



# ADC<sup>24</sup> Bristol

## REAL-TIME INFERENCE OF NEURAL NETWORKS

*A PRACTICAL APPROACH FOR DSP ENGINEERS — PART II*

**FARES SCHULZ  
& VALENTIN ACKVA**



**Audio Developer Conference**  
Bristol 2024



# **Real-Time Inference of Neural Networks**

**A Practical Approach for DSP Engineers – Part II**

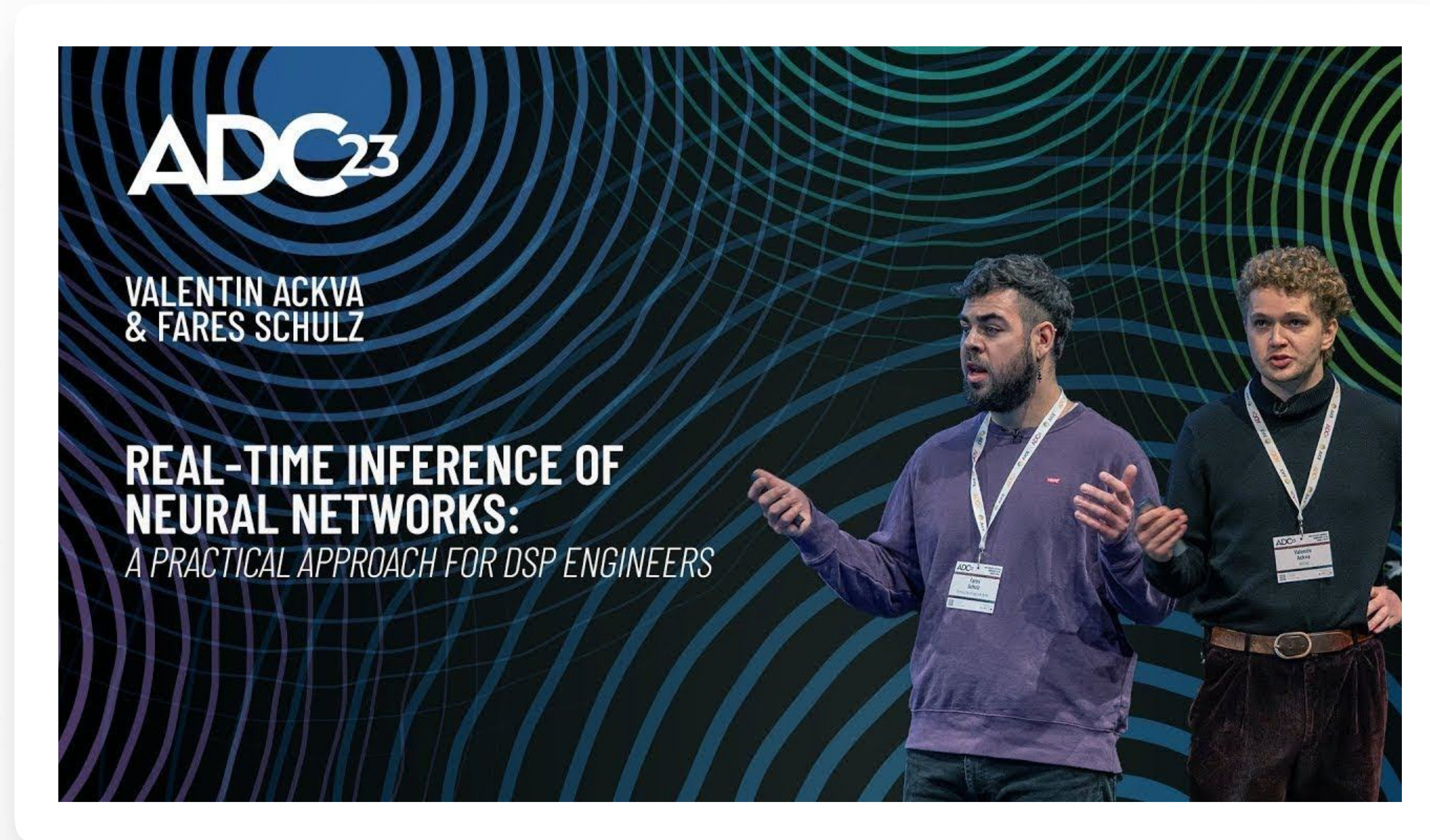
Fares Schulz

Valentin Ackva

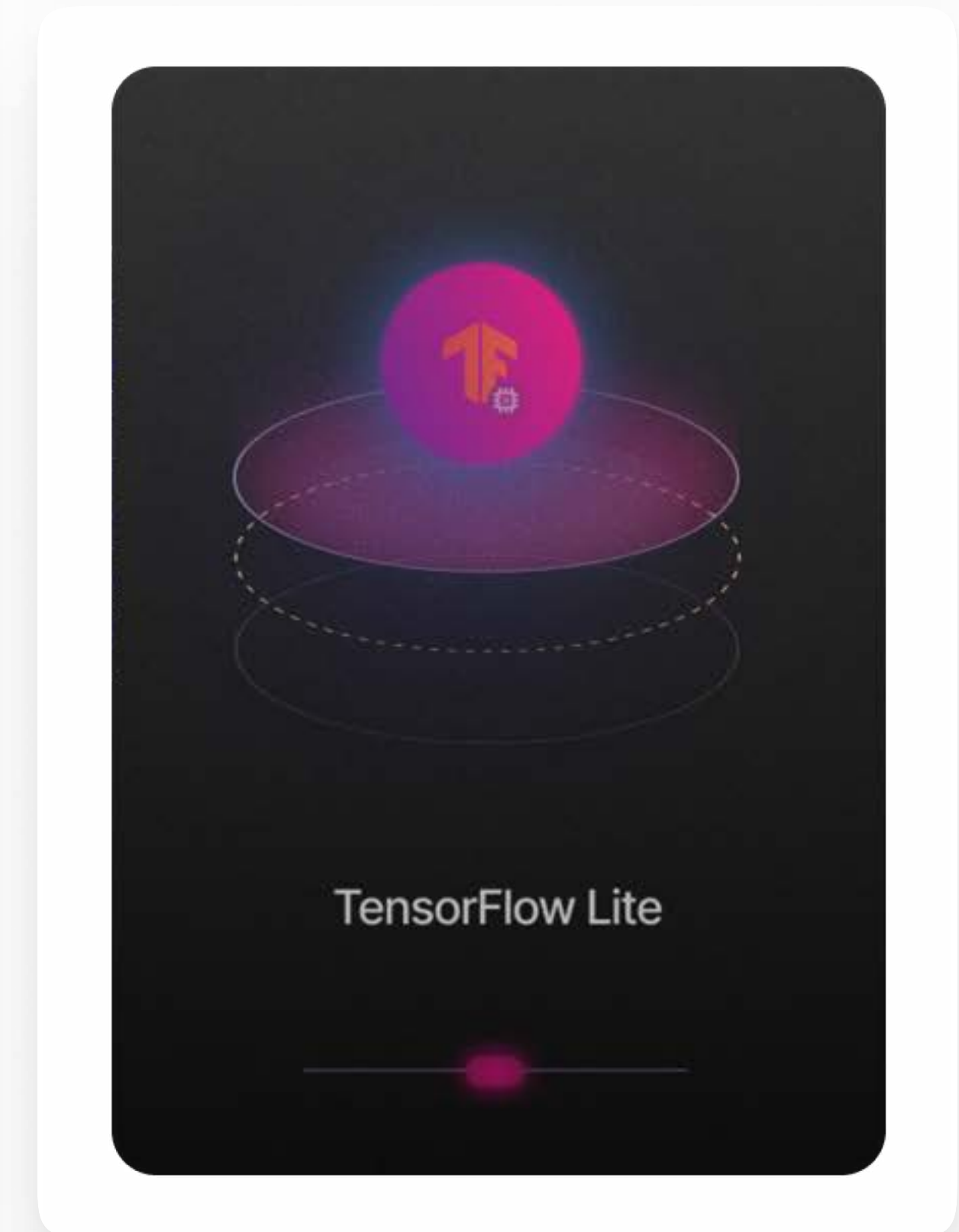




# Background



London, ADC 2023



nn-inference-template



# About us



**Fares Schulz**

**Researcher / Head of TU Studio**  
Technische Universität Berlin, Germany



**Valentin Ackva**

**Audio Software Developer**  
INSONE GmbH in Leipzig, Germany



# Table of Content

1. Introduction
2. Examine Real-Time Violations
3. Library Architecture
4. Deep Dive Thread Pool and Latency
5. Impact on Inference Runtimes
6. Conclusion



Chapter I

# Introduction

Recap Part I, Relevance









# How to Implement the Inference?

**Write inference  
yourself**

using  
std::lib, Eigen, SIMD

**Use specialized  
libraries**

known as  
inference engines



# How to Implement the Inference?

**Write inference yourself**

using  
std::lib, Eigen, SIMD

**Use specialized libraries**

known as  
inference engines





# Real-time Principles

- These inference engines favor average execution times
- None of the libraries gives real-time safety guaranties
- Confusion on real-time safety of major inference engines
  - Chowdhury finds no real-time safe (2021)
  - Stefani et al. conclude real-time safety after first inference (2022)
- Noteworthy: RTNeural – a real-time safe inference engine
  - Fast for small networks
  - Very limited layer support

Chowdhury, J. (2021). Rtnatural: Fast neural inferencing for real-time systems. *arXiv preprint arXiv:2106.03037*.

Stefani, D., Peroni, S., & Turchet, L. (2022). A comparison of deep learning inference engines for embedded real-time audio classification. In *Proceedings of the International Conference on Digital Audio Effects, DAFx* (Vol. 3, pp. 256-263).





# In the Last Talk

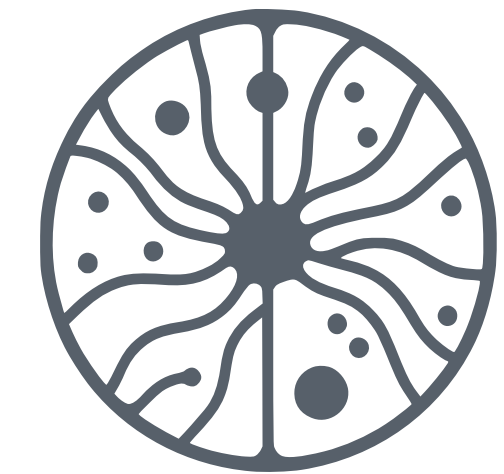
- Outline of the pipeline for implementing neural networks in audio plug-ins
- Overview of major inference engines
- Presentation of a real-time safe plug-in template
- Basic benchmark for inference engines inferring a neural network model
- Discussion about continuous signals / streamability of neural networks



Part I - ADC23

# In This Talk

- Quantification of real-time violations of inference engines
- Cross-platform library - ANIRA: An Architecture for Neural Network Inference in Real-Time Audio Applications
  - Streamlines the use of neural networks in any real-time audio environment
  - Significant improvements for the use of multiple instances
  - Refined latency calculation
  - Built-in benchmarking capabilities
- Performance impact of various factors on inference runtimes



**anira**



Chapter II

# **Examine Real-Time Violations**

## **Methods, Results**



# Real-Time Violation Checks

How can we check real-time violations by external libraries?

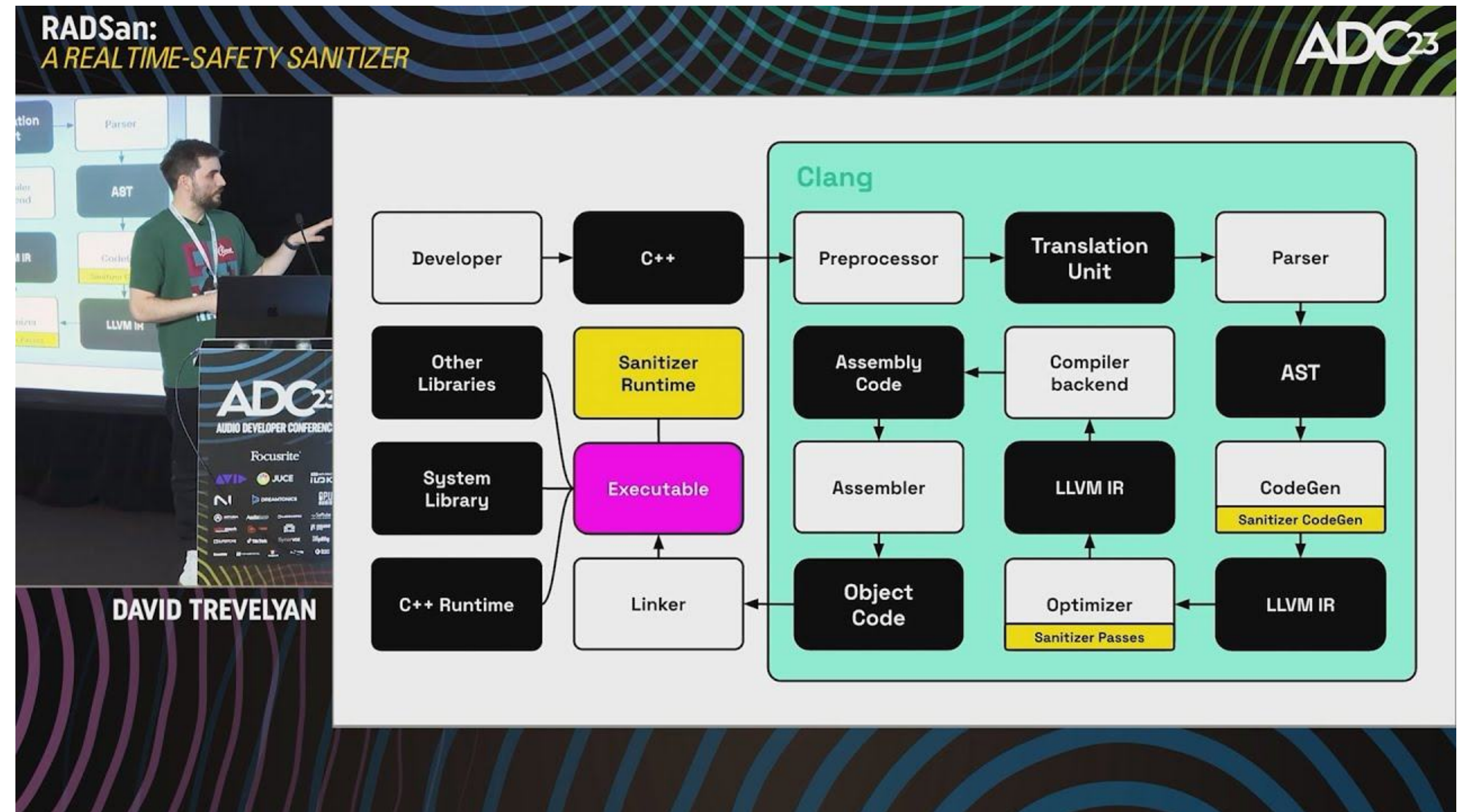




# Verification Method

Real Time Sanitizer - RTSan (prev. RADSan):

