

Leveraging Modern Data Processing & Engineering Techniques for Cosmological Simulations

Soumyodeep Mukherjee
Associate Director, Commercial Data Engineering,
Genmab, Plainsboro, NJ, US1

ABSTRACT

Cosmological simulations play a crucial role in understanding the formation and evolution of the universe. As these simulations generate and process vast amounts of data, applying modern data engineering & processing techniques becomes essential to manage, analyze, and derive meaningful insights from this information. This paper explores how these techniques can optimize the performance, scalability, and accuracy of cosmological simulations. The focus is on the integration of distributed computing, real-time data processing, and advanced storage solutions to enhance simulations. Furthermore, examination was done to determine initial boundary conditions from observational data & discussion has been included in the paper on popular models used to simulate the universe's evolution and consideration of methods for tuning simulation parameters to balance accuracy with manageable data growth. Estimations of data generation and computational requirements are also provided, emphasizing the role of cloud computing in handling these challenges.

General Terms

Data Processing on Cosmological Simulation leveraging Data Engineering at scale

Keywords

Cosmological Simulations; Modern Data Engineering; Distributed Computing; High-Performance Computing; Cloud Computing

1. INTRODUCTION

The study of cosmology seeks to understand the large-scale structure and evolution of the universe. To achieve this, cosmologists rely on simulations that model the formation of galaxies, clusters, and other cosmic structures over billions of years. These simulations generate petabytes of data, necessitating advanced modern data engineering & processing techniques to manage and process this information effectively.

2. INITIAL BOUNDARY CONDITIONS

Accurate initial boundary conditions are essential for realistic cosmological simulations. These conditions are typically derived from:

2.1 Cosmic Microwave Background (CMB)

Data

The CMB provides a snapshot of the early universe, offering insights into temperature fluctuations and density variations. This data, obtained from missions like WMAP and Planck, is used to set the initial density perturbations in simulations.

2.2 Large-Scale Structure Surveys

Observations from surveys such as SDSS and DES help determine the large-scale distribution of matter, guiding the initial conditions for simulations.

2.3 Primordial Power Spectrum

The primordial power spectrum $P(k)$ is inferred from CMB data and describes the initial density fluctuations as a function of the wavenumber k . It influences the distribution of matter in simulations, where $P(k) \propto k^n$ with $n \approx 1$ (scale-invariant spectrum).

2.4 N-Body Simulations

These simulations evolve an initial particle distribution under gravitational forces according to the N-body equations:

$$\frac{d^2 \mathbf{r}_i}{dt^2} = -G \sum_{j \neq i} \frac{m_j (\mathbf{r}_i - \mathbf{r}_j)}{|\mathbf{r}_i - \mathbf{r}_j|^3}$$

where \mathbf{r}_i and \mathbf{r}_j are the positions of the i -th and j -th particles, respectively, and G is the gravitational constant.

3. POPULAR MODELS FOR COSMOLOGICAL SIMULATIONS

Several models are used to simulate the universe's evolution, each with specific strengths

3.1 Λ CDM Model

The standard cosmological model, Λ CDM, incorporates dark energy (Λ) and cold dark matter (CDM). The evolution of the universe in this model is governed by the Friedmann equations.

$$\left(\frac{\dot{a}}{a}\right)^2 = \frac{8\pi G}{3} \rho - \frac{k}{a^2} + \frac{\Lambda}{3}$$

where $a(t)$ is the scale factor, ρ is the total energy density, k is the curvature parameter, and Λ is the cosmological constant

3.2 Hydrodynamical Simulations

These simulations solve the equations of hydrodynamics along with gravity to model the evolution of baryonic matter. The Navier-Stokes equations, along with the equation of state for the gas, are solved numerically:

$$\frac{\partial \rho}{\partial t} + \nabla \cdot (\rho \mathbf{v}) = 0$$
$$\frac{\partial (\rho \mathbf{v})}{\partial t} + \nabla \cdot (\rho \mathbf{v} \otimes \mathbf{v} + p \mathbf{I}) = -\rho \nabla \Phi$$

where ρ is gas density, \mathbf{v} is the velocity, p is the pressure and Φ is the gravitational potential.

