

# Manycore algorithms for batch scalar and block tridiagonal solvers

ENDRE LÁSZLÓ, University of Oxford

MIKE GILES, University of Oxford

JEREMY APPELEYARD, NVIDIA Corporation

Engineering, scientific and financial applications often require the simultaneous solution of a large number of independent tridiagonal systems of equations with varying coefficients. Since the number of systems is large enough to offer considerable parallelism on many-core systems, the choice between different tridiagonal solution algorithms, such as Thomas, CR (Cyclic Reduction) or PCR (Parallel Cyclic Reduction) needs to be re-examined. This work investigates the optimal choice of tridiagonal algorithm for CPU, Intel MIC and NVIDIA GPU with a focus on minimizing the amount of data transfer to and from the main memory using novel algorithms and register blocking mechanism, and maximizing the achieved bandwidth. It also considers block tridiagonal solutions which are sometimes required in CFD (Computational Fluid Dynamic) applications. A novel work-sharing and register blocking based Thomas solver is also presented.

CCS Concepts: • **Mathematics of computing** → **Solvers; Mathematical software performance;**

Additional Key Words and Phrases: scalar tridiagonal solver, block tridiagonal solver, CPU, GPU, MIC, Xeon Phi, CUDA, vectorization

## ACM Reference Format:

Endre László, Mike Giles and Jeremy Appleyard, 2014. Manycore algorithms for batch scalar and block tridiagonal solvers. *ACM Trans. Math. Softw.* 0, 0, Article 0 (November 2014), 39 pages.

DOI: 0000001.0000001

## INTRODUCTION

The numerical approximation of multi-dimensional PDE problems on regular grids often requires the solution of multiple tridiagonal systems of equations. In computational finance such problems arise frequently, due to the ADI (alternating direction implicit) time discretization favored by many in the community, see [Dang et al. 2010]. The ADI method requires the solution of multiple tridiagonal systems of equations in each dimension of the problem, see [Peaceman and Rachford 1955; Craig and Sneyd 1988; Douglas and Rachford 1956; Douglas and Gunn 1964]. The requirement also arises when using line-implicit smoothers as part of a multi-grid solver [Douglas et al. 1998], and when using high-order compact differencing [Dring et al. 2014; Karaa and Zhang 2004]. In most cases the tridiagonal systems are scalar, with one unknown per grid point, but this is not always the case. For example, computational fluid dynamics

---

The research at the Oxford e-Research Centre has been partially supported by the ASEArch project on Algorithms and Software for Emerging Architectures, funded by the UK Engineering and Physical Sciences Research Council. The research has also been supported by the TÁMOP-4.2.1./B-11/2/KMR-2011-002, TÁMOP - 4.2.2./B-10/1-2010-0014 projects at PPCU - University of National Excellence. The authors would like to acknowledge the use of the University of Oxford Advanced Research Computing (ARC) facility in carrying out this work. We are thankful for the support of Péter Szolgay at PPCU-FIT, Hungary.

Author's addresses: Endre László, e-Research Centre, University of Oxford, (Current address) Faculty of Information Technology and Bionics, Pázmány Pter Catholic University; Mike Giles, Mathematical Institute, University of Oxford; Jeremy Appleyard, NVIDIA Corporation. E-mail: laszlo.endre@itk.ppke.hu, mike.giles@maths.ox.ac.uk, jappleyard@nvidia.com

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

© 2014 ACM. 0098-3500/2014/11-ART0 \$15.00

DOI: 0000001.0000001

applications often have systems with block-tridiagonal structure up to 8 unknowns per grid point [Pulliam 1986]. The paper is structured into two parts. The first part of the paper considers scalar tridiagonal system of equations and the second part considers block tridiagonal system of equations. Both parts cover methods suitable for parallel, multi and many-core architectures. The conclusions for the scalar and block tridiagonal solvers are given separately at the end of the each part.

## SCALAR TRIDIAGONAL SOLVERS ON SIMD AND SIMT ARCHITECTURES

### 1. INTRODUCTION

Let us start the discussion of the scalar tridiagonal solvers with a specific example of a regular 3D grid of size  $256^3$ . An ADI time-marching algorithm will require the solution of a tridiagonal system of equations along each line in the first coordinate direction, then along each line in the second direction, and then along each line in the third direction. There are two important observations to make here. The first is that there are  $256^2$  separate tridiagonal solutions in each phase, and these can be performed independently and in parallel, ie. there is plenty of natural parallelism to exploit on many-core processors. The second is that a data layout which is optimal for the solution in one particular direction, might be far from optimal for another direction. This will lead us to consider using different algorithms and implementations for different directions. Due to this natural parallelism in batch problems, the improved parallelism of CR (Cyclic Reduction) and PCR (Parallel Cyclic Reduction) are not necessarily advantageous and the increased work-load of these methods might not necessarily pay off. If we take the example of one of the most modern accelerator cards, the NVIDIA Tesla K40 GPU, we may conclude that the 12GB device memory is suitable to accommodate a  $N^d$  multi-dimensional problem domain with  $N^d = 12 \text{ GB} / 4 \text{ arrays} = 0.75 \times 10^9$  single precision grid points. This means that the length  $N$  along each dimension is  $N = \sqrt[d]{0.75 \times 10^9}$  for dimensions  $d = 2 - 4$  as shown on Table I.

d	N	# parallel systems
2	27386	27386
3	908	824464
4	165	4492125

Table I: Number of parallel systems increases rapidly as dimension  $d$  is increased.  $N$  is chosen to accommodate an  $N^d$  single precision problem domain on the available 12 GB of an NVIDIA Tesla K40 GPU.

Before discussing the specific implementations on different architectures, we review a number of different algorithms for solving tridiagonal systems, and discuss their properties in terms of the number of floating point operations and the amount of memory traffic generated. Research has been conducted by [Zhang et al. 2010; Chang et al. 2012] to solve tridiagonal system of equations on GPUs. The Recursive Doubling algorithm has been developed by [Stone 1973]. Earlier work by [Wang 1981; van der Vorst 1987; Bondeli 1991; Mattor et al. 1995; Spaletta and Evans 1993] has decomposed a larger tridiagonal system into smaller ones that can be solved independently, in parallel. The algorithms described in the following subsections are the key building blocks of a high performance implementation of a tridiagonal solver. We introduce the tridiagonal system of equations as shown on Eq. (1) or in its matrix form shown on Eq. (2).

$$a_i u_{i-1} + b_i u_i + c_i u_{i+1} = d_i, \quad i = 0, 1, \dots, N-1 \quad (1)$$

$$\begin{pmatrix} b_0 & c_0 & 0 & 0 & \cdots & 0 \\ a_1 & b_1 & c_1 & 0 & \cdots & 0 \\ 0 & a_2 & b_2 & c_2 & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \cdots & a_{N-1} & b_{N-1} \end{pmatrix} \begin{pmatrix} u_0 \\ u_1 \\ u_2 \\ \vdots \\ u_{N-1} \end{pmatrix} = \begin{pmatrix} d_0 \\ d_1 \\ d_2 \\ \vdots \\ d_{N-1} \end{pmatrix} \quad (2)$$

where  $a_0 = c_{N-1} = 0$ .

## 2. TRIDIAGONAL ALGORITHMS

In this section the standard algorithms like Thomas, PCR (Parallel Cyclic Reduction) and CR (Cyclic Reduction) are presented along with a new hybrid Thomas-PCR algorithm which is the combination of the Thomas and the PCR algorithms.

### 2.1. Thomas algorithm

The Thomas algorithm [Thomas 1949] is a sequential algorithm which is described in [Press et al. 2007]. It is simply a customization of Gaussian elimination to the case in which the matrix is tridiagonal. The algorithm has a forward pass in which the lower diagonal elements  $a_i$  are eliminated by adding a multiple of the row above. This is then followed by a reverse pass to compute the final solution using the modified  $c_i$  values. In both passes the  $d_i$  is also modified according to the row operations.

---

#### ALGORITHM 1: Thomas algorithm

---

**Require:**  $thomas(a, b, c, d)$

```

1:  $d_0^* := d_0/b_0$ 
2:  $c_0^* := c_0/b_0$ 
3: for  $i = 1, \dots, N-1$  do
4:    $r := 1 / (b_i - a_i c_{i-1}^*)$ 
5:    $d_i^* := r (d_i - a_i d_{i-1}^*)$ 
6:    $c_i^* := r c_i$ 
7: end for
8: for  $i = N-2, \dots, 0$  do
9:    $d_i := d_i^* - c_i^* d_{i+1}$ 
10: end for
11: return  $d$ 
```

---

The full Thomas algorithm with in-place solution (RHS vector is overwritten) is given in Algorithm 1. Note that this does not perform any pivoting; it is assumed the matrix is diagonally dominant, or at least sufficiently close to diagonal dominance so that the solution is well-conditioned. Floating point multiplication and addition, as well as FMA (Fused Multiply and Add) have the same instruction throughput on almost every hardware. Also, the use of FMA is considered to be a compiler optimization feature and therefore optimistic calculations based on this instruction give a good lower estimate on the number of actual instruction uses. In the following discussion the FMA instruction is used as the basis of estimating the work complexity of the algorithm. The computational cost per row is approximately three FMA operations, one reciprocal and two multiplications. If we treat the cost of the reciprocal as being equivalent to five FMAs (which is the approximate cost on a GPU for a double precision reciprocal), then the total cost is equivalent to 10 FMAs per grid point.

On the other hand, the algorithm requires at the minimum the loading of  $a_i, b_i, c_i, d_i$ , and the storing of the final answer  $d_i$ . Worse still, it may be necessary to store the  $c_i^*, d_i^*$  computed in the forward pass, and then read them back during the reverse pass.

This shows that the implementation of the algorithm is likely to be memory-bandwidth limited, because there are very few operations per memory reference. Thus, when using the Thomas algorithm it will be critical to ensure the maximum possible memory bandwidth.

## 2.2. PCR and CR Algorithms

PCR (Parallel Cyclic Reduction) [Gander and Golub 1997] is an inherently parallel algorithm which is ideal when using multiple execution threads to solve each tridiagonal system.

---

### ALGORITHM 2: PCR algorithm

---

**Require:**  $pcr(\mathbf{a}, \mathbf{b}, \mathbf{c}, \mathbf{d})$

```

1: for  $p = 1, \dots, P$  do
2:    $s := 2^{p-1}$ 
3:   for  $i = 0, \dots, N-1$  do
4:      $r := 1 / (1 - a_i^{(p-1)} c_{i-s}^{(p-1)} - c_i^{(p-1)} a_{i+s}^{(p-1)})$ 
5:      $a_i^{(p)} := -r a_i^{(p-1)} a_{i-s}^{(p-1)}$ 
6:      $c_i^{(p)} := -r c_i^{(p-1)} c_{i+s}^{(p-1)}$ 
7:      $d_i^{(p)} := r (d_i^{(p-1)} - a_i^{(p-1)} d_{i-s}^{(p-1)} - c_i^{(p-1)} d_{i+s}^{(p-1)})$ 
8:   end for
9: end for
10: return  $\mathbf{d}^{(P)}$ 
```

---

If we start with the same tridiagonal system of equations, but normalized so that  $b_i = 1$ ,

$$a_i u_{i-1} + u_i + c_i u_{i+1} = d_i, \quad i = 0, 1, \dots, N-1,$$

with  $a_0 = c_{N-1} = 0$ , then subtracting the appropriate multiples of rows  $i \pm 1$ , and re-normalising, gives

$$a'_i u_{i-2} + u_i + c'_i u_{i+2} = d'_i, \quad i = 0, 1, \dots, N-1,$$

with  $a'_0 = a'_1 = c'_{N-2} = c'_{N-1} = 0$ . Repeating this by subtracting the appropriate multiples of rows  $i \pm 2$  gives

$$a''_i u_{i-4} + u_i + c''_i u_{i+4} = d''_i, \quad i = 1, 2, \dots, N,$$

with  $a''_j = c''_{N-1-j} = 0$  for  $0 \leq j < 4$ .

After  $P$  such steps, where  $P$  is the smallest integer such that  $2^P \geq N$ , then all of the modified  $a$  and  $c$  coefficients are zero (since otherwise the value of  $u_i$  would be coupled to the non-existent values of  $u_{i \pm 2^P}$ ) and so we have the value of  $u_i$ .

The PCR algorithm is given in Algorithm 2. Note that any reference to a value with index  $j$  outside the range  $0 \leq j < N$  can be taken to be zero; it is multiplied by a zero coefficient so as long as it is not an IEEE exception value (such as a NaN) then any valid floating point value can be used. If the computations within each step are performed simultaneously for all  $i$ , then it is possible to reuse the storage so that  $a^{(p+1)}$  and  $c^{(p+1)}$  are held in the same storage (e.g. the same registers) as  $a^{(p)}$  and  $c^{(p)}$ . The computational cost per row is approximately equivalent to 14 FMAs in each step, so the total cost per grid point is  $14 \log_2 N$ . This is clearly much greater than the cost of the Thomas algorithm with 10 FMA, but the data transfer to/from the main memory may be lower if there is no need to store and read back the various intermediate values of the  $a$  and  $c$  coefficients.

Cyclic Reduction (CR) is a slightly different algorithm in which the  $p^{th}$  pass of the above algorithm is performed only for those  $i$  for which  $i \bmod 2^p = 0$ . This gives the

forward pass of the PCR algorithm; there is then a corresponding reverse pass which performs the back solve. This involves fewer floating point operations overall, but there is less parallelism on average and it requires twice as many steps so it is slower than PCR when there are many active threads. Hence, we have not found it to be helpful in the work presented here.

### 2.3. Hybrid Thomas-PCR Algorithm

The hybrid algorithm is a combination of a modified Thomas and PCR algorithms. Some hybrid algorithms have been developed over the years with different application aims, like [Wang 1981] who decomposed tridiagonal systems in an upper tridiagonal form to be solved multiple RHS. The slides of [Sakharnykh 2010] show that the PCR algorithm can be used to divide a larger system into smaller ones which can be solved using the Thomas algorithm. In the present paper it is shown that a larger system can be divided into smaller ones using a modified Thomas algorithm and the warp level PCR can be used to solve the remaining system. The complete computation with the presented algorithms can be done in registers using the `_shfl()` instructions introduced in the NVIDIA Kepler GPU architecture. In the case of the Thomas solver intermediate values  $c_i^*$  and  $d_i^*$  had to be stored and loaded in main memory. With the hybrid algorithm the input data is read once and output is written once without the need to move intermediate values in global memory. Therefore this algorithm allows an optimal implementation.

---

#### ALGORITHM 3: First phase of hybrid algorithm

---

**Require:**  $hybridFW(\mathbf{a}, \mathbf{b}, \mathbf{c}, \mathbf{d})$

```

1:  $d_0 := d_0/b_0$ 
2:  $a_0 := a_0/b_0$ 
3:  $c_0 := c_0/b_0$ 
4: for  $i = 1, \dots, M-2$  do
5:    $r := 1 / (b_i - a_i c_{i-1})$ 
6:    $d_i := r (d_i - a_i d_{i-1})$ 
7:    $a_i := -r a_i a_{i-1}$ 
8:    $c_i := r c_i$ 
9: end for
10: for  $i = M-3, \dots, 1$  do
11:    $d_i := d_i - c_i u_{i+1}$ 
12:    $a_i := a_i - c_i a_{i+1}$ 
13:    $c_i := -c_i c_{i+1}$ 
14: end for
```

---

**ALGORITHM 4:** Middle phase: modified PCR algorithm.

---

**Require:** *pcr2(a, b, c, d)*

```

1: for  $i = 1 : 2 : N-1$  do
2:    $r = 1 / (1 - a_i^{(0)} c_{i-1}^{(0)} - c_i^{(0)} a_{i+1}^{(0)})$ 
3:    $d_i^{(1)} = r (d_i^{(0)} - a_i^{(0)} d_{i-1}^{(0)} - c_i^{(0)} d_{i+1}^{(0)})$ 
4:    $a_i^{(1)} = -r a_i^{(0)} a_{i-1}^{(0)}$ 
5:    $c_i^{(1)} = -r c_i^{(0)} c_{i+1}^{(0)}$ 
6: end for
7: for  $p = 1, \dots, P$  do
8:    $s := 2^p$ 
9:   for  $i = 1 : 2 : N-1$  do
10:     $r = 1 / (1 - a_i^{(p-1)} c_{i-s}^{(p-1)} - c_i^{(p-1)} a_{i+s}^{(p-1)})$ 
11:     $d_i^{(p)} = r (d_i^{(p-1)} - a_i^{(p-1)} d_{i-s}^{(p-1)} - c_i^{(p-1)} d_{i+s}^{(p-1)})$ 
12:     $a_i^{(p)} = -r a_i^{(p-1)} a_{i-s}^{(p-1)}$ 
13:     $c_i^{(p)} = -r c_i^{(p-1)} c_{i+s}^{(p-1)}$ 
14:   end for
15: end for
16: for  $i = 1 : 1 : N-1$  do
17:    $d_i^{(P)} = d_i^{(P-1)} - a_i^{(P-1)} d_i^{(P-1)} - c_i^{(P-1)} d_{i+1}^{(P-1)}$ 
18: end for
19: return  $d^{(P)}$ 

```

---

**ALGORITHM 5:** Last phase: solve for the remaining unknowns.

---

**Require:** *hybrid1b(a, b, c, d)*

```

1: for  $i = 1, \dots, M-1$  do
2:    $d_i := d_i - a_i r d_0 - c_i d_{M-1}$ 
3: end for

```

---

Suppose the tridiagonal system is broken into a number of sections of size  $M$  (show in Figure 1), each of which will be handled by a separate thread. Within each of these pieces, using local indices ranging from 0 to  $M-1$ , a slight modification to the Thomas algorithm (shown in Algorithm 3) operating on rows 1 to  $M-2$  enables one to obtain an equation of the form

$$a_i u_0 + u_i + c_i u_{M-1} = d_i, \quad i = 1, 2, \dots, M-2 \quad (3)$$

expressing the central values as a linear function of the two end values,  $u_0$  and  $u_{M-1}$ . The forward and backward phases of the modified Thomas algorithm are shown on Figure 2 and 3. Using equation 3 for  $i = 1, M-2$  to eliminate the  $u_1$  and  $u_{M-2}$  entries in the equations for rows 0 and  $M-1$ , leads to a reduced tridiagonal system (shown in Figure 3) of equations involving the first and last variables within each section. This reduced tridiagonal system can be solved using PCR using Algorithm 4, and then Algorithm 5 gives the interior values.

