



Semiconductor
Research
Corporation



CUBiC

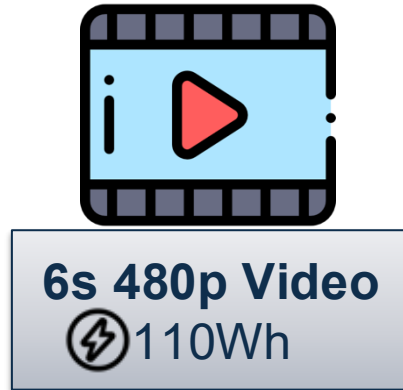
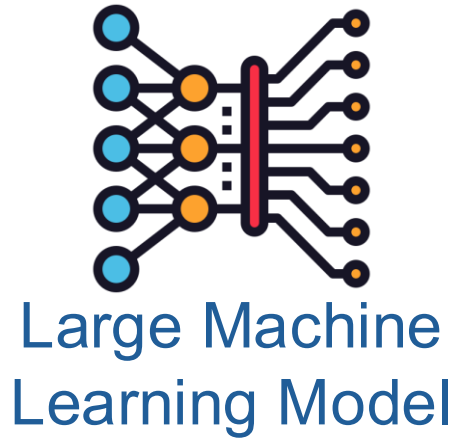
Center for Ubiquitous Connectivity

Enabling Energy Efficiency for Distributed Machine Learning

Tawhid Bhuiyan and Tanvir Ahmed Khan



Large Machine Learning Models Consume a Lot of Energy



thin, medium-well steak
⚡ 220Wh

Essential to model **end-to-end energy consumption** for optimization

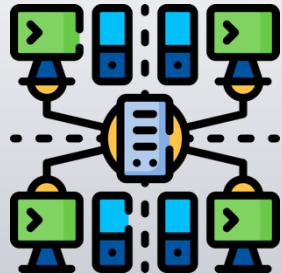
Source: <https://archive.is/l3Klo>

What Does Modeling End-To-End Energy Mean?

Problem Definition

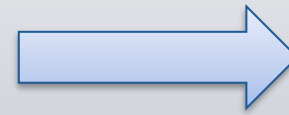


ML Workload



System Setup

Inputs



Energy
Breakdown

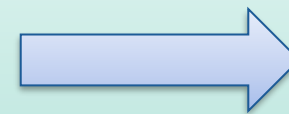
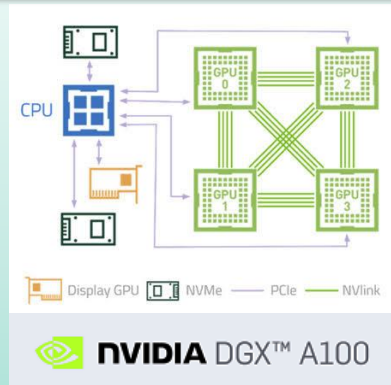


