

Block-asynchronous multigrid smoothers for GPU-accelerated systems

H. Anzt, J. Dongarra, M. Gates, S. Tomov

No. 2011-15

Preprint Series of the Engineering Mathematics and Computing Lab (EMCL)





Preprint Series of the Engineering Mathematics and Computing Lab (EMCL)
ISSN 2191-0693
No. 2011-15

Impressum

Karlsruhe Institute of Technology (KIT)
Engineering Mathematics and Computing Lab (EMCL)

Fritz-Erler-Str. 23, building 01.86
76133 Karlsruhe
Germany

KIT – University of the State of Baden Wuerttemberg and
National Laboratory of the Helmholtz Association

Published on the Internet under the following Creative Commons License:
<http://creativecommons.org/licenses/by-nc-nd/3.0/de> .



www.emcl.kit.edu

Block-asynchronous multigrid smoothers for GPU-accelerated systems

Hartwig Anzt^{*1}, Stanimire Tomov^{†2}, Mark Gates^{‡2}, and Jack Dongarra^{§2}

¹*Karlsruhe Institute of Technology (KIT), Institute for Applied and Numerical Mathematics 4, Karlsruhe, Germany*

²*University of Tennessee, Innovative Computing Lab (ICL), Knoxville, USA*

Abstract This paper explores the need for asynchronous iteration algorithms as smoothers in multigrid methods. The hardware target for the new algorithms is top-of-the-line, highly parallel hybrid architectures – multicore-based systems enhanced with GPGPUs. These architectures are the most likely candidates for future high-end supercomputers. To pave the road for their efficient use, challenges related to the established notion that “data movement, not FLOPS, is the bottleneck to performance” must be resolved. Our work is in this direction – we designed block-asynchronous multigrid smoothers that perform more flops in order to reduce synchronization, and hence data movement. We show that the extra flops are done for “free,” while synchronization is reduced and the convergence properties of multigrid with classical smoothers like Gauss-Seidel are preserved.

Keywords Block-asynchronous Iteration, Multigrid Smoothers, GPU

1 Introduction

Classical relaxation methods such as Gauss-Seidel and Jacobi require a synchronization between each iteration. This implies a severe restriction for parallel implementations. An asynchronous iteration method removes this synchronization barrier, updating components using the latest available values. It allows a large freedom in the update order and the number of updates per component, while every component update uses the latest available values for the other components. In the end the obtained algorithm is neither deterministic nor does it imply convergence for all systems that can be solved by the classical Jacobi approach, in fact it requires the linear equation system to fulfill additional conditions. While due to the poor convergence rate they may seem to be very unattractive from the mathematical point of view, the asynchronous iteration is, in contrast to most other iterative methods, able to exploit the high computational power of modern hardware platforms. The high parallelization potential and the high tolerance to communication and synchronization latencies make them perfect candidates for computing systems that

consist of a high number of parallel computing cores, that are eventually even located in different devices, and provide excellent performance as long as they run independently, but severely suffer from data transfers and synchronization points. For linear equation systems fulfilling the required additional conditions, the asynchronous iteration is often able to overcompensate the inferior convergence rate by leveraging the computational power of the available hardware resources. This leads to the fact, that for many problems the asynchronous iteration outperforms the classical synchronized relaxation methods like Gauss-Seidel and Jacobi [2].

While iterative methods based on matrix splitting are nowadays seldom applied directly to solve a problem, they are often used in multigrid methods, to smooth the error terms related to the large eigenvalues on the distinct grid levels [11] [17]. The superiority of an asynchronous iteration method has been shown for many linear equation problems, but it is an open question whether it is suitable to replace the classical smoothers in multigrid methods. Since the parallelization of smoothers is usually crucial when optimizing multigrid solvers [7], it may have a huge impact on the overall multigrid performance.

Targeting this topic we split this paper into different parts: We start with an introductory section providing the mathematical background of geometric multigrid methods and asynchronous iteration. In the second part we give details about the implementation we used for the numerical tests and the hardware configuration. We then report numerical results of multigrid implementations using asynchronous iteration smoothers and compare them to Gauss-Seidel schemes. In the last section we conclude and provide ideas about which topics could be interesting to address in this research field.

2 Mathematical Background

2.1 Multigrid Methods

Multigrid methods are a type of error correction methods that attempt to find a solution approximation by using a sequence of problems that are similar with respect to their structure, but differ in the successively decreasing dimen-

^{*}Email: hartwig.anzt@kit.edu

[†]Email: tomov@eecs.utk.edu

[‡]Email: mgates3@eecs.utk.edu

[§]Email: dongarra@eecs.utk.edu

sion. [17] [20] [21] They are usually applied to problems occurring in the field of finite element or finite difference discretizations of partial differential equations. In this case, different discretizations with decreasing dimension of the continuous problem can be used for the sequence of discretized problems, as in Figure 1. This enables splitting the approximation error into high and low frequency terms that can then be treated with different efficiency on the distinct grid levels. A basic multigrid algorithm is given in Algorithm 1.

