

Carbon Aware Serverless System - Group 6

Jaykumar Patel
patel.jay4802@utexas.edu
jnp2369

Afnan Mir
afnanmir@utexas.edu
amm23523

Nidhi Dubagunta
nidhi.dubagunta@utexas.edu
nsd632

Abstract

Serverless computing is a method that enables developers to build and deploy applications without having to provision resources or manage infrastructure. However, existing serverless systems do not consider the carbon footprint of function invocations when allocating or configuring resources, such as memory and vCPU. In this paper, we propose a carbon-aware serverless system that focuses on minimizing the carbon emissions of function invocations while continuing to meet service level objectives (SLOs) by finding an optimal configuration for vCPU count and operational frequency of the vCPU cores. We first perform a measurement study by running a variety of workloads with different resource configurations to understand the relationship between resource usage, energy consumption, and performance. We also collect measurements for energy consumption during different stages in the serverless function invocation lifecycle. Leveraging these insights, we aim to incorporate a Bayesian Optimization algorithm into OpenWhisk, an open-source serverless cloud platform, to build a serverless system that performs carbon-aware dynamic frequency scaling (DFS) and vCPU allocation.

1 Introduction

The rise of cloud computing has transformed the space of digital infrastructure, offering higher amounts of scalability, flexibility, and cost-effectiveness compared to on-premise solutions. Within cloud computing, the concept of serverless computing emerged, which further increased efficiency in application development by abstracting the complexities of infrastructure management and resource provisioning.

In a serverless system, when a function is invoked, the system checks for a pre-existing warmed container to host the function. If no warmed container is available, a cold start is initiated, during which the runtime environment is set up, necessary dependencies are installed, and the function code is loaded. After the function is executed, the con-

tainer can remain idle for a period of time for future invocations before it is killed.

Several existing serverless systems, such as AWS and GCP, couple memory and CPU allocation, assigning a CPU share proportional to the memory requested by the user, or using a preset resource allocation option [2]. Serverless systems today also do not incorporate dynamic frequency scaling (DFS), in which the operating frequency of the vCPU cores can be dynamically modified. These policies used for resource allocation are carbon-agnostic, meaning they do not consider the carbon emissions of the function execution and idle containers.

In this paper, we propose a serverless framework that provisions resources based on workload characteristics to minimize the carbon emissions associated with the function invocation, while continuing to meet SLOs, such as total runtime. We introduce the frequency at which vCPUs run as a new resource to allocate per invocation. Our system will optimize the number of vCPUs allocated, as well as the frequency at which the vCPU cores operate.

2 Motivation

According to a study by Microsoft and the Green Software Foundation, data centers and data transmission networks account for almost 1% of energy-related greenhouse gas emissions, which contribute to rising global temperatures and climate change [3]. By adopting carbon-aware computing, these data centers can play a significant role in decarbonization efforts.

The boxplot in Figure 1, shows the variability in energy consumption (in joules (J)) for different functions with a fixed input, while changing the amount of vCPUs allocated and the frequency at which the cores run. The large range of energy consumption for the same function due to changes in vCPU count and frequency underscores the importance of understanding the relationship between these factors and energy usage. This motivates further study into optimizing vCPU allocation and frequency at a function-level to minimize energy consumption and therefore carbon footprint.

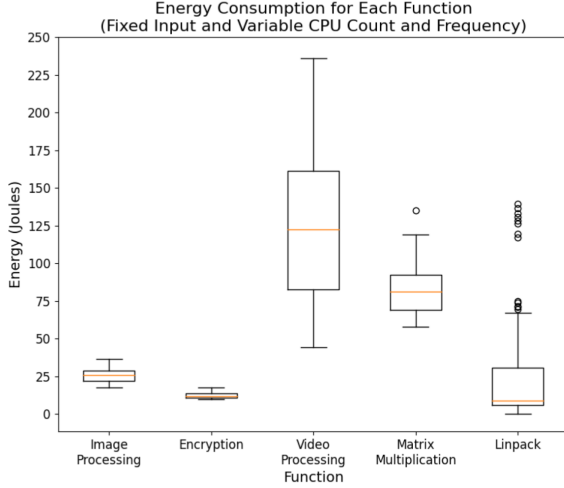


Figure 1: Boxplot showing energy variability across functions with fixed input, while changing vCPU count and frequency.

3 Related Work

CherryPick [1] is a framework that was developed to find optimal or near-optimal cloud configurations (vCPU and memory) for repetitive big data analytics jobs, addressing challenges of system cost, performance, and adaptivity. CherryPick leverages Bayesian Optimization (BO) to build performance models that are able to determine optimal or near-optimal configurations that satisfy performance constraints with a relatively low number of samples.

