

MDBenchmark: a toolkit to optimize the performance of molecular dynamics simulations

Michael Gecht,^{1, a)} Marc Siggel,¹ Max Linke,¹ Gerhard Hummer,^{1, 2} and Jürgen Köfinger^{1, b)}

¹⁾*Department of Theoretical Biophysics, Max Planck Institute of Biophysics, Max-von-Laue-Straße 3, 60438 Frankfurt am Main, Germany*

²⁾*Institute for Biophysics, Goethe University, Max-von-Laue-Straße 9, 60438 Frankfurt am Main, Germany*

(Dated: 22 June 2020)

Despite the impending flattening of Moore’s law, the system size, complexity and length of molecular dynamics (MD) simulations keeps on increasing, thanks to effective code parallelization and optimization combined with algorithmic developments. Going forward, exascale computing poses new challenges to the efficient execution and management of MD simulations. The diversity and rapid developments of hardware architectures, software environments, and MD engines make it necessary that users can easily run benchmarks to optimally set up simulations, both with respect to time-to-solution and overall efficiency. To this end, we have developed the software MDBenchmark to streamline the setup, submission, and analysis of simulation benchmarks and scaling studies. The software design is open and as such not restricted to any specific MD engine or job queuing system. To illustrate the necessity and benefits of running benchmarks and the capabilities of MDBenchmark, we measure the performance of a diverse set of 22 MD simulation systems using GROMACS 2018. We compare the scaling of simulations with the number of nodes for CPU-only and mixed CPU-GPU nodes and optimize the numbers of MPI ranks and OpenMP threads. Our results demonstrate the importance of benchmarking for finding optimal system and hardware specific simulation parameters. Running MD simulations with optimized settings leads to a significant performance increase that reduces the monetary, energetic, and environmental costs of MD simulations.

I. INTRODUCTION

Molecular dynamics (MD) simulations have become an integral part of the molecular life sciences and material sciences. Their predictive power has been continuously increasing thanks to methodological advances and the exponential growth of compute power. The latter is captured by Moore’s law for the number of transistors of an integrated circuit and, across different technologies, by the law of accelerated returns^{1,2}. This continuing growth translates into a similar growth of the temporal and spatial scales that can be assessed in MD simulations³. Thus, MD simulations are becoming more powerful in extending and connecting the different scales that are accessible to experimental methods⁴. As so-called computational microscopes^{5,6}, MD simulations are widely used to make predictions and to analyze, design, and validate experiments.

The tools to perform MD simulations have reached a level of sophistication which allows non-expert users to set up, run, and analyze simulations. Various software packages for performing MD simulations, so called MD engines, have been developed, e.g., ACEMD⁷, Amber⁸, CHARMM⁹, Desmond¹⁰, GROMACS¹¹, LAMMPS¹², NAMD¹³, and OpenMM¹⁴. These engines, each with a unique set of features, were designed to efficiently compute the time evolution of particles in a simulation box and take advantage of different hardware architectures and parallel computation to varying degrees.

Optimal settings for MD simulations aim to decrease time-to-solution (TTS) and to increase the throughput of simulations and energy efficiency. To run an MD simulation on a compute cluster, the user has to specify numerous parameters to control the behavior of the underlying hardware and software. Moreover, optimal parameters might vary between different versions of the same MD engine and depend on the details of the molecular system under consideration. Users new to MD simulations might miss the details that are required to run efficient simulations on high-performance compute clusters. A poor choice of parameters can notably degrade the simulation performance, usually measured in simulated time per day, e.g., ns/day, or increase it only by a small margin while wasting resources. Such inefficient simulations decrease overall simulation throughput, lead to a higher electricity demand and operating cost, and ultimately increase the carbon footprint. It is therefore in the interest of every user to optimize their usage of hardware resources, and, at the same time, keep their TTS to a minimum. This daunting task of finding the optimal parameters for running MD simulations can only be tackled by thorough benchmarking.

Therefore, we need to enable individual users to run benchmarks conveniently and efficiently for their given settings. This complements the goals of systematic benchmark studies performed by experts^{15–18}. For example, Kutzner *et al.* have performed extensive benchmarks using GROMACS to determine the best performance-to-price ratio for a variety of MD systems and numerous hardware architectures^{17,18}. They provided valuable guidelines for configuring and purchasing new clusters and for choosing optimal parameters. However, the rapid development of hardware, algorithms, and software and the variety of MD engines and simulation systems requires that users themselves are able to run benchmark studies effi-

^{a)}Electronic mail: michael.gecht@biophys.mpg.de

^{b)}Electronic mail: juergen.koefinger@biophys.mpg.de

ciently.

We developed MDBenchmark, a standalone software package, implemented as a command-line interface (CLI), to conveniently set up, run and analyze benchmarks of MD simulations. With this tool users can run benchmarks and scaling studies for their specific molecular system, MD engine, and compute cluster. MDBenchmark was developed to streamline and simplify the process of finding the optimal run parameters and settings for any simulation and hardware stack. It takes care of submitting simulations to the queuing system, performs scaling studies by varying the number of nodes, automatically toggles the usage of CPUs and/or GPUs, and scans the numbers of processes used for parallelization (MPI ranks, OpenMP threads) if applicable. The package was designed for ease of use not only by expert users, but also by researchers without prior detailed knowledge of the ins and outs of high-performance computing (HPC).

To illustrate the application and the capabilities of MDBenchmark and to highlight the value and necessity of running benchmarks, we report on an extensive scaling study of 22 MD simulations of varying sizes ($1 \times 10^5 - 4 \times 10^6$ atoms) and system compositions. We identify numbers of MPI ranks and OpenMP threads that produce the best performance for a range of system sizes, study the benefits of hyperthreading, and analyze when it is beneficial to use CPU-only or mixed CPU-GPU nodes. For this study, we use the GROMACS software suite, as it is widely used, freely available, and highly optimized for different kinds of hardware. However, MDBenchmark has been designed such that different MD engines and job queuing systems can easily be added.

II. BACKGROUND

Current compute clusters are composed of compute nodes, each containing at least one central processing unit (CPU), an optional graphics processing unit (GPU), as well as gigabytes of dedicated random-access memory (RAM). These nodes are connected in a network, such that data can be exchanged between nodes and calculations can be performed in parallel on multiple nodes.

Modern CPUs contain dozens of physical cores, where each core can perform computations independent from the others. In addition, a single physical core can often perform two computations at the same time, a feature called “hyperthreading”¹⁹. When enabled, the number of physical cores is virtually doubled, i.e., for each physical core two “logical cores” are introduced.

To use these heterogeneous resources efficiently and run calculations in parallel, two interfaces are widely used: *message passing interface* (MPI) and *open multi-processing* (OpenMP). MPI spawns processes, which we will refer to as ranks. A single MPI rank can comprise all cores of a single node or only a subset of them. In contrast, OpenMP creates computational threads, where each is composed only of a subset of cores available inside a MPI rank. OpenMP threads share the same memory.

Running a simulation on a computer cluster requires the user to submit a compute job to a queuing system. The user must configure a submission script that launches the MD engine over its CLI. Users have to define the correct parameters for the specific queuing system. Activating hyperthreading on the “Sun Grid Engine” (SGE) differs, for example, from activating it on the “Simple Linux Utility for Resource Management” (SLURM) queuing system.

The numbers of MPI ranks to OpenMP threads influence the performance of GROMACS simulations^{17,18}. To take advantage of the parallel compute infrastructure, a simulation box is first divided into separate domains in a process called domain decomposition²⁰. Each domain is regarded in an isolated manner and information at the borders is communicated with the other domains. In a hybrid MPI-OpenMP approach, the calculations of a single domain are managed by a single MPI rank. This rank spawns multiple OpenMP threads, which then perform the actual calculations in each domain. Each MPI rank communicates with the ranks responsible for its neighbouring domains. The number of ranks n_{ranks} per node times the number of threads n_{threads} per node gives the number of logical cores per node. If hyperthreading is deactivated then the number of logical cores is equal to the number of physical cores. If activated, it is equal to twice the number of physical cores.

III. THE MDBENCHMARK SOFTWARE

The CLI of MDBenchmark provides access to four main functions (Fig. 1). In the first step, all parameters for the benchmark(s) are defined using the `GENERATE` command [Fig. 1(a)]. Here, the user chooses the MD engine, the numbers of nodes to perform scaling on, whether to use GPUs and the numbers of MPI ranks. In addition, a run input file of an equilibrated MD simulation must be provided. A `.TPR` file is sufficient for GROMACS. Different MD engines require different input files. For example, the `.NAMD`, `.PSF` and `.PDB` files have to be provided for NAMD. MDBenchmark automatically checks the availability of the requested MD engine and its installed version using the “Environment Modules” system²¹. This feature was put in place to safeguard against typos in the module name. If a module environment is not used on the compute cluster, the user can skip this availability check. MDBenchmark will prompt the user to confirm the action, before proceeding to create the folder structure.

The folder structure was intentionally designed with a nested hierarchy to allow users to access files themselves, if needed. Each requested MD engine is put into its own folder, with a subfolder denoting the engine’s version and whether MDBenchmark is going to request CPU-only or mixed CPU-GPU nodes. The last subfolder layer separates the benchmarks by the numbers of nodes, MPI ranks, and OpenMP threads, and by whether hyperthreading is enabled or disabled. Each of these subfolders then contains a copy of the run input file, as provided by the user, the job submission script containing all parameters and commands to run the benchmark, as well as a hidden folder holding all metadata in JSON for-

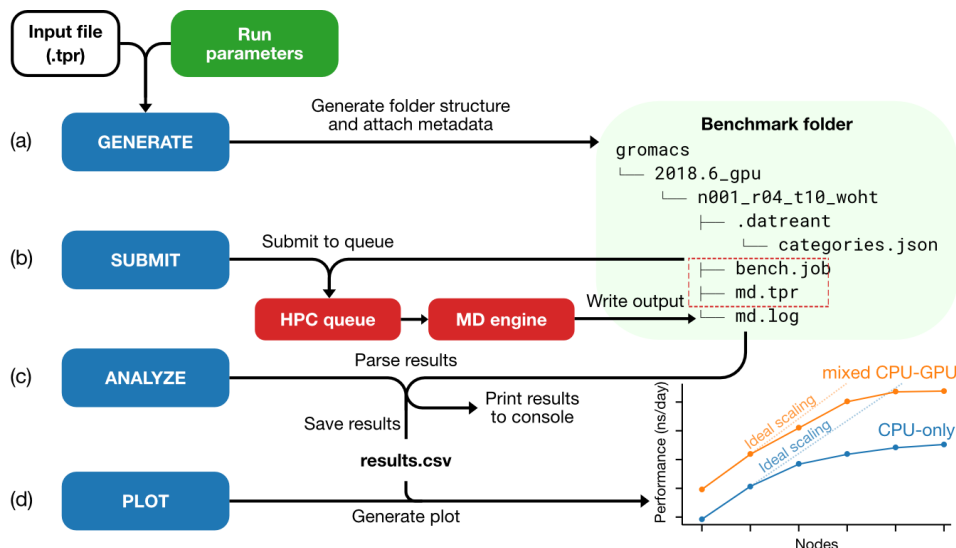


FIG. 1. Schematic representation of the implementation and workflow of MDBenchmark. (a) Run input files and benchmark parameters are supplied to the GENERATE command. For every benchmark, a folder with attached metadata is created. (b) The SUBMIT command sends the benchmark job to the HPC queue, starting the MD engine on the requested resources. (c) The log file, written by the engine, contains all necessary results and is parsed through the ANALYZE command. Results are shown in the console or saved to a CSV file. (d) The PLOT command can visualize results in a plot using the CSV file as input. It produces a scaling plot with the performance (ns/day) as a function of the number of nodes, with every benchmark setting as a separate line.

mat. This metadata is managed through the datreant Python package²². We use it to define each benchmark as an entity and add our parameters as metadata. The package can be used to search and filter benchmarks through their Python API. This way, benchmarks can easily be grouped by distinct parameters for additional customized analysis by the user.

