

CS249r: Generative AI



Nov 27

Course Logistics

Assignment Schedule Updates

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

Scribing

- Peer review - Generally, 1 detailed review
 - Security and Privacy is in PR mode
 - This week
 - i. Responsible AI
 - ii. Sustainable AI
 - Only 1 detailed review
 - i. 25/33 are done! 👍 🙌
- ES91r
 - Opportunity to work with me and my students in the Edge Computing Lab
 - Deepen your knowledge of TinyML and more broadly ML systems.

Project Check-Ins

- Any issues?
 - Debugging, profiling, optimization...
- Feel free to contact us on slack, so that we can be more responsive
- Office hours dedicated to projects
 - Please check in with the TAs if there are any issues or you want feedback on presentation/papers etc.

RE: Responsible AI Class

Discussion of sensitive topics

Mindful delivery of these topics

Open to ideas and suggestions

Course Topics

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

Today's Schedule

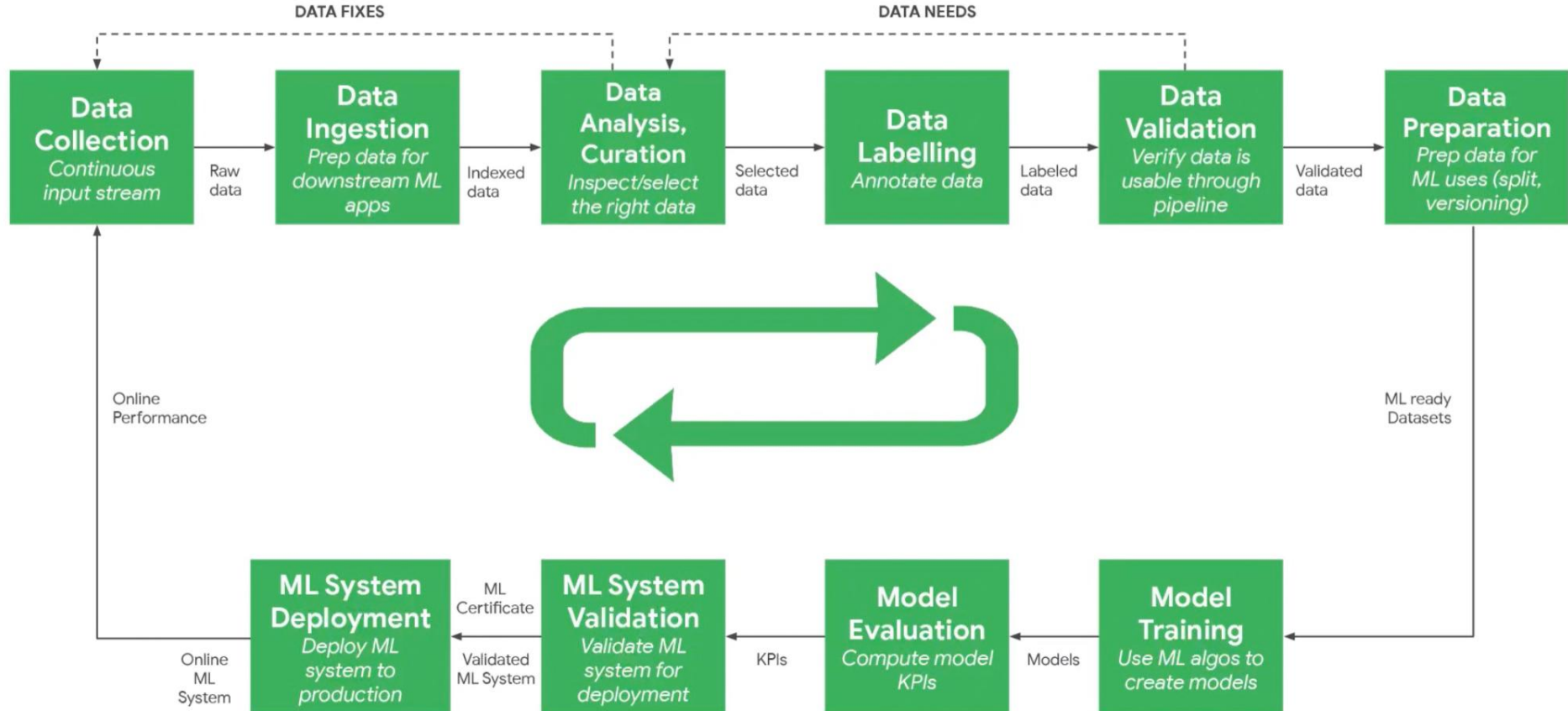
12:45 pm to ~2:00pm - Lecture material

2:00 pm to 2:30pm - Paper discussions

2:30 pm to 3:30pm - Jason Wei / Open AI

Course Logistics

Life cycle of ML



Recap: Data Engineering

Contributors: Oishi Banerjee, Shreya Johri, Itai Shapira

Understand the **importance of clearly defining the problem statement** and objectives

Recognize **various data sourcing techniques** like web scraping, crowdsourcing, and synthetic data generation, along with their advantages and limitations.

Appreciate the need for **thoughtful data labeling, using manual or AI-assisted approaches**, to create high-quality training datasets.

Methods for **storing & managing data** – databases, data warehouses, and data lakes.

Comprehend the role of **transparency through metadata and dataset documentation**

Understand how **licensing protocols** govern legal data access and usage

Recognize key **challenges in data engineering**, including privacy risks, representation gaps, legal restrictions around data access, and balancing competing priorities.

Models Cards for Model Reporting (paper)

Multilingual Spoken Words Corpus (paper)

The screenshot shows a web browser displaying the '6 Data Engineering' chapter from the 'MACHINE LEARNING SYSTEMS' book. The page layout includes a table of contents on the right, a list of learning objectives on the left, and a central diagram illustrating the data engineering process. The diagram shows a flow from 'Raw Data Sources' through 'Data Sourcing', 'Data Cleaning', 'Data Processing', 'Data Storage', and 'Data Labeling' to 'Data Analysis'. The 'Data Engineering' box is highlighted in the center. The 'Learning Objectives' section lists several key points:

- Understand the importance of clearly defining the problem statement and objectives when embarking on a ML project.
- Recognize various data sourcing techniques like web scraping, crowdsourcing, and synthetic data generation, along with their advantages and limitations.
- Appreciate the need for thoughtful data labeling, using manual or AI-assisted approaches, to create high-quality training datasets.
- Briefly learn different methods for storing and managing data such as databases, data warehouses, and data lakes.
- Comprehend the role of transparency through metadata and dataset documentation, as well as tracking data provenance to facilitate ethics, auditing, and reproducibility.
- Understand how licensing protocols govern legal data access and usage, necessitating careful compliance.

Recap: Embedded ML Frameworks

Contributors: Henry Bae, Divya Amirtharaj, Sophia Cho, Emeka Ezike

Understand the evolution and capabilities of **major machine learning frameworks**. This includes graph execution models, programming paradigms, hardware acceleration support, and how they have expanded over time.

Learn the **core components and functionality of frameworks** like computational graphs, data pipelines, optimization algorithms, training loops, etc. that enable efficient model building.

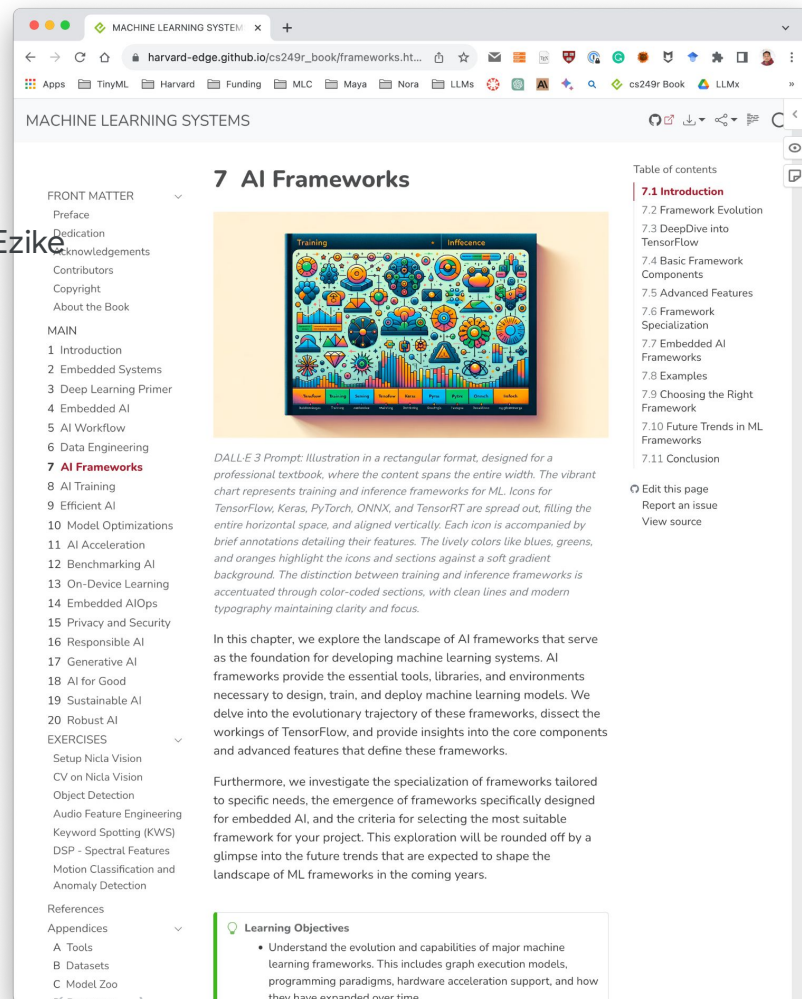
Compare **frameworks across different environments like cloud, edge, and tinyML**. Learn how frameworks specialize based on computational constraints and hardware.

Dive deeper into embedded and tinyML focused frameworks like **TensorFlow Lite Micro**, **CMSIS-NN**, **TinyEngine** etc. and how they optimize for microcontrollers.

Explore **model conversion and deployment considerations** when choosing a framework, including aspects like latency, memory usage, and hardware support.

Evaluate **key factors in selecting the right framework** like performance, hardware compatibility, community support, ease of use, etc. based on the specific project needs and constraints.

Understand the **limitations of current frameworks** and potential future trends like using ML to improve frameworks, decomposed ML systems, and high performance compilers.



TensorFlow Lite Micro: Embedded Machine Learning on TinyML Systems (paper)
MCUNet: Tiny Deep Learning on IoT Devices (paper)

Recap: Representation & Compression

Contributors: Jeffrey Ma, Aghyad Deeb, Costin Oncescu, Jayson Lin

Learn **techniques like pruning, knowledge distillation and specialized model architectures** to represent models more efficiently

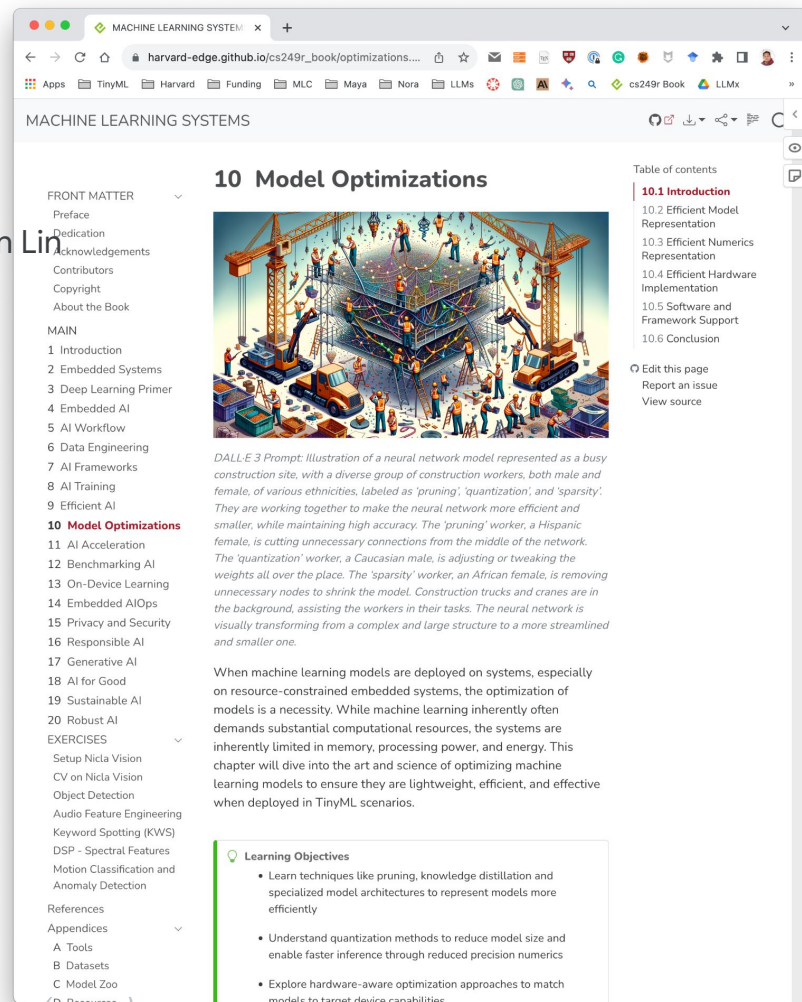
Understand **quantization methods** to reduce model size and enable faster inference through reduced precision numerics

Explore **hardware-aware optimization approaches** to match models to target device capabilities

Discover **software tools like frameworks and model conversion platforms** that enable deployment of optimized models

Develop **holistic thinking to balance tradeoffs in model complexity, accuracy, latency, power** etc. based on application requirements

Gain **strategic insight into selecting and applying model optimizations** based on use case constraints and hardware targets



A Survey of Quantization Methods for Efficient Neural Network Inference (paper)
The Lottery Ticket Hypothesis: Finding Sparse, Trainable Neural Networks (paper)

Recap: Benchmarking of ML Systems

Contributors: Colby Banbury, Mark Mazumder

Understand the **purpose and goals of benchmarking AI systems**, including performance assessment, resource evaluation, validation, and more.

Learn about the **different types of benchmarks** - micro, macro, and end-to-end - and their role in evaluating different aspects of an AI system.

Become familiar with the **key components of an AI benchmark**, including datasets, tasks, metrics, baselines, reproducibility rules, and more.

Understand the **distinction between training and inference**, and how each phase warrants specialized ML systems benchmarking.

Learn about **system benchmarking concepts** like throughput, latency, power, and computational efficiency.


Appreciate the **evolution of model benchmarking from accuracy to more holistic metrics** like fairness, robustness and real-world applicability.

Recognize the **growing role of data benchmarking** in evaluating issues like bias, noise, balance and diversity.

Understand the **limitations of evaluating models, data, and systems in isolation**, and the emerging need for integrated benchmarking.

MACHINE LEARNING SYSTEMS

12 Benchmarking AI



DALL-E 3 Prompt: Photo of a podium set against a tech-themed backdrop. On each tier of the podium, there are AI chips with intricate designs. The top chip has a gold medal hanging from it, the second one has a silver medal, and the third has a bronze medal. Banners with 'AI Olympics' are displayed prominently in the background.

Benchmarking is a critical part of developing and deploying machine learning systems, especially for tinyML applications. Benchmarks allow developers to measure and compare the performance of different model architectures, training procedures, and deployment strategies. This provides key insights into which approaches work best for the problem at hand and the constraints of the deployment environment.

This chapter will provide an overview of popular ML benchmarks, best practices for benchmarking, and how to use benchmarks to improve model development and system performance. It aims to provide developers with the right tools and knowledge to effectively benchmark and optimize their systems, especially for tinyML systems.

Learning Objectives

- Understand the purpose and goals of benchmarking AI systems, including performance assessment, resource evaluation, validation, and more.
- Learn about the different types of benchmarks - micro, macro, and end-to-end - and their role in evaluating different aspects of an AI system.
- Become familiar with the key components of an AI benchmark, including datasets, tasks, metrics, baselines, reproducibility rules, and more.

FRONT MATTER

12.1 Introduction

12.2 Historical Context

12.3 AI Benchmarks: System, Model, and Data

12.4 System Benchmarking

12.5 Model Benchmarking

12.6 Data Benchmarking

12.7 The Trifecta

12.8 Benchmarks for Emerging Technologies

12.9 Conclusion

14

Recap: Learning on the Edge

Contributors: Michael Schnebly, Alex Rodriguez, Aditi Raju, Jared Ni

Understand **on-device learning** and how it differs from cloud-based training

Recognize the **benefits and limitations of on-device learning**

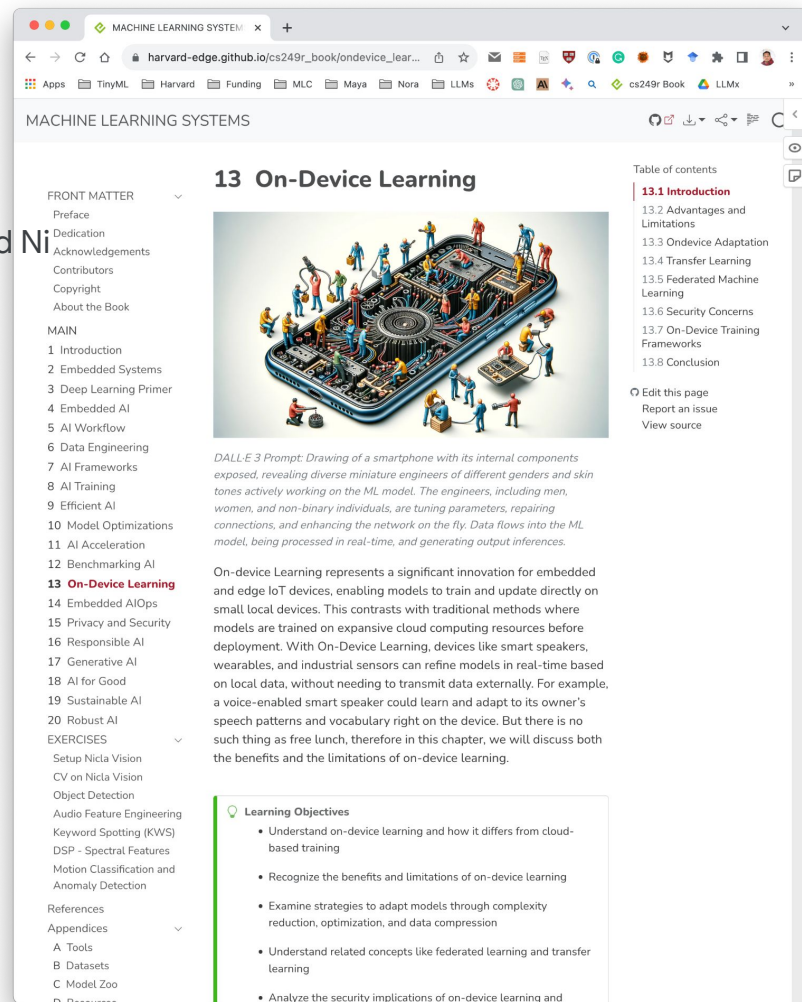
Examine **strategies to adapt models** through complexity reduction, optimization, and data compression

Understand **related concepts like federated learning and transfer learning**

Analyze the **security implications of on-device learning** and mitigation strategies

Flower: A Friendly Federated Learning Framework ([paper](#))

On-Device Training Under 256KB Memory ([paper](#))



Recap: AI Acceleration

Contributors: Jennifer Zhou, Eric Dong, Arnau Marin, Pong Trairatvorakul

Understand **why hardware acceleration is needed** for AI workloads

Survey key accelerator options like **GPUs, TPUs, FPGAs, and ASICs** and their tradeoffs

