



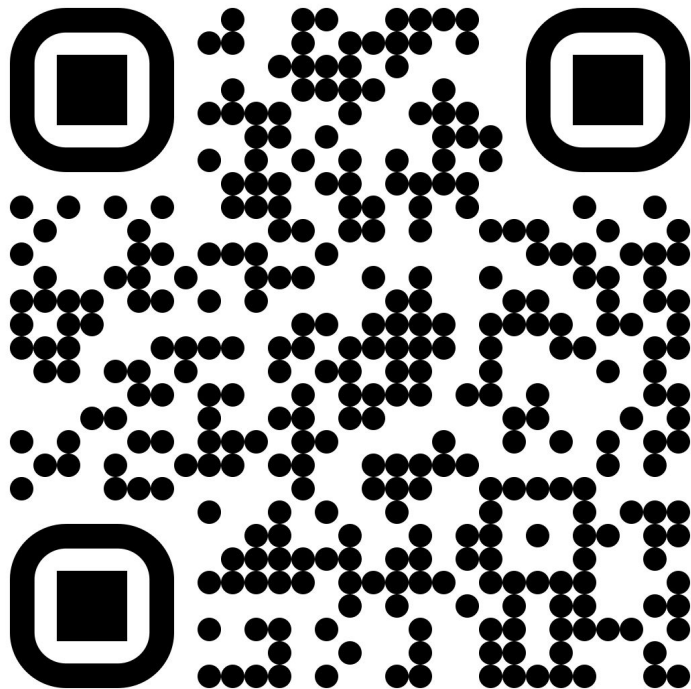
# Edge AI Workshop

Rust Workshop





# Materials

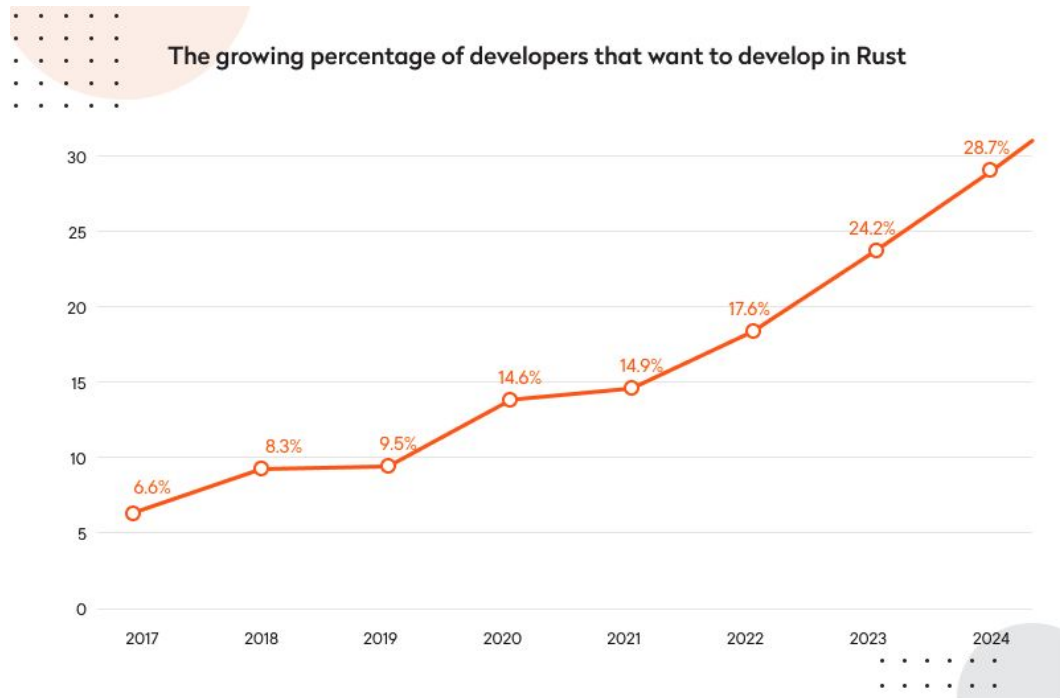




# Interest in Rust

- Memory Safe
- Strongly Typed
- Extensive Ecosystem  
([crates.io](https://crates.io))