3.3 N-body Simulations:

These focus on the evolution of dark matter particles under gravity. The collisionless Boltzmann equation describes the phase-space distribution function $f(\mathbf{r}, \mathbf{v}, t)$

3.4 Modified Gravity Models

These models explore alternatives to general relativity, modifying the Poisson equation for gravity, for example, in $f(R)$ theories:

$$f'(R)R_{\mu\nu} - \frac{1}{2}f(R)g_{\mu\nu} + (g_{\mu\nu} \square - \nabla_\mu \nabla_\nu)f'(R) = 8\pi G T_{\mu\nu}$$

3.5 Cosmic Reionization Simulations

These simulations incorporate radiative transfer equations to model the propagation of ionizing radiation and its impact on the intergalactic medium (IGM)

4. MODERN DATA ENGINEERING & PROCESSING TECHNIQUES FOR COSMOLOGICAL SIMULATIONS

To handle the complexity and scale of cosmological simulations, advanced modern data engineering & processing techniques are employed.

4.1 Scalability and Performance

Distributed computing platforms, such as Apache Spark, enable parallel processing of the data generated by cosmological simulations. For example, an N-body simulation might be divided into spatial domains, each processed on a different node.

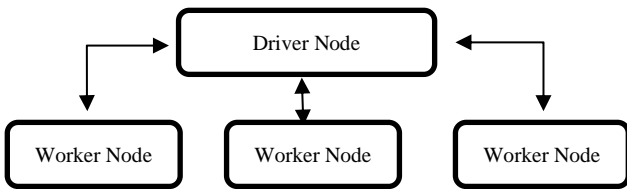


Fig 1: Apache Spark (distributed architecture)

4.2 Data Storage and Management

With data volumes exceeding petabytes, distributed storage solutions like HDFS provide the necessary scalability. Techniques such as data sharding and partitioning are essential to manage and retrieve simulation data efficiently.

4.3 Real-Time Data Processing

Frameworks like Apache Kafka allow for real-time streaming and analysis of simulation data. For instance, changes in the density field can be monitored in real time to adjust parameters dynamically.

4.4 Data Visualization and Interpretation

Visualization tools, such as ParaView or yt, are integrated with simulation frameworks to render large-scale 3D visualizations of cosmic structures, making use of data reduction techniques like PCA to handle the vast amounts of data.

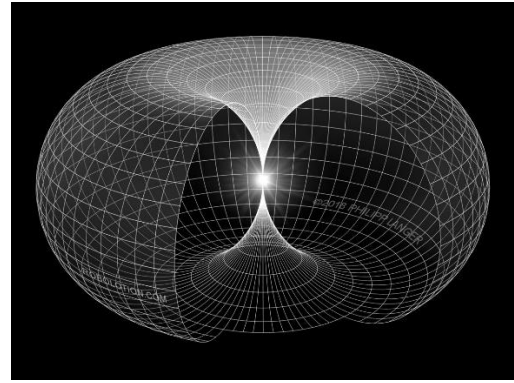


Fig 2: Representation of Cosmic Simulation

5. TUNING SIMULATION PARAMETERS TO MANAGE DATA GROWTH

One of the significant challenges in cosmological simulations is the exponential growth of data as the complexity and resolution of models increase. To manage this

5.1 Parameter Tuning

Consider a simulation with N particles. The computational complexity typically scales as $O(N \log N)$ for tree-based algorithms and $O(N^2)$ for direct methods. Reducing N by a factor of 2 can reduce the computational cost by a factor of 4 in direct methods, illustrating the importance of parameter tuning.

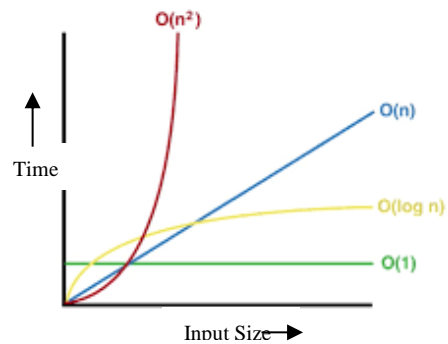


Fig 3: Comparison of time complexity for various Big O representations

5.2 Balancing Accuracy and Efficiency

Adaptive mesh refinement (AMR) techniques adjust the resolution dynamically based on local density variations, thereby reducing the total number of cells required while maintaining accuracy. For example, if the resolution scales as h^{-3} where h is the cell size, then using AMR reduces the number of high-resolution cells significantly.

5.3 Data Compression

The application of Fourier transforms to compress density fields, or wavelet transforms to capture multiscale structures, can reduce the storage requirements without losing critical information.

