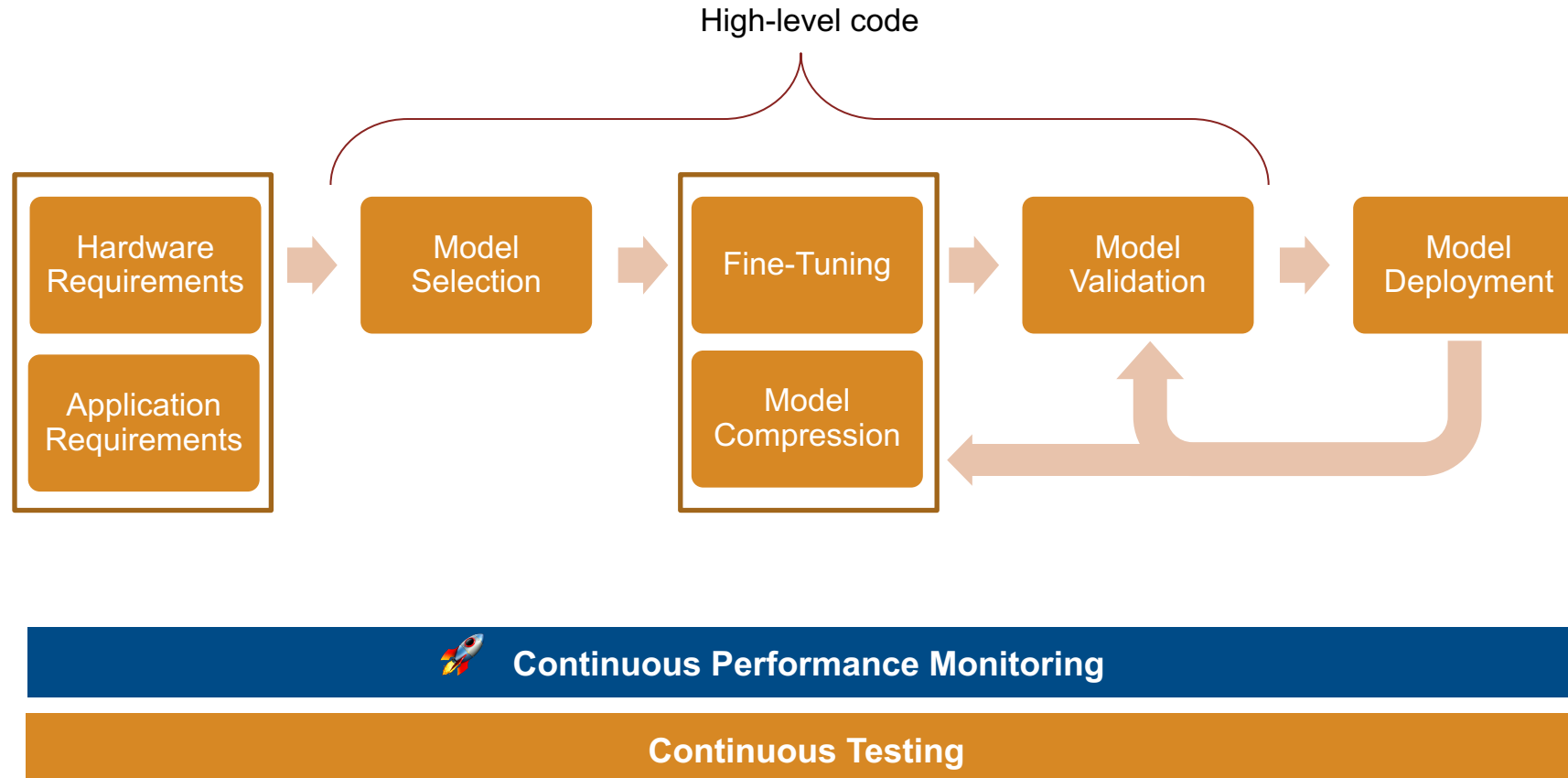# Accelerating Model Optimization on the Edge Through Automated Performance Benchmarking and End-to-End Profiling

