

MolSieve: A Progressive Visual Analytics System for Molecular Dynamics Simulations

Rostyslav Hnatyshyn, Jieqiong Zhao, Danny Perez, James Ahrens, Ross Maciejewski

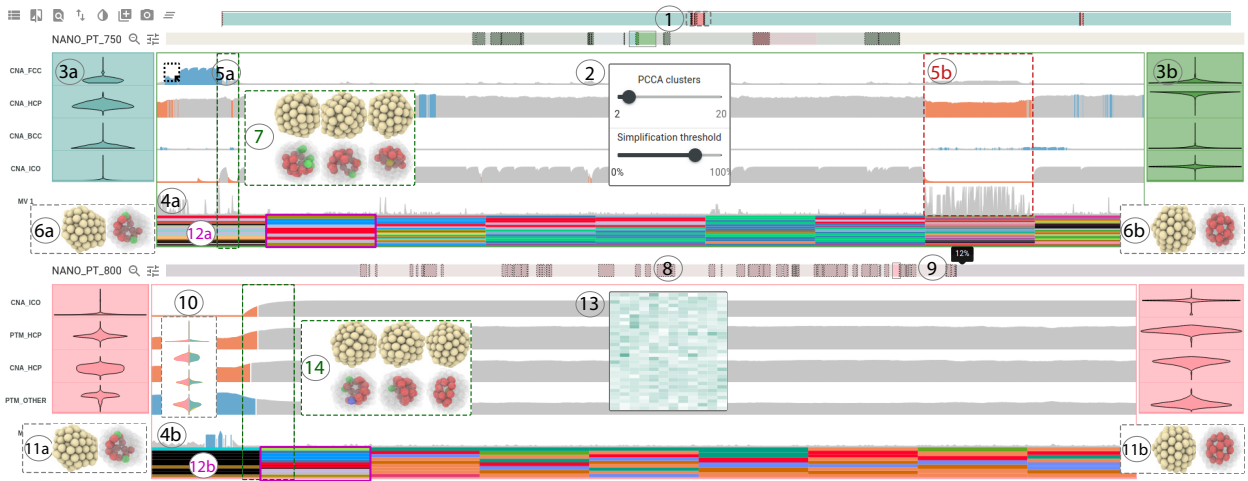


Fig. 1: A sample analysis performed in MolSieve. Two long-duration simulations of nano-particles at 750 (top) and 800 kelvins (bottom) are shown. (1) is the Timeline View of the top trajectory before zooming. (2) shows a menu for adjusting exploratory parameters of each trajectory. (3a) and (3b) are Super-State Views; (6a) and (6b) are their characteristic states. (4a) and (4b) are multi-variate control charts that were dynamically added during the analysis. (5a) (green dashed outline) and (5b) are regions of interest within a Transition Region View. (8) shows the visible extents of the 800K trajectory. (9) shows the results of using the **Find similar regions** button on the transition region within the 750K simulation. (10) is the Region Comparison Widget for comparing the teal and pink super-states. (12a) and (12b) show State Space Charts; the **boxes** display similarities between the two regions. (13) is a Sub-Sequence Comparison Widget comparing the sequences (7) and (14). An analyst has found a structural change in common between the simulations: (6a), (7), and (6b) detail the change in the 750K simulation and (11a), (14), and (11b) detail the change at the 800K simulation. (6a) and (11a) are the same state, and (6b) and (11b) are rotations of each other.

Abstract—Molecular Dynamics (MD) simulations are ubiquitous in cutting-edge physio-chemical research. They provide critical insights into how a physical system evolves over time given a model of interatomic interactions. Understanding a system's evolution is key to selecting the best candidates for new drugs, materials for manufacturing, and countless other practical applications. With today's technology, these simulations can encompass millions of unit transitions between discrete molecular structures, spanning up to several milliseconds of real time. Attempting to perform a brute-force analysis with data-sets of this size is not only computationally impractical, but would not shed light on the physically-relevant features of the data. Moreover, there is a need to analyze simulation ensembles in order to compare similar processes in differing environments. These problems call for an approach that is analytically transparent, computationally efficient, and flexible enough to handle the variety found in materials-based research. In order to address these problems, we introduce MolSieve, a progressive visual analytics system that enables the comparison of multiple long-duration simulations. Using MolSieve, analysts are able to quickly identify and compare regions of interest within immense simulations through its combination of control charts, data-reduction techniques, and highly informative visual components. A simple programming interface is provided which allows experts to fit MolSieve to their needs. To demonstrate the efficacy of our approach, we present two case studies of MolSieve and report on findings from domain collaborators.

Index Terms—Molecular dynamics, time-series analysis, visual analytics

1 INTRODUCTION

Molecular dynamics (MD) simulations allow scientists to observe how systems of atoms evolve over time using a potential energy function

- R. Hnatyshyn, J. Zhao and R. Maciejewski are with Arizona State University. E-mail: {rhnatysh,jzhao,rmacieje}@asu.edu.
- D. Perez and J. Ahrens are with Los Alamos National Laboratory. E-mail: {ahrens,danny_perez}@lanl.gov

Manuscript received 31 March 2023; revised 1 July 2023; accepted 8 August 2023.
Date of publication 8 November 2023; date of current version 21 December 2023.
This article has supplementary downloadable material available at <https://doi.org/10.1109/TVCG.2023.3326584>, provided by the authors.
Digital Object Identifier no. 10.1109/TVCG.2023.3326584

that calculates interatomic forces. Understanding the nanoscale behavior of matter has widespread applications, from guiding protein mutations in bio-medical research [24] to validating the robustness of a material in engineering contexts [32]. A large family of software packages have been developed in order to generate MD simulations, such as GROMACS [8] for biological simulations, LAMMPS [46] for materials modeling, as well as countless others, e.g. [27, 37]. A recently introduced simulation management tool called ParSplice [34] has enabled MD simulations to span time-scales reaching into the hundreds of thousands of nano-seconds (milliseconds), two orders of magnitude larger than simulations typically performed with biological systems. The time-scales ParSplice is able to simulate typically contain millions of discrete transitions between molecular configurations. Some of these

systems suffer from the heterogeneous energy barrier problem [34], a prevalent issue in long MD trajectories [35, 36].

Trajectories with a heterogeneous energy barrier distribution are difficult to analyze since relevant regions within a trajectory are buried amongst a myriad of repetitive transitions within so-called super-states. While calculating all of the energy paths between every state and visualizing them seems like a solution at first glance, not only are the computational costs involved impractical, but the results generated by this method are impossible to sift through manually. To further compound the problem, ParSplice generates trajectories as ensembles because MD is inherently a stochastic process; attempting to generalize the behavior of a system from an individual simulation could lead to brittle conclusions. These issues dictate the need to develop an analysis tool that highlights the essential components of a trajectory (i.e., its transition regions), while understating the parts of a trajectory where there is little to no change in the structure of the system (i.e., its super-states), as well as facilitating comparisons between trajectories.

A number of visual analytics systems enable the exploration of molecular dynamics simulations, e.g. [11, 14, 23, 31, 49]. However, most existing systems focus on biological simulations, which typically do not involve the same time-scales as their inorganic counterparts, rendering them impractical for analyzing the data-sets produced by ParSplice. To address this gap, we worked closely with domain experts to develop MolSieve, a visual analytics system that aggressively reduces molecular dynamics simulations to their essential components (super-states and transition regions) to facilitate their analysis and comparison. To evaluate the efficacy of MolSieve, we performed two case studies with materials science experts on data-sets from their daily workflows. They demonstrate that our system is not only efficient in extracting insight but is also adaptable to an expert's needs. This work contributes:

- A novel combination of coordinated multiple views consisting of temporal charts for examining long sequences by distinguishing regions of interest and uninteresting regions;
- A novel state space chart for visualizing discrete temporal events in a limited screen space while outlining their general trend;
- An efficient, scalable, and customizable progressive visual analytics system that supports analyzing large materials MD trajectory ensembles in real-time with the aforementioned visual designs.

2 RELATED WORK

In this section, we review various methods to analyze long-duration molecular dynamics simulations. We also discuss the visualization techniques and analytical methods that inspired our system.

2.1 Molecular Dynamics Analysis Approaches

Many approaches exist for exploring long-duration molecular dynamics trajectories which utilize various methods of reducing the data-set to a size tractable for real-time analysis. We found that these approaches are typically tailored for specific analyses of biological systems. For example, PyContact [41] enables the exploration of non-covalent interactions within molecular dynamics trajectories. It aims to provide access to points of interest within the trajectory by filtering on the amount of contact molecules within the simulation at any given time-step. However, PyContact requires the calculation of every molecular contact before the data-set can be analyzed, which can be time-consuming. VIA-MD [44] allows the exploration of long duration biological molecular systems through a combination of linked 2D and 3D views, which work together to highlight events of interest in both the spatial and temporal domains. Our proposed solution differs in locating regions of interest due to the difference in scale – VIA-MD was tested on a biological simulation that spanned twenty-three nano-seconds, while our case studies average five thousand nano-seconds. To extract insights from data-sets of this size, we developed a unique data simplification scheme based on the internal dynamics of the simulation. To the best of our knowledge, this simplification scheme has not yet been explored. ExaViz [14] enables the in-situ analysis of biological molecular systems. This in-situ approach reduces the data-set by allowing experts to decide what portions of the trajectory are relevant before saving them for long-term

storage, which requires a tremendous amount of computing power and tedious manual analysis. Byška et al. [11] built a focus+context visual analytics system that tied statistical properties of simulations to their 3D renders. Building on this work, sMolBoxes [49] utilized a data-flow model embedded in CAVER [23] to identify important snapshots within long duration bio-molecular simulations. sMolBoxes identifies important snapshots (states) within a trajectory by relying on domain specific information provided by analysts, e.g., using the root-mean-square deviation (RMSD) between states to identify abnormal structural changes in proteins. Analysts are able to select individual parts of a protein to track throughout the trajectory. Unfortunately, this powerful interaction is inherently coupled with the spatial dimension of the data, which reduces its scope to biological systems. Duran et al. [15] explore building a similar system using traditional statistical charts and linking them to a 3D visualization of the protein being studied. Non-biological systems do not behave in the same manner as proteins, reducing the effectiveness of these approaches as a general solution to identifying regions of interest within a molecular dynamics simulation. Chae et al. [12] used a deep learning model to reduce the dimensionality of a molecular dynamics simulation to a 3D space for easier exploration, using multiple views to display the original data alongside the 3D embedding. LaSCA [47] is a visual analytics system which identifies crystalline structures within large molecular systems in great detail; however, the system does not support analyzing these structures within the context of a MD trajectory. Wu et al. [52] proposed a visualization pipeline to identify point defects in nuclear materials – as with LaSCA, this approach does not consider the trajectory as a whole. To the best of our knowledge, the visual analytics systems currently available do not offer an efficient method to identify and compare analyst-defined regions of interest within MD simulations of materials.

