# BLAS

Christoph Ortner

Stef Salvini

# The *BLASics*

- **B**asic **L**inear **A**lgebra **S**ubroutines
  - Building blocks for more complex computations
  - Very widely used

- Level means "number of operations"
  - Level 1: vector-vector operations – O($n$) operations
  - Level 2: matrix vector operations – O($n^2$) operations
  - Level 3: Matrix-Matrix operations  - O($n^3$) operations

- A *Flop* is any numerical operation
  - Adds, Mults, divisions, square roots (!!!!), etc
    - Of course divisions & square roots are more expensive …
  - Loads/stores are not taken into account (history … )

- BLAS provide a good basis to understand performance issues

# A Fistful of Flops



π BLAS take off with **Vector** processors (70s – 80s)

π Level 1 first, then level 2 BLAS

  ν Encapsulate expert knowledge

  ν Efficient building blocks

  ν "Local" optimisation of code enough to increase performance

# Level 1 BLAS

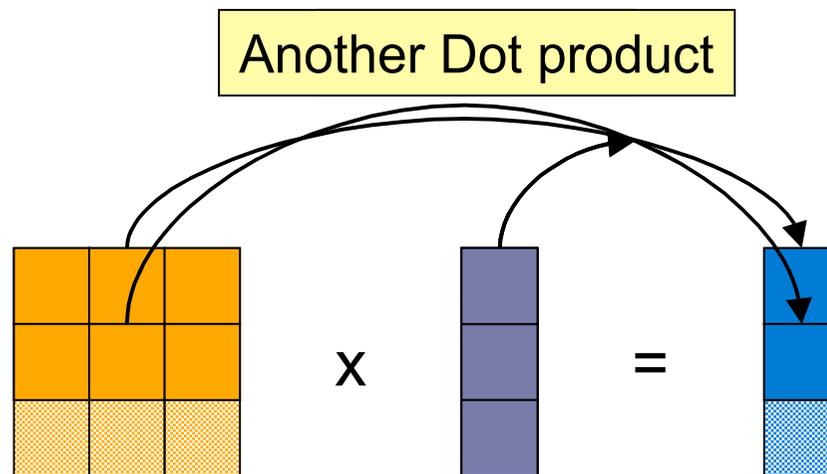- π  O($n$) operands (I.e. load/stores), operations O($n$) (flops)

- π  Vector operations (loved by vector processors)
  - ν  Ratio between load/stores and operations: O(1)
  - ν  E.g. "axpy" : $\alpha\, x + y \lozenge y$

- π  Reduction operations (hated by vector processors)
  - ν  Ratio between load/stores and operations: O($n$)
  - ν  E.g. dot product: $\alpha = x^T\, y$

- π  Available:
  - ν  Single & double precision, real and complex
    - π  Names beginning with S, D, C and Z, respectively
    - π  Axpy: SAXPY, DAXPY, CAXPY, ZAXPY
    - π  Dot Products (SDOT, DDOT, CDOTC, ZDOTC)

