

---

# Shadows of Deception: Unveiling AI-generated Images Through Inconsistencies in Scene Lighting

**Muhammad Abdullah**

**Supervisors: PD. Dr. Christian Riess, Prof. Dr. Bernhard Egger**

# Light in real scenes follows the laws of physics. Can data-driven AI-generated images learn these laws?



AI-generated image

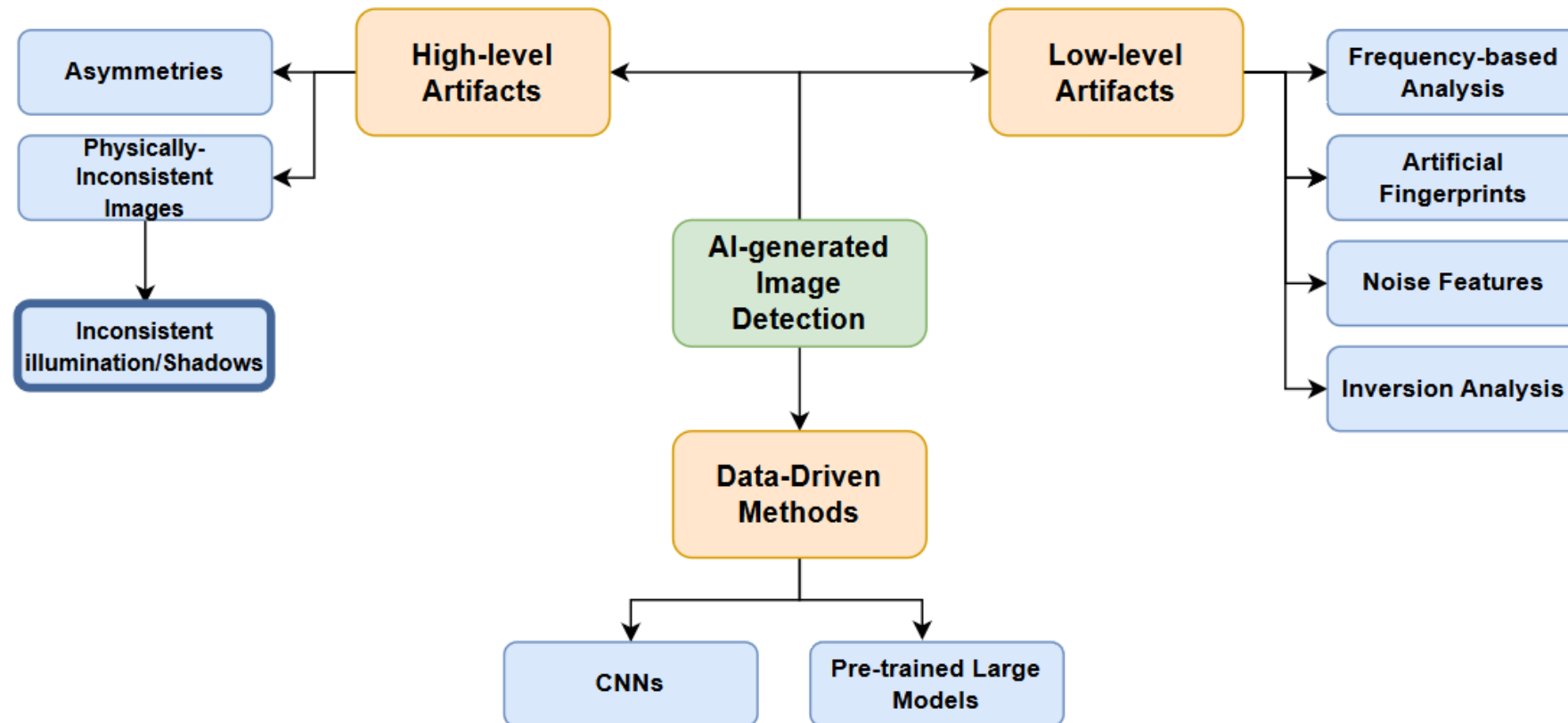


AI-generated image

- 
- 01 Related work
  - 02 Modeling scene illumination
  - 03 Inverse rendering
  - 04 Methodology
  - 05 Baseline methods
  - 06 Dataset creation
  - 07 Experiments and results
  - 08 Qualitative analysis
  - 09 Conclusion and future directions

# Related Work: Detecting AI-generated Images

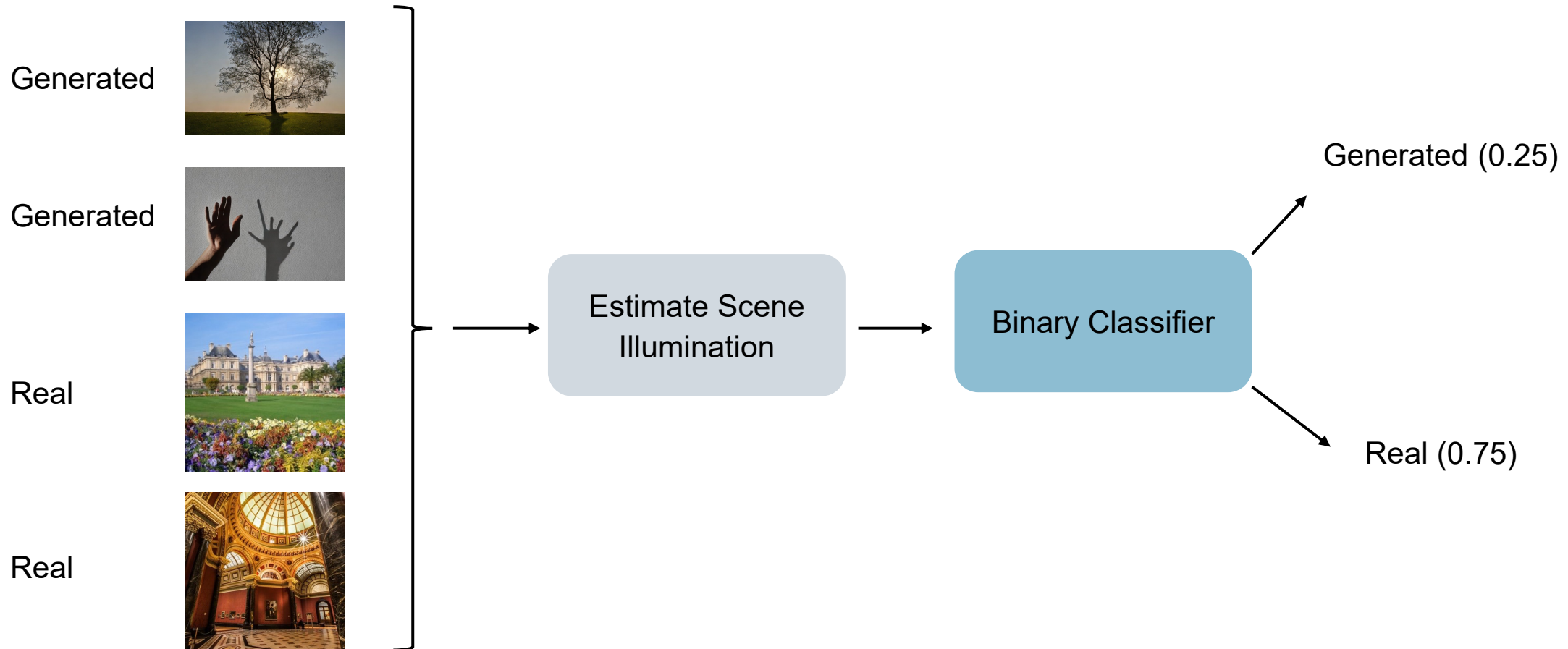
Classification of different approaches used for detecting AI-generated images



Classification of AI-generated images detection methods<sup>1</sup>

# Scene Illumination for Detecting AI-generated Images

Estimated scene illumination for real and generated images is passed to a binary classifier to classify real and generated images



---

# Modeling Scene Illumination



### The Rendering Equation<sup>1</sup>

$$L(\mathbf{x}, \vec{\omega}_o) = L_e(\mathbf{x}, \vec{\omega}_o) + \int_S f_r(\mathbf{x}, \vec{\omega}_i \rightarrow \vec{\omega}_o) L(\mathbf{x}', \vec{\omega}_i) G(\mathbf{x}, \mathbf{x}') V(\mathbf{x}, \mathbf{x}') d\omega_i$$

where

$L(\mathbf{x}, \vec{\omega}_o)$  = the intensity reflected from position  $\mathbf{x}$  in direction  $\omega_o$

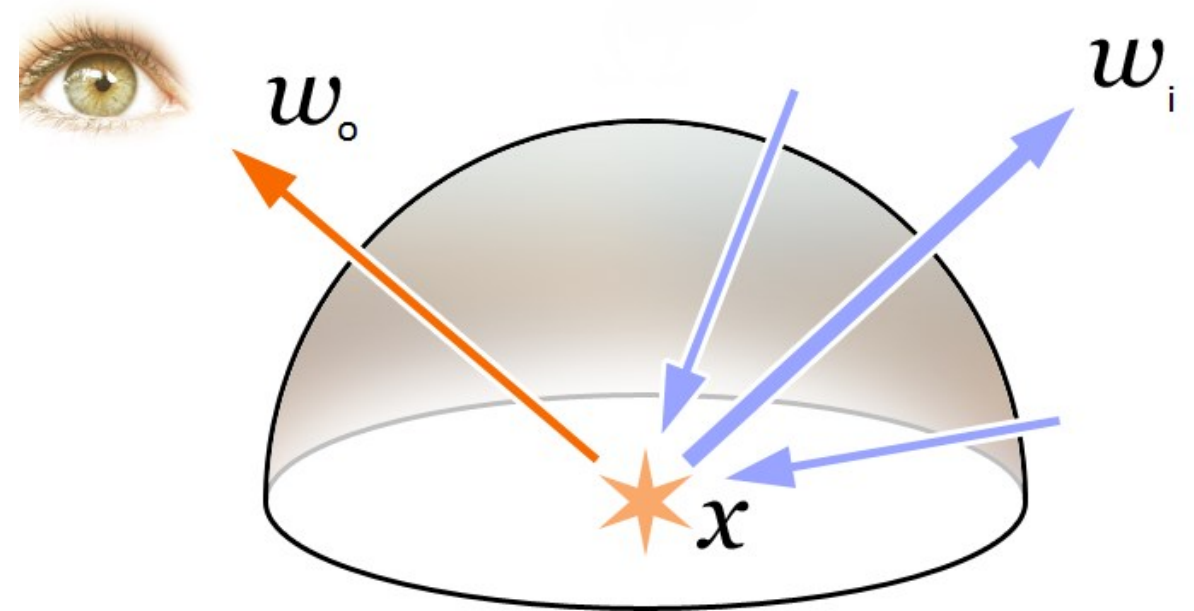
$L_e(\mathbf{x}, \vec{\omega}_o)$  = the light emitted from  $\mathbf{x}$  by this object itself