A number of programming tool-kits also provide solutions for MD trajectories [10, 28, 38, 45]. However, these tool-kits cannot identify regions of interest within a trajectory without being integrated into a larger framework. Blindly applying these tool-kits to long simulations without a scheme to filter and organize their output will simply produce large bodies of data that are difficult to interpret.

2.2 Visual Analytics Methods for Time-series Exploration

In this section, we discuss several works that directly inspired views in MolSieve. Tominski et al. [48] developed a multi-attribute temporal view for a spatial trajectory by stacking horizon charts representing each attribute. This stacked trajectory chart is then rendered on top of 3D map data to facilitate a spatio-temporal analysis of the data-set. DQNViz [50] also took a similar approach to visualize multi-variate sequence data by stacking line charts, bar charts, and area charts on top of each other to provide a multi-dimensional view of the behavior of a machine learning model. Additionally, their approach provides methods to identify and compare patterns within the trajectory using segment mining and dynamic time warping. MolSieve does not use dynamic time warping for comparing sequences, as the structure of a system is far too complex to be modeled by dynamic time warping; instead, we use domain-specific methods to compare analyst-defined regions. Our approach combines the visual elements of the aforementioned systems and uses trajectory information to generate and arrange charts based on the detected importance of a region. SignalLens [26] uses a distorted scale where interesting parts of an electronic signal are magnified while uninteresting regions are minimized in their sequence view. Regardless of the level of distortion, context is maintained, which is essential to navigating long time-series on a screen limited by size. MolSieve distorts the trajectory's sequence to emphasize transition regions while minimizing super-states. For a comprehensive review of time series visualization techniques, we refer to Aigner et al [5].

3 ANALYTICAL TASKS, REQUIREMENTS, AND DEFINITIONS

In this section, we define tasks for MD analysis, the requirements for an analytical tool, and domain-specific definitions.

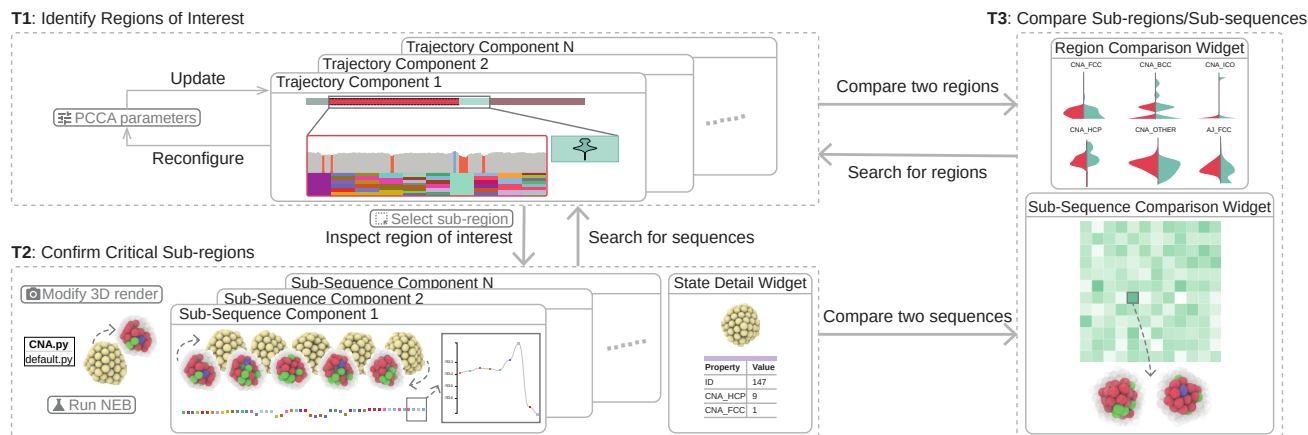


Fig. 2: MolSieve is designed to extract insight from MD simulations in three stages. First, an analyst uses a modal window to set the system's exploratory parameters. When the initial simplification completes, analyst-defined properties are mined on portions of the trajectory and progressively rendered in a Trajectory Component. While the properties are being calculated, the analyst uses the Trajectory Component's embedded views to interactively identify regions of interest (T1). If no regions of interest are found, the exploratory parameters can be reconfigured. If the analyst finds a sub-region of interest, they can select and examine it in detail using the Sub-Sequence Component and the State Detail Widget (T2). The analyst can use the comparison interactions provided by MolSieve to explore other trajectories in the context of their new discovery (T3).

3.1 Definitions

ParSplice simulations typically generate tens to hundreds of thousands of unique configurations of the system being simulated, with each discrete configuration being referred to as a *state*. Each configuration has its own state ID. These states contain meta-data about the system being simulated at a given point in time, such as the positions, chemical species, velocities, etc., of its atoms. This meta-data can be used to calculate properties that characterize its structure and geometry.

A *trajectory* is a sequence of states, and a single trajectory describes one of the many possible ways a system can evolve. To *transition* between two states, the system must overcome the *energy barrier* between them; therefore, state transitions that have a low energy barrier tend to occur exponentially more frequently than state transitions associated with larger barriers. This causes states to repeat throughout a trajectory, since structurally similar states are easier to transition to than radically different ones. Each transition has a discrete time-step associated with it to organize it temporally; transitions take a variable amount of time, but usually they usually occur in the span of hundreds of picoseconds.

The frequency of low energy transitions causes trajectories to often get trapped in so-called *super-states*, subsets of states connected together by low energy barriers, separated from outside regions by high energy barriers. Parsplice simulations tend to visit these super-states for long periods of time before transitioning to another super-state. These movements between super-states are referred to as *transition regions*, which typically contain the most important kinetic information of a system because it controls its long-term behavior. Transition regions are often comparatively short compared to the time spent trapped within super-states, while intra-super-state transitions occur very frequently.

When analyzing the structure of molecules, experts often investigate the neighbors of each atom and determine the shapes that these neighborhoods form in order to characterize a system. Mutations in the shape and crystalline structure of a system have a strong influence on its properties. There are seven main types of crystalline structures commonly found in materials, and our case studies are focused around analyzing cubic (face-centered cubic – FCC, body-centered cubic – BCC) and icosahedral (ICO) structures as they commonly occur in nano-particles; please refer to Misra [33] for a thorough discussion.

3.2 Analytical Tasks

We adopted an iterative design process to develop MolSieve with two domain experts who work in computational materials science; one of them has over twenty years of experience, and the other has more than six. We met bi-weekly for two years, using the feedback from these meetings to refine MolSieve's functionality and visual design. Through the design process, we identified a set of analytical tasks that are essential for gaining insight into long duration molecular dynamics

simulations. Simplifying these tasks became one of the core design objectives of MolSieve (Figure 2).

T1: Classify super-states and transition regions in individual trajectories. The first step in analyzing large simulations is to identify super-states and the transition regions that separate them, which are not known *a priori*. Transition regions are critical because they control how rapidly the system will experience significant changes that could affect its properties. This separation reduces the data-set to a manageable size and allows experts to concentrate their analysis on transition regions.

T2: Identify critical sub-regions, relevant patterns and motifs within transition regions. There are a number of patterns and motifs to be discovered within the transition regions of a trajectory. Patterns of state transitions often signify the presence of a structural change, but they can also be misleading due to the nature of long duration simulations, where repeated behavior is often due to the system making rapid low-energy transitions between states. The challenge lies in identifying patterns and sub-regions within transition regions where meaningful changes occur while ignoring low information density portions. The analysis of these sub-regions is the crux of molecular dynamics research; understanding how the structure of a material changes allows domain experts to make decisions on whether or not to use a certain material in an engineering application.

T3: Compare regions of interest between trajectories. MD trajectories are generated in a stochastic manner, so it is unlikely that two trajectories will contain the same behavior and physical structures. Therefore, there is a need to develop flexible methods that can differentiate robust features of the dynamics that are common to many simulations.

3.3 Requirements

After identifying the primary tasks found in MD analysis, we derived the following set of requirements for a visual analytics system.

R1: Guide the analyst to transition regions. Analysts should be guided to regions that are most likely to reveal significant changes in a system's structure.

R2: Automatic calculation of analyst-defined properties. The trajectory should be populated with automatically calculated properties that can be defined by an analyst. Time should only be spent computing properties for regions that are potentially interesting. The results should be stored in a data-base for future use.