After benchmarks have been set up, they can be submitted to the queuing system with the SUBMIT command in the second step [Fig. 1(b)]. When called, it will traverse all subfolders and gather information on each benchmark. The user will be shown all benchmarks that are to be submitted and they will be prompted to submit or cancel. If a benchmark was already submitted, it will be excluded from further submissions. The user can ask MDBenchmark to submit all benchmarks, ignoring their submission status. Currently, MDBenchmark can submit jobs to the queuing systems SGE, SLURM, and IBM’s LoadLeveler. The available queuing system will be automatically determined before submission. The submission files are conveniently implemented using the template syntax provided by the Jinja2 Python package, which facilitates the easy adaption to other HPC resources and different requirements (further details are available in the documentation).

After submission to the queuing system, the status of benchmarks can be checked with the ANALYZE command in the third step [Fig. 1(c)]. MDBenchmark will print all available information on each benchmark in a tabular format to the console. Every call to the ANALYZE command will parse the log files produced by GROMACS and extract all important information on the corresponding run. If a benchmark has already finished, its performance will be printed in the last column of the table. All available benchmark data can be saved in a CSV format to a file for future analysis.

In a fourth step, the name of the file containing the benchmark results can be provided as an argument to the PLOT command, which will produce a scaling plot [Fig. 1(d)]. The plot shows the performance (ns/day) as a function of the number of nodes, with each line representing a different combination of parameters (CPU-only/mixed CPU-GPU nodes, numbers of MPI ranks and OpenMP threads, and hyperthreading). A linear fit through the first two data points is added to indicate ideal scaling. The PLOT command also accepts multiple CSV files from separate benchmarks, to allow for straightforward comparison of different benchmark parameters. In addition, the CSV file can be easily read with the pandas Python package²³ for customized analysis.

Further detailed explanations of all available options can be found online in the MDBenchmark documentation (<https://mdbenchmark.readthedocs.io/en/version-3/>).

IV. METHODS

We used MDBenchmark to run benchmarks of atomistic MD simulations for 22 biomolecular systems. The system sizes range from 10^5 to 4×10^6 atoms [Fig. 2(a)]. 15 out of the 22 tested systems (68 %) have less than 300,000 atoms. Due to the nature of biomolecular simulations, the 22 simulation systems studied here are mostly composed of solvent: 55 – 99 % of the systems total atoms belong to water molecules [Fig. 2(b)]. If present, lipid bilayers account for 8 – 45 % of all remaining atoms. Other solutes, i.e., proteins, nucleic acids, small molecules and ions, make up only 0.1 – 27 %. These systems cover a broad range of force fields, water models, compositions, box geometries and sizes (Ta-

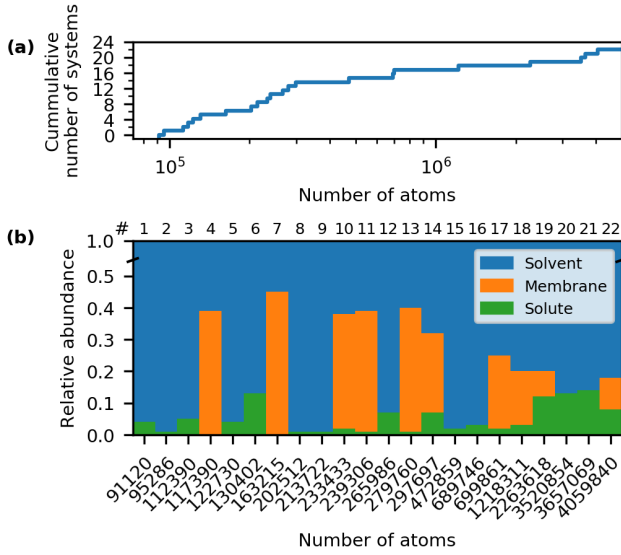


FIG. 2. Size and composition of simulation systems. (a) Cumulative number of systems as a function of the atom number. Most MD systems (68 %) contain less than 300,000 atoms. (b) Relative abundance of system constituents as a function of the number of atoms. Systems are composed of 55 – 99 % solvent (blue). Lipid bilayers, if present, make up 8 – 45 % (orange) of the systems total number of atoms. All other solutes account for, at most 27 %, but usually less than 10 % (green). Note that the numbers of atoms comprise all interaction sites, i.e., also the additional interaction sites for TIP4P-D water.

ble I). Their exact setup can be found in the corresponding references, if already published.

For these 22 systems, we performed scaling studies in which we determine the performance $P(N)$ as a function of the number of nodes N . We vary the number of MPI ranks, n_{ranks} , which also determine the number of OpenMP threads per rank, n_{threads} , for activated and deactivated hyperthreading.

We use Amdahl's law as a simple model to summarize the results of our scaling studies²⁴. This law describes the speed-up $S(N)$ of parallelized computations as a function of the number of nodes N , i.e.,

$$S(N) = \frac{1}{1 - p + \frac{p}{N}}, \quad (1)$$

where p is the fraction of the code that benefits from parallelization.

We use Amdahl's law to express the performance $P(N) = S(N)P(1)$ as

$$P(N) = \frac{P^{\max}}{1 + \frac{p}{N(1-p)}}, \quad (2)$$

where the maximum performance $P^{\max} = \lim_{N \rightarrow \infty} P(N)$ is given by

$$P^{\max} = \frac{P(1)}{1 - p}. \quad (3)$$

The ideal scaling, determined by the speed-up going from $N = 1$ to $N = 2$, is then given by

$$P^{\text{id}}(N) = P(1) + (N - 1)(P(2) - P(1)), \quad (4)$$

which becomes

$$P^{\text{id}}(N) = P^{\max} \left(1 + (N - 1) \frac{p}{2 - p} \right). \quad (5)$$

We use Amdahl's law to estimate the performance that can be achieved by increasing the number of nodes, while being reasonably efficient, i.e., close to ideal scaling as given by Eq. (4). For parallelization to be efficient, we demand that the performance is a fraction f of the ideal performance, i.e.,

$$P(N) = f P^{\text{id}}(N) \quad (6)$$

and solve for N . Introducing $q = 1 - p$ for brevity, we obtain

$$N(f) = \frac{1}{2pqf} \left[q + 1 - f(p^2 + 2q^2) + \sqrt{(q + 1 - f(p^2 + 2q^2))^2 - 8p^2q^2f^2} \right] \quad (7)$$

The performance corresponding to a fraction f of the ideal scaling is then given by $P(N(f))$.

We performed benchmarks with GROMACS 2018¹¹ without the built-in hybrid thread-MPI library, but with external MPI libraries and enabled OpenMP support. AVX_512 SIMD instructions were enabled at compile time with GCC 8.3 and CUDA 10.1. We used nodes with two Intel® Xeon® Gold 6148 CPUs (2.40 GHz) for all benchmarks, with additionally two NVIDIA Quadro RTX 5000 graphic cards in mixed CPU-GPU nodes. Benchmarks were run on specific numbers of MPI ranks and OpenMP threads for a total wall time of 15 minutes, if not mentioned otherwise. We used SLURM as the queuing system for all benchmarks on the available supercomputer. Example submission files for both CPU-only and mixed CPU-GPU nodes can be found in the supporting information.

All data generated in this study was managed and analyzed using datreant²², MDAnalysis^{25,26}, NumPy²⁷, Pandas²³, SciPy²⁸ and IPython²⁹. Plots were generated with Matplotlib³⁰.

V. RESULTS

Using MDBenchmark, we first examine the dependence of the performance estimates on the run time of the benchmarks. We then use two exemplary systems to discuss in detail the performance scaling and its dependence on the number of MPI ranks and on hyperthreading. We then present results for the complete data set of 22 systems and show how optimal parameters vary with system size.

A. Run time of benchmarks

Ideally, a benchmark run provides the most accurate performance estimate (ns/day) in the least amount of time. Short

TABLE I. Details of the fully atomistic simulation systems, for which we perform scaling benchmarks. The systems vary in size, composition, and box geometry, as well as in force fields and water models: CHARMM36³¹, CHARMM36m³², AMBER99SB*-ILDN^{33–35}, AMBER99SB*-ILDN-Q^{33–37}, TIP3P³⁸, and TIP4P-D³⁹.