$f_r(\mathbf{x}, \vec{\omega}_i \rightarrow \vec{\omega}_o)$  = the BRDF of the surface at point  $\mathbf{x}$ ,  
transforming incoming light  $\omega_i$  to reflected light  $\omega_o$

$L(\mathbf{x}', \vec{\omega}_i)$  = light from  $\mathbf{x}'$  on another object arriving along  $\omega_i$

$G(\mathbf{x}, \mathbf{x}')$  = the geometric relationship between  $\mathbf{x}$  and  $\mathbf{x}'$

$V(\mathbf{x}, \mathbf{x}')$  = a visibility test, returns 1 if  $\mathbf{x}$  can see  $\mathbf{x}'$ , 0 otherwise



[https://en.wikipedia.org/wiki/Rendering\\_equation](https://en.wikipedia.org/wiki/Rendering_equation)

We make some assumptions: A diffused surface is illuminated by a distant light source  $L$ . We can rewrite the rendering equation:

$$L = f_r(x) \int_S L(w_i)(n \cdot w_i) dw_i$$

$$\text{Irradiance} = E(n) = \int_S L(w_i)(n \cdot w_i) dw_i$$

$$L = f_r(x) E(n)$$

Irradiance is parameterized by the surface normal only. It can be considered as a function defined over a sphere.



Environment map



Irradiance  
environment map



# Modeling Scene Illumination

## Irradiance as spherical harmonics

Each function defined on the surface of sphere can be written as a sum of spherical harmonics that form orthonormal basis.

Assuming Lambertian surface, we can write the irradiance as

$$E(n) = \sum_{l=0}^{\infty} \sum_{m=-l}^l c_{lm} Y_{lm}(n(x)) = Bi$$

$$L = f_r \odot Bi$$

$$L = Mi \longrightarrow \text{Solve using least squares method}$$

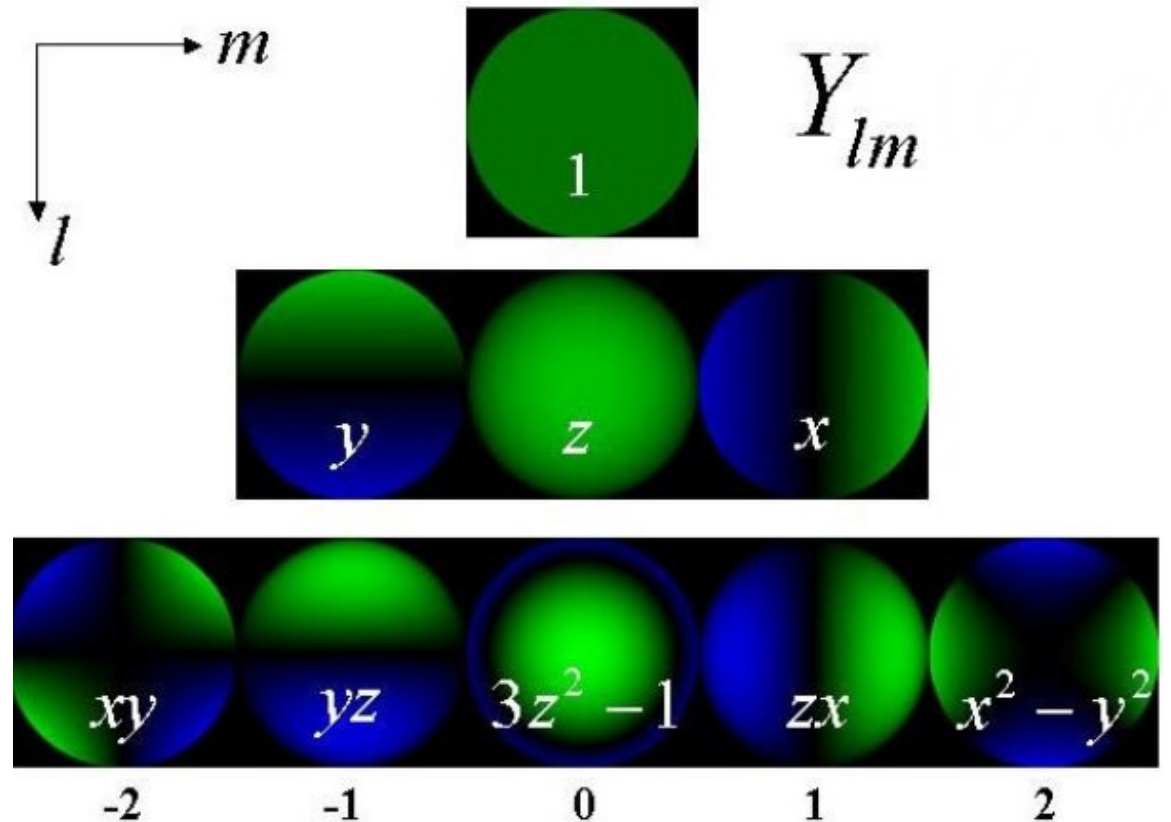
$c_{lm}$  = Spherical harmonic coefficients

$Y_{lm}$  = Spherical harmonic basis functions

$n(x)$  = Surface normal at point x

$i$  = Vector containing spherical harmonic coefficients

$M$  = Matrix containing spherical harmonic basis



The first 3 orders of spherical harmonics<sup>1</sup>

**Irradiance for a lambertian object can be well approximated by up to 2 orders of spherical harmonics<sup>1</sup>.**

- Low frequency estimation but fewer parameters
- Total 9 spherical harmonics for 1-channel images
- 27 for 3-channel images (9 for each channel)
- Spherical harmonics coefficients can be used as a proxy for scene lighting
- On Right, a lambertian sphere is lightened using spherical harmonic coefficients



An exemplary SH illumination

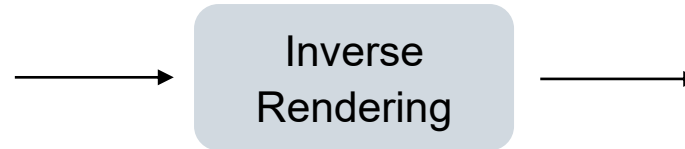
# Inverse Rendering

# Inverse Rendering

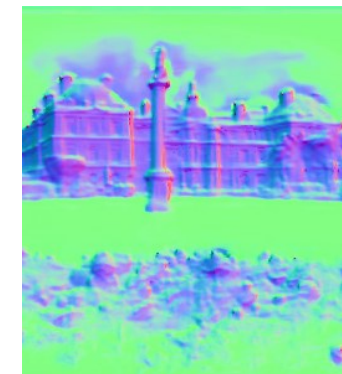
Decompose an image into albedo, surface normal, scene illumination, and shadows



