# Literature Review

## Computational Constraints in Astrophysical Simulation

Over the past several decades, progress in computational astrophysics has closely followed the performance improvements predicted by Moore's Law, with central processing unit (CPU) speeds increasing at a near-exponential rate. Once an algorithm was implemented, substantial performance gains could often be achieved simply by running existing code on newer hardware, with minimal additional development effort. However, as single-core CPU performance has plateaued, this implicit scaling model has begun to break down. As a result, continued advances in computational astrophysics increasingly depend on exploiting parallelism and adapting algorithms to emerging computing architectures. @

One of the most significant architectural shifts has been the widespread adoption of graphics processing units (GPUs), which offer massive parallelism and high memory bandwidth. While GPUs provide the potential for orders-of-magnitude performance improvements, realizing these gains typically requires substantial algorithmic reformulation. Algorithms that scale well on serial or modestly parallel CPUs may perform poorly on highly parallel architectures if they involve irregular control flow or memory access patterns.

## Gravitational N-Body Simulations and Quadratic Scaling

Among the most computationally demanding problems in astrophysics are gravitational N-body simulations, which form the foundation of many studies in galactic dynamics and structure formation. In such simulations, a galaxy is modeled as a system of particles (representing stars or dark matter) whose evolution is governed by their mutual gravitational interactions.

According to Newton's law of universal gravitation, the force exerted on a particle of mass $m_i$ by another particle of mass $m_j$ is given by

$$\vec{F}_{ij} = G\frac{m_i m_j (\vec{r}_j - \vec{r}_i)}{|\vec{r}_j - \vec{r}_i|^3}$$

The total force acting on particle $i$ is obtained by summing the contributions from all other particles in the system,

$$\vec{F}_i \sum_{j\neq i}^{N} G\frac{m_i m_j (\vec{r}_j - \vec{r}_i)}{|\vec{r}_j - \vec{r}_i|^3}$$

Evaluating this expression requires computing $N-1$ pairwise interactions per particle. Consequently, a direct implementation of gravitational force evaluation scales as $O(N^2)$ per simulation timestep. This quadratic scaling rapidly becomes computationally prohibitive as the number of particles increases, particularly when long integration times are required to observe large-scale dynamical phenomena such as spiral structure or halo evolution.

## Hierarchical Approximation and the Barnes–Hut Algorithm

To address the computational limitations imposed by direct force evaluation, a variety of approximation techniques have been developed to reduce the cost of N-body simulations while preserving essential physical behavior. One of the most influential of these methods is the Barnes–Hut algorithm, which exploits the hierarchical spatial structure of particle distributions to approximate gravitational interactions.

Originally introduced by Barnes and Hut in 1986, the algorithm organizes particles into a tree structure—typically a quadtree in two dimensions or an octree in three dimensions. During force

evaluation, distant groups of particles are approximated as a single effective mass located at their center of mass, allowing subsets of the force summation to be replaced by a single interaction. An opening-angle parameter controls the accuracy of this approximation, providing a tunable tradeoff between computational efficiency and force accuracy.

By replacing many distant pairwise interactions with aggregate approximations, Barnes–Hut reduces the computational complexity of force evaluation from $O(N^2)$ to approximately $O(N \log N)$. This reduction makes the algorithm particularly well-suited for simulating large, collisionless systems such as galaxies, where global dynamical behavior is often of greater interest than exact short-range interactions.

## Numerical Integration and Stability Considerations

In addition to force evaluation, the accuracy and long-term stability of N-body simulations depend critically on the choice of numerical integration scheme. Simple methods such as forward Euler integration are computationally inexpensive but suffer from poor energy conservation, leading to unphysical behavior over long time scales. As a result, symplectic integrators such as Verlet or leapfrog methods are commonly employed in astrophysical simulations due to their improved conservation properties.

In Barnes–Hut simulations, integration error interacts with approximation error introduced by hierarchical force evaluation. Prior studies have shown that, for many galactic-scale simulations, integration error can dominate approximation error if inappropriate time-stepping schemes are used. Consequently, careful selection of integration methods remains essential even when approximate force models are employed.

