

Confidential Execution of Deep Learning Inference at the Untrusted Edge with ARM TrustZone

Md Shihabul Islam, Mahmoud Zamani, Chung Hwan Kim, Latifur Khan, Kevin W. Hamlen

**ACM Conference on Data and Application Security and Privacy (CODASPY '23),
April 24–26, 2023, Charlotte, NC, USA**

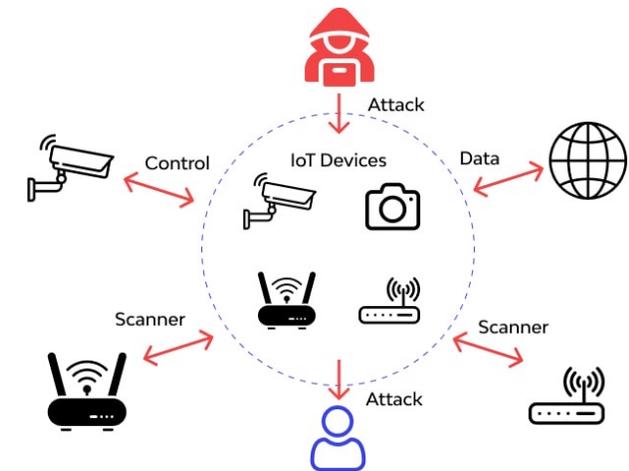
MOTIVATION

- ❑ Ubiquitousness of Internet-of-Things (IoT) devices
- ❑ On-device Machine Learning
 - Performance of edge/IoT applications
 - Low bandwidth
 - Reducing communication cost
 - Privacy of user data



CHALLENGES

- Protection of **user data** on untrusted and resource-constrained edge/IoT devices
 - Model Inversion Attack
 - Membership Inference Attack
- Unfeasible existing techniques for edge/IoT devices
 - Homomorphic encryption
 - Differential privacy



Solution:

- ✓ Trusted Execution Environment (TEE) for edge/IoT devices
 - ARM TrustZone

ARM TrustZone

❖ ARM: Pioneer in embedded device processors

❖ TrustZone

- Optional hardware security extension
- Ensures the integrity and confidentiality of an application's data on a device
- Two architectures:
 - Cortex-A
 - Cortex-M

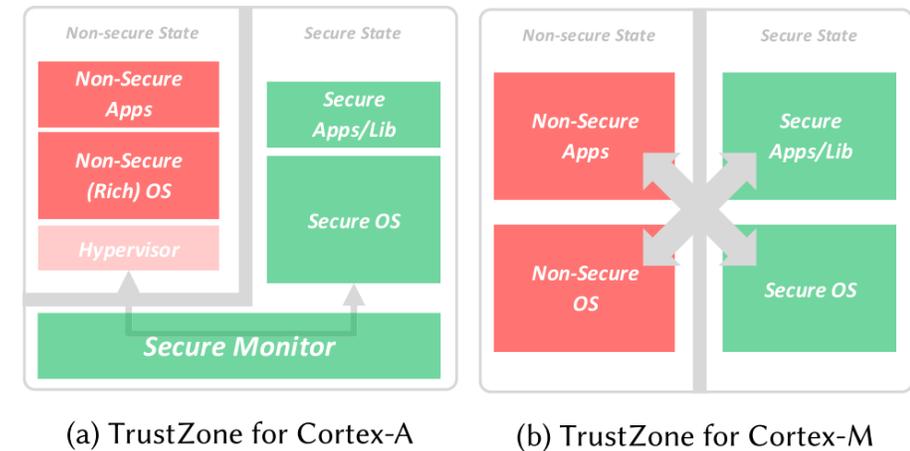


Fig. 1. TrustZone technology.

