

## **TUGAS AKHIR – EF234801**

# **PENGELOLAAN PENGGUNAAN INFRASTRUKTUR GPU UNTUK PENGGUNA BERBASIS DOCKER CONTAINER MENGUNAKAN JUPYTERLAB**

**Gloriyano Cristho Daniel Pepuho**  
NRP 5025201121

Dosen Pembimbing 1

**Ir. Ary Mazharuddin Shiddiqi, S.Kom., M.Comp.Sc., Ph.D.**  
NIP 19810620 200501 1 003

Dosen Pembimbing 2

**Royyana Muslim Ijtihadie, S.Kom., M.Kom., Ph.D.**  
NIP 19770824 200604 1 001

**Program Studi Strata 1 (S1) Teknik Informatika**

Departemen Teknik Informatika

Fakultas Fakultas Teknologi Elektro dan Informatika Cerdas

Institut Teknologi Sepuluh Nopember

Surabaya

2025



**TUGAS AKHIR – EF234801**

**PENGELOLAAN PENGGUNAAN INFRASTRUKTUR  
GPU UNTUK PENGGUNA BERBASIS DOCKER  
CONTAINER MENGGUNAKAN JUPYTERLAB**

**Gloriyano Cristho Daniel Pepuho**

NRP 5025201121

Dosen Pembimbing 1

**Ir. Ary Mazharuddin Shiddiqi, S.Kom., M.Comp.Sc., Ph.D.**

NIP 19810620 200501 1 003

Dosen Pembimbing 2

**Royyana Muslim Ijtihadie, S.Kom., M.Kom., Ph.D.**

NIP 19770824 200604 1 001

**Program Studi Strata 1 (S1) Teknik Informatika**

Departemen Teknik Informatika

Fakultas Fakultas Teknologi Elektro dan Informatika Cerdas

Institut Teknologi Sepuluh Nopember

Surabaya

2025

*[Halaman ini sengaja dikosongkan]*



**FINAL PROJECT - EF234801**

***Managing Distributed GPU Infrastructure Usage for  
Users Based on Docker Containers Using JupyterLab***

**Gloriyano Cristho Daniel Pepuho**

NRP 5025201121

Advisor

**Ir. Ary Mazharuddin Shiddiqi, S.Kom., M.Comp.Sc., Ph.D.**

NIP 19810620 200501 1 003

**Royyana Muslim Ijtihadie, S.Kom., M.Kom., Ph.D.**

NIP 19770824 200604 1 001

**Undergraduate Study Program of Department of Informatics Engineering**

Department of Department of Informatics Engineering

Faculty of Faculty of Intelligent Electrical and Informatics Technology

Sepuluh Nopember Institute of Technology

Surabaya

2025

*[Halaman ini sengaja dikosongkan]*

# LEMBAR PENGESAHAN

## PENGELOLAAN PENGGUNAAN INFRASTRUKTUR GPU UNTUK PENGGUNA BERBASIS DOCKER CONTAINER MENGGUNAKAN JUPYTERLAB

### TUGAS AKHIR

Diajukan untuk memenuhi salah satu syarat  
memperoleh gelar Sarjana Teknik pada  
Program Studi S-1 Departemen Teknik Informatika  
Departemen Departemen Teknik Informatika  
Fakultas Fakultas Fakultas Teknologi Elektro dan Informatika Cerdas  
Institut Teknologi Sepuluh Nopember

Oleh: **Gloriyano Cristho Daniel Pepuhu**  
NRP. 5025201121

Disetujui oleh Tim Penguji Tugas Akhir:

Ir. Ary Mazharuddin Shiddiqi, S.Kom., M.Comp.Sc., Ph.D.  
NIP: 19810620 200501 1 003

(Pembimbing I)

.....

Royyana Muslim Ijtihadie, S.Kom., M.Kom., Ph.D.  
NIP: 19770824 200604 1 001

(Pembimbing II)

.....

Dr. Galileo Galilei, S.T., M.Sc.  
NIP: 18560710 194301 1 001

(Penguji I)

.....

Friedrich Nietzsche, S.T., M.Sc.  
NIP: 18560710 194301 1 001

(Penguji II)

.....

Alan Turing, ST., MT.  
NIP: 18560710 194301 1 001

(Penguji III)

.....

Mengetahui,  
Kepala Departemen Departemen Teknik Informatika FTEIC - ITS

Prof. Albus Percival Wulfric Brian Dumbledore, S.T., M.T.  
NIP. 18810313 196901 1 001

**SURABAYA**  
**Juni, 2025**

*[Halaman ini sengaja dikosongkan]*

# APPROVAL SHEET

## *Managing Distributed GPU Infrastructure Usage for Users Based on Docker Containers Using JupyterLab*

### FINAL PROJECT

Submitted to fulfill one of the requirements  
for obtaining a degree Bachelor of Engineering at  
Undergraduate Study Program of Department of Informatics Engineering  
Department of Department of Informatics Engineering  
Faculty of Faculty of Intelligent Electrical and Informatics Technology  
Sepuluh Nopember Institute of Technology

By: **Gloriyano Cristho Daniel Pepuho**  
NRP. 5025201121

Approved by Final Project Examiner Team:

Ir. Ary Mazharuddin Shiddiqi, S.Kom., M.Comp.Sc., Ph.D.  
NIP: 19810620 200501 1 003

(Advisor I)

.....

Royyana Muslim Ijtihadie, S.Kom., M.Kom., Ph.D.  
NIP: 19770824 200604 1 001

(Co-Advisor II)

.....

Dr. Galileo Galilei, S.T., M.Sc.  
NIP: 18560710 194301 1 001

(Examiner I)

.....

Friedrich Nietzsche, S.T., M.Sc.  
NIP: 18560710 194301 1 001

(Examiner II)

.....

Alan Turing, ST., MT.  
NIP: 18560710 194301 1 001

(Examiner III)

.....

Acknowledged,  
Head of Department of Informatics Engineering Department FTEIC - ITS

Prof. Albus Percival Wulfric Brian Dumbledore, S.T., M.T.  
NIP. 18810313 196901 1 001

**SURABAYA**  
**June, 2025**



*[Halaman ini sengaja dikosongkan]*

## PERNYATAAN ORISINALITAS

Yang bertanda tangan dibawah ini:

Nama Mahasiswa / NRP : Gloriyano Cristho Daniel Pepuho / 5025201121  
Departemen : Departemen Teknik Informatika  
Dosen Pembimbing / NIP : Ir. Ary Mazharuddin Shiddiqi, S.Kom., M.Comp.Sc., Ph.D. /  
19810620 200501 1 003

Dengan ini menyatakan bahwa Tugas Akhir dengan judul "PENGELOLAAN PENGGUNAAN INFRASTRUKTUR GPU UNTUK PENGGUNA BERBASIS DOCKER CONTAINER MENGGUNAKAN JUPYTERLAB" adalah hasil karya sendiri, berfsifat orisinal, dan ditulis dengan mengikuti kaidah penulisan ilmiah.

Bilamana di kemudian hari ditemukan ketidaksesuaian dengan pernyataan ini, maka saya bersedia menerima sanksi sesuai dengan ketentuan yang berlaku di Institut Teknologi Sepuluh Nopember.

Surabaya, June 2025

Mengetahui  
Dosen Pembimbing

Mahasiswa

Ir. Ary Mazharuddin Shiddiqi, S.Kom., M.Comp.Sc., Ph.D.  
NIP. 19810620 200501 1 003

Gloriyano Cristho Daniel Pepuho  
NRP. 5025201121

*[Halaman ini sengaja dikosongkan]*

## STATEMENT OF ORIGINALITY

The undersigned below:

Name of student / NRP : Gloriyano Cristho Daniel Pepuho / 5025201121  
Department : Department of Informatics Engineering  
Advisor / NIP : Ir. Ary Mazharuddin Shiddiqi, S.Kom., M.Comp.Sc., Ph.D. /  
19810620 200501 1 003

Hereby declared that the Final Project with the title of "*Managing Distributed GPU Infrastructure Usage for Users Based on Docker Containers Using JupyterLab*" is the result of my own work, is original, and is written by following the rules of scientific writing.

If in future there is a discrepancy with this statement, then I am willing to accept sanctions in accordance with provisions that apply at Sepuluh Nopember Institute of Technology.

Surabaya, June 2025

Acknowledged  
Advisor

Student

Ir. Ary Mazharuddin Shiddiqi, S.Kom., M.Comp.Sc., Ph.D.      Gloriyano Cristho Daniel Pepuho  
NIP. 19810620 200501 1 003      NRP. 5025201121

*[Halaman ini sengaja dikosongkan]*

## ABSTRAK

Nama Mahasiswa : Gloriyano Cristho Daniel Pepuho  
Judul Tugas Akhir : PENGELOLAAN PENGGUNAAN INFRASTRUKTUR GPU UNTUK PENGGUNA BERBASIS DOCKER CONTAINER MENGGUNAKAN JUPYTERLAB  
Pembimbing : 1. Ir. Ary Mazharuddin Shiddiqi, S.Kom., M.Comp.Sc., Ph.D.  
2. Royyana Muslim Ijtihadie, S.Kom., M.Kom., Ph.D.

Dalam era teknologi yang semakin maju, kebutuhan akan komputasi berbasis GPU menjadi sangat penting, khususnya dalam bidang kecerdasan buatan (AI) dan analisis data skala besar. GPU memungkinkan pemrosesan paralel yang cepat dan efisien, sehingga sering digunakan untuk melatih model deep learning dan menjalankan tugas-tugas komputasi intensif. Namun, pengelolaan GPU di lingkungan multi-pengguna menghadapi tantangan besar, seperti alokasi sumber daya yang tidak merata dan potensi penurunan efisiensi sistem. Untuk mengatasi masalah ini, penelitian ini bertujuan untuk mengembangkan mekanisme penjadwalan GPU yang efisien dengan memanfaatkan teknologi Docker Container dan antarmuka JupyterLab. Docker digunakan untuk menciptakan lingkungan kerja yang terisolasi bagi setiap pengguna, sementara JupyterLab menyediakan platform interaktif yang memudahkan pengguna dalam mengakses dan menjalankan tugas berbasis GPU secara simultan. Penelitian ini dibagi kedalam beberapa tahap yang meliputi analisis kebutuhan, desain sistem, serta perancangan metode evaluasi. Rancangan sistem yang diusulkan akan diimplementasikan pada kluster GPU di lingkungan laboratorium atau institusi pendidikan. Evaluasi direncanakan mencakup pengujian efisiensi alokasi sumber daya, kemudahan akses pengguna, dan skalabilitas sistem dalam mendukung banyak pengguna secara bersamaan. Penelitian ini diharapkan dapat memberikan kontribusi terhadap pengelolaan sumber daya GPU dalam lingkungan komputasi terdistribusi, mendukung efisiensi dan keadilan alokasi, serta meningkatkan pengalaman pengguna dalam mengakses sumber daya GPU untuk kebutuhan komputasi modern.

**Kata Kunci:** *Kluster GPU, Docker Container, JupyterLab, Pengelolaan pengguna*

*[Halaman ini sengaja dikosongkan]*

## ABSTRACT

*Name* : Gloriyano Cristho Daniel Pepuho  
*Title* : *Managing Distributed GPU Infrastructure Usage for Users Based on Docker Containers Using JupyterLab*  
*Advisors* : 1. Ir. Ary Mazharuddin Shiddiqi, S.Kom., M.Comp.Sc., Ph.D.  
2. Royyana Muslim Ijtihadie, S.Kom., M.Kom., Ph.D.

