

# Optimizing Resource Allocation and Time Completion in Data Centers Using Integer Linear Programming

Daniel Safavisohi<sup>1</sup> and Dr. Melkonian<sup>2</sup>

<sup>1</sup>Graduate Student, Mathematics Department, Ohio University

<sup>2</sup>Associate Professor, Mathematics Department, Ohio University

December 13, 2024

## Abstract

This study presents optimization models for task assignment and resource allocation in data centers, with a focus on minimizing task completion time, energy consumption, and load imbalance. Two distinct models are developed: one leveraging batching with dynamic sizes and durations, and another incorporating multi-resource allocation and energy-aware scheduling. Experimental evaluations on synthetic datasets demonstrate that the first model effectively assigns tasks to servers within predefined batches, while the second model optimizes task prioritization and load balancing across resources. The results highlight the computational challenges associated with nonlinear constraints and the advantages of linearized models for scalability and efficiency. Future work includes integrating the strengths of both models and applying them to real-world datasets from industry leaders to enhance their applicability and performance in large-scale environments.

## 1 Introduction

In the era of big data and cloud computing, data centers play a pivotal role in processing and storing vast amounts of information. Efficient resource allocation and minimizing task completion time are critical for optimizing performance and reducing operational costs in data centers. Integer Linear Programming (ILP) offers a mathematical approach to model and solve such optimization problems. This project focuses on formulating an ILP model to optimize resource allocation and task scheduling in data centers, aiming to minimize total completion time while adhering to resource constraints.

## 2 Overview of Data Centers

Data centers are rapidly expanding across the globe to meet the escalating demand for digital services. The United States leads this growth, hosting approximately 2,670 data centers—the highest number of any country—contributing to a global total of over 8,000 facilities spread across regions like Europe, Asia-Pacific, and Latin America. Constructing and operating these centers involve substantial investments; building costs can range from \$ 10 million for small to medium-sized facilities to over \$1 billion for hyperscale data centers developed by tech giants such as Amazon, Google, and Facebook. Operational expenses, including energy consumption, maintenance, staffing, and technological upgrades, can amount to millions of dollars annually. Notable projects like Google's \$2.5 billion data center in Iowa and Microsoft's investments exceeding \$1 billion in various locations underscore the significant financial commitment required. The construction timeline for a data center typically spans 18 to 24 months, influenced by factors such as the facility's scale, technological requirements, regulatory approvals, and site-specific challenges.

### 2.1 Types of Projects Handled by Data Centers

Data centers are the backbone of modern digital infrastructure, supporting a wide array of projects critical to today's business operations and services. They provide the essential foundation for cloud computing platforms like Amazon Web Services (AWS), Microsoft Azure, and Google Cloud Platform, enabling scalable and flexible computing resources. In the realm of big data and analytics,

data centers process massive datasets for industries such as finance, healthcare, and retail, driving informed decision-making and strategic insights. They are instrumental in powering artificial intelligence (AI) and machine learning (ML) applications, hosting the computational resources required to train complex models. Additionally, data centers facilitate content delivery networks (CDNs), ensuring the rapid and efficient global distribution of media content and web applications. They manage data from the Internet of Things (IoT), overseeing interconnected devices and sensors used in smart cities, industrial automation, and consumer electronics. In financial services, data centers support high-frequency trading platforms, online banking, and transaction processing systems, maintaining the robustness and security of financial transactions. Furthermore, they run critical enterprise applications like Enterprise Resource Planning (ERP), Customer Relationship Management (CRM), and Supply Chain Management (SCM) systems, underpinning the operational efficiency of businesses worldwide.

## 2.2 Main Challenges Facing Data Centers Today

Data centers today face several significant challenges that impact their efficiency, cost-effectiveness, and sustainability. Energy consumption and environmental impact are major concerns, with data centers accounting for about 1 percent of global electricity use. This has led to a concerted effort to reduce carbon footprints through the adoption of renewable energy sources and more efficient technologies. Heat dissipation and cooling present another critical issue; managing the heat generated by high-density computing equipment is essential, and traditional cooling methods are energy-intensive. Innovations like liquid cooling and free-air cooling are being explored to address this. Cybersecurity threats have become increasingly sophisticated, necessitating robust security measures to protect sensitive data and ensure compliance with data protection regulations. The need for scalability and flexibility is a constant challenge due to rapid technological advancements, requiring data centers to implement modular designs and scalable infrastructure. Regulatory compliance adds another layer of complexity, as data centers must navigate complex regulations related to data sovereignty, privacy laws like GDPR, and industry-specific compliance standards. Lastly, supply chain disruptions caused by global events such as pandemics and geopolitical tensions can disrupt the supply of critical components, affecting both construction and maintenance.

## 2.3 Optimization Techniques in Data Centers

To enhance resource utilization and minimize completion times, data centers employ various optimization techniques that are essential for efficient operations. Load balancing is a fundamental strategy that distributes workloads evenly across servers, preventing both overloading and underutilization, which in turn enhances performance and reliability. Dynamic resource allocation further refines this process by adjusting resources in real-time based on current workload demands, thereby improving efficiency and responsiveness to fluctuating needs. Another critical technique is task consolidation, which groups smaller tasks to run concurrently on the same server, freeing up resources and reducing overall energy consumption. Additionally, data centers leverage predictive analytics, utilizing historical data and machine learning algorithms to forecast workload patterns. This allows for proactive adjustments in resource allocation, ensuring that resources are optimally used and ready to meet future demands.

## 2.4 Resource Scheduling Algorithms in Data Centers

Data centers utilize various resource scheduling algorithms to optimize task allocation and enhance overall efficiency. First-Come, First-Served (FCFS) processes tasks strictly in the order they arrive. While this method is straightforward, it may not be optimal for resource utilization or for handling time-sensitive tasks. Shortest Job First (SJF) prioritizes tasks with the shortest expected execution time, effectively reducing the average waiting time for all tasks. However, this can lead to the starvation of longer tasks if shorter ones continue to arrive. Priority Scheduling assigns tasks based on predetermined priority levels, ensuring that critical tasks receive immediate attention. This approach, though, can result in lower-priority tasks being neglected. Lastly, the Round Robin (RR) algorithm allocates fixed time slices to each task in a cyclic order, promoting fairness by giving all tasks equal opportunity to utilize resources. This method can increase overhead due to the frequent context switching between tasks, potentially impacting performance.

## 2.5 Batch Processing vs. Continuous Processing

Understanding the nature of computational tasks is essential for effective resource allocation in data centers, where tasks are generally categorized into batch processing and continuous processing. Batch processing involves executing a series of tasks collectively without manual intervention, making it suitable for non-interactive, time-insensitive tasks such as data analysis, report generation, and large-scale computations. This approach offers advantages like resource efficiency—allowing scheduling during off-peak hours to optimize utilization—and cost reduction by timing operations to coincide with lower energy rates or peak availability of renewable energy sources. It also simplifies management by reducing the complexity of real-time resource allocation. However, batch processing is not suitable for tasks requiring immediate results due to latency tolerance and requires careful scheduling to maximize resource utilization without impacting other operations.

In contrast, continuous processing handles tasks that require immediate processing, such as real-time data analytics, online transaction processing, and streaming services. Its advantages include low latency, providing immediate responses essential for user-facing applications; scalability through dynamic resource allocation to handle fluctuating workloads; and high availability to ensure services remain accessible at all times. The considerations for continuous processing involve higher resource demands, as it requires constant resource availability, potentially increasing operational costs. It also necessitates complex management with sophisticated scheduling and monitoring systems to maintain optimal performance.

## 2.6 Energy Management in Data Centers

Energy consumption constitutes a significant operational cost for data centers and has considerable environmental implications, making effective energy management strategies crucial for sustainable operations. A fundamental metric for assessing data center energy efficiency is Power Usage Effectiveness (PUE), defined as the ratio of total facility energy consumption to the energy consumed by IT equipment alone; a PUE value closer to 1 indicates higher efficiency. To optimize energy usage, data centers employ several techniques. Dynamic Voltage and Frequency Scaling (DVFS) adjusts the voltage and frequency of processors based on workload demands, reducing energy consumption during periods of low activity. Server virtualization consolidates multiple virtual servers onto a single physical server, optimizing hardware utilization and reducing the number of active physical servers required. Implementing efficient cooling systems further enhances energy efficiency: free cooling utilizes external environmental conditions, such as cool air or water, to lessen reliance on energy-intensive cooling systems, while hot/cold aisle containment organizes server racks to separate hot and cold air flows, improving cooling efficiency. Additionally, integrating renewable energy sources like solar or wind power reduces the carbon footprint and dependence on non-renewable energy sources. Effective workload management also plays a pivotal role in energy efficiency. Energy-aware scheduling incorporates energy consumption metrics into scheduling algorithms to balance performance with energy usage. Workload shifting transfers computational tasks to data centers in regions with lower energy costs or cooler climates, leveraging geographical advantages. Lastly, idle resource management powers down or places idle servers into low-power states to minimize unnecessary energy consumption. Together, these strategies contribute to more sustainable and cost-effective data center operations.

In the following section, we present three models along with various scenarios. We then select two models for evaluation using two synthetic datasets. Lastly, we conclude the study and provide appendices for additional analysis.

# 3 Models

## 3.1 First Model

The primary objective is to minimize the total completion time of tasks while efficiently allocating resources. The ILP model will focus on the following objective function:

$$\sum_{i \in T} C_i, \quad \forall i \in T$$

### Decision Variables

Let:

$x_{ij}$  be a binary decision variable where:

$$x_{ij} = \begin{cases} 1, & \text{if task } i \text{ is assigned to server } j, \\ 0, & \text{otherwise.} \end{cases}$$

$C_i$  is a variable and represents the completion time of task  $i$ .

### Parameters

- $T$  is the set of tasks,  $i \in T$
- $S$  is the set of servers,  $j \in S$
- $p_{ij}$  is the processing time required to complete task  $i$  on server  $j$ . This is a fixed parameter based on server capabilities and task requirements. For instance, we have an estimate that task A takes 5 hours on my laptop (Server 1) but only 1 hour on a powerful university PC (Server 2).
- $r_j$  is the resource capacity of server  $j$
- $d_i$  is the resource demand of task  $i$

### 3.2 Scenarios and Constraints

#### Scenario 1: Single Resource Type

In this scenario, we consider a data center where tasks require a single type of resource, such as CPU or GPU.

#### Constraints

**Assignment Constraint:** Each task must be assigned to exactly one server.

$$\sum_{j \in S} x_{ij} = 1, \quad \forall i \in T$$

**Resource Capacity Constraint:** The total resource demand on each server must not exceed its capacity.

$$\sum_{i \in T} d_i x_{ij} \leq r_j, \quad \forall j \in S$$

**Completion Time Calculation:** The completion time for each task is determined by its processing time on the assigned server.

$$C_i = \sum_{j \in S} p_{ij} x_{ij}, \quad \forall i \in T$$

#### Scenario 2: Multiple Resource Types

Here, tasks require multiple types of resources (such as CPU, memory, and storage).

#### Additional Parameters

- $R$  is the set of resource types,  $k \in R$
- $d_{ik}$  is the demand of resource  $k$  by task  $i$
- $r_{jk}$  is the capacity of resource  $k$  on server  $j$

#### New Constraints

**Resource Capacity Constraints:** For each resource type, the total demand must not exceed the server's capacity.

$$\sum_{i \in T} d_{ik} x_{ij} \leq r_{jk}, \quad \forall j \in S, \forall k \in R$$

### Scenario 3: Batch Scheduling with Dynamic Batch Sizes and Durations

159

This model allows batches to have variable sizes and durations, optimizing the schedule based on task requirements. The model has a non-linear constraint (constraint 6).

160

161

#### Sets and Indices

162

- $T$ : Set of tasks, indexed by  $i$ .
- $S$ : Set of servers, indexed by  $j$ .
- $K$ : Set of batches (cycles), indexed by  $k$ .

163

164

165

#### Parameters

166

- $p_{ij}$ : Processing time required to complete task  $i$  on server  $j$ .
- $d_i$ : Resource demand of task  $i$ .
- $r_j$ : Resource capacity of server  $j$ .

167

168

169

#### Decision Variables

170

- $x_{ijk} \in \{0, 1\}$ : Binary variable equal to 1 if task  $i$  is assigned to server  $j$  in batch  $k$ ; 0 otherwise.
- $s_k \geq 0$ : Continuous variable representing the start time of batch  $k$ .
- $D_k \geq 0$ : Continuous variable representing the duration of batch  $k$ .
- $C_i \geq 0$ : Continuous variable representing the completion time of task  $i$ .

