# Communication-Efficient Cluster Scalable Genomics Data Processing Using Apache Arrow Flight

**Tanveer Ahmad**

**TU**Delft

- Background
- Technologies introduction
- Implementation
- Results
  - Performance evaluation
  - Comparison with MPI and Apache Spark
- Conclusion

**T**U Delft

# Background

- High throughput sequencing (HTS) technologies

  - Short reads (Illumina), Long reads (ONT, PacBio)

- Single-node to clouds and HPC systems

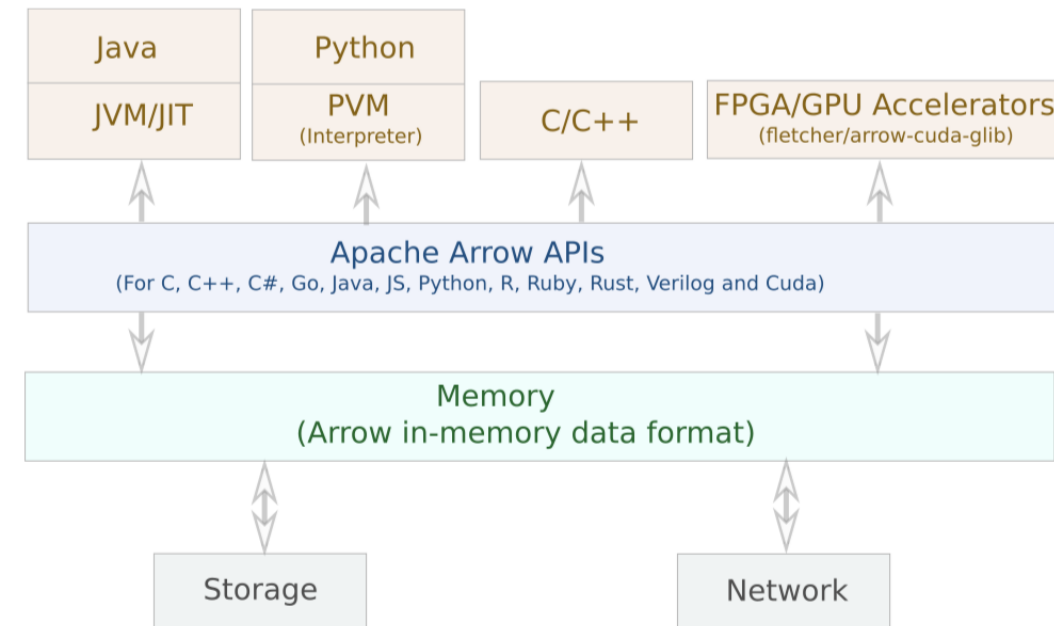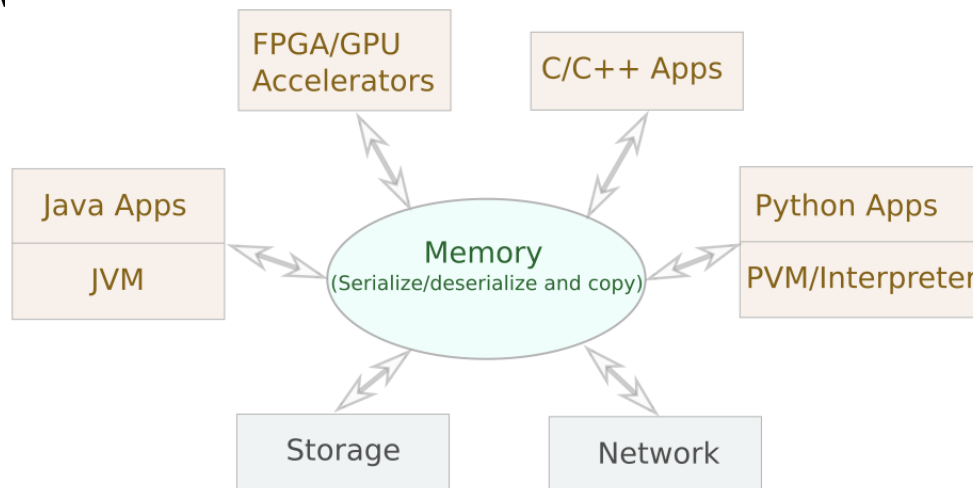- Genomics computational workflows


