

# RNNPool : Efficient Non-linear Pooling for RAM Constrained Inference

Code: <https://github.com/Microsoft/EdgeML>



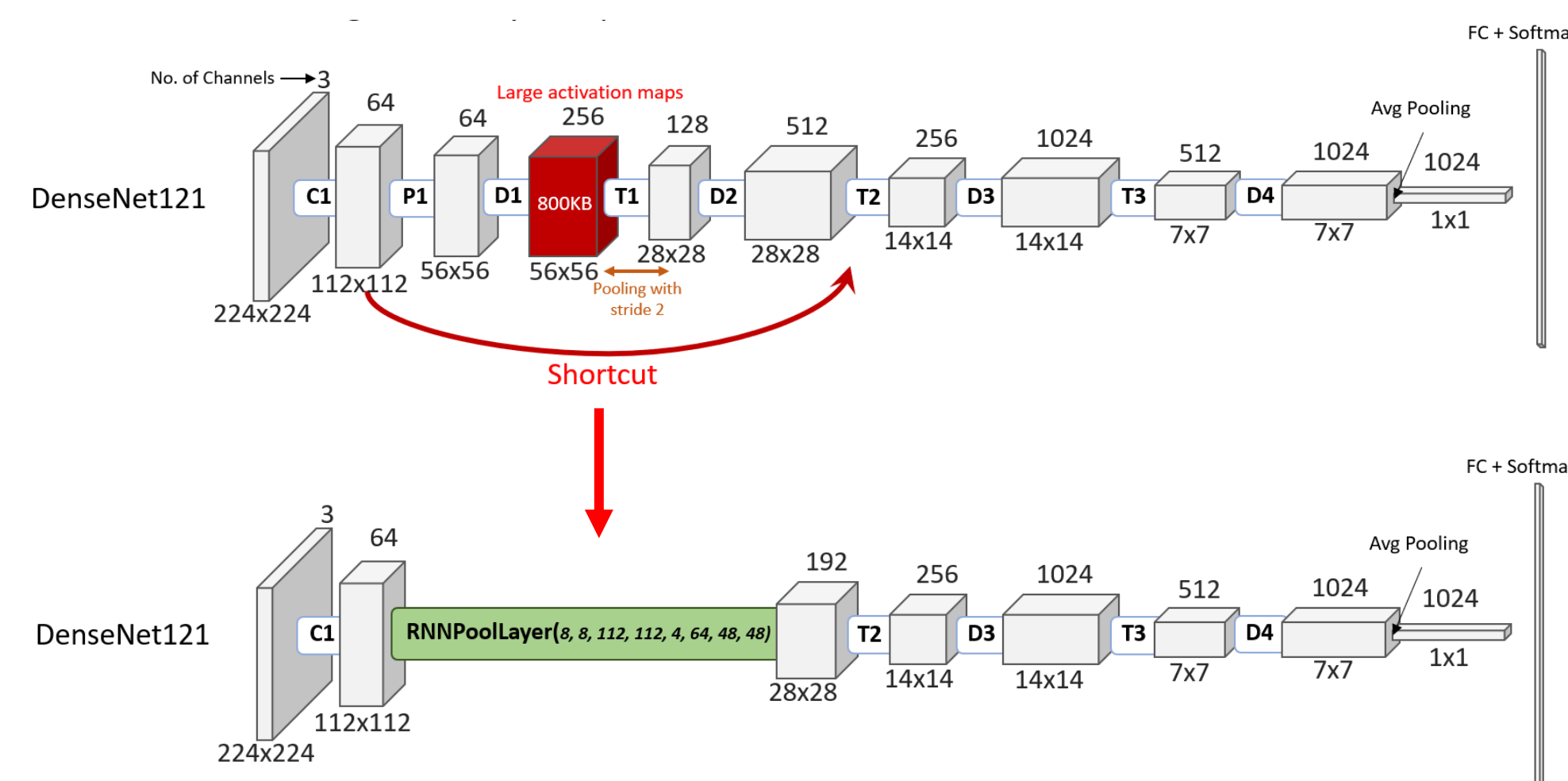
Oindrila Saha<sup>1</sup>, Aditya Kusupati<sup>2</sup>, Harsha Vardhan Simhadri<sup>1</sup>, Manik Varma<sup>1</sup> and Prateek Jain<sup>1</sup>

