

Widening Access to Applied Machine Learning with TinyML

Vijay Janapa Reddi* Brian Plancher* Susan Kennedy* Laurence Moroney†
 Pete Warden† Anant Agarwal*‡ Colby Banbury* Massimo Banzi§ Matthew Bennett*
 Benjamin Brown* Sharad Chitlangia¶ Radhika Ghosal* Sarah Grafman* Rupert Jaeger||
 Srivatsan Krishnan* Maximilian Lam* Daniel Leiker|| Cara Mann* Mark Mazumder*
 Dominic Pajak§ Dhilan Ramaprasad* J. Evan Smith* Matthew Stewart* Dustin Tingley*

*Harvard University

†Google

ABSTRACT

Broadening access to both computational and educational resources is critical to diffusing machine-learning (ML) innovation. However, today, most ML resources and experts are siloed in a few countries and organizations. In this paper, we describe our pedagogical approach to increasing access to applied ML through a massive open online course (MOOC) on Tiny Machine Learning (TinyML). We suggest that TinyML, ML on resource-constrained embedded devices, is an attractive means to widen access because TinyML both leverages low-cost and globally accessible hardware, and encourages the development of complete, self-contained applications, from data collection to deployment. To this end, a collaboration between academia (Harvard University) and industry (Google) produced a four-part MOOC that provides application-oriented instruction on how to develop solutions using TinyML. The series is openly available on the edX MOOC platform, has no prerequisites beyond basic programming, and is designed for learners from a global variety of backgrounds. It introduces pupils to real-world applications, ML algorithms, data-set engineering, and the ethical considerations of these technologies via hands-on programming and deployment of TinyML applications in both the cloud and their own microcontrollers. To facilitate continued learning, community building, and collaboration beyond the courses, we launched a standalone website, a forum, a chat, and an optional course-project competition. We also released the course materials publicly, hoping they will inspire the next generation of ML practitioners and educators and further broaden access to cutting-edge ML technologies.

1 INTRODUCTION

The past two decades have seen machine learning (ML) progress dramatically from a purely academic discipline to a widespread commercial technology that serves a range of sectors. ML allows developers to improve business processes and human productivity through data-driven automation. Given applied ML's ubiquity and success, its commercial use should only increase. Existing ML applications cover a wide spectrum that includes digital assistants [1, 2],

§Arduino, ¶BITS Pilani, work done as Harvard intern, ||CreativeClass.ai, *edX, ‡MIT

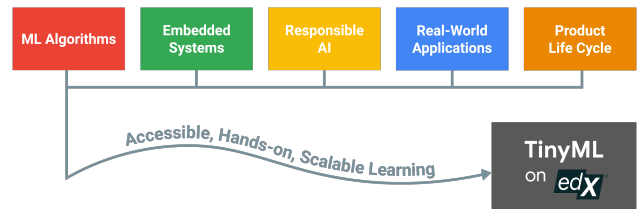


Figure 1: We designed a new applied-ML course motivated by real-world applications, covering not only the software (algorithms) and hardware (embedded systems) but also the product life cycle and responsible AI considerations needed to deploy these applications. To make it globally accessible and scalable, we focused on the emerging TinyML domain and released the course as a MOOC on edX.

autonomous vehicles [3, 4], robotics [5], health care [6], transportation [7, 8], security [9], and education [10, 11], with new application use cases continuously emerging every few days.

The proliferation of this technology and associated jobs have great potential to improve society and uncover new opportunities for technological innovation, societal prosperity, and individual growth. But it all rests on the assumption that everyone, globally, has unfettered access to ML technologies, which isn't the case.

Expanding access to applied ML faces three challenges. First is a shortage of ML educators at all levels [12, 13]. Second is insufficient resources, as training and running ML models often requires costly, high-performance hardware, especially as data sets continue to balloon. Third is a growing gap between industry and academia, as even the best academic institutions and research labs struggle to keep pace with industry's rapid progress. Addressing these critical issues requires innovative education and workforce training to prepare the next generation of applied-ML engineers.

This paper presents a pedagogical approach, developed as an academic and industry collaboration led by Harvard University and Google, to address these challenges and thereby increase global access to applied ML. The resulting course, TinyML on edX, focuses not only on teaching the topic by exploring real-world TinyML

applications running on low-cost embedded systems, but it also considers the ethical and life-cycle challenges of industrial product development and deployment (see Figure 1).

To improve accessibility we employ both cloud computing and low-cost hardware. We leverage Google’s free, open-source TensorFlow and Colaboratory tools along with globally accessible inexpensive embedded devices from Arm and Arduino. We believe hands-on learning that transcends the underlying ML equations is essential. To this end, we focus our approach on TinyML.

Tiny Machine Learning (TinyML), a rapidly growing subfield of applied ML, is a prime candidate for enabling hands-on education globally. This budding area focuses on deploying simple yet powerful models on extremely low-power, low-cost microcontrollers at the network edge. TinyML models require relatively small amounts of data, and their training can employ simple procedures. Furthermore, as TinyML can run on microcontroller development boards with extensive hardware abstraction, such as Arduino products, deploying an application onto hardware is easy. TinyML enables a variety of always-on applications for battery-powered devices—for instance, environmental monitoring and industrial predictive-maintenance analytics. Moreover, the same cost and efficiency benefits open the door to distributed TinyML systems working in concert at the “edge” of the cloud-computing network.

Since TinyML systems are becoming powerful enough for many commercial tasks, learners can acquire skills that can directly apply to their professional careers and future job prospects. The lessons of complete-TinyML-application design, development, deployment and management are also transferable to large-scale ML systems and applications, such as those in data centers and mobile devices. This technology thus provides an attractive entry into applied ML.

Our approach to applied ML through the lens of TinyML provides experience with the complete industrial ML workflow, and it also explores the ethics of software deployment—crucial knowledge for the applied-ML workforce. Creating an ML system is a high-stakes endeavor since inaccurate or unpredictable model performance can erode consumer trust and reduce the chance of success. So understanding ethical reasoning is a crucial skill for ML engineers. To this end, we collaborated with the Harvard Embedded EthICS program to develop and integrate a responsible-AI curriculum into each course, providing opportunities to practice identifying ethical challenges and thinking through potential solutions to concrete problems, many of which are based on real-world case studies.

To broaden and widen access, we aimed to provide TinyML on a globally available platform that lets users benefit at no cost from instructional resources. We therefore deployed our pedagogical approach on edX, a MOOC provider created by Harvard and MIT that hosts university-level courses in many disciplines. Notably, professionals can choose to prove their newly earned skills with a certificate, available for a small fee, once they satisfy the testing requirements. To foster collaboration and continued learning beyond this edX course, we developed a standalone website, a Discourse forum, a Discord chat, and an optional course-project competition.

We launched the core TinyML edX series, comprising three sequential courses, between October 2020 and March 2021; an optional fourth course is under development. On average, more than

1,000 new students enroll each week. After eight months, over 40,000 have enrolled from 170 countries. They come from diverse backgrounds and experiences, ranging from complete novices to experts who want to master an emerging field. Feedback suggests this strong enrollment may owe to the unique collaborative structure we foster between students, teachers, and industry leaders. Shared ownership between Harvard faculty and staff and Google instructors and engineers appears to give participants confidence they are gaining skills that industry needs both today and tomorrow. Moreover, we recognize that opportunities to interact with experts is both encouraging and validating.

In summary, our effort to expand access to and participation in applied ML reflects five guiding principles:

- (1) Focus on application-based pedagogy that covers all ML aspects. Instead of isolated theory and ML-model training, show how to physically design, develop, deploy, and manage trained ML models.
- (2) Work with industry and academic leaders to aid learners in developing the skills that industry requires today and will require in the foreseeable future.
- (3) Raise awareness of the ethical challenges associated with ML and familiarize learners with ethical-reasoning skills to identify and address these challenges.
- (4) Prioritize open access to students worldwide by teaching TinyML at a global scale through a MOOC platform using low-cost hardware that is available anywhere.
- (5) Build community by providing a variety of platforms so participants can learn collaboratively and showcase their work no matter where they live.

We hope the approach we devised brings ML to more people. As such we have open sourced our courseware materials which can be found at <https://github.com/tinyMLx/courseware>, and note that this paper is part of a broader effort to enable activities such as TinyML4 Developing Countries (TinyML4D) and TinyML4 Science, Technology, Engineering, and Mathematics (TinyML4STEM).

We have organized the rest of the paper as follows. It begins with a discussion of the criteria for increasing access to applied ML (Section 2). We then explain both why TinyML is a useful entry to practical ML and how our courses meet those criteria (Section 3). Next, the paper describes our series (Section 4) and how we integrated ethics throughout (Section 5). We then detail how we quickly and efficiently deployed TinyML by innovating in both multimedia production and the use of MOOCs as well as other online platforms (Section 6), in addition to analyzing early data on our courses’ impact (Section 7). Finally, we introduce the TinyML Open Education Initiative: our effort to further broaden the courses’ impact through activities such as TinyML4STEM (Section 8). To provide a balanced viewpoint, we discuss some limitations of our approach and suggest alternatives (Section 9). Finally, we conclude the paper with the main takeaways and the lessons learned (Section 10).

2 CHALLENGES AND OPPORTUNITIES

We propose three criteria to empower applied-ML practitioners. First, no one size fits all with regard to interest, experience, and motivation, especially when broadening participation. Second, given ML innovation’s breakneck pace, academic/industrial collaboration on cutting-edge technologies is paramount. Third, learners who wish to prepare for ML careers need experience with the entire development process from data collection to deployment, and they must understand the ethical implications of their designs before deploying them.

2.1 Student Background Diversity

A major challenge in expanding ML access is that participants begin applied-ML courses with diverse background knowledge, work experience, learning styles, and goals. Hence, we must provide multiple on-ramps to meet the needs of a varied population.

Participants include high-school and university students who want to learn about AI for the first time. Not only will this knowledge empower them to develop cutting-edge applications, but it will also give them an edge in their careers, as many employers expect new hires to have some ML background.

Other participants are industry veterans looking to either pivot their careers toward ML or study the landscape of the TinyML field. For example, some are computer-systems engineers who want to learn about ML in general. Others are ML engineers and data scientists who want to expand their skills by applying ML. Yet others are doctors, scientists, or environmentalists who are curious about how TinyML technology could transform their fields.

Other participants are self-taught, makers, tinkerers, and hobbyists who want to build smart “things” based on emerging technologies. This group typically operates at the systems level, drawing on prior art, but they want to understand how different components or functional blocks fit together to create intelligent ML devices.

Given this broad spectrum, we have a unique opportunity to enable inclusive learning for all despite differing backgrounds and expertise. But we must provide multiple on-ramps. Specifically, we chose to structure the course in a spiral that sequentially addresses the same concepts with increasing complexity [14]. Doing so ensures that not only do participants reinforce fundamentals while picking up new details, but they also master important objectives at every stage. This approach has been shown to improve learning while meeting each individual’s objectives [15].

2.2 Need for Academia/Industry Collaboration

Expanding ML access requires the expertise of academia and industry. Academia is strong in structured teaching: it creates in-depth, rigorous curricula to impart a deep understanding of a field. Conversely, industry is more pragmatic, developing the skills necessary for employment. These approaches are complementary.

Also, ML is moving rapidly thanks largely to industry’s access to rich data. Analysis of ML-research publications at the NeurIPS ML conference suggests industry leads ML innovation [16]. As such, industry has essential domain-specific knowledge that helps ground ML pedagogy in practical skills and real-world applications.

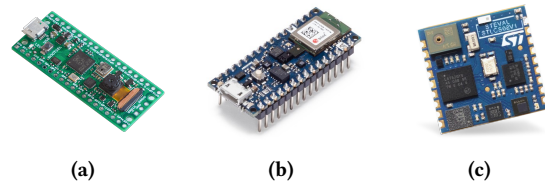


Figure 2: Example TinyML devices: (a) Pico4ML, (b) Arduino Nano 33 BLE Sense, and (c) STMicroelectronics Sensor Tile.

We believe that academia and industry must work in tandem to deliver high-quality, accessible, foundational, and skills-based ML content. Joining a strong academic institution and a industry leader in technology innovation, with a history of releasing free and accessible resources, makes students confident that they are learning the best skills from the best teachers.

2.3 Demand for Full-Stack ML Expertise

In ML, the “full stack”¹ approach to building and using ML models is the core skill that will define future engineers. The engineers who bring long-term value to their industry are those who have the in-depth knowledge to innovate beyond well-known applications and scenarios. In fact, full-stack developers are now more numerous than all other developers combined, with 55% of developers identifying as full-stack in a 2020 report [17].

Our academia and industry collaboration can ensure the course series imparts the full-stack abilities that industry demands. Doing so requires content beyond the narrow, well-lit path of ML-model training, optimization, and inference. We therefore also focus on acquiring and cleansing data, deploying models in hardware, and managing continuous model updates on the basis of field results. Our hope is that learners will gain a whole new set of applied-ML skills and unlock new ideas.

3 ML’S FUTURE IS TINY AND BRIGHT

We employ Tiny Machine Learning (TinyML), a cutting-edge applied-ML field that brings the potential of ML to low-cost, low-performance, and power-constrained embedded systems and thereby enables hands-on learning. TinyML lets us impart ML-application design, development, deployment, and life-cycle-management skills.

