

Design, Usage and Features of (T)LAPACK

github.com/tlapack/tlapack

Julien Langou^{1,2} Weslley Pereira³ Johnathan Rhyne¹
Thijs Steel⁴

¹University of Colorado Denver, Denver - CO, USA

²Inria, France

³National Renewable Energy Laboratory, Golden - CO, USA

⁴KU Leuven, Belgium

SIAM LA Minitutorial
May 14, 2024

Agenda

- General introduction to <T>LAPACK (T. Steel)
- Experience of a new developer (J. Rhyne)
- Level 0 (J. Langou)
- Mixed precision and fp8 (J. Langou)
- Rotation sequences (T. Steel)
- Some hands on exercises

What is $\langle T \rangle$ LAPACK ?

github.com/tlapack/tlapack

- C++ Template-based Linear Algebra Package.
- Routines are free functions using high-level abstraction.
- Interoperability with: SLATE, mdspan, Eigen3, StarPU.

```
1 // User's code looks like:  
2  
3 auto A = std::mdspanA( A_ptr, 100, 100 );  
4 auto B = Eigen::MatrixXd::Random(100,100).eval();  
5  
6 int infoA = getrf( A, piv );  
7 int infoB = getrf( B, piv );
```

```
1 // Inside potrf we find commands like:  
2  
3 A(piv[0], 0) = A(0, 0);  
4 auto A00 = slice(A, range(0, k0), range(0, k0));
```

Contributors

Senior Team Members

- Julien Langou, CU Denver
- Weslley Pereira, CU Denver
- Thijs Steel, KU Leuven

Undergrads, Grads, Postdocs (CU Denver)

- Heidi Meier, Lindsay Slager, Naomi Wilson, Natnael Tebeje, Yuxin Cai (Summer 2022, undergraduate)
- Ali Lotfi (Fall 2022, postdoc)
- Michael Pultorak, Racheal Asamoah, Skylar Johns (Spring 2023, undergraduate)
- Jonhathan Rhyne (AY 22-23, graduate)
- Alexander Hatle, Aslak Djupskas, Dominic Lioce, Hugh Kadhem, Seok-Ho Yoon, Sudhanva Kulkarni (Fall 2023, undergraduate)

Why $\langle T \rangle$ LAPACK ?

LAPACK

- provides reference implementation for a wide set of state-of-the-art algorithms.
- is largely used for its efficiency and reliability.
- has been incorporated into several packages.

BLAS / LAPACK

- are not templated, currently support 4 data types.
- are not templated, replicate a single algorithm on 4 files.
- work on a few data layouts.

$\langle T \rangle$ LAPACK

- uses templates that encompass data types and layouts.
- interoperates with SLATE, mdspan and Eigen3.
- uses modern software standards.

High-level design

Templates

- Any type works
- Any matrix layout works (but need to write plugin)
 - Access through ()
 - Submatrix through slice()
 - sizes through nrows, ncols
 - ...
- float, double, Eigen::Half, std::complex
- Your own matrix, Eigen, mdspan, ...
- All the same interface!

Power of templates: example

- Magma (GPU) has: gehrd and orghr, but not ormhr
(Hessenberg factorization and applying the orthogonal factor)
- ormhr is just a wrapper around ormqr

```
1      CALL DORMQR( SIDE, TRANS, MI, NI, NH, A( ILO+1,  
2      ILO ), LDA, TAU( ILO ), C( I1, I2 ), LDC, WORK,  
      LWORK, IINFO )
```

2

- In ⟨T⟩LAPACK , ormhr would not need to be rewritten for GPU!

High-level design

Optional arguments

- Workspace (allocated if not provided)
- Blocksize
- Algorithm selection
- Many parameters can be runtime or compile time

```
1 PotrfOpts opts;
2 opts.variant = PotrfVariant::Blocked;
3 opts.nb = 64;
4 potrf(Uplo::Lower, A, opts);
```

High-level design

Dimension errors

- Matrix and vector objects know their dimension
- Bounds checks on element access
- Checks for compatible dimensions in routines
- Bounds checks on submatrix/subvector
- Runtime checks can be disabled for performance

```
1 // A is a 5x5 matrix, x and y are size 4 vectors
2 // This throws an error
3 A(5,2) += 1.5;
4 // This throws an error
5 x[4] = 3;
6 // This throws an error
7 gemv( Op::NoTrans, 1.0, A, x, 0.0, y);
8 // This throws an error
9 auto A2 = slice(A, pair{0,5}, pair{3,6});
```

High-level design

Testing and debugging

- Let's be honest, the LAPACK test suite is not great
- $\langle T \rangle$ LAPACK offers:
 - Unit tests with an actual framework! (Catch2)
 - Synergy with bounds checks
 - Visualizer for debugging
- Less developer time spent on simple issues, more time spent on cool algorithms!

High-level design

Optimized BLAS

- Paradigm of LAPACK: use highly optimized implementations of matrix-matrix multiplication for portable high-performance.
- ⟨T⟩LAPACK uses BLAS++ and LAPACK++ from SLATE to access optimized BLAS and LAPACK.

High-level design

Open question: naming scheme

- In Fortran 77: limited amount of characters for routine names.
- Old names are now commonly used: geqrf, geev, ...
- New names would make much more sense for the future

High-level design

Open question: naming scheme

- In Fortran 77: limited amount of characters for routine names.
- Old names are now commonly used: `geqr`, `geev`, ...
- New names would make much more sense for the future

Audience participation what do these routines do?:

- `dpotrf`
- `dgghrd`
- `dgeqp3`

High-level design

Open question: naming scheme

- In Fortran 77: limited amount of characters for routine names.
- Old names are now commonly used: `geqr`, `geev`, ...
- New names would make much more sense for the future

Audience participation what do these routines do?:

- `dpotrf` Cholesky factorization
- `dgghrd`
- `dgeqp3`

High-level design

Open question: naming scheme

- In Fortran 77: limited amount of characters for routine names.
- Old names are now commonly used: `geqr`, `geev`, ...
- New names would make much more sense for the future

Audience participation what do these routines do?:

- `dpotrf` Cholesky factorization
- `dgghrd` Hessenberg-triangular reduction
- `dgeqp3`

High-level design

Open question: naming scheme

- In Fortran 77: limited amount of characters for routine names.
- Old names are now commonly used: `geqr`, `geev`, ...
- New names would make much more sense for the future

Audience participation what do these routines do?:

- `dpotrf` Cholesky factorization
- `dgghrd` Hessenberg-triangular reduction
- `dgeqp3` Pivoted QR

High-level design

Generality has a price

- ⟨T⟩LAPACK needs a full reimplementation. Extremely large amount of developer time!
- C++ is not slower, but also not faster.
- Near-term benefit is limited. Most of the benefits are down the line.

Case studies: Cholesky and Nonsymmetric eigenvalues

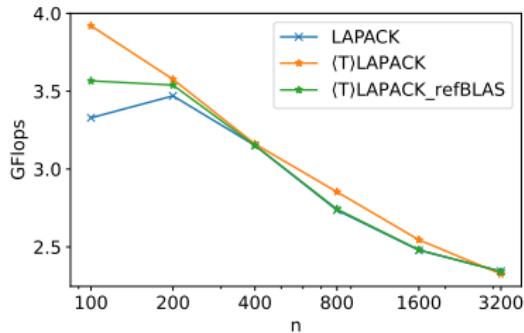
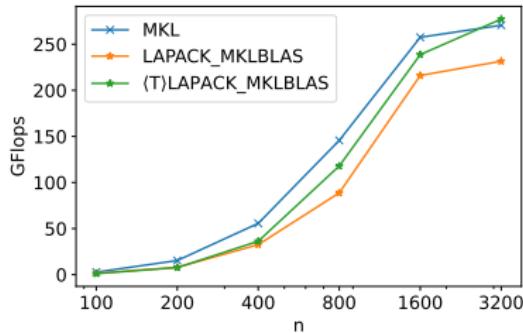
Wrappers to optimized BLAS

$\langle T \rangle$ LAPACK uses BLAS++ and LAPACK++ from SLATE to access optimized BLAS and LAPACK libraries.

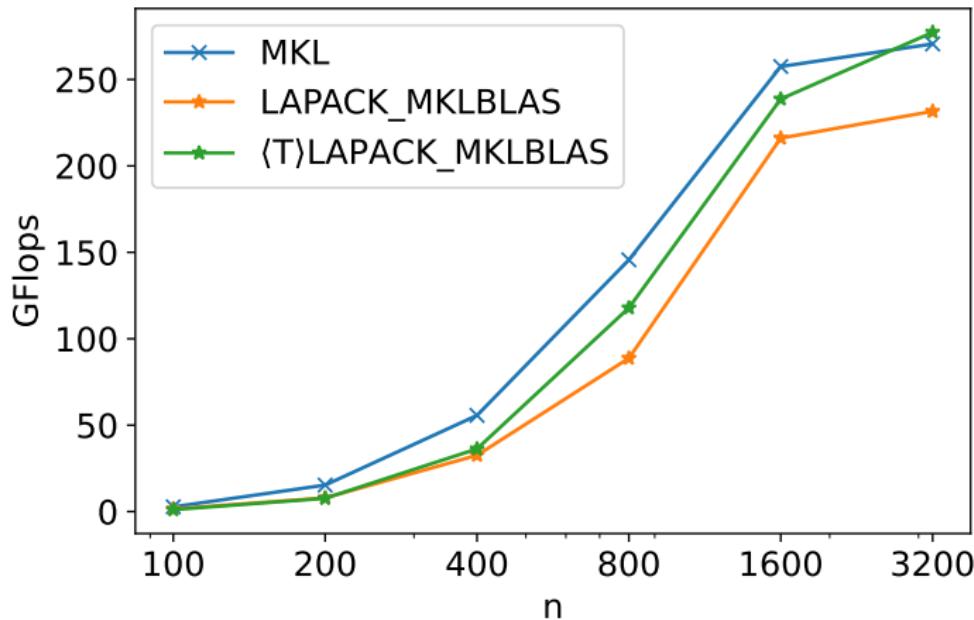
```
1 //>LAPACK implementation
2 template<TLAPACK_VECTOR vector_t,
3           disable_if_allow_optblas_t< vector_t > = 0>
4 auto nrm2( const vector_t& x )
5 { ... }
6
7 #ifdef TLAPACK_USE_LAPACKPP
8
9     /// Wrapper to BLAS++
10    template<TLAPACK_LEGACY_VECTOR vector_t,
11              enable_if_allow_optblas_t< vector_t > = 0>
12    auto nrm2( const vector_t& x )
13    {
14        auto x_ = legacy_vector(x);
15        const auto& n = x_.n;
16        return ::blas::nrm2( n, x_.ptr, x_.inc );
17    }
18
19 #endif
```

Performance test with recursive Cholesky

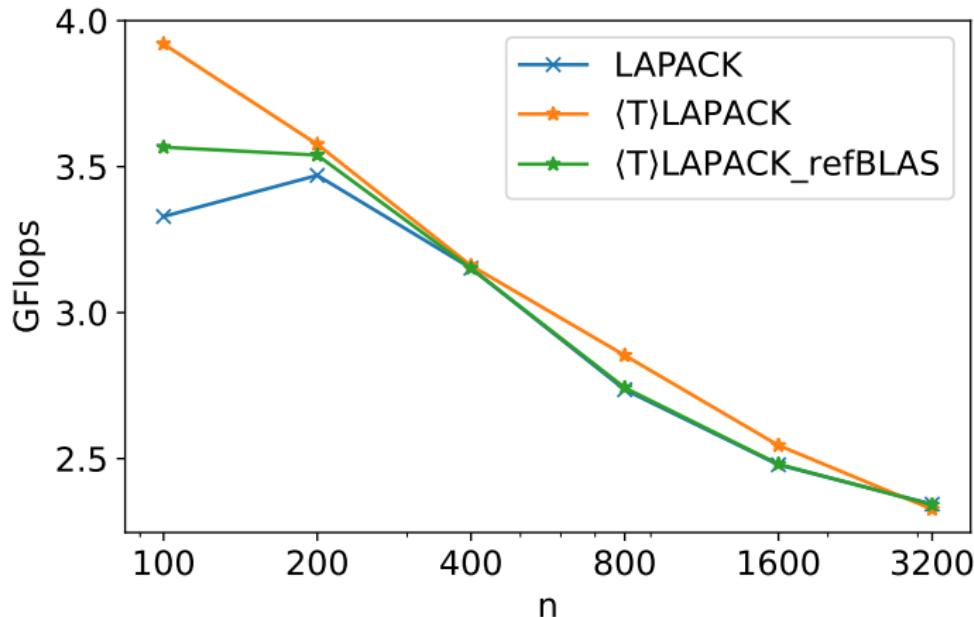
- 11th Gen i7-11800H @ 2.30GHz, 16 cores
- gcc version 9.3.0 (Ubuntu 9.3.0-17ubuntu1 20.04)
- LAPACK version 3.9.0
- oneAPI/MKL version 2022.0.2
- Main code: github.com/weslleyspereira/tlapack/blob/try-performance-example/examples/potrf/example_potrf.cpp
- Double matrix of size n -by- n



Performance test with recursive Cholesky



Performance test with recursive Cholesky



Integration with StarPU

In a few lines:

- `starpu_data_handle_t`: handle to manage a piece of data.
Splits data in a grid of nx-by-ny tiles.
- `tlapack::starpu::Matrix<T>`: matrix abstraction to `starpu_data_handle_t`. Allows for blocking that doesn't match with the nx-by-ny grid.
- `<T>LAPACK plugin`: provides `slice`, `nrows`, `ncols`, etc.

```
1 int main() {
2     /* initialize StarPU here */
3     float *A_;
4     starpu_malloc((void**)&A_, 90 * 90 * sizeof(float));
5     genHermitianMatrix(A_, 90, 90); // 90x90 matrix
6     Matrix<float> A(A_, 90, 90, 20, 20); // 20x20 grid
7
8     potrf_opts_t<size_t> opts; opts.nb = 9; // 9x9 blocks
9     int info = potrf(upperTriangle, A, opts);
10    /* terminate StarPU here */
11 }
12 }
```

Integration with StarPU

- **Key aspect:** No need to modify top-level routines!
- Blocked Cholesky calls gemm, herk, trsm and potf2.
- Those four routines have a StarPU (tiled) implementation.

Preliminary results with $n = 7680$ in single precision:

- ① Using only CPU workers:

MKL potrf: 589.0 GFLOPS,
StarPU Cholesky example: 519.3 GFLOPS,
(T)LAPACK potrf: 477.9 GFLOPS.

- ② Using GPU and CPU workers:

StarPU Cholesky example: 1835.1 GFLOPS,
(T)LAPACK potrf: 1445.6 GFLOPS.

Ubuntu 20.04.6 LTS x86_64. 11th Gen Intel(R) Core(TM) i7-11800H @ 2.30GHz (8 cores). NVIDIA GeForce RTX 3050 (2048 cores, 4GB). Driver version: 530.30.02. CUDA 11.0 (cuBLAS 11.2.0 cuSOLVER 10.6.0). MKL 2023.1.0. GNU compilers, version 9.4.0. StarPU scheduler: dmdas.

