

## Notes for Math 221, Lecture 5, Feb 8, 2022

Last time we spoke about dense matrix multiplication, and how to make it go fast, by doing as little communication as possible, or by doing asymptotically fewer floating point operations too (eg Strassen's method). Our goals for the rest of dense linear algebra (solving  $Ax=b$ , least squares, etc) will include these. But we have more goals:

Backward stability: get the right answer for a slightly different problem  $A + E$  where  $\text{norm}(E) = O(\text{macheps}) * \text{norm}(A)$

Exploit structure: go faster when matrix is symmetric, positive definite, or "sparse". Here sparse means that matrix depends on  $\ll n^2$  parameters, either because most matrix entries are zero, or if it is dense but the entries depend on few parameters.

For example, the Vandermonde matrix  $V(i,j) = x(i)^{(j-1)}$  is dense but depends on just  $n$  parameters  $x(1:n)$ , and there is an  $O(n^2)$  algorithm to solve  $V*y=b$  (hint: it is just polynomial interpolation in disguise).

One of the linear algebra problems will be computing the SVD  $A = U*\Sigma*V^T$ , which is one example of a matrix factorization of  $A$  into a product of simpler matrices. All the other algorithms we discuss will also explicitly or implicitly compute matrix factorizations:

Gaussian elimination:  $A = P*L*U$  where  $P$  is a permutation,  $L$  lower triangular,  $U$  upper triangular (with variations if  $A$  is symmetric, or positive definite too)

Least Squares:  $A = Q*R$ ,  $Q$  orthogonal,  $R$  upper triangular

Eigenvalue problem:  $A = Q*T*Q^H$  where  $Q$  unitary,  $T$  (nearly) triangular

And there are others, for example for computing modes of vibration (eigenvalues and eigenvectors) of a system with two matrices (mass and stiffness) or three (damping too), or deciding whether a linear control system with feedback can be made stable.

We begin with Gaussian elimination for a general (rectangular) matrix.

Def: Permutation matrix: identity with permuted rows

Facts: let  $P, P1$  etc be permutation matrices

$P*X$  = row permutation of  $X$  and  $X*P$  = column permutation of  $X$

$P1*P2$  also a permutation matrix

$\text{inv}(P) = P^T$ , i.e.  $P$  is orthogonal (enough to check that  $\text{diag}(P*P^T) = \text{all ones}$ )

$\det(P) = \pm 1$

to store and multiply by  $P$ , just keep track of locations of ones (cheap)

Thm (LU decomposition): Given any  $m \times n$  full rank matrix  $A$ , with  $m \geq n$ , there exist

$m \times m$  permutation matrix  $P$

$m \times n$  unit lower triangular matrix  $L$  (unit means  $L(i,i)=1$ )

$n \times n$  nonsingular upper triangular matrix  $U$   
such that  $A = P * L * U$

In the proof below, we will use induction, which is just Gaussian elimination.

Cor: If  $A$   $n \times n$  and nonsingular, there exist an  $n \times n$  permutation  $P$ , unit lower triangular  $L$ , and nonsingular upper triangular  $U$  such that  $A = P * L * U$

To solve  $A * x = b$ :

- (1) factor  $A = P * L * U$  (expensive part, cost =  $(2/3)n^3$ )
- (2) Solve  $P * L * U * x = b$  for  $L * U * x = P^T * b$  by permuting  $b$   
... cost =  $O(n)$
- (3) Solve  $L * U * x = P^T * b$  for  $U * x = \text{inv}(L) * P^T * b$  using forward substitution with  $L$  ... cost =  $n^2$  flops
- (4) Solve  $U * x = \text{inv}(L) * P^T * b$  for  $x = \text{inv}(U) * \text{inv}(L) * P^T * b = \text{inv}(A) * b$  using back substitution with  $U$  ... cost =  $n^2$  flops

If we are given another  $b'$  vector, we can solve  $A * x' = b'$  in  $2n^2$  more flops.

