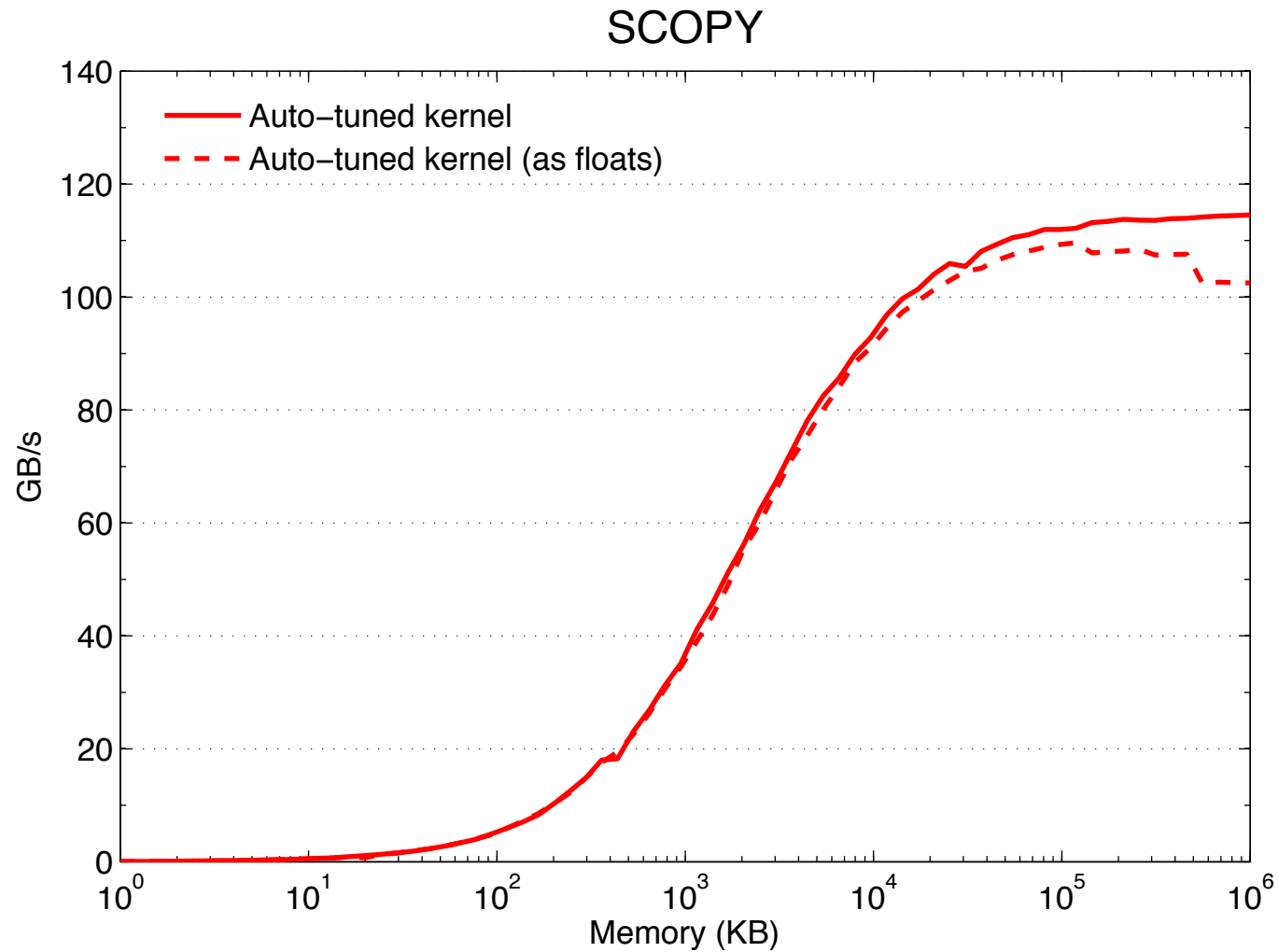


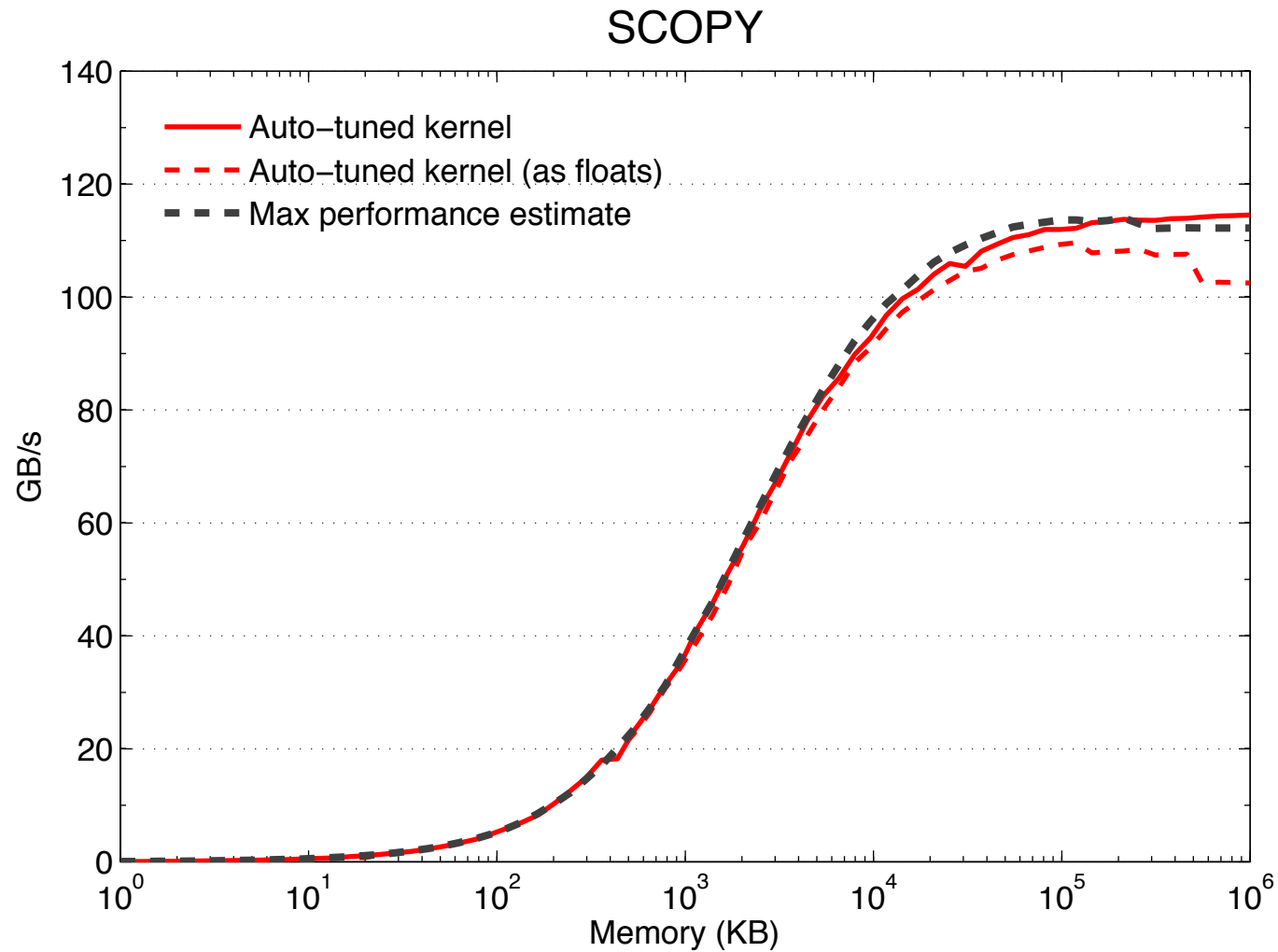
