# INVESTIGATING ECO-FRIENDLY RENDERING SOLUTIONS FOR 3D MODELS ON MULTI-CORE PROCESSORS

*V.V.S. Rao, Mukesh Ranjan and Prof. V. Kunwar*

**ABSTRACT:**

*Graphics Processing Units (GPUs) are vital for handling computationally intensive tasks, particularly in real-time 3D applications. Designed with additional processing power, GPUs are optimized to efficiently render complex scenes. Shaders, which are small programs running on the GPU, control how each pixel is processed and rendered. They affect key aspects of the rendering pipeline, such as geometric complexity, texture resolution, and per-pixel lighting. However, the extensive use of shaders can lead to increased power consumption, especially when managing complex 3D scenes. This study investigates GPU power usage during the rendering process and proposes a dynamic power prediction model. Initially, the power consumption of various rendering configurations is analyzed to understand the impact of shader parameters. Based on this analysis, a power-aware, dynamically reconfigurable rendering model is introduced, designed to optimize power consumption through shader calls while maintaining performance and visual quality. The implementation integrates TensorFlow management libraries and dynamic voltage and frequency scaling (DVFS), a technique that adjusts the GPU's voltage and frequency based on workload demands. Shader parameters such as geometric complexity, texture resolution, per-pixel lighting, and filtering are evaluated for their power usage during common GPU tasks. By continuously monitoring the framerate, the system dynamically adjusts the GPU frequency to optimize power efficiency. A lower framerate limit of 30 frames per second (fps) and an upper limit of 60 fps are established. If the framerate exceeds 60 fps, the GPU frequency is reduced to conserve power, while the frequency is increased if the framerate drops below 30 fps to maintain smooth performance. This dynamic adjustment ensures the GPU operates efficiently without sacrificing the quality of 3D rendering. The experimental results demonstrate that the proposed model offers around 40% power savings compared to previous methods. Additionally, it improves computational speed while preserving the quality of the rendered scenes. This approach highlights the potential for achieving a balance between power efficiency and rendering performance in GPUs, making it ideal for power-sensitive applications.*

*Index Terms: GPU; 3D graphics rendering; Shader parameters; TensorFlow; DVFS; Workloads.*

***Biographical notes:***

**V.V.S. Rao-** *M. Tech student Department of AI &ML, VIT, Vellore*

**Mukesh Ranjan-** *M. Tech student Department of AI &ML, VIT Vellore*

**Prof. V. Kunwar –** *Associate Professor, Department of AI&ML, VIT Vellore*

## 1. INTRODUCTION

Rendering 3D models in real-time is a crucial and challenging task in computer graphics, requiring a balance between visual quality and hardware constraints such as power consumption and computational overhead. There are two main categories of 3D model rendering: real-time dynamic rendering and pre-rendered scenes. Real-time rendering is particularly demanding because it involves the immediate synthesis of 3D scenes on raw graphics hardware, ensuring flexibility, interactivity, and high-quality visuals. This technique is essential for applications like video games, virtual reality (VR), and augmented reality (AR), where responsiveness and user interaction are critical. On the other hand, pre-rendered scenes are typically used in areas like 3D movie production, architectural visualization, and simulations, where accuracy and visual fidelity take precedence over speed. The algorithms involved in these cases require more time for computation, as they simulate complex lighting physics, including reflection, refraction, and shading. This computational intensity makes these techniques unsuitable for real-time applications, as their focus is on producing detailed and visually precise results, even if the rendering process is time-consuming. One of the biggest challenges in real-time rendering, especially in mobile devices and battery-powered systems, is minimizing power consumption. GPUs (Graphics Processing Units) play a critical role in rendering, but they are known for being power-hungry due to the complex calculations involved in processing high-quality 3D scenes. This power usage directly affects battery life and system efficiency, making power optimization a top priority for developers. In this field, research efforts focus on reducing computational overhead while maintaining visual fidelity. The objective is to optimize rendering algorithms and hardware configurations to produce high-quality scenes without excessively increasing power consumption or compromising system performance. As 3D models grow more complex, it becomes harder to maintain a balance between rendering quality and power efficiency, especially when real-time applications require natural, dynamic responses from the GPU. There is a significant relationship between rendering performance and power consumption. As performance improves, typically through higher frame rates and more detailed scenes, power usage increases substantially. Real-time applications, including 3D gaming, simulations, and interactive virtual environments, tend to consume large amounts of power, reducing battery life and impacting other system components. Several factors play a role in evaluating power consumption in rendering applications, including rendering configurations, power efficiency, application characterization, and throughput. These metrics help developers understand the interplay between GPU workloads and power consumption. Although advanced GPUs are equipped with more energy-efficient technologies, the challenge remains in finding an optimal balance between performance and power usage.

The goal of this paper is to render high-quality, photorealistic 3D models as quickly as possible, while staying within a defined power budget. Achieving this objective requires a thorough understanding of GPU architecture, rendering techniques, and power optimization strategies. During the rendering process, the GPU accelerates tasks such as determining the correct colors for each pixel on the screen, after loading each 3D scene into the rendering pipeline. This pipeline comprises several stages, each contributing to the final output. Key steps in the rendering pipeline include handling geometrical complexity, which involves processing 3D models made up of vertices, polygons, and surfaces. Texture resolution is another important factor, as it directly affects the visual quality of the final image. Higher texture resolutions require more memory and processing power, leading to increased power consumption. Other parameters, such as viewport settings (which determine the visible area on the screen), lighting (how objects are illuminated), and shading (how light interacts with objects), also influence the quality of the rendered scene. The complexity of handling these parameters can significantly increase the overall rendering challenge. Modern GPUs adapt dynamically during runtime to optimize performance while adhering to power limitations. One technique commonly used is dynamic voltage and frequency scaling (DVFS), which adjusts the GPU's voltage and clock speed based on current workloads. This allows the system to save power when full performance isn't needed. For instance, when rendering a less demanding scene, the GPU can lower its clock speed, reducing power

consumption without sacrificing visual quality. Conversely, when more computational power is required, the GPU can increase its clock speed to maintain performance. Both GPU core frequency and voltage are affected by these power-saving strategies. As the available bandwidth for processing changes, the GPU adjusts accordingly to maintain optimal performance within the given power constraints. Power limitations are especially critical for mobile and portable devices, where battery life is a key factor. In these situations, the GPU must strike a balance between maintaining high performance and conserving power to extend battery life and ensure efficiency. As GPU technology and rendering techniques continue to advance, new opportunities for power optimization emerge. Power-aware rendering models can help systems adjust dynamically to application demands, reducing unnecessary power usage while still delivering high-quality visuals. Techniques such as adaptive shading, selective rendering, and efficient memory management contribute to minimizing power consumption in 3D rendering. In conclusion, the ongoing challenge in real-time 3D rendering is to deliver high-quality visuals while efficiently managing power consumption. By employing techniques such as dynamic voltage and frequency scaling, shader optimization, and intelligent workload distribution, developers can enhance the efficiency of GPU rendering. This balance between performance and power usage is essential for the future of interactive 3D applications, particularly as the demand for high-quality visuals continues to grow.

In the field of computer graphics, rendering 3D models in real time is a complex process that demands a lot of computational power, especially from the Graphics Processing Unit (GPU). As the demand for high-quality images in 3D games and applications grows, so does the need to find ways to reduce power consumption without compromising on the visual quality or performance. The goal of this research is to develop a strategy that optimizes GPU settings to save power during real-time 3D rendering, while still producing the best possible image quality.