3.1 Introduction to TinyML

TinyML refers to the deployment of ML resources on small, resource-constrained devices (Figure 2). It starkly contrasts with traditional ML, which increasingly focuses on large-scale implementations that are often confined to the cloud. TinyML is neither a specific technology nor a method per se, but it acts in many ways as a proto-engineering discipline that combines machine learning, embedded systems, and performance engineering. Similar to how chemical

¹The term *full-stack* comes from historic career growth in web technologies and the Internet. It began as a series of loosely linked skills but now encompasses web development from the lowest level, the server, to the highest level, the web browser or mobile app.

Table 1: Cloud & Mobile ML systems versus TinyML systems.

Platform	Architecture	Memory	Storage	Power	Price
Cloud E.g., Nvidia V100	GPU Nvidia Volta	HBM 16GB	SSD/disk TB-PB	250W	~\$9,000
Mobile E.g., cellphone	CPU Arm Cortex-A78	DRAM 4GB	Flash 64GB	~8W	~\$750
Tiny E.g., Arduino Nano 33 BLE Sense	MCU Arm Cortex-M4	SRAM 256KB	eFlash 1MB	0.05W	\$3

engineering evolved from chemistry and how electrical engineering evolved from electromagnetism, TinyML has evolved from machine learning in cloud and mobile computing systems.

The TinyML approach dispels the barriers of traditional ML, such as the high cost of suitable computing hardware and the availability of data. As Table 1 shows, TinyML systems are nearly two to three orders of magnitude cheaper and more power efficient than traditional ML systems. As such, this approach can serve in embedded devices at little to no cost and can handle tasks that go beyond traditional ML. The TinyML approach also makes it easy to emphasize the importance of responsible AI (Section 5).

TinyML supports large-scale, distributed, and local ML tasks. Inference on low-cost embedded devices allows scalability, and their low power consumption enables distribution even to remote locations far from the electric grid. The number of tiny devices in the wild far exceeds the number of traditional cloud and mobile systems [18]. The ubiquity of tiny embedded devices makes TinyML a candidate for local ML tasks that were once prohibitively expensive, such as distributed sensor networks and predictive maintenance systems in industrial manufacturing settings.

TinyML applications are broad and continue to expand as the field gains traction. The approach’s unique value stems primarily from bringing ML close to the sensor, right where the data stream originates. Therefore, TinyML permits a wide range of new applications that traditional ML cannot deliver because of bandwidth, latency, economics, reliability, and privacy (BLERP) limitations.

Common TinyML applications include keyword spotting, visual wake words, and anomaly detection. Keyword spotting generally refers to identification of words that typically act as part of a cascade architecture to kick-start or control a system, such as a mobile phone responding to voice commands [19, 20]. Visual wake words involve parsing image data to find an individual (human or animal) or object. This task can potentially serve in security systems [9], intelligent lighting [21], wildlife conservation [22, 23], and more. Anomaly detection looks for abnormalities in persistent activities [24]. It has many applications in both consumer and commercial markets, such as checking for abnormal vibrations [25] or temperatures [26] to provide early warnings of potential failures and to enable preventive maintenance [27, 28].

3.2 TinyML for Applied ML

An applied-ML engineer should have this full-stack experience to appreciate the impact of the various ML-development stages on the end user. In prototypical ML, such as training large neural-network

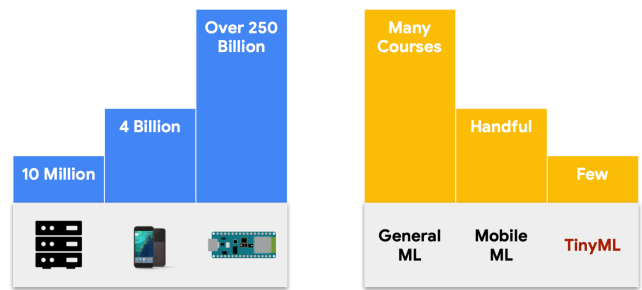


Figure 3: The number of available ML courses is disproportionate to the number of systems in the field.

models in the cloud, learners are unable to participate locally in end-to-end ML development. For example, it is impossible to require them to collect millions of images (akin to ImageNet [29]) for large and complex tasks, such as general image classification. Even more difficult is asking all learners to buy the computational resources to train a complex ML model and then evaluate its performance in the real world.²

By contrast, the small form factor and domain-specific tasks of TinyML enable the full ML workflow, starting from data collection and ending with model deployment on embedded devices. Students thereby gain a unique experience. For example, to implement keyword spotting in their native language, course participants learn to collect their own speech data (e.g., by saying “yá’át’ééh,” which is Navajo for “hello”), train a model on that data, deploy it in an embedded device, and test the device in their community.

Such activities create an immersive learning experience, and they are feasible with TinyML because they only require about 30–40 samples of spoken keywords—easy to collect (only from people with their explicit consent) using a laptop with a web browser and microphone. Learners can then train the model using Google’s free Colab environment [30] and deploy it in a TinyML device using TensorFlow Lite for Microcontrollers [31] or another open-source software technology. This approach allows small keyword-spotting models (about 16KB) to run efficiently on low-cost, highly constrained hardware (less than 256KB of RAM).

3.3 TinyML for Expanding Access

The most difficult task in expanding applied-ML access is making low-cost hardware available anywhere. Cloud-ML technologies cost thousands of dollars, and their physical power, scale, and operational requirements limit their accessibility. Mobile-ML devices are more affordable and pervasive, but their availability is still limited because of network-infrastructure requirements and other factors.

Research shows that although smartphones have become more affordable, their cost remains a barrier in many low- and middle-income countries (LMICs) [32]. Statista estimates only 59.5% of the world’s population has Internet access, with large offline populations residing in both India and China [33]. According to Pew

²Just because a trained model performs well on a test data set does not automatically mean it will perform well in the real world.

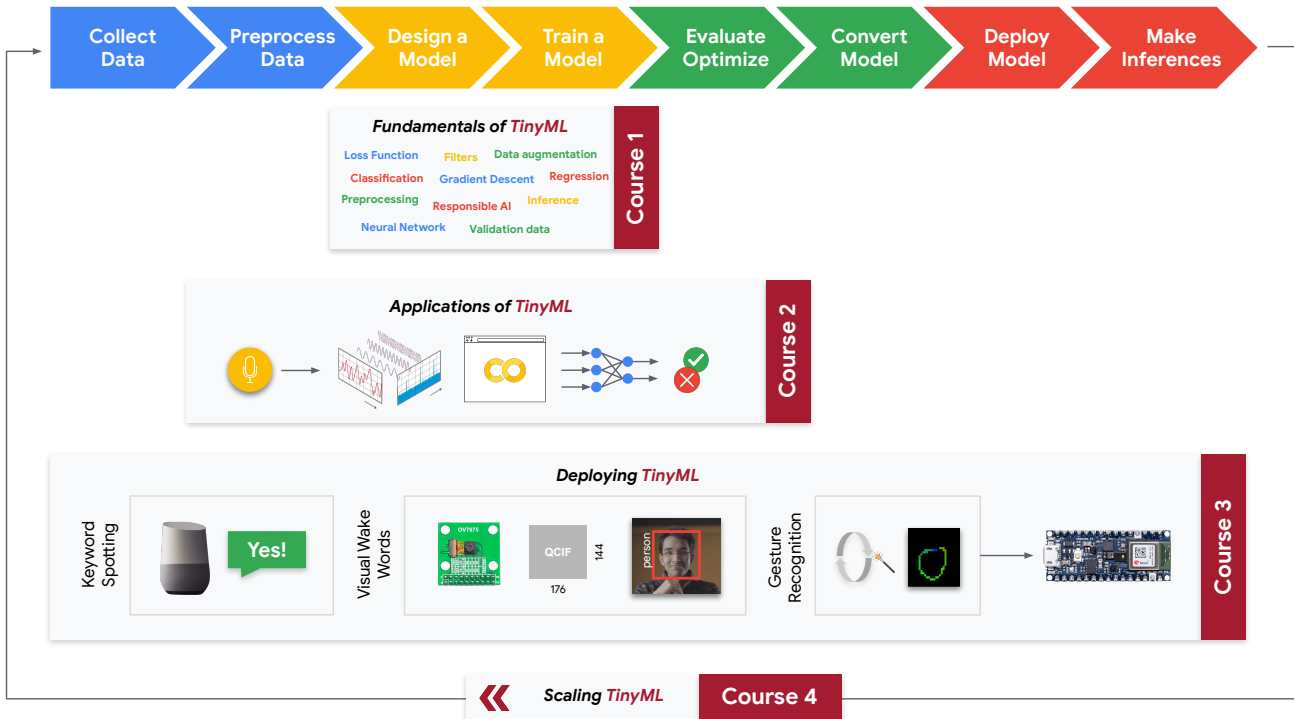


Figure 4: The ML workflow from data collection to model training to inference. The spiral course design focuses on the neural-network model in Course 1, model application in Course 2, application deployment in Course 3, and, finally, TinyML-model management and scaled deployment in Course 4.

Research, 76% of individuals in advanced economies have smartphones compared with 45% in emerging economies. Last in the latter group is India, where only 24% of the population has a smartphone [34]. Students and teachers in many developing countries lack the resources necessary to learn and use traditional ML.

In contrast, TinyML devices are low cost and pervasive. They are readily accessible, enabling hands-on learning anywhere in the world, and their portability eases demonstration of the complete applied-ML workflow in a realistic setting. Furthermore, TinyML applications are more numerous and easier to deploy than mobile-ML and cloud-ML applications. However, despite the wide availability of tiny devices, there is little material for teaching TinyML (see Figure 3). The number of general-ML courses far exceeds the number of TinyML courses (or, more generally, embedded-ML courses).

4 APPLIED-TINYML SPECIALIZATION

We developed an applied-ML course specialization focusing on TinyML. Our specialization provides multiple on-ramps to enable a diverse learner population. Moreover, because TinyML is easy to deploy on hardware and test in the real world, it allows us to systematically explore applied ML’s vast design space (algorithms, optimization techniques, etc.). It also lets us incorporate responsible AI in all four ML stages: design, development, deployment, and

management at scale, which we discuss in greater depth in Section 5. We hope our description of this applied-ML specialization serves as a roadmap for anyone wishing to adopt the program.

4.1 A Four-Course Spiral Design

The TinyML specialization comprises three foundational courses and one advanced course, which we consider optional. Participants would ideally start with the first course and work through the natural progression, but we allow them to go in any order they choose. Depending on their background, they can skip some courses and take the one most relevant to their knowledge and expertise.

As we mentioned earlier, our application-focused spiral design covers the complete ML workflow, going outward from the middle. The curriculum begins with neural networks for TinyML in Course 1, expands to cover the details of TinyML applications in Course 2, deploys full TinyML applications in Course 3, and application management and scaled deployment in Course 4 (Figure 4). Our application focus increases learner engagement and enthusiasm [35], and our spiral design increases the technical depth over time while reinforcing the main concepts, providing multiple on-ramps and eventually enabling students to create their own TinyML application and deploy it on a physical microcontroller.

Course 1: Fundamentals of TinyML	Course 2: Applications of TinyML	Course 3: Deploying TinyML
1.1. Course 1 Overview	2.1. Course 2 Overview	3.1. Course 3 Overview
1.2. The Future of ML Is Tiny and Bright	2.2. AI Life Cycle and ML Workflow	3.2. Getting Started
1.3. Tiny Machine Learning Challenges	2.3. ML on Mobile and Edge Devices (Pt. 1)	3.3. Embedded Hardware and Software
1.4. Getting Started With ML	2.4. ML on Mobile and Edge Devices (Pt. 2)	3.4. TensorFlow Lite Micro
1.5. The ML Paradigm	2.5. Keyword Spotting (KWS)	3.5. Deploying Keyword Spotting
1.6. The Elements of Deep Learning	2.6. Data Engineering	3.6. KWS Custom-Data-Set Engineering
1.7. Exploring ML Scenarios	2.7. Visual Wake Words (VWW)	3.7. Deploying Visual Wake Words
1.8. Building a Computer-Vision Model	2.8. Anomaly Detection	3.8. Gesturing Magic Wand
1.9. Responsible AI Design	2.9. Responsible AI Development	3.9. Responsible AI Deployment
1.10. Summary	2.10. Summary	3.10. Summary
Course 4: Scaling TinyML		
4.1. Course 4 Overview	4.5. Neural Architecture Search (NAS)	4.9. Responsible AI Management
4.2. Profiling TinyML Systems	4.6. Machine Learning Operations (MLOps)	4.10. Summary
4.3. Benchmarking TinyML Systems	4.7. TinyML as a Service (TinyMLaaS)	
4.4. Micro NPUs & Hardware Acceleration	4.8. Federated Learning for TinyML	

Table 2: A breakdown of topics in the four TinyML courses. Each one has several activities, including videos, colabs, hands-on labs, quizzes, readings, assignments, tests, and discussion-forum participation. For a detailed overview of the program as well as links to the course materials, visit our courseware Github: <https://github.com/tinyMLx/courseware>

Table 2 shows a breakdown of the courses. Roughly, each one takes five or six weeks to complete. For a more detailed and up-to-date overview and links to all course materials, visit our courseware Github at <https://github.com/tinyMLx/courseware>.