R3: Highlight potentially interesting sub-regions. Once the expert-defined properties are rendered, the analyst should be guided towards sub-regions within transition regions that potentially express a change in the system's behavior. While **R2** focuses on calculating properties, guided visual exploration is another crucial aspect that accelerates the

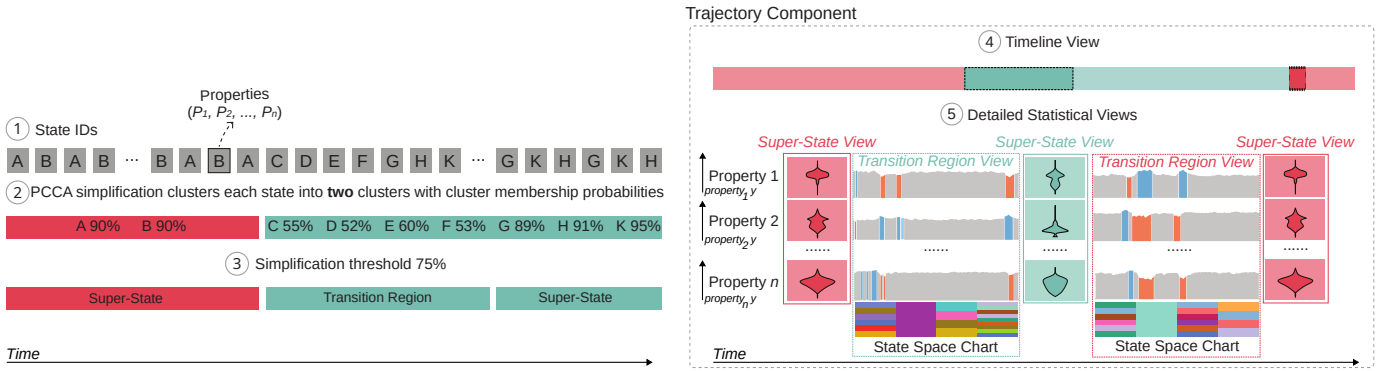


Fig. 3: The simplification scheme employed by MolSieve and its relation to the visual components in the system. (1) displays a portion of a sample trajectory's sequence, where each rectangle represents a state, the capital letter represents its state ID, and P_1 to P_n represent its analyst-defined properties. These properties are not required for the simplification and can be calculated and assigned afterward. (2) PCCA is performed on this sequence and it yields the maximum cluster membership probability for each state. Then, the simplification (3) is applied using an analyst-defined threshold (75% by default). States with a maximum cluster membership probability **above** this percentage are rendered as **super-states**, and states **below** are rendered as **transition regions**. These regions are mapped to views in a Trajectory Component which consists of the Timeline View and statistical views. The Timeline View (4) provides temporal context for the statistical views below. Regions drawn with dashed outlines are transition regions, while regions without outlines are super-states. The statistical views (5) are arranged temporally and split vertically into partitions, with each partition corresponding to a single property; we include axes to indicate the relative scale for each property. Super-State Views display small multiples of violin plots that outline the distributions of each property within a super-state. Transition Region Views are small multiples of control charts for each property that are accompanied by a State Space Chart. These charts collaborate to describe the most frequently occurring states within evenly divided segments; the number of states in a segment directly correlates to the number of unique states visited. Segments with large numbers of states usually indicate a structural change is occurring.

discovery process.

R4: Select, compare, and inspect regions of interest in detail. Integrating **R1–3** should enable the analyst to effectively select and refine regions of interest in a responsive manner, as well as allowing them to inspect a set of customized properties through expressive visualizations.

R5: On-demand calculation of detailed analyses. The selection process detailed in **R4** generates sub-regions that may include states that express behaviors of interest. Understanding their behavior requires physically grounded analyses which can be computationally expensive. The analyst should be able to request these analyses on demand and be able to continue exploring the trajectory.

R6: Extensibility. An intuitive extension of **R2** is the ability to define new properties. The solution should accommodate a broad spectrum of simulation types, enabling analysts to provide customized scripts for calculating system-specific properties. By providing this amount of flexibility, analysts can define properties which typically denote changes in a system. They can then use the visualizations and interactions provided by the solution to quickly identify regions of interest based on these properties.

R7: Ease of use and performance. The analyst should be able to easily navigate and discern patterns within trajectories. Additionally, the proposed solution must remain responsive during computationally intensive tasks and progressively render partially calculated data while waiting for results. The analyst should receive feedback regarding the progress of complex calculations as well as any errors that may occur, with the ability to adjust or cancel them as needed.

4 MOLSIEVE

MolSieve is a visual analytics system implemented using a FastAPI [1] back-end, and an interface powered by D3 [9], React [2], and Redux [3]. The back-end provides a powerful method for simplifying dense MD trajectories; its results are mapped to the views in the interface (Figure 1). The interface is designed to quickly guide analysts to potential regions of interest within MD trajectories (**T1**) and provides tools to interactively verify (**T2**) and compare (**T3**) multiple data-sets. Due to the tremendous amount of data that needs to be processed and stored on the fly, we designed our approach based on the progressive visual analytics paradigm [16].

To support a wide range of simulations, MolSieve automatically executes, stores, and renders the results of analyst-defined Python

scripts (**R2, R6**). This feature enables analysts to specify properties that indicate a region of interest for the simulation they are studying. These scripts are provided access to Atomic Simulation Environment (ASE) [28] representations of each state, which can be leveraged to calculate physically relevant properties of dynamic systems, e.g., the Common Neighbor Analysis (CNA) [19] counts for atomic structures. These n properties are calculated and assigned to each state within the trajectory (Figure 3.1). To further accelerate the process of discovery, these properties are calculated and rendered progressively, allowing analysts to gather insights throughout the data-set without having to wait for computations to finish (**R7**).

Background - Trajectory Simplification We used Perron Cluster Cluster Analysis (PCCA) [13] as implemented by pyGPCCA [39] as the basis for MolSieve's simplification scheme. PCCA has been proven to accurately simplify MD trajectories by clustering together groups of kinetically linked states [20, 21]. PCCA can be applied to simulations where transitions are modeled as a Markov chain.

MolSieve simplifies the trajectory by dividing it into tentative transition regions and super-states. This is achieved by first running PCCA on the trajectory, which divides it into N meta-stable sets of states, referred to as *clusters*. PCCA assigns a vector of N *cluster membership probabilities* to each individual state which describes how strongly it belongs to each cluster (Figure 3.2). Then, each individual state's membership probability is compared to a threshold set by the analyst (Figure 3.3); if its maximum membership probability is **above** the threshold, it is considered part of a *super-state*; otherwise, it is considered to be part of a possible *transition region* (i.e., it occurs in regions where the trajectory moves between clusters). If the simplification threshold is set to its maximum value of 1.0, no portion of the trajectory will be simplified, and every state will be considered a transition region.



When initially loading a trajectory, analysts have the opportunity to set a range for the PCCA clusterings they are interested in, as PCCA is not guaranteed to yield results for all numbers of clusters. The back-end uses the range to determine and return the optimal PCCA clustering for the trajectory and then simplifies it using the simplification threshold. Simultaneously, analyst-defined properties (P_1 to P_n) are calculated and assigned to each state within the trajectory.

The optimal clustering may not always reveal the best possible splits between transition regions and super-states, so analysts are free to adjust the PCCA cluster counts as well as the simplification threshold within the interface. The simplification threshold is set to a default

value of 0.75 and the PCCA clustering range to 2-20 which provides a reasonable starting point for exploration. A simplification threshold value of 0.75 tends to reveal sets of states that are weakly clustered, regardless of the PCCA cluster count. The default PCCA clustering range is set wide enough to ensure a clustering is found. Once a trajectory is simplified, its results are directly mapped to MolSieve's Trajectory Components (Figure 3 right).


4.1 Trajectory Components

Trajectory Components adopt a focus+context approach [17] to assist analysts in identifying regions of interest through the use of a variable number of Transition Region and Super-State Views (Figure 2.T1). Each trajectory belongs to a separate component, organized on the main area of the screen. The Timeline View (Figure 3.4) provides temporal context and a means of control for the statistical views (Figure 3.5).

Timeline View The Timeline View (Figure 3.4) displays the regions that are currently being rendered as statistical views (Figure 3.5) and allows experts to adjust which regions are visible to focus their analysis. Regions are colored according to the PCCA cluster they are assigned; transition regions are rendered with a dashed outline and super-states with no outline and a slightly lighter color in order to differentiate between them. We colored the clusters with a color scheme adapted from ColorBrewer's [18] qualitative set. Hovering over either type of statistical view highlights its corresponding region in the Timeline View. Brushing the view adjusts the visible extent of the trajectory, saturating regions that are outside of the brush's extent and reorganizing the statistical views. Double clicking the view zooms it in on the currently brushed region, which allows analysts to view regions that may have been rendered too small initially. There are two additional interactions provided by buttons next to each Timeline View. Clicking the  **PCCA parameters** button shows a menu containing two sliders that adjust the number of PCCA clusters and the simplification threshold (R1). The  **Reset timeline** button resets the Trajectory Component to show the entire trajectory.

The detailed statistical views (Figure 3.5) are arranged in temporal order from left to right and are drawn with a scale that exaggerates the size of transition regions to direct the analyst's attention (R1). Transition regions are exaggerated because they contain details on how the system evolved within a critical region, which demands more screen space; meanwhile, super-states are small multiples of violin plots which remain legible at small sizes. As a result, super-states are only offered a maximum of 10% of the total screen space unless there are no transition regions within the visible extents of the trajectory, in which case they occupy the entire width of the screen. This exaggerated scale was inspired by SignalLens [26].

Each view is bordered by the color of the cluster it is associated with, which assists in finding transition regions between clusters. The slices marked Property₁ to Property_n in Figure 3.5 display how views within Trajectory Components are vertically split into partitions. Each property is assigned a partition per trajectory, and each property is consistently rendered with its own scale in order to facilitate comparison. Partitioning the views in this manner allows experts to follow the evolution of multiple properties simultaneously.