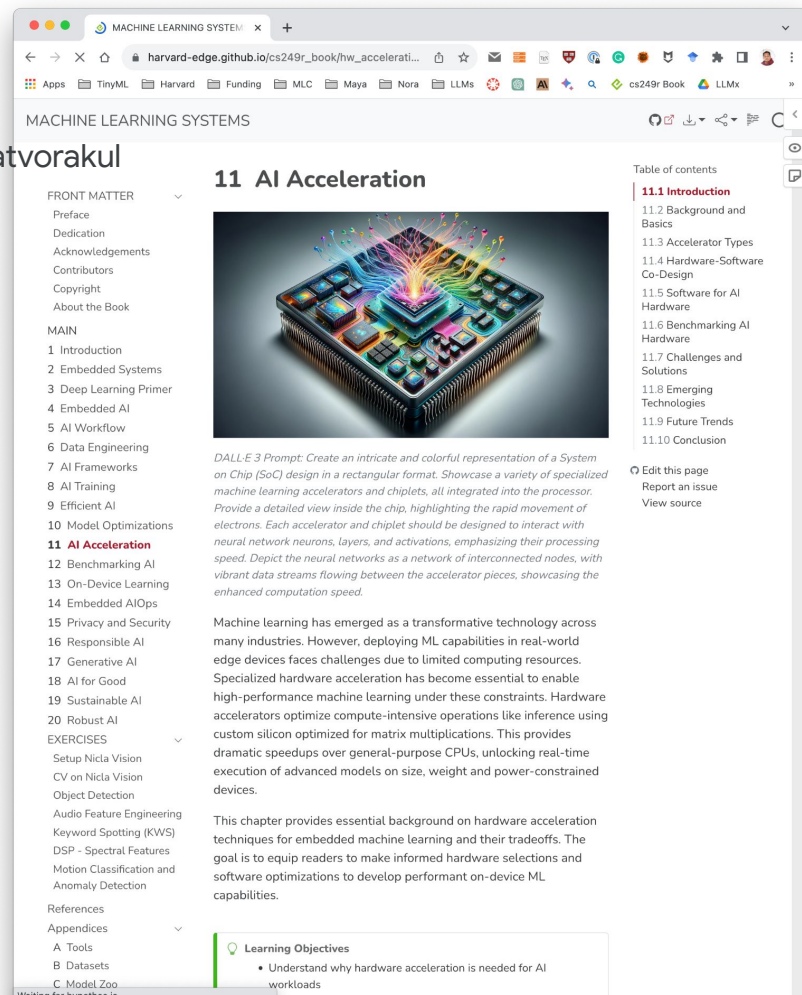
Learn about **programming models, frameworks, compilers** for AI accelerators

Appreciate the **importance of benchmarking and metrics** for hardware evaluation

Recognize the **role of hardware-software co-design** in building efficient systems

Gain exposure to cutting-edge **research directions like neuromorphics** and quantum computing

Understand **how ML is beginning to augment and enhance hardware design**



CFU Playground: Full-Stack Open-Source Framework for Tiny Machine Learning (tinyml) Acceleration on FPGAs (paper)
An Evaluation of Edge TPU Accelerators for Convolutional Neural Networks (paper)

Recap: MLOps

Contributors: Andrew Bass, Annie Landefeld, Vijay Edupuganti, Curren Iyer

Understand **what is MLOps** and why it is needed

Learn the **architectural patterns** for traditional MLOps

Contrast **traditional vs. embedded MLOps** across the ML lifecycle

Identify **key constraints of embedded environments**

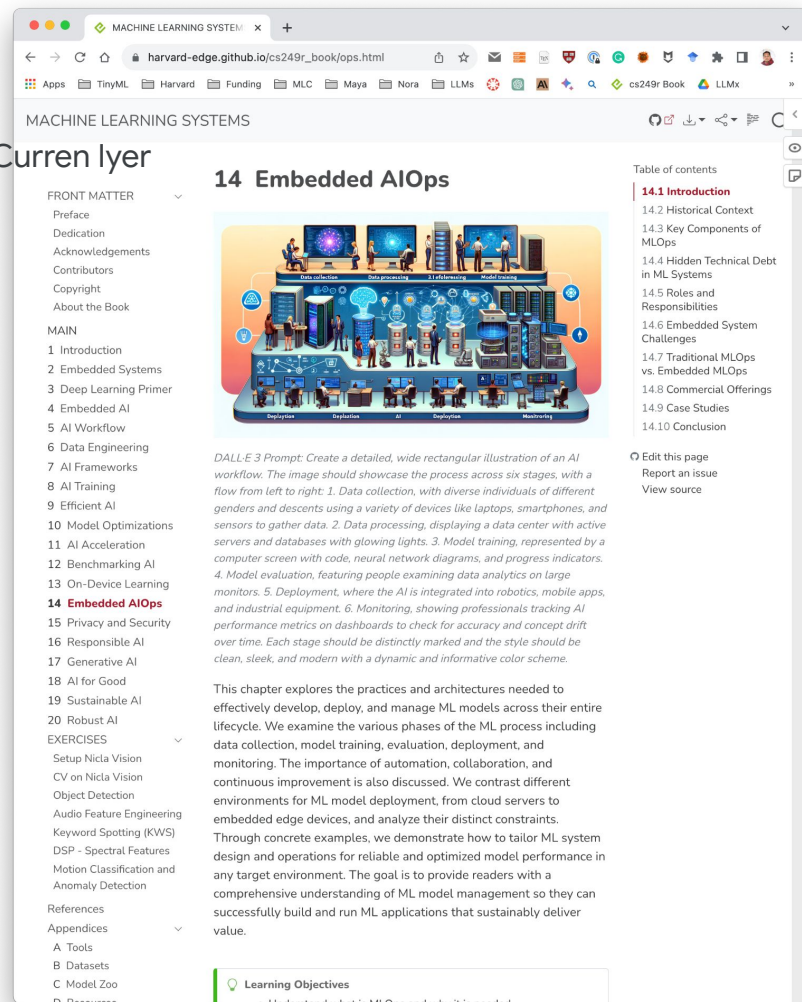
Learn **strategies to mitigate embedded ML challenges**

Examine **real-world case studies** demonstrating embedded MLOps principles

Appreciate **the need for holistic technical and human approaches**

Edge Impulse (paper)

Hidden Technical Debt in Machine Learning (paper)



Recap: Security & Privacy

Contributors: Elizabeth Sutor, Eliza Kimball, Jothi Ramaswamy, Elias Nuwara

Understand **key ML privacy and security risks** like data leaks, model theft, adversarial attacks, bias, and unintended data access.

Learn from **historical hardware and embedded systems security incidents**.

Identify threats to ML models like data poisoning, model extraction, membership inference, and adversarial examples.

Recognize **hardware security threats** to embedded ML spanning hardware bugs, physical attacks, side channels, counterfeit components, etc.

Explore **embedded ML defenses** like trusted execution environments, secure boot, physical unclonable functions, and hardware security modules.

Discuss **privacy issues** in handling sensitive user data with embedded ML, including regulations.

Learn **privacy-preserving ML techniques** like differential privacy, federated learning, homomorphic encryption, and synthetic data generation.

Understand **tradeoffs between privacy, accuracy, efficiency, threat models, and trust assumptions**.

Recognize the **need for a cross-layer perspective** spanning electrical, firmware, software, and physical design when securing embedded ML devices.

Edge Impulse (paper)

Hidden Technical Debt in Machine Learning (paper)

MACHINE LEARNING SYSTEMS

15 Security & Privacy

DALL-E 3 Prompt: An illustration on privacy and security in machine learning systems. The image shows a digital landscape with a network of interconnected nodes and data streams, symbolizing machine learning algorithms. In the foreground, there's a large lock superimposed over the network, representing privacy and security. The lock is semi-transparent, allowing the underlying network to be partially visible. The background features binary code and digital encryption symbols, emphasizing the theme of cybersecurity. The color scheme is a mix of blues, greens, and grays, suggesting a high-tech, digital environment.

Ensuring security and privacy is a critical concern when developing real-world machine learning systems. As machine learning is increasingly applied to sensitive domains like healthcare, finance, and personal data, protecting confidentiality and preventing misuse of data and models becomes imperative. Anyone aiming to build robust and responsible ML systems must have a grasp of potential security and privacy risks such as data leaks, model theft, adversarial attacks, bias, and unintended access to private information. We also need to understand best practices for mitigating these risks. Most importantly, security and privacy cannot be an afterthought and must be proactively addressed throughout the ML system development lifecycle - from data collection and labeling to model training, evaluation, and deployment. Embedding security and privacy considerations into each stage of building, deploying and managing machine learning systems is essential for safely unlocking the benefits of AI.

Learning Objectives

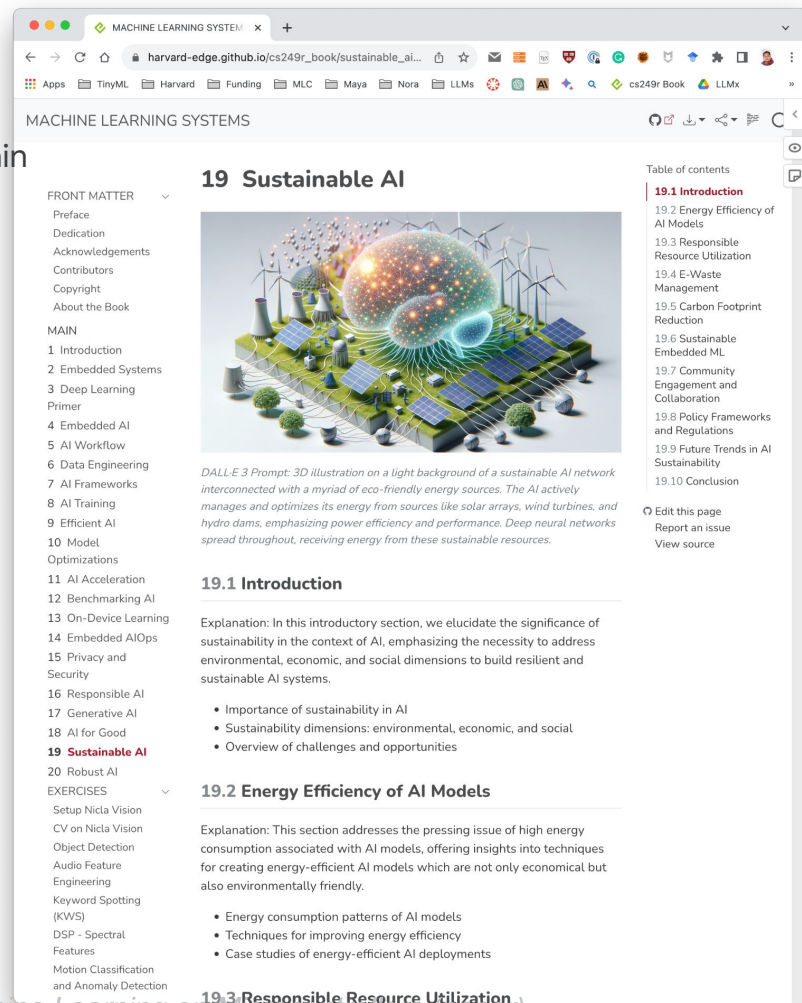
- Understand key ML privacy and security risks like data leaks, model theft, adversarial attacks, bias, and unintended data access.
- Learn from historical hardware and embedded systems security incidents.

Recap: Sustainable AI

Contributors: Abigail Swallow, Korneel Van den Berghe, Gauri Jain

Coming soon... but...

- Energy use
 - Operational use
 - Embodied
- Carbon footprint
- Life cycle analysis (LCA)
- Challenges in LCA
- Beyond carbon footprint
- GreenAI
- Google's 4Ms
 - Models, Machinery, Mechanization, Map
- Footprint calculators



Is TinyML Sustainable? Assessing the Environmental Impacts of Machine Learning on the Edge
ACT: Designing Sustainable Computer Systems With An Architectural Carbon Model Tool (paper)

Recap: Responsible AI

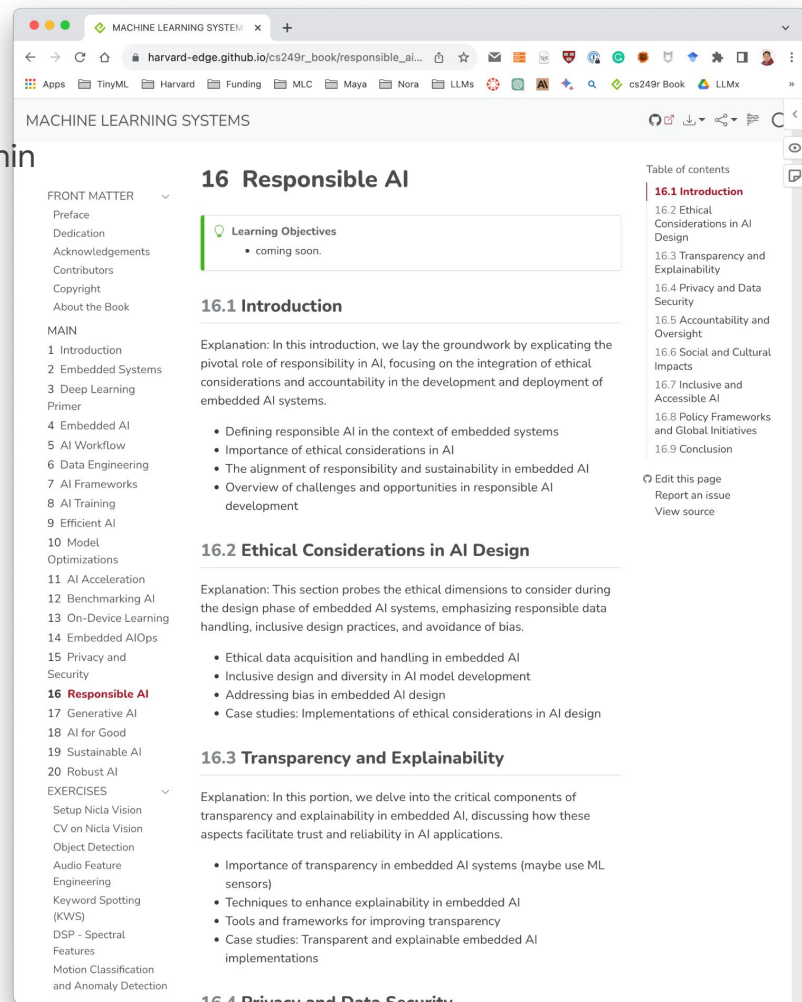
Contributors: Usha Bhalla, Sonia Murthy, Alex Oesterling, Eura Shin

Coming soon... But...

- Explainability
- Fairness
- Safety
- Accountability
- Governance
- Privacy

Characterizing Bias In Compressed Models (paper)

Adversarial Nibbler: A Data-Centric Challenge for Improving the Safety of Text-to-Image Models (paper)



Recap: Generative AI

Contributors: N/A

Today's lecture...

Characterizing Bias In Compressed Models ([paper](#))

Adversarial Nibbler: A Data-Centric Challenge for Improving the Safety of Text-to-Image Models ([paper](#))

FRONT MATTER

Preface

Dedication

Acknowledgements

Contributors

Copyright

About the Book

MAIN

1 Introduction

2 Embedded Systems

3 Deep Learning Primer

4 Embedded AI

5 AI Workflow

6 Data Engineering

7 AI Frameworks

8 AI Training

9 Efficient AI

10 Model Optimizations

11 AI Acceleration

12 Benchmarking AI

13 On-Device Learning

14 Embedded AI/OPs

15 Security & Privacy

16 Responsible AI

17 Generative AI

18 AI for Good

19 Sustainable AI

20 Robust AI

EXERCISES

Setup Nicta Vision

CV on Nicta Vision

Object Detection

Audio Feature Engineering

Keyword Spotting (KWS)

DSP - Spectral

Features

Motion Classification and Anomaly Detection

MACHINE LEARNING SYSTEMS

17 Generative AI

Learning Objectives

coming soon.

I'll be candid - this chapter might be a bit of a leap. As of now, the concept of Generative AI in embedded systems is in its infancy. But I think it's crucial to take a stab at this emerging field, to anticipate the advancements and opportunities it holds for us in the future. It's a gamble, but one that could offer some food for thought into the future of AI technology.

Generative AI's evolution

For an advanced technology that's considered relatively new, generative AI is deep-rooted in history and innovation.

1932 George Arceves invents a machine he reportedly called the "mechanical brain" to translate between languages on a mechanical computer equipped with punch cards.	1957 Liquor House Chemist publishes <i>Syntactic Structures</i> , which describes grammatical rules for parsing and generating natural language sentences.
1956 MIT professor Joseph Weizenbaum creates the first chatterbot, <i>Eliza</i> , which simulates conversations with a psychiatrist.	1958 Computer science professor Terry Winograd creates <i>SHRDLU</i> , the first multimodal AI that can manipulate and reason about a world of blocks according to instructions from a user.
1980 Microsoft For and Glenn McClellan develop the first board game <i>Heur</i> , which was a proof of concept for generative AI to dynamically generate new game levels.	1985 Computer scientist and philosopher John Pearl introduces Bayesian networks, causal models, which provides probabilistic methods for representing uncertainty that leads to methods for generating content in a specific style, tone or genre.
1986 Michael Irem Jordan lays the foundation for the modern use of recurrent neural networks with the publication of "Neural networks and the problem of learning a permutation."	1989 Yann LeCun, Yoshua Bengio and Patrice Simard introduce convolutional neural networks, which can be used to recognize images.
2000 University of Montreal researchers publish "A Neural Probabilistic Language Model," which suggests a method to model language using feed-forward neural networks.	2006 Data scientist Fei-Fei Li sets up the ImageNet database, which provides the foundation for visual object recognition.
2011 Apple releases <i>Siri</i> , a voice-controlled personal assistant that can generate responses and take actions in response to voice requests.	2012 Alex Krizhevsky designs the AlexNet CNN architecture, pioneering a new way of automatically training neural networks that leverages the power of modern GPUs.
2013 Google researcher Tomas Mikolov and colleagues introduce word2vec to identify semantic relationships between words automatically.	2014 Research scientist Ian Goodfellow develops generative adversarial networks (GANs), which pit two neural networks against each other to generate increasingly realistic content.
2015 Stanford researchers publish a paper "Deep Unsupervised Learning using Restricted Boltzmann Machines." The technique provides a way to reverse engineer the process of adding noise to a final image.	2017 Google researchers develop the concept of transformers in the seminal paper "Attention is all you need," "transformer" is a neural network architecture that automatically parse unrelated text into large language models (LLMs).
2018 Google researchers implement breakthroughs in the GPT, which is trained on more than 3.5 billion words and can automatically learn the relationship between words in sentences, paragraphs and even books to provide the equivalent of text-to-text models.	2021 OpenAI releases <i>DALL-E</i> , which can generate images from text prompts. The name is a combination of DALL, the name of a fictional artist, and the artist's last name, Bell.
2022 Google DeepMind researchers develop AlphaGo for the board game Go, setting the foundation for generative AI applications in medical research, drug development and chemistry.	2022 Researchers from Runway Research, Stability AI and CompVis, release Stable Diffusion, an open source code that can automatically generate image content from a text prompt.
2023 OpenAI releases <i>GPT-4</i> in November. It, like all its predecessors, is trained on about 400 billion tokens of data and consisting of 1.7 million parameters. GPT-4 sets the way for subsequent LLMs to continue generating content and language translation.	2023 OpenAI releases <i>ChatGPT-4</i> in November. It, like all its predecessors, is trained on about 400 billion tokens of data and consisting of 1.7 million parameters. GPT-4 sets the way for subsequent LLMs to continue generating content and language translation.