4.2 Fundamentals of TinyML (Course 1)

Course 1 is titled Fundamentals of TinyML. Its objective is to ensure students understand the “language” of (tiny) ML so they can dive into future courses. TinyML differs from mainstream (e.g., cloud-based) ML in that it requires not only software expertise but also embedded-hardware expertise. It sits at the intersection of embedded-ML applications, algorithms, hardware, and software, so we cover each of these topics. As Figure 4 shows, the course focuses on a portion of the complete ML workflow. Moving to subsequent courses, we progressively expand participants’ understanding of the rest of that workflow.

The course introduces students to basic concepts of embedded systems (e.g., latency, memory, embedded operating systems, and software libraries) and ML (e.g., gradient descent and convolution). The first portion emphasizes the relevance of embedded systems to TinyML. It describes embedded-system concepts through the lens of TinyML, exploring the memory, latency, and portability tradeoffs of deploying ML models in resource-constrained devices versus deploying them in cloud- and mobile-based systems.

The second portion goes deeper by focusing on the theory and practice of ML and deep learning, ensuring all students gain the requisite ML knowledge necessary for later courses. Through hands-on coding exercises, students explore central ML concepts, training

their own ML models to perform classification using Python and the TensorFlow library in Google’s Colaboratory programming environment.

We provide an overview of embedded systems and ML to ensure students recognize that the topics we cover in the specialization are relevant to their lives and careers, boosting motivation and retention [36, 37]. For those with sufficient ML and embedded-systems experience, Course 1 is optional. By designing the series with these multiple on-ramps, we can meet participants wherever they are, regardless of their background and expertise.

4.3 Applications of TinyML (Course 2)

The objective of the second course is to give learners the opportunity to see practical (tiny) ML applications. Nearly all such applications differ from traditional ML because TinyML is all about real-time processing of time-series data that comes directly from sensors. As Figure 4 shows, we help students understand the complete end-to-end ML workflow by including additional stages, such as data preprocessing and model optimization. Moreover, when we revisit the same stages (e.g., model design and training), we employ spiral design to broach advanced concepts that build on Course 1.

Course 2 examines ML applications in embedded devices. Participants study the code behind common TinyML use cases, such as keyword spotting (e.g., “OK Google”), in addition to how such front-end, user-facing, technologies integrate with more-complex smartphone functions, such as natural-language processing (NLP). They also examine other industry applications and full-stack topics,

including visual wake words, anomaly detection, data-set engineering, and responsible AI.

We take an application-driven approach to teaching the technical components. For example, we use the keyword-spotting (KWS) example to demonstrate the importance of preprocessing sensor inputs, showing the power of FFTs [38] and MFCCs [39] through coding exercises. We additionally explore the importance of holistic architecture by discussing the QoS metrics that evaluate KWS applications and the “cascade architecture” (i.e., ML models staged one after another for efficiency) for deploying them [40]. As another example, through the lens of the visual wake words (VWW) application, we introduce transfer learning [41], teaching students to develop their neural-network models without voluminous training data and expensive hardware. Supplementing the theoretical concepts is a coding exercise that employs transfer learning on a pretrained MobileNet [42] model to detect whether an individual is wearing a mask—a real-world application that will resonate with learners in light of Covid-19. As a final example, we use anomaly detection (AD), in the context of predictive maintenance for manufacturing, to demonstrate the power (and limitations) of supervised learning and deep neural networks by exploring k nearest neighbors [43], an unsupervised traditional-ML technique, and comparing it with autoencoders [44], an unsupervised neural-network technique.

The course not only teaches students about TinyML applications and their technical components, but also how to run and test these applications using TensorFlow Lite in Google’s Colaboratory programming environment. This step completes the second learning spiral, providing a hands-on opportunity to explore full TinyML applications. The inclusion of TensorFlow Lite lets participants explore important TinyML topics (e.g., neural-network quantization), preparing them to add the next layer in Course 3: physical hardware. We intentionally avoided introducing microcontroller hardware until the third course so students could complete the first two entirely for free. They can then make an informed decision on whether to buy the low-cost hardware that Course 3 requires.

4.4 Deploying TinyML (Course 3)

Most instructional ML material focuses on models and algorithms, failing to provide hands-on experience in gathering input or training data, making decisions on the basis of a model’s output, and testing models in the real world. Therefore, Course 3 explicitly focuses on demonstrating the complete ML workflow (Figure 4).

We have found that learners are excited about using the knowledge they gain to solve real problems. But absent guidance in how to build an entire system, many become frustrated because they are unable to apply their knowledge. This issue arises in the form of questions that traditional courses leave out, such as “How can I find the training data for my problem?” and “What threshold should I use to decide whether a classification score is high enough for my application?” and “How do I go from an RGB-camera-image byte array to the float-tensor image my model needs?”

Providing universally applicable answers to all such questions is impossible. But alerting early learners to them and offering a set

of comparable guided practical experiences reduces the frustration before these questions arise in their projects, hobbies, or careers. This is critical, as frustration squelches the desire to master the skills a field requires. Instead, we want to give students the support and confidence they need to overcome challenges and solve problems.

“Deploying TinyML” mixes computer science and electrical engineering. It gives participants fundamental knowledge and hands-on experiences with ML training, inference, and deployment on embedded devices. Following the spiral pattern, it builds atop many of the techniques and applications from previous courses and adds new technical topics and extensions. Students develop and deploy full applications, such as KWS, person detection, audio/visual reaction, and gesture detection on their own microcontrollers.

The course introduces TensorFlow Lite for Microcontrollers [31], an embedded-ML software library that eases the task of efficiently running ML models on embedded hardware. Students learn how the library works, helping them appreciate the challenges that an embedded-ML-framework engineer faces in the real world. They also examine the library’s APIs as they deploy applications such as KWS, VWW, and magic wand to their microcontrollers.

Building on the concepts from the first two courses, we introduce new concepts such as multitenancy—that is, running more than one ML model at a time—when we revisit certain stages of the ML workflow (Figure 4). We present the tradeoffs between using multimodal learning that fuses sensor data versus using two separate models to make inferences. The former stresses the first half of the ML workflow (training), while the latter stresses the second half (deployment).

A unique benefit of this course is the exercises that involve the entire ML process. Before they know it, students are implementing an entire TinyML application from beginning to end on a physical device they can hold in their hands. This approach gives our course the unique value of allowing participants to fully develop and use their own TinyML projects at home. This type of hands-on project-based learning is proven to enhance learning, motivation, and retention [47, 48]. For instance, participants collect their own custom keyword data for training a KWS model, giving them first-hand experience with the challenges of getting ML models to work accurately. Some find that the *tinyConv* [49] KWS model works well; others find that they must collect more data or adjust the preprocessing. A few individuals in this latter category are perplexed that even those improvements may fail to dramatically increase accuracy, finding that the 16 KB KWS model is too small. The point of the exercises is not necessarily to increase model accuracy, but to instead understand the challenges of applying ML models to the real world.

The course uses the Tiny Machine Learning Kit, which we co-designed with Arduino for hands-on, low-cost, accessible, project-based learning. This kit, shown in Figure 5, is globally accessible, and includes an Arduino board containing numerous sensors (microphone, temperature, humidity, pressure, vibration, orientation, color, brightness, proximity, gesture, etc.) that enable a wide range of TinyML applications. More importantly, it has a popular Arm Cortex-M-class microcontroller [50] that binds the learning experience to reality. The kit provides everything a student needs to build

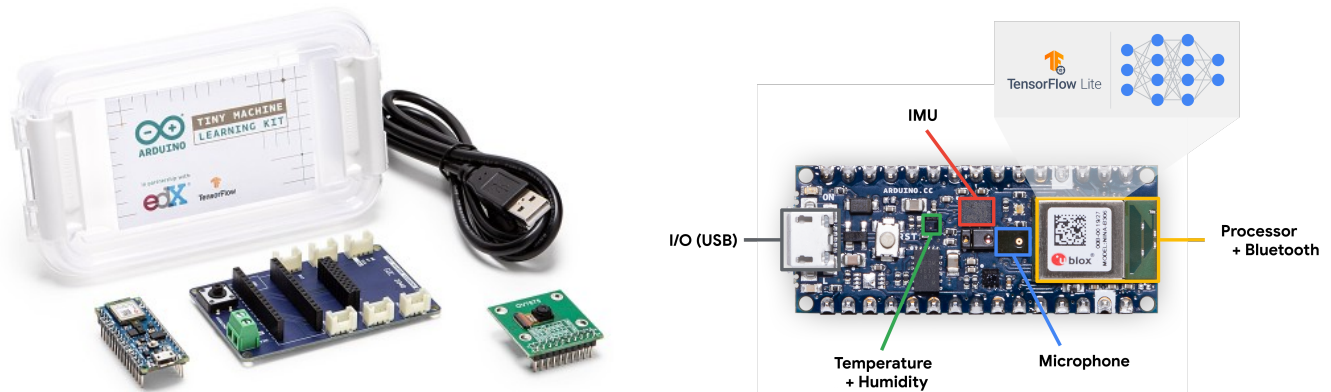


Figure 5: Shown at left, the TinyML Kit includes the Arduino Nano 33 BLE Sense [45], an OV7675 camera module [46], and a TinyML shield that simplifies sensor integration. Shown at right, the Nano 33 BLE includes a host of onboard integrated sensors (e.g., temperature/humidity, IMU, and microphone), a Bluetooth Low Energy (BLE) module, and an Arm Cortex-M microcontroller that can run neural-network models using TensorFlow Lite for Microcontrollers.

TinyML applications for image recognition, audio processing, and gesture detection.

After completing Courses 1, 2, and 3, students are eligible to receive the HarvardX/edX Tiny Machine Learning Certificate, testifying they are trained as full-stack TinyML engineers. We offer the certificate because many professional learners desire such awards to enhance their resumes and prove to potential employers that they have mastered particular skills. At this point in the series, participants have not only explored the technical and societal challenges that TinyML poses, but they have also gained hands-on experience with the complete TinyML-application pipeline: collecting data, developing and training models in the cloud using TensorFlow, testing them in the cloud using TensorFlow Lite, and deploying them in hardware with TensorFlow Lite for Microcontrollers.

4.5 Scaling TinyML (Course 4)

The first three courses bring learners up to speed on designing, developing, and deploying TinyML applications on a device. Course 4 builds on this foundation and considers scaled management of TinyML-application deployments. This advanced course covers two aspects of scaling. The first half focuses on “scaling up” the effectiveness of individual TinyML applications through performance benchmarking, model optimization, and hardware/software co-design. The second half focuses on “scaling out” TinyML applications from one device to thousands.

In the scaling-up portion of the course, we start by introducing the concept of system-performance profiling through TinyMLPerf [51] and other open-source, industry-standard benchmarks. Because embedded systems deal with sensor data in real time, the TinyML device must be able to keep up with the data rates. In safety-critical systems, such as automobiles, a slow response to new sensor inputs can be life threatening. Knowledge of benchmarking principles is therefore essential; it enables applied-ML engineers to compare ML

systems in a fair and useful manner and to make informed decisions when selecting a device for a particular task.

Next, given a suitable system, we discuss the art of picking a suitable model, emphasizing when and whether to be more code-centric (improve the model) or data-centric (improve the data set). The most accessible means of decreasing model latency is to change the neural-network architecture. Students can accelerate inference without changing any code by designing a new model that is sufficiently accurate yet requires fewer calculations. But designing ML-model architectures is difficult and time consuming because of the many decisions that affect model quality and latency: what type of neural network to choose, what size to make it, how many hidden layers and neurons to include, how to best initialize the network, and so forth. Fortunately, services have emerged to help design network architectures automatically. Cloud services such as AutoML [52] and techniques such as neural-architecture search (NAS) [53] allow even developers with limited ML expertise to quickly train high-quality models that meet their needs. In TinyML, such services are essential because achieving efficiency means co-designing the MCU hardware, software, and models, a challenging task for humans. The course therefore introduces concepts such as AutoML and NAS, explaining the fundamentals so students can employ high-level automation tools with confidence.

The second half of the course focuses on “scaling out” TinyML applications from one device to thousands. Applied-ML engineers must know how to manage such deployments as a production ecosystem may involve hundreds or thousands of devices. We thus offer a preview through the TinyML lens. We start by leveraging “ML operations” (MLOps) to develop, monitor, and improve a TinyML application. MLOps automates the complete workflow (Figure 4) as Figure 6 shows. We discuss ways to automatically manage and process data, train ML models, version them, and evaluate, compare, and deploy them, all from the perspective of complete MLOps

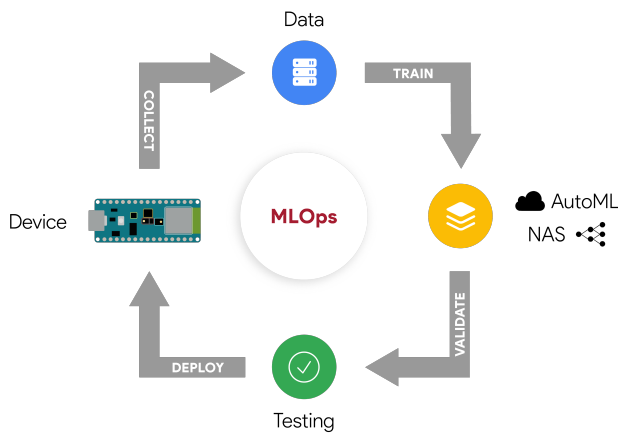


Figure 6: Scaling TinyML through MLOps.

platforms. This approach introduces the advantages of an automatic ML workflow, which include managing the overwhelming complexity of ML deployments, reducing the burden of maintaining in-house ML knowledge, being more scientific, easing long-term maintenance, and ultimately improving the model’s performance in the field.

