

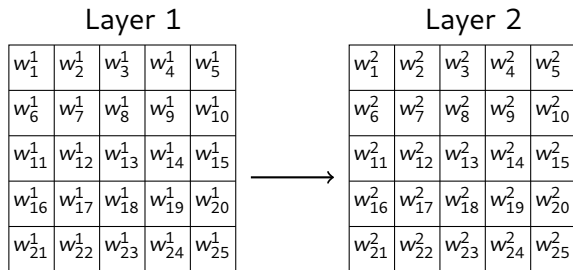
# Attention Based Pruning for Shift Networks

Ghouthi BOUKLI HACENE, Carlos Lassance, Vincent Gripon,  
Matthieu Courbariaux and Yoshua Bengio



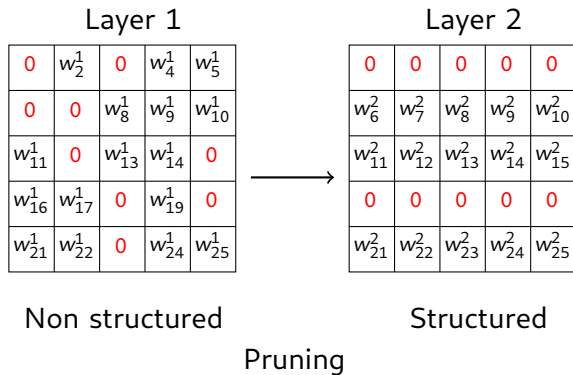
January 13th, 2021

# Pruning



Baseline

# Pruning

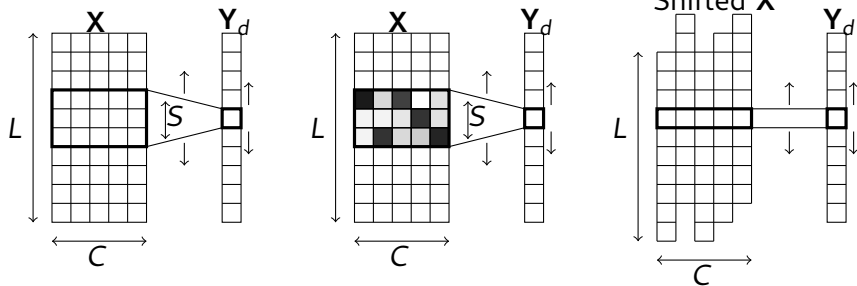


- Evaluate the importance of neurons and eliminate the least important ones to reduce neural network size.
- Non structured pruning: eliminate neurons independently, **only exploitable for very large levels of sparsity.**
- Structured pruning: eliminate kernels, filters or even layers, **exploitable for even low levels of sparsity.**

- Evaluate the importance of neurons and eliminate the least important ones to reduce neural network size.
- Non structured pruning: eliminate neurons independently, **only exploitable for very large levels of sparsity**.
- Structured pruning: eliminate kernels, filters or even layers, **exploitable for even low levels of sparsity**.

- Evaluate the importance of neurons and eliminate the least important ones to reduce neural network size.
- Non structured pruning: eliminate neurons independently, **only exploitable for very large levels of sparsity**.
- Structured pruning: eliminate kernels, filters or even layers, **exploitable for even low levels of sparsity**.

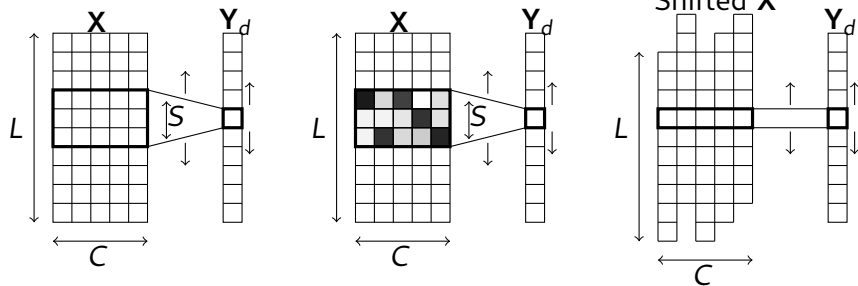
# Structured pruning and shift layers



## Shift Attention Layer (SAL)

- Simplified operations,
- Reduced number of parameters,
- Fully exploitable technique.

