

Advanced Automatic Code Generation For Multiple Relaxation-Time Lattice Boltzmann Methods

Frederick Hennig, Markus Holzer, and Ulrich Rüde

Presented by Russell Bentley

Stony Brook

2024



Stony Brook University

What are we even talking about?

ADVANCED AUTOMATIC CODE GENERATION FOR MULTIPLE RELAXATION-TIME LATTICE BOLTZMANN METHODS

- ▶ Present my graphics final project
 - ▶ Computational Fluid Dynamics (CFD)
 - ▶ MRT Lattice Boltzmann Methods (LBM)
 - ▶ Theory & Implementation
 - ▶ We'll understand the title (hopefully)
- ▶ Present paper
 - ▶ waLBerla, pystencils, and lbmpy
 - ▶ Implementation & Symbolic Simplification,
- ▶ Future Work



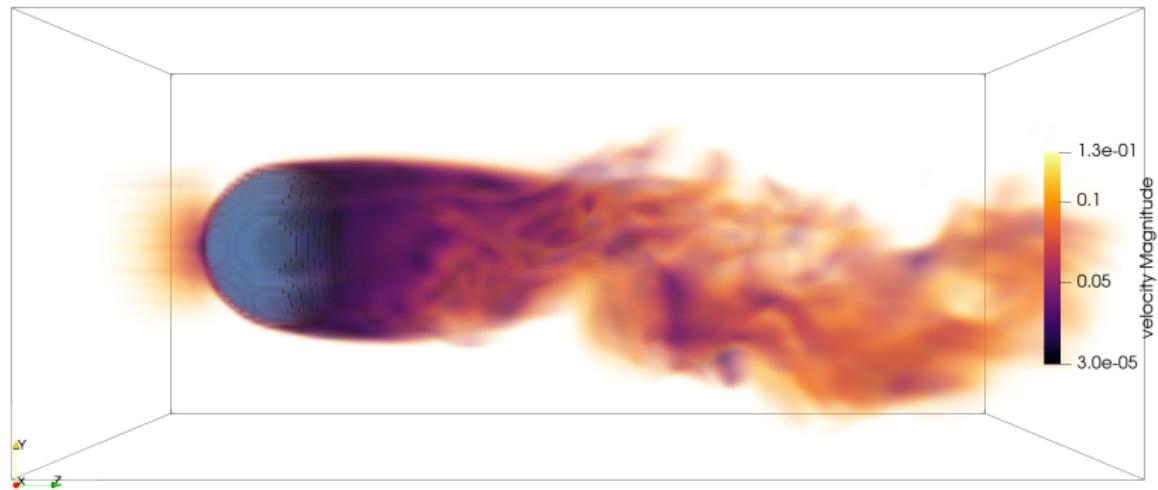
Stony Brook University

Graphics Final Introduction

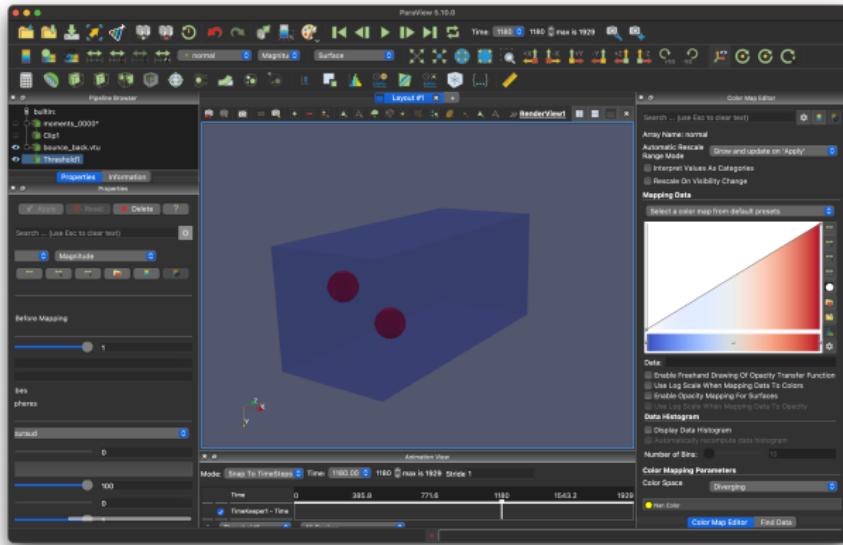
- ▶ Final Project for Computer Graphic
 - ▶ Rough Goal to implement [11, 12]
- ▶ Lattice Boltzmann method (LBM)
 - ▶ Numerically model fluid dynamics
 - ▶ Known to struggle with turbulence
- ▶ CM-MRT: A new collision operator
 - ▶ Based on non-orthogonal moment space [5]
 - ▶ Adoption in graphics research [9–12]