171

172

173

174

#### Objective Function

175

Minimize the total completion times of all tasks:

176

$$\min \sum_{i \in T} C_i \quad (1)$$

#### Constraints

177

**1. Assignment Constraint** Each task must be assigned to exactly one server in one batch:

178

$$\sum_{j \in S} \sum_{k \in K} x_{ijk} = 1, \quad \forall i \in T \quad (2)$$

**2. Resource Capacity Constraints** For each server in each batch, the total resource demand must not exceed its capacity:

179

180

$$\sum_{i \in T} d_i x_{ijk} \leq r_j, \quad \forall j \in S, \forall k \in K \quad (3)$$

**3. Batch Duration Constraints** The duration of each batch must cover the processing times of the tasks assigned to it:

181

182

$$D_k \geq p_{ij} x_{ijk}, \quad \forall i \in T, \forall j \in S, \forall k \in K \quad (4)$$

**4. Batch Sequencing Constraints** Batches are processed sequentially; the start time of the next batch begins after the previous one ends:

183

184

$$s_{k+1} \geq s_k + D_k, \quad \forall k \in K \quad (5)$$

**5. Batch Start Time** The first batch starts at or after time zero:

185

$$s_1 \geq 0 \quad (6)$$

**6. Completion Time Calculation (non-linear)** The completion time of each task is the start time of its batch plus its processing time:

186

187

$$C_i = \sum_{j \in S} \sum_{k \in K} (s_k + p_{ij}) x_{ijk}, \quad \forall i \in T \quad (7)$$

#### Scenario 4: Batch Scheduling with Fixed Batch Sizes and Durations

188

This model uses predetermined batch sizes and durations, providing a structured scheduling framework.

189

190

#### Additional Parameters

191

- $D_k$ : Fixed duration of batch  $k$  (given).
- $m$ : Maximum number of tasks per batch (fixed batch size).

192

193

#### Decision Variables

194

- $x_{ijk} \in \{0, 1\}$ : Binary variable equal to 1 if task  $i$  is assigned to server  $j$  in batch  $k$ ; 0 otherwise.
- $s_k \geq 0$ : Continuous variable representing the start time of batch  $k$ .
- $C_i \geq 0$ : Continuous variable representing the completion time of task  $i$ .

195

196

197

#### Objective Function

198

Minimize the total completion times of all tasks:

199

$$\min \sum_{i \in T} C_i \quad (8)$$

#### Constraints

200

**1. Assignment Constraint** Each task must be assigned to exactly one server in one batch:

201

$$\sum_{j \in S} \sum_{k \in K} x_{ijk} = 1, \quad \forall i \in T \quad (9)$$

**2. Resource Capacity Constraints** For each server in each batch, the total resource demand must not exceed its capacity:

202

203

$$\sum_{i \in T} d_i x_{ijk} \leq r_j, \quad \forall j \in S, \forall k \in K \quad (10)$$

**3. Batch Duration Constraints** The processing time of any task assigned to a batch must not exceed the fixed duration of that batch:

204

205

$$p_{ij} x_{ijk} \leq D_k, \quad \forall i \in T, \forall j \in S, \forall k \in K \quad (11)$$

**4. Batch Size Constraints** The number of tasks assigned to each batch must not exceed the maximum batch size  $m$ :

206

207

$$\sum_{i \in T} \sum_{j \in S} x_{ijk} \leq m, \quad \forall k \in K \quad (12)$$

**5. Batch Sequencing Constraints** Batches are processed sequentially; the start time of the next batch begins after the fixed duration of the current batch:

208

209

$$s_{k+1} \geq s_k + D_k, \quad \forall k \in K \quad (13)$$

**6. Batch Start Time** The first batch starts at or after time zero:

210

$$s_1 \geq 0 \quad (14)$$

**7. Completion Time Calculation** The completion time of each task is the start time of its batch plus its processing time:

211

212

$$C_i = \sum_{j \in S} \sum_{k \in K} (s_k + p_{ij}) x_{ijk}, \quad \forall i \in T \quad (15)$$

### 3.3 Second Model

The original objective function aims to minimize the total completion time of tasks. We can augment this by adding:

- **Energy Consumption Minimization:** Reduce the total energy consumption of servers.
- **Load Balancing:** Distribute tasks evenly across servers to prevent overloading.
- **Task Prioritization:** Prioritize critical tasks by assigning weights.

The new objective function becomes:

$$\text{Minimize } \sum_{i \in T} w_i C_i + \alpha \sum_{j \in S} E_j + \beta \sum_{j \in S} L_j$$

Where:

- $w_i$  is the priority weight of task  $i$ .
- $E_j$  is the energy consumption of server  $j$ .
- $L_j$  is a load balancing term for server  $j$ .
- $\alpha$  and  $\beta$  are scaling coefficients.

#### New Decision Variables and Parameters

##### Decision Variables

- $y_j$  is a binary variable where:

$$y_j = \begin{cases} 1, & \text{if server } j \text{ is active,} \\ 0, & \text{otherwise.} \end{cases}$$

##### Parameters

- $e_j$  is the energy consumption rate of server  $j$ .
- $w_i$  is the priority weight of task  $i$ .
- $M$  is a large constant used in constraints.

#### Additional Constraints

**Server Activation Constraint:** A server must be activated if any task is assigned to it.

$$\sum_{i \in T} x_{ij} \leq M y_j, \quad \forall j \in S$$

**Load Balancing Constraints:** Ensure that the number of tasks assigned to each server is within a certain range.

$$L_j = \left| \sum_{i \in T} x_{ij} - \frac{|T|}{|S|} \right|, \quad \forall j \in S$$

**Task Precedence Constraints:** Some tasks must be completed before others can start.

$$C_i + S_{ij} \leq C_k, \quad \forall (i, k) \in P$$

Where:

- $P$  is the set of task pairs with precedence relations.
- $S_{ij}$  is the setup time between tasks  $i$  and  $k$ .

**Time Window Constraints:** Tasks must start and finish within specific time windows.

$$s_i \leq C_i \leq f_i, \quad \forall i \in T$$

Where:

- $s_i$  is the earliest start time of task  $i$ .
- $f_i$  is the latest finish time of task  $i$ .

### 3.4 Third Model

To create a more accurate and efficient resource allocation model for data centers, we can include additional separate load balancing terms for each critical resource, incorporating utilization ratios. This enhancement allows us to balance the load of each resource type (such as CPU, GPU, Memory, Storage) across servers, ensuring that no single resource becomes a bottleneck. The updated model builds upon the original objective function, which aims to minimize the total completion time of tasks, energy consumption, and load imbalance.

#### New Objective Function

The new objective function is:

$$\text{Minimize } Z = \sum_{i \in T} w_i \sum_{j \in S} p_{i,j} x_{i,j} + \alpha \sum_{j \in S} e_j \left( \sum_{i \in T} (p_{i,j} + S_{i,j}) x_{i,j} \right) + \beta \sum_{j \in S} \sum_{k \in R} L_{j,k}$$

To better show it:

$$\text{Minimize } Z = \sum_{i \in T} w_i C_i + \alpha \sum_{j \in S} E_j + \beta \sum_{j \in S} \sum_{k \in R} L_{j,k}$$

Where:

- $w_i$ : Priority weight of task  $i$ .
- $C_i$ : Completion time of task  $i$ .
- $E_j$ : Energy consumption of server  $j$ .
- $e_j$ : is the energy consumption rate of server  $j$
- $p_{ij}$ : Processing time required to complete task  $i$  on server  $j$ .
- $L_{j,k}$ : Load imbalance term for server  $j$  for each resource  $k$  in the set of resources  $R$ .
- $S_{ij}$ : is the setup time between tasks  $i$  and  $k$ .
- $\alpha, \beta$ : Scaling coefficients.
- $T$ : Set of tasks.
- $S$ : Set of servers.
- $R$ : Set of resources.

#### Definitions

- **Resources** ( $R$ ): The set of key resources (GPU, CPU, Memory, Storage).
- **Utilization Ratio** ( $U_{j,k}$ ): The utilization of resource  $k$  on server  $j$ .
- **Average Utilization** ( $U_{\text{avg},k}$ ): The average utilization of resource  $k$  across all servers.
- **Load Imbalance** ( $L_{j,k}$ ): The absolute difference between  $U_{j,k}$  and  $U_{\text{avg},k}$ .

#### Parameters

- $d_{i,k}$ : Demand of resource  $k$  by task  $i$ .
- $r_{j,k}$ : Capacity of resource  $k$  on server  $j$ .
- $e_j$ : Energy consumption rate of server  $j$ .
- $p_{i,j}$ : Processing time of task  $i$  on server  $j$ .
- $w_i$ : Priority weight of task  $i$ .
- $M$ : A large constant for server activation constraints.
- $\alpha, \beta$ : Scaling coefficients.



## Decision Variables

- $x_{i,j} \in \{0, 1\}$ : Assignment of task  $i$  to server  $j$ .
- $y_j \in \{0, 1\}$ : Server activation indicator.
- $C_i \geq 0$ : Completion time of task  $i$ .
- $E_j \geq 0$ : Energy consumption of server  $j$ .
- $U_{j,k} \geq 0$ : Utilization of resource  $k$  on server  $j$ .
- $L_{j,k} \geq 0$ : Load imbalance of resource  $k$  on server  $j$ .

## Constraints

### 1. Assignment Constraint

Each task must be assigned to exactly one server:

$$\sum_{j \in S} x_{i,j} = 1, \quad \forall i \in T$$

### 2. Server Activation Constraint

A server must be activated if any task is assigned to it:

$$\sum_{i \in T} x_{i,j} \leq M y_j, \quad \forall j \in S$$

### 3. Resource Capacity Constraints

For each resource  $k$ , the total demand on a server cannot exceed its capacity:

$$\sum_{i \in T} d_{i,k} x_{i,j} \leq r_{j,k} y_j, \quad \forall j \in S, \forall k \in R$$

### 4. Completion Time Calculation

The completion time for each task is determined by its processing time on the assigned server:

$$C_i = \sum_{j \in S} p_{i,j} x_{i,j}, \quad \forall i \in T$$

### 5. Energy Consumption Calculation

The total energy consumption of a server over the period processing a set of sequential tasks (setup time included):

$$E_j = e_j \sum_{i \in T} (p_{i,j} + S_{i,j}) x_{i,j}, \quad \forall j \in S$$

### 6. Utilization Ratio Calculation

Calculate the utilization of each resource on each server:

$$U_{j,k} = \frac{\sum_{i \in T} d_{i,k} x_{i,j}}{r_{j,k}}, \quad \forall j \in S, \forall k \in R$$

### 7. Average Utilization Calculation

Compute the average utilization for each resource across all servers:

$$U_{\text{avg},k} = \frac{\sum_{j \in S} \sum_{i \in T} d_{i,k} x_{i,j}}{\sum_{j \in S} r_{j,k}}, \quad \forall k \in R$$

## 8. Load Imbalance Constraints

Calculate the load imbalance for each resource on each server:

$$L_{j,k} \geq |U_{j,k} - U_{\text{avg},k}|, \quad \forall j \in S, \forall k \in R$$

Alternatively:

$$L_{j,k} \geq U_{j,k} - U_{\text{avg},k}, \quad \forall j \in S, \forall k \in R$$

$$L_{j,k} \geq U_{\text{avg},k} - U_{j,k}, \quad \forall j \in S, \forall k \in R$$

## 4 Assessment of Optimization Models

In this study, we examined Model One, Scenario Three to evaluate the impact of batching on the optimal completion time for each server. Our objective was to determine whether the model could effectively assign tasks with varying demands to servers with different resource capacities across a set of predetermined batches. In this scenario, both the duration and size of each batch are unknown and are determined by the model itself. The model's sole issue is a non-linear constraint, which imposes certain limitations on solver selection.

Additionally, we analyzed Model Three to investigate how a diverse range of resources, energy efficiency, and load balancing influence the model's task assignment strategy. This model operates without the complexities associated with batching, allowing us to assess task assignments in a more straightforward context.

We applied these optimization models to two distinct datasets to evaluate their effectiveness in enhancing data center operations. Each model requires a specific dataset, and given the computational resources available for this project, we carefully considered the size of each dataset accordingly.

### First Experiment

The dataset for Model One, Scenario Three comprises 60 tasks, 4 servers, and 4 batches. Each server's capacity is a randomly assigned integer between 20 and 30, while each task's demand is a randomly assigned value between 3 and 8. Predicted processing time of task on each server is also provided. This setup allows the model to determine the optimal assignment of tasks to servers within the specified batches.

For this specific application, each solver presents its own advantages and disadvantages. Through a trial-and-error approach, we ultimately selected the solver that offered the lowest computation time. The solver chosen for this model is FilMINT, which is based on the LP/NLP algorithm developed by Quesada and Grossmann and implemented within a branch-and-cut framework. FilMINT was developed by a group of scientists, Kumar Abhishek, Sven Leyffer, and Jeff Linderoth. We utilized NEOS Server to run these experiments.

Here is a brief summary of model 1 scenario 3 statistics.

Table I: Model 1 Scenario 3 Statistics	
Parameter	Value
Number of Constraints	1,040
Number of Variables	1,029
Number of Continuous Variables	69
Number of Binary Variables	960

The FilMINT solver employed a robust Branch and Cut methodology integrated with Sequential Quadratic Programming (filterSQP) to efficiently identify the optimal solution for the given Mixed-Integer Nonlinear Programming (MINLP) problem. Initially, the solver performed presolve operations to eliminate redundant constraints, thereby simplifying the problem structure. FilMINT strategically handled the 60 nonlinear constraints by solving a single NLP relaxation, leveraging filterSQP to manage nonlinearity while maintaining a linear objective function for enhanced computational efficiency. The solver generated seven lifted knapsack covers as cutting planes to tighten the feasible region, effectively reducing the solution space without extensive cut generation. Additionally, active primal heuristics and bound improvement techniques facilitated rapid convergence

by quickly identifying feasible solutions and refining objective bounds. Notably, the Branch and Bound process concluded at the root node with a tree depth of zero, indicating that the optimal solution was attained without further branching. Overall, FilMINT demonstrated high efficiency by solving the problem within 0.44 seconds, underscoring its capability to handle large-scale MINLPs through a combination of advanced preprocessing, selective cut generation, and effective heuristic strategies.

