

# Robustness as Architecture: Designing IQA Models to Withstand Adversarial Perturbations

Igor Meleshin<sup>1</sup>, Anna Chistyakova<sup>1,2</sup>, Anastasia Antsiferova<sup>2,3,4</sup>, Dmitriy Vatolin<sup>1,2,3</sup>

<sup>1</sup>Lomonosov Moscow State University, Moscow, Russia

<sup>2</sup>ISP RAS Research Center for Trusted Artificial Intelligence, Moscow, Russia

<sup>3</sup>MSU Institute for Artificial Intelligence, Moscow, Russia

<sup>4</sup>Laboratory of Innovative Technologies for Processing Video Content, Innopolis University, Innopolis, Russia

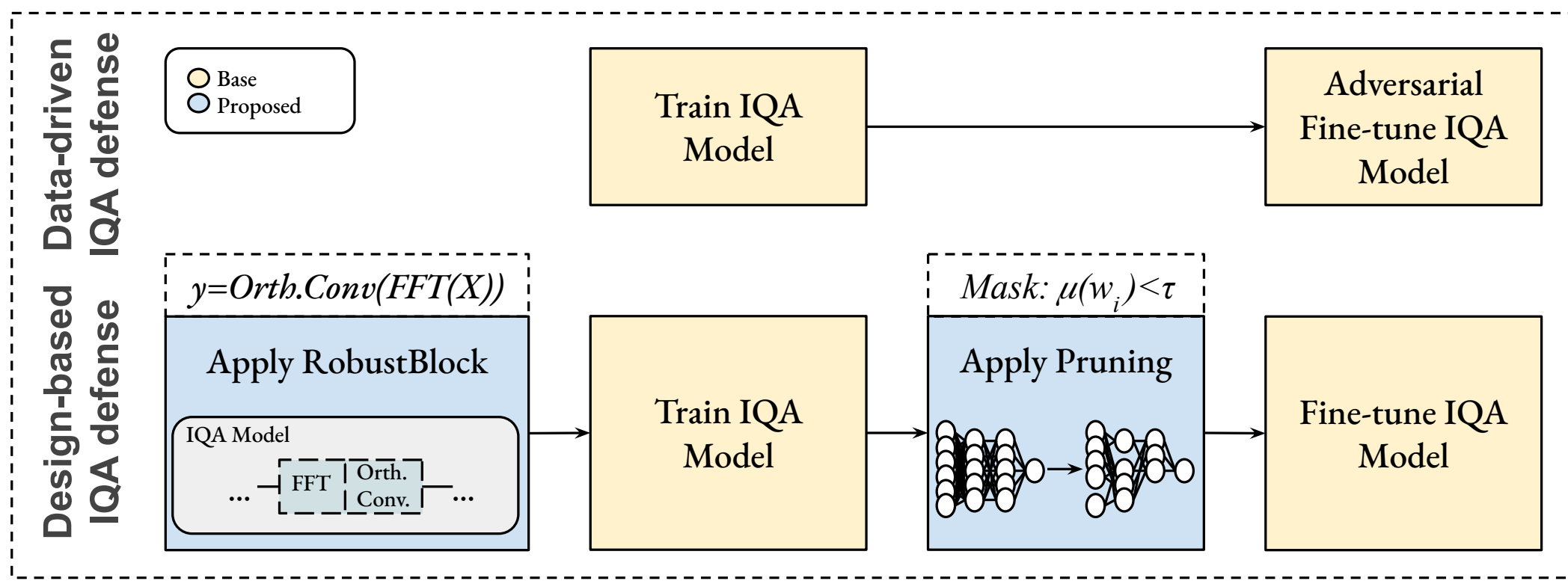


## MOTIVATION

Image Quality Assessment (IQA) models are critical for modern vision systems, from compression and enhancement to generation and streaming. However, they suffer from fundamental instability and can be easily compromised by minimal, visually imperceptible adversarial perturbations. This vulnerability leads to unreliable quality scores in critical applications and raises concerns about the trustworthiness of automated perceptual evaluation.

