



HAL
open science

CGCompiler: Automated Coarse-Grained Molecule Parametrization via Noise-Resistant Mixed-Variable Optimization

Kai Steffen Stroh, Paulo C T Souza, Luca Monticelli, Herre Jelger Risselada

► **To cite this version:**

Kai Steffen Stroh, Paulo C T Souza, Luca Monticelli, Herre Jelger Risselada. CGCompiler: Automated Coarse-Grained Molecule Parametrization via Noise-Resistant Mixed-Variable Optimization. *Journal of Chemical Theory and Computation*, In press, 10.1021/acs.jctc.3c00637 . hal-04296767

HAL Id: hal-04296767

<https://hal.science/hal-04296767v1>

Submitted on 20 Nov 2023

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

CGCompiler: Automated coarse-grained molecule parameterization via noise-resistant mixed-variable optimization

Kai Steffen Stroh,^{†,‡} Paulo Cesar Telles de Souza,[¶] Luca Monticelli,[¶] and Herre Jelger Risselada^{*,†,‡,§}

[†]*Department of Physics, Technische Universität Dortmund, Dortmund, Germany*

[‡]*Institute for Theoretical Physics, Georg-August University Göttingen, Göttingen, Germany*

[¶]*Molecular Microbiology and Structural Biochemistry (MMSB, UMR 5086), CNRS & University of Lyon, Lyon, France*

[§]*Leiden Institute of Chemistry, Leiden University, Leiden, The Netherlands*

E-mail: jelger.risselada@tu-dortmund.de

Abstract

Coarse-grained force-fields (CG FF) such as the Martini model entail a predefined, fixed set of Lennard-Jones parameters (building blocks) to model virtually all possible non-bonded interactions between chemically relevant molecules. Owing to its universality and transferability, the building block coarse-grained approach has gained a tremendous popularity over the last decade. The parameterization of molecules can be highly complex and often involves the selection and fine tuning of a large number of parameters (e.g., bead types and bond lengths) to optimally match multiple relevant targets simultaneously. The parameterization of a molecule within the building block CG approach is a mixed-variable optimization problem: The non-bonded interactions

are discrete variables whereas the bonded interactions are continuous variables. Here, we pioneer the utility of mixed-variable particle swarm optimization in automatically parameterizing molecules within the Martini 3 coarse-grained force-field by matching both structural (e.g., RDFs) as well as thermodynamic data (phase-transition temperatures). For sake of demonstration, we parameterize the linker of the lipid sphingomyelin. The important advantage of our approach is that both bonded- and non-bonded interactions are simultaneously optimized while conserving the search efficiency of vector guided particle swarm optimization (PSO) methods over other metaheuristic search methods such as genetic algorithms. In addition, we explore noise-mitigation strategies in matching the phase transition temperatures of lipid membranes, where nucleation and concomitant hysteresis introduces a dominant noise term within the objective function. We propose that noise-resistant mixed-variable PSO methods can both improve as well as automate parameterization of molecules within building block CG FFs, such as Martini.

1 Introduction

Atomically detailed molecular dynamics (MD) simulations provide great insights into the structure and dynamics of biomolecular and other soft matter systems, but larger time- and length scales often require a coarse-grained (CG) description. In coarse-graining a group of atoms is mapped into one bead or supra-atom. Coarse-grained descriptions achieve computational efficiency by reducing degrees of freedom while preserving relevant aspects. This not only allows for bridging larger time and length scales but also enhances our understanding of the fundamental physics underlying molecular processes within biological cells. For example, it can enable fundamental insights into phenomena like the self-organization of lipid membranes and the formation of characteristic thermodynamic phases, including liquid-ordered, liquid-disordered, and gel phases.¹⁻³ Systematic coarse-graining approaches such as inverse Boltzmann and inverse Monte-Carlo approaches^{4,5} as well as force-matching approaches^{6,7} parameterize coarse-grained force-fields by reproducing the structural part of the partition function of the fine-grained system by either matching relevant radial distribution functions or (combined) forces within the fine-grained system. However, because the partition function only describes a single thermodynamic state point at equilibrium, i.e., a unique combination of pressure & temperature values, systematically parameterized 'bottom-up' coarse-grained force-fields are not suited to describe phase transitions over a wider temperature range. Phase-transitions or phase-diagrams can, however, be optimally modeled using coarse-grained force-fields based on the alternative Statistical Associating Fluid Theory (SAFT) parameterization approach, which uses a scaled Lennard-Jones interaction potential whose functional form (the exponent) is uniquely adapted for each interaction type.^{8,9} However, the main practical problem of all of these coarse-grained force-fields is their lack of chemical transferability, i.e. inclusion of a new molecule (interaction type) within the system would require reparameterization of all the existing interaction parameters.

The Martini coarse-grained force-field^{10,11} is a building block force-field, i.e., common chemical groups are parameterized as basic building blocks, which can be combined to build

up any existing molecule. These basic building blocks of Martini, the beads, are parameterized top-down and reproduce the thermodynamic properties of the chemical groups they model, such as partitioning free energies in liquid-liquid systems, while complete molecules are parameterized with a combination of top-down (experimental data) and bottom-up (atomistic simulation). Such a parameterization enables the qualitative simulation of phase transitions as well as phase segregation in lipid membranes while simultaneously conserving molecular compatibility (transferability) by describing all non-bonded interactions with the same 12-6 Lennard-Jones potential form. However, a major drawback compared to other systematic coarse-grained approaches is that parameterization of molecules in Martini can be highly complex and often involves the selection and fine tuning of a large number of parameters (e.g., bead types and bond lengths) to optimally match multiple relevant targets simultaneously. A task that is time consuming when done by human labor. Additionally, it is not always obvious which parameters have to be changed in what manner to enhance a certain behavior, particularly when cooperative processes are involved. While the choice of individual bead types can be made using chemical intuition, still a sizable subset of combined possibilities exists. Importantly, parameterization of bonded and non-bonded parameters should be optimally performed simultaneously since bonded and non-bonded interactions are not independent – they are directly influencing each other via the density of interactions.^{12,13} Recent versions of the Martini force-fields such as Martini 3 rebalanced the density of interactions by introducing an even larger number of possible interaction types, thereby rendering the parameterization of molecules often a non-tractable problem to common users. Automation of coarse-graining is thus critical, especially when constructing large databases of molecules. Automation offers a solution to address the challenge of force-field development, which typically involves collaboration among multiple researchers working on interdependent parameters. By automating the process, a clear and structured flowchart-based hierarchy is established, providing an overview of how the parameterization is conducted and which objectives are targeted.

This automation approach facilitates collaborations by allowing researchers to focus on selecting a set of relevant objectives and assigning importance or weights to each objective. These objectives, along with their individual weights, define the force-field's philosophy. Furthermore, automation empowers collaborations to prioritize two key aspects: the generation and provision of reference data for the objectives at hand, and the design of analysis tools to quantitatively assess how each objective is addressed within the automation pipeline. By automating the parameterization process, collaborators can allocate their efforts towards obtaining high-quality reference data that accurately represents the desired objectives. Simultaneously, they can focus on developing comprehensive analysis tools that enable thorough quantitative evaluation, ensuring the effectiveness of the automation pipeline in achieving the defined objectives. This collaborative approach maximizes the efficiency and reliability of the parameterization process while facilitating a deeper understanding of the force-field's performance.

Particle swarm optimization (PSO) is a powerful computational method used to optimize problems by iteratively improving candidate solutions based on a defined objective function. Compared to evolutionary optimization methods like genetic algorithms, PSO offers advantages in efficiently finding global optima within high-dimensional continuous spaces due to its vectorial search direction. PSO has been successfully employed in various coarse-grained (CG) parameterization tasks, as demonstrated in previous studies.^{14–18}

PSO is primarily designed for continuous variables, making it well-suited for optimizing structure-based coarse-grained (CG) models where bonded and non-bonded parameters can be chosen from a continuum of values. However, in building block models like Martini, the non-bonded parameters are predefined and discrete, representing different interaction levels. Consequently, the parameterization of molecules in a building block CG force field becomes a mixed-variable optimization problem.

When using PSO for parameterization in building block models, a transformation from the continuous space to the discrete space of force field parameters is necessary. This trans-

formation introduces cumulative rounding errors, which can potentially affect the quality of the parameterization, especially in larger molecules. Therefore, additional evaluation and reparameterization steps are often required to ensure the optimal performance of the force field.

It is crucial to parameterize both bonded and non-bonded interactions simultaneously since they are not independent and their optimization should be performed in a coordinated manner.¹³ By considering their interplay during the parameterization process, the resulting force field can better capture the complex behavior of molecules in the system.

To address the limitations of existing PSO approaches, we employ a mixed-variable PSO scheme (mv-PSO) for parameterization. This approach allows for the simultaneous optimization of both discrete parameters (representing non-bonded interactions) and continuous parameters (representing bonded interactions), enhancing the accuracy and reliability of the parameterization process.

Furthermore, due to the chaotic nature of MD simulations, observables measured in MD simulations are subject to noise. Since standard PSO was designed for deterministic objective functions, straightforward application to noisy optimization problems is error prone, because the algorithm can no longer correctly identify global and personal best solutions when noise levels are similar to differences between objective function values.¹⁹ Noise-mitigation strategies are particularly important when utilizing thermodynamic data as targets, as these are notoriously expensive to estimate accurately in MD simulations, even when employing CG models. Particularly problematic is the targeting of phase transition temperatures, which involve a first order phase transition and are thus subject to nucleation and concomitant hysteresis.

In this paper, we pioneer the application of mixed-variable particle swarm optimization in automated parameterization of molecules within the Martini 3 coarse-grained force-field by matching both structural (e.g., RDFs) as well as thermodynamic data (phase-transition temperatures). The important advantage of this approach is that both bonded- and non-

bonded interactions are simultaneously optimized while conserving the search efficiency of vector guided particle swarm methods over other metaheuristic search methods such as genetic algorithms. In addition, we explore noise-mitigation strategies in matching the phase transition temperatures, where nucleation and concomitant hysteresis introduces a dominant noise term within the objective function. To the best of our knowledge, the impact of noisy objective function values has not been previously addressed in the context of applying PSO for CG parameterization. The manuscript is structured in the following way: Section 2 describes the mixed-variable PSO algorithm and parameterization procedure. As an example, we parameterized the linker region of sphingolipids, a biological highly relevant class of lipid molecules, that constitutes approximately 30 mol% of the plasma membrane lipids,²⁰ but has not been updated for Martini 3, yet. Details of the simulated molecules, systems and observables are given in Section 3. Results are presented in Section 4, followed by conclusions in Section 5.