*Java Styles at C/C++ speed*

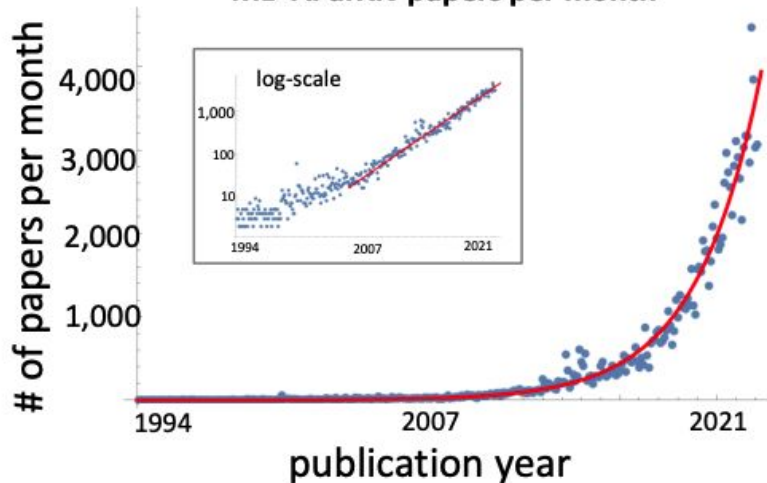




# Interest in AI - It keeps growing, but for how long?

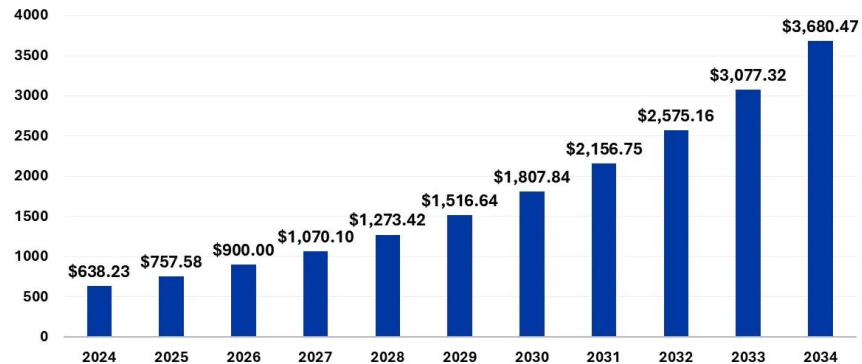
- Exponential number of new papers each year!
  - NeurIPS (“best” AI conference) 2025: a record of **25,000** submissions
- More funding, investments for AI companies

ML+AI arXiv papers per month



Precedence  
RESEARCH

Artificial Intelligence (AI) Market Size 2024 to 2034 (USD Billion)

