

Supplementary Material for DUO

Xin Yao*, Yu Zhan*, Yimin Chen^{†§}, Fengxiao Tang*, Ming Zhao*, Enlang Li*, Yanchao Zhang^{†§}

*School of Computer Science and Engineering, Central South University, Changsha, Hunan, 410082, China

[†]Miner School of Computer & Information Sciences, University of Massachusetts Lowell, Lowell, MA 01854, USA

[‡]School of Electrical, Computer and Energy Engineering, Arizona State University, Tempe, AZ 85287, USA

[§]Corresponding authors

I. MORE DETAILS OF SURROGATE MODEL BUILDING

The architecture of the surrogate model is also a typical DNN (denoted by $\mathcal{S}(\cdot)$) with model parameters ρ . Similar to [1], $\mathcal{S}(\cdot)$ consists of a stacked convolution neural network and multiple fully connected layers to extract the spatial features of video frames and the temporal features of the video and outputs flattened and continuous video features. Note that $\mathcal{S}(\cdot)$ can be in other forms such as C3D [1] and our attack is applicable to them as well. Given $\mathcal{S}(\cdot)$, we move forward to denote the projected features of a given video $\mathbf{v} \in \mathbb{R}^{N \times W \times H \times C}$ by $\text{Fea}_\rho(\mathbf{v})$ and consequently the distance between the features of \mathbf{v} and \mathbf{v}_i by $D(\mathbf{v}, \mathbf{v}_i) = \|\text{Fea}_\rho(\mathbf{v}) - \text{Fea}_\rho(\mathbf{v}_i)\|_2^2$. The smaller $D(\mathbf{v}, \mathbf{v}_i)$ is, the higher similarity \mathbf{v} and \mathbf{v}_i share.

The training process of the surrogate model is very straightforward. Utilizing $\mathbf{T} = \{(\mathbf{v}, \mathbf{v}_i, \mathbf{v}_j) | \mathbf{v}_i, \mathbf{v}_j \in \mathcal{R}^n(\mathbf{v}), i < j \in [1, n]\}$ as the training set, we use the following loss function to update ρ , similar to [2]:

$$\arg \max_{\rho} \sum_{j>i} [D(\mathbf{v}, \mathbf{v}_j) - D(\mathbf{v}, \mathbf{v}_i) + \gamma]_+.$$

The intuition behind is to find the optimal ρ to maximize the feature distance between the query video \mathbf{v} and the returned video \mathbf{v}_j while minimizing that between the query video \mathbf{v} and the returned video \mathbf{v}_i . The parameter γ is a margin constant to prevent a negative loss value. Apparently, we can use the popular adaptive moment estimation (ADAM) [3] optimizer to obtain the updated ρ . In the end, we can obtain a trained surrogate model for the victim system.

II. MORE DETAILS OF AE GENERATION ON $\mathcal{S}(\cdot)$

We formulate AE generation on $\mathcal{S}(\cdot)$ as the following:

$$\begin{aligned} \arg \min_{\theta, \mathcal{I}, \mathcal{F}} \mathcal{L}(\text{Fea}_\rho(\mathbf{v}_{adv}), \text{Fea}_\rho(\mathbf{v}_t)) + \lambda \|\theta \odot \mathcal{I} \odot \mathcal{F}\|_2^2 \\ \text{s.t. } \mathbf{1}^\top \mathcal{I} = k, \|\mathcal{F}\|_{2,0} = n, \|\theta\|_\infty \leq \tau. \end{aligned} \quad (1)$$

Considering that θ is a continuous variable while \mathcal{I} and \mathcal{F} are binary ones, the problem in Eq. (1) is mixed integer programming [4], making it difficult to find the binary pixel mask \mathcal{I} and the binary frame mask \mathcal{F} . To tackle this problem, we follow [5] to loosen the binary constraint with two strong continuous constraints. For instance, the binary variable \mathcal{I} can be represented as $\mathcal{S}_b \cup \mathcal{S}_p$, where $\mathcal{S}_b = [0, 1]^{N \times W \times H \times C}$ is a box function and $\mathcal{S}_p = \{\mathcal{I} : \|\mathcal{I} - \frac{1}{2}\|_2^2 = \frac{N \times W \times H \times C}{4}\}$ is a l_2 -sphere constraint. Similarly, it is very challenging to