Nayara Aguiar, PhD
Performance Engineer
MathWorks

# Deploying AI Models: A Tale on Performance

High-level code

| Hardware Requirements | Model Selection | Fine-Tuning | Model Validation | Model Deployment |
| Application Requirements | | Model Compression | | |

🚀 **Continuous Performance Monitoring**

**Continuous Testing**

*When does performance become a concern in your AI deployment pipeline?*

# Evaluating performance throughout the development process enables early detection of bottlenecks

Tracking performance makes it easier to pinpoint the source of issues
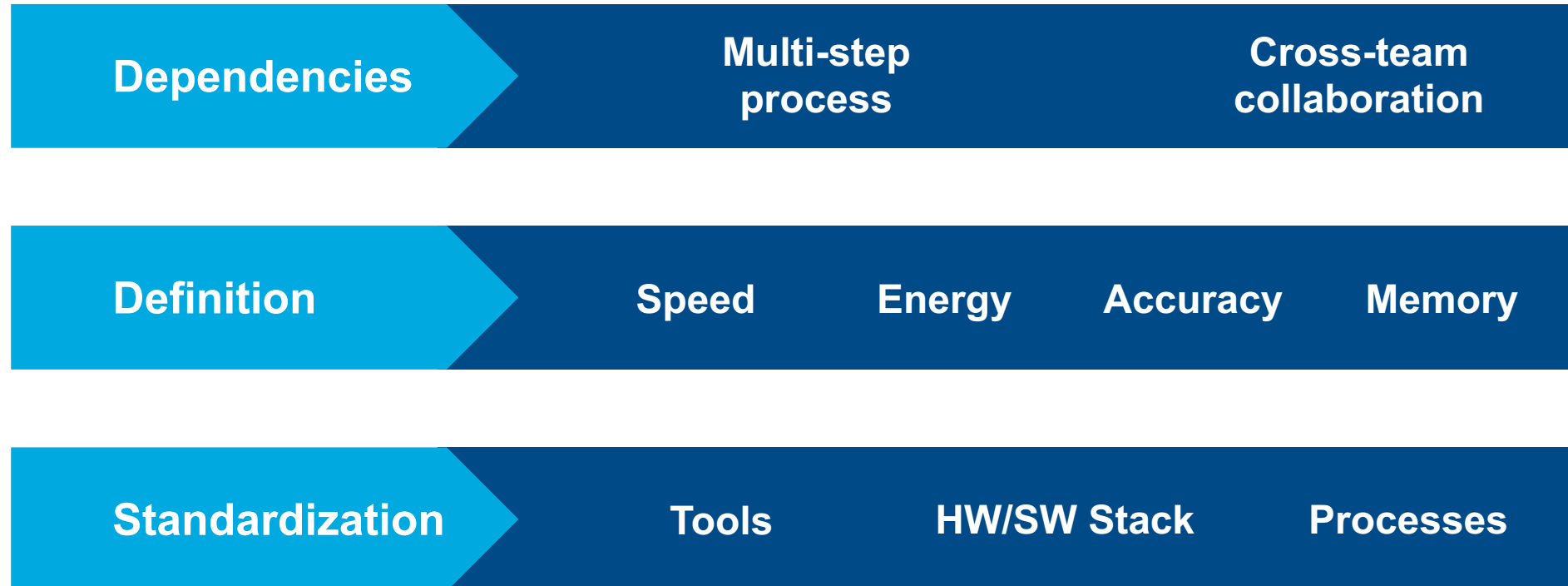
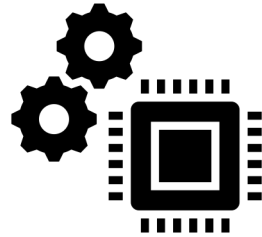Mitigation of performance issues is less costly with early detection