- Big data frameworks-based solutions

  - Halvade

  - ADAM and SparkGA2

Serialization,
Memory overhead
and
Scalability issues

# Technologies introduction
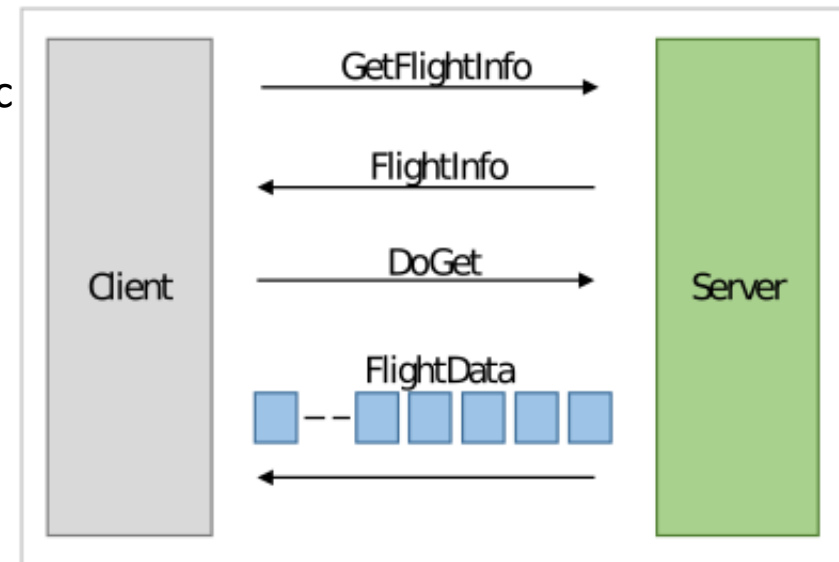
- ## Apache Arrow

  - Apache Arrow is an in-memory standard columnar data format

  - Columnar data storage enables efficient vectorized operations

  - Better cache locality can be exploited using this format

  - Arrow provides cross-language interoperability
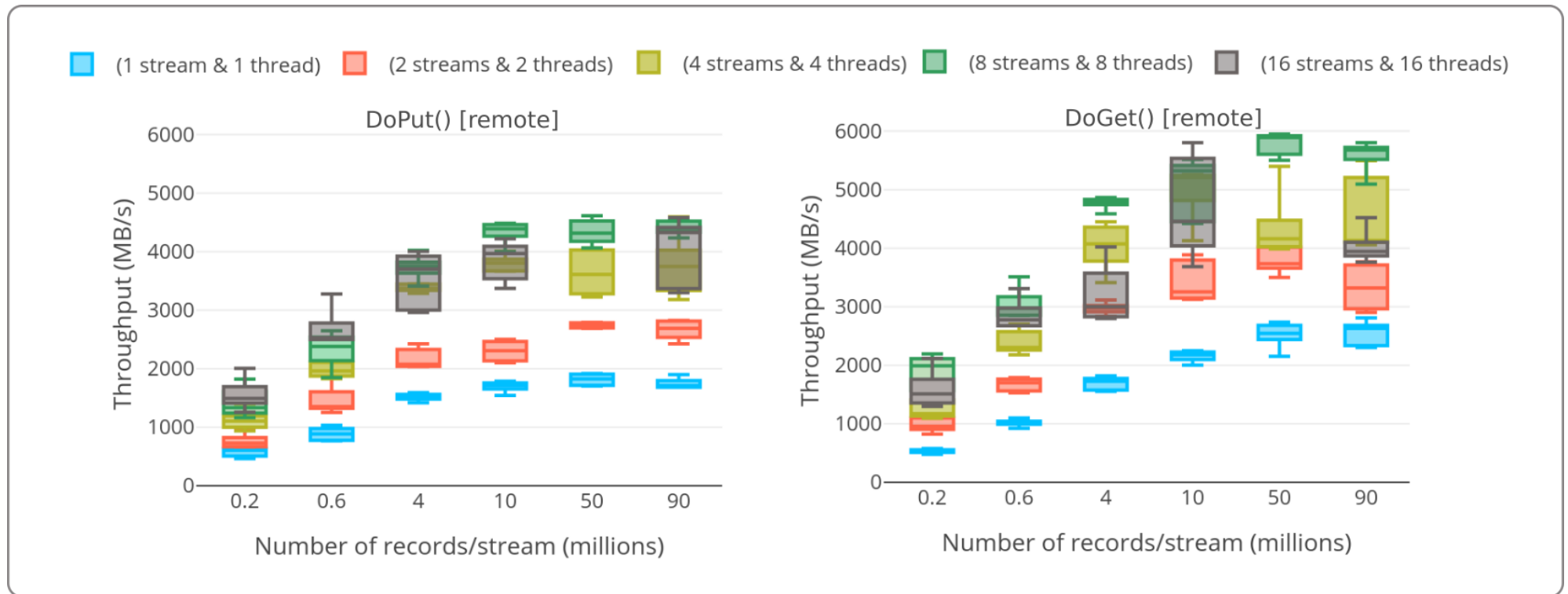
  - GPU/FPGA

# Technologies introduction

- ## Apache Arrow Flight

  - Arrow Flight is a submodule in the Apache Arrow project

  - Arrow Flight provides a high performance, secure, parallel and c
    platform language support

  - Apache Spark integration
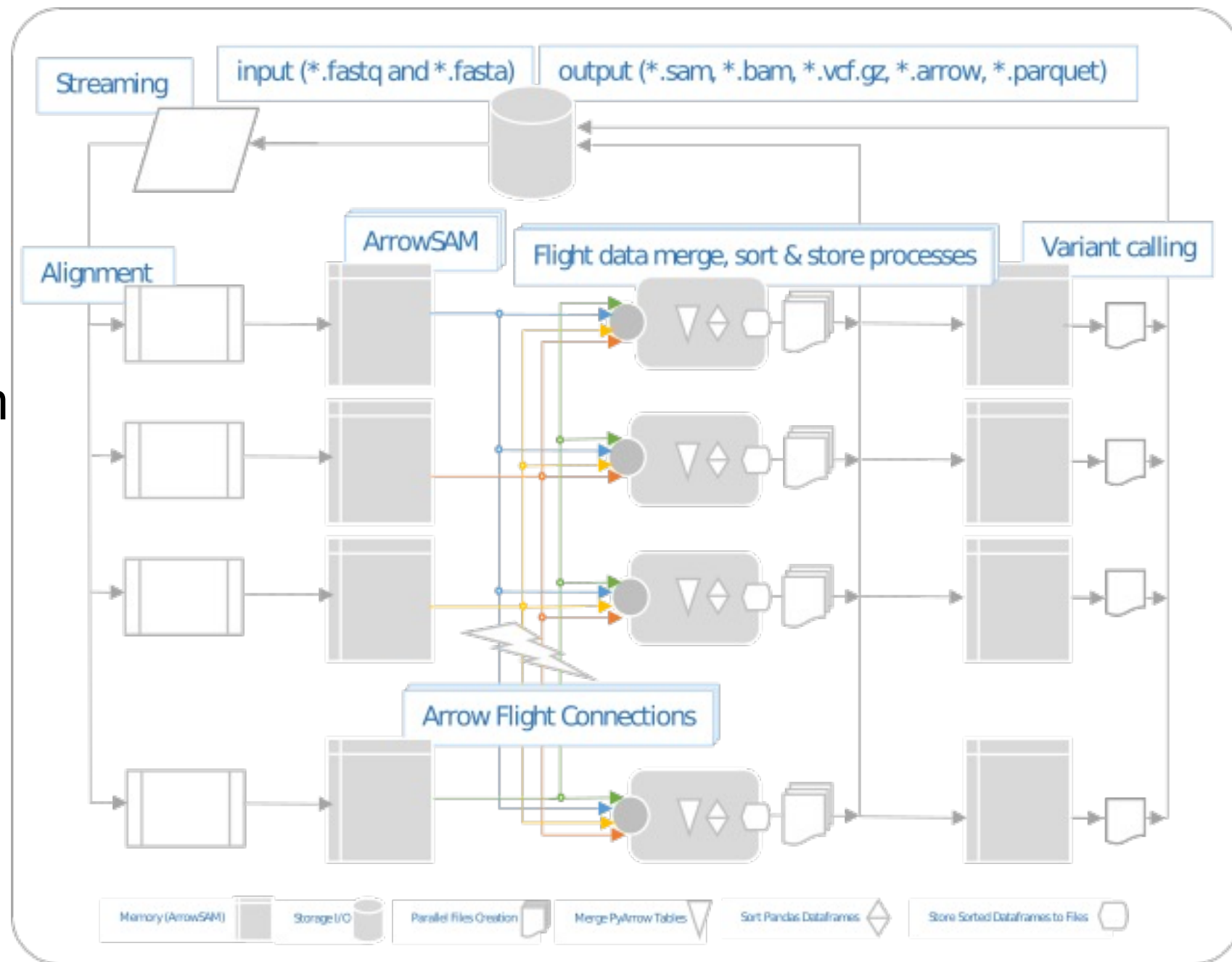
  - Arrow Flight client-server setup