Using an optimal combination of problem sequence and operators, one can obtain a solver with optimal complexity $\mathcal{O}(n)$. Hence, multigrid solvers are among the most efficient solvers for the discretized partial differential equations.

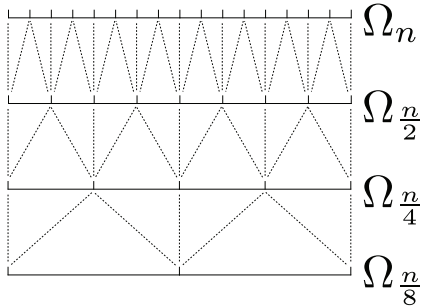


Figure 1: Sequence of successively coarser grids.

- 1: $\text{MG}(x_l, b_l, l)$
- 2: **if** $l = 0$ **then**
- 3: Solve $A_l x_l = b_l$ {exact solution on coarsest grid}
- 4: **else**
- 5: $x_l = \mathcal{S}_l^{v_1}(x_l, b_l)$ {pre-smoothing}
- 6: $r_{l-1} = r_l^{l-1}(b_l - A_l x_l)$ {restriction}
- 7: $v_{l-1} = 0$
- 8: **for** $j = 0; j < \gamma; j++$ **do**
- 9: $\text{MG}(v_{l-1}, r_{l-1}, l-1)$ {coarse grid correction}
- 10: **end for**
- 11: $x_l = x_l + p_{l-1}^l v_{l-1}$ {prolongation of coarse grid correction}
- 12: $x_l = \mathcal{S}_l^{v_2}(x_l, b_l)$ {post smoothing}
- 13: **end if**

Algorithm 1: Basic multigrid method [17].

In the case of using only two grids and solving the error correction equation exact on the grid $\Omega_{\frac{n}{2}}$, the method is called the “Two-Grid Iteration Method.” Usually, the process is recursively applied to successively coarser grids, creating a “Multigrid Method.” In this case there exist different schemes for how to organize the multigrid process. The structure is determined by the number R of error correction computations on every grid level. Implementations usually choose $R = 1$ or $R = 2$. The resulting multigrid iteration schemes are called “V-Cycle” and “W-Cycle,” respectively. While the V-cycle is usually very efficient in

terms of computational cost, it can be unstable with respect to the properties of the problem. The W-cycle is computationally more expensive, but at the same time more robust in terms of the problem. For $R \geq 3$, the multigrid iteration method becomes inefficient. A trade off between V-cycle and W-cycle is the “F-Cycle,” shown in Figure 2.

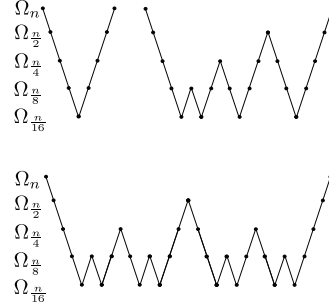


Figure 2: Visualizing V-cycle, W-cycle and F-cycle.

2.2 Smoother in Multigrid Methods

One critical component of multigrid methods is the smoother. Usually, a simple relaxation method such as Gauss-Seidel or Jacobi is used for pre- and post-smoothing the solution approximations on the distinct grid levels. The idea of applying the smoother is to make the underlying error smooth so that it can be approximated efficiently on a coarser grid. From the analytical point of view, if the error is expressed in terms of the eigenvectors of the system, the smoother must eliminate the error associated with the eigenvectors having large eigenvalues, while the coarse-grid correction eliminates the remaining error associated with eigenvectors having small eigenvalues [17].

The best smoothers, such as Gauss Seidel, usually do not parallelize well. Therefore, much effort is put into developing parallel smoothers that scale on multicore architectures. The main approach is to use a set of local smoothers that exchange boundary values in a Jacobi-like manner [12], [3]. The performance of these hybrid smoothers may then be enhanced furthermore by using weights [22].

Still, the synchronization necessary to exchange boundary values may be detrimental to the performance on highly parallel architectures. Therefore, it is worthwhile to consider a block-asynchronous iteration, which lacks any synchronization and therefore scales optimally on any architecture, for the smoother in multigrid methods.

2.3 Asynchronous Iteration

The motivation for an asynchronous iteration is modern hardware, which provides a large number of cores that achieve excellent performance when running in parallel, but suffer when synchronizing or exchanging data. Therefore, algorithms that lack any synchronization would achieve outstanding performance on these devices, while most of the numerical algorithms are poorly parallel and

require regular data exchange. For computing the next iteration in relaxation methods, one usually requires the latest values of all components. For some algorithms, e.g., Gauss-Seidel [13], even the already computed values of the current iteration step are used. This requires a strict order of the component updates, limiting the parallelization potential to a stage, where a component cannot be updated several times before all the other components are updated.

