# Hardware-Accelerated Phase-Averaging for Cavitating Bubbly Flows

Diego Vaca-Revelo[a], Benjamin Wilfong[b], Spencer H. Bryngleson[b,c,d], Aswin Gnanaskandan[a]

[a]*Mechanical and Materials Engineering, Worcester Polytechnic Institute, Worcester, MA 01609, USA*

[b]*School of Computational Science & Engineering, Georgia Institute of Technology, Atlanta, GA 30332, USA*

[c]*Daniel Guggenheim School of Aerospace Engineering, Georgia Institute of Technology, Atlanta, GA 30332, USA*

[d]*George W. Woodruff School of Mechanical Engineering, Georgia Institute of Technology, Atlanta, GA 30332, USA*

## Abstract

We present a comprehensive validation, performance characterization, and scalability analysis of a hardware-accelerated phase-averaged multiscale solver designed to simulate acoustically driven dilute bubbly suspensions. The carrier fluid is modeled using the compressible Navier–Stokes equations. The dispersed phase is represented through two distinct subgrid formulations: a volume-averaged model that explicitly treats discrete bubbles within a Lagrangian framework, and an ensemble-averaged model that statistically represents the bubble population through a discretized distribution of bubble sizes. For both models, the bubble dynamics are modeled via the Keller–Miksis equation. For the GPU cases, we use OpenACC directives to offload computation to the GPUs. The volume-averaged model is validated against the analytical Keller-Miksis solution and experimental measurements, showing excellent agreement with root-mean-squared errors of less than $8\,\%$ for both single-bubble oscillation and collapse scenarios. The ensemble-averaged model is validated by comparing it to volume-averaged simulations. On an NCSA Delta node with 4 NVIDIA A100 GPUs, we observe a speedup 16-fold compared to a 64-core AMD Milan CPU. The ensemble-averaged model offers additional reductions in computational cost by solving a single set of averaged equations, rather than multiple stochastic realizations. However, the volume-averaged model enables the interrogation of individual bubble dynamics, rather than the averaged statistics of the bubble dynamics. Weak and strong scaling tests demonstrate good scalability across both CPU and GPU platforms. These results show the proposed method is robust, accurate, and efficient for the multiscale simulation of acoustically driven dilute bubbly flows.

*Keywords:* Cavitation, Bubble dynamics, Subgrid bubble models, Hardware acceleration

## 1. Introduction

Acoustic cavitation refers to the process in which sound waves induce rapid pressure fluctuations within a liquid, causing the formation, growth, shrinkage, and eventual collapse of small gas-filled bubbles [1, 2]. These bubbles oscillate in response to the alternating phases of compression and rarefaction in the acoustic field. Under suitable conditions, the oscillations can become unstable,

Code available at https://github.com/MFlowCode/MFC.

leading to inertial collapse that concentrates energy into a small region. This collapse can generate extreme local temperatures, high-pressure shock waves, and high-velocity microjets [3–5]. Acoustic cavitation has been widely applied in fields such as biomedicine [6–8] and engineering [9–11], where the unique properties of bubble dynamics are exploited for various purposes.

While acoustic cavitation can occur in a wide range of bubbly environments, a particularly important case is that of dilute bubbly suspensions. In such systems, small gas bubbles are dispersed within a liquid, with a low volume fraction (typically less than 0.01) of bubbles. The bubble distribution, bubble size, acoustic bubble-bubble interaction, and the surrounding liquid typically influence the behavior of these suspensions under the influence of an acoustic wave. Even at low concentrations, these small bubbles strongly influence sound propagation, introducing dispersion and attenuation which becomes prominent when the bubbles oscillate near their resonance frequencies [12, 13]. These characteristics make dilute bubbly suspensions a versatile medium for both theoretical and experimental validation of acoustic models. Applications of dilute bubbly suspensions span multiple disciplines. In medical ultrasound, microbubble suspensions are used as contrast agents to enhance diagnostic imaging and as carriers for targeted drug delivery, where controlled cavitation can release therapeutic compounds at precise locations [7, 8]. It is also central to shock wave lithotripsy, a non-invasive medical procedure that uses acoustic cavitation-induced shock waves to break down kidney stones [6]. Furthermore, microbubble-enhanced high-intensity focused ultrasound therapy leverages acoustic cavitation to enhance the therapeutic effects of ultrasound in medical treatments, such as tumor ablation [14–17]. It has also been utilized in materials processing, including in the enhancement of chemical reactions through sonochemistry [9, 10]. Industrial processes make use of their cleaning capabilities, emulsification potential, and effectiveness in food processing, while environmental engineering applications include wastewater treatment, where enhanced aeration and oxidation accelerate contaminant breakdown [18–20]. Therefore, elucidating a better understanding of the dynamics of gas bubbles dispersed in a liquid medium under the influence of acoustic waves is of vital importance.

The length scales in dilute bubbly flows span centimeter-scale acoustic wave propagation and micrometer-scale bubbles, and fully resolving both would require an impractically fine mesh across the entire domain. Subgrid bubble models are therefore essential for representing such multiscale bubbly systems efficiently. In this framework, the acoustic field is resolved on a relatively coarse Eulerian grid, while the unresolved microbubbles are modeled at the subgrid scale. Two phase-averaged subgrid-scale approaches are commonly employed: ensemble averaging (Euler–Euler) and volume averaging (Euler–Lagrange). In the ensemble-averaged approach, the aim is to capture the mean collective response of a bubble population characterized by a known probability distribution of radii. This method assumes the presence of many stochastically distributed bubbles within each computational cell and evaluates the statistically averaged mixture dynamics. Individual bubbles are not explicitly resolved, though this approach retains the essential acoustic and dynamical effects of the unresolved bubbles, offering good computational efficiency [21, 22]. Volume-averaging employs a Lagrangian strategy, where individual bubbles are treated as discrete entities [23, 24], and their volumetric oscillations are typically modeled using one of the standard bubble dynamics equations. While this approach retains the individual size dynamics of each bubble, it is usually expensive to obtain the mean behavior of a bubble population due to the need for several independent ensembles of bubble distributions within a cloud. The computationally advantageous approach depends on the number of bubbles within the bubbly suspension and the available computational resources. Typically, the ensemble-averaged models are preferable for a large number of bubbles, and volume-averaged models are preferred when representing individual bubble dynamics is important [25].

Although these models are well-established, improving their computational performance remains a priority for enabling large-scale parametric studies that can close knowledge gaps in cavitation physics. A common acceleration strategy is distributed memory parallelization via MPI [26]. With a purely grid-based decomposition, each CPU core performs separate calculations, and the MPI protocol exchanges data with its neighboring cores. This strategy is effective when the computational workload is nearly evenly distributed, such as when solving the compressible Navier–Stokes equations in a finite volume framework, where each subdomain typically requires a similar amount of computations if the number of cells is distributed evenly. Simulations involving subgrid bubbles, on the other hand, can introduce load balancing challenges [27]. Load imbalance occurs when bubbles are localized in specific regions of the domain, resulting in an uneven computational workload across processors. This issue may appear in both the phase-averaged subgrid bubble models, being more prejudicial for the volume-averaged model, since for the discrete phase, the number of equations being solved is proportional to the number of bubbles present. Subdomains with high bubble concentrations require significantly more processing time, while others with fewer or no bubbles finish their computations sooner. The underutilized processors must then wait for the heavily loaded ones to complete before MPI synchronization can proceed. This unproductive time reduces parallel efficiency and slows the overall simulation, introducing an opportunity to improve existing solvers. The present work addresses this by optimizing solver performance and mitigating the effects of load imbalance.

Another well-established approach for accelerating CFD solvers is the use of offloading strategies, particularly to modern graphics processing units (GPUs) [28–31]. GPUs can offer substantial performance gains; for example, Sweet et al. [30] reported speedups of up to 14x in simulations of particle-laden turbulent flows compared to CPU-only implementations. Piscaglia and Ghioldi [32] used OpenFOAM and offloaded computationally intensive calculations to GPUs and reported a 10-fold speedup. The acceleration reported in the work of Jespersen [33] goes from 2.5x to 3.0x when comparing a GPU against a single CPU. Nevertheless, GPU acceleration may be constrained by limited memory and increased communication overhead. To overcome these issues, Radhakrishnan et al. [31] developed a directive-based offloading strategy for multiphase compressible flow solvers, achieving high memory reuse and efficient computation. This strategy was implemented in the open-source CFD solver Multi-Component Flow Code (MFC) [34, 35], where a speedup of 40x was reported on a single node using NVIDIA V100 GPUs over IBM POWER9 CPUs. Building on this foundation, the present work introduces a GPU-based hardware acceleration strategy for phase-averaged subgrid models to enhance computational efficiency.

This work explores the computational speedup achievable with the state-of-the-art subgrid bubble models on CPU and GPU architectures, identifies the conditions under which each phase-averaged formulation is most computationally efficient, and evaluates how effectively modern GPU hardware can mitigate load imbalance arising from localized bubble clusters. To address these questions, we develop and assess a GPU-accelerated framework for both volume- and ensemble-averaged subgrid bubble models, quantify the speedup relative to CPU-based execution, and examine the performance of each approach across a range of bubble populations. The analysis identifies the scenarios in which each model offers the greatest computational advantage and shows how GPU acceleration enables simulations over a broader range of parameters and large-scale systems that would not be feasible with CPU-only computations.

In the following sections, we present the computational methodology in detail, validate the solver against several representative test cases, highlight the substantial performance gains achieved with GPU-based computations compared to CPU-only simulations, and demonstrate code scalability through strong and weak scaling studies on both CPU and GPU architectures.