Note: We do not compute  $\text{inv}(A)$  and multiply  $x = \text{inv}(A) * b$  because it is  
(1) 3x more expensive in dense case, can be much worse in sparse case ( $O(n^2)$  times!)  
(2) not as numerically stable

Proof of Theorem:

If  $A$  is full rank, the first column is nonzero, so there is a permutation  $P$  such that  $(P * A)(1,1)$  is nonzero. Write

$$P * A = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & I \end{bmatrix} * \begin{bmatrix} A_{11} & A_{12} \\ 0 & A_{22} - A_{21} * A_{12} / A_{11} \end{bmatrix}$$

where  $A_{11}$  is  $1 \times 1$ ,  $A_{21}$  is  $(m-1) \times 1$ ,  $A_{12}$  is  $1 \times (n-1)$  and  $A_{22}$  is  $(m-1) \times (n-1)$ .

Now  $A$  full (column) rank  $\Rightarrow P * A$  full rank  $\Rightarrow S = A_{22} - A_{21} * A_{12} / A_{11}$  is full rank. (Otherwise, if some nonzero linear combination of columns of  $S$  were  $0$ , say  $S * x = 0$ , then a linear combination of columns of  $A$  would be zero,  $A * [-A_{12} * x / A_{11}; x] = 0$ , contradicting  $A$  being full column rank. More simply in the square case:

$$\det(A) \text{ nonzero} \Rightarrow 0 \neq \det(\text{first factor}) * \det(\text{second factor}) \\ = 1 * A_{11} * \det(S) \Rightarrow \det(S) \text{ nonzero}$$

Notation:  $S$  called Schur complement

Now apply induction:  $S = P' * L' * U'$ , so

$$P * A = \begin{bmatrix} 1 & 0 \\ A_{21}/A_{11} & I \end{bmatrix} * \begin{bmatrix} A_{11} & A_{12} \\ 0 & P' * L' * U' \end{bmatrix} \\ = \begin{bmatrix} 1 & 0 \\ A_{21}/A_{11} & P' * L' \end{bmatrix} * \begin{bmatrix} A_{11} & A_{12} \\ 0 & U' \end{bmatrix} \\ = \begin{bmatrix} 1 & 0 \\ 0 & P' \end{bmatrix} * \begin{bmatrix} 1 & 0 \\ P'^T * (A_{21}/A_{11}) & L' \end{bmatrix} * \begin{bmatrix} A_{11} & A_{12} \\ 0 & U' \end{bmatrix} \\ = P'' * L * U$$

= permutation \* unit\_lower\_triangular \* upper\_triangular

where  $U$  is nonsingular.

So  $A = P^T P^{-1} * L * U = \text{permutation} * L * U$  QED.

(We pause the recorded lecture here.)

Expressing this as an algorithm, we get the following (ignore permutations first)

```
for i = 1 to n
  ... i = 1 performs algebra shown above,
  ... i>1 applies the same algorithm recursively to the Schur
  ... complement
  L(i,i) = 1, L(i+1:n,i) = A(i+1:n,i)/A(i,i)
  U(i,i:n) = A(i,i:n)
  if (i < n) A(i+1:n,i+1:n) = A(i+1:n,i+1:n) - L(i+1:n,i)*U(i,i+1:n)
```

Add permutations: after "for i=1 to n", add:

```
  if A(i,i) zero and A(j,i) nonzero, swap rows i and j of L and A;
  record swap
  ... choosing A(j,i) is called "pivoting", more on this below
```

Don't waste space:

```
row i of U overwrites row of i A: omit U(i,i:n) = A(i,i:n)
col i of L (below diagonal) overwrites same entries of A, which are
zeroed out: change first line to A(i+1:n,i) = A(i+1:n,i)/A(i,i)
only need to loop from i = 1 to n-1, and change last line from
  A(i+1:n,i+1:n) = A(i+1:n,i+1:n) - L(i+1:n,i)*U(i,i+1:n)
to
  A(i+1:n,i+1:n) = A(i+1:n,i+1:n) - A(i+1:n,i)*A(i,i+1:n)
```