# Level 2 BLAS

- π  $O(n^2)$ operands, $O(n^2)$ operations
- π  Performance can be understood in terms of Level 1 cache
- π  Matrix-vector product, matrix updates, solution of a triangular system of equations, etc
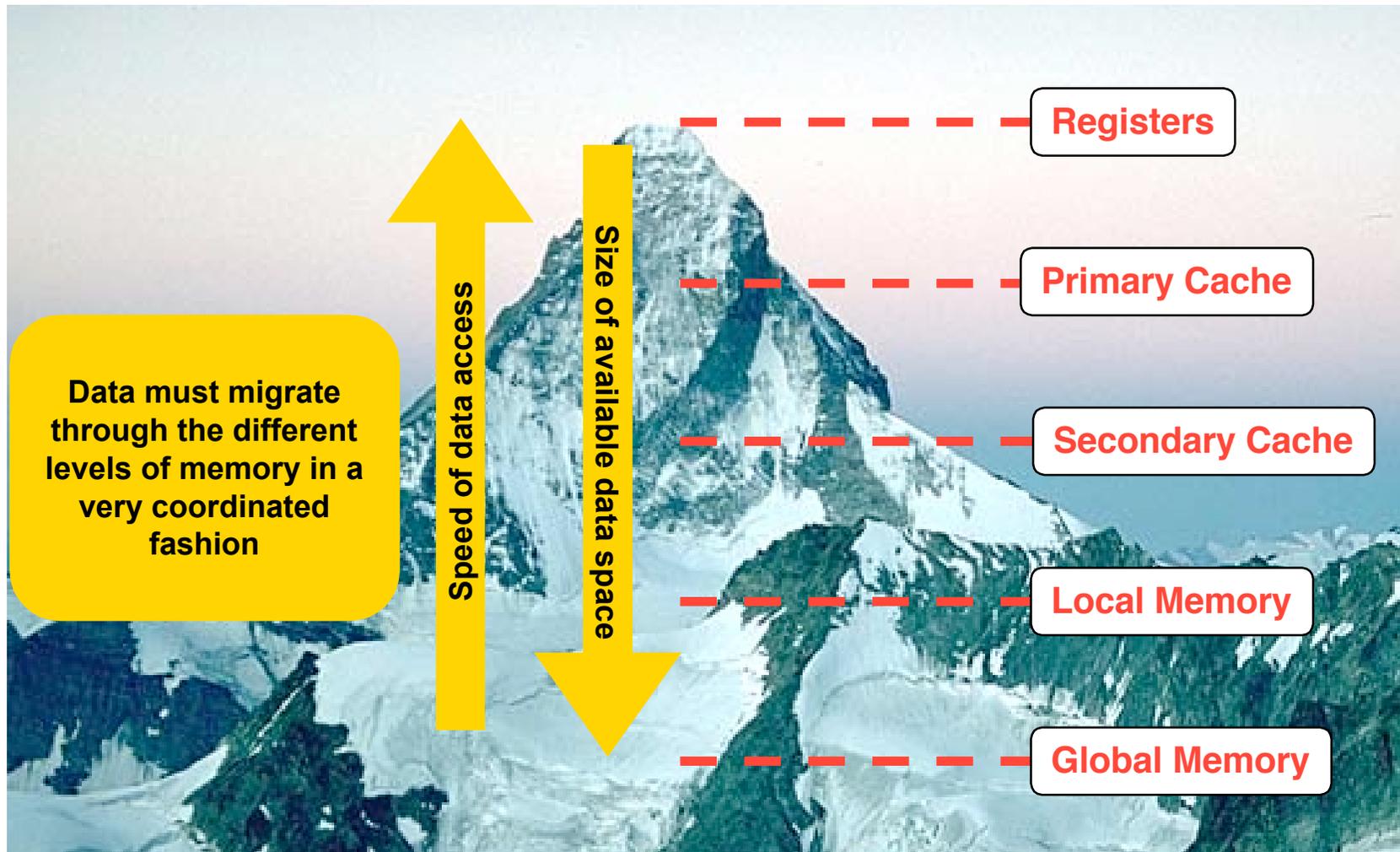
Another Dot product

X  =

# Superscalar takes over

- π Technology dictated by
  - ν Cost
  - ν Widespread use
  - ν Relatively small HPC market

- π Superscalar here means more than one operation per cycle

- π All supercalar architecture (give-or-take)
  - ν No direct access to memory
    - π Hierarchical memory layout
    - π Use of caches to make use of any data locality
  - ν **Rule-of-thumb for performance:**
    - π *"Thou shalt not have as many operands as operations"*
    - π In fact: poor performance of Level 1 and 2 BLAS (sometimes horrifyingly so)
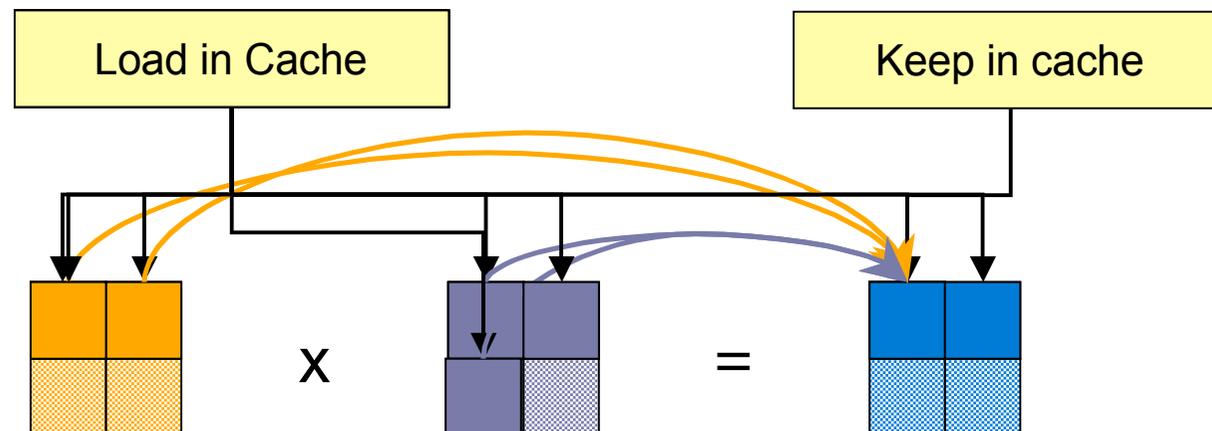    - π Poor performance for indirect addressing
    - π FFTs very difficult

# The Data View of an SS Architecture



Data must migrate through the different levels of memory in a very coordinated fashion

Speed of data access

Size of available data space

Registers

Primary Cache

Secondary Cache

Local Memory

Global Memory

# For A Few Flops More

$\pi$ Level 3 BLAS (matrix-matrix operation) to the rescue!