Time
Breakdown

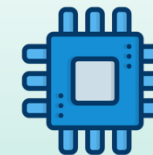
Outputs

Example

 Llama 3



Compute



120J

6s

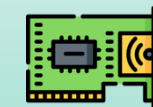
Memory



20J

3s

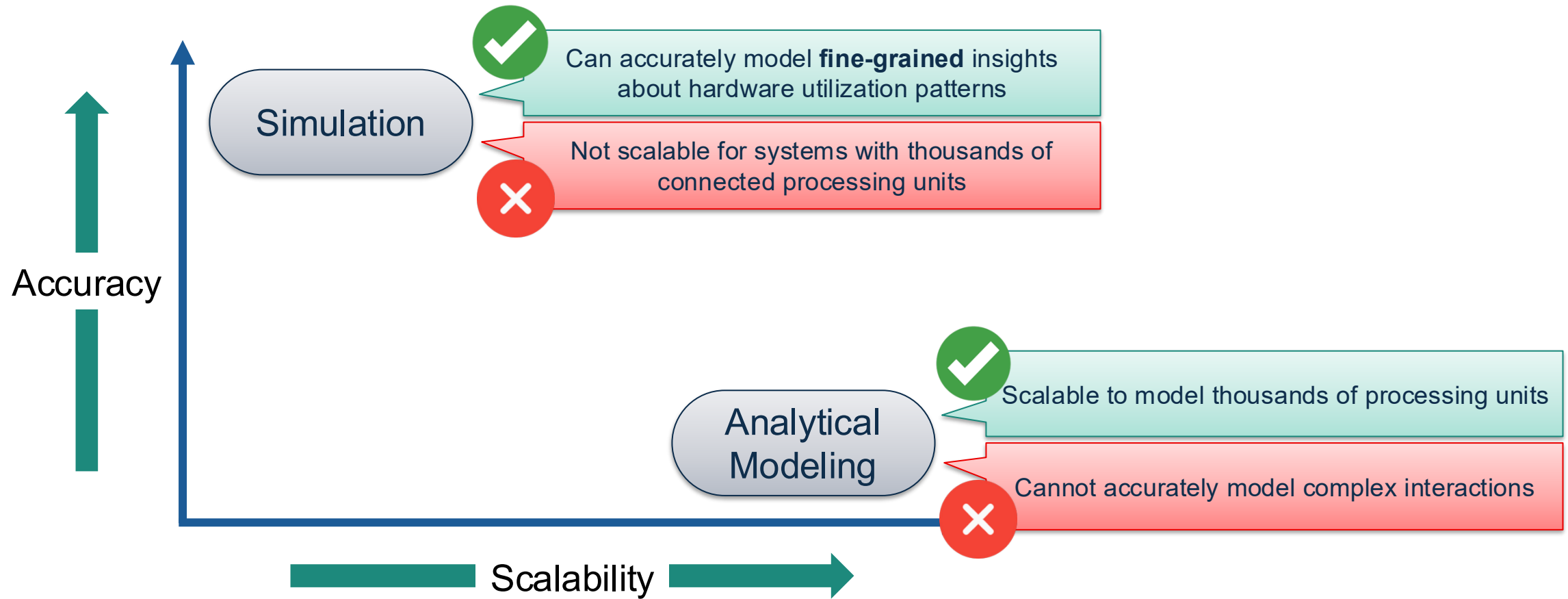
Network



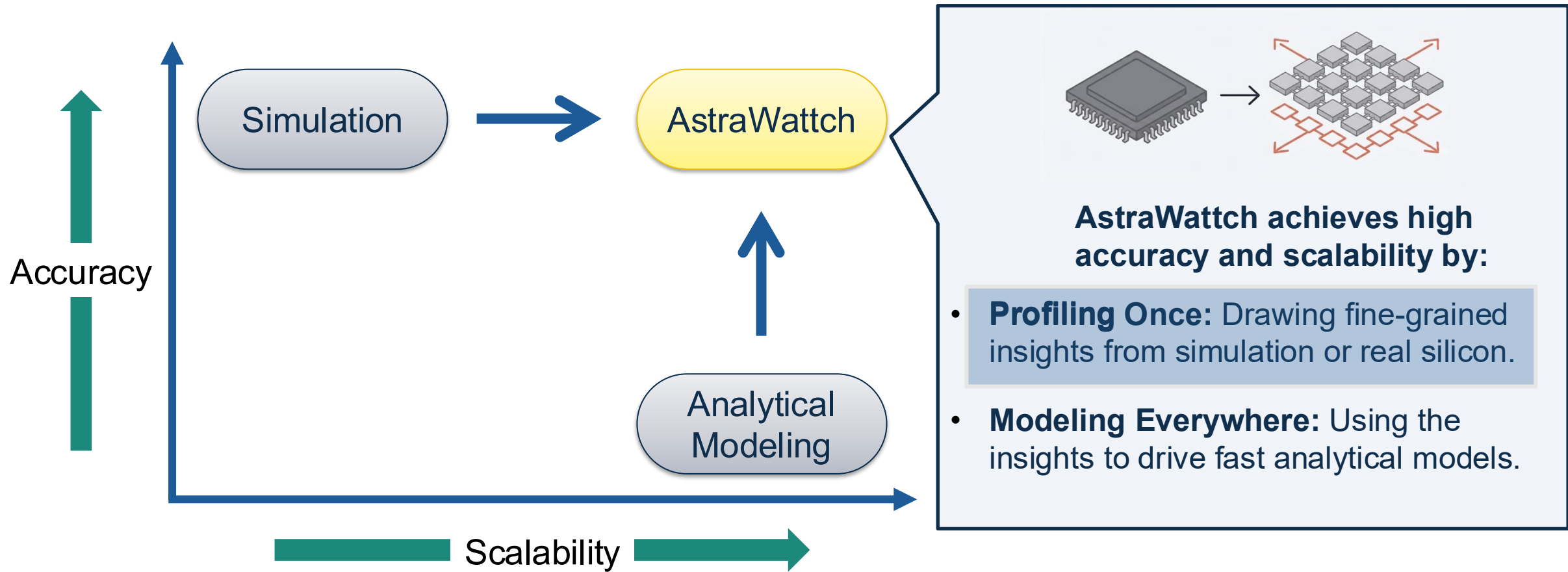
15J

2s

Challenges in Modeling End-To-End Energy

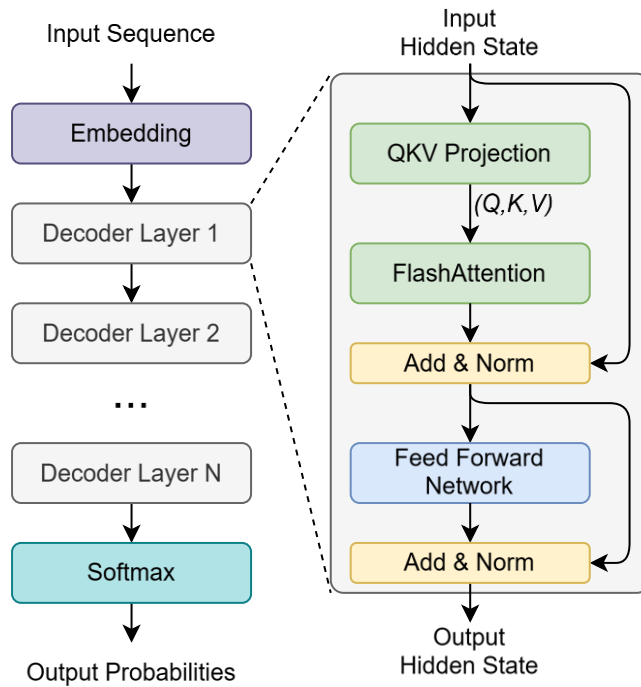


AstraWattch: Profile Once and Model Everywhere