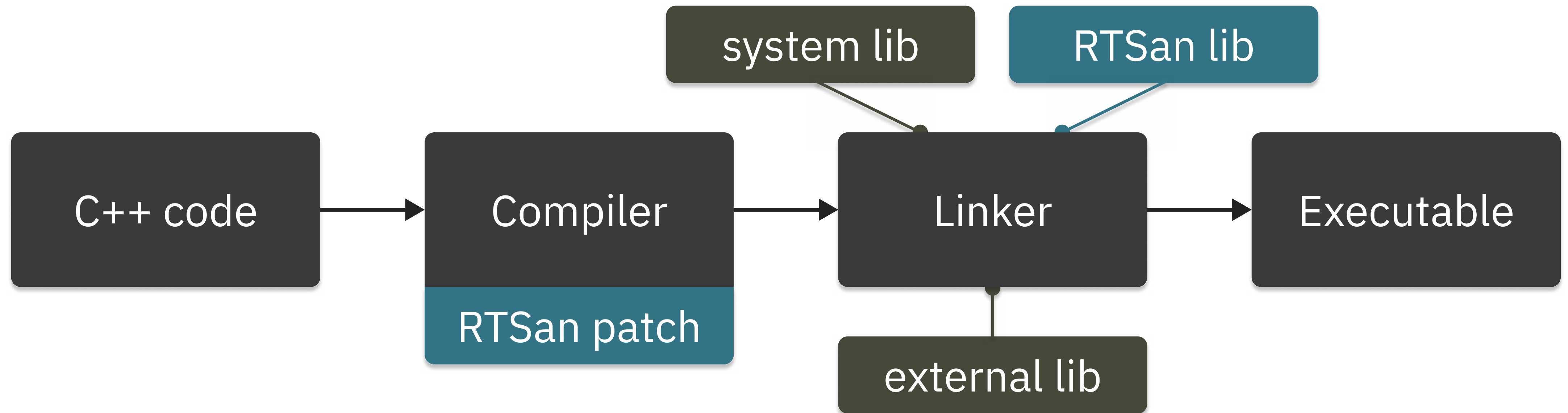
- Part of clang compiler and runtime library



# Verification Method

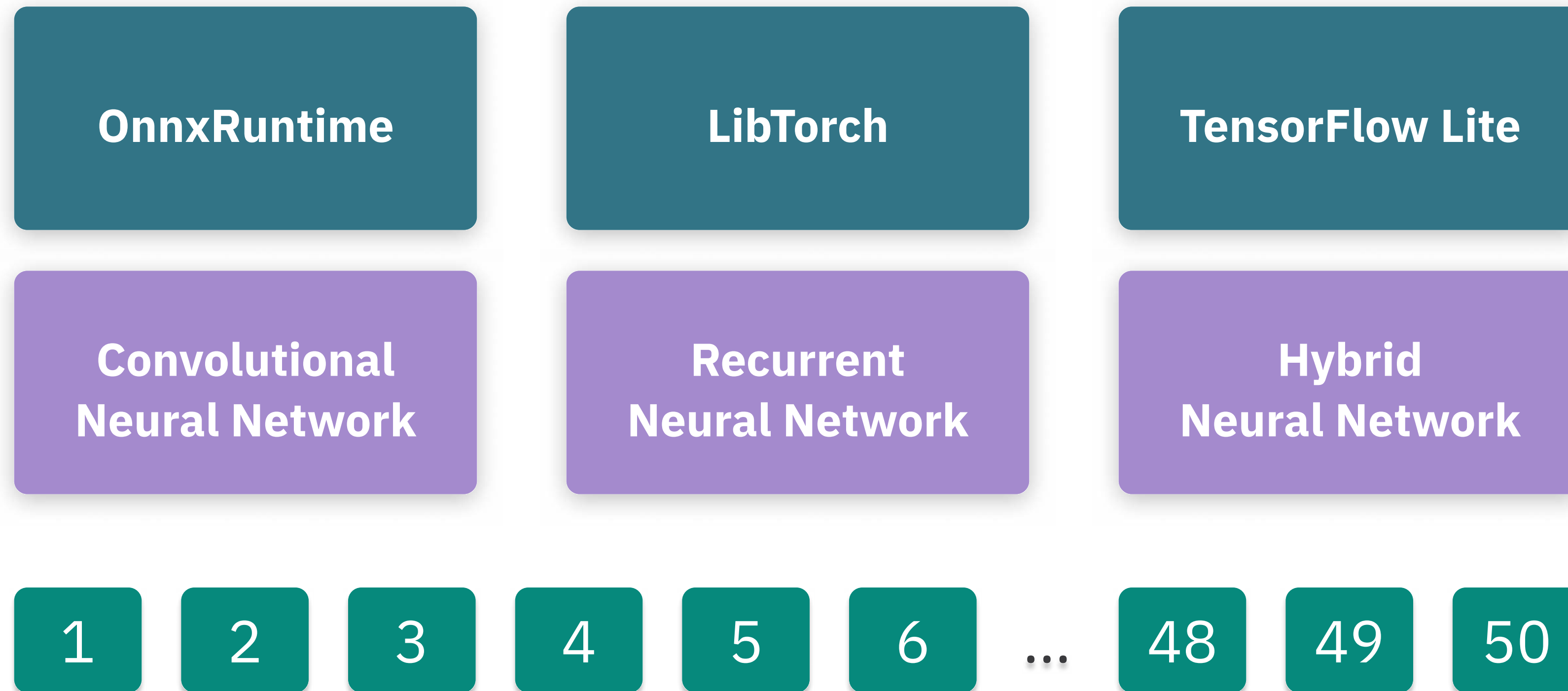
## How does RTSan work?

- Introduces real-time context and intercepts for non-real-time safe operations like `malloc` and `free`

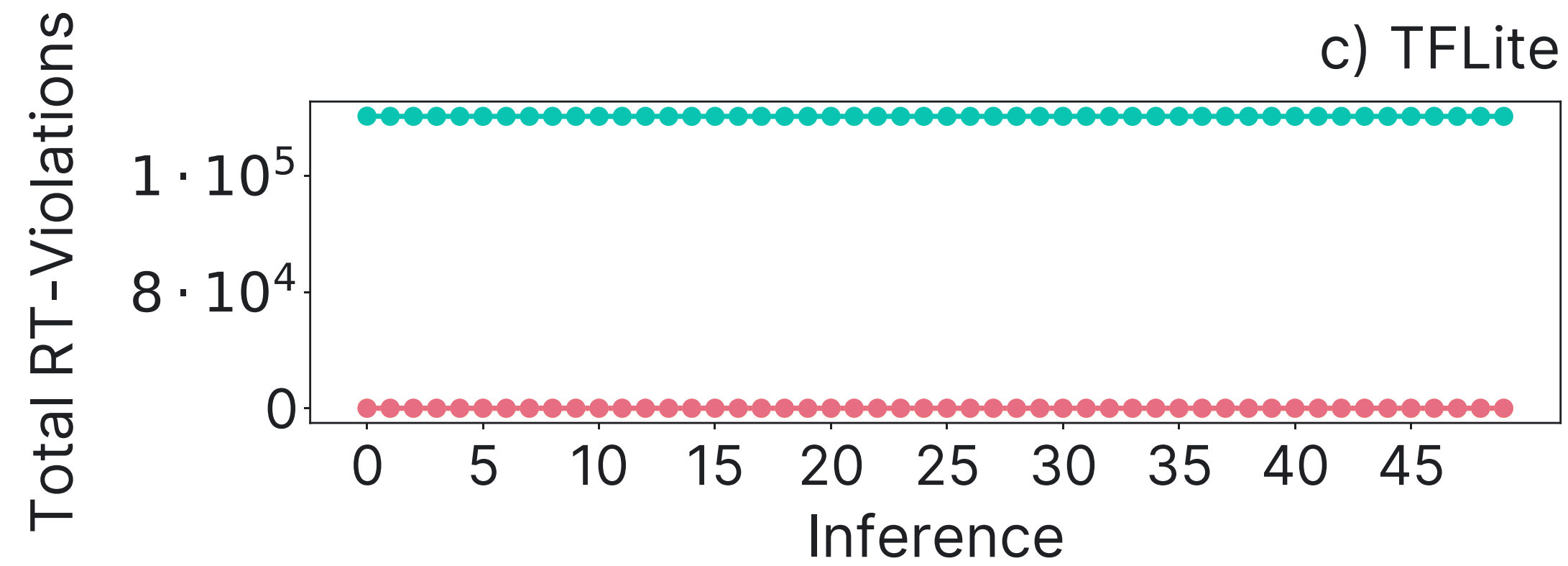
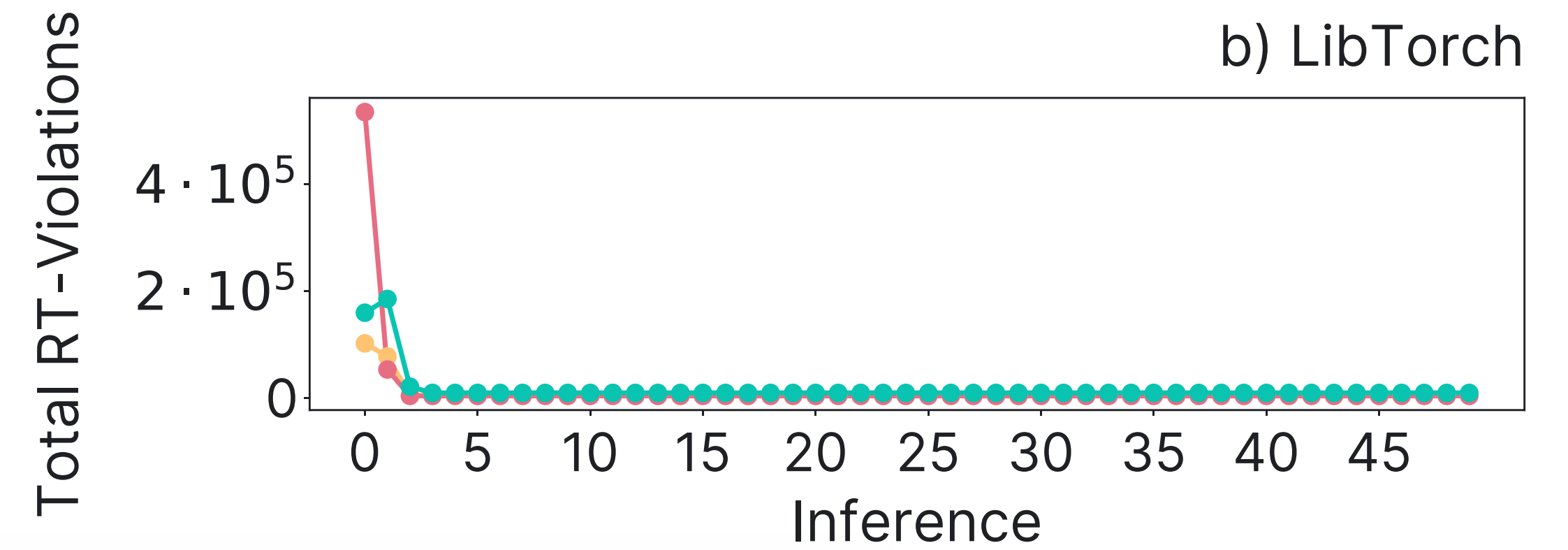
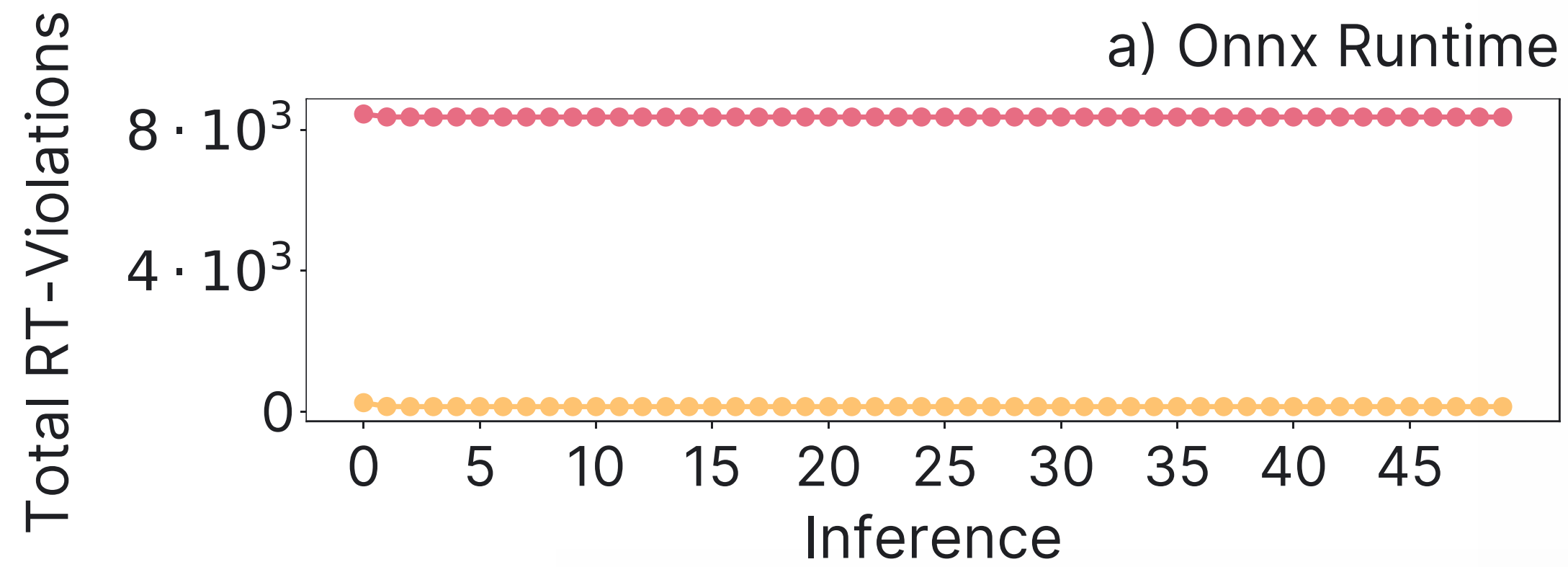




