

# Foundations of High Performance Computing

## Lecture 11: Using High Performance Libraries

2021-2022 Stefano Cozzini



“Foundation of HPC” course  
**DATA SCIENCE &  
SCIENTIFIC COMPUTING**

# Agenda

- Intro:
  - What are we using HPC for ?
  - Where should we start optimizing ?
- High Performance Libraries
- Linear Algebra libraries
- Using HP libraries: some examples.

# The seven dwarfs of HPC



- Phil Colella (LBL) identified **7 kernels** of which most simulation and data analysis program are composed:
- Dense Linear Algebra
  - Ex: solve  $Ax=B$  or  $Ax=\lambda x$  where  $A$  is a dense matrix
- Sparse Linear Algebra
  - solve  $Ax=B$  or  $Ax=\lambda x$  where  $A$  is a sparse matrix (mostly zero)
- Operation on structured Grids:
  - $ANEW_j() = 4 * (A(i,j) - A(i-1,j) - A(i+1,j) - A(i,j-1) - A(i,j+1))$
- Operation on unstructured Grids:
  - similar but list of neighbours varies from entry to entry
- Spectral Methods
  - Fast Fourier Transform (FFT)
- Particle Methods
  - Compute electrostatic forces on  $n$ -particles
- Monte Carlo
  - many independent simulation using different inputs

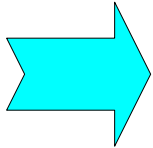
# Is this the real picture in 2021 ?

- We are missing ALL data analytics working load..
- But...
- A lot of them relies on highly optimized libraries.

Where should you start  
optimizing your application ?

# Optimization techniques

- There are basically three different categories:
  - Improve memory performance (the most important)
  - Improve CPU performance
  - Use already highly optimized libraries/subroutines



The easiest and more efficient way..

# What are High Performance Libraries ?

- Routines for common (math) functions written in a specific way to take advantage of all capabilities of the CPU.
- Each CPU type normally has its own version of the library specifically written or compiled to maximally exploit that architecture

# Why using High Performance Libraries ?

- Compilers can optimize code only to a certain point. Effective programming needs deep knowledge of the platform
- Performance libraries are designed to use the CPU in the most efficient way, which is not necessarily the most straightforward way.
- It is normally best to use the libraries supplied by or recommended by the CPU vendor
- On modern hardware they are hugely important, as they most efficiently exploit caches, special instructions and parallelism
- **Parallelism (at least on single node) comes for free..**



# Any other reason apart from performance ?

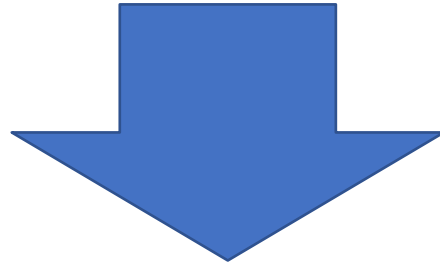
- Usage of libraries makes coding easier.  
Complicated math operations can be used from existing routines
- Increase portability of code as standard (and well optimized) libraries exist for ALL computing platforms.
- Lego approach: build your own code using already available bricks..

# What is available ?

- Linear Algebra:
  - BLAS/LAPACK/SCALAPACK
- FFT:
  - FFTW
- ODE/PDE
  - PETSC
- Machine Learning:
  - Tensorflow / Caffe etc..

# Should I write my own algorithm for L. A. ?

- 99.99% of time: NO !
- Tons of libraries out there
- Well tested
- Extremely efficient in 99.99% of the case
- With some “de facto” standard implemented



PORTABILITY IS COMING (almost) FOR FREE

# Why Linear Algebra ?

- Basic Linear Algebra Problems:
  - Linear Equations: Solve  $Ax=b$  for  $x$
  - Least Squares: Find  $x$  that minimizes  $\|r\|_2 \equiv \sqrt{\sum r_i^2}$  where  $r=Ax-b$ 
    - Statistics: Fitting data with simple functions
  - Eigenvalues: Find  $\lambda$  and  $x$  where  $Ax = \lambda x$ 
    - Vibration analysis, Quantum Simulations, etc..
  - Singular Value Decomposition:  $A^T Ax = \sigma^2 x$ 
    - Data fitting, Information retrieval

Lots of variations depending on structure of  $A$

# Why dense Linear Algebra ?

- Many large matrices are sparse, but ...
- Dense algorithms easier to understand
- Some applications yields large dense matrices
- Large sparse matrix algorithms often yield smaller (but still large) dense problems
- LINPACK Benchmark ([www.top500.org](http://www.top500.org))

“ How fast is your computer?”