2 CG molecule parameterization via mixed-variable particle swarm optimization

With CGCompiler we present a Python package that streamlines CG molecule parameterization. It employs mixed-variable particle swarm optimization to simultaneously optimize categorical (beadtype) and continuous (bonds, angles, dihedrals, ...) variables. Therefore, CGCompiler is particularly well suited for, but not limited to, parameterization tasks in CG FFs that follow a building block approach. To enable the application of the building block approach also to larger molecular fragments, consisting of more than one CG bead, the method allows for optimization of shared building blocks in different molecules, e.g. the headgroup, linker, or tails of lipids.

Molecule parameterization in Martini 3 follows three steps: i) Choice of mapping and bead sizes ii) Assignment of chemical bead types iii) Choice of bonded terms and assignment of bonded parameters.¹¹ While a mapping from atomistic to CG model and the set of bonded terms have to be predefined, the here-presented algorithm optimizes bead size, chemical bead type and bonded parameters simultaneously.

The parameterization workflow is shown in Figure 1. For a given parameterization task, the user provides or generates the target data, and creates a set of CG training systems, that allows measurement of the target observables. In the initial iteration, the optimization algorithm generates a number N_p , i.e., the swarm size, of candidate solutions with random FF parameters, and runs MD simulations for each candidate solution and each training system. Candidate solutions are then scored by how well the parameterization targets are reproduced. By utilizing the swarm's knowledge of the fitness landscape, candidate solutions are updated, and a new cycle of MD simulations, analyses, and fitness evaluations, starts. This is repeated until a termination criterion is fulfilled. Due to noise in the objective function evaluation, the selection of the true best parameters can only be done with a certain probability. Therefore, the set of the best, statistically equal candidate solutions undergoes a screen-to-the-best

procedure, which either provides one solution that is significantly better than the rest, or reduces the field of viable candidate solutions further, on which more expensive evaluation simulations would be performed.

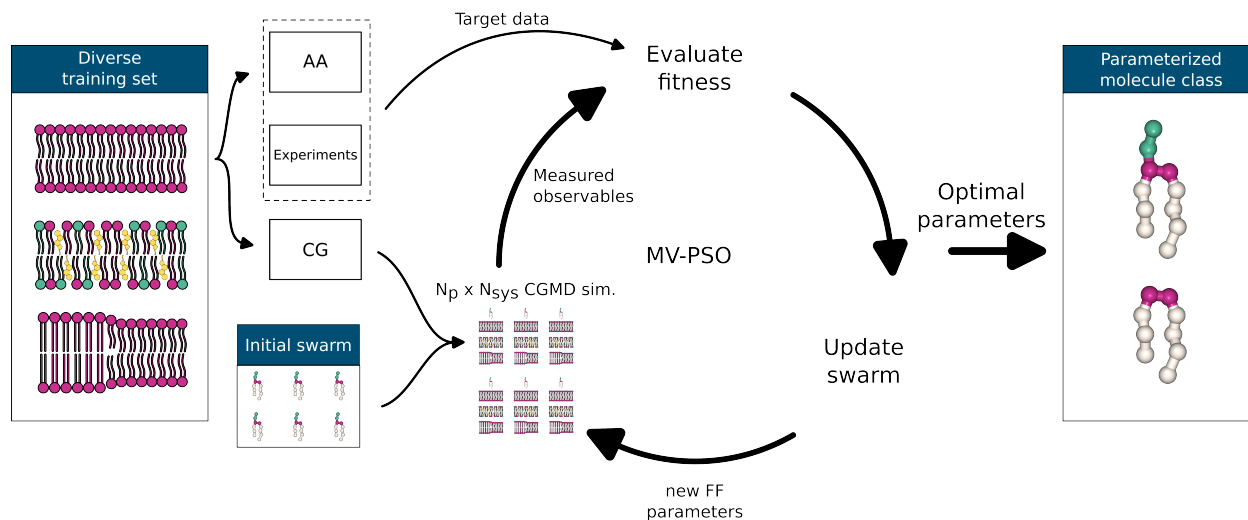


Figure 1: Parameterization workflow. i) A set of training systems from which the target properties can be extracted. ii) Target data is acquired from atomistic simulations and experiments. iii) An initial swarm is generated with FF parameters randomly selected from a predefined range of feasible parameters. iv) All candidate solutions are simulated in all training systems, the target observables are measured and compared to the target data, i.e., the fitness of the candidate solutions is estimated. New candidate solutions are generated by utilizing the swarm’s knowledge of the fitness landscape. v) Step iv is repeated until a termination criterion is fulfilled. vi) A screen-to-the best procedure yields the optimized set of FF parameters.

2.1 Mixed-variable particle swarm optimization

In the original PSO algorithm for continuous optimization problems in a D -dimensional parameter space, particle i has a position vector $X_i = (x_i^1, \dots, x_i^D)$ and a velocity $V_i =$

(v_i^1, \dots, v_i^D) .²¹ At each iteration t the velocity and position are updated by

$$\begin{aligned} V_i(t+1) &= w * V_i(t) \\ &+ c_1 \mathbf{r}_1 (\text{pbest}_i(t) - X_i(t)) \\ &+ c_2 \mathbf{r}_2 (\text{gbest}(t) - X_i(t)) \end{aligned} \quad (1)$$

$$X_i(t+1) = X_i(t) + V_i(t+1) \quad (2)$$

Where $\text{pbest}_i(t)$ is the personal best position of particle i and $\text{gbest}(t)$ is the best position found by the whole swarm. w is an inertia weight, which balances global vs. local search. The coefficients c_1 and c_2 are balancing personal vs. social experience. \mathbf{r}_1 and \mathbf{r}_2 are vectors of random numbers. In the mv-PSO algorithm, that is utilized in our work, the position vector of a particle takes a hybrid form, where Z dimensions encode continuous variables and V dimensions encode categorical variables.²²

$$X_i = \underbrace{(x_i^1, x_i^2, \dots, x_i^Z)}_{\text{continuous}}, \underbrace{(x_i^{Z+1}, x_i^{Z+2}, \dots, x_i^{Z+V})}_{\text{categorical}} \quad (3)$$

The continuous and categorical parts of the position vector are updated separately.

2.1.1 Continuous reproduction method

In classical PSO the swarm can get trapped in local optima and therefore prematurely converge.²² To promote diversity while maintaining good convergence efficiency Wang et al. proposed an altered continuous reproduction scheme, where particle i learns from the best position of a randomly selected particle.²² In order to guide the swarm towards improved solutions, the pool of pbest to choose from, only consists of solutions whose fitness is superior to $\text{pbest}_i(t)$.

$$V_i(t+1) = w \cdot V_i(t) + c \cdot \mathbf{r} \cdot (\text{pbest}_r(t) - X_i(t)) \quad (4)$$

Algorithm 1 Continuous reproduction method

- 1: **Input:** sorted swarm, particle i , parameter w_i
- 2: **for** $j = 1..Z$ **do**
- 3: Randomly choose r , $i \leq r \leq N$
- 4: $v_i^j(t+1) = w_i \cdot v_i^j(t) + c \cdot r \cdot (pbest_r^j - x_i^j)$
- 5: $x_i^j(t+1) = x_i^j(t) + v_i^j(t+1)$
- 6: **end for**
- 7: **return** $(x_i^1, x_i^2, \dots, x_i^Z)$

2.1.2 Categorical reproduction method

Values of categorical variables are assigned according to a probability. Initial probabilities are given by

$$Prob_{j,n}(0) = \frac{1}{n_j} \quad (5)$$

where n_j is the number of available values for the j th variable. To leverage the swarm's knowledge of good solutions, only the superior half of the sorted swarm is utilized in updating the probabilities of available categorical values. To avoid premature extinction of available values, a lower limit is assigned for $Prob_{j,n}$. If $Prob_{j,n}$ falls below that lower limit, $Prob_{j,n}$ is set to that threshold value, and all probabilities are renormalized such that $\sum_n Prob_{j,n} = 1$. The categorical update method is shown in Algorithm 2.

Algorithm 2 Categorical reproduction method

- 1: **Input:** sorted swarm, particle i , parameter α_i
- 2: **for** $j = 1..V$ **do**
- 3: **for** each available value n , $n = 1$ to n_j **do** $Count_{j,n} = 0$
- 4: **for** each personal best $pbest_i$, $i = N/2$ to N **do**
- 5: **if** $pbest_{i,j} == Values_{j,n}$ **then**
- 6: $Count_{j,n} + = 1$
- 7: **end if**
- 8: **end for**
- 9: $Prob_{j,n}(t+1) = \alpha_i \cdot Prob_{j,n}(t) + (1 - \alpha_i) \cdot \frac{Count_{j,n}}{N/2}$
- 10: **end for**
- 11: **end for**
- 12: **for** $j = 1..V$ **do**
- 13: Assign an available value to x_i^{Z+j} according $Prob_j$
- 14: **end for**
- 15: **return** $(x_i^{Z+1}, x_i^{Z+2}, \dots, x_i^{Z+V})$

2.1.3 Cost function

Molecule parameterization is typically a multiobjective optimization problem (MOP). A simple way to scalarize an MOP is by linear weighting. The scalarized optimization problem is solved by minimizing the cost, which is given by

$$\text{cost} = \sum_o w_o f_o(\mathbf{x}) \quad (6)$$

Where w_o is an objective weight, f_o the objective cost function, and \mathbf{x} the parameter vector. The objective weights can be used to balance the importance of the utilized parameterization targets. The weights are set by the user. Setting weights might require some intuition about the parameterized molecule, quality of target data, etc.

Each objective can have a different objective cost function f_o . New objective cost functions can be added by the user easily. In its present form, the parameterization algorithm uses two distinct objective cost functions. For *single valued observables*, such as area per lipid, membrane thickness, melting temperature, solvent accessible surface area (SASA) the objective cost function is defined as

$$f_o(\mathbf{x}) = \frac{1}{\sum_s^{N_s} w_{o,s}} \left(\sum_s^{N_s} w_{o,s} \frac{1}{N_{\text{types},s}} \sum_t^{N_{\text{types},s}} \max(0, SAE(y_{s,t}(\mathbf{x}), \hat{y}_{s,t}) - E_{o,s}^{\text{tol}}) \right). \quad (7)$$

$y_s(\mathbf{x})$ is the observed value, given the FF parameters \mathbf{x} . \hat{y}_s is the target value. N_s is the number of training systems that is used for the current parameterization objective. N_{types} is the number of bond or angle types being parameterized. The deviation from the target is calculated by the scaled absolute error $SAE(y, \hat{y}) = \left| \frac{\hat{y} - y}{\hat{y}} \right|$. With the error tolerance $E_{o,s}^{\text{tol}}$, uncertainties in target data can be accounted for. Each training system has an additional weight $w_{o,s}$, which can be used in case of differences in target data quality or similar cases. Generally these are set to 1.

