



THE UNIVERSITY *of* EDINBURGH
informatics

Future Efficient Distributed AI Systems

Luo Mai

Large-scale AI Systems Group

University of Edinburgh



THE UNIVERSITY *of* EDINBURGH
INFORMATICS FORUM

What are the key challenges in AI systems?

Computational requirement of AI computing



Gap to be filled by **Distributed AI Systems**

Examples of distributed AI systems:

- Uber Horovod
- TensorFlow / PyTorch Distributed
- Microsoft DeepSpeed
- DeepMind Acme
- ...



Why should ARM care?

- AI becomes critical workload in **data centres**
 - Gigantic reasoning models: GPT-3, AlphaFold2
 - Neural recommendation services: Facebook DLRM
 - Real-time gaming services: AlphaStar
- Many AI services are deployed **at the edge** and **endpoints**
 - Intelligent edge: traffic engineering, content caching
 - Intelligent endpoints: personal assistant (phones), scene understanding (autonomous vehicles)

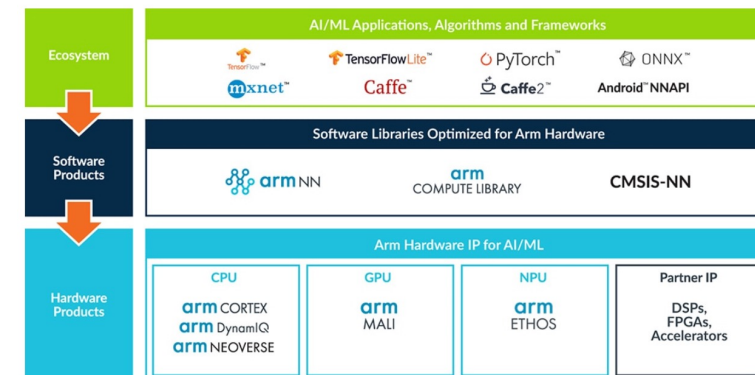
What are the opportunities for ARM?

Existing distributed AI systems exhibits **extremely low energy efficiency**

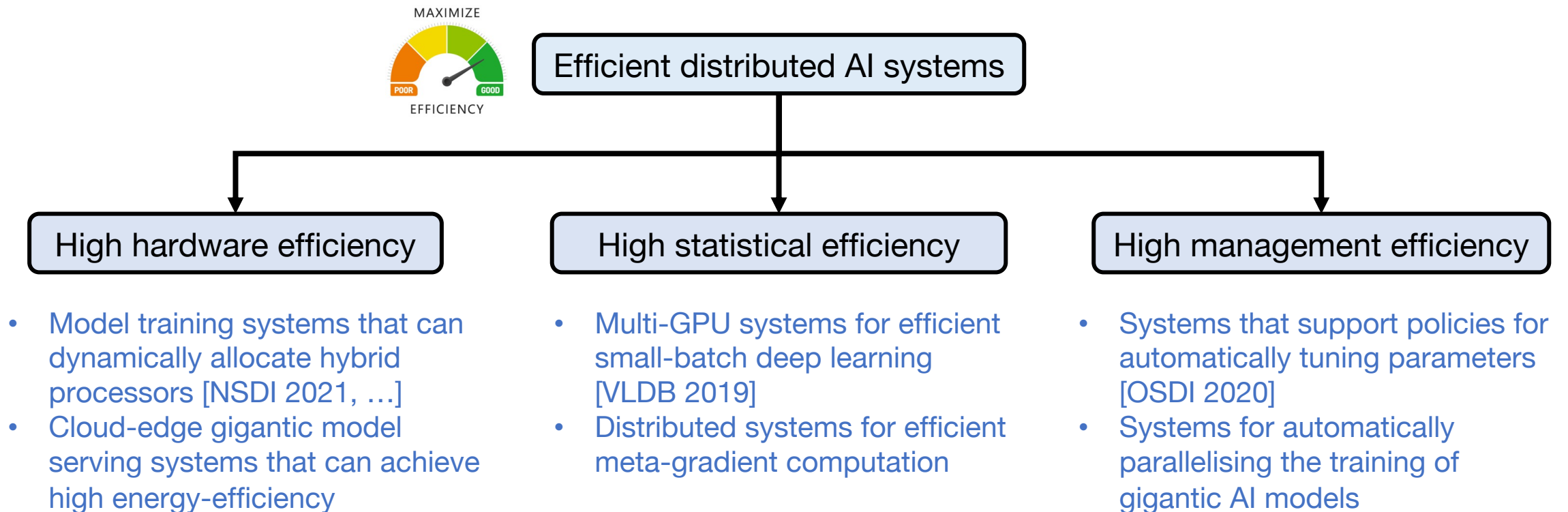
- Example: training a GPT-3 consumes energy equivalent to drive a car from the earth to the moon
- Strong demands for **sustainable AI infrastructure**

