# lbmpy: Automatic code generation for efficient parallel lattice Boltzmann methods

Martin Bauer
*Chair for System Simulation,*
*FAU Erlangen-Nürnberg*
martin.bauer@fau.de

Harald Köstler
*Chair for System Simulation,*
*FAU Erlangen-Nürnberg*
harald.koestler@fau.de

Ulrich Rüde
*Chair for System Simulation,*
*FAU Erlangen-Nürnberg*
*CERFACS, 31057 Toulouse Cedex 1, France*
ulrich.ruede@fau.de

*Abstract*—Lattice Boltzmann methods are a popular mesoscopic alternative to classical computational fluid dynamics based on the macroscopic equations of continuum mechanics. Many variants of lattice Boltzmann methods have been developed that vary in complexity, accuracy, and computational cost. Extensions are available to simulate multi-phase, multi-component, turbulent, and non-Newtonian flows. In this work we present *lbmpy*, a code generation package that supports a wide variety of different lattice Boltzmann methods. Additionally, *lbmpy* provides a generic development environment for new schemes. A high-level domain-specific language allows the user to formulate, extend and test various lattice Boltzmann methods. In all cases, the lattice Boltzmann method can be specified in symbolic form. Transformations that operate on this symbolic representation yield highly efficient compute kernels. This is achieved by automatically parallelizing the methods, and by various application-specific automatized steps that optimize the resulting code. This pipeline of transformations can be applied to a wide range of lattice Boltzmann variants, including single- and two-relaxation-time schemes, multi-relaxation-time methods, as well as the more advanced cumulant methods, and entropically stabilized methods. *lbmpy* can be integrated into high-performance computing frameworks to enable massively parallel, distributed simulations. This is demonstrated using the WALBERLA multiphysics package to conduct scaling experiments on the SuperMUC-NG supercomputing system on up to 147 456 compute cores.

## I. INTRODUCTION

Computational science and engineering is an interdisciplinary field. The workflow of creating computational models starts at (physical) reality and ends with the production of efficiently executable computer code [47]. However, on the route from physical phenomena to machine code lies the formulation of mathematical models and their discretization, the construction of time stepping schemes and solution methods, the design and analysis of parallel algorithms, the realization of complex software systems, and finally the transformation to code that can be executed on a given hardware. The target computer architecture may be a massively parallel system, possibly heterogeneous and using accelerators that can only be exploited by special programming techniques. During the development, many choices must be taken and alternatives considered. Thus, creating computational science software is a work-intensive, time-consuming, and error prone task, whose complexity is easily underestimated, despite its fundamental relevance for extracting reliable predictions from scientific principles. In this article, we will present progress towards the systematic design

of scientific software based on the automatic derivation of methods including the automatic design of efficient software. When designing scientific software with conventional programming techniques, it is often difficult to find the right balance between the flexibility of an approach and its performance, since these are often conflicting goals. Often the specialization to a restricted class of problems would permit special optimizations and can thus lead to more efficient codes. However, a flexible software design may lead to more extensible and more generally usable software. Additionally, using general-purpose libraries for subtasks, such as for the solution of linear systems, may reduce development time, but may also lead to overheads when special structures that would lead to faster algorithms can not be exploited.

Furthermore, it is typical that computational software outlives the computer systems that it was originally designed for. This leads to the problem of performance portability. The choice of a specific algorithm and a specific software design may have been good for the efficient execution on older computer architectures, but these design choices may turn out to be a major bottleneck on the accelerator-based architectures that dominate high-end computing today. A complete rewrite would be necessary, but since this is too time-consuming and expensive. Legacy code today may stay in use though it severely underperforms on modern hardware.

A generic approach to overcome these difficulties are more advanced abstractions. For example, libraries such as Kokkos [13] or programming systems such as OpenACC [15] present abstractions of fine-granular concurrent execution on modern hardware that can help to alleviate the problems of performance portability. An alternative to such libraries can be metaprogramming techniques and the usage of domain specific languages (DSL). Here machine optimized code is generated utilizing program transformations and compiler technology. Using DSLs opens additional possibilities based on application-specific abstractions. A prominent and successful example for a finite element specific DSL is UFL (unified form language) [1] which is embedded in python. It is e.g. used in FeniCS [40] or Firedrake [45]. These automated computing platforms permit the generation of computational models based on a wide variety of partial differential equations that can be expressed in a DSL. For stencil based computations there exist a number of different DSLs e.g. [54], [31], [29], [27] that work on

structured grids and e.g. the ExaSlang DSL that can also work on block-structured grids [38]. These DSLs succeed in supporting a wide variety of mathematical models with high flexibility. Additionally, most of them also support parallel program execution. Such approaches can exploit application knowledge represented in the DSL design and can thus use optimization techniques that are more powerful than optimizing compilers for general-purpose languages.

In the present article, we will focus alternatively on the lattice Boltzmann method (LBM), as a promising mesoscopic modeling paradigm. Our scope of modeling thus supports kinetic schemes as an alternative to continuum mechanics for fluid dynamics. Kinetic schemes are often performance-hungry, but they also offer a high degree of parallelism. Therefore the scalability and performance on modern computer architectures is a central goal of our work. In particular, we will attempt to reach the performance of the best manually optimized LBM codes by developing a new application-specific code generation technology. To achieve this, we will leverage long term experience in manually optimizing general stencil codes and systematic performance engineering for LBM methods [59], [57], [58], [22], [17]. In particular, we will build on WALBERLA as a state-of-the-art LBM software framework with excellent performance characteristics [5], [39].

However, the current article goes beyond developing tools that help to optimize given kinetic simulation algorithms for a variety of computer architectures. While this alone is also useful, we here extend the code transformation approach to the earlier stages of the computational science workflow. In particular, we point out that the design and derivation of lattice Boltzmann models follows a complex but systematic methodology. The development of a specific LBM for a given application is characterized by various options and choices that the method designer must make. This involves selecting a lattice model, defining momenta and relaxation rates, and many more, as will be elaborated in detail below. Based on these choices, a very wide variety of LBMs can be derived. While any such model could be constructed manually and then optimized using program generation technology, this article goes an essential step further. The manual development of advanced LBM can also be time-consuming and error-prone. We will present here, how the derivation of the models themselves can be conceptualized so that it becomes amenable to automation. In *lbmpy*, the tedious mechanical steps of the LBM development can be performed by automatic symbolic manipulations, saving precious developer time and making the methods more reliable and reproducible.

In summary, *lbmpy* jointly with WALBERLA becomes a computing platform for LBMs that is equivalent to what FEniCS is for finite elements. *lbmpy* helps to realize highly efficient LBM implementations and makes it easy for the developer to experiment with different variants of the methods. Fully grown implementations of many different LBMs can be generated with a single mouse click. We emphasize here specifically that our code generation technology has the unique capability to generate highly efficient parallel code that is ready to run with optimal scalability on the largest supercomputers. The LBM is based on concepts from statistical mechanics. The fluid is modeled on the mesoscopic level using distribution functions that represent the statistical behavior of the particles constituting the fluid. Compared to traditional computational fluid dynamics (CFD), that describes the fluid macroscopically with the Navier-Stokes equations, the mesoscopic description permits higher modeling expressivity, leading to a variety of LBMs, e.g., for porous media or multi-phase flows. Its local data access pattern makes the method very well suited for modern hardware architectures that draw their computational power of ever-increasing concurrency. To run code efficiently on these architectures, parallelism on different levels must be exploited, starting from single-instruction-multiple-data (SIMD) instruction sets, utilizing multiple cores per node with OpenMP, to distributed memory parallelism with MPI. Optimizing LB compute kernels to get the best possible performance requires hardware-specific adaptation. This process costs significant development time. Unfortunately, it must be repeated for each new hardware platform. Furthermore, the optimization process often leads to code that is hard to read and maintain. In practice, this can lead to the effect that for prototyping and development, a slow but flexible code is used, and only a few proven methods are re-implemented in a highly optimized way.

Over time a large variety of LBMs have been proposed [34]. Starting from the simple, but widely used, BGK single-relaxation-time (SRT) operator that relaxes the current state linearly to equilibrium with a single relaxation rate, over two-relaxation-time methods (TRT) [21] to general multi-relaxation-time (MRT) [11] methods. All these methods relax to equilibrium in moment space and can be viewed as special cases of MRT. Then, there are multiple advanced methods like cumulant [20] or entropic LBMs [10], where a different set of statistical quantities is chosen for relaxation or the discrete equilibrium is modified. All LB versions have in common that they are parameterized by a set of relaxation rates, which can either be chosen constant or adapted locally. Computing relaxation parameters from local quantities, e.g., shear rates, is used to implement turbulence models or to simulate non-Newtonian fluids. Entropic KBC-type models [10], [2], [18] can also be placed into this group since they are constructed based on an MRT method with two relaxation rates, where one relaxation rate is chosen subject to a local entropy condition. Additionally, LBMs could be extended with different forcing schemes [53], [26], and the equilibrium could be approximated up to different orders either in the so-called incompressible [28] or the standard compressible version. During the implementation phase, the space of options is growing still. Different storage patterns for the distribution functions can be chosen [3], [19]. Hardware-dependent optimizations, like loop splitting and non-temporal stores, further increase the performance on recent CPU architectures.

Faced with this space of physical models and implementation/optimization options, developers of LB frameworks are faced with two options: Either pick a small set of LB methods and optimize them for a specific target architecture,

repeating the full process when a new LB method or hardware platform must be supported or, trying to abstract and automate the development task. In this work, we show that it is indeed possible to automate many tedious development tasks, like reformulating equations to save floating-point operations, splitting loops for better memory access behavior, or fusing stream and collision kernels. We present the lattice Boltzmann code generation package *lbmpy* that solves these problems by automating large parts of the LBM development and optimization process. Its source code and documentation are available open-source under the GNU AGPLv3 license [1].

*lbmpy* is a system for development of computational fluid dynamics codes based on the LBM. In this regard, it is comparable to other LBM frameworks such as OpenLB [30], [42], Palabos [36], [43], elbe [41], [14], LB3D [24], [50], [37], and HemeLB [25]. Some of these LB frameworks like Sailfish [48] and TCLB [55] also use metaprogramming techniques, mainly for portability to GPUs. While these frameworks support many LB methods, systematic performance evaluations and optimizations are predominantly available for simple LB collision operators like single- and two-relaxation-time methods [60], [61], [58], [22].

In this work, we first give an overview over the specific LBM design workflow and the associated code generation pipeline in section II. Then we describe the formalism for LB method specification in section III that comprises the two upper abstraction layers. The transformation to an algorithmic description and its optimization is discussed in section IV. Finally, we present performance and scaling results in section V.

## II. OVERVIEW: LB METAPROGRAMMING PIPELINE

In this section, we first give an overview of the various abstraction layers of *lbmpy*'s metaprogramming pipeline. In the following sections, we then discuss each layer in detail. All abstraction layers are implemented based on the computer algebra system *sympy*. The code generation system itself is implemented in Python, but the generated code is produced in C/C++/LLVM for CPUs and in CUDA or OpenCL for GPUs. As illustrated in fig. 1 the most abstract layer represents a LBM in $q$-dimensional collision space for a DdQq stencil. This collision space is either based on moments or cumulants [20]. To specify an LB scheme, the user defines a basis of the collision space as well as an equilibrium value and a relaxation parameter for each component. Relaxation parameters do not have to be constants but can be chosen as a symbolic expression that depends on local quantities like shear rates. This permits the formulation of models for non-Newtonian fluids, turbulence models, or entropic stabilization.

*lbmpy* transforms this high-level representation into a symbolic description of the collision operator. The collision operator is stored as a symbolic function $\Omega : \mathbb{R}^q \to \mathbb{R}^q$, where $q$ is the number of discrete distribution functions.

In the next step, the collision operator is mapped to a computational kernel. This stage allows the generation of pure collision kernels as well as fused stream-collide kernels.
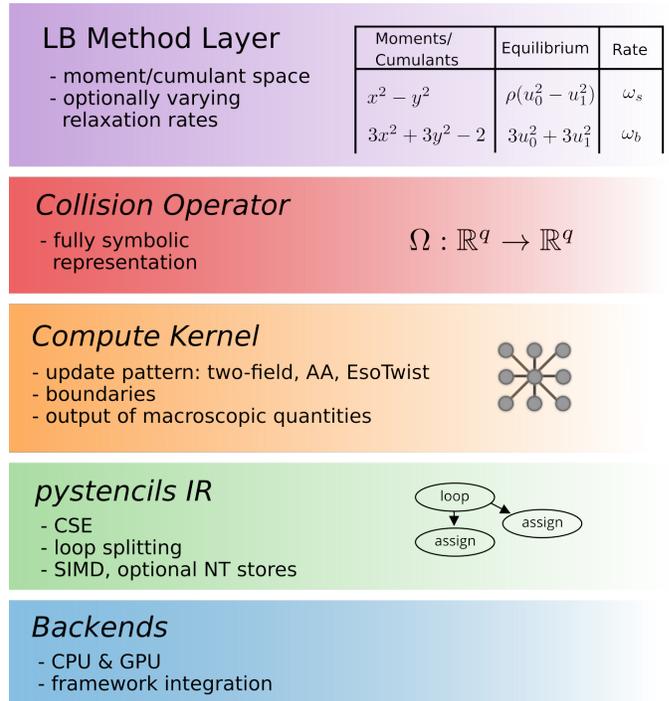
Figure 1: Abstraction layers of metaprogramming pipeline.

