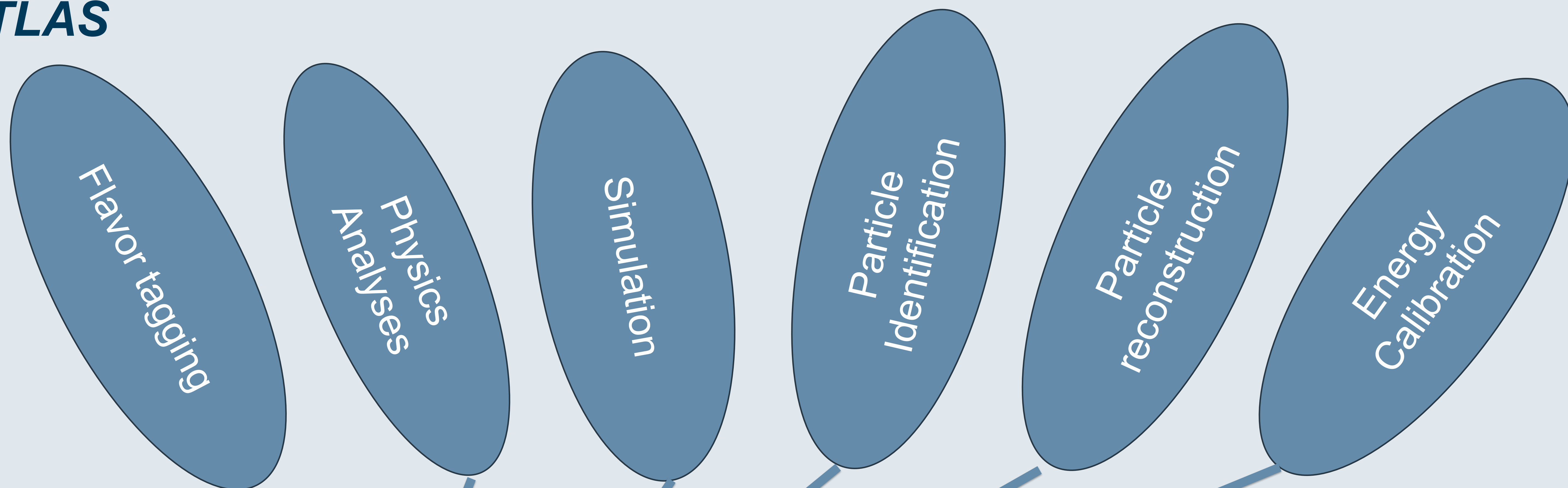
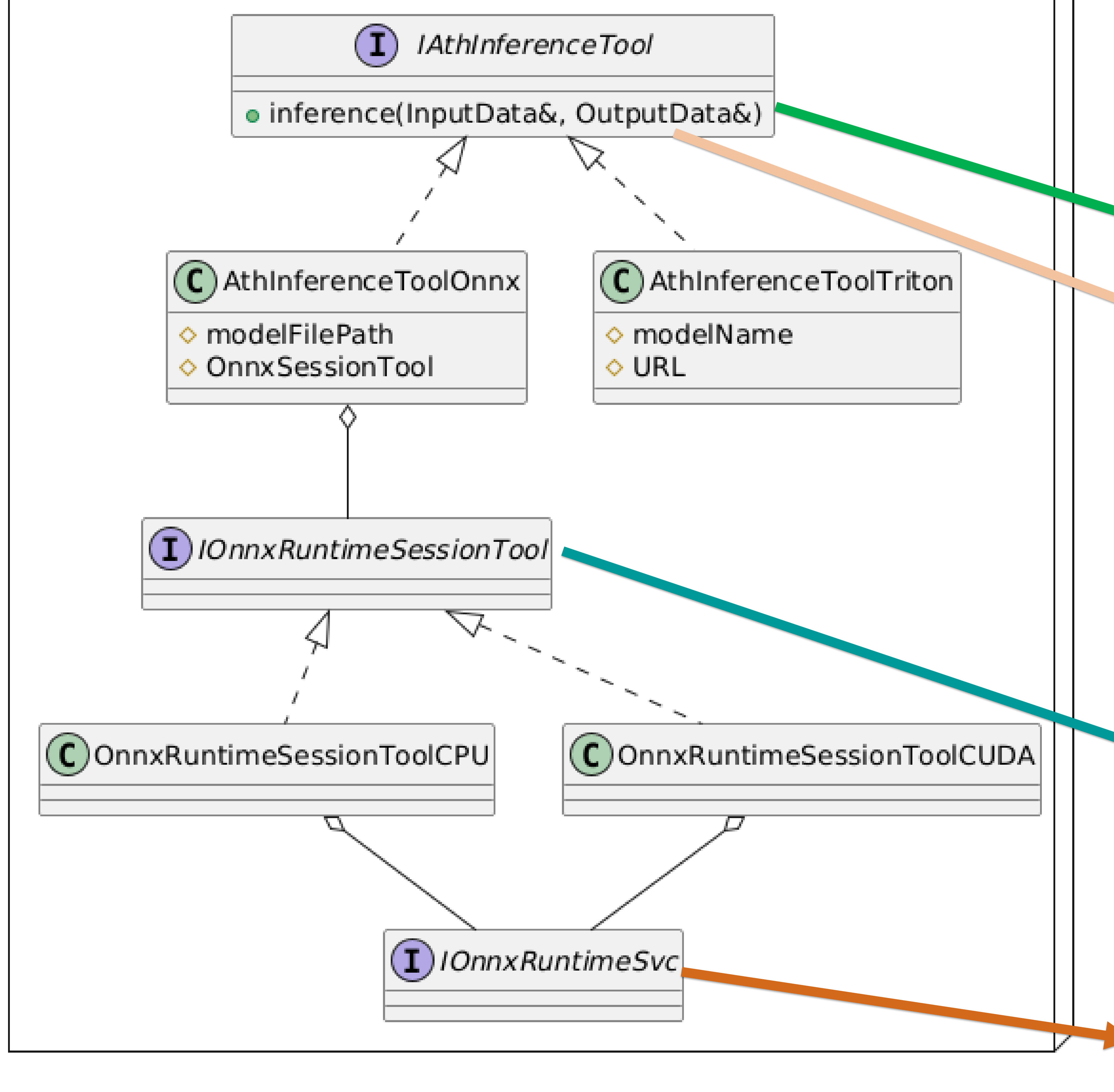


ML applications in ATLAS



Control

AthInferenceTool



OnnxRuntime [1] is a high-performance inference engine designed to accelerate machine learning models across various platforms and hardware, supporting the ONNX (Open Neural Network Exchange) format for interoperability.

A **common interface** for Machine Learning inference in Athena that are implemented by OnnxRuntime and Inference as a Service [2].

Input and output data is a **dictionary** that maps input and output names to a data object that contains the tensor shape and data. The dictionary allows dynamic number of input and output data.

Tensor data is represented by **std::variant**, allowing different data types.

Lightweight **Athena Tool**, each for a specific backend such as CPU or NVIDIA GPU, allows the possibility that other backends can be used.

One ORT::Env as **Athena Service**, shared by all Athena Algorithms. Its functionality includes 1) logging, 2) managing thread pools for running ONNX Sessions (To-be-done), 3) managing memory allocations for ONNX models (To-be-done)