Chair of Aerodynamics and Fluid Mechanics
Department of Mechanical Engineering
Technical University of Munich

# Probabilistic Programming for Scientific Discovery

Lecture 4

Ludger Paehler
*Lviv Data Science Summer School*

July 30, 2020

# Table of Contents

# **Outline**

Interacting with Scientific Simulators: The Engineering Challenge

    Merlin

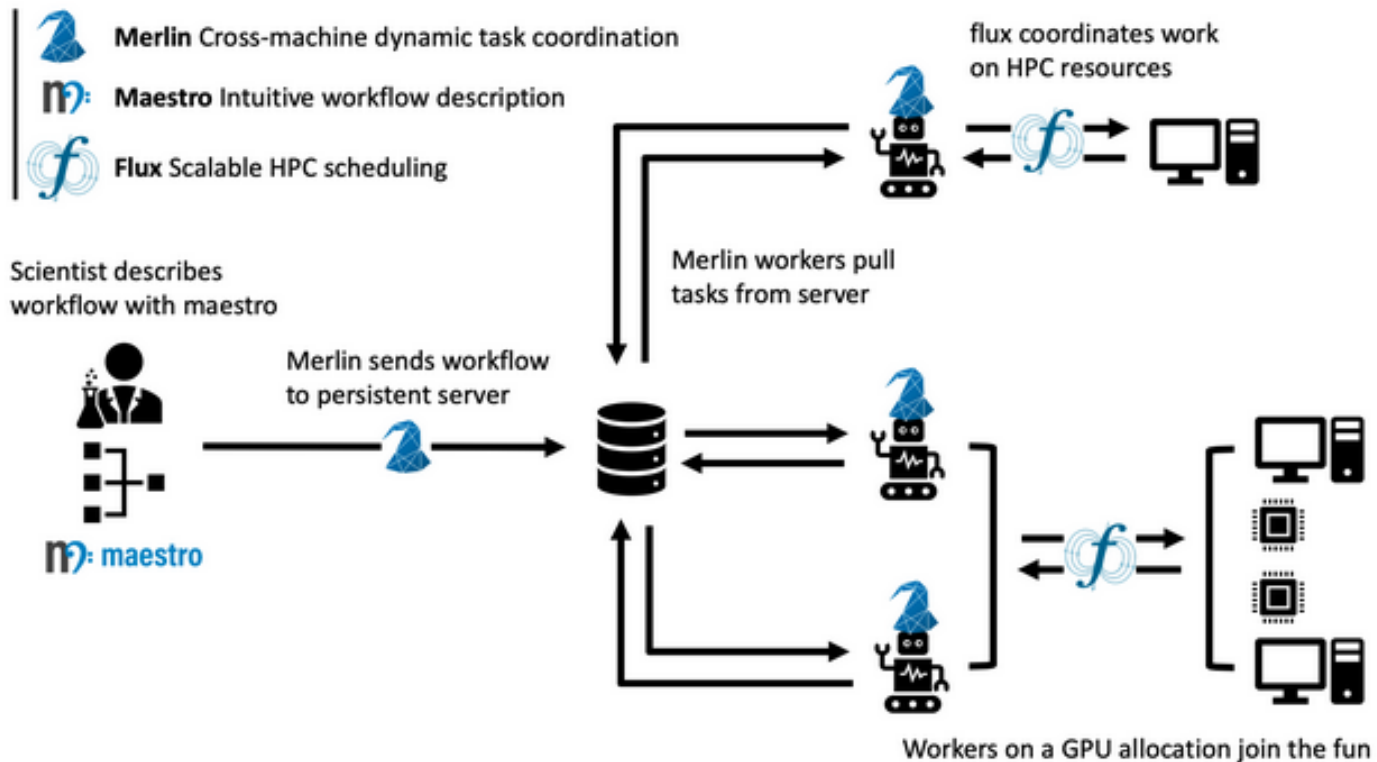    Reverb

    Apache Arrow

    XLA

# The Challenge

- When scaling our application/workflows to scientific applications we are faced with a multitude of problems:
  - Data movement
  - Amount of data
  - Concurrent processes, which might be prone to crashing
  - Application portability
- This is the point where probabilistic programming turns into an HPC engineering problem, where we have choose our frameworks wisely to get the most performance out of the hardware
- Offshot: Will also run much more efficient/faster on small computers

