Compositional Targeted Multi-Label Universal Perturbations

Supplementary Material

7. Derivation of CMLU_{α} and CMLU_{β}

In this section, we provide more details of the derivation from the main paper. In particular, we focus on deriving Eq. (7). We begin by considering Eq. (5)

$$\max_{\boldsymbol{U}} \sum_{\boldsymbol{\Omega} \in \mathcal{P}(\mathcal{C}) k \in \Omega} \sum_{k \in \Omega} \left(\mathcal{H}_{k}(\boldsymbol{u}_{k} + \boldsymbol{U}\boldsymbol{\lambda}_{\Omega \setminus k}) - \frac{\gamma}{|\Omega|} \sum_{q \in \bar{\Omega}} \mathcal{R}_{q}(\boldsymbol{u}_{k} + \boldsymbol{U}\boldsymbol{\lambda}_{\Omega \setminus k}) \right),$$

s.t. $\hat{\boldsymbol{U}}^{\top} \hat{\boldsymbol{U}} = \boldsymbol{I}, \|\boldsymbol{U}\|_{1,p} \leq \epsilon$ (12)

where \hat{U} is column-normalized U and $||U||_{1,p} = \max_{1 \le i \le |\mathcal{C}|} ||u_i||_p$, is the maximum norm of the columns of U. Notice that we wrote $u_{\Omega} = u_k + U\lambda_{\Omega\setminus k}$ for a k that belongs to Ω , where $\Omega\setminus k$ denotes removing k from Ω . The sum in Eq. (12) is across the power set of \mathcal{C} .

7.1. Powerset Reformulation

In our derivation, we have two key steps (corresponding to two observations) that involve expressing the powerset of a set in a specific way. We explain those steps below.

Observation 1: For any $k \in C$, the subsets of powerset $\mathcal{P}(C)$ that include k can be written as:

$$\mathcal{P}_{k} \triangleq \bigcup_{\Omega \in \mathcal{P}(\mathcal{C} \setminus k)} \left\{ \{k\} \cup \Omega \right\} \triangleq \left\{ \{k\} \cup \Omega \mid \Omega \in \mathcal{P}(\mathcal{C} \setminus \{k\}\}$$
(13)

 \mathcal{P}_k in Eq. (13) gives all subsets of \mathcal{C} that include k, constructed by taking all the possible subsets of $\mathcal{C} \setminus \{k\}$ and adding k to them. Based on this, the key observation is that we can combine these contributions (\mathcal{P}_k) from all $k \in \mathcal{C}$ to reconstruct the entire powerset $\mathcal{P}(\mathcal{C})$ (excluding the empty set) as

$$\mathcal{P}(\mathcal{C}) \triangleq \bigcup_{k \in \mathcal{C}} \mathcal{P}_k \triangleq \bigcup_{k \in \mathcal{C}} \bigcup_{\Omega \in \mathcal{P}(\mathcal{C} \setminus k)} \left\{ \{k\} \cup \Omega \right\}$$
(14)

Eq. (14) allows us to leverage the property of decomposing the powerset into a union of sets. This structure reveals the inherent symmetry and compositional nature of the powerset. Therefore, using Eq. (14), we can reconstruct the summation across the power set in Eq. (12). However, (14) introduced redundancy, as each subset is repeated $|\Omega| + 1$ time in the formulation. To account for this, one could multiply the loss expression by a factor of $1/(|\Omega|+1)$, which just scales the loss. So instead, we directly optimize

$$\max_{\boldsymbol{U}} \sum_{k \in \mathcal{C}} \sum_{\Omega \in \mathcal{P}(\mathcal{C} \setminus k)} \left(\mathcal{H}_{k}(\boldsymbol{u}_{k} + \boldsymbol{U}\boldsymbol{\lambda}_{\Omega \setminus k}) - \eta \sum_{q \in \{\bar{\Omega} \setminus k\}} \mathcal{R}_{q}(\boldsymbol{u}_{k} + \boldsymbol{U}\boldsymbol{\lambda}_{\Omega \setminus k}) \right),$$

s.t. $\hat{\boldsymbol{U}}^{\top} \hat{\boldsymbol{U}} = \boldsymbol{I}, \|\boldsymbol{U}\|_{1,p} \leq \epsilon,$ (15)

where $\eta = \gamma/|\Omega|$. Note the change in the indexing of the second loss term in (15). This adjustment ensures consistency with (12) by enforcing that the sum of non-targeted loss excludes k in our new formulation.

