



Tuning of BLAS level 1 and 2

Sørensen, Hans Henrik Brandenburg

Publication date:
2011

Document Version
Publisher's PDF, also known as Version of record

[Link back to DTU Orbit](#)

Citation (APA):
Sørensen, H. H. B. (2011). Tuning of BLAS level 1 and 2 [Sound/Visual production (digital)]. Workshop on "GPU Computing Today and Tomorrow", Kongens Lyngby, Denmark, 01/01/2011, <http://gpulab.imm.dtu.dk/courses.html>

DTU Library

Technical Information Center of Denmark

General rights

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the public portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the public portal

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

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 **Whenever possible, use Level 3**
 - $(2 \times) 4 \times N$ **BLAS in your (GPU) applications!**
 - Compute bound for large N

BLAS (Basic Linear Algebra Subprograms)

- Level 1 BLAS: (xAXPY, xAMAX, xDOT, xNRM2, etc.)
 - Vector of length N (4xN bytes)
 - 4xN bytes
 - Memory bound
- Level 2 BLAS: (xGEMV, xSYMV, etc.)
 - Matrix-vector multiply (4xN^2 bytes)
 - 4xN^2 bytes
 - Memory bound
- Level 3 BLAS: (xGEMM, xSYMM, xHEMM, xTRSM, etc.)
 - 1 or 2 matrices of size NxN (4xNxN bytes)
 - (2x)4xN^2 bytes : O(N^3) flops
 - Compute bound for large N

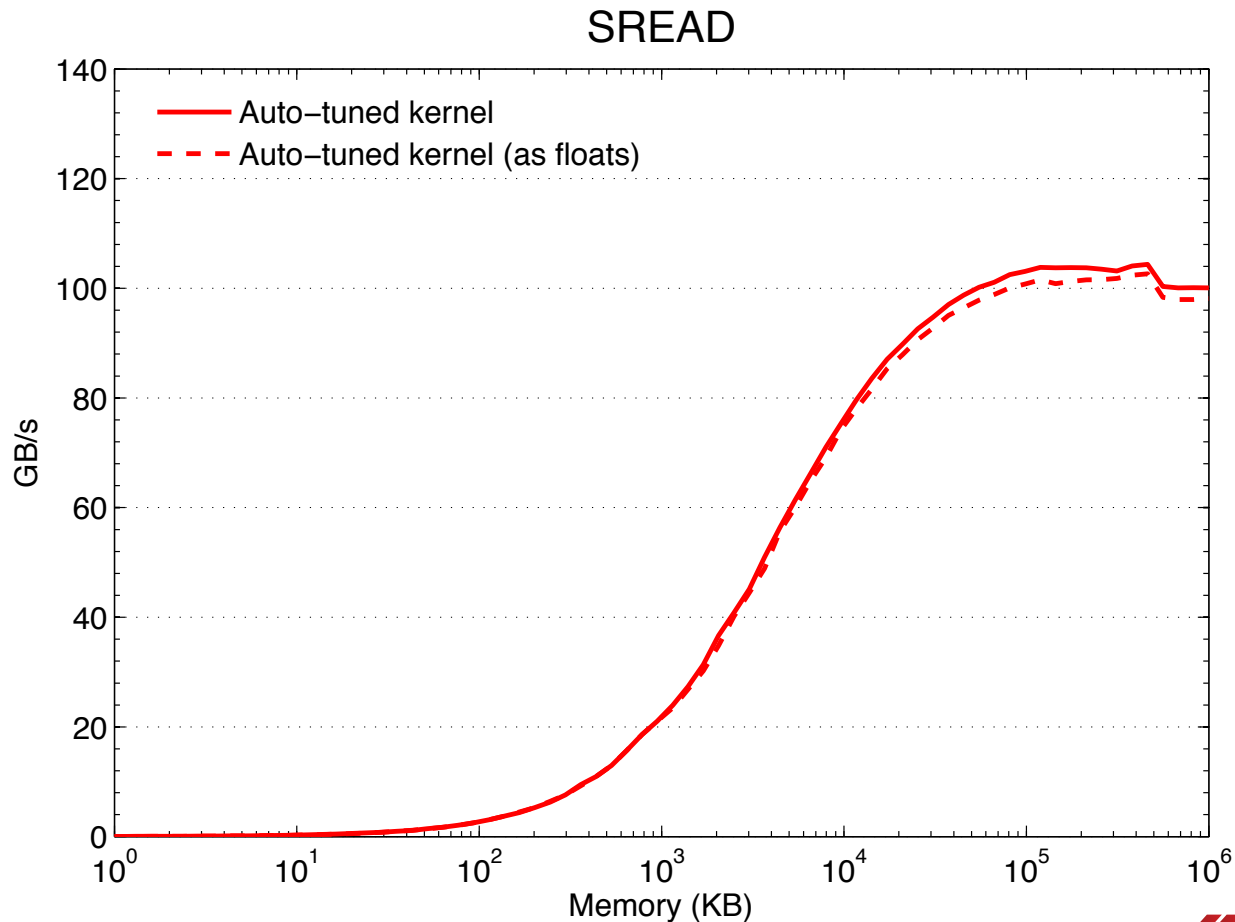
Why consider Level 1 and Level 2:

- Matvecs, orthogonalizations, norms, triangular solves, LAPACK building blocks, e.g., factorizations.
- Still very little attention from GPU community.

