# SwapNet: Efficient Swapping for DNN Inference on Edge AI Devices Beyond the Memory Budget

Kun Wang, Jiani Cao, Zimu Zhou, Zhenjiang Li
City University of Hong Kong

香港城市大學
**City University of Hong Kong**

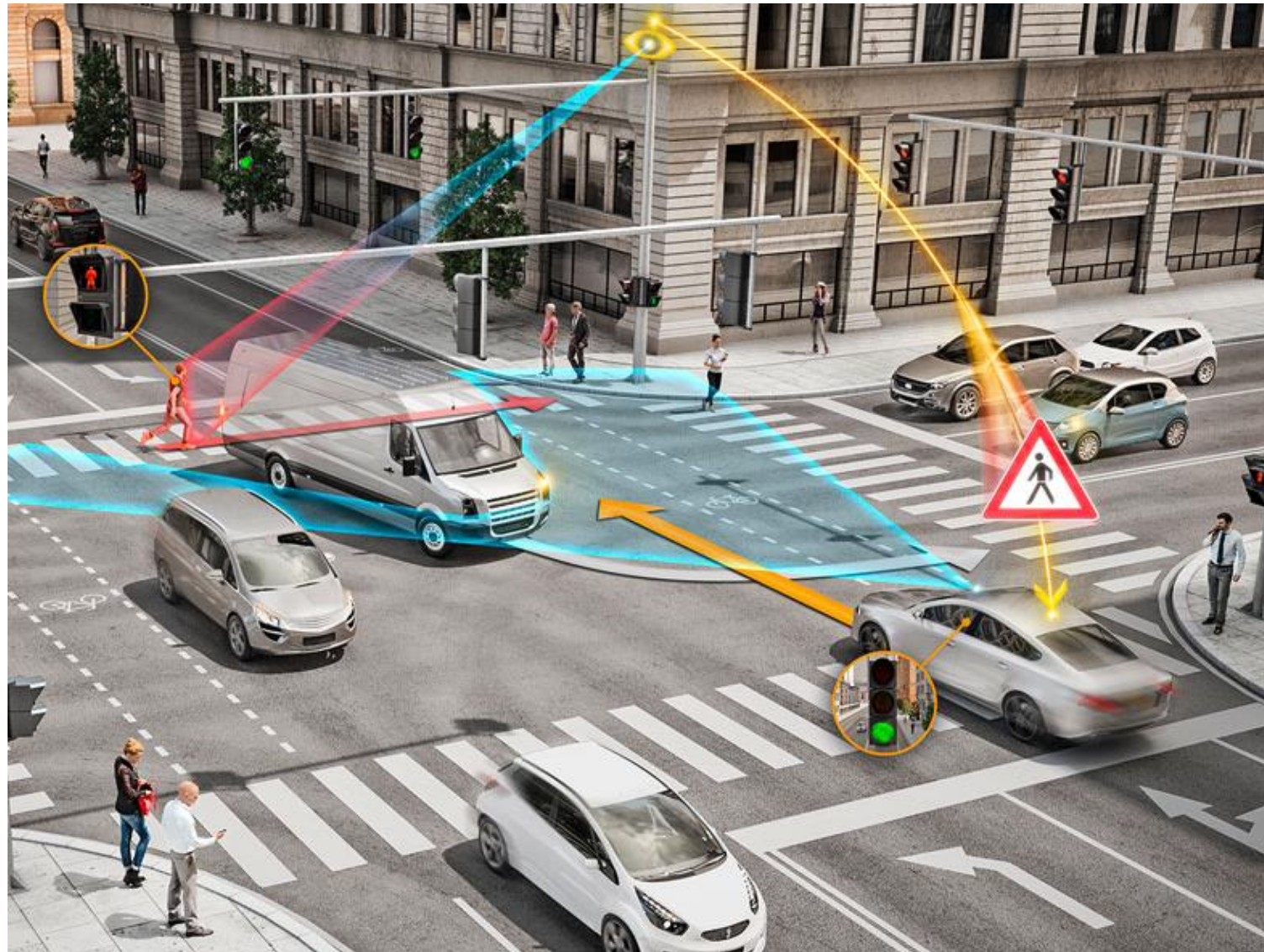# Traffic condition in Hong Kong

Very **complex** road conditions

Very **high** traffic densities
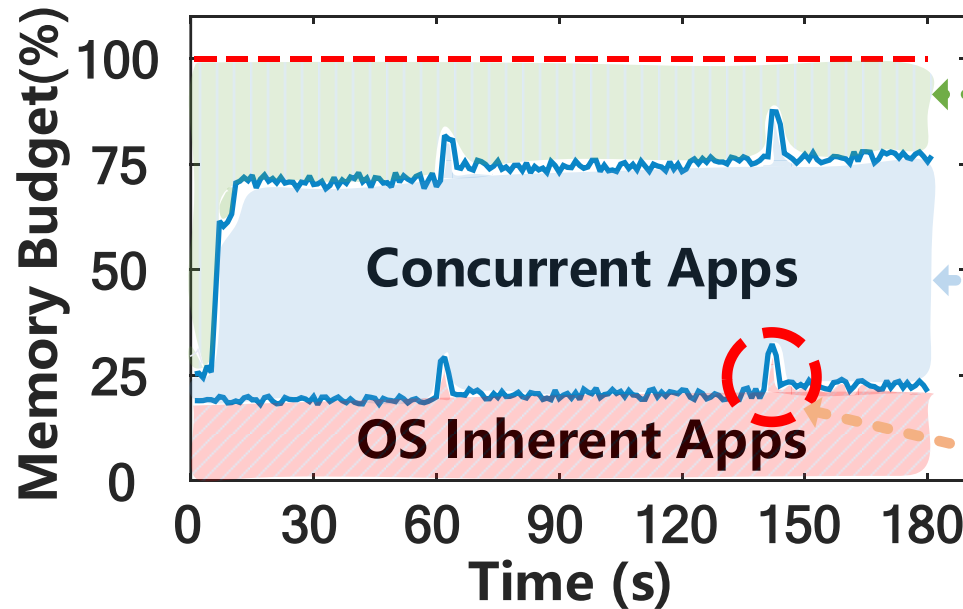
# Vehicle-to-Everything (V2X)

# Roadside Unit - Edge AI Device



NVIDIA Jetson Board

HUAWEI Atlas Board

Memory Budget(%)

Concurrent Apps

OS Inherent Apps

Time (s)

Memory Budget for AI Apps