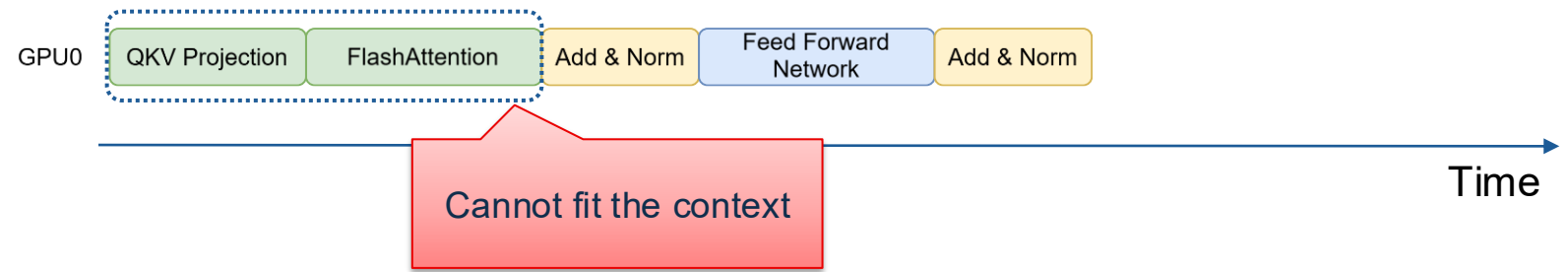


Preliminary results: AstraWattch enables accurate energy modeling in minutes

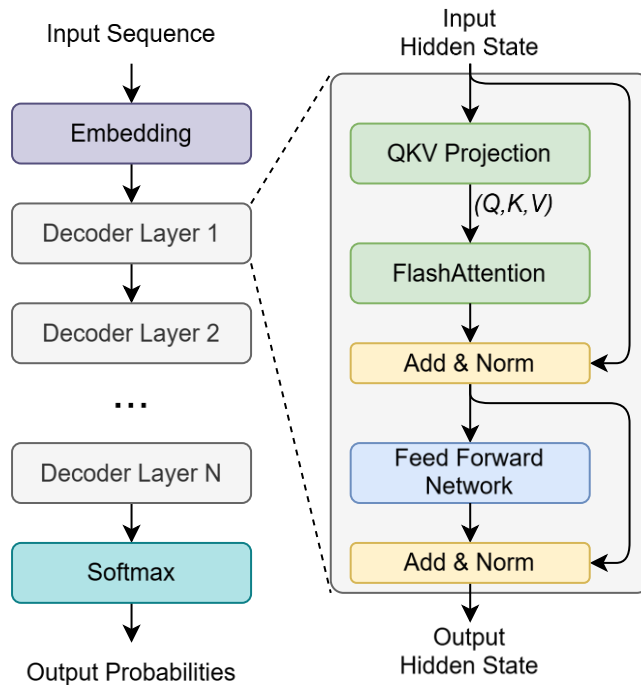
Profiling Machine Learning



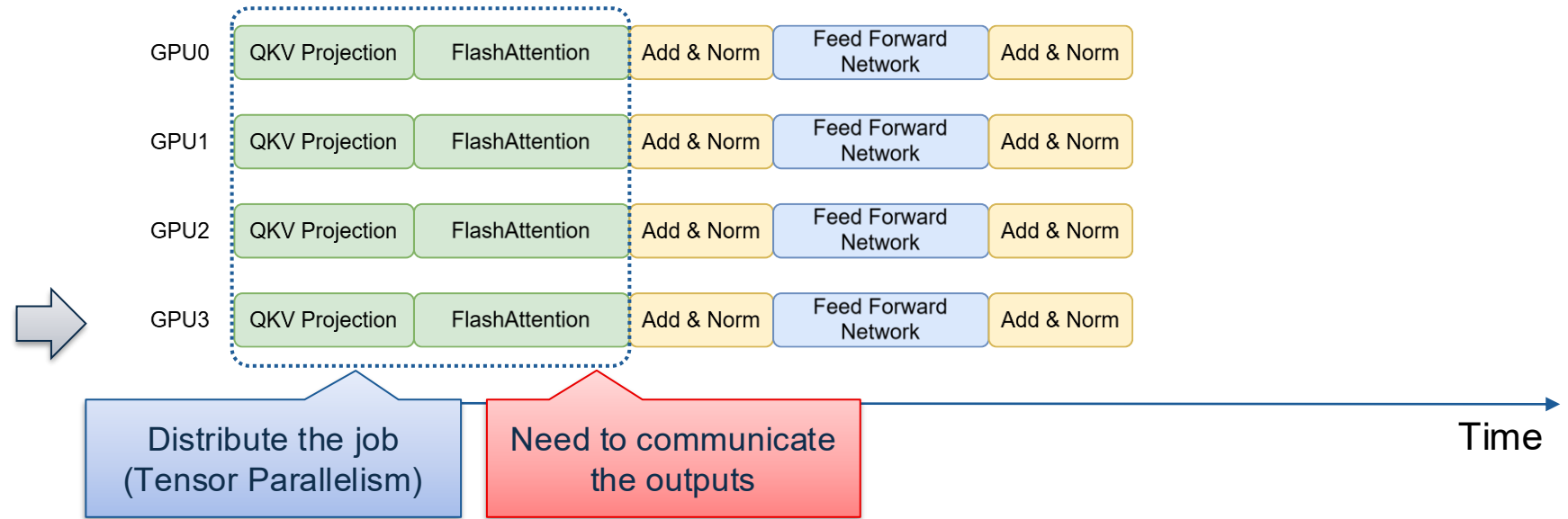
LLM Architecture



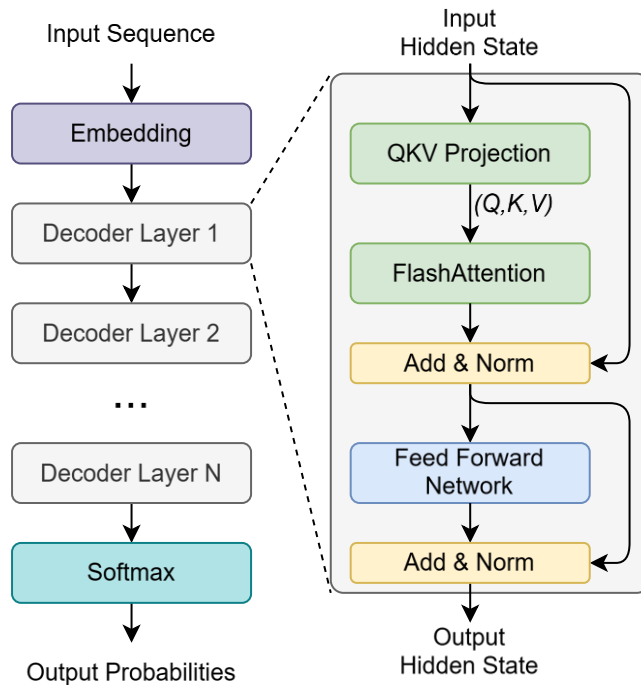
Profiling Distributed Machine Learning



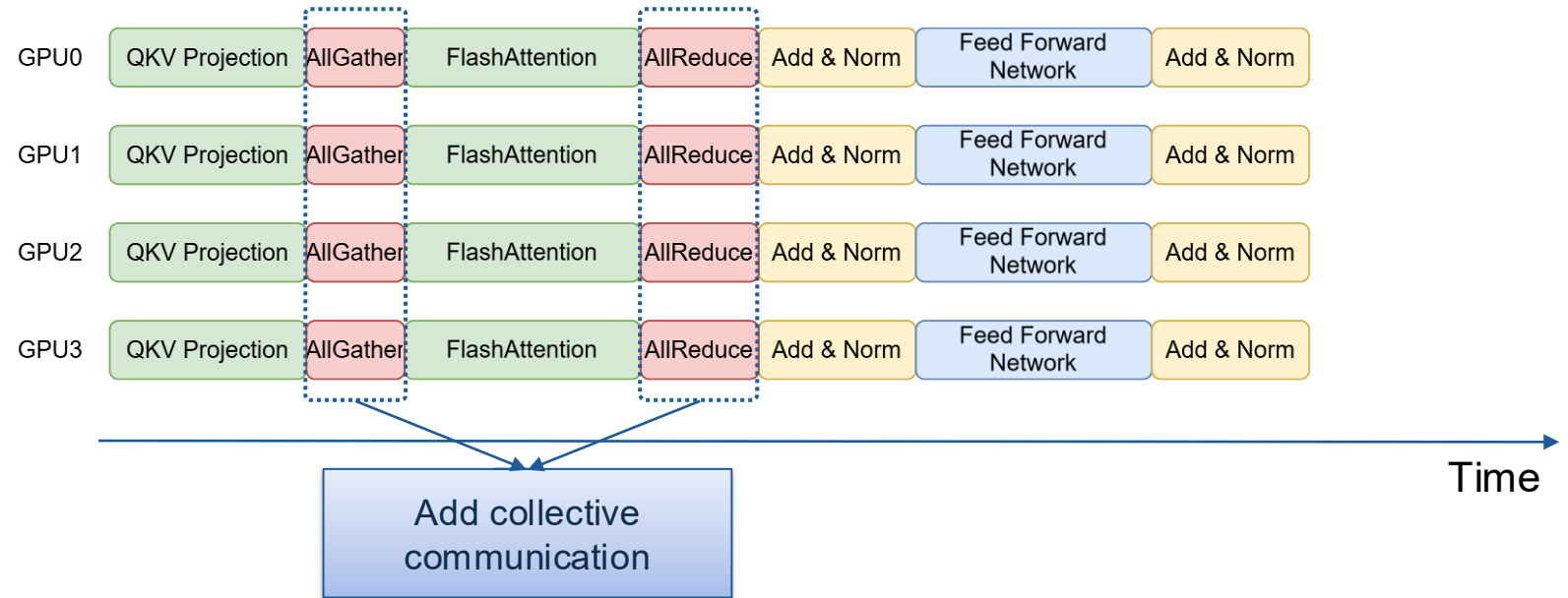