Performance of the real eigenvalue solver

Schur decomposition:

$$A = ZTZ^H$$

Two input scenarios:

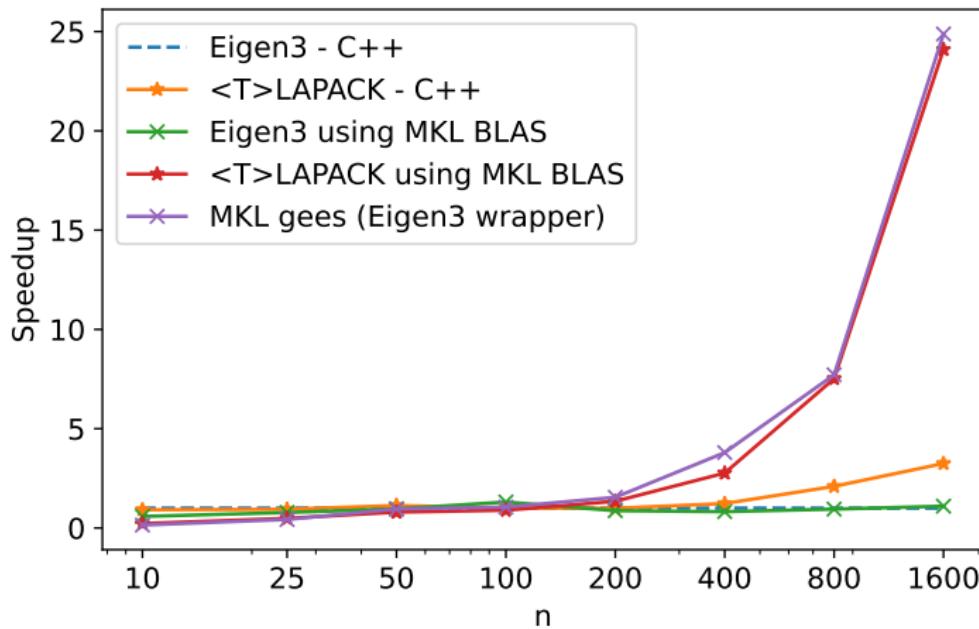
- ① Random A , with entries taken uniformly in $[0, 1]$.
- ② $A = (Q_1 D Q_2) \Lambda (Q_2^H D^{-1} Q_1^H)$, where
 - Λ is diagonal with entries taken uniformly in $[0, 1]$.
 - D is diagonal with $D_{ii} := 2^{(i \text{ mod } 54)}$.
 - Given A_i is random full rank, generate Q_i the Q-factor of the QR factorization of A_i , $i = 1, 2$.

Configurations in the following slides:

- Experiments in double precision.
- Package Eigen3 v3.3.7.
- Package mdsppan v0.4.0.
- oneAPI/MKL version 2022.2.1.

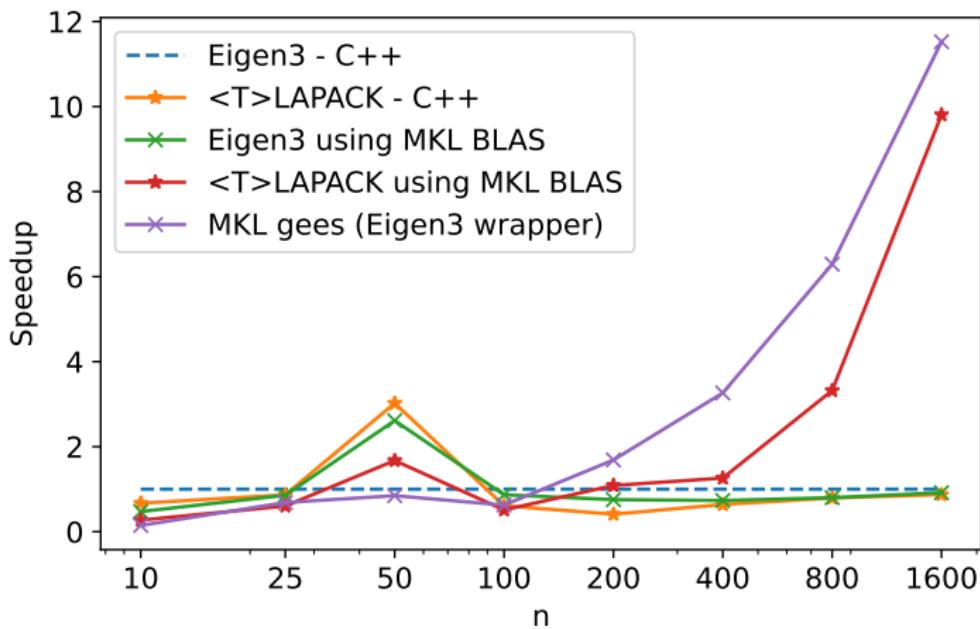
Performance in Scenario 1

Time to compute the eigenvalues compared to Eigen3.



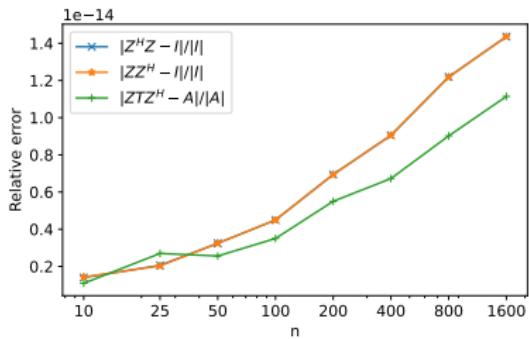
Performance in Scenario 2

Time to compute the eigenvalues compared to Eigen3.

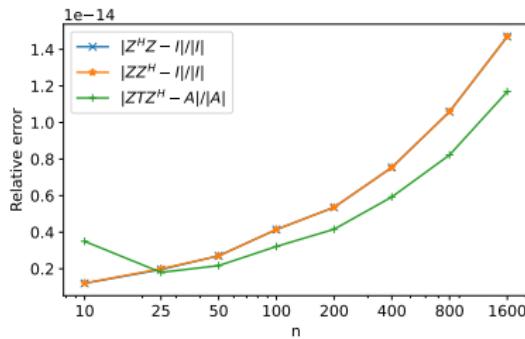


Stability analysis of the C++ codes

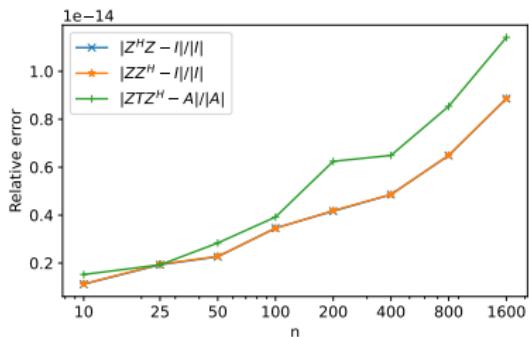
(T)LAPACK scenario 1



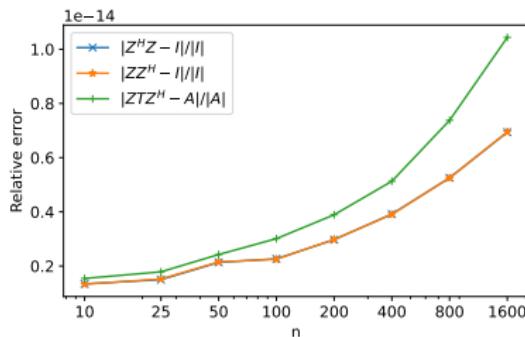
Eigen3 scenario 1



(T)LAPACK scenario 2

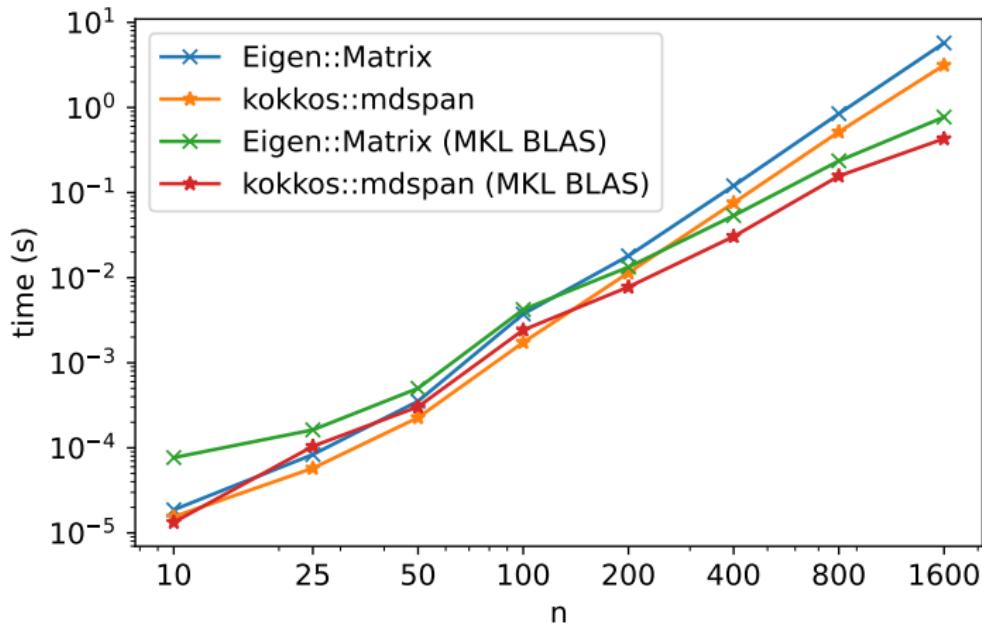


Eigen3 scenario 2



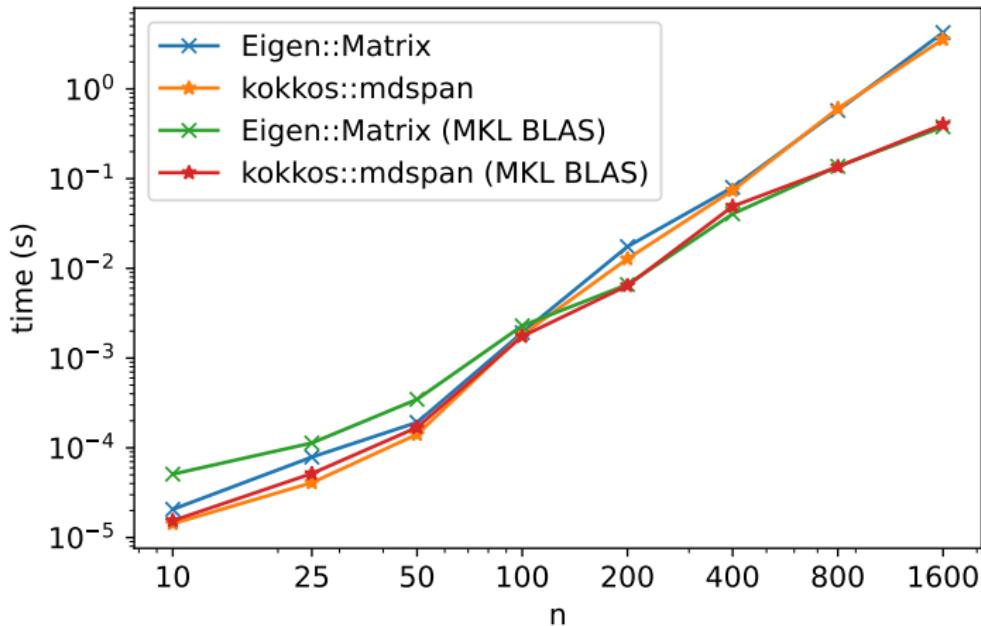
$\langle T \rangle$ LAPACK using mdspan from kokkos

Scenario 1 using different data structures with $\langle T \rangle$ LAPACK



$\langle T \rangle$ LAPACK using mdspan from kokkos

Scenario 2 using different data structures with $\langle T \rangle$ LAPACK



Code in the last example

Step-by-step ($\langle T \rangle$ LAPACK) \times Straight-forward (Eigen3)

Example using $\langle T \rangle$ LAPACK implementation:

```
1 // Hessenberg reduction
2 Eigen::Matrix<T,-1, 1> tau(n);
3 tlapack::gehrd( 0, n, A, tau );
4
5 // Discard reflectors
6 for (idx_t j = 0; j < n; ++j)
7     for (idx_t i = j + 2; i < n; ++i)
8         A(i,j) = T(0);
9
10 // Compute eigenvalues and Schur decomposition
11 Eigen::Matrix<std::complex<T>,-1,1> lambda(n);
12 Eigen::Matrix<T,-1,-1> Z(n,n); Z.setIdentity();
13 tlapack::multishift_qr( false, false, 0, n, A, lambda, Z );
```

Example using Eigen3 implementation:

```
1 Eigen::EigenSolver<decltype(A)> solver( A, false );
2 auto lambda = solver.eigenvalues();
```

Developing with <T>LAPACK from a student perspective

Developing for <T>LAPACK

As was said before, the codebase for <T>LAPACK is written in header files. This allows us to take advantage of templates!

My project within ⟨T⟩LAPACK

I worked on some software that was for a fork of ⟨T⟩LAPACK .
The code can be found here:

<https://github.com/jprhyne/tlapack/tree/hqr>, which was investigating porting EISPACK hqr into the ⟨T⟩LAPACK framework.

My project within <T>LAPACK

```
template <class matrix_t,
          class vector_t,
          enable_if_t<
              is_complex<type_t<vector_t>>::value,
              int> = 0
          >

int eispack_hqr(
    bool want_T,
    bool want_Q,
    size_type<matrix_t> low,
    size_type<matrix_t> igh,
    matrix_t &A,
    vector_t &eigs,
    matrix_t &Q,
    real_type<type_t<matrix_t>> &norm )
```

My workflow

- ① Write in a language I am comfortable with (C)
- ② Segment the code (create more functions!)
- ③ Transcribe a function at a time along with white-box testing
- ④ Transcribe the main driver

Examples of calling code

```
//Call hqr
real_t norm = real_t(zero);
retCode = tlapack::eispack_hqr(true, true, ilo, igh, H, s, Q, norm);
```

Looking at the testing files, we don't need to do anything special in order to call our templated functions. In fact, it is very similar to how you would use an interface like CBLAS or LAPACKE. It's possible to omit the namespace dereference in the function calls, but it is overall not too cumbersome to leave it this way and makes it much more clear where functions you are using come from.

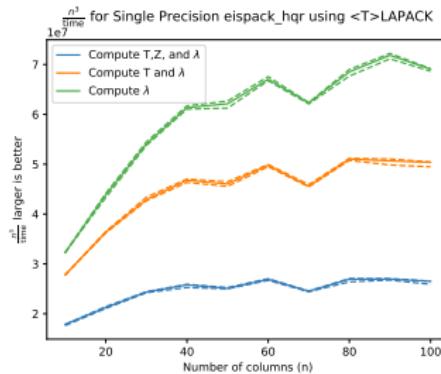
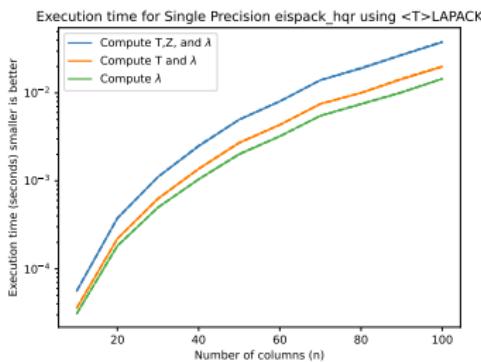
Testing our code

```
/* **** * \n/* Single precision section */\n\\ **** *\nstd::cout << "Single precision" << std::endl;\ncomputeTimes<legacyMatrix<float, std::size_t, Layout::ColMajor>>(n);\n/* **** *\n/* Double precision section */\n\\ **** *\nstd::cout << "Double precision" << std::endl;\ncomputeTimes<legacyMatrix<double, std::size_t, Layout::ColMajor>>(n);
```

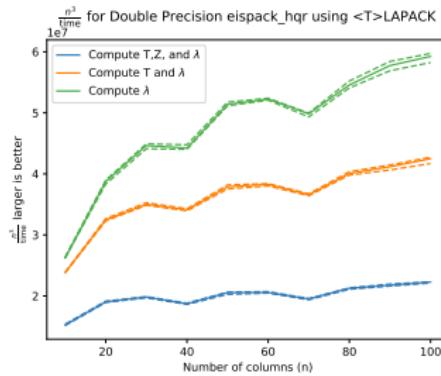
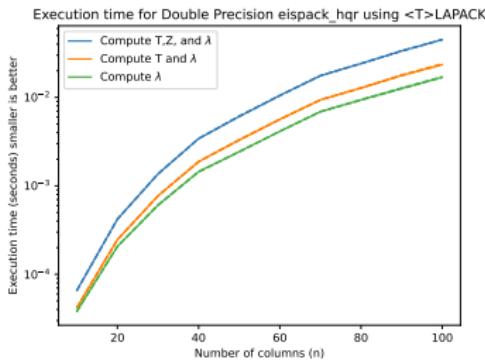
As we can see, when we template, the code that is called can be much simpler and if we are being given our data from another source, we don't need to implement any kind of checking to determine the correct call!