In the last phase the solution of the interior unknowns using the outer  $i = 0$  and  $i = M-1$  values needs to be completed with Algorithm 5.

When there are 32 CUDA threads and 8 unknowns per thread the cost is approximately 14 FMAs per point, about 40% more than the cost of the Thomas algorithm. The data transfer cost depends on whether these intermediate coefficients  $a_i, c_i, d_i$  need to be stored in the main memory. If they do, then it is again more costly than the Thomas algorithm, but if not then it is less costly. We need to add this to the cost of the PCR

solution. The relative size of this depends on the ratio  $(\log_2 N)/M$ . We will discuss this in more detail later.

The hybrid Thomas-PCR solver is validated against the Thomas solver by computing the MSE (Mean Square Error) between the two solutions. The results show MSE error in the order of  $10^{-9}$  in single precision and  $10^{-18}$  in double precision.

$$\left( \begin{array}{ccc|ccc|ccc} b_0 & c_0 & & & & & & & \\ a_1 & b_1 & c_1 & & & & & & \\ & a_2 & b_2 & c_2 & & & & & \\ & & a_3 & b_3 & c_3 & & & & \\ \hline & & & a_4 & b_4 & c_4 & & & \\ & & & & a_5 & b_5 & c_5 & & \\ & & & & & a_6 & b_6 & c_6 & \\ & & & & & & a_7 & b_7 & c_7 \\ \hline & & & & & & & a_8 & b_8 & c_8 \\ & & & & & & a_9 & b_9 & c_9 \\ & & & & & & & a_{10} & b_{10} & c_{10} \\ & & & & & & & & a_{11} & b_{11} \end{array} \right) \begin{pmatrix} u_0 \\ u_1 \\ u_2 \\ u_3 \\ u_4 \\ u_5 \\ u_6 \\ u_7 \\ u_8 \\ u_9 \\ u_{10} \\ u_{11} \end{pmatrix} = \begin{pmatrix} d_0 \\ d_1 \\ d_2 \\ d_3 \\ d_4 \\ d_5 \\ d_6 \\ d_7 \\ d_8 \\ d_9 \\ d_{10} \\ d_{11} \end{pmatrix}$$

Fig. 1: Starting state of the equation system. 0 values outside the bar are part of the equation system and they are needed for the algorithm and implementation.

$$\left( \begin{array}{ccc|ccc|ccc} 1 & c_0 & & & & & & & \\ a_1^* & 1 & c_1 & & & & & & \\ & a_2^* & & 1 & c_2 & & & & \\ & & a_3^* & & 1 & c_3 & & & \\ \hline & & & a_4^* & 1 & c_4 & & & \\ & & & & a_5^* & 1 & c_5 & & \\ & & & & & a_6^* & 1 & c_6 & \\ & & & & & & a_7^* & 1 & c_7 \\ \hline & & & & & & & a_8^* & 1 & c_8 \\ & & & & & & a_9^* & 1 & c_9 \\ & & & & & & & a_{10}^* & 1 & c_{10} \\ & & & & & & & & a_{11}^* & 1 \end{array} \right) \begin{pmatrix} u_0 \\ u_1 \\ u_2 \\ u_3 \\ u_4 \\ u_5 \\ u_6 \\ u_7 \\ u_8 \\ u_9 \\ u_{10} \\ u_{11} \end{pmatrix} = \begin{pmatrix} d_0^* \\ d_1^* \\ d_2^* \\ d_3^* \\ d_4^* \\ d_5^* \\ d_6^* \\ d_7^* \\ d_8^* \\ d_9^* \\ d_{10}^* \\ d_{11}^* \end{pmatrix}$$

Fig. 2: State of the equation system after the forward sweep of the modified Thomas algorithm.

### 3. DATA LAYOUT IN MEMORY

Knowledge of the data layout and the access patterns of the executed code with various algorithms is essential to understand the achieved performance. Taking the example of a 3D application shown on Figure 4, each of the coefficients arrays is a 3D indexed array which is mathematically of the form  $u_{i,j,k}$ . Ie. each arrays a, b, c, d and u are stored in a separate cubic data array.

However this is stored in memory as a linear array and so we require a mapping from the index set  $(i, j, k)$  to a location in memory. If the array has dimensions  $(N_x, N_y, N_z)$ , we choose to use the mapping:

$$\text{index} = i + j N_x + k N_x N_y.$$

$$\begin{pmatrix}
 \mathbf{1} & & & & & & & & & & & \\
 a_1^* & \mathbf{1} & & & & & & & & & & \\
 a_2^* & & \mathbf{1} & & & & & & & & & \\
 \mathbf{a_3^*} & & & \mathbf{1} & & & & & & & & \\
 & & & & \mathbf{a_4^*} & & & & & & & \\
 & & & & & \mathbf{1} & & & & & & \\
 & & & & & & \mathbf{a_5^*} & & & & & \\
 & & & & & & & \mathbf{1} & & & & \\
 & & & & & & & & \mathbf{a_6^*} & & & \\
 & & & & & & & & & \mathbf{1} & & \\
 & & & & & & & & & & \mathbf{a_7^*} & \\
 & & & & & & & & & & & \mathbf{1} \\
 & & & & & & & & & & & & \mathbf{a_8^*} \\
 & & & & & & & & & & & & & \mathbf{1} \\
 & & & & & & & & & & & & & & \mathbf{a_9^*} \\
 & & & & & & & & & & & & & & & \mathbf{1} \\
 & & & & & & & & & & & & & & & & \mathbf{a_{10}^*} \\
 & & & & & & & & & & & & & & & & & \mathbf{1} \\
 & & & & & & & & & & & & & & & & & & \mathbf{a_{11}^*}
 \end{pmatrix}
 \begin{pmatrix}
 \mathbf{u_0} \\
 u_1 \\
 u_2 \\
 \mathbf{u_3} \\
 \mathbf{u_4} \\
 u_5 \\
 u_6 \\
 \mathbf{u_7} \\
 \mathbf{u_8} \\
 u_9 \\
 u_{10} \\
 \mathbf{u_{11}}
 \end{pmatrix}
 =
 \begin{pmatrix}
 \mathbf{d_0^*} \\
 d_1^* \\
 d_2^* \\
 \mathbf{d_3^*} \\
 \mathbf{d_4^*} \\
 d_5^* \\
 d_6^* \\
 \mathbf{d_7^*} \\
 \mathbf{d_8^*} \\
 d_9^* \\
 d_{10}^* \\
 \mathbf{d_{11}^*}
 \end{pmatrix}$$

Fig. 3: State of the equation system after the backward sweep of the modified Thomas algorithm. Bold variables show the elements of the reduced system which is to be solved with the PCR algorithm.

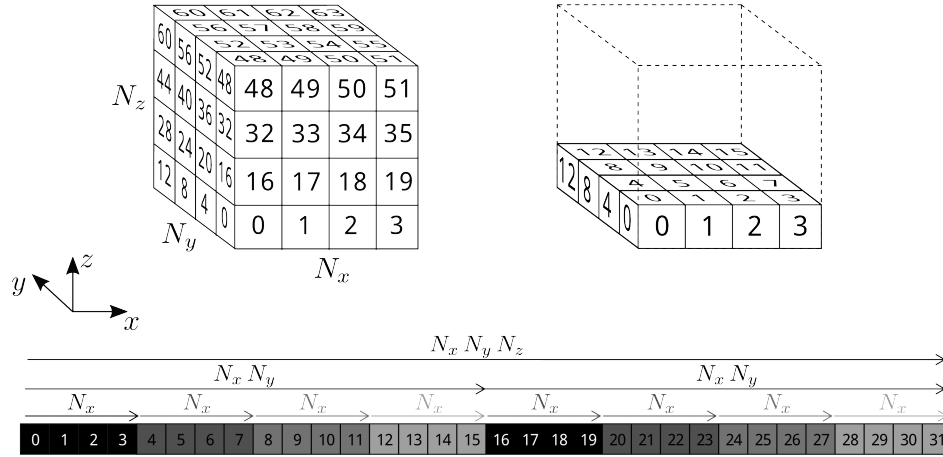


Fig. 4: Data layout of 3D cube data-structure.

This means that if we perform a Thomas algorithm solution in the first coordinate direction, then consecutive elements of the same tridiagonal system are contiguous in memory. On the other hand, in the second coordinate direction they are accessed with a stride of  $N_x$ , and in the third direction the stride is  $N_x N_y$ . This data layout extends naturally to applications with higher number of dimensions.

To understand the consequences of this layout and access pattern, it is also necessary to understand the operation of the cache memory hierarchy. The memory bus which transfers data from the main memory (or graphics memory) to the CPU (or GPU) does so in cache lines which have a size usually in the range of 32-128 bytes. These cache lines are held in the cache hierarchy until they are displaced to make room for a new cache line. One difference between CPUs and GPUs is that CPUs have a lot more cache memory per thread, and so cache lines are usually resident for longer.

The effectiveness of caches depends on two kinds of locality: 1) *temporal locality*, which means that a data item which has been loaded into the cache is likely to be used again in the near future; 2) *spatial locality*, which means that when a cache line is loaded to provide one particular piece of data, then it is likely that other data items



in the same cache line will be used in the near future. An extreme example of spatial locality can occur in a GPU when a set of 32 threads (known as a *warp*) each load one data item. If these 32 items form a contiguous block of data consisting of one or more cache lines then there is perfect spatial locality, with the entire line of data being used. This is referred to as a *coalesced* read or write. See [NVIDIA 2015] for further details.

Aligned memory access for efficient data transfer [Intel 2012a], [Intel 2012b], [NVIDIA 2015], TLB (Translation Lookaside Buffer) [Intel 2012a], [Wong et al. 2010] hit rate reduction for lower cache latency are among the techniques used to reduce execution time. Avoiding cache and TLB cache thrashing [Intel 2012a] can be done by proper padding which is the responsibility of the user of the library.

#### 4. DETAILS OF CPU, MIC AND GPU HARDWARE

Processors used in the present study include a high-end, dual socket Intel Xeon server CPU, Intel Xeon Phi coprocessor and NVIDIA K40m GPU. Processor specifications are listed in the Appendix A. CPUs operate with fast, heavy-weight cores, with large cache, out-of-order execution and branch prediction. GPUs on the other hand use light weight, in-order, hardware threads with moderate, but programmable caching capabilities and without branch prediction.

The CPU cores are very complex out-of-order, shallow-pipelined, low latency execution cores, in which operations can be performed in a different order to that specified by the executable code. This on-the-fly re-ordering is performed by the hardware to avoid delays due to waiting for data from the main memory, but is done in a way that guarantees the correct results are computed. Each CPU core also has an AVX vector unit. This 256-bit unit can process 8 single precision or 4 double precision operations at the same time, using vector registers for the inputs and output. For example, it can add or multiply the corresponding elements of two vectors of 8 single precision variables. To achieve the highest performance from the CPU it is important to exploit the capabilities of this vector unit, but there is a cost associated with gathering data into the vector registers to be worked on. Each core owns an L1 (32KB) and an L2 (256KB) level cache and they also own a portion of the distributed L3 cache. These distributed L3 caches are connected with a special ring bus to form a large, usually 15-35 MB L3 or LLC (Last Level Cache). Cores can fetch data from the cache of another core.

The Xeon Phi or MIC (Many Integrated Core) is Intel's HPC coprocessor aimed to accelerate compute intensive problems. The architecture is composed of 60 individual simple, in-order, deep-pipelined, high latency, high throughput CPU cores equipped with 512-bit wide floating point vector units which can process 16 single precision or 8 double precision operations at the same time, using vector registers for inputs and output. This architecture faces mostly the same programmability issues as the Xeon CPU, although due to the in-order execution the consequences of inefficiencies in the code can result in more server performance deficit. The MIC coprocessor uses a similar caching architecture as the Xeon CPU but has only two levels: L1 with 32KB and the distributed L2 with 512KB/core.

The Kepler generation of NVIDIA GPUs have a number of SMX (Streaming Multiprocessor – eXtended) functional units. Each SMX has 192 relatively simple in-order execution cores. These operate effectively in warps of 32 threads, so they can be thought of as vector processors with a vector length of 32. To avoid poor performance due to the delays associated with accessing the graphics memory, the GPU design is heavily multi-threaded, with up to 2048 threads (or 64 warps) running on each SMX simultaneously; while one thread in a warp waits for data, execution can switch to a different thread in the warp which has the data it requires for its computation. Each thread has its own share (up to 255 32 bit registers) of the SMX's registers (65536 32 bit registers) to avoid the context-switching overhead usually associated with multi-

threading on a single core. Each SMX has a local L1 cache, and also a shared memory which enables different threads to exchange data. The combined size of these is 64kB, with the relative split controlled by the programmer. There is also a relatively small 1.5MB global L2 cache which is shared by all SMX units.

## 5. CPU AND MIC SOLUTION

In the present section the CPU and MIC solutions are described. The MIC implementation is essentially the same as the CPU implementation with some minor differences in terms of the available ISA (Instruction Set Architecture). The ZMM vector registers are 512 bit wide therefore they allow handling 8x8 double precision or 16x16 single precision data transposition.

To have an efficient implementation of the Thomas algorithm running on a CPU, compiler auto-vectorization and single thread, sequential execution is not sufficient. To exploit the power of a multi-core CPU the vector instruction set together with the multi-threaded execution needs to be fully utilized. Unfortunately, compilers today are not always capable of making full use of the vector units of a CPU even though the instruction set would make it feasible. Often, algorithmic data dependencies and execution flow can prevent the vectorization of a loop. Usually these conditions are related to the data alignment, layout and data dependencies in the control flow graph. Such conditions are live-out variables, inter-loop dependency, non-aligned array, non-contiguous data-access pattern, array aliasing, etc. If the alignment of arrays cannot be proven at compile time, then special vector instructions with unaligned access will be used in the executable code, and this leads to a significant performance deficit. In general, the largest difference between the vector level parallelism of the CPU and GPU are the differences how the actual parallelism is executed. In the case of the CPU and MIC code the compiler decides how a code segment can be vectorized with the available instruction set – this is called SIMD (Single Instruction Multiple Data) parallelism. The capabilities of the ISA influences the efficiency of the compiler to accomplish this task. GPUs on the other hand leave the vector level parallelization for the hardware – this is SIMT (Single Instruction Multiple Threads). It is decided in run-time by the scheduling unit whether an instruction can be executed in parallel or not. If not, then the threads working in a group (called a warp in CUDA terms) are divided into a set of idle and active threads. When the active threads complete the task, another set of idle threads are selected for execution. This sequential scheduling goes on while there are unaccomplished threads left. A more detailed study on this topic in the case of unstructured grids can be found in [Reguly et al. 2014].

Since the Thomas algorithm needs huge data traffic due to the varying coefficients, one would expect the implementation of the algorithm to be limited by memory bandwidth. If an implementation is bandwidth limited and can be implemented with a non-branching computation stream the GPU is supposed to be more efficient than the CPU due to the 2-4 times larger memory bandwidth. But, this is not true for CPUs with out-of-order execution, large cache and sophisticated caching capabilities. A single socket CPU in Appendix A has 20 MB of LLC/L3 (Last Level Cache) per socket, 8 x 256 KB L2 and 8 x 32KB L1 Cache per socket, which is in total 40 MB L3, 4 MB L2 and 512 KB L1 cache total. This caching/execution mechanism makes the CPU efficient when solving a tridiagonal algorithm with the Thomas algorithm. The two temporary arrays ( $c^*$  and  $d^*$ ) can be held in this cache hierarchy and reused in the backward sweep. Therefore input data ( $a$ ,  $b$ ,  $c$  and  $d$ ) only passes through the DDR3 memory interface once when it is loaded and result array ( $u$ ) passes once when it is stored.

