

AI Factory Benchmarking Framework

Technical Report

Team 6 – EUMaster4HPC 2025-2026

Supervisor: Dr. Farouk Mansouri, LuxProvide

Institution: LuxProvide – MeluXina Supercomputer

January 12th, 2026

Abstract

The emergence of AI Factories in the European Union demands robust evaluation methodologies for AI infrastructure components. This report presents a comprehensive benchmarking framework designed to evaluate end-to-end performance for critical AI Factory components on the MeluXina supercomputer, one of Europe’s most powerful supercomputers.

The AI Factory Benchmarking Framework provides unified benchmarking capabilities across diverse AI services. Supported services include Large Language Model (LLM) inference servers (Ollama, vLLM, Triton), vector databases (ChromaDB, Faiss, Weaviate), relational databases (PostgreSQL), object storage (S3-compatible MinIO), and file storage systems. The framework integrates production-grade monitoring using Prometheus and Grafana. It also supports SLURM-based HPC cluster orchestration and utilizes Apptainer containerization for reproducible deployments.

Key contributions include: (1) a modular Python-based architecture enabling extensible service benchmarking, (2) real-time monitoring with comprehensive metrics collection and visualization, (3) YAML-driven configuration for reproducible benchmark scenarios, (4) parallel execution support achieving up to 10× faster benchmarking, and (5) detailed reporting with statistical analysis and percentile calculations.

This framework equips researchers and practitioners with the tools necessary to design and operate future AI Factory workflows at scale, providing comparative insights, performance guidelines, and scalability assessments for AI infrastructure planning.

Keywords: High-Performance Computing, AI Benchmarking, LLM Inference, Vector Databases, Prometheus, Grafana, SLURM, Apptainer, GPU Acceleration, MeluXina

Contents

Abstract	1
1 Introduction	7
1.1 Background and Motivation	7
1.2 Project Objectives	7
1.3 Challenge Overview	8
1.4 Report Structure	8
2 System Architecture	10
2.1 Architectural Overview	10
2.2 Design Principles	10
2.3 Data Flow	11
3 Core Components	12
3.1 Main Orchestrator	12
3.1.1 Execution Modes	12
3.2 Middleware Interface	13
3.3 Server Manager	13
3.3.1 SLURM Script Generation	13
3.4 Unified Client Factory	14
3.5 Monitor	14
3.6 Metrics Interceptor	15
3.7 Reporter	15
3.7.1 Report Format	15
4 Supported Services	17
4.1 LLM Inference Services	17
4.1.1 Ollama	17
4.1.2 vLLM	17
4.1.3 Triton Inference Server	18

4.2	Vector Databases	18
4.2.1	ChromaDB	18
4.2.2	Faiss	18
4.2.3	Weaviate	18
4.3	Relational Databases	19
4.3.1	PostgreSQL	19
4.4	Storage Systems	19
4.4.1	Object Storage (MinIO)	19
4.4.2	File Storage	19
5	Monitoring and Visualization	20
5.1	Monitoring Architecture	20
5.2	Prometheus Integration	20
5.2.1	Key Metrics	21
5.3	Pushgateway	21
5.4	Grafana Dashboards	21
5.5	Accessing Monitoring from Local Machine	22
6	HPC Integration	23
6.1	MeluXina Supercomputer	23
6.2	SLURM Integration	23
6.2.1	Configuration Parameters	23
6.2.2	Job Lifecycle	24
6.3	Apptainer Containerization	24
6.3.1	Container Management	24
6.3.2	Bind Mounts for Persistent Data	24
6.4	GPU Acceleration	24
7	Implementation Details	26
7.1	Technology Stack	26
7.2	Dependencies	26
7.3	Recipe Configuration System	27
7.4	Installation and Setup	28
7.4.1	HPC Installation (MeluXina)	28
7.5	Usage Examples	28
7.5.1	Running Benchmarks	28
7.5.2	SLURM Batch Execution	29
7.5.3	Interactive Session	29
8	Benchmark Methodology and Results	30

8.1	Benchmark Methodology	30
8.1.1	Warm-up Phase	30
8.1.2	Statistical Rigor	30
8.2	Performance Metrics	31
8.2.1	Latency Metrics	31
8.2.2	Throughput Metrics	31
8.2.3	Reliability Metrics	31
8.3	Result Analysis	31
8.3.1	Querying Metrics Database	31
8.3.2	Comparing Benchmarks	32
8.4	Expected Performance Characteristics	32
9	Conclusion and Future Work	33
9.1	Summary of Achievements	33
9.2	Lessons Learned	33
9.3	Future Work	34
9.4	Acknowledgments	34
A	Apptainer Container Recipes	35
B	Grafana Dashboard Configuration	36
C	Prometheus Configuration	37
D	Troubleshooting Guide	38
D.1	Common Issues	38
D.1.1	Docker Not Found on HPC	38
D.1.2	Job Stays in PENDING	38
D.1.3	Service Connection Failed	38
D.1.4	Out of Memory	39