Demo

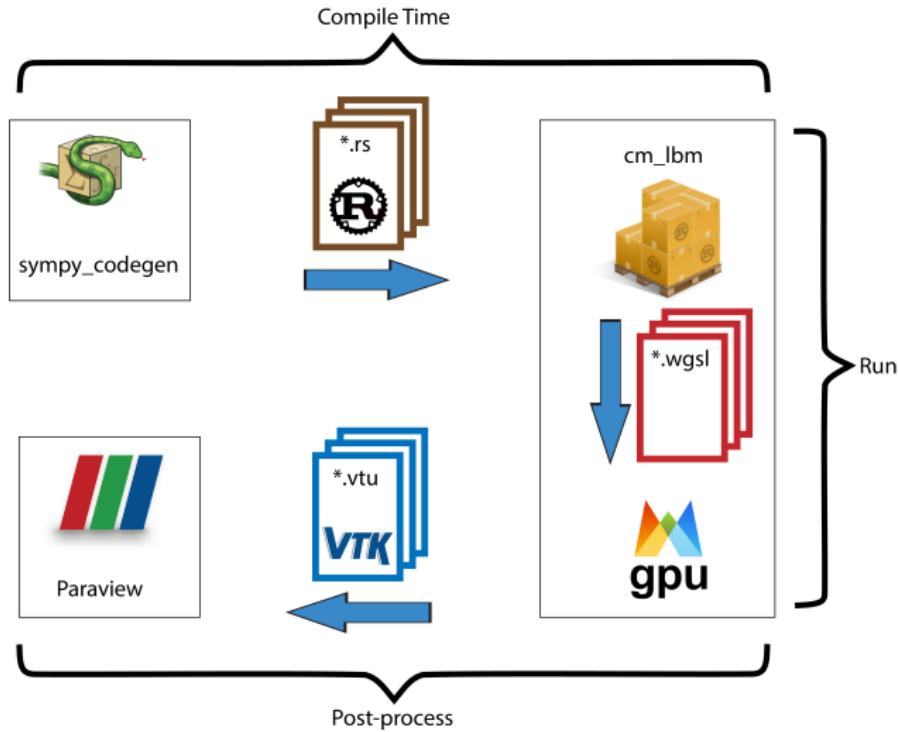


Paraview + VTK



<https://www.paraview.org>
<https://vtk.org>

Architecture



Stony Brook University

Classic CFD

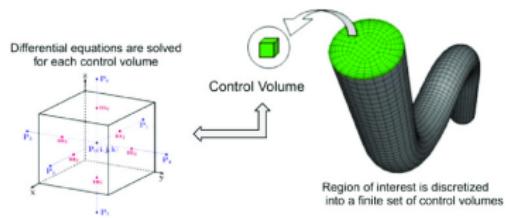
- ▶ Discretize space into control Volumes
- ▶ Solve for macroscopic fluid values
- ▶ Common to require global solve for pressure term

Conservation of Mass

$$\frac{\partial \rho}{\partial t} + \nabla \cdot (\rho \mathbf{u}) = 0$$

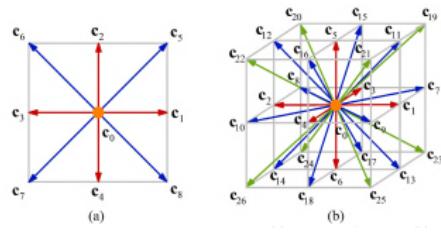
Conservation of Momentum

$$\frac{\partial}{\partial t}(\rho \mathbf{u}) + \nabla \cdot (\rho \mathbf{u} \times \mathbf{u}) = -\nabla p + \nabla \cdot \sigma$$