If this order is not adhered to, i.e., the individual components are updated independently and without consideration of the current state of the other components, the resulting algorithm is called a chaotic or asynchronous iteration method. In the 1970s, Chazan and Miranker analyzed some basic properties of these methods and established convergence theory [9] [18] [4] [8] [10]. For the last 30 years, these algorithms went out of favor due to the superior convergence properties of synchronized iteration methods like the Krylov subspace methods. Today, due to the complexity of heterogeneous hardware platforms and the large number of computing units in parallel devices like GPUs, these schemes may become interesting again for applications like multigrid methods, where highly parallel smoothers are required on the distinct grid levels. While traditional smoothers like the sequential Gauss-Seidel obtain their efficiency from their fast convergence, an asynchronous iteration scheme may compensate its inferior convergence behavior by superior scalability.

The chaotic or asynchronous relaxation scheme defined by Chazan and Miranker [9] can be characterized by two functions, an update function $u(\cdot)$ and a shift function $s(\cdot, \cdot)$. For each non-negative integer v , the component of the solution approximation x that is updated at step v is given by $u(v)$. For the update at step v , the m^{th} component used in this step is $s(v, m)$ steps back. All the other components are kept the same. This can be expressed as:

$$x_l^{v+1} = \begin{cases} \sum_{m=1}^N b_{l,m} x_m^{v-s(v,m)} + d_l & \text{if } l = u(v) \\ x_l^v & \text{if } l \neq u(v). \end{cases} \quad (1)$$

Furthermore, the following conditions can be defined to guarantee the well-posedness of the algorithm [16]:

1. The update function $u(\cdot)$ takes each of the values l for $1 \leq l \leq N$ infinitely often.
2. The shift function $s(\cdot, \cdot)$ is bounded by some \bar{s} such that $0 \leq s(v, m) \leq \bar{s} \forall v \in \{1, 2, \dots\}, \forall m \in \{1, 2, \dots, N\}$. For the initial step, we additionally require $s(v, m) \leq v$.
3. The shift function $s(\cdot, \cdot)$ is independent of m .

If these conditions are satisfied and $\rho(|M|) < 1$ (i.e., the spectral radius of the iteration matrix, taking the absolute values of its elements, is less than one), the convergence of the asynchronous method is guaranteed [16].

3 Numerical Experiments

3.1 Experimental Setup

The numerical problem we target is the finite difference discretization of the differential equation

$$-\Delta u + \varepsilon u = f,$$

which for $\varepsilon = 0$ becomes the Laplace equation. For Dirichlet boundary condition equal to zero, the 1D discretization for this problem on a grid of size h can be written as a system of linear equations of the form $Ax = b$ where $b = h^2 f$ and

$$A = \begin{pmatrix} 2+h^2\varepsilon & -1 & 0 & \dots & 0 \\ -1 & 2+h^2\varepsilon & \ddots & \ddots & \vdots \\ 0 & \ddots & 2+h^2\varepsilon & \ddots & 0 \\ \vdots & \ddots & \ddots & \ddots & -1 \\ 0 & \dots & 0 & -1 & 2+h^2\varepsilon \end{pmatrix}. \quad (2)$$

Although this may seem to be a very basic problem, it contains many essential aspects necessary to analyze the convergence behavior of the multigrid method. We vary $\varepsilon \in [10^{-6}, 10^{-1}]$ to control the condition number κ of the linear equation system. Using the Gerschgorin circle theorem [19], we can estimate the respective condition numbers by $\kappa \approx 4 \cdot \frac{1}{h^2\varepsilon}$. The first experiment in section 3.3 will analyze the impact of this additive component determining the condition number.

The geometric multigrid method we apply to this system is implemented according to the Algorithm 1, where we use the Conjugate Gradient method for the solution of the coarse grid system. To analyze the performance of a GPU-based block-asynchronous iteration as smoother, we compare it with a CPU implementation of Gauss-Seidel performing smoothing iteration.

For all smoothers, we use a stencil implementation of the corresponding linear equation system, updating the distinct components by using the adjacent components. This reduces the computational cost, since we do not have to perform a sparse matrix vector multiplication, as well as the memory requirements, which are usually daunting when performing GPU-based kernels. Still, we utilize the explicit matrix to compute the error term on each grid level. Utilizing stencils for this may be beneficial for the overall performance as well, but we refrain from doing so since this is not the main target of this paper.

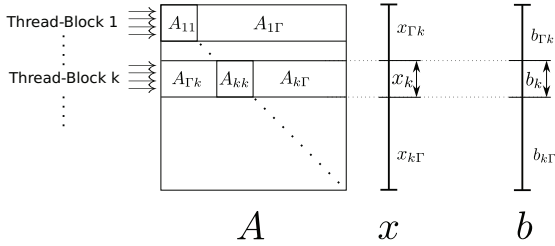
In general, it is very difficult to analyze the performance of a smoother within a multigrid framework. The reason is that a superior smoother will cut the work of the direct solver on the coarsest grid, while using an inferior smoother, the direct solver has a higher workload. This leads to a trade-off between solver workload and smoother workload. Depending on the linear system characteristics, the applied multigrid scheme, and the solver used on the coarsest grid, it may be beneficial to put more or less work-

load on the exact solution process. To enable a comparison, we use a block-asynchronous iteration implementation with smoothing properties similar to the Gauss-Seidel reference method. We furthermore split the numerical tests into two parts, where we first analyze the two-grid iteration and then extend it to a complete multilevel V-cycle.

