

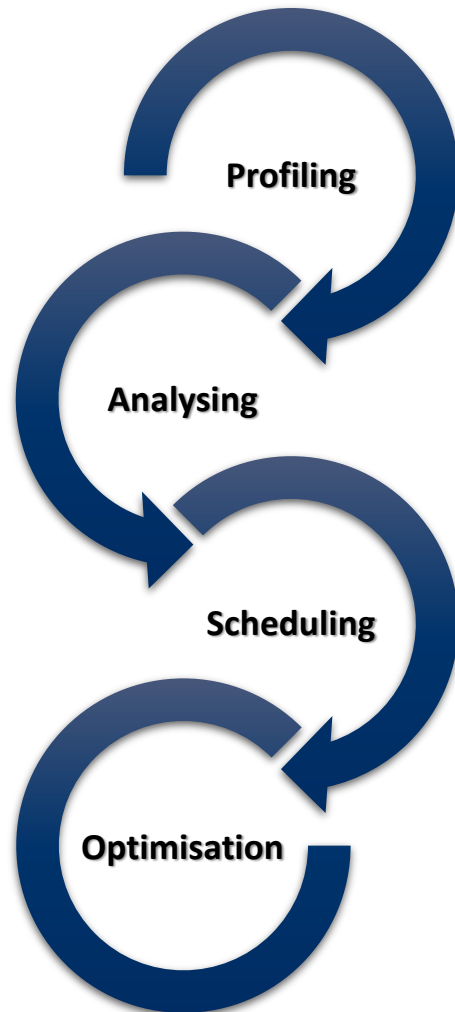
# Deployment on Edge Devices

Profiling, Analysing, Scheduling & Optimisation (PASO)

# Motivation - The Need for Optimisation

Aspect	Model-Centric Optimisation	System-Level Optimisation
Focus Area	Neural network architecture and parameters	Hardware, OS, memory, and data movement
Objective	Reduce model size, complexity, and computation	Maximize hardware utilization and system efficiency
Techniques	Pruning, Quantization, Knowledge Distillation, NAS	I/O Optimization, Memory Management, Scheduling
Benefits	Lower memory and computation costs	Faster and more efficient inference
Challenges	Maintaining accuracy after optimization	Balancing resource contention and system overhead

# Deployment Optimisation Process



Aspect	Techniques	Benefits
Advanced Profiling	Fine-grained performance analysis, Memory access patterns	Identifies bottlenecks and guides optimizations
Analysing Bottlenecks	CPU/GPU utilization, Memory/I/O, OS-level	Maximizes hardware utilization
Dynamic Scheduling	Task scheduling, Model partitioning, Adaptive precision	Reduces latency and improves resource efficiency
Advanced Optimization	Hardware acceleration, Code optimization, Graph optimization	Accelerates computation and minimizes power usage

# Advanced Profiling Techniques



Technique	Description	Tools/Examples	Benefits
<b>Model Performance Analysis</b>	Analyzes execution time of individual layers, operations, and functions. Identifies bottlenecks like convolutions or activation functions.	<ul style="list-style-type: none"><li>- <b>TensorBoard Profiler</b>: Layer-wise execution time and memory usage.</li><li>- <b>NVIDIA Nsight Systems</b>: GPU-specific profiling for CUDA operations.</li></ul>	Pinpoints the most time-consuming layers for targeted optimizations.
<b>Memory Analysis</b>	Examines memory access patterns to detect cache misses and inefficient data movement.	<ul style="list-style-type: none"><li>- <b>Valgrind's Cachegrind</b>: Analyses cache performance.</li><li>- <b>ARM Streamline</b>: Memory usage profiling on ARM devices.</li></ul>	Optimizes data locality and cache usage, reducing memory bottlenecks.
<b>Power Profiling</b>	Measures power consumption to identify energy-intensive operations. Essential for battery-powered edge devices.	<ul style="list-style-type: none"><li>- <b>ARM Energy Probe</b></li><li>- <b>Intel Power Gadget</b></li></ul>	Optimizes power usage and scheduling strategies for energy efficiency.
<b>System-Level Profiling</b>	Analyzes low-level metrics such as cycles, instructions, cache misses and branch mispredictions.	<ul style="list-style-type: none"><li>- <b>perf</b>: Profiling tool for Linux systems, e.g. <code>perf stat -e cycles, instructions my_program</code></li></ul>	Provides deep insights into CPU and hardware behavior (system-level bottlenecks) beyond model-specific issues.
<b>Network Profiling</b>	Captures and analyzes network traffic to detect latency, bandwidth limitations, and packet loss.	<ul style="list-style-type: none"><li>- <b>tcpdump, Wireshark</b></li><li>- Network latency and bandwidth usage diagrams.</li></ul>	Identifies network-related performance issues in edge-cloud communication.