Lattice Boltzmann method (LBM): Discretization

- ▶ Model microscopic particles
 - ▶ Recover macroscopic properties
 - ▶ D dimensional grid of points x
 - ▶ Particle position distribution
 - ▶ Q dimensional “lattice” at each point
 - ▶ Particle velocity distribution
 $f(x, \xi, t)$
 - ▶ Directions c_i streaming and collision



(a) D2Q9 (b) D3Q27

$$c_i = \begin{cases} (0, 0, 0) & i = 0 \\ (\pm 1, 0, 0), (0, \pm 1, 0), (0, 0, \pm) & i = 1, \dots, 6 \\ (\pm 1, \pm 1, 0), (\pm 1, 0, \pm 1), (0, \pm 1, \pm 1) & i = 7, \dots, 18 \\ (\pm 1, \pm 1, \pm 1) & i = 19, \dots, 26 \end{cases}$$

LBM Theory Primer

- ▶ Lattice Boltzmann equation
- ▶ $f(x, t)$ approximates ξ
- ▶ Hermite Polynomial Basis [5]
 - ▶ Quadrature Rule with c_i .
- ▶ Collision operator, Ω
 - ▶ Relax towards equilibrium f_{eq}
- ▶ Moments of f are macroscopic quantities
- ▶ Advance time with two operators, *streaming* and *collision*

$$\partial_t f(x, \xi, t) + \xi \cdot \nabla f(x, \xi, t) = \Omega(f)$$

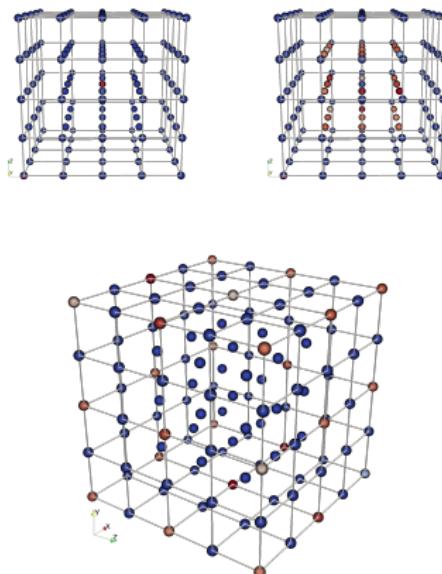
$$\rho(x, t) = \sum_{i=0}^{Q-1} f_i(x, t)$$

$$\rho(x, t)u(x, t) = \sum_{i=0}^{Q-1} c_i f_i(x, t)$$

Streaming

- ▶ Non-local operator (Transport)
- ▶ Particles stream to next point in grid
- ▶ Information travels at speed of sound!

$$f_i^*(x, t) = f_i(x - c_i, t)$$



Stony Brook University

Collision - Single Relaxation Time (SRT)

- ▶ Collision is local
- ▶ Relax distributions towards equilibrium
- ▶ SRT is standard collision operator
 - ▶ Relaxes all distributions at constant rate
 - ▶ Also known as Bhatnagar-Gross-Krook (BGK)
- ▶ First order approximation
- ▶ Equilibrium from continuous Maxwell-Boltzmann distribution

$$\text{SRT}(f) = \frac{1}{\tau} (f_{\text{eq}} - f)$$

$$f_{\text{eq}} = w_i \rho \left(1 + 3c_i \cdot u + \frac{9(c_i \cdot u)^2}{2} - \frac{3(u \cdot u)^2}{2} \right)$$



Stony Brook University

Multiple Relaxation Time (MRT)