Enhancing the asynchronous iteration method by local component updates in every thread block leads to the design of a block-asynchronous iteration [2]. While the motivation for this scheme is the design of graphics processing units and the CUDA programming language, it is equivalent to a two-staged asynchronous iteration [5]. We split the linear system into blocks of rows, and then assign the computations for each block to one thread block on the GPU. Between these thread blocks, an asynchronous iteration method is used, while within each thread block, multiple Jacobi-like iterations are performed, instead of a

single iteration. During these local iterations, the x values used from outside the block are kept constant, equal to their values at the beginning of the global iteration. After the local iterations, the updated values are communicated. This approach is inspired by the well know hybrid relaxation schemes [7] [6]. In other words, using domain-decomposition terminology, our blocks correspond to subdomains and thus we iterate locally on every subdomain. We denote this scheme by *async-(i)*, where the index i indicates that we use i Jacobi-like updates on the subdomain. As the subdomains are relatively small and the data needed largely fits into the multiprocessor's cache, these additional iterations on the subdomains come for almost free. The obtained algorithm, visualized in Figure 3, can be written as a component-wise update of the solution approximation:

$$x_k^{(m+1)} = \frac{1}{a_{kk}} \left(b_k - \underbrace{\sum_{j=1}^{T_S} a_{kj} x_j^{(m-v(m+1,j))}}_{\text{global part}} - \underbrace{\sum_{j=T_S}^{T_E} a_{kj} x_j^{(m)}}_{\text{local part}} - \underbrace{\sum_{j=T_E}^n a_{kj} x_j^{(m-v(m+1,j))}}_{\text{global part}} \right), \quad (3)$$



$$x_k = D_{kk}^{-1} (b_k - A_{\Gamma k} x_{\Gamma k} - A_{kk} x_k - A_{k\Gamma} x_{k\Gamma})$$

Figure 3: Visualizing the asynchronous iteration in block description used for the GPU implementation.

where T_S and T_E denote the starting and the ending indexes of the matrix/vector part in the thread block. Furthermore, for the local components, the antecedent values are always used, while for the global part, the values from the beginning of the iteration are used. The shift function $v(m+1, j)$ denotes the iteration shift for the component j — this can be positive or negative, depending on whether the respective other thread block has already conducted more or less iterations. Note that this gives a block Gauss-Seidel flavor to the updates. It should also be mentioned that the shift function may not be the same in different thread blocks.

Using the design of the block-asynchronous iteration, there exist two different parameters that can be used to adjust the smoother: The number of global iterations that

correspond to the number of iterations of a synchronized iterative method like Jacobi or Gauss-Seidel, and the number of local iterations on the subdomains. While the number of global iterations is usually daunting concerning the execution time of a block-asynchronous iteration on GPUs, which is still small compared to synchronized Gauss-Seidel on the CPU, due to the data locality and the used GPU architecture, the latter ones basically come for free [2]. But at the same time, adding local iterations may not trigger the same improvement to the solution approximation. Since the factor between the convergence rate of Gauss-Seidel and Jacobi — which is the fundamental idea of block-asynchronous iteration — equals 2, we always merge 2 global block-asynchronous iterations into one smoothing step. The local iterations may then be used to compensate for the convergence loss due to the chaotic behavior. Without investigating the trade-off between global and local iterations, we set the latter one to the fixed number of 5, and denote the obtained block-asynchronous iteration with *async-(5)*. We also want to neglect the issue of the non-deterministic behavior of *async-(5)*, and refer to all further results as average.

In the second part of the numerical experiment section we then extend the Two-Grid iteration to a full V-cycle. We analyze the impact of adding grid levels, and report the smoother run times for different problem sizes. Finally, we provide a detailed time-to-solution comparison between block-asynchronous iteration and Gauss-Seidel smoothed multigrid for a 10-level implementation using

different numbers of smoothing steps.

3.2 Hardware and Software Issues

The experiments were conducted on a heterogeneous GPU-accelerated multicore system located at the Karlsruhe Institute of Technology, Germany. The system is equipped with two Intel XEON E5540 @ 2.53GHz and 4 Fermi C2070 (14 Multiprocessors x 32 CUDA cores @ 1.15GHz, 6 GB memory). The GPUs are connected to the host through a PCI-e x16.

On the CPU, the synchronous Gauss-Seidel implementation runs on 4 cores. Intel compiler version 11.1.069 [1] is used with optimization flag “-O3”. The GPU implementation is based on CUDA [14], while the respective libraries used are from CUDA 4.0.17 [15]. The component updates were coded in CUDA, using thread blocks of size 512. The kernels are then launched through different streams. The thread block size, the number of streams, along with other parameters, were determined through empirically based tuning.