We implemented a dynamic ranking system for each partition, which reduces the amount of irrelevant data on the screen in order to streamline intra-trajectory comparisons. Properties that change dynamically throughout a simulation are more likely to be relevant to analysts than properties that stay constant, so each property partition is ranked vertically based on the magnitude of its difference throughout the trajectory. This difference is calculated by performing a statistical z-test between each pair of adjacent super-states and summing them over the trajectory to get the final score. Each set of rankings is individual to a trajectory, as certain properties can be highly dynamic in one trajectory and stagnant in another. Since the data is being loaded progressively, the rank is calculated based on the information currently available to MolSieve. By default, only 4 properties are loaded, but there is a  **Property control** button in the main toolbar to adjust the number of properties shown.

Transition Region View To aid analysts in discovering sub-regions of



Fig. 4: The importance of being able to dynamically adjust the moving average period for the Transition Region View. (1) has a moving average period that is too short and contains a lot of false positives for anomalies. (2) has a moving average period that is appropriately chosen for the region being studied, showing only two major anomalous events within the given time-frame. (3) has a period that is too large to capture the interesting events occurring within the region.


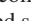
interest within a transition region, we designed the Transition Region View (Figure 3.5) to leverage control charts for detecting statistical anomalies [43]. Each view contains a small multiple of control charts which correspond to analyst-defined properties which work to highlight regions that have drastic changes in value. At the bottom of each view, a chart displays the state space within the region, which provide an overview of the most frequently occurring states.

Control Charts Each control chart displays the moving average of a property inside a transition region and it is colored based on the distance of the moving average from the mean. If the moving average moves one standard deviation **above** the mean, it is colored **blue**. If the moving average moves one standard deviation **below** the mean, it is colored **orange**, and if it stays **within** the control limits, it is colored **light gray**. This coloring scheme, inspired by ColorBrewer [18], draws attention to sequences within the transition region where a change is occurring, and allows analysts to quickly determine what sub-regions are of special interest, fulfilling R3 and R4. Hovering over a control chart displays a tooltip with the current value of the property and associated time-step.

By default, the moving average time period for each control chart is set to one tenth of the length of its Transition Region. However, if the analyst finds that the moving average time period is not capturing regions of interest, they can adjust the moving average time period for all of the control charts within the view (Figure 4).

State Space Chart A state ID vs time-step pixel plot is a familiar way of visualizing ParSplice trajectories [20]. Each time-step within a ParSplice simulation corresponds to one state. Rendering states this way allows analysts to quickly determine regions of interest. Through our iterative design process, we found that rendering each state within a transition region would lead to highly cluttered and cumbersome graphs, since transition regions often consist of thousands of states (Figure 5 top). In order to address these issues, we devised and implemented an aggregate version of this plot, the State Space Chart (Figure 5 bottom).

The State Space Chart highlights changes within the transition region by splitting it into ten evenly divided segments and calculating which states occur most frequently within each segment. To be considered part of a segment, a state's actual distribution value needs to be greater than its expected value. This is defined as 1 divided by the number of unique states within the segment; i.e., a state occurs equally likely as all of its neighbors. Segments with many colors indicate that the simulation is rapidly moving between many unique atomic configurations. Using the control charts coupled with this view allows analysts to quickly estimate the visit frequency within the region and identify sub-sequences worthy of a detailed inspection (T1).

Transition Region View Interactions Upon hovering over a Transition Region View, a toolbar with multiple controls is displayed. The  **Select sub-region** button toggles a brush to select sub-regions of interest; completing the selection generates a corresponding Sub-Sequence component (Section 4.2; R4; Figure 6a). To facilitate making fine-grained selections within a region, the  **Zoom into region** button enlarges the transition region so it occupies the entirety of the screen.

Super-State View Super-states revealed by the simplification algorithm tend to constitute the majority of ParSplice simulations. To maximize performance, we elected to use aggregate statistical charts [51] when representing super-states. A Super-State View (Fig-

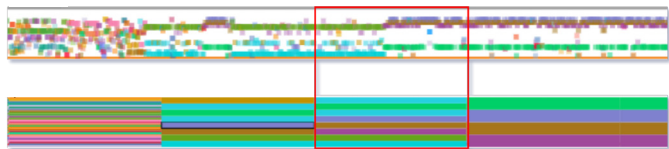


Fig. 5: The original design (top) for state ID vs. time-step charts was flawed, as it attempted to render a large amount of data in a small space. States would often occlude each other, and each render would be computationally costly, as each rectangle is an SVG element. Our final design (bottom) underlines the behavior of a sub-region succinctly and makes it easier to read and interpret when many states are in a region. Note the highlighted segment, which captures the transition between two minor repetitive sub-regions.

ures 1.3a, 1.3b, Figure 3.5) is a small multiple of violin plots that describes the overall distribution of each property. They highlight the evolution of each property throughout a simulation in a compact manner. We originally used box-plots to display the distributions of these properties, but we found that they were highly cluttered due to the small amount of screen space allotted to them and they did not capture the variance of each distribution as well as violin plots.

Each violin plot is constructed using the property values from the so-called **dominant** states of the region plus a randomly selected 1%. Dominant states within the super-state are states that occur with larger than median frequency. Using a small randomly sampled portion of the region provides a reasonable overview without having to compute a prohibitive amount of data. In order to ensure that important details about states are not hidden from analysts, we implemented an expansion feature that allows experts to explore super-states in more detail. Double clicking any of the control charts causes the Transition Region View to “expand” (Figure 6b), revealing its state space and the moving averages of its neighbors. Expansion occurs 100 time-steps at a time to avoid loading unnecessary data.

Interactions The toolbar above all Trajectory Components provides interactions that enhance the analyst’s ability to examine a trajectory in detail (Figure 1, top left). Analysts are able to construct multi-variate control charts [30] with the properties they provided by clicking the **Add multi-variate chart** button which opens a modal window (Figure 6c). These multi-variate charts (Figures 1.4a, 1.4b) are dynamically added to each Transition Region View, allowing the analyst to combine various properties to generate more powerful control charts that highlight synchronized movements across property values that are difficult to detect using single variable charts, fulfilling **R5**. The **Swap trajectory** button allows analysts to swap the vertical positions of Trajectory Components to facilitate direct comparisons. The **Clear selection** button allows experts to undo a selection they are currently making if they decide they want to abort the process.

4.2 Sub-Sequence Component

Since the State Space charts within Transition Region Views only provide an overview, there is a need to look at sub-regions in more detail. Sub-Sequence Components (Figure 2.T2) are added to the bottom of the screen once an analyst completes a selection in a Transition Region View using the **Select sub-region** button (Figure 6a). They are designed to fulfill **R4** and **R5**, as they allow experts to glean additional insight from regions that they deem to be interesting, and correspond to the abstract/elaborate interaction category in Yi et al. [53]. Each Sub-Sequence Component provides a small multiple of 3D state visualizations, which serves as an overview of the structural changes occurring within the selection. To generate the overview, we developed a greedy search algorithm that uses the *Frobenius norm* (provided by ASE [28]) of the spatial distance between all atomic coordinates. A high distance between states indicates that they are structurally different. The algorithm iterates over the selection and takes the distance between the state being queried and the rest. To find states that are as different as possible, we start at the initial state of the selection, find its most dissimilar counterpart, and start the search again at this state until we reach a maximum iteration count or the end of the selection.

At the bottom of each Sub-Sequence Component is a traditional state ID vs. time-step plot of the selection’s constituent states (Figure 5 top). The Sub-Sequence Component also supports running the Nudged Elastic Band calculation [22] using the **Run NEB** button. Clicking the button (Figure 2.T2) opens a modal window that allows analysts to adjust the parameters of the calculation and make a selection on the sub-sequence that will be used in the calculation. The results from the calculation are used to generate a potential energy graph which shows the minimum energy pathways for the selection they made, fulfilling **R5**. Potential energy graphs are commonly used by analysts to determine if a sequence of states constitutes a structural change in the simulation. An exceedingly high potential energy barrier between any two pairs of states in the sequence followed by any number of low energy barriers, usually indicates a transition. This is because particles are known to move towards their lowest energy configurations.

State Detail Widget Whenever a state is clicked throughout the UI (e.g., within State Space Charts, Sub-Sequence Components etc.), the State Detail Widget (Figure 2.T2) is updated. It displays a static 3D visualization of the state, inspired by guidelines outlined in Byška et al. [11] that suggest linking 3D visualizations of a system to its properties. Additionally, a table is shown below the 3D render displaying the properties of the state that was selected. Since states are all colored consistently throughout the visual interface, we included a bar under all 3D renders that displays the selected state’s color, making it easy to visually link the state to other visualizations. The **Modify 3D render** button in the trajectory toolbar allows analysts to change the way states are rendered in 3D throughout the interface by Python scripts inside the `vis_scripts` folder in the source code (**R6**). Experts pick the visualization script they want to use with a pop-up menu that is populated with the contents of the `vis_scripts` folder. Analysts are expected to define a function that takes an OVITO [45] rendering pipeline object as a parameter which they can modify to suit their needs. Figure 2.T2 demonstrates an example: the default view is swapped for a visualization of crystalline structure neighborhoods where each atom is colored according to its structural classification (see Section 3.1). Customizing the visualization gives analysts an additional method to verify their conclusions made from the 2D charts in MolSieve and is integral to certain types of analyses (Section 5.2).

4.3 Comparison Widgets and Interactions

MD ensembles are practically impossible to analyze due to the amount of data that needs to be compared. To address this, we included a variety of comparison interactions that quantify the difference between regions of interest from multiple trajectories (Figure 2.T3 and Figure 6).

The **Compare regions/selections** button allows experts to select regions or sub-sequences they want to compare directly. When two are selected, a Region Comparison Widget is placed at the bottom of the screen which contains asymmetrical violin plots that compare the distributions of each property (Figure 6d). Transition Region Views can also be selected with this interaction, making it easy to compare them with Super-State Views; MolSieve uses the properties from the dominant states in transition regions to compute the distribution of each property, allowing for a fair comparison. Comparing regions this way reduces the cognitive load of having to look back and forth between two distributions that are visually separated. When two Sub-Sequence Components are selected, a Sub-Sequence Comparison Widget is generated, which displays a state similarity heat-map (Figure 6e). State similarity is defined as the inverse of the distance used in the 3D overview for Sub-Sequence Components, see Section 4.2.

