

ms2: A Molecular Simulation Tool for Thermodynamic Properties

Stephan Deublein^a, Bernhard Eckl^b, Jürgen Stoll^b, Sergey V. Lishchuk^c, Gabriela Guevara-Carrion^a, Colin W. Glass^d, Thorsten Merker^a, Martin Bernreuther^d, Hans Hasse^a, Jadran Vrabec^{e,*}

^a*Lehrstuhl für Thermodynamik, Universität Kaiserslautern, 67653 Kaiserslautern, Germany*

^b*Institut für Technische Thermodynamik und Thermische Verfahrenstechnik, Universität Stuttgart, 70550 Stuttgart, Germany*

^c*Department of Mathematics, University of Leicester, Leicester LE1 7RH, United Kingdom*

^d*Höchstleistungsrechenzentrum Universität Stuttgart (HLRS), 70550 Stuttgart, Germany*

^e*Lehrstuhl für Thermodynamik und Energietechnik, Universität Paderborn, 33098 Paderborn, Germany*

Abstract

This work presents the molecular simulation program *ms2* that is designed for the calculation of thermodynamic properties of bulk fluids in equilibrium consisting of small electro-neutral molecules. *ms2* features the two main molecular simulation techniques, molecular dynamics (MD) and Monte-Carlo. It supports the calculation of vapor-liquid equilibria of pure fluids and multi-component mixtures described by rigid molecular models on the basis of the grand equilibrium method. Furthermore, it is capable of sampling various classical ensembles and yields numerous thermodynamic properties. To evaluate the chemical potential, Widom's test molecule method and gradual insertion are implemented. Transport properties are determined by equilibrium MD simulations following the Green-Kubo formalism. *ms2* is designed to meet the requirements of academia and industry, particularly achieving short response times and straightforward handling. It is written in Fortran90 and optimized for a fast execution on a broad range of computer architectures, spanning from single processor PCs over PC-clusters and vector computers to high-end parallel machines. The standard Message Passing Interface (MPI) is used for parallelization and *ms2* is therefore easily portable onto a broad range of computing platforms. Auxiliary feature tools facilitate the interaction with the code and the interpretation of input and output files. The accuracy and reliability of *ms2* has been shown for a large variety of fluids in preceding work.

Keywords: Molecular simulation, molecular dynamics, Monte-Carlo, grand equilibrium method, vapor-liquid equilibrium, transport properties

1. Program summary

Manuscript: *ms2*: A Molecular Simulation Program for Thermodynamic Properties

Authors: Stephan Deublein, Bernhard Eckl, Jürgen Stoll, Sergey V. Lishchuk, Gabriela Guevara-Carrion, Colin W. Glass, Thorsten Merker, Martin Bernreuther, Hans Hasse, Jadran Vrabec (jadran.vrabec@upb.de)

Title of program: *ms2*

Operating system: Unix/Linux, Windows

*Corresponding author: Jadran Vrabec, Warburger Str. 100, 33098 Paderborn, Germany, Tel.: +49-5251/60-2421, Fax: +49-5251/60-3522, Email: jadran.vrabec@upb.de

Computer: The simulation program *ms2* is usable on a wide variety of platforms, from single processor machines over PC-clusters and vector computers to vector-parallel architectures. (Tested with Fortran compilers, gfortran, Intel compiler, PathScale compiler, Portland Group compiler and Sun Studio compiler)

Memory: *ms2* runs on single processors with 512 MB RAM. The memory demand rises with increasing number of processors used per node and increasing number of molecules.

Distribution format: tar.gz

Keywords: Molecular simulation, molecular dynamics, Monte-Carlo, grand equilibrium method, vapor-liquid equilibrium, transport properties, parallel algorithms

Programming language used: Fortran90

External: Message passing interface (MPI).

Classification: 7.7, 7.9, 12

vectorised: Yes. Message Passing Interface (MPI) protocol

Scalability: Excellent scalability up to 16 processors for molecular dynamics and > 512 processors for Monte-Carlo simulations.

Nature of problem: Calculation of application oriented thermodynamic properties for rigid electro-neutral molecules: vapor-liquid equilibria of pure fluids and multi-component mixtures, thermal and caloric data as well as transport properties.

Method of solution: Molecular dynamics, Monte-Carlo, various classical ensembles, grand equilibrium method, Green-Kubo formalism

Restrictions: None. The system size is user-defined. Typical problems addressed by *ms2* can be solved by simulating systems containing typically 2000 molecules or less.

Unusual Features: Auxiliary feature tools are available for creating input files, analyzing simulation results and visualizing molecular trajectories.

Additional comments: Sample makefiles for multiple operation platforms are provided.

Documentation: Documentation is provided with the installation package and is available at <http://www.ms-2.de>.

Typical running time: The running time of *ms2* depends on the specified problem, the system size and the number of processes used in the simulation. Running four processes on a "Nehalem" processor, simulations calculating vapor-liquid equilibrium data take between two and 12 hours, calculating transport properties between six and 24 hours.

2. Introduction

Due to the advances in computing power, methodological efficiency and the development of accurate force fields, it is understood that "molecular modeling and simulation will become a breakthrough technology that is widely accepted in the chemical industry and applied in conjunction with other predictive methods to meet the industry's evolving fluid property data needs" [1]. At the core of this prospect lies the sound physical basis of molecular modeling and simulation. It allows to adequately describe structure, energetics and dynamics on the microscopic level, which govern material properties on the macroscopic level. Therefore it provides convenient access to thermodynamic properties, particularly if they are difficult or expensive to obtain by experiment, e.g. at high temperatures and pressures, or may not be measured at all, e.g. when new substances are not available in sufficient quantity. Furthermore, for toxic, explosive or in any other way hazardous substances, measuring thermodynamic properties experimentally can be unfeasible even at moderate thermodynamic conditions, while molecular modeling and simulation offers a reliable route [2].

Presently, the chemical industry extends the scope of thermodynamic models that are regularly used from phenomenological equations of state (EOS) or excess Gibbs energy models to more advanced types, such as statistical EOS like SAFT [3, 4] or continuum solvation models like COSMO-RS [5, 6]. It should be noted that all these advanced models are developed using a bottom-up approach, i.e. bridging molecular properties to the macroscopic level. These methods mainly offer access to aggregated properties, such as the Helmholtz energy of bulk fluids and its derivatives. Regarding transport properties, they are not useful.

Molecular modeling and simulation is applicable with very few constraints as both static and dynamic thermodynamic data may be calculated, be it in bulk fluids or in confinements. Both, equilibrium or non-equilibrium conditions, may be studied. Additionally, detailed insight into the mechanisms on the microscopic level is provided. This versatility, however, is associated with a substantial computational effort which is orders of magnitude larger than needed for the methods mentioned above. Traditionally, molecular simulations were carried out in computing centers of national institutions or universities on powerful computing equipment, which is usually unavailable even in large chemical companies.

However, with the increase in computer power, molecular simulations become feasible even off the shelf, if suitable simulations programs are available. The scope of data accessible within reasonable times even for industrial workflows is increasing rapidly. Using molecular modeling and simulation can therefore contribute for a reduction of process development time and production costs [7].

Molecular modeling and simulation is rewarding, because it provides reliable predictions of essentially all thermodynamic properties in a consistent manner [8]. Molecular models yield predictions for any property at any condition. This is of particular interest for industrial applications, where a wide variety of properties needs to be known.

To further stimulate this issue, the Industrial Fluid Properties Simulation Collective [1] has organized six Simulation Challenges to date [8, 9, 10, 11, 12, 1]. The goal of these Simulation Challenges is to assess the capabilities of molecular methods regarding typical industrial tasks, where classical methods are insufficient.

The present paper presents a molecular simulation program, named *ms2*, which successfully competed in two Simulation Challenges, finishing second place 2006 [11] and first place the following year [8]. *ms2* is aimed at the calculation of thermodynamic properties that are needed for applications in the chemical industry, without requiring expert knowledge by the user. It was developed for rigid molecular models with a focus on condensed phases, covering mixtures containing an arbitrary number of small molecular species.

The simulation program *ms2* is optimized for a fast execution on a broad range of computer architectures, spanning from single processor PCs over PC-clusters and vector computers to high end parallel machines. It is a standalone Fortran90 code that does not require any libraries or have any other software prerequisites. *ms2* is provided as a source code. It only needs to be compiled. The structure of *ms2* is modular and object-oriented, except at the very core of energy and force calculations, where some structural compromises were made to achieve a better performance.

ms2 is aimed at homogeneous bulk systems in equilibrium. For VLE calculations, the coexisting phases are simulated independently and subsequently. Thus, it is usually sufficient to sample molecular systems containing in the order of 10^3 molecules to obtain statistically reliable results [13]. Applying standard cut-off radii, the intermolecular interactions are evaluated roughly up to the maximal distance, which is given by the edge length of the cubic simulation volume. Therefore, spatial decomposition schemes for parallelization are not rewarding. Instead, Plimpton's method [14] was implemented for parallel execution of molecular dynamics (MD), which scales reasonably well up to 16 or 32 processors, depending on the particular architecture. For Monte-Carlo (MC), a optimal parallelization scheme was introduced, taking advantage of the stochastic nature of such simulations. The computing resources are splitted into single runs each on one node, sampling the phase space independently. Both parallelization methods allow for short response times¹. A typical molecular simulation of a vapor liquid equilibrium state point takes roughly six hours on a current workstation. The main expenses for determining thermodynamic data *in silico* are installation and maintenance of computing equipment and/or computing time, which is presently below US\$ 0.35/CPU-hour [15].

The simulation program *ms2* presented here was developed in a long-time cooperation of engineers and computer scientists. With the release of *ms2*, it is intended to facilitate the transfer of a state of the art and user-friendly molecular modeling and simulation package to academic institutions and the chemical industry. The release package consists of the simulation program *ms2* itself and auxiliary feature tools for setup and analysis, including a simple 3D visualization tool to monitor the molecular trajectories. Furthermore, accurate molecular models for more than 100 small molecules are provided.

In sections 3 and 4 of this paper, *ms2* and the implemented thermodynamic properties are outlined. In section 5, the structure of the code is explained, while section 6 compares its performance to similar simulation programs and section 7 describes the auxiliary feature tools. Section 8 of this paper draws a conclusion and offers a brief outlook on future developments and improvements.

¹Under "response time" we understand the real time difference between submission and termination of a simulation run.

3. Simulation program *ms2*

The simulation program *ms2* is capable of sampling the phase space for rigid electro-neutral molecules by applying the two most fundamental molecular simulation techniques, i.e. MD and MC. MD simulations rely on the numerical solution of Newton's equations of motion: for a point in time, the intermolecular interactions, in particular the resulting forces and torques, are evaluated and treated as constant for a specified time step. They are the driving forces of the molecular motion. The displacements for the time step are calculated on that basis, resulting in a new configuration. This process is repeated in a loop. The chronologically ordered configurations are a time discretized approximation of a molecular process. Both static and dynamic thermodynamic properties are determined via time averages. MC simulation explores the phase space stochastically for a given molecular system. Molecules in the simulation volume are displaced randomly. The probability of accepting the displacement is chosen such that a representative set of configurations is obtained. The Markov chain of configurations generated in this way allows for a rigorous calculation of static thermodynamic properties via ensemble averages.

3.1. Overview

The simulation program *ms2* allows for the determination of static and dynamic thermodynamic properties in equilibrium. The implemented static properties are:

- Thermal and caloric properties
- Chemical potential
- Vapor-liquid equilibria
- Henry's law constant
- Second virial coefficient

The transport properties can be calculated on the fly during a MD simulation with a reasonable additional computational effort using the Green-Kubo formalism [16, 17]. The implemented dynamic properties are:

- Self-diffusion coefficient
- Maxwell-Stefan diffusion coefficient
- Shear viscosity
- Bulk viscosity

The model class that is supported by *ms2* covers rigid multi-center Lennard-Jones (LJ) 12-6 interaction sites with an arbitrary number of superimposed electrostatic sites [18, 19, 20, 21]. The supported electrostatic models are point charges, point dipoles and point quadrupoles, which can be positioned anywhere within the molecule. Currently, *ms2* is designed for electro-neutral species.

The quality of thermodynamic properties calculated by molecular simulation is basically determined by two factors: first, the employed molecular model, i.e. the force field, which fully defines the thermodynamic properties

and deviates to some degree from the behavior of the real fluid; second, the sampling of the phase space during simulation, which is associated with statistical uncertainties.

Molecular models have been investigated with *ms2* in numerous cases in the past [22, 23, 24, 25, 26, 27, 28]. For many of these models, the geometric and electrostatic interaction parameters were passed on from ab-initio calculations. The remaining parameters were adjusted to reproduce the vapor pressure and the saturated liquid density of the regarded pure substance. These molecular models, combined with *ms2* and its analysis methods, allow for time-efficient high quality simulations.

The second factor can be influenced for a given simulation method by the number of molecules and the number of sampled configurations. The more data, the lower the statistical uncertainties.

3.2. Methods

3.2.1. Ensembles

Following ensembles are currently supported by *ms2*:

- canonical ensemble (NVT) - MD and MC
- micro-canonical ensemble (NVE) - MD
- isobaric-isothermal ensemble (NpT) - MD and MC
- grand equilibrium method (pseudo- μVT) - MC

Most of these ensembles are well known and widely in use [29]. Therefore, the discussion is restricted here to the grand equilibrium method, since it was developed recently [30] and is probably new to many readers.

The grand equilibrium method is a technique to determine the VLE of pure substances or mixtures. It is a two-step procedure, where the coexisting phases are simulated independently and subsequently. The specified thermodynamic variables for the VLE are the temperature T and the composition \boldsymbol{x} of the liquid. In the first step, one NpT simulation of the liquid phase is performed at T , \boldsymbol{x} and some pressure p_0 to determine the chemical potentials μ_i^l and the partial molar volumes v_i^l of all components, which corresponds to the molar volume in case of a pure substance. If the entropic properties are determined by Widom's test molecule method [31], both MD or MC can be used to sample the phase space. A more advanced technique for these properties is implemented in combination with MC, i.e. gradual insertion [32, 33] (see below). On the basis of the chemical potentials and partial molar volumes at p_0 , first order Taylor expansions can be made for the pressure dependence

$$\mu_i^l(T, \boldsymbol{x}, p) \approx \mu_i^l(T, \boldsymbol{x}, p_0) + v_i^l(T, \boldsymbol{x}, p_0) \cdot (p - p_0) . \quad (1)$$

Note that *ms2* yields $\mu_i(T, \boldsymbol{x}, p) - \mu_i^{\text{id}}(T)$, where $\mu_i^{\text{id}}(T)$ is the temperature dependent part of the ideal gas contribution to the chemical potential of the pure component. $\mu_i^{\text{id}}(T)$ does not need to be determined for VLE calculations, because it cancels out when Eq. (1) is equated to the corresponding expression for the vapor.

In the second step, one pseudo- μVT simulation [30] is performed for the vapor phase on the basis of Eq. (1) that yields the saturated vapor state point of the VLE. This simulation takes place in a pseudo ensemble in the

sense that the specified chemical potentials are not constant, but dependent on the actual pressure of the vapor phase. However, experience based on hundreds of systems [25, 28] shows that the vapor phase simulation rapidly converges to the saturated vapor state point during equilibration so that effectively the equilibrium chemical potentials are specified via the attained vapor pressure. The number of molecules varies during the vapor phase simulation of the grand equilibrium method. Starting from a specified number of molecules in an arbitrarily chosen gaseous state, *ms2* adjusts the extensive volume V after an equilibration period so that the vapor density does not have to be known in advance.

The grand equilibrium method has been widely used and compares favorably to other methods used for vapor liquid equilibrium simulations, as shown in section 6.

3.2.2. Integrators, thermostat and barostat

In *ms2*, two integrators are implemented to solve Newton's equations of motion during MD simulation: Leapfrog and Gear predictor-corrector. These integrators are well known and described in literature [29]. The leapfrog integrator is a second order integrator that requires little computational effort, while being robust for many applications. The error of the method scales quadratically with the length of the chosen time step. The Gear predictor-corrector integrator is implemented with fifth order and is more accurate for small time steps compared to the Leapfrog integration [29]. The computational demand for both integration schemes is similar as implemented in *ms2*. The Gear integration, though being of higher order, is only 0.3% slower than the Leapfrog algorithm for the total calculation of one MD time step with a molecular model composed of three LJ sites, having three rotational degrees of freedom.

For MC, the Markov chain is generated by repeatedly disturbing the system by translational or rotational motion of a molecule and evaluating the resulting configurations with respect to energy using Metropolis acceptance criterion. The thermostat incorporated in *ms2* is velocity scaling. Here, the velocities are scaled such that the actual kinetic energy matches the specified temperature. The scaling is applied equally over all molecular degrees of freedom. The pressure is kept constant using Andersen's barostat in MD, and random volume changes evaluated according to Metropolis acceptance criterion in MC, respectively.