=

“How fast can you solve dense  $Ax=b$ ?”

BLAS

Basic Linear Algebra Subprograms

# BLAS history

- In the beginning it was libraries like **EISPACK** (for eigenvalue problems)
- Then the **BLAS-1** were invented (1973-1977)
  - Create a standard library of 15 operations (mostly) on vectors “AXPY” ( $y = \alpha \cdot x + y$ ), dot product, scale ( $x = \alpha \cdot x$ ), etc
  - Up to 4 versions of each (S/D/C/Z), 46 routines, 3300 LOC
  - Language: FORTRAN
- GOALS
  - Common “pattern” to ease programming, readability
  - Robustness, via careful coding (avoiding over/underflow) --> Accuracy Portability (common interface)
  - Efficiency via machine specific implementations
  - Maintainability
- Why **BLAS-1**? They do  $O(n)$  ops on  $O(n)$  data
  - Used in libraries like LINPACK (for linear systems)
  - Source of the name “LINPACK Benchmark” (not the code!)

# BLAS history

- But the BLAS-1 weren't enough
  - Consider AXPY (  $y = \alpha \cdot x + y$  ):  $2n$  flops on  $3n$  read/writes
  - Computational intensity =  $(2n)/(3n) = 2/3$
  - Too low to run near peak speed (read/write dominates)
- So the BLAS-2 were developed (1984-1986)
  - Standard library of 25 operations (mostly) on matrix/vector pairs
  - “GEMV”:  $y = \alpha \cdot A \cdot x + \beta \cdot x$ , “GER”:  $A = A + \alpha \cdot x \cdot y^T$ ,  $x = T^{-1} \cdot x$
  - Up to 4 versions of each (S/D/C/Z), 66 routines, 18K LOC
- Why BLAS-2 ?
  - They do  $O(n^2)$  ops on  $O(n^2)$  data
  - So computational intensity still just  $\sim (2n^2)/(n^2) = 2$
  - OK for **vector machines**, but not for machine with caches



# BLAS history

- The next step: BLAS-3 (1987-1988)
  - Standard library of 9 operations (mostly) on matrix/matrix pairs
  - “GEMM”:  $C = \alpha \cdot A \cdot B + \beta \cdot C$ ,  $C = \alpha \cdot A \cdot A^T + \beta \cdot C$ ,  $C = T^{-1} \cdot B$
  - Up to 4 versions of each (S/D/C/Z), 30 routines, 10K LOC
- Why BLAS 3 ?
  - They do  $O(n^3)$  ops on  $O(n^2)$  data
  - So computational intensity  $(2 n^3)/(4 n^2) = n/2$  – big at last!
  - Good for machines with caches, other mem. hierarchy levels
  - Performing implementations left to others..

# Where can I get BLAS?

[www.netlib.org/blas](http://www.netlib.org/blas)

- Source: 142 routines, 31K LOC,
- Testing: 28K LOC
- Reference (unoptimized) implementation only !
- [http://www.netlib.org/blas/#\\_reference\\_blas\\_version\\_3\\_8\\_0](http://www.netlib.org/blas/#_reference_blas_version_3_8_0)
- Ex: 3 nested loops for GEMM

