

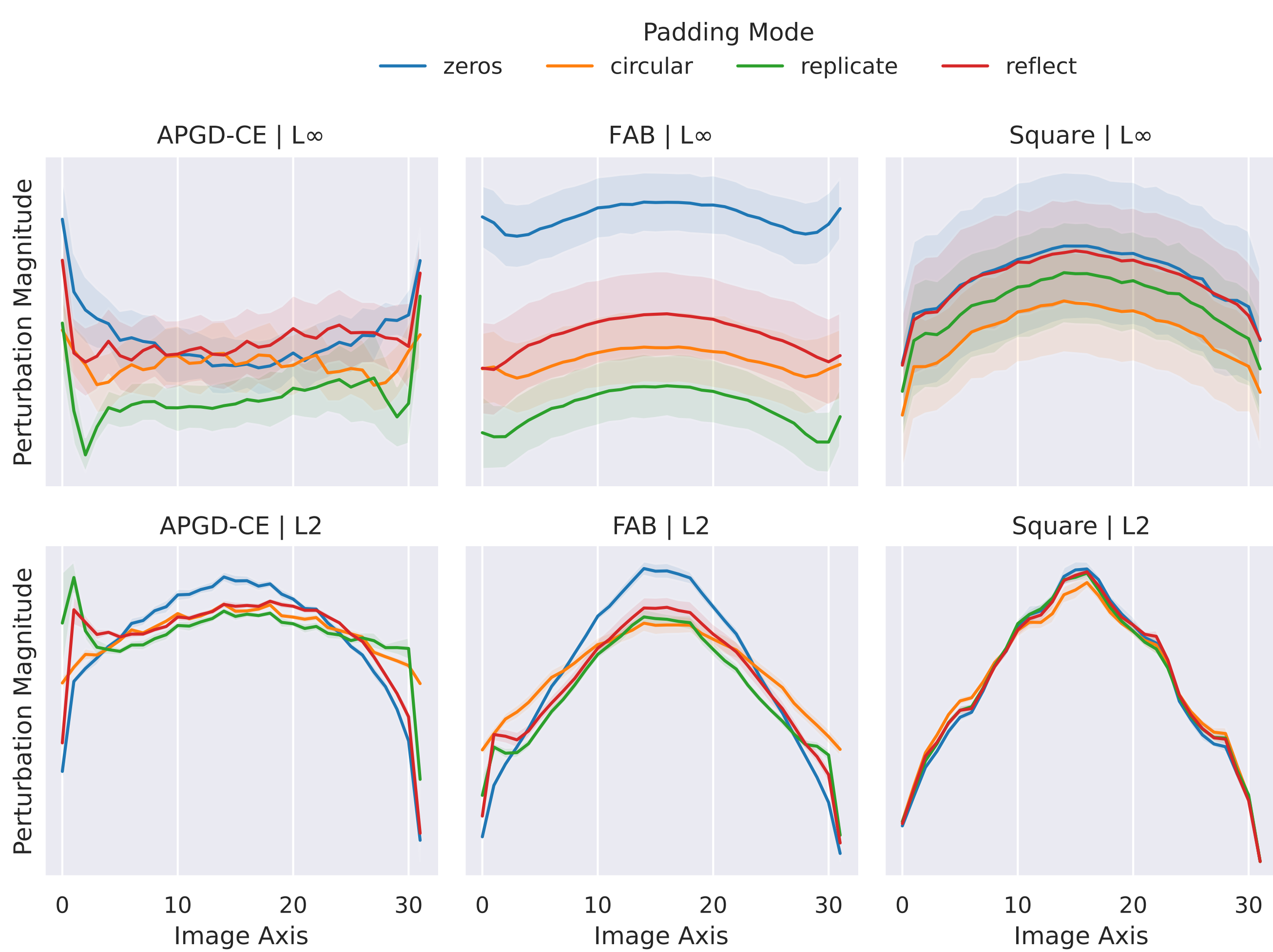


The padding mode is an essential yet rarely tuned CNN hyperparameter. How does its choice affect robustness?