$\pi$ Why Level 3 BLAS are *good* (example of matrix-matrix product)
  - $\nu$ Use blocked algorithms to maximise cache reuse
  - $\nu$ $O(b^2)$ loads/stores – $O(b^3)$ flops
  - $\nu$ Enough operations to "hide" memory traffic costs

# The Good the Bad and the Ugly



- π Lots of Packages depend on and benefit from BLAS
  - ν LAPACK (Linear Algebra)
  - ν Many Sparse Solvers (using local dense blocks of sparse matrices, such as SuperLU, MUMPS, etc)
- π A Myth
  - ν BLAS are parallelised by vendors, hence all LAPACK etc is parallel and scalable – NOT TRUE!
    - π Level 1 BLAS: NEVER parallelised
    - π Level 2 BLAS: SOME parallelised
    - π Level 3 BLAS: ALL parallelised
- π Most codes do not contain the nice packed aligned data that BLAS require (indirect addressing on SS architectures very tough!)
- π What about SSE & SSE2 on Intel & AMD architectures
  - ν They are for multimedia!
    - π Pack several words (numbers) in register
    - π Operate simultaneous on all words in register
    - π Operations crossing low & high in register very expensive! (What about complex numbers: well, they do not exist for vendors)

# Sparse BLAS

**Slide Intentionally left (almost) Blank**

Zillions of sparse formats

Efforts to generate sparse BLAS automatically (performance poor – indirect addressing)

# Getting the BLAS

- π "Model" BLAS
  - ν Model implementation in Fortran
  - ν No optimization in source
  - ν Some compilers can block Level 3-BLAS approaching level of more sophisticated implementations (only DGEMM)
  - ν C interface is available
- π Vendor BLAS
  - ν Hand-optimized by vendors (IESSL/IBM, MKL/Intel, ACML/AMD, …)
  - ν Achieves highest performance on vendors' platforms.
  - ν **YOU SHOULD USE THIS!**
- π ATLAS
  - ν "Automatically Tuned Linear Algebra Software"
  - ν Brute force optimization
    - π trying out all possible combinations of memory layout, loop reordering, etc.
  - ν Competitive performance on Level 3 BLAS
  - ν Can be generated for virtually all platforms

# The "Mythical" Goto BLAS

- π  BLAS designed by **Kazushige Goto**
- π  Optimizes all memory traffic in a very clever way
- π  Currently beats most commercial libraries
- π  Only few non-threaded BLAS
- π  http://www.tacc.utexas.edu/resources/software/

# Measuring Performance

π  Performance is measured in Mflops/s

π  E.g.: multiplication of two square N x N matrices (DGEMM)

  ν  $N^3$ multiplications and

  ν  $N^3$ additions

  ν  = 2 $N^3$ flops

π  t seconds for m `dgemm` calls gives

$$\frac{2n^3 m}{10^6 t} \text{ Mflops}$$

# Benchmarks

π 1. Loop to determine number **m** of calls to run for **T** seconds (say **T = 1.0**)

  ν This loop does many timings which is a significant overhead

π 2. Time **m** calls

π 3. Repeat step 2 several times and take best timing

# The henrici Systems

- $\pi$   Intel® Xeon™ CPU 3.20 GHz
- $\pi$   512 KB L2 cache
- $\pi$   2 Processors with hyperthreading
- $\pi$   2GB main memory