Performance Prediction

C2050: Theoretical bandwidth 144 GB/s

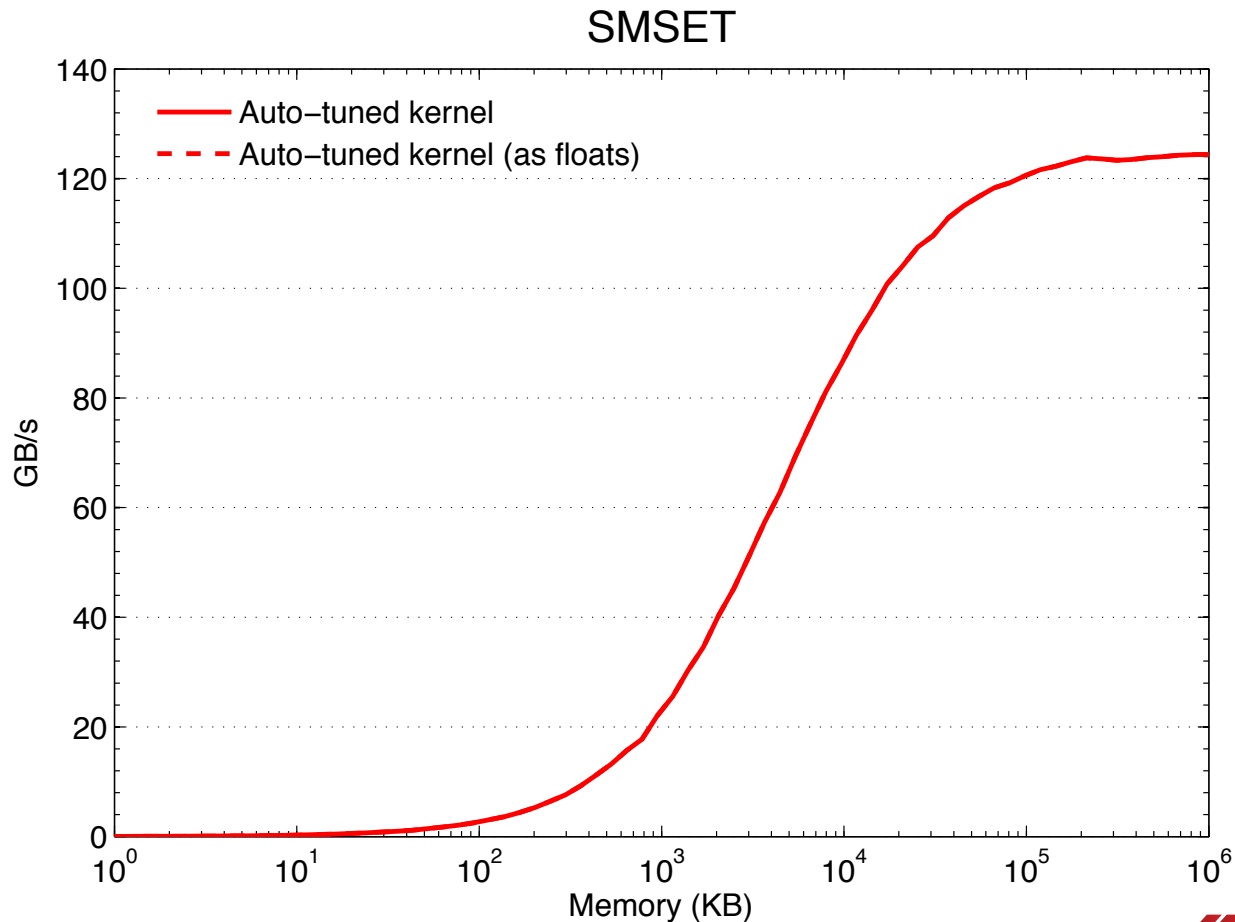
What is the effective bandwidth?