Quality targets for final product can be met while enhancing development process
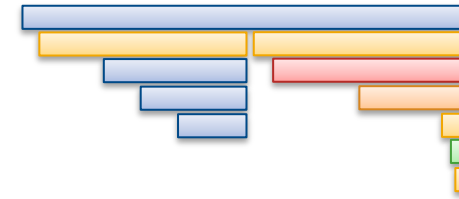
# Optimizing performance of the AI deployment pipeline presents multiple challenges

| Dependencies | Multi-step process | | Cross-team collaboration |
|---|---|---|---|

| Definition | Speed | Energy | Accuracy | Memory |
|---|---|---|---|---|

| Standardization | Tools | HW/SW Stack | Processes |
|---|---|---|---|

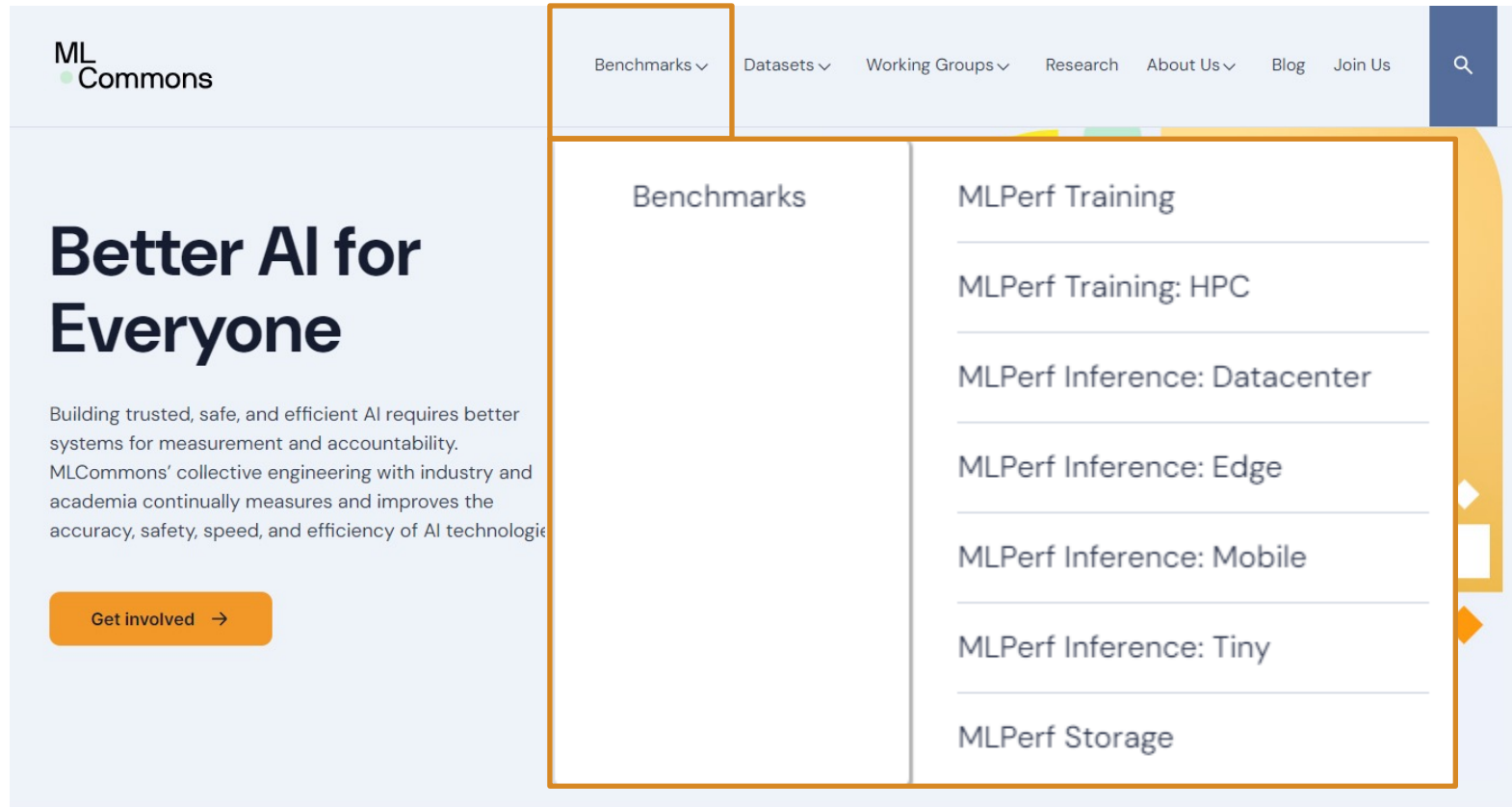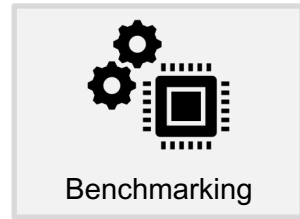# We will focus on our ongoing efforts to enhance the process for performance benchmarking and profiling