For observables that are given in the form of *distributions*, such as bond lengths, angles,

or radial distribution functions (RDF), the objective cost function is given by:

$$f_o(\mathbf{x}) = \frac{1}{\sum_s w_{o,s}} \left(\sum_s w_{o,s} \frac{1}{N_{\text{types},s}} \sum_t^{N_{\text{types},s}} EMD(\phi(\mathbf{x}_{s,t}), \hat{\phi}_{s,t}) \right) \quad (8)$$

Where $\phi(\mathbf{x})$ is the observed distribution, given the FF parameters \mathbf{x} . $\hat{\phi}$ is the target distribution. The earth mover's distance $EMD(\phi(\mathbf{x}_{s,t}), \hat{\phi}_{s,t})$ is a measure of the distance between the two distributions.²³

2.2 Noise mitigation strategies for PSO

PSO was designed for deterministic objective functions. Due to the chaotic nature of MD simulations hereby measured observables are subject to noise. With noise in objective functions, selection of the true best solutions is not guaranteed. Since solutions, that are identified as the best, attract the swarm toward regions of interest in parameter space, noise can misguide the swarm and therefore deteriorate PSO performance.

2.2.1 Resampling

Resampling is a widely applied strategy for noise mitigation within the objective function. Relatively simple resampling methods are *equal resampling* (PSO-ER), *extended equal resampling* PSO-EER, and *equal resampling* with allocation to top-N solutions PSO-ERN.²⁴ These simpler methods are regularly outperformed by state-of-the-art resampling methods, such as *optimal computing budget allocation* PSO-OCBA,²⁵ but the quality of results depends on the specific optimization problem and noise levels.^{19,24} OCBA aims to maximize the probability of correctly selecting good solutions. This is done by first allocating a primary computational budget equally to all current solutions to estimate their cost means and variances. A secondary budget is then sequentially allocated to solutions with lower means and higher variances to improve the fitness estimations of potentially good solutions. For efficient secondary budget allocation at least 5 primary evaluations should be executed for mean and

variance estimation.²⁶ This might make application of OCBA prohibitively expensive for regular CG molecule parameterization tasks. Based on the observation that most observables utilized in the multiobjective optimization of the sphingomyelin linker region have a low variance and only a few suffer from a larger variance (cf. Figure S5), we hypothesize that in the molecule parameterization task at hand, one primary objective function evaluation is sufficient to differentiate potentially good solutions from bad solutions, but to maximize the probability of correctly selecting the true best solution, the accuracy of the fitness estimates has to be increased. Therefore, we propose a somewhat pragmatic approach, that salvages the core idea of OCBA, i.e., allocate additional computational budgets to where it is the most useful (low mean and high variance). At each iteration, our resampling method involves one full objective function evaluation of the current solutions. The current solutions are then ranked by their fitness, and for the best N solutions only the observables that have significant variance are reevaluated.

2.2.2 Set of statistically equivalent solutions

Even with noise mitigation, at the end of an optimization run, there will be a number of solutions with very similar scores. While in a deterministic setting, the global best position is determined by

$$\mathbf{gbest} = \arg \min_{\mathbf{x} \in \mathcal{P}_t} f(\mathbf{x}), \quad (9)$$

where \mathcal{P}_t is the set of all positions that have been visited by the swarm up to iteration t , with noise in the objective function no solution can be declared the best with 100% certainty.¹⁹ With the *screen-to-the-best* procedure of Boesel et al.²⁷ a set of positions $\mathcal{P}_t^g \subseteq \mathcal{P}_t$ can be selected, such that the true global best solution \mathbf{gbest} is contained in \mathcal{P}_t^g with probability of at least $1 - \alpha$ (with $0 < \alpha < 1$).¹⁹

For solutions $i, j \in \mathcal{P}_t$, \bar{f}_i and S_i^2 denote the sample mean and sample variance of objective function values. The elementary steps of the screen-to-the-best procedure are:

1. Compute W_{ij} ,

$$W_{ij} = \left(\frac{t_i S_i^2}{n_i} + \frac{t_j S_j^2}{n_j} \right)^{1/2}, \forall i \neq j \in \mathcal{P}_t \quad (10)$$

where $t_i = t_{(1-\alpha)^{1/|\mathcal{P}_t|-1}, n_i-1}$ and $t_{\beta, \nu}$ is the β quantile of the t distribution with ν degrees of freedom

2. Set $\mathcal{P}_t^g = \{i : i \in \mathcal{P}_t, \bar{f}_i \leq \bar{f}_j + W_{ij}, \forall i \neq j \in \mathcal{P}_t\}$

3. Return \mathcal{P}_t^g

W_{ij} is the half-width of pooled t-confidence intervals on the difference between the scores of solutions i and j .¹⁹ Therefore, the procedure entails a pair-wise comparison of solutions and determines if differences of the sample averaged scores are statistically significant.¹⁹

3 Example application: Sphingolipid linker parameterization

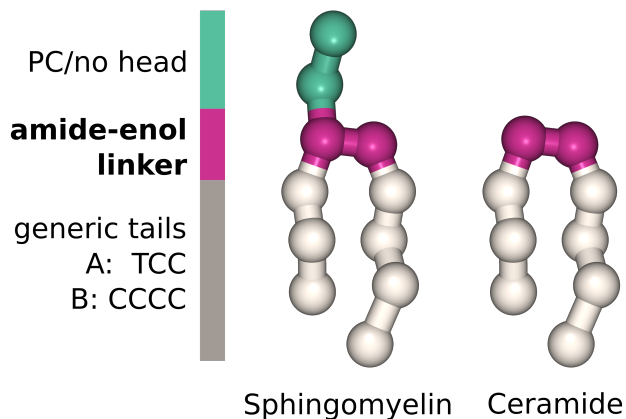


Figure 2: CG description of sphingomyelin and ceramide.

3.1 Simulation details

The Python package is based on evo-MD.²⁸ All simulations were performed with GROMACS 2020.4 and 2021.4²⁹ and analyzed with in-house Python scripts that are utilizing MDAnalysis,^{30,31} LiPyphilic,³² SciPy,³³ and pyemd, which is a Python wrapper for Pele and Werman’s EMD implementation.^{34,35} Visualization was done with NGLview.³⁶

3.1.1 Atomistic models

All atomistic models were simulated using the CHARMM36^{37–39} force field. Table 1 provides details about the atomistic target systems. Initial configurations of the membrane systems were generated with the CHARMM-GUI membrane builder.^{40–42} Following energy minimization and equilibration, all systems were simulated with a 2 fs time step. Bonds of hydrogen atoms were constrained employing the LINCS algorithm.⁴³ Van der Waals forces were gradually switched off between 1.0 nm and 1.2 nm. The PME algorithm⁴⁴ was used for electrostatic interactions. Temperature coupling was done via the velocity rescale algo-

rithm⁴⁵ with a coupling time $\tau_t = 1.0$ ps. System pressures were held at 1 bar by using the Parinello-Rahman barostat⁴⁶ with a coupling time $\tau_p = 5.0$ ps. Pressure coupling was applied isotropically for aqueous solutions and semi-isotropically for membrane systems.

Table 1: Atomistic target system details. In the naming scheme of the CHARMM FF, SSM and CHL1 denote sphingomyelin (18:0) and cholesterol, respectively.

system	lipids	# TIP3P	# NA	# CL	T / K	sim. time / ns
DPSM128 328K	128 SSM	5120	-	-	328.15	150
POPC SSM CHOL	100 POPC 100 SSM 100 CHL1	9000	18	18	321.15	300

3.1.2 Coarse-grained models

All coarse-grained models were simulated using the Martini 3¹¹ force field. Beta version 14 of the Martini 3 cholesterol parameters was used.^{47,48} Initial configurations of membrane systems were generated with the Python script insane.⁴⁹ Details of the employed training systems are listed in Table 2. All systems were energy minimized and equilibrated with the current version of DPSM, that made the Martini 2 model of sphingomyelin compatible with Martini 3. During the particle swarm optimization each system was equilibrated with the candidate FF parameters in two stages, with time steps of 2 fs and 20 fs, respectively. For all coarse-grained production simulations a time step of 20 fs was used. Non-bonded interactions were cut off at 1.1 nm. For electrostatic interactions the reaction-field method was used with a dielectric constant of 15 and the reaction-field dielectric constant was set to infinity.

Temperature coupling was obtained via the velocity rescale algorithm⁴⁵ with a coupling time $\tau_t = 1.0$ ps. System pressures were held at 1 bar by using the Parinello-Rahman barostat⁴⁶ with a coupling time $\tau_p = 12.0$ ps. Pressure coupling was applied isotropically for aqueous solutions and semi-isotropically for membrane systems. In simulations for melting temperature estimation anisotropic pressure coupling was employed, using the Berendsen barostat⁵⁰ with a coupling time $\tau_p = 4.0$ ps.

Table 2: Coarse-grained training system details.

system	lipids	# W	# NA	# CL	T / K
DPSM128 328K	128 DPSM	1177	-	-	328.15
DPSM256 biphasic	256 DPSM half gel/half liquid	2300	26	26	286, 291, 296, 301, 303, 305, 307, 308, 309, 310, 311, 316, 321, 326
POPC SSM CHOL	96 POPC 96 DPSM 96 CHOL	2124	23	23	321.15

4 Results

Our aim was the development of an automatization framework for molecule parameterization in building-block force fields. As an example we parameterized the sphingolipid linker region. Section 4.1 shows the results of the parameterization with CGCompiler using a simple noise-mitigation strategy. Since noise-mitigation strategies can only reduce the effects of noise when selecting the true best solution, the best statistically equivalent solutions generated during the mv-PSO run are subsequently screened-to-the-best, as described in Section 2.2.2.

4.1 Parameterization of the sphingolipid linker region

Table 3 shows the observables and their weights used in the parameterization. The swarm size was 64. Noise-mitigation was done by reevaluating the melting temperature of the 16 best candidate solutions of the current iteration 12 times, i.e., results were obtained with noise-mitigation setting mv-PSO-R16 (cf. Section 4.2). As T_m is the major contribution to cost variance, but the employed T_m estimation method is good for differentiating good from bad solutions, i.e., it has an accuracy of a few K. Other observables were only evaluated once, area per lipid (APL) fluctuations were the second largest cause of cost variance. For more details on noise-mitigation efficacy see Section 4.2.

All results shown include the complete set of the best statistically equivalent candidate solutions \mathcal{P}^g that remained after two rounds of the screen-to-the-best procedure (cf. Section 2.2.2). This set contains 18 candidate solutions.

