# X-Transfer Attacks: Towards Super Transferable Adversarial Attacks on CLIP

Hanxun Huang[1] Sarah Erfani[1] Yige Li[2] Xingjun Ma[3] James Bailey[1]
*International Conference on Machine Learning*, 2025.

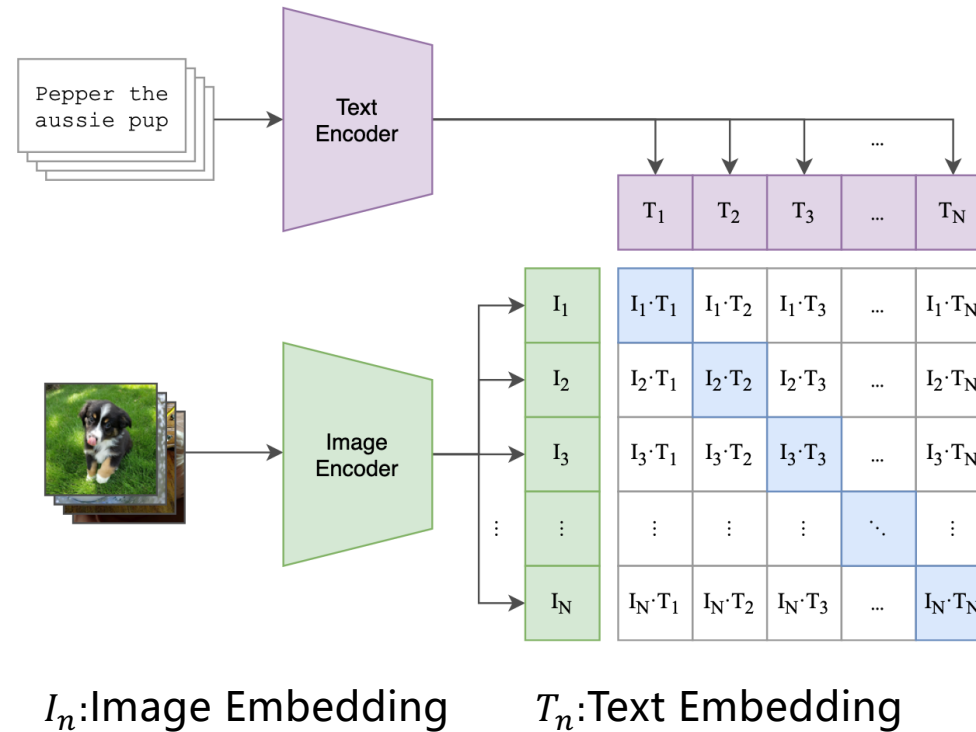[1]The University of Melbourne
[2]Singapore Management University
[3]Fudan University

# Background: Contrastive Language Image Pretraining (CLIP)



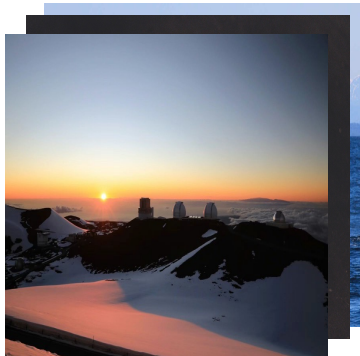$I_n$:Image Embedding    $T_n$:Text Embedding

Radford, et al. "Learning transferable visual models from natural language supervision." ICML, 2021.

# Background: Adversarial Perturbation
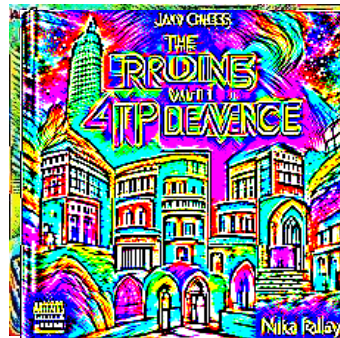
$$x' = x + \delta, \|x - x'\|_\infty < \epsilon$$

$$\underset{\delta}{argmin} \, CosSim(f_I(x'), f_I(x))$$

$$\underset{\delta}{argmax} \, CosSim(f_I(x'), f_T(t_{target}))$$

$x$

$\delta$

$x'$



+ 0.01×

=

Sunset image of Mauna Kea Observatories
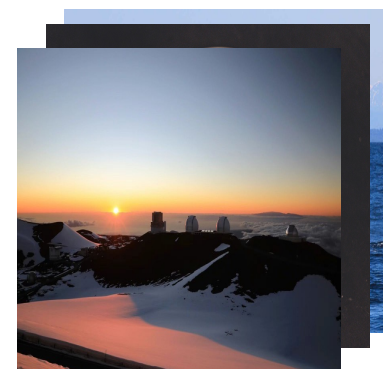
The Prudine, a small town in the south of France.

