

SASA: Sequence-Aware Shadow Attacks via Attention Alignment for Traffic Sign Recognition

Anonymous CVPR submission

Paper ID 21

Abstract

001 We propose **SASA** (*Sequence-Aware Shadow Attack*), a black-
002 box adversarial framework that uses physically realistic, dif-
003 ferentiable shadow patterns to deceive traffic sign recogni-
004 tion systems. Unlike prior image-based attacks, **SASA** targets
005 video sequences—common in real-world driving—by gen-
006 erating smooth, temporally consistent shadows that remain
007 visually plausible and imperceptible to humans. Guided by
008 attention maps from frozen vision transformers, **SASA** aligns
009 shadow placement with semantically salient regions without
010 querying the target model. Evaluated on the *GTSRB* dataset,
011 **SASA** reduces classification accuracy by up to 86% and
012 sequence-level accuracy by over 90% on black-box models,
013 including CNNs and ViTs. The method generalizes across
014 architectures, preserves perceptual quality, and reveals a
015 novel vulnerability in sequential vision systems.

016 1. Introduction

017 Adversarial examples have exposed critical vulnerabilities
018 in deep neural networks, raising significant safety concerns
019 in domains such as autonomous driving. While most adver-
020 sarial research has focused on imperceptible perturbations
021 applied to static images [9, 21], real-world perception sys-
022 tems typically operate over temporally structured inputs like
023 video streams. In such scenarios, frame-level perturbations
024 often fail to generalize due to temporal filtering, camera
025 dynamics, and perceptual inconsistencies [3, 7].

026 In the context of traffic sign recognition (TSR), physical-
027 world attacks have demonstrated real-world feasibility using
028 printed posters [8], adversarial stickers [19], and more re-
029 cently, light projection [25] and shadow casting [26]. These
030 techniques highlight a shift toward stealthy, naturally plausi-
031 ble perturbations that can deceive models without arousing
032 human suspicion. Among these, cast shadows are particu-
033 larly compelling: they are common in real driving scenes,
034 inherently persistent across frames, and do not require modi-
035 fying the object itself. Figure 1 presents visual comparisons

of different light-based physical attacks, including our own
shadow variants.

In this paper, we propose **SASA** (*Semantically Aligned Shadow Attack*), a novel black-box, video-based adversarial attack that leverages physically plausible shadows to induce misclassification in TSR systems. Unlike prior work that either targets single images [26] or assumes white-box access [15], **SASA** generates a temporally consistent shadow mask that aligns with semantically important regions, guided by frozen transformer attention maps (DINO and DeiT), without requiring any access to the target model.

Our approach integrates three components: (1) a compact, physically inspired shadow generator with spatiotemporal control, (2) a differentiable attention-guided alignment strategy, and (3) an optimization routine that ensures realism and temporal consistency. We evaluate **SASA** on both convolutional (CNN, STN) and transformer-based (EffNet, ViT) TSR models. Under a strict black-box setting, **SASA** achieves up to **90% sequence-level accuracy drop** on ViT models and maintains high effectiveness across all architectures—even without direct model access. Among all variants, our DeiT-guided StripShadow delivers the strongest results, lowering average frame-level accuracy by over **30 percentage points** on modern TSR systems. Crucially, the generated shadows remain visually natural and consistent across frames, maintaining realism while significantly degrading model performance.

Contributions. Our main contributions are as follows:

- We introduce **SASA**, the first fully black-box, temporally coherent shadow-based adversarial attack tailored for sequential TSR systems.
- We design a differentiable and physically grounded shadow generation module that allows compact control over shape, position, and opacity.
- We guide perturbation placement using saliency maps from frozen vision transformers (DINO and DeiT), aligning shadows with semantically meaningful regions.
- We demonstrate that **SASA** generalizes across architectures and achieves stealthy, transferable attacks under real-world constraints, reducing TSR accuracy by up to **86%**



Figure 1. Visual comparison of light-based physical attacks including shadow casting [26], reflected laser [23], spotlight projection [25], and natural illumination [10].

076 on ViT and over 60% on EffNet models.

077 2. Related Work

078 **Adversarial Attacks on Vision Models.** Deep neural net-
079 works are vulnerable to adversarial attacks in both digital and
080 physical settings. Traditional white-box [9, 14] and black-
081 box [11, 16] strategies generate subtle perturbations at the
082 pixel level. However, these often fail under real-world trans-
083 formations [2]. Recent evaluations [17] stress the need to
084 move beyond simplistic baselines (e.g., LISA-CNN, GTSRB-
085 CNN) and assess attacks under realistic, sequential, and
086 physically plausible conditions, especially for traffic sign
087 recognition (TSR).

088 **Physical-World and Light-Based Attacks.** A growing
089 body of work focuses on physical adversarial examples, in-
090 cluding posters [18], stickers [8], and camouflage [4]. More
091 recently, light-based attacks like RFLA [23] and AdvSL [25]
092 introduced adversarial perturbations using sunlight and spot-
093 lights, demonstrating high stealth and feasibility in outdoor
094 settings. These methods avoid object modification and offer
095 attacker-controlled timing and adaptability. Natural light
096 attacks [10] further reveal that everyday illumination varia-
097 tions alone can degrade TSR reliability.