This means that a system in single precision with 3 coefficient arrays (a RHS array and 2 temporary arrays) can stay in L1 cache up to system length 1365 if no hardware hyper-threading (HTT) is used or half the size, namely 682 if HTT is used. The effi-

ciency of the solver still remains very good above this level and a gradual decrease can be observed as L2 and L3 get more significant role in caching.

*Thomas solver in X dimension.* The data layout described in Section 3 doesn't allow for natural parallelism available with AVX vector instructions. The loading of an 8-wide SP (Single Precision) or 4-wide DP (Double Precision) chunk of array into a register can be accomplished with the use of `_mm256_load_p{s,d}` intrinsic or the `vmovap{s,d}` assembly instruction if the first element of the register is on an aligned address, otherwise a `_mm256_loadu_p{s,d}` intrinsic or the `vmovup{s,d}` assembly instruction needs to be used.

Unfortunately, since the data in the vector register belongs to one tridiagonal system, it needs to be processed sequentially. When the algorithm acquires a value of a coefficient array, the consecutive 7 SP (or 3 DP) values will also be loaded into the L1 cache. The L1 cache line size is the same as the vector register width and therefore no saving on the efficiency of the main memory traffic can be achieved. According to Appendix A this memory can be accessed with 4 clock cycle latency. On a Xeon server processor this latency is hidden by the out-of-order execution unit if data dependencies allow it. When the instructions following the memory load instruction are independent from the loaded value – eg. the instructions to calculate the index of the next value to be read – the processor can skip waiting for the data to arrive and continue on with the independent computations. When the data arrives or the execution flow reaches an instruction that depends on the loaded value, the processor enters in idle state until the data arrives. L1 caching and out-of-order execution on the CPU therefore enables the Thomas algorithm in the  $X$  dimension to run with high efficiency, although some performance is still lost due to non-vectorized code.

Vectorization and efficient data reuse can be achieved by applying local data transposition. This local transposition combines the register or cache blocking optimization with butterfly transposition. If the code is written using intrinsic operations, the compiler decides where the temporary data is stored, either in registers or cache. The advantage of this approach is that it allows for the use of vector load/store and arithmetic instructions and vector registers (or L1 cache) which leads to highly efficient execution.

The process in the case in double precision coefficient arrays is depicted on Figure 6. As the AVX YMM registers are capable holding 4 DP values, 4 tridiagonal systems can be solved in parallel. The first 4 values of the 4 systems are loaded and stored in 4 separate vector registers for each coefficient array  $a$ ,  $b$  etc. This procedure allows for perfect cache line utilization and it is depicted on Figure 6 for the first  $4 \times 4$  array case. The same procedure is applied for the  $8 \times 8$  single precision case with one additional transposition of  $4 \times 4$  blocks. The process in the case of  $4 \times 4$  of array  $a$  is the following. Load the first 4 values of the first tridiagonal system into YMM0 register, load the first 4 values of the second tridiagonal system into YMM1 register etc. Once data is in the register perform transposition according to Figure 6 and perform 4 steps of the Thomas algorithm. Do the same for the next 4 values of the system. Repeat this procedure until the end of the systems are reached.

*Thomas solver in Y and Z dimension.* The data layout described in Section 3 suggests natural parallel execution with `_mm256_load{u}_p{s,d}` loads. Unfortunately, the compiler is not able to determine how to vectorize along these dimensions, since it can not prove that the neighboring solves will access neighboring elements in the nested *for* loops of the forward and backward phases. Not even Intel's elemental function and array notation is capable of handling this case correctly.

Since solving tridiagonal problems is inherently data parallel and data reads fit the AVX load instructions, manual vectorization with AVX intrinsic functions are used to

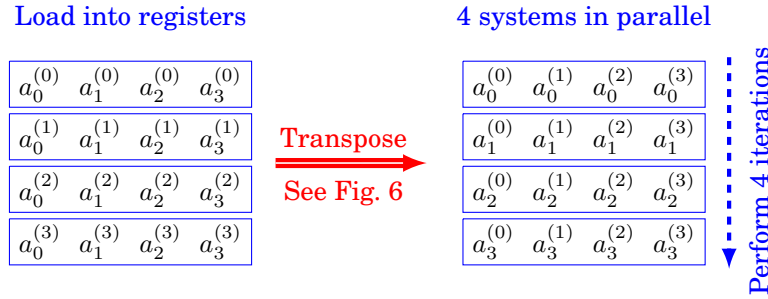


Fig. 5: Vectorized Thomas solver on CPU and MIC architectures. Data is loaded into  $M$  wide vector registers in  $M$  steps. On AVX  $M = 4$  for double precision and  $M = 8$  for single precision. On IMCI  $M = 8$  for double precision and  $M = 16$  for single precision. Data is then transposed within registers as described in Figure 6 and after that  $M$  number of iterations of the Thomas algorithm is performed. When this is done the procedure repeats until the end of the systems is reached.

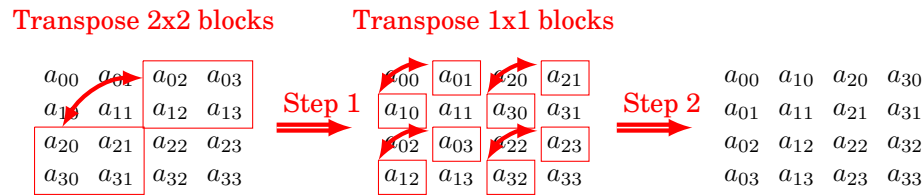


Fig. 6: Local data transposition of an  $M \times M$  two dimensional array of  $M$  consecutive values of  $M$  systems, where  $a_{ti}$  is the  $i$ th coefficient of tridiagonal system  $t$ . Transposition can be performed with  $M \times \log_2 M$  AVX or IMCI shuffle type instructions such as swizzle, permute and align or their masked counter part. On AVX  $M = 4$  for double precision and  $M = 8$  for single precision. On IMCI  $M = 8$  for double precision and  $M = 16$  for single precision.

vectorize the method. Compared to a scalar, multi-threaded implementation, using the explicit vectorization gives a 2.5 times speedup in dimension  $Y$  in single precision.

The vectorized code reads 8/4 consecutive SP/DP scalar elements into YMM registers (`_mm256` or `_mm256d`) from the multidimensional coefficient array. The implementation is straightforward with `_mm256_load_p{s,d}` intrinsic and the `vmovap{s,d}` assembly instruction if the first element of the register is on aligned address or with `_mm256_loadu_p{s,d}` intrinsic and the `vmovup{s,d}` assembly instruction when data is not aligned. In case the length along the  $X$  dimension is not a multiple of 8/4 padding is used to achieve maximum performance. Arithmetic operations are carried out using `_mm256_{add,sub,mul,div}_p{s,d}` intrinsics.

When stepping in the  $Z$  dimension to solve a system, one may hit the NUMA penalty. The consequences of ccNUMA (cache coherent NUMA) are two-fold. First, the LLC miss rate increases since the values of  $d$  are still in the cache where they have been used the last time. The coherent cache is implemented through the QPI bus which introduces an extra access latency when the cache line needs to be fetched from a different socket. LLC cache miss rate can go up to a level of 84.7% of all LLC-cache access in the  $Z$  dimension. The second consequence is the TLB (Translation Lookaside Buffer) miss penalty which occurs when elements in an array are accessed with large stride. A TLB page only covers a 4KB range of the dynamically allocated memory. Once the data acquired is outside this range, a new page frame pointer needs to be loaded into the TLB cache. The latency of doing this is the time taken to access the main

memory and time taken by the Linux kernel to give the pointer and permission to access that particular page.

## 6. GPU SOLUTION

The GPU implementation of the tridiagonal solver is not as straightforward as the CPU solver with the Thomas algorithm. Regardless of the underlying algorithms (Thomas or Thomas-PCR hybrid), solving tridiagonal systems is usually considered to be memory bandwidth limited, since the ratio of the minimal amount of data that needs to be loaded ( $4N$ ) and stored ( $N$ ) and the minimal amount of FLOPs need to be carried out (with the Thomas algorithm  $9N$ ) is  $\frac{9NFLOP}{(4+1)N\text{ values}}$ . This figure is called the compute ratio and implies that for every loaded value this amount of compute needs to be performed. This figure also depends on the SP/DP floating point data and processing unit throughput. Thus the required compute ratio for single and double precision is  $0.45 \frac{FLOP}{byte}$  and  $0.225 \frac{FLOP}{byte}$  respectively. An algorithm is theoretically expected to be memory bandwidth limited if the compute ratio capability of the specific hardware exceeds the compute ratio of the algorithm. For the K40 GPU these figures are  $16.92 \frac{FLOP}{byte}$  and  $14.89 \frac{FLOP}{byte}$  respectively which suggest that solving the problem with the most compute-efficient algorithm – the Thomas algorithm – is memory bound. In the forthcoming discussion the aim is to improve this bottleneck and achieve high memory utilization with suitable access patterns and reduced memory traffic.

Global memory load on the GPU poses strict constraints on the access patterns used to solve systems of equations on a multidimensional domain. The difference with regards to CPUs comes from the way SIMT thread level parallelism provides access to data in the memory. Solvers using unit-stride access pattern (along the  $X$  dimension) and the long-strided access pattern (along  $Y$  and higher dimensions) need to be distinguished. The former is explained in Section 6.1 while the latter is explained in Section 6.2.

The memory load instructions are issued by threads within the same warp. In order to utilize the whole 32 (or 128) byte cache line threads need to use all the data loaded by that line. At this point we need to distinguish between the two efficient algorithms discussed in this paper for solving tridiagonal problems along the  $X$  dimension, because the two needs different optimization strategies. These algorithms are

- (1) Thomas algorithm with optimized memory access pattern, detailed in Sections 6.1
- (2) Thomas-PCR hybrid algorithm, detailed in Section 6.3

### 6.1. Thomas algorithm with optimized memory access pattern

The optimization is performed on the Thomas algorithm, which is detailed in Section 2.1. The problem with the presented solver is two-fold: 1) the access pattern along dimension  $X$  is different from the pattern along dimensions  $Y$  and higher ; 2) in  $X$  dimension the actual data transfer is more than the theoretically required lower limit. In this section it is shown how to overcome problem 1) and how to optimize the effect of 2) on a GPU. In order to make the discussion unambiguous, the focus of the following discussion is put on the single precision algorithm. The same argument applies for the access pattern of double precision.

The first problem becomes obvious when one considers execution time along the three dimensions of an ADI solution step, see Figure 7. The  $X$  dimensional execution time is more than one order of magnitude lower than the solution along the  $Y$  and  $Z$  dimension. This is due to two factors: 1) high cache under-utilization and 2) high TLB thrashing.

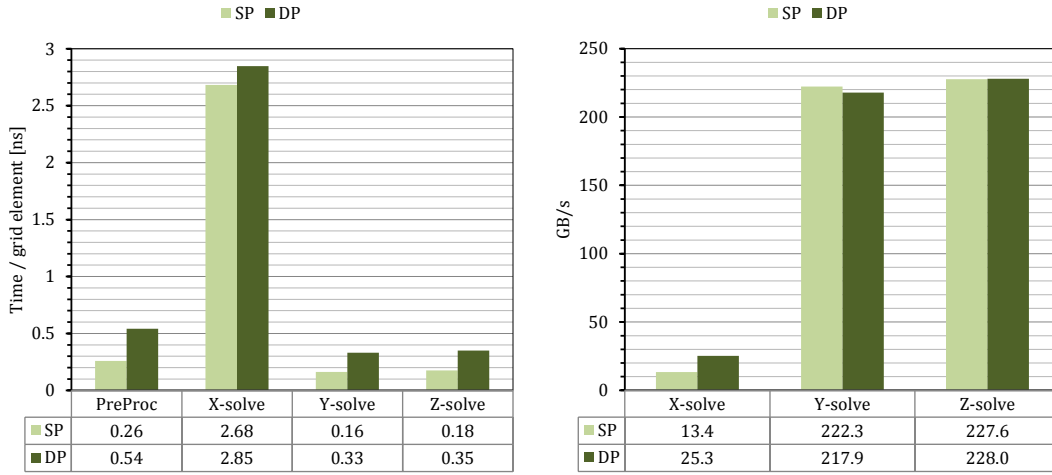


Fig. 7: ADI kernel execution times on K40m GPU, and corresponding bandwidth figures for  $X$ ,  $Y$  and  $Z$  dimensional solves based on the Thomas algorithm. Single precision  $X$  dimensional solve is slower  $\times 16.75$  than the  $Y$  dimensional solve. High bandwidth in the case of  $Y$  and  $Z$  dimensional solve is reached due to local memory caching and read-only (texture) caching.

The high cache line underutilization can be explained with the access pattern of coefficient  $a_i$  in Algorithm 1. A whole 32 byte cache line with 8 SP values is loaded into the L2/L1 or non-coherent cache, but temporarily only one value is used by the algorithm before the cache line is subsequently evicted – this results in poor cache line utilization. The other values within the cache line are used in a different time instance of the solution process, but at that point the cache line will already be flushed from the cache. The same applies to arrays  $b$ ,  $c$ ,  $d$ ,  $u$  as well. To overcome this issue one needs to do local data transposition with cache or register blocking similar to that discussed in Section 5. Two solutions to perform this optimization are introduced in Section 6.1.1 and Section 6.1.2. The idea of these optimizations is to increase cache line utilization by reading consecutive memory locations with consecutive threads into a scratch memory (shared memory or registers) and then redistributing them as needed.

The issue with TLB thrashing in the  $X$  dimensional solution relates to the mapping of threads to the problem. One might map a 2 dimensional thread block on the  $Y$ - $Z$  plane of the 3 dimensional ADI problem, and process the systems along the  $X$  dimension. In this case the start index ( $ind$  in  $a[ind]$ ) calculation is easy, according to the following formula:

---

**ALGORITHM 6:** The common way of mapping CUDA threads on a 3D problem.

---

```

1:  $y = \text{threadIdx.x}$ 
2:  $z = \text{threadIdx.y}$ 
3:  $ind = z * NX * NY + y * NX$ 

```

---

The problem with this mapping is that it introduces significant TLB (Translation Lookaside Buffer) thrashing. TLB is an important part of the cache hierarchy in every processor with virtual memory. A virtual memory page covers a certain section of the memory, usually 4 kB in the case of systems with relatively low memory and 2 MB "huge memory page" on systems with large memory. Since NVIDIA does not reveal the details of its cache and TLB memory subsystem, further details on the supposed TLB structure for the older GT200 architecture can be found in [Wong et al. 2010]. Note

that significant architecture changes have been made since GT200, but the latency is expected to be dominated by DRAM access latency, which did not change significantly in the last few years. The problem with TLB thrashing arises as  $z$  changes and when the different coefficient arrays are accessed. This access pattern induces reads from  $NX \times NY \times 4$  bytes distance and even larger when reading different arrays, which are  $NX \times NY \times NZ$  apart. For sufficiently large domain this requires the load of a new TLB page table pointer for every  $z$ . The TLB cache can only handle a handful of page pointers, thus in such a situation thrashing is more significant. Depending on the level of TLB page misses the introduced latency further increases. Explaining the in-depth effects of TLB is out of the scope of this paper and the reader is referred to [Wong et al. 2010].

The solution for this problem is simple. One needs to preserve data locality when accessing the coefficients of the systems, by mapping the threads to always solve the closest neighboring system. One may map the threads with a 1 dimensional thread block according to:

---

**ALGORITHM 7:** Mapping threads to neighboring systems.

---

```
1: tid = threadIdx.x + blockIdx.x*blockDim.x
2: ind = tid*NX
```

---

The inefficiency due to non-coalesced memory access has also been shown on the slides of [Sakharnykh 2009] and a global transposition solution has been given. This essentially means that the data has been transposed before performing the execution of the tridiagonal solve and therefore data is read and written unnecessarily during the transposition. In the present work we give two local data transposition solutions both which keep the data in registers for the time of the tridiagonal solution. These solutions work in a register-blocking manner and therefore avoid the need of the load and store of the global transposition resulting in much higher efficiency.

*6.1.1. Thomas algorithm with local transpose in shared memory.* Local data transpose can be performed in shared memory available on most GPU devices. The optimization is based on warp-synchronous execution, therefore there is no need for explicit synchronization. Although a warp size of 32 is assumed in the following, the implementation uses macros to define warp size and thus it allows for changes in the future. Threads within a warp cooperate to read data into shared memory. The data is then transposed via shared memory and stored back in registers. The threads perform 8 steps of the Thomas algorithm with the data in registers and then read the next 8 values, and so on. Meanwhile the updated coefficients and intermediate values  $c^*, d^*$  are stored in the local memory, which automatically creates coalesced and cached load/store memory transactions. The algorithm is shown in Algorithm 8 and further detailed in Figure 8.