3.2.3. Basics on simulations with *ms2*

The basic molecular simulation techniques employed in *ms2* are well described in the literature [29, 34, 35] and are thus not repeated here. All non-standard methods that are implemented in *ms2* are introduced in section 4 or the Supplementary Material. The following paragraph briefly describes the most significant assumptions and techniques employed in *ms2*.

Simulations with *ms2* are performed in a quasi-infinite Cartesian space, using periodic boundary conditions [36] and the minimum image convention [37]. The molecular interactions are considered to be pairwise additive, like in most other simulation packages. Including three- and many-body interactions leads to an increase in computational effort while the benefits are questionable. As usual, the computational effort in *ms2* is reduced by the introduction of a cut-off radius, up to which the intermolecular interactions are explicitly evaluated. The contributions of interactions with molecules beyond the cut-off radius are accounted for by correction schemes.

For the electrostatic models, the reaction field method is included, while for the dispersive and repulsive interactions, a homogeneous fluid beyond the cut-off radius is assumed. The long range contributions are added to all thermodynamic data that are calculated in *ms2*.

Many thermodynamic properties calculated with *ms2* are residual quantities, indicated by the superscript "res". To compare them to data that are measured on the macroscopic level, the solely temperature dependent contributions of the ideal gas have to be added. These ideal properties are accessible e.g. by quantum chemical methods and can often be found in data bases [38].

3.3. Thermodynamic properties

The simulation program *ms2* calculates the thermodynamic properties from the trajectories (MD) or Markov chains (MC) on the fly. The results are written to file with a specified frequency during the course of the simulation. The statistical uncertainties of all results are estimated using the block averaging method according to Flyvbjerg and Petersen [39] and the error propagation law.

3.3.1. Density, pressure, internal energy and enthalpy

The static thermodynamic properties accessible via *ms2* depend on the chosen ensemble. At constant temperature and volume (*NVT* ensemble), e.g. the pressure is determined by the virial expression W . Other important thermodynamic properties are accessible at constant temperature and pressure. In the *NpT* ensemble, the volume is a fluctuating parameter that on average yields the volume corresponding to the specified temperature and pressure. In both ensembles, the residual internal energy is the sum of all pairwise interaction energies u_{ij} and the appropriate long range correction. The residual enthalpy is linked to the residual internal energy, pressure and volume of the system using the thermodynamic definition.

3.3.2. Second derivatives

The implemented second derivatives vary with the employed ensemble. In the *NVT* ensemble, the residual isochoric heat capacity c_v^{res} is determined by fluctuations of the residual potential energy U^{res} . The partial derivative of the potential energy with respect to the volume at constant temperature $(\partial U^{\text{res}}/\partial V)_T$ is determined by fluctuations of the residual potential energy U^{res} and the virial W .

In the *NpT* ensemble, the residual isobaric heat capacity c_p^{res} , the isothermal compressibility β_T and the volume expansivity α_p are functions of ensemble fluctuations [29]. The residual isobaric heat capacity c_p^{res} is related to fluctuations of the residual enthalpy H^{res} and the isothermal compressibility β_T is calculated from the volume fluctuations. The partial derivative of the residual enthalpy with respect to the pressure at constant temperature $(\partial H^{\text{res}}/\partial p)_T$ is determined by fluctuations of volume and residual internal energy and the volume expansivity α_p again to volume and residual enthalpy fluctuations.

The speed of sound c is the velocity of a sound wave traveling through an elastic medium. It is defined by the isothermal compressibility β_T , the volume expansivity α_p , the isobaric heat capacity c_p and the temperature T . In *ms2*, the speed of sound is calculated both for pure components and mixtures in the *NpT* ensemble.

3.3.3. Chemical potential

The chemical potential of a given component can be separated into the temperature dependent part $\mu_i^{\text{id}}(T)$ of the ideal gas contribution to the chemical potential of the pure component and the remaining contribution $\mu_i(T, \mathbf{x}, p) - \mu_i^{\text{id}}(T)$. The solely temperature dependent ideal gas contribution cancels out in the calculation of phase equilibria and is therefore not implemented in *ms2*. The chemical potential depends on the real substance behavior due to the molecular interactions and can be determined with *ms2* using two different techniques: Widom’s test molecule method and gradual insertion.

Widom. A conceptually straightforward approach to calculate the chemical potential of a component i was presented by Widom [31]. It allows for a determination of the chemical potential with low computational cost, both for pure substances and mixtures. For a three-center LJ fluid, the computational demand is shown in Table 1. The execution time for the determination of the chemical potential increases roughly linearly with the specified number of test molecules that are sampled.

The accuracy of the calculation varies with the number of test molecules and the density of the investigated fluid. For very dense fluids, the results are subject to poor statistics and the method may even fail. Within limits, lower statistical uncertainties of the chemical potential can be achieved by inserting a large number of test molecules into the simulation volume.

Gradual insertion. A more advanced method to determine the chemical potential, which is reliable also at very high densities, is gradual insertion. It is briefly described here, further details including all parameters are discussed in literature [32, 33, 40, 41].

Instead of inserting complete test molecules, a fluctuating molecule is introduced into the simulation, which can appear in different states of coupling with the other molecules. In the decoupled state, the fluctuating molecule does not interact at all with the other molecules, while in the fully coupled state, it acts like a real molecule of the specified component i . Between these states, a set of partially coupled states has to be defined, each with a larger fraction of the real molecule interaction, cf. Figure 1.

The $N-1$ real molecules plus the fluctuating molecule π_l in the state l form a set of sub-ensembles, which can be depicted by the following scheme

$$[N + \pi_0] \leftrightarrow [N + \pi_1] \leftrightarrow \dots \leftrightarrow [N + \pi_l] \leftrightarrow \dots \leftrightarrow [N + \pi_{k-1}] \leftrightarrow [N + \pi_k]. \quad (2)$$

To switch between neighboring sub-ensembles, an additional move is introduced in a standard MC simulation. The probability of accepting a change of the fluctuating molecule from a state of coupling l to a state of coupling m is given by

$$P_{\text{acc}}(l \rightarrow m) = \min \left(1, \frac{\omega_m}{\omega_l} \exp \left(- \frac{\psi_m - \psi_l}{k_{\text{B}}T} \right) \right), \quad (3)$$

where ψ_l denotes the interaction energy of the fluctuating molecule in the state l with all other $N - 1$ real molecules and $k_{\text{B}}T$ is the Boltzmann constant multiplied by the temperature. The states of coupling are weighted by the weighting factors ω_l to avoid an unbalanced sampling of the different states. If specified, the weighting

factors are adjusted during simulation, depending on the number of times N_s the fluctuating molecule appeared in state l , according to

$$\omega_l^{new} = \omega_l^{old} \frac{N_s(k)}{N_s(l)}. \quad (4)$$

The fully coupled state k serves as a reference state for the weighting factors. Local relaxation around the fluctuating molecule is enhanced by biased translational and rotational moves in the vicinity of the fluctuating molecule throughout the simulation [33]. The chemical potential $\mu_i(T, \mathbf{x}, p) - \mu_i^{\text{id}}(T)$ of component i is then determined by

$$\mu_i(T, \mathbf{x}, p) - \mu_i^{\text{id}}(T) = k_B T \ln \left\langle \frac{N_i \omega_k \text{Prob}[N + \pi_0]}{V \omega_0 \text{Prob}[N + \pi_k]} \right\rangle, \quad (5)$$

where $\text{Prob}[N + \pi_0]$ and $\text{Prob}[N + \pi_k]$ are the probabilities to observe an ensemble with the fluctuating molecule in the fully decoupled and fully coupled state, respectively.

The gradual insertion method yields good results for the chemical potential even in cases where Widom's test molecule method fails. Disadvantages of the method are the extended simulation time and the additional effort needed to define the fluctuating states.

3.3.4. Henry's law constant

The solubility of solutes in solvents is characterized by the Henry's law constant. In solvents of a given composition, it is solely a function of temperature. The Henry's law constant H_i is related to the residual chemical potential of the solute i at infinite dilution $\mu_i^{\text{res}, \infty}$ as defined in e.g. [42] by

$$H_i = \rho k_B T \exp \left(\frac{\mu_i^{\text{res}, \infty}}{k_B T} \right), \quad (6)$$

where ρ denotes the solvent density. The residual chemical potential of a component at infinite dilution can often be calculated using Widom's test molecule method. In this case, a simulation of the solvent is performed in the NpT ensemble and the saturated vapor pressure of the solvent, while the solute is introduced as an additional component of the mixture with a molar fraction of zero. Then, the solute is only added in form of test molecules to determine its chemical potential at infinite dilution. For dense phases, the chemical potential and thus the Henry's law constant can be calculated with the gradual insertion method, if Widom's method fails. The number of solute molecules in the simulation is reduced to one, representing the infinitely diluted molecule.

3.3.5. Second virial coefficient

The second virial coefficient is related to the intermolecular potential by [43]

$$B = -2\pi \int_0^\infty dr_{ij} \left\langle \exp \left(-\frac{u_{ij}(r_{ij}, \boldsymbol{\omega}_i, \boldsymbol{\omega}_j)}{k_B T} \right) - 1 \right\rangle_{\boldsymbol{\omega}_i, \boldsymbol{\omega}_j} r_{ij}^2, \quad (7)$$

where the $\langle \dots \rangle$ brackets indicate the average over the orientations $\boldsymbol{\omega}_i$ and $\boldsymbol{\omega}_j$ of two molecules i and j separated by a center of mass distance r_{ij} . The integrand, called Mayer's f -function, is evaluated at numerous distances in an appropriate range. *ms2* evaluates Mayer's f -function by averaging over numerous random orientations at each radius.

3.3.6. Transport properties

In *ms2*, transport properties are calculated using equilibrium MD. Here, the fluctuations of a system around its equilibrium state are evaluated as a function of time. The Green-Kubo formalism relates these microscopic fluctuations to the respective transport properties.

Diffusion coefficients. The self-diffusion coefficient D_i is related to the mass flux of single molecules within a fluid. Therefore, the relevant Green-Kubo expression is based on the individual molecule velocity autocorrelation function [44]. Since all molecules contribute to the self-diffusion coefficient, the autocorrelation function is averaged over all N_i molecules of component i in the ensemble to achieve better statistics. In binary mixtures, the Maxwell-Stefan diffusion coefficient \mathcal{D}_{ij} is defined by [45]

$$\mathcal{D}_{ij} = \frac{x_j}{x_i} \Lambda_{ii} + \frac{x_i}{x_j} \Lambda_{jj} - \Lambda_{ij} - \Lambda_{ji}, \quad (8)$$

where $x_i = N_i/N$ and Λ_{ij} can be written in terms of the center of mass velocity

$$\Lambda_{ij} = \frac{1}{3N} \int_0^\infty dt \left\langle \sum_{k=1}^{N_i} \mathbf{v}_{i,k}(0) \cdot \sum_{l=1}^{N_j} \mathbf{v}_{j,l}(t) \right\rangle. \quad (9)$$

From the expressions above, the collective character of the Maxwell-Stefan diffusion coefficient is evident. This leads to significantly less data for a given system size and time step and therefore to larger statistical uncertainties than in case of the self-diffusion coefficient. Note that equations for the Maxwell-Stefan diffusion coefficient are implemented in *ms2* both for binary and ternary mixtures [45].

Shear viscosity. The shear viscosity η , as defined by Newton's "law" of viscosity, is a measure of the resistance of a fluid to a shearing force [46]. It is associated with the momentum transport under the influence of velocity gradients. Hence, the shear viscosity can be related to the time autocorrelation function of the off-diagonal elements of the stress tensor \mathbf{J}_p [44]

$$\eta = \frac{1}{Vk_B T} \int_0^\infty dt \left\langle J_p^{xy}(t) \cdot J_p^{xy}(0) \right\rangle. \quad (10)$$

Averaging over all three independent elements of the stress tensor, i.e. J_p^{xy} , J_p^{xz} and J_p^{yz} , improves the statistics. The component J_p^{xy} of the microscopic stress tensor \mathbf{J}_p is given in terms of the molecular positions and velocities by [46]

$$J_p^{xy} = \sum_{i=1}^{N-1} m v_i^x v_i^y - \sum_{i=1}^{N-1} \sum_{j=i+1}^N \sum_{k=1}^n \sum_{l=1}^n r_{ij}^x \frac{\partial u_{ij}}{\partial r_{kl}^y}. \quad (11)$$

Here, the lower indices l and k indicate the n interaction sites of a molecule and the upper indices x and y denote the spatial vector components, e.g. for velocity v_i^x or site-site distance r_{ij}^x . The first term of Eq. (11) is the kinetic energy contribution and the second term is the potential energy contribution to the shear viscosity. Consequently, the Green-Kubo integral (10) can be decomposed into three parts, i.e. the solely kinetic energy contribution, the solely potential energy contribution and the mixed kinetic-potential energy contribution [46].

Bulk viscosity. The bulk viscosity η_v refers to the resistance to dilatation of an infinitesimal volume element at constant shape [47]. The bulk viscosity can be calculated by integration of the time-autocorrelation function of the diagonal elements of the stress tensor and an additional term that involves the product of pressure p and volume V , which does not appear in the shear viscosity, cf. Eq. (10). In the NVE ensemble, the bulk viscosity is given by [44, 48]

$$\eta_v = \frac{1}{Vk_{\text{B}}T} \int_0^\infty dt \langle (J_p^{xx}(t) - pV(t)) \cdot (J_p^{xx}(0) - pV(0)) \rangle. \quad (12)$$

The diagonal component J_p^{xx} of the microscopic stress tensor \mathbf{J}_p is defined analogous to Eq. (11). The statistics of the ensemble average in Eq. (12) can be improved using all three independent diagonal elements of the stress tensor J_p^{xx} , J_p^{yy} and J_p^{zz} . Eqs. (11) and (12) can directly be applied for mixtures.

4. Simulation program *ms2*: Detailed description

This section describes *ms2* in detail and important options for its application are introduced. It is intended to allow for a better understanding of the simulation program and to facilitate access to the program features.

4.1. Input and output

ms2 was designed to be an easily applicable simulation program. Therefore, the input files are restricted to one file for the definition of simulation scenario and one file for each of the molecular species that are used in the simulation. The output files contain structured information on the simulation. All calculated thermodynamic properties are summarized in one output file, which is straight forwardly readable and self-explaining. For a more detailed evaluation of the simulation, the instantaneous and running averages of the most important thermodynamic properties are written to other files. In total, the current status of the simulation and many more details are written to six files and can be accessed during execution.

4.2. Reduced quantities

The simulation program *ms2* internally uses reduced quantities for its calculations. All quantities are reduced by a reference length σ_{R} , a reference energy ϵ_{R} and a reference mass m_{R} , respectively. These reference values are input variables and may, in principle, be chosen arbitrarily. However, it is recommended to use a reference length σ_{R} in the order of 3 Å, a reference energy $\epsilon_{\text{R}}/k_{\text{B}}$ in the order of 100 K and a reference mass m_{R} in the order of 50 atomic units. The reduction scheme for the most important physical quantities is listed in the Supplementary Material. From these properties, the reduced form of all other quantities can be derived. An exception in the reduction scheme is the chemical potential, which is normalized in *ms2* by $k_{\text{B}}T$ instead of ϵ_{R} .

4.3. Molecular positions and orientations

The simulation program *ms2* updates only the positions of the centers of mass and the orientations of the molecules. From these values, the coordinates of the sites are derived on the fly. The center of mass of each molecule is stored in Cartesian coordinates. The absolute positions are reduced by the reference length σ_{R} and scaled by the reduced edge length of the cubic simulation volume, leading to position values in the range of

$-0.5 < x, y, z < 0.5$. The advantage of this scaling scheme lies mainly in the efficient application of the periodic boundary condition [36] and the minimum image convention [37].

The molecular orientations are stored using normalized quaternions [29]. Quaternions are a biunique representation that avoid divergence problems at low angles. Their use allows for an efficient calculation of site positions, while being less demanding concerning execution time and memory than updating and storing all site positions.

4.4. Initial configuration

It is important to define a stable and physically reasonable starting configuration for a simulation in a reliable way. In *ms2*, the molecules are initially placed on a face-centered cubic lattice in order to avoid overlaps between molecules. For mixtures, the positions of the different molecule species are distributed randomly on the lattice, ensuring a homogeneous distribution of all components in the simulation volume.

After the initial molecule placement, the configuration can be relaxed by translational and rotational MC moves. This step is part of the initialization process and can therefore be performed regardless of the simulation technique used. The number of relaxation moves is user defined and should be chosen large enough to achieve a physically reasonable starting configuration. For MD, the molecules are subsequently assigned with initial velocities such that the temperature is specified and no net translational and rotational moment is present. It is recommended to continue with a MD equilibration until a physically reasonable configuration and distribution of velocities is achieved. The MC loops relax possible overlaps in the initial configuration, whereas the MD equilibration drives the system into a physically reasonable dynamic microstate. Note that the equilibration process does not contribute to the calculation of the thermodynamic properties.

