# Performance Engineering for a Tall & Skinny Matrix Multiplication Kernel on GPUs

Dominik Ernst[1], Georg Hager[1], Jonas Thies[2], and Gerhard Wellein[1]

[1] Erlangen Regional Computing Center (RRZE), 91058 Erlangen, Germany
`dominik.ernst@fau.de`
[2] German Aerospace Center (DLR), Simulation and Software Technology

**Abstract.** General matrix-matrix multiplications (GEMM) in vendor-supplied BLAS libraries are best optimized for square matrices but often show bad performance for tall & skinny matrices, which are much taller than wide. Nvidia's current CUBLAS implementation delivers only a fraction of the potential performance (as given by the roofline model) in this case. We describe the challenges and key properties of an implementation that can achieve perfect performance. We further evaluate different approaches of parallelization and thread distribution, and devise a flexible, configurable mapping scheme. A code generation approach enables a simultaneously flexible and specialized implementation with autotuning. This results in perfect performance for a large range of matrix sizes in the domain of interest, and at least 2/3 of maximum performance for the rest on an Nvidia Volta GPGPU.

## 1 Introduction

### 1.1 Tall & Skinny Matrix Multiplications (TSMM)

The general matrix-matrix multiplication (GEMM) is such an essential linear algebra operation used in many numerical algorithms that hardware vendors usually supply an implementation that is perfectly optimized for their hardware. In case of Nvidia, this is part of CUBLAS ([6]). However, since these implementations are focused on mostly square matrices, they often perform poorly for matrices with unusual shapes.

In this paper, we cover the operation $\mathbf{C} = \mathbf{A}^T \mathbf{B}$, with matrices $\mathbf{A}$ and $\mathbf{B}$ being tall & skinny, i.e., much taller than they are wide. If $\mathbf{A}$ and $\mathbf{B}$ are of size $K \times M$ and $K \times N$, the shared dimension $K$ is long (of the order of $10^6$), whereas the dimensions $M$ and $N$ are short, which we define here as the range $[1, 64]$. $\mathbf{A}$ and $\mathbf{B}$ are both stored in row-major order. We are interested in a highly efficient implementation of this operation on the Nvidia Volta GPGPU.

### 1.2 Application

Row-major tall & skinny matrices are the result of combining several vectors to block vectors. *Block Vector Algorithms* are linear algebra algorithms that compute on multiple vectors simultaneously for improved performance. For instance,
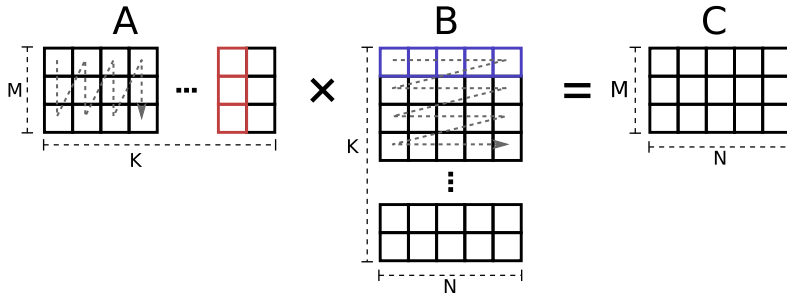
Fig. 1: Illustration of $\mathbf{A}^T\mathbf{B} = \mathbf{C}$ with $\mathbf{A}$ and $\mathbf{B}$ being tall & skinny matrices. Note that $\mathbf{A}$ is transposed in the illustration.

by combining multiple, consecutive *Sparse Matrix Vector* (SpMV) multiplications to a *Sparse Matrix Multiple Vector* (SpMMV) multiplication, the matrix entries are loaded only once and used for the multiple vectors, which reduces the overall memory traffic and consequently increases performance of this memory bound operation. This has first been analytically shown in [3] and is used in many applications such as described in [8].

The simultaneous computation on multiple vectors can also be used to gain numerical advantages. This has been shown for block vector versions of the Lanzcos algorithm [1], of the biconjugate gradient algorithm [7], and of the Jacobi-Davidson Method [8], each of which use block vectors to compute multiple eigenvectors simultaneously. Many such algorithms require multiplications of block vectors. For example, the tall & skinny matrix matrix multiplication $\mathbf{A}^T\mathbf{B}$ occurs in classical Gram-Schmidt orthogonalization of a number of vectors represented by $\mathbf{B}$ against an orthogonal basis $\mathbf{A}$.