Original Image



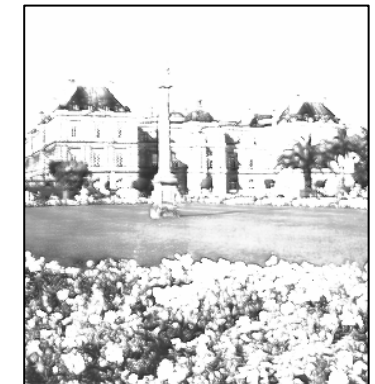
Albedo map



Surface normal map



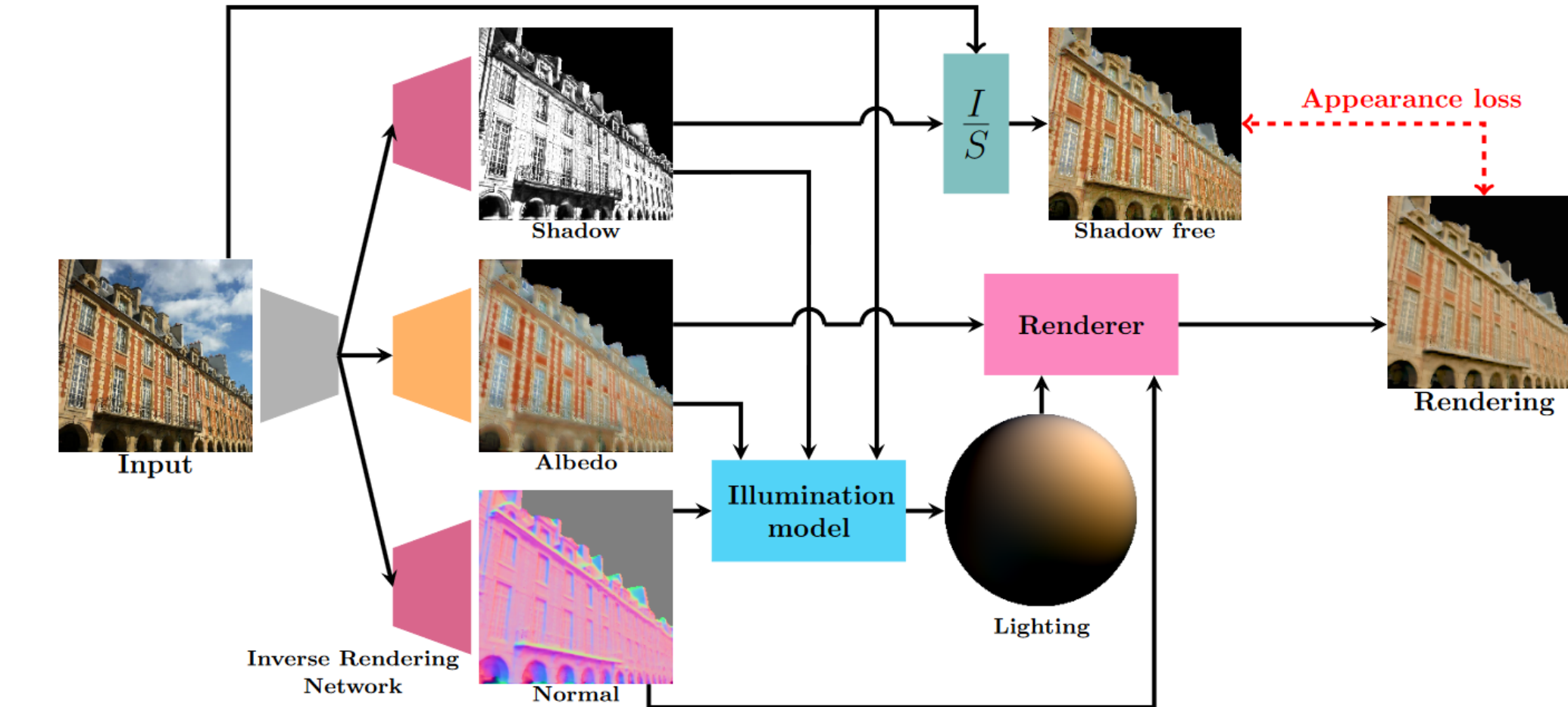
Scene illumination



Shadow map

# Inverse Rendering Network (IRN)

Directly regress albedo, surface normal and shadow maps from image using autoencoder



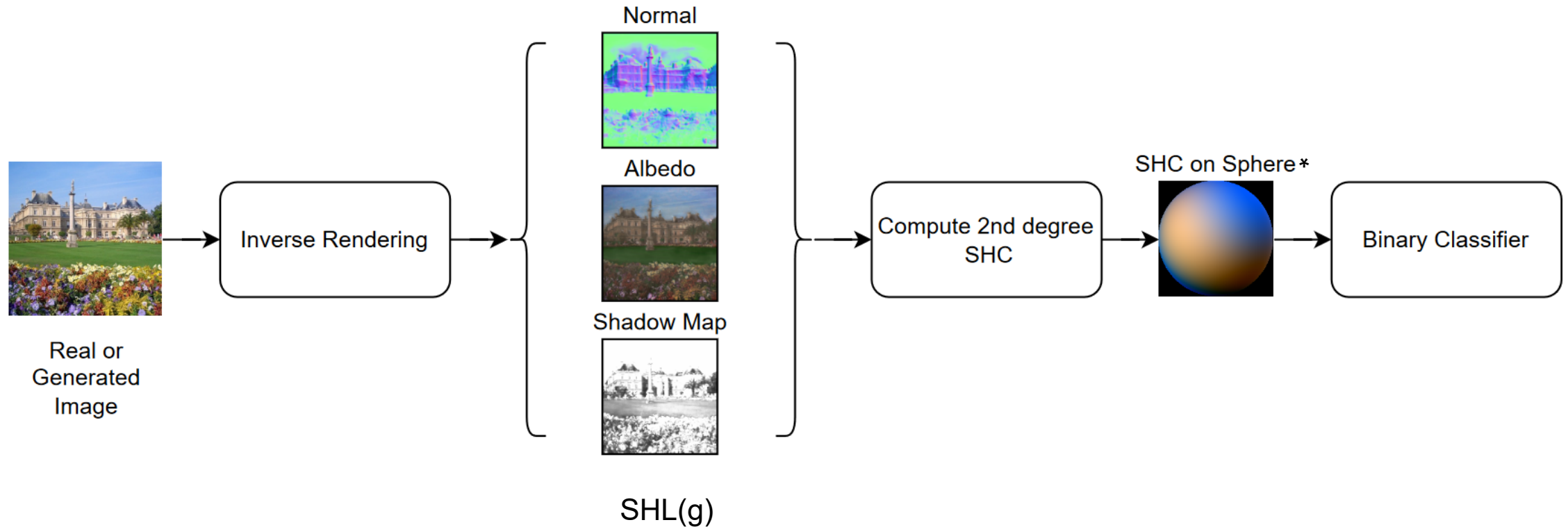
Inverse rendering network architecture<sup>1</sup>

---

# Methodology

# Proposed Pipeline

Compute 27 spherical harmonics coefficients for the whole image and pass them to different binary classifiers

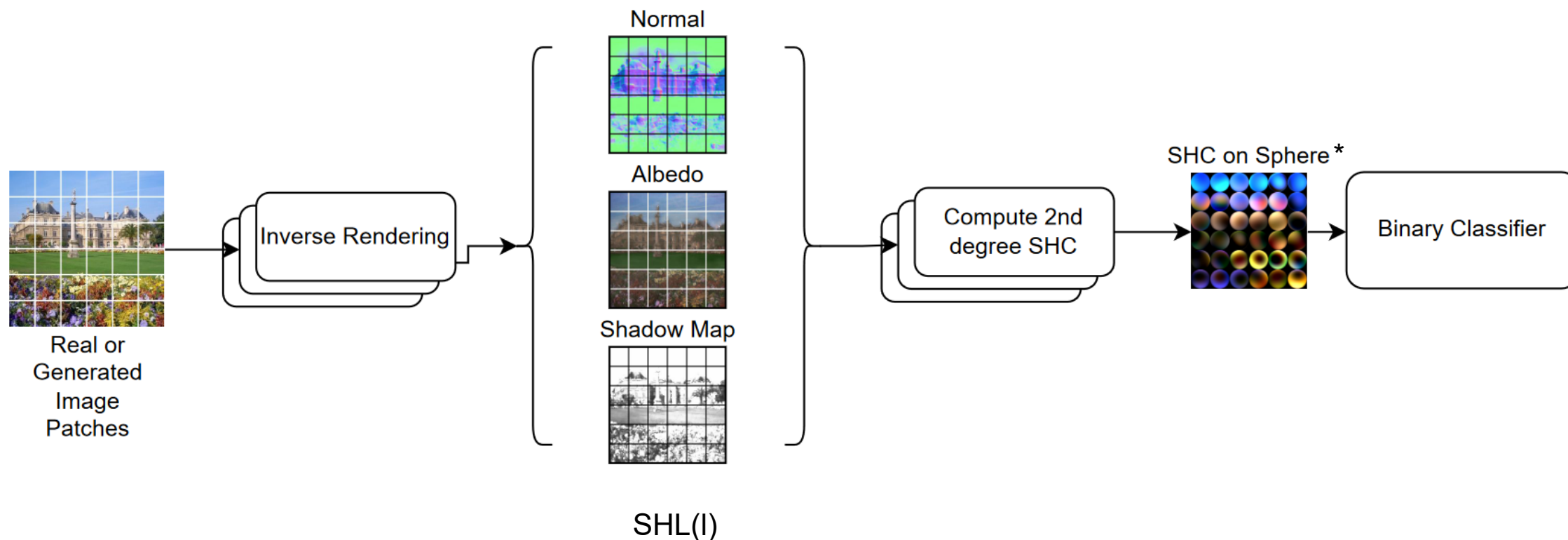


\* This image is for visualisation only. We directly pass 27 spherical harmonics coefficients to the binary classifier.



# Proposed Pipeline

Divide image into patches and compute spherical harmonics coefficients for each patch and pass them to the binary classifier



\* This image is for visualisation only. We accumulate 27 spherical harmonics coefficients for each patch and pass them to the binary classifier.

# Binary Classifiers

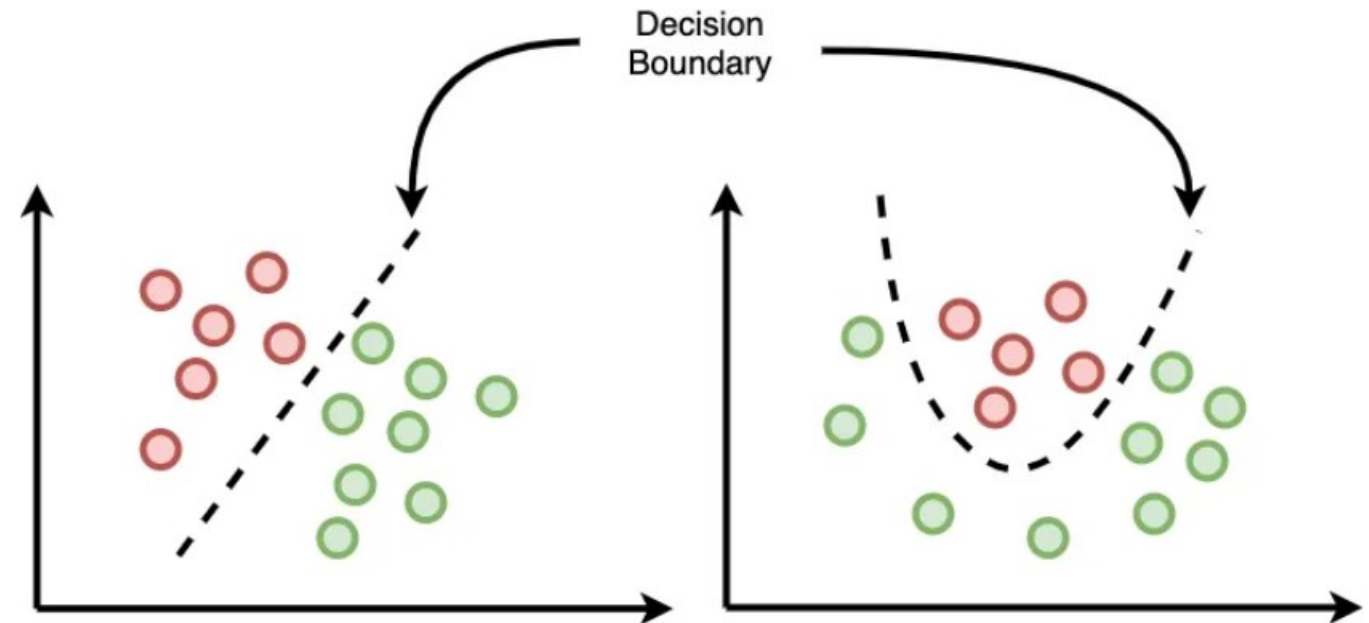
Different linear and non-linear binary classifiers are trained for classification

