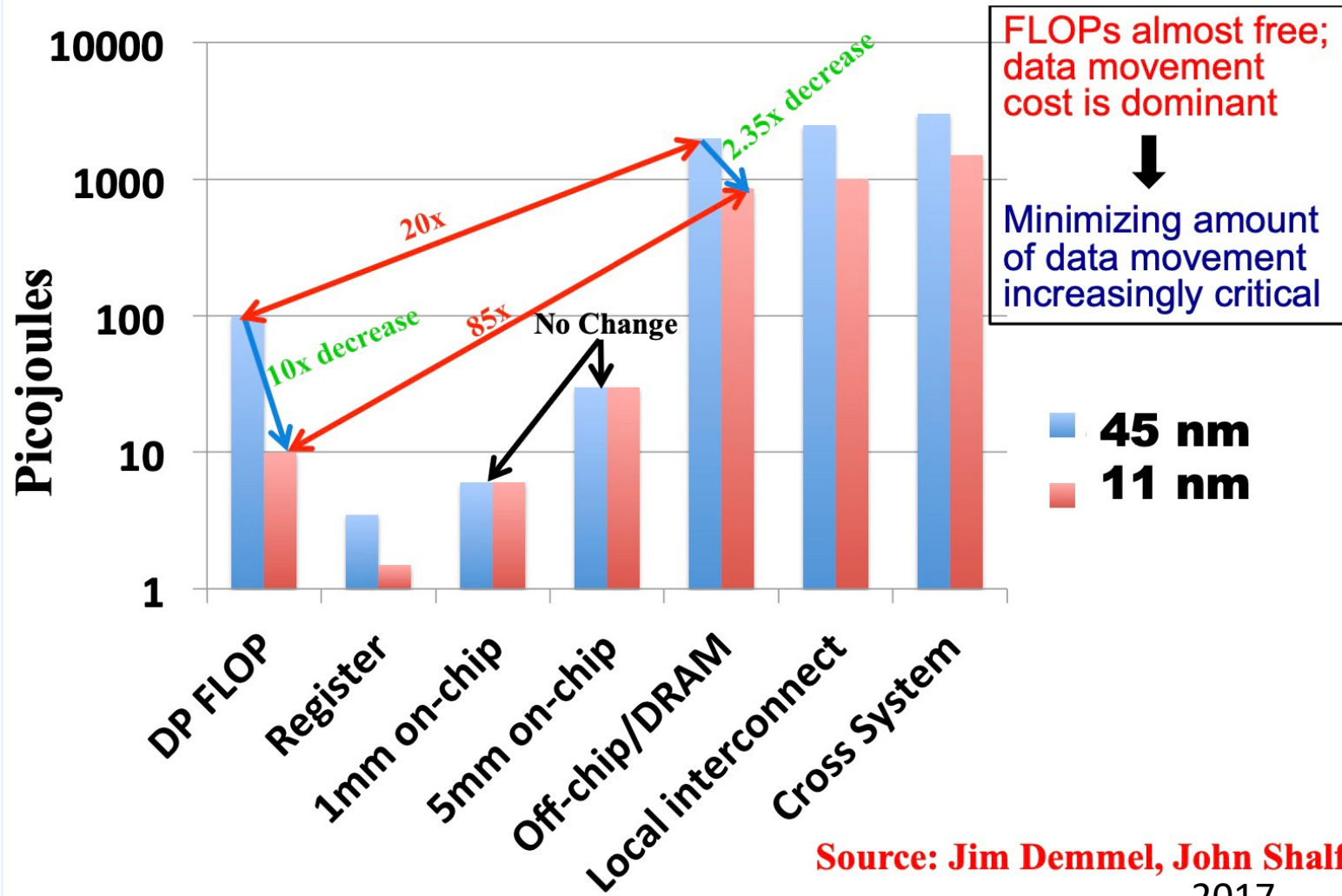# The I/O Requirements of Various Numerical Linear Algebra Kernels

Julien Langou

Wednesday June 26th 2024

# Data Movement Cost: Energy Trends



Source: Jim Demmel, John Shalf
2017

H. T. Kung

# I/O complexity: The red-blue pebble game

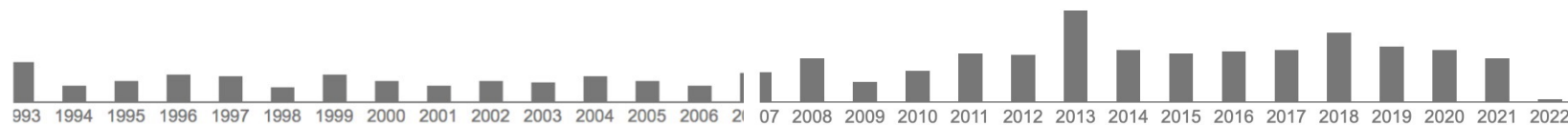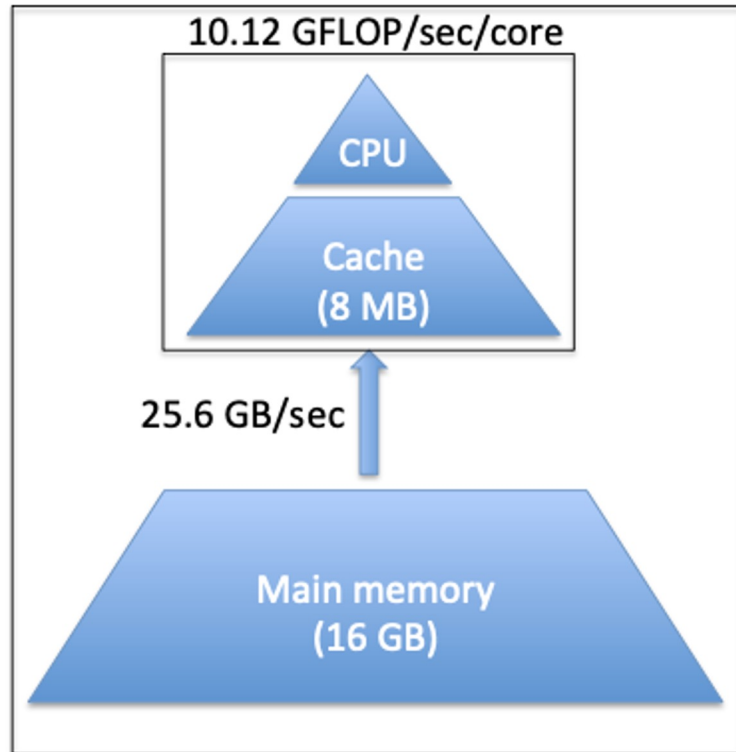| | |
|---|---|
| Authors | Hong Jia-Wei, HT Kung |
| Publication date | 1981/5/11 |
| Conference | Proceedings of the thirteenth annual ACM symposium on Theory of computing |
| Pages | 326-333 |
| Publisher | ACM |
| Description | Abstract In this paper, the red-blue pebble game is proposed to model the input-output complexity of algorithms. Using the pebble game formulation, a number of lower bound results for the I/O requirement are proven. For example, it is shown that to perform the n-point FFT or the ordinary n× n matrix multiplication algorithm with O (S) memory, at least &Ohgr;(n log n/log S) or &Ohgr;(n 3/@@@@ S), respectively, time is needed for the I/O. Similar results are obtained for algorithms for several other problems. All of the lower ... |
| Total citations | Cited by 401 |

993 1994 1995 1996 1997 1998 1999 2000 2001 2002 2003 2004 2005 2006 2( 07 2008 2009 2010 2011 2012 2013 2014 2015 2016 2017 2018 2019 2020 2021 2022

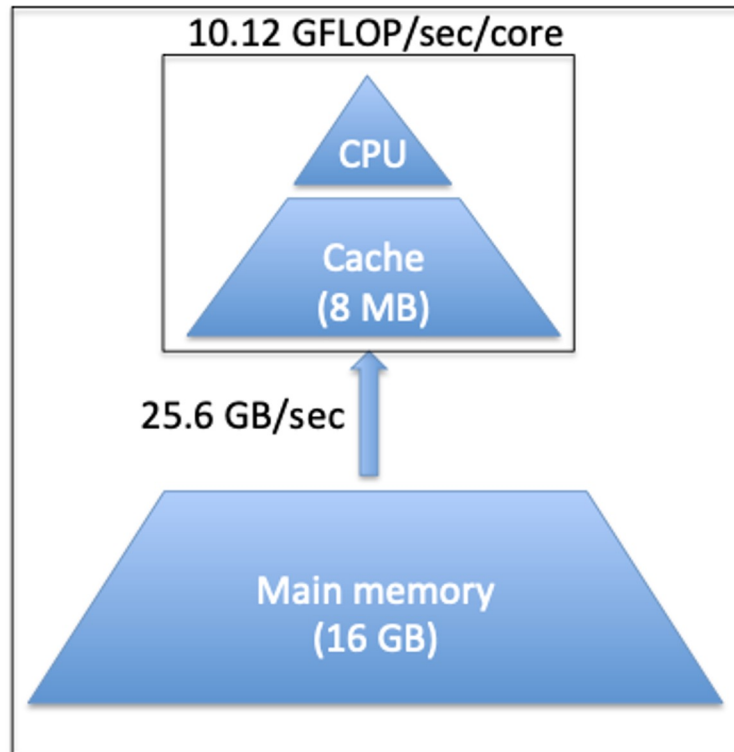| | |
|---|---|
| Scholar articles | I/O complexity: The red-blue pebble game |
| | H Jia-Wei, HT Kung - Proceedings of the thirteenth annual ACM symposium ..., 1981 |
| | Cited by 401 - Related articles - All 3 versions |

# IO lower bound for computer programs

Intel Xeon Processor E5520 (Nehalem)



10.12 GFLOP/sec/core

CPU

Cache
(8 MB)

25.6 GB/sec

Main memory
(16 GB)

# IO lower bound for computer programs
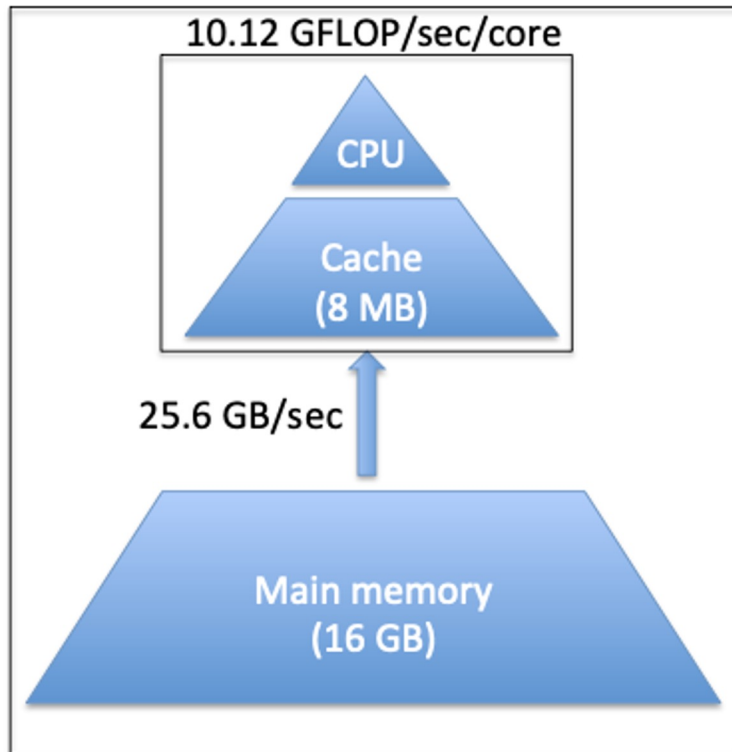
Intel Xeon Processor E5520 (Nehalem)



For MM

The number of words transferred between slow and fast memory is at least

$$2 \left( \frac{n^3}{\sqrt{M}} \right) - 2M.$$

where n is matrix size and M is size of cache.

See: A Tight I/O Lower Bound for Matrix Multiplication. T. M. Smith, B. Lowery, J. Langou, R. A. van de Geijn https://arxiv.org/abs/1702.02017

# IO lower bound for computer programs

Intel Xeon Processor E5520 (Nehalem)



One core Intel Xeon Processor E5520 (nehalem)
$\beta^{-1} = 580 \cdot 10^6$ words/sec    $\gamma^{-1} = 10.12 \cdot 10^9$ flops/sec    $M = 10^6$ words

DGEMM on one core Intel Xeon Processor E5520 (nehalem)

For MM

The number of words transferred between slow and fast memory is at least
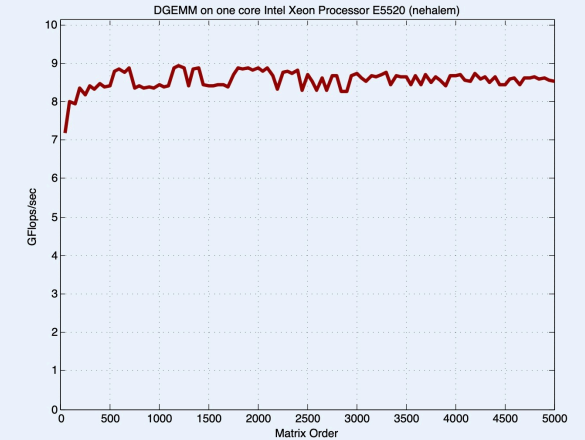
$$2\left(\frac{n^3}{\sqrt{M}}\right) - 2M.$$

where n is matrix size and M is size of cache.

See: A Tight I/O Lower Bound for Matrix Multiplication. T. M. Smith, B. Lowery, J. Langou, R. A. van de Geijn https://arxiv.org/abs/1702.02017

```
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];
```

**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];
```

## 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];
    }
  }
}
```

## 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]);
}
```

## 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]);
}
```

## cgs – Classical Gram-Schmidt

```c
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

```c
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

```c
for(k = 0; k < N; k++){
  norma2 = 0.e+00;
  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 ) ;

  tau[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[k][k]= ( A[k][k] > 0 ) ? ( - norma ) : ( norma ) ;

  for(j = k+1; j < N; j++){
    tau[j] = A[k][j];
    for(i = k+1; i < M; i++){
      tau[j] += A[i][k] * A[i][j];
    }
    tau[j] = tau[k] * tau[j];
    A[k][j] = A[k][j] - tau[j];
    for(i = k+1; i < M; i++){
      A[i][j] = A[i][j] - A[i][k] * tau[j];
    }
  }
}
```

## v2q – orqr2 – Construction of Q from Householder

```c
for(k = N-1; k > -1; k--){
  for(j = k+1; j < N; j++){
    tau[j] = 0.e+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(j = k+1; j < N; j++){
    A[k][j] = -tau[j];
  }
  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];
  }
}
```

| gebd2 − reduction to bidiagonalization | gehd2 − reduction to Hessenberg form |
|---|---|

```
for(k = 0; k < N; k++){
  norma2 = 0.e+00;
  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[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+00;
  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 ) ;
  taup[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][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 * taup[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] ;
    }
  }
}
```

```
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 ( i = j+2; i < n; i++ ) {
    A[i][j] /= A[j+1][j] ;
  }
  A[j+1][j] = ( A[j+1][j] > 0 ) ? ( - norma ) : ( norma ) ;
  for (i = j+1; i < n; i++) {
    tmp[i] = A[j+1][i] ;
    for (k = j+2; k < n; k++) {
      tmp[i] += A[k][j] * A[k][i];
    }
  }
  for (i = j+1; i < n; i++) {
    tmp[i] *= tau ;
  }
  for (i = j+1 ; i < n ; i++) {
    A[j+1][i] -= tmp[i] ;
  }
  for (i = j+2; i < n; i++) {
    for (k = j+1; k < n; k++) {
      A[i][k] -= A[i][j] * tmp[k];
    }
  }
  for (i = 0; i < n; i++) {
    tmp[i] = A[i][j+1];
    for (k = j+2; k < n; k++) {
      tmp[i] += A[i][k] * A[k][j];
    }
  }
  for (i = 0; i < n; i++) {
    tmp[i] *= tau ;
  }
  for (i = 0; i < n; i++) {
    A[i][j+1] -= tmp[i] ;
  }
  for (i = 0; i < n; i++) {
    for (k = j+2; k < n; k++) {
      A[i][k] -= tmp[i] * A[k][j] ;
    }
  }
}
```

| sytd2 − reduction to symmetric tridiagonal form | gghd2 − reduction to triangular Hessenberg 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 ) ? (A[i+1][i] + norma) : (A[i+1][i] - norma) ;
    taui = 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[k][j] * A[j][i];
        }
        for(j = k+1; j < n; j++){
            tau[k-1] += A[j][k] * A[j][i];
        }
        tau[k-1] *= taui;
    }
    alpha = tau[i];
    for(k = i+2; k < n; k++){
        alpha += tau[k-1] * A[k][i];
    }
    alpha = -0.5e+00 * taui * 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 * tau[i] ;
    for(j = i+2; 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] = taui;
}
```

```
for (j = 0; j < _PB_N-2; j++) {
    for (i = _PB_N-2; i > j; i--) {
        nrm = SQRT_FUN ( A[i][j] * A[i][j] + A[i+1][j] * A[i+1][j] );
        c = A[i][j] / nrm;
        s = A[i+1][j] / nrm;
        A[i][j] = nrm;
        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;
        }
        nrm = SQRT_FUN ( B[i+1][i+1] * B[i+1][i+1] + B[i+1][i] * B[i+1][i] );
        c = B[i+1][i+1] / nrm;
        s = B[i+1][i] / nrm;
        B[i+1][i+1] = nrm;
        B[i+1][i] = SCALAR_VAL(0.0);
        for (k = 0; k <= i; k++) {
            tmp = c * B[k][i] - s * B[k][i+1];
            B[k][i+1] = s * B[k][i] + c * B[k][i+1];
            B[k][i] = 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;
        }
    }
}
```
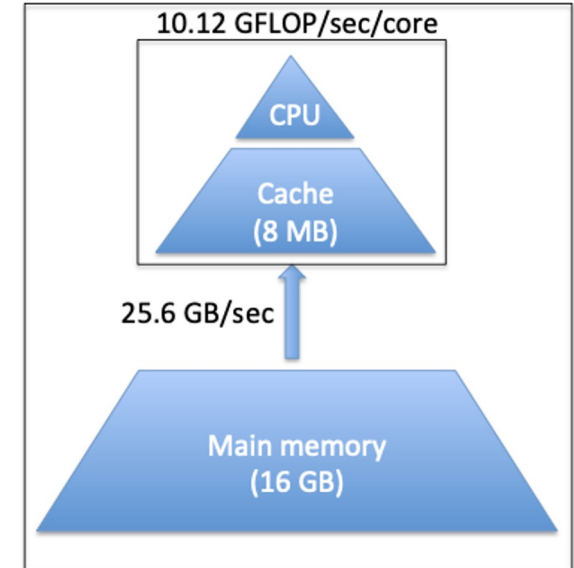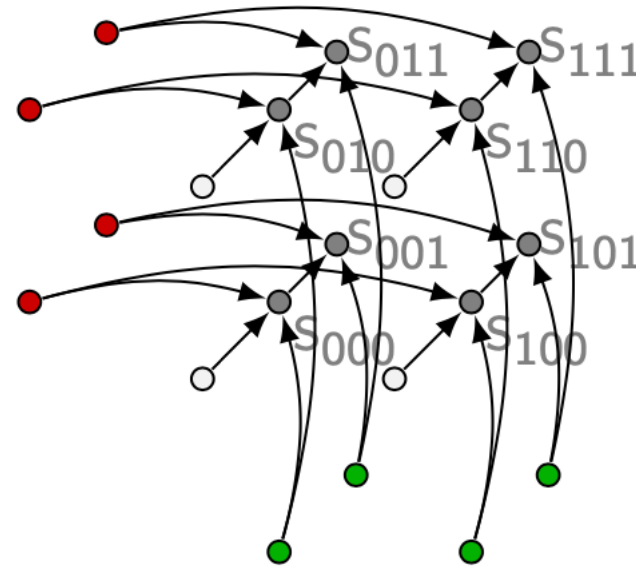
```
for (i = 0; i < m; i++)
  for (j = 0; j < n; j++)
    for (k = 0; k < p; k++)
      C[i][j] += A[i][k] * B[k][j];
```

C[0][0] += A[0][0] * B[0][0];  // S000
C[0][0] += A[0][1] * B[1][0];  // S001
C[0][0] += A[0][2] * B[2][0];  // S002
C[0][0] += A[0][3] * B[3][0];  // S003
...
C[0][1] += A[0][0] * B[0][1];  // S010
C[0][1] += A[0][1] * B[1][1];  // S011
C[0][1] += A[0][2] * B[2][1];  // S012
C[0][1] += A[0][3] * B[3][1];  // S013

Intel Xeon Processor E5520 (Nehalem)



gemm

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 )

**input** is an affine code

```
for (i = 0; i < _PB_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_FUN(A[i][i]);
}
```

**cholesky code**

Cholesky

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 )

**input** is an affine code    with parameter **N**
and parameter **S** for
cache size

```
for (i = 0; i < _PB_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_FUN(A[i][i]);
}
```

**cholesky code**

Unit is for **N** for, **S** for and IO lower bound is "number".

Cholesky

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

**input** is an affine code   with parameter **N**
and parameter **S** for
cache size

IO lower bound
with parameter **N** and parameter **S** for cache size

asymptotic
IO lower bound

```
for (i = 0; i < _PB_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_FUN(A[i][i]);
}
```

**cholesky code**

Unit is for **N** for, **S** for and IO lower bound is "number".

Cholesky

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

**input** is an affine code

```
for (i = 0; i < _PB_N; i++) {
    for (j = 0; j < i; j++) {
        for (k = 0; k < 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_FUN(A[i][i]);
}
```

cholesky code

**Output #2** is a proof

Cholesky

We consider a segmentation of the execution every T IOs. For each segment, we have at most K = S + T "input" data to perform the statements in the segment.

We consider a segmentation of the execution every T IOs. For each segment, we have at most K = S + T input data to perform the statements in the segment.

the maximum number of statements in a segment

We consider a segmentation of the execution every T IOs. For each segment, we have at most K = S + T input data to perform the statements in the segment.

$$2 \cdot (K/3)^{3/2}$$

the maximum number of statements in a segment

Cholesky

the total number of statements in the execution

the IO of a segment

We consider a segmentation of the execution every T IOs. For each segment, we have at most K = S + T input data to perform the statements in the segment.

$$Q \geq T \times \left\lfloor \frac{N^3/6}{2 \cdot (K/3)^{3/2}} \right\rfloor$$

The minimum number of segments in a schedule

the maximum number of statements in a segment

is greater than or equal to

the IO of any schedule

Cholesky

We consider a segmentation of the execution every T IOs. For each segment, we have at most K = S + T input data to perform the statements in the segment.

$$2 \cdot (K/3)^{3/2}$$

↑
the maximum number of statements in a segment

Cholesky

# What do we want to compute?

We want to compute the $n^3$

$$c_{ijk} = a_{ik} b_{kj}.$$

The computation of $c_{ijk}$ requires $a_{ik}$ and $b_{kj}$ to be in cache.

## Lemma (Loomis-Whitney Inequality)

Let $V \in Z[3]$ be a finite set, and let $V_x$, $V_y$, and $V_z$ be orthogonal projections of $V$ onto the coordinate planes. The cardinality of $V$, $|V|$, satisfies

$$|V| \leq \sqrt{|V_x| \cdot |V_y| \cdot |V_z|}.$$

## Lemma (Loomis-Whitney Inequality)

Let $V \in Z[3]$ be a finite set, and let $V_x$, $V_y$, and $V_z$ be orthogonal projections of V onto the coordinate planes. The cardinality of V, $|V|$, satisfies

$$|V| \leq \sqrt{|V_x| \cdot |V_y| \cdot |V_z|}.$$

We consider a segmentation of the execution every T IOs. For each segment, we have at most K = S + T input data to perform the statements in the segment.

We will use a generalization: the Brascamp-Lieb inequality.