Table 3: Weights of observables w_o and system specific observable weights $w_{o,s}$ for optimization run 1.

observable	w_o	$w_{o, \text{DPSM128}}$	$w_{o, \text{DPSM256}}$	$w_{o, \text{POPC SSM CHOL}}$
bond length dist.	1	1	0	1
angle dist.	100	1	0	1
d_{HH}	500	1	0	0.25
APL	1000	1	0	0.25
T_m	250	0	1	0
RDF COM DPSM-CHOL	1	0	0	1

4.1.1 Improved reproduction of membrane properties

Figure 3 shows thickness, average area per lipid and melting temperature of pure DPSM membranes for the set of statistically equal candidate solutions that remained after the second screen-to-the-best procedure performed after reevaluating the initial set 20 times. All new candidate solutions outperform the current DPSM model regarding thickness. The average area per lipid of the current model is closer to the target value, but most of the candidate solutions are within the tolerance of 1.5% deviation. In general, thickness and APL are inversely correlated, increasing one will always result in decreasing the other, therefore, with both values inside the tolerance, the new models represent a better balance of thickness and APL. It is important to note that in the comparison, SM(18:0) was used as the atomistic target. The current tail model of the Martini FF represents both SM(16:0) and SM(18:0). The CHARMM model for SM(16:0) exhibits a reduced thickness when compared to SM(18:0).³⁸ It is therefore not unexpected that the Martini DPSM models show a reduced thickness compared to SM(18:0).

While the melting temperatures estimated with the biphasic approach, that is used during optimization for performance reasons, are not within the specified tolerance regime of 2 K but $\approx 5 - 6$ K below the target value and $\approx 3 - 4$ K below the lower target threshold, the new models are greatly improved compared to the current model, which was 20 K off target. Notably, the estimation of T_m is approach dependent. Estimations using the alternative, reversible melting approach with slow melting rates, based on Kowalik et al.⁵¹ and Sun and Böckmann⁵² (see SI for further details), which requires a very large computational budget (as done here, total simulation time for one T_m estimation $> 90 \mu s$) show an even better agreement with the experimental melting temperature.

The here-performed biphasic approach utilizes a bilayer that is half gel and half liquid. The gel phase is fabricated by quenching to a temperature well below the melting temperature, and the gel phase system is combined with a preequilibrated liquid system. The combined system is then equilibrated with thermostats set to different temperatures for

the two phases. As quenching and equilibration can take up to several hundreds of ns, reconstructing the starting structure for every candidate solution would significantly increase computational cost of a PSO run. Therefore, starting structures for this procedure were generated with the current DPSM parameters beforehand and equilibrated using the parameters of each candidate solution. While equilibration of the fluid phase is generally fast, this certainly is not the case for the gel phase. Considering that an unequilibrated phase is inherently less stable, the presence of an equilibrated liquid phase alongside an unequilibrated gel phase may lead to a slight systematic underestimation of the melting temperature (T_m).⁵³ However, this potential underestimation can be anticipated and taken into account during the analysis.

The equilibrium melting rate approach does not suffer from this potential problem of unequally equilibrated phases. To minimize bias caused by the quenched starting structures used in this approach, for each validated candidate solution eight different starting conformations were generated.

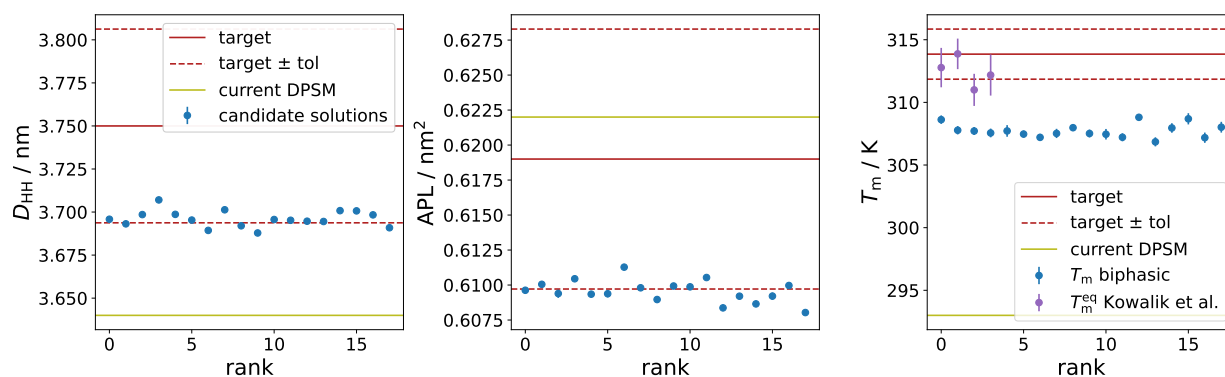


Figure 3: Thickness, average area per lipid and melting temperature for the set of statistically equal candidate solutions that remained after the second screen-to-the-best procedure performed after reevaluating the initial set 20 times.

4.1.2 Structural properties of the parameterized sphingomyelin models

Figure 4 shows the distributions of the newly parameterized bonds and angles for the candidate solutions in \mathcal{P}^g . The atomistic target distributions are matched reasonably well in

all cases. Some finer details of the atomistic model, like double peaks or extensive shoulders cannot be matched in the CG model. The parameterization philosophy of Martini 3 adopts a size-shape concept, where bond lengths are determined based on the molecular volume of the atomistic fragment mapped by the beads, rather than simply center of masses. This complication further underscores the necessity of employing multi-objective optimization algorithms to achieve effective molecule parameterization.

The solvent accessible surface area (SASA) is commonly used to further compare the molecular volumes and shapes between CG and AA models.^{11,54} Figure 5 shows the SASA values of \mathcal{P}^g in comparison to the AA and current CG DPSM models. The SASAs are computed for the linker beads AM1 and AM2, as well as all supra-atoms that are directly connected to the linker, i.e., beads PO4, T1A, and C1B, as these connections are also parameterized. With SASA values of $\approx 6.24 \text{ nm}^2$ all newly parameterized CG models show a better reproduction of the AA value (5.24 nm^2) compared to the current model (6.45 nm^2), but with discrepancy of $\approx 19\%$ all SASA values remain grossly too high. It appears that solely reparameterizing the linker region is not enough to fix this issue. Furthermore, using SASA directly as a target in the high-throughput optimization scheme is not necessarily beneficial, since a specific SASA value is not a unique representation of a certain shape. Therefore, comparisons of solvent accessible surface areas between AA and CG models are most helpful when done by simultaneous visual inspection. For automated parameterization, however, more detailed shape descriptors should be used.

4.1.3 Force field parameters

Non-bonded interactions: Due to the polar nature of the linker region of sphingolipids, only the chemical types of the P-block of the Martini 3 FF were eligible. As groups of 3 or 4 heavy atoms were combined into supra-atoms in the specified mapping, bead sizes small (S) and regular (default) could be chosen by the algorithm. Both bead sizes were permitted for both interaction sites, to allow for some wiggle room, even though 4 heavy atoms are

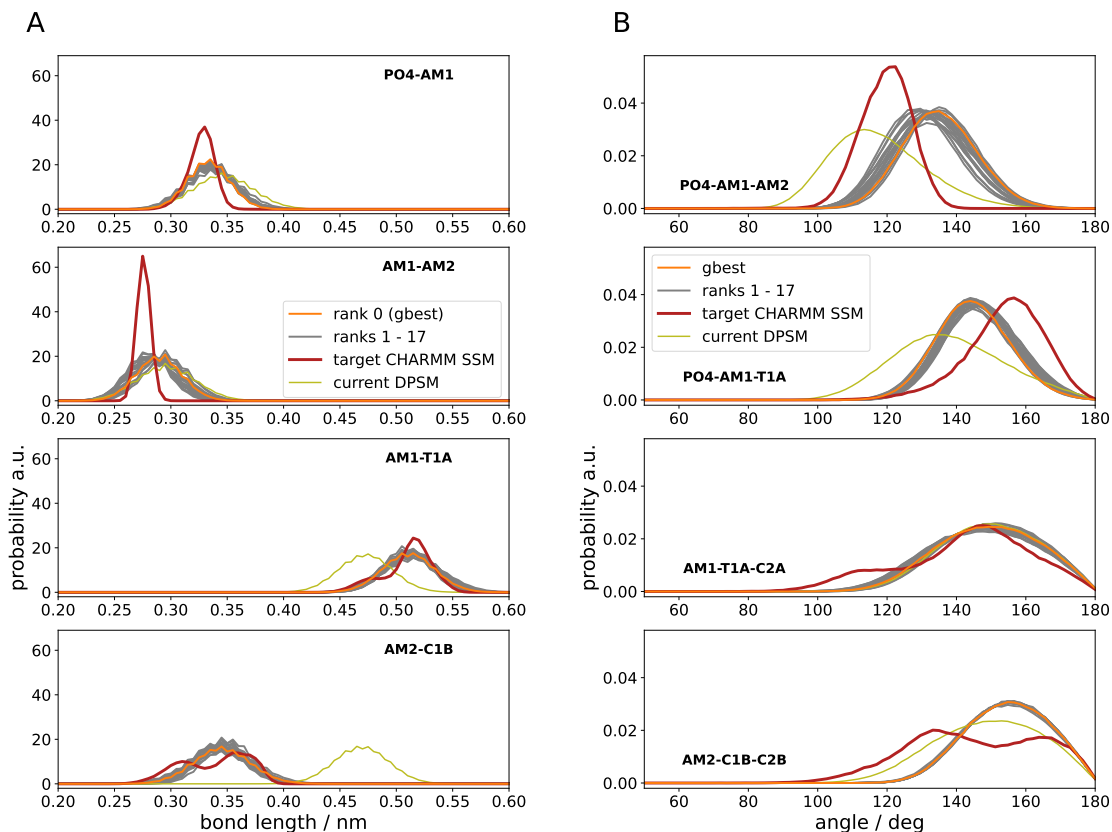


Figure 4: Validation of targets from rerun simulations for the set \mathcal{P}^S . **A)** Bond length distributions. **B)** Angle distributions.