# Technologies introduction

- ## Apache Arrow Flight
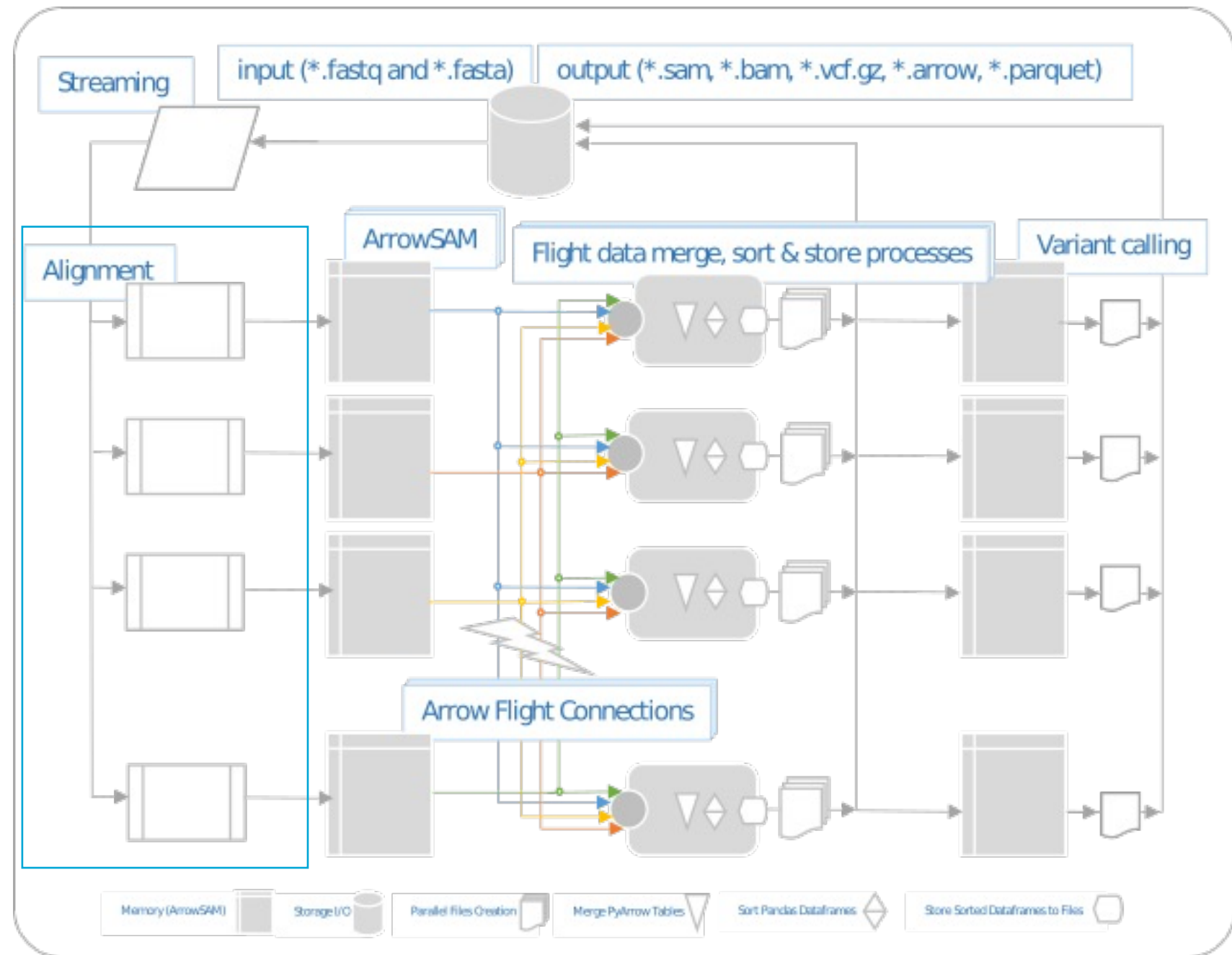
  - Remote client-server throughput