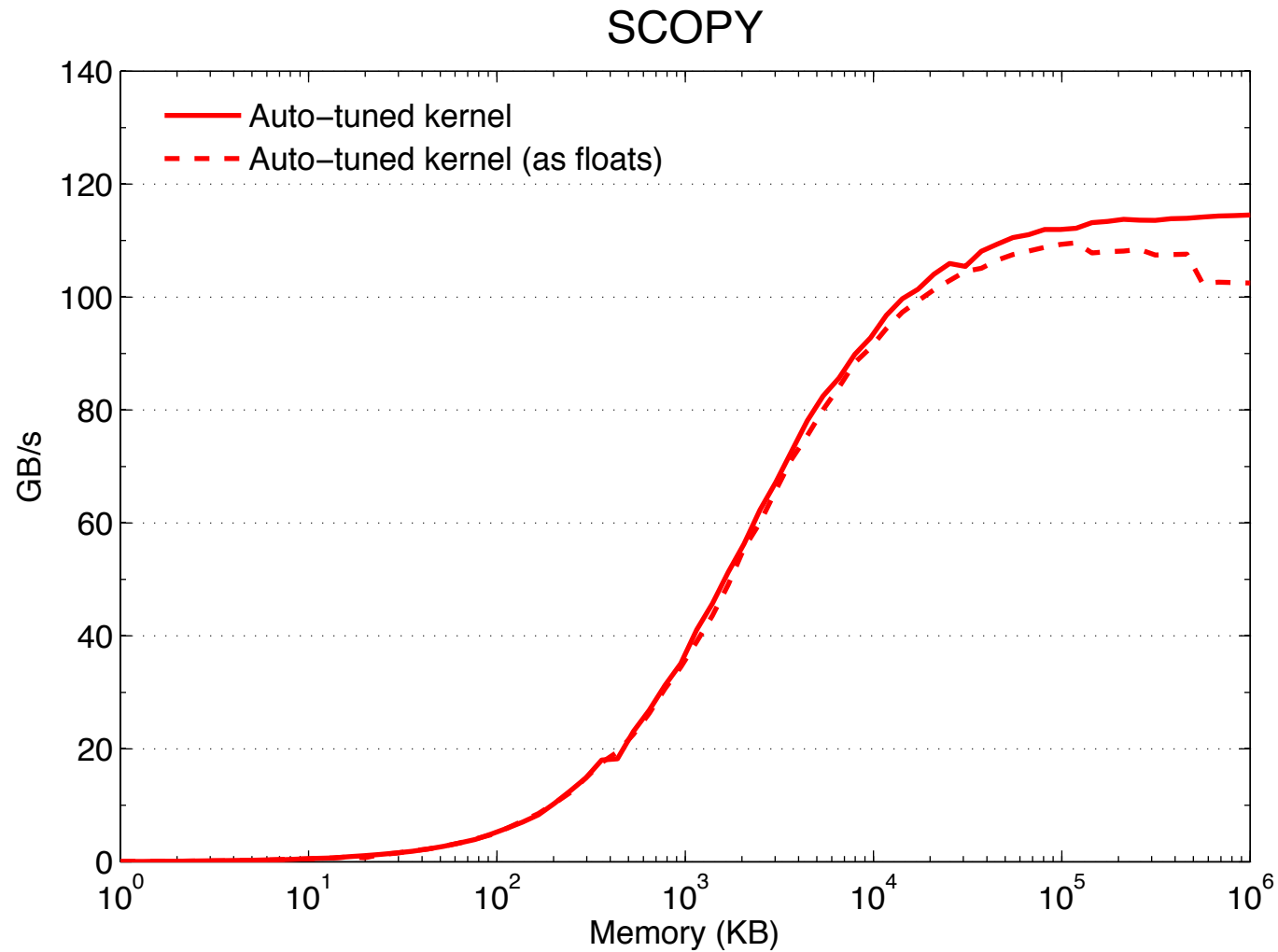
Performance Prediction

C2050: Theoretical bandwidth 144 GB/s

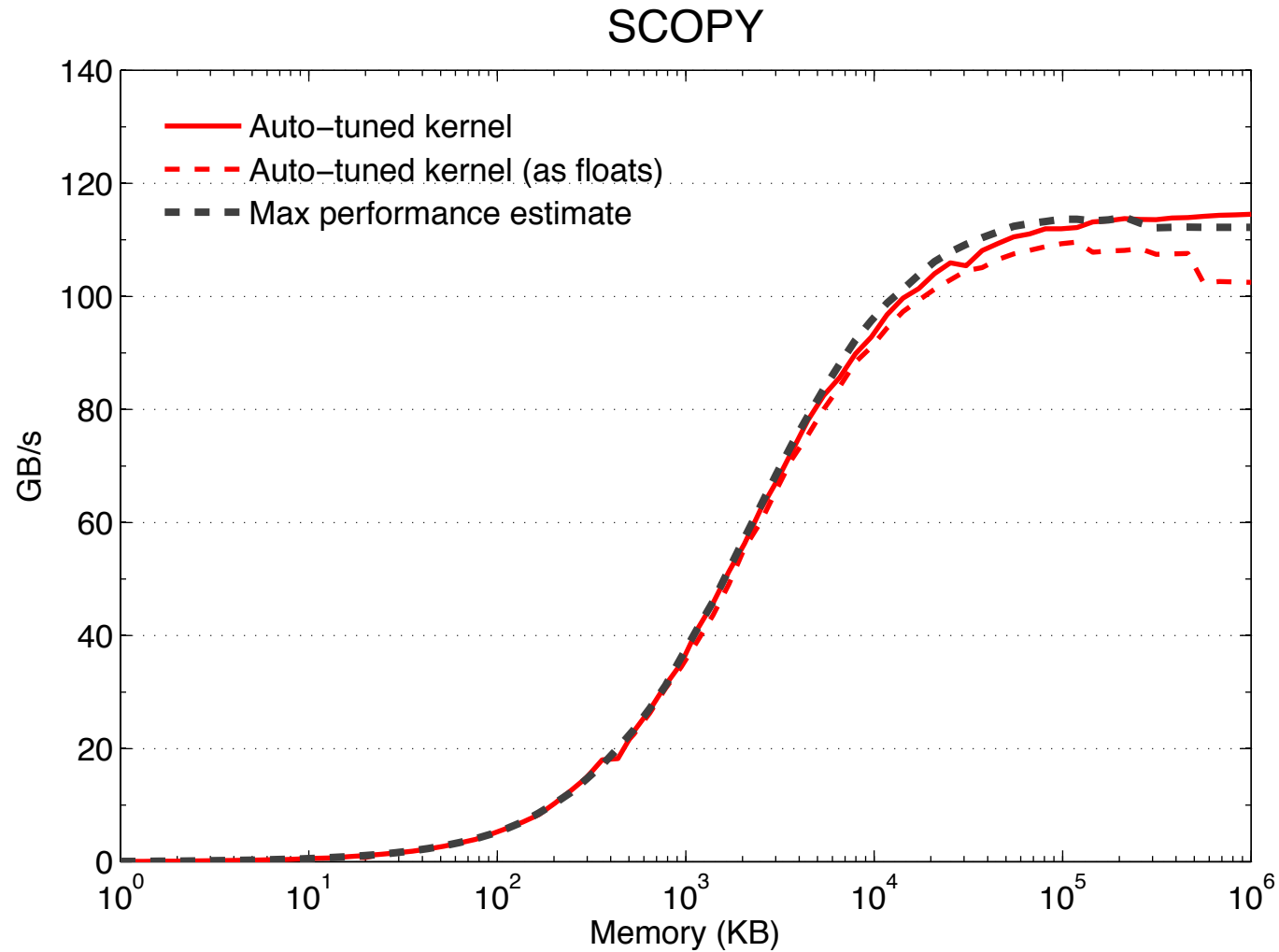
What is the effective bandwidth?

