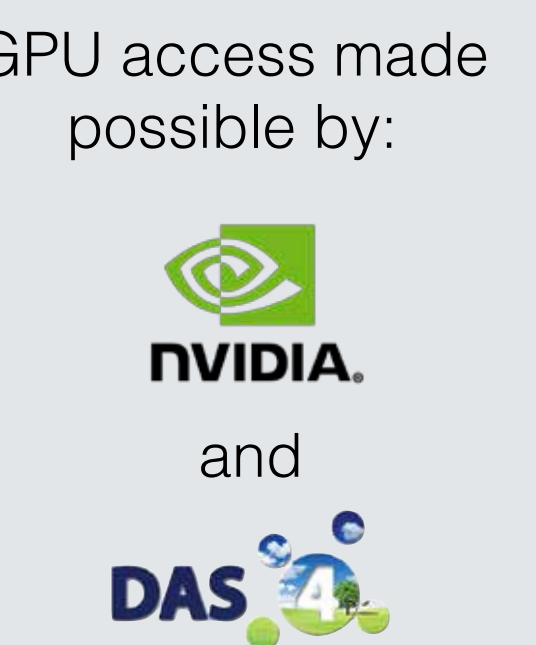


Cedric Nugteren

cedric.nugteren@surfsara.nl
SURFsara Supercomputer Centre
Amsterdam, The Netherlands

Auto-Tuning OpenCL Matrix-Multiplication: K40 versus K80



Generic OpenCL Auto-Tuning Framework

- Specify the OpenCL kernel **source-file**
- Define tuning **parameters** and their **values**
- Give kernel **arguments**
- Optional: specify a **reference** kernel
- Run** the tuner

Annotations:

- OpenCL kernel with pre-processor parameters
- Tuner iterates over all legal permutations
- Built-in correctness verification

Now: full-search; Future: smart-search strategies

Tuner available on GitHub (open-source): <https://github.com/cnugteren/cltune>

OpenCL Matrix-Multiplication (SGEMM/DGEMM)

Per thread tiling

Per workgroup tiling

Unrolled by a factor K_{wi}

$N_{wi} \times N_{dimC} = N_{wg}$

Tuneable parameters

- Per workgroup tiling: M_{wg}, N_{wg}, K_{wg}
- Workgroup sizes: M_{dimC}, N_{dimC}
- Re-shaped tiles: $M_{dimA} \times K_{dimA}, K_{dimB} \times N_{dimB}$
- Tile in local memory: $\$localA, \$localB$
- Stride within a thread: $M_{stride} (A\&C), N_{stride} (B)$
- Vector widths: $M_{vec} (A\&C), N_{vec} (B)$
- Unroll factor K_{wg} loop: K_{wi}

SGEMM/DGEMM code and parameters based on:

- Performance Tuning of Matrix Multiplication in OpenCL on Different GPUs and CPUs. K. Matsumoto, N. Nakasato, S.G. Sedukhin. In: *SC-Companion '12*. IEEE.
- OpenCL SGEMM tuning tutorial. C. Nugteren. <http://www.cedricnugteren.nl/tutorial.php>

Best-Case Matrix-Multiplication Parameters

Tested ~100K permutations!

Dual GPU board

Same chip as K40m but twice the number of registers and L1 cache

Tile sizes	NVIDIA Tesla K40m		½ NVIDIA Tesla K80		AMD Radeon HD7970	
	SGEMM	DGEMM	SGEMM	DGEMM	SGEMM	DGEMM
M_{wg}, N_{wg}, K_{wg}	128, 128, 16	64, 64, 8	128, 128, 8	64, 128, 8	128, 128, 16	128, 64, 8
M_{dimC}, N_{dimC}	16, 16	8, 16	8, 16	8, 16	16, 16	16, 16
M_{dimA}, N_{dimB}	16, 16	32, 16	32, 16	16, 16	16, 16	32, 32
$\$localA, \$localB$	yes, yes	yes, yes	yes, yes	yes, yes	yes, yes	no, no
M_{stride}, N_{stride}	yes, yes	yes, yes	yes, yes	yes, yes	yes, yes	yes, no
M_{vec}, N_{vec}	2, 2	1, 1	4, 4	2, 2	2, 8	2, 1
Unroll factor	8	8	8	8	16	4

Best vector width **x2** for K80

Best vector width **x½** for doubles

Only slightly slower without local memory or strides (definitely not the case for K40/K80)

Average execution time after a choice e.g. $M_{vec} = 4$

Visualization of Tested Parameter Permutations

50ms, 100ms, 150ms, 200ms, 250ms, 300ms

SGEMM $M=N=K=4096$ execution time

Subset with $\$localA = \$localB = M_{str} = N_{str} = yes$

corresponding colours per GPU

top: Tesla K40m

bottom: ½ Tesla K80 (w/o boost)

Best performance at low occupancy

Parameters with highest impact at the bottom

More registers

Performance in GFLOPS

Assembly gap (LDS, register bank conflicts, ...)

peak 4291

SGEMM $M=N=K=4096$

peak 4368

peak 1430

DGEMM $M=N=K=2048$

peak 935

peak 1126

Legend:

- cBLAS 2.2.0
- This work (best-case)
- cuBLAS 6.5

The AMD GPU has significantly better performing source-level code

cBLAS also uses auto-tuning, but uses an inferior and less tuneable kernel

Both using non-IEEE-754 "mad()"

all K80 results +55% at max clock

Conclusions

- Generic **OpenCL auto-tuning** framework ensures **portability**
- Fastest** source-level **matrix-multiplication** implementation
- Tree **visualization** shows architectural and algorithm properties
- K80's** extra registers and cache yield **higher efficiency** (vs. K40)