

Enabling Efficient Execution of In Situ Workflows

Tu Mai Anh Do Information Sciences Institute, University of Southern California

UTK Seminar, September 3rd, 2021

1





Acknowledgements





M. Taufer



E. Deelman



R. Ferreira da Silva



H. Weinstein



M. Cuendet



T. Estrada



S. Caino-Lores



L. Pottier



E. Kots



I. Lumsden

This work is funded by NSF contracts #1741040 and #1741057; and DOE contract #DE-SC0012636 슃







Outline



- 1. Motivations and Background
- 2. Contributions :
 - a. Modeling framework for in situ workflows
 - b. Performance indicators to determine performance of workflow ensembles
- 3. Conclusion



Molecular dynamics (MD) Protein backbone



- Molecular dynamics (MD) is a simulation model computing the atomic states of a molecular system evolving over time by observing interactions between atoms
- MD serves as a productive method to:
 - Control the configurations of the molecular systems
 - Observe important processes at atomic resolution



Razavi et al, 2017





Post-hoc Analysis





I/O limitations





- The increase in computing capability helps the MD simulations generate more data that needs to be analyzed (150,000 atoms + 500,000 snapshots would generate ~ 1.8TB data)
- However, the I/O bandwidth does not grow at the same pace \rightarrow I/O bottleneck





In situ analysis





In situ Analysis

- Data is analyzed as soon as generated
- Decoupling analysis from the simulation to interleave their executions → Reduce time-to-solution
- Study insights into phenomena of the molecular system in a timely fashion





Comparison of capacity, latency for most important technologies of the memory and storage hierarchies. (Lüttgau et al., 2018)



Information Sciences Institute



Page 8

• We focus on memory-to-memory for the lowest latency



Placement variants





In situ framework design



Problem: How to enable in situ execution of simulations and analyses?

Challenges/complexities:

- Decoupling analyses from the simulation to perform in situ requires to manage data staging to additional concurrent components
- Orchestrating data coupling over iterations
- Data incompatibility between decoupled components (simulations, analyses)

We model/design a framework that allows to decouple in situ analyses from the simulation to address these complexities:

- Data staging, coupling \rightarrow Data transport layer
- Data compatibility \rightarrow Data abstraction



In situ framework implementation





- Two main components:
 - Data transport layer (DTL) provides interfaces to support different storage tiers
 - DTL plugins interact with the DTL via data chunk abstraction for compatibility
- In this work
 - The DTL is implemented with the help of DIMES (Zhang et al., 2017) to enable in-memory data staging
 - The DTL plugin is integrated with Plumed (Tribello et al., 2014) to offer non-intrusive approach







In situ Workflow Ensembles





Performance indicators for in situ workflow ensembles



Problem: How to evaluate performance of an in situ workflow ensembles?

Challenges:

- Using traditional metrics and capturing them separately are not straightforward to synthesize performance of entire workflow ensemble
- Need to take into account resources to characterize performance under resource cost

We introduce multi-stage performance indicators that capture performance of the entire in situ workflow ensembles in terms of multiple resource perspectives

- **Resource usage**: How efficiently the resource is utilized?
- Resource allocation: How efficiently the simulations, analyses are placed on allocated resources?
- **Resource provisioning**: How many resources are provisioned to execute efficiently?



Multi-stage performance indicators



Information Sciences Institute

Page 15



Experiment setup





Medium-scale all-atom system containing the GltPh transporter protein (Akyuz 2015) implemented in GROMACS (P Bjelkmar et al., 2010)



Collective variable (largest eigenvalue of bipartite distance matrices between two substructures) (Barducci 2011, Johnston 2017)

- Execution platform: Cori, a Cray XC40 supercomputer at NERSC. Each compute node is equipped with
 - 2 Intel Xeon E5-2698 v3 (16 cores each)
 - 128 GB of DRAM
- TAU (Shende et al., 2006) is leveraged to collect performance information



One analysis per simulation







Takeaways



- The proposed indicators can be leveraged for evaluating scheduling decision of in situ ensemble under resource constraints
- Provide hints to improve effectiveness of resource usage → optimizing simulation exploration by running many MD simulations at a time
- Future work will consider leveraging the proposed indicators for scheduling in situ components of a workflow ensemble to enable high-throughput ensemble of simulations







Information Sciences Institute





References



- Simakov et al., 2018. A Workload Analysis of NSF's Innovative HPC Resources Using XDMoD. arXiv preprint arXiv:1801.04306 (2018)
- Lüttgau et al., 2018. Survey of Storage Systems for High-Performance Computing. Supercomputing Frontiers And Innovations, 5(1), 31-58.
- Chelli et al., 2012. Serial Generalized Ensemble Simulations of Biomolecules with Self-Consistent Determination of Weights. Journal of Chemical Theory and Computation 8, 3.
- Sanchez-Martinez et al., 2015. Enzymatic minimum free energy path calculations using swarms of trajectories. J Phys Chem B. 2015;119(3):1103–13
- Ballard et al., 2009. Replica exchange with nonequilibrium switches. Proc Natl Acad Sci. 2009;106(30):12224–9
- Comer et al., 2014. Multiple-Replica Strategies for Free-Energy Calculations in NAMD: Multiple-Walker Adaptive Biasing Force and Walker Selection Rules. Journal of Chemical Theory and Computation 10, 12 (2014).
- Bowman et al., 2010. Enhanced Modeling via Network Theory: Adaptive Sampling of Markov State Models. Journal of Chemical Theory and Computation 2010 6 (3), 787-794
- P Bjelkmar et al., 2010. Implementation of the CHARMM Force Field in GROMACS: Analysis of Protein Stability Effects from Correction Maps, Virtual Interaction Sites, and Water Models. J. Chem. Theory Comput.6, 2 (2010).
- Tribello et al., 2014. Plumed 2: New feathers for an old bird. Comput. Phys. Commun., 185(2):604–613.