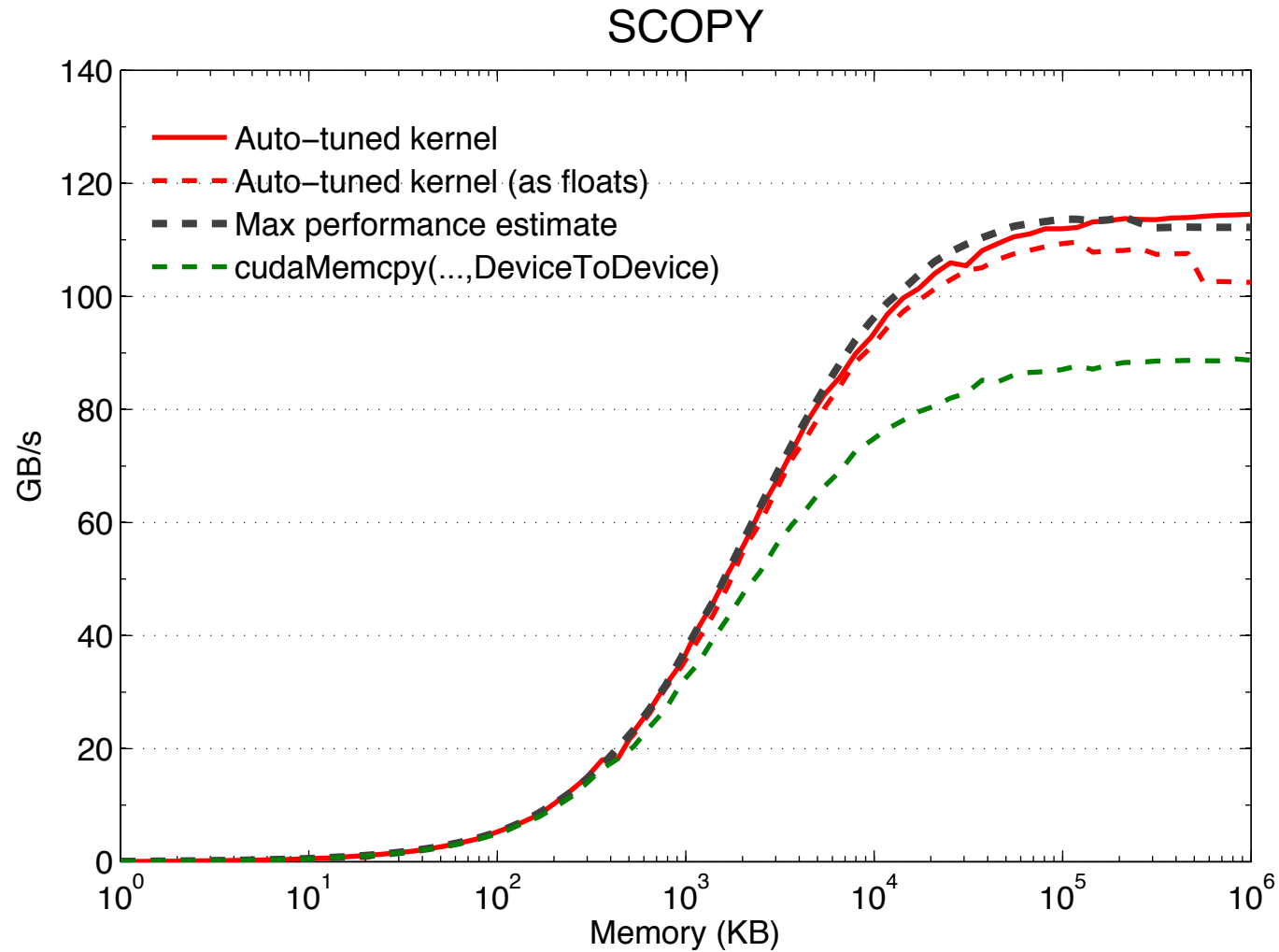


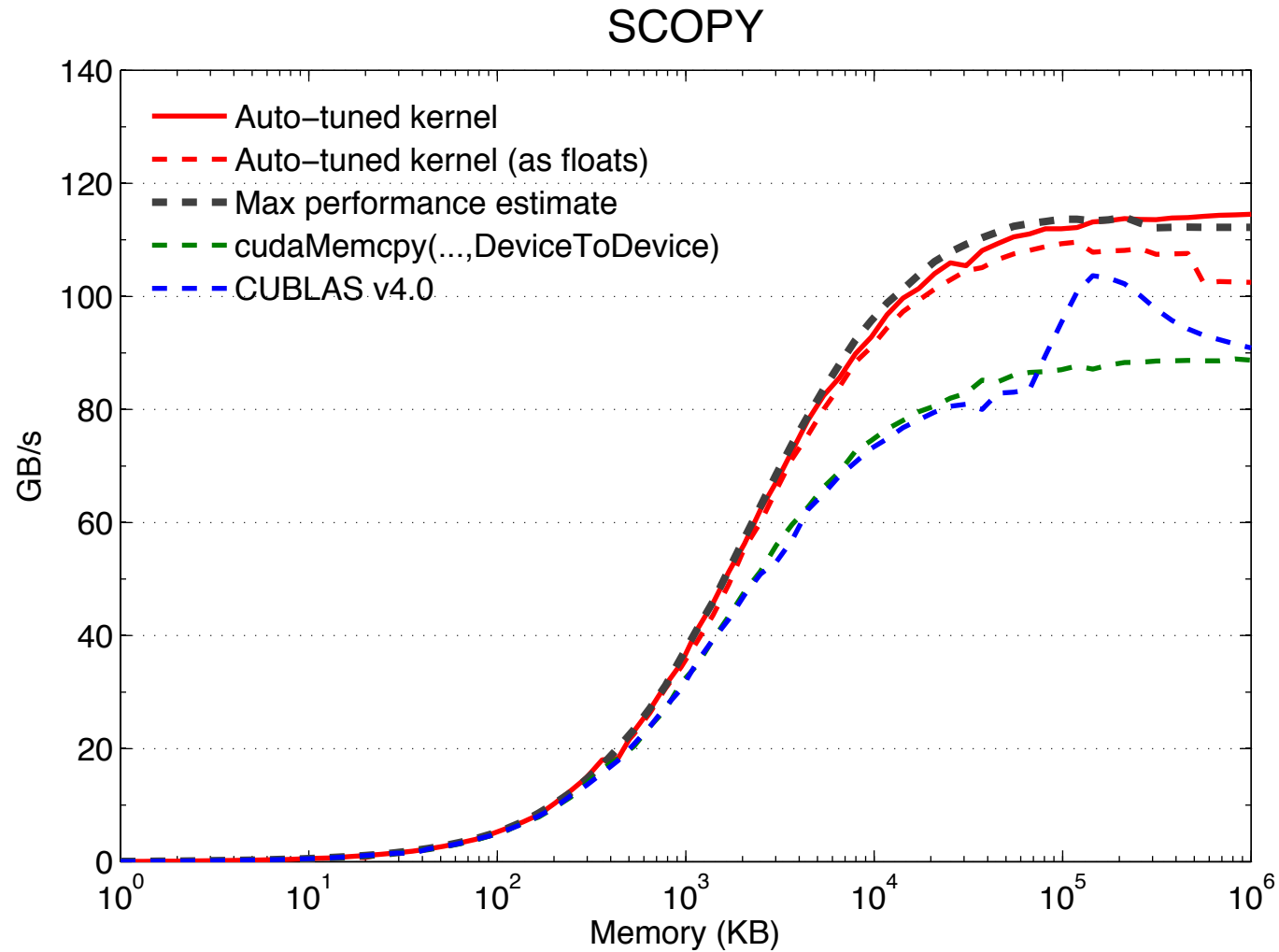
SCOPY



SCOPY

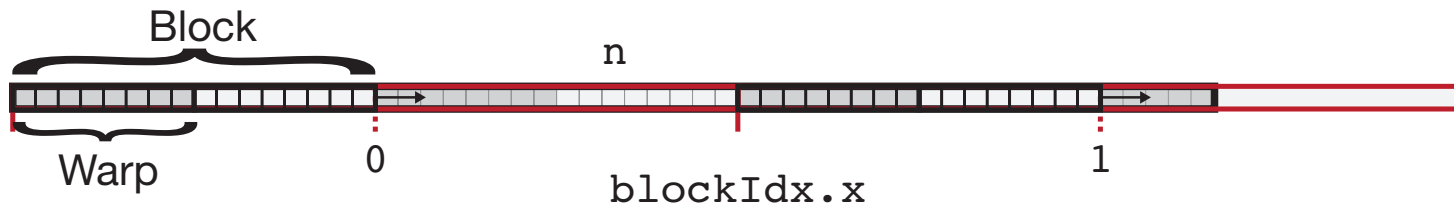