Performance Prediction

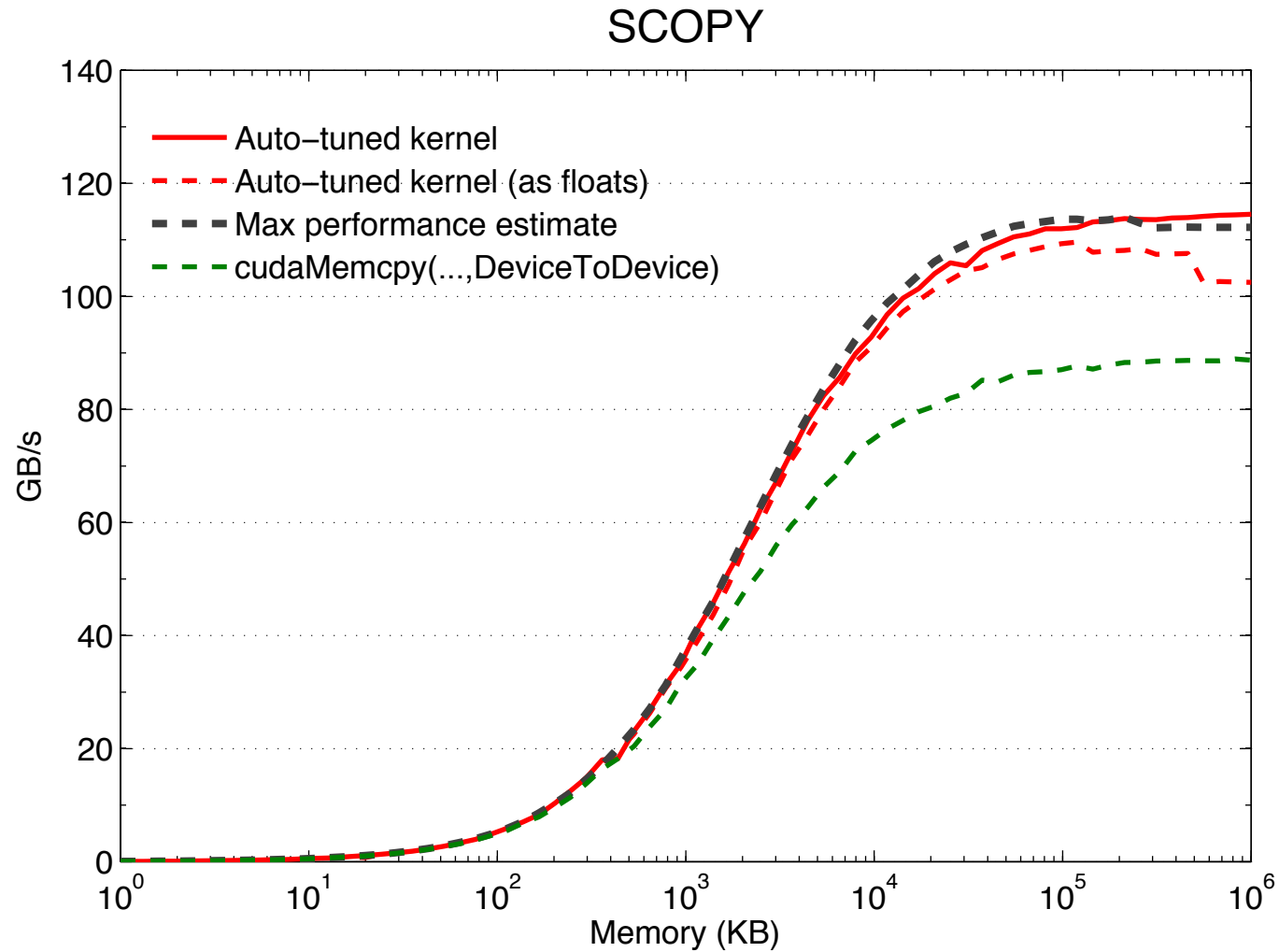
C2050: Theoretical bandwidth 144 GB/s