# BLAS list

## Level 1 BLAS

	dim	scalar	vector	vector	scalars	5-element array
SUBROUTINE xROTG (					A, B, C, S )	
SUBROUTINE xROTMG(					D1, D2, A, B,	PARAM )
SUBROUTINE xROT ( N,			X, INCX, Y, INCY,		C, S )	PARAM )
SUBROUTINE xROTM ( N,			X, INCX, Y, INCY,			
SUBROUTINE xSWAP ( N,			X, INCX, Y, INCY )			
SUBROUTINE xSCAL ( N,		ALPHA,	X, INCX )			
SUBROUTINE xCOPY ( N,			X, INCX, Y, INCY )			
SUBROUTINE xAXPY ( N,		ALPHA,	X, INCX, Y, INCY )			
FUNCTION xDOT ( N,			X, INCX, Y, INCY )			
FUNCTION xDOTU ( N,			X, INCX, Y, INCY )			
FUNCTION xDOTC ( N,			X, INCX, Y, INCY )			
FUNCTION xxDOT ( N,			X, INCX, Y, INCY )			
FUNCTION xRNRM2 ( N,			X, INCX )			
FUNCTION xASUM ( N,			X, INCX )			
FUNCTION IxAMAX( N,			X, INCX )			

Generate plane rotation  
 Generate modified plane rotation  
 Apply plane rotation  
 Apply modified plane rotation  
 $x \leftrightarrow y$   
 $x \leftarrow \alpha x$   
 $y \leftarrow x$   
 $y \leftarrow \alpha x + y$   
 $dot \leftarrow x^T y$   
 $dot \leftarrow x^T y$   
 $dot \leftarrow x^H y$   
 $dot \leftarrow \alpha + x^T y$   
 $norm2 \leftarrow ||x||_2$   
 $asum \leftarrow ||re(x)||_1 + ||im(x)||_1$   
 $amax \leftarrow 1^{*}k \ni |re(x_k)| + |im(x_k)|$   
 $= \max\{|re(x_i)| + |im(x_i)|\}$

prefixes  
 S, D  
 S, D  
 S, D  
 S, D  
 S, D, C, Z  
 S, D, C, Z, CS, ZD  
 S, D, C, Z  
 S, D, C, Z  
 S, D, DS  
 C, Z  
 C, Z  
 SDS  
 S, D, SC, DZ  
 S, D, SC, DZ  
 S, D, C, Z

## Level 2 BLAS

options	dim	b-width	scalar	matrix	vector	scalar	vector
xGEMV ( TRANS,	M, N,		ALPHA,	A, LDA,	X, INCX,	BETA,	Y, INCY )
xGEMV ( TRANS,	M, N,	KL, KU,	ALPHA,	A, LDA,	X, INCX,	BETA,	Y, INCY )
xHEMV ( UPLO,	N,		ALPHA,	A, LDA,	X, INCX,	BETA,	Y, INCY )
xHEMV ( UPLO,	N, K,		ALPHA,	A, LDA,	X, INCX,	BETA,	Y, INCY )
xHPMV ( UPLO,	N,		ALPHA,	AP,	X, INCX,	BETA,	Y, INCY )
xSYMV ( UPLO,	N,		ALPHA,	A, LDA,	X, INCX,	BETA,	Y, INCY )
xSBMV ( UPLO,	N, K,		ALPHA,	A, LDA,	X, INCX,	BETA,	Y, INCY )
xSPMV ( UPLO,	N,		ALPHA,	AP,	X, INCX,	BETA,	Y, INCY )
xTRMV ( UPLO, TRANS, DIAG,	N,		A, LDA,	X, INCX )			
xTRMV ( UPLO, TRANS, DIAG,	N, K,		A, LDA,	X, INCX )			
xTRMV ( UPLO, TRANS, DIAG,	N,		AP,	X, INCX )			
xTRSV ( UPLO, TRANS, DIAG,	N,		A, LDA,	X, INCX )			
xTRSV ( UPLO, TRANS, DIAG,	N, K,		A, LDA,	X, INCX )			
xTPSV ( UPLO, TRANS, DIAG,	N,		AP,	X, INCX )			
options	dim	scalar	vector	vector	matrix		
xGER (	M, N,	ALPHA,	X, INCX,	Y, INCY,	A, LDA )		
xGERU (	M, N,	ALPHA,	X, INCX,	Y, INCY,	A, LDA )		
xGERC (	M, N,	ALPHA,	X, INCX,	Y, INCY,	A, LDA )		
xHER ( UPLO,	N,	ALPHA,	X, INCX,		A, LDA )		
xHPR ( UPLO,	N,	ALPHA,	X, INCX,		AP )		
xHER2 ( UPLO,	N,	ALPHA,	X, INCX,	Y, INCY,	A, LDA )		
xHPR2 ( UPLO,	N,	ALPHA,	X, INCX,	Y, INCY,	AP )		
xSYR ( UPLO,	N,	ALPHA,	X, INCX,		A, LDA )		
xSPR ( UPLO,	N,	ALPHA,	X, INCX,		AP )		
xSYR2 ( UPLO,	N,	ALPHA,	X, INCX,	Y, INCY,	A, LDA )		
xSPR2 ( UPLO,	N,	ALPHA,	X, INCX,	Y, INCY,	AP )		