Traditional classifiers:

- Random Forest (RF)
- Support Vector Machine (SVM)
- Logistic Regression (LR)

Neural network-based classifiers:

- MLP
- Modified Vision Transformer (ViT)



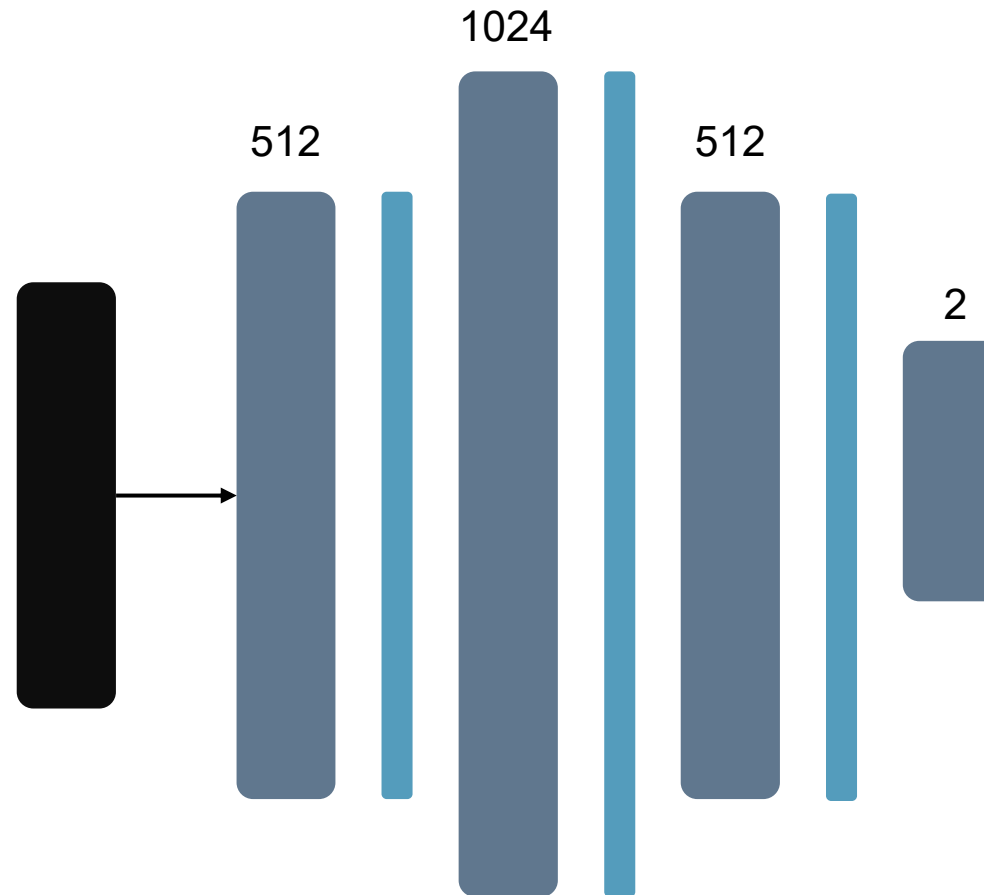
We use linear/non-linear binary classifiers to explore different decision boundaries.

<https://towardsdatascience.com/logistic-regression-and-decision-boundary-eab6e00c1e8>

# Binary Classification with MLP

4-layer MLP

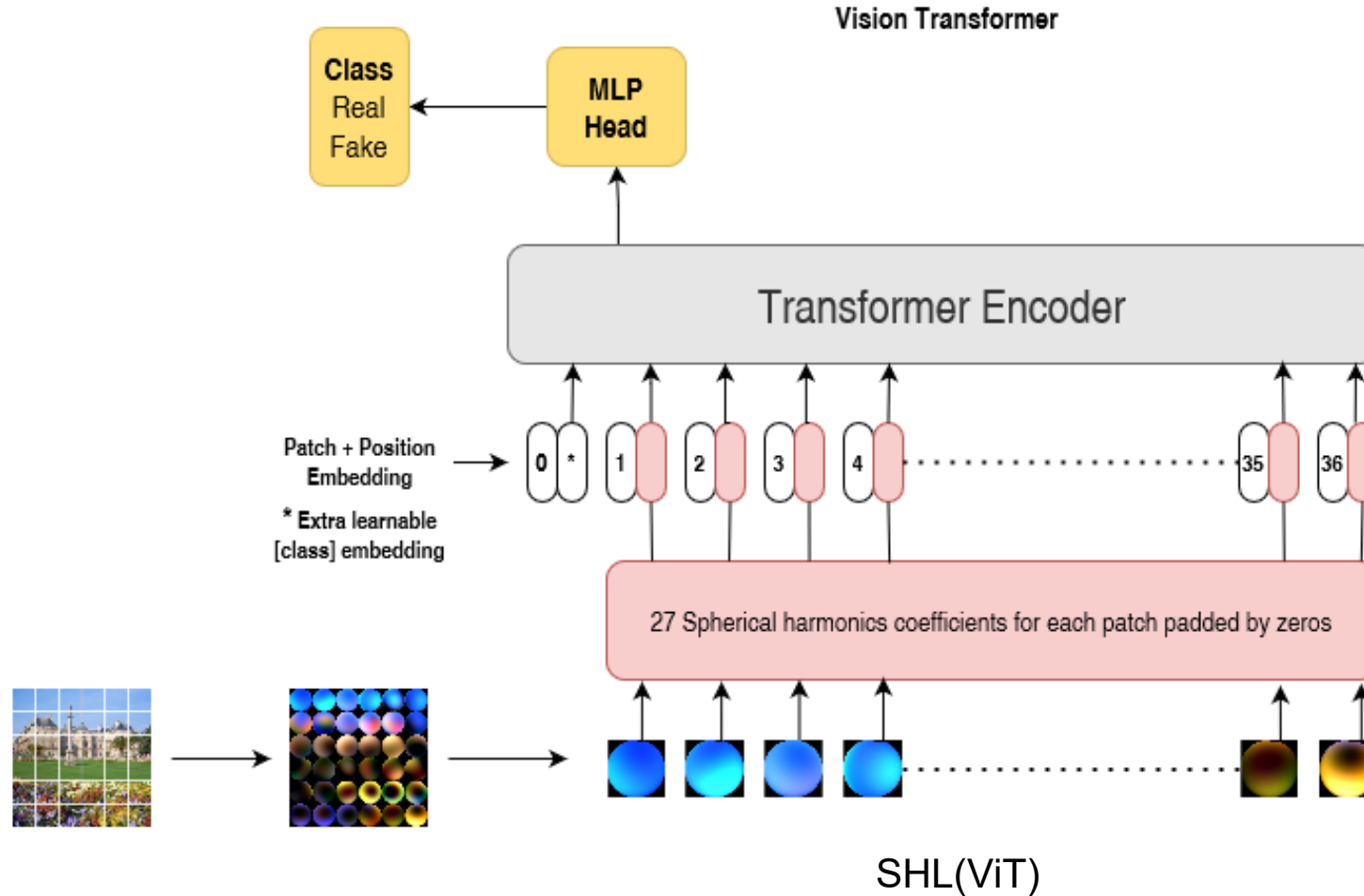
- Input features
- Fully-connected + ReLU
- Dropout ( $p=0.5$ )