Additionally, a concrete storage pattern for the distribution functions is selected, e.g., two-array push or pull patterns or more elaborate single array storage schemes. Optionally one can integrate boundary handling or output of macroscopic quantities at this layer as well. The output of this stage is a symbolic stencil representation of an LB compute kernel. This stencil representation is passed to the *pystencils* package [7] that produces the actual code for CPUs or GPUs. *pystencils* is extended with custom optimization passes to extend the code-generation pass with domain-specific knowledge. In this work we focus on the CPU backend, GPU-specific optimizations and performance results are part of a second publication.

## III. MODEL DESCRIPTION IN COLLISION SPACE

In this section we first give a short overview over the theory that is used to define LBMs on the highest abstraction level. Then we describe the "LB method layer" in more detail followed by examples showing how common collision operators can be specified on this layer.

### A. Collision in moment space

All LBMs considered here discretize the computational domain using a Cartesian, uniformly spaced grid. The discretization is specified in stencil notation as $DdQq$, describing a $d$-dimensional domain where each cell contains $q$ particle distribution functions (PDFs) labeled $f_q(x_i, t)$ with $q \in \{1, ..., q\}$. The PDF $f_q$ represents the mass fraction of particles moving along a lattice velocity $\mathbf{c}_q$. In the following, we use lattice units, i.e., the positions $x_i$ and the time $t$ are integers.

The most basic and probably also the still most commonly used collision operator for LBMs is the single-relaxation-time (SRT)

or BGK operator. It relaxes each population to its equilibrium distribution using a single relaxation parameter $\omega$. The SRT LBM can be succinctly written as

$$f_q(\mathbf{x} + \mathbf{c}_q, t + 1) = \omega\, f_q^{(eq)}(\mathbf{x}, t) + (1 - \omega)\, f_q(\mathbf{x}, t), \quad (1)$$

describing a single time step of the LBM, evolving the system from time $t$ to $t + 1$. It can be split into a local collision step that computes a convex combination of the current state with the equilibrium state and a non-local streaming step that copies PDFs to neighboring cells. The collision is formulated using a relaxation rate $\omega$ which is the inverse of the relaxation time $\tau$, i.e., $\omega = 1/\tau$. The local density $\rho$ and velocity $\mathbf{u}$ are computed from the distribution function as

$$\rho = \sum_q f_q \quad \text{and} \quad u_i = \frac{1}{\rho_0} \sum_q c_{qi} f_q. \quad (2)$$

With these macroscopic values the equilibrium is given as

$$f_q^{(eq)} = w_q \rho + w_q \rho_0 \Big[ \frac{c_{q\alpha} u_\alpha}{c_s^2} + \underbrace{\frac{1}{2 c_s^4} u_\alpha u_\beta (c_{q\alpha} c_{q\beta} - c_s^2 \delta_{\alpha\beta})}_{\text{2nd order}}$$
$$+ \underbrace{\frac{1}{2 c_s^6} \left( u_\alpha u_\beta u_\gamma (c_{q\alpha} c_{q\beta} c_{q\gamma} - c_s^2 c_{q\alpha} \delta_{\beta\gamma}) \right)}_{\text{3rd order}} \Big].$$
$$(3)$$

The reference density $\rho_0$ can either be chosen as 1 to obtain a so-called incompressible LBM [28], for standard compressible LBM set $\rho_0 = \rho$. Note that both versions approximate the incompressible Navier-Stokes equations (NSE) [34]. Equation (3) shows a third-order equilibrium approximation. To obtain the NSE in the macroscopic limit, the equilibrium is required only up to second order in $\mathbf{u}$ [34].

While the BGK collision operator is still widely used in practice, we aim for a more generic description of LBMs. To develop a high-level description of LBMs that is used as the input to the metaprogramming pipeline, we rely on the formalism of multi-relaxation-time methods. It includes the BGK operator and the popular two-relaxation-time method as special cases. With this approach *lbmpy* is able to generate all LB schemes whose collision operator can be written in linear matrix form. This covers the majority of LBMs, with the most notable exception being the cumulant collision operator that will be treated in section III-B.

The MRT formalism also allows us to derive the discrete equilibrium (3) from its continuous counterpart instead of manually specifying it. The collision operator of MRT methods first transforms the PDFs from population space into moment space via a moment matrix $\mathbf{M}$. The components of this moment vector $\mathbf{m} = \mathbf{M} \mathbf{f}$ are relaxed to equilibrium values $\mathbf{m}^{(eq)}$ using a diagonal matrix of relaxation rates $\mathbf{S}$. One collide-stream step of an MRT method then reads

$$f_q(\mathbf{x} + \mathbf{c}_q, t + 1) = \mathbf{M}^{-1} \Big[ \mathbf{S}\, \mathbf{m}^{(eq)} + (\mathbf{I} - \mathbf{S})\, \mathbf{M} \mathbf{f} \Big] \quad (4)$$

with $\mathbf{I}$ denoting the identity matrix. To fully specify the method, we have to define a concrete moment matrix $\mathbf{M}$, a vector

of corresponding equilibrium values $\mathbf{m}^{(eq)}$, and a diagonal relaxation matrix $\mathbf{S}$. We now discuss how each of these three ingredients can be specified in *lbmpy*.

*1) Moment Space:* We begin with the transformation from population space to moment space via the moment matrix $\mathbf{M}$. To derive an invertible transformation $\mathbf{M}$, a set of $q$ independent moments is required. Moments can be identified with polynomials in the lattice velocities $\mathcal{P}(\mathbf{c_q}) : \mathbb{R}^d \to \mathbb{R}$. For example, the zeroth moment, i.e., the density, is given by the constant polynomial 1. The first moments, i.e., the x, y, and z momentum densities, are represented by $c_{qx}$, $c_{qy}$, and $c_{qz}$. Second order polynomials describe viscosity modes, for example in 2D, $c_{qx}^2 + c_{qy}^2$ is a mode related to bulk viscosity, whereas $c_{qx} \cdot c_{qy}$ and $c_{qx}^2 - c_{qy}^2$ are modes related to shear viscosity. The moment value is then computed as

$$\Pi_\mathcal{P}(\mathbf{f}) = \sum_{\mathbf{c_q} \in \mathcal{S}} \mathcal{P}(\mathbf{c}_q) f_q \quad (5)$$

with a stencil $\mathcal{S}$ that is given by a sequence of directions with integer components. Given a sequence of moment polynomials $\mathcal{P}_1 ... \mathcal{P}_k$ the elements of the moment matrix are computed as $M_{kq} = \mathcal{P}_k(\mathbf{c}_q) f_q$. So we have to select $q$ moment polynomials that yield an invertible moment matrix $\mathbf{M}$. In *lbmpy*, there are various options to provide these moment polynomials. The most basic but also most flexible option is to list them explicitly. Then *lbmpy* automatically computes the moment matrix and checks that it is invertible. This option is useful if code for a given MRT method from literature has to be generated.

One central goal of *lbmpy* is to derive LB methods automatically, and not having to pass in, for example, the discrete equilibrium or even the moment basis. Thus, our system additionally offers routines to construct the moment basis for first neighborhood stencils automatically. First, these routines have to find $q$ monomials that lead to an invertible moment matrix $\mathbf{M}$, which can then be orthogonalized in a second step. To illustrate this procedure, we first consider the D3Q27 stencil. Since the velocity vector components only contain the values $\{-1, 0, 1\}$, moments with velocity powers larger than 2 alias a lower order moment. For example

$$\sum_{\mathbf{c}_q \in \mathcal{S}} c_{qi}^3 f_q = \sum_{\mathbf{c}_q \in \mathcal{S}} c_{qi} f_q \quad \text{if } c_{qi} \in \{-1, 0, 1\}. \quad (6)$$

Similarly, moments with even exponents larger or equal than 4 are aliases by corresponding moments with exponent 2. Potentially non-aliased moments are thus $c_{q0}^{e_0} c_{q1}^{e_1} c_{q2}^{e_2}$ with $e_i \in \{0, 1, 2\}$. These systematically constructed 27 monomial moments can be used for the D3Q27 stencil to construct an invertible moment matrix. Similarly, in 2D, this strategy also yields an invertible moment matrix for the D2Q9 stencil. For D3Q15 and D3Q19 the situation is slightly more complex. For D3Q19, we start with the 27 possible monomial moments and discard moments that produce a zero line in the moment matrix, like e.g. the moments defined by the polynomials $c_{q0}^2 c_{q1} c_{q2}$ or $c_{q0} c_{q1} c_{q2}$. For this stencil, there are in total 8 out of the 27 possible monomial moments leading to zero lines, leaving 19 independent rows. For the D3Q15 stencil this procedure

has to be further refined. There, some monomial moments produce the same non-zero row in the moment matrix, which are trivially linear dependent. *lbmpy* groups moments together that yield the same row, resulting in 15 groups. One group of moments, for example is $[c_{q0}^2 c_{q1}, c_{q1} c_{q2}^2, c_{q0}^2 c_{q1} c_{q2}^2]$. In each group, we keep only the lowest order moments. If there is more than one moment remaining, their sum is used. In case of above example this leads to $c_{q0}^2 c_{q1} + c_{q1} c_{q2}^2$. This systematic procedure constructs moment matrices that span the same space as MRT matrices reported in [11], [12], [49].

However, the constructed $q$ independent moments are not orthogonal yet. For MRT methods, typically, an orthogonal moment set is required. Using a symbolic Gram-Schmidt procedure, *lbmpy* can orthogonalize the moments, either utilizing the standard scalar product, or a scalar product weighted by the lattice weights. The exact outcome of the Gram-Schmidt orthogonalization depends on the ordering of the non-orthogonal moments that are put in. For reproducible results *lbmpy* sorts the input by moment order and within each order lexicographically. Before the orthogonalization, also the second-order moments are manually split into bulk and shear part.

*2) Equilibrium State:* The second element necessary for the construction of an MRT method is the equilibrium. It can either be given in population space (3) or directly as a vector of equilibrium moments $\mathbf{m}^{(eq)}$. In this mode, the user has full control over the equilibrium values and can create LBMs not only for the Navier-Stokes equations but also for other partial differential equations.

Beyond this, our meta-programming approach attempts to derive as much as possible from a more general formulation. Therefore, we provide functionality in *lbmpy* to derive equilibrium values for hydrodynamic LBMs automatically. One way to obtain a hydrodynamic discrete equilibrium is to compute the equilibrium moments directly from the continuous Maxwell-Boltzmann distribution

$$f^{(MB)}(\rho, \boldsymbol{u}, \boldsymbol{\xi}) = \frac{\rho}{(2\pi c_s^2)^{\frac{D}{2}}} \exp\left(-\frac{||\boldsymbol{\xi} - \boldsymbol{u}||^2}{2c_s^2}\right) \quad (7)$$

with

$$\int P(\boldsymbol{\xi}) \, f^{(MB)}(\rho, \boldsymbol{u}, \boldsymbol{\xi}) \, d\boldsymbol{\xi}. \quad (8)$$

In this case, the user first chooses a value for the speed of sound $c_s$, typically $c_s = 1/\sqrt{3}$ for first neighborhood stencils, then the integral (8) is evaluated symbolically with the help of *sympy*. The resulting continuous moments of the Maxwellian are used in the equilibrium moment vector $\mathbf{m}^{(eq)}$. Optionally the moments can be truncated to a given order in the macroscopic velocity $\mathbf{u}$. If the equilibrium is required in population space, it can be easily transformed using the assembled moment matrix with $\mathbf{M}^{-1} \mathbf{m}^{(eq)}$. For the cartesian product stencils D2Q9 and D3Q27, this method yields exactly the standard equilibrium (3). For the D3Q15 and D3Q19 velocity sets, however, a different equilibrium is obtained. A comparison of the standard equilibrium and the equilibrium obtained with this moment-matching technique can be found in [9]. The standard

equilibrium for D3Q15 and D3Q19, including the weights, can also be derived by *lbmpy* using a Hermite projection of (7).

*3) Relaxation rates:* The third building block to fully define the method are the relaxation parameters. In *lbmpy*, the user can specify a relaxation rate separately for each of the previously selected moments. Each relaxation rate can either be a constant value, a symbol, or an arbitrary expression of local or neighboring values. In the simplest case, the relaxation rate is a compile-time constant and equal for all time steps and lattice cells. This case allows the computer algebra system to pre-evaluate and simplify expressions containing constants only, thus leading to significant savings. If the relaxation parameter is chosen as a symbol, it becomes a run-time parameter of the generated kernel function. In this case it can be changed, e.g., in a configuration file of the final application without requiring re-compilation. Of course this comes possibly at the cost of executing more FLOPs as compared to the pre-evaluation. The third option, where the relaxation rate is given as a symbolic expression, gives the most modeling power and flexibility. The expression may contain any local or neighboring quantities like equilibrium or non-equilibrium moments. This allows the formulation of a wide range of turbulence models, where the relaxation rate needs to be adapted depending on shear rates. Entropically stabilized schemes like the KBC-type models [10] can also be described in this way. The relaxation rate expression may also contain values of other arrays, allowing for easy coupling of multiple LB schemes, e.g., for multiphase or thermal flows.

### B. Collision in cumulant space

Recently an alternative collision space has been proposed in [20]. Before collision, cumulants of the distribution function are calculated that are relaxed against their respective equilibrium value. Conceptually, cumulant collision operators are realized in *lbmpy* similar to collision operators in moment space. The user specifies a set of cumulants, together with relaxation rates. The cumulant equilibrium values are obtained from the continuous Maxwellian. This allows the formulation of cumulant methods not only for D3Q27 and D2Q9 but also for D3Q19 and D3Q15 stencils.