$y \leftarrow \alpha Ax + \beta y, y \leftarrow \alpha A^T x + \beta y, y \leftarrow \alpha A^H x + \beta y, A - m \times n$   
 $y \leftarrow \alpha Ax + \beta y, y \leftarrow \alpha A^T x + \beta y, y \leftarrow \alpha A^H x + \beta y, A - m \times n$   
 $y \leftarrow \alpha Ax + \beta y$   
 $y \leftarrow \alpha Ax + \beta y$   
 $y \leftarrow \alpha Ax + \beta y$   
 $y \leftarrow \alpha Ax + \beta y$   
 $y \leftarrow \alpha Ax + \beta y$   
 $y \leftarrow \alpha Ax + \beta y$   
 $x \leftarrow Ax, x \leftarrow A^T x, x \leftarrow A^H x$   
 $x \leftarrow Ax, x \leftarrow A^T x, x \leftarrow A^H x$   
 $x \leftarrow Ax, x \leftarrow A^T x, x \leftarrow A^H x$   
 $x \leftarrow A^{-1} x, x \leftarrow A^{-T} x, x \leftarrow A^{-H} x$   
 $x \leftarrow A^{-1} x, x \leftarrow A^{-T} x, x \leftarrow A^{-H} x$   
 $x \leftarrow A^{-1} x, x \leftarrow A^{-T} x, x \leftarrow A^{-H} x$

S, D, C, Z  
 S, D, C, Z  
 C, Z  
 C, Z  
 C, Z  
 S, D  
 S, D  
 S, D  
 S, D, C, Z  
 S, D, C, Z  
 S, D, C, Z  
 S, D, C, Z  
 S, D, C, Z  
 S, D, C, Z  
 S, D, C, Z

$A \leftarrow \alpha xy^T + A, A - m \times n$   
 $A \leftarrow \alpha xy^T + A, A - m \times n$   
 $A \leftarrow \alpha xy^H + A, A - m \times n$   
 $A \leftarrow \alpha xx^H + A$   
 $A \leftarrow \alpha xx^H + A$   
 $A \leftarrow \alpha xy^H + y(\alpha x)^H + A$   
 $A \leftarrow \alpha xy^H + y(\alpha x)^H + A$   
 $A \leftarrow \alpha xx^T + A$   
 $A \leftarrow \alpha xx^T + A$   
 $A \leftarrow \alpha xy^T + \alpha yx^T + A$   
 $A \leftarrow \alpha xy^T + \alpha yx^T + A$