Designing **efficient distributed AI systems with ARM technologies**

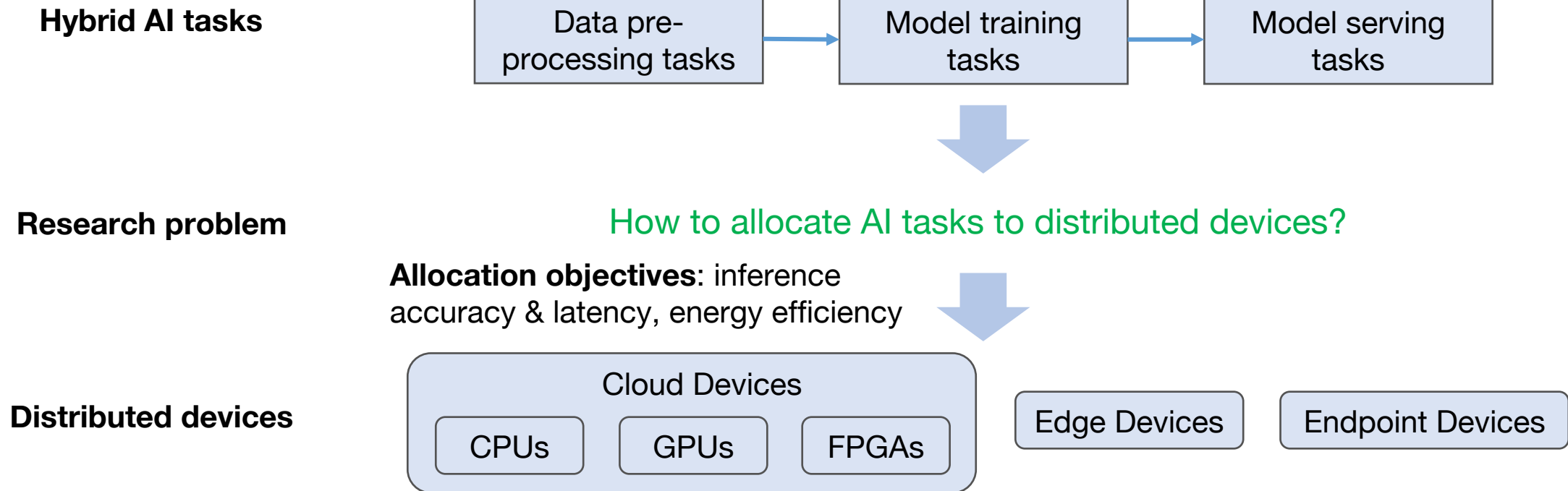
- AI systems optimised for ARM chips:
 - ARM Desktop Chips, ARM Server CPUs, ARM GPUs
- Distributed AI systems optimised for ARM clouds:
 - ARM AI Platform for Machine Learning



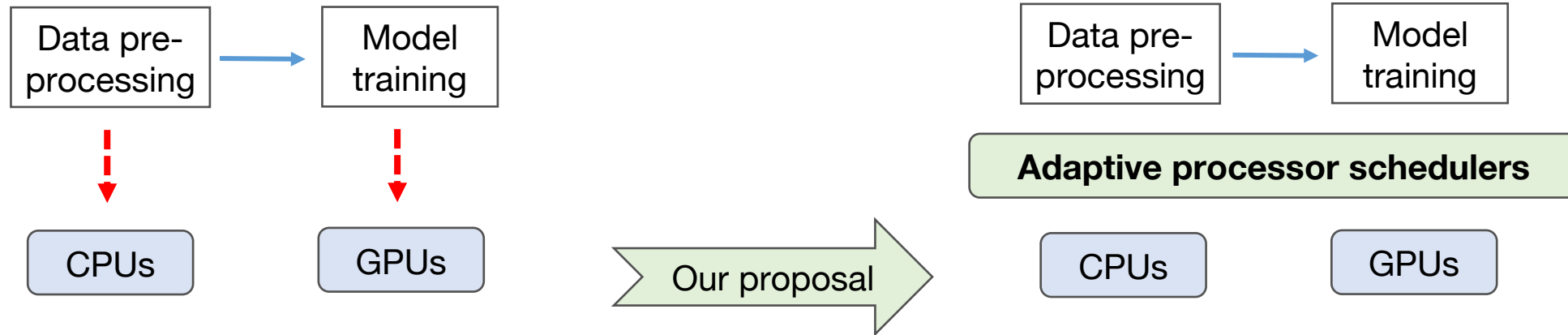
My research towards efficient distributed AI systems



Key problem in optimising hardware efficiency



How to allocate devices for **training** models?