The **Find similar regions** button lets an expert select a Transition Region View to quickly compare to all other Transition Region Views that are currently selected using the Timeline View’s brush, which corresponds to a Connect interaction in Yi et al. [53]. Once the selection is complete, MolSieve computes the difference between their state distributions and then displays the result with a tooltip rendered above each region (Figure 6f). This computation provides a crude preview of similarities between two transition regions, which can be used to narrow down which regions require an in-depth comparison.

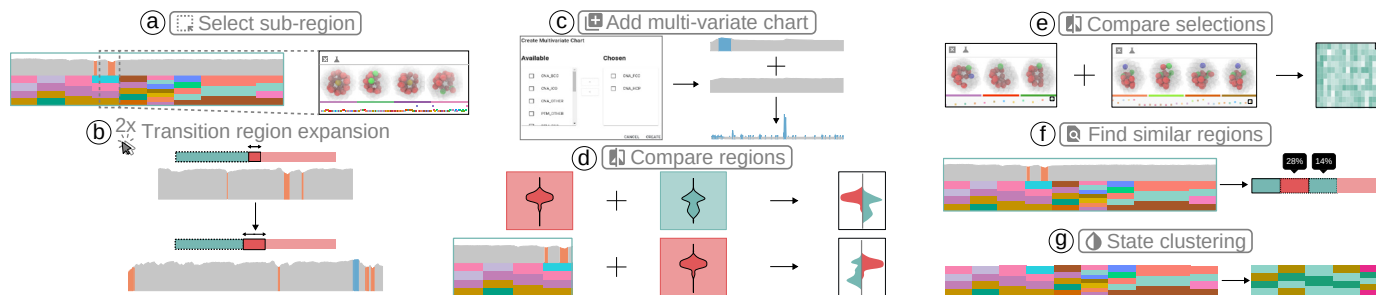


Fig. 6: MolSieve's unique interactions. (a) lets analysts select sub-regions of interest within a Transition Region View to create a Sub-Sequence Component. (b) Double clicking a Transition Region View causes it to expand into its neighbors, which makes it possible to view parts of super-states in detail. (c) allows experts to create multi-variate charts. The **Compare regions/selections** button can compare any Transition Region or Super-State View by creating a Region Comparison Widget (d) which is a small multiple of asymmetrical violin plots detailing the distributions from each selected region. The button also works with Sub-Sequence Components, creating a Sub-Sequence Comparison Widget (e) that contains a heat-map detailing the similarities between the selections. (f) allows experts to select a single Transition Region View, which MolSieve uses to compare with all other visible transition regions, automatically highlighting their similarity on the Timeline View. (g) recolors all of the states in the interface according to the OPTICS clustering algorithm using analyst-defined properties.

Clicking the **State clustering** button clusters all of the states in the transition regions visible on the screen based on their properties (Figure 6g). MolSieve uses the OPTICS clustering algorithm [6] to generate clusters to color the states by. Clustering states together based on their properties provides a slow, but flexible method to directly compare trajectories.

These interactions were designed to replace one of the inherent visual features of chord diagrams in an earlier design, where regions could be rendered as arcs on a circle and linked together based on similarity. When we attempted to implement chord diagrams in MolSieve, we found them to be cluttered and confusing when trying to interpret the temporal structure of the data.

Clicking anywhere in the Super-State View updates the State Detail Widget with the state that occurs most frequently within the region, referred to as the region's *characteristic state*. Characteristic states describe the general properties of these regions [42]. Thus, this interaction allows analysts to quickly determine if any structural changes occurred between super-states. Once a change has been identified, analyst can seek more detailed information about the change within the Transition Region View between the two differing super-states.

5 CASE STUDIES

We demonstrate the efficacy of MolSieve by presenting two case studies in which we conducted pairwise analysis [7] with our domain experts E1 and E2. The first case study involves analyzing two long-duration trajectories of platinum nano-particles, first by determining sub-regions in each trajectory where the particle undergoes a structural change and then comparing them. The second case study focuses on atom vacancy analysis, where a reference atomic configuration is compared to states within the trajectory. In atom vacancy analyses, experts typically look for regions within a simulation where the “missing” atoms begin to displace in tandem.

5.1 Platinum Nano-particles

Our nano-particle expert (E1) aimed to identify and characterize significant fluctuations in the shape of a platinum nano-particle subjected to high temperatures. E1 began the case study by loading a simulation of a platinum nano-particle at 750 kelvins, which consists of approximately eighteen million transitions and twenty-five thousand unique states (Table 1), with each state representing different configurations of a nano-particle with 147 platinum atoms.

Based on a prior study of nano-particles [21], the analyst decided that the best properties to analyze this simulation were the Common Neighbor Analysis (CNA) [19], Ackland-Jones [4] (AJ), and Polyhedral Template Matching [29] (PTM) atom characterization counts. These analyses attempt to characterize the structure of a nano-particle based on descriptors of the local environment around each component atom and have been found to be strong indicators of transition regions. The

analyst wrote a script that used OVITO [45] to compute these properties and loaded them into MolSieve (R2). Since it was difficult to tell what was occurring to the nano-particle from the default 3D render, our analyst wrote a visualization script that highlights CNA counts within states (Figure 1.6a, 1.6b, 1.7, 1.11a, 1.11b, 1.14). The CNA visualization script renders HCP atoms as red, ICO atoms as yellow, and FCC atoms as green.

Identify Transition Regions (T1): E1 decided to load the trajectory with a PCCA clustering range of 2-20 and a simplification threshold of 0.75. PCCA split the trajectory into two clusters, yielding a small red cluster in between a dominant teal cluster (Figure 1.1) which E1 zoomed in on using the Timeline View. This revealed a busy region with many possible transitions; however, the Super-State Views showed that the super-state distributions did not vary greatly between each other, so the analyst increased the number of clusters to 4, hoping to reveal more fine-grained super-states (Figure 1.2). Once the simplification was rendered, they found that there were a number of transition regions between super states where the ICO and HCP counts of the nano-particle were rising (R1; Figure 1.3a and 1.3b).

Analyze Transition Patterns (T2): The analyst added a multi-variate control chart using the ICO counts from all three analyses to see if they would all point towards the same regions (Figure 1.4a and 1.4b). The analyst then found two sub-regions within a Transition Region View where the control charts indicated that the structure of the nano-particle changed (Figures 1.5a, 1.5b; R3).

Next, E1 clicked on the Super-state Views (Figure 1.3a, 1.3b) surrounding that Transition Region View to get an understanding of how the nano-particle changed from the first super-state to the second; the characteristic states of each super-state are shown in Figures 1.6a and 1.6b. Since it was difficult to tell what was occurring to the nano-particle from the default 3D render, the analyst changed the 3D view to highlight CNA counts. This revealed a sudden change in the ICO count, where the two green atoms disappear in Figure 1.3a disappear.



To verify that the sudden change in ICO count was not a random event, they double-clicked the Transition Region View to expand it. This confirmed that the nano-particle stays in the same configuration for some time before suddenly undergoing a drastic change in the FCC and ICO counts. Satisfied, they made a selection in the region where the ICO count suddenly changed from zero to one (Figure 1.5a), which rendered a Sub-Sequence component. Then, they clicked through the states in the Sub-Sequence Component to get a detailed look at what was occurring to the particle (R4). This revealed that the trajectory was undergoing a transformation (Figure 1.7) within the region the analyst selected (Figure 1.5a); the nano-particle started the transition with two FCC atoms (Figure 1.6a) and lost them (Figure 1.6b). They ignored the other sub-region where the ICO count dropped (Figure 1.5b), stating that “This is normal behavior in simulations with a heterogeneous energy barrier: the system tries to escape its configuration but is not able


to, causing it to change before returning to its previous configuration; this is why I wanted to check the region on the left."

Once the transition was found, they decided to run Nudged Elastic Band (NEB) calculations on the both ends of the suspected transition region (**R5**). The NEBs confirmed that the transition to and from the suspected transition region took a large amount of energy, thus demonstrating that our system is a significant improvement in terms of detecting regions of interest in large molecular dynamics simulations.

Ensemble Analysis (T3): Once the transition was confirmed, the analyst decided to load another platinum nano-particle trajectory at 800K. E1 aimed to determine if the structural changes they observed in the particle at 750K were similar to the ones observed at 800. The 800K simulation contains thirteen million transitions and fifty-three thousand unique states (Table 1).

Once the simulation was loaded, the analyst used a similar workflow to determine where the transition regions occurred in the trajectory by carefully adjusting the simplification threshold until a suitable number of possible transition regions were displayed. Starting at the simplification threshold's default value of 0.75 did not yield any transition regions; however, it led the analyst to zoom into a sequence of super-states where the ICO count was changing from zero to one. Increasing the simplification threshold to 0.85 revealed super-states undergoing a transition similar to the 750K trajectory (Figure 1.8). The state IDs overlapped between the two trajectories, and many regions that contained the same 4 states that the 750K simulation spent large amounts of time in (Figure 1.12a and 1.12b; **R1, R3**).

The analyst then decided to click the  **Find similar regions** button and select the transition region they discovered in the 750K trajectory. This revealed many regions shared a large portion of states with the selection, a region which scored 12% similarity based on the set of unique states present in each region, which can be seen on the Timeline View (Figure 1.9) (**R4**). A similarity of this magnitude is significant due to the fact that simulations are unlikely to contain the same states in a small temporal region. They then used the  **Compare regions/selections** button to examine the difference in distributions between the regions that scored highest on similarity and the original transition region they discovered (Figure 1.10). While the region that was 12% similar did not have the same transition characteristics, the analyst found a region that had a similar shift in its ICO and HCP count. Moreover, when the analyst clicked on the two Super-State Views surrounding the region, they found that the first super-state had the same characteristic state as the first in the previously found region, and the second super-state was a rotation of the previous trailing super-state (Figure 1.12a, 1.12b).

