

# Boosting Adversarial Transferability with Spatial Adversarial Alignment



Zhaoyu Chen<sup>1,\*</sup>, Haijing Guo<sup>2,\*</sup>, Kaixun Jiang<sup>1</sup>, Jiyuan Fu<sup>2</sup>, Xinyu Zhou<sup>2</sup>, Dingkan Yang<sup>1</sup>, Hao Tang<sup>3</sup>, Bo Li<sup>4</sup>, Wenqiang Zhang<sup>1,2,†</sup>

<sup>1</sup>College of Intelligent Robotics and Advanced Manufacturing, Fudan University

<sup>2</sup>College of Computer Science and Artificial Intelligence, Fudan University

<sup>3</sup>School of Computer Science, Peking University <sup>4</sup>vivo Mobile Communication Co., Ltd

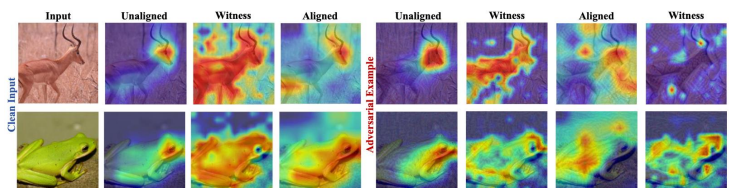
## Motivation

Cross-model transferability has been extensively studied for CNNs. However, few works explore adversarial transferability on ViT and the performance of existing work extending CNN to ViT is poor due to significant structural differences.

We argue that **unique structural features** are critical to cross-architecture adversarial transferability. If we can obtain a surrogate model whose features are similar to those of models with different architectures, then the resulting adversarial perturbation can be transferable across different models.

However, directly applying this idea to black-box attacks may lead to the degradation of cross-architecture transferability. The main reasons are:

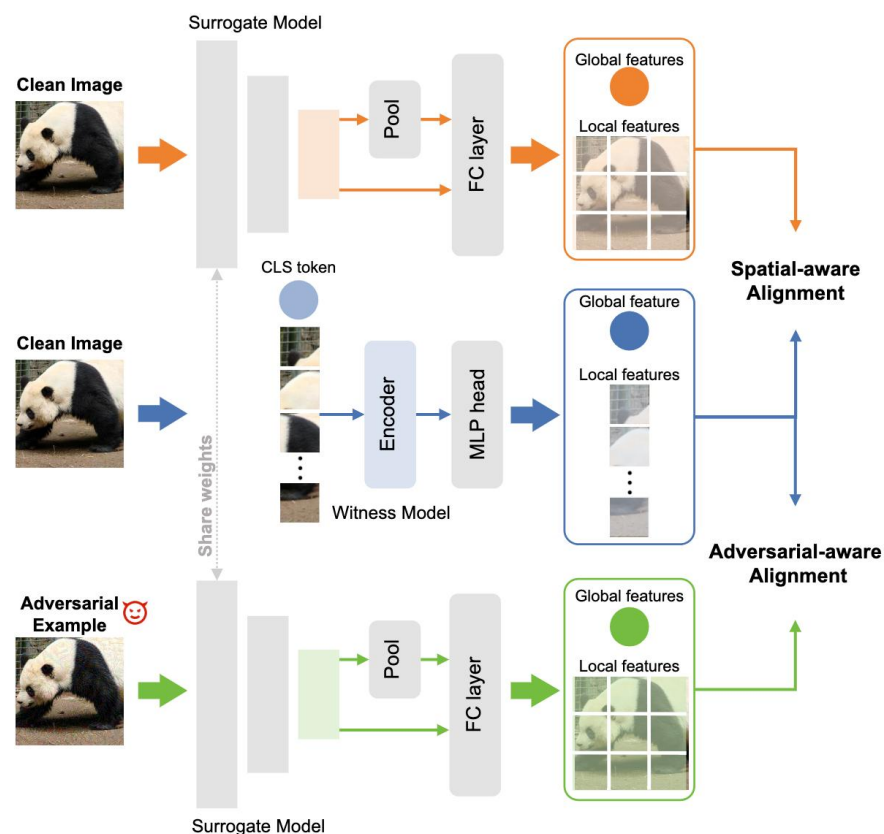
1. Features are not aligned in space. It is hard to directly constrain the similarity of features only by the final logits.
2. Features are not aligned from the perspective of adversarial features. The features of adversarial examples also have similarities across models and need to be considered.