List of Figures

2.1	High-level system architecture of the AI Factory Benchmarking Framework showing the main orchestrator coordinating middleware, server management, client factory, and monitoring components.	10
5.1	HPC monitoring architecture with SSH tunneling for visualization access	20

List of Tables

4.1	Supported services and their capabilities	17
7.1	Technology stack	26
8.1	Expected performance ranges by service type	32

Chapter 1

Introduction

1.1 Background and Motivation

The European Union is actively investing in AI Factories—high-performance computing facilities designed to accelerate artificial intelligence research, development, and deployment across member states. These AI Factories harness the power of next-generation HPC and AI systems to revolutionize data processing, analytics, and model deployment. As these facilities come online, there is a critical need for comprehensive benchmarking methodologies to evaluate their performance characteristics and optimize their configurations.

The MeluXina supercomputer, operated by LuxProvide in Luxembourg, serves as an ideal platform for developing and testing such benchmarking frameworks. With its heterogeneous architecture combining CPU and GPU resources, MeluXina provides the computational capabilities necessary to evaluate the diverse workloads typical of AI Factory operations.

1.2 Project Objectives

The primary objectives of this project are:

1. **Design a unified benchmarking framework** to evaluate end-to-end performance for critical AI Factory components
2. **Implement modular benchmark components** covering:
 - LLM inference servers (Ollama, vLLM, Triton)
 - Vector databases (ChromaDB, Faiss, Weaviate)
 - Relational databases (PostgreSQL)

- Object storage (S3-compatible systems)
 - File storage (POSIX filesystem)
3. **Integrate production-grade monitoring** using Prometheus and Grafana
 4. **Enable reproducible benchmarking** through YAML-based configuration and SLURM orchestration
 5. **Provide comparative insights** with detailed performance metrics and statistical analysis

1.3 Challenge Overview

This project was conducted as part of the Student Challenge 2025-2026, spanning four months with the following phases:

Phase 1 – Onboarding: Introduction to MeluXina, exploration of tools and methodologies, design specification

Phase 2 – Prototyping: Development of applications, monitoring dashboards, and benchmarking scripts with iterative testing

Phase 3 – Evaluation: Deployment at realistic scales, performance measurements, resource profiling, and documentation

Phase 4 – Defense: Presentation of results and findings

1.4 Report Structure

This report is organized as follows:

- **Chapter 2:** Presents the system architecture and design decisions
- **Chapter 3:** Describes the core framework components in detail
- **Chapter 4:** Documents the supported services and their configurations
- **Chapter 5:** Explains the monitoring and visualization infrastructure
- **Chapter 6:** Details the HPC integration with SLURM and Apptainer
- **Chapter 7:** Covers implementation details and usage
- **Chapter 8:** Presents benchmark methodology and performance metrics

- **Chapter 9:** Summarizes conclusions and future work

Chapter 2

System Architecture

2.1 Architectural Overview

The AI Factory Benchmarking Framework follows a modular, layered architecture designed for extensibility, maintainability, and scalability. The system is composed of six core components orchestrated by a main controller, as illustrated in [Figure 2.1](#).

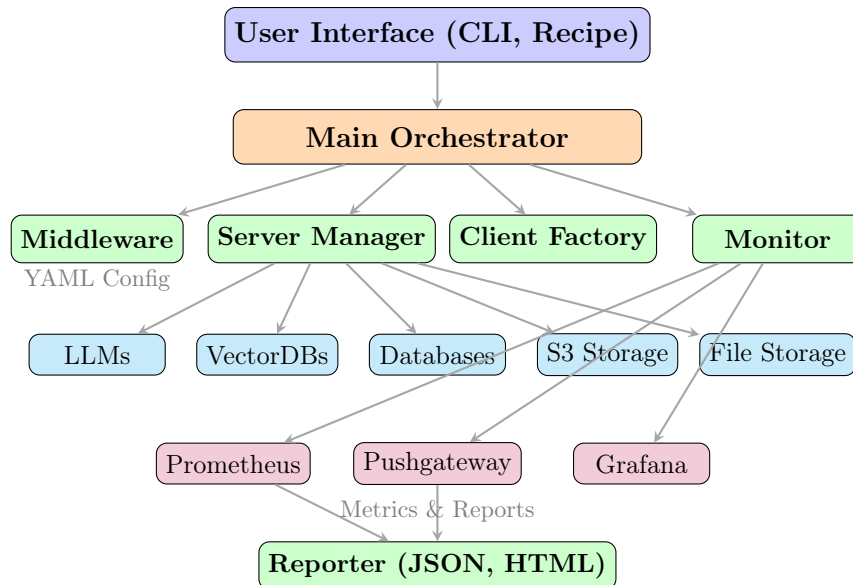


Figure 2.1: High-level system architecture of the AI Factory Benchmarking Framework showing the main orchestrator coordinating middleware, server management, client factory, and monitoring components.