# Merlin



- Manager for HPC-focuses simulation workflow, which can include multiple supercomputers, cloud components, as well as heterogeneous clusters.
- Made to scale to multiple hundred million simulations at a time
  - (Broken down into individual batches of simulations)
- Can automatically take advantage of spare computing capacity
- Persistent in case of individual servers failing
- Can be likened to a microservices-focussed architecture for HPC simulations
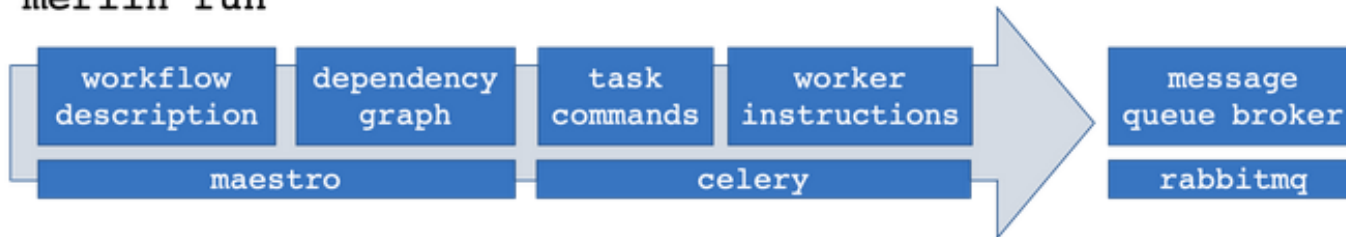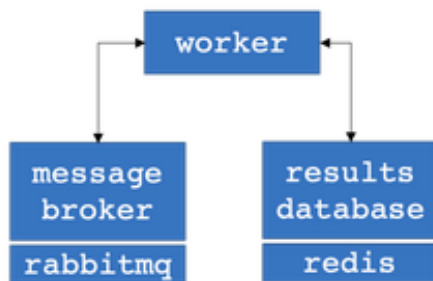
# Merlin

# Merlin

- Asynchronous task queuing library (Merlin)

- Dynamic task queuing (Celery)

- Description of our model is contained in a central workflow file

- Resilience and support for multiple users and queues (RabbitMQ)

- Scalability and fast database access (Redis)