- ▶ M transforms f into moment space
- ▶ 26 of them!
- ▶ Relax moments separately
- ▶ M^{-1} transforms result back to distribution

$$m_{0,j} = \{1\},$$

$$m_{1,j} = \{c_{j,x}\},$$

$$m_{18,j} = \{c_{j,x}^2 c_{j,y}^2 + c_{j,x}^2 c_{j,z}^2 - c_{j,y}^2 c_{j,z}^2\},$$

$$m_{19,j} = \{c_{j,x}^2 c_{j,y}^2 - c_{j,x}^2 c_{j,z}^2\},$$

$$m_{26,j} = \{c_{j,x}^2 c_{j,y}^2 c_{j,z}^2\},$$

$$\text{MRT}(f) = M^{-1} \cdot R \cdot M(f_{\text{eq}} - f)$$



Stony Brook University

Relaxation Rates, pt 1

- ▶ Independent relaxation rates for each moment
- ▶ Zeroth moment, density, should be conserved.
- ▶ Next three are momentum, $r_i = 2$ is common
- ▶ Next six moments are related to momentum flux
 - ▶ Relaxation rate is related to viscosity, ν

$$R = \begin{bmatrix} r_0 & 0 & \cdots & 0 \\ 0 & r_1 & \cdots & 0 \\ \vdots & & & \vdots \\ 0 & & \cdots & r_{26} \end{bmatrix}$$

$$r_0 = 0,$$

$$r_i = 2, \text{ for } i = 1, 2, 3$$

$$r_i = (3\nu + 1/2)^{-1}, \text{ for } i = 4, \dots, 9$$



Relaxation Rates, pt 2

- ▶ Higher-order relaxation rates are an area of active research
- ▶ [11] Describes an adaptive approach
- ▶ I implemented fixed rates as described in [9].

$$R = \begin{bmatrix} r_0 & 0 & \cdots & 0 \\ 0 & r_1 & \cdots & 0 \\ \vdots & & & \vdots \\ 0 & \cdots & & r_{26} \end{bmatrix}$$

$$r_i = (3v_i + 1/2)^{-1}$$

$$v_i = 0.005, \text{ for } i = 9, \dots, 16$$

$$v_i = 0.007, \text{ for } i = 17, \dots, 22$$

$$v_i = 0.009, \text{ for } i = 26$$



Stony Brook University

Central Moment MRT (CM-MRT)

- ▶ MRT violates Galilean invariance
 - ▶ u changes collision outcome
- ▶ Center Moments on u [4]
- ▶ $\bar{m}_\rho = M_u f_{eq}(u, \rho)$
 - ▶ Equilibrium in CM space
 - ▶ Only depends on ρ !

$$\bar{c}_j = c_j - u$$

$$m_{0,j} = \{1\},$$

$$m_{1,j} = \{\bar{c}_{j,x}\},$$

$$m_{18,j} = \{\bar{c}_{j,x}^2 \bar{c}_{j,y}^2 + \bar{c}_{j,x}^2 \bar{c}_{j,z}^2 - \bar{c}_{j,y}^2 \bar{c}_{j,z}^2\},$$

$$m_{19,j} = \{\bar{c}_{j,x}^2 \bar{c}_{j,y}^2 - \bar{c}_{j,x}^2 \bar{c}_{j,z}^2\},$$

$$m_{26,j} = \{\bar{c}_{j,x}^2 \bar{c}_{j,y}^2 \bar{c}_{j,z}^2\},$$

$$\text{CM-MRT}_{u,\rho}(f) = M_u^{-1} \cdot R \cdot (\bar{m}_\rho - M_u f)$$

Central Moment MRT (CM-MRT)

- ▶ MRT violates Galilean invariance
- ▶ Center Moments on u [4]
- ▶ f_{eq} only depends on ρ in CM space

$$\bar{c}_j = c_j - u$$

$$m_{0,j} = \{1\},$$

$$m_{1,j} = \{\bar{c}_{j,x}\},$$

$$m_{18,j} = \{\bar{c}_{j,x}^2 \bar{c}_{j,y}^2 + \bar{c}_{j,x}^2 \bar{c}_{j,z}^2 - \bar{c}_{j,y}^2 \bar{c}_{j,z}^2\},$$

$$m_{19,j} = \{\bar{c}_{j,x}^2 \bar{c}_{j,y}^2 - \bar{c}_{j,x}^2 \bar{c}_{j,z}^2\},$$