ARM TrustZone Limitations

Limitations:

- Resource-intensive DL methods
- Limited trusted memory and resources in TrustZone

Possible Solutions:

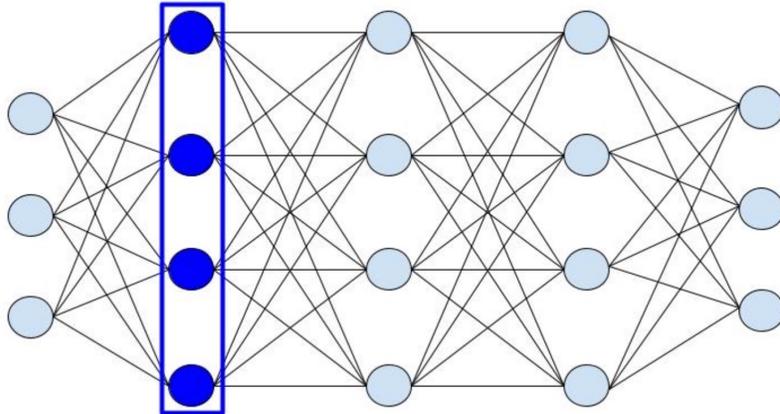
- Quantization
- Model pruning

But affects model's prediction accuracy



Common Practice: Partitioning

Layer-base Partitioning



Model	# Layers	Pre-trained Model Size (MB)	Peak Mem. Usage (MB)
LeNet	10	0.2	7
VGG-7	13	0.3	7
CIFAR	18	30.7	45
Tiny	22	4.2	71
Darknet	16	29.3	88
Extraction	27	93.8	163
Alexnet	14	249.5	272
Darknet53	78	159	273
Inception-v3	145	95.5	448
Yolov3	107	237	840
VGG-16	24	528	923

Too
Large!

Typical Trusted Memory \approx 16 MB

❑ How to solve?

- ❑ Run only a few layers in the TrustZone

- Model Inversion Attack
- Membership Inference Attack

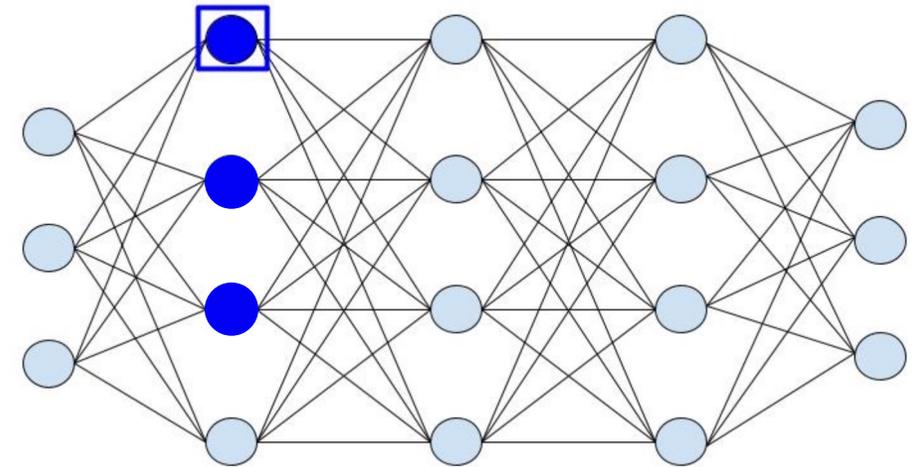
T-Slices

Overview:

- Utilizes ARM TrustZone with limited trusted memory to protect the entire DL execution
- Does not sacrifice original prediction accuracy