# Test Setup Overview



# Results



# Inference Engine Integration

**We should not run these engines on the real-time thread!**





Chapter III

# Library Architecture

Solving Limitations, Interface

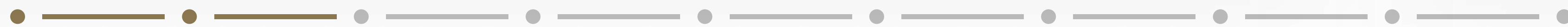
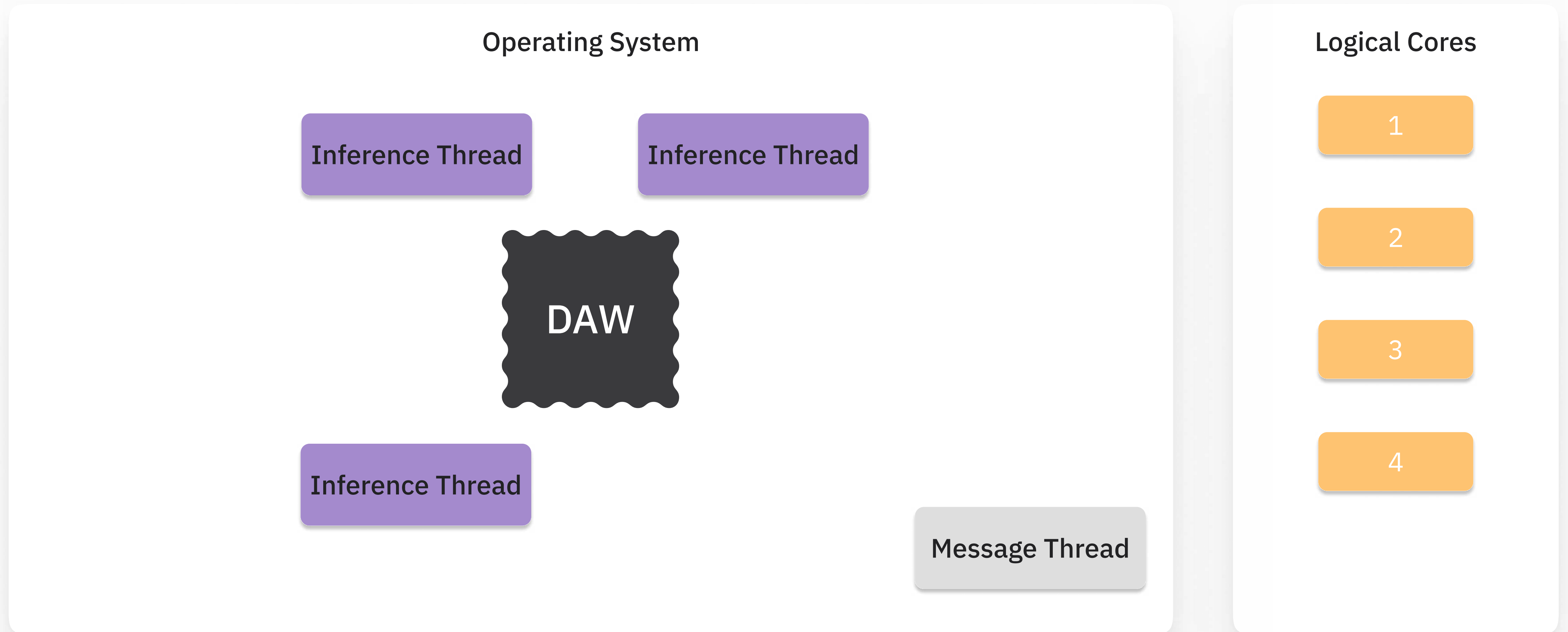








# Architecture Limitation





# Architecture Limitation

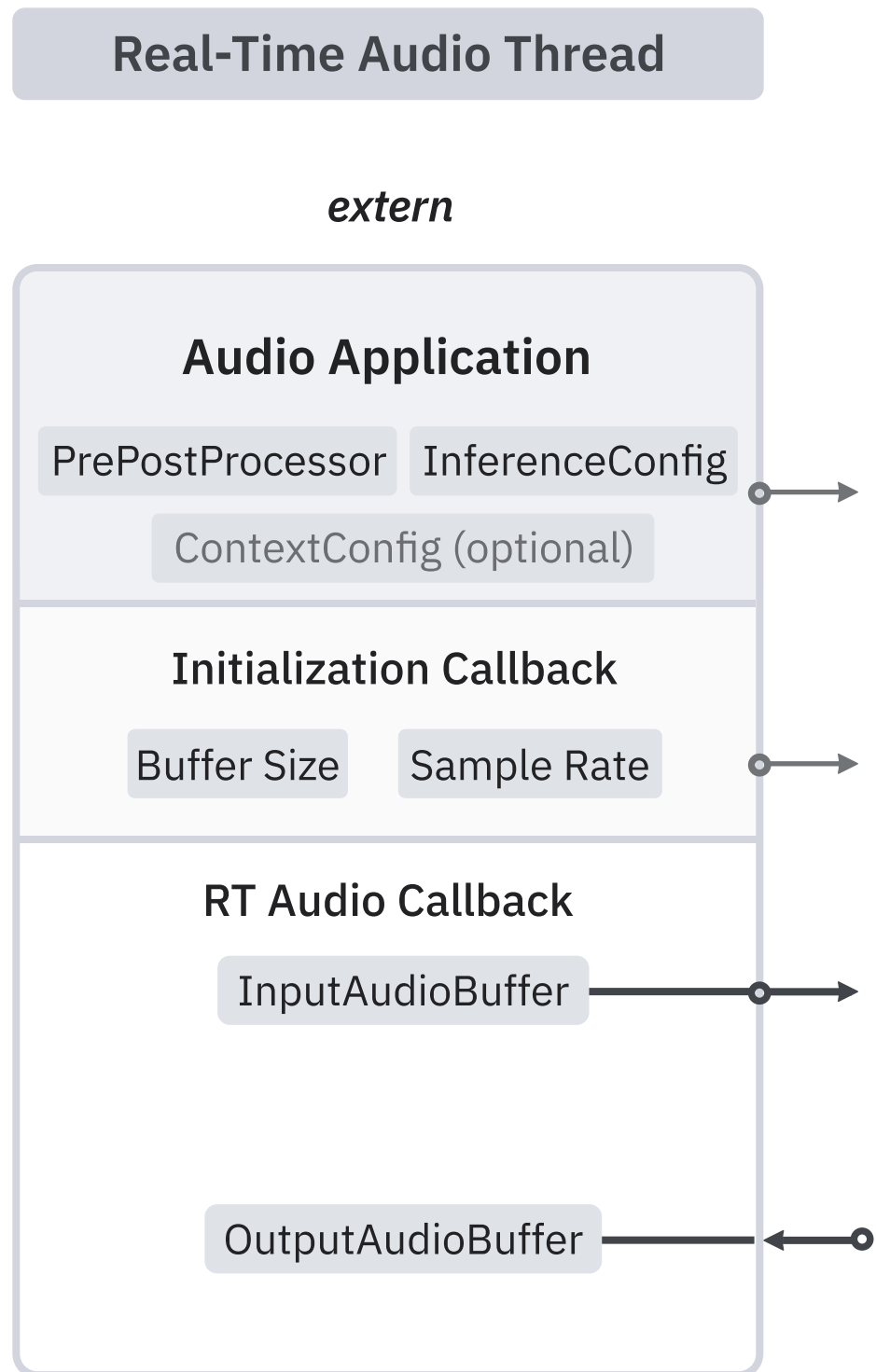
## Oversubscription Problem:

- If number of high-priority threads  $>$  number of logical cores
- Causes threads to compete for the same cores
- Can be especially problematic for real-time processes

## Solution - Static Thread Pool Design:

- Shares inference threads across instances (e.g. different plugins)
- Can enable parallel execution for faster inference
- Implementation requires developing a **dedicated library**