	Type	# atoms	Box geometry	System size (nm)	Force field	Water model	Reference
1	Protein in solution	91,120	cuboid	9.9x9.9x9.9	CHARMM36	TIP3P	Unpublished
2	Protein in solution	95,286	cuboid	10.0x10.0x10.0	CHARMM36	TIP3P	Unpublished
3	Protein in solution	112,390	cuboid	10.0x10.0x10.0	AMBER99SB*-ILDN-Q	TIP3P	Unpublished
4	Protein in membrane	117,390	cuboid	11.0x11.0x10.0	CHARMM36	TIP3P	Hofbauer <i>et al.</i> ⁴⁰
5 ^a	Protein in solution	122,730	cuboid	10.0x10.0x10.0	AMBER99SB*-ILDN-Q	TIP4P-D	Unpublished
6 ^a	Dense protein solution	130,402	cuboid	9.9x9.9x9.9	AMBER99SB*-ILDN-Q	TIP4P-D	Bülow <i>et al.</i> ⁴¹
7	Empty membrane	163,215	cuboid	13.1x13.1x9.1	CHARMM36m	TIP3P	Unpublished
8 ^a	Protein in solution	202,512	rhomb. dodec.	13.0	AMBER99SB*-ILDN-Q	TIP4P-D	Unpublished
9	dsDNA in solution	213,722	rhomb. dodec.	14.5x14.5x10.3	AMBER99SB*-ILDN	TIP3P	Unpublished
10	Protein and membrane	233,433	cuboid	14.6x14.6x10.5	CHARMM36m	TIP3P	Wu <i>et al.</i> ⁴²
11	Protein and membrane	239,306	cuboid	14.6x14.6x10.9	CHARMM36	TIP3P	Unpublished
12 ^a	Dense protein solution	265,986	cuboid	12.6x12.6x12.6	AMBER99SB*-ILDN-Q	TIP4P-D	Bülow <i>et al.</i> ⁴¹
13	Protein and membrane	279,760	cuboid	15.9x15.9x10.8	CHARMM36	TIP3P	Unpublished
14	Protein in membrane	297,697	cuboid	13.8x13.8x15.4	CHARMM36m	TIP3P	Hofmann <i>et al.</i> ⁴³
15 ^a	Protein in solution	472,859	cuboid	15.0x15.0x15.0	AMBER99SB*-ILDN-Q	TIP4P-D	Unpublished
16 ^a	Dense protein solution	689,746	cuboid	17.3x17.3x17.3	AMBER99SB*-ILDN-Q	TIP4P-D	Bülow <i>et al.</i> ⁴¹
17	Protein in membrane	699,861	cuboid	19.5x19.5x18.1	CHARMM36m	TIP3P	Unpublished
18	Protein and membrane	1,219,446	cuboid	23.0x23.0x23.0	CHARMM36m	TIP3P	Unpublished
19	Protein and membrane	2,263,618	cuboid	38.2x21.6x31.5	CHARMM36	TIP3P	Unpublished
20 ^a	Dense protein solution	3,520,854	cuboid	30.0x30.0x30.0	AMBER99SB*-ILDN-Q	TIP4P-D	Bülow <i>et al.</i> ⁴¹
21 ^a	Dense protein solution	3,657,069	cuboid	30.4x30.4x30.4	AMBER99SB*-ILDN-Q	TIP4P-D	Bülow <i>et al.</i> ⁴¹
22 ^b	Protein and membrane	4,059,840	hexameric	34.9x34.9x37.9	CHARMM36m	TIP3P	Unpublished

^a Systems use TIP4P-D water model.

^b The timestep of this system was to 4 fs.

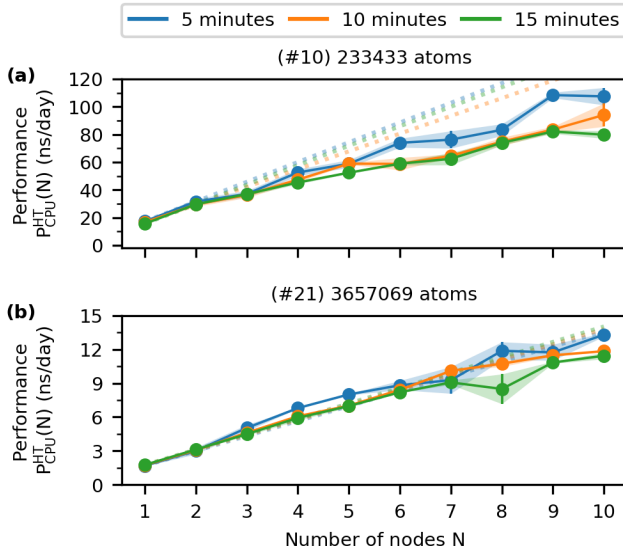


FIG. 3. Dependence of the performance estimates on the run times of the benchmarks. MD simulations with (a) 233k and (b) 3.6M atoms scaled from 1 to 10 nodes and run for 5 (blue), 10 (orange) and 15 (green) minutes. Each data point shows the average of 5 independent runs with error bars and filled-in area showing the standard deviation. Transparent dashed lines show ideal scaling as estimated by linearly extrapolating the performance of one and two nodes [Eq. (4)].

benchmarks facilitate a higher throughput, a more efficient use of the limited computing resources, and thus allow us to sample a broader range of parameters.

To set the length of a benchmark run, we can either use a fixed number of steps, as previously done by Kutzner *et al.*^{17,18}, or set the run time explicitly. In a GROMACS simulation, the first hundreds to thousands of steps can be used to balance computational load between different ranks using dynamic load balancing. This auto-tuning is enabled by default and beneficial for overall simulation performance. However, a benchmark will be aborted if this process has not finished after half of the available compute time. For larger systems, the auto-tuning process takes longer. Thus, we decided to set the run time explicitly as jobs with fixed, short run times can be given higher priority by the queuing systems.

To determine the minimum run time of benchmarks, we have performed scaling studies of two systems using $N = 1 - 10$ nodes and used run times of 5, 10, and 15 minutes (Fig. 3) using CPU-only nodes. Results for the shortest run times of 5 minutes can deviate significantly from the results for 10 and 15 minutes. The shortest run time can be used for a first screening while longer run times seem to be necessary to obtain accurate results. The benchmark results can be influenced, for example, by the amount of traffic handled by the network infrastructure at run time. In the following, all reported benchmarks were run for 15 minutes independent of the system size. Interestingly, the performance estimates from multiple runs with fixed benchmark time vary more for CPU-

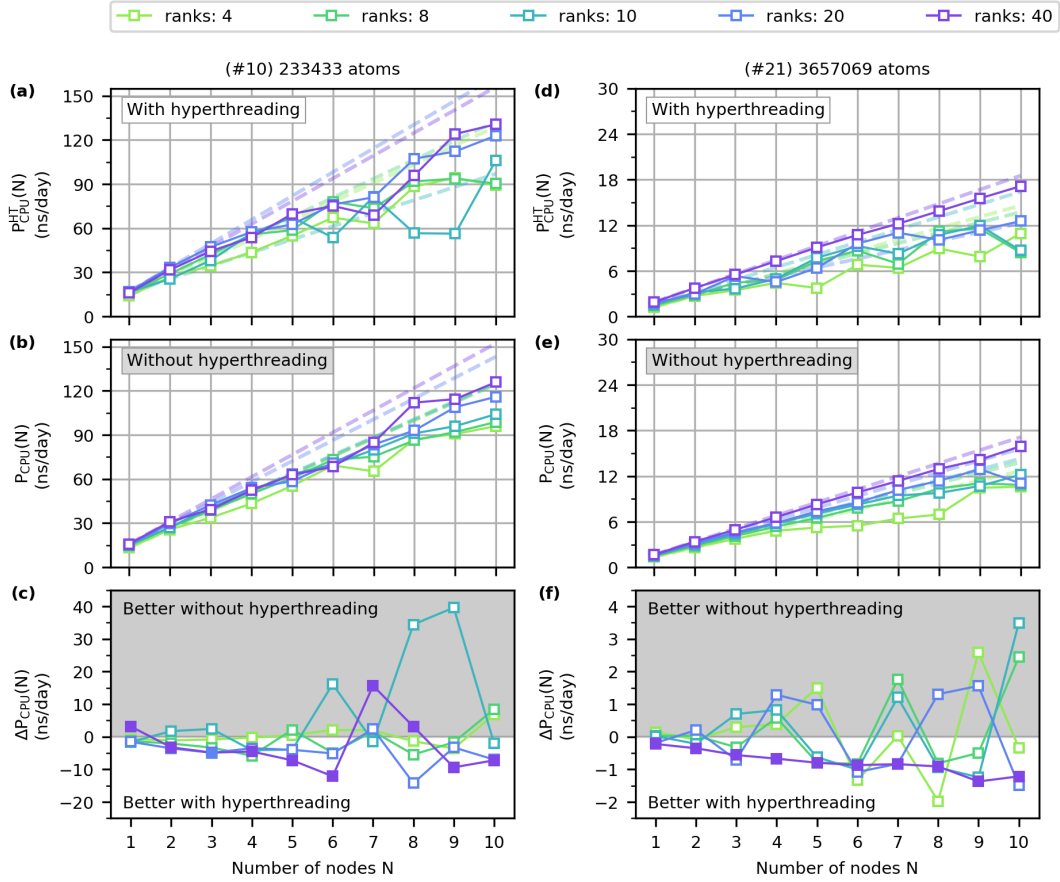


FIG. 4. Scaling of the performance $P_{\text{CPU}}(N)$ with the number N of CPU-only nodes for different numbers of MPI ranks n_{ranks} (colors) with (a,d) and without (b,e) hyperthreading. Ideal scaling was estimated by a linear fit to the performances of one and two nodes [Eq. (4), dashed lines]. (a,b) Prototypical membrane protein system with 233k atoms and (d,e) dense protein solution with 3.6M atoms. Both systems scale the best with $n_{\text{ranks}} = 40$. (c,f) Differences $\Delta P_{\text{CPU}}(N)$ between the performances with and without hyperthreading [Eq. (8)]. For the best performing rank settings (filled squares) of $n_{\text{ranks}} = 40$, both systems generally benefit from running simulations with hyperthreading.

only nodes than for mixed CPU-GPU nodes, as can be seen, for example, below in Fig. 6.

B. Optimizing performance for two exemplary systems

We evaluated how the performance scales with different numbers of nodes using either CPU-only or mixed CPU-GPU nodes. Using MDBenchmark, we ran benchmarks using 1 to 10 nodes and scanned the values of the MPI ranks with and without hyperthreading. In the following, we present detailed results for two exemplary systems of different size and composition (systems 10 and 21 in Table I): system #10, a prototypical membrane protein system with 233k atoms⁴² and system #21, a large 3.6M atom system of a dense protein solution using TIP4P-D as the water model⁴¹.

Our results show that the dependence of the performance on the number of MPI ranks is different for CPU-only nodes and for mixed CPU-GPU nodes. For CPU-only nodes with hyperthreading activated, both systems show the best performance for $n_{\text{ranks}} = 40$, consistently for all node numbers [Figs. 4(a)

and 4(d)]. Thus, for CPU-only nodes the optimal number of ranks is independent of the system size. In contrast, for mixed CPU-GPU nodes with hyperthreading activated, the optimal choice of n_{ranks} depends on the system size [Figs. 5(a) and 5(d)]. For the default parameters set by GROMACS 2018, $n_{\text{ranks}} = 40$ and hyperthreading activated, the medium sized system with 233k atoms shows the worst performance [Fig. 5(a)]. We find that the optimal numbers of ranks yielding the highest performances are given by $n_{\text{ranks}} = 8$ as well as by $n_{\text{ranks}} = 10$. These two settings can scale up to 5 nodes with 80% efficiency (dashed lines). This example illustrates that blindly trusting the default settings can decrease performance by more than half. For the larger systems of 3.6M atoms, we observe that a higher number of ranks n_{ranks} results in better performance and that the optimal value is actually the default setting of $n_{\text{ranks}} = 40$ and $n_{\text{threads}} = 2$ [Fig. 5(d)].