Real-time Memory Budget Fluctuation

Available Memory for AI Apps

Concurrent Apps Memory footprint

OS Inherent Apps

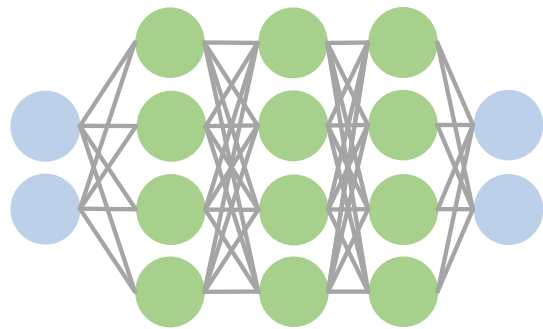Only Around **25% memory** remain in NVIDIA Jetson Nano

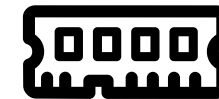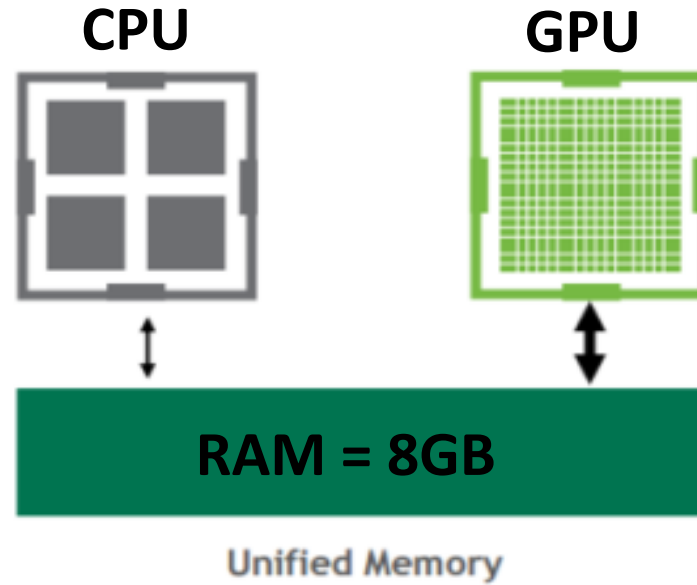# Challenge #1 - Memory Scarcity Problem

**CPU**

**GPU**

**RAM = 8GB**

Unified Memory

**AI apps need to compete memory with other apps**

**Model Size: > 2GB**

\>

**Available RAM: < 2GB**

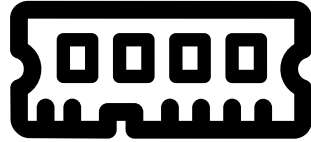**Q** How to run large model in small memory?