The main focus of this research is to create a method for power savings that doesn't negatively impact the performance or visual quality of 3D models. The approach involves dynamically adjusting different parameters during the rendering process to minimize power consumption. We aim to find the best balance between high-quality rendering and efficient power usage, allowing us to maintain computational efficiency without overusing the GPU or other system resources.

*Key 3D Model Rendering Parameters -*

First, we take a closer look at the main parameters involved in the rendering process. Understanding how these parameters contribute to power consumption helps us identify the areas where we can save power without compromising the visual quality. Some of the key parameters we analyzed include geometric complexity, texture resolution, level of detail, and per-pixel lighting.

**- Geometric Complexity** refers to the number of polygons or vertices used to build the 3D model. More complex models with a higher polygon count require more computational power to render, which in turn increases power consumption.

**- Texture Resolution** plays a significant role in the quality of the rendered image. Higher resolution textures create more detailed and realistic images but demand more memory and processing power from the GPU. This increased load contributes to greater energy usage.

**- Per-Pixel Lighting** is a technique that calculates light interactions at the individual pixel level, which improves the visual realism of the scene. However, this approach is computationally intensive and leads to higher power consumption.

Two additional factors that play a prominent role in power consumption during rendering are frame rate and GPU frequency. The frame rate, or frames per second (fps), refers to the number of frames displayed per second. A higher frame rate ensures smoother motion in the visual output, with 30 fps to 60 fps being the ideal range for most applications. GPU frequency, on the other hand, refers to the clock speed

at which the GPU operates. By controlling these two factors, we can manage power usage more effectively.

In our proposed method, the GPU's frequency is continuously adjusted based on the current frame rate. If the frame rate exceeds 60 fps, the GPU frequency is lowered to save power, since the extra frames are not needed. Similarly, if the frame rate drops below 30 fps, the GPU frequency is increased to ensure smooth performance. This dynamic adjustment of GPU settings helps reduce unnecessary power consumption while maintaining an optimal visual experience.

To better understand power consumption, we conducted a series of experiments using various GPU settings to explore how different configurations affect power usage during rendering. These experiments involved testing different combinations of GPU settings, labeled from C1 to C8, each with specific adjustments to factors like texture resolution, lighting, and geometric complexity. In addition to adjusting the settings, we also used TensorFlow, a widely-used machine learning library, to help manage the GPU's performance during the rendering process. TensorFlow provides important information about the GPU, such as the type and number of GPUs in use, the kind of application running, and the memory usage. TensorFlow also has a built-in pipeline mechanism that allows for parallel data processing, speeding up the rendering process and reducing computational time. TensorFlow's pipeline mechanism is a key part of our power-saving approach. By enabling large-scale parallelism, TensorFlow allows the system to handle more data in less time, which contributes to faster rendering without increasing power consumption. Additionally, TensorFlow helps manage GPU frequency and voltage adjustments, which are crucial to minimizing power usage during the rendering of 3D models.

By integrating TensorFlow into our proposed power-saving strategy, we have created an efficient rendering pipeline that dynamically adjusts the GPU's settings based on real-time data. This dynamic adjustment allows the system to maintain the highest possible image quality while reducing power consumption. During the rendering of 3D models in real time, TensorFlow continuously monitors the GPU's performance and power metrics. If the frame rate is higher than needed (above 60 fps), TensorFlow lowers the GPU frequency to conserve power. If the frame rate drops too low (below 30 fps), TensorFlow raises the frequency to maintain a smooth and high-quality visual output. This constant balancing ensures that the GPU is only using the necessary power for the rendering task at hand.

Our experimental results show that the proposed power-saving strategy effectively reduces GPU power consumption without negatively impacting the quality of the 3D rendered images. When compared to other rendering methods, our approach achieved up to 40% power savings. This improvement is due to the dynamic tuning of GPU parameters, such as frequency and voltage, in response to real-time rendering conditions. In addition to saving power, the integration of TensorFlow's parallel processing capabilities allowed for faster rendering times, further enhancing the system's overall efficiency. By combining power-saving techniques with TensorFlow's ability to manage large-scale data processing, we have created a rendering method that optimizes both performance and power consumption. In conclusion, this research demonstrates that it is possible to achieve high-quality real-time rendering while significantly reducing power consumption. Our proposed power-saving methodology offers an efficient solution for rendering 3D models, with potential applications in areas such as gaming, animation, and virtual reality. Future work in this area could further explore other dynamic power-saving techniques and optimize the rendering pipeline even further, potentially leading to even greater power efficiency and longer battery life for devices running graphics-intensive applications.

The key objectives of this research aim to develop an optimized, power-efficient rendering framework for GPUs.

**Objective 1** focuses on investigating the power consumption of GPUs in real-time, specifically analyzing how rendering configurations and shader parameters (such as geometric complexity, texture

resolution, and lighting) impact energy use during 3D scene rendering. This exploration helps define power-heavy parameters.

**Objective 2** aims to propose a power prediction model that dynamically adjusts based on the complexity of the scene and the bandwidth occupied by different shader calls. This model will predict power usage during rendering to enable better power management.

**Objective 3** involves implementing an algorithm on the GPU accelerator that adjusts the rendering configuration in real-time. This algorithm will ensure an optimal balance between visual quality and power consumption, enabling more efficient 3D rendering while minimizing energy usage.

Together, these objectives create a framework for rendering 3D models that is both high-performance and energy-efficient.

## 2. RELATED WORK

Real-time 3D graphics rendering is a critical component in many multimedia applications today, ranging from gaming to simulations and virtual environments. The demand for high-quality 3D scenes, combined with the need for immediate dynamic responsiveness, drives system developers to continually enhance the performance of rendering processes. However, one significant challenge in this area is managing the balance between performance and power consumption. In existing research [1][3], the role of integrated and discrete Graphics Processing Units (GPUs) is discussed in relation to workload distribution. Poor distribution of rendering tasks across GPU cores often results in excessive power consumption, which is a major issue for computing resources dedicated to 3D model rendering [21]. This is particularly problematic when the load is not evenly shared among the GPU cores, leading to inefficiencies and unnecessary energy usage. To address this issue, our proposed work leverages TensorFlow libraries, which offer a solution by identifying individual GPU cores based on their allocated workloads. TensorFlow enables the efficient distribution of tasks across the cores, ensuring that rendering processes are executed with greater speed and precision. This approach helps in optimizing the rendering workload, thereby reducing power consumption while maintaining high performance. Researchers in this field often argue that achieving low-power consumption designs for processors and accelerators is one of the most significant challenges in modern technology. As advancements in 3D rendering continue, power optimization will remain a critical limitation for future technological developments [6]. In the design process, several metrics must be evaluated to ensure efficiency, including optimized power dissipation, effective throughput, and chip area [4][11]. These are crucial factors that influence the overall performance of GPUs and other accelerators used in 3D rendering tasks. Our work addresses these concerns by proposing an optimized power-aware framework that ensures better workload distribution and power savings while preserving the quality and speed of 3D scene rendering.

