# CS249r: Learning on the Edge



### Goals for today

- 1. Course logistics
- 2. Lecture Challenges with on-device training/learning/updates
- 3. Paper discussions
- 4. Guest speaker

# **Course Logistics**

#### Schedule

## Please check the website for the latest and most up to date information



•••	😵 CS2	49r: Tir	ny Machine l	.earnin ×	+															
$\leftarrow \   \rightarrow \   G$	۵	â site	es.google.	com/g.ha	rvard.edu/	cs249-tiny	ml-202	23	Ô	☆			Tex 💘	0	Ŷ	C	Ŭ	+	*	3
TinyML	Harva	ard 🗎	Funding	MLC	🗎 Nora	🗎 LLMs	$\odot$	0	A\	+.	Q	📀 cs	249r Boo	k 🛆	LLMx	۵	CS2	49r		
= 🕲	CS24	9r: Ti	iny Machi	ine Lear	ning															C
SCH	IED	ULE																		
= Wh	nen: E	very	Monday																	
• Tir	ne <sup>,</sup> 13	2.45 P	PM - 3.30	PM EST	-															

Required · Flower: A Friendly Federated Learning Framework (paper) · On-Device Training Under 256KB Memory (paper) Nic Lane, Oct Professor at Optional Learning on the Edge None None 16 University of Cambridge CS · TinvTL: Reduce Memory, Not Parameters for Efficient On-Device Learning (paper) · TinyTrain: Deep Neural Network Training at the Extreme Edge (paper) Required · CFU Playground: Full-Stack Open-Source Framework for Tiny Machine Learning (tinyML) Acceleration on FPGAs (paper) · An Evaluation of Edge TPU Accelerators for Convolutional Neural Networks (paper) Optional Hardware Acceleration Hardik Sharma, · In-Datacenter Performance Oct Assignment 2 for ML: GPUs, TPUs Silicon Architect at Assignment 3 23 Analysis of a Tensor Due and FPGAs Googla

#### **GENERAL INFORMATION**

#### RESOURCES

Where: SEC LL2.229

## Assignments Goals

Four assignments

- Assignment 1
- Assignment 2
- Assignment 3
- Assignment 4

Each assignment focuses on different aspects

Want to connect the dots on why we picked what we picked and how you do it

## Assignment - Why, What, How

- Assignment 1
  - **Why**: Vision (Camera) use cases + Dataset generation
  - What: End-to-end ML workflow
  - **How**: Frameworks (Edge Impulse) + Embedded Hardware (Nicla)
- Assignment 2
  - Why: Audio (Microphone) use cases
  - What: ML Pipeline (Data preprocessing + Model optimizations)
  - **How**: NAS (Edge Impulse) + Quantization/Pruning (TFLite)
- Assignment 3 Model development
  - **Why**: Time series (IMU) use cases
  - What: ML frameworks programming (Tensorflow/TFLite/TFLM)
  - How: Bare metal implementation of end-to-end ML workflow
- Assignment 4 Responsible Al
  - Why: Sustainability + Responsibility
  - What: Sustainability-aware system design and safe AI
  - **How**: TinyML Footprint + Adversarial Nibbler

## Assignment Schedule Updates

- Assignment 2
  - Due: October 23rd (Monday)
- Mid-Project Review
  - Due: October 30th (Monday)
- Assignment 3
  - Due: November 6th (Monday)
- Assignment 4 Part 1
  - Due: November 20th (Monday)
- Assignment 4 Part 2
  - Due: November 27th (Monday)
- Project Presentations
  - Due: December 4th (Monday)
- Final Report
  - Due: December 11th (Monday)

## **Project Logistics**

Deadlines

- Project Proposal Oct 10th
- Mid-Project Presentation (5 slides/3 mins) Oct 30th
- Final Project Presentation Dec 4th
- Final Project Report Dec 11th

Reviewed the projects

- More application centric, than TinyML engineering centric that's OK!
- Generally, groups of 1 ~ 3/4