With hyperthreading activated, the previously determined optimal values of n_{ranks} remain unchanged, both for CPU-only and mixed CPU-GPU nodes [Figs. 4(b), 4(e), 5(b) and 5(e)]. To quantify the effect of hyperthreading, we calculate the dif-

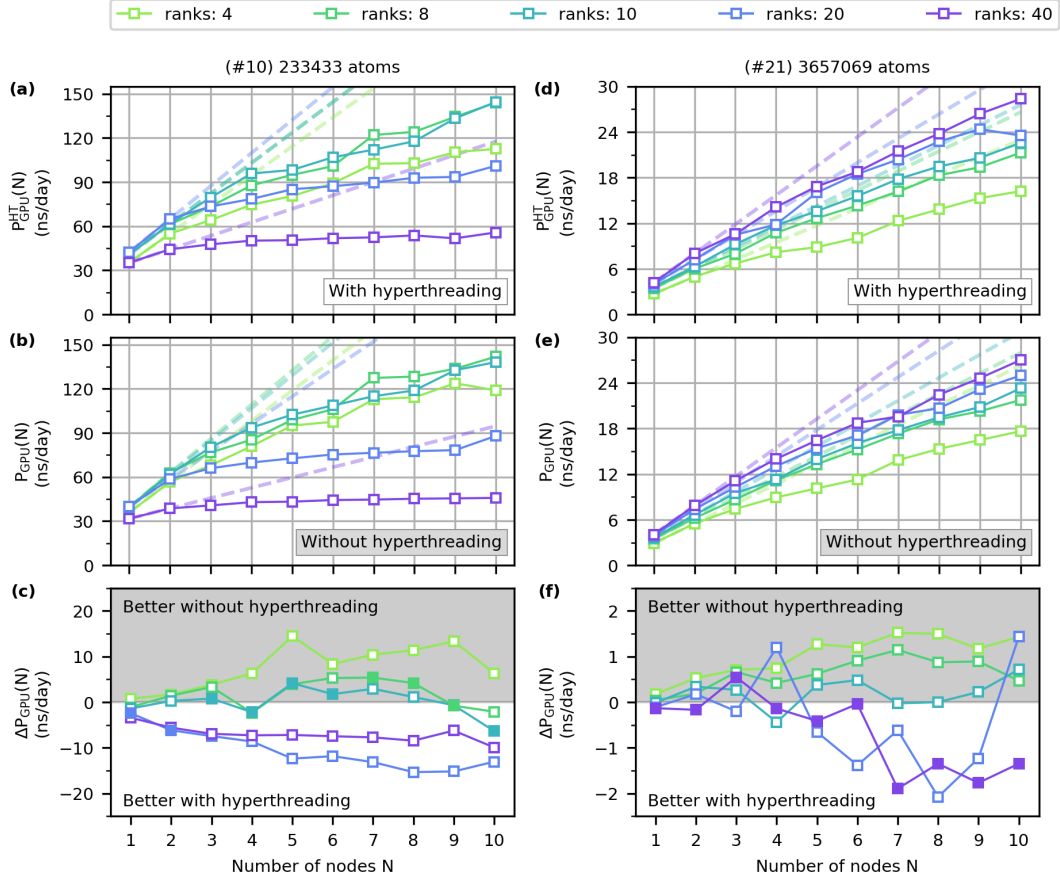


FIG. 5. Scaling of the performance $P_{\text{GPU}}(N)$ with the number N of mixed CPU-GPU nodes for different numbers of MPI ranks n_{ranks} (colors) with (a,d) and without (b,e) hyperthreading. Ideal scaling was estimated by a linear fit to the performances of one and two nodes [Eq. (4), dashed lines]. For the membrane protein system with 233k atoms, we obtain the best performances for $n_{\text{ranks}} = 8$, (a) with and (b) without hyperthreading. For the dense protein solution with 3.6M atoms, we obtain the best performance for $n_{\text{ranks}} = 40$, (c) with and (d) without hyperthreading. (c,f) Differences $\Delta P_{\text{GPU}}(N)$ between the performances with and without hyperthreading [Eq. (8)]. (c) We find that for the 233k atom system, the performance for the optimal rank settings (filled symbols) does not benefit from hyperthreading. (f) For the 3.6M atom system, a significant increase in performance can be achieved by activating hyperthreading for the optimal rank settings (filled squares).

ference of the absolute performances as

$$\Delta P_x(N) = P_x(N) - P_x^{\text{HT}}(N) \quad (8)$$

where $P_x(N)$ is the performance without hyperthreading and $P_x^{\text{HT}}(N)$ is the performance with hyperthreading. Here and in the following, subscript $x = \text{CPU}$ denotes CPU-only and $x = \text{GPU}$ denotes mixed CPU-GPU nodes. The superscript HT indicates that hyperthreading is activated, its absence that it is deactivated.

We find that simulations on CPU-only nodes benefit from hyperthreading [Figs. 4(c) and 4(f)] whereas simulations on mixed CPU-GPU nodes do not in general [Figs. 5(c) and 5(f)]. The differences in performance for activated and deactivated hyperthreading $P_x(N)$ for a given value of n_{ranks} can change signs between different node numbers. However, the overall tendencies become clear if we consider the optimal rank setting for each node number individually. For CPU-only nodes, performances generally benefit from hyperthreading. Only for the 233k atom system, we see deviations from this behav-

ior for three node numbers ($N=1,7,8$). For mixed CPU-GPU nodes, we find that there is no significant benefit of activating hyperthreading for the 233k atom system and all node numbers and for the 3.6M atom systems and smaller node numbers ($N \lesssim 6$). Only for the larger system and larger node numbers, the activation of hyperthreading leads to a $\sim 5\%$ performance increase for the optimal rank settings.

We now show that the scaling behavior, i.e., the dependence of the performance on the number of nodes, is well captured by Amdahl's law in the form of Eq. (2) for both hardware architectures (Figs. 4 and 5). This simple law captures the linear increase of the performance with small node numbers and the convergence to a plateau for larger node numbers. If node numbers become too large, the scaling breaks down and Amdahl's law cannot be applied. Using the optimal rank and hyperthreading settings for the two systems and hardware architectures considered here, we performed scaling studies up to 64 nodes. We fit Amdahl's law to the scaling curves using the single fit parameter p . We find that Amdahl's law fits the

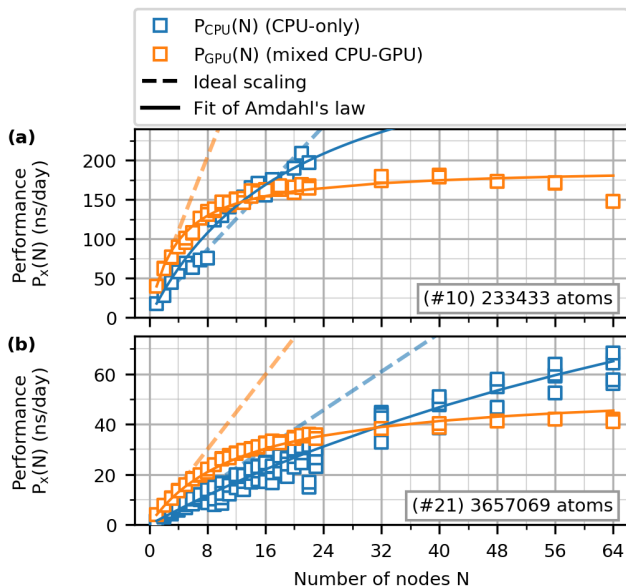


FIG. 6. Amdahl’s law captures the scaling behaviour. Performance $P_x(N)$ of two MD simulations with varying sizes as a function of nodes N . Each panel shows the best performing n_{ranks} values on mixed CPU-GPU nodes without hyperthreading (orange) and CPU-only nodes with 40 n_{ranks} and hyperthreading (blue) from Figs. 4 and 5. Each combination of nodes N was run in multiple independent simulations of 15 minutes each. All replicates are shown. Solid colored line shows the fits to the data using Eq. (2). Transparent dashed lines show ideal scaling as estimated from the performances of one and two nodes.

scaling curves reasonably well and that it can thus be used to summarize the results, as is done in the next section.

A comparison of the performance scaling for the two hardware architectures shows that the mixed CPU-GPU nodes perform better than CPU-only nodes for small node numbers only (Fig. 6). For the membrane protein system, mixed CPU-GPU nodes perform better up to $N = 14$ [Fig. 6(a)], for the dense protein solution up to 32 nodes. For larger N the performance increases only slowly for mixed CPU-GPU nodes such that CPU-only nodes achieve higher absolute performances. This behavior is consistent for all of the 22 systems considered here and the point of equal performance shifts to higher N with increasing system size (Figs. S1 and S2). The location of this point of equal performance is also determined by the relative computational power of CPU and GPU. Note that both node architectures contain the same CPU. A less powerful CPU would thus shift the point of equal performance to higher node numbers.

C. Size dependence of optimal parameters

As we have shown above for two exemplary systems, the optimal values for n_{ranks} depend on the number of atoms in a system for mixed CPU-GPU nodes. We further validate these observations with additional scaling results of the remaining

20 systems on mixed CPU-GPU nodes without hyperthreading, varying the n_{ranks} values for different numbers of nodes N (Fig. S1). We also present scaling benchmarks for each system with the best performing settings on CPU-only nodes ($n_{\text{ranks}} = 40$ with hyperthreading).

Our benchmark results for the 22 systems consistently show that for mixed CPU-GPU nodes the optimal numbers of n_{ranks} decreases with increasing node number and increases with system size (Fig. 7). Note that the number of OpenMP threads is given by the number of logical cores divided by the number of MPI ranks. The observed trends are clear and consistent and provide guidelines for the optimal rank settings. However, the deviations from these monotonic trends also show that for a specific system it is beneficial to run benchmarks for numbers of ranks close to the optimal values.

The absolute performances for the 22 systems follow the trends as exemplified above for the 233k and 3.6M atom systems (Fig. 8). Mixed CPU-GPU nodes [Fig. 8(a)] perform better than CPU-only nodes [Fig. 8(b)] only for fairly small numbers of nodes, as indicated by their performance ratio [Fig. 8(c)]. With the exception of systems 6 and 12, the speed-up factor is larger than two for all simulations on a single node. For two nodes, a speed-up larger than two is achieved for nearly all system sizes larger than 236k atoms, with system 12 and 16 being the exception. As we discuss in the following, these exceptions are due to aggressive simulation settings.

The overall scaling of the performance with the number of nodes and system size is fairly smooth and quite monotonic. However, the dense protein solutions (systems 6, 12, 16, 20, 21) systematically deviate from the overall trends for CPU-only nodes [Fig. 8(b)], which also becomes noticeable in the performance ratios [Fig. 8(c)]. The reason is that in these simulations the cutoff-distance in the real space interactions has been reduced and the Fourier space cut-off in the PME calculation increased compared to the other systems. With these settings, performance could be increased significantly on CPU-only nodes. However, the performance on mixed CPU-GPU nodes remains largely unaffected. Note that in addition to the dense protein solutions listed here also systems 5, 8, and 15 use the TIP4P-D³⁹ water model. However, these systems do not show deviations from the overall performance trends.

