Exploiting mesh structure to improve multigrid performance for saddle-point problems

Journal Title XX(X):1-16 ©The Author(s) 0000 Reprints and permission: sagepub.co.uk/journalsPermissions.nav DOI: 10.1177/ToBeAssigned www.sagepub.com/



Lukas Spies¹, Luke Olson¹, and Scott MacLachlan²

Abstract

In recent years, solvers for finite-element discretizations of linear or linearized saddle-point problems, like the Stokes and Oseen equations, have become well established. There are two main classes of preconditioners for such systems: those based on a block-factorization approach and those based on monolithic multigrid. Both classes of preconditioners have several critical choices to be made in their composition, such as the selection of a suitable relaxation scheme for monolithic multigrid. From existing studies, some insight can be gained as to what options are preferable in low-performance computing settings, but there are very few fair comparisons of these approaches in the literature, particularly for modern architectures, such as GPUs. In this paper, we perform a comparison between a Block-Triangular preconditioner and monolithic multigrid methods with the three most common choices of relaxation scheme - Braess-Sarazin, Vanka, and Schur-Uzawa. We develop a performant Vanka relaxation algorithm for structured-grid discretizations, which takes advantage of memory efficiencies in this setting. We detail the behavior of the various CUDA kernels for the multigrid relaxation schemes and evaluate their individual arithmetic intensity, performance, and runtime. Running a preconditioned FGMRES solver for the Stokes equations with these preconditioners allows us to compare their efficiency in a practical setting. We show that monolithic multigrid can outperform Block-Triangular preconditioning, and that using Vanka or Braess-Sarazin relaxation is most efficient. Even though multigrid with Vanka relaxation exhibits reduced performance on the CPU (up to 100% slower than Braess-Sarazin), it is able to outperform Braess-Sarazin by more than 20% on the GPU, making it a competitive algorithm, especially given the high amount of algorithmic tuning needed for effective Braess-Sarazin relaxation.

Keywords

Monolithic multigrid, GPU performance, Relaxation scheme, Vanka, Braess-Sarazin, Schur-Uzawa, Block-Triangular preconditioners

Introduction

Finite-element discretizations are a popular choice for coupled systems such as magnetohydrodynamics (MHD), or the Stokes or Navier-Stokes equations. Even though solvers for finite-element discretizations of such saddlepoint problems are well established, designing efficient and scalable solvers on emerging computing architectures for such systems remains an ongoing challenge (Dou and Liang 2023; Nataf and Tournier 2022; Ershkov et al. 2021).

Here, we focus on preconditioned Krylov methods for the linear or linearized systems of equations that arise in solving such problems. There are two main classes of preconditioners for such systems: preconditioners based on block-factorization approaches (Elman et al. 2005; Benzi et al. 2005; Notay 2019) and those based on monolithic multigrid principles (Brandt and Dinar 1979; Voronin et al. 2022). Within each class there is considerable variability in their building blocks, such as the choice of relaxation scheme in monolithic multigrid, including Braess-Sarazin (Braess and Sarazin 1997; Zulehner 2000), Vanka (Vanka 1986), and Schur-Uzawa (Maitre et al. 1984) relaxation.

From existing studies (Adler et al. 2016; Farrell et al. 2021; Voronin et al. 2022; Adler et al. 2023) some insight can be gained into which algorithms are preferable in serial (low-performance) computing settings, but there are

relatively few fair comparisons of these approaches in literature (John and Tobiska 2000b; Paisley and Bhatti 1998; Larin and Reusken 2008a; Adler et al. 2017), in particular for geometric multigrid on modern architectures, such as GPUs. This is particularly important given changes in prevailing HPC architectures in the past two decades. In (Greif and He 2021; Farrell et al. 2021) it is observed that monolithic multigrid with Vanka relaxation leads to scalable performance mathematically and that additive variants of Vanka are well-suited for implementation on modern GPUs but, to our knowledge, no performance studies support this claim. As we discuss below, one main difficulty in getting good performance out of the Vanka algorithm is the high cost of memory movement for forming the various Vanka patches and updating the global solution, which requires a careful approach to achieve good performance. Similar issues have recently been considered using related additive

Corresponding author:

¹University of Illinois Urbana-Champaign, USA

²Memorial University of Newfoundland, Canada

Lukas Spies, Siebel Center for Computer Science, University of Illinois Urbana-Champaign, 201 N. Goodwin Ave, Urbana, IL 61801, USA Email: Ispies@illinois.edu

Schwarz relaxation schemes within multigrid applied to the Poisson equation (Munch and Kronbichler 2023), where it was found that, with optimization of memory caching and reducing communication between patches, additive Schwarz methods built around cell-centric patches are capable of outperforming optimized point-Jacobi-based relaxation schemes.

More broadly, the parallel scalability of multigrid algorithms on modern architectures faces many challenges related to indirection and increased coarse-grid communication costs (Bienz et al. 2016, 2020). This makes data locality and the cost of data movement a central issue, but one that can be solved by carefully exploiting structure in the problem as is done, e.g., in black box multigrid (BoxMG) algorithms (Dendy 1982; Reisner et al. 2018, 2020). In this work, we use highly structured meshes that allow us to encode various information about the data and how it is accessed in the structure itself, similarly to the BoxMG paradigm. Notably, we work with a structured matrix representation and implement algorithms that take full advantage of this, in order to minimize memory accesses and maximize floating point operations (arithmetic intensity).

In this paper, we first introduce the Stokes equations as our model problem and provide an overview of their structure and resulting discretization. We then introduce two different preconditioners for the FGMRES algorithm used to solve such equations, monolithic multigrid and the upper Block-Triangular preconditioner. For monolithic multigrid, we introduce three different relaxation schemes, Braess-Sarazin, Vanka, and Schur-Uzawa. As Braess-Sarazin and Vanka are two common choices for relaxation schemes, we focus our performance analysis on these, noting that Schur-Uzawa can be implemented with the same kernels as Braess-Sarazin. After illustrating which kernels are the biggest contributors to each algorithm, we then break the performance analysis into two parts: common kernels (matrix-vector and array operations) and Vanka-specific kernels (forming patches, updating the global solution, solving patch systems). For each part, the arithmetic intensity, performance, and runtime are analyzed in order to gain a full understanding of the algorithms and how they compare. This leads to a roofline model to investigate how much better the kernels could be doing, if at all. To finish our performance analysis, we compare a full solve of the Stokes equations using FGMRES preconditioned with both a Block-Triangular preconditioner and a multigrid V-cycle preconditioner with Braess-Sarazin, Vanka, and Schur-Uzawa as relaxation schemes. We show that Vanka is, indeed, a competitive algorithm on modern architectures with careful design. We also show that simply porting a performant CPU implementation of Vanka directly to the GPU is not sufficient for achieving competitive performance.

The Stokes equations and their discretization

Problem setup

Fluid flow where viscous forces are much greater than advective inertia is called Stokes flow. In nature, flow with such properties occurs in many places, e.g., in geodynamics or in the swimming movement of microorganisms. The equations of motions arising from this flow are called the Stokes equations and can be viewed as a simplification of the steady-state Navier-Stokes equations in the limit of small Reynolds number, $Re \ll 1$. They are not only well-suited to be solved by iterative solvers (ur Rehman et al. 2011), but they also serve as a suitable prototype for a wide range of models that lead to saddle-point structure.