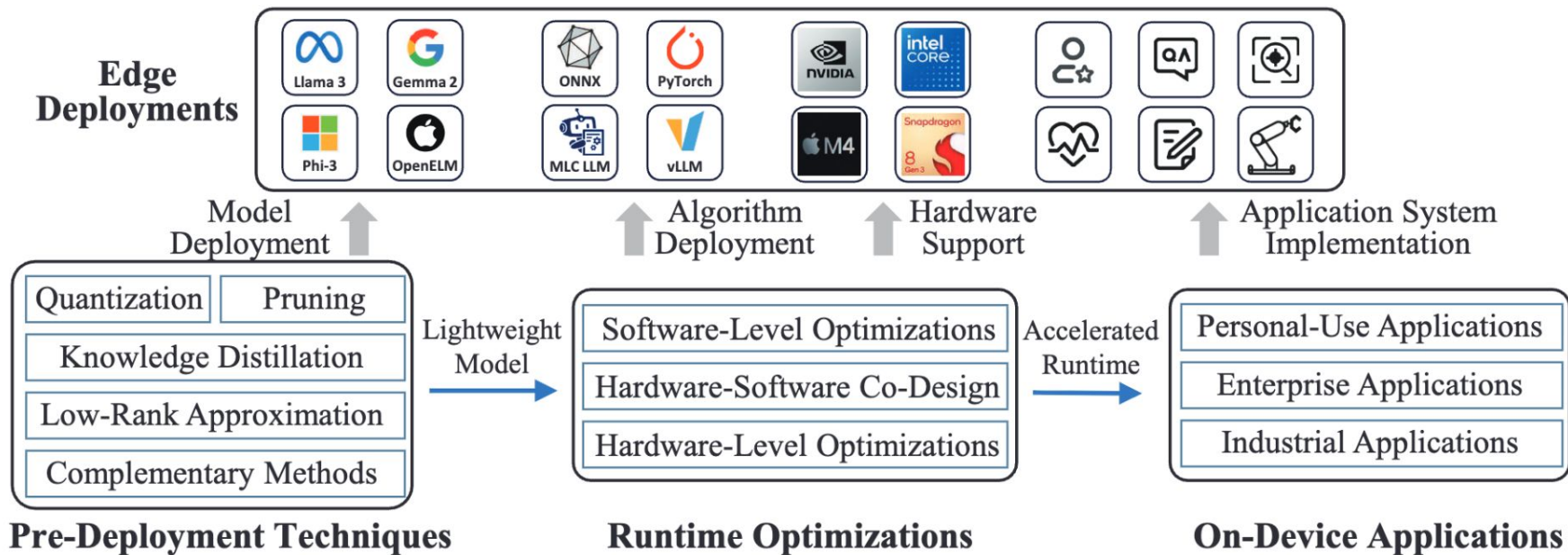


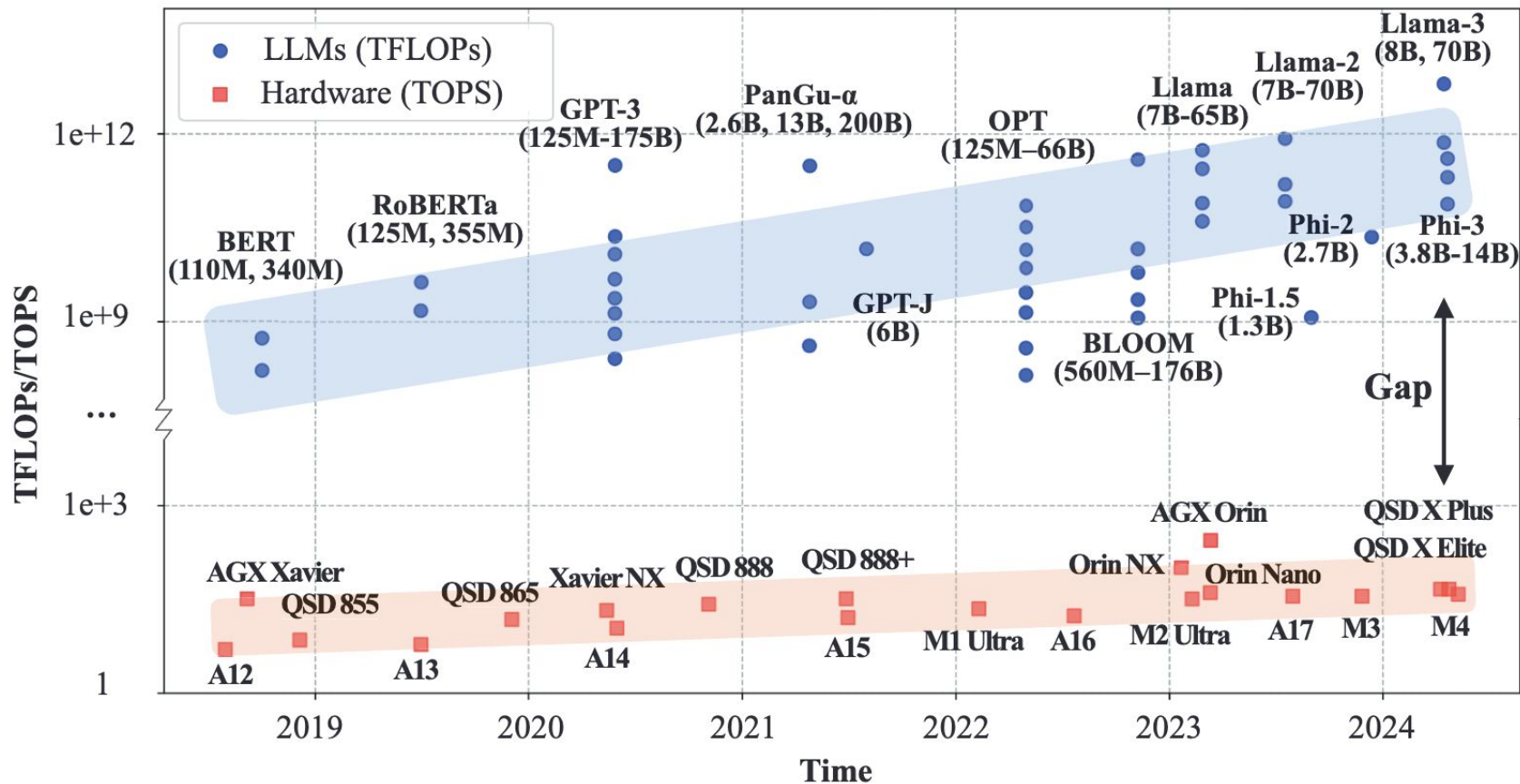
Source: <https://www.precedenceresearch.com/artificial-intelligence-market>



# Edge AI

- Deployment of AI models directly on local devices or “edge devices”
  - like sensors, cameras, IoT gadgets, smartphones, industrial machinery
  - Raspberry Pi
- Self-driving cars
- Drones
- Smartphones
- etc.







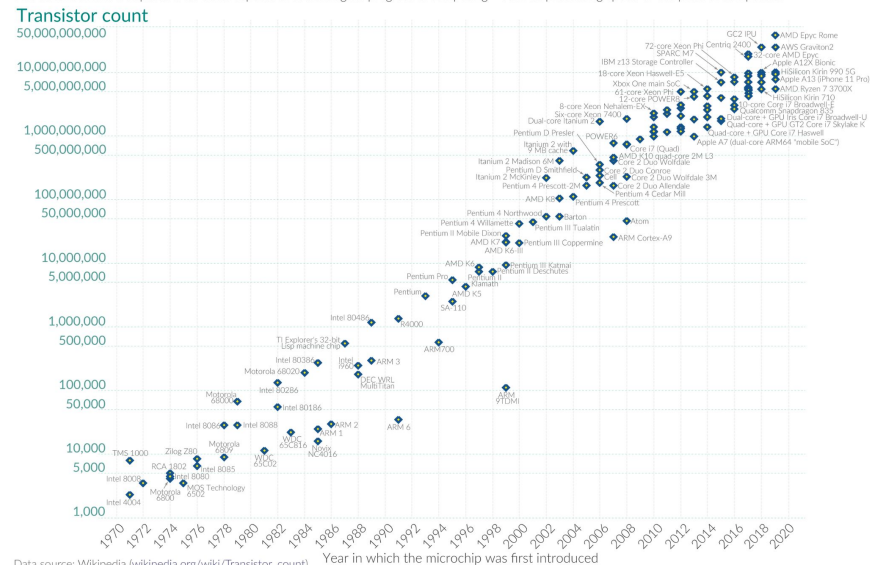
# Edge AI - Why is it possible?

## - Moore's Law

- Just 6-7 years ago, running a real-time ~50M parameter model on edge devices was a struggle
- Now, we can run 1B parameter models
  - ~10x more computing power

## Moore's Law: The number of transistors on microchips doubles every two years

Moore's law describes the empirical regularity that the number of transistors on integrated circuits doubles approximately every two years. This advancement is important for other aspects of technological progress in computing – such as processing speed or the price of computers.





- Moore's Law

- Just 6-7 years ago, running a real-time ~50M parameter model on edge devices was a struggle
  - Now, we can run 1B parameter models
    - ~10x more computing power
- Moore's Law: The number of transistors on microchips doubles every two years.  
Moore's law describes the empirical regularity that the number of transistors on integrated circuits doubles approximately every two years.  
This advancement is important for other aspects of technological progress in computing – such as processing speed or the amount of data that can be stored.  
Transistor count