For all system sizes, simulations on CPU-only nodes scale much better with the number of nodes than mixed CPU-GPU nodes (Fig. 9). We quantify the computational efficiency of choosing N nodes by calculating the ratio of the actual performance to the performance we would get for ideal scaling. We determine ideal scaling using Eq. (5), using values of p from fits of Amdahl’s law to the performance data. Instead of fitting to the scaling curve for the n_{ranks} value which gives the overall best performance, we generate an optimal scaling curve by first identifying for each node number N the maximum performance over all n_{ranks} values. We then identify the rank giving the optimal performance given N , as shown in Fig. 7. Finally, we fit Amdahl’s law to these optimal scaling curves.

Note that with this definition of ideal scaling, we ignore that for mixed CPU-GPU nodes, the performance increase from one to two nodes is usually small compared to the perfor-

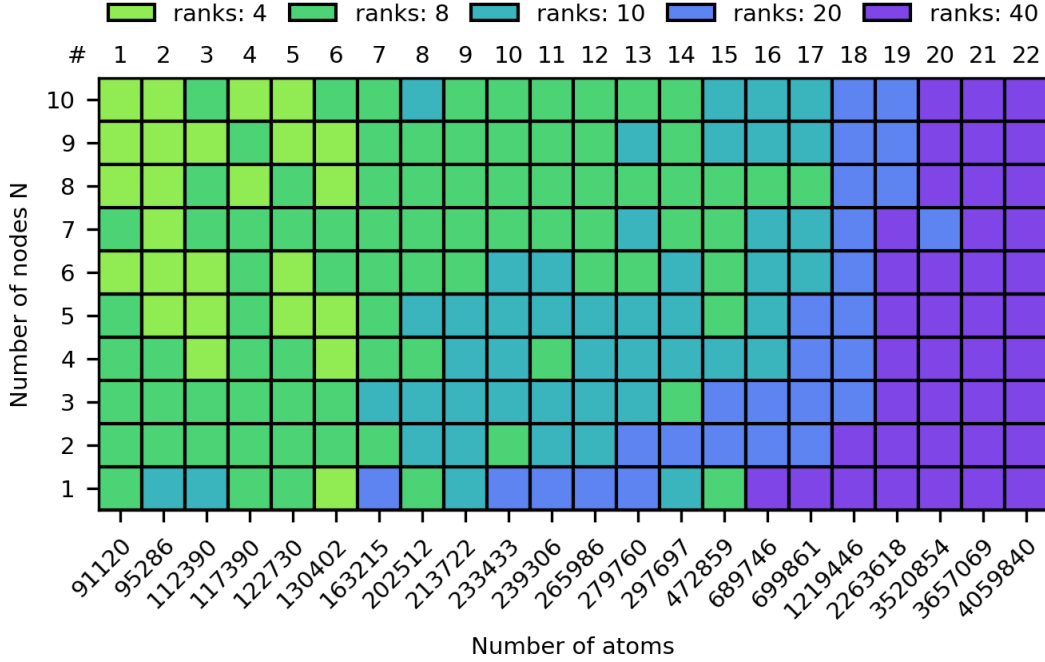


FIG. 7. On mixed CPU-GPU nodes, the optimal number of MPI ranks (colors) increases with system size and decreases with increasing node numbers N . The same systems (columns) were scaled over different numbers of nodes (rows), using different n_{ranks} values. The top axis indicates the system numbers according to Table I.

mance on one node. Only for systems larger than 2M atoms, the performance increase from one to two nodes is above 90% the performance of a single node. For the smallest systems, this increase is only about 35% (see Fig. 8a). In contrast, for CPU-only nodes this increase is closer to 100% for all system sizes (see Fig. 8b).

We next use the fits of Amdahl's law to the performance data to summarize the overall trends and the differences between the two hardware architectures (Fig. 10). We extract the values of the single fit parameter p and the maximum performance (Fig. 10a and b). We find that for both architectures the values of p increase with increasing system size. For CPU-only nodes, $p \gtrsim 0.9$ for all sizes. For mixed CPU-GPU nodes, p increases from ~ 0.5 for 10^5 atoms to > 0.9 for $\sim 4 \times 10^6$ atoms. Note that in all cases the values of p for CPU-only nodes are larger than for mixed CPU-GPU nodes. Although we cannot expect the estimates for the maximum performance to be highly accurate, they are useful to summarize the observed trends (Fig. 10b). For CPU-only nodes the maximum performance decreases from ~ 545 ns/day to ~ 157 ns/day with increasing system size. The maximum performance of mixed CPU-GPU nodes is always smaller than the maximum performance for CPU-only nodes and decreases from ~ 226 ns/day for the smallest system considered here to ~ 50 -80 ns/day for the largest system sizes.

The number of nodes, at which the performance is 70% of the ideal performance, increases with system size much more quickly for CPU-only nodes than for mixed CPU-GPU nodes (Fig. 10c). The estimates for these number are quite sensi-

tive to the quality of the fit. However, the performance values at these node numbers are rather robust estimates, as the performance changes only slowly with the number of nodes in this regime (Fig. 10d). These estimates show that the highest performance, i.e., shortest TTS, for a reasonable efficiency of $f = 0.7$ is achieved with CPU-only nodes in all cases.

Note that the relative performance increase going from one to two nodes, which determines ideal scaling, is only $\sim 35\%$ for mixed CPU-GPU nodes for the smallest system sizes (Fig. 10e). In contrast, for CPU-only nodes it is always larger than $\sim 70\%$. This relative performance increase calculated directly from the performance values agrees well with results from the fits of Amdahl's law.

VI. CONCLUSION

High-performance computing in general and MD simulations in particular are fast growing and highly dynamic fields. In a rapidly changing environment of hardware, software and systems, running MD simulations efficiently thus requires continuous benchmarking and monitoring of the simulation performance. The MDBenchmark toolkit presented here has been designed to simplify the benchmarking process. Its design is open to different MD engines and queuing systems, acknowledging the fact that it is becoming common practice that a single user uses different MD engines on various high-performance compute clusters.

The performance of an MD simulation depends on many

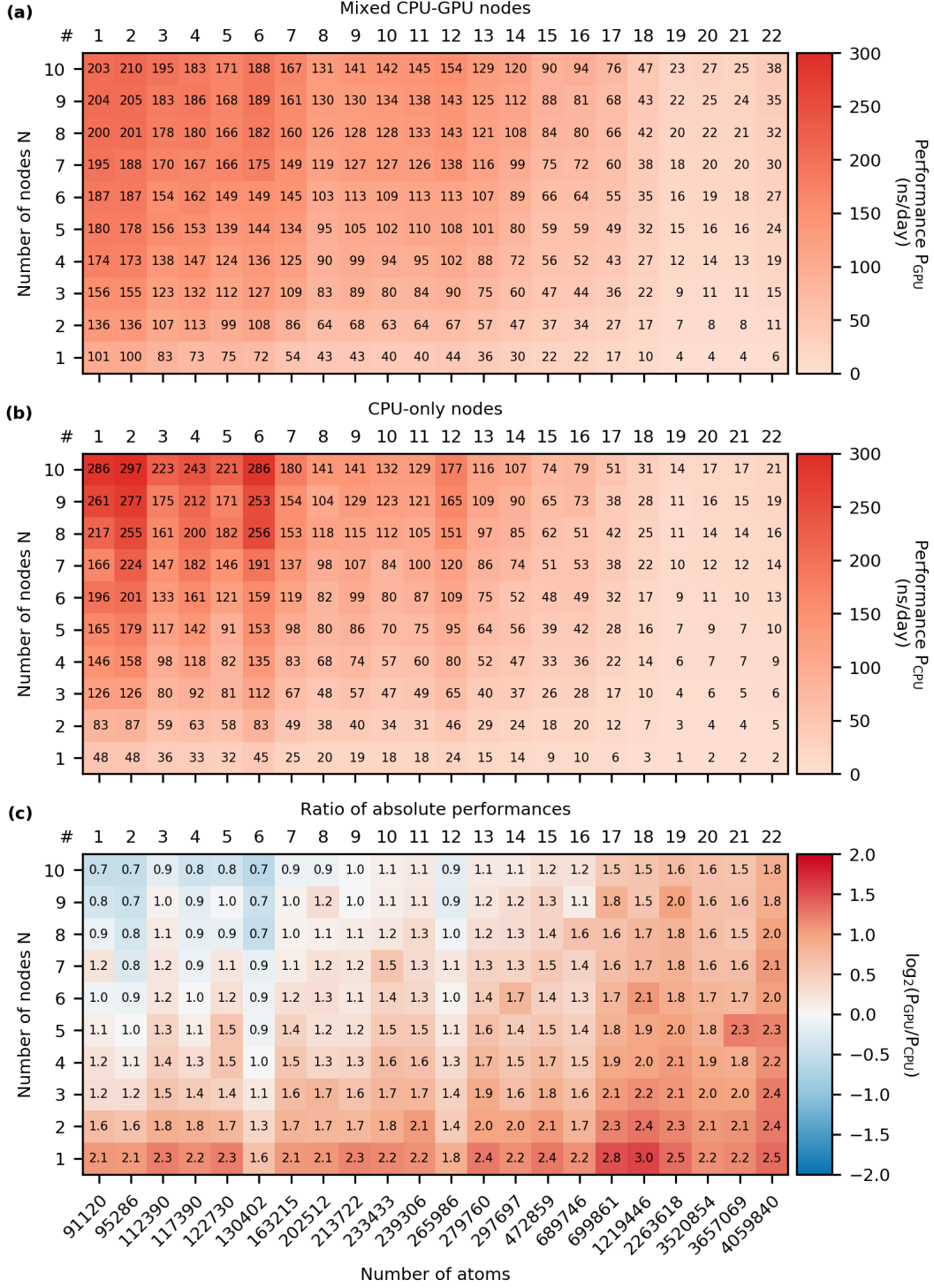


FIG. 8. Comparison of performance on (a) mixed CPU-GPU nodes and (b) CPU-only nodes. The numbers in the fields show the performance values in ns/day. (c) The performance ratio $P_{\text{GPU}}(N)/P_{\text{CPU}}(N)$ is shown on a logarithmic (base 2) color scale. The color scale is centered at equal performance ($P_{\text{GPU}}(N)/P_{\text{CPU}}(N) = 1$, white). The numbers in the cells are the actual values of the ratios. Performance data was taken for the best performing n_{ranks} values on mixed CPU-GPU nodes without hyperthreading (Fig. 7) for each node individually. For CPU-only nodes, $n_{\text{ranks}} = 40$ and hyperthreading was activated. The top axis indicates the system numbers according to Table I. Note that for system 22 we used a time step of 4 fs and doubled the hydrogen mass.