In the era of advancing technology, the demand for GPU-based computing has become increasingly critical, particularly in the fields of artificial intelligence (AI) and large-scale data analysis. GPUs enable fast and efficient parallel processing, making them widely used for training deep learning models and performing computationally intensive tasks. However, managing GPUs in multi-user environments presents significant challenges, such as uneven resource allocation and potential system inefficiencies. To address these issues, this study aims to develop an efficient GPU scheduling mechanism utilizing Docker container technology and the JupyterLab interface. Docker creates isolated work environments for each user, while JupyterLab provides an interactive platform that simplifies simultaneous GPU-based task execution. The research consists of several phases, including requirement analysis, system design, and evaluation method planning. The proposed system design will be implemented on a GPU cluster in a laboratory or educational institution environment. Evaluation will include testing resource allocation efficiency, user accessibility, and system scalability in supporting multiple concurrent users. This study is expected to make a significant contribution to GPU resource management in distributed computing environments, promoting efficiency and fairness in resource allocation while enhancing the user experience in accessing GPU resources for modern computational needs.

***Keywords:*** *GPU Cluster, Docker Container, JupyterLab, User Management*



*[Halaman ini sengaja dikosongkan]*

## KATA PENGANTAR

Puji dan syukur kehadiran Tuhan Yang Maha Esa yang memberikan karunia, rahmat, dan pertolongan sehingga penulis dapat menyelesaikan penelitian tugas akhir yang berjudul 'PENGELOLAAN PENGGUNAAN INFRASTRUKTUR GPU UNTUK PENGGUNA BERBASIS DOCKER CONTAINER MENGGUNAKAN JUPYTERLAB'. Melalui kata pengantar ini, penulis mengucapkan terima kasih sebesar-besarnya kepada seluruh pihak yang telah membantu dan mendukung penulis selama mengerjakan penelitian tugas akhir ini, diantaranya adalah:

1. Tuhan Yang Maha Esa, atas karunia dan rahmat-Nya sehingga penulis dapat mencapai titik akhir perkuliahan strata satu di Departemen Teknik Informatika, Institut Teknologi Sepuluh Nopember.
2. Kedua orang tua yang telah mendukung penulis selama berkuliah di Departemen Teknik Informatika, Institut Teknologi Sepuluh Nopember.
3. Bapak Ir. Ary Mazharuddin Shiddiqi, S.Kom., M.Comp.Sc., Ph.D. dan Bapak Royyana Muslim Ijtihadie, S.Kom., M.Kom., Ph.D. sebagai dosen pembimbing yang telah membimbing, memberi arahan, dan masukan kepada penulis selama mengerjakan tugas akhir ini.
4. Dosen dan tenaga pendidik di Departemen Teknik Informatika, Institut Teknologi Sepuluh Nopember yang telah memberikan pengetahuan, wawasan, dan pengalaman yang sangat berarti selama masa studi.
5. Pihak-pihak lain yang tidak dapat disebutkan satu persatu yang telah membantu penulis dalam pelaksanaan penelitian tugas akhir ini.

Akhir kata, semoga penelitian tugas akhir ini dapat memberikan kontribusi yang bermanfaat. Terima kasih dan permohonan maaf atas kekurangan dan kesalahan dalam pelaksanaan tugas akhir ini.

Surabaya, Juni 2025

Gloriyano Cristho Daniel Pepuho

*[Halaman ini sengaja dikosongkan]*

# DAFTAR ISI

<b>ABSTRAK</b>	<b>i</b>
<b>ABSTRACT</b>	<b>iii</b>
<b>KATA PENGANTAR</b>	<b>v</b>
<b>DAFTAR ISI</b>	<b>vii</b>
<b>DAFTAR GAMBAR</b>	<b>ix</b>
<b>DAFTAR TABEL</b>	<b>xi</b>
<b>1 PENDAHULUAN</b>	<b>1</b>
1.1 Latar Belakang . . . . .	1
1.2 Rumusan Masalah . . . . .	2
1.3 Batasan Masalah atau Ruang Lingkup . . . . .	2
1.4 Tujuan . . . . .	2
1.5 Manfaat . . . . .	2
1.6 Sistematika Penulisan . . . . .	3
<b>2 TINJAUAN PUSTAKA</b>	<b>5</b>
2.1 Hasil penelitian/perancangan terdahulu . . . . .	5
2.1.1 Containerisation for High Performance Computing Systems: Survey and Prospects . . . . .	5
2.1.2 An accessible infrastructure for artificial intelligence using a Docker-based JupyterLab in Galaxy . . . . .	5
2.1.3 Syndeo: Portable Ray Clusters with Secure Containerization . . . . .	6
2.1.4 Ray: A Distributed Framework for Emerging AI Applications . . . . .	6
2.2 Teori/Konsep Dasar . . . . .	7
2.2.1 Klaster GPU . . . . .	7
2.2.2 Docker . . . . .	7
2.2.3 Penjadwalan GPU dengan Ray . . . . .	9
2.2.4 JupyterLab . . . . .	9
2.2.5 JupyterHub . . . . .	9

2.2.6	Ray Framework . . . . .	11
<b>3</b>	<b>DESAIN DAN IMPLEMENTASI</b>	<b>13</b>
3.1	Perancangan Arsitektur Sistem . . . . .	13
3.1.1	Service Discovery . . . . .	13
3.1.2	JupyterHub . . . . .	14
3.1.3	Ray Cluster . . . . .	15
3.2	Implementasi Sistem . . . . .	15
3.2.1	Discovery Service . . . . .	15
3.2.2	Agent Service . . . . .	20
3.2.3	Implementasi JupyterHub . . . . .	23
3.3	Peralatan Pendukung . . . . .	25
<b>4</b>	<b>PENGUJIAN DAN ANALISIS</b>	<b>29</b>
4.1	Skenario Pengujian . . . . .	29
4.1.1	Skenario 1: Pemilihan Node dengan Beban Terendah . . . . .	29
4.1.2	Skenario 2: Multi-User Concurrent . . . . .	29
4.1.3	Skenario 3: Simulasi Beban Tinggi . . . . .	29
4.1.4	Skenario 4: Validasi Profil GPU . . . . .	29
4.1.5	Skenario 5: TTL Redis dan Node Tidak Aktif . . . . .	30
4.2	Evaluasi Pengujian . . . . .	30
<b>5</b>	<b>PENUTUP</b>	<b>31</b>
5.1	Kesimpulan . . . . .	31
5.2	Saran . . . . .	31
	<b>DAFTAR PUSTAKA</b>	<b>33</b>
	<b>BIOGRAFI PENULIS</b>	<b>35</b>

## DAFTAR GAMBAR

2.1	Arsitektur Docker (Sumber: Team, 2024) . . . . .	8
2.2	Komponen <i>RAY</i> (Sumber: Moritz et al., 2018) . . . . .	12

*[Halaman ini sengaja dikosongkan]*

## DAFTAR TABEL

3.1	Struktur Direktori Proyek Discovery Service . . . . .	15
3.2	Struktur Direktori Proyek JupyterHub . . . . .	24
3.3	Spesifikasi Peralatan Pendukung . . . . .	26
3.4	Daftar Perangkat Lunak Pendukung . . . . .	27
4.1	Ringkasan Evaluasi Pengujian Sistem . . . . .	30



*[Halaman ini sengaja dikosongkan]*

# BAB I

## PENDAHULUAN

### 1.1 Latar Belakang

GPU telah menjadi elemen krusial dalam komputasi modern, dengan penggunaan GPU untuk deep learning workloads meningkat secara eksponensial dalam dekade terakhir. Artikel "Deep Learning Workload Scheduling in GPU Datacenters: A Survey" mengidentifikasi bahwa traditional approaches yang dirancang untuk big data atau HPC workloads tidak dapat mendukung deep learning workloads untuk fully utilize GPU resources, sehingga memerlukan pendekatan scheduling yang khusus dirancang untuk karakteristik unik dari AI workloads.

Teknologi seperti Docker Container telah menjadi solusi inovatif untuk meningkatkan efisiensi dalam pengelolaan aplikasi, terutama di lingkungan komputasi modern. Dengan mengisolasi aplikasi dan dependensinya, *container* memungkinkan penyebaran yang cepat dan konsisten. Artikel "*Containerisation for High Performance Computing Systems: Survey and Prospects*" menjelaskan bahwa teknologi kontainerisasi tidak hanya relevan dalam komputasi awan, tetapi juga memiliki potensi besar untuk sistem *High Performance Computing (HPC)*. Dalam konteks *HPC*, kontainer mampu mengemas pustaka yang dioptimalkan untuk perangkat keras tertentu, meskipun tantangan seperti ukuran yang besar dan kebutuhan orkestrasi yang kompleks tetap perlu diatasi. Penggunaan mekanisme orkestrasi yang efisien dapat membantu memaksimalkan pemanfaatan sumber daya GPU dalam lingkungan *HPC* yang terdistribusi.

Namun, tantangan seperti ukuran kontainer yang besar dan kebutuhan akan mekanisme orkestrasi yang kompleks tetap perlu diatasi. Artikel ini juga menyoroti bahwa penggunaan teknologi orkestrasi yang efisien, seperti Kubernetes, dapat membantu memaksimalkan pemanfaatan sumber daya GPU dalam lingkungan *HPC* yang terdistribusi, menjadikan kontainerisasi solusi yang menjanjikan untuk pengelolaan sumber daya multi-pengguna (Zhou et al., 2022).

Integrasi teknologi kontainer seperti Docker dengan JupyterLab telah membuka peluang baru dalam pengelolaan infrastruktur komputasi berbasis GPU untuk proyek kecerdasan buatan (AI). Artikel "*An accessible infrastructure for artificial intelligence using a Docker-based JupyterLab in Galaxy*" menunjukkan bahwa pendekatan berbasis kontainer ini dapat menyediakan lingkungan komputasi yang terisolasi namun fleksibel, memungkinkan *training model deep learning* yang cepat dan aman melalui akses GPU yang teroptimalkan. Selain itu, JupyterLab yang berjalan di atas Docker mendukung kolaborasi antar peneliti melalui *notebook* interaktif yang mudah diakses. Infrastruktur ini tidak hanya meningkatkan efisiensi penggunaan GPU tetapi juga memungkinkan pengelolaan sumber daya secara lebih adil, sebagaimana diilustrasikan dalam studi kasus pelatihan model untuk analisis gambar dan prediksi struktur protein. Dengan menggunakan pendekatan serupa, penelitian ini bertujuan untuk mengembangkan mekanisme penjadwalan GPU yang efektif dalam lingkungan terdistribusi, guna mendukung kebutuhan komputasi AI modern yang semakin kompleks dan kolaboratif.

Mekanisme penjadwalan pengguna pada lingkungan GPU terdistribusi menjadi aspek krusial dalam mendukung efisiensi dan keadilan alokasi sumber daya. Dalam konteks ini, JupyterLab menawarkan antarmuka interaktif yang memungkinkan pengguna menjalankan tugas se-

cara bersamaan dengan akses GPU yang terorkestrasi oleh Docker Container. Hal ini relevan dalam skenario penelitian atau pengembangan model AI oleh kami yang memanfaatkan sumber daya GPU secara kolektif. Dengan demikian, penelitian ini bertujuan untuk mengembangkan solusi penjadwalan yang tidak hanya meningkatkan efisiensi, tetapi juga memaksimalkan pengalaman pengguna dalam mengakses sumber daya GPU pada kluster terdistribusi.