### 1.3   Roofline Model

We use the roofline model to obtain an upper limit for the performance of this operation. Under the assumption that the three matrices are transferred just once from memory and that $2MNK$ floating point operations are performed, the arithmetic intensity assuming $K \gg M, N$ and $M = N$ is

$$I_C = \frac{2MNK}{(MK + NK + MN) \times 8} \frac{\text{flop}}{\text{byte}} \stackrel{K \gg M,N}{\approx} \frac{2MN}{(M + N) \times 8} \frac{\text{flop}}{\text{byte}} \stackrel{M = N}{=} \frac{M}{8} \frac{\text{flop}}{\text{byte}} .$$
(1)

In this symmetric case, the arithmetic intensity grows linearly with $M$. This paper will show measurements only for the symmetric case, although the non-symmetric case is not fundamentally different, with the intensity being proportional to the harmonic mean of both dimensions and consequently dominated by the smaller number. With the derived intensity, the model predicts $P = \max\left(M/8 \times B_s, P_{peak}\right)$ as the "perfect" performance yardstick.

Usually the GEMM is considered a classic example for a compute-bound problem with high arithmetic intensity. However, at $M, N = 1$, the arithmetic

intensity is just $1/8$ flop/byte, which is far below the roofline knee of modern compute devices and therefore strongly memory bound. This is not surprising given that a matrix multiplication with $M, N = 1$ is the same as a scalar product. At the other endpoint of the considered spectrum, at $M, N = 64$, the arithmetic intensity is $8$ flop/byte, which is close to the inverse machine balance of a V100 GPU (see below). Therefore the performance character of the operation changes from extremely memory bound at $M, N = 1$ to simultaneously memory and compute bound at $M, N = 64$. An implementation with perfect performance thus needs to fully utilize the memory bandwidth at all sizes and reach peak floating point performance for the large sizes. The very different performance characteristics make it hard to write an implementation that fits well for both ends of the spectrum. Different optimizations are required for both cases.

With the performance as given by the roofline model, it is possible to judge the quality of an implementation's performance as the percentage of the roofline limit. This is plotted in Figure 2 for CUBLAS. The graph shows a potential performance headroom of $5\times$ to $100\times$.

### 1.4   Contribution

This paper shows the necessary implementation techniques to achieve near-perfect performance for double precision tall & skinny matrix-matrix multiplications on an Nvidia V100 GPGPU.

Two different parallel reduction schemes are implemented and analyzed as to their suitability for matrices with lower row counts.

A code generator is implemented that generates code for specific matrix sizes and tunes many configuration options specifically to that size. This allows to exploit regularity, where size parameters allow it, while still generating the least possible overhead where they do not.

### 1.5   Related Work

*CUBLAS* is NVIDIA's BLAS implementation. The GEMM function interface in BLAS only expects column-major matrices. Treating the matrices as transposed column major matrices and executing $\mathbf{AB}^T$ is an equivalent operation. *CUTLASS* is a collection of primitives for multiplications especially of small matrices, which can be composed in different ways to form products of larger matrices. One of these is the `splitK` kernel, which additionally parallelizes the inner summation of the matrix multiplication. Square matrix multiplications usually do not do this, but it is what is required for sufficient parallelism. An adapted version of the "`06_splitK_gemm`" example is used for benchmarking.

### 1.6   Hardware

In this work we use Nvidia's V100-PCIe-16GB GPGPU (Volta architecture) with CUDA 10.0. The hardware data was collected with our own CUDA microbenchmarks, available at [2] together with more detailed data.

| % of occupancy | ILP, Gbyte/s | | |
| --- | --- | --- | --- |
| | 1 | 4 | 16 |
| 1 block, 4 warps | 3.0 | 10.1 | 16.3 |
| 6.25% | 228 | 629 | 815 |
| 12.5% | 419 | 824 | 877 |
| 25% | 681 | 872 | 884 |
| 50% | 834 | 884 | 887 |
| 100% | 879 | 891 | 877 |

Table 1: Measured memory bandwidth on a Tesla V100-PCIe-16GB of a read only kernel with different amount of load parallelism (ILP) and occupancies
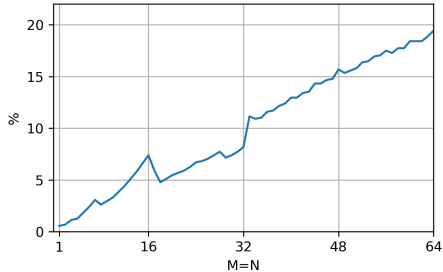