## 2. Governing equations

At the macroscale, acoustic cavitation problems are resolved using a fixed-grid Eulerian framework, while the influence of unresolved bubbles is incorporated via phase-averaged subgrid-scale models. In dilute bubbly suspensions, the propagation of acoustic waves can be described by a fully compressible continuum model of the liquid–bubble mixture. Any mixture property, denoted by $(\cdot)$ is defined as $(\cdot) = (1 - \alpha)(\cdot)_l + \alpha(\cdot)_g$, where $\alpha$ is the volume fraction of the gas contributed by the bubbles, and the subscripts $l$ and $g$ denote the liquid and gas phases, respectively. The governing conservation equations for mass, momentum, and energy take the following form:

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

$$\frac{\partial (\rho \boldsymbol{u})}{\partial t} + \nabla \cdot (\rho \boldsymbol{u} \otimes \boldsymbol{u} + p\mathcal{I} - \mathcal{T}) = \boldsymbol{0}, \tag{1}$$

$$\frac{\partial E}{\partial t} + \nabla \cdot [(E + p)\boldsymbol{u} - \mathcal{T} \cdot \boldsymbol{u}] = 0,$$

where $\rho$ is the density, $\boldsymbol{u}$ is the velocity vector, $p$ is the pressure, and $E$ is the total energy. The term $\mathcal{T}$ denotes the effective viscous stress tensor of the mixture. In dilute suspensions, the characteristic low void fraction, up to $O(10^{-2})$, allows us to treat the liquid density as significantly greater than that of the gas $\rho_l >> \rho_g$ [22, 23]. Then, we can approximate the mixture density to be $\rho \approx (1 - \alpha)\rho_l$. Additionally, we assume zero slip velocity between the phases $\boldsymbol{u} \approx \boldsymbol{u}_l = \boldsymbol{u}_g$. Thus, there is effectively no momentum transfer across the gas–liquid interface, allowing us to approximate the effective viscous stress as that of the continuous phase:

$$\mathcal{T} \equiv \mathcal{T}_l = \mu_l \left( \nabla \boldsymbol{u} + \nabla \boldsymbol{u}^\top - \tfrac{2}{3}(\nabla \cdot \boldsymbol{u})\mathcal{I} \right) \tag{2}$$
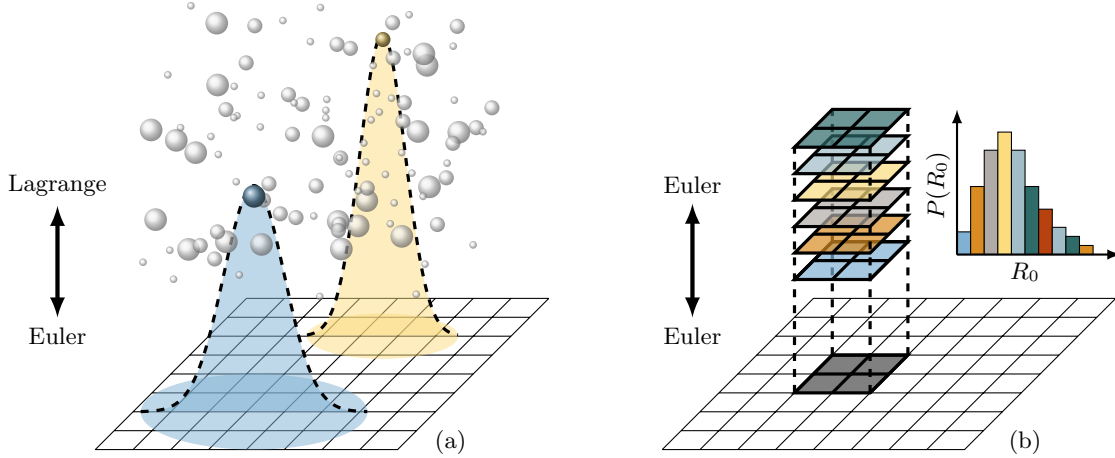
where $\mu_l$ is the liquid viscosity.

To account for the presence of the bubbles and their interaction with the surrounding liquid, we use the volume- and ensemble-averaged models. Each of them is derived from a distinct set of governing equations, based on physically justified simplifying assumptions. Detailed formulations are provided in the subsequent subsections of the manuscript. Both phase-averaged models assume that the bubbles are spherical and sufficiently separated so that collisions, bouncing, and coalescence can be ignored. Additionally, bubble–bubble interaction occurs only through their effect on the liquid-bubble mixture.

In both subgrid models, the volumetric oscillations of the bubbles in response to pressure variations in the surrounding liquid are described by the Keller–Miksis equation, which incorporates the liquid's compressibility effect, and is given by:

$$\left( R \left( 1 - \frac{\dot{R}}{c} \right) \right) \ddot{R} + \frac{3}{2} \dot{R}^2 \left( 1 - \frac{\dot{R}}{3c} \right) = \frac{p_{bw} - p_\infty}{\rho} \left( 1 + \frac{\dot{R}}{c} \right) + \frac{R\dot{p}_{bw}}{\rho c},$$

$$p_{bw} = p_b - \frac{4\mu_l \dot{R}}{R} - \frac{2\sigma}{R}, \tag{3}$$

where $R$, $\dot{R}$, and $\ddot{R}$ are the radius, interface velocity, and interface acceleration of the bubble. The term $p_{bw}$ is the pressure at the bubble wall, $p_b$ is the pressure inside the bubble, and $p_\infty$ is the pressure that forces the radial oscillations of the bubble. $\mu_l$ is the dynamic viscosity of the background medium, $\sigma$ is the surface tension, and $c$ is the liquid's speed of sound. We assume that

4

**Figure 1:** Schematics of the (a) volume-averaged and (b) ensemble-averaged subgrid bubble models.

the bubble contains both non-condensable gas and vapor, and we adopt the reduced-order models with constant heat and mass transfer coefficients at the bubble wall, as described by Preston et al. [36]. These models account for the effects of vapor and heat diffusion through the interface. Here, the vapor mass transfer rate is given by

$$\dot{m}_v = \frac{\mathcal{D}\rho_{bw}}{1 - \chi_{vw}} \frac{\partial \chi_{vw}}{\partial r}\bigg|_{r=R}, \tag{4}$$

where $\chi_v$ is the vapor mass fraction, $\mathcal{D}$ is the binary diffusion coefficient, and subscript $w$ denotes properties at the bubble wall. The internal pressure $p_b$ evolves following the model of Ando et al. [22] as follows:

$$\dot{p}_b = \frac{3\gamma_b}{R} \left( -\dot{R}p_b + \mathcal{R}_v T_{bw} \dot{m}_v + \frac{\gamma_b - 1}{\gamma_b} k_{bw} \frac{\partial T}{\partial r}\bigg|_{r=R} \right), \tag{5}$$

where $\gamma_b$ is the specific-heat ratio of the bubble contents, $\mathcal{R}_v$ is the gas constant of vapor, $T_{bw}$ is the bubble-wall temperature and $k_{bw}$ is the thermal conductivity for the bubble contents.

In dilute bubbly suspensions, external acoustic perturbations are necessary to excite bubble oscillations and sustain cavitation dynamics. In this work, we use a source-term approach to generate one-way acoustic waves [37]. This method enables the injection of unidirectional acoustic waves from an arbitrarily shaped source surface by introducing appropriate forcing terms into the mass, momentum, and energy equations on that surface within the computational domain.

### 2.1. Volume-averaged (EL) subgrid model

This model is formulated within an Euler–Lagrange (EL) framework, where the bubbles are treated as discrete entities in three-dimensional space and their interaction with the carrier fluid is represented through a two-way coupling strategy. Our volume-averaged model (or EL model) follows the formulation of Maeda and Colonius [23]. In cavitation problems, a bubble's behavior is influenced not only by its radial oscillations but also by its translational motion. However, the formulation of Maeda and Colonius [23] simplifies the problem by assuming that bubbles remain fixed in space. As a result, bubble translation, which would require solving an additional equation of motion, is neglected. This assumption provides a reasonable approximation because the timescale is typically much longer than that of the radial dynamics. Under moderate acoustic forcing and in the absence

of additional perturbations, the induced liquid velocities remain small, making bubble motion negligible compared to the rapid radial oscillations of the bubble.

The two-way Euler–Lagrange coupling is established as follows: the far-field pressure driving the bubble oscillations, $p_\infty$, is obtained from the Eulerian pressure field, while the bubbles' influence is transferred back by smearing their effect on the background grid. Applying the model assumptions to eq. (1), the following inhomogeneous hyperbolic system is derived:

$$\frac{\partial \rho_l}{\partial t} + \nabla \cdot (\rho_l \boldsymbol{u}_l) = \frac{\rho_l}{1-\alpha}\left[\frac{\partial \alpha}{\partial t} + \boldsymbol{u}_l \cdot \nabla \alpha\right],$$

$$\frac{\partial (\rho_l \boldsymbol{u}_l)}{\partial t} + \nabla \cdot (\rho_l \boldsymbol{u}_l \otimes \boldsymbol{u}_l + p_l \mathcal{I} - \mathcal{T}_l) = \frac{\rho_l \boldsymbol{u}_l}{1-\alpha}\left[\frac{\partial \alpha}{\partial t} + \boldsymbol{u}_l \cdot \nabla \alpha\right] - \frac{\alpha \nabla \cdot (p_l \mathcal{I} - \mathcal{T}_l)}{1-\alpha}, \quad (6)$$

$$\frac{\partial E_l}{\partial t} + \nabla \cdot ((E_l + p)\boldsymbol{u}_l - \mathcal{T}_l \cdot \boldsymbol{u}_l) = \frac{E_l}{1-\alpha}\left[\frac{\partial \alpha}{\partial t} + \boldsymbol{u}_l \cdot \nabla \alpha\right] - \frac{\alpha \nabla \cdot (p \boldsymbol{u}_l - \mathcal{T}_l \cdot \boldsymbol{u}_l)}{1-\alpha}.$$

The left-hand side corresponds to the conservation equations of the liquid phase, while the right-hand side represents source terms that incorporate the effects of the bubbles. The liquid pressure $p_l$ follows the stiffened-gas equation of state

$$p_l = (\gamma_l - 1)\rho_l \varepsilon_l - \gamma_l \pi_{\infty,l}, \quad (7)$$

where $\varepsilon_l$ is the liquid's internal energy, and $\gamma_l$ and $\pi_{\infty,l}$ are the specific heat ratio and the stiffness of the liquid, respectively.

The instantaneous bubble sizes are communicated to the background flow solver through the local void fraction, $\alpha$. This coupling is achieved by computing an effective void fraction that distributes each bubble's volume contribution to the surrounding computational cells, as illustrated in fig. 1 (a). At the bubble center $\boldsymbol{x}_n$ of bubble $n$, the volume is smeared into the continuous void fraction field by a regularization kernel $\delta$, giving

$$\alpha(\boldsymbol{x}_n) = \sum_{n=1}^{N} V_n \delta = \sum_{n=1}^{N}\left(\frac{4}{3}\pi R_n^3\right)\delta, \quad (8)$$

where $N$ is the total number of bubbles and $V_n$ is the volume. We use a continuous, second-order truncated Gaussian kernel:

$$\delta(d_n, h) = \begin{cases} \frac{1}{h^3 (2\pi)^{3/2}} e^{-\frac{d_n^2}{2h^2}} & 0 \le \frac{d_n}{h} < 3, \\ 0, & 3 \le \frac{d_n}{h}, \end{cases} \quad (9)$$

where $h$ is the kernel width and $d_n = |\boldsymbol{x} - \boldsymbol{x}_n|$ is the distance to the bubble center. The evolution of $\alpha$ follows as:

$$\frac{\partial \alpha(\boldsymbol{x})}{\partial t} = \sum_{n=1}^{N}\frac{\partial V_n}{\partial t}\delta + \sum_{n=1}^{N}V_n\frac{\partial \delta}{\partial t}, \quad (10)$$

with

$$\frac{\partial V_n}{\partial t} = 4\pi R_n^2 \dot{R}_n, \quad \frac{\partial \delta}{\partial t} = -\boldsymbol{u}_l \cdot \nabla \delta. \quad (11)$$

Maeda and Colonius [23] impose a constraint on the kernel support width $h$ to ensure that the model resolves the small-scale dynamics inside the bubble cloud:

$$\begin{Bmatrix} R_b \\ \Delta \end{Bmatrix} \le h < L_b, \quad (12)$$

where $R_b$ is the characteristic bubble radius, $\Delta$ is the Eulerian grid spacing, and $L_b$ is the characteristic inter-bubble distance. This condition prevents kernel supports from overlapping and ensures that the physics at the inter-bubble scale are accurately captured. Because the maximum bubble radius $R_b$ is not known a priori, the solver dynamically adjusts $h$ to satisfy eq. (12). As long as this condition holds, the model is valid even when bubbles grow larger than the grid size. Accordingly, the solver sets $h = \Delta$ when $R_b < \Delta$, and $h = R_b$ otherwise. A more detailed description of the volume-averaged model can be found in [23].

