

GIO: Generating Efficient Matrix and Frame Readers for Custom Data Formats by Example^[1]

SAEED FATHOLLAHZADEH, Graz University of Technology & Know-Center GmbH, Austria
MATTHIAS BOEHM, Technische Universität Berlin, Germany

This guide provides instructions for replicating the paper’s findings. To ensure successful execution, it is essential to meet specific setup requirements.

1 SOURCE CODE INFO

- **Repository:** Apache SystemDS (<https://github.com/apache/systemds>, commit: `82d9d130861be8e36d37a08c22cdd8d3231de6c2`)
- **Reproducibility Repository:** <https://github.com/damslab/reproducibility/tree/master/sigmod2023-GIO-p454>
- **Programming Language:** Java, Clang++10, Python 3.8, SystemDS
- **Packages/Libraries Needed:** JDK 11, cmake, RapidJSON, Clang++, Python, Git, Maven, pdflatex, unzip, unrar, xz-utils

2 EXPERIMENT INDEX STRUCTURE

```
1  .
2  |— baselines
3      |— JavaBaselines
4      |— PythonPandas
5      +— RapidJSONCPP
6  |— datagen
7  |— explocal
8      |— exp1_micor_bench_identification
9      |— exp2_identification_1k_10k
10     |— exp3_early
11     |— exp4_micro_bench
12     |— exp5_systematic
13     |— exp6_end_to_end
14     +— plotting
15 |— load-had3.3-java11.sh
16 |— plots
17 |— run0LoadConfig.sh
18 |— run1SetupDependencies.sh
19 |— run2SetupBaseLines.sh
20 |— run3DownloadData.sh
21 |— run4GenerateData.sh
22 |— run5LocalExperiments.sh
23 |— run6PlotResults.sh
24 +— runAll.sh
```

3 INSTALL AND CONFIG

Hardware and Software Info: We ran all experiments on a server node an AMD EPYC 7302 CPU @ 3.0-3.3, GHz (16 physical/32 virtual cores) with 128 GB, two 480 GB SATA SSDs (system/home), and twelve 2 TB SATA HDDs. All reader experiments utilize a single SSD.

Setup and Experiments: The repository is pre-populated with the paper's experimental results ("results"), individual plots ("plots"), and SystemDS source code.

To acquire installation and config, follow these steps:

```
1 ./run0LoadConfig.sh
2 ./run1SetupDependencies.sh;
3 ./run2SetupBaselines.sh;
```

Note 1. "run0LoadConfig.sh" is crucial to execute before proceeding with other steps. It loads JVM memory settings and configures Hadoop settings. Please refrain from commenting out this line when running the steps

4 DATA PREPARATION

The datasets will be downloaded and generated automatically, follow these steps:

```
1 ./run3DownloadData.sh;
2 ./run4GenerateData.sh;
```

5 RUN EXPERIMENTS

To execute all experiments, simply run "run5LocalExperiments.sh". The experiments will commence, collecting output results.

```
1 ./run5LocalExperiments.sh;
```

Note 2. Our configuration entails running all experiments five times, making it a more time-consuming process.

6 PLOTTING

For plotting, simply run "run6PlotResults.sh":

```
1 ./run6PlotResults.sh;
```

Since the experiments are run five times, at the end of the experiments, we will merge all of them and create average results.

7 RUN ALL

The entire experimental evaluation can be run via "runAll.sh", which deletes the results and plots and performs setup, dataset download, dataset preparation, dataset generating, local experiments, and plotting. However, for a more controlled evaluation, we recommend running the individual steps separately.

```
1 ./run0LoadConfig.sh
2 ./run1SetupDependencies.sh;
3 ./run2SetupBaselines.sh;
4 ./run3DownloadData.sh;
5 ./run4GenerateData.sh;
6 ./run5LocalExperiments.sh;
7 ./run6PlotResults.sh;
```

Please be aware that this process may take a considerable amount of time, and you may prefer to have better control over the timing of each experiment.

Note 3. we are clearing the system cache using the command `"echo 3 >/proc/sys/vm/drop_caches && sync"`. Therefore, it is necessary to use `"sudo"` when executing `"runAll.sh"` (e.g., `"sudo runAll.sh"`).

REFERENCES

- [1] Saeed Fathollahzadeh and Matthias Boehm. 2023. GIO: Generating Efficient Matrix and Frame Readers for Custom Data Formats by Example. *Proc. ACM Manag. Data* 1, 2 (2023), 120:1–120:26. <https://doi.org/10.1145/3589265>