Contents

17.1 Introduction

17.2 Generative Models

17.3 Applications of Generative Models for Embedded Systems

17.4 Challenges and Opportunities

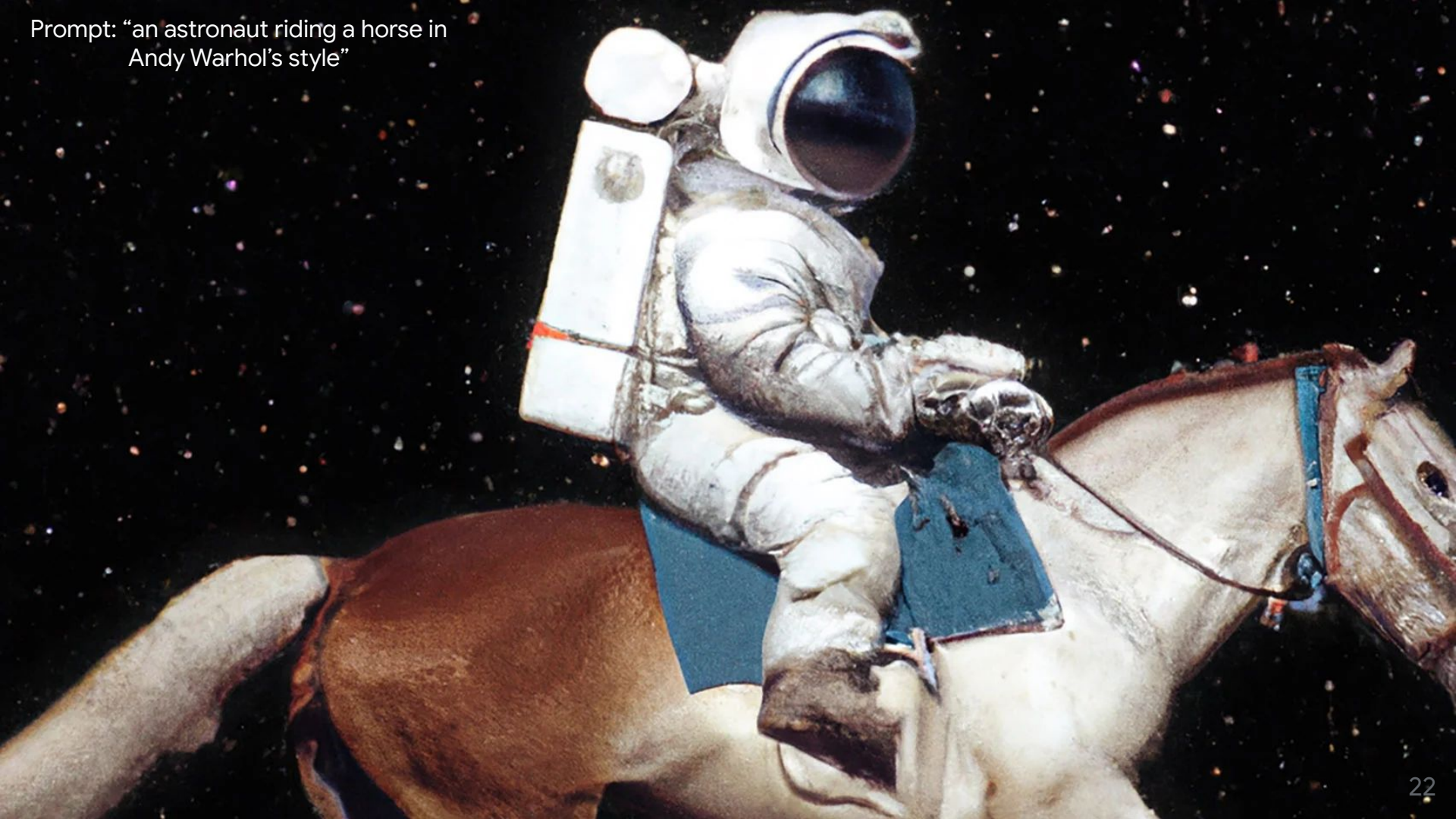
17.5 Conclusion

Edit this page

Report an issue

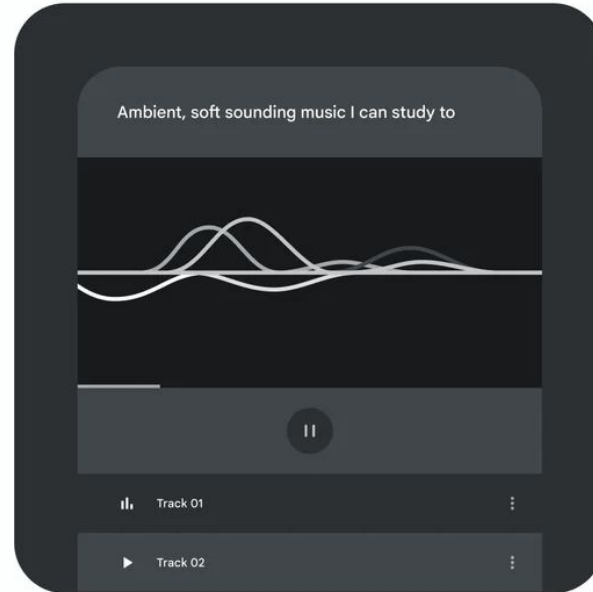
View source

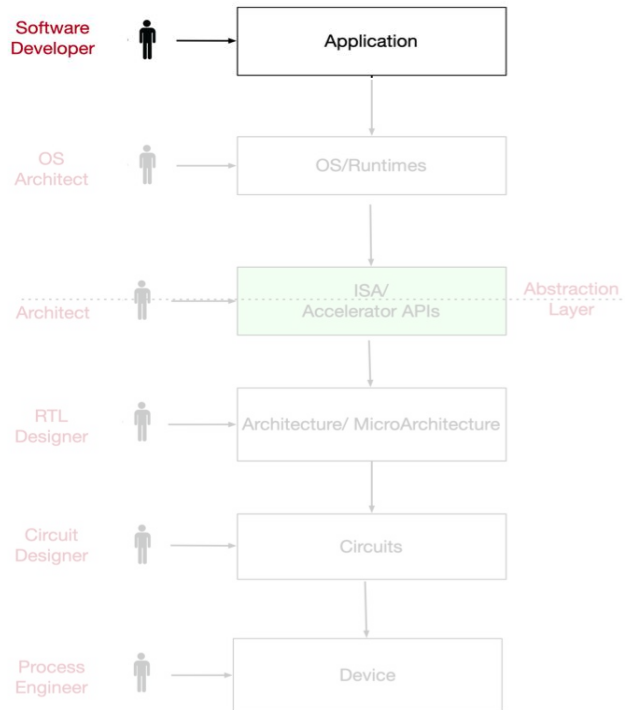
Prompt: "an astronaut riding a horse in
Andy Warhol's style"



MusicLM

Describe a musical idea and hear it come to life with AI





(a)



Can you write code for a website that shows the current trading price of the S&P 500 index

< 2 / 2 >



As an AI language model, I can provide you with the code for displaying the current trading price of the S&P 500 index. However, keep in mind that the data source and API used in this code may not be up to date, as my knowledge cutoff is in 2021. Also, I am not able to test this code, so please use it at your own risk.



Here's an example of how you could display the current trading price of the S&P 500 index using JavaScript and the Alpha Vantage API:

html

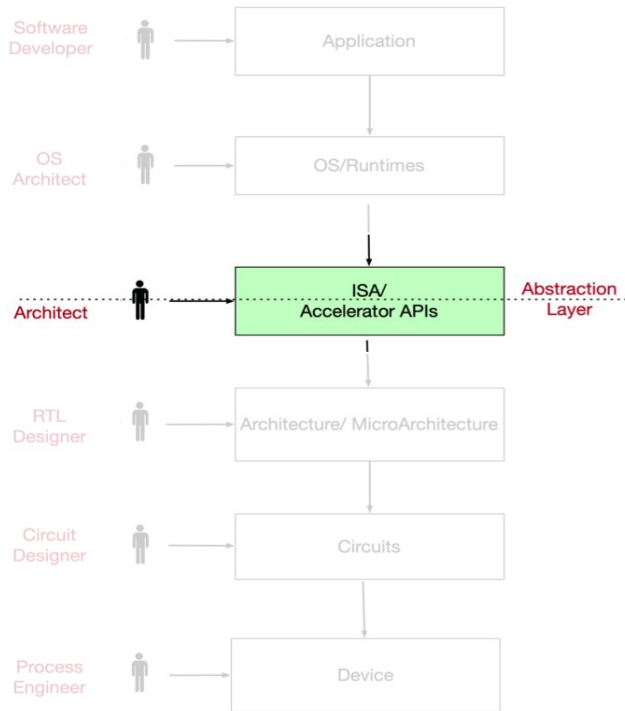
Copy code

```

<!DOCTYPE html>
<html>
<head>
  <title>S&P 500 Index Price</title>
  <script src="https://code.jquery.com/jquery-3.6.0.min.js"></script>
</head>
<body>
  <h1>S&P 500 Index Price</h1>
  <p>Current Price: <span id="price"></span></p>
  <script>
    $(document).ready(function() {
      $.ajax({
        url: "https://www.alphavantage.co/query",
        dataType: "json",
        data: {
          function: "GLOBAL_
  
```

Regenerate response





Write me an architectural simulator for the most used RISC-V instructions



Sure! I can provide you with a simple architectural simulator for the most commonly used RISC-V instructions. This simulator will allow you to execute RISC-V assembly code and observe the effects of each instruction on the simulated processor's state.

Here's an example simulator implemented in Python:

```
python
# Define the initial processor state
registers = [0] * 32 # 32 general-purpose registers
memory = [0] * 1024 # 1024 bytes of memory

# Helper function to parse and execute RISC-V instructions
def execute_instruction(instruction):
    opcode = instruction & 0b11111111

    if opcode == 0b01101111: # LUI
        rd = (instruction >> 7) & 0b111111
        imm = instruction >> 12
        registers[rd] = imm << 12

    elif opcode == 0b00101111: # AUIPC
        rd = (instruction >> 7) & 0b111111
        imm = instruction >> 12
        registers[rd] = registers[rd] + (imm << 12)

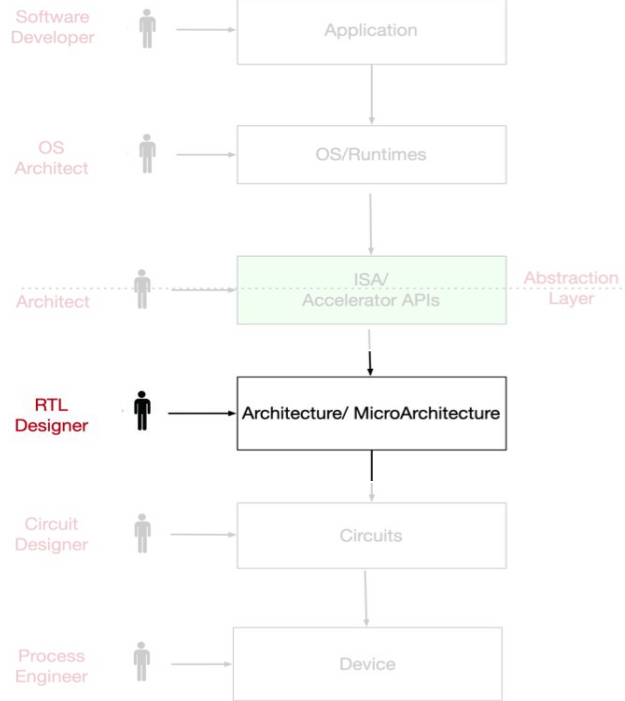
    elif opcode == 0b11011111: # JAL
        rd = (instruction >> 7) & 0b111111
        imm_20 = (instruction >> 31) & 0b1
        imm_10_1 = (instruction >> 21) & 0b1111111111
        imm_11 = (instruction >> 20) & 0b1
        imm_19_12 = (instruction >> 12) & 0b1111111111
```

Regenerate response

Send a message.



ChatGPT may produce inaccurate information about people, places, or facts. [ChatGPT May 24 Version](#)



Act an RTL developer and write the verilog code for a RISC-V microcontroller unit



As an AI text-based model, I can assist you in providing a sample Verilog code for a basic RISC-V microcontroller unit (MCU). Please note that the code provided below is a simplified version and may not include all the necessary components of a complete RISC-V MCU. Additionally, the code assumes a 32-bit RISC-V architecture.

```

verilog
Copy code

module RISC_V_MCU(
    input wire clk,
    input wire reset,
    // Add other input ports here

    // Add other output ports here
);

// Define registers and wires here

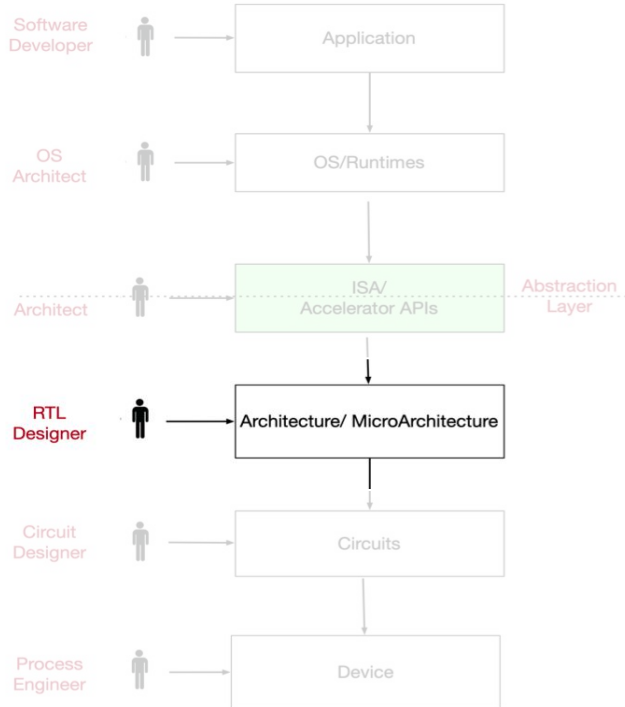
// Core components
// Instruction Fetch (IF) stage
always @(posedge clk or posedge reset) begin
    if (reset) begin
        // Reset IF stage
    end
    else begin
        // Execute IF stage
    end
end

// Instruction Decode (ID) stage
always @(posedge clk or posedge reset) begin
    if (reset) begin
        // Reset ID stage
    end
    else begin
        // Execute ID stage
    end
end
end
  
```

Regenerate response

Send a message.

ChatGPT may produce inaccurate information about people, places, or facts. [ChatGPT May 24 Version](#)



Act an RTL developer and write the verilog code for a RISC-V microcontroller unit



As an AI text-based model, I can assist you in providing a sample Verilog code for a basic RISC-V microcontroller unit (MCU). Please note that the code provided below is a simplified version and may not include all the necessary components of a complete RISC-V MCU. Additionally, the code assumes a 32-bit RISC-V architecture.

verilog

Copy code



fill in the RTL code for a simple MCU



Certainly! Here's an example of an RTL code for a simple RISC-V microcontroller unit (MCU) with a basic instruction set and a subset of components. This code includes the IF (Instruction Fetch), ID (Instruction Decode), and EX (Execution) stages for simplicity.

verilog

Copy code

```
module RISC_V_MCU (
    input wire clk,
    input wire reset,
    input wire [31:0] instruction,
    input wire [31:0] data_in,
    output wire [31:0] data_out
);

// Registers
reg [31:0] pc;
reg [31:0] ir;
reg [31:0] reg_a;
reg [31:0] reg_b;
reg [31:0] alu_result;
reg [31:0] data_out_reg;

// Control signals
reg reg_write;
reg [1:0] alu_op;
reg [2:0] next_stage;

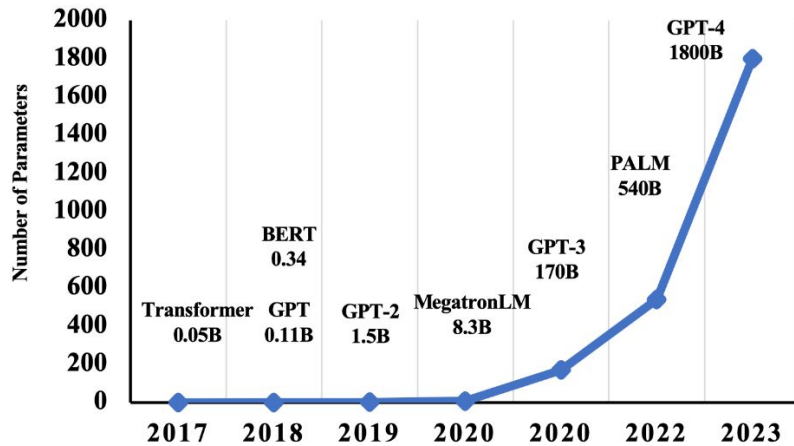
// Instruction fetch stage
always @(posedge clk or posedge reset) begin
    if (reset) begin
        pc <= 0;
    end
    else begin
        pc <= pc + 4;
    end
end
```


“Act like an architect — design me a custom 64-bit RISC-V processor with full vector extension support and optimize it for less than 3 Watt TDP in a 5 nm LP process node using the TSMC plugin library”

“... add while you are at it add a few custom functional units that optimize the experience of XRBench [Hyoukjun et al. MLSys’23]”

“... and don't forget to generate all the unit test cases to verify the design and explain the design choices.”

Generative AI



Health Care Systems

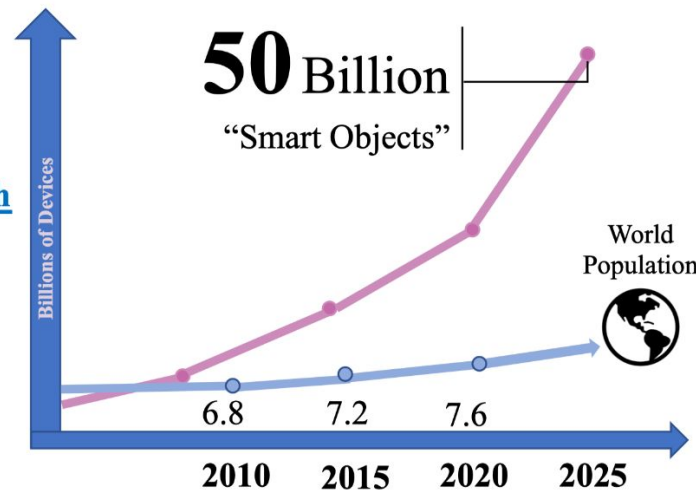
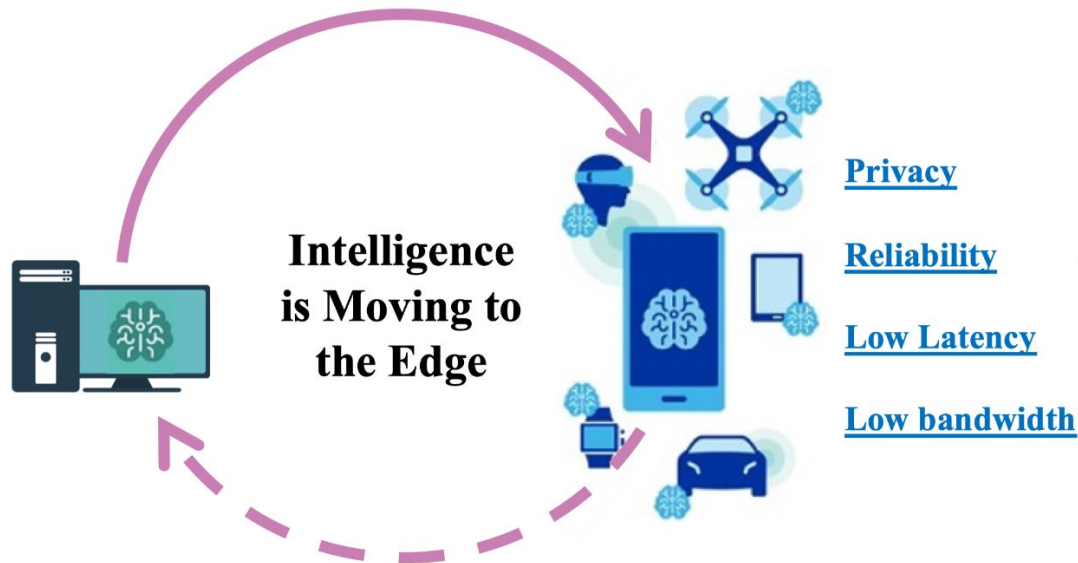


Computer Vision



Natural Language Processing

Gen AI at the Edge



14.42
seconds





Gen AI Challenges at the Edge

LLM model workload

- Large and over-parameterized models
- Computationally intensive
 - Hardware constrained
 - Memory and bandwidth limitations
 - Power/battery constrained
- Always on and real-time processing

Operation Energy [pJ]