```
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];
```

C[0][0] += A[0][0] * B[0][0];  // S000
C[0][0] += A[0][1] * B[1][0];  // S001
C[0][0] += A[0][2] * B[2][0];  // S002
C[0][0] += A[0][3] * B[3][0];  // S003
...
C[0][1] += A[0][0] * B[0][1];  // S010
C[0][1] += A[0][1] * B[1][1];  // S011
C[0][1] += A[0][2] * B[2][1];  // S012
C[0][1] += A[0][3] * B[3][1];  // S013



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];
```

C[0][0] += A[0][0] * B[0][0];  // S000
C[0][0] += A[0][1] * B[1][0];  // S001
C[0][0] += A[0][2] * B[2][0];  // S002
C[0][0] += A[0][3] * B[3][0];  // S003
…
C[0][1] += A[0][0] * B[0][1];  // S010
C[0][1] += A[0][1] * B[1][1];  // S011
C[0][1] += A[0][2] * B[2][1];  // S012
C[0][1] += A[0][3] * B[3][1];  // S013

gemm

```
for(k = 0; k < n; k++) {
        A[k][k] = sqrt(A[k][k]);                //S1
    for(i = k+1; i < n; i++)
        A[i][k] /= A[k][k];                     // S2
    for(i = k+1; i < n; i++)
        for(j = k+1; j <= i; j++)
            A[i][j] -= A[i][k] * A[j][k];        // S3
}
```

Cholesky

```
for(k = 0;  k < n;  k++) {
        A[k][k] = sqrt(A[k][k]);                    //S1
    for(i = k+1;  i < n;  i++)
      A[i][k] /= A[k][k];                           // S2
    for(i = k+1;  i < n;  i++)
        for(j = k+1;  j <= i;  j++)
            A[i][j] -= A[i][k] * A[j][k];            // S3
}
```



Cholesky

```
for(k = 0; k < n; k++) {
        A[k][k] = sqrt(A[k][k]);                    //S1
    for(i = k+1; i < n; i++)
        A[i][k] /= A[k][k];                         // S2
    for(i = k+1; i < n; i++)
        for(j = k+1; j <= i; j++)
            A[i][j] -= A[i][k] * A[j][k];           // S3
}
```



$$R_{e_1} = \{S_3[k-1,i,j] \rightarrow S_3[k,i,j]:\ 1 \leq k < N\ \wedge\ k+1 \leq i < N\ \wedge\ k+1 \leq j \leq i\}$$
$$R_{e_2} = \{S_2[k,j] \rightarrow S_3[k,i,j]:\ 0 \leq k < N\ \wedge\ k+1 \leq i < N\ \wedge\ k+1 \leq j \leq i\}$$
$$R_{e_3} = \{S_2[k,i] \rightarrow S_3[k,i,j]:\ 0 \leq k < N\ \wedge\ k+1 \leq i < N\ \wedge\ k+1 \leq j \leq i\}$$
$$R_{e_4} = \{S_3[k-1,i,k] \rightarrow S_2[k,i]:\ 1 \leq k < N\ \wedge\ k+1 \leq i < N\}$$
$$R_{e_5} = \{S_1[k] \rightarrow S_2[k,i]:\ 0 \leq k < N\ \wedge\ k+1 \leq i < N\}$$
$$R_{e_6} = \{S_1[k-1,k,k] \rightarrow S_1[k]:\ 1 \leq k < N\ \wedge\ k+1 \leq i < N\}$$

Cholesky

```
for(k = 0; k < n; k++) {
        A[k][k] = sqrt(A[k][k]);                    //S1
    for(i = k+1; i < n; i++)
        A[i][k] /= A[k][k];                         // S2
    for(i = k+1; i < n; i++)
        for(j = k+1; j <= i; j++)
            A[i][j] -= A[i][k] * A[j][k];           // S3
}
```

$$\operatorname{Dom}(P_1) = \{S_3[k,i,j]: \; 1 \le k < N \; \wedge \; k+1 \le i < N \; \wedge \; k+1 \le j \le i\}$$
$$\operatorname{Dom}(P_2) = \{S_3[k,i,j]: \; 0 \le k < N \; \wedge \; k+1 \le i < N \; \wedge \; k+1 \le j \le i\}$$
$$\operatorname{Dom}(P_3) = \{S_3[k,i,j]: \; 0 \le k < N \; \wedge \; k+1 \le i < N \; \wedge \; k+1 \le j \le i\}$$
$$D_S := \operatorname{Dom}(P_1) \cap \operatorname{Dom}(P_2) \cap \operatorname{Dom}(P_3)$$
$$= \{S_3[k,i,j]: \; 1 \le k < N \; \wedge \; k+1 \le i < N \; \wedge \; k+1 \le j \le i\}$$



$$R_{e_1} = \{S_3[k-1,i,j] \to S_3[k,i,j]: \; 1 \le k < N \; \wedge \; k+1 \le i < N \; \wedge \; k+1 \le j \le i\}$$
$$R_{e_2} = \{S_2[k,j] \to S_3[k,i,j]: \; 0 \le k < N \; \wedge \; k+1 \le i < N \; \wedge \; k+1 \le j \le i\}$$
$$R_{e_3} = \{S_2[k,i] \to S_3[k,i,j]: \; 0 \le k < N \; \wedge \; k+1 \le i < N \; \wedge \; k+1 \le j \le i\}$$
$$R_{e_4} = \{S_3[k-1,i,k] \to S_2[k,i]: \; 1 \le k < N \; \wedge \; k+1 \le i < N\}$$
$$R_{e_5} = \{S_1[k] \to S_2[k,i]: \; 0 \le k < N \; \wedge \; k+1 \le i < N\}$$
$$R_{e_6} = \{S_1[k-1,k,k] \to S_1[k]: \; 1 \le k < N \; \wedge \; k+1 \le i < N\}$$

$$P_1 = (e_1) \qquad\qquad \text{is a chain path}$$
$$P_2 = (e_2) \qquad\qquad \text{is a broadcast path}$$
$$P_3 = (e_3) \qquad\qquad \text{is a broadcast path}$$

Cholesky

```
for(k = 0; k < n; k++) {
        A[k][k] = sqrt(A[k][k]);              //S1
    for(i = k+1; i < n; i++)
        A[i][k] /= A[k][k];                   // S2
    for(i = k+1; i < n; i++)
        for(j = k+1; j <= i; j++)
            A[i][j] -= A[i][k] * A[j][k];     // S3
}
```

$\text{Dom}(P_1) = \{S_3[k,i,j]: \; 1 \leq k < N \; \wedge \; k+1 \leq i < N \; \wedge \; k+1 \leq j \leq i\}$

$\text{Dom}(P_2) = \{S_3[k,i,j]: \; 0 \leq k < N \; \wedge \; k+1 \leq i < N \; \wedge \; k+1 \leq j \leq i\}$

$\text{Dom}(P_3) = \{S_3[k,i,j]: \; 0 \leq k < N \; \wedge \; k+1 \leq i < N \; \wedge \; k+1 \leq j \leq i\}$

$D_S := \text{Dom}(P_1) \cap \text{Dom}(P_2) \cap \text{Dom}(P_3)$

$\quad = \{S_3[k,i,j]: \; 1 \leq k < N \; \wedge \; k+1 \leq i < N \; \wedge \; k+1 \leq j \leq i\}$



$\phi_1(k,i,j) = \text{proj}_{(1,0,0)}(k,i,j) = (0,i,j)$ $\qquad$ kernel $k_1 = \text{Ker}(\phi_1) = \langle(1,0,0)\rangle$

$\phi_2(k,i,j) = (k,j)$ $\qquad$ kernel $k_2 = \text{Ker}(\phi_2) = \langle(0,1,0)\rangle$

$\phi_3(k,i,j) = (k,i)$ $\qquad$ kernel $k_3 = \text{Ker}(\phi_3) = \langle(0,0,1)\rangle$



$R_{e_1} = \{S_3[k-1,i,j] \rightarrow S_3[k,i,j]: \; 1 \leq k < N \; \wedge \; k+1 \leq i < N \; \wedge \; k+1 \leq j \leq i\}$

$R_{e_2} = \{S_2[k,j] \rightarrow S_3[k,i,j]: \; 0 \leq k < N \; \wedge \; k+1 \leq i < N \; \wedge \; k+1 \leq j \leq i\}$

$R_{e_3} = \{S_2[k,i] \rightarrow S_3[k,i,j]: \; 0 \leq k < N \; \wedge \; k+1 \leq i < N \; \wedge \; k+1 \leq j \leq i\}$

$R_{e_4} = \{S_3[k-1,i,k] \rightarrow S_2[k,i]: \; 1 \leq k < N \; \wedge \; k+1 \leq i < N\}$

$R_{e_5} = \{S_1[k] \rightarrow S_2[k,i]: \; 0 \leq k < N \; \wedge \; k+1 \leq i < N\}$

$R_{e_6} = \{S_1[k-1,k,k] \rightarrow S_1[k]: \; 1 \leq k < N \; \wedge \; k+1 \leq i < N\}$

$P_1 = (e_1)$ $\qquad\qquad$ is a chain path

$P_2 = (e_2)$ $\qquad\qquad$ is a broadcast path

$P_3 = (e_3)$ $\qquad\qquad$ is a broadcast path

Cholesky

```
for(k = 0; k < n; k++) {
        A[k][k] = sqrt(A[k][k]);              //S1
    for(i = k+1; i < n; i++)
        A[i][k] /= A[k][k];                   // S2
    for(i = k+1; i < n; i++)
        for(j = k+1; j <= i; j++)
            A[i][j] -= A[i][k] * A[j][k];     // S3
}
```
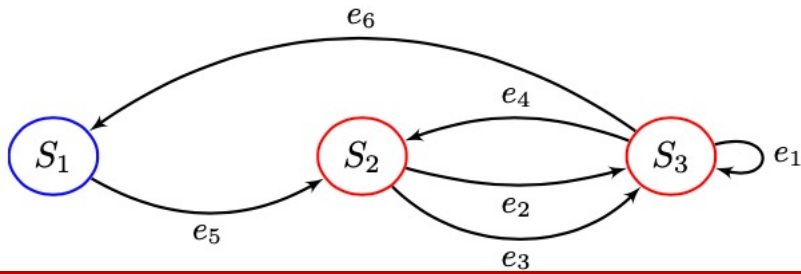
$$\text{Dom}(P_1) = \{S_3[k,i,j]: \; 1 \le k < N \; \wedge \; k+1 \le i < N \; \wedge \; k+1 \le j \le i\}$$

$$\text{Dom}(P_2) = \{S_3[k,i,j]: \; 0 \le k < N \; \wedge \; k+1 \le i < N \; \wedge \; k+1 \le j \le i\}$$

$$\text{Dom}(P_3) = \{S_3[k,i,j]: \; 0 \le k < N \; \wedge \; k+1 \le i < N \; \wedge \; k+1 \le j \le i\}$$

$$D_S := \text{Dom}(P_1) \cap \text{Dom}(P_2) \cap \text{Dom}(P_3)$$

$$= \{S_3[k,i,j]: \; 1 \le k < N \; \wedge \; k+1 \le i < N \; \wedge \; k+1 \le j \le i\}$$

$$\phi_1(k,i,j) = \text{proj}_{(1,0,0)}(k,i,j) = (0,i,j) \qquad \text{kernel } k_1 = \text{Ker}(\phi_1) = \langle(1,0,0)\rangle$$

$$\phi_2(k,i,j) = (k,j) \qquad\qquad\qquad\qquad \text{kernel } k_2 = \text{Ker}(\phi_2) = \langle(0,1,0)\rangle$$

$$\phi_3(k,i,j) = (k,i) \qquad\qquad\qquad\qquad \text{kernel } k_3 = \text{Ker}(\phi_3) = \langle(0,0,1)\rangle$$



$$R_{e_1} = \{S_3[k-1,i,j] \to S_3[k,i,j]: \; 1 \le k < N \; \wedge \; k+1 \le i < N \; \wedge \; k+1 \le j \le i\}$$

$$R_{e_2} = \{S_2[k,j] \to S_3[k,i,j]: \; 0 \le k < N \; \wedge \; k+1 \le i < N \; \wedge \; k+1 \le j \le i\}$$

$$R_{e_3} = \{S_2[k,i] \to S_3[k,i,j]: \; 0 \le k < N \; \wedge \; k+1 \le i < N \; \wedge \; k+1 \le j \le i\}$$

$$R_{e_4} = \{S_3[k-1,i,k] \to S_2[k,i]: \; 1 \le k < N \; \wedge \; k+1 \le i < N\}$$

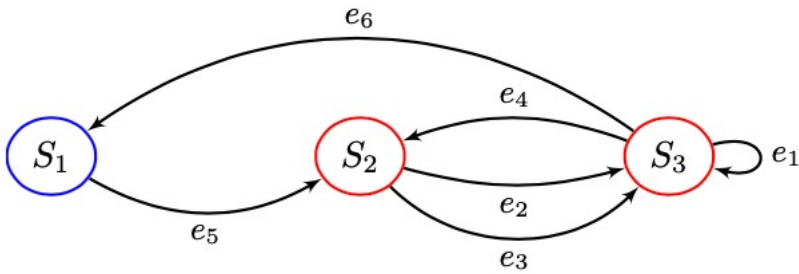$$R_{e_5} = \{S_1[k] \to S_2[k,i]: \; 0 \le k < N \; \wedge \; k+1 \le i < N\}$$

$$R_{e_6} = \{S_1[k-1,k,k] \to S_1[k]: \; 1 \le k < N \; \wedge \; k+1 \le i < N\}$$

THEOREM 3.10 (BRASCAMP-LIEB INEQUALITY, DISCRETE CASE [CHRIST ET AL. 2013]). *Let d and $d_j$ be nonnegative integers and $\phi_j : \mathbb{Z}^d \mapsto \mathbb{Z}^{d_j}$ be group homomorphisms for $1 \le j \le m$. Let $0 \le s_1, s_2, \ldots, s_m \le 1$. Suppose that:*

$$\text{rank}(H) \le \sum_{j=1}^{m} s_j \cdot \text{rank}(\phi_j(H)) \text{ for all subgroups } H \text{ of } \mathbb{Z}^d \qquad (2)$$

*Then:*

$$|E| \le \prod_{j=1}^{m} |\phi_j(E)|^{s_j} \text{ for all nonempty finite sets } E \subseteq \mathbb{Z}^d. \qquad (3)$$

$$P_1 = (e_1) \qquad\qquad \text{is a chain path}$$

$$P_2 = (e_2) \qquad\qquad \text{is a broadcast path}$$

$$P_3 = (e_3) \qquad\qquad \text{is a broadcast path}$$

Cholesky

```
for(k = 0; k < n; k++) {
        A[k][k] = sqrt(A[k][k]);              // S1
    for(i = k+1; i < n; i++)
        A[i][k] /= A[k][k];                   // S2
    for(i = k+1; i < n; i++)
        for(j = k+1; j <= i; j++)
            A[i][j] -= A[i][k] * A[j][k];     // S3
}
```
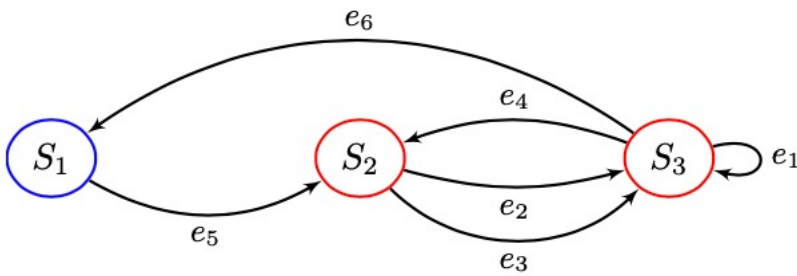
$\text{Dom}(P_1) = \{S_3[k,i,j] : \ 1 \le k < N \ \land \ k+1 \le i < N \ \land \ k+1 \le j \le i\}$

$\text{Dom}(P_2) = \{S_3[k,i,j] : \ 0 \le k < N \ \land \ k+1 \le i < N \ \land \ k+1 \le j \le i\}$

$\text{Dom}(P_3) = \{S_3[k,i,j] : \ 0 \le k < N \ \land \ k+1 \le i < N \ \land \ k+1 \le j \le i\}$

$D_S := \text{Dom}(P_1) \cap \text{Dom}(P_2) \cap \text{Dom}(P_3)$

$\quad\quad = \{S_3[k,i,j] : \ 1 \le k < N \ \land \ k+1 \le i < N \ \land \ k+1 \le j \le i\}$



$\phi_1(k,i,j) = \text{proj}_{(1,0,0)}(k,i,j) = (0,i,j)$ \quad\quad $\text{kernel } k_1 = \text{Ker}(\phi_1) = \langle(1,0,0)\rangle$

$\phi_2(k,i,j) = (k,j)$ \quad\quad\quad\quad\quad\quad\quad\quad\quad $\text{kernel } k_2 = \text{Ker}(\phi_2) = \langle(0,1,0)\rangle$

$\phi_3(k,i,j) = (k,i)$ \quad\quad\quad\quad\quad\quad\quad\quad\quad $\text{kernel } k_3 = \text{Ker}(\phi_3) = \langle(0,0,1)\rangle$

THEOREM 3.10 (BRASCAMP-LIEB INEQUALITY, DISCRETE CASE [CHRIST ET AL. 2013]). *Let $d$ and $d_j$ be nonnegative integers and $\phi_j : \mathbb{Z}^d \mapsto \mathbb{Z}^{d_j}$ be group homomorphisms for $1 \le j \le m$. Let $0 \le s_1, s_2, \ldots, s_m \le 1$. Suppose that:*

$$\text{rank}(H) \le \sum_{j=1}^{m} s_j \cdot \text{rank}(\phi_j(H)) \text{ for all subgroups } H \text{ of } \mathbb{Z}^d \quad\quad (2)$$

*Then:*

$$|E| \le \prod_{j=1}^{m} |\phi_j(E)|^{s_j} \text{ for all nonempty finite sets } E \subseteq \mathbb{Z}^d. \quad\quad (3)$$

$R_{e_1} = \{S_3[k-1,i,j] \to S_3[k,i,j] : \ 1 \le k < N \ \land \ k+1 \le i < N \ \land \ k+1 \le j \le i\}$

$R_{e_2} = \{S_2[k,j] \to S_3[k,i,j] : \ 0 \le k < N \ \land \ k+1 \le i < N \ \land \ k+1 \le j \le i\}$

$R_{e_3} = \{S_2[k,i] \to S_3[k,i,j] : \ 0 \le k < N \ \land \ k+1 \le i < N \ \land \ k+1 \le j \le i\}$

$R_{e_4} = \{S_3[k-1,i,k] \to S_2[k,i] : \ 1 \le k < N \ \land \ k+1 \le i < N\}$

$R_{e_5} = \{S_1[k] \to S_2[k,i] : \ 0 \le k < N \ \land \ k+1 \le i < N\}$

$R_{e_6} = \{S_1[k-1,k,k] \to S_1[k] : \ 1 \le k < N \ \land \ k+1 \le i < N\}$

$0 \le s_1, s_2, s_3 \le 1$

$1 \le s_2 + s_3$

$1 \le s_1 + s_3$

$1 \le s_1 + s_2$

$P_1 = (e_1)$ \quad\quad\quad\quad\quad\quad is a chain path

$P_2 = (e_2)$ \quad\quad\quad\quad\quad\quad is a broadcast path

$P_3 = (e_3)$ \quad\quad\quad\quad\quad\quad is a broadcast path

Cholesky

```
for(k = 0; k < n; k++) {
        A[k][k] = sqrt(A[k][k]);          // S1
    for(i = k+1; i < n; i++)
        A[i][k] /= A[k][k];               // S2
    for(i = k+1; i < n; i++)
        for(j = k+1; j <= i; j++)
            A[i][j] -= A[i][k] * A[j][k];  // S3
}
```

$$\text{Dom}(P_1) = \{S_3[k,i,j] : \boxed{1} \le k < N \ \wedge \ k+1 \le i < N \ \wedge \ k+1 \le j \le i\}$$
$$\text{Dom}(P_2) = \{S_3[k,i,j] : 0 \le k < N \ \wedge \ k+1 \le i < N \ \wedge \ k+1 \le j \le i\}$$
$$\text{Dom}(P_3) = \{S_3[k,i,j] : 0 \le k < N \ \wedge \ k+1 \le i < N \ \wedge \ k+1 \le j \le i\}$$
$$D_S := \text{Dom}(P_1) \cap \text{Dom}(P_2) \cap \text{Dom}(P_3)$$
$$= \{S_3[k,i,j] : 1 \le k < N \ \wedge \ k+1 \le i < N \ \wedge \ k+1 \le j \le i\}$$

$$\phi_1(k,i,j) = \text{proj}_{(1,0,0)}(k,i,j) = (0,i,j) \qquad \text{kernel } k_1 = \text{Ker}(\phi_1) = \langle(1,0,0)\rangle$$
$$\phi_2(k,i,j) = (k,j) \qquad\qquad\qquad\qquad\qquad \text{kernel } k_2 = \text{Ker}(\phi_2) = \langle(0,1,0)\rangle$$
$$\phi_3(k,i,j) = (k,i) \qquad\qquad\qquad\qquad\qquad \text{kernel } k_3 = \text{Ker}(\phi_3) = \langle(0,0,1)\rangle$$