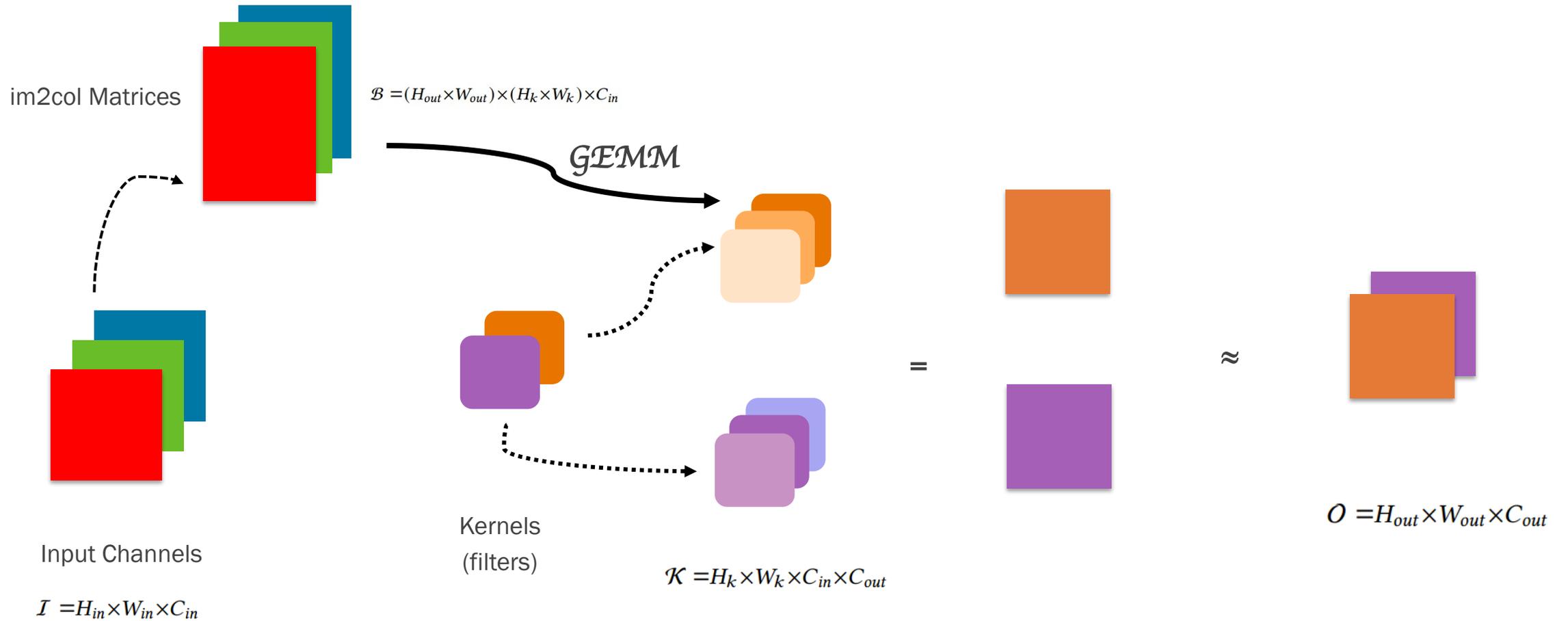
T-Slices

- Partitions DNN layer into smaller independent segments called **Slices**
- Follows an optimized Memory Management plan with on-demand parameter loading scheme
 - Calculated from Hyperparameters
- Dynamically determines a set of **Slices** based on the available trusted memory buffer in TrustZone