# Architecture

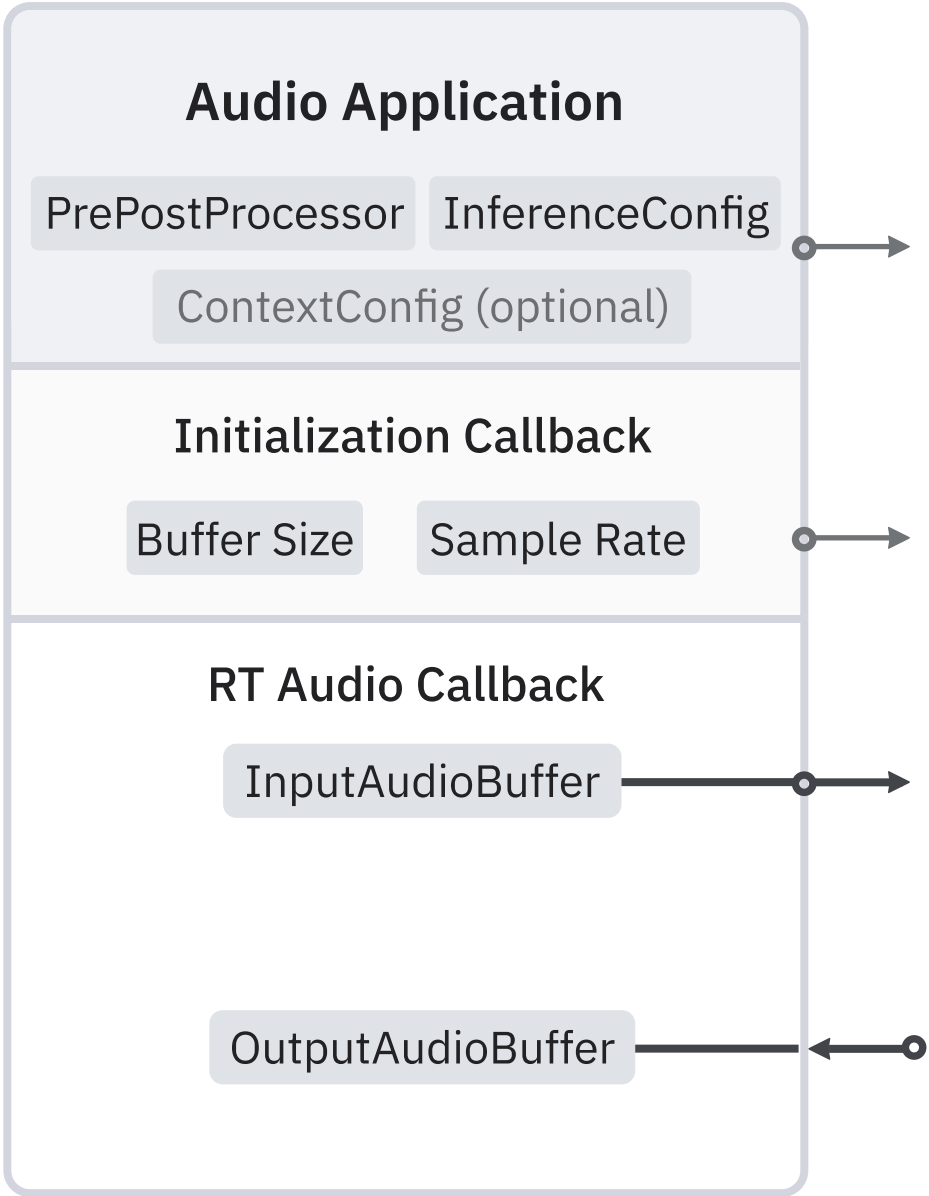




# Architecture

## Real-Time Audio Thread

*extern*



## Inference Threads

*private anira*

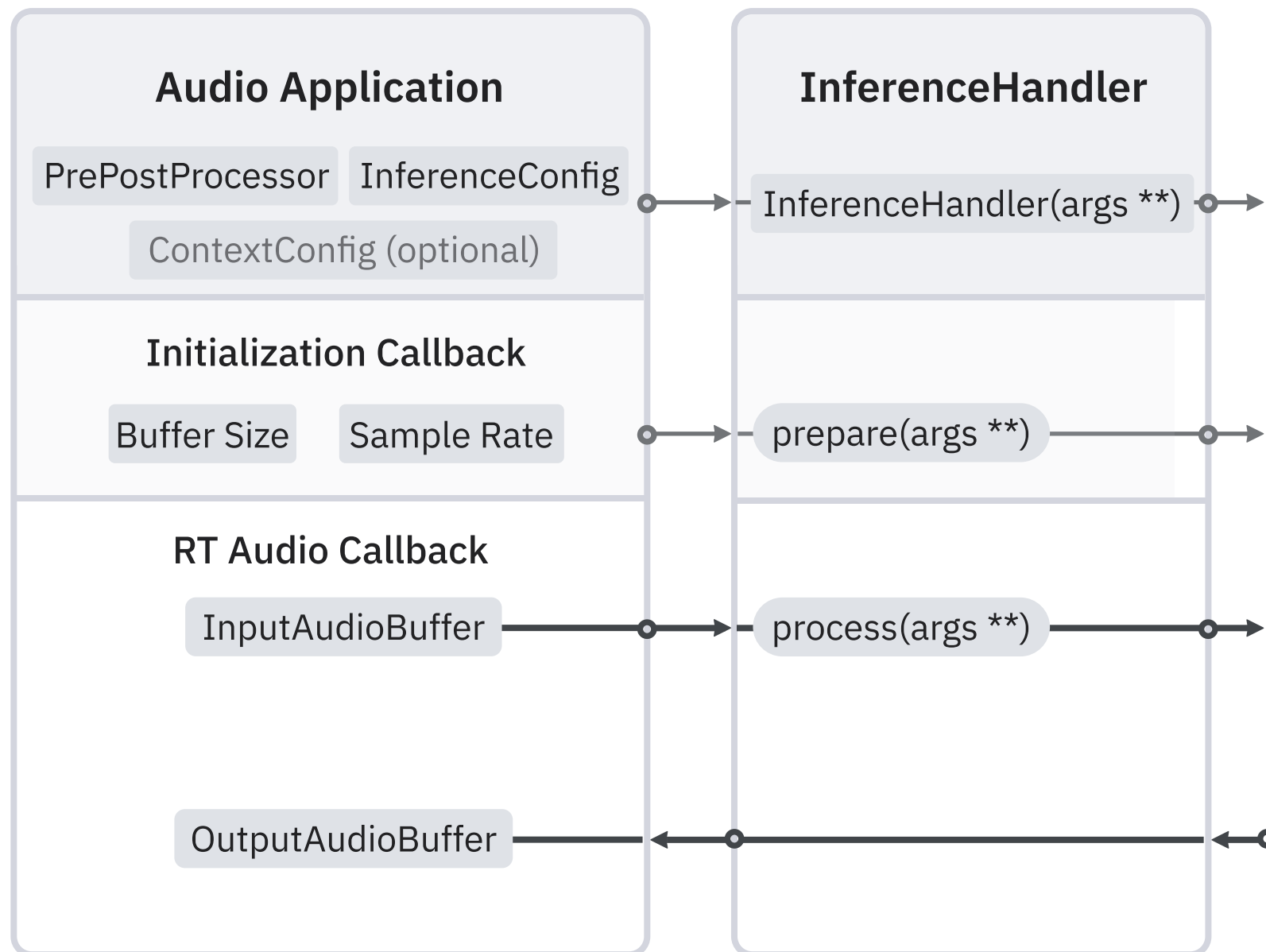


# Architecture

## Real-Time Audio Thread

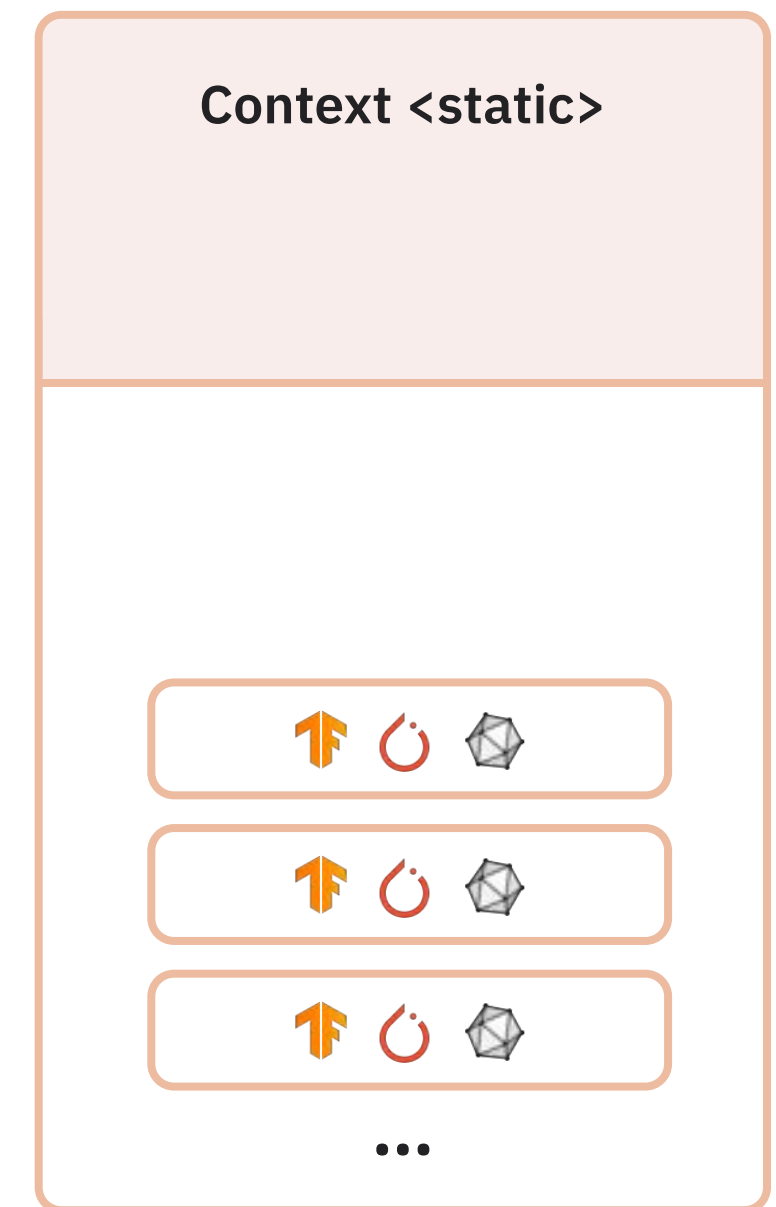
*extern*

*public anira API*

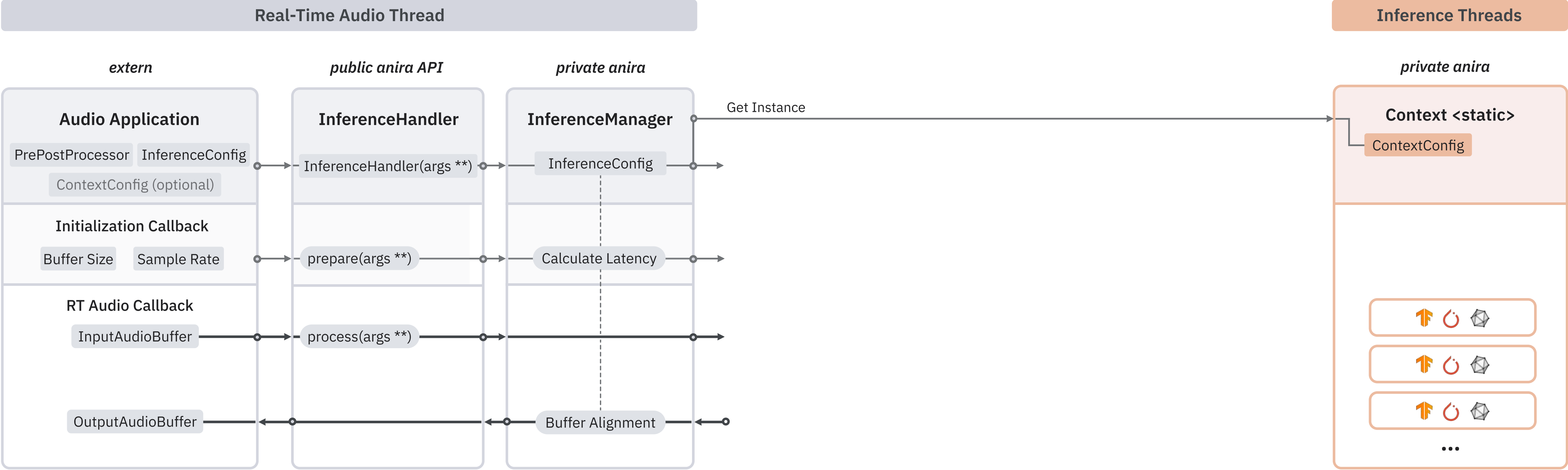


## Inference Threads

*private anira*

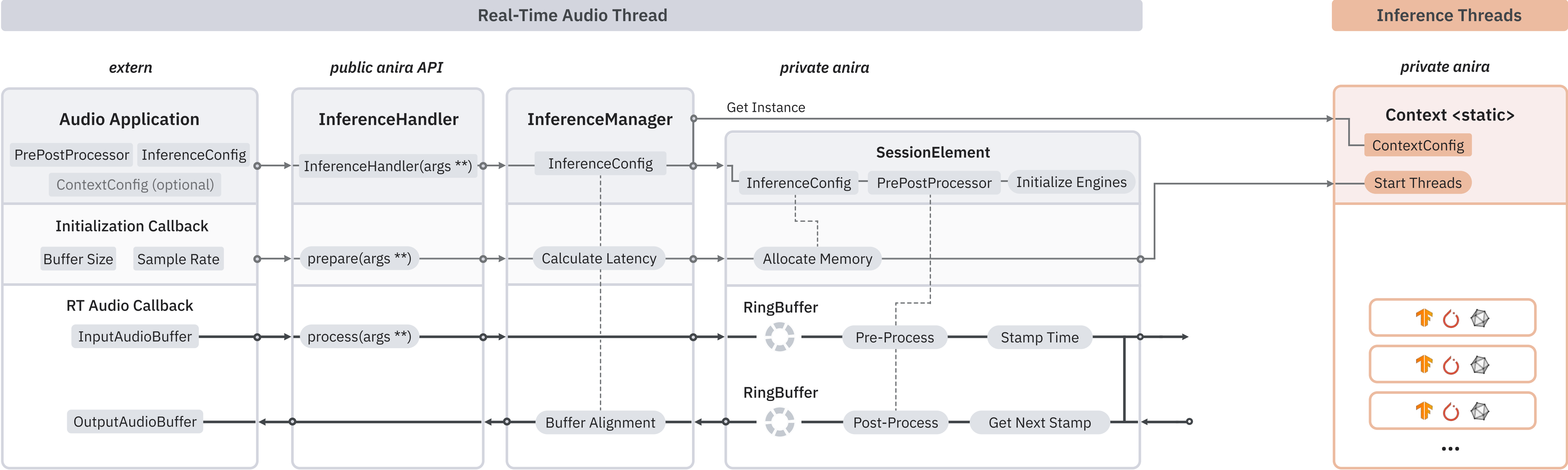


# Architecture

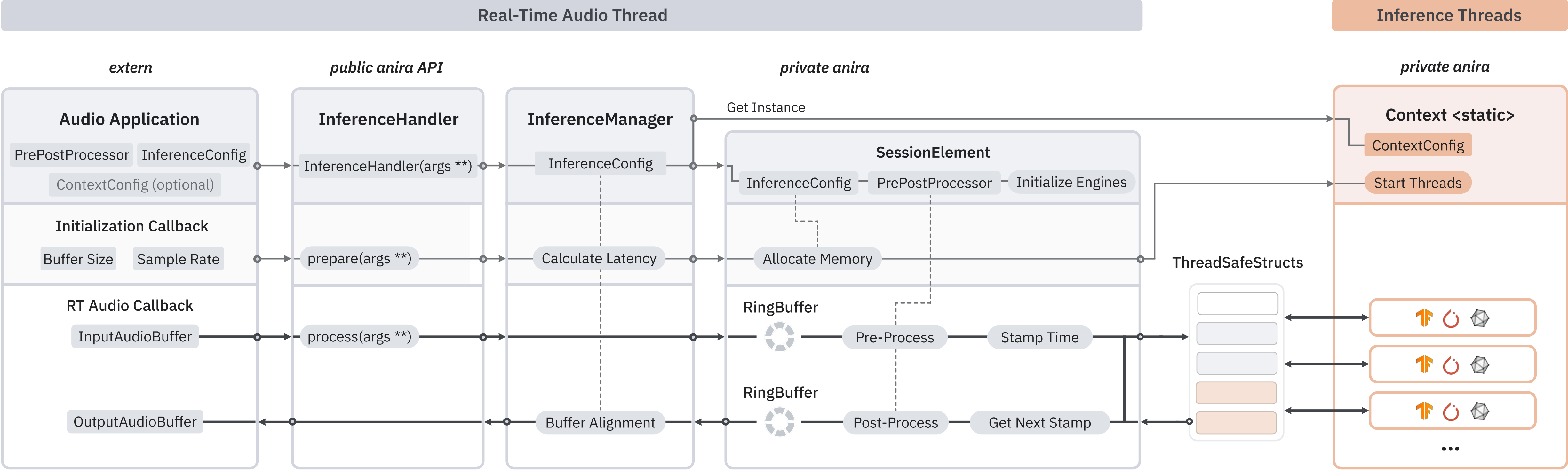




# Architecture



# Architecture



# Library Integration

How does the interaction with the library work?







# Interface

## Necessary Parameters:

- Model data path\*
- Model shapes\*
- Max inference time

\* definable for multiple inference engines

## Context (optional, shared):

- Number of threads

## Optional Parameters:

- Model latency
- Warm-up inference
- Number of channels
- Bind session to processor
- Number of parallel processors
- ...