Specifically, we consider the incompressible Stokes equation with constant viscosity ν in the unit-square domain $\Omega = [0, 1]^2 \in \mathbb{R}^2$. The equations are given by

$$-\nu\nabla^2 \mathbf{u} + \nabla p = \mathbf{f},\tag{1}$$

$$\nabla \cdot \mathbf{u} = 0. \tag{2}$$

Dirichlet boundary conditions on velocity are enforced on all edges of the domain, but no boundary conditions are imposed on pressure. Here, we consider no-flux boundary conditions,

$$\mathbf{u} \cdot \mathbf{n} = 0 \text{ on } \partial\Omega, \tag{3}$$

where **n** is the outward pointing normal vector. Thus, we define the Hilbert space $\mathbf{H}_{0}^{1}(\Omega)$ as

$$\mathbf{H}_{0}^{1}(\Omega) = \{ \mathbf{v} \in \mathbf{H}^{1}(\Omega) : \mathbf{v} \cdot \mathbf{n} = 0 \text{ on } \partial \Omega \}$$
(4)

and $L^2(\Omega)/\mathbb{R}$ as the quotient space of equivalence classes of functions in $L^2(\Omega)$ differing by a constant. The weak form is then defined as: Find $(\mathbf{u}, p) \in \mathbf{H}_0^1(\Omega) \times L^2(\Omega)/\mathbb{R}$ such that

$$a(\mathbf{u}, \mathbf{v}) + b(\mathbf{v}, p) = (\mathbf{f}, \mathbf{v}) \quad \forall \mathbf{v} \in \mathbf{H}_0^1(\Omega)$$
(5)

$$b(\mathbf{u},q) = 0 \qquad \forall q \in L^2(\Omega)/\mathbb{R}$$
 (6)

where

$$u(\mathbf{u}, \mathbf{v}) = \nu \int_{\Omega} \nabla \mathbf{u} : \nabla \mathbf{v}, \tag{7}$$

$$b(\mathbf{v}, p) = \int_{\Omega} p \nabla \cdot \mathbf{v}.$$
 (8)

We work with a manufactured solution (Ayuso de Dios et al. 2014; Adler et al. 2017) that satisfies the properties above, given by

$$\mathbf{u}(x,y) = \begin{cases} x(1-x)(2x-1)(6y^2 - 6y + 1) \\ y(y-1)(2y-1)(6x^2 - 6x + 1) \end{cases}$$
(9)

$$p(x,y) = x^2 - 3y^2 + \frac{8}{3}xy,$$
(10)

with f computed to satisfy (1). A visualization of this solution is found in Figure 1.

Discretization

We consider the standard Q2–Q1 Taylor-Hood mixed finiteelement discretization on a uniform grid for discretizing the system in Equations (1) and (2). For the velocity, this uses Q2 elements, with biquadratic polynomials for each component on each element as a basis. For the pressure, this uses Q1 elements, with bilinear polynomials on each element as a basis. Both velocities and pressures are required to be continuous across element boundaries. An illustration of the degrees of freedom in these elements is shown in Figure 2.



Figure 1. Visualization of three components of the manufactured solution.



Figure 2. Illustration of the degrees of freedom for a Q2 and Q1 element, with different types of degrees of freedom identified by different shapes.

The resulting Q2 and Q1 elements have a total of nine and four degrees of freedom per element, respectively. Such a discretization of the Stokes equations directly relates back to the weak form as shown in Equations (5) and (6), defining the matrices L and B by

$$L_{i,j} = a(\psi_j, \psi_i) \tag{11}$$

$$B_{k,j} = b(\psi_j, \phi_k). \tag{12}$$

It is important to note that the introduction of a basis gives us two "views" on the finite-element approximations, writing $\mathbf{u} = \sum_i u_i \psi_i$ and $p = \sum_k p_k \phi_k$, so we can consider the functions \mathbf{u} and p directly, or think about their coefficients in the basis expansion, $\{u_i\}$ and $\{p_k\}$. In what follows, we follow the usual convention of overloading the notation \mathbf{u} and p to denote both the functions themselves and the vectors of basis coefficients, with the distinction typically clear from context, in that $L\mathbf{u}$ is the matrix acting on the basis coefficients, while $a(\mathbf{u}, \mathbf{v})$ is the bilinear form evaluated on the function. With this matrix representation, the solution of the weak form can be expressed as that of the linear system

$$\begin{bmatrix} L & B^T \\ B & 0 \end{bmatrix} \begin{bmatrix} \mathbf{u} \\ p \end{bmatrix} = \begin{bmatrix} \mathbf{f} \\ 0 \end{bmatrix}$$
(13)

where L is the discretized Laplacian, B and B^T are the discretized divergence of u and gradient of p, respectively, u and p are the discretized velocity and pressure components of the solution, and f is the velocity component of the right-hand side. In what follows, to save space, we will write the 2×2 block matrix in Equation (13) as A.

Such a system is challenging to solve, as it is symmetric but indefinite due to the zero block in the lower right-hand corner of the system matrix. This causes many common iterative methods (e.g., simple stationary methods like Jacobi and Gauss-Seidel) to not work as they typically involve inverting the diagonal of the system matrix. Study of numerical methods for solution of such saddle-point systems is a well-established discipline (Benzi et al. 2005). One possible solution to these challenges is the use of block preconditioners, based on the block LU factorization of the coupled system. This approach has been well-developed for discretizations of the Stokes equations (Benzi et al. 2005; Elman et al. 2005). However, existing studies (Adler et al. 2017) suggest that monolithic preconditioners can be more efficient. Thus, we also consider monolithic multigrid as preconditioner for FGMRES. We note that monolithic multigrid preconditioners are generally not symmetric and positive definite and, so, they cannot be used directly as preconditioners for MINRES; however, the added computational work for orthogonalization in FGMRES is more than made up for by the quick convergence of the monolithic multigrid approach.

Structured matrix representation

Any iterative solver naturally depends on calculations of matrix-vector products for the block-structured matrix in Equation (13). In general, such calculations require indirect addressing, when arbitrary numbers of elements can be adjacent to each node of the mesh, leading to irregular communication patterns. However, when we restrict the mesh to have logically rectangular structure (meaning that each node is at the intersection of four edges, and is adjacent to four elements), then applying the discretization matrix can be done in a stencil-wise fashion, where each degree of freedom requires information from at most a 2×2 element patch. This allows storing the system matrices in a highly efficient data structure by numbering the degrees of freedom in lexicographic order. For the Q1 discretization, this ordering is natural, since the only degrees of freedom occur at the nodes in the mesh, that can be labeled lexicographically by (x, y)-indices. For the Q2 discretization, we separate the degrees of freedom into four sets, given by those associated with the nodes of the mesh, then those at midpoints of the x and y edges, and, finally, those associated with the cell centers. Figure 3 shows a local numbering of those degrees of freedom on the 2×2 patch around a node.



Figure 3. Local numbering of degrees of freedom around nodal degree of freedom 5.

With this ordering, the system matrix is stored as an array of arrays, where the first (outer) array index corresponds to the row in the matrix associated with that degree of freedom, and the inner arrays stored at each outer index contain a fixed number of entries corresponding to the number of degrees of freedom surrounding the respective degree of freedom, with the respective columns implicitly stored by the positions in the inner array. Thus, the inner arrays have a fixed length of 9 for Q1 elements and a fixed length of 25 for the nodal degree of freedoms for Q2 elements (with the ordering given in Figure 3). Such a matrix format corresponds to a hybrid of the ELLPACK and the SELL-n matrix formats (Rice and Boisvert 1985; Anzt et al. 2014), as we use padding to always store the full stencil and, thus, have a consistent row length, but we also use a (variable) blocksize n with each type of degree of freedom forming its own block. Exploiting the structure of our setup enables us to avoid having to store any indices or pointers, as this information is encoded in the structure itself. Thus we are able to minimize the memory storage needed and maximize the performance, as only a small amount of memory needs to be loaded and read for any operation. Figure 4 illustrates how the matrix entries for the different types of degrees of freedom are stored in memory, with each box corresponding to one inner array.





This approach translates directly to higher dimensions. For a 3D discretization of this type on "brick" elements, we would extend the above to $2 \times 2 \times 2$ element patches with $5 \times 5 \times 5 = 125$ nodal degrees of freedom (5 in each dimension). These can be labeled in an analogous way and,

consequently, stored similarly in a single array of arrays, with outer index corresponding to the rows in the system matrix, and inner index corresponding to the local numbering around each degree of freedom.

Multigrid

Multigrid methods are based on the notion that standard (but slow-to-converge) iterative methods are generally effective at reducing oscillatory errors in a discrete approximation, but that the subspace of smooth and slow-to-converge errors of such an iteration is better treated by a complementary process (Briggs et al. 2000; Trottenberg et al. 2001). A natural approach to reduce the smooth errors is with correction from a coarser-grid realization of the same discretized problem, where those smooth modes can be accurately resolved by recursively applying the same iterations to problems with fewer degrees of freedom, until some suitably coarse version of the problem is found where a sparse direct solver can be effectively applied. For higherorder discretizations and systems of PDEs, such an error classification breaks down (He and MacLachlan 2020), but the multigrid principle remains effective, in that we can define relaxation schemes that effectively damp a large portion of the error in a given approximation, and the remaining error can be effectively corrected from a coarse grid.