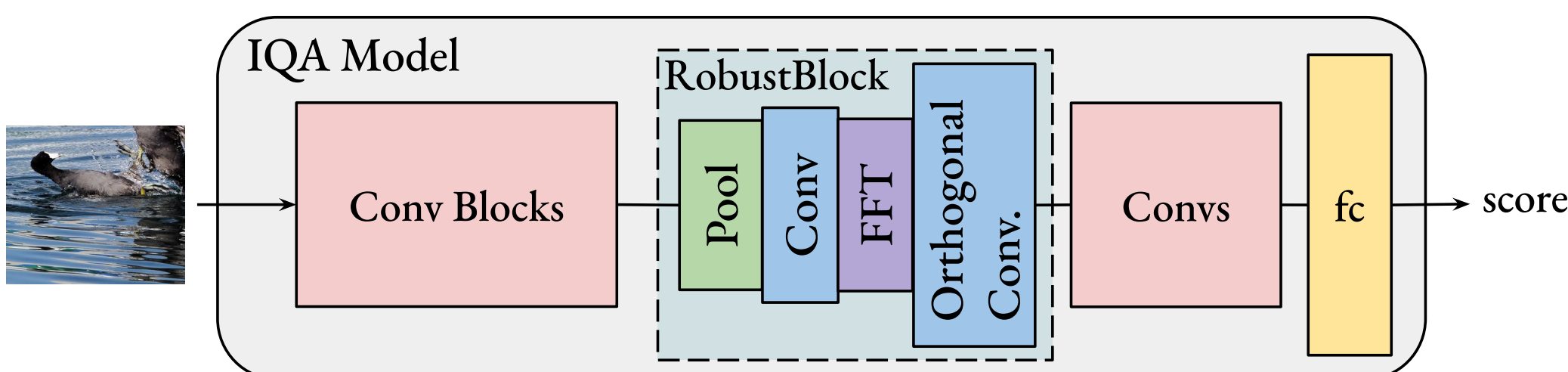
## IDEA

Traditional defenses rely on expensive and attack-specific data-driven retraining or purification. We propose an architectural modification to increase the robustness of IQA models against adversarial attacks.



## ARCHITECTURE

We introduce a lightweight, modular pipeline built around the RobustBlock: FFT-based feature transformation to frequency domain, orthogonal convolution for norm-preserving mappings, and pruning. RobustBlock receives a squared feature matrix as input, so we use pooling and convolutions to transform the input tensor.