$$m_{26,j} = \{\bar{c}_{j,x}^2 \bar{c}_{j,y}^2 \bar{c}_{j,z}^2\},$$

$$\text{CM-MRT}_{u,\rho}(f) = M_u^{-1} \cdot R \cdot (\bar{m}_\rho - M_u f)$$

Code Example: SymPy

$$f_i^* = f_i + \Omega_i(f, \rho, u).$$

```
u = cm_mrt.u()
f = cm_mrt.f()
density = Symbol("density")
omega = Symbol("omega")
eq_op = cm_mrt.f_eq(density, u)
bgk_op = f + omega * (eq_op - f)

(source_body, debug_raw) = util.rust_generate_op(bgk_op.evalf())
```

```
u = cm_mrt.u()
f = cm_mrt.f()
density = Symbol("density")
riv = Symbol("riv")
m = cm_mrt.M(u)
m1 = m.inv()
r = cm_mrt.R(riv)
f = cm_mrt.f()
mbar = cm_mrt.MBar(density)
cm_mrt_op = f + m1 * r * (mbar - (m * f))

(source_body, debug_raw) = util.rust_generate_op(cm_mrt_op.evalf())
```

- ▶ Using SymPy was inspired by [7]
- ▶ Looks like a compiler



Stony Brook University

Code Example: Shader

```
result[0] = -0.0698812019566736*density + 1.88679245283019*q0*ux * ux*uy*uz*(1.0 - ux) * (1.0 - ux) + 1.88679245283019*q10*uz * uz*(-.0) + 1.88679245283019*q11*uy * uy*(1.0 - ux) * (1.0 - ux)*(1.0 - uz) * (1.0 - ux)*(-uz - 1.0) * (-uz - 1.0) + 1.88679245283019*q12*(-ux - 1.0) + 1.88679245283019*q14*uy * uy*(-ux - 1.0) * (-ux - 1.0), 15*ux * ux*(1.0 - uy) * (1.0 - uy)*(1.0 - uz) * (1.0 - uz) + 1.88679245283019*q15*uz * (-uz - 1.0) + 1.88679245283019*q17*ux * ux*(1.0 - uz) * (1.0 - ux)*q18*ux * ux*(-uy - 1.0) * (-uy - 1.0)*(-uz - 1.0) * (-uz - 1.0) + 1.88679245283019*q19*x*(1.0 - uy) * (1.0 - uy)*(1.0 - uz) * (1.0 - uz) + 1.88679245283019*q20*(1.0 - ux) * (1.0 - ux)*(1.0 - uy) * (1.0 - uy)*21*(1.0 - ux) * (1.0 - ux)*(1.0 - uz) * (1.0 - uz)*(-uy - 1.0) * (-uy - 1.0)*(-uz - 1.0) * (-uz - 1.0) + 1.88679245283019*q24*(1.0 - uy - 1.0) * (-uy - 1.0) + 1.88679245283019*q25*(1.0 - uy) * (1.0 - uy) * (1.0 - 1.0) + 1.88679245283019*q26*(-ux - 1.0) * (-ux - 1.0)*(-uy - 1.0) * 9245283019*q3*ux * ux*uz * uz*(1.0 - uy) * (1.0 - uy) + 1.88679245283019*q5*ux * ux*uy * uy*(1.0 - uz) * (1.0 - uz) + 1.88679245283019*q7*uz * uz*(1.0 - ux) * (1.0 - ux)*(1.0 - ux)
```

- ▶ $\approx 5\text{kb}$ for moments shader
- ▶ 1.1mb for CM-MRT shader
- ▶ Cheap to compute
- ▶ Challenging to debug



Code Example: Rust

```
let q25 = f[25]; => f32
let q26 = f[26]; => f32
result[0] = -0.0698812019566736 * density
    + 1.88679245283019 * q0 * ux * ux * uy * uy * uz * uz
    + q0
    + 1.88679245283019 * q1 * uy * uy * uz * uz * (1.0 - ux) * (1.0 - ux)
    + 1.88679245283019
        * q10
        * uz
        * uz
        * (-ux - 1.0)
        * (-ux - 1.0)
        * (-uy - 1.0)
        * (-uy - 1.0)
    + 1.88679245283019
        * q11
        * uy
        * uy
        * (1.0 - ux)
        * (1.0 - ux)
        * (1.0 - uz)
        * (1.0 - uz)
```

- ▶ With formatting, $\approx 45k$ lines
- ▶ Easier to test and debug
- ▶ Tricky parts match shader exactly