The growing need for optimized power consumption in GPUs has become a critical concern in various fields, particularly in the development of algorithms and hardware-integrated systems for rendering 3D graphics. As the demand for higher-quality scenes increases, so do the computational requirements, resulting in longer processing times and higher resource usage. This inevitably leads to significant power consumption, especially as GPUs must frequently access memory, causing substantial bandwidth utilization. A key focus of current research is identifying shader parameters that are closely tied to memory access within GPUs. Efficient management of these parameters can greatly reduce power consumption. The studies referenced in papers [2] and [7] present a power model that addresses this issue, which is based on a traditional stage-based rendering approach that includes batch processing, texture handling, and the use of vertex and fragment shaders. The algorithm developed in this research revolves around two key components: a power consumption prediction model and a real-time quality error estimation strategy. These two mechanisms work in tandem to identify an optimal rendering

configuration that minimizes power usage while remaining transparent to the user. The research emphasizes the importance of managing shader parameters to achieve significant GPU power savings during rendering. One of the main challenges addressed is the time required to compute scene complexity. Precomputation, a necessary step for assessing scene complexity, is time-consuming and requires numerous memory accesses, which can degrade performance, particularly in dynamic scenes. The research explores shader parameters like geometric complexity, texture resolution, per-pixel lighting, and filtering, all of which contribute significantly to power consumption. The objective of the proposed approach is to identify and optimize these parameters, reducing their impact on power usage. Another challenge arises when a scene's complexity exceeds the GPU's memory capacity. In such cases, the GPU may need to access system RAM or even swap data to disk, significantly slowing the rendering process. To address this, the proposed solution integrates the TensorFlow library, which supports massive parallelism. TensorFlow cleans up disk space and improves execution speed, making the rendering process faster and more efficient.

Further research, as discussed in [5], proposes a real-time power budget methodology to balance power consumption with rendered image quality. This approach provides a generalized implementation for various platforms and 3D scene types. The power-saving configurations developed as part of this research are tested across six different scenes, showcasing their ability to reduce power consumption while maintaining rendering quality. The GPU-based rendering pipeline, as illustrated in Fig. 1, includes multiple stages, such as vertex and pixel shaders, which perform complex computations. Larger textures consume more memory and require additional data, directly affecting GPU performance. In the vertex shader stage, each vertex is processed sequentially, and the manipulation algorithm is executed before being passed to subsequent stages in the pipeline. The next stage in the pipeline is the rasterizer, where fragments are generated and sent to the pixel shader. The pixel shader manipulates pixels according to the parameters, and the results are stored in the frame buffer. This pipeline plays a critical role in rendering 3D scenes efficiently, but it also underscores the need to manage power consumption throughout the process. The experiments conducted in [17] provide insights into factors such as vertex sharing, vertex count, texture mapping, and Level of Detail (LOD), all of which impact GPU performance during 3D rendering. These findings offer a foundation for optimizing GPU power consumption and enhancing performance in the rendering of 3D graphics. The research highlights that more complex scenes lead to increased power consumption, primarily due to higher memory access and bandwidth utilization. When more detailed textures and lighting effects are required, power usage rises, which can be problematic for applications that demand real-time rendering, such as video games or virtual reality environments. To address these issues, the proposed solution focuses on optimizing shader parameters and adjusting GPU configurations to reduce power consumption while maintaining high-quality rendering. By modifying factors like geometric complexity and texture resolution, the overall energy usage of the GPU can be minimized without affecting the visual quality of the rendered scene. Additionally, TensorFlow libraries help distribute workloads more efficiently across GPU cores, resulting in faster execution and reduced power consumption. The results of this research indicate that by fine-tuning shader parameters and implementing power-saving techniques, significant energy savings can be achieved in GPU-based rendering. These optimizations can be applied across a variety of 3D scenes and platforms, offering a versatile approach to reducing power consumption in multimedia applications. In conclusion, the growing demand for high-quality 3D graphics rendering has increased the need for power-efficient GPU solutions. By focusing on optimizing shader parameters and refining the rendering process, power consumption can be reduced without compromising image quality or performance. TensorFlow libraries further enhance the efficiency of the rendering process, enabling faster execution and reducing energy usage. This research provides critical insights into the factors influencing GPU power consumption and offers practical strategies for achieving power savings in 3D graphics rendering.
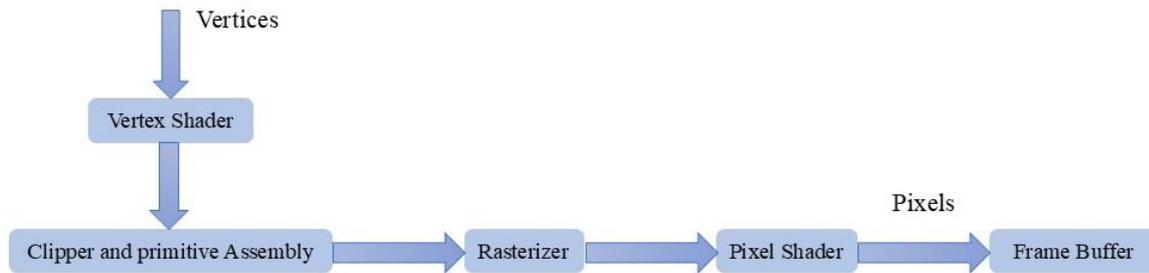
**Fig. 1.** Graphics pipeline with programming

When bandwidth (BW) occupancy is too high, the GPU may underperform due to multiple threads competing for limited space in its memory caches, leading to cache thrashing [7][8]. This inefficiency can cause pixels to be shaded without appearing on screen, resulting in excessive bandwidth usage and power wastage [11]. Overdraw, a pixel rendering issue, exacerbates this by increasing the fill rate (the number of pixels rendered per second) and bandwidth (the data transferred to and from the GPU per second), consuming valuable resources and power even during idle periods [12][13]. In such cases, overdraw creates unnecessary rendering layers, which intensifies both fill rate and bandwidth demand, straining the GPU's capabilities and causing it to consume more power without improving performance. This inefficient resource management leads to heightened power usage while the GPU remains idle or underutilized. Our proposed methodology addresses these challenges by optimizing power consumption through careful analysis of shader parameters. By dynamically adjusting the GPU core frequency based on frame rate monitoring, the rendering process is fine-tuned to maintain optimal performance without unnecessary power drain. This strategy helps prevent bandwidth overuse and reduces overall power consumption while preserving the quality of the rendered scenes.

## 3. PROPOSED METHODOLOGY

The proposed design architecture for a power-aware rendering framework is detailed in Figure 2, and it includes three key components: (a) the TensorFlow library, (b) the Graphics Rendering Parameters Control (GRPC), and (c) the Dynamic Voltage and Frequency Scheduling (DVFS) technique. Each of these components plays a crucial role in optimizing the power consumption of GPUs during 3D model rendering, ensuring that high-quality rendering is maintained while minimizing power usage.
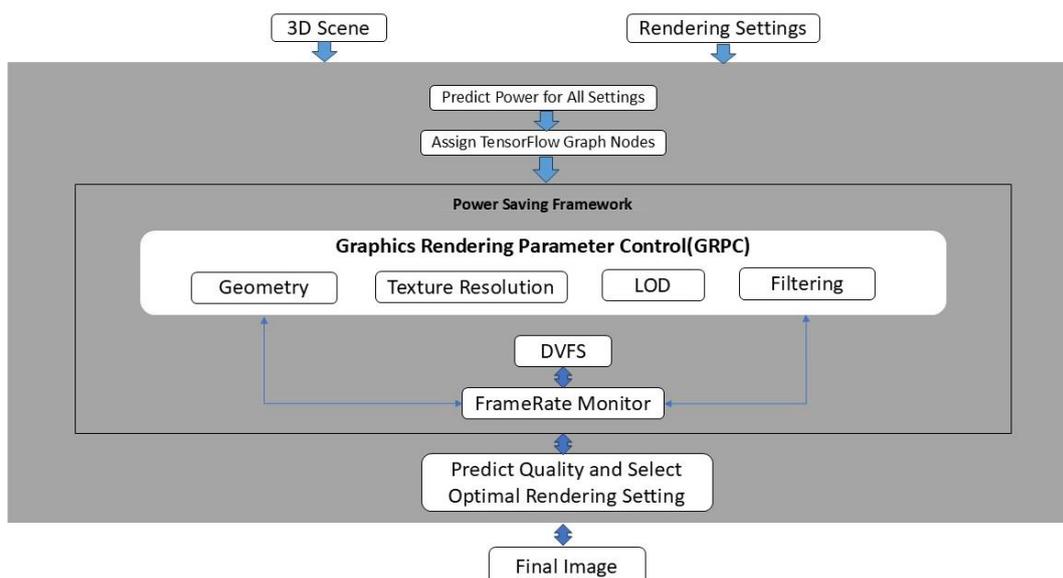
**Fig. 2.** The suggested methodology's architecture