While the nature of the transition was similar based on the control charts, E1 wondered if the states were truly structurally similar, so they went to use the  **State clustering** button. Recoloring the states based on their structural cluster revealed that the region shared many states, particularly around the sub-regions where the analyst believed a transition was occurring (Figures 1.12a and 1.12b). The analyst also used the Sub-Sequence Comparison Widget to compare the two selections (Figure 1.13), which verified that the transitions were similar in nature as the states were rotational analogs of each other. The combination of these comparisons reassured the analyst in their conclusion that these transitions were of a similar nature (**R4**). While not identical to the one found in the 750K simulation, the sub-region found by the analyst also describes how the nano-particle loses FCC atoms and gains an ICO atom (Figure 1.14); this slight difference is to be expected due to the fact that MD simulations are stochastic by nature. This discovery demonstrates that our system is effective in not only detecting regions of interest in one long-duration simulation but is also capable of detecting similar physical occurrences in multiple simulations.

5.2 Bulk Tungsten Defect Analysis

The goal of a defect analysis is to understand the way the point defects in a crystalline structure evolve over the course of a simulation; these defects determine the properties of a given material. Typically, analysts use the Wigner-Seitz cell method [54] to visualize the difference between a state in a defect simulation and a reference structure that does not have any defects. Our analyst, who specializes in cell defects (E2),

provided a reference Tungsten lattice with 2,000 atoms, which represented a perfect, defect-free crystalline structure, as well as a Python script from his daily workflow that compares a state and the reference structure using the Wigner-Seitz analysis. The script they provided outputs the defective atoms in each state and displays them, which the analyst used as the state view for the case study, seen in the renders for Figures 7.A and 7.B. Additionally, the analyst used the output from the script to create three properties which described the center of mass of the atoms that were defective (**R6**).

To begin the case study, the analyst loaded their scripts and a simulation of a Tungsten crystalline lattice being subject to various deformations at 1000 kelvin. This data-set was considerably smaller than the nano-particle case study, having only approximately 800 transitions and only 50 unique states (Table 1). However, the size of each state was considerably larger, as each state represented a Tungsten lattice with 1996 atoms.

Analyze Transition Patterns (T2): MolSieve initially classified the entire trajectory as 3 super-states, which meant that the PCCA simplification was not useful for this data-set. This prompted E2 to set the simplification threshold to 1.0, and rendering all of the PCCA clusters as transition regions, allowing the analyst to see the control charts for each property. Once it was re-rendered, the analyst noticed that the moving average time period for each Transition Region View was very high, obscuring potentially interesting transitions, so they set the moving average time period for each transition region to 10. Once the system was configured properly, the control charts exposed regions where the center of mass changed rapidly in all three dimensions. MolSieve immediately identified diffusive transitions (Figure 7.A), highlighting them among the numerous repetitive thermal vibrational motions (Figure 7.B) that were composed of single vacancies moving back and forth (**R1, R3**). E2 then selected several regions highlighted by the control charts and was able to identify and follow the chain of events for several diffuse transitions. This case study was able to demonstrate that MolSieve is effective in finding regions of interest in diverse analysis scenarios.

5.3 Domain Expert Feedback

To evaluate MolSieve, we conducted an hour-long semi-structured interview session with E1 and E2. During the interview, we asked them to compare their daily workflow to using MolSieve and solicited their suggestions on improving the system.

A typical workflow for a molecular dynamics analyst consists of running scripts for several days on simulation data and sifting through states manually. They typically visualize the states in OVITO [45] and then click frame-by-frame to get an idea of what changes the system is going through. The greatest challenge in analyzing simulations this way is the amount of data that needs to be processed which makes it difficult to keep track of transitions and one's temporal context within the trajectory. E1, our nano-particle expert, remarked that "The overall layout of MolSieve makes it easy to analyze these data-sets. It is very easy to understand where you are in the trajectory, just by looking at the Timeline View. This helps me think about what is going on in the simulation as a whole, and I don't feel like I have tunnel vision while examining data." They continued their reflection on the system by comparing the experience of examining regions of interest in MolSieve with their daily workflow, specifically praising the 3D overview within Sub-Sequence Components, saying that "The 3D overview [within the Sub-Sequence Component] provides a very nice, pictorial, visual effect that gives a preview of what the particle is going through. I don't have to waste time clicking back and forth between states to get an idea of what I'm looking at." E2, our defect analysis expert, reflected on MolSieve's visual design by saying, "The combination of the control charts and the aggregate state space chart make it easy to find regions of interest within a transition region. The aggregate state chart also tells me which regions to avoid selecting, since it's so easy to see where the simulation gets stuck jumping between a small set of states."

The experts found that MolSieve was efficient, providing a massive productivity increase over their accustomed workflows. E1 said, "The system is exciting, as it takes an unimaginable amount of data and makes it interpretable. The nano-particle simulations we examined

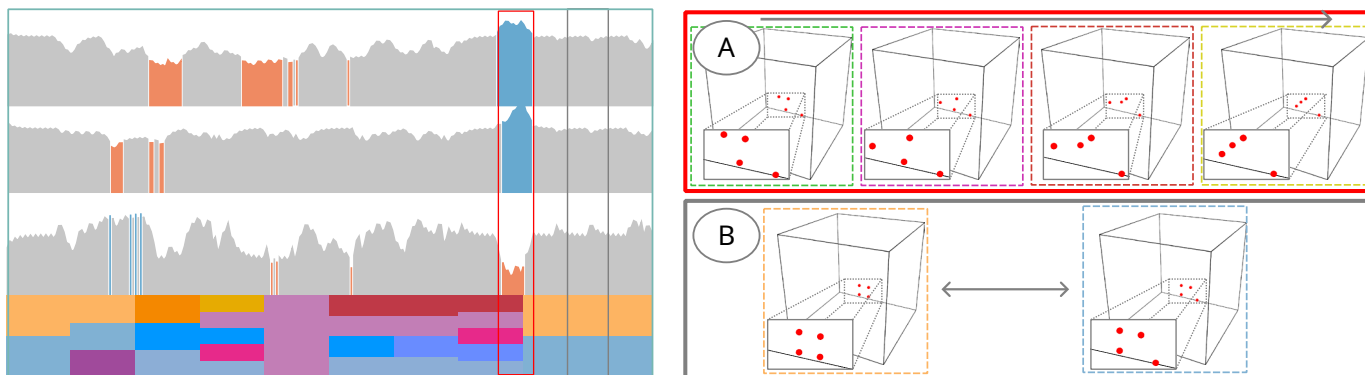


Fig. 7: Results of the defect analysis case study. (A) Demonstrates an example of the diffusive transitions discovered within the simulation, which are a set of unique structural changes occurring to the defective region within the Tungsten crystalline lattice. (B) Demonstrates an example of "fluttering", where the defects within the lattice move back and forth between two configurations, one atom at a time. These kinds of transitions are the predominant transformations occurring to the lattice throughout the trajectory. The dashed rectangles represent the colors of each state's ID, and demonstrate that the State Space Charts were effective in capturing the regions of interest.

Simulation	Generation time (s)	Simulation time (ns)	# Timesteps	# States	Time to load (cached) (s)	Total preprocessing time (s)
nano-pt-700	35,994	62,857.99	6,711,821	16,631	5.876	459.156
nano-pt-750	35,992	49,869.45	18,463,872	24,457	7.714	1507.284
nano-pt-800	35,993	43,152.76	13,348,978	53,018	10.489	1150.753
nano-pt-900	57,586	31,636.00	7,721,529	490,226	14.979	9172.017
tungsten	10,800	10,000.00	866	241	2.000	6.800

Table 1: Several simulations that were tested on MolSieve are presented here. This table displays the total time it took to generate each simulation in ParSplice, the length of time the simulation represents in nanoseconds, the number of discrete timesteps in the simulation, the number of unique states, the time it takes to load the simulation when cached, and how long it takes to load each simulation.

could take several lifetimes to sift through, and MolSieve manages to make it look trivial, with near real-time performance." E2 added, "The amount of data I was able to comb through with MolSieve would have normally taken a few weeks to do, and I managed to do this in just a few minutes," which indicates that we fulfilled R7.

The customizability of the system was a major selling point, as E2 stated, "That is what really makes it come to life - this makes it applicable for a wide array of applications and will save us a considerable amount of time in the future." E2 continued the discussion by suggesting that analysts should be able to customize the simplification algorithm. This idea stems from the results of the atom vacancy case study (Section 5.2), where the simplification algorithm failed to produce transition regions. To get around this, E2 increased the simplification threshold to include the entire trajectory. E2 warned that, "In principle, the simplification scheme in MolSieve should work on most data-sets but molecular dynamics simulations are often analyzed in various modalities, some of which are not captured by dividing the trajectory using PCCA." Thus, allowing experts to customize how the trajectory is simplified could make it easier to find relevant regions for various analyses. E2 added that the distance metric used in both the overview in the Sub-Sequence Component (Section 4.2) as well as the heatmap in the Sub-Sequence Comparison Widget did not effectively describe the difference between two states. This was due to the fact that we were studying the **absence** of atoms within a state. To make these comparisons more useful, they suggested that the distance functions in MolSieve should also be customizable.

Finally, E1 felt that the MolSieve was lacking a feature for comparing multiple individual states. We focused on comparing sub-regions within trajectories and did not consider the importance of being able to easily compare two or more states. The Sub-Sequence Comparison Widget supports this to a limited extent, but E1 suggested an interaction that could "save" states and show them on demand.