# Existing Methods

① **Upgrade Memory**

**8GB RAM**

**$200**

**32GB RAM**

② **Model Compression**

before pruning

after pruning

pruning synapses

pruning neurons

**Accuracy Drop**

**90.33** **90.25** **89.97** **89.69** **89.12** **88.95** **85.30** **79.32** **71.27**

**Over Compression**

Accuracy(%)

95 85 75 65

Parameters(M)

138 100 90 85 80 75 70 65 60
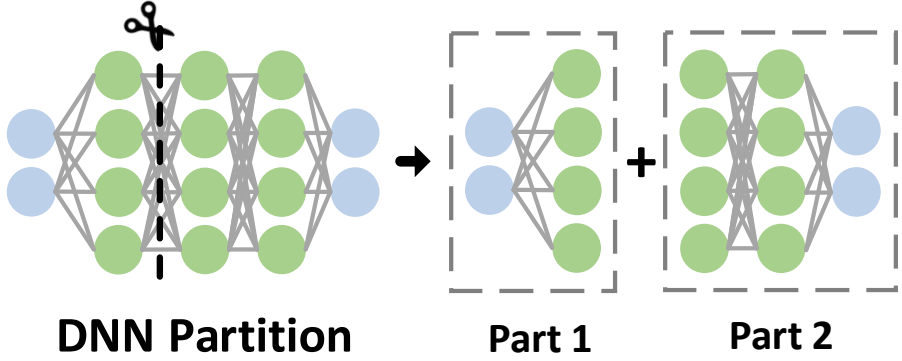
# Existing Methods

③ **Computing Offloading**

**DNN Partition**

**Part 1** + **Part 2**

**Edge Side**

**Transmit**

**5G**

**Server Side** → **man**

**Attacks**

**Hacker**



Latency(ms) vs Bandwidth(MB/s)

- 289.35
- 466.37
- 589.48
- 784.84
- 1503.97
- 2892.37

Offloading

Edge only

# Thinking



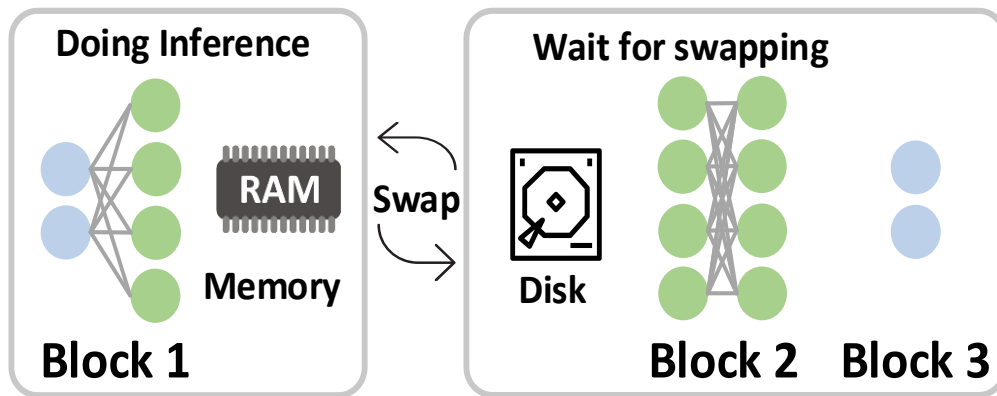DNN Partition → Part 1 + Part 2
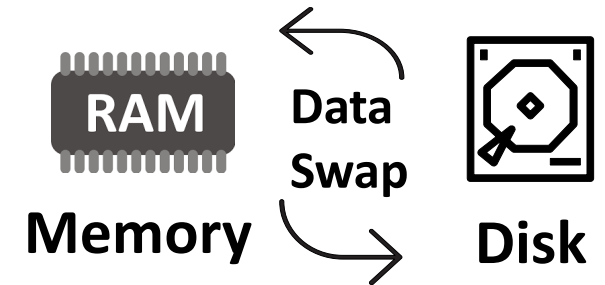
Keep Accuracy