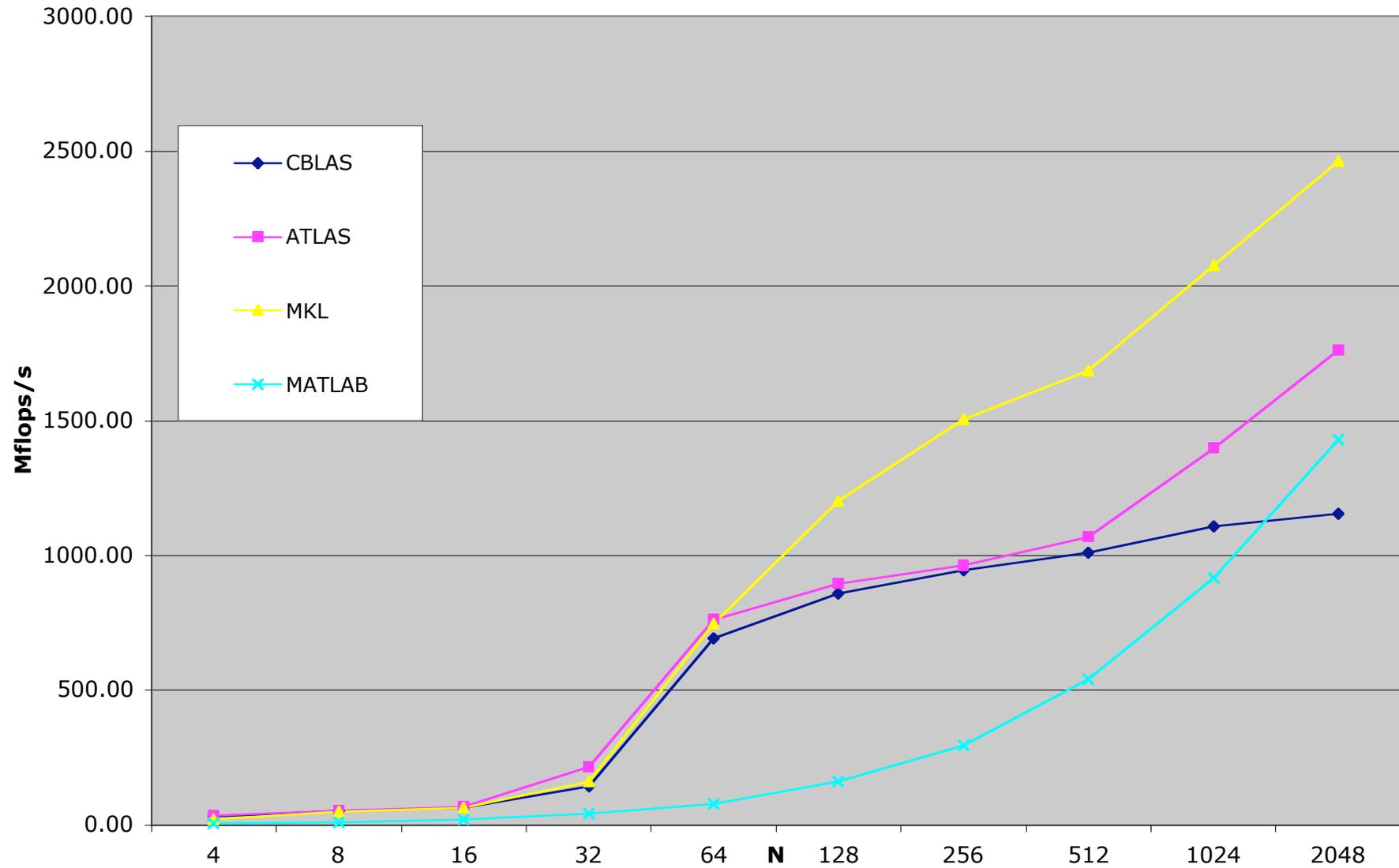- $\pi$   (theoretical **serial** peak peformance: 6400 MFLOPs/s)

# DDOT

π  `a = x' * y;`

π  double **cblas_ddot**(const int N, const double *X,
        const int incX, const double *Y, const int incY);