directly obtain the binary frame mask \mathcal{F} with the $\ell_{2,0}$ -norm constraint. To this end, we introduce a continuous variable $\mathcal{C} \in \mathbb{R}^{N \times W \times H \times C}$ to replace \mathcal{F} , obtain it with [6], and utilize an indicator function to select more sensitive frames. Eq. (1) is a non-convexity problem, so we cannot directly solve it (i.e., obtain the solution) but calculate the parameters $\mathcal{P} = \{\theta, \mathcal{F}, \mathcal{I}\}$ in an iterative fashion using the following steps until they converge.

Step 1: Given \mathcal{F} and \mathcal{I} , the sub-problem of \mathcal{P}_θ can be formulated as:

$$\begin{aligned} \arg \min_{\theta} \mathbb{L}(\mathcal{P}_\theta) = \mathcal{L}(\text{Fea}_\rho(\mathbf{v}_{adv}), \text{Fea}_\rho(\mathbf{v}_t)) + \lambda \|\theta \odot \mathcal{I} \odot \mathcal{F}\|_2^2 \\ \text{s.t. } \|\theta\|_\infty \leq \tau. \end{aligned} \quad (2)$$

This step is similar to a typical adversarial example attack [7] and therefore can be readily solved by the gradient descent method as

$$\begin{aligned} \theta = \theta - \mu_1 \cdot \frac{\mathbb{L}(\mathcal{P}_\theta)}{\partial \theta} \\ = \theta - \mu_1 \cdot \left[\frac{\partial \mathcal{L}(\text{Fea}_\rho(\mathbf{v}_{adv}), \text{Fea}_\rho(\mathbf{v}_t))}{\partial \theta} + 2\lambda \theta \odot (\mathcal{I} \odot \mathcal{F})^2 \right], \end{aligned}$$

where μ_1 is the update step size. To satisfy the constraint of $\|\theta\|_\infty \leq \tau$, we clip each perturbation magnitude of $\theta_{i,w,h,c}$ with the function of $\max(-\tau, \min(\theta_{i,w,h,c}, \tau))$ after each iteration.

Step 2: Given θ and \mathcal{F} , the sub-problem on $\mathcal{P}_\mathcal{I}$ can be formulated as:

$$\begin{aligned} \arg \min_{\mathcal{I}} \mathcal{L}(\text{Fea}_\rho(\mathbf{v}_{adv}), \text{Fea}_\rho(\mathbf{v}_t)) + \lambda \|\theta \odot \mathcal{I} \odot \mathcal{F}\|_2^2 \\ \text{s.t. } \mathcal{I} \in \{1, 0\}^{N \times W \times H \times C}, \mathbf{1}^\top \mathcal{I} = k. \end{aligned} \quad (3)$$

Since θ and \mathcal{F} are a continuous variable and a binary pixel mask, respectively, the optimization problem of Eq. (3) is also mixed integer programming. We follow [5] to replace the binary constraint \mathcal{I} with $\mathcal{I} = \mathbf{Y}_b \in \mathcal{S}_b$ and $\mathcal{I} = \mathbf{Y}_p \in \mathcal{S}_p$. As a result, the problem of Eq. (3) can be further formulated as

$$\begin{aligned} \arg \min_{\mathcal{I}} \mathcal{L}(\text{Fea}_\rho(\mathbf{v}_{adv}), \text{Fea}_\rho(\mathbf{v}_t)) + \lambda \|\theta \odot \mathcal{I} \odot \mathcal{F}\|_2^2 \\ \text{s.t. } \mathcal{I} = \mathbf{Y}_b, \mathcal{I} = \mathbf{Y}_p, \mathbf{1}^\top \mathcal{I} = k. \end{aligned} \quad (4)$$

