

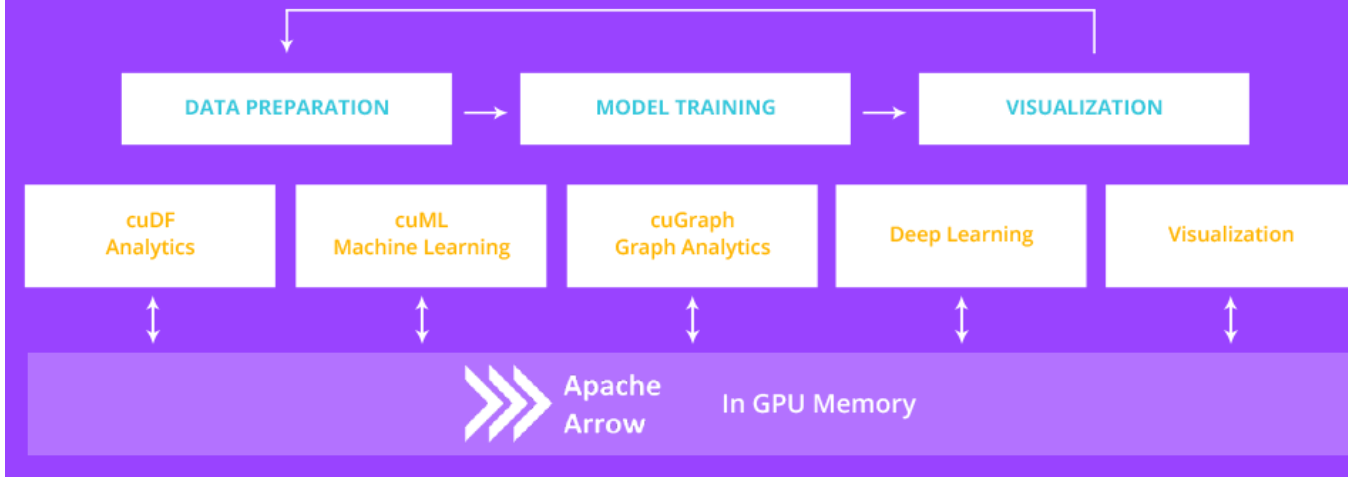
RAPIDS: end-to-end data science

RAPIDS in a nutshell:

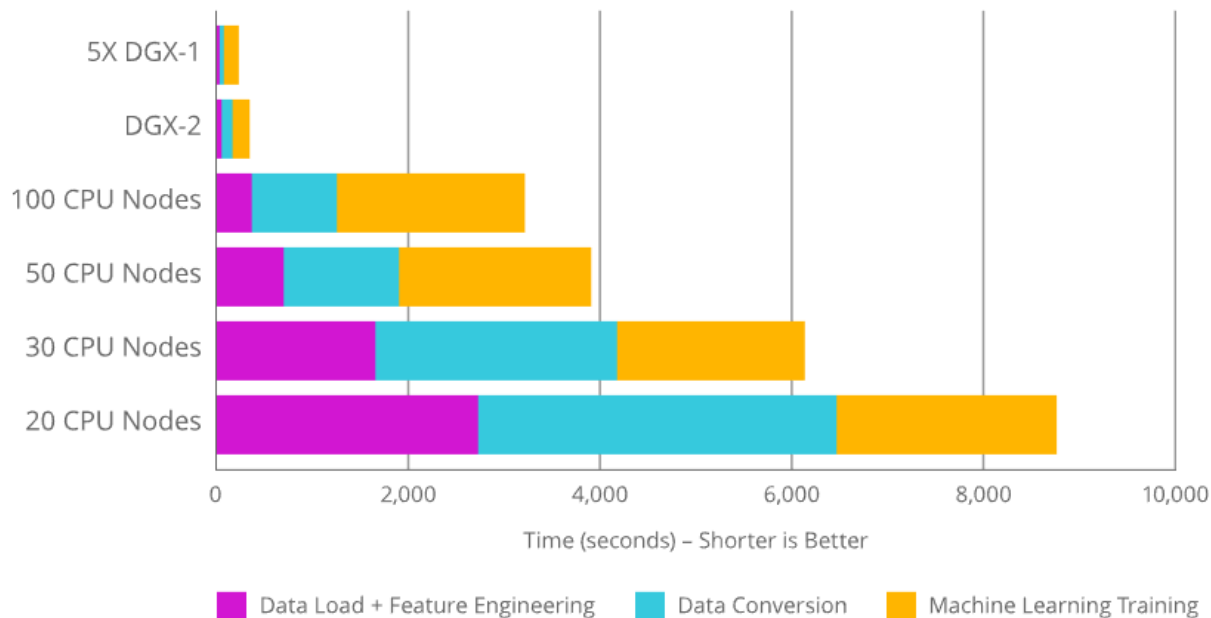
RAPIDS is a suite of open source software libraries aim to enable execution of end-to-end data science and analytics pipelines entirely on GPUs.