Observation 2: The power set of C can be written as the union of subsets of C of fixed size:

$$\mathcal{P}(\mathcal{C}) \triangleq \bigcup_{t=0}^{|\mathcal{C}|} \binom{\mathcal{C}}{t} \triangleq \bigcup_{t=0}^{|\mathcal{C}|} \bigcup_{S \in \binom{\mathcal{C}}{t}} \{S\}$$
(16)

where, each $S \in \binom{\mathcal{C}}{t}$ is a subset of \mathcal{C} of size t.

In our case, any $S \in \binom{\mathcal{C}}{t}$ would represent the indices of the columns of U. Using (16), we write the summation across the power set $\mathcal{P}(\mathcal{C}\setminus k)$ in (15) as

$$\max_{\boldsymbol{U}} \sum_{k \in \mathcal{C}} \sum_{t=0}^{|\mathcal{C}|-1} \sum_{S \in \binom{\mathcal{C} \setminus k}{t}} \left(\mathcal{H}_{k}(\boldsymbol{u}_{k} + \boldsymbol{U}\boldsymbol{\lambda}_{S}) - \eta \sum_{q \in \{\bar{S} \setminus k\}} \mathcal{R}_{q}(\boldsymbol{u}_{k} + \boldsymbol{U}\boldsymbol{\lambda}_{S}) \right),$$

s.t. $\hat{\boldsymbol{U}}^{\top} \hat{\boldsymbol{U}} = \boldsymbol{I}, \|\boldsymbol{U}\|_{1,p} \leq \epsilon,$ (17)

Please refer to the main paper for the rest of the derivation.

8. Analysis of Universal Perturbations

In Fig. 6, we visualize label-wise universal perturbations generated using the Oracle, $CMLU_{\alpha}$, and $CMLU_{\beta}$ methods on the NUS-WIDE dataset for the ASL model, along with some of their compositions. Note that the perturbations produced by the Oracle method exhibit significant texture variations across labels, and combining these perturbations often causes the key adversarial features of individual classes to diminish or fade away. In contrast, $CMLU_{\alpha}$ and $CMLU_{\beta}$ learn label-wise perturbations with subtle variations, enabling their composition to preserve shared features across labels while de-emphasizing conflicting textures. This demonstrates the effectiveness of $CMLU_{\beta}$ for attacking multiple classes.

In Tab. 7, we show the pairwise dot products of universals computed using Oracle and CMLU_{β} . The purpose of this experiment is to evaluate the difference between the vulnerability directions found through each of the methods. From the table, we can see that CMLU_{β} finds unique directions, different from Oracle.

In Tab. 4 in the main paper, we showed the fooling success rate to evaluate the transferability of universal attacks across model. In Tab. 5, we show non-target flip rate (NT_R) values. From both tables, we can see that when CMLU_{β}

Table 5. Non-target flip rate (NT_R) of Transferability Experiments using the ASL model.

Transferability of Universal Attacks											
		N	US-WI	DE	MS-COCO						
Models	Method	1	3	5	1	3	5				
$\mathbf{A} ightarrow \mathbf{D}$	$\begin{array}{c} \mathbf{Oracle} \\ \mathbf{CMLU}_{\beta} \end{array}$	3.81 4.21	3.88 3.80	3.65 3.29	5.29 7.39	8.03 8.39	9.45 9.04				
$\mathbf{D} ightarrow \mathbf{A}$	$\begin{array}{c} \mathbf{Oracle} \\ \mathbf{CMLU}_{\beta} \end{array}$	4.22 6.02	4.86 0.85	-	7.54 10.6	9.07 10.9	10.07 13.3				

achieves better FR than Oracle, it also gets higher NT_R than Oracle. In the case, when it performs worse than Oracle, it has lower NT_R .