098 **Shadow-Based and Temporally Coherent Attacks.**
099 Shadow perturbations have emerged as a subtle, physically
100 grounded threat. Zhong et al. [26] showed triangular shad-
101 ows can achieve over 90% attack success in TSR. Building
102 on this, MohajerAnsari et al. [15] proposed a black-box
103 video attack using a temporally scaled shadow with fixed
104 shape and opacity. Their work incorporates saliency supervi-
105 sion from DINO transformers and introduces the Sequence-
106 Level Attack Success Rate (SL-ASR) to assess robustness
107 across video frames.

108 **Attention and Saliency-Guided Attacks.** Attention
109 maps derived from vision transformers (ViTs) have been
110 used to guide adversarial perturbations in a more inter-
111 pretable manner [13, 24]. The ShadowSeq framework [15]
112 extends this to traffic sign videos, using a dual-loss objective
113 that targets both classification and attention misalignment.
114 These attacks disrupt not only outputs but also internal visual
115 reasoning.

3. Methodology

We propose a physically grounded adversarial attack that
constructs temporally coherent shadow patterns across video
sequences of traffic signs. Our method, SASA (Semantically
Aligned Shadow Attack), comprises three core components:
(1) a differentiable shadow generator, (2) attention-guided
placement informed by pretrained vision transformers, and
(3) an optimization strategy that aligns shadows with seman-
tically salient regions—without requiring access to the target
classifier.

3.1. Problem Formulation

Conventional adversarial attacks in vision often target single-
frame inputs with unconstrained pixel-level perturbations.
However, real-world systems like traffic sign recognition
(TSR) process temporally structured video streams, rely-
ing on frame-to-frame consistency to reject transient noise.
Frame-wise perturbations typically lead to visual flickering
or temporal artifacts, which are easily suppressed by post-
processing or tracking modules. SASA instead introduces
a physically plausible attack that applies a single, tempo-
rally consistent shadow transformation across all frames. By
mimicking natural lighting effects—such as occlusion from
poles or foliage—these shadows maintain spatial coherence
while avoiding perceptual artifacts.

Formally, let $\{x_1, \dots, x_T\}$, where $x_t \in [0, 1]^{3 \times H \times W}$,
denote a video sequence. Let $f: \mathbb{R}^{3 \times H \times W} \rightarrow \mathcal{Y}$ be a TSR
classifier. We seek a shared shadow function $\mathcal{S}(x_t; \theta, \gamma)$
such that each perturbed frame $\tilde{x}_t = \mathcal{S}(x_t)$ is misclassified:
 $f(\tilde{x}_t) \neq f(x_t)$ for all t . Here, θ defines geometric param-
eters of the shadow, and $\gamma \in [0, 1]$ controls its intensity. This
formulation permits efficient, differentiable optimization of
realistic shadows across time.

3.2. Threat Model

We assume a strict black-box threat model: the attacker has
no access to model weights, architecture, gradients, or predic-
tions. Unlike white-box methods that rely on backpropaga-
tion through the target, SASA utilizes frozen attention maps
extracted from transformer models—specifically DINO and
DeiT—fine-tuned on the TSR dataset. These models are used
solely to guide the shadow placement. Importantly, the final
adversarial shadows are evaluated against entirely unseen

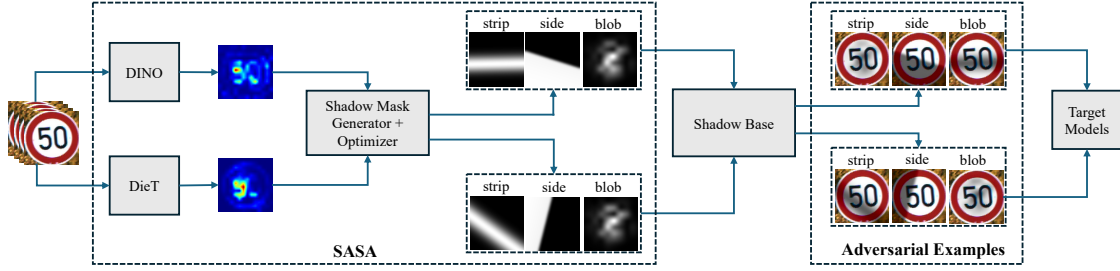


Figure 2. Overview of the SASA pipeline. A clean video sequence is processed by a frozen attention model (DINO/DeiT) to produce a fused saliency map. This map guides the optimization of a shared, differentiable shadow mask. The mask is applied uniformly across all frames, and the resulting perturbed sequence is evaluated against unseen TSR classifiers (e.g., CNNs, ViTs) in a strict black-box setting.

157 target models, ensuring a realistic and rigorous black-box
158 setting.

159 3.3. The SASA Framework

160 SASA produces a single shadow mask that is spatially
161 aligned and temporally consistent across an entire video
162 sequence. The mask is optimized via a compact set of differ-
163 entiable parameters, guided by pretrained attention maps.