Timing of my project

We first present the first two “obvious” options for our types, single and double precision

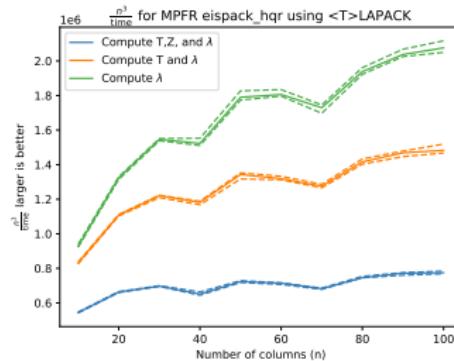
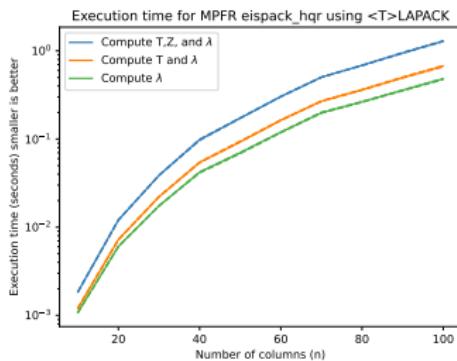


Timing of my project



Timing of my project

Using multi precision was also just a drag and drop!



Level 0

- easier for “compilers” to work with
- try to have more “versions” of algorithm and not necessarily one algorithm
- might offer interesting speedup for small sizes
- easier for humans to work with, (batched version, etc.)

Related to for example “User Guide for LDL, a concise sparse Cholesky package”, Tim Davis, version 3.3.2, Mar 22, 2024.

User Guide for LDL, a concise sparse Cholesky package - Tim Davis

```

void ldl_symmetric
{
    int n;           /* A and L are n-by-n, where n >= 0 */
    int Ap [ ],      /* input of size n!, not modified */
    int Lx [ ],      /* output of size n!, not modified */
    int Lxg [ ],     /* output of size n!, not modified */
    int Ip [ ],      /* output of size n!, not defined on input */
    int Px [ ],      /* output of size n, not defined on input */
    int P [ ],       /* optional input of size n */
    int Plse [ ]     /* optional output of size n (used if P is not NULL) */

    {
        int i, k, p, kk, y0 [ ];

        if (P)
        {
            /* If P is present then compute Plse, the inverse of P */
            for (Ct = 0; Ctx < n; Ct++)
            {
                for (Ip [k] = 1; k < n; k++)
                {
                    /* Lx(k,:); pattern: all nodes reachable in tree from m in A(0:k-1,k) */
                    Parent [k] = -1;          /* parent of k is not yet known */
                    Flag [k] = -1;            /* mark node k as visited */
                    Lx [k] = 0;               /* zero out Lx(k,:); column k of L */
                    kk = (k * Ct) + (k * Ct); /* kth original, or parent, column */
                    p0 = Ap [kk];
                    for (p = p0; p < p0 + p0; p++)
                    {
                        if (A (k,p) > 0) y0 [p] = y0 [p];
                    }
                    /* A (k,k) is nonzero (original or permuted k) */
                    i = (k * Ct) + (Flag [k] * Ct) + (Aii [k]);
                    if (i < k)
                    {
                        /* fall through from k to root of tree, stop at flagged node */
                        for (l : Flag [l] >= h : l < Parent [l])
                        {
                            /* find parent of l if not yet determined */
                            if (Parent [l] == -1) Parent [l] = i;
                            Lx [l] = 0;              /* zero out L (l,:); column l of L */
                            Flag [l] = k;             /* mark l as visited */
                        }
                    }
                }
            }
            /* construct Ip index array from Lx column counts */
            Ip [0] = 0;
            for (k = 0; k < n; k++)
            {
                Ip [k+1] = Ip [k] + Lx [k];
            }
        }
        /* returns n if successful, N if D (x,k) is zero */
        lnl ldl_symmetric
        {
            int a,           /* A & L are n-by-n, where n >= 0 */
            int Ap [ ],      /* input of size n!, not modified */
            int Lx [ ],      /* output of size n!, not modified */
            int Lxg [ ],     /* output of size n!, not modified */
            int Ip [ ],      /* input of size n!, not modified */
            int Px [ ],      /* output of size n!, not defined on input */
            int Lx [ ],      /* output of size n!, not defined on input */
            int Lxg [ ],     /* output of size n!, not defined on input */
            int Lxg_ip [ ],  /* output of size n!, not defined on input */
            double D [ ],    /* optional input of size n */
            double P [ ],    /* optional input of size n */
            int Plse [ ]     /* optional input of size n */

            double y1, y2, z1, z2, p0, len, top ;
            for (k = 0; k < n; k++)
            {
                /* compute nonzero Pattern of kth row of L in topological order */
                Lx [k] = 0;               /* Lx(k,:); column k of L */
                top = k;                 /* top of stack of all rows w/ */
                Ap [k] = 0;               /* stack for pattern in row k */
                Flag [k] = k;             /* mark node k as visited */
                Lxg_ip [k] = 0;           /* zero out Lxg_ip(k,:); column k of L */
                z1 = (D * k) + (P * k);  /* D (k,k) + P (k,k) = D (k,k) */
                z2 = (D * k) + (P * k) - (Aii [k]); /* with arithmetic, or permuted, column */
                p0 = Ap [kk];
                for (p = p0; p < p0 + p0; p++)
                {
                    if (A (k,p) > 0) y1 = (Ap [p] * z1) + (Aii [p]);
                    else y1 = 0;
                    if (y1 > 0)
                    {
                        y1 = (A (k,p) * z1) / Ap [p]; /* get A (k,p) */
                        if (y1 > 0)
                        {
                            y1 = -y1; /* As (p) */
                            Ap [p] = Ap [p] + Aii [p]; /* update A (p,p) into T (sum duplication) */
                            Lxg_ip [k] += k; /* k = k + 1 */
                            if (Lxg_ip [k] >= n) break; /* k = k + 1 */
                            /* Pattern (k,:) = 1;   /* Lx (k,:); column k of L */
                            Flag [k] = k;             /* mark k as visited */
                            Lxg_ip [k] = k;
                            while (len > 0) Pattern [-top] = Pattern [-len];
                        }
                    }
                    /* compute numerical values kth row of L (square triangular solve) */
                    D [k] = y1 * D [k];           /* get D (k,k) and clear Y (k,k) */
                    Y [k] = 0.0;
                    for (l : top <= k : top++)
                    {
                        l = Pattern [top];
                        y1 = (Lxg_ip [l] * y1) / Ap [l]; /* get and clear Y (l,k) */
                        Y [l] = 0.0;
                        p0 = Lxg_ip [l];
                        for (p = (p0 + 1); p < p0 + p0; p++)
                        {
                            if (A (l,p) > 0) y1 = (Lxg_ip [p] * y1) / Ap [p];
                            Lxg_ip [p] = y1;
                        }
                        Lxg_ip [l] = y1 / D [k]; /* the nonzero entry L (k,l) */
                        D [k] = D [k] * (1.0 - y1 * y1);
                        Lxg_ip [k] = 1.0;
                        Lxg_ip [k] += 1;
                        len = Lxg_ip [k] - 1;
                        /* increment count of nonzeros in col i */
                    }
                    if (D [k] == 0.0) return (k); /* failure, D (k,k) is zero */
                }
                return (n); /* success, diagonal of D is all nonzero */
            }
        }
    }
}

void ldl_leftsel
{
    int n;           /* L is n-by-n, where n >= 0 */
    double Ip [ ],   /* size n: right-hand-side on input, mln: on output */
    int Lx [ ],      /* input of size n!, not modified */
    int Lxg [ ],     /* input of size n!, not modified */
    double Lxg_ip [ ], /* input of size n!, not modified */
    int Lxg_ip_ip [ ] /* input of size n!, not modified */

    {
        int j, p, q0 ; /* L is n-by-n, where n >= 0 */
        for (j = 0; j < n; j++)
        {
            p = Ip [j];
            for (q = Ip [j] ; p < p0 + p0; p++)
            {
                X [j] = Lx [p] + X [j];
            }
        }
    }
}

void ldl_desel
{
    int a,           /* D is n-by-n, where n >= 0 */
    double Ip [ ],   /* size n: right-hand-side on input, mln: on output */
    double Lx [ ],   /* input of size n!, not modified */
    double Lxg [ ],  /* input of size n!, not modified */
    double Lxg_ip [ ] /* input of size n!, not modified */

    {
        int j ; /* L is n-by-n, where n >= 0 */
        for (j = 0; j < n; j++)
        {
            X [j] = D [j];
        }
    }
}

void ldl_ltrisolve
{
    int n;           /* L is n-by-n, where n >= 0 */
    double Ip [ ],   /* size n: right-hand-side on input, mln: on output */
    int Lx [ ],      /* input of size n!, not modified */
    int Lxg [ ],     /* input of size n!, not modified */
    double Lxg_ip [ ] /* input of size n!, not modified */

    {
        int l, j ; /* L is n-by-n, where n >= 0 */
        for (l = 0; l < n; l++)
        {
            for (j = 0; j < n; j++)
            {
                X [j] = D [j];
            }
        }
    }
}

void ldl_trisolve
{
    int n;           /* L is n-by-n, where n >= 0 */
    double Ip [ ],   /* size n: right-hand-side on input, mln: on output */
    int Lx [ ],      /* input of size n!, not modified */
    int Lxg [ ],     /* input of size n!, not modified */
    double Lxg_ip [ ] /* input of size n!, not modified */

    {
        int l, j ; /* L is n-by-n, where n >= 0 */
        for (l = 0; l < n; l++)
        {
            for (j = 0; j < n; j++)
            {
                X [j] = D [j];
            }
        }
    }
}

int ldl_ltrisolve
{
    int n;           /* L is n-by-n, where n >= 0 */
    double Ip [ ],   /* size n: right-hand-side on input, mln: on output */
    int Lx [ ],      /* input of size n!, not modified */
    int Lxg [ ],     /* input of size n!, not modified */
    double Lxg_ip [ ] /* input of size n!, not modified */

    {
        int l, j ; /* L is n-by-n, where n >= 0 */
        for (l = 0; l < n; l++)
        {
            for (j = 0; j < n; j++)
            {
                X [j] = D [j];
            }
        }
    }
}

int ldl_trisolve
{
    int n;           /* L is n-by-n, where n >= 0 */
    double Ip [ ],   /* size n: right-hand-side on input, mln: on output */
    int Lx [ ],      /* input of size n!, not modified */
    int Lxg [ ],     /* input of size n!, not modified */
    double Lxg_ip [ ] /* input of size n!, not modified */

    {
        int l, j ; /* L is n-by-n, where n >= 0 */
        for (l = 0; l < n; l++)
        {
            for (j = 0; j < n; j++)
            {
                X [j] = D [j];
            }
        }
    }
}

```

Level 0

gemm

```
for (i = 0; i < m; i++)  
    for (j = 0; j < n; j++)  
        for (k = 0; k < p; k++)  
            C[i][j] += A[i][k] * B[k][j];
```

syrk

```
for (i = 0; i < n; i++)  
    for (k = 0; k < m; k++)  
        for (j = 0; j <= i; j++)  
            C[i][j] += A[i][k] * A[j][k];
```

lu

```
for (i = 0; i < n; i++) {  
    for (j = 0; j < i; j++) {  
        for (k = 0; k < j; k++) {  
            A[i][j] -= A[i][k] * A[k][j];  
        }  
        A[i][j] /= A[j][j];  
    }  
    for (j = i; j < n; j++) {  
        for (k = 0; k < i; k++) {  
            A[i][j] -= A[i][k] * A[k][j];  
        }  
    }  
}
```

apply_givens_rot_sequence

```
for (p = 0; p < k; p++) {  
    for (j = 0; j < n - 1; j++) {  
        for (i = 0; i < m; i++) {  
            temp = C[j][p] * A[i][j] + S[j][p] * A[i][j + 1];  
            A[i][j + 1] = -S[j][p] * A[i][j] + C[j][p] * A[i][j + 1];  
            A[i][j] = temp;  
        }  
    }  
}
```

cholesky

```
for (i = 0; i < n; i++) {  
    for (j = 0; j < i; j++) {  
        for (k = 0; k < j; k++) {  
            A[i][j] -= A[i][k] * A[j][k];  
        }  
        A[i][j] /= A[j][j];  
    }  
    for (k = 0; k < i; k++) {  
        A[i][i] -= A[i][k] * A[i][k];  
    }  
    A[i][i] = sqrt(A[i][i]);  
}
```

Level 0

cgs – Classical Gram-Schmidt

```

for (j = 0; j < N; j++) {
    for (i = 0; i < j; i++) {
        R[i][j] = 0.0e+00;
        for (k = 0; k < M; k++)
            R[i][j] += A[k][i] * A[k][j];
    }
    for (i = 0; i < j; i++)
        for (k = 0; k < M; k++)
            A[k][j] -= A[k][i] * R[i][j];
    R[j][j] = 0.0e+00;
    for (k = 0; k < M; k++)
        R[j][j] += A[k][j] * A[k][j];
    R[j][j] = sqrt(R[j][j]);
    for (k = 0; k < M; k++)
        A[k][j] /= R[j][j];
}

```

mgs – Modified Gram-Schmidt

```

for (j = 0; j < N; j++) {
    for (i = 0; i < j; i++) {
        R[i][j] = 0.0e+00;
        for (k = 0; k < M; k++)
            R[i][j] += A[k][i] * A[k][j];
        for (k = 0; k < M; k++)
            A[k][j] -= A[k][i] * R[i][j];
    }
    R[j][j] = 0.0e+00;
    for (k = 0; k < M; k++)
        R[j][j] += A[k][j] * A[k][j];
    R[j][j] = sqrt(R[j][j]);
    for (k = 0; k < M; k++)
        A[k][j] /= R[j][j];
}

```

a2v – geqr2 – Householder QR factorization

```

norm2 = 0; k=0; k++){
    for(i = k; i < N; i++){
        for(l = k+1; l < M; l++){
            norm2 += A[l][k] * A[l][k];
        }
        norm = sqrt(A[k][k] + A[k][k] + norm2);
        A[k][k] = (A[k][k] > 0) ? (A[k][k] + norm) : (A[k][k] - norm);
        tau[k] = 2.0 / (1.0 + norm2 / (A[k][k] + A[k][k]));
        for(l = k+1; l < M; l++){
            A[l][k] /= A[l][k];
        }
        A[k][k] = (A[k][k] > 0) ? (-norm) : (norm);
        for(j = k+1; j < N; j++){
            tau[j] = 0.0;
            for(i = k+1; i < M; i++){
                tau[j] += A[i][k] * A[i][j];
            }
            tau[j] = tau[k] + tau[j];
            A[i][j] = A[i][j] - tau[j];
            for(l = k+1; l < M; l++){
                A[i][l] = A[i][l] - A[i][k] * tau[j];
            }
        }
    }
}

```

v2q – orqr2 – Construction of Q from Householder

```

for(k = N-1; k > -1; k--){
    for(j = k+1; j < N; j++){
        tau[j] = 0.0e+00;
        for(i = k+1; i < M; i++){
            tau[j] += A[i][k] * A[i][j];
        }
    }
    for(j = k+1; j < N; j++){
        tau[j] *= tau[k];
    }
    A[k][k] = 1.0e+00 - tau[k];
    for(i = k+1; i < N; i++){
        A[k][i] = -tau[i];
    }
    for(j = k+1; j < N; j++){
        for(i = k+1; i < M; i++){
            A[i][j] := A[i][k] * tau[j];
        }
    }
    for(i = k+1; i < M; i++){
        A[i][k] = -A[i][k] * tau[k];
    }
}