# Implementation

- Alignment
- ArrowSAM
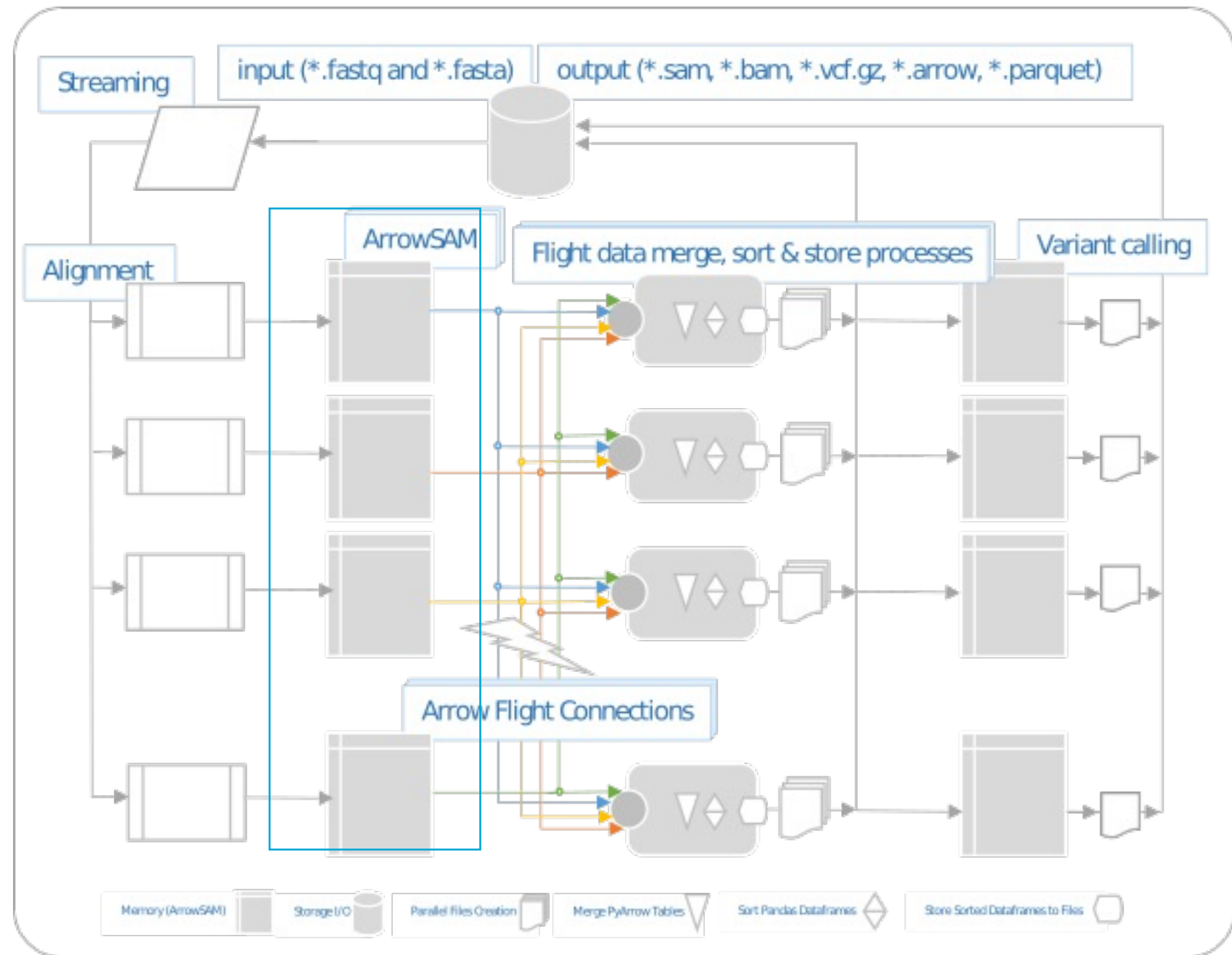- Arrow Flight comm
- Pre-processing
- Variant calling