# Case Study: CPU & Memory Bottleneck Analysis (*perf*)

- ***perf*** is a powerful profiling tool for Linux systems that provides detailed insights into CPU and system performance. It tracks low-level metrics such as *CPU cycles*, *instructions executed*, *cache misses* and *branch mispredictions*.
- **Scenario:** Profiling an object detection model running on a Raspberry Pi 4 to analyze CPU utilization, cache misses, and branch mispredictions.
- **Objective:** Identify CPU and memory bottlenecks to optimize model inference speed.

```
sudo taskset -c 2 python3 object_detection.py
```

```
sudo perf stat -C 2 -e cycles,instructions,cache-references,cache-misses,branch-misses -t <process_id>
```

```
sudo perf stat -e cycles,instructions,cache-references,cache-misses,branch-misses python3 object_detection.py
```

```
Performance counter stats for 'python3 object_detection.py':
      2,345,678,901      cycles                  #    1.50 GHz
      3,456,789,012      instructions             #    1.47  insns per cycle
      234,567,890       cache-references
      12,345,678        cache-misses                #    5.27% of all cache refs
       1,234,567        branch-misses              #    1.02% of all branches
```

## Key Insights:

- **Low IPC:** Indicates a CPU bottleneck or poor instruction parallelism.
- **Cache Misses:** Suggests memory bottlenecks due to poor data locality.
- **Branch Misses:** Indicates control flow changes that can affect pipeline efficiency.

## Optimization Actions:

- **Improve Cache Usage:** Optimize data structures and memory access patterns for better cache locality. Use smaller model weights or quantisation techniques to fit in the cache.
- **Reduce Branch Mispredictions:** Minimize complex conditional statements in the critical path. Refactor loops and branches for more predictable control flow.

# Case Study: CPU & Memory Bottleneck Analysis (*perf*)

```
sudo perf record -F 99 -g python3 object_detection.py
```

Overhead	Command	Shared Object	Symbol
30.15%	python3	/usr/lib/arm-linux-gnueabi/libc-2.28.so	[.] memcpy
20.23%	python3	/usr/lib/python3.7/site-packages/numpy/core/umath.cpython-37m-arm-linux-gnueabi.so	[.] ufunc_generic_call
15.10%	python3	/usr/lib/arm-linux-gnueabi/libc-2.28.so	[.] __memcpy_arm
10.12%	python3	/usr/lib/python3.7/site-packages/tensorflow/python/_pywrap_tensorflow_internal.so	[.] _tensorflow::CopyData
8.05%	python3	/usr/lib/arm-linux-gnueabi/libpthread-2.28.so	[.] __pthread_mutex_lock
7.02%	python3	/usr/lib/arm-linux-gnueabi/libc-2.28.so	[.] malloc
4.18%	python3	/usr/lib/python3.7/site-packages/tensorflow/python/_pywrap_tensorflow_internal.so	[.] _tensorflow::Conv2DCompute
3.15%	python3	/usr/lib/python3.7/site-packages/numpy/core/umath.cpython-37m-arm-linux-gnueabi.so	[.] reduce_sum
2.00%	python3	/usr/lib/arm-linux-gnueabi/libc-2.28.so	[.] free

