

BRUSLEATTACK: A QUERY-EFFICIENT SCORE- BASED BLACK-BOX SPARSE ADVERSARIAL ATTACK

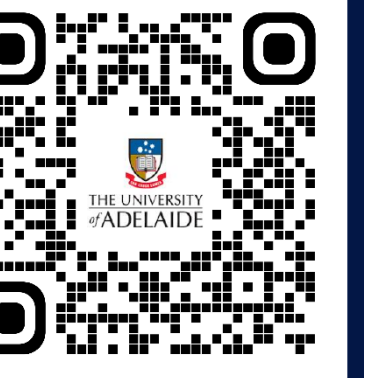


Viet Quoc Vo, Ehsan Abbasnejad, Damith C. Ranasinghe



ICLR
2024

SCAN ME

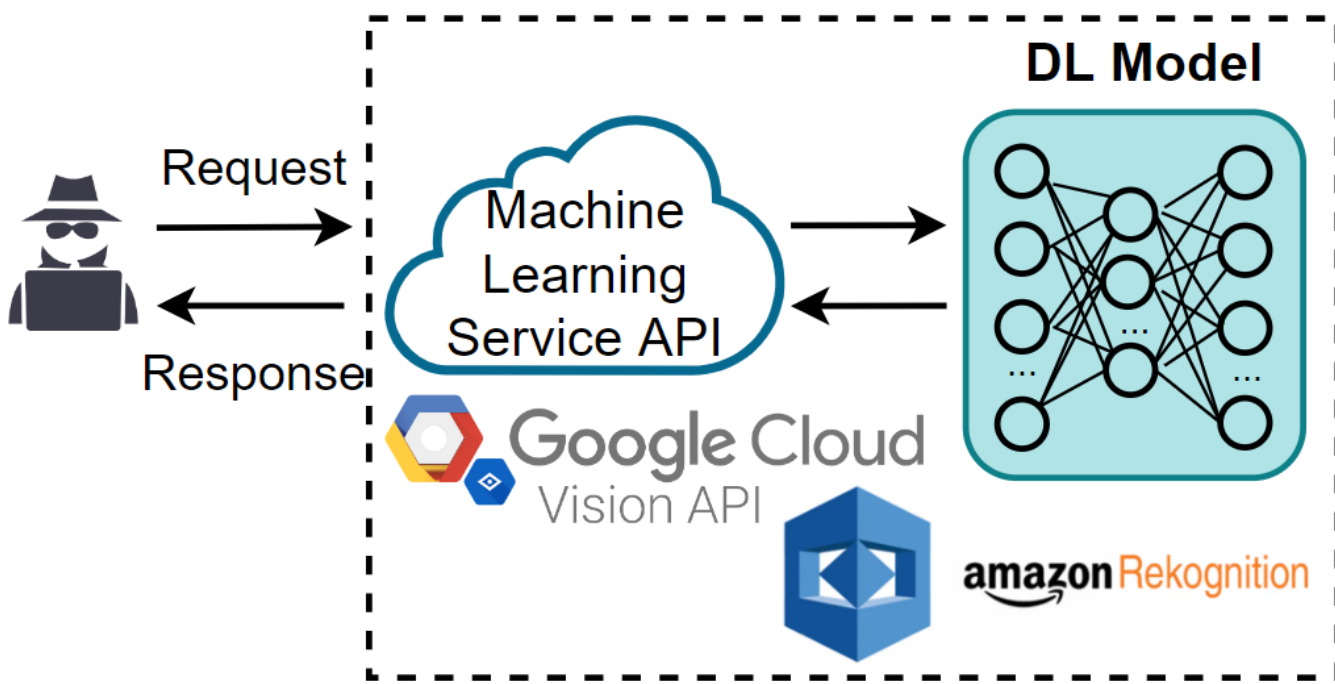


Project page: <https://brusliattack.github.io/>

Introduction

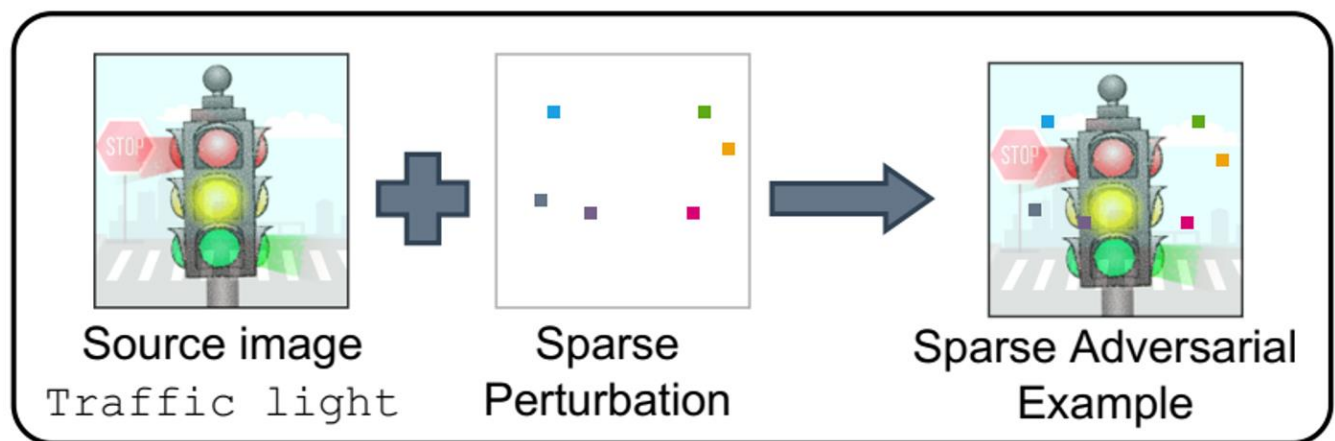
Threat Model

- In real-world systems, the model is hidden from users except for the access to the model's outputs (**a black-box**), e.g. confidence scores.
- This is a pragmatic threat scenario for operational systems, as it underscores the potential for adversaries to exploit even the minimal information available for launching attacks. Thus, we are interested in exploring the vulnerability of models in this context.