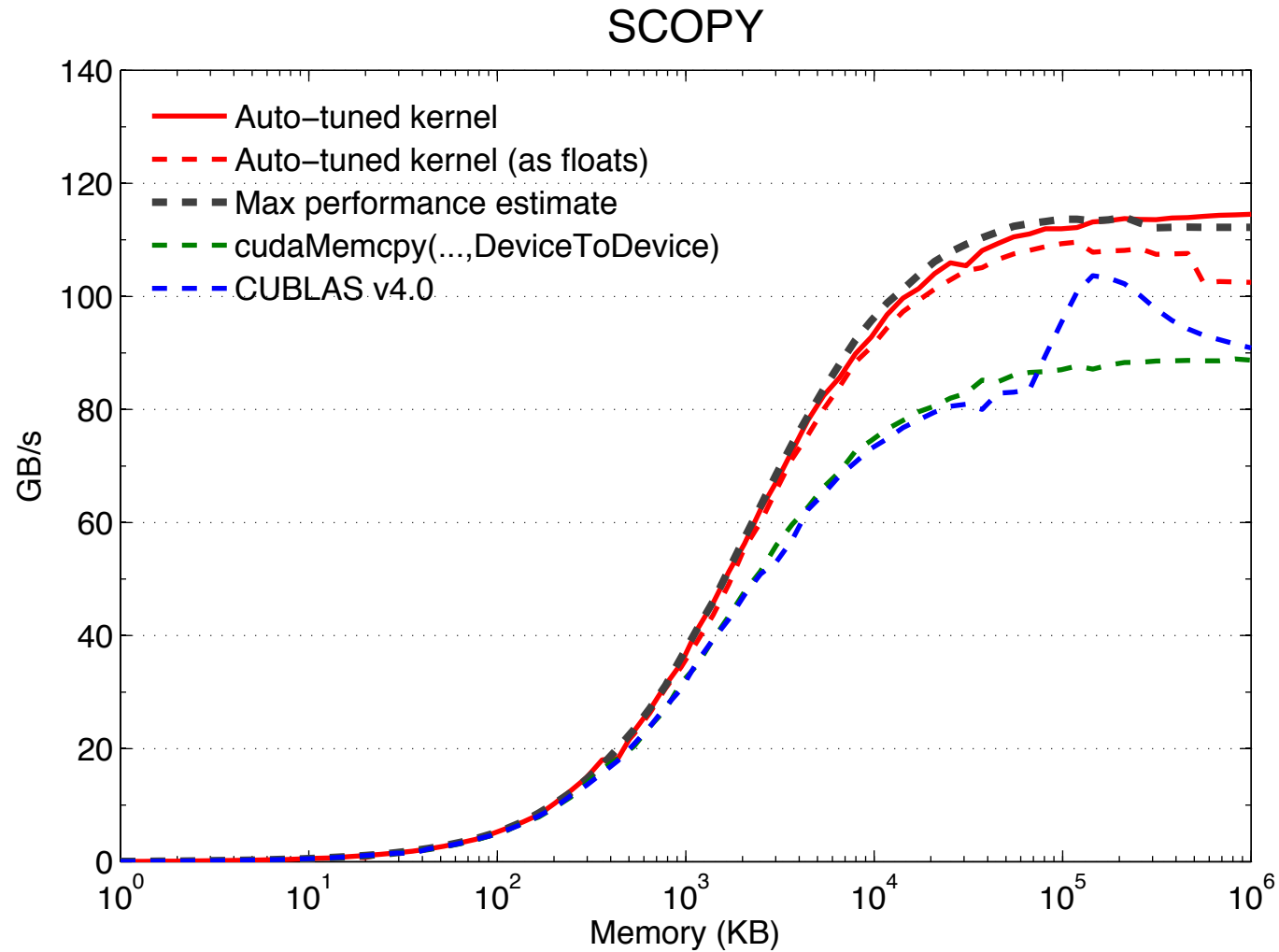
What is the effective bandwidth?





