

An Introduction to Modern Deep Learning Systems and Frameworks

Lecture 1 for Advanced Deep Learning Systems

Aaron Zhao, Imperial College London, a.zhao@imperial.ac.uk

Table of contents

1. Introduction
2. Dates and Structures
3. An Introduction to Modern Deep Learning Systems

Introduction

Introduction - Myself

My name is Aaron, and my research looks at the intersections between algorithm, hardware and security in Deep Learning (DL) Systems.

My email is a.zhao@imperial.ac.uk, and my office is 903.

The course is also supported by the following GTAs:

- Cheng Zhang
- Pedro Gimenez
- Mingzhu Shen

Introduction - Where to find stuff?

Everything is online, you simply need a single url

- The course webpage
(<https://aaron-zhao123.github.io/teaching/adls>)
- Bugs should be reported on Github Issues (if it is lab related)
- There are also in-person lab sessions for questions
- Email me (a.zhao@imperial.ac.uk)

Introduction - Why this course?

Learning pure DL theories and model architectures is **superficial**

You need to understand:

- What is actually running **in hardware** when you do "torch.matmul"?
- How do these ML/DL frameworks (eg. Tensorflow, Pytorch) work?
- How to deploy your model on different backends?
 - eg. CPUs, GPUs, etc.
- How to make them run better or faster? What are the **software** and **hardware** tricks for these optimisations?

Hopefully you will have an idea of the answers of these questions.

Introduction - Assessments

The course and the coursework is designed around an in-house compilation framework that we have built, named MASE.

The final score is composed of two parts:

- A **lab oral** with a GTA to show your progress on the four lab practicals (20%).
- A **team project** (2 people) to pick or self-propose a topic, implement and evaluate the idea on MASE (80%).

This course contains both software and hardware aspects, this is a **CHALLENGING** course, and you will be expected to write a lot of code!

The lab orals and assessments for MSc ADIC (more hardware focused) and MSc AML (more software focused) students are different. MEngs, feel free to pick your path.

The purpose of the labs is for you to learn what is involved in a DL framework, we have four lab sessions supported by GTAs:

- Lab 1: Getting familiar with MASE
- Lab 2: Treating optimization as Passes
- Lab 3: Automated Machine Learning a Pass search
- Lab 4 (Software): Network Architecture Search
- Lab 4 (Hardware): Implement an integer matrix multiplication
- Questions about MASE functionalities, apart from asking the GTAs in Lab sessions, can also be queried through Github Issues.

Introduction - Assessments - Lab oral

The labs are intended to assist you in navigating and comprehending a substantial codebase. To ensure your understanding, we will facilitate a one-on-one oral (each 15-min) lab session with you that aims to:

- You finish it individually.
- Verify that you have grasped the material thoroughly.
- Ensure that you are engaging with the content rather than merely replicating commands or code.
- Show us your code **with your comments**.

Introduction - Assessments - Team Project

You will need to find a partner to make up a two-person team. If you cannot find a partner yourself, this is fine, we will find you one.

- A **roadmap** meeting with a GTA (not assessed).
 - We expect to see an **outline of functionalities** that you would like to implement, and **an arrangement of who is doing what**.
 - This is to provide you feedback with your planned work.
- Coursework **submission as a pull request** on Github (assessed).
 - Your code (correctness and quality)
 - Your augmented testing script (coverage)
 - A short Readme file and a few commands to verify it works (clarity)
- A **show me your code** session (assessed).
- Documentation on the functionalities you have implemented and also an **8-page report** per group (assessed).

Introduction - Assessments - Team Project

- We will provide a list of projects, each project would be assigned with a GTA supervisor, each project would be assigned with a GTA supervisor.
- We allow **max 3 groups** on a single project, first come first serve.
- You can also propose your own group project.
- Your **outline of functionalities** should describe clearly who is doing what, this would affect your individual mark.
- If your project is merged to the upstream, we will add your name to our annual tech report.

Introduction - Assessments - Team project

With the support of MASE, you can do many fancy stuff, we encourage you to self-propose a team project.

- MASE has a fairly strict CI, this means when you push on the PR, there is **auto-checking**, you need to pass all the tests, obviously.
- Our GTAs may ask questions on your code in the PR.
- In this way, I expect you to learn how to work on an open-source project.

The framework is completely open-source, and valuable team projects will be integrated by our team in the upstream, so others can see your contributions on Github. If your code is integrated, the report would be included in an annual technical report that we plan to write on the toolchain.

Introduction - My assumptions

I assume you know

- Basic terminal commands (eg. cp, mkdir ...).
- Basic version control knowledge (eg. git, github ...)
- Python and Object Orientated Programming
- System Verilog (for hardware-stream students)

If you never heard of these words, and are not willing to pick them up yourself through Google/Youtube/ChatGPT, think twice before you take this course.

Dates and Structures

Important Dates

- Team project sign-up/proposal starts (29th Jan)
- Team project sign-up ends (9th Feb)
- Midterm lab oral date (TBC, but around 5th Feb)
- Show me your code session (TBC, but around 18th Mar)
- Report submission and pull request finalization (28th Mar)

You have in total 8-hour lectures.

The lectures do not necessarily cover all aspects of MASE. The lectures aim to deliver theoretical knowledge, background and an overview of recent progresses in the field of Deep Learning Systems.

I believe you **learn coding by doing it** (that's why labs exist), and you do good research/engineering by understanding the field and code it up (that's why there are lectures and team project).