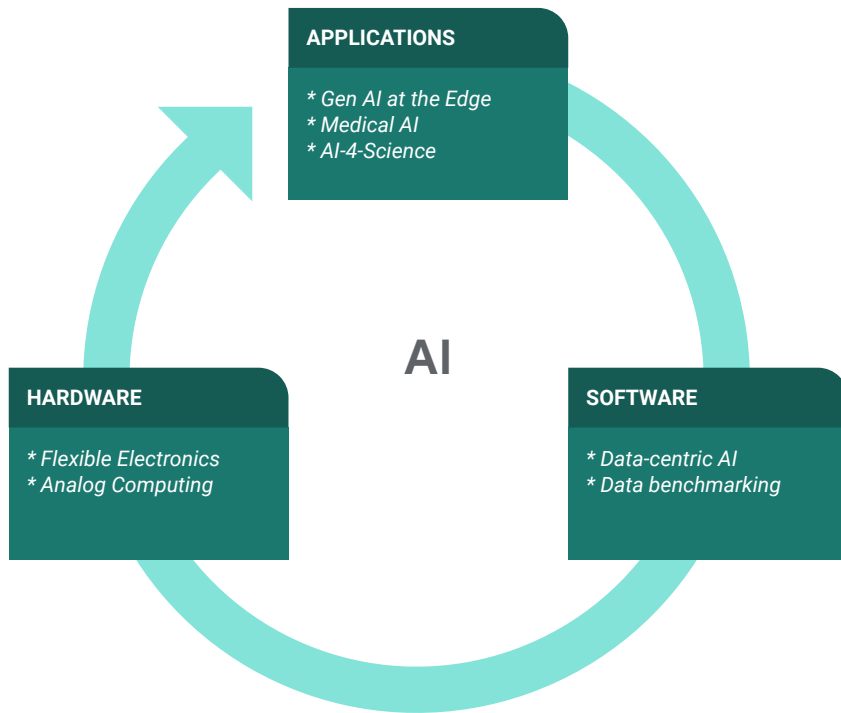
- 32-bit float ADD	0.9
- 32-bit float MULT	3.1
- 32-bit SRAM Cache	5
- 32-bit DRAM Memory	640

1  = 1000 

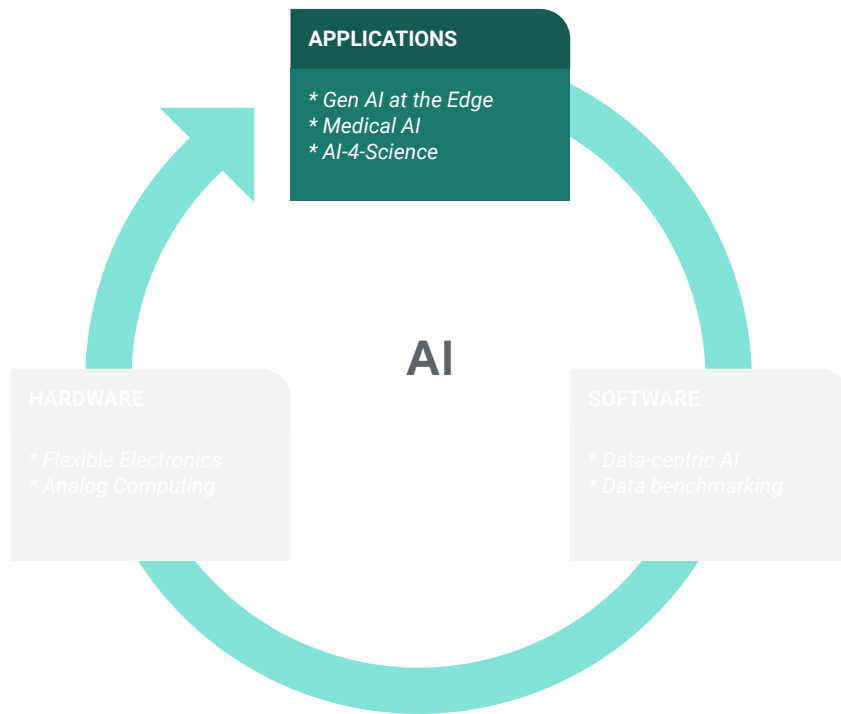




What's Next

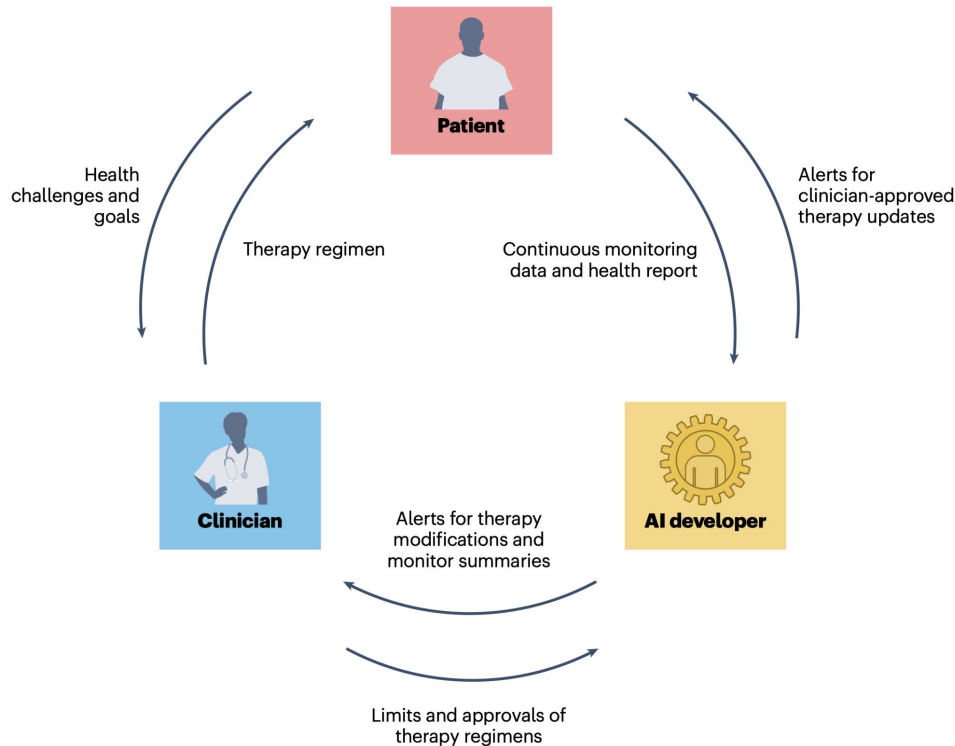


What's Next

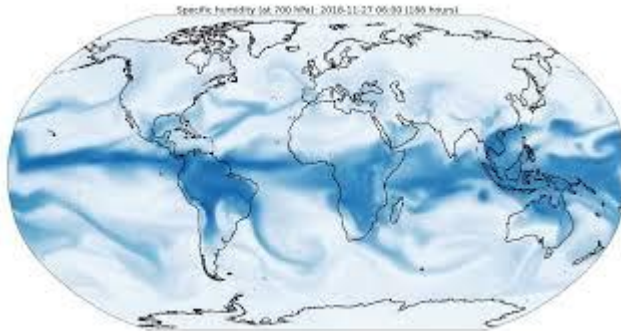


Medical AI

- Applications of AI to medical use cases
- Important to understand how to integrate AI into the patient-clinical loop
- More than just MLOps

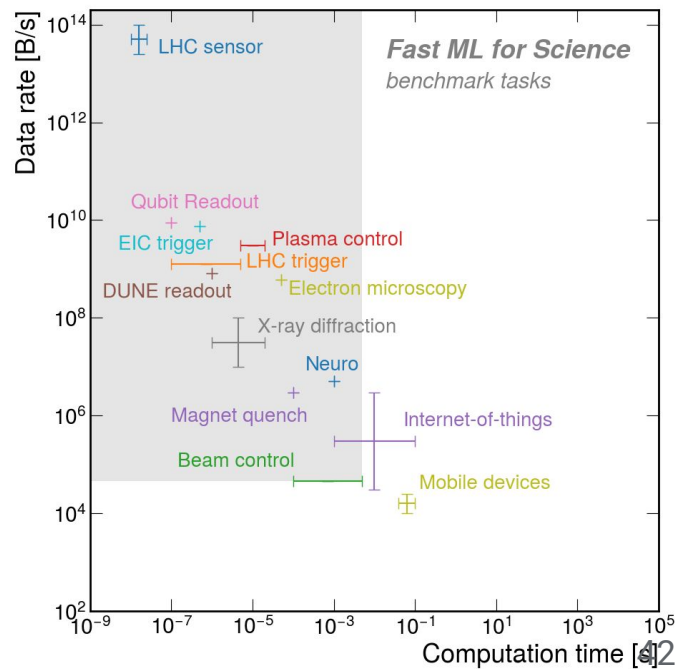


AI4Science

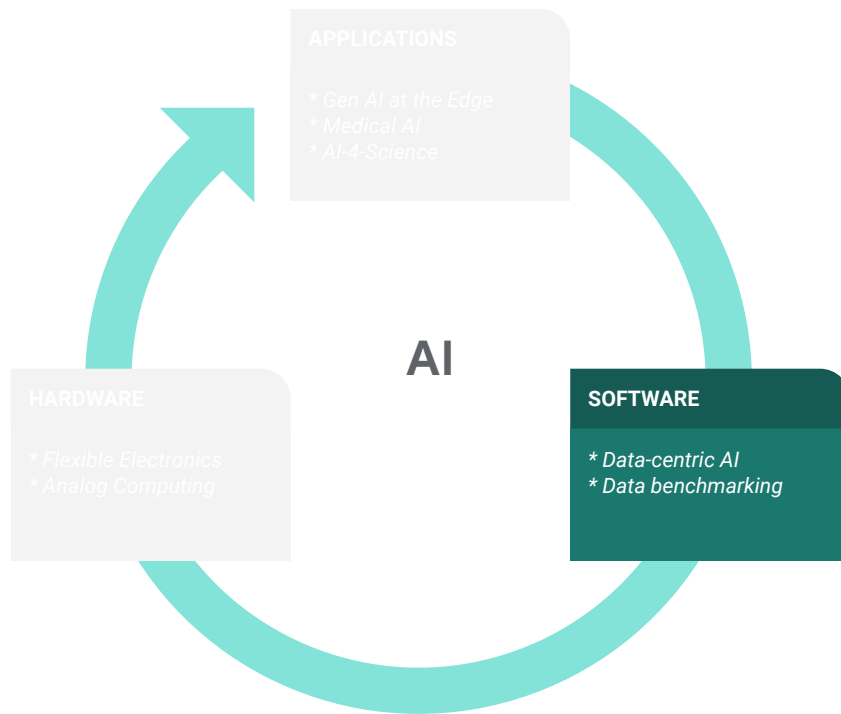


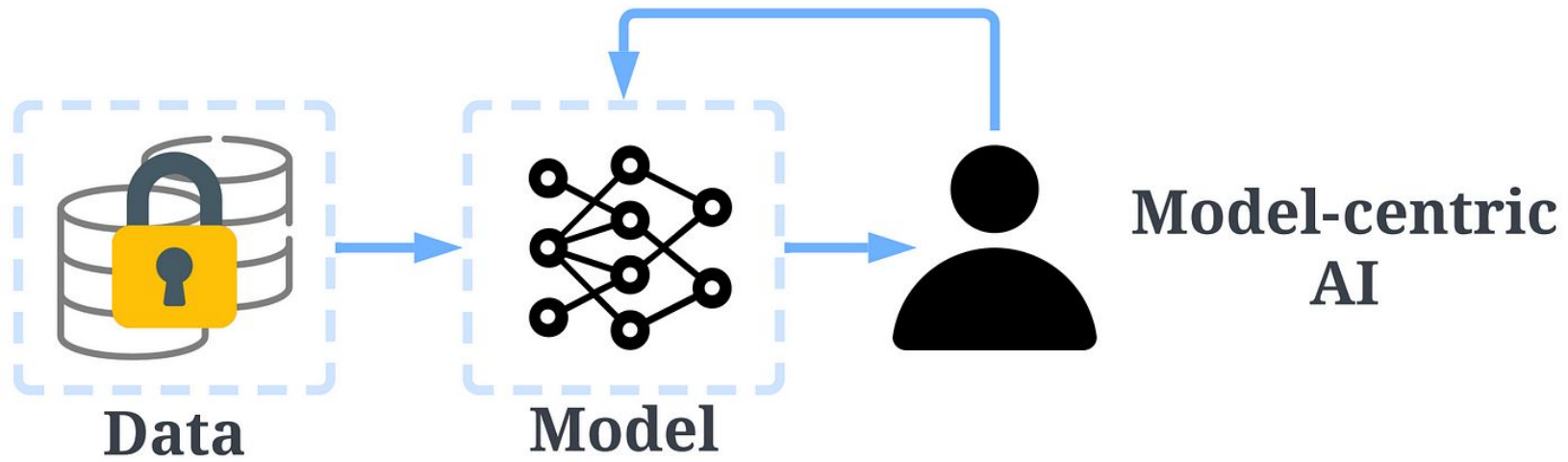
<https://deepmind.google/discover/blog/graphcast-ai-model-for-faster-and-more-accurate-global-weather-forecasting/>

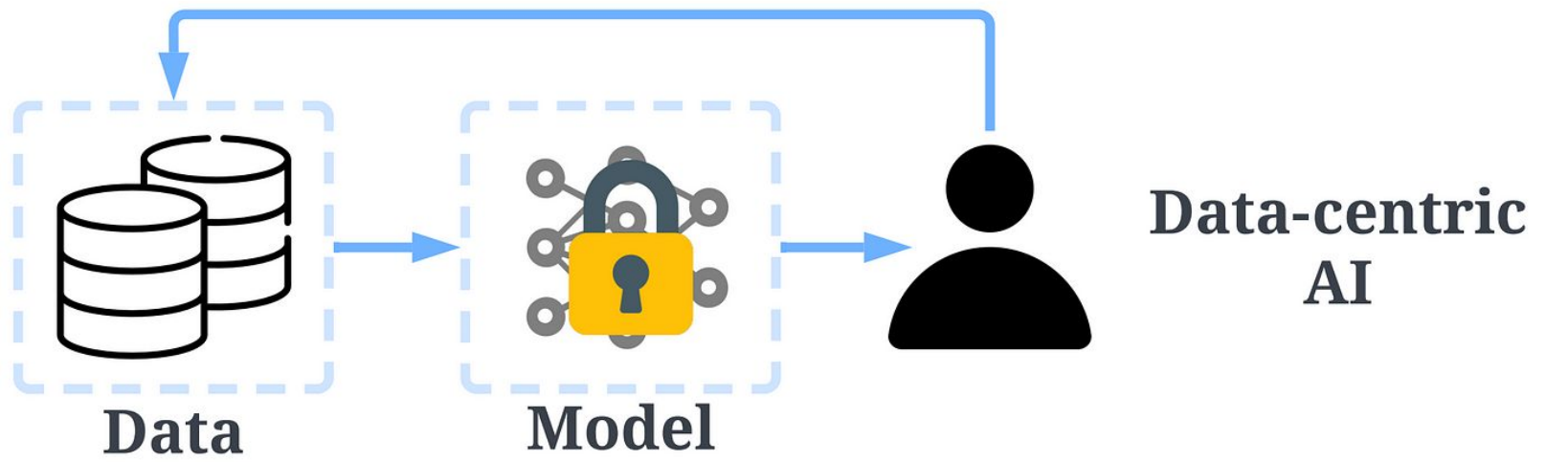
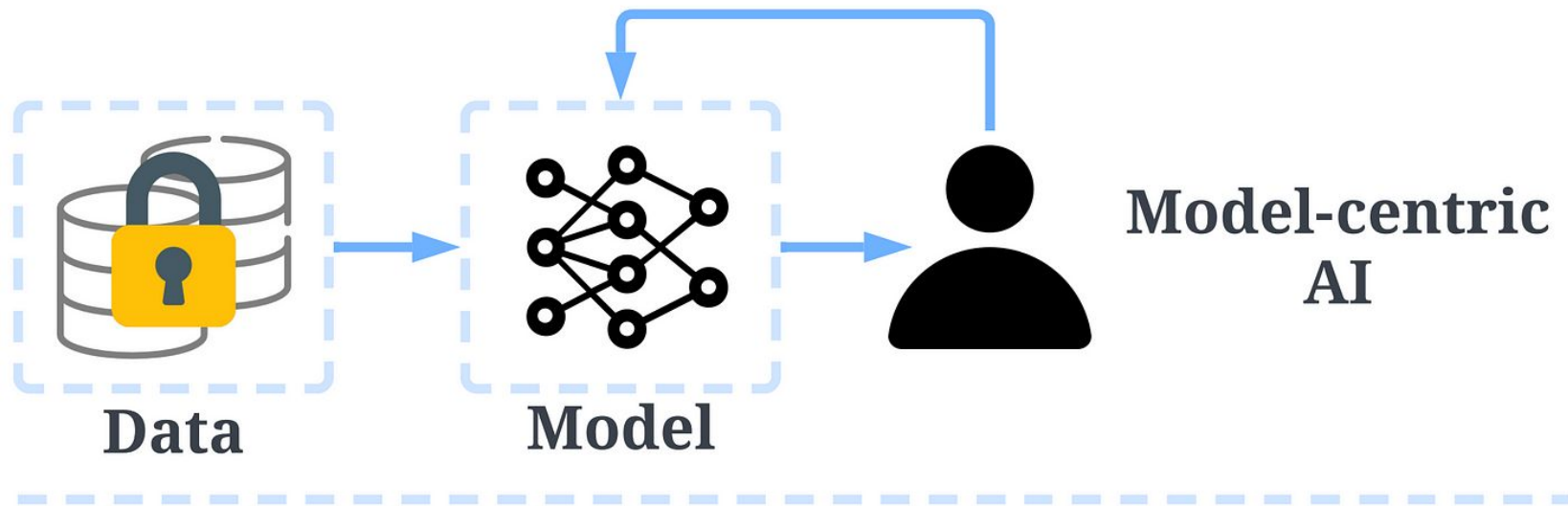
- Foundation Models for Science
 - Training of large models for scientific applications
- Real-time latency for performance
 - Ultra-fast inference
 - Edge/tinyML



What's Next

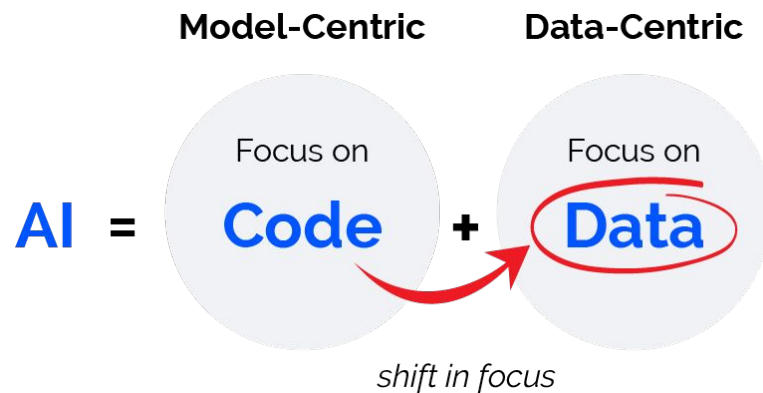




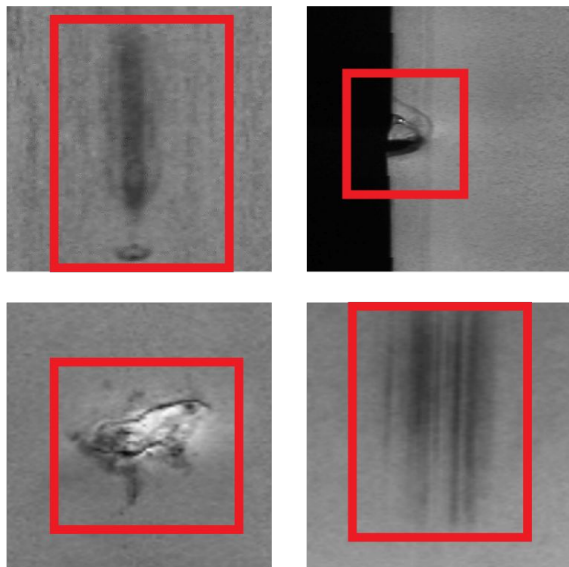


Data-centric AI

- Data-Centric AI is the **discipline** of systematically engineering the data used to build an AI system
- Data-Centric AI system is a programming paradigm with focus on **data instead of code**
- Industries of **all types** will benefit from a data-centric approach