164 The pipeline comprises four key stages:

- 165 1. **Attention Aggregation.** A frozen DeiT or DINO trans-
166 former extracts spatial attention maps from each frame x_t .
167 These are fused via temporal max-pooling into a single
168 2D heatmap.
- 169 2. **Shadow Generation.** A shadow generator (*BlobShadow*,
170 *StripShadow*, or *SideShadow*) generates $M \in [0, 1]^{H \times W}$
171 via a differentiable function $\text{Gen}(\theta)$, parameterized by
172 compact geometric variable θ . The same mask is applied
173 to every frame using the transformation $\mathcal{S}(x_t; \theta, \gamma)$.
- 174 3. **Attention-aligned Optimization.** The mask M is opti-
175 mized to match the fused attention map A , using a top- k
176 alignment loss $\mathcal{L}_{\text{align}}$. Crucially, this optimization does
177 not require target model access.
- 178 4. **Black-box Evaluation.** The resulting adversarial shadow
179 is tested on multiple held-out TSR classifiers (e.g., CNNs,
180 ViTs), across a range of intensities γ , with no further
181 tuning.

182 This approach decouples attack generation from
183 model-specific gradients, yielding transferable, physically
184 grounded, and temporally smooth adversarial sequences. An
185 overview of the complete pipeline is shown in Figure 2. We
186 now detail the shadow generation process and the physical
187 priors embedded in our design.

188 3.4. Differentiable Shadow Generation

189 SASA employs a physically motivated shadow generator
190 that simulates real-world occlusions (e.g., from trees, poles,
191 overpasses). Rather than perturbing individual pixels, SASA
192 learns interpretable, low-dimensional parameters that define
193 soft shadow masks $M \in [0, 1]^{H \times W}$.

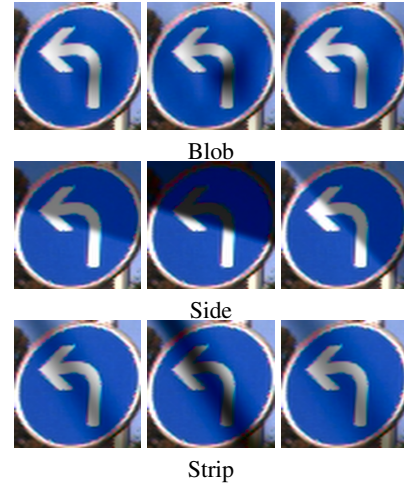


Figure 3. Visual comparison of different shadow styles: Blob, Side, and Strip Shadow. Each sequence shows shadow behavior under varying parameters.

Masks are applied uniformly across the sequence to pre- 194
serve temporal consistency. Each shadow is rendered in 195
CIELAB space, darkening only the luminance channel: 196

$$L' = L \cdot (1 - \gamma M), \quad a' = a, \quad b' = b \quad 197$$

This formulation darkens the luminance channel proportional 198
to the mask while preserving color consistency. The image 199
is then converted back to RGB and clamped to a valid range. 200
To ensure perceptual realism, we apply a Gaussian falloff 201
around the shadow boundaries: 202

$$F(d) = \exp\left(-\left(\frac{d}{r}\right)^2 \cdot \tau\right) \quad 203$$

where d is distance from the mask center, r is radius, and τ 204
controls softness. 205

SASA supports three shadow styles—Blob, Strip, and 206
Side—described in detail next. 207

Algorithm 1 Shadow Mask Generation: BLOBSHADOW and SIDESHADOW**Require:** Image size (H, W) , type $t \in \{\text{blob, side}\}$, parameters θ , blur kernel size

```

1: if  $t = \text{blob then}$ 
Require:  $\theta = (c_x, c_y, s)$ , seed resolution  $k$ 
2: Sample noise tensor  $N \in \mathbb{R}^{1 \times 1 \times k \times k}$ 
3: Upsample  $N \rightarrow \mathbb{R}^{1 \times 1 \times H \times W}$  via bilinear interpolation
4: Apply Gaussian blur; normalize to  $[0, 1]$ 
5: Generate radial falloff at  $(c_x, c_y)$  with radius  $r = s \cdot H$ 
6: Multiply blob by radial falloff and normalize
7: else if  $t = \text{side then}$ 
Require:  $\theta = (c_x, c_y, \alpha, f)$ , blend sharpness  $\kappa$ 
8: Create coordinate grid centered at  $(c_x, c_y)$ ; rotate by  $\alpha$  to get  $y_{\text{rot}}$ 
9: Define half-plane masks:  $M_1 = \mathcal{H}[y_{\text{rot}} > 0]$ ,  $M_2 = \mathcal{H}[y_{\text{rot}} < 0]$ 
10:  $M \leftarrow \sigma((f-0.5)\kappa) \cdot M_2 + (1 - \sigma((f-0.5)\kappa)) \cdot M_1$ 
11: Apply Gaussian blur and distance-based falloff
12: end if
13: return Shadow mask  $M \in [0, 1]^{H \times W}$ 

```

208 **BlobShadow: Radial Shadows.** *BlobShadow* generates
209 soft, amorphous occlusions resembling natural phenomena
210 such as tree cover or cloud shadows. It is parameterized by
211 $\theta = (c_x, c_y, s)$, where $(c_x, c_y) \in [0, 1]^2$ defines the mask
212 center, and $s \in [0.1, 0.8]$ controls spatial extent. The result
213 is a diffuse, circular shadow with smooth falloff, well-suited
214 for simulating irregular, organic occlusions (see Figure 3,
215 left).