*1. TensorFlow Library and Data Parallelism*

The first component, the TensorFlow library, is primarily responsible for managing device placement for 3D rendering in the initial stage. By assigning specific tasks to each GPU based on their capabilities and availability, TensorFlow helps ensure that the workload is distributed effectively across multiple GPUs. This device placement strategy allows for the efficient handling of 3D graphics tasks by maximizing data parallelism, which significantly speeds up the rendering process. The TensorFlow library plays a pivotal role in defining the frame rates at which the 3D applications will run. By dynamically adjusting these frame rates, TensorFlow ensures that the GPU operates within a range that has been pre-configured by the user, typically between 30 and 60 frames per second (fps). This range is carefully selected to provide an optimal balance between performance and power consumption, with 30 fps representing the lower limit for acceptable performance and 60 fps serving as the upper limit for most real-time 3D applications. TensorFlow's ability to notify the system about the current load on the GPU is another critical function. By monitoring the GPU's workload, TensorFlow can identify periods when the GPU is idle or underutilized. During these periods, it may be possible to reduce the GPU's voltage and frequency, thereby saving power without affecting the overall performance of the rendering process.

This research emphasizes the use of the TensorFlow library for dynamic resource management through a system known as Dynamic Resource Management System (DRMS), which generates graph nodes for computation. In this setup, mathematical operations are executed by the nodes, while the edges connecting them are referred to as tensors, represented by multidimensional arrays. The flexible and portable architecture of TensorFlow enables it to distribute computations across one or more GPUs based on specific requirements. TensorFlow operates as a graph, which is a robust data structure that leverages parallel execution to significantly enhance execution speed. One of its notable advantages is its ability to identify duplicate operations and suboptimal graphs, replacing them with more efficient alternatives to optimize processing time. These features contribute to improved computation efficiency, resulting in reduced execution times and enhanced power savings.

*2. Graphics Rendering Parameters Control (GRPC)*

The second component, the Graphics Rendering Parameters Control (GRPC), serves as an intermediary in the proposed architecture, manipulating key rendering parameters to reduce power consumption during runtime. The GRPC is responsible for adjusting several parameters that affect the rendering of a 3D scene, including geometric complexity, texture resolution, level of detail, and filtering options. These parameters are closely tied to the computational resources required for rendering. By adjusting them dynamically, GRPC ensures that the GPU uses only as much power as needed for the current task, without overloading the system. For example, geometric complexity refers to the number of polygons that make up the 3D models in the scene. By reducing the number of polygons for distant objects or simplifying complex shapes in real-time, the system can reduce the computational load on the GPU. Similarly, texture resolution can be adjusted based on the distance of objects from the camera, with lower resolution textures being used for distant objects. Level of detail (LOD) refers to the process of simplifying a model as it moves further from the camera, ensuring that only essential details are rendered when the object is far away. Filtering parameters, such as anisotropic filtering, can also be dynamically adjusted based on the available GPU resources and the scene's rendering requirements.

The Graphics Rendering Parameters Control (GRPC) serves as an intermediary between the power-optimized frame rate configuration and rendering parameters, ensuring optimal power budget savings. There is potential for the GRPC to be enhanced by incorporating additional parameters, as it is not limited to the current study and can be adapted to integrate various other factors. In practice, the GRPC is implemented as software that independently conducts power analysis of shader parameters, including

geometric complexity, level of detail management, texture size, and filtering selection. The proposed methodology can be further tested with more parameters to improve the predictive accuracy of the framework.

*3. Dynamic Voltage and Frequency Scheduling (DVFS)*

The third and final component is the Dynamic Voltage and Frequency Scheduling (DVFS) technique. DVFS is responsible for adjusting the frequency and voltage of the GPU based on the current frame rate. During the rendering process, the GPU frequency does not need to remain at its maximum value if the current frame rate is lower than the user-defined upper limit. By lowering the frequency when the frame rate is within acceptable limits, DVFS ensures that the GPU does not waste power when it is not needed. Conversely, if the frame rate drops below the minimum acceptable level (30 fps), the GPU frequency is increased to ensure smoother rendering and prevent any visible performance issues. The DVFS technique is applied during the experimentation phase by monitoring the GPU core frequency through the frame monitor. When the frame rate exceeds the upper limit of 60 fps, the GPU core frequency is reduced to the next lower level. This helps conserve power without affecting the user's experience since the visual quality remains the same. On the other hand, if the frame rate falls below 30 fps, the GPU core frequency is increased to maintain an acceptable rendering speed. This dynamic adjustment process ensures that power consumption is minimized while maintaining optimal performance for real-time 3D rendering.

The third and final component of the framework is the Dynamic Voltage and Frequency Scheduling (DVFS) technique. DVFS adjusts the GPU's frequency and voltage according to the current frame rate. During the rendering process, the GPU does not need to operate at its maximum frequency if the frame rate remains below the user-defined upper limit. By lowering the frequency when the frame rate is within acceptable bounds, DVFS prevents unnecessary power consumption. Conversely, if the frame rate dips below the minimum acceptable level of 30 frames per second (fps), the GPU frequency is increased to ensure smoother rendering and avoid noticeable performance degradation. During the experimentation phase, the DVFS technique is implemented by continuously monitoring the GPU core frequency through the frame monitor. When the frame rate exceeds the upper limit of 60 fps, the GPU core frequency is reduced to the next lower level. This adjustment conserves power without compromising the user experience, as visual quality remains consistent. Conversely, if the frame rate drops below 30 fps, the GPU core frequency is raised to maintain a satisfactory rendering speed. This dynamic adjustment process optimizes power consumption while ensuring optimal performance for real-time 3D rendering.

