

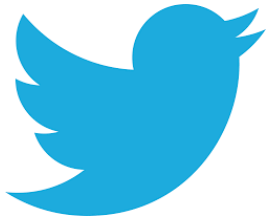
# High-throughput conversion of Apache Parquet files to Apache Arrow in-memory format using FPGAs

Lars van Leeuwen (l.t.j.vanleeuwen@student.tudelft.nl)

Johan Peltenburg, Jian Fang, Zaid Al-Ars, Peter Hofstee

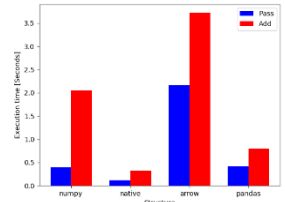
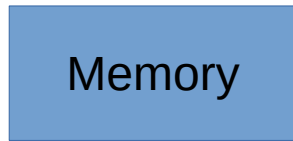
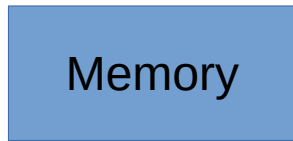


# Analyzing some tweets

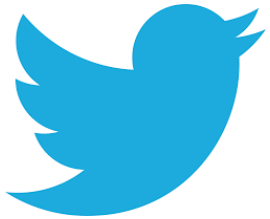


Decompression  
Deserialization

Compute

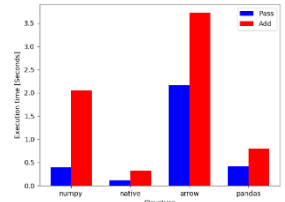
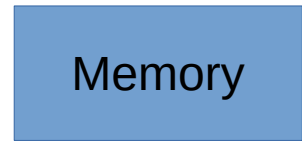
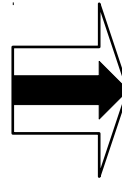
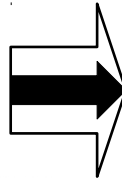
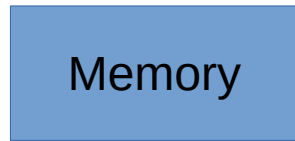