Cumulants can be succinctly defined through the cumulant-generating function

$$K(\boldsymbol{\xi}) = \ln\left(\sum_{\mathbf{c}_q \in \mathcal{S}} f_q \exp(\mathbf{c_q} \cdot \boldsymbol{\xi})\right). \quad (9)$$

The cumulants are computed by multi-differentiation of (9) and evaluating the derivative at zero. For example, the "bulk cumulant", that we associated with the polynomical $c_{qx}^2 + c_{qy}^2$ is computed as

$$\left.\frac{\partial^2 K(\boldsymbol{\xi})}{\partial \xi_0^2}\right|_{\boldsymbol{\xi}=0} + \left.\frac{\partial^2 K(\boldsymbol{\xi})}{\partial \xi_1^2}\right|_{\boldsymbol{\xi}=0}. \quad (10)$$

Originally, we implemented the cumulant transformation with this approach in *lbmpy*, however, the resulting expressions get very elaborate, especially for large stencils. *sympy*'s common

subexpression evaluation capabilities then run an unpracticable long time manipulating these expressions. Thus we have developed an alternative multi-step transformation, where the populations are first transformed to moment space and then to cumulants. The moments are intermediate quantities that serve as common subexpressions. In *lbmpy*, we use Faà di Bruno's formula [46] to derive the transformation of raw moments to cumulants and vice versa.

### C. Collision Model Examples

In this section, we demonstrate how LBMs can be formulated in *lbmpy* by constructing collision operators of varying complexity.

*1) SRT, TRT:* We start with the single- and two-relaxation-time collision operators. Even if these methods are typically not derived in moment space, we use the MRT formalism for these operators as well to not introduce special cases. The challenge with this general approach is, however, that the simplification system needs to be able to reduce the resulting expressions to their short form.

The following code example shows the definition of a D2Q9 TRT method. Stencils are represented by a tuple of discrete directions with integer components. Common stencils, like the D2Q9, can be obtained by their name. This stencil is used to construct a set of independent raw moments using the algorithm described above.

```
d2q9 = get_stencil("D2Q9")
moments = independent_raw_moments(d2q9)
ω_e, ω_o = symbols("ω_e, ω_o")
ωs = [ω_e if is_even_moment(m) else ω_o
      for m in moments]
m_eq = maxwellian_moments(moments, dim=2,
                          c_s=1/sqrt(3))
trt = create_method(d2q9, moments, ωs, m_eq)
```

In this example, the moment equilibrium values are computed from the continuous Maxwellian, and the relaxation rates are defined for each moment. *lbmpy* offers various classification functions for moments like the is_even_moment function, used here. Other functions can determine the order of a moment, or if it is related to shear or bulk viscosity. Putting these elements together, the method is fully defined and can be displayed to the user in a Jupyter notebook [33] in tabular form, as shown below.

| Moment | Equilibrium | Relaxation rate |
|--------|-------------|-----------------|
| $1$ | $\rho$ | $\omega_e$ |
| $x$ | $\rho u_0$ | $\omega_o$ |
| $y$ | $\rho u_1$ | $\omega_o$ |
| $x^2$ | $\rho u_0^2 + \frac{\rho}{3}$ | $\omega_e$ |
| $y^2$ | $\rho u_1^2 + \frac{\rho}{3}$ | $\omega_e$ |
| $xy$ | $\rho u_0 u_1$ | $\omega_e$ |
| $x^2 y$ | $\frac{\rho u_1}{3}$ | $\omega_o$ |
| $xy^2$ | $\frac{\rho u_0}{3}$ | $\omega_o$ |
| $x^2 y^2$ | $\frac{\rho u_0^2}{3} + \frac{\rho u_1^2}{3} + \frac{\rho}{9}$ | $\omega_e$ |

For better readability we denote moment polynomials using variables $x, y$ and $z$ instead of $c_{qx}, c_{qy}$ and $c_{qz}$. Note that no explicit equilibrium formulation similar to (3) was necessary to construct this method. Only the stencil, the continuous Maxwellian, and a systematically constructed set of independent raw moments have been used to derive this method.

*2) MRT:* Next, we show how to construct a generic MRT method in *lbmpy*. We stick with the D2Q9 stencil to keep the listing of the method tableaus short. Similar to the TRT method above, we start with a set of independent raw moments. For MRT methods, the moments have to be orthogonalized. In *lbmpy* we provide an orthogonalization routine based on the Gram-Schmidt procedure. This routine either uses the standard or a weighted scalar product. A common choice is to use a scalar product weighted with the lattice weights, which we demonstrate in the code example below. If we want to control bulk and shear viscosities using different relaxation rates, the second-order moments must be modified before the orthogonalization. This is handled by the split_shear_bulk_moments function. We pass in the list of all raw moments, containing the second-order moments $x^2, y^2$ and $xy$. These are split into the bulk moment $x^2 + y^2$ and the remaining $xy$ and $x^2 - y^2$ moments. In 3D, this function works analogously.

```
moments = independent_raw_moments(d2q9)
moments = split_shear_bulk_moments(moments)
moments = gram_schmidt(moments, d2q9,
                       weights=get_weights(d2q9))
ω = symbols("ω_:4")
ωs = [0    if get_order(m) < 2   else
      ω[0] if is_shear_moment(m) else
      ω[1] if is_bulk_moment(m)  else
      ω[get_order(m)-1]
      for m in moments]
m_eq = maxwellian_moments(moments, dim=2,
                          c_s=1/sqrt(3))
mrt = create_method(d2q9, moments, ωs,
                    to_incompressible(m_eq))
```

The Gram-Schmidt orthogonalization step then produces the moments listed in the first column of the following table. Then a list is constructed that defines the relaxation rate for each moment. Moments of order less than two are conserved and the relaxation rate can be chosen arbitrarily. In this example, the relaxation rate is set to zero for these moments. Having split up the second order bulk and shear moments, we can pick separate relaxation rates $\omega_0$ and $\omega_1$ for these. In this example, we choose a common relaxation rate for the third- and fourth-order moments.

| Moment | Equilibrium | Relaxation rate |
|--------|-------------|-----------------|
| $1$ | $\rho$ | $0$ |
| $x$ | $u_0$ | $0$ |
| $y$ | $u_1$ | $0$ |
| $x^2 - y^2$ | $u_0^2 - u_1^2$ | $\omega_0$ |
| $xy$ | $u_0 u_1$ | $\omega_0$ |
| $3x^2 + 3y^2 - 2$ | $3u_0^2 + 3u_1^2$ | $\omega_1$ |
| $3x^2 y - y$ | $0$ | $\omega_2$ |
| $3xy^2 - x$ | $0$ | $\omega_2$ |
| $9x^2 y^2 - 3x^2 - 3y^2 + 1$ | $0$ | $\omega_3$ |

Compared to the TRT example, we have done another modification here. The equilibrium moments are modified to yield a so-called incompressible equilibrium [28]. The incompressible equilibrium moments are obtained by writing them as polynomial in the velocity **u** and substituting $\rho = 1$ in all terms that contain at least one velocity component, e.g., $\rho + \rho u_0 \to \rho + u_0$.

Having full information about an LB method in the form of the moment table, as shown above, enables us to analyze the method using a Chapman-Enskog procedure symbolically, as long as relaxation rates are chosen constant. The primary input for this analysis are the moment equilibrium values. The automated analysis can show the user the approximated PDE as well as higher-order error terms. Additionally, it can derive the connection between relaxation parameters and macroscopic parameters, e.g., viscosities. The following snippet shows the analysis of the MRT method defined here.

```
>>> ce = ChapmanEnskogAnalysis(mrt)
>>> ce.get_bulk_viscosity()
-1/9 - 1/(3*ω_1) + 5/(9*ω_0)
>>> ce.get_macroscopic_equations()[0]
∂_t ρ + ∂_0 u_0 + ∂_1 u_1
```

*3) Boundary Conditions:* Similar to the collision operator, boundary conditions are also described in symbolic form. Boundary conditions have to specify the value of a population that is streamed in from a boundary lattice cell. Here is an example of a velocity-bounce-back boundary that models a moving wall.

```
def vel_bounce_back(f, c, method, vel):
    c_s = method.speed_of_sound
    w_q = method.weights[method.stencil.idx(c)]
    vel_term = 2 / c_s**2 * c * v * w_q
    return f.center(c) - vel_term
```

In the boundary definition, the user has access to the method definition, that offers properties like speed of sound or lattice weights. The lattice direction `c` is an integer vector pointing from the fluid to the boundary cell. With this information, an expression for the missing population is constructed. The population field `f` and macroscopic properties can be accessed using relative addressing, where the center is the fluid cell.

Additional information, like in this example, the velocity of the moving wall, can be used. This data can be supplied by various sources. In the simplest case, it is a compile-time constant value. It can also be an expression that depends on spatial coordinates, time step, local population values, or macroscopic quantities. It can also be supplied at runtime. In this case the data is read from a field or a sparse list data structure that stores this information for every connected boundary cell. More details will be covered in the section on the algorithmic treatment of boundary conditions. However, in all these cases, the boundary definition, as shown in the above example, does not change at all. The definition and implementation are strictly separated. Boundary conditions are defined per lattice link, not by lattice cell. That means, for example, that for each link a different velocity can be prescribed.

*4) Turbulence models:* Up to now, we have presented methods with constant relaxation rates. The modeling power of LBMs stems partially from the ability to vary relaxation rates on a cell-by-cell basis, depending on local quantities. With this technique, one can for example, model non-Newtonian fluids or implement turbulence models. To provide this modeling power to the user, *lbmpy* does not only allow for compile- and run-time constants as relaxation rates. It can also take arbitrary expressions of neighboring distribution functions or macroscopic quantities as relaxation rates. We illustrate this for the example of a Smagorinsky subgrid turbulence model. This model adds an eddy viscosity $\nu_t$ to represent energy damping on unresolved scales [32]. The eddy viscosity is calculated from the local strain rate tensor as

$$\nu_t = \underbrace{(C_S\Delta)^2|S|}_{\nu_t} \tag{11}$$

where $C_S$ is a constant and $\Delta$ is a filter length chosen as 1 in lattice coordinates. $|S| = \sqrt{2S_{ij}S_{ij}}$ is the Frobenius norm of the local strain rate tensor

$$S_{ij} = \frac{1}{2}\left(\partial_i u_j + \partial_j u_i\right) = -\frac{3\omega}{2\rho}\Pi_{ij}^{(neq)}. \tag{12}$$

This uses the fortunate property of LBMs that the strain rate tensor can be computed from local quantities, only using the second non-equilibrium moment [35]

$$\Pi_{ij}^{(neq)} = \sum_q c_{qi}c_{qj}\left(f_q - f_q^{(eq)}\right). \tag{13}$$

Equation (12) contains the total relaxation rate $\omega$ which is computed from the total viscosity $\nu = \nu_0 + \nu_t$ which again depends on the eddy viscosity $\nu_t$ that we want to determine

$$\omega = \frac{2}{6C_S^2|S| + 6\nu_0 + 1}. \tag{14}$$

Thus we have a system of two equations in $\omega$ and $|S|$ that we now like to solve for $\omega$. Here it pays off that *lbmpy* is based on the computer algebra system *sympy* where these steps can be performed automatically:

```
S, ω = symbols("|S|, ω", positive=True)

f_neq = pre_collision_symbols() - equilibrium_symbols()
Π = frobenius_norm(second_order_moment_matrix(f_neq))

eqs = [ Eq(ω, ω_from_ν( ν_from_ω(ω_0) + C_S**2 * S )),
        Eq(S, 3 * ω / 2 * Π) ]

effective_ω = solve(eqs, [ω, S])[ω]
```

The resulting symbolic expression for the effective $\omega$ value can be be used in all places where in previous examples a constant has been used. Thus, one can construct MRT or cumulant methods where some or all relaxation rates vary locally, using potentially different expressions for different relaxation rates. For brevity, we have shown here only the construction of a simple turbulence model. The possibility to employ arbitrary expressions as relaxation rates in *lbmpy* can be used to realize also more advanced turbulence models with

only little programming effort. Additionally, *lbmpy* also comes with several pre-defined turbulence models. Thus the user does not have to perform the steps outline above manually, if only a common turbulence model is required.

*5) Entropic KBC Models:* Another important class of models, where relaxation rates are varied locally, are entropic LB schemes. In this section, we show how entropic MRT methods, labeled KBC models by the authors in [10], are realized in *lbmpy*. The central idea of these methods is to maximize a discrete entropy measure $S$ of the post-collision state. The free variable that is tuned to obtain maximum entropy is a relaxation rate associated with higher-order moments. Single relaxation time entropic methods also change the effective viscosity by varying this single rate to maximize entropy. KBC models present an improvement by using two relaxation rates: One rate for the shear moments called $\omega_s$, and a second relaxation rate $\omega_h$ that controls higher order moments. Only $\omega_h$ is changed according to the entropy condition. The shear relaxation rate $\omega_s$ is not altered, and thus also the viscosity remains constant. By choosing which moment is relaxed by which relaxation rate, one obtains different KBC variants, that are labeled by the authors as KBC-N1 up to KBC-N4.