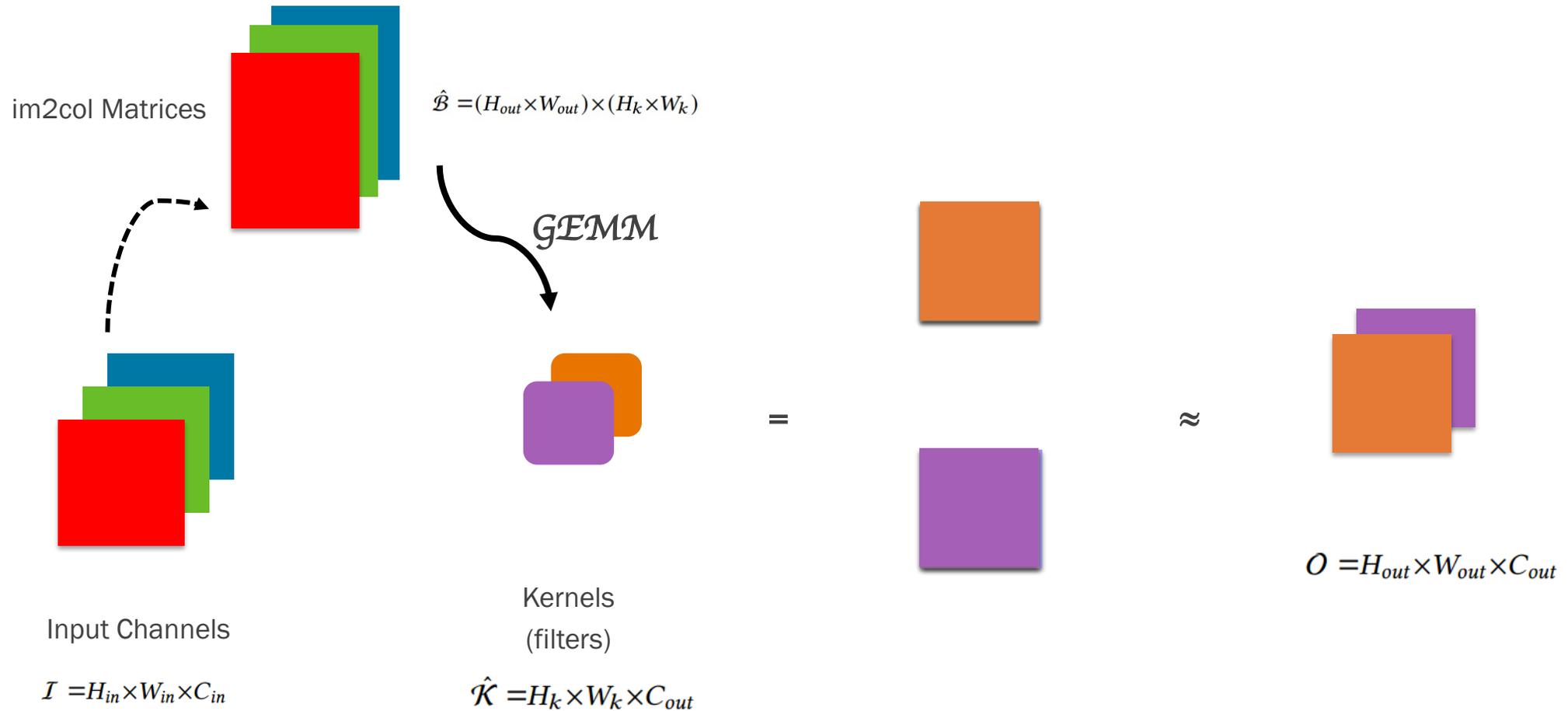


Background

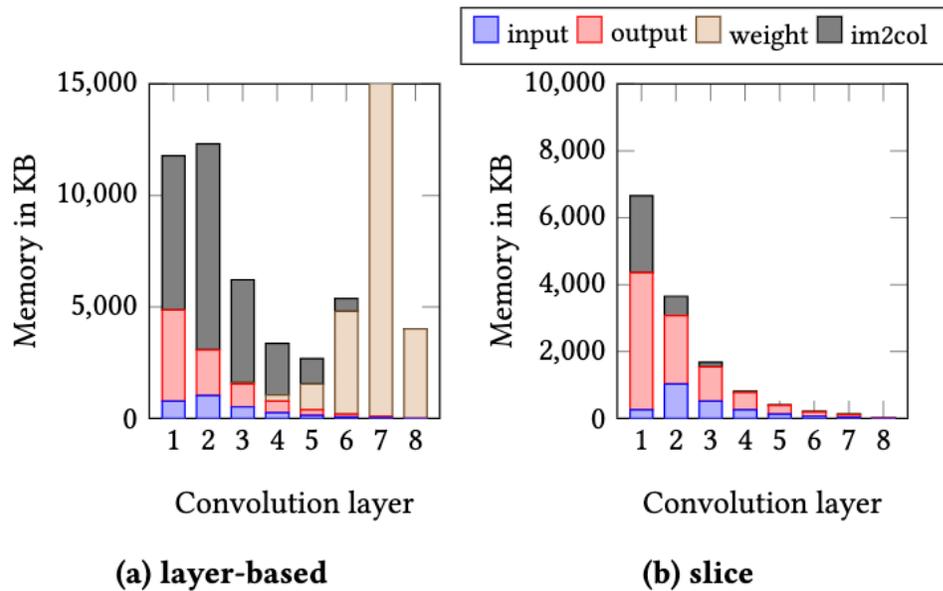
Convolution Operation



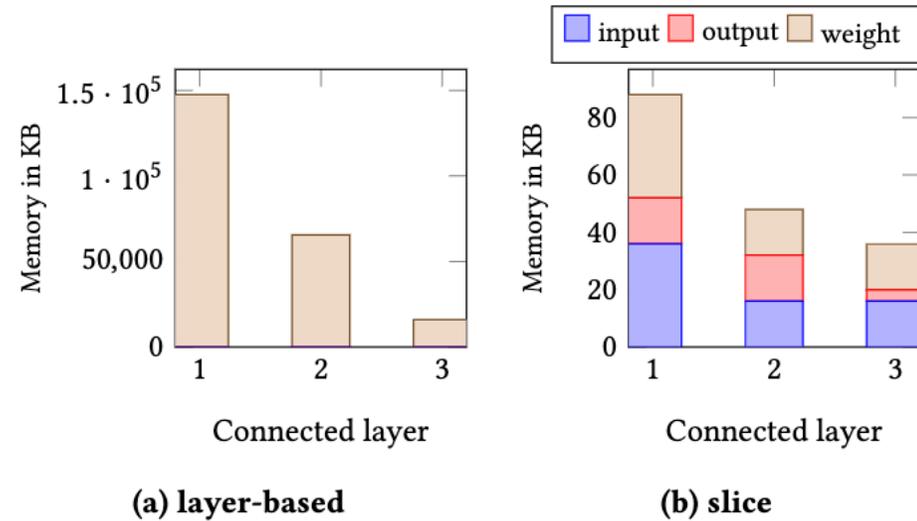