Fig. 2: Percentage of roofline predicted performance achieved by CUBLAS for the range $M = N \in [1, 64]$ on a Tesla V100-PCIe-16GB

*Memory Bandwidth.* The STREAM benchmarks ([5]) all contain a write stream, while the TSMM does not. We thus use a thread-local sum reduction to estimate the achievable memory bandwith (see Table 1). The maximum is above 880 Gbyte/s; but this is only possible with sufficient parallelism, either through high occupancy or instruction level parallelism (ILP) in the form of multiple read streams, achieved here through unrolling.

*Floating Point Throughput.* The V100 can execute one 32-wide double precision (DP) floating point multiply add (FMA) per cycle on each of its 80 streaming multiprocessors (SMs) and runs at a clock speed of 1.38 GHz for a DP peak of $80 \times 32 \times 2 \times 1.38$ Gflop/s = 7066 Gflop/s. One SM quadrant can process a 32 warp lanes-wide instruction every four cycles at a latency of eight cycles. Full throughput can already be achieved with a single warp per quadrant, if instructions are independent.

*L1 Cache.* The L1 cache plays an instrumental role in achieving the theoretically possible arithmetic intensity. The per-SM private L1 cache can transfer a full 128-byte cache line per cycle. Therefore a 32-wide, unit-stride DP load takes two cycles but this number increases by a cycle for each 128-byte cache line that is affected. This amounts to a rate of one load for two FMA instructions.

## 2   Implementation Strategies

The arithmetic intensity (1) was derived assuming perfect caching, i.e. that **A** and **B** are transferred from memory just once. The fastest way to reuse values is to use a register and have the thread, the register belongs to, perform all required operations on this data. Data used by multiple threads can preferably be shared in the L1 cache for threads in the same thread block or in the L2 cache otherwise. This works best only with some spatial and temporal access locality in place.

```
for m = 0...M:
  for n = 0...N:
    for k = 0...K:
      C[m][n] += A[k][m] * B [k][n]
```

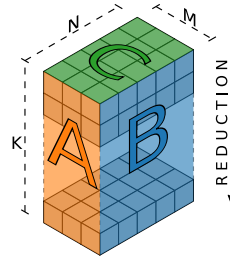Listing 1.1: Naive MMM pseudo code. Note that A is transposed.

Fig. 3: Illustration of the iteration space

## 2.1   Thread Mapping Options

The parallelization scheme, i.e., the way in which work is mapped to GPU threads, plays an important role for data flow in the memory hierarchy. The canonical formulation of an MMM is the three-level loop nest shown in Listing 1.1.

The two outer loops, which are completely independent and therefore well parallelizable, are usually the target of an implementation focused on square matrices. For skinny matrices, these loops are much too short to yield enough parallelism for a GPU. In consequence, the loop over the long K dimension has to be parallelized as well, which also involves parallelizing the sum inside the loop. There are many more terms in the parallel reduction than threads, so that each thread can first serially compute a partial sum, which is afterwards reduced to a total sum.

The iteration space of an MMM can be visualized as the cuboid spanned by the outer product of $\mathbf{A}$ and $\mathbf{B}$ (see Figure 3), where each cell contains a multiplication. A sum reduction over the $K$ dimension yields the result matrix $\mathbf{C}$. Each horizontal slice requires to load a row of $\mathbf{A}$ and $\mathbf{B}$ of length $M$ and $N$, respectively, and computes a rectangle with $M \times N$ FMAs in it.

For data locality, the two small loops have to be moved into the $K$ loop. This creates $MN$ intermediate sums that have to be stored. Depending on whether and how the two small loops are parallelized, each thread computes only some of these intermediates. Figures 4 to 6 visualize this by showing a slice of the multiplication cube and which values a single thread would compute. The number of loads that each thread has to do are the affected values in the row of $\mathbf{A}$ and $\mathbf{B}$, also visible in the illustrations, while the number of FMA operations is the number of highlighted cells in the slice. It is favorable to have a high FMA/loads ratio, as the L1 cache is not as fast as the FP units. This can be achieved by maximizing the area and the squareness of the area that is computed by a single thread. At the same time, more intermediate results per thread increase the register count, which decreases occupancy and eventually leads to spilling.