LLM Architecture



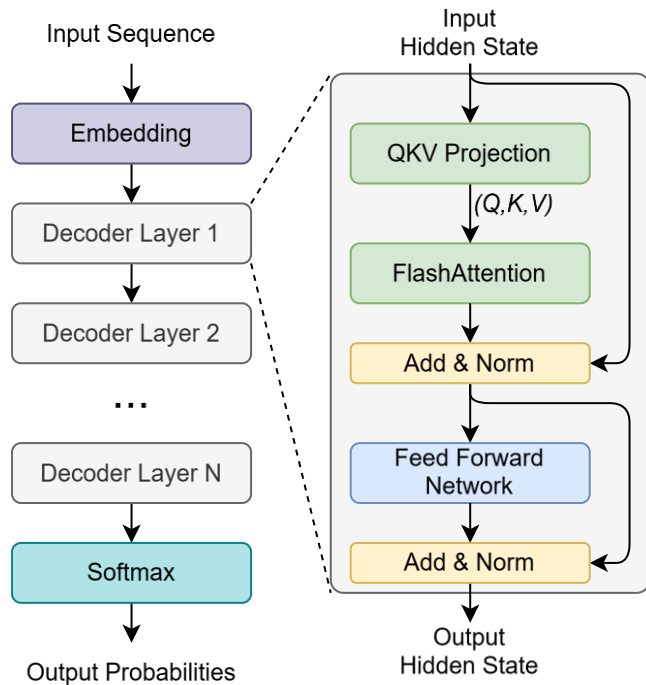
Profiling Distributed Machine Learning



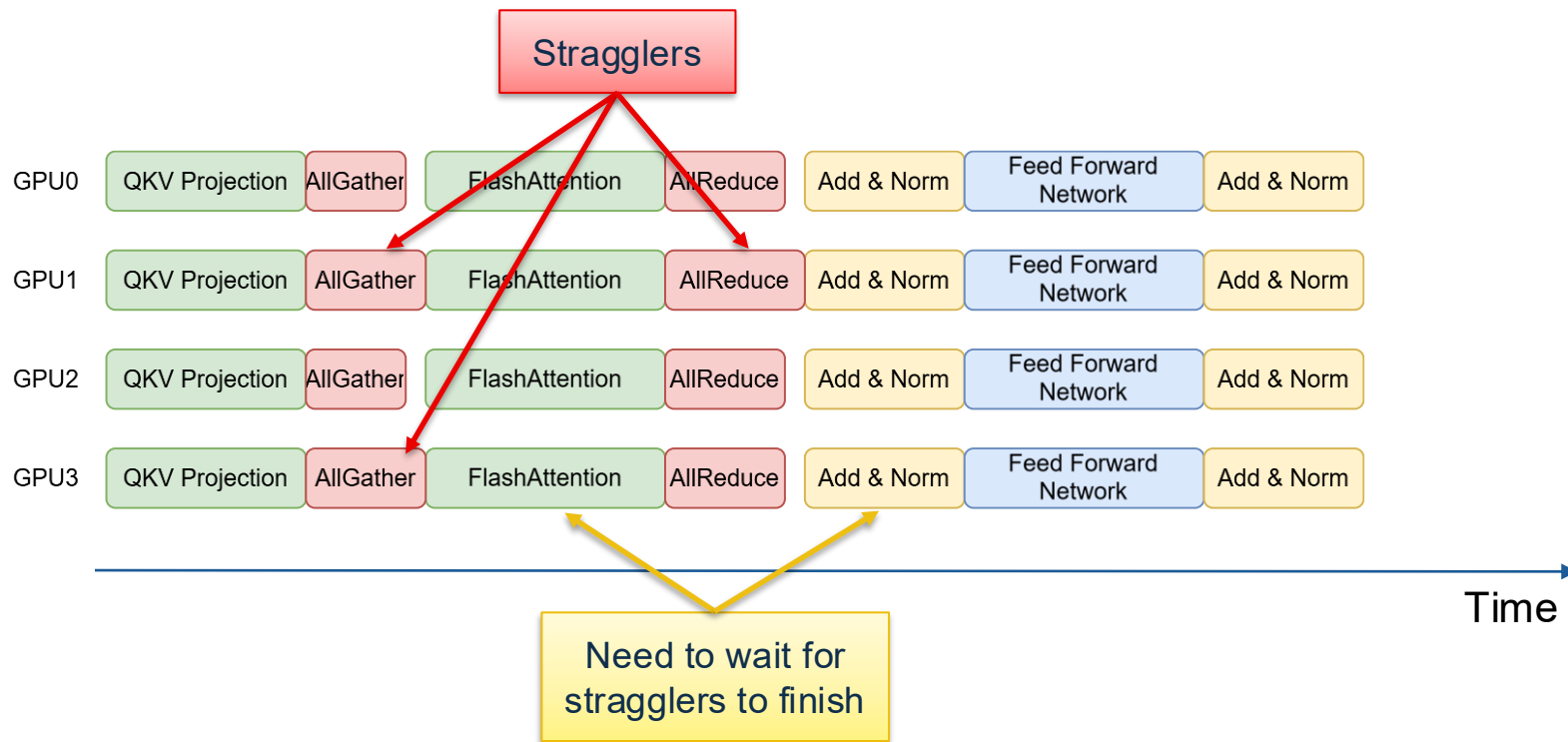
LLM Architecture



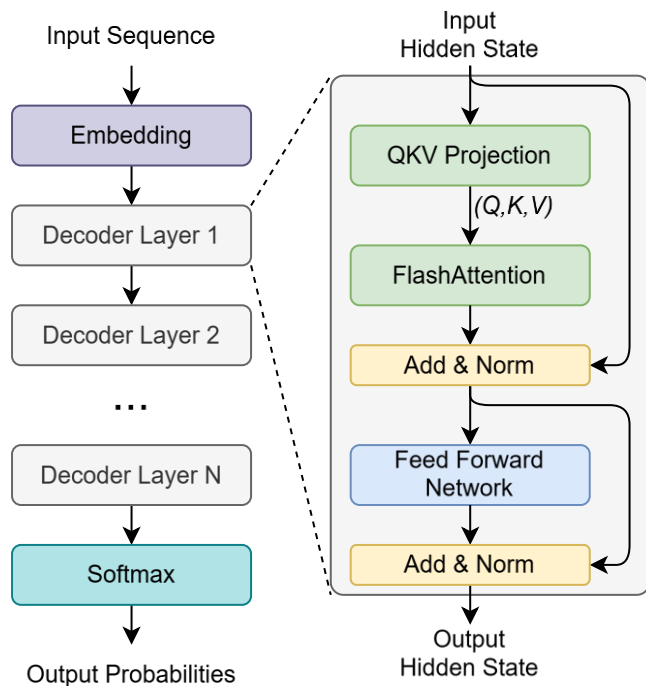
Profiling Insights



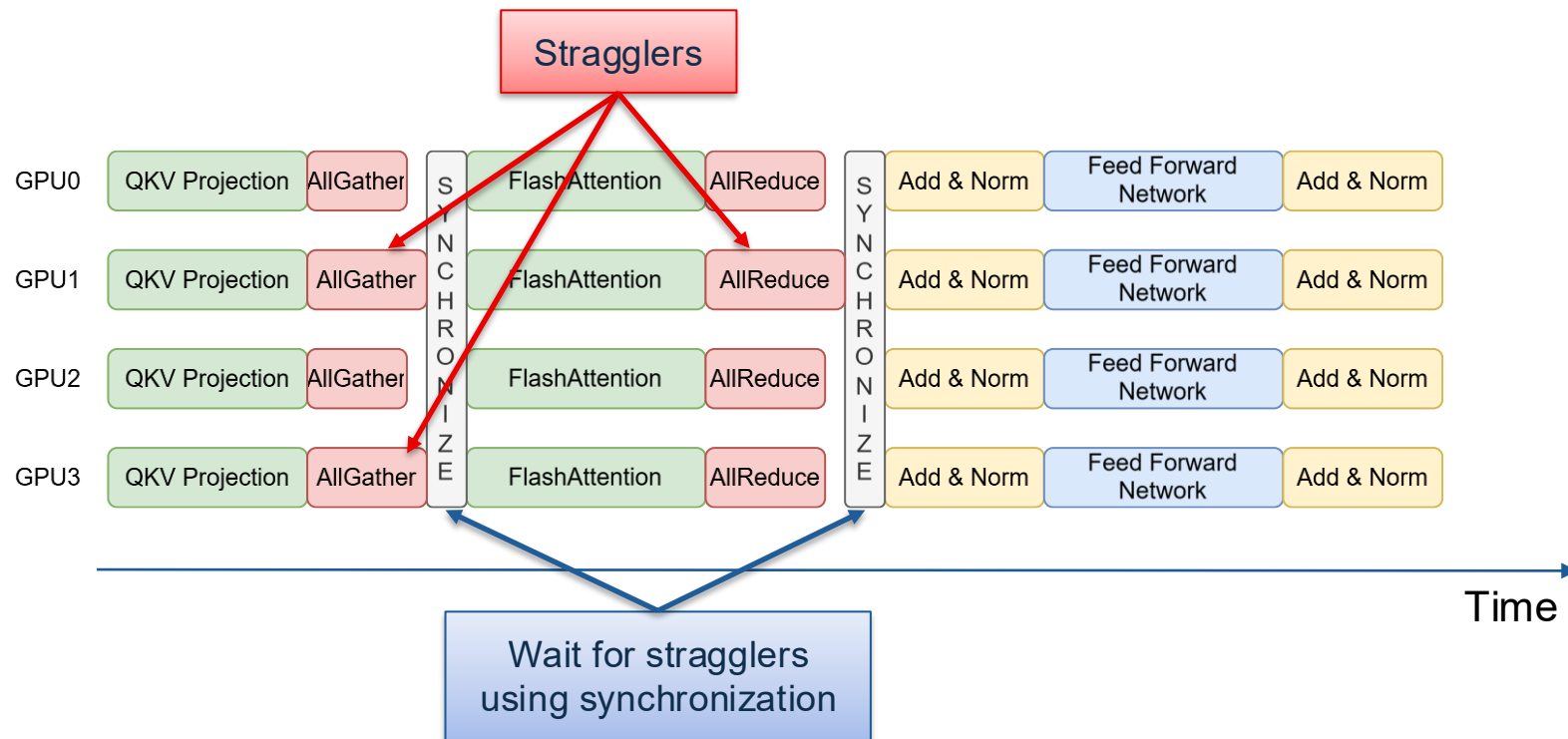
LLM Architecture



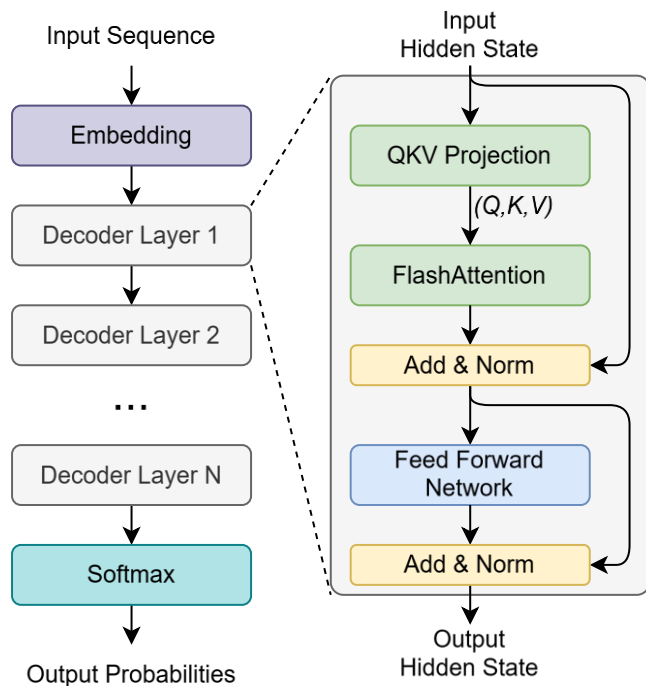
Profiling Insights



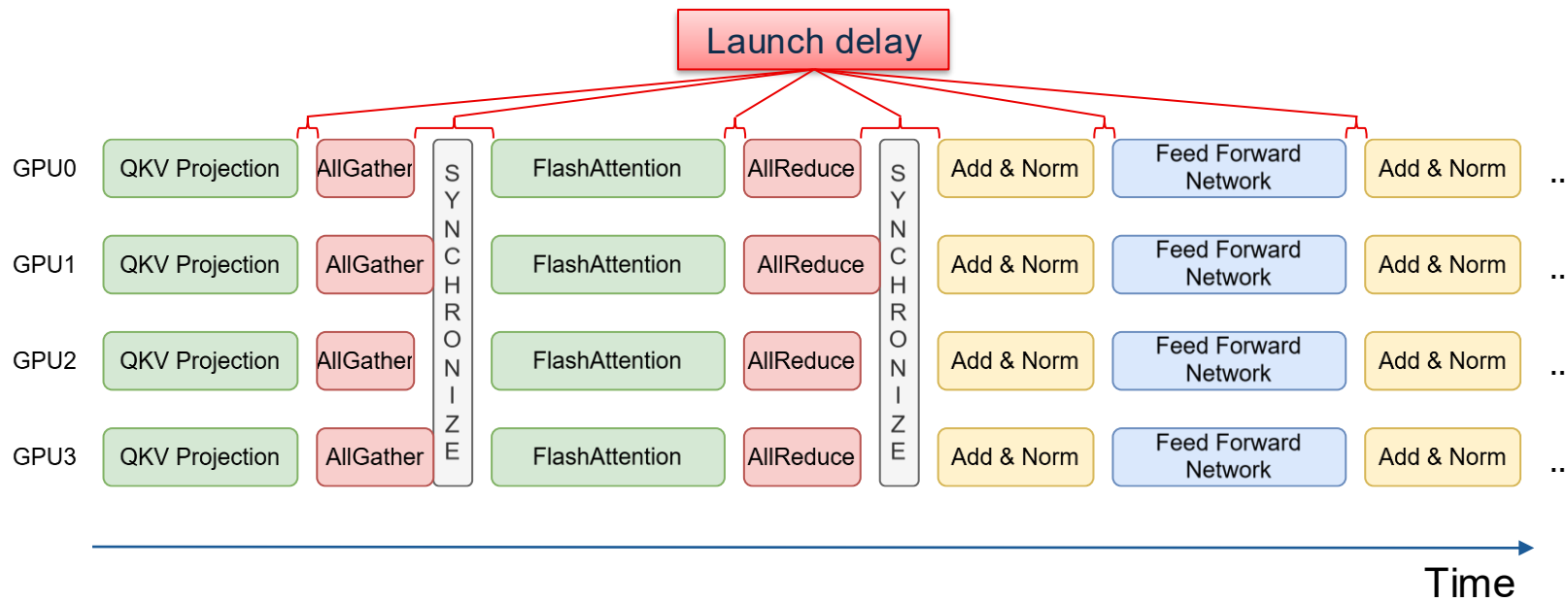