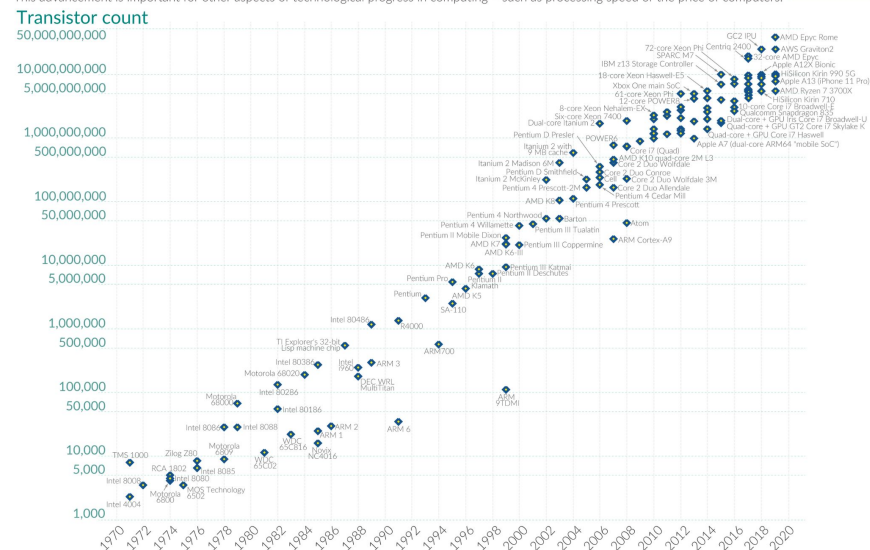
- Not just due to hardware advancements!

- We have better inference techniques:

- Quantization, pruning
- **Distillation!**



Moore's law describes the empirical regularity that the number of transistors on integrated circuits doubles approximately every two years. This advancement is important for other aspects of technological progress in computing – such as processing speed or the price of computers.

Our World  
in Data

Data source: Wikipedia ([wikipedia.org/wiki/Transistor\\_count](http://wikipedia.org/wiki/Transistor_count))

OurWorldinData.org – Research and data to make progress against the world's largest problems

Licensed under CC-BY by the authors Hannah Ritchie and Max Roser.



# Ingredients for Edge AI



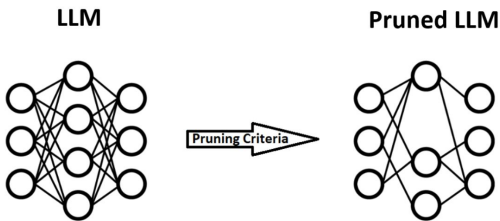
- Good, cheap hardware
- If you don't have it, just wait a couple of years.



# Ingredients for Edge AI



- Good, cheap hardware
- If you don't have it, just wait a couple of years.
- Model pruning
- Get a big model and make it smaller by deleting neurons / layers

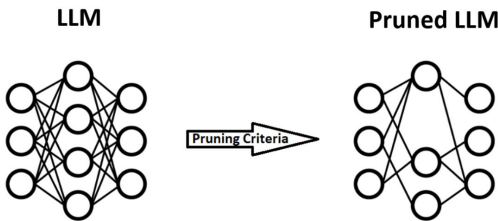




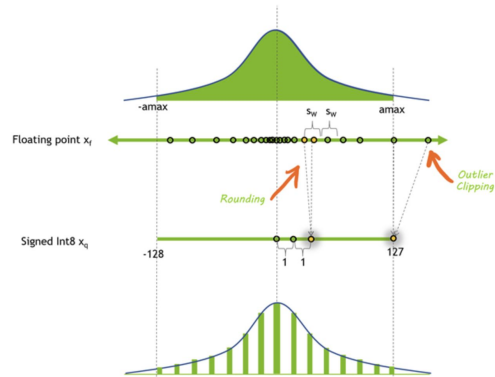
# Ingredients for Edge AI



- Good, cheap hardware
- If you don't have it, just wait a couple of years.



- Model pruning
- Get a big model and make it smaller by deleting neurons / layers

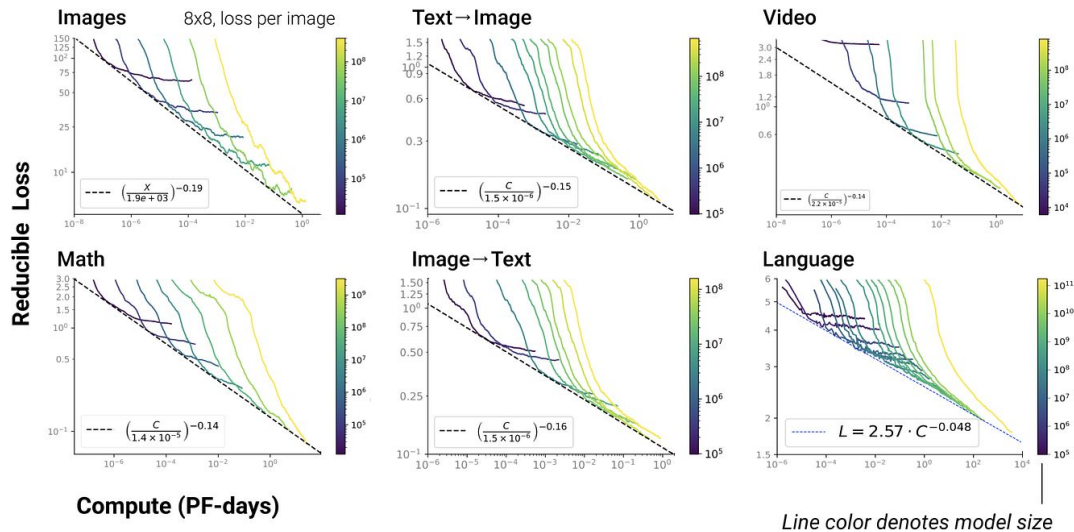


- Reduce precision of weights to save memory / inference time



# The real reason: We scaled-up

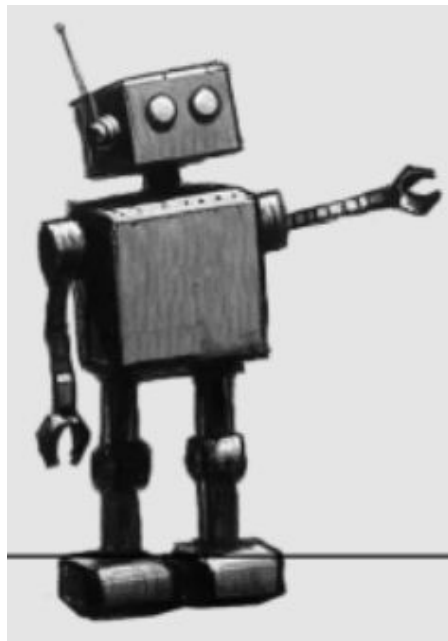
- Consistent in AI: **scale** leads to predictable improvements in capability
  - Model size, dataset size, compute
- Moore's Law enables us to scale faster
- Paradoxically, scaling up also facilitates scaling down!