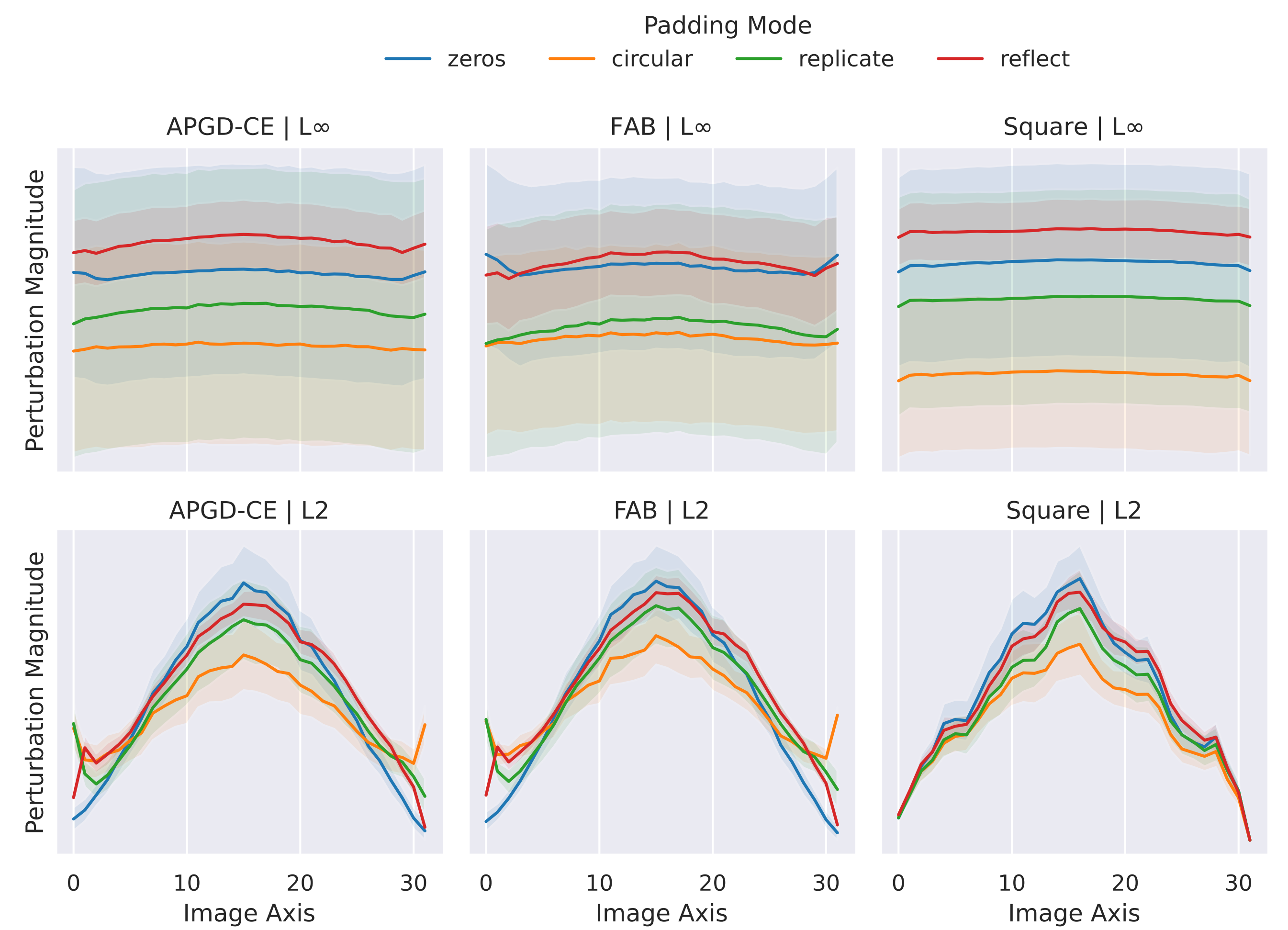


## Normal Training

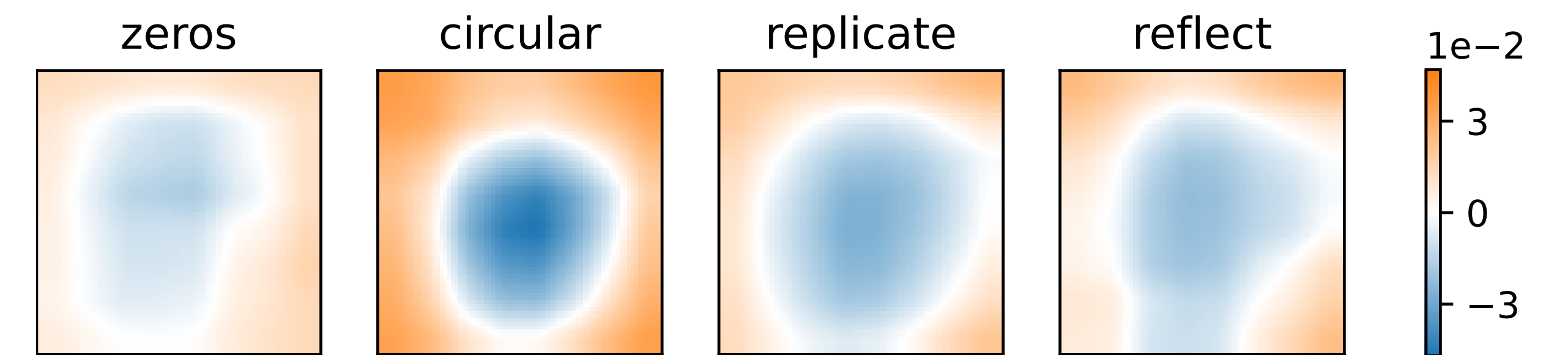
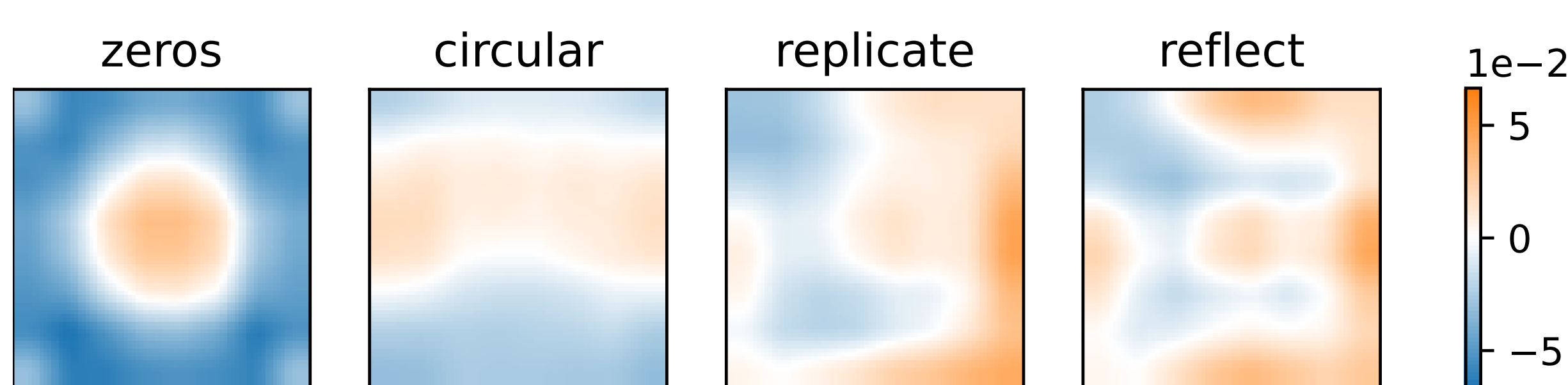
### Anomalies in Perturbations in Padding Regions.