The standard multigrid solution algorithm is known as the V-cycle, since it traverses a given hierarchy of meshes from the given finest grid to the coarsest, then back to the finest. In the "downward" sweep of the traversal (from fine-to-coarse), on each level, an initial approximation (generally a zero vector) is improved by a specified relaxation scheme. Then, the residual associated with that improved approximation is calculated and restricted to the next coarsest grid, where the scheme recurses. On the "upward" sweep, the current approximation is improved by interpolating a correction back from the next coarser grid, then running relaxation again, before proceeding to the next finer grid. For transferring residuals and corrections between grids, we define a single interpolation operator that maps from a coarse grid to the next finer grid, and use its transpose as a restriction operator. In this work, we follow the standard geometric multigrid approach of using the finite-element interpolation operators, that naturally map from coarse-grid versions of the Q2 and Q1 spaces to their fine-grid analogues. Algorithm 1 presents an algorithmic overview of the multigrid V-cycle for the Stokes equations, following the convention that level 0 is the coarsest grid in the hierarchy, and we are interested in the solution on some given fine grid, for fixed l > 0.

While traditional relaxation schemes, such as (weighted) Jacobi or Gauss-Seidel are effective for elliptic problems, they generally cannot be applied directly to saddle-point systems, due to the zero block in the matrix. Thus, specialized relaxation schemes are commonly developed and analyzed for the Stokes equations. In this work, we consider four different preconditioning approaches, comparing monolithic multigrid with Braess-Sarazin (Braess and Sarazin 1997; Zulehner 2000), Vanka (Vanka 1986), and Schur-Uzawa (Maitre et al. 1984) relaxation with an upper Block-Triangular preconditioner. We focus in particular on

Algorithm 1 Multigrid V-cycle for Stokes equations

1:	1: function MG(A_l , \mathbf{u}_l , p_l , \mathbf{f}_l , g_l , l)			
2:	Relax on \mathbf{u}_l and p_l			
3:	Compute residual: $\begin{bmatrix} \mathbf{r}_{\mathbf{u},l} \\ r_{p,l} \end{bmatrix} = \begin{bmatrix} \mathbf{f}_l \\ g_l \end{bmatrix} - A_l \begin{bmatrix} \mathbf{u}_l \\ p_l \end{bmatrix}$			
4:	Restriction: $\begin{vmatrix} \mathbf{r}_{\mathbf{u},l-1} \\ r_{p,l-1} \end{vmatrix} = P_{l-1}^T \begin{vmatrix} \mathbf{r}_{\mathbf{u},l} \\ r_{p,l} \end{vmatrix}$			
5:	if l is 1 then			
6:	$\begin{bmatrix} \mathbf{e}_{\mathbf{u},0} \\ e_{p,0} \end{bmatrix} = A_0^{-1} \begin{bmatrix} \mathbf{r}_{\mathbf{u},0} \\ r_{p,0} \end{bmatrix}$			
7:	else			
8:	$\begin{vmatrix} \mathbf{e}_{\mathbf{u},l-1} \\ e_{p,l-1} \end{vmatrix} = \mathrm{MG}(A_{l-1}, 0, 0, \mathbf{r}_{\mathbf{u},l-1}, r_{p,l-1}, l-1)$			
9:	end if			
10:	Correction: $\begin{bmatrix} \mathbf{u}_l \\ p_l \end{bmatrix} = \begin{bmatrix} \mathbf{u}_l \\ p_l \end{bmatrix} + P_{l-1} \begin{bmatrix} \mathbf{e}_{\mathbf{u},l-1} \\ e_{p,l-1} \end{bmatrix}$			
11:	Relax on \mathbf{u}_l and p_l			
12:	end function			

the former two, Braess-Sarazin and Vanka, as these are known to lead to effective monolithic multigrid methods, but also expose key kernels that are reused in the implementation of the latter two. We next provide an overview of all four algorithms, before focusing on aspects of implementation and performance when implementing these approaches on the GPU.

Braess-Sarazin relaxation scheme

The Braess-Sarazin iteration is based on an approximation of the block factorization of the system matrix in (13),

$$\begin{bmatrix} L & B^T \\ B & 0 \end{bmatrix} = \begin{bmatrix} L & 0 \\ B & \hat{S} \end{bmatrix} \begin{bmatrix} I & L^{-1}B^T \\ 0 & I \end{bmatrix}, \quad (14)$$

for $\hat{S} = -BL^{-1}B^T$. The original algorithm (Braess and Sarazin 1997) proposed replacing the matrix, L, in the above by a scaled version of its diagonal, tD, for scalar t, and updating the current approximation by an under-relaxed solve of the saddle-point system with this replacement,

$$\begin{bmatrix} \mathbf{u} \\ p \end{bmatrix}^{new} = \begin{bmatrix} \mathbf{u} \\ p \end{bmatrix}^{old} + \omega_{BS} \begin{bmatrix} tD & B^T \\ B & 0 \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{r}_{\mathbf{u}} \\ r_p \end{bmatrix}^{old}$$
(15)

where $\mathbf{r}_{\mathbf{u}}$ and r_p are the respective residuals. For simplicity, we fix $\omega_{BS} = 1$ in what follows. This is equivalent to computing the (unweighted) updates, $\delta \mathbf{u}$ and δp , as approximate solutions of the block-factorized approximation to the system matrix,

$$\begin{bmatrix} tD & 0\\ B & S \end{bmatrix} \begin{bmatrix} I & \frac{1}{t}D^{-1}B^T\\ 0 & I \end{bmatrix} \begin{bmatrix} \delta \mathbf{u}\\ \delta p \end{bmatrix} = \begin{bmatrix} \mathbf{r}_{\mathbf{u}}\\ r_p \end{bmatrix}$$
(16)

where $S = -\frac{1}{t}BD^{-1}B^{T}$ is the Schur complement of the approximated system. Equation (16) can be rewritten as two equations

$$S\delta p = r_p - \frac{1}{t}BD^{-1}\mathbf{r_u},\tag{17}$$

$$\delta \mathbf{u} = \frac{1}{t} D^{-1} (\mathbf{r}_{\mathbf{u}} - B^T \delta p).$$
(18)

The *inexact* variant of Braess-Sarazin (Zulehner 2000) only approximately computes the solution to Equation (17),

using standard weighted Jacobi (or other algorithms) to approximate the inverse of S (see also (He and MacLachlan 2019). The full algorithm is given in Algorithm 2.

Algorithm 2 Braess-Sarazin relaxation

- 1: Approximately solve $S\delta p = r_p \frac{1}{t}BD^{-1}\mathbf{r_u}$ for δp by relaxation.
- 2: Compute $\delta \mathbf{u} = \frac{1}{t} D^{-1} (\mathbf{r}_{\mathbf{u}} B^T \delta p).$
- 3: Update $p^{new} = p^{old} + \omega_{BS}\delta p$.
- 4: Update $\mathbf{u}^{new} = \mathbf{u}^{old} + \omega_{BS} \delta \mathbf{u}$.

Vanka relaxation scheme

Vanka relaxation, in contrast, applies a block overlapping Schwarz iteration to the global saddle-point system. In this approach, we define sets of "patches" (or "subdomains" in the usual Schwarz notation) corresponding to 2×2 blocks of elements, where we take a single pressure degree of freedom at the central vertex and all velocity degrees of freedom on the associated (neighboring) elements, see Figure 5. For each patch, we define a restriction operator, V_i , that extracts degrees of freedom from the global matrix to those present on local patch *i*, and use this to restrict the system matrix to the *i*th patch, as

$$A_i = V_i A V_i^T \tag{19}$$

The Vanka algorithm is defined by looping over the patches and solving

$$A_{i} \begin{bmatrix} \delta \mathbf{u}_{i} \\ \delta p_{i} \end{bmatrix} = V_{i} \begin{bmatrix} \mathbf{r}_{\mathbf{u}} \\ r_{p} \end{bmatrix}$$
(20)

exactly for $\delta \mathbf{u}_i$ and δp_i , which are then used to update the global approximate solution in a weighted additive manner. The Vanka algorithm is given in Algorithm 3. We note



Figure 5. Illustration of overlapping 2×2 Vanka patches

that we solve the patch systems exactly by inverting each patch matrix ahead of time, as they do not change between iterations.

Algorithm 3 Vanka relaxation (additive)			
1:	for $i = 1, \ldots, N$ do		
2:	Solve $A_i \begin{bmatrix} \delta \mathbf{u}_i \\ \delta p_i \end{bmatrix} = V_i \begin{bmatrix} \mathbf{r}_\mathbf{u} \\ r_p \end{bmatrix}$ for $\begin{bmatrix} \delta \mathbf{u}_i \\ \delta p_i \end{bmatrix}$.		
3:	end for		
4:	Update $\begin{bmatrix} \mathbf{u}^{new} \\ p^{new} \end{bmatrix} = \begin{bmatrix} \mathbf{u}^{old} \\ p^{old} \end{bmatrix} + \sum_{i=1}^{N} V_i^T W_i \begin{bmatrix} \delta \mathbf{u}_i \\ \delta p_i \end{bmatrix}$. where W_i is the matrix with the weights.		