## Explanation of perf Report:

- **Overhead:** % of total CPU time spent on a specific function or library.
- **Command:** The executable that triggered the event.
- **Shared Object:** The shared library or binary file where the event occurred.
- **Symbol:** The specific function or symbol being executed

## Some Insights:

Event	Description	Possible Optimization Actions
30.15% memcpy	Frequent data copying between memory buffers.	- <b>Reduce Memory Copy Operations:</b> Optimize tensor data structures to minimize copying. - <b>Use In-Place Operations:</b> Modify data in place to avoid unnecessary copies.
15.10% __memcpy_arm	ARM-specific memory copy operation, confirming memory-bound bottlenecks.	- <b>Optimize Data Layout:</b> Reorder data structures to improve cache locality. - <b>Use NumPy Broadcasting:</b> Leverage broadcasting to minimize memory operations.
7.02% malloc	frequent memory allocations due to dynamic creation of temporary buffers in NumPy and TensorFlow operations.	Pre-allocating memory outside loops and reusing memory buffers to minimise heap allocations.

# Analysing System-Level Bottlenecks



Aspect	Description	Tools/Examples	Optimization Approach
<b>CPU/GPU Utilization</b>	Identifies under-utilization or overloading of processing units.	<ul style="list-style-type: none"><li>- <b>Perf</b> (Linux Performance Profiler)</li><li>- <b>NVIDIA Visual Profiler (nvvp)</b></li></ul>	<ul style="list-style-type: none"><li>• Load balancing between CPU and GPU.</li><li>• Utilize heterogeneous computing (e.g., NPUs).</li></ul>
<b>Memory and I/O</b>	Analyzes latency due to memory access or data transfer operations.	<ul style="list-style-type: none"><li>- <b>gprof, Valgrind</b></li><li>- I/O solutions: Faster storage (e.g., SSD), buffering data</li></ul>	<ul style="list-style-type: none"><li>• Memory Pooling: Reuse memory buffers.</li><li>• Layer Fusion: Reduce intermediate data movement.</li><li>• Optimize data access patterns.</li></ul>
<b>OS-Level</b>	Examines the impact of OS scheduling, resource allocation, and context switching.	<ul style="list-style-type: none"><li>- <b>Zephyr RTOS</b> for real-time constraints</li><li>- Lightweight runtime (e.g., TensorFlow Lite Micro)</li></ul>	<ul style="list-style-type: none"><li>• Real-Time Scheduling: Prioritize ML tasks.</li><li>• Minimize OS overhead with lightweight runtimes.</li></ul>
<b>Concurrency Issues</b>	Addresses thread contention, locking, and synchronization challenges.	<ul style="list-style-type: none"><li>- <b>Thread Pools</b></li><li>- <b>Asynchronous Operations</b></li></ul>	<ul style="list-style-type: none"><li>• Use thread pools and asynchronous operations.</li><li>• Implement lock-free data structures.</li></ul>