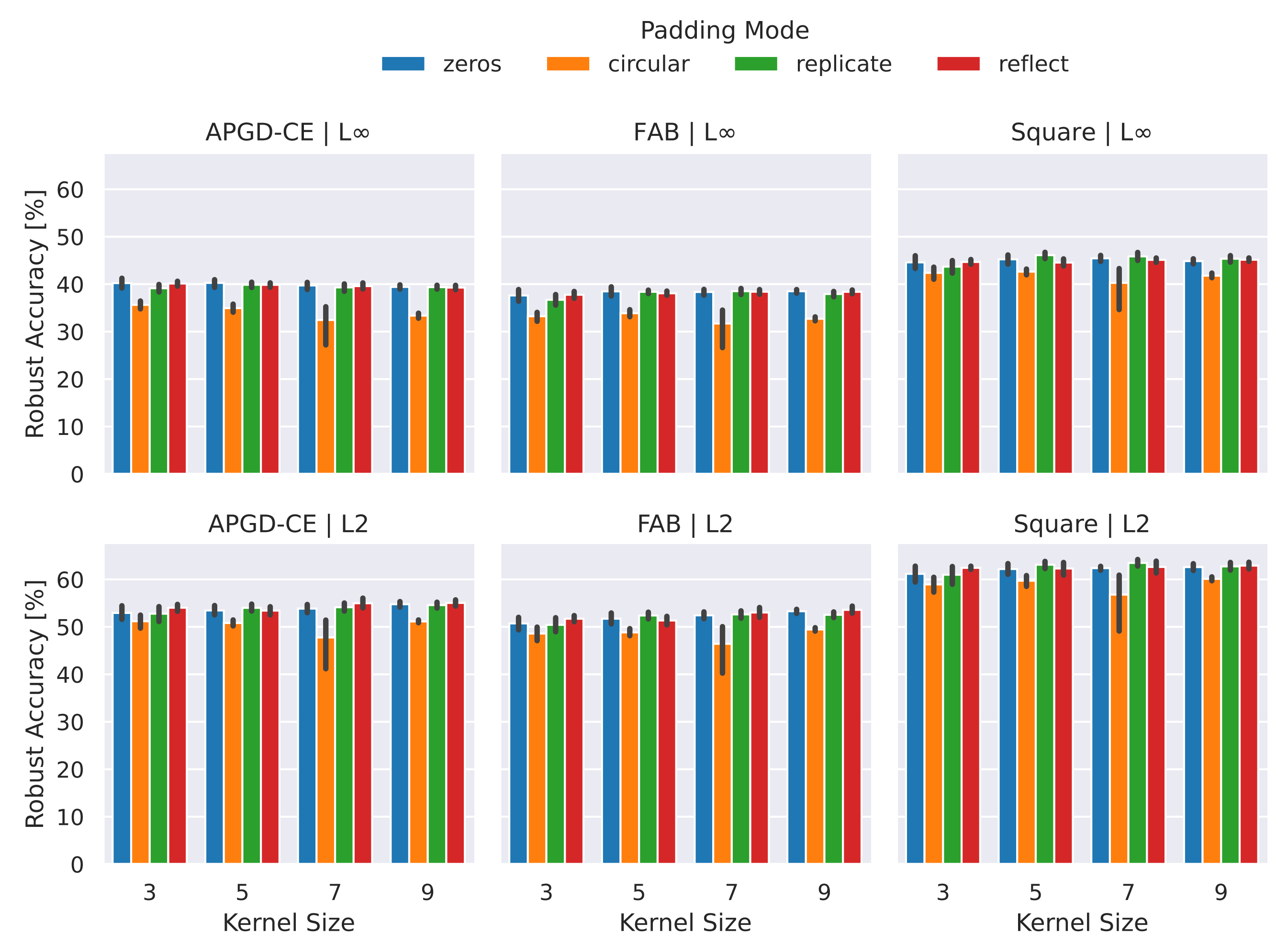
