

P1417: Lessons learned for linear algebra library standardization

Mark Hoemmen (mhoemme@sandia.gov)

WG21, Kona, Feb. 2019

Outline

- Over 40 years of standard linear algebra interfaces
- P1417 shares some history & lessons learned
 - Incomplete & can never be fully complete
 - Meant to grow & include your contributions
- This talk: Lessons from BLAS standardization
 - Standardize in layers, bottom up
 - Layer based on expertise & performance portability
 - Start w/ fundamentals least specific to linear algebra
- Lowest-level fundamentals include
 - Multidimensional arrays
 - Multidimensional parallel iteration
 - Explicit SIMD types & ops

Basic Linear Algebra Subprograms



- Standard published 2002
 - 1995 workshop
 - 1996-99 meetings 3x/year (w/ minutes online)
 - Fortran & C interfaces
 - Dense matrix & vector ops(*)
- Developed in levels (1,2,3):
 - Vector-vector (BLAS 1): 1979 (least data reuse)
 - Matrix-vector (BLAS 2): 1988 (constant factor more reuse)
 - Matrix-matrix (BLAS 3): 1990 (most data reuse)

(Fortran) BLAS quick reference:

<http://www.netlib.org/blas>

(See also Jack Dongarra's oral history)

(*) Sparse etc. too; ask me later

BLAS 1-3 coevolved w/ computers

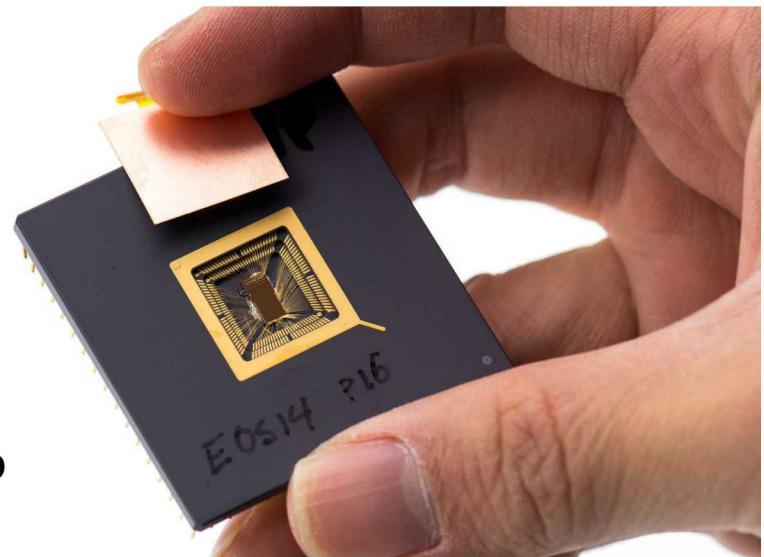


Seymour Cray w/
Cray 1, circa 1976
(when LINPACK
funding started)

Vector (Cray, NEC)

- Low flop/byte ratio
- Favor long, dependence-free loops & regular data access
- BLAS 1 target

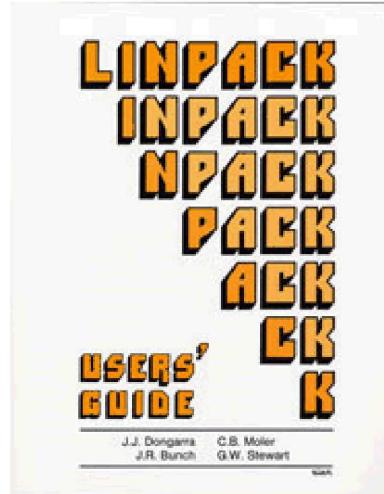
Yunsup Lee holding
RISC V prototype, 2013