*Rendering Workflow*

The proposed methodology is designed to achieve power savings by dynamically adjusting GPU settings based on the real-time performance of the system. The workflow begins by analyzing the power consumption of the GPU and related computational resources required to render a 3D scene. The TensorFlow library is used to manage device placement and data parallelism, which significantly improves processing speed. During the rendering process, parameters such as geometric complexity, texture resolution, level of detail, and filtering are modified by the GRPC to optimize the power consumption of the GPU. Additionally, the DVFS technique continuously monitors the frame rate and adjusts the GPU core frequency accordingly. If the frame rate is within the user-defined limits, the GPU operates at a lower frequency to save power. If the frame rate falls outside of these limits, the frequency is increased or decreased as necessary to ensure that the 3D application runs smoothly. The key advantage of this approach is that it allows the system to maintain high-quality rendering while reducing power consumption. The dynamic nature of the TensorFlow library, combined with the flexibility of the GRPC and DVFS techniques, enables the system to respond quickly to changes in workload and adjust GPU settings in real-time.
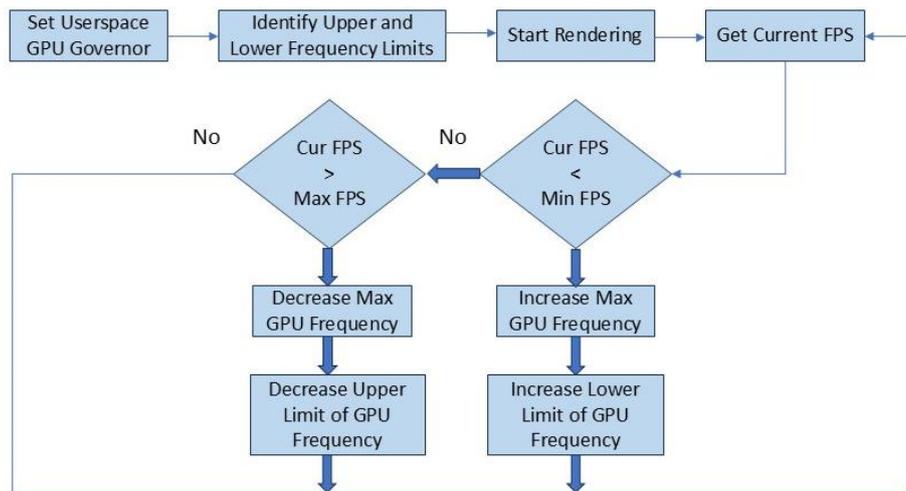
**Fig. 3.** An overview of the algorithm for dynamic voltage and frequency scheduling

In our experimental tests, we set a lower frame rate limit of 30 fps and an upper limit of 60 fps. This range was chosen because it provides acceptable performance for most real-time 3D applications while keeping power consumption to a minimum. The experiments revealed that when the frame rate exceeded 60 fps, the GPU core frequency was reduced, resulting in significant power savings. Conversely, when the frame rate dropped below 30 fps, the GPU core frequency was increased, ensuring that the rendering process remained smooth. Figure 3 shows the dynamic voltage and frequency scheduling algorithm used in the proposed methodology. It provides an overview of how the GPU core frequency is adjusted based on the frame rate. The results of the experiments demonstrate that this approach is highly effective at balancing power consumption and rendering quality in real-time 3D applications. By dynamically adjusting the GPU core frequency, the system can achieve power savings without sacrificing performance. Overall, the proposed methodology offers a comprehensive solution to the problem of high power consumption in GPUs during 3D rendering. By leveraging the TensorFlow library for device placement and data parallelism, the GRPC for rendering parameter control, and the DVFS technique for dynamic voltage and frequency scheduling, the system is able to optimize power usage while maintaining high-quality rendering performance. This approach has the potential to significantly reduce the power consumption of GPUs in real-time 3D applications, making it a valuable tool for developers and system designers.

To minimize the additional overhead related to frequency adjustments, our proposed methodology restricts the GPU's frequency switching to a maximum of once per second. This strategy ensures that the rendering process remains stable while allowing for effective power management. During the real-time rendering of a 3D scene, the system samples the frame rate within a user-defined window, typically consisting of five samples. The average frame rate from these samples is calculated and compared to the current frame rate to determine if any adjustments are necessary. When TensorFlow assesses the GPU's workload, the Dynamic Voltage and Frequency Scheduling (DVFS) technique is activated. This technique adjusts the voltage and frequency levels of the GPU cores based on the identified workload. For example, if the average frame rate indicates that the GPU is underutilized, the DVFS algorithm will reduce the GPU frequency and voltage to conserve energy. Conversely, if the average frame rate falls below the desired threshold—suggesting that the GPU is struggling to keep pace with rendering demands—the system will increase the voltage and frequency to enhance performance. The primary objective of this research is to develop a power-optimized rendering configuration that allows for the efficient rendering of the final image while significantly reducing power consumption. This is accomplished through the careful control and modification of various

rendering parameters, including geometric complexity, texture resolution, level of detail, and filtering techniques. By dynamically adjusting these parameters in real time, the system minimizes the computational resources needed for effective 3D scene rendering. A critical aspect of this approach is recognizing that maintaining a high GPU frequency is not always necessary, particularly if the current frame rate is stable and meets user-defined expectations. If the frame rate remains consistent, the GPU can operate at a lower frequency, thus saving energy without compromising rendering quality. The proposed algorithm emphasizes a rapid, computationally efficient design for frame rendering while adhering to an optimal power budget. It achieves this by simplifying complexities within the 3D scene, primarily through adjustments in geometry and texture details. By alleviating unnecessary computational burdens, the system enhances overall rendering efficiency while ensuring a visually appealing output. This balance between performance and power savings is crucial for real-time applications, particularly in contexts where energy consumption is a significant concern. Ultimately, the implementation of these strategies provides a comprehensive framework for efficient GPU utilization in 3D rendering tasks, showcasing substantial potential for various applications.

$$\textbf{Power} = \textbf{Capacitance} \times \textbf{Voltage}^2 \times \textbf{Frequency}$$

P: the actual power consumed by the GPU
C: capacitance
V: voltage level supplied to GPU
F: clock frequency of the GPU

## 4. IMPLEMENTATION

The rendering of 3D applications relies on several critical parameters, including GPU core and memory frequencies, bandwidth, clock speeds, the number of shaders, and texture attributes. These factors significantly influence both the performance and power consumption of GPUs. The proposed methodology is versatile and can be implemented across various platforms, particularly those supporting NVIDIA GeForce GTX and RTX graphics cards, paired with processors as basic as the Intel Core i5. For the experimentation phase, the setup includes a desktop PC equipped with an Intel Core i7-7700 processor and NVIDIA GeForce RTX 2080 Ti graphics card, along with an NVIDIA Quadro P4000. This configuration boasts approximately 4,000 CUDA cores, a GPU core clock speed ranging from 1,200 MHz to 1,350 MHz, and memory bandwidth spanning from 192 GB/s to 616 GB/s, with a texture fill rate reaching up to 420.2. These specifications enable robust performance, particularly when executing the proposed framework, which adjusts rendering configurations dynamically. The selected frequency configuration plays a pivotal role, as it triggers each process to identify a new optimal rendering configuration for the specified number of frames following each previous setting. This adjustment allows for rapid detection of any alterations in the rendered scene, effectively minimizing the impact of associated computations on overall performance. By dynamically optimizing the rendering parameters based on real-time requirements, the proposed work ensures efficient use of computational resources while maintaining high-quality visuals. This not only enhances the rendering speed but also contributes to significant power savings, making it an effective solution for real-time 3D rendering tasks across a variety of applications. The adaptability of the methodology to different hardware setups further broadens its potential for use in diverse scenarios, including gaming, simulation, and virtual reality environments.

*4.1 Power Measurement*

The NVIDIA Management Library (NVML) API is utilized to accurately assess the power consumption of graphics cards in both laptops and desktop computers. NVML offers a comprehensive framework for monitoring and managing the status of NVIDIA GPU devices, enabling direct evaluation of power

usage alongside the GPU's supporting circuitry. This capability is crucial for understanding the performance and efficiency of the graphics card during demanding tasks like 3D rendering. In the experimental setup, power consumption is averaged over 30 frames to provide a reliable measurement that reflects the GPU's performance under typical operating conditions. The average Thermal Design Power (TDP) for the system is around 270 watts, accompanied by an impressive floating-point performance of 5,304 gigaflops (GFLOPS). This performance level indicates the GPU's capability to effectively handle resource-intensive computational tasks. The Dynamic Voltage and Frequency Scaling (DVFS) mechanism is implemented to further enhance power efficiency by allowing the GPU frequency to be dynamically adjusted every second. This real-time frequency modulation helps maintain the frame rate within the range of 30 to 60 frames per second (fps), ensuring smooth rendering while minimizing power consumption. Through this NVML-configured setup, the methodology effectively demonstrates optimal power usage patterns, balancing performance with efficiency. By continuously monitoring power metrics and adjusting operational parameters, the approach not only enhances the rendering experience but also contributes to sustainable GPU resource utilization across various applications.