In the original work [10], the notion of mirror states and according relaxation parameters is used. Here we use a different notation that is closer to the formalism of MRT methods. We start with an arbitrary MRT method that uses two symbolic relaxation rates $\omega_s$ and $\omega_h$. The rate $\omega_s$ must include the shear moments if shear viscosity should remain constant. The collision operator in population space is then of the form

$$f_q' = f_q - \omega_s \Delta s_q - \omega_h \Delta s_h, \tag{15}$$

where $f_q'$ is the post-collision state, $f_q$ the pre-collision state, and $\Delta s_q, \Delta h_q$ being the coefficients multiplying the relaxation rates. Then we need to maximize the entropy

$$S(\mathbf{f}') = -\sum_q f_q'(\omega_h) \ln\left(\frac{f_q'(\omega_h)}{f_q^{(eq)}}\right) \tag{16}$$

in every cell at every time step by varying $\omega_h$. Taking the first derivative of (16) w.r.t. $\omega_h$ we get the optimality condition

$$\sum_q \Delta h \left[\ln\left(\frac{f_q'(\omega_h)}{f_q^{(eq)}}\right) + 1\right] = 0. \tag{17}$$

This condition could be solved numerically in every cell using Newton's method. However, in this case, a more efficient way can be devised. We expect $f'$ to be close to $f^{(eq)}$ and approximate the logarithm around $1$ up to first order with $\ln(x) \approx x - 1$. t The optimality condition then simplifies to

$$\sum_q \Delta h \, \frac{f_q'(\omega_h)}{f_q^{(eq)}} = 0. \tag{18}$$

Inserting the post-collision value as $f_q' = f_q - \omega_s \Delta s_q - \omega_h \Delta s_h$, and introducing the entropic scalar product $\langle a, b \rangle_E := \sum_q a_q b_q \left[f_q^{(eq)}\right]^{-1}$, we can solve for $w_h$ and obtain

$$\omega_h = 1 + (1 - \omega_s) \frac{\langle \Delta s, \Delta h \rangle_E}{\langle \Delta h, \Delta h \rangle_E}. \tag{19}$$

To obtain this result, one has to replace $f_q = f_q^{(eq)} + \Delta s_q + \Delta h_q$ and use $\sum_q \Delta h_q = 0$, which holds because of the mass conservation property of the collision operator. All steps leading to (19), are implemented using *sympy*, to obtain an automatic derivation of KBC methods from high-level principles.

Using this technique, we can construct a wide range of entropically stabilized methods, not only for D3Q27 stencils as in [10] but for D3Q19 and D3Q15 stencils as well. Furthermore, we offer a more costly but also more general numerical maximization procedure for the post-collision state entropy, which is based on Newton's method. This can e.g. be used for cumulant methods where the update is not linear in the relaxation rates any more as in (15), but has the quadratic form $f_q' = f_q - a_1\omega_s - a_2\omega_s^2 - b_1\omega_h - b_2\omega_h^2$.

## IV. COMPUTE KERNEL GENERATION

All steps described up to now produce a symbolic representation of the collision operator $\Omega : \mathbb{R}^q \to \mathbb{R}^q$. Together with the stencil, represented as a list of $q$ discrete velocities with integer components, an efficient LB compute kernel must be generated for various hardware platforms. This process is discussed in the following section.

### A. Simplification

To obtain an efficient formulation of the resulting compute kernel, the symbolic collision operator must be rewritten in a form where as few as possible floating point operations (FLOPs) are required to compute post-collision values. This is a very challenging task, since the automatic operator derivation yields a highly inefficient formulation by default. Consider, for example, the case of an SRT model that is derived by transforming populations to moment space, relaxing with a single rate, and transforming back. The matrix products produce lengthy expressions, that are mathematically equivalent to the usual SRT formulation but are now expressed using many more FLOPs. With standard mathematical techniques, like expanding or factoring, terms can already be simplified considerably. However, the most significant reduction in the number of FLOPs is achieved with common subexpression elimination (CSE). General CSE algorithms implemented in computer algebra systems are not guaranteed to find the global optimum and have to rely on heuristics to find a reasonably good solution. These algorithms do not just identify common subtrees as it is typically done as a compiler optimization, but they also try to rewrite the expressions in a form where they have more common subtrees.

For an illustration of the optimization possible in *lbmpy* we present here results for the D3Q19 BGK method. The first row of Table I shows the number of FLOPs in the expressions as they are produced by the automatic derivation. In total, this initial automatically generated code version needs 1263 operations. To reduce cost, we first employ the simplification and CSE capabilities of the *sympy* computer algebra system directly. The results are labeled "*Only CSE*" in the table. This reduces the number of operations significantly, down to only 261. However, a manually optimized implementation of the

|  | Additions | Muls | Divs | Total |
|---|---|---|---|---|
| *Only CSE:* | | | | |
| initial | 686 | 574 | 3 | 1263 |
| sympy CSE | 199 | 61 | 1 | **261** |
| | | | | |
| *Custom:* | | | | |
| initial | 686 | 574 | 3 | 1263 |
| expand | 173 | 423 | 3 | 599 |
| quadratic velocity prod. | 203 | 447 | 3 | 653 |
| expand | 179 | 423 | 3 | 605 |
| factor $\omega$'s | 179 | 305 | 3 | 487 |
| common quadratic term | 131 | 161 | 3 | 295 |
| substitute existing subexpr. | 119 | 119 | 3 | 241 |
| sympy CSE | 119 | 73 | 1 | **193** |

Table I: Detailed simplification results for compressible D3Q19 SRT

BGK method developed by the authors requires only 204 FLOPs. A value around 200 FLOPs is also reported by Wellein et al. [58]. The default simplification and CSE of *sympy* thus is unable to produce code as good as manually tuned, since it needs about 30% more FLOPs than the best solutions known. We also tested the simplification capabilities of other computer algebra systems, including Maple and Mathematica which also could not find simplifications competitive with hand-tuned code. Therefore, it was necessary to develop a new set of custom transformations to rewrite the equations before they are passed to the CSE function of *sympy*. These transformations are listed in the order of application in the lower part of table I as "*Custom*" transformations. Next to the name of the transformation we display the number of FLOPs after the transformation has been applied. Some of these transformations use LBM application knowledge, e.g., they treat density and velocity symbols differently.

We now study these transformations one by one and describe them in detail. The initial formulation is first expanded, i.e., transformed into a sum of products using a function provided by *sympy*. The next transformation called "quadratic velocity products" is specifically developed for LBM simplification. It picks out mixed quadratic terms in macroscopic velocity components $u_i u_j$ and replaces them by $(e^2 - u_i^2 - u_j^2)/2$, with a new subexpression $e := u_i + u_j$. This transformation may seem counterintuitive since it increases the number of FLOPs. However, it helps the following transformations to obtain better results. This is an example of a transformation that requires domain knowledge since this replacement may only be applied to the velocity symbols. Next, a standard expansion is performed. The previously introduced subexpression $e$ prevents this expansion from undoing the previous transformation. The next transformation uses a generic *sympy* function to factor out relaxation rates. The "common quadratic term" step introduces a subexpression that is obtained by taking the expression for the center point, setting all pre-collision values to zero, and relaxation rates to one. For the TRT method this yields $\rho - 3/2\rho(u_0^2 + u_1^2 + u_2^2)$. After this transformation, already existing subexpressions like density and velocity are searched in the equations, and finally, a CSE from *sympy* is performed. With

this custom simplification pipeline, we eventually arrive at a kernel that costs only 193 FLOPs. Note that this is even slightly better than the previously known and carefully hand-optimized version. The improvement compared to using only the generic simplification is 35%. Note also, that the same automatic simplification and optimization steps can now be performed for other LBM methods.

Table II shows the total number of FLOPs for a selection of LB schemes that *lbmpy* can generate and optimize. We compare methods that use the so-called compressible and incompressible equilibrium [28]. We also compare SRT, TRT, MRT, and the SRT with Smagorinsky turbulence model. The table does not display in detail which type of FLOPs each method is composed of. However, all methods only require additions and multiplications with the following exceptions: All compressible models have one division by the density, and the Smagorinsky methods additionally require two square root operations. The results in the first column are obtained by using only the *sympy* CSE. The third column uses the custom simplification strategy introduced above. The second column also uses the custom simplification strategy, with a modified CSE step at the end. In this CSE step, we first search for subexpressions in terms that update opposing lattice directions. By construction, these contain terms that differ in sign only and are good common subexpression candidates. This step is then again followed by a global CSE. This approach is labeled "direction CSE" since lattice directions are taken into account. The lowest FLOP count is marked for each method in boldface.

Let us first discuss the results for all methods with one or two relaxation rates. We see that for SRT and TRT methods the custom simplification pipeline consistently leads to better results. It is generic enough to work not only for the SRT method it was designed for, but leads to good results for TRT methods as well. Also turbulence models built on top of these collision operators are simplified better by the custom strategy, as shown in the table with the Smagorinsky example. Whether the direction-aware CSE is beneficial depends on the stencil. For the D3Q27, it gives the best or equal result across methods, for D2Q9 and D3Q19 it is helpful only for the SRT operators.

We also have chosen two example MRT methods. One, that uses the standard scalar product for moment orthogonalization, and one with moments that are orthogonal w.r.t. to the weighted scalar product. Second order shear and bulk moments are relaxed with different rates, and for each order larger than 2 a separate relaxation rate is chosen. Table II shows that the custom simplification pipeline cannot handle MRT methods. A straight application of CSE obtains much better results for all tested MRT methods, regardless of the stencil.

Currently we employ these three simplification options for each method, and then automatically select the best one. In the future we plan to also use machine learning techniques to optimize the application order of transformations or for finding new transformations.

|  | Only CSE | Custom with direction CSE | Custom, default CSE |
|---|---|---|---|
| *D2Q9* | | | |
| compr. SRT | 113 | **90** | 90 |
| incompr. SRT | 107 | **75** | 75 |
| compr. Smag. | 137 | **122** | 122 |
| compr. TRT | 114 | 110 | **101** |
| incompr. TRT | 108 | 103 | **94** |
| compr. MRT | **150** | 349 | 317 |
| compr. weighted MRT | **153** | 350 | 325 |
| *D3Q19* | | | |
| compr. SRT | 261 | **193** | 193 |
| incompr. SRT | 252 | **162** | 162 |
| compr. Smag. | 306 | **251** | 251 |
| compr. TRT | 262 | 233 | **214** |
| incompr. TRT | 253 | 225 | **206** |
| compr. MRT | **444** | 1098 | 962 |
| compr. weighted MRT | **406** | 947 | 903 |
| *D3Q27* | | | |
| compr. SRT | 444 | **293** | 389 |
| incompr. SRT | 435 | **289** | 346 |
| compr. Smag. | 510 | **370** | 370 |
| compr. TRT | 446 | **379** | 516 |
| incompr. TRT | 437 | **374** | 482 |
| compr. MRT | **651** | 3155 | 4054 |
| compr. weighted MRT | **786** | 3290 | 3984 |

Table II: Total number of FLOPs for different LB schemes. The "Only CSE" columns runs only a CSE from *sympy*. The "Custom with direction CSE" runs the custom simplification pipeline, then a PDF direction-aware CSE followed by a standard CSE. "Custom with default CSE" is the similar, but without the direction-aware CSE.

### B. Collision Operator to Stencil

*1) Streaming and Collision:* After simplification we have the collision operator given as a function $\mathbb{R}^q \to \mathbb{R}^q$. It is represented by a list of $q$ symbolic expressions for the post-collision population values accompanied by a set of subexpressions. The next stage transforms this formulation into a stencil representation.

The stencil representation and all following low-level transformations are part of the *pystencils* package[2] that is also developed by the authors [7]. *pystencils* generates stencil kernels, i.e., routines that iterate over arrays, applying the same operation on every cell. It distinguishes between spatial and index dimensions. Only spatial dimensions are iterated over, while index dimensions are used to address values stored inside a cell, e.g., the $q$ populations for a PDF array or components of a vector field. The central concept of *pystencils* are fields, and field accesses. A field is defined by a name and its number of spatial and index coordinates. Fields are indexed relatively, so the field access `f[1][0](q)`, for example, refers to the $q$'th population value of the east neighbor cell in a 2D setup. *pystencils* is built on top of *SymPy*, and field accesses can be used just like a built-in symbol. The collision operator can be transformed into a stencil representation by replacing the pre- and post-collision symbols by field accesses. Two additional pieces of information are required for this process. The user has to choose the data layout of the population array and a
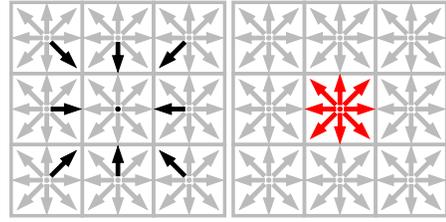
[2]https://i10git.cs.fau.de/pycodegen/pystencils



Figure 2: Visualization of stream-pull-collide update pattern using two arrays. The left part encodes the reads of pre-collision values, the right part shows where post-collision values are written to.

kernel type that describes the operations done inside a kernel call.

*lbmpy* supports three different population storage options. The simplest approach is to have two arrays, where, during one kernel invocation, one array is read-only and the second array is write-only. For this storage pattern the system can generate a fused stream-pull-collide, a fused collide-stream-push, or a collision-only kernel. For a pure LBM simulation that is not coupled to other simulation models, typically a stream-pull-collide kernel gives the best performance. In *lbmpy* the data access patterns are encoded by the field accesses where pre-collision values are loaded from, and field accesses where post-collision values are written to. These are visualized in fig. 2 for a stream-pull-collide kernel. This mechanism cleanly separates the LB method definition from algorithmic- and data structure aspects, avoiding any code duplication.

Besides the simple two array swapping technique, *lbmpy* supports also more advanced storage patterns that operate on a single PDF array and thus require only half the memory. Supported single-array schemes are the AA pattern [3] and the esoteric twist (EsoTwist) update scheme [19]. Single array storage schemes that introduce a data dependency between cell updates like [44] are not supported.
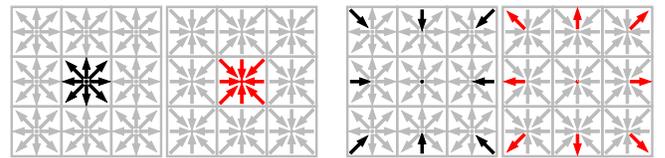


Figure 3: AA update pattern. The two leftmost schematics show the even time step consisting of an in-place collision with inversed storage of populations. The odd time step (right) is a fused stream-pull, collide, stream-push step.