Automating Performance
Benchmarking

Enabling End-to-End
Performance Profiling

# Integrating MLPerf™ to our internal performance benchmarking strategy helped us standardize tooling

# We used automation to enhance the benchmarking process



Benchmarking

**Network**
resnet50

**BatchSize**
128

**TargetConfiguration**
TargetLanguage: C++
Hardware: RaspberryPi

**BenchmarkMode**
PerformanceOnly
SingleStream

**Prepare High-Level Code** → **Generate C++ Code** → **Deploy Code to Device** → **Invoke MLPerf™ Benchmark**

**Results**
MLPerf™ Logs

# This benchmark can also be wrapped in our testing infrastructure for further automation

Benchmarking

**Testing Framework**
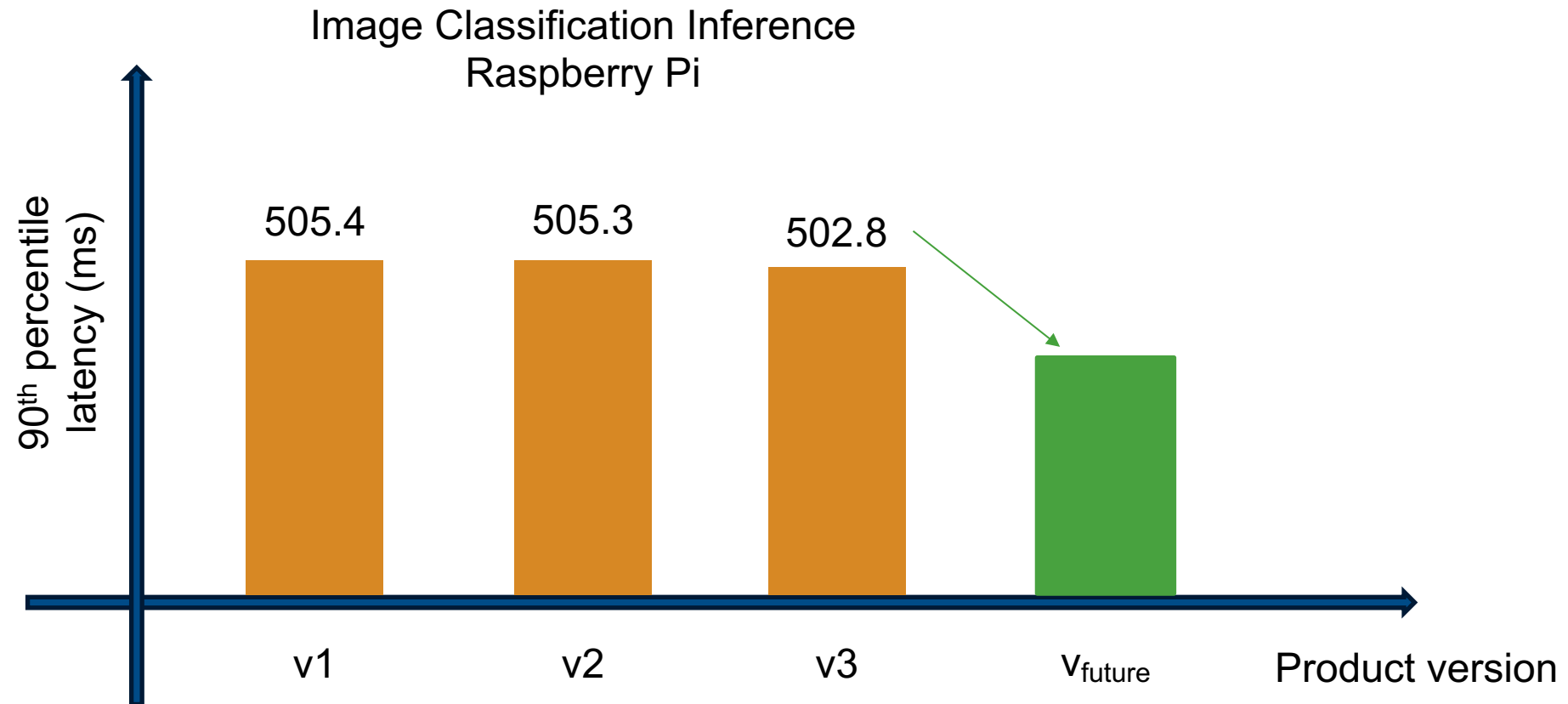


Results

Network = {networkA, networkB}
BatchSize = {64, 128}
TargetConfiguration = {configA, configB}
BenchmarkMode = {mode1, mode2}

1. Select pre-trained models:
   - mobilenetv2
   - resnet50
2. Create configurations for Raspberry Pi runs:
   - Using original network (FP32)
   - Using quantized network (INT8)
   - Using network equalized before quantization (INT8)
3. Batch size:
   - 1024
4. Select benchmark mode:
   - AccuracyOnly, SingleStream

|  | Original (FP32) | Quantized (INT8) | Equalized (INT8) |
|---|---|---|---|
| mobilenetv2 | 70.3% | 0.2% | 60.3% |
| resnet50 | 72.2% | 69.3% | 68.9% |

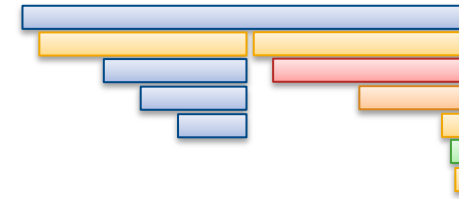# Automated benchmarks allow for performance monitoring
## … and profiling helps investigating performance bottlenecks

Image Classification Inference
Raspberry Pi



90th percentile latency (ms)