The files included in Appendix A are model1.mod, data1.dat, and job1.run, all written in AMPL. Additionally, the data generator file, written in Python (datagen1.py), is attached in Appendix A for further analysis.

The results show that the model works fine and correctly assign tasks in various batches to each server based on the optimum processing time. The result table is shown in Figure 1. The complete table is attached in Appendix A. Figure 2 is also provided to show the completion time for each task. By allowing the model to determine batch sizes and durations dynamically, we observe that tasks are grouped in batches that minimize idle times and reduce total completion time. Figure 3 shows the distribution of tasks across batches, highlighting how the model leverages batching to enhance scheduling efficiency. The model considers the varying processing times of tasks on different servers, assigning tasks to servers where they can be completed more quickly. This strategic assignment contributes to the minimization of total completion time, as depicted in Figure 5, which shows the assigned processing time by server and batch.

Task	t1	t10	t11	t12	t13	t14	t15	t16	t17	t18	...	t55	t56	t57	t58	t59	t6	t60	t7	t8	t9
Server	s3	s4	s3	s1	s1	s3	s3	s4	s4	s4	...	s1	s1	s1	s2	s2	s3	s1	s2	s2	s2
Batch	2	1	3	4	4	2	1	1	3	1	...	2	1	1	3	2	1	1	1	2	1

Figure 1: A sample of tasks assigned to each server in each batch

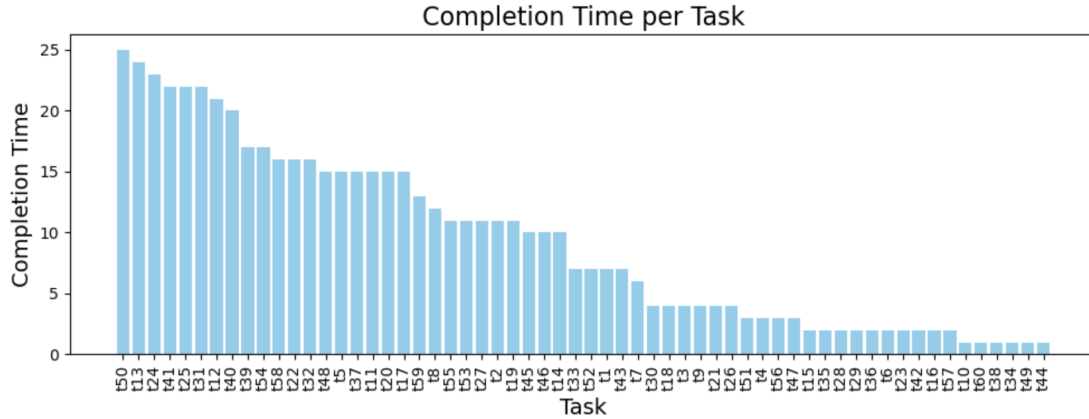


Figure 2: Completion time for each task

## Second Experiment

The dataset for Model 3 includes 60 tasks, 4 servers, and 4 resources: CPU, GPU, Memory, and Storage. The processing time for tasks is randomly assigned as an integer between 1 and 10, while setup time ranges from 0 to 3. The resource demand for each task ranges from 12 to 25, and the available resources for each server fall between 50 and 100. Task priority weights are randomly assigned as integers between 1 and 5. Additionally, the energy consumption rate for each server is randomly selected within the range of 5 to 15.

To address this model, we employed Cplex which is a well known model for solving mixed integer linear programming. We used IBM ILOG CPLEX Optimizer provided by NEOS to implement this experiment. A summary of model 3 statistics is provided in table 2.

The CPLEX solver version 22.1.1.0 efficiently tackled the Mixed-Integer Programming (MIP) problem by utilizing its advanced multi-threaded capabilities, employing four threads to enhance computational performance. The solver implemented a primal simplex algorithm, executing 154 MIP simplex iterations to methodically explore the feasible region and optimize the objective

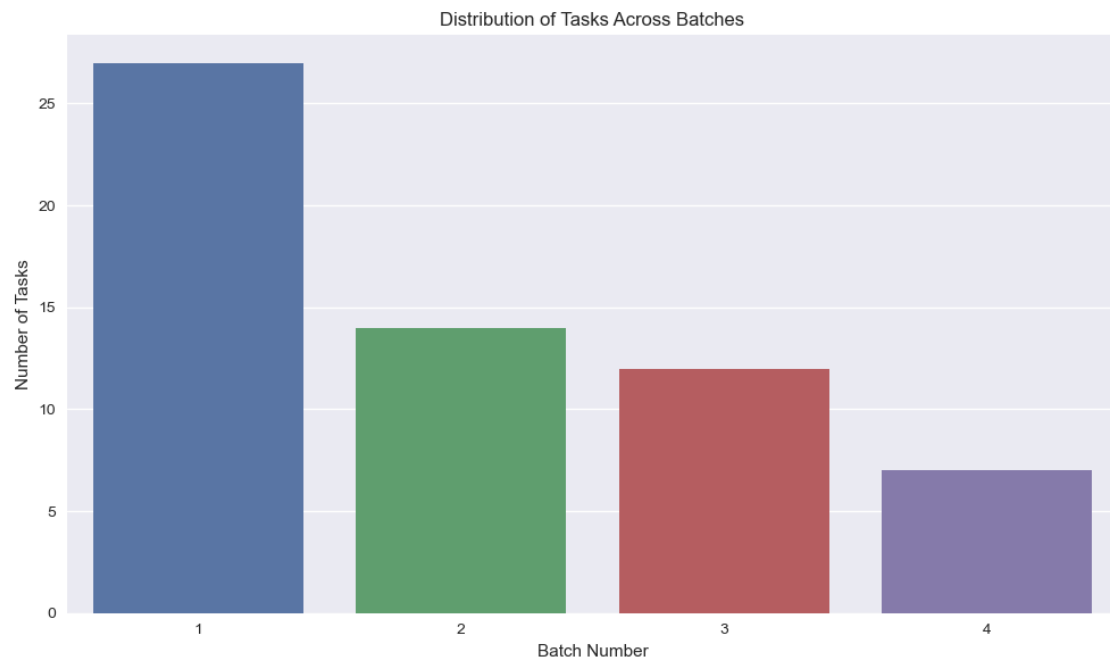


Figure 3: Distribution of tasks across batches

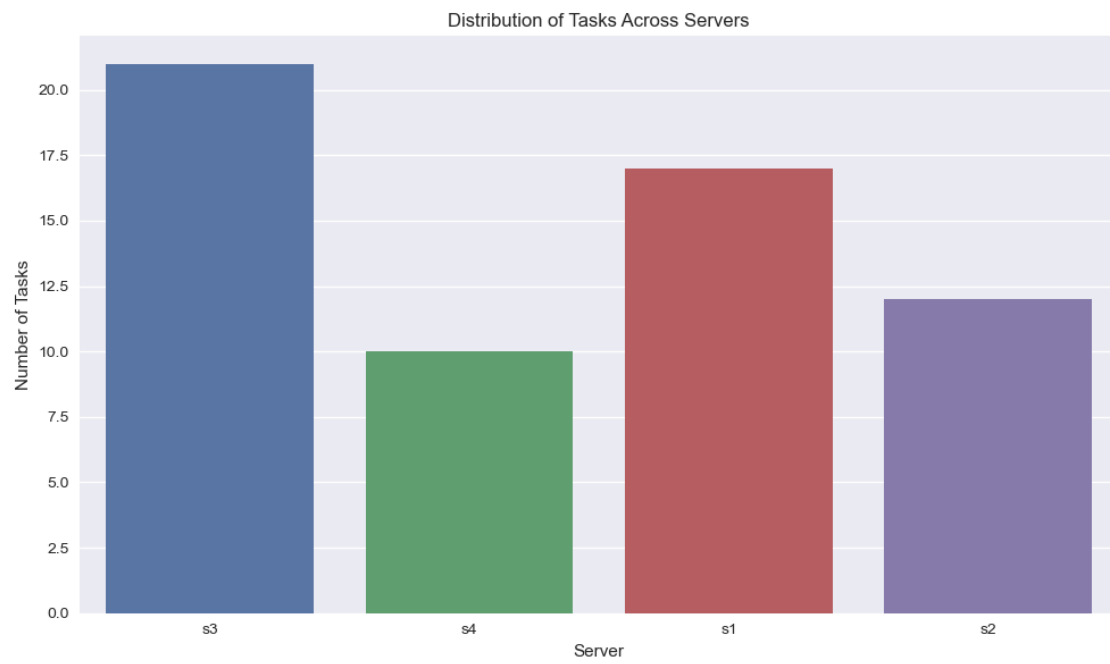


Figure 4: Distribution of tasks across servers

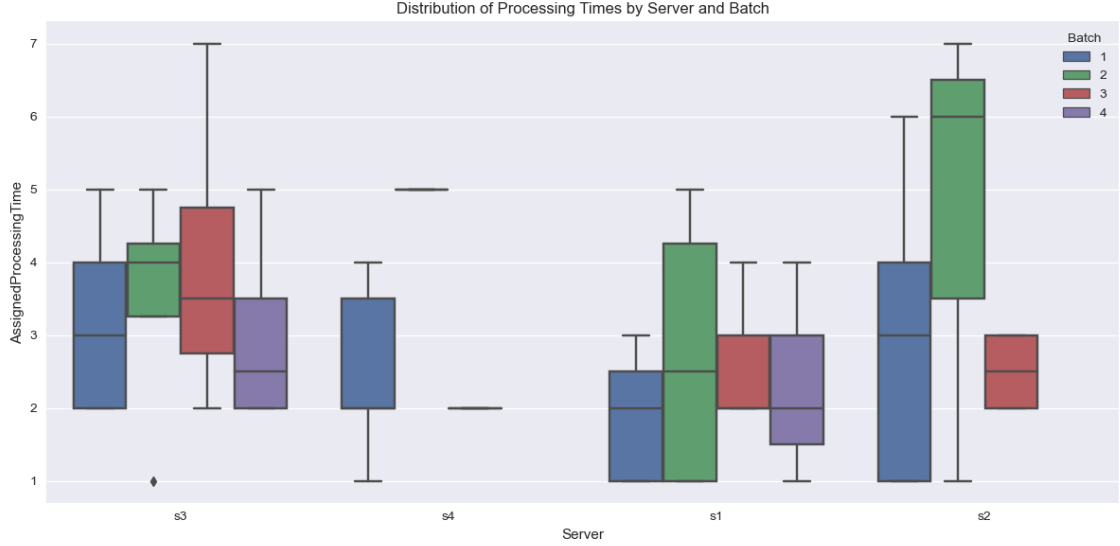


Figure 5: Distribution of assigned processing time by server and batch

Table II: Model 3 Statistics

Parameter	Value
Number of Equality Constraints	140
Number of Inequality Constraints	52
Number of Linear Variables	96
Number of Binary Variables	244

function, ultimately achieving an optimal integer solution with an objective value of 1728.301497 in less than two seconds. Notably, the optimization process concluded without initiating any branch-and-bound nodes, indicating that the linear programming relaxation of the problem was either inherently integer-feasible or that CPLEX’s sophisticated presolving and cutting-plane techniques were sufficiently effective in identifying the optimal solution without the need for further branching. The essential files associated with this model, including model2.mod, data2.dat, and job2.run, are provided in Appendix B. These files are all written in AMPL, ensuring consistency and ease of use within the modeling environment. Additionally, we have included the Python-based data generator script, datagen2.py, in Appendix B to support further analysis and replication of our results. Figure 6 presents the results for Model 3, highlighting the outcomes of our optimization efforts and offering insights into the performance and efficiency of task assignments across servers based on task priorities. The model achieves a balanced utilization of critical resources (CPU, GPU, Memory, and Storage) across all servers. Figure 7 demonstrates that the load imbalance for each resource is minimized, preventing any single resource from becoming a bottleneck. By incorporating energy consumption into the objective function, the model assigns tasks in a way that reduces the total energy usage. Servers with lower energy consumption rates are preferred for tasks with higher processing times, leading to an overall reduction in energy expenditure. The inclusion of task priority weights ensures that critical tasks are prioritized in the scheduling process. This is reflected in Figure 8, where high-priority tasks are assigned to servers capable of completing them more efficiently.

Comparing the results from both experiments, several insights emerge:

- **Model Complexity versus Scalability:** The nonlinear constraints in Model One (constraint 6), while allowing for dynamic batching, increase computational complexity and limit scalability. In contrast, the linearized Model Three handles larger datasets more efficiently, making it more suitable for real-world applications where computational resources may be limited.
- **Flexibility in Scheduling:** Model One offers greater flexibility in scheduling through dynamic batching, which can be advantageous in environments with highly variable workloads. However, this flexibility comes at the cost of increased computational overhead.
- **Comprehensive Resource Management:** Model Three’s incorporation of multiple resource

types and energy considerations provides a more holistic approach to resource management. This model effectively balances the utilization of different resources and aligns with sustainability objectives by minimizing energy consumption.