*2.2. Ensemble-averaged (EE) subgrid model*

This model can be viewed as an Euler–Euler (called EE) formulation, as illustrated in fig. 1 (b). Here, instead of solving for the dynamics of individual bubbles, it evaluates the statistically-averaged mixture dynamics by assuming a large number of stochastically scattered bubbles dispersed within each computational grid cell. Our ensemble-averaged model follows the description of Zhang and Prosperetti [21] and Bryngelson et al. [25].

The equilibrium radii of the bubbles, $\boldsymbol{R}_0$, are assumed to follow a log-normal distribution, which is further discretized using $N_{\text{bin}}$ number of bins. The bubble population is represented statistically through the variables $\boldsymbol{R}$, $\dot{\boldsymbol{R}}$, $\boldsymbol{p}_b$, and $\boldsymbol{m}_v$, corresponding to the instantaneous bubble radii, interface velocities, bubble pressures, and vapor mass. Each of these variables contains $N_{\text{bin}}$ components. The mixture-averaged pressure in eq. (1) is given by

$$p = (1 - \alpha)p_\ell + \alpha \left( \frac{\overline{\boldsymbol{R}^3 \boldsymbol{p}_{bw}}}{\overline{\boldsymbol{R}^3}} - \rho \frac{\overline{\boldsymbol{R}^3 \dot{\boldsymbol{R}}^2}}{\overline{\boldsymbol{R}^3}} \right), \tag{13}$$

where $\boldsymbol{p}_{bw}$ is the associated bubble wall pressure defined in eq. (3). The liquid pressure $p_\ell$ is computed using the stiffened-gas equation of state given in eq. (7). The bubble number density per unit volume $n_{\text{bub.}}(\boldsymbol{x}, t)$ is conserved as

$$\frac{\partial n_{\text{bub.}}}{\partial t} + \nabla \cdot (n_{\text{bub.}} \boldsymbol{u}) = 0. \tag{14}$$

For the spherical bubbles considered here, $n_{\text{bub.}}$ is related to the void fraction, $\alpha$, via the conservation of the number density function:

$$\alpha(\boldsymbol{x}, t) = \frac{4}{3} \pi \overline{\boldsymbol{R}^3} n_{\text{bub.}}(\boldsymbol{x}, t), \tag{15}$$

and so the void fraction $\alpha(\boldsymbol{x}, t)$ transports as

$$\frac{\partial \alpha}{\partial t} + \nabla \cdot (\alpha \boldsymbol{u}) = 3\alpha \frac{\overline{\boldsymbol{R}^2 \dot{\boldsymbol{R}}}}{\overline{\boldsymbol{R}^3}}, \tag{16}$$

where the right-hand side represents the change in void fraction resulting from the growth and collapse of bubbles. The bubble dynamics are evaluated as

$$\frac{\partial n_{\text{bub.}} \boldsymbol{\phi}}{\partial t} + \nabla \cdot (n_{\text{bub.}} \boldsymbol{\phi} \boldsymbol{u}) = n_{\text{bub.}} \dot{\boldsymbol{\phi}}, \tag{17}$$

where $\boldsymbol{\phi} = \{\boldsymbol{R}, \dot{\boldsymbol{R}}, \boldsymbol{p}_b, \boldsymbol{m}_v\}$ contains the bubble dynamic variables. The over-barred terms in the above equations denote averages computed across the bubble dispersion, confining all $N_{\text{bin}}$ bubble groups in each control volume. All over-barred terms require a numerical closure, which is accomplished by

distributing the bubble equilibrium radii $R_o$ in $N_{\text{bin}}$ bins with a log-normal probability distribution function (PDF). The integration of these terms follows from Simpson's rule, though more advanced techniques are available in the limit of small bubble oscillation amplitude [38, 39]. The Euler–Euler ensemble averaging technique can be readily extended to a population balance formulation that accounts for distributions in all independent bubble coordinates, $R$, $\dot{R}$, and $R_o$ [40, 41]. More details of the ensemble-averaged model can be found in [21, 42].

## 3. Numerical method

The conservative form of the set of governing equations eq. (1) can be generalized as:

$$\frac{\partial \boldsymbol{q}}{\partial t} + \nabla \cdot \boldsymbol{f}(\boldsymbol{q}) = \boldsymbol{s}, \tag{18}$$

where, $\boldsymbol{q}$ is the vector of conservative variables, $\boldsymbol{f}(\boldsymbol{q})$ are the fluxes, and $\boldsymbol{s}$ contains any source terms. To numerically solve eq. (18), MFC employs the finite volume method. This equation can be spatially discretized in a Cartesian framework as

$$\frac{\partial \boldsymbol{q}}{\partial t} + \frac{\partial \boldsymbol{f}^x(\boldsymbol{q})}{\partial x} + \frac{\partial \boldsymbol{f}^y(\boldsymbol{q})}{\partial y} + \frac{\partial \boldsymbol{f}^z(\boldsymbol{q})}{\partial z} = \boldsymbol{s}, \tag{19}$$

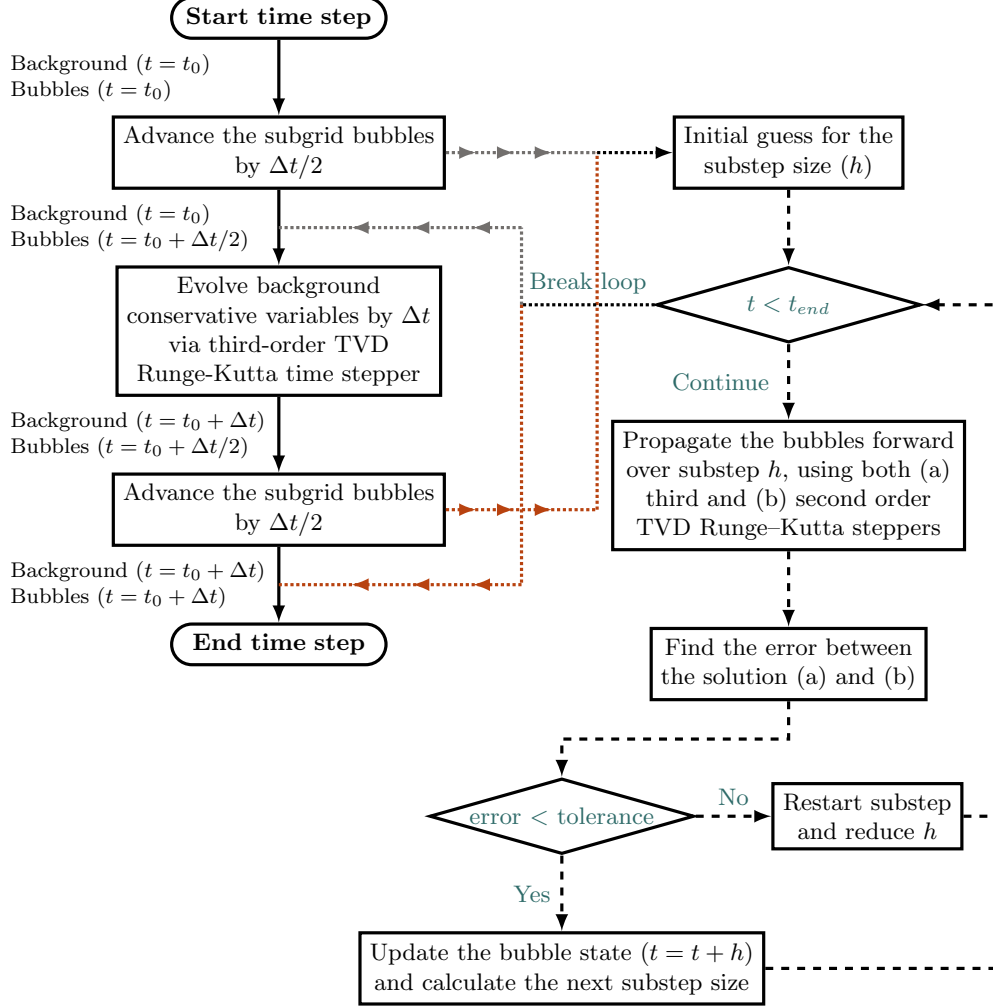where $\boldsymbol{f}^x, \boldsymbol{f}^y$ and $\boldsymbol{f}^z$ are vectors of fluxes in the $x$, $y$, and $z$ directions. The above equation is integrated over each finite volume grid cell $(i, j, k)$ in the three coordinate directions. The resultant equation in semi-discrete form is

$$\begin{aligned}
\frac{\mathrm{d}\boldsymbol{q}_{i,j,k}}{\mathrm{d}t} = \frac{1}{\Delta x_i}\left[\boldsymbol{f}^x_{i-1/2,j,k} - \boldsymbol{f}^x_{i+1/2,j,k}\right] + \frac{1}{\Delta y_j}\left[\boldsymbol{f}^y_{i,j-1/2,k} - \boldsymbol{f}^y_{i,j+1/2,k}\right] + \\
\frac{1}{\Delta z_k}\left[\boldsymbol{f}^z_{i,j,k-1/2} - \boldsymbol{f}^z_{i,j,k+1/2}\right] + \boldsymbol{s}_{i,j,k}.
\end{aligned} \tag{20}$$

We use the HLLC approximate Riemann solver to compute the fluxes of the primitive variables across the cell faces. The right and left states for the Riemann problem follow from a fifth-order accurate WENO reconstruction, which is robust to grid-scale phase and state discontinuities [25].