---

**ALGORITHM 8:** Thomas algorithm with shared memory transpose

---

Forward pass:

- 1: Wrap a warp (32 threads) into  $4 \times 8$  blocks to perform non-caching (32byte) loads
- 2: Load  $32 \times 8$  size tiles into shared memory: in 8 steps of  $4 \times 8$  block loads
- 3: Transpose data by putting values into *float a[8]* register arrays;
- 4: Perform Thomas calculation with the 8 values along the  $X$  dimension
- 5: Repeat from step 2 until end of  $X$  dimension is reached

Backward pass:

- 1: Compute backward stage of Thomas in chunks of 8
  - 2: Transpose results and store in global memory
-

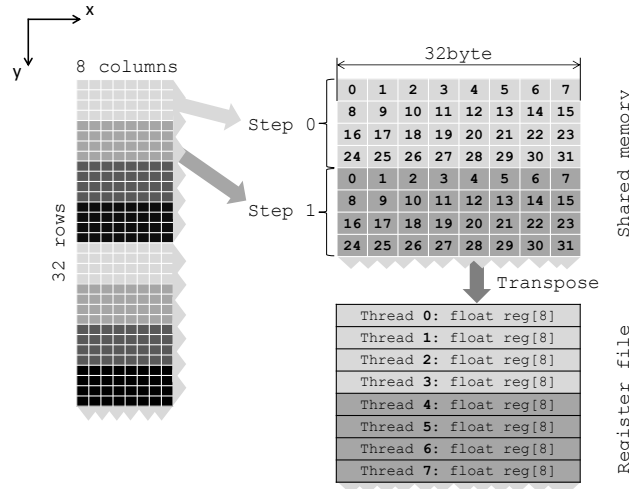


Fig. 8: Local transpose with shared memory.

The shared memory on an NVIDIA GPU can be accessed through 32 banks of width 32 or 64 bits – latter only applicable from NVIDIA Compute Capability 3.0 and higher. Since the width of the block that is stored in shared memory is 8, this leads to shared memory bank conflicts in the 3rd step of Algorithm 8. In the first iteration the first 4 threads will access banks 0,8,16,24, the next 4 threads will again access 0,8,16,24 and so on. To overcome this problem the leading dimension (stride) in shared memory block is padded to 9 instead of 8. This common trick helps to get rid of most of the bank conflicts. The first 4 threads will access banks 0,9,18,27, the next 4 threads will access banks 1,10,19,28 and so on. The amount of shared memory utilized for the 4 register arrays is  $9 * 32 * 4 * 4 \text{ bytes/warp} = 4608 \text{ bytes/warp}$ .

**6.1.2. Thomas algorithm with local transpose using register shuffle.** Local data transposition can also be performed in registers exclusively. This solution fits recent NVIDIA GPU architectures, starting with the Kepler architecture with Compute Capability 3.0. The optimization is again based on warp-synchronous execution with the use of `__shfl()` register shuffle intrinsic instructions. Although a warp size of 32 is assumed in the following, the implementation uses C macros definitions to set warp size and thus it allows for future architecture change. Threads within a warp cooperate to read a  $32 \times 16$  SP block or  $32 \times 8$  DP block of the x-y plane, ie. threads cooperate to read 16 SP or 8 DP values of 32 neighboring systems. Every 16 or 8 threads within the warp initiate a read of two cache-lines ( $2 \times 32$  bytes). They store the data into their register arrays of size 16 for SP and 8 for DP. These local arrays will be kept in registers for the same reasons that are discussed in Section 6.1.1. Once data is read, threads cooperatively distribute values with the XOR (Butterfly) transpose algorithm using the `__shfl_xor()` intrinsic in two steps. Once every thread has the data they perform 16 SP or 8 DP steps of the Thomas algorithm with the data in registers and then read the next 16 SP or 8 DP long array, and so on. Meanwhile the updated coefficients and intermediate values  $c^*$ ,  $d^*$  are stored in the local memory, which automatically creates coalesced, cached load/store memory transactions. The algorithm is shown in Algorithm 9 and further detailed in Figure 9.



**ALGORITHM 9:** Thomas algorithm with register shuffle transpose

Forward pass:

- 1: Wrap 32 threads into  $8 \times 4$  blocks to perform  $4 \times \text{float4}$  vector loads
- 2: Load  $32 \times 16$  size tiles into registers:
- 3: 4 threads read 4 consecutive *float4* vectors = 64 bytes
- 4: Do this 4 times for rows under each other
- 5: Transpose data within 4 threads:
- 6: 4 threads exchange data on a  $4 \times 4$  2D array with *\_shfl(float4)*
- 7: Each element in the 2D array is a *float4* vector
- 8: Perform Thomas calculation with the 16 values along  $X$  dimension
- 9: Repeat from 2 until end of  $X$  dimension is reached

Backward pass:

- 1: Compute backward stage of Thomas in chunks of 8
- 2: Transpose results and store in global memory

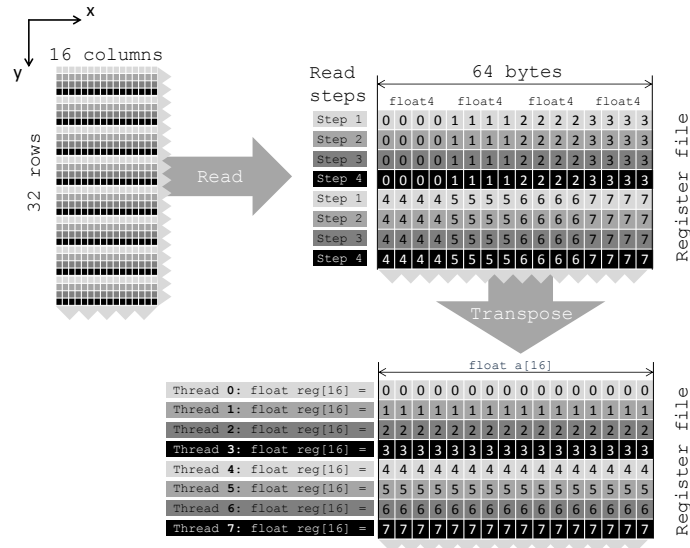


Fig. 9: Local transpose with registers.

The use of the register shuffle instruction increases the register pressure compared to the shared memory version, but since the data is kept in registers, all the operations within a 16 SP or 8 DP solve are performed in the fastest possible way, thus the overall gain is significant.

## 6.2. Thomas algorithm in higher dimensions

The Thomas algorithm gives itself to natural parallelization and efficient access pattern in the dimensions  $Y$  and above. The memory access patterns fit the coalesced way of reading and storing data. Therefore, the Thomas algorithm is efficient in these directions and the only limiting factor is the bandwidth throughput requirement posed by reading 3 coefficients and 1 RHS arrays, writing and later reading the 2 temporary arrays and finally writing out the solution array.

### 6.3. Thomas-PCR hybrid algorithm

The Thomas-PCR hybrid algorithm introduced in Section 2.3 is implemented using either shared memory or register shuffle instructions. The advantage of the hybrid algorithm is that it allows for storing the entire system in the registers of a warp. Every thread in the warp is responsible for part of the system of size  $M$ , eg. if the system size is 256, then  $M = 8$  and each thread stores an 8 long chunk of the 3 coefficients, the RHS and the solution vector in its registers. This storage method puts a high pressure on register allocation, and it stays efficient only up to a system size which allows the sub-arrays to be stored in registers. Every thread individually performs the first phase of the modified Thomas algorithm and then switches to PCR to cooperatively solve the second phase of the hybrid algorithm. Meanwhile the intermediate, modified coefficients and RHS values are kept in registers. Once the PCR finished, its solution can be used to solve the subsystems with the recently computed boundary values. Then follows the third phase where the back-substitution of the modified Thomas algorithm is performed.

In the case of the  $X$  dimensional solve, threads of a warp need contiguous chunks of size  $M$  of the whole array. This poses the same need for reading and transposing data in the  $X$  dimension, just like in Section 6.1.1 and 6.1.2.

Data transposition can be performed both using *shared memory* or register shuffles. To achieve good efficiency warp-synchronous implementation is used in the solver. Threads of a warp read consecutive values in a coalesced way, in multiple steps of 32 values and store data in shared memory. Then each thread reads its own  $M$  values from shared memory.

---

#### ALGORITHM 10: Thomas-PCR hybrid in $X$ dimension

---

- 1: Read  $N$  values with 32 threads in a coalesced way
  - 2: Transpose array data so that every thread has its own  $M$  interior values
  - 3: For each thread, compute interior  $M$  elements in terms of outer 2
  - 4: The new tridiagonal system of size  $2*32$  is solved by PCR, using shuffles or data exchange through shared memory
  - 5: Perform back-substitution for interior elements in terms of outer 2 values
- 

In  $Y$  and higher dimensions the access pattern fits the algorithm better, just like in the case of the standard Thomas algorithm. In this case there is no need for data transposition, but there is need for thread block synchronization. The solution along these dimensions is the most efficient with few restrictions. Algorithm 11 is similar to Algorithm 10 in terms of the underlying hybrid algorithm, with the significant difference that warps cooperate to solve multiple systems. Each thread within a warp solves an  $M$  size chunk. The chunks within a warp don't necessarily contribute to the same system. At a certain phase warps share the necessary information to connect the matching systems. Sharing data is done through shared memory with thread block synchronization, but the solution of the reduced system is done with PCR the same way as in Algorithm 10.

## 7. PERFORMANCE COMPARISON

The performance of the solvers presented are discussed in terms of scalability on different system sizes and along different dimensions on all three architectures – namely on an NVIDIA GPU, Intel Xeon CPU and Xeon Phi coprocessor. Further in the paper the GPU will also be referred as SIMT (Single Instruction Multiple Thread) architecture and the CPU and Xeon Phi collectively will also be referred to as SIMD (Single Instruction Multiple Data) architecture.

**ALGORITHM 11:** Thomas-PCR hybrid in dimension  $Y$  and above

- 
- 1: Threads within a warp are grouped into  $(W/G) \times G$  blocks, where  $G$  is the number of systems solved by a warp.
  - 2: The first  $G$  threads read the first  $M$  values of  $G$  systems, second  $G$  threads read the second  $M$  values and so on.
  - 3: Data of  $G$  systems are now loaded into registers
  - 4: Every thread computes an independent  $M$  size system: thread 0 computes the first  $M$  values of problem 0, thread  $G$  computes the second  $M$  values, thread  $2G$  computes the second  $M$  values etc.
  - 5: Threads cooperate through shared memory to each reduced system into one warp.
  - 6: Every warp computes an independent reduced system independently.
  - 7: Reduced solution is propagated back to threads to solve interior systems.
  - 8: Data is stored the same way as it was read.
- 

SIMD (CPU and Xeon Phi) measurement results are heavily contaminated with system noise. Therefore significantly longer integration time had to be used for averaging execution time. There is however a significant recurring spike in the execution time of the Xeon Phi results, that is worth discussing. In single precision almost every system size with a multiple of 128 and almost every system size in double precision with a multiple of 64 results in a significant, more than  $\times 2$  slowdown. This is the consequence of cache-thrashing. Cache-thrashing happens when cache-lines on  $2^n$  bytes boundaries with same  $n$  are accessed frequently. These memory addresses have different tag ID-s, but the same address within a tag. This means that threads accessing the two cache-lines of such boundaries contaminate the shared L3 (or LLC) level cache, i.e. they thrash the shared cache. This problem can be overcome if the cache architecture uses set-associative cache with high associativity – at least as many threads are sharing the cache. This is true in the case of a modern processor. The size of L2 cache is so low in the case of the Xeon Phi coprocessor, that only  $512KB/4threads = 128KB/thread$  is available and that might not allow for an efficient set-associativity.

The presented scaling measurements in the following subsections are run with fixed,  $256^2 = 65536$  number of systems. The large number of systems ensures enough parallelism even for the Thomas solver – which by nature contains the most sequential computations – in the case of the GPU. Since GPUs need a tremendous amount of work to be done in parallel to hide memory latency, the 65536 parallel work units are enough to saturate the CUDA cores and more importantly the memory controllers. The length of the systems along the different dimensions are changed by extruding the dimension under investigation from 64 up to 1024 with steps of size 4. The resulting execution times are averaged over 100-200 iterations to integrate out system noise. The execution time reported is the proportion of the total execution time to one element of the three dimensional cube and does not include the data transfer to the accelerator card. This gives good reference to compare the different solvers, since it is independent of the number and the length of the systems to be solved. It is expected that the execution time per element be constant for a solver, since the execution time is limited by the memory transfer bottleneck. The execution time differences of the algorithms presented in this section only relate to the memory transfer and memory access pattern of the particular algorithm. All implementations presented here utilize a whole cache line except the SIMD  $X$  solvers and the naïve GPU solver. The dependence on system size relates to running out of scratch memory (registers caching, L1/L2/L3 cache, TLB cache) for large systems and having enough workload for efficiency in the case of small systems. These dependences are discussed in the following in the discussion of the corresponding solver. The results are also compared against the *dtsub()* function

of the Intel Math Kernel Library 11.2 Update 2 for Linux library [Intel 2015] and the *gtsv()* function of the NVIDIA's cuSPARSE library [NVIDIA 2015]. The *dtsvb()* solver is a diagonally dominant solver which according to the Intel documentation is two times faster than the partial pivoting based *gtsv()* solver. Details of the hardware used are discussed in Appendix A. Since the implementations for  $X$  and higher dimensions differ we also need separate discussion for these cases.

In a realistic scenario in high dimensions the data transfer between the NUMA nodes is unavoidable since an NUMA-efficient layout for the  $X$  dimensional solve is not efficient for the  $Y$  dimensional solve and vice versa. Therefore no special attention was made to handle any NUMA related issues in the 2 socket CPU implementation for any dimensions.

### 7.1. Scaling in the $X$ dimension

Figures 10 and 11 show the performance scaling in the  $X$  dimension for single and double precision for all the architectures studied in this paper. Figure 13 compares the solvers for a specific setup.

**7.1.1. SIMD solvers.** The SIMD solvers rely on the regular Thomas algorithm with or without local data transposition. The detailed description of these solvers is in Section 5. The MKL *dtsvb()* solver for diagonally dominant tridiagonal systems was chosen as the baseline solver. In the following comparison all the presented SIMD implementations take advantage of multithreading with OpenMP. Threads using *KMP\_AFFINITY=compact* were pinned to the CPU and MIC cores to avoid unnecessary cache reload when threads are scheduled to another core. Other run-time optimization approaches using *numactl* and *membind* were also considered, but they did not provide any significant speedup.

As seen in Figure 10 and 11 the naïve Thomas algorithm with OpenMP outperforms the MKL *dtsvb()* solver and the transposition based Thomas solver further increases the speedup on a 2 socket Xeon processor. The naïve Thomas solver in the  $X$  dimension is not capable of utilizing any AVX instruction due to the inherent nature of vector units on CPUs, however the transposition based solvers are capable of taking advantage of the vector units as it is described in Section 5. The execution time declines and saturates in both single and double precision as the system size increases, due to the workload necessary to keep threads busy on the CPU. Forking new threads with OpenMP has a constant overhead which becomes less significant when threads have more work. Above a certain system size the CPU provides stable execution time. The efficiency of the CPU relies on the cache and out-of-order execution performance of the architecture. The size of the temporary arrays during the solution is small enough to fit into the L1 cache. The out-of-order execution is capable of hiding some of the 4 clock cycle latency of accessing these temporary, cached values which in total results in high efficiency.