SOTA: Existing AI systems (e.g., TensorFlow, PyTorch) **statically** allocate CPUs for pre-processing and GPUs for training

Problem: Data pre-processing often become bottleneck in emerging AI workloads (e.g., GNNs, RL)

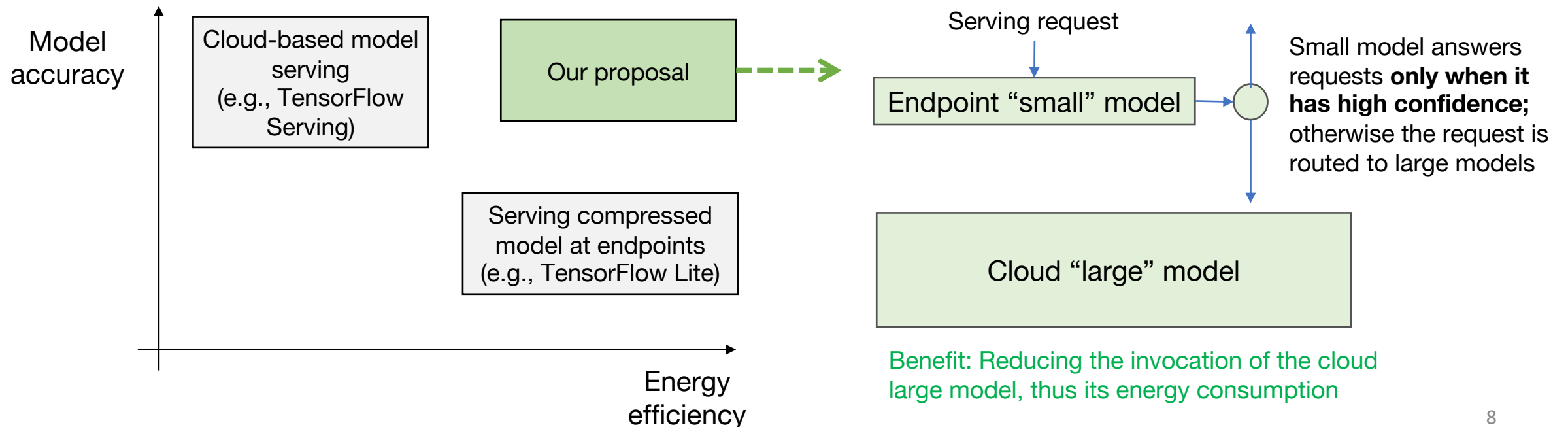
Idea: Designing schedulers that can **dynamically** allocate pre-processing and training tasks to CPUs and GPUs based on **monitored metrics** [NSDI 2021, ...]

Benefits: Up to 10x performance improvement in GNN / Streaming applications

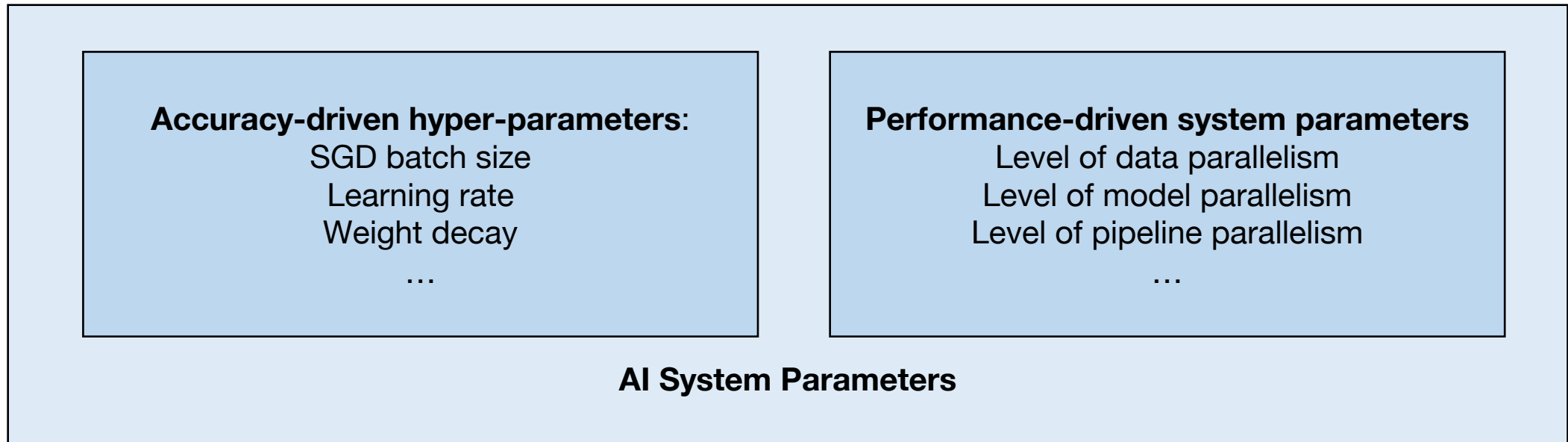
How to allocate devices for **-serving** models?