6 CONCLUSION

In this work, we present MolSieve, a visual analytics system for long-duration molecular dynamics simulations modeled by discrete Markov

chains. Through the use of multiple coordinated visualizations powered by a data simplification scheme unique to MD simulations, MolSieve makes it possible to analyze previously unexplored simulation data-sets. The comparison interactions offered by the system provide support for analyzing simulation ensembles. Additionally, MolSieve's Python programming interface lets it accommodate a wide variety of simulations. To demonstrate the effectiveness of MolSieve's design, we analyzed three simulations alongside our domain experts: two nano-particle simulations and one atom vacancy simulation. Table 1 provides a detailed look at the efficiency of the system.

However, it became apparent that some of its components need to support further customization. We found that the simplification algorithm would sometimes return many regions in a trajectory, which led to the screen being highly cluttered. This would require the analyst to zoom in using the Timeline View to get a better idea of the general trend within the trajectory. This can be mitigated by reformulating the way regions are rendered to only show large regions until the zoom level is appropriate. We also found that some visual design elements must be adjusted; these issues are particularly prevalent in the color encodings of the interface. The analysts found that coloring states by their IDs made it difficult to distinguish them from one another once a large number of states were rendered on the screen, which we attempted to remedy by implementing the state clustering function. However, the state clustering function would sometimes also have color overlap, which could be reduced by mapping the number of clusters to a set of salient colors. Alternatively, we could explore using different visual encodings to distinguish a large number of classes. Another limitation is the inability to view a list of the most frequently occurring states within a Super-State. This can be addressed by adding a widget that shows all of the most frequently occurring states in a region.

In the future, we plan to address some of the limitations of the system, including the cramped visual encoding space and the need for extra customization. Providing additional support for exploring biological simulations would be of particular interest, as this could lead to a truly general MD region-of-interest visual analytics system. To continue scaling, we plan to switch the rendering engine to use WebGL instead of SVG, allowing MolSieve to take advantage of the current innovations in consumer graphics technology. Moreover, a number of techniques have yet to be integrated into our system - improving the 3D rendering pipeline will allow MolSieve to support a number of novel analyses (e.g., [25, 47, 52] and rendering techniques [40]). Future work will also include a method to recall expert selections, a direct state comparison view, and better 3D rendering support.

7 SUPPLEMENTARY MATERIALS

We included a demo video that showcases the first case study and an instruction manual for MolSieve as supplementary material. MolSieve's source code is available at <https://github.com/rostyhn/MolSieve>.

ACKNOWLEDGMENTS

This project was supported by the Exascale Computing Project (17-SC-20-SC), a collaborative effort of the US Department of Energy Office of Science and the National Nuclear Security Administration and the U.S. Department of Homeland Security under Grant Award Number 17STQAC00001-06-04. Los Alamos National Laboratory is operated by Triad National Security LLC, for the National Nuclear Security Administration of the U.S. DOE under Contract No. 89233218CNA0000001. We graciously acknowledge computing resources from the Los Alamos National Laboratory. The views and conclusions contained in this document are those of the authors and should not be interpreted as necessarily representing the official policies, either expressed or implied, of the U.S. Department of Homeland Security. We also would like to acknowledge Jiayi Hong and Andrew Garmon for their discussions and contributions to the paper.

REFERENCES

- [1] FastAPI. <https://fastapi.tiangolo.com>, 2019. Accessed: 2023-03-22. 4
- [2] React. <https://reactjs.org>, 2019. Accessed: 2023-03-22. 4
- [3] Redux. <https://redux.js.org>, 2019. Accessed: 2023-03-22. 4
- [4] G. J. Ackland and A. P. Jones. Applications of local crystal structure measures in experiment and simulation. *Physical Review B*, 73(5):054104:1–054104:7, Feb. 2006. doi: 10.1103/PhysRevB.73.054104 7
- [5] W. Aigner, S. Miksch, W. Müller, H. Schumann, and C. Tominski. Visualizing time-oriented data—A systematic view. *Computers & Graphics*, 31(3):401–409, June 2007. doi: 10.1016/j.cag.2007.01.030 2
- [6] M. Ankerst, M. M. Breunig, H.-P. Kriegel, and J. Sander. OPTICS: Ordering points to identify the clustering structure. In *Proceedings of the ACM SIGMOD International Conference on Management of Data*, p. 49–60. ACM, New York, June 1999. doi: 10.1145/304182.304187 7
- [7] R. Arias-Hernandez, L. T. Kaastra, T. M. Green, and B. Fisher. Pair analytics: Capturing reasoning processes in collaborative visual analytics. In *Proceedings of the Hawaii International Conference on System Sciences*, p. 10 pages. IEEE, New York, Jan. 2011. doi: 10.1109/HICSS.2011.339 7
- [8] H. Bekker, H. J. C. Berendsen, E. J. Dijkstra, S. Achterop, R. van Drunen, D. van der Spoel, A. Sijbers, H. Keegstra, and M. K. R. Renardus. GRO-MACS: A parallel computer for molecular dynamics simulations. In R. A. de Groot and J. Nadrchal, eds., *Proceedings of the International Conference on Computational Physics*, pp. 252–256. World Scientific Publishing, Singapore, Apr. 1993. 1
- [9] M. Bostock. D3. <https://d3js.org/>, 2011. Accessed: 2023-03-22. 4
- [10] M. Brehm, M. Thomas, S. Gehrke, and B. Kirchner. TRAVIS—A free analyzer for trajectories from molecular simulation. *The Journal of Chemical Physics*, 152(16):164105:1–164105:20, Apr. 2020. doi: 10.1063/5.0005078 2
- [11] J. Byška, T. Trautner, S. M. Marques, J. Damborský, B. Kozlíková, and M. Waldner. Analysis of long molecular dynamics simulations using interactive focus+context visualization. *Computer Graphics Forum*, 38(3):441–453, July 2019. doi: 10.1111/cgf.13701 2, 6
- [12] J. Chae, D. Bhowmik, H. Ma, A. Ramanathan, and C. Steed. Visual analytics for deep embeddings of large scale molecular dynamics simulations. In *Proceedings of the IEEE International Conference on Big Data*, pp. 1759–1764. IEEE, New York, Dec. 2019. doi: 10.1109/BigData47090.2019.9006048 2
- [13] P. Deufhard and M. Weber. Robust Perron cluster analysis in conformation dynamics. *Linear Algebra and its Applications*, 398:161–184, Mar. 2005. doi: 10.1016/j.laa.2004.10.026 4
- [14] M. Dreher, J. Prevotau-Jonquet, M. Trellet, M. Piuze, M. Baaden, B. Raffin, N. Ferey, S. Robert, and S. Limet. ExaViz: A flexible framework to analyse, steer and interact with molecular dynamics simulations. *Faraday Discussions*, 169:119–142, Feb. 2014. doi: 10.1039/C3FD00142C 2
- [15] D. Duran, P. Hermosilla, T. Ropinski, B. Kozlíková, Vinacua, and P.-P. Vázquez. Visualization of large molecular trajectories. *IEEE Transactions on Visualization and Computer Graphics*, 25(1):987–996, Jan. 2019. doi: 10.1109/TVCG.2018.2864851 2
- [16] J.-D. Fekete, D. Fisher, A. Nandi, and M. Sedlmair. Progressive data analysis and visualization. *Dagstuhl Reports*, 8(10):40 pages, Apr. 2019. doi: 10.4230/DagRep.8.10.1 4
- [17] G. W. Furnas. A fisheye follow-up: Further reflections on focus + context. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, pp. 999–1008. ACM, New York, Apr. 2006. doi: 10.1145/1124772.1124921 5
- [18] M. Harrower and C. A. Brewer. Colorbrewer.org: An online tool for selecting colour schemes for maps. *The Cartographic Journal*, 40(1):27–37, June 2003. doi: 10.1179/000870403235002042 5
- [19] J. D. Honeycutt and H. C. Andersen. Molecular dynamics study of melting and freezing of small Lennard-Jones clusters. *Journal of Physical Chemistry*, 91(19):4950–4963, Sept. 1987. doi: DOI: 10.1021/j100303a014 4, 7
- [20] R. Huang, L.-T. Lo, Y. Wen, A. F. Voter, and D. Perez. Cluster analysis of accelerated molecular dynamics simulations: A case study of the decahedron to icosahedron transition in Pt nanoparticles. *The Journal of Chemical Physics*, 147(15):152717:1–152717:6, Oct. 2017. doi: 10.1063/1.4996922 4, 5
- [21] R. Huang, Y. Wen, A. F. Voter, and D. Perez. Direct observations of shape fluctuation in long-time atomistic simulations of metallic nanoclusters. *Physical Review Materials*, 2(12):126002:1–126002:9, Dec. 2018. doi: 10.1103/PhysRevMaterials.2.126002 4, 7
- [22] H. Jónsson, G. Mills, and K. W. Jacobsen. *Nudged Elastic Band Method for Finding Minimum Energy Paths of Transitions*, pp. 385–404. World Scientific, Singapore, June 1998. doi: 10.1142/9789812839664_0016 6
- [23] A. Jurčík, D. Bednar, J. Byška, S. M. Marques, K. Furmanová, L. Daniel, P. Kokkonen, J. Brezovský, O. Strnad, J. Štourač, A. Pavelka, M. Maňák, J. Damborský, and B. Kozlíková. CAVER analyst 2.0: Analysis and visualization of channels and tunnels in protein structures and molecular dynamics trajectories. *Bioinformatics*, 34(20):3586–3588, Oct. 2018. doi: 10.1093/bioinformatics/bty386 2
- [24] M. Karplus and G. A. Petsko. Molecular dynamics simulations in biology. *Nature*, 347:631–639, Oct. 1990. doi: 10.1038/347631a0 1
- [25] M. Kern, S. Jaeger-Honz, F. Schreiber, and B. Sommer. APL@voro—interactive visualization and analysis of cell membrane simulations. *Bioinformatics*, 39(2):1–3, Feb. 2023. doi: 10.1093/bioinformatics/btad083 9
- [26] R. Kincaid. SignalLens: Focus+context applied to electronic time series. *IEEE Transactions on Visualization and Computer Graphics*, 16(6):900–907, Nov. 2010. doi: 10.1109/TVCG.2010.193 2, 5
- [27] T. D. Kühne, M. Iannuzzi, M. D. Ben, V. V. Rybkin, P. Seewald, F. Stein, T. Laino, R. Z. Khaliullin, O. Schütt, F. Schiffrmann, D. Golze, J. Wilhelm, S. Chulkov, M. H. Bani-Hashemian, V. Weber, U. Borštnik, M. Taillefumier, A. S. Jakobovits, A. Lazzaro, H. Pabst, T. Müller, R. Schade, M. Guidon, S. Andermatt, N. Holmberg, G. K. Schenter, A. Hehn, A. Bussy, F. Belleflamme, G. Tabacchi, A. Glöß, M. Lass, I. Bethune, C. J. Mundy, C. Plessl, M. Watkins, J. VandeVondele, M. Krack, and J. Hutter. CP2K: An electronic structure and molecular dynamics software package – Quickstep: Efficient and accurate electronic structure calculations. *The Journal of Chemical Physics*, 152(19):194103:1–194103:47, May 2020. doi: 10.1063/5.0007045 1
- [28] A. H. Larsen, J. J. Mortensen, J. Blomqvist, I. E. Castelli, R. Christensen, M. Dulak, J. Friis, M. N. Groves, B. Hammer, C. Hargus, E. D. Hermes, P. C. Jennings, P. B. Jensen, J. Kermode, J. R. Kitchin, E. L. Kolsbjerg, J. Kubal, K. Kaasbjerg, S. Lysgaard, J. B. Maronsson, T. Maxson, T. Olsen, L. Pastewka, A. Peterson, C. Rostgaard, J. Schiøtz, O. Schütt, M. Strange, K. S. Thygesen, T. Vegge, L. Vilhelmsen, M. Walter, Z. Zeng, and K. W. Jacobsen. The atomic simulation environment—a Python library for working with atoms. *Journal of Physics: Condensed Matter*, 29(27):273002:1–273002:30, June 2017. doi: 10.1088/1361-648X/aa680e 2, 4, 6
- [29] P. M. Larsen, S. Schmidt, and J. Schiøtz. Robust structural identification via polyhedral template matching. *Modelling and Simulation in Materials Science and Engineering*, 24(5):055007:1–055007:18, May 2016. doi: 10.1088/0965-0393/24/5/055007 7
- [30] C. A. Lowry and D. C. Montgomery. A review of multivariate control charts. *IIE Transactions*, 27(6):800–810, Apr. 1995. doi: 10.1080/07408179508936797 6
- [31] X. Martinez, M. Chavent, and M. Baaden. Visualizing protein structures — Tools and trends. *Biochemical Society Transactions*, 48(2):499–506, Mar. 2020. doi: 10.1042/BST20190621 2
- [32] C. Massobrio, J. Du, M. Bernasconi, and P. S. Salmon, eds. *Molecular Dynamics Simulations of Disordered Materials*, vol. 215. Springer, New York, 2015. 1
- [33] P. K. Misra. *Physics of Condensed Matter*, chap. 17, pp. 567–597. Academic Press, Boston, Jan. 2012. doi: 10.1016/B978-0-12-384954-0.00017-7 3