The easiest approach to achieve this goal would be to only parallelize the $K$ loop, as visualized in Figure 4. While this maximizes the arithmetic intensity
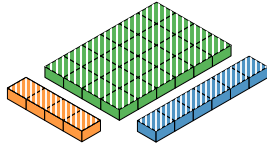
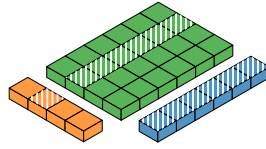Fig. 4:     Parallelization over $K$ only



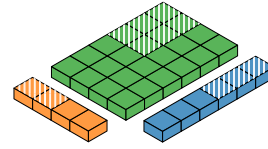Fig. 5:     Parallelization over $K$ loop and an inner loop



Fig. 6:     Parallelization over $K$ and tiling of the two inner loops, here with tile size $2 \times 3$

already in the L1 cache, the $MN$ intermediate results occupy $2MN$ registers, so the maximum of 256 registers per thread is already exceeded at $M, N > 11$, causing spilling and poor performance.

One of the inner loops could be parallelized as well, leading to the pattern in Figure 5. The amount of registers required is only $M$ or $N$, so there is no spilling even at $M, N = 64$. However, the narrow shape results in a `FMA`/loads ratio below 1 and therefore a low arithmetic intensity in the L1 cache.

A better approach that combines managable register requirements with a more square form is to subdivide the two smaller loops into tiles like in Figure 6. This mapping also allows for much more flexibility, as the tile sizes can be chosen small enough to avoid spilling or reach a certain occupancy goal but also large enough to create a high `FMA`/loads ratio.

## 2.2   Global Reduction

After each thread has serially computed its partial, thread-local result, a global reduction is required, which is considered overhead. It depends only on the thread count, though, whereas the time spent in the serial summation grows linearly with the row count and therefore becomes marginal for large enough row counts. However, the authors of [9] argue that the performance at small row counts is still relevant, as the available GPU memory is shared by more data structures than just the two tall & skinny matrices, which limits the data set size.

Since the *Pascal* architecture, atomic add operations are available for global memory, making global reductions more efficient than on older systems. Each thread can just use an `atomicAdd` of its partial value to the final results. The throughput of `atomicAdd` operations is limited by the amount of contention, which grows for smaller matrix sizes. To cut down on the amount of `atomicAdd` operations and therefore contention, it is also possible to first compute thread block-wide partial results, which are then added to the global results. Inside of a thread block, shared memory atomics can be used to compute the block result. Additionally, we opportunistically reduce the amount of launched threads for small row counts.

### 2.3   Leapfrogging and Unrolling

On Nvidia's GPU architectures, load operations can overlap with each other. The execution will only stall at an instruction that requires an operand from an outstanding load. The compiler maximizes this overlap by moving all loads to the beginning of the loop body, followed by the FP instructions that consume the loaded values. At least one or two of the loads come from memory, which take longer to complete than queueing all load operations, so that execution stalls at the first FP instruction. A way to circumvent this stall is to load the inputs one loop iteration ahead into a separate set of *next* registers, while the computations still happen on the *current* values. At the end of the loop, the *next* values become the *current* values of the next loop iteration by assignment. These assignments are the first instructions that depend on the loads and thus the computations can happen while the loads are still in flight. A similar effect can be achieved by unrolling the $K$ loop. The compiler can move the loads of multiple iterations to the front, therefore creating more overlap.

### 2.4   Code Generation

A single implementation cannot be suitable for all matrix sizes. In order to engineer the best code for each size, some form of metaprogramming is required. C++ templates allow some degree of metaprogramming but are limited in their expressivity or require convoluted constructs. Usually the compiler unrolls and eliminates short loops with known iteration count in order to reduce loop overhead, combine address calculations, avoid indexed loads from arrays for the thread-local results, deduplicate and batch loads, and much more. A direct manual generation of the code offers more control, however. For example, in order to enable tile sizes that are not divisors of the matrix dimensions, guarding statements have to be added around computations that could exceed the matrix size. These should be omitted wherever it is safe in order to not compromise performance for dividing tile sizes. We therefore use a code generating script in python, which allows to prototype new techniques much quicker and with more control. Many different parameters can be configured easily, for example whether leapfrogging and unrolling is used, how the reduction is performed, and what tile sizes to set.