THEOREM 3.10 (BRASCAMP-LIEB INEQUALITY, DISCRETE CASE [CHRIST ET AL. 2013]). *Let d and* $d_j$ *be nonnegative integers and* $\phi_j : \mathbb{Z}^d \mapsto \mathbb{Z}^{d_j}$ *be group homomorphisms for* $1 \le j \le m$*. Let* $0 \le s_1, s_2, \ldots, s_m \le 1$*. Suppose that:*

$$\text{rank}(H) \le \sum_{j=1}^{m} s_j \cdot \text{rank}(\phi_j(H)) \quad \text{for all subgroups } H \text{ of } \mathbb{Z}^d \tag{2}$$
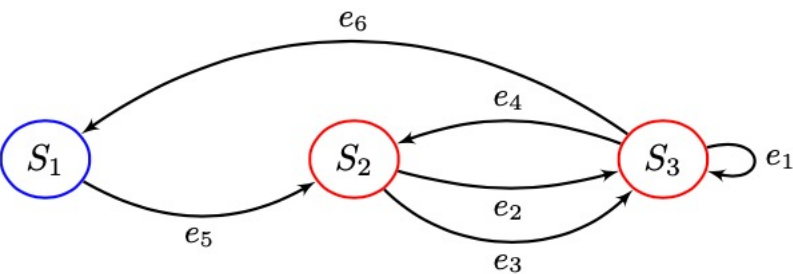
*Then:*

$$|E| \le \prod_{j=1}^{m} |\phi_j(E)|^{s_j} \quad \text{for all nonempty finite sets } E \subseteq \mathbb{Z}^d. \tag{3}$$

$$R_{e_1} = \{S_3[k-1,i,j] \to S_3[k,i,j] : 1 \le k < N \ \wedge \ k+1 \le i < N \ \wedge \ k+1 \le j \le i\}$$
$$R_{e_2} = \{S_2[k,j] \to S_3[k,i,j] : 0 \le k < N \ \wedge \ k+1 \le i < N \ \wedge \ k+1 \le j \le i\}$$
$$R_{e_3} = \{S_2[k,i] \to S_3[k,i,j] : 0 \le k < N \ \wedge \ k+1 \le i < N \ \wedge \ k+1 \le j \le i\}$$
$$R_{e_4} = \{S_3[k-1,i,k] \to S_2[k,i] : 1 \le k < N \ \wedge \ k+1 \le i < N\}$$
$$R_{e_5} = \{S_1[k] \to S_2[k,i] : 0 \le k < N \ \wedge \ k+1 \le i < N\}$$
$$R_{e_6} = \{S_1[k-1,k,k] \to S_1[k] : 1 \le k < N \ \wedge \ k+1 \le i < N\}$$

$$0 \le s_1, s_2, s_3 \le 1 \qquad \boxed{s_1 = s_2 = s_3 = \tfrac{1}{2}}$$
$$1 \le s_2 + s_3$$
$$1 \le s_1 + s_3$$
$$1 \le s_1 + s_2$$

$$P_1 = (e_1) \qquad\qquad \text{is a chain path}$$
$$P_2 = (e_2) \qquad\qquad \text{is a broadcast path}$$
$$P_3 = (e_3) \qquad\qquad \text{is a broadcast path}$$

Cholesky

```
for(k = 0; k < n; k++) {
        A[k][k] = sqrt(A[k][k]);               //S1
    for(i = k+1; i < n; i++)
        A[i][k] /= A[k][k];                    // S2
    for(i = k+1; i < n; i++)
        for(j = k+1; j <= i; j++)
            A[i][j] -= A[i][k] * A[j][k];      // S3
}
```

$$\text{Dom}(P_1) = \{S_3[k,i,j] : \boxed{1} \le k < N \ \wedge \ k+1 \le i < N \ \wedge \ k+1 \le j \le i\}$$

$$\text{Dom}(P_2) = \{S_3[k,i,j] : 0 \le k < N \ \wedge \ k+1 \le i < N \ \wedge \ k+1 \le j \le i\}$$

$$\text{Dom}(P_3) = \{S_3[k,i,j] : 0 \le k < N \ \wedge \ k+1 \le i < N \ \wedge \ k+1 \le j \le i\}$$

$$D_S := \text{Dom}(P_1) \cap \text{Dom}(P_2) \cap \text{Dom}(P_3)$$
$$= \{S_3[k,i,j] : 1 \le k < N \ \wedge \ k+1 \le i < N \ \wedge \ k+1 \le j \le i\}$$



$$\phi_1(k,i,j) = \text{proj}_{(1,0,0)}(k,i,j) = (0,i,j) \qquad \text{kernel } k_1 = \text{Ker}(\phi_1) = \langle(1,0,0)\rangle$$
$$\phi_2(k,i,j) = (k,j) \qquad\qquad\qquad\qquad \text{kernel } k_2 = \text{Ker}(\phi_2) = \langle(0,1,0)\rangle$$
$$\phi_3(k,i,j) = (k,i) \qquad\qquad\qquad\qquad \text{kernel } k_3 = \text{Ker}(\phi_3) = \langle(0,0,1)\rangle$$

THEOREM 3.10 (BRASCAMP-LIEB INEQUALITY, DISCRETE CASE [CHRIST ET AL. 2013]). *Let d and $d_j$ be nonnegative integers and $\phi_j : \mathbb{Z}^d \mapsto \mathbb{Z}^{d_j}$ be group homomorphisms for $1 \le j \le m$. Let $0 \le s_1, s_2, \ldots, s_m \le 1$. Suppose that:*

$$\text{rank}(H) \le \sum_{j=1}^{m} s_j \cdot \text{rank}(\phi_j(H)) \ \text{for all subgroups H of } \mathbb{Z}^d \qquad (2)$$

*Then:*

$$|E| \le \prod_{j=1}^{m} |\phi_j(E)|^{s_j} \ \text{for all nonempty finite sets } E \subseteq \mathbb{Z}^d. \qquad (3)$$
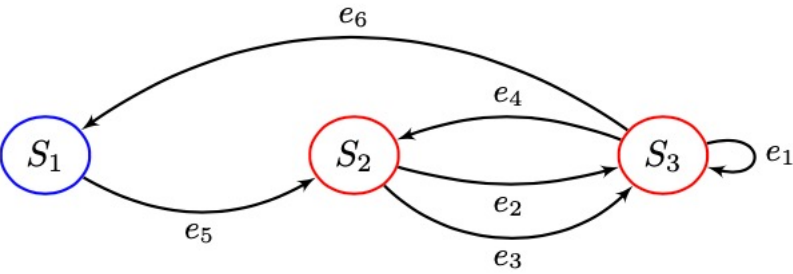
$$R_{e_1} = \{S_3[k-1,i,j] \to S_3[k,i,j] : 1 \le k < N \ \wedge \ k+1 \le i < N \ \wedge \ k+1 \le j \le i\}$$
$$R_{e_2} = \{S_2[k,j] \to S_3[k,i,j] : 0 \le k < N \ \wedge \ k+1 \le i < N \ \wedge \ k+1 \le j \le i\}$$
$$R_{e_3} = \{S_2[k,i] \to S_3[k,i,j] : 0 \le k < N \ \wedge \ k+1 \le i < N \ \wedge \ k+1 \le j \le i\}$$
$$R_{e_4} = \{S_3[k-1,i,k] \to S_2[k,i] : 1 \le k < N \ \wedge \ k+1 \le i < N\}$$
$$R_{e_5} = \{S_1[k] \to S_2[k,i] : 0 \le k < N \ \wedge \ k+1 \le i < N\}$$
$$R_{e_6} = \{S_1[k-1,k,k] \to S_1[k] : 1 \le k < N \ \wedge \ k+1 \le i < N\}$$

$$0 \le s_1, s_2, s_3 \le 1 \qquad s_1 = s_2 = s_3 = \frac{1}{2}$$
$$1 \le s_2 + s_3$$
$$1 \le s_1 + s_3$$
$$1 \le s_1 + s_2$$

$$\boxed{|E| \le |\phi_1(E)|^{\frac{1}{2}} \cdot |\phi_2(E)|^{\frac{1}{2}} \cdot |\phi_3(E)|^{\frac{1}{2}}}$$

$P_1 = (e_1)$      is a chain path

$P_2 = (e_2)$      is a broadcast path

$P_3 = (e_3)$      is a broadcast path
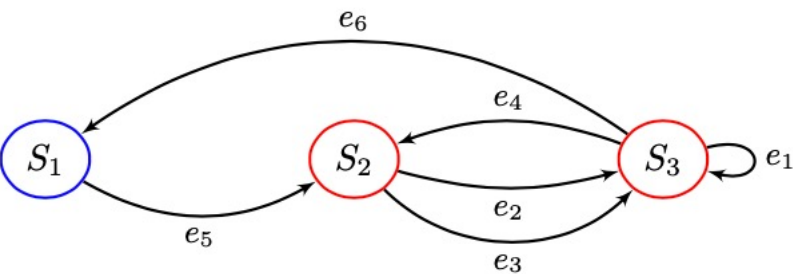
Cholesky

```
for(k = 0; k < n; k++) {
        A[k][k] = sqrt(A[k][k]);              //S1
    for(i = k+1; i < n; i++)
        A[i][k] /= A[k][k];                   // S2
    for(i = k+1; i < n; i++)
        for(j = k+1; j <= i; j++)
            A[i][j] -= A[i][k] * A[j][k];     // S3
}
```

$$\text{Dom}(P_1) = \{S_3[k,i,j] : \; 1 \le k < N \; \wedge \; k+1 \le i < N \; \wedge \; k+1 \le j \le i\}$$
$$\text{Dom}(P_2) = \{S_3[k,i,j] : \; 0 \le k < N \; \wedge \; k+1 \le i < N \; \wedge \; k+1 \le j \le i\}$$
$$\text{Dom}(P_3) = \{S_3[k,i,j] : \; 0 \le k < N \; \wedge \; k+1 \le i < N \; \wedge \; k+1 \le j \le i\}$$

$$D_S := \text{Dom}(P_1) \cap \text{Dom}(P_2) \cap \text{Dom}(P_3)$$
$$= \{S_3[k,i,j] : \; 1 \le k < N \; \wedge \; k+1 \le i < N \; \wedge \; k+1 \le j \le i\}$$

$$\phi_1(k,i,j) = \text{proj}_{(1,0,0)}(k,i,j) = (0,i,j) \qquad \text{kernel } k_1 = \text{Ker}(\phi_1) = \langle (1,0,0) \rangle$$
$$\phi_2(k,i,j) = (k,j) \qquad\qquad\qquad\qquad \text{kernel } k_2 = \text{Ker}(\phi_2) = \langle (0,1,0) \rangle$$
$$\phi_3(k,i,j) = (k,i) \qquad\qquad\qquad\qquad \text{kernel } k_3 = \text{Ker}(\phi_3) = \langle (0,0,1) \rangle$$

THEOREM 3.10 (BRASCAMP-LIEB INEQUALITY, DISCRETE CASE [CHRIST ET AL. 2013]). *Let $d$ and $d_j$ be nonnegative integers and $\phi_j : \mathbb{Z}^d \mapsto \mathbb{Z}^{d_j}$ be group homomorphisms for $1 \le j \le m$. Let $0 \le s_1, s_2, \ldots, s_m \le 1$. Suppose that:*

$$\text{rank}(H) \le \sum_{j=1}^{m} s_j \cdot \text{rank}(\phi_j(H)) \text{ for all subgroups } H \text{ of } \mathbb{Z}^d \qquad (2)$$

*Then:*

$$|E| \le \prod_{j=1}^{m} |\phi_j(E)|^{s_j} \text{ for all nonempty finite sets } E \subseteq \mathbb{Z}^d. \qquad (3)$$



$$R_{e_1} = \{S_3[k-1,i,j] \to S_3[k,i,j] : \; 1 \le k < N \; \wedge \; k+1 \le i < N \; \wedge \; k+1 \le j \le i\}$$
$$R_{e_2} = \{S_2[k,j] \to S_3[k,i,j] : \; 0 \le k < N \; \wedge \; k+1 \le i < N \; \wedge \; k+1 \le j \le i\}$$
$$R_{e_3} = \{S_2[k,i] \to S_3[k,i,j] : \; 0 \le k < N \; \wedge \; k+1 \le i < N \; \wedge \; k+1 \le j \le i\}$$
$$R_{e_4} = \{S_3[k-1,i,k] \to S_2[k,i] : \; 1 \le k < N \; \wedge \; k+1 \le i < N\}$$
$$R_{e_5} = \{S_1[k] \to S_2[k,i] : \; 0 \le k < N \; \wedge \; k+1 \le i < N\}$$
$$R_{e_6} = \{S_1[k-1,k,k] \to S_1[k] : \; 1 \le k < N \; \wedge \; k+1 \le i < N\}$$

$$0 \le s_1, s_2, s_3 \le 1 \qquad s_1 = s_2 = s_3 = \tfrac{1}{2}$$
$$1 \le s_2 + s_3$$
$$1 \le s_1 + s_3$$
$$1 \le s_1 + s_2$$

$P_1 = (e_1)$ \qquad is a chain path
$P_2 = (e_2)$ \qquad is a broadcast path
$P_3 = (e_3)$ \qquad is a broadcast path

$$|E| \le |\phi_1(E)|^{\frac{1}{2}} \cdot |\phi_2(E)|^{\frac{1}{2}} \cdot |\phi_3(E)|^{\frac{1}{2}}$$
$$|\phi_i(E)| \le K, \; i = 1, 2, 3$$
$$|E| \le K^{3/2}$$

Cholesky

```
for(k = 0; k < n; k++) {
        A[k][k] = sqrt(A[k][k]);              //S1
    for(i = k+1; i < n; i++)
        A[i][k] /= A[k][k];                   // S2
    for(i = k+1; i < n; i++)
        for(j = k+1; j <= i; j++)
            A[i][j] -= A[i][k] * A[j][k];     // S3
}
```

$$\text{Dom}(P_1) = \{S_3[k,i,j] : \ 1 \le k < N \ \wedge \ k+1 \le i < N \ \wedge \ k+1 \le j \le i\}$$
$$\text{Dom}(P_2) = \{S_3[k,i,j] : \ 0 \le k < N \ \wedge \ k+1 \le i < N \ \wedge \ k+1 \le j \le i\}$$
$$\text{Dom}(P_3) = \{S_3[k,i,j] : \ 0 \le k < N \ \wedge \ k+1 \le i < N \ \wedge \ k+1 \le j \le i\}$$
$$D_S := \text{Dom}(P_1) \cap \text{Dom}(P_2) \cap \text{Dom}(P_3)$$
$$= \{S_3[k,i,j] : \ 1 \le k < N \ \wedge \ k+1 \le i < N \ \wedge \ k+1 \le j \le i\}$$



$$\phi_1(k,i,j) = \text{proj}_{(1,0,0)}(k,i,j) = (0,i,j) \qquad \text{kernel } k_1 = \text{Ker}(\phi_1) = \langle(1,0,0)\rangle$$
$$\phi_2(k,i,j) = (k,j) \qquad \text{kernel } k_2 = \text{Ker}(\phi_2) = \langle(0,1,0)\rangle$$
$$\phi_3(k,i,j) = (k,i) \qquad \text{kernel } k_3 = \text{Ker}(\phi_3) = \langle(0,0,1)\rangle$$

THEOREM 3.10 (BRASCAMP-LIEB INEQUALITY, DISCRETE CASE [CHRIST ET AL. 2013]). *Let d and $d_j$ be nonnegative integers and $\phi_j : \mathbb{Z}^d \mapsto \mathbb{Z}^{d_j}$ be group homomorphisms for $1 \le j \le m$. Let $0 \le s_1, s_2, \ldots, s_m \le 1$. Suppose that:*

$$\text{rank}(H) \le \sum_{j=1}^{m} s_j \cdot \text{rank}(\phi_j(H)) \text{ for all subgroups } H \text{ of } \mathbb{Z}^d \qquad (2)$$

*Then:*

$$|E| \le \prod_{j=1}^{m} |\phi_j(E)|^{s_j} \text{ for all nonempty finite sets } E \subseteq \mathbb{Z}^d. \qquad (3)$$

$$R_{e_1} = \{S_3[k-1,i,j] \to S_3[k,i,j] : \ 1 \le k < N \ \wedge \ k+1 \le i < N \ \wedge \ k+1 \le j \le i\}$$
$$R_{e_2} = \{S_2[k,j] \to S_3[k,i,j] : \ 0 \le k < N \ \wedge \ k+1 \le i < N \ \wedge \ k+1 \le j \le i\}$$
$$R_{e_3} = \{S_2[k,i] \to S_3[k,i,j] : \ 0 \le k < N \ \wedge \ k+1 \le i < N \ \wedge \ k+1 \le j \le i\}$$
$$R_{e_4} = \{S_3[k-1,i,k] \to S_2[k,i] : \ 1 \le k < N \ \wedge \ k+1 \le i < N\}$$
$$R_{e_5} = \{S_1[k] \to S_2[k,i] : \ 0 \le k < N \ \wedge \ k+1 \le i < N\}$$
$$R_{e_6} = \{S_1[k-1,k,k] \to S_1[k] : \ 1 \le k < N \ \wedge \ k+1 \le i < N\}$$

$$0 \le s_1, s_2, s_3 \le 1 \qquad s_1 = s_2 = s_3 = \frac{1}{2}$$
$$1 \le s_2 + s_3$$
$$1 \le s_1 + s_3 \qquad\qquad |E| \le |\phi_1(E)|^{\frac{1}{2}} \cdot |\phi_2(E)|^{\frac{1}{2}} \cdot |\phi_3(E)|^{\frac{1}{2}}$$
$$1 \le s_1 + s_2$$

$$P_1 = (e_1) \qquad\qquad \text{is a chain path}$$
$$P_2 = (e_2) \qquad\qquad \text{is a broadcast path}$$
$$P_3 = (e_3) \qquad\qquad \text{is a broadcast path}$$

$$\boxed{|\phi_i(E)| \le K, \ i = 1, 2, 3}$$

$$|E| \le K^{3/2}$$

$$\boxed{\begin{array}{l} |\phi_1(E)| + |\phi_2(E)| \le K \\ |\phi_1(E)| + |\phi_3(E)| \le K \qquad |\phi_1(E)| + \frac{1}{2}|\phi_2(E)| + \frac{1}{2}|\phi_3(E)| \le K \end{array}}$$

Cholesky

```c
for(k = 0; k < n; k++) {
        A[k][k] = sqrt(A[k][k]);              //S1
    for(i = k+1; i < n; i++)
        A[i][k] /= A[k][k];                   // S2
    for(i = k+1; i < n; i++)
        for(j = k+1; j <= i; j++)
            A[i][j] -= A[i][k] * A[j][k];     // S3
}
```

$$\text{Dom}(P_1) = \{S_3[k,i,j] : \boxed{1} \le k < N \ \wedge \ k+1 \le i < N \ \wedge \ k+1 \le j \le i\}$$
$$\text{Dom}(P_2) = \{S_3[k,i,j] : \ 0 \le k < N \ \wedge \ k+1 \le i < N \ \wedge \ k+1 \le j \le i\}$$
$$\text{Dom}(P_3) = \{S_3[k,i,j] : \ 0 \le k < N \ \wedge \ k+1 \le i < N \ \wedge \ k+1 \le j \le i\}$$
$$D_S := \text{Dom}(P_1) \cap \text{Dom}(P_2) \cap \text{Dom}(P_3)$$
$$= \{S_3[k,i,j] : \ 1 \le k < N \ \wedge \ k+1 \le i < N \ \wedge \ k+1 \le j \le i\}$$



$$\phi_1(k,i,j) = \text{proj}_{(1,0,0)}(k,i,j) = (0,i,j) \qquad \text{kernel } k_1 = \text{Ker}(\phi_1) = \langle(1,0,0)\rangle$$
$$\phi_2(k,i,j) = (k,j) \qquad\qquad\qquad\qquad \text{kernel } k_2 = \text{Ker}(\phi_2) = \langle(0,1,0)\rangle$$
$$\phi_3(k,i,j) = (k,i) \qquad\qquad\qquad\qquad \text{kernel } k_3 = \text{Ker}(\phi_3) = \langle(0,0,1)\rangle$$

THEOREM 3.10 (BRASCAMP-LIEB INEQUALITY, DISCRETE CASE [CHRIST ET AL. 2013]). *Let d and $d_j$ be nonnegative integers and $\phi_j : \mathbb{Z}^d \mapsto \mathbb{Z}^{d_j}$ be group homomorphisms for $1 \le j \le m$. Let $0 \le s_1, s_2, \ldots, s_m \le 1$. Suppose that:*

$$\text{rank}(H) \le \sum_{j=1}^{m} s_j \cdot \text{rank}(\phi_j(H)) \ \text{for all subgroups } H \text{ of } \mathbb{Z}^d \qquad (2)$$

*Then:*

$$|E| \le \prod_{j=1}^{m} |\phi_j(E)|^{s_j} \ \text{for all nonempty finite sets } E \subseteq \mathbb{Z}^d. \qquad (3)$$

$$R_{e_1} = \{S_3[k-1,i,j] \to S_3[k,i,j] : \ 1 \le k < N \ \wedge \ k+1 \le i < N \ \wedge \ k+1 \le j \le i\}$$
$$R_{e_2} = \{S_2[k,j] \to S_3[k,i,j] : \ 0 \le k < N \ \wedge \ k+1 \le i < N \ \wedge \ k+1 \le j \le i\}$$
$$R_{e_3} = \{S_2[k,i] \to S_3[k,i,j] : \ 0 \le k < N \ \wedge \ k+1 \le i < N \ \wedge \ k+1 \le j \le i\}$$
$$R_{e_4} = \{S_3[k-1,i,k] \to S_2[k,i] : \ 1 \le k < N \ \wedge \ k+1 \le i < N\}$$
$$R_{e_5} = \{S_1[k] \to S_2[k,i] : \ 0 \le k < N \ \wedge \ k+1 \le i < N\}$$
$$R_{e_6} = \{S_1[k-1,k,k] \to S_1[k] : \ 1 \le k < N \ \wedge \ k+1 \le i < N\}$$

$$0 \le s_1, s_2, s_3 \le 1 \qquad s_1 = s_2 = s_3 = \frac{1}{2}$$

$$1 \le s_2 + s_3$$
$$1 \le s_1 + s_3 \qquad |E| \le |\phi_1(E)|^{\frac{1}{2}} \cdot |\phi_2(E)|^{\frac{1}{2}} \cdot |\phi_3(E)|^{\frac{1}{2}}$$
$$1 \le s_1 + s_2 \qquad |\phi_i(E)| \le K, \ i = 1,2,3$$

$P_1 = (e_1)$ is a chain path

$P_2 = (e_2)$ is a broadcast path

$P_3 = (e_3)$ is a broadcast path

$$|E| \le K^{3/2}$$

Cholesky

$$|\phi_1(E)| + |\phi_2(E)| \le K$$
$$|\phi_1(E)| + |\phi_3(E)| \le K \qquad |\phi_1(E)| + \frac{1}{2}|\phi_2(E)| + \frac{1}{2}|\phi_3(E)| \le K \qquad \boxed{|E| \le 2 \cdot (K/3)^{3/2}}$$
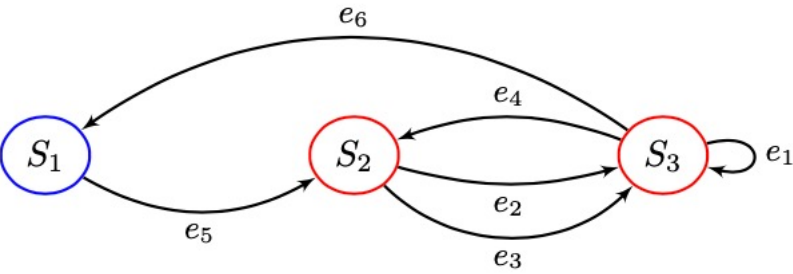
```
for(k = 0; k < n; k++) {
        A[k][k] = sqrt(A[k][k]);               //S1
   for(i = k+1; i < n; i++)
      A[i][k] /= A[k][k];                      // S2
   for(i = k+1; i < n; i++)
       for(j = k+1; j <= i; j++)
            A[i][j] -= A[i][k] * A[j][k];      // S3
}
```