References

- Razavi et al, 2017. A Markov State-based Quantitative Kinetic Model of Sodium Release from the Dopamine Transporter. Sci Rep 7, 40076
- Pronk et al., 2011. Copernicus: a new paradigm for parallel adaptive molecular dynamics. In Proceedings of 2011 International Conference for High Performance Computing, Networking, Storage and Analysis (SC '11). Association for Computing Machinery, New York, NY, USA, Article 60, 1–10.
- Hruska et al., 2020. Extensible and Scalable Adaptive Sampling on Supercomputers. Journal of Chemical Theory and Computation 2020 16 (12), 7915-7925
- Johnston et al. "In situ data analytics and indexing of protein trajectories", Journal of Computational Chemistry, 38 (16), 1419-1430, (2017)
- Akyuz et al. 2015. Transport domain unlocking sets the uptake rate of an aspartate transporter. Nature 518, 7537(2015)
- Barducci et al. 2011. Meta-dynamics.WIREs Computational Molecular Science 1, 5 (2011)
- Johnston et al. In situ data analytics and indexing of protein trajectories. Journal of Computational Chemistry, 38 (16), 1419-1430, (2017)
- Fan Zhang et al., 2017. In-memory staging and data-centric task placement for coupled scientific simulation workflows.Concurrency and Computation: Practice and Experience 29, 12 (2017)
- Agelastos et al., 2014. The Lightweight Distributed Metric Service: A Scalable Infrastructure for Continuous Monitoring of Large Scale Computing Systems and Applications. International Conference for High Performance Computing, Networking, Storage and Analysis (SC '14). 154–165.



References



- Chad Wood et al., 2016. A Scalable Observation System for Introspection and in Situ Analytics. The 5th Workshop on Extreme-Scale Programming Tools (ESPT '16). IEEE Press, Salt Lake City, Utah, 42–49.
- Taufer et al., 2019.Characterizing In Situ and In Transit Analytics of Molecular Dynamics Simulations for Next-Generation Supercomputers. 15th International Conference on eScience (eScience)
- Zou et al., 2014. FlexAnalytics: A flexible data analytics framework for big data applications with I/O performance improvement. Big Data Research 1, 4 13.
- Zhang et al., 2016. WOWMON: A Machine Learning-based Profiler for Self-adaptive Instrumentation of Scientific Workflows. Procedia Computer Science 80 (2016), 1507–1518. International Conference on Computational Science 2016, ICCS.
- Malakar et al., 2015. Optimal scheduling of in-situ analysis for large-scale scientific simulations. International Conference for High Performance Computing, Networking, Storage and Analysis, pp. 1-11
- Shende et al., 2006. The Tau parallel performance system. The International Journal of High Performance Computing Applications 20, 287–311.
- Lofstead et al., 2008. Flexible io and integration for scientific codes through the adaptable io system (adios),6th international workshop on Challenges of large applications in distributed environments
- Docan et al., 2012.Dataspaces: an interaction and coordination framework for coupled simulation workflows. Cluster Computing .
- Dreher, M., et al., 2017. Decaf: Decoupled dataflows for in situ high-performance workflows



In-line in situ analysis



Stride

- Data is analyzed as soon as generated
- The simulation and analysis interchangeably execute
- Analysis needs to be embedded in simulation code
- Less robust







Efficiency model



- Computing in situ step is lightweight and can be performed online
- *In situ step* is leveraged to estimate:
 - 1. Makespan: duration of an in situ step x number of in situ steps overlapped part between in situ steps
 - 2. Useful computation: measured by time for computation except idle time during an in situ step
- Evaluate computational efficiency:

Maximize
$$E = \frac{\text{Estimation of useful computation}}{\text{Estimation of makespan}}$$
 \rightarrow Minimize idle time \rightarrow Minimize makespan \rightarrow Minimize makespan



Characterization challenges



- Evaluating each ensemble component/member exclusively is hard to:
 - Generate a full picture of the workflow ensemble performance
 - Compare between different executions/configurations
- Without taking into account resources, the performance could be misleading







Component placement



- The simulation is co-located with the analysis, iff $|s| = |s \cup a|$
- The simulation and analysis are assigned to different nodes, iff $|s| < |s \cup a|$

Set of node indexes where a simulation is executed



Set of node indexes where the coupled analysis is executed

Placement indicator of ensemble member i with K_i analyses

$$CP_i = rac{1}{K_i} \sum_{j=1}^{K_i} rac{|s_i|}{|s_i \cup a_i^j|}$$

Mean of ratios forming by all (simulation, analysis) pairs

Maximize placement indicator prioritizes placements that minimize the number of computing resources (number of compute nodes) used by that ensemble member.

> USC Viterbi School of Engineering

Performance indicator P_i of ensemble member i









Synthesis of performance indicators

• P_i can be either $P_i^{U}, P_i^{U,A}, P_i^{U,P}, P_i^{U,A,P} (= P_i^{U,P,A})$

• The objective function of N ensemble members (the higher the better)

Maximize
$$F(P_i) = \overline{P} - \sqrt{\frac{1}{N} \sum_{i=1}^{N} (P_i - \overline{P})^2}$$
 where $\overline{P} = \frac{1}{N} \sum_{i=1}^{N} P_i$
Maximize average performance Mean Standard deviation \rightarrow Minimize variability among ensemble members