Schur-Uzawa relaxation scheme

The Schur-Uzawa iteration is derived in a similar way to the Braess-Sarazin iteration. It is again based on an approximation of the factorization of the system matrix in (13),

$$\begin{bmatrix} L & B^T \\ B & 0 \end{bmatrix} = \begin{bmatrix} L & 0 \\ B & \hat{S} \end{bmatrix} \begin{bmatrix} I & L^{-1}B^T \\ 0 & I \end{bmatrix}, \quad (21)$$

for $\hat{S} = -BL^{-1}B^T$. As in Braess-Sarazin, we again replace the matrix, L, by some approximation that is easy to invert, again denoted tD, but also drop the upper-triangular term from this factorization, resulting in an inexact system computing updates, $\delta \mathbf{u}$ and δp , as approximate solutions of the block system,

$$\begin{bmatrix} tD & 0\\ B & \hat{S} \end{bmatrix} \begin{bmatrix} \delta \mathbf{u}\\ \delta p \end{bmatrix} = \begin{bmatrix} \mathbf{r}_{\mathbf{u}}\\ r_{p} \end{bmatrix}.$$
 (22)

The system in (22) can be rewritten as two equations

$$tD\delta \mathbf{u} = \mathbf{r}_{\mathbf{u}} \tag{23}$$

$$S\delta p = B\delta \mathbf{u} - r_p \tag{24}$$

that are solved for $\delta \mathbf{u}$ by directly inverting tD and for δp by standard weighted Jacobi to approximate the inverse of $S = -\frac{1}{t}BD^{-1}B^T$. The full algorithm is given in Algorithm 4.

Algorithm 4 Schur-Uzawa relaxation

1: Compute $\delta \mathbf{u} = \frac{1}{4} D^{-1} \mathbf{r}_{\mathbf{u}}$.

- 2: Approximately solve $S\delta p = B\delta \mathbf{u} r_p$ for δp by relaxation.
- 3: Update $p^{new} = p^{old} + \delta p$.
- 4: Update $\mathbf{u}^{new} = \mathbf{u}^{old} + \delta \mathbf{u}$.

Block-Triangular preconditioner

The Block-Triangular preconditioner is also based on an approximation of the system matrix in (13), but we now consider an alternate form with unit block diagonal for the lower-triangular factor,

$$\begin{bmatrix} L & B^T \\ B & 0 \end{bmatrix} = \begin{bmatrix} I & 0 \\ BL^{-1} & I \end{bmatrix} \begin{bmatrix} L & B^T \\ 0 & \hat{S} \end{bmatrix}, \quad (25)$$

with $\hat{S} = -BL^{-1}B^T$. While Braess-Sarazin and Schur-Uzawa relaxation use simple approximations to L and \hat{S} to approximate the inverse for relaxation within a multigrid cycle, it is more common to use multigrid on the blocks when using the block preconditioner directly. Here, as is common (Elman et al. 2005), we first approximate \hat{S} by a mass matrix, -M, on the pressure space, then apply multigrid to this approximation. With this, we compute updates, δu and δp , as approximate solutions of the uppertriangular approximation to the system matrix,

$$\begin{bmatrix} L & B^T \\ 0 & -M \end{bmatrix} \begin{bmatrix} \delta \mathbf{u} \\ \delta p \end{bmatrix} = \begin{bmatrix} \mathbf{r}_{\mathbf{u}} \\ r_p \end{bmatrix}, \qquad (26)$$

using multigrid V-cycles to approximately invert L and M. The system in (26) can be rewritten as two equations

$$-M\delta p = r_p \tag{27}$$

$$L\delta \mathbf{u} = \mathbf{r}_{\mathbf{u}} - B^T \delta p, \qquad (28)$$

leading to the full algorithm given in Algorithm 5.

Algorithm 5 Block-Triangular preconditioner

- 1: Approximately solve $M\delta p = -r_p$ for δp using multigrid on M.
- 2: Approximately solve $L\delta \mathbf{u} = \mathbf{r}_{\mathbf{u}} B^T \delta p$ for $\delta \mathbf{u}$ using multigrid on L.
- 3: Update $p^{new} = p^{old} + \delta p$.
- 4: Update $\mathbf{u}^{new} = \mathbf{u}^{old} + \delta \mathbf{u}$.

Our implementation

We have implemented the outer FGMRES iteration and a multigrid V-cycle with the three relaxation schemes, Vanka, Schur-Uzawa, and Braess-Sarazin, as well as the Block-Triangular preconditioner, in C++, with custom data structures that provide structured matrix implementations of the required matrix and vector operations, as discussed above. This is achieved by using operator overloading to allow the optimization of the code for different architectures while preserving a clean implementation of the highlevel algorithms. Underlying the custom data structures are standard STL vectors of double data type. Additionally, support for CUDA and OpenCL requires only a switch of the backend implementation while the high-level algorithm implementation remains largely untouched.

The resulting implementation yields an efficient solution algorithm for the incompressible Stokes equations in two dimensions on both the CPU and the GPU. We limit our attention to optimizing implementations for a single CPU node or single GPU and focus on comparing performance using the different algorithms, taking advantage of the underlying structure. In principle, similar performance is expected for other discretizations of Stokes and other saddlepoint problems, in two and three dimensions, but studying performance in these contexts is left for future work. We also do not consider extending this work to MPI-based parallelism or multi-GPU systems.

Existing work

John and Tobiska (2000a) investigate the performance of multigrid paired with Braess-Sarazin for solving the Stokes equations using P1-P0 finite elements. In the case of a Wcycle, the improvement in error reduction is approximately linear with the number of smoothing steps. For a V-cycle, convergence increases in general with increasing level of refinement. The work shows that multigrid paired with Braess-Sarazin is indeed a robust and reliable preconditioner.

Larin and Reusken (2008b) compares use of a coupled multigrid method with Vanka and Braess-Sarazin type relaxation schemes, along with preconditioned MINRES and an inexact Uzawa method. The focus is on solving the Stokes equations on the then-current hardware and architectures. The results show that all four methods are robust with respect to variations in parameters. The conclusion is that a multigrid W-cycle paired with diagonal Vanka results in the most efficient solver in terms of CPU time.

More recently, Adler et al. (2017) compare of a fullycoupled monolithic multigrid paired with Braess-Sarazin or Vanka as relaxation scheme, and a block-factorization preconditioner similar to the one we presented here. On CPU-only systems, multigrid paired with Vanka results in the best scaling and lowest iteration count. Yet, these solvers require significantly more work per iteration than the other preconditioners. As a result, multigrid paired with Braess-Sarazin yields the best time-to-solution on CPUs for the problems studied in that work.

Performance Analysis

By far the most costly component of the monolithic multigrid-preconditioned FGMRES solver is the relaxation scheme within the multigrid V-cycle. Thus, we focus our performance analysis on the implementations of two of the relaxation schemes, Vanka and Braess-Sarazin, alone, noting that multigrid with Schur-Uzawa relaxation reuses only components from that using Braess-Sarazin, while the Block-Triangular preconditioner also reuses primarily kernels from Braess-Sarazin relaxation as well.

Studying the performance of these methods requires careful analysis of memory movement and access. The standard metric for this is *arithmetic intensity*, which quantifies the relationship between floating-point operations and memory reads and writes. Another important quantity is the *FLOP rate*, which describes how many floating-point operations are performed in a given time frame. The runtime of the kernels (and how they relate to one-another) indicates the importance of each kernel when it comes to studying the performance. We first study and optimize the component kernels individually, before comparing performance of our four preconditioners for the FGMRES solver for the Stokes equations. For measuring the various metrics to evaluate and compare implementations on the GPU, we use NVIDIA's Nsight Compute* and Nsight Systems[†] software.

Test System

The system we use for each test is the Delta supercomputer[‡] located at the National Center for Supercomputing Applications (NCSA). It is equipped with NVIDIA A100 GPUs that have a measured peak double-precision floating-point performance of 9472.34 GFLOP/s, 80GB on-chip memory, and a measured GPU memory bandwidth of 1264.42 GB/s, measured using the CS roofline toolkit[§]. For our final comparison of algorithms, we also run on the CPU nodes of the Delta supercomputer, that carry dual 64-core AMD 7763 processors with a base frequency of 2.45 GHz (max boost frequency of 3.5 GHz) and a per socket memory bandwidth of 204.8 GB/s, although we only consider serial runs on a single core here.

Kernels

To start the analysis of the different kernels, we first present an overview of each kernel and get a sense of how much they contribute to the overall runtime. Although the problem itself and the final parameter choices can have an important impact on the performance of a kernel, we compare kernels for a generic case here to provide a baseline.