9. Wall Clock Time

In Tab. 6, we present the computational complexity and wall-clock times for each method, illustrating their practical efficiency. It is important to highlight that the total time column (T) in the table reflects only a few target sizes (1, 3, 5) and accounts for only a small subset $S(|\Omega|)$ of all possible combinations of a fixed size $|\Omega|$. For methods like Oracle and Or-C, these times would increase significantly with larger target sizes.

From the table, we make the following observations:

- The computational time of Oracle and Or-C continues to increases as additional target sizes are considered, making them less practical for large datasets and target sizes. In contrast, our methods, CMLU_α and CMLU_β, require training only for |Ω| = 1, and their computational costs remain constant regardless of the number of target sizes.
- Note that CMLU_α is more computationally expensive than CMLU_β. This cost is due to the inner minimization step involving N_α iterations.
- Or-S and CMLU_β have the same time complexity and wall clock time. But as we have shown in the main paper, CMLU_β achieves performance close to exponential time Oracle method while Or-S scales poorly to large target sizes.

10. Hyperparameters and Experiment Setting

In algorithm 1, we mainly have five hyperparameters: $\epsilon, \eta = \frac{\kappa}{|\Omega|}, \alpha, N_{\alpha}$, and ξ . In our experiments, we set $\kappa = 20$ for NUS-WIDE and OpenImages experiments and $\kappa = 10$ for MS-COCO. In all of our main experiments, we use $N_{\alpha} = 10, \alpha = 1e2, \epsilon = 0.05$, and $\xi = 0.002$. We set the random seed for pytorch and numpy to 999. We use the officially provided train/val split to train the models and test split to evaluate the attacks. We used batch size of 40 in our experiments. All images were resized and center-cropped to 448 x 448.

11. Generating Convex Subcones

The purpose of Fig. 5 in the main paper is to demonstrate why our assumption (4) works with CMLU but fails with other methods. The figure compares the input regions that Oracle and CMLU_{β} have learnt. In the figure, Oracle plot shows the class-specific universals learnt using Oracle for $|\Omega| = 1$.

Consider a multi-label classifier $\mathcal{F} : \mathbb{R}^d \to \mathbb{R}^{|\mathcal{C}|}$ where \mathcal{C} is the set of all labels. The classifier $\mathcal{F} = \{\mathcal{F}_1, \mathcal{F}_2, \dots, \mathcal{F}_{|\mathcal{C}|}\}$ consists of $|\mathcal{C}|$ binary classifiers (one for each label), where $\mathcal{F}_c(\mathbf{x}) \in (-\infty, +\infty)$ represents the *logit* of label c.

- Assume U_{Or} represents $|\mathcal{C}|$ universals learnt using Oracle Eq. (1), where each universal vector $u_k \in U_{\text{Or}}$ attacks the specific class k.
- Assume U_{β} represents $|\mathcal{C}|$ universals learnt using CMLU_{β} Eq. (10), where each universal vector $u_k \in U_{\beta}$ attacks the specific class k.

First, we plot the average model confidence for all possible pairs of vectors from U_{Or} using (18) and then we plot the average model confidence for all possible pairs of vectors from U_{β} . We consider a set $Q = \binom{C}{2}$ i.e., all possible pairs of classes. Assume \mathcal{I}_{Ω} are the clean images that contain classes Ω . Then, the value at each point (i, j) in the 2D plot in Fig. 5 is computed as:

$$\mathbb{C}(i,j) = \frac{1}{|\mathcal{Q}|} \sum_{\Omega \in \mathcal{Q}} \frac{1}{|\mathcal{I}_{\Omega}|} \sum_{x \in \mathcal{I}_{\Omega}} \frac{1}{|\Omega|} \sum_{k \in \Omega} \sigma(\mathcal{F}_k(x + v_{\Omega}(i,j)))$$
(18)

where, $v_{\Omega}(i,j) = \prod_{\epsilon} (i * u_{\Omega_1} + j * u_{\Omega_2})$

for $i, j \in [0, 1]$. Π_{ϵ} projects the input on ϵ infinity norm ball and σ is the sigmoid function. Note that the value $\mathbb{C}(i, j)$ is computed across all target classes Q and all images \mathcal{I}_{Ω} for a point (i, j). This renders the cross-section of the input space to visualize the universal adversarial regions. Note that the vectors are not mutually sparse. We clip the cone at 0.05 infinity norm.