# Structured pruning and shift layers



## Shift Attention Layer (SAL)

- Simplified operations,
- Reduced number of parameters,
- Fully exploitable technique.



# Attention Shift Layer

$$\begin{array}{|c|c|c|} \hline w_1 & w_2 & w_3 \\ \hline w_4 & w_5 & w_6 \\ \hline w_7 & w_8 & w_9 \\ \hline \end{array} \times \begin{array}{|c|c|c|} \hline a_1 & a_2 & a_3 \\ \hline a_4 & a_5 & a_6 \\ \hline a_7 & a_8 & a_9 \\ \hline \end{array} = \begin{array}{|c|c|c|} \hline \tilde{w}_1 & \tilde{w}_2 & \tilde{w}_3 \\ \hline \tilde{w}_4 & \tilde{w}_5 & \tilde{w}_6 \\ \hline \tilde{w}_7 & \tilde{w}_8 & \tilde{w}_9 \\ \hline \end{array}$$

Figure: Using attention tensor A to select which weight should be kept

# Conv Operation

8	2	4	6	1	1
3	2	0	9	7	4
0	6	4	0	0	3
5	1	2	0	3	4
1	3	4	0	2	1
0	2	3	1	2	2

 $*$ 

2	0	0
0	0	0
0	0	0

 $=$ 

16	4	8	12
6	4	0	18
0	12	8	0
10	2	4	0

# Shift Operation

8	2	4	6	1	1
3	2	0	9	7	4
0	6	4	0	0	3
5	1	2	0	3	4
1	3	4	0	2	1
0	2	3	1	2	2

 $\times$ 

2
---

 $=$ 

16	4	8	12
6	4	0	18
0	12	8	0
10	2	4	0

# Experimental Results

**Table:** Comparison of accuracy, number of parameters and FLOPs between a standard CNN, SAL and vanilla Shiftnet on ImageNet ILSVRC 2012.

		Top-1	Params	FLOPs
Large budget	ResNet-w24 (CLs)	63.47%	<b>3.2M</b>	664M
	ShiftNet-A	70.1%	4.1M	1.4G
	ResNet-w64 + SAL	<b>71%</b>	3.3M	<b>538M</b>
Small budget	ResNet-w16 (CLs)	56.6%	1.4M	295M
	ShiftNet-B	61.2%	1.1M	371M
	ResNet-w32 + SAL	<b>62.7%</b>	<b>0.97M</b>	136M
Mobile Architecture	MobileNetV2	56.6%	1.76M	57M
	ShuffleNetV2	60.7%	1.3M	<b>41M</b>

# Conclusion and Future Work

## Conclusion

- We introduced novel attention-based pruning method.
- The pruning method aims at replacing convolutional layers by shift layers.
- We showed SAL outperformed other existing methods.

## Future Work

- Extend SAL to all kernel shapes, and to other domains than classification.
- Work on reducing complexity of the training process.

# Conclusion and Future Work

## Conclusion

- We introduced novel attention-based pruning method.
- The pruning method aims at replacing convolutional layers by shift layers.
- We showed SAL outperformed other existing methods.

## Future Work

- Extend SAL to all kernel shapes, and to other domains than classification.
- Work on reducing complexity of the training process.

# Thank you

Thank you for watching this presentation. I will be glad to answer any questions you have via e-mail

[ghouthi.boukliyacene@imt-atlantique.fr](mailto:ghouthi.boukliyacene@imt-atlantique.fr).

## References

Wu, Bichen, et al. "Shift: A zero flop, zero parameter alternative to spatial convolutions." Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. 2018.

Jeon, Yunho, and Junmo Kim. "Constructing fast network through deconstruction of convolution." Advances in Neural Information Processing Systems. 2018.