

**Department of Artificial Intelligence & Data Science****Vision of the Department**

To be a well-known centre for pursuing computer education through innovative pedagogy, value-based education and industry collaboration.

Mission of the Department

To establish learning ambience for ushering in computer engineering professionals in core and multidisciplinary area by developing Problem-solving skills through emerging technologies.

Session 2025-2026**Vision:** Dream of where you want.**Mission:** Means to achieve Vision

Program Educational Objectives of the program (PEO): (broad statements that describe the professional and career accomplishments)

| | | | |
|------|---------------------------------|--|--|
| PEO1 | Preparation | P: Preparation | Pep-CL abbreviation pronounce as Pep-si-IL easy to recall |
| PEO2 | Core Competence | E: Environment (Learning Environment) | |
| PEO3 | Breadth | P: Professionalism | |
| PEO4 | Professionalism | C: Core Competence | |
| PEO5 | Learning Environment | L: Breadth (Learning in diverse areas) | |

Program Outcomes (PO):

1. Understand and Apply Parallel Programming Concepts
2. Analyse and Improve Program Performance.
3. Demonstrate Practical Skills in HPC Tools and Environments.

Keywords of POs:

Engineering knowledge, Problem analysis, Design/development of solutions, Conduct Investigations of Complex Problems, Engineering Tool Usage, The Engineer and The World, Ethics, Individual and Collaborative Team work, Communication, Project Management and Finance, Life-Long Learning

PSO Keywords: Cutting edge technologies, Research

“I am an engineer, and I know how to apply engineering knowledge to investigate, analyse and design solutions to complex problems using tools for entire world following all ethics in a collaborative way with proper management skills throughout my life.” to contribute to the development of cutting-edge technologies and Research.

Integrity: I will adhere to the Laboratory Code of Conduct and ethics in its entirety.

Name and Signature of Student and Date

Sakshi Gokhale

04/11/25



Department of Artificial Intelligence & Data Science

Vision of the Department

To be a well-known centre for pursuing computer education through innovative pedagogy, value-based education and industry collaboration.

Mission of the Department

To establish learning ambience for ushering in computer engineering professionals in core and multidisciplinary area by developing Problem-solving skills through emerging technologies.

| | | | |
|-----------------|---------------|------------------------|----------------|
| Session | 2025-26 (ODD) | Course Name | HPC Lab |
| Semester | 7 AIDS | Course Code | 22ADS706 |
| Roll No | 16 | Name of Student | Sakshi Gokhale |

| | |
|-------------------------------------|---|
| Practical Number | 9 (Mini Project) |
| Course Outcome | 1. Understand and Apply Parallel Programming Concepts 2. Analyse and Improve Program Performance |
| Aim | Performance Comparison of Data Serialization Formats |
| Problem Definition | Performance Comparison of Data Serialization Formats |
| Theory (100 words) | In High Performance Computing (HPC), huge volumes of structured numerical data are generated from simulations, scientific experiments, sensors, and analytical workloads. Efficient storage and fast data movement are essential to reduce execution time. Data serialization formats define how data gets converted into byte stream for storage or transfer. Different formats provide different performance characteristics. Text-based formats like CSV and JSON are human-readable but very slow in I/O operations and consume more disk space. Pickle is fast but it is Python specific and not suitable for cross-platform HPC environments. Parquet is a modern columnar binary format which stores data in compressed, vectorizable columns. This improves caching, minimizes storage size and supports parallel read operations. Therefore selection of the right serialization format is important for HPC pipelines. Parquet achieves high throughput and low memory usage, making it suitable for large scale data analysis and distributed computational systems. |
| Procedure and Execution (100 Words) | Code: import pandas as pd, numpy as np, time df = pd.DataFrame({ "A": np.random.rand(5_000_000), "B": np.random.rand(5_000_000), "C": np.random.rand(5_000_000) }) |



Department of Artificial Intelligence & Data Science

Vision of the Department

To be a well-known centre for pursuing computer education through innovative pedagogy, value-based education and industry collaboration.