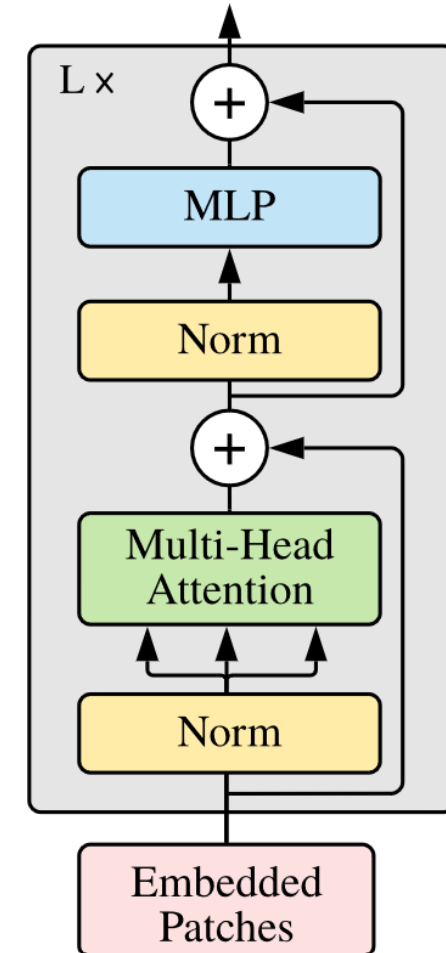
MLP architecture

# Binary Classification with ViT

Modified ViT: Instead of image patched, pass spherical harmonics for each patch



## Transformer Encoder



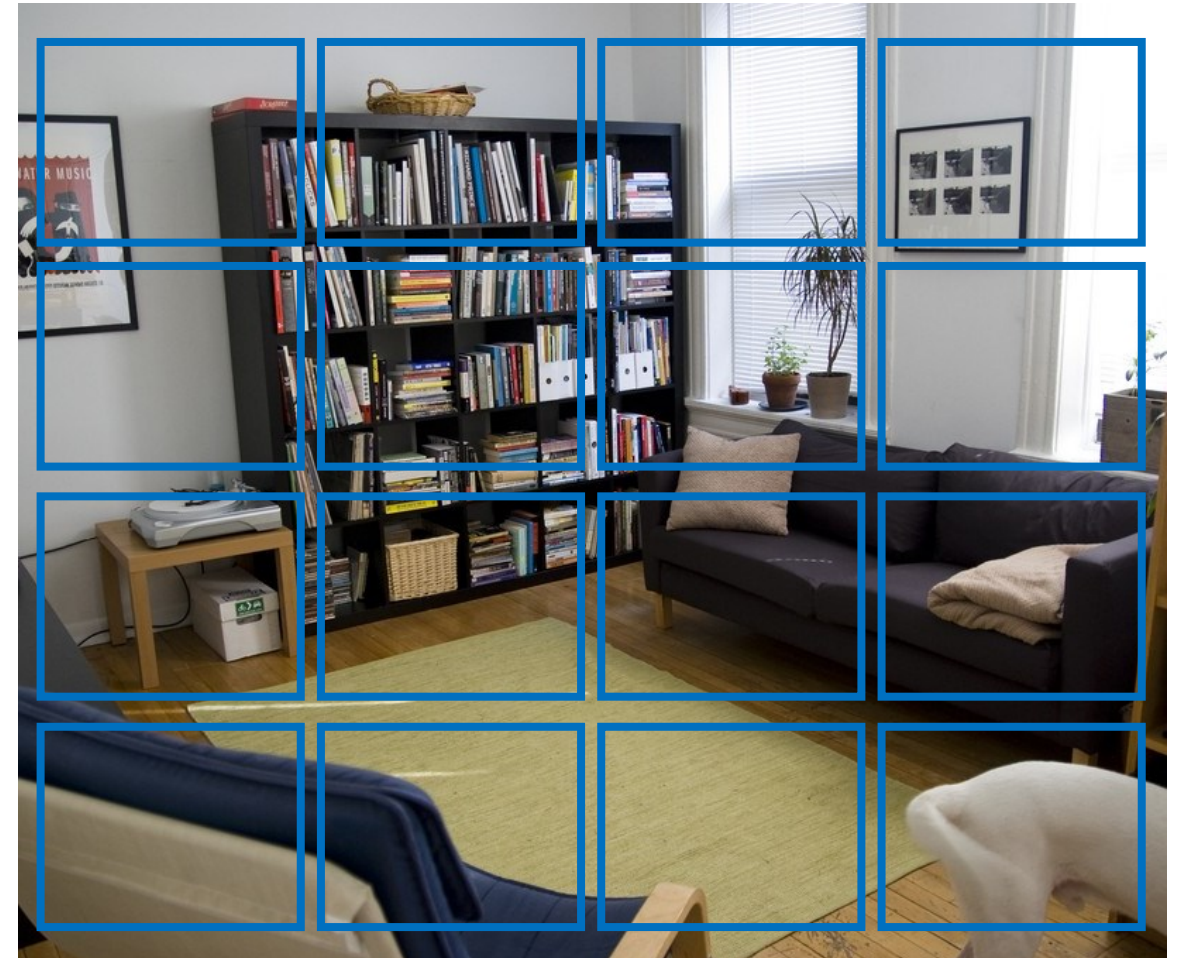
# Baseline Methods

# GIST Features as a Baseline

Learn scene-level representation instead of segmentation or object-level representations

## Estimate structure of a scene by computing spectral signature using a few perceptual scene properties<sup>1</sup>

- Scene properties include openness, roughness, expansion, ruggedness
- Spectral templates for different values of scene properties that are learned from real images
- 32 spectral templates (8 for each scene property)
- Divide image into 4x4 patches
- Average filter response for each filter and patch
- Total 512 features for a single image

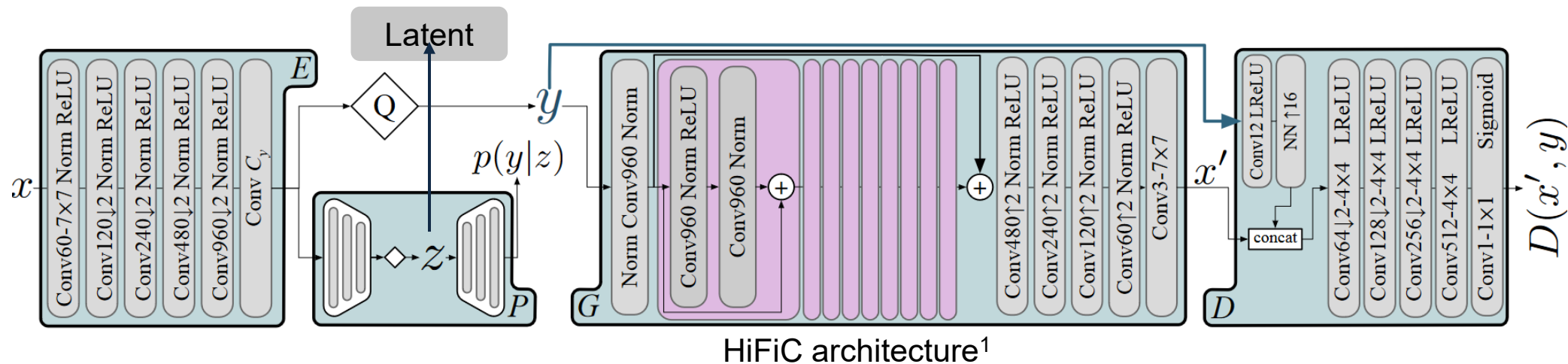


# Latent Space of a Compression Network

Real and generated images have different latent space representation

Compression network project input image into a very compact latent representation. HiFiC is one such compression network

- Reconstructs perceptually similar images at extremely low bit-rate
- An input image is projected into a latent representation using an encoder
- A conditional-GAN reconstruct the image from the latent representation
- We only use latent representation



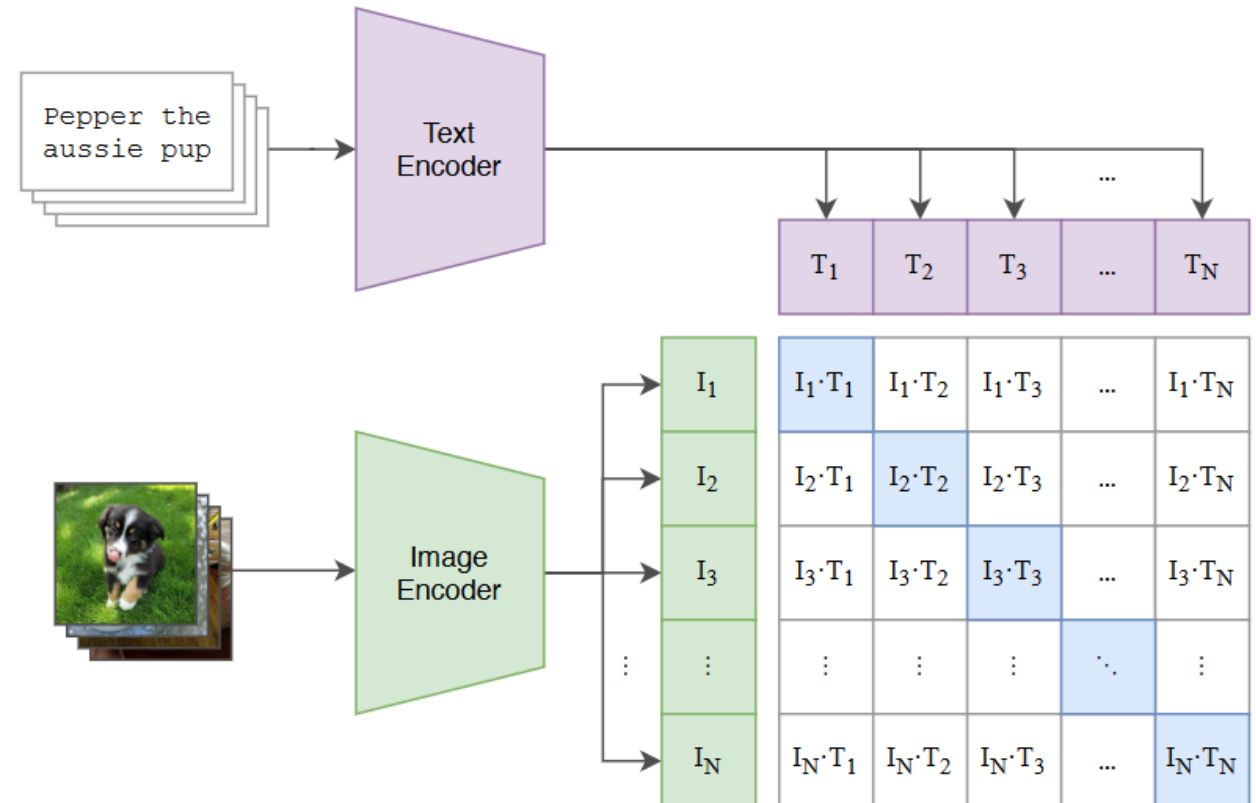


# Contrastive Language-Image Pretraining (CLIP)

CLIP features with linear classifiers have better cross-generalization performance

CLIP is trained on large amount of real image-caption pairs collected from the internet.

- Separate text and image encoder
- Minimize cosine similarity between two embedding
- Contrastive loss function
- CLIP image features have good cross-generalization performance<sup>1</sup>.