505.4          505.3          502.8

v1          v2          v3          $v_{future}$          Product version

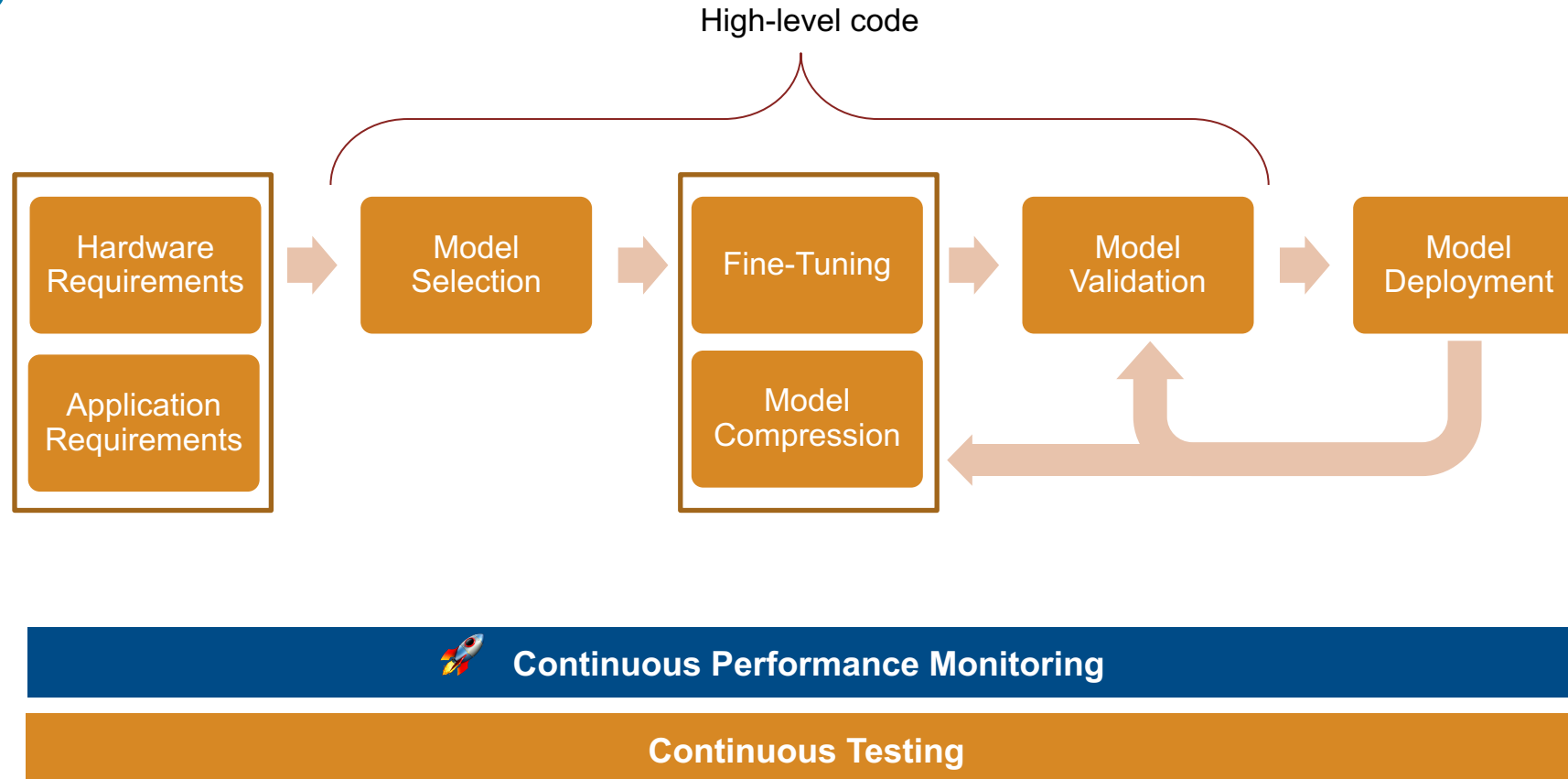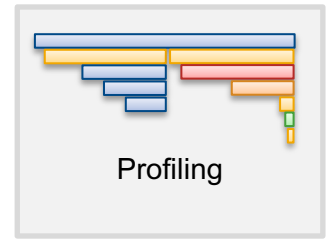# We will focus on our ongoing efforts to enhance the process for performance benchmarking and profiling

Automating Performance
Benchmarking

Enabling End-to-End
Performance Profiling

# We develop tools that cover the entire AI deployment pipeline

Profiling

High-level code

| Hardware Requirements | Model Selection | Fine-Tuning | Model Validation | Model Deployment |
| Application Requirements | | Model Compression | | |

🚀 **Continuous Performance Monitoring**

**Continuous Testing**

*To facilitate end-to-end performance investigations, we developed the Unified Timeline*