Aquatope [8] also utilizes Bayesian models to optimize for resource allocation for serverless function invocations. This framework also utilizes BO to learn the mapping from resource configuration to performance and cost, while considering noise and uncertainty in the cloud environment. Aquatope also uses batch sampling to reduce the cost of exploration when finding a resource configuration.

Our work aims to draw from the ML techniques utilized in CherryPick and Aquatope. We plan to use BO to find the near-optimal configuration of vCPU and operational frequency with the objective of minimizing carbon footprint for serverless function invocations. We will also set performance constraints to ensure user-specified SLOs are met.

In the space of carbon-aware computing, GreenCourier [4] is a scheduling framework that enables runtime scheduling of serverless functions across geographically distributed regions based on their carbon efficiencies. This work focuses on minimizing carbon footprint based on the location at which a function is run, whereas our work aims to optimize the host environment through proper resource allocation and configuration to minimize carbon emissions using function-level characteristics.

4 Measuring Carbon Emission

Before conducting the measurement study, we need to define how to measure carbon emissions. According to Gupta et al. [5], the carbon footprint of a server can be calculated using the following formula where OP_{CF} is the operational carbon footprint, CI_{use} is the carbon intensity of the energy source, and $Energy$ is the operational energy:

$$OP_{CF} = CI_{use} \times Energy \quad (1)$$

The carbon intensity refers to how *clean* the energy source is. For example, coal has a higher carbon intensity than solar. We assume that the energy source of the serverless system is constant, and therefore, the carbon intensity is constant. We can then use the operational energy as a proxy for the carbon footprint.

4.1 Measuring Energy

We first need a way to measure the energy consumption of a serverless function invocation. EnergAt [6] is a tool that measures fine-grained energy consumption at the thread-level while taking into account NUMA effects. An EnergAt process continuously samples a target container, where each sample returns energy readings over an interval of at least 50ms. Therefore, EnergAt will get continuous back-to-back samples of energy readings over time. The energy consumption of a target container is the sum of the energy readings over the duration of the function invocation. One limitation of EnergAt is that every target container for which we want to measure energy consumption must be instrumented with its own EnergAt process. This is not feasible as the number of EnergAt processes and resource contention scales with the number of target containers.

To mitigate this issue, we created EnergyLiteDaemon. EnergyLiteDaemon is a daemon that samples multiple target containers in a round-robin fashion, where each sample returns energy readings over an interval of at least 50ms. This allows us to measure the energy consumption of multiple target containers with a single, lightweight, process. Due to the round-robin fashion, however, the sampling frequency per target container decreases as a function of the number of target containers. If there are multiple containers, the samples are no longer back-to-back – there may be some time delay between two samples. To solve this, we complete the following steps to interpolate the total energy consumption of a target container:

Power Calculation: We compute the average power for every sample by dividing the energy consumed over the interval of the sample by the duration of that interval. This allows us to get a power-time scatter plot that contains n points, where n is the number of samples.

Variable	Values
Function Type	Float Matrix Multiplication, Image Processing, Video Processing, Encryption, Linpack
vCPU Count	1, 2, 3, 5, 7, 10, 13, 16, 19, 22, 25, 28, 31
vCPU Frequency (GHz)	1.0 to 2.4 in 0.2 increments

Table 1: Measurement Study Configurations – Energy Variation with vCPU and Frequency Allocation

Variable	Values
Image	main-python, video-ow, mobilenet-ow, sentiment-ow, audio-ow
vCPU Count	1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 12, 14, 16, 18, 20, 22, 24, 26, 28, 30, 32
Memory (MB)	512 to 10240 in 512 increments

Table 2: Measurement Study Configurations – Energy Variation with vCPU and Frequency Allocation

Trapezoidal Integration: Given a power-time scatter-plot, we compute the total energy consumption of a target container using trapezoidal integration. All of the points in the power-time scatter plot are connected by straight lines, and the area under the curve is calculated. This area is the total energy consumption of the target container.

A toy example of energy measurement using EnergAt and EnergyLiteDaemon is shown in Appendix A.1.

Due to EnergyLiteDaemon’s round-robin sampling approach, which reduces its overhead, it may not capture all energy readings for a container. Even with trapezoidal integration, the energy consumption will not be exact. To determine the accuracy of EnergyLiteDaemon, we ran 20 stress docker containers, with varying CPU usage and duration. Figure 2 shows the energy consumption of the containers as measured by EnergAt and EnergyLiteDaemon. The absolute percentage error of EnergyLiteDaemon is less than 15%, which is acceptable for our purposes.