S, D  
 C, Z  
 C, Z  
 C, Z  
 C, Z  
 C, Z  
 C, Z  
 S, D  
 S, D  
 S, D  
 S, D

## Level 3 BLAS

options	dim	scalar	matrix	matrix	scalar	matrix
xGEMM ( TRANSA, TRANSB,	M, N, K,	ALPHA,	A, LDA,	B, LDB,	BETA,	C, LDC )
xSYMM ( SIDE, UPLO,	M, N,	ALPHA,	A, LDA,	B, LDB,	BETA,	C, LDC )
xHEMM ( SIDE, UPLO,	M, N,	ALPHA,	A, LDA,	B, LDB,	BETA,	C, LDC )
xSYRK ( UPLO, TRANS,	N, K,	ALPHA,	A, LDA,		BETA,	C, LDC )
xHERK ( UPLO, TRANS,	N, K,	ALPHA,	A, LDA,		BETA,	C, LDC )
xSYR2X ( UPLO, TRANS,	N, K,	ALPHA,	A, LDA,	B, LDB,	BETA,	C, LDC )
xHER2X ( UPLO, TRANS,	N, K,	ALPHA,	A, LDA,	B, LDB,	BETA,	C, LDC )
xTRMM ( SIDE, UPLO, TRANSA,	DIAG, M, N,	ALPHA,	A, LDA,	B, LDB )		
xTRSM ( SIDE, UPLO, TRANSA,	DIAG, M, N,	ALPHA,	A, LDA,	B, LDB )		

$C \leftarrow \alpha op(A)op(B) + \beta C, op(X) = X, X^T, X^H, C - m \times n$   
 $C \leftarrow \alpha AB + \beta C, C \leftarrow \alpha BA + \beta C, C - m \times n, A = A^T$   
 $C \leftarrow \alpha AB + \beta C, C \leftarrow \alpha BA + \beta C, C - m \times n, A = A^H$   
 $C \leftarrow \alpha AA^T + \beta C, C \leftarrow \alpha A^T A + \beta C, C - n \times n$   
 $C \leftarrow \alpha AA^H + \beta C, C \leftarrow \alpha A^H A + \beta C, C - n \times n$   
 $C \leftarrow \alpha AB^T + \delta BA^T + \beta C, C \leftarrow \alpha A^T B + \delta B^T A + \beta C, C - n \times n$   
 $C \leftarrow \alpha AB^H + \delta BA^H + \beta C, C \leftarrow \alpha A^H B + \delta B^H A + \beta C, C - n \times n$   
 $B \leftarrow \alpha op(A)B, B \leftarrow \alpha Bop(A), op(A) = A, A^T, A^H, B - m \times n$   
 $B \leftarrow \alpha op(A^{-1})B, B \leftarrow \alpha Bop(A^{-1}), op(A) = A, A^T, A^H, B - m \times n$

S, D, C, Z  
 S, D, C, Z  
 C, Z  
 S, D, C, Z  
 C, Z  
 S, D, C, Z  
 C, Z  
 S, D, C, Z  
 S, D, C, Z  
 S, D, C, Z

# BLAS summary

## Basic Linear Algebra Subroutines

Name	Description	Examples
Level-1 BLAS	Vector Operations	$C = \sum X_i Y_i$
Level-2 BLAS	Matrix-Vector Operations	$B_i = \sum_k A_{ik} X_k$
Level-3 BLAS	Matrix-Matrix Operations	$C_{ij} = \sum_k A_{ik} B_{kj}$

# Why BLAS are important ?

- Because the BLAS are **efficient**, **portable**, **parallel**, and **widely available**, they are commonly used in the development of high quality linear algebra software.
- Performance of lot of applications depends a lot on the performance of the underlying BLAS
- Lot of applications include ML/DL stuff as well....
  - <https://petewarden.com/2015/04/20/why-gemm-is-at-the-heart-of-deep-learning/>