3.3 Numerical Experiments

In the first experiment, we analyze the impact of the condition number of the linear equation system on the performance of multigrid methods smoothed by block-asynchronous iteration and Gauss-Seidel, respectively. We choose a dimension of $n = 10,000,000$ and compare with respect to the iterations. In Figure 4, we use 1 smoothing step of Gauss-Seidel or block-asynchronous iteration, where we apply 2 smoothing steps in Figure 5. First, we observe that the number of necessary multigrid steps to convergence can be considerably decreased by performing 2 instead of one smoothing iteration. Second, the block-asynchronous smoother has smoothing properties similar to the Gauss-Seidel, so the convergence behavior of the multigrid is not affected by the respective smoother. Only for very small condition numbers, the block-asynchronous iteration performs even better than the Gauss-Seidel smoother. The only difference is the accuracy of the final solution: The Gauss-Seidel method allows a higher approximation quality than the block-asynchronous iteration. But the variations are small and the more crucial factor determining the accuracy of the final solution approximation are the limitations of the used floating point format. Still, if very accurate solution approximations are requested, it may also be reasonable to switch to a Gauss-Seidel method for the last V-cycles.

In the next experiment we investigate the impact of the problem size on the finest grid level. For this purpose we choose problem sizes between 10^4 and 10^8 and analyze the convergence behavior. The results shown in Figure 6 reveal that the problem size has almost no influence on the convergence rate.

While the convergence rate with respect to iteration number is interesting from the theoretical point of view,

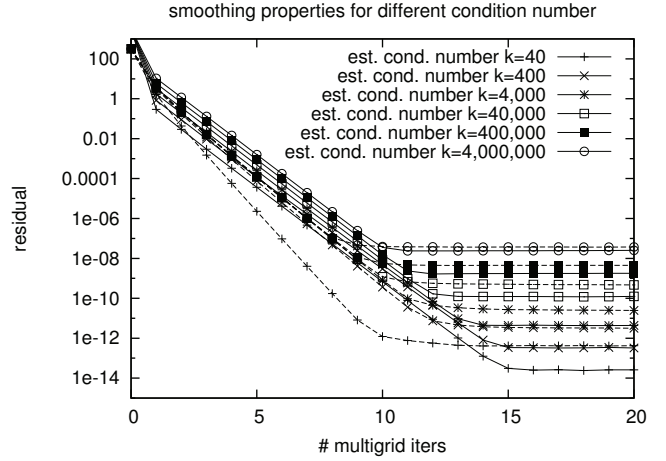


Figure 4: Two-level convergence using one smoothing step for $n = 10,000,000$ and different condition numbers. Dashed lines are block-asynchronous and solid lines are Gauss-Seidel.

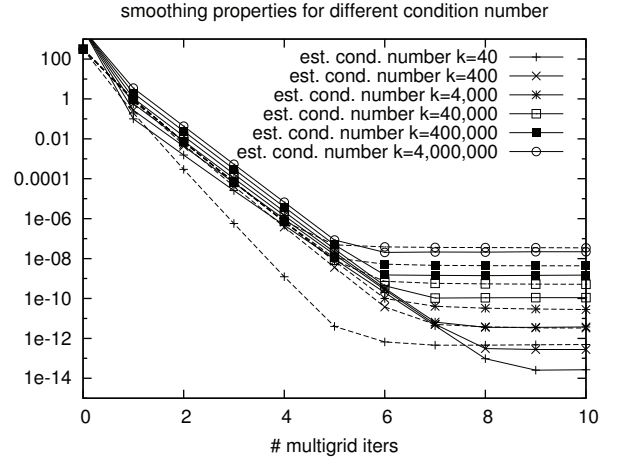


Figure 5: Two-level convergence using two smoothing steps for $n = 10,000,000$ and different condition numbers. Dashed lines are block-asynchronous and solid lines are Gauss-Seidel.

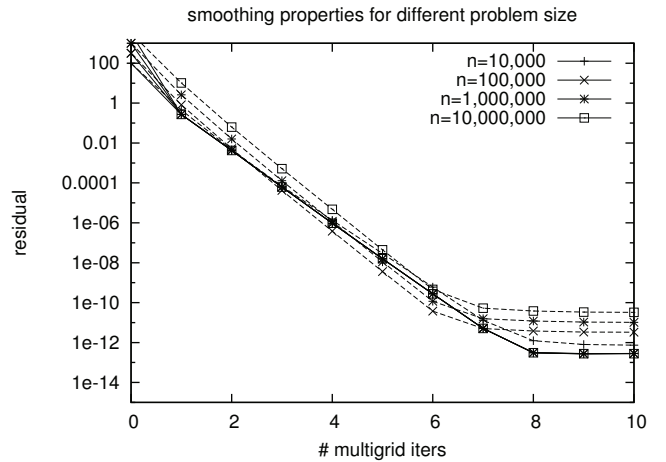


Figure 6: Two-level convergence using two smoothing steps for different problem sizes. Dashed lines are block-asynchronous and solid lines are Gauss-Seidel.

the more relevant factor is the convergence with respect to time. This depends not only on the convergence rate, but also on the efficiency of the respective algorithm on the available hardware resources. While the Gauss-Seidel smoother requires strict update order and therefore allows only sequential implementations, the block-asynchronous iteration is tolerant to varying update orders and synchronization latencies, and therefore enables a highly parallel and synchronization free GPU implementation. In Table 1, we report run times of the respective smoothers for different problem sizes. We also extend the analysis from a two-grid method to multiple levels.