Sparse Attacks

Sparse attacks pose an insidious threat as they aim to perturb only a **minimal number of pixels** in an image (l_0 norm) to **mislead a model**. However, sparse attacks have received relatively limited attention in research, we are motivated to explore highly query-efficient sparse attack strategies.



Challenges

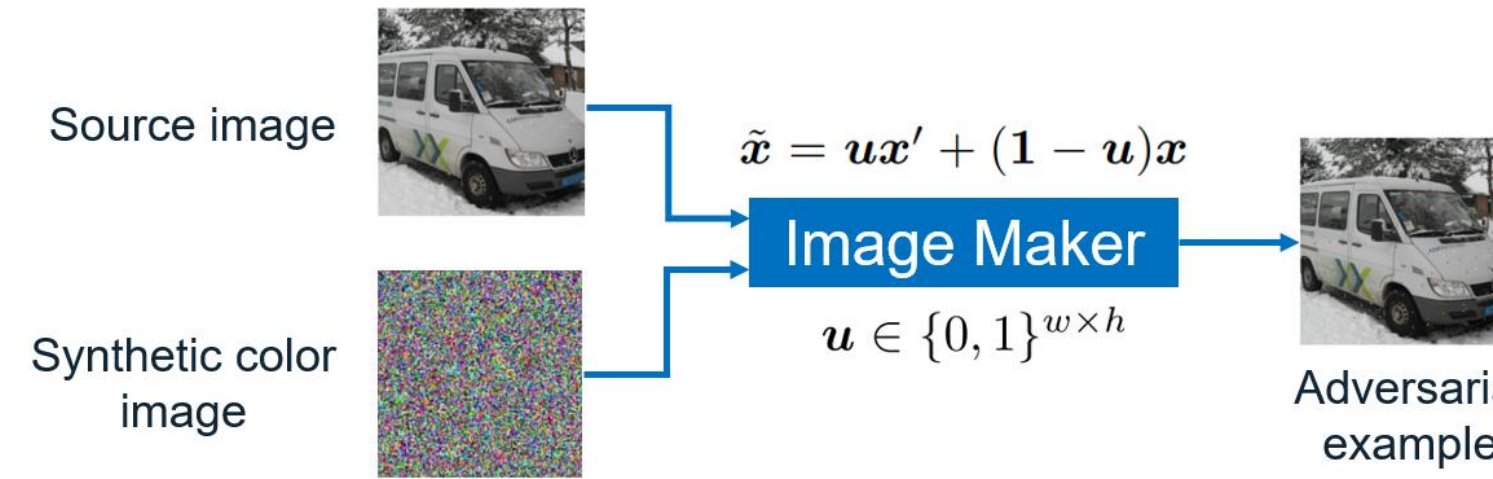
- The search space (Width \times Height \times Channels \times Color values) is incredibly **enormous** to discover a minimum number of pixels. Hence, with limited information from the model's responses, attacks require excessive queries, which is undesirable.
- The **NP-hard problem** [1, 2]
• **Discrete** and **non-differentiable** search space (mixed discrete and continuous) [3].