The list of team projects would be centered around lectures' contents.

You have in total 10-hour labs.

In the first 2 lab sessions, we will cover Lab1 - Lab4.

The other lab sessions are for you to consult with the GTAs for your team projects.

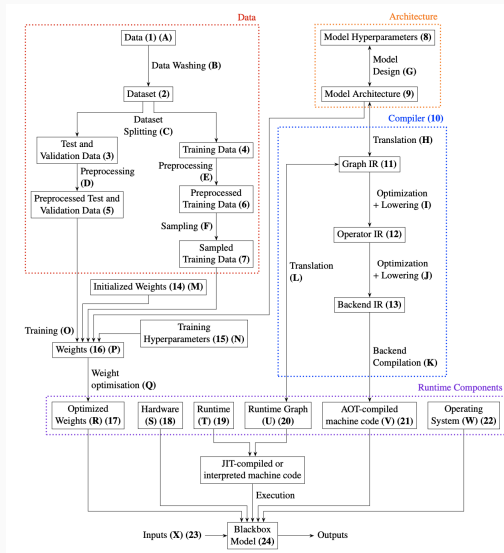
Additional independent study and lab hours are needed.

Dates - Lectures and Labs

- 15th Jan: Lecture 1 and 2
- 22nd Jan: Lab 1
- 29th Jan: Lab 2 (Team project sign-up start)
- 5th Feb: Lecture 2 (Team project sign-up ddl, Mid-term oral)
- 12th Feb: Lecture 3 and 4
- 19th Feb: Lecture 5 and 6
- 26th Feb: Lab 3
- 4th Mar: Lab 4
- 11th Mar: Lab 5
- 18th Mar: Show me your code

An Introduction to Modern Deep Learning Systems

An Overview of an ML System



The Limitations on Hardware Designs

- The end (or slow-down) of Moore's law
- Wire delays
- Power/Cooling is a key limitation (dark silicon)

Trends in ML model designs

It seems like it is going from domain-specific models to general (or foundation) models.

- The rise of foundation models (eg. GPTs, CLIP).
- The idea of handling multi-modality inputs (eg. Google Pathway).

We are also designing systems that adapts to different use scenarios:

- Federated learning, because of privacy concerns (or not).
- On-device learning.
- the list really goes on

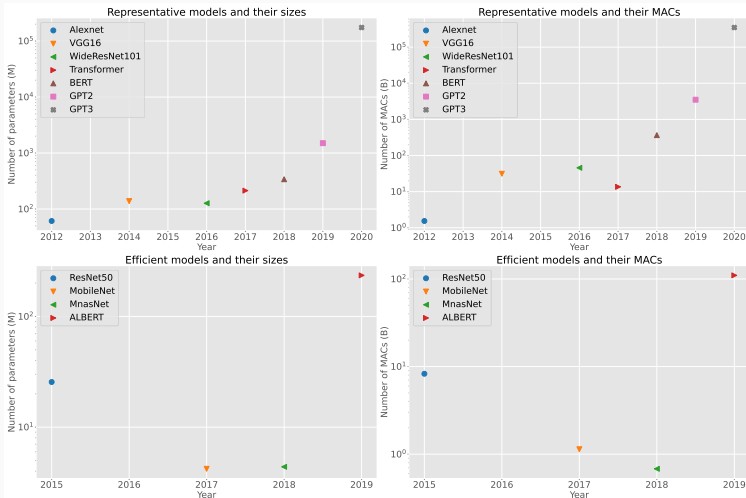
Two key deployment scenarios

We consider mainly two systems: the **cloud system** and the **edge system**:

- Difference in scale
- Difference in budgeting (power, speed, reliability, input rate ...)
- Difference in optimization techniques

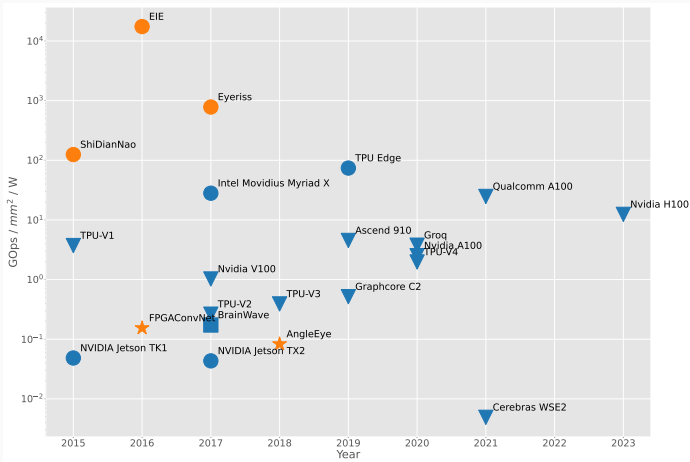
Why ML is a very challenging workload?

- Models are getting larger and larger
- We are developing efficient models, but it requires time



Why ML is a very challenging workload?

- Hardware scaling is a lot slower!
- Basically just Moore's law...



How should we deal with it?

- Software level optimizations (eg. pruning, quantization), typically offers you 2 to 50 \times performance gains.
- Compiler level optimizations (eg. loop tiling, loop unrolling)
 - How can we better map the workload to existing hardware, typically offers another 2 to 5 \times performance gain.
- Hardware level optimizations (eg. PE design, memory system design), typically 2 to 5 \times .

We will explore these optimizations in this course.