# Merlin

# Merlin
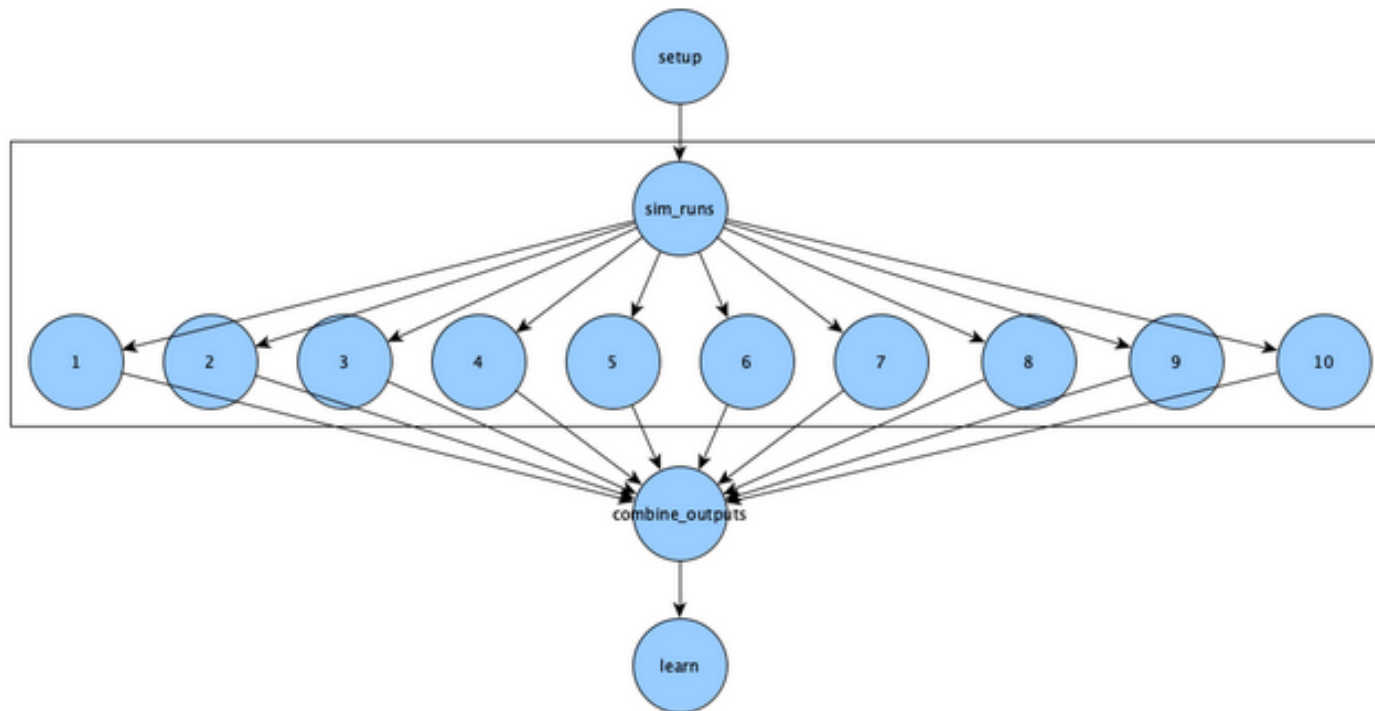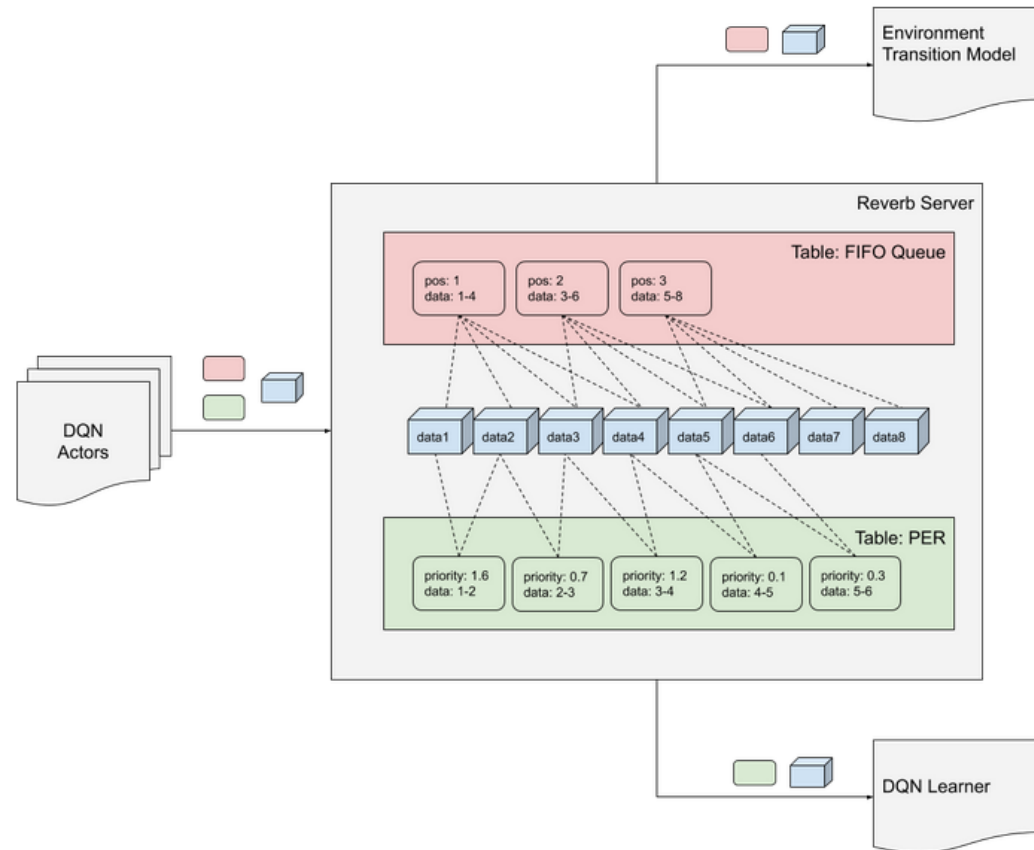


**Figure:** Any workflow we can define as a directed-acyclic graph can be performed.

# Reverb

- Experience replay system originally built for off-policy reinforcement learning
  - Can be adapted to probabilistic inference as well
