

# Text2Mesh: Text-Driven Neural Stylization for Meshes

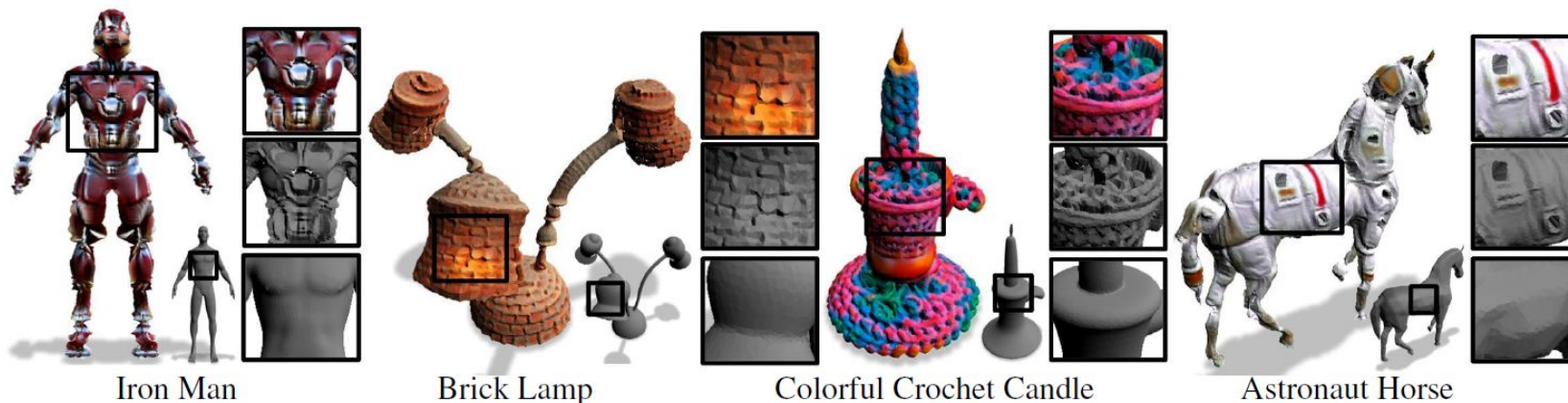


Figure 1. Text2Mesh produces color and geometric details over a variety of source meshes, driven by a target text prompt. Our stylization results coherently blend unique and ostensibly unrelated combinations of text, capturing both global semantics and part-aware attributes.

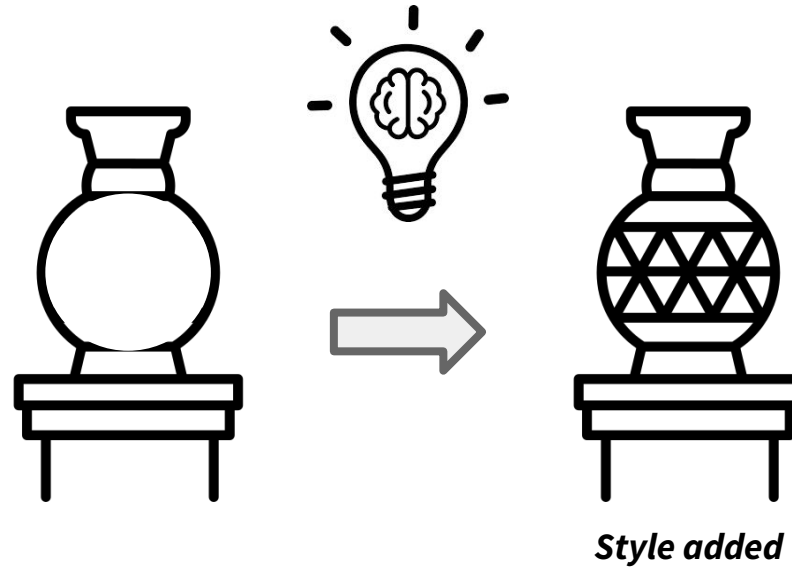


**Project ID: 14**  
**Mooyeol Oh & Muhammad Izaaz Inhar Ramahdani**  
**Development Track**  
[Michel et al., arXiv 2021]



# Motivation

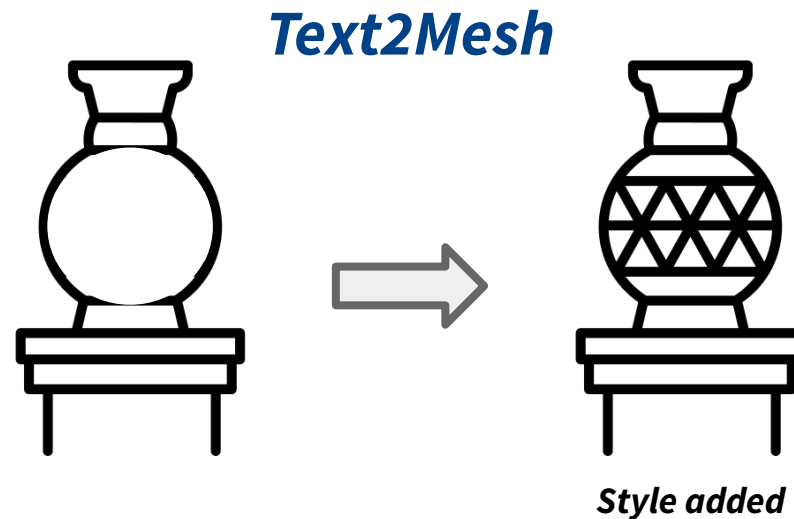
*Intuitive control the 3D