Finally, we get

```
for i=1 to n-1
  if A(i,i) zero and A(j,i) nonzero, swap rows i and j of A;
  record swap
  A(i+1:n,i) = A(i+1:n,i)/A(i,i) ... call to BLAS1 routine scal
  A(i+1:n,i+1:n) = A(i+1:n,i+1:n) - A(i+1:n,i)*A(i,i+1:n)
  ... call BLAS2 ger
```

(draw picture of intermediate step)

When we're done:

```
U overwrites upper triangle and diagonal of A
L (below diagonal) overwrites A (below diagonal)
The unit diagonal entries L(i,i)=1 are not stored.
```

To see that this is the same Gaussian Elimination you learned long ago, start from

```
for i=1 to n-1 ... for each column i
  for j=i+1 to n ... add a multiple of row i to row j to zero out
  ... entry (j,i) below diagonal
  m = A(j,i)/A(i,i)
  A(j,i:n) = A(j,i:n) - m*A(i,i:n)
```

"Optimize" this by

- (1) not bothering to compute the entries below the diagonal you know are zero:  
change last line to  $A(j,i+1:n) = A(j,i+1:n) - m \cdot A(i,i+1:n)$
- (2) compute all the multipliers  $m$  first, store them in zeroed-out locations:  
for  $i=1$  to  $n-1$   
  for  $j=i+1$  to  $n$   
     $A(j,i) = A(j,i)/A(i,i)$   
  for  $j=i+1$  to  $n$   
     $A(j,i+1:n) = A(j,i+1:n) - A(j,i) \cdot A(i,i+1:n)$
- (3) combine loops into single expressions to get same as before:  
for  $i=1$  to  $n-1$   
   $A(i+1:n,i) = A(i+1:n,i)/A(i,i)$   
   $A(i+1:n,i+1:n) = A(i+1:n,i+1:n) - A(i+1:n,i) \cdot A(i,i+1:n)$

The cost is  $\sum_{i=1}^{n-1} 2(n-i)^2 = (2/3)n^3 + O(n^2)$  multiplies and adds.

(We pause the recorded lecture here.)

We need to do pivoting carefully, i.e. choosing which nonzero to put on the diagonal, even if  $A(i,i)$  is nonzero, to get a backward stable result, i.e.  $P \cdot L \cdot U = A + E$  where  $E$  is small compared to  $A$ . For example, if we run in single precision, so with 7 decimal digits, and take

$$A = \begin{bmatrix} 1e-8 & 1 \\ 1 & 1 \end{bmatrix}, \quad \text{inv}(A) \sim \begin{bmatrix} -1 & 1 \\ 1 & -1e-8 \end{bmatrix},$$

then  $\text{cond}(A) \sim 2.6$ , really small, so we expect a good answer. But

$$L = \begin{bmatrix} 1 & 0 \\ 1e8 & 1 \end{bmatrix} \quad U = \begin{bmatrix} 1e-8 & 1 \\ 0 & \text{fl}(1-1e8 \cdot 1) = -1e8 \end{bmatrix}$$

$$\text{and so } L \cdot U = \begin{bmatrix} 1e-8 & 1 \\ 1 & 0 \end{bmatrix}$$

which is very different from  $A$  in the  $(2,2)$  entry. In fact, we'd get the same (wrong)  $L$  and  $U$  if  $A(2,2)$  were  $.5$ ,  $-1$ , etc. because the operation  $\text{fl}(A(2,2)-1e8 \cdot 1)$  "forgets"  $A(2,2)$  if  $A(2,2)$  is small enough, say  $O(1)$ . So going on to solve  $A \cdot x = b$  using  $L$  and  $U$  will give very wrong answers.

If we instead pivot (swap rows 1 and 2) so  $A(1,1)=1$ , then we get full accuracy. The intuition is that we want to pick a large entry of  $A$  to be the "pivot"  $A_{11}$ , and repeat this at each step, for reasons we formalize below. To motivate the statement of the result, recall HW 1.10, where you showed  $C = \text{fl}(A \cdot B) = A \cdot B + E$  where  $|E| \leq n \cdot \text{macheps} \cdot |A| \cdot |B|$ .