- Reverb server consists of one or more tables, where tables also define sample and removal strategies, maximum item capacity, and a rate limiter.
- Multiple item items can refer to the same data element, as entries in tables only contain references to the data
- Data is deleted if there is no reference to it anymore

# Reverb

# Apache Arrow

- OSS Community initiative conceived in 2015
- Sits right at the intersection of big data, database systems, and data science tools
- Core of the project: Language agnostic open standards to accelerate in-memory computing
- Components have public APIs
- Arrow is front-end agnostic
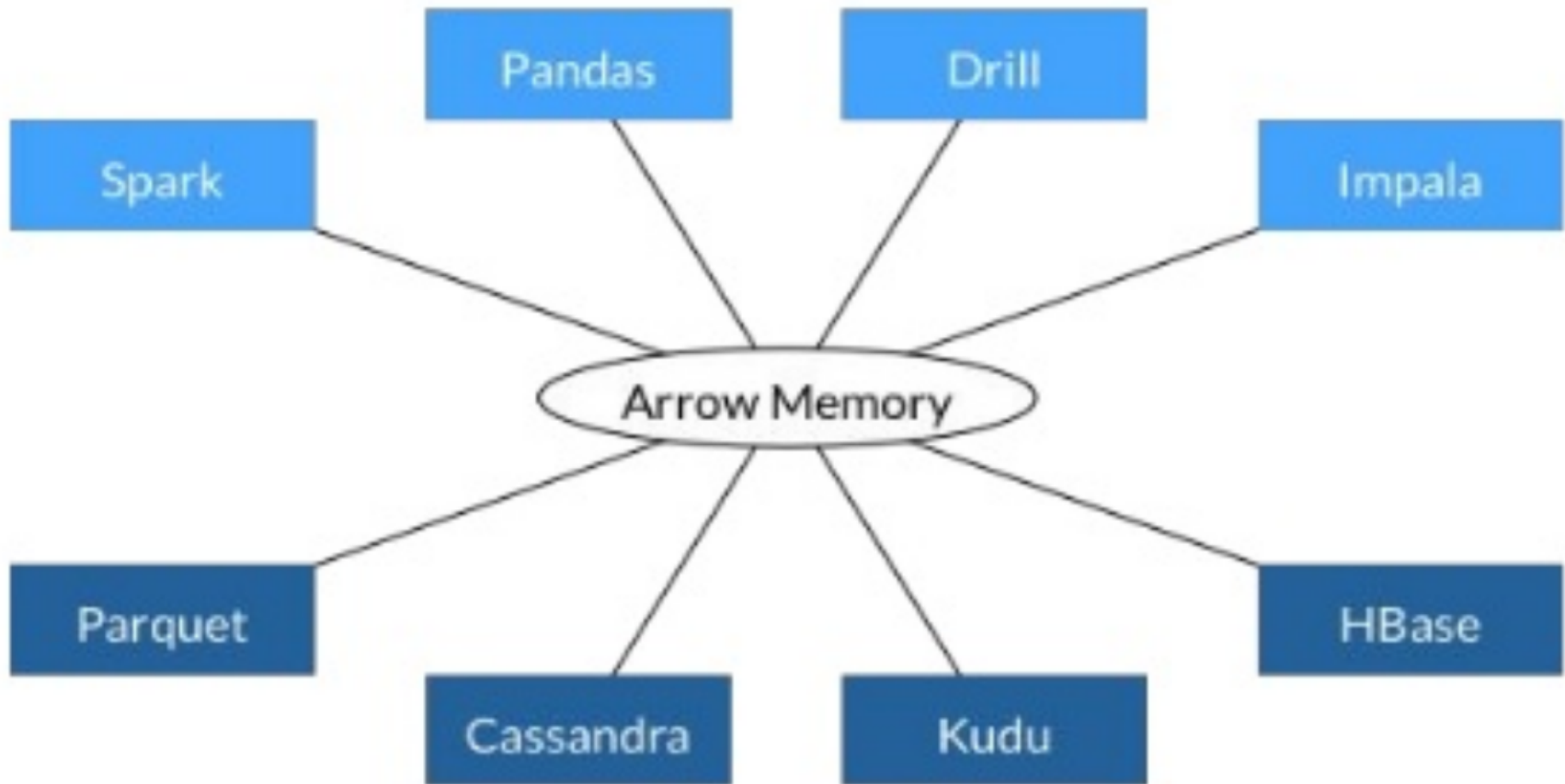- Dataframe backend for NVIDIA Rapids