```

Level 0

gebd2 – reduction to bidiagonalization

```

for(k = 0; k < N; k++){
    norma2 = 0.e+0;
    for(i = k+1; i < M; i++){
        norma2 += A[i][k] * A[i][k];
    }
    norma = sqrt( A[k][k] + A[k][k] + norma2 );
    A[k][k] = ( A[k][k] > 0 ) ? ( A[k][k] + norma ) : ( A[k][k] - norma );
    tauq[k] = 2.0 / ( 1.0 + norma2 / ( A[k][k] * A[k][k] ) );
    for(i = k+1; i < M; i++){
        A[i][k] /= A[k][k];
        A[i][k] *= A[i][k];
    }
    A[k][k] = ( A[k][k] > 0 ) ? ( - norma ) : ( norma );
    for(j = k+1; j < N; j++){
        ttmp = A[k][j];
        for(i = k+1; i < M; i++){
            ttmp += A[i][k] * A[i][j];
        }
        ttmp = tauq[k] * ttmp;
        A[k][j] = A[k][j] - ttmp;
        for(i = k+1; i < M; i++){
            A[i][j] = A[i][j] - A[i][k] * ttmp;
        }
    }
    norma2 = 0.e+0;
    for(j = k+2; j < N; j++){
        norma2 += A[k][j] * A[k][j];
    }
    norma = sqrt( A[k][k+1] * A[k][k+1] + norma2 );
    A[k][k+1] = ( A[k][k+1] > 0 ) ? ( A[k][k+1] + norma ) : ( A[k][k+1] - norma );
    tauq[k] = 2.0 / ( 1.0 + norma2 / ( A[k][k+1] * A[k][k+1] ) );
    for(j = k+2; j < N; j++){
        A[k][j] /= A[k][k+1];
        A[k][j] *= A[k][k+1];
    }
    A[k][k+1] = ( A[k][k+1] > 0 ) ? ( - norma ) : ( norma );
    for(i = k+1; i < M; i++){
        ttmp = A[i][k+1];
        for(j = k+2; j < N; j++){
            ttmp += A[i][j] * A[k][j];
        }
        ttmp = ttmp * tauq[k];
        A[i][k+1] = A[i][k+1] - ttmp;
        for(j = k+2; j < N; j++){
            A[i][j] = A[i][j] - ttmp * A[k][j];
        }
    }
}

```

gehd2 – reduction to Hessenberg form

```

for(j = 0; j < n-2; j++) {
    norma2 = 0.0;
    for(i = j+2; i < n; i++) {
        norma2 += A[i][j] * A[i][j];
    }
    norma = sqrt( A[j+1][j] + A[j+1][j] + norma2 );
    A[j+1][j] = ( A[j+1][j] > 0 ) ? ( A[j+1][j] + norma ) : ( A[j+1][j] - norma );
    tau = 2.0 / ( 1.0 + norma2 / ( A[j+1][j] * A[j+1][j] ) );
    for(k = j+2; k < n; k++) {
        A[i][j] /= A[i+1][j];
    }
    A[j+1][j] = ( A[j+1][j] > 0 ) ? ( - norma ) : ( norma );
    for(i = j+1; i < n; i++) {
        ttmp[i] = A[i+1][j];
        for(k = j+2; k < n; k++) {
            ttmp[i] += A[k][j] * A[k][i];
        }
        ttmp[i] /= A[i+1][j];
    }
    for(i = j+1; i < n; i++) {
        ttmp[i] *= tau;
    }
    for(i = j+1; i < n; i++) {
        A[i+1][i] = ttmp[i];
    }
    for(i = j+2; i < n; i++) {
        for(k = j+1; k < n; k++) {
            A[i][k] -= A[i][j] * ttmp[k];
        }
    }
    for(i = 0; i < n; i++) {
        ttmp[i] = A[i][i+1];
        for(k = j+2; k < n; k++) {
            ttmp[i] += A[i][k] * A[k][j];
        }
    }
    for(i = 0; i < n; i++) {
        ttmp[i] *= tau;
    }
    for(i = 0; i < n; i++) {
        A[i][j+1] = ttmp[i];
    }
    for(i = 0; i < n; i++) {
        for(k = j+2; k < n; k++) {
            A[i][k] -= ttmp[i] * A[k][j];
        }
    }
}

```

Level 0

sytd2 – reduction to symmetric tridiagonal form

```

for(i = 0; i < n-2; i++){
    norma2 = 0.0e+00 ;
    for ( k = i+2; k < n ; k++ ) {
        norma2 += A[k][i] * A[k][i];
    }
    norma = sqrt( A[i+1][i] * A[i+1][i] + norma2 ) ;
    A[i+1][i] = ( A[i+1][i] > 0.0 ) ? ( A[i+1][i] + norma ) : ( A[i+1][i] - norma ) ;
    tau1 = 2.0e+00 / ( 1.0e+00 + norma2 / ( A[i+1][i] * A[i+1][i] ) ) ;

    for (k = i+2; k < n ; k++) {
        A[k][i] /= A[i+1][i];
    }
    A[i+1][i] = ( A[i+1][i] > 0.0e+00 ) ? ( - norma ) : ( norma ) ;
    for(k = i+1; k < n; k++){
        tau[k-1] = A[k][i+1];
        for(j = i+2; j <= k; j++){
            tau[k-1] *= A[j][i];
        }
        for(j = k+1; j < n; j++){
            tau[k-1] *= A[j][k] * A[j][i];
        }
        tau[k-1] *= tau1;
    }
    alpha = tau1;
    for(k = i+2; k < n; k++){
        alpha -= tau[k-1] * A[k][i];
    }
    alpha = -0.5e+00 * tau1 * alpha;
    tau[i] += alpha ;
    for(k = i+2; k < n; k++){
        tau[k-1] += alpha * A[k][i];
    }
    A[i+1][i+1] -= 2.0e+00 * tau1[1];
    for(j = i+1; j < n; j++){
        A[j][i+1] -= tau[j-1];
        A[j][i+1] -= tau[i] * A[j][i];
        for(k = i+2; k <= j; k++){
            A[j][k] -= tau[j-1] * A[k][i];
            A[j][k] -= tau[k-1] * A[j][i];
        }
    }
    tau[i] = tau1;
}

```

gghd2 – reduction to triangular Hessenberg form

```

for (j = 0; j < _PB_N-2; j++) {
    for (i = _PB_N-2; i > j; i--) {
        nm = SORT_FUN ( A[i][i] * A[i][j] + A[i+1][j] * A[i+1][j] );
        c = A[i][j] / nm;
        s = A[i+1][j] / nm;
        A[i][j] = nm;
        A[i+1][j] = SCALAR_VAL(0.0);
        for (k = j+1; k < _PB_N; k++) {
            tmp = c * A[i][k] + s * A[i+1][k];
            A[i+1][k] = - s * A[i][k] + c * A[i+1][k];
            A[i][k] = tmp;
        }
        for (k = i; k < _PB_N; k++) {
            tmp = c * B[i][k] + s * B[i+1][k];
            B[i+1][k] = - s * B[i][k] + c * B[i+1][k];
            B[i][k] = tmp;
        }
        for (k = 0; k < _PB_N; k++) {
            tmp = c * Q[i][k] + s * Q[i+1][k];
            Q[i+1][k] = - s * Q[i][k] + c * Q[i+1][k];
            Q[i][k] = tmp;
        }
        nm = SORT_FUN ( B[i+1][i+1] * B[i+1][i+1] + B[i+1][1] * B[i+1][1] );
        c = B[i+1][i+1] / nm;
        s = B[i+1][1] / nm;
        B[i+1][i+1] = nm;
        B[i+1][1] = SCALAR_VAL(0.0);
        for (k = 0; k <= i; k++) {
            tmp = c * B[k][1] - s * B[k][i+1];
            B[k][i+1] = s * B[k][1] + c * B[k][i+1];
            B[k][1] = tmp;
        }
        for (k = 0; k < _PB_N; k++) {
            tmp = c * A[k][i] - s * A[k][i+1];
            A[k][i+1] = s * A[k][i] + c * A[k][i+1];
            A[k][i] = tmp;
        }
        for (k = i-j; k < _PB_N; k++) {
            tmp = c * Z[k][i] - s * Z[k][i+1];
            Z[k][i+1] = s * Z[k][i] + c * Z[k][i+1];
            Z[k][i] = tmp;
        }
    }
}

```

Level 0

a2v – gqr2 – Householder QR factorization	v2q – orqr2 – Construction of Q from Householder	sytd2 – reduction to symmetric tridiagonal form
<pre> for(k = 0; k < N; k++){ norm2 = 0.0e+00; for(l = k+1; l < N; l++) norm2 += A[l][k] * A[l][k]; norma = sqrt(norm2); A[k][k] = 1.0e+00 / norma; for(l = k+1; l < N; l++) A[l][k] = A[l][k] / norma; tauw[k] = 2.0 * (A[k][k] - 1.0) * A[k][k]; for(k = N-1; k > -1; k--){ for(i = k+1; i < N; i++){ tauw[i] = tauw[k] + A[i][k] * A[i][k]; } } } </pre>	<p>v2q – orqr2 – Construction of Q from Householder</p> <pre> for (i = 0; i < n; i++) for (k = 0; k < m; k++) for (j = 0; j < i; j++) C[i][j] += A[i][k] * A[j][k]; A[0][0] * R[i][j] = -tauw[i]; for(j = k+1; j < N; j) for (j = 0; j < N; j++) { for(l = k+1; l < m; l) A[l][j] -= A[i][l] * R[i][j]; R[i][j] = 0.0e+00; } } </pre>	<pre> for(i = 0; i < n-2; i++){ norma2 = 0.0e+00; for (k = i+2; k < n; k++) norma2 += A[k][i] * A[k][i]; norma = sqrt(A[i][i]); A[i][i] = 1.0e+00 / norma; for(k = i+1; k < N; k++) norma2 += A[k][i] * A[i][k]; norma = sqrt(A[i][i] + norma2); A[i][i] = 2.0 / (1.0 + norma2 / (A[i][i] * A[i][i])); for(i = k+1; i < N; i++) A[i][i] = 1.0e+00 / norma; } </pre>
syrk	gemm	gebd2 – reduction to bidiagonalization
<pre> for (i = 0; i < n; i++) for (k = 0; k < m; k++) for (j = 0; j < i; j++) C[i][j] += A[i][k] * A[j][k]; A[0][0] * R[i][j] = -tauw[i]; for(j = k+1; j < N; j) for (j = 0; j < N; j++) { for(l = k+1; l < m; l) A[l][j] -= A[i][l] * R[i][j]; R[i][j] = 0.0e+00; } } </pre>	<p>gemm</p> <pre> for (i = 0; i < m; i++) for (j = 0; j < n; j++) for (k = 0; k < p; k++) C[i][j] += A[i][k] * B[k][j]; tauw[k-1] = A[i][k]; for(j = i+2; j < N; j) tauw[k-1] += A[i][j]; for(j = k+1; j < N; j) tauw[k-1] += A[i][j]; } </pre>	<pre> for(k = 0; k < N; k++){ if (for[i = k+1; i < M; i++] norma2 = 0.0e+00; A[i][k] = (A[i][k] > 0) ? (A[i][k] + norma2) : (A[i][k] - tauw[k]) / (1.0 + norma2 / (A[i][k] * A[i][k])); for(i = k+1; i < M; i++) A[i][k] = 1.0e+00 / norma2; } </pre>
mgs – Modified Gram-Schmidt	apply_givens_rot_sequence	geh2d – reduction to Hessenberg
<pre> for (j = 0; j < N; j++) { for (i = 0; i < j; i++) { R[i][j] = 0.0e+00; for (k = 0; k < M; k++) R[i][j] += A[k][i] * A[k][j]; for (k = 0; k < M; k++) A[k][j] -= A[k][i] * R[i][j]; } R[j][j] = 0.0e+00; for (k = 0; k < M; k++) R[j][j] += A[k][j] * A[k][j]; R[j][j] = sqrt(R[j][j]); for (k = 0; k < M; k++) A[k][j] /= R[j][j]; } </pre>	<p>apply_givens_rot_sequence</p> <pre> += A for (p = 0; p < k; p++) { for (j = 0; j < n-1; j++) { for (i = 0; i < m; i++) { temp = C[i][p] * A[i][j] + S[i][p] * A[i][j+1]; A[i][j+1] = -S[i][p] * A[i][j] + C[i][p] * A[i][j+1]; j+2; i < n; i++) += A[i][j] * A[i][i]; } } </pre>	<pre> int (A[j+1][j] * A[j+2][j] + norma2) = (A[j+1][j] > 0) ? (A[j+1][j] * A[j+2][j] + norma2) : (1.0 + norma2 / (A[j+1][j] * A[j+2][j])); A[j+1][j] = 1.0e+00 / norma2; A[j+1][j] = (A[j+1][j] > 0) ? (-norma) : (+1; i < n; i++) = A[j+1][i]; = j+2; k < n; k++) tmp[i] += A[k][j] * A[k][i]; } for (i = j+1; i < n; i++) for (j = i; j < n; j++) for (k = 0; k < i; k++) A[i][j] = -A[i][j] * A[j][k]; A[i][j] /= A[j][j]; } for (k = 0; k < i; k++) A[i][i] = A[i][i] * A[i][k] * </pre>

Level 0

A. Olivry, J. Langou, L-N Pouchet, P. Sadayappan, and F. Rastello. **Automated derivation of parametric data movement lower bounds for affine programs**. In the Proceedings of PLDI 2020: Proceedings of the 41st ACM SIGPLAN Conference on Programming Language Design and Implementation, page 808–822, June 2020. ([link](#))

Output #1 is an IO lower bound

cholesky	$\max\left(\frac{1}{2}N(N+1), \frac{1}{6}\frac{1}{\sqrt{S}}(N-1)(N-2)(N-3) + \frac{1}{2\sqrt{S}}\frac{1}{S}(N-1)(N-2) - (N-2)(N-7) - 4\sqrt{2}S\right)$	$\frac{1}{6}\frac{1}{\sqrt{S}}N^3$
----------	---	------------------------------------

input is an affine code

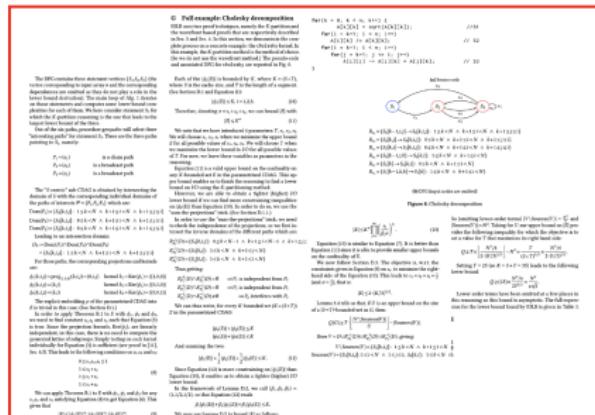
```

for (i = 0; i < _PB_N; i++) {
    for (j = 0; j < i; j++) {
        for (k = 0; k < k; k++) {
            A[i][j] -= A[i][k] * A[j][k];
        }
        A[i][j] /= A[j][j];
    }
    for (k = 0; k < i; k++) {
        A[i][i] -= A[i][k] * A[i][k];
    }
    A[i][i] = SQRT_FUN(A[i][i]);
}