In addition, we also introduce TinyML as a service (TinyMLaaS), which allows production ecosystems to easily manage and integrate heterogeneous TinyML devices. Because embedded TinyML devices are specialized—from the ML compilers to the operating system and ML hardware—to achieve ultra-low-power energy efficiency, a highly fragmented ecosystem can result. Fragmentation limits a precompiled ML inference model’s portability when the hardware changes. The model then requires recompilation for a particular device, leading to deployment complexity. To reap the efficiency benefits of hardware heterogeneity while coping with the fragmented ecosystem, we need a new “as-a-service” abstraction. To this end, we introduce learners to TinyMLaaS [54], a general method for tailoring an ML inference model to a specific device. It is a software abstraction layer that gathers information about the target device—such as the CPU type, RAM and ROM sizes, available peripherals, and underlying software—to generate the correct compiled inference model. The designated device then automatically downloads this generated inference model, and the process repeats for all other devices. Because TinyMLaaS enables firmware-over-the-air (FOTA) and software-over-the-air (SOTA) updates, it introduces privacy and security concerns for both the model and the data. Hence, we also cover how large device networks can train models while maintaining user privacy through federated learning [55].

Participants who complete the four courses will have learned all the fundamentals of ML-model design, development, deployment, and management through the TinyML lens. This knowledge is invaluable for career advancement in this quickly emerging field.

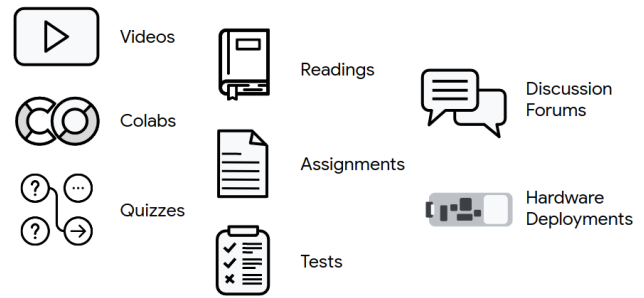


Figure 7: We employ a wide array of learning activities to give students an immersive, self-paced online experience.

4.6 Student Activities

People learn differently [56]. To support many learning styles, we implemented proven strategies [57] and a variety of methods (Figure 7). Our approach mixes video lectures, short readings, and coding exercises in Google’s Colaboratory programming environment to teach and reinforce the course’s main technical components. Thus, visual, auditory, and experiential participants all learn by their preferred method.

We keep the videos short (4–10 minutes). Research shows that people learn better from numerous short content modules than from a few long ones [58]. Because some students will find that many of the coding components are new, we provide walk-throughs of major sections, numerous comments describing the code, and introductory text to explain the purpose of each code snippet.

In the first two courses, each section builds toward a coding assignment in Colaboratory to encourage project-based exploration and creativity. The assignments in Course 3 expand to full-on hardware deployment that lets students hold their own designed, trained, and deployed model in their hands and test it in the real world.

The activities grow in complexity and detail as students progress through the courses, following our spiral-design principle. Students thus gain confidence throughout, as complete application deployments can be challenging. Finally, we sought the input of industry veterans on our course staff to ensure the hands-on activities build relevant full-stack skills.

The courses also include formative multiple-choice quizzes throughout, focusing on the main concepts so students see their progress even if they do not understand every line of code. The quizzes also reinforce the importance of high-level tradeoffs and applied-ML concepts, which will be relevant to ML careers even if the technical stack changes. For those pursuing a paid certificate, we also included summary tests at the end of each section.

Finally, we provide many discussion forums that allow students to ask questions and get answers. Forums allow the course staff to support all participants regardless of their location. They also serve the dual function of building a community around the course. Our forums encourage students to ask any and all questions and to answer them for one another.

Each activity includes a strong ethics component, which we describe in detail later (Section 5). Briefly, however, we ask many open-ended questions to elicit student opinions on the opportunities and challenges of responsible TinyML-application design, development, and deployment. As the literature predicts [57], many of these questions have led to conversations and debates between our online participants, despite their different geographic locations, ages, and technical backgrounds.

4.7 Accessible, Hands-on Learning

To enable hands-on learning anywhere in the world, we need a low-cost, self-contained yet extensible, approachable yet representative, flexibly abstracted, and globally accessible TinyML platform. Once again, microcontrollers are promising because they are inexpensive and widely available. So, to provide an easily accessible out-of-the-box experience, we custom designed the Tiny Machine Learning Kit (Figure 5) with Arduino. This section describes the kit and its development.

Systematic selection. The range of TinyML hardware and software options is wide, but we believe an ideal solution is fully self-contained yet extensible, approachable yet representative, and flexibly abstracted. As such, we searched for one that not only made it simple to integrate the sensors required for the course but also supported easy integration of additional sensors for future study.

Putting application-level software aside leaves two fundamental elements: hardware, from the microprocessor to peripherals and discrete circuitry, and software, from the application layer (our focus) to the silicon layer. An initial constraint on the field of potential microprocessors is the need to support TensorFlow Lite for Microcontrollers [31], which is written in C++ and requires that the microcontroller support 32-bit computing.

We developed criteria for compatible microcontroller development boards, recognizing that an integrated off-the-shelf product would greatly increase accessibility. These criteria include a small form factor (it is *tiny* ML, after all), a low power budget (efficiency is critical to edge computing), a small system memory (some controllers have large memories, making them less accessible and limiting their range of application), sufficient clock speeds, wireless-communication capability (to enable periodic reporting and/or distributed systems), select sensor integration, and serial channels for extensibility. We defined similar criteria for the accompanying software, comprising the development environment, embedded framework, and logistics (fast, reliable distribution). Next, we added weights to the selection criteria and compiled the candidates in a Pugh matrix [59]. We ranked a field of about two dozen hardware products, giving some preference to controllers that had undergone more-extensive testing—in particular, Arm’s Cortex-M series [50] and Espressif Systems’ products (namely, the ESP32t) [60]. Both of these embedded systems are widely popular.

The TinyML kit. Ultimately we selected the Arduino Nano 33 BLE Sense [45] because it uniquely blends expert embedded-systems engineering and remarkable isolation of the application developer from many low-level hardware details [61]. Furthermore, the Arduino framework and its software APIs (“cores”) fit naturally

with our spiral design. Arduino’s many libraries and simple IDE are easy for inexperienced students to learn, yet it typically permits those interested in the “bare metal” to work their way down the embedded-software stack. Moreover, the Nordic nRF52840 Cortex-M4-based controller [62] on the Nano 33 BLE Sense development board, along with its Mbed real-time OS [63], represent industry-level hardware and software.

We also developed the Tiny Machine Learning Shield to enable plug-and-play integration of sensors that the Nano 33 BLE Sense lacks. In particular, it eliminates the need for users to make 18 individual connections between the microcontroller and the low-cost camera module we selected for the course—the OV7675 [46], which typically sells for about US\$2. A series of Grove connectors [64] line each side of the shield for connection to numerous additional sensors, which students can purchase for their own projects and integrate without soldering or low-level circuit design.

We bundled the Nano 33 BLE Sense with the shield, the OV7675 camera module, and a USB cable to form the Tiny Machine Learning Kit (Figure 5); learners can purchase a single item and be fully prepared for Course 3 for US\$49.99. To accommodate those few who prefer to purchase elements individually, we provide wiring diagrams and a custom Arduino software library so they can readily swap the OV7675 for the related OV7670 camera module.

Alternatives. In the months after we developed and announced our TinyML kit, similar boards emerged to provide alternative options. For example, the Pico4ML by ArduCam [65] is a notable single-board example that comes complete with a microphone, inertial measurement unit (IMU), and camera module, and is suitable for the course exercises. We are working to support some of these new and exciting hardware platforms to give students more flexibility with their projects.

5 ETHICAL & RESPONSIBLE AI

Ethical and Responsible AI is about putting people, social benefit, and safety first. More specifically, ethical AI emphasizes the need for ML engineers to safeguard user privacy and security, mitigate algorithmic bias and discrimination, and ensure ML models perform reliably after deployment. It also extends to developing consumer trust. In this section, our goal is to shift learners from thinking about which ML technology is feasible to which is useful, with an understanding of how it will influence users and society.

5.1 Ethical Consideration of Ubiquitous ML

TinyML offers many helpful features, ranging from data privacy and security to low latency and high availability. Coupled with low-cost embedded hardware, these features make it a pervasive technology that can enable ML everywhere. TinyML sensors will monitor the environment in which they are deployed, be it mechanical or human, around the clock. With the prospect of ML everywhere comes a pressing need to address privacy, drift, bias, and other ethical issues.

Fortunately, TinyML allows us to incorporate responsible AI into all four ML stages: design (Course 1), development (Course 2), deployment (Course 3), and scaling (Course 4). By embedding

ethics into each TinyML course, we communicate the technology’s ethical and social dimensions in a personal and practical manner.

To achieve deep integration, we follow the Embedded EthiCS pedagogy at Harvard [66], where philosophers participate directly in computer-science courses to teach students how to think through the ethical and social implications of their work. We collaborated with a philosopher from this program to co-develop and include such material in our curriculum. Her commitment to learning the technical aspects of TinyML enabled us to customize the ethical content to meet the unique course needs of TinyML.

By distributing responsible AI throughout the series, covering the entire ML workflow, students discover how ethical issues permeate all aspects of their work. Our aim is to introduce them to the conceptual tools for navigating these issues, in hopes they will view responsible AI as an active enterprise. Next, we describe our pedagogical goals for each responsible-AI unit, some examples we covered, and the exercises that reinforce the concepts.

5.2 Designing AI Responsibly (Course 1)

Access to, adoption of, and use of ML products is inequitably distributed. According to Pew Research, 64% of Americans believe technology companies create products and services that benefit people who are already advantaged, and 65% believe these companies fail to anticipate the societal impact of those offerings [67].

To enable more-widespread, safer, and more-secure ML, we must raise awareness of its capabilities. Thanks to the low cost and accessibility of TinyML hardware, our students are diverse, and they will probably have to address different social and cultural factors when designing ML applications. To ensure they all can anticipate the effects of ML products and ensure equitable access, our approach to responsible AI focuses on forming a vision of both the problem to be solved and the people a solution will affect.

We believe that by taking an active role in responsible ML design, students will be better able to address ethical challenges such as bias, fairness, and security. We therefore cover real-world examples, such as a Winterlight Labs auditory test for Alzheimer’s disease. In this case, research revealed that nonnative English speakers are more likely to be mistakenly flagged as having Alzheimer’s [68]. In a discussion forum, students reflected on what the product designers could have done differently to avoid this failure. Such activities reinforce the importance of considering diverse user perspectives during the design phase, as doing so can inform data-collection decisions that mitigate ML bias.

In a subsequent forum, learners practice ethical reasoning about the consequences of a KWS model’s failure in terms of Type I (false positive) and Type II (false negative) errors. In this case, a false positive would result in audio being recorded, unbeknownst to the user, and sent to the cloud. A false negative means the device failed to activate when the user spoke the wake word. Students must justify their decision to optimize the model for high precision, thereby minimizing false positives, or to optimize for high recall, thereby minimizing false negatives.

For the KWS activity, nearly all participants chose to optimize for high precision to minimize the risk of privacy violations. Interestingly, one provided a justification based on sustainability concerns related to unnecessary data transmission and storage in the cloud. Those who decided to optimize for high recall cited a variety of reasons. One noted that although people claim to value privacy, they tend to prioritize convenience. In contrast, another suggested enacting privacy measures elsewhere to offset the potential harm of optimizing for high recall. Lastly, yet another student prioritized model performance to meet user expectations. That student claimed the burden of preserving privacy should fall on the user, who has the ability to decide whether to purchase the product. There is no right or wrong answer. Our desire is to spur self-reflection and foster constructive discussion among learners from different backgrounds and cultures.

5.3 Developing AI Responsibly (Course 2)

Any developer employing ML must be aware of how data-collection bias and fairness affect application behavior. Our courses use public data sets, including Speech Commands [69], Mozilla Common Voice [70], ImageNet [29], and Visual Wake Words [71], for nearly all of the programming assignments. Most data sets, however, have demographic-representation problems [72]. For example, despite crowdsourcing efforts to increase diversity, the Common Voice data set lacks equal gender representation (only 24% of English-data-set contributors who revealed their gender are female) [73].

Our goal is for students to see how data collection, bias, and fairness intertwine, as well as to equip them to mitigate the problems, because they are working with KWS models, we cover real-life biases relevant to this kind of ML application. For instance, research shows that voice-recognition tools struggle to identify African American Vernacular English, causing popular voice assistants to work less well for black individuals [74]. Similarly, research shows that voice recognition struggles to identify nonnative English speakers and those with speech impairments [75]. To acquaint learners with recent work in mitigating bias, we discuss Project Euphonia, a initiative that launched in 2019 with the goal of collecting more data from individuals with speech impairments or heavy accents to fill the gaps in voice data sets [76].