4.5. Intermolecular interactions

4.5.1. Dispersive and repulsive interactions

In *ms2*, the dispersive and repulsive interactions between molecules are reduced to pairwise interactions u_{ij}^{LJ} of the different molecule sites i and j , which are modeled by the 12-6-LJ potential [49]

$$u_{ij}^{\text{LJ}} = 4\varepsilon \left(\left(\frac{\sigma}{r_{ij}} \right)^{12} - \left(\frac{\sigma}{r_{ij}} \right)^6 \right). \quad (13)$$

Here, the site-site distance between two interacting LJ sites i and j is denoted by r_{ij} , while σ and ε are the LJ size and energy parameters, respectively. The LJ potential is widely used and allows for a fast computation of these basic interactions. It has only two parameters, which facilitates the parameterization of molecular models.

For pure components, the interactions between two different LJ sites are described by the Lorentz-Berthelot combination rules [50, 51]. For mixtures, the combination rules are extended to the modified Lorentz-Berthelot rules, which include two additional parameters η and ξ to describe the interactions between LJ sites of unlike molecules [52]

$$\sigma_{ij} = \eta \frac{\sigma_i + \sigma_j}{2}, \quad (14)$$

$$\varepsilon_{ij} = \xi \sqrt{\varepsilon_i \varepsilon_j}. \quad (15)$$

The two parameters scale all LJ interactions between molecules of different components equally. Arbitrary modifications of the combination rules are possible. *ms2* allows the specification of η and ξ for every molecule pair independently and thus the free parameterization of the combination rules.

4.5.2. Electrostatic interactions

ms2 considers electrostatic interactions between electro-neutral molecules. Three different electrostatic site type models are available: point charge, point dipole and linear point quadrupole. The two higher order electrostatic interaction sites integrate characteristic arrangements of several single partial charges on a molecule. The simulative advantage of higher order polarities is faster execution and a better description of the electrostatic potential of the molecule for a given number of electrostatic sites [53]. Combining two partial charges to one point dipole reduces the computational effort with *ms2* by about 30%, combining three point charges to one point quadrupole reduces the computational demand by even 60%. An additional advantage is the simplification of the molecular model, which reduces the number of molecular model parameters. In the following, the implemented electrostatic interactions between sites of the same type are briefly described. The interaction potentials between unlike electrostatic site types is presented in the Supplementary Material.

Point charge. Point charges are first order electrostatic interaction sites. The electrostatic interaction between two point charges q_i and q_j is given by Coulomb’s law [29]

$$u_{ij}^{qq}(r_{ij}, q_i, q_j) = \frac{1}{4\pi\epsilon_0} \frac{q_i q_j}{r_{ij}}. \quad (16)$$

This interaction decays with the inverse of r_{ij} and is therefore significant up to very large distances. For computational efficiency, the Coulombic contributions to the potential energy are explicitly evaluated up to a specified cut-off radius r_c . The long range interactions with charges beyond r_c are corrected for by the reaction field method. Its application to point charges is discussed below.

Point dipole. A point dipole describes the electrostatic field of two point charges with equal magnitude, but opposite sign at a mutual distance $a \rightarrow 0$. Its moment μ is defined by $\mu = qa$. The electrostatic interaction between two point dipoles with the moments μ_i and μ_j at a distance r_{ij} is given by [29, 43]

$$u_{ij}^{DD}(r_{ij}, \theta_i, \theta_j, \phi_{ij}, \mu_i, \mu_j) = \frac{1}{4\pi\epsilon_0} \frac{\mu_i \mu_j}{r_{ij}^3} (\sin \theta_i \sin \theta_j \cos \phi_{ij} - 2 \cos \theta_i \cos \theta_j), \quad (17)$$

with θ_i being the angle between the dipole direction and the distance vector of the two interacting dipoles and ϕ_{ij} being the azimuthal angle of the two dipole directions, cf. Figure 2. In *ms2*, the interaction between two dipoles is explicitly evaluated up to the specified cut-off radius r_c . The long range contributions are considered in *ms2* by the reaction field method [29].

Linear point quadrupole. A linear point quadrupole describes the electrostatic field induced either by two collinear point dipoles with the same moment, but opposite orientation at a distance $a \rightarrow 0$, or by at least three collinear point charges with alternating sign ($q, -2q, q$). The resulting quadrupole moment Q is defined by

$Q = 2aq$. The electrostatic interaction between two linear point quadrupoles with the moments Q_i and Q_j at a distance r_{ij} is given by [43, 54]

$$u_{ij}^{QQ}(r_{ij}, \theta_i, \theta_j, \phi_{ij}, Q_i, Q_j) = \frac{1}{4\pi\epsilon_0} \frac{3}{4} \frac{Q_i Q_j}{r_{ij}^5} \left[1 - 5((\cos \theta_i)^2 + (\cos \theta_j)^2) - 15(\cos \theta_i)^2(\cos \theta_j)^2 + 2(\sin \theta_i \sin \theta_j \cos \phi_{ij} - 4 \cos \theta_i \cos \theta_j)^2 \right], \quad (18)$$

where the angles θ_i , θ_j and ϕ_{ij} indicate the relative angular orientation of the two point quadrupoles, as discussed above. Note that no long range correction is necessary for this interaction type if the fluid is isotropic.

Reaction field method. The truncation of interactions between first and second order electrostatic sites leads to errors that need to be corrected. The reaction field method [34, 55] is implemented in *ms2* for this task, being widely used and well accepted [56, 57]. Its advantages are accuracy and stability, while requiring little computational effort compared to other techniques like Ewald summation [58]. However, it is limited to electro-neutral systems, thus *ms2* is currently restricted to electro-neutral molecules. The basic assumption of the reaction field method implemented in *ms2* is that the system is sufficiently large so that tinfoil boundary conditions ($\epsilon_s \rightarrow \infty$) are applicable without a loss of accuracy. This is the case for $N \geq 500$ [59, 60, 61, 62].

4.6. Cut-off modes

ms2 supports two different cut-off modes: the site-site cut-off and the center of mass cut-off. The site-site cut-off mode explicitly considers the interactions between all sites that are within a distance of r_c . Beyond r_c , the long-range contributions to the energy and pressure are estimated by analytical functions assuming a homogeneous fluid [29]. A disadvantage of the site-site cut-off arises, if molecular models contain point charges. In many cases close to r_c , molecules are only partially considered in the explicit calculation, cf. Figure 3. The point charges within the cut-off radius may be unbalanced so that an overall charge within the cut-off sphere might occur. For this condition, the reaction field is not valid.

A more robust alternative is the center of mass cut-off mode. It considers all interaction sites of different molecules explicitly, if their molecular centers of mass are within the cut-off radius. The long-range contributions beyond r_c due to LJ interactions are approximated by the formulations of Lustig [63]. Note that the computational advantage of the center of mass cut-off mode increases with the number of sites per molecule.

4.7. Monte-Carlo algorithm

Thermodynamic properties are determined by MC simulation via a Markov chain of molecular configurations. In *ms2*, this Markov chain is generated by executing a loop of N_{MC} moves per configuration, where N_{MC} is defined by

$$N_{MC} = \frac{1}{3} \sum_{i=1}^N N_{i,DGF}. \quad (19)$$

Here, N represents the number of molecules in the system and $N_{i,DGF}$ is the number of degrees of freedom of molecule i . For MC simulations, three different moves are implemented in *ms2*, translational and rotational

displacement of a single molecule and fluctuation of the simulation volume. All three moves are well known and widely in use [29] and thus not further discussed here.

The acceptance of MC moves is associated with energy differences before and after that perturbation. To speed up execution, it is avoided in *ms2* to calculate the old energy for each attempted move. Instead, the pairwise interactions of all molecules with all the remaining $N - 1$ molecules in the system are stored in a $N \times N$ matrix. The sum of each column and row of that matrix, respectively, equals the potential energy of one particular molecule. If a molecule i is assigned to be moved, its old energy is determined by simply summing up all contributions in column i of the energy matrix. After the move, its new energy is determined by calculating the pairwise interactions with the remaining molecules and the individual contributions are stored in a vector of dimensions $N \times 1$. If the move is accepted, the vector replaces column i as well as row i of the energy matrix. If the move is rejected, the vector is discarded. The same technique is used for the virial contributions.

5. Implementation

In this section, the code design and implementation issues of *ms2* are discussed. These explanations of the source code are intended to help understanding the program and to encourage the reader to further develop and improve the code. New developers should get a smooth entry and the possibility to make fast adaptations to new problems. Due to the experience of the core developers and the suitability for an efficient numerical code, Fortran90 was chosen as programming language. The program makes extensive use of the concepts introduced with Fortran90, notably modular programming. A wide variety of different simulation setups is possible, each requiring different computations. This variety leads to a need for a highly modular structure, where modules have clear-cut interfaces. Given such modularity, code parts that are not needed for the user defined simulation setup, can be simply skipped within the calculation. Besides avoiding if-statements at computationally demanding points and therefore leading to a better data flow, this allows for an efficient implementation of new functionalities, as existing modules can be used whenever feasible. In *ms2*, a stringent philosophy regarding modularization was followed, which is discussed in section 5.1. The runs need to be fast and scale well with increasing number of processors in order to reduce response times. *ms2* is parallelized with the Message Passing Interface (MPI) [64]. While MD needs synchronization after every time step due to its deterministic nature, MC can be executed embarrassingly parallel, requiring synchronization only at the very end of a run. Furthermore, the by far most expensive part of the code, calculating forces and potential energies, is highly vectorized and optimized.

5.1. Modular structure

Goals.

The simulation program *ms2* is intended to be flexible, expandable and easy to get familiar with. Still, the program has to be efficient and effective in its calculations, such that the employed computer resources are efficiently used and therefore, the response times are minimized. The introduction of Abstract Data Types (ADT) is a popular approach to reduce the complexity of a software system, decomposing it into smaller subunits, which are easier to maintain. The ADT model, with its data and code abstraction, forms the basis of Object-Oriented

Programming (OOP). It should be mentioned that OOP holds challenges regarding efficiency. Distributing the data between objects can lead to data fragmentation and indirections can occur, leading to a negative impact on the performance. One objective was to introduce OOP concepts, while still maintaining efficiency, cf. section 6.

Fortran issues.

Fortran90 modules are a further development of Fortran77 common blocks. A module may not only contain data, but also define subroutines and user derived data types. Due to the fact that module data and associated module subroutines implement singletons from the start, grouping module-defined user derived data types and associated functions also allows an OOP-like programming style. A small example will demonstrate this. The module `ms_simulation`, defined in the file `ms_simulation.F90`, contains the type `TSimulation` as well as the subroutines

- `TSimulation_Construct (this)`
- `TSimulation_Destruct (this)`
- `TSimulation_RunSteps (this , StepStart , StepEnd)`

among others. The strict naming convention helps avoiding clashes in name space. All the subroutines above receive a reference to an instance of `TSimulation` through the first parameter, named *this*, and they serve as (default) constructor, destructor or ordinary member function, respectively. In the following, constructs like *this* will be referred to as a class. `ms2` does not strictly follow an OOP approach, because all data elements can be accessed directly, therefore there is no information hiding through data encapsulation.

Class structure.

In order to keep the class hierarchy flat, `ms2` does not make use of inheritance and hence polymorphism. E.g. the module `ms2_site` contains the classes `TSiteLJ126`, `TSiteCharge`, `TSiteDipole` and `TSiteQuadrupole`, which have data and functionalities in common, but are not derived from a common class. The class structure and the relations between the classes are organized in six levels, as depicted in Figure 4. Every class is located on a certain level and hence has data and modules relevant to its level. The six levels are intuitive:

1. Global – all data and functions of global use are handled. Scope: whole code
2. Simulation – setup and control of the simulation. Scope: simulation framework
3. Ensemble – the required ensemble is set up and initialized. Scope: molecular ensemble
4. Component – data and functions dealing with all molecules of a given molecule type. Scope: one molecule species
5. Molecule – data regarding the basic structure of a molecule. Scope: one molecule
6. Site – position of sites with respect to a molecule. Scope: one site of one molecule

All subroutines on every level are task specific and implemented as modules. Depending on the simulation setup, e.g. MD/MC or ensemble, the appropriate modules are engaged individually during simulation. The modules of `ms2` can coarsely be divided into following tasks:

- initialization
- simulation
- accumulation
- output

If a task spans more than one level, the corresponding module on the topmost level is called, which in turn calls modules on lower hierarchical levels and so on.

First level - global. The first level contains subroutines and variables that are needed globally in all parts of the simulation. The variables are mainly natural constants like the Boltzmann constant k_B , the Avogadro constant N_A etc. Additionally, variables defining the simulation setup for the given run are stored, e.g. the simulation technique and the ensemble type. These variables determine the type of calculation and are assigned according to user specifications to the levels below in *ms2*. Subroutines that are stored on this level are also of global use in *ms2*. These are basic routines for input into the simulation program and output into files, as well as more advanced routines for automated communication between hardware and program via signals. The latter routines are of particular interest for running *ms2* on high performance computers, where possible data loss can be avoided by automated communication between cluster nodes and the program.

Second level - simulation. On the second level, the simulation flow is controlled. The simulation is initiated, started and sustained. This includes controlling the length of the simulation and its output. User specifications concerning the simulation setup are read and the corresponding global variables are assigned accordingly. Furthermore, all simulation quantities and averages are analyzed and written to file.

Third level - ensemble. On the third level, the simulation is organized according to the specified ensemble. The respective ensemble is initialized, by characterizing the ensemble setup and defining e.g. the volume and the composition of the molecular species in the simulation volume. In addition, the initial positions and velocities of all molecules are distributed. Furthermore, the routines for the subsequent simulation are assigned on this level. The time integration or the Markov chain, respectively, is performed and the ensemble averages of thermodynamic properties are calculated during the simulation run.

Fourth level - component. On this level, the molecule species (component) and their interactions are addressed. The class structure has three distinct branches. While the first branch deals with issues regarding structure and properties of one component at a time, the second branch deals with interactions between molecules of one or two components. The third branch contains all routines for the accumulation of data, determining the average value of the data as well as their statistical uncertainties. The first two branches will be discussed in more detail. The first branch deals with the structure of one individual component. An instance of the class component stores center of mass position, orientation, velocity etc. of every molecule of one species, e.g. H_2O . The contained modules deal with calculations of all molecules of this species, e.g. modules for solving Newton's equations of

motion, kinetic energy calculation, atom to molecule transformation and vice versa. Furthermore, this branch is responsible for the initialization. It reads the molecular model for every component, calculates positions, introduces the reduced quantities etc.

The second branch deals with interactions between molecules of the same or different type. Therefore, it contains all force and potential energy calculations, which are the computationally most expensive parts by far. This branch is highly vectorized to achieve a fast calculation, in particular if executed in vector-parallel. The modules in this branch are simulation technique specific, i.e. MD or MC. In a MD simulation, first, the interaction partners of every molecule are determined, i.e. other molecules or sites within the cut-off radius. Then the interaction energies and virial contributions as well as the resulting forces and torques are calculated and summed up for every molecule.

In a MC simulation, the modules are designed to calculate only the energy and the virial. Forces are not evaluated and the corresponding algorithms are therefore entirely omitted in these modules. For MD, all interactions of one component-component pair are calculated in one call of the module. For MC, the modules are finer-grained, calculating all interactions of one component-component pair by determining every molecule-molecule interaction individually. The corresponding modules are thus called more frequently than in case of MD.

Fifth level - molecule. The basic information on a molecular model is stored on this level, such as mass, moment of inertia tensor, rotational degrees of freedom etc. Higher level modules access these data via the class molecule.

Sixth level - site. Here, the assembly of a molecule is stored. The data are used by higher level modules to determine the site positions from the molecular positions and orientations. This allows for the calculation of site-site interactions. It should be noted that *ms2* only integrates (MD) or accepts (MC) center of mass positions and orientations. Site specific data are calculated for every simulation step on the fly.

5.2. Parallelism

The source code is implemented for an efficient parallelism with distributed memory, using the MPI [64] standard for communication. The program uses the following MPI calls:

- `MPI_Init`, `MPI_Finalize`, `MPI_Comm_rank`, `MPI_Comm_size` to set up basic MPI functionality
- `MPI_Abort` to stop the program in case of an error
- `MPI_Barrier`, `MPI_Wtime` and `MPI_Wtick` within the stopwatch class
- `MPI_Bcast`, `MPI_Reduce`, `MPI_Allreduce` for simulation data exchange

Exclusively collective communication is used, an approach suitable for molecular simulations, where the cut-off radius is in the order of the simulation volume edge length. For MD simulations with *ms2*, only the force calculation is parallelized using molecule decomposition according to Plimpton [14]. The interaction matrix is rearranged such that the number of interaction partners is almost equally distributed in the matrix. This is achieved by calculating the interactions between molecules 1 and N as interactions between molecules N and

1, cf. Figure 5. The parallelization scheme then distributes the calculation among N_P processors with an almost equal work load, such that every processor executes N/N_P rows of the interaction matrix. The master process reduces then all the force components to sum up the resulting molecular forces on each molecule.