216 **StripShadow: Elongated Directional Shadows.** *Strip-*
217 *Shadow* emulates shadows cast by upright objects like poles,
218 signposts, or fences. It is parameterized by $\theta = (c_x, c_y, \alpha, s)$,
219 where $\alpha \in [0, \pi]$ controls orientation and s specifies the
220 width of the shadow band. The resulting mask supports
221 sharp directional control, making it effective for adversarial
222 occlusion of vertically aligned semantic features (see Figure
223 3, right).

224 **SideShadow: Lateral Projection Shadows.** *SideShadow*
225 produces oblique or asymmetric shadows cast from lateral
226 occluders such as roadside barriers or parked vehicles. It is
227 defined by $\theta = (c_x, c_y, \alpha, f)$, where $f \in [0, 1]$ is a continu-
228 ous flip factor interpolating between left- and right-projected
229 masks. Two directional masks are blended using sigmoid-
230 weighted interpolation, followed by Gaussian smoothing.
231 This design allows nuanced, asymmetric shading and cap-
232 tures edge cases often overlooked by uniform shadow models
233 (see Figure 3, middle). Implementation pseudocode for both
234 generators is included in Algorithm 1. To ensure that these
235 shadows align with semantically meaningful features, we
236 rely on saliency signals from pretrained transformers.

3.5. Attention-Guided Shadow Alignment

238 To ensure shadows occlude semantically critical content,
239 SASA leverages attention maps extracted from pretrained
240 vision transformers. We adopt two complementary mod-
241 els—DINO [5] and DeiT [22]—which capture different

forms of visual saliency. DINO, trained via self-supervised
242 contrastive learning, produces object-centric maps, while
243 DeiT, trained in a supervised manner, tends to localize class-
244 discriminative regions. Both models are fine-tuned on the
245 TSR dataset and remain frozen during optimization. 246

Attention Rollout and Temporal Fusion. We use the at-
247 tention rollout method [1] to extract dense saliency maps
248 from transformer layers. At each layer l , we compute
249 residual-aware attention: 250

$$\bar{A}^{(l)} = \alpha A^{(l)} + (1 - \alpha)I, \quad 251$$

where $\alpha = 0.9$ balances the raw attention matrix $A^{(l)} \in$
252 $\mathbb{R}^{N \times N}$ with identity flow I . Cumulative attention is obtained
253 via matrix multiplication across layers: 254

$$\tilde{A} = \bar{A}^{(1)} \bar{A}^{(2)} \dots \bar{A}^{(L)}. \quad 255$$

We extract the first row (class-to-token attention), discard the
256 CLS token, and reshape the remaining values into a 2D map
257 $A_t \in \mathbb{R}^{h \times w}$. Each A_t is upsampled to the input resolution
258 $H \times W$. To consolidate information over time, we perform
259 temporal max-pooling: 260

$$A = \max_{t \in \{1, \dots, T\}} A_t. \quad 261$$

Optionally, we apply a 2D Gaussian prior centered in the
262 image to bias attention toward regions where traffic signs are
263 commonly located: 264

$$G(x, y) = \exp \left(-\lambda \left(\frac{(x - c_x)^2}{W^2} + \frac{(y - c_y)^2}{H^2} \right) \right), \quad 265$$

$$A(x, y) \leftarrow A(x, y) \cdot G(x, y). \quad (1) \quad 266$$

Top- k Alignment Loss. To focus optimization on the most
267 influential pixels, we define a top- k alignment loss. Let \mathcal{T}_k
268 be the indices of the top- k values in the attention map A . The
269 shadow mask M is optimized to overlap with these salient
270 regions: 271

$$\mathcal{L}_{\text{align}} = \frac{1}{k} \sum_{i \in \mathcal{T}_k} |M_i - A_i|^p, \quad 272$$

where $p = 1$ corresponds to L1 loss and $p = 2$ to mean
273 squared error. This design encourages sparse but targeted
274 occlusion, yielding compact shadows that effectively disrupt
275 classification. 276

3.6. Full Optimization Objective

The final loss function combines semantic alignment with
277 optional regularization: 278

$$\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{align}} + \lambda_{\text{reg}} \cdot \mathcal{L}_{\text{reg}}. \quad 280$$

281 Here, \mathcal{L}_{reg} may include total variation, sparsity, or spatial
282 smoothness penalties to enhance physical plausibility. Importantly, SASA requires no gradients or queries from the target
283 classifier; all optimization is driven by attention-derived signals, preserving the black-box threat model while supporting
284 efficient and transferable adversarial attacks.

285 Together, these components form a compact, transferable, and physically realizable adversarial attack framework
286 suitable for black-box, real-world TSR systems.

290 4. Experiments

291 We evaluate SASA in a strict black-box setting on the traffic
292 sign recognition (TSR) task. Our experimental goals are
293 fourfold: (1) quantify attack effectiveness across a range
294 of model architectures, (2) assess temporal consistency and
295 perceptual realism of shadow-based perturbations, (3) evaluate
296 transferability to unseen models without adaptation, and
297 (4) study the impact of key design choices through targeted
298 ablations.

299 Our analysis spans three shadow types, multiple shadow
300 intensities, and two pretrained attention backbones (DeiT
301 and DINO), offering a comprehensive view into SASA’s
302 effectiveness, generalization, and physical plausibility.