To be able to process all cells in parallel, while only having a single array for PDF storage, the AA pattern needs two different access patterns for even and odd time steps (fig. 3). This also leads to two different kernels that have to be run in an alternating fashion. The different data layout after even and odd steps may complicate boundary handling and coupling the LBM

to other solvers, when traditional implementation techniques are used. With our code generation approach this additional complexity can be handled automatically. The symbolic, high level formulation of method, boundaries, and update scheme is sufficient to e.g. generate boundary handling for even and odd steps automatically.
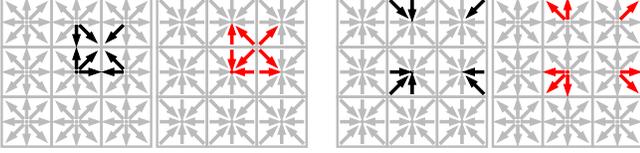


Figure 4: Esoteric twist split in even (left) and odd (right) time step kernels. Black arrows indicate reads, red arrows writes.

The esoteric twist pattern also requires only a single array. In contrast to the AA pattern, it was designed to not require an even and odd time step. If the populations for different lattice velocities are stored in separate arrays, a pointer swapping technique can be used for streaming. In *lbmpy*, however, we do not use this technique, in order to keep a common kernel interface with a single population array. Instead we also use an even and an odd time step for the EsoTwist pattern as well (fig. 4).

*2) Boundary Conditions:* In this section we discuss in more detail how boundary conditions are realized algorithmically. One option is to leave the LB kernel unchanged, and run separate boundary handling kernels before. These kernels prepare the population array by writing values that will be streamed in from boundary cells. *lbmpy* can take symbolic boundary definitions, as shown above, and generate one kernel per boundary. In the simplest case, these boundary kernels operate on a rectangular subdomain, e.g., at the borders of the computational domain. This very simple, but also very common case can thus be handled in the most efficient way possible.

For more general boundary shapes a flag field is used. The flags store a bitmask in every cell that encodes the type of boundary. The flag field is initialized by the user using image/voxel data or with the help of surface meshes. The boundary condition kernel could iterate over the full domain, masking out cells, but especially if boundary conditions are static, this would be rather inefficient. Thus, *lbmpy* also offers an alternative approach, where instead of iterating over all cells, a pre-processing step extracts boundary cell coordinates from the flag field and creates a list of indices for each boundary condition. This index list contains the spatial coordinate of the boundary cell together with the lattice direction to the neighboring fluid cell. Consequently, this list has one entry per boundary link. For each boundary link, custom boundary data can be stored additionally, e.g., the wall velocity for a velocity bounce-back boundary. The flag field then only acts as a convenient way to setup boundaries. During the simulation itself, it is not required

any more since all necessary information is stored in the list data structure to accelerate the processing.

So far we have studied time steps where the boundaries are treated in a separate kernel. With *lbmpy*, boundaries can also be compiled directly into the LB compute kernel. Then a conditional is added to the kernel that determines the cell type either by using a flag field or by a boolean expression that depends on the spatial coordinates.

Using the information about the data access pattern, the symbolic formulation is transformed to obtain a concrete boundary assignment. We point out that it is particularly beneficial to automate the error prone implementation of boundary conditions for the single-array patterns AA or EsoTwist.

All boundary treatment options discussed above are not new since they have already been implemented in existing frameworks or applications. The new contribution here is that code for all these options can be automatically generated avoiding tedious manual coding and debugging. Since all versions can be generated easily, this makes it possible to benchmark all versions and to choose the fastest version for a specific setup. Furthermore, there is no trade-off between flexibility and performance any more. It is now possible to compile a specific boundary treatment into a kernel to get an application-specific implementation with the best performance. Different from a hand-tuned version of the same method, the *lbmpy* approach keeps the system maintainable and extensible. The separation of concerns is realized on the symbolic abstraction level.

### C. Transformations in the Intermediate Representation

In this section we describe low level optimizations executed on the intermediate representation of *pystencils* with the goal to further accelerate the LB compute kernels. *pystencils* is designed as a modular package that allows the user to write custom code transformations specific to the application.

*1) Splitting inner loop:* For LB kernels we expect that the memory interface to be the performance limiting factor, if domains exceed the capacity of the outer level cache. The first optimization we discuss here, aims to increase the maximum attainable bandwidth of the kernel by reducing the number of parallel load/store streams from/to memory. A standard LB kernel iterates over all cells, loads all $q$ pre-collision values at once, computes the post-collision values and stores all $q$ of them. This leads to $q$ parallel load and store streams. Reducing the number of parallel streams to memory can increase the obtained bandwidth [57]. Therefore we develop an automatic transformation that splits the innermost loop into multiple smaller loops. To avoid the re-computation of common subexpressions in every inner loop, buffer arrays are introduced. The first inner loop computes density and velocity and writes them to the buffer arrays. The following loops then handle only two lattice direction updates and have only two parallel load and store streams. Algorithm 1 shows the state after the transformation in pseudo-code. It assumes a simple two-field storage pattern with source and destination array. There are different options on how to exactly perform

**Algorithm 1** Stream-collide kernel with split inner loops

---

**for** all slices $y, z$ **do**
 $\rho\_\text{arr} \leftarrow$ array[x-size]
 u\_arr $\leftarrow$ array[x-size]
 **for** line $x$ **do**
  $f \leftarrow \text{src}[x, y, z]$
  $\rho\_\text{arr}[x] \leftarrow \rho(\mathbf{f})$
  u\_arr[x] $\leftarrow \mathbf{u}(\mathbf{f})$
  dst[x, y, z, center] $\leftarrow \Omega(f, \rho\_\text{arr}[x], \text{u\_arr}[x])$
 **end for**

 **for** line $x$ **do**
  $f \leftarrow \text{src}[x, y, z]$
  dst[x, y, z, east] $\leftarrow \Omega(f, \rho\_\text{arr}[x], \text{u\_arr}[x])$
  dst[x, y, z, west] $\leftarrow \Omega(f, \rho\_\text{arr}[x], \text{u\_arr}[x])$
 **end for**
 **for** line $x$ **do**
  $f \leftarrow \text{src}[x, y, z]$
  dst[x, y, z, north west] $\leftarrow \Omega(f, \rho\_\text{arr}[x], \text{u\_arr}[x])$
  dst[x, y, z, south east] $\leftarrow \Omega(f, \rho\_\text{arr}[x], \text{u\_arr}[x])$
 **end for**
 *... (more loops for remaining directions)*
**end for**

---

this transformation. One free parameter is the number of directions that are updated in the inner loops. In the example, we update two opposing directions at once, since these updates share many common subexpressions. One could also create a separate inner loop for each direction, or group more than two directions together. This transformation is also parametrized by the common subexpressions that are pre-computed in temporary arrays. *lbmpy* can introduce additional temporary arrays for other subexpressions besides density and velocity as well. A heuristic is used to determine subexpressions that are compute intensive enough to justify introducing a temporary array for them. It is important that all temporary arrays fit into the inner level cache so that they do not generate additional pressure on the memory interface. In the simple example, as shown in algorihm 1, the temporary arrays grow with the domain size in $x$-direction. To make this optimization work for arbitrary domain sizes, the inner loop is blocked before splitting it up. The chunk size can be selected such that the arrays fit into L1 cache.

*2) OpenMP and SIMD vectorization:* All LB kernels are designed in a way that cells can be updated in parallel. *pystencils* uses the fact that loop iterations are independent to automatically parallelize the kernel with OpenMP. By default the outer loop is parallelized using a static scheduling strategy. If the domain size is known at compile time, and the outer dimension is very small, *pystencils* uses OpenMP `collapse` to increase the number of parallel interations.

Knowing that iterations are independent, allows *pystencils* to vectorize the code. To have full control over the vectorization process, we do not rely on compiler auto vectorization or pragma-based approaches but generate C code with SIMD

intrinsics. *pystencils* currently supports SSE, AVX, AVX2 and AVX512 vector instruction sets. If the data layout and alignment of the population array is known at compile time, we generate aligned load/store instructions where possible.

*3) Non-temporal stores:* The intrinsics-based vectorization allows us to explicitly use *non-temporal* (NT) stores, also called *streaming stores*, in kernels that use two population arrays. This optimization reduces the total amount of data that has to be transferred from/to memory. By default, modern CPUs have a "write-allocate" or "read for ownership" cache policy [61]. This means that a store operation causes the respective cache line to be read into cache and thus generates twice the memory traffic that is actually required. This actually causes $3q$ values to be transferred over the memory interface per lattice cell. The custom vectorization allows us to change the store instructions from default to streaming stores that bypass the cache. Then only $2q$ values have to be loaded and stored per cell. So this optimization can increase the performance of two-array LB kernels by a factor of 1.5, assuming they are memory-bound and the PDF array does not fit into the outer level cache.

### D. Framework Integration

The intermediate representation of the compute- and boundary kernels is finally transformed by a backend to either C, CUDA or OpenCL code. For each kernel a C function with a well-defined interface is generated. Arrays are passed in as raw pointer, together with shape and stride information that define the memory layout of the arrays. Symbolic quantities that have not been replaced during the code generation process automatically become parameters to the generated C function, e.g., values for relaxation rates or constant external forces.

This simple interface was chosen, such that the generated kernels can be called from a variety of different languages and can be easily integrated into existing frameworks. In this section we describe different ways of utilizing *lbmpy*. The first option allows the user to completely work in a Python environment, preferably an interactive Jupyter notebook for a convenient display of symbolic expressions. There, the user derives the LBM symbolically and passes the method definition to *lbmpy*. After automatic simplification and optimization the generated C/CUDA/OpenCL code is automatically compiled and dynamically loaded as a Python module. The compilation process is fully transparent to the user. The optimized, shared-memory parallel kernel can then be directly called from Python. Data is stored in *numpy* for CPU simulation or in *gpuarray*'s from the *pycuda* package for GPU simulations. In this mode *lbmpy* offers a flexible and fast prototyping environment for LB methods, where simulations can be run on a single node or a single GPU.

For distributed memory parallelization we use the WAL-BERLA framework [16], a multiphysics software system that is optimized for massively parallel simulations with stencil codes. The distributed memory parallelization uses a block-structured domain partitioning based on a forest of octrees and is characterized by excellent scalabilty since it uses no central nor globally shared data structures so that no

global communication is necessary [51]. This fully parallel data structure enables adaptive grid refinement and dynamic load balancing between MPI processes [52], [5]. wALBERLA has Python bindings [8] that allow for simple distributed simulations with *lbmpy* generated kernels directly from Python on a uniform grid. For advanced use cases, e.g., those that require grid refinement, the user has to switch to C++ as the driving language. Integrations of *lbmpy* into the CMake build system of wALBERLA control the generation of LB compute kernels, boundary kernels and packing/unpacking kernels for distributed memory MPI communication. The reason why we also generate communication kernels are the single-field AA and EsoTwist storage patterns. Manually determining what values have to be sent to neighboring processes is tedious and error-prone in these cases. Since the compute kernels are available in symbolic form, we can extract that information and generate the necessary communication routines automatically.

## V. Performance Results

In this section we present benchmark results using the automatically generated LB kernels and compare them to manual implementations with different optimization level.

### A. Single Node Benchmark

*1) Hardware:* We first investigate the single-node performance on two test systems. We scale all kernels on one socket of two Intel Xeon processors with different microarchitecture. The first system is an Intel Xeon E5-2695v3 Haswell system with 14 physical cores per socket. For benchmarking, we deactivate the turbo mode of this processor and set the frequency to a fixed value of 2.3 GHz using the likwid tool suite [56]. The second system is an Intel Xeon Gold 6148 CPU Skylake that has 20 cores per socket with a fixed frequency of 2.4 GHz. The Sub-NUMA clustering features of both systems are switched off to have one NUMA domain per socket. We use transparent huge pages and disable automatic NUMA balancing in the Linux kernel. For all benchmarks we use a domain size of $300 \times 100 \times 100$ that is too large to fit in the outer level cache of any of the tested systems.

To find an upper bound for the possible performance, assuming kernels are memory-bound, we use bandwidth measurements for both systems from [60]. Table III shows the measured `copy` bandwidth for both machines where write-allocates have already been taken into account. Additionally, we use bandwidth measurements of scenarios that closely mimic the memory access behaviour of the D3Q19 LB kernels. Kernels with two-array population storage are compared to a stream benchmark with 19 parallel streams and non-temporal stores, labeled `copy-19-nt-sl`. Kernels with a single-array update pattern are compared to a benchmark that updates 19 values in place called `update-19`. The upper bound for the kernels, as measured in million lattice updates per second (MLUP/s), the bandwidth is divided by the number of bytes that have to be transferred per lattice cell. We assume double precision for all kernels. Thus each cell update requires $2 \cdot q \cdot 8$ bytes per cell.

| processor | Xeon E5-2695v3 | Xeon Gold 6148 |
|---|---|---|
| micro architecure | Haswell | Skylake |
| cores per socket | 14 | 20 |
| frequency | 2.3 GHz | 2.4 GHz |
| *Measured Bandwidths* | | |
| copy | 52.0 GB/s | 102.8 GB/s |
| copy-19-nt-sl | 47.1 GB/s | 92.4 GB/s |
| update-19 | 44.0 GB/s | 93.6 GB/s |

Table III: Test system specification with measured bandwidths from [60]. `copy` uses one load and store stream, write-allocate is already taken into account. `copy-19-nt-sl` uses 19 load and store streams and non-temporal stores. `update-19` updates 19 values in-place.