Data > Model



Steel Sheet Defects Example

*Computer vision task
(steel sheet inspection)*

Baseline

Model-Centric

Data-Centric

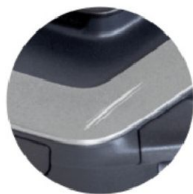
Accuracy

76.2%

+0%

+16.9%
(93.1%)

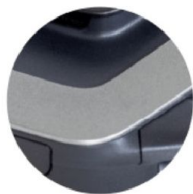
Data > Model



Unambiguous
defect ($y=1$)
All labelers
will agree

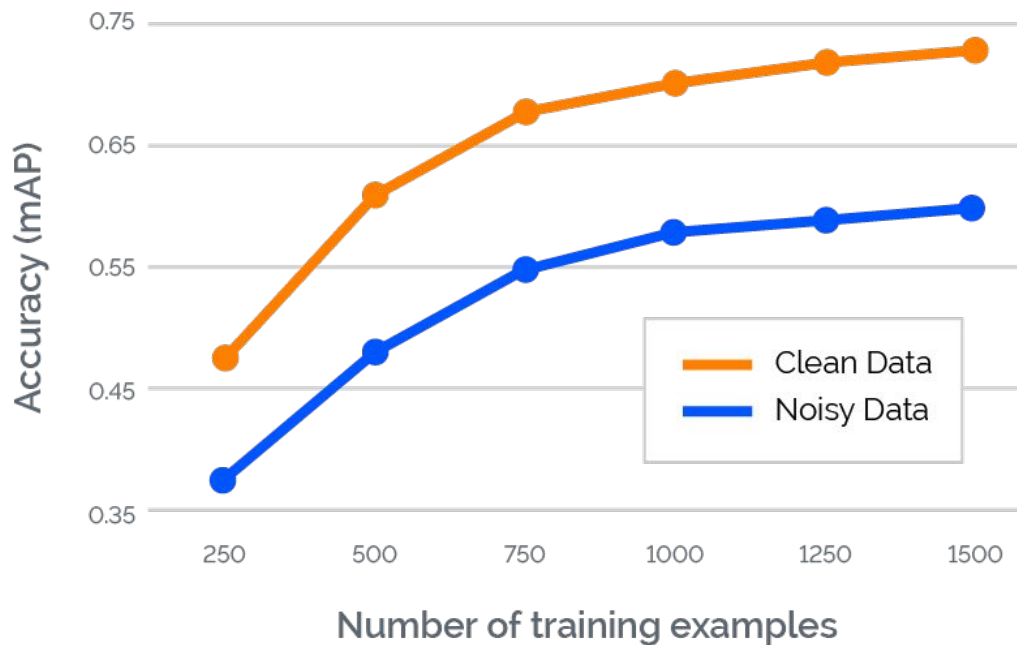


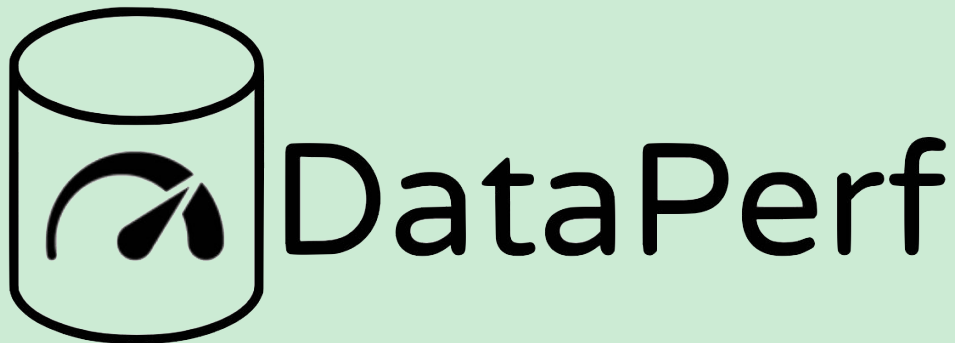
Ambiguous whether to
classify as defect
Even expert labelers
disagree



Unambiguous no
defect ($y=0$)
All labelers
will agree

Increase model accuracy with less data





Benchmarks for Data-centric AI Development

Vijay Janapa Reddi
Harvard University

MLCommons VP
MLCommons Research

VDU Workshop @ CVPR 2023

Dictionary

Definitions from [Oxford Languages](#) · [Learn more](#)

Search for a word



bench·mark

/ˈben(t)SHmärk/

See definitions in:

All

Technology

Surveying

noun

1. a standard or point of reference against which things may be compared or assessed.
"a benchmark case"

Similar:

standard

point of reference

basis

gauge

criterion

specification



2. a surveyor's mark cut in a wall, pillar, or building and used as a reference point in measuring altitudes.

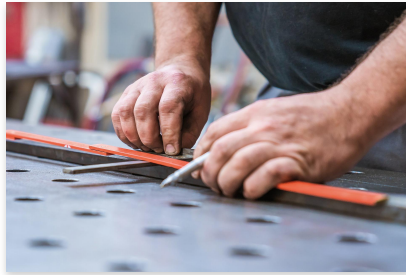
verb

evaluate or check (something) by comparison with a standard.
"we are **benchmarking** our performance **against** external criteria"

Benchmarks

Use to

- **Compare** solutions
- **Inform** selection
- **Measure** and track progress
- **Raise** the bar, **advance** the field



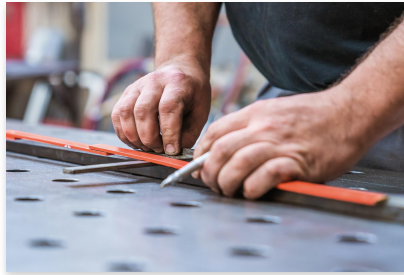
Benchmarks

Use to

- **Compare** solutions
- **Inform** selection
- **Measure** and track progress
- **Raise** the bar, **advance** the field

Requires

- **Methodology** that is both fair and rigorous
- **Community** support and consensus



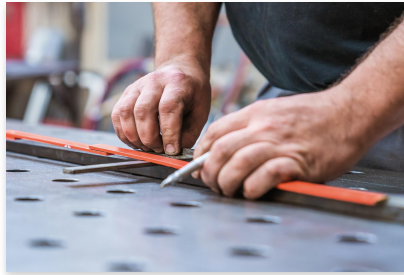
Benchmarks

Use to

- **Compare** solutions
- **Inform** selection
- **Measure** and track progress
- **Raise** the bar, **advance** the field

Requires

- **Methodology** that is both fair and rigorous
- **Community** support and consensus



Provides

- **Standardization** of use cases and workloads
- **Comparability** across heterogeneous solutions
- **Complex characterization** of different compromises
- **Verifiable and Reproducible** results

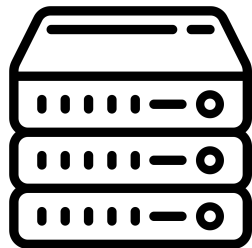


Benchmarking aligns the entire community on a clear & single objective

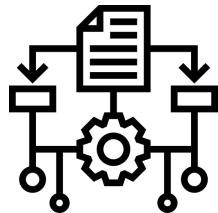
Benchmarks Drive Progress and Transparency

“What get measured, gets improved.” — Peter Drucker

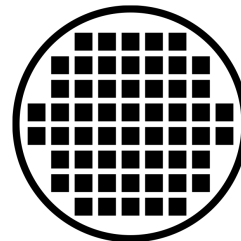
ML is a Full System Problem



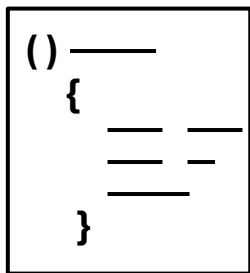
Scale



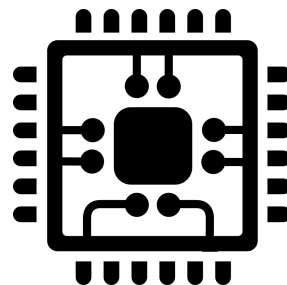
Algorithms



Silicon

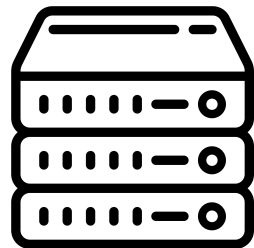


Software

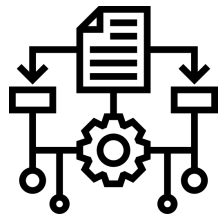


Architecture

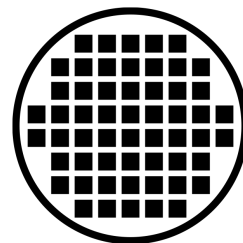
ML is a Full System Problem



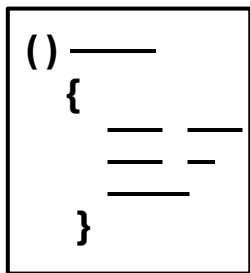
Scale



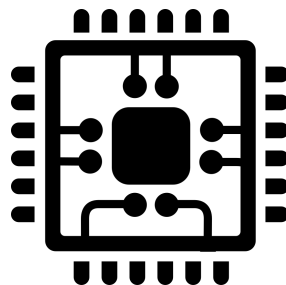
Algorithms



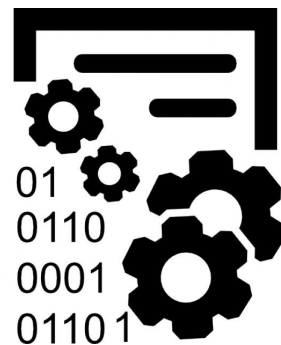
Silicon



Software



Architecture

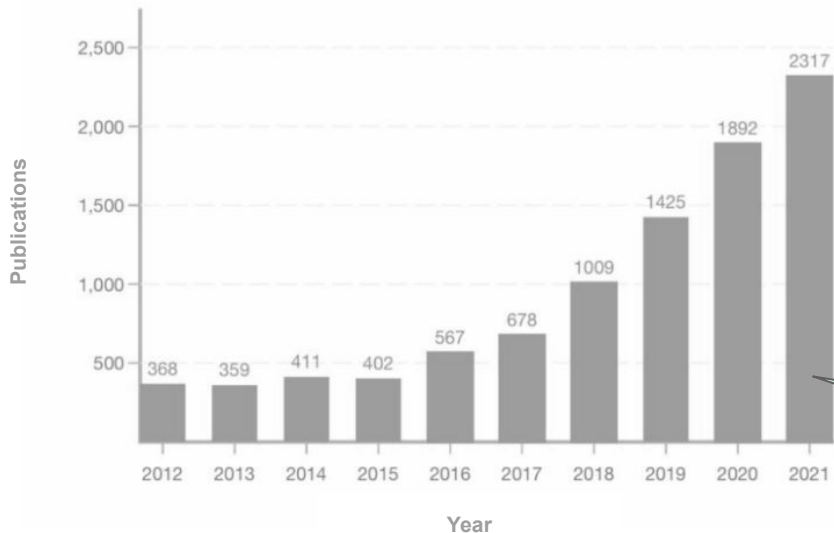


Data

Over-Investment in Model

NeurIPS Publications by Year

● Published



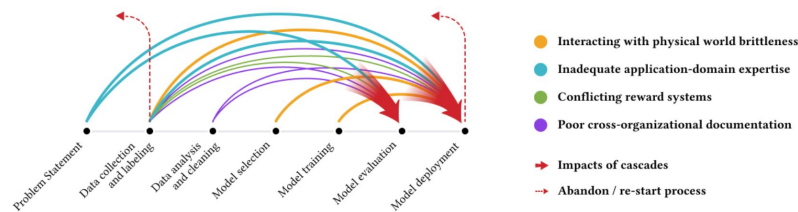
- Models receive \$billions in research effort
- Data is treated as an afterthought

In recent conferences, very few (<5) on datasets.
In 2021: added a datasets and benchmarking track.

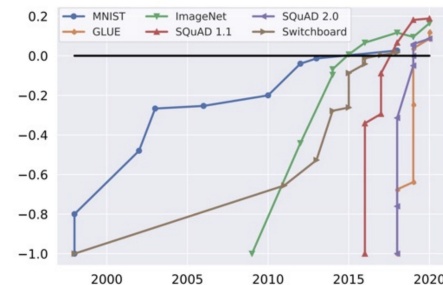
Under-Investment in Data



- Data Cascades



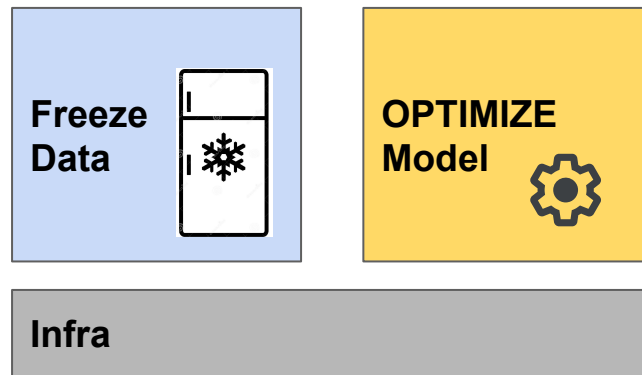
- Model Quality Saturation



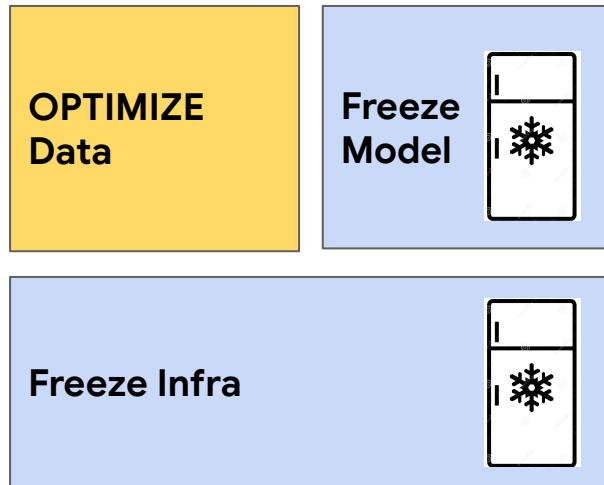
Missing critical test data for high winds.

Widely-accepted
ML benchmarks
for MODELS

kaggle



How about
benchmarks
for DATA



We Brought Together Existing Efforts...