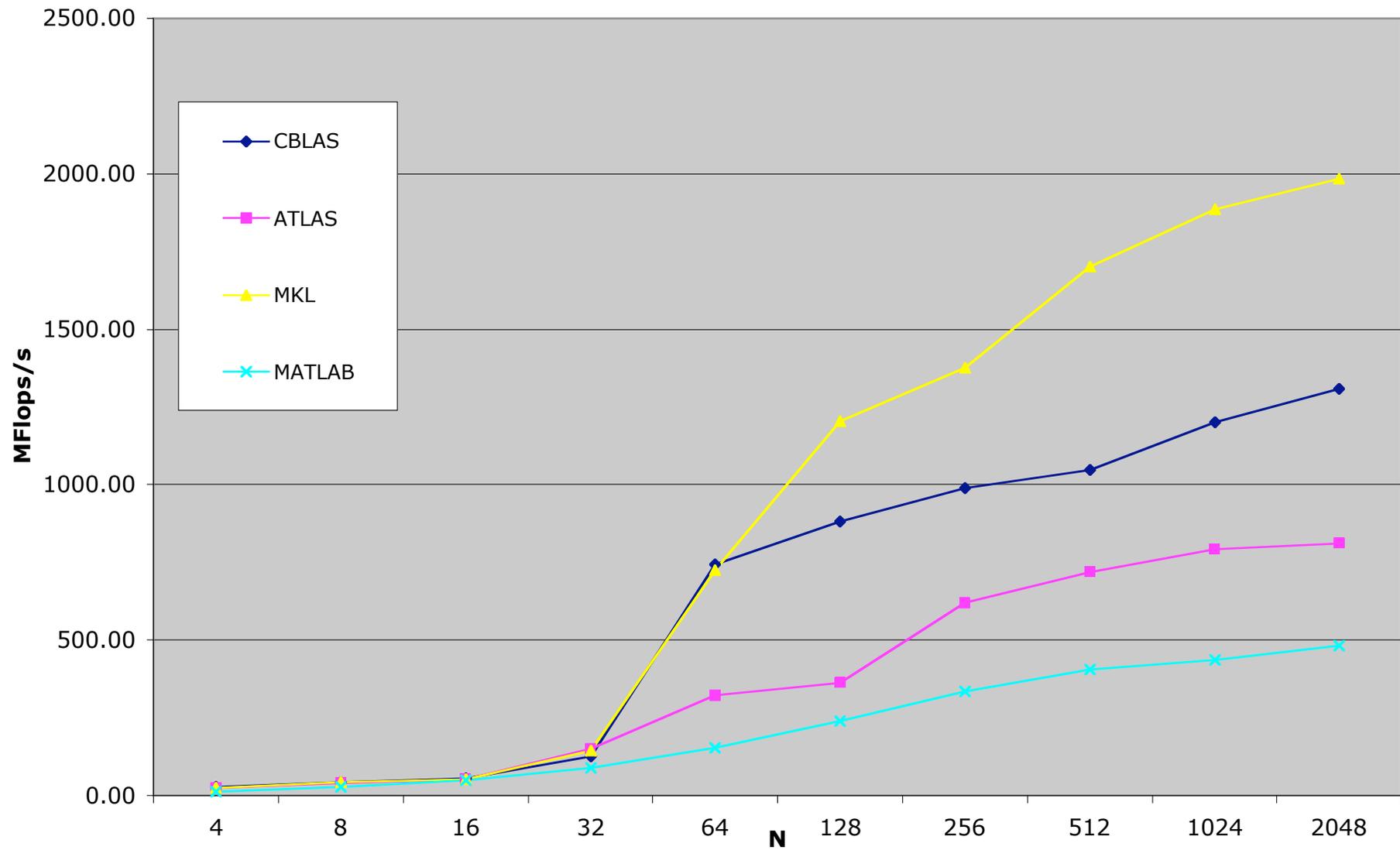π  2 N operations + 2 N memory accesses

# DDOT Performance

# DAXPY

π **`y = y + alpha * x;`**

π void **`cblas_daxpy`**`(const int N, const double alpha,`
`                const double *X, const int incX,`
`                double *Y, const int incY);`