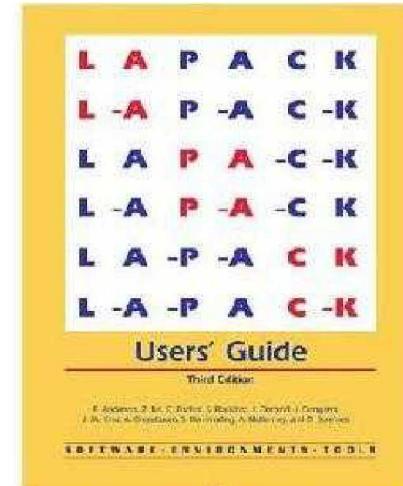
Cache based

- “Killer micros” (see 1991 NYT article)
- High flop/byte ratio
- Favor data reuse
- BLAS 3 target

BLAS codesigned w/ algorithms



- LINPACK library: 1979
 - Linear systems (LU, Cholesky)
 - Linear least squares (QR)
 - General dense, symmetric, & banded
 - Designed to use BLAS (1), for good performance on many different computers
- LAPACK: 1990
 - Combines functionality of LINPACK + EISPACK ({eigen,singular} value problems)
 - “Coreleased” w/ BLAS 3, w/ common authors
 - Algorithms that better exploit data reuse
 - BLAS 3 designed for those algorithms



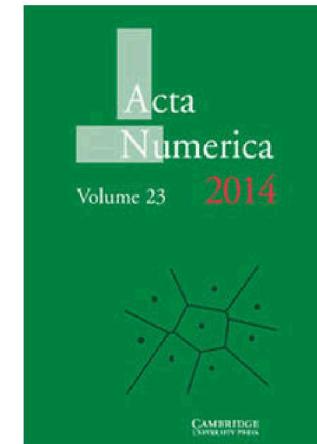
Imitate BLAS' std-ization approach



- How it got standardized, not necessarily the interface itself
 - I do not necessarily advocate a “C++ BLAS binding”
 - Fullest development of Ranges may reduce BLAS 1’s value
 - BLAS optimized for single large problems, not small & many
 - Larger, different customer base, solving different problems
- Process: Standardize in layers, bottom up
- Layer by expertise
 - BLAS: hardware / performance; LAPACK: numerical analysis
 - Kazushige Goto (talented BLAS optimizer) is self-taught [1]
 - Numerical analysis not typical background for Std Lib developers
- Layer for performance portability
 - Value added in multiple implementations, tuned for specific hardware

Counter & rebuttal

- Counter: BLAS & LAPACK codesigned by numerical analysts
 - Will future algorithms need different fundamental operations?
 - Would a bottom-up approach risk overspecializing the lowest level?
 - Radical changes in computer architecture?
- Rebuttal
 - Fundamentals not specific to linear algebra
 - Last 2 decades: Many theoretical results on algorithm optimality (matrix & tensor), so unlikely to need radical changes like BLAS 1-3
 - Participation by more communities (graphics, machine learning, embedded, ...) will reduce risk of overspecification
 - Many of us plugged into SG1 etc., so we're watching architecture changes



See “Communication lower bounds & optimal algorithms for numerical linear algebra,” Ballard, Carson, Demmel, Hoemmen, Knight, & Schwarz, 2014

Alternate procedural approaches

- What if WG21 rejects / delays standardization?
- Join the BLAS Technical Forum?
 - (+) Meeting minutes show they planned a C++ interface [2]
 - (o) C++ & Java excluded: “time constraints” & possible Java changes
 - (+) CBLAS exists & has support from multiple vendors
 - (-) Invading small (< 20) focused committee, inactive nearly 20 years
 - (-) May reject C++ “binding” if too different (see MPI C++ binding)
- Form our own standard?
 - GraphBLAS: Express graph algorithms in the language of linear algebra
 - Batched BLAS: Do many small dense ops at once, expose parallelism
 - Non-language standards more useful when it’s worth having multiple implementations of the same interface
 - Multiple Engines / back-ends permit this, also encourage interoperability

Fundamental or batteries included?