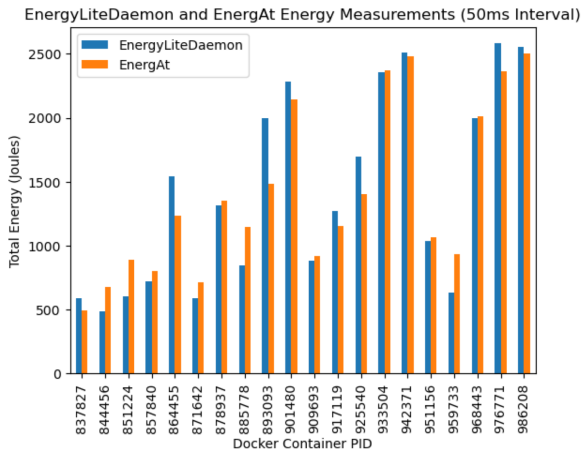


Figure 2: Energy measurements of EnergAt and EnergyLiteDaemon for 20 stress docker containers.

5 Measurement Study

We conduct two measurement studies to answer the following questions: **1. How does energy vary with changes in vCPU and frequency allocation?** **2. What is the energy consumption over the lifecycle of a serverless function invocation?**

To perform the measurement study, we use OpenWhisk, which is an open-source serverless platform, to run function invocations. We use EnergyLiteDaemon to measure the energy consumption of the function invocations.

5.1 Energy Variation with vCPU and Frequency Allocation

A variety of workloads were run on different vCPU and frequency configurations. We collected the vCPU utilization, energy consumption, and performance in terms of duration. The resource configurations are shown in Table 2.

We then analyzed the data to understand how energy consumption varies with changes in vCPU and frequency allocation. For float matrix multiplication, energy usage initially decreases and then flat-lines as the vCPU allocation increases. This is because matrix multiplication is highly parallelizable and has high vCPU utilization, leading to a decrease in duration and energy usage. For image processing, energy usage has no apparent correlation with vCPU allocation, likely because most of the image processing tasks are non-parallelizable. For encryption, energy usage, duration, and vCPU utilization do not change with vCPU allocation because encryption is also non-parallelizable and does not benefit from additional vCPUs.

For both types of functions, energy consumption as a function of frequency displays a convex curve, where energy consumption decreases with frequency until a certain point (which varies between different functions and function inputs), and then begins to increase again. This is because the energy consumption per time unit of a vCPU is

proportional to the frequency at which it runs, but the overall duration of the function is inversely proportional to the frequency at which the vCPU runs. More details on the results can be found in Appendix A.2.

These trends support the idea that we can use machine learning to learn the relationship between vCPU, frequency, and energy consumption, and use this model to minimize energy consumption while meeting SLOs.

5.2 Energy Consumption Over the Lifecycle of a Serverless Function Invocation

Additionally, we ran a variety of docker images, each of a different size, with a variety of vCPU and memory configurations to investigate carbon emissions of different stages of the container lifecycle. For each (image, vCPU, memory) configuration, we measured the energy consumption at different stages of the function invocation lifecycle. More specifically, we measured the energy consumed during the container spin-up, the idle time, and the container spin-down.

After collecting and analyzing the data, we first found that for each stage of the lifecycle we measured, when changing the vCPU and memory configurations, the energy consumption did not vary significantly. This was true for all image sizes. This is likely due to the fact that these resources will not actually be allocated to the container until the compute inside the container requires them. We did notice that if we varied the image size, the energy consumption of the container spin-up process significantly increased as the image size increased. However, it is worth noting that most of the energy consumption of the spin-up is accredited to the ‘docker pull’ process, which would only be a one-time operation. The ‘docker run’ does not consume a significant amount of energy and does not show any statistical correlation. For the idle time, we noticed a strange negative correlation trend with the image size, and for the spin-down, we noticed a somewhat positive correlation with a few outlier points. Since we only used a few image sizes, we cannot definitively conclude that these trends will extrapolate.

Due to our findings of little to no correlations, we decided to not pursue a carbon-aware keep-alive policy further. More details on the study can be found in Appendix A.3.