```

cholesky code

Cholesky



Output #2 is a proof

Level 0

kernel	# input data	#ops	ratio	OI_{up}	OI_{manual}	ratio
2mm	$N_i N_k + N_k N_j$ $+ N_j N_l + N_i N_l$	$N i N_j N_k$ $+ N_i N_j N_l$	-	\sqrt{S}	\sqrt{S}	1 ✓
3mm	$N_i N_k + N_k N_j$ $+ N_j N_m + N_m N_l$	$N i N_j N_k + N_j N_l N_m$ $+ N_i N_j N_l$	-	\sqrt{S}	\sqrt{S}	1 ✓
cholesky	$\frac{1}{2} N^2$	$\frac{1}{3} N^3$	$\frac{2}{3} N$	$2\sqrt{S}$	\sqrt{S}	2
correlation	MN	$M^2 N$	M	$2\sqrt{S}$	\sqrt{S}	2
covariance	MN	$M^2 N$	M	$2\sqrt{S}$	\sqrt{S}	2
doitgen	$N_p N_q N_r$	$2N_q n_r N_p^2$	$2N_p$	\sqrt{S}	\sqrt{S}	1 ✓
fdtd-2d	$3N_x N_y$	$11N_x N_y T$	$\frac{11}{3} T$	$22\sqrt{2}\sqrt{S}$	$\frac{11}{24} \sqrt{3}\sqrt{S}$	$\frac{48\sqrt{2}}{\sqrt{3}}$
floyd-warshall	N^2	$2N^3$	$2N$	$2\sqrt{S}$	\sqrt{S}	2
gemm	$N_i N_j + N_j N_k + N_i N_k$	$2N_i N_j N_k$	-	\sqrt{S}	\sqrt{S}	1 ✓
heat-3d	N^3	$30N^3 T$	$30T$	$\frac{160}{3\sqrt{3}} \sqrt[3]{S}$	$\frac{5}{2} \sqrt[3]{S}$	$\frac{64}{3\sqrt{3}}$
jacobi-1d	N	$6NT$	$6T$	$24S$	$\frac{3}{2} S$	16
jacobi-2d	N^2	$10N^2 T$	$10T$	$15\sqrt{3}\sqrt{S}$	$\frac{5}{4}\sqrt{S}$	$12\sqrt{3}$
lu	N^2	$\frac{2}{3} N^3$	$\frac{2}{3} N$	\sqrt{S}	\sqrt{S}	1 ✓
ludcmp	N^2	$\frac{2}{3} N^3$	$\frac{2}{3} N$	\sqrt{S}	\sqrt{S}	1 ✓
seidel-2d	N^2	$9N^2 T$	$9T$	$\frac{27\sqrt{3}}{2}\sqrt{S}$	$\frac{9}{4}\sqrt{S}$	$6\sqrt{3}$
symm	$\frac{1}{2} M^2 + 2MN$	$2M^2 N$	-	\sqrt{S}	\sqrt{S}	1 ✓
syr2k	$\frac{1}{2} N^2 + 2MN$	$2MN^2$	-	$2\sqrt{S}$	\sqrt{S}	2
syrk	$\frac{1}{2} N^2 + MN$	MN^2	-	$2\sqrt{S}$	\sqrt{S}	2
trmm	$\frac{1}{2} M^2 + MN$	$M^2 N$	-	\sqrt{S}	\sqrt{S}	1 ✓
atax	MN	$4MN$	4	4	4	1 ✓
bicg	MN	$4MN$	4	4	4	1 ✓
deriche	HW	$32HW$	32	32	$\frac{16}{3}$	6
gemver	N^2	$10N^2$	10	10	5	2
gesummv	$2N^2$	$4N^2$	2	2	2	1 ✓
mvt	N^2	$4N^2$	4	4	4	1 ✓
trisolv	$\frac{1}{2} N^2$	N^2	2	2	2	1 ✓
adi	N^2	$30N^2 T$	$30T$	30	5	6
durbin	N	$2N^2$	$2N$	4	$\frac{2}{3}$	6
gramschmidt	MN	$2MN^2$	$2N$	$2\sqrt{S}$	1	$2\sqrt{S}$
nussinov	$\frac{1}{2} N^2$	$\frac{1}{3} N^3$	$\frac{2}{3} N$	$2\sqrt{S}$	1	$2\sqrt{S}$

Integration and Interaction

LAPACK 3.11 (Nov 2022)

criterion for deflation in the QZ algorithm by @thijssteel

- New “strict criterion” and back to the normwise criterion for detecting infinite eigenvalues (as opposed to elementwise criterion) ([PR#698](#)). See also <https://arxiv.org/abs/2208.02057>. Consider a pencil (H, T) in Hessenberg-triangular form:

$$(H, T) = \left(\begin{bmatrix} h_{11} & \dots & \dots & h_{1,n} \\ h_{21} & \ddots & & \vdots \\ & \ddots & \ddots & \vdots \\ 0 & & h_{n,n-1} & h_{nn} \end{bmatrix}, \begin{bmatrix} t_{11} & \dots & \dots & t_{1,n} \\ & \ddots & & \vdots \\ & & \ddots & \vdots \\ 0 & & & t_{nn} \end{bmatrix} \right).$$

- normwise (original): $|h_{i,i-1}| \leq u\|H\|_F$
- elementwise (since 3.9.1): $|h_{i,i-1}| \leq u(|h_{i-1,i-1}| + |h_{i,i}|)$
- strict (since 3.11.0): $|h_{i,i-1}| \leq u(|h_{i-1,i-1}| + |h_{i,i}|)$ and $|h_{i-1,i}t_{i,i} - h_{i,i}t_{i-1,i}| |h_{i,i-1}| \leq u|h_{i,i}| |h_{i-1,i-1}t_{i,i} - h_{i,i}t_{i-1,i-1}|$

LAPACK 3.11 (Nov 2022)

criterion for deflation in the QZ algorithm by @thijssteel

In single precision (32-bit), consider the pencil:

$$(H, T) = \left(\begin{bmatrix} 1 & c & 0 \\ \eta & (1+d) & 1 \\ 0 & \eta & (1+2d)c^{-1} \end{bmatrix}, \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & c^{-1} \end{bmatrix} \right),$$

with $\eta = 1.1 \cdot 10^{-8}$, $c = 1.1 \cdot 10^5$ and $d = 10^{-2}$.

- The eigenvalues of this pencil are approximately:

0.95371503, 1.0261424, and 1.0501426.

LAPACK 3.11 (Nov 2022)

criterion for deflation in the QZ algorithm by @thijssteel

In single precision (32-bit), consider the pencil:

$$(H, T) = \left(\begin{bmatrix} 1 & c & 0 \\ \eta & (1+d) & 1 \\ 0 & \eta & (1+2d)c^{-1} \end{bmatrix}, \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & c^{-1} \end{bmatrix} \right),$$

with $\eta = 1.1 \cdot 10^{-8}$, $c = 1.1 \cdot 10^5$ and $d = 10^{-2}$.

- The eigenvalues of this pencil are approximately:
[0.95371503](#), [1.0261424](#), and [1.0501426](#).
- Criteria applied only to H (the [normwise criterion](#) (3.9 and before), the [elementwise criterion](#) (3.10) or even the criterion by Ahues and Tisseur 1997) would consider $h_{3,2}$ to be small enough in IEEE single precision.
The returned eigenvalues are then

[1.](#), [1.02](#) and [1.04](#),

resulting in less than 2 digits of accuracy.

LAPACK 3.11 (Nov 2022)

criterion for deflation in the QZ algorithm by @thijssteel

In single precision (32-bit), consider the pencil:

$$(H, T) = \left(\begin{bmatrix} 1 & c & 0 \\ \eta & (1+d) & 1 \\ 0 & \eta & (1+2d)c^{-1} \end{bmatrix}, \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & c^{-1} \end{bmatrix} \right),$$

with $\eta = 1.1 \cdot 10^{-8}$, $c = 1.1 \cdot 10^5$ and $d = 10^{-2}$.

- The eigenvalues of this pencil are approximately:
[0.95371503](#), [1.0261424](#), and [1.0501426](#).
- Criteria applied only to H (the [normwise criterion](#) (3.9 and before), the [elementwise criterion](#) (3.10) or even the criterion by Ahues and Tisseur 1997) would consider $h_{3,2}$ to be small enough in IEEE single precision.

The returned eigenvalues are then

[1.](#), [1.02](#) and [1.04](#),

resulting in less than 2 digits of accuracy.

- When using the [strict criterion](#) (3.11), the QZ algorithm performs 2 more iterations, resulting in smaller subdiagonal entries. The returned eigenvalues are:

[0.95371503](#), [1.0261424](#), and [1.0501425](#),

which are almost accurate to the working precision.

LAPACK 3.11 (Nov 2022)

level-3 BLAS solvers for the triangular system and triangular Sylvester equation by @angsch

The triangular Sylvester equation has been recognized to be prone to overflow. For that purpose, *TRSYL utilizes a scaling factor to represent the solution as and solve the scaled equation . Due to the scaling factor, there is some flexibility in the representation of the solution. The proposed level-3 BLAS version, *TRSYL3, computes the scaling factors based on the upper bounds of blocks to enable level-3 BLAS. The scaling is typically slightly more aggressive so that an alternatively scaled final solution is computed. This is no problem as long as the scaling factor does not get flushed to zero ([PR#651](#)). The same upper bound calculation was used to write the level-3 BLAS solver for the triangular system, *LATRS3.

LAPACK 3.11 (Nov 2022)

level-3 BLAS solvers for the triangular system and triangular Sylvester equation by @angsch

Carl Christian Kjelgaard Mikkelsen, Angelika Beatrix Schwarz and Lars Karlsson. [Parallel robust solution of triangular linear systems.](#)
In Concurrency and Computation: Practice and Experience, e5064, Wiley Online Library, 2018.

Dimensions		Good Systems						Bad Systems			
m	n	DTRS{V,M}		Kiya		DLATRS		Kiya		DLATRS	
10000	1	0.00572	20	0.0183	5.6	0.0791	1.3	0.0183	5.3	0.0780	1.3
20000	1	0.0175	23.5	0.0494	8.2	0.328	1.2	0.0496	8.2	0.327	1.2
40000	1	0.0625	26.7	0.143	11.2	1.47	1.1	0.158	10.1	1.47	1.1
10000	4000	0.510	784	0.971	412	23.9	16.7	0.983	407	23.9	16.7
20000	1000	0.483	828	0.932	429	53.6	7.46	0.932	429	57.3	6.98
40000	250	0.665	601	1.04	385	226	1.8	1.05	381	229	1.7

Comparison of different solvers. Kiya is *LATRS3 introduced in 3.11. Timings are given in seconds with three significant digits and are followed by the performance given in GFLOPS/s. Timings of DLATRS for $n > 1$ right hand sides are extrapolated from measurements with 112 right-hand sides.

LAPACK 3.11 (Nov 2022)

Givens rotations generated with less accumulation errors by
@weslleyespereira

New algorithms for computing Givens rotations in complex arithmetic that reduce the accumulation errors for computing each of the outputs, c , s , r . The new algorithms are, on average, more accurate than both the algorithms from LAPACK 3.9.1 and LAPACK 3.10.0 ([PR#631](#)). See also <https://arxiv.org/abs/2211.04010>.

Motivated by a bug report about the algorithm in LAPACK 3.10 (See [Issue #629](#)). The author reported that the new Givens rotations may have lower accuracy than the ones that were in LAPACK up to release 3.9. This could be easily verified by noticing that, after applying several rotations to a unitary 2×2 matrix, the departure from unitarity was larger (in average) for 3.10 than 3.9. We performed additional experiments where we couldn't find the referred problem nor verify that the algorithm in LAPACK 3.10 was less accurate.

LAPACK 3.11 (Nov 2022)

Givens rotations generated with less accumulation errors by
@weslleyespereira

Issue #629 – In the test 5000 of plane rotations are applied to the vector of length 2 with the initial 2-norm of 10.

LAPACK 3.10

```
-bash-4.2$ ./a.out
...
The 2-norm of the first vector =
0.10001301E+02
...
Max of all errors printed above
0.13008118E-02
```

LAPACK 3.9

```
-bash-4.2$ ./a.out
...
The 2-norm of the first vector =
0.99999943E+01
...
Error abs(norm of the first column - initial norm )
0.57220459E-05
...
Max of all errors printed above
0.57220459E-05
```

LAPACK 3.11 (Nov 2022)

Givens rotations generated with less accumulation errors by
@weslleyespereira

We run each code several (10^6) times with random input data.

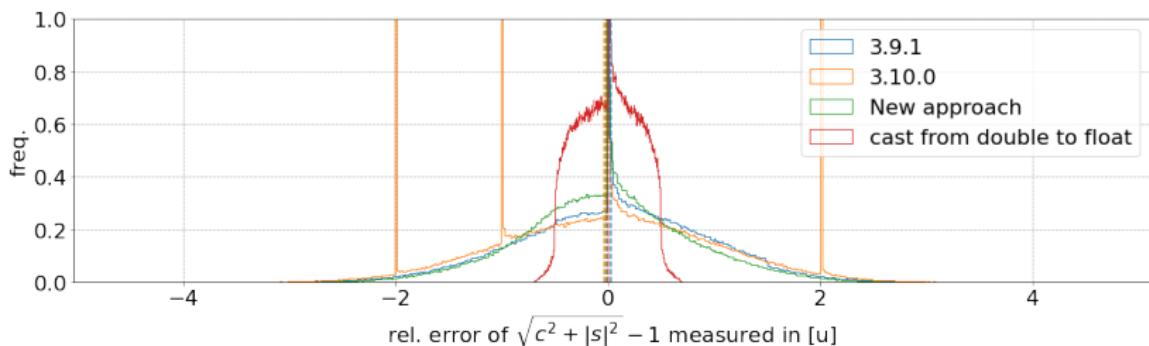


Figure: Histogram of the error in the singular values,
 $Err := \sqrt{c^2 + |s|^2} - 1$, for each algorithm measured in unit roundoff.

LAPACK 3.11 (Nov 2022)

Givens rotations generated with less accumulation errors by
@weslleyespereira

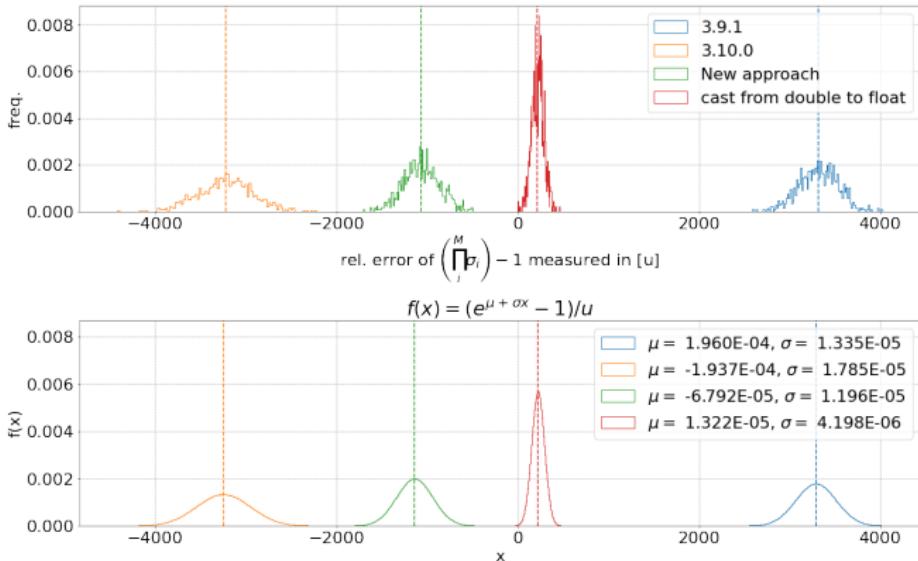


Figure: Top: Histogram of $\left(\prod_{i=1}^M \sigma_i\right) - 1$, $M = 10^5$, for each algorithm measured in unit roundoff. Repeat $N = 10^3$ times. Bottom: Estimates from formulae.

LAPACK 3.11 (Nov 2022)

Faster algorithms for Least Squares by @scr2016

- The new algorithms, *GELST, are similar to *GELS. *GELST avoids computing triangular blocks twice as in *GELS, which means *GELST runs faster ([PR#739](#)).