## **Project Topics**

- 1. Cache 4 Trash
- 2. Coloring our Thoughts
- 3. Context Dependent Notification System
- 4. Decoding Sleep States through Neural Models in Wearable Devices
- 5. Does regularization during training improve performance of quantized models?
- 6. Early Wildfire Mapping via TFLite and Image Detection
- 7. Hand-Tracking for BrilliantLabs Monocle
- 8. Improving waste disposal practices through smart trash classification
- 9. Quant-Eyes: Fast and Efficient Quantized Models for On-Device Real-Time Eye Tracking
- 10. Quantization and Fine-Tuning in IMU-based Gesture Recognition
- 11. Real-time Visual Narration with tinyCLIP for the Visually Impaired
- 12. Saving power in ultrasound sensing for tracking neuromuscular function
- 13. WickWatcher

#### Waste Management and Environmental Sustainability

- Cache 4 Trash
- Improving waste disposal practices through smart trash classification

#### **Visualization and Digital Art**

- Coloring our Thoughts
- WickWatcher

#### **Notifications and Context-aware Systems**

Context Dependent Notification System

#### Health and Biometric Monitoring

- Decoding Sleep States through Neural Models in Wearable Devices
- Saving power in ultrasound sensing for tracking neuromuscular function
- Hand-Tracking for BrilliantLabs Monocle

#### **Quantization and Model Efficiency**

- Does regularization during training improve performance of quantized models
- Quant-Eyes: Fast and Efficient Quantized Models for On-Device Real-Time Eye Tracking:
- Quantization and Fine-Tuning in IMU-based Gesture Recognition

#### Image Recognition and Processing

- Early Wildfire Mapping via TFLite and Image Detection
- Real-time Visual Narration with tinyCLIP for the Visually Impaired

## Projects

- Create a project channel
- Add all staff to the channel
- Update description of the channel w/ project summary

#### Projects

- # project-decoding-sleep-states
- # project-quanteyes
- 合 project-regularization-quantiza...
- # project-tinyclip
- 合 project-trashy-ml
- 合 project-ultrasound-power
- 合 project-wildfiredetection

## Scribing Updates

#### Chapters

- Data engineering 🗸
- ML frameworks 🗸
- Model optimizations
- On-device learning (today!)
  - \circ Come say hi! 👋

**New Contributors** 

- Itai, Oishi, Shreya (Data Engineering) ......
- Marcelo NICLA

••• • 📀 M	ACHINE LEARNING S	YSTEM ×	+																~
$\leftrightarrow$ $\rightarrow$ $C$ $\Delta$	â harvard-edg	e.github.io	/cs249r_b	ook/contri	ibutors.I	h 🖞	☆	M		718 🐨	0	Ļ	٢	M	•	* 1		2	:
🗎 TinyML 🗎 Har	vard 🗎 Funding	MLC	🗎 Nora	🗎 LLMs	0	A\	+.	Q	📀 cs	249r Book	۵	LLMx	۵	CS24	9r				39
MACHINE LEA	ARNING SYS	TEMS											C	00.	↓. •	≪° ▼	50	С	<
																			0

FRONT MATTER Preface Dedication Acknowledgements Contributors Copyright About the Book MAIN 1 Introduction 2 Embedded Systems 3 Deep Learning Primer 4 Embedded Al 5 AI Workflow 6 Data Engineering 7 AI Frameworks 8 Al Training 9 Efficient Al 10 Model Optimizations 11 Al Acceleration 12 Benchmarking Al 13 On-Device Learning 14 Embedded AlOps 15 Privacy and Security 16 Responsible AI 17 Generative Al 18 Al for Good 19 Sustainable Al EXERCISES Setup Nicla Vision CV on Nicla Vision References

Appendices A Tools B Datasets

C Model Zoo D Resources E Communities F Case Studies

#### Contributors

Jessica.Ouz

ш

We extend our sincere thanks to the diverse group of individuals who have generously contributed their expertise, insights, and time to enhance both the content and codebase of this project. Below you will find a list of all contributors. If you would like to contribute to this project, please see our GitHub 🖄 page.