The performance of the naïve Thomas algorithm on the Xeon Phi coprocessor is far from the expected. Due to the difficulty in vectorizing in the  $X$  dimension described in Section 5, the coprocessor needs to process the data with a scalar code. Since the scalar code is more compute limited than a vectorized code, and since the clock rate of the Xeon Phi is almost a third of the Xeon processor, the efficiency of the code drops to about the third of the vectorized code. However significant x2 increase in speedup is reached on the MIC architecture both in single and double precision when using the transposition based Thomas solver. The performance increase due to vectorization would imply a 8 or 16 times speedup in single or double precision respectively, but the underlying execution architecture, compiler and caching mechanism is not capable of providing this speedup. The overall performance of the Xeon Phi with the transposition

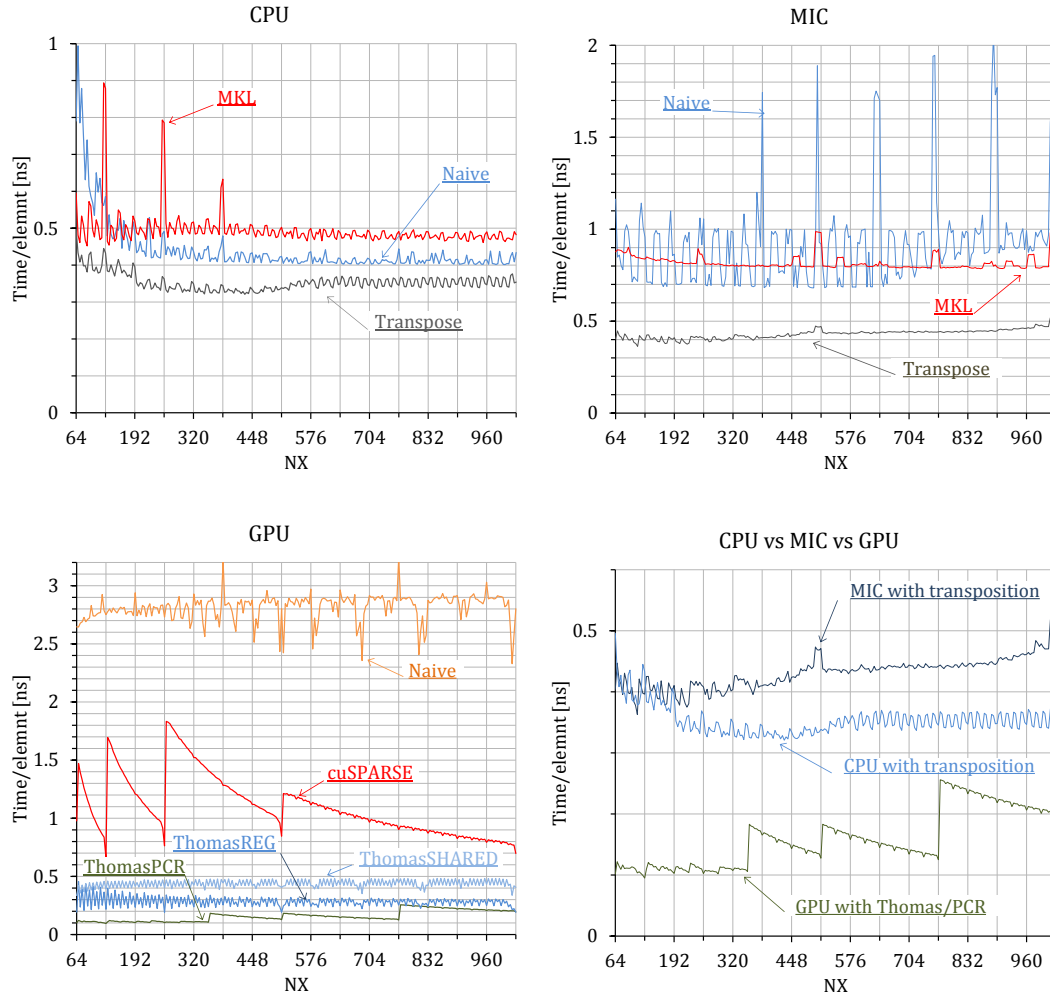


Fig. 10: Single precision execution time per grid element. 65536 tridiagonal systems with varying  $NX$  length along the X dimension are solved. CPU, MIC and GPU execution times along with the comparison of the best solver on all three architectures are show.

based Thomas algorithm is comparable to the CPU for small system size and becomes slower than the CPU as the system size increases above 200-300. Significant spikes in execution time of the naïve solver can be seen in both single and double precision on almost every 512 byte steps, either in steps of 128 in single precision or in steps of 64 in double precision. Thread pinning with *KMP\_AFFINITY=compact* option prevents thread migration and improves the performance on the coprocessor significantly. To understand the difference between the CPU and the MIC architecture the reader is suggested to study the ISA manuals [Intel 2012a] and [Intel 2012b] of the two architectures. The two architectures are radically different. The CPU works at high clock rates with complex control logic, out-of-order execution, branch prediction, low latency instructions and large, low latency distributed cache per thread. The MIC (Knights Corner) architecture is the opposite in many of these properties. It works at low clock rate, has simple control logic, in-order execution, no branch prediction, higher latency

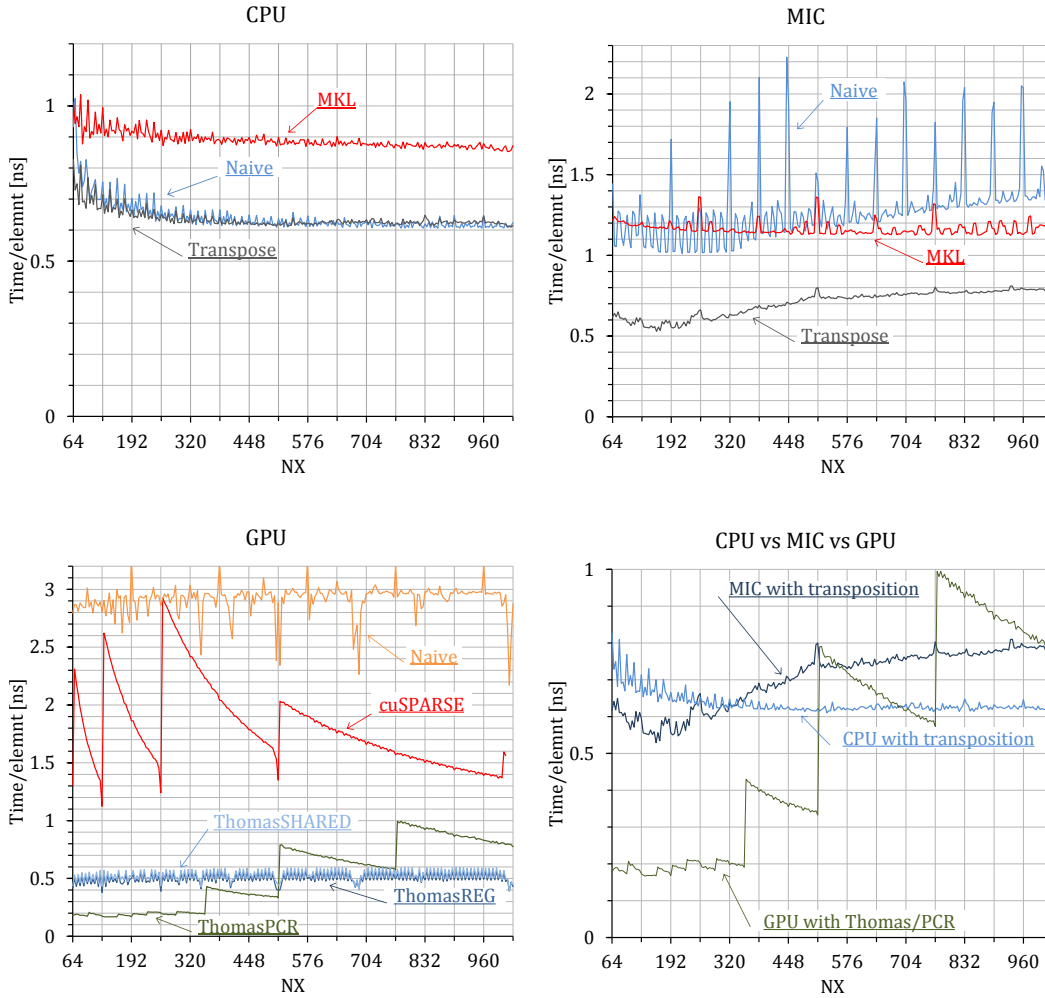


Fig. 11: Double precision execution time per grid element. 65536 tridiagonal systems with varying  $NX$  length along the X dimension are solved. CPU, MIC and GPU execution times along with the comparison of the best solver on all three architectures are show.

instructions and small, higher latency cache per thread. These fact suggest that when the input data is fetched into the cache, the heavy weight CPU cores utilize these data much efficiently than the light weight MIC cores. One may note that the upcoming Intel Knights Landing architecture with its out-of-order execution unit may significantly improve the MIC performance.

**7.1.2. SIMT solvers.** One may notice that the worst performance is achieved by the naïve GPU solver – even worse than any SIMD implementation. The inefficiency comes from the poor cache line utilizations discussed in Section 6. The measurements are contaminated with deterministic, non-stochastic noise that may come from cache efficiency and cache thrashing etc. This noise can not be attenuated by longer integration time. Due to the low cache line utilization the implementation is fundamentally latency limited. This is supported by the fact that the execution time is almost the same for the single and double precision cases. Every step of the solver requires the load of

a new cache line and since cache lines on GPU architectures are not prefetched, both single and double precision versions move the same amount of data (32 byte cache line) through the memory bus.

The cuSPARSE v7.0 solver provides better performance than the naïve solver, but it also has a significant recurring dependence on system size. The *cusparse*{*S,D*}*gtsvStridedBatch()* solver performs the best when the system size is a power of two in both single and double precision cases. The performance can vary approximately by two times speed difference in both single and double precision. This dependence is due to the internal implementation of the hybrid CR/PCR algorithm that is used inside the solver [NVIDIA 2015]. The performance of this solver is even worse than the MKL *dtsvb()* library function on the CPU in both single and double precision. The slowdown of cuSPARSE versus the MKL solver is only valid for the regime of short systems in the order of thousands of length. For larger systems this tendency changes for the advantage of the cuSPARSE library. Since the practical applications detailed in the introductory Section involves the solution of systems below the thousand size regime, systems above this limit are not the scope of the paper and we don't elaborate on these differences any deeper.

The transposition based Thomas solvers perform better than cuSPARSE by a factor of 2.5 – 3.1 in the case of transposition in shared memory and by 4.3 – 3.3 in the case of register shuffle depending on floating point precision and system size. The advantage of making extra effort for improving the cache line utilization is obvious. The achieved speedup compared to the naïve Thomas solver is 6.1 – 7.2 in the case of single and double precision respectively. The efficiency of the Thomas solver with transposition remains constant as the system size increases, ie. there is no significant fluctuation in performance.

One may notice that the padding is important for the register transposition so that aligned *float4* or *double2* loads may be done, which can not always be ensured. There is a factor 2 speed difference between the shared memory transposition and the register shuffle transposition in single precision. This difference is negligible in double precision since 64bit wide shared memory bank access is used to access the scratch memory instead of 32bit in the case of single precision. The wider bank access has been introduced in the Kepler architectures and it effectively doubles the bandwidth of the shared memory when 64bit access can be ensured. In case of shared memory a read throughput of 1211 GB/s and store throughput of 570 GB/s is measured with NVVP (NVIDIA Visual Profiler) on a  $256^3$  single precision cube domain. Also, the transposition using shared memory happens by reading 32bytes at a time, whereas in the case of the register shuffle based transposition 64 bytes are read from the global memory at a time.

The hybrid Thomas-PCR algorithm outperforms every solver in the X dimension in single precision. In double precision however the performance drops significantly beyond system size 512. The efficiency is due to the register blocking nature of the algorithm. Each system that is solved is completely resident in registers and therefore only the input data is read once and the output is stored once. This results in the minimum amount of data transfers and leads to the best possible performance.

## 7.2. Scaling in higher dimensions

Figure 12 shows the performance scaling in the Y dimension for single and double precision for all the architectures studied in this paper. Since the solution of tridiagonal systems with non-unit stride is not implemented in any library up to date, the higher dimensional execution time benchmarks do not contain any standard library execution times. Transposing the data structure would allow for the use of the standard solvers, but it has not been done for two reasons: 1) the efficiency of transposing the whole data

structure would involve further optimization of implementations attached to standard library code and the overall performance would be influenced by the extra programming effort; 2) even with a highly optimized transposition the overall execution time would be higher then in the X dimension.

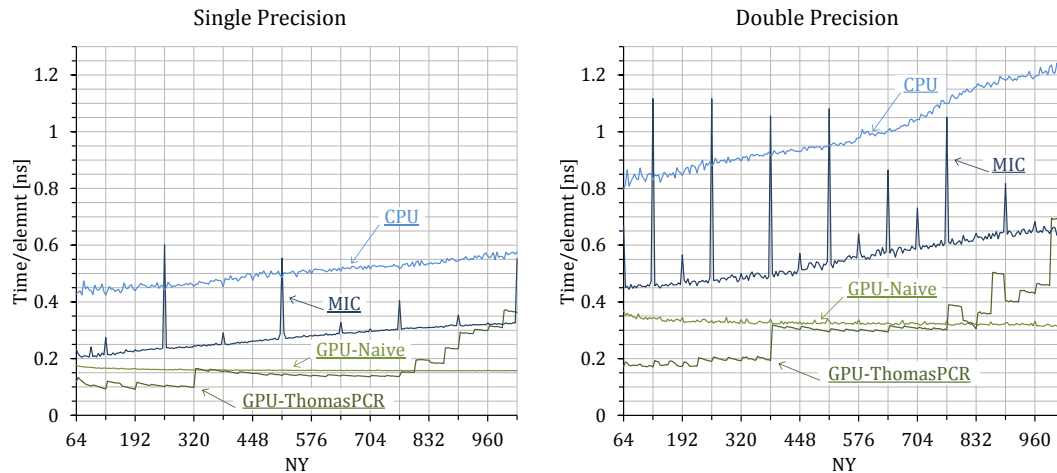


Fig. 12: Single and double precision execution time per grid element. 65536 tridiagonal systems with varying  $NY$  length along the Y dimension are solved. CPU, MIC and GPU execution times along with the comparison of the best solver on all three architectures are show.

**7.2.1. SIMD solvers.** The SIMD solvers rely on the regular Thomas algorithm. The detailed description is in Section 5. Both the CPU and the MIC SIMD architectures are able to utilize the 256 byte and 512 byte wide AVX and IMCI vector units to solve a tridiagonal system. This means that 8(4) and 16(8) systems can be solved by a single thread in parallel in single(double) precision on CPU and MIC respectively. The efficiency of the CPU relies on the cache and out-of-order execution performance of the architecture. The size of the temporary arrays for the systems solved in parallel is small enough to fit into the L1 cache and the L2 cache. The out-of-order execution is capable of hiding some of the 4 clock cycle latency of accessing these temporary in the L1 cache. On Figure 12 it can be seen that the CPU and MIC performance changes in parallel. Although the scaling should be constant, it is changing almost linearly in the case of single precision and super linearly in double precision. The reason for the latter is that the cores are running out of L1 cache and the L2 cache performance starts to dominate. Two arrays ( $c^*$  and  $d^*$ ) need to be cached for good performance. For instance, system size  $N = 1024$  needs  $1024 \text{ element} * 8 \text{ bytes/element} * 2 \text{ arrays} * 4 \text{ parallel systems} = 64 \text{ KB}$  to be cached, which is twice the size of the L1 cache. The single precision solver remains linear in this regime, because even 1024 long system fits into the L1 cache. The order of the two for loops (X and Z) iterating on the 65526 systems have set so that better data locality is achieved and thus better TLB hit rate is reached. The Thomas solver on the dual socket Xeon CPU is the slowest among all the solvers. The MIC architecture is by a factor 1.8 – 2 faster then the CPU. The SIMD architectures require aligned loads for best performance which can be ensured with padding otherwise the performance is hit by unaligned data loads and stores. The vector operations can still be performed, the non-alignment only hits data transfer.



**7.2.2. SIMT solvers.** The naïve GPU solver provides stable execution time and it is up to 3.6(3.8) times faster than the dual socket Xeon CPU and 2.1(2.5) times faster than the Xeon Phi in single(double) precision. The GPU implementation is capable of solving 32 systems within a warp using coalesced memory access. The performance is therefore predictable and very efficient. The only drawback of the solver is that it is moving more data than the Thomas-PCR hybrid solver (detailed in Section 2.3) which caches the data in registers. Therefore the Thomas-PCR hybrid algorithm is up to 1.5(1.8) faster than the naïve GPU solver in case of single(double) precision. Compared to the highly optimized two socket CPU solver the Thomas-PCR solver is 4.3(4.6) faster in single(double) precision. Compared to the highly optimized MIC the speedup is 2.2 – (2.5). These are significant differences which can be maintained until there is enough register to hold the values of the processed systems. Once the compiled code runs out of register, the effect of register spill becomes obvious, since the execution time jumps by more the 50% – this happens with system size 320 in single precision and 384 in double precision. In case of register spill the advantage of the Thomas-PCR over the naïve solver is negligible and above certain system size it is even worse.

## 8. CONCLUSION FOR SCALAR TRIDIAGONAL SOLVERS

In the past many algorithms have been introduced to solve a system of tridiagonal equations, but only a few of them took advantage of the fact that in certain cases (eg. an ADI solver) the problem to be solved contains a large number of small systems to solve and the access pattern of the system elements might change in data structures with 2 dimensions and above. It has been shown that in the  $X$  dimension the standard Thomas algorithm with modified access pattern using local data transposition on both CPU, MIC and GPU can outperform library quality, highly optimized code, such as: 1) *dtsvb()* MKL diagonally dominant solver running on multiple threads on a dual socket Xeon processor and 2) the PCR-CR hybrid algorithm implemented in the cuSPARSE library provided by NVIDIA. If the system size allows caching in registers, then a new proposed Thomas-PCR hybrid algorithm on the GPU can be used to solve the problem even more efficiently with a speedup of about 2 compared to the Thomas algorithm with local data transposition. It has been shown that in higher dimensions ( $Y$ ,  $Z$  etc.) the naïve solver allows for coalesced (or almost coalesced) access pattern and therefore there is no need for transposition and the performance is high. The Thomas-PCR for  $X$  dimension is modified to handle systems in higher dimensions by using more warps to load and store the data required in the computation. The register and shared memory pressure is higher in these cases and register spills occur above system size 320 in single and 384 in double precision. The Thomas algorithm performs better for above these system sizes. Figure 13 summarizes the execution times for the  $X$  and  $Y$  dimensions in the case of a  $240 \times 256 \times 256$  domain.