6 Final System

6.1 Machine Learning Model

To implement a machine learning model to minimize carbon emissions, we plan on using BO model similar to

[1] and [8]. We will optimize vCPU allocation and frequency scaling to minimize carbon with meeting SLOs as our constraint. Since we have collected data on 5 different functions, we will have a separate BO model for each function. Each function and its input will be encoded into an input vector for the BO model, either using feature engineering or a neural encoder similar to [8]. Then, the vCPU count and frequency will be concatenated to this input vector, and the BO model will search the space of vCPU count and frequency to minimize energy usage. The output of the BO model will be the vCPU count and frequency that minimizes energy usage while meeting SLOs. A diagram of the machine learning model is shown in Figure 3.

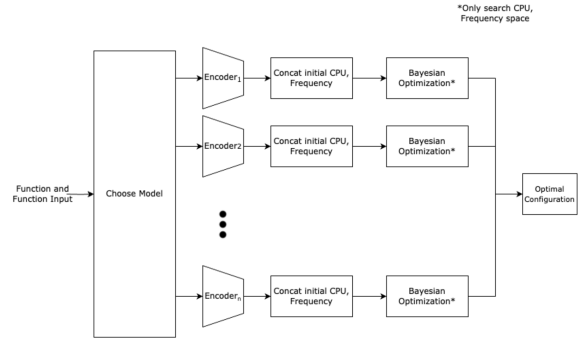


Figure 3: Bayesian Optimization Model for Carbon-Aware Serverless System

6.2 Adapting OpenWhisk

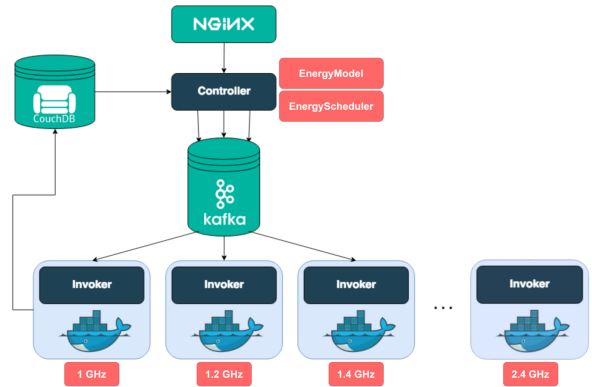


Figure 4: Adapted OpenWhisk System Overview

We adapted OpenWhisk to incorporate Bayesian Optimization to drive a carbon-aware DFS and vCPU allocation model. Figure 4 shows the system overview.

The Bayesian Optimization model is contained in a shim layer, called the EnergyModel. Whenever a user invokes a

function, the EnergyModel determines the optimal vCPU and frequency allocation (using Bayesian Optimization). The EnergyModel is trained online periodically.

We created a custom scheduler, called EnergyScheduler, that replaces OpenWhisk’s current scheduler. There are multiple invoker nodes, each running at a different frequency. The EnergyScheduler uses the output of EnergyModel to determine which invoker to send the function invocation to. We base the EnergyScheduler’s algorithm off of Shabri’s scheduler [7]. EnergyScheduler is designed to mitigate cold starts by:

1. Prioritizing routing the invocation to warm containers: The scheduler will receive the invocation and the predicted vCPU and frequency allocation from the EnergyModel. The scheduler will then iterate through the invokers running at at least the predicted frequency and attempt to find a warm container that is closest to the predicted vCPU allocation. If a warm container is found, the invocation is sent to that container. If no warm container is found, the invocation is sent to a randomly selected invoker.

2. Preemptively spinning up containers: In the case that there is no warm container with the exact predicted vCPU allocation on the invoker running at the predicted frequency, the scheduler will preemptively spin up a container on the invoker running at the predicted frequency with the predicted vCPU allocation. This is done to mitigate future cold starts and underutilization of resources. The underlying assumption is that function invocations exhibit temporal locality, meaning that in a short time frame, the same function is likely to be invoked again. Thus, preemptively spinning up a container will reduce the likelihood of cold starts for future invocations.

7 Evaluation

7.1 Methodology

Workloads. We evaluate our system using functions shown in Table 3. These are the same functions used in the measurement study. These functions vary in their parallelizability and resource requirements, making them suitable for evaluating the effectiveness of our system. We utilize the same trace used in Shabri [7] to simulate function invocations, which is a scaled down sample of the Azure trace. We utilize a requests per second (RPS) of 1 for the evaluation.

Metrics.

Baselines. We compare our system to two baselines Parrotfish and Aquatope. Parrotfish utilizes parametric regression to predict resource configuration, while Aquatope uses Bayesian Optimization to predict resource configuration. We compare our system to these baselines in terms of energy consumption and SLO violations.