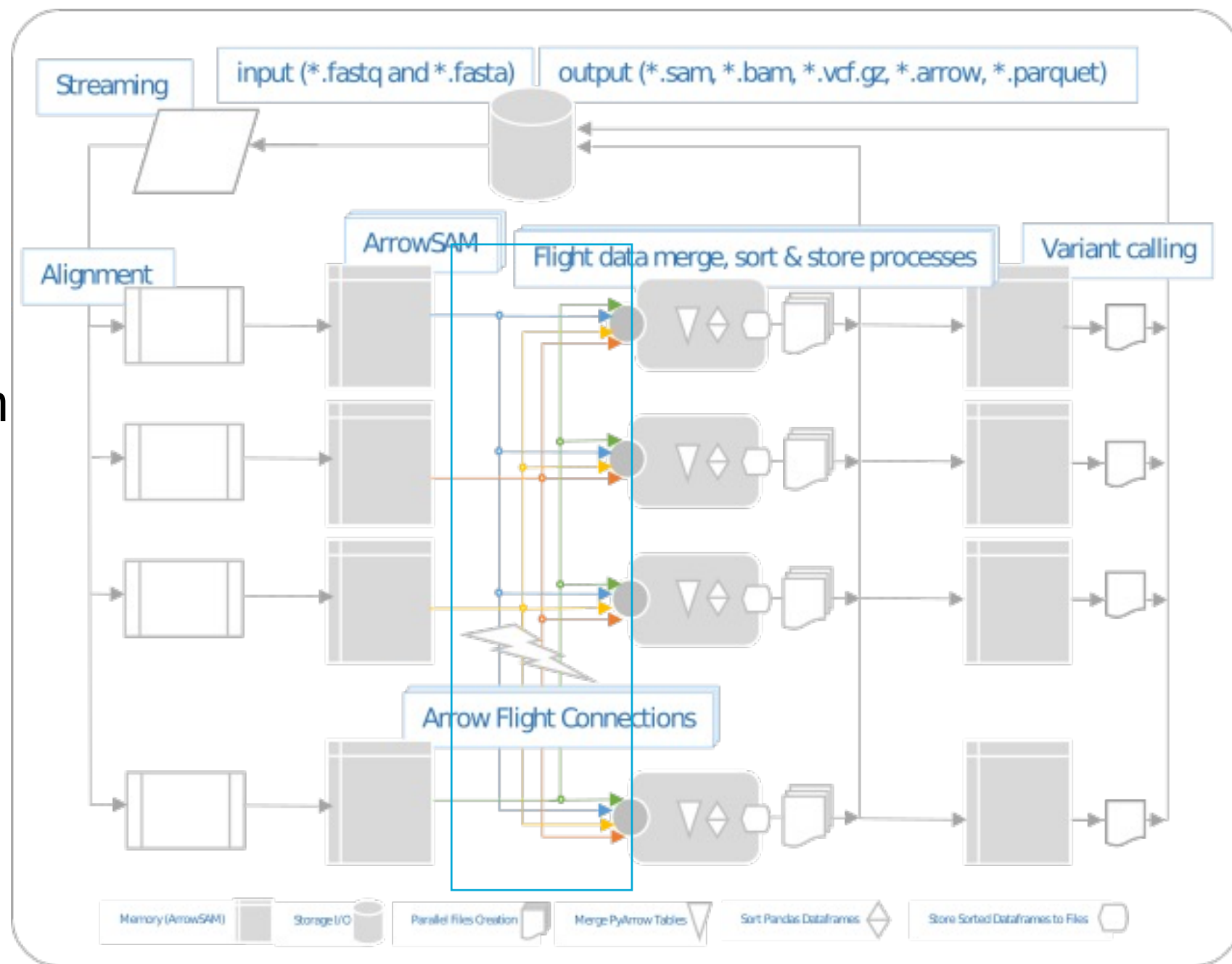# Implementation

- Alignment
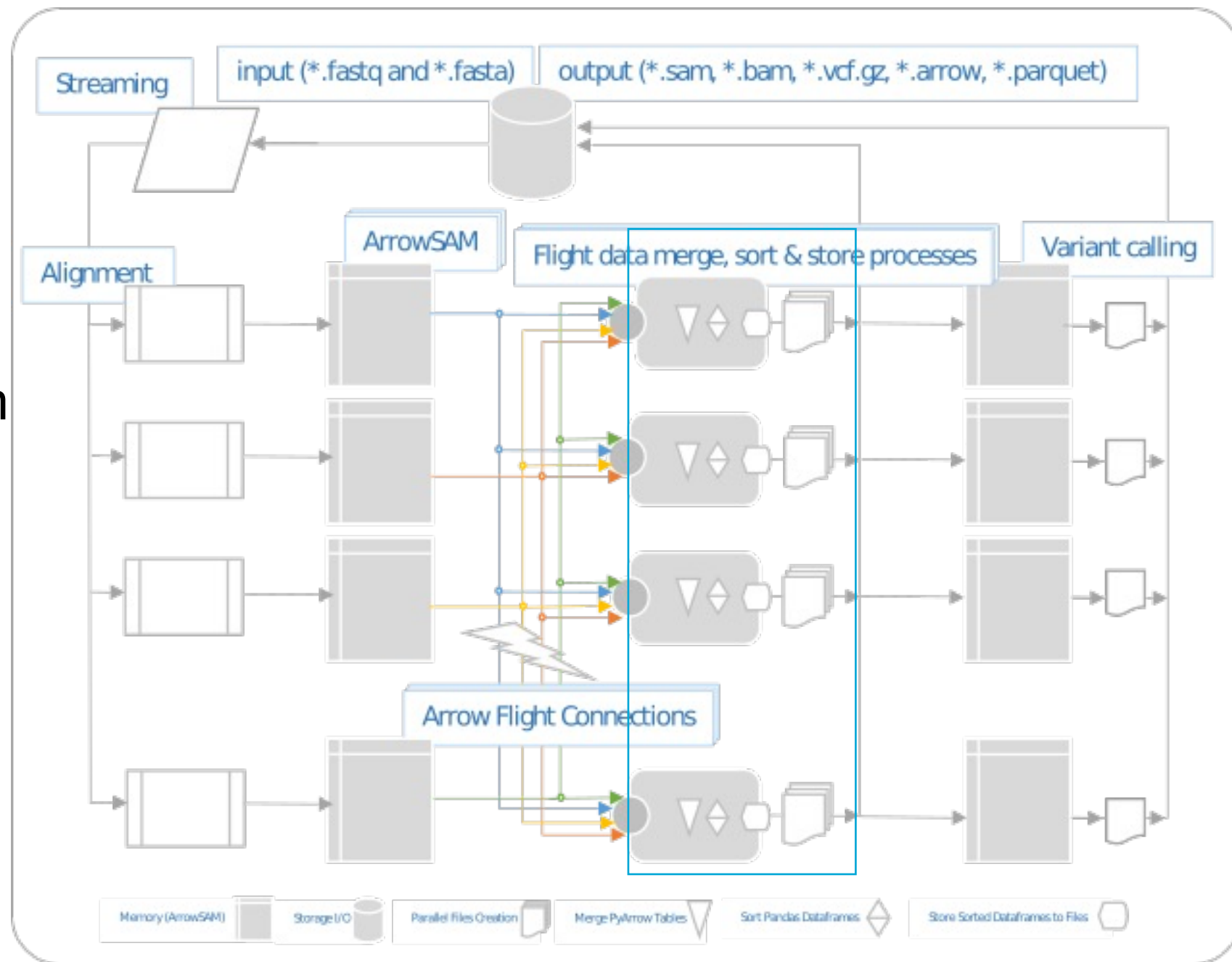
# Implementation

- Alignment
- ArrowSAM

# Implementation