We created a Colab activity that uses Google’s What-If Tool (WIT) [77], based on its Responsible AI tool kit [78]. The WIT is one of the company’s many open-source, interactive visualization tools for investigating trained-ML-model behavior with minimal coding. In this exercise, participants practiced ethical reasoning by exploring a real-life data set, identifying sources of bias, and evaluating threshold-optimization strategies for fairness. For the WIT activity, students noted how the visual representations fostered a deeper understanding of issues pertaining to fairness. One claimed that the focus on confusion matrices in particular was an effective way to clearly distinguish between the fairness metrics. In general, learners appreciated the opportunity to try out the WIT.

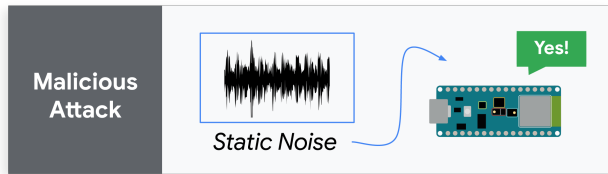


Figure 8: Students attack a pre-trained KWS model with malicious static noise and trigger a spotting of the keyword "yes," showing them the importance of security and privacy.

5.4 Deploying AI Responsibly (Course 3)

Even after designing and developing an ML model, deployment raises a new set of ethical challenges. For example, TinyML systems are often touted for preserving privacy. When an embedded system processes data locally (close to the sensor) rather than transmitting it to the cloud, we tend to believe it protects user privacy. But user interaction with the model raises new privacy and security concerns [79]. Moreover, ML interacting with a dynamic real-world environment using sensors raises concerns about model drift³ over a product’s lifetime.

To the extent that TinyML enables ML everywhere, the privacy, security, and even model-drift risks could be more widespread compared with traditional ML. To familiarize students with these risks, we cover real-life examples, such as doorbells that share data with law enforcement [80] and fitness devices that leak user information [81].

Our goal is to equip students with strategies to mitigate these risks when deploying trained models in embedded devices. Importantly, the mitigation strategies available to traditional ML systems are sometimes unattractive or infeasible for TinyML. For instance, the resource constraints of an embedded device, such as low power and small memories, complicate implementation of robust security systems and model retraining. Therefore, we acquaint course participants with a wide array of strategies, such as minimizing the transmitted and stored data to preserve user privacy, minimizing hardware design to limit vulnerability to attackers, and running supervised experiments in the real world before releasing ML models.

Inspired by research showing we can use inaudible ultrasonic noise to trigger or eavesdrop on KWS models [82], we created an exercise that gives students hands-on experience attacking a KWS model. They trigger a false positive—“yes”—with seemingly innocent but adversarial static noise (Figure 8), which in a real application would cause the system to constantly record and transmit the audio. This experience builds on our videos and readings and makes the security threat real—a crucial part of any major security-awareness program at large [83]. At the same time, it is also a cautionary tale of ML’s limitations, a lesson all applied-AI engineers should learn.

³Model drift generally refers to prediction-accuracy degradation owing to environmental changes.

To further reinforce this point, a subsequent discussion forum allows course participants to practice ethical reasoning to determine when malicious triggering of a false positive can cause serious harm. Some have noted that this vulnerability would be most likely to cause harm where security is a paramount value, such as using KWS to grant access to a secure space or to initiate a financial transaction. One student drew a connection to the practice of ethical hacking, or penetration testing, and the possibility of developing adversarial data for retraining the model to be more resilient. Interestingly, another noted that since users lack the ability to fix security issues, their only option is to stop using the device. But this choice ultimately depends on whether the company informs customers about the vulnerability. Lastly, one student claimed the adversarial example was more reliable than that student’s own voice for triggering a “yes.” The course staff then responded, prompting a discussion of data-set bias and the likelihood that American male accents are overrepresented in the data.

5.5 Scaling AI Responsibly (Course 4)

Many ethical implications require consideration when applying technology. Even minor biases, which can be difficult to detect in the proof of concept, can have a major impact when appearing in thousands or millions of devices.

This problem highlights the need to treat responsible ML as an iterative process. Rather than introduce entirely new ethical considerations, we revisit and expand on previous ones. For instance, to guide students in cleaning up a test data set before they conduct benchmarks, we revisit the ethical issues of data collection and bias. Similarly, we revisit privacy in depth once participants become acquainted with federated learning.

We are incorporating an active learning exercise using Google’s Model Card Toolkit (MCT) [84]. Model cards are a reporting mechanism that can increase model transparency and facilitate the exchange of information between model creators, users, and others. This exercise requires that students practice using the model-card framework to document information relating to the model’s development, performance, and ethical considerations.

We additionally discuss the environmental impact of large-scale TinyML networks, as the production and maintenance of billions of MCUs can have led to substantial carbon emissions. Beyond the ethical pitfalls of scaling TinyML, we cover the potential positive social impact this technology can have in domains such as environmental sustainability, public health, and AI equity.

6 ACCESS VIA MOOCS

In this section we describe how we leveraged technology to make the TinyML specialization broadly accessible and highlight important considerations made to ensure we supported our remote learners.

6.1 Massive Open Online Course

Our goal was to reach a global audience. We therefore chose to employ massive-open-online-course (MOOC) platforms. Examples such as edX and Coursera are ideal for making the content globally

accessible; students need not travel to a different country to learn. These platforms host a wide variety of university-level courses and are generally cheaper than equivalent academic and professional training thanks to the economics of scale [85]. We deployed the TinyML specialization on edX through HarvardX.

Participants can audit the course for free or pay to earn a professional certificate. Since they can “upgrade” to a professional certificate at any point during the course, both students and professionals can try before they buy, encouraging more to enroll. Although the professional certification includes summary tests that are absent from the audit version, we designed the curriculum so individuals who are just auditing learn the same crucial principles. Thus, all can attend the entire class, developing their skills for free. At the time of this writing, the number of auditors far outweighs the number of paying students by more than order of magnitude.

The course is asynchronous and self-paced rather than instructor led. Students progress through the material at whatever speed they find comfortable. But unlike in-person courses, interaction between students and staff is minimal, forcing staff to develop high-quality, self-explanatory, and self-sufficient materials that rely heavily on media (which we describe in greater detail in Section 6.2).

Unlike most MOOCs, Course 3 employs the TinyML kit (Section 4.7) for hands-on learning. To maximize hardware accessibility, we worked with Arduino to make a custom all-in-one kit globally available for purchase through either that company’s website [86] or one of its many distributors. We also provided a detailed bill of materials for students who wish to buy individual components instead. The main benefit of this approach is that it improves the efficiency for the host institutions (Harvard and Google) by reducing the burden on them for managing inventory and shipping logistics (taxes, international shipping rates, etc.).

6.2 Accelerated Remote Media Production

The typical development timeline for a series of online courses, such as the ones we described in Section 4, is about two years—far too long to keep up with changing ML technology. Applied ML, especially in the context of TinyML, remains a nascent yet quickly developing field. Therefore, media production for the online curriculum must be rapid to ensure the material is timely and relevant and to ensure broad access.

We compressed the media-production time greatly, achieving an average development cadence of 6–8 weeks per course. To maintain this cadence, we created a custom remote-media-production workflow. We produced the TinyML course under the specter of Covid-19, but regardless of the safety limitations, a remote production strategy would still have been the only way to achieve these quick results. Remote production methods offer flexibility and allow an international crew to make contributions, meaning the process continues around the clock. Regardless, no matter when and how it is done, creating a flexible workflow requires a principled content-design approach, and advanced technology is necessary for rapid progress. The following breakdown can serve as a roadmap for others attempting to follow a similar approach.

Production design. To expand access such that our effort meets the needs of a global audience (Section 2), we built our media-production strategy around five critical ingredients: compelling instructional narrative, best media practices for online learning, a diverse and skilled production team, prioritized use of production equipment, and willingness to innovate.

A *compelling instructional narrative* that whets student interest is critical, as all great media experiences unite around a good story. TinyML offers a sound narrative because it provides an accessible, hands-on introduction to ML (Section 3). We aimed to communicate with a global audience and provide the practical knowledge for building complete, relevant TinyML applications and tools.

Effectively communicating that narrative requires *best media practices for online learning*. Decisions made in postproduction often hold more weight than any others. For example, the decision of when to show the instructor, slides, or both in a picture-in-picture cut versus when to display graphics or other visual/auditory information can affect the viewer’s cognitive load and overall learning [87, 88]. When in doubt, Mayer’s “12 principles of multimedia learning” [89] is an excellent place to discover such general practices for enhancing the student’s experience.

From the start, we determined the primary media types we would produce to hold students’ attention and maintain their cognitive load balance [87]. We chose picture-in-picture and split-screen formats, allowing us to show the instructor or other imagery in full-screen mode to focus on the most important aspects of the presentation (Figures 9). We emphasized instructor screen time, however, to improve student learning [90].

A *geographically distributed and responsive team* is necessary to quickly produce highly sophisticated content, especially for an emerging technical field. Our media team included a producer and director to establish a creative vision and ensure media delivery, a senior editor to assemble and craft the videos, a motion-graphics designer to provide custom graphical elements for our brand, and a production assistant to wrangle data, review content, and integrate the final videos into the platform. This team was relatively lean. One additional advantage was that contributors were scattered across 12 time zones (San Francisco to Boston to London to Mumbai), meaning at all times someone was awake and working on the project.

A crucial ingredient to quickly producing content is *prioritized use of production equipment*. The remote nature of the production and the Covid-19 precautions only heightened this need. For example, webcams and audio supplies were sold out or on back order because people were setting up home offices so they could continue to work. Fortunately, we were able to make acceptable compromises and buy equipment in a way that ensured the greatest impact. We prioritized production-equipment purchases as follows: 1) audio, because it is more important to retaining viewer attention than video [91]; 2) lighting, as it can improve even a nonideal camera to draw the student’s eye; and 3) video, which we mention last because it is the most expensive in a context where higher production value does not necessarily imply a better learning experience [58].

The final ingredient was *willingness to innovate*. Course 3 (Deploying TinyML) involves hands-on learning. Typically, in-person

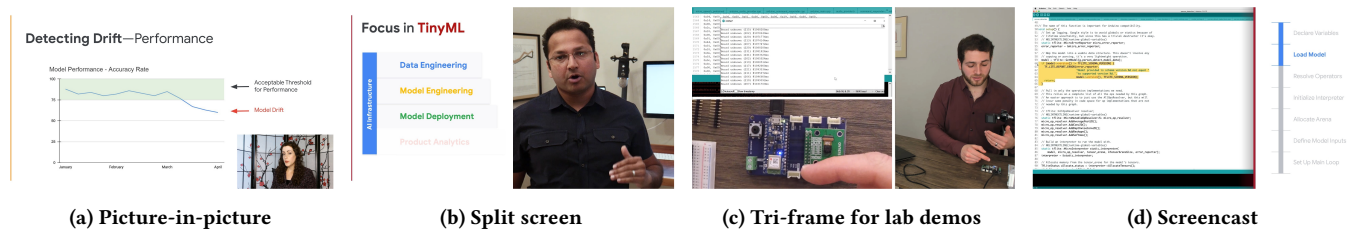


Figure 9: We used various video-production strategies throughout the course to maximize learning efficiency. (a) Picture-in-picture places a video clip in a small frame on top of another frame, playing them simultaneously. It enhances the perception of instructor presence while showing the student relevant content. (b) Split-screen, a slight variation, also improves the perception of instructor presence. (c) Tri-frames are useful for lab demos to enable “hands-on” instruction from teaching assistants. (d) Screencasts coupled with associated lecture material reinforce concepts with code.

teaching assistants (TAs) demonstrate labs to show students the goals and scope and to preemptively troubleshoot common errors. Doing the same online is extremely difficult. We developed a three-way split-screen medium (Figure 9c) that displays the device assembly, device testing, and TinyML lab exercises. We assembled a new film location to (remotely) support the teaching staff with the lab exercises, adding an overhead camera and additional lighting. Furthermore, we enhanced our visuals for the three-way split screen with a custom motion-graphics layer. This setup reached completion and underwent rapid testing without disturbing the production timeline. From start to finish, Course 3 took only eight weeks despite involving five hours of produced-video time, which includes short lectures, screencasts, and lab videos.

Technology. Without globally accessible technology and services, remote media production at the level and pace we achieved would have been impossible. Cloud storage was the backbone of our strategy. It allowed contributors to ingest and manage footage globally. It was also the heart of our production workflow, giving us the ability to sync media project files instantly. Videotelephony services such as Zoom and Google Meet aided in assessing home-studio setups in addition to serving as a virtual rehearsal stage and writers’ room. Amazon supplied 90% of our equipment. Frame.io streamlined our video-quality review and revision [92].

Copyright. Although on-camera presence was a major focus of remote production, video lectures are just one part of the students’ activities. At the same time, a multidisciplinary team of content experts, graphic designers, and web developers at HarvardX rapidly designed and formatted readings and coding exercises. A major challenge in quickly producing course materials was ensuring each illustration, photo, and code library met strict licensing requirements to avoid copyright infringement. Given our project’s more than two thousand graphics and tight timelines, we trained all content developers on proper sourcing for course materials. In-house custom graphics—necessary for a nascent field—predominated, and copyright specialists at HarvardX evaluated each piece as it arrived to cite all external creators.

6.3 Building Community

A common and well-known pitfall of MOOC platforms is the difficulty of developing community and fostering peer learning among a geographically distributed population. Students often struggle to discuss and collaborate after completing the course and even during the course. We therefore developed the TinyMLx community, which welcomes everyone beyond the edX platform.