The run times are the aggregated smoother times for the multigrid to converge. Since we usually obtain higher accuracy approximations for the Gauss-Seidel smoother, we choose a stopping criterion for the multigrid iteration that can be achieved for both methods. Additionally we provide the data transfer time for the GPU implementation of the block-asynchronous iteration smoother. This also contains the overhead of the GPU initialization. Note that we report only the run times for the smoother, which increase for multiple levels due to additional smoother calls. The total runtime for the multigrid iteration may still decrease with more levels due to the smaller linear equation system solved on the lowest grid level.

We observe in Table 1 that for Gauss-Seidel on the CPU, linearly increasing the problem size on the finest grid corresponds to a linear run time increase for the smoother. This is also true for the block-asynchronous iteration on the GPU, except for small problem sizes, where calling the GPU kernels triggers some overhead. Also the data transfer is also influenced by the GPU initialization. For all problem sizes and grid sequences, the `async-(5)` smoother outperforms the Gauss-Seidel smoother. While for small problems the improvement is at least a factor of three, it rises to 7 for larger dimension. Since the multigrid framework is often implemented on the host of the system, we should also take the data transfer time into account. Then, for small problem sizes, the `async-(5)` smoother suffers from this overhead due to the GPU initialization and expensive data transfer. For larger problem sizes, the `async-(5)` smoother outperforms the Gauss-Seidel smoother at least by a factor of two, independent of the number of grid levels.

The question is how this corresponds to an acceleration of the multigrid method, since the smoothers usually accounts for a small part in the overall execution time. To investigate this issue, we apply a 10-level multigrid method to a very ill-conditioned problem of size 10,000,000 and provide detailed analysis on the execution time of the smoother, the grid operations like restriction, prolongation and residual computation, and the direct solver on the coarsest grid level.

Analyzing the results, we realize that applying more smoothing steps reduces the number of V-cycles in the multigrid method, which again reduces the number of

solver calls on the coarsest grid level and the number of grid operations. Considering the trade-off between V-cycles and smoothing steps, there is a point for maximal performance. This can usually be determined only heuristically. In our case, it even differs for the block-asynchronous smoother and the Gauss-Seidel smoother. The reason is that the block-asynchronous iteration is not only considerably faster than Gauss-Seidel, but is also dominated by the data transfers between host and CPU. The component updates using a stencil for the block-asynchronous iteration come almost for free. Therefore, increasing the number of iterations does not cause a linear increase of the computation time. Hence, for large numbers of smoothing steps, the speedup factor between Gauss-Seidel and the block-asynchronous iteration smoothed multigrid rises.

Since the smoother accounts for a high percentage of the overall multigrid time, replacing the Gauss-Seidel smoother leads to considerable acceleration. While for more complex problems the ratio between smoother and overall multigrid time is often smaller, the necessity of more smoothing steps rises at the same time, making the block-asynchronous smoother even more attractive. Note that this analysis also takes the data transfer time into account, implementing the multigrid framework on the GPU would trigger even higher speedups.

4 Conclusion

We have shown that asynchronous iteration may be a suitable replacement for Gauss-Seidel or Jacobi smoothers in multigrid methods when targeting highly parallel implementations. Not only for two-grid methods but also for multigrid methods the convergence of the multigrid solver is not affected by the chaotic behavior of the asynchronous smoother. For most test problems we were able to outperform the CPU-based Gauss-Seidel smoother by an asynchronous iteration smoother that utilized the computing power of a graphics processor. Only for small problem sizes, where the overhead of the GPU calls is crucial, was the Gauss-Seidel smoother superior. While for sequential CPU-based smoothers like Gauss-Seidel, a large number of smoothing steps directly corresponds to increased computation time, there is not this linear trade-off for GPU-based block-asynchronous iteration. Hence, choosing a larger number of smoothing steps may be reasonable and improve the overall multigrid performance. Also, adjusting the number of local iterations may have a positive impact. But since the improvement for local and global iterations depends on characteristics of the respective system, more research is necessary at this point. Another research topic may be the field of algebraic multigrid methods. There, depending on the type, asynchronous iteration can even be adapted to the structure of the multigrid method.