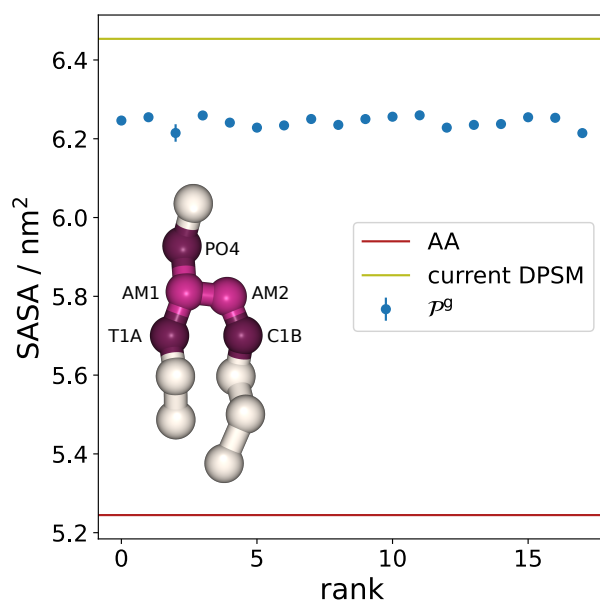


Figure 5: Solvent accessible surface area of the linker and beads connected directly to it. Beads involved in SASA calculation are highlighted.

grouped together into supra-atom AM1 and 3 into AM2. A slight miscount of mapped atoms is not uncommon in Martini, e.g., the mapping of the NC3 bead is actually 6-to-1.¹⁰

One feature of the mixed-variable approach is that the optimization procedure directly yields a probability distribution of bead types, cf. Figure 6A. While for the interaction site AM2 there is clear consensus on the bead type, for AM1 only the size (small) is clearly determined, but there is some ambiguity regarding the interaction strength. The reduced size of one of the beads seems to be warranted, given the still too high SASA values shown above, and is also inline with the new Martini 3 models of glycerolipids.¹¹ It is also worth mentioning that the chemical bead types chosen by our algorithm match the expected assignment suggested by Martini 3.

A converged "degenerate" probability distribution of bead types is the result of two or more bead types having indistinguishable effects on fitness. This can be caused by noise levels being larger than the fitness differences or the employed set of observables and training systems is lacking the necessary discriminatory power. Both issues can be remedied in post-optimization screening, but should optimally be addressed during optimization. As the former option would merely improve selection from the pool of generated candidate solutions, the later would potentially allow the generation of truly better solutions.

Additionally, for both, non-bonded and bonded FF parameters, diversity can be caused by the fact that the objective cost function for single valued observables (Eq. 7) has an error tolerance to accommodate for uncertainties in target data. With respect to these observables, different parameterizations with different "phenotypes" can have the same objective cost, as long as they are within the specified tolerances.

Bonded interactions: Table 4 lists the range of permitted bond parameters used in the optimization. The resulting bonded parameters of \mathcal{P}^g are shown in Figure 6. For equilibrium bond lengths b_0 there is little variation between different candidate solutions. This strong consensus suggests that the optimization has converged and that small changes in equilibrium bond length are linked to significant cost changes. The situation for the force constants is

quite different. The values fluctuate over a relatively large range, compared to the predefined domain of permitted values. The measured bond length distributions (Figure 4A) show that these seemingly substantial differences in force constant values have only minor effects on the molecule’s behavior.

The situation for the angle FF parameters is similar. The equilibrium values show smaller variances than the force constants, compared to their respective domain sizes of applicable values. Again, the differences in FF parameters have little effect on the observed distributions (cf. Figure 4B). Notably, the optimal force constants for the angles PO4-AM1-T1A and AM2-C1B-C2B were close to or at the maximum of their permitted ranges. Further optimization was therefore likely hindered, and a wider range should have been chosen.

In a similar vein to the discussion surrounding non-bonded parameters, the relatively wide range of force constants in \mathcal{P}^g indicates that additional metrics or training systems could be employed to further optimize the overall performance of candidate solutions while maintaining the quality of the employed observables. For instance, exploring lipids in environments other than a bilayer, which induce different lipid conformations, could benefit from a candidate solution with a lower angle force constant to allow for increased conformational variation.

Table 4: Bonded interactions. GROMACS function type; permitted parameter ranges for equilibrium bond length / angle, and corresponding force constants.

bond	GROMACS		
	bond func. type	b_0 / nm	fc / kJ/mol/nm ²
PO4-AM1	1	0.25 – 0.40	1000 – 9000
AM1-AM2	1	0.20 – 0.35	1000 – 9000
AM1-T1A	1	0.40 – 0.55	1000 – 9000
AM2-C1B	1	0.25 – 0.50	1000 – 9000
angle	GROMACS		
	angle func. type	a_0 / deg	fc / kJ/mol
PO4-AM1-AM2	2	90 – 180	5 – 100
PO4-AM1-T1A	2	90 – 180	5 – 100
AM1-T1A-C2A	2	180	5 – 100
AM2-C1B-C2B	2	180	5 – 100

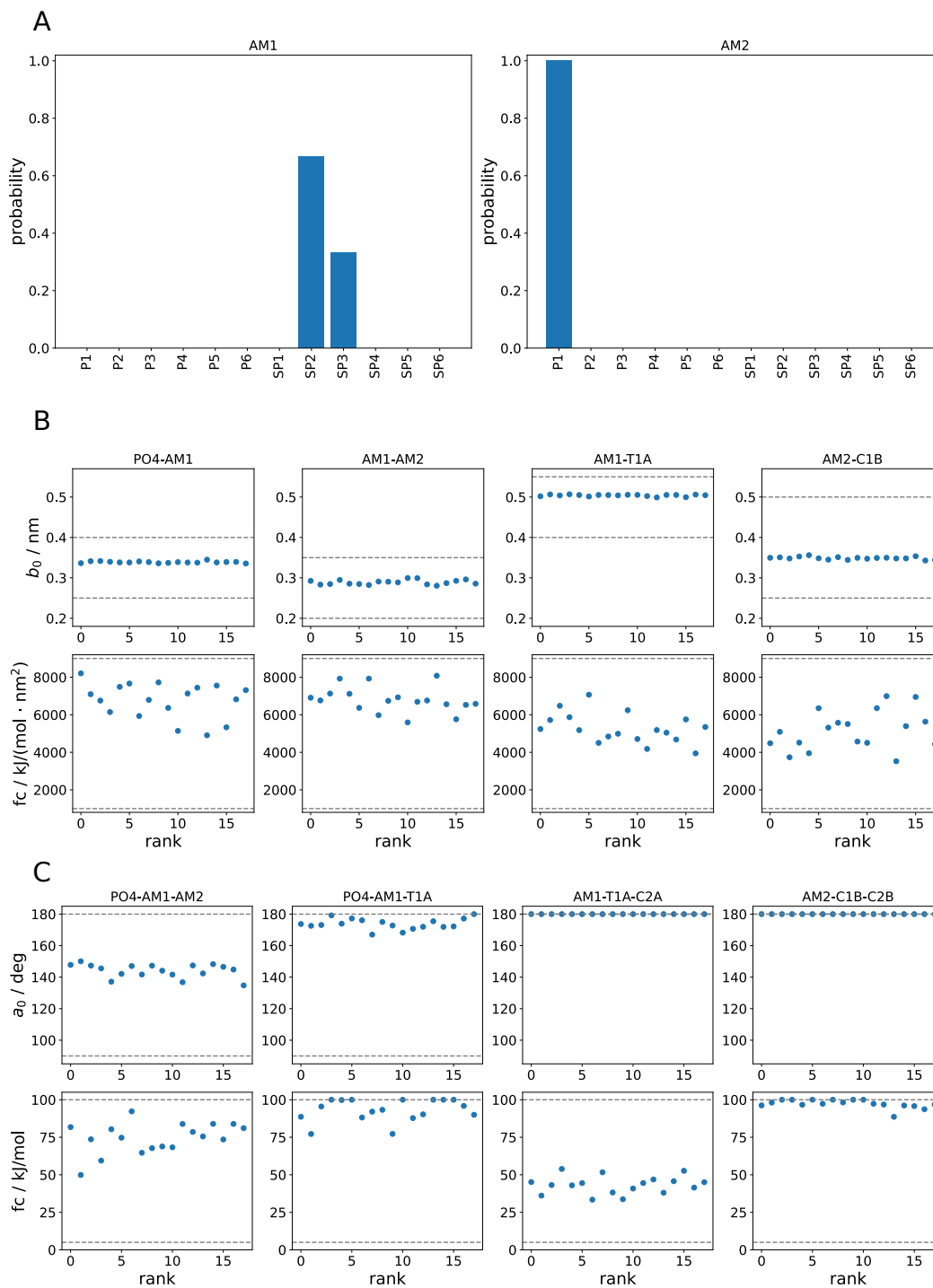


Figure 6: Force field parameters of the set of statistically equivalent solutions \mathcal{P}^g for the sphingolipid linker region. **A)** Bead probability distributions. **B)** Bond parameters. Dashed lines are upper and lower parameter limits. **C)** Angle parameters. Dashed lines are upper and lower parameter limits. The equilibrium angles of AM1-T1A-C2A and AM2-C1B-C2B are not varied during optimization. They are fixed at 180° .

4.2 Noise-mitigation improves quality of parameterized models

We investigated whether the simple noise-mitigation strategy described in Section 2.2.1 can improve the quality of solutions found by the algorithm. The swarm size, training systems, observables and weights are the same as in Section 4.1. We tested three different resampling allocation settings and compared these to the mv-PSO without noise-mitigation. Each optimization run was given a fixed computational budget of 16128 MD simulation slots. With a swarm size of 64 particles, and 3 training systems required for one full objective function evaluation, this amounts to 84 iterations for the mv-PSO without resampling (named mv-PSO-R0). In the optimization runs with resampling an initial computational budget of $64 \cdot 3 = 192$ MD simulation slots was used for one full objective function evaluation of each particle, and a second equally sized computational budget was allocated to reevaluate the melting temperature (the target observable with the largest variance) of the best 16, best 32, or all 64 candidate solutions of the current iteration. For brevity we will refer to these as mv-PSO-R16, mv-PSO-R32, and mv-PSO-R64. Due to the fixed computational budget, for each particle involved in resampling, T_m was reevaluated 12, 6, or 3 times. As half of the total computational budget was used for resampling, the number of iterations was set to 42 in these runs.

From the literature on PSO noise-mitigation^{24,55} we draw the expectation that which of the resampling, or no resampling, strategies is the best, depends on the level of noise. If noise levels are very low, the additional number of possible iterations, when forgoing resampling, could lead to better solutions. For intermediate noise levels, initial fitness evaluation results in a sufficient differentiation of good and bad solutions, i.e., overall sorting is roughly correct, and the focus on improving sorting of the very best solutions is most helpful. In case of even higher noise levels initial sorting would be vastly incorrect and a larger fraction of the swarm needs to be resampled to achieve satisfactory overall sorting. As a consequence, the sorting quality of the very top would be degraded, as there is less computational budget allocated here.