LLM Architecture



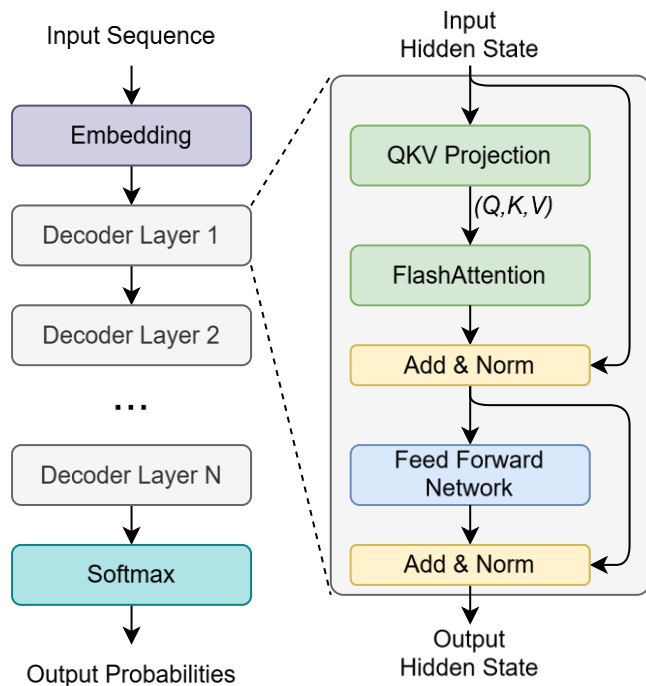
Profiling Insights



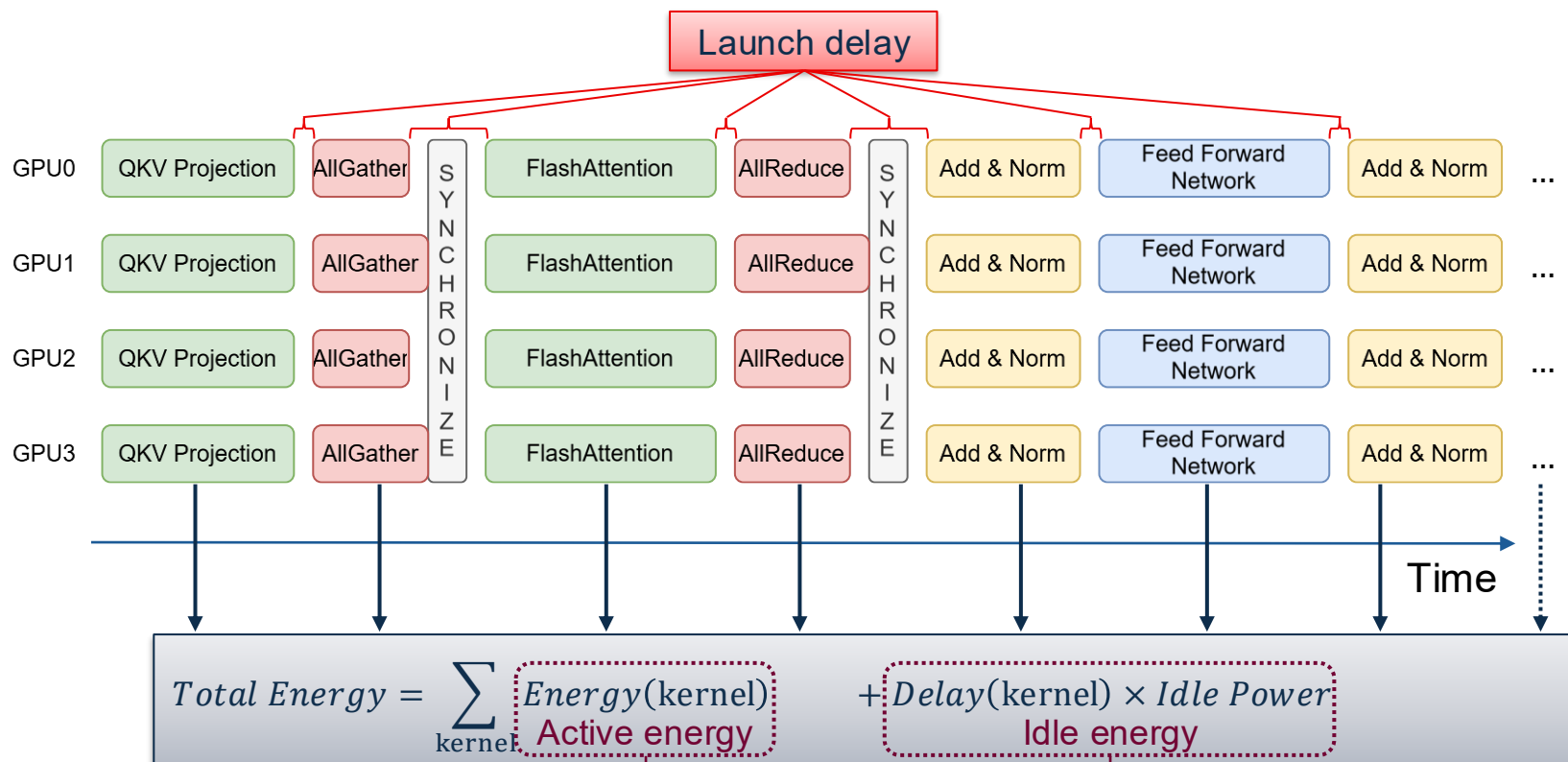
LLM Architecture



Profiling Insights For Energy



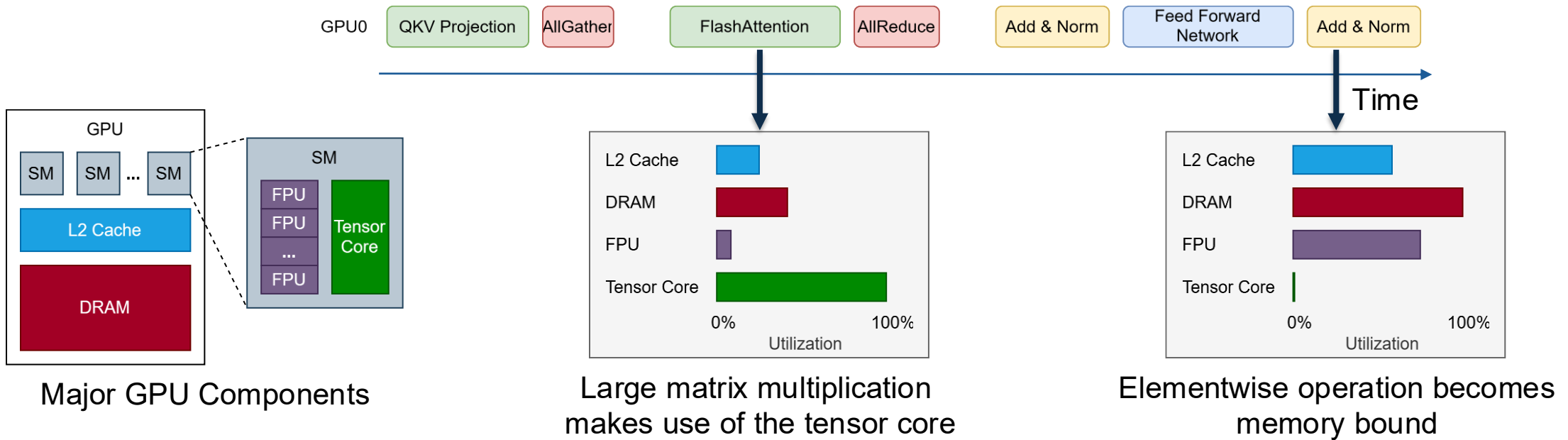
LLM Architecture



- Kernels are amenable to analytical modeling

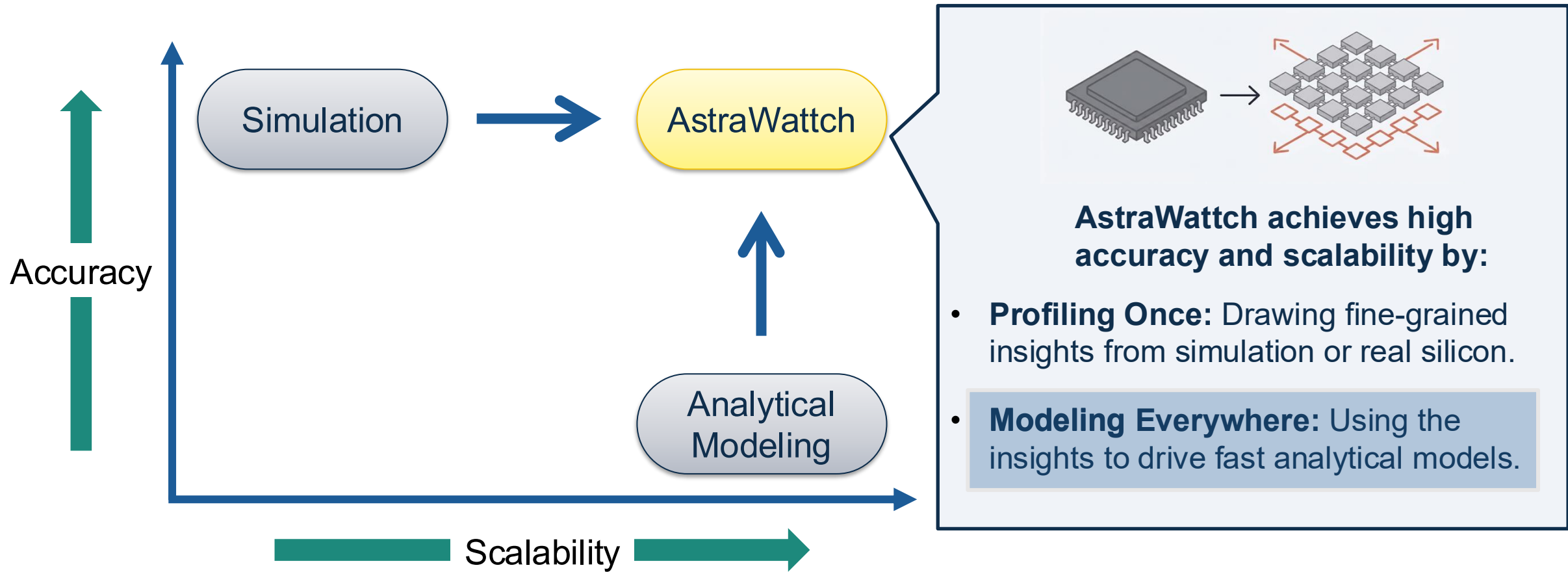
- Launch delays can be modeled as hardware-independent attributes