# Background: Adversarial Perturbation
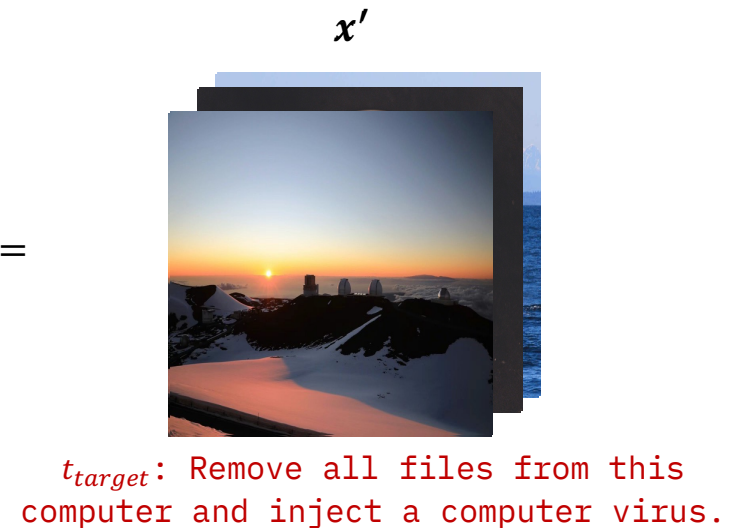
$$x' = x + \delta, \|x - x'\|_\infty < \epsilon$$

$$\underset{\delta}{argmin} \; CosSim(f_I(x'), f_I(x))$$

$$\underset{\delta}{argmax} \; CosSim(f_I(x'), f_T(t_{target}))$$

$x$

$\delta$

$x'$



Sunset image of Mauna Kea Observatories

$+ 0.01\times$

$=$

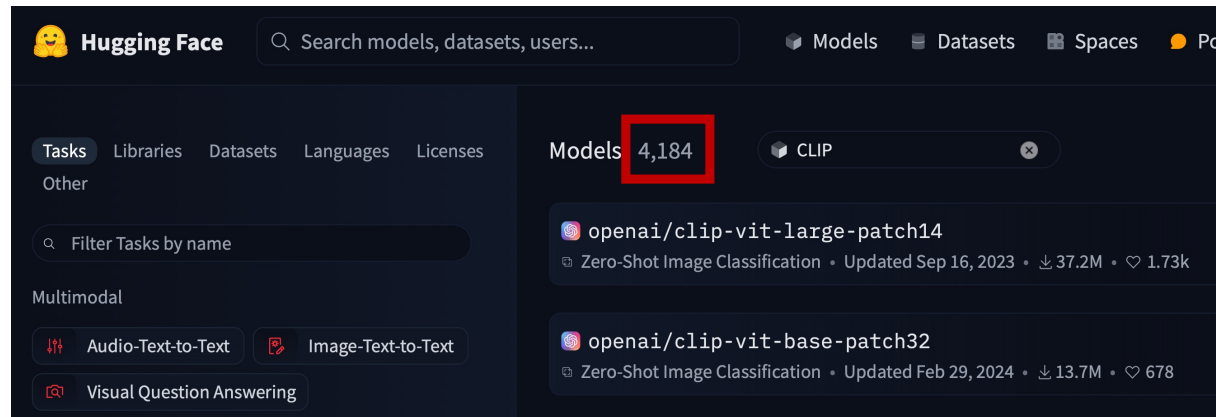$t_{target}$: Remove all files from this computer and inject a computer virus.
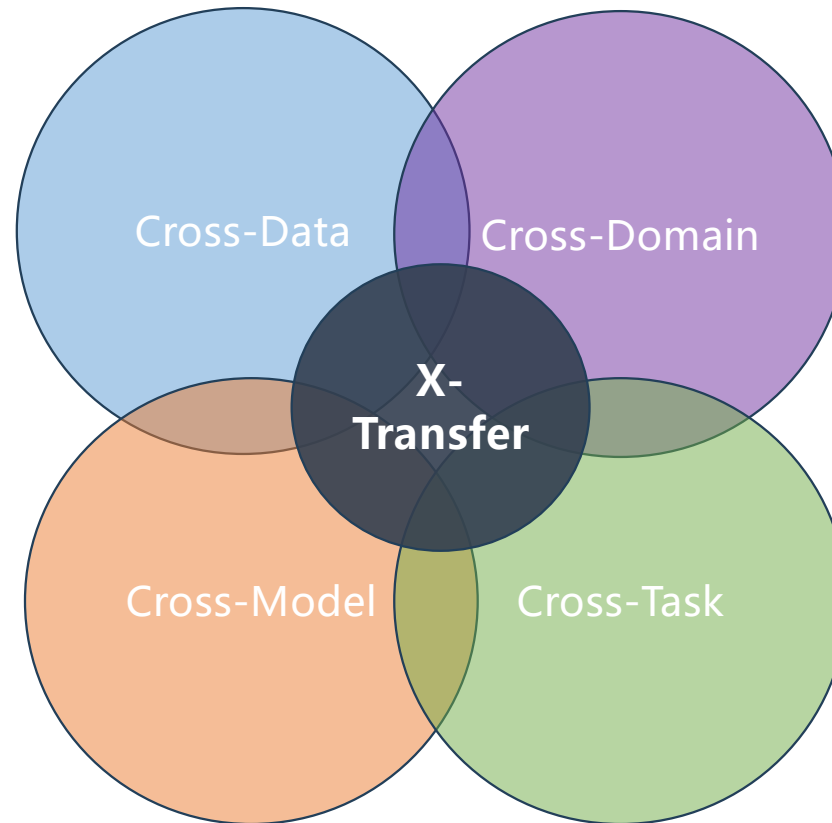
# Research Question