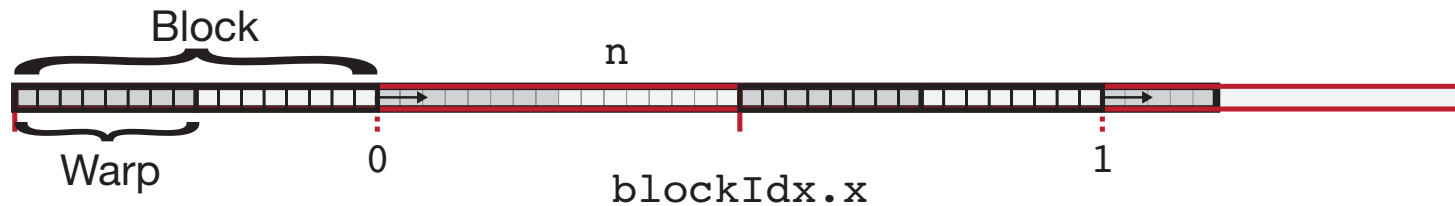




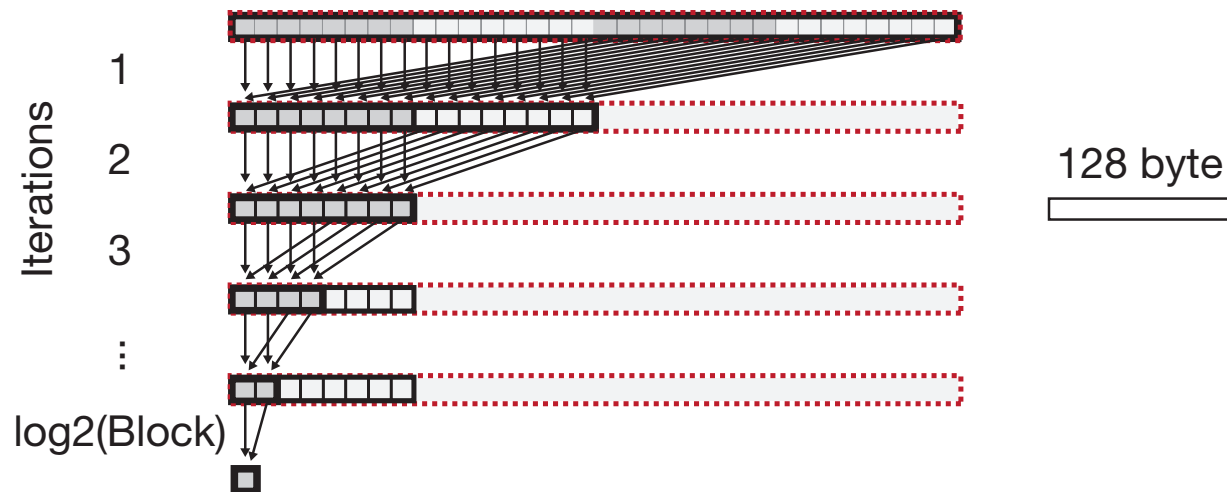


Level 1 BLAS: Vector Operations

Elementwise Vector Operation - Coalesced



Reduction Operation in Shared Memory