# Standardization: BLAS example

- Each BLAS Subroutines have a standardized layout
- BLAS is documented in the source code
- Man pages exist
- Vendor supplied docs
- Different BLAS implementations have the same calling sequence

```
SUBROUTINE DGEMM ( TRANSA, TRANSB, M, N, K, ALPHA, A, LDA, B, LDB,
+ BETA, C, LDC )
*
* .. SCALAR ARGUMENTS ..
* CHARACTER*1      TRANSA, TRANSB
* INTEGER          M, N, K, LDA, LDB, LDC
* DOUBLE PRECISION ALPHA, BETA
* .. ARRAY ARGUMENTS ..
* DOUBLE PRECISION A( LDA, = ), B( LDB, = ), C( LDC, = )
*
* ..
*
* PURPOSE
* =====
*
* DGEMM PERFORMS ONE OF THE MATRIX-MATRIX OPERATIONS
*
*   C := ALPHA*OP( A )*OP( B ) + BETA*C,
*
* WHERE OP( X ) IS ONE OF
*
*   OP( X ) = X   OR   OP( X ) = X'.
*
* ALPHA AND BETA ARE SCALARS, AND A, B AND C ARE MATRICES, WITH OP( A )
* AN M BY K MATRIX, OP( B ) A K BY N MATRIX AND C AN M BY N MATRIX.
*
* PARAMETERS
* =====
*
* TRANSA - CHARACTER*1,
* ON ENTRY, TRANSA SPECIFIES THE FORM OF OP( A ) TO BE USED IN
* THE MATRIX MULTIPLICATION AS FOLLOWS:
*
*   TRANSA = 'N' OR 'n',  OP( A ) = A,
*   TRANSA = 'T' OR 't',  OP( A ) = A',
*   TRANSA = 'C' OR 'c',  OP( A ) = A'.
*
* UNCHANGED ON EXIT.
*
* TRANSB - CHARACTER*1,
* ON ENTRY, TRANSB SPECIFIES THE FORM OF OP( B ) TO BE USED IN
* THE MATRIX MULTIPLICATION AS FOLLOWS:
*
*   TRANSB = 'N' OR 'n',  OP( B ) = B,
*   TRANSB = 'T' OR 't',  OP( B ) = B'.
```

# Vendor/Optimized BLAS libraries

- ACML
  - The AMD Core Math Library, supporting the AMD processors
- ATLAS
  - Automatically Tuned Linear Algebra, an open source implementation of BLAS APIs for C and Fortran 77
- Intel MKL
  - The Intel Math Kernel Library supporting x86 32-bits and 64-bits. Includes optimizations for Intel Pentium, Core and Intel Xeon CPUs and Intel Xeon Phi; support for Linux, Windows and Mac OS X
- cuBLAS
  - Optimized BLAS for NVIDIA based GPU cards
- ESSL
  - IBM's Engineering and Scientific Subroutine Library, supporting the PowerPC architecture under AIX and Linux
- GotoBLAS
  - Kazushige Goto's BSD-licensed implementation of BLAS, tuned in particular for Intel, VIA Nanoprocessor, AMD Opteron
- OpenBLAS
  - Optimized BLAS based on Goto BLAS hosted at GitHub, supporting Intel platform and other
- And many others...

# What about C/C++ program?

- BLAS routines are Fortran-style, when calling them from C language programs, follow the Fortran-style calling conventions:
  - Pass variables by address, not by value.
  - Store your data in Fortran style, that is, column-major rather than row-major order.
- Be aware that because the Fortran language is case-insensitive, the routine names can be both upper-case or lower-case, with or without the trailing underscore.
- For example, the following names are equivalent:
  - dgemm, DGEMM, dgemm\_, and DGEMM\_



# Use CBLAS !

- C-style interface to the BLAS routines

[www.netlib.org/blas/blast-forum/cblas.tgz](http://www.netlib.org/blas/blast-forum/cblas.tgz)

- You can call CBLAS routines using regular C- style calls.
- The header file specifies enumerated values and prototypes of all the functions.
- Details and examples here:

<https://software.intel.com/en-us/mkl-tutorial-c-multiplying-matrices-using-dgemm>

# Q parameter: aka computational efficiency...

**Table 2: Basic Linear Algebra Subroutines (BLAS)**

Operation	Definition	Floating point operations	Memory references	$q$
<b>saxpy</b>	$y_i = \alpha x_i + y_i, i = 1, \dots, n$	$2n$	$3n + 1$	$2/3$
<b>Matrix-vector mult</b>	$y_i = \sum_{j=1}^n A_{ij}x_j + y_i$	$2n^2$	$n^2 + 3n$	$2$
<b>Matrix-matrix mult</b>	$C_{ij} = \sum_{k=1}^n A_{ik}B_{kj} + C_{ij}$	$2n^3$	$4n^2$	$n/2$