CLIP architecture

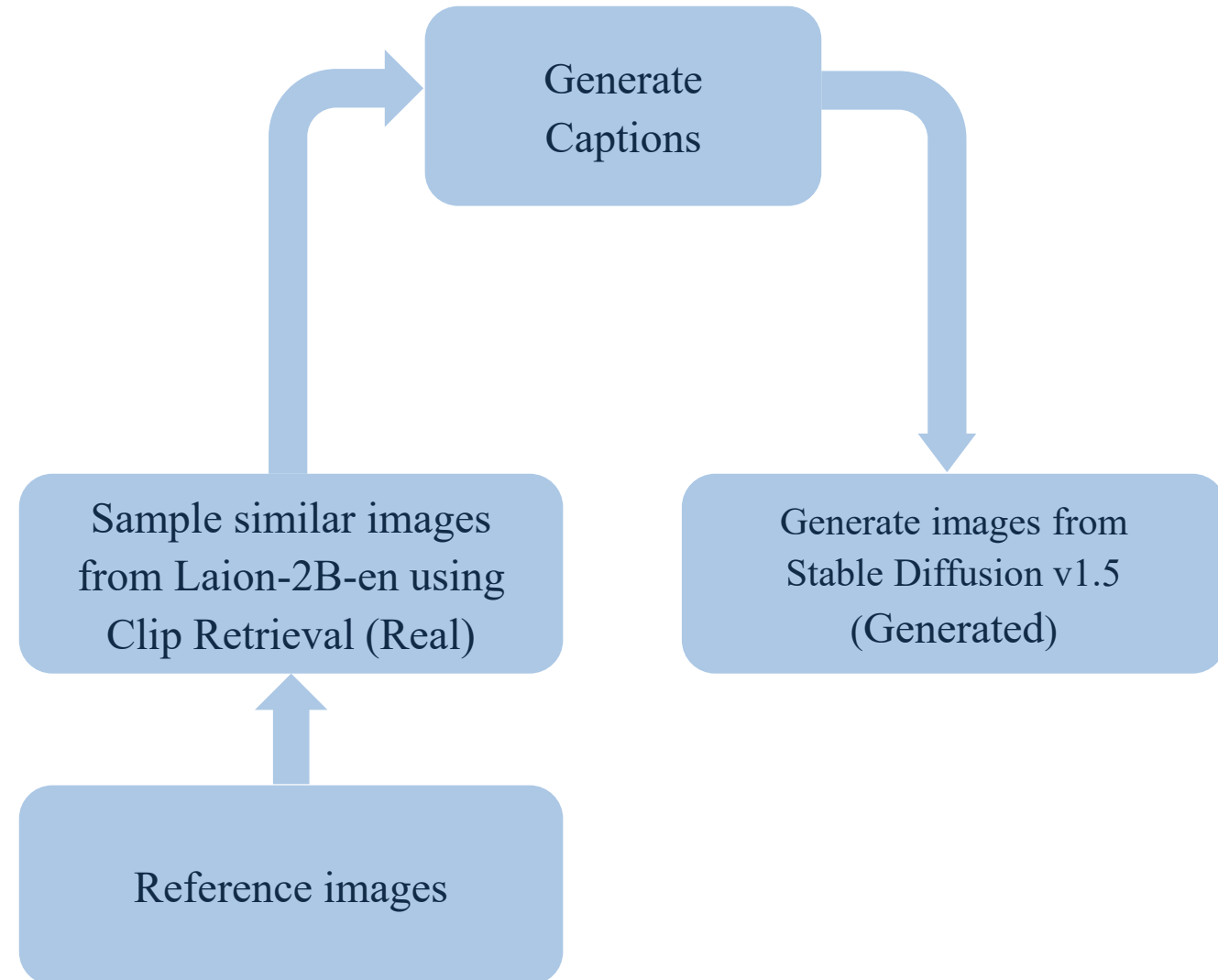
# Dataset Creation

# Dataset Creation

Real and stable diffusion generated images are conditioned on common captions

## Dataset creation pipeline

- Select few reference images manually
- Search similar images based on clip embedding of reference images from LAION-2B-en dataset ([rom1504.github.io/clip-retrieval/](https://rom1504.github.io/clip-retrieval/))
- Collected images are assumed **real**. Generate captions using **microsoft/git-large-r-coco** image captioning model
- Use captions to generate images from SD1.5 model. Append special keywords for photorealistic generation
- All images are JPEG compressed
- Remove white background images



# A Few Reference Images

Reference images are used to collect real images from the LAION-2B-en dataset



All images are real.

# Real and Generated Samples

Real and stable diffusion generated images are conditioned on common captions

*A white chair in a white room*



An image sampled from LAION-2B-en dataset (assumed real scene)



A stable diffusion v1.5 generated image



# Real and Generated Samples

Real and stable diffusion generated images are conditioned on common captions

*A bunch of pears sitting on top of a wooden table.*



An image sampled from LAION-2B-en dataset (assumed real scene)



A stable diffusion v1.5 generated image

# Data Splits For Training and Evaluation

Distribution of real and generated images (Stable Diffusion v1.5) across data splits

Split	Real Images	Generated Images	Total Images
Train	5123*	5501	10624
Validation	665	679	1344
Test	651	719	1370

\* Some of the real images were removed from the dataset as they contained objects with backgrounds removed.



# Cross-generator Evaluation Dataset

Distribution of real and generated images in cross-generator dataset

Dataset	Images
RAISE-1k	1000

Real Images

Dataset	Images	Architecture
DALL-E2 <sup>*</sup>	1000	open
DALL-E3 <sup>*</sup>	1000	closed
Firefly <sup>*</sup>	1000	closed
Midjourney <sup>*</sup> v5	1000	closed
SD v1.3 <sup>*</sup>	1000	open
SD v1.4 <sup>*</sup>	1000	open
SD v1.5	1000	open
SD v2 <sup>*</sup>	1000	open
SD XL <sup>*</sup>	1000	open

AI-Generated Images

Open: Architectural  
detail shared

Closed: Architectural  
detail are not fully  
disclosed

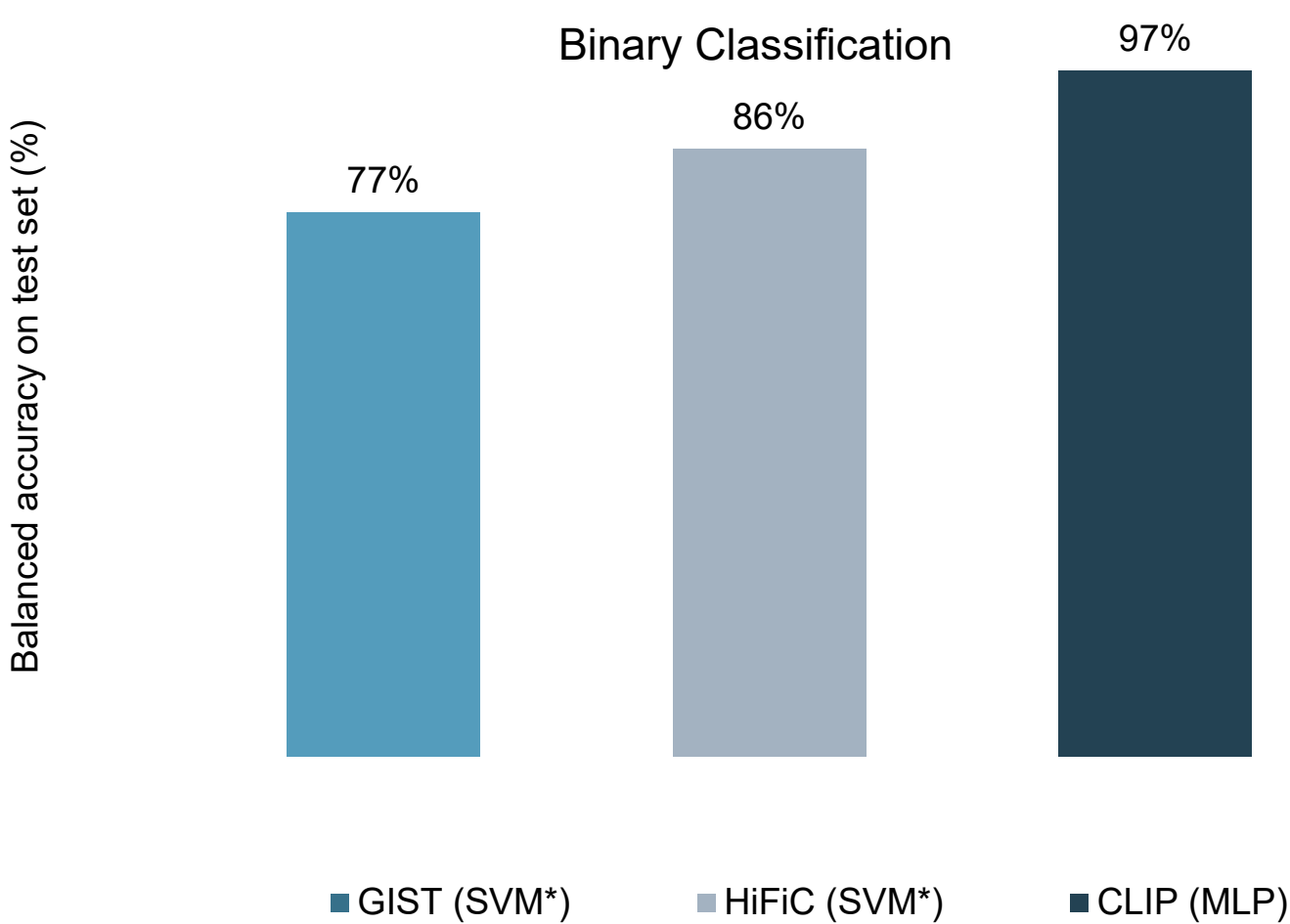
<sup>\*</sup> Images taken from  
the Synthbuster<sup>1</sup>  
dataset