6. DATA GENERATION AND RESOURCE ESTIMATIONS

To illustrate the data generation and computational requirements, let us assume a simulation with $N=10^{12}$ particles representing dark matter and baryonic matter over a cosmological volume of approximately 1 cubic gigaparsec:

- **Data Generation:** Assume each particle has position $\mathbf{r}_i = (x_i, y_i, z_i)$ and velocity $\mathbf{v}_i = (v_{x,i}, v_{y,i}, v_{z,i})$ and additional properties such as mass m_i . If each property requires 16 bytes

(double precision), each particle requires approximately 100 bytes. The total data per snapshot is:

$$\text{Data per snapshot} = 10^{12} \times 100 \text{ bytes} = 100 \text{ TB}$$

Given M snapshots, the total data generated is 100M TB. If M = 1000, the total data is approximately 100 PB.

- **Computational Resources:** The computational load T can be estimated by considering the number of operations per timestep. For an $O(N \log N)$ algorithm over K timesteps, the computational time scales as:

$$T \sim \alpha N \log N \times K$$

where α is a constant that depends on the specifics of the hardware and algorithm used. If $N=10^{12}$ and $K=10^6$, a rough estimate of the computational load can be made:

$$T \sim \alpha \times 10^{12} \times \log(10^{12}) \times 10^6$$

Assuming $\log(10^{12}) \approx 27.6$:

$$T \sim \alpha \times 2.76 \times 10^{18}$$

For a given computational resource, such as a supercomputer with p petaflops (floating-point operations per second), the total wall-time W required for the simulation can be estimated as:

$$W = T / P$$

If p=1 petaflops:

$$W \sim 2.76 \times 10^{18} / 10^{15} \text{ seconds} = 2.76 \times 10^3 \text{ seconds} \\ \approx 46 \text{ minutes per timestep}$$

Thus, for a million timesteps, the total simulation time would be approximately 46 million minutes, or around 87 years, underscoring the importance of optimization and parallelization in cosmological simulations.

7. METHODOLOGY

This study focuses on integrating modern data engineering and processing techniques into cosmological simulations to address challenges related to scale, complexity, and data management. The methodology employed in this research involves the following steps:

7.1. Framework Selection:

- **Distributed Computing:** Apache Spark and Hadoop Distributed File System (HDFS) were selected to process and manage the vast datasets generated by cosmological simulations. These frameworks enable the partitioning and parallel processing of data to improve scalability and performance.
- **Real-Time Data Processing:** Apache Kafka was utilized to process streaming data, allowing for dynamic adjustments in simulation parameters based on real-time changes in the density fields.

7.2. Simulation Initialization:

- **Initial Boundary Conditions:** Data from Cosmic Microwave Background (CMB) and large-scale structure surveys were analyzed to establish accurate initial conditions. The primordial power spectrum was

used to set the density perturbations, ensuring physical consistency.

- **Parameter Optimization:** Key parameters, such as particle counts and resolution, were optimized using adaptive mesh refinement (AMR) techniques to balance computational efficiency and accuracy.

7.3. Algorithm Implementation:

- **Gravitational N-Body Simulations:** These simulations employed a tree-based algorithm to reduce computational complexity $O(N \log N)$. Additional optimizations, such as hierarchical time-stepping, ensured computational efficiency.
- **Hydrodynamical Simulations:** The equations of hydrodynamics were solved using finite volume methods, and numerical stability was maintained through adaptive time-stepping.

7.4. Data Management and Compression:

- **Data Storage:** Simulation data exceeding petabyte scales were managed using distributed storage solutions like HDFS, with data partitioning to ensure efficient retrieval.
- **Data Compression:** Techniques such as Fourier transforms and wavelet compression were applied to reduce storage requirements while preserving critical structural information.

7.5. Visualization and Analysis:

- **High-resolution visualization tools,** including ParaView and yt, were employed to render the large-scale structure of the universe. Dimensionality reduction techniques, such as Principal Component Analysis (PCA), were applied to facilitate interpretation of high-dimensional data.

7.6. Evaluation:

- The methodology was evaluated using multiple datasets, representing diverse cosmological scenarios. Key metrics included computational efficiency, data throughput, and the accuracy of simulated structures compared to observational data.

8. CONCLUSION

Modern data engineering and processing techniques have demonstrated their indispensability in managing the immense computational and data demands of cosmological simulations. These advancements enable researchers to create accurate and scalable models of the universe's evolution, tackling challenges such as exponential data growth, computational bottlenecks, and dynamic parameter tuning. By leveraging distributed computing, real-time data processing, and optimized storage solutions, cosmologists can enhance the granularity and fidelity of simulations without compromising efficiency.