303 4.1. Experimental Setup

304 **Dataset and Preprocessing** We use the **German Traffic**
305 **Sign Recognition Benchmark (GTSRB)** [20], a real-world
306 dataset consisting of over 50,000 traffic sign images cate-
307 gorized into 43 classes. Each sample includes bounding box
308 annotations and metadata that support constructing tempo-
309 rally adjacent sequences of the same sign.

310 To simulate sequential video input, we group images
311 with the same class and sample ID to form short sequences.
312 Each sequence is standardized to a fixed length of $T = 30$
313 frames using a sliding window approach. Images are cropped
314 based on bounding boxes and resized according to model-
315 specific input dimensions: 128×128 for CNN-based models
316 and 224×224 for transformer-based models. Color and
317 luminance statistics (e.g., average LAB and HSV values)
318 are computed for each frame to enable filtering and ablation
319 control.

320 We curate a high-visibility evaluation subset by discard-
321 ing sequences with low average luminance ($L < 120$ in
322 CIELAB space), ensuring that generated shadows remain
323 perceptually distinguishable. All train/test splits are class-
324 stratified and fixed across experiments to support repro-
325 ducibility.

326 **Target Models** We evaluate SASA against a diverse set
327 of TSR classifiers spanning convolutional and transformer-
328 based architectures.

329 **CNN-based models.** We include two widely used base-
330 lines: **GTSRB-CNN** [6] and **GTSRB-STN** [12], reimple-

mented and trained from scratch on the GTSRB training set.
These models process inputs at 128×128 resolution and
serve as classical vision backbones for TSR.

Transformer-based models. We fine-tune ImageNet-
pretrained **EfficientNet-B0** and **ViT-Small** models on GT-
SRB, adapting them to traffic sign classification via end-to-
end supervised training. Both models use inputs resized to
 224×224 . This setup allows us to evaluate SASA’s per-
formance on modern attention-based architectures under a
black-box constraint.

Attention models. **DeiT** [22] and **DINO** [5] are used
exclusively as frozen saliency extractors. These models are
fine-tuned on GTSRB to improve attention map fidelity, but
are not used as targets during evaluation. Their outputs guide
shadow placement during optimization without accessing the
target model, thus strictly adhering to the black-box threat
model.

Evaluation Metrics We evaluate SASA’s effectiveness
using two complementary metrics that reflect both frame-
level performance and sequence-level robustness.

Frame-wise Accuracy. This metric computes the aver-
age classification accuracy across all individual frames in a
sequence. It serves as a fine-grained measure of how con-
sistently a model classifies perturbed inputs on a per-frame
basis.

Sequence-level Accuracy at 50% (seq@50). This
stricter metric assesses whether a majority of frames in a
sequence are correctly classified. A sequence is considered
correct if at least 50% of its frames yield the correct label.
seq@50 captures temporal robustness and provides a more
task-relevant evaluation for video-based TSR, where tran-
sient misclassifications may be tolerable, but sustained errors
are not.

Together, these metrics allow us to assess not only the raw
effectiveness of shadow-based perturbations but also their
practical impact on real-world TSR pipelines.

4.2. Evaluation Results

We evaluate SASA across four black-box TSR models:
GTSRB-CNN, GTSRB-STN, EfficientNet-B0, and ViT. Table 1
reports performance drops in frame-wise accuracy and
sequence-level accuracy (seq@50) at three shadow intensi-
ties ($\gamma = \{0.2, 0.5, 0.8\}$), relative to the clean baseline.

Note on Comparability. We do not include direct com-
parisons with prior shadow-based attacks such as [26], as
these methods are designed to attack single-frame images
using per-image optimization. In contrast, SASA targets
full video sequences by optimizing a single, temporally
consistent shadow mask. This fundamental difference in
formulation makes direct comparisons methodologically in-
compatible.

Table 1. Classification accuracy (%) and sequence accuracy at 50% threshold ($\text{seq}@50$, %) for clean sequences and relative drop (\downarrow) under each shadow variant. For each intensity level ($\gamma = 0.2, 0.5, 0.8$), the highest drop per model (i.e., strongest attack) is **bolded**.

Intensity Variant		CNN		STN		EffB0		ViT	
		Acc \downarrow	seq@50 \downarrow	Acc \downarrow	seq@50 \downarrow	Acc \downarrow	seq@50 \downarrow	Acc \downarrow	seq@50 \downarrow
–	Clean Sequence	96.6	100.0	95.8	100.0	99.9	100.0	95.1	100.0
0.2	BlobShadow (DeiT)	0.6	0.0	2.4	2.5	1.1	0.0	4.6	5.0
	BlobShadow (DINO)	0.6	0.0	2.7	2.5	1.0	0.0	3.8	2.5
	SideShadow (DeiT)	2.2	0.0	2.4	2.5	0.5	0.0	5.6	2.5
	SideShadow (DINO)	0.9	0.0	1.9	2.5	0.3	0.0	2.7	2.5
	StripShadow (DeiT)	1.3	0.0	2.6	2.5	2.5	0.0	4.0	2.5
	StripShadow (DINO)	0.5	0.0	2.6	2.5	0.7	0.0	1.4	0.0
0.5	BlobShadow (DeiT)	9.0	7.5	12.4	10.0	8.4	5.0	20.3	12.5
	BlobShadow (DINO)	8.6	5.0	12.9	7.5	6.5	2.5	21.0	17.5
	SideShadow (DeiT)	11.8	7.5	14.3	15.0	9.1	7.5	32.5	35.0
	SideShadow (DINO)	8.8	2.5	10.5	7.5	8.2	2.5	24.2	20.0
	StripShadow (DeiT)	12.4	5.0	20.2	20.0	30.8	27.5	31.4	37.5
	StripShadow (DINO)	7.0	2.5	15.8	15.0	17.8	15.0	27.6	27.5
0.8	BlobShadow (DeiT)	22.8	22.5	30.1	25.0	21.7	12.5	75.3	80.0
	BlobShadow (DINO)	21.6	22.5	31.7	25.0	22.1	15.0	75.8	90.0
	SideShadow (DeiT)	35.1	35.0	44.2	45.0	33.3	32.5	80.6	87.5
	SideShadow (DINO)	35.3	40.0	43.2	45.0	35.5	27.5	75.8	80.0
	StripShadow (DeiT)	40.6	47.5	46.9	47.5	62.0	62.5	86.3	92.5
	StripShadow (DINO)	27.3	25.0	44.7	45.0	40.1	35.0	78.8	85.0