Since  $A = P \cdot L \cdot U$  in exact arithmetic, it is perhaps no surprise that

Thm (backward error analysis of LU decomposition):

If  $P$ ,  $L$  and  $U$  are computed by the above algorithm, then

$$A - E = P \cdot L \cdot U \text{ where } |E| \leq n \cdot \text{macheps} \cdot |P| \cdot |L| \cdot |U|.$$

Before proving this, we state and prove a corollary, bounding the backward error of solving  $Ax=b$ :

Corollary: Suppose we solve  $Ax=b$  by Gaussian elimination, following by forward and back substitution with  $L$  and  $U$  as described above. Then the computed results  $\hat{x}$  satisfies

$$(A-F)\hat{x} = b \text{ where } |F| \leq 3n \text{macheps} P \|L\| \|U\|.$$

Proof: Here, and in the proof of the Thm, we assume  $P=I$  for simplicity of notation; just imagine running without pivoting on the matrix  $P^T A$ . In HW 1.11, you showed that the computed solution  $\hat{y}$  of  $L*y=b$  satisfied

$$(L+dL)*\hat{y} = b \text{ with } |dL| \leq n \text{macheps} \|L\|,$$

and that the computed solution  $\hat{x}$  of  $U*x=\hat{y}$  satisfied

$$(U+dU)*\hat{x} = \hat{y} \text{ with } |dU| \leq n \text{macheps} \|U\|.$$

Thus

$$\begin{aligned} b &= (L+dL)*\hat{y} = (L+dL)*(U+dU)*\hat{x} \\ &= (L*U + dL*U + L*dU + dL*dU)*\hat{x} \\ &= (A - E + dL*U + L*dU + dL*dU)*\hat{x} \quad \text{by the Thm} \\ &= (A - F)*\hat{x} \end{aligned}$$

where  $F = E - dL*U - L*dU - dL*dU$

$$\begin{aligned} \text{so } |F| &\leq |E| + |dL*U| + |L*dU| + |dL*dU| \\ &\leq |E| + |dL|*|U| + |L|*|dU| + |dL|*|dU| \\ &\leq n \text{macheps} \|L\| \|U\| \quad \dots \text{ by the Thm} \\ &\quad + n \text{macheps} \|L\| \|U\| \quad \dots \text{ from the bound on } |dL| \\ &\quad + n \text{macheps} \|L\| \|U\| \quad \dots \text{ from the bound on } |dU| \\ &\quad + (n \text{macheps})^2 \|L\| \|U\| \quad \dots \text{ from the bounds on } |dL| \text{ and } |dU| \\ &= (3n \text{macheps} + O(\text{macheps}^2)) \|L\| \|U\| \end{aligned}$$

What the Thm (proven below) and Corollary tell us is that we need

$$\text{norm}(F) \leq 3 \text{norm}(\text{macheps} \|L\| \|U\|) = O(\text{macheps} \text{norm}(A))$$

for the algorithm to be backward stable, i.e.

$$\text{norm}(\|L\| \|U\|) = O(\text{norm}(A)).$$

This in turn depends on pivoting carefully.

Proof of Thm: Recall that for simplicity we assume  $P=I$ . If we trace through the algorithm and ask how  $U(i,j)$  is computed, we see that we start with  $A(i,j)$  when  $i \leq j$  and repeatedly subtract  $L(i,k)*U(k,j)$  for  $k=1,2,\dots,i-1$  until we get

$$U(i,j) = A(i,j) - \sum_{k=1}^{i-1} L(i,k)*U(k,j)$$

which, not surprisingly, is what you get when you take the  $(i,j)$  entry of  $A = L*U$  and solve for  $U(i,j)$ . So  $U(i,j)$  is essentially computed by a dot product of the (previously computed)  $i$ -th row  $L$  and  $j$ -th column of  $U$ , and using the same error analysis approach for dot products as in HW 1.10, we get the result.