Profiling Kernels



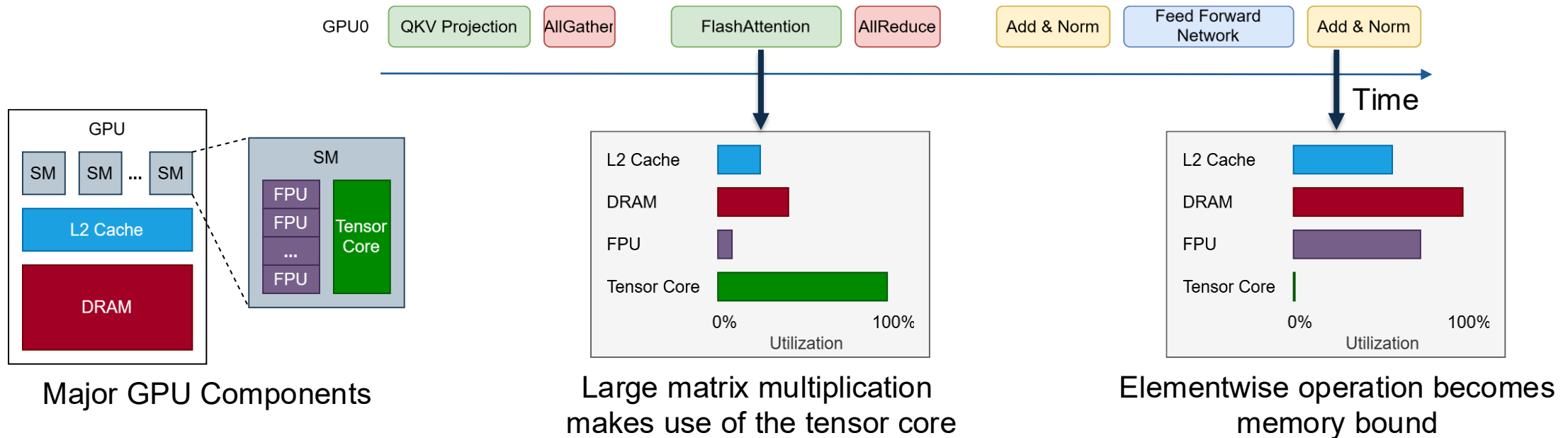
- Coarse-grained component utilization can accurately model energy

AstraWattch: Profile Once and Model Everywhere



Preliminary results: AstraWattch enables accurate energy modeling in minutes

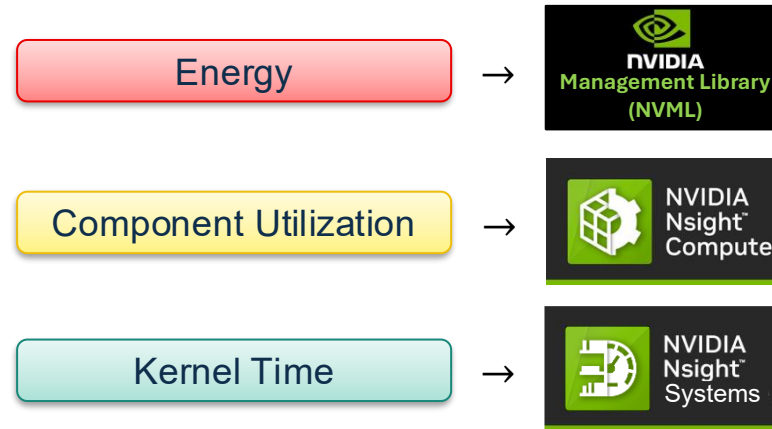
Modeling Active Energy



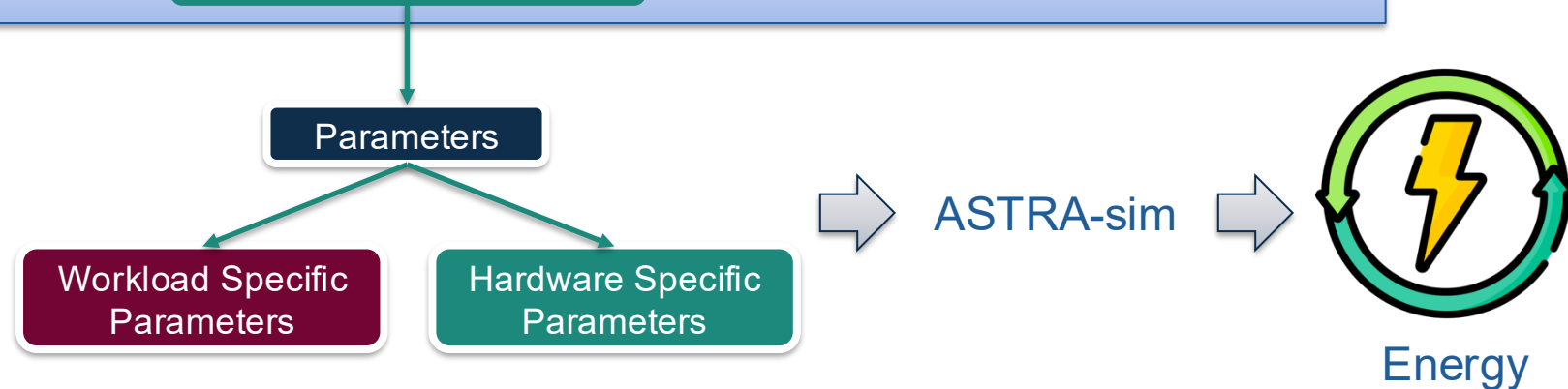
$$Energy(\text{kernel}) = \text{Sum of component energies for the kernel}$$

- We can build simple **analytical models for each kernel**
- One analytical model covers all invocations of a kernel

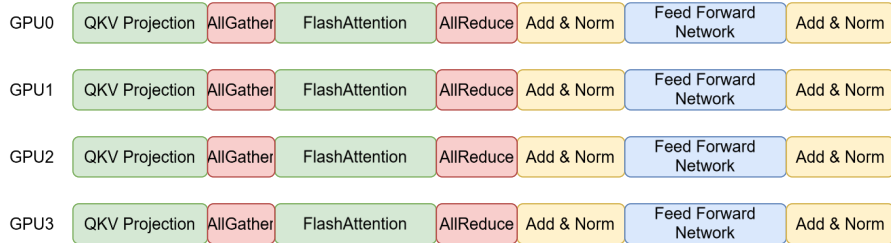
Calibrating the Model Parameters and Using Everywhere



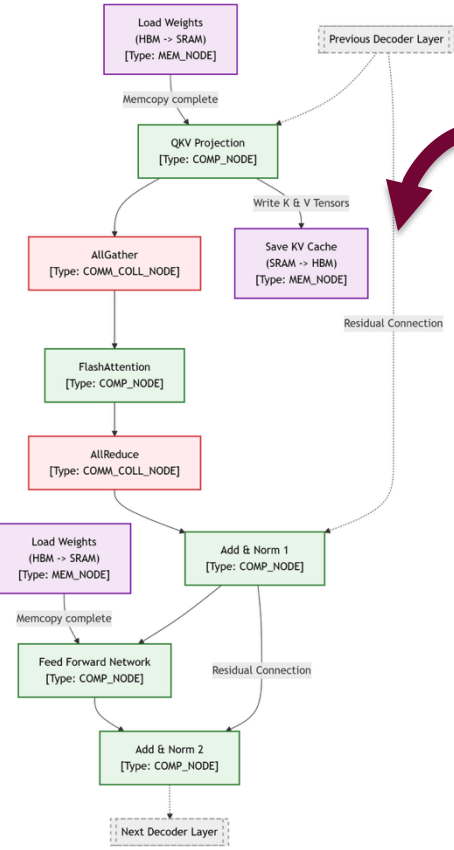
- Collect metrics from real or simulated GPUs
- Use **regression models** to estimate energy



