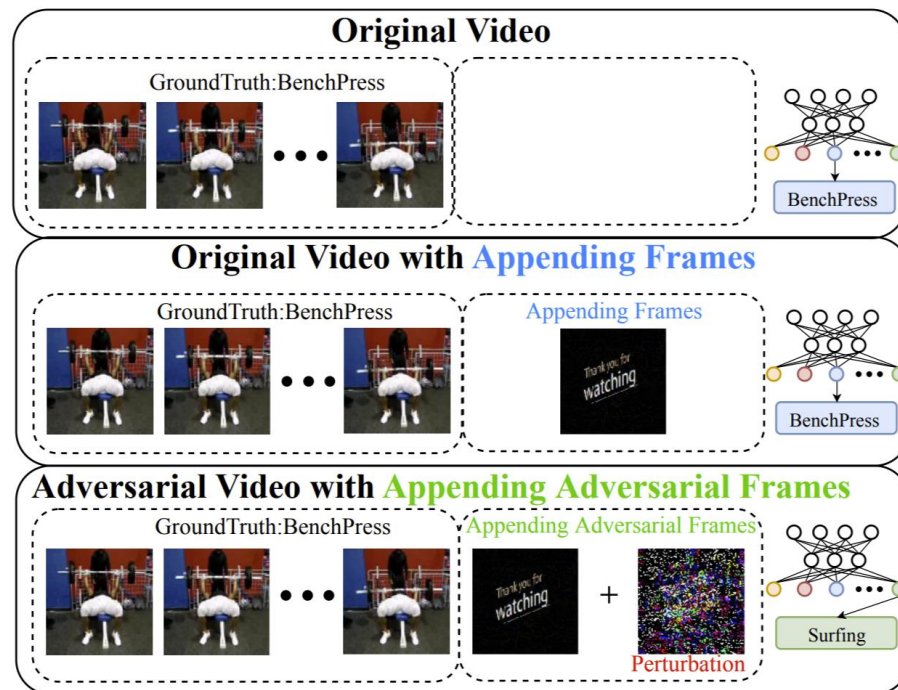


# **Appending Adversarial Frames for Universal Video Attack**

# AAF



# Basic Attack Methods

- **Model:**

$$X = \{f_1, f_2 \dots, f_T\}$$

$$E = \{p_1, p_2 \dots p_T\}$$

$$\hat{X} = \{f_1 \oplus p_1, f_2 \oplus p_2 \dots, f_T \oplus p_T\}$$

- **Weakness:**

High authority needed.

Unsafe. ( $f_1, f_2 \dots$  are related while  $p_1, p_2$  are not)

High perturbation rate. (every frame)

Weak transferability.

# Appending Adversarial Frames Method

- **Model:**

$$X \in R^{T \times W \times H \times C} \quad (\text{T: number of frames } W, H, C: \text{width, height, and channel of each frame})$$

$$\Delta \in R^{\Delta T \times W \times H \times C} \quad (\text{adversarial frames without perturbations})$$

$$\hat{\Delta} \in R^{\Delta T \times W \times H \times C} \quad (\text{adversarial frames with perturbations})$$

$$\hat{X} \in R^{(T+\Delta T) \times W \times H \times C} \quad (\text{adversarial video})$$

- **Optimization Function**

$$\arg \min_{\mathbf{E}} \lambda ||\mathbf{E}||_p - \ell(\mathbf{1}_y, \mathbb{J}(\hat{\mathbf{X}}; \boldsymbol{\theta}))$$

$$\arg \min_{\mathbf{E}} \lambda ||\mathbf{E}||_p + \ell(\mathbf{1}_{y^*}, \mathbb{J}(\hat{\mathbf{X}}; \boldsymbol{\theta}))$$

# Variants of AAFM

- **Across Videos :**

$$\arg \min_{\mathbf{E}} \lambda \|\mathbf{E}\|_p - \sum_{n=1}^N \alpha_n \ell(\mathbf{1}_{y_n}, \mathbb{J}(\hat{\mathbf{X}}_n; \boldsymbol{\theta}))$$

- **Across Models:**

$$\arg \min_{\mathbf{E}} \lambda \|\mathbf{E}\|_p - \sum_{k=1}^K \beta_k \ell(\mathbf{1}_{y_k}, \mathbb{J}(\hat{\mathbf{X}}; \boldsymbol{\theta}_k))$$

- **Feature Similarity:**

$$\begin{aligned} \arg \min_{\mathbf{E}} \lambda \|\mathbf{E}\|_p - \ell(\mathbf{1}_y, \mathbb{J}(\hat{\mathbf{X}}; \boldsymbol{\theta})) \\ + \lambda_l \|\phi_l(\boldsymbol{\Delta}_s) - \phi_l(\hat{\boldsymbol{\Delta}})\|_p \end{aligned}$$

# Experiments

Table 2. Comparison of BAM and A<sup>2</sup>FM with different video classification models.