---

# Experiments and Results

# Baseline Binary Classification Results

Balanced classification accuracy for different baseline methods on test set



## GIST Confusion Matrix

	Pred Real	Pred Gen
Real	496	151
Gen	167	546

## HiFiC

	Pred Real	Pred Gen
Real	517	130
Gen	65	648

## CLIP

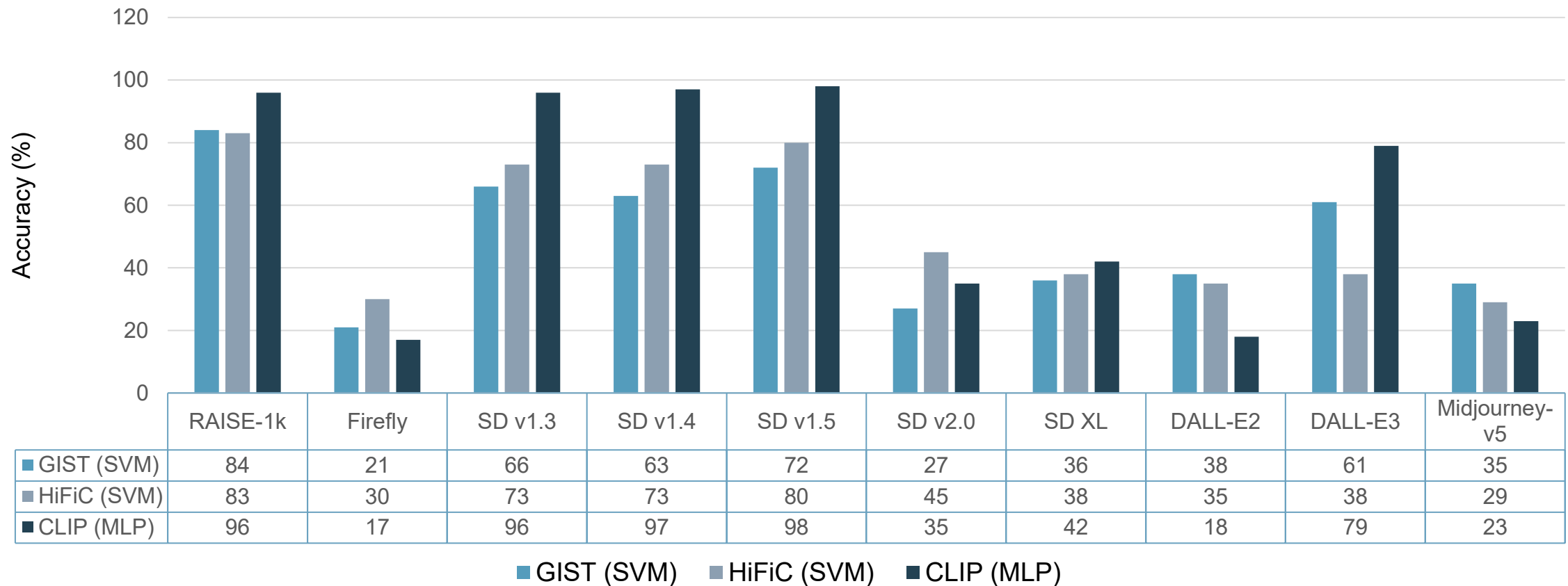
	Pred Real	Pred Gen
Real	617	21
Gen	16	690