https://github.com/apache/arrow

| Front end API |
| :---: |

| Computation Engine |
| :---: |

| In-memory storage |
| :---: |

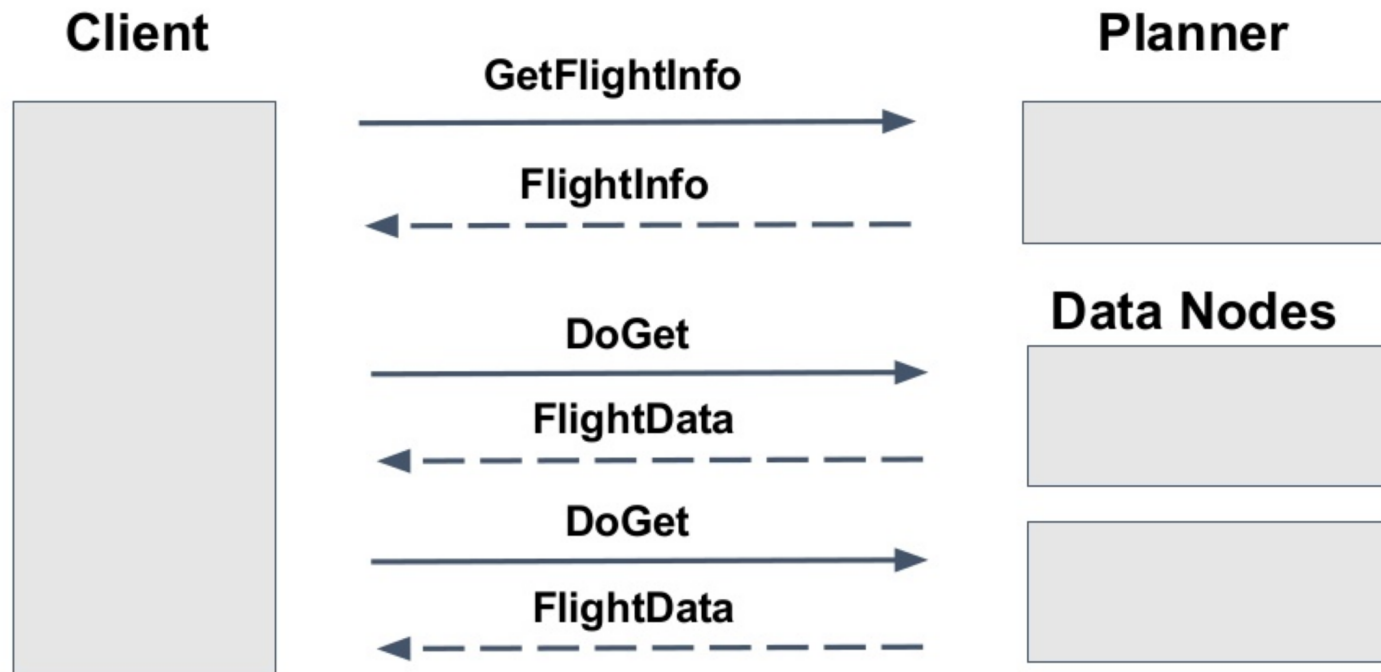| IO and Deserialization |
| :---: |

# Apache Arrow

# Apache Arrow

- Applied to eliminate serialization overhead in data interchange
- Improve CPU/GPU in-memory processing efficiency
- Simplify architectures
- Organized for cache-efficient access on CPUs/GPUs
- Optimized for data locality, SIMD, and parallel processing
- In-memory encoding, compression, and sparseness are in the works
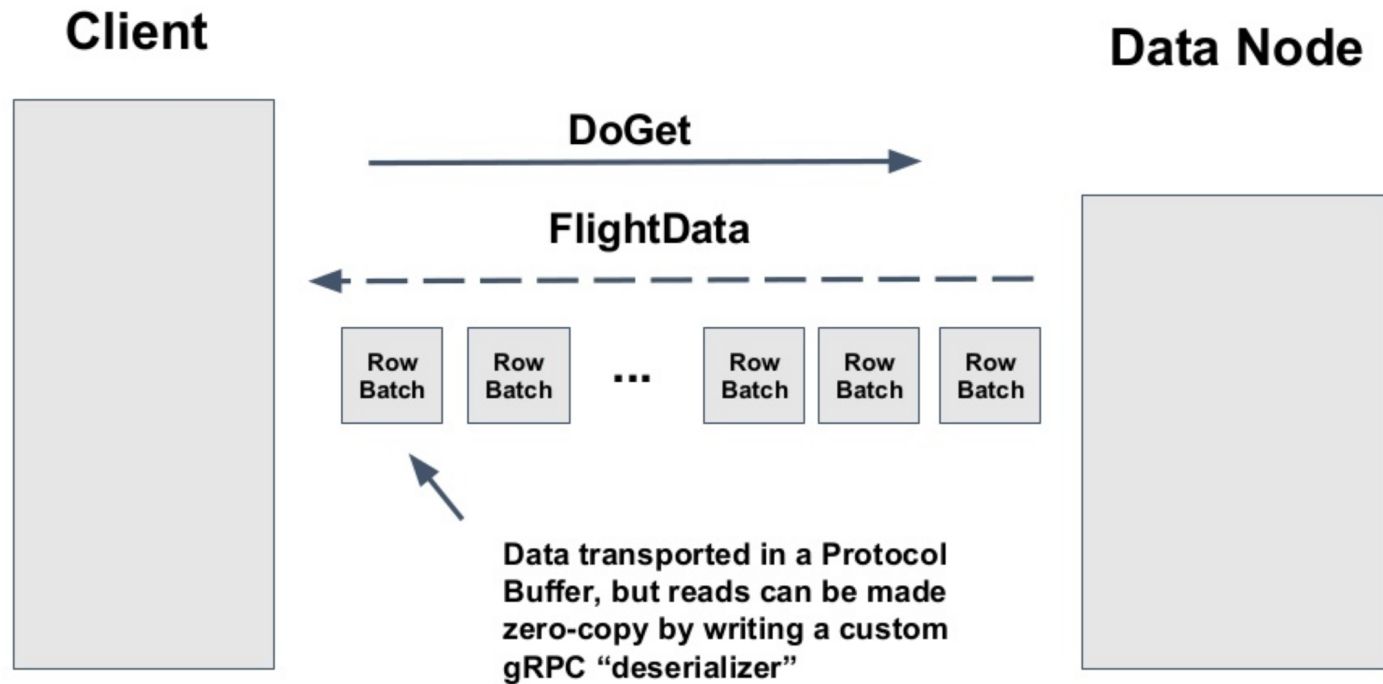
# Apache Arrow