Target Model	Methods	UCF-101		HMDB-51	
		FR (%)	AAP	FR (%)	AAP
I3D-ResNet	BAM	<b>100</b>	0.22	<b>100</b>	0.31
	A <sup>2</sup> FM	<b>100</b>	<b>0.05</b>	<b>100</b>	<b>0.06</b>
I3D-Inception	BAM	<b>99.5</b>	0.20	<b>100</b>	0.28
	A <sup>2</sup> FM	<b>99.5</b>	<b>0.08</b>	<b>100</b>	<b>0.07</b>
CNN+LSTM	BAM	<b>100</b>	0.20	<b>100</b>	0.28
	A <sup>2</sup> FM	<b>100</b>	<b>0.02</b>	<b>100</b>	<b>0.02</b>
C3D	BAM	<b>99.5</b>	0.24	<b>100</b>	0.30
	A <sup>2</sup> FM	97.3	<b>0.14</b>	96.8	<b>0.16</b>
ResNet3D	BAM	<b>97.8</b>	0.25	<b>100</b>	0.30
	A <sup>2</sup> FM	95.1	<b>0.09</b>	<b>100</b>	<b>0.07</b>
P3D	BAM	<b>100</b>	0.20	<b>100</b>	0.28
	A <sup>2</sup> FM	<b>100</b>	<b>0.02</b>	<b>100</b>	<b>0.02</b>

fewer perturbations

Table 3. Comparison of BAM and A<sup>2</sup>FM-AV in transferability across different videos.

Target Model	Methods	UCF-101		HMDB-51	
		FR (%)	AAP	FR (%)	AAP
I3D-ResNet	BAM	95.4	0.62	93.0	0.70
	A <sup>2</sup> FM-AV	<b>98.1</b>	<b>0.52</b>	<b>97.8</b>	<b>0.60</b>
I3D-Inception	BAM	2.6	<b>0.34</b>	2.0	<b>0.25</b>
	A <sup>2</sup> FM-AV	<b>69.3</b>	1.25	<b>2.3</b>	0.84
CNN+LSTM	BAM	18.1	<b>0.09</b>	<b>69.6</b>	<b>0.13</b>
	A <sup>2</sup> FM-AV	<b>47.1</b>	0.16	45.7	0.21
C3D	BAM	97.9	<b>0.75</b>	<b>98.0</b>	<b>0.68</b>
	A <sup>2</sup> FM-AV	<b>98.1</b>	1.21	96.9	1.75
ResNet3D	BAM	45.2	<b>0.65</b>	58.6	<b>0.49</b>
	A <sup>2</sup> FM-AV	<b>96.6</b>	1.21	<b>94.1</b>	0.79
P3D	BAM	20.7	<b>0.11</b>	46.9	0.16
	A <sup>2</sup> FM-AV	<b>98.4</b>	0.25	<b>97.4</b>	<b>0.15</b>

Better transferability  
AAP is not guaranteed

# Experiments

Table 4. Comparison of BAM and A<sup>2</sup>FM-AM in transferability across models on UCF-101 dataset. The first column indicates we use the Leave-One-Out ensemble method that excludes one model to produce perturbations. For instance, ‘–I3D-ResNet’ means the corresponding ensemble model excludes I3D-ResNet. The numbers in the 3-8 columns are the fooling rates (%) for each attacked model.

Models	Method	I3D-ResNet	ResNet3D	P3D	I3D-Inception	C3D	CNN+LSTM
–I3D-ResNet	BAM	0	78.7	84.6	87.8	70.8	56.2
	A <sup>2</sup> FM-AM	<b>39.5</b>	68.1	97.4	42.9	85.4	81.6
–ResNet3D	BAM	100	0	84.6	87.8	70.8	38.9
	A <sup>2</sup> FM-AM	89.5	<b>6.4</b>	97.4	52.2	85.4	71.4
–P3D	BAM	100	80.9	15.4	87.8	72.9	58.8
	A <sup>2</sup> FM-AM	86.8	74.5	<b>59.0</b>	50.0	85.4	83.7
–I3D-Inception	BAM	100	83.0	97.4	0	73.0	61.1
	A <sup>2</sup> FM-AM	86.8	78.7	100	<b>2.0</b>	85.4	50.0
–C3D	BAM	100	83.0	100	90.0	0	64.7
	A <sup>2</sup> FM-AM	92.1	80.9	100	60.0	<b>20.8</b>	79.6
–CNN+LSTM	BAM	100	80.9	97.4	97.8	72.9	35.7
	A <sup>2</sup> FM-AM	89.5	74.5	100	55.6	85.4	<b>77.6</b>