The augmented Lagrangian function for the objective function in Eq. (4) is given by

$$\begin{aligned} \mathbb{L}(\mathcal{P}_{\mathcal{I}}) = & \mathcal{L}(\text{Fea}_{\rho}(\mathbf{v}_{adv}), \text{Fea}_{\rho}(\mathbf{v}_t)) + \lambda \|\boldsymbol{\theta} \odot \mathcal{I} \odot \mathcal{F}\|_2^2 \\ & + \mathbf{w}_1^\top (\mathcal{I} - \mathbf{Y}_b) + \mathbf{w}_2^\top (\mathcal{I} - \mathbf{Y}_p) + \mathbf{w}_3^\top (\mathbf{1}^\top \mathcal{I} - k) \\ & + \frac{\rho_1}{2} \|\mathcal{I} - \mathbf{Y}_b\|_2^2 + \frac{\rho_2}{2} \|\mathcal{I} - \mathbf{Y}_p\|_2^2 + \frac{\rho_3}{2} (\mathbf{1}^\top \mathcal{I} - k)^2 \\ & + h_1(\mathbf{Y}_b) + h_2(\mathbf{Y}_p), \end{aligned} \quad (5)$$

where $\{\mathbf{w}_1, \mathbf{w}_2\} \in \mathbb{R}^{N \times W \times H \times C}$, $\mathbf{w}_3 \in \mathbb{R}^N$ are dual variables, and $\{\rho_1, \rho_2, \rho_3\}$ are penalty parameters. $h_1(\mathbf{Y}_b) = \mathbb{I}_{\{\mathbf{Y}_b \in \mathcal{S}_b\}}$ and $h_2(\mathbf{Y}_p) = \mathbb{I}_{\{\mathbf{Y}_p \in \mathcal{S}_p\}}$ are “indicator functions”, meaning $\mathbb{I}_{\{a\}} = 0$ if a is true and $+\infty$ otherwise. We follow the general procedure of ADMM to update the above dual variables and primal variables iteratively as follows.

Step 2.1: We independently update \mathbf{Y}_b and \mathbf{Y}_p in parallel as

$$\begin{cases} \mathbf{Y}_b^{t+1} = \arg \min_{\mathbf{Y}_b \in \mathcal{S}_b} \mathbf{w}_1^\top (\mathcal{I} - \mathbf{Y}_b) + \frac{\rho_1}{2} \|\mathcal{I} - \mathbf{Y}_b\|_2^2 \\ \quad = \mathcal{P}_b(\mathcal{I} + \frac{1}{\rho_1} \mathbf{w}_1), \\ \mathbf{Y}_p^{t+1} = \arg \min_{\mathbf{Y}_p \in \mathcal{S}_p} \mathbf{w}_2^\top (\mathcal{I} - \mathbf{Y}_p) + \frac{\rho_2}{2} \|\mathcal{I} - \mathbf{Y}_p\|_2^2 \\ \quad = \mathcal{P}_p(\mathcal{I} + \frac{1}{\rho_2} \mathbf{w}_2), \end{cases} \quad (6)$$

where $\mathcal{P}_b(\mathbf{a}) = \min(\mathbf{1}, \max(\mathbf{0}, \mathbf{a}))$ with $\mathbf{a} \in \mathbb{R}^{N \times W \times H \times C}$ denotes the projection onto the box constraint \mathcal{S}_b , and $\mathcal{P}_p(\mathbf{a}) = \frac{\sqrt{N \times W \times H \times C}}{2} \frac{\bar{\mathbf{a}}}{\|\bar{\mathbf{a}}\|} + \frac{1}{2} \mathbf{1}$ with $\bar{\mathbf{a}} = \mathbf{a} - \frac{1}{2} \mathbf{1}$ and $\mathbf{a} \in \mathbb{R}^{N \times W \times H \times C}$ denotes the projection onto the ℓ_2 -sphere constraint. Both of the constraints have been proved in [5].