To advance the conservative variables in time, we adopt a third-order total variation diminishing (TVD) Runge–Kutta time-stepping scheme [43]. This method strikes an effective balance between numerical accuracy and stability, particularly in the presence of sharp gradients. Unlike the relatively smooth evolution of the grid-resolved background flow, subgrid-scale bubbles often undergo transient and nonlinear dynamics, particularly during events such as bubble collapse and rebound. These phenomena involve rapid variations in pressure and volume, necessitating significantly smaller time steps to maintain numerical stability and accurately capture the physics. However, applying such fine time resolution uniformly across the bubbles and background flow would drastically increase computational cost, making the simulations infeasible. To address this challenge, we employ a Strang splitting algorithm [44], which decouples the time step of the background flow from that of the bubble dynamics. Figure 2 schematically illustrates the implementation of the Strang splitting method within our framework. This approach enables us to treat the rapid, localized dynamics of subgrid bubbles independent of the slower evolution of the background flow. In practice, the bubble equations are solved using smaller time steps nested within each coarser background time step. To ensure accuracy and stability, the solver only advances each substep if the relative errors in the bubble radius $R$ and interface velocity $\dot{R}$ calculated using third- and second-order TVD Runge–Kutta

**Start time step**

Background ($t = t_0$)
Bubbles ($t = t_0$)

Advance the subgrid bubbles by $\Delta t/2$

Initial guess for the substep size ($h$)

Background ($t = t_0$)
Bubbles ($t = t_0 + \Delta t/2$)

Evolve background conservative variables by $\Delta t$ via third-order TVD Runge-Kutta time stepper

Break loop

$t < t_{end}$

Continue

Background ($t = t_0 + \Delta t$)
Bubbles ($t = t_0 + \Delta t/2$)

Advance the subgrid bubbles by $\Delta t/2$

Propagate the bubbles forward over substep $h$, using both (a) third and (b) second order TVD Runge–Kutta steppers

Background ($t = t_0 + \Delta t$)
Bubbles ($t = t_0 + \Delta t$)

**End time step**

Find the error between the solution (a) and (b)

error < tolerance

No

Restart substep and reduce $h$

Yes

Update the bubble state ($t = t + h$) and calculate the next substep size

**Figure 2:** General schematic of the Strang splitting algorithm for time integration.

steppers remain below a tolerance of $10^{-4}$. The algorithm alternates between advancing the states of the subgrid bubbles and the conservative variables of the flow field, ensuring consistent coupling between them. By manipulating the timescale requirements of the two systems, this strategy enables the efficient simulation of complex multiscale phenomena without incurring prohibitive computational costs.

## 4. Hardware acceleration strategy

The multiscale algorithm is hardware-accelerated using graphics processing units (GPUs). The grid-resolved background flow follows the acceleration strategy developed by Radhakrishnan et al. [31], which is implemented in MFC. Building upon this framework, we extended the GPU acceleration to our previously described phase-averaged subgrid bubble models. In the volume-averaged (Euler–Lagrange) formulation, the acceleration targets the evolution of discrete bubbles and the two-way coupling routines within the EL framework. The ensemble-averaged (Euler–Euler) formulation focuses on accelerating the transport of bubble number density, void fraction, and averaged bubble dynamics. This integrated approach enables full exploitation of GPU-based parallelism across both the micro- and macro-scale components of the simulation.

Once the initial conditions are established on the CPU, the state variables are transferred to the GPU, which subsequently performs the majority of the computational workload. OpenACC directives are used to offload all computationally intensive tasks, ensuring that parallel regions and loop structures are automatically mapped to the GPU architecture. After identifying independent loops and defining their levels of parallelization, the OpenACC runtime automatically selects optimal kernel configurations, maximizing performance and parallel efficiency. This automated tuning process ensures that the algorithm is well-optimized for the target GPU hardware. A key advantage of this directive-based offloading approach is maintaining a unified codebase for both CPU and GPU execution. Rather than maintaining separate implementations for different hardware platforms, the OpenACC directives allow the compiler to generate architecture-specific code at compile time via simple compiler flags. By adhering to standard OpenACC syntax, the implementation remains portable across multiple computing environments, as these directives are supported by compilers such as NVHPC, GNU, and Cray (CCE). This flexibility ensures compatibility with both NVIDIA and AMD GPUs. In this study, GPU results were obtained using the NVHPC 24.1 SDK, and CPU-based runs used GNU 11.4. No meaningful performance differences were observed for different compiler versions. MFC also uses metaprogramming techniques enabled by the Fypp preprocessor to enhance GPU kernel performance [31]. User-defined inputs are treated as compile-time constants, allowing the compiler to allocate fixed-size thread-local arrays and optimize memory access through register utilization. This approach also removes conditional branching and redundant kernel duplication across spatial dimensions. Additional optimizations are employed, which are described in full in Radhakrishnan et al. [31]. These optimizations contribute to MFC's high computational performance and serve as the foundation for the acceleration strategy adopted in this work.

An example of the OpenACC kernel used to incorporate the influence of Lagrangian bubble volume within the Eulerian framework is shown in listing 1. The `parallel loop` construct is augmented with a `gang vector` clause, which enables the compiler to automatically determine the optimal number of gangs and vectors, maximizing resource usage [45]. The outer loop iterates over all bubbles, and each bubble is subsequently processed through three nested loops spanning the spatial coordinates. These loops are collapsed into a single memory-coalesced loop via the `loop collapse(3)` clause, improving loop-level parallelization across the three dimensions. In cases where bubbles are located in proximity, multiple threads may attempt to update the same computational cell concurrently, leading to potential race conditions and data corruption. To address this, OpenACC provides the `atomic update` clause, which enforces thread-safe access to memory. When multiple threads update the same cell, their operations are serialized to ensure that all contributions are correctly accumulated without data loss or overwriting. In our implementation, the `atomic update` clause is applied selectively to ensure accurate updates to the Eulerian grid induced by the bubbles, while preserving parallel efficiency.

Similarly, listing 2 presents the OpenACC kernel used to advance the time evolution of the averaged bubble state in the EE model. As before, the `gang vector` clause allows the compiler to automatically select the optimal number of gangs and vectors to maximize utilization of the available GPU resources. Three nested spatial loops are merged using the `loop collapse(3)` clause. In contrast, an additional `loop seq` clause is employed to iterate over the discrete bins of the log-normal probability distribution, which represents the initial bubble size distribution. Within this loop, the Strang Splitting subroutine is invoked to apply the algorithm illustrated in fig. 2. Finally, the scalar fields in the Eulerian framework where the subgrid bubbles reside are updated accordingly. In this step, the use of an `atomic update` clause is unnecessary, as each GPU thread operates on a unique entry (`j`, `k`, `l`, `q`) defined by its $x$-, $y$-, and $z$-coordinate indices and bin index, respectively.

**Listing 1:** EL Model: OpenACC directives to smear the bubble volume.

```
!$acc parallel loop gang vector
!$acc default(present) private(..)
do q = 1, nBubs
  cell = f_find_cell(q)
  stdev = f_kernel_width(q)

  !$acc loop collapse(3)
  do j = 0,delt; do k = 0,delt; do l = 0,delt
  ! Coordinate directions of local smearing
    aux(1) = cell(1) + j - delt
    aux(2) = cell(2) + k - delt
    aux(3) = cell(3) + l - delt
    call s_check_outside(aux)
    if (.not. outside) then
      call s_gaussian(stdev, fun)
    end if

    !$acc atomic update
    updVar%sf(aux) += fun*bubVolume
  end do; end do; end do
end do
!$acc end parallel loop
```

**Listing 2:** EE Model: OpenACC directives to evolve the bubble state.

```
!$acc parallel loop collapse(3)
!$acc gang vector default(present)
!$acc private(..)
do j = 0,Nx; do k = 0,Ny; do l = 0,Nz
! Loop over grid

  ! Bubble number density
  nbub = f_obtain_nbub()

  !$acc loop seq
  do q = 1,nBins
    bub = f_bubState(q)
    if (adap_dt) then !time advance
      ! Strang Split
      call s_advance(bub, R_loc, Rdot_loc)
      ! Update variables
      q(r(q))%sf(j,k,l)= nbub*R_loc
      q(v(q))%sf(j,k,l)= nbub*Rdot_loc
    end if
  end do
end do; end do; end do
!$acc end parallel loop
```
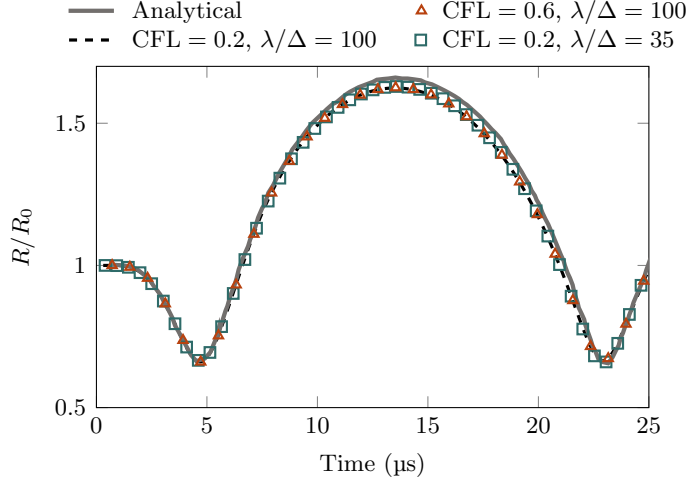
## 5. Results and discussion

### 5.1. Validation

The accuracy of our volume-averaged multiscale solver is rigorously validated through two distinct test cases, each designed to assess its ability to accurately predict the dynamic behavior of oscillating and spherically collapsing bubbles.
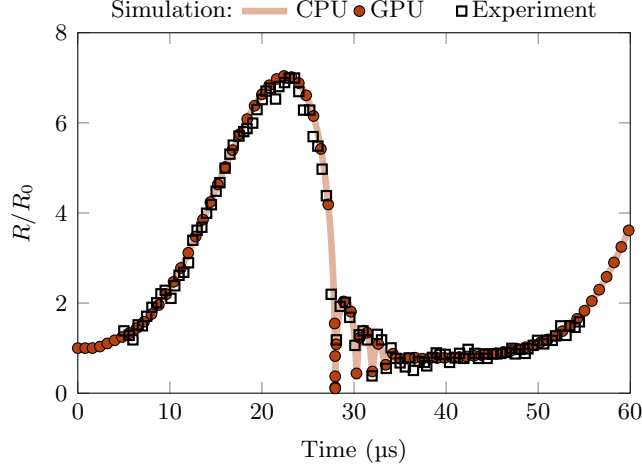
In the first validation scenario, we consider an isolated gas bubble in water. Its initial radius is $50\,\mu m$, and is positioned at a fixed location $(0, 0, 0)$ at the center of the computational domain. This bubble is exposed to a planar sinusoidal acoustic wave with an amplitude of $0.2\,MPa$ and a frequency of $150\,kHz$. To prevent acoustic reflections through the boundaries, we use the non-reflective boundary conditions described by Thompson [46]. The temporal evolution of the bubble radius, shown in fig. 3, is compared against the analytical solution of the Keller–Miksis equation, as reported by Maeda and Colonius [23]. We performed a series of simulations to quantify the solver's variability with respect to the Courant–Friedrichs–Lewy (CFL) number and background grid size. Because the high-order TVD Runge–Kutta scheme is explicit, the numerical stability criteria require a CFL number below unity. MFC can automatically adjust the time step to ensure that the CFL number remains below a user-defined maximum throughout the domain. In our simulations, we observe negligible variability as the CFL number is increased from 0.2 to 0.6. We also examined the influence of grid resolution while holding the CFL number fixed at 0.2. The results again show negligible sensitivity, with comparable accuracy when resolving the acoustic wavelength using 35 or 100 cells. Overall, this parametric study highlights the temporal and spatial robustness of the solver. The numerical predictions exhibit excellent agreement with the analytical solution, with a maximum root-mean-squared error (RMSE) of $2.76\,\%$.