\* SVM with RBF Kernel

# Cross-generator Generalization

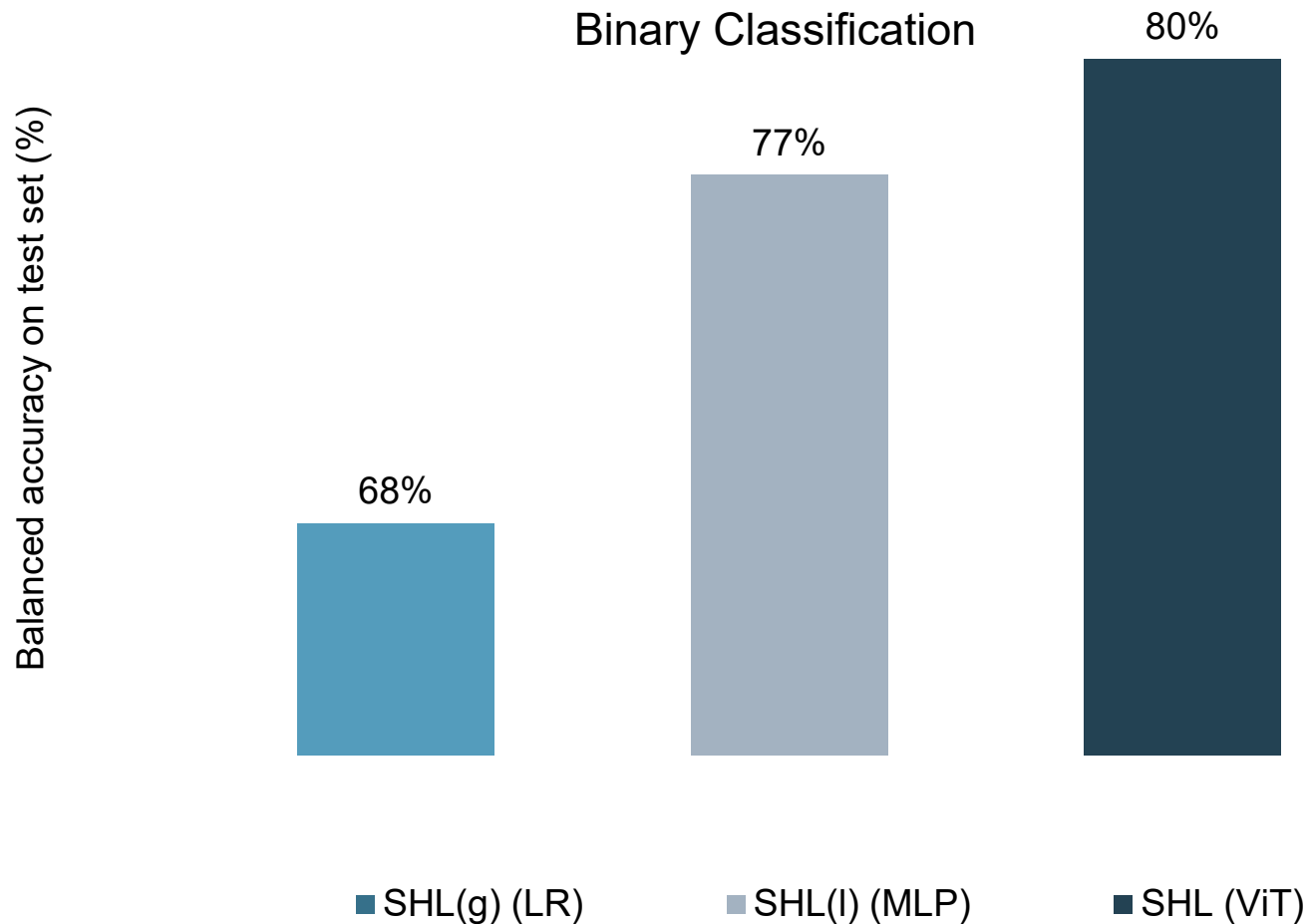
Evaluation on cross-generator dataset

## Cross-generator performance of baseline methods



# Binary Classification Results for Proposed Method

Balanced classification accuracy for proposed methods on test set



## SHL(g) Confusion Matrix

	Pred Real	Pred Gen
Real	401	246
Gen	189	524

## SHL(I)

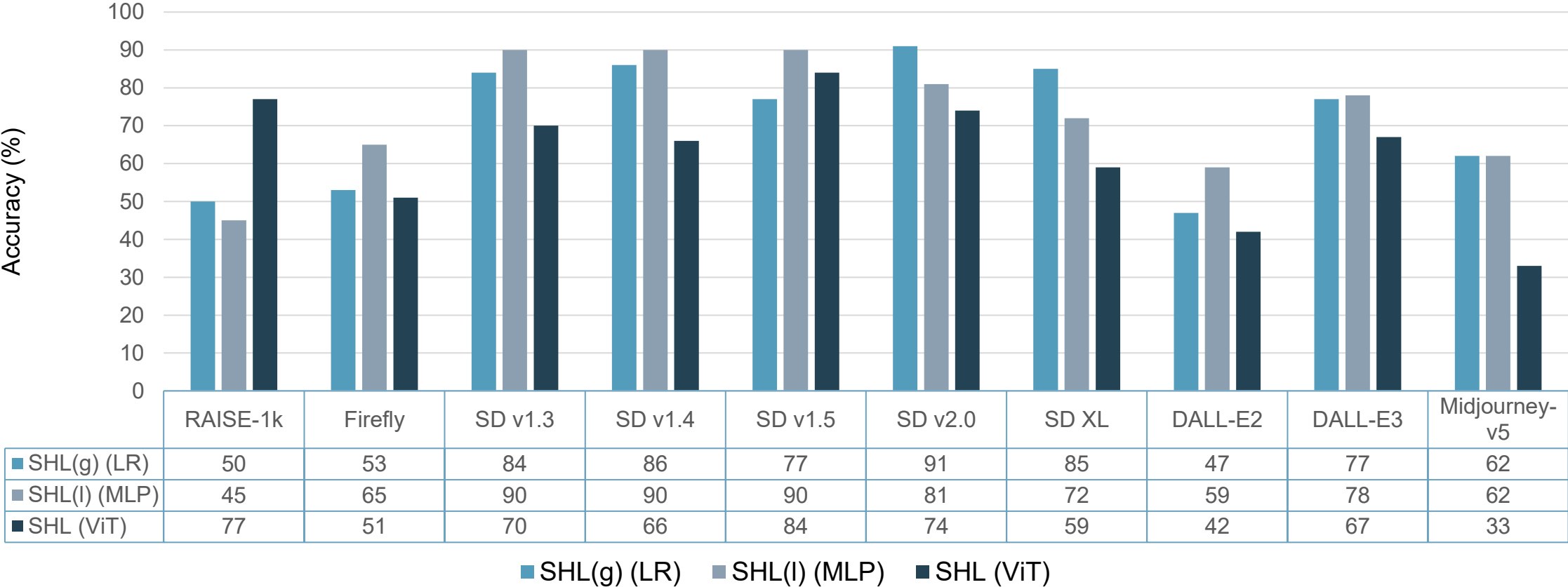
378	260
52	654

## SHL (ViT)

469	181
68	642

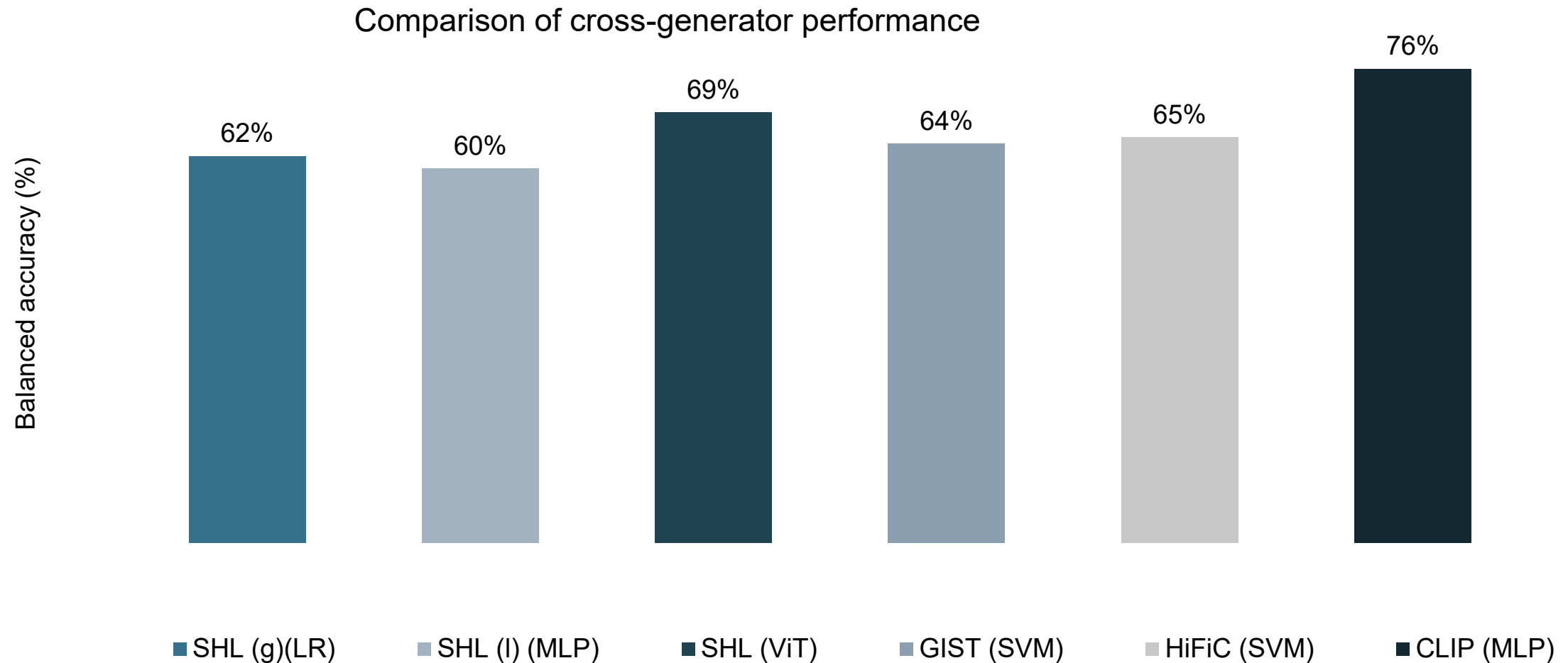
\* SVM with RBF Kernel

Cross-generator performance of our proposed methods



# Comparison Cross-generator Performance

Comparison between baseline methods and proposed method for cross-generator performance





# Qualitative Analysis of SHL (ViT)

# Some Misclassified Examples

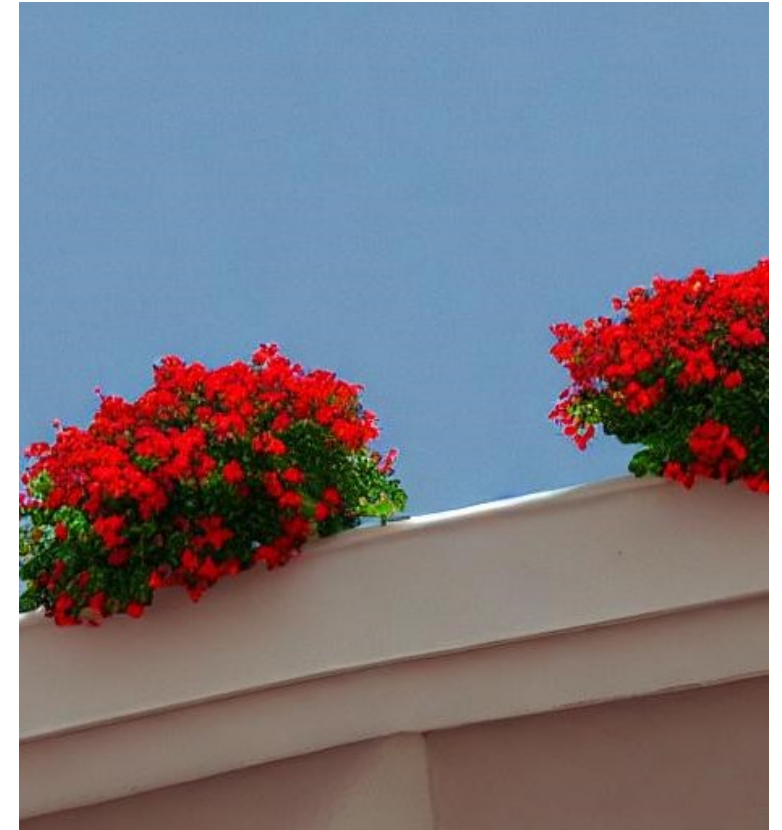
A few generated images taken from a smaller subset (50 images, manually picked)  
of test set generated images



Generated but predicted real



Generated but predicted real



Generated but predicted real



# Some Misclassified Examples

A few real images taken from a smaller subset (50 images, manually picked) of test set real images



Real but predicted generated



Real but predicted generated

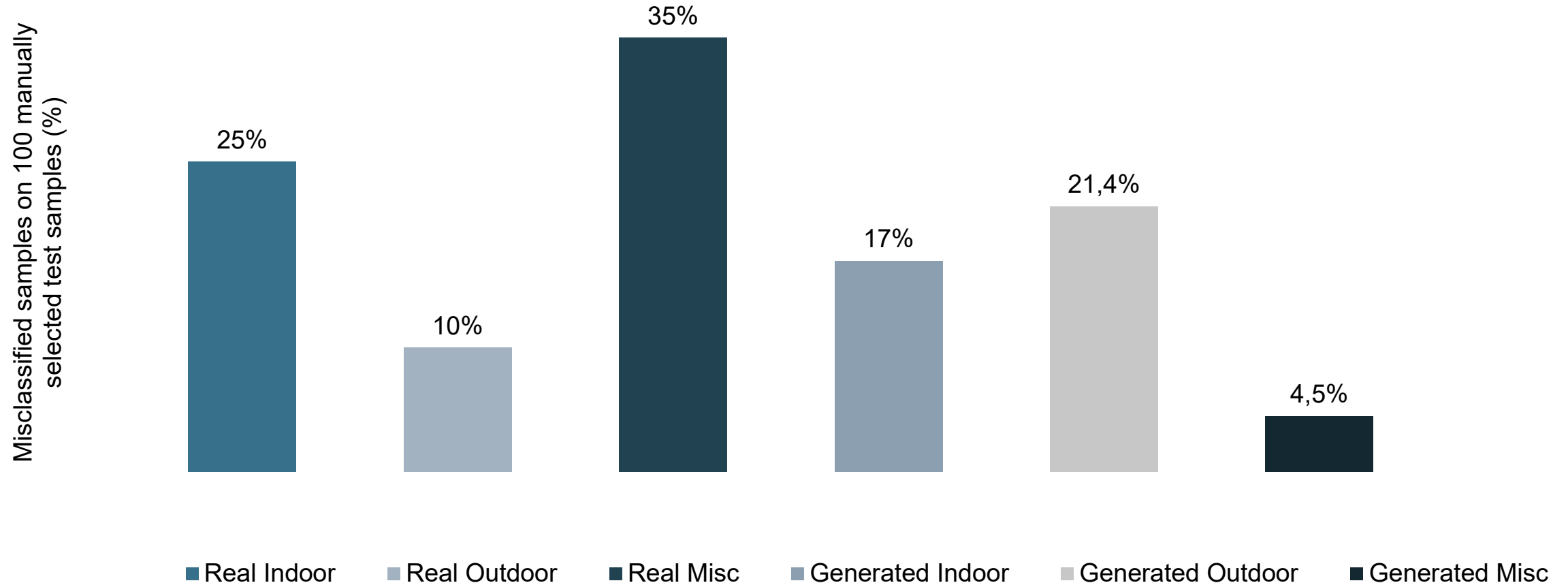


Real but predicted generated

# Categorization of Misclassified Samples

Evaluation on manually picked 100 images from test set (Nearly 17 images for each category)

Misclassification rate for different categories using SHL (ViT)



---

# Conclusion and Future Directions

---

## Conclusion

- Spherical harmonics representation of scene illumination can be used to detect real and generated images
- ViT architecture can capture illumination inconsistencies
- Our proposed method has ability to generalize on unseen image generators.

## Future outlook:

- Increase the training set
- Include more diverse generators in cross-generator evaluation dataset
- Experiment with recently proposed better inverse rendering pipelines

**Thank you  
for your attention!**