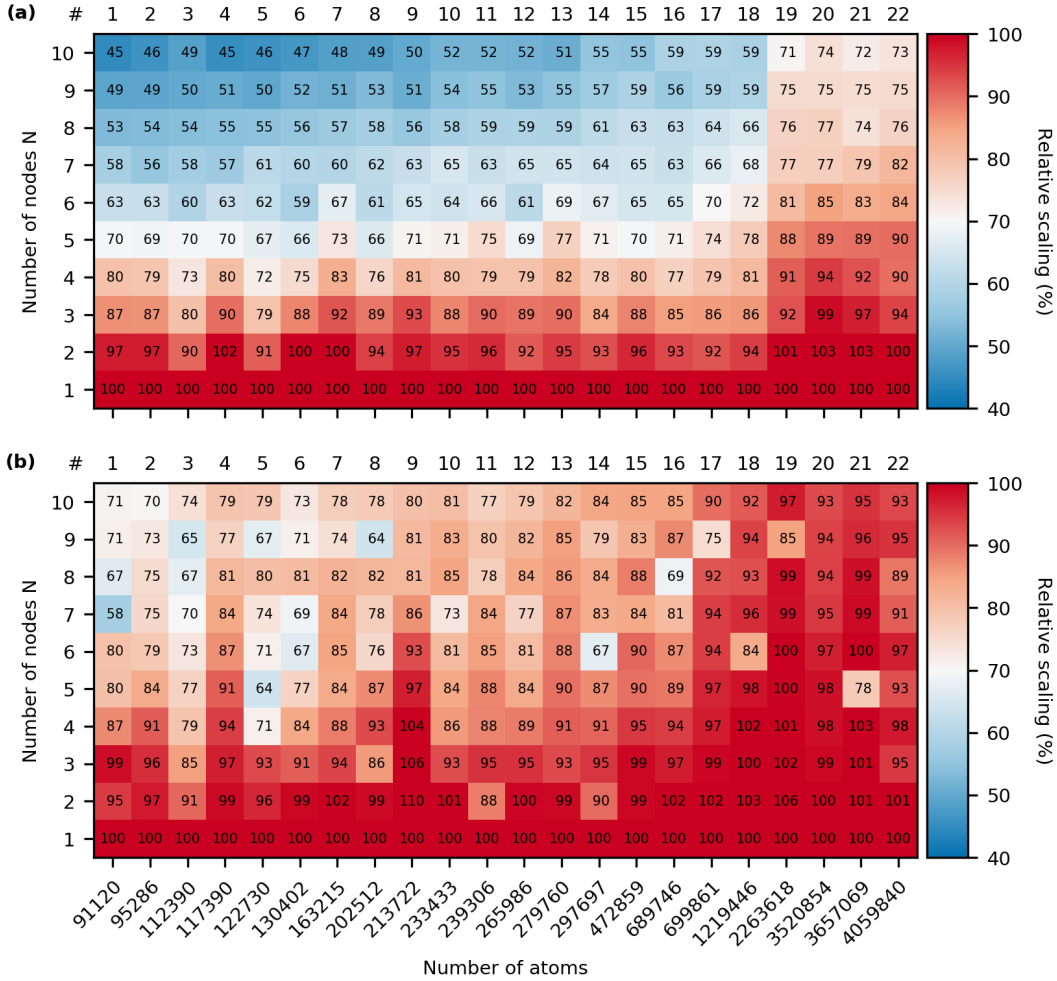


FIG. 9. Simulation performance scales differently on different architectures. Relative scaling of (a) mixed CPU-GPU nodes and (b) CPU-only nodes to their best performing settings. (a) Systems scale up to 70% of the ideal scaling as the number of atoms and nodes increase. (b) Systems scale to more than 70% of the ideal scaling as the number of atoms and nodes increase on CPU-only nodes. The numbers in each cell shows the scaling efficiency to the ideal scaling of the best performing setting, as shown in Fig. 7. The top axis indicates the system numbers according to Table I. Note that for system 22 we used a time step of 4 fs and doubled the hydrogen mass.

factors, of which only some are controlled by the user. Even for a given MD engine and hardware configuration, the performance is sensitive to the choice of the underlying algorithms and the simulation parameters. An example for the latter presented here is the increase in performance of the dense protein solutions compared to similarly sized systems, which was achieved by aggressively tuning simulation parameters. Activating an enhanced sampling method can affect the performance and the scaling with the number of nodes dramatically. While efficiency usually changes with the version of the MD engine itself, variations in the hardware drivers, compilers, and interfaces for parallelization can also have huge effects on the performance. For example, we observed that an update of the NVIDIA driver for the GPU increased performance on mixed CPU-GPU nodes by up to $\sim 20\%$.

Thus, our extensive performance scaling study surveyed only a small region in a high-dimensional parameter space.

For example, we have not explored the effects of offloading specific calculations, i.e., PME, to separate GPU ranks or systematically investigated how the tuning the cut-off parameters affects the performance. Nevertheless, our extensive quantification of the performance scaling provides guidelines for GROMACS users, reveals general trends, and serves as a point of reference for performance comparison, also for users of other MD engines.

Our results illustrate that benchmarking is necessary to find optimal parameters and to identify inefficiencies due to singular deviations from observed scaling trends. In the case of GROMACS 2018, for example, the proper choice of the number of MPI ranks is crucial. For mixed CPU-GPU nodes, this choice depends on system size and the number of nodes. Hyperthreading is generally beneficial for CPU-only runs. In contrast, for mixed CPU-GPU nodes, it depends on system size and node number whether hyperthreading leads to per-

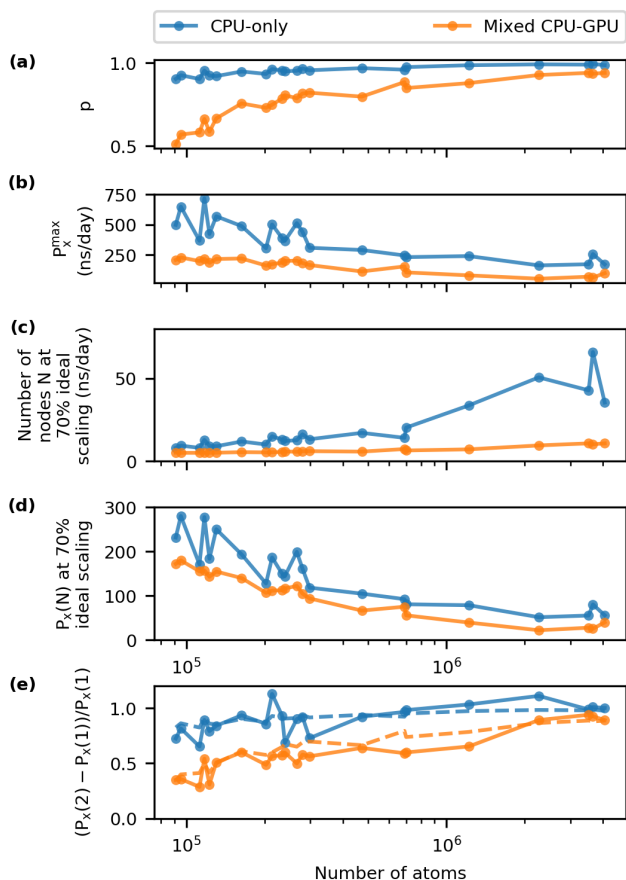


FIG. 10. Summary of the dependence on the performance on the atom number for CPU-only nodes (blue) and for mixed CPU-GPU nodes (orange). We fit Amdahl’s law to the scaling data for 1 to 10 nodes using the single fit parameter p . (a) The fitted values of p increase with increasing atom numbers. (b) The estimates for the maximum performance decrease with increasing systems size. CPU-only nodes always have larger values than mixed CPU-GPU nodes. (c) The number of nodes at which we reach 70% of the ideal performance and (d) the performance at 70% of the ideal scaling provide guidance for setting up benchmarks and simulations. (e) The relative performance increase when going from a single node to two nodes, directly calculated from the respective performance values. Dashed lines indicate the relative performance obtained from the fits of Amdahl’s law. Note that for the largest system (#22) we used a time step of 4 fs and doubled the hydrogen mass.

formance gains. Even though the results presented here serve as guidelines, confirming these settings for the specific simulation system and the resources available is inevitable.

In its current version, MDBenchmark can scan parameters like number of nodes, the numbers of MPI ranks and OpenMP threads, and the activation of hyperthreading, which are set when submitting the job. However, to achieve the best performance it is also necessary to tune parameters that are set in the configuration files of the respective MD engines. Currently, these configuration files have to be provided by the user. Ideally, future versions of MDBenchmark perform scans over simulation parameters specified in the configuration files

automatically and validate the results.

The monetary and environmental costs of molecular simulations are significant and even small relative performance improvements have large absolute effect on overall cost and efficiency. For many, MDBenchmark might be a first step to start continuously monitoring and evaluating the efficient use of their hardware resources. It is fair to assume that within a typical research group, with a mix of members with essentially no experience and members who are experts in running simulations, resources can be easily wasted if insufficient attention is paid to monitoring simulation efficiency. Running benchmarks as a rule for any new system and the comparison with existing benchmarks could greatly reduce the risk of wasteful use of resources.

Thus, ideally, a tool like MDBenchmark additionally collects the benchmark information generated by the users in a single database accessible to all. This information should be supplemented by the actual performance data of production runs. In principle, the MDBenchmark toolkit could be already used to perform production runs and thus automatically collect performance information. With this kind of information, inefficiencies can be identified quickly and the database can provide accurate guidelines for setting up simulations. Such guidelines also serve to keep the number of necessary benchmarks to a minimum.

The design of MDBenchmark embraces the philosophy that we should always choose the best tool for the task at hand. Ideally, we can easily switch between different MD engines to take advantage of their unique features^{7–14}. MDBenchmark is open to all MD engines (and queuing systems). We hope that the community will appreciate the design and capabilities of the provided framework to run and analyze benchmarks such that they contribute their expert knowledge by adding their favorite MD engines and queuing systems. Ultimately, running simulations more efficiently translates into doing better science.

The source code of MDBenchmark is freely available under the GPLv3 license at <https://github.com/bio-phys/mdbenchmark>. It can be installed either via the *pip* or *conda* package managers using the PyPI or conda-forge repository, respectively. The code is accompanied by an extensive documentation that is hosted at <https://mdbenchmark.readthedocs.io/en/version-3/>. Detailed instructions for adding currently unsupported MD engines can also be found in the documentation.