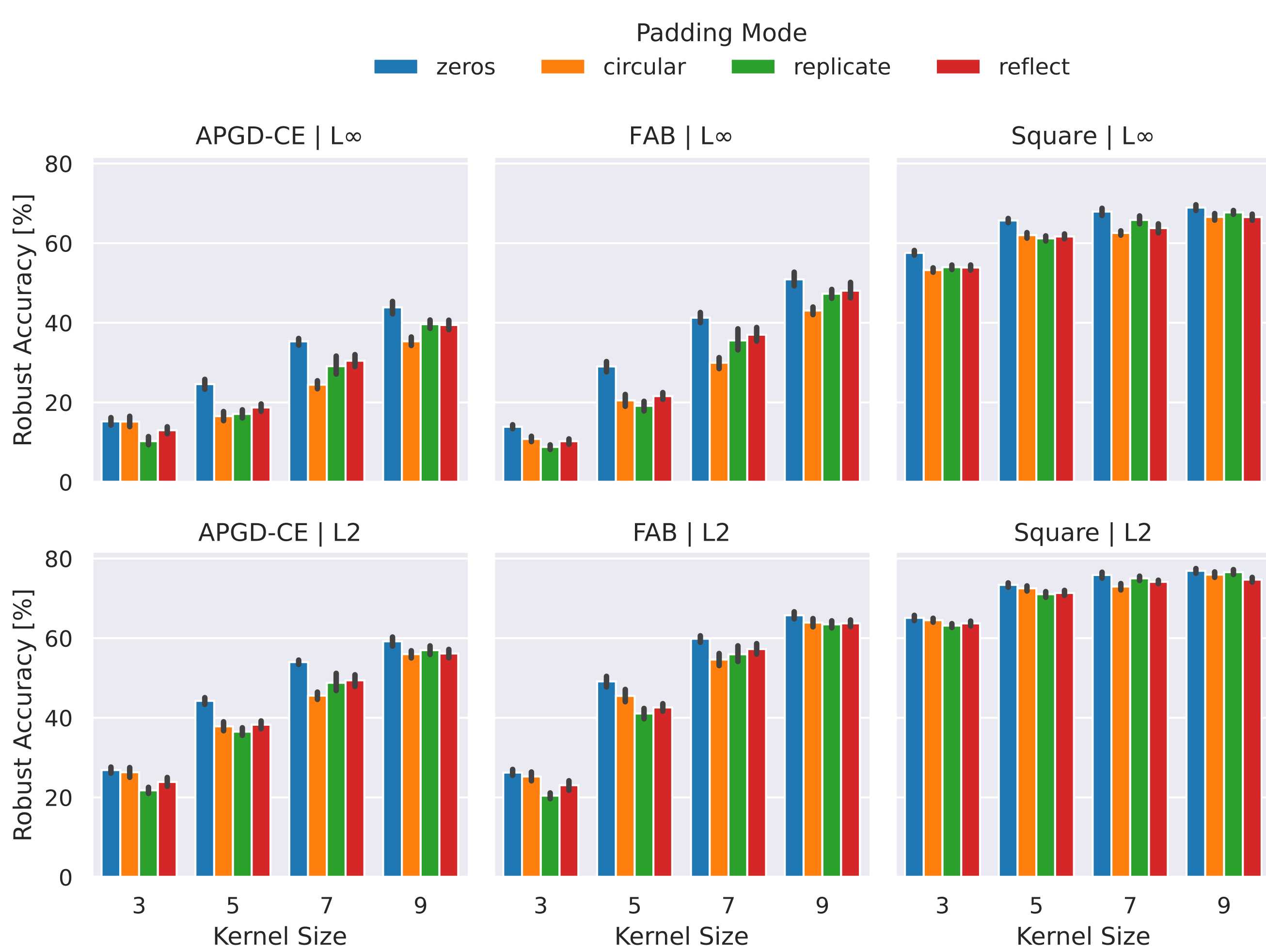
## Adversarial Training