7.2 Results

References

- [1] O. Alipourfard, H. H. Liu, J. Chen, S. Venkataraman, M. Yu, and M. Zhang. CherryPick: Adaptively unearthing the best cloud configurations for big data analytics. In *14th USENIX Symposium on Networked Systems Design and Implementation (NSDI 17)*, pages 469–482, Boston, MA, Mar. 2017. USENIX Association.
- [2] M. Bilal, M. Canini, R. Fonseca, and R. Rodrigues. With great freedom comes great opportunity: Rethinking resource allocation for serverless functions. In *Proceedings of the Eighteenth European Conference on Computer Systems, EuroSys ’23*, page 381–397, New York, NY, USA, 2023. Association for Computing Machinery.
- [3] W. Buchanan, J. Foxon, D. Cooke, S. Iyer, E. Graham, B. DeRusha, C. Binder, K. Chiu, L. Corso, H. Richardson, V. Knight, A. Hussain, A. Allison, and N. Mathews. Carbon-aware computing: Measuring and reducing the carbon intensity associated with software in execution. White Paper, 2023.
- [4] M. Chadha, T. Subramanian, E. Arima, M. Gerndt, M. Schulz, and O. Abboud. Greencourier: Carbon-aware scheduling for serverless functions. In *WoSC ’23: Proceedings of the 9th International Workshop on Serverless Computing*, pages 18–23. ACM Press, Dec. 2023.
- [5] U. Gupta, M. Elgamal, G. Hills, G.-Y. Wei, H.-H. S. Lee, D. Brooks, and C.-J. Wu. Act: designing sustainable computer systems with an architectural carbon modeling tool. In *Proceedings of the 49th Annual International Symposium on Computer Architecture, ISCA ’22*, page 784–799, New York, NY, USA, 2022. Association for Computing Machinery.
- [6] H. Hè, M. Friedman, and T. Rekatsinas. Energat: Fine-grained energy attribution for multi-tenancy. In *Proceedings of the 2nd Workshop on Sustainable Computer Systems, Hot-Carbon ’23*. ACM, July 2023.
- [7] P. Sinha, K. Kaffes, and N. J. Yadwadkar. Shabri: Delayed decision-making for faster and efficient serverless functions, 2024.
- [8] Z. Zhou, Y. Zhang, and C. Delimitrou. Aquatope: QoS-and-uncertainty-aware resource management for multi-stage serverless workflows. In *ASPLOS 2023: Proceedings of the 28th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 1*, pages 1–14. ACM Press, Mar. 2023.

A Appendix

A.1 Energy Measurement

A toy example of the energy measurement of a container using EnergAt and EnergyLiteDaemon is shown in Figure 5. Figure 5a shows the energy readings of EnergAt. Figure 5b shows the energy readings of EnergyLiteDaemon. EnergyLiteDaemon doesn’t capture all energy readings for the

Function	Input Type	# Runs	Sizes	Size Range
floatmatmult	square matrix	1	1	1
linpack	square matrix	1	1	1
encrypt	string	1	1	1
imageprocess	image	1	1	1
videoprocess	video	1	1	1

Table 3: Summary of serverless functions studied.

container due to the round-robin sampling approach. Figure 5c shows EnergyLiteDaemon’s power calculation of each energy sample. Figure 5d shows the energy interpolation of EnergyLiteDaemon using trapezoidal integration.

A.2 Energy Variation with vCPU and Frequency Allocation

The energy usage, duration, and CPU usage for Float Matrix Multiplication, Image Processing, and Encryption as a function of vCPU allocation are shown in Figure 6.

The energy usage and duration for Float Matrix Multiplication, Image Processing, and Encryption as a function of vCPU frequency are shown in Figure 7.

A.3 Energy Variation for Different Stages of a Serverless Function Invocation

We can see here the energy consumption of the container spin-up, idle time, and container spin-down for different image sizes in Figure 8.