## 1.2 Rumusan Masalah

Rumusan masalah yang diangkat dalam tugas akhir ini adalah sebagai berikut:

1. Bagaimana cara mengelola penggunaan GPU agar dapat digunakan secara adil dan efisien oleh banyak pengguna?
2. Bagaimana cara memanfaatkan kontainerisasi untuk mengisolasi lingkungan kerja setiap pengguna dan meningkatkan efisiensi serta skalabilitas?
3. Apakah sistem yang diusulkan dapat mengakomodasi kebutuhan penggunaan GPU oleh banyak pengguna secara bersamaan?
4. Bagaimana cara mempermudah pengguna dalam mengakses dan memanfaatkan GPU melalui antarmuka yang interaktif?

## 1.3 Batasan Masalah atau Ruang Lingkup

Batasan dalam pengerjaan tugas akhir ini adalah sebagai berikut:

1. Infrastruktur GPU yang digunakan adalah infrastruktur berbasis kluster komputer yang tersedia pada laboratorium atau institusi pendidikan, dengan konfigurasi spesifik seperti *node* berbasis Docker.
2. Implementasi sistem difokuskan pada integrasi Docker dengan JupyterLab untuk orkestrasi akses GPU.

## 1.4 Tujuan

Tujuan dari pembuatan Tugas Akhir ini adalah sebagai berikut:

1. Mengembangkan sistem yang mampu meningkatkan efisiensi penggunaan GPU, khususnya dalam konteks lingkungan multi-pengguna.
2. Menyediakan antarmuka berbasis JupyterLab yang memungkinkan akses GPU dengan mudah, interaktif, dan terintegrasi dengan baik.

## 1.5 Manfaat

Manfaat dari penelitian ini adalah sebagai berikut:

1. Sistem ini dapat meningkatkan efisiensi pemanfaatan infrastruktur GPU yang terbatas, mendukung penelitian, dan pengembangan berbasis komputasi AI.
2. Sistem ini memberikan akses GPU yang lebih mudah, terstruktur, dan aman, sehingga mendukung produktivitas dalam pengembangan aplikasi berbasis GPU.

## 1.6 Sistematika Penulisan

Laporan penelitian tugas akhir ini terbagi menjadi Lorem ipsum dolor sit amet, consectetur adipiscing elit. Ut purus elit, vestibulum ut, placerat ac, adipiscing vitae, felis. Curabitur dictum gravida mauris. yaitu:

### 1. **BAB I Pendahuluan**

Bab ini berisi latar belakang penelitian yang menjelaskan pentingnya pengelolaan infrastruktur GPU terdistribusi. rumusan masalah yang dihadapi dalam penggunaan GPU multi-user, batasan masalah dan ruang lingkup penelitian, tujuan yang ingin dicapai, manfaat penelitian, serta sistematika penulisan laporan.

### 2. **BAB II Tinjauan Pustaka**

Bab ini berisi tinjauan terhadap penelitian-penelitian terdahulu yang relevan dengan topik penelitian, teori dan konsep dasar yang meliputi kluster GPU, teknologi Docker, Ray Framework, penjadwalan GPU, JupyterLab, dan JupyterHub. Bab ini menjadi landasan teoritis untuk melakukan pengembangan sistem.

### 3. **BAB III Desain dan Implementasi Sistem**

Bab ini berisi perancangan arsitektur sistem yang mencakup service discovery, integrasi JupyterHub, dan konfigurasi Ray cluster. Selain itu bab ini membahas peralatan apa saja yang digunakan pada saat penelitian serta setiap detail implementasi komponen yang dikembangkan.

### 4. **BAB IV Pengujian dan Analisa**

Bab ini berisi bab ini dirancang untuk memvalidasi fungsionalitas sistem, evaluasi perform dalam berbagai kondisi beban, analisis efisiensi penggunaan resource, serta pembahasan hasil pengujian terhadap tujuan penelitian yang telah ditetapkan

### 5. **BAB V Penutup**

Bab ini berisi kesimpulan dari penelitian yang merangkum pencapaian tujuan penelitian, kontribusi yang diberikan, serta saran untuk pengembangan dan penelitian lebih lanjut yang dapat dilakukan berdasarkan penelitian ini.

*[Halaman ini sengaja dikosongkan]*

## **BAB II**

### **TINJAUAN PUSTAKA**

#### **2.1 Hasil penelitian/perancangan terdahulu**

Dalam melakukan penelitian ini, penulis akan menggunakan beberapa penelitian terdahulu sebagai pedoman dan referensi dalam mengerjakan tugas akhir ini.

##### **2.1.1 Containerisation for High Performance Computing Systems: Survey and Prospects**

Pada artikel ini, peneliti melakukan survei tentang penggunaan *container* dalam sistem *High Performance Computing (HPC)*. Fokus utama adalah bagaimana *container*, seperti Docker, dapat meningkatkan portabilitas, efisiensi, dan isolasi lingkungan di *HPC*. Artikel ini juga mengkaji kelebihan *container* dibandingkan *virtual machine*, termasuk *overhead* yang lebih rendah dan waktu *startup* yang lebih cepat. Survei ini dilakukan dengan mengkaji literatur dari berbagai penelitian terbaru tentang penggunaan *container* di sistem *HPC*, termasuk studi kasus implementasi *container* dalam kluster GPU dan simulasi berbasis AI.

Hasil survei menunjukkan bahwa *container* memungkinkan *workload HPC* dijalankan di berbagai platform dengan efisiensi yang tinggi, menjadikannya solusi populer untuk komputasi terdistribusi. Selain itu, peneliti juga mengidentifikasi tantangan seperti integrasi dengan sistem manajemen kluster dan penjadwalan yang optimal. Kontribusi utama dari artikel ini adalah menyajikan analisis perbandingan yang mendalam antara *container* dan *virtual machine* dalam konteks *HPC*, serta menyarankan pendekatan penjadwalan yang lebih optimal untuk *container* di lingkungan *multi-user*.

Temuan ini relevan dengan penelitian ini, terutama dalam konteks penggunaan Docker untuk mengelola pengguna pada kluster GPU terdistribusi. Konsep efisiensi yang diangkat dalam artikel ini memberikan dasar teoritis untuk pengembangan mekanisme penjadwalan GPU yang akan digunakan dalam penelitian ini, terutama dalam hal mengurangi *overhead* dan memastikan alokasi sumber daya yang adil. Selain itu, contoh aplikasi *container* di sistem *HPC* yang disebutkan dalam artikel ini memberikan inspirasi untuk implementasi praktis dalam pengelolaan lingkungan kerja berbasis *container* (Zhou et al., 2022).

##### **2.1.2 An accessible infrastructure for artificial intelligence using a Docker-based JupyterLab in Galaxy**

Pada artikel ini, peneliti mengembangkan infrastruktur yang dapat diakses untuk kecerdasan buatan dengan memanfaatkan JupyterLab berbasis Docker di dalam platform *Galaxy*. Infrastruktur ini dirancang untuk mempermudah pengguna dalam mengakses alat komputasi AI melalui antarmuka berbasis web.

Hasilnya menunjukkan bahwa pendekatan ini meningkatkan portabilitas dan aksesibilitas bagi pengguna. Penggunaan *container* Docker memungkinkan pengelolaan lingkungan komputasi yang konsisten dan meminimalkan konfigurasi manual. Artikel ini juga menyoroti man-

faat JupyterLab dalam menyediakan antarmuka yang intuitif bagi pengguna. Temuan ini relevan dengan penelitian ini, terutama dalam konteks penggunaan JupyterLab berbasis Docker untuk mengelola sumber daya komputasi GPU. Artikel ini memberikan wawasan tentang bagaimana desain antarmuka berbasis *container* dapat meningkatkan efisiensi dan aksesibilitas sistem (Kumar et al., 2023).

### 2.1.3 Syndeo: Portable Ray Clusters with Secure Containerization

Pada paper ini, peneliti memperkenalkan *Syndeo*, sebuah framework untuk mengelola dan mengorkestrasi cluster RAY secara portable menggunakan *container*. Fokus utama dari penelitian ini adalah bagaimana *Syndeo* dapat memanfaatkan kontainerisasi untuk meningkatkan portabilitas, keamanan, dan skalabilitas dalam menjalankan *workload* RAY di berbagai platform cloud, seperti AWS, Azure, dan Google Cloud. Framework ini dirancang untuk mendukung komputasi *throughput* tinggi *multi-node* dan memastikan keamanan dengan membatasi hak istimewa pengguna, sehingga administrator memiliki kontrol penuh atas akses sistem. *Syndeo* juga memungkinkan implementasi *workflow paralell* Ray pada sistem manajemen kluster seperti *Slurm*, yang sebelumnya tidak didukung secara *native*.

Temuan dalam paper ini relevan dengan penelitian ini, terutama dalam konteks penggunaan Docker dan Ray untuk mengelola sumber daya GPU dalam kluster terdistribusi. *Syndeo* memberikan wawasan tentang pentingnya portabilitas, keamanan, dan orkestrasi yang efisien dalam lingkungan multi-pengguna, yang dapat menjadi inspirasi untuk pengelolaan pengguna dan alokasi sumber daya dalam sistem berbasis Kluster yang digunakan pada penelitian ini (Li et al., 2024).

### 2.1.4 Ray: A Distributed Framework for Emerging AI Applications

Dalam paper ini, peneliti memperkenalkan Ray, sebuah framework terdistribusi yang dirancang untuk mendukung aplikasi AI modern, seperti *reinforcement learning* dan *deep learning*. Framework ini menawarkan antarmuka terpadu yang mampu mengekspresikan komputasi berbasis tugas (*task-parallel*) dan aktor (*actor-based*), didukung oleh mesin eksekusi dinamis tunggal. Ray mengimplementasikan penjadwalan terdistribusi dan penyimpanan yang toleran terhadap kesalahan untuk mengelola status kontrol sistem. Eksperimen yang dilakukan menunjukkan bahwa Ray dapat menskalakan hingga lebih dari 1,8 juta tugas per detik dan memberikan kinerja yang lebih baik dibandingkan sistem khusus lainnya.

Temuan dari paper ini relevan dengan penelitian ini dalam beberapa aspek. Pertama, Ray sebagai framework terdistribusi untuk aplikasi AI sesuai dengan kebutuhan penelitian untuk memanfaatkan teknologi tersebut dalam pengelolaan infrastruktur GPU terdistribusi. Kedua, kemampuan Ray dalam mendukung model komputasi *task-parallel* dan *actor-based* memberikan fleksibilitas yang diperlukan dalam penjadwalan dan alokasi sumber daya GPU di lingkungan multi-pengguna. Ketiga, fitur penjadwalan terdistribusi dan toleransi kesalahan pada Ray dapat meningkatkan efisiensi dan keandalan sistem yang dikembangkan. Keempat, skalabilitas tinggi Ray, yang mampu menangani jutaan tugas per detik, relevan untuk mendukung penggunaan GPU oleh banyak pengguna secara simultan (Moritz et al., 2018).

## 2.2 Teori/Konsep Dasar

### 2.2.1 Klaster GPU

Klaster GPU adalah kumpulan unit pemrosesan grafis (GPU) yang terhubung dalam satu sistem untuk mendukung komputasi paralel intensif. Klaster ini sering digunakan untuk mempercepat pemrosesan aplikasi dengan kebutuhan komputasi tinggi, seperti pelatihan model *deep learning* dan simulasi ilmiah. Efisiensi komputasi dicapai dengan membagi beban kerja antar GPU secara terdistribusi. Setiap GPU bekerja secara paralel untuk menyelesaikan bagian tertentu dari tugas besar, memungkinkan pengurangan waktu pemrosesan dan penggunaan sumber daya secara optimal. Manajemen sumber daya yang baik diperlukan agar alokasi beban kerja berjalan efisien dan terkoordinasi. (Shikai Wang and Shang, 2024).