# Case Study: Smart Camera for Object Detection

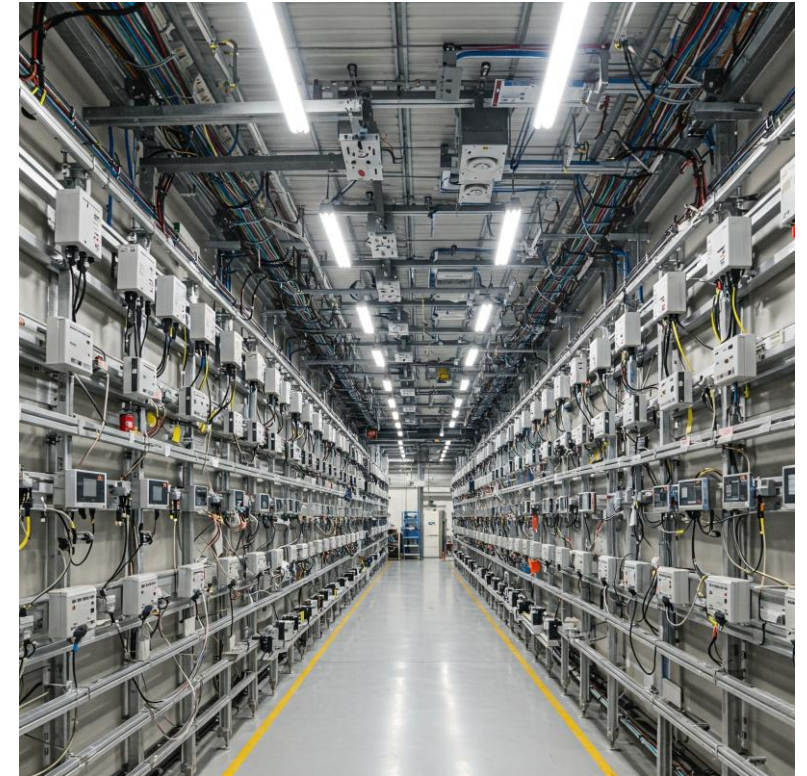
- **Application:** An edge-based smart camera performing real-time object detection & alerting.
- **Bottleneck:** Camera transmits results with high latency to the cloud for processing.
- **Tools:** *tcpdump* and *Wireshark*, network traffic analysis revealed high transmission delays due to limited bandwidth and intermittent connectivity.
- **Possible Solution:**
  - *Edge-Cloud Partitioning:* Critical detection logic was moved to the edge device, sending only summarised alerts to the cloud, reducing data.
  - *Protocol Optimisation:* Switching from HTTP to MQTT for more efficient, lightweight communication.





# Case Study: Industrial Sensor Network

- **Application:** A large-scale sensor network monitoring environmental parameters in a manufacturing facility.
- **Bottleneck:** The gateway device experienced delays in aggregating and processing data from hundreds of sensors.
- **Tools:** *perf* and *Valgrind* profiling showed CPU under usage due to inefficient data handling and I/O operations.
- **Possible Solution:**
  - *Memory Pooling:* Implemented buffer reuse to reduce memory allocation overhead.
  - *Layer Fusion:* Combined data pre-processing layers to minimise intermediate data movement.
  - *Concurrency Optimization:* Introduced asynchronous data processing and lock-free data structures.

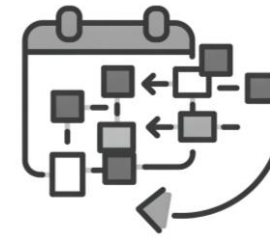


# Case Study: Autonomous Drone Navigation

- **Application:** An edge-based drone using deep learning for obstacle detection and navigation.
- **Bottleneck:** The GPU was overloaded while the CPU was underutilised, leading to navigation lag.
- **Tools:** Profiling with *NVIDIA Nsight Systems* showed GPU saturation during inference while CPU resources were idle.
- **Possible Solution:**
  - *Load Balancing:* Offloaded pre-processing and sensor fusion tasks to the CPU, optimising GPU utilisation for model inference only.
  - *Heterogeneous Computing:* Integrated an NPU to accelerate lightweight tasks, balancing the workload across CPU, GPU, and NPU.