The conclusion is, that the Thomas algorithm with modified access pattern is advantageous up to relatively large system sizes in the order of thousands. The Thomas-PCR hybrid gives better performance in the  $X$  dimension in this regime. In the  $Y$  dimensions and above the Thomas-PCR is the best performing up to system size 320(384) in single(double) precision, but above this size the Thomas regains its advantage due to its simplicity.

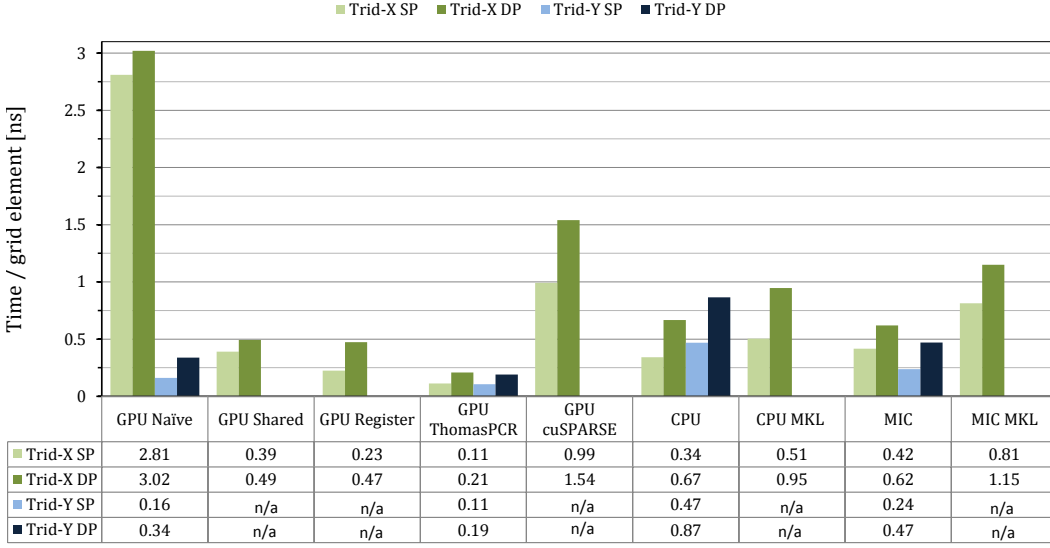


Fig. 13: Grid element execution times per grid element on a  $240 \times 256 \times 256$  cube domain.  $X$  dimension is chosen to be 240 to alleviate the effect of caching and algorithmic inefficiencies which occurs on sizes of power of 2 as it can be seen on Figures 10,11,12. Missing bars and n/a values are due to non-feasible or non-advantagous implementations in the  $Y$  dimension.

## BLOCK TRIDIAGONAL SOLVER ON SIMD AND SIMT ARCHITECTURES

### 1. INTRODUCTION

In many real-world CFD and financial applications the multidimensional PDEs have interdependent state variables. The state variable dependence creates a block structure in the matrix used in the implicit solution of the PDE. In certain cases these matrices are formed to be tridiagonal with block matrices in the diagonal and off-diagonals. The  $M^2$  block sizes are usually in the range of  $M = 2, \dots, 8$  and therefore the tridiagonal matrix takes the forms shown in Equation 4 and 5.

$$\mathbf{A}_i \mathbf{u}_{i-1} + \mathbf{B}_i \mathbf{u}_i + \mathbf{C}_i \mathbf{u}_{i+1} = \mathbf{d}_i \quad (4)$$

$$\begin{pmatrix} \mathbf{B}_0 & \mathbf{C}_0 & \mathbf{0} & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{A}_1 & \mathbf{B}_1 & \mathbf{C}_1 & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{0} & \mathbf{A}_2 & \mathbf{B}_2 & \mathbf{C}_2 & \cdots & \mathbf{0} \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \cdots & \mathbf{A}_{N-1} & \mathbf{B}_{N-1} \end{pmatrix} \begin{pmatrix} \mathbf{u}_0 \\ \mathbf{u}_1 \\ \mathbf{u}_2 \\ \vdots \\ \mathbf{u}_{N-1} \end{pmatrix} = \begin{pmatrix} \mathbf{d}_0 \\ \mathbf{d}_1 \\ \mathbf{d}_2 \\ \vdots \\ \mathbf{d}_{N-1} \end{pmatrix} \quad (5)$$

where  $\mathbf{A}_i, \mathbf{B}_i, \mathbf{C}_i \in \mathbb{R}^{M \times M}$ ,  $\mathbf{u}_i, \mathbf{d}_i \in \mathbb{R}^M$  and  $M \in [2, 8]$ .

Algorithms for solving block-tridiagonal system of equations have been previously developed in [Hirshman et al. 2010; Seal et al. 2013; van der Vorst 1987]. [Stone et al. 2011] gave a GPU based solution using the block PCR algorithm motivating their choice by the inherent parallelism given by the algorithm and the demand for high parallelism by the GPU. The overall arithmetic complexity of the PCR is known to be higher than that of the Thomas algorithm as it is detailed in Section 2. As many CFD applications consider block sizes of  $M = 2, \dots, 8$  the motivation of our work is to

make use of the computationally less expensive algorithm, namely the Thomas algorithm and exploit parallelism in the block-matrix operations. In other words, the overall arithmetic complexity is kept low by the Thomas algorithm and the parallelism is increased by the work-sharing threads which solve the block-matrix operations.

[Stone et al. 2011] used a highly interleaved SOA (Structure of Arrays) storage format to store the coefficients of the systems in order to achieve coalesced memory access pattern suitable for the GPU. In that approach the data on the host is stored in AOS (Array of Structures) format which had to be converted into SOA format to suit the needs of the GPU. In our approach the work-sharing and register blocking solution alleviates the need for AOS-SOA conversion, ie. the data access efficiency remains high.

## 2. BLOCK-THOMAS ALGORITHM

---

### ALGORITHM 12: Block Thomas algorithm

---

**Require:** *block\_thomas*(A, B, C, d)

```

Forward pass
1:  $C_0^* = B_0^{-1}C_0$ 
2:  $d_0^* = B_0^{-1}d_0$ 
3: for  $i = 1, \dots, N-1$  do
4:    $C_i^* = (B_i - A_i C_{i-1}^*)^{-1}C_i$ 
5:    $d_i^* = (B_i - A_i C_{i-1}^*)^{-1}(d_i - A_i d_{i-1}^*)$ 
6: end for
Backward pass
7:  $u_{N-1} = d_{N-1}^*$ 
8: for  $i = N-2, \dots, 0$  do
9:    $u_i = d_i^* - C_i^* u_{i+1}$ 
10: end for
11: return u

```

---

The block structure introduces matrix operations such as 1) block matrix inversion with  $O(M^3)$ ; 2) block matrix-matrix multiplication with  $O(M^3)$  and addition with  $O(M^2)$ ; 3) block matrix-vector multiplication with  $O(M^2)$ . On the other hand the data transfer is only  $O(M^2)$  per block matrix. As the matrix operations are dominated by  $O(M^3)$  versus the  $O(M^2)$  data transfer, the solution of the block tridiagonal system becomes compute limited as the block size increases. Once the data is read and stored in scratch memory, the cost of making the matrix operations is the bottleneck, both because arithmetic complexity and control flow complexity are significant. Let us define the  $O$  complexity in terms of block matrix operations with the arithmetic complexity states above, and define the total work as the complexity of solving a single system times the number of systems to be solved on a given number of parallel processors. When solving  $N$  long systems on  $N$  processors the Thomas algorithm has  $N O(N)$  total work complexity versus the  $N O(N \log N)$  total work complexity of the PCR algorithm. This significant difference establishes the use of the Thomas algorithm over the PCR in the block tridiagonal case.

The block Thomas algorithm is essentially the same as the scalar algorithm, assuming that scalar operations are exchanged with matrix operations. The lack of commutative property of the matrix multiplication, the order of these matrix operations have to be maintained throughout the implementation. See Algorithm 12 for details.

## 3. DATA LAYOUT FOR SIMT ARCHITECTURE

The data layout is a critical point of the solver, since this influences the effectiveness of the memory access pattern. Coalesced memory access is needed to achieve the best theoretical performance, therefore SOA data structures are used in many GPU appli-

cations. In this section an AOS data storage is presented in which blocks of distinct tridiagonal system are interleaved. Block coefficients  $A_n^p, B_n^p, C_n^p, d_n^p$  are stored in separate arrays. The data layout of  $A, B, C$  and  $C^*$  coefficient block arrays are the same. Within the array of blocks the leading dimension is the row of a block, ie. blocks are stored in row major format. The block of system  $p$  is followed by the block of system block  $p + 1$  in the array, ie. the blocks in array  $A$  are stored by interleaving the  $n$ th blocks of problems  $p = 0, \dots, P - 1$  in the way shown on Figure 14.

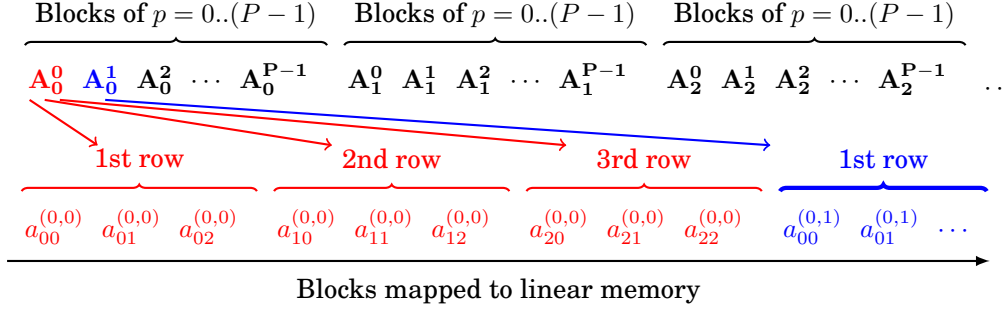


Fig. 14: Data layout within coefficient array  $A$ . This layout allows nearly coalesced access pattern and high cache hit rate when coalescence criteria is not met. Here  $A_n^p$  is the  $n$ th block (ie.  $n$ th block-row in the coefficient matrix) of problem  $p$  and  $p \in [0, P - 1]$ ,  $n \in [0, N - 1]$ . Notation for the scalar elements in the  $n$ th block of problem  $p$  are shown on Eq. (6). Bottom part of the figure shows the mapping of scalar values to the linear main memory.

$$A_n^p = \begin{pmatrix} a_{00}^{(n,p)} & a_{01}^{(n,p)} & a_{02}^{(n,p)} \\ a_{10}^{(n,p)} & a_{11}^{(n,p)} & a_{12}^{(n,p)} \\ a_{20}^{(n,p)} & a_{21}^{(n,p)} & a_{22}^{(n,p)} \end{pmatrix} \quad (6)$$

Vectors  $d, d^*$  and  $u$  are stored in a similar, interleaved fashion as depicted on Figure 14. Subvectors in array  $d$  are stored by interleaving the  $n$ th subvector of problems  $p = 0, \dots, P - 1$  in the way shown on Figure 15. The notation of the scalar values of  $d_n^p$  can be seen on Eq. (7).

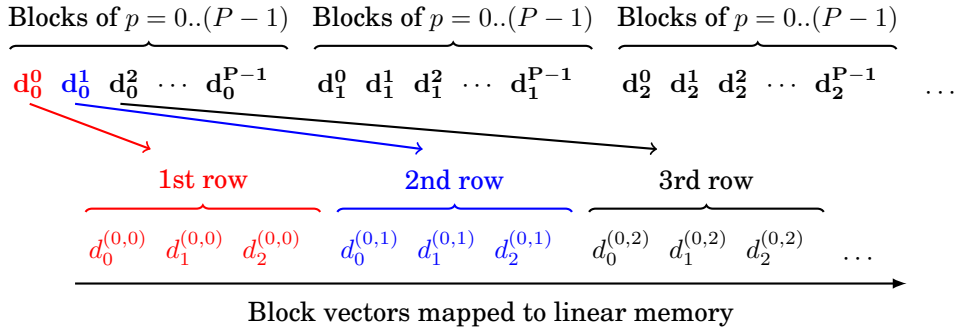


Fig. 15: Data layout of  $d$  as block vectors. This layout allows nearly coalesced access pattern and high cache hit rate when coalescence criteria is not met. Here  $d_n^p$  is the  $n$ th block (ie.  $n$ th block-row in the coefficient matrix) of problem  $p$  and  $p \in [0, P - 1]$ ,  $n \in [0, N - 1]$ . Notation for the scalar elements in the  $n$ th block of problem  $p$  are shown on Eq. (7). Bottom part of the figure shows the mapping of scalar values to the linear main memory.

$$\mathbf{d}_n^p = \left( d_0^{(n,p)} \quad d_1^{(n,p)} \quad d_2^{(n,p)} \right)^T \quad (7)$$

This data layout allows the threads to load data through texture cache from global memory with efficient access pattern and cache line utilization. A block is read in  $M$  steps. In each step a row of the block is read. The scalar values are stored in the registers as shown of Figure 16. In the second step threads read the second row in the same manner, and so on.

This storage format allows for high texture cache hit rate by assuring that most of the data that is read in within a 32 bytes cache line. Let us take the case of  $M = 3$  in single precision when reading  $A_0^0$ . In the first step only the first row of the block is read, ie. 12 bytes of the 32 bytes cache line is utilized. Once a texture cache line is loaded it will be reused immediately in the following few instructions or by the neighboring thread. In the next step, when the second row of the block is read, the cache line is already in the texture cache and this time 12 bytes are directly read from the texture cache. In the third step, when the third row is read, the 12 bytes are again in the cache, since 8 bytes out of the 12 bytes are in the same cache line as the first two rows, and the remaining 4 bytes are in the cache line read by the next group of 3 threads which read the block  $A_0^1$ . All this is done in parallel by  $\lfloor 32/3 \rfloor = 30$  threads. The total amount of data that is read by the warp (ie. 10 groups of threads) is  $10 \times 9 \times 4 = 360\text{bytes}$ , which fits into 12 cache lines. The probability of a cache line being evicted is low, since the cache lines are reused by the threads in the same warp. Since the sequence of instructions of loading a block doesn't contain data dependency, there is no reason for the scheduler to idle the active threads which started loading the data.

#### 4. COOPERATING THREADS FOR INCREASED PARALLELISM ON SIMT ARCHITECTURE

SIMT architectures are inherently sensitive to parallelism, ie. they need enough running threads to hide the latency of accessing the memory and filling up the arithmetic pipelines efficiently. Also, the implemented solver has to consider the constraints of the processor architecture being used. The Kepler GK110b architecture has 48KiB shared memory and 64K 32bit registers available for use within a thread block with up to 255 registers/thread. Computing block matrix operations with a single thread is not efficient due to the following reasons: 1) the limited number of registers and shared memory doesn't allow for temporary storage of block matrices and reloading data from global memory is costly; 2) in order to utilize coalesced memory access a high level of input data interleaving would be required, which is not useful in real application environment. As a consequence the problem would become memory bandwidth limited rather than compute limited. The workload of computing a single systems therefore needs to be shared among threads, so that block matrix data is distributively stored in the registers (and shared memory) of threads. This means that both workload and data storage is distributed among the cooperating threads.

We propose to use  $M$  number of threads to solve the  $M$  dimensional block matrix operations, so that every thread stores one column of a block and one scalar value of a vector that is being processed, see Figure 16 for details. Subsequent  $M$  threads are computing a system and a warp computes  $\lfloor 32/M \rfloor$  number of systems. This means that  $\lfloor 32/M \rfloor * M$  threads are active during the computation. The rest are idle. With this work distribution in the  $M = 2, \dots, 8$  range the worst case is  $M = 7$  when 4 out of 32 threads are inactive and thus 87.5% is the actual computation performance that can be reached.

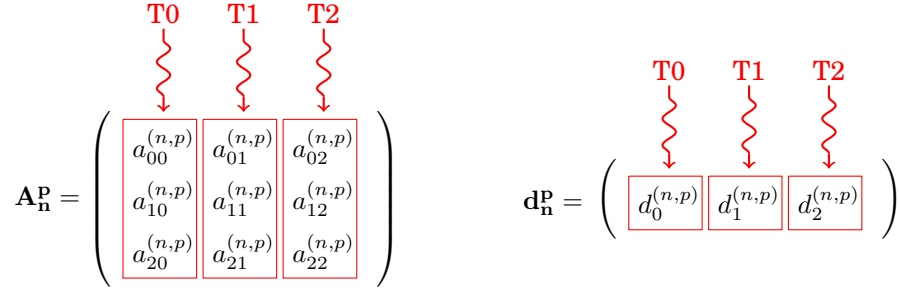


Fig. 16: Shared storage of blocks in registers. Arrows show load order. Thread  $T0, T1$  and  $T2$  stores the first, second and third columns of block  $A_n^p$  and first, second and third values of  $d_n^p$ .

The communication to perform matrix operations is carried out using either shared memory or register shuffles (`_shfl()` intrinsic). Just like in the case of the scalar solver, register and shared memory is not enough to store the intermediate  $C^*$ ,  $d^*$  block arrays and for this purpose local memory is used – this is an elegant way of getting coalesced memory access.