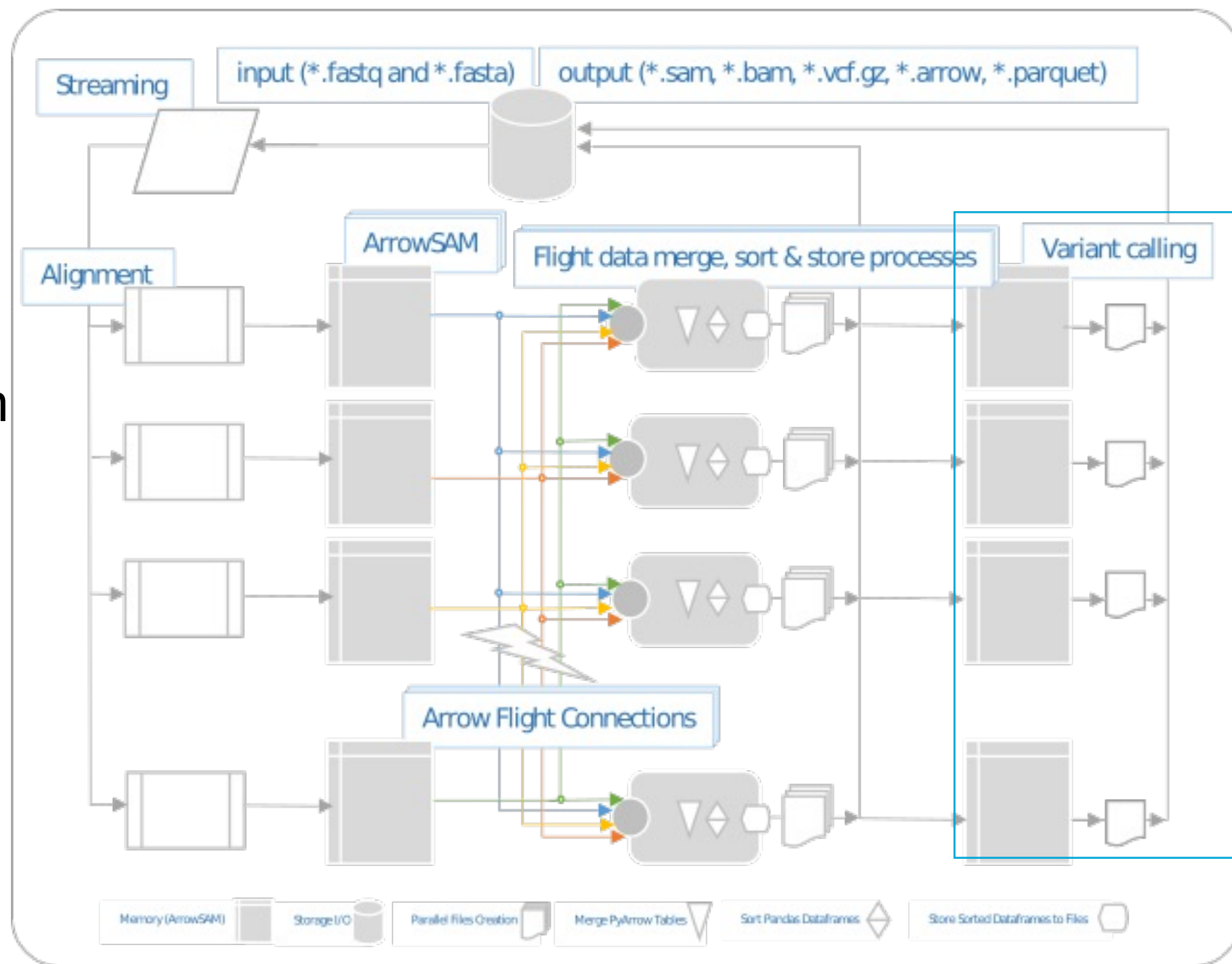- Alignment
- ArrowSAM
- Arrow Flight comm

# Implementation

- Alignment
- ArrowSAM
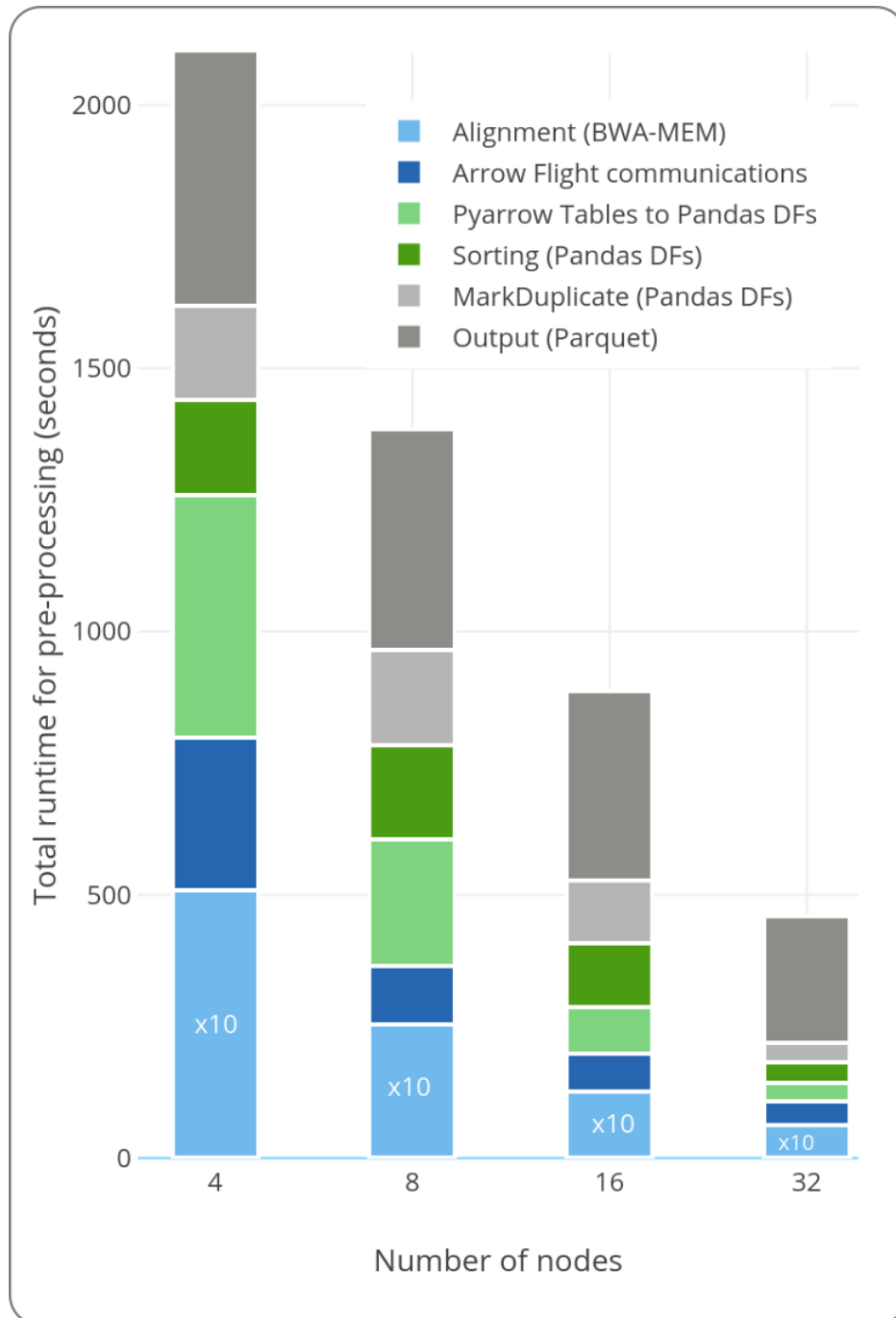- Arrow Flight comm
- Pre-processing