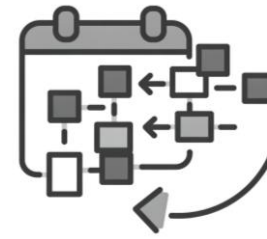


# Dynamic and Adaptive Scheduling



Technique	Description	Approach	Benefits	Examples
<b>Dynamic Task Scheduling</b>	Dynamically assigns tasks to processing units (CPU, GPU, NPU) based on resource availability.	<ul style="list-style-type: none"><li>- <b>Work Stealing Algorithms:</b> Idle processors take over tasks from busy ones.</li><li>- <b>Priority-Based Scheduling:</b> Assigns priority levels to tasks (e.g., high, medium, low).</li></ul>	<ul style="list-style-type: none"><li>- Maximizes hardware utilization.</li><li>- Reduces latency by prioritizing critical tasks.</li></ul>	<ul style="list-style-type: none"><li>- <b>TVM Runtime:</b> Dynamic scheduling for heterogeneous devices.</li><li>- <b>Preemptive Scheduling:</b> Interrupts lower-priority tasks for high-priority ones.</li></ul>
<b>Adaptive Model Partitioning</b>	Splits model execution across multiple devices (e.g., edge and cloud) and adapts partitioning based on runtime conditions like network latency.	<ul style="list-style-type: none"><li>- <b>Edge-Cloud Offloading:</b> Executes latency-critical layers on the edge and offloads other layers to the cloud.</li></ul>	<ul style="list-style-type: none"><li>- Balances latency, energy, and computational load.</li><li>- Improves performance under varying network conditions.</li></ul>	<ul style="list-style-type: none"><li>- <b>Neurosurgeon Framework:</b> Real-time model partitioning between edge and cloud.</li></ul>
<b>Adaptive Resource Allocation</b>	Dynamically allocates system resources (CPU, memory) based on current system load and task requirements.	<ul style="list-style-type: none"><li>- Allocates more resources to high-priority or critical tasks during peak load.</li><li>- Minimizes resource allocation for background tasks during idle periods.</li></ul>	<ul style="list-style-type: none"><li>- Enhances system responsiveness and efficiency.</li><li>- Ensures optimal resource usage without overloading.</li></ul>	<ul style="list-style-type: none"><li>- <b>Real-Time Operating Systems (RTOS):</b> Adaptive scheduling in time-critical applications.</li></ul>

# Honourable mentions



Technique	Description	Approach	Benefits	Examples
<b>Adaptive Precision and Quantization</b>	Dynamically adjusts model precision (e.g., FP32 to INT8) based on available resources and accuracy requirements.	- <b>Hardware-Aware Quantization:</b> Adapts precision to match hardware constraints, optimizing for both accuracy and efficiency.	- Maintains a balance between accuracy and computational efficiency. - Reduces power consumption and memory usage.	- <b>HAQ (Hardware-Aware Quantization):</b> Adjusts precision based on hardware limitations.
<b>Containerization for Edge</b>	Packages applications and dependencies into containers for portability and isolated resource management on edge devices.	- <b>Docker and Containers:</b> Ensures consistent environment across edge devices. - Simplifies deployment and scaling of ML workloads.	- Enhances portability and isolation. - Efficient resource management and scheduling.	- <b>Docker Containers:</b> Edge deployment with consistent environments and efficient scaling.



# Case Study: Smart Security Camera with Motion-Activated Object Detection

**Imagine:** You have a security camera to detect people or objects only when there's movement, saving power and resources when nothing is happening.

## Smart Scheduling:

- **Movement Detected:** The camera immediately runs object detection to see what's happening.
- **No Movement:** Object detection pauses or runs in the background at a lower priority, freeing it to do other things or saving power.

## What You Need:

- **Motion Sensor:** A sensor that detects movement.
- **Camera:** To capture video.
- **Software:**
  - **OpenCV:** Helps the camera understand the video stream.
  - **Multiprocessing:** The MCU can perform object detection as a separate task.
  - **Object Detection Model:** A pre-trained model like MobileNet SSD, which is small and fast.

## Why This is Cool:

- **Saves Power:** The camera doesn't need to run object detection all the time, so it uses less power.
- **Faster Response:** When something happens, the camera reacts instantly.
- **Efficient:** The Raspberry Pi can do other things when it's not busy detecting objects.

*Think of it like a guard dog that only wakes up when it hears something!*





## Case Study: Smart Home Camera Performing Real-time Face Recognition

**Imagine:** You have a home security camera that can recognise faces. It can tell if a family member is at the door or a stranger is approaching. The camera can quickly alert you if an unknown person is at your door, even if the internet is down.

**Challenge:** High latency and power consumption when running complex face recognition models entirely on the edge device.

### Adaptive Model Partitioning:

- The camera's face recognition "brain" is split into two parts.
- **Local:** The camera uses its own "brain" to quickly analyse faces and see if they're familiar.
- **Cloud:** If the camera needs extra brainpower to recognise someone, it sends the video to a powerful computer in the cloud for help. This happens when the internet connection is strong, and the camera isn't too busy.