dimension	2 levels			3 levels			4 levels			5 levels		
	G.-S.	async-(5)	transfer	G.-S.	async-(5)	transfer	G.-S.	async-(5)	transfer	G.-S.	async-(5)	transfer
10,000	0.00398	0.00085	0.01193	0.00621	0.00181	0.02243	0.00805	0.00242	0.03289	0.00957	0.00338	0.04102
100,000	0.03613	0.00602	0.02571	0.05539	0.00913	0.04267	0.06621	0.01044	0.05673	0.07181	0.01156	0.06866
1,000,000	0.36815	0.05397	0.10779	0.54973	0.08433	0.16522	0.64393	0.09654	0.20734	0.69662	0.10531	0.23634
10,000,000	3.62690	0.53082	0.85057	5.60806	0.80394	1.29404	6.48827	0.92531	1.53345	6.95943	1.00858	1.64843

Table 1: Average smoother runtime [s] for Gauss-Seidel (G.-S.) on CPU, block-asynchronous iteration (async-(5)) on GPU and the data transfer time performing 2 Pre- and 2 Post-smoothing steps on the respective grid levels.

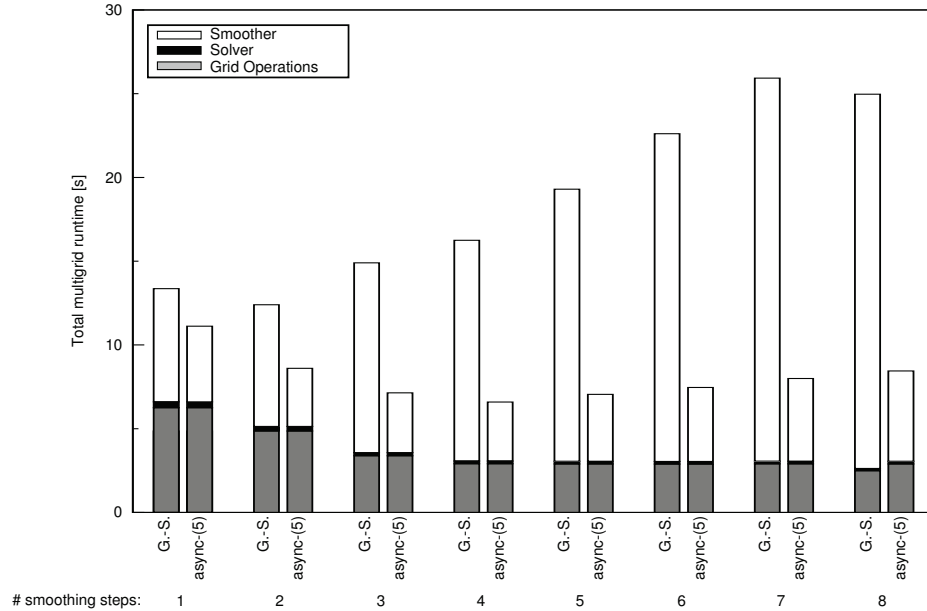


Figure 7: 10-level multigrid v-cycle runtime analysis for different numbers of smoothing steps using Gauss-Seidel and async-(5), respectively. Problem size $n=10,000,000$, condition number estimate 10^6 . The async-(5) smoother includes the data transfer times from and to the GPU.

References

- [1] Intel C++ Compiler Options. Intel Corporation. Document Number: 307776-002US.
- [2] H. Anzt, S. Tomov, J. Dongarra, and V. Heuveline. A block-asynchronous relaxation method for graphics processing units. Technical report, Innovative Computing Laboratory, University of Tennessee, UT-CS-11-687, 2011.
- [3] A. T. Are Magnus Bruaset, editor. *Numerical solution of partial differential equations on parallel computers*. Birkhäuser, 2006.
- [4] U. Aydin and M. Dubois. Sufficient conditions for the convergence of asynchronous iterations. *Parallel Computing*, 10(1):83–92, 1989.
- [5] Z.-Z. Bai, V. Migallón, J. Penadés, and D. B. Szyld. Block and asynchronous two-stage methods for mildly nonlinear systems. *Numerische Mathematik*, 82:1–20, 1999.
- [6] A. H. Baker, R. D. Falgout, T. Gamblin, T. V. Kolev, S. Martin, and U. Meier Yang. Scaling algebraic multigrid solvers: On the road to exascale. *Proceedings of Competence in High Performance Computing CiHPC 2010*.
- [7] A. H. Baker, R. D. Falgout, T. V. Kolev, and U. Meier Yang. Multigrid smoothers for ultra-parallel computing, 2011. LLNL-JRNL-435315.
- [8] D. P. Bertsekas and J. Eckstein. Distributed asynchronous relaxation methods for linear network flow problems. *Proceedings of IFAC '87*, 1986.
- [9] D. Chazan and W. Miranker. Chaotic Relaxation. *Linear Algebra and Its Applications*, 2(7):199–222, 1969.
- [10] A. Frommer and D. B. Szyld. On asynchronous iterations. *Journal of Computational and Applied Mathematics*, 123:201–216, 2000.
- [11] W. Hackbusch. *Multigrid Methods and Applications*. Wolfgang Hackbusch: Multigrid Methods and Applications, Springer, 1985, 1985.
- [12] V. E. Henson and U. M. Yang. Boomeramg: a parallel algebraic multigrid solver and preconditioner. *Applied Numerical Mathematics*, 41:155–177, 2000.
- [13] C. T. Kelley. *Iterative Methods for Linear and Nonlinear Equations*. SIAM, 1995.
- [14] NVIDIA Corporation. *NVIDIA CUDA Compute Unified Device Architecture Programming Guide*, 2.3.1 edition, August 2009.
- [15] NVIDIA Corporation. *CUDA TOOLKIT 4.0 READINESS FOR CUDA APPLICATIONS*, 4.0 edition, March 2011.
- [16] J. C. Strikwerda. A convergence theorem for chaotic asynchronous relaxation. *Linear Algebra and its Applications*, 253(1-3):15–24, Mar. 1997.