# Implementation

- Alignment
- ArrowSAM
- Arrow Flight comm
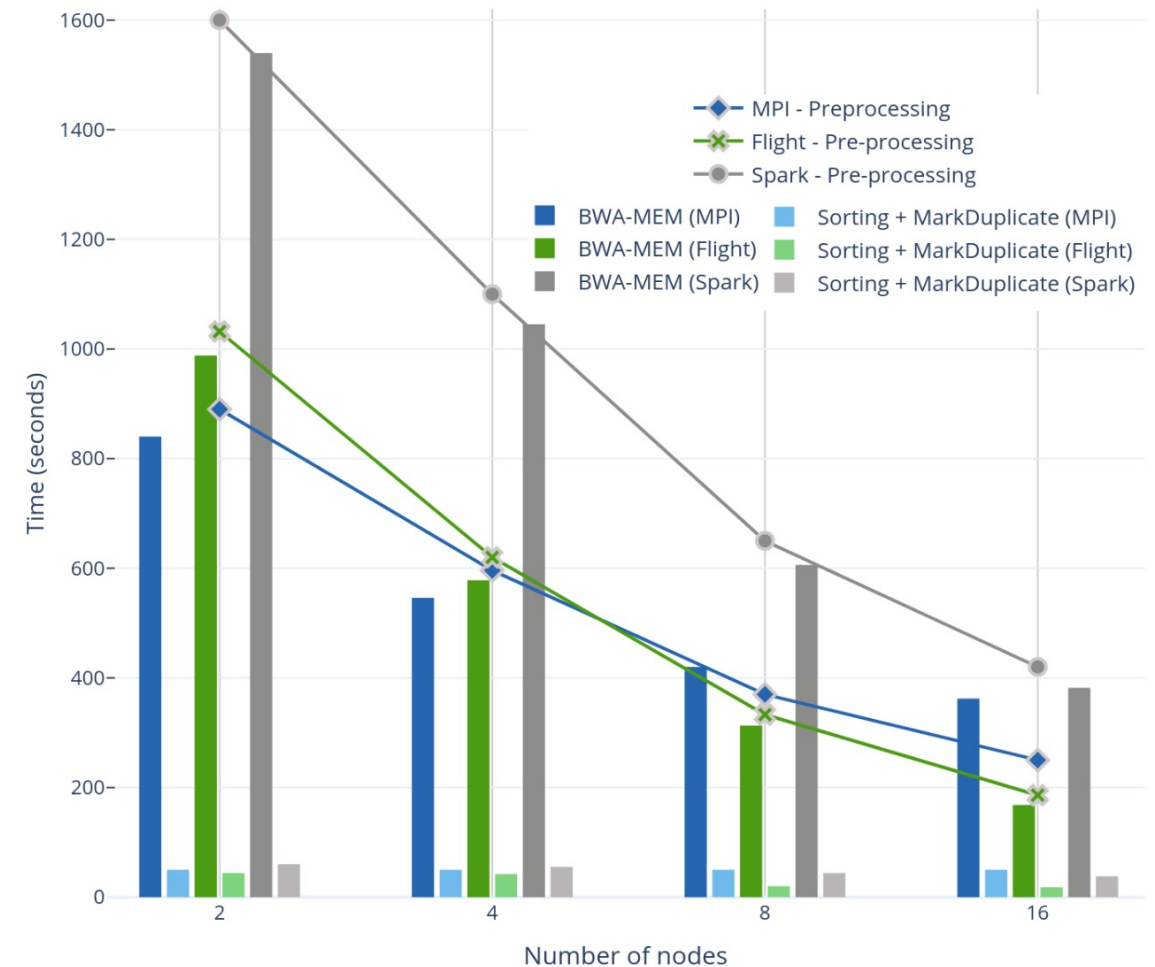- Pre-processing
- Variant calling

# Results

- **Performance evaluation**

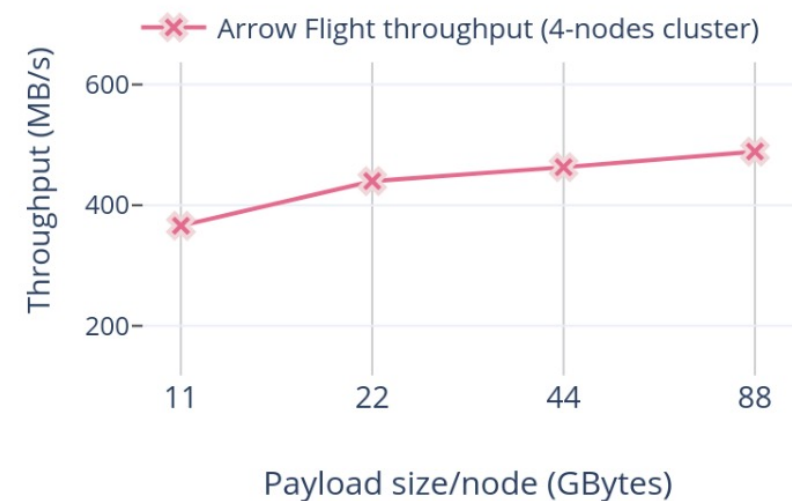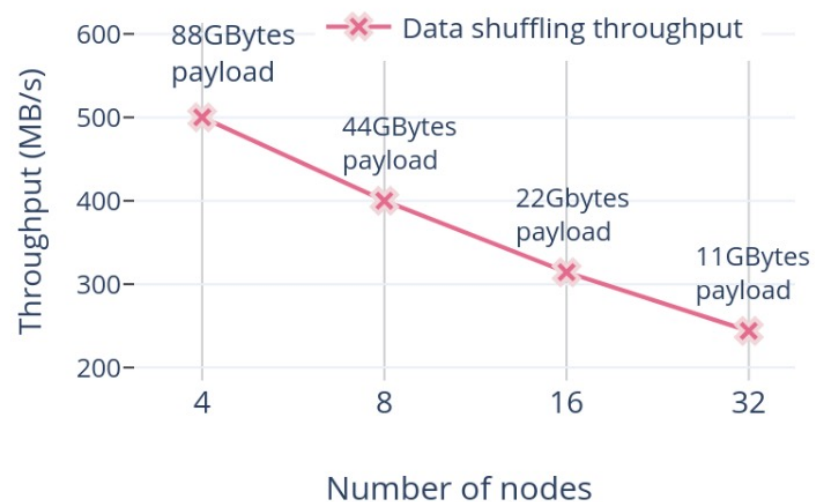- Runtime speedup