Edge Side

5G

Transmit

Server Side

→ man
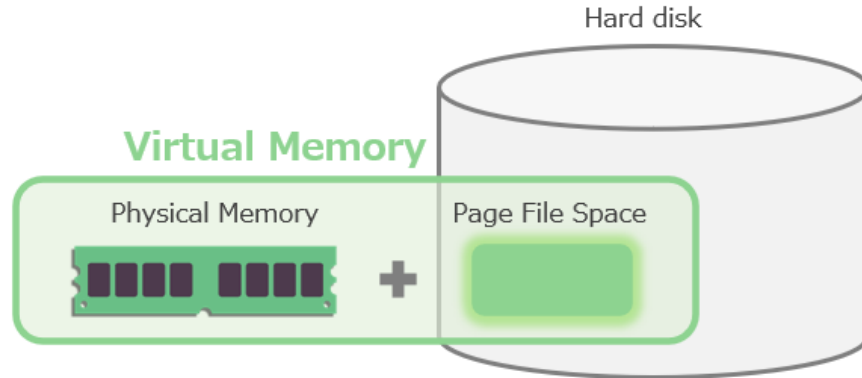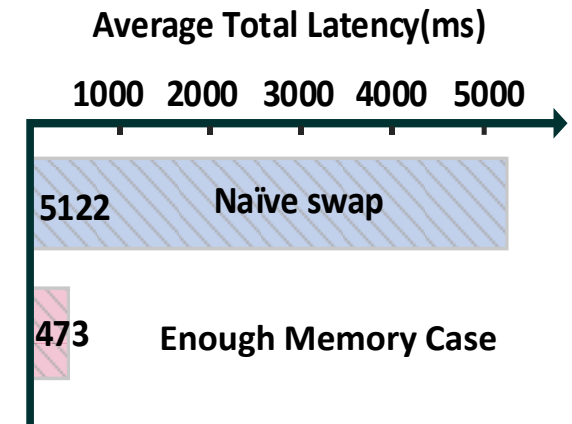
① Transmission safety concern

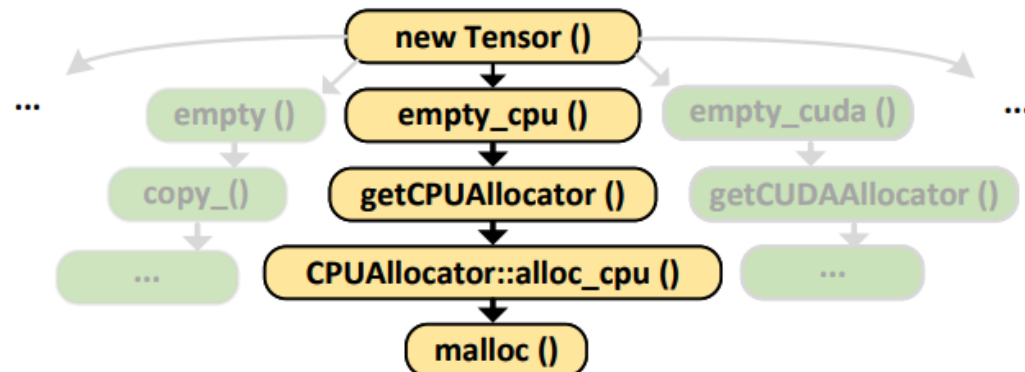② Transmission rate concern

# Main Idea - Virtual Memory



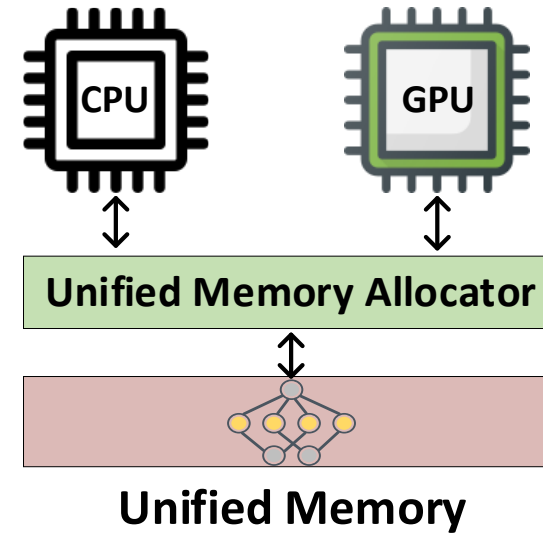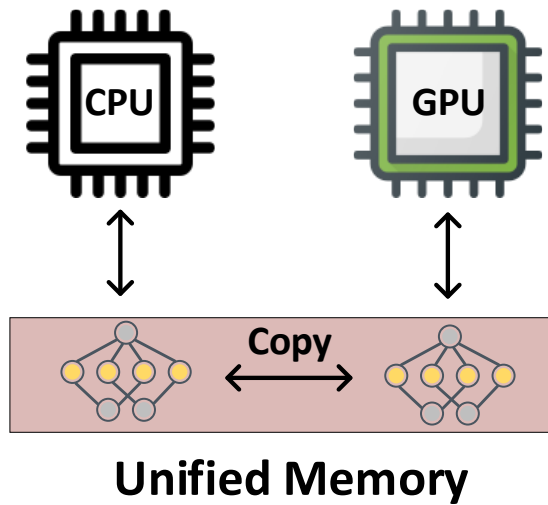**Model Block Swap**