### 2.2.2 Docker

Docker merupakan *tool open-source* yang mengotomatisasi proses penyebaran aplikasi ke dalam wadah (*container*). Docker dikembangkan oleh tim di Docker, Inc (sebelumnya dikenal sebagai dotCloud Inc), salah satu pelopor di pasar *Platform-as-a-Service* atau (PAAS), dan dirilis di bawah lisensi Apache 2.0. Apa yang membuat Docker istimewa? Docker menyediakan platform untuk penyebaran aplikasi yang dibangun di atas lingkungan eksekusi *container* yang tervirtualisasi. Teknologi ini dirancang untuk menghadirkan lingkungan yang ringan dan cepat bagi pengembangan serta eksekusi aplikasi, sekaligus menyederhanakan alur kerja distribusi kode—mulai dari perangkat pengembang, lingkungan pengujian, hingga tahap produksi. Dengan kemudahan yang ditawarkannya, Docker memungkinkan pengguna memulai hanya dengan host minimal yang memiliki kernel Linux yang kompatibel dan biner Docker (Turnbull, 2014).

Docker memiliki beberapa komponen penting, seperti berikut:

- **Docker Client**  
Docker Client adalah antarmuka utama yang digunakan oleh pengguna untuk berinteraksi dengan Docker. Melalui Docker Client, pengguna dapat mengirim perintah seperti membangun, mendistribusikan, dan menjalankan *container*. Perintah-perintah ini kemudian diteruskan ke Docker Daemon untuk diproses. Docker Client mendukung penggunaan antarmuka command line (CLI) yang intuitif, sehingga memudahkan pengelolaan infrastruktur container.
- **Docker Daemon**  
Docker Daemon adalah proses latar belakang yang bertanggung jawab untuk menangani perintah yang diterima dari Docker Client. Fungsinya meliputi pembuatan dan pengelolaan berbagai objek Docker, seperti *images*, *containers*, *networks*, dan *volumes*. Docker Daemon memastikan *container* berjalan dengan stabil dan memonitor aktivitasnya. Ia juga berperan penting dalam komunikasi dengan *registry* untuk *push* atau *pull* Docker images.
- **Docker Container**  
Docker Container adalah unit eksekusi yang ringan dan mandiri. *Container* ini berisi semua komponen yang diperlukan untuk menjalankan aplikasi, termasuk kode aplikasi, pustaka, dependensi, dan konfigurasi. Karena sifatnya yang terisolasi, *container* memberikan lingkungan konsisten untuk aplikasi, terlepas dari perbedaan konfigurasi sistem di berbagai host.
- **Docker Images**  
Docker Images adalah *template read-only* yang menjadi dasar untuk membangun Docker



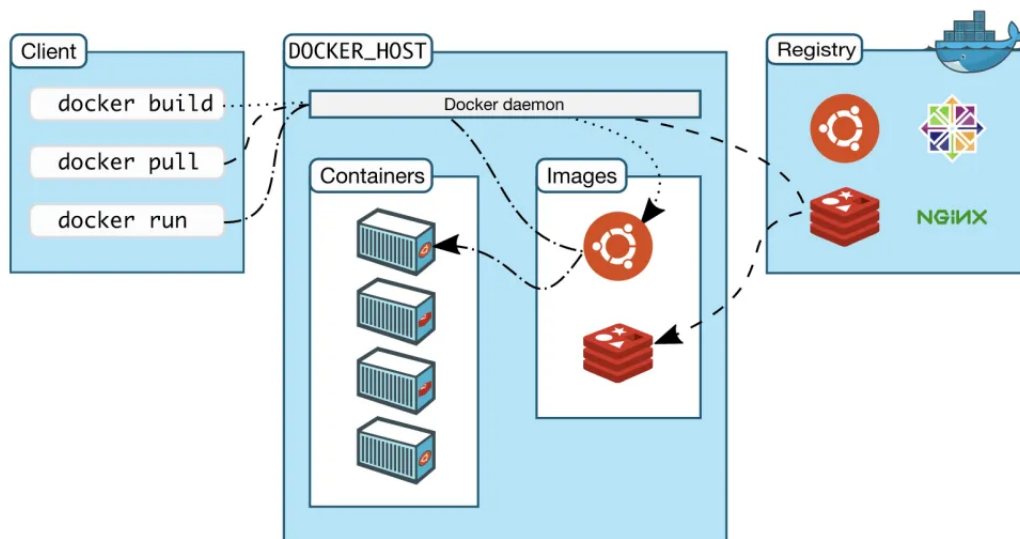
Container. Image ini mencakup semua dependensi, pustaka, dan file yang diperlukan untuk menjalankan aplikasi dalam container. Pengguna dapat membuat *images* dari *Dockerfile* atau *pull images* yang sudah ada dari Docker Hub atau *registry* lainnya. *Images* bersifat modular dan dapat *diupdate* atau digunakan kembali untuk berbagai kebutuhan.

- Registry

Registry adalah layanan penyimpanan dan distribusi Docker Images. Docker menyediakan *registry* publik seperti Docker Hub, tempat pengguna dapat mengunggah, menyimpan, dan berbagi *images*. Selain itu, pengguna juga dapat mengatur *registry* privat untuk kebutuhan spesifik organisasi. Registry mempermudah pengelolaan *images* dalam pengembangan kolaboratif dan siklus hidup *container*.

Arsitektur Docker dapat dilihat pada gambar 2.2. Dalam arsitektur ini, Docker Client berfungsi sebagai jembatan antara pengguna dan sistem Docker, di mana setiap perintah yang dikirimkan oleh pengguna akan diteruskan ke Docker Daemon yang berjalan pada sistem Docker Host.

Docker Daemon kemudian akan menjalankan proses yang dibutuhkan, mulai dari menarik *image (pull)* dari Docker Registry, membangun *container* dari *image* tersebut, hingga menjalankan dan mengelola siklus hidup *container*. Docker Registry sendiri berperan sebagai tempat penyimpanan dan distribusi Docker Image, baik melalui *registry* publik seperti Docker Hub, maupun *registry* privat yang disiapkan secara internal.



Gambar 2.1: Arsitektur Docker (Sumber: Team, 2024)

### 2.2.3 Penjadwalan GPU dengan Ray

Penjadwalan GPU merupakan komponen penting dalam sistem komputasi terdistribusi, terutama ketika sumber daya GPU terbatas harus dibagi ke banyak pengguna atau tugas. Ray menyediakan mekanisme penjadwalan tugas berbasis sumber daya yang fleksibel dan dinamis untuk mengelola alokasi GPU secara efisien dalam lingkungan multi-pengguna.

### 2.2.4 JupyterLab

JupyterLab adalah antarmuka pengguna berbasis web untuk Project Jupyter yang menyediakan lingkungan pengembangan interaktif yang fleksibel dan modular. JupyterLab memungkinkan pengguna untuk bekerja dengan *notebook*, file, *terminal*, dan editor teks dalam satu antarmuka terpadu yang dapat disesuaikan.

JupyterLab berperan sebagai antarmuka utama yang memungkinkan pengguna mengakses sumber daya GPU secara interaktif. Setiap pengguna akan mendapatkan instance JupyterLab yang berjalan dalam container Docker terisolasi, memberikan lingkungan kerja yang konsisten dan aman. Integrasi dengan Ray framework memungkinkan pengguna menjalankan komputasi terdistribusi langsung dari notebook tanpa konfigurasi manual yang kompleks.

JupyterLab sendiri dipilih karena kemudahannya dalam lingkungan multi-pengguna dan kompatibilitasnya dengan container Docker, sehingga cocok untuk implementasi sistem penjadwalan GPU yang diusulkan dalam penelitian ini.

#### Komponen Utama JupyterLab:

1. **Notebook Interface:** Menyediakan lingkungan interaktif untuk menjalankan kode Python, R, atau bahasa pemrograman lainnya dengan dukungan visualisasi data yang kaya.
2. **File Browser:** Memungkinkan navigasi dan manajemen file dalam sistem, termasuk upload dan download file secara langsung melalui antarmuka web.
3. **Extension System:** Akses terminal penuh yang terintegrasi dalam antarmuka web, memungkinkan eksekusi perintah sistem langsung dari browser.
4. **Terminal:** Mendukung plugin dan ekstensi untuk memperluas fungsionalitas sesuai kebutuhan spesifik pengguna.

### 2.2.5 JupyterHub

JupyterHub adalah platform open-source yang memungkinkan banyak pengguna untuk mengakses dan menjalankan lingkungan Jupyter Notebook secara terisolasi melalui antarmuka web. Dirancang untuk mendukung skenario multi-pengguna, JupyterHub sangat cocok digunakan dalam lingkungan pendidikan, penelitian, dan industri yang memerlukan akses bersama ke sumber daya komputasi.

#### Arsitektur Utama JupyterHub

JupyterHub terdiri dari tiga komponen utama yang bekerja secara sinergis:

1. **Hub:** Komponen inti yang bertanggung jawab atas manajemen akun pengguna, proses autentikasi, dan koordinasi peluncuran server notebook individu melalui mekanisme yang disebut Spawner.
2. **Proxy:** Berfungsi sebagai gerbang utama yang menerima semua permintaan HTTP dari pengguna dan meneruskannya ke Hub atau server notebook pengguna yang sesuai. Secara default, JupyterHub menggunakan configurable-http-proxy yang dibangun di atas node-http-proxy.
3. **Single-User Notebook:** Server Jupyter Notebook yang dijalankan secara terpisah untuk setiap pengguna setelah proses autentikasi berhasil. Server ini memungkinkan pengguna untuk menjalankan kode dan berinteraksi dengan lingkungan Jupyter secara pribadi.

### Alur Kerja JupyterHub

Proses interaksi pengguna dengan JupyterHub dapat dijelaskan sebagai berikut:

1. Akses Awal: Pengguna mengakses JupyterHub melalui browser web dengan mengunjungi alamat IP atau nama domain yang telah dikonfigurasi.
2. Autentikasi: Data login yang dimasukkan oleh pengguna dikirim ke komponen Authenticator untuk validasi. Jika valid, pengguna akan dikenali dan diizinkan untuk melanjutkan.
3. Peluncuran Server Notebook: Setelah autentikasi berhasil, JupyterHub akan meluncurkan instance server notebook khusus untuk pengguna tersebut menggunakan Spawner.
4. Konfigurasi Proxy: Proxy dikonfigurasi untuk meneruskan permintaan dengan URL tertentu (misalnya, `/user/[username]`) ke server notebook pengguna yang sesuai.
5. Penggunaan Lingkungan Jupyter: Pengguna diarahkan ke server notebook pribadi mereka, di mana mereka dapat mulai bekerja dengan lingkungan Jupyter seperti biasa.

### Melakukan kustomisasi dan menambah extension

JupyterHub dirancang dengan fleksibilitas tinggi, memungkinkan kustomisasi melalui dua komponen utama:

- **Authenticator:** Mengelola proses autentikasi pengguna. JupyterHub mendukung berbagai metode autentikasi, termasuk PAM (Pluggable Authentication Modules), OAuth, LDAP, dan lainnya.
- **Spawner:** Mengontrol cara peluncuran server notebook untuk setiap pengguna. Beberapa jenis Spawner yang umum digunakan antara lain:
  - **LocalProcessSpawner:** Meluncurkan server notebook sebagai proses lokal di mesin yang sama.
  - **DockerSpawner:** Menjalankan server notebook dalam container Docker, memberikan isolasi lingkungan yang lebih baik.
  - **KubeSpawner:** Menggunakan Kubernetes untuk mengelola dan menskalakan server notebook di lingkungan kluster.