## Method

**Spatial Adversarial Alignment** consists of two key parts:

1. In **spatial-aware alignment**, in addition to aligning on the final global features, we also focus on the features of local regions.
2. In **adversarial-aware alignment**, we introduce a self-adversarial strategy, which constructs adversarial examples so that the model can learn the differences between different architectures in adversarial features, thereby enabling the model to further capture more common features.



## Experiments

SAA has stronger adversarial transferability with transfer attacks.

Attack	Target Model										Avg. ASR (%)
	CNNs					ViTs					
	Res18	Res50	Res101	VGG19	DN121	Inc-v3	ViT-B	Swin-B	PVT-v2	MobViT	
MI	57.7	<b>99.9</b>	58.1	54.2	55.1	39.0	9.4	33.0	38.0	35.7	42.2
MI-SAA	<b>84.1</b>	99.6	<b>74.7</b>	<b>80.3</b>	<b>81.8</b>	<b>65.7</b>	<b>24.3</b>	<b>48.7</b>	<b>52.9</b>	<b>62.5</b>	<b>63.9</b>
NI	58.9	<b>100.0</b>	63.2	59.3	61.4	40.0	9.6	37.4	41.8	38.1	45.5
NI-SAA	<b>86.1</b>	<b>99.9</b>	<b>76.3</b>	<b>82.2</b>	<b>83.7</b>	<b>67.6</b>	<b>24.0</b>	<b>50.6</b>	<b>55.7</b>	<b>64.8</b>	<b>65.7</b>
GI	57.3	<b>100.0</b>	62.3	60.5	60.5	40.7	12.5	36.8	40.7	39.6	45.7
GI-SAA	<b>86.5</b>	99.7	<b>78.8</b>	<b>83.9</b>	<b>84.8</b>	<b>70.4</b>	<b>27.6</b>	<b>52.7</b>	<b>55.8</b>	<b>66.3</b>	<b>67.4</b>
DI	44.1	<b>95.8</b>	41.7	56.1	44.2	26.1	5.6	30.9	36.7	35.1	35.6
DI-SAA	<b>74.6</b>	94.4	<b>61.1</b>	<b>81.1</b>	<b>73.2</b>	<b>53.1</b>	<b>10.6</b>	<b>44.0</b>	<b>50.8</b>	<b>63.6</b>	<b>56.9</b>
TI	38.5	<b>99.9</b>	<b>37.8</b>	33.9	36.1	24.2	5.4	21.0	29.0	21.1	27.4
TI-SAA	<b>59.7</b>	94.9	35.2	<b>50.6</b>	<b>54.5</b>	<b>40.6</b>	<b>9.5</b>	<b>22.8</b>	<b>29.3</b>	<b>33.6</b>	<b>37.3</b>
SSA	75.8	<b>99.9</b>	<b>78.6</b>	76.0	77.8	57.0	16.5	<b>48.5</b>	55.0	50.5	59.5
SSA-SAA	<b>91.5</b>	99.5	77.8	<b>85.7</b>	<b>88.4</b>	<b>74.9</b>	<b>23.4</b>	46.7	<b>57.1</b>	<b>66.0</b>	<b>67.9</b>
DI-MI	65.5	97.0	65.0	74.7	65.7	49.1	16.4	49.0	54.9	57.9	55.4
DI-MI-SAA	<b>91.9</b>	<b>98.7</b>	<b>84.9</b>	<b>94.1</b>	<b>90.5</b>	<b>78.3</b>	<b>34.1</b>	<b>69.8</b>	<b>76.1</b>	<b>86.8</b>	<b>78.5</b>
TI-MI	61.4	<b>99.9</b>	60.5	60.9	60.9	44.3	15.2	37.4	42.3	41.8	47.2
TI-MI-SAA	<b>84.8</b>	99.3	<b>71.9</b>	<b>79.0</b>	<b>81.8</b>	<b>69.1</b>	<b>27.0</b>	<b>45.0</b>	<b>52.8</b>	<b>62.4</b>	<b>63.8</b>
SSA-MI	89.6	<b>99.9</b>	92.2	89.5	91.0	77.6	39.2	<b>74.4</b>	76.4	76.3	78.5
SSA-MI-SAA	<b>96.3</b>	99.8	<b>95.6</b>	<b>96.5</b>	<b>97.2</b>	<b>91.5</b>	<b>46.3</b>	74.1	<b>80.4</b>	<b>88.4</b>	<b>85.1</b>
SSA-DI-TI-MI	93.5	98.5	92.3	95.0	93.7	85.5	<b>55.9</b>	<b>83.2</b>	87.1	89.8	86.2
SSA-DI-TI-MI-SAA	<b>97.5</b>	<b>98.8</b>	<b>93.8</b>	<b>97.6</b>	<b>96.7</b>	<b>94.3</b>	53.6	81.6	<b>84.4</b>	<b>94.7</b>	<b>88.2</b>