Profiling

# The *Unified Timeline* facilitates performance analysis across the AI deployment pipeline

C++
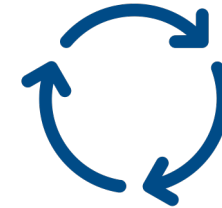MATLAB
JavaScript
…

Generic
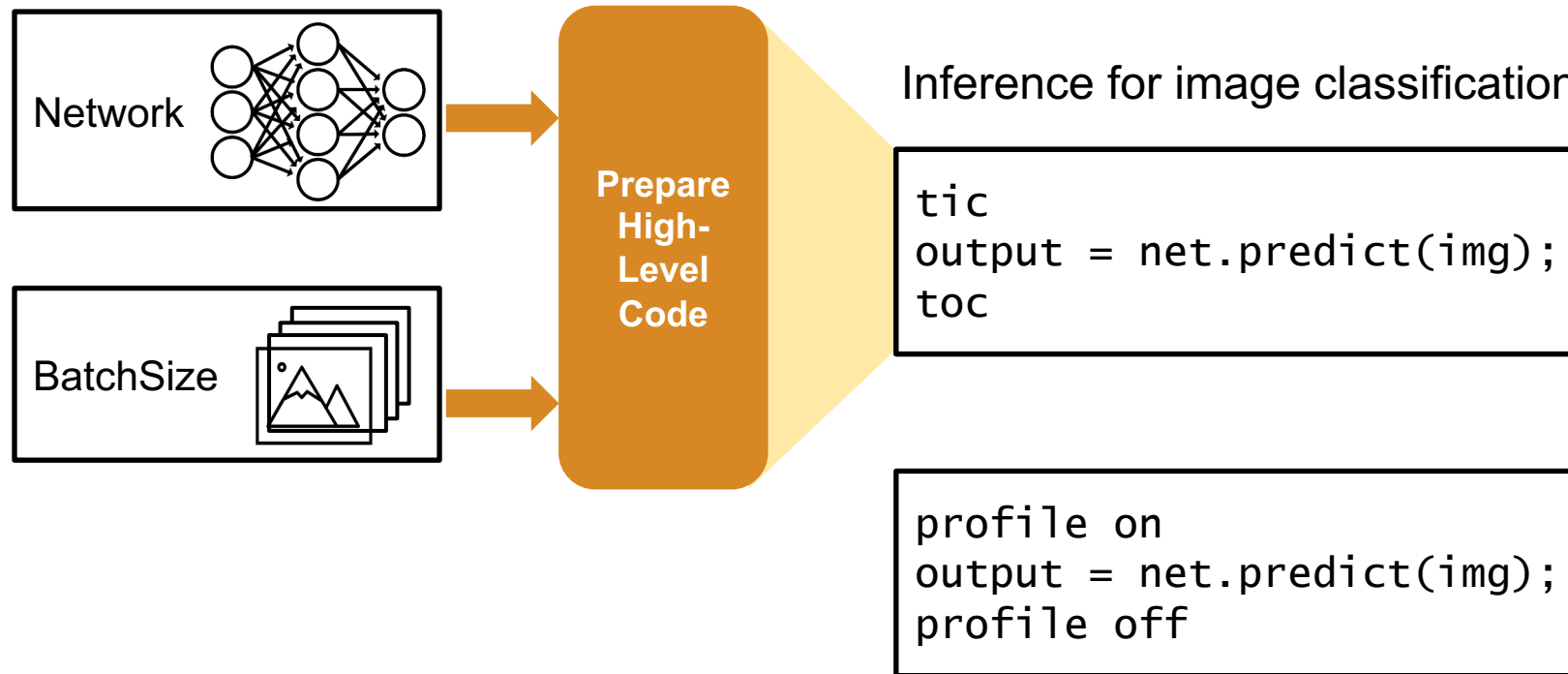
Easy to use

Continuously enhanced

# We can use our in-house profiling tool to investigate performance bottlenecks

Network

BatchSize

**Prepare High-Level Code**

Inference for image classification

```
tic
output = net.predict(img);
toc
```

```
profile on
output = net.predict(img);
profile off
```

Total time

~ 700 ms

Time breakdown

# We can visualize the timeline for this code in a browser and zoom in stacks of interest

Profiling

Inference for image classification



Can this be optimized?

# Make performance a core aspect for AI applications on the edge



- Automate performance benchmarking
- Explore performance profiling tools

**… accelerate model optimization on the edge**

# Copyright Notice

This presentation in this publication was presented at the tinyML® Summit 2024. The content reflects the opinion of the author(s) and their respective companies. The inclusion of presentations in this publication does not constitute an endorsement by tinyML Foundation or the sponsors.

There is no copyright protection claimed by this publication. However, each presentation is the work of the authors and their respective companies and may contain copyrighted material. As such, it is strongly encouraged that any use reflect proper acknowledgement to the appropriate source. Any questions regarding the use of any materials presented should be directed to the author(s) or their companies.

tinyML is a registered trademark of the tinyML Foundation.

## www.tinyml.org