Kemampuan untuk menyesuaikan dan memperluas JupyterHub melalui Authenticator dan Spawner memungkinkan integrasi yang mulus dengan berbagai infrastruktur dan kebutuhan spesifik pengguna.

## 2.2.6 Ray Framework

Ray merupakan sebuah *framework open-source* yang dirancang untuk membangun dan menjalankan aplikasi komputasi paralel dan terdistribusi secara efisien. Framework ini menyediakan abstraksi tingkat tinggi yang memungkinkan pengembang untuk membuat aplikasi yang skalabel dan mudah dijalankan pada kluster komputasi yang terdiri dari banyak node (Moritz et al., 2018).

Ray mendukung dua paradigma pemrograman utama:

### 1. Task-based Computing (Stateless)

Task-based computing memungkinkan pengguna untuk menjalankan fungsi Python secara paralel menggunakan dekorator `@ray.remote`. Model ini bersifat stateless dan cocok digunakan untuk proses komputasi yang dapat dibagi menjadi unit-unit kecil independen.

### 2. Actor-based Computing (Stateful)

Paradigma ini memungkinkan pengguna untuk membuat komponen yang mempertahankan state selama siklus hidupnya. Cocok untuk aplikasi dengan shared state atau layanan yang berjalan terus-menerus.

## Arsitektur Ray

Arsitektur Ray terbagi menjadi dua lapisan utama, yaitu **Application Layer** dan **System Layer (Backend)**.

### Application Layer

Pada lapisan aplikasi, terdapat beberapa komponen utama:

- **Driver:** Komponen yang memulai eksekusi program Ray dan bertanggung jawab dalam mengatur tasks dan actors.
- **Worker:** Unit eksekusi stateless yang digunakan untuk menjalankan fungsi-fungsi remote.
- **Actor:** Unit eksekusi stateful yang mempertahankan data selama proses berjalan.

### System Layer (Backend)

Lapisan sistem ini mencakup manajemen koordinasi, scheduling, dan penyimpanan objek:

**1. Object Store** Object Store adalah penyimpanan berbasis memori untuk menyimpan objek hasil eksekusi:

- Mendukung *zero-copy* antar proses
- Manajemen memori otomatis
- Penyimpanan efisien untuk objek besar

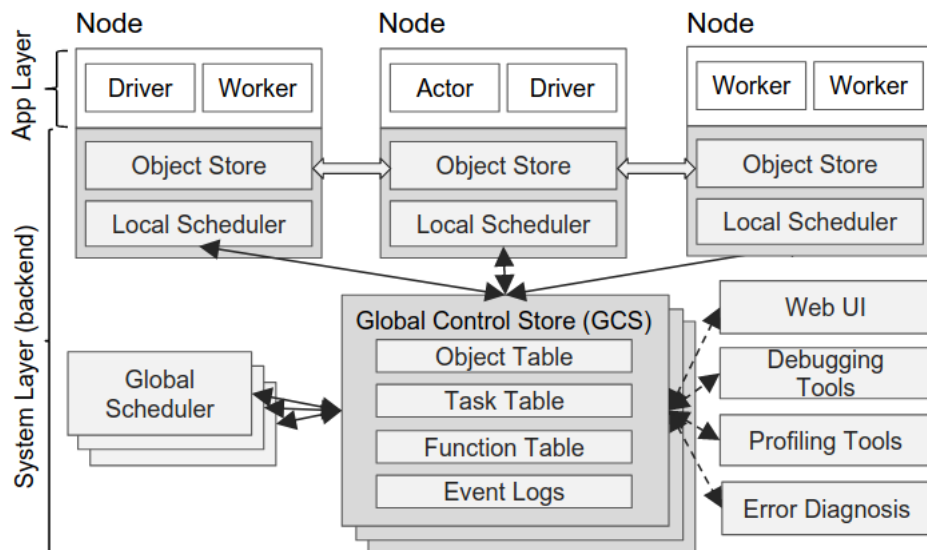
## 2. Raylet (Local Scheduler + Object Manager)

- Menangani penjadwalan tasks lokal
- Mengelola transfer objek antar node
- Berinteraksi dengan Global Control Store (GCS)

## 3. Global Control Store (GCS) Komponen pusat metadata dan koordinasi antar node:

- Menyimpan metadata objek, status tasks, actors, dan nodes
- Menyediakan event logs untuk debugging

**4. Global Scheduler** Mengatur penjadwalan global saat node lokal tidak memiliki resource yang mencukupi.



Gambar 2.2: Komponen RAY (Sumber: Moritz et al., 2018)

## BAB III

### DESAIN DAN IMPLEMENTASI

Penelitian ini dilaksanakan sesuai Lorem ipsum dolor sit amet, consectetur adipiscing elit. Ut purus elit, vestibulum ut, placerat ac, adipiscing vitae, felis. Curabitur dictum gravida mauris. Nam arcu libero, nonummy eget, consectetur id, vulputate a, magna. Donec vehicula augue eu neque.

#### 3.1 Perancangan Arsitektur Sistem

Dalam tahap ini, dilakukan perancangan arsitektur sistem yang mencakup proses penjadwalan pengguna, isolasi lingkungan komputasi, pemilihan node GPU, serta orkestrasi komputasi terdistribusi. Sistem dirancang untuk melayani banyak pengguna secara bersamaan dengan memanfaatkan infrastruktur GPU multi-node dan container Docker yang terintegrasi dengan JupyterHub. Perancangan sistem ini mencakup beberapa tahapan utama, yaitu:

##### 3.1.1 Service Discovery

Discovery Service merupakan komponen inti dari sistem pengelolaan infrastruktur GPU terdistribusi yang berfungsi sebagai pusat koordinasi dan penyimpanan informasi status untuk seluruh node dalam klaster. Implementasi dilakukan menggunakan Flask REST API dengan arsitektur hybrid yang mengintegrasikan PostgreSQL sebagai penyimpanan utama dan Redis sebagai cache layer untuk optimalisasi performa sistem.

##### Arsitektur Penyimpanan Data

Sistem menggunakan arsitektur hybrid storage yang terdiri dari:

- **Redis:** Berfungsi sebagai cache layer untuk data real-time dengan TTL (time-to-live) yang pendek. Redis menyimpan informasi status node yang sering diakses seperti CPU usage, RAM usage, GPU utilization, dan jumlah container aktif. Data ini diperbarui setiap 30 detik oleh agen monitoring di setiap node.
- **PostgreSQL:** Berperan sebagai persistent storage untuk menyimpan data historis, konfigurasi node, metadata pengguna, dan log aktivitas sistem. Database ini memungkinkan analisis tren penggunaan resource dan audit trail untuk keperluan monitoring jangka panjang

##### Mekanisme Pengumpulan Data:

Setiap node dalam klaster menjalankan agen monitoring yang secara berkala mengumpulkan informasi sistem menggunakan library psutil dan gpustat untuk metrics CPU, RAM, dan GPU. Data yang dikumpulkan meliputi:

- Persentase penggunaan CPU, 'RAMdisk'.

- Status dan utilisasi GPU (memory usage, temperature, processes)
- Jumlah container Docker yang sedang berjalan

### REST API Endpoints:

Service discovery menyediakan REST API berbasis Flask dengan endpoint utama:

- **POST /register-node:** Menerima data status dari agen node dan menyimpannya ke Redis dengan PostgreSQL sebagai backup
- **GET /available-nodes:** Mengembalikan daftar node yang tersedia untuk alokasi resource
- 

### Integrasi dengan JupyterHub:

Data dari *service discovery* digunakan oleh custom spawner JupyterHub untuk melakukan pemilihan node (*node selection*). Algoritma pemilihan node mempertimbangkan faktor-faktor seperti *available GPU memory*, *CPU load*, dan *proximity* ke data untuk memastikan distribusi beban yang optimal dan pengalaman pengguna yang responsif.

#### 3.1.2 JupyterHub

JupyterHub bertindak sebagai komponen inti untuk manajemen pengguna dan orkestrasi container JupyterLab dalam sistem GPU terdistribusi. Setiap pengguna mendapatkan lingkungan komputasi terisolasi yang dapat mengakses sumber daya GPU melalui mekanisme penjadwalan yang terintegrasi dengan service discovery.

Sistem menggunakan NativeAuthenticator yang menyediakan fitur registrasi dan autentikasi pengguna built-in tanpa ketergantungan pada sistem eksternal. NativeAuthenticator menyimpan informasi pengguna langsung dalam database PostgreSQL dengan enkripsi password menggunakan bcrypt.

Sistem juga menggunakan *custom spawner* yang terintegrasi dengan *service discovery* untuk pemilihan *node* atau komputer (*node selection*). Setiap spawn request menghasilkan Docker container dengan konfigurasi otomatis yang mencakup:

- **Volume Mounting:** Persistent storage untuk user home directory
- **Network Configuration:** Isolated network dengan port forwarding yang aman
- **Resource Limits:** CPU dan memory limits berdasarkan user quota

JupyterHub mengelola lifecycle container pengguna dengan fitur:

- **Auto-shutdown:** Container otomatis dihentikan setelah periode idle tertentu
- **Resource Monitoring:** Tracking penggunaan resource real-time per user
- **Logging:** Comprehensive logging untuk debugging dan audit

### 3.1.3 Ray Cluster

Ray digunakan untuk mengatur workload komputasi paralel. Dalam perancangannya, setiap container JupyterLab pengguna akan menjalankan Ray Worker secara otomatis dan terhubung ke Ray Head Node. Dengan cara ini, pengguna dapat langsung menggunakan fitur komputasi terdistribusi seperti `ray.remote()` tanpa konfigurasi manual. Ray menjembatani antar-node agar task berat bisa dijalankan dengan efisien di GPU atau CPU sesuai kapasitas.

## 3.2 Implementasi Sistem

Implementasi sistem terdiri dari beberapa komponen utama yang saling terintegrasi, sebagai berikut:

### 3.2.1 Discovery Service

Berdasarkan arsitektur yang telah dijelaskan pada bagian 3.1.1, implementasi Discovery Service dilakukan menggunakan Flask REST API dengan menerapkan arsitektur *hybrid storage* sesuai dengan desain yang telah ditetapkan. Bagian ini menjelaskan detail implementasi kode dari konsep yang telah dirancang sebelumnya.

#### A. Arsitektur Aplikasi

Implementasi layanan mengikuti *application factory pattern*, yang memungkinkan konfigurasi fleksibel berdasarkan lingkungan yang berbeda seperti *development*, *testing*, maupun *production*. Pada tahap inisialisasi, konfigurasi sistem dimuat dari objek konfigurasi yang disesuaikan, kemudian seluruh *extensions* seperti database, migrasi skema, dan *CORS* diaktifkan agar layanan dapat berjalan lintas domain.