## 5. RESULTS AND DISCUSSION

To establish the rendering workload for our experiments, we utilized a complex synthetic 3D scene that accurately depicts a specific scenario typical in 3D gaming, as shown in Figure 4. This scene serves as the basis for our tests, providing a rich and dynamic environment to evaluate the performance and power consumption of various rendering techniques. The 3D scene is meticulously crafted, featuring over 300,000 polygons, which contribute to its complexity and realism. We included 70 different model types, primarily statues, each with three levels of geometric detail. This variety enables comprehensive testing across different rendering configurations.



**Fig. 4.** 3D scene with power profile

The camera path within the 3D scene was predetermined to ensure a consistent flow throughout our experimental tests. This systematic approach allows us to analyze the rendering process effectively, eliminating redundancy in our tests. The geometric models are programmatically rotated to optimize

the rendering engine's calculations, minimizing the workload associated with caching and maximizing efficiency. To validate the proposed methodology, we conducted a series of systematic experimental tests, encompassing around eight distinct tuned settings, as outlined in Table 1. These settings were developed by examining four key parameters: geometric complexity, texture resolution, level of detail (LOD), and filtering features, along with the application of Dynamic Voltage and Frequency Scheduling (DVFS). Each of these parameters significantly influences the power consumption of GPUs during rendering.

**Table 1.** Power consumption of each parameter in the experimental setup separately

| Tuned Setting | Texture Resolution | Filtering | LOD management | DVFS |
|---|---|---|---|---|
| C1 | $2048^2$ | Yes | No | No |
| C2 | $1024^2$ | Yes | No | No |
| C3 | $512^2$ | Yes | No | No |
| C4 | $2048^2$ | No | No | No |
| C5 | $2048^2$ | Yes | Yes | No |
| C6 | $2048^2$ | Yes | No | Yes |
| C7 | $2048^2$ | No | Yes | Yes |
| C8 | $512^2$ | No | Yes | Yes |

For texture resolution, we considered three different configurations: 2048x2048 pixels, 1024x1024 pixels, and 512x512 pixels. Each texture setting corresponds to a specific configuration, allowing us to assess how changes in resolution impact power usage and rendering quality. The geometric complexity is adjusted using a Level of Detail manager, which dynamically alters the detail of the models based on their distance from the camera. This adjustment enhances performance by reducing the rendering load without compromising visual fidelity. All configurations were loaded programmatically at the beginning of the rendering process, ensuring that the correct settings were consistently applied across each test. We utilized standard filtering techniques throughout the experiments to enhance image quality. This filtering process requires additional computations over the fragment shader but ultimately improves visual output by providing more detail and smoother transitions in textures. To evaluate the power consumption associated with each configuration, we established a standard reference condition known as C1. This condition is assumed to represent the maximum expected power consumption among all the tuned settings outlined in Table 1. In this scenario, we aimed to achieve the highest possible quality while considering all available factors. Each tuned setting was tested with frame rates set between 30 fps and 60 fps, carefully balancing the output close to user-defined limits. This flexibility ensures that we can simulate real-world scenarios where users have varying performance expectations. Our observations revealed that the lowest power consumption values were associated with settings C7 and C8. These configurations are enabled with DVFS, which dynamically adjusts the GPU's voltage and frequency based on rendering demands. As a result, they achieve substantial power savings while maintaining an acceptable frame rate. Conversely, configuration C4 yielded a significantly high frame rate; however, this setting did not utilize filtering, LOD adjustments, or DVFS, indicating that while it performs well in terms of speed, it does not optimize power consumption effectively. From the data presented in Table 1, it is clear that the reference condition C1 is the most power-intensive setting. This condition serves as a benchmark against which we can measure the

efficiency of the other configurations. The power consumption results for C7 and C8 demonstrate that enabling DVFS leads to significant power savings without sacrificing performance. This finding highlights the importance of adaptive power management techniques in modern rendering frameworks. Furthermore, our experiments emphasized the trade-offs between performance and power consumption. While achieving high frame rates is critical for user satisfaction, it is equally important to consider the energy costs associated with different rendering settings. The ability to dynamically adjust the GPU's power usage based on rendering demands provides a viable solution to mitigate these costs while maintaining high-quality visuals. As we analyzed the results, we recognized that integrating LOD and filtering techniques could further enhance our framework's efficiency. By intelligently managing the level of detail and applying appropriate filtering methods, we can optimize power consumption without sacrificing the quality of the rendered scenes. This optimization is particularly beneficial in applications where battery life is a concern, such as mobile gaming or laptop usage. In conclusion, our research demonstrates that using a complex 3D synthetic scene enables effective evaluation of various rendering techniques and their impact on power consumption. The careful selection of parameters, such as geometric complexity, texture resolution, and filtering features, is crucial for optimizing the rendering process. By employing a reference condition and systematically analyzing different tuned settings, we have established a clear understanding of how each configuration affects GPU power usage. The implementation of DVFS is particularly noteworthy, as it offers a powerful mechanism to achieve significant power savings without compromising rendering performance. As the gaming and graphics industries continue to evolve, incorporating such adaptive techniques will be essential for developing efficient rendering solutions that meet both performance expectations and power constraints. Future research can build on these findings by exploring additional parameters and optimizing rendering workflows to further enhance the efficiency and quality of 3D rendering in real-time applications.

*5.1. Textures*

Tuned settings C1 through C3 were analyzed to assess the power consumed due to the increased workload from changing the texture resolution of the test scene. When varying the texture resolution in the tested scene, we examined the power consumption for these tuned settings. However, no significant benefits were observed for tuned settings C2 and C3. Specifically, reducing the texture resolution to 1024x1024 pixels resulted in a marginal improvement in GPU power consumption, yielding a negligible reduction of just 2%.

*5.2. Filtering*

To assess the power consumption associated with filtering on the GPU, we define tuned setting C4. In this configuration, filtering is turned off, resulting in a minor reduction in image quality. All other parameters remain unchanged from the reference setting, C1. Disabling the filtering process allows the GPU to save around 35% of energy in comparison to the reference condition setting.

*5.3. Level-of-Detail (LOD)*

The impact of geometric complexity is evaluated using tuned setting C5, where Level of Detail (LOD) management is activated. With LOD management in place, the GPU achieves energy savings of approximately 33% compared to the reference tuned setting C1, resulting in an overall performance improvement of 15%.

*5.4. Dynamic Voltage and Frequency Scheduling*

The implementation of the Dynamic Voltage and Frequency Scaling (DVFS) mechanism is examined in tuned setting C6. By enabling DVFS, the GPU frequency is dynamically adjusted every second during runtime, ensuring that the frame rate remains between 30 and 60 fps. This approach results in a significant reduction in GPU power consumption, achieving an improvement of approximately 74%. Under tuned setting C6, the effects of DVFS on power consumption are evident, with the GPU

frequency being managed dynamically and maintaining the specified frame rate range. Overall, more than 50% improvement in GPU power consumption is observed in this DVFS state.