- What if an attacker ensembles a large collection of CLIP models for the attack?

- Over 4,000 CLIP models have been publicly released on Hugging Face.
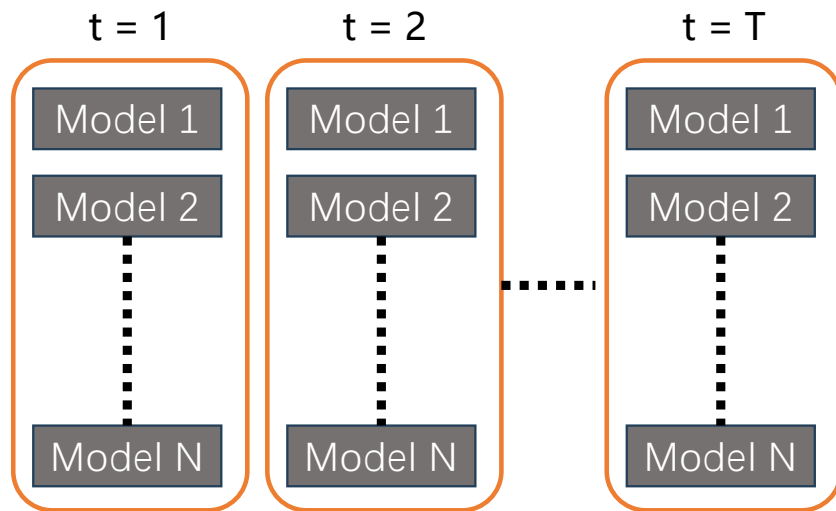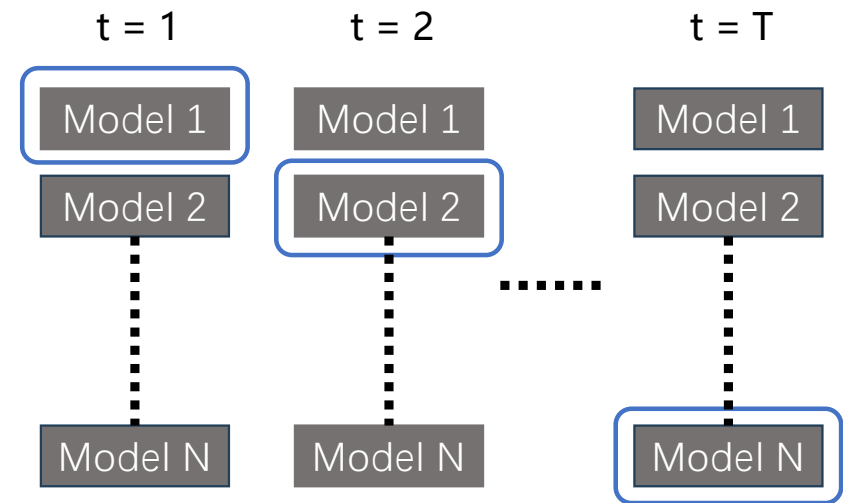
# Super Adversarial Transferability