Mission of the Department

To establish learning ambience for ushering in computer engineering professionals in core and multidisciplinary area by developing Problem-solving skills through emerging technologies.

```
})

# write timings
for fmt, fn, save in [
    ("CSV", "test.csv", lambda: df.to_csv("test.csv")),

    ("JSON", "test.json", lambda: df.to_json("test.json", orient='records')),
    ("Pickle", "test.pkl", lambda: df.to_pickle("test.pkl")),

    ("Parquet", "test.parquet", lambda: df.to_parquet("test.parquet", engine='pyarrow'))
]:
    t=time.time(); save(); print(fmt, "write:", time.time()-t)

print("-----READING-----")

# read timings
for fmt, load in [
    ("CSV", lambda: pd.read_csv("test.csv")),
    ("JSON", lambda: pd.read_json("test.json")),
    ("Pickle", lambda: pd.read_pickle("test.pkl")),
    ("Parquet", lambda: pd.read_parquet("test.parquet"))
]:
    t=time.time(); load(); print(fmt, "read:", time.time()-t)
```

Output:

```
[lab1@localhost ~]$ python3 -m ensurepip --user
Looking in links: /tmp/tmpwjw5ydlk
Requirement already satisfied: setuptools in /usr/lib/python3.9/site-packages (53.0.0)
Processing /tmp/tmpwjw5ydlk/pip-21.3.1-py3-none-any.whl
Installing collected packages: pip
Successfully installed pip-21.3.1
[lab1@localhost ~]$ python3 -m pip install --user pandas pyarrow
Collecting pandas
  Downloading pandas-2.3.3-cp39-cp39-manylinux_2_24_x86_64.manylinux_2_28_x86_64.whl (12.8 MB)
    |#####| 12.8 MB 6.0 kB/s
Collecting pyarrow
  Downloading pyarrow-21.0.0-cp39-cp39-manylinux_2_28_x86_64.whl (42.7 MB)
    |#####| 42.7 MB 256 kB/s ^[[24~
Collecting python-dateutil<=2.8.2
  Downloading python_dateutil-2.9.0.post0-py2.py3-none-any.whl (229 kB)
    |#####| 229 kB 271 kB/s
Collecting tzdata>=2022.7
  Downloading tzdata-2025.2-py2.py3-none-any.whl (347 kB)
    |#####| 347 kB 280 kB/s
Collecting pytz>=2020.1
  Downloading pytz-2025.2-py2.py3-none-any.whl (509 kB)
    |#####| 509 kB 250 kB/s
Collecting numpy>=1.22.4
  Downloading numpy-2.0.2-cp39-cp39-manylinux_2_17_x86_64.manylinux2014_x86_64.whl (19.5 MB)
    |#####| 19.5 MB 14 kB/s
Requirement already satisfied: six>=1.5 in /usr/lib/python3.9/site-packages (from python-dateutil>=2.8.2->pandas) (1.15.0)
Installing collected packages: tzdata, pytz, python-dateutil, numpy, pyarrow, pandas
Successfully installed numpy-2.0.2 pandas-2.3.3 pyarrow-21.0.0 python-dateutil-2.9.0.post0 pytz-2025.2 tzdata-2025.2
```



Department of Artificial Intelligence & Data Science

Vision of the Department

To be a well-known centre for pursuing computer education through innovative pedagogy, value-based education and industry collaboration.