# The real reason: We scaled-up and distilled the knowledge

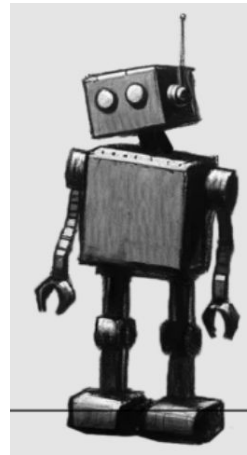
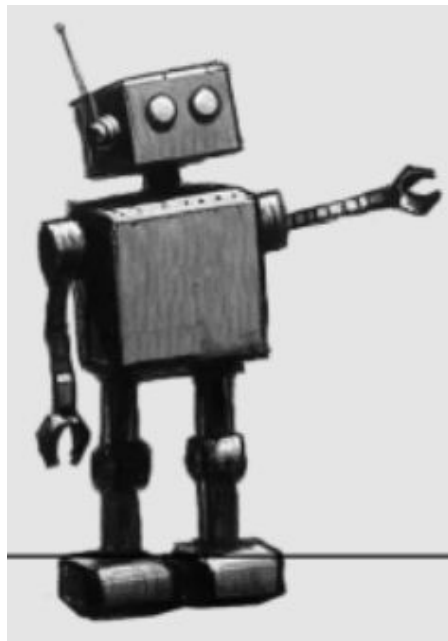
- We trained larger and larger and more capable models
  - GPT-2, GPT-3, GPT-4, GPT-5, ...
  - Deepseek-R1 (~400B params)
  - Qwen-480B





# The real reason: We scaled-up and distilled the knowledge

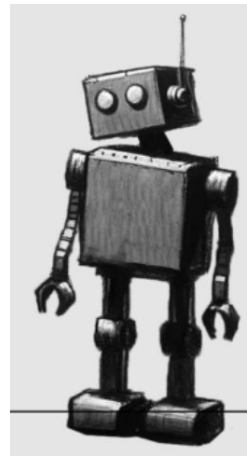
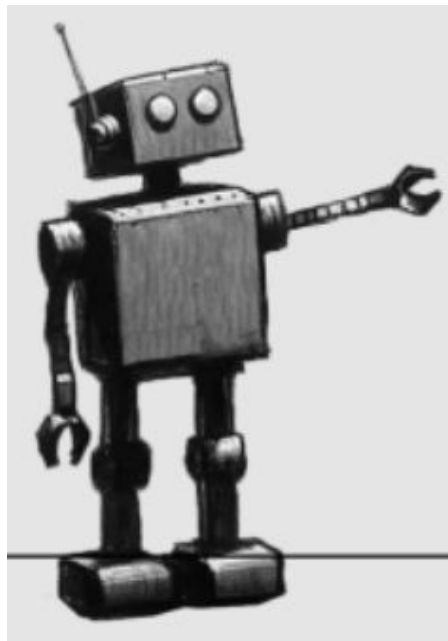
- We trained larger and larger and more capable models
  - GPT-2, GPT-3, GPT-4, GPT-5, ...
  - Deepseek-R1 (~400B params)
  - Qwen-480B
- **Large models** can then teach **smaller faster** versions of themselves
  - gpt-mini, gpt-nano
  - deepseek -7B
  - qwen-1B, qwen-3B





# The real reason: We scaled-up and distilled the knowledge

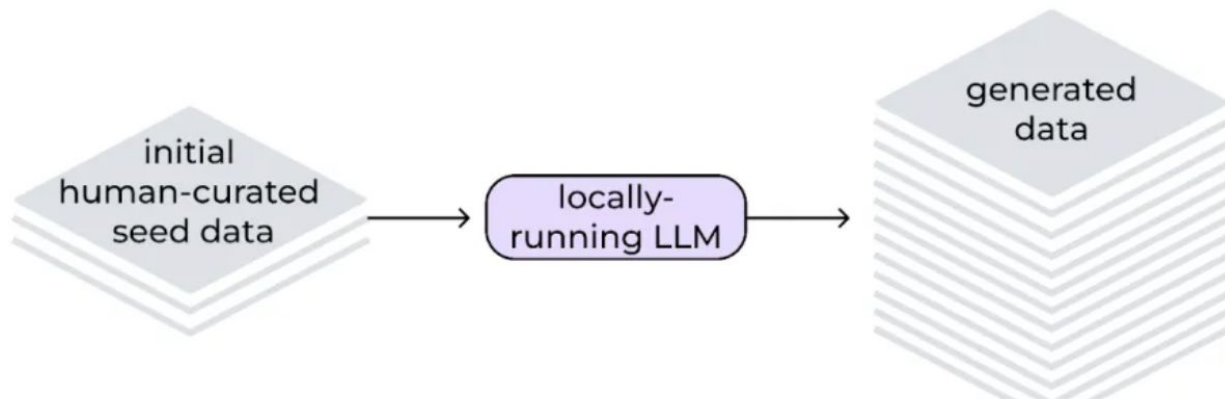
- We trained larger and larger and more capable models
  - GPT-2, GPT-3, GPT-4, GPT-5, ...
  - Deepseek-R1 (~400B params)
  - Qwen-480B
- **Large models** can then teach **smaller faster** versions of themselves
  - gpt-mini, gpt-nano
  - deepseek -7B
  - qwen-1B, qwen-3B
- Small models reach a performance level that cannot be otherwise obtained