The MC code is parallelized by exploiting the stochastic nature of the simulation method. Starting from an equilibrated state in the simulation volume, the molecular configuration is copied multiple times into different volumes. Then, each copy runs independently in parallel, using a different random number seed, to calculate the thermodynamic properties. At the end, the data from all copies are gathered and averaged.

5.3. Vector-based structure

The program *ms2* specifically accounts for vector-parallel machines. All the main information needed in the computationally costly inner loops of the calculation, like the positions of the molecules and their sites, are stored in vectors that are accessed sequentially. This allows for an efficient use of memory on vector-parallel machines.

For the evaluation of one configuration in *ms2*, e.g. the position vectors are loaded once and then distributed among all processors. Then, additional information like interaction partners of each individual molecule is computed, tabulated in vectors and communicated. The order of this data follows the order of all other vectors, e.g. the position vector. The appropriate order allows for a serial access to the data in all vectors in the random access memory. Therefore, the calculations can be vectorized, increasing the efficiency.

6. Benchmarks

A subset of the *ms2* examples, which come along with the code release, was used for profiling and runtime tests that are presented here. The number of time steps (MD) or loops (MC), respectively, considered was lower than for a normal production run, since only comparisons were carried out. The MD and MC test case simulates the equimolar liquid mixture of methanol and ethanol at 298.15 K and 0.1 MPa in the NpT ensemble [65]. Methanol was modeled by two LJ sites and three point charges and ethanol by three LJ sites and three point charges [66]. This test case was chosen, because it is a typical application. Similar results are expected for a wide class of problems. However, note that the actual run times will differ with varying thermodynamic conditions that are simulated. Run times increase with higher density of the system, among others.

6.1. Sequential version

Compiler. *ms2* is distributed as a source code and can be compiled with a wide variety of compilers. However, the performance is significantly influenced by the compiler and linker as well as the used options. Figure 6 shows the runtime of the test case on different platforms using different compilers. The binaries were generated by

- GNU gfortran² with `”-fdefault-real-8 -O3”`
- Intel ifortran³ with `”-r8 -fast”`

²<http://gcc.gnu.org/fortran/>

³<http://software.intel.com/en-us/intel-compilers/>

- PGI pgf95⁴ with `”-r8 -fastsse”`
- Sun Studio sunf90⁵ with `”-r8const -fast”`
- Pathscale pathf90⁶ with `”-r8 -Ofast”`

with options activated to enforce the use of double precision floating point numbers. The chosen optimization flags represent a common choice.

The best combination of compiler and platform was the Intel ifortran compiler and the Intel Xeon E5560 `”Nehalem”` processor⁷. An Intel ifortran compiled binary of *ms2* significantly outperforms binaries compiled by GNU gfortran and PGI pgf90, independent on the computing platform.

Profiling. The test case was chosen for profiling studies running the sequential version of *ms2* with valgrind using the callgrind tool⁸. The binary was built by the Intel ifortran 11.1 compiler and optimized for SSSE3 to work with valgrind 3.5.0. From the resulting estimated CPU cycles, the computational hot spots are presented for both MD (Figure 7a) and MC (Figure 7b) simulations. The calculation of forces and energies consumes most of the CPU cycles. For the MD run, 54.33% of the CPU time is spent for charge-charge interactions and 36.09% for LJ-LJ interactions, which adds up to 90.42%. Further program routines, e.g. determining the interaction partners (7.21%), add up to a total of 99.70% of the CPU time spent for the entire calculation of potential energies and forces. In MC calculations, 95.68% of the time is spent for energy calculations, including the routine for determining the interaction partners (4.28%).

With valgrind set up to use a 32768 Byte, 4-way associative level 1 data cache with 64 Byte cachelines and prefetch enabled, the cache miss rate for reading (D1r) was below 5%, cf. Table 2. The misses occur mainly in the force calculation for MD (99.19%) and the energy calculation for MC (99.07%).

The ratio of read and write access rates is larger than six, therefore, level 1 reading cache misses (D1mr) have a larger impact on the performance than level 1 writing cache misses (D1mw) and level 2 cache misses. Instruction cache misses are negligible here.

6.2. Vectorization

ms2 MD simulations also run efficiently on vector machines like the NEC SX architectures. Table 3a shows profiling data running the MD test case on a NEC SX8R using `”ftrace”`. The force calculation routines reached an average vector length of about 170 and 1892 MFLOPS, which is good, especially considering that the test case is not favorable regarding vectorization. The MC algorithm, however, is not well suited for vectorization, cf. Table 3b. The NEC SX8 shows comparable results, proportional to the lower processor clock speed, whereas results for the NEC SX9 are below expectation with respect to the higher hardware peak performance and memory access speed. For comparison, the same test case was executed on a mainstream Intel Xeon E5440 `”Harpertown”`

⁴<http://www.pgroup.com/products/>

⁵<http://developers.sun.com/sunstudio/>

⁶<http://www.pathscale.com/>

⁷<http://www.hlrs.de/systems/platforms/nec-nehalem-cluster/>

⁸<http://valgrind.org/info/tools.html#callgrind>

2.83 GHz system, reading hardware counters with PAPI. The MD version achieved an average of 1461 MFLOPS and the MC version 1374 MFLOPS on this platform.

6.3. Parallelization

ms2 supports parallel systems with distributed memory using the MPI standard for communication. Taking a look at the MPI routines used to exchange data, excluding the *ms2_global* and *ms2_stopwatch* modules, solely three different MPI calls are employed: `MPI_Bcast` (42x), `MPI_Reduce` (13x) and `MPI_Allreduce` (12x). As a result of the profiling in section 6.1, only the force calculations were parallelized for MD and therefore a master process has to reduce all force components to sum up the resulting molecular forces. MC calculations use this kind of parallelization only in the equilibration phase. During production, the phase space is sampled with fully independent random walks on each processing unit, generating independent Markov chains, which have to be consolidated only once at the end of the simulation.

The *ms2* parallelization capabilities were tested on a "Nehalem" dual-Quadcore-CPU cluster with Infiniband, using the Intel ifortran compiler and OpenMPI V1.4 as well as Intel MPI V4.0. The test case runtime results for a range of commonly used numbers of processors are shown in Figure 8. Using more than a single node, i.e. 8 processes here, Intel MPI showed a better scalability than OpenMPI. While the MD scaling behavior is good for a decent number of processors, MC production steps can be characterized as optimally parallel. In contrast to this, the MC equilibration steps, parallelized using a different strategy, show the expected scaling characteristics of MC simulations. A closer look at the time spent for communication shows the differences between equilibration and production for a parallel MC run with 32 processors. After the equilibration, which consumes roughly 23 seconds in this case, no communication takes place, cf. Figure 9a.

For MD, MPI communication is necessary throughout the whole simulation, with communication phases in each time step. Figure 9b shows the communication needed for a single time step, where the barrier induced by the first collective call results in a waiting time for faster processes and particularly for the last process, which might be notably faster than the others, cf. Figure 9c. This imbalance of the last process is because the force calculation for N molecules is divided among the N_p processes, where each processor is assigned with $\lceil N/N_p \rceil$ molecules, except the last one to which the rest is assigned. Regarding the MD test case with 1372 molecules in Figure 9, 31 processes each calculate the forces for 43 molecules, while the last process only deals with 39 molecules. This load imbalance is not a serious issue, since also with a different molecule distribution, where the number of molecules on each processor only differs by one, there is still a processor with 43 molecules, which dictates the overall execution time. The fraction of the time spent in MPI routines increased to about 10% for 32 processes from 5% for 16 processes for the test case. Within the collective communication, process 0 is the master and therefore exhibits a different communication pattern, cf. Figure 9.

6.4. Memory

The memory demand of *ms2* is low, e.g. for a simulation of 1372 molecules with five interaction sites, roughly 200 MB of RAM is required. For parallel execution, the concept of distributed memory is used. Here, the

memory scales linearly with the number of processors that are used. For most applications of *ms2* a total RAM of 2 GB is sufficient.

6.5. Comparison to other codes

Monte-Carlo. The performance of *ms2* was evaluated against the simulation program *MCCCS Towhee V6.2.7* [67], which is widely in use and well accepted by the scientific thermodynamic community. In a first step, the comparison was restricted to a simple *NVT* simulation of the test case. *ms2* and *MCCCS Towhee*, both compiled with the Intel ifortran compiler, executed the MC test case sequentially, i.e. with one process only, on a "Nehalem"/Infiniband cluster with an equal number of MC moves. A sequential run was executed since *MCCCS Towhee* is not yet prepared for parallelization. The performance of *ms2* was faster than the *MCCCS Towhee* program for this simulation by a factor of around 20.

An important application of *ms2* is the determination of VLE data. The present comparison was restricted to two other programs: *MCCCS Towhee* [67] and the *Errington* Gibbs ensemble Monte-Carlo code [68]. Both codes use the Gibbs ensemble approach for the determination of VLE data, a simulation scenario that is not yet supported by *ms2*. Instead, the grand equilibrium method was used in *ms2* for determining the VLE.

The test case was again a mixture of methanol and ethanol at $T = 393$ K, using $N_l = 1372$ molecules for the liquid phase and $N_v = 500$ molecules for the gas phase. These numbers were also used as starting values for the Gibbs ensemble calculations. For *MCCCS Towhee* and the *Errington* Code, three different types of moves were allowed, translational and rotational displacement, volume exchange and molecule transition between the two simulation volumes. The moves were chosen with probabilities of 79%, 1% and 20% respectively, which are typical numbers for simulations including dense liquid phases. In *ms2*, the number of moves was restricted to two in the liquid phase, translational and rotational displacements (N_l times) and one volume fluctuation each loop. For the gas phase simulation, a move for insertion and deletion of molecules was additionally invoked twice in each loop. The overall execution time was set to 96 hours execution time, which is about the maximal time that is currently provided for a single simulation in high performance computing centers. The runs were performed sequentially, since *MCCCS Towhee* and the *Errington* code are not suited for parallel execution. For this test case, *ms2* showed the best performance. Within the 96 h, *ms2* calculated 52000 loops in the liquid phase simulation, and 50000 loops in the vapor phase simulation. The *Errington* code performed second best in this test, running 22000 loops, followed by *Towhee*, which ran 6000 loops.

The performance of the codes influenced the quality of the simulation results drastically. For the VLE test case, the reference phase equilibrium data taken from the literature comprise a vapor pressure of 0.58 MPa, a saturated liquid density of 17.22 mol/l and a saturated vapor density of 0.19 mol/l. *ms2* reproduced the experimental results well in the given time. The calculated vapor pressure of *ms2* was 0.53 MPa, while the saturated liquid density of the mixture was predicted to be 17.54 mol/l and the saturated vapor density to be 0.18 mol/l. The statistical uncertainties for the quantities was acceptably low, being 0.01 MPa for the vapor pressure and 0.03 mol/l and 0.003 mol/l for the saturated liquid density and saturated vapor density, respectively. The *Errington* code calculated in the given time a vapor pressure of 0.19 MPa well as a saturated liquid density of 19.93 mol/l and

a vapor density of 0.443 mol/l. Using *Towhee*, the calculation yielded a vapor pressure of 1.81 MPa as well as densities of 13.62 mol/l and 0.64 mol/l, respectively. Here, the standard deviations of the results were significantly higher. These simulation results showed a drift, e.g. for the saturated liquid density by an average of 1.50 mol/l over 1000 steps. The results for the *Errington* Code and *MCCCS Towhee* indicate that the simulations with these programs had not reached the phase equilibrium in 96 h. This effect is much more pronounced for *MCCCS Towhee*, having executed significantly less loops than the *Errington* code.

Molecular dynamics. The performance of the MD part of *ms2* was evaluated against the simulation program *Gromacs V4.0.3* [69], which is designed for simulations of biological systems. The comparison was restricted to a simple *NVT* simulation of the equimolar liquid mixture of ethanol and methanol at $T = 298$ K and $p = 0.1$ MPa. The runs were executed sequentially on a "Harpertown"/Infiniband cluster with an equal number of time steps and the same interaction cut-off radius of 21 Å. Both programs, *ms2* and *Gromacs*, were compiled with the Intel ifortran compiler. *Gromacs* was faster by almost a factor of two than *ms2* and also the scaling behavior for the parallel version was superior. This shows that there is still an optimization potential in *ms2*, which will be exploited in the future.

A comparison between *Gromacs* and *ms2* for the VLE test case was not carried out, since *Gromacs* does not support VLE simulations.

6.6. Computational demand for the calculation of transport properties

The computational demand for transport properties following the Green-Kubo formalism was evaluated. This investigation was performed for the test case equimolar liquid mixture of methanol and ethanol at a temperature $T = 298$ K and $p = 0.1$ MPa with an autocorrelation length of 13.82 ps. The time period between two autocorrelation functions was specified to be 197.4 fs. A total of 100 autocorrelation functions was explicitly written to file in order to check the results. The transport properties were evaluated every four autocorrelation functions. The reference case was a MD run performed with the same technical details, except for the calculation of the autocorrelation functions. Calculating transport properties increased the CPU time of a simulation step by about 78% compared to the reference case. Note that all autocorrelation functions were evaluated each time step, which is an extreme frequency that can be reduced without much loss of accuracy.

7. Features

The feature programs described in the following are intended to facilitate the handling of *ms2*. They should allow for an easy start of molecular simulations with *ms2* and give access to all output data generated by *ms2*.

7.1. Simulation setup: *ms2par*

The GUI-based feature tool *ms2par* allows for a convenient generation of input (*.par) files according to user specifications, cf. Figure 10. It assigns the molecular model, ensemble, thermodynamic state point, time step, number of equilibration steps etc. Furthermore, the tool proposes a cut-off radius based on the input data. After the user completes the specifications, the program generates the *.par-file in ASCII format, to be read by *ms2*.

Applying *ms2par* is simple and should be intuitive even for new users. *ms2par* is a java application, thus it runs on all operating systems.

7.2. Simulation analysis: *ms2chart*

ms2chart is a GUI-based java applet for evaluating the simulation results from *ms2*, cf. Figure 11. This tool plots the evolution of properties calculated by *ms2*. The properties can be shown for different graph axes that are chosen from a drop down menu and the plot is shown directly in the GUI. For a better analysis, *ms2chart* allows various features: plotting block averages as well as simulation averages into the same plot, changing the design of the plot, individual labeling of the axes and adjusting the frame detail. All plots can be exported in png format.

The analysis program *ms2chart* can be executed at any time of the simulation, i.e. also while the simulation is still running. There is no loss of data, if it is run on the fly. The tool *ms2chart* is easy to understand and can be handled intuitively. It allows for an easy first analysis of the *ms2* simulation results. Users will employ *ms2chart* for a quick check of their simulation results, such as whether the equilibration time was appropriate, which is important if extensive series of simulation runs are performed.

7.3. Visualization tool: *ms2molecules*

The program *ms2molecules* is a visualization tool for *ms2*. The program displays the molecular trajectories that are stored by *ms2* as a series of configurations in the *.vim file. *ms2molecules* visualizes molecular sites by colored spheres. The colors are user defined. The size of a sphere is by default proportional to the LJ size parameter σ , but can be changed manually in the first lines of the *.vim file. This feature facilitates monitoring a component of interest by reducing the size of other components, e.g. solvents. Other features like zooming into or out of the simulation volume as well as rotating the simulation volume also facilitate the analysis of simulations. The visualization can be exported via snapshots in the jpg format.

The program is based on OpenGL and written in C. The handling of the program is simple and console based. Requirements for the feature tool are OpenGL in a Windows or Linux environment. Figure 12 shows a snapshot of a ternary mixture, taken with the program *ms2molecules*, which is convenient to gain insight into the trajectory of the system or the state of the system, respectively.

8. Conclusions

The molecular simulation program *ms2* was designed for the calculation of thermodynamic properties of bulk fluids. Special care was given to a minimization of the response time. The capabilities of *ms2* are broad, ranging from basic static thermodynamic properties, like thermal and caloric data, over vapor-liquid equilibria to transport properties, like diffusion coefficient, viscosity and thermal conductivity. These data are accessible for pure substances and mixtures. The accuracy of the simulation data generated by *ms2* is high, while consuming a reasonable computational effort. Molecular models of more than 100 pure fluids are supplied with *ms2* that accurately describe their thermodynamic properties. Despite the fact that *ms2* is a sophisticated Fortran90 program, new developers benefit from its modular structure and object-orientation. The application of *ms2* is

straightforward, because the input files are well structured and auxiliary feature tools help to create input files, analyze simulation results and visualize molecular trajectories. The code was optimized for the current hardware technology and achieves a high efficiency.