Slicing for Convolution Operation



Memory Buffer Size Comparison

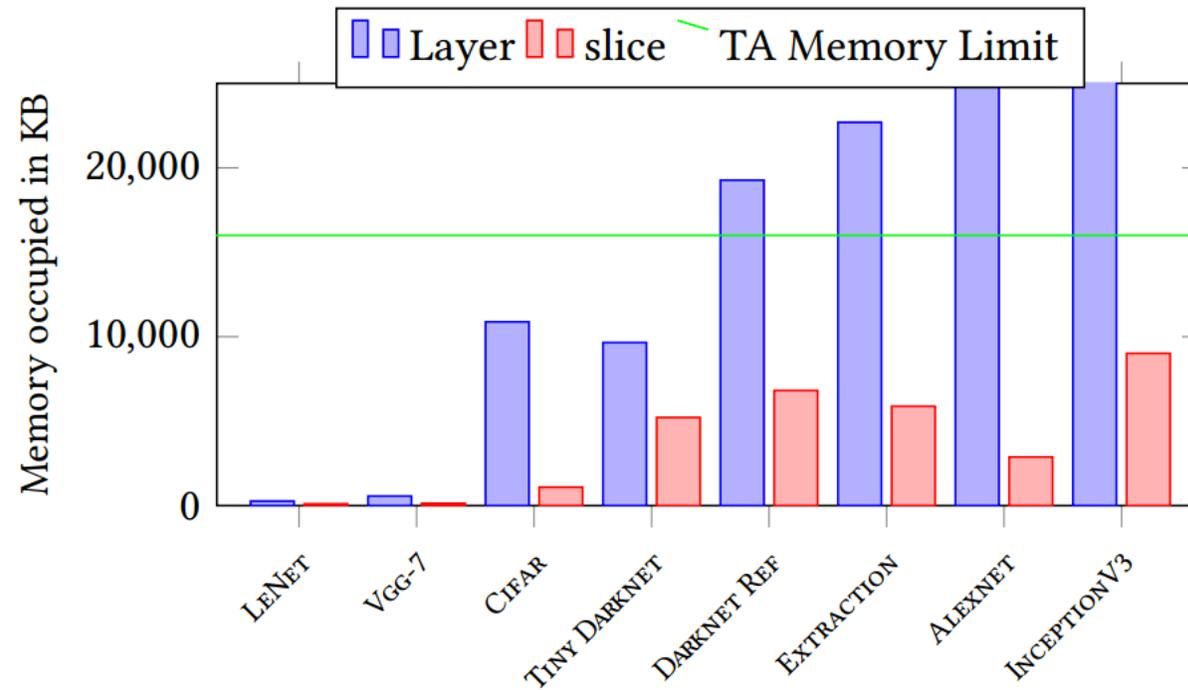


Darknet Reference Model



Alexnet Model

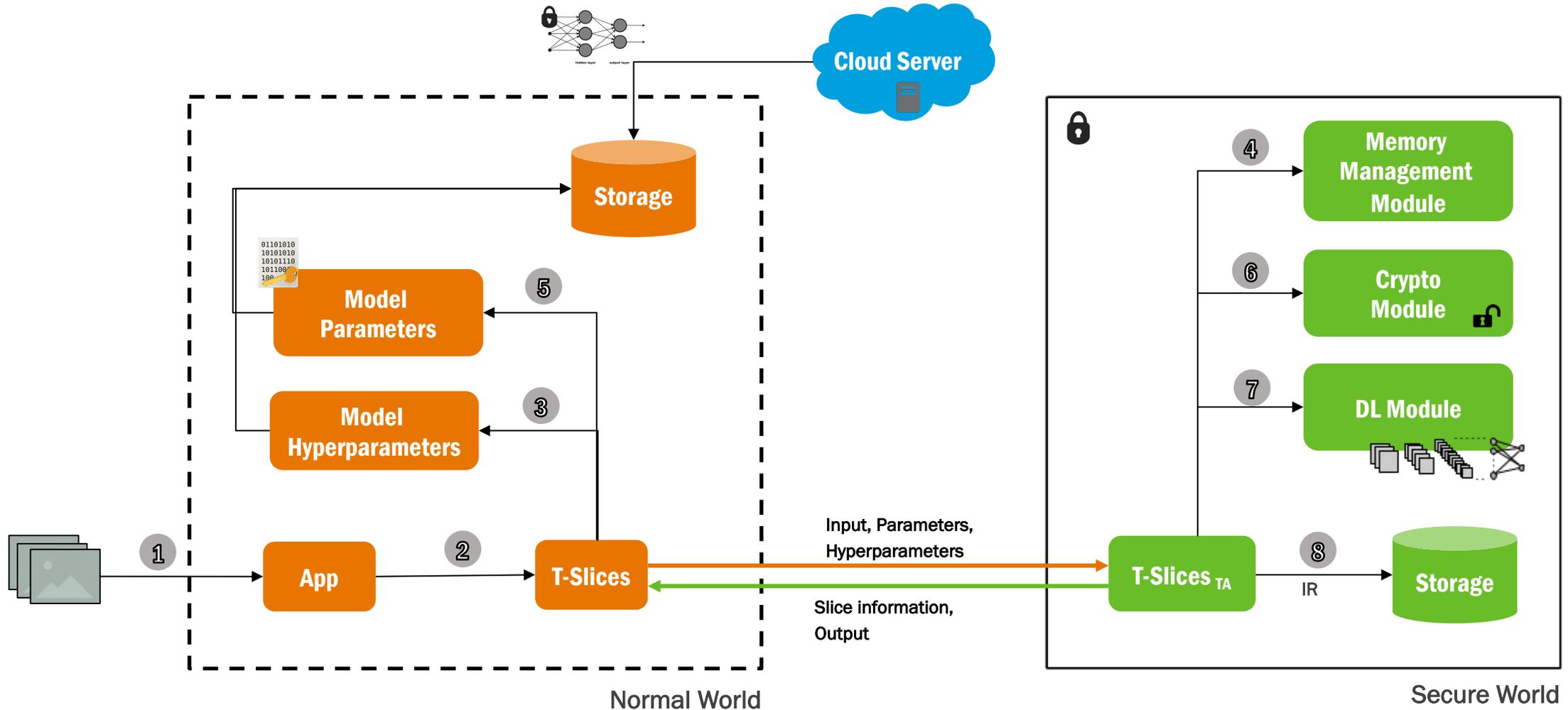
Memory Buffer Size Comparison



Peak memory required to execute any convolution/connected layer in different CNN architectures. Trusted memory limit considered as 16 MB.