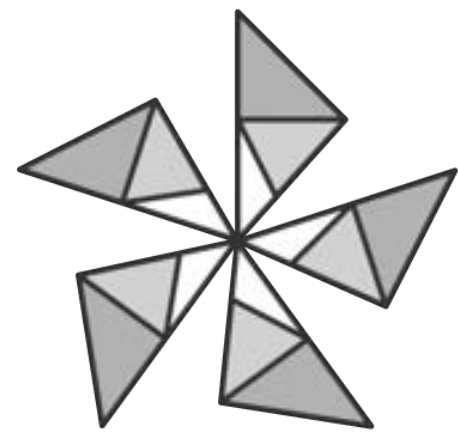
# Code Example

```
1  #include <anira/anira.h>
2
3  anira::InferenceConfig inference_config(
4      {"path/to/model.onnx", anira::InferenceBackend::ONNX} // Model data
5      {{{1, 1, 512}}, {{1, 1, 512}}}, // Input- Output-Tensor shapes
6      5.0f, // Maximum inference time in ms
7  );
8
9  anira::PrePostProcessor pp_processor; // Create default pre- and post-processor
10 anira::InferenceHandler inference_handler(pp_processor, inference_config); // InferenceHandler
11
12 inference_handler.prepare({buffer_size, sample_rate}); // Allocate memory
13 int latency_in_samples = inference_handler.get_latency(); // Get latency of the inference process
14
15 inference_handler.set_inference_backend(anira::InferenceBackend::ONNX); // Select the backend
16
17 process_block(float** audio_data, int num_samples) {
18     inference_handler.process(audio_data, num_samples); // Real-time safe audio processing
19 }
20 }
```



# Supported Inference Engines

**Onnx Runtime**



**LibTorch**



**TensorFlow Lite**



**Custom**

Any Backend /  
Customization





Chapter IV

# Deep Dive Thread Pool and Latency

Static Thread Pool Design, Minimum Latency















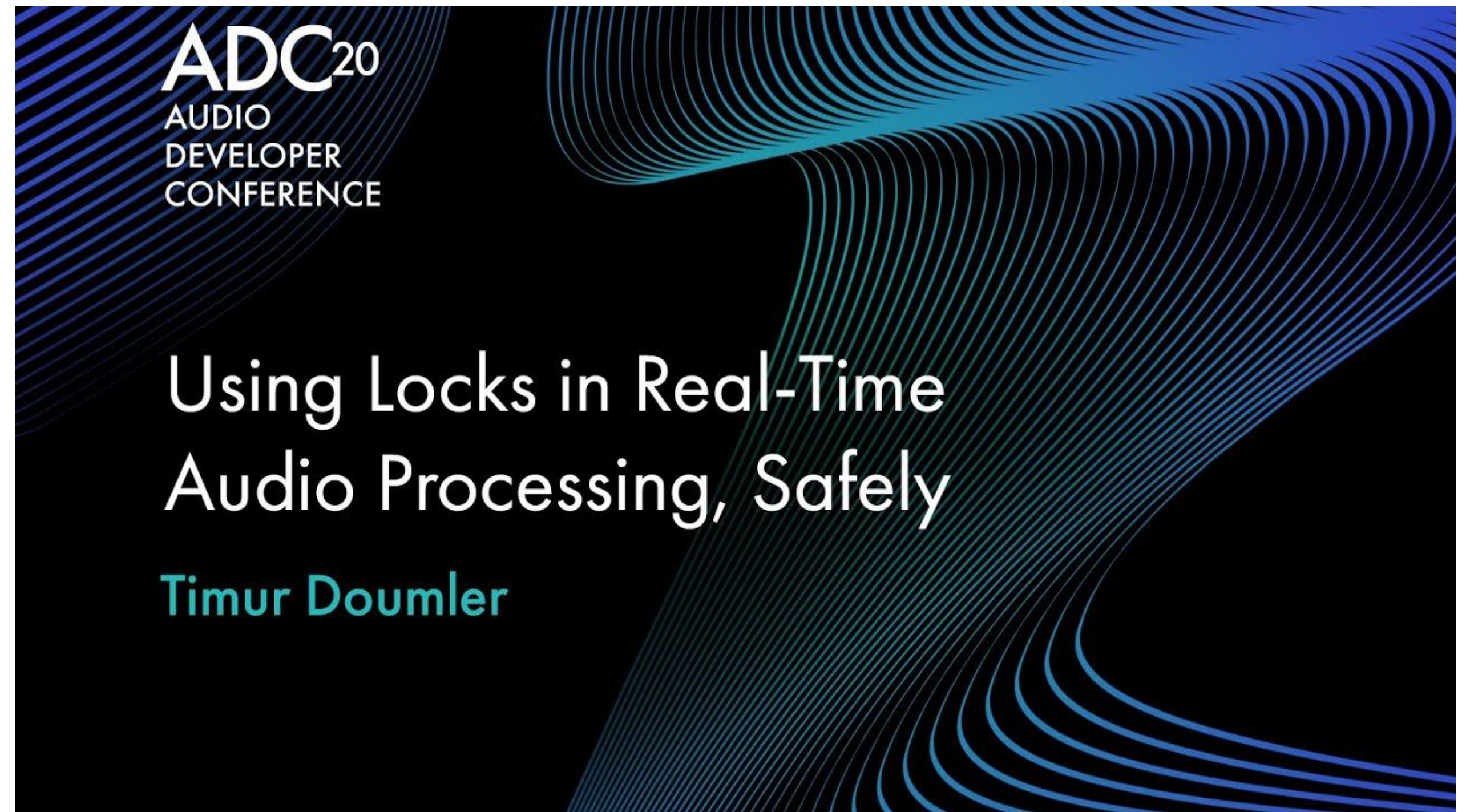




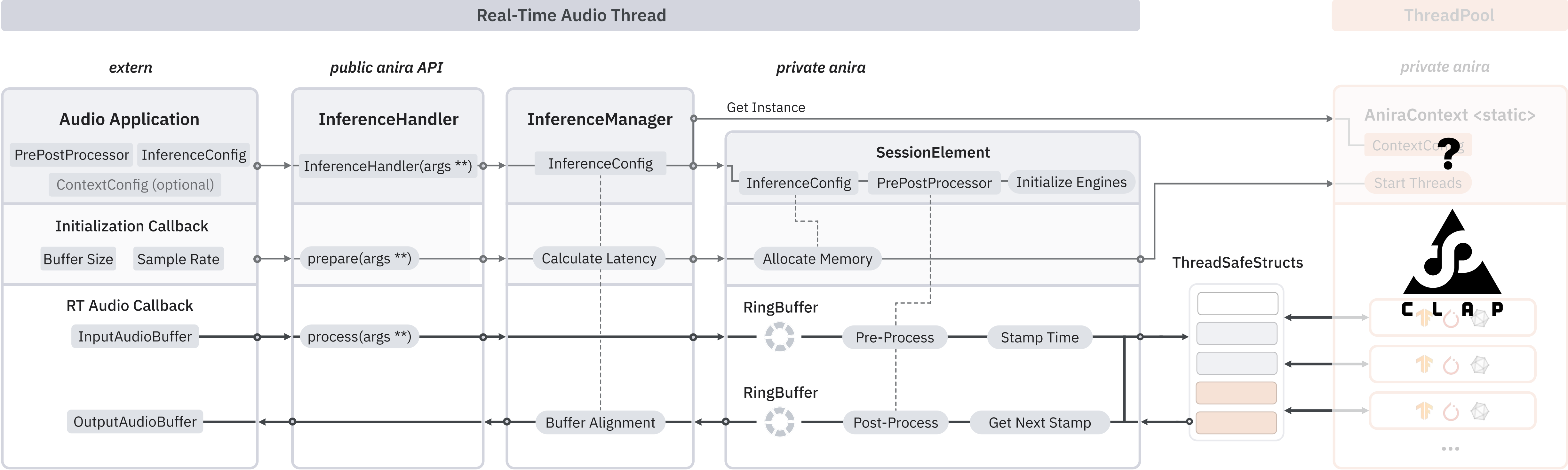


# Safe Spin Locks

- Keeping threads alive without using all CPU resources
- The magic keyword is exponential backoff
- Great talk on this topic:

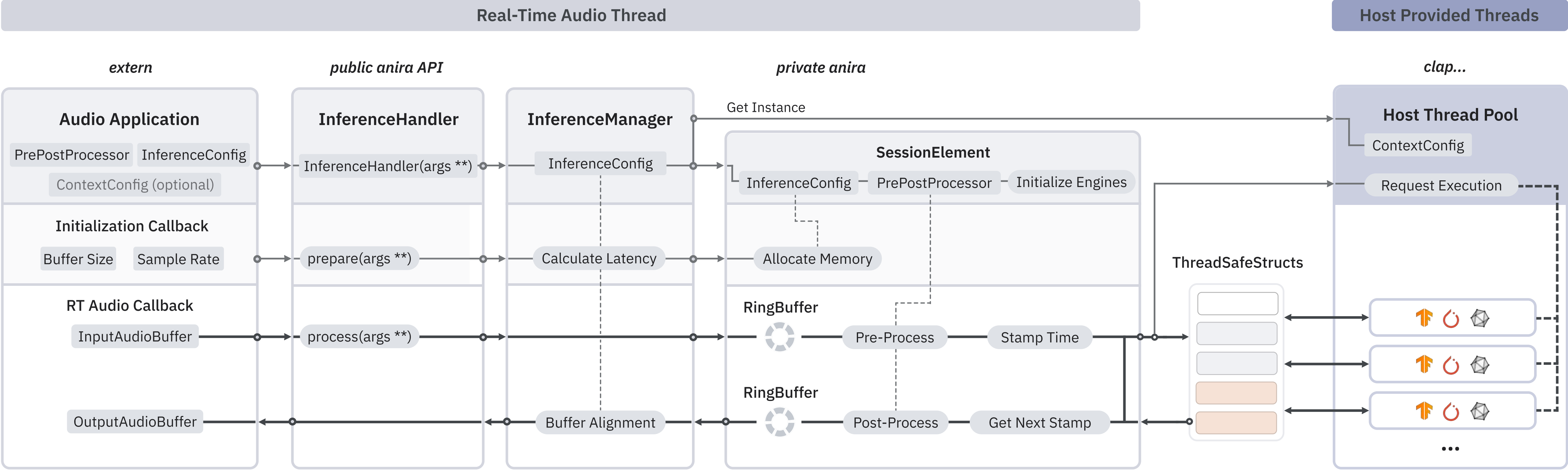


# Host Provided Threads ?





# Host Provided Threads ?



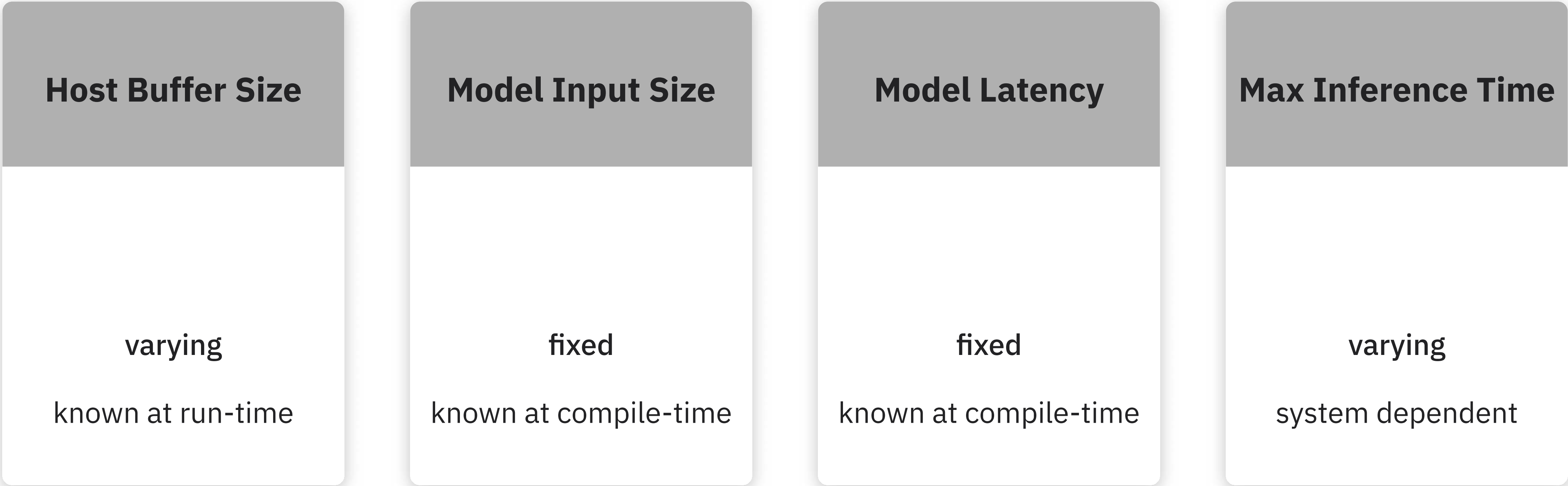


# Host Provided Threads ?

```
1  typedef struct clap_host_thread_pool {
2      // Schedule num_tasks jobs in the host thread pool.
3      // Will block until all the tasks are processed.
4      // This must be used exclusively for realtime processing within the process call.
5      // It can't be called concurrently or from the thread pool.
6      // Returns true if the host did execute all the tasks, false if it rejected the request.
7      // The host should check that the plugin is within the process call, and if not,
8      // reject the exec request.
9      // [audio-thread]
10
11     bool(CLAP_ABI *request_exec)(const clap_host_tQ*host, uint32_t num_tasks);
12 } clap_host_thread_pool_t;
```

bad idea!