$R_{e_1} = \{S_3[k-1, i, j] \rightarrow S_3[k, i, j] : \ 1 \leq k < N \ \wedge \ k+1 \leq i < N \ \wedge \ k+1 \leq j \leq i\}$

$R_{e_2} = \{S_2[k, j] \rightarrow S_3[k, i, j] : \ 0 \leq k < N \ \wedge \ k+1 \leq i < N \ \wedge \ k+1 \leq j \leq i\}$

$R_{e_3} = \{S_2[k, i] \rightarrow S_3[k, i, j] : \ 0 \leq k < N \ \wedge \ k+1 \leq i < N \ \wedge \ k+1 \leq j \leq i\}$

$R_{e_4} = \{S_3[k-1, i, k] \rightarrow S_2[k, i] : \ 1 \leq k < N \ \wedge \ k+1 \leq i < N\}$

$R_{e_5} = \{S_1[k] \rightarrow S_2[k, i] : \ 0 \leq k < N \ \wedge \ k+1 \leq i < N\}$

$R_{e_6} = \{S_1[k-1, k, k] \rightarrow S_1[k] : \ 1 \leq k < N \ \wedge \ k+1 \leq i < N\}$

$P_1 = (e_1)$      is a chain path

$P_2 = (e_2)$      is a broadcast path

$P_3 = (e_3)$      is a broadcast path

Cholesky

$\mathrm{Dom}\,(P_1) = \{S_3[k, i, j] : \ 1 \leq k < N \ \wedge \ k+1 \leq i < N \ \wedge \ k+1 \leq j \leq i\}$

$\mathrm{Dom}\,(P_2) = \{S_3[k, i, j] : \ 0 \leq k < N \ \wedge \ k+1 \leq i < N \ \wedge \ k+1 \leq j \leq i\}$

$\mathrm{Dom}\,(P_3) = \{S_3[k, i, j] : \ 0 \leq k < N \ \wedge \ k+1 \leq i < N \ \wedge \ k+1 \leq j \leq i\}$

$D_S := \mathrm{Dom}\,(P_1) \cap \mathrm{Dom}\,(P_2) \cap \mathrm{Dom}\,(P_3)$

$\quad = \{S_3[k, i, j] : \ 1 \leq k < N \ \wedge \ k+1 \leq i < N \ \wedge \ k+1 \leq j \leq i\}$

$\phi_1(k, i, j) = \mathrm{proj}_{(1,0,0)}(k, i, j) = (0, i, j)$      kernel $k_1 = \mathrm{Ker}(\phi_1) = \langle(1, 0, 0)\rangle$

$\phi_2(k, i, j) = (k, j)$      kernel $k_2 = \mathrm{Ker}(\phi_2) = \langle(0, 1, 0)\rangle$

$\phi_3(k, i, j) = (k, i)$      kernel $k_3 = \mathrm{Ker}(\phi_3) = \langle(0, 0, 1)\rangle$

THEOREM 3.10 (BRASCAMP-LIEB INEQUALITY, DISCRETE CASE [CHRIST ET AL. 2013]). *Let d and $d_j$ be nonnegative integers and $\phi_j : \mathbb{Z}^d \mapsto \mathbb{Z}^{d_j}$ be group homomorphisms for $1 \leq j \leq m$. Let $0 \leq s_1, s_2, \ldots, s_m \leq 1$. Suppose that:*

$$\mathrm{rank}\,(H) \leq \sum_{j=1}^{m} s_j \cdot \mathrm{rank}\,(\phi_j(H)) \ \text{for all subgroups } H \text{ of } \mathbb{Z}^d \qquad (2)$$

*Then:*

$$|E| \leq \prod_{j=1}^{m} |\phi_j(E)|^{s_j} \ \text{for all nonempty finite sets } E \subseteq \mathbb{Z}^d. \qquad (3)$$

$0 \leq s_1, s_2, s_3 \leq 1$      $s_1 = s_2 = s_3 = \frac{1}{2}$

$1 \leq s_2 + s_3$

$1 \leq s_1 + s_3$      $|E| \leq |\phi_1(E)|^{\frac{1}{2}} \cdot |\phi_2(E)|^{\frac{1}{2}} \cdot |\phi_3(E)|^{\frac{1}{2}}$

$1 \leq s_1 + s_2$      $|\phi_i(E)| \leq K, \ i = 1, 2, 3$

     $|E| \leq K^{3/2}$

$|\phi_1(E)| + |\phi_2(E)| \leq K$

$|\phi_1(E)| + |\phi_3(E)| \leq K$      $|\phi_1(E)| + \frac{1}{2}|\phi_2(E)| + \frac{1}{2}|\phi_3(E)| \leq K$      $\boxed{|E| \leq 2 \cdot (K/3)^{3/2}}$

the total number of statements in the execution

the IO of a segment

We consider a segmentation of the execution every T IOs. For each segment, we have at most K = S + T input data to perform the statements in the segment.

$$Q \geq T \times \left\lfloor \frac{N^3/6}{2 \cdot (K/3)^{3/2}} \right\rfloor$$

The minimum number of K-partition in a schedule

the maximum number of statements in a segment

is greater than or equal to

the IO of any schedule

Cholesky

the total number of statements in the execution

the IO of a segment

We consider a segmentation of the execution every $T$ IOs. For each segment, we have at most $K = S + T$ input data to perform the statements in the segment.

$$Q \geq T \times \left\lfloor \frac{N^3/6}{2 \cdot (K/3)^{3/2}} \right\rfloor$$

The minimum number of K-partition in a schedule

the maximum number of statements in a segment

is greater than or equal to

the IO of any schedule

Setting $T = 2S$ (so $K = S + T = 3S$)

$$Q \geq (2S) \times \frac{N^3/6}{2S^{3/2}} = \boxed{\frac{N^3}{6\sqrt{S}}}$$

Cholesky

```
for(k = 0; k < n; k++) {
        A[k][k] = sqrt(A[k][k]);                    //S1
    for(i = k+1; i < n; i++)
        A[i][k] /= A[k][k];                         // S2
    for(i = k+1; i < n; i++)
        for(j = k+1; j <= i; j++)
            A[i][j] -= A[i][k] * A[j][k];            // S3
}
```

$$\text{Dom}(P_1) = \{S_3[k,i,j] :\ 1 \leq k < N \ \wedge \ k+1 \leq i < N \ \wedge \ k+1 \leq j \leq i\}$$
$$\text{Dom}(P_2) = \{S_3[k,i,j] :\ 0 \leq k < N \ \wedge \ k+1 \leq i < N \ \wedge \ k+1 \leq j \leq i\}$$
$$\text{Dom}(P_3) = \{S_3[k,i,j] :\ 0 \leq k < N \ \wedge \ k+1 \leq i < N \ \wedge \ k+1 \leq j \leq i\}$$
$$D_S := \text{Dom}(P_1) \cap \text{Dom}(P_2) \cap \text{Dom}(P_3)$$
$$= \{S_3[k,i,j] :\ 1 \leq k < N \ \wedge \ k+1 \leq i < N \ \wedge \ k+1 \leq j \leq i\}$$

$$\phi_1(k,i,j) = \text{proj}_{(1,0,0)}(k,i,j) = (0,i,j) \qquad \text{kernel } k_1 = \text{Ker}(\phi_1) = \langle(1,0,0)\rangle$$
$$\phi_2(k,i,j) = (k,j) \qquad \text{kernel } k_2 = \text{Ker}(\phi_2) = \langle(0,1,0)\rangle$$
$$\phi_3(k,i,j) = (k,i) \qquad \text{kernel } k_3 = \text{Ker}(\phi_3) = \langle(0,0,1)\rangle$$

THEOREM 3.10 (BRASCAMP-LIEB INEQUALITY, DISCRETE CASE [CHRIST ET AL. 2013]). *Let d and $d_j$ be nonnegative integers and $\phi_j : \mathbb{Z}^d \mapsto \mathbb{Z}^{d_j}$ be group homomorphisms for $1 \leq j \leq m$. Let $0 \leq s_1, s_2, \ldots, s_m \leq 1$. Suppose that:*

$$\text{rank}(H) \leq \sum_{j=1}^{m} s_j \cdot \text{rank}(\phi_j(H)) \ \text{for all subgroups } H \text{ of } \mathbb{Z}^d \qquad (2)$$

*Then:*



$S_1$ $S_2$ $S_3$ graph with edges $e_1, e_2, e_3, e_4, e_5, e_6$

$$R_{e_1} = \{S_3[k-1,i,j] \rightarrow S_3[k,i,j] :\ 1 \leq k < N \ \wedge \ k+1 \leq i < N \ \wedge \ k+1 \leq j \leq i\}$$
$$R_{e_2} = \{S_2[k,j] \rightarrow S_3[k,i,j] :\ 0 \leq k < N \ \wedge \ k+1 \leq i < N \ \wedge \ k+1 \leq j \leq i\}$$
$$R_{e_3} = \{S_2[k,i] \rightarrow S_3[k,i,j] :\ 0 \leq k < N \ \wedge \ k+1 \leq i < N \ \wedge \ k+1 \leq j \leq i\}$$
$$R_{e_4} = \{S_3[k-1,i,k] \rightarrow S_2[k,i] :\ 1 \leq k < N \ \wedge \ k+1 \leq i < N\}$$
$$R_{e_5} = \{S_1[k] \rightarrow S_2[k,i] :\ 0 \leq k < N \ \wedge \ k+1 \leq i < N\}$$
$$R_{e_6} = \{S_1[k-1,k,k] \rightarrow S_1[k] :\ 1 \leq k < N \ \wedge \ k+1 \leq i < N\}$$

$0 \leq s_1, s_2, s_3 \leq 1$

$1 \leq s_2 + s_3$

$1 \leq s_1 + s_3$

$1 \leq s_1 + s_2$

| | |
|---|---|
| $P_1 = (e_1)$ | is a chain path |
| $P_2 = (e_2)$ | is a broadcast path |
| $P_3 = (e_3)$ | is a broadcast path |

$|\phi_1(E)| + |\phi_2(E)| \leq K$

$|\phi_1(E)| + |\phi_3(E)| \leq K$

$|\phi_1$

Cholesky

We consider a segmentation of the execution every T IOs. For each segment, we have at most K = S + T input data to perform the statements in the segment.

the total number of statements in the execution

the IO of a segment

$$Q \geq T \times \left\lfloor \frac{N^3/6}{2 \cdot (K/3)^{3/2}} \right\rfloor$$

The minimum number of K-partition in a schedule

the maximum number of statements in a segment
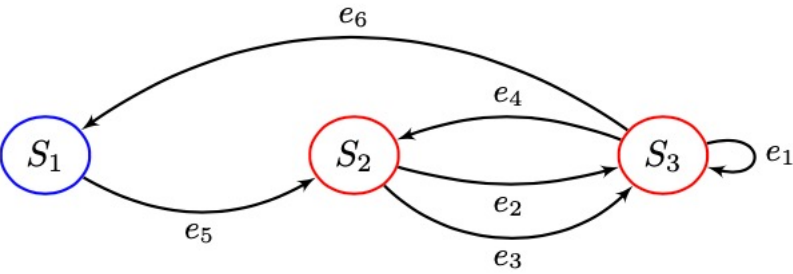
is greater than or equal to

the IO of any schedule

Setting $T = 2S$ (so $K = S + T = 3S$)

$$Q \geq (2S) \times \frac{N^3/6}{2S^{3/2}} = \boxed{\frac{N^3}{6\sqrt{S}}}$$

```
for(k = 0; k < n; k++) {
        A[k][k] = sqrt(A[k][k]);                      //S1
    for(i = k+1; i < n; i++)
        A[i][k] /= A[k][k];                           // S2
    for(i = k+1; i < n; i++)
        for(j = k+1; j <= i; j++)
            A[i][j] -= A[i][k] * A[j][k];             // S3
}
```

$\text{Dom}(P_1) = \{S_3[k, i, j]: \; 1 \leq k < N \; \wedge \; k+1 \leq i < N \; \wedge \;$
$1 \leq i < N \; \wedge \;$
$< N \; \wedge \;$
$j \leq i\}$

$\phi_1(k, i, j) = \text{proj}_{(1,0,0)}(k, i, j) = (0, i, j)$  $(1, 0, 0)\rangle$

$\phi_2(k$  $1, 0)\rangle$

$\phi_3(k$  $\text{ker}$  $0, 1)\rangle$

THE... NEQUALITY, DISCRETE CASE [CHRIS...
$d_j$ be n... $\mathbb{Z}^d \mapsto \mathbb{Z}^{d_j}$ be group homomorph...
$0 \leq s_1, $

$\text{rank}(H) \leq \sum_{j=1} s_j \cdot \text{rank}(\phi_j(H)) \; \text{for all subgroups } H \; o$

Then:



the total number of statements in the execution
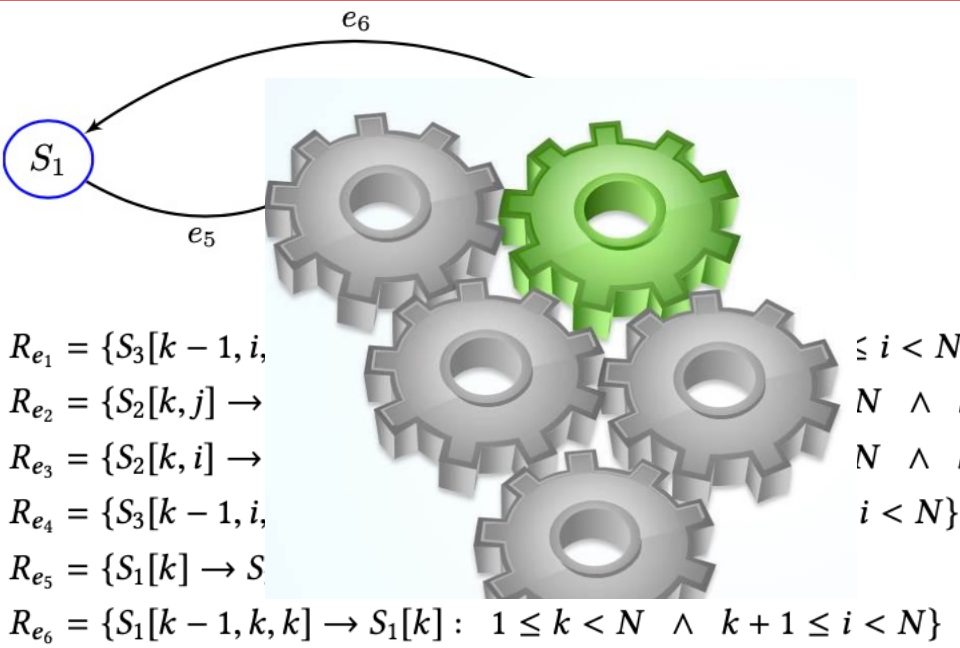the IO of a segment

We consider a... T IOs. For
each segment... a to perform
the statement...

$0 \leq s_1, s_2, s_3 \leq 1$
$1 \leq s_2 + s_3$
$1 \leq s_1 + s_3$
$1 \leq s_1 + s_2$

$Q \geq T \times \left\lfloor \dfrac{N^3/6}{2 \cdot (K/3)^{3/2}} \right\rfloor$

the maximum number of statements in a segment
is greater than or equal to
the IO of any schedule

$e_6$
$S_1$
$e_5$

$R_{e_1} = \{S_3[k-1, i,$  $\leq i < N \; \wedge \; k+1 \leq j \leq i\}$
$R_{e_2} = \{S_2[k, j] \to$  $N \; \wedge \; k+1 \leq j \leq i\}$
$R_{e_3} = \{S_2[k, i] \to$  $N \; \wedge \; k+1 \leq j \leq i\}$
$R_{e_4} = \{S_3[k-1, i,$  $i < N\}$
$R_{e_5} = \{S_1[k] \to S$
$R_{e_6} = \{S_1[k-1, k, k] \to S_1[k]: \; 1 \leq k < N \; \wedge \; k+1 \leq i < N\}$

IOLB: I/O Lower Bound
https://iocomplexity.corse.inria.fr/iolb

Cholesky

Setting $T = 2S$ (so $K = S + T = 3S$)

$|\phi_1(E)| + |\phi_2(E)| \leq K$
$|\phi_1(E)| + |\phi_3(E)| \leq K$   $|\phi_1($

$Q \geq (2S) \times \dfrac{N^3/6}{2S^{3/2}} = \boxed{\dfrac{N^3}{6\sqrt{S}}}$

# IOLB: I/O Lower Bound
## https://iocomplexity.corse.inria.fr/iolb

- PET: parse the C code (PET = Polyhedral Extraction Tool )
- DFG (in house) Data Flow Graph
- Piplib: for the ILP (PIP = Parametric Integer Programming )
- GINAC: manipulate symbolic expressions
- Lib ISL: (Integer Set Library) Barvinok: omputes cardinalities of Z polyhedron. Main tool for polyhedral compiling.

```
for(k = 0; k < n; k++) {
    for(i = k+1; i < n; i++)
        A[i][k] /= A[k][k];                    // S1
    for(i = k+1; i < n; i++)
        for(j = k+1; j < n; j++)
            A[i][j] += A[i][k] * A[k][j];       // S2
}
```



$$Q \geq (2S)\frac{N^3/3}{S^{3/2}} = \boxed{\frac{2N^3}{3\sqrt{S}}}$$

Operational intensity of LU is at most sqrt(S)

LU

# LU

similar algorithm as Béreux's Cholesky (SIMAX 2008)

```
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];
    }
  }
}
```

Reorder statements to add blocking in the update

explicit cache directives

Operational intensity of LU is at least sqrt(S)

```
for( j = 0; j < n; j+=nb ){
    jb = (( nb < n-j ) ? ( nb ) : ( n-j ));
    for( k= j+jb; k < n; k+=nb ){
        kb = (( nb < n-k ) ? ( nb ) : ( n-k ));
//      read A[ k:k+kb, j:j+jb ]
        for( kk= 0; kk < j; kk++ ){
//          read A[ k:k+kb, kk ]
//          read A[ kk, j:j+jb ]
            for( ii= k; ii < k+kb; ii++ )
                for( jj= j; jj < j+jb; jj++ )
                    A[ ii ][ jj ] -= A[ ii ][ kk ] * A[ kk ][ jj ];
//          erase A[ k:k+kb, kk ]
//          erase A[ kk, j:j+jb ]
        }
//      write A[ k:k+kb, j:j+jb ]
    }
    for( k= j+jb; k < n; k+=nb ){
        kb = (( nb < n-k ) ? ( nb ) : ( n-k ));
//      read A[ j:j+jb, k:k+kb ]
        for( kk= 0; kk < j; kk++ ){
//          read A[ j:j+jb, kk ]
//          read A[ kk, k:k+kb ]
            for( ii= j; ii < j+jb; ii++ )
                for( jj= k; jj < k+kb; jj++ )
                    A[ ii ][ jj ] -= A[ ii ][ kk ] * A[ kk ][ jj ];
//          erase A[ j:j+jb, kk ]
//          erase A[ kk, k:k+kb ]
        }
//      write A[ j:j+jb, k:k+kb ]
    }

//  read A[ 1:jb, 1:jb ]
    for( kk= 0; kk < j; kk++ ){
//      read A[ 1:jb, kk ]
//      read A[ kk, 1:jb ]
        for( ii= j; ii < j+jb; ii++ )
            for( jj= j; jj < j+jb; jj++ )
                A[ ii ][ jj ] -= A[ ii ][ kk ] * A[ kk ][ jj ];
//      erase A[ 1:jb, kk ]
//      erase A[ kk, 1:jb ]
    }
    for( jj = j; jj < j+jb; jj++ ){
        for(ii = j; ii < jj; ii++)
            for(kk = j; kk < ii; kk++)
                A[ ii ][ jj ] -= A[ ii ][ kk ] * A[ kk ][ jj ];
        for(ii = jj; ii < j+jb; ii++)
            for(kk = j; kk < jj; kk++)
                A[ ii ][ jj ] -= A[ ii ][ kk ] * A[ kk ][ jj ];
        for( ii = jj+1; ii < j+jb; ii++ )
            A[ ii ][ jj ] /= A[ jj ][ jj ];
    }
    for(ii = j+jb; ii < n; ii++){
        // read A[ ii, 1:jb ]
        for( jj = j; jj < j+jb; jj++ ){
            for(kk = j; kk < jj; kk++)
                A[ ii ][ jj ] -= A[ ii ][ kk ] * A[ kk ][ jj ];

            A[ ii ][ jj ] /= A[ jj ][ jj ];
        }
        // write A[ ii, 1:jb ]
    }
    for(jj = j+jb; jj < n; jj++){
        // read A[ 1:jb, jj ]
        for( ii = j; ii < j+jb; ii++ ){
            for(kk = j; kk < ii; kk++){
                A[ ii ][ jj ] -= A[ ii ][ kk ] * A[ kk ][ jj ];
            }
        }
        // write A[ 1:jb, jj ]
    }
//  write A[ 1:jb, 1:jb ]
}
```

| kernel | # input data | #ops | ratio | $OI_{up}$ | $OI_{manual}$ | ratio |
|---|---|---|---|---|---|---|
| 2mm | $N_iN_k + N_kN_j$ $+N_jN_l + N_iN_l$ | $NiN_jN_k$ $+N_iN_jN_l$ | – | $\sqrt{S}$ | $\sqrt{S}$ | 1 ✓ |
| 3mm | $N_iN_k + N_kN_j$ $+N_jN_m + N_mN_l$ | $NiN_jN_k + N_jN_lN_m$ $+N_iN_jN_l$ | – | $\sqrt{S}$ | $\sqrt{S}$ | 1 ✓ |
| cholesky | $\frac{1}{2}N^2$ | $\frac{1}{3}N^3$ | $\frac{2}{3}N$ | $2\sqrt{S}$ | $\sqrt{S}$ | 2 |
| correlation | $MN$ | $M^2N$ | $M$ | $2\sqrt{S}$ | $\sqrt{S}$ | 2 |
| covariance | $MN$ | $M^2N$ | $M$ | $2\sqrt{S}$ | $\sqrt{S}$ | 2 |
| doitgen | $N_pN_qN_r$ | $2N_qn_rN_p^2$ | $2N_p$ | $\sqrt{S}$ | $\sqrt{S}$ | 1 ✓ |
| fdtd-2d | $3N_xN_y$ | $11N_xN_yT$ | $\frac{11}{3}T$ | $22\sqrt{2}\sqrt{S}$ | $\frac{11}{24}\sqrt{3}\sqrt{S}$ | $\frac{48\sqrt{2}}{\sqrt{3}}$ |
| floyd-warshall | $N^2$ | $2N^3$ | $2N$ | $2\sqrt{S}$ | $\sqrt{S}$ | 2 |
| gemm | $N_iN_j + N_jN_k + N_iN_k$ | $2N_iN_jN_k$ | – | $\sqrt{S}$ | $\sqrt{S}$ | 1 ✓ |
| heat-3d | $N^3$ | $30N^3T$ | $30T$ | $\frac{160}{3\sqrt[3]{3}}\sqrt[3]{S}$ | $\frac{5}{2}\sqrt[3]{S}$ | $\frac{64}{3\sqrt[3]{3}}$ |
| jacobi-1d | $N$ | $6NT$ | $6T$ | $24S$ | $\frac{3}{2}S$ | 16 |
| jacobi-2d | $N^2$ | $10N^2T$ | $10T$ | $15\sqrt{3}\sqrt{S}$ | $\frac{5}{4}\sqrt{S}$ | $12\sqrt{3}$ |
| 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 | $N^2$ | $9N^2T$ | $9T$ | $\frac{27\sqrt{3}}{2}\sqrt{S}$ | $\frac{9}{4}\sqrt{S}$ | $6\sqrt{3}$ |
| symm | $\frac{1}{2}M^2 + 2MN$ | $2M^2N$ | – | $\sqrt{S}$ | $\sqrt{S}$ | 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 | $\frac{1}{2}M^2 + MN$ | $M^2N$ | – | $\sqrt{S}$ | $\sqrt{S}$ | 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 | $N^2$ | $30N^2T$ | $30T$ | 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 | $2\sqrt{S}$ |

polybench test suite

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 )