While summing universals learnt using (1) fails to successfully attack (no red regions), $CMLU_{\beta}$ learns universals whose sums effectively attack target labels (red regions). This validates that our framework can learn compositional universal perturbations. Furthermore, the figure provides intuitive 2D insights into MLL model behavior.

12. Limitation

A limitation of our work is the assumption of label independence. In MLL, several labels can be correlated, such as "Cat and Animal" or "Beach and Ocean". This correlation results in a tightly coupled FR and NT_R (reducing NT_R would also reduce FR e.g., if we want to attack Cat but keep Animal unchanged). This was also shown in our ablation experiments. Similarly, for high *FR*, we see an increased NT_R for larger target sizes. To address this, future

Table 6. Complexity and wall clock time of the methods. We show the approximate average time (in hours) for a single training iteration of the methods on an NVIDIA-V100 32GB GPU. The values are rounded off to the nearest integer. For a given target set size $|\Omega|$, $S(|\Omega|)$ is the number of selected sets or label combinations of size $|\Omega|$. The column **T** is the sum of the columns ($|\Omega| = 1, 3, 5$) to show the total compute time across different targets.

	Compl	NUS-WIDE				MS-COCO				OpenImages				
$ \Omega $	Space Time		1	3	5	Т	1	3	5	Т	1	3	5	Т
$\mathcal{S}(\Omega)$			25	40	40	105	40	50	50	140	100	100	100	300
Oracle	$\mathcal{O}(d \ 2^{ C })$	$\mathcal{O}(E \ 2^{ C })$	11	18	18	47	18	23	23	64	47	47	47	141
Or-S	$\mathcal{O}(d C)$	$\mathcal{O}(E C)$	11	-	-	11	18	-	-	18	47	-	-	47
Or-C	$\mathcal{O}(d C +2^{ C })$	$\mathcal{O}(E \; 2^{ C })$	11	17	17	45	18	22	22	62	47	43	43	133
SGA	$\mathcal{O}(d C)$	$\mathcal{O}(E C)$	23	-	-	23	31	-	-	31	112	-	-	112
NAG	$\mathcal{O}(G_p)$	$\mathcal{O}(E C)$	23	-	-	23	37	-	-	37	92	-	-	92
\mathbf{CMLU}_{α}	$\mathcal{O}(d C)$	$\mathcal{O}(EN_{\alpha} C)$	28	-	-	28	38	-	-	38	137	-	-	137
\mathbf{CMLU}_{β}	$\mathcal{O}(d C)$	$\mathcal{O}(E C)$	11	-	-	11	18	-	-	18	48	-	-	48

Table 7. The table shows mean dot product values of the normalized perturbation vectors computed across all target class combinations for each target size $|\Omega|$. The perturbations are learnt on NUS-WIDE, using ASL and ML-Decoder models.

			Tar	get 1			Tar	get 3		Target 5				
		ASL		ML-Dec		ASL		ML-Dec		ASL		ML-Dec		
		Oracle	\mathbf{CMLU}_{β}	Oracle	\mathbf{CMLU}_{β}	Oracle	\mathbf{CMLU}_{β}	Oracle	\mathbf{CMLU}_{β}	Oracle	\mathbf{CMLU}_{β}	Oracle	\mathbf{CMLU}_{β}	
ASL	Oracle	1.0	0.005	0.001	0.0	1.0	0.004	0.001	0.0	1.0	0.004	0.002	0.0	
	\bigcup_{β}	0.005	1.0	0.001	0.0	0.004	1.0	0.0	0.0	0.004	1.0	0.001	0.0	
ML-Dec	$\begin{array}{c c} \mathbf{Oracle} \\ \mathbf{CMLU}_{\beta} \end{array}$	0.002 0.0	0.001 0.0	1.0 0.005	0.005 1.0	0.001 0.0	0.0 0.0	1.0 0.003	0.003 1.0	0.002	0.001 0.0	1.0 0.004	0.004 1.0	



Figure 6. Visualizing the universal perturbations and their compositions, generated using various methods on the NUS-WIDE dataset for the ASL model. The arrows show which label-wise perturbations have been composed together to construct a new universal vector.

works could investigate the integration of label correlation and using a soft loss penalty on non-targeted labels.