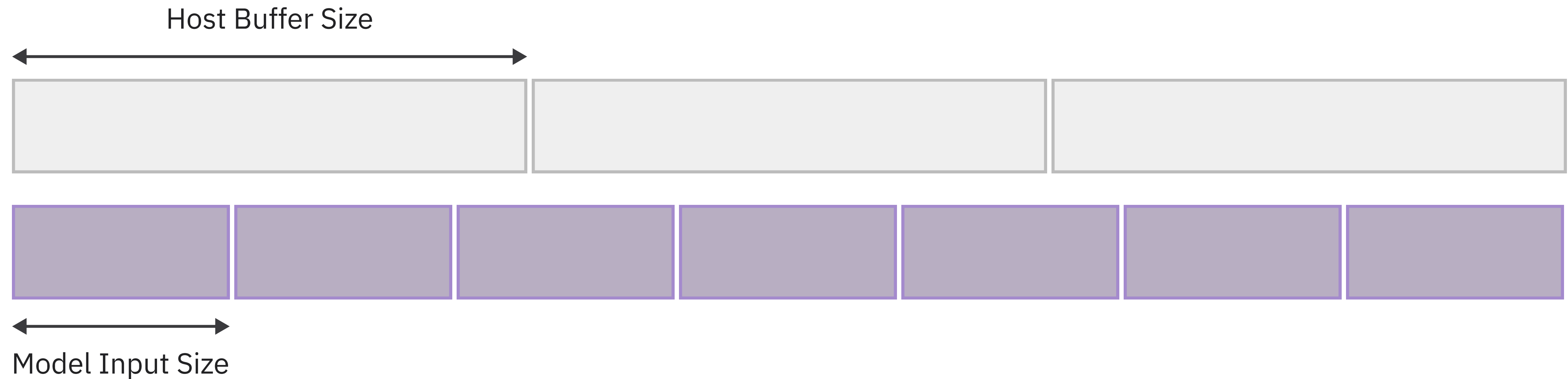
# Latency Influences



# Latency

We often face situations where the host and model input size don't match

→ leaving a remainder when divided









# Latency

$$L_{\text{total}} = H_{\text{adapt}} + \left\lceil \frac{I_{\text{max}}}{H_{\text{host}}} \right\rceil \cdot H_{\text{host}} + M_{\text{int}}$$

$L_{\text{total}}$  Total Latency       $H_{\text{adapt}}$  Largest Remainder

$I_{\text{max}}$  Max. Inference Time       $H_{\text{host}}$  Host Buffer Size       $M_{\text{int}}$  Model Latency