P

## **Scribing Updates**

General flow for how things are working  $\rightarrow$ 

- Collaboratively generate content
- Peer review + feedback + comments
- Resolve comments w/ staff
- Staff merge update

Comments on Implementation

- Put links to figures in Google Doc
- Make sure chapter renders correctly on quarto
- Check figure caption format is the same as other book chapters
- Make sure references render correctly and update refs.bib file—do not have a separate references section at the end of your chapter



## **Course Topics**

- 1. Overview and Introduction to Embedded Machine Learning
- 2. Data Engineering
- 3. Embedded Machine Learning Frameworks
- 4. Efficient Model Representation and Compression
- 5. Performance Metrics and Benchmarking of ML Systems
- 6. Learning on the Edge
- 7. Hardware Acceleration for Edge ML: GPUs, TPUs and FPGAs
- 8. Embedded MLOps
- 9. Secure and Privacy-Preserving On-Device ML
- 10. Responsible Al
- 11. Sustainability at the Edge
- 12. Generative AI at the Edge

#### **Attendance Sheet**

To be filled in during class.

Will be passed around each class.











#### **GBoard Example**



### **GBoard Example**



#### How do we realize GBoard with MLaaS



#### How do we realize GBoard with MLaaS



## **Traditional** ML

• Data is **aggregated** from different sources at the server



## **Traditional** ML

- Data is **aggregated** from different sources at the server
- **Central server** builds the machine learning model



## **Traditional** ML

- Data is **aggregated** from different sources at the server
- **Central server** builds the machine learning model
- Central server **distributes the** global model to everyone



 Send all of the raw clients' user data to the server



- Send all of the raw clients' user data to the server
- All the ML model training is done in the remote cloud datacenters



• Key concern with sending raw data to the server: **Privacy** 



- Key concern with sending raw data to the server: **Privacy**
- Exposes user's raw data to the central server, which may potentially be compromised risking the loss of private data



## How can **privacy be preserved**?

#### Minimize

• Avoid collecting unnecessary data, and dispose or delete data periodically

#### Protect

• Use encryption techniques to protect data

#### Map the flow of information

• Context, the type of information, and who has access

#### Informed consent

• Be transparent with users about how their data is being collected and used

#### Cloud

#### Embedded



#### Cloud

Embedded

ML-dedicated hardware: CPU, GPU, TPU ML-dedicated software: many tools ML Tasks → Data collection and preprocessing, data transformation, model training, model deployment, inference



#### Cloud

Embedded

ML-dedicated hardware: CPU, GPU, TPU ML-dedicated software: many tools ML Tasks → Data collection and preprocessing, data transformation, model training, model deployment, inference

No ML-dedicated Hardware No ML-dedicated software ML Tasks  $\rightarrow$  Inference

## Alternative approach?

- The key flaw behind the prior approach to training is the sending of **raw client data** to central server
  - Server has access to raw client data, exposing clients to intrusion of privacy by central server






#### Partially trained models $\rightarrow$ server





Partially trained models → server



Server combines and makes federated model





# Federated ML

• Data is kept local to the endpoint device (data does not ever leave the device)





# Federated ML

- Data is kept local to the endpoint device (data does not ever leave the device)
- Only local model updates are set to the central server



# Federated Learning

• A method to **collaboratively** learn a shared model while keeping data on device



# Federated ML

- Data is kept local to the endpoint device (data does not ever leave the device)
- Only local model updates are set to the central server
- Server creates a **global model** and sends it to the endpoints





**Hyper-Personalized** 



Local



**Hyper-Personalized** 



Minimum Latencies



Local



**Hyper-Personalized** 



Low Cloud Infra Overheads



Minimum Latencies





Global Model ---Loca Model Local

Local

Unbalanced training samples

Need a high # of endpoint devices/clients

Unbalanced training samples

Need a high # of endpoint devices/clients

Unbalanced training samples

Need a high # of endpoint devices/clients

Unbalanced training samples

Need a high # of endpoint devices/clients

## **Federated Learning**

**Development Environment** 

**Production Environment** 









## Computation vs. Communication

- Data movement is extremely important to keep in mind for efficient system engineering
- Data movement is more costly than computation itself





# Continuous Monitoring for TinyML

- Monitoring may **not always** be a **feasible** option
  - Low power communication protocol
  - Device isn't wifi-enabled
- Monitoring opens up security and privacy risks

# Continuous Monitoring for TinyML

- Monitoring may **not always** be a **feasible** option
  - Low power communication protocol
  - Device isn't wifi-enabled
- Monitoring opens up **security and privacy risks**
- How can we enable **Continuous Monitoring** to enable **Continuous Training** without moving the data off the endpoint tiny ML device?