Apache Arrow Flight

- Framework to define custom data services and receive Arrow data natively
  - Avoids any copying or deserializatio
- Low-level optimizations
- Made for continous extreme-scale datastreams, i.e. also useful for HPC applications
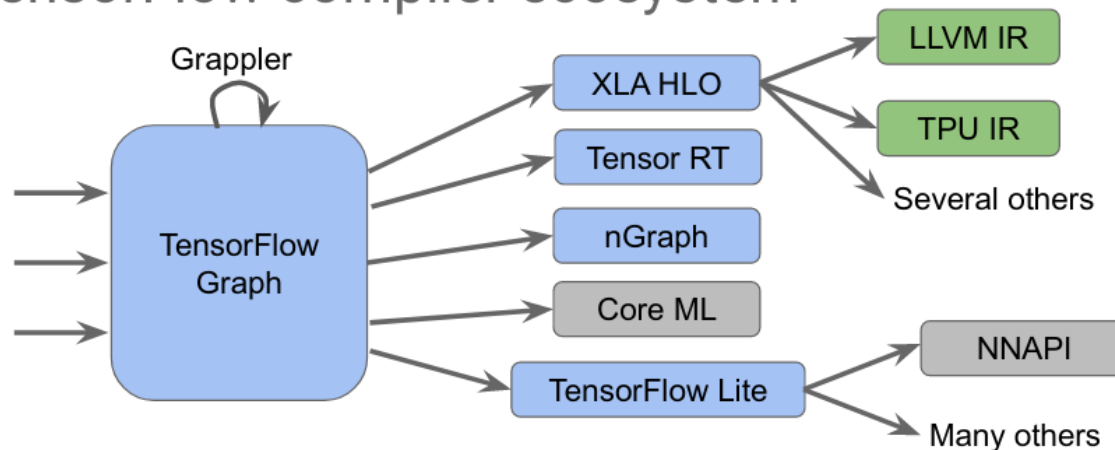
# Apache Arrow

# Apache Arrow

# Machine Learning Compilers

The TensorFlow compiler ecosystem

Grappler

TensorFlow Graph → XLA HLO → LLVM IR, TPU IR, Several others

TensorFlow Graph → Tensor RT

TensorFlow Graph → nGraph

TensorFlow Graph → Core ML

TensorFlow Graph → TensorFlow Lite → NNAPI, Many others

Many "Graph" IRs, each with challenges:

- Similar-but-different proprietary technologies: not going away anytime soon
- Fragile, poor UI when failures happen: e.g. poor/no location info, or even crashes
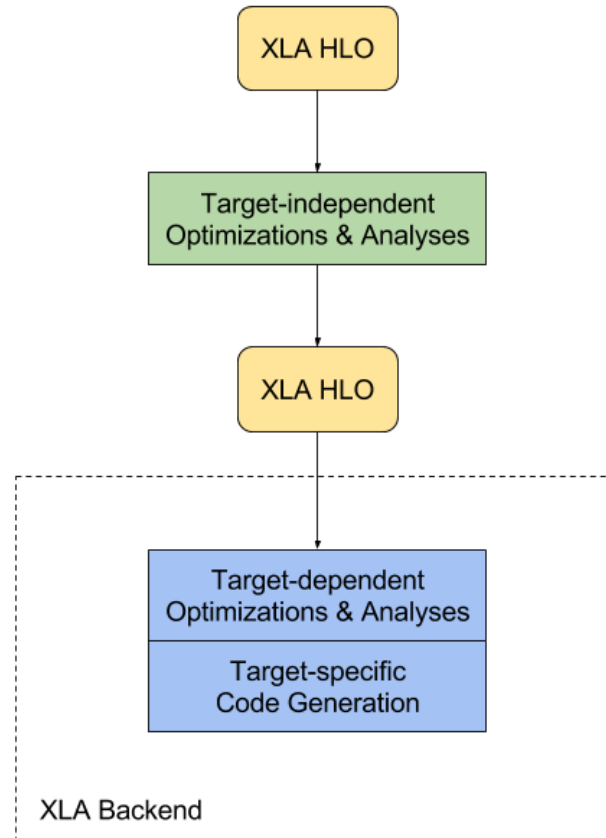- Duplication of infrastructure at all levels

Google

# XLA

Purpose behind XLA:

- Improve the execution speed
- Improve memory usage
- Reduce reliance on custom operations
- Reduce mobile footprint
- Improve portability
  - Accelerator-producers just need to implement the XLA API, and our code will run seamlessly on accelerators of all types

# XLA

# The Approach

- Data movement $\longrightarrow$ Apache Arrow

- Amount of data $\longrightarrow$ Apache Arrow + Compression

- Concurrent processes, which might be prone to crashing $\longrightarrow$ Merlin

- Application portability $\longrightarrow$ XLA or MLIR