**Figure 3:** Evolution of an isolated bubble in response to a single cycle of a sinusoidal pressure wave using different CFL numbers and grid sizes.

In the second validation case, we replicate the experimental observations reported by Ohl et al. [47], who investigated the dynamics of a trapped bubble in a water–glycerine mixture undergoing spherical collapse. The liquid host has a density of $1000\,\text{kg/m}^3$, viscosity of $0.006\,\text{N\,s/m}^2$, and the surface tension of the liquid-bubble interface is $0.07\,\text{N/m}$. A single bubble with an initial radius of $R_0 = 8\,\mu\text{m}$ is fixed in space and subjected to a sinusoidal acoustic wave with an amplitude of $1.32\,\text{bar}$ and a frequency of $21.4\,\text{kHz}$. The grid spacing is uniform across the domain with $\Delta = 100\,\mu\text{m}$, and the CFL number is kept at 0.2. For this particular simulation, we disabled the mass transfer model given in eq. (4). The reduced-order formulation developed by Preston et al. [36] assumes that the vapor pressure remains in equilibrium at the gas–liquid interface. Under strong bubble collapses, however, this assumption becomes less certain, as the rapid dynamics may prevent the interface from maintaining thermodynamic equilibrium. Consequently, we chose not to include this model here to maintain consistency with the experimental observations. The expansion part of the pressure wave makes the bubble grow, and we observe a maximum bubble radius of $7.04 R_0$, as depicted in fig. 4. Then, the compressive part of the wave provokes a sudden collapse of the bubble, followed by a rebound cycle. The simulated bubble radius evolution during growth, collapse, and rebound cycles shows close agreement with the experimental measurements, yielding an RMSE of $7.46\,\%$. Furthermore, the results obtained using GPU-based simulations are identical to those from CPU-based computations, confirming the correct implementation of GPU offloading and data-parallelism within our solver.

To evaluate the accuracy of the ensemble-averaged subgrid model, we simulate the interaction between a dilute bubble screen and a single sinusoidal planar acoustic wave. Both the volume-averaged and ensemble-averaged subgrid models are employed to enable a systematic comparison. Using our previously validated volume-averaged model as a reference, we test the performance of the ensemble-averaged formulation. The computational domain is a square prism, shown in fig. 5, with spatial extents defined as $x \in [-20, 20]$ mm and $y, z \in [-2.5, 2.5]$ mm. The bubble distribution is cubic, centered at the origin $(0, 0, 0)$, extends 5 mm along the $x$-direction, and spans the full extent of the domain in the $y$- and $z$-directions. The screen consists of a cloud of bubbles with an initial void fraction of $\alpha_0 = 4 \times 10^{-5}$, ensuring the flow remains in the dilute regime, and all bubbles are initialized in their equilibrium states.
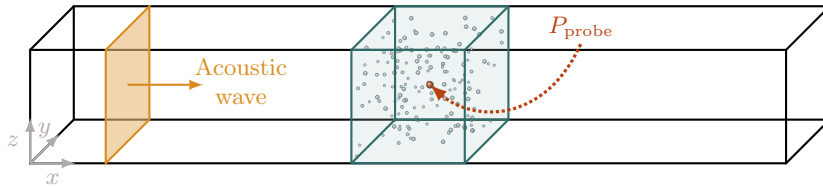
**Figure 4:** Radius evolution of the spherical collapse of an isolated bubble in response to a sinusoidal acoustic wave.
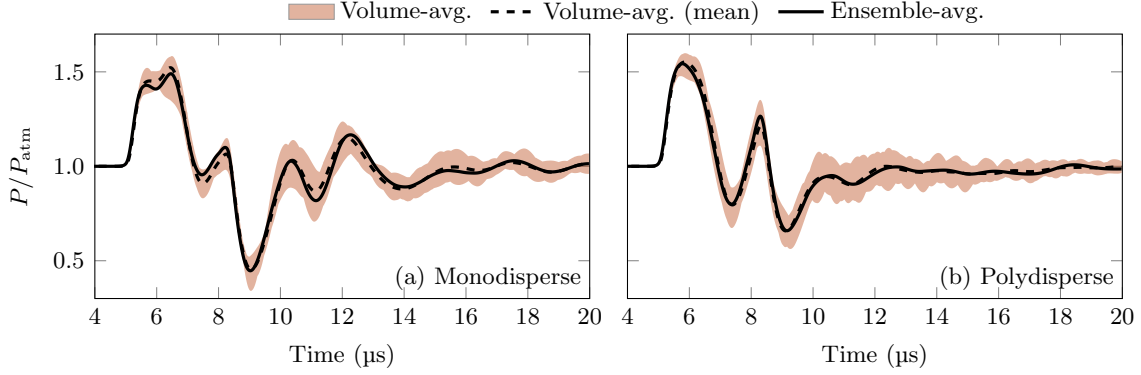
Two test cases are considered for the bubble size distribution: a monodisperse case and a polydisperse case. In the monodisperse scenario, all bubbles have the same initial radius of $10\,\mu m$. In the polydisperse scenario, a more realistic distribution of bubble sizes is introduced by assigning radii according to a log-normal distribution, also centered at $10\,\mu m$, with a log-normal shape parameter of $\sigma_p = 0.3$. To excite the bubble screen, we impose a planar sinusoidal acoustic wave that is generated at $x = -7.5\,mm$ and propagates in the positive $x$-direction toward the bubble cloud. The wave has a frequency of $300\,kHz$ and an amplitude of $0.1\,MPa$. The surrounding fluid is water, and the gas inside the bubbles is air, representing a typical two-phase system in acoustically active environments such as underwater acoustics or biomedical ultrasound. The domain boundaries are configured with non-reflective boundary conditions to prevent acoustic reflections.

A 3D Cartesian grid is used for spatial discretization. The grid is uniformly spaced, with 400 cells in the $x$-direction and 50 cells in the $y$- and $z$-directions. This configuration results in a uniform cell size of $\Delta_x = \Delta_y = \Delta_z = 100\,\mu m$. The choice of grid resolution is guided by the need to resolve the acoustic wavelength with sufficient accuracy. Specifically, the acoustic wave is resolved using 50 cells per wavelength, which is adequate to prevent numerical dissipation that could otherwise attenuate the wave during its propagation through the medium [24]. We maintain a CFL number of 0.2 throughout all simulations.

For simulations using the volume-averaged model, bubbles are randomly placed within the defined cubic bubble screen. To account for statistical variations introduced by the random placement of bubbles and to extract a meaningful average behavior, we perform 40 independent simulations, which are sufficient for statistical convergence, as shown in a previous study [25]. The mean inter-bubble



**Figure 5:** Schematic of the dilute bubble screen configuration (not to scale).

**Figure 6:** Pressure profiles at the origin of the bubble screen using the volume-averaged and ensemble-averaged subgrid models. (a) Monodisperse and (b) polydisperse bubble cloud.
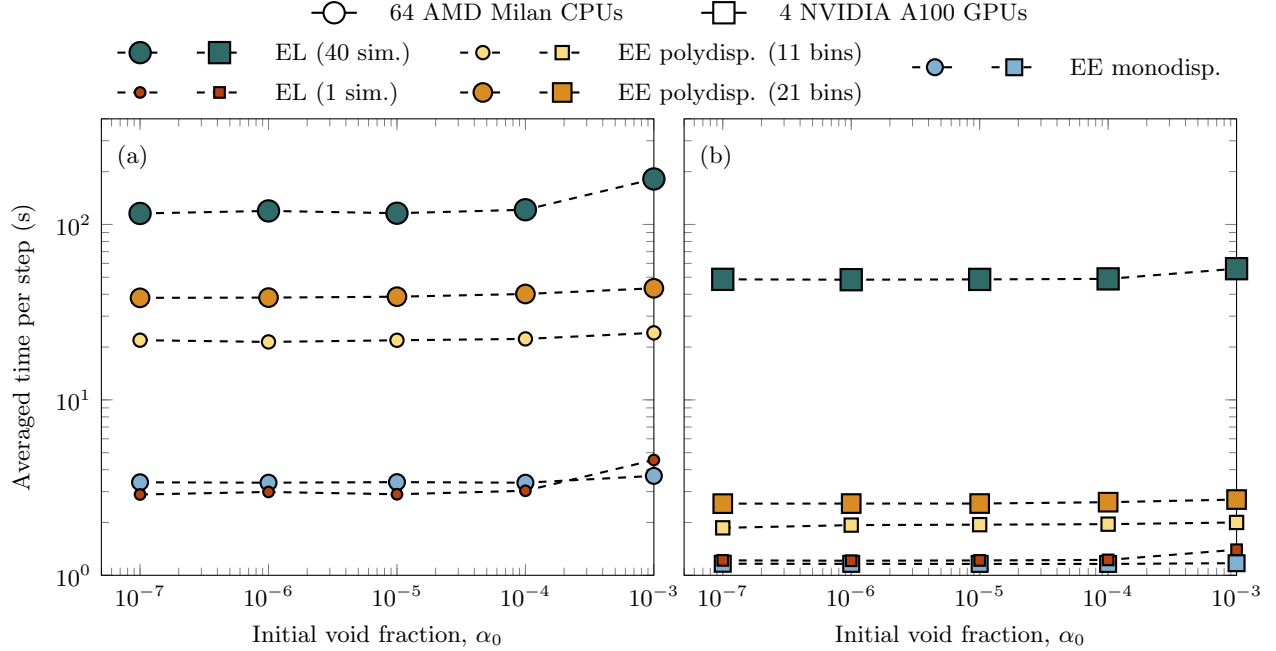
distance for the given void fraction and bubble radius is approximately $265\,\mu\text{m}$, which is more than twice the grid spacing. This satisfies the requirement of eq. (12).