Table 2. Classification accuracy (%) and sequence accuracy at 50% ($\text{seq}@50$, %) averaged over all shadow intensities ($\gamma \in [0.1, 1.0]$). Gray-box evaluations (DeiT and DINO classifiers) are shown in the last columns. The lowest accuracy (i.e., strongest attack) for each model is highlighted in green.

Variant	CNN		STN		EffB0		ViT		DeiT		DINO	
	Acc	seq@50	Acc	seq@50	Acc	seq@50	Acc	seq@50	Acc	seq@50	Acc	seq@50
Clean (baseline)	96.6	100.0	95.8	100.0	99.9	100.0	95.1	100.0	98.8	100.0	92.8	100.0
BlobShadow (DeiT)	83.5	86.8	77.2	82.8	87.1	91.3	55.3	58.8	59.4	61.5	68.6	74.0
BlobShadow (DINO)	83.9	87.8	76.3	83.3	87.5	91.3	55.5	57.5	59.9	62.3	69.6	74.5
SideShadow (DeiT)	74.4	80.5	68.1	71.5	77.7	80.3	49.8	53.5	46.4	48.0	62.6	67.5
SideShadow (DINO)	76.2	79.8	70.7	74.3	79.2	81.5	54.9	58.5	51.3	52.8	65.7	70.3
StripShadow (DeiT)	73.5	78.5	66.3	69.8	64.7	65.3	49.4	52.0	43.3	44.8	54.8	61.0
StripShadow (DINO)	81.1	86.5	69.2	74.3	76.6	79.0	53.4	55.8	49.2	50.8	61.4	64.8

381 At low intensity ($\gamma = 0.2$), shadows are subtle, and no
382 single configuration dominates. SideShadow (DeiT) yields
383 the strongest effect on CNN and ViT, while other variants per-
384 form comparably. This suggests that at low visibility, effec-
385 tiveness depends more on attention alignment than shadow
386 shape.

387 At medium intensity ($\gamma = 0.5$), StripShadow (DeiT) con-
388 sistentlly outperforms other variants, achieving the highest
389 drop on all models except ViT, where SideShadow (DeiT) is
390 slightly better. Its long pattern appears well-suited to overlap
391 with high-attention regions, mostly on structured signs. At
392 high intensity ($\gamma = 0.8$), StripShadow (DeiT) stays domi-
393 nant, with the largest accuracy and seq@50 drops across
394 all models—exceeding 60% on EfficientNet and 90% on
395 ViT. SideShadow and BlobShadow also remain effective
396 but fall short of StripShadow’s disruption capacity. Overall,
397 DeiT-guided variants outperform their DINO-based variants,
398 confirming that class-discriminative attention better informs

shadow placement. The results validate SASA’s core design: structured shadows guided by transformer attention yield strong, transferable black-box attacks.

4.3. Average Performance Across Intensities

To evaluate the overall impact of shadow-based attacks, we compute the average classification accuracy and sequence-level accuracy ($\text{seq}@50$) across the full intensity range ($\gamma \in [0.1, 1.0]$). Table 2 reports these results for each shadow variant. The table includes both black-box settings, where the attention model differs from the classifier, and gray-box cases, where the same transformer (DeiT or DINO) is used for both guiding attention and optimizing shadow. This offers a unified view of attack efficacy across shadow types, models, and threat settings.

Among all variants, StripShadow with DeiT attention consistently yields the lowest average accuracy, particularly on transformer-based targets like ViT and DeiT, suggesting that

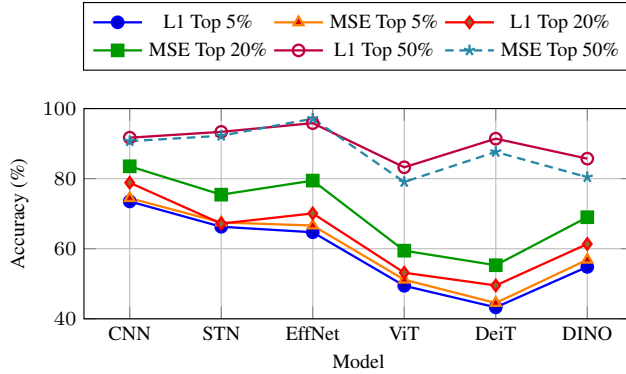


Figure 4. Final classification accuracy (lower is better) across six target models under different loss types and top- k attention coverage. Lower accuracy indicates stronger attacks.