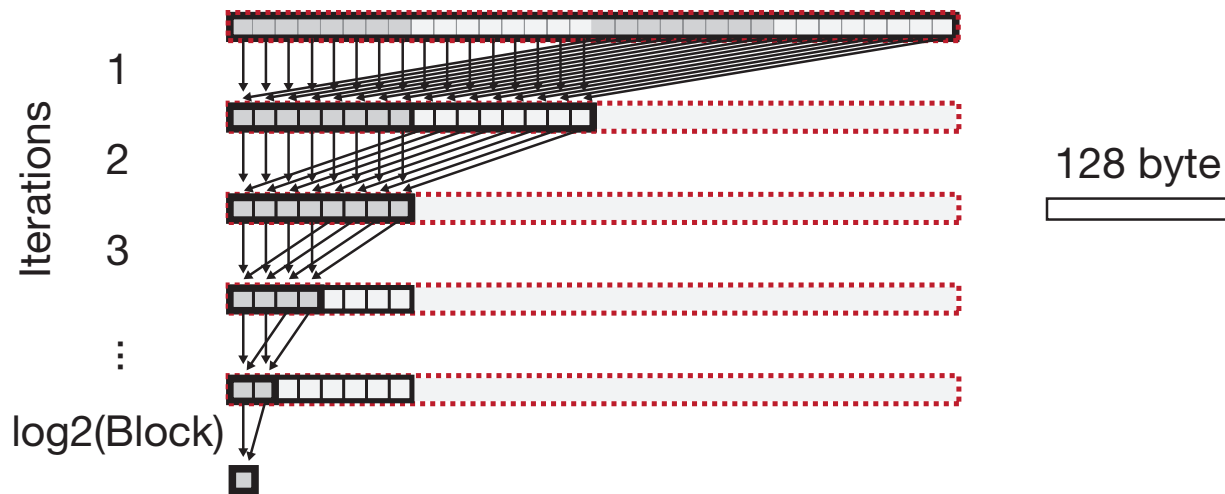


Level 1 BLAS: Vector Operations

Elementwise Vector Operation - Coalesced



Reduction Operation in Shared Memory