In the ensemble-averaged modeling framework, the polydispersity of bubble sizes introduces additional complexity. Instead of simulating the full spatial distribution of individual bubbles, the model resolves the average behavior across a spectrum of equilibrium bubble radii. To capture this variability, the radius space is discretized into a finite number of bins, $N_{\text{bin}}$, each representing a subset of the log-normal distribution. Bryngelson et al. [25] conducted a parametric study to quantify the sensitivity of the model predictions for different $N_{\text{bin}}$, showing that the error decreases as the number of bins increases. For example, they report an $8\,\%$ error for $\sigma_p = 0.3$ when using $N_{\text{bin}} = 11$. To reduce this error and enhance the model's predictive capability, we increase the resolution of the bubble size spectrum to $N_{\text{bin}} = 21$ in our polydisperse simulations. This choice provides a practical balance between accuracy and computational cost.

Figure 6 shows the pressure measurements at the center of the bubble cloud, $P_{\text{probe}}(0,0,0)$, for both bubble distributions. The shaded region represents the variability range observed across all Euler–Lagrange realizations. In both monodisperse and polydisperse configurations, the pressure profiles computed using the ensemble-averaged model closely follow the mean behavior of the 40 volume-averaged simulations. The RMSE error between the ensemble-averaged and volume-averaged models is $2.10\,\%$ for the monodisperse case and $1.53\,\%$ for the polydisperse case. These results confirm the accuracy and validity of the ensemble-averaged subgrid model for dilute bubbly flows.

*5.2. Computational cost of the bubbles*

In this section, we evaluate the computational expense associated with our phase-averaged bubble models and demonstrate the performance benefits of the hardware-accelerated solver utilizing GPUs. We consider a representative bubble screen configuration with varying initial void fraction values, ranging from $\alpha_0 = 10^{-7}$ to $\alpha_0 = 10^{-3}$. The computational domain is a cubic volume of water with dimensions $x, y, z \in [-7.5, 7.5]$ mm, within which a bubble cloud is initialized. The bubble radii follow a log-normal distribution centered at $R_0 = 10\,\mu\text{m}$, with a shape parameter of $\sigma_p = 0$ for monodisperse and $\sigma_p = 0.3$ for polydisperse cases. The bubbles are initialized out of equilibrium; a non-zero radial velocity $\dot{R}_0 = 0.015\,\text{m/s}$ is prescribed at the gas–liquid interface, which in the polydisperse case varies linearly with bubble size. This initialization promotes bubble oscillations in the absence of external excitation sources such as traveling pressure waves. The computational

**Figure 7:** Computational cost of the volume-averaged (EL) and ensemble-averaged (EE) models on (a) CPU cores and (b) GPUs.

cost of the EL model primarily depends on the total number of oscillating discrete bubbles in the domain and is independent of the bubble size distribution. Consequently, we restrict our EL tests to the monodisperse bubble cloud configuration. In contrast, the cost of the ensemble-averaged model depends on the discretization of the bubble size distribution. Specifically, the number of equations solved scales with the number of bins ($N_{bin}$) used to represent the log-normal distribution of the bubble radii. In our simulations, the EE monodisperse model corresponds to $N_{bin} = 1$, while the EE polydisperse cases employ $N_{bin} = 11$ and $N_{bin} = 21$.

The Eulerian background flow is discretized using a uniform grid of $\Delta_x = \Delta_y = \Delta_z = 75\,\mu\text{m}$. To enable a fair comparison between the subgrid models, this background mesh is kept fixed across all simulations. We use a constant time step of $\Delta t = 0.03\,\mu\text{s}$, and each simulation runs for a total duration of $1\,\mu\text{s}$ using 33 time steps. After discarding the first three start-up time steps, we compute the averaged wall-clock time per step, as reported in fig. 7. All simulations were executed on the NCSA Delta supercomputer; detailed hardware specifications are available in NCSA [48]. We use the 4-way NVIDIA A100 GPU compute nodes, which contain 4 NVIDIA A100 GPUs and 64 AMD Milan CPU cores per node. Based on this, we use one node to run each of our simulations. Thus, CPU-based simulations use 64 CPU cores and GPU simulations use 4 GPUs, each linked to one CPU core. Figure 7 summarizes the simulation outcomes. The EE model demonstrates a clear advantage in computational efficiency, as it requires only one simulation to capture the mean dynamics of the dilute bubbly flow. In contrast, the EL model requires multiple independent realizations (40 in this case, following Bryngelson et al. [25]) to achieve statistically converged results, due to the stochastic distribution of bubble locations, which significantly increases its overall cost. The computational expense of the EE model scales with the number of bins. The 21-bin configuration is the most expensive, followed by the 11-bin and monodisperse (1-bin) cases. Notably, the performance of EE simulations remains largely unaffected by variations in void fraction across both CPU and GPU architectures.

| $\alpha_0$ | (a) 64 AMD Milan CPU cores | | (b) 4 NVIDIA A100 GPUs | |
|---|---|---|---|---|
| | max(N. cells) | max(N. bubbles) | max(N. cells) | max(N. bubbles) |
| $10^{-7}$ | 119458 | 6 | 1980050 | 81 |
| $10^{-6}$ | 119458 | 40 | 1980050 | 806 |
| $10^{-5}$ | 119458 | 255 | 1980050 | 8058 |
| $10^{-4}$ | 119458 | 2435 | 1980050 | 80573 |
| $10^{-3}$ | 119458 | 23891 | 1980050 | 805722 |

**Table 1:** Euler–Lagrange simulation configurations, showing the maximum number of cells and discrete bubbles per processor on CPU and GPU architectures.

GPU acceleration substantially reduces computational time in all tested configurations. We observe maximum speedups of 16x for the EE polydisperse case with 21 bins, 12x for 11 bins, 3.1x for the EE monodisperse case, and 3.2x for the EL model. On CPUs, the EL model's cost increases significantly with void fraction, especially for $\alpha_0 = 1 \times 10^{-3}$, as a larger number of discrete bubbles intensifies the per-core workload. Despite full parallelization, each CPU core handles bubble dynamics in a largely serialized fashion, in which its computational time scales directly with the number of bubbles assigned to it, which are shown in table 1. This leads to potential load imbalance if the bubble distribution is uneven. Conversely, each NVIDIA A100 GPU contains approximately 7K CUDA cores, enabling an additional layer of parallelization that accelerates bubble dynamics computations and mitigates imbalances. The computational cost when performing a single simulation of the EL model is comparable to that of the EE monodisperse model.

These computational cost assessments provide a comprehensive understanding of the expected performance of each subgrid model in both CPU and GPU architectures. Beyond the performance metrics, the choice between the phase-averaged subgrid models depends largely on the underlying flow characteristics, the degree of spatial and temporal variability, and the level of statistical fidelity required for the analysis. The EL model is particularly advantageous in configurations where resolving the spatial distribution and individual dynamics of discrete bubbles is required to represent or infer physical features. Its Lagrangian framework naturally accommodates bubble–flow coupling, making it the most accurate choice when bubble motion or specific spatial distribution play a critical role. Still, the requirement for multiple statistically independent realizations to obtain converged ensemble-averaged quantities significantly increases the computational expense of the EL approach, particularly for large void fractions or bubble counts. In contrast, the EE model is more appropriate for statistically homogeneous bubbly flows where averaged properties are sufficient to describe the overall behavior. Its formulation allows for a continuous representation of the dispersed phase through a set of discrete size bins, efficiently capturing both monodisperse and polydisperse distributions. This discrete-bin treatment results in substantial reductions in computational cost while preserving the essential dynamics of the bubble population, including size evolution and interaction with the carrier fluid. Another important consideration is memory usage, which becomes critical in large-scale simulations. Among the models considered, the polydisperse EE approach is the most memory-intensive, particularly when a large number of bins is used to resolve the bubble size spectrum. The total memory footprint scales with both the number of bins and the resolution of the Eulerian mesh, imposing practical constraints on simulations at high spatial resolutions or when a wide size distribution is represented. Consequently, the balance between computational speed, memory requirements, and model fidelity must be carefully evaluated when selecting the appropriate subgrid model for a given problem.

The simulations discussed so far were carried out on a single compute node. To assess their efficiency as the problem size changes, we proceed with a scaling study, described in the following section.
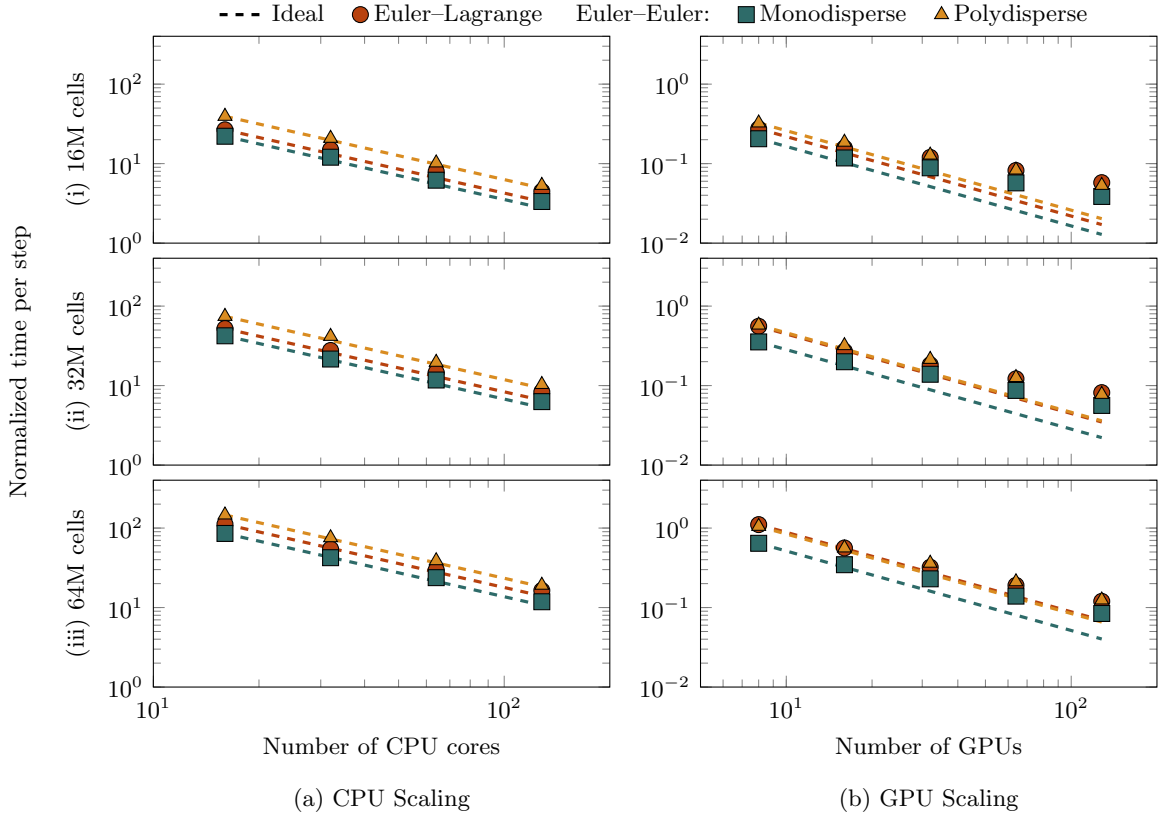
### 5.3. Scaling study

Here, we evaluate the scalability of the subgrid bubble models on both CPU- and GPU-based architectures through strong and weak scaling analyses. The computational time is averaged over twenty time steps, which we find sufficient for the performance metrics to reach a steady state. The physical configuration consists of a liquid domain filled with water serving as the base medium, within which a population of gas bubbles is dispersed at an initial void fraction of $\alpha_0 = 1 \times 10^{-3}$. Consistent with our earlier study cases, the bubble sizes follow a log-normal distribution centered at $10\,\mu m$, with $\sigma_p = 0$ for monodisperse and $\sigma_p = 0.3$ for polydisperse clouds. As in the previous study, the bubbles are initialized in a non-equilibrium state with an interface velocity of $\dot{R}_0 = 0.015\,m/s$, which in the polydisperse case changes linearly with bubble size. For the EE polydisperse scaling tests, we employ $N_{\mathrm{bin}} = 3$, as increasing the number of bins significantly raises memory requirements and leads to out-of-memory issues during scaling studies. For the EL tests, we focus on monodisperse bubble clouds, as in the previous section, with bubble positions randomly initialized to mimic realistic conditions. The background mesh resolution is held constant across both weak and strong scaling studies, with $\Delta_x = \Delta_y = \Delta_z = 80\,\mu m$.