2.2 Design Principles

The framework is built upon the following design principles:

Modularity: Each component is self-contained with well-defined interfaces, enabling independent development and testing

Extensibility: New service types can be added by implementing the appropriate client interface

Reproducibility: YAML-based configuration ensures benchmark scenarios can be exactly reproduced

Scalability: Parallel execution and SLURM integration support benchmarking at scale

Observability: Comprehensive metrics collection and real-time monitoring provide insights into system behavior

2.3 Data Flow

The benchmark execution follows a structured workflow:

1. **Configuration Loading:** The middleware interface loads and validates the recipe YAML configuration
2. **Service Initialization:** The server manager starts required services via SLURM and Apptainer
3. **Endpoint Discovery:** Service endpoints are retrieved from SLURM job outputs
4. **Benchmark Execution:** The client factory creates appropriate clients and executes benchmark operations
5. **Metrics Collection:** The monitor component records performance metrics in SQLite
6. **Report Generation:** The reporter generates comprehensive reports with statistical analysis
7. **Cleanup:** Services are stopped and resources are released

Chapter 3

Core Components

3.1 Main Orchestrator

The Main Orchestrator (`src/main.py`) serves as the central coordinator for the entire benchmarking workflow. It is responsible for:

- Initializing all framework components
- Managing the benchmark lifecycle
- Coordinating parallel execution of benchmarks
- Handling errors and cleanup procedures
- Generating final reports

3.1.1 Execution Modes

The orchestrator supports two execution modes:

Parallel Mode (Default): Uses Python's `ThreadPoolExecutor` to run benchmarks concurrently, achieving up to 10× faster execution

Sequential Mode: Executes benchmarks one at a time, useful for debugging and resource-constrained environments

```
1 # Parallel execution (default)
2 python src/main.py --recipe recipe.yml --max-workers 10
3
4 # Sequential execution
5 python src/main.py --recipe recipe.yml --sequential
```

Listing 3.1: Command-line execution options

3.2 Middleware Interface

The Middleware Interface (`src/middleware/interface.py`) provides the configuration management layer. It handles:

- Loading and parsing YAML recipe files
- Validating configuration parameters
- Providing service metadata to other components
- Managing benchmark identifiers

3.3 Server Manager

The Server Manager (`src/server/server.py`) manages the lifecycle of benchmark services running in Apptainer containers via SLURM. Key capabilities include:

- Automatic SLURM batch script generation
- Container image management (pulling from Docker registries)
- Service health checking
- Endpoint discovery and URL construction
- Graceful shutdown via SLURM job cancellation

3.3.1 SLURM Script Generation

The server manager dynamically generates SLURM submission scripts based on service configuration:

```
1 #!/bin/bash -l
2 #SBATCH --job-name=ollama-llama2
3 #SBATCH --time=01:00:00
4 #SBATCH --partition=gpu
5 #SBATCH --account=p200981
6 #SBATCH --nodes=1
7 #SBATCH --ntasks=1
8 #SBATCH --output=logs/ollama-llama2_%j.out
9
10 module add Apptainer
11 export COMPUTE_NODE=$(hostname -s)
12 export SERVICE_PORT=11434
13
14 # Save endpoint for discovery
```

```
15 echo "$COMPUTE_NODE:$SERVICE_PORT" > \  
16     $HOME/.cache/service_endpoints/ollama-llama2_endpoint.txt  
17  
18 # Run service in Apptainer container  
19 apptainer exec --nv containers/ollama.sif ollama serve &  
20 wait $!
```

Listing 3.2: Example generated SLURM script structure

3.4 Unified Client Factory

The Unified Client Factory (`src/client/unified_client.py`) creates specialized benchmark clients based on service type. It provides a unified interface while delegating to appropriate backend implementations:

- **BenchmarkClient:** HTTP-based clients for LLM inference (Ollama, vLLM, Triton)
- **PostgreSQLBenchmarkClient:** Database operations client
- **S3BenchmarkClient:** Object storage client
- **FileStorageBenchmarkClient:** POSIX filesystem client
- **VectorDBBenchmarkClient:** Vector database client (ChromaDB, Faiss, Weaviate)