TUNING PARAMETERS: Block Size, Work Size per Thread, Unroll level

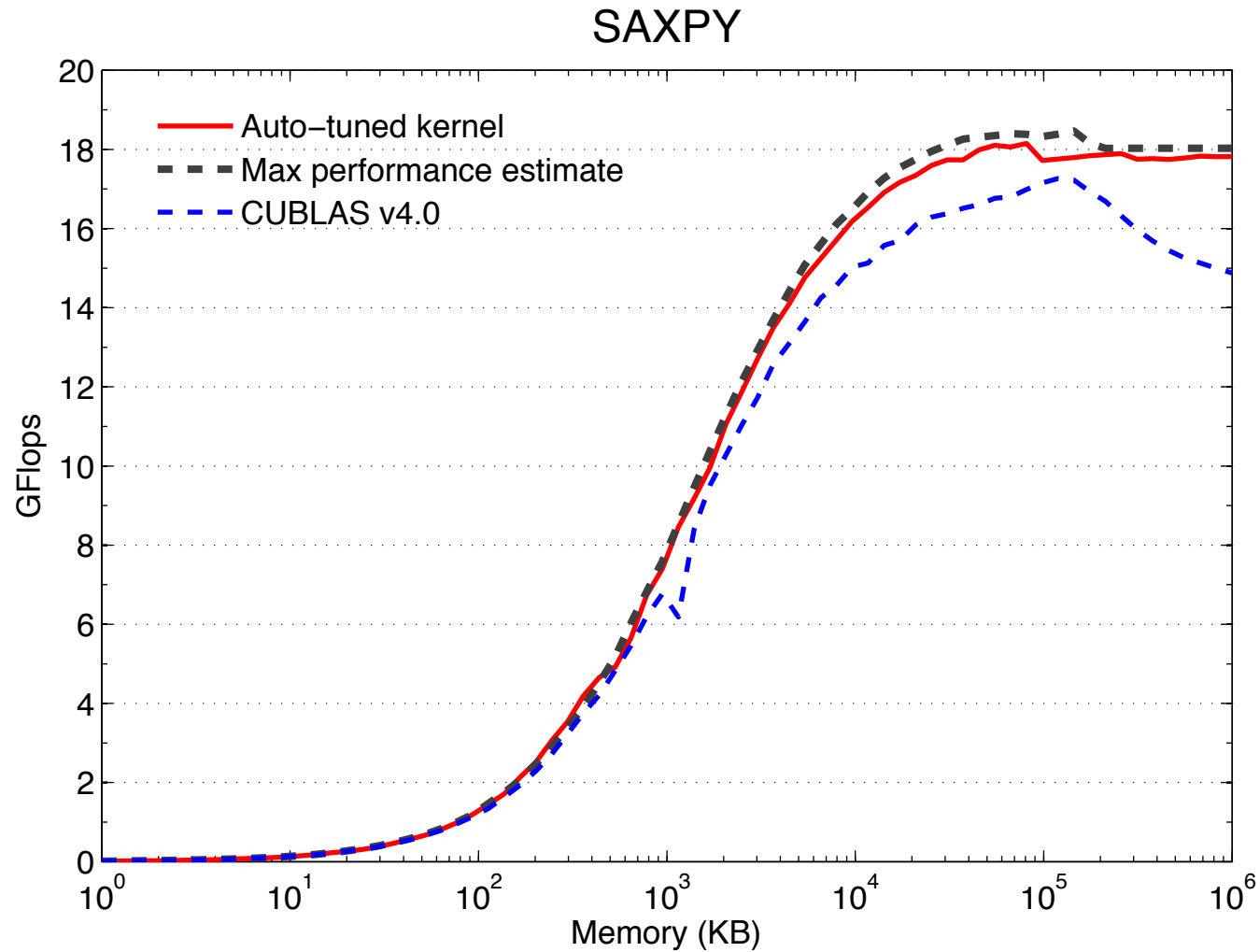
Auto-tuning for High Performance

- Heuristic search of parameter space:

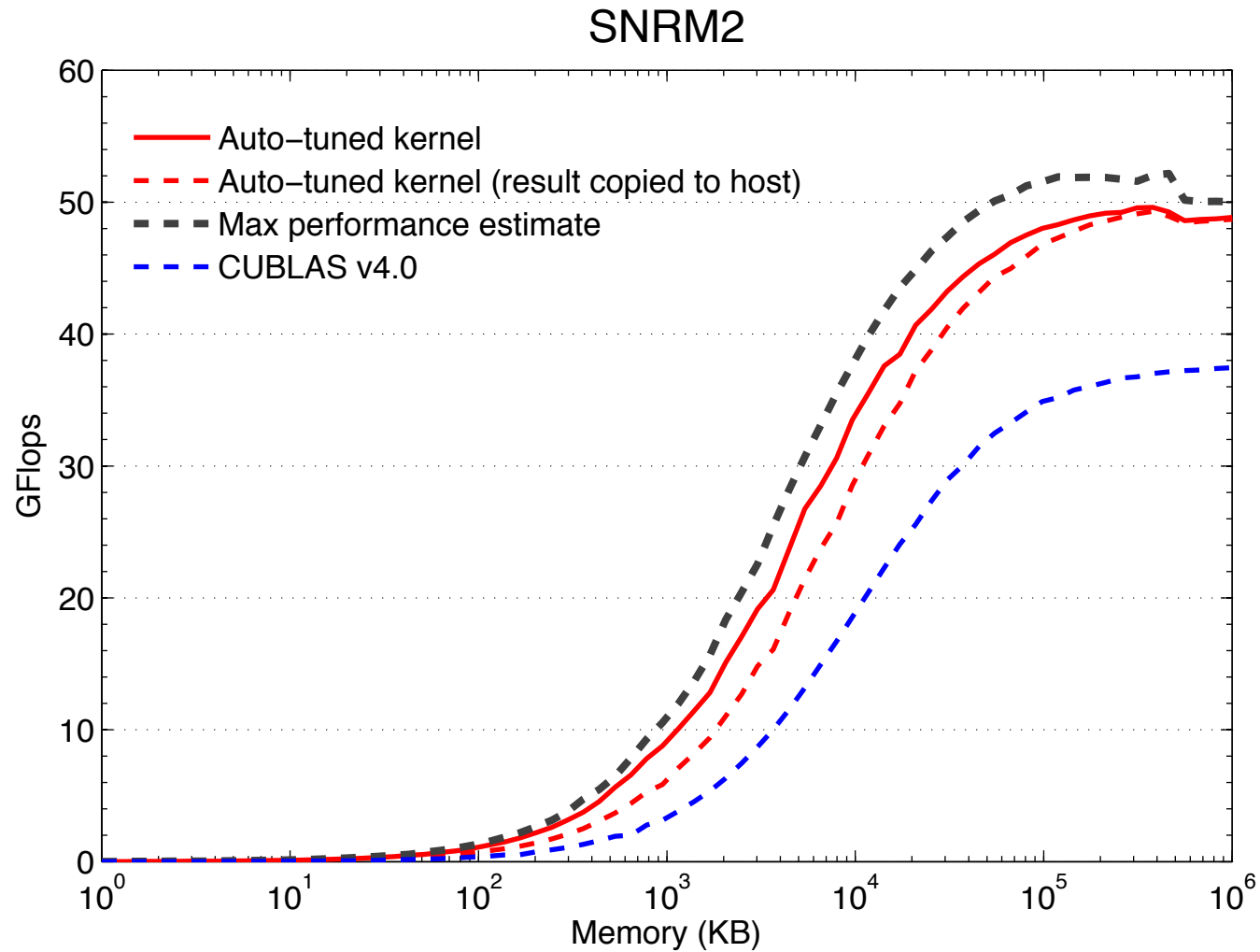
$$\text{BLOCKSIZE} \in \{32, 64, 96, 128, 160, 192, 224, 256\}$$