- [17] U. Trottenberg. *Multigrid*. Academic Press, 2000.
- [18] A. Üresin and M. Dubois. Generalized asynchronous iterations. In *CONPAR*, pages 272–278, 1986.
- [19] R. S. Varga. *Geršgorin and his circles*. Springer, 2004.
- [20] P. Wesseling. *An introduction to multigrid methods*. John Wiley & Sons Inc, 1992.
- [21] S. F. M. William L. Briggs, Van Emden Henson. *A multigrid tutorial*. SIAM, 2000.
- [22] U. M. Yang. On the use of relaxation parameters in hybrid smoothers. *Numerical Linear Algebra With Applications*, 11:155–172., 2004.

Preprint Series of the Engineering Mathematics and Computing Lab

recent issues

- No. 2011-14 Hartwig Anzt, Jack Dongarra, Vincent Heuveline, Stanimire Tomov: A Block-Asynchronous Relaxation Method for Graphics Processing Units
- No. 2011-13 Vincent Heuveline, Wolfgang Karl, Fabian Nowak, Mareike Schmidtobreck, Florian Wilhelm: Employing a High-Level Language for Porting Numerical Applications to Reconfigurable Hardware
- No. 2011-12 Vincent Heuveline, Gudrun Thäter: Proceedings of the 4th EMCL-Workshop Numerical Simulation, Optimization and High Performance Computing
- No. 2011-11 Thomas Gengenbach, Vincent Heuveline, Mathias J. Krause: Numerical Simulation of the Human Lung: A Two-scale Approach
- No. 2011-10 Vincent Heuveline, Dimitar Lukarski, Fabian Oboril, Mehdi B. Tahoori, Jan-Philipp Weiss: Numerical Defect Correction as an Algorithm-Based Fault Tolerance Technique for Iterative Solvers
- No. 2011-09 Vincent Heuveline, Dimitar Lukarski, Nico Trost, Jan-Philipp Weiss: Parallel Smoothers for Matrix-based Multigrid Methods on Unstructured Meshes Using Multicore CPUs and GPUs
- No. 2011-08 Vincent Heuveline, Dimitar Lukarski, Jan-Philipp Weiss: Enhanced Parallel $ILU(p)$ -based Preconditioners for Multi-core CPUs and GPUs – The Power(q)-pattern Method
- No. 2011-07 Thomas Gengenbach, Vincent Heuveline, Rolf Mayer, Mathias J. Krause, Simon Zimny: A Preprocessing Approach for Innovative Patient-specific Intranasal Flow Simulations
- No. 2011-06 Hartwig Anzt, Maribel Castillo, Juan C. Fernández, Vincent Heuveline, Francisco D. Igual, Rafael Mayo, Enrique S. Quintana-Ortí: Optimization of Power Consumption in the Iterative Solution of Sparse Linear Systems on Graphics Processors
- No. 2011-05 Hartwig Anzt, Maribel Castillo, José I. Aliaga, Juan C. Fernández, Vincent Heuveline, Rafael Mayo, Enrique S. Quintana-Ortí: Analysis and Optimization of Power Consumption in the Iterative Solution of Sparse Linear Systems on Multi-core and Many-core Platforms
- No. 2011-04 Vincent Heuveline, Michael Schick: A local time-dependent Generalized Polynomial Chaos method for Stochastic Dynamical Systems
- No. 2011-03 Vincent Heuveline, Michael Schick: Towards a hybrid numerical method using Generalized Polynomial Chaos for Stochastic Differential Equations
- No. 2011-02 Panagiotis Adamidis, Vincent Heuveline, Florian Wilhelm: A High-Efficient Scalable Solver for the Global Ocean/Sea-Ice Model MPIOM
- No. 2011-01 Hartwig Anzt, Maribel Castillo, Juan C. Fernández, Vincent Heuveline, Rafael Mayo, Enrique S. Quintana-Ortí, Björn Rucker: Power Consumption of Mixed Precision in the Iterative Solution of Sparse Linear Systems
- No. 2010-07 Werner Augustin, Vincent Heuveline, Jan-Philipp Weiss: Convey HC-1 Hybrid Core Computer – The Potential of FPGAs in Numerical Simulation