3.5 Monitor

The Monitor component (`src/monitor/monitor.py`) collects and stores performance metrics in a SQLite database. Features include:

- Batch metric storage for efficiency
- Indexed queries by benchmark ID
- Raw data export to JSON
- Thread-safe operations for parallel benchmarking

The database schema captures:

- Timestamp
- Service name
- Client ID

- Request duration
- Success/failure status
- Status code
- Error messages
- Response size

3.6 Metrics Interceptor

The Metrics Interceptor (`src/interceptor/interceptor.py`) provides real-time metrics recording during benchmark execution. It captures:

- Request timing with microsecond precision
- Success/failure tracking
- Error classification
- Response size measurements

3.7 Reporter

The Reporter component (`src/reporter/reporter.py`) generates comprehensive reports from collected metrics. Report features include:

- Summary statistics (total, successful, failed requests)
- Timing analysis (average, median, min, max, standard deviation)
- Percentile calculations (P50, P90, P95, P99)
- Per-service breakdown
- Per-client analysis
- Throughput calculations (requests per second)

3.7.1 Report Format

Reports are generated in JSON format with the following structure:

```
1 {  
2   "benchmark_id": "benchmark_1234567890",  
3   "generated_at": "2025-01-04T12:00:00",
```



```
4  "summary": {
5    "total_requests": 3000,
6    "successful_requests": 2950,
7    "failed_requests": 50,
8    "success_rate": 98.33,
9    "avg_duration": 0.523,
10   "median_duration": 0.501
11 },
12 "services": {
13   "service-name": {
14     "timing": {...},
15     "percentiles": {"p50": 0.501, "p95": 0.812},
16     "throughput": {"requests_per_second": 9.98}
17   }
18 }
19 }
```

Listing 3.3: Example report structure

Chapter 4

Supported Services

The AI Factory Benchmarking Framework supports comprehensive benchmarking across five major service categories, as summarized in [Table 4.1](#).

Table 4.1: Supported services and their capabilities

Service Category	Backends	GPU Support	Protocol
LLM Inference	Ollama, vLLM, Triton	✓	HTTP/gRPC
Vector Databases	ChromaDB, Faiss, Weaviate	✓	HTTP/Native
Relational Databases	PostgreSQL	✓	psycopg2
Object Storage	MinIO (S3-compatible)	✓	S3 API
File Storage	POSIX Filesystem	✓	Direct I/O

4.1 LLM Inference Services

4.1.1 Ollama

Ollama provides a simple HTTP API for running large language models locally. The framework supports:

- Model selection (Llama2, Mistral, etc.)
- Configurable prompt generation
- Response streaming support

4.1.2 vLLM

vLLM is a high-performance inference engine optimized for large language models. Features include:

- OpenAI-compatible API
- PagedAttention for efficient memory management
- Continuous batching for high throughput

4.1.3 Triton Inference Server

NVIDIA Triton provides a platform for deploying AI models at scale:

- Multi-framework support (TensorRT, PyTorch, TensorFlow)
- Dynamic batching
- Model versioning

4.2 Vector Databases

4.2.1 ChromaDB

ChromaDB is an AI-native open-source embedding database:

- HTTP REST API
- Collection-based organization
- Configurable embedding dimensions

4.2.2 Faiss

Facebook AI Similarity Search (Faiss) provides efficient similarity search:

- In-memory operation (no server required)
- GPU-accelerated indexing
- Multiple index types (Flat, IVF, HNSW)

4.2.3 Weaviate

Weaviate is a cloud-native vector search engine:

- GraphQL and REST APIs
- Semantic search capabilities
- Hybrid search support

4.3 Relational Databases

4.3.1 PostgreSQL

PostgreSQL benchmarking supports standard database operations:

- SELECT, INSERT, UPDATE, DELETE operations
- Configurable operation mix
- Connection pooling support
- Batch operations

4.4 Storage Systems

4.4.1 Object Storage (MinIO)

MinIO provides S3-compatible object storage:

- PUT, GET, LIST, DELETE operations
- Configurable object sizes
- Bucket management
- Multi-part upload support

4.4.2 File Storage

POSIX filesystem benchmarking supports:

- Read, Write, Stat, Delete operations
- Configurable file sizes
- Direct I/O option
- Sync mode for durability testing

Chapter 5

Monitoring and Visualization

5.1 Monitoring Architecture

The framework integrates a production-grade monitoring stack based on Prometheus and Grafana, specifically designed for HPC environments. The architecture, shown in [Figure 5.1](#), supports deployment on MeluXina compute nodes using Apptainer containers.