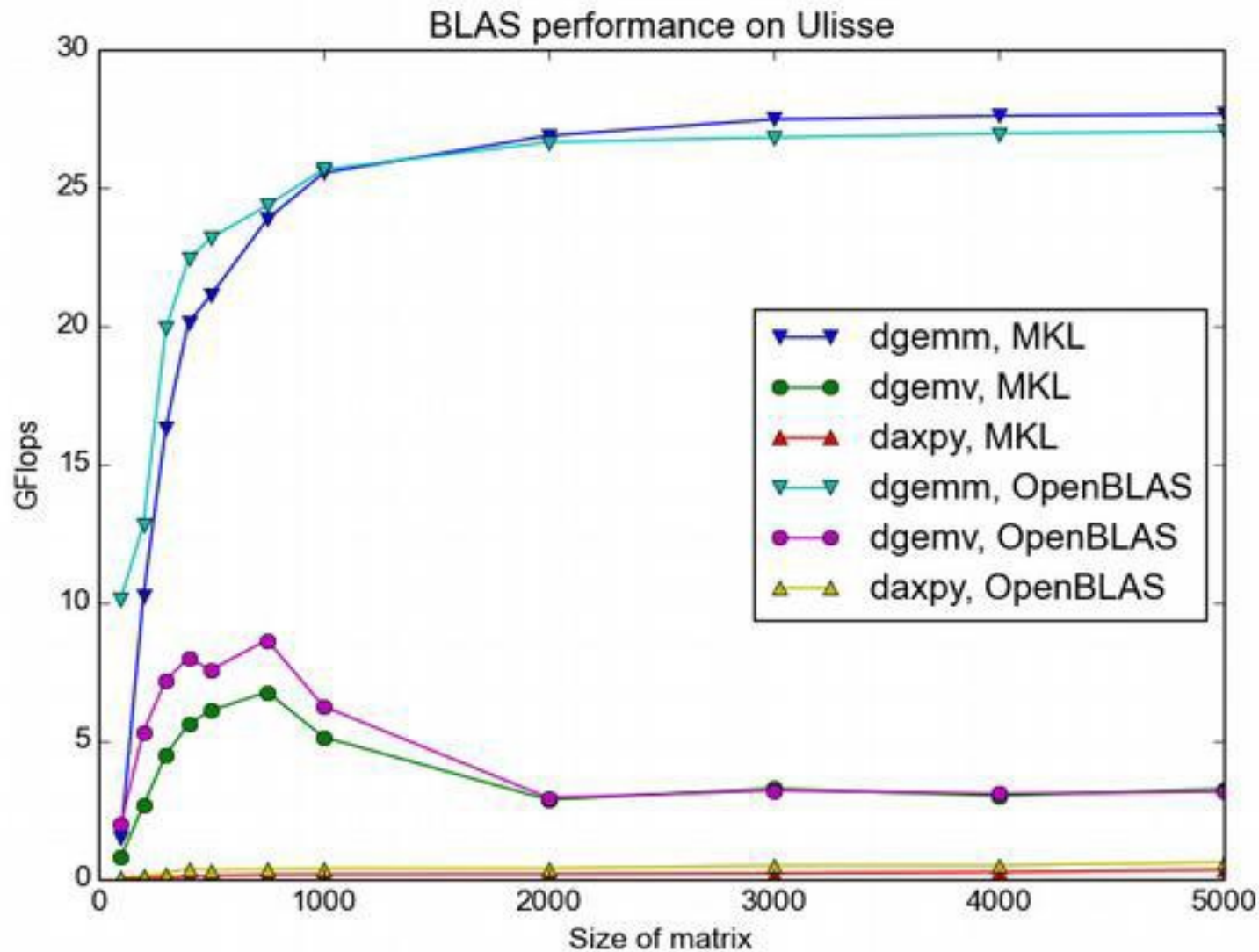
The parameter  $q$  is the ratio of flops to memory references.  
Generally:

1. Larger values of  $q$  maximize useful work to time spent moving data.
2. The higher the level of the BLAS, the larger  $q$ .