Problem Formulation and *BruSLeAttack* Methodology

1 Reduce the search space

Problem formulation, $\mathbf{x}^* = \operatorname{argmin}_{\tilde{\mathbf{x}}} L(f(\tilde{\mathbf{x}}), y_{\text{target}})$ s. t. $\|\mathbf{x} - \tilde{\mathbf{x}}\|_0 < B$, leads to an *enormous search space*.

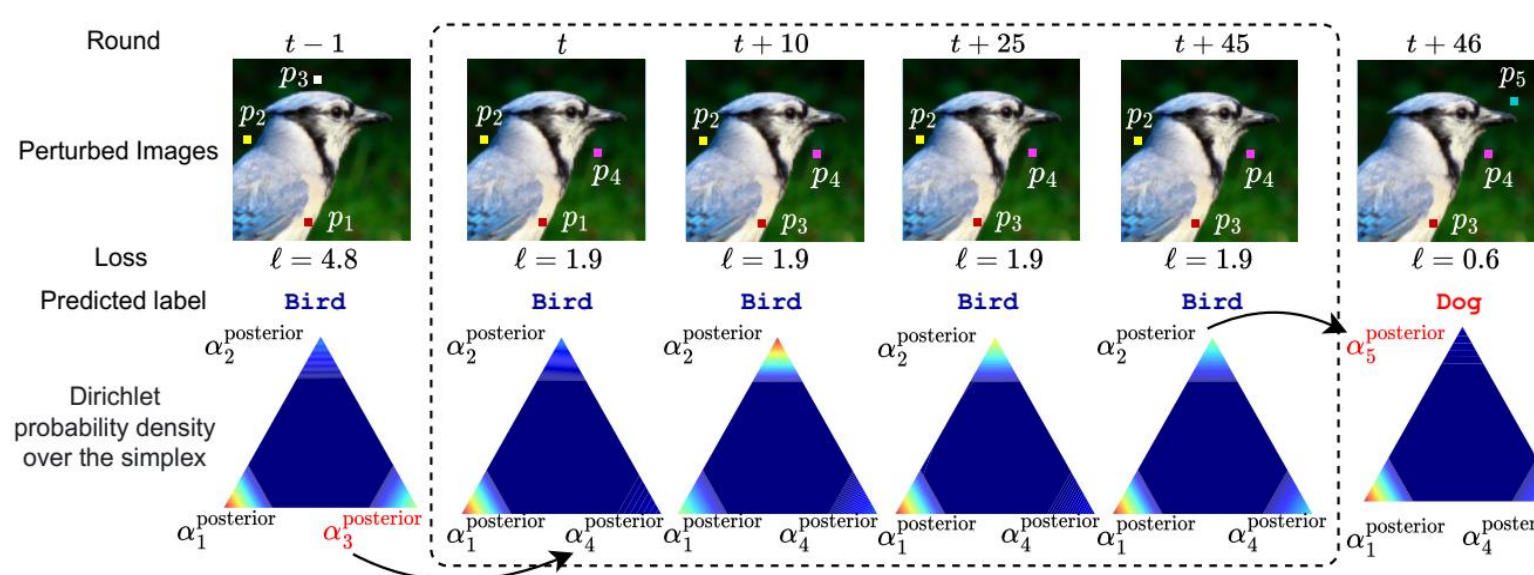
Our **IDEA** is to search and replace pixels in a **source image** with their *corresponding pixels* in a **synthetic color image** so that the malicious loss is reduced, in contrast to attempting to search for pixel values and locations.



Then, we introduce a **new formulation**: $\mathbf{u}^* = \operatorname{argmin}_{\mathbf{u}} \ell(\mathbf{u})$ s. t. $\|\mathbf{u}\|_0 < B$,

where $\ell(\mathbf{u}) = L(f(\tilde{\mathbf{x}}), y_{\text{target}})$, a binary matrix \mathbf{u} is used to *encode an adversarial example*. Elements 0 and 1 denote pixels from the source and synthetic color images, respectively.