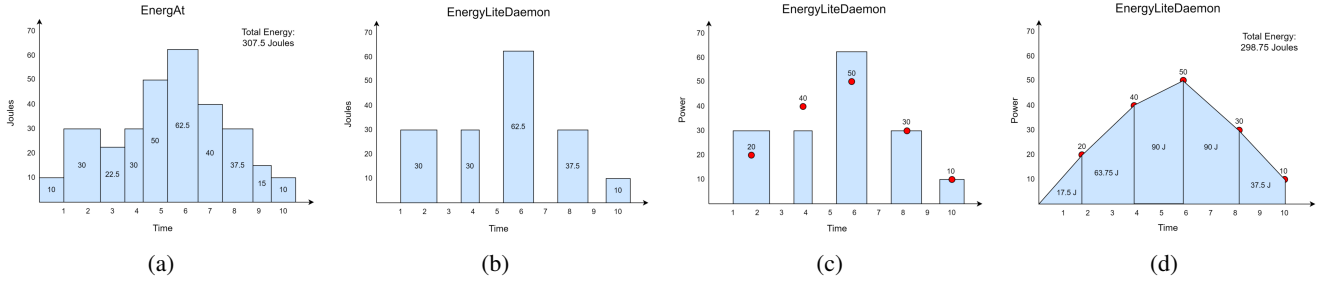


Figure 5: Toy example of energy measurement of a container using EnergyAt and EnergyLiteDaemon. (a) EnergyAt energy readings. (b) EnergyLiteDaemon energy readings. (c) EnergyLiteDaemon power calculation (red dots). (d) EnergyLiteDaemon energy interpolation using trapezoidal integration.

Float Matrix Multiplication

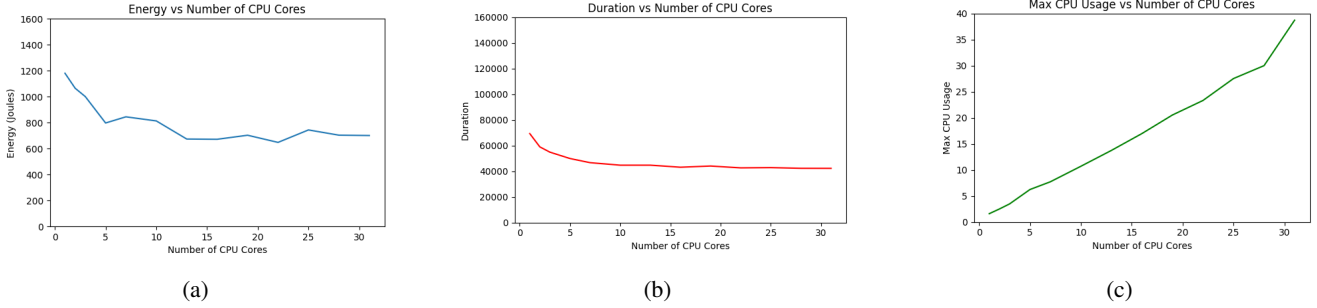
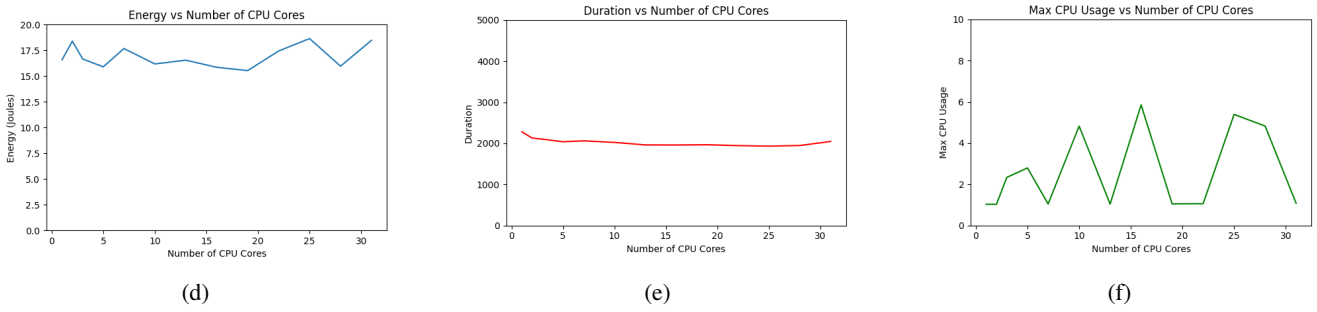