Algorithm 6 shows the Braess-Sarazin algorithm in slightly different form than Algorithm 2, to focus on the kernels involved for the various steps. These kernels are color-coded for ease of comparing with the cost breakdown for a single iteration of Braess-Sarazin shown in Figure 6,

7

indicating how much each kernel contributes to the overall runtime. (Noting that the percentages may not sum to 100%, due to rounding.) As is seen (and expected), most of the runtime is consumed by the matrix-vector operations. Many of these kernels are reused in the other solvers. We note in particular that the weighted Jacobi kernel for the pressure solve, used both for Braess-Sarazin and Schur-Uzawa relaxation (and in the Block-Triangular preconditioner) contributes very little to the overall runtime (less than 2%).

Algorithm 6 Braess-Sarazin with kernel breakdowns		
1:	Compute current residuals, $\mathbf{r}_{\mathbf{u}}$ and r_p .	
	Q2 matrix * Q2 vector	
	Q2Q1 matrix * Q2 vector	
	Q2Q1 matrix * Q1 vector	
	array operations	
2:	Form right hand side of Equation (17).	
	Q2 matrix * Q2 vector	
	Q2Q1 matrix * Q2 vector	
	array operations	
3:	Use Jacobi to compute approximation of δp in Equa-	
	tion (17).	
	weighted Jacobi	
4:	Use δp to compute δu in Equation (18).	
	Q2 matrix * Q2 vector	
	Q2Q1 matrix * Q1 vector	

array operations

5: Update global solution with $\delta \mathbf{u}$ and δp .



Figure 6. Braess-Sarazin: Kernels and their proportion of runtime

A similar kernel-focused restatement of Algorithm 3 is given in Algorithm 7, with the various kernels color-coded to correspond to timing breakdown for a single iteration shown in Figure 7. The left figure in Figure 7 shows the runtime for what we call "simple Vanka", where we do not take advantage of the fact that, for many settings, the Vanka submatrices are identical for most patches and can, thus, be stored once and used many times. This approach results in more than 75% of the runtime being spent applying the patch matrix inverses as each patch needs to load its own Vanka submatrix from global memory. The right figure in Figure 7

^{*}NVIDIA Nsight Compute: https://developer.nvidia.com/ nsight-compute

[†]NVIDIA Nsight Systems: https://developer.nvidia.com/ nsight-systems

[‡]Delta supercomputer: https://delta.ncsa.illinois.edu/

[§]CS Roofline Toolkit: https://bitbucket.org/berkeleylab/ cs-roofline-toolkit



Figure 7. Vanka: Kernels and their proportion of runtime, both (a) simple and (b) tuned Vanka.

_
id Fig
I IU

2415 $\overline{23}$ Q 13 19 2521 22 20 3 4 18 $\overline{7}$ $\mathbf{2}$ 175 6

Figure 8. Vanka: groups of shared patch matrices

shows the results using a "tuned Vanka" implementation, where patches that have identical submatrices take advantage of fast shared memory to optimize memory accesses and, in turn, improve performance. For a uniform grid as considered here, Figure 8 sketches the grouping of patches into those that have a submatrix in common. Here, there are special cases for patches adjacent to the edges or corners of the mesh, including those associated with nodes on the boundary and those distance one from the boundary (where some degrees of freedom in the patch have Dirichlet boundary conditions applied), and a general case for all patches at nodes at least distance two from the boundary. This approach results in only about 40% of the overall runtime being taken up by applying patch matrix inverses. In total, just over three quarters of the runtime in the tuned approach is used for the four unique-to-Vanka operations. The other portion is contributed by the same simple matrix-vector operations as in the analysis of Braess-Sarazin. In what follows, we focus on tuned Vanka in our performance analysis and show a comparison of tuned and simple Vanka as part of our final comparison of relaxation schemes within multigridpreconditioned FGMRES.

Vanka patch matrices

update global solution

GPUs are built around multithreaded streaming multiprocessors. Whenever a kernel is launched from the host, all threads are grouped together into smaller thread blocks, which are then enumerated and distributed to available multiprocessors. All threads within a thread block are executed concurrently, and all blocks can be executed concurrently. Threads within a block are able to access local shared memory and can be synchronized. Additionally, all threads are able to access global memory. Choosing the right size of block is essential for good performance, as a too small block size leads to streaming multiprocessors that remain partially idle, whereas a too large block size leads to an imbalanced load over all of the streaming multiprocessors. CUDA is designed with a maximum of 1024 possible threads per thread block.

Making use of shared and local memory as much as possible within a CUDA thread block allows us to optimize memory accesses further, as data that is used repeatedly can be cached in memory that is faster than global memory. This is of particular importance for the Vanka algorithm, as many of the Vanka patches share the same patch matrices, as discussed above. In total, there are 25 different patch matrices, as depicted in Figure 8, for a constant-coefficient Stokes problem on a uniform mesh, independent of the number of elements. The shaded orange areas in Figure 8 denote single patches that have their own unique patch matrix. The normal orange areas are one-dimensional areas along element edges that share the same patch matrix, and the blue area is the two-dimensional interior region of the domain, where all patches share the same submatrix. Within any one of these regions, we can load the patch matrix into shared memory once, to be used by all threads in the block.

kernel	reads	writes	flops
Q2 array plus/minus array	2n	n	n
Q2 matrix * Q2 vector	n+1 $50+160\ell+128\ell^2$	$\frac{n}{n}$	n^{n} 25 + 80 ℓ + 64 ℓ^{2}
Q2Q1 matrix * Q2 vector Q2Q1 matrix * Q1 vector	$50 + 100\ell + 50\ell^2 18 + 60\ell + 50\ell^2$	m	$25 + 50\ell + 25\ell^2$ $9 + 30\ell + 25\ell^2$
Braess-Sarazin: weighted Jacobi	2m + 1	m	2m
Vanka: Form Patch RHS	$51 + 102\ell + 51\ell^2$	$51 + 102\ell + 51\ell^2$	0
Vanka: Apply matrix inverse Vanka: Update global solution	$52 + 5304\ell + 52\ell^2$	$51 + 102\ell + 51\ell^{-1}$ $52 + 104\ell + 52\ell^{2}$	$\frac{2601 + 5202\ell + 2601\ell^2}{52 + 104\ell + 52\ell^2}$

Table 1. Theoretical reads [double], writes [double] and flops of the various operations with $\ell + 1$ as the number of nodal degrees of freedom in one dimension, n as the total number of velocity degrees of freedom, and m as the total number of pressure degrees of freedom.

Arithmetic Intensity

The *arithmetic intensity* of a kernel is defined as the ratio of how many floating point operations (flops) are performed per byte read/written. The algorithms for both Vanka and Braess-Sarazin involve various general vector and matrix-vector operations. In addition, Braess-Sarazin requires a weighted Jacobi application, and Vanka requires operations to extract the current residuals, to apply the patch matrix inverses, and to update the global solution. Table 1 denotes the counts of all reads, writes, and floating points operations (flops) for the various kernels, obtained by counting the operations in the algorithms. Note that the Q1 array operations are similar to the Q2 array operations leading to equivalent arithmetic intensity and performance values.

Based on the values in Table 1, we compute the arithmetic intensity of the various kernels. On the GPU (using CUDA), we use the following formula,

$$AI = \frac{flops}{32(sectors read + sectors written)}.$$
 (29)

where the reads and writes are counted per sector. One sector consists of a total of 32 bytes and, thus, we multiply by that constant in order to recover the byte count. The performance of each operation is then computed by

$$perf = \frac{flops}{\max\left(\frac{32(\text{sectors read + written})}{\text{bandwidth}}, \frac{flops}{\text{peak perf}}\right)}$$
(30)

The theoretical arithmetic intensity and performance computed this way is shown in Table 2. This analysis, however, has its limitations. In practice, we expect the actual arithmetic intensity and performance to be better, as values are typically not read from or written to memory one-byone. Instead, a memory range is typically loaded all at once, allowing us to reuse values. Additional strategies, like using shared memory for Vanka patches with the same patch matrix, further optimize the memory accesses, increasing both the arithmetic intensity and performance. Similarly, varying the size of CUDA blocks also has an effect on these quantities.