The scaling simulations were conducted on the NCSA Delta supercomputer [48]. CPU-based simulations were executed on Delta's CPU nodes, each equipped with dual AMD EPYC 7764 (Milan) processors, providing a total of 128 cores per node across 132 nodes. GPU-based simulations were performed on Delta's 4-way NVIDIA A100 GPU compute nodes, comprising a total of 100 nodes. Each GPU node includes four NVIDIA A100 GPUs and 64 AMD Milan CPU cores. All Delta nodes are interconnected via the HPE/Cray Slingshot 11 high-performance interconnect, which supports data transfer rates of up to $200\,Gbit/s$ and enables high-bandwidth communication across the system. In the GPU configuration, each GPU is directly coupled to a dedicated CPU core.

### 5.3.1. Strong scaling

In strong scaling, the total problem size is fixed, and performance is evaluated as the number of computational partitions increases. Accordingly, the total number of Eulerian grid cells and, for the EL model, the total number of discrete bubbles, remains constant across simulations. To comprehensively assess strong scaling behavior, three (fig. 8 (i)–(iii)) increasingly large problem sizes were considered, consisting of 16 million, 32 million, and 64 million 3D Eulerian grid cells. The corresponding EL simulations contained 2 million, 4 million, and 8 million discrete bubbles, respectively.

Figure 8 presents the results of the strong scaling tests, and table 2 reports the corresponding minimum parallel efficiencies. For each problem size, the ideal scaling trend was defined using the average time per simulation step obtained at the smallest processor and GPU count, respectively. Across both architectures, the computational cost trends are consistent with those observed in the previously presented subgrid cost assessment, where the EE polydisperse model exhibits the highest computational demand, followed by the single EL and EE monodisperse simulations. On CPUs, the solver maintains high parallel efficiency as the problem size increases. Even at the smallest problem size, efficiencies remain above $74\%$ for the EL model and exceed $82\%$ for both EE tests. At 64M cells, the EE monodisperse and EE polydisperse cases reach efficiencies above $90\%$, indicating that the CPU implementation scales effectively and continues to benefit from additional cores. For GPU-based tests, the solver closely follows the ideal scaling behavior up to intermediate configurations, after which the performance trend deviates as the number of GPUs increases, as
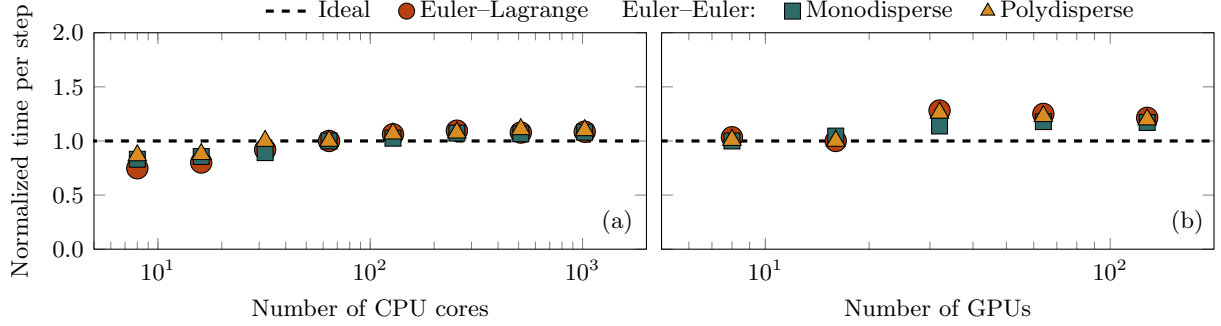
**Figure 8:** Strong scaling on (a) AMD CPUs and (b) NVIDIA A100 GPUs for different problem sizes: (i) 16M, (ii) 32M, (iii) 64M grid cells.

| | (a) AMD Milan CPU cores | | | (b) NVIDIA A100 GPUs | | |
|---|---|---|---|---|---|---|
| Grid Cells | 16M | 32M | 64M | 16M | 32M | 64M |
| EL | 74.21 % | 78.64 % | 85.31 % | 29.69 % | 42.70 % | 57.68 % |
| EE (monodisperse) | 82.76 % | 84.12 % | 89.98 % | 33.26 % | 39.35 % | 47.73 % |
| EE (polydisperse) | 92.43 % | 89.19 % | 94.02 % | 38.13 % | 46.52 % | 52.77 % |

**Table 2:** Minimum parallel efficiency of the strong scaling tests on CPU and GPU architectures. The grid cells are given in millions.

shown in fig. 8 (b). Efficiency is lowest for the smallest problem and recovers as the problem size increases. For the 64M problem size, efficiencies rise to 58 % for EL and exceed 47 % for the EE cases. This upward trend suggests that the GPU solver achieves better scaling once each device has a sufficiently large computational workload. The observed drop in efficiency for configurations with more than 16 GPUs is likely due to each GPU being underutilized, combined with the increasing cost of MPI communication.

These results demonstrate that the phase-averaged multiscale solver scales well on CPUs and efficiently utilizes the available computational resources. They also suggest that GPUs are used most effectively when the hardware is fully occupied.

**Figure 9:** Weak scaling on (a) AMD CPUs and (b) NVIDIA A100 GPUs for Euler–Lagrange (EL) and Euler–Euler (EE) models as labeled.

### 5.3.2. Weak scaling

In weak scaling, the problem size assigned to each processor remains constant while the total number of compute devices increases. Within the EE framework, the number of bins in the polydisperse tests is held constant, and the problem size is therefore only determined by the number of grid cells in the domain. In contrast, for the EL model framework, maintaining a consistent problem size per processor requires fixing both the number of Eulerian grid cells, representing the continuous liquid phase, and the number of Lagrangian bubbles, representing the dispersed phase, assigned to each CPU or GPU. To evaluate weak scaling behavior, we define a representative problem size in which each processor handles 100K 3D grid cells, and, in the EL model, an additional 20K discrete bubbles.

Figure 9 presents the weak scaling performance of the solver. The vertical axis indicates the average time per time step, normalized by the baseline cases using 64 CPU cores and 8 GPUs, respectively. In an ideal scenario, a perfectly scalable solver would align with the dashed reference line. The CPU-based simulations, depicted in fig. 9 (a), exhibit excellent weak scaling, maintaining minimum parallel efficiencies of 91.31 %, 92.51 %, and 90.29 % for the EL, EE monodisperse, and EE polydisperse models, respectively, for cases with 64 or more cores. Each AMD socket comprises 64 cores organized into 4 NUMA (Non-Uniform Memory Access) domains. Simulations executed with 8 or 16 cores are confined to a single NUMA domain, enabling low-latency access to local memory and improved performance. When 32 cores are used, additional cross-NUMA communication introduces latency, and the 64 core case saturates the intra-socket communication bandwidth. Beyond this point, performance stabilizes as all cores within the socket are fully utilized. Conversely, the GPU-based simulations, shown in fig. 9 (b), achieve minimum parallel efficiencies of 78.05 %, 84.81 %, and 79.40 % for the EL, EE monodisperse, and EE polydisperse models, respectively. Performance improves up to 16 GPUs, after which a gradual degradation is observed. From 32 to 128 GPUs, the performance plateaus, likely due to inter-node communication reaching a saturation point. The results confirm that the solver maintains strong parallel efficiency and scalability across both CPU and GPU architectures as computational resources increase, while keeping a constant problem size per compute device.

## 6. Summary

We present a hardware-accelerated phase-averaged multiscale solver for acoustically driven bubbly flows. The solver integrates both volume-averaged (EL) and ensemble-averaged (EE) models to capture the effects of the bubbles within a liquid host. Its limitations are determined by the justified

assumptions underlying each model. In particular, we primarily assume bubble–bubble interaction occurs only through its effect on the liquid–bubble mixture, and the bubble radial oscillations obey the Keller–Miksis formulation. The solver's physical accuracy was verified through three test cases. The oscillation of a single bubble was compared with the analytical Keller–Miksis solution, achieving an RMSE of 2.76 %. We also considered the simulated spherical collapse of a gas bubble in a water-glycerine mixture which closely matched the experimental data from Ohl et al. [47] with an RMSE of 7.46 %, demonstrating the solver's ability to capture highly transient, nonlinear bubble dynamics. Additionally, both CPU and GPU implementations produced identical results, confirming numerical consistency across architectures. The EE model's ability to reproduce ensemble-scale pressure dynamics was verified by comparing its results with averaged data from 40 EL simulations. For monodisperse and polydisperse bubble distributions, the RMSE between the EE and EL results was 2.10 % and 1.53 %, respectively, demonstrating that the EE formulation accurately represents the averaged behavior in statistically homogeneous bubbly systems.

The computational performance of both models was analyzed across varying void fractions. The EL model's cost scales with the number of discrete bubbles, while the EE model's cost scales with the number of bubble size bins. On CPUs, the EL model exhibited sharp cost increases at higher void fractions, resulting from serial bubble updates per core, which can lead to potential load imbalance. In contrast, GPU-based EL simulations achieved up to 3.2 times speedups compared to 64-core AMD Milan CPU runs, due to the parallelism available via the NVIDIA A100 for bubble computations. The EE model proved even more computationally efficient, as a single simulation can represent ensemble dynamics.