LAPACK 3.11 (Nov 2022)

Improvements on the building system and the Continuous Testing framework

- Regression test for illegal modification of Y in *GEMV by @matcross in [PR#622](#)
- Test the Fortran intrinsic ABS and complex divisions and report failures during build by @weslleyspereira in [PR#623](#)
- Appveyor is working in the Continuous Testing once more by @martin-frbg in [PR#627](#)
- More tests and fixes for type mismatches by @angsch in [PR#675](#)
- Solve build issues with IBM XLF by @friedc in [PR#677](#)
- Explicit type casts to INT by @angsch in [PR#684](#)
- Fix an out-of-bounds write in the tests by @angsch in [PR#685](#)
- Adding support for NAG Fortran compiler (nagfor) by @ACSimon33 in [PR#686](#)
- Fix time check flag propagation and default value by @ACSimon33 in [PR#696](#)
- Fixed format warnings in 64 bit integer builds by @ACSimon33 in [PR#700](#)
- Explicit type casts and more rigorous CI by @weslleyspereira in [PR#703](#)
- Fix type mismatches in function calls of testing code by @mjacobse in [PR#706](#)
- Fix lower triangular band matrix generation by @angsch in [PR#728](#)

LAPACK 3.11 (Nov 2022)

Other improvements

- Improves the computation of norms in S,DORBDB6 and C,ZUNBDB6 by @christoph-conrads in [PR#647](#)
- Uses a Newton step to compute the initial guess for ETA in S,DLAED4 to reduce the number of iterations by @wesleleyspereira, suggestion from Ren-Cang Li, in [PR#655](#)
- Return immediately when scaling with one by @angschn in [PR#674](#)
- Rearrange the application of the Householder reflectors in *LAQR5, which lowers the instruction count when FMA is available by @angschn in [PR#681](#)
- Add C,ZROTG, CS,ZDROT, S,DCABS1 to CBLAS by @angschn in [PR#721](#)
- Add *LANGB to LAPACKE by @ACSimon33 in [PR#725](#)
- NaN check for trapezoidal matrices on LAPACKE by @ACSimon33 in [PR#738](#) and [PR#742](#)

LAPACK 3.11 (Nov 2022)

Uses a Newton step to compute the initial guess for ETA in S,DLAED4 to reduce the number of iterations by @weslleyspereira, suggestion from Ren-Cang Li

- Uses a Newton step to compute the initial guess for ETA in S,DLAED4 to reduce the number of iterations by @weslleyspereira in [PR#655](#)

line 333 in xLAED4 in 3.10

```
ETA = DLTUB - TAU
```

line 336 in xLAED4 in 3.11

```
ETA = -W / ( DPSI+DPHI )
```

See as well [Issue#99](#)

Bug fixes

- Fix LAPACKE_*tpmqrt_work for row-major matrices by @weslleyspereira in [PR#540](#)
- Fix leading dimension check in LAPACKE's *geesv[x] and *gges[x] by @angsch in [PR#655](#)
- Fix the left-looking variant of GEQRF by @weslleyspereira in [PR#690](#)
- Fix workspace query for *SYEVD and *HEEVD routines by @neil-lindquist in [PR#691](#)
- Fix behavior of SCALE in *LATBS and *LATRS, and avoids NaN generation if entries in CNORM exceed the overflow threshold by @angsch in [PR#712](#)

LAPACK 3.11 (Nov 2022)

New Contributors

- @PhiliHS made their first contribution in [PR#601](#)
- @mtowara made their first contribution in [PR#663](#)
- @dbakshee made their first contribution in [PR#671](#)
- @ACSSimon33 made their first contribution in [PR#686](#)
- @friedc made their first contribution in [PR#677](#)
- @neil-lindquist made their first contribution in [PR#691](#)
- @ivan-pi made their first contribution in [PR#699](#)
- @mjacobse made their first contribution in [PR#702](#)
- @MonicaLiu0311 made their first contribution in [PR#720](#)
- @eltociear made their first contribution in [PR#735](#)

LAPACK 3.12 (June/July 2023?)

- [xGEDMD](#) and [xGEDMDQ](#) – Nonsymmetric DMD by Zlatko Drmač and Daniel Bielich and AIMdyn Inc. Santa Barbara, CA – PR#736
- [xGEQP3RK](#) – Truncated QR With Column Pivoting by Igor Kozachenko and Jim Demmel.
Igor's branch
- [xRSCL](#) – Reciprocal Scaling by Weslley Pereira
PR#749
- Better Organization for the Documentation by Mark Gates
PR#802
- (maybe) [IxAMAX](#) – exception handling “compliant” MAX search by Weslley Pereira with input from Jim Demmel, Greg Henry, Jason Riedy, etc.

LAPACK updates

J. Demmel, J. Dongarra, M. Gates, G. Henry, J. Langou, X. Li, P. Luszczek, W. Pereira, Wesleley, J. Riedy, and C. Rubio-Gonzalez.

[Proposed Consistent Exception Handling for the BLAS and LAPACK](#). In Sixth International Workshop on Software Correctness for HPC Applications (Correctness 2022), a workshop of ACM/IEEE SC 2022 Conference (SC'22), Dallas, TX, USA, November 13-18, 2022 10.1109/Correctness56720.2022.00006

- [rcond is 0 for NaN matrix in lapack-3.11.](#)

<https://github.com/Reference-LAPACK/lapack/issues/763>

<https://savannah.gnu.org/bugs/?63384>

- proper scaling function (reciprocal scaling)
- NaN propagation (IxAMAX)

[Related to BLAS:](#)

- [testBLAS](#) to test and find inconsistencies in the use of the BLAS standard. Testing MKL, BLIS, OpenBLAS, and also BLAS++.

NaN propagation in BLAS

```
alpha = 1  
beta = 0  
call GEMM for C := alpha * A * B + beta * C
```

Do we want the NaNs in “input C” to propagate to “output C”?

testBLAS – compiler part

```
./test_zcomplexabs 2> test_zcomplexabs.err
# All tests pass for ABS(a+b*I)
./test_zcomplexdiv 2> test_zcomplexdiv.err
!! Some (x+x*I)/(x+x*I) differ from 1
!! Some (x+x*I)/(x-x*I) differ from I
#           12602  tests out of      12606  pass for complex division,
# Please check the failed divisions in [stderr]
./test_zcomplexmult 2> test_zcomplexmult.err
# All tests pass for complex multiplication.
./test_zminMax 2> test_zminMax.err
[i8] MIN(  NaN,    0.) =    0.
[i8] MAX(  NaN,    0.) =    0.
#           14  tests out of      16  pass for intrinsic MIN and MAX
```

Continuous Testing for <T>LAPACK

testBLAS: test corner cases and exception handling in the BLAS.

- Started in Jun 20, 2021.
- Used in Demmel, J. et.al. (2022). Proposed Consistent Exception Handling for the BLAS and LAPACK

BLAS++ and LAPACK++ testers for comparing with legacy code.

- Tests compare output and execution time with LAPACKE.

Tests inside the <T>LAPACK project.

- Test different data types, layouts and matrix classes, e.g., legacyMatrix, Eigen::Matrix and mspan.

Impact on other software

Related to Eigen:

- Commented on bug <https://gitlab.com/libeigen/eigen/-/issues/408>
- GCC complex was incompatible to Eigen::half. We proposed a PR that, we believe, is being evaluated: <https://github.com/gcc-mirror/gcc/pull/8>.

Related to StarPU:

- StarPU does not have the means to declare static constant codelets.

<https://github.com/starpu-runtime/starpu/pull/20>.

- Bug on one of their examples: <https://github.com/starpu-runtime/starpu/issues/18>.
- Bug in cuSOLVER 11: <https://github.com/starpu-runtime/starpu/issues/12>. (Still need to report it)

- Makes StarPU compatible with CUDA 12:

<https://github.com/starpu-runtime/starpu/pull/10>

- Task on elements: This is something outside StarPU in our opinion. They provide the means to do that. But they do not define a matrix class.

Related to SLATE, BLAS++, LAPACK++, and LAPACK:

- Test suite for Inf and NaN propagation and corner cases that works on both $\langle T \rangle$ LAPACK and BLAS++: <https://github.com/tlapack/testBLAS>.

- Failures in the compilation of $\langle T \rangle$ LAPACK with SLATE helps improving the latter, e.g., <https://bitbucket.org/icl/lapackpp/pull-requests/14>,
<https://github.com/icl-utk-edu/blaspp/pull/11>,
<https://bitbucket.org/icl/blaspp/pull-requests/38>.

- Reciprocal scaling: <https://github.com/Reference-LAPACK/lapack/pull/749>

- Related to BLAS:**
- testBLAS to test and find inconsistencies in the use of the BLAS standard: <https://arxiv.org/abs/2207.09281>. Testing MKL, BLIS, OpenBLAS, and also BLAS++.

Happening Now

- Thijs Steel working on the symmetric eigenvalue solver. See github.com/tlapack/tlapack/pulls.
- Julien Langou, Johnathan Rhyne, Thijs Steel on SIAM LA'24. See siam.org/conferences/cm/program/minitutorials/la24-minitutorials
- Strategies for more sustainable software at NREL, running $\langle T \rangle$ LAPACK in low precision. See linkedin.com/jobs/search/?currentJobId=3826121686. (Graduate Summer Internship)

Mixed precision in <T>LAPACK

New data types come as plugins, e.g.,

- float16, gnuquad, mpreal.

Steps to add a new type:

- ① Create a plugin file based on other plugins and on concepts.hpp.

It is actually just one step! But you can do better:

- ② Add tests for the new type.
- ③ Add an example of usage.
- ④ Maybe, add overload a routine for better use with that type.

See ongoing addition of bfloat16 at

<https://github.com/tlapack/tlapack/pull/382>.

Mixed precision in <T>LAPACK

Functions in <T>LAPACK work with matrices of different types.
They also handle mixed precision. How?

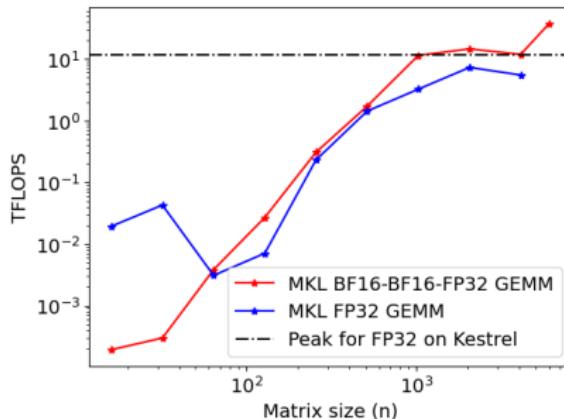
- `matrixA_t` and `matrixB_t` have entry types `TA` and `TB`, resp.
- Operations between $A(i,j)$ and $B(k,l)$ may have inputs in different precision.
- Any accumulator have the type common type between `TA` and `TB`, e.g.,
`double` for `pair<float,double>`,
`complex<double>` for `pair<complex<float>,double>`.

What if I want to operate in a different precision?

Mixed precision in <T>LAPACK : An ongoing experiment

Intel AMX¹:

- TMUL GEMM operates in mixed precision.
- `cblas_gemm_bf16bf16f32()` available in MKL (2023-2024).
- Less memory + Less time (up to 4x) + Lower accuracy.



Mixed precision in <T>LAPACK : TRMM blocked (and mixed)

```
1 // trmm's interface:  
2 void trmm_blocked_mixed(..., const matrixA_t& A, matrixB_t&  
    B, work_t& work)
```

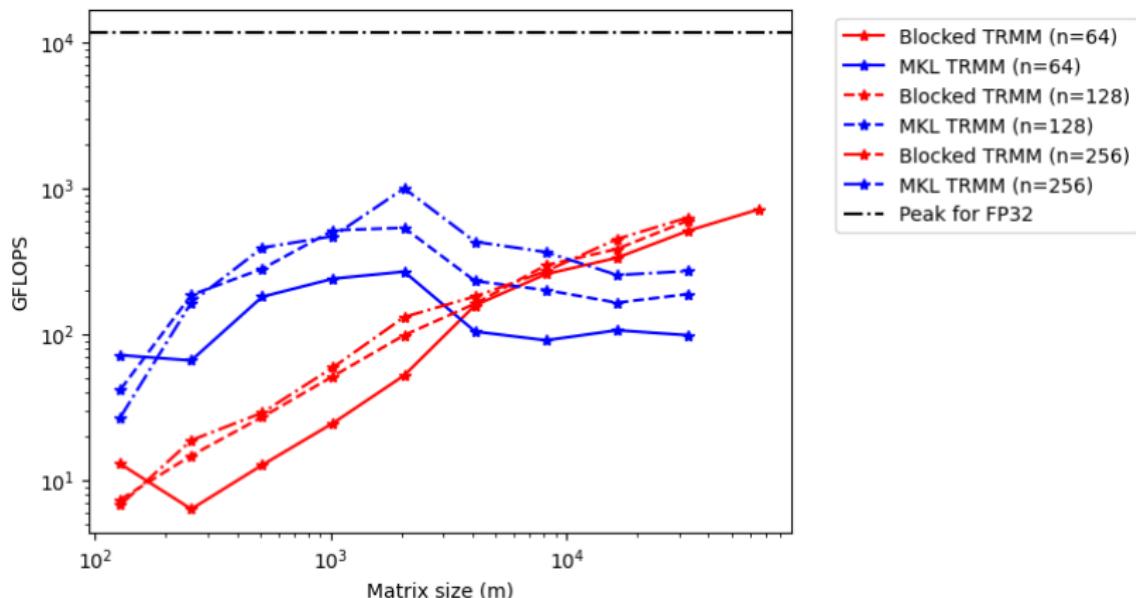
- `matrixA_t` and `work_t` in low precision (bfloating16).
- `matrixB_t` in high precision (float).

```
1 // Inside trmm's iteration:  
2 lacpy(GENERAL, Bi, BiLowPrecision);  
3 gemm(NO_TRANS, NO_TRANS, alpha, AOi, BiLowPrecision,  
4       real_t(1), B0);  
5 trmm(side, uplo, trans, diag, alpha, Aii, Bi);
```

- `lacpy()` does the type cast.
- `gemm()` has an overload to `cblas_gemm_bf16bf16f32()`.
- `trmm()` in mixed precision in <T>LAPACK , native C++.

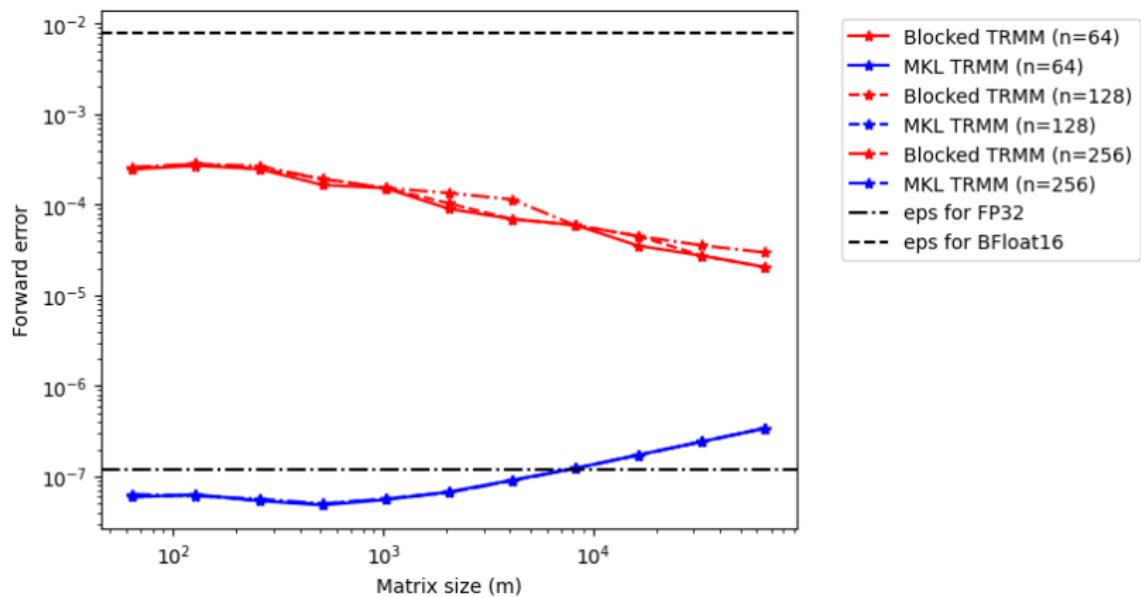
Mixed precision in $\langle T \rangle$ LAPACK : TRMM blocked (and mixed)

- Random double-precision matrices with values in $[0,1)$.
- $n_B = 110$ for blocked TRMM (Not necessarily optimal).