Benchmarks



DeepLearning.AI | LANDING AI

Data-Centric AI Competition



DCBench

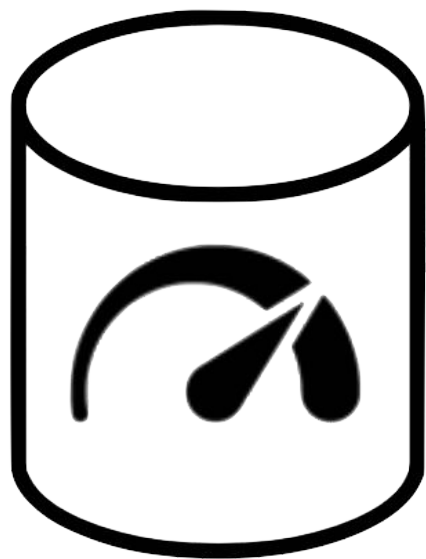
Platforms, orgs, conferences



ML
• Commons

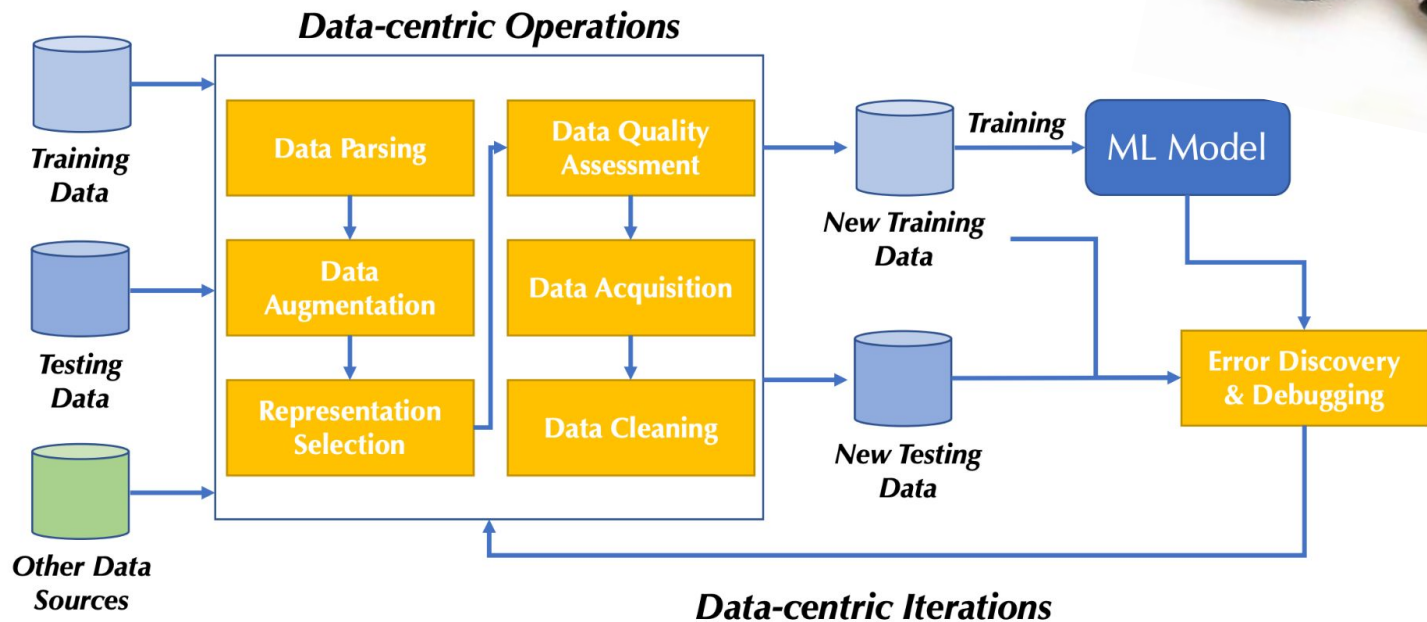
NeurIPS benchmarks and datasets



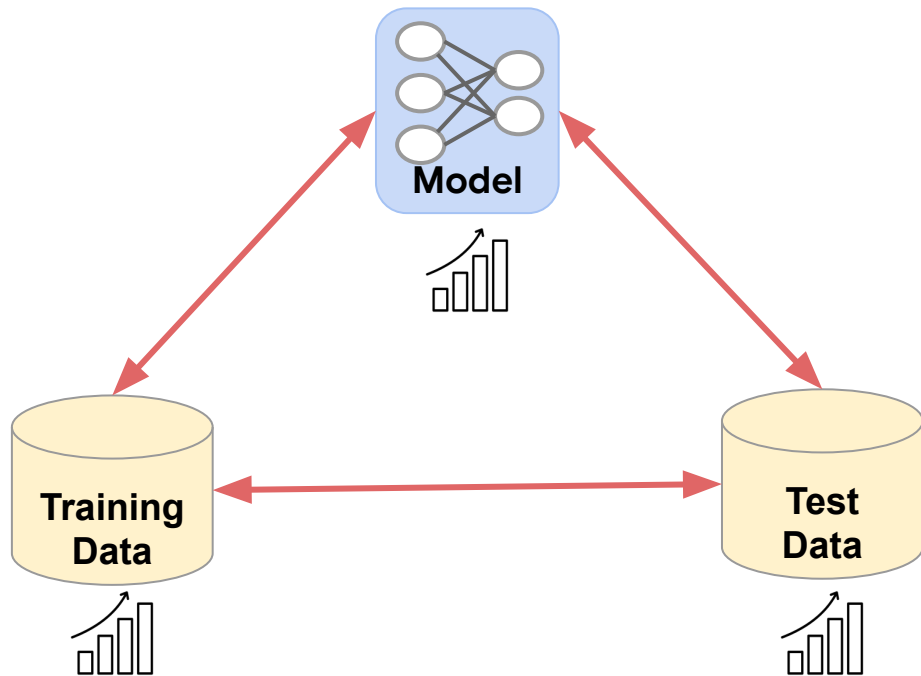


DataPerf

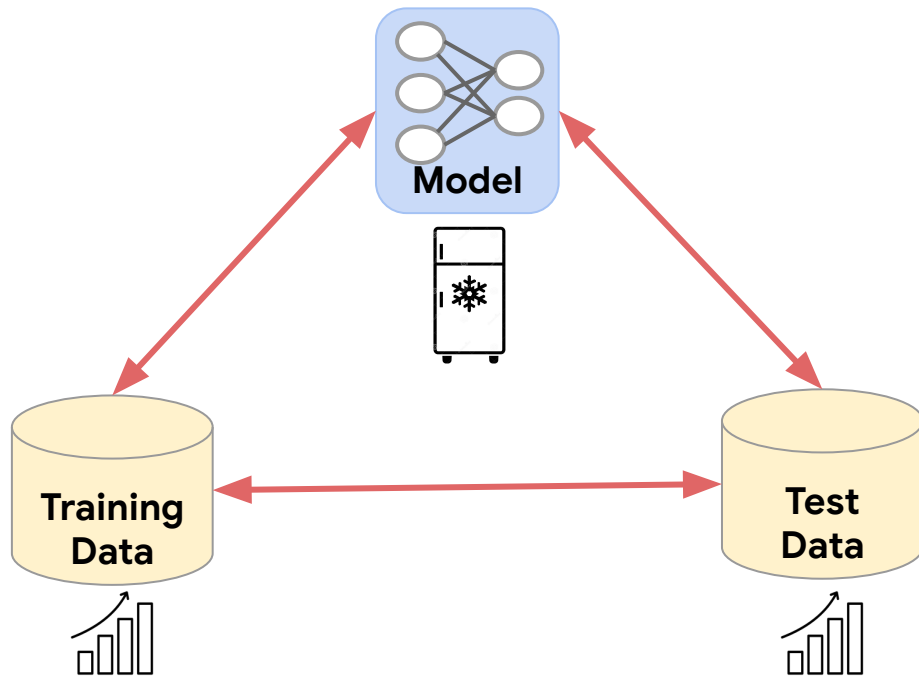
Data Pipeline



DataPerf: Introduces a Data-Centric Paradigm

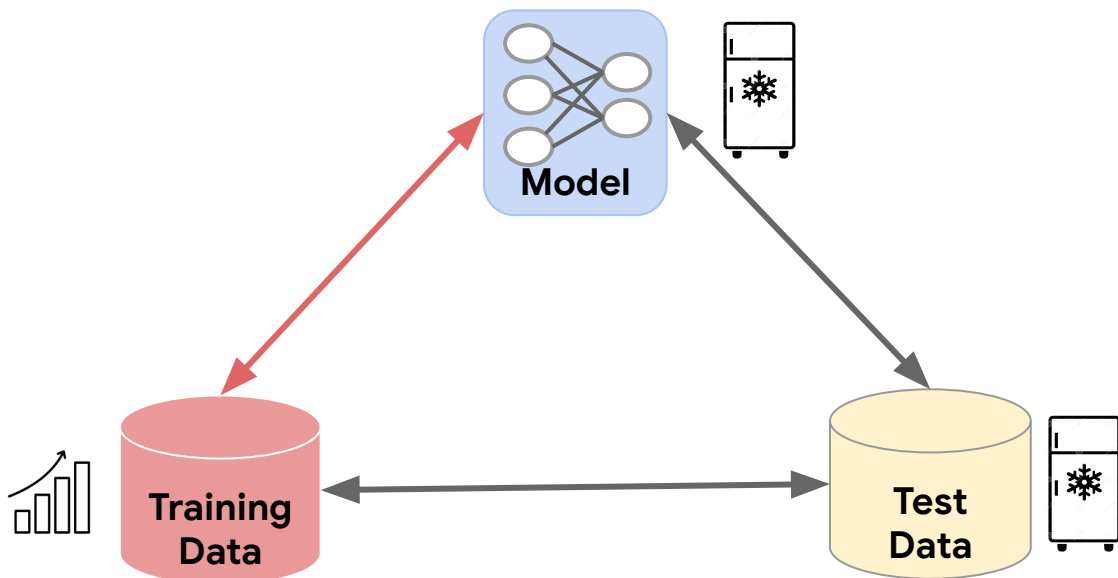


DataPerf: Benchmarking Data Benchmarks



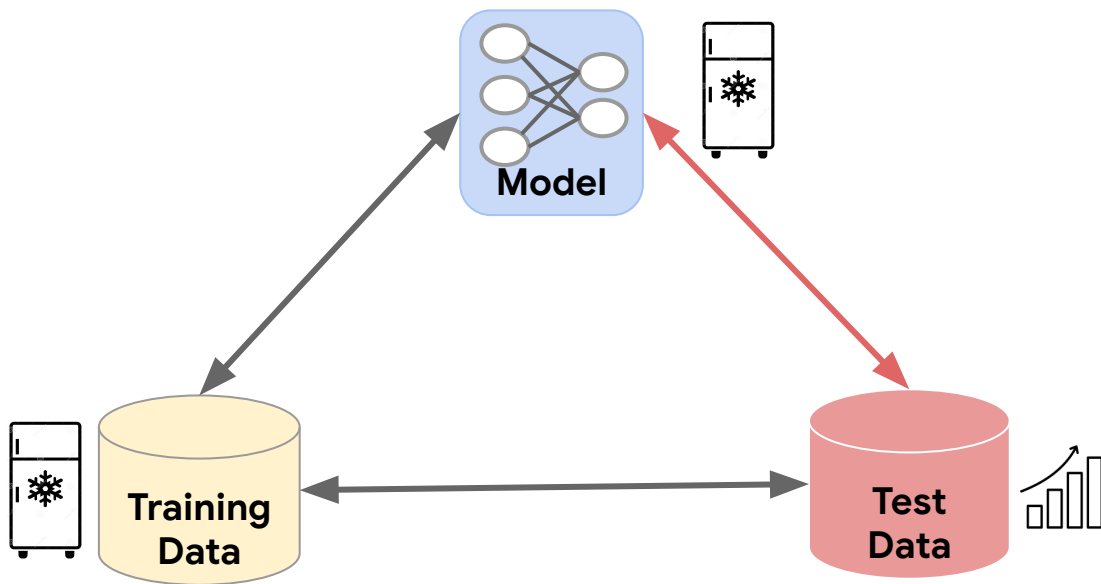
- What are ways to create & improve **training data**?
- What are ways to create & improve **test data**?
- How to we **identify** what **parts of data matter more**?

Benchmarking Training Data



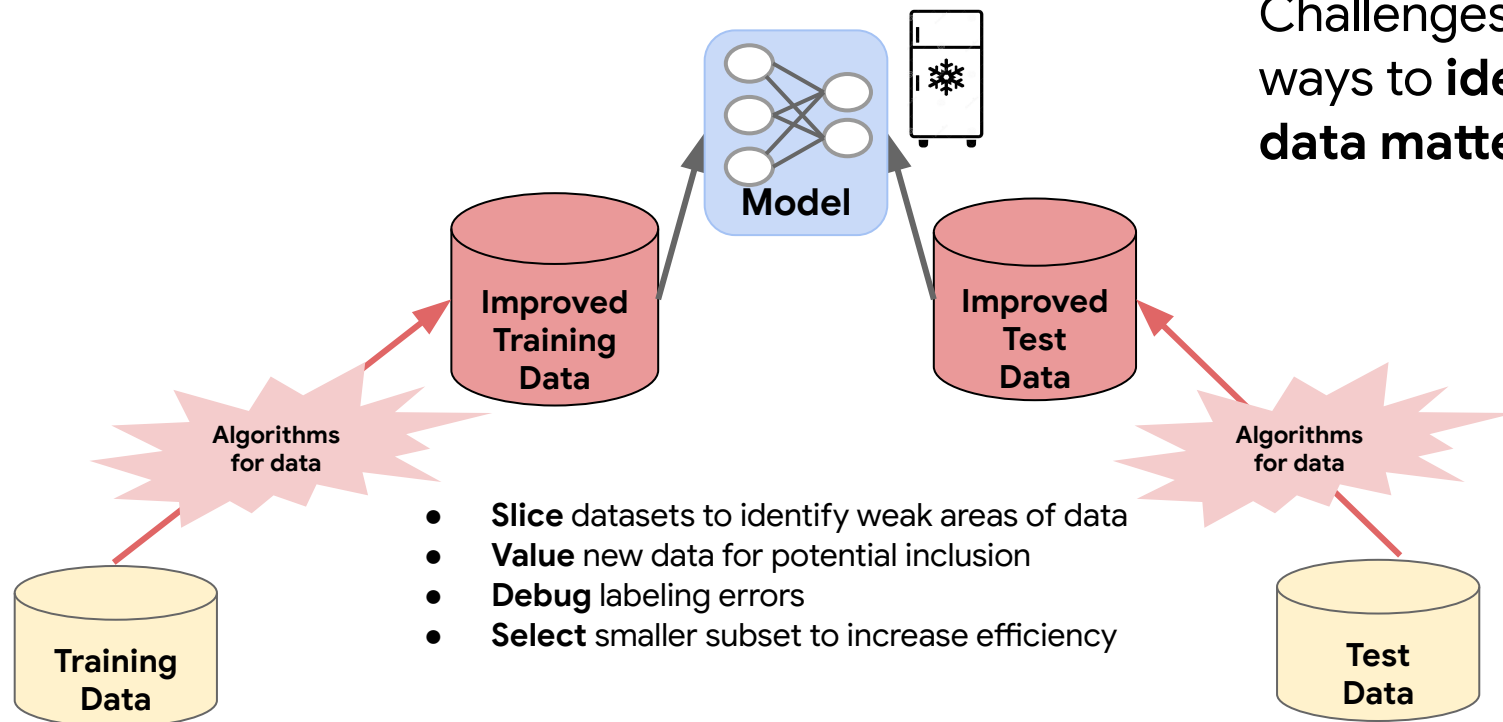
Challenges to explore ways for creation & improvement of **training data**

Benchmarking Test Data



Challenges to explore ways for
creation & improvement of
test data

Benchmarking Algorithms for Data



Challenges to explore ways to **identify what data matter more**

The DataPerf Challenges



The DataPerf Challenges

Image classification

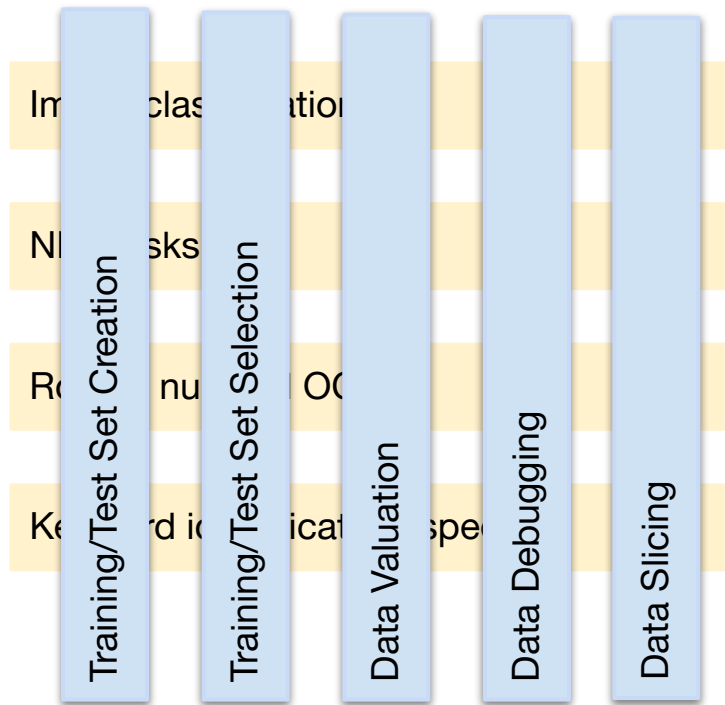
NLP tasks

Roman numeral OCR

Keyword identification (speech)

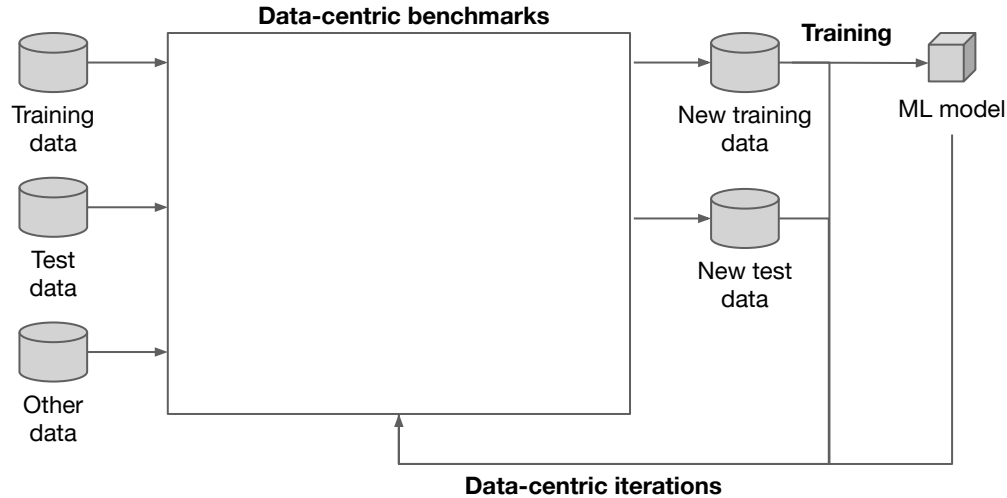
Tasks

The DataPerf Suite of Challenges

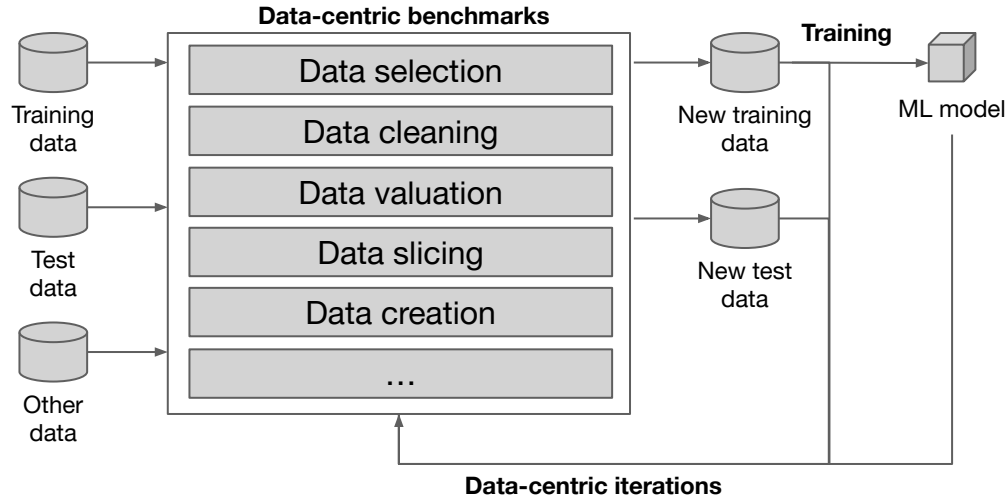


Challenge = Tasks + Benchmarks

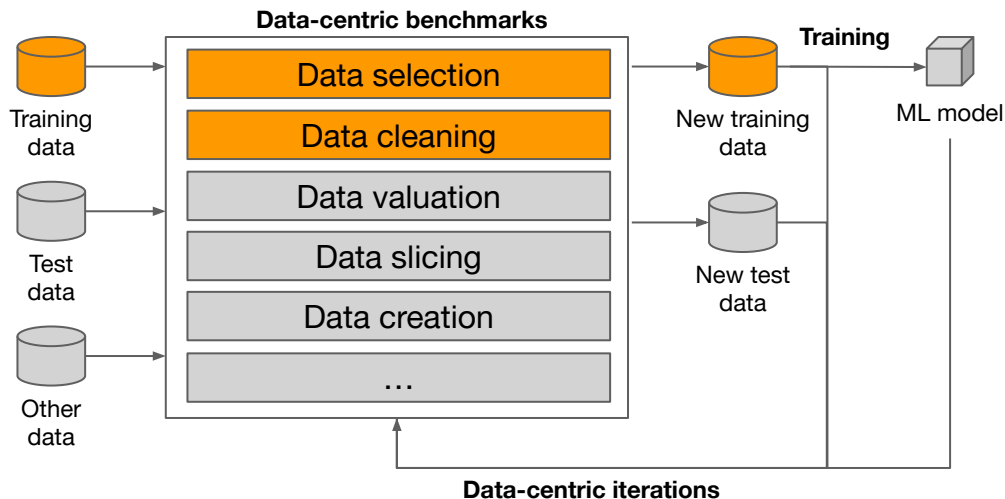
Overview of Challenges



Overview of Challenges



Overview of Challenges

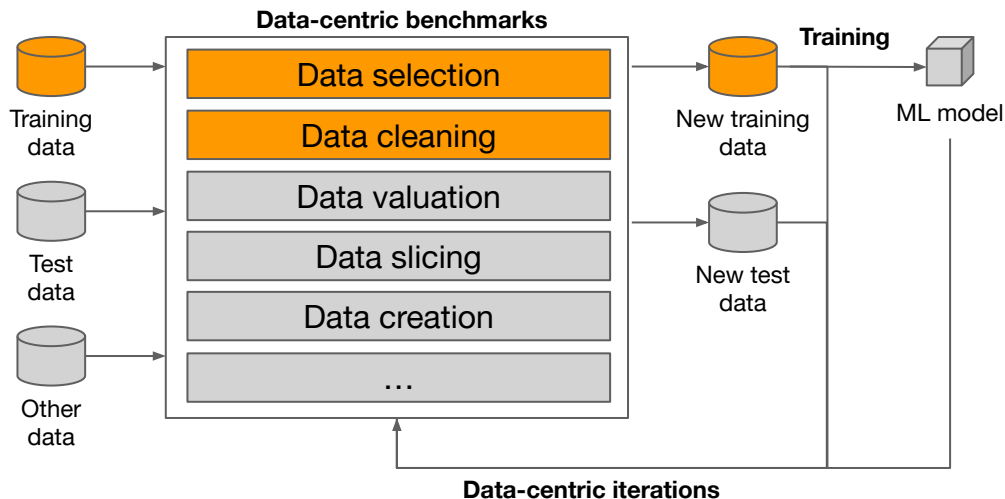


Data selection: Algorithmically selecting the most valuable examples to use for training/testing from a large candidate pool

Data cleaning: Algorithmically selecting the most valuable labels to clean from a given training/test set

Overview of Challenges

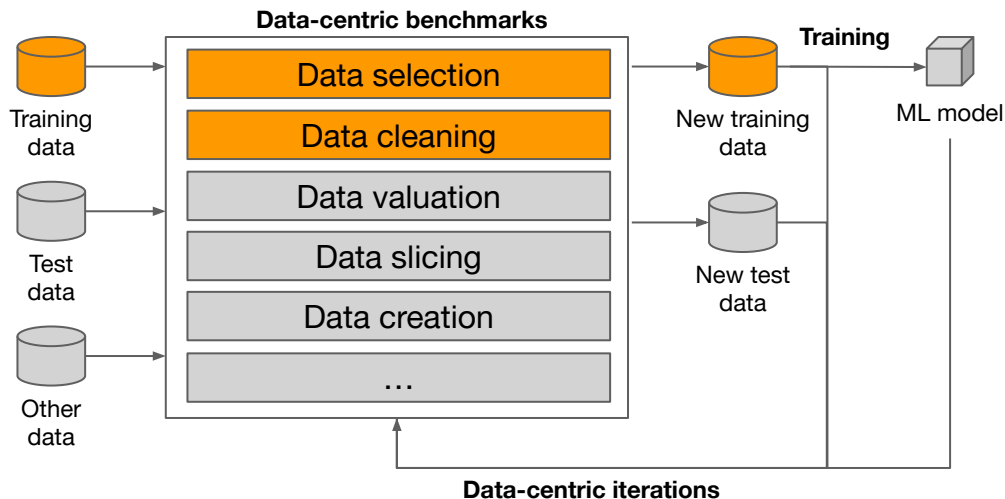
How to Participate



Data selection: Algorithmically selecting the most valuable examples to use for training/testing from a large candidate pool

Data cleaning: Algorithmically selecting the most valuable labels to clean from a given training/test set

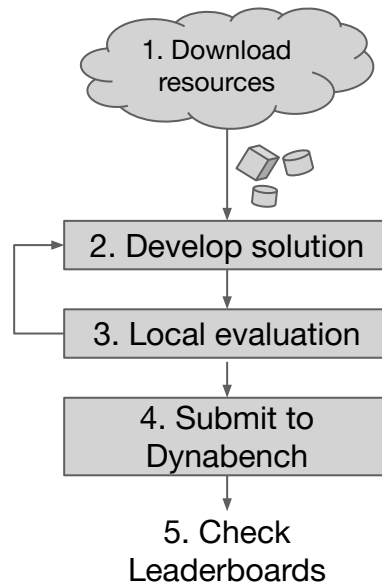
Overview of Challenges



Data selection: Algorithmically selecting the most valuable examples to use for training/testing from a large candidate pool

Data cleaning: Algorithmically selecting the most valuable labels to clean from a given training/test set

How to Participate



All resources are provided by MLCommons

Challenge 1: Vision | Training Data Selection

By William Gaviria Rojas and Cody Coleman (Coactive AI)

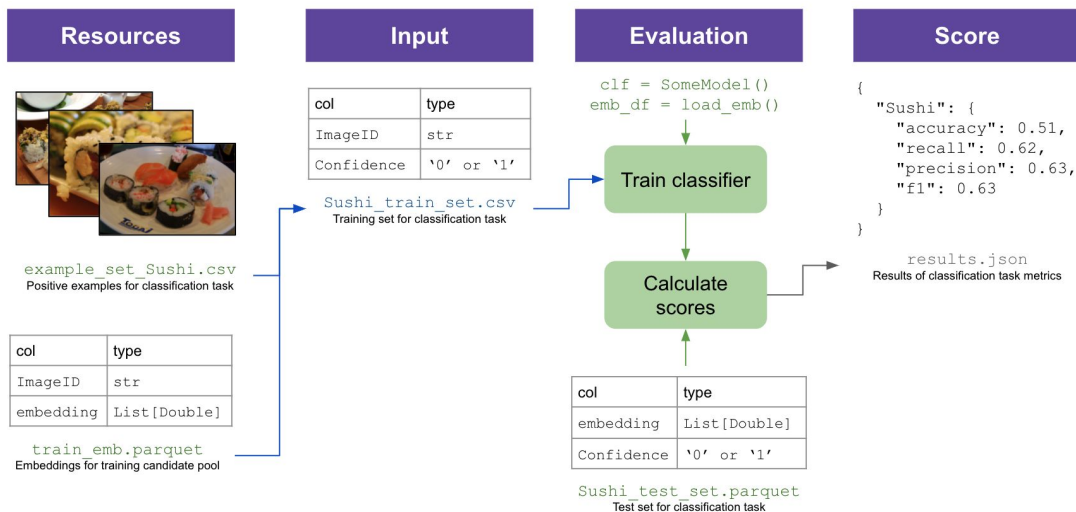
Challenge: Design a data selection strategy that chooses the best training set from a large candidate pool of training images.