The benchmark codes are compiled with Intel compiler 19.0.2 if not specified otherwise. Where explicitly noted, the GCC in version 7.4.0 is used. For both compilers we set the optimization flag `-O3`, enable AVX512 on Skylake and AVX2 on Haswell, and switch on fast math flags that allow the compiler to reorder floating point operations. For the intel compiler these are `-fp-model fast=2 -no-prec-sqrt -no-prec-div` and for GCC `-ffast-math`. All kernels are parallelized with OpenMP.

*2) Two-array kernels and comparison to manual implementations:* As discussed before, there is a trade-off between code quality and performance when developing LB kernels manually in C/C++. Code is optimized by specializing it for a particular scenario. To illustrate this trade-off we compare *lbmpy* generated kernels with two manual implementations. The code of the manual implementations can be accessed at [4]. For this test we restrict ourselves to a SRT collision operator. The first manually implemented code is written in a stencil-agnostic way, where lattice velocities and stencil weights are abstracted away through template meta-programming. Theoretically, the compiler should be able to resolve these indirections fully at compile-time. Relaxation rates are also passed in via a templated functor, to enable a flexible integration of turbulence models. While the main aim of this kernel is to be as generic as possible it is still restricted to a single collision operator and a single two-array population storage pattern. But it is easily readable and extensible.

The second manual implementation is written specifically for a D3Q19 stencil. All loops over lattice directions are manually unrolled, expressions are simplified by leaving out multiplications with zero lattice direction components, and common subexpressions are eliminated. These steps lead to code duplication and decreased readability, but may lead to better performance. These optimization steps should not be necessary since the compiler should be able to do them automatically. However, the unrolled stencil-specific version is the basis of further optimization like loop splitting. Figure 5 shows benchmark results on the Skylake system for both manually implemented kernels using GCC and the Intel compiler. We can see that indeed the Intel compiler (right) was able to resolve the compile-time abstractions, such that the generic version is as fast as the stencil-specific one. However, GCC
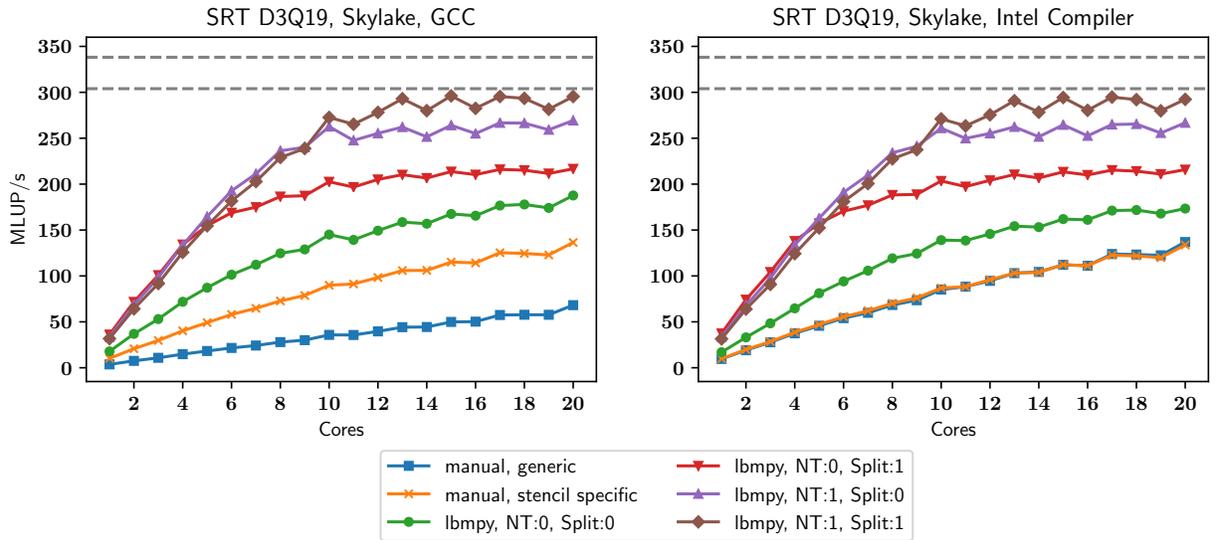
Figure 5: Comparison of kernels with different optimization level on Skylake using a BGK method with two-array population storage on a $300 \times 100 \times 100$ domain. Horizontal lines indicate the roofline estimate using the measured `copy-19-nt-sl` (lower) and `copy` bandwidth (higher).

cannot optimize the generic code automatically, only obtaining about half the performance. The manual implementations scale perfectly, but are far from utilizing the available bandwidth on the system. The generated kernel without loop splitting and non-temporal stores already performs better than the manual implementations. This kernel is explicitly vectorized with AVX512 SIMD intrinsics and uses pointer arithmetic to access the population arrays, whereas the manual implementations use getter/setter methods of an array class. Splitting the inner loops lets the kernel saturate at about 200 MLUP/s. Due to the write-allocate strategy in total 1.5 times more data is moved across the memory interface than necessary. As can be seen in the plot, the performance of this kernel is consequently also about a factor of 1.5 worse than the best kernel with NT stores. Activating NT-stores results in the expected performance of about 300 MLUP/s on this system, very close to the maximal 304 MLUP/s predicted by the roofline estimate obtained with the `copy-19-nt-sl` bandwidth. So the loop splitting and non-temporal stores optimizations are indeed necessary to obtain best possible performance on this system. All manual implementations, where these optimizations have been applied, are lengthy and hard to read, due to the unrolled loops and the usage of SIMD intrinsics. These hand-optimized codes are not only difficult and time-consuming to develop, but also their maintainability and flexibility have been sacrificed for performance. With code generation it is possible to resolve these conflicting goals. Figure 5 also shows, that the generated code performs consistently across different compilers, since all abstractions are already transformed to perform with best possible efficiency by the code generation system, leaving only standard optimizations to the back end compiler.

Figure 6 shows the corresponding results for the Haswell system. Overall the behavior of this older system is similar to
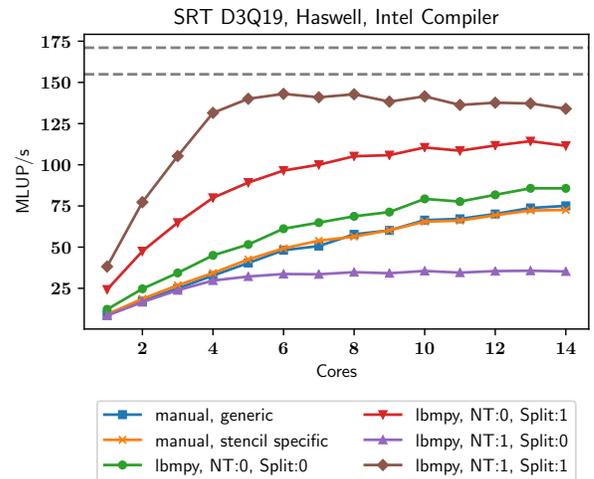


Figure 6: Two-field BGK D3Q19 kernels on Haswell. Configuration and roofline indicators as in fig. 5.

Skylake, with the exception that the system is apparently not able to handle 19 parallel non-temporal store streams, as the version with NT-stores without loop splitting performs very poorly.

*3) Kernels with AA pattern and boundary handling:* Next, we show performance results for single-array kernels that use the AA update pattern. We use the TRT collision operator for these benchmarks. SRT and TRT kernels have very similar performance characteristics, because they have about the same number of FLOPs. In fig. 7 we compare the best two-field version with split loops and non-temporal stores to the corresponding AA kernel. The two-array kernel saturates at about 13 cores, the AA version achieves the highest performance
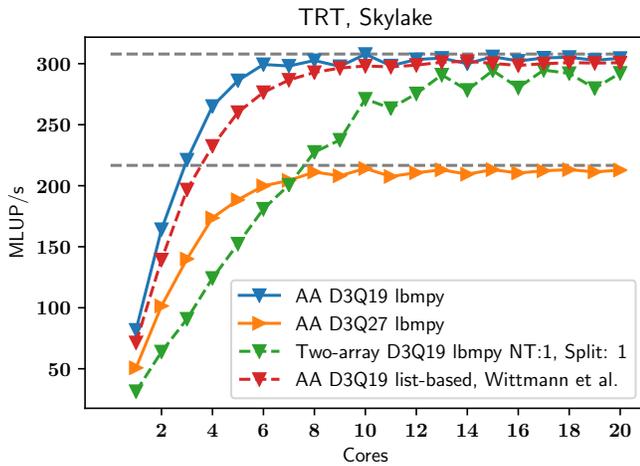
Figure 7: AA pattern, TRT on Skylake. Roofline estimate uses `update-19` bandwidth.



Figure 8: TRT collision operator on Skylake using AA pattern with different boundary options.

already with 6 cores. The additional development effort that is required for the AA pattern pays off not only in half the memory consumption but also in single core performance.

For the AA kernels the NT-store optimization is not applicable, since all values are updated in-place. The inner loop splitting, however, may be beneficial. Since there are two different kernels for even and odd time steps, there are in total four options, where the splitting transformation has been applied to none, only one, or both kernels. We find, that loop splitting does not help in this case. All four options yield almost equal performance results. Thus, fig. 7 only shows the version where neither of the two kernels has been split.

We can also see, that the roofline limit based on the measured `copy-19` bandwidth is a very good model for the performance on the full socket. The D3Q27 version saturates at very close to the expected value, that is by a factor of $19/27$ lower than that of the D3Q19 stencil.

Figure 7 also shows performance results of the TRT LB benchmark kernel by Wittmann et al. [60]. From this benchmark we use the fastest kernel `list-aa-pv-soa` on a channel geometry. It also uses the AA pattern in a SoA layout. In contrast to the *lbmpy* kernels, it operates on a sparse list data structure, such that only populations in fluid cells have to be stored. Also, the benchmark kernels have boundary handling built in, while the *lbmpy* results in fig. 7 show the performance of the compute kernel only.

The boundary handling performance of *lbmpy* is investigated in fig. 8. It shows the *lbmpy* kernels where two different boundary handling approaches are used for a channel scenario, where non-periodic boundaries are set on all sides. The simplest option is to generate separate, external kernels that handle boundaries. Since cache lines containing population at the border must be loaded twice during a time step, the final performance obtained at the full socket is decreased by about 15%. The second approach introduces conditionals in the compute kernel. With this approach we can obtain about the same performance on
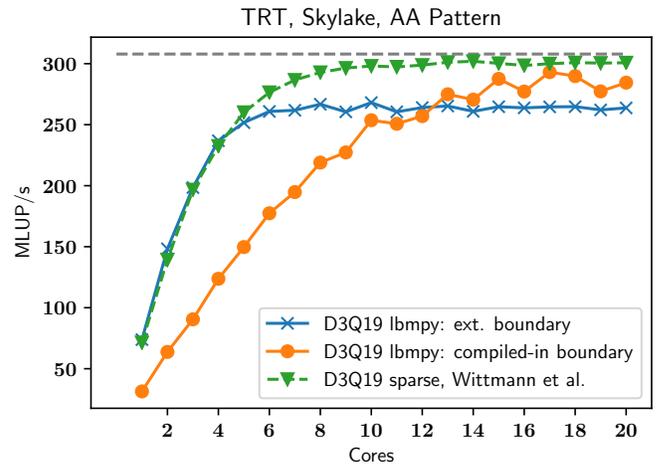
the full socket as the pure compute kernel, but the performance on a few cores is much lower. The intrinsics-based SIMD vectorization in *pystencils* cannot handle the conditionals in an optimal way yet. This limitation is expected to be remedied in future work.

*4) Advanced collision operators:* The goal of code generation in *lbmpy* is not to make a single collision operator fast, but provide a framework that is can obtain good performance for a wide range of different LBMs. Figure 9 shows results for different D3Q19 LB schemes on the test systems using the AA update pattern. The TRT results, we have seen above are included for reference again. A slightly more complex scheme is the BGK operator with included Smagorinsky turbulence model. In this kernel, the relaxation rate is determined on a cell-by-cell basis. The computation of the adapted rate is done on the fly inside the kernel, to not introduce additional memory accesses. Schemes with variable relaxation rates are oftentimes implemented in a way where the rates are computed in a separate kernel and stored into an additional array, for flexibility reasons. This is not necessary in *lbmpy*, so that the Smagorinsky kernel obtains identical performance as the TRT kernel on the full socket. On Skylake also the performance on small core counts is almost identical to the TRT kernel, whereas on Haswell the additional computational complexity, like e.g., the two sqrt operators per cell, lead to lower single core performance compared to TRT.