Mixed precision in <T>LAPACK : TRMM blocked (and mixed)

- Error: $\|B_{double} - \tilde{B}\|_1 / \|B_{double}\|_1$.



Thank you!

Questions?

Wesley.daSilvaPereira@nrel.gov

Github repo:

<https://github.com/tlapack/tlapack>

This work was authored in part by the National Renewable Energy Laboratory, operated by Alliance for Sustainable Energy, LLC, for the U.S. Department of Energy (DOE) under Contract No. DE-AC36-08GO28308. The views expressed in the article do not necessarily represent the views of the DOE or the U.S. Government. The U.S. Government retains and the publisher, by accepting the article for publication, acknowledges that the U.S. Government retains a nonexclusive, paid-up, irrevocable, worldwide license to publish or reproduce the published form of this work, or allow others to do so, for U.S. Government purposes.

testing fp8 linear algebra in <T>LAPACK

Sudhanva Kulkarni

adding fp8 in $\langle T \rangle$ LAPACK

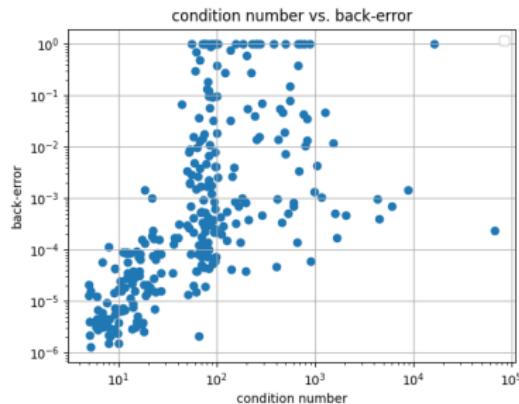
- IEEE P3109 format (8-bit floats) currently in development.
Software implementation easy to use in $\langle T \rangle$ LAPACK.
- Constraints on new fp types enforced by the use of "Scalar" concept in concepts.hpp .
- arithmetic for fp8 simulated by using higher precision datatypes like fp32 (operator+, operator-, etc). Can include stochastic rounding by changing the fp8 constructor.
Generates random bitstring with Mersenne Twister PRNG (std::mt19937)
- Simulate mixed precision BLAS by accumulating in fp32/fp16/bfloat vector and casting down

LAPACK in fp8 - getrf

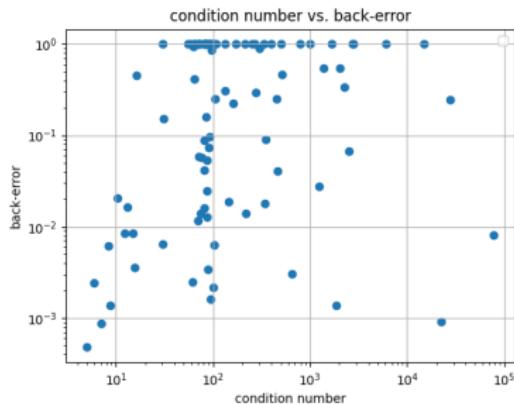
- Example routine - getrf for PLU decomposition
- Modifications made to reduce error include mixed precision BLAS/accumulation in fp32, equilibration in fp32, LU based iterative refinement.
- Experiments seem to convey that fp8 works well only for extremely well conditioned small matrices. Could use GMRES based IR for stronger refinement.
- Simple error analysis using Doolittle form gives element-wise worst case error bound $\sim |\Delta A| \leq ((1 + \epsilon_8)^2 - 1)|L||U|$ for LU on small matrices.
- need to account for underflow due to large underflow thresholds for fp8 types!