Our Contribution

T-Slices Architecture/Flow



Experimental Setting

➤ Device Configuration

- STM32MP157C-DK2 with Cortex-A7 32-bit and Cortex-M4 32-bit MPUs
- Raspberry Pi 3 Model B (RPi3B)

➤ Experiment

- Image classification with CNN models
- Compare with Baseline DarknetZ^γ

➤ Performance Metric

- Trusted Memory Consumption
- Prediction Time Overhead
- Case Studies against prevalent privacy attacks

Dataset and Models

Model	# Layers	# Conv. Layers	Dataset	Pre-trained Model Size (MB)
LeNET	10	2	MNIST	0.2
CIFAR_SMALL	12	7	CIFAR10	0.08
VGG-7	13	6	CIFAR10	0.26
VGG-7	13	6	CIFAR100	0.3
CIFAR	18	10	CIFAR10	30.7
TINY DARKNET	22	16	ImageNet1k	4.2
EXTRACTION	27	21	ImageNet1k	93.8
DARKNET REF	16	8	ImageNet1k	29.3
ALEXNET	14	5	ImageNet1k	249.5
INCEPTIONV3	145	94	ImageNet1k	95.5

Trusted Memory Consumption

- T-Slices on average achieves 72% reduction in peak memory consumption

Model	DARKNETZ per Layer	DARKNETZ* per Layer	T-SLICES per Slice	% Decrease[†]
LENET	7	0.25	0.1	60
VGG-7	7	0.7	0.2	71
CIFAR	45	10.5	1.25	88
TINY DARKNET	71	9.5	5	47
DARKNET REF	88	18.5	6.5	65
EXTRACTION	163	22.6	5.6	75
ALEXNET	272	144	2.75	98
INCEPTIONV3	337	33	9	73

* with *on-demand* parameter loading scheme

† decrease from DARKNETZ* to T-SLICES

Prediction Time Overhead

- T-Slices on average achieves 29% improvement in execution time

STM32MP157C-DK2

CNN	Dataset	DARKNETZ*	T-SLICES	% Improvement
LeNET	MNIST	2.44	2.10	14
CIFAR_SMALL	CIFAR10	3.49	3.24	7
VGG-7	CIFAR10	11.93	6.38	47
CIFAR	CIFAR10	608.04	285.07	53
TINY DARKNET	ImageNet1k	874.58	859.34	2
EXTRACTION	ImageNet1k	1244.84	615.56	51
DARKNET REF	ImageNet1k	1175.69	815.55	31
ALEXNET	ImageNet1k	✗	1219.31	✗
INCEPTIONV3	ImageNet1k	✗	1928.41	✗

*with *on-demand* parameter loading

✗: Unable to execute due to not enough trusted memory

RPI3B

CNN	Dataset	DARKNETZ*	T-SLICES	% Improvement
LeNET	MNIST	0.092	0.092	0
CIFAR_SMALL	CIFAR10	0.19	0.19	0
VGG-7	CIFAR10	0.309	0.307	1
CIFAR	CIFAR10	30.43	30.26	1
TINY DARKNET	ImageNet1k	14.72	14.71	0
EXTRACTION	ImageNet1k	116.57	116.24	0
DARKNET REF	ImageNet1k	18.81	18.78	0
ALEXNET	ImageNet1k	✗	44.19	✗
INCEPTIONV3	ImageNet1k	✗	468.1	✗

*with *on-demand* parameter loading

✗: Unable to execute due to not enough trusted memory

Security Analysis

- **Model Inversion Attack** [1]
 - Reconstruct/recover the training data or any sensitive attributes from the trained ML model
- **Membership Inference Attack** [2]
 - Discover whether a given data sample is a part of the training dataset for the trained ML model

[1] Model Inversion Attacks That Exploit Confidence Information and Basic Countermeasures, ACM CCS 2015

[2] Membership inference attacks against machine learning models, IEEE S&P 2017

Limitations & Future Work

- ❑ Investigate vast DL models unsuitable for memory-constrained edge/IoT devices
 - ❑ Peak memory of vgg-16 ~ 923 MB, Yolov3 ~ 840 MB
 - ❑ Parallel processing using multiple TZ devices
- ❑ Investigate other DL architectures (RNNs)
- ❑ Investigate the capability of side-channel attacks on T-Slices



Thank you

Contact information

md.shihabul.islam@utdallas.edu