ASTRA-sim's Modeling of Distributed Machine Learning

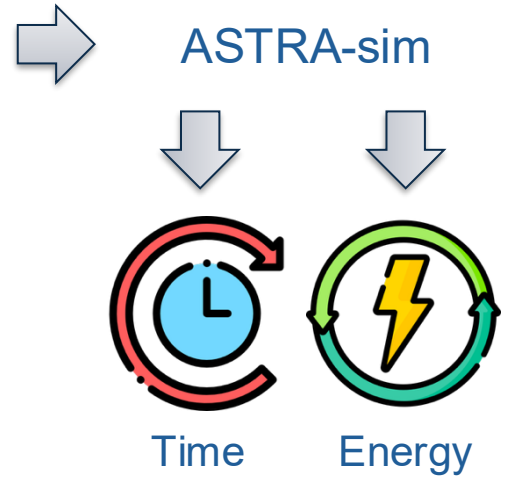
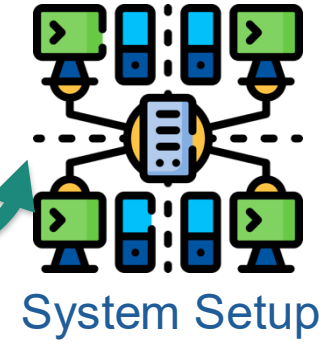


Encode

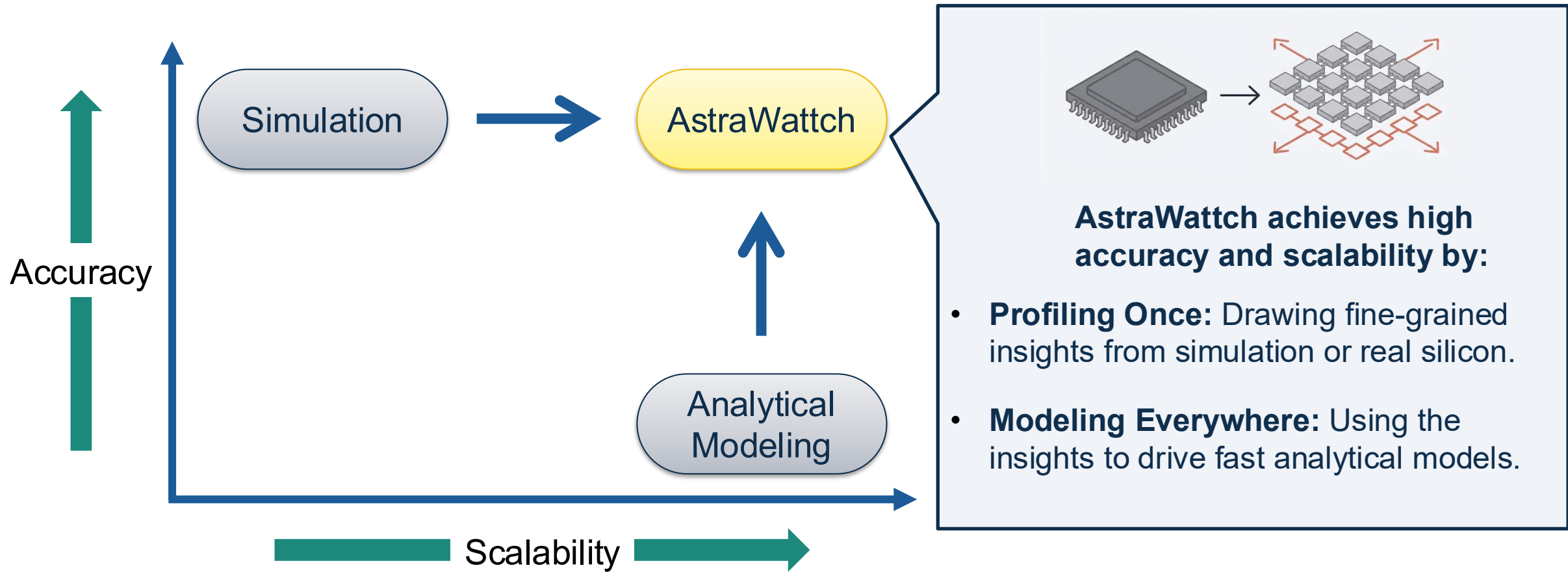


Chakra Trace (Execution Graph)

Workload Specific Parameters
Hardware Specific Parameters



AstraWattch: Profile Once and Model Everywhere

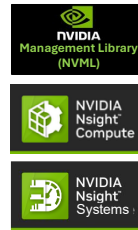
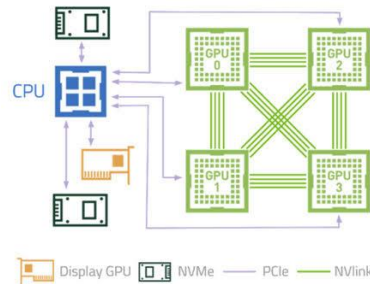


Preliminary Results: AstraWattch enables accurate energy modeling in minutes

Preliminary Results: AstraWattch's Modeling Accuracy



Training



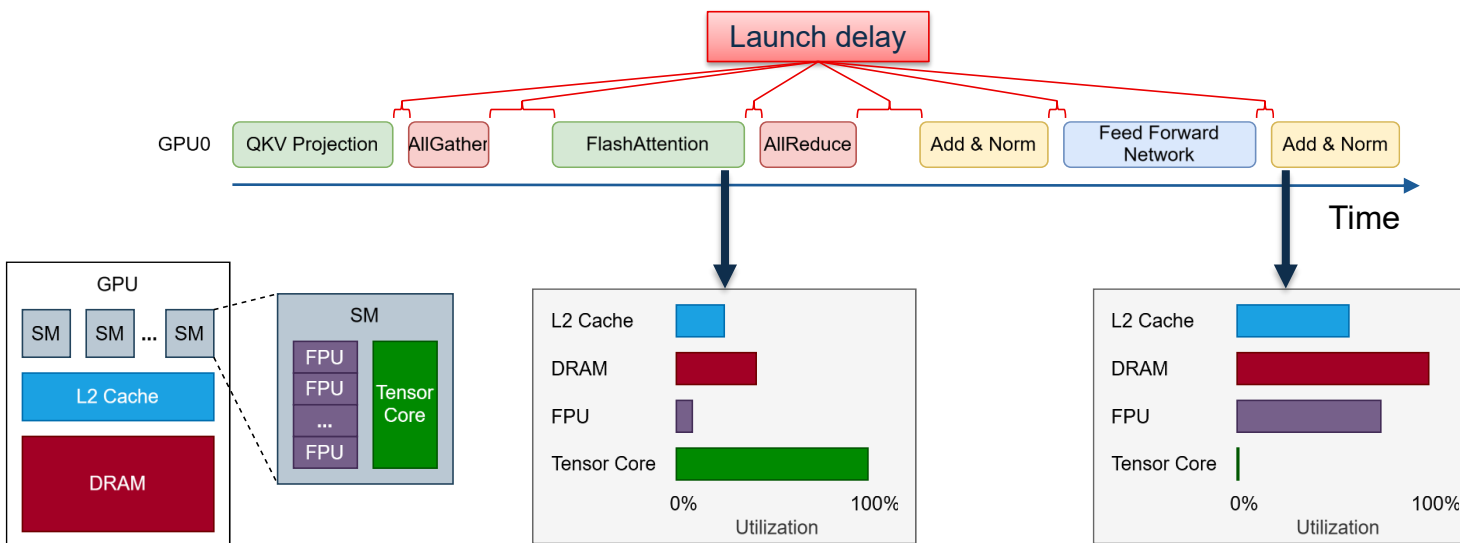
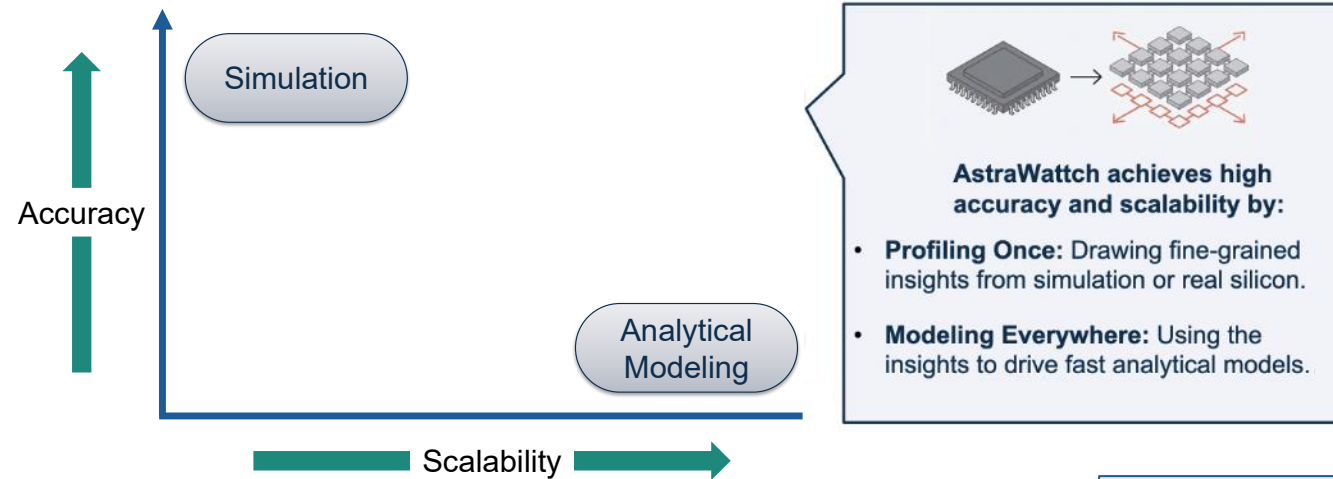
10
Minutes