Next, we investigate performance characteristics of an MRT kernel with weighted orthogonal moments. It has four relaxation rates, two for controlling shear and bulk viscosity separately, one for third, and one for forth order moments. All relaxation rates remain symbolic at compile time and become run time parameters. *lbmpy* is capable to optimize this model so that it almost runs as fast as the TRT model on both test architectures. This is true also for other MRT models that the system can generate, e.g. weighted/unweighted moment orthogonalization or compressible/incompressible equilibria.
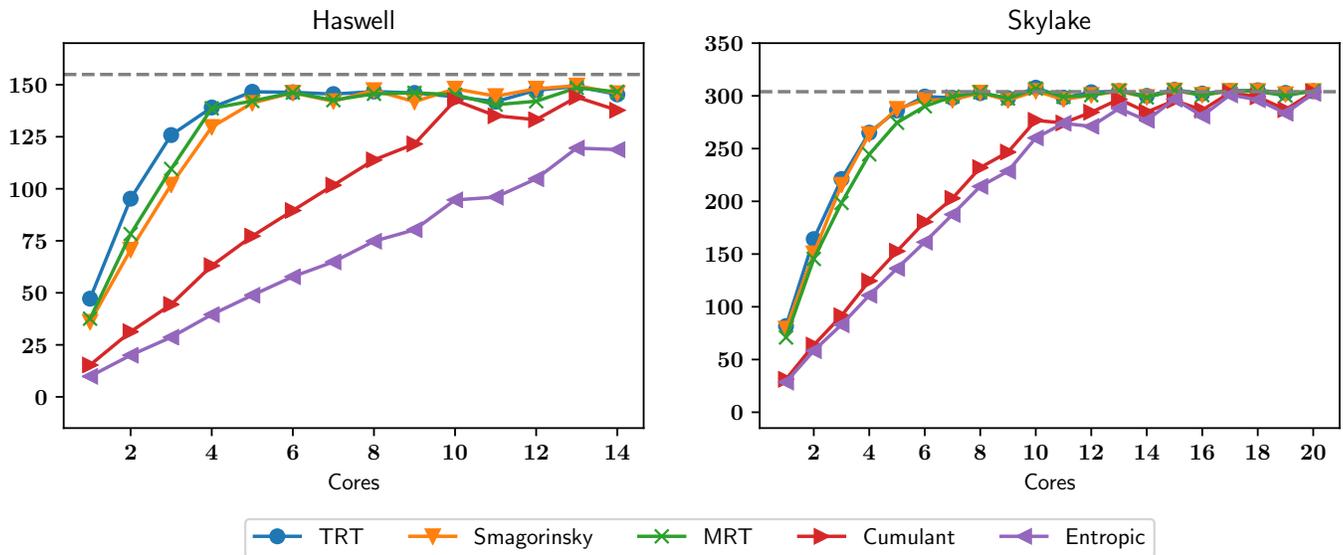
Figure 9: Comparison of different LB collision operators. All kernels use the AA pattern and a D3Q19 stencil.

Figure 9 also contains measurements for a D3Q19 cumulant method. The non-linear transformation to cumulant-space makes this collision operator more compute intensive than MRT methods. Nonetheless, *lbmpy* can optimize the cumulant kernel such that it saturates the available memory bandwidth on both systems. Additionally, we try an entropic method of KBC type. Shear and bulk viscosity is kept fix, the relaxation rate for higher order moments is chosen adaptively to maximize entropy. This method is too compute intensive to be memory-bound on Haswell, but on a full Skylake socket it achieves performance similar to the simpler methods.

Summarizing our findings, after careful optimization there is no performance penalty using complex LB collision operators. On modern CPU architectures all LBM implementaions are memory bound when they are properly optimized. This optimization, however, is only achievable by tedious manual coding by experts or by using automatic code generation with tools like *lbmpy*.

## B. Scaling Benchmark

Integrating the generated *lbmpy* kernels into the WALBERLA framework allows us to run large scale simulations on distributed memory systems. We use the MPI communication capabilities of WALBERLA together with generated serialization/deserialization kernels to run a large parallel simulations. In contrast to previous work [23], [6], where scaling results for manual implementations of TRT two-field kernels have been shown, we demonstrate the performance of a more complex MRT kernel with AA pattern here. As we have shown above, this kernel runs as fast as a SRT or TRT collision operator on the full node when properly optimized.

As test system the SuperMUC-NG supercomputer in Munich is used. It consists out of Intel Xeon Platinum 8174 processors with Skylake architecture. Each node has two sockets with 24
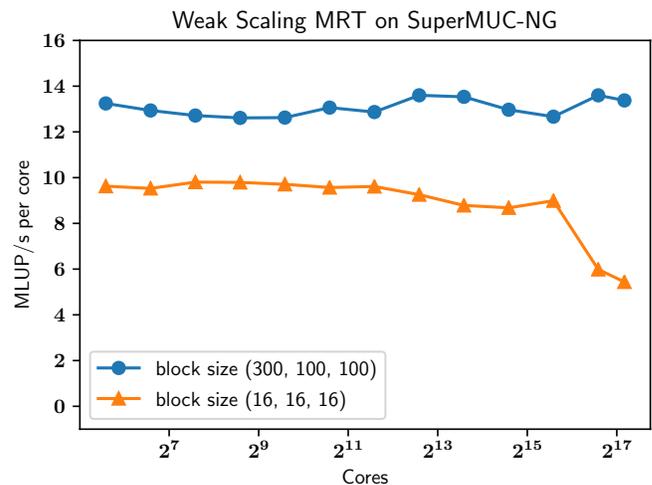


Figure 10: Weighted orthogonal MRT method with AA update pattern scaled on SuperMUC-NG up to 147,456 cores (3072 nodes) using a channel geometry.

physical cores each. We run a weak scaling setup up to the 3072 compute nodes that we have access to. This is about half of the full machine size.

Figure 10 shows the scaling results of a channel geometry. The domain is partitioned into equally sized blocks and each block is assigned to a physical core. The machine is best utilized when choosing large block sizes, since the communication overhead is kept small in this case. For a block size of $(300, 100, 100)$, we observe perfect scalability up to all 147,456 cores used. In this configuration we obtain about 1972 GLUP/s on half of SuperMUC-NG. Besides this weak scaling scenario demonstrating maximal GLUP/s rates, we may alternatively want to maximize the number of time steps per second. To

illustrate the capabilities of *lbmpy* with waLBerla we also study a scaling scenario with much smaller block size. fig. 10 shows the scaling behavior for a block size of $16^3$. Here we achieve still good scalability up to 1024 nodes, then the performance in GLUP/s drops to approximately 40% of the performance that was observed for the large block size. Note however, that in this configuration we can execute 1327 time steps per second. Note also that on the 3072 nodes on SuperMUC-NG this is still for an LBM grid consisting of $6 \times 10^8$ LBM cells.

## VI. Conclusion and Outlook

In this article, we presented a programming system named *lbmpy* that supports the flexible creation of highly optimized parallel LBMs. The scope of *lbmpy* are moment-based MRT methods plus cumulant and entropically stabilized collision operators. All methods can be created with locally varying relaxation parameters so that various turbulence models can be realized or also models for non-Newtonian fluids. *lbmpy* automatically optimizes the compute kernels with domain-specific transformations and it can produce codes that employ memory-efficient single-array population storage. Even for complex LBM models, the kernels generated by *lbmpy* can reach the same performance as manually optimized state-of-the-art TRT implementations of [60]. After automatic optimization, all methods are memory-bound on a recent Skylake system. Thus, also advanced collision operators can be used without performance penalty, provided that large domain sizes are used and memory bandwidth is the bottleneck. Through integration into the HPC framework waLBerla, the *lbmpy*-generated methods can be executed on large scale distributed-memory systems with excellent scalability.

The *lbmpy*/waLBerla programming system supports the computational science workflow. It automizes the tedious and error prone development steps. The support of *lbmpy* is *vertically integrated*: it starts with the development of advanced kinetic schemes, it assists the development of scalable parallel codes, and it includes advanced hardware-specific code optimizations. Note that the derivation of modern LB schemes will usually require expertise in mathematics and physics by specialists in LBMs, while the nodel-level kernel optimization for modern CPU microarchitecture will require a detailed understanding of CPU microarchitecture. In this sense *lbmpy*/waLBerla is an excercise in interdisciplinary co-design to create advanced simulation software for future extreme scale computing. The code-generation paradigm permits a higher level of abstraction than it can be realized by conventional software engineering methods. These new, application-specific methods of abstractions can only be realized with automatic code generation. They help to resolve the fundamental conflict between flexibility of software and the need for hardware specific optimizations. The approach taken in *lbmpy* offers a road to performance portability and thus to improve the sustainability of scientific software.

## References

[1] ALNÆS, M. S. ; LOGG, A. ; ØLGAARD, K. B. ; ROGNES, M. E. ; WELLS, G. N.: Unified form language: A domain-specific language for weak formulations of partial differential equations. In: *ACM Trans. Math. Software* 40 (2014), Nr. 2, S. 9

[2] ANSUMALI, S. ; KARLIN, I. V. ; ÖTTINGER, H. C.: Minimal entropic kinetic models for hydrodynamics. In: *Europhysics Letters* 63 (2003), Nr. 6, 798–804. http://dx.doi.org/10.1209/epl/i2003-00496-6. – DOI 10.1209/epl/i2003–00496–6. – ISBN 0295–5075

[3] BAILEY, P. ; MYRE, J. ; WALSH, S. D. ; LILJA, D. J. ; SAAR, M. O.: Accelerating lattice Boltzmann fluid flow simulations using graphics processors. In: *2009 international conference on parallel processing* IEEE, 2009, S. 550–557

[4] BAUER, M. : *Implementation manual LB kernels*. https://github.com/lssfau/walberla/blob/48bd19800b0c46030ae7f5e510e896a9154d78b8/apps/benchmarks/UniformGridGenerated/ManualKernels.h, 2020. – [Online; accessed 24-January-2020]

[5] BAUER, M. ; EIBL, S. ; GODENSCHWAGER, C. ; KOHL, N. ; KURON, M. ; RETTINGER, C. ; SCHORNBAUM, F. ; SCHWARZMEIER, C. ; THÖNNES, D. ; KÖSTLER, H. u.a.: waLBerla: A block-structured high-performance framework for multiphysics simulations. In: *Computers & Mathematics with Applications* (2020)

[6] BAUER, M. ; EIBL, S. ; GODENSCHWAGER, C. ; KOHL, N. ; KURON, M. ; RETTINGER, C. ; SCHORNBAUM, F. ; SCHWARZMEIER, C. ; THÖNNES, D. ; KÖSTLER, H. ; RÜDE, U. : waLBerla: A block-structured high-performance framework for multiphysics simulations. In: *Computers & Mathematics with Applications* (2020). http://dx.doi.org/https://doi.org/10.1016/j.camwa.2020.01.007. – DOI https://doi.org/10.1016/j.camwa.2020.01.007. – ISSN 0898–1221

[7] BAUER, M. ; HÖTZER, J. ; ERNST, D. ; HAMMER, J. ; SEIZ, M. ; HIERL, H. ; HÖNIG, J. ; KÖSTLER, H. ; WELLEIN, G. ; NESTLER, B. u.a.: Code generation for massively parallel phase-field simulations. In: *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis* ACM, 2019, S. 59

[8] BAUER, M. ; SCHORNBAUM, F. ; GODENSCHWAGER, C. ; MARKL, M. ; ANDERL, D. ; KÖSTLER, H. ; RÜDE, U. : A Python extension for the massively parallel multiphysics simulation framework waLBerla. In: *International Journal of Parallel, Emergent and Distributed Systems* 31 (2016), Nr. 6, 529–542. http://dx.doi.org/10.1080/17445760.2015.1118478. – DOI 10.1080/17445760.2015.1118478

[9] BAUER, M. ; SILVA, G. ; RÜDE, U. : Truncation errors of the D3Q19 lattice model for the lattice Boltzmann method. In: *Journal of Computational Physics* 405 (2020), 109111. http://dx.doi.org/https://doi.org/10.1016/j.jcp.2019.109111. – DOI https://doi.org/10.1016/j.jcp.2019.109111. – ISSN 0021–9991

[10] BÖSCH, F. ; CHIKATAMARLA, S. S. ; KARLIN, I. V.: Entropic multirelaxation lattice Boltzmann models for turbulent flows. In: *Physical Review E - Statistical, Nonlinear, and Soft Matter Physics* (2015). http://dx.doi.org/10.1103/PhysRevE.92.043309. – DOI 10.1103/PhysRevE.92.043309. – ISSN 15502376

[11] D'HUMIÈRES, D. ; GINZBURG, I. ; KRAFCZYK, M. ; LALLEMAND, P. ; LUO, L.-S. : Multiple-relaxation-time lattice Boltzmann models in three dimensions. In: *Philosophical transactions. Series A, Mathematical, physical, and engineering sciences* 360 (2002), Nr. 1792, S. 437–451. http://dx.doi.org/10.1098/rsta.2001.0955. – DOI 10.1098/rsta.2001.0955. – ISBN 1364503X

[12] DÜNWEG, B. ; SCHILLER, U. D. ; LADD, A. J. C.: Statistical mechanics of the fluctuating lattice Boltzmann equation. In: *Physical Review E - Statistical, Nonlinear, and Soft Matter Physics* 76 (2007), Nr. 3, S. 1–10. http://dx.doi.org/10.1103/PhysRevE.76.036704. – DOI 10.1103/PhysRevE.76.036704. – ISSN 15393755

[13] EDWARDS, H. C. ; TROTT, C. R. ; SUNDERLAND, D. : Kokkos: Enabling manycore performance portability through polymorphic memory access patterns. In: *Journal of Parallel and Distributed Computing* 74 (2014), Nr. 12, S. 3202–3216

[14] ELBE: *https://www.tuhh.de/elbe/home.html*. – accessed on 2020-01-21

[15] FARBER, R. : *Parallel programming with OpenACC*. Cambridge : Morgan Kaufmann, 2017

[16] FEICHTINGER, C. ; DONATH, S. ; KÖSTLER, H. ; GÖTZ, J. ; RÜDE, U. : WaLBerla: HPC software design for computational engineering simulations. In: *Journal of Computational Science* 2 (2011), Nr. 2, S. 105–112

[17] FEICHTINGER, C. ; HABICH, J. ; KÖSTLER, H. ; RÜDE, U. ; AOKI, T. : Performance modeling and analysis of heterogeneous lattice Boltzmann simulations on CPU–GPU clusters. In: *Parallel Computing* 46 (2015), S. 1–13