Graphics Final Conclusion

- ▶ How are we feeling about the paper title?
 - ▶ MRT Lattice Boltzmann methods?
 - ▶ Code generation using `sympy`?
- ▶ LBM + `sympy`
 - ▶ “A Literate Lattice Boltzmann Code” [8]

**ADVANCED AUTOMATIC CODE GENERATION FOR MULTIPLE
RELAXATION-TIME LATTICE BOLTZMANN METHODS**



Stony Brook University

Paper Context

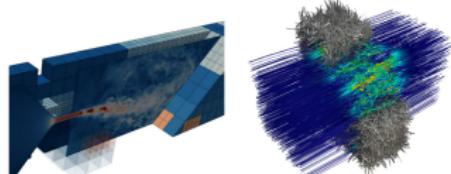
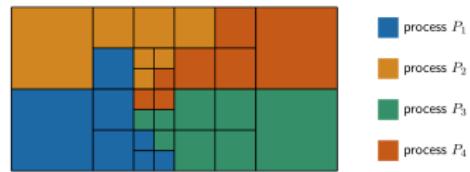
- ▶ Today's paper is part of a series
- ▶ Collaborative HPC work at German University
 - ▶ Friedrich-Alexander-Universität
- ▶ Also part of MultiXScale
 - ▶ European joint effort on exascale simulation



Stony Brook University

waLBerla

- ▶ waLBerla framework
 - ▶ Introduced in [3] (2021)
 - ▶ HPC-scale simulations
 - ▶ Stencil problems specifically
 - ▶ Bring your Compute Kernels
- ▶ Provides:
 - ▶ Multi-node domain partitioning
 - ▶ Communication patterns
 - ▶ Checkpoints

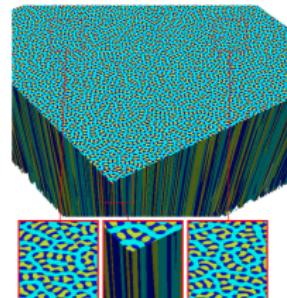
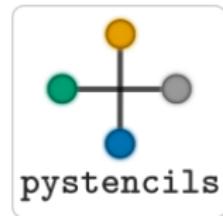


(a) Adaptive grid refinement for turbulent flow inside a vocal fold [40].

(b) Fluid flow through a porous crystal geometry [75].

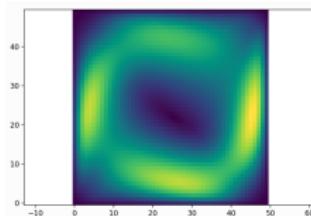
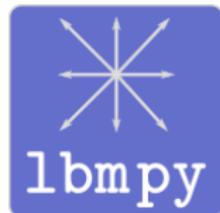
pystencils

- ▶ Introduced in [2] (2019)
- ▶ Stencil Compiler
 - ▶ Generate C or CUDA source
 - ▶ sympy based IR
 - ▶ Assignments + Instructions
 - ▶ Static Single Assignment (SSA)
- ▶ Kernels can be
 - ▶ Compiled and used with numpy
 - ▶ Integrated with other code (e.g. waLBerla)



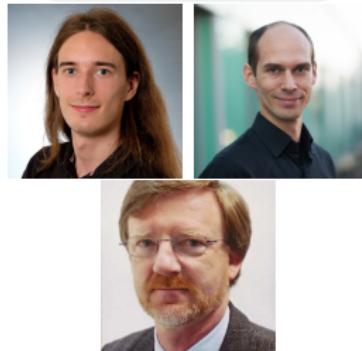
Stony Brook University

- ▶ Compiler for LBM Kernels
 - ▶ Introduced in [1] (2021)
 - ▶ Built on pystencils
- ▶ Supports
 - ▶ Many grid configurations ($D3Q27$)
 - ▶ Arbitrary governing equations (f_{eq})
- ▶ Originally only supported SRT, MRT
- ▶ Deployable to exa-scale machines with waLBerla



Introduction

- ▶ Presenting [7] (2023)
- ▶ Re-architecture of lbmpy
 - ▶ Zero-centered Storage
 - ▶ Central Moment (CM)
 - ▶ Cumulants (K)
 - ▶ Support more generic governing eqns
- ▶ Minimizes arithmetic operations for CM and K methods
- ▶ Demonstrate small relative cost compared to SRT