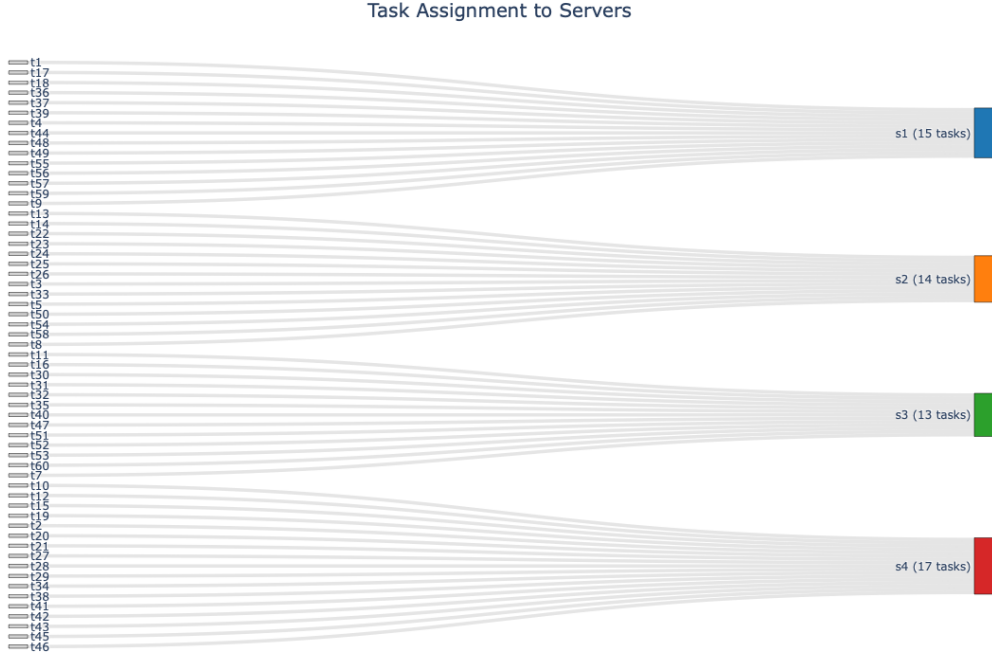


Figure 6: Task assignment to servers

## 5 Conclusion

In this project, we developed several optimization models aimed at enhancing task assignments within data center operations. Specifically, we proposed a couple of distinct models and applied two of them to synthetic datasets to evaluate their performance and effectiveness. Through these applications, we observed that the inherent nonlinearity within the models had a detrimental impact on both computation time and overall efficiency. The complex nonlinear constraints not only increased the computational burden, making the models less practical for larger datasets, but also resulted in out-of-disk or out-of-memory errors when scaling up, thereby further limiting their applicability in real-world, large-scale environments.

We analyzed various models designed for the first scenario, which involved nonlinear constraints. Our analysis revealed that the lack of proper documentation often creates confusion, limiting researchers' ability to fully leverage these models. Among the models tested, FilmINT demonstrated the best performance for the first scenario. However, we remain uncertain whether the limitations we encountered arose from the data generation process or from resource constraints on the NEOS server—particularly limited disk space for less commonly used solvers like FilmINT—which prevented us from testing larger datasets. In contrast, the model used in the second experiment, which was linear, posed fewer challenges when applied to large-scale datasets.

For future projects, we recommend adopting a strategic approach that combines the strengths of the two models explored. By linearizing these models, we can simplify complex relationships and reduce computational overhead, thereby improving both speed and efficiency. Linearized models are typically more tractable and faster to solve, making them suitable for real-time applications and larger datasets. Integrating the two models would address resource limitations on servers, enhance energy efficiency, and accommodate task prioritization. Additionally, distributing tasks across different batches would pave the way for more efficient models and improved overall performance.

Additionally, we propose applying these refined models to real-world datasets sourced from industry leaders such as Google and Microsoft. Utilizing actual data from these organizations will

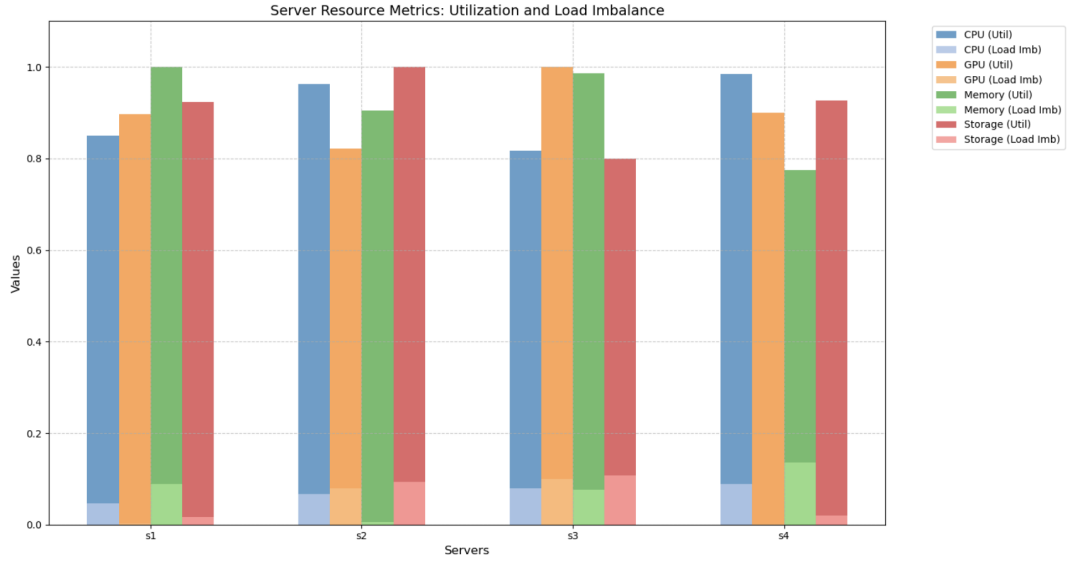


Figure 7: Server utilization and load imbalance for each server

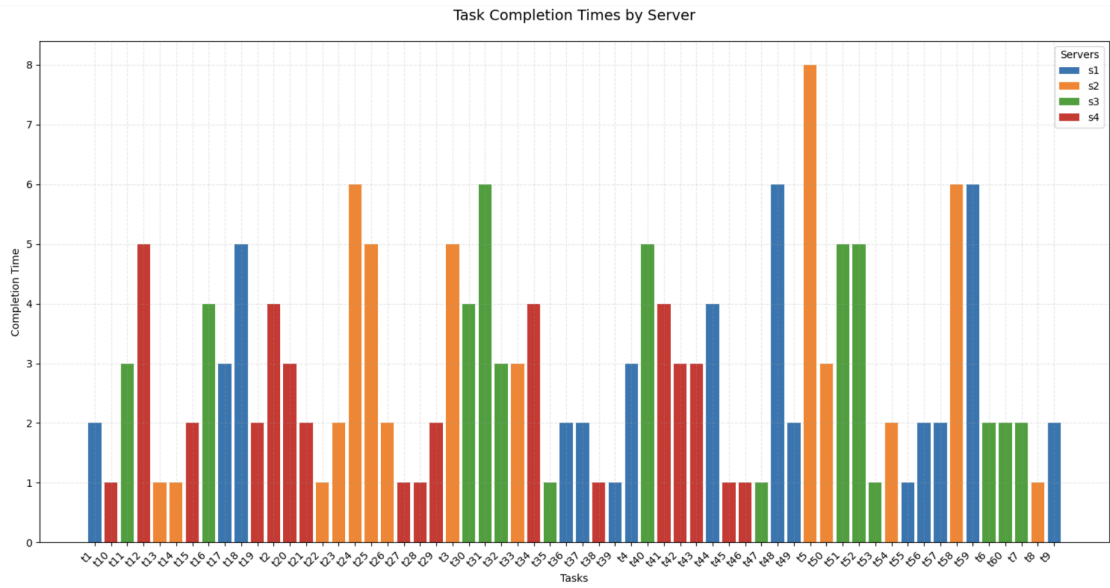


Figure 8: Completion time of assigned tasks on servers

provide a more accurate assessment of the models' applicability and performance in practical, large- 433  
scale environments. This real-world testing is essential for validating the models' effectiveness and 434  
ensuring they can meet the operational demands of modern data centers. Ultimately, this approach 435  
will contribute to the development of more robust and efficient task assignment solutions, leading 436  
to improved data center management and resource utilization. 437



## 6 References

- Barroso, L. A., Hölzle, U. (2009). The Datacenter as a Computer: An Introduction to the Design of Warehouse-Scale Machines. Morgan and Claypool Publishers. 438
- Beloglazov, A., Abawajy, J., Buyya, R. (2012). "Energy-aware resource allocation heuristics for efficient management of data centers for Cloud computing." Future Generation Computer Systems, 28(5), 755-768. 439
- Abhishek, K., Leyffer, S., and Linderoth, J. T. 2010. FilMINT: An Outer-Approximation-Based Solver for Nonlinear Mixed Integer Programs. INFORMS Journal on Computing 22: 555-567. DOI:10.1287/ijoc.1090.0373. 440
- Quesada, I. and I. E. Grossmann. 1992. An LP/NLP based branch-and-bound algorithm for convex MINLP optimization problems. Computers and Chemical Engineering 16: 937-947. 441
- IBM. 2023. IBM ILOG CPLEX Optimization Studio. Version 22.1.1.0. IBM. [www.ibm.com/products/ilog-cplex-optimization-studio](http://www.ibm.com/products/ilog-cplex-optimization-studio). 442
- Koomey, J. G. (2011). "Growth in data center electricity use 2005 to 2010." Analytics Press. 443
- Kliazovich, D., Bouvry, P., Khan, S. U. (2010). "GreenCloud: a packet-level simulator of energy-aware cloud computing data centers." The Journal of Supercomputing, 62(3), 1263-1283. 444
- Garg, S. K., Yeo, C. S., Anandasivam, A., Buyya, R. (2011). "Energy-efficient scheduling of HPC applications in cloud computing environments." Computing, 91(9), 1199-1219. 445
- OpenAI. (2024). Generative Pre-trained Transformer (November 7 Version). Retrieved from <https://chat.openai.com> 446
- Anthropic. (2024). Claude [Large Language Model]. Retrieved from <https://www.anthropic.com/claude> 447
- Overleaf. (2024). [Online LaTeX Editor]. Retrieved from <https://www.overleaf.com> 448
- Google Inc. (2019). Google Cluster Data V3. 449
- Fourer, R., Gay, D. M., & Kernighan, B. W. (2002). AMPL: A Modeling Language for Mathematical Programming (2nd ed.). Duxbury Press/Brooks/Cole Publishing Company. 450
- AMPL Optimization Inc. (2024). AMPL [Mathematical Programming Software]. Retrieved from <https://ampl.com> 451
- NEOS Server. (2024). Wisconsin Institute for Discovery at the University of Wisconsin in Madison. Retrieved from <https://neos-server.org> 452
- Vidyarthi, D. P., & Bhattacharya, B. (2008). Scheduling in distributed computing systems: Analysis, design and models. Springer. 453
- Marinescu, D. C. (2013). Cloud computing: Theory and practice. Morgan Kaufmann. 454
- Hunter, J. D. (2007). Matplotlib: A 2D graphics environment. Retrieved from <https://matplotlib.org/> 455
- The pandas development team. (2023). Zenodo. <https://doi.org/10.5281/zenodo.3509134> 456
- Plotly Technologies Inc. (2023). Plotly: Open-source graphing library for Python. Retrieved from <https://plotly.com/python/> 457
- Microsoft. (2021). Azure Data Center Workload Dataset. Microsoft Research. Retrieved from <https://github.com/Azure/azure-datacenter-workload-dataset> 458
- GitHub. (2024). GitHub Copilot. GitHub. Retrieved from <https://github.com/features/copilot> 459
- CBRE Group, Inc. (2022). Data Center Solutions Market Update. Retrieved from <https://www.cbre.com/> 460
- Uptime Institute. (2022). Global Data Center Survey. Retrieved from <https://uptimeinstitute.com/> 461