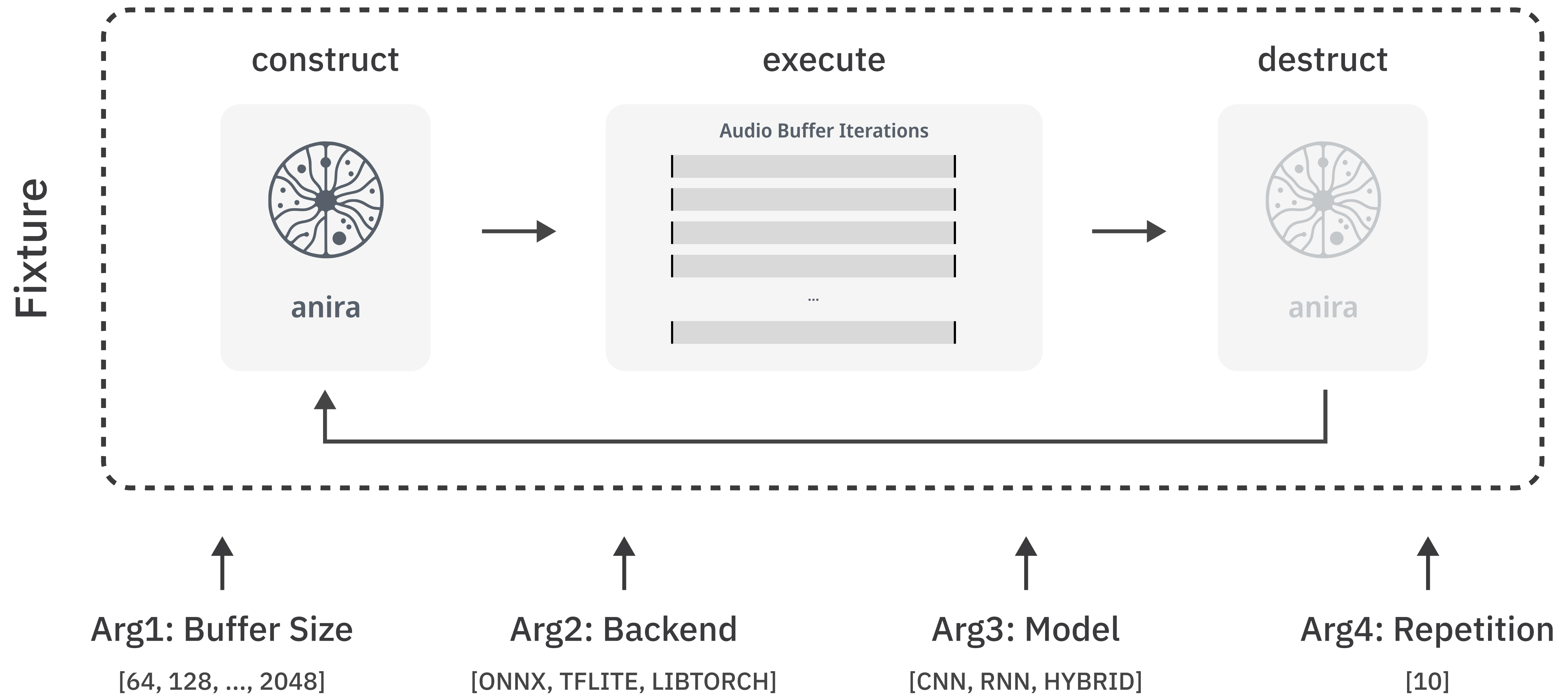
Chapter V

# **Impact on Inference Runtimes**

**Various Factors Affecting Performance**

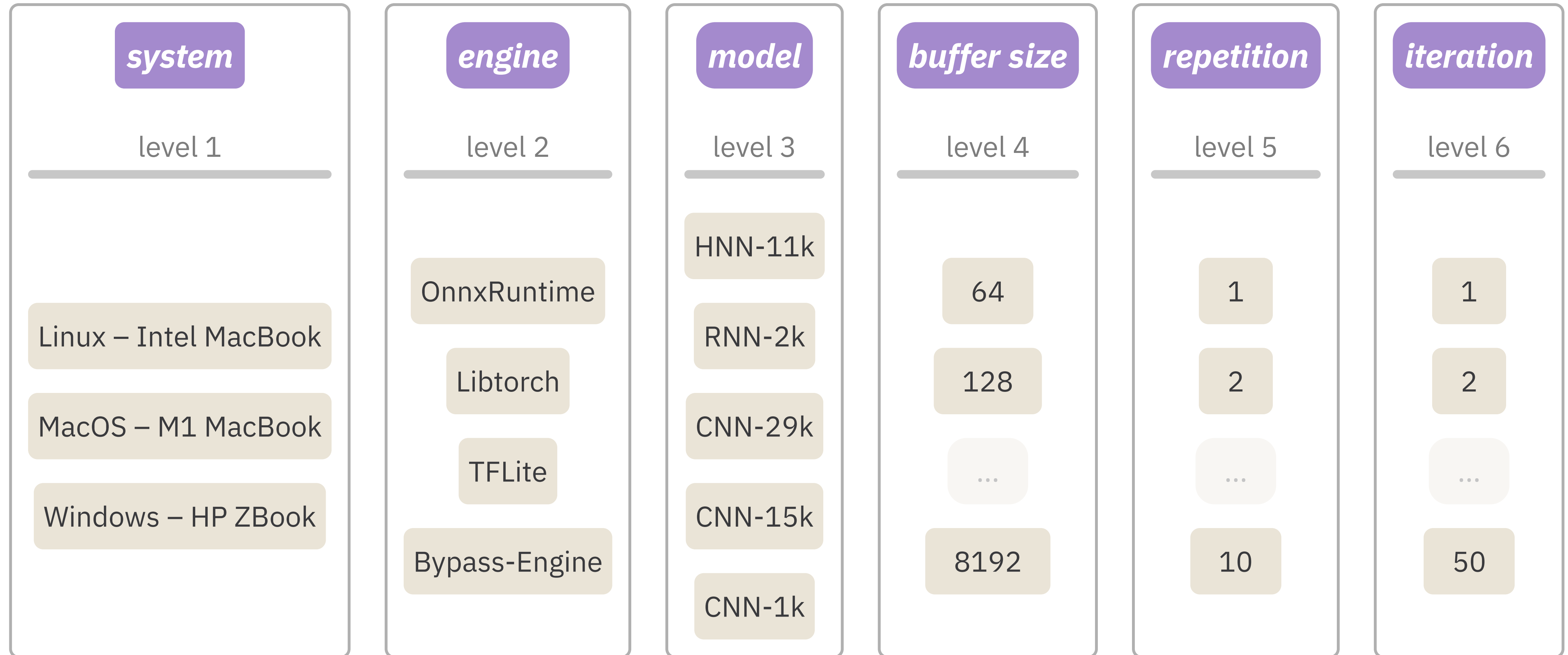


# Benchmarking



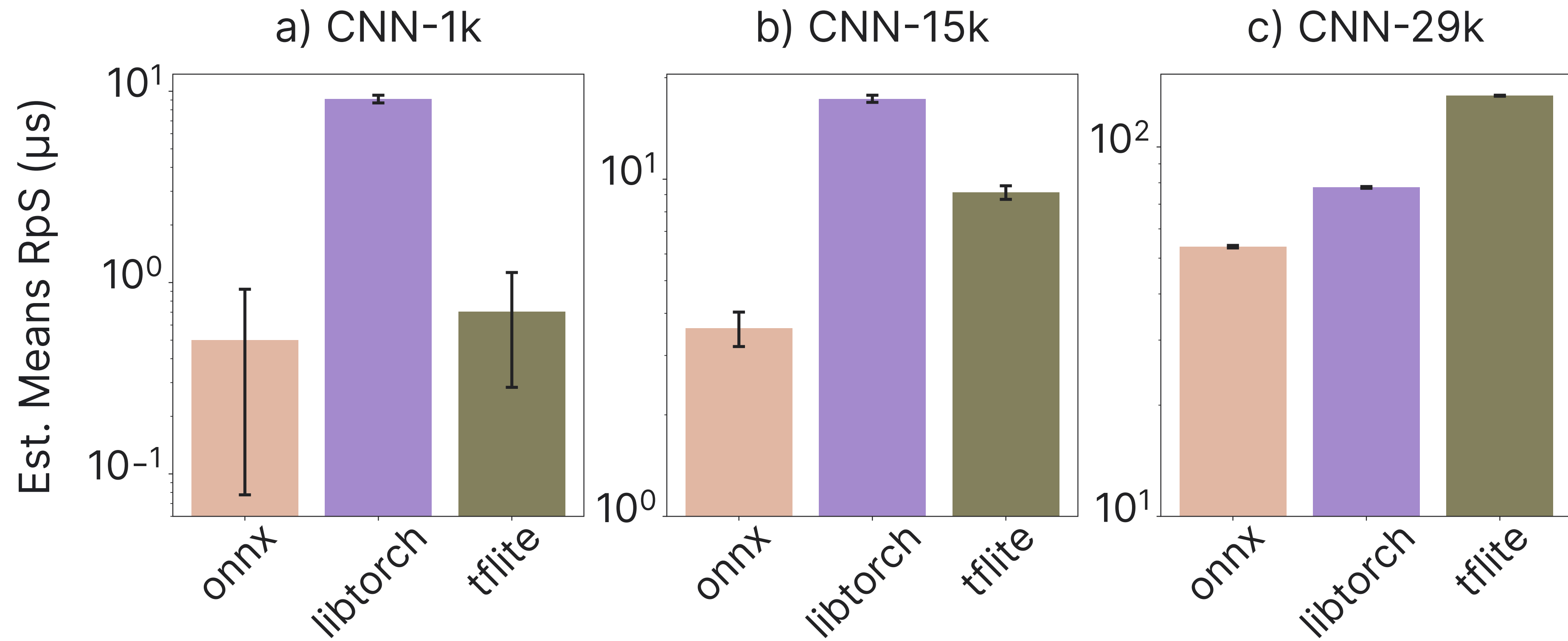


# Measurements

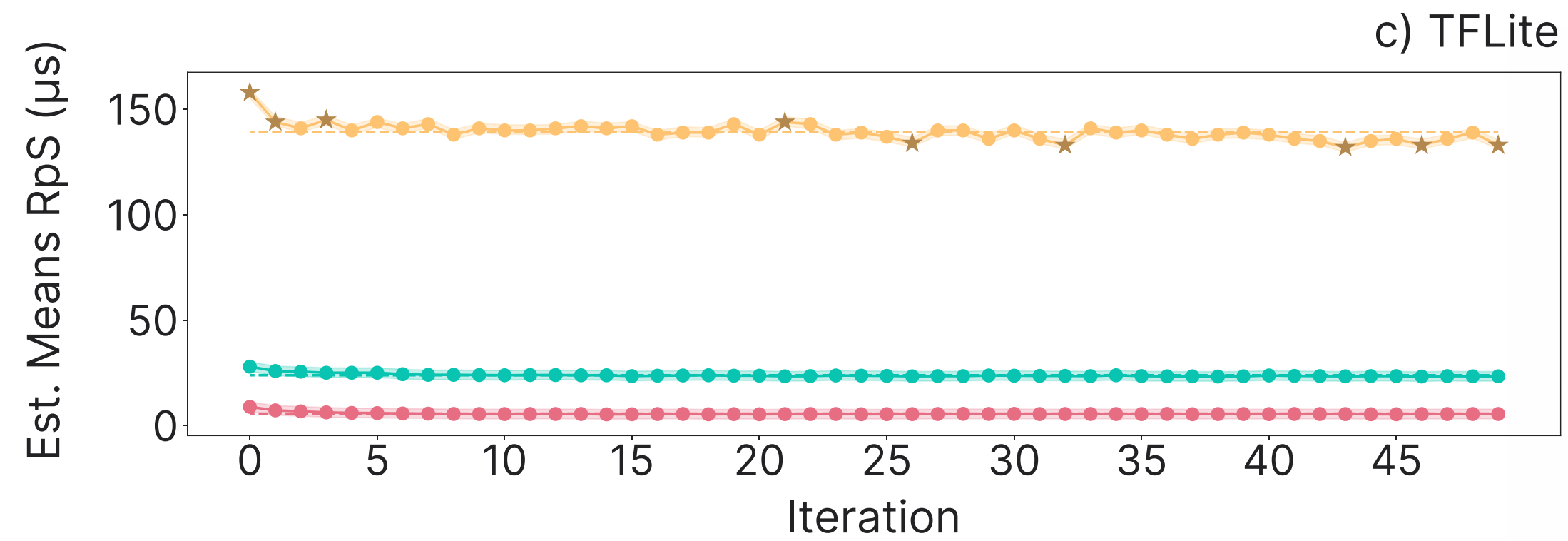
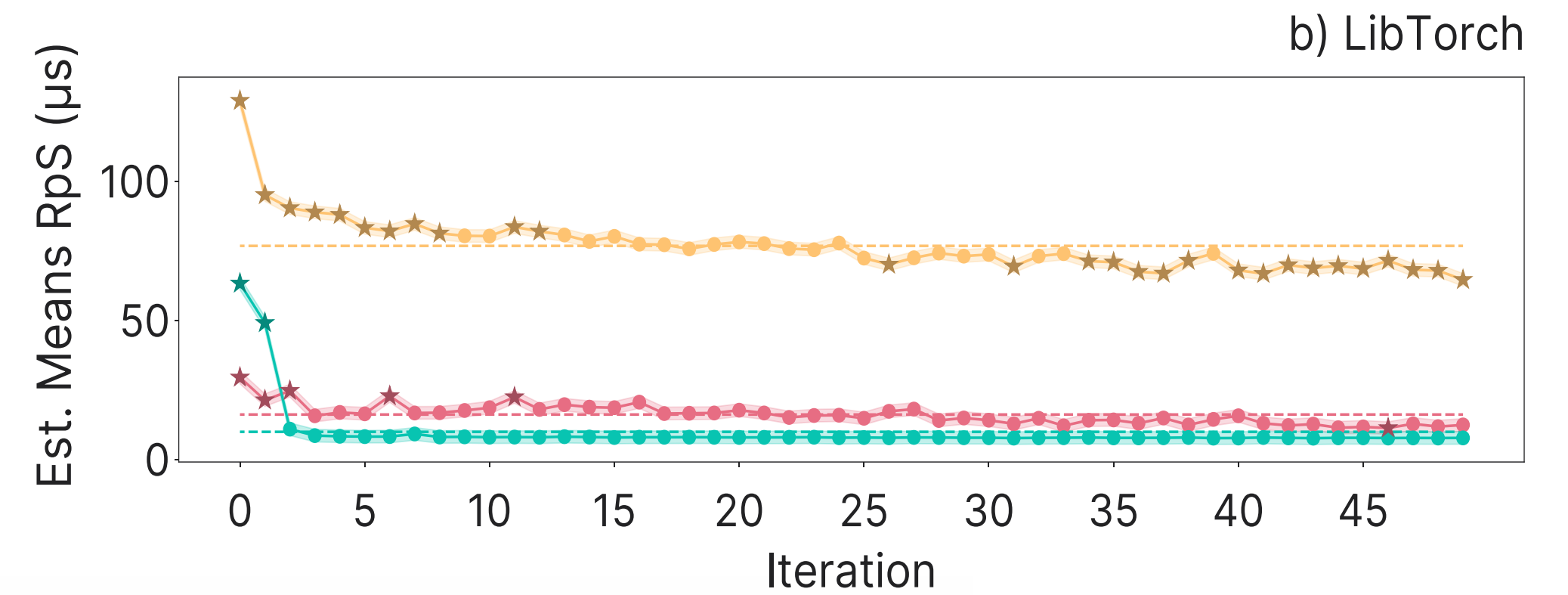
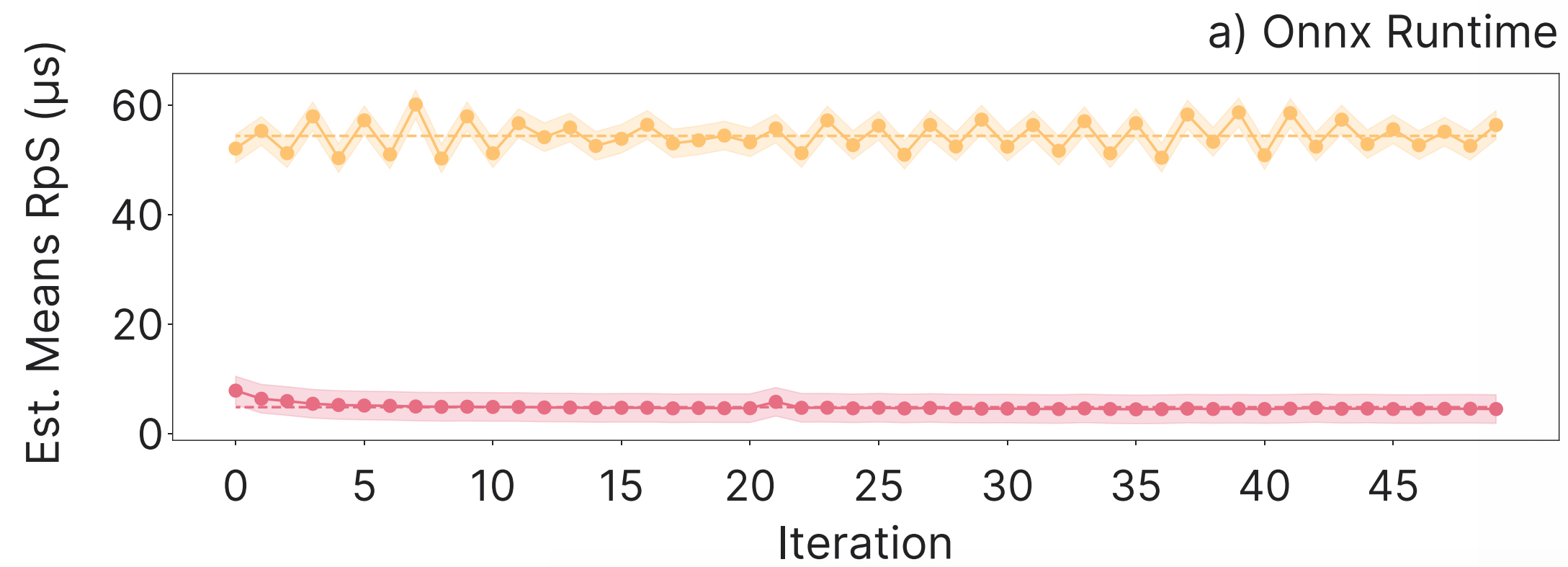




# CNN Hyperparameter Comparison



# Influence of the Iteration

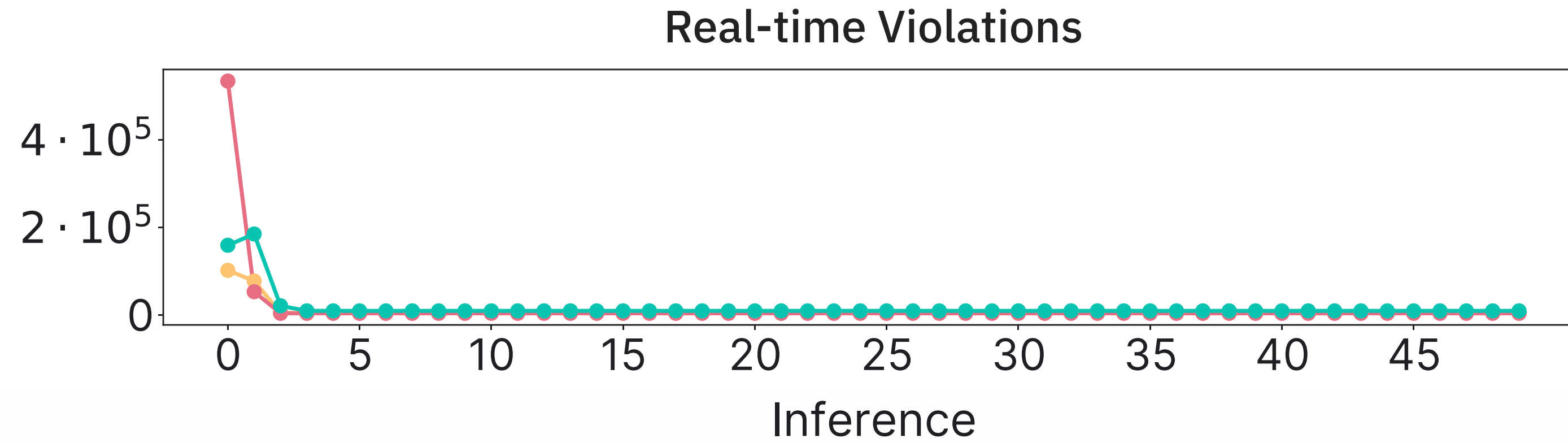


- HNN-11k
- CNN-29k
- RNN-2k

- Average
- Significant Deviation

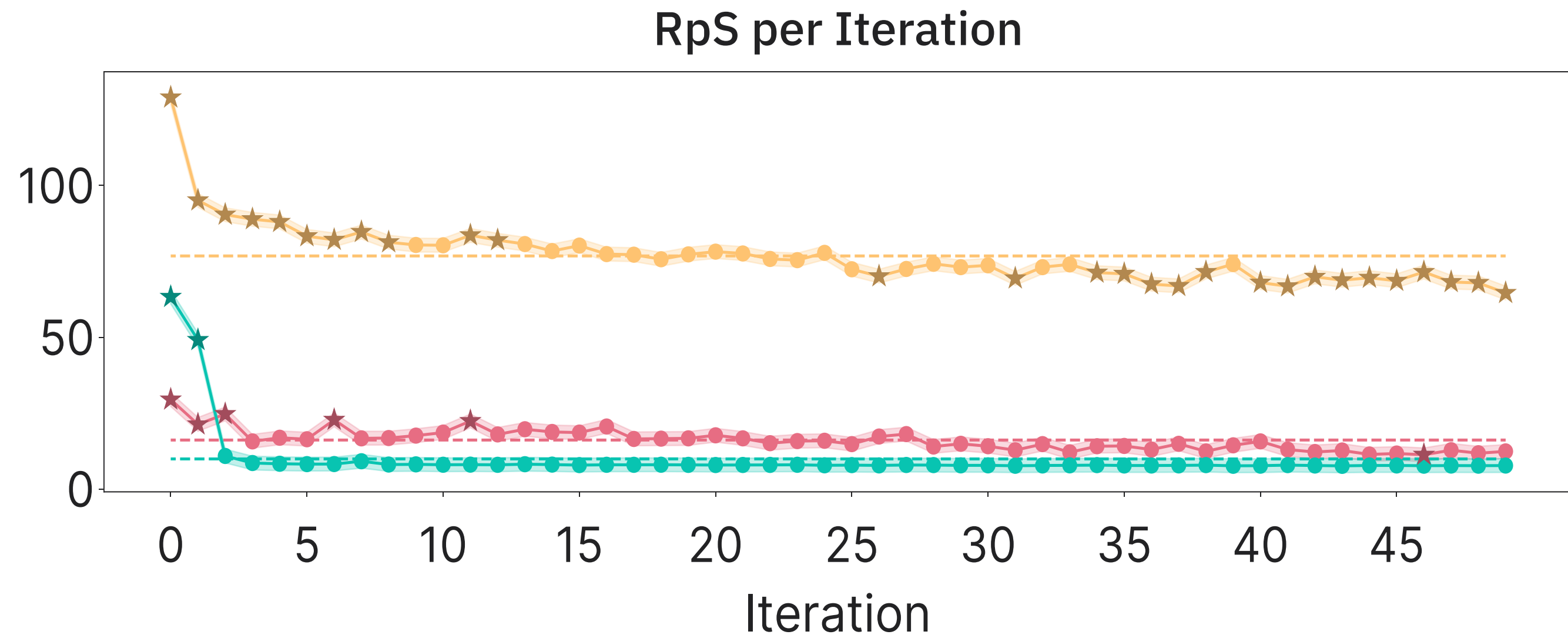
# Violations vs. Runtime Performance - LibTorch