Tabel 3.1: Struktur Direktori Proyek Discovery Service

Berkas/Folder	Deskripsi
<code>app.py</code>	Entry point utama untuk Flask API Discovery
<code>config.py</code>	Konfigurasi environment dan koneksi
<code>database.py</code>	Inisialisasi SQLAlchemy dan Redis
<code>database_service.py</code>	Wrapper fungsi query dan insert PostgreSQL
<code>models.py</code>	Model SQLAlchemy untuk Node, GPU, dan Statistik
<code>load_balancer.py</code>	Implementasi algoritma pemilihan node
<code>redis.conf</code>	Konfigurasi Redis lokal
<code>init.sql</code>	Skrip inisialisasi awal database PostgreSQL
<code>Dockerfile</code>	Dockerfile untuk build Discovery API
<code>docker-compose.yml</code>	Orkestrasi Redis, PostgreSQL, dan API
<code>env-agent</code>	File <code>.env</code> untuk menjalankan agent dari container
<code>requirements.txt</code>	Dependensi Python proyek



---

```

1 def create_app(config_name=None):
2     app = Flask(__name__)
3     config_name = config_name or os.environ.get('FLASK_ENV', 'development')
4     config_class = config_mapping.get(config_name, config_mapping['default'])
5     app.config.from_object(config_class)
6
7     db.init_app(app)
8     migrate = Migrate(app, db)
9     CORS(app, origins="*")

```

---

Listing 3.1: Application Factory untuk Discovery Service

## B. Inisialisasi Layanan

Setelah konfigurasi dasar dilakukan, sistem menginisialisasi komponen-komponen layanan seperti penyambung database, algoritma *load balancing*, dan mekanisme seleksi node untuk JupyterHub.

---

```

1 # Initialize services
2 db_service = DatabaseService(app)
3 load_balancer = LoadBalancer()
4 jupyterhub_load_balancer = JupyterHubLoadBalancer(
5     max_cpu=app.config['JUPYTERHUB_MAX_CPU'],
6     max_memory=app.config['JUPYTERHUB_MAX_MEMORY'],
7     max_containers=app.config['JUPYTERHUB_MAX_CONTAINERS']
8 )

```

---

Listing 3.2: Inisialisasi Service dan Load Balancer

## C. Koneksi ke Redis

Redis digunakan untuk menyimpan informasi status node secara *real-time* dengan *time-to-live (TTL)* tertentu. Koneksi ke Redis menggunakan *connection pooling* agar efisien dan tangguh terhadap kegagalan jaringan.

---

```

1 try:
2     pool = ConnectionPool(
3         host=app.config['REDIS_HOST'],
4         port=app.config['REDIS_PORT'],
5         password=app.config['REDIS_PASSWORD'],
6         decode_responses=True
7     )
8     redis_client = Redis(connection_pool=pool)
9     redis_client.ping()
10    logger.info(f"Connected to Redis at {app.config['REDIS_HOST']}:{app.config['REDIS_PORT']}")
11 except Exception as e:
12    logger.error(f"Failed to connect to Redis: {e}")

```

---

Listing 3.3: Koneksi Redis dengan Error Handling

## D. Endpoint Health Check

Layanan ini menyediakan endpoint `/health-check` untuk memeriksa status koneksi terhadap Redis dan PostgreSQL, serta status *background task* pembersihan otomatis dari data node.

---

```
1 @app.route("/health-check")
2 def health_check():
3     redis_status = "connected" if app.redis_client else "disconnected"
4     try:
5         db.session.execute(db.text('SELECT 1'))
6         db_status = "connected"
7     except Exception:
8         db_status = "disconnected"
9
10    cleanup_status = app.db_service.get_cleanup_status()
11
12    return jsonify({
13        "status": "ok",
14        "message": "Hello, from [DiscoveryAPI]",
15        "redis_status": redis_status,
16        "database_status": db_status,
17        "cleanup_task_active": cleanup_status.get('cleanup_thread_active', ←
18        , False),
19        "node_summary": {
20            "active_nodes": cleanup_status.get('nodes', {}).get('active', ←
21            0),
22            "total_nodes": cleanup_status.get('nodes', {}).get('total', ←
23            0)
24        }
25    }), 200
```

---

Listing 3.4: Endpoint Health Check

## E. Endpoint Registrasi Node

Endpoint utama `/register-node` digunakan oleh *agent* di setiap node untuk mengirimkan data status terkini, seperti pemakaian CPU, memori, GPU, serta jumlah kontainer yang sedang berjalan. Data ini disimpan ke PostgreSQL untuk keperluan jangka panjang dan Redis untuk kecepatan akses oleh sistem pemilihan node di JupyterHub.

---

```
1 @app.route("/register-node", methods=["POST"])
2 def register_node():
3     data = request.get_json()
4     hostname = data.get("hostname")
5     ip = data.get("ip")
6
7     if not hostname:
8         return jsonify({"error": "hostname is required"}), 400
9
10    try:
11        node = app.db_service.register_or_update_node(data)
12
13        if app.redis_client:
```

---

```

14         app.redis_client.set(
15             f"node:{hostname}:info",
16             json.dumps(data),
17             ex=app.config['REDIS_EXPIRE_SECONDS']
18         )
19         if ip:
20             app.redis_client.set(
21                 f"node:{hostname}:ip",
22                 ip,
23                 ex=app.config['REDIS_EXPIRE_SECONDS']
24             )
25
26         return jsonify({
27             "status": "ok",
28             "stored": True,
29             "node_id": node.id,
30             "hostname": node.hostname,
31             "load_score": node.load_score
32         })
33     except Exception as e:
34         logger.error(f"Node registration error: {e}")
35         return jsonify({"error": f"Registration failed: {str(e)}"}), 500

```

---

Listing 3.5: Endpoint Registrasi Node

## F. Layanan Abstraksi Basis Data

Untuk memisahkan logika bisnis dari interaksi langsung dengan database, sistem ini menggunakan lapisan abstraksi yang disebut `DatabaseService`. Kelas ini bertanggung jawab atas pengolahan data node ke dalam PostgreSQL, baik saat registrasi, pembaruan metrik, maupun pengambilan informasi untuk analisis dan seleksi node.

Cuplikan berikut menunjukkan salah satu metode penting yaitu `register_or_update_node()`, yang digunakan dalam endpoint `/register-node`:

---

```

1 def register_or_update_node(self, data: Dict[str, Any]) -> Node:
2     hostname = data.get("hostname")
3     node = db.session.query(Node).filter_by(hostname=hostname).first()
4
5     if not node:
6         node = Node(hostname=hostname, ip=data.get("ip"))
7         db.session.add(node)
8
9     node.update_from_dict(data)
10    db.session.commit()
11    return node

```

---

Listing 3.6: Metode Register/Update Node

## G. Penutup Integrasi Komponen

Seluruh komponen yang telah dijelaskan sebelumnya, mulai dari konfigurasi awal, koneksi ke Redis dan PostgreSQL, hingga definisi endpoint dan pemrosesan data, membentuk fondasi

utama dari layanan Discovery. Proses seleksi node untuk keperluan JupyterHub tidak dilakukan langsung di endpoint, melainkan dipisahkan ke dalam modul LoadBalancer yang mendukung berbagai strategi pemilihan berdasarkan kondisi aktual node.

Dengan desain modular dan terintegrasi ini, layanan Discovery mampu menyediakan informasi node yang akurat dan *real-time*, sekaligus mempertahankan fleksibilitas dan skalabilitas sistem.

### 3.2.1.1 Load Balancer untuk Seleksi Node

Komponen *Load Balancer* merupakan bagian integral dari Discovery Service yang bertugas melakukan seleksi node secara efisien berdasarkan kondisi sumber daya dan kriteria tertentu. Meskipun diimplementasikan dalam satu layanan yang sama, fungsionalitas pemilihan node ini dipisahkan dalam modul terpisah bernama `load_balancer.py` untuk mendukung prinsip pemisahan tanggung jawab (*separation of concerns*).

#### A. Struktur dan Algoritma

Kelas LoadBalancer menyediakan beberapa algoritma seleksi seperti *round robin*, *least loaded*, *weighted round robin*, dan *random*. Pemilihan node didasarkan pada nilai skor beban (`load_score`) yang dikalkulasi dari penggunaan CPU, memori, dan jumlah kontainer aktif. Untuk kasus khusus seperti permintaan dari JupyterHub, digunakan kelas turunan JupyterHubLoadBalancer dengan kriteria seleksi yang lebih ketat.

---

```
1 class LoadBalancer:
2     def __init__(self):
3         self.round_robin_counter = 0
4         self.counter_lock = threading.Lock()
5         self.algorithms = {
6             'round_robin': self._round_robin_select,
7             'least_loaded': self._least_loaded_select,
8             'weighted_round_robin': self._weighted_round_robin_select,
9             'random': self._random_select
10        }
11
12    def select_node(self, nodes: List[Node], algorithm: str = '↵
round_robin',
13                  requested_count: int = 1, **kwargs) -> Dict[str, Any]:
14        if not nodes:
15            return {'selected_nodes': [], 'error': 'No nodes available'}
16        sorted_nodes = sorted(nodes, key=lambda n: n.load_score or 0)
17        selector = self.algorithms.get(algorithm, self.↵
_round_robin_select)
18        ... # Pemilihan node dan update counter
```

---

Listing 3.7: Contoh Seleksi Node dengan Round Robin

#### B. Seleksi Khusus untuk JupyterHub

Kelas JupyterHubLoadBalancer mengimplementasikan metode `select_jupyterhub_node` yang menyaring node berdasarkan batas maksimum CPU, memori, dan kontainer aktif. Hal ini bertujuan untuk memastikan kualitas layanan pengguna JupyterHub tetap terjaga.

---

```
1 class JupyterHubLoadBalancer(LoadBalancer):
```

---

```

2     def select_jupyterhub_node(self, nodes: List[Node], algorithm: str = ↵
      'round_robin',
3                                     requested_count: int = 1) -> Dict[str, Any↵
      ]:
4         filtered_nodes = []
5         for node in nodes:
6             if (node.cpu_usage_percent < self.max_cpu and
7                 node.memory_usage_percent < self.max_memory and
8                 (node.active_jupyterlab + node.active_ray) < self.↵
                  max_containers):
9                 filtered_nodes.append(node)
10        return self.select_node(filtered_nodes, algorithm, ↵
            requested_count)

```

---

Listing 3.8: Seleksi Node JupyterHub dengan Kriteria Khusus

### C. Penerapan dalam Endpoint Discovery

Load balancer ini digunakan secara langsung dalam endpoint `/available-nodes` dan `/jupyterhub-nodes`, yang masing-masing akan memanggil fungsi seleksi sesuai algoritma yang diminta oleh klien, baik itu pengguna biasa maupun sistem JupyterHub.

Dengan arsitektur ini, sistem dapat mendistribusikan beban secara lebih merata, menghindari bottleneck, dan meningkatkan utilisasi sumber daya komputasi dalam kluster.

## 3.2.2 Agent Service

Setelah layanan Discovery Service diimplementasikan, sistem memerlukan komponen tambahan yang berjalan secara periodik di setiap node. Komponen ini disebut sebagai *Agent Service*. Agent bertanggung jawab untuk mengumpulkan informasi sistem dan mengirimkannya secara berkala ke endpoint `/register-node` pada Discovery API. Informasi tersebut mencakup pemanfaatan CPU, memori, disk, deteksi GPU, serta jumlah container yang sedang aktif.

### 3.2.2.1 Arsitektur dan Fungsi Agent

Agent dikembangkan sebagai skrip Python mandiri yang berjalan sebagai proses daemon ringan pada setiap node. Agent ini dirancang agar:

- Mengirimkan data sistem setiap 15 detik.
- Menangkap informasi hardware dan aktivitas container.
- Tetap ringan dan tidak membebani node secara signifikan.