- Cushman & Wakefield. (2022). Global Data Center Market Comparison. Retrieved from <https://www.cushmanwakefield.com/>. 482  
483
- Statista. (2023). Number of data centers worldwide from 2015 to 2022, with forecasts until 2025. Retrieved from <https://www.statista.com/>. 484  
485
- Turner & Townsend. (2022). Data Center Cost Index. Retrieved from <https://www.turnerandtownsend.com/>. 486  
487
- Google. (2023). Google Data Centers. Retrieved from <https://www.google.com/about/datacenters/>. 488  
489
- Microsoft. (2023). Inside Microsoft’s Datacenters. Retrieved from <https://www.microsoft.com/en-us/datacenters>. 490  
491
- Amazon Web Services (AWS). (2023). AWS Cloud Products. Retrieved from <https://aws.amazon.com/products/>. 492  
493
- Microsoft Azure. (2023). Azure Products by Category. Retrieved from <https://azure.microsoft.com/>. 494  
495
- Google Cloud. (2023). Google Cloud Services. Retrieved from <https://cloud.google.com/products/>. 496  
497
- IDC (International Data Corporation). (2021). Worldwide Big Data and Analytics Software Forecast. Retrieved from <https://www.idc.com/>. 498  
499
- Gartner. (2022). Emerging AI and ML Use Cases in Data Centers. Retrieved from <https://www.gartner.com/>. 500  
501
- Akamai Technologies. (2023). Content Delivery Network (CDN) Services. Retrieved from <https://www.akamai.com/>. 502  
503
- Cisco. (2023). Internet of Things (IoT). Retrieved from <https://www.cisco.com/>. 504
- Bank for International Settlements (BIS). (2021). Technology in Financial Services. Retrieved from <https://www.bis.org/>. 505  
506
- International Energy Agency (IEA). (2021). Data Centres and Data Transmission Networks. Retrieved from <https://www.iea.org/reports/data-centres-and-data-transmission-networks>. 507  
508  
509
- Koomey, J. (2011). Growth in Data Center Electricity Use 2005 to 2010. Analytics Press. Retrieved from <http://www.analyticspress.com/datacenters.html>. 510  
511
- ASHRAE Technical Committee 9.9. (2021). Thermal Guidelines for Data Processing Environments. Retrieved from <https://www.ashrae.org/>. 512  
513
- ENISA (European Union Agency for Cybersecurity). (2022). Cyber Threat Landscape Report. Retrieved from <https://www.enisa.europa.eu/>. 514  
515
- European Commission. (2020). General Data Protection Regulation (GDPR). Retrieved from <https://gdpr.eu/>. 516  
517
- Supply Chain Management Review. (2021). Data Center Supply Chain Challenges. Retrieved from <https://www.scmr.com/>. 518  
519
- Cardoso, M., Singh, A., Mirhoseini, A., & Bruno, J. (2009). Exploiting Dynamic Resource Allocation for Efficient Parallel Data Processing in the Cloud. IEEE Cloud. 520  
521
- Beloglazov, A., Buyya, R., Lee, Y. C., & Zomaya, A. (2011). A Taxonomy and Survey of Energy-Efficient Data Centers and Cloud Computing Systems. Advances in Computers, Elsevier. 522  
523  
524
- Gmach, D., Rolia, J., Cherkasova, L., & Kemper, A. (2009). Resource Pool Management: Reactive versus Proactive or Let’s be Friends. IEEE Computer Society. 525  
526
- Silberschatz, A., Galvin, P. B., & Gagne, G. (2018). Operating System Concepts (10th ed.). Wiley. 527  
528

- Xu, J., & Fortes, J. A. B. (2010). Multi-Objective Virtual Machine Placement in Virtualized Data Center Environments. IEEE/ACM International Conference on Green Computing and Communications. 529 530 531
- Mishra, M., & Sahoo, A. (2011). On Theory of VM Placement: Anomalies in Existing Methodologies and Their Mitigation Using a Novel Vector Based Approach. IEEE Cloud. 532 533
- Dean, J., & Ghemawat, S. (2008). MapReduce: Simplified Data Processing on Large Clusters. Communications of the ACM, 51(1), 107-113. 534 535
- Kumar, V., Grama, A., Gupta, A., & Karypis, G. (2003). Introduction to Parallel Computing. Addison-Wesley. 536 537
- The Green Grid. (2020). Green Grid Data Center Power Efficiency Metrics: PUE and DCiE. Retrieved from <https://www.thegreengrid.org/>. 538 539
- Fan, X., Weber, W.-D., & Barroso, L. A. (2007). Power Provisioning for a Warehouse-sized Computer. ACM SIGARCH Computer Architecture News, 35(2), 13-23. 540 541
- Beloglazov, A., & Buyya, R. (2012). Optimal Online Deterministic Algorithms and Adaptive Heuristics for Energy and Performance Efficient Dynamic Consolidation of Virtual Machines in Cloud Data Centers. Concurrency and Computation: Practice and Experience, 24(13), 1397-1420. 542 543 544 545
- ASHRAE Datacom Series. (2015). Liquid Cooling Guidelines for Datacom Equipment Centers. American Society of Heating, Refrigerating and Air-Conditioning Engineers. 546 547
- Shehabi, A., Smith, S. J., Masanet, E., & Koomey, J. (2018). Data center growth in the United States: decoupling the demand for services from electricity use. Energy & Environmental Science, 11(3), 623-635. 548 549 550
- Google Sustainability. (2023). Our Data Centers: Efficiency and Sustainability. Retrieved from <https://sustainability.google/projects/data-centers/>. 551 552

## 7 Appendix A

553

- model 1 scenario 3 result

554

Table III: Model 1 Scenario 3 Results

Task	Demand	Server	Resource	Batch
t1	8	s3	30	2
t10	6	s4	28	1
t11	7	s3	30	3
t12	6	s1	24	4
t13	7	s1	24	4
t14	7	s3	30	2
t15	3	s3	30	1
t16	3	s4	28	1
t17	6	s4	28	3
t18	5	s4	28	1
t19	5	s4	28	2
t2	8	s4	28	2
t20	6	s2	22	3
t21	3	s3	30	1
t22	6	s3	30	3
t23	4	s4	28	1
t24	7	s3	30	4
t25	8	s1	24	4
t26	3	s3	30	1
t27	8	s4	28	2
t28	4	s4	28	1
t29	3	s3	30	1
t3	3	s4	28	1
t30	6	s3	30	1
t31	8	s3	30	4
t32	6	s2	22	3
t33	5	s1	24	2
t34	3	s1	24	1
t35	4	s3	30	1
t36	3	s3	30	1
t37	8	s1	24	3
t38	3	s2	22	1
t39	8	s3	30	3
t4	4	s2	22	1
t40	7	s3	30	3
t41	7	s3	30	4
t42	3	s1	24	1
t43	4	s2	22	2
t44	3	s1	24	1
t45	7	s1	24	2
t46	8	s3	30	2
t47	4	s1	24	1
t48	5	s2	22	3
t49	5	s2	22	1
t5	6	s1	24	3
t50	8	s3	30	4
t51	4	s3	30	1
t52	3	s1	24	2
t53	6	s3	30	2
t54	6	s1	24	3
t55	8	s1	24	2

Continued on next page

Table III: Model 1 Scenario 3 Results

Task	Demand	Server	Resource	Batch
t56	3	s1	24	1
t57	3	s1	24	1
t58	5	s2	22	3
t59	6	s2	22	2
t6	3	s3	30	1
t60	3	s1	24	1
t7	5	s2	22	1
t8	8	s2	22	2
t9	4	s2	22	1

• model1.mod

```

1  # Scenario 3
2  # Batch Scheduling with Dynamic Batch Sizes and Durations
3
4  # Sets and Indices
5  set T; # Set of tasks
6  set S; # Set of servers
7  param N integer > 0; # Number of batches
8  set K; # Set of batches
9
10 # Parameters
11 param p {T, S} >= 0; # Processing time required to complete task i
    on server j
12 param d {T} >= 0; # Resource demand of task i
13 param r {S} >= 0; # Resource capacity of server j
14
15 # Decision Variables
16 var x {T, S, K} binary; # 1 if task i is assigned to server j in
    batch k
17 var s {K} >= 0; # Start time of batch k
18 var D {K} >= 0; # Duration of batch k
19 var C {T} >= 0; # Completion time of task i
20
21 # Objective Function
22 minimize TotalCompletionTime:
23     sum {i in T} C[i];
24
25 # Constraints
26
27 # 1. Assignment Constraint
28 subject to Assignment {i in T}:
29     sum {j in S, k in K} x[i,j,k] = 1;
30
31 # 2. Resource Capacity Constraints
32 subject to ResourceCapacity {j in S, k in K}:
33     sum {i in T} d[i] * x[i,j,k] <= r[j];
34
35 # 3. Batch Duration Constraints
36 subject to BatchDuration {i in T, j in S, k in K}:
37     D[k] >= p[i,j] * x[i,j,k];
38
39 # 4. Batch Sequencing Constraints
40 subject to BatchSequencing {k in K: k < N}:
41     s[k+1] >= s[k] + D[k];
42
43 # 5. Batch Start Time
44 subject to BatchStartTime:
45     s[1] >= 0;
46
47 # 6. Completion Time Calculation
48 subject to CompletionTime {i in T}:
49     C[i] = sum {j in S, k in K} (s[k] + p[i,j]) * x[i,j,k];

```

• data1.dat

```

1  set T := t1 t2 t3 t4 t5 t6 t7 t8 t9 t10 t11 t12 t13 t14 t15 t16 t17
    t18 t19 t20 t21 t22 t23 t24 t25 t26 t27 t28 t29 t30 t31 t32 t33
    t34 t35 t36 t37 t38 t39 t40 t41 t42 t43 t44 t45 t46 t47 t48 t49
    t50 t51 t52 t53 t54 t55 t56 t57 t58 t59 t60;
2  set S := s1 s2 s3 s4;
3  param N := 4;
4  set K := 1 2 3 4;
5
6  param p : s1 s2 s3 s4 :=
7  t1 7 7 1 5
8  t2 9 8 7 5

```

```

9   t3 8 6 10 4
10  t4 9 3 5 3
11  t5 2 10 5 9
12  t6 10 3 5 2
13  t7 2 6 8 9
14  t8 2 6 7 6
15  t9 10 4 9 8
16  t10 8 9 5 1
17  t11 9 1 2 7
18  t12 1 10 8 6
19  t13 4 6 2 4
20  t14 10 4 4 3
21  t15 9 8 2 2
22  t16 6 9 8 2
23  t17 5 9 5 2
24  t18 9 6 9 4
25  t19 10 9 10 5
26  t20 8 2 10 7
27  t21 6 10 4 5
28  t22 3 4 3 1
29  t23 10 5 8 2
30  t24 2 3 3 1
31  t25 2 9 7 9
32  t26 5 9 4 4
33  t27 10 7 10 5
34  t28 8 8 6 2
35  t29 6 10 2 8
36  t30 10 6 4 4
37  t31 1 5 2 4
38  t32 6 3 6 7
39  t33 1 2 3 4
40  t34 1 10 9 10
41  t35 2 1 2 4
42  t36 10 10 2 7
43  t37 2 6 2 1
44  t38 10 1 4 3
45  t39 2 8 4 1
46  t40 1 9 7 10
47  t41 2 5 2 4
48  t42 2 5 6 7
49  t43 3 1 9 8
50  t44 1 10 2 7
51  t45 4 5 6 8
52  t46 10 3 4 1
53  t47 3 3 6 9
54  t48 5 2 10 8
55  t49 3 1 8 7
56  t50 10 9 5 6
57  t51 7 5 3 9
58  t52 1 8 2 6
59  t53 1 9 5 3
60  t54 4 8 6 10
61  t55 5 6 10 10
62  t56 3 5 7 7
63  t57 2 1 10 4
64  t58 6 3 4 4
65  t59 8 7 10 7
66  t60 1 7 10 7
67  ;
68
69  param d :=
70  t1 8
71  t2 8
72  t3 3
73  t4 4
74  t5 6

```

```

622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687

```

75	t6 3	688
76	t7 5	689
77	t8 8	690
78	t9 4	691
79	t10 6	692
80	t11 7	693
81	t12 6	694
82	t13 7	695
83	t14 7	696
84	t15 3	697
85	t16 3	698
86	t17 6	699
87	t18 5	700
88	t19 5	701
89	t20 6	702
90	t21 3	703
91	t22 6	704
92	t23 4	705
93	t24 7	706
94	t25 8	707
95	t26 3	708
96	t27 8	709
97	t28 4	710
98	t29 3	711
99	t30 6	712
100	t31 8	713
101	t32 6	714
102	t33 5	715
103	t34 3	716
104	t35 4	717
105	t36 3	718
106	t37 8	719
107	t38 3	720
108	t39 8	721
109	t40 7	722
110	t41 7	723
111	t42 3	724
112	t43 4	725
113	t44 3	726
114	t45 7	727
115	t46 8	728
116	t47 4	729
117	t48 5	730
118	t49 5	731
119	t50 8	732
120	t51 4	733
121	t52 3	734
122	t53 6	735
123	t54 6	736
124	t55 8	737
125	t56 3	738
126	t57 3	739
127	t58 5	740
128	t59 6	741
129	t60 3	742
130	;	743
131		744
132	param r :=	745
133	s1 24	746
134	s2 22	747
135	s3 30	748
136	s4 28	749
137	;	750
		751

• job1.run

752  
753



1	# Load the model file	754
2	#model model1_scenario3.mod;	755
3		756
4	# Load the data file	757
5	#data data1.dat;	758
6		759
7		760
8	# Solve the model	761
9	solve;	762
10		763
11	# Display the decision variables and parameters	764
12	display x;	765
13	display s;	766
14	display D;	767
15	display C;	768
		769

