

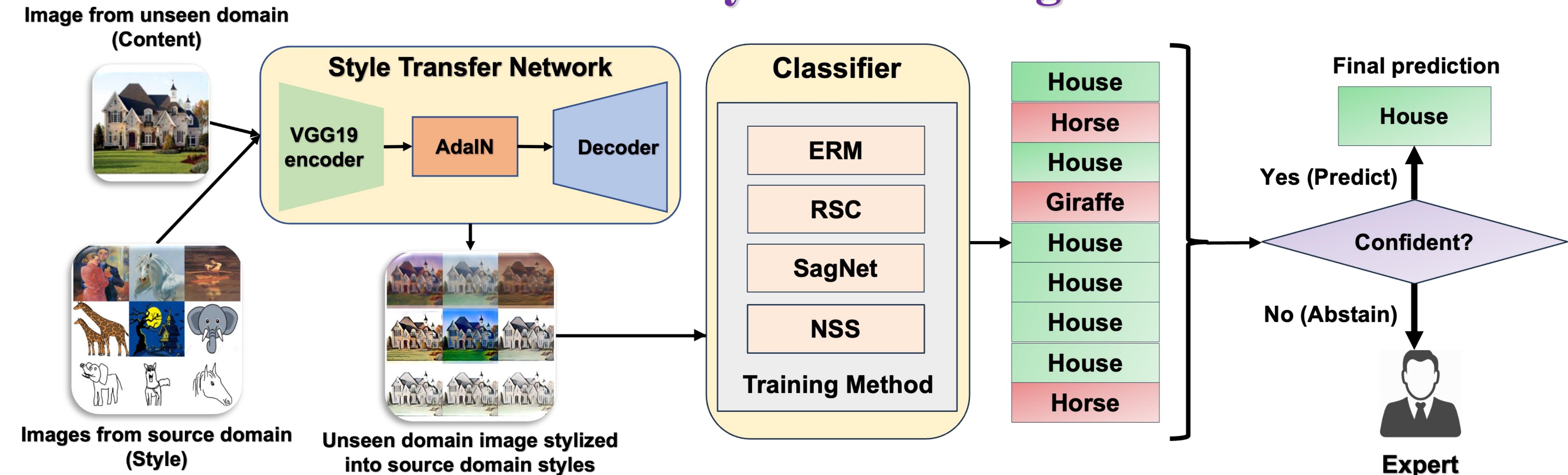
Motivation

- Machine learning (ML) models suffer a significant accuracy degradation on data from unseen domains.
- Recent works demonstrated the use of domain-specific information (e.g., image styles) for predictions as a potential cause for this accuracy degradation.
- This makes ML models unreliable for deployment in high-risk scenarios where misclassifications can have catastrophic consequences.

Contributions

- We propose an efficient inference procedure to obtain risk-averse predictions on the fly from a model by requiring only black-box access to it.
- We propose a training procedure relying on losses that enforce prediction consistency on random stylization of the training data.
- We demonstrate the effectiveness of our inference and training methods on benchmark datasets and their variations.

Neural Style Smoothing



Let \mathcal{X}, \mathcal{Y} be the data domain/labels, f be the classifier, and x_c/x_s be the content/style images, g/h be the encoder/decoder, then $h(t)$ is the AdaIn image with $t := AI(x_c, x_s) = \sigma(g(x_s)) \left(\frac{g(x_c) - \mu(g(x_c))}{\sigma(g(x_c))} \right) + \mu(g(x_s))$.

Given a classifier $f: \mathcal{R}^d \rightarrow \mathcal{Y}$, prediction of a style smoothed classifier $\psi: \mathcal{R}^d \rightarrow \mathcal{Y}$ on a sample x is given by

$$\psi(x) := \arg \max_{y \in \mathcal{Y}} P(f(h(t)) = y),$$

where $t = AI(x, x_s)$, $x_s \sim P_S$ and P_S is the source distribution.

```
def get_counts(x, f, g, h):
    set cls_counts to zeros.
    for i = 1, ..., n do
        Sample x_s^i from P_S
        t = AI(x, x_s^i; g)
        x_stylized = h(t)
        pred = f(x_stylized)
        cls_counts[pred] += 1
    return cls_counts
```

```
def TT_NSS(x, f, g, h, n, alpha):
    cls_counts = get_counts(x, f, g, h)
    c_max = arg max_y cls_counts
    n_max = cls_counts[c_max]
    if n_max / n < alpha then
        return ABSTAIN
    else:
        return c_max
```

Test-Time Neural Style Smoothing (TT-NSS)

- NSS-based losses to improve risk-averse predictions.
 - Style Augmentation:** $\mathbb{E}_{x_s \sim P_S} [\ell(f(h(t)), y)]$.
 - Style Consistency:** $\mathbb{E}_{x_s \sim P_S} [KL(\bar{F}(x)) || F(h(t))] + H(\bar{F}(x), y)$, where $F: \mathcal{R}^d \rightarrow \Delta^{K-1}$ is the soft classifier of f and Δ^{K-1} is the probability simplex in \mathcal{R}^K , $H(\cdot)$ is the entropy and $\bar{F}(x) = \mathbb{E}_{x_s \sim P_S} [F(h(t))]$.

Key Results

TT-NSS produces better risk-averse predictions than the confidence-based abstaining mechanism.

Effect of the number of re-stylizations, n , in TT-NSS.

NSS trained models achieve better risk averse predictions than models trained with ERM.