<sup>1</sup>Microsoft Research

<sup>2</sup>University of Washington

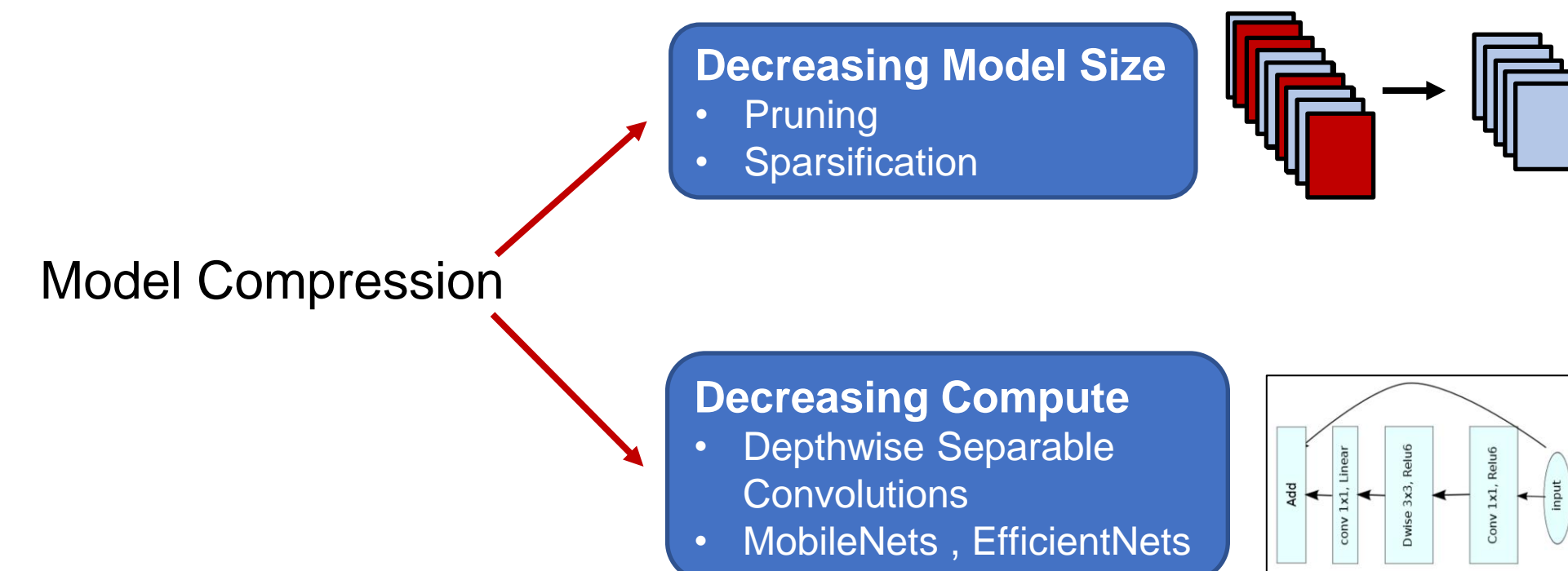
**Our goal:** Accurate computer vision models that can be deployed on tiny devices

- Barriers:**
- CNN models have many layers with large activations
  - Large memory footprint, the most constrained resource on microcontrollers