- Referencing debate over Graphics API [3,4]
 - For this talk, I'm not taking a side
 - Fundamental-ness may influence WG21 prioritization
 - But: just as BLAS was codesigned w/ LINPACK & LAPACK, we should design the linear algebra we want, & standardize what we can
- Lowest level of linear algebra: Clearly fundamental, including
 - Multidimensional arrays: P0009r9, mspan, at LWG level
 - Multidimensional parallel iteration: Papers & ideas by authors here
 - Explicit SIMD types & ops: Voted into N4744 (Parallelism TS 2)
- What makes them fundamental?
 - Need / would benefit from compiler support
 - Vocabulary
 - Multiple use cases outside linear algebra

Multidimensional arrays: `mdspan`

- P0009r9, currently in LWG wording review
- Christian Trott will present on this
- Meant as a zero-overhead abstraction
 - Like Fortran arrays, complete w/ slices
 - Can mix compile- & run-time dimensions
 - No implicit temporaries / allocation (9x)
 - Compiler could optimize indexing
- Polymorphic data layout
 - Needed for performance with tensors & batched small dense
 - Performance portability (tune layout for architecture)
 - Compatibility with different libraries in other languages

Is mdspan a “matrix”?

- mdspan is nonowning
 - Generalizes span [views]
 - Can't create (allocate) matrix
 - Expressions may need temps
 - We use owning Kokkos::View in BLAS-like kokkos-kernels
- mdspan views memory
 - For small dense SIMD-based classes, may prefer values
 - Ask us about combining Kokkos::View & SIMD types
- Not a linear algebra library
 - No NumPy-style broadcast
 - No expressions, arithmetic, ...



If we want “Matlab in C++”:

- Owning tensor / mat / vec classes
- Dispatch to mdspan kernels
- mdspan may support expression templates through AccessorPolicy

Multidimensional parallel iteration



- Who needs this?
 - Linear & (especially) tensor algebra
 - Representations of physical space (structured grids)
- Why not just 1-D iterators?
 - Iterating 1 dim at a time ignores data locality
 - Problems most naturally expressed in multiple dimensions; “flattening” the iteration space by hand complicates code
 - All these must be optimized together:
 - Algorithm & data layout (hence mdspan’s polymorphic layout)
 - Iteration order & parallelization strategy
- Interface options
 - For-each i... over multidimensional index domain (ask us)
 - Einstein tensor index notation (ask other folks here)

Conclusions

- If we want to standardize a C++ linear algebra library,
- we should learn from BLAS Standard's successes:
 - Standardize in layers, from the bottom up
 - Layer by expertise: hardware / performance vs. numerical analysis
- Identify fundamentals for any linear algebra / tensor library
 - Explicit SIMD types & ops: Voted into N4744 (Parallelism TS 2)
 - Multidimensional arrays: P0009r9, in LWG wording review
 - Multidimensional parallel iteration: ????, OPPORTUNITY
- Thanks to P1417 coauthors & all of you!

References

- [1] John Markoff, “Writing the Fastest Code, by Hand, for Fun: A Human Computer Keeps Speeding Up Chips,” *New York Times*, Nov. 28, 2005.
- [2] BLAST Forum MINUTES, NIST, Washington, D.C., Oct. 8-9, 1998. Available online: <http://www.netlib.org/blas/blast-forum/blast-forum-minutes.oct8-9.98.html> [last accessed Feb. 14, 2019].
- [3] Guy Davidson, “Batteries not included: what should go in the C++ standard library?”, *World of hatcat*, Feb. 16, 2018. Available online: <https://hatcat.com/?p=16> [last accessed Feb. 10, 2019].
- [4] Titus Winters, “What Should Go Into the C++ Standard Library,” *Abseil Blog*, Feb. 27, 2018. Available online: <https://abseil.io/blog/20180227-what-should-go-stdlib> [last accessed Feb. 10, 2019].