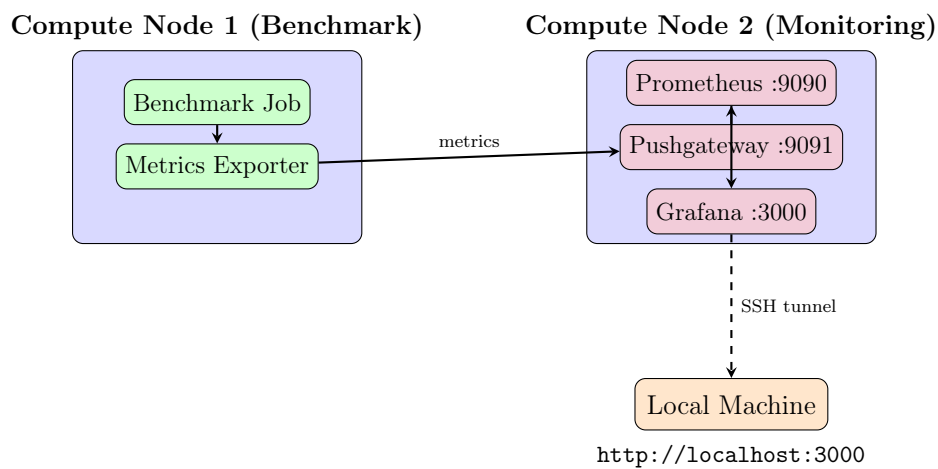


Figure 5.1: HPC monitoring architecture with SSH tunneling for visualization access

5.2 Prometheus Integration

Prometheus serves as the time-series database for metrics storage:

- Scrapes metrics from Pushgateway
- Stores time-series data for historical analysis
- Supports PromQL queries for data exploration

- Configurable retention and alerting

5.2.1 Key Metrics

The framework exports the following Prometheus metrics:

```
1 # Request rate (throughput)
2 rate(benchmark_requests_total[1m])
3
4 # Average latency
5 benchmark_request_duration_seconds_sum /
6 benchmark_request_duration_seconds_count
7
8 # Success rate
9 (benchmark_requests_successful /
10 benchmark_requests_total) * 100
11
12 # P95 latency
13 histogram_quantile(0.95,
14   rate(benchmark_request_duration_seconds_bucket[5m]))
```

Listing 5.1: Prometheus metric queries

5.3 Pushgateway

The Pushgateway serves as an intermediary for batch and ephemeral jobs to push metrics to Prometheus. This is essential for HPC environments where benchmark jobs have limited lifetimes:

```
1 python src/monitoring/prometheus_exporter.py benchmark_id \
2   --push \
3   --gateway-url http://mel2110:9091 \
4   --job ai_factory_benchmark
```

Listing 5.2: Pushing metrics to Pushgateway

5.4 Grafana Dashboards

Grafana provides visualization and analytics capabilities. The pre-configured dashboard includes:

- **Request Rate (Throughput):** Time-series graph of requests per second
- **Success Rate Gauge:** Visual indicator with configurable thresholds

- **Response Time:** Latency graph with percentile overlays
- **Total Requests Counter:** Stat panel for overall volume
- **Failed Requests Counter:** Error tracking
- **Requests by Service:** Pie chart for workload distribution
- **Success vs Failure:** Stacked bar chart

5.5 Accessing Monitoring from Local Machine

Since compute nodes are not directly accessible from external networks, SSH tunneling is used to access Grafana:

```
1 # Create SSH tunnel to compute node
2 ssh -L 3000:mel2345:3000 \
3     -L 9090:mel2345:9090 \
4     -L 9091:mel2345:9091 \
5     <username>@login.lxp.lu -p 8822
6
7 # Access Grafana at http://localhost:3000
8 # Login: admin / admin
```

Listing 5.3: SSH tunnel setup for monitoring access

Chapter 6

HPC Integration

6.1 MeluXina Supercomputer

MeluXina is Luxembourg's petascale supercomputer, operated by LuxProvide. It provides:

- Heterogeneous architecture with CPU and GPU nodes
- SLURM workload manager for job scheduling
- High-speed interconnect for distributed computing
- Shared filesystems for data management

6.2 SLURM Integration

The framework provides deep integration with SLURM for job management:

6.2.1 Configuration Parameters

```
1 slurm:
2   partition: gpu           # or 'cpu'
3   account: p200981        # Project account
4   time: "01:00:00"       # Maximum runtime
5   nodes: 1                # Number of nodes
6   ntasks: 1               # Tasks per job
7   qos: default            # Quality of service
8   gres: gpu:1             # GPU resources
```

Listing 6.1: SLURM configuration in recipe.yml

6.2.2 Job Lifecycle

1. **Script Generation:** Dynamic SLURM scripts based on service configuration
2. **Job Submission:** Automatic `sbatch` submission
3. **Status Monitoring:** `squeue` integration for job tracking
4. **Endpoint Discovery:** Reading node hostname from job output
5. **Job Cancellation:** `scancel` for cleanup