In the following the essential block matrix operations are discussed using Figure 16. Important to notice that all the implementation with register shuffle is written in a way that the local arrays storing block columns are held in registers. Two criteria need to be satisfied to achieve this, 1) local array indexing needs to be known at compile time and 2) local array size can not exceed a certain size defined internally in the compiler.

*Matrix-vector multiplication.* When performing matrix-vector multiplication, threads have the column of a block and the corresponding scalar value to perform the first step of matrix-vector multiplication, namely scalar-vector multiplication – this is done in parallel, independently by the threads. The second step is to add up the result vectors of the previous step. In this case threads need to share data, which can either be done through shared memory or using register shuffle instructions. In the case of *shared memory* the result is stored in shared memory. It is initialized to 0. In the  $m$ th step (where  $m = 0, \dots, M - 1$ ) thread  $m$  adds its vector to the shared vector and thread block synchronizes. In the case of *register shuffle* the multiplication is also done in  $M$  steps. In the  $m$ th step the  $M$  threads compute a scalar product of the  $m$ th row and shuffle the computed values around (in round-robin) to accumulate these values. The actual addition of the accumulation is done by the  $m$ th thread.

*Matrix-matrix multiplication.* Matrix-matrix multiplication needs to communicate the  $M \times M$  values of one of the blocks. This can either be done through shared memory or register shuffle. In the case of shared memory this approach uses  $M^2$  number of `_syncthreads()` which would suggest heavy impact on performance, but the results are still satisfying. The register shuffle approach doesn't require synchronization, thus it is expected to work faster than the shared memory approach.

*Gauss-Jordan elimination.* The solution of block systems in line 4 and 5 in Algorithm 12 is done by immediately solving these systems when performing a Gauss-Jordan elimination, ie. the systems are solved without composing the explicit inverse of matrix  $B_i - A_i C_{i-1}^*$ . The Gauss-Jordan elimination is computed cooperatively by using either shared memory or register shuffle instructions. Both versions have high compute throughput.

## 5. OPTIMIZATION ON SIMD ARCHITECTURE

The approach for efficient solution in the case of SIMD architectures is different in terms of data storage and execution flow. The data layout needs to be changed to get better data locality (for better cache performance) and second, each individual HT (Hyper Thread) thread computes an independent system without sharing the workload of the solution of one system. The latter consideration is the natural way of doing computation on multi-core systems. As each thread solves an independent system, best cache locality is reached when the blocks of array  $A$  are reordered to store the blocks  $A_0^p, A_1^p, \dots, A_N^p$  next to each other. Figure 17 depicts this in more details.

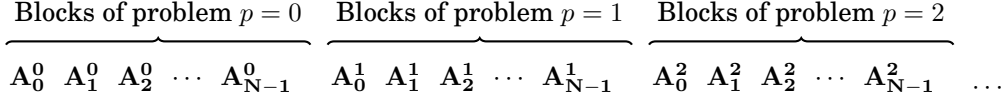


Fig. 17: Data layout within coefficient array  $A$  suited for better data locality on CPU and MIC. Here  $A_n^p$  is the  $n$ th block (ie.  $n$ th block-row in the coefficient matrix) of problem  $p$  and  $p \in [0, P - 1]$ ,  $n \in [0, N - 1]$ . Notation for the scalar elements in the  $n$ th block of problem  $p$  are shown on Eq. (6).

On these systems threads are heavy weight with out-of-order execution and have enough L1 cache to efficiently solve a block tridiagonal system. The C++ code is optimized to take advantage of the vector units. Loops are written so that the requirements for auto-vectorization are satisfied. Where needed `#pragma ivdep` or `#pragma simd` is used to help and force the compiler to vectorize loops. Dimension sizes are passed through template variables to make better use of compile time optimizations such as loop unrolling, vectorization and static indexing. Thread level parallelism is provided by OpenMP and threads solving independent systems are scheduled with guided scheduling to balance work load and give an overall better performance.

## 6. PERFORMANCE COMPARISON AND ANALYSIS

Performances of the implemented solvers are compared in terms of memory bandwidth, computation throughput in the case of GPUs and speedups of GPU and CPU implementations compared to the banded solver `gbsv()`. The size of the problem is always chosen to be such that it saturates both the CPU and the GPU with enough work, so that these architectures can provide the best performance, ie. as the block size is increased the length of the system to be solved is decreased so that the use of the available memory is kept close to maximum. Table III. and IV. show the selected length of a system ( $N$ ) and the number of systems to be solved ( $P$ ) respectively.

The Intel Math Kernel Library 11.2 Update 2 library [Intel 2015] was chosen with its `LAPACKE_{s,d}gbsv_work()` function to perform the banded solution. The library function was deploy using OpenMP to achieve the best performance MKL can provide. According to the MKL manual [Intel 2015] this routine solves for  $X$  the linear equations  $AX = B$ , where  $A \in \mathbb{R}^{n \times n}$  band matrix with  $kl$  subdiagonals and  $ku$  superdiagonals. The columns of matrix  $B$  are individual right-hand sides (RHS), and the columns of  $X$  are the corresponding solutions. This function is based on LU decomposition with partial pivoting and row interchanges. The factored form of  $A$  is then used to solve the system of equations  $AX = B$ . The solution of a system is done in two steps. First, a partial solve is done with the upper-triangular  $U$  matrix and then a solve with the lower-triangular  $L$  matrix is performed. This is an efficient approach when many right hand side exist. In the present case there is always one RHS, ie.  $X \in \mathbb{R}^{n \times 1}$  and  $B \in \mathbb{R}^{n \times 1}$ . As the systems which are solved are defined with diagonally

dominant matrices, pivoting is not performed during execution time. Moreover, the *work* version of the solver neglects any data validity check and thus provides a fair comparison. The scalar elements of diagonal block matrix arrays  $A_n^p$ ,  $B_n^p$  and  $C_n^p$  are mapped to band matrix  $A$  and the scalar elements of diagonal block vector arrays of  $d_n^p$  and  $u_n^p$  are mapped to  $X$  and  $B$  accordingly. The performance of the routine is expected to be high as the triangular sub-matrices are small enough to reside in L1, L2 or L3 cache. Comparing the solutions of the banded solver with the block tridiagonal solver by taking the differences between the corresponding scalar values shows MSE (Mean Square Error) in the order of  $10^{-4}$ . MSE does varies but stays in the order of  $10^{-4}$  as floating point precision, block size, system size or the number of systems is changed.

Figure 18 and Figure 19 show the effective bandwidth and computation throughput. The term effective is used to emphasize the fact that these performance figures are computed on basis of the actual data needed to be transferred and actual floating point arithmetic needed to be performed. Any caching and FMA influencing these figures are therefore implicitly included. In general the register shuffle based GPU solver outperforms the shared memory version, with one major exception in double precision with  $M = 8$  block size. In this case the register pressure is too high and registers get spilled into local memory.

One may notice that, as the block size increases the effective bandwidth decreases on Figure 18 and the effective compute throughput increases at the same time on Figure 19. Therefore it is implied that the problem is becoming compute limited rather than bandwidth limited as it is discussed in Section 2 due to the increasing difference between the  $O(M^3)$  compute and  $O(M^2)$  memory complexity of a single block.

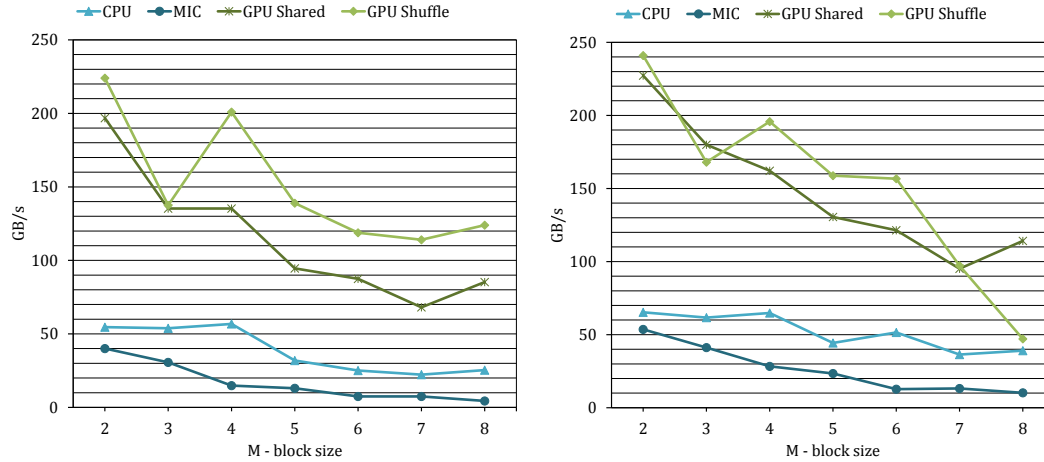


Fig. 18: Single (left) and double (right) precision effective bandwidth when solving block tridiagonal systems with varying  $M$  block sizes. Bandwidth is computed based on the amount of data to be transferred. Caching in L1 and registers can make this figure higher than the achievable bandwidth with a bandwidth test.

Execution time per block row are shown of Figure 20. The relative execution time measures the efficiency which is independent of problem size. The total execution time of the solver is divided by  $NP$ , where  $N$  is the length of a system and  $P$  is the number system that are solved. Execution time drastically increases for  $M = 8$  in double precision shuffle version. This is due to register spilling and local memory allocation, i.e. data can no longer fit into registers, therefore it is put into local memory, also the small local arrays that supposed to be allocated in registers due to compiler considerations



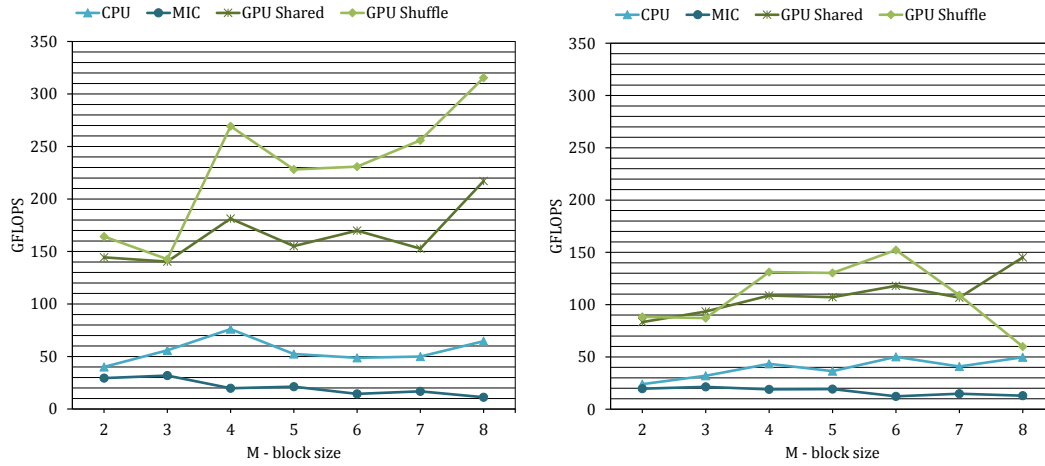


Fig. 19: Single (left) and double (right) precision effective computational throughput when solving block tridiagonal systems with varying  $M$  block sizes. GFLOPS is computed based on the amount of floating point arithmetic operations needed to accomplish the task. Addition, multiplication and division is considered as floating point arithmetic operations, ie. FMA (Fused Multiply and Add) is considered as two floating point instructions.

are put into local memory as well. This shows that the presented GPU approaches have very high efficiency for  $M = 2, \dots, 8$  block sizes.

An NVIDIA Tesla C2070 GPU card has been used to compare the results against the PCR solver presented in [Stone et al. 2011] where an NVIDIA Tesla C2050 was used. The two cards contain identical GPUs with different amount of global memory, see [NVIDIA 2010]. 3GB is available on the C2050 and 6 GB is available on the C2070. The execution times were measured for the single precision case on block size  $M = 4$ , on  $P = 512$  number of problems each of which is  $N = 128$  long. The block PCR based solver completes in 10.5ms and the (shared memory based) block Thomas solver completes in 2.42ms. This is a  $\times 4.3$  speed difference for the sake of the Thomas solver.  $\times 8.3$  and  $\times 9.8$  improvement is achieved if the execution time per block metrics are compared and the number of systems to be solved is increased to  $P = 4096$  or  $P = 32768$ .

The SIMD solution presented in the paper performs well on the CPU, but performs poorly on the MIC architecture. On both architectures the compute intensive loops were vectorized as it is reported by the Intel compiler. Both the MKL banded solver and the presented block tridiagonal solver run more efficiently on the CPU.

Figure 21 presents the speedup of the block-tridiagonal solvers on GPU over the MKL banded solvers and the CPU and MIC based block tridiagonal solvers. This proves the benefit of the presented GPU based solutions. Also, the highly efficient CPU and MIC implementations show the benefit of using a block tridiagonal solver over a banded solver for the range of block sizes  $M = 2, \dots, 8$ .

Significant speedup against the CPU MKL banded solver is reached with the GPU-based solver, up to  $\times 27$  in single and  $\times 12$  in double precision. The multi-threaded CPU code provides  $\times 2 - 6$  speedup in single and  $\times 1.5 - 3$  speedup in double precision. The multi-threaded MIC performance of the block solver is better than the CPU MKL, but the MKL banded solver perform poorly on the MIC. The efficiency of the CPU relies on the cache and out-of-order execution performance of the architecture. The size of the temporary blocks and arrays of blocks is small enough to fit into the L1 and L2 cache. The out-of-order execution is capable of hiding some of the 4 clock cycle latency of accessing these temporary data structures in the L1 cache. As the MIC lacks the

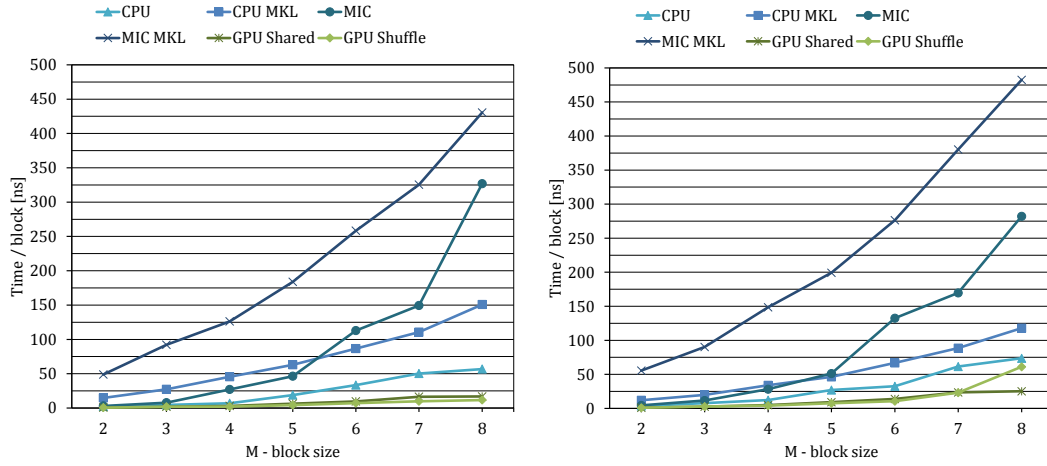


Fig. 20: Single (left) and double (right) precision block-tridiagonal solver execution time per block matrix row for varying block sizes.

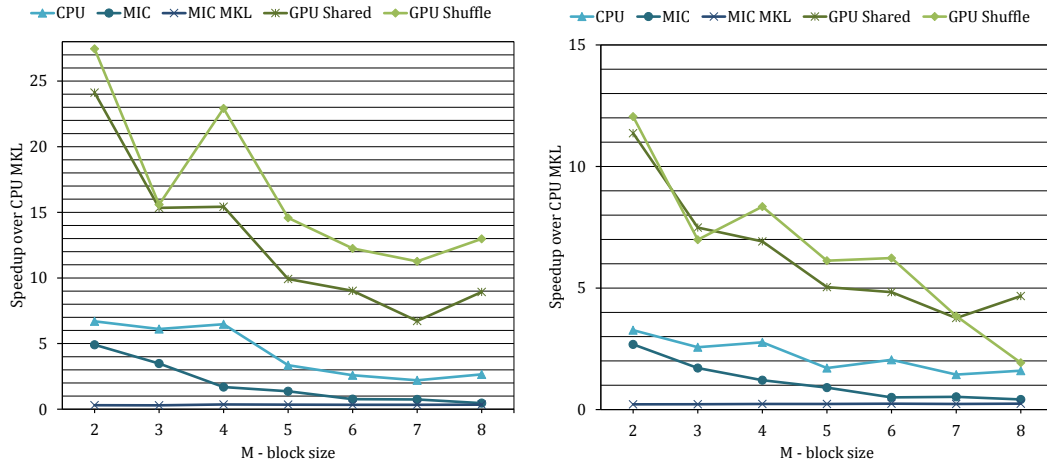


Fig. 21: Single (left) and double (right) precision block tridiagonal solver speedup over the CPU MKL banded solver.

out-of-order execution and the cache size per thread is much smaller the performance is also worse than the CPU. Moreover, the *gbsv()* banded solver of the MKL library does not perform well on the MIC architecture as it shows 3 times lower performance than the CPU MKL version.