First, we created a Discourse forum ([discuss.TinyMLx.org](https://discuss.tinymlx.org)) to provide both a communication platform for students and a home for future initiatives. It has been successful, garnering over 3,500 user visits and over 58,500 page views in its first five months. We also conducted two live Q&A sessions for the TinyML community. For each session, between 100 and 200 learners joined live from around the world, and many more have since watched the recording. We received dozens of questions leading up to the events and dozens more during, with topics including how to best teach TinyML material, how to improve diversity in TinyML, and many others in between. Participants enjoyed the events, e.g., 90% of respondents to our first post-event poll said they would like to attend another. Finally, based on learner feedback we recently created a Discord chat to further enable easy collaboration, communication, and community building.

To challenge our Course 3 students, who went on to deploy TinyML models on their microcontrollers, we developed an optional “capstone-project” competition. We believe this competition reinforces the value and usefulness of the technical skills that students are gaining. A prize will go to the individual (or individuals) whose project demonstrates technical mastery, is most creative in its implementation, and has the most potential to improve society. This initiative has already spawned collaborative-learning groups.

To increase the impact of these projects and further reinforce the real-world applicability of the knowledge students gain through this course series, we are working with the Arribada Initiative [93] to create larger advised projects. This partnership will allow students to contribute their newly acquired TinyML skills to real-world conservation efforts, such as human/elephant conflict mitigation and sea-turtle monitoring, while receiving advice and support from both industry professionals and course staff. Finally, we are asking the community to continuously improve the course, since it is as

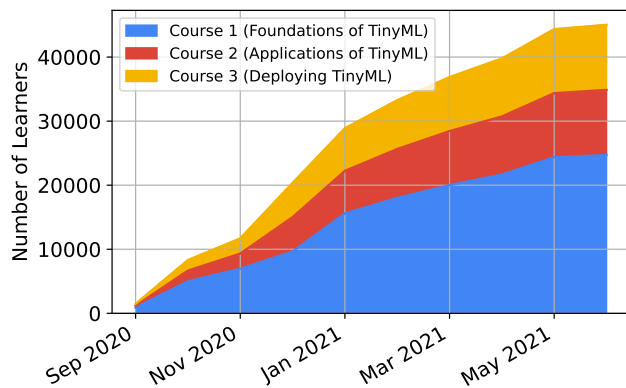


Figure 10: Course-enrollment metrics for Courses 1 through 3. Over 43,000 students are currently enrolled across all three courses. Course 4 has yet to open for enrollment. We expect another enrollment spike with “Scaling TinyML.”

much theirs as ours. As a result, we’ve seen many forum posts and GitHub pull requests offering typo corrections, bug fixes, and even content-improvement suggestions.

7 BROADER IMPACT

Our goal is to expand global access to applied ML through the lens of TinyML. In this section, we assess our work’s initial impact by presenting data from edX Insights, a service that provides course statistics to instructors and staff. It is merely an initial impact assessment, as the first cohort of participants have just begun graduating from the core TinyML series (Courses 1-3), and Course 4 (optional) remains in development. As such, our early analysis considers enrollment in the first three courses by geography, background, age, and gender.

7.1 Course Enrollment

At the time of this writing, the total course enrollment stands at 43,000. Figure 10 shows the daily enrollment data, starting from the opening date. We announced Courses 1, 2, and 3 together in early October 2020 and launched them on October 27, 2020; December 16, 2020; and March 2, 2021, respectively. Students could enroll in any or all courses at the same time but could only start after each course’s launch date.

TinyML is a young field, so the first useful metric is interest in the topic (i.e., acquiring applied-ML skills via TinyML). Figure 10 shows the strong and steady increase over time. On average, ~1,000 new students enroll in at least one course each week. Interest in Courses 2 and 3 continues to grow—a phenomenon we attribute to participants promoting them through social media such as LinkedIn, Twitter, and Facebook as they earn their course-completion certificates. The sharp increases around the first week of October, third week of December, and third week of February align with course-announcement dates or major social-media activity. For instance, on January 24, Mashable handpicked “Fundamentals of TinyML” as

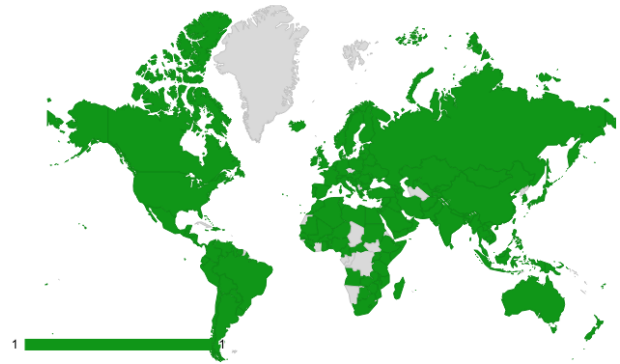


Figure 11: Global access to TinyML courses. At the time of this writing, people from more than 171 of the 193 United Nations member states have participated in TinyML.

one of the 10 best free Harvard courses to learn something new [94]. TinyML ranked at the top of the STEM-courses listed.

Figure 11 shows that TinyML students come from more than 171 countries. Because edX reaches a wide audience, our learners come from nearly all continents. Today, the top 10 countries by participant activity are the US, India, Turkey, the UK, Canada, Pakistan, Germany, Brazil, Australia, and Indonesia.

7.2 Completion Rates

People take online courses for a wide variety of reasons. Some are curious about the topic and want to get their feet wet; they may audit a course but not complete it. Others would like to master the program and earning a certificate of completion, assuming they can afford it. Therefore, enrollment numbers alone are insufficient.

We assessed how many verified enrollees complete the course. We have access only to the percentage who have earned a passing grade among those officially enrolled in the courses (i.e., the paid-certificate program). This number is constantly changing. At the time of writing, the completion rates are 59%, 55%, and 44% for Courses 1, 2, and 3, respectively. We believe Course 3’s number is slightly lower because it is more challenging than Courses 1 and 2, which do not have a hands-on component. The average completion rate for most MOOCs is somewhere between 5% and 15% [95], so the TinyML courses appear to be faring well. Although these results are preliminary (we need more data to make better quantitative comparisons), they shed a positive light on our design approach.

7.3 Learner Demographics

We conducted a demographic analysis of students’ age, educational background, and gender. They volunteer this information to edX, so it covers only a fraction of the numbers in Figure 10. Nonetheless, the data is extensive enough that we can draw general conclusions. At the start of each course, a forum post asks students to introduce themselves and summarize what they hope to get out of the edX series. We derived additional qualitative analysis from these responses. So far we have a good distribution across age groups and

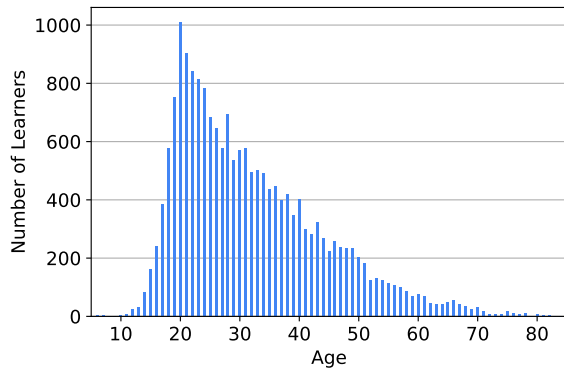


Figure 12: Age demographics across all courses based on voluntarily provide information. We have learners who are still in high school to learners who are retired and learning TinyML to understand its impact on our global society.

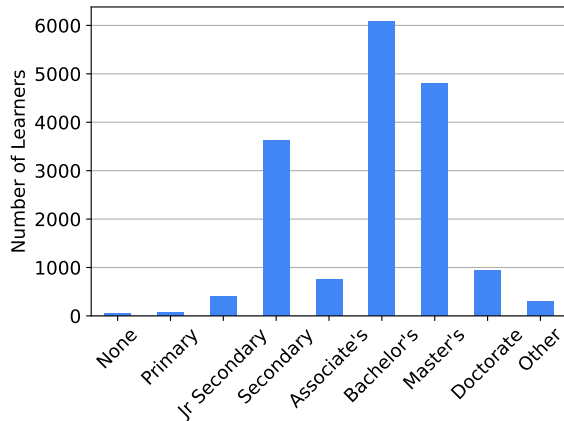


Figure 13: Education demographics based on voluntarily provide information. Many of our learners indicated an interest in TinyML to understand applied ML technologies to either pivot or grow further in their current positions.

educational backgrounds. Our gender diversity is lacking, however, but we are working to address it (Section 8).

Age. Figure 12 shows the age distribution for all three courses combined. The median is 30. Some participants are high-school students as young as 15 and wish to pursue an ML career. Others are over 60 and wish to understand the latest technological innovations as well as their societal implications. This age diversity was one of our objectives (Section 2).

Education. Figure 13 shows that nearly all our learners have either just a secondary (high-school) diploma or a bachelor’s/master’s degree. A few others have doctorate degrees. Judging from the forum discussions, we gather that individuals with a bachelor’s or master’s degree are trying to advance or shift their careers by adding an ML focus. Most participants with a doctorate want to apply (tiny)

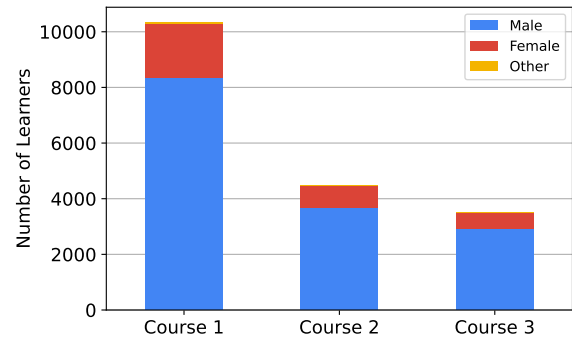


Figure 14: Gender demographics based on voluntarily provide information. Collectively, we are working on engaging a more diverse population of learners with the aid of working groups that are part of the tinyML Foundation.

ML in their research. Many students expressed enthusiasm about enrolling in a career-advancing course backed by both Harvard and Google. This variety of educational backgrounds and career focuses also meets our expectations and objectives and further emphasizes the importance of our academia/industry partnership (Section 2).

Gender. Figure 14 depicts the gender diversity across all three courses. It weighs heavily toward men; on average, across all three courses, 20% of our learners are women. We are working to change that ratio through our open education initiative (Section 8). More specifically, we are putting together a TinyML4Everyone working group to encourage more women to learn about TinyML.

8 FUTURE DIRECTIONS

TinyML can dramatically transform applied-ML education and development at many levels, far beyond what we achieved with the edX specialization. To this end, we launched the Tiny Machine Learning Open Education Initiative (TinyMLx) [96] to sponsor a wide variety of initiatives, such as TinyML4D (for applied-ML education via TinyML for developing countries), TinyML4STEM (for nurturing creative research in science, technology, engineering, and math), TinyML4Everyone (for building a shared identity and breaking stereotypes), and TinyML4x (for your favorite topic x).

We are currently running the TinyML4D and TinyML4Everyone working groups that are looking for ways to broaden TinyML participation, access, and belonging. One way is to provide TinyML materials in a student’s native language. For instance, we already have two projects for developing course content and instructional materials in Spanish and Portuguese [97]. Additionally, the makers of TinyML on edX, along with students and faculty of Navajo Technical University in New Mexico, plan to conduct a workshop in June 2021 that teaches Navajo students the basics of hardware programming and how to employ ML for their communities by creating voice-activated applications trained on the Navajo language.

TinyML can be instrumental for inspiring youth, as it offers a superb introduction to programming and ML for K–12 students.

Deployment of TinyML applications on physical embedded devices intrigues students by allowing them to interact with actual technologies, not just on-screen representations. Our first step in this direction was to publicly release all course materials on our GitHub: <https://github.com/tinyMLx/courseware>. We are working with STEM teachers worldwide to help us refine our tools for the classroom. Our aim is to develop ready-to-go project-based lessons and accompanying lesson plans to further increase ML access by reaching younger children. One possible project is to enable the use of visual programming abstractions (e.g., Microsoft MakeCode editor for the Arduino Nano BLE Sense 33) so people of all ages can apply ML without learning a programming language.

In addition, we are working with various organizations to assist teachers in learning applied ML. The 2021 Backyard Brains AI Fellowship [98], for example, is an early opportunity for teachers to help design TinyML projects for classrooms.

9 LIMITATIONS

We believe TinyML is an effective means to widen access to applied ML. Indeed, it is one way but not the only way. To provide a more balanced viewpoint, we describe some limitations of our approach and suggest alternative methods that may be more suitable.

Hardware cost. TinyML requires the purchase of embedded hardware to acquire the full-stack ML-development experience. The TinyML kit we developed costs US\$49.99. In some developing countries, this exceeds the average income in a week, in some rare cases, even a month. Although this price is considered reasonable in some countries, it may still be too high in others. We have found that the cost of shipping to distant parts of the world depends heavily on the presence of nearby distribution centers that carry the device. If none exist, the kit's cost, including shipping, can sometimes double the original kit price.

Ideally, TinyML would require no physical hardware, making the hardware cost zero. We are experimenting with open-source emulation platforms such as Renode.io from Antmicro [99]. Renode is an open-source framework that allows users to build and test embedded (ML) software without physical embedded hardware. It will enable developers to run their original code, which would have run on the hardware, unmodified in an emulated environment. Although this approach eliminates the hardware cost, students miss the opportunity and excitement of interacting with a device.