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 )

**Table 2.** Complete Lower Bound Formulae for PolyBench obtained with the current version of IOLB

gemm — $\max\left(N_i N_j + N_j N_k + N_i N_k + 2, \frac{2}{\sqrt{S}} N_i N_j (N_k - 1) + 2N_i + 2N_j + N_k - 4\sqrt{2S}\right)$ — $2\frac{1}{\sqrt{S}} N_i N_j N_k$

gemm

```
for (i = 0; i < _PB_NI; i++) {
    for (j = 0; j < _PB_NJ; j++)
        C[i][j] *= beta;
    for (k = 0; k < _PB_NK; k++) {
        for (j = 0; j < _PB_NJ; j++)
            C[i][j] += alpha * A[i][k] * B[k][j];
    }
}
```

Operational intensity of GEMM is at sqrt(S)

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 )

| kernel | Complete Lower Bound Formulae for PolyBench obtained with the current version of IOLB | asymptotic simplified formula |
|---|---|---|
| 2mm | $\max(N_iN_j+N_jN_k+N_iN_k+N_jN_l+2,$ $\left(\frac{2}{\sqrt{S}}N_iN_j(N_k-1)+2N_i+2N_j+N_k-4\sqrt{2S}\right)$ $+\left(\frac{2}{\sqrt{S}}N_iN_l(N_j-1)+2N_i+2N_l+N_j-4\sqrt{2S}\right)$ $-2N_iN_j-2)$ | $\frac{2}{\sqrt{S}}N_iN_jN_k$ $+\frac{2}{\sqrt{S}}N_iN_lN_j$ |
| 3mm | $\max(N_iN_k+N_jN_k+N_jN_m+N_lN_m,$ $\left(\frac{2}{\sqrt{S}}N_iN_j(N_k-1)+2N_i+2N_j+N_k-4\sqrt{2S}\right)$ $+\left(\frac{2}{\sqrt{S}}N_iN_l(N_j-1)+2N_i+2N_l+N_j-4\sqrt{2S}\right)$ $+\left(\frac{2}{\sqrt{S}}N_jN_l(N_m-1)+2N_j+2N_l+N_m-4\sqrt{2S}\right)$ $-2N_jN_i-2N_jN_l-N_iN_l-6)$ | $\frac{2}{\sqrt{S}}N_iN_jN_k$ $+\frac{2}{\sqrt{S}}N_iN_lN_j$ $+\frac{2}{\sqrt{S}}N_jN_lN_m$ |
| adi | $4N^2+\max(0,(N^2-4N-S+5)(T-2))$ | $N^2T$ |
| atax | $MN+N+\max(0,\frac{1}{8}\frac{1}{S}((2M-1-8S)(2N-1-8S)-1)-10S+2)$ | $MN$ |
| bicg | $MN+M+N+\max(0,\frac{1}{8}\frac{1}{S}((2M-1-8S)(2N-1-8S)-1)-10S+2)$ | $MN$ |
| cholesky | $\max(\frac{1}{2}N(N+1),\frac{1}{6}\frac{1}{\sqrt{S}}(N-1)(N-2)(N-3)+\frac{1}{2\sqrt{2}}\frac{1}{S}(N-1)(N-2)$ $-(N-2)(N-7)-4\sqrt{2S})$ | $\frac{1}{6}\frac{1}{\sqrt{S}}N^3$ |
| correlation | $\max(MN+2,\frac{1}{2}\frac{1}{\sqrt{S}}M(M-1)(N-1+\frac{\sqrt{2}}{2}\frac{1}{\sqrt{S}})-\frac{1}{2}(M-3)(M+2N-2)+2-4S\sqrt{2})$ | $\frac{1}{2}\frac{1}{\sqrt{S}}M^2N$ |
| covariance | $\max(MN+2,\frac{1}{2}\frac{1}{\sqrt{S}}M(M-1)(N-1+\frac{\sqrt{2}}{2}\frac{1}{\sqrt{S}})-\frac{1}{2}(M-3)(M+2N-2)+1-4S\sqrt{2})$ | $\frac{1}{2}\frac{1}{\sqrt{S}}M^2N$ |
| deriche | $HW+1$ | $HW$ |
| doitgen | $\max(N_p^2+N_pN_qN_r,\frac{2}{\sqrt{S}}N_qN_rN_p(N_p-1+\frac{1}{\sqrt{2}}\frac{1}{\sqrt{S}})-N_qN_r(N_p-1)+2N_p-8\sqrt{2S}-1)$ | $2\frac{1}{\sqrt{S}}N_qN_rN_p^2$ |
| durbin | $2N+\max(0,\frac{1}{2}(N-3)(N-2-2S))$ | $\frac{1}{2}N^2$ |
| fdtd-2d | $\max(3N_xN_y-N_y+T-1,\frac{1}{2\sqrt{2}}\frac{1}{\sqrt{S}}(N_x-2)(N_y-2)(T-1)+2(N_x+2)(N_y+2)$ $-T(N_x+N_y-6)-N_y-S-23)$ | $\frac{1}{2\sqrt{2}}\frac{1}{\sqrt{S}}N_xN_yT$ |
| floyd-warshall | $\max(N^2,\frac{1}{\sqrt{S}}(N-1)^3-(6N-19)(N-2)-8\sqrt{2S})$ | $\frac{1}{\sqrt{S}}N^3$ |
| gemm | $\max(N_iN_j+N_jN_k+N_iN_k+2,\frac{2}{\sqrt{S}}N_iN_j(N_k-1)+2N_i+2N_j+N_k-4\sqrt{2S})$ | $2\frac{1}{\sqrt{S}}N_iN_jN_k$ |
| gemver | $N^2+8N+2+\max(0,\frac{1}{4}\frac{1}{S}(3N-2)(N-8S)-3S+1)$ | $N^2$ |
| gesummv | $2N^2+N+2+\max(0,\frac{1}{2}\frac{1}{S}(N-1)(N-8S)-2S)$ | $2N^2$ |
| gramschmidt | $\max(MN,\frac{1}{\sqrt{S}}MN(N-3)-M(N-5-\frac{2}{\sqrt{S}})-\frac{1}{2}(N-1)(N-6)-4\sqrt{2S}-3)$ | $\frac{1}{\sqrt{S}}MN^2$ |
| heat-3d | $\max((N-10)(N+2)^2,\frac{9}{16}\frac{\sqrt[3]{5}}{\sqrt{S}}\frac{1}{\sqrt{S}}(T-1)(N-3)^3-3(T-7)(N-3)(N-4)$ $+42N-T-\frac{9}{4}\frac{\sqrt[3]{5}}{\sqrt[4]{S}}S-111)$ | $\frac{9}{16}\frac{\sqrt[3]{5}}{\sqrt{S}}N^3T$ |
| jacobi-1d | $\max(2+n,\frac{1}{4}\frac{1}{S}(T-1)(N-3)-T-S+7)$ | $\frac{1}{4}\frac{1}{S}NT$ |
| jacobi-2d | $\max((N-2)(N+6),\frac{2}{3\sqrt{3}}\frac{1}{\sqrt{S}}(N-3)^2(T-1)-\frac{4\sqrt{2}}{3\sqrt{3}}S-(T-7)(2N-7)+14)$ | $\frac{2}{3\sqrt{3}}\frac{1}{\sqrt{S}}N^2T$ |
| lu | $\max(N^2,\frac{2}{3}\frac{1}{\sqrt{S}}(N-2)(N^2-4N+6)-2(N^2-10N+18)-8\sqrt{2S})$ | $\frac{2}{3}\frac{1}{\sqrt{S}}N^3$ |
| ludcmp | $\max(N^2+N,\frac{1}{3}\frac{1}{\sqrt{S}}(2N-3)(N-1)(N-2)\sqrt{2}\frac{1}{S}(N-1)(N-2)-(2N^2-15N+19)-16\sqrt{2S})$ | $\frac{2}{3}\frac{1}{\sqrt{S}}N^3$ |
| mvt | $N^2+4N+\max(0,\frac{1}{6}\frac{1}{S}N(N-1)-2S-4N+4)$ | $N^2$ |
| nussinov | $\frac{1}{2}N^2+\frac{5}{2}N-1+\max(0,\frac{1}{6}\frac{1}{\sqrt{S}}(N-3)(N-4)(N-5)+\frac{1}{4}\frac{1}{S}\sqrt{2}(3N^2-19N+6)$ $-(N^2-13N+22)-8\sqrt{2S})$ | $\frac{1}{6}\frac{1}{\sqrt{S}}N^3$ |
| seidel-2d | $\max(N^2,\frac{2}{3\sqrt{3}}\frac{1}{\sqrt{S}}(N-3)^2(T-1)-(2N-7)(T-5)-\frac{4\sqrt{2}}{3\sqrt{3}}S+12)$ | $\frac{2}{3\sqrt{3}}\frac{1}{\sqrt{S}}N^2T$ |
| symm | $\max(\frac{1}{2}M(M+1)+2MN+2,2\frac{1}{\sqrt{S}}(M-1)(M-2)N-\frac{1}{2}((4N+M)(M-5))$ $+5(M-2)-8\sqrt{2S})$ | $2\frac{1}{\sqrt{S}}M^2N$ |
| syr2k | $\max(2+2MN+\frac{1}{2}N(N+1),\frac{1}{\sqrt{S}}(M-1)(N+1)N+M+4N-4\sqrt{2S})$ | $\frac{1}{\sqrt{S}}MN^2$ |
| syrk | $\max(MN+\frac{1}{2}(N+1)N+2,\frac{1}{2}\frac{1}{\sqrt{S}}(M-1)(N+1)N-(M-4)(N-1)-2\sqrt{2S}+4)$ | $\frac{1}{2}\frac{1}{\sqrt{S}}MN^2$ |
| trisolv | $\frac{1}{2}N(N+1)+N+\max(0,\frac{1}{8}\frac{1}{S}(N-1)(N-2)-2N-S+5)$ | $\frac{1}{2}N^2$ |
| trmm | $\max(\frac{1}{2}M(M-1)+MN+1,\frac{1}{\sqrt{S}}(M-2+\frac{\sqrt{2}}{\sqrt{S}})(M-1)N-(M-4)(N-2)-8\sqrt{2S}+5)$ | $\frac{1}{\sqrt{S}}M^2N$ |

**Table 2.** Complete Lower Bound Formulae for PolyBench obtained with the current version of IOLB

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

cholesky

```
for (i = 0; i < _PB_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_FUN(A[i][i]);
}
```

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 )

| kernel | Complete Lower Bound Formulae for POLYBENCH obtained with the current version of IOLB | asymptotic simplified formula |
|---|---|---|
| 2mm | $\max(N_iN_j+N_jN_k+N_iN_k+N_jN_l+2,$ $\left(\frac{2}{\sqrt{S}}N_iN_j(N_k-1)+2N_i+2N_j+N_k-4\sqrt{2}S\right)$ $+\left(\frac{2}{\sqrt{S}}N_iN_l(N_j-1)+2N_l+2N_l+N_j-4\sqrt{2}S\right)$ $-2N_iN_j-2)$ | $\frac{2}{\sqrt{S}}N_iN_jN_k$ $+\frac{2}{\sqrt{S}}N_iN_lN_j$ |
| 3mm | $\max(N_iN_k+N_jN_k+N_jN_m+N_lN_m,$ $\left(\frac{2}{\sqrt{S}}N_iN_j(N_k-1)+2N_i+2N_j+N_k-4\sqrt{2}S\right)$ $+\left(\frac{2}{\sqrt{S}}N_iN_l(N_j-1)+2N_i+2N_l+N_j-4\sqrt{2}S\right)$ $+\left(\frac{2}{\sqrt{S}}N_jN_l(N_m-1)+2N_j+2N_l+N_m-4\sqrt{2}S\right)$ $-2N_jN_i-2N_jN_l-N_iN_l-6)$ | $\frac{2}{\sqrt{S}}N_iN_jN_k$ $+\frac{2}{\sqrt{S}}N_iN_lN_j$ $+\frac{2}{\sqrt{S}}N_jN_lN_m$ |
| adi | $4N^2+\max(0,(N^2-4N-S+5)(T-2))$ | $N^2T$ |
| atax | $MN+N+\max(0,\frac{1}{4}\frac{1}{S}((2M-1-8S)(2N-1-8S)-1)-10S+2)$ | $MN$ |
| bicg | $MN+M+N+\max(0,\frac{1}{4}\frac{1}{S}((2M-1-8S)(2N-1-8S)-1)-10S+2)$ | $MN$ |
| cholesky | $\max(\frac{1}{2}N(N+1),\frac{1}{6}\frac{1}{\sqrt{S}}(N-1)(N-2)(N-3)+\frac{1}{2\sqrt{2}}\frac{1}{S}(N-1)(N-2)$ $-(N-2)(N-7)-4\sqrt{2}S)$ | $\frac{1}{6}\frac{1}{\sqrt{S}}N^3$ |
| correlation | $\max(MN+2,\frac{1}{2}\frac{1}{\sqrt{S}}M(M-1)(N-1+\frac{\sqrt{2}}{2}\frac{1}{\sqrt{S}})-\frac{1}{2}(M-3)(M+2N-2)+2-4S\sqrt{2})$ | $\frac{1}{2}\frac{1}{\sqrt{S}}M^2N$ |
| covariance | $\max(MN+2,\frac{1}{2}\frac{1}{\sqrt{S}}M(M-1)(N-1+\frac{\sqrt{2}}{2}\frac{1}{\sqrt{S}})-\frac{1}{2}(M-3)(M+2N-2)+1-4S\sqrt{2})$ | $\frac{1}{2}\frac{1}{\sqrt{S}}M^2N$ |
| deriche | $HW+1$ | $HW$ |
| doitgen | $\max(N_p^2+N_pN_qN_r,\frac{2}{\sqrt{S}}N_qN_rN_p(N_p-1+\frac{1}{\sqrt{2}}\frac{1}{\sqrt{S}})-N_qN_r(N_p-1)+2N_p-8\sqrt{2}S-1)$ | $2\frac{1}{\sqrt{S}}N_qN_rN_p^2$ |
| durbin | $2N+\max(0,\frac{1}{2}(N-3)(N-2-2S))$ | $\frac{1}{2}N^2$ |
| fdtd-2d | $\max(3N_xN_y-N_y+T-1,\frac{1}{2\sqrt{2}}\frac{1}{\sqrt{S}}(N_x-2)(N_y-2)(T-1)+2(N_x+2)(N_y+2)$ $-T(N_x+N_y-6)-N_y-S-23)$ | $\frac{1}{2\sqrt{2}}\frac{1}{\sqrt{S}}N_xN_yT$ |
| floyd-warshall | $\max(N^2,\frac{1}{\sqrt{S}}(N-1)^3-(6N-19)(N-2)-8\sqrt{2}S)$ | $\frac{1}{\sqrt{S}}N^3$ |
| gemm | $\max(N_iN_j+N_jN_k+N_iN_k+2,\frac{2}{\sqrt{S}}N_iN_j(N_k-1)+2N_i+2N_j+N_k-4\sqrt{2}S)$ | $2\frac{1}{\sqrt{S}}N_iN_jN_k$ |
| gemver | $N^2+8N+2+\max(0,\frac{1}{4}\frac{1}{S}(3N-2)(N-8S)-3S+1)$ | $N^2$ |
| gesummv | $2N^2+N+2+\max(0,\frac{1}{2}\frac{1}{S}(N-1)(N-8S)-2S)$ | $2N^2$ |
| gramschmidt | $\max(MN,\frac{1}{\sqrt{S}}MN(N-3)-M(N-5-\frac{2}{\sqrt{S}})-\frac{1}{2}(N-1)(N-6)-4\sqrt{2}S-3)$ | $\frac{1}{\sqrt{S}}MN^2$ |
| heat-3d | $\max((N-10)(N+2)^2,\frac{9}{16}\frac{\sqrt[3]{3}}{\sqrt{S}}\frac{1}{\sqrt{S}}(T-1)(N-3)^3-3(T-7)(N-3)(N-4)$ $+42N-T-\frac{9}{4}\frac{\sqrt[3]{3}}{\sqrt[4]{4}}S-111)$ | $\frac{9}{16}\frac{\sqrt[3]{3}}{\sqrt{S}}N^3T$ |
| jacobi-1d | $\max(2+n,\frac{1}{4}\frac{1}{S}(T-1)(N-3)-T-S+7)$ | $\frac{1}{4}\frac{1}{S}NT$ |
| jacobi-2d | $\max((N-2)(N+6),\frac{2}{3\sqrt{3}}\frac{1}{\sqrt{S}}(N-3)^2(T-1)-\frac{4\sqrt{2}}{3\sqrt{3}}S-(T-7)(2N-7)+14)$ | $\frac{2}{3\sqrt{3}}\frac{1}{\sqrt{S}}N^2T$ |
| lu | $\max(N^2,\frac{2}{3}\frac{1}{\sqrt{S}}(N-2)(N^2-4N+6)-2(N^2-10N+18)-8\sqrt{2}S)$ | $\frac{2}{3}\frac{1}{\sqrt{S}}N^3$ |
| ludcmp | $\max(N^2+N,\frac{1}{3}\frac{1}{\sqrt{S}}(2N-3)(N-1)(N-2)\sqrt{2}\frac{1}{S}(N-1)(N-2)-(2N^2-15N+19)-16\sqrt{2}S)$ | $\frac{2}{3}\frac{1}{\sqrt{S}}N^3$ |
| mvt | $N^2+4N+\max(0,\frac{1}{6}\frac{1}{S}N(N-1)-2S-4N+4)$ | $N^2$ |
| nussinov | $\frac{1}{2}N^2+\frac{5}{2}N-1+\max(0,\frac{1}{4}\frac{1}{\sqrt{S}}(N-3)(N-4)(N-5)+\frac{1}{4}\frac{1}{S}\sqrt{2}(3N^2-19N+6)$ $-(N^2-13N+22)-8\sqrt{2}S)$ | $\frac{1}{6}\frac{1}{\sqrt{S}}N^3$ |
| seidel-2d | $\max(N^2,\frac{2}{3\sqrt{3}}\frac{1}{\sqrt{S}}(N-3)^2(T-1)-(2N-7)(T-5)-\frac{4\sqrt{2}}{3\sqrt{3}}S+12)$ | $\frac{2}{3\sqrt{3}}\frac{1}{\sqrt{S}}N^2T$ |
| symm | $\max(\frac{1}{2}M(M+1)+2MN+2,2\frac{1}{\sqrt{S}}(M-1)(M-2)N-\frac{1}{2}((4N+M)(M-5))$ $+5(M-2)-8\sqrt{2}S)$ | $2\frac{1}{\sqrt{S}}M^2N$ |
| syr2k | $\max(2+2MN+\frac{1}{2}N(N+1),\frac{1}{\sqrt{S}}(M-1)(N+1)N+M+4N-4\sqrt{2}S)$ | $\frac{1}{\sqrt{S}}MN^2$ |
| syrk | $\max(MN+\frac{1}{2}(N+1)N+2,\frac{1}{2}\frac{1}{\sqrt{S}}(M-1)(N+1)N-(M-4)(N-1)-2\sqrt{2}S+4)$ | $\frac{1}{2}\frac{1}{\sqrt{S}}MN^2$ |
| trisolv | $\frac{1}{2}N(N+1)+N+\max(0,\frac{1}{8}\frac{1}{S}(N-1)(N-2)-2N-S+5)$ | $\frac{1}{2}N^2$ |
| trmm | $\max(\frac{1}{2}M(M-1)+MN+1,\frac{1}{\sqrt{S}}(M-2+\frac{\sqrt{2}}{\sqrt{S}})(M-1)N-(M-4)(N-2)-8\sqrt{2}S+5)$ | $\frac{1}{\sqrt{S}}M^2N$ |