ACKNOWLEDGMENTS

We thank Drs. Markus Rampp, Klaus Reuter, and Sebastian Kehl for technical support and useful discussions. We thank Dr. Florian Blanc, Sören von Bülow, Daniel Chavez Rojas, Dr. Roberto Covino, Sergio Cruz, Dr. Sonya Hanson, Dr. Ahmadreza Mehdipour, Laura Schulz, and Jan Stuke for providing molecular dynamics systems. This study used the high-performance computing resources of the Max Planck Computing and Data Facility (MPCDF). We acknowledge financial support by the Max Planck Society and the Landes-Offensive

zur Entwicklung Wissenschaftlich-ökonomischer Exzellenz (LOEWE) DynaMem program of the state of Hesse (MG, MS and GH).

DATA AVAILABILITY STATEMENT

The data that support the findings of this study are available from the corresponding author upon reasonable request.

- ¹G. E. Moore, *Electronics* **38** (1965).
- ²R. Kurzweil, *The Age of Spiritual Machines: When Computers Exceed Human Intelligence*, A Penguin Book. Science/Technology (Penguin Books, 2000).
- ³M. Vendruscolo and C. M. Dobson, *Curr. Biol.* **21**, R68 (2011).
- ⁴R. O. Dror, M. Ø. Jensen, D. W. Borhani, and D. E. Shaw, *J. Gen. Physiol.* **135**, 555 (2010).
- ⁵E. H. Lee, J. Hsin, M. Sotomayor, G. Comellas, and K. Schulten, *Structure* **17**, 1295 (2009).
- ⁶R. O. Dror, R. M. Dirks, J. Grossman, H. Xu, and D. E. Shaw, *Annu. Rev. Biophys.* **41**, 429 (2012).
- ⁷M. J. Harvey, G. Giupponi, and G. D. Fabritiis, *J. Chem. Theory Comput.* **5**, 1632 (2009).
- ⁸R. Salomon-Ferrer, A. W. Götz, D. Poole, S. Le Grand, and R. C. Walker, *J. Chem. Theory Comput.* **9**, 3878 (2013).
- ⁹B. R. Brooks, C. L. Brooks, A. D. Mackerell, L. Nilsson, R. J. Petrella, B. Roux, Y. Won, G. Archontis, C. Bartels, S. Boresch, A. Caflisch, L. Caves, Q. Cui, A. R. Dinner, M. Feig, S. Fischer, J. Gao, M. Hodosek, W. Im, K. Kuczera, T. Lazaridis, J. Ma, V. Ovchinnikov, E. Paci, R. W. Pastor, C. B. Post, J. Z. Pu, M. Schaefer, B. Tidor, R. M. Venable, H. L. Woodcock, X. Wu, W. Yang, D. M. York, and M. Karplus, *J. Comput. Chem.* **30**, 1545 (2009).
- ¹⁰K. J. Bowers, D. E. Chow, H. Xu, R. O. Dror, M. P. Eastwood, B. A. Gregersen, J. L. Klepeis, I. Kolossvary, M. A. Moraes, F. D. Sacerdoti, J. K. Salmon, Y. Shan, and D. E. Shaw, in *SC '06: Proceedings of the 2006 ACM/IEEE Conference on Supercomputing* (Tampa, FL, 2006) pp. 43–43.
- ¹¹M. J. Abraham, T. Murtola, R. Schulz, S. Páll, J. C. Smith, B. Hess, and E. Lindahl, *SoftwareX* **1-2**, 19 (2015).
- ¹²W. M. Brown, A. Kohlmeier, S. J. Plimpton, and A. N. Tharrington, *Comput. Phys. Commun.* **183**, 449 (2012).
- ¹³J. C. Phillips, R. Braun, W. Wang, J. Gumbart, E. Tajkhorshid, E. Villa, C. Chipot, R. D. Skeel, L. Kalé, and K. Schulten, *J. Comput. Chem.* **26**, 1781 (2005).
- ¹⁴P. Eastman, M. S. Friedrichs, J. D. Chodera, R. J. Radmer, C. M. Bruns, J. P. Ku, K. A. Beauchamp, T. J. Lane, L.-P. Wang, D. Shukla, T. Tye, M. Houston, T. Stich, C. Klein, M. R. Shirts, and V. S. Pande, *J. Chem. Theory Comput.* **9**, 461 (2013).
- ¹⁵A. M. Bonvin, A. E. Mark, and W. F. van Gunsteren, *Comput. Phys. Commun.* **128**, 550 (2000).
- ¹⁶C. C. Gruber and J. Pleiss, *J. Comput. Chem.* **32**, 600 (2011).
- ¹⁷C. Kutzner, S. Páll, M. Fechner, A. Esztermann, B. L. de Groot, and H. Grubmüller, *J. Comput. Chem.* **36**, 1990 (2015).
- ¹⁸C. Kutzner, S. Páll, M. Fechner, A. Esztermann, B. L. de Groot, and H. Grubmüller, *J. Comput. Chem.* **40**, 2418 (2019).
- ¹⁹D. T. Marr, F. Binns, D. L. Hill, G. Hinton, D. A. Koufaty, J. A. Miller, and M. Upton, *Intel Technol. J.* **6** (2002).
- ²⁰B. Hess, C. Kutzner, D. van der Spoel, and E. Lindahl, *J. Chem. Theory Comput.* **4**, 435 (2008).
- ²¹J. L. Furlani, in *Proceedings of the fifth large installation systems administration conference (LISA V)* (San Diego, CA, 1991) pp. 141–152.
- ²²D. L. Dotson, S. L. Seyler, M. Linke, R. J. Gowers, and O. Beckstein, in *Proceedings of the 15th Python in Science Conference*, edited by Sebastian Benthall and Scott Rostrop (Austin, TX, 2016) pp. 51–56.
- ²³W. McKinney, in *Proceedings of the 9th Python in Science Conference*, edited by S. van der Walt and J. Millman (Austin, TX, 2010) pp. 51–56.
- ²⁴G. M. Amdahl, *SSCS* **12**, 19 (2007).
- ²⁵N. Michaud-Agrawal, E. J. Denning, T. B. Woolf, and O. Beckstein, *J. Comput. Chem.* **32**, 2319 (2011).
- ²⁶R. J. Gowers, M. Linke, J. Barnoud, T. J. E. Reddy, M. N. Melo, S. L. Seyler, J. Domański, D. L. Dotson, S. Buchoux, I. M. Kenney, and O. Beckstein, in *Proceedings of the 15th Python in Science Conference*, edited by Sebastian Benthall and Scott Rostrop (Austin, TX, 2016) pp. 98–105.
- ²⁷S. van der Walt, S. C. Colbert, and G. Varoquaux, *Comput Sci Eng* **13**, 22 (2011).
- ²⁸P. Virtanen, R. Gommers, T. E. Oliphant, M. Haberland, T. Reddy, D. Cournapeau, E. Burovski, P. Peterson, W. Weckesser, J. Bright, S. J. van der Walt, M. Brett, J. Wilson, K. Jarrod Millman, N. Mayorov, A. R. J. Nelson, E. Jones, R. Kern, E. Larson, C. Carey, Í. Polat, Y. Feng, E. W. Moore, J. VanderPlas, D. Laxalde, J. Perktold, R. Cimrman, I. Henriksen, E. A. Quintero, C. R. Harris, A. M. Archibald, A. H. Ribeiro, F. Pedregosa, P. van Mulbregt, and SciPy 1.0 Contributors, *Nat. Methods* (2020).
- ²⁹F. Perez and B. E. Granger, *Comput Sci Eng* **9**, 21 (2007).
- ³⁰J. D. Hunter, *Comput Sci Eng* **9**, 90 (2007).
- ³¹R. B. Best, X. Zhu, J. Shim, P. E. Lopes, J. Mittal, M. Feig, and A. D. MacKerell Jr, *J. Chem. Theory Comput.* **8**, 3257 (2012).
- ³²J. Huang, S. Rauscher, G. Nawrocki, T. Ran, M. Feig, B. L. de Groot, H. Grubmüller, and A. D. MacKerell, *Nat. Methods* **14**, 71 (2017).
- ³³J. Wang, P. Cieplak, and P. A. Kollman, *J. Comput. Chem.* **21**, 1049 (2000).
- ³⁴V. Hornak, R. Abel, A. Okur, B. Strockbine, A. Roitberg, and C. Simmerling, *Proteins* **65**, 712 (2006).
- ³⁵K. Lindorff-Larsen, S. Piana, K. Palmo, P. Maragakis, J. L. Klepeis, R. O. Dror, and D. E. Shaw, *Proteins* **78**, 1950 (2010).
- ³⁶R. B. Best and G. Hummer, *J. Phys. Chem. B* **113**, 9004 (2009).
- ³⁷R. B. Best, D. de Sancho, and J. Mittal, *Biophys. J.* **102**, 1462 (2012).
- ³⁸W. L. Jorgensen, J. Chandrasekhar, J. D. Madura, R. W. Impey, and M. L. Klein, *J. Chem. Phys.* **79**, 926 (1983).
- ³⁹S. Piana, A. G. Donchev, P. Robustelli, and D. E. Shaw, *J. Phys. Chem. B* **119**, 5113 (2015).
- ⁴⁰H. F. Hofbauer, M. Gecht, S. C. Fischer, A. Seybert, A. S. Frangakis, E. H. K. Stelzer, R. Covino, G. Hummer, and R. Ernst, *J. Cell Biol.* **217**, 3109 (2018).
- ⁴¹S. v. Bülow, M. Siggel, M. Linke, and G. Hummer, *PNAS* **116**, 9843 (2019).
- ⁴²X. Wu, M. Siggel, S. Ovchinnikov, W. Mi, V. Svetlov, E. Nudler, M. Liao, G. Hummer, and T. A. Rapoport, *Science* **368** (2020).
- ⁴³S. Hofmann, D. Janulienė, A. R. Mehdipour, C. Thomas, E. Stefan, S. Brüchert, B. T. Kuhn, E. R. Geertsma, G. Hummer, R. Tampé, and A. Moeller, *Nature* **571**, 580 (2019).