Image Process



Encryption

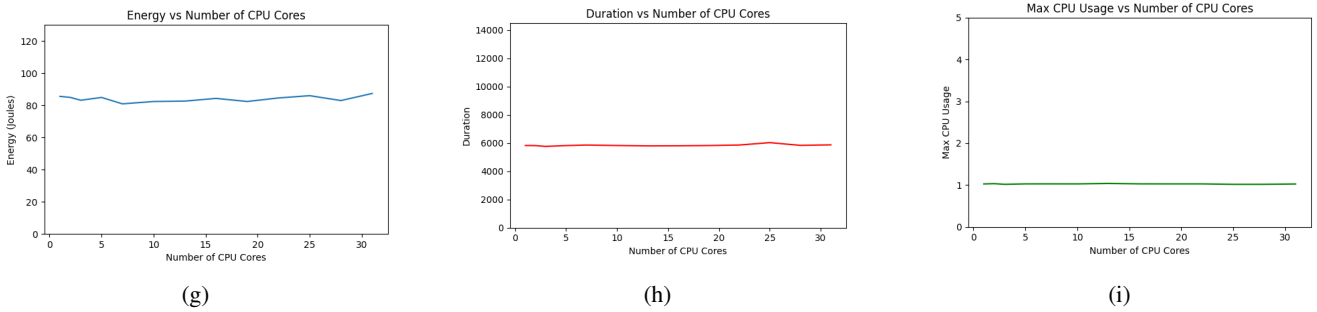
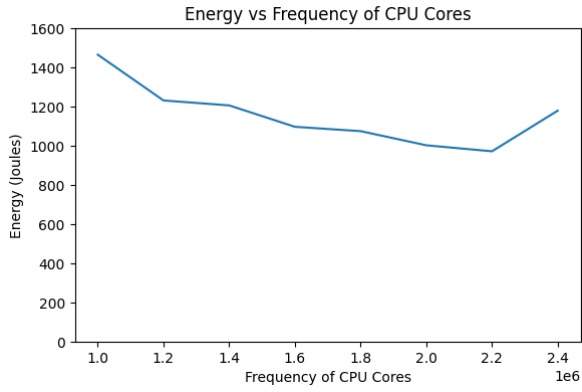
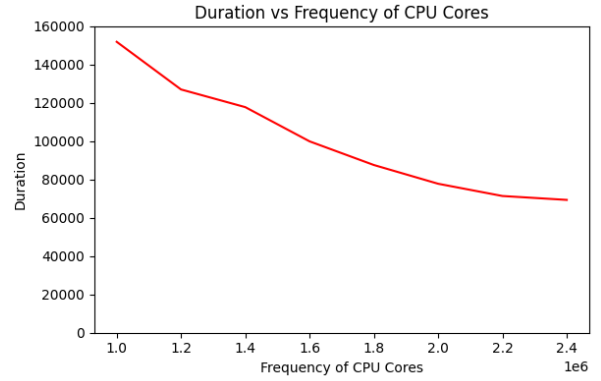


Figure 6: Energy Usage (a), Duration (b), and CPU Usage (c) for Float Matrix Multiplication, Image Processing, and Encryption as a function of vCPU allocation. Frequency was fixed at 2.4 GHz.

Float Matrix Multiplication

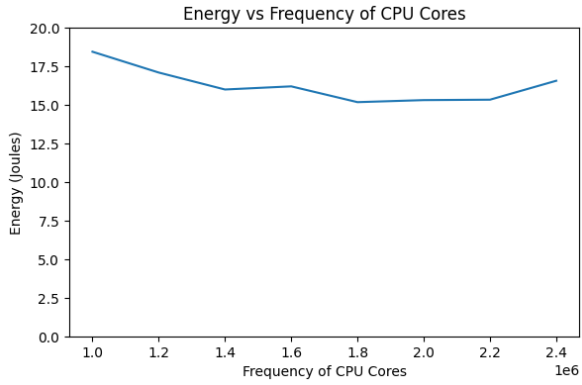


(a)

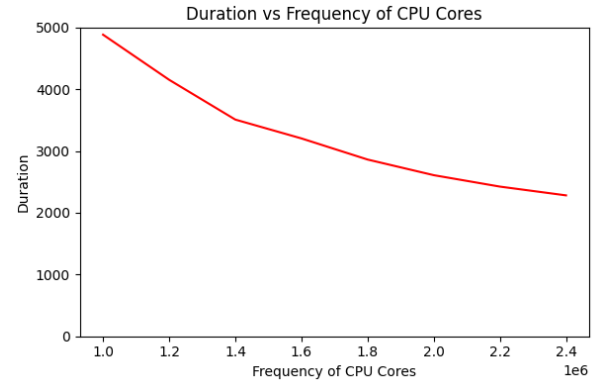


(b)

Image Process

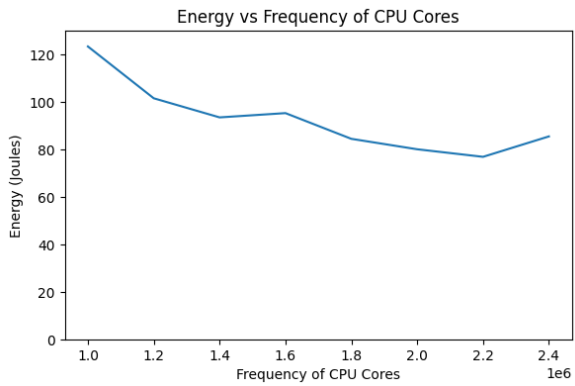


(c)