• datagen1.py

1		771
2		772
3		773
4	import random	774
5	import math	775
6		776
7	# Set random seed for reproducibility	777
8	random.seed(0)	778
9		779
10	num_tasks = 60	780
11	num_servers = 4	781
12	#num_batches = 8	782
13	batch_duration = 10	783
14		784
15		785
16	tasks = [f"t{i}" for i in range(1, num_tasks + 1)]	786
17	servers = [f"s{j}" for j in range(1, num_servers + 1)]	787
18		788
19		789
20	# Generate processing times	790
21	processing_times = {}	791
22	for t in tasks:	792
23	processing_times[t] = {}	793
24	for s in servers:	794
25	processing_times[t][s] = random.randint(1, 10)	795
26		796
27	# Generate resource demands	797
28	resource_demands = {t: random.randint(3, 8) for t in tasks}	798
29		799
30	# Generate resource capacities	800
31	resource_capacities = {s: random.randint(20, 30) for s in servers}	801
32		802
33	# Compute total processing time (minimum across servers for each task)	803
34	total_processing_time = 0	804
35	for t in tasks:	805
36	min_p = min(processing_times[t][s] for s in servers)	806
37	total_processing_time += min_p	807
38	# Compute total processing capacity per batch	808
39	processing_capacity_per_batch = len(servers) * batch_duration	809
40		810
41	# Estimate minimum number of batches based on processing time	811
42	min_batches_processing = math.ceil(total_processing_time /	812
43	processing_capacity_per_batch)	813
44		814
45	# Compute total resource demand	815
46	total_resource_demand = sum(resource_demands[t] for t in tasks)	816
47	# Compute total resource capacity per batch	817
		818
		819

```

48 total_resource_capacity_per_batch = sum(resource_capacities[s] for s
    in servers)
49
50 # Estimate minimum number of batches based on resource capacities
51 min_batches_resource = math.ceil(total_resource_demand /
    total_resource_capacity_per_batch)
52
53 # Compute the overall minimum number of batches required
54 min_batches_required = max(min_batches_processing,
    min_batches_resource)
55
56 # Print the minimum number of batches required
57 #print(f"Minimum number of batches required: {min_batches_required}")
58
59 batches = [str(k) for k in range(1, min_batches_required + 1)]
60 # Write data to AMPL data file
61 with open('data_synthetic_dec1.dat', 'w') as f:
62     # Write sets
63     f.write('set T := ' + ' '.join(tasks) + ';\n')
64     f.write('set S := ' + ' '.join(servers) + ';\n')
65     f.write(f'param N := {min_batches_required};\n')
66     f.write('set K := ' + ' '.join(batches) + ';\n\n')
67     #f.write(f'param D0 := {batch_duration};\n')
68
69     # Write processing times
70     f.write('param p : ' + ' '.join(servers) + ' :=\n')
71     for t in tasks:
72         f.write(t + ' ')
73         for s in servers:
74             f.write(str(processing_times[t][s]) + ' ')
75         f.write('\n')
76     f.write(';\n\n')
77
78     # Write resource demands
79     f.write('param d :=\n')
80     for t in tasks:
81         f.write(t + ' ' + str(resource_demands[t]) + '\n')
82     f.write(';\n\n')
83
84     # Write resource capacities
85     f.write('param r :=\n')
86     for s in servers:
87         f.write(s + ' ' + str(resource_capacities[s]) + '\n')
88     f.write(';\n')

```

## 8 Appendix B

- model2.mod

```

1 # -----
2 # Optimization Model3
3 # Objective: Minimize weighted completion time, energy consumption,
    and load imbalance
4 # -----
5
6 # Sets
7 set T;          # Set of Tasks
8 set S;          # Set of Servers
9 set R;          # Set of Resources
10
11 # Parameters
12 param w {T};    # Priority weight of task i
13 param e {S};    # Energy consumption rate of server
    j

```

```

14 param p {T, S}; # Processing time of task i on server j
15 param S_time {T, S}; # Setup time between task i and server j
16 param d {T, R}; # Demand of resource k by task i
17 param r {S, R}; # Capacity of resource k on server j
18 param alpha; # Scaling coefficient for energy consumption
19 param beta; # Scaling coefficient for load imbalance
20 param M; # Large constant for server activation constraints
21
22 # Precomputed Parameters
23 param sum_r_k {k in R} := sum {j in S} r[j, k]; # Sum of capacities for each resource
24
25 # Decision Variables
26 var x {T, S} binary; # Assignment of task i to server j
27 var y {S} binary; # Server activation indicator
28 var C {T} >= 0; # Completion time of task i
29 var E {S} >= 0; # Energy consumption of server j
30 var U {S, R} >= 0; # Utilization of resource k on server j
31 var L {S, R} >= 0; # Load imbalance of resource k on server j
32
33 # Objective Function
34 minimize Z:
35     sum {i in T} w[i] * C[i]
36     + alpha * sum {j in S} E[j]
37     + beta * sum {j in S, k in R} L[j, k];
38
39 # Constraints
40
41 # 1. Assignment Constraint
42 subject to Assignment {i in T}:
43     sum {j in S} x[i, j] = 1;
44
45 # 2. Server Activation Constraint
46 subject to ServerActivation {j in S}:
47     sum {i in T} x[i, j] <= M * y[j];
48
49 # 3. Resource Capacity Constraints
50 subject to ResourceCapacity {j in S, k in R}:
51     sum {i in T} d[i, k] * x[i, j] <= r[j, k] * y[j];
52
53 # 4. Completion Time Calculation
54 subject to CompletionTime {i in T}:
55     C[i] = sum {j in S} p[i, j] * x[i, j];
56
57 # 5. Energy Consumption Calculation
58 subject to EnergyConsumption {j in S}:
59     E[j] = e[j] * sum {i in T} (p[i, j] + S_time[i, j]) * x[i, j];
60
61 # 6. Utilization Ratio Calculation
62 subject to UtilizationRatio {j in S, k in R}:
63     U[j, k] = sum {i in T} d[i, k] * x[i, j] / r[j, k];
64
65 # 7. Load Imbalance Constraints
66 subject to LoadImbalanceUpper {j in S, k in R}:
67     L[j, k] >= U[j, k] - (sum {j2 in S} sum {i in T} d[i, k] * x[i, j2] / sum_r_k[k]);
68
69 subject to LoadImbalanceLower {j in S, k in R}:
70     L[j, k] >= (sum {j2 in S} sum {i in T} d[i, k] * x[i, j2] /

```

```
sum_r_k[k]) - U[j, k];
```

949  
950

• data2.dat

951

```

1  # -----
2  # Data File for model 3
3  # Generated by datagen2.py
4  # -----
5
6  # Sets
7  set T := t1 t2 t3 t4 t5 t6 t7 t8 t9 t10 t11 t12 t13 t14 t15 t16 t17
      t18 t19 t20 t21 t22 t23 t24 t25 t26 t27 t28 t29 t30 t31 t32 t33
      t34 t35 t36 t37 t38 t39 t40 t41 t42 t43 t44 t45 t46 t47 t48 t49
      t50 t51 t52 t53 t54 t55 t56 t57 t58 t59 t60;
8  set S := s1 s2 s3 s4;
9  set R := GPU CPU Memory Storage;
10
11 # Parameters
12 # Priority weight of each task
13 param w :=
14     t1 5
15     t2 2
16     t3 3
17     t4 2
18     t5 2
19     t6 3
20     t7 5
21     t8 2
22     t9 5
23     t10 2
24     t11 1
25     t12 4
26     t13 3
27     t14 3
28     t15 1
29     t16 2
30     t17 5
31     t18 1
32     t19 5
33     t20 4
34     t21 5
35     t22 3
36     t23 1
37     t24 4
38     t25 3
39     t26 1
40     t27 5
41     t28 1
42     t29 2
43     t30 4
44     t31 3
45     t32 3
46     t33 5
47     t34 2
48     t35 4
49     t36 5
50     t37 4
51     t38 3
52     t39 3
53     t40 4
54     t41 2
55     t42 2
56     t43 3
57     t44 2
58     t45 2
59     t46 3

```

952  
953  
954  
955  
956  
957  
958  
959  
960  
961  
962  
963  
964  
965  
966  
967  
968  
969  
970  
971  
972  
973  
974  
975  
976  
977  
978  
979  
980  
981  
982  
983  
984  
985  
986  
987  
988  
989  
990  
991  
992  
993  
994  
995  
996  
997  
998  
999  
1000  
1001  
1002  
1003  
1004  
1005  
1006  
1007  
1008  
1009  
1010  
1011  
1012  
1013  
1014

60	t47 4	1015
61	t48 2	1016
62	t49 2	1017
63	t50 5	1018
64	t51 2	1019
65	t52 2	1020
66	t53 3	1021
67	t54 5	1022
68	t55 2	1023
69	t56 3	1024
70	t57 3	1025
71	t58 4	1026
72	t59 1	1027
73	t60 2	1028
74	;	1029
75		1030
76	# Energy consumption rate of each server	1031
77	param e :=	1032
78	s1 12	1033
79	s2 10	1034
80	s3 8	1035
81	s4 10	1036
82	;	1037
83		1038
84	# Processing time of each task on each server	1039
85	param p :=	1040
86	t1 s1 2	1041
87	t1 s2 10	1042
88	t1 s3 7	1043
89	t1 s4 5	1044
90	t2 s1 5	1045
91	t2 s2 9	1046
92	t2 s3 7	1047
93	t2 s4 4	1048
94	t3 s1 10	1049
95	t3 s2 5	1050
96	t3 s3 5	1051
97	t3 s4 5	1052
98	t4 s1 3	1053
99	t4 s2 5	1054
100	t4 s3 8	1055
101	t4 s4 10	1056
102	t5 s1 7	1057
103	t5 s2 8	1058
104	t5 s3 9	1059
105	t5 s4 7	1060
106	t6 s1 6	1061
107	t6 s2 8	1062
108	t6 s3 2	1063
109	t6 s4 6	1064
110	t7 s1 7	1065
111	t7 s2 5	1066
112	t7 s3 2	1067
113	t7 s4 6	1068
114	t8 s1 6	1069
115	t8 s2 1	1070
116	t8 s3 5	1071
117	t8 s4 7	1072
118	t9 s1 2	1073
119	t9 s2 3	1074
120	t9 s3 4	1075
121	t9 s4 5	1076
122	t10 s1 3	1077
123	t10 s2 2	1078
124	t10 s3 5	1079
125	t10 s4 1	1080

126	t11 s1 8	1081
127	t11 s2 2	1082
128	t11 s3 3	1083
129	t11 s4 3	1084
130	t12 s1 6	1085
131	t12 s2 7	1086
132	t12 s3 9	1087
133	t12 s4 5	1088
134	t13 s1 10	1089
135	t13 s2 1	1090
136	t13 s3 3	1091
137	t13 s4 4	1092
138	t14 s1 5	1093
139	t14 s2 1	1094
140	t14 s3 6	1095
141	t14 s4 2	1096
142	t15 s1 8	1097
143	t15 s2 4	1098
144	t15 s3 9	1099
145	t15 s4 2	1100
146	t16 s1 9	1101
147	t16 s2 6	1102
148	t16 s3 4	1103
149	t16 s4 7	1104
150	t17 s1 3	1105
151	t17 s2 10	1106
152	t17 s3 10	1107
153	t17 s4 10	1108
154	t18 s1 5	1109
155	t18 s2 8	1110
156	t18 s3 10	1111
157	t18 s4 10	1112
158	t19 s1 8	1113
159	t19 s2 3	1114
160	t19 s3 6	1115
161	t19 s4 2	1116
162	t20 s1 4	1117
163	t20 s2 10	1118
164	t20 s3 7	1119
165	t20 s4 3	1120
166	t21 s1 2	1121
167	t21 s2 7	1122
168	t21 s3 4	1123
169	t21 s4 2	1124
170	t22 s1 1	1125
171	t22 s2 1	1126
172	t22 s3 3	1127
173	t22 s4 8	1128
174	t23 s1 3	1129
175	t23 s2 2	1130
176	t23 s3 4	1131
177	t23 s4 6	1132
178	t24 s1 5	1133
179	t24 s2 6	1134
180	t24 s3 8	1135
181	t24 s4 6	1136
182	t25 s1 7	1137
183	t25 s2 5	1138
184	t25 s3 8	1139
185	t25 s4 9	1140
186	t26 s1 10	1141
187	t26 s2 2	1142
188	t26 s3 7	1143
189	t26 s4 6	1144
190	t27 s1 4	1145
191	t27 s2 4	1146

192	t27 s3 3	1147
193	t27 s4 1	1148
194	t28 s1 5	1149
195	t28 s2 3	1150
196	t28 s3 7	1151
197	t28 s4 1	1152
198	t29 s1 4	1153
199	t29 s2 4	1154
200	t29 s3 9	1155
201	t29 s4 2	1156
202	t30 s1 4	1157
203	t30 s2 8	1158
204	t30 s3 4	1159
205	t30 s4 7	1160
206	t31 s1 8	1161
207	t31 s2 5	1162
208	t31 s3 6	1163
209	t31 s4 10	1164
210	t32 s1 8	1165
211	t32 s2 7	1166
212	t32 s3 3	1167
213	t32 s4 5	1168
214	t33 s1 5	1169
215	t33 s2 3	1170
216	t33 s3 3	1171
217	t33 s4 9	1172
218	t34 s1 5	1173
219	t34 s2 8	1174
220	t34 s3 7	1175
221	t34 s4 4	1176
222	t35 s1 7	1177
223	t35 s2 3	1178
224	t35 s3 1	1179
225	t35 s4 2	1180
226	t36 s1 2	1181
227	t36 s2 3	1182
228	t36 s3 5	1183
229	t36 s4 5	1184
230	t37 s1 2	1185
231	t37 s2 8	1186
232	t37 s3 7	1187
233	t37 s4 7	1188
234	t38 s1 5	1189
235	t38 s2 8	1190
236	t38 s3 3	1191
237	t38 s4 1	1192
238	t39 s1 1	1193
239	t39 s2 6	1194
240	t39 s3 7	1195
241	t39 s4 9	1196
242	t40 s1 7	1197
243	t40 s2 9	1198
244	t40 s3 5	1199
245	t40 s4 7	1200
246	t41 s1 3	1201
247	t41 s2 9	1202
248	t41 s3 2	1203
249	t41 s4 4	1204
250	t42 s1 6	1205
251	t42 s2 7	1206
252	t42 s3 7	1207
253	t42 s4 3	1208
254	t43 s1 9	1209
255	t43 s2 8	1210
256	t43 s3 9	1211
257	t43 s4 3	1212