Ongoing efforts focus on the implementation of new algorithms that extend the applicability of the program as well as its analysis tools. The two major current developments are the implementation of internal degrees of freedom to allow for the application of *ms2* to larger molecules and Ewald summation to allow for the application of *ms2* to charged molecules. The source code is available to the scientific community at <http://www.ms-2.de>.

9. Acknowledgments

The authors gratefully acknowledge financial support by the BMBF "01H08013A - Innovative HPC-Methoden und Einsatz für hochskalierbare Molekulare Simulation" and computational support by the Steinbruch Centre for Computing under the grant LAMO and the High Performance Computing Center Stuttgart (HLRS) under the grant MMHBF. The present research was conducted under the auspices of the Boltzmann-Zuse Society of Computational Molecular Engineering (BZS).

Table 1: Computational demand for the calculation of the chemical potential for a three-center LJ fluid using Widom’s test molecule insertion in MD or MC simulation. The calculations were performed on a single processor. The simulation time per time step (MD) was 0.043 s and 0.300 s per loop (MC), respectively. The simulations were performed with 500 molecules at a density of 0.23 mol/l. A cut-off radius of 60 Å was employed.

# test molecules	Additional CPU time per MD time step or MC loop / s
1	0.0015
500	0.1925
1000	0.3865
1500	0.5810
2000	0.7750

Table 2: *ms2* data cache behavior for the test case investigated with callgrind (D: data, 1m: level 1 cache misses, r/w: read/write access). The data misses are normalized by the total amount of data access.

	D1mr/Dr [%]	D1mw/Dw [%]	D2mr/Dr [%]	D2mw/Dw [%]	Dr/Dw
MD	2.41	0.34	0.0001	0.0002	6.18
MC	4.98	7.48	0.006	2.04	10.35

Table 3: NEC SX8R "ftrace" profiling. The five routines with the largest computational demand and their vectorization for the test case are shown for MD (a) and MC (b).

(a)

PROG.UNIT EXCLUSIVE TIME[sec]	(%)	AVER.TIME [msec]	FREQUENCY MOPS	MFLOPS	V.OP RATIO	AVER. V.LEN	VECTOR TIME	I-CACHE MISS	O-CACHE MISS	BANK CONF CPU	NETWORK
ms2_potential.tpotcc_force		4.713	126900								
598.058	(53.4)		8661.0	2164.7	99.53	169.3	579.369	0.012	3.383	3.487	13.433
ms2_potential.tpotljj_force		4.184	89300								
373.605	(33.4)		8039.8	1804.6	99.60	170.0	360.055	0.010	2.775	2.317	8.112
ms2_interaction.tinteraction_force		8.395	14100								
118.367	(10.6)		1749.6	429.2	80.68	172.6	36.326	0.025	39.243	0.092	0.961
ms2_interaction.tinteraction_calcpartners		1.929	14100								
27.200	(2.4)		14262.1	3412.8	99.28	191.2	25.795	0.021	0.091	0.169	0.830
ms2_component.tcomponent_restartsave		151.159	2								
0.302	(0.0)		341.9	7.7	0.14	3.3	0.003	0.048	0.010	0.000	0.001
total			648179								
1119.044	(100.0)	1.726	7855.1	1891.6	99.09	170.5	1002.287	0.265	45.647	6.157	23.568

(b)

PROG.UNIT EXCLUSIVE TIME[sec]	(%)	AVER.TIME [msec]	FREQUENCY MOPS	MFLOPS	V.OP RATIO	AVER. V.LEN	VECTOR TIME	I-CACHE MISS	O-CACHE MISS	BANK CONF CPU	NETWORK
ms2_interaction.tinteraction_energy		0.498	10155544								
5055.675	(98.3)		2388.1	563.9	69.23	162.3	635.244	22.679	827.141	6.305	45.971
ms2_component.tcomponent_mol2atom1		0.004	5604792								
25.066	(0.5)		305.0	41.3	0.00	0.0	0.000	3.765	7.560	0.001	0.453
ms2_ensemble.tensemble_energy1		0.004	3430000								
13.135	(0.3)		1271.2	492.5	87.92	237.8	2.329	1.990	4.316	0.001	0.163
ms2_ensemble.tensemble_move		0.006	1714383								
9.818	(0.2)		144.8	7.1	7.24	3.0	0.395	2.692	2.778	0.000	0.216
ms2_ensemble.tensemble_rotate		0.005	1715617								
8.658	(0.2)		141.9	9.9	0.00	0.0	0.000	2.491	2.313	0.000	0.199
total			44555359								
5142.345	(100.0)	0.115	2366.6	559.1	69.33	162.7	646.942	38.509	849.895	6.401	47.877

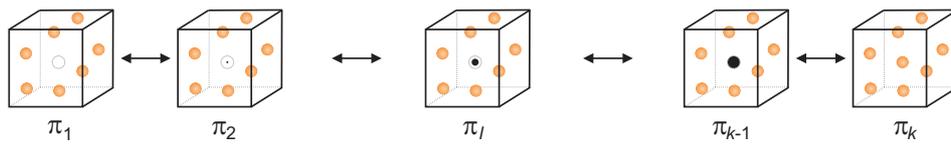


Figure 1: Schematic of the gradual insertion of a molecule. The molecule (dark gray) in state π_l fluctuates between state $l = 0$ (fully decoupled) and state $l = k$ (fully coupled).

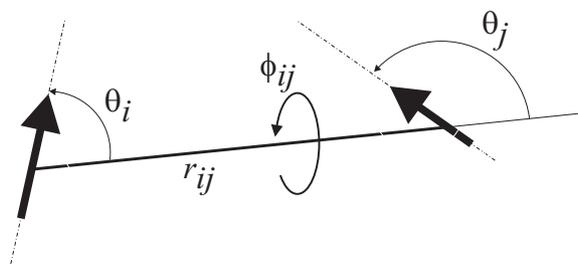


Figure 2: Schematic of the angles between two point dipoles i and j indicated by the arrows, which are situated on different molecules at a distance r_{ij} .

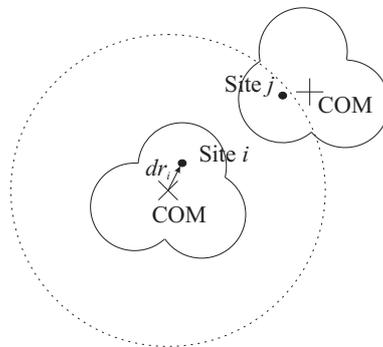


Figure 3: Schematic of the different cut-off modes. Using the site-site cut-off mode, the interaction between the two sites marked by a dot is explicitly evaluated, since they are within the cut-off radius indicated by the dotted line. The distance between the centers of mass (COM), marked by a cross, is irrelevant. Using the center of mass cut-off mode, none of the site-site interactions of the two molecules are explicitly evaluated, since the centers of mass are not within the cut-off radius. The distance between the center of mass and one particular site of molecule i is defined as dr_i .

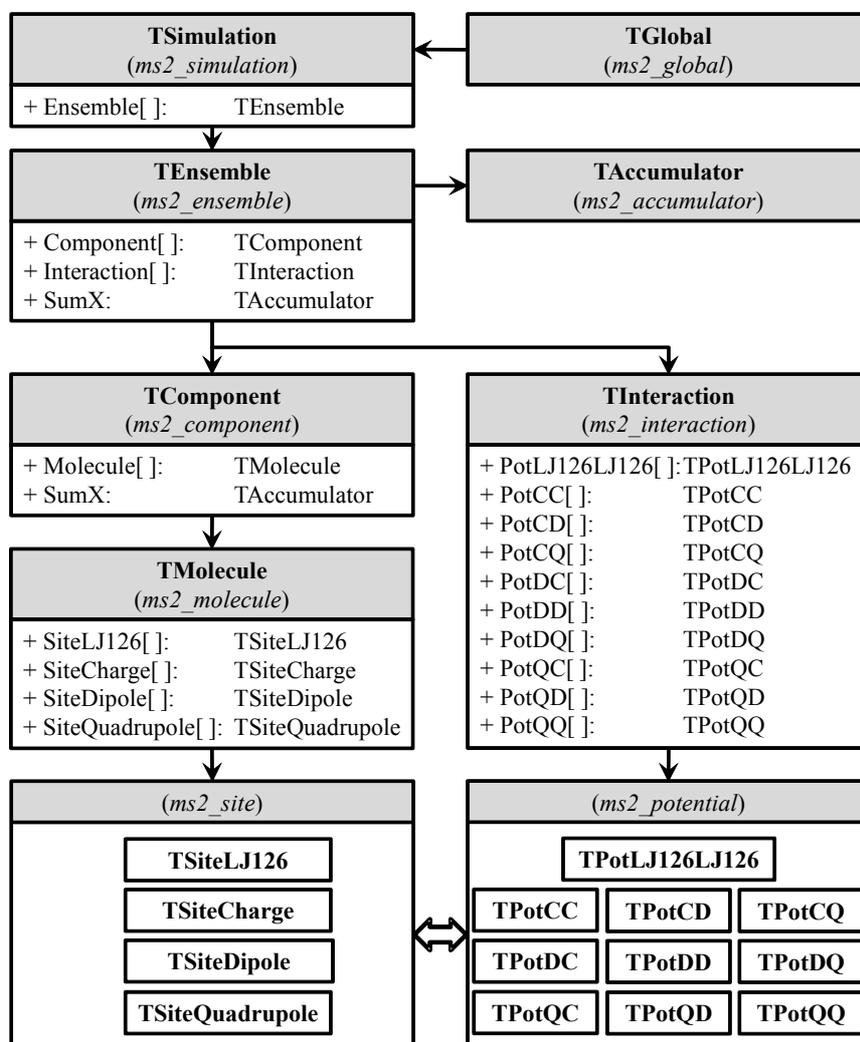


Figure 4: UML class diagram of *ms2*, showing object attributes. The bold headers specify the class name, while the italic names correspond to the program files, the class is implemented. The arrows indicate the command structure, while the bold arrows indicate information transfer. For a better understanding of the structure, following abbreviations are introduced into the diagram: CC: ChargeCharge; CD: ChargeDipole; CQ: ChargeQuadrupole; DC: DipoleCharge; DD: DipoleDipole; DQ: DipoleQuadrupole; QC: QuadrupoleCharge; QD: QuadrupoleDipole; QQ: QuadrupoleQuadrupole

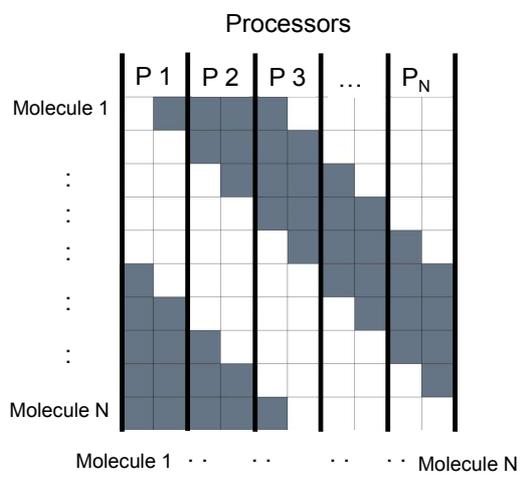


Figure 5: Schematic of the molecule decomposition parallelization method by Plimpton [14]. The colored boxes signify the pair interactions that need to be determined in order to proceed the simulation. The black lines define the range of molecules each processor is assigned to calculate their interactions.

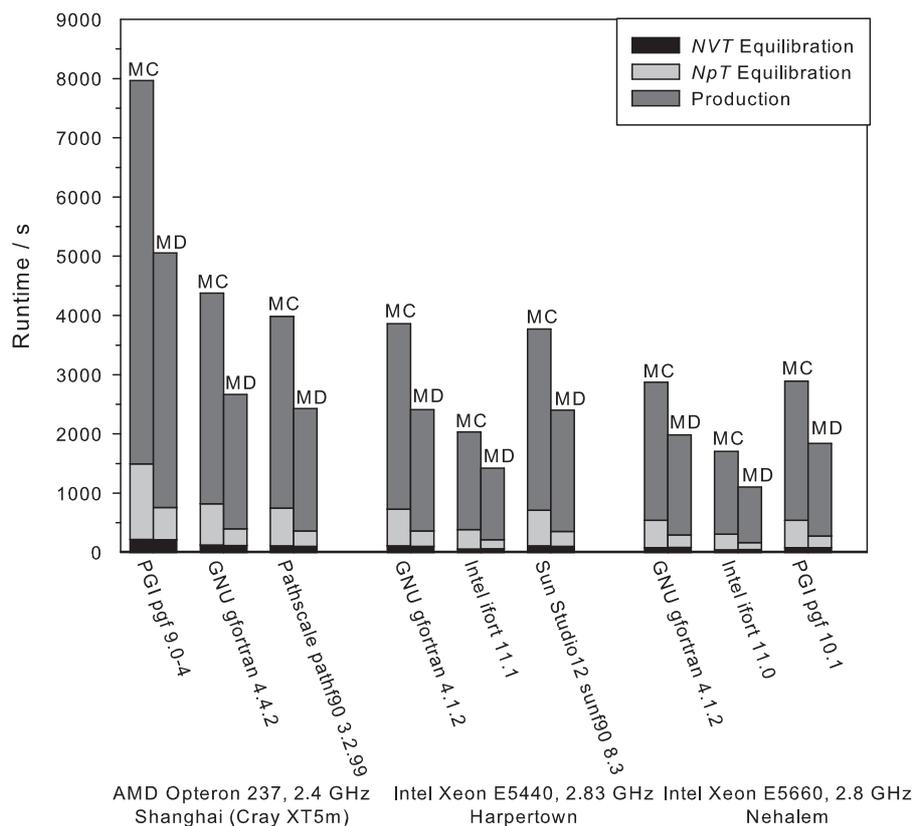
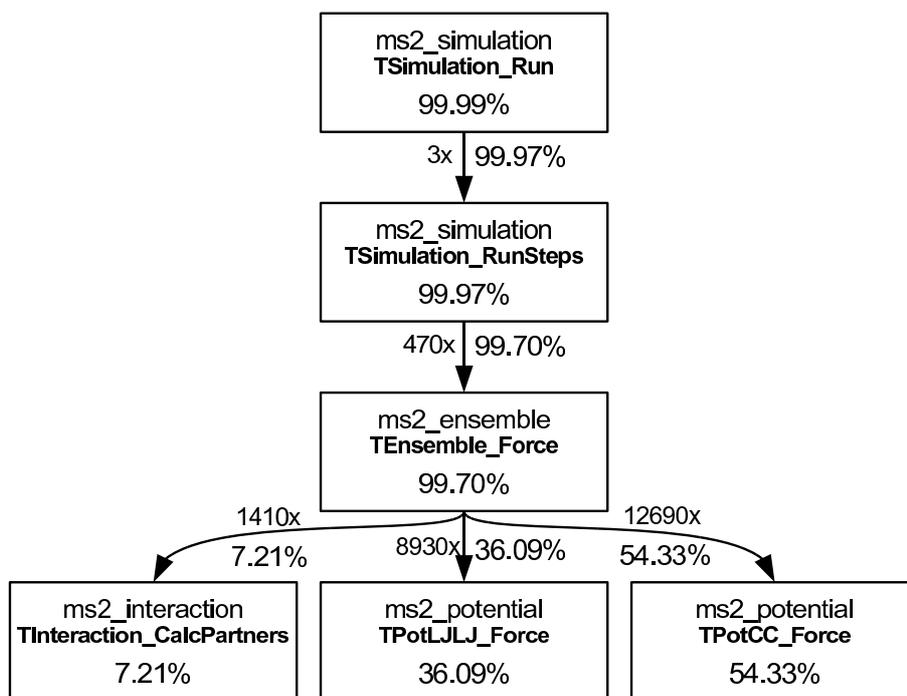
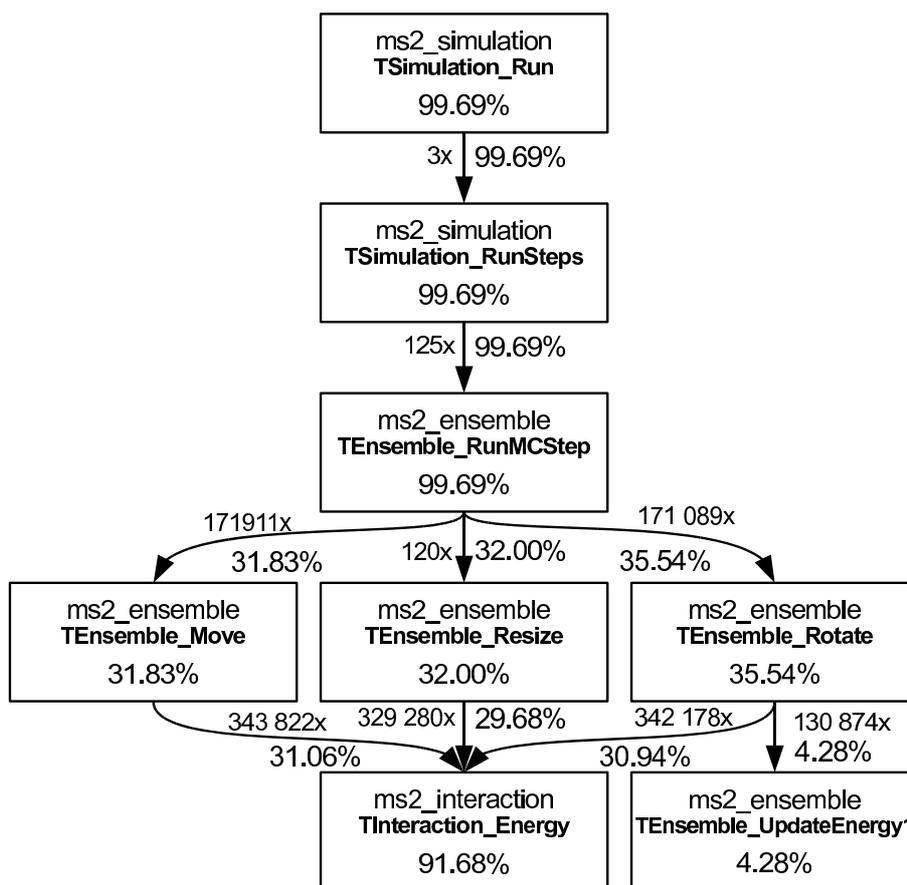


Figure 6: Runtime for the equimolar liquid mixture of methanol and ethanol at 298.15 K and 0.1 MPa in the NpT ensemble using different compilers and processors. The simulation was executed with 1372 molecules with MD for 200 NVT equilibration time steps, 500 NpT equilibration time steps and 4000 production time steps. For MC, 50 NVT equilibration loops, 200 NpT equilibration loops and 1000 production loops.