*5.5. Overall Observations*

Figure 5 illustrates the frame rates achieved across all experimental tuned settings, revealing performance differences resulting from the applied optimizations. As indicated in Table 1, settings C7 and C8 show considerable power savings, particularly when all optimization configurations are enabled and texture resolutions of 512x512 and 2048x2048 pixels are used. Specifically, under condition C7, a notable 50% power saving for the GPU is recorded. The majority of these savings are attributed to the Dynamic Voltage and Frequency Scaling (DVFS) technique, while a small portion comes from Level of Detail (LOD) management and filtering. This indicates that enabling DVFS while disabling LOD and filtering produces the most significant GPU power savings. In terms of texture resolution, even with the reduction to 512x512 pixels in tuned setting C8, the expected power savings were not as significant. This suggests that simply lowering texture resolution has limitations in its impact on power consumption, indicating that other factors play a crucial role in overall GPU performance. Additionally, the proposed methodology investigates the effects of per-pixel lighting and shading calculations, commonly referred to as vertex-based calculations, during the rendering of 3D scenes on the GPU. The per-pixel lighting scenario is particularly computationally demanding, resulting in higher power requirements during the rendering process. Consequently, both fill rate and bandwidth can become limited resources, further contributing to power usage. To address these challenges, the proposed methodology includes an analysis of per-pixel lighting in relation to GPU power dissipation. By reducing compute time and eliminating duplicate pixels, this approach decreases switching activity and optimizes the processing of relevant triangles, ultimately improving power efficiency.
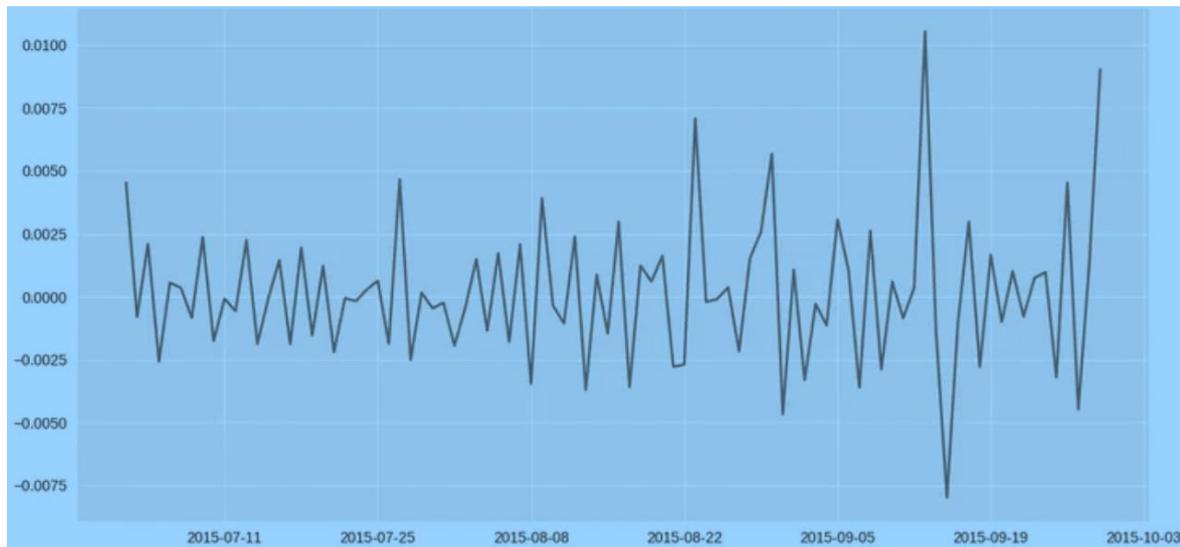


**Fig. 5.** Achieved frame rates across all experimentally adjusted parameters

Table 2 shows enhancements in rendering performance, with an increase from an existing rate of 1.235× to 1.548×. This reflects significant overall power savings, as the average power consumption is reduced compared to previous studies, such as [5] and [7], which typically require extensive precomputation and complex implementation mechanisms for rendering scenes. The proposed methodology demonstrates a remarkable improvement in power efficiency, ranging from 1.812× to 4.826×, depending on the size of the triangles being rendered. Smaller triangles need less precision, correlating with the higher power savings achieved with this methodology. Even when larger triangles, up to a size of 255 pixels in both width and height, are rendered, a notable improvement of 1.876× in power

efficiency is still observed. Further analysis reveals that when the GPU operates at a frequency of 250 MHz, the minimal average power consumed is around 142.0 mW. Rendering a single pixel requires an average of 600.8 picojoules (pJ), leading to an overall average efficiency of the rendering pipeline of approximately 2.143 million pixels per millijoule (MPixels/mJ). Implementing Dynamic Voltage and Frequency Scheduling (DVFS) can be especially advantageous during times when performance configurations do not demand maximum processing power. This adaptive approach not only enhances power efficiency but also contributes to the overall sustainability of the rendering process in 3D applications.

**Table 2.** Comparison of power using and without using the suggested methods

| Triangle Size | Without proposed methodology | | | | Without proposed methodology | | | | Efficiency improvement |
|---|---|---|---|---|---|---|---|---|---|
| | Bus Bandwidth Occupancy | Avg Power | Pixel Energy | Efficiency | Bus Bandwidth Utilization | Average Power | Pixel Energy | Efficiency | |
| (Width in pixels) | (%) | (pJ) | (MPixels/mJ) | (mW) | (%) | (pJ) | (MPixels/mJ) | (mW) | |
| 255(Large) | 93 | 248.343 | 518.989 | 1.045 | 94 | 137.534 | 282.454 | 3.888 | 1.812x |
| 171 | 90 | 256.867 | 570.542 | 1.812 | 92 | 137.623 | 290.214 | 3.675 | 1.856x |
| 128 | 87 | 268.443 | 614.324 | 1.753 | 89 | 137.745 | 300.099 | 3.402 | 1.948x |
| 102 | 83 | 274.676 | 617.136 | 1.612 | 85 | 138.098 | 310.785 | 3.217 | 1.892 |
| 85 | 80 | 302.554 | 700.087 | 1.587 | 82 | 138.435 | 320.498 | 3.056 | 2.300 |
| 51(Middle) | 63 | 317.825 | 864.231 | 1.291 | 66 | 139.786 | 390.213 | 2.975 | 2.203 |
| 25 | 48 | 340.774 | 1303.645 | 0.843 | 52 | 140.897 | 480.899 | 2.587 | 2.8 |
| 10 | 30 | 377.432 | 2260.065 | 0.554 | 35 | 142.567 | 760.012 | 2.067 | 2.998 |
| 5 | 20 | 473.988 | 3622.978 | 0.398 | 24 | 145.544 | 1080.034 | 1.560 | 3.401 |
| 2(small) | 14 | 324.433 | 6417.825 | 0.256 | 15 | 150.091 | 1350.067 | 1.989 | 4.826 |
| Average | - | 318.533 | 1748.982 | 1.115 | - | 140.832 | 556.527 | 2.841 | |

These findings showcase the effectiveness of the proposed methodology in optimizing GPU power consumption while maintaining rendering performance. By focusing on critical aspects of power management, including DVFS and the handling of geometric complexity, texture resolutions, and lighting calculations, this research successfully addresses the inherent challenges of rendering in 3D environments. The results indicate that further exploration into additional parameters, such as fill rates and bandwidth occupancy, could provide even deeper insights and improvements in future research. A comprehensive understanding of the interactions among these various factors is essential for enhancing the efficiency of GPU rendering processes, making significant strides toward more sustainable graphics performance in both gaming and professional applications.