$$\text{WORKSIZE}_n \in \{1, 2, 3, 4, 5, 6, 7, 8\} \times \text{BLOCKSIZE}$$

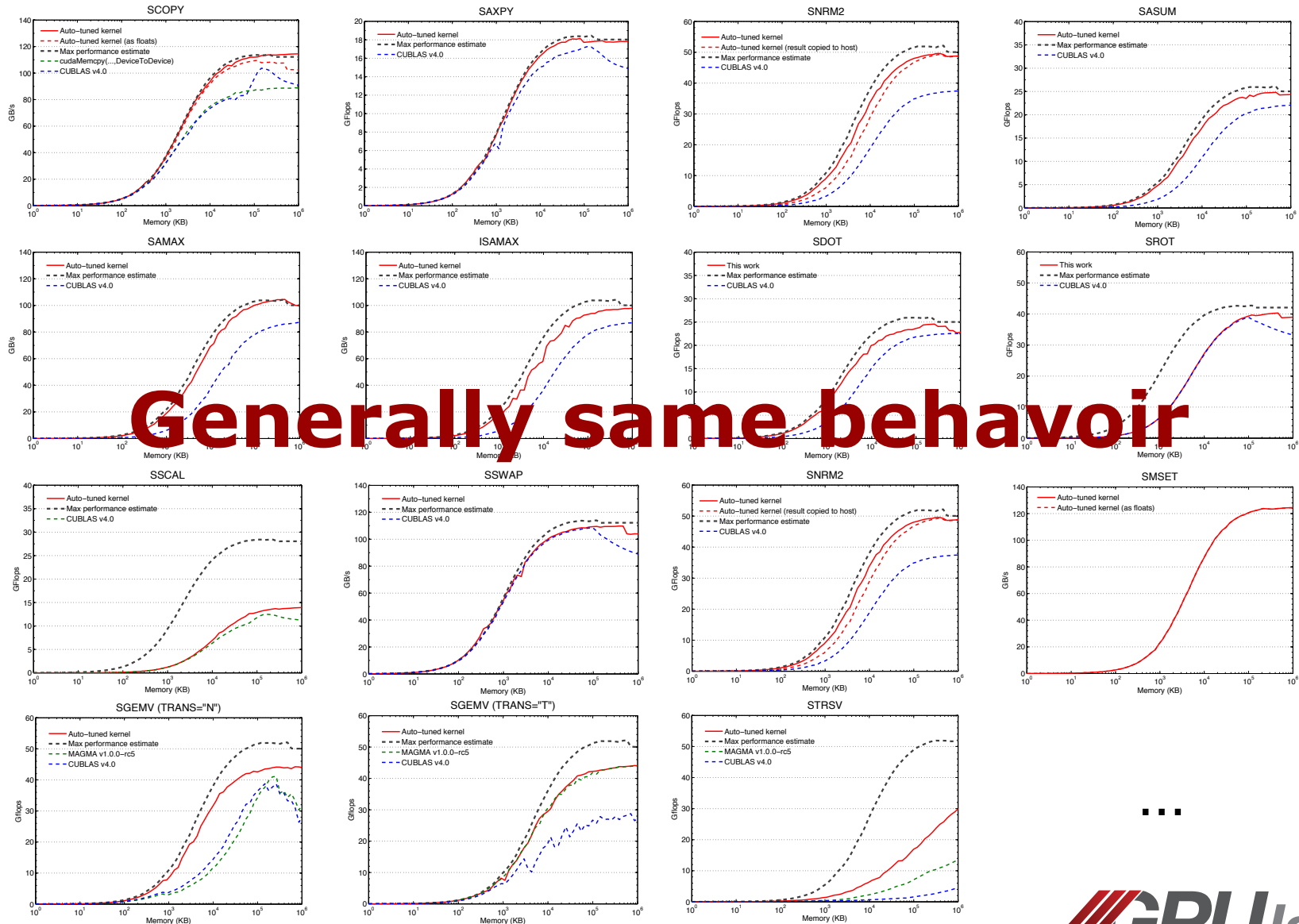
$$\text{UNROLL_LEVEL} \in \{1, 2, 3, 4, 5, 6, 7, 8\}$$