**Table 2.** Complete Lower Bound Formulae for POLYBENCH obtained with the current version of IOLB

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

Off by a factor of 2 from the lowest known upper bound

cholesky

```
for (i = 0; i < _PB_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_FUN(A[i][i]);
}
```

| kernel | # input data | #ops | ratio | $OI_{up}$ | $OI_{manual}$ | ratio |
|---|---|---|---|---|---|---|
| 2mm | $N_iN_k + N_kN_j$ $+N_jN_l + N_iN_l$ | $N_iN_jN_k$ $+N_iN_jN_l$ | – | $\sqrt{S}$ | $\sqrt{S}$ | 1 ✓ |
| 3mm | $N_iN_k + N_kN_j$ $+N_jN_m + N_mN_l$ | $N_iN_jN_k + N_jN_lN_m$ $+N_iN_jN_l$ | – | $\sqrt{S}$ | $\sqrt{S}$ | 1 ✓ |
| cholesky | $\frac12 N^2$ | $\frac13 N^3$ | $\frac23 N$ | $2\sqrt{S}$ | $\sqrt{S}$ | 2 |
| correlation | $MN$ | $M^2N$ | $M$ | $2\sqrt{S}$ | $\sqrt{S}$ | 2 |
| covariance | $MN$ | $M^2N$ | $M$ | $2\sqrt{S}$ | $\sqrt{S}$ | 2 |
| doitgen | $N_pN_qN_r$ | $2N_qn_rN_p^2$ | $2N_p$ | $\sqrt{S}$ | $\sqrt{S}$ | 1 ✓ |
| fdtd-2d | $3N_xN_y$ | $11N_xN_yT$ | $\frac{11}{3}T$ | $22\sqrt2\sqrt{S}$ | $\frac{11}{24}\sqrt3\sqrt{S}$ | $\frac{48\sqrt2}{\sqrt3}$ |
| floyd-warshall | $N^2$ | $2N^3$ | $2N$ | $2\sqrt{S}$ | $\sqrt{S}$ | 2 |
| gemm | $N_iN_j + N_jN_k + N_iN_k$ | $2N_iN_jN_k$ | – | $\sqrt{S}$ | $\sqrt{S}$ | 1 ✓ |
| heat-3d | $N^3$ | $30N^3T$ | $30T$ | $\frac{160}{3\sqrt[3]{3}}\sqrt[3]{S}$ | $\frac52\sqrt[3]{S}$ | $\frac{64}{3\sqrt[3]{3}}$ |
| jacobi-1d | $N$ | $6NT$ | $6T$ | $24S$ | $\frac32 S$ | 16 |
| jacobi-2d | $N^2$ | $10N^2T$ | $10T$ | $15\sqrt3\sqrt{S}$ | $\frac54\sqrt{S}$ | $12\sqrt3$ |
| lu | $N^2$ | $\frac23N^3$ | $\frac23N$ | $\sqrt{S}$ | $\sqrt{S}$ | 1 ✓ |
| ludcmp | $N^2$ | $\frac23N^3$ | $\frac23N$ | $\sqrt{S}$ | $\sqrt{S}$ | 1 ✓ |
| seidel-2d | $N^2$ | $9N^2T$ | $9T$ | $\frac{27\sqrt3}{2}\sqrt{S}$ | $\frac94\sqrt{S}$ | $6\sqrt3$ |
| symm | $\frac12M^2 + 2MN$ | $2M^2N$ | – | $\sqrt{S}$ | $\sqrt{S}$ | 1 ✓ |
| syr2k | $\frac12N^2 + 2MN$ | $2MN^2$ | – | $2\sqrt{S}$ | $\sqrt{S}$ | 2 |
| syrk | $\frac12N^2 + MN$ | $MN^2$ | – | $2\sqrt{S}$ | $\sqrt{S}$ | 2 |
| trmm | $\frac12M^2 + MN$ | $M^2N$ | – | $\sqrt{S}$ | $\sqrt{S}$ | 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 | $\frac12N^2$ | $N^2$ | 2 | 2 | 2 | 1 ✓ |
| adi | $N^2$ | $30N^2T$ | $30T$ | 30 | 5 | 6 |
| durbin | $N$ | $2N^2$ | $2N$ | 4 | $\frac23$ | 6 |
| gramschmidt | $MN$ | $2MN^2$ | $2N$ | $2\sqrt{S}$ | 1 | $2\sqrt{S}$ |
| nussinov | $\frac12N^2$ | $\frac13N^3$ | $\frac23N$ | $2\sqrt{S}$ | 1 | $2\sqrt{S}$ |

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 )

| kernel | # input data | #ops | ratio | $OI_{up}$ | $OI_{manual}$ | ratio |
|---|---|---|---|---|---|---|
| 2mm | $N_iN_k + N_kN_j$ $+N_jN_l + N_iN_l$ | $N_iN_jN_k$ $+N_iN_jN_l$ | – | $\sqrt{S}$ | $\sqrt{S}$ | 1 ✓ |
| 3mm | $N_iN_k + N_kN_j$ $+N_iN_m + N_mN_l$ | $N_iN_jN_k + N_jN_lN_m$ $+N_iN_jN_l$ | – | $\sqrt{S}$ | $\sqrt{S}$ | 1 ✓ |
| cholesky | $\frac12 N^2$ | $\frac13 N^3$ | $\frac23 N$ | $2\sqrt{S}$ | $\sqrt{S}$ | 2 |
| correlation | $MN$ | $M^2N$ | $M$ | $2\sqrt{S}$ | $\sqrt{S}$ | 2 |
| covariance | $MN$ | $M^2N$ | $M$ | $2\sqrt{S}$ | $\sqrt{S}$ | 2 |
| doitgen | $N_pN_qN_r$ | $2N_qn_rN_p^2$ | $2N_p$ | $\sqrt{S}$ | $\sqrt{S}$ | 1 ✓ |
| fdtd-2d | $3N_xN_y$ | $11N_xN_yT$ | $\frac{11}{3}T$ | $22\sqrt2\sqrt{S}$ | $\frac{11}{24}\sqrt3\sqrt{S}$ | $\frac{48\sqrt2}{\sqrt3}$ |
| floyd-warshall | $N^2$ | $2N^3$ | $2N$ | $2\sqrt{S}$ | $\sqrt{S}$ | 2 |
| gemm | $N_iN_j + N_jN_k + N_iN_k$ | $2N_iN_jN_k$ | – | $\sqrt{S}$ | $\sqrt{S}$ | 1 ✓ |
| heat-3d | $N^3$ | $30N^3T$ | $30T$ | $\frac{160}{3\sqrt3}\sqrt[3]{S}$ | $\frac{5}{2}\sqrt[3]{S}$ | $\frac{64}{3\sqrt3}$ |
| jacobi-1d | $N$ | $6NT$ | $6T$ | $24S$ | $\frac32 S$ | 16 |
| jacobi-2d | $N^2$ | $10N^2T$ | $10T$ | $15\sqrt3\sqrt{S}$ | $\frac54\sqrt{S}$ | $12\sqrt3$ |
| lu | $N^2$ | $\frac23N^3$ | $\frac23 N$ | $\sqrt{S}$ | $\sqrt{S}$ | 1 ✓ |
| ludcmp | $N^2$ | $\frac23N^3$ | $\frac23 N$ | $\sqrt{S}$ | $\sqrt{S}$ | 1 ✓ |
| seidel-2d | $N^2$ | $9N^2T$ | $9T$ | $\frac{27\sqrt3}{2}\sqrt{S}$ | $\frac94\sqrt{S}$ | $6\sqrt3$ |
| symm | $\frac12M^2 + 2MN$ | $2M^2N$ | – | $\sqrt{S}$ | $\sqrt{S}$ | 1 ✓ |
| syr2k | $\frac12N^2 + 2MN$ | $2MN^2$ | – | $2\sqrt{S}$ | $\sqrt{S}$ | 2 |
| syrk | $\frac12N^2 + MN$ | $MN^2$ | – | $2\sqrt{S}$ | $\sqrt{S}$ | 2 |
| trmm | $\frac12M^2 + MN$ | $M^2N$ | – | $\sqrt{S}$ | $\sqrt{S}$ | 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 | $\frac12N^2$ | $N^2$ | 2 | 2 | 2 | 1 ✓ |
| adi | $N^2$ | $30N^2T$ | $30T$ | 30 | 5 | 6 |
| durbin | $N$ | $2N^2$ | $2N$ | 4 | $\frac23$ | 6 |
| gramschmidt | $MN$ | $2MN^2$ | $2N$ | $2\sqrt{S}$ | 1 | $2\sqrt{S}$ |
| nussinov | $\frac12N^2$ | $\frac13N^3$ | $\frac23N$ | $2\sqrt{S}$ | 1 | $2\sqrt{S}$ |

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 )

(1) Olivier Beaumont, Lionel Eyraud-Dubois, Julien Langou, and Mathieu Vérité. I/O-optimal algorithms for symmetric linear algebra kernels. In the 34th ACM Symposium on Parallelism in Algorithms and Architectures (SPAA '22), Philadelphia, PA, USA, July 11–14, 2022. DOI information: 10.1145/3490148.3538587.

(2) Olivier Beaumont, Philippe Duchon, Lionel Eyraud-Dubois, Julien Langou, and Mathieu Vérité. Symmetric Block-Cyclic Distribution: Fe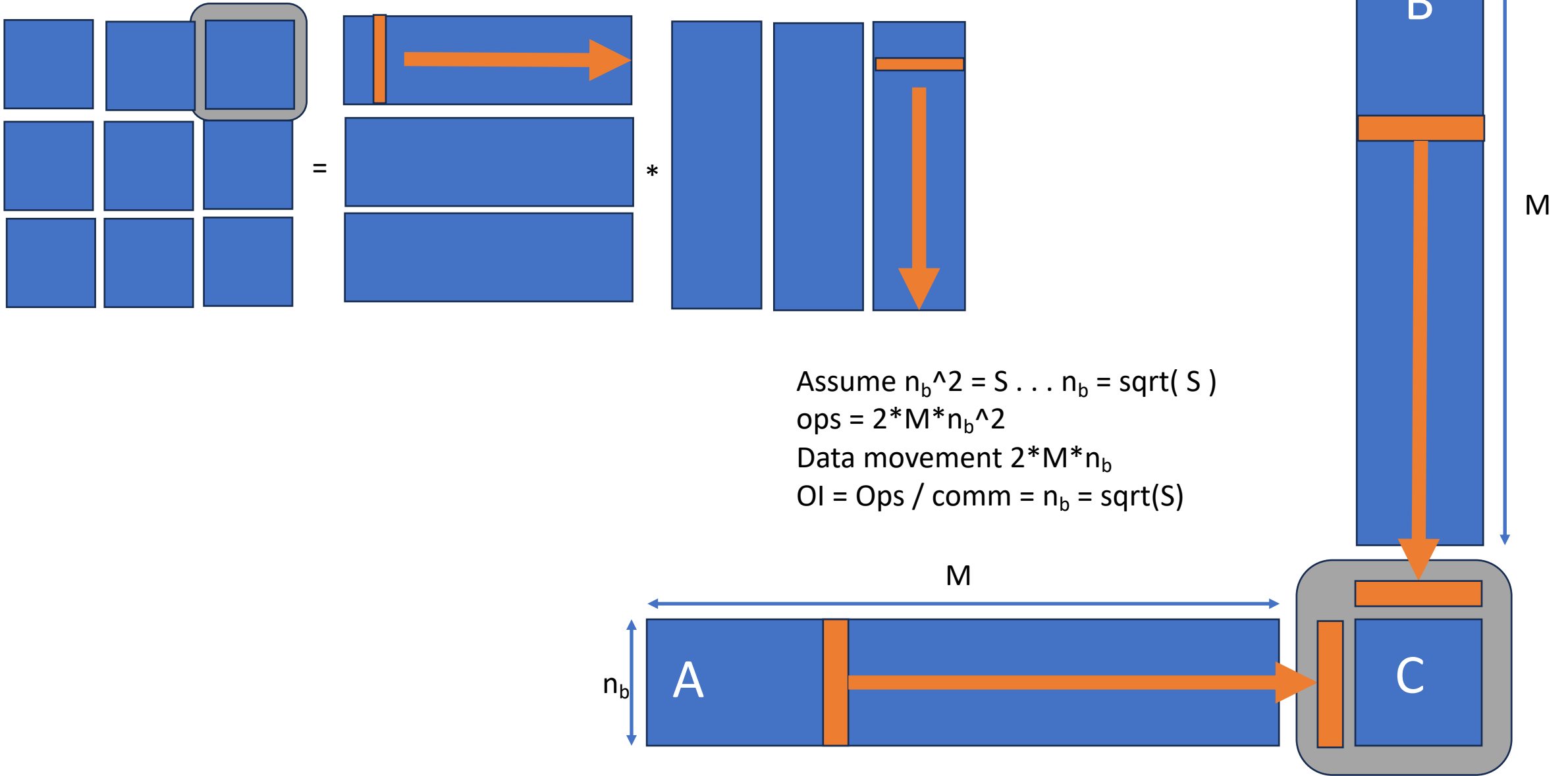wer Communications leads to Faster Dense Cholesky Factorization. In ACM/IEEE SC 2022 Conference (SC'22), Dallas, TX, USA, November 13-18, 2022. DOI information: 10.1109/SC41404.2022.00034. Nominee for Best Paper Award.

(3) Emmanuel Agullo, Alfredo Buttari, Olivier Coulaud, Lionel Eyraud-Dubois, Mathieu Faverge, Alain Franc, Abdou Guermouche, Antoine Jego, Romain Peressoni, and Florent Pruvost. On the Arithmetic Intensity of Distributed-Memory Dense Matrix Multiplication Involving a Symmetric Input Matrix (SYMM). In 2023 IEEE International Parallel and Distributed Processing Symposium (IPDPS'23), St. Petersburg, FL, USA, 15-19 May 2023. DOI information: 10.1109/IPDPS54959.2023.00044
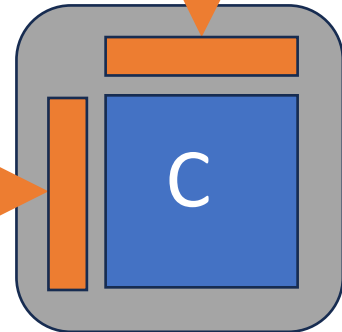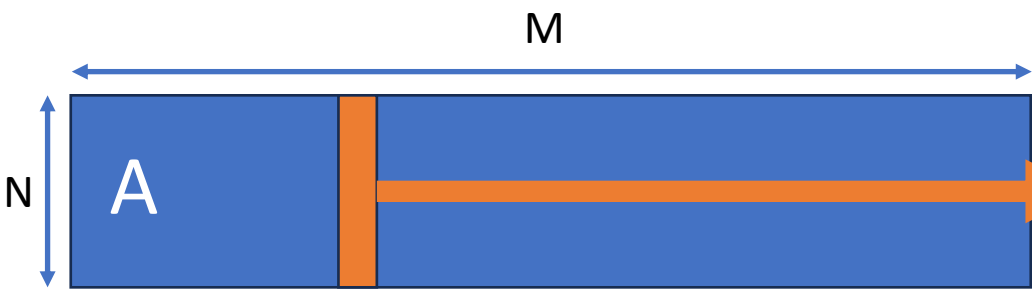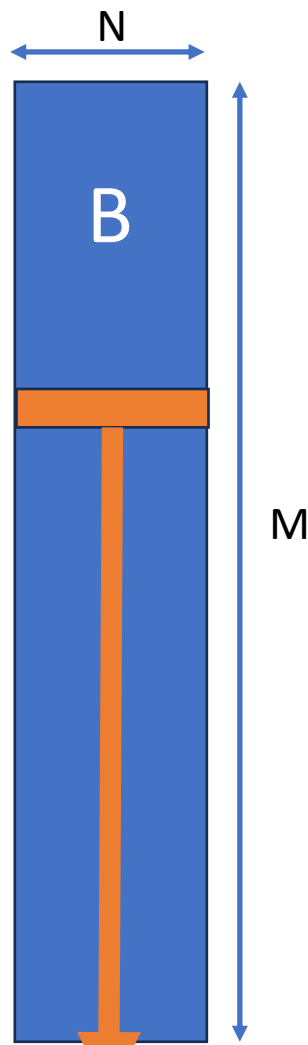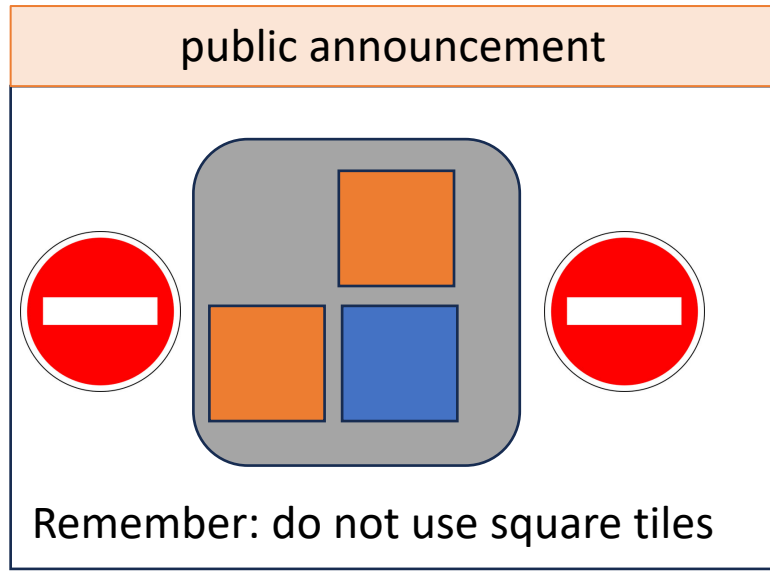
Operational intensity of LU is at sqrt(2) * sqrt(S)

# GEMM



Assume $n_b^2 = S \ldots n_b = \sqrt{S}$

ops $= 2*M*n_b^2$

Data movement $2*M*n_b$

OI $=$ Ops $/$ comm $= n_b = \sqrt{S}$

# GEMM



public announcement

Remember: do not use square tiles

# SYRK

Assume $(N^2)/2 = S \ldots N = \sqrt{S} * \sqrt{2}$
ops = $M*N^2$
Data movement $M*N$
OI = Ops / comm = $N = \sqrt{S}*\sqrt{2}$
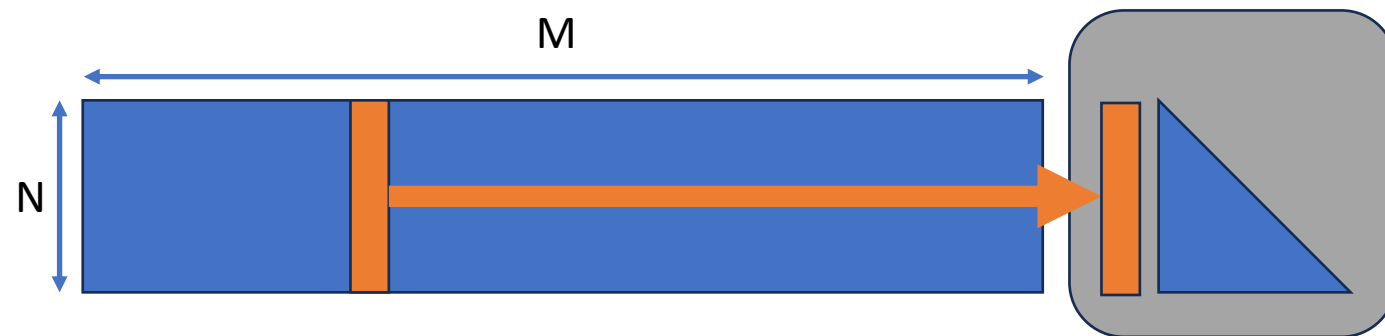
Observation:
SYRK has an OI sqrt(2) better
than MM when a "triangle fit in cache"

# SYRK



Assume $(N^2)/2 = S \ldots N = \sqrt{S} * \sqrt{2}$
ops = $M*N^2$
Data movement $M*N$
OI = Ops / comm = N = $\sqrt{S}*\sqrt{2}$

SYRK has an OI sqrt(2) better than MM when a "triangle fit in cache"

- Building an optimal GEMM is easy.
- But how to build an optimal SYRK?

# SYRK

- IOLB says OI is less than 2 sqrt(S)

- Béreux's algorithm has an OI of sqrt(S)

- This begs the question what is <u>the</u> OI of SYRK?

# SYRK

- IOLB says OI is less than 2 sqrt(S)

- Béreux's algorithm has an OI of sqrt(S)

- The OI upper bound from IOLB of 2 sqrt(S) is too optimistic, we can find a lower OI upper bound: sqrt(2) * sqrt(S).

- And we can improve Béreux's algorithm. Instead of an OI of sqrt(S), we can improve to sqrt(2) * sqrt(S).

| kernel | # input data | #ops | ratio | $OI_{up}$ | $OI_{manual}$ | ratio |
|---|---|---|---|---|---|---|
| 2mm | $N_iN_k + N_kN_j$ $+N_jN_l + N_iN_l$ | $N_iN_jN_k$ $+N_iN_jN_l$ | – | $\sqrt{S}$ | $\sqrt{S}$ | 1 ✓ |
| 3mm | $N_iN_k + N_kN_j$ $+N_iN_m + N_mN_l$ | $N_iN_jN_k + N_jN_lN_m$ $+N_iN_jN_l$ | – | $\sqrt{S}$ | $\sqrt{S}$ | 1 ✓ |
| cholesky | $\frac{1}{2}N^2$ | $\frac{1}{3}N^3$ | $\frac{2}{3}N$ | $2\sqrt{S}$ | $\sqrt{S}$ | 2 |
| correlation | $MN$ | $M^2N$ | $M$ | $2\sqrt{S}$ | $\sqrt{S}$ | 2 |
| covariance | $MN$ | $M^2N$ | $M$ | $2\sqrt{S}$ | $\sqrt{S}$ | 2 |
| doitgen | $N_pN_qN_r$ | $2N_qn_rN_p^2$ | $2N_p$ | $\sqrt{S}$ | $\sqrt{S}$ | 1 ✓ |
| fdtd-2d | $3N_xN_y$ | $11N_xN_yT$ | $\frac{11}{3}T$ | $22\sqrt{2}\sqrt{S}$ | $\frac{11}{24}\sqrt{3}\sqrt{S}$ | $\frac{48\sqrt{2}}{\sqrt{3}}$ |
| floyd-warshall | $N^2$ | $2N^3$ | $2N$ | $2\sqrt{S}$ | $\sqrt{S}$ | 2 |
| gemm | $N_iN_j + N_jN_k + N_iN_k$ | $2N_iN_jN_k$ | – | $\sqrt{S}$ | $\sqrt{S}$ | 1 ✓ |
| heat-3d | $N^3$ | $30N^3T$ | $30T$ | $\frac{160}{3\sqrt[3]{3}}\sqrt[3]{S}$ | $\frac{5}{2}\sqrt[3]{S}$ | $\frac{64}{3\sqrt[3]{3}}$ |
| jacobi-1d | $N$ | $6NT$ | $6T$ | $24S$ | $\frac{3}{2}S$ | 16 |
| jacobi-2d | $N^2$ | $10N^2T$ | $10T$ | $15\sqrt{3}\sqrt{S}$ | $\frac{5}{4}\sqrt{S}$ | $12\sqrt{3}$ |
| 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 | $N^2$ | $9N^2T$ | $9T$ | $\frac{27\sqrt{3}}{2}\sqrt{S}$ | $\frac{9}{4}\sqrt{S}$ | $6\sqrt{3}$ |
| symm | $\frac{1}{2}M^2 + 2MN$ | $2M^2N$ | – | $\sqrt{S}$ | $\sqrt{S}$ | 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 | $\frac{1}{2}M^2 + MN$ | $M^2N$ | – | $\sqrt{S}$ | $\sqrt{S}$ | 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 | $N^2$ | $30N^2T$ | $30T$ | 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 | $2\sqrt{S}$ |

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 )
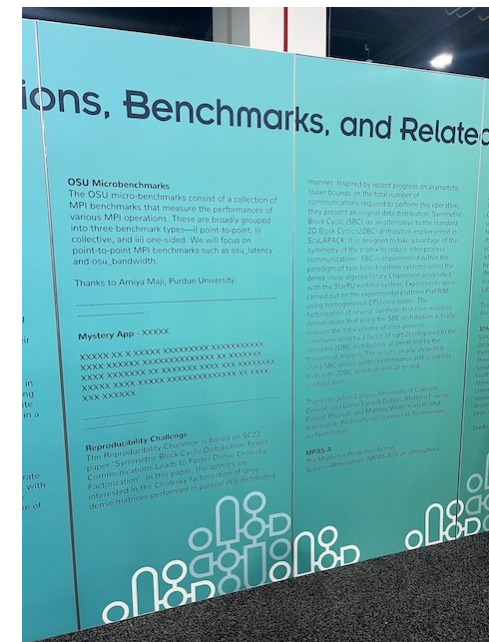
(1) Olivier Beaumont, Lionel Eyraud-Dubois, Julien Langou, and Mathieu Vérité. I/O-optimal algorithms for symmetric linear algebra kernels. In the 34th ACM Symposium on Parallelism in Algorithms and Architectures (SPAA '22), Philadelphia, PA, USA, July 11–14, 2022. DOI information: 10.1145/3490148.3538587.