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

Auto-tuning for High Performance

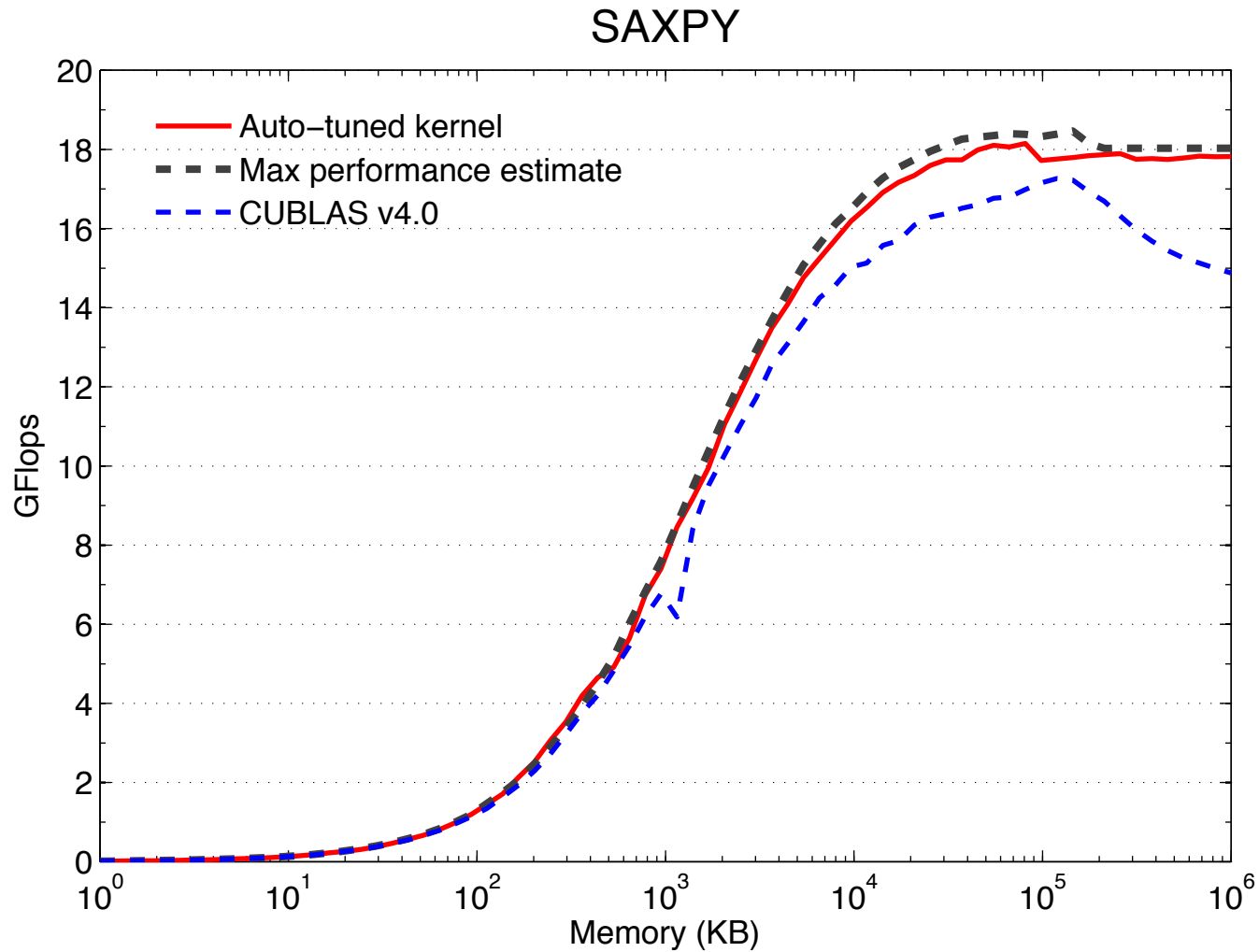
- Heuristic search of parameter space:

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

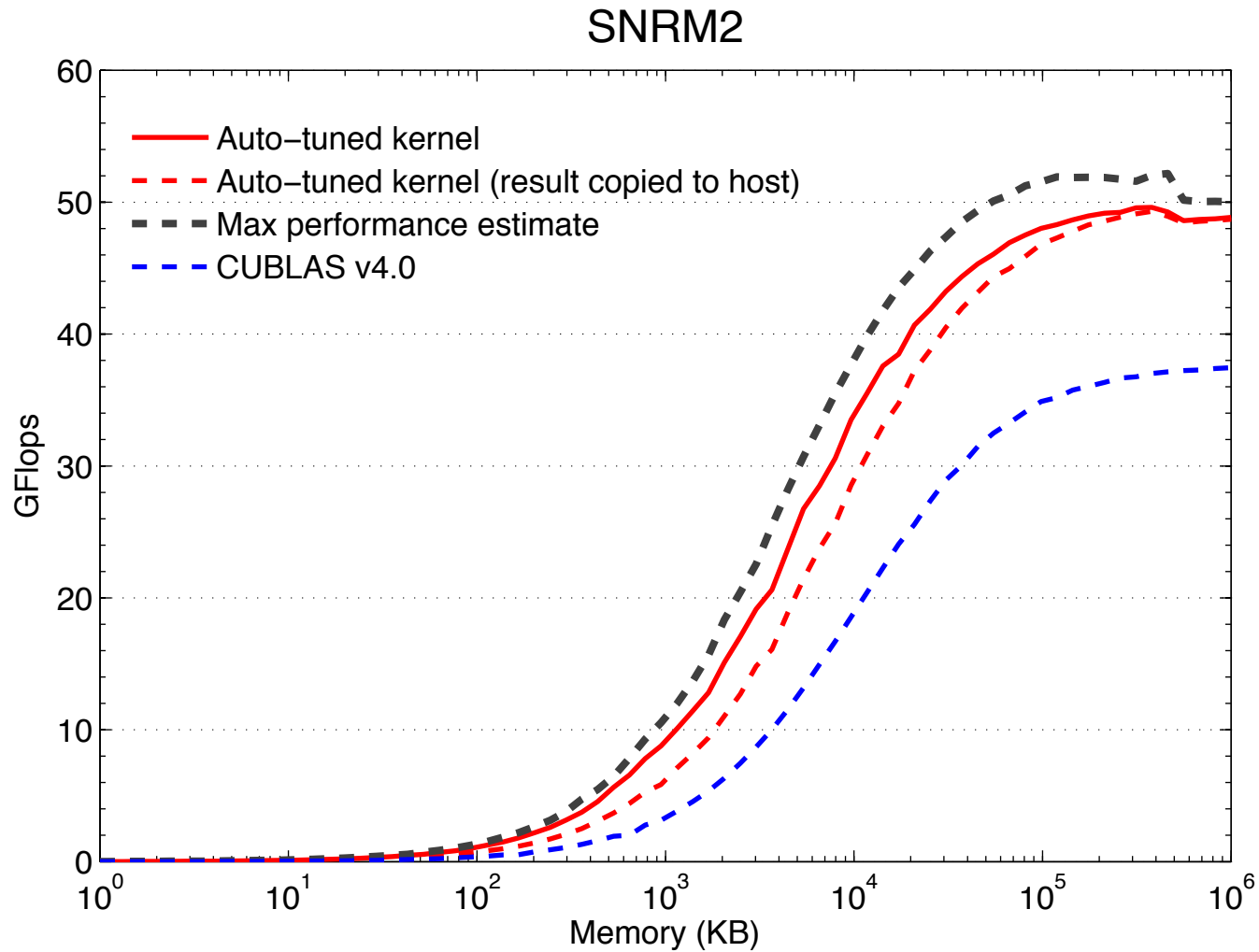
$\text{WORKSIZE}_n \in \{1, 2, 3, 4, 5, 6, 7, 8\} \times \text{BLOCKSIZE}$