Similarly, when  $i > j$  we get  $L(i,j)$  by

starting with  $A(i,j)$  and subtracting  $L(i,k)*U(k,j)$  for  $k=1,2,\dots,j-1$  dividing the resulting sum by  $A(i,i)$

or

$L(i,j) = (A(i,j) - \sum_{k=1}^{j-1} L(i,k)U(k,j))/A(i,i)$   
 so again we have a dot product, followed by a division, and a similar approach works.

Now we discuss how to pivot so that  $\text{norm}(|L|*|U|) = O(\text{norm}(A))$ , or not much larger.

Def: we call  $g = \text{norm}(|L|*|U|)/\text{norm}(A)$  the pivot growth factor. (Note: this is defined somewhat differently in the literature, but is very similar).

The unstable 2x2 example above, where  $L(2,1) = 1/1e-8 = 1e8$ , suggests that we choose the pivot  $A_{11}$  to be as large as possible, so entries of  $L$  are as small as possible.

(1) Simplest, and standard, approach, used in most libraries:

"Partial pivoting" (also called GEPP).

At each step, permute rows so  $A_{11}$  is the largest (in absolute value) entry in column

Then  $L_{21} = A_{21}/A_{11}$  has  $|L_{21}| \leq 1$

Theorem (easy): with GEPP,  $|L| \leq 1$  and  $\max(|U(:,i)|) \leq 2^{(n-1)} \max(|A(:,i)|)$

Bad news: worst case is terrible; even for  $n=24$  in singular precision, all wrong

Good news: hardly ever happens (only very small family of matrices where this occurs)

Empirical observation, with some justification:  $g_{PP} < n^{(2/3)}$

If all entries of matrix were "random", this would be true; as you perform pivoting, they seem to get more random

(2) Complete pivoting (GECP): permute rows and columns so that  $A_{11}$  largest entry in the whole matrix; again repeat at every step.

Get  $A = P_r * L * U * P_c$  where  $P_r$  and  $P_c$  are both permutations.

Theorem:  $g_{CP} < n^{(\log n / 4)}$

Empirical:  $g_{CP} < n^{(1/2)}$

Long-standing Conjecture:  $g_{CP} < n$  (false, but nearly true)

More expensive, hardly used, not in most libraries

(3) Tournament pivoting - something new, needed to minimize communication (#messages): will present it later.

(4) Threshold pivoting: this and similar schemes try to preserve sparsity while maintaining stability

Figures 2.1 and 2.2 in the book show empirical results that demonstrate that these bounds for GEPP and GECP are "worst case", and pessimistic in general.

Altogether, our worst-case error bound is

$$\begin{aligned} \text{norm}(x-\hat{x})/\text{norm}(x) &\leq \text{cond}(A) * (\text{backward error}) \\ &\leq \text{cond}(A) * 3 * n * \text{macheps} * \text{pivot\_growth} \end{aligned}$$

where we can estimate  $\text{cond}(A)$  with  $O(n^2)$  flops after LU decomposition, and can bound pivot growth in  $O(n^2)$  work too.

(We pause the recorded lecture here.)

What if this error is too large for your application, or too slow? We can run iterative refinement, aka Newton's method, using mixed precision, doing most of the work (the  $O(n^3)$  part) in low (fast) precision, and a little more ( $O(n^2)$ ) in high precision. In the following algorithm, low/high precision could mean single/double, half/single, bfloat16/single, double/quad, or other combinations:

```

Do GEPP to solve  $Ax=b$  in low precision, call solution  $x(1)$ 
... this is the  $O(n^3)$  part
i = 1
repeat
  r =  $A*x(i) - b$     ... in high precision, but costs just  $O(n^2)$ ;
                    ... round final result r to low precision
  solve  $A*d = r$     ... in low precision using existing LU
                    factors, costs just  $O(n^2)$ 
  update  $x(i+1) = x(i) - d$  ... in low precision, costs  $O(n)$ 
until "convergence"

```

We compute  $r$  in high precision, because otherwise the computed residual  $r$  may be mostly roundoff error, so the correction  $d$  is mostly noise, and there is no guarantee of progress (though some benefits have been proven, and so both versions are in LAPACK).