Future Scope: The ongoing evolution of computational resources and methodologies offers exciting opportunities to further advance this field:

- **Integration with Quantum Computing:** As quantum computing matures, its potential to handle complex computations at unprecedented speeds could revolutionize cosmological simulations. This could include solving N-body gravitational equations more efficiently or simulating quantum effects in the early universe.
- **Multi-Scale Modeling:** The development of techniques that seamlessly integrate simulations across vastly

different spatial and temporal scales will allow for more comprehensive models that bridge the gap between local phenomena and large-scale cosmic structures.

- **AI-Driven Insights:** Machine learning and artificial intelligence could be employed for automated parameter tuning, anomaly detection, and predictive modeling, further optimizing simulations, and uncovering new insights from data.
- **Real-Time Observational Feedback:** The integration of real-time data from next-generation telescopes (e.g., JWST, SKA) with simulations can allow for adaptive models that evolve in sync with observational discoveries.
- **Improved Data Accessibility:** Efforts to democratize access to simulation data through cloud-based platforms and open datasets will enable broader collaboration and innovation, allowing researchers worldwide to contribute to and benefit from these simulations.
- **Sustainability in Computing:** As the scale of simulations grows, optimizing energy consumption and developing sustainable practices for high-performance computing will become increasingly important.

By addressing these future directions, researchers can continue to push the boundaries of our understanding of the cosmos, unraveling the mysteries of its formation and evolution while laying the foundation for interdisciplinary innovations. These advancements will not only deepen our comprehension of the universe but also pave the way for breakthroughs in data science, computational physics, and high-performance computing.

9. REFERENCES

- [1] Planck Collaboration. "Planck 2018 results. VI. Cosmological parameters." *Astronomy & Astrophysics*, 641 (2020): A6.
- [2] Eisenstein, D. J., et al. "SDSS-III: Massive Spectroscopic Surveys of the Distant Universe, the Milky Way Galaxy, and Extra-Solar Planetary Systems." *The Astronomical Journal* 142.3 (2011): 72.
- [3] Navarro, J. F., Frenk, C. S., & White, S. D. M. "A Universal Density Profile from Hierarchical Clustering." *The Astrophysical Journal* 490.2 (1997): 493-508.
- [4] Springel, V., et al. "Simulations of the formation, evolution and clustering of galaxies and quasars." *Nature* 435.7042 (2005): 629-636.
- [5] WMAP Collaboration, Hinshaw, G., et al. "Five-Year Wilkinson Microwave Anisotropy Probe (WMAP) Observations: Data Processing, Sky Maps, and Basic Results." *The Astrophysical Journal Supplement Series* 180.2 (2009): 225.
- [6] Springel, V., et al. "The large-scale structure of the Universe." *Nature* 435.7042 (2005): 629-636.
- [7] Bond, J. R., et al. "Excursion set mass functions for hierarchical Gaussian fluctuations." *The Astrophysical Journal* 379 (1991): 440-460.
- [8] Davis, M., et al. "The evolution of large-scale structure in a universe dominated by cold dark matter." *The Astrophysical Journal* 292 (1985): 371-394.
- [9] Pfrommer, C., et al. "Simulating cosmic rays in clusters of galaxies – III. Non-thermal scaling relations and comparison to observations." *Monthly Notices of the Royal Astronomical Society* 378.2 (2007): 385-402.
- [10] Weinberg, S. "Gravitation and Cosmology: Principles and Applications of the General Theory of Relativity." Wiley, 1972.
- [11] Klypin, A., et al. "Resolving the Structure of Cold Dark Matter Halos." *The Astrophysical Journal* 522.1 (1999): 82-92.
- [12] Navarro, J. F., Frenk, C. S., & White, S. D. M. "The Structure of Cold Dark Matter Halos." *The Astrophysical Journal* 462 (1996): 563.
- [13] White, S. D. M., & Rees, M. J. "Core condensation in heavy halos: A two-stage theory for galaxy formation and clustering." *Monthly Notices of the Royal Astronomical Society* 183.3 (1978): 341-358.
- [14] Peebles, P. J. E. "The large-scale structure of the universe." Princeton University Press, 1980.
- [15] Mo, H. J., van den Bosch, F. C., & White, S. D. M. "Galaxy Formation and Evolution." Cambridge University Press, 2010.