6.3 Apptainer Containerization

Apptainer (formerly Singularity) provides containerization support for HPC environments:

6.3.1 Container Management

```
1 # Pull containers from Docker registry
2 apptainer pull ollama.sif docker://ollama/ollama
3 apptainer pull vllm.sif docker://vllm/vllm-openai:latest
4 apptainer pull chromadb.sif docker://chromadb/chroma:latest
5
6 # Run with GPU support
7 apptainer exec --nv chromadb.sif chroma run --port 8000
```

Listing 6.2: Container operations

6.3.2 Bind Mounts for Persistent Data

Different services require specific bind mounts for data persistence:

- **ChromaDB:** `-bind $DATA_DIR:/chroma/data`
- **PostgreSQL:** `-bind $PGDATA:/var/lib/postgresql/data`
- **MinIO:** `-bind $MINIO_DATA:/data`
- **Weaviate:** `-bind $WEAVIATE_DATA:/var/lib/weaviate`

6.4 GPU Acceleration

All services support GPU acceleration through NVIDIA containers:

- Use `-nv` flag with Apptainer for GPU passthrough

-
- Request GPU resources via SLURM `gres` parameter
 - LLM inference benefits significantly from GPU acceleration
 - Vector databases can use GPU for accelerated similarity search

Chapter 7

Implementation Details

7.1 Technology Stack

The framework is built using the following technologies:

Table 7.1: Technology stack

Category	Technologies
Programming Language	Python 3.12
Orchestration	SLURM, Python ThreadPoolExecutor
Containerization	Apptainer (Singularity)
Monitoring	Prometheus, Pushgateway, Grafana
Data Storage	SQLite (metrics), JSON (reports)
GPU Support	NVIDIA CUDA via Apptainer <code>-nv</code>
Configuration	YAML (PyYAML)
HTTP Client	requests
CLI	rich, argparse

7.2 Dependencies

Core Python dependencies are managed via `requirements.txt`:

```
1 PyYAML>=6.0          # Configuration
2 requests>=2.31.0      # HTTP client
3 numpy>=1.24.0         # Numerical operations
4 pandas>=2.0.0         # Data analysis
5 psycopg2-binary>=2.9.0 # PostgreSQL
6 boto3>=1.28.0         # S3 client
7 chromadb>=0.4.0       # ChromaDB client
8 faiss-cpu>=1.7.4      # Faiss
9 weaviate-client>=3.24.0 # Weaviate
```

```
10 prometheus-client>=0.18.0 # Metrics export
11 rich>=13.0.0             # CLI output
12 matplotlib>=3.7.0        # Visualization
```

Listing 7.1: Core dependencies

7.3 Recipe Configuration System

The YAML-based recipe system enables reproducible benchmark configurations:

```
1 benchmark:
2   name: "AI Factory Comprehensive Benchmark"
3   description: "Full-scale benchmarking of AI components"
4   duration: 3600
5
6 global:
7   log_level: INFO
8   metrics_db: metrics.db
9   reports_dir: reports
10  pushgateway_url: http://localhost:9091
11
12 services:
13   # LLM Inference
14   - service_name: ollama-llama2
15     service_type: ollama
16     container_image: docker://ollama/ollama
17     port: 11434
18     client_count: 5
19     requests_per_second: 10
20     duration: 300
21     model: llama2
22     slurm:
23       partition: gpu
24       time: "01:00:00"
25       gres: gpu:1
26
27   # Vector Database
28   - service_name: chromadb
29     service_type: vectordb
30     backend: chromadb
31     port: 8000
32     dimension: 384
33     operation_mix:
34       insert: 0.3
35       search: 0.5
36       update: 0.1
```

```
37 delete: 0.1
```

Listing 7.2: Comprehensive recipe example

7.4 Installation and Setup

7.4.1 HPC Installation (MeluXina)

```
1 # Connect to MeluXina
2 ssh <username>@login.lxp.lu -p 8822
3
4 # Navigate to project directory
5 cd $PROJECT/p200981
6
7 # Load required modules
8 module purge
9 module add Apptainer
10 module load Python/3.12.3-GCCcore-13.3.0
11
12 # Clone repository
13 git clone https://github.com/Valegr1/Team6_EUMASTER4HPC2526.git
14 cd Team6_EUMASTER4HPC2526
15
16 # Create virtual environment
17 python -m venv venv
18 source venv/bin/activate
19 pip install -r requirements.txt
20
21 # Pull service containers
22 ./scripts/pull_containers.sh
```