2 Remedy the NP-hard problem and handle a non-differentiable search space



Bayesian Framework

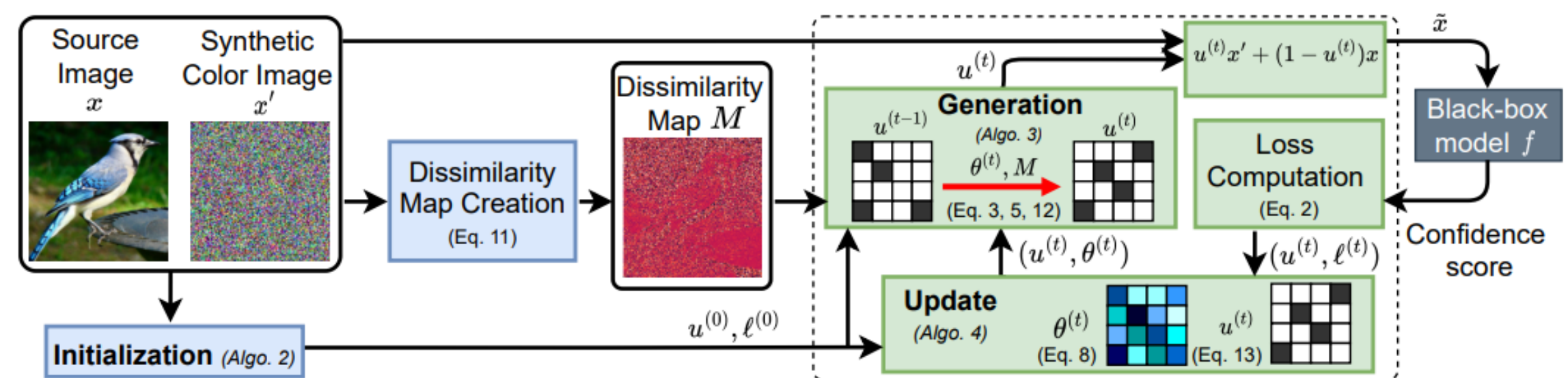
1. **Prior**: $P(\theta, \alpha) = \text{Dir}(\alpha)$

2. **Sampling** $\mathbf{u}^{(t)}$:

$$\begin{aligned} v_1^{(t)}, \dots, v_b^{(t)} &\sim \text{Cat}(v|\theta^{(t)}, u^{(t-1)} = 1) \\ q_1^{(t)}, \dots, q_{B-b}^{(t)} &\sim \text{Cat}(q|\theta^{(t)}, u^{(t-1)} = 0) \\ \mathbf{u}^{(t)} &= [\mathbf{V}_{k=1}^b v_k^{(t)}] \mathbf{V}[\mathbf{V}_{r=1}^{B-b} v_r^{(t)}] \end{aligned}$$

3. **Update** $\theta^{(t)}$:

$$\begin{aligned} \alpha_{i,j}^{\text{posterior}} &= \alpha_{i,j}^{\text{prior}} + s_{i,j}^{(t)} \\ P(\theta|\alpha, u^{(t-1)}, \ell^{(t-1)}) &= \text{Dir}(\alpha^{\text{posterior}}) \\ \theta^{(t)} &= \mathbb{E}_{\theta \sim P(\theta|\alpha, u^{(t-1)}, \ell^{(t-1)})}[\theta] \end{aligned}$$