### Why This is Cool:

- **Fast Response:** Critical layers run on the edge, ensuring quick response.
- **Saves Energy:** Heavy layers offloaded to the cloud, reducing power consumption on the edge device.
- **Reduces Cost:** Minimizes cloud usage by running lightweight layers locally.

***Think of it like a security guard with a walkie-talkie. They can handle most situations themselves, but they can call for backup from the central office when needed!***

## Case Study: Smart Farm with Adaptive Crop Monitoring

**Imagine:** A network of cameras and sensors in a greenhouse monitors the health and growth of crops. The system needs to detect pests, diseases, and nutrient deficiencies, while also optimizing water and fertilizer usage.

### How it Works:

- **Smart Priorities:** The system's "brain" shifts focus based on the plants' needs.
  - **Young Plants:** More attention to pest and disease detection.
  - **Mature Plants:** Focus on optimising water and nutrients.
- **Quick Response:** If a problem is found, the system focuses more resources on that area.

**Challenge:** Continuously running all analysis tasks at full power would consume excessive energy and processing capacity.

**Adaptive Resource Allocation:** The system dynamically adjusts how much processing power and sensor sampling frequency is dedicated to each task based on the current needs of the crops and the presence of any issues.

*Think of it like a farmer tending to their crops. They pay close attention to young seedlings, providing extra care and attention. As the plants mature, the farmer adjusts their focus to ensure optimal growing conditions and address any specific needs!*

### Why This is Cool:

- **Healthy Crops:** By focusing on the most critical tasks at each growth stage, the system helps ensure healthy crop development and maximises yield.
- **Resource Efficiency:** The system conserves energy and processing power by allocating resources only where and when they are most needed.
- **Early Problem Detection:** By dynamically increasing resource allocation to potential problem areas, the system enables early intervention, preventing widespread issues.

# Advanced Optimisation Techniques



Focus Area	Description	Techniques	Examples	Challenges
Hardware Acceleration	Utilizes specialized hardware components (e.g., NPUs, TPUs, FPGAs) to accelerate ML computations and reduce latency.	<ul style="list-style-type: none"><li>- <b>Operator Fusion:</b> Combines multiple operations into one kernel to minimize memory access.</li><li>- <b>Domain-Specific Accelerators:</b> Custom accelerators for specific tasks (e.g., convolution, matrix multiplication).</li></ul>	<ul style="list-style-type: none"><li>- <b>Google Edge TPU:</b> Low-latency inference on edge devices.</li><li>- <b>NVIDIA TensorRT:</b> GPU acceleration for deep learning models.</li></ul>	<ul style="list-style-type: none"><li>- <b>Hardware-Specific Tuning:</b> Requires optimization for specific hardware.</li><li>- <b>Limited Portability:</b> Custom accelerators may not be compatible across devices.</li></ul>
Code Optimization	Refines code for faster execution, reduced memory usage, and lower power consumption.	<ul style="list-style-type: none"><li>- <b>Kernel Optimization:</b> Improves cache locality and reduces memory access.</li><li>- <b>Loop Unrolling and Tiling:</b> Enhances parallelism and data locality.</li><li>- <b>Vectorization:</b> Utilizes SIMD (Single Instruction, Multiple Data) for faster computation.</li></ul>	<ul style="list-style-type: none"><li>- <b>LLVM and GCC Compilers:</b> Low-level code optimization.</li><li>- <b>XLA (Accelerated Linear Algebra) Compiler:</b> Efficient kernel generation in TensorFlow.</li></ul>	<ul style="list-style-type: none"><li>- <b>Complexity vs. Maintainability:</b> Balancing performance gains with code maintainability.</li><li>- <b>Deep Hardware Knowledge Required:</b> Effective optimization demands understanding of underlying hardware architecture.</li></ul>
Graph-Level Optimization	Rearranges computation graphs for improved memory and computation efficiency.	<ul style="list-style-type: none"><li>- <b>Common Subexpression Elimination:</b> Avoids redundant computations.</li><li>- <b>Constant Folding:</b> Pre-computes constant expressions during compilation.</li></ul>	<ul style="list-style-type: none"><li>- <b>TensorFlow XLA and TVM:</b> Perform graph-level optimizations for efficient execution.</li></ul>	<ul style="list-style-type: none"><li>- <b>Graph Dependency Management:</b> Requires careful management of graph dependencies to avoid runtime errors.</li></ul>