**Big latency**

# Design #1 - Unified Memory Allocator



```
1:  // In Copy.cu
2:  // data_ptr pointed to existing CPU tensor.
3:      void* src = iter.data_ptr(1);
4:  // Original method needs to allocate GPU Memory
5:      and copy data to it.
6:  // void* dst = iter.data_ptr(0);
7:  // cudaMemcpyAsync(dst, src, size, kind, stream);
8:      void* dst = src;
9:      cudaDeviceSynchronize();
10:     return dst;
```

# Design #2 - Weights restoration optimization

# Challenge #2 - Inefficiency of Sequential Swap



**(a)**

**(b)**



**Unavoidable Swap latency of Sequential Inference**

# Design #3 - Partition Module: Parallel Inference



Case 1    $T_{SO1}+T_{SI3} \leqslant T_{I2}$

Case 2    $T_{SO1}+T_{SI3} > T_{I2}$

Total Latency

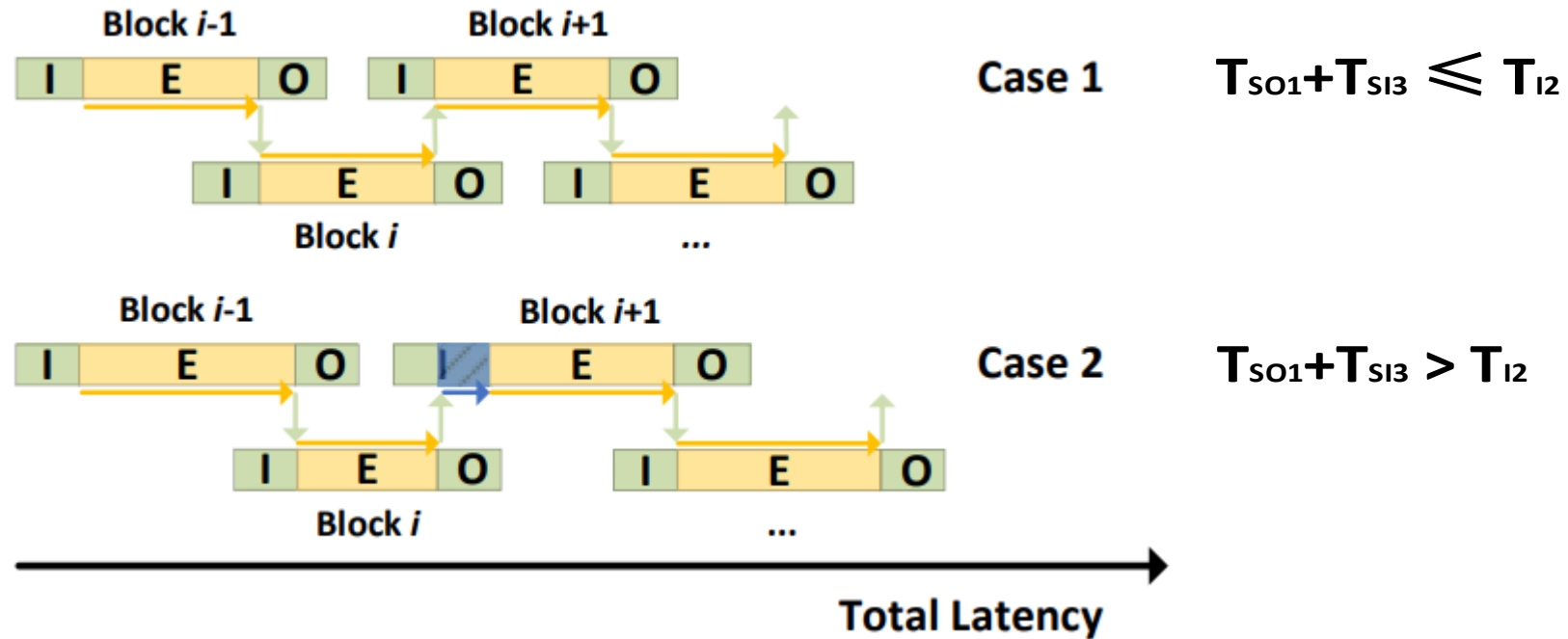$$T_{overlap}(i) = T_{swap\_out}(i-1) + T_{swap\_in}(i+1) - T_{inference}(i)$$