416 this combination best aligns perturbations with the model’s
 417 discriminative regions. Despite its simplicity, StripShadow
 418 is especially effective when high-intensity shadows are
 419 applied. ViT emerges as the most vulnerable target, likely due
 420 to its strong reliance on attention-based localization, which
 421 aligns with SASA’s strategy. In contrast, GTSRB-CNN
 422 shows higher resilience across shadow types, highlighting
 423 the relative robustness of convolutional architectures under
 424 such physically inspired perturbations.

4.4. Loss Function Ablation

426 We ablate the effect of loss type and attention sparsity
 427 on SASA’s attack effectiveness. Specifically, we com-
 428 pare six configurations: L1 and MSE alignment losses ap-
 429 plied over the top- k % most salient attention values, with
 430 $k \in \{5, 20, 50\}$. All experiments use the StripShadow gen-
 431 erator guided by DeiT attention, and results are averaged over
 432 shadow intensities $\gamma \in [0.1, 1.0]$.

433 Figure 4 reports final classification accuracy for each
 434 configuration across six target models. Since lower accuracy
 435 indicates more effective attacks, several trends emerge:

436 **Larger top- k regions result in stronger attacks.** The top-
 437 50% variants consistently yield the lowest accuracy across
 438 all models, confirming that broader attention coverage leads
 439 to more effective occlusion of discriminative regions.

440 **MSE slightly outperforms L1** in low- k settings (e.g.,
 441 top-5%), likely due to smoother gradient propagation during
 442 optimization. However, this difference diminishes as the
 443 attention mask grows denser.

444 **Transformer models are more vulnerable.** ViT and
 445 DeiT show the lowest final accuracies, suggesting greater
 446 susceptibility to attention-aligned perturbations. In contrast,
 447 CNN-based models retain higher accuracy, consistent with
 448 their reduced reliance on global saliency cues.

449 These results validate the design of SASA’s top- k align-
 450 ment loss and highlight the importance of both attention

coverage and loss smoothness. Even simple losses like L1
 451 and MSE can yield strong black-box attacks when guided by
 452 semantically meaningful saliency maps. 453

4.5. Qualitative Attention Alignment

454 To visualize how SASA targets model saliency, we present
 455 qualitative examples for each shadow type. Figure 5
 456 shows three representative sequences—one per shadow mod-
 457 ule—highlighting the effect of attention-guided shadow
 458 placement and its impact on model focus. 459

460 Each row includes: (1) the original input frame, (2) the
 461 DeiT-derived attention map on the clean frame, (3) the opti-
 462 mized shadow mask, and (4) the attention map after applying
 463 the shadow. These visualizations demonstrate how shadows
 464 occlude high-saliency regions, leading to degraded or dif-
 465 fused attention in the perturbed frames.

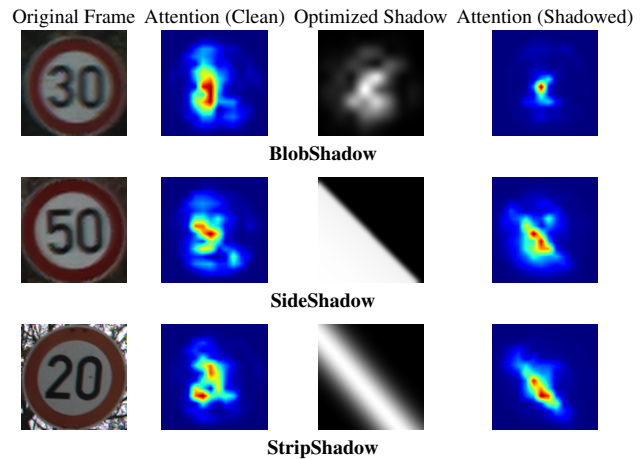


Figure 5. Qualitative attention analysis across shadow types. Each row shows: (1) the original frame, (2) the DeiT attention map before attack, (3) the optimized shadow mask, and (4) attention map after applying the shadow.

5. Conclusion

466 We introduced SASA, a black-box adversarial framework
 467 that generates physically plausible, temporally coherent shad-
 468 ows to mislead traffic sign recognition systems. By combin-
 469 ing differentiable shadow generators with transformer-based
 470 attention maps, SASA aligns perturbations with semanti-
 471 cally critical regions—without accessing model gradients
 472 or predictions. Our experiments on the GTSRB benchmark
 473 demonstrate that SASA significantly degrades classification
 474 performance across CNN and transformer models, while
 475 maintaining high perceptual realism and transferability. Ab-
 476 lation studies confirm the importance of attention-guided
 477 alignment and shadow diversity. Overall, SASA offers a
 478 scalable, interpretable, and physically grounded approach
 479 for evaluating real-world vulnerabilities in sequential vision
 480 systems. 481