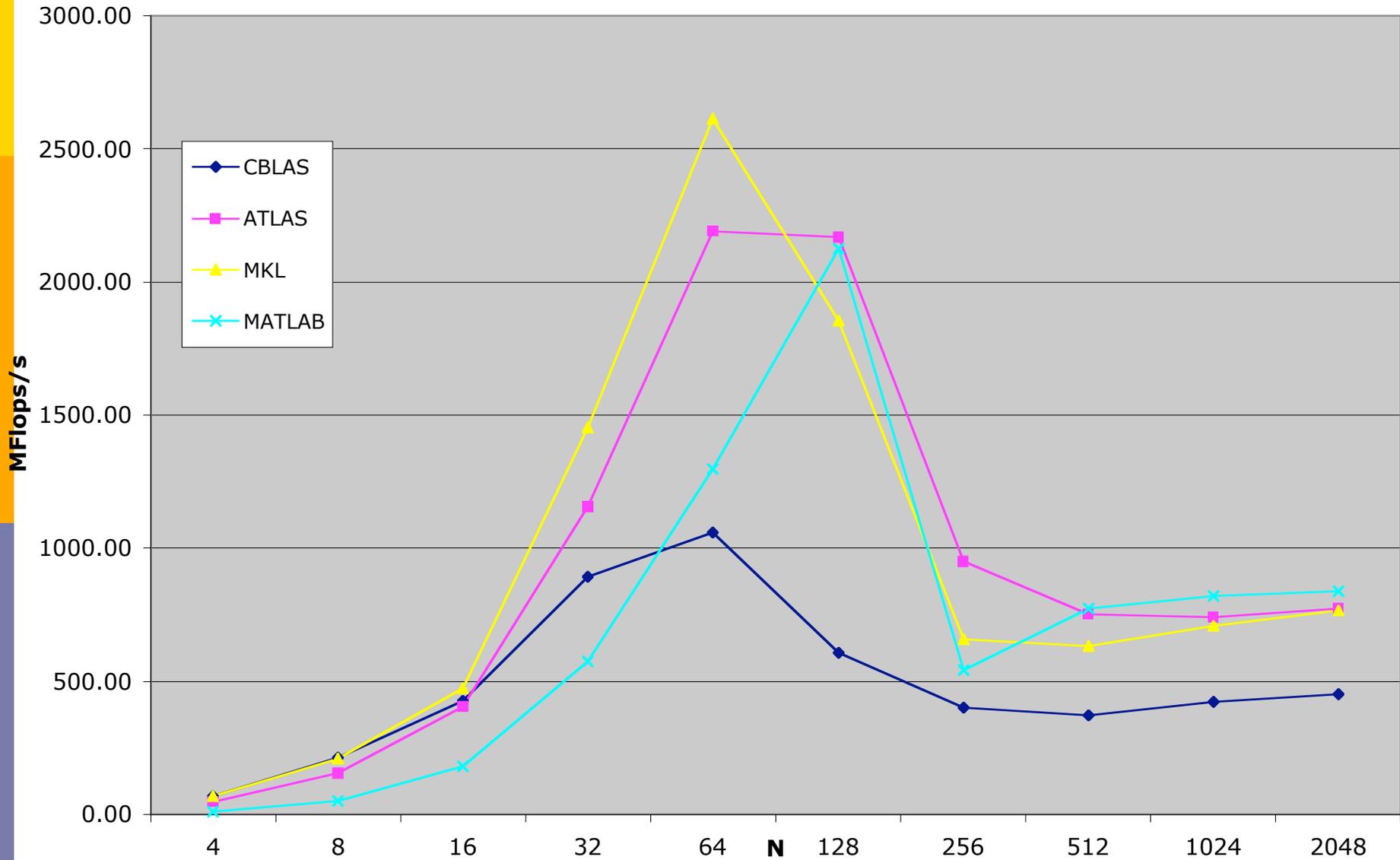π 2 N operations + 3 N memory accesses

# DAXPY Performance

# DGEMV

π **y = alpha * A * x + y;**

π void **cblas_dgemv**(const enum CBLAS_ORDER Order,
         const enum CBLAS_TRANSPOSE TransA, const int M,
         const int N, const double alpha, const double *A,
         const int lda, const double *X, const int incX,
         const double beta, double *Y, const int incY);

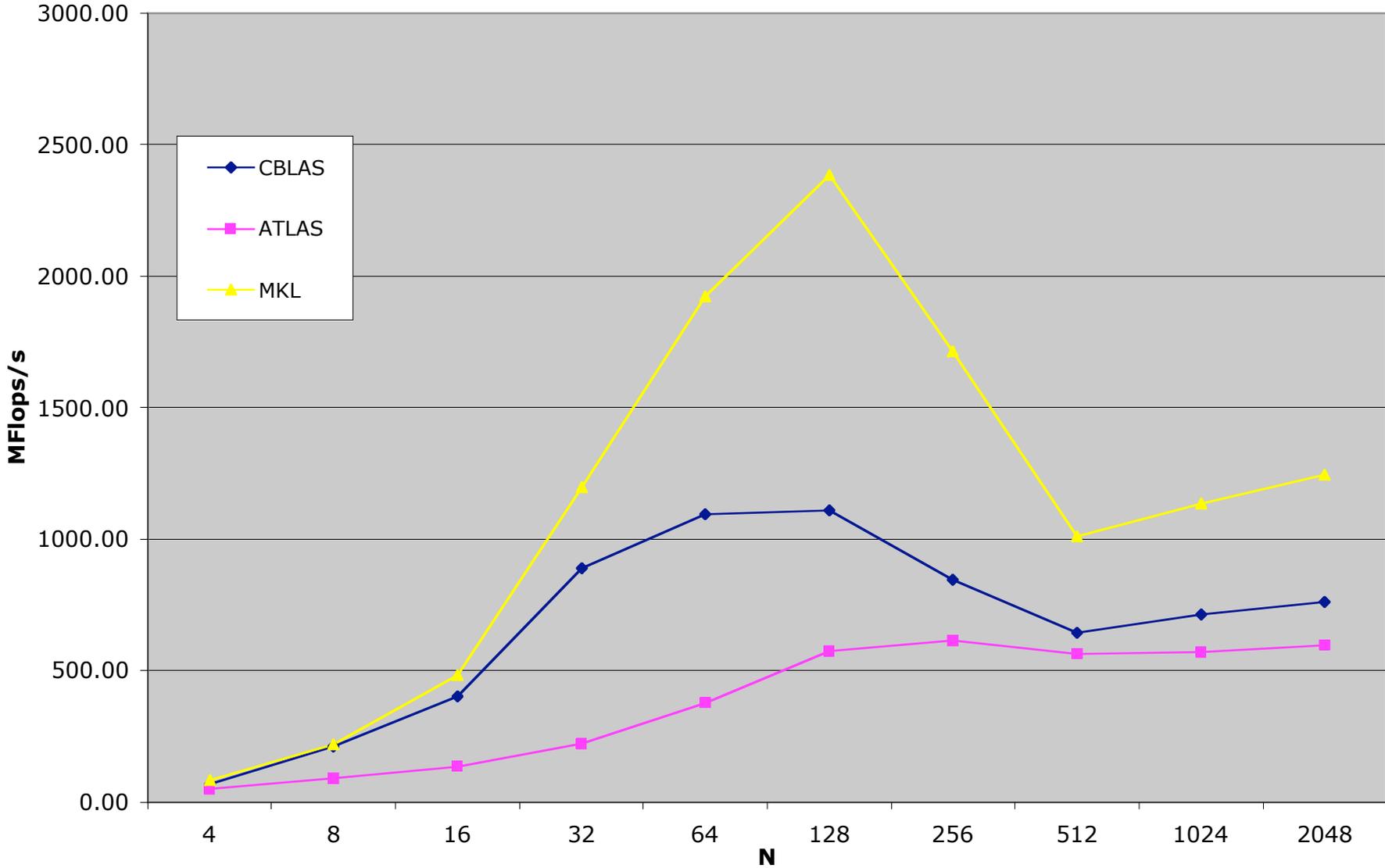π 2 $N^2$ operations + $N^2$ memory accesses

# DGEMV Performance

# DSYMV

π  **y = alpha * A * x + y;     % A symmetric**

π  void **cblas_dsymv**(const enum CBLAS_ORDER Order,
          const enum CBLAS_UPLO Uplo, const int M,
        const int N, const double alpha, const double *A,
        const int lda, const double *X, const int incX,
        const double beta, double *Y, const int incY);