Common Kernels

For simplicity, we group the kernels into two classes. First, we examine those kernels that are common to both Vanka and Braess-Sarazin relaxation, involving matrix-vector products

kernel	AI p	erformance
array plus/minus array	0.0417	52.684
	0.0025	79.020
Q2 matrix $*$ Q2 vector	0.0606	76.633
Q2Q1 matrix * $Q2$ vector Q2Q1 matrix * $Q1$ vector	0.0613 0.0578	(1.4)(1.47)
	0.0018	13.113
Braess-Sarazin: Jacobi	0.0833	105.368
Vanka: Form Patches	0.0	0.0
Vanka: Apply matrix Vanka: Undate solution	0.120 0.125	152.088 157 744
Vanka: Apply matrix Vanka: Update solution	$0.120 \\ 0.125$	152.088 157.744

Table 2. Theoretical AI [flops/byte] and performance [GFLOP/s], calculated for a 512×512 element patch.

and array operations. Following this, we analyze the Vankaspecific kernels.

The common kernels are listed at the top of Figure 9, where we break down the matrix-vector products producing Q2 vectors into those that compute values at the nodes, denoted by n, the x- and y-edge midpoints, denoted by xand y, respectively, and the cell centers, denoted by c. The color-coding of these kernels matches that in Figure 6. For these kernels, we can choose the CUDA block size in an attempt to improve performance. Figure 9 shows how the arithmetic intensity, performance, and runtime vary for the various common kernels with varying CUDA block size.

We first note that the measured arithmetic intensities are indeed better than the theoretical values described in Table 2, albeit generally not far off. Additionally, the measured arithmetic intensity does not vary much (or at all) with varying CUDA block size. This is due to the nature of the underlying memory operations, as the structured matrix data structures used here already optimize the loading and writing of memory. Due to the global nature of the kernels, they are not able to take advantage of shared memory on the GPU. The performance and runtime, however, are impacted by the CUDA block size, with increases in the performance leading to decreases in runtime. Over all results, we see differences in performance of up to a factor of 5 as we vary the block size. Choosing the best overall parameter comes down to selecting the best block size for the kernels that contribute the most to each algorithm.

For Braess-Sarazin, the Q2 matrix by Q2 vector multiplication makes up more than 50% of the overall runtime and, thus, choosing the best parameter for the 4



Figure 9. Common kernels: CUDA block size vs. AI, performance, and runtime. 256^2 elements in left column, 1024^2 elements in right column.

kernels within this operation has the largest impact on the overall runtime of the algorithm. For both problem sizes, the best (or near-best, within 2% of the best) runtime for these 4 kernels is achieved for a CUDA block size of 12×12 . Analyzing the other common kernels yields a very similar picture. Thus, all of the common kernels achieve peak (or near-peak) performance for a CUDA block size of 12×12 , which we choose for the Braess-Sarazin algorithm for which these kernels dominate the cost. We confirmed that these are the best choices by timing a full iteration of the algorithm for both problem sizes. Table 3 provides a concise overview of the best parameters for both algorithms.

Vanka-Specific Kernels

For the Vanka-specific kernels, we perform a similar analysis as for the common kernels. For all four kernels, we vary the thread block size from 4×4 to 16×16 . Figure 10 shows how the arithmetic intensity, performance, and runtime varies with this parameter, again matching the color-coding used in Figure 7. Here, we notice that the measured arithmetic intensity is higher than the theoretical analysis in Table 2, in particular for the kernels applying the patch inverses. This is expected, as we take advantage of fast shared memory for storing the shared patch matrices, which is not accounted for in that analysis. We note, however, that the arithmetic intensity does not vary much with block size, remaining largely constant. The kernel updating the global solution has a comparatively low arithmetic intensity, as it



Figure 10. Vanka-specific kernels: thread-block size vs. Al, performance, and runtime. 256^2 elements in left column, 1024^2 elements in right column.

consists largely of memory movements and only very few floating-point operations. Similarly, the kernel forming the various Vanka patches does not contain any floating-point operations, resulting in zero arithmetic.

Analyzing the performance of the four kernels shows a rather similar picture, with the thread block size causing only small variations in the performance. Even though the kernel for the exterior patches and the kernel for the interior patches have a very similar arithmetic intensity, they differ widely in terms of performance, by up to 2 orders of magnitude. This is due to the comparatively high amount of work to be done for the interior patches. Once again, the kernel for forming the Vanka patches has a performance of 0 GFLOP/s, as it does not contain any floating point operations.

Both of these metrics, the arithmetic intensity and performance, are useful for evaluating the different kernels, but the effective runtime is the defining criteria for which any set of parameters is, ultimately, the best choice. Even though the kernels applying the patch matrices to the exterior patches (A) has a much lower performance than the kernel applying the patch matrices to the interior patches (B), the runtime of (A) for the smaller problem size is only about a factor of 3 larger than that for (B). For the larger problem size, the runtime of (A) is much lower than that for (B), by a factor of about 8. This is due to the overall relatively small amount of computations required for (A), as the exterior regions only grow linearly with the grid size in each dimension, whereas the interior region grows quadratically with (one-dimensional) grid size. Here, we can also see that the proportional runtime for the kernels (B), (C), and (D) is very much comparable, as already indicated in the kernel runtime breakdown in Figure 7.

Next, we investigate the effect of "grouping" computational threads, by passing more than one Vanka patch off to single CUDA thread within any one of the regions where the Vanka patches share the same patch matrix. This reduces the number of overall threads that need to be launched, while potentially further improving the memory accesses required. Figure 11 shows the runtime of the two sets of kernels applying the patch matrices for the four different thread block sizes, grouping patches together in groups of 1 to 64 patches per thread. From Figure 11, we see that increasing the number of patches per thread typically does not lead to a faster runtime; at best, the performance remains relatively constant. Thus, we do not pursue this any further and remain using one thread per patch.

The four Vanka-specific kernels make up more than 75% of the overall runtime of a Vanka iteration (see Figure 7), with each kernel taking up roughly the same proportion of overall runtime. To avoid unnecessary complexity in the code, we choose a single CUDA block size to use for the entire algorithm and all connected kernels. For the smaller problem size, a CUDA block size of 8×8 is not the optimal choice for many of the individual kernels, but all four kernels exhibit near-peak performance for this CUDA block size. For the larger problem size, the best choice of CUDA block size is 12×12 . We have also confirmed that these are the best choices by timing a full iteration of the algorithm for both problem sizes. Table 3 provides a concise overview of the best parameters for both algorithms.

Algorithm	# elements	threads/block
Braess-Sarazin	256^2 1024^2	$\begin{array}{c} 12\times12\\ 12\times12 \end{array}$
Vanka	256^2 1024^2	$\begin{array}{c} 8\times8\\ 12\times12 \end{array}$

Table 3. Best parameter choices for both algorithms.

Roofline model

Having analyzed the kernels above and selected the optimal thread block size, we now consider a roofline model to measure for how efficient the kernels are on a given GPU. Such models tell us whether an operation is memory or compute bound, and whether all theoretically available computing power is used. Figure 12 shows two roofline models, one for each of the two problem sizes, showing measured performance vs. arithmetic intensity for each kernel in a Braess-Sarazin or Vanka relaxation sweep. Here, we very clearly see that, for the larger problem size, most kernels lie right on or very close to the performance bound, meaning that they are running as fast as possible given their arithmetic intensity. In order to improve their performance, we would need to find ways to increase their arithmetic intensity. However, given the nature of these kernels and the underlying structured matrix data structures, there is not an obvious avenue to do this.

Even though most of the kernels (all the matrix-vector and array operations) are clustered together, there are four outliers in this data that we want to highlight:

1. The first outlier is the kernel corresponding to forming the Vanka patch submatrices, which is not visible in the plot, as it consists entirely of memory movements and no floating-point operations. Its arithmetic intensity and floating-point performance are, thus, 0.

- 2. The second outlier is the kernel applying the patch matrix inverse to the exterior patches of the domain. This has a low peak performance of only 6 GFLOP/s, as it consists of mostly small operations (16 unique patch matrices, with 8 one-dimensional regions sharing a patch matrix). It also acts on little enough data that, even for the large problem size, we do not achieve the performance expected from the roofline model.
- 3. The third outlier is a kernel that we mostly ignored in our analysis, the weighted Jacobi kernel. This kernel achieves a higher performance than all but one other kernel, with a peak performance of 753 GFLOP/s. However, it contributes less than 2% to the overall runtime of Braess-Sarazin relaxation, with similar percentages of runtime expected for the other algorithms that use it. Thus, even though its performance is rather high, it has barely any measurable effect on the algorithm runtime.
- 4. The final outlier is the kernel applying the patch matrix inverse to the interior patches of the domain. Its peak performance is roughly 2145 GFLOP/s, almost three times as high as the next highest kernel. With this high performance, it still makes up about 20% of the overall runtime of Vanka. Thus, achieving this performance on this single kernel results in the overall Vanka algorithm achieving much better performance.