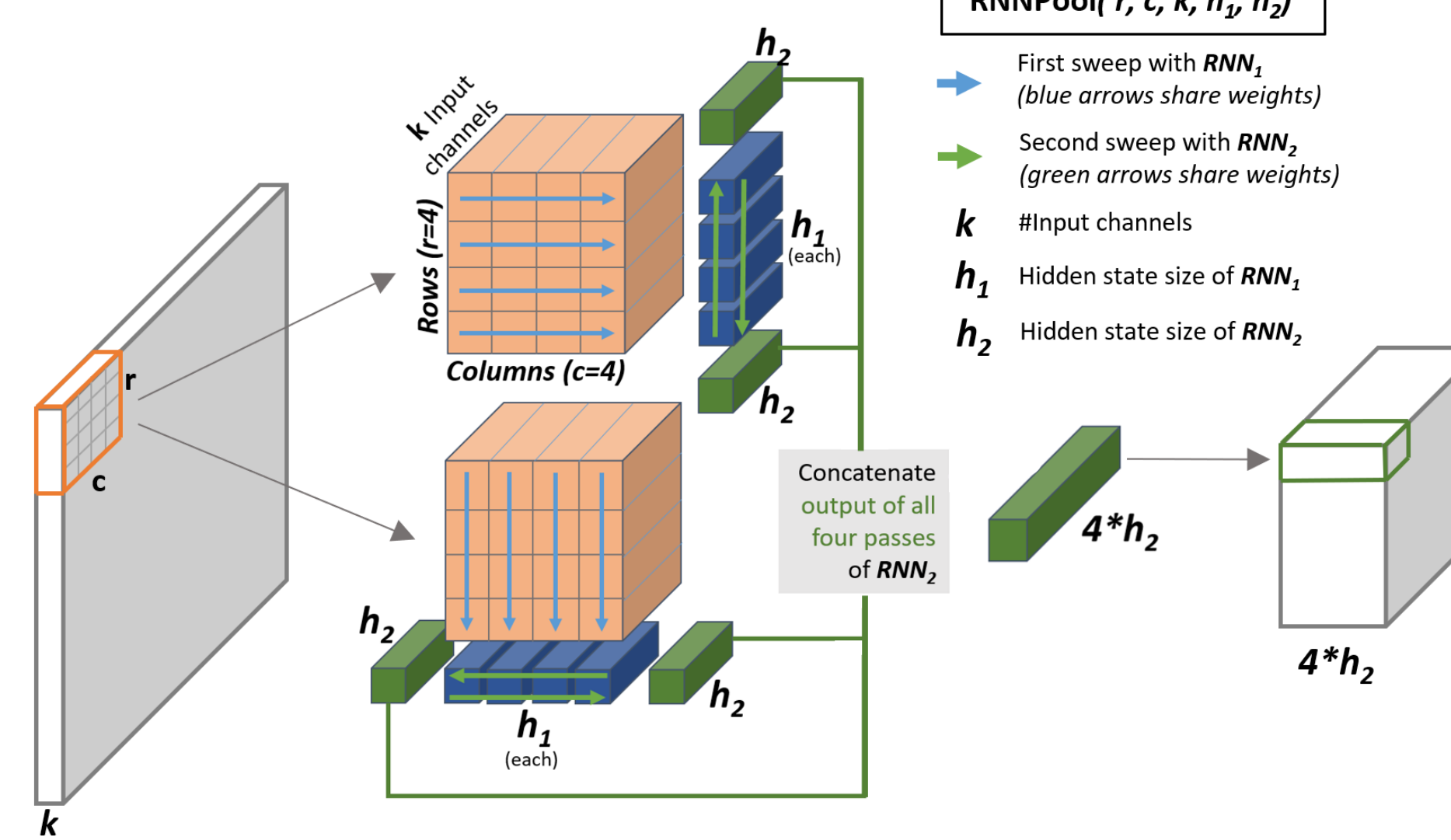
- Standard pooling operators (e.g. max pool) are gross aggregators, thus are only used with a maximum stride of 2
- RNNPool can reduce intermediate feature maps significantly (up to 16x) with small loss in accuracy
  - Ensures heavy convolution blocks run on smaller activation maps