Testing "convergence" depends on one's goals. We mention two:

(1) Getting a small backward error in high precision:  
 $\text{norm}(A*x_{\text{computed}} - b) = O(\text{macheps}_{\text{high}})*\text{norm}(A)*\text{norm}(x_{\text{computed}})$   
 or a warning that the matrix is too ill-conditioned to converge.  
 This is straightforward to implement, since we need to compute the residual anyway. This is motivated by the availability of hardware accelerators for machine learning (from Google, Nvidia, Intel, etc.) that can do 16-bit arithmetic much faster than higher precisions. Some of them also accumulate dot-products internally in 32-bit precision, giving us the high precision residual automatically. There is also recent work on alternatives to the simple Newton's iteration above, with better convergence behavior, see the links to "A Survey of Numerical Methods Utilizing Mixed Precision Arithmetic" on the class webpage, or the recent Nvidia blog post.  
 (The algorithm used, GMRES, will be discussed in Chap. 6.)

(2) Getting a small relative error in low precision:  
 $\text{norm}(x_{\text{computed}} - x_{\text{true}})/\text{norm}(x_{\text{true}}) = O(\text{macheps}_{\text{low}})$   
 or a warning that the matrix is too ill-conditioned to do this.  
 This is tricky, because we have to avoid getting fooled by a very ill-conditioned matrix that appears to "accidentally" converge; for details see [www.netlib.org/lapack/lawnspdf/lawn165.pdf](http://www.netlib.org/lapack/lawnspdf/lawn165.pdf). Slides 33-35 of [www.cs.berkeley.edu/~demmel/Future\\_Sca-LAPACK\\_v7.ppt](http://www.cs.berkeley.edu/~demmel/Future_Sca-LAPACK_v7.ppt) show empirical results for millions of randomly generated test cases, using single/double, with condition numbers ranging from

1 to well beyond  $1/\text{macheps\_low}$ : The relative error of LU without refinement is indeed usually close to  $\text{condition\_number} * \text{macheps\_low}$ , but LU with refinement is much better: the relative error is  $O(\text{macheps\_low})$  as long as the condition number is less than about  $1/\text{macheps\_low}$ , which the algorithm reports to the user; when the condition number is larger, the relative error can rise as high as 1, and this lack of convergence is also reported to the user. This algorithm is available in LAPACK as `sgesvxx`, `dgesvxx`, etc.

(We pause the recorded lecture here.)

Now we return to minimizing communication. Historically, GEPP was rewritten to do most of its work by calling the matrix multiplication routine in the BLAS3, which often led to high performance in libraries like LAPACK and ScaLAPACK. The idea is based on a similar induction proof as for classical LU decomposition, but instead of working on one column at a time, the algorithm works  $b$  columns at a time, where  $b$  is a block size analogous to one used in matrix multiplication. Ignoring pivoting for simplicity we write

$$A = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix}$$

... where  $A_{11}$  is  $b$ -by- $b$ ,  $A_{21}$   $(n-b)$ -by- $b$ ,  $A_{12}$   $b$ -by- $(n-b)$  and  
 ...  $A_{22}$   $(n-b)$ -by- $(n-b)$

$$= \begin{bmatrix} L_{11} * U_{11} & A_{12} \\ L_{21} * U_{11} & A_{22} \end{bmatrix}$$

... where we have performed LU decomposition on the  $n$ -by- $b$  matrix  
 ...  $\begin{bmatrix} A_{11} \\ A_{21} \end{bmatrix} = \begin{bmatrix} L_{11} \\ L_{21} \end{bmatrix} * U_{11}$   
 ... by using the previous "unblocked" algorithm

$$= \begin{bmatrix} L_{11} * U_{11} & L_{11} * U_{12} \\ L_{21} * U_{11} & A_{22} \end{bmatrix}$$