- Cluster scalability

# Results

- Comparison with MPI and Apache Spark

# Conclusion

- Arrow Flight Throughput

# Conclusion

- Accuracy

### TABLE I

Accuracy evaluation of small variants of HG002 (NA24385 with 50x coverage taken from PrecisionFDA challenge V2 datasets) against GIAB HG002 v4.2 benchmarking set. This table shows the SNP and INDEL results for "Chr1" on a single node (default) run.

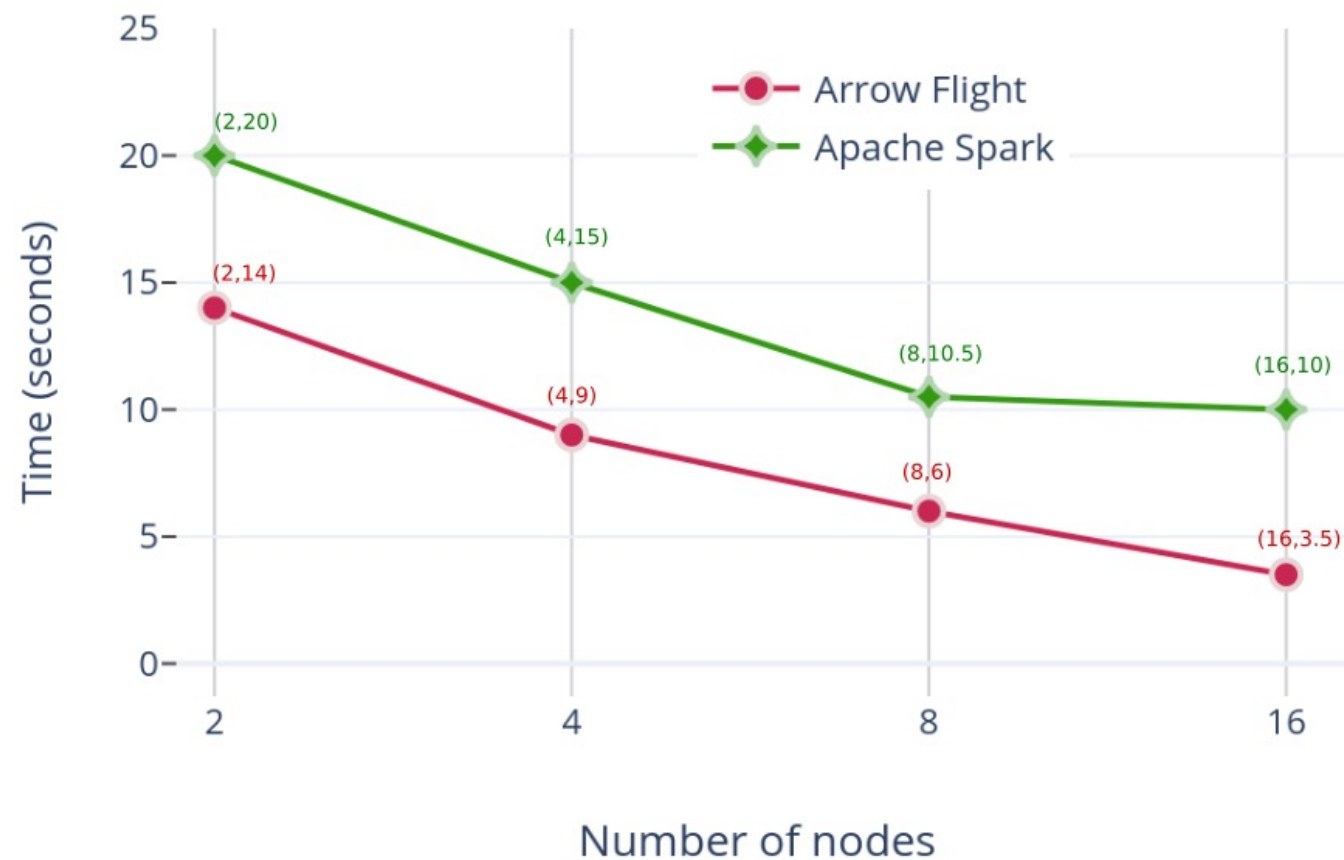| Variant type | Truth total | True positives | False negatives | False positives | Recall | Precision | F1-Score |
|---|---|---|---|---|---|---|---|
| INDEL | 42689 | 42390 | 299 | 131 | **0.992996** | 0.997053 | 0.995020 |
| SNP | 264143 | 262367 | 1776 | 351 | **0.993276** | **0.998665** | **0.995963** |

### TABLE II

Accuracy evaluation of small variants of HG002 (NA24385 with 50x coverage taken from PrecisionFDA challenge V2 datasets) against GIAB HG002 v4.2 benchmarking set. This table shows the SNP and INDEL results for "Chr1" on a cluster scaled (distributed) implementation. "Chr1" has been chunked into ten parts.

| Variant type | Truth total | True positives | False negatives | False positives | Recall | Precision | F1-Score |
|---|---|---|---|---|---|---|---|
| INDEL | 42689 | 42390 | 299 | 127 | **0.992996** | **0.997142** | **0.995065** |
| SNP | 264143 | 262365 | 1778 | 355 | 0.993269 | 0.998649 | 0.995952 |

TUDelft

# Conclusion

- Cluster scalability

Thank you for your attention

TU Delft