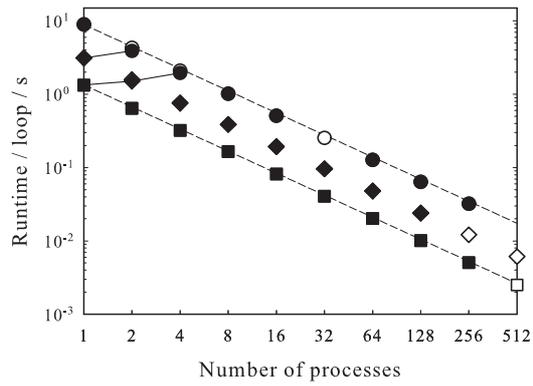


(a)

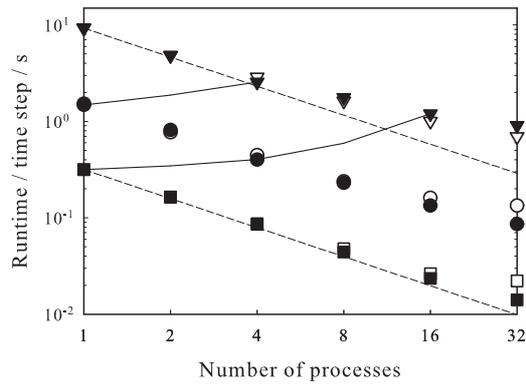


(b)

Figure 7: CPU cycles for one MD step (a) and one MC loop (b) of the test case estimated by valgrind. All functions with less than 1% CPU time share are not shown.



(a)

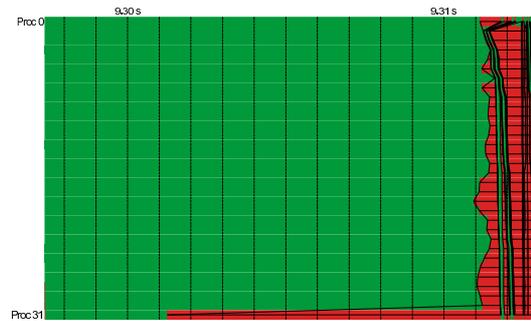


(b)

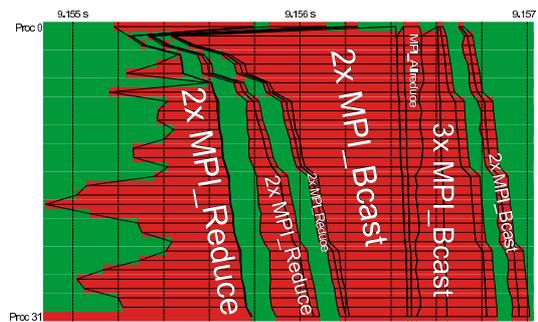
Figure 8: *ms2* runtime of the test case with MC (a) and MD (b) on a "Nehalem"/Infiniband cluster for Amdahl and Gustafson scaling. Full symbols: Intel MPI V4.0, empty symbols: OpenMPI V1.4.1, \square : $N=1372$, \diamond : $N=2744$, \circ : $N=5488$, ∇ : $N=21952$. The dashed lines indicate perfect strong scaling, while the full lines indicate the achieved weak scaling.



(a)



(b)



(c)

Figure 9: Communication between 32 processes executing the test case. The green color indicates calculations, the red color indicates communication. Figure (a) shows the summary time line for MC simulation, while the lower figures show the summary time line for MD simulation. Figure (b) exhibits the communication for one entire MD step, while (c) gives a magnified view on the communication pattern at the end of the time step.

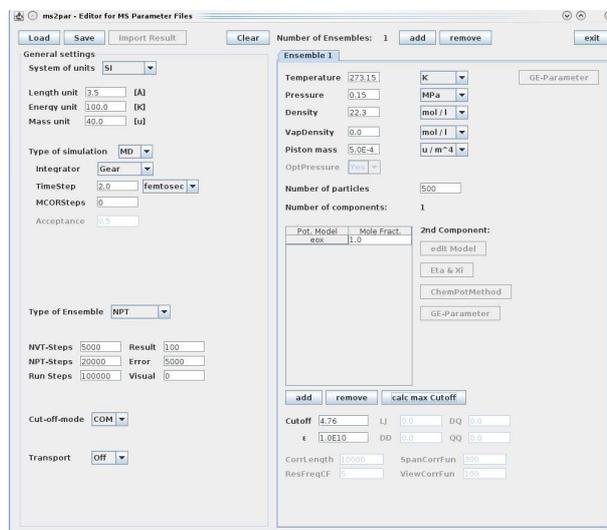


Figure 10: Snapshot of *ms2par* generating an input (*.par) file for a MD simulation of pure liquid ethylene oxide.



Figure 11: Snapshot of the tool *ms2chart*. The picture shows the analysis of a simulation of liquid cyclohexane.

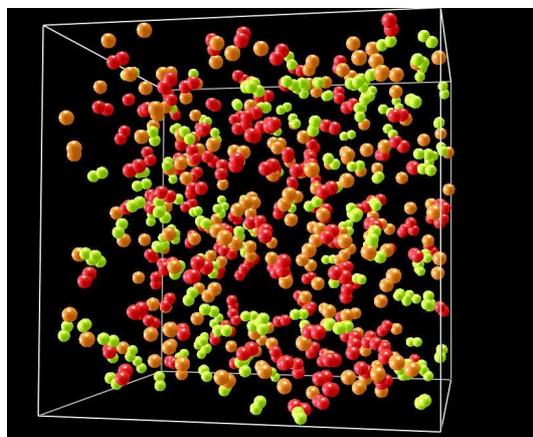


Figure 12: Snapshot of a ternary mixture taken with the visualization tool *ms2molecules*. The mixture consists of ethane (red), methane (orange), and carbon dioxide (green). The white frame illustrates the size of the simulation volume.

References

- [1] <http://www.fluidproperties.org>, Industrial fluid properties simulation collective, 2010.
- [2] Gubbins, K. E. and Moore, J. D., *The Journal of Chemical Physics* **49** (2010) 3026.
- [3] Tan, S. P., Adidharma, H., and Radosz, M., *Industrial & Engineering Chemistry Research* **47** (2008) 8063.
- [4] Kleiner, M., Turnakaka, F., and Sadowski, G., *Molecular Therodynamics of Complex Systems* **131** (2009) 75.
- [5] Klamt, A., *The Journal of Physical Chemistry* **99** (1995) 2224.
- [6] Klamt, A. and Eckert, F., *Fluid Phase Equilibria* **172** (2000) 43.
- [7] Gupta, S. and Olson, J., *Industrial & Engineering Chemistry Research* **42** (2003) 6359.
- [8] Case, F. H. et al., *Fluid Phase Equilibria* **274** (2008) 2.
- [9] Case, F. et al., *Fluid Phase Equilibria* **217** (2004) 1.
- [10] Case, F. et al., *Fluid Phase Equilibria* **236** (2005) 1.
- [11] Case, F. H. et al., *Fluid Phase Equilibria* **260** (2007) 153.
- [12] Case, F. H. et al., *Fluid Phase Equilibria* **285** (2009) 1.
- [13] Lyubartsev, A. P. and Laaksonen, A., *Computer Physics Communications* **128** (2000) 565.
- [14] Plimpton, S., *The Journal of Computational Physics* **117** (1993) 1.
- [15] *The Economist* **March 13th** (2010) 68.
- [16] Green, M., *The Journal of Chemical Physics* **22** (1954) 398.
- [17] Kubo, R., *The Journal of The Physical Society of Japan* **12** (1957) 570.
- [18] Poncela, A., Rubio, A. M., and Freire, J. J., *Molecular Physics* **91** (1997) 189.
- [19] Kristof, T., Vorholz, J., Liszi, J., Rumpf, B., and Maurer, G., *Molecular Physics* **97** (1999) 1129.
- [20] Ketko, M. H., Rafferty, J., Siepmann, J. I., and Potooff, J. J., *Fluid Phase Equilibria* **274** (2008) 44.
- [21] Bourasseau, E., Haboudou, M., Boutin, A., Fuchs, A. H., and Ungerer, P., *The Journal of Chemical Physics* **118** (2003) 3020.
- [22] Eckl, B., Vrabec, J., and Hasse, H., *Chemie Ingenieur Technik* **80** (2008) 25.
- [23] Eckl, B., Vrabec, J., and Hasse, H., *Molecular Physics* **106** (2008) 1039.

- [24] Eckl, B., Vrabec, J., and Hasse, H., *Fluid Phase Equilibria* **274** (2008) 16.
- [25] Huang, Y.-L., Vrabec, J., and Hasse, H., *Fluid Phase Equilibria* **287** (2009) 62.
- [26] Eckl, B., Horsch, M., Vrabec, J., and Hasse, H., *High Performance Computing In Science And Engineering '08* (2009) 119.
- [27] Merker, T., Guevara-Carrion, G., Vrabec, J., and Hasse, H., *High Performance Computing In Science And Engineering '08* (2009) 529.
- [28] Vrabec, J., Huang, Y.-L., and Hasse, H., *Fluid Phase Equilibria* **279** (2009) 120.
- [29] Allen, M. and Tildesley, D., *Computer Simulation of Liquids*, Clarendon Press, Oxford, 1987.
- [30] Vrabec, J. and Hasse, H., *Molecular Physics* **100** (2002) 3375.
- [31] Widom, B., *The Journal of Chemical Physics* **39** (1963) 2808.
- [32] Lyubartsev, A. P., Martinovski, A. A., Shevkunov, S. V., and Vorontsov-Velyaminov, P. N., *The Journal of Chemical Physics* **96** (1992) 1776.
- [33] Nezbeda, I. and Kolafa, J., *Molecular Simulation* **5** (1991) 391.
- [34] Frenkel, D. and Smith, B., *Understanding Molecular Simulation*, Academic Press, Elsevier, San Diego, 1993.
- [35] Rappaport, D., *The Art of Molecular Dynamics Simulation*, Cambridge University Press, Cambridge, 2004.
- [36] Born, M. and von Karman, T., *Physikalische Zeitschrift* **13** (1912) 297.
- [37] Metropolis, N., Rosenbluth, A., Rosenbluth, M. N., Teller, A. H., and Teller, E., *The Journal of Chemical Physics* **21** (1953) 1087.
- [38] Rowley, R. et al., *DIPPR Data Compilation of Pure Compound Properties*, Design Institute for Physical Properties, AIChE, 2003.
- [39] Flyvberg, H. and Petersen, H., *The Journal of Chemical Physics* **91** (1989) 461.
- [40] Lyubartsev, A., Laaksonen, A., and Vorontsov-Velyaminov, P., *Molecular Physics* **82** (1994) 455.
- [41] Vrabec, J., Kettler, M., and Hasse, H., *Chemical Physics Letters* **356** (2002) 431.
- [42] Shing, K., Gubbins, K., and Lucas, K., *Molecular Physics* **65** (1988) 1235.
- [43] Gray, C. and Gubbins, K., *Theory of molecular fluids, Volume 1: Fundamentals*, Clarendon Press, Oxford, 1984.
- [44] Gubbins, K., *Statistical Mechanics*, volume 1, The Chemical Society Burlington House, London, 1972.

- [45] Krishna, R. and van Baten, J. M., *Industrial & Engineering Chemistry Research* **44** (2005) 6939.
- [46] Hoheisel, C., *Physics Reports* **245** (1994) 111.
- [47] Barton, A., *The dynamic liquid state*, Longman, London, 1974.
- [48] Steele, W., *Transport Phenomena in fluids*, Marcel Dekker, New York, 1969.
- [49] Lennard-Jones, J., *Proceedings of the Physical Society* **43** (1931) 461.
- [50] Berthelot, D., *Comptes Rendues Hebdomadaires des Seances de l'Academie des Sciences* **126** (1898) 1703.
- [51] Lorentz, H., *Annalen der Physik* (1881) 127.
- [52] Fischer, J., Möller, D., Chialvo, A., and Halle, J. M., *Fluid Phase Equilibria* **48** (1989) 161.
- [53] Stone, A. J., *Science* **321** (2008) 787.
- [54] Hirschfelder, J., Curtiss, C., and Bird, R., *Molecular Theory of Gases and Liquids*, J. Wiley & Sons, 1954.
- [55] Nymand, T. and Linse, P., *The Journal of Chemical Physics* **112** (2000) 6152.
- [56] Steihauser, O., *Molecular Physics* **45** (1982) 335.
- [57] van der Spoel, D., van Maaren, P., and Berendsen, H., *The Journal of Chemical Physics* **108** (1998) 10220.
- [58] Ewald, P., *Annalen der Physik* **64** (1921) 253.
- [59] Lisal, M., Budinsky, R., and Vacek, V., *Fluid Phase Equilibria* **135** (1997) 193.
- [60] Garzon, B., Lago, S., Vega, C., and Rull, L., *The Journal of Chemical Physics* **102** (1995) 7204.
- [61] Lisal, M., William, W., and Nezbeda, I., *Fluid Phase Equilibria* **181** (2001) 127.
- [62] Jedlovsky, P. and Mezei, M., *The Journal of Chemical Physics* **110** (1999) 2991.
- [63] Lustig, R., *Molecular Physics* **65** (1988) 175.
- [64] MPI-Forum, *MPI: A Message-Passing Interface Standard, Version 2.2*, High Performance Computing Center Stuttgart (HLRS), 2009.
- [65] Guevara-Carrion, G., Nieto-Draghi, C., Vrabec, J., and Hasse, H., *The Journal of Physical Chemistry B* **112** (2008) 16664.
- [66] Schnabel, T., Srivastava, A., Vrabec, J., and Hasse, H., *The Journal of Physical Chemistry B* **111** (2007) 9871.
- [67] Towhee, <http://www.towhee.sourceforge.org>, 2008.
- [68] Errington, J. and Panagiotopoulos, A. Z., <http://kea.princeton.edu/jerring/gcmc/index.html>.
- [69] Hess, B., Kutzner, C., van der Spoel, D., and Lindahl, E., *The Journal of Chemical Theory and Computation* **4** (2008) 435.

Supplementary Material to *ms2*: A Molecular Simulation Tool for Thermodynamic Properties

Stephan Deublein^a, Bernhard Eckl^b, Jürgen Stoll^b, Sergey V. Lishchuk^c, Gabriela Guevara-Carrion^a, Colin W. Glass^d, Thorsten Merker^a, Martin Bernreuther^d, Hans Hasse^a, Jadran Vrabec^{e,*}

^aLehrstuhl für Thermodynamik, Universität Kaiserslautern, 67653 Kaiserslautern, Germany

^bInstitut für Technische Thermodynamik und Thermische Verfahrenstechnik, Universität Stuttgart, 70550 Stuttgart, Germany

^cDepartment of Mathematics, University of Leicester, Leicester LE1 7RH, United Kingdom

^dHöchstleistungsrechenzentrum Universität Stuttgart (HLRS), 70550 Stuttgart, Germany

^eLehrstuhl für Thermodynamik und Energietechnik, Universität Paderborn, 33098 Paderborn, Germany

Abstract

This supplementary material includes detailed descriptions of the equations implemented in *ms2* as well as detailed information on the input and output files of the program.