# Analyzing some tweets

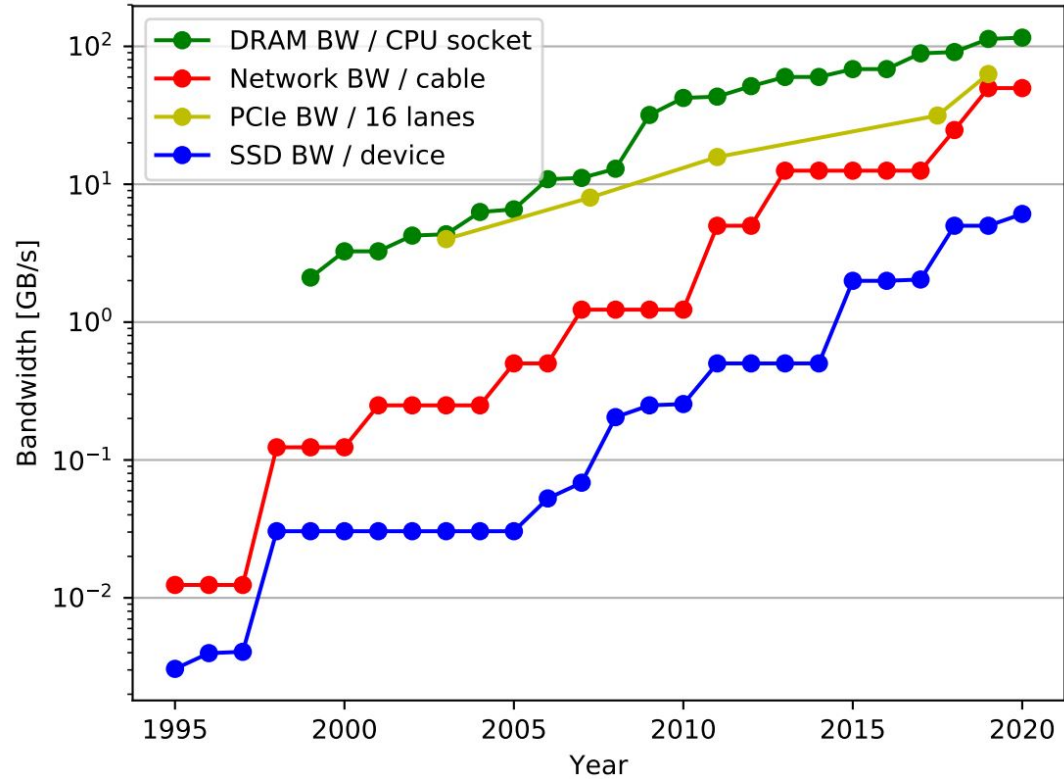


Decompression  
Deserialization

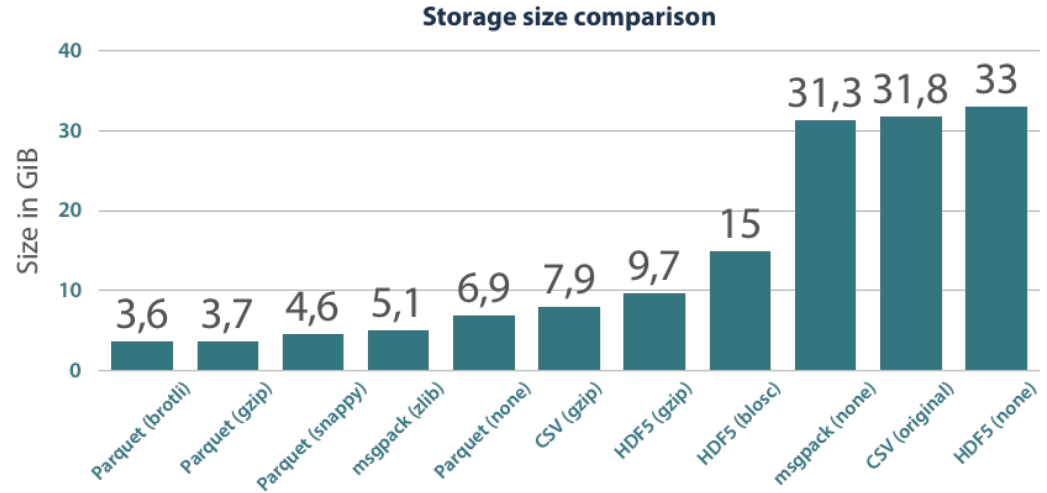
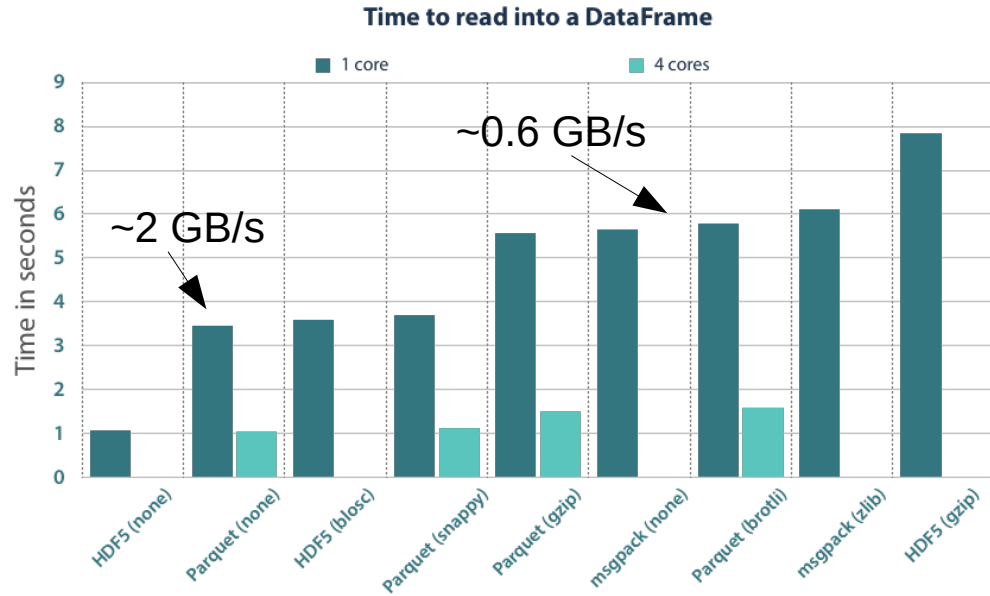
Compute



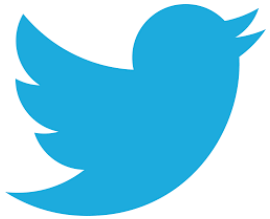
# Trends in bandwidth



# Reading data into Pandas

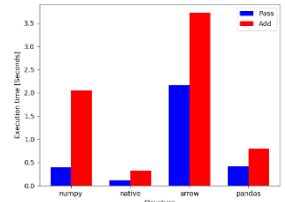


# Analyzing some tweets

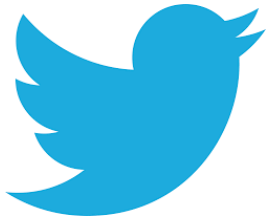


Decompression  
Deserialization

Compute

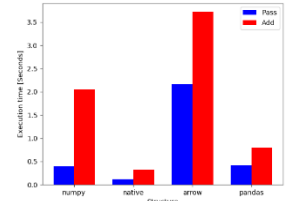
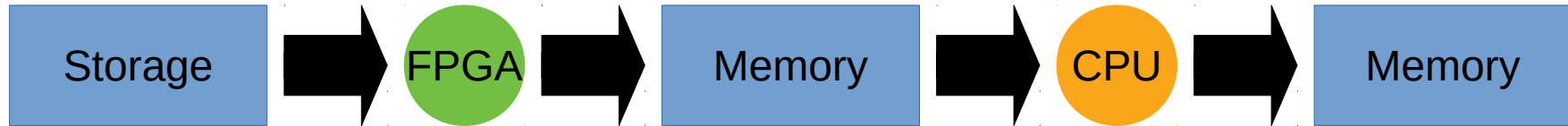