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

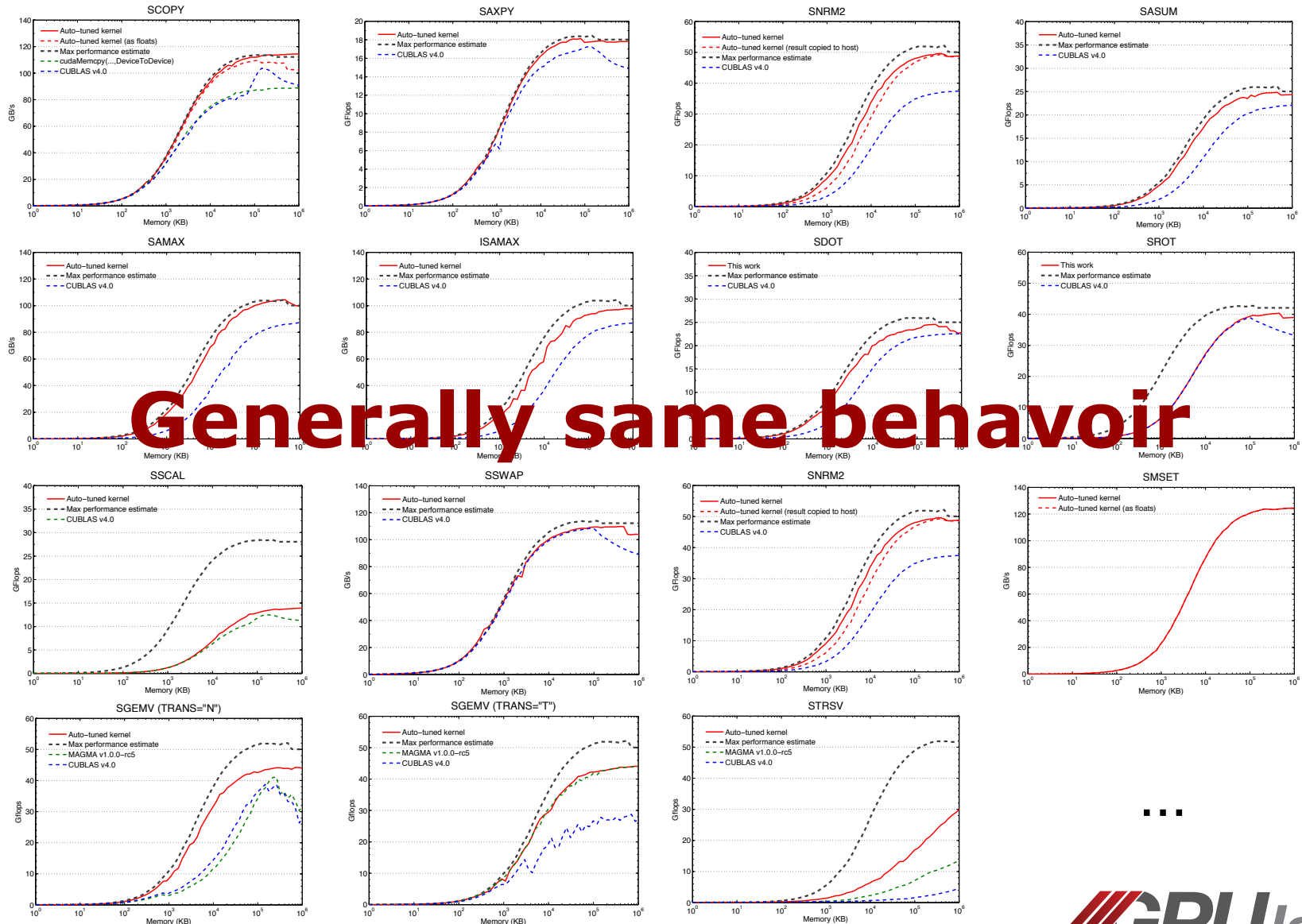
SAXPY



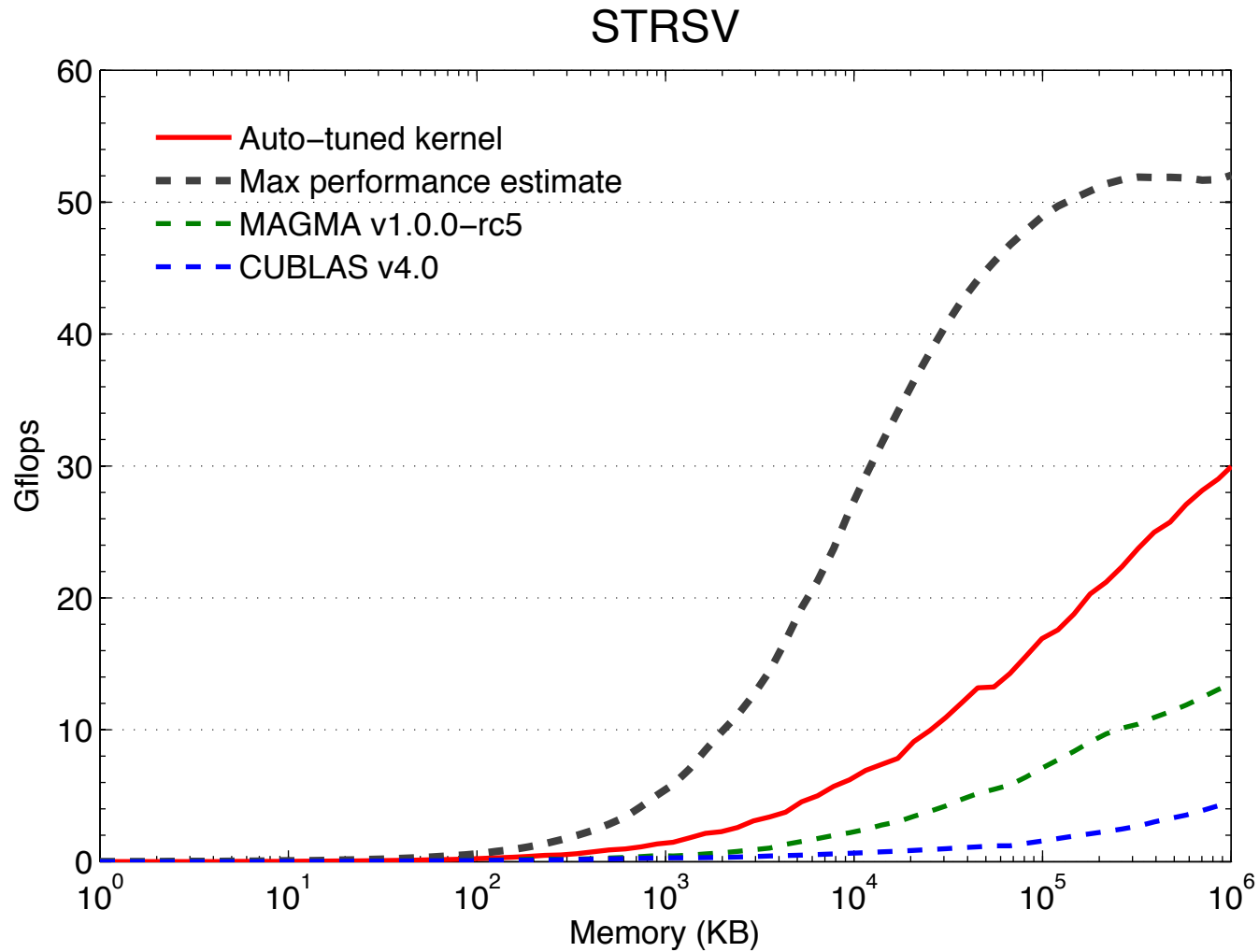
SNRM2



GLAS on Nvidia C2050



STRSV

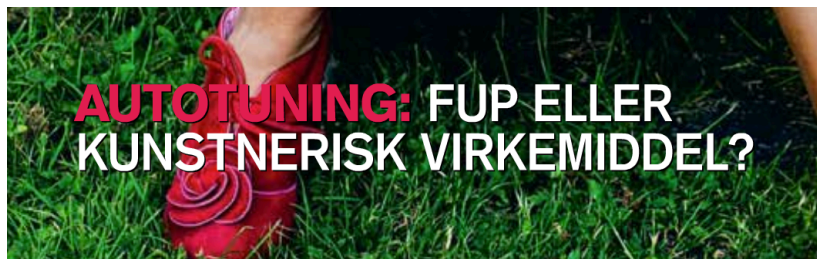


Summary

- 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?