(2) Olivier Beaumont, Philippe Duchon, Lionel Eyraud-Dubois, Julien Langou, and Mathieu Vérité. Symmetric Block-Cyclic Distribution: Fewer Communications leads to Faster Dense Cholesky Factorization. In ACM/IEEE SC 2022 Conference (SC'22), Dallas, TX, USA, November 13-18, 2022. DOI information: 10.1109/SC41404.2022.00034. Nominee for Best Paper Award.

(3) Emmanuel Agullo, Alfredo Buttari, Olivier Coulaud, Lionel Eyraud-Dubois, Mathieu Faverge, Alain Franc, Abdou Guermouche, Antoine Jego, Romain Peressoni, and Florent Pruvost. On the Arithmetic Intensity of Distributed-Memory Dense Matrix Multiplication Involving a Symmetric Input Matrix (SYMM). In 2023 IEEE International Parallel and Distributed Processing Symposium (IPDPS'23), St. Petersburg, FL, USA, 15-19 May 2023. DOI information: 10.1109/IPDPS54959.2023.00044

# SIGHPC Reproducibility Award
# SCC @ SC'23







Thanks to Mathieu Faverge and Florent Pruvost
Thanks to Thomas Herault, George Bosilca, and Weslley Pereira

Olivier Beaumont, Philippe Duchon, Lionel Eyraud-Dubois, Julien Langou, and Mathieu Vérité.
Symmetric Block-Cyclic Distribution: Fewer Communications leads to Faster Dense Cholesky Factorization.
In ACM/IEEE SC 2022 Conference (SC'22), Dallas, TX, USA, November 13-18, 2022. DOI information:
10.1109/SC41404.2022.00034. Nominee for Best Paper Award.

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

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 )

## Classical Gram-Schmidt (CGS)

**Input:** $a_1$, $a_2$,…, $a_{n-1}$, $a_n$

**Output:** $q_1$, $q_2$,…, $q_{n-1}$, $q_n$

**for** j=1:n,

$\quad$ w = ( I - $Q_{1:j-1}$ $Q_{1:j-1}{}^{\mathsf{T}}$ ) $a_j$

$\quad$ $q_j$ = w / || w ||$_2$

**end**

## Modified Gram-Schmidt (MGS)

**Input:** $a_1$, $a_2$,…, $a_{n-1}$, $a_n$

**Output:** $q_1$, $q_2$,…, $q_{n-1}$, $q_n$

**for** j=1:n,

$\quad$ w = ( I - $q_{j-1}$ $q_{j-1}{}^{\mathsf{T}}$ ) .. ( I - $q_1$ $q_1{}^{\mathsf{T}}$ ) $a_j$

$\quad$ $q_j$ = w / || w ||$_2$

**end**

```matlab
function [ Q, R ] = cgs( A )

  n = size(A,2);

  Q = zeros(size(A));

  R = zeros (n);

  for j=1:n,

    Q(:,j) = A(:,j) ;

    for i=1:j-1,

        R(i,j) =  Q(:,i)' * A(:,j) ;
        Q(:,j)  = Q(:,j)  - Q(:,i) * R(i,j);

    end

    R(j,j) = norm( Q(:,j)  );

    Q(:,j) = Q(:,j)  / R(j,j);

  end

end
```

```matlab
function [ Q, R ] = mgs( A )

  n = size(A,2);

  Q = zeros(size(A));

  R = zeros (n);

  for j=1:n,

    Q(:,j) = A(:,j) ;

    for i=1:j-1,

        R(i,j) =  Q(:,i)' * Q(:,j) ;
        Q(:,j)  = Q(:,j)  - Q(:,i) * R(i,j);

    end

    R(j,j) = norm( Q(:,j)  );

    Q(:,j) = Q(:,j)  / R(j,j);

  end

end
```

```c
#include <math.h>
void qr_cgs_ll (int M, int N, double A[M][N], double R[N][N] )
{
int i, j, k;
#pragma scop
if (M>=N) {
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];
}
}
#pragma endscop
}
```
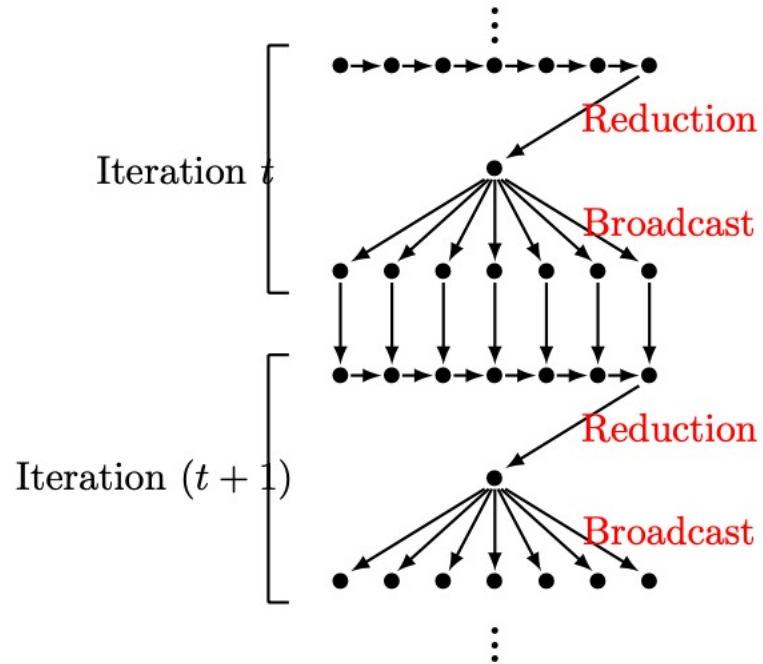
```c
#include <math.h>
void qr_mgs_ll (int M, int N, double A[M][N], double R[N][N] )
{
int i, j, k;
#pragma scop
if (M>=N) {
for (j = 0; j < N; j++) {
    for (i = 0; i < j; i++) {
        R[i][j] = 0.0e+00;
        for (k = 0; k < M; k++)
            R[i][j] += A[k][i] * A[k][j];
        for (k = 0; 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];
}
}
#pragma endscop
}
```

CGS

MN input data (+ N^2 / 2 for R in output)

2 MN^2 operations

IOLB: operational intensity is at most $2\sqrt{S}$

MGS

MN data (+ N^2 / 2 for R in output)

2 MN^2 operations

IOLB: operational intensity is at most $2\sqrt{S}$

# Hourglass pattern



Standard "MM-way" of doing business

$$|E| \leq |\phi_{i,j}(E)|^{\frac{1}{2}} \cdot |\phi_{i,k}(E)|^{\frac{1}{2}} \cdot |\phi_{k,j}(E)|^{\frac{1}{2}} \leq \boxed{K^{3/2}}$$

If Hourglass pattern is detected then

$$|I'| \leq |\phi_i(I')| \times |\phi_j(I')| \times |\phi_k(I')|.$$

$$|I'| \leq M \times \frac{K}{M} \times \frac{K}{M} = \boxed{\frac{K^2}{M}}$$

# Modified Gram-Schmidt

**Theorem 5** (Lower bounds for MGS). *The communication volume $Q$ for the MGS algorithm on a $M \times N$ matrix can be bounded as follows:*

$$\frac{M^2 N(N-1)}{8(S+M)} \leq Q$$

*Furthermore, if $S \leq M$, we also have:*

$$\frac{(M-S)N(N-1)}{4} \leq Q$$

Interpretation: (with approximations)

If M < N then IO is at least     $(1/8) * M^2 N^2 / S$
or in another words              $(1/8) * ( M / S ) * M N^2$

If M > N then IO is at least     $(1/4) * M N^2$

Interpretation:
Since M < sqrt(S), $M^2 N^2 / S$ is a much better lower bound on IO than $M N^2 / $ sqrt(S) – (larger is better)

| Kernel | Old bound | New bound (hourglass) | |
|---|---|---|---|
| MGS | $\Omega\left(\dfrac{MN^2}{\sqrt{S}}\right)$ | $\Omega\left(\dfrac{M^2N(N-1)}{S+M}\right)$ | |
| QR HH A2V | $\Omega\left(\dfrac{MN^2}{\sqrt{S}}\right)$ | $\Omega\left(\dfrac{MN^2(N-M)}{N-M-S}\right)$ | GEQR2 |
| QR HH V2Q | $\Omega\left(\dfrac{MN^2}{\sqrt{S}}\right)$ | $\Omega\left(\dfrac{MN^2(N-M)}{N-M-S}\right)$ | ORG2R |

| Kernel | Old bound [17] | New bound (hourglass) |
|---|---|---|
| MGS | $\dfrac{2M+3MN+MN^2}{\sqrt{S}}+5M-MN+\dfrac{7N-N^2}{2}-S-6$ | $\dfrac{N^2M^2+2M^2-3NM^2}{8(M+S)}+5M-MN+\dfrac{7N-N^2}{2}-S-6$ |
| QR HH A2V | $\dfrac{3MN^2+6M+7N-N^3-9MN-6}{3\sqrt{S}}+5M-MN+5N-S-13$ | $\dfrac{3MN^2-9MN+7N+6M-6-N^3}{24\left(1-\frac{S}{N-M}\right)}+5M-MN+5N-S-13$ |
| QR HH V2Q | $\dfrac{3MN^2-N^3+6M+7N-9MN-6}{3\sqrt{S}}+2M+2N+\dfrac{N-N^2}{2}-S-4$ | $\dfrac{3MN^2-N^3+6M+7N-9MN-6}{24\left(1+\frac{S}{M-N}\right)}+2M+2N+\dfrac{N-N^2}{2}-S-4$ |

| Kernel | Old bound [17] | New bound (hourglass) |
|---|---|---|
| MGS | $\Omega\left(\dfrac{MN^2}{\sqrt{S}}\right)$ | $\Omega\left(\dfrac{M^2N(N-1)}{S+M}\right)$ |
| QR HH A2V | $\Omega\left(\dfrac{MN^2}{\sqrt{S}}\right)$ | $\Omega\left(\dfrac{MN^2(N-M)}{N-M-S}\right)$ |
| QR HH V2Q | $\Omega\left(\dfrac{MN^2}{\sqrt{S}}\right)$ | $\Omega\left(\dfrac{MN^2(N-M)}{N-M-S}\right)$ |
| GEBD2 | $\Omega\left(\dfrac{MN^2}{\sqrt{S}}\right)$ | $\Omega\left(\dfrac{MN^2(M-N+1)}{8(S+M-N+1)}\right)$ |
| GEHD2 | $\Omega\left(\dfrac{N^3}{\sqrt{S}}\right)$ | $\Omega\left(\dfrac{N^4}{N+2S}\right)$ |

| Kernel | Old bound [17] | New bound (hourglass) |
|---|---|---|
| MGS | $\dfrac{2M+3MN+MN^2}{\sqrt{S}} + 5M - MN + \dfrac{7N-N^2}{2} - S - 6$ | $\dfrac{N^2M^2+2M^2-3NM^2}{8(M+S)} + 5M - MN + \dfrac{7N-N^2}{2} - S - 6$ |
| QR HH A2V | $\dfrac{3MN^2+6M+7N-N^3-9MN-6}{3\sqrt{S}} + 5M - MN + 5N - S - 13$ | $\dfrac{3MN^2-9MN+7N+6M-6-N^3}{24(1-\frac{S}{N-M})} + 5M - MN + 5N - S - 13$ |
| QR HH V2Q | $\dfrac{3MN^2-N^3+6M+7N-9MN-6}{3\sqrt{S}} + 2M + 2N + \dfrac{N-N^2}{2} - S - 4$ | $\dfrac{3MN^2-N^3+6M+7N-9MN-6}{24(1+\frac{S}{M-N})} + 2M + 2N + \dfrac{N-N^2}{2} - S - 4$ |
| GEBD2 | $\dfrac{3MN^2-N^3-9MN+6M+7N-6}{3\sqrt{S}} + 5N + 5M - MN - S - 13$ | $\dfrac{3MN^2-N^3+3N^2-15MN+4N+18M-12}{24(1+\frac{S}{1+M-N})} + 5N + 7M - MN - S - 18$ |
| GEHD2 | $\dfrac{5N^3-30N^2+55N-30}{3\sqrt{S}} + \dfrac{69N-9N^2}{2} - 3*S - 56$ | $\dfrac{N^3-6N^2+11N-6}{12(1+\frac{S}{N-M-1})} - N^2 + 12N - S - 19$ |

./iolb-affine-indocker.sh examples/lin-alg/gehd2_splitted.c

$-1+n^2+\max(0,-107+5/24*(S-(2+m-n)^{\wedge}(-1)*S^{\wedge}2)^{\wedge}(-1)*S*n^{\wedge}3-11/12*((1+m-n)^{\wedge}(-1)*S^{\wedge}2-S)^{\wedge}(-1)*S*n-8*m*S^{\wedge}(-1/2)*n^{\wedge}2-5/4*(S-(2+m-n)^{\wedge}(-1)*S^{\wedge}2)^{\wedge}(-1)*S*n^{\wedge}2-23*S^{\wedge}(-1/2)-27*m*S^{\wedge}(-1/2)-1/12*((1+m-n)^{\wedge}(-1)*S^{\wedge}2-S)^{\wedge}(-1)*S*n^{\wedge}3+55/24*(S-(2+m-n)^{\wedge}(-1)*S^{\wedge}2)^{\wedge}(-1)*S*n-2*m^{\wedge}2-5/4*(S-(2+m-n)^{\wedge}(-1)*S^{\wedge}2)^{\wedge}(-1)*S+1/2*((1+m-n)^{\wedge}(-1)*S^{\wedge}2-S)^{\wedge}(-1)*S*n^{\wedge}2+31*m*S^{\wedge}(-1/2)*n-9*m+235/6*S^{\wedge}(-1/2)*n+7*m^{\wedge}2*S^{\wedge}(-1/2)*n-13*m^{\wedge}2*S^{\wedge}(-1/2)+5*m*n-25/2*n^{\wedge}2-41/2*S^{\wedge}(-1/2)*n^{\wedge}2-10*S+145/2*n+1/2*((1+m-n)^{\wedge}(-1)*S^{\wedge}2-S)^{\wedge}(-1)*S+10/3*S^{\wedge}(-1/2)*n^{\wedge}3-2*m^{\wedge}3*S^{\wedge}(-1/2))$

# Modified Gram-Schmidt

**Theorem 5** (Lower bounds for MGS). *The communication volume $Q$ for the MGS algorithm on a $M \times N$ matrix can be bounded as follows:*

$$\frac{M^2 N(N-1)}{8(S+M)} \leq Q$$

*Furthermore, if $S \leq M$, we also have:*

$$\frac{(M-S)N(N-1)}{4} \leq Q$$

Interpretation:

If M < N then IO is at least (1/8) * $M^2 N^2$ / S or (1/8) * ( M / S ) M $N^2$

If M > N then IO is at least (1/4) * M $N^2$

Interpretation:

Much better lower bound on IO than M $N^2$ / sqrt(S)

Communication-optimal parallel and sequential QR and LU factorizations
Demmel, Grigori, Hoemmen and Langou, 2008
https://arxiv.org/abs/0808.2664



**N**

**K**

**M**

Finished part of Q    panel    trailing matrix

Modified Gram-Schmidt, GEQR2, ORGQR2

We need at least a few columns to fit in cache!

Take K such that

$$M * K + M < S$$

panel fits in cache and stays in cache

At start of step, panel+trailing are updated and ready to go
Load panel
Do MGS on panel (panel factorization)
Update trailing matrix by loading one column at a time

Analysis:
K is about M / S
Load about (less than) N columns of size M every N / K steps
So IO is about $M^2 N^2 / S$

Which is a factor of S/M better than $M N^2$
So a factor of K better than $M N^2$

This is right looking, you can do a left looking variant too.

A2V (geqr2)
bora@plafrim
Intel CascadeLake

Variant
— Left Looking
— QR
— QR2
— Recursive
— Right Looking
— Tiled LL
— Tiled RL

A2V (geqr2)
bora@plafrim
Intel CascadeLake

Variant
Left Looking
QR
QR2
Recursive
Right Looking
Tiled LL
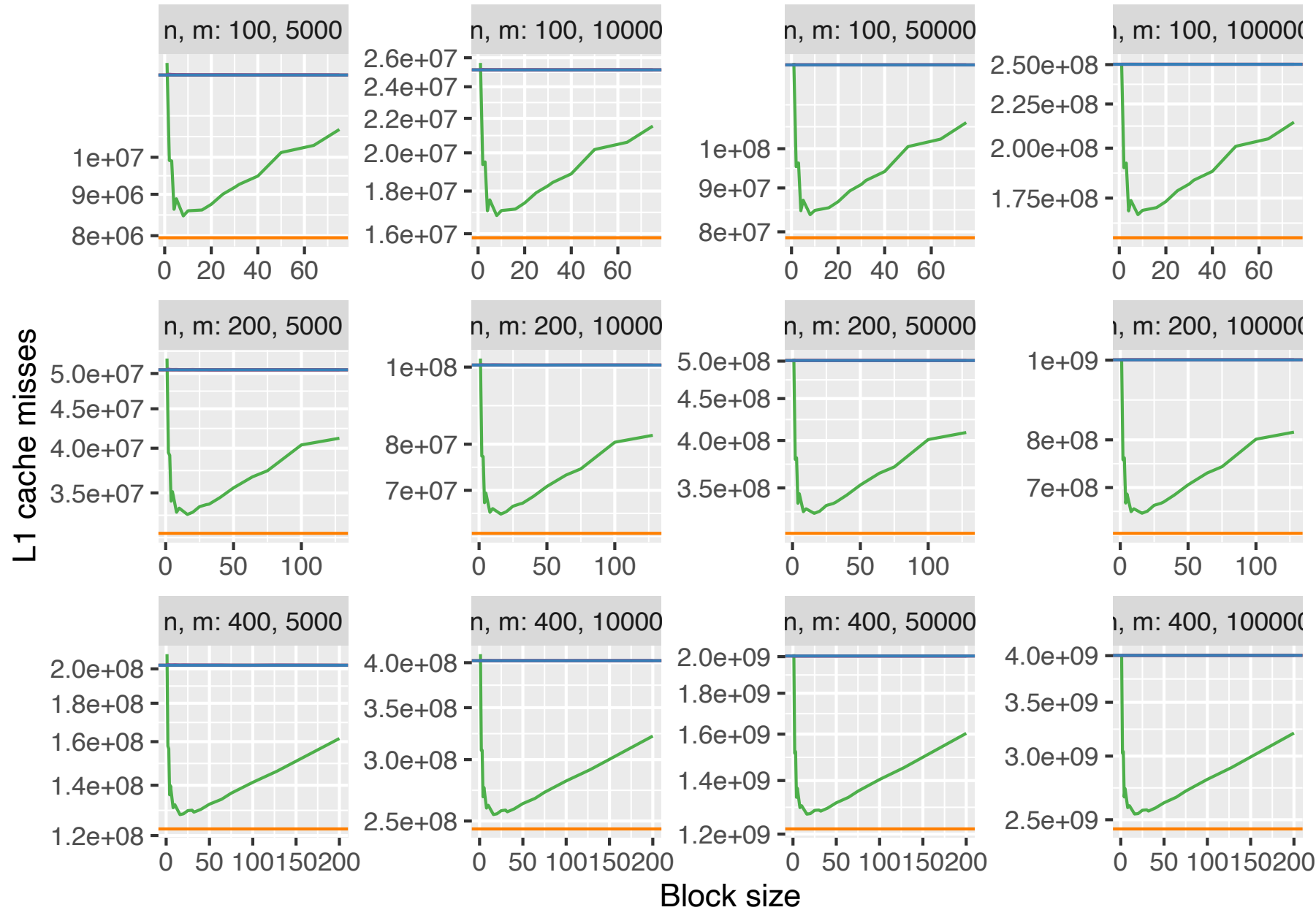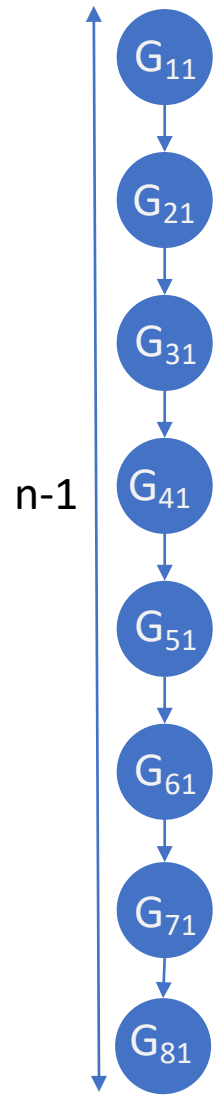Tiled RL

A2V (geqr2)
bora@plafrim
Intel CascadeLake

Variant: Left Looking, QR, QR2, Recursive, Right Looking, Tiled LL, Tiled RL

MGS
bora@plafrim
Intel CascadeLake

Variant — Left Looking — Recursive — Right Looking — Tiled LL — Tiled RL

MGS
bora@plafrim
Intel CascadeLake

L1 cache misses

Block size

Variant — Left Looking — Recursive — Right Looking — Tiled LL — Tiled RL

# Current work: chains of Givens rotations with Thijs Steel

```
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;
        }
    }
}
```