# Analyzing some tweets

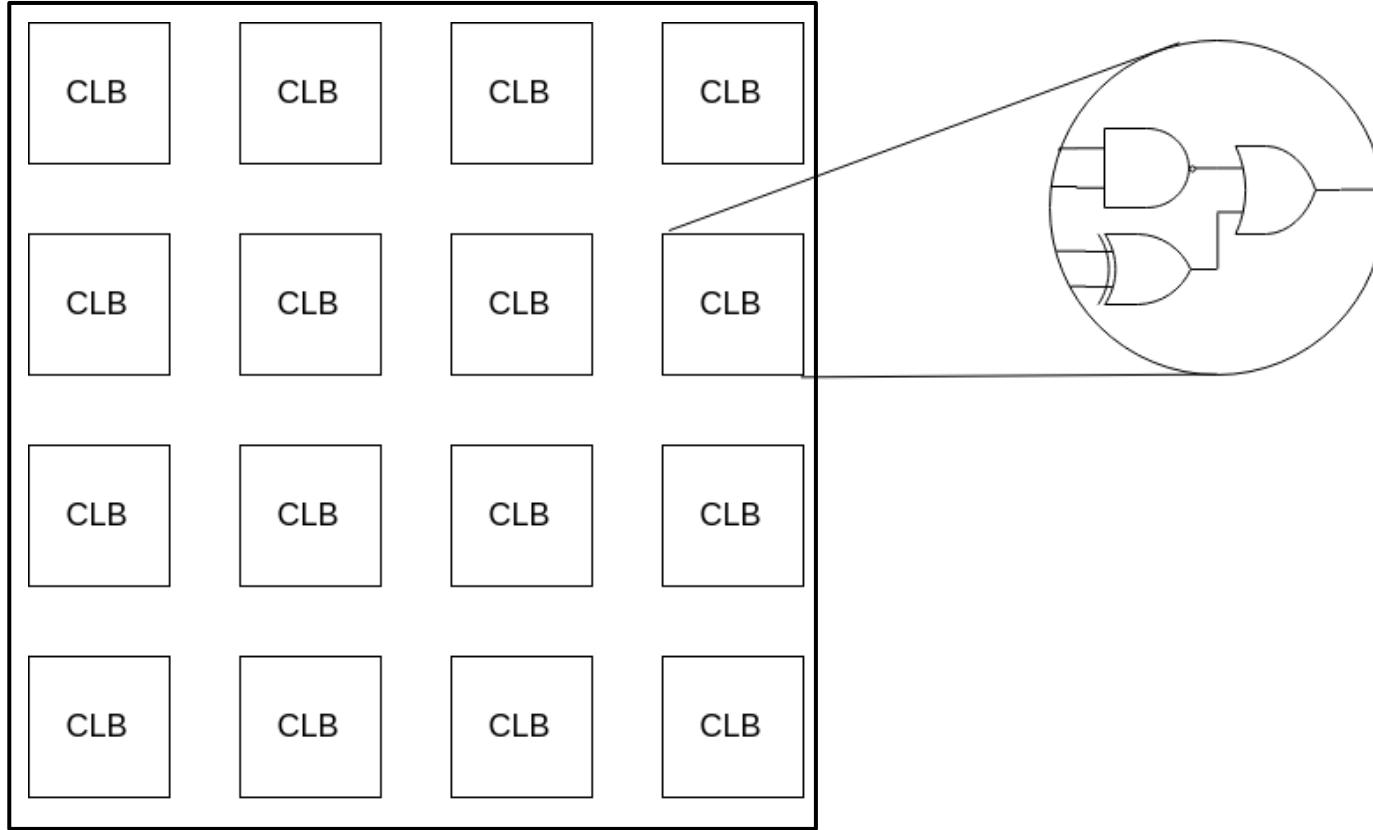


Decompression  
Deserialization

Compute

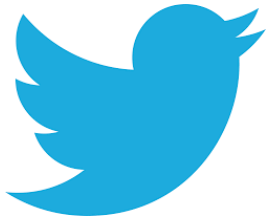


# Field Programmable Gate Array?



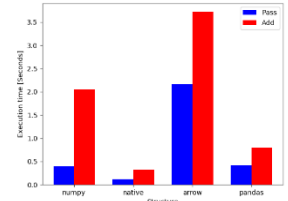
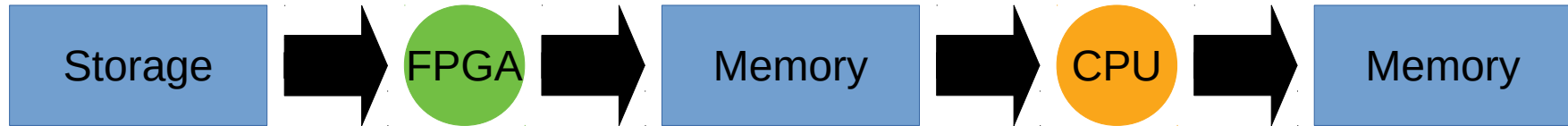


# Analyzing some tweets

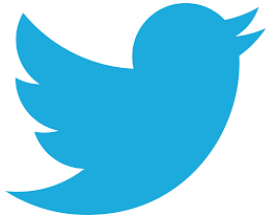


Decompression  
Deserialization

Compute

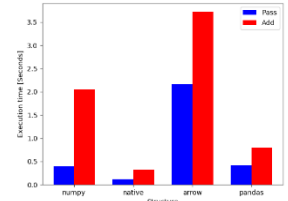
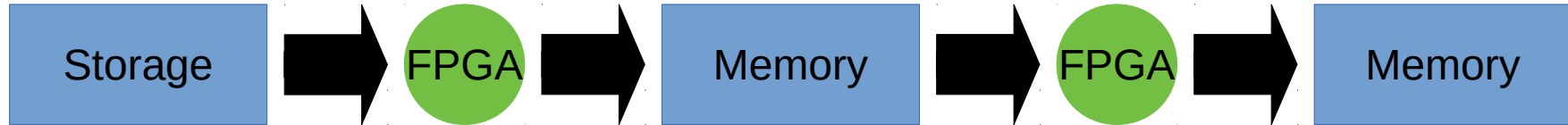


# Analyzing some tweets

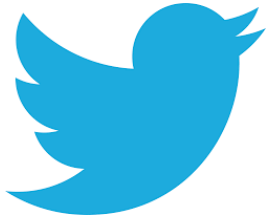


Decompression  
Deserialization

Compute



# Analyzing some tweets



Storage

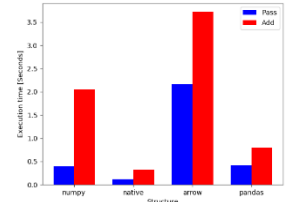


Decompression  
Deserialization  
Compute

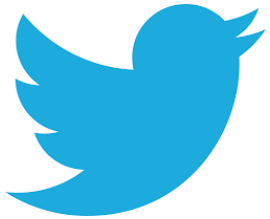
FPGA



Memory

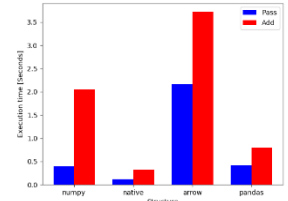
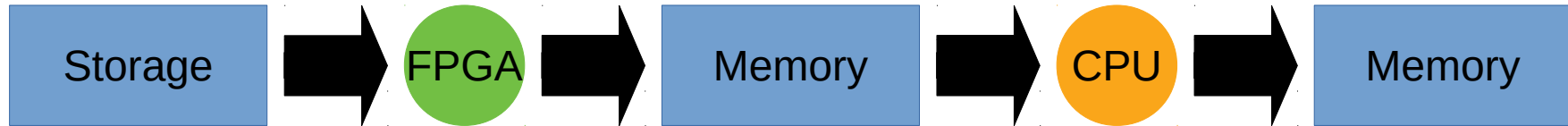