The true quality of a candidate solution is not necessarily reflected by the cost estimated during an optimization run, as there is some uncertainty in estimates of target observables other than T_m , and the confidence level of the T_m estimation with different resampling settings differs vastly. Therefore, validation is required. As we are mostly interested in the quality verification of the best solutions, the first step of the screen-to-the-best procedure from Boesel et al.²⁷ can be used to select the statistically equivalent set of candidate solutions. For mv-PSO-R16 the set \mathcal{P}_t^g contains 69 candidate solutions. Due to the increased uncertainty in mv-PSO-R32 and mv-PSO-R64, their respective sets \mathcal{P}_t^g contain hundreds of candidate solutions. To keep the computational cost for validation manageable, we selected only the 72 best solutions of these optimization runs for validation. As there are no variance estimates in the optimization run without resampling, the selection procedure is not applicable. Again, the 72 best solutions from the optimization run were selected for validation. All candidate solutions chosen for validation were fully (all training systems, all observables) reevaluated 20 times. The resulting rerun cost vs. the originally estimated cost is shown in Figure 7. Clearly, mv-PSO-R16 gave the best results, while the quality of the best solutions in the three other cases does not differ much. Furthermore, the fact that for all selected candidate solutions of mv-PSO-R0 the rerun cost estimate is substantially higher than the original cost estimate indicates that these original estimate are particularly favorable. While there are also candidate solutions with substantial differences in original and rerun cost for the resampling systems – in this case mostly caused by APL fluctuations – these are much less frequent and there is much better correlation between original and rerun cost (Pearson correlation coefficient 0.21 vs. 0.64, for mv-PSO-R0 and mv-PSO-R16, respectively).

Our interpretation of these results is the following: The noise level is low enough, so that even without noise-mitigation, the sorting of candidate solutions is correct in a coarser sense and the swarm is guided towards the "correct" vicinity in parameter space. Yet, noise levels are substantial enough, so that resolution of finer cost differences is impeded. Only the concentrated allocation of the resampling budget on the top 16 solutions lowers the cost

estimation errors sufficiently such that improved candidate solutions can be found.

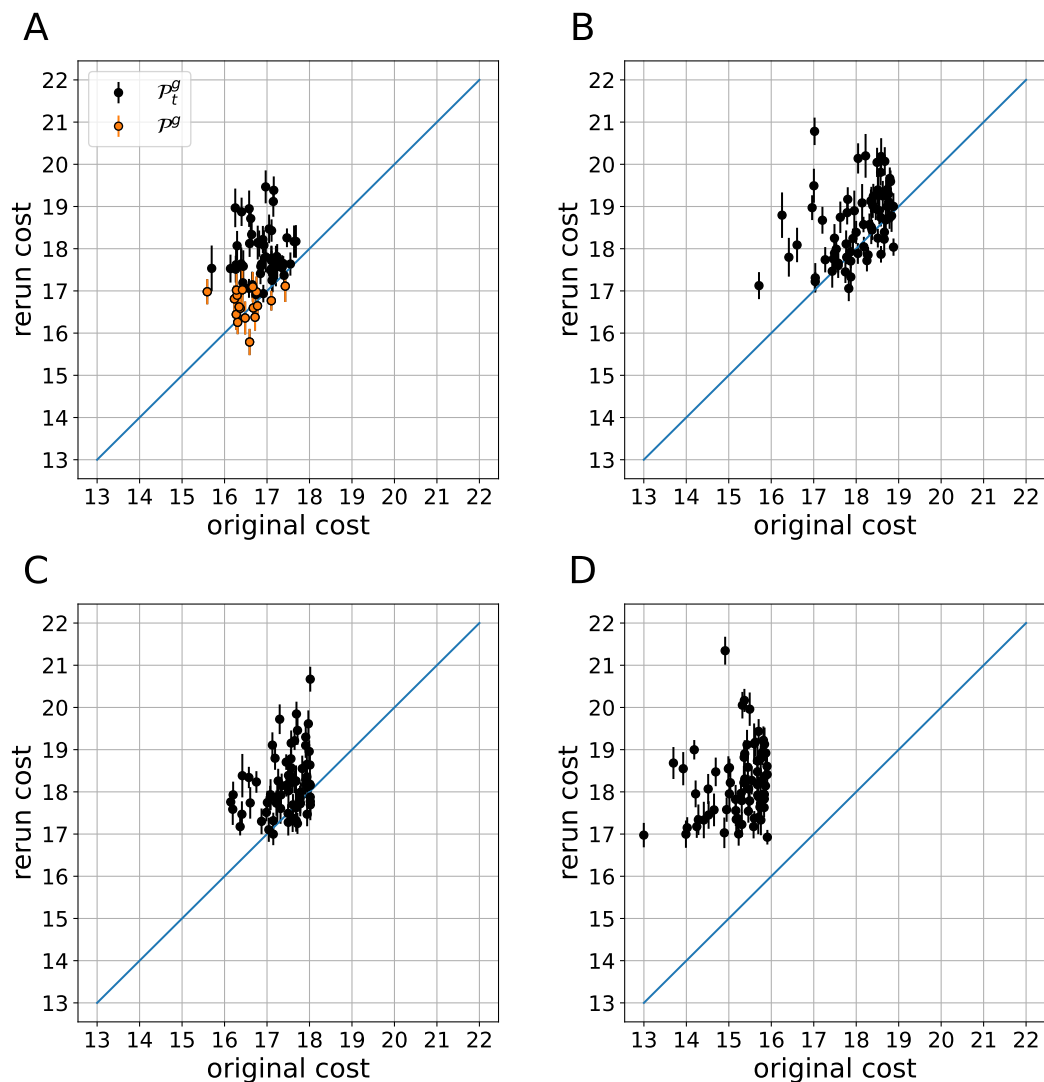


Figure 7: Comparison of cost estimated during the optimization run and average cost estimated from repeated reruns of \mathcal{P}_t^g in (A) and the 72 best candidate solutions in (B)-(D). Error bars are standard errors. (A) Original cost: 16 particles resampled, 1+12 T_m samples. (B) Original cost: 32 particles resampled, 1+6 T_m samples. (C) Original cost: 64 particles resampled, 1+3 T_m samples. (D) No resampling during optimization, but twice as many iterations

5 Discussion & Conclusion

We have illustrated how to apply mixed-variable particle swarm optimization for automated CG molecule parameterization. As an example application, we parameterized the sphingolipid linker region for the Martini 3 FF. The newly parameterized sphingomyelin model reproduces important target observables accurately, including the melting temperature, which was ≈ 20 K off target before and is now within ≈ 2 K of the experimental reference. Notably, reproduction of experimental melting temperatures had been historically problematic in Martini lipid models.⁵⁶

The mixed-variable approach offers a major advantage when parameterizing molecules for building-block force fields. Due the explicit use of building blocks, every candidate model is a valid parameterization in the given FF. Otherwise, changing non-bonded interaction parameters of the FF's building blocks breaks the validity of their parameterization. Candidate solutions generated by a continuous treatment of non-bonded interactions have to be converted to a valid FF model, followed by additional validation of this model.

A drawback of the mixed-variable treatment is that some advanced improvements to PSO, such as fuzzy parameter tuning of Nobile et al.,⁵⁷ are not directly applicable to mv-PSO, because in the categorical representation there is no similarity metric, which is utilized in the PSO parameter tuning. This could be overcome by using discrete ordered representation for non-bonded interactions instead of the categorical treatment.

One of the great benefits of automated parameterization algorithms is the simultaneous optimization against multiple structural and thermodynamic target data. As thermodynamic observables can be expensive to estimate accurately in MD simulations, the formal consideration of noise in objective function values is an important conceptual improvement. As demonstrated, optimization with applied noise-mitigation produced significantly better solutions and the utilized screen-to-the-best procedure provides a systematic approach to the post-optimization selection of the best model.

Although we have demonstrated the adverse effects of objective function value noise on

the sorting and performance of PSO, it is important to note that the non-deterministic nature of particle swarm optimization necessitates multiple repetitions of full optimization runs to confidently determine the most effective noise-mitigation setting. Achieving a high level of confidence in identifying the optimal approach would require a significant number of iterations. Furthermore, the ‘ground truth’, i.e., the true score of a candidate parameterization, is unknown, hence a large amount of validation simulations would be required. This is not feasible, due to a high computational cost. Rigorous development and testing of noise-mitigation strategies should not be done with objective function evaluations that require costly MD simulations, and are therefore beyond the scope of this paper. Moreover, the additionally gained insight, would only be of moderate value. The PSO literature has shown that under significant noise PSO performance is degraded and performance differences between resampling methods for noise-mitigation are problem and noise-level dependent. Generally, noise-mitigation methods employing OCBA perform the best under various circumstances,^{24,58} but its sequential secondary budget allocation puts constraints on the parallelization of the parameterization algorithm. Still, its integration into the parameterization pipeline should be explored in the future.

Together with the general benefits of automation, the here-presented conceptual advantages will further facilitate rigorous CG molecule parameterization. The CGCompiler Python package that comes with our method is tailor-made for parameterization tasks in building-block FFs, such as Martini. Also larger building blocks, i.e., a molecule class with shared regions can be parameterized simultaneously. Our approach is not limited to lipid parameterization, but can be applied to any kind of molecule. CGCompiler can be easily adapted to the needs of a specific parameterization task. Implementing new observables is not much different from writing Python functions for analyzing MD data. Importantly, our automation platform eases collaborations between individual researchers since a clear overview of the parameterization flow is provided. This also renders force-field reproducibility as well as retrospective force-field corrections, such as corrections to the targets (e.g., improved atomistic

force-fields or simulation settings) or inclusion of additional targets rather straightforward.

The here-presented study focuses on method development and the sphingolipid linker parameterization was merely a test case. The parameters of the head group and lipid tails, predefined in our study, are still actively improved/(re)parameterized by the core developers.¹¹ Once these final parameters are released, reparameterization of the linker may be necessary, ideally with an even broader set of training systems, including liquid ordered-disordered phase behavior.

In order to achieve fully automated molecule parameterization in high-throughput applications, the development of an automated mapping and selection of bonded terms remains a crucial component. Currently, mapping and parameter optimization are separate tasks, but integrating an automated mapping scheme into the parameterization pipeline could be facilitated prior to employing mixed-variable particle swarm optimization, utilizing CGCompiler. The choice of bonded parameters not only influences the accuracy of the model but also impacts simulation stability. Various strategies, such as the use of virtual sites, restricted bending potentials, hinge and "divide and conquer" constructions,^{59,60} have been previously described to address instability. Additionally, careful consideration of constraints is necessary to ensure simulation stability and prevent artificial temperature gradients.^{61,62} These aspects should be incorporated as essential steps in a future fully automated parameterization pipeline.

Another future prospect is the advancement of true non-scalarized multi-objective optimization, which eliminates the need for user-defined weights on the targets within the objective function. However, it can also be argued that these user-defined weights, which reflect the importance of targets based on intuition, experience, or additional knowledge, along with the predefined set of relevant structural and thermodynamic targets for the CG force-field, encompass what is commonly known as the "force-field's philosophy". In this sense, the user-defined weights embody the guiding principles that shape the force-field.