## 3   Results

### 3.1   Impact of Reductions

Figure 7 shows the relative performance of our TSMM implementation versus row count with respect to a baseline without any reduction for a selection of inner matrix sizes and tile sizes, choosing either of the two reduction methods described in Section 2.2. As expected, the impact of the reduction generally decreases with increasing row count. The method with only global atomics is especially slow for the narrower matrices ($M, N = 4$). Many threads writing to a
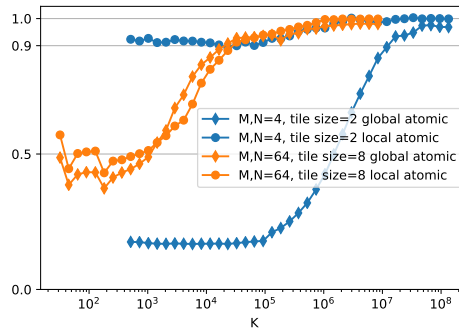
Fig. 7: Ratio of achieved performance of four kernels with two matrix sizes and two different final reduction methods vs. the long dimension $K$ compared to the performance of a kernel that lacks the final reduction

small amount of result values leads to contention and causes a noticeable impact even for a device memory filling matrix ($K = 10^8$). The local atomic variant drastically reduces the number of writing threads, resulting in less than 10% overhead even for the smallest sizes and near perfect performance for $K > 10^6$. For the wider matrices, the difference is smaller. The global atomic version is not as slow because writes spread out over more result values and the local atomic variant is not as fast because the larger tile size requires more work in the local reduction. Both variants incur less than 10% overhead just above $10^4$, a point where only about 0.2% of the GPU memory is used.

## 3.2   Tile Sizes

Figure 8 shows the dependence of performance on tile sizes $T_M$ and $T_N$ for the case $M, N = 32$ with leapfrogging. Performance drops off sharply if the tile sizes become too large and too many registers are used. The number of registers can be approximated by $2 \times (T_M T_N + 2(T_M + T_N))$, which accounts for the thread-local sums ($T_M T_N$) and the loaded values ($T_M + T_N$). Leapfrogging requires two registers for each of the latter and double precision doubles the overall number of needed registers. The graph shows the isolines of 112 and 240 registers, which, accounting for some registers used otherwise, represent the occupancy drop from 25% to 12.5% at 128 registers and the start of spilling at 256 registers.

The best-performing tile sizes generally sit on these lines, maximizing the area of the tile without requiring too many registers. The dimensions are largely symmetric but not perfectly, as threads are mapped to tiles in $M$ direction first. While there are some patterns, where some tile sizes would generally be faster than others, there is no clear trend that either dividers of the matrix width or powers of two are favored in any way.

## 3.3   Best Performance

An exhaustive search was used to find the best tile size for each matrix size. Figure 9 shows the results with and without leapfrogging. The performance meets the roofline prediction (gray dashed line) perfectly until $M, N = 20$. Until
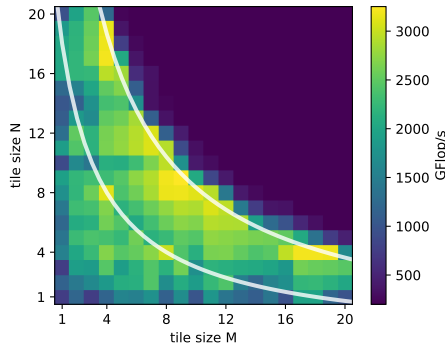
Fig. 8: Performance for different tile sizes, $M, N = 32$, $K = 2^{24}$, leapfrogging enabled. The two white lines are $2 \times (T_M T_N + 2(T_M + T_N)) = R$, with $R = 112,240$ to mark approximate boundaries of register usage.
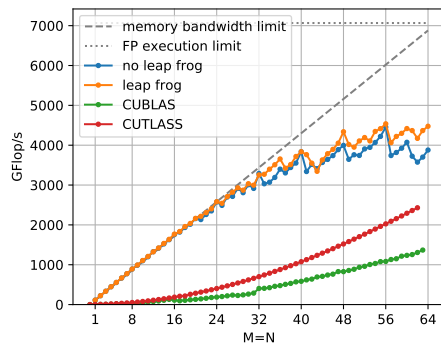
Fig. 9: Best achieved performance for each matrix size with $M = N$ in comparison with the roofline limit, CUBLAS and CUTLASS, with $K = 2^{23}$

$M, N = 36$, the best performance stays within 95% of the limit. Beyond that, the growing arithmetic intensity does not translate into a proportional speedup anymore, although we are still about a factor of two away from peak. The best performance appears to plateau at about 4500 Gflop/s, or ⅔ of peak. Leapfrogging gives about 10–15% advantage for the large sizes. Note that the best tile size changes when leapfrogging is used as it requires more registers.