# Nowadays, “distillation” has be rebranded

- “**Synthetic Data**”
- Use an (very) large, highly capable model to generate more (high-quality) data / clean existing data
  - Paraphrasing
  - Generate instructions
  - **Enrich dataset** by automatic annotations





# Sounds good but ...

Bitter lesson: progress of AI in the past 70 years boils down to

- Develop progressively more general methods with weaker modeling assumptions
- Add more data and computation (i.e. scale up)

## The Bitter Lesson

**Rich Sutton**

**March 13, 2019**

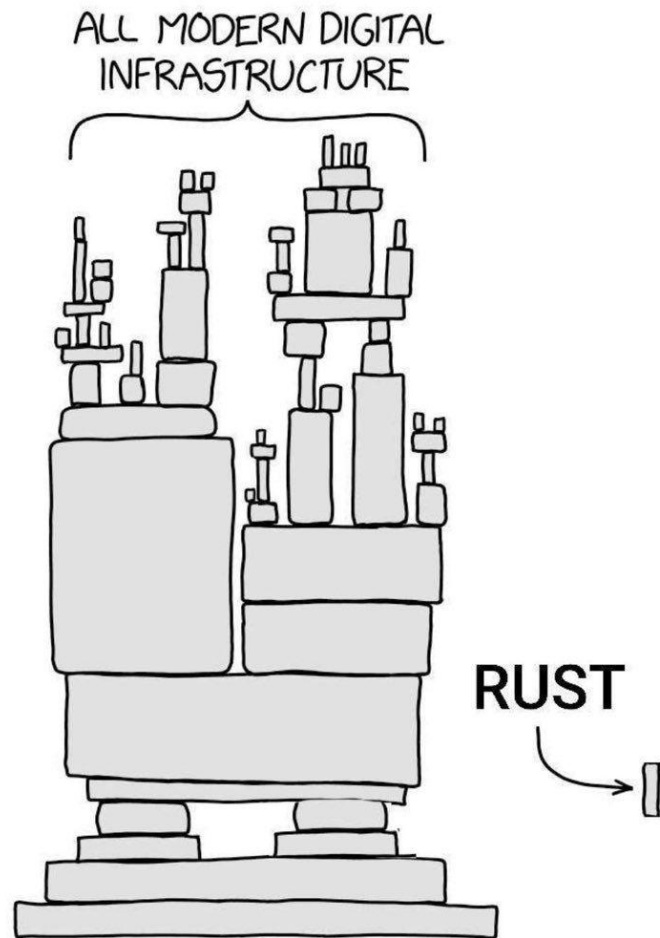
The biggest lesson that can be read from 70 years of AI research is that general methods that leverage computation are ultimately the most effective, and by a large margin. The ultimate reason for this is Moore's law, or rather its generalization of continued exponentially falling cost per unit of computation. Most AI research has been conducted as if the computation available to the agent were constant (in which case leveraging human knowledge would be one of the only ways to improve performance) but, over a slightly longer time than a typical research project, massively more computation inevitably becomes available. Seeking an improvement that makes a difference in the shorter term, researchers seek to leverage their human knowledge of the domain, but the only thing that matters in the long run is the leveraging of computation. These two need not run counter to each other, but in practice they tend to. Time spent on one is time not spent on the other. There are





# AI + Rust

- Edge AI
  - Tokenizers
  - WASM
  - Raspberry PIs
- 
- In practice, Rust is useful for building high-performance data pipelines
- 
- For example: 🤗 **HuggingFace** tokenizers library is built on top of Rust





# Common ML Frameworks (non Rust)

## Training

- Pytorch (Meta)
- Tensorflow (Google)
- JAX (Google)

## Inference

- GGML
- ONNX Runtime (Microsoft)
- LiteRT (Google)
- Executorch (Meta)
- TVM (Apache)
- TF Micro (Google)



# Why different Frameworks for Inference and Training?

## Training

- Focus on flexibility
- Feature-rich
- GPU first

(most in Python)

## Inference

- Focus on speed and size
- Compiled into apps / devices
- Cross platform
- GPU, CPU, NPU, MCU

(most in C++)



# Specialized Inference Frameworks

- Cadence HiFi NN Lib
- ARM-NN
- Apple Core ML
- Rockchip NPU
- Huawei CANN
- Android NNAPI
- Intel OpenVino
- Xilinx Vitis-AI
- Qualcomm QNN

Everyone builds their own inference engine...