# Design #3 - Partition Module: Select Optimal Solution



$$\min \sum_{n=1}^{B} (\alpha * WS_n + (\beta + \gamma) * PD_n - \theta * FLOPs_n)$$

$T_{I/O} \propto Weight\_Size$

$T_{restore} \propto param\_depth$

$T_{inference} \propto FLOPs$

$T_{remove} \propto param\_depth$

# Design #3 - Partition Module: Select Optimal Solution



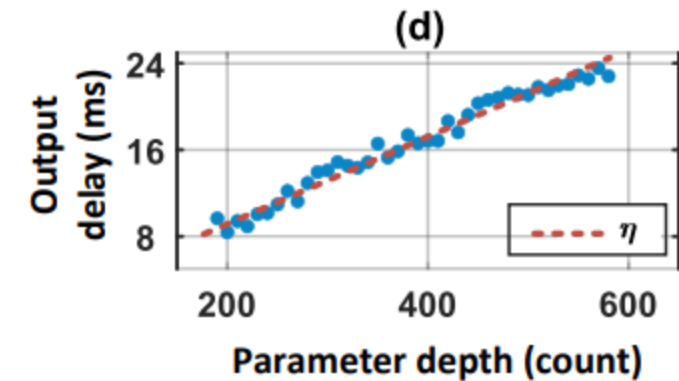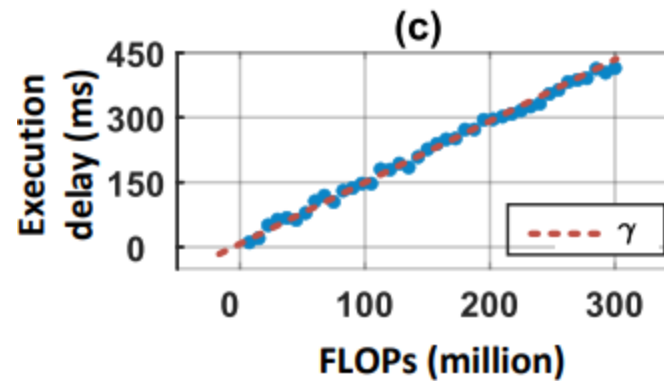$$\min \sum_{n=1}^{B} (\alpha * WS_n + (\beta + \gamma) * PD_n - \theta * FLOPs_n)$$
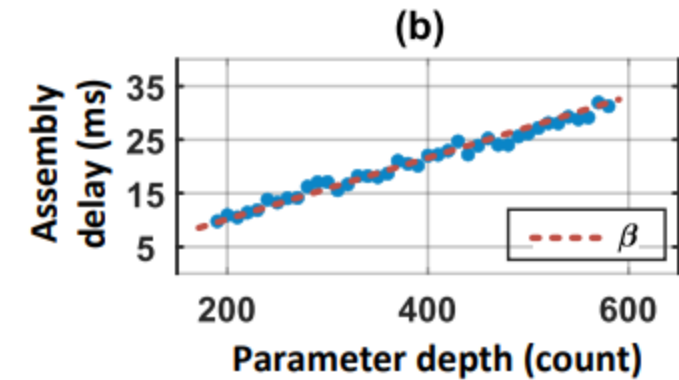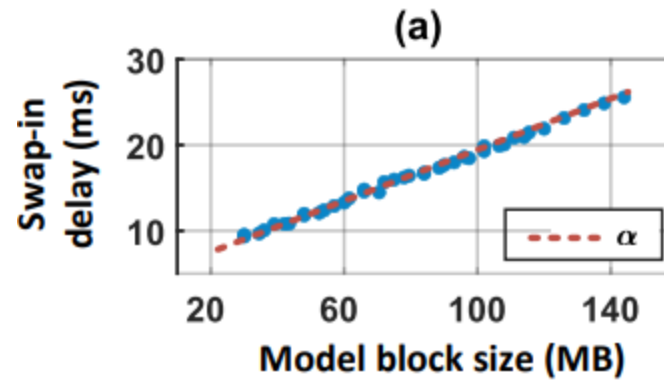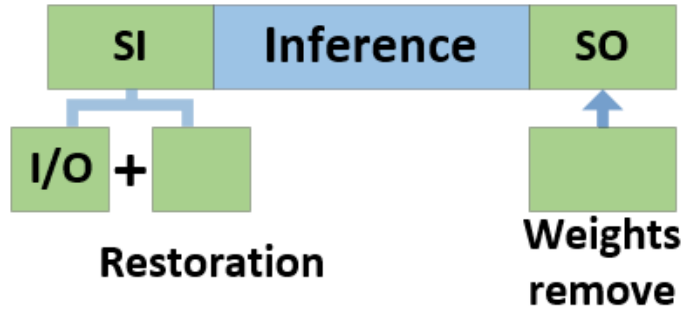
$T_{I/o} \propto Weight\_Size$