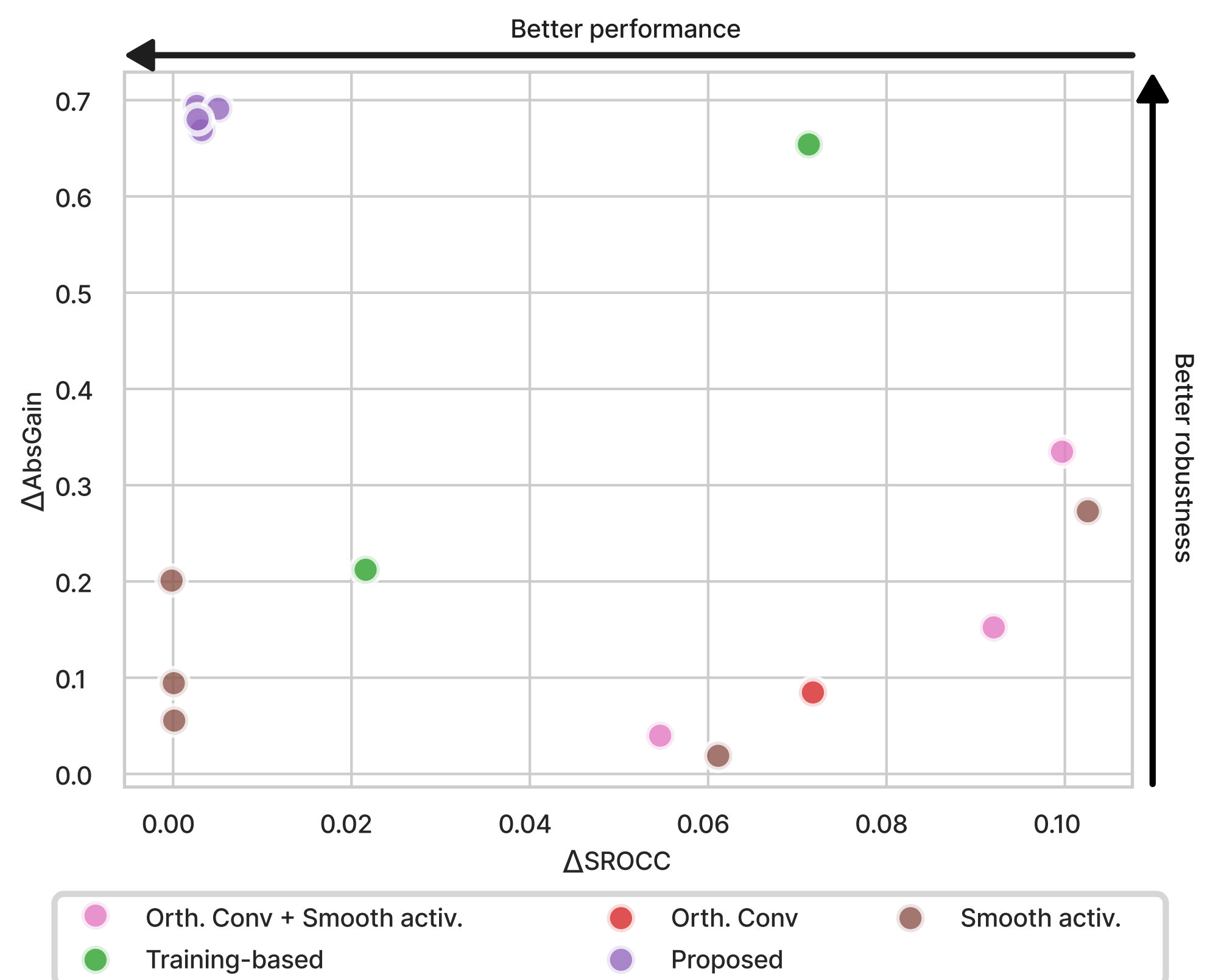
## BEST POSITION

We compared a few positions of the RobustBlock in the architecture. The optimal placement of the RobustBlocks is near fully connected layers where perturbation amplification is minimal.

Modification	InputDim	SROCC	PLCC	AbsGain	RScore	Time	
						train epoch	test
---	---	<b>0.926</b>	<b>0.936</b>	0.409	0.012	<b>100.727</b>	<b>31.735</b>
AT	---	0.855	0.877	<b>0.209</b>	0.025	183.682	<b>32.216</b>
NT	---	0.904	0.921	<b>0.265</b>	<b>0.026</b>	229.612	33.028
RobustBlock 1	3×498×664	0.766	0.794	0.415	0.018	131.172	36.787
RobustBlock 2	64×125×166	0.793	0.823	0.508	0.006	<b>129.313</b>	63.864
RobustBlock 4	128×63×83	0.854	0.876	0.402	0.016	199.581	103.655
RobustBlock 6	1024×32×42	<b>0.923</b>	<b>0.936</b>	0.289	<b>0.031</b>	181.826	89.703

## COMPARISON

Trade-off between robustness and quality correlation on the NIPS2017 dataset. Our architectural modification (purple) consistently reduces vulnerability to adversarial attacks (lower AbsGainAUC) while maintaining high perceptual alignment (minimal drop in SROCC).



## RESULTS

Performance and robustness metrics of IQA models with different architectural variations. AbsGain and R-Score are shown for UAP and stAdv, while area-under-curve (AUC) values are provided for PGD attacks with 1 and 8 iterations. The best result for each metric within each model group is highlighted in bold.

	SROCC	PLCC	AbsGainAUC		RScoreAUC		AbsGain		RScore	
			PGD-1		UAP		stAdv		stAdv	
Linearity	<b>0.926</b>	<b>0.936</b>	0.409	0.013	0.465	0.463	0.025	1.732		
Linearity+NT	0.904	0.921	<b>0.305</b>	<b>0.026</b>	<b>0.190</b>	<b>0.825</b>	<b>0.015</b>	<b>1.937</b>		
Linearity+AT	0.855	0.877	0.365	<b>0.026</b>	<b>0.322</b>	0.491	<b>0.013</b>	<b>1.985</b>		
Linearity+our	<b>0.921</b>	<b>0.935</b>	<b>0.261</b>	<b>0.034</b>	0.360	<b>0.668</b>	0.022	1.862		
KonCept	<b>0.915</b>	<b>0.929</b>	0.359	0.015	1.276	0.293	<b>0.006</b>	2.176		
KonCept+NT	0.794	0.808	<b>0.122</b>	0.036	<b>0.419</b>	<b>0.608</b>	<b>0.001</b>	<b>3.873</b>		
KonCept+AT	0.815	0.850	<b>0.062</b>	<b>0.046</b>	0.633	<b>0.488</b>	0.009	<b>2.386</b>		
KonCept+our	<b>0.884</b>	<b>0.908</b>	0.244	<b>0.040</b>	<b>0.396</b>	0.413	<b>0.006</b>	2.324		
TReS	<b>0.918</b>	<b>0.929</b>	<b>0.300</b>	<b>0.027</b>	0.472	0.181	<b>0.103</b>	0.862		
TReS+NT	0.887	0.901	<b>0.307</b>	<b>0.029</b>	<b>0.273</b>	<b>0.395</b>	0.122	<b>1.427</b>		
TReS+AT	0.886	<b>0.904</b>	0.357	0.026	<b>0.426</b>	<b>0.188</b>	0.117	0.846		
TReS+our	<b>0.919</b>	<b>0.929</b>	0.314	<b>0.029</b>	0.503	0.178	<b>0.051</b>	<b>0.903</b>		

## FUTURE WORK

Although our FFT-based robustness method provides strong gains, it introduces trade-offs like moderate computational overhead and limitations to convolutional networks with square inputs into RobustBlock. We view these constraints as opportunities for future work, focusing on faster spectral approximations and extending the design to new architectures.