258	t44 s1 4	1213
259	t44 s2 4	1214
260	t44 s3 5	1215
261	t44 s4 6	1216
262	t45 s1 7	1217
263	t45 s2 1	1218
264	t45 s3 1	1219
265	t45 s4 1	1220
266	t46 s1 4	1221
267	t46 s2 5	1222
268	t46 s3 3	1223
269	t46 s4 1	1224
270	t47 s1 8	1225
271	t47 s2 4	1226
272	t47 s3 1	1227
273	t47 s4 6	1228
274	t48 s1 6	1229
275	t48 s2 6	1230
276	t48 s3 7	1231
277	t48 s4 9	1232
278	t49 s1 2	1233
279	t49 s2 5	1234
280	t49 s3 7	1235
281	t49 s4 3	1236
282	t50 s1 4	1237
283	t50 s2 3	1238
284	t50 s3 5	1239
285	t50 s4 9	1240
286	t51 s1 9	1241
287	t51 s2 7	1242
288	t51 s3 5	1243
289	t51 s4 9	1244
290	t52 s1 10	1245
291	t52 s2 5	1246
292	t52 s3 5	1247
293	t52 s4 8	1248
294	t53 s1 9	1249
295	t53 s2 10	1250
296	t53 s3 1	1251
297	t53 s4 10	1252
298	t54 s1 1	1253
299	t54 s2 2	1254
300	t54 s3 10	1255
301	t54 s4 1	1256
302	t55 s1 1	1257
303	t55 s2 6	1258
304	t55 s3 8	1259
305	t55 s4 5	1260
306	t56 s1 2	1261
307	t56 s2 9	1262
308	t56 s3 8	1263
309	t56 s4 4	1264
310	t57 s1 2	1265
311	t57 s2 6	1266
312	t57 s3 6	1267
313	t57 s4 7	1268
314	t58 s1 10	1269
315	t58 s2 6	1270
316	t58 s3 8	1271
317	t58 s4 8	1272
318	t59 s1 6	1273
319	t59 s2 9	1274
320	t59 s3 5	1275
321	t59 s4 10	1276
322	t60 s1 10	1277
323	t60 s2 8	1278



```

324         t60 s3 2
325         t60 s4 6
326     ;
327
328     # Setup time between tasks on servers
329     param S_time :=
330         t1 s1 0
331         t1 s2 2
332         t1 s3 0
333         t1 s4 0
334         t2 s1 1
335         t2 s2 3
336         t2 s3 1
337         t2 s4 3
338         t3 s1 0
339         t3 s2 0
340         t3 s3 2
341         t3 s4 2
342         t4 s1 0
343         t4 s2 0
344         t4 s3 0
345         t4 s4 2
346         t5 s1 2
347         t5 s2 1
348         t5 s3 3
349         t5 s4 3
350         t6 s1 0
351         t6 s2 1
352         t6 s3 2
353         t6 s4 0
354         t7 s1 1
355         t7 s2 0
356         t7 s3 0
357         t7 s4 1
358         t8 s1 0
359         t8 s2 1
360         t8 s3 1
361         t8 s4 2
362         t9 s1 1
363         t9 s2 1
364         t9 s3 2
365         t9 s4 3
366         t10 s1 2
367         t10 s2 3
368         t10 s3 0
369         t10 s4 2
370         t11 s1 2
371         t11 s2 2
372         t11 s3 0
373         t11 s4 3
374         t12 s1 1
375         t12 s2 3
376         t12 s3 3
377         t12 s4 1
378         t13 s1 3
379         t13 s2 0
380         t13 s3 2
381         t13 s4 2
382         t14 s1 3
383         t14 s2 0
384         t14 s3 0
385         t14 s4 0
386         t15 s1 2
387         t15 s2 3
388         t15 s3 2
389         t15 s4 0

```

```

1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295
1296
1297
1298
1299
1300
1301
1302
1303
1304
1305
1306
1307
1308
1309
1310
1311
1312
1313
1314
1315
1316
1317
1318
1319
1320
1321
1322
1323
1324
1325
1326
1327
1328
1329
1330
1331
1332
1333
1334
1335
1336
1337
1338
1339
1340
1341
1342
1343
1344

```

390	t16 s1 1	1345
391	t16 s2 3	1346
392	t16 s3 1	1347
393	t16 s4 0	1348
394	t17 s1 0	1349
395	t17 s2 3	1350
396	t17 s3 2	1351
397	t17 s4 2	1352
398	t18 s1 0	1353
399	t18 s2 1	1354
400	t18 s3 2	1355
401	t18 s4 3	1356
402	t19 s1 1	1357
403	t19 s2 2	1358
404	t19 s3 1	1359
405	t19 s4 1	1360
406	t20 s1 3	1361
407	t20 s2 1	1362
408	t20 s3 0	1363
409	t20 s4 2	1364
410	t21 s1 3	1365
411	t21 s2 1	1366
412	t21 s3 2	1367
413	t21 s4 0	1368
414	t22 s1 2	1369
415	t22 s2 0	1370
416	t22 s3 0	1371
417	t22 s4 3	1372
418	t23 s1 2	1373
419	t23 s2 0	1374
420	t23 s3 2	1375
421	t23 s4 3	1376
422	t24 s1 1	1377
423	t24 s2 1	1378
424	t24 s3 0	1379
425	t24 s4 2	1380
426	t25 s1 3	1381
427	t25 s2 2	1382
428	t25 s3 1	1383
429	t25 s4 1	1384
430	t26 s1 0	1385
431	t26 s2 3	1386
432	t26 s3 2	1387
433	t26 s4 3	1388
434	t27 s1 0	1389
435	t27 s2 1	1390
436	t27 s3 1	1391
437	t27 s4 3	1392
438	t28 s1 2	1393
439	t28 s2 0	1394
440	t28 s3 2	1395
441	t28 s4 2	1396
442	t29 s1 3	1397
443	t29 s2 1	1398
444	t29 s3 2	1399
445	t29 s4 3	1400
446	t30 s1 2	1401
447	t30 s2 3	1402
448	t30 s3 2	1403
449	t30 s4 2	1404
450	t31 s1 1	1405
451	t31 s2 3	1406
452	t31 s3 0	1407
453	t31 s4 0	1408
454	t32 s1 3	1409
455	t32 s2 3	1410

456	t32 s3 2	1411
457	t32 s4 0	1412
458	t33 s1 2	1413
459	t33 s2 0	1414
460	t33 s3 2	1415
461	t33 s4 2	1416
462	t34 s1 1	1417
463	t34 s2 1	1418
464	t34 s3 1	1419
465	t34 s4 0	1420
466	t35 s1 2	1421
467	t35 s2 2	1422
468	t35 s3 2	1423
469	t35 s4 2	1424
470	t36 s1 3	1425
471	t36 s2 3	1426
472	t36 s3 1	1427
473	t36 s4 2	1428
474	t37 s1 0	1429
475	t37 s2 3	1430
476	t37 s3 0	1431
477	t37 s4 2	1432
478	t38 s1 1	1433
479	t38 s2 1	1434
480	t38 s3 3	1435
481	t38 s4 0	1436
482	t39 s1 3	1437
483	t39 s2 1	1438
484	t39 s3 1	1439
485	t39 s4 1	1440
486	t40 s1 2	1441
487	t40 s2 3	1442
488	t40 s3 3	1443
489	t40 s4 1	1444
490	t41 s1 3	1445
491	t41 s2 0	1446
492	t41 s3 1	1447
493	t41 s4 0	1448
494	t42 s1 1	1449
495	t42 s2 3	1450
496	t42 s3 2	1451
497	t42 s4 2	1452
498	t43 s1 0	1453
499	t43 s2 0	1454
500	t43 s3 3	1455
501	t43 s4 0	1456
502	t44 s1 2	1457
503	t44 s2 3	1458
504	t44 s3 2	1459
505	t44 s4 3	1460
506	t45 s1 1	1461
507	t45 s2 0	1462
508	t45 s3 3	1463
509	t45 s4 3	1464
510	t46 s1 1	1465
511	t46 s2 0	1466
512	t46 s3 0	1467
513	t46 s4 2	1468
514	t47 s1 1	1469
515	t47 s2 3	1470
516	t47 s3 3	1471
517	t47 s4 0	1472
518	t48 s1 0	1473
519	t48 s2 0	1474
520	t48 s3 2	1475
521	t48 s4 0	1476

522	t49 s1 1	1477
523	t49 s2 2	1478
524	t49 s3 1	1479
525	t49 s4 2	1480
526	t50 s1 2	1481
527	t50 s2 2	1482
528	t50 s3 3	1483
529	t50 s4 3	1484
530	t51 s1 2	1485
531	t51 s2 2	1486
532	t51 s3 2	1487
533	t51 s4 2	1488
534	t52 s1 1	1489
535	t52 s2 0	1490
536	t52 s3 2	1491
537	t52 s4 0	1492
538	t53 s1 2	1493
539	t53 s2 3	1494
540	t53 s3 3	1495
541	t53 s4 1	1496
542	t54 s1 3	1497
543	t54 s2 1	1498
544	t54 s3 3	1499
545	t54 s4 0	1500
546	t55 s1 0	1501
547	t55 s2 3	1502
548	t55 s3 0	1503
549	t55 s4 3	1504
550	t56 s1 1	1505
551	t56 s2 2	1506
552	t56 s3 0	1507
553	t56 s4 1	1508
554	t57 s1 2	1509
555	t57 s2 1	1510
556	t57 s3 3	1511
557	t57 s4 0	1512
558	t58 s1 1	1513
559	t58 s2 3	1514
560	t58 s3 0	1515
561	t58 s4 1	1516
562	t59 s1 1	1517
563	t59 s2 0	1518
564	t59 s3 2	1519
565	t59 s4 3	1520
566	t60 s1 3	1521
567	t60 s2 0	1522
568	t60 s3 1	1523
569	t60 s4 3	1524
570	;	1525
571		1526
572	# Demand of each resource by each task	1527
573	param d :=	1528
574	t1 GPU 1	1529
575	t1 CPU 3	1530
576	t1 Memory 4	1531
577	t1 Storage 5	1532
578	t2 GPU 4	1533
579	t2 CPU 1	1534
580	t2 Memory 1	1535
581	t2 Storage 4	1536
582	t3 GPU 7	1537
583	t3 CPU 8	1538
584	t3 Memory 1	1539
585	t3 Storage 6	1540
586	t4 GPU 9	1541
587	t4 CPU 2	1542

588	t4 Memory 1	1543
589	t4 Storage 3	1544
590	t5 GPU 8	1545
591	t5 CPU 4	1546
592	t5 Memory 1	1547
593	t5 Storage 3	1548
594	t6 GPU 8	1549
595	t6 CPU 5	1550
596	t6 Memory 4	1551
597	t6 Storage 1	1552
598	t7 GPU 9	1553
599	t7 CPU 4	1554
600	t7 Memory 5	1555
601	t7 Storage 3	1556
602	t8 GPU 8	1557
603	t8 CPU 1	1558
604	t8 Memory 9	1559
605	t8 Storage 8	1560
606	t9 GPU 2	1561
607	t9 CPU 4	1562
608	t9 Memory 8	1563
609	t9 Storage 8	1564
610	t10 GPU 9	1565
611	t10 CPU 1	1566
612	t10 Memory 9	1567
613	t10 Storage 3	1568
614	t11 GPU 7	1569
615	t11 CPU 5	1570
616	t11 Memory 2	1571
617	t11 Storage 9	1572
618	t12 GPU 9	1573
619	t12 CPU 1	1574
620	t12 Memory 7	1575
621	t12 Storage 2	1576
622	t13 GPU 5	1577
623	t13 CPU 7	1578
624	t13 Memory 8	1579
625	t13 Storage 5	1580
626	t14 GPU 4	1581
627	t14 CPU 1	1582
628	t14 Memory 8	1583
629	t14 Storage 1	1584
630	t15 GPU 1	1585
631	t15 CPU 3	1586
632	t15 Memory 1	1587
633	t15 Storage 8	1588
634	t16 GPU 1	1589
635	t16 CPU 5	1590
636	t16 Memory 1	1591
637	t16 Storage 1	1592
638	t17 GPU 1	1593
639	t17 CPU 8	1594
640	t17 Memory 3	1595
641	t17 Storage 4	1596
642	t18 GPU 2	1597
643	t18 CPU 1	1598
644	t18 Memory 10	1599
645	t18 Storage 7	1600
646	t19 GPU 7	1601
647	t19 CPU 6	1602
648	t19 Memory 4	1603
649	t19 Storage 9	1604
650	t20 GPU 2	1605
651	t20 CPU 9	1606
652	t20 Memory 8	1607
653	t20 Storage 3	1608