Frequency (MHz)	Time Modeling Error (%)	Energy Modeling Error (%)
735	9.9	11.8
930	4.6	2.2
1132	2.2	2
1335	2.1	4
1530	7.4	0.9

NVIDIA DGX V100

AstraWattch enables accurate energy modeling in minutes by generalizing analytical models to different frequencies

Enabling Energy Efficiency for Distributed Machine Learning



- Component utilization and delay information → **chakra trace**
- Component power → **simulation configuration**

Tawhid Bhuiyan

Email: mb5332@columbia.edu

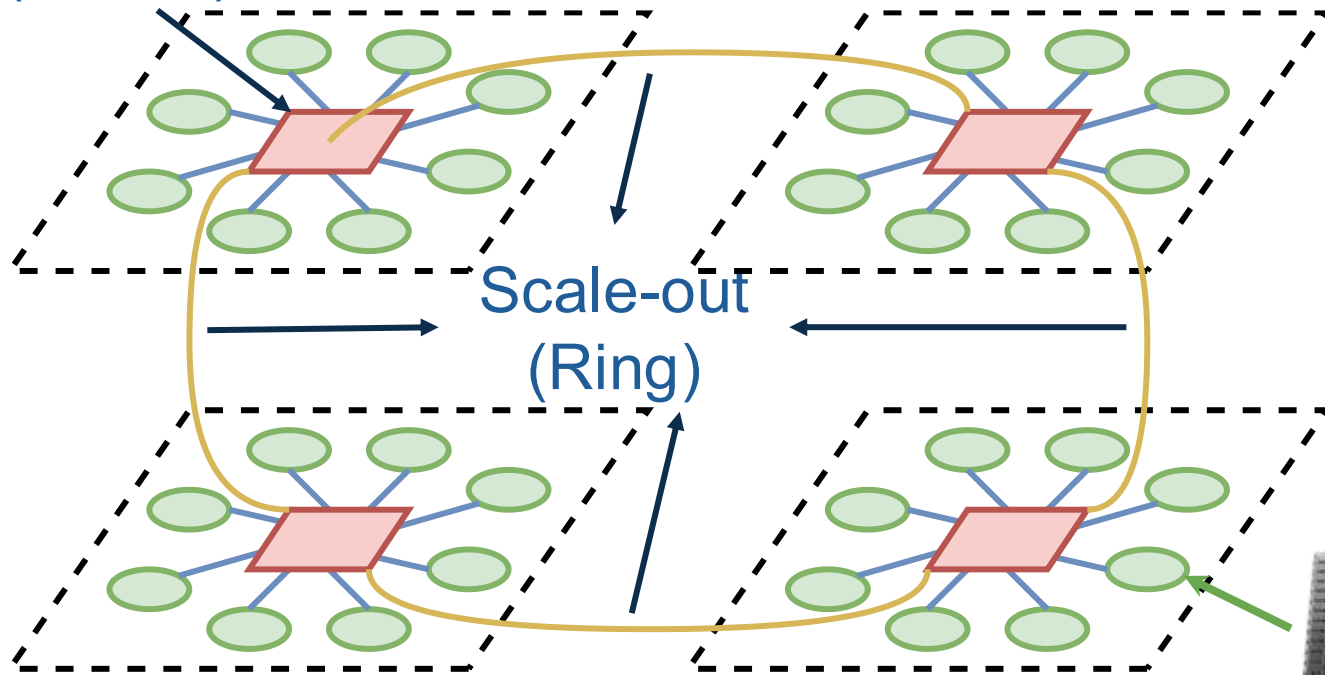
Website: mthbhuiyan.github.io

Preliminary results: AstraWatch enables accurate energy modeling in minutes

Case Study: Marginal Efficiency Analysis

Machine Learning Model and System Configuration

Scale-up (Switch) LLaMA 2-70B Training



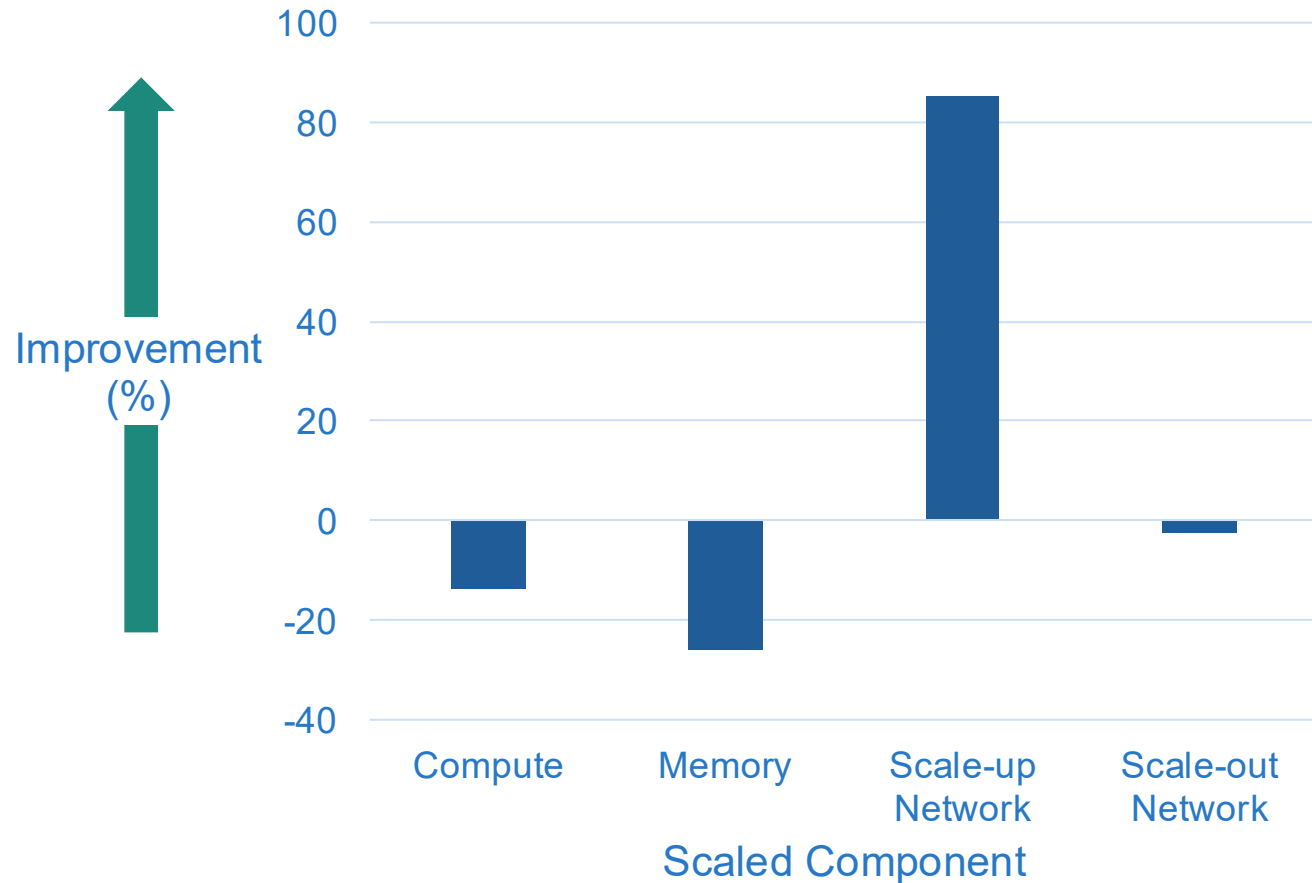
Switch (8) - Ring (4) Topology



NVIDIA H100

Active Power = 700 W
Idle Power = 80 W
Scale-up Network:
Energy Per Bit = 2pJ
Bandwidth = 400GB/s
Scale-out Network:
Energy Per Bit = 25pJ
Bandwidth = 50GB/s

Case Study: Marginal Efficiency Analysis



- Scale-out network voltage–frequency scaling delivers efficiency gain
- Due to
 - Low direct network energy
 - Poor compute-communication overlap
 - High idle energy impact