Listing 7.3: Installation steps for MeluXina

7.5 Usage Examples

7.5.1 Running Benchmarks

```
1 # Basic execution
2 python src/main.py --recipe recipe.yml
3
4 # With Pushgateway for monitoring
5 python src/main.py --recipe recipe.yml \
6     --pushgateway http://mel2110:9091
7
8 # Parallel execution with custom workers
```

```
9 python src/main.py --recipe recipe.yml --max-workers 10
10
11 # Sequential execution for debugging
12 python src/main.py --recipe recipe.yml --sequential
```

Listing 7.4: Benchmark execution examples

7.5.2 SLURM Batch Execution

```
1 # Submit benchmark job
2 sbatch scripts/run_benchmark.sh
3
4 # Monitor job status
5 squeue -u $USER
6
7 # View job output
8 tail -f logs/benchmark_JOBID.out
```

Listing 7.5: SLURM job submission

7.5.3 Interactive Session

```
1 # Request interactive session
2 salloc --partition=gpu --account=p200981 \
3       --time=01:00:00 --nodes=1
4
5 # Load modules and run
6 module add Apptainer
7 module load Python/3.12.3-GCCcore-13.3.0
8 source venv/bin/activate
9 python src/main.py --recipe recipe.yml
```

Listing 7.6: Interactive SLURM session

Chapter 8

Benchmark Methodology and Results

8.1 Benchmark Methodology

The framework follows established benchmarking best practices:

8.1.1 Warm-up Phase

Each benchmark begins with a warm-up period to:

- Initialize service caches
- Establish connection pools
- Load models into GPU memory
- Stabilize system performance

8.1.2 Statistical Rigor

To ensure meaningful results:

- Multiple runs (3-5) are recommended for averaging
- Sufficiently long durations (300-600 seconds) capture variance
- Control variables by changing one parameter at a time
- Report confidence intervals when appropriate

8.2 Performance Metrics

8.2.1 Latency Metrics

- **Average Duration:** Mean request processing time
- **Median Duration (P50):** 50th percentile latency
- **P90/P95/P99:** Tail latency percentiles
- **Min/Max:** Extreme values for outlier detection
- **Standard Deviation:** Latency consistency measure

8.2.2 Throughput Metrics

- **Requests per Second:** Overall throughput rate
- **Total Requests:** Volume of operations completed
- **Duration:** Total benchmark runtime

8.2.3 Reliability Metrics

- **Success Rate:** Percentage of successful requests
- **Failed Requests:** Count of failed operations
- **Error Classification:** Types of failures encountered

8.3 Result Analysis

8.3.1 Querying Metrics Database

```
1  -- Summary statistics by service
2  SELECT
3      service_name ,
4      COUNT(*) as total ,
5      SUM(CASE WHEN success = 1 THEN 1 ELSE 0 END) as successful ,
6      AVG(request_duration) as avg_duration ,
7      MIN(request_duration) as min_duration ,
8      MAX(request_duration) as max_duration
9  FROM metrics
10 WHERE benchmark_id = 'benchmark_1234567890'
11 GROUP BY service_name;
12
13 -- Percentile calculation
```



```

14 SELECT
15     service_name ,
16     request_duration
17 FROM metrics
18 WHERE success = 1
19 ORDER BY request_duration;

```

Listing 8.1: SQLite query examples

8.3.2 Comparing Benchmarks

```

1 # Compare multiple runs
2 for report in reports/benchmark*_report.json; do
3     echo "=== $(basename $report) ==="
4     jq '.summary.success_rate, .timing.avg_duration' $report
5 done

```

Listing 8.2: Comparing multiple benchmark runs

8.4 Expected Performance Characteristics

Based on the architecture and typical HPC deployments, expected performance characteristics include:

Table 8.1: Expected performance ranges by service type

Service Type	Latency (P50)	Throughput	GPU Impact
LLM Inference (Ollama)	100-500ms	10-50 req/s	5-10×
LLM Inference (vLLM)	50-200ms	20-100 req/s	10-20×
Vector DB (ChromaDB)	5-50ms	100-500 req/s	2-5×
Vector DB (Faiss)	1-10ms	1000+ req/s	5-10×
PostgreSQL	1-10ms	500-2000 req/s	Minimal
S3 Storage	10-100ms	100-500 req/s	N/A
File Storage	0.1-10ms	1000+ req/s	N/A

Chapter 9

Conclusion and Future Work

9.1 Summary of Achievements

This project has successfully delivered a comprehensive AI Factory Benchmarking Framework with the following key achievements:

1. **Unified Framework:** A single, cohesive framework supporting 10+ AI infrastructure services
2. **Production-Grade Monitoring:** Integration with Prometheus and Grafana for real-time visualization
3. **HPC Integration:** Seamless SLURM orchestration and Apptainer containerization
4. **Reproducible Benchmarks:** YAML-driven configuration for exact reproduction of experiments
5. **GPU Acceleration:** Support for NVIDIA GPUs across all service types
6. **Parallel Execution:** ThreadPoolExecutor enables up to 10× faster benchmarking
7. **Comprehensive Reporting:** Detailed statistical analysis with percentile calculations

9.2 Lessons Learned

Key insights gained during development:

- **Container Compatibility:** Apptainer provides excellent Docker compatibility for HPC environments

- **Endpoint Discovery:** Dynamic service endpoint discovery is essential for SLURM-based deployments
- **Monitoring Architecture:** Compute node-based monitoring eliminates SSH tunneling complexity for metrics collection
- **Concurrency Benefits:** Thread-based parallelism proved effective for I/O-bound benchmark workloads
- **Configuration Management:** YAML-based recipes greatly simplify reproducibility

9.3 Future Work

Several enhancements are planned for future development:

Multi-Node Benchmarking: Extend support for distributed benchmarks across multiple compute nodes

Advanced Triton Models: Add pre-configured model repositories for Triton Inference Server

Scaling Analysis: Implement automated scaling studies with varying client counts and node configurations

EU AI Factory Integration: Align with emerging European AI Factory standards and requirements

Machine Learning Workloads: Add benchmarks for training workloads alongside inference

Network Analysis: Include network bandwidth and latency measurements between services

Cost Modeling: Integrate resource utilization with cost estimation for cloud deployments

9.4 Acknowledgments

We thank Dr. Farouk Mansouri of LuxProvide for his supervision and guidance throughout this project. We also acknowledge LuxProvide for providing access to the MeluXina supercomputer and the EUMASTER4HPC program for the opportunity to work on this challenge.

Appendix A

Apptainer Container Recipes

The framework includes Apptainer definition files for building custom containers when needed:

- `chromadb.def` – ChromaDB vector database
- `minio.def` – MinIO S3-compatible storage
- `ollama.def` – Ollama LLM inference
- `postgres.def` – PostgreSQL database
- `triton.def` – NVIDIA Triton Inference Server
- `vllm.def` – vLLM high-performance inference
- `weaviate.def` – Weaviate vector search

Appendix B

Grafana Dashboard Configuration

The pre-configured Grafana dashboard is available at:

`monitoring/grafana/dashboards/ai-factory-benchmark.json`

Dashboard panels include:

- Request Rate (RPS) time series
- Success Rate gauge
- Response Time graph with percentiles
- Total Requests stat
- Failed Requests stat
- Requests by Service pie chart
- Success vs Failure stacked bar

Appendix C

Prometheus Configuration

Prometheus configuration is located at:

`monitoring/prometheus/prometheus.yml`

The configuration includes:

- Scrape configurations for Pushgateway
- Alerting rules in `alerts.yml`
- Retention settings

Appendix D

Troubleshooting Guide

D.1 Common Issues

D.1.1 Docker Not Found on HPC

Issue: “docker: command not found” on MeluXina

Solution: Docker is not available on HPC systems. Use Apptainer instead:

```
1 module add Apptainer
2 apptainer run container.sif
```

D.1.2 Job Stays in PENDING

Solution:

```
1 # Check reason
2 squeue -u $USER --start
3
4 # Check partition availability
5 sinfo -p gpu
```

D.1.3 Service Connection Failed

Solution:

```
1 # Check service logs
2 grep -i error logs/benchmark_*.log
3
4 # Verify container is running
5 singularity instance list
```

D.1.4 Out of Memory

Solution:

```
1 # Request more memory
2 sbatch --mem=64G scripts/run_benchmark.sh
```


Bibliography

- [1] LuxProvide. MeluXina Documentation. <https://docs.lxp.lu/>
- [2] Prometheus Authors. Prometheus - Monitoring system & time series database. <https://prometheus.io/>
- [3] Grafana Labs. Grafana - The open observability platform. <https://grafana.com/>
- [4] SchedMD. Slurm Workload Manager. <https://slurm.schedmd.com/>
- [5] Apptainer Community. Apptainer Documentation. <https://apptainer.org/docs/>
- [6] Ollama. Ollama - Get up and running with large language models locally. <https://ollama.ai/>
- [7] vLLM Team. vLLM: Easy, Fast, and Cheap LLM Serving. <https://vllm.ai/>
- [8] NVIDIA. NVIDIA Triton Inference Server. <https://developer.nvidia.com/nvidia-triton-inference-server>
- [9] Chroma. ChromaDB - AI-native open-source embedding database. <https://www.trychroma.com/>