$T_{restore} \propto param\_depth$

$T_{inference} \propto FLOPs$

$T_{remove} \propto param\_depth$

| Layer | Size | Depth | FLOPs |
|---|---|---|---|
| Layer1 | 0.38 MB | 1 | 26.2 M |
| Layer2 | 1.49 MB | 5 | 0.8 K |
| Layer3 | 1.12 MB | 1 | 123.9 M |
| Layer4 | 5.93MB | 5 | 4.2 K |
| Layer5 | 4.38MB | 6 | 316.7 M |
| ... | ... | ... | ... |
| Layer100 | 23.6 MB | 1 | 30 K |
| Layer101 | 17.45 MB | 1 | 5 K |

**Model Info Table**

| Partition Points | Maximum Memory | Predicted Latency |
|---|---|---|
| 1,2 | exceed | null |
| 1,3 | exceed | null |
| ... | ... | ... |
| 30,66 | 105 MB | 496 ms |
| 30,67 | 109 MB | 488 ms |
| ... | ... | ... |
| 98,100 | exceed | null |
| 99,100 | exceed | null |

**Decision Table**

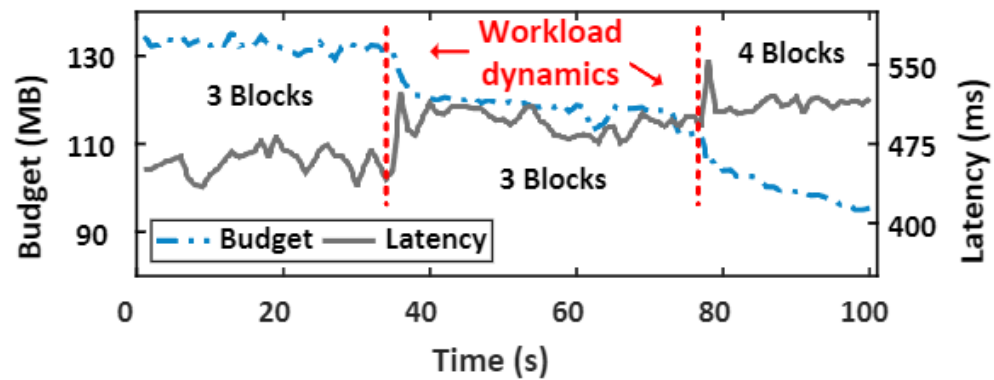**Overlap latency can be computed through the model info table**
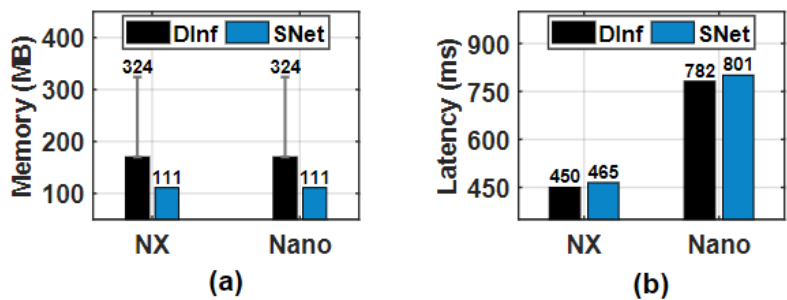
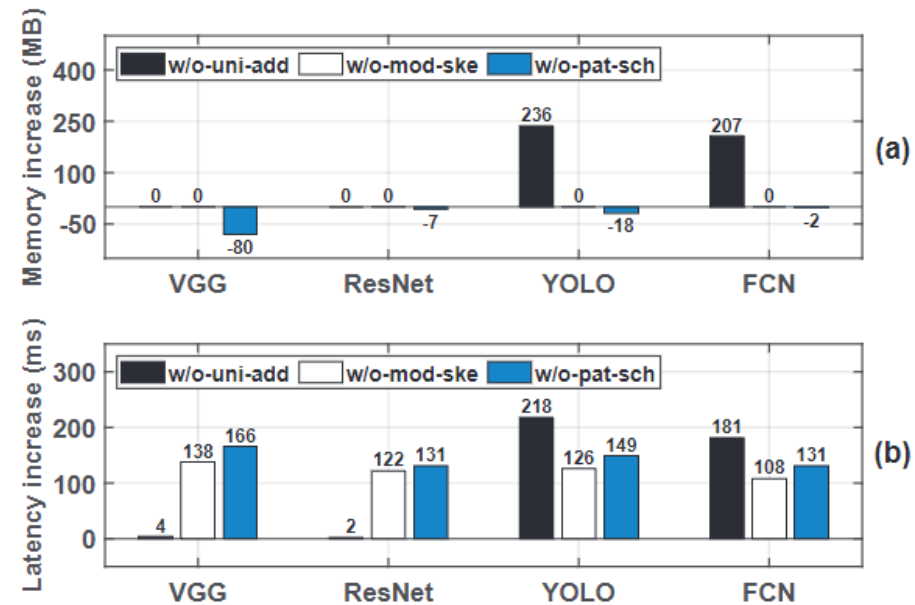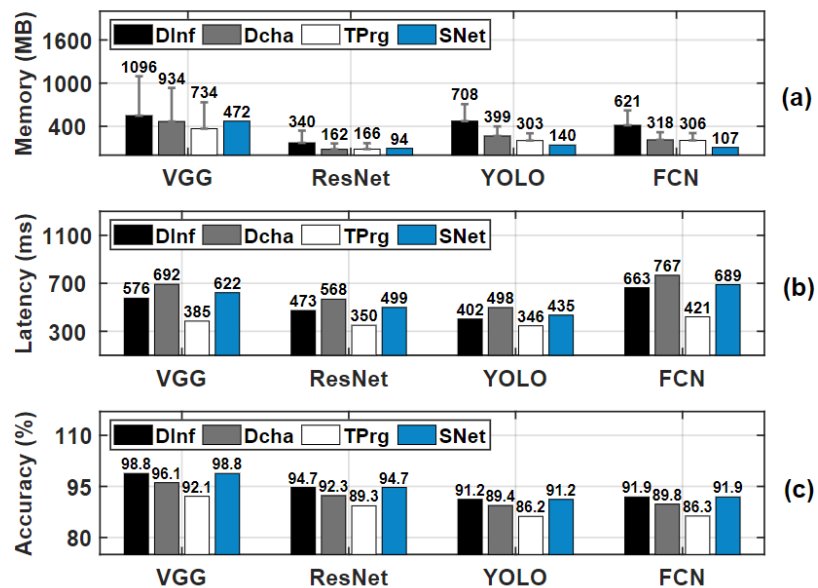# Implementation



**Proposed SwapNet Framework**



**Test Bed**

1. **Self-driving(4 tasks):** lane detection, object detection, segmentation, traffic sign classification
2. **Road-Side Unit (5 tasks):** 2 object detection, 2 natural scenes classification and traffic light classification
3. **UAV surveillance(2 tasks):** fire source detection, wild animal recognition

**Scenarios**

# Evaluation

# Conclusion

- We introduce SwapNet, a middleware that logically executes large DNN models on a small memory budget. SwapNet partitions large DNN models into blocks for execution by swapping them between the memory and the external storage in order.

- Our main contribution is a transparent design that eliminates the substantial latency and memory overhead occurred during block swapping while remaining compatible with the DNN development tool chains for edge AI devices.

- Extensive evaluations show the promising performance gains of SwapNet in combination with parallel optimization for efficient execution.