# It follows...

- BLAS1 are **memory bounded** !(for each computation a memory transfer is required)
- BLAS2 are not so memory bounded (can have good performance on super-scalar architecture)
- BLAS3 can be very efficient on super-scalar computers because **not memory bounded**

# BLAS performance on SandyBridge (1core)



# Proposed exercise/tutorial

- Create the same graph for ORFEO cores using MKL and OpenBLAS
- STEPS:
  - Install OpenBLAS libraries in your directory
  - Write a small program to call the three routines
    - Dgemm/dgemv/daxypi
  - Write a script to collect all sizes of interest
  - Make nice plots

# Linking optimized libraries...

- OpenBLAS:

- It comes with cblas bundled in so no problem with C/C++
- Automatically includes lapack reference implementation
- Compilation is straightforward:

```
gcc -o test test.c -I  
/your_path/OpenBLAS/include/  
-L/your_path/OpenBLAS/lib -lopenblas
```

- MKL

- Generally complex and highly dependent on version and/or HW/SW implementation

- <https://software.intel.com/en-us/articles/intel-mkl-link-line-advisor>

# Are these libraries multithreaded ?

- MKL
  - Both sequential and multithreaded version available

Intel® oneAPI Math Kernel Library (oneMKL) Link Line Advisor v6.13

Reset

Select Intel® product:	Intel(R) Parallel Studio XE 2020 ▼
Select OS:	Linux* ▼
Select compiler:	GNU C/C++ ▼
Select architecture:	Intel(R) 64 ▼
Select dynamic or static linking:	Dynamic ▼
Select interface layer:	C/Fortran API with 32-bit integer ▼
Select threading layer:	OpenMP threading ▼
Select OpenMP library:	<Select threading> OpenMP threading Sequential TBB threading
Enable OpenMP offload feature to GPU:	

# MKL: how to control number of threads?

- OpenMP threading ? → OMP\_NUM\_THREADS
- Other threading ? → MKL\_NUM\_THREADS
- Define yourself the number of threads:
  - Place `mkl_set_num_thread( N )` routine in your code.
- All MKL routines call takes precedence over any environment variables. and MKL environment Variables will take precedence over the OpenMP\* environments.

More details

here: <https://software.intel.com/content/www/us/en/develop/articles/recommended-settings-for-calling-intel-mkl-routines-from-multi-threaded-applications.html>



# OpenBLAS:

- By default : multithreaded version, maximum number of threads established by cores available on the machine when compilation is performed;
- In our case:

```
OpenBLAS build complete. (BLAS CBLAS LAPACK LAPACKE)
```

```
OS          ... Linux
Architecture ... x86_64
BINARY      ... 64bit
C compiler  ... GCC (cmd & version : cc (GCC) 4.8.5 20150623 (Red Hat 4.8.5-39))
Fortran compiler ... GFORTRAN (cmd & version : GNU Fortran (GCC) 9.3.0)
Library Name ... libopenblas_haswellp-r0.3.13.a (Multi-threading; Max num-threads is 48)
```

```
To install the library, you can run "make PREFIX=/path/to/your/installation install".
```

# OpenBLAS: caveat

- If your application is already multi-threaded, **it will conflict** with OpenBLAS multi-threading. Thus, you must set OpenBLAS to use single thread as following.
  - export OPENBLAS\_NUM\_THREADS=1 in the environment variables.
  - call `openblas_set_num_threads(1)` in the application on runtime.
- You can compile the library itself in sequential mode:
  - make `USE_THREAD=0 USE_LOCKING=1`\*(see comment below)
  - If your application is parallelized by OpenMP, please build OpenBLAS with `USE_OPENMP=1`

\* the thread management provided by OpenMP is not sufficient to prevent race conditions when OpenBLAS was built single-threaded by `USE_THREAD=0` and there are concurrent calls from multiple threads to OpenBLAS functions. In this case, it is vital to also specify `USE_LOCKING=1`