# How Rust can help

Everything that makes Rust great in other use-cases also applies to ML:

- Easy cross-compilation
- Great optimization of the same code on different architectures
- Memory Safety
- Awesome Tooling
- Great abstractions of complex patterns



# Rust ML Frameworks

**Burn**

(Tracel AI)



**tract**

(Sonos)



(Pyke.io)

**Candle**



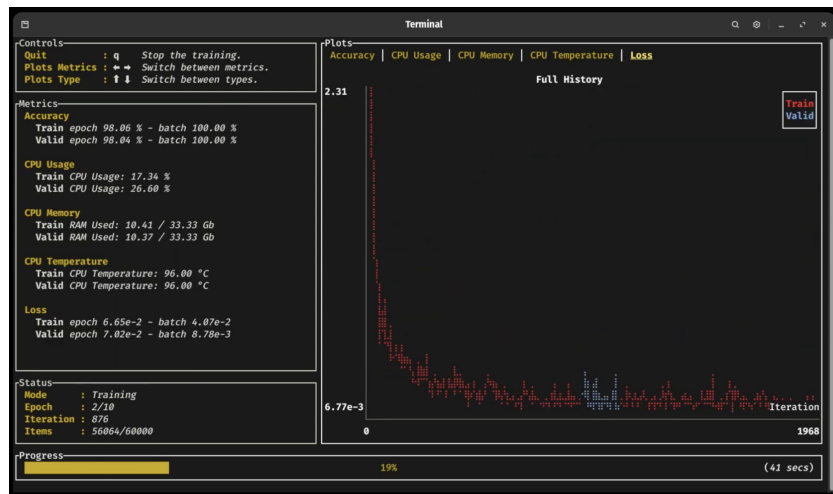
**Hugging Face**



# Burn (Tracel AI)

- Training & Inference
- GPU, CPU & MCU (no\_std)
- Own GPU compute language (CubeCL)
- GPU works on NVIDIA, AMD, Intel & Web
- Only basic ONNX import
- GPU and larger model focus
- Especially slow for embedded platforms
- No Accelerator support (NPU)

# Burn





# Burn ONNX Import

| ONNX OP                            | Import Support | Burn Support |
|------------------------------------|----------------|--------------|
| <a href="#">Abs</a>                | ✓              | ✓            |
| <a href="#">Acos</a>               | ✗              | ✗            |
| <a href="#">Acosh</a>              | ✗              | ✗            |
| <a href="#">Add</a>                | ✓              | ✓            |
| <a href="#">And</a>                | ✓              | ✓            |
| <a href="#">ArgMax</a>             | ✓              | ✓            |
| <a href="#">ArgMin</a>             | ✓              | ✓            |
| <a href="#">Asin</a>               | ✗              | ✗            |
| <a href="#">Asinh</a>              | ✗              | ✗            |
| <a href="#">Atan</a>               | ✗              | ✗            |
| <a href="#">Atanh</a>              | ✗              | ✗            |
| <a href="#">Attention</a>          | ✓              | ✓            |
| <a href="#">AveragePool1d</a>      | ✓              | ✓            |
| <a href="#">AveragePool2d</a>      | ✓              | ✓            |
| <a href="#">BatchNormalization</a> | ✓              | ✓            |
| <a href="#">Bernoulli</a>          | ✓              | ✓            |
| <a href="#">BitShift</a>           | ✓              | ✓            |

|                                   |   |   |
|-----------------------------------|---|---|
| <a href="#">RNN</a>               | ✗ | ✓ |
| <a href="#">RoiAlign</a>          | ✗ | ✗ |
| <a href="#">Round</a>             | ✓ | ✓ |
| <a href="#">Scan</a>              | ✗ | ✗ |
| <a href="#">Scatter</a>           | ✗ | ✓ |
| <a href="#">ScatterElements</a>   | ✗ | ✗ |
| <a href="#">ScatterND</a>         | ✗ | ✗ |
| <a href="#">Selu</a>              | ✗ | ✗ |
| <a href="#">SequenceAt</a>        | ✗ | ✗ |
| <a href="#">SequenceConstruct</a> | ✗ | ✗ |
| <a href="#">SequenceEmpty</a>     | ✗ | ✗ |
| <a href="#">SequenceErase</a>     | ✗ | ✗ |
| <a href="#">SequenceInsert</a>    | ✗ | ✗ |
| <a href="#">SequenceLength</a>    | ✗ | ✗ |
| <a href="#">SequenceMap</a>       | ✗ | ✗ |
| <a href="#">Shape</a>             | ✓ | ✓ |
| <a href="#">Shrink</a>            | ✗ | ✗ |
| <a href="#">Sigmoid</a>           | ✓ | ✓ |
| <a href="#">Sign</a>              | ✓ | ✓ |
| <a href="#">Sin</a>               | ✓ | ✓ |
| <a href="#">Sinh</a>              | ✓ | ✓ |

[...]



# Candle (Hugging Face)



**Hugging Face**

- Good selection of already implemented models
- Many examples
- Support for CUDA, Metal, Intel MKL, Apple Accelerate
- Programmatically first
- Only very basic ONNX import
- No universal GPU support (AMD, Intel, Web)
- No embedded accelerators

```
use candle_core::{Device, Tensor};

fn main() -> Result<(), Box<dyn std::error::Error>> {
    let device = Device::Cpu;

    let a = Tensor::randn(0f32, 1., (2, 3), &device)?;
    let b = Tensor::randn(0f32, 1., (3, 4), &device)?;

    let c = a.matmul(&b)?;
    println!("{c}");
    Ok(())
}
```



# Tract (Sonos)

- Pure Rust
- Support of most ONNX ops
- Super easy import
- Optimized for smaller models on embedded devices (Raspberry Pi)
- Great CLI Test and Benchmarking Tool
  
- CPU only
- No no\_std support





# ORT (pyke.io)

- Wrapper around ONNX Runtime (C++)
  - Support 100% of ONNX operators
  - Super easy to use with prebuilt static libraries
  - Supports all GPUs and many accelerators
  - Super fast on CPUs and embedded
  - Battle tested and optimized
- 
- C++ wrapping can make things harder
  - Builds for accelerators can be tricky