## GPU Acceleration of Hierarchical N-Body Methods

The computational demands of gravitational N-body simulations have made them a natural target for graphics processing units (GPUs), which provide massive data parallelism and high memory bandwidth. Modern GPUs are designed with thousands of small, efficient cores capable of executing the same operation on many data elements simultaneously. This architecture is particularly well-suited to the independent, particle-wise force computations inherent in N-body simulations.

To leverage GPUs effectively, computational tasks are expressed in terms of shaders—small programs that run on the GPU cores. Originally, shaders were designed to compute visual effects for rendering pipelines, including vertex transformations, fragment coloring, and texture operations. Over time, GPU programming languages such as CUDA and OpenCL have generalized this model to allow general-purpose computation, enabling the parallel execution of tasks that are not directly related to graphics, including N-body simulations.

Despite these advancements, mapping hierarchical algorithms like Barnes–Hut onto GPU architectures is challenging. Tree construction and traversal involve irregular memory access and control flow divergence, which reduce the efficiency of parallel execution. Optimizations such as linearized trees, stackless traversal, and space-filling curves have been proposed to improve memory coherence and parallel efficiency. While these approaches have demonstrated substantial speedups in native GPU implementations, they require low-level memory control and recursive data structures that are not directly available in web environments.

## Web-Based Simulation and WebGL

Web-based platforms offer the advantage of accessibility, allowing simulations to run on any device with a modern browser. Early web-based N-body simulations relied primarily on WebGL, a web standard for rendering interactive 2D and 3D graphics. WebGL exposes the GPU through a

rendering pipeline that is largely designed for graphics applications, relying on vertex and fragment shaders to manipulate visual elements.

In this model, computation that is not directly related to rendering—such as calculating gravitational forces—is often forced into the graphics pipeline, using hacky tricks such as encoding particle data into textures and interpreting colors as numerical values. While this enables real-time visualization and interactive particle systems, the approach has several limitations:

- Shaders in WebGL are designed for graphics, not physics computation
- Recursive or tree-based algorithms like Barnes–Hut are difficult to implement
- Numerical precision is limited, often constraining simulations to visually plausible but physically approximate behavior

As a result, most prior web-based particle simulations have focused on graphics and educational interactivity, rather than accurate physical modeling of gravitational systems.

## WebGPU: General-Purpose GPU Programming in the Browser

WebGPU is a new web standard that introduces general-purpose GPU compute to browser environments. Unlike WebGL, WebGPU is not restricted to the graphics pipeline and provides:

- Compute shaders: GPU programs that perform arbitrary calculations on buffers of data, independent of rendering
- Structured buffers: Typed memory storage for particle data, enabling efficient hierarchical computation
- Explicit memory control and synchronization: Allowing parallel computations with predictable performance

WebGPU's architecture closely mirrors modern native GPU APIs such as Vulkan, Metal, and Direct3D 12, enabling more complex algorithms like Barnes–Hut to be implemented in-browser. Compute shaders in WebGPU can perform hierarchical tree traversal and force accumulation without being forced through the rendering pipeline, making accurate particle physics simulations feasible in a web context.

## Positioning of This Work

While native GPU implementations of Barnes–Hut have been extensively studied, prior web-based simulations have largely emphasized visualization and graphics rather than physics-accurate N-body computation. Existing WebGL approaches typically rely on approximations that simplify or bypass hierarchical force computation, limiting their scientific validity.

WebGPU presents a unique opportunity to bridge this gap: by enabling true general-purpose GPU computation in the browser, it allows for hierarchical algorithms such as Barnes–Hut to be implemented without forcing the computation through the graphics pipeline. WebGPU provides the parallel performance necessary to handle large particle counts efficiently, making physically accurate simulations feasible in a lightweight, interactive environment.

This thesis investigates how Barnes–Hut hierarchical N-body simulations can be adapted to WebGPU, aiming to demonstrate that accurate galaxy simulations can be performed in-browser without compromising performance. By doing so, this work seeks to democratize astrophysical simulation, allowing users to explore large-scale galactic dynamics interactively without requiring high-end hardware, complex software setup, or native GPU access.