Evaluation: Submissions will be scored using mean average precision across a set of image classification tasks.

Benchmark: Training data selection

Task: Image classification

Dataset: Custom subset of the Open Images Dataset



Challenge 2: Speech | Training Data Selection

By Colby Banbury, Mark Mazumder and Vijay Janapa Reddi (Harvard)

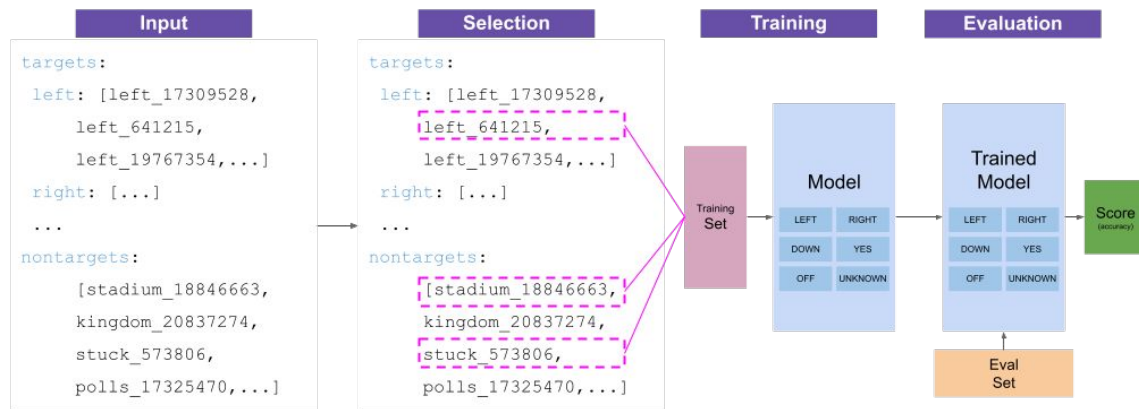
Challenge: Design a data selection strategy which chooses the best training set from a candidate pool of spoken words.

Evaluation: Submissions will be scored using classification accuracy across a limited set of keywords.

Benchmark: Training data selection

Task: Keyword spotting

Dataset: The Multilingual Spoken Words Corpus



Challenge 3: Vision | Training Data Cleaning

By Xiaozhe Yao and Ce Zhang (ETH Zürich)

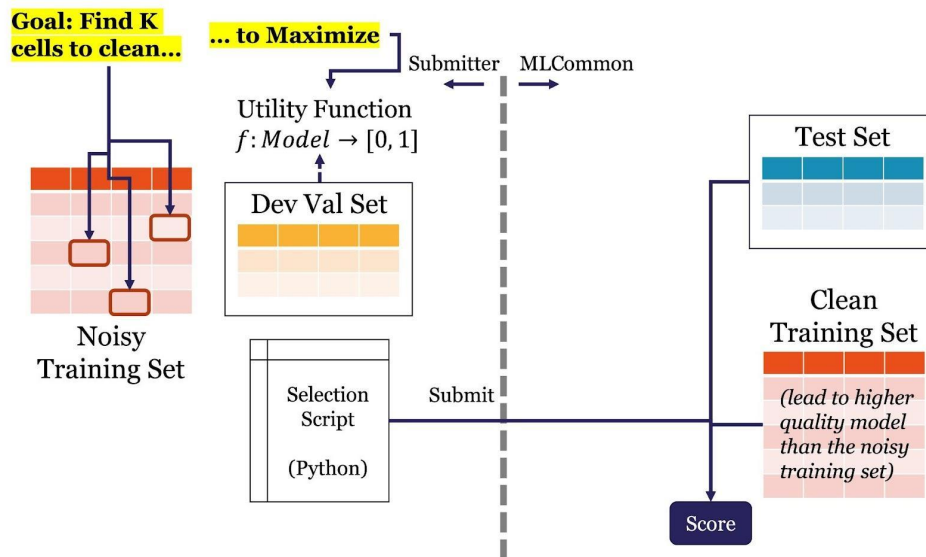
Challenge: Design a data cleaning strategy that chooses samples to relabel from a noisy training set.

Evaluation: Submissions will be scored on the minimum number of cleaned samples needed to achieve an accuracy threshold across a set of image classification tasks.

Benchmark: Training data label cleaning

Task: Image classification

Dataset: Custom subset of the Open Images Dataset with noisy labels



Challenge 4: Nibbler | Safe AI

By Xiaozhe Yao and Ce Zhang (ETH Zürich)

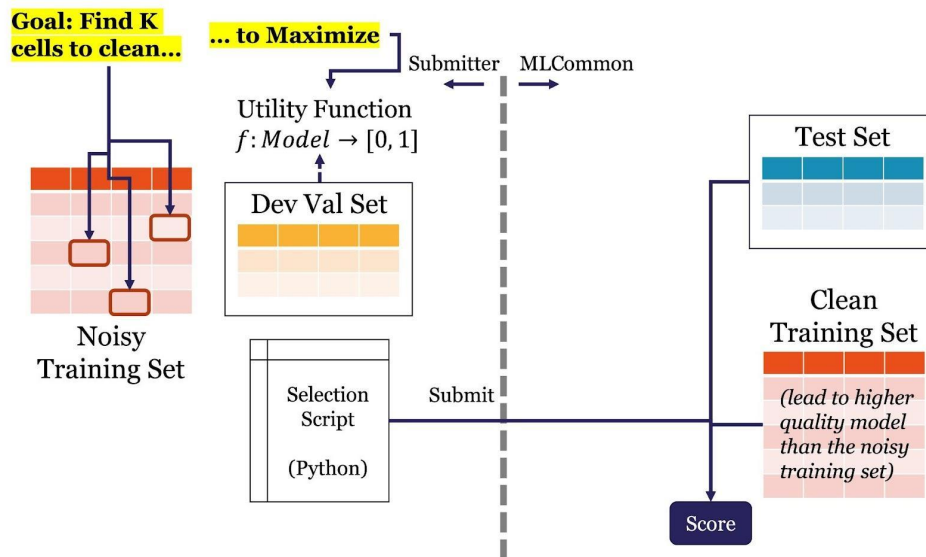
Challenge: collect prompts that are likely to cause a generative text-to-image model to fail in an unsafe manner (i.e., safety policy violations)

Evaluation:

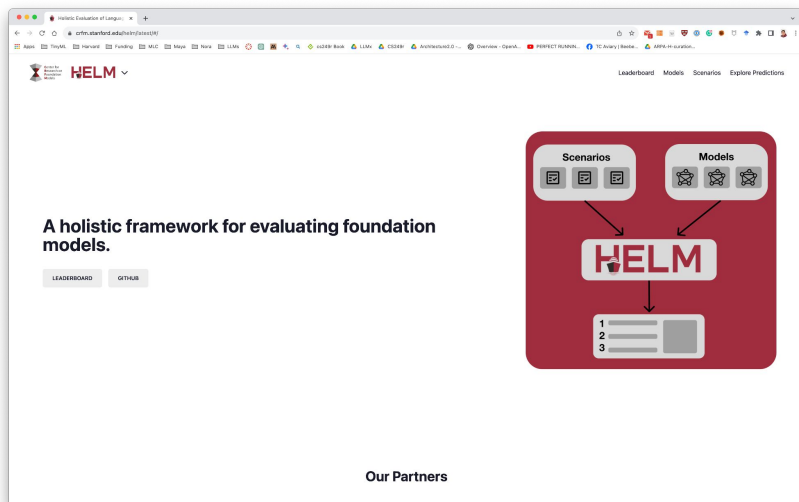
Benchmark: Training data label cleaning

Task: Image classification

Dataset: Custom subset of the Open Images Dataset with noisy labels



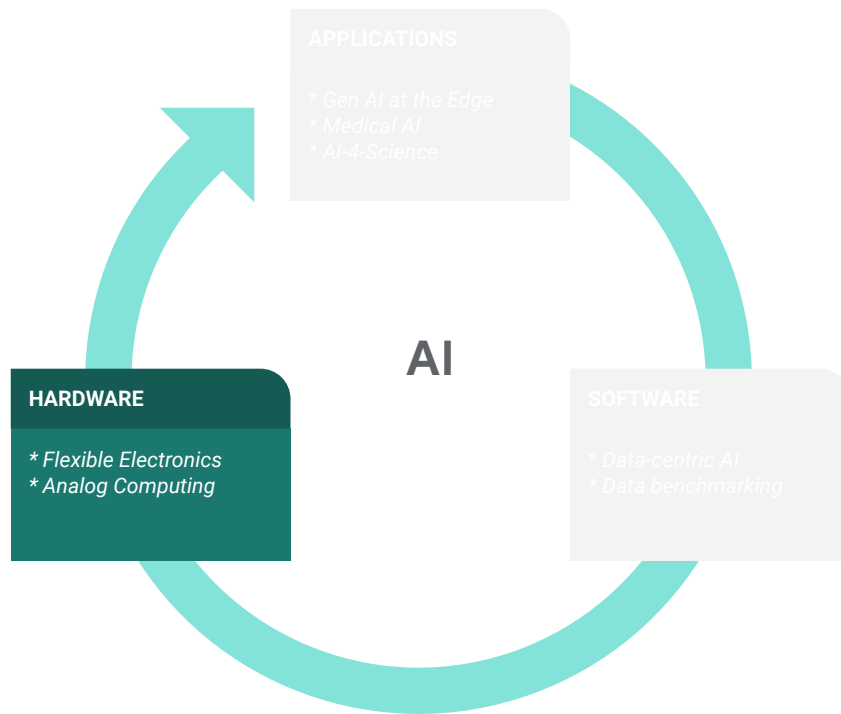
Challenges with Benchmarking GenAI



The screenshot shows the HELM Leaderboard page. The header includes the HELM logo and navigation links: Leaderboard, Models, Scenarios, and Explore Predictions. The main content area features the text "Leaderboard" and a description: "The leaderboard shows how the various models (with particular adaptation procedures) perform across different groups of scenarios and different metrics." Below this is a table with columns for Accuracy, Calibration, Robustness, Fairness, Efficiency, General information, Bias, Toxicity, and Summarization metrics. The table lists various models and their performance metrics across different scenarios.

Model(adapter)	Mean win rate	MMLU - EM	BoolQ - EM	NarrativeQA - F1	NaturalQuestions (closed-book)
Llama 2 (70B)	0.944	0.582	0.886	0.77	0.458
LLaMA (65B)	0.908	0.584	0.871	0.755	0.431
text-davinci-002	0.905	0.568	0.877	0.727	0.383
Mistral v0.1 (7B)	0.884	0.572	0.874	0.716	0.365
Cohere Command beta (52.4B)	0.874	0.452	0.856	0.752	0.372
text-davinci-003	0.872	0.569	0.881	0.727	0.406
Jurassic-2 Jumbo (178B)	0.824	0.48	0.829	0.733	0.385
Llama 2 (13B)	0.823	0.507	0.811	0.744	0.376
TNLG v2 (530B)	0.787	0.469	0.809	0.722	0.384
gpt-3.5-turbo-0613	0.783	0.391			

What's Next



Introduction & Motivation



Introduction & Motivation



Introduction & Motivation



Introduction & Motivation



Introduction & Motivation



Introduction & Motivation

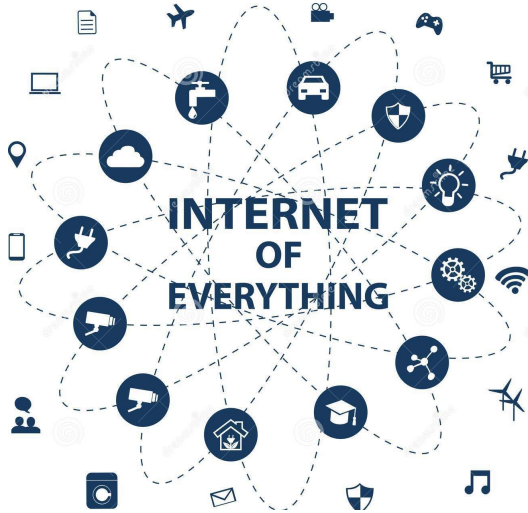


Internet-of-*Everything*

TinyML Connection

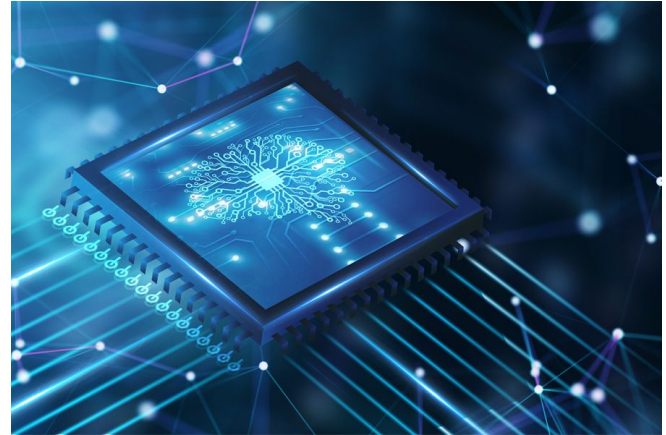


TinyML Connection



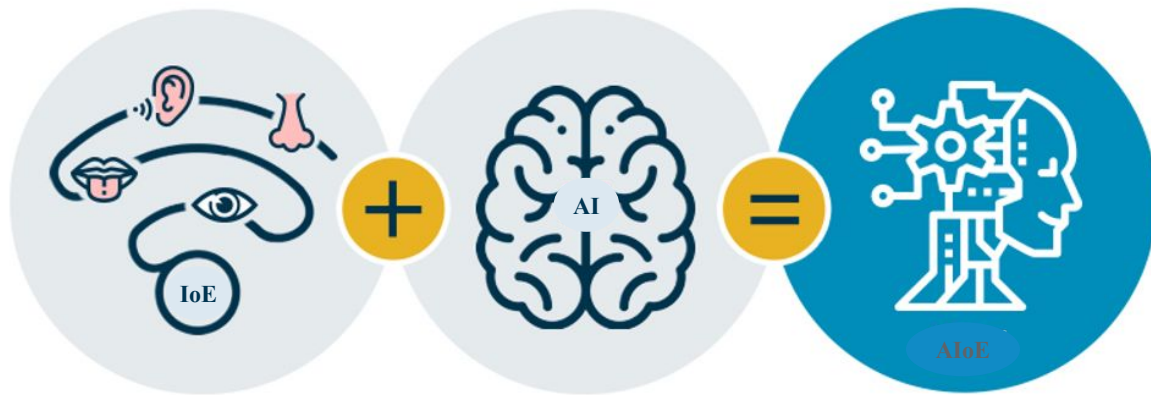
IoE

+



TinyML

TinyML Connection

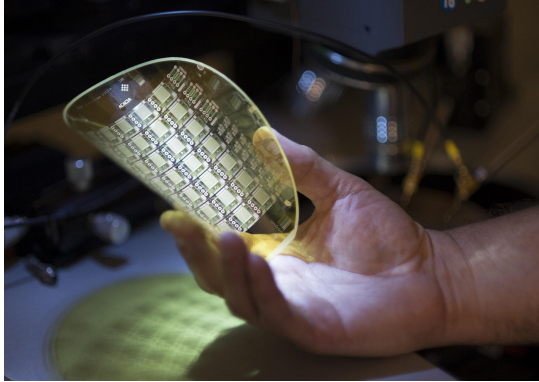


IoE (sensory organs) eyes, nose, ears, skin, tongue + AI (brain) = AIoE

Intelligent Internet-of-Everything

Flexible Electronics

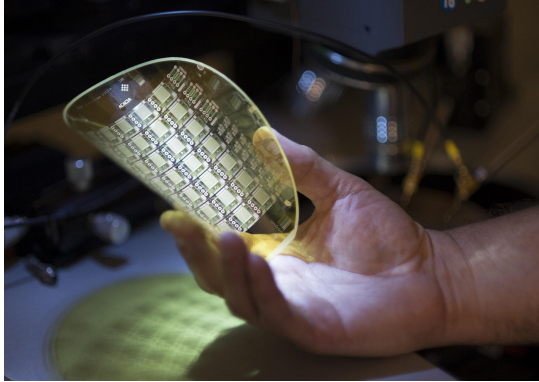
Potential of Flexible Integrated Circuits (FlexICs)



Conformability

Flexible Electronics

Potential of Flexible Integrated Circuits (FlexICs)



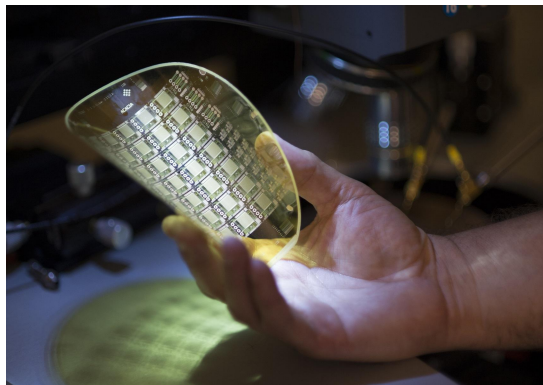
Conformability



Cost

Flexible Electronics

Potential of Flexible Integrated Circuits (FlexICs)