Better transferability

# Experiments

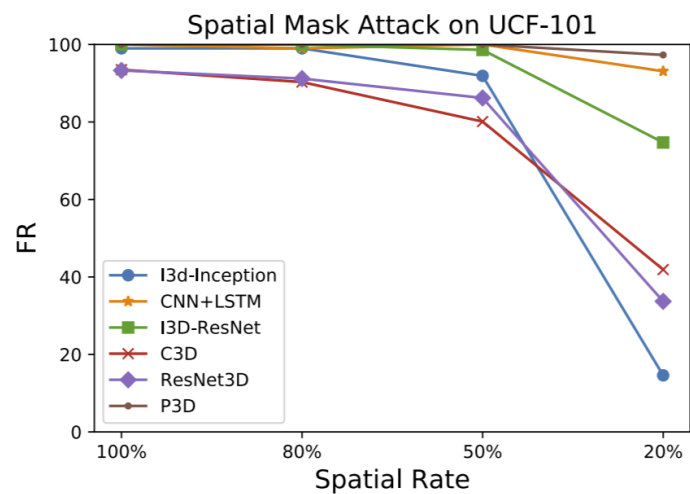
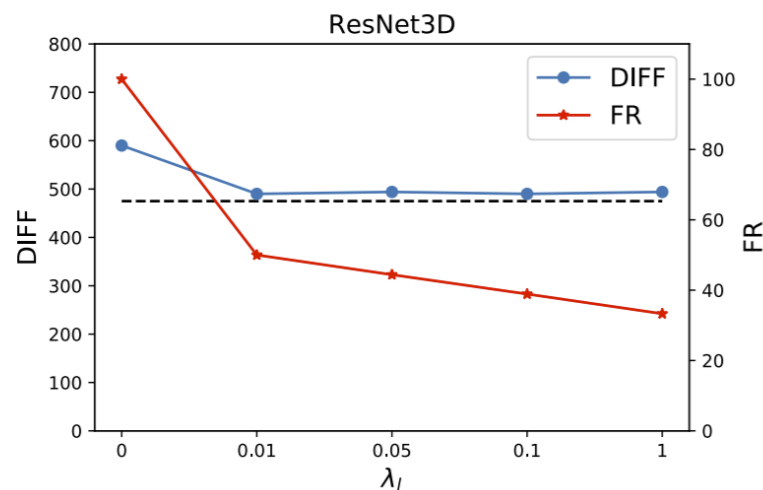


Table 5. Comparison of BAM and A<sup>2</sup>FM for targeted attack.

Target Model	Methods	UCF-101		HMDB-51	
		FR (%)	AAP	FR (%)	AAP
I3D-ResNet	BAM	97.6	0.29	<b>97.8</b>	0.31
	A <sup>2</sup> FM	<b>97.7</b>	<b>0.17</b>	<b>97.8</b>	<b>0.14</b>
I3D-Inception	BAM	<b>84.6</b>	0.23	<b>96.8</b>	0.27
	A <sup>2</sup> FM	27.4	<b>0.08</b>	40.2	<b>0.08</b>
CNN+LSTM	BAM	<b>61.6</b>	0.23	<b>55.8</b>	0.27
	A <sup>2</sup> FM	53.2	<b>0.07</b>	42.4	<b>0.07</b>
C3D	BAM	<b>97.9</b>	0.30	<b>97.8</b>	0.31
	A <sup>2</sup> FM	83.8	<b>0.26</b>	95.0	<b>0.22</b>
Resnet3D	BAM	<b>98.1</b>	0.28	<b>98.0</b>	0.30
	A <sup>2</sup> FM	<b>98.1</b>	<b>0.15</b>	<b>98.0</b>	<b>0.13</b>
P3D	BAM	<b>98.0</b>	0.22	<b>97.8</b>	0.26
	A <sup>2</sup> FM	97.8	<b>0.07</b>	<b>97.8</b>	<b>0.08</b>

Bad performance on targeted attack



# Conclusion

The paper present an adversarial video attack method, which appends a few dummy frames with adversarial perturbations to the original video.

Comparing with the basic adversarial video attack method, AAF has a better performance on perturbation rate and transferability between different videos and classifier modules, while a worse performance on targeted attack.