	LAN Bluetooth, ZigBee, WiFi	Cellular 2G, 3G, 4G	LPWAN LoRa, Sigfox, NB-loT	
Data Rate	~100kbps - 100mbps	~100kbps - 100mbps	10kbps	
Range	Short	Long	Long Range (10 miles)	
Battery Life	Varies	Medium	Long Battery Life (10 yrs)	
Cost	Expensive	Very Expensive	Best Price	
Use Cases	Smart TV, WiFi Network, Bluetooth Speakers	Smart Grid, CCTV, Personal Communication (smartphones)	Monitoring, Metering, Temperature, Asset Tracking, Weather, Location	

	ка	In r	e rej
 			90

Deep indoor coverage (multi-floor buildings)

Star topology network design

### Long Battery Life

Low-power optimized

Up to 10 year lifetime

Up to 10x versus Cellular M2M

#### **High Capacity**

Millions of messages per gateway

Multi-tenant interoperability

Public/Private network deployments

#### Low Cost

Minimal infrastructure

Low cost end-node

Open source software

### Geolocation

Indoor/outdoor

Accurate without GPS

No battery life impact

#### Updates

Firmware updates over-the-air for apps and LoRaWAN stack

#### Roaming

Seamless handoffs from one network to another

#### Security

Embedded end-to-end AES-128 encryption

Unique ID

Application

Network

• Vending machines could alert distributors when a product requires **maintenance** 



- Vending machines could alert distributors when a product requires **maintenance**
- Animal lovers can **track** pets or study migration patterns over longer distances



- Vending machines could alert distributors when a product requires **maintenance**
- Animal lovers can **track** pets or study migration patterns over longer distances
- Oil companies could **receive alerts** when home oil tanks are running low



- Vending machines could alert distributors when a product requires **maintenance**
- Animal lovers can **track** pets or study migration patterns over longer distances
- Oil companies could **receive alerts** when home oil tanks are running low
- Logistics providers could **track** cargo containers on trucks, ships and trains



- Natural **Disaster Prevention**
- Smart Agriculture Monitoring
- Animal Production Monitoring
- Endangered Species **Protection**
- Smart Industry Control
- Smart Cities, Homes, Buildings & Offices
- Supply Chain Logistics, Asset Tracking & Quality Management
- Smart Metering Facilities (Water, Gas, Electricity)





**AES Secured Payload**


Network Communication Latency

Battery Life



Network Communication Latency

Battery Life



Network Communication Latency











# Network Advantages for Class A Devices

- No complicated network planning is required. Gateways can be added anywhere at any time
- Accurate message delivery is robust, since multiple gateways receive the same data packet during each uplink. This is called uplink spatial diversity
- There is **no need to plan for different frequencies for each gateway**, or to reallocate frequencies when the number of gateways change. All gateways are constantly listening to all frequencies of the network
- Mobile devices can operate at ultra low power thanks to the fact that any gateway can receive messages from any device at any time



### **Cross Comparison**

- Class A device communication is initiated only by the end device
- Class B devices have regularly-scheduled windows, in addition to those that open when a Class A-style uplink is sent to server
- **Class C** devices achieve the lowest latency among all operating modes







### Paper Discussions

### **Discussion Questions**

Discuss

- Why did we read the two papers?
- What are the differences between them?
- <lead discussion>

Summary

• Discuss as a group

## Guest Speaker

#### Nicholas Lane

Nic is a professor of machine learning systems at Cambridge University. His research focuses on the design, architecture and algorithms of scalable and robust end-to-end machine learning systems. He is also the co-founder and Chief Scientific Officer of Flower Labs, a venture-backed AI company behind the Flower federated learning framework.



Personal Website

#### **Google Scholar**