# Scoping the project



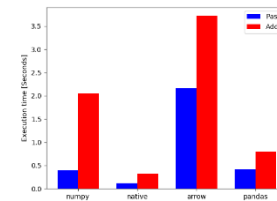
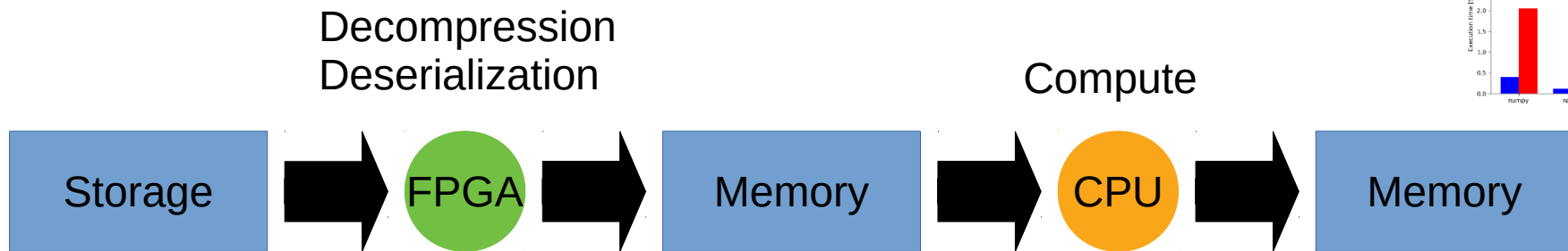
Decompression  
Deserialization