The cause of the observed performance limit is insufficient memory parallelism. In a large tile, like the $10 \times 8$ which performs best for $M, N = 64$, only 1–2 of the 18 loads go to memory and are not backed by a cache. At 12.5% occupancy and 1–2 read streams, Table 1 shows an achieved memory bandwidth between 400 and 700 Gbyte/s, far below the 880 Gbyte/s needed to meet the roofline expectation.

Both CUBLAS' and CUTLASS' performance is far below the potential performance, especially for the small sizes. Their peformance increases with wider matrices and it can be expected that CUTLASS would overtake the presented implementation at around $M, N = 100$, and CUBLAS somewhat later.

## 4 Conclusion and Outlook

We have shown how to get perfect performance on a V100 GPU for the multiplication of tall & skinny matrices narrower than 32 columns, and at least ⅔ of the potential performance for the rest of the skinny range until width 64. This was achieved using a code generator on top of a suitable thread mapping pattern, which enabled an exhaustive parameter space search. Two different ways to achieve fast, parallel device-wide reductions have been devised in order to ensure a fast ramp up of performance already for shorter matrices.

In future work, in order to push the limits of the current implementation, shared memory could be integrated into the mapping scheme to speed up the many loads, especially scattered ones, that are served by the L1 cache. Another point for improvement would be a transposition of tiles and threads by interleaving threads instead of blocking them, which is a strategy used by CUTLASS. However, our current setup has somewhat lower overhead for nondividing tile sizes and configurations, where a warp would work on multiple slices.

The presented performance figures were obtained by parameter search. An advanced performance model, currently under development, could be fed with code characteristics such as load addresses and instruction counts generated with the actual code and then used to eliminate bad candidates much faster. It will also support a better understanding of performance limiters.

Prior work by us in this area is already part of the sparse matrix toolkit *GHOST* ([4]) and we plan to integrate the presented work as well.

# Bibliography

[1] Cullum, J., Donath, W.E.: A block lanczos algorithm for computing the *q* algebraically largest eigenvalues and a corresponding eigenspace of large, sparse, real symmetric matrices. In: 1974 IEEE Conference on Decision and Control including the 13th Symposium on Adaptive Processes. pp. 505–509 (Nov 1974)

[2] Ernst, D.: CUDA Microbenchmarks. `http://tiny.cc/cudabench` (????)

[3] Gropp, W.D., Kaushik, D.K., Keyes, D.E., Smith, B.F.: Towards realistic performance bounds for implicit cfd codes. In: Proceedings of Parallel CFD'99. pp. 233–240. Elsevier (1999)

[4] Kreutzer, M., Thies, J., Röhrig-Zöllner, M., Pieper, A., Shahzad, F., Galgon, M., Basermann, A., Fehske, H., Hager, G., Wellein, G.: GHOST: Building blocks for high performance sparse linear algebra on heterogeneous systems. International Journal of Parallel Programming pp. 1–27 (2016), `http://dx.doi.org/10.1007/s10766-016-0464-z`

[5] McCalpin, J.D.: Memory bandwidth and machine balance in current high performance computers. IEEE Computer Society Technical Committee on Computer Architecture (TCCA) Newsletter pp. 19–25 (Dec 1995)

[6] NVIDIA: CUBLAS reference. `https://docs.nvidia.com/cuda/cublas` (2019), [Online; accessed 05-May-2019]

[7] O'Leary, D.P.: The block conjugate gradient algorithm and related methods. Linear Algebra and its Applications 29, 293 – 322 (1980), `http://www.sciencedirect.com/science/article/pii/0024379580902475`, special Volume Dedicated to Alson S. Householder

[8] Röhrig-Zöllner, M., Thies, J., Kreutzer, M., Alvermann, A., Pieper, A., Basermann, A., Hager, G., Wellein, G., Fehske, H.: Increasing the performance of the jacobi–davidson method by blocking. SIAM Journal on Scientific Computing 37(6), C697–C722 (2015), `https://doi.org/10.1137/140976017`

[9] Thies, J., Röhrig-Zöllner, M., Overmars, N., Basermann, A., Ernst, D., Wellein, G.: Phist: a pipelined, hybrid-parallel iterative solver toolkit. ACM Transactions on Mathematical Software (01 2019)