π  2 $N^2$ operations + $N^2$ / 2 memory accesses
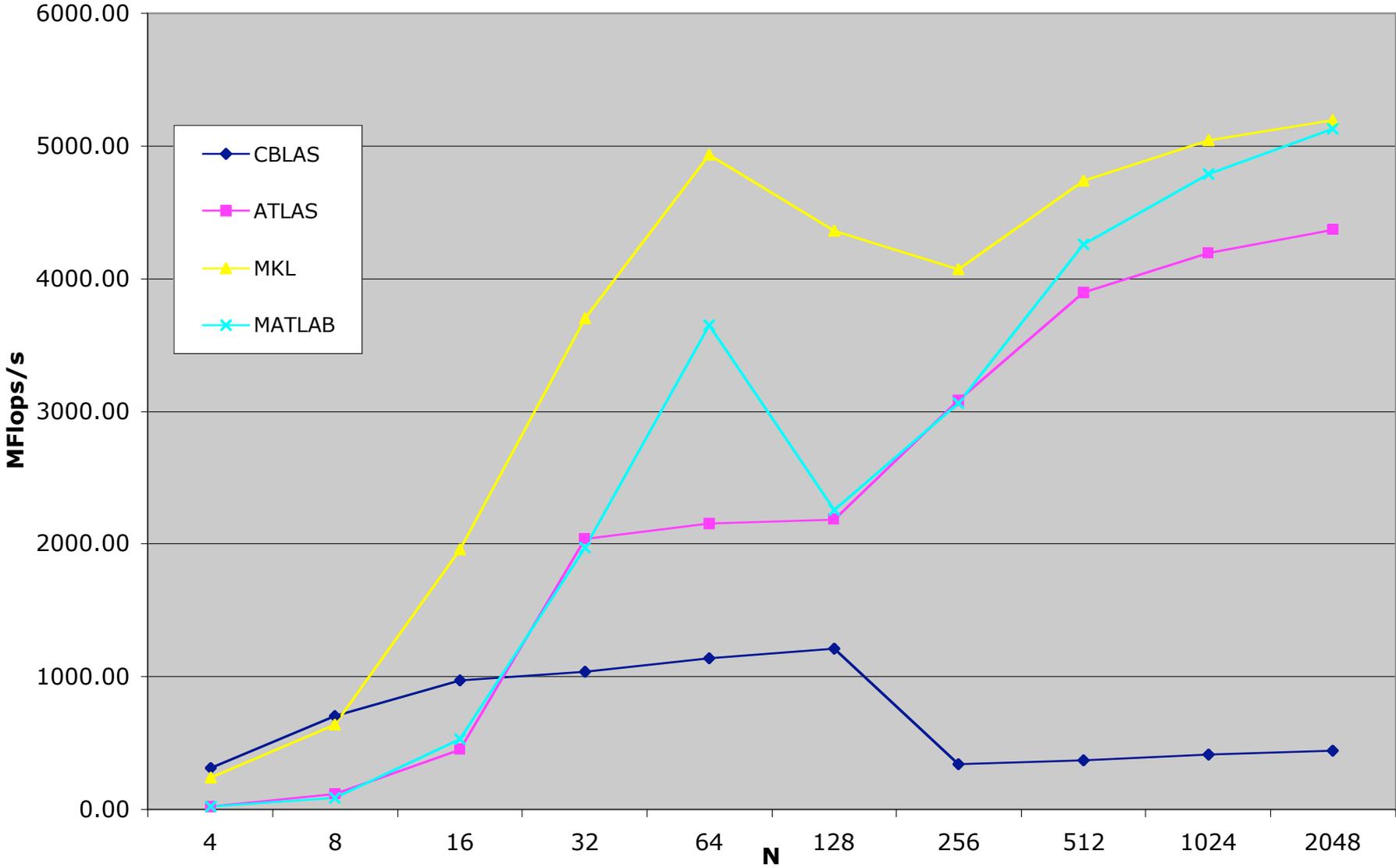
# DSYMV Performance

# DGEMM

$\pi$   **`C = alpha * A * B + C;`**

$\pi$   void **cblas_dgemm**(const enum CBLAS_ORDER Order,
    const enum CBLAS_TRANSPOSE TransA,
    const enum CBLAS_TRANSPOSE TransB,
    const int M, const int N, const int K,
    const double alpha, const double *A, const int lda,
    const double *B, const int ldb, const double beta,
    double *C, const int ldc);

$\pi$   $2\ N^3$ operations + $3\ N^2$ memory accesses

# DGEMM Performance

# Blitz++

- π C++ Array Library
- π Fully object oriented and templated
- π Supports operator overloading
- π Level 1 BLAS via **Expression Templates**
- π 'Lightweight' classes for small vectors or matrices
- π Current Version: 0.9
  http://www.oonumerics.org/blitz/
- π Similar Package: uBLAS (boost Library)

# Expression Templates

π   daxpy: **y = y + a x**

π   Expression Templates:
**for i = 1:N
    y(i) = y(i) + a * x(i);**

π   Expression Templates create such loop for any possible operation of any possible Object (e.g. Arrays of short vectors for explicit codes)

# DDOT Performance

# DAXPY Performance

# TinyVector, TinyMatrix

π    Fully Templated: e.g.