Acknowledgement

K.S.S. and H.J.R. thank the NWO Vidi Scheme, The Netherlands, (project number: 723.016.005) for funding this work. H.J.R. thanks the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) for funding this work under Germany's Excellence Strategy - EXC 2033 - 390677874 - RESOLV. K.S.S. and H.J.R. gratefully acknowledge the computing time granted by the Resource Allocation Board and provided on the supercomputer Lise and Emmy at NHR@ZIB and NHR@Göttingen as part of the NHR infrastructure. The calculations for this research were conducted with computing resources under the project nip00054. K.S.S. and H.J.R. gratefully acknowledge the Gauss Centre for Supercomputing e.V. (www.gauss-centre.eu) for funding this project by providing computing time through the John von Neumann Institute for Computing (NIC) on the GCS Supercomputer JUWELS at Jülich Supercomputing Centre (JSC). P.C.T.S. acknowledges the support of the French National Center for Scientific Research (CNRS) and the funding from a research collaboration agreement with PharmCADD. L.M. acknowledges funding by the Institut National de la Santé et de la Recherche Médicale (INSERM)

Supporting Information Available

The following files are available free of charge.

- SI.pdf: Supporting Information including
 - methods to estimate melting temperature
 - additional plots regarding noise
 - Sphingomyelin-cholesterol 2d-center-of-mass radial distribution functions
 - DPSM Gromacs topology file

Data availability

The CGCompiler Python package is available upon request from the authors and will be made publicly available on github at a later stage.

A Gromacs topology file of the final sphingomyelin (and ceramide) can be download from the Martini Database server⁶³ (<https://mad.ibcp.fr>).

References

- (1) Risselada, H. J.; Marrink, S. J. The molecular face of lipid rafts in model membranes. *Proceedings of the National Academy of Sciences* **2008**, *105*, 17367–17372.
- (2) Marrink, S. J.; Risselada, J.; Mark, A. E. Simulation of gel phase formation and melting in lipid bilayers using a coarse grained model. *Chem .Phys. Lipids* **2005**, *135*, 223–244.
- (3) Risselada, H. J.; Marrink, S. J. The freezing process of small lipid vesicles at molecular resolution. *Soft Matter* **2009**, *5*, 4531–4541.
- (4) Lyubartsev, A. P.; Laaksonen, A. Calculation of effective interaction potentials from radial distribution functions: A reverse Monte Carlo approach. *Phys Rev E* **1995**, *52*, 3730.
- (5) Lyubartsev, A. P. Multiscale modeling of lipids and lipid bilayers. *Eur. Biophys. J.* **2005**, *35*, 53–61.
- (6) Izvekov, S.; Parrinello, M.; Burnham, C. J.; Voth, G. A. Effective force fields for condensed phase systems from ab initio molecular dynamics simulation: A new method for force-matching. *J. Chem. Phys.* **2004**, *120*, 10896–10913.
- (7) Izvekov, S.; Voth, G. A. A multiscale coarse-graining method for biomolecular systems. *J. Phys. Chem. B* **2005**, *109*, 2469–2473.

- (8) Lafitte, T.; Apostolakou, A.; Avendaño, C.; Galindo, A.; Adjiman, C. S.; Müller, E. A.; Jackson, G. Accurate statistical associating fluid theory for chain molecules formed from Mie segments. *J. Chem. Phys.* **2013**, *139*, 154504.
- (9) Papaioannou, V.; Lafitte, T.; Avendaño, C.; Adjiman, C. S.; Jackson, G.; Müller, E. A.; Galindo, A. Group contribution methodology based on the statistical associating fluid theory for heteronuclear molecules formed from Mie segments. *J. Chem. Phys.* **2014**, *140*, 054107.
- (10) Marrink, S. J.; Risselada, H. J.; Yefimov, S.; Tieleman, D. P.; de Vries, A. H. The MARTINI Force Field: Coarse Grained Model for Biomolecular Simulations. *The Journal of Physical Chemistry B* **2007**, *111*, 7812–7824.
- (11) Souza, P. C. T.; Alessandri, R.; Barnoud, J.; Thallmair, S.; Faustino, I.; Grünewald, F.; Patmanidis, I.; Abdizadeh, H.; Bruininks, B. M. H.; Wassenaar, T. A.; Kroon, P. C.; Melcr, J.; Nieto, V.; Corradi, V.; Khan, H. M.; Domański, J.; Javanainen, M.; Martinez-Seara, H.; Reuter, N.; Best, R. B.; Vattulainen, I.; Monticelli, L.; Periolo, X.; Tieleman, D. P.; de Vries, A. H.; Marrink, S. J. Martini 3: a general purpose force field for coarse-grained molecular dynamics. *Nat. Methods* **2021**, *18*, 382–388.
- (12) Alessandri, R.; Souza, P. C. T.; Thallmair, S.; Melo, M. N.; de Vries, A. H.; Marrink, S. J. Pitfalls of the Martini Model. *J. Chem. Theory Comput.* **2019**, *15*, 5448–5460.
- (13) Risselada, H. J. Martini 3: a coarse-grained force field with an eye for atomic detail. *Nat. Methods* **2021**, *18*, 342–343.
- (14) Bejagam, K. K.; Singh, S.; An, Y.; Deshmukh, S. A. Machine-Learned Coarse-Grained Models. *J. Phys. Chem. Lett.* **2018**, *9*, 4667–4672.
- (15) Bejagam, K. K.; Singh, S.; An, Y.; Berry, C.; Deshmukh, S. A. PSO-Assisted Devel-

- opment of New Transferable Coarse-Grained Water Models. *J. Phys. Chem. B* **2018**, *122*, 1958–1971.
- (16) Empereur-Mot, C.; Pesce, L.; Doni, G.; Bochicchio, D.; Capelli, R.; Perego, C.; Pavan, G. M. Swarm-CG: Automatic Parametrization of Bonded Terms in MARTINI-Based Coarse-Grained Models of Simple to Complex Molecules via Fuzzy Self-Tuning Particle Swarm Optimization. *ACS Omega* **2020**, *5*, 32823–32843.
- (17) Empereur-mot, C.; Capelli, R.; Perrone, M.; Caruso, C.; Doni, G.; Pavan, G. M. Automatic multi-objective optimization of coarse-grained lipid force fields using SwarmCG. *J. Chem. Phys.* **2022**, *156*, 024801.
- (18) Empereur-mot, C.; Pedersen, K. B.; Capelli, R.; Crippa, M.; Caruso, C.; Perrone, M.; Souza, P. C. T.; Marrink, S. J.; Pavan, G. M. On the Automatic Optimization of Lipid Models in the Martini Force Field using SwarmCG. *ChemRxiv* **2023**,
- (19) Taghiyeh, S.; Xu, J. A new particle swarm optimization algorithm for noisy optimization problems. *Swarm Intell* **2016**, *10*, 161–192.
- (20) van Meer, G.; Voelker, D. R.; Feigenson, G. W. Membrane lipids: where they are and how they behave. *Nat. Rev. Mol. Cell Biol.* **2008**, *9*, 112–124.
- (21) Eberhart, R.; Kennedy, J. A new optimizer using particle swarm theory. MHS'95. Proceedings of the Sixth International Symposium on Micro Machine and Human Science. 1995.
- (22) Wang, F.; Zhang, H.; Zhou, A. A particle swarm optimization algorithm for mixed-variable optimization problems. *Swarm Evol. Comput.* **2021**, *60*, 100808.
- (23) Rubner, Y.; Tomasi, C.; Guibas, L. J. The Earth Mover's Distance as a Metric for Image Retrieval. *Int J Comput Vision* **2000**, *40*, 99–121.

- (24) Rada-Vilela, J.; Johnston, M.; Zhang, M. Population statistics for particle swarm optimization: Resampling methods in noisy optimization problems. *Swarm Evol. Comput.* **2014**, *17*, 37–59.
- (25) Pan, H.; Wang, L.; Liu, B. Particle swarm optimization for function optimization in noisy environment. *Appl Math Comput* **2006**, *181*, 908–919.
- (26) Chen, C.-H.; Lin, J.; Yücesan, E.; Chick, S. E. Simulation Budget Allocation for Further Enhancing the Efficiency of Ordinal Optimization. *Discrete Event Dyn Syst* **2000**, *10*, 251–270.
- (27) Boesel, J.; Nelson, B. L.; Kim, S.-H. Using Ranking and Selection to "Clean up" after Simulation Optimization. *Oper. Res.* **2003**, *51*, 814–825.
- (28) Methorst, J.; van Hilten, N.; Risselada, H. J. Inverse design of cholesterol attracting transmembrane helices reveals a paradoxical role of hydrophobic length. *bioRxiv* **2021**,
- (29) Abraham, M. J.; Murtola, T.; Schulz, R.; Páll, S.; Smith, J. C.; Hess, B.; Lindahl, E. GROMACS: High performance molecular simulations through multi-level parallelism from laptops to supercomputers. *SoftwareX* **2015**, *1-2*, 19–25.
- (30) Michaud-Agrawal, N.; Denning, E. J.; Woolf, T. B.; Beckstein, O. MDAAnalysis: A toolkit for the analysis of molecular dynamics simulations. *J. Comput. Chem.* **2011**, *32*, 2319–2327.
- (31) Gowers, R.; Linke, M.; Barnoud, J.; Reddy, T.; Melo, M.; Seyler, S.; Domański, J.; Dotson, D.; Buchoux, S.; Kenney, I.; Beckstein, O. MDAAnalysis: A Python Package for the Rapid Analysis of Molecular Dynamics Simulations. Proceedings of the 15th Python in Science Conference. 2016.
- (32) Smith, P.; Lorenz, C. D. LiPyphilic: A Python Toolkit for the Analysis of Lipid Membrane Simulations. *J. Chem. Theory Comput.* **2021**, *17*, 5907–5919.