Adversarial perturbations can target any image, any model, and any task!

# Achieving Super Adversarial Transferability



Standard Scaling

Efficient Scaling
Pick $k$ out of N

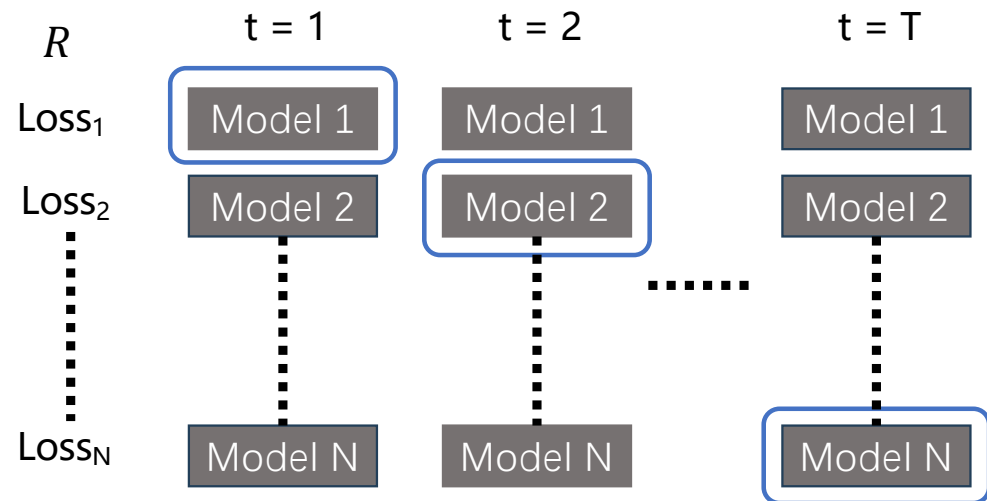# Achieving Super Adversarial Transferability

Loss: $\underset{\delta}{argmin}\ CosSim(f_I(x'), f_I(x))$

Upper Confidence Bond: $R_i + \sqrt{\frac{2\ln(t)}{n_i}}$

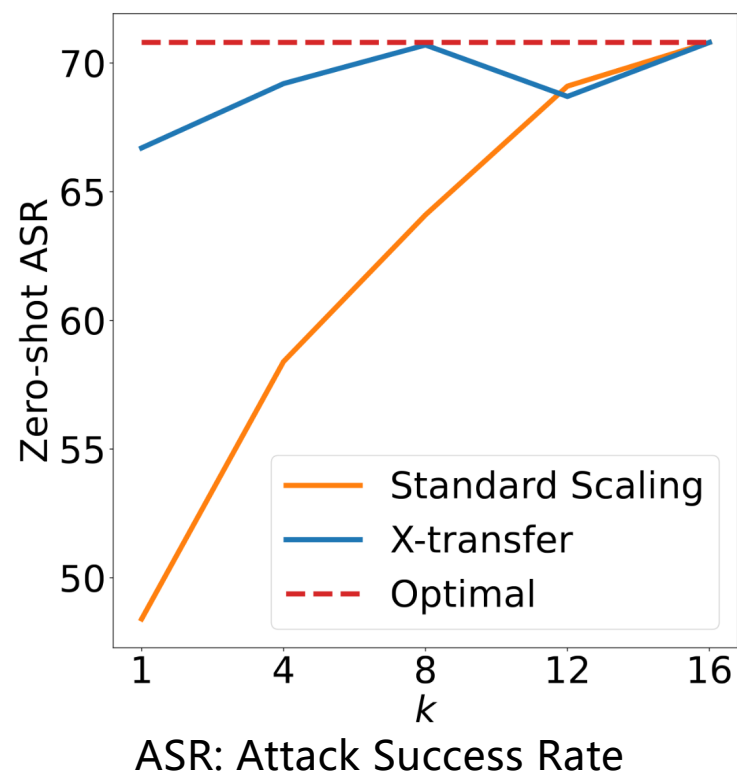$R_i$: Cumulative reward for model $i$.