Agent memanfaatkan pustaka seperti `psutil` untuk informasi sistem, `gpustat` untuk deteksi GPU NVIDIA, dan pustaka `docker` untuk menghitung container aktif. Semua data dikirim dalam format JSON ke Discovery API.

---

```

1 def register():
2     print("[DEBUG] adding node...")
3     payload = collect_node_info()
4     if payload:
5         try:
6             resp = requests.post(DISCOVERY_URL, json=payload)

```

```

7         print(f"[AGENT] Registered: {payload['hostname']} ({payload['↵
            ip']}) {resp.status_code}")
8     except Exception as e:
9         print(f"[AGENT] Failed to register node: {e}")

```

---

Listing 3.9: Pengumpulan dan Pengiriman Informasi oleh Agent

Fungsi `collect_node_info()` bertanggung jawab mengumpulkan data sistem, termasuk deteksi GPU dan status container Docker. Cuplikan fungsi berikut menjelaskan proses ini:

```

1 def collect_node_info():
2     hostname = socket.gethostname()
3     ip_address = os.popen("hostname -I").read().strip().split()[0]
4     ram_gb = round(psutil.virtual_memory().total / 1e9, 2)
5     cpu_usage = psutil.cpu_percent(interval=1)
6     memory = psutil.virtual_memory()
7     disk = psutil.disk_usage("/")
8     container_info = get_container_info()
9     gpu_stats = get_gpu_stats()
10
11     return {
12         "hostname": hostname,
13         "ip": ip_address,
14         "cpu": os.cpu_count(),
15         "gpu": gpu_stats,
16         "has_gpu": len(gpu_stats) > 0,
17         "ram_gb": ram_gb,
18         "cpu_usage_percent": round(cpu_usage, 2),
19         "memory_usage_percent": round(memory.percent, 2),
20         "disk_usage_percent": round(disk.percent, 2),
21         "active_jupyterlab": container_info["jupyterlab_count"],
22         "active_ray": container_info["ray_count"],
23         "total_containers": container_info["total_count"],
24         "last_updated": datetime.now().isoformat() + "Z"
25     }

```

---

Listing 3.10: Fungsi Pengumpulan Informasi Sistem

### 3.2.2.2 Mekanisme Deteksi GPU dan Kontainer

Deteksi GPU dilakukan melalui pustaka `gpustat`, yang mampu membaca informasi memori, suhu, dan utilisasi GPU. Untuk menghitung kontainer aktif, agent memeriksa nama atau image dari setiap container menggunakan pustaka Docker.

```

1 def get_gpu_stats():
2     try:
3         stats = gpustat.GPUStatCollection.new_query()
4         return [{
5             "name": gpu.name,
6             "index": gpu.index,
7             "uuid": gpu.uuid,
8             "memory_total_mb": gpu.memory_total,
9             "memory_used_mb": gpu.memory_used,
10            "memory_util_percent": round(gpu.memory_used / gpu.↵
                memory_total * 100, 2)
11            if gpu.memory_total else 0,
12            "utilization_gpu_percent": gpu.utilization,

```

```

13         "temperature_gpu": gpu.temperature
14     } for gpu in stats.gpus]
15 except Exception as e:
16     print(f"[GPUSTAT] NVIDIA GPU not available or error: {e}")
17     return []

```

---

Listing 3.11: Fungsi Deteksi GPU dan Container Aktif

### 3.2.2.3 Integrasi dengan Discovery Service

Data yang dikirim oleh Agent akan diproses oleh endpoint `/register-node` pada Discovery API, disimpan ke Redis untuk akses cepat, dan ke PostgreSQL untuk histori dan audit jangka panjang. Detail lengkap terkait penyimpanan, struktur model database, serta pencatatan statistik dapat dilihat pada bagian 3.2.1.

### 3.2.2.4 Integrasi dengan Database PostgreSQL

Sistem Service Discovery dan Load Balancer pada platform ini mengandalkan PostgreSQL sebagai penyimpanan jangka panjang untuk informasi node. Basis data ini menyimpan status dan atribut setiap node, termasuk informasi penggunaan CPU, memori, GPU, dan histori alokasi.

#### A. Model Node dan Relasi GPU

Model utama dalam database adalah entitas Node, yang merepresentasikan setiap mesin komputasi yang terdaftar dalam sistem. Tiap Node memiliki sejumlah atribut seperti hostname, IP address, jumlah CPU, RAM, penggunaan sumber daya, dan total kontainer aktif. Untuk mendukung kebutuhan GPU, terdapat relasi ke entitas NodeGPU menggunakan relasi satu ke banyak.

---

```

1 class Node(db.Model):
2     __tablename__ = 'nodes'
3     id = db.Column(db.Integer, primary_key=True)
4     hostname = db.Column(db.String(255), unique=True, nullable=False)
5     ip = db.Column(db.String(45), nullable=False)
6     ram_gb = db.Column(db.Float, default=0.0)
7     cpu_usage_percent = db.Column(db.Float, default=0.0)
8     ...
9     gpus = db.relationship('NodeGPU', backref='node', lazy='dynamic', ←
        cascade='all, delete-orphan')
10
11 class NodeGPU(db.Model):
12     __tablename__ = 'node_gpus'
13     id = db.Column(db.Integer, primary_key=True)
14     node_id = db.Column(db.Integer, db.ForeignKey('nodes.id'), nullable=↵
        False)
15     name = db.Column(db.String(255), nullable=False)
16     memory_total_mb = db.Column(db.Integer, default=0)
17     ...

```

---

Listing 3.12: Model Node dan Relasi GPU

#### B. Pembaruan Status Node

Ketika agent melakukan pendaftaran ke endpoint `/register-node`, informasi node disimpan

ke dalam PostgreSQL menggunakan metode `update_from_dict`. Proses ini memastikan data node diperbarui secara berkala dan historinya tercatat dengan baik.

---

```
1 def update_from_dict(self, data):
2     self.cpu = data.get('cpu', self.cpu)
3     self.ram_gb = data.get('ram_gb', self.ram_gb)
4     self.cpu_usage_percent = data.get('cpu_usage_percent', self.cpu_usage_percent)
5     self.active_jupyterlab = data.get('active_jupyterlab', self.active_jupyterlab)
6     self.last_seen = datetime.now(timezone.utc)
7     ...
```

---

Listing 3.13: Pembaruan Status Node ke Database

### C. Pencatatan Riwayat Alokasi dan Statistik Load Balancer

Untuk kebutuhan analisis, sistem juga mencatat riwayat penugasan node dalam entitas `NodeAssignment` dan performa algoritma pemilihan node ke dalam `LoadBalancerStats`. Hal ini memungkinkan audit dan evaluasi terhadap distribusi beban.

---

```
1 class NodeAssignment(db.Model):
2     __tablename__ = 'node_assignments'
3     node_id = db.Column(db.Integer, db.ForeignKey('nodes.id'))
4     session_id = db.Column(db.String(255), index=True)
5     algorithm_used = db.Column(db.String(50))
6     status = db.Column(db.String(50))
7     assigned_at = db.Column(db.DateTime)
8
9 class LoadBalancerStats(db.Model):
10     __tablename__ = 'load_balancer_stats'
11     algorithm = db.Column(db.String(50))
12     avg_cpu_usage = db.Column(db.Float)
13     recorded_at = db.Column(db.DateTime, index=True)
```

---

Listing 3.14: Model Assignment dan Statistik

Melalui integrasi ini, sistem tidak hanya dapat mengambil keputusan secara real-time melalui Redis, tetapi juga menyimpan informasi penting untuk pemantauan jangka panjang, pelaporan, dan peningkatan sistem di masa depan.

## 3.2.3 Implementasi JupyterHub

JupyterHub merupakan komponen utama dalam sistem ini yang bertanggung jawab atas autentikasi pengguna, manajemen sesi, dan peluncuran container JupyterLab secara terdistribusi. Implementasi dilakukan menggunakan pendekatan modular dengan direktori terpisah untuk konfigurasi, spawner, kernel, dan form antarmuka pengguna.

### 3.2.3.1 A. Struktur Proyek

Struktur proyek JupyterHub ditunjukkan pada Tabel 3.2, yang mencakup konfigurasi utama, custom spawner, image container pengguna, dan kernel Ray.



Tabel 3.2: Struktur Direktori Proyek JupyterHub

Folder	Deskripsi
config/	Berisi konfigurasi modular untuk spawner, proxy, autentikasi
spawner/	Implementasi MultiNodeSpawner yang mewarisi DockerSpawner
form/	Form HTML interaktif untuk pemilihan profil dan konfigurasi node
singleuser/	Dockerfile dan skrip inisialisasi untuk container JupyterLab
ray-kernel/	Implementasi custom kernel untuk integrasi Ray pada JupyterLab
data/	Menyimpan data runtime JupyterHub seperti cookie secret dan SQLite
docker-compose.yml	Orkestrasi service JupyterHub dan dependency lain

### 3.2.3.2 B. Spawner dengan Service Discovery

Spawner kustom bernama MultiNodeSpawner dikembangkan dalam direktori spawner/. Spawner ini memperluas fungsi dari DockerSpawner untuk dapat memilih node dari Discovery API secara dinamis berdasarkan profil pengguna.

```

1 c.JupyterHub.spawner_class = 'spawner.multinode.MultiNodeSpawner'
2 c.MultiNodeSpawner.discovery_api_url = 'http://discovery-api:15002'
3 c.MultiNodeSpawner.enable_multi_node = True
4 c.MultiNodeSpawner.allowed_images = {
5     "cpu": "danielcrith0/jupyterlab:cpu",
6     "gpu": "danielcrith0/jupyterlab:gpu"
7 }
```

Listing 3.15: Contoh Konfigurasi MultiNodeSpawner

### 3.2.3.3 C. Form Interaktif dan Profil

Antarmuka pemilihan konfigurasi pengguna dibuat dalam berkas form/form.html yang di-load saat pengguna login. Form ini menampilkan daftar node dan profil yang dapat dipilih oleh pengguna, termasuk opsi CPU/GPU dan batas memori.

Form tersebut diproses menggunakan pre\_spawn\_hook dan options\_form dalam konfigurasi spawner untuk menyesuaikan container yang akan dijalankan.

### 3.2.3.4 D. Integrasi Ray Kernel

JupyterLab yang dijalankan pada setiap container telah terintegrasi dengan kernel Ray, yang memungkinkan pengguna untuk langsung menjalankan kode terdistribusi. Kernel tersebut diimplementasikan dalam direktori ray-kernel/ dan disisipkan dalam image container melalui Dockerfile.

```

1 FROM jupyter/base-notebook
2 RUN pip install -r /singleuser/requirements.txt
3 COPY ray-kernel/ /usr/local/share/ray-kernel/
4 RUN python -m ray_kernel.install_ray_kernel
```

Listing 3.16: Dockerfile untuk Container CPU

#### **3.2.3.5 E. Isolasi dan Manajemen Resource**

Setiap container JupyterLab dijalankan secara terisolasi di node terpilih, dengan batasan sumber daya yang ditentukan oleh profil. Auto-culling diterapkan untuk menghentikan container idle, dan penyimpanan dilakukan di volume lokal masing-masing container.

### **3.3 Peralatan Pendukung**

Perangkat yang digunakan untuk pengerjaan tugas akhir ini merupakan sebuah komputer dengan spesifikasi sebagai berikut.