482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537**References**

- [1] Samira Abnar and Willem Zuidema. Quantifying attention flow in transformers. *arXiv preprint arXiv:2005.00928*, 2020. 4
- [2] Anish Athalye, Nicholas Carlini, and David Wagner. Obfuscated gradients give a false sense of security: Circumventing defenses to adversarial examples. In *International conference on machine learning*, pages 274–283. PMLR, 2018. 2
- [3] Anish Athalye, Logan Engstrom, Andrew Ilyas, and Kevin Kwok. Synthesizing robust adversarial examples. In *International conference on machine learning*, pages 284–293. PMLR, 2018. 1
- [4] Tom B Brown, Dandelion Mané, Aurko Roy, Martín Abadi, and Justin Gilmer. Adversarial patch. *arXiv preprint arXiv:1712.09665*, 2017. 2
- [5] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 9650–9660, 2021. 4, 5
- [6] Dan Ciregan, Ueli Meier, and Jürgen Schmidhuber. Multicolumn deep neural networks for image classification. In *2012 IEEE conference on computer vision and pattern recognition*, pages 3642–3649. IEEE, 2012. 5
- [7] Yinpeng Dong, Tianyu Pang, Hang Su, and Jun Zhu. Evading defenses to transferable adversarial examples by translation-invariant attacks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4312–4321, 2019. 1
- [8] Kevin Eykholt, Ivan Evtimov, Earlene Fernandes, and et al. Robust physical-world attacks on deep learning visual classification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1625–1634, 2018. 1, 2
- [9] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. *arXiv preprint arXiv:1412.6572*, 2014. 1, 2
- [10] Teng-Fang Hsiao, Bo-Lun Huang, Zi-Xiang Ni, Yan-Ting Lin, Hong-Han Shuai, Yung-Hui Li, and Wen-Huang Cheng. Natural light can also be dangerous: Traffic sign misinterpretation under adversarial natural light attacks. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 3915–3924, 2024. 2
- [11] Andrew Ilyas, Logan Engstrom, Anish Athalye, and Jessy Lin. Black-box adversarial attacks with limited queries and information. In *International conference on machine learning*, pages 2137–2146. PMLR, 2018. 2
- [12] Max Jaderberg, Karen Simonyan, Andrew Zisserman, et al. Spatial transformer networks. *Advances in neural information processing systems*, 28, 2015. 5
- [13] Hsueh-Ti Derek Liu, Michael Tao, Chun-Liang Li, et al. Beyond pixel norm-balls: Parametric adversaries using an analytically differentiable renderer. In *International Conference on Machine Learning (ICML)*, pages 224–233, 2018. 2
- [14] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. In *International Conference on Learning Representations (ICLR)*, 2018. 2
- [15] Pedram MohajerAnsari, Amir Salarpour, David Fernandez, Cigdem Kokenoz, Bing Li, and Mert D Pesé. Attention-aware temporal adversarial shadows on traffic sign sequences. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2025. 1, 2
- [16] Nicolas Papernot, Patrick McDaniel, Ian Goodfellow, Somesh Jha, Z Berkay Celik, and Ananthram Swami. Practical black-box attacks against deep learning systems using adversarial examples. *arXiv preprint arXiv:1602.02697*, 1(2):3, 2016. 2
- [17] Svetlana Pavlitska, Nico Lambing, and J Marius Zöllner. Adversarial attacks on traffic sign recognition: A survey. In *2023 3rd International conference on electrical, computer, communications and mechatronics engineering (ICECCME)*, pages 1–6. IEEE, 2023. 2
- [18] Mahmood Sharif, Sruti Bhagavatula, Lujo Bauer, and Michael K Reiter. Accessorize to a crime: Real and stealthy attacks on state-of-the-art face recognition. In *Proceedings of the 2016 acm sigsac conference on computer and communications security*, pages 1528–1540, 2016. 2
- [19] Chawin Sitawarin, Arjun Nitin Bhagoji, Arsalan Mosenia, Mung Chiang, and Prateek Mittal. Darts: Deceiving autonomous cars with toxic signs. *arXiv preprint arXiv:1802.06430*, 2018. 1
- [20] Johannes Stalldkamp, Marc Schlipf, Jan Salmen, and Christian Igel. Man vs. computer: Benchmarking machine learning algorithms for traffic sign recognition. *Neural networks*, 32: 323–332, 2012. 5
- [21] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus. Intriguing properties of neural networks. *arXiv preprint arXiv:1312.6199*, 2013. 1
- [22] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. Training data-efficient image transformers & distillation through attention. In *International conference on machine learning*, pages 10347–10357. PMLR, 2021. 4, 5
- [23] Donghua Wang, Wen Yao, Tingsong Jiang, Chao Li, and Xiaoqian Chen. Rfla: A stealthy reflected light adversarial attack in the physical world. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 4455–4465, 2023. 2
- [24] Jie Wang, Zhaoxia Yin, Jing Jiang, and Yang Du. Attention-guided black-box adversarial attacks with large-scale multi-objective evolutionary optimization. *International Journal of Intelligent Systems*, 37(10):7526–7547, 2022. 2
- [25] LI Yufeng, YANG Fengyu, LIU Qi, LI Jiangtao, and CAO Chenhong. Light can be dangerous: Stealthy and effective physical-world adversarial attack by spot light. *Computers & Security*, 132:103345, 2023. 1, 2
- [26] Yiqi Zhong, Xianming Liu, Deming Zhai, Junjun Jiang, and Xiangyang Ji. Shadows can be dangerous: Stealthy and effective physical-world adversarial attack by natural phenomenon. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 15345–15354, 2022. 1, 2, 5