- (33) Virtanen, P.; Gommers, R.; Oliphant, T. E.; Haberland, M.; Reddy, T.; Cournapeau, D.; Burovski, E.; Peterson, P.; Weckesser, W.; Bright, J.; van der Walt, S. J.; Brett, M.; Wilson, J.; Millman, K. J.; Mayorov, N.; Nelson, A. R. J.; Jones, E.; Kern, R.; Larson, E.; Carey, C. J.; Polat, İ.; Feng, Y.; Moore, E. W.; VanderPlas, J.; Laxalde, D.; Perktold, J.; Cimrman, R.; Henriksen, I.; Quintero, E. A.; Harris, C. R.; Archibald, A. M.; Ribeiro, A. H.; Pedregosa, F.; van Mulbregt, P.; Vijaykumar, A.; Bardelli, A. P.; Rothberg, A.; Hilboll, A.; Kloeckner, A.; Scopatz, A.; Lee, A.; Rokem, A.; Woods, C. N.; Fulton, C.; Masson, C.; Häggström, C.; Fitzgerald, C.; Nicholson, D. A.; Hagen, D. R.; Pasechnik, D. V.; Olivetti, E.; Martin, E.; Wieser, E.; Silva, F.; Lenders, F.; Wilhelm, F.; Young, G.; Price, G. A.; Ingold, G.-L.; Allen, G. E.; Lee, G. R.; Audren, H.; Probst, I.; Dietrich, J. P.; Silterra, J.; Webber, J. T.; Slavič, J.; Nothman, J.; Buchner, J.; Kulick, J.; Schönberger, J. L.; de Miranda Cardoso, J. V.; Reimer, J.; Harrington, J.; Rodríguez, J. L. C.; Nunez-Iglesias, J.; Kuczynski, J.; Tritz, K.; Thoma, M.; Newville, M.; Kümmerer, M.; Bolingbroke, M.; Tartre, M.; Pak, M.; Smith, N. J.; Nowaczyk, N.; Shebanov, N.; Pavlyk, O.; Brodtkorb, P. A.; Lee, P.; McGibbon, R. T.; Feldbauer, R.; Lewis, S.; Tygier, S.; Sievert, S.; Vigna, S.; Peterson, S.; More, S.; Pudlik, T.; Oshima, T.; Pingel, T. J.; Robitaille, T. P.; Spura, T.; Jones, T. R.; Cera, T.; Leslie, T.; Zito, T.; Krauss, T.; Upadhyay, U.; Halchenko, Y. O.; and, Y. V.-B. SciPy 1.0: fundamental algorithms for scientific computing in Python. *Nat. Methods* **2020**, *17*, 261–272.
- (34) Pele, O.; Werman, M. Fast and robust Earth Mover's Distances. 2009 IEEE 12th International Conference on Computer Vision. 2009.
- (35) Pele, O.; Werman, M. A Linear Time Histogram Metric for Improved SIFT Matching. Computer Vision – ECCV 2008. Berlin, Heidelberg, 2008; pp 495–508.
- (36) Nguyen, H.; Case, D. A.; Rose, A. S. NGLview–interactive molecular graphics for Jupyter notebooks. *Bioinformatics* **2017**, *34*, 1241–1242.

- (37) Klauda, J. B.; Venable, R. M.; Freites, J. A.; O'Connor, J. W.; Tobias, D. J.; Mondragon-Ramirez, C.; Vorobyov, I.; MacKerell, A. D.; Pastor, R. W. Update of the CHARMM All-Atom Additive Force Field for Lipids: Validation on Six Lipid Types. *J. Phys. Chem. B* **2010**, *114*, 7830–7843.
- (38) Venable, R. M.; Sodt, A. J.; Rogaski, B.; Rui, H.; Hatcher, E.; MacKerell, A. D.; Pastor, R. W.; Klauda, J. B. CHARMM All-Atom Additive Force Field for Sphingomyelin: Elucidation of Hydrogen Bonding and of Positive Curvature. *Biophys. J.* **2014**, *107*, 134–145.
- (39) Wang, E.; Klauda, J. B. Molecular Dynamics Simulations of Ceramide and Ceramide-Phosphatidylcholine Bilayers. *J. Phys. Chem. B* **2017**, *121*, 10091–10104.
- (40) Jo, S.; Lim, J. B.; Klauda, J. B.; Im, W. CHARMM-GUI Membrane Builder for Mixed Bilayers and Its Application to Yeast Membranes. *Biophys. J.* **2009**, *97*, 50–58.
- (41) Wu, E. L.; Cheng, X.; Jo, S.; Rui, H.; Song, K. C.; Dávila-Contreras, E. M.; Qi, Y.; Lee, J.; Monje-Galvan, V.; Venable, R. M.; Klauda, J. B.; Im, W. CHARMM-GUI Membrane Builder toward realistic biological membrane simulations. *J. Comput. Chem.* **2014**, *35*, 1997–2004.
- (42) Lee, J.; Cheng, X.; Swails, J. M.; Yeom, M. S.; Eastman, P. K.; Lemkul, J. A.; Wei, S.; Buckner, J.; Jeong, J. C.; Qi, Y.; Jo, S.; Pande, V. S.; Case, D. A.; Brooks, C. L.; MacKerell, A. D.; Klauda, J. B.; Im, W. CHARMM-GUI Input Generator for NAMD, GROMACS, AMBER, OpenMM, and CHARMM/OpenMM Simulations Using the CHARMM36 Additive Force Field. *J. Chem. Theory Comput.* **2015**, *12*, 405–413.
- (43) Hess, B.; Bekker, H.; Berendsen, H. J. C.; Fraaije, J. G. E. M. LINCS: A linear constraint solver for molecular simulations. *J. Comput. Chem.* **1997**, *18*, 1463–1472.
- (44) Darden, T.; York, D.; Pedersen, L. Particle mesh Ewald: An $N \cdot \log(N)$ method for Ewald sums in large systems. *J. Chem. Phys.* **1993**, *98*, 10089–10092.

- (45) Bussi, G.; Donadio, D.; Parrinello, M. Canonical sampling through velocity rescaling. *J. Chem. Phys.* **2007**, *126*, 014101.
- (46) Parrinello, M.; Rahman, A. Polymorphic transitions in single crystals: A new molecular dynamics method. *J Appl Phys* **1981**, *52*, 7182–7190.
- (47) Borges-Araújo, L.; Borges-Araújo, A.; Ozturk, T.; Ramirez-Echemendia, D. P.; Fábíán, B.; Carpenter, T. S.; Thallmair, S.; Barnoud, J.; Ingólfsson, H. I.; Hummer, G.; Tieleman, D. P.; Marrink, S. J.; Souza, P. C. T.; Melo, M. N. Martini 3 Coarse-Grained Force Field for cholesterol. *ChemRxiv* **2023**,
- (48) Borges-Araújo, L.; Borges-Araújo, A.; Ozturk, T.; Ramirez-Echemendia, D. P.; Fábíán, B.; Carpenter, T. S.; Thallmair, S.; Barnoud, J.; Ingólfsson, H. I.; Hummer, G.; Tieleman, D. P.; Marrink, S. J.; Souza, P. C. T.; Melo, M. N. Parameterization of cholesterol for the Martini 3 coarse grained force field. 2023; <https://github.com/Martini-Force-Field-Initiative/M3-Sterol-Parameters>.
- (49) Wassenaar, T. A.; Ingólfsson, H. I.; Böckmann, R. A.; Tieleman, D. P.; Marrink, S. J. Computational Lipidomics with insane: A Versatile Tool for Generating Custom Membranes for Molecular Simulations. *J. Chem. Theory Comput.* **2015**, *11*, 2144–2155.
- (50) Berendsen, H. J. C.; Postma, J. P. M.; van Gunsteren, W. F.; DiNola, A.; Haak, J. R. Molecular dynamics with coupling to an external bath. *J. Chem. Phys.* **1984**, *81*, 3684–3690.
- (51) Kowalik, B.; Schubert, T.; Wada, H.; Tanaka, M.; Netz, R. R.; Schneck, E. Combination of MD Simulations with Two-State Kinetic Rate Modeling Elucidates the Chain Melting Transition of Phospholipid Bilayers for Different Hydration Levels. *J. Phys. Chem. B* **2015**, *119*, 14157–14167.
- (52) Sun, L.; Böckmann, R. A. Membrane phase transition during heating and cooling: molecular insight into reversible melting. *Eur. Biophys. J.* **2017**, *47*, 151–164.

- (53) Coppock, P. S.; Kindt, J. T. Determination of Phase Transition Temperatures for Atomistic Models of Lipids from Temperature-Dependent Stripe Domain Growth Kinetics. *J. Phys. Chem. B* **2010**, *114*, 11468–11473.
- (54) Borges-Araújo, L.; Souza, P. C. T.; Fernandes, F.; Melo, M. N. Improved Parameterization of Phosphatidylinositide Lipid Headgroups for the Martini 3 Coarse-Grain Force Field. *J. Chem. Theory Comput.* **2021**, *18*, 357–373.
- (55) Rada-Vilela, J.; Johnston, M.; Zhang, M. Deception, blindness and disorientation in particle swarm optimization applied to noisy problems. *Swarm Intell* **2014**, *8*, 247–273.
- (56) Marrink, S. J.; Risselada, J.; Mark, A. E. Simulation of gel phase formation and melting in lipid bilayers using a coarse grained model. *Chem. Phys. Lipids* **2005**, *135*, 223–244.
- (57) Nobile, M. S.; Cazzaniga, P.; Besozzi, D.; Colombo, R.; Mauri, G.; Pasi, G. Fuzzy Self-Tuning PSO: A settings-free algorithm for global optimization. *Swarm Evol. Comput.* **2018**, *39*, 70–85.
- (58) Rada-Vilela, J.; Johnston, M.; Zhang, M. Population statistics for particle swarm optimization: Hybrid methods in noisy optimization problems. *Swarm Evol. Comput.* **2015**, *22*, 15–29.
- (59) Bulacu, M.; Goga, N.; Zhao, W.; Rossi, G.; Monticelli, L.; Periole, X.; Tieleman, D. P.; Marrink, S. J. Improved Angle Potentials for Coarse-Grained Molecular Dynamics Simulations. *J. Chem. Theory Comput.* **2013**, *9*, 3282–3292.
- (60) Alessandri, R.; Barnoud, J.; Gertsen, A. S.; Patmanidis, I.; de Vries, A. H.; Souza, P. C. T.; Marrink, S. J. Martini 3 Coarse-Grained Force Field: Small Molecules. *Adv. Theory Simul.* **2021**, *5*, 2100391.
- (61) Fábíán, B.; Thallmair, S.; Hummer, G. Optimal Bond Constraint Topology for Molec-

ular Dynamics Simulations of Cholesterol. *J. Chem. Theory Comput.* **2023**, *19*, 1592–1601.

- (62) Thallmair, S.; Javanainen, M.; Fábíán, B.; Martinez-Seara, H.; Marrink, S. J. Non-converged Constraints Cause Artificial Temperature Gradients in Lipid Bilayer Simulations. *J. Phys. Chem. B* **2021**, *125*, 9537–9546.
- (63) Hilpert, C.; Beranger, L.; Souza, P. C. T.; Vainikka, P. A.; Nieto, V.; Marrink, S. J.; Monticelli, L.; Launay, G. Facilitating CG Simulations with MAD: The MArtini Database Server. *J. Chem. Inf. Model.* **2023**, *63*, 702–710.

Graphical TOC Entry