Step 2.2: Given $\mathbf{Y}_b, \mathbf{Y}_p, \mathbf{w}_1, \mathbf{w}_2$, and \mathbf{w}_3 , we update \mathcal{I} by gradient descent as below:

$$\begin{aligned} \mathcal{I} &= \mathcal{I} - \mu_2 \cdot \frac{\mathbb{L}(\mathcal{P}_{\mathcal{I}})}{\partial \mathcal{I}} \\ &= \mathcal{I} - \mu_2 \cdot \left[\frac{\partial \mathcal{L}(\text{Fea}_{\rho}(\mathbf{v}_{adv}), \text{Fea}_{\rho}(\mathbf{v}_t))}{\partial \mathcal{I}} + 2\lambda \mathbf{I} \odot (\boldsymbol{\theta} \odot \mathcal{F})^2 \right. \\ &\quad + (\mathbf{w}_1 + \rho_1(\mathcal{I} - \mathbf{Y}_b)) + (\mathbf{w}_2 + \rho_2(\mathcal{I} - \mathbf{Y}_p)) \\ &\quad \left. + (\mathbf{w}_3 + \rho_3(\mathbf{1}^\top \mathcal{I} - k)) \right], \end{aligned} \quad (7)$$

where μ_2 is the update step size.

Step 2.3: Given $\mathcal{I}, \mathbf{Y}_b$, and \mathbf{Y}_p , we update the dual variables $\mathbf{w}_1, \mathbf{w}_2$, and \mathbf{w}_3 using the general ascent method as follows:

$$\begin{cases} \mathbf{w}_1^{t+1} = \mathbf{w}_1^t + \rho_1(\mathcal{I}^{t+1} - \mathbf{Y}_b^{t+1}) \\ \mathbf{w}_2^{t+1} = \mathbf{w}_2^t + \rho_2(\mathcal{I}^{t+1} - \mathbf{Y}_p^{t+1}) \\ \mathbf{w}_3^{t+1} = \mathbf{w}_3^t + \rho_3(\mathbf{1}^\top \mathcal{I}^{t+1} - k) \end{cases} \quad (8)$$

Step 3: Given $\boldsymbol{\theta}$ and \mathcal{I} , the sub-problem of $\mathcal{P}_{\mathcal{F}}$ can be formulated below:

$$\begin{aligned} \arg \min_{\mathcal{F}} & \mathcal{L}(\text{Fea}_{\rho}(\mathbf{v}_{adv}), \text{Fea}_{\rho}(\mathbf{v}_t)) + \lambda \|\boldsymbol{\theta} \odot \mathcal{I} \odot \mathcal{F}\|_2^2 \\ \text{s.t.} & \mathcal{F} \in \{0, 1\}^{N \times W \times H \times C}, \|\mathcal{F}\|_{2,0} = n. \end{aligned} \quad (9)$$

Similarly, it is very challenging to obtain the binary frame mask \mathcal{F} . Thus, we introduce a continuous variable $\mathcal{C} \in \mathbb{R}^{N \times W \times H \times C}$ to replace \mathcal{F} and formulate Eq. (9) as

$$\begin{aligned} \arg \min_{\mathcal{C}} & \mathcal{L}(\text{Fea}_{\rho}(\mathbf{v}_{adv}), \text{Fea}_{\rho}(\mathbf{v}_t)) + \lambda \|\boldsymbol{\theta} \odot \mathcal{I} \odot \mathcal{C}\|_2^2 \\ \text{s.t.} & \mathcal{C} \in \mathbb{R}^{N \times W \times H \times C}, \|\mathcal{C}\|_{2,0} = n. \end{aligned}$$

\mathcal{C} can be easily solved following [6]. Subsequently, we update \mathcal{F} as

$$\mathcal{F}_{\pi(i)} = \begin{cases} \mathbf{1}, i \leq n, \\ \mathbf{0}, i > n, \end{cases}$$

where $\mathcal{F}_{\pi(i)}$ is the $\pi(i)$ -th row of \mathcal{F} and $\pi(\cdot)$ is an indicator function such that $\|\mathcal{C}_{\pi(1)}\|_2 \geq \|\mathcal{C}_{\pi(2)}\|_2 \geq \dots \geq \|\mathcal{C}_{\pi(N)}\|_2$. Note that $\|\mathcal{C}_{\pi(i)}\| \geq \|\mathcal{C}_{\pi(j)}\|$ means that the $\pi(i)$ -th frame is more sensitive than the $\pi(j)$ -th frame for the same amount of perturbations and thus should be selected with a higher probability.