Compute

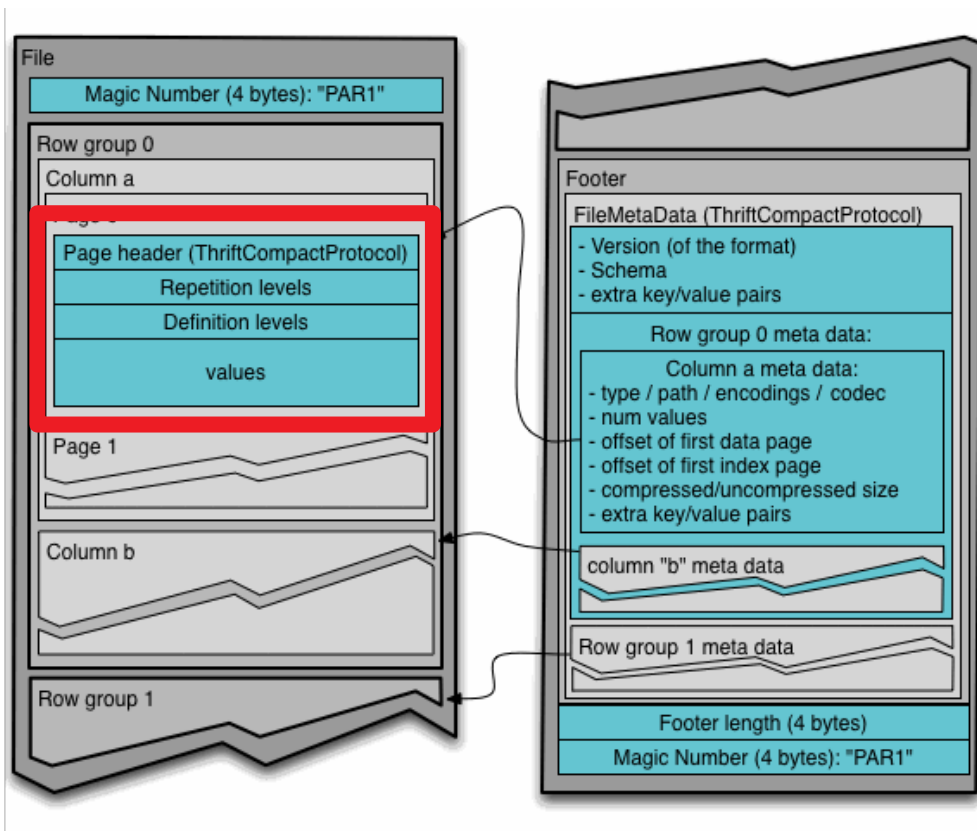


# Scoping the project

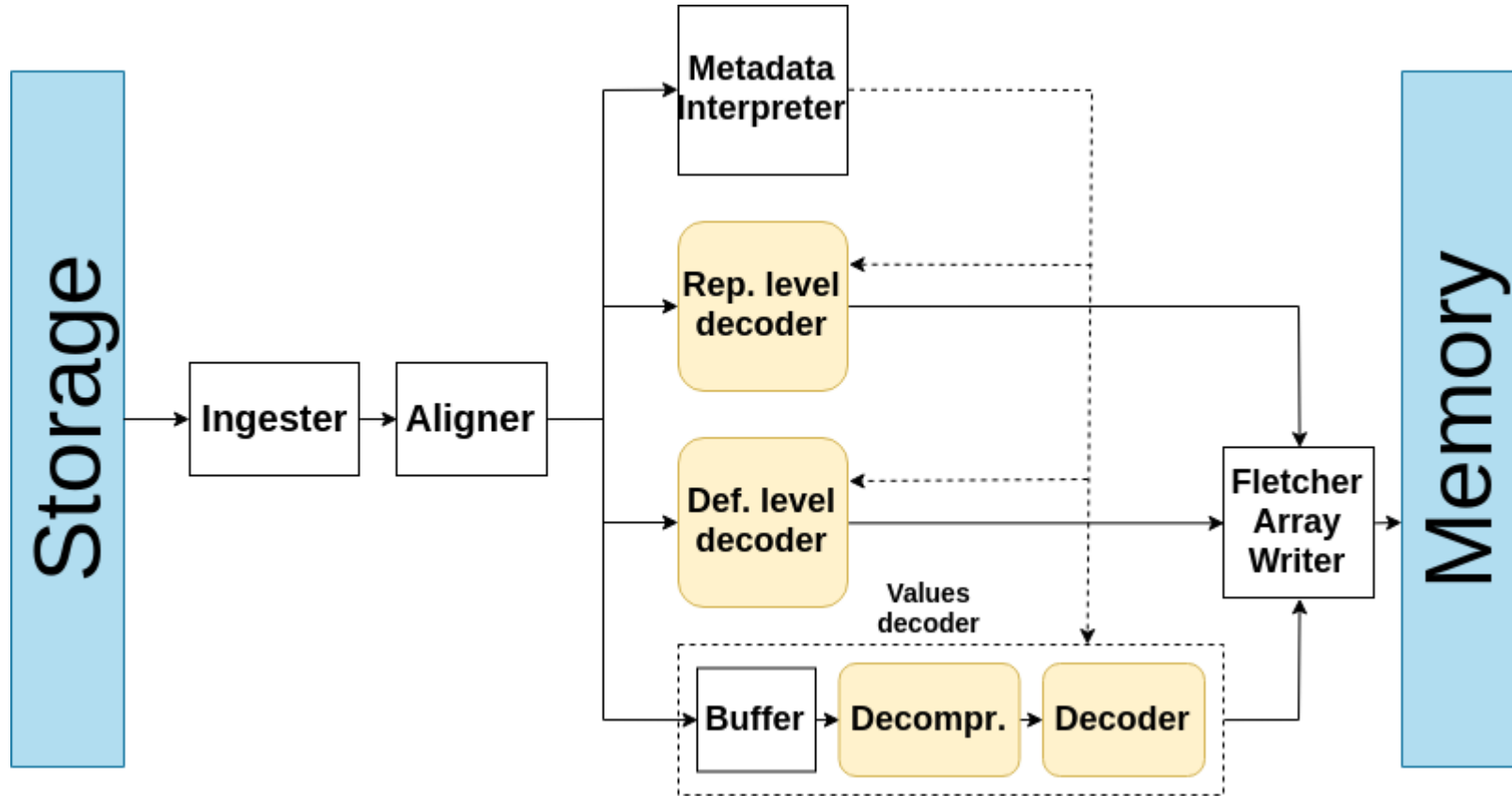
- 1) Both Parquet and Arrow are columnar formats
- 2) Parquet files can be split into parts
- 3) Arrow data is accessible to multiple language runtimes