Tabel 3.3: Spesifikasi Peralatan Pendukung

No.	Komponen	Spesifikasi
1	<b>Laptop</b>	
	<i>Brand</i>	Asus
	<i>Processor</i>	AMD Ryzen 3
	<i>Operating System</i>	Ubuntu 22.04 LTS
	<i>GPU</i>	AMD Radeon vega 3 graphics
	<i>Memory</i>	18 GB
	<i>Storage</i>	512 GB
2	<b>Komputer</b>	
	<i>Brand</i>	Asus
	<i>Processor</i>	12th Gen Intel i9-12900K
	<i>Operating System</i>	Ubuntu 24.04 LTS
	<i>GPU</i>	NVIDIA GeForce RTX 3080 Ti
	<i>Memory</i>	64 GB
	<i>Storage</i>	723 GB
3	<b>Virtual Machine 1</b>	
	<i>Brand</i>	Asus
	<i>Processor</i>	12th Gen Intel i9-12900K
	<i>Operating System</i>	Ubuntu 24.04 LTS
	<i>GPU</i>	NVIDIA GeForce RTX 3080 Ti
	<i>Memory</i>	64 GB
	<i>Storage</i>	723 GB
4	<b>Virtual Machine 2</b>	
	<i>Brand</i>	Asus
	<i>Processor</i>	12th Gen Intel i9-12900K
	<i>Operating System</i>	Ubuntu 24.04 LTS
	<i>GPU</i>	NVIDIA GeForce RTX 3080 Ti
	<i>Memory</i>	64 GB
	<i>Storage</i>	723 GB

Selain perangkat keras, terdapat juga perangkat lunak pendukung seperti berikut.

Tabel 3.4: Daftar Perangkat Lunak Pendukung

<b>Nama Perangkat Lunak</b>	<b>Versi</b>	<b>Keterangan</b>
Docker Engine	24.x	Digunakan untuk menjalankan container JupyterLab dan Ray Worker secara terisolasi. Memungkinkan lingkungan komputasi tiap pengguna berjalan secara independen dan mudah didistribusikan ke berbagai node.
Docker Compose	2.x	Membantu mendefinisikan dan mengatur layanan multi-container (seperti JupyterHub dan Redis) dalam satu berkas konfigurasi. Memudahkan manajemen dan replikasi layanan.
JupyterHub	5.3.0	Menangani autentikasi pengguna serta spawn container JupyterLab ke node terpilih berdasarkan data dari service discovery.
Ray	2.46	Framework komputasi paralel dan terdistribusi. Setiap pengguna dapat langsung menjalankan task terdistribusi secara otomatis dari dalam JupyterLab.
Redis	7.0	Database key-value in-memory yang digunakan untuk menyimpan status sistem (CPU, RAM, GPU) dan log aktivitas pengguna secara real-time.
Flask	3.1.0	Framework web Python yang digunakan untuk membangun service discovery berupa REST API yang menerima dan menyediakan data status node.
Python	3.11	Bahasa pemrograman utama yang digunakan untuk seluruh komponen sistem, seperti konfigurasi JupyterHub, pengembangan REST API, Ray, serta skrip monitoring.
PostgreSQL	14	Bahasa pemrograman utama yang digunakan untuk seluruh komponen sistem, seperti konfigurasi JupyterHub, pengembangan REST API, Ray, serta skrip monitoring.

*[Halaman ini sengaja dikosongkan]*

## **BAB IV**

### **PENGUJIAN DAN ANALISIS**

Bab ini membahas proses pengujian dan hasil analisis terhadap sistem yang telah dibangun. Tujuan utama dari pengujian ini adalah untuk mengevaluasi kinerja Service Discovery dalam memilih node yang optimal untuk menjalankan container JupyterLab, serta memastikan bahwa integrasi antar komponen (JupyterHub, Discovery API, Agent, dan Docker) berjalan sesuai ekspektasi.

#### **4.1 Skenario Pengujian**

Pengujian dilakukan dengan lima skenario yang dirancang untuk mewakili berbagai kondisi nyata dalam penggunaan sistem. Setiap skenario difokuskan pada aspek tertentu seperti seleksi node, distribusi user, validasi profil GPU, serta penanganan beban dan kegagalan node.

##### **4.1.1 Skenario 1: Pemilihan Node dengan Beban Terendah**

- **Tujuan:** Memastikan sistem memilih node dengan CPU usage terendah.
- **Langkah:** Jalankan 1 user (profil: single-cpu) dan catat node terpilih.
- **Hasil:** Sistem memilih node dengan 24.7% CPU, lebih rendah dari kandidat lainnya.

##### **4.1.2 Skenario 2: Multi-User Concurrent**

- **Tujuan:** Menguji distribusi kontainer saat 5 user masuk secara paralel.
- **Profil:** 2 user GPU, 3 user CPU.
- **Hasil:** Node GPU digunakan optimal, node CPU terdistribusi merata.

##### **4.1.3 Skenario 3: Simulasi Beban Tinggi**

- **Tujuan:** Memastikan sistem menghindari node yang sedang overload.
- **Langkah:** Menjalankan stress-ng di satu node.
- **Hasil:** Node tersebut tidak dipilih oleh sistem.

##### **4.1.4 Skenario 4: Validasi Profil GPU**

- **Tujuan:** User GPU hanya boleh dialokasikan ke node dengan GPU.
- **Hasil:** Semua user GPU dialokasikan ke node dengan NVIDIA A100.

#### 4.1.5 Skenario 5: TTL Redis dan Node Tidak Aktif

- **Tujuan:** Node yang tidak update status akan dihapus otomatis.
- **Langkah:** Matikan agent dan tunggu TTL (45s), lalu spawn user.
- **Hasil:** Node tidak aktif diabaikan oleh Discovery API.

## 4.2 Evaluasi Pengujian

Tabel berikut merangkum hasil pengujian dan status keberhasilan sistem dalam setiap skenario:

Tabel 4.1: Ringkasan Evaluasi Pengujian Sistem

No	>Skenario Pengujian	Berhasil
1	Pemilihan node dengan CPU terendah untuk profil single-cpu	Ya
2	Distribusi 5 user secara paralel (profil campuran CPU dan GPU)	Ya
3	Sistem menghindari node yang overload (stress test)	Ya
4	Alokasi profil GPU hanya ke node yang memiliki GPU	Ya
5	Node dengan agent mati tidak terpilih setelah TTL Redis habis	Ya

Evaluasi menunjukkan bahwa sistem berhasil melakukan alokasi kontainer sesuai dengan desain arsitektur dan kriteria seleksi node. Load balancing bekerja optimal, dan fitur failover serta integrasi antar komponen telah diuji dengan baik.

## **BAB V**

### **PENUTUP**

#### **5.1 Kesimpulan**

Berdasarkan hasil pengujian yang Lorem ipsum dolor sit amet, consectetur adipiscing elit. Ut purus elit, vestibulum ut, placerat ac, adipiscing vitae, felis. Curabitur dictum gravida mauris. sebagai berikut:

1. Pembuatan Nam dui ligula, fringilla a, euismod sodales, sollicitudin vel, wisi. Morbi auctor lorem non justo. Nam lacus libero, pretium at, lobortis vitae, ultricies et, tellus.
2. Donec aliquet, tortor sed accumsan bibendum, erat ligula aliquet magna, vitae ornare odio metus a mi. Morbi ac orci et nisl hendrerit mollis. Suspendisse ut massa.
3. Cras nec ante. Pellentesque a nulla. Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus. Aliquam tincidunt urna.

#### **5.2 Saran**

Untuk pengembangan lebih lanjut pada Lorem ipsum dolor sit amet, consectetur adipiscing elit. Ut purus elit, vestibulum ut, placerat ac, adipiscing vitae, felis. Curabitur dictum gravida mauris. antara lain:

1. Memperbaiki Nam dui ligula, fringilla a, euismod sodales, sollicitudin vel, wisi. Morbi auctor lorem non justo. Nam lacus libero, pretium at, lobortis vitae, ultricies et, tellus.
2. Donec aliquet, tortor sed accumsan bibendum, erat ligula aliquet magna, vitae ornare odio metus a mi. Morbi ac orci et nisl hendrerit mollis. Suspendisse ut massa.
3. Cras nec ante. Pellentesque a nulla. Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus. Aliquam tincidunt urna.



*[Halaman ini sengaja dikosongkan]*

## DAFTAR PUSTAKA

- Kumar, A., Cuccuru, G., Grüning, B., & Backofen, R. (2023). An accessible infrastructure for artificial intelligence using a docker-based jupyterlab in galaxy [Published: 26 April 2023]. *GigaScience*, 12. <https://doi.org/10.1093/gigascience/giad028>
- Li, W., Lafuente Mercado, R. S., Pena, J. D., & Allen, R. E. (2024). Syndeo: Portable ray clusters with secure containerization [arXiv:2409.17070v1 [cs.DC]]. *arXiv preprint arXiv:2409.17070*. <https://arxiv.org/abs/2409.17070>
- Moritz, P., Nishihara, R., Wang, S., Tumanov, A., Liaw, R., Liang, E., Elibol, M., Yang, Z., Paul, W., Jordan, M. I., & Stoica, I. (2018). Ray: A distributed framework for emerging ai applications. *Proceedings of the 13th USENIX Symposium on Operating Systems Design and Implementation (OSDI '18)*, 561–577. <https://www.usenix.org/conference/osdi18/presentation/moritz>
- Shikai Wang, X. W., Haotian Zheng, & Shang, F. (2024). Distributed high-performance computing methods for accelerating deep learning training. *jklst*. <https://jklst.org/index.php/home/article/view/230>
- Team, S. D. (2024). *What is docker architecture?* [Accessed: December 24, 2024]. <https://sysdig.com/learn-cloud-native/what-is-docker-architecture>
- Turnbull, J. (2014). *The docker book: Containerization is the new virtualization*.
- Zhou, N., Zhou, H., & Hoppe, D. (2022). Containerisation for high performance computing systems: Survey and prospects. *arXiv*. <https://arxiv.org/abs/2212.08717>

*[Halaman ini sengaja dikosongkan]*

## BIOGRAFI PENULIS



Gloriyano Cristho Daniel Pepuho, lahir pada Lorem ipsum dolor sit amet, consectetur adipiscing elit. Ut purus elit, vestibulum ut, placerat ac, adipiscing vitae, felis. Curabitur dictum gravida mauris. Nam arcu libero, nonummy eget, consectetur id, vulputate a, magna. Donec vehicula augue eu neque. Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Mauris ut leo. Cras viverra metus rhoncus sem. Nulla et lectus vestibulum urna fringilla ultrices. Phasellus eu tellus sit amet tortor gravida placerat. Integer sapien est, iaculis in, pretium quis, viverra ac, nunc. Praesent eget sem vel leo ultrices bibendum. Aenean faucibus. Morbi dolor nulla, malesuada eu, pulvinar at, mollis ac, nulla. Curabitur auctor semper nulla. Donec varius orci eget risus. Duis nibh mi, congue eu, accumsan eleifend, sagittis quis, diam. Duis eget orci sit amet orci dignissim rutrum.

Nam dui ligula, fringilla a, euismod sodales, sollicitudin vel, wisi. Morbi auctor lorem non justo. Nam lacus libero, pretium at, lobortis vitae, ultricies et, tellus. Donec aliquet, tortor sed accumsan bibendum, erat ligula aliquet magna, vitae ornare odio metus a mi. Morbi ac orci et nisl hendrerit mollis. Suspendisse ut massa. Cras nec ante. Pellentesque a nulla. Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus. Aliquam tincidunt urna. Nulla ullamcorper vestibulum turpis. Pellentesque cursus luctus mauris.

*[Halaman ini sengaja dikosongkan]*