... where we have solved the triangular system of equations  
 ...  $A_{12} = L_{11} * U_{12}$  for  $U_{12}$  by calling BLAS3 routine TRSM

$$= \begin{bmatrix} L_{11} & 0 \\ L_{21} & I \end{bmatrix} * \begin{bmatrix} U_{11} & U_{12} \\ 0 & A_{22} - L_{21} * U_{12} \end{bmatrix}$$

... where we compute  $A_{22hat} = A_{22} - L_{21} * U_{12}$  by calling BLAS3 GEMM

The inductive step then applies the same procedure to  $A_{22hat}$ . While this often works well in practice, since most of the work is done by the BLAS3, it turns out that there are some combinations of dimension  $n$  and cache size  $M$  for which  $b$  cannot be chosen to attain the lower bound.

Just as there was a recursive, cache oblivious version of `matmul` that minimized communication without explicitly depending on the cache size  $M$ , we can do the same for LU (due to Toledo, 1997):

Here is a high level description of the algorithm:

```
"Do LU on left half of matrix
  Update right half (U at top, Schur complement at bottom)
```

Do LU on Schur complement"

```
function [L,U] = RLU(A) ... RLU = Recursive LU
... assume A is n by m with n >= m, m a power of 2
if m=1 ... one column
    pivot so largest entry on diagonal, update rest of matrix
    L = A/A11, U = A11
else
    ... write A = [ A11 A12 ], L1 = [ L11 ]
    ...           [ A21 A22 ]     [ L12 ]
    ... where A11, A12, L11, U1 and U2 are m/2 by m/2
    ...       A21, A22 and L12 are n-m/2 by m/2
    [L1,U1] = RLU([ A11 ]) ... LU of left half of A
                [ A21 ]
    Solve A12=L11*U12 ... update U, in upper part of right half of A
    A22 = A22 - L21*U12 ... update Schur complement
    [L2,U2] = RLU( A22 } ... LU on Schur complement
    L = [ L1, [0;L2] ] ... return complete nxm L factor
    U = [ U1 U12 ]     ... return complete mxm U factor
        [ 0 U2 ]
```

Correctness follows by induction on m, as does showing when m=n:

A(n) = #arithmetic operations =  $(2/3)n^3 + O(n^2)$ , and

W(n) = #words moved =  $O(n^3/\sqrt{M})$

As stated, this algorithm only minimizes the #words moved, not the #messages. To minimize #messages, we either need to use a different pivoting scheme (tournament pivoting, to be discussed later, see "CALU: A communication optimal LU Factorization Algorithm") or a more complicated data structure, (see "Communication efficient Gaussian elimination with Partial Pivoting using a Shape Morphing Data Layout"); both papers are at [bebop.cs.berkeley.edu](http://bebop.cs.berkeley.edu).

Fact: if we do  $L21*U12$  by Strassen, and  $L11 \setminus A12$  by

(1) inverting L11 by divide-and-conquer:

$$\text{inv} \begin{pmatrix} T11 & T12 \\ 0 & T22 \end{pmatrix} = \begin{pmatrix} \text{inv}(T11) & -\text{inv}(T11)*T12*\text{inv}(T22) \\ 0 & \text{inv}(T22) \end{pmatrix}$$

(2) multiplying  $\text{inv}(L11)*A12$  etc by Strassen

then the RLU algorithm costs  $O(n^{(\log_2 7)})$  like Strassen, but can be slightly less stable than usual  $O(n^3)$  version of GEPP (see [arxiv.org/abs/math.NA/0612264](http://arxiv.org/abs/math.NA/0612264)).

Where to find all this implemented (all "blocked" algorithms, unless marked as "recursive"):

Matlab: A\b or [P,L,U]=lu(A), rcond or condest

LAPACK:

xGETRF just for GEPP where x = S/D/C/Z

xGETRF2 recursive

xGESV to solve  $A*x=b$

xGESVX for condition estimation, iterative refinement with

no extra precision  
xGESVXX for condition estimation, iterative refinement with  
extra precision  
xGECOM for condition estimation alone  
ScaLAPACK: PxGETRF etc