    `TinyVector<double, 3> x, y, z;`
    creates 3D-vectors of type double.
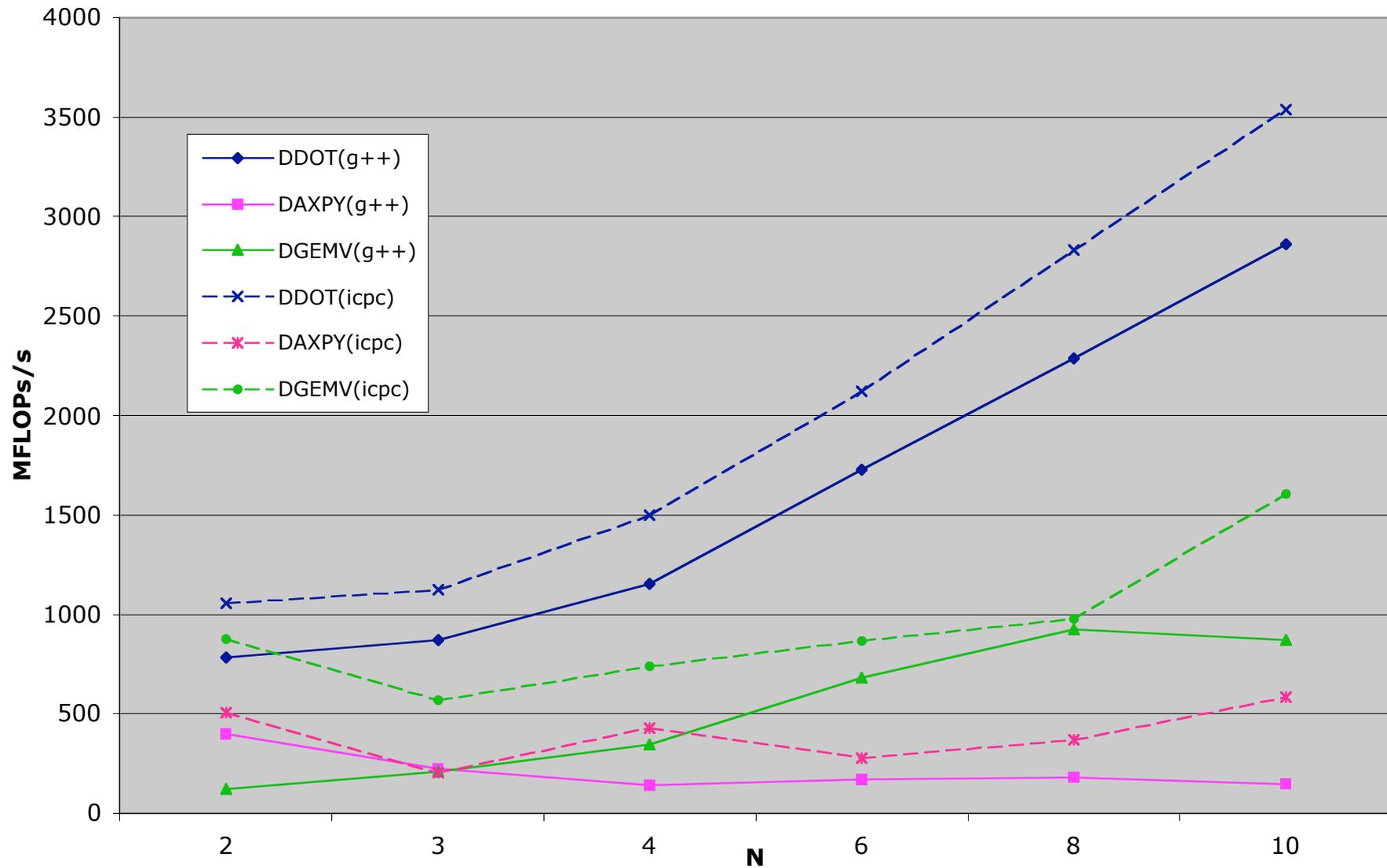
π    Vector length known at compile time.

π    Fully unroll all loops: `z = x + y;` becomes

```
z(1) = x(1) + y(1);
z(2) = x(2) + y(2);
z(3) = x(3) + y(3);
```

# TinyVector/Matrix Performance

# Conclusion

- π BLAS are essential items in scientific computing
- π Standardized interface to basic matrix and vector operations
- π Highly optimized BLAS are available
- π Many applications/packages/libraries depend dramatically on BLAS (e.g. dense and sparse solvers, both direct and iterative)

- π **We recommend you use VENDOR BLAS!**