# Future Trends

Trend	Description	Benefits	Challenges	Example Use Cases
<b>AI-Driven Optimization</b>	ML algorithms dynamically optimize system parameters and model architectures.	<ul style="list-style-type: none"><li>- Automated tuning</li><li>- Real-time adaptation</li><li>- Hardware-aware optimization</li></ul>	<ul style="list-style-type: none"><li>- High complexity</li><li>- Privacy concerns</li></ul>	<ul style="list-style-type: none"><li>- Google AutoML</li><li>- Alibaba RL Cluster Optimizer</li></ul>
<b>Dynamic and Adaptive Allocation</b>	AI-driven scheduling adjusts resources based on real-time demands.	<ul style="list-style-type: none"><li>- Optimized utilization</li><li>- Reduced latency</li><li>- Cost efficiency</li></ul>	<ul style="list-style-type: none"><li>- Complex scheduling</li><li>- Coordination overhead</li></ul>	<ul style="list-style-type: none"><li>- Alibaba Adaptive Resource Manager</li><li>- TVM Runtime</li></ul>
<b>Cross-Layer Optimization</b>	Jointly optimizes hardware, runtime, compiler, and application layers.	<ul style="list-style-type: none"><li>- Holistic efficiency</li><li>- Maximum performance gains</li><li>- Scalability</li></ul>	<ul style="list-style-type: none"><li>- High complexity</li><li>- Compatibility issues</li></ul>	<ul style="list-style-type: none"><li>- TVM Compiler</li><li>- Google Edge TPU Compiler</li></ul>
<b>Collaborative &amp; Federated Opt.</b>	Distributed optimization preserving data privacy and security.	<ul style="list-style-type: none"><li>- Data privacy</li><li>- Personalized models</li><li>- Scalable learning</li></ul>	<ul style="list-style-type: none"><li>- Non-IID data</li><li>- Communication overhead</li></ul>	<ul style="list-style-type: none"><li>- Google Federated Learning</li><li>- OpenMined PySyft</li></ul>

# Summary

- **Advanced Profiling Techniques** provide detailed insights into model and system performance, guiding further optimizations.
- **System-Level Bottlenecks** address hardware, memory, and OS constraints, ensuring efficient resource usage.
- **Dynamic and Adaptive Scheduling** enables real-time adjustments, maximizing performance under changing conditions.
- **Advanced Optimization Techniques** exploit hardware accelerators and low-level code enhancements for peak efficiency.

These strategies collectively enhance the deployment of ML models on edge devices, achieving high performance and low latency while minimizing power consumption.

---



## Some Relevant References

1. J. Zhang, C. Li, and Y. Wang, "**EdgeShard: Efficient LLM Inference via Collaborative Edge Computing**," arXiv preprint arXiv:2405.14371v1, 2024. [Online]. Available: <https://arxiv.org/html/2405.14371v1>.
  2. H. Kim, S. Park, and J. Lee, "**Inferencing on Edge Devices: A Time- and Space-aware Co-scheduling Framework**," Proceedings of the ACM on Embedded Systems, vol. 2, no. 1, pp. 1-13, 2023. [Online]. Available: <https://dl.acm.org/doi/10.1145/3576197>.
  3. X. Liu, M. Zhang, and K. Chen, "**Boosting DNN Cold Inference on Edge Devices**," arXiv preprint arXiv:2206.07446, 2022. [Online]. Available: <https://arxiv.org/abs/2206.07446>.
  4. Y. Sun, X. Wang, and Z. Li, "**BCEdge: SLO-Aware DNN Inference Services with Adaptive Batching on Edge Platforms**," arXiv preprint arXiv:2305.01519, 2023. [Online]. Available: <https://arxiv.org/abs/2305.01519>.
-