Problem: Allocating cloud and endpoint devices for serving large AI models (e.g., GPT-3, AlphaFold2)

SOTA: (1) Cloud-based model serving shows **low energy-efficiency**; (2) Compressing models for endpoint deployment **hurts model accuracy**



Key problem in improving management efficiency



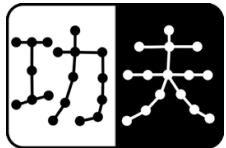
How to find optimal parameters?

How to find optimal hyper-parameters?

SOTA: Hyper-parameters are statically configured according to empirical experience

Problem: Users must frequently re-configure hyper-parameters whenever update models

Our proposal



KungFu Framework
[OSDI 2020]

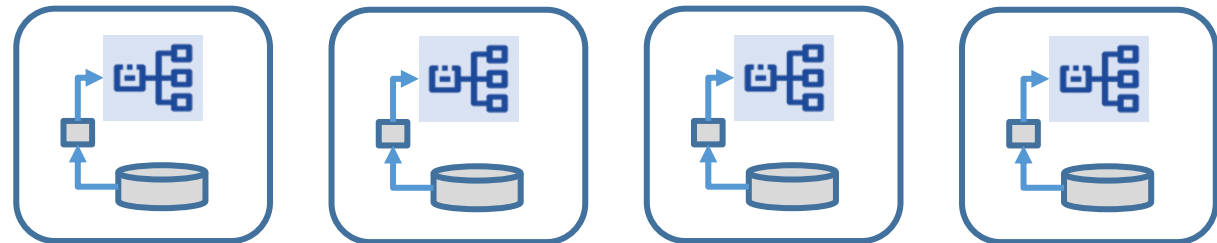
<https://github.com/llds/KungFu>

Benefits:

- Up to 80% improvement in model training time
- Elastic resource usage

High-level parameter adaptation policies

Monitoring training metrics  Adapting hyper-parameters 



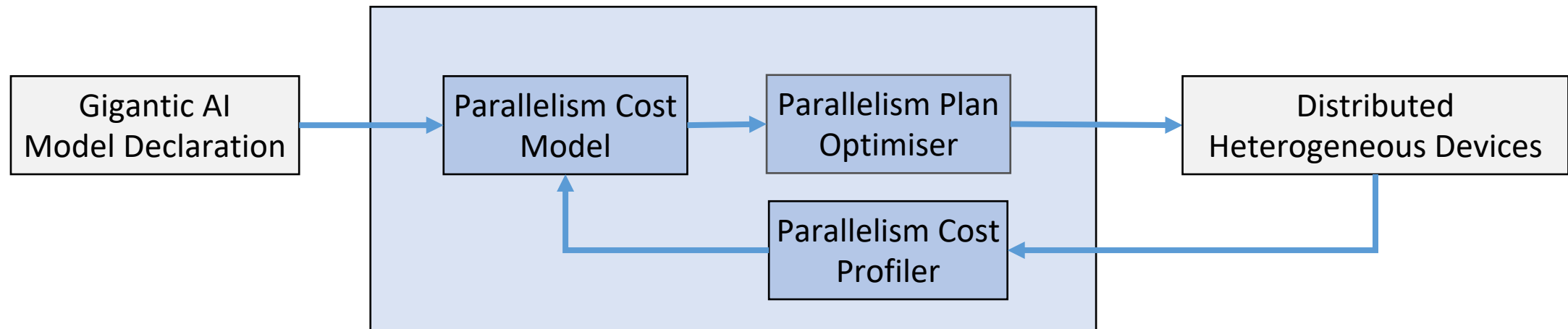
TensorFlow/PyTorch/Keras Workers

Scalable policy execution runtime

How to find optimal system parameters?

SOTA: System parameters (e.g., data parallel, model parallel, pipeline parallel) are hard-coded in system implementation

Problem: Users must frequently re-configure system parameters whenever change hardware or environments



Our proposal: Automatic Parallelism Compiler



Summary

- ARM technologies are keys to design efficient distributed AI systems
- At Edinburgh, we are designing distributed AI systems that can improve
 - Hardware efficiency [NSDI'21, ...]
 - Statistical efficiency [VLDB'19, ...]
 - Management efficiency [OSDI'20, ...]



Large-scale AI Systems Group
University of Edinburgh

Luo Mai
luo.mai@ed.ac.uk
<https://luomai.github.io>