Overall, we note that most of the kernels achieve their maximum possible performance, as they lie right on the performance limit in the roofline model for the larger problem size. Due to the nature of their operations, increasing their arithmetic intensity is not possible and, thus, the performance of these algorithms cannot reasonably be expected to be increased. One avenue to consider to improve the performance of the kernels that require a disproportionately large volume of memory movement would be to try to "trade" some memory movement for increasing numbers of floating-point operations; this will be a subject for future research.

Overall solver performance

Finally, having analyzed and optimized the performance of the Vanka and Braess-Sarazin relaxation schemes, we now look to see how they compare in practice, when used as relaxation schemes inside of a multigrid V-cycle that is used as preconditioner for FGMRES applied to the Stokes equations. We will also compare their performances to the performance of FGMRES preconditioned with a multigrid V-cycle with Schur-Uzawa and preconditioned with a Block-Triangular preconditioner with multigrid Vcycles used to approximate the block inverses. The additional parameters needed for Schur-Uzawa and the Block-Triangular preconditioner have been determined through further experiment. The optimal value of t in the Schurcomplement scheme is 1 with an optimal Jacobi weight of $\omega = 0.4$. For the Block-Triangular preconditioner, we determined that a total of 3 V-cycles are necessary for both the pressure update solve and velocity update solve, and the two respective weights for the weighted Jacobi relaxation are



Figure 11. Vanka-specific kernels: Group size vs. runtime for 256² elements (left) and 1024² elements (right).



Figure 12. Roofline Model for all kernels

 $\omega = 0.6$ for the pressure update solve, and $\omega = 1.0$ for the velocity update solve.

The choice of outer Krylov method for a linear solve requires considering many factors. Here, because the preconditioners are not guaranteed to be symmetric and positive definite, we must use a general Krylov method instead of a specialized technique like CG or MINRES. We choose to use FGMRES for two reasons. First, we find that right preconditioning is a preferable framework to left preconditioning, since it does not change the norm of the underlying minimization. Secondly, all of the chosen components lead to preconditioned FGMRES algorithms that converge in tens of iterations, so the additional memory costs for vector storage in FGMRES are feasible (even on the GPU) and preferable to the added computational cost of an extra preconditioner application that is needed in classical GMRES. We note, however, that none of the conclusions from this study would be substantially changed by using classical right-preconditioned GMRES.

We use our own implementation of FGMRES, making use of our structured data structures, and use a multigrid V(1,1) cycle as preconditioner, and a V(3,3) cycle as part of the Block-Triangular preconditioner. At each level of the multigrid algorithm, we use a sweep of either Braess-Sarazin, Vanka, or Schur-Uzawa relaxation. With the Block-Triangular preconditioner, we use three sweeps of weighted Jacobi. At the coarsest level, we use either an exact solve on the CPU or three sweeps of the relaxation scheme on the GPU.

The first comparison we consider is a comparison of Braess-Sarazin to both our tuned Vanka implementation described in this paper and a simple Vanka implementation, shown in Figure 13. All of the results are for problems of size 1024^2 , with the exception of the simple Vanka runs. Due to its higher memory requirements, the largest problem size that successfully ran was a problem of size 768^2 elements. However, even though simple Vanka has just over half as many elements as the other approaches, it is still not able match their performance. On the CPU, we see that multigrid



Figure 13. Comparing Vanka to Braess-Sarazin for a problem of size 1024^2 elements (768^2 elements for simple Vanka)

with Braess-Sarazin relaxation strongly outperforms the use of Vanka relaxation. Even though multigrid with Vanka relaxation typically requires one fewer iteration to reach convergence, the work required per iteration is significantly larger than for Braess-Sarazin, resulting in multigrid with Vanka taking about twice as long. On the GPU, however, we are able to take advantage of the throughput of Vanka, resulting in a runtime that is more than 23 times smaller than on the CPU, whereas the runtime for Braess-Sarazin is only reduced by a factor of about 11. Overall, on the GPU, tuned Vanka outperforms Braess-Sarazin by about 10%.

Next, we compare Braess-Sarazin and tuned Vanka to the other two preconditioning strategies, monolithic multigrid with Schur-Uzawa and the Block-Triangular preconditioner, shown in Figure 14. We can see that both the multigrid



Figure 14. Comparing all four preconditioning strategies for a problem of size 1024^2 elements

preconditioner with Schur-Uzawa relaxation and the Block-Triangular preconditioner are not able to match the performance of both Braess-Sarazin and Vanka. Initially they perform very well, in particular the Block-Triangular preconditioner, but they soon slow down requiring up to more than 3 times as long as Braess-Sarazin and Vanka (on the GPU).

The third metric to consider is the performance of our tuned Vanka implementation for two problem sizes when normalized per element on the CPU and per row of elements on the GPU; this is shown in Figure 15. We observe that the time for the tuned Vanka implementation (per element) remains the same no matter the problem size on the CPU, requiring about 0.03ms per element. Thus, there is no

runtime [s] **Figure 15.** Showing the work per element (in ms) of tuned Vanka for 256^2 and 1024^2 elements.

additional overhead introduced by the size of the problem. On the GPU, we are able expose the fine-grained parallelism in the tuned Vanka implementation, resulting in a constant work per row of elements at just over 1 ms. In fact, we are able to perform about 10% faster per row of elements for the larger problem size.

These results show that a careful implementation of Vanka on the GPU not only results in the fastest time to convergence, but it also does so without requiring additional parameters to be set. In addition, the parallelism of Vanka makes it a clear favorite in distributed memory settings.

Conclusions and Future Work

tigrid Several preconditioners for FGMRES are well-known to yield scalable solution algorithms for saddle-point problems, such as the Stokes equations, including both monolithic Schur-U**zawk**tigrid and multigrid-based block-factorization precon-

Block-Tditigutars. Here, we consider their implementation, performance, and optimization, on modern CPU and GPU architectures. Different metrics were presented and analyzed, including arithmetic intensity, performance, and runtime, for the various kernels making up these algorithms. Given a highly structured setup, we show that multigrid with Vanka relaxation can be very performant on the GPU, leading to faster convergence than when using Braess-Sarazin, both in terms of iterations (saving just 1 iteration) and runtime (up to

10% faster). This shows that using Vanka relaxation is both mathematically and computationally competitive, although a careful design of the algorithm is warranted. This also highlights the benefit of using GPUs for such algorithms, as multigrid with Vanka on the GPU is up to 23 times faster than on the CPU, while multigrid with Braess-Sarazin is up to 11 times faster.

We also presented two other preconditioning strategies, multigrid preconditioner with Schur-Uzawa relaxation, and Block-Triangular preconditioner with multigrid and weighted Jacobi within. Both of these have been shown to not be able to compete with multigrid with Braess-Sarazin or Vanka, in particular on the GPU. In addition, they introduce additional parameters that need to be carefully chosen.

Future work includes extending this work to cases where the tuned Vanka approach is not applicable, such as for linearizations of the Navier-Stokes equations. It is also not clear how well these results generalize to other discretizations of saddle-point systems (including both higher-order discretizations using generalized Taylor-Hood elements and other discretizations, such as using discontinuous Galerkin methods). Extensions to threedimensional incompressible flow problems and other saddlepoint systems are also important future work. One such system of interest, for example, that combines some of these difficulties is the Reynolds-Averaged Navier-Stokes (RANS) equations, in the context of wind-turbine simulations.

Acknowledgments

This work used the Delta system at the National Center for Supercomputing Applications through allocation CIS230037 from the Advanced Cyberinfrastructure Coordination Ecosystem: Services & Support (ACCESS) program, which is supported by National Science Foundation grants #2138259, #2138286, #2138307, #2137603, and #2138296. The work of Scott MacLachlan was partially supported by an NSERC Discovery Grant.

Resources

The implementation of FGMRES and multigrid with Braess-Sarazin, Vanka, and Schur-Uzawa relaxation, and the Block-Triangular preconditioner, as well as problem setup and all data structures is available at: https://gitlab.com/luspi/gmres.