654	t21 GPU 5	1609
655	t21 CPU 2	1610
656	t21 Memory 1	1611
657	t21 Storage 1	1612
658	t22 GPU 7	1613
659	t22 CPU 7	1614
660	t22 Memory 6	1615
661	t22 Storage 3	1616
662	t23 GPU 4	1617
663	t23 CPU 3	1618
664	t23 Memory 9	1619
665	t23 Storage 1	1620
666	t24 GPU 8	1621
667	t24 CPU 4	1622
668	t24 Memory 4	1623
669	t24 Storage 1	1624
670	t25 GPU 2	1625
671	t25 CPU 9	1626
672	t25 Memory 9	1627
673	t25 Storage 9	1628
674	t26 GPU 2	1629
675	t26 CPU 8	1630
676	t26 Memory 2	1631
677	t26 Storage 7	1632
678	t27 GPU 5	1633
679	t27 CPU 4	1634
680	t27 Memory 9	1635
681	t27 Storage 6	1636
682	t28 GPU 2	1637
683	t28 CPU 3	1638
684	t28 Memory 1	1639
685	t28 Storage 8	1640
686	t29 GPU 6	1641
687	t29 CPU 1	1642
688	t29 Memory 1	1643
689	t29 Storage 3	1644
690	t30 GPU 9	1645
691	t30 CPU 1	1646
692	t30 Memory 8	1647
693	t30 Storage 4	1648
694	t31 GPU 2	1649
695	t31 CPU 7	1650
696	t31 Memory 10	1651
697	t31 Storage 6	1652
698	t32 GPU 1	1653
699	t32 CPU 7	1654
700	t32 Memory 4	1655
701	t32 Storage 2	1656
702	t33 GPU 1	1657
703	t33 CPU 10	1658
704	t33 Memory 1	1659
705	t33 Storage 3	1660
706	t34 GPU 9	1661
707	t34 CPU 10	1662
708	t34 Memory 8	1663
709	t34 Storage 1	1664
710	t35 GPU 4	1665
711	t35 CPU 9	1666
712	t35 Memory 6	1667
713	t35 Storage 6	1668
714	t36 GPU 4	1669
715	t36 CPU 6	1670
716	t36 Memory 6	1671
717	t36 Storage 2	1672
718	t37 GPU 8	1673
719	t37 CPU 8	1674

720	t37 Memory 4	1675
721	t37 Storage 5	1676
722	t38 GPU 1	1677
723	t38 CPU 10	1678
724	t38 Memory 7	1679
725	t38 Storage 1	1680
726	t39 GPU 2	1681
727	t39 CPU 9	1682
728	t39 Memory 8	1683
729	t39 Storage 2	1684
730	t40 GPU 9	1685
731	t40 CPU 3	1686
732	t40 Memory 1	1687
733	t40 Storage 8	1688
734	t41 GPU 9	1689
735	t41 CPU 1	1690
736	t41 Memory 3	1691
737	t41 Storage 9	1692
738	t42 GPU 3	1693
739	t42 CPU 5	1694
740	t42 Memory 1	1695
741	t42 Storage 1	1696
742	t43 GPU 7	1697
743	t43 CPU 4	1698
744	t43 Memory 3	1699
745	t43 Storage 9	1700
746	t44 GPU 9	1701
747	t44 CPU 1	1702
748	t44 Memory 9	1703
749	t44 Storage 3	1704
750	t45 GPU 9	1705
751	t45 CPU 1	1706
752	t45 Memory 8	1707
753	t45 Storage 7	1708
754	t46 GPU 2	1709
755	t46 CPU 3	1710
756	t46 Memory 4	1711
757	t46 Storage 1	1712
758	t47 GPU 3	1713
759	t47 CPU 1	1714
760	t47 Memory 8	1715
761	t47 Storage 8	1716
762	t48 GPU 9	1717
763	t48 CPU 9	1718
764	t48 Memory 5	1719
765	t48 Storage 7	1720
766	t49 GPU 3	1721
767	t49 CPU 9	1722
768	t49 Memory 9	1723
769	t49 Storage 2	1724
770	t50 GPU 1	1725
771	t50 CPU 1	1726
772	t50 Memory 7	1727
773	t50 Storage 5	1728
774	t51 GPU 7	1729
775	t51 CPU 3	1730
776	t51 Memory 9	1731
777	t51 Storage 2	1732
778	t52 GPU 1	1733
779	t52 CPU 9	1734
780	t52 Memory 5	1735
781	t52 Storage 2	1736
782	t53 GPU 5	1737
783	t53 CPU 5	1738
784	t53 Memory 5	1739
785	t53 Storage 2	1740

786	t54 GPU 5	1741
787	t54 CPU 9	1742
788	t54 Memory 1	1743
789	t54 Storage 1	1744
790	t55 GPU 2	1745
791	t55 CPU 1	1746
792	t55 Memory 8	1747
793	t55 Storage 7	1748
794	t56 GPU 5	1749
795	t56 CPU 9	1750
796	t56 Memory 7	1751
797	t56 Storage 5	1752
798	t57 GPU 5	1753
799	t57 CPU 8	1754
800	t57 Memory 4	1755
801	t57 Storage 6	1756
802	t58 GPU 7	1757
803	t58 CPU 5	1758
804	t58 Memory 1	1759
805	t58 Storage 2	1760
806	t59 GPU 8	1761
807	t59 CPU 7	1762
808	t59 Memory 2	1763
809	t59 Storage 7	1764
810	t60 GPU 4	1765
811	t60 CPU 3	1766
812	t60 Memory 7	1767
813	t60 Storage 6	1768
814	;	1769
815		1770
816	# Capacity of each resource on each server	1771
817	param r :=	1772
818	s1 GPU 78	1773
819	s1 CPU 100	1774
820	s1 Memory 88	1775
821	s1 Storage 79	1776
822	s2 GPU 84	1777
823	s2 CPU 80	1778
824	s2 Memory 74	1779
825	s2 Storage 55	1780
826	s3 GPU 70	1781
827	s3 CPU 82	1782
828	s3 Memory 76	1783
829	s3 Storage 75	1784
830	s4 GPU 100	1785
831	s4 CPU 66	1786
832	s4 Memory 98	1787
833	s4 Storage 82	1788
834	;	1789
835		1790
836	# Scaling coefficients and large constant	1791
837	param alpha := 0.5;	1792
838	param beta := 0.3;	1793
839	param M := 1000;	1794
840		1795
841	‘	1796

• job2.run

1	# Load the model file	1798
2	#model model2.mod;	1799
3		1800
4	# Load the data file	1801
5	#data data2.dat;	1802
6		1803
7		1804



8	# Solve the model	1807
9	solve;	1808
10		1809
11	# Display the decision variables and parameters	1810
12	display x;	1811
13	display C;	1812
14	display E;	1813
15	display U;	1814
16	display L;	1815

• datagen2.py

1	import random	1817
2		1818
3	def generate_ampl_data(num_tasks=60, num_servers=4,	1819
	output_file='dataset_model3_nov30.dat'):	1820
4	# Define sets	1821
5	tasks = [f't{i}' for i in range(1, num_tasks + 1)]	1822
6	servers = [f's{j}' for j in range(1, num_servers + 1)]	1823
7	resources = ['GPU', 'CPU', 'Memory', 'Storage']	1824
8		1825
9	# Initialize data structures	1826
10	w = {}	1827
11	e = {}	1828
12	p = {}	1829
13	S_time = {}	1830
14	d = {}	1831
15	r = {}	1832
16		1833
17	# Priority weights	1834
18	for task in tasks:	1835
19	w[task] = random.randint(1, 5)	1836
20		1837
21	# Energy consumption rates	1838
22	for server in servers:	1839
23	e[server] = random.randint(5, 15)	1840
24		1841
25	# Resource capacities	1842
26	r = {server: {} for server in servers}	1843
27	for server in servers:	1844
28	for resource in resources:	1845
29	r[server][resource] = random.randint(50, 100)	1846
30		1847
31	# Resource demands	1848
32	d = {task: {} for task in tasks}	1849
33	for task in tasks:	1850
34	for resource in resources:	1851
35	# Generate demands that are significantly less than	1852
	capacities	1853
36	# This helps ensure that tasks can fit on servers	1854
37	max_demand = min(r[server][resource] for server in	1855
	servers) // 4	1856
38	d[task][resource] = random.randint(1, max_demand)	1857
39		1858
40	# Processing times and setup times	1859
41	p = {}	1860
42	S_time = {}	1861
43	for task in tasks:	1862
44	p[task] = {}	1863
45	S_time[task] = {}	1864
46	for server in servers:	1865
47	p[task][server] = random.randint(1, 10)	1866
48	S_time[task][server] = random.randint(0, 3)	1867
49		1868
50	# Feasibility Checks	1869
51	# 1. Total Demand      Total Capacity for each resource	1870

```

52 total_demand = {resource: 0 for resource in resources}
53 total_capacity = {resource: 0 for resource in resources}
54 for resource in resources:
55     total_demand[resource] = sum(d[task][resource] for task in
56         tasks)
57     total_capacity[resource] = sum(r[server][resource] for server
58         in servers)
59     if total_demand[resource] > total_capacity[resource]:
60         print(f"Adjusting demands for resource {resource} to
61             ensure feasibility.")
62         # Scale down demands proportionally
63         scaling_factor = total_capacity[resource] /
64             total_demand[resource]
65         for task in tasks:
66             d[task][resource] = max(1, int(d[task][resource] *
67                 scaling_factor))
68         # Recalculate total demand
69         total_demand[resource] = sum(d[task][resource] for task
70             in tasks)
71
72 # 2. Ensure each task can be assigned to at least one server
73 for task in tasks:
74     assignable = False
75     for server in servers:
76         can_assign = all(d[task][resource] <= r[server][resource]
77             for resource in resources)
78         if can_assign:
79             assignable = True
80             break
81     if not assignable:
82         print(f"Adjusting demands for task {task} to ensure it
83             can be assigned.")
84         # Adjust demands to fit the smallest capacity server
85         for resource in resources:
86             min_capacity = min(r[server][resource] for server in
87                 servers)
88             d[task][resource] = min(d[task][resource],
89                 min_capacity)
90
91 # Write data to file
92 with open(output_file, 'w') as f:
93     f.write("#\n")
94     f.write("-----\n")
95     f.write("# Data File for Task Assignment to Servers Model\n")
96     f.write("# Generated by generate_data.py\n")
97     f.write("#\n")
98     f.write("-----\n\n")
99
100 # Sets
101 f.write("# Sets\n")
102 f.write(f"set_T:= {' '.join(tasks)};\n")
103 f.write(f"set_S:= {' '.join(servers)};\n")
104 f.write(f"set_R:= {' '.join(resources)};\n\n")
105
106 # Parameters
107
108 # Priority weights
109 f.write("# Parameters\n")
110 f.write("# Priority weight of each task\n")
111 f.write("param_w:=\n")
112 for task in tasks:
113     f.write(f" {task} {w[task]}\n")
114 f.write(";\n\n")
115
116 # Energy consumption rates
117 f.write("# Energy consumption rate of each server\n")

```

106	f.write("param_e:=\n")	1939
107	for server in servers:	1940
108	f.write(f"      {server}_{e[server]}\n")	1941
109	f.write(";\n\n")	1942
110		1943
111	# Processing times	1944
112	f.write("#_Processing_time_of_each_task_on_each_server\n")	1945
113	f.write("param_p:=\n")	1946
114	for task in tasks:	1947
115	for server in servers:	1948
116	f.write(f"      {task}_{server}_{p[task][server]}\n")	1949
117	f.write(";\n\n")	1950
118		1951
119	# Setup times	1952
120	f.write("#_Setup_time_between_tasks_on_servers\n")	1953
121	f.write("param_S_time:=\n")	1954
122	for task in tasks:	1955
123	for server in servers:	1956
124	f.write(f"      {task}_{server}_	1957
	{S_time[task][server]}\n")	1958
125	f.write(";\n\n")	1959
126		1960
127	# Resource demands	1961
128	f.write("#_Demand_of_each_resource_by_each_task\n")	1962
129	f.write("param_d:=\n")	1963
130	for task in tasks:	1964
131	for resource in resources:	1965
132	f.write(f"      {task}_{resource}_	1966
	{d[task][resource]}\n")	1967
133	f.write(";\n\n")	1968
134		1969
135	# Resource capacities	1970
136	f.write("#_Capacity_of_each_resource_on_each_server\n")	1971
137	f.write("param_r:=\n")	1972
138	for server in servers:	1973
139	for resource in resources:	1974
140	f.write(f"      {server}_{resource}_	1975
	{r[server][resource]}\n")	1976
141	f.write(";\n\n")	1977
142		1978
143	# Scaling coefficients and large constant	1979
144	f.write("#_Scaling_coefficients_and_large_constant\n")	1980
145	f.write("param_alpha:=0.5;\n")	1981
146	f.write("param_beta:=0.3;\n")	1982
147	f.write("param_M:=1000;\n\n")	1983
148		1984
149	f.write("#_	1985
	-----\n")	1986
150	f.write("#_End_of_Data_File\n")	1987
151	f.write("#_	1988
	-----\n")	1989
152		1990
153	print(f"Data_file_{output_file}_generated_successfully_with_	1991
	{num_tasks}_tasks_and_{num_servers}_servers.")	1992
154		1993
155	if __name__ == "__main__":	1994
156	generate_ampl_data()	1995
		1996