Conformability



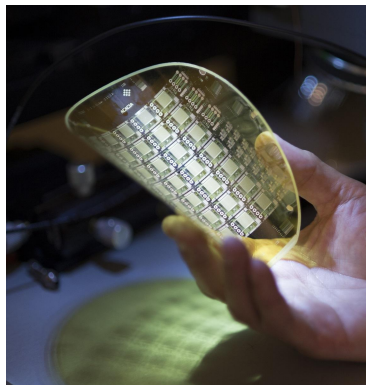
Cost



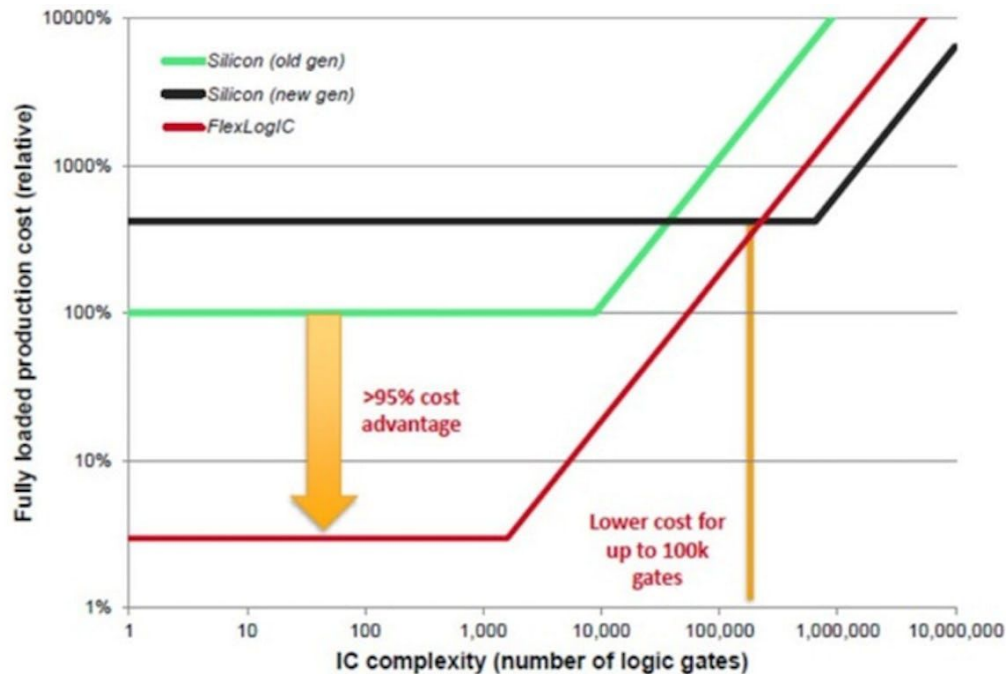
Footprint

Flexible Electronics

Potential of Flexible Integrated Circuits (FlexICs)



Conformab

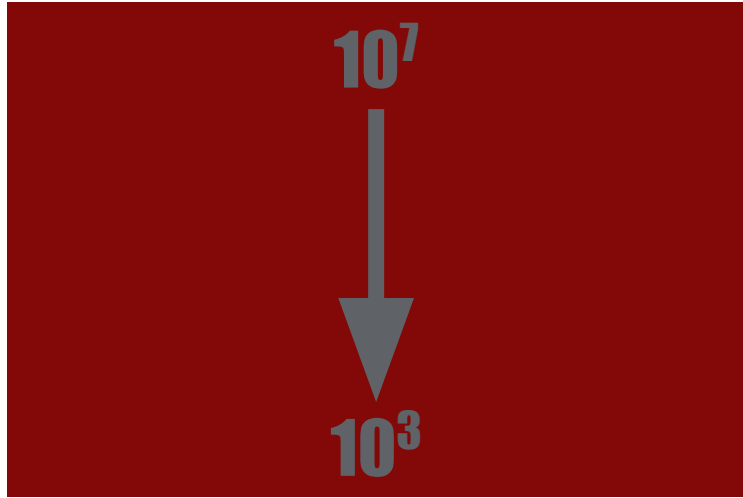


footprint

NB: analysis assumes footprint is limited by complexity not I/O pads (more pads would increase inflection & crossover point)

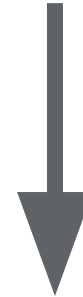
Flexible Electronics

Challenges of Flexible Integrated Circuits (FlexICs)



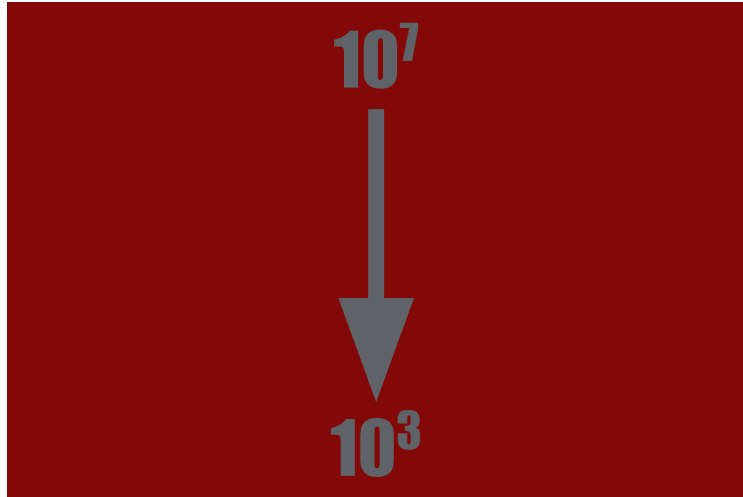
Low Yield □ *Small Designs*

MHz-GHz



Flexible Electronics

Challenges of Flexible Integrated Circuits (FlexICs)



Low Yield □ *Small Designs*

Low Performance



Existing Applications for FlexICs

Low Performance Use-Cases

Table 1: Example applications and their performance/precision requirements

Application	Sample Rate (Hz)	Prec. (bits)	Duty Cycle Period	Application	Sample Rate (Hz)	Prec. (bits)	Duty Cycle Period
Blood Pressure Sensor [70]	< 100	< 8	Hours [19]	Body Temperature Sensor [70]	< 1	< 8	Minutes [44]
Odor Sensor [70]	16-25	< 8	Minutes [73]	Smart Bandage [65]	< 0.01	< 8	Continuous to Hours [23]
Heart Beat Sensor [70]	< 4	1	Seconds [90]	Tremor Sensor [33]	< 25	16	Seconds [20]
Pressure Sensor [31]	1-5.5	12	Continuous to Hours [81]	Oral-Nasal Airflow [70]	< 25	< 8	Seconds
Light Level Sensor [70]	< 1	< 8	Continuous to Hours [22]	Perspiration Sensor [47]	< 25	< 8	Minutes [99]
Trace Metal Sensor [47]	25	16	Minutes	Pedometer [72]	< 25	1	Seconds [72]
Food Temp. Sensor [70]	< 1	< 8	5 minutes [82]	Timer [40]	1	1	Single Use
Alcohol Sensor [48]	1	< 8	Single Use [64]	POS Computation [63]	< 100	< 8	Single Use [63]
Humidity Sensor [34]	10	16	Continuous to Hours [80]	Smart Labels [7]	1	< 8	Seconds
Pseudo-RNG	n/a	< 8	Seconds	Error Detection Coding	< 100	< 8	Continuous to Hours

Tiny Processor

Small Design for High Yield

- **SERV:** World's Smallest RISC-V Processor

Tiny Processor

Small Design for High Yield

- **SERV:** World's Smallest RISC-V Processor
- **Bit-serial** MCU

Tiny Processor

Small Design for High Yield

- **SERV**: World's Smallest RISC-V Processor
- **Bit-serial** MCU
- Supports **RV32I** (~40 instructions)

Tiny Processor

Small Design for High Yield

- **SERV**: World's Smallest RISC-V Processor
 - **Bit-serial** MCU
 - Supports **RV32I** (~40 instructions)
 - **M (multiplication)** extension added

Tiny Processor

Small Design for High Yield

- **SERV**: World's Smallest RISC-V Processor
 - **Bit-serial** MCU
 - Supports **RV32I** (~40 instructions)
 - **M (multiplication)** extension added

Simple Example

$$2 + 3 = 5$$

Tiny Processor

Small Design for High Yield

- **SERV**: World's Smallest RISC-V Processor
 - **Bit-serial** MCU
 - Supports **RV32I** (~40 instructions)
 - **M (multiplication)** extension added

Simple Example

$$2 + 3 = 5$$

$$00000010 + 00000011 = 00000101$$

Tiny Processor

Small Design for High Yield

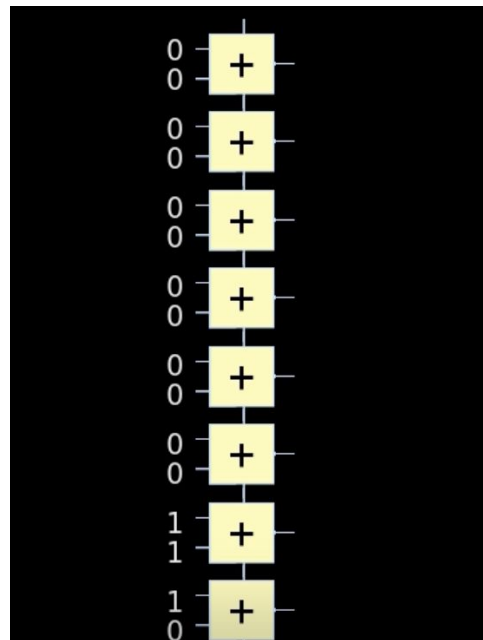
- **SERV**: World's Smallest RISC-V Processor
- **Bit-serial** MCU
- Supports **RV32I** (~40 instructions)
- **M (multiplication)** extension added

Simple Example

$$2 + 3 = 5$$

$$00000010 + 00000011 = 00000101$$

Modern
Bit-Parallel Design



Tiny Processor

Small Design for High Yield

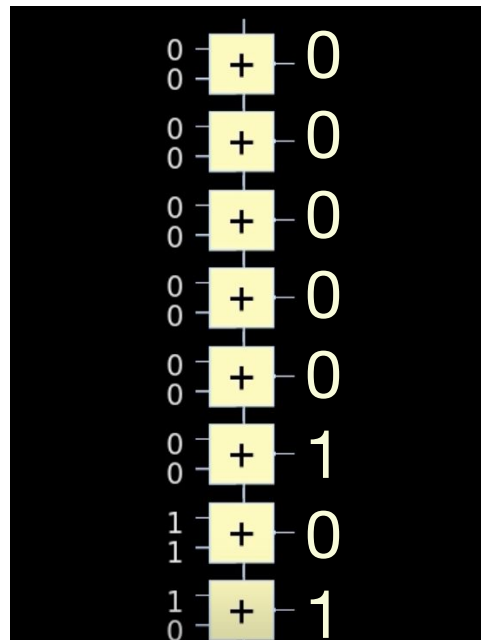
- **SERV**: World's Smallest RISC-V Processor
- **Bit-serial** MCU
- Supports **RV32I** (~40 instructions)
- **M (multiplication)** extension added

Simple Example

$$2 + 3 = 5$$

$$00000010 + 00000011 = 00000101$$

Modern
Bit-Parallel Design



Tiny Processor

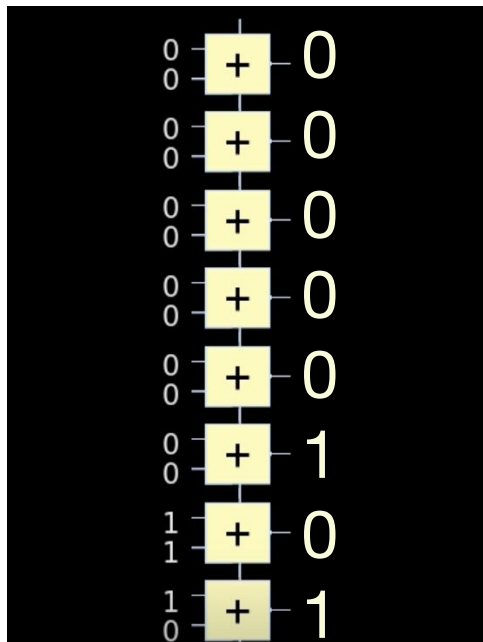
Small Design for High Yield

- **SERV**: World's Smallest RISC-V Processor
- **Bit-serial** MCU
- Supports **RV32I** (~40 instructions)
- **M (multiplication)** extension added

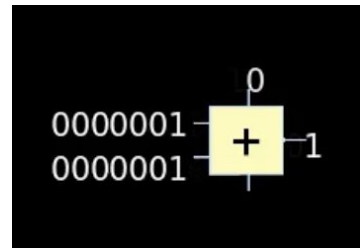
Simple Example

$$2 + 3 = 5$$

$$00000010 + 00000011 = 00000101$$



Less Hardware, Longer Latency

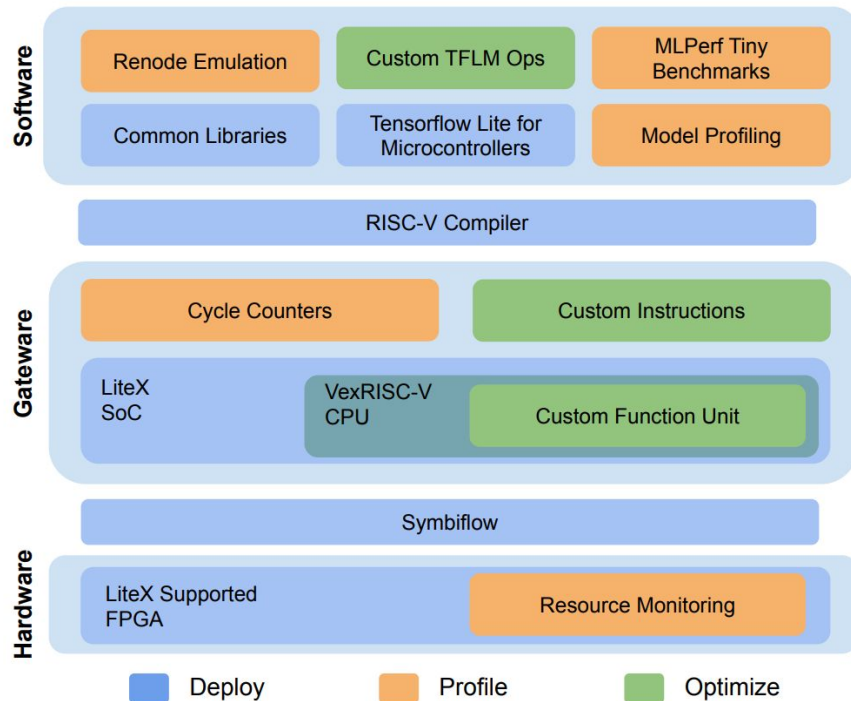
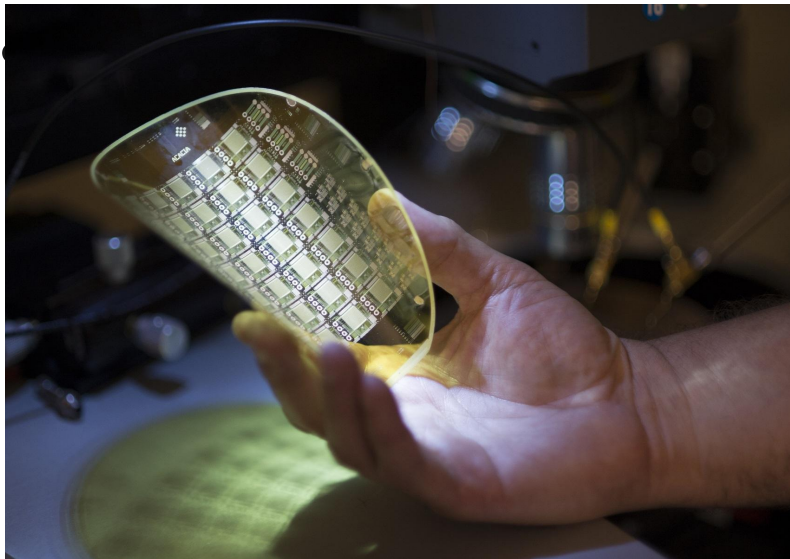


Bit-Serial Design

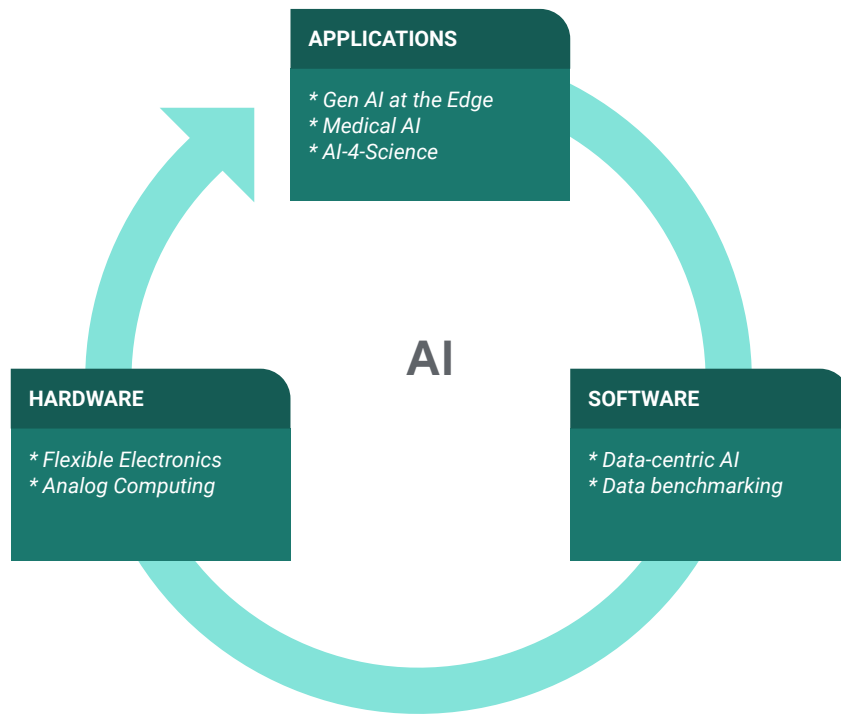
FlexICs for TinyML Updates & Status

Developer Tool Flow

- CFU Playground Integration

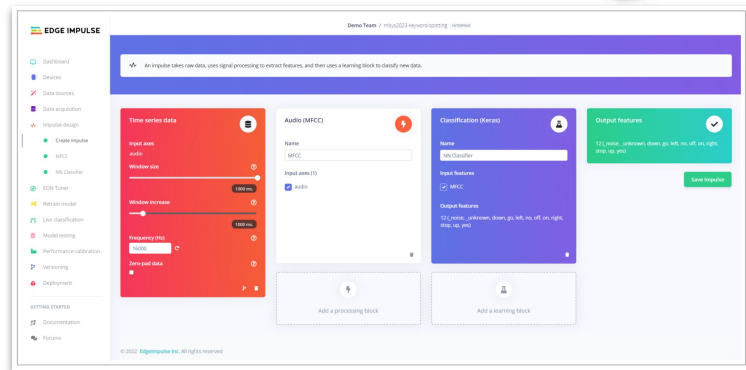
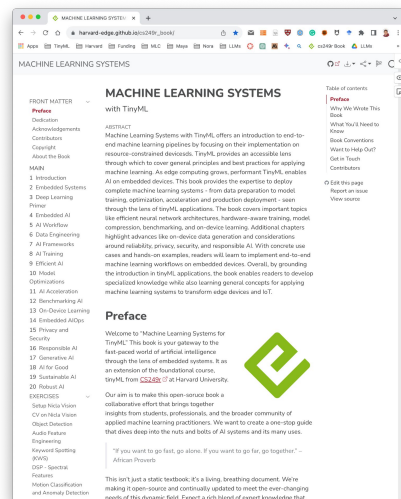


What's Next



Course Recap

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



Guest Speaker: Jason Wei

Jason Wei is an AI researcher living in San Francisco. He currently works at OpenAI on the ChatGPT team. Previously, he was a senior research scientist at Google Brain, where he popularized chain-of-thought prompting, co-led the first efforts on instruction tuning, and wrote about emergence in large language models.

Chain-of-thought prompting was presented by Google CEO Sundar Pichai at the Google I/O press event in 2022.



[Google Scholar](#)
[Personal Website](#)