### Effect on Model Decisions.



### Robust Performance.



## Comparison

AT	k	Clean Test [%] (↑)				AutoAttack [%] (↑)			
		zeros	circular	replicate	reflect	zeros	circular	replicate	reflect
x	3	<b>90.26</b>	90.10	90.13	90.15	<b>8.52</b>	4.69	4.90	<u>5.79</u>
	5	<b>90.14</b>	89.66	89.82	89.67	<b>17.69</b>	10.44	11.12	<u>12.33</u>
	7	<b>89.36</b>	88.49	88.52	88.47	<b>29.06</b>	17.86	<u>24.55</u>	24.35
	9	<b>88.22</b>	87.50	87.03	87.25	<b>39.18</b>	30.52	<u>36.39</u>	34.81
✓	3	71.84	69.17	70.79	<b>73.11</b>	<b>36.88</b>	32.09	35.91	<u>36.82</u>
	5	73.72	71.34	<b>74.02</b>	73.08	<b>37.48</b>	32.34	37.30	37.12
	7	73.86	67.33	<b>73.89</b>	73.10	<b>37.42</b>	30.16	37.08	<u>37.26</u>
	9	<u>73.51</u>	71.53	72.24	<b>73.90</b>	<b>37.49</b>	31.09	36.89	<u>37.25</u>

## Take-Home Messages

- Padding results in anomalies in the spatial distribution of adversarial attacks.
- Increasing the kernel size (and padding) natively improves robustness without adversarial training.
- Zero padding performs best in, both, clean and adversarial evaluation with normal training.
- Adversarial training balances the robust performance under different padding modes (except *circular*) and kernel sizes.
- When using adversarial training, *replicate/reflect* notably improves clean performance with marginal impairments in robust performance compared to *zero* padding.
- Padding is an essential operation. Removing padding results in deteriorated performance in clean and adversarial settings.
- Limitation: We only studied image classification on CIFAR-10 with ResNet-20. As with many “toy datasets”, objects in question are usually perfectly centered in the images → not clear if the results transfer to real-world scenarios.