1. Definitions of the thermodynamic properties accessible in *ms2*

1.1. Density, pressure, internal energy and enthalpy

At constant temperature T and volume V , the pressure p is determined by [1]

$$p = \frac{k_B T}{V} + \frac{\langle W \rangle}{V} + \Delta p^L = \frac{k_B T}{V} + \frac{1}{3V} \left\langle \sum_{i=1}^{N-1} \sum_{\substack{j=i+1 \\ r_{ij} < r_c}}^N \mathbf{r}_{ij} \mathbf{f}_{ij} \right\rangle + \Delta p^L, \quad (1)$$

where k_B is the Boltzmann constant and the brackets $\langle \dots \rangle$ indicate the ensemble average. W denotes the virial, which is defined by the force vector \mathbf{f}_{ij} acting between two molecules at a separation vector \mathbf{r}_{ij} between their centers of mass. The contribution Δp^L considers the long-range interactions with molecules beyond the cut-off radius r_c .

In the isothermal-isobaric ensemble, the volume is a fluctuating parameter that on average yields the volume corresponding to the specified pair of temperature T and pressure p_0 . In MD simulations with Andersen's barostat [2], which is implemented in *ms2*, the fluctuations of the volume are damped by a fictive piston mass Q_p according to the equation of motion

$$\ddot{V} = \frac{p - p_0}{Q_p}. \quad (2)$$

In MC simulations, the volume is changed randomly and the new volume accepted by applying the Metropolis acceptance criterion

$$P_{acc} = \min\left(1, \exp\left(\frac{\Delta E}{k_B T}\right)\right), \quad (3)$$

*Corresponding author: Jadran Vrabec, Warburger Str. 100, 33098 Paderborn, Germany, Tel.: +49-5251/60-2421, Fax: +49-5251/60-3522, Email: jadran.vrabec@upb.de

where P_{acc} is the probability of accepting the volume change and ΔE is the difference between the old state with energy U_{old} for the volume V_{old} and the new state with energy U_{new} for the volume V_{new} according to

$$\Delta E = p_0(V_{new} - V_{old}) + (U_{new} - U_{old}) + Nk_B T \ln \frac{V_{old}}{V_{new}} . \quad (4)$$

The residual enthalpy H^{res} is directly linked to the residual internal energy, pressure and volume of the system

$$H^{res} = U^{res} + pV - Nk_B T , \quad (5)$$

where U^{res} is the residual potential energy. It is defined as the sum of all pairwise interaction energies u_{ij} and the appropriate long range correction ΔU^L , cf. section 2.4 of this Supplementary Material

$$U^{res} = \sum_{i=1}^{N-1} \sum_{\substack{j=i+1 \\ r_{ij} < r_c}}^N u_{ij} + \Delta U^L . \quad (6)$$

1.2. Second derivatives

The accessible second derivatives vary with the employed ensemble. In the NVT ensemble, the residual isochoric heat capacity c_v^{res} is determined by fluctuations of the residual potential energy U^{res}

$$c_v^{res} = \frac{1}{N} \left(\frac{\partial U^{res}}{\partial T} \right)_v = \frac{1}{k_B (NT)^2} (\langle (U^{res})^2 \rangle - \langle U^{res} \rangle^2) . \quad (7)$$

The partial derivative of the potential energy with respect to the volume at constant temperature $(\partial U^{res} / \partial V)_T$ is determined by fluctuations of the residual potential energy U^{res} and the virial W

$$\left(\frac{\partial U^{res}}{\partial V} \right)_T = \frac{N}{3V} \left(\frac{1}{k_B T} (\langle W \rangle \langle U^{res} \rangle - \frac{1}{N} \langle W U^{res} \rangle) + \langle W \rangle \right) . \quad (8)$$

In the NpT ensemble, the second derivatives of the Gibbs energy, namely the residual isobaric heat capacity c_p^{res} , the isothermal compressibility β_T and the volume expansivity α_p are functions of ensemble fluctuations [1]. The residual isobaric heat capacity c_p^{res} is related to fluctuations of the residual enthalpy H^{res}

$$c_p^{res} = \frac{1}{N} \left(\frac{\partial H^{res}}{\partial T} \right)_p = \frac{1}{k_B (NT)^2} (\langle (H^{res})^2 \rangle - \langle H^{res} \rangle^2) . \quad (9)$$

To obtain the total isobaric heat capacity c_p , the solely temperature dependent ideal gas contribution $c_p^{id}(T)$ has to be added. This ideal property is accessible e.g. by quantum chemical methods and can often be found in data bases [3].

An analogous relationship links the isothermal compressibility β_T to volume fluctuations in the NpT ensemble

$$\beta_T = -\frac{1}{V} \left(\frac{\partial V}{\partial p} \right)_T = \frac{1}{k_B T \langle V \rangle} (\langle V^2 \rangle - \langle V \rangle^2) . \quad (10)$$

The partial derivative of the residual enthalpy with respect to the pressure at constant temperature $(\partial H^{\text{res}}/\partial p)_T$ is linked to fluctuations of volume and residual internal energy

$$\left(\frac{\partial H^{\text{res}}}{\partial p}\right)_T = V - \frac{N}{T}(\langle U^{\text{res}}V \rangle - \langle U^{\text{res}} \rangle \langle V \rangle + p(\langle V^2 \rangle - \langle V \rangle^2)), \quad (11)$$

and the volume expansivity α_p again to volume and residual enthalpy fluctuations

$$\alpha_p = \frac{1}{V} \left(\frac{\partial V}{\partial T}\right)_p = \frac{N}{T^2 \langle V \rangle} (\langle H^{\text{res}}V \rangle - \langle H^{\text{res}} \rangle \langle V \rangle). \quad (12)$$

Note that Eqs. (7) to (12) are valid for mixtures as well.

1.3. Speed of sound

The speed of sound c is defined by the isothermal compressibility β_T , the volume expansivity α_p , the isobaric heat capacity c_p and the temperature T by

$$c = \left(\frac{1}{M(\beta_T \rho - T \alpha_p^2 / c_p)}\right)^{0.5}, \quad (13)$$

where M is the molar mass. In *ms2*, the speed of sound is calculated both for pure components and mixtures in the NpT ensemble.

1.4. Chemical potential - Widom test molecule insertion

For the calculation of the chemical potential of component i according to Widom [4], a so-called "test" molecule l of component i is inserted into the simulation volume at a random position with a random orientation. At constant temperature and pressure its potential energy ψ_l , due to its interactions with all other "real" molecules, is related to the chemical potential according to

$$\mu_i - \mu_i^{\text{id}}(T) = -k_B T \ln \frac{\langle V \exp(-\psi_l / k_B T) \rangle}{\langle N_i \rangle}. \quad (14)$$

The test molecule is removed immediately after the calculation of its potential energy, thus it does not influence the time evolution of the system or the Markov chain, respectively.

The value of ψ_l is highly dependent on the random position of the test molecule. In addition, the density of the system has a significant influence on the accuracy of the calculation. For very dense fluids, test molecules almost always overlap with some real molecules, which leads to a potential energy $\psi_l \rightarrow \infty$ and thus to no contribution to Eq. (14) resulting in poor statistics for the chemical potential or even complete failure of the sampling. Within limits, lower statistical uncertainties of the chemical potential can be achieved by inserting a large number of test molecules into the simulation volume, which leads to an increasing computational demand.

1.5. Self-diffusion coefficients

The self-diffusion coefficient D_i is related to the mass flux of single molecules within a fluid. Therefore, the relevant Green-Kubo expression is based on the individual molecule velocity autocorrelation function [5]

$$D_i = \frac{1}{3N_i} \int_0^\infty dt \langle \mathbf{v}_k(t) \cdot \mathbf{v}_k(0) \rangle, \quad (15)$$

where $\mathbf{v}_k(t)$ is the center of mass velocity vector of molecule k at some time t . Eq. (15) is an average over all N_i molecules of component i in the ensemble, since all contribute to the self-diffusion coefficient.

2. Details of *ms2*

This section provides a closer insight into *ms2* and the implemented equations.

2.1. Reduced quantities

The simulation program *ms2* internally uses reduced quantities for its calculations. All quantities are reduced by a reference length σ_R , a reference energy ϵ_R and a reference mass m_R , respectively.

Table 1: Important physical quantities in their reduced form. Note that ϵ_0 indicates the permittivity of the vacuum $\epsilon_0 = 8.854187 \times 10^{-12} \text{ A}^2 \text{ s}^4 \text{ kg}^{-1} \text{ m}^{-3}$.

length	$l^* = \frac{l}{\sigma_R}$
energy	$u^* = \frac{u}{\epsilon_R}$
mass	$m^* = \frac{m}{m_R}$
time	$t^* = \frac{t}{\sigma_R} \sqrt{\frac{\epsilon_R}{m_R}}$
piston mass	$Q_p^* = \frac{Q_p \sigma_R^4}{m_R}$
temperature	$T^* = \frac{T k_B}{\epsilon_R}$
pressure	$p^* = \frac{p \sigma_R^3}{\epsilon_R}$
density	$\rho^* = \rho \sigma_R^3 N_A$
volume	$V^* = \frac{V}{\sigma_R^3}$
chemical potential	$\tilde{\mu} = \frac{\mu}{k_B T}$
point charge	$q^* = \frac{1}{\sqrt{4\pi\epsilon_0}} \frac{q}{\sqrt{\epsilon_R \sigma_R}}$
dipole moment	$\mu^* = \frac{1}{\sqrt{4\pi\epsilon_0}} \frac{\mu}{\sqrt{\epsilon_R \sigma_R^3}}$
quadrupole moment	$Q^* = \frac{1}{\sqrt{4\pi\epsilon_0}} \frac{Q}{\sqrt{\epsilon_R \sigma_R^5}}$
diffusion coefficient	$D^* = \frac{D}{\sigma \sqrt{m_R / \epsilon_R}}$
viscosity	$\eta^* = \frac{\eta \sigma_R^2}{\sqrt{m_R \epsilon_R}}$
thermal conductivity	$\lambda^* = \frac{\lambda \sigma_R^2}{k_B \sqrt{m_R / \epsilon_R}}$

2.2. Intermolecular interactions

Point dipole interactions. The interaction between a point dipole with moment μ_i and a point charge q_j at a distance r_{ij} is given by [6, 7]

$$u_{ij}^{Dq}(r_{ij}, \theta_i, \mu_i, q_j) = -\frac{1}{4\pi\epsilon_0} \frac{\mu_i q_j}{r_{ij}^2} \cos \theta_i. \quad (16)$$

Here, θ_i is the angle between the distance vector of the point charge and the orientation vector of the point dipole, as illustrated in Figure 2 in the associated paper.

Point quadrupole interactions. The interaction potential of a linear point quadrupole Q_i with a point charge q_j at a distance r_{ij} is given by [1, 8]

$$u_{ij}^{Qq}(r_{ij}, \theta_i, Q_i, q_j) = \frac{1}{4\pi\epsilon_0} \frac{Q_i q_j}{4r_{ij}^3} (3 \cos^2 \theta_i - 1), \quad (17)$$

where θ_i is the angle between the distance vector of the point charge and the orientation vector of the point quadrupole, as illustrated in Figure 2 in the associated paper.

The interaction between a linear point quadrupole with moment Q_i and a point dipole with moment μ_j is given by [6, 8]

$$u_{ij}^{Q\mu}(r_{ij}, \theta_i, \theta_j, \phi_{ij}, Q_i, \mu_j) = \frac{1}{4\pi\epsilon_0} \frac{3}{2} \frac{Q_i \mu_j}{r_{ij}^4} (\cos \theta_i - \cos \theta_j) \quad (18)$$

$$(1 + 3 \cos \theta_i \cos \theta_j - 2 \cos \phi_{ij} \sin \theta_i \sin \theta_j),$$

where the angles θ_i , θ_j and ϕ_{ij} indicate the relative angular orientation of the point dipole j and the point quadrupole i , as shown in Figure 2 in the associated paper.

2.3. Reaction field method

The truncation of electrostatic interactions of first and second order are corrected for with the reaction field method [9, 10]. Here, all dipoles within the cut-off radius r_c polarize the fluid surrounding the cut-off sphere, which is modeled as a dielectric continuum with a relative permittivity ϵ_s . The polarization gives rise to a homogeneous electric field within the cut-off radius, called the reaction field \mathbf{E}_i^{RF} of magnitude

$$\mathbf{E}_i^{\text{RF}} = \frac{1}{4\pi\epsilon_0} \frac{2(\epsilon_s - 1)}{2\epsilon_s + 1} \frac{1}{r_c^3} \sum_{b=1}^m \sum_{\substack{j=1 \\ r_{ij} < r_c}}^{N_b} \boldsymbol{\mu}_{b,j}. \quad (19)$$

Note that all m dipoles with moment $\boldsymbol{\mu}$ on all N_b molecules within the cut-off sphere have to be summed up. The reaction field acting outside of the cut-off sphere interacts with a dipole $\boldsymbol{\mu}_i$ at the center of the cut-off sphere by [1, 11]

$$u_i^{\mu\text{RF}} = -\frac{1}{2} \boldsymbol{\mu}_i \mathbf{E}_i^{\text{RF}}, \quad (20)$$

where $u_i^{\mu\text{RF}}$ is its contribution to the potential energy. In Eq. (20), it is assumed that the system is sufficiently large so that tinfoil boundary conditions ($\varepsilon_s \rightarrow \infty$) are applicable without a loss of accuracy, e.g. $N \geq 500$ [12, 13, 14, 15].

The reaction field method can also be applied to correct for sets of point charges, as long as they add up to a total charge of zero. Therefore, a point charge distribution is reduced to a dipole vector $\boldsymbol{\mu}^q$ according to

$$\boldsymbol{\mu}^q = \sum_{i=1}^n \mathbf{r}_i q_i, \quad (21)$$

where n denotes the total number of point charges and \mathbf{r}_i the position vector of charge q_i . The resulting dipole $\boldsymbol{\mu}^q$ is transferred into the correction term according to Eqs. (19) and (20).

2.4. Cut-off mode

In the site-site cut-off mode, the LJ contributions of molecules beyond the cut-off radius r_c to the internal energy are estimated by [1]

$$\Delta U^{L*} = \frac{8}{9} N \rho^* \pi \left(\left(\frac{1}{r_c^*} \right)^9 - 3 \left(\frac{1}{r_c^*} \right)^3 \right). \quad (22)$$

The contributions to the pressure are considered by [1]

$$\Delta p^{L*} = \frac{32}{9} (\rho^*)^2 \pi \left(\left(\frac{1}{r_c^*} \right)^9 - \frac{3}{2} \left(\frac{1}{r_c^*} \right)^3 \right). \quad (23)$$

Using the center of mass cut-off mode, the LJ contributions of molecules beyond the cut-off radius r_c are estimated on the basis of the correction terms of Lustig [16]. The correction Δu^{L*} for the residual internal energy can be divided into three contributions, the contributions by interactions between two molecule centers (TICC), between one molecule center and one site (TICS), and between two molecule sites (TISS). The first contributions describe the long range contributions between molecular interaction sites that are positioned in the center of the mass of the individual molecules. The TICS terms describe the contributions, where one site is located in the center of mass of its molecule, while the other site is not. The TISS terms correct for contributions between molecular sites that are not positioned in the center of mass of their molecules. Using these formulations, the correction term for the residual internal energy is defined by

$$\Delta u^{L*} = 2\pi\rho \sum_{i=1}^{N_C} \sum_{j=1}^{N_C} \sum_{\alpha=1}^{N_{LJ,i}} \sum_{\beta=1}^{N_{LJ,j}} x_i x_j 4\varepsilon \begin{cases} (\text{TICCu}(-6) - \text{TICCu}(-3)) \\ (\text{TICSu}(-6) - \text{TICSu}(-3)) \\ (\text{TISSu}(-6) - \text{TISSu}(-3)) \end{cases} \quad (24)$$

Here, N_C is the number of components in the system and $N_{LJ,i}$ the number of LJ sites of the molecule of component i . The TIXXu functions are calculated by the following equations, where the argument in the brackets defines the exponent n

$$\text{TICCu}(n) = \frac{r_c^{2n+3}}{\sigma^{2n}(2n+3)}, \quad (25)$$

$$\text{TICSu}(n) = \frac{(r_c + \tau)^{2n+3} - (r_c - \tau)^{2n+3}}{4\sigma^{2n}\tau(n+1)(2n+3)}r_c + \frac{(r_c + \tau)^{2n+4} - (r_c - \tau)^{2n+4}}{4\sigma^{2n}\tau(n+1)(2n+3)(2n+4)}, \quad (26)$$

$$\begin{aligned} \text{TISSu}(n) = & -\frac{(r_c + \tau_+)^{2n+4} - (r_c + \tau_-)^{2n+4} - (r_c - \tau_-)^{2n+4} + (r_c - \tau_+)^{2n+4}}{8\sigma^{2n}\tau_1\tau_2(n+1)(2n+3)(2n+4)}r_c + \\ & \frac{(r_c + \tau_+)^{2n+5} - (r_c + \tau_-)^{2n+5} - (r_c - \tau_-)^{2n+5} + (r_c - \tau_+)^{2n+5}}{8\sigma^{2n}\tau_1\tau_2(n+1)(2n+3)(2n+4)(2n+5)}. \end{aligned} \quad (27)$$

The terms r_+ and r_- are defined by

$$\tau_+ = \tau_1 + \tau_2, \quad (28)$$

$$\tau_- = \tau_1 - \tau_2, \quad (29)$$

where τ_1 and τ_2 define the distance between sites 1 and 2, respectively, to the center of mass of the molecules they belong to.