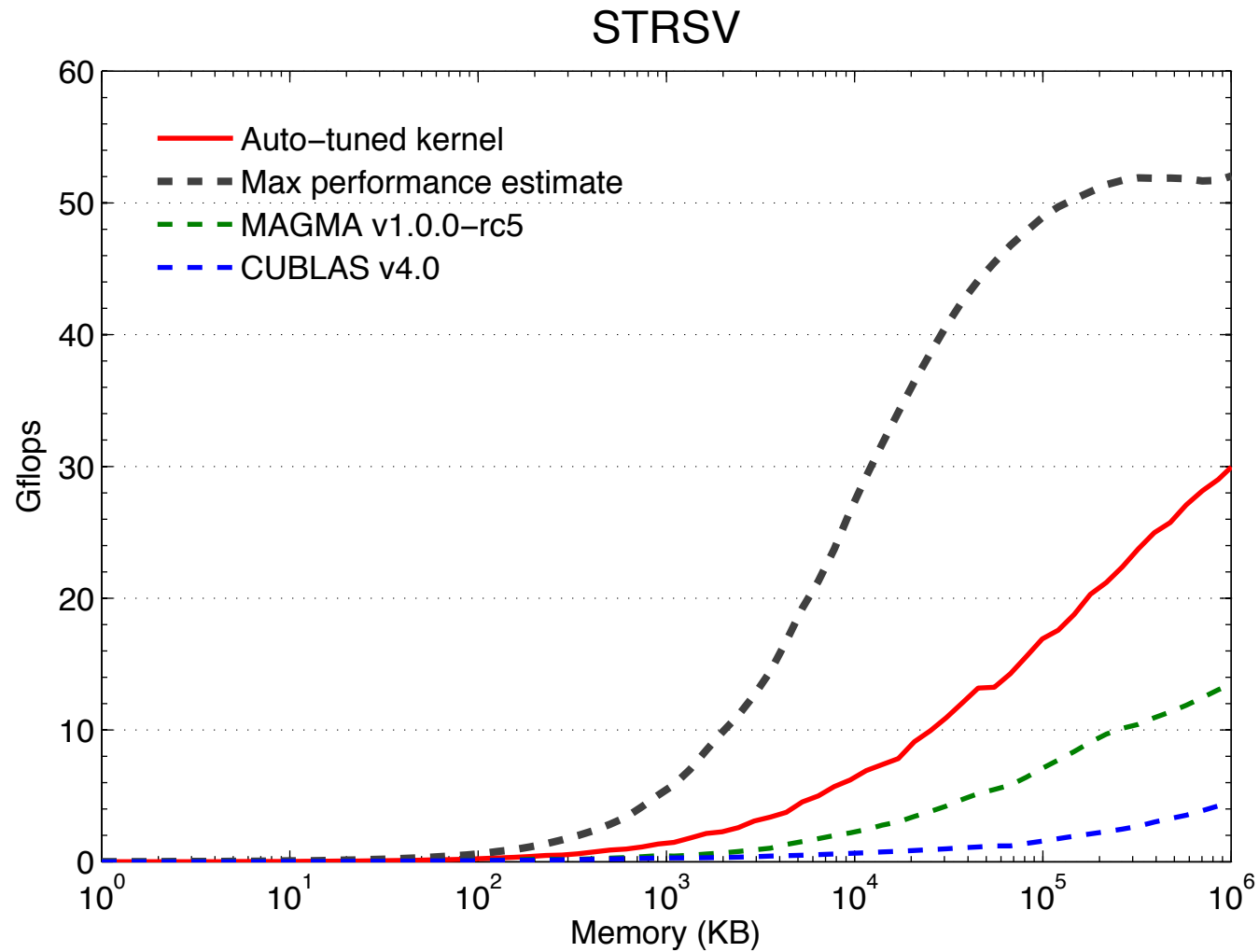
SNRM2



GLAS on Nvidia C2050



STRSV

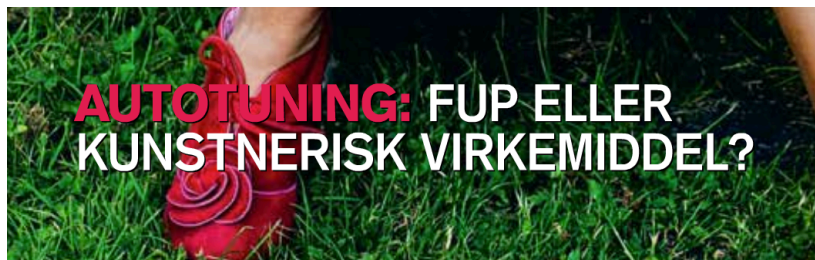


Summary

- BLAS level 1 and 2 kernels are totally memory bound.
- GPU's effective bandwidth sets the max performance.
- Simple auto-tuning can facilitate high-performance BLAS kernels for all input sizes and shapes.

- GLAS for single precision is available for download at `gpulab.imm.dtu.dk` (Open Source MIT License):

- `glas_v0.2_C2050_cuda_4.0_linux.tar.gz`
- `glas_v0.2_GTX590_cuda_4.0_linux.tar.gz`



Autotuning: Scam or artistic instrument?