Device accessibility. Globally, the number of embedded devices far exceeds the number of cloud and mobile devices (as Figure 3 shows). But individuals must procure the necessary embedded hardware, such as the TinyML kit that we have developed with Arduino, to learn. By comparison, devices such as laptop and desktop computers connected to the web benefit from easier access. Students can use a regular computer to gain access to the online course materials. Even if they lack immediate access to computers in their homes, they can access the online resources from Internet cafés that provide web access for a nominal fee. A crucial shortcoming of this approach, however, is that learners will have difficulty experiencing the complete ML workflow (Figure 4), since they will be unable to deploy in a device the models they train in the cloud.

Smartphones may be a suitable compromise. They are highly accessible, even though they can be an order of magnitude more costly than the TinyML kit. Nevertheless, they enable students to experience the complete TinyML design, development, deployment, and management workflow. Also, an average smartphone has more than 10 sensors—many more than the Arduino Nano 33 BLE Sense we use in Course 3, enabling additional applications. Learners can hold the smartphone in their hands, much like the TinyML device. That said, conveying the significance of ML's future being tiny and bright (Section 3) is more challenging (though not impossible) because mobile devices have far more resources (compute power, memory, bandwidth, etc.) than TinyML devices (Table 1). Students may therefore miss the fundamental issue of embedded-resource constraints. But if the goal is ultimately to expand access to applied ML, mobile devices may be a fair compromise.

Programming background. Building ML models for mobile devices (using TensorFlow Lite [100]) or the web (using TensorFlow.js [101]) is possible using high-level programming languages such as Python and JavaScript, respectively. These languages are easy to learn and far more accessible to beginners than C/C++, which is necessary to program embedded hardware (similar to Course 3). So although TinyML creates an opportunity to showcase the full-stack ML experience using embedded hardware, and we leverage the Arduino IDE and heavily scaffolded code with video walkthroughs to minimize the lift to C/C++, it may also narrow access in some regards. The additional necessary programming skills and associated education can be a roadblock.

In the future, we believe that end-to-end developer platforms such as Edge Impulse [102, 103] that lower the entry barrier into TinyML will likely become mainstream and an essential part of the future developer ecosystem. Not every embedded ML engineer must know and understand all of the inner workings of TensorFlow Lite Micro or how an ML compiler works or how to extract the best performance from a highly customized ML hardware accelerator etc. Instead, learners need the right level of abstraction that allows them to focus on what matters most. Platforms such as Edge Impulse make it easy for learners, software developers, engineers and other domain experts to solve real-world problems using ML on the edge and TinyML devices without advanced degrees in ML or embedded systems. We therefore expose learners to the end-to-end MLOps platforms in Course 4, but note that more focus on such platforms in future courses could enable even more accessibility.

In summary, there are many paths to broaden applied-ML access. The correct approach—or, better, the most suitable approach—depends on the situation. We, therefore, hope this discussion clarifies the pros and cons of approaching applied ML through TinyML.

10 CONCLUSION

Expanding access to high-quality educational content, especially for machine learning, is important to ensuring that expertise diffuses beyond just a few prominent organizations. But doing so in a way that is both accessible and affordable to many different people is a difficult task. The four-part TinyML edX series we present here aims to tackle these challenges by providing application-driven

content that covers the entire ML life cycle, giving students hands-on experience guided by world experts and developing their ML skills regardless of their background. The forums, chats, optional project, and online discussions with the class creators promote community development and continued learning. The early impact of this approach is demonstrable: numerous participants from a variety of locations and demographics have signed up. We have also begun initiatives to further increase access by helping develop courses that target K–12 students and teachers, as well as courses in other languages.

ACKNOWLEDGEMENTS

Our approach to broadening access to applied machine learning using TinyML is based on input from many individuals at various organizations. We thank Rod Crawford, Tomas Edsö, and Felix Thomasmathibalan from **Arm**; Joe Lynas and Sacha Krstulović from **Audio Analytic**; Joshua Meyer from **Coqui**; Adam Benzio, Jenny Plunkett, Daniel Situnayake and Zach Shelby from **Edge Impulse**; Tulsee Doshi, Josh Gordon, Alex Gruenstein, Prateek Jain, and Nat Jeffries from **Google**; Marco Zennaro from **International Centre for Theoretical Physics (ICTP)**; Sek Chai from **Latent.AI**; Jane Polak Scowcroft from **Mozilla Common Voice**; Thiery Moreau from **OctoML**; Evgeni Gousev and Erich Plondke from **Qualcomm**; and Danilo Pau from **STMicroelectronics** for their valuable feedback. We are also grateful to the **Google TensorFlow Lite Micro team**, which includes Robert David, Jared Duke, Advait Jain, Vijay Janapa Reddi, Nat Jeffries, Jian Li, Nick Kreeger, Ian Nappier, Meghna Natraj, Shlomi Regev, Rocky Rhodes, and Tiezhen Wang, without whom we would have been unable to deploy models in microcontrollers, and **Arduino**—Jose Garcia Dotel and Martino Facchin—who helped us with global distribution of the TinyML kit. We also thank the tinyML Foundation for nurturing activity around embedded ML, providing guidance and supporting educational and outreach activities around TinyML.

REFERENCES

- [1] Tom M Mitchell, Rich Caruana, Dayne Freitag, John McDermott, David Zabowski, et al. Experience with a learning personal assistant. *Communications of the ACM*, 37(7), 1994.
- [2] Alexander Maedche, Christine Legner, Alexander Benlian, Benedikt Berger, Henner Gimpel, Thomas Hess, Oliver Hinz, Stefan Morana, and Matthias Söllner. AI-based digital assistants. *Business & Information Systems Engineering*, 61(4), 2019.
- [3] Ürün Dogan, Johann Edelbrunner, and Ioannis Iossifidis. Autonomous driving: A comparison of machine learning techniques by means of the prediction of lane change behavior. In *2011 IEEE International Conference on Robotics and Biomimetics*. IEEE, 2011.
- [4] Chenyi Chen, Ari Seff, Alain Kornhauser, and Jianxiong Xiao. Deepdriving: Learning affordance for direct perception in autonomous driving. In *Proceedings of the IEEE international conference on computer vision*, 2015.
- [5] Andy Zeng, Shuran Song, Johnny Lee, Alberto Rodriguez, and Thomas A. Funkhouser. TossingBot: Learning to Throw Arbitrary Objects with Residual Physics. *CoRR*, abs/1903.11239, 2019.
- [6] Muhammad Aurangzeb Ahmad, Carly Eckert, and Ankur Teredesai. Interpretable machine learning in healthcare. In *Proceedings of the 2018 ACM international conference on bioinformatics, computational biology, and health informatics*, 2018.
- [7] Arash Jahangiri and Hesham A Rakha. Applying machine learning techniques to transportation mode recognition using mobile phone sensor data. *IEEE transactions on intelligent transportation systems*, 16(5), 2015.
- [8] Fotios Zantalis, Grigorios Koulouras, Sotiris Karabetsos, and Dionisis Kandris. A review of machine learning and IoT in smart transportation. *Future Internet*, 11(4), 2019.
- [9] Anna L Buczak and Erhan Guven. A survey of data mining and machine learning methods for cyber security intrusion detection. *IEEE Communications surveys & tutorials*, 18(2), 2015.
- [10] Vijaya B Kolachalama and Priya S Garg. Machine learning and medical education. *NPJ digital medicine*, 1(1), 2018.
- [11] Hadeel S Alenezi and Maha H Faisal. Utilizing crowdsourcing and machine learning in education: Literature review. *Education and Information Technologies*, 2020.
- [12] Terry Brown. The AI Skills Shortage. <https://itchronicles.com/artificial-intelligence/the-ai-skills-shortage/>, October 2019.
- [13] Jean-François Gagné. Global AI Talent Report 2019. <https://jfgagne.ai/talent-2019/>.
- [14] Ronald M Harden. What is a spiral curriculum? *Medical teacher*, 21(2), 1999.
- [15] Yoram Neumann, Edith Neumann, and Shelia Lewis. The robust learning model with a spiral curriculum: Implications for the educational effectiveness of online master degree programs. *Contemporary Issues in Education Research*, 10(2), 2017.
- [16] Gleb Chuvpilo. Who's ahead in AI research AT NeurIPS 2020? <https://chuvpilo.medium.com/whos-ahead-in-ai-research-at-neurips-2020-bf2a40a54325>, December 2020.
- [17] Stack Overflow Developer Survey 2020. <https://insights.stackoverflow.com/survey/2020#developer-roles>.
- [18] IC Insights. *MCUs Expected to Make Modest Comeback after 2020 Drop*, 2020.
- [19] Yundong Zhang, Naveen Suda, Liangzhen Lai, and Vikas Chandra. Hello edge: Keyword spotting on microcontrollers. *arXiv preprint arXiv:1711.07128*, 2017.
- [20] Seungwoo Choi, Seokjun Seo, Beomjun Shin, Hyeongmin Byun, Martin Kersner, Beomsu Kim, Dongyoung Kim, and Sungjoo Ha. Temporal convolution for real-time keyword spotting on mobile devices. *arXiv preprint arXiv:1904.03814*, 2019.
- [21] Aravind Kota Gopalakrishna, Tanir Özçelebi, Antonio Liotta, and Johan J Lukkien. Exploiting machine learning for intelligent room lighting applications. In *2012 6th IEEE International Conference Intelligent Systems*. IEEE, 2012.
- [22] Clement Duhart, Gershon Dublon, Brian Mayton, Glorianna Davenport, and Joseph A Paradiso. Deep learning for wildlife conservation and restoration efforts. In *36th International Conference on Machine Learning, Long Beach*, volume 5, 2019.
- [23] Enrico Di Minin, Christoph Fink, Henrikki Tenkanen, and Tuomo Hiippala. Machine learning for tracking illegal wildlife trade on social media. *Nature ecology & evolution*, 2(3), 2018.
- [24] Varun Chandola, Arindam Banerjee, and Vipin Kumar. Anomaly detection: A survey. *ACM computing surveys (CSUR)*, 41(3), 2009.
- [25] Alan Turnbull, James Carroll, and Alasdair McDonald. Combining SCADA and vibration data into a single anomaly detection model to predict wind turbine component failure. *Wind Energy*, 24(3), 2021.
- [26] XJ Zeng, M Yang, and YF Bo. Gearbox oil temperature anomaly detection for wind turbine based on sparse Bayesian probability estimation. *International Journal of Electrical Power & Energy Systems*, 123, 2020.
- [27] Sudha Gupta, Faruk Kazi, Sushama Wagh, and Ruta Kambli. Neural network based early warning system for an emerging blackout in smart grid power networks. In *Intelligent Distributed Computing*. Springer, 2015.
- [28] Takehisa Yairi, Yoshinobu Kawahara, Ryohei Fujimaki, Yuichi Sato, and Kazuo Machida. Telemetry-mining: a machine learning approach to anomaly detection and fault diagnosis for space systems. In *2nd IEEE International Conference on Space Mission Challenges for Information Technology (SMC-IT'06)*. IEEE, 2006.
- [29] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*. Ieee, 2009.
- [30] Ekaba Bisong. Google colabatory. In *Building Machine Learning and Deep Learning Models on Google Cloud Platform*. Springer, 2019.
- [31] Robert David, Jared Duke, Advait Jain, Vijay Janapa Reddi, Nat Jeffries, Jian Li, Nick Kreeger, Ian Nappier, Meghna Natraj, Shlomi Regev, et al. Tensorflow lite micro: Embedded machine learning on tinyml systems. *arXiv preprint arXiv:2010.08678*, 2020.
- [32] Kalvin Bahia and Stefano Suardi. The state of mobile internet connectivity 2920. *GSMA Connected Society: London*, 2020.
- [33] Joseph Johnson. Global digital population as of January 2021. <https://www.statista.com/statistics/617136/digital-population-worldwide/#:~:text=Global%20internet%20usage&text=The%20global%20internet%20penetration%20rate,penetration%20rate%20among%20the%20population,> January 2021.
- [34] Laura Silver. Smartphone Ownership Is Growing Rapidly Around the World, but Not Always Equally. <https://www.pewresearch.org/global/2019/02/05/smartphone-ownership-is-growing-rapidly-around-the-world-but-not->