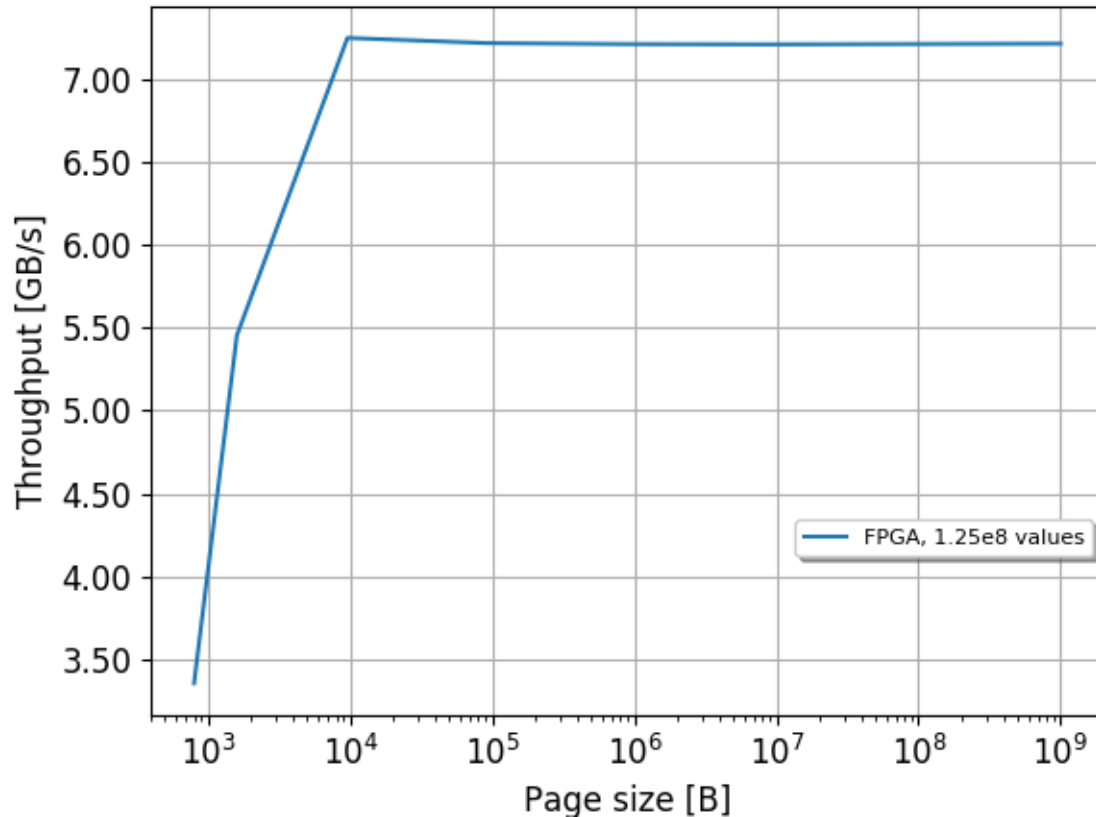
# Apache Parquet



# High-level architecture



# Performance for plain encoding (uncompressed)

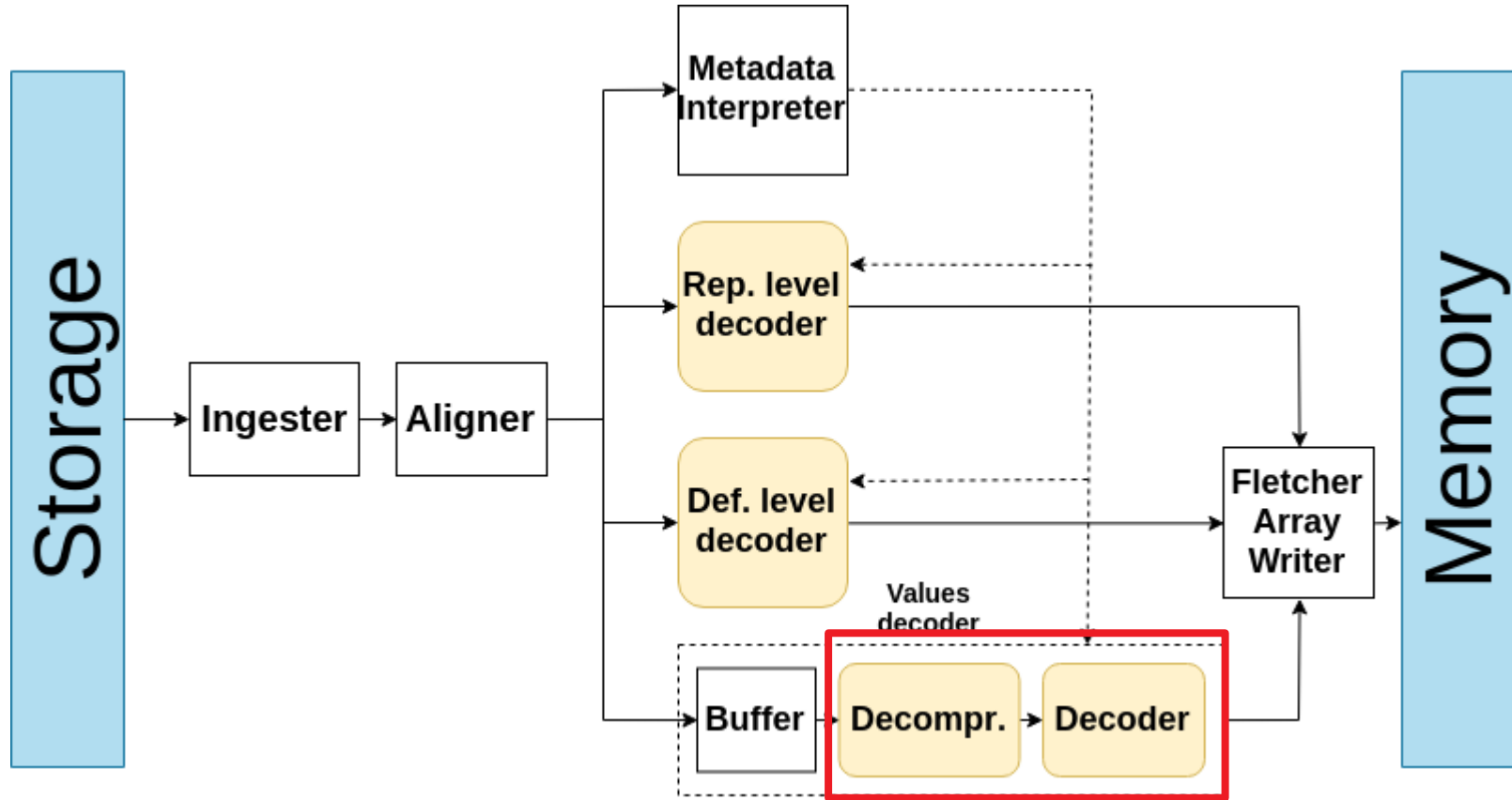


Single ParquetReader  
Implemented on XCVU9P

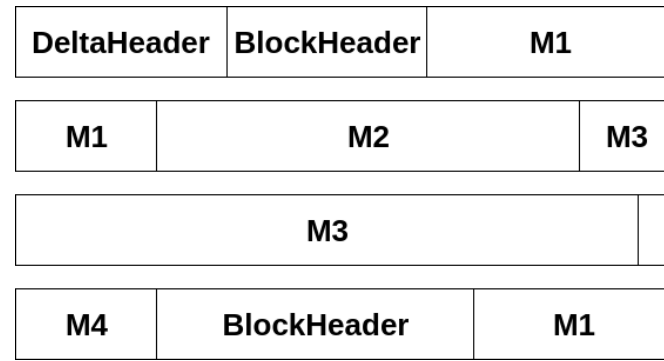