Appendix A: Appendix

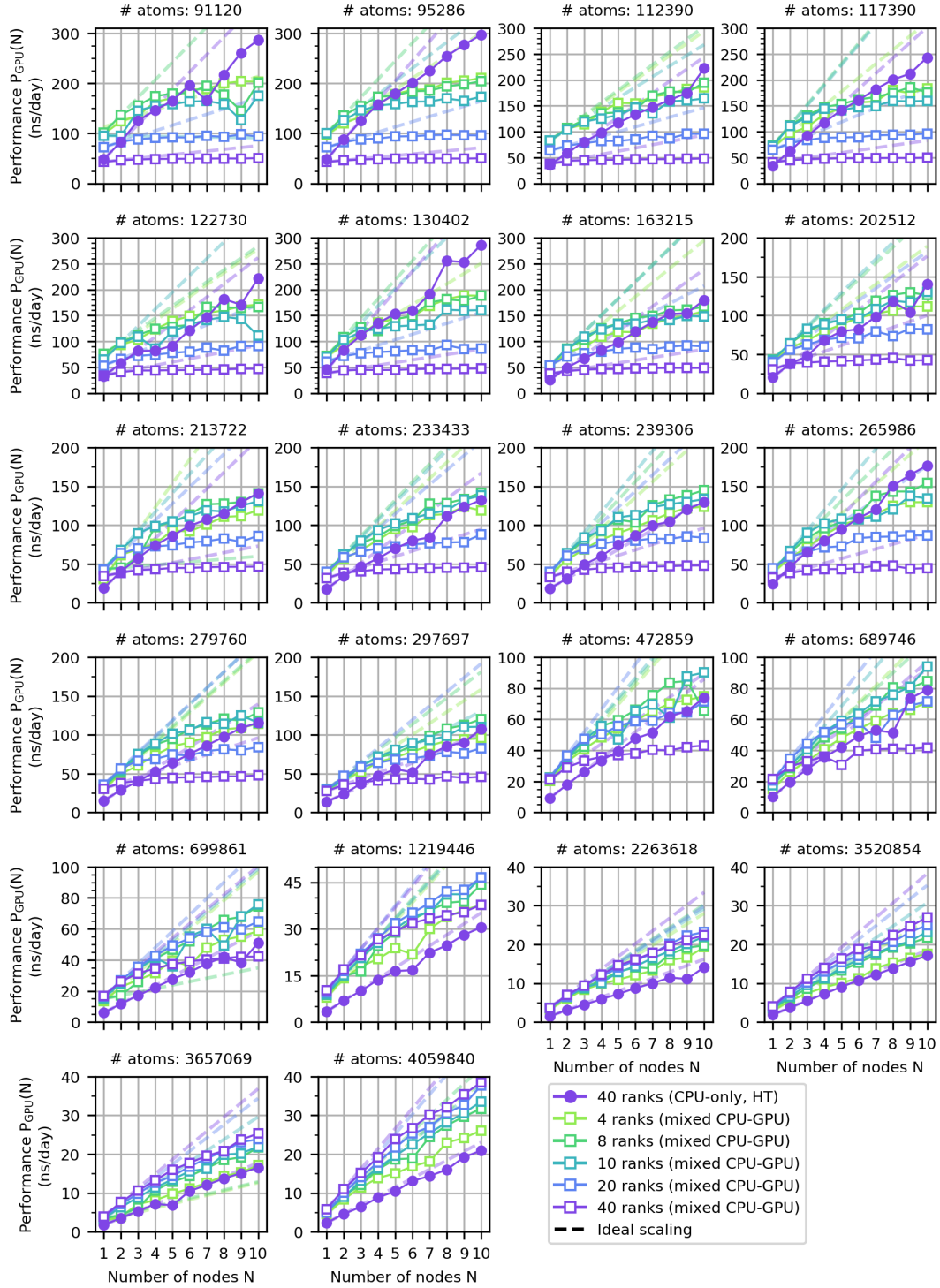


FIG. S1. For mixed CPU-GPU nodes, the optimal choice of n_{ranks} depends on system size. Performance $P_{\text{GPU}}(N)$ of 22 MD simulations with varying sizes as a function of nodes N for different numbers of MPI ranks, n_{ranks} on mixed CPU-GPU nodes without hyperthreading (colored squares). The best performing setting with CPU-only nodes is shown as reference with circles. Each data point in the performance plot shows one independent run of 15 minutes each. Transparent dashed lines show ideal scaling as estimated by the performance of one and two nodes. Note that for system 22 we used a time step of 4 fs and doubled the hydrogen mass.

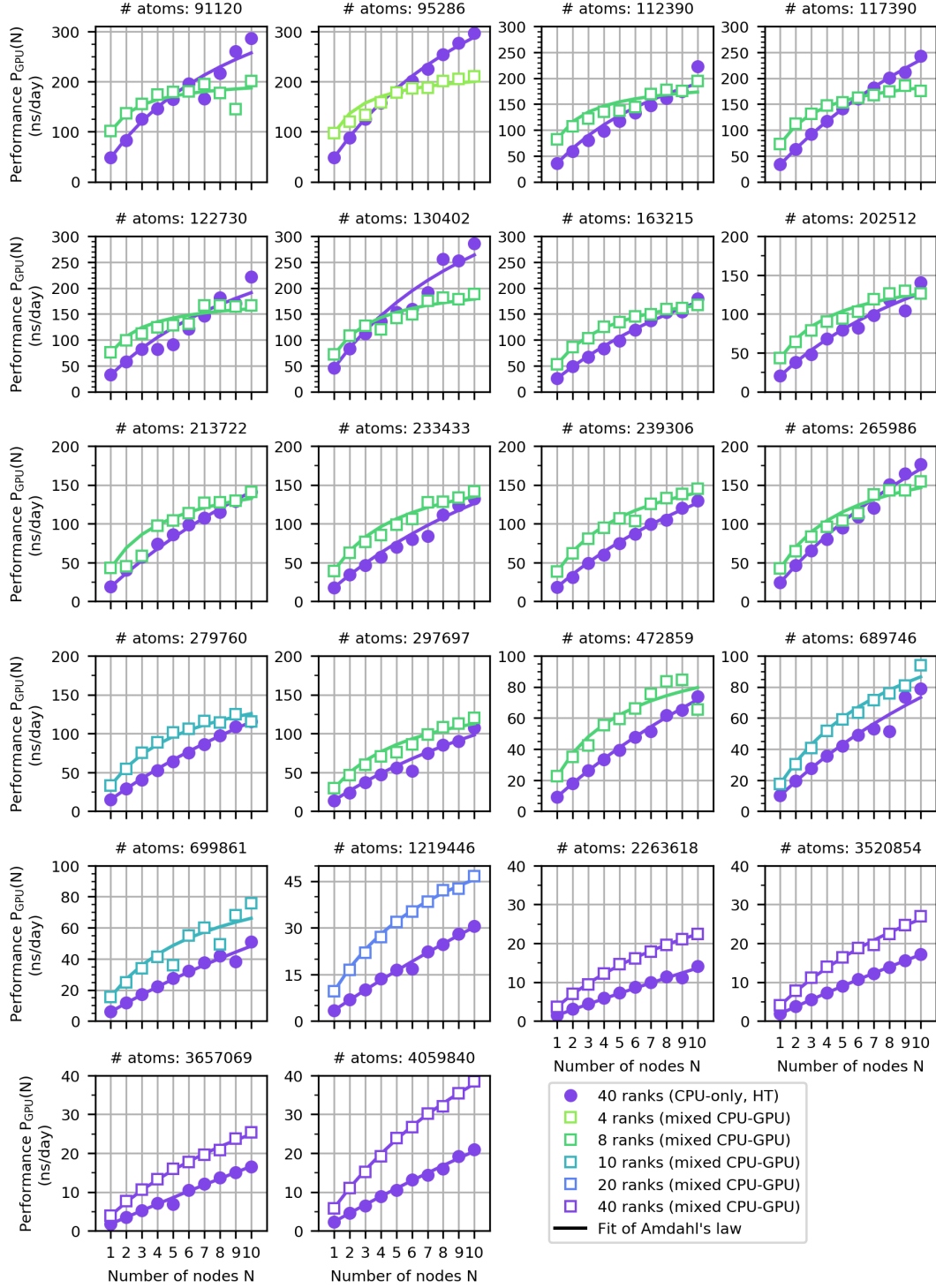


FIG. S2. GPUs scale to higher N with bigger systems. Performance $P(N)$ of 22 MD simulations with varying sizes as a function of nodes N . Each panel shows the best performing number of MPI ranks, n_{ranks} on mixed CPU-GPU nodes without hyperthreading (colored squares) and on CPU-only nodes with 40 n_{ranks} and hyperthreading (filled circles) from Fig. S1. The solid colored line shows the interpolated fit to the data using Eq. (2). The intersection of both curves increases to higher N with increasing number of atoms. Note that for system 22 we used a time step of 4 fs and doubled the hydrogen mass.

```

1 #!/bin/bash -l
2 #SBATCH -o ./benchjob.out.%j
3 #SBATCH -e ./benchjob.err.%j
4 #SBATCH -D ./
5 #SBATCH -J n001_r40_t02_wht
6 #
7 #SBATCH --nodes=1
8 # Set the number of tasks per node (=MPI ranks)
9 #SBATCH --ntasks-per-node=40
10 # Set the number of threads per rank (=OpenMP threads)
11 #SBATCH --cpus-per-task=2
12 # Enable hyperthreading
13 #SBATCH --ntasks-per-core=2
14 #SBATCH --time=00:17:00
15
16 module purge
17 module load gcc
18 module load impi
19 module load cuda
20 module load gromacs/2018.8
21
22 # Set number of OpenMP threads and proper core pinning with hyperthreading
23 export OMP_NUM_THREADS=$SLURM_CPUS_PER_TASK
24 export OMP_PLACES=threads
25 export SLURM_HINT=multithread
26
27 # Run gromacs/2018.8 for 15 minutes
28 srun gmx_mpi mdrun -v -ntomp $OMP_NUM_THREADS -maxh 0.25 -resethway -deffnm md -noconfout

```

Listing 1. SLURM submission script for a 15 minutes run on a CPU-only single node with 40 MPI ranks, 2 OpenMP threads and hyperthreading enabled.

```

1 #!/bin/bash -l
2 #SBATCH -o ./benchjob.out.%j
3 #SBATCH -e ./benchjob.err.%j
4 #SBATCH -D ./
5 #SBATCH -J n001_r20_t02_woht
6 #
7 #SBATCH --constraint="gpu"
8 #SBATCH --gres=gpu:rtx5000:2
9 #
10 #SBATCH --nodes=1
11 # Set the number of tasks per node (=MPI ranks)
12 #SBATCH --ntasks-per-node=20
13 # Set the number of threads per rank (=OpenMP threads)
14 #SBATCH --cpus-per-task=2
15 #SBATCH --time=00:17:00
16
17 module purge
18 module load gcc
19 module load impi
20 module load cuda
21 module load gromacs/2018.8
22
23 # Set number of OpenMP threads and proper core pinning without hyperthreading
24 export OMP_NUM_THREADS=$SLURM_CPUS_PER_TASK
25 export OMP_PLACES=cores
26
27 # Run gromacs/2018.8 for 15 minutes
28 srun gmx_mpi mdrun -v -ntomp $OMP_NUM_THREADS -maxh 0.25 -resethway -deffnm md -noconfout

```

Listing 2. SLURM submission script for a 15 minutes run on a single mixed CPU-GPU node with 20 MPI ranks, 2 OpenMP threads and hyperthreading disabled.