A Francis step

work on columns 1 and 2

work on columns 2 and 3

work on columns 3 and 4
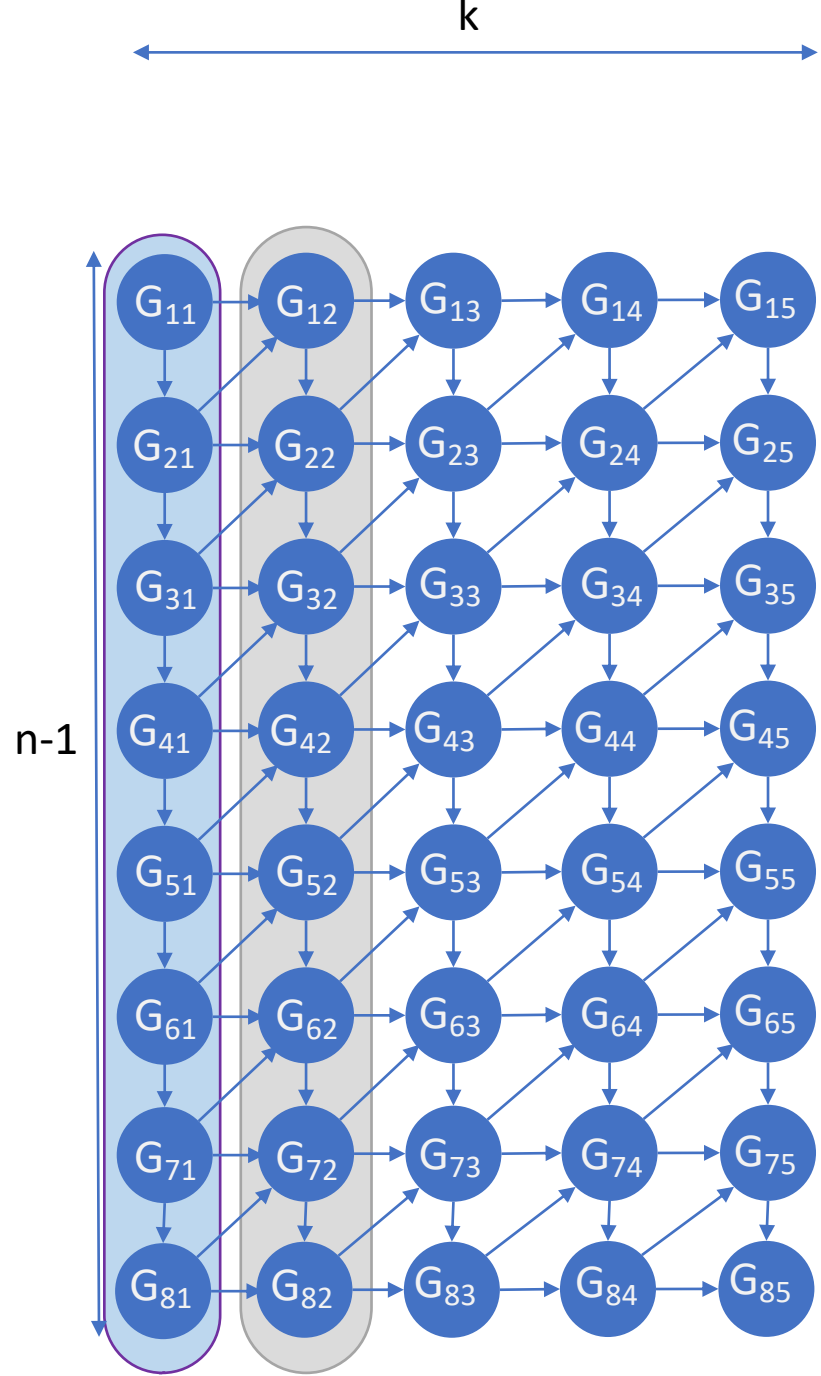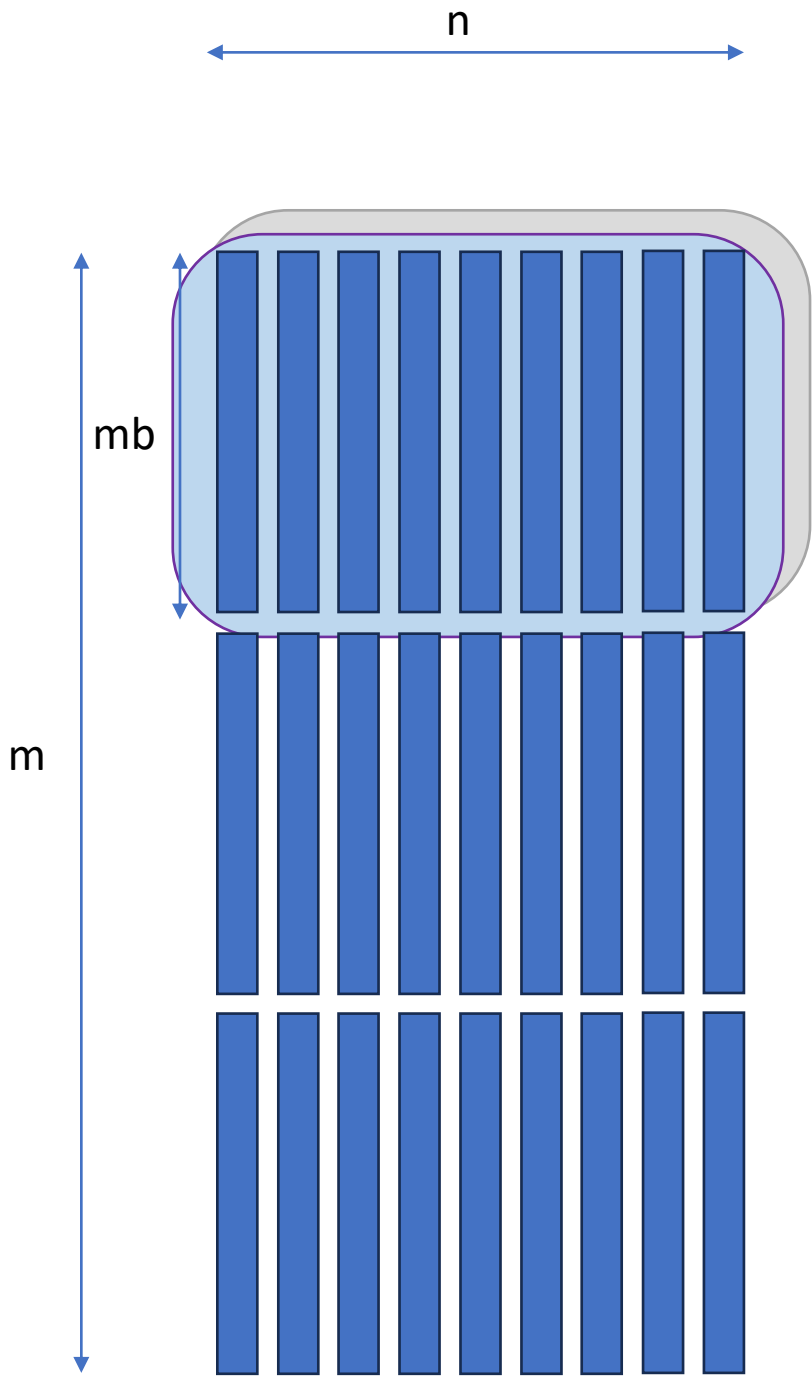
work on columns 4 and 5

work on columns 5 and 6

work on columns 6 and 7

work on columns 7 and 8

work on columns 8 and 9

n

k

A waves of five Francis steps

work on columns 1 and 2

work on columns 2 and 3
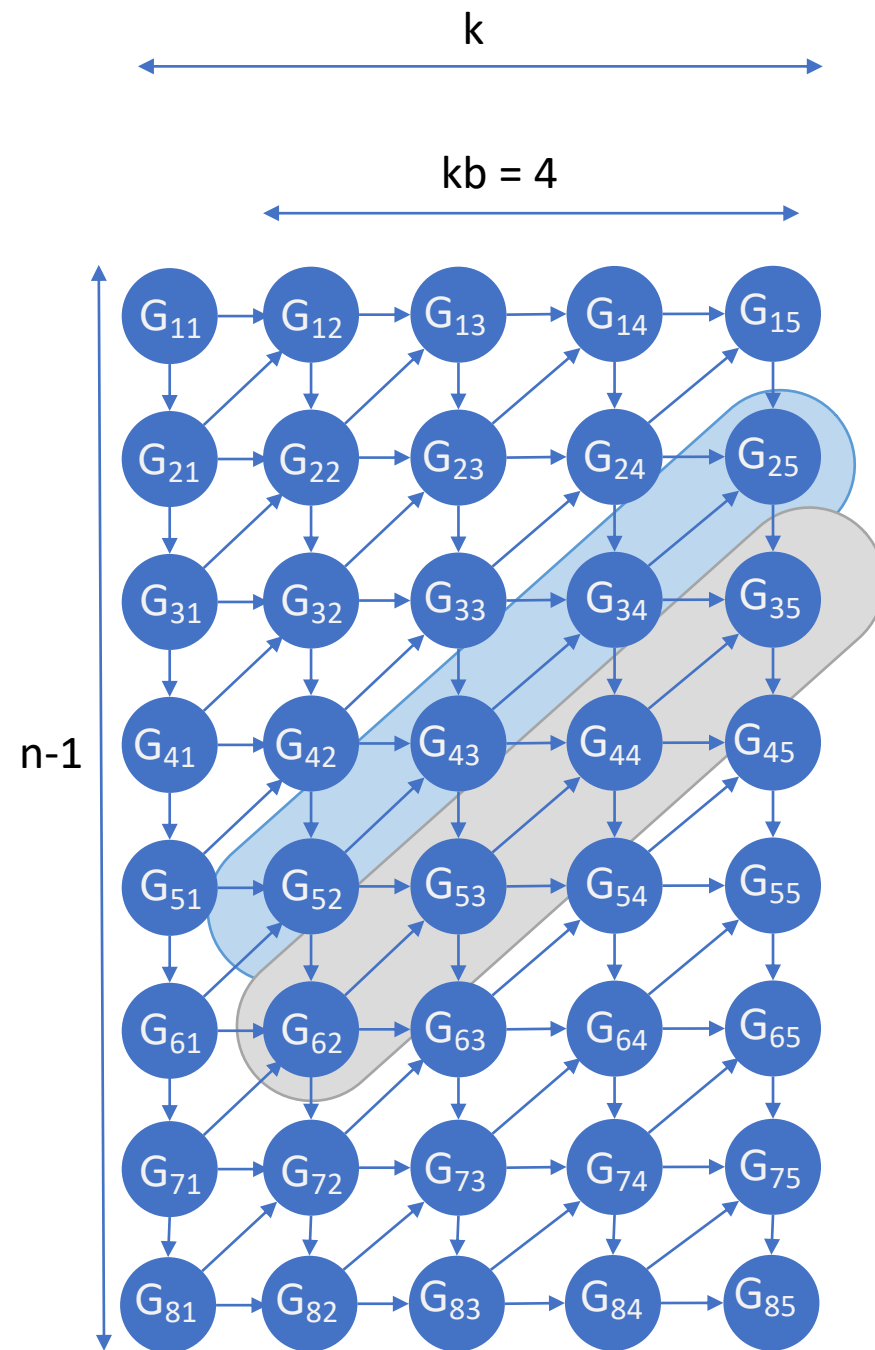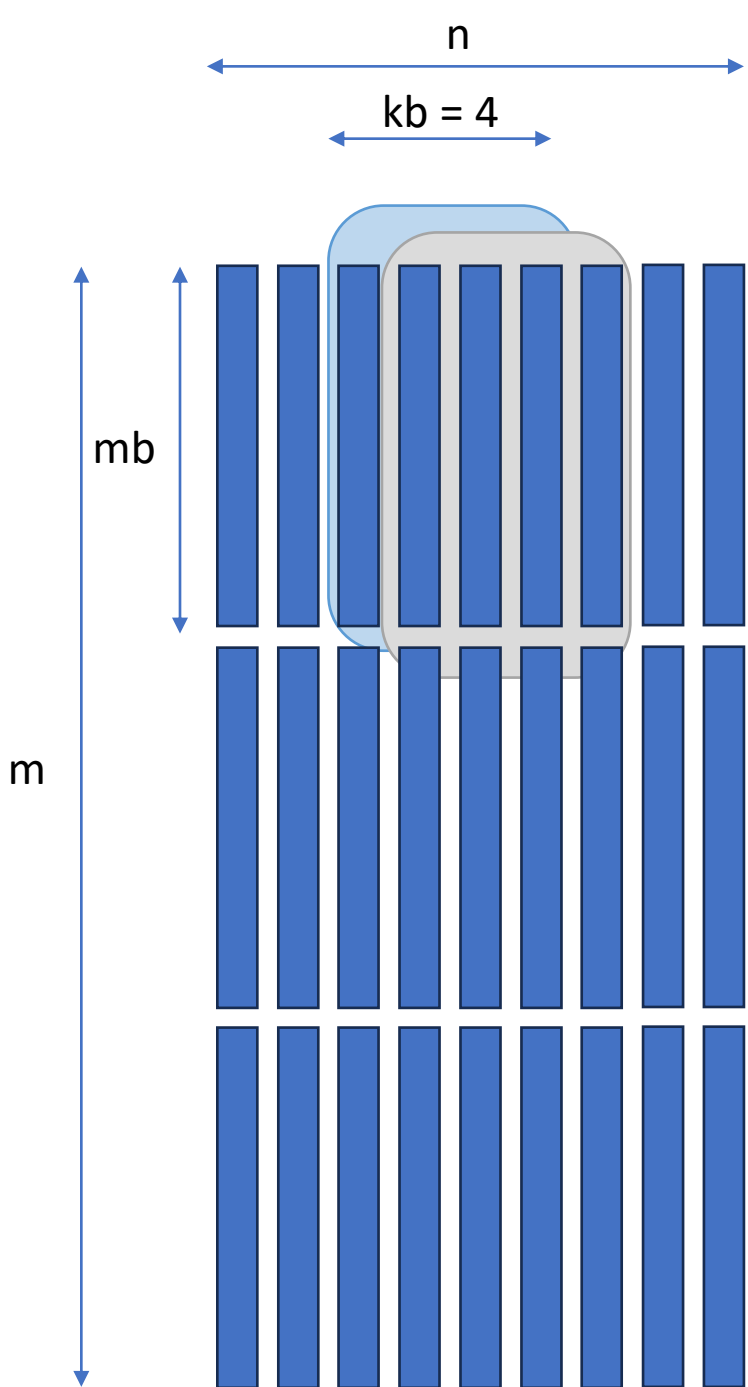
work on columns 3 and 4

work on columns 4 and 5

work on columns 5 and 6

work on columns 6 and 7

work on columns 7 and 8

work on columns 8 and 9

work on columns 1 and 2

work on columns 2 and 3

work on columns 3 and 4

work on columns 4 and 5

work on columns 5 and 6

work on columns 6 and 7

work on columns 7 and 8

work on columns 8 and 9

work on columns 1 and 2

work on columns 2 and 3

work on columns 3 and 4

work on columns 4 and 5

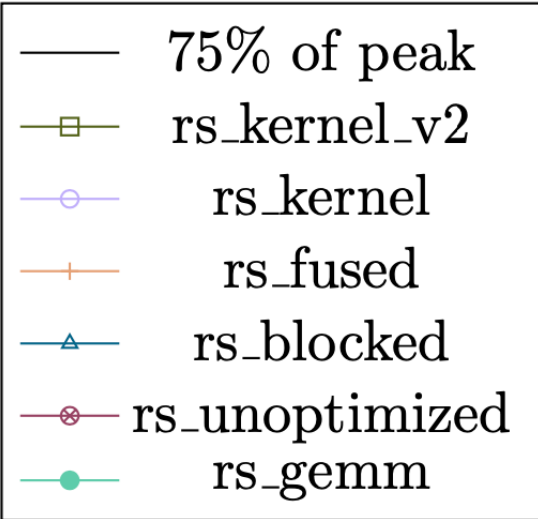work on columns 5 and 6

work on columns 6 and 7

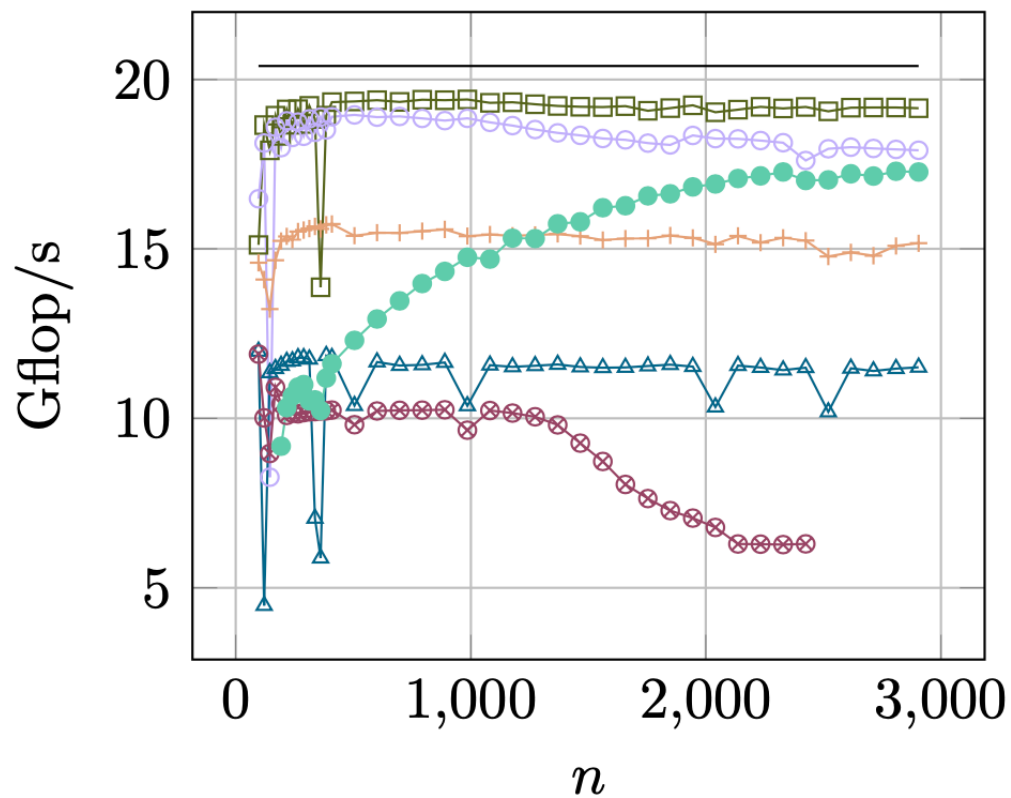work on columns 7 and 8

work on columns 8 and 9

# Analysis

- Going fast.
- We need to apply mnk Givens rotations.
- We assume that $k_b * m_b \leq M$
- We move on piece of column of size $m_b$, write back one piece of size $m_b$, and load $2k_b$ rotations
- And we will do $k_b * m_b$ rotations in a segment
- Forgetting start and end clean up code, we need to do
  $mnk / m_b / k_b$ segments
- And so the volume of communication of the algorithm is
  $( mnk / m_b / k_b ) * ( 2m_b + 2k_b )$
- For $k_b = sqrt(M)$ and $m_b = sqrt(M)$ we get
  $4 mnk / sqrt( M )$
- Since the total number of operations of our code is 6mnk, this means that the operational intensity for the wavefront algorithm is **(3/2) sqrt{S}**
- IOLB returns **6 sqrt{S}** as an upper bound for the operational intensity
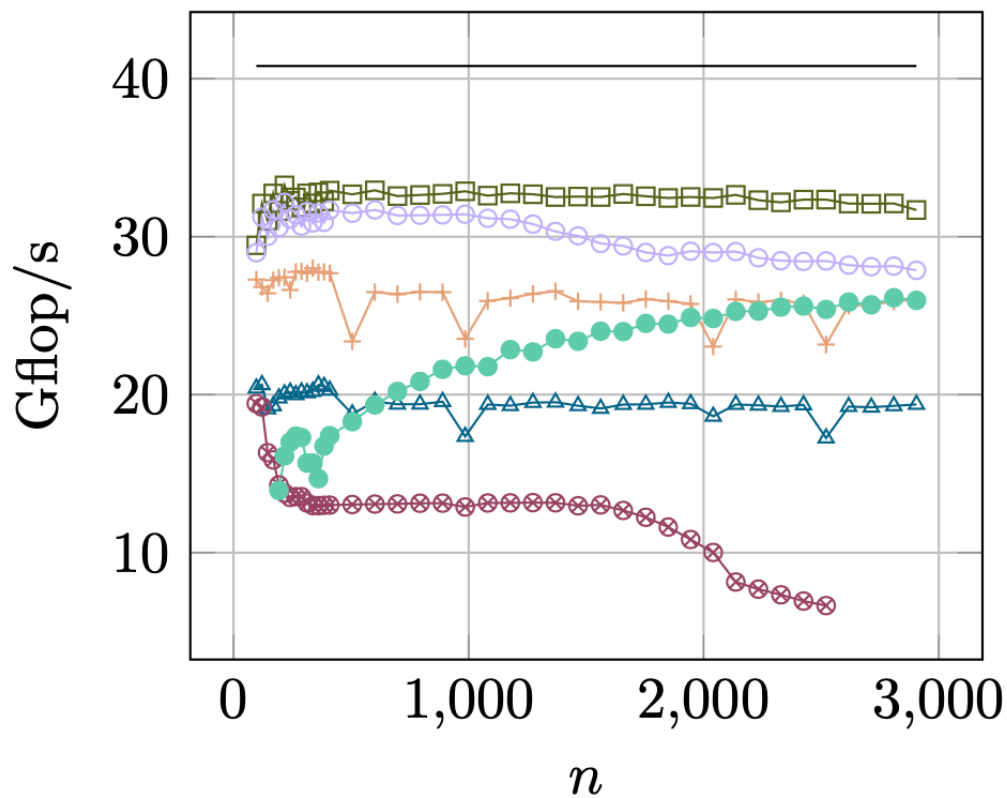- So we are a factor of 4 off.

apply_givens( )
k = 180, varying n and m = n

"fused" is a 2x2 kernel
"kernel" is a 16x2 kernel
"kernel_v2" uses packed format for A

(a) Xeon V2 flop rate

(b) Xeon V3 flop rate

## a2v – geqr2 – Household QR factorization

```c
for(k = 0; k < N; k++){
  norma2 = 0.e+00;
  for(i = k+1; i < M; i++){
    norma2 += A[i][k] * A[i][k];
  }
  norma =
  A[k][k]
  tau[k] =
```

## v2q – orqr2 – Construction of Q from Householder

```c
2. for(k = N-1; k > -1; k--){
    for(j = k+1; j < N; j++){
```

## syrk

```c
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];
```

## gemm

```c
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];
```

## sytd2 – reduction to symmetric tridiagonal form

```c
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[
  A[i+1][i] = ( A[
```

## gebd2 – reduction to bidiagonalization

```c
( 1.0e+00 + norma2 / ( A[i+1][i] * A[i+1][i] ) ) ;
for(k = 0; k < N; k++){
  norma2 = 0.e+00;
  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] -
  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]=
  for(j = k+1;
    tau[k-1] +=
  }
```

## cgs – Classical Gram-Schmidt

```c
for (j = 0; j < N; j++) {
  for (i = 0; i < j; i++) {
    R[i][j] = 0.0e+00;
    for (k = 0; k <
      += A
```

## apply_givens_rot_sequence

```c
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;
```

## mgs – Modified Gram-Schmidt

```c
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];
}
```

## cholesky

```c
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++) {
```

## gehd2 – reduction to Hessenb

```c
for ( j = 0; j < n-2; j++ ) {
  j+2; i < n; i++ ) {
    += A[i][j] * A[i][j] ;
  rt ( A[j+1][j] * A[j+1][j] + norma2 )
  = ( A[j+1][j] > 0 ) ? ( A[j+1][j] + no
  / ( 1.0 + norma2 / ( A[j+1][j] * A[j+1
  j+2; i < n; i++ ) {
    /= 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];
  }
}
tmp[i]
```

## lu

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

# The I/O Requirements of Various Numerical Linear Algebra Kernels

Julien Langou

Tuesday June 26th 2024