These libraries rely on NVIDIA® CUDA® primitives for low-level compute optimization, while concurrently exposing GPU parallelism and high-bandwidth memory speed through user-friendly Python interfaces.

The New GPU Data Science Pipeline



End-to-End Faster Speeds on RAPIDS



<https://rapids.ai/index.html>

KEY FEATURES:

cuDF:

DataFrame
manipulation for data
preparation

cuDF currently has 2
APIs: C++ and a Python
API that mimics Pandas

cuML:

GPU-accelerated
machine learning
libraries.

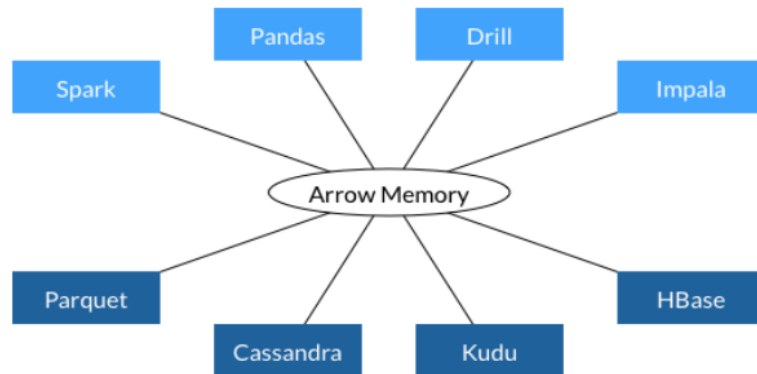
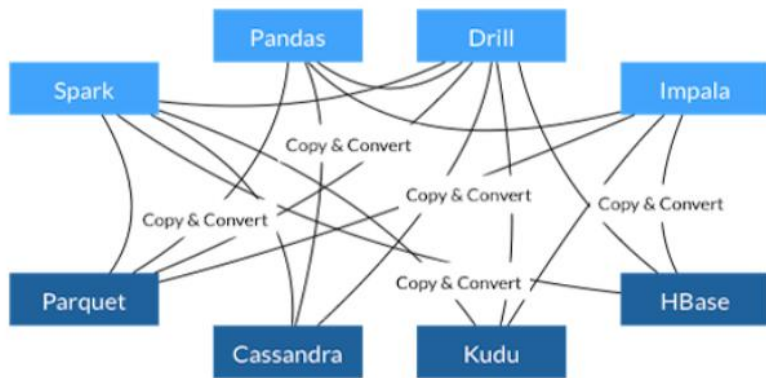
Apache Arrow:

Columnar in-memory
data structure that
delivers efficient and
fast data interchange
with flexibility to
support complex data
models.

Apache Arrow

Apache Arrow is a cross-language development platform for in-memory data.

Standard language-independent columnar memory format for flat and hierarchical data, organized for efficient analytic operations on modern hardware.



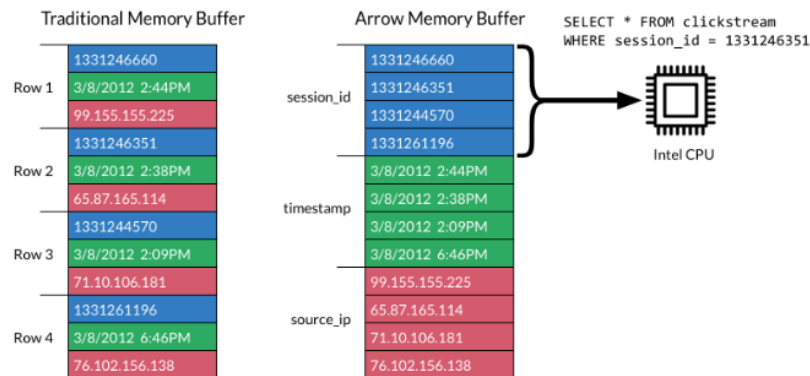
Apache Arrow is an open standard for DataFrame-like data manipulation

Main draw: **zero-copy** reads between interfaces

Additional advantages:

- Execution engines can take advantage of SIMD while processing DataFrames
- Cache misses minimized
- Leverages columnar compression techniques
- Allows for the deconstruction of traditionally vertically integrated analytic database architectures

	session_id	timestamp	source_ip
Row 1	1331246660	3/8/2012 2:44PM	99.155.155.225
Row 2	1331246351	3/8/2012 2:38PM	65.87.165.114
Row 3	1331244570	3/8/2012 2:09PM	71.10.106.181
Row 4	1331261196	3/8/2012 6:46PM	76.102.156.138



Example 1:

```
import cudf
gdf = cudf.read_csv('path/to/file.csv')
for column in gdf.columns:
    print(gdf[column].mean())
```

Example 2: Kaggle competition