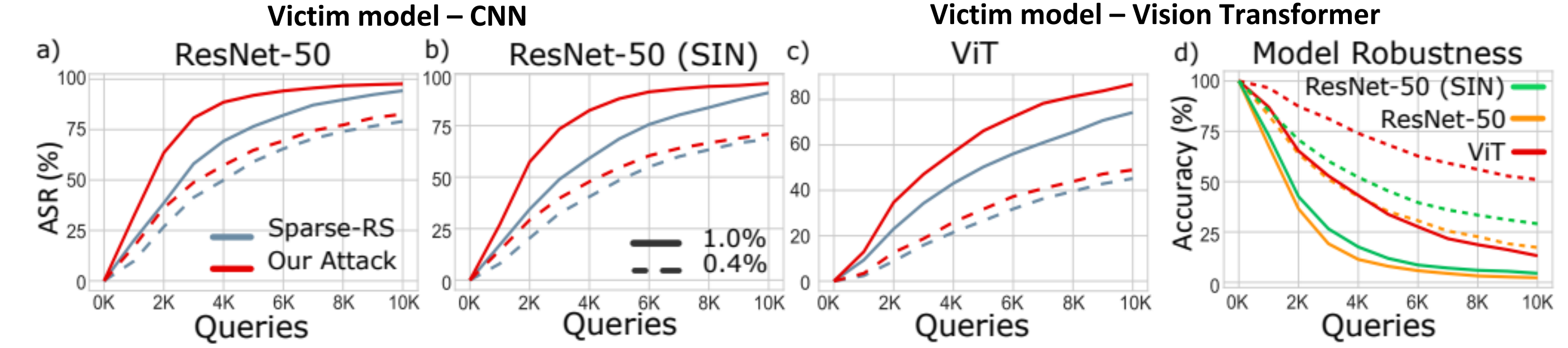


BruSLeAttack Algorithm

- Create a **dissimilarity map** \mathbf{M} between the source and the synthetic color images.
- Initialize some solutions randomly and choose the best for $\mathbf{u}^{(0)}$.
- Sample new $\mathbf{u}^{(t)}$ based on $\theta^{(t)}$, $\mathbf{u}^{(t-1)}$ and \mathbf{M} . Then craft an adversarial image $\tilde{\mathbf{x}}$ from $\mathbf{u}^{(t)}$, \mathbf{x} and \mathbf{x}' .
- Query a black-box model f and calculate loss $\ell^{(t)}$.
- Update both $\theta^{(t)}$ and $\mathbf{u}^{(t)}$ based on the change in the loss and the current solution $\mathbf{u}^{(t)}$.

Results

Attack Transformers & Convolutional Nets



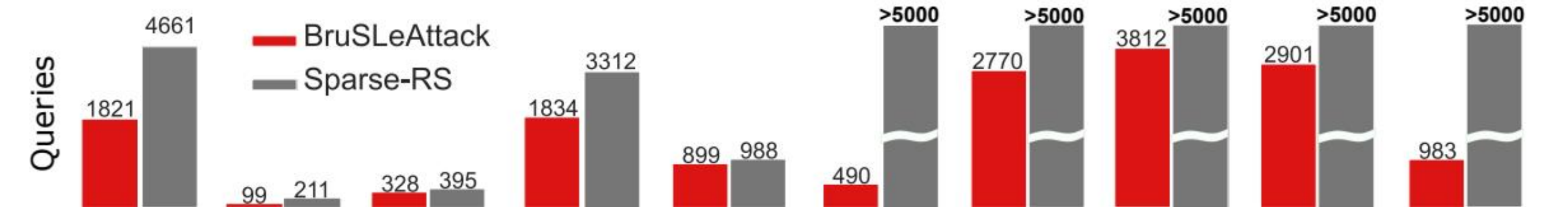
- Query Efficiency**: Within 10K queries, **BruSLeAttack** outperforms state-of-the-art Sparse-RS [4].
- Attack Success Rate (ASR, up to 10K queries)**: **BruSLeAttack** can achieve a much higher ASR than Sparse-RS across different query budgets.

Attack Defended Models

| Sparsity | Undefended Model | | l_∞ -AT | | l_2 -AT | | RND | |
|----------|------------------|--------------|----------------|--------------|-----------|--------------|-----------|--------------|
| | SPARSE-RS | BRUSLEATTACK | SPARSE-RS | BRUSLEATTACK | SPARSE-RS | BRUSLEATTACK | SPARSE-RS | BRUSLEATTACK |
| 0.04% | 33.6% | 24.0% | 43.8% | 42.2% | 89.8% | 88.4% | 90.8% | 85.0% |
| 0.08% | 13.2% | 6.8% | 26.8% | 24.4% | 81.2% | 79.2% | 82.2% | 72.6% |
| 0.12% | 7.6% | 2.6% | 19.0% | 18.4% | 75.8% | 73.8% | 73.6% | 61.0% |
| 0.16% | 5.2% | 1.0% | 16.6% | 14.8% | 71.4% | 69.2% | 64.8% | 51.4% |
| 0.2% | 4.6% | 1.0% | 12.2% | 11.8% | 68.4% | 66.4% | 56.8% | 42.6% |

BruSLeAttack consistently outweighs Sparse-RS against different defense methods and sparsity levels.

Attack Real-world System - Google Cloud Vision



BruSLeAttack is more query efficient than State-of-the-art Sparse-RS.

Conclusions

- BruSLeAttack** is able to remedy the NP-hard problem
- BruSLeAttack** is capable of handling a discrete and non-differentiable search space.
- BruSLeAttack** is more query-efficient than Sparse-RS.

References

- [1] Modas and P. Moosavi-Dezfooli, S. Frossard. Sparsefool: a few pixels make a big difference. CVPR 2019.
- [2] X. Dong, D. Chen, J. Bao, C. Qin, L. Yuan, W. Zhang, N. Yu, and D. Chen. GreedyFool: DistortionAware Sparse Adversarial Attack, NeurIPS, 2020.
- [3] N. Carlini and D. Wagner. Towards evaluating the robustness of neural networks. IEEE SSP, 2017.
- [4] Croce F., Andriushchenko M., Singh N. D., Flammarion N., and Hein M. Sparse-RS: A Versatile Framework for Query-Efficient Sparse Black-Box Adversarial Attacks, AAAI 2022.