The correction of the pressure for the center of mass cut-off is given by

$$\Delta p^{L*} = \frac{-2}{3}\pi\rho^2 \sum_{i=1}^{N_C} \sum_{j=1}^{N_C} \sum_{\alpha=1}^{N_{L,j,i}} \sum_{\beta=1}^{N_{L,j}} x_i x_j 4\varepsilon \begin{cases} (\text{TICCP}(-6) - \text{TICCP}(-3)) \\ (\text{TICSp}(-6) - \text{TICSp}(-3)) \\ (\text{TISSp}(-6) - \text{TISSp}(-3)). \end{cases} \quad (30)$$

These expressions follow the same naming scheme as for the internal energy correction. The functions TICCP , TICSp and TISSp are defined by

$$\text{TICCP}(n) = 2n \cdot \text{TICCu}(n), \quad (31)$$

$$\text{TICSp}(n) = -\frac{(r_c + \tau)^{2n+2} - (r_c - \tau)^{2n+2}}{4\sigma^{2n}\tau(n+1)}r_c^2 - 3 \cdot \text{TICSu}(n), \quad (32)$$

$$\text{TISSp}(n) = -\frac{(r_c + \tau_+)^{2n+3} - (r_c + \tau_-)^{2n+3} - (r_c - \tau_-)^{2n+3} + (r_c - \tau_+)^{2n+3}}{8\sigma^{2n}\tau_1\tau_2(n+1)(2n+3)}r_c^2 - 3 \cdot \text{TISSu}(n). \quad (33)$$

3. Modular structure - example

The modular structure of *ms2* is discussed by highlighting a MD simulation calculation process of a pure fluid modeled by two LJ sites and two point charges in the *NVT* ensemble. During the initialization process, the modules for memory allocation and definition of the simulation scenario are executed. This process is not unique for any of the hierarchical levels, but takes place on all levels, cf. Figure 4 in the associated paper.

On level one, the simulation is delimited to "molecular dynamics", therefore all modules concerning MC simulation are entirely neglected. On level two, the "initialization" modules define the simulation environment by setting ensemble specific properties, like simulation volume, temperature as well as the thermostat and integration algorithm etc. While the use of the first modules mentioned above is required in each simulation, independently of MD or MC, e.g. the trajectory calculation is specific for MD. However, the modular structure retains the efficiency of the program by initiating only the molecule propagation that is specific for the simulation technique. The same effect of modules can be found in the initiation process on lower levels. Here, the composition of the molecular species in the simulation volume (module component), the description of the molecule (module molecule) and the storage of the positions, velocities and forces of all LJ sites and charges (module sites) are determined. All other modules, e.g. treating additional components for mixtures or containing characteristic properties of dipoles or quadrupoles, are fully omitted by *ms2* in this case, because they are not needed in this simulation. Once the simulation has been initialized, a set of modules is invoked, which deals with maintaining the simulation run. On level one, this includes modules of global use for any simulation, e.g. limiting the execution time etc. On level two, Newton's equations of motion are numerically integrated, requiring information about the fluid's composition, the positions and orientations of the molecules etc. provided by modules of the respective level. The thermostat is applied and the calculation of thermodynamic properties is performed. The forces, torques and energies are calculated in the second branch on the level component. For a MD simulation, this covers routines for the calculation of energies as well as forces at the same time, whereas for a MC simulation, the modules only calculate energies and virial contributions. A further set of modules deals with the accumulation of data. These modules are designed to store data and evaluate the statistics of the produced data. The actual thermodynamic properties are calculated on the ensemble level. The last set of modules deals with the output, where the results of the molecular simulation are written to file. These modules are called independently on the user settings in all simulation scenarios.

4. Input and output

ms2 was designed to be an easily applicable simulation program. The structure of the input files as well as the output files is shown in Figure 1.

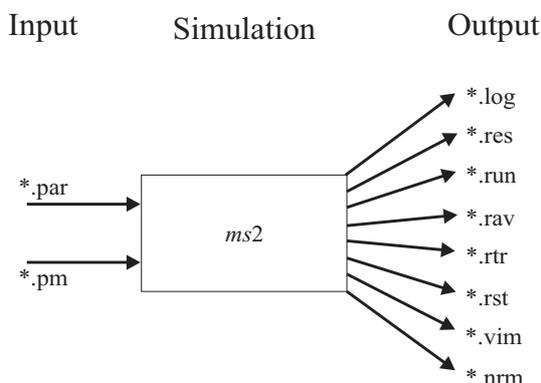


Figure 1: File structure needed and generated by *ms2*.

4.1. Input file *.par

The simulation program *ms2* requires one input file (*.par) to specify the simulation parameters and one molecular model file (*.pm) for every molecular species considered. The *.par file contains all input variables for the simulation process, such as simulation type, ensemble, number of equilibration and production steps, time step length etc. Furthermore, the user has to specify temperature, density and fluid composition. Table 2 lists all input parameters and options to be specified in the *.par file. An example for a complete *.par file is given below for a MC simulation of pure ethylene oxide in the NpT ensemble, where the chemical potential is calculated by gradual insertion.

Table 2: Parameters and options specified in the *.par file.

Parameter	Option	Explanation	Recommended value
Units	SI	Physical properties in the *.par file are given in SI units	SI
	reduced	Physical properties in the *.par file are given in reduced units with respect to the reference values of length σ_R , energy ϵ_R and mass m_R	
LengthUnit		Reference length σ_R in Å	3.0
EnergyUnit		Reference energy ϵ_R/k_B in K	100.0
MassUnit		Reference mass m_R in atomic units $u = 1.6605 \times 10^{-27}$ kg	100.0
Simulation	MD	Molecular dynamics simulation	
	MC	Monte-Carlo simulation	
Integrator	Gear	Gear predictor-corrector integrator (MD only)	Gear
	Leapfrog	Leapfrog integrator (MD only)	

Table 2 continued

Parameter	Option	Explanation	Recommended value
TimeStep		Time step of one MD step in fs (MD only)	~1
Acceptance		Acceptance rate for MC moves (MC only)	0.5
Ensemble	NVT	Canonical ensemble	
	NVE	Micro-canonical ensemble	
	NPT	Isobaric-isothermal ensemble	
	GE	Grand equilibrium method (pseudo- μVT)	
MCORSteps		MC relaxation loops for pre-equilibration	100
NVTSteps		Number of equilibration time steps (MD) or loops (MC) in the <i>NVT</i> ensemble	20000
NPTSteps		Number of equilibration time steps (MD) or loops (MC) in the <i>NpT</i> ensemble (optional)	50000
RunSteps		Number of production time steps (MD) or loops (MC)	300000
ResultFreq		Size of block averages in time steps or loops	100
ErrorFreq		Frequency of writing the *.res file in time steps or loops	5000
VisualFreq		Frequency of saving configurations in the *.vim file for visualization in time steps (MD) or loops (MC)	5000
CutoffMode	COM	Center of mass cut-off	COM
	Site	Site-site cut-off	
NEnsemble		Number of ensembles in the simulation	1
CorrfunMode	yes	Calculation of autocorrelation functions enabled	
	no	Calculation of autocorrelation functions disabled	
Temperature		Specified temperature	
Pressure		Specified pressure	
Density		Specified density	
PistonMass		Piston mass for simulations at constant pressure	
NParticles		Total number of molecules	864 - 4000
Liqdensity		Simulation result density	
VarLiqDensity		Statistical uncertainty of density	
LiqEnthalpy		Simulation result residual enthalpy	
VarLiqEnthalpy		Statistical uncertainty of residual enthalpy	
LiqBetaT		Simulation result isothermal compressibility	
VarLiqBetaT		Statistical uncertainty of isothermal compressibility	
LiqHdP		Simulation result $(dh^{\text{res}}/dp)_T$	
VarHdP		Statistical uncertainty of $(dh^{\text{res}}/dp)_T$	
NComponents		Number of components in the simulation	

Table 2 continued

Parameter	Option	Explanation	Recommended value
Corrlength		Length of the autocorrelation in time steps	
SpanCorrFun		Time steps separating subsequent autocorrelation functions	
ViewCorrFun		Output frequency of the full autocorrelation functions into the *.rtr file	
ResultFreqCF		Output frequency of transport properties into the *.res file	
PotModel		Potential model *.pm file of a component	
MolarFract		Molar fraction of a component	
ChemPotMethod	none	No calculation of the chemical potential for this component	none
	Widom	Calculation of the chemical potential for this component using Widom's test molecule method	
	GradIns	Calculation of the chemical potential for this component using gradual insertion	
NTest		Number of test molecules for Widom's test molecule method	2000
WeightFactors	Guess	For gradual insertion: use user defined initial values for the weight factors with optimization of these factors during simulation	Guess
	OptSet	For gradual insertion: user defined values for the weight factors without adjustment during simulation	
Cutoff		Cut-off radius for center of mass cut-off	
CutoffLJ		Cut-off radius for LJ interactions (site-site cut-off)	
CutoffDD		Cut-off radius for dipole-dipole interactions (site-site cut-off)	
CutoffDQ		Cut-off radius for dipole-quadrupole interactions (site-site cut-off)	
CutoffQQ		Cut-off radius for quadrupole-quadrupole interactions (site-site cut-off)	
Epsilon		Dielectric constant	1.0E+10

An example for a *.par file is given in Table 3. The scenario is a MC simulation in the NpT ensemble for ethylene oxide.

Table 3: Parameters and options specified in the *.par file.

Sim_EOX.par		
Units	=	SI
LengthUnit	=	3.0
EnergyUnit	=	100.0
MassUnit	=	100.0
Simulation	=	MC
Acceptance	=	0.5
Ensemble	=	NVT
MCORSteps	=	100
NVTSteps	=	2000
NPTSteps	=	10000
RunSteps	=	50000
ResultFreq	=	100
ErrorsFreq	=	2000
VisualFreq	=	10
CutoffMode	=	COM
NEnsembles	=	1
Temperature	=	400.0
Pressure	=	0.79355
Density	=	21.09227
PistonMass	=	0.00001
NParticles	=	500
NComponents	=	1
PotModel	=	eox.pm
MolarFract	=	1.0
ChemPotMethod	=	GradIns
WeightFactors	=	Guess
		1.00
		2.00
		3.00
		4.00
		5.00
		10.00
		20.00
		40.00
		60.00
Cutoff	=	5.0
Epsilon	=	1.E+10

4.2. Input file *.pm

A *.pm file contains the molecular model for a given substance. It contains the relative positions and the parameters of all sites. The potential model file for methanol is shown in Table 4. Methanol was modeled by two LJ sites and three point charges [17]. All positions and distances in the *.pm file are given in Å, the LJ parameters σ and ε/k_B are given in Å and K, respectively, while the mass is given in atomic units ($u = 1.6605 \times 10^{-27}$ kg). The magnitudes of the charges are specified in electronic charges ($e = 1.602 \times 10^{-19}$ C), while the dipole moments and quadrupole moments are given in Debye ($D = 3.33564 \times 10^{-30}$ Cm) and Buckingham ($B = 3.33564 \times 10^{-40}$ Cm²), respectively. The orientations of the dipole and quadrupole are represented by spherical coordinates, where the azimuthal angle ϕ specifies the angle to the positive x -axis and the polar angle θ defines the angle to the positive z -axis. Both angles are specified in degrees. Molecular models can be orientated arbitrarily in the *.pm file. All site positions are transformed into a principal axes coordinate system at the beginning of each simulation with *ms2*. The normalized site positions are written to a *.nrm file for each component.

Table 4: Parameters and options specified in the *.pm file for a molecular model of methanol.

MeOH.pm		
NSiteTypes	=	2
SiteType	=	LJ126
NSites	=	2
x	=	7.660331E-01
y	=	1.338147E-02
z	=	0.0
sigma	=	3.754348
epsilon	=	120.591759
mass	=	15.034
x	=	-6.564695E-01
y	=	-6.389332E-02
z	=	0.0
sigma	=	3.030
epsilon	=	87.879094
mass	=	16.00
SiteType	=	Charge
NSites	=	3
x	=	7.660331E-01
y	=	1.338147E-02
z	=	0.0
charge	=	0.247461
mass	=	0.0
shielding	=	0.1
x	=	-6.564695E-01
y	=	-6.389332E-02
z	=	0.0
charge	=	-0.678742
mass	=	0.0
shielding	=	0.1
x	=	-1.004989E+00
y	=	8.145993E-01
z	=	0.0
charge	=	0.431281
mass	=	1.008
shielding	=	0.05
NRotAxes	=	auto

4.3. Output files

ms2 yields seven output files:

- *.log file - stores a complete summary of all execution steps taken by *ms2*.
- *.res file - contains the results of the simulation in an aggregated form. The data is written to file in reduced quantities as well as in SI units, along with the statistical uncertainties of the calculated properties. The *.res file is created during simulation and updated every specified number of time steps or loops.
- *.run file - contains the calculated properties of the simulation for a specified time step or loop interval. The file is in tabular form, where the data is given in reduced units. The file is subsequently updated according to the user specification, which is set in the *.par file.
- *.rav file - contains the block averages of the calculated properties. The file is in tabular form, where the data is given in reduced units. The file is subsequently updated according to the user specification, which is set in the *.par file.
- *.rtr file - stores the final values of the autocorrelation functions and their integrals. The number of output lines has to be defined in the *.par file.
- *.rst file - is the restart file of the simulation. It contains all molecular positions, velocities, orientations, forces, torques and block averages for the thermodynamic properties. It is written once at the end of a simulation or immediately after having received a termination signal of the operating system. The *.rst file allows for a stepwise execution of the simulation, necessary e.g. in case of an early interruption of the simulation, time limits on a queuing system or unexpected halts.
- *.vim file - is the trajectory visualization file. It contains the positions and orientations of all molecules in an aggregated ASCII format. The configurations are written to file after a user-specified interval of time steps or loops. The *.vim file is readable by the visualization tool *ms2molecules*, which is also part of the simulation package.
- *.nrm file - stores the normalized coordinates of a potential model after a principal axes transformation.

References

- [1] Allen, M. and Tildesley, D., *Computer Simulation of Liquids*, Clarendon Press, Oxford, 1987.
- [2] Andersen, H., *The Journal of Chemical Physics* **72** (1980) 2384.
- [3] Rowley, R. et al., *DIPPR Data Compilation of Pure Compound Properties*, Design Institute for Physical Properties, AIChE, 2003.
- [4] Widom, B., *The Journal of Chemical Physics* **39** (1963) 2808.
- [5] Gubbins, K., *Statistical Mechanics*, volume 1, The Chemical Society Burlington House, London, 1972.
- [6] Hirschfelder, J., Curtiss, C., and Bird, R., *Molecular Theory of Gases and Liquids*, J. Wiley & Sons, 1954.
- [7] Lorrain, P., Corson, D., and Lorrain, F., *Elektromagnetische Felder und Wellen*, Walter de Gruyter, Berlin, 1995.
- [8] Gray, C. and Gubbins, K., *Theory of molecular fluids, Volume 1: Fundamentals*, Clarendon Press, Oxford, 1984.
- [9] Frenkel, D. and Smith, B., *Understanding Molecular Simulation*, Academic Press, Elsevier, San Diego, 1993.
- [10] Nymand, T. and Linse, P., *The Journal of Chemical Physics* **112** (2000) 6152.
- [11] Boettcher, C., van Belle, O., Bordewijk, P., and Rip, A., *Theory of electric polarization, Vol. 1: Dielectrics in static fields*, Elsevier, 1973.
- [12] Lisal, M., Budinsky, R., and Vacek, V., *Fluid Phase Equilibria* **135** (1997) 193.
- [13] Garzon, B., Lago, S., Vega, C., and Rull, L., *The Journal of Chemical Physics* **102** (1995) 7204.
- [14] Lisal, M., William, W., and Nezbeda, I., *Fluid Phase Equilibria* **181** (2001) 127.
- [15] Jedlovsky, P. and Mezei, M., *The Journal of Chemical Physics* **110** (1999) 2991.
- [16] Lustig, R., *Molecular Physics* **65** (1988) 175.
- [17] Schnabel, T., Srivastava, A., Vrabec, J., and Hasse, H., *The Journal of Physical Chemistry B* **111** (2007) 9871.