## Existing Works



However, peak RAM requirement still remains high

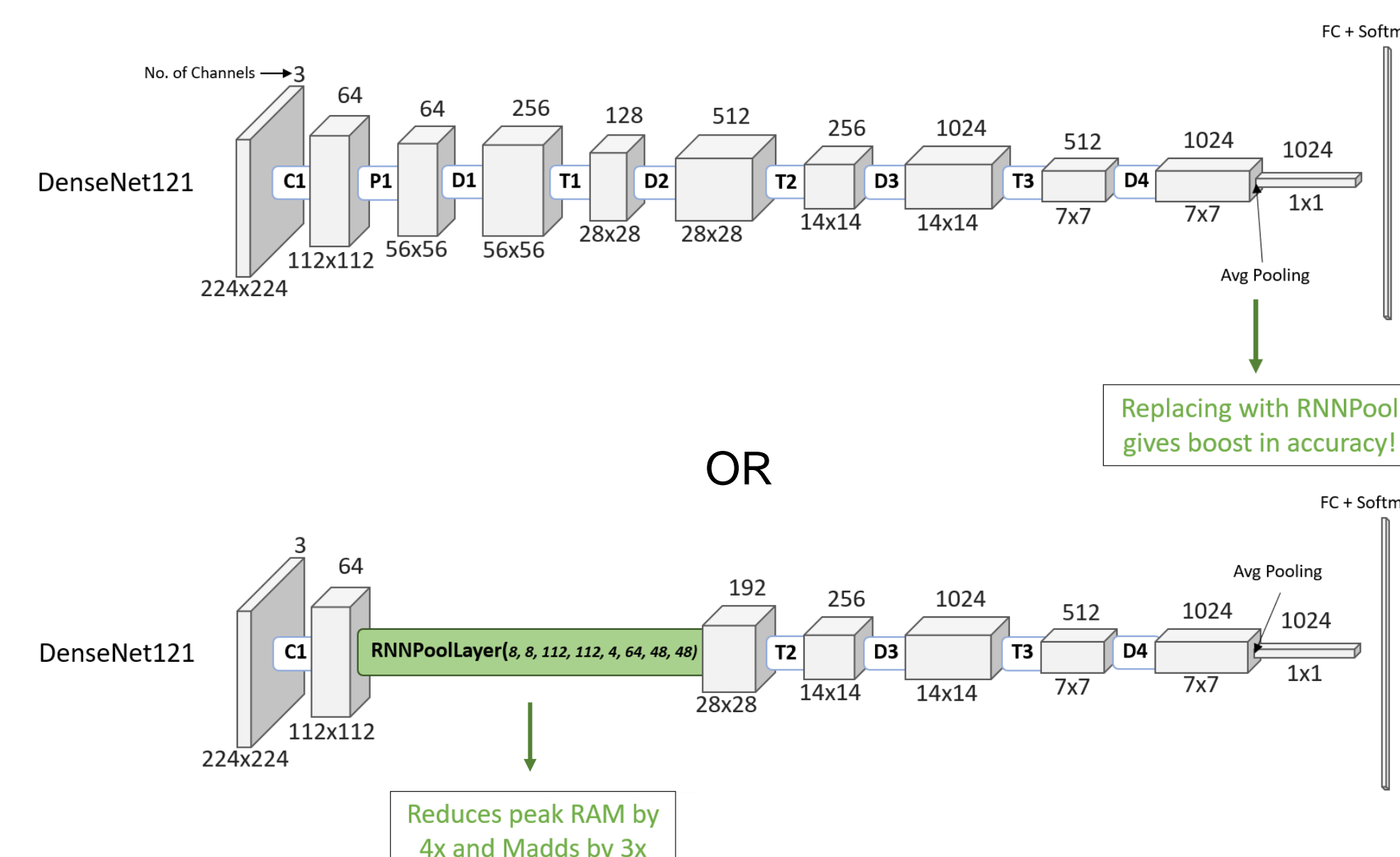
## RNNPool



- Takes a patch of the input and produces a 1x1 summary
- Patches having overlap of size = PatchSize - Stride
- For each patch, 4 RNN runs produce the pooled feature vector