- always-equally/, August 2020.
- [35] Dazhi Yang. Instructional strategies and course design for teaching statistics online: perspectives from online students. *International Journal of STEM Education*, 4(1), 2017.
- [36] Nadia Rahbek Dyrberg and Henriette Tolstrup Holmegaard. Motivational patterns in STEM education: a self-determination perspective on first year courses. *Research in Science & Technological Education*, 37(1), 2019.
- [37] Claire Wladis, Alyse C Hachey, and Katherine Conway. An investigation of course-level factors as predictors of online STEM course outcomes. *Computers & Education*, 77, 2014.
- [38] E Oran Brigham. *The fast Fourier transform and its applications*. Prentice-Hall, Inc., 1988.
- [39] Vibha Tiwari. MFCC and its applications in speaker recognition. *International journal on emerging technologies*, 1(1), 2010.
- [40] Alexander Gruenstein, Raziel Alvarez, Chris Thornton, and Mohammadali Ghodrati. A cascade architecture for keyword spotting on mobile devices. *arXiv preprint arXiv:1712.03603*, 2017.
- [41] Jason Yosinski, Jeff Clune, Yoshua Bengio, and Hod Lipson. How transferable are features in deep neural networks? *arXiv preprint arXiv:1411.1792*, 2014.
- [42] Andrew G. Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand, Marco Andreetto, and Hartwig Adam. MobileNets: Efficient Convolutional Neural Networks for Mobile Vision Applications. *CoRR*, abs/1704.04861, 2017.
- [43] Gongde Guo, Hui Wang, David Bell, Yaxin Bi, and Kieran Greer. KNN model-based approach in classification. In *OTM Confederated International Conferences "On the Move to Meaningful Internet Systems"*. Springer, 2003.
- [44] Pierre Baldi. Autoencoders, unsupervised learning, and deep architectures. In *Proceedings of ICML workshop on unsupervised and transfer learning*. JMLR Workshop and Conference Proceedings, 2012.
- [45] Arduino Nano 33 BLE. <https://store.arduino.cc/usa/nano-33-ble>. (Accessed on 04/02/2021).
- [46] OV7675. <https://www.arducam.com/products/camera-breakout-board/0-3mp-ov7675/>. (Accessed on 04/02/2021).
- [47] Joseph S Krajcik and Phyllis C Blumenfeld. *Project-based learning*. na, 2006.
- [48] Petri Vesikivi, Minna Lakkala, Jaana Holvikivi, and Hanni Muukkonen. The impact of project-based learning curriculum on first-year retention, study experiences, and knowledge work competence. *Research Papers in Education*, 35(1), 2020.
- [49] TensorFlow. TinyConv. https://github.com/tensorflow/tensorflow/blob/master/tensorflow/examples/speech_commands/models.py, 2021. [Online; accessed 15-Mar-2021].
- [50] Trevor Martin. *The designer's guide to the Cortex-M processor family*. Newnes, 2016.
- [51] Colby R Banbury, Vijay Janapa Reddi, Max Lam, William Fu, Amin Fazel, Jeremy Holleman, Xinyuan Huang, Robert Hurtado, David Kanter, Anton Lokhmotov, et al. Benchmarking TinyML systems: Challenges and direction. *arXiv preprint arXiv:2003.04821*, 2020.
- [52] Cloud automl custom machine learning models | google cloud. <https://cloud.google.com/automl>. (Accessed on 05/30/2021).
- [53] Barret Zoph, Vijay Vasudevan, Jonathon Shlens, and Quoc V Le. Learning transferable architectures for scalable image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 8697–8710, 2018.
- [54] Hiroshi Doyu, Roberto Morabito, and Martina Brachmann. A tinyml ecosystem for machine learning in iot: Overview and research challenges. In *2021 International Symposium on VLSI Design, Automation and Test (VLSI-DAT)*, pages 1–5. IEEE, 2021.
- [55] Jakub Konečný, H Brendan McMahan, Felix X Yu, Peter Richtárik, Ananda Theertha Suresh, and Dave Bacon. Federated learning: Strategies for improving communication efficiency. *arXiv preprint arXiv:1610.05492*, 2016.
- [56] Harold Pashler, Mark McDaniel, Doug Rohrer, and Robert Bjork. Learning styles: Concepts and evidence. *Psychological science in the public interest*, 9(3), 2008.
- [57] Alison S Lockman and Barbara R Schirmer. Online Instruction in Higher Education: Promising, Research-Based, and Evidence-Based Practices. *Journal of Education and E-Learning Research*, 7(2), 2020.
- [58] Philip J Guo, Juho Kim, and Rob Rubin. How video production affects student engagement: An empirical study of MOOC videos. In *Proceedings of the first ACM conference on Learning scale conference*, 2014.
- [59] H Frank Cervone. Applied digital library project management: Using Pugh matrix analysis in complex decision-making situations. *OCLC Systems & Services: International digital library perspectives*, 2009.
- [60] Espressif Systems. ESP32. <https://www.espressif.com/en/products/socs/esp32>.
- [61] Peter Jamieson. Arduino for teaching embedded systems. are computer scientists and engineering educators missing the boat? In *Proceedings of the international conference on frontiers in education: computer science and computer engineering (FECS)*. The Steering Committee of The World Congress in Computer Science, Computer ... , 2011.
- [62] nRF52840. <https://www.nordicsemi.com/-/media/Software-and-other-downloads/Product-Briefs/nRF52840-SoC-product-brief.pdf?la=en&hash=EDF4C48A053E7943AD3C9DD3963B626D768B5885>. (Accessed on 04/02/2021).
- [63] Mbed OS. <https://os.mbed.com/mbed-os/>. (Accessed on 04/02/2021).
- [64] Thomas W Schubert, Alessandro D'Ausilio, and Rosario Canto. Using Arduino microcontroller boards to measure response latencies. *Behavior research methods*, 45(4), 2013.
- [65] Pico4ML by Arducam. <https://www.arducam.com/pico4ml-an-rp2040-based-platform-for-tiny-machine-learning/>. (Accessed on 04/02/2021).
- [66] Barbara J. Grosz, David Gray Grant, Kate Vredenburg, Jeff Behrends, Lily Hu, Alison Simmons, and Jim Waldo. Embedded ethicsCS: Integrating ethics across CS curriculum. *Communications of the ACM*, 62(8), 2019. <https://cacm.acm.org/magazines/2019/8/238345-embedded-ethics/fulltext>.
- [67] Aaron Smith. Public Attitudes Toward Technology Companies. <https://www.pewresearch.org/internet/2018/06/28/public-attitudes-toward-technology-companies/>.
- [68] Kathleen C. Fraser, J. Meltzer, and F. Rudzicz. Linguistic Features Identify Alzheimer's Disease in Narrative Speech. *Journal of Alzheimer's disease : JAD*, 49 2, 2016.
- [69] Pete Warden. Speech commands: A dataset for limited-vocabulary speech recognition. *arXiv preprint arXiv:1804.03209*, 2018.
- [70] Rosana Ardia, Megan Branson, Kelly Davis, Michael Henretty, Michael Kohler, Josh Meyer, Reuben Morais, Lindsay Saunders, Francis M. Tyers, and Gregor Weber. Common Voice: A Massively-Multilingual Speech Corpus, 2020.
- [71] Aakanksha Chowdhery, Pete Warden, Jonathon Shlens, Andrew Howard, and Rocky Rhodes. Visual wake words dataset. *arXiv preprint arXiv:1906.05721*, 2019.
- [72] Joy Buolamwini and Timnit Gebru. Gender Shades: Intersectional Accuracy Disparities in Commercial Gender Classification. In Sorelle A. Friedler and Christo Wilson, editors, *Proceedings of the 1st Conference on Fairness, Accountability and Transparency*, volume 81 of *Proceedings of Machine Learning Research*, New York, NY, USA, 23–24 Feb 2018. PMLR.
- [73] Discourse Discussion Forum: How will the lack of female voices be handled. <https://discourse.mozilla.org/t/how-will-the-lack-of-female-voices-be-handled/40551>.
- [74] Allison Koencke, Andrew Nam, Emily Lake, Joe Nudell, Minnie Quartey, Zion Mengesha, Connor Toups, John R. Rickford, Dan Jurafsky, and Sharad Goel. Racial disparities in automated speech recognition. *Proceedings of the National Academy of Sciences*, 117(14), 2020.
- [75] Yunhan Wu, Daniel Rough, Anna Bleakley, Justin Edwards, Orla Cooney, Philip R. Doyle, Leigh Clark, and Benjamin R. Cowan. See What I'm Saying? Comparing Intelligent Personal Assistant Use for Native and Non-Native Language Speakers. In *22nd International Conference on Human-Computer Interaction with Mobile Devices and Services*, MobileHCI '20, New York, NY, USA, 2020. Association for Computing Machinery.
- [76] Joel Shor, Dotan Emanuel, Oran Lang, Omry Tuval, Michael Brenner, Julie Cattiau, Fernando Vieira, Maeve McNally, Taylor Charbonneau, Melissa Nollstadt, Avinatan Hassidim, and Yossi Matias. Personalizing ASR for Dysarthric and Accented Speech with Limited Data, 2019.
- [77] Google. <https://pair-code.github.io/what-if-tool/>, 2020.
- [78] Google. Responsible AI Toolkit & TensorFlow. <https://www.tensorflow.org/responsible%5Fai>, 2020.
- [79] Kartik Prabhu, Brian Jun, Pan Hu, Zain Asgar, Sachin Katti, and Pete Warden. Privacy-Preserving Inference on the Edge: Mitigating a New Threat Model. In *Research Symposium on Tiny Machine Learning*, 2021.
- [80] Drew Harwell. Doorbell-camera firm Ring has partnered with 400 police forces, extending surveillance concerns. <https://www.washingtonpost.com/technology/2019/08/28/doorbell-camera-firm-ring-has-partnered-with-police-forces-extending-surveillance-reach/>.
- [81] Zack Whittaker. Fitness app PumpUp leaked health data, private messages. <https://www.zdnet.com/article/fitness-app-pumpup-leaked-health-data-private-messages/>.
- [82] Guoming Zhang, Chen Yan, Xiaoyu Ji, Tianchen Zhang, Taimin Zhang, and Wenyuan Xu. Dolphinattack: Inaudible voice commands. In *Proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security*, 2017.
- [83] Benjamin Reinheimer, Lukas Aldag, Peter Mayer, Mattia Mossano, Reyhan Duezgen, Bettina Lofthouse, Tatiana von Landesberger, and Melanie Volkamer. An investigation of phishing awareness and education over time: When and how to best remind users. In *Sixteenth Symposium on Usable Privacy and Security (SUSOUPS) 2020*, 2020.
- [84] Margaret Mitchell, Simone Wu, Andrew Zaldivar, Parker Barnes, Lucy Vasseran, Ben Hutchinson, Elena Spitzer, Inioluwa Deborah Raji, and Timnit Gebru. Model Cards for Model Reporting. *Proceedings of the Conference on Fairness*,

- Accountability, and Transparency*, January 2019.
- [85] Paul Belleflamme and Julien Jacqmin. An Economic Appraisal of MOOC Platforms: Business Models and Impacts on Higher Education. *CESifo Economic Studies*, 62, 09 2014.
- [86] Arduino Tiny Machine Learning Kit. <https://store.arduino.cc/usa/tiny-machine-learning-kit>, April 2021.
- [87] Chih-Ming Chen and Chung-Hsin Wu. Effects of different video lecture types on sustained attention, emotion, cognitive load, and learning performance. *Computers & Education*, 80, 2015.
- [88] Richard E Mayer, Logan Fiorella, and Andrew Stull. Five ways to increase the effectiveness of instructional video. *Educational Technology Research and Development*, 68(3), 2020.
- [89] Richard Mayer and Richard E Mayer. *The Cambridge handbook of multimedia learning*. Cambridge university press, 2005.
- [90] Jiahui Wang, Pavlo Antonenko, and Kara Dawson. Does visual attention to the instructor in online video affect learning and learner perceptions? An eye-tracking analysis. *Computers & Education*, 146, 2020.
- [91] Rain Dance Canada. Audio Quality vs. Video Quality. <https://www.youtube.com/watch?v=-PLMiA18tBc>.
- [92] Frame.io. Video Review and Collaboration Software. <https://www.frame.io/>.
- [93] Arribada Initiative: Open Source Conservation Technology. <https://arribada.org/>, March 2021.
- [94] Amy-Mae Turner. 10 free online classes from Harvard to learn something new. <https://mashable.com/article/free-harvard-classes-online/>, January 2021.
- [95] F Hollands and A Kazi. Benefits and Costs of MOOC-Based Alternative Credentials: 2018-2019 Results from End-of-Program Surveys. *Center for Benefit-Cost Studies of Education, Teachers College, Columbia University*, 2019.
- [96] Tiny Machine Learning Open Education Initiative (TinyMLx). <https://tinymlx.org/>.
- [97] Marcelo Rovai. TinyML - Machine Learning for Embedding Devices. <https://github.com/Mjrovai/UNIFEI-IESTI01-T01-2021.1>.
- [98] Backyard Brains 2021 AI Fellowship. <https://blog.backyardbrains.com/2021/03/backyard-brains-2021-ai-fellowship/>, March 2021.
- [99] Antmicro. Renode.io - Antmicro's virtual development framework for complex embedded systems. <https://renode.io/>.
- [100] Google. Tensorflow lite | ml for mobile and edge devices. <https://www.tensorflow.org/lite>. (Accessed on 06/02/2021).
- [101] Daniel Smilkov, Nikhil Thorat, Yannick Assogba, Ann Yuan, Nick Kreeger, Ping Yu, Kangyi Zhang, Shanqing Cai, Eric Nielsen, David Soergel, Stan Bileschi, Michael Terry, Charles Nicholson, Sandeep N. Gupta, Sarah Sirajuddin, D. Sculley, Rajat Monga, Greg Corrado, Fernanda B. Viégas, and Martin Wattenberg. Tensorflow.js: Machine learning for the web and beyond. *arXiv preprint arXiv:1901.05350*, 2019.
- [102] Edge impulse. <https://www.edgeimpulse.com/>. (Accessed on 06/02/2021).
- [103] Introduction to embedded machine learning | coursera. <https://www.coursera.org/learn/introduction-to-embedded-machine-learning>. (Accessed on 06/09/2021).