## 6. CONCLUSION AND FUTURE WORK

This paper discusses the results of experimental tests aimed at optimizing power consumption in 3D rendering. Our proposed methodology integrates three key components: the TensorFlow library, Graphics Rendering Parameter Control (GRPC), and Dynamic Voltage and Frequency Scheduling (DVFS). These elements work together to improve power efficiency while maintaining the visual quality of rendered 3D scenes. The heart of our experimental framework focuses on essential

parameters of 3D scenes, such as geometry resolution, filtering, level of detail (LOD), texture resolution, and texture size. Each of these parameters plays a significant role in the rendering process, and we carefully evaluated their impact on both power usage and the quality of the images produced. By fine-tuning these parameters, we ensured that we could consistently maintain a frame rate of 30 frames per second (fps). This frame rate allows users to explore the rendered 3D scenes in detail without interruptions. An important aspect of our work was to keep the scene resolution high. Lowering the resolution can lead to noticeable quality issues, so we focused on maintaining high-quality scenes while also optimizing power consumption. Our results indicate that it is possible to achieve this balance effectively. The experimental findings show notable power savings when examining each parameter independently. More impressively, when we combine all the parameters, we achieved overall power savings of up to 40%. This substantial reduction in power consumption demonstrates the effectiveness of our proposed methodology in optimizing GPU performance while still delivering high-quality rendering results. The framework we developed has several key advantages that make it suitable for real-time 3D rendering applications. One significant benefit is that it does not rely on extensive precomputation, which is common in traditional rendering techniques. Instead, our framework can smoothly handle dynamic scenes without any decrease in frame rates. This ability to adapt quickly is crucial for interactive applications, where user experience depends heavily on smooth visuals. Additionally, the methodology excels in performing per-pixel lighting analysis. By integrating this analysis into our framework, we were able to significantly reduce power consumption while ensuring that the quality of the rendered scenes remains high. This capability is essential for applications where realistic lighting effects are important for enhancing the overall visual experience.

The power settings within our framework are individually configured for each demonstration, allowing for precise monitoring of optimal rendering configurations. This flexibility helps us understand how different settings affect power usage and performance, contributing to our overall goal of power optimization. The use of the TensorFlow library is also vital to the success of our methodology. TensorFlow enhances the parallelism of computations, which leads to faster rendering processes. Its ability to manage and distribute computational tasks across multiple processing units increases overall system efficiency. By using TensorFlow's dynamic computation graph, our framework can adapt to varying workloads in real-time, optimizing both performance and power consumption effectively. Moreover, the implementation of DVFS technology is critical for maximizing power savings. This technology allows us to dynamically adjust the GPU's frequency and voltage based on real-time frame rate measurements. By lowering the GPU's power usage when full performance is not required, we can conserve energy while still delivering high-quality graphics. This dynamic adjustment is particularly beneficial in scenarios where frame rates fluctuate, allowing the GPU to operate at lower power levels without sacrificing visual fidelity. While our research demonstrates the effectiveness of the proposed framework, there remains significant potential for further development. The predefined camera paths, scene exploration techniques, and various measurements of quality under different configurations show the promise of our approach. However, future work can expand on this foundation by exploring additional parameters related to complex 3D models. For example, analyzing bandwidth occupancy and fill rates among pixels could provide deeper insights into the computational demands of different rendering tasks. These factors often require substantial power during the rendering process and are worth investigating alongside our proposed methodology. In addition to these considerations, we can also explore resource balancing techniques that focus on managing the workloads of both the CPU and GPU. This approach could yield valuable insights into how to allocate resources dynamically, ensuring that both processing units work together effectively to deliver optimal performance and efficiency. In conclusion, this paper highlights the importance of our proposed methodology for power optimization in 3D rendering applications. By integrating the TensorFlow library, GRPC, and DVFS techniques, we have developed a framework that reduces power consumption while maintaining high-quality rendering. Our experimental results demonstrate a promising 40% reduction in power usage, showcasing the potential of our approach in real-world applications. By continuing to explore additional parameters and resource balancing techniques, we can further enhance the framework's capabilities and contribute to advancements in power-efficient 3D rendering solutions.

**REFERENCES**

[1] Abraham, F., Celes, W., Cerqueira, R., Campos, J.L.: A load-balancing strategy for sort-first distributed rendering. In: Proceedings. 17th Brazilian Symposium on Computer Graphics and Image Processing, pp. 292–299 (2004)

[2] Bhaniramka, P., Robert, P.C., Eilemann, S.: OpenGL multipipe SDK: a toolkit for scalable parallel rendering. In: VIS 05. IEEE Visualization, pp. 119–126. IEEE (2005)

[3] Yunjin Zhang, Marta Ortin, Victor Arellano, Rui Wang, Diego Gutierrez, Hujun Bao, On-the-Fly Power-Aware Rendering, Eurographics Symposium on Rendering 37 (2018) Number 4.

[4] Xiao Lei, Vetria L. Byrd, Real-Time Rendering With Heterogeneous Gpus.", in: International Conferences Computer Graphics, Visualization, Computer Vision and Image Processing, 2020, ISBN: 978-989-8704-21-4 © 2020.

[5] J. Ma et al., "An analytical framework for estimating scale-out and scale-up power efficiency of heterogeneous manycores," IEEE Transactions on Computers, vol. 65, no. 2, pp. 367–381, 2016.

[6] M. Cavalcante et al., "Ara: A 1-GHz+ scalable and energy-efficient RISC-V vector processor with multiprecision floating-point support in 22-nm FD-SOI," IEEE Transactions on Very Large Scale Integration, vol. 28, no. 2, pp. 530–543, 2020.

[7] Jaekyu Lee, Hyesoon Kim (2012) TAP: A TLP-aware cache management policy for a CPU-GPU heterogeneous architecture IEEE International Symposium on High-Performance Comp Architecture

[8] An-Chow Lai, Cem Fide, B. Falsafi (2001) Dead-block prediction & dead-block correlating prefetchers Proceedings 28th Annual International Symposium on Computer Architecture

[9] M. B. Taylor, "Is dark silicon useful? Harnessing the four horsemen of the coming dark silicon apocalypse", Proc. 49th Annu. Des. Autom. Conf., pp. 1131-1136, 2012.

[10] J. Guo, J. Meng, Q. Yi, V. Morozov and K. Kumaran, "Analytically modeling application execution for software-hardware co-design", Proc. IEEE 28th Int. Symp. Parallel Distrib. Process, pp. 468-477, 2014.

[11] William V Baxter, Avneesh Sud, Naga K Govindaraju and Dinesh Manocha, "Gigawalk: Interactive walkthrough of complex environments" in Rendering Techniques, pp. 203-214, 2002.

[12] Enrico Gobbetti and Fabio Marton, "Far voxels: a multiresolution framework for interactive rendering of huge complex 3D models on commodity graphics platforms", ACM Transactions on Graphics (TOG), vol. 24, no. 3, pp. 878-885, 2005.

[13] Abe Stephens, Solomon Boulos, James Bigler, Ingo Wald and Steven Parker, "An application of scalable massive model interaction using shared-memory systems", Proceedings of the 6th Eurographics conference on Parallel Graphics and Visualization, pp. 19-27, 2006.

[14] Lukefahr, S. Padmanabha, R. Das, R. Dreslinski, T. F. Wenisch and S. Mahlke, "Heterogeneous microarchitectures trump voltage scaling for low-power cores", Proc. Parallel Archit. Compilation Techn., pp. 237-250, 2014.

[15] S. Sarma, T. Muck, L. A. Bathen, N. Dutt and A. Nicolau, "SmartBalance: A sensing-driven linux load balancer for energy efficiency of heterogeneous MPSoCs", Proc. Des. Autom. Conf., pp. 1-6, 2015.

[16] J. Kong, S.W. Chung and K. Skadron, "Recent thermal management techniques for microprocessors", *ACM Comput. Survey*, vol. 44, no. 3, 2012.