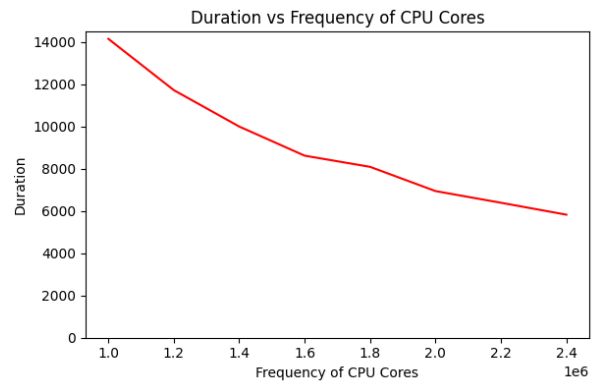


(d)

Encryption



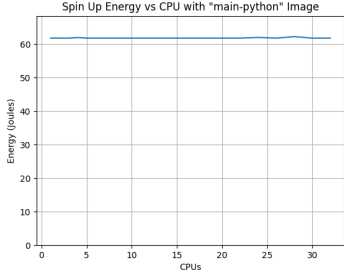
(e)



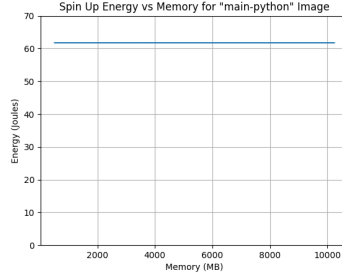
(f)

Figure 7: Energy Usage (a), Duration (b), and CPU Usage (c) for Float Matrix Multiplication, Image Processing, and Encryption as a function of vCPU frequency. vCPU allocation was fixed at 1.

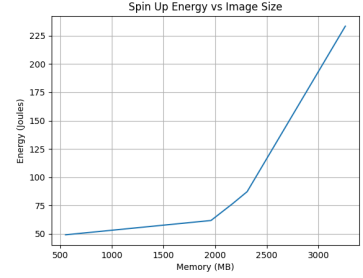
Spin Up Stage



(a)

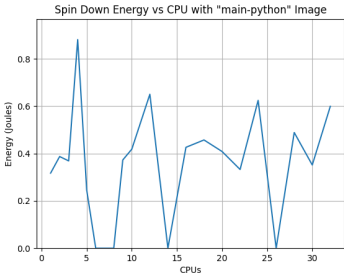


(b)

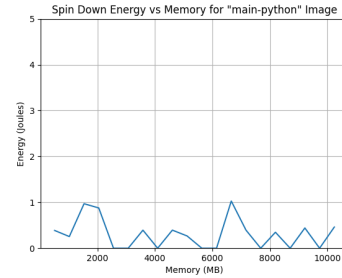


(c)

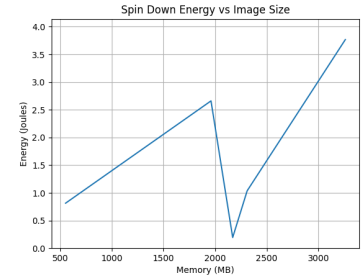
Spin Down Stage



(d)

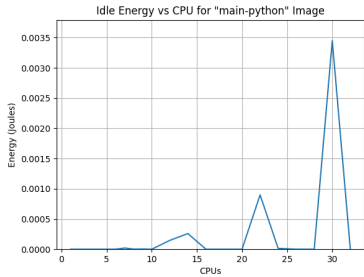


(e)

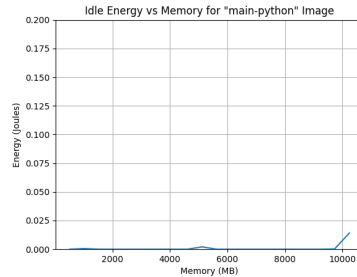


(f)

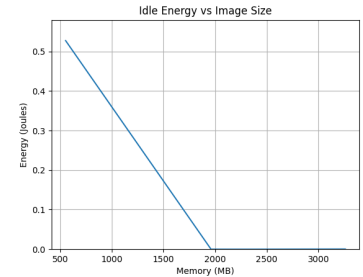
Idle Stage



(g)



(h)



(i)

Figure 8: Energy Usage vs CPU (a, d, g), Memory (b, e, h), and Image Size (c, f, i) for the Spin Up, Spin Down, and Idle stages of a serverless function invocation.