`xw0 = 0_0000_00 +Ob1.00*2^-15 = 0`
`xw1 = 0_0000_01 +Ob1.01*2^-15 = 7.62939e-06`
`xw2 = 0_0000_10 +Ob1.01*2^-15 = 1.52388e-05`
`xw3 = 0_0000_11 +Ob1.01*2^-15 = 2.28882e-05`
`xw4 = 0_0000_01 +Ob1.01*2^-15 = 3.05176e-05`
`xw5 = 0_0000_01 +Ob1.01*2^-15 = 3.8147e-05`
`xw6 = 0_0000_01 +Ob1.01*2^-15 = 4.57764e-05`
`xw7 = 0_0000_01 +Ob1.11*2^-15 = 5.34058e-05`
`xw8 = 0_0000_01 +Ob1.01*2^-14 = 6.10325e-05`
`xw9 = 0_0000_01 +Ob1.01*2^-14 = 7.62939e-05`
`xw10 = 0_0000_01 +Ob1.01*2^-14 = 9.15527e-05`
`xw11 = 0_0000_01 +Ob1.11*2^-14 = 0.00010812`
`xw12 = 0_0000_01 +Ob1.01*2^-13 = 0.00012207`
`xw13 = 0_0000_01 +Ob1.01*2^-13 = 0.000135288`
`xw14 = 0_0000_01 +Ob1.11*2^-13 = 0.000138105`
`xw15 = 0_0001_11 +Ob1.11*2^-13 = 0.000213623`
`xw16 = 0_0001_11 +Ob1.11*2^-13 = 0.0002244141`
`xw17 = 0_0001_10 +Ob1.01*2^-12 = 0.000303176`
`xw18 = 0_0001_10 +Ob1.11*2^-12 = 0.000363211`
`xw19 = 0_0001_11 +Ob1.11*2^-12 = 0.000427246`
`xw20 = 0_0001_09 +Ob1.01*2^-11 = 0.000488281`
`xw21 = 0_0001_09 +Ob1.01*2^-11 = 0.0010352`
`xw22 = 0_0001_10 +Ob1.10*2^-10 = 0.00073422`
`xw23 = 0_0001_11 +Ob1.11*2^-10 = 0.000854492`
`xw24 = 0_0001_10 +Ob1.01*2^-10 = 0.000975656`
`xw25 = 0_0001_10 +Ob1.11*2^-10 = 0.0012207`
`xw26 = 0_0001_10 +Ob1.01*2^-9 = 0.00146484`
`xw27 = 0_0001_11 +Ob1.11*2^-9 = 0.00170988`
`xw28 = 0_0001_10 +Ob1.11*2^-9 = 0.00193532`
`xw29 = 0_0001_11 +Ob1.01*2^-9 = 0.00244141`
`xw30 = 0_0001_10 +Ob1.11*2^-9 = 0.0029269`
`xw31 = 0_0001_11 +Ob1.11*2^-9 = 0.00341797`
`xw32 = 0_0100_00 +Ob1.00*2^-12 = 0.00398425`
`xw33 = 0_0100_01 +Ob1.01*2^-12 = 0.00488281`
`xw34 = 0_0100_10 +Ob1.10*2^-12 = 0.00585938`
`xw35 = 0_0100_11 +Ob1.11*2^-12 = 0.00683394`
`xw36 = 0_0100_00 +Ob1.00*2^-11 = 0.0078125`
`xw37 = 0_0100_01 +Ob1.01*2^-11 = 0.007976562`
`xw38 = 0_0100_10 +Ob1.10*2^-11 = 0.017188`
`xw39 = 0_0100_11 +Ob1.11*2^-11 = 0.01867719`
`xw40 = 0_0100_00 +Ob1.00*2^-10 = 0.015625`
`xw41 = 0_0100_01 +Ob1.01*2^-10 = 0.0195332`
`xw42 = 0_0100_10 +Ob1.10*2^-10 = 0.02347125`
`xw43 = 0_0100_11 +Ob1.11*2^-10 = 0.02737438`
`xw44 = 0_0101_00 +Ob1.00*2^-9 = 0.03125`
`xw45 = 0_0101_01 +Ob1.01*2^-9 = 0.0390625`
`xw46 = 0_0101_10 +Ob1.10*2^-9 = 0.046875`
`xw47 = 0_0101_11 +Ob1.11*2^-9 = 0.0546875`
`xw48 = 0_0110_00 +Ob1.00*2^-8 = 0.0623`
`xw49 = 0_0110_01 +Ob1.01*2^-8 = 0.078125`
`xw50 = 0_0110_10 +Ob1.10*2^-8 = 0.09375`
`xw51 = 0_0110_11 +Ob1.11*2^-8 = 0.109375`
`xw52 = 0_0110_00 +Ob1.00*2^-7 = 0.125`
`xw53 = 0_0110_01 +Ob1.01*2^-7 = 0.15625`
`xw54 = 0_0110_10 +Ob1.10*2^-7 = 0.1875`
`xw55 = 0_0110_11 +Ob1.11*2^-7 = 0.21875`
`xw56 = 0_0110_00 +Ob1.00*2^-6 = 0.25`
`xw57 = 0_0110_01 +Ob1.01*2^-6 = 0.3125`
`xw58 = 0_0110_10 +Ob1.10*2^-6 = 0.375`
`xw59 = 0_0110_11 +Ob1.11*2^-6 = 0.4375`
`xw60 = 0_0110_00 +Ob1.00*2^-5 = 0.5`
`xw61 = 0_0110_01 +Ob1.01*2^-5 = 0.625`
`xw62 = 0_0110_10 +Ob1.10*2^-5 = 0.75`
`xw63 = 0_0110_11 +Ob1.11*2^-5 = 0.875`
`xw64 = 0_1000_00 +Ob1.00*2^-12 = 4096`
`xw65 = 0_1000_01 +Ob1.01*2^-12 = 5120`
`xw66 = 0_1000_10 +Ob1.10*2^-12 = 6144`
`xw67 = 0_1000_11 +Ob1.11*2^-12 = 7168`
`xw68 = 0_1000_00 +Ob1.00*2^-11 = 8192`
`xw69 = 0_1000_01 +Ob1.01*2^-11 = 10240`
`xw70 = 0_1000_10 +Ob1.10*2^-11 = 12288`
`xw71 = 0_1000_11 +Ob1.11*2^-11 = 14336`
`xw72 = 0_1000_00 +Ob1.00*2^-10 = 16384`
`xw73 = 0_1000_01 +Ob1.01*2^-10 = 20480`
`xw74 = 0_1000_10 +Ob1.10*2^-10 = 24576`
`xw75 = 0_1000_11 +Ob1.11*2^-10 = 28672`
`xw76 = 0_1000_00 +Ob1.00*2^-9 = 32768`
`xw77 = 0_1000_01 +Ob1.01*2^-9 = 40960`
`xw78 = 0_1000_10 +Ob1.10*2^-9 = 49152`
`xw79 = 0_1000_11 +Ob1.11*2^-9 = -0.5`
`xw80 = 0_1000_00 +Ob1.00*2^-8 = -0.625`
`xw81 = 0_1000_01 +Ob1.01*2^-8 = -0.875`
`xw82 = 0_1000_10 +Ob1.10*2^-8 = -1.0`
`xw83 = 0_1000_11 +Ob1.11*2^-8 = -1.25`
`xw84 = 0_1000_00 +Ob1.00*2^-7 = -1.5`
`xw85 = 0_1000_01 +Ob1.01*2^-7 = -1.75`
`xw86 = 0_1000_10 +Ob1.10*2^-7 = -2.0`
`xw87 = 0_1000_11 +Ob1.11*2^-7 = -2.25`
`xw88 = 0_1000_00 +Ob1.00*2^-6 = -2.5`
`xw89 = 0_1000_01 +Ob1.01*2^-6 = -2.75`
`xw90 = 0_1000_10 +Ob1.10*2^-6 = -3.0`
`xw91 = 0_1000_11 +Ob1.11*2^-6 = -3.25`
`xw92 = 0_1000_00 +Ob1.00*2^-5 = -3.5`
`xw93 = 0_1000_01 +Ob1.01*2^-5 = -3.75`
`xw94 = 0_1000_10 +Ob1.10*2^-5 = -4.0`
`xw95 = 0_1000_11 +Ob1.11*2^-5 = -4.25`
`xw96 = 0_1000_00 +Ob1.00*2^-4 = -4.5`
`xw97 = 0_1000_01 +Ob1.01*2^-4 = -4.75`
`xw98 = 0_1000_10 +Ob1.10*2^-4 = -5.0`
`xw99 = 0_1000_11 +Ob1.11*2^-4 = -5.25`
`xw100 = 0_1000_00 +Ob1.00*2^-3 = -5.5`
`xw101 = 0_1000_01 +Ob1.01*2^-3 = -5.75`
`xw102 = 0_1000_10 +Ob1.10*2^-3 = -6.0`
`xw103 = 0_1000_11 +Ob1.11*2^-3 = -6.25`
`xw104 = 0_1000_00 +Ob1.00*2^-2 = -6.5`
`xw105 = 0_1000_01 +Ob1.01*2^-2 = -6.75`
`xw106 = 0_1000_10 +Ob1.10*2^-2 = -7.0`
`xw107 = 0_1000_11 +Ob1.11*2^-2 = -7.25`
`xw108 = 0_1000_00 +Ob1.00*2^-1 = -7.5`
`xw109 = 0_1000_01 +Ob1.01*2^-1 = -7.75`
`xw110 = 0_1000_10 +Ob1.10*2^-1 = -8.0`
`xw111 = 0_1000_11 +Ob1.11*2^-1 = -8.25`
`xw112 = 0_1000_00 +Ob1.00*2^0 = -8.5`
`xw113 = 0_1000_01 +Ob1.01*2^0 = -8.75`
`xw114 = 0_1000_10 +Ob1.10*2^0 = -9.0`
`xw115 = 0_1000_11 +Ob1.11*2^0 = -9.25`
`xw116 = 0_1000_00 +Ob1.00*2^-1 = -9.5`
`xw117 = 0_1000_01 +Ob1.01*2^-1 = -9.75`
`xw118 = 0_1000_10 +Ob1.10*2^-1 = -10.0`
`xw119 = 0_1000_11 +Ob1.11*2^-1 = -10.25`
`xw120 = 0_1000_00 +Ob1.00*2^-2 = -10.5`
`xw121 = 0_1000_01 +Ob1.01*2^-2 = -10.75`
`xw122 = 0_1000_10 +Ob1.10*2^-2 = -11.0`
`xw123 = 0_1000_11 +Ob1.11*2^-2 = -11.25`
`xw124 = 0_1000_00 +Ob1.00*2^-3 = -11.5`
`xw125 = 0_1000_01 +Ob1.01*2^-3 = -11.75`
`xw126 = 0_1000_10 +Ob1.10*2^-3 = -12.0`
`xw127 = 0_1000_11 +Ob1.11*2^-3 = -12.25`
`xw128 = 0_1000_00 +Ob1.00*2^-4 = -12.5`
`xw129 = 0_1000_01 +Ob1.01*2^-4 = -12.75`
`xw130 = 0_1000_10 +Ob1.10*2^-4 = -13.0`
`xw131 = 0_1000_11 +Ob1.11*2^-4 = -13.25`
`xw132 = 0_1000_00 +Ob1.00*2^-5 = -13.5`
`xw133 = 0_1000_01 +Ob1.01*2^-5 = -13.75`
`xw134 = 0_1000_10 +Ob1.10*2^-5 = -14.0`
`xw135 = 0_1000_11 +Ob1.11*2^-5 = -14.25`
`xw136 = 0_1000_00 +Ob1.00*2^-6 = -14.5`
`xw137 = 0_1000_01 +Ob1.01*2^-6 = -14.75`
`xw138 = 0_1000_10 +Ob1.10*2^-6 = -15.0`
`xw139 = 0_1000_11 +Ob1.11*2^-6 = -15.25`
`xw140 = 0_1000_00 +Ob1.00*2^-7 = -15.5`
`xw141 = 0_1000_01 +Ob1.01*2^-7 = -15.75`
`xw142 = 0_1000_10 +Ob1.10*2^-7 = -16.0`
`xw143 = 0_1000_11 +Ob1.11*2^-7 = -16.25`
`xw144 = 0_1000_00 +Ob1.00*2^-8 = -16.5`
`xw145 = 0_1000_01 +Ob1.01*2^-8 = -16.75`
`xw146 = 0_1000_10 +Ob1.10*2^-8 = -17.0`
`xw147 = 0_1000_11 +Ob1.11*2^-8 = -17.25`
`xw148 = 0_1000_00 +Ob1.00*2^-9 = -17.5`
`xw149 = 0_1000_01 +Ob1.01*2^-9 = -17.75`
`xw150 = 0_1000_10 +Ob1.10*2^-9 = -18.0`
`xw151 = 0_1000_11 +Ob1.11*2^-9 = -18.25`
`xw152 = 0_1000_00 +Ob1.00*2^-10 = -18.5`
`xw153 = 0_1000_01 +Ob1.01*2^-10 = -18.75`
`xw154 = 0_1000_10 +Ob1.10*2^-10 = -19.0`
`xw155 = 0_1000_11 +Ob1.11*2^-10 = -19.25`
`xw156 = 0_1000_00 +Ob1.00*2^-11 = -19.5`
`xw157 = 0_1000_01 +Ob1.01*2^-11 = -19.75`
`xw158 = 0_1000_10 +Ob1.10*2^-11 = -20.0`
`xw159 = 0_1000_11 +Ob1.11*2^-11 = -20.25`
`xw160 = 0_1000_00 +Ob1.00*2^-12 = -20.5`
`xw161 = 0_1000_01 +Ob1.01*2^-12 = -20.75`
`xw162 = 0_1000_10 +Ob1.10*2^-12 = -21.0`
`xw163 = 0_1000_11 +Ob1.11*2^-12 = -21.25`
`xw164 = 0_1000_00 +Ob1.00*2^-13 = -21.5`
`xw165 = 0_1000_01 +Ob1.01*2^-13 = -21.75`
`xw166 = 0_1000_10 +Ob1.10*2^-13 = -22.0`
`xw167 = 0_1000_11 +Ob1.11*2^-13 = -22.25`
`xw168 = 0_1000_00 +Ob1.00*2^-14 = -22.5`
`xw169 = 0_1000_01 +Ob1.01*2^-14 = -22.75`
`xw170 = 0_1000_10 +Ob1.10*2^-14 = -23.0`
`xw171 = 0_1000_11 +Ob1.11*2^-14 = -23.25`
`xw172 = 0_1000_00 +Ob1.00*2^-15 = -23.5`
`xw173 = 0_1000_01 +Ob1.01*2^-15 = -23.75`
`xw174 = 0_1000_10 +Ob1.10*2^-15 = -24.0`
`xw175 = 0_1000_11 +Ob1.11*2^-15 = -24.25`
`xw176 = 0_1000_00 +Ob1.00*2^-16 = -24.5`
`xw177 = 0_1000_01 +Ob1.01*2^-16 = -24.75`
`xw178 = 0_1000_10 +Ob1.10*2^-16 = -25.0`
`xw179 = 0_1000_11 +Ob1.11*2^-16 = -25.25`
`xw180 = 0_1000_00 +Ob1.00*2^-17 = -25.5`
`xw181 = 0_1000_01 +Ob1.01*2^-17 = -25.75`
`xw182 = 0_1000_10 +Ob1.10*2^-17 = -26.0`
`xw183 = 0_1000_11 +Ob1.11*2^-17 = -26.25`
`xw184 = 0_1000_00 +Ob1.00*2^-18 = -26.5`
`xw185 = 0_1000_01 +Ob1.01*2^-18 = -26.75`
`xw186 = 0_1000_10 +Ob1.10*2^-18 = -27.0`
`xw187 = 0_1000_11 +Ob1.11*2^-18 = -27.25`
`xw188 = 0_1000_00 +Ob1.00*2^-19 = -27.5`
`xw189 = 0_1000_01 +Ob1.01*2^-19 = -27.75`
`xw190 = 0_1000_10 +Ob1.10*2^-19 = -28.0`
`xw191 = 0_1000_11 +Ob1.11*2^-19 = -28.25`
`xw192 = 0_1000_00 +Ob1.00*2^-20 = -28.5`
`xw193 = 0_1000_01 +Ob1.01*2^-20 = -28.75`
`xw194 = 0_1000_10 +Ob1.10*2^-20 = -29.0`
`xw195 = 0_1000_11 +Ob1.11*2^-20 = -29.25`
`xw196 = 0_1000_00 +Ob1.00*2^-21 = -29.5`
`xw197 = 0_1000_01 +Ob1.01*2^-21 = -29.75`
`xw198 = 0_1000_10 +Ob1.10*2^-21 = -30.0`
`xw199 = 0_1000_11 +Ob1.11*2^-21 = -30.25`
`xw200 = 0_1000_00 +Ob1.00*2^-22 = -30.5`
`xw201 = 0_1000_01 +Ob1.01*2^-22 = -30.75`
`xw202 = 0_1000_10 +Ob1.10*2^-22 = -31.0`
`xw203 = 0_1000_11 +Ob1.11*2^-22 = -31.25`
`xw204 = 0_1000_00 +Ob1.00*2^-23 = -31.5`
`xw205 = 0_1000_01 +Ob1.01*2^-23 = -31.75`
`xw206 = 0_1000_10 +Ob1.10*2^-23 = -32.0`
`xw207 = 0_1000_11 +Ob1.11*2^-23 = -32.25`
`xw208 = 0_1000_00 +Ob1.00*2^-24 = -32.5`
`xw209 = 0_1000_01 +Ob1.01*2^-24 = -32.75`
`xw210 = 0_1000_10 +Ob1.10*2^-24 = -33.0`
`xw211 = 0_1000_11 +Ob1.11*2^-24 = -33.25`
`xw212 = 0_1000_00 +Ob1.00*2^-25 = -33.5`
`xw213 = 0_1000_01 +Ob1.01*2^-25 = -33.75`
`xw214 = 0_1000_10 +Ob1.10*2^-25 = -34.0`
`xw215 = 0_1000_11 +Ob1.11*2^-25 = -34.25`
`xw216 = 0_1000_00 +Ob1.00*2^-26 = -34.5`
`xw217 = 0_1000_01 +Ob1.01*2^-26 = -34.75`
`xw218 = 0_1000_10 +Ob1.10*2^-26 = -35.0`
`xw219 = 0_1000_11 +Ob1.11*2^-26 = -35.25`
`xw220 = 0_1000_00 +Ob1.00*2^-27 = -35.5`
`xw221 = 0_1000_01 +Ob1.01*2^-27 = -35.75`
`xw222 = 0_1000_10 +Ob1.10*2^-27 = -36.0`
`xw223 = 0_1000_11 +Ob1.11*2^-27 = -36.25`
`xw224 = 0_1000_00 +Ob1.00*2^-28 = -36.5`
`xw225 = 0_1000_01 +Ob1.01*2^-28 = -36.75`
`xw226 = 0_1000_10 +Ob1.10*2^-28 = -37.0`
`xw227 = 0_1000_11 +Ob1.11*2^-28 = -37.25`
`xw228 = 0_1000_00 +Ob1.00*2^-29 = -37.5`
`xw229 = 0_1000_01 +Ob1.01*2^-29 = -37.75`
`xw230 = 0_1000_10 +Ob1.10*2^-29 = -38.0`
`xw231 = 0_1000_11 +Ob1.11*2^-29 = -38.25`
`xw232 = 0_1000_00 +Ob1.00*2^-30 = -38.5`
`xw233 = 0_1000_01 +Ob1.01*2^-30 = -38.75`
`xw234 = 0_1000_10 +Ob1.10*2^-30 = -39.0`
`xw235 = 0_1000_11 +Ob1.11*2^-30 = -39.25`
`xw236 = 0_1000_00 +Ob1.00*2^-31 = -39.5`
`xw237 = 0_1000_01 +Ob1.01*2^-31 = -39.75`
`xw238 = 0_1000_10 +Ob1.10*2^-31 = -40.0`
`xw239 = 0_1000_11 +Ob1.11*2^-31 = -40.25`
`xw240 = 0_1000_00 +Ob1.00*2^-32 = -40.5`
`xw241 = 0_1000_01 +Ob1.01*2^-32 = -40.75`
`xw242 = 0_1000_10 +Ob1.10*2^-32 = -41.0`
`xw243 = 0_1000_11 +Ob1.11*2^-32 = -41.25`
`xw244 = 0_1000_00 +Ob1.00*2^-33 = -41.5`
`xw245 = 0_1000_01 +Ob1.01*2^-33 = -41.75`
`xw246 = 0_1000_10 +Ob1.10*2^-33 = -42.0`
`xw247 = 0_1000_11 +Ob1.11*2^-33 = -42.25`
`xw248 = 0_1000_00 +Ob1.00*2^-34 = -42.5`
`xw249 = 0_1000_01 +Ob1.01*2^-34 = -42.75`
`xw250 = 0_1000_10 +Ob1.10*2^-34 = -43.0`
`xw251 = 0_1000_11 +Ob1.11*2^-34 = -43.25`
`xw252 = 0_1000_00 +Ob1.00*2^-35 = -43.5`
`xw253 = 0_1000_01 +Ob1.01*2^-35 = -43.75`
`xw254 = 0_1000_10 +Ob1.10*2^-35 = -44.0`
`xw255 = 0_1000_11 +Ob1.11*2^-35 = -44.25`
`xw256 = 0_1000_00 +Ob1.00*2^-36 = -44.5`
`xw257 = 0_1000_01 +Ob1.01*2^-36 = -44.75`
`xw258 = 0_1000_10 +Ob1.10*2^-36 = -45.0`
`xw259 = 0_1000_11 +Ob1.11*2^-36 = -45.25`
`xw260 = 0_1000_00 +Ob1.00*2^-37 = -45.5`
`xw261 = 0_1000_01 +Ob1.01*2^-37 = -45.75`
`xw262 = 0_1000_10 +Ob1.10*2^-37 = -46.0`
`xw263 = 0_1000_11 +Ob1.11*2^-37 = -46.25`
`xw264 = 0_1000_00 +Ob1.00*2^-38 = -46.5`
`xw265 = 0_1000_01 +Ob1.01*2^-38 = -46.75`
`xw266 = 0_1000_10 +Ob1.10*2^-38 = -47.0`
`xw267 = 0_1000_11 +Ob1.11*2^-38 = -47.25`
`xw268 = 0_1000_00 +Ob1.00*2^-39 = -47.5`
`xw269 = 0_1000_01 +Ob1.01*2^-39 = -47.75`
`xw270 = 0_1000_10 +Ob1.10*2^-39 = -48.0`
`xw271 = 0_1000_11 +Ob1.11*2^-39 = -48.25`
`xw272 = 0_1000_00 +Ob1.00*2^-40 = -48.5`
`xw273 = 0_1000_01 +Ob1.01*2^-40 = -48.75`
`xw274 = 0_1000_10 +Ob1.10*2^-40 = -49.0`
`xw275 = 0_1000_11 +Ob1.11*2^-40 = -49.25`
`xw276 = 0_1000_00 +Ob1.00*2^-41 = -49.5`
`xw277 = 0_1000_01 +Ob1.01*2^-41 = -49.75`
`xw278 = 0_1000_10 +Ob1.10*2^-41 = -50.0`
`xw279 = 0_1000_11 +Ob1.11*2^-41 = -50.25`
`xw280 = 0_1000_00 +Ob1.00*2^-42 = -50.5`
`xw281 = 0_1000_01 +Ob1.01*2^-42 = -50.75`
`xw282 = 0_1000_10 +Ob1.10*2^-42 = -51.0`
`xw283 = 0_1000_11 +Ob1.11*2^-42 = -51.25`
`xw284 = 0_1000_00 +Ob1.00*2^-43 = -51.5`
`xw285 = 0_1000_01 +Ob1.01*2^-43 = -51.75`
`xw286 = 0_1000_10 +Ob1.10*2^-43 = -52.0`
`xw287 = 0_1000_11 +Ob1.11*2^-43 = -52.25`
`xw288 = 0_1000_00 +Ob1.00*2^-44 = -52.5`
`xw289 = 0_1000_01 +Ob1.01*2^-44 = -52.75`
`xw290 = 0_1000_10 +Ob1.10*2^-44 = -53.0`
`xw291 = 0_1000_11 +Ob1.11*2^-44 = -53.25`
`xw292 = 0_1000_00 +Ob1.00*2^-45 = -53.5`
`xw293 = 0_1000_01 +Ob1.01*2^-45 = -53.75`
`xw294 = 0_1000_10 +Ob1.10*2^-45 = -54.0`
`xw295 = 0_1000_11 +Ob1.11*2^-45 = -54.25`
`xw296 = 0_1000_00 +Ob1.00*2^-46 = -54.5`
`xw297 = 0_1000_01 +Ob1.01*2^-46 = -54.75`
`xw298 = 0_1000_10 +Ob1.10*2^-46 = -55.0`
`xw299 = 0_1000_11 +Ob1.11*2^-46 = -55.25`
`xw300 = 0_1000_00 +Ob1.00*2^-47 = -55.5`
`xw301 = 0_1000_01 +Ob1.01*2^-47 = -55.75`
`xw302 = 0_1000_10 +Ob1.10*2^-47 = -56.0`
`xw303 = 0_1000_11`

Plots



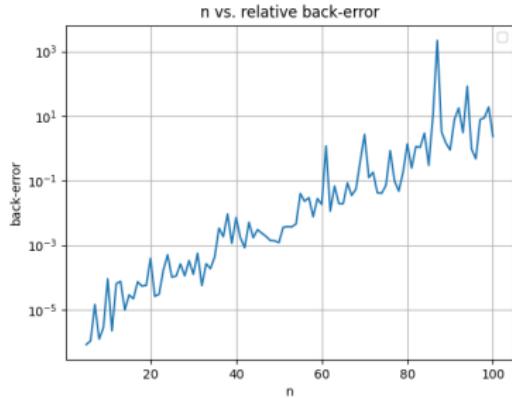
(a) $n = 30$ in binary8p4; $\epsilon_{mach} = 0.125$



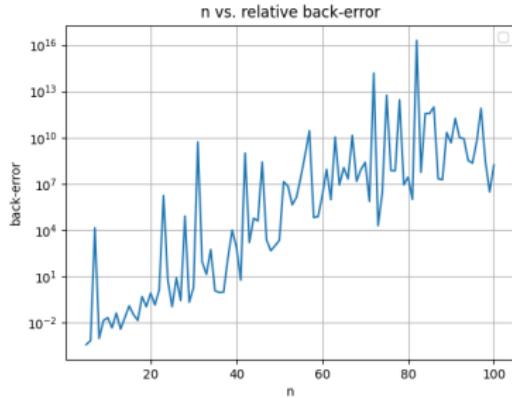
(b) $n = 30$ in binary8p3; $\epsilon_{mach} = 0.25$

Figure: Relative Backwards error vs condition number for different precisions after 5 iterations of iterative refinement in fp32

More plots



(a) $\text{cond} \approx 20$ in binary8p4;
 $\epsilon_{\text{mach}} = 0.125$



(b) $\text{cond} \approx 20$ in binary8p3;
 $\epsilon_{\text{mach}} = 0.25$

Figure: Relative Backwards error vs n for different precisions after 5 iterations of iterative refinement in fp32

Some hands-on experiences

Open your laptops!

<https://github.com/tlapack/tutorials>

We will make:

- LU in multi-precision
- Recursive matrix-matrix multiplication