SAA further improves the adversarial transferability on ViTs.

Attack	Target Model										Avg. ASR (%)
	CNNs						ViTs				
	Res18	Res50	Res101	VGG19	DN121	Inc-v3	ViT-B	Swin-B	PVT-v2	MobViT	
SGM	82.9	67.6	59.4	81.2	75.4	71.3	99.7	83.3	72.7	78.8	78.3
SGM-SAA	91.1	79.8	73.3	87.5	87.3	80.9	99.5	90.5	82.6	86.3	86.6
PatchOut	45.6	27.4	20.3	45.5	36.1	33.9	93.0	41.0	34.2	40.5	43.3
PatchOut-SAA	76.5	72.4	70.3	79.4	78.1	71.3	94.7	83.6	77.2	76.8	78.7
PNA	61.2	45.0	38.1	60.8	54.8	49.0	99.6	66.3	55.8	56.8	60.3
PNA-SAA	82.7	78.1	73.4	85.6	84.0	75.1	97.4	89.3	80.3	82.1	83.3
TGR	74.0	55.6	48.4	73.2	66.6	59.0	99.3	74.5	61.6	69.6	69.6
TGR-SAA	85.9	78.1	71.5	87.4	85.6	79.6	99.3	89.1	81.0	86.2	85.1

SAA improves adversarial transferability against defenses.

Attack	HGD	R&P	NIPS-r3	JPEG	FD	RS	Bit-Red	NRP	Diffuse	Inc-v3 <sub>ens3</sub>	Inc-v3 <sub>ens4</sub>	IncRes-v2 <sub>ens</sub>	Avg. ASR (%)
MI	42.5	21.9	25.3	33.9	42.4	23.6	29.3	6.7	13.8	33.3	31.3	23.3	27.3
MI-SAA	<b>73.3</b>	<b>57.8</b>	<b>60.4</b>	<b>69.9</b>	<b>65.6</b>	<b>39.6</b>	<b>42.3</b>	<b>12.0</b>	<b>22.4</b>	<b>68.9</b>	<b>65.5</b>	<b>57.9</b>	<b>53.0</b>
SSA-DI-TI-MI	93.7	89.6	90.2	91.9	89.7	82.3	81.7	14.8	71.1	92.5	91.1	89.8	81.5
SSA-DI-TI-MI-SAA	<b>96.0</b>	<b>93.2</b>	<b>94.8</b>	<b>95.1</b>	<b>94.0</b>	<b>89.8</b>	<b>85.5</b>	<b>20.2</b>	<b>78.8</b>	<b>95.7</b>	<b>94.5</b>	<b>93.3</b>	<b>85.9</b>

Grad-CAM on target models.