Total RT-Violations



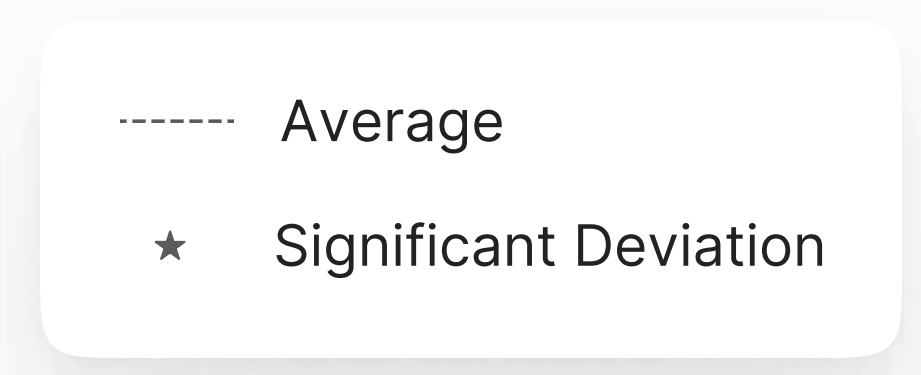
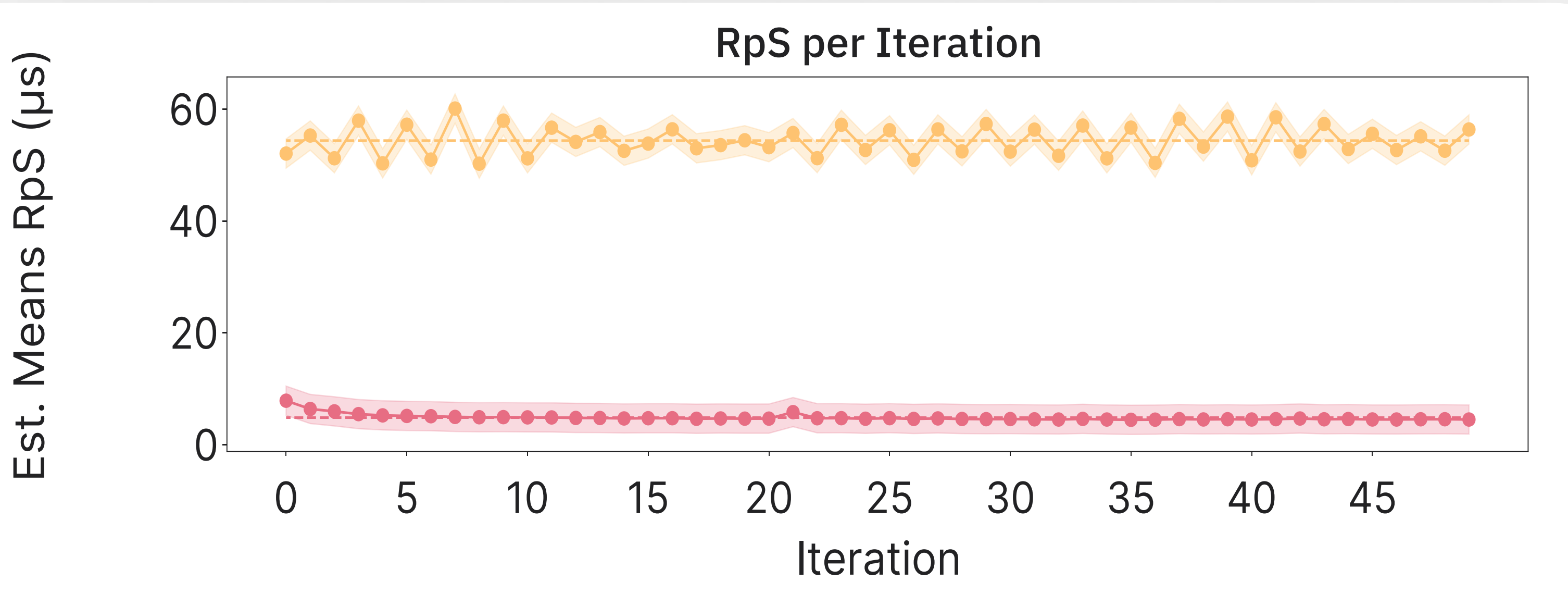
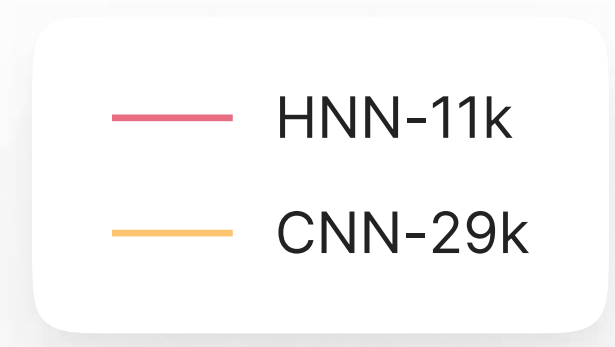
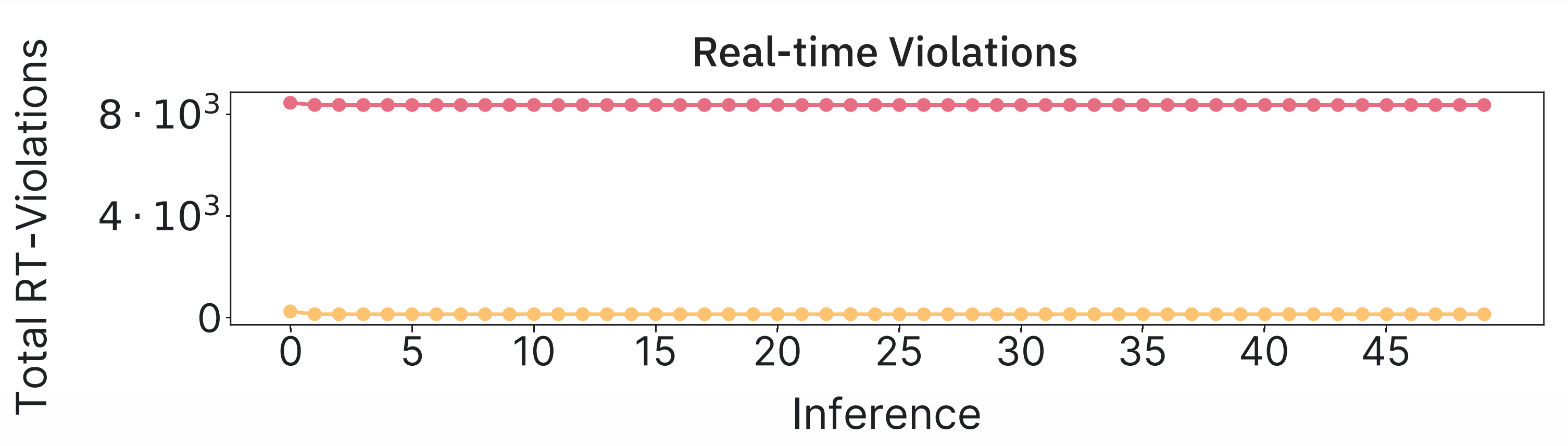
- HNN-11k
- CNN-29k
- RNN-2k

Est. Means RpS ( $\mu$ s)



- Average
- Significant Deviation

# Violations vs. Runtime Performance - Onnx Runtime







Chapter VI

## **Conclusion**

**Summary, Getting Started, Questions**



# Library Feature Overview

## Thread-Safe Inference Engine Wrapper Library:

- Supports major inference engines
- Compatible with various neural network types
- Enables handling of multiple session
  - Across different plugins (and formats) with multiple instances
  - Across various neural network types

## C++ Cross Platform Library:

- available on Windows (x64), Linux (x64, aarch64, armv7l), macOS (x64, arm64)

**Permissive License: Apache 2.0**

# Objectives of the Library

Implementing real-time safe inference requires expertise across various domains

**The library aims to simplify real-time safe inference implementation**



# Getting Started

We provide five example neural networks, each featuring:

- ML implementation
- Training code and data
- Pre-trained models

Every example can be adapted into a real-time application, available as:

- Audio plugin using the JUCE framework
- Native CLAP plugin, Bela embedded application
- Soon: JACK client, Max external

Additionally, we include example benchmarks to validate runtimes



# Paper

## ANIRA: An Architecture for Neural Network Inference in Real-Time Audio Applications

Valentin Ackva \*  
Audio Communication Group  
Technische Universität Berlin  
Berlin, Germany  
valentin.ackva@gmail.com

Fares Schulz \*  
Audio Communication Group  
Technische Universität Berlin  
Berlin, Germany  
fares.schulz@tu-berlin.com

**Abstract**—Numerous tools for neural network inference are currently available, yet many do not meet the requirements of real-time audio applications. In response, we introduce *anira*, an efficient cross-platform library. To ensure compatibility with a broad range of neural network architectures and frameworks, *anira* supports ONNX Runtime, LibTorch, and TensorFlow Lite as backends. Each inference engine exhibits real-time violations, which *anira* mitigates by decoupling the inference from the audio callback to a static thread pool. The library incorporates built-in latency management and extensive benchmarking capabilities, both crucial to ensure a continuous signal flow. Three different neural network architectures for audio effect emulation are then subjected to benchmarking across various configurations

and then export and load the trained parameters from the high-level framework. While this approach provides the ability to optimize the inference code for the neural network being used, it is often time-consuming and lacks the flexibility to integrate different neural networks.

To facilitate this process, a variety of inference engines have been developed in high-performance languages. These engines optimize the execution of neural networks on different hardware platforms and provide implementations of the most common neural network layers. In this paper, we focus on the three most common inference engines: ONNX Run-



<https://doi.org/10.1109/IS262782.2024.10704099>



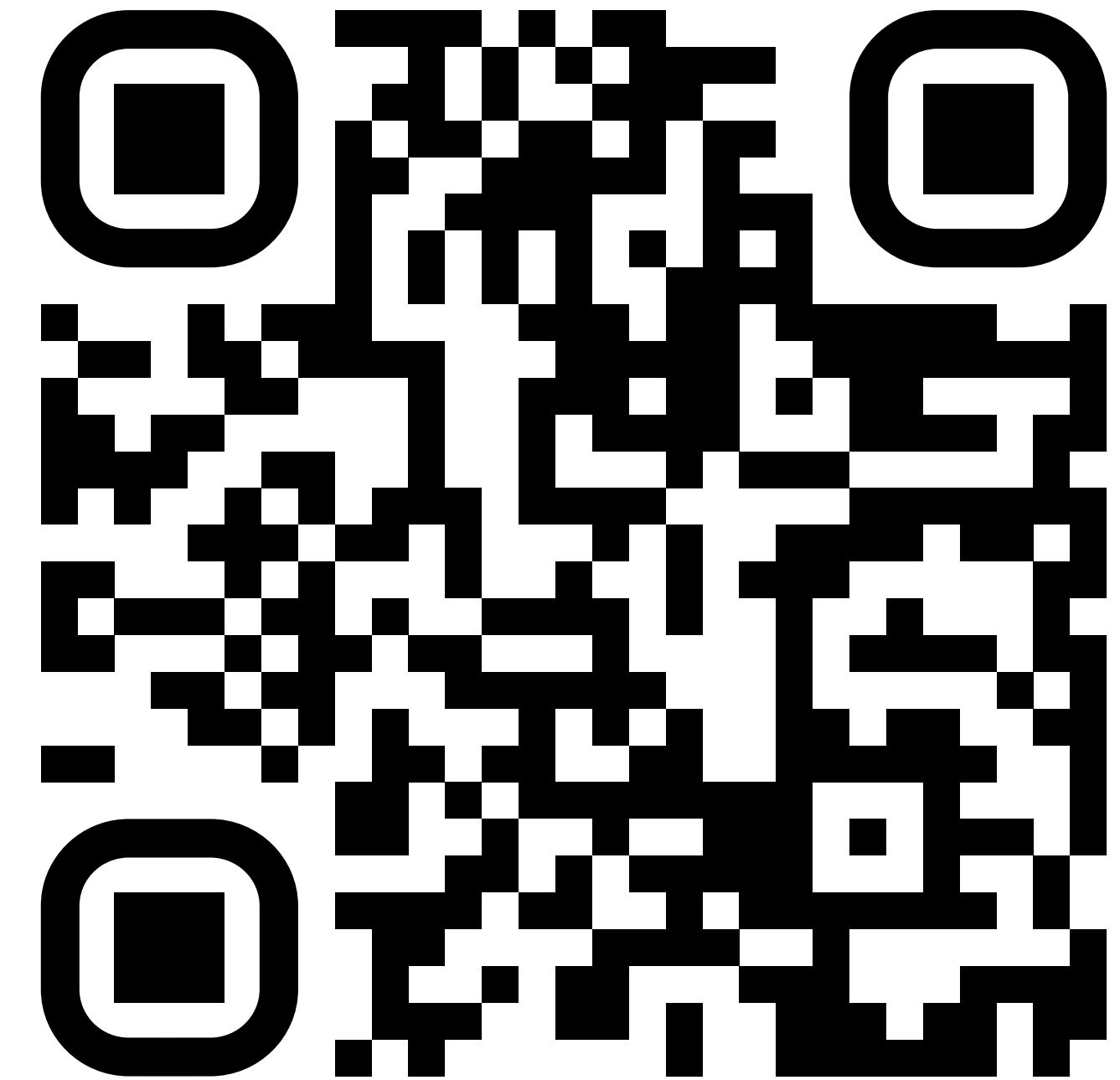


# Repository



**anira**

an architecture for neural network inference  
in real-time audio applications



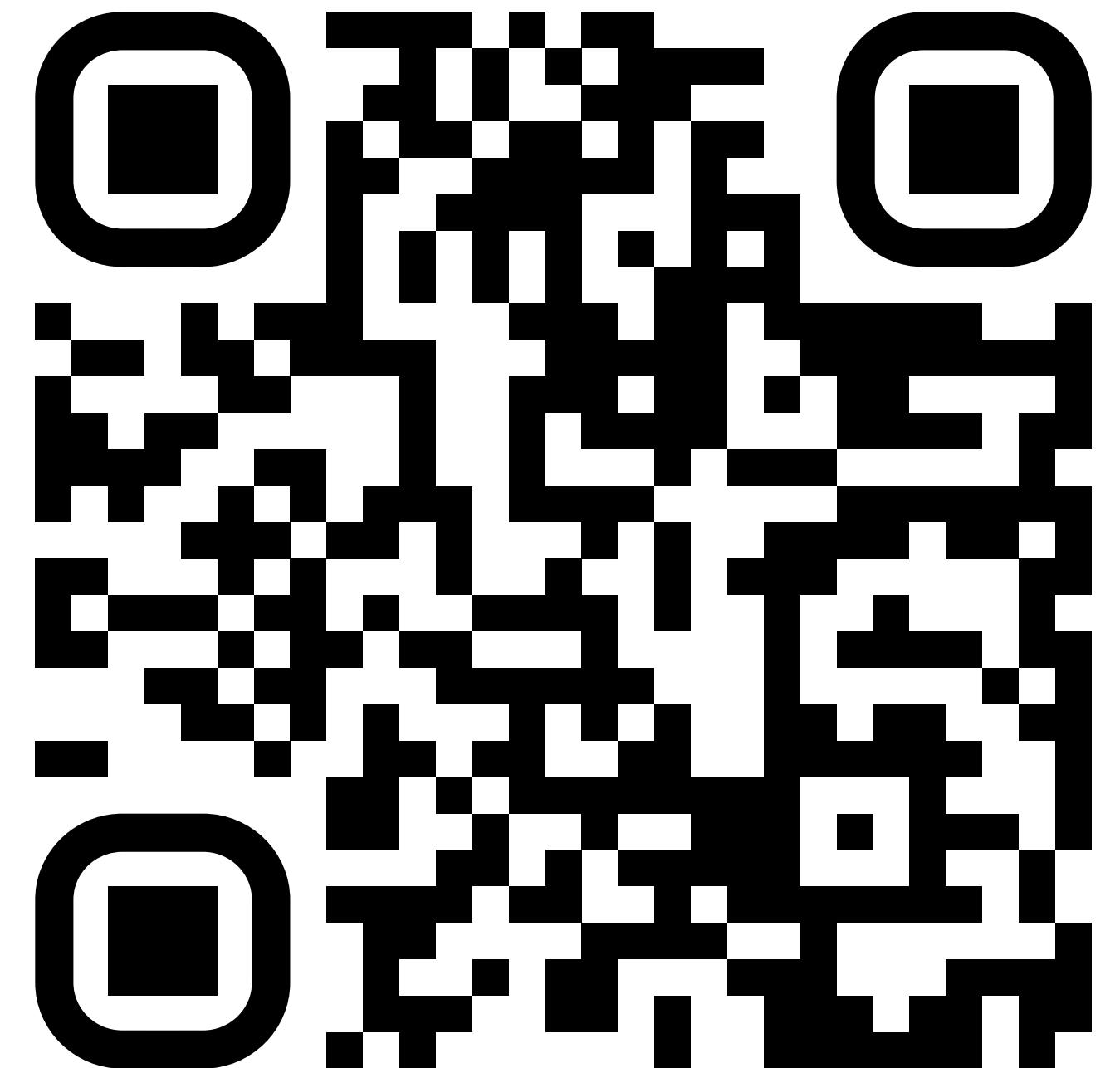
<https://github.com/anira-project/anira>



# Thank you for listening

Do you have questions?

Repository:



Audio Communication Group  
Technische Universität Berlin