References

- Adler J, He Y, Hu X, MacLachlan S and Ohm P (2023) Monolithic multigrid for a reduced-quadrature discretization of poroelasticity. *SIAM J. Sci. Comput.* 45(3): S54–S81.
- Adler JH, Benson TR, Cyr EC, MacLachlan SP and Tuminaro RS (2016) Monolithic multigrid methods for two-dimensional resistive magnetohydrodynamics. *SIAM J. Sci. Comput.* 38(1): B1–B24. DOI:10.1137/151006135. URL http://dx.doi. org/10.1137/151006135.
- Adler JH, Benson TR and MacLachlan SP (2017) Preconditioning a mass-conserving discontinuous Galerkin discretization of the Stokes equations. *Numerical Linear Algebra with Applications* 24(3): e2047. DOI:10.1002/nla.2047.
- Anzt H, Tomov S and Dongarra JJ (2014) Implementing a sparse matrix vector product for the SELL-C/SELL-C- σ formats on NVIDIA GPUs. URL https://api.semanticscholar.org/CorpusID:61878276.
- Ayuso de Dios B, Brezzi F, Marini LD, Xu J and Zikatanov L (2014) A simple preconditioner for a discontinuous Galerkin method for the Stokes problem. J. Sci. Comput. 58(3): 517–547. DOI: 10.1007/s10915-013-9758-0. URL http://dx.doi.org/ 10.1007/s10915-013-9758-0.
- Benzi M, Golub G and Liesen J (2005) Numerical solution of saddle point problems. Acta Numer. 14: 1–137.
- Bienz A, Falgout RD, Gropp W, Olson LN and Schroder JB (2016) Reducing parallel communication in algebraic multigrid through sparsification. SIAM Journal on Scientific Computing 38(5): S332–S357.
- Bienz A, Gropp WD and Olson LN (2020) Reducing communication in algebraic multigrid with multi- step node aware communication. *The International Journal of High Performance Computing Applications* 34(5): 547–561.

- Braess D and Sarazin R (1997) An efficient smoother for the Stokes problem. *Appl. Numer. Math.* 23(1): 3–19. DOI: 10.1016/S0168-9274(96)00059-1. URL http://dx.doi. org/10.1016/S0168-9274(96)00059-1. Multilevel methods (Oberwolfach, 1995).
- Brandt A and Dinar N (1979) Multigrid solutions to elliptic flow problems. In: PARTER SV (ed.) Numerical Methods for Partial Differential Equations. Academic Press. ISBN 978-0-12-546050-7, pp. 53–147. DOI: https://doi.org/10.1016/B978-0-12-546050-7.50008-3. URL https://www.sciencedirect.com/science/ article/pii/B9780125460507500083.
- Briggs WL, Henson VE and McCormick SF (2000) A Multigrid Tutorial. Philadelphia: SIAM Books. Second edition.
- Dendy JE (1982) Black box multigrid. J. Comput. Phys. 48: 366– 386.
- Dou Y and Liang ZZ (2023) A class of block alternating splitting implicit iteration methods for double saddle point linear systems. *Numerical Linear Algebra with Applications* 30(1): e2455. DOI:https://doi.org/10.1002/nla. 2455. URL https://onlinelibrary.wiley.com/ doi/abs/10.1002/nla.2455.
- Elman H, Silvester D and Wathen A (2005) Finite elements and fast iterative solvers: with applications in incompressible fluid dynamics. Numerical Mathematics and Scientific Computation. New York: Oxford University Press. ISBN 978-0-19-852868-5; 0-19-852868-X.
- Ershkov SV, Prosviryakov EY, Burmasheva NV and Christianto V (2021) Towards understanding the algorithms for solving the Navier–Stokes equations. *Fluid Dynamics Research* 53(4): 044501. DOI:10.1088/1873-7005/ac10f0. URL https:// dx.doi.org/10.1088/1873-7005/ac10f0.
- Farrell PE, He Y and MacLachlan SP (2021) A local Fourier analysis of additive Vanka relaxation for the Stokes equations. *Numerical Linear Algebra with Applications* 28(3): e2306. DOI:https://doi.org/10.1002/nla. 2306. URL https://onlinelibrary.wiley.com/ doi/abs/10.1002/nla.2306.
- Greif C and He Y (2021) A closed-form multigrid smoothing factor for an additive Vanka-type smoother applied to the Poisson equation.
- He Y and MacLachlan S (2019) Local Fourier analysis for mixed finite-element methods for the Stokes equations. *Journal of Computational and Applied Mathematics* 357: 161–183.
- He Y and MacLachlan S (2020) Two-level Fourier analysis of multigrid for higher-order finite-element discretizations of the Laplacian. *Numer. Linear Alg. Appl.* 27(3): e2285.
- John V and Tobiska L (2000a) A coupled multigrid method for nonconforming finite element discretizations of the 2D-Stokes equation. *Computing* 64(4): 307–321. DOI:10. 1007/s006070070027. URL http://dx.doi.org/10. 1007/s006070070027. International GAMM-Workshop on Multigrid Methods (Bonn, 1998).
- John V and Tobiska L (2000b) Numerical performance of smoothers in coupled multigrid methods for the parallel solution of the incompressible Navier-Stokes equations. *International Journal For Numerical Methods In Fluids* 33(4): 453–473.
- Larin M and Reusken A (2008a) A comparative study of efficient iterative solvers for generalized Stokes equations. *Numer.*

Linear Algebra Appl. 15(1): 13–34. DOI:10.1002/nla.561. URL http://dx.doi.org/10.1002/nla.561.

- Larin M and Reusken A (2008b) A comparative study of efficient iterative solvers for generalized Stokes equations. *Numer. Linear Algebra Appl.* 15(1): 13–34. DOI:10.1002/nla.561. URL http://dx.doi.org/10.1002/nla.561.
- Maitre JF, Musy F and Nignon P (1984) A fast solver for the Stokes equations using multigrid with a UZAWA smoother. In: Braess D, Hackbusch W and Trottenberg U (eds.) Advances in Multi–Grid Methods, Notes on Numerical Fluid Mechanics, volume 11. Braunschweig: Vieweg, pp. 77–83.
- Munch P and Kronbichler M (2023) Cache-optimized and low-overhead implementations additive Schwarz of methods for high-order FEM multigrid computations. The International Journal of High Performance Computing **Applications** 10943420231217221DOI: 10.1177/10943420231217221. URL https://doi. org/10.1177/10943420231217221.
- Nataf F and Tournier PH (2022) Recent advances in domain decomposition methods for large-scale saddle point problems. *Comptes Rendus. Mécanique* DOI:10.5802/crmeca.130. Online first.
- Notay Y (2019) Convergence of some iterative methods for symmetric saddle point linear systems. SIAM Journal on Matrix Analysis and Applications 40(1): 122–146. DOI:10. 1137/18M1208836. URL https://doi.org/10.1137/ 18M1208836.
- Paisley M and Bhatti N (1998) Comparison of multigrid methods for neutral and stably stratified flows over two-dimensional obstacles. J. Comput. Phys. 142(2): 581–610. DOI:10.1006/ jcph.1998.5915. URL https://doi.org/10.1006/ jcph.1998.5915.
- Reisner A, Berndt M, Moulton JD and Olson LN (2020) Scalable line and plane relaxation in a parallel structured multigrid solver. *Parallel Computing* 100: 102705. DOI:https://doi.org/10.1016/j.parco.2020.102705. URL https://www.sciencedirect.com/science/ article/pii/S0167819120300922.
- Reisner A, Olson LN and Moulton JD (2018) Scaling structured multigrid to 500K+ cores through coarse-grid redistribution. *SIAM Journal on Scientific Computing* 40(4): C581–C604. DOI:10.1137/17M1146440. URL https://doi.org/10. 1137/17M1146440.
- Rice JR and Boisvert RF (1985) Solving elliptic problems using *ELLPACK*. Berlin, Heidelberg: Springer-Verlag. ISBN 0387909109.
- Trottenberg U, Oosterlee CW and Schüller A (2001) *Multigrid*. London: Academic Press.
- ur Rehman M, Geenen T, Vuik C, Segal G and MacLachlan SP (2011) On iterative methods for the incompressible Stokes problem. *International Journal for Numerical Methods in Fluids* 65(10): 1180–1200. DOI:https://doi.org/10.1002/fld. 2235. URL https://onlinelibrary.wiley.com/ doi/abs/10.1002/fld.2235.
- Vanka S (1986) Block-implicit multigrid solution of Navier-Stokes equations in primitive variables. Journal of Computational Physics 65(1): 138-158. DOI: https://doi.org/10.1016/0021-9991(86)90008-2. URL https://www.sciencedirect.com/science/ article/pii/0021999186900082.

- Voronin A, He Y, MacLachlan S, Olson LN and Tuminaro R (2022) Low-order preconditioning of the Stokes equations. Numerical Linear Algebra with Applications 29(3): e2426. DOI:https://doi.org/10.1002/nla.2426. URL https://onlinelibrary.wiley.com/doi/abs/ 10.1002/nla.2426.
- Zulehner W (2000) A class of smoothers for saddle point problems. *Computing* 65(3): 227–246.