[18] FRAPOLLI, N. ; KARLIN, I. V.: Entropic Lattice Boltzmann Models for Thermal and Compressible Flows. (2017). http://dx.doi.org/10.3929/ethz-a-010782581. – DOI 10.3929/ethz–a–010782581. ISBN 8610828378018

[19] GEIER, M. ; SCHÖNHERR, M. : Esoteric Twist: An Efficient in-Place Streaming Algorithmus for the Lattice Boltzmann Method on Massively Parallel Hardware. In: *Computation* 5 (2017), Nr. 2, 19. http://dx.doi.org/10.3390/computation5020019. – DOI 10.3390/computation5020019. – ISSN 2079–3197

[20] GEIER, M. ; SCHÖNHERR, M. ; PASQUALI, A. ; KRAFCZYK, M. : The cumulant lattice Boltzmann equation in three dimensions: Theory and validation. In: *Computers and Mathematics with Applications* (2015). http://dx.doi.org/10.1016/j.camwa.2015.05.001. – DOI 10.1016/j.camwa.2015.05.001. – ISSN 08981221

[21] GINZBURG, I. ; VERHAEGHE, F. ; D'HUMIERES, D. : Two-relaxation-time lattice Boltzmann scheme: About parametrization, velocity, pressure and mixed boundary conditions. In: *Communications in Computational Physics* 3 (2008), Nr. 2, S. 427–478

[22] GODENSCHWAGER, C. ; SCHORNBAUM, F. ; BAUER, M. ; KÖSTLER, H. ; RÜDE, U. : A framework for hybrid parallel flow simulations with a trillion cells in complex geometries. In: *Proceedings of the International Conference on High Performance Computing, Networking, Storage and Analysis* ACM, 2013, S. 35

[23] GODENSCHWAGER, C. ; SCHORNBAUM, F. ; BAUER, M. ; KÖSTLER, H. ; RÜDE, U. : A Framework for Hybrid Parallel Flow Simulations with a Trillion Cells in Complex Geometries. In: *Proceedings of the International Conference on High Performance Computing, Networking, Storage and Analysis*. New York, NY, USA : ACM, 2013 (SC '13). – ISBN 978–1–4503–2378–9, S. 35:1–35:12

[24] GROEN, D. ; HENRICH, O. ; JANOSCHEK, F. ; COVENEY, P. ; HARTING, J. : Lattice-Boltzmann methods in fluid dynamics: Turbulence and complex colloidal fluids. In: *Jülich Blue Gene/P Extreme Scaling Workshop*, 2011, S. 17

[25] GROEN, D. ; HETHERINGTON, J. ; CARVER, H. B. ; NASH, R. W. ; BERNABEU, M. O. ; COVENEY, P. V.: Analysing and modelling the performance of the HemeLB lattice-Boltzmann simulation environment. In: *Journal of Computational Science* 4 (2013), Nr. 5, S. 412–422. http://dx.doi.org/10.1016/j.jocs.2013.03.002. – DOI 10.1016/j.jocs.2013.03.002. – ISSN 1877–7503

[26] GUO, Z. ; ZHENG, C. ; SHI, B. : Discrete lattice effects on the forcing term in the lattice Boltzmann method. In: *Physical Review E* 65 (2002), Nr. 4, S. 046308

[27] GYSI, T. ; OSUNA, C. ; FUHRER, O. ; BIANCO, M. ; SCHULTHESS, T. C.: STELLA: a domain-specific tool for structured grid methods in weather and climate models. In: *Proc. of the International Conference for High Performance Computing, Networking, Storage and Analysis* ACM, 2015, S. 41

[28] HE, X. ; LUO, L.-S. : Lattice Boltzmann Model for the Incompressible Navier–Stokes Equation. In: *Journal of Statistical Physics* 88 (1997), Nr. 3/4, 927–944. http://dx.doi.org/10.1023/B:JOSS.0000015179.12689.e4. – DOI 10.1023/B:JOSS.0000015179.12689.e4. – ISBN 0022–4715

[29] HENRETTY, T. ; VERAS, R. ; FRANCHETTI, F. ; POUCHET, L.-N. ; RAMANUJAM, J. ; SADAYAPPAN, P. : A stencil compiler for short-vector SIMD architectures. In: *Proc. of the 27th international ACM conference on International conference on supercomputing* ACM, 2013, S. 13–24

[30] HEUVELINE, V. ; LATT, J. : THE OPENLB PROJECT: AN OPEN SOURCE AND OBJECT ORIENTED IMPLEMENTATION OF LATTICE BOLTZMANN METHODS. In: *International Journal of Modern Physics C* 18 (2007), Nr. 04, S. 627–634. http://dx.doi.org/10.1142/S0129183107010875. – DOI 10.1142/S0129183107010875

[31] HOLEWINSKI, J. ; POUCHET, L.-N. ; SADAYAPPAN, P. : High-performance code generation for stencil computations on GPU architectures. In: *Proc. of the 26th ACM international conference on Supercomputing* ACM, 2012, S. 311–320

[32] HOU, S. ; STERLING, J. ; CHEN, S. ; DOOLEN, G. D.: A Lattice Boltzmann Subgrid Model for High Reynolds Number Flows. In: *Fields Institute Communications* 6 (1996), 151–166. http://arxiv.org/abs/comp-gas/9401004

[33] KLUYVER, T. ; RAGAN-KELLEY, B. ; PÉREZ, F. ; GRANGER, B. ; BUSSONNIER, M. ; FREDERIC, J. ; KELLEY, K. ; HAMRICK, J. ; GROUT, J. ; CORLAY, S. ; IVANOV, P. ; AVILA, D. ; ABDALLA, S. ; WILLING, C. : Jupyter Notebooks – a publishing format for reproducible computational workflows. In: LOIZIDES, F. (Hrsg.) ; SCHMIDT, B. (Hrsg.) ; IOS Press (Veranst.): *Positioning and Power in Academic Publishing: Players, Agents and Agendas* IOS Press, 2016, S. 87 – 90

[34] KRÜGER, T. ; KUSUMAATMAJA, H. ; KUZMIN, A. ; SHARDT, O. ; SILVA, G. ; VIGGEN, E. M.: *The Lattice Bolzmann Method - Principles and Practics*. Springer, 2017. – ISBN 978–3–319–44649–3

[35] KRÜGER, T. ; VARNIK, F. ; RAABE, D. : Shear stress in lattice Boltzmann simulations. In: *Physical Review E - Statistical, Nonlinear, and Soft Matter Physics* 79 (2009), Nr. 4, S. 1–14. http://dx.doi.org/10.1103/PhysRevE.79.046704. – DOI 10.1103/PhysRevE.79.046704. – ISBN 1539–3755

[36] LAGRAVA, D. ; MALASPINAS, O. ; LATT, J. ; CHOPARD, B. : Advances in multi-domain lattice Boltzmann grid refinement. In: *Journal of Computational Physics* 231 (2012), Nr. 14, S. 4808 – 4822. http://dx.doi.org/10.1016/j.jcp.2012.03.015. – DOI 10.1016/j.jcp.2012.03.015. – ISSN 0021–9991

[37] LB3D: *http://ccs.chem.ucl.ac.uk/lb3d*. – accessed on 2020-01-21

[38] LENGAUER, C. ; APEL, S. ; BOLTEN, M. ; CHIBA, S. ; RÜDE, U. ; TEICH, J. ; GRÖSSLINGER, A. ; HANNIG, F. ; KÖSTLER, H. ; CLAUS, L. u. a.: ExaStencils–Advanced Multigrid Solver Generation. In: *Software for Exascale Computing-SPPEXA 2nd phase* (2020)

[39] LIU, Z. ; CHU, X. ; LV, X. ; MENG, H. ; SHI, S. ; HAN, W. ; XU, J. ; FU, H. ; YANG, G. : Sunwaylb: Enabling extreme-scale lattice boltzmann method based computing fluid dynamics simulations on sunway taihulight. In: *2019 IEEE International Parallel and Distributed Processing Symposium (IPDPS)* IEEE, 2019, S. 557–566

[40] LOGG, A. ; MARDAL, K.-A. ; WELLS, G. : *Automated solution of differential equations by the finite element method: The FEniCS book*. Bd. 84. Springer Science & Business Media, 2012

[41] MIERKE, D. ; JANSSEN, C. ; RUNG, T. : An efficient algorithm for the calculation of sub-grid distances for higher-order LBM boundary conditions in a GPU simulation environment. In: *Computers & Mathematics with Applications*. (2018). http://dx.doi.org/10.1016/j.camwa.2018.04.022. – DOI 10.1016/j.camwa.2018.04.022. – ISSN 0898–1221

[42] OPENLB: *https://www.openlb.net/*. – accessed on 2020-01-21

[43] PALABOS: *http://www.palabos.org/*. – accessed on 2020-01-21

[44] POHL, T. ; KOWARSCHIK, M. ; WILKE, J. ; IGLBERGER, K. ; RÜDE, U. : Optimization and profiling of the cache performance of parallel lattice Boltzmann codes. In: *Parallel Processing Letters* 13 (2003), Nr. 04, S. 549–560

[45] RATHGEBER, F. ; HAM, D. A. ; MITCHELL, L. ; LANGE, M. ; LUPORINI, F. ; MCRAE, A. T. ; BERCEA, G.-T. ; MARKALL, G. R. ; KELLY, P. H.: Firedrake: automating the finite element method by composing abstractions. In: *ACM Transactions on Mathematical Software (TOMS)* 43 (2016), Nr. 3, S. 1–27

[46] ROMAN, S. : The formula of Faa di Bruno. In: *The American Mathematical Monthly* 87 (1980), Nr. 10, S. 805–809

[47] RÜDE, U. ; WILLCOX, K. ; MCINNES, L. C. ; STERCK, H. D.: Research and education in computational science and engineering. In: *Siam Review* 60 (2018), Nr. 3, S. 707–754

[48] SAILFISH: *https://github.com/sailfish-team/sailfish*. – accessed on 2020-01-21

[49] SCHILLER, U. D.: Thermal fluctuations and boundary conditions in the lattice Boltzmann method. (2008), Nr. November. http://www2.mpip-mainz.mpg.de/theory/Theses/Thesis{_}Documents/50/thesis

[50] SCHMIESCHEK, S. ; SHAMARDIN, L. ; FRIJTERS, S. ; KRÜGER, T. ; SCHILLER, U. D. ; HARTING, J. ; COVENEY, P. V.: LB3D: A parallel implementation of the Lattice-Boltzmann method for simulation of interacting amphiphilic fluids. In: *Computer Physics Communications*

217 (2017), S. 149–161. http://dx.doi.org/10.1016/j.cpc.2017.03.013. –
DOI 10.1016/j.cpc.2017.03.013

[51] SCHORNBAUM, F. ; RÜDE, U. : Massively parallel algorithms for the
lattice Boltzmann method on nonuniform grids. In: *SIAM Journal on
Scientific Computing* 38 (2016), Nr. 2, S. C96–C126

[52] SCHORNBAUM, F. ; RÜDE, U. : Extreme-scale block-structured adaptive
mesh refinement. In: *SIAM Journal on Scientific Computing* 40 (2018),
Nr. 3, S. C358–C387

[53] SILVA, G. ; SEMIAO, V. : First-and second-order forcing expansions
in a lattice Boltzmann method reproducing isothermal hydrodynamics
in artificial compressibility form. In: *Journal of fluid mechanics* 698
(2012), S. 282–303

[54] TANG, Y. ; CHOWDHURY, R. A. ; KUSZMAUL, B. C. ; LUK, C.-K.
; LEISERSON, C. E.: The pochoir stencil compiler. In: *Proc. of the
twenty-third annual ACM symposium on Parallelism in algorithms and
architectures* ACM, 2011, S. 117–128

[55] TCLB: *https://github.com/CFD-GO/TCLB*. – accessed on 2020-01-21

[56] TREIBIG, J. ; HAGER, G. ; WELLEIN, G. : Likwid: A lightweight
performance-oriented tool suite for x86 multicore environments. In:
*2010 39th International Conference on Parallel Processing Workshops*
IEEE, 2010, S. 207–216

[57] WELLEIN, G. ; LAMMERS, P. ; HAGER, G. ; DONATH, S. ; ZEISER, T.
: Towards optimal performance for lattice Boltzmann applications on
terascale computers. In: *Parallel Computational Fluid Dynamics 2005*.
Elsevier, 2006, S. 31–40

[58] WELLEIN, G. ; ZEISER, T. ; HAGER, G. ; DONATH, S. : On the single
processor performance of simple lattice Boltzmann kernels. In: *Computers
and Fluids* 35 (2006), Nr. 8-9, S. 910–919. http://dx.doi.org/10.1016/j.
compfluid.2005.02.008. – DOI 10.1016/j.compfluid.2005.02.008. – ISBN
0045–7930

[59] WILKE, J. ; POHL, T. ; KOWARSCHIK, M. ; RÜDE, U. : Cache
performance optimizations for parallel lattice Boltzmann codes. In:
*European Conference on Parallel Processing* Springer, 2003, S. 441–450

[60] WITTMANN, M. ; HAAG, V. ; ZEISER, T. ; KÖSTLER, H. ; WELLEIN,
G. : Lattice Boltzmann benchmark kernels as a testbed for performance
analysis. In: *Computers & Fluids* 172 (2018), S. 582–592

[61] WITTMANN, M. ; ZEISER, T. ; HAGER, G. ; WELLEIN, G. : Com-
parison of different propagation steps for lattice Boltzmann methods.
In: *Computers and Mathematics with Applications* 65 (2013), Nr.
6, 924–935. http://dx.doi.org/10.1016/j.camwa.2012.05.002. – DOI
10.1016/j.camwa.2012.05.002. – ISBN 0898–1221