- [34] D. Perez, E. D. Cubuk, A. Waterland, E. Kaxiras, and A. F. Voter. Long-time dynamics through parallel trajectory splicing. *Journal of Chemical Theory and Computation*, 12(1):18–28, Jan. 2016. doi: 10.1021/acs.jctc.5b00916 1, 2
- [35] D. Perez, B. P. Uberuaga, Y. Shim, J. G. Amar, and A. F. Voter. Accelerated molecular dynamics methods: Introduction and recent developments. In R. A. Wheeler, ed., *Annual Reports in Computational Chemistry*, vol. 5, chap. 4, p. 79–98. Elsevier, Amsterdam, Netherlands, Jan. 2009. doi: 10.1016/S1574-1400(09)00504-0 2
- [36] D. Perez, B. P. Uberuaga, and A. F. Voter. The parallel replica dynamics method – Coming of age. *Computational Materials Science*, 100:90–103, Apr. 2015. doi: 10.1016/j.commatsci.2014.12.011 2
- [37] E. F. Pettersen, T. D. Goddard, C. C. Huang, G. S. Couch, D. M. Greenblatt, E. C. Meng, and T. E. Ferrin. UCSF chimera—A visualization system for exploratory research and analysis. *Journal of Computational Chemistry*, 25(13):1605–1612, July 2004. doi: 10.1002/jcc.20084 1
- [38] S. Redona. SAMSON: Software for adaptive modeling and simulation of nanosystems. <https://samson-connect.net>, 2016. Accessed: 2023-03-22. 2
- [39] B. Reuter, K. Fackeldey, and M. Weber. Generalized Markov modeling of nonreversible molecular kinetics. *The Journal of Chemical Physics*, 150(17):174103:1–174103:12, May 2019. doi: 10.1063/1.5064530 4
- [40] K. Schatz, M. Krone, J. Pleiss, and T. Ertl. Interactive visualization of biomolecules’ dynamic and complex properties. *The European Physical Journal Special Topics*, 227(14):1725–1739, Mar. 2019. doi: 10.1140/epjst/e2019-800162-y 9
- [41] M. Scheurer, P. Rodenkirch, M. Siggel, R. C. Bernardi, K. Schulten, E. Tajkhorshid, and T. Rudack. PyContact: Rapid, customizable, and visual analysis of noncovalent interactions in MD simulations. *Biophysical Journal*, 114(3):577–583, Feb. 2018. doi: 10.1016/j.bpj.2017.12.003 2
- [42] D. Shalloway. Macrostates of classical stochastic systems. *The Journal of Chemical Physics*, 105(22):9986–10007, Sept. 1996. doi: 10.1063/1.472830 7
- [43] W. A. Shewhart. *Statistical Method from the Viewpoint of Quality Control*. Dover Publications, New York, 1986. 5
- [44] R. Skånberg, M. Linares, C. König, P. Norman, D. Jönsson, I. Hotz, and A. Ynnerman. VIA-MD: Visual interactive analysis of molecular dynamics. In *Proceedings of the Workshop on Molecular Graphics and Visual Analysis of Molecular Data*, pp. 19–27. The Eurographics Association, Eindhoven, June 2018. doi: 10.2312/molva.20181102 2
- [45] A. Stukowski. Visualization and analysis of atomistic simulation data with OVITO—The open visualization tool. *Modelling and Simulation in Materials Science and Engineering*, 18(1):015012:1–015012:7, Dec. 2009. doi: 10.1088/0965-0393/18/1/015012 2, 6, 7, 8
- [46] A. P. Thompson, H. M. Aktulga, R. Berger, D. S. Bolintineanu, W. M. Brown, P. S. Crozier, P. J. in ’t Veld, A. Kohlmeyer, S. G. Moore, T. D. Nguyen, R. Shan, M. J. Stevens, J. Tranchida, C. Trott, and S. J. Plimpton. LAMMPS - A flexible simulation tool for particle-based materials modeling at the atomic, meso, and continuum scales. *Computer Physics Communications*, 271:108171:1–108171:34, Feb. 2022. doi: 10.1016/j.cpc.2021.108171 1
- [47] Z. Tian, Z. Zhang, X. Jiang, F. Wei, S. Ping, and F. Wu. LaSCA: A visualization analysis tool for microstructure of complex systems. *Metals*, 13(2):415:1–415:10, Feb. 2023. doi: 10.3390/met13020415 2, 9
- [48] C. Tominski, H. Schumann, G. Andrienko, and N. Andrienko. Stacking-based visualization of trajectory attribute data. *IEEE Transactions on Visualization and Computer Graphics*, 18(12):2565–2574, Dec. 2012. doi: 10.1109/TVCG.2012.265 2
- [49] P. Ulbrich, M. Waldner, K. Furmanová, S. M. Marques, D. Bednář, B. Kozlíková, and J. Byška. sMolBoxes: Dataflow model for molecular dynamics exploration. *IEEE Transactions on Visualization and Computer Graphics*, 29(1):581–590, Jan. 2023. doi: 10.1109/TVCG.2022.3209411 2
- [50] J. Wang, L. Gou, H.-W. Shen, and H. Yang. DQNViz: A visual analytics approach to understand Deep Q-Networks. *IEEE Transactions on Visualization and Computer Graphics*, 25(1):288–298, Jan. 2019. doi: 10.1109/TVCG.2018.2864504 2
- [51] H. Wickham and L. Stryjewski. 40 years of boxplots. Technical report, Nov. 2011. 5
- [52] G. Wu, D. Lin, H. Wang, and L. Liu. Visual analysis of defect clustering in 3D irradiation damage simulation data. *Journal of Visualization*, 25(1):31–45, Feb. 2022. doi: 10.1007/s12650-021-00769-9 2, 9
- [53] J. S. Yi, Y. ah Kang, J. Stasko, and J. A. Jacko. Toward a deeper understanding of the role of interaction in information visualization. *IEEE Transactions on Visualization and Computer Graphics*, 13(6):1224–1231, Nov. 2007. doi: 10.1109/TVCG.2007.70515 6
- [54] P. Zou and R. Bader. A topological definition of a Wigner–Seitz cell and the atomic scattering factor. *Acta Crystallographica Section A: Foundations and Advances*, 50(6):714–725, Apr. 1994. doi: 10.1107/S0108767394003740 8