*Principles of Machine Learning Systems*

*A: Class Administration*

v1.6

Nicholas D. Lane

# Principles of ML Systems

10 Lectures (FW26) covering:

- ML Systems Landscape
- Mapping to Hardware
- Model Compression
- Accelerators: GPUs, NPUs
- Frameworks and Run-times
- Single/Multi GPU Training

- Scalable Inference Serving
- Deep Learning Compilers
- Automated ML
- Federated Learning
- Development Practices
- MLOps related

# Principles of ML Systems

3 Labs (SW02), covering:

- MCU and Model Compression: Speech Recognition
  - Lab scheduled for Oct 20th  -- Submit material: Nov 2nd

- Single/Multi-GPU Training
  - Lab scheduled for Nov 10th  -- Submit material: Nov 23rd

- Federated Learning: Experiments and Deployment
  - Lab scheduled for Nov 17th -- Submit material: Nov 30th

(each lab counts for 10% of your final grade; total is 30%)
  - Moodle submission. zip file of Google Colab + txt file of Colab URL
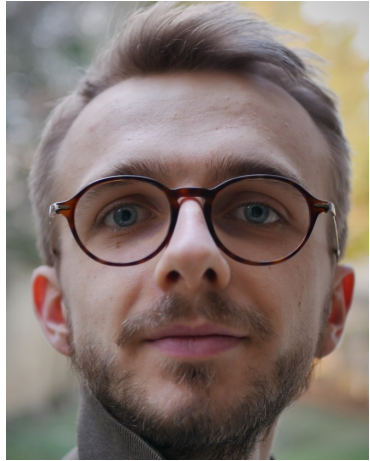
# Principles of ML Systems

- Lectures vs. Labs
- Lectures recorded
  - Available on Moodle (maybe posted on YouTube); subtitles available
  - Past lectures also available (very similar content)
  - **You might be heard/seen in a recording**, especially if you ask questions
  - if uncomfortable with this feel free to approach us via an alt. channel
- TA sessions, and Office Hours
  - TA General+Project (weekly w/ Filip; 11am Wed mix of online/in-person)
  - TA Labs (upon appointment)
  - Office Hours (upon appointment)
- No required textbook: Lectures have optional reading lists
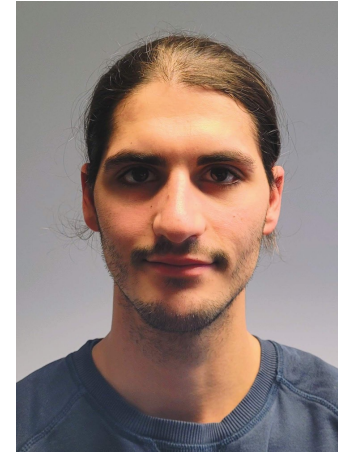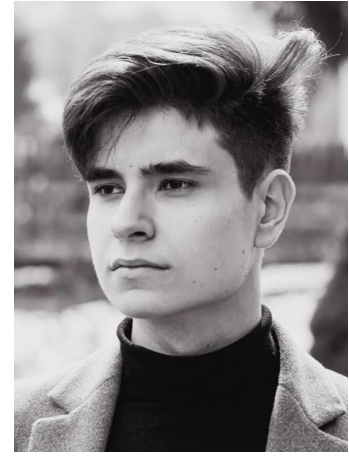
# The Team



Titouan Parcollet

Filip Svoboda

Dr **Thomas** Ploetz
Georgia Institute of Technology, USA
thomas.ploetz@gatech.edu

Dr **Yu** Guan
Hongxiang Fan
Open Lab
Newcastle University, UK
yu.guan@newcstle.ac.uk

Nic Lane
University of Oxford
Nokia Bell Labs, UK
niclane@acm.org

Dr **Sourav** Bhattacharya
Lorenzo Sani
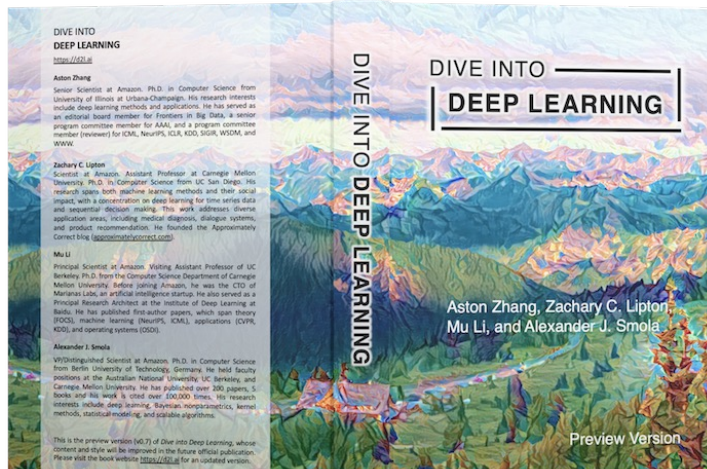Nokia Bell Labs, Cambridge, UK
sourav.bhattacharya@nokia-bell-labs.com

Alex Iacob

# Pre-requisites

1. Comfortable with programming in Python

2. Undergraduate level: operating systems; computer architecture; machine learning

3. Familiarity with Deep Neural Networks (Intro Level)



**Dive into Deep Learning**
https://d2l.ai/

# Assessment: Class Project

*Primary* assessment of the class (70%)

- Individual or Teams (sorry but Part III and MPhil students can't be in the same group)
- Repository for code, decision planning, results and write up
- Repository access provided to assessors *(myself + team)*
- Primary output assessed: written report *(submitted by moodle)*
- Repository will used understand team contributions and process
- Assume your report will be public *(unless otherwise requested)*
- Report in format of NeurIPS conference paper *(8 pages)*
- Due: Start of Lent term *(16/1/24 at 12:00 noon)*

# Assessment: Class Project

Example Projects *(select projects from last year will be posted)*

- Build an application – applying some ML Sys ideas
- Explore a new direction *(novelty is not a must)*
- Leverage interesting hardware or architecture
- Detail an investigation to examine a direction of interest
- Replicating paper results
- Replicating expected system behavior *(textbook etc.)*