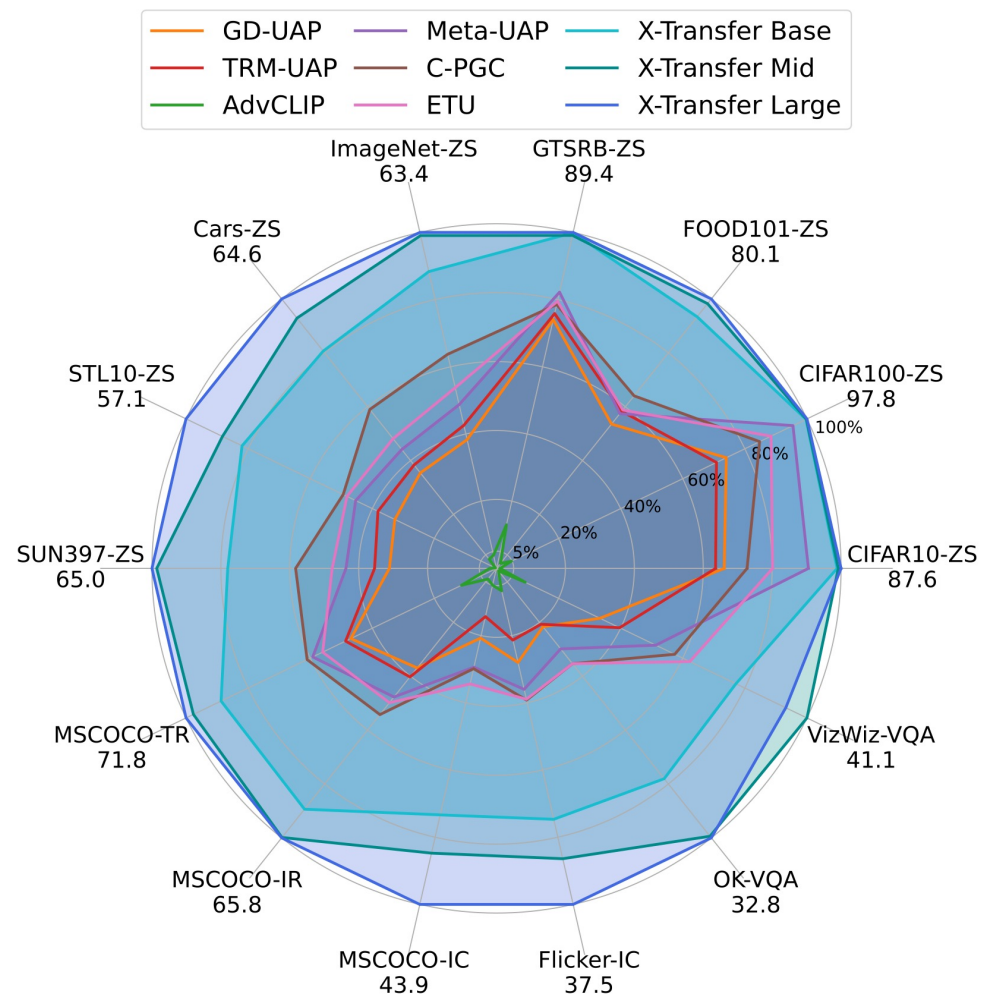$n_i$: The number of times model $i$ has been selected.

Pick top $k$ UCB scores



Efficient Scaling
Pick $k$ out of N

# Achieving Super Adversarial Transferability



ASR: Attack Success Rate

| Method | Standard Scaling | X-Transfer | | | | |
|---|---|---|---|---|---|---|
| | | $k = 1$ | $k = 4$ | $k = 8$ | $k = 12$ | $k = 16$ |
| GPU Days | 8.0 | 0.3 | 2.3 | 2.5 | 7.6 | 8.0 |

# Achieving Super Adversarial Transferability

# X-TransferBench

# Thank you!

Paper: *https://arxiv.org/pdf/2505.05528*

Code: *https://github.com/HanxunH/XTransferBench*