## Usage

- Semantically equivalent to pooling, so can be used to replace any pooling operator
- But key usage in reducing image size in beginning of network to save RAM requirement



## Evaluation

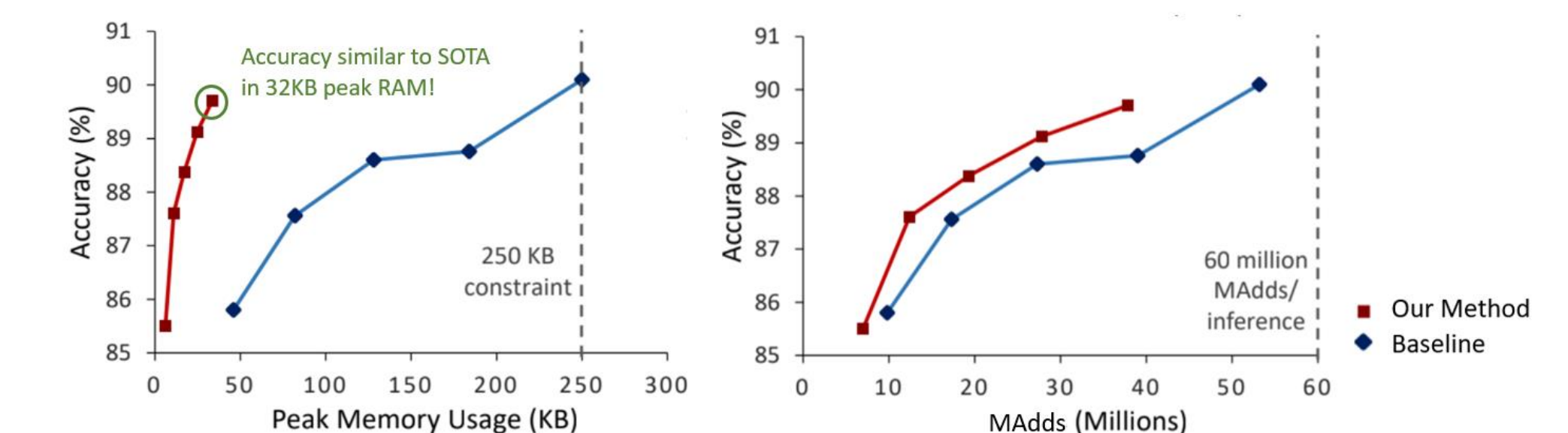
### Task 1: Image Classification

Results on ImageNet-10 dataset formed by subsampling 10 classes from ImageNet1K and MobileNetV2 as the base architecture

Model	Accuracy	MAdds	#Params	RAM
Base Network	94.2	0.300G	2.2M	2.29MB
Average Pooling	90.8	0.334G	2.0M	0.24MB (10x Lower RAM!)
Max Pooling	92.8	0.200G	2.0M	
Strided Conv	93.0	0.200G	2.1M	
ReNet	92.2	0.296G	2.3M	
<b>RNNPool</b>	<b>94.4</b>	0.226G	<b>2.0M</b>	

### Task 2: Visual Wakeword

- 8x less RAM
- 40% less compute



### Task 3: Face Detection

Comparison with SOTA for very low MAdds category on WIDER FACE validation subset

Method	RAM	#Params	MAdds	MAP		
				E	M	H
EagleEye	1.17MB	0.23M	0.1G	0.74	0.70	0.44
<b>Rpool-Face-Quant</b>	<b>225KB</b>	<b>0.07M</b>	0.1G	<b>0.80</b>	<b>0.78</b>	<b>0.53</b>

**Face Detection M4 Deployment**

- 188KB peak RAM
- 70M MAdds
- 160KB Model Size

10.45 sec/image on STM32F439-M4 device clocked at 168 MHz

SCUT Head Dataset: <https://github.com/HCILAB/SCUT-HEAD-Dataset-Release>

Image from <https://researchdesignlab.com/stm32-arm-cortex-m4-development-board-stm32f407vet6.html>