Zero-centered Storage

- ▶ For incompressible hydrodynamics:
- ▶ f is generally close to f^0 .
- ▶ So lets only store $\partial f = f - f^0$
- ▶ Improve numerical accuracy
- ▶ Can also use smaller float format!
- ▶ Adds code-gen complexity

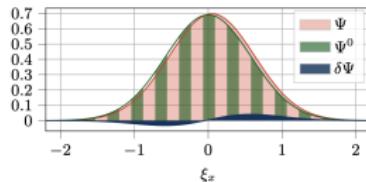


FIGURE 2.1. Plot of Ψ , Ψ^0 and $\delta\Psi$ in one dimension at $\rho = 1.01$, $u_x = 0.05$.



Stony Brook University

Moment Transform, pt. 1

- ▶ Transform populations to Raw Moments (Mf)
- ▶ Use Chimera transform to compute moments
 - ▶ Introduced in [6]
- ▶ More efficient than Matrix transform
- ▶ Inverse uses trick to split populations

$$f_i^* = f_i^+ + f_i^-, \quad f_{\bar{i}}^* = f_i^+ - f_i^-$$

$$m_{xy|\gamma} := \sum_{z \in \{-1,0,1\}} f_{xyz} \cdot z^\gamma$$

$$m_{x|\beta\gamma} := \sum_{y \in \{-1,0,1\}} m_{xy|\gamma} \cdot y^\beta$$

$$m_{\alpha\beta\gamma} := \sum_{x \in \{-1,0,1\}} m_{x|\beta\gamma} \cdot x^\alpha$$



Stony Brook University

Moment Transform, pt. 2

- ▶ To compute central moments
- ▶ First compute raw moments $m_{\alpha\beta\gamma}$
- ▶ Then transform them to central moments $\kappa_{\alpha\beta\gamma}$

$$\kappa_{ab|\gamma} := \sum_{c=0}^{\gamma} \binom{\gamma}{c} (-u_z)^{\gamma-c} m_{abc}$$

$$\kappa_{a|\beta\gamma} := \sum_{c=0}^{\gamma} \binom{\gamma}{c} (-u_z)^{\gamma-c} m_{abc}$$

$$\kappa_{\alpha\beta\gamma} = \sum_{a,b,c=0}^{\alpha\beta\gamma} \binom{\alpha}{a} \binom{\beta}{b} \binom{\gamma}{c} (-u_x)^{\alpha-a} (-u_y)^{\beta-b} (-u_z)^{\gamma-c} m_{abc}$$



Stony Brook University

Simplification Passes

- ▶ Conserved Quantity Rewriting
 - ▶ Use moments instead of ρ, \mathbf{u}
- ▶ Propagation of Logarithms
 - ▶ Relevant for Cumulants
- ▶ Common Subexpression Elimination (CSE)
 - ▶ Make new assignments for them
- ▶ Expression Propagation
 - ▶ For assignments with constant RHS
 - ▶ Copy into use-sites
- ▶ Unused Subexpression Elimination
 - ▶ Remove unused assignments

$$\kappa_{000} = \rho$$

$$\kappa_{100} = -F_x/2$$

$$\kappa_{000}^* = \kappa_{000}$$

$$\kappa_{100}^* = \kappa_{100} + F_x$$

$$m_{10|0}^* = \kappa_{000}^* \mathbf{u}_x + \kappa_{100}^*$$

To

$$m_{10|0}^* = \rho \mathbf{u}_x + F_x/2$$



Stony Brook University

Symbolic Simplification Results

TABLE 4.1

Arithmetic operation counts in the symbolic representation of compute kernels generated by lbmpy for several method definitions on the D2Q9, D3Q19 and D3Q27 stencils. Numbers for kernels derived without simplification (N), standard simplification (S) and simplification plus CSE (S+CSE) are listed.