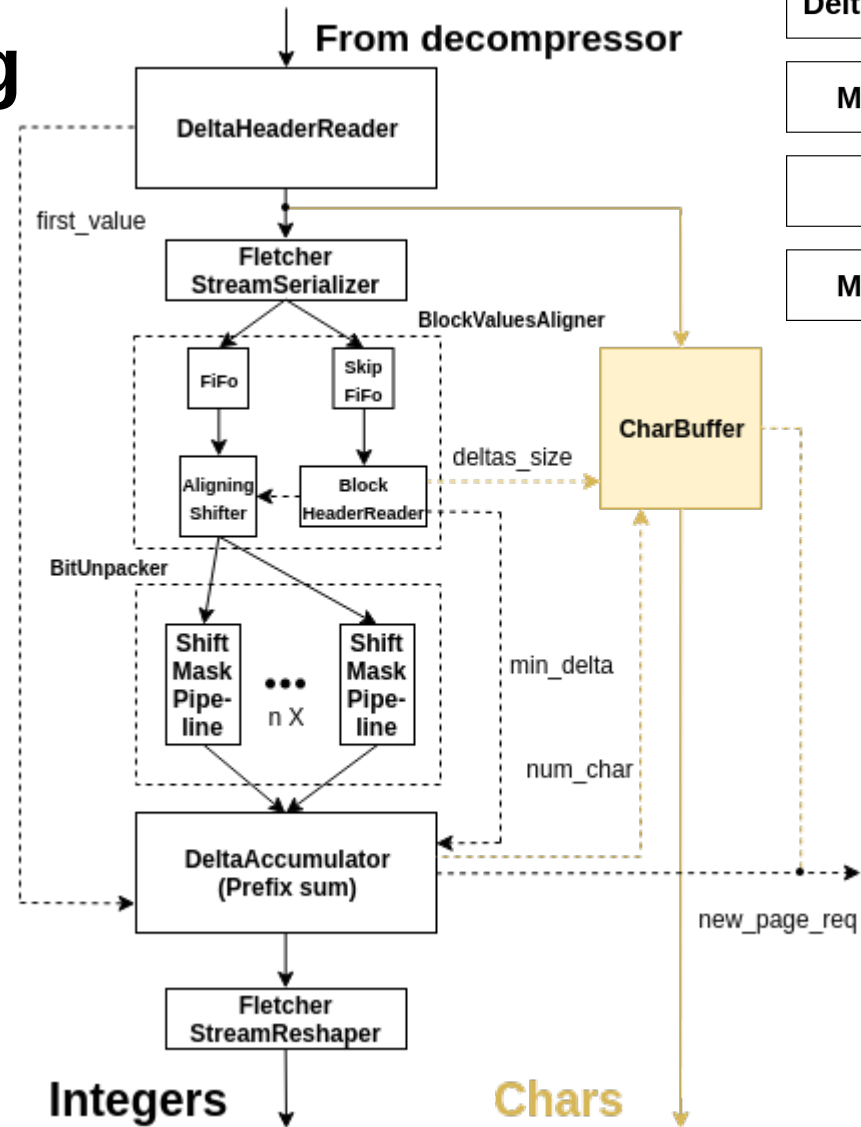
- Area utilization
  - CLB 3.14%
  - BRAM 1.78%
- Estimated power usage
  - ~1 W



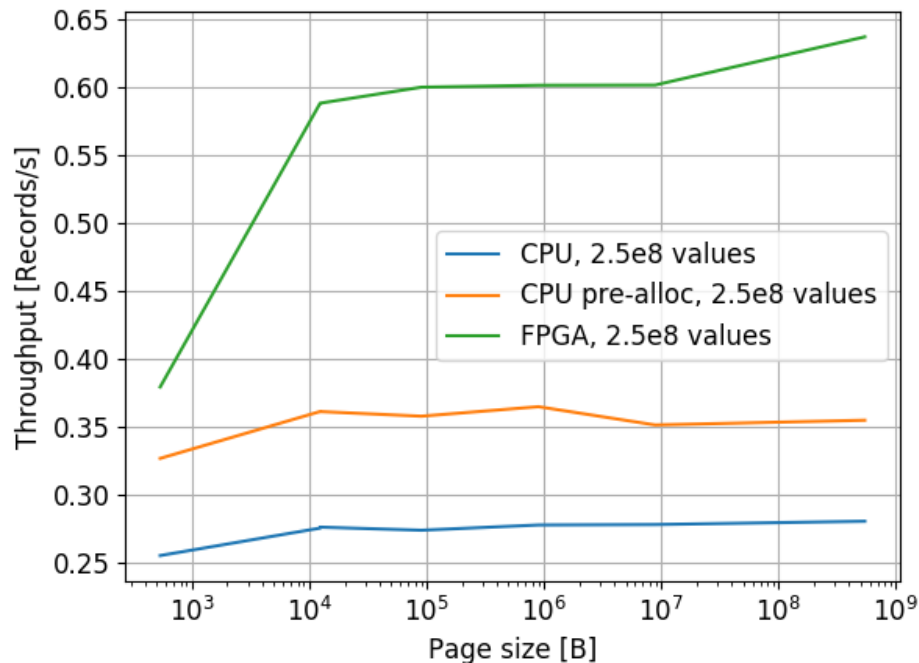
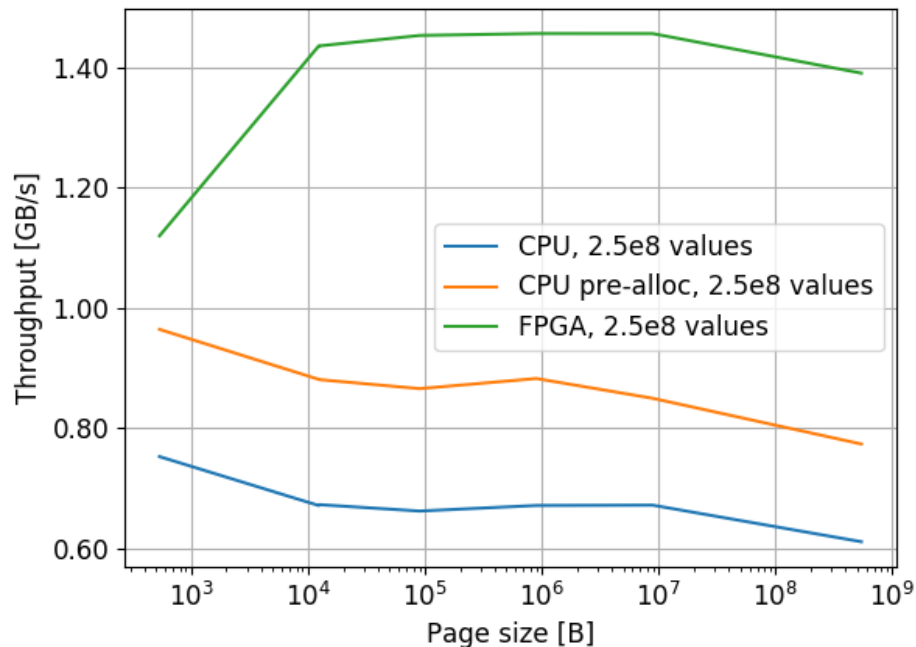
# Giving the FPGA something to do