# ORT - Execution Providers

| EP                                   | Cargo feature | Supported | Binaries |
|--------------------------------------|---------------|-----------|----------|
| <a href="#">NVIDIA CUDA</a> ↗        | cuda          | ◆         | ✓        |
| <a href="#">NVIDIA TensorRT</a> ↗    | tensorrt      | ◆         | ✓        |
| <a href="#">Microsoft DirectML</a> ↗ | directml      | ◆         | ✓        |
| <a href="#">Apple CoreML</a> ↗       | coreml        | ◆         | ✓        |
| <a href="#">AMD ROCm</a> ↗           | rocm          | ◆         | ✗        |
| <a href="#">Intel OpenVINO</a> ↗     | openvino      | ◆         | ✗        |
| <a href="#">Intel oneDNN</a> ↗       | onednn        | ◆         | ✗        |
| <a href="#">XNNPACK</a> ↗            | xnnpack       | ◆         | ✓        |
| <a href="#">Qualcomm QNN</a> ↗       | qnn           | ◆         | ✗        |

|                                   |          |   |   |
|-----------------------------------|----------|---|---|
| <a href="#">Huawei CANN</a> ↗     | cann     | ◆ | ✗ |
| <a href="#">Android NNAPI</a> ↗   | nnapi    | ◆ | ✗ |
| <a href="#">Apache TVM</a> ↗      | tvm      | ◆ | ✗ |
| <a href="#">Arm ACL</a> ↗         | acl      | ◆ | ✗ |
| <a href="#">ArmNN</a> ↗           | armnn    | ◆ | ✗ |
| <a href="#">AMD MIGraphX</a> ↗    | migraphx | ◆ | ✗ |
| <a href="#">AMD Vitis AI</a> ↗    | vitis    | ◆ | ✗ |
| <a href="#">Rockchip RKNPU</a> ↗  | rknpu    | ◆ | ✗ |
| WebGPU                            | webgpu   | ◆ | ✓ |
| <a href="#">Microsoft Azure</a> ↗ | azure    | ◆ | ✗ |



# ORT - Upcoming features (v2.0.0-rc.10)

- Alternative Backends (Tract / Candle)
- Model Editor API (build models programmatically)
- no\_std support
- WebGPU support with prebuilt-libraries
- Static linking of CUDA and TensorRT libraries





# Workshop Overview

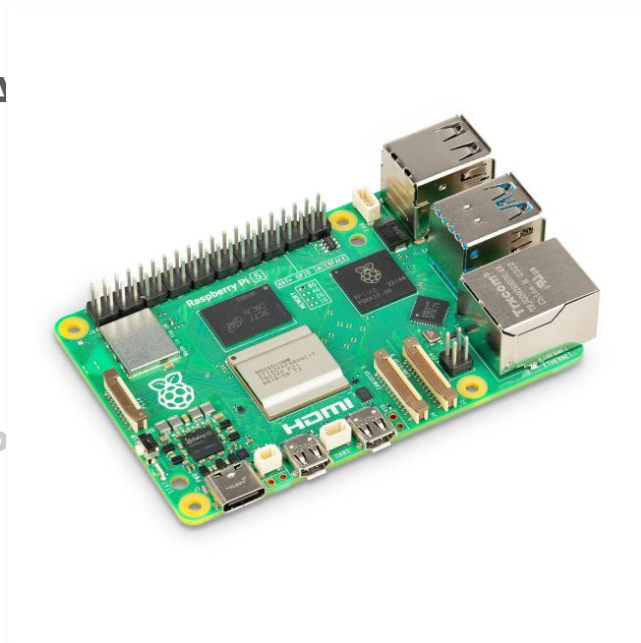
1. **Part 0: Introduction to the Team & Rust for Edge AI**
2. **Part I: Lecture on Computer Vision**
  - a. Main problems in Computer Vision
  - b. What exactly is a neural network? (CNNs / Transformers)
  - c. What exactly is an image embedding?
  - d. Computer vision on the Edge
3. **Hands-On I: Air-gapped Face recognition on the Pi**
4. **Part II: Lecture on Natural Language Processing**
  - a. A bit of history & development of modern LLMs
  - b. How does an LLM work? Tokenizers, pretraining, post-training
  - c. Context Engineering: Tool calling, RAG
  - d. Libraries: tokenizers-rs, llama.cpp
5. **Hands-On II: Chat with a LLM on Pi**





# Workshop Overview

1. **Part 0: Introduction to the Team & Rust for Edge AI**
2. **Part I: Lecture on Computer Vision**
  - a. Main problems in Computer Vision
  - b. What exactly is a neural network? (CNNs / Transformers)
  - c. What exactly is an image embedding?
  - d. Computer vision on the Edge
3. **Hands-On I: Air-gapped Face recognition on the Pi**
4. **Part II: Lecture on Natural Language Processing**
  - a. A bit of history & development of modern LLMs
  - b. How does an LLM work? Tokenizers, pretraining, post-training
  - c. Context Engineering: Tool calling, RAG
  - d. Libraries: tokenizers-rs, llama.cpp
5. **Hands-On II: Chat with a LLM on Pi**





# Workshop Overview

1. Part 0: Introduction to the Team & Rust for Edge AI
- 2. Part I: Lecture on Computer Vision**
  - a. Main problems in Computer Vision
  - b. What exactly is a neural network? (CNNs / Transformers)
  - c. What exactly is an image embedding?
  - d. Computer vision on the Edge
3. Hands-On I: Air-gapped Face recognition on the Pi
- 4. Part II: Lecture on Natural Language Processing**
  - a. A bit of history & development of modern LLMs
  - b. How does an LLM work? Tokenizers, pretraining, post-training
  - c. Context Engineering: Tool calling, RAG
  - d. Libraries: tokenizers-rs, llama.cpp
- 5. Hands-On II: Chat with a LLM on Pi**