Method	D2Q9			D3Q19			D3Q27		
	N	S	S+CSE	N	S	S+CSE	N	S	S+CSE
SRT	432	113	91	1156	312	204	3842	428	285
TRT	437	191	111	1161	480	239	3847	668	343
O-MRT	196	142	122	554	397	324	928	596	488
R-O-MRT	196	105	89	554	290	232	928	407	329
WO-MRT	178	117	102	507	339	283	843	484	415
R-WO-MRT	178	87	75	507	244	196	843	316	265
CM	236	156	132	638	415	344	1140	747	600
R-CM	236	108	95	638	255	216	1140	417	342
K	676	167	142	1854	454	376	7623	1035	820
R-K	676	114	100	1854	272	231	7623	489	397

For reference, my had CM implementation (slight over estimate) 220,702 arithmetic operations, 676 for SRT



Stony Brook University

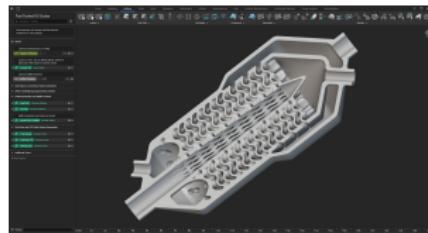
Paper Conclusion

- ▶ lbmpy now has SOTA collision spaces
 - ▶ Far more efficient than other implementations
- ▶ High level symbolic representation
 - ▶ Easy to experiment with variant methods
 - ▶ Easy to target alternate architectures
- ▶ Easy to scale with waLBerla



Future Work (For Me)

- ▶ Tracer Particles
- ▶ Boundary Conditions
 - ▶ Generate code for [12]?
 - ▶ Implicit Geometry Support?
- ▶ Problem Setup
 - ▶ Unit conversions,
 - ▶ Adaptive relaxation from [11]
- ▶ Alternative Compute Targets (Metal)



Stony Brook University

Conclusion

- ▶ Today we covered:
 - ▶ Lattice Boltzmann Methods
 - ▶ Central Moments (CM) collision Space
 - ▶ Code generation with sympy
- ▶ Questions?



Stony Brook University

References I

- [1] Martin Bauer, Harald Köstler, and Ulrich Rüde. “Ibmpy: Automatic code generation for efficient parallel lattice Boltzmann methods”. In: *Journal of Computational Science* 49 (2021), p. 101269.
- [2] Martin Bauer et al. “Code generation for massively parallel phase-field simulations”. In: *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis*. 2019, pp. 1–32.
- [3] Martin Bauer et al. “waLBerla: A block-structured high-performance framework for multiphysics simulations”. In: *Computers & Mathematics with Applications* 81 (2021), pp. 478–501.
- [4] Alessandro De Rosis. “Nonorthogonal central-moments-based lattice Boltzmann scheme in three dimensions”. In: *Physical Review E* 95.1 (2017), p. 013310.



References II

- [5] Alessandro De Rosis and Kai H Luo. "Role of higher-order Hermite polynomials in the central-moments-based lattice Boltzmann framework". In: *Physical Review E* 99.1 (2019), p. 013301.
- [6] Martin Geier et al. "The cumulant lattice Boltzmann equation in three dimensions: Theory and validation". In: *Computers & Mathematics with Applications* 70.4 (2015), pp. 507–547.
- [7] Frederik Hennig, Markus Holzer, and Ulrich Rüde. "Advanced automatic code generation for multiple relaxation-time lattice Boltzmann methods". In: *SIAM Journal on Scientific Computing* 45.4 (2023), pp. C233–C254.
- [8] Adrian Kummerländer. *A literate Lattice Boltzmann Code*. <http://literatelb.org>. 2020.



References III

- [9] Wei Li, Kai Bai, and Xiaopei Liu. "Continuous-scale kinetic fluid simulation". In: *IEEE transactions on visualization and computer graphics* 25.9 (2018), pp. 2694–2709.
- [10] Wei Li, Kui Wu, and Mathieu Desbrun. "Kinetic Simulation of Turbulent Multifluid Flows". In: *ACM Transactions on Graphics (TOG)* 43.4 (2024), pp. 1–17.
- [11] Wei Li et al. "Fast and scalable turbulent flow simulation with two-way coupling.". In: *ACM Trans. Graph.* 39.4 (2020), p. 47.
- [12] Chaoyang Lyu et al. "Fast and versatile fluid-solid coupling for turbulent flow simulation". In: *ACM Transactions on Graphics* 40.6 (2021), p. 201.