III. MORE DETAILS OF $\mathbb{H}(\cdot, \cdot)$ AND \mathbf{Q}

To evaluate the similarities of two retrieval lists, e.g., $\mathcal{R}^m(\mathbf{v}) = \{\mathbf{v}_1, \dots, \mathbf{v}_i, \dots, \mathbf{v}_m\}$ and $\mathcal{R}^m(\mathbf{v}') = \{\mathbf{v}'_1, \dots, \mathbf{v}'_j, \dots, \mathbf{v}'_m\}$, we introduce the NDCG-based function [2] $\mathbb{H}(\mathcal{R}^m(\mathbf{v}), \mathcal{R}^m(\mathbf{v}'))$. In effect, for each retrieval video $\mathbf{v}_i \in \mathcal{R}^m(\mathbf{v})$, we first calculate the prior sampling probability ω_i that \mathbf{v}_i appears in $\mathcal{R}^m(\mathbf{v})$ as:

$$\omega_i = \frac{\log r_i - 1}{\sum_{i=1}^m (\log r_i - 1)},$$

where r_i represents the relevance between \mathbf{v}_i and \mathbf{v} and is equivalent to $m - i + 1$. Note that the smaller i , the higher the similarity between \mathbf{v}_i and \mathbf{v} . Subsequently, for each retrieval video $\mathbf{v}_i \in \mathcal{R}^m(\mathbf{v})$, we calculate the conditioned attack failure probability ϑ_i as:

$$\vartheta_i = \begin{cases} \omega_j, & \mathbf{v}_i \in \mathcal{R}^m(\mathbf{v}') \text{ and } \mathbf{v}_i = \mathbf{v}'_j; \\ 0, & \mathbf{v}_i \notin \mathcal{R}^m(\mathbf{v}'). \end{cases}$$

Therefore, the similarity function $\mathbb{H}(\mathcal{R}^m(\mathbf{v}), \mathcal{R}^m(\mathbf{v}'))$ can be defined as:

$$\sum_{i=1}^m \omega_i \cdot \vartheta_i.$$

Note that the more $\mathcal{R}^m(\mathbf{v})$ and $\mathcal{R}^m(\mathbf{v}')$ intersect, the larger $\mathbb{H}(\mathcal{R}^m(\mathbf{v}), \mathcal{R}^m(\mathbf{v}'))$, and vice versa.

In addition, for the κ -th iteration, the random matrix $\mathbf{q}^{(\kappa)} \in \mathbb{R}^{N, W, H, C}$ is selected from an orthonormal candidate vector set (e.g., the Cartesian basis) without replacement to ensure that the vector $\mathbf{q}^{(\kappa)}$ will not be cancelled out across different iterations.

REFERENCES

- [1] D. Tran, L. Bourdev, R. Fergus, L. Torresani, and M. Paluri, “Learning spatiotemporal features with 3d convolutional networks,” in *ICCV’15*, Santiago, Chile, Dec. 2015.
- [2] X. Li, J. Li, Y. Chen, S. Ye, Y. He, S. Wang, H. Su, and H. Xue, “Qair: Practical query-efficient black-box attacks for image retrieval,” in *CVPR’21*, Nashville, TN, June 2021.
- [3] D. Kingma and J. Ba, “Adam: A method for stochastic optimization,” in *ICLR’15*, San Diego, CA, May 2015.

- [4] B. Fan, Y. Wu, H. Li, Y. Zhang, Y. Li, F. Li, and J. Yang, "Sparse adversarial attack via perturbation factorization," in ECCV'20, Glasgow, UK, Aug. 2020.
- [5] B. Wu and B. Ghanem, " ℓ_p -box ADMM: A Versatile Framework for Integer Programming," IEEE Transactions on Pattern Analysis and Machine Intelligence, vol. 41, no. 7, pp. 1695–1708, Jan. 2018.
- [6] W. Gu, X. Gu, J. Gu, B. Li, Z. Xiong, and W. Wang, "Dependence guided unsupervised feature selection," in AAAI'18, New Orleans, LA, Feb. 2018.
- [7] N. Carlini and D. Wagner, "Towards evaluating the robustness of neural networks," in S&P'17, San Jose, CA, Jun. 2017.