On GPUs, the EE formulation achieved speedups of 3.1x (1 bin), 12x (11 bins), and 16x (21 bins) relative to 64-core CPU runs. Although the 21-bin EE configuration was the most computationally demanding among the EE cases, it remained substantially cheaper than performing the 40 EL realizations required for ensemble convergence. The computational cost of the EE model was largely insensitive to variations in void fraction, indicating robust scalability with respect to dispersed-phase volume fraction. Memory usage, however, became a key consideration, particularly for the polydisperse EE model, where the cost scales with both the number of bins and Eulerian grid resolution. For large-scale problems, the EE model's memory footprint can exceed that of the EL model despite shorter runtimes, emphasizing the need to balance accuracy, resolution, and hardware constraints when selecting between models.

Overall, GPU acceleration proved highly effective for both subgrid bubble models, reducing runtimes by more than an order of magnitude and enabling simulations that would otherwise be computationally prohibitive on CPU-only systems. The demonstrated accuracy, scalability, and efficiency establish the solver as a robust tool for large-scale, high-fidelity simulations of multiphase systems relevant to biomedical ultrasound, underwater acoustics, and cavitating flows.

## References

[1] E. Neppiras, Acoustic cavitation, Physics Reports **61** (1980) 159–251.

[2] K. Yasui, Acoustic cavitation, in: Acoustic Cavitation and Bubble Dynamics, Springer, 2017, pp. 1–35.

[3] E. B. Flint, K. S. Suslick, The temperature of cavitation, Science **253** (1991) 1397–1399.

[4] W. Lauterborn, T. Kurz, Physics of bubble oscillations, Reports on Progress in Physics **73** (2010) 106501.

[5] S. Beig, B. Aboulhasanzadeh, E. Johnsen, Temperatures produced by inertially collapsing bubbles near rigid surfaces, Journal of Fluid Mechanics **852** (2018) 105–125.

[6] T. Ikeda, S. Yoshizawa, N. Koizumi, M. Mitsuishi, Y. Matsumoto, Focused ultrasound and lithotripsy, Therapeutic Ultrasound (2016) 113–129.

[7] K. Kooiman, S. Roovers, S. A. Langeveld, R. T. Kleven, H. Dewitte, M. A. O'Reilly, J.-M. Escoffre, A. Bouakaz, M. D. Verweij, K. Hynynen, et al., Ultrasound-responsive cavitation nuclei for therapy and drug delivery, Ultrasound in Medicine & Biology **46** (2020) 1296–1325.

[8] C. C. Coussios, R. A. Roy, Applications of acoustics and cavitation to noninvasive therapy and drug delivery, Annual Review of Fluid Mechanics **40** (2008) 395–420.

[9] J. H. Bang, K. S. Suslick, Applications of ultrasound to the synthesis of nanostructured materials, Advanced Materials **22** (2010) 1039–1059.

[10] J. Luo, Z. Fang, R. L. Smith Jr, Ultrasound-enhanced conversion of biomass to biofuels, Progress in Energy and Combustion Science **41** (2014) 56–93.

[11] K. S. Suslick, Y. Didenko, M. M. Fang, T. Hyeon, K. J. Kolbeck, W. B. McNamara III, M. M. Mdleleni, M. Wong, Acoustic cavitation and its chemical consequences, Philosophical Transactions of the Royal Society of London. Series A: Mathematical, Physical and Engineering Sciences **357** (1999) 335–353.

[12] D. Voronin, G. Sankin, V. Teslenko, R. Mettin, W. Lauterborn, Secondary acoustic waves in a polydisperse bubbly medium, Journal of Applied Mechanics and Technical Physics **44** (2003) 17–26.

[13] V. Leroy, A. Strybulevych, J. H. Page, M. G. Scanlon, Sound velocity and attenuation in bubbly gels measured by transmission experiments, The Journal of the Acoustical Society of America **123** (2008) 1931–1940.

[14] Y. Kaneko, T. Maruyama, K. Takegami, T. Watanabe, H. Mitsui, K. Hanajiri, H. Nagawa, Y. Matsumoto, Use of a microbubble agent to increase the effects of high intensity focused ultrasound on liver tissue, European Radiology **15** (2005) 1415–1420.

[15] K. Kajiyama, K. Yoshinaka, S. Takagi, Y. Matsumoto, Micro-bubble enhanced HIFU, Physics Procedia **3** (2010) 305–314.

[16] D. J. Chung, S. H. Cho, J. M. Lee, S.-T. Hahn, Effect of microbubble contrast agent during high intensity focused ultrasound ablation on rabbit liver in vivo, European Journal of Radiology **81** (2012) e519–e523.

[17] E. K. Juang, L. H. De Koninck, K. S. Vuong, A. Gnanaskandan, C.-T. Hsiao, M. A. Averkiou, Controlled hyperthermia with high-intensity focused ultrasound and ultrasound contrast agent microbubbles in porcine liver, Ultrasound in Medicine & Biology **49** (2023) 1852–1860.

[18] J. Canselier, H. Delmas, A. Wilhelm, B. Abismail, Ultrasound emulsification—An overview, Journal of Dispersion Science and Technology **23** (2002) 333–349.

[19] O. Krasulya, V. Bogush, V. Trishina, I. Potoroko, S. Khmelev, P. Sivashanmugam, S. Anandan, Impact of acoustic cavitation on food emulsions, Ultrasonics Sonochemistry **30** (2016) 98–102.

[20] M. Dular, T. Griessler-Bulc, I. Gutierrez-Aguirre, E. Heath, T. Kosjek, A. K. Klemenčič, M. Oder, M. Petkovšek, N. Rački, M. Ravnikar, et al., Use of hydrodynamic cavitation in (waste) water treatment, Ultrasonics Sonochemistry **29** (2016) 577–588.

[21] D. Zhang, A. Prosperetti, Ensemble phase-averaged equations for bubbly flows, Physics of Fluids **6** (1994) 2956–2970.

[22] K. Ando, T. Colonius, C. E. Brennen, Numerical simulation of shock propagation in a polydisperse bubbly liquid, International Journal of Multiphase Flow **37** (2011) 596–608.

[23] K. Maeda, T. Colonius, Eulerian–Lagrangian method for simulation of cloud cavitation, Journal of Computational Physics **371** (2018) 994–1017.

[24] A. Gnanaskandan, C.-T. Hsiao, G. Chahine, Modeling of microbubble-enhanced high-intensity focused ultrasound, Ultrasound in Medicine & Biology **45** (2019) 1743–1761.

[25] S. H. Bryngelson, K. Schmidmayer, T. Colonius, A quantitative comparison of phase-averaged models for bubbly, cavitating flows, International Journal of Multiphase Flow **115** (2019) 137–143.

[26] M. Snir, MPI–The Complete Reference: The MPI core, volume 1, MIT Press, 1998.

[27] D. Böhme, Characterizing load and communication imbalance in parallel applications, volume 23, Forschungszentrum Jülich, 2014.

[28] P. Wang, T. Abel, R. Kaehler, Adaptive mesh fluid simulations on gpu, New Astronomy **15** (2010) 581–589.

[29] F. Salvadore, M. Bernardini, M. Botti, Gpu accelerated flow solver for direct numerical simulation of turbulent flows, Journal of Computational Physics **235** (2013) 129–142.

[30] J. Sweet, D. H. Richter, D. Thain, GPU acceleration of Eulerian–Lagrangian particle-laden turbulent flow simulations, International Journal of Multiphase Flow **99** (2018) 437–445.

[31] A. Radhakrishnan, H. Le Berre, B. Wilfong, J.-S. Spratt, M. Rodriguez Jr., T. Colonius, S. H. Bryngelson, Method for portable, scalable, and performant GPU-accelerated simulation of multiphase compressible flow, Computer Physics Communications **302** (2024) 109238.

[32] F. Piscaglia, F. Ghioldi, GPU acceleration of CFD simulations in OpenFOAM, Aerospace **10** (2023) 792.

[33] D. C. Jespersen, Acceleration of a CFD code with a GPU, Scientific Programming **18** (2010) 193–201.

[34] S. H. Bryngelson, K. Schmidmayer, V. Coralic, J. C. Meng, K. Maeda, T. Colonius, MFC: An open-source high-order multi-component, multi-phase, and multi-scale compressible flow solver, Computer Physics Communications **266** (2021) 107396.

[35] B. Wilfong, H. Le Berre, A. Radhakrishnan, A. Gupta, D. Vaca-Revelo, D. Adam, H. Yu, H. Lee, J. R. Chreim, M. Carcana Barbosa, Y. Zhang, E. Cisneros-Garibay, A. Gnanaskandan, M. Rodriguez Jr., R. D. Budiardja, S. Abbott, T. Colonius, S. H. Bryngelson, MFC 5.0: An exascale many-physics flow solver, arXiv preprint arXiv:2503.07953 (2025).

[36] A. Preston, T. Colonius, C. Brennen, A reduced-order model of diffusive effects on the dynamics of bubbles, Physics of Fluids **19** (2007).

[37] K. Maeda, T. Colonius, A source term approach for generation of one-way acoustic waves in the Euler and Navier–Stokes equations, Wave Motion **75** (2017) 36–49.

[38] S. H. Bryngelson, Fast integration method for averaging polydisperse bubble population dynamics, Computers & Fluids **304** (2026) 106877.

[39] A. Sinha, S. H. Bryngelson, Neural networks can be FLOP-efficient integrators of 1D oscillatory integrands, Transactions on Machine Learning Research (2024).

[40] S. H. Bryngelson, R. O. Fox, T. Colonius, Conditional moment methods for polydisperse cavitating flows, Journal of Computational Physics **477** (2023) 111917.

[41] A. Charalampopoulos, S. H. Bryngelson, T. Colonius, T. P. Sapsis, Hybrid quadrature moment method for accurate and stable representation of non-Gaussian processes and their dynamics, Philosophical Transactions of the Royal Society A **380** (2022).

[42] T. Colonius, R. Hagmeijer, K. Ando, C. E. Brennen, Statistical equilibrium of bubble oscillations in dilute bubbly flows, Physics of Fluids **20** (2008).

[43] S. Gottlieb, C.-W. Shu, Total variation diminishing Runge-Kutta schemes, Mathematics of Computation **67** (1998) 73–85.

[44] G. Strang, On the construction and comparison of difference schemes, SIAM Journal on Numerical Analysis **5** (1968) 506–517.

[45] S. Chandrasekaran, G. Juckeland, OpenACC for programmers: Concepts and strategies, Addison-Wesley Professional, 2017.

[46] K. W. Thompson, Time dependent boundary conditions for hyperbolic systems, Journal of Computational Physics **68** (1987) 1–24.

[47] C.-D. Ohl, T. Kurz, R. Geisler, O. Lindau, W. Lauterborn, Bubble dynamics, shock waves and sonoluminescence, Philosophical Transactions of the Royal Society of London. Series A: Mathematical, Physical and Engineering Sciences **357** (1999) 269–294.

[48] NCSA, Delta Architecture User Guide, 2023. URL: https://docs.ncsa.illinois.edu.