# Delta decoding



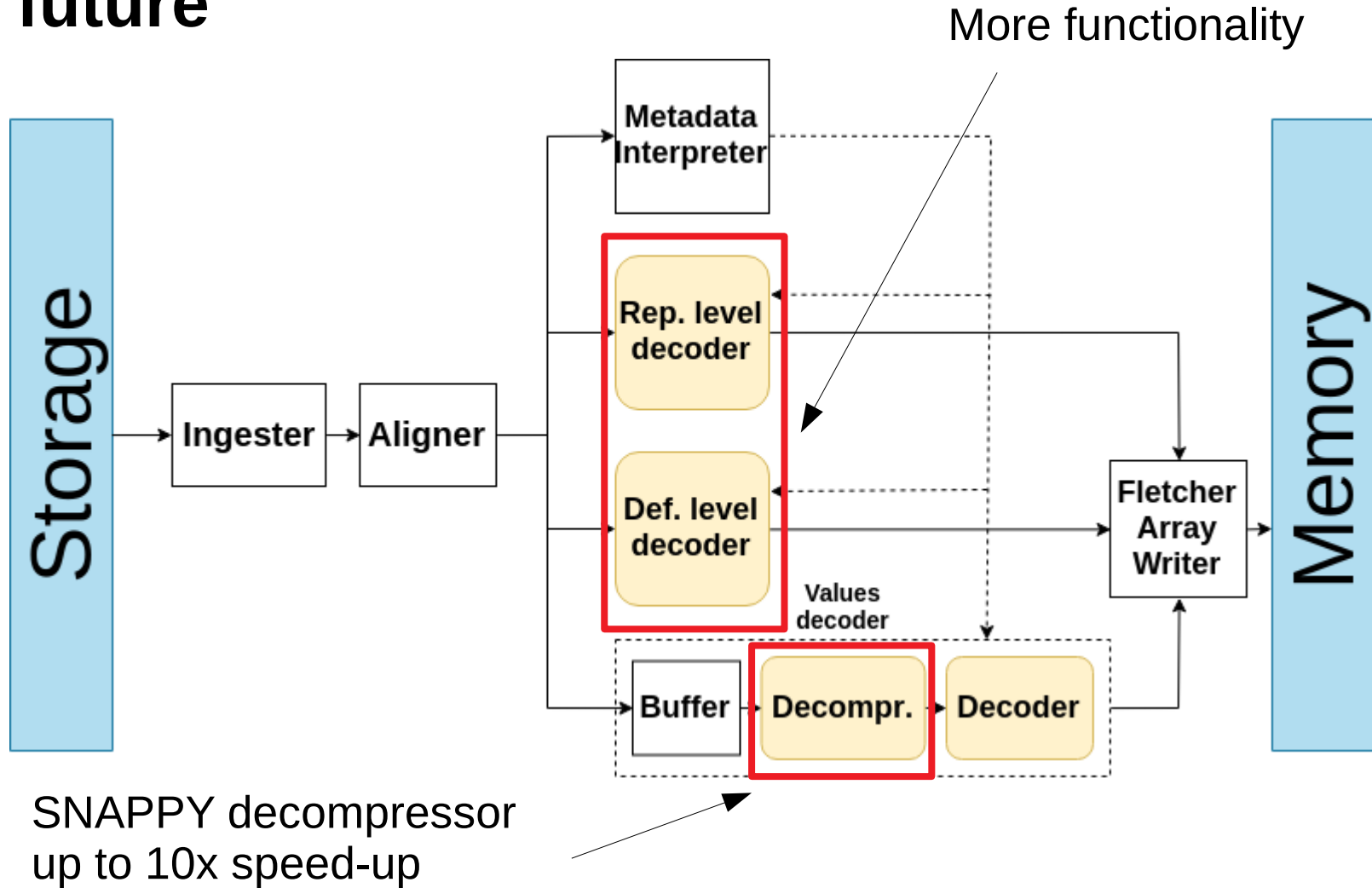
# Performance with delta decoding



Single ParquetReader implemented on XCVU9P

- Area utilization
  - CLB 3.20%
  - BRAM 2.85%
- Estimated power usage
  - ~1 W

# The future



# Conclusion

- 1.75x speed up shown for delta encoding
- Single ParquetReader uses only 3.2% of CLBs
- Promising future with decompression hardware

# High-throughput conversion of Apache Parquet files to Apache Arrow in-memory format using FPGAs

Lars van Leeuwen (l.t.j.vanleeuwen@student.tudelft.nl)

Johan Peltenburg, Jian Fang, Zaid Al-Ars, Peter Hofstee