The advantage of doing block tridiagonal solve in the range of  $M = 2, \dots, 8$  instead of a banded solve is obvious. It is important to note that, as the block size  $M$  increases, the computations involved in performing block-matrix operations make the problem compute limited instead of memory bandwidth limited. The total execution time of a computing a systems is composed of loading the blocks of data and performing matrix linear algebra on the blocks. The former one scales as  $O(NM^2)$  and the latter one as  $O(NM^3)$ , where  $N$  is the system size and  $M$  is the block size. As the block size increases the computational part becomes more dominant in the execution time and the memory access time becomes less significant. This can be read from the Figures 18 and 19 – as  $M$  increases, the bandwidth decreases and the GFLOPS increases.

## 7. CONCLUSION FOR BLOCK TRIDIAGONAL SOLVERS

In the present paper it has been shown that solving block-tridiagonal systems with  $M = 2, \dots, 8$  block sizes with the Thomas algorithm pays off over the Intel MKL banded solver. It has been shown that the advantage of the block Thomas algorithm is the computational complexity over the CR or PCR algorithm and that parallelism can be increased by exploiting the properties of the block matrix operations. The superior performance of the GPU relies on the low arithmetic complexity of the Thomas algorithm and the efficiency of the parallel block matrix operations allowed by the work sharing and the register blocking capabilities of the GPU threads. Since the work complexity (ie. number of block-matrix operations) of the CR/PCR algorithms are significantly higher than the Thomas algorithm, the CR/PCR has no advantage in block-tridiagonal solvers. Significant speedup is reached with the GPU-based solver with up to  $\times 27$  in single and  $\times 12$  in double precision. The multi-threaded CPU code provides  $\times 2 - 6$  speedup in single and  $\times 1.5 - 3$  times speedup in double precision against the MKL banded solver.

## APPENDIX

### A. HARDWARE AND SOFTWARE

The most salient properties of the CPU, MIC and GPU used in the present study are shown in Table II. The Intel Composer XE Suite 2015.2.164 with compiler version 15.0.2 compiler and Intel Math Kernel Library version 11.2 Update 2 for Linux was used to perform the benchmarks. For the NVIDIA GPU benchmark the CUDA 7.0 runtime and driver version 346.46 was used.

Table II: Details of Intel Xeon Sandy Bridge server processor, Intel Xeon Phi coprocessor and the NVIDIA Tesla GPU card. \*Estimates are based on [Saini et al. ], [Wong et al. 2010] and [amd 2013]. Both Intel architectures have 8-way, shared LLC with ring-bus topology. HT - Hyper Thread, MM - Multi Media, RO - Read Only

	Intel Xeon E5-2680	Intel Xeon Phi 5110P	NVIDIA K40m
Microarchitecture	Sandy Bridge	Knights Corner	GK110b
No. sockets	2	1	1
CPU Core / SMX unit	2x8	59 (+1 for OS)	15
HT / core or CUDA core/SMX	2	4	192
MM register width	256 bits	512 bits	32 bits
Registers / thread	16 YMM	32 ZMM	255
L1 data cache / core	32 KB	32 KB	64 KB + 48 KB RO
L2 data cache / core	256 KB	30 MB	1.5 MB
L3 data cache / socket	2x20 MB	-	-
Cache line	64 bytes	64 bytes	32 bytes
Cache latency L1/L2/L3*	4/11/21*	3/22/-*	38/-/ - *
Virtual page size	4KB-2MB	2MB	4KB*
Clock rate	2.7 (3.5 Turbo) GHz	1053 MHz	745 MHz
fp32 perf.	2x172.8 GFLOPS	2.022 TFLOPS	4.29 TFLOPS
fp64 perf.	2x86.4 GFLOPS	1.011 TFLOP	1.43 TFLOPS
Installed memory	64 GB DDR3	8 GB GDDR5	12 GB GDDR5
Memory bandwidth	2x51.2 GB/s	320 GB/s	288 GB/s
PCI bus	PCI-E x40 Gen3	PCI-E x16 Gen2	PCI-E x16 Gen3
TDP / Cooling	2x130 W / Active	225W / Passive	235 W / Passive
Recommended price	2x1727 USD	2649.00	5499 USD

### B. SYSTEM SIZE CONFIGURATION FOR BENCHMARKING BLOCK TRIDIAGONAL SOLVERS.

Table III: Parameter  $N$  - Length of a system used for benchmarking a processor architecture and solver. The length of the system is chosen such that the problem fits into the memory of the selected architecture.

	Block size	2	3	4	5	6	7	8
SP	CPU	128	128	128	128	128	96	96
	CPU MKL 11.2.2	128	128	128	128	128	96	96
	MIC	128	128	128	128	128	96	96
	MIC MKL 11.2.2	128	128	128	128	128	96	96
	GPU Shuffle	128	128	128	128	128	96	96
	GPU Shared	128	128	128	128	128	96	96
DP	CPU	128	128	128	128	128	96	96
	CPU MKL 11.2.2	128	128	128	128	128	96	96
	MIC	128	128	128	128	128	64	64
	MIC MKL 11.2.2	64	64	64	64	64	48	48
	GPU Shuffle	128	128	128	128	128	96	96
	GPU Shared	128	128	128	128	128	96	96

Table IV: Parameter  $P$  - Number of systems used for benchmarking a processor architecture and solver. The number of systems is chosen such that the problem fits into the memory of the selected architecture.

	Block size	2	3	4	5	6	7	8
SP	CPU	65536	65536	65536	65536	65536	32768	32768
	CPU MKL 11.2.2	65536	65536	65536	65536	65536	32768	32768
	MIC	32768	32768	32768	32768	32768	16384	16384
	MIC MKL 11.2.2	32768	32768	32768	32768	32768	16384	16384
	GPU Shuffle	65536	65536	65536	65536	65536	32768	32768
	GPU Shared	65536	65536	65536	65536	65536	32768	32768
DP	CPU	65536	65536	65536	65536	65536	32768	32768
	CPU MKL 11.2.2	65536	65536	65536	65536	65536	32768	32768
	MIC	32768	32768	32768	32768	32768	16384	16384
	MIC MKL 11.2.2	32768	32768	32768	32768	32768	16384	16384
	GPU Shuffle	65536	65536	65536	65536	65536	32768	32768
	GPU Shared	65536	65536	65536	65536	65536	32768	32768

## ELECTRONIC APPENDIX

The electronic appendix for this article can be accessed in the ACM Digital Library.

## ACKNOWLEDGMENTS

The research at the Oxford e-Research Centre has been partially supported by the ASEArch project on Algorithms and Software for Emerging Architectures, funded by the UK Engineering and Physical Sciences Research Council. The research has also been supported by the TÁMOP-4.2.1./B-11/2/KMR-2011-002, TÁMOP - 4.2.2./B-10/1-2010-0014 projects at PPCU - University of National Excellence. The authors would like to acknowledge the use of the University of Oxford Advanced Research Computing (ARC) facility in carrying out this work. We are thankful for the support of Péter Szolgay at PPCU-FIT, Hungary.

## REFERENCES

2013. AMD64 Architecture Programmer's Manual Volumes 1-5. <http://developer.amd.com/resources/documentation-articles/developer-guides-manuals/#manuals>. (2013).
- Stefan Bondeli. 1991. Divide and conquer: a parallel algorithm for the solution of a tridiagonal linear system of equations. *Parallel Comput.* 17, 45 (1991), 419 – 434. DOI: [http://dx.doi.org/10.1016/S0167-8191\(05\)80145-0](http://dx.doi.org/10.1016/S0167-8191(05)80145-0)
- Li-Wen Chang, John A. Stratton, Hee-Seok Kim, and Wen-Mei W. Hwu. 2012. A Scalable, Numerically Stable, High-performance Tridiagonal Solver Using GPUs. In *Proceedings of the International Conference on High Performance Computing, Networking, Storage and Analysis (SC '12)*. IEEE Computer Society Press, Los Alamitos, CA, USA, Article 27, 11 pages. <http://dl.acm.org/citation.cfm?id=2388996.2389033>
- I.J.D. Craig and A.D. Sneyd. 1988. An alternating-direction implicit scheme for parabolic equations with mixed derivatives. *Computers and Mathematics with Applications* 16, 4 (1988), 341 – 350. DOI: [http://dx.doi.org/10.1016/0898-1221\(88\)90150-2](http://dx.doi.org/10.1016/0898-1221(88)90150-2)
- Duy M. Dang, Christina Christara, and Kenneth R. Jackson. 2010. Parallel Implementation on GPUs of ADI Finite Difference Methods for Parabolic PDEs with Applications in Finance. *Social Science Research Network Working Paper Series* (03 April 2010). <http://ssrn.com/abstract=1580057>
- Craig C. Douglas, Sachit Malhotra, and Martin H. Schultz. 1998. Parallel Multigrid with ADI-like Smoothers in Two Dimensions. (1998).
- J. Douglas and H. H. Rachford. 1956. On the numerical solution of heat conduction problems in two and three space variables. *Transaction of the American Mathematical Society* 82 (1956), 421–489.
- Jr. Douglas, Jim and JamesE. Gunn. 1964. A general formulation of alternating direction methods. *Numer. Math.* 6, 1 (1964), 428–453. DOI: <http://dx.doi.org/10.1007/BF01386093>
- B. Dring, M. Fourni, and A. Rigal. 2014. High-Order ADI Schemes for Convection-Diffusion Equations with Mixed Derivative Terms. In *Spectral and High Order Methods for Partial Differential Equations - ICOSAHOM 2012*, Mejdi Azaez, Henda El Fekih, and Jan S. Hesthaven (Eds.). Lecture Notes

- in Computational Science and Engineering, Vol. 95. Springer International Publishing, 217–226. DOI: [http://dx.doi.org/10.1007/978-3-319-01601-6\\_17](http://dx.doi.org/10.1007/978-3-319-01601-6_17)
- Walter Gander and Gene H. Golub. 1997. Cyclic Reduction - History and Applications. In *Proceedings of the Workshop on Scientific Computing*.
- S.P. Hirshman, K.S. Perumalla, V.E. Lynch, and R. Sanchez. 2010. BCYCLIC: A parallel block tridiagonal matrix cyclic solver. *J. Comput. Phys.* 229, 18 (2010), 6392 – 6404. DOI: <http://dx.doi.org/10.1016/j.jcp.2010.04.049>
- Intel 2012a. *Intel 64 and IA-32 Architectures Optimization Reference Manual*. Intel. <http://www.intel.com/content/dam/doc/manual/64-ia-32-architectures-optimization-manual.pdf>
- Intel 2012b. *Intel Xeon Phi Coprocessor Instruction Set Architecture Reference Manual*. Intel. <https://software.intel.com/sites/default/files/forums/278102/327364001en.pdf>
- Intel. 2015. *Math Kernel Library*. <http://software.intel.com/en-us/articles/intel-mkl/>
- Samir Karaa and Jun Zhang. 2004. High order {ADI} method for solving unsteady convectiondiffusion problems. *J. Comput. Phys.* 198, 1 (2004), 1 – 9. DOI: <http://dx.doi.org/10.1016/j.jcp.2004.01.002>
- Nathan Mattor, Timothy J. Williams, and Dennis W. Hewett. 1995. Algorithm for solving tridiagonal matrix problems in parallel. *Parallel Comput.* 21, 11 (1995), 1769 – 1782. DOI: [http://dx.doi.org/10.1016/0167-8191\(95\)00033-0](http://dx.doi.org/10.1016/0167-8191(95)00033-0)
- NVIDIA 2010. *TESLA C2050 / C2070 GPU Computing Processor*. NVIDIA. [http://www.nvidia.com/docs/IO/43395/NV\\_DS\\_Tesla\\_C2050\\_C2070\\_jul10\\_lores.pdf](http://www.nvidia.com/docs/IO/43395/NV_DS_Tesla_C2050_C2070_jul10_lores.pdf)
- NVIDIA 2015. *CUDA C Programming Guide*. NVIDIA. <http://docs.nvidia.com/cuda/cuda-c-programming-guide/#axzz3aTPUq4Jo>
- NVIDIA 2015. *CUSPARSE LIBRARY v7.0*. NVIDIA. <http://docs.nvidia.com/cuda/cusparse/index.html#axzz3aTPUq4Jo>
- D. W. Peaceman and Jr. Rachford, H. H. 1955. The Numerical Solution of Parabolic and Elliptic Differential Equations. *J. Soc. Indust. Appl. Math.* 3, 1 (1955), pp. 28–41. <http://www.jstor.org/stable/2098834>
- William H. Press, Saul A. Teukolsky, William T. Vetterling, and Brian P. Flannery. 2007. *Numerical Recipes 3rd Edition: The Art of Scientific Computing* (3 ed.). Cambridge University Press, New York, NY, USA.
- Thomas H. Pulliam. 1986. Implicit solution methods in computational fluid dynamics. *Applied Numerical Mathematics* 2, 6 (1986), 441 – 474. DOI: [http://dx.doi.org/10.1016/0168-9274\(86\)90002-4](http://dx.doi.org/10.1016/0168-9274(86)90002-4)
- I. Reguly, E. Laszlo, G. Mudalige, and M. Giles. 2014. Vectorizing Unstructured Mesh Computations for Many-core Architectures. (2014).
- Subhash Saini, Johnny Chang, and Haoqiang Jin. *Performance Evaluation of the Intel Sandy Bridge Based NASA Pleiades Using Scientific and Engineering Applications*. Technical Report. NASA Advanced Supercomputing Division, NASA Ames Research Center, Moffett Field, California 94035-1000, USA.
- Nikolai Sakharlykh. 2009. Tridiagonal Solvers on the GPU and Applications to Fluid Simulation. Presented at the GPU Technology Conference, San Jose, CA. [http://www.nvidia.com/content/gtc/documents/1058\\_gtc09.pdf](http://www.nvidia.com/content/gtc/documents/1058_gtc09.pdf)
- Nikolai Sakharlykh. 2010. Efficient Tridiagonal Solvers for ADI methods and Fluid Simulation. Presented at the GPU Technology Conference, San Jose, CA. <http://on-demand.gputechconf.com/gtc/2010/presentations/S12015-Tridiagonal-Solvers-ADI-Methods-Fluid-Simulation.pdf>
- Sudip K. Seal, Kalyan S. Perumalla, and Steven P. Hirshman. 2013. Revisiting parallel cyclic reduction and parallel prefix-based algorithms for block tridiagonal systems of equations. *J. Parallel and Distrib. Comput.* 73, 2 (2013), 273 – 280. DOI: <http://dx.doi.org/10.1016/j.jpdc.2012.10.003>
- G Spaletta and D.J Evans. 1993. The parallel recursive decoupling algorithm for solving tridiagonal linear systems. *Parallel Comput.* 19, 5 (1993), 563 – 576. DOI: [http://dx.doi.org/10.1016/0167-8191\(93\)90006-7](http://dx.doi.org/10.1016/0167-8191(93)90006-7)
- Christopher P Stone, Earl PN Duque, Yao Zhang, David Car, John D Owens, and Roger L Davis. 2011. GPGPU parallel algorithms for structured-grid CFD codes. In *Proceedings of the 20th AIAA Computational Fluid Dynamics Conference*, Vol. 3221.
- Harold S. Stone. 1973. An Efficient Parallel Algorithm for the Solution of a Tridiagonal Linear System of Equations. *J. ACM* 20, 1 (Jan. 1973), 27–38. DOI: <http://dx.doi.org/10.1145/321738.321741>
- L. H. Thomas. 1949. *Elliptic Problems in Linear Differential Equations over a Network*. Technical Report. Columbia University.
- Henk A van der Vorst. 1987. Large tridiagonal and block tridiagonal linear systems on vector and parallel computers. *Parallel Comput.* 5, 12 (1987), 45 – 54. DOI: [http://dx.doi.org/10.1016/0167-8191\(87\)90005-6](http://dx.doi.org/10.1016/0167-8191(87)90005-6)
- Proceedings of the International Conference on Vector and Parallel Computing-Issues in Applied Research and Development.
- H. H. Wang. 1981. A Parallel Method for Tridiagonal Equations. *ACM Trans. Math. Softw.* 7, 2 (June 1981), 170–183. DOI: <http://dx.doi.org/10.1145/355945.355947>

- H. Wong, M.-M. Papadopoulou, M. Sadooghi-Alvandi, and A. Moshovos. 2010. Demystifying GPU microarchitecture through microbenchmarking. In *Performance Analysis of Systems Software (ISPASS), 2010 IEEE International Symposium on*. 235–246. DOI: <http://dx.doi.org/10.1109/ISPASS.2010.5452013>
- Yao Zhang, Jonathan Cohen, and John D. Owens. 2010. Fast Tridiagonal Solvers on the GPU. *SIGPLAN Not.* 45, 5 (Jan. 2010), 127–136. DOI: <http://dx.doi.org/10.1145/1837853.1693472>

Received October 2014; revised April 2015; accepted September 2015