Mission of the Department

To establish learning ambience for ushering in computer engineering professionals in core and multidisciplinary area by developing Problem-solving skills through emerging technologies.

```
lab1@localhost:~  
Successfully installed numpy-2.0.2 pandas-2.3.3 pyarrow-21.0.0 python-dateutil-2.9.0.post0 pytz-2025.2 tzdata-2025.2  
WARNING: You are using pip version 21.3.1; however, version 25.3 is available.  
You should consider upgrading via the '/usr/bin/python3 -m pip install --upgrade pip' command.  
[lab1@localhost ~]$ ~python3  
bash: ~python3: command not found...  
Similar command is: 'python3'  
[lab1@localhost ~]$ python3  
Python 3.9.23 (main, Aug 19 2025, 00:00:00)  
[GCC 11.5.0 20240719 (Red Hat 11.5.0-11)] on linux  
Type "help", "copyright", "credits" or "license()" for more information.  
>>> import pandas as pd, numpy as np, time  
>>>  
>>> df = pd.DataFrame({  
...     "A": np.random.rand(5_000_000),  
...     "B": np.random.rand(5_000_000),  
...     "C": np.random.rand(5_000_000)  
... })  
>>>  
>>> # write timings  
>>> for fmt, fn, save in [  
...     ("CSV", "test.csv", lambda: df.to_csv("test.csv")),  
...     ("JSON", "test.json", lambda: df.to_json("test.json", orient='records')),  
...     ("Pickle", "test.pkl", lambda: df.to_pickle("test.pkl")),  
...     ("Parquet", "test.parquet", lambda: df.to_parquet("test.parquet", engine='pyarrow'))  
... ]:  
...     t=time.time(); save(); print(fmt, "write:", time.time()-t)  
...
```

```
lab1@localhost:~  
CSV write: 14.345051527023315  
JSON write: 2.7007336616516113  
Pickle write: 0.06801319122314453  
Parquet write: 0.46419453620910645  
>>> print("-----READING-----")  
-----READING-----  
>>>  
>>> # read timings  
>>> for fmt, load in [  
...     ("CSV", lambda: pd.read_csv("test.csv")),  
...     ("JSON", lambda: pd.read_json("test.json")),  
...     ("Pickle", lambda: pd.read_pickle("test.pkl")),  
...     ("Parquet", lambda: pd.read_parquet("test.parquet"))  
... ]:  
...     t=time.time(); load(); print(fmt, "read:", time.time()-t)  
...  
... Unnamed: 0      A      B      C  
0      0  0.668127  0.465055  0.140999  
1      1  0.861483  0.352243  0.218870  
2      2  0.901892  0.687647  0.122469  
3      3  0.085235  0.298829  0.630436  
4      4  0.949967  0.265607  0.700942  
...      ...      ...      ...  
4999995 4999995 0.262878  0.301076  0.566506  
4999996 4999996 0.151081  0.017528  0.088710  
4999997 4999997 0.938092  0.572218  0.287541  
4999998 4999998 0.899285  0.138958  0.140682  
4999999 4999999 0.767493  0.809369  0.129159  
[5000000 rows x 4 columns]  
CSV read: 1.6857192516326904  
...      A      B      C  
0      0  0.668127  0.465055  0.140999  
1      1  0.861483  0.352243  0.218870  
2      2  0.901892  0.687647  0.122469  
3      3  0.085235  0.298829  0.630436  
4      4  0.949967  0.265607  0.700942
```



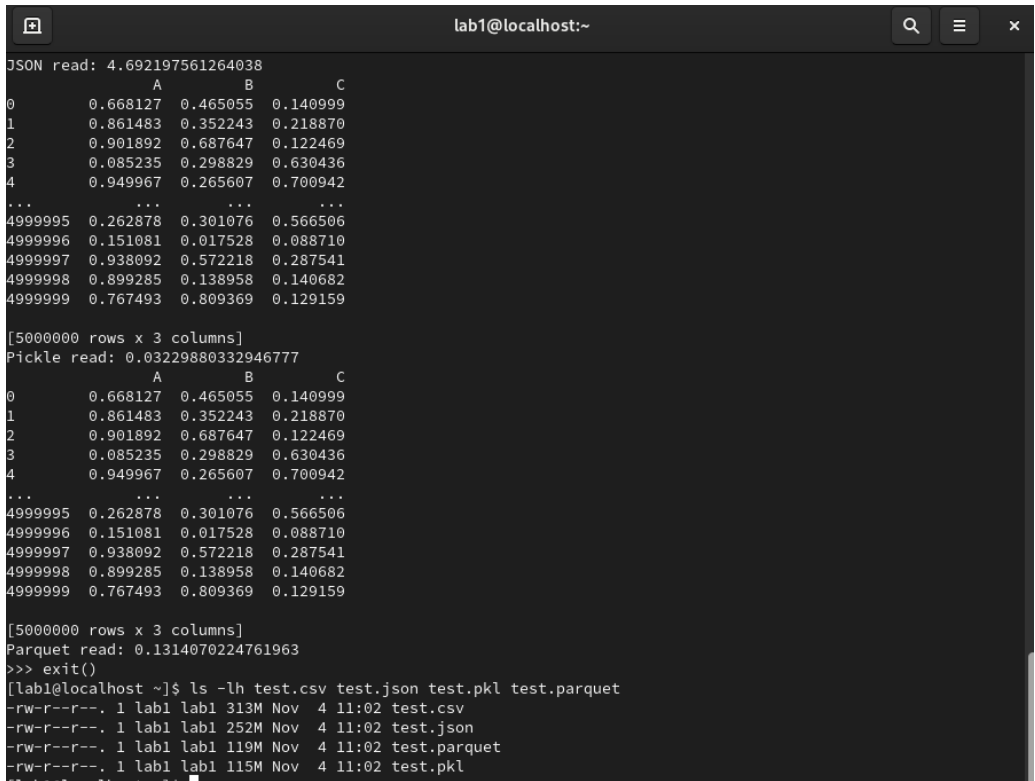
Department of Artificial Intelligence & Data Science

Vision of the Department

To be a well-known centre for pursuing computer education through innovative pedagogy, value-based education and industry collaboration.

Mission of the Department

To establish learning ambience for ushering in computer engineering professionals in core and multidisciplinary area by developing Problem-solving skills through emerging technologies.

| | |
|---|--|
| |  <pre>lab1@localhost:~ JSON read: 4.692197561264038 A B C 0 0.668127 0.465055 0.140999 1 0.861483 0.352243 0.218870 2 0.901892 0.687647 0.122469 3 0.085235 0.298829 0.630436 4 0.949967 0.265607 0.700942 4999995 0.262878 0.301076 0.566506 4999996 0.151081 0.017528 0.088710 4999997 0.938092 0.572218 0.287541 4999998 0.899285 0.138958 0.140682 4999999 0.767493 0.809369 0.129159 [5000000 rows x 3 columns] Pickle read: 0.03229880332946777 A B C 0 0.668127 0.465055 0.140999 1 0.861483 0.352243 0.218870 2 0.901892 0.687647 0.122469 3 0.085235 0.298829 0.630436 4 0.949967 0.265607 0.700942 4999995 0.262878 0.301076 0.566506 4999996 0.151081 0.017528 0.088710 4999997 0.938092 0.572218 0.287541 4999998 0.899285 0.138958 0.140682 4999999 0.767493 0.809369 0.129159 [5000000 rows x 3 columns] Parquet read: 0.1314070224761963 >>> exit() [lab1@localhost ~]\$ ls -lh test.csv test.json test.pkl test.parquet -rw-r--r--. 1 lab1 lab1 313M Nov 4 11:02 test.csv -rw-r--r--. 1 lab1 lab1 252M Nov 4 11:02 test.json -rw-r--r--. 1 lab1 lab1 119M Nov 4 11:02 test.parquet -rw-r--r--. 1 lab1 lab1 115M Nov 4 11:02 test.pkl</pre> |
| Output Analysis | CSV and JSON formats produced the largest file sizes and slowest read/write times. Pickle was comparatively faster but it is language dependent. Parquet generated the smallest file size and showed the best performance in both write and read operations. Therefore, the experimental observation proves that columnar optimized formats perform significantly better for large HPC datasets. |
| Link of student Github profile where lab assignment | https://github.com/sakshi-gokhale/Lab-HPC |

**Department of Artificial Intelligence & Data Science****Vision of the Department**

To be a well-known centre for pursuing computer education through innovative pedagogy, value-based education and industry collaboration.

Mission of the Department

To establish learning ambience for ushering in computer engineering professionals in core and multidisciplinary area by developing Problem-solving skills through emerging technologies.

| | | | | | | | | | | | | | |
|--------------------------------------|---|-------|-----|------------|------|-----------|---|------------|----|-----------|-------------|------------|-------------|
| has been uploaded | | | | | | | | | | | | | |
| Conclusion | Parquet format is the most efficient serialization format for HPC because it is columnar, compressed and supports very fast reading operations. CSV and JSON are slower and not ideal for HPC scale data. Pickle is fast but not portable. Hence, Parquet is recommended for large scale scientific and HPC workflows. | | | | | | | | | | | | |
| Plag Report (Similarity index < 12%) | <div>Plagiarism Scan Report <p>3% Plagiarism</p><p>3% Exact Match</p><p>0% Partial Match</p><p>97% Unique</p><table><tr><td>Words</td><td>226</td></tr><tr><td>Characters</td><td>1461</td></tr><tr><td>Sentences</td><td>9</td></tr><tr><td>Paragraphs</td><td>14</td></tr><tr><td>Read Time</td><td>2 minute(s)</td></tr><tr><td>Speak Time</td><td>2 minute(s)</td></tr></table><p>Content Checked For Plagiarism</p><p>In High Performance Computing (HPC), huge volumes of structured numerical data are generated from simulations, scientific experiments, sensors, and analytical workloads. Efficient storage and fast data movement are essential to reduce execution time. Data serialization formats define how data gets converted into byte stream for storage or transfer. Different formats provide different</p></div> | Words | 226 | Characters | 1461 | Sentences | 9 | Paragraphs | 14 | Read Time | 2 minute(s) | Speak Time | 2 minute(s) |
| Words | 226 | | | | | | | | | | | | |
| Characters | 1461 | | | | | | | | | | | | |
| Sentences | 9 | | | | | | | | | | | | |
| Paragraphs | 14 | | | | | | | | | | | | |
| Read Time | 2 minute(s) | | | | | | | | | | | | |
| Speak Time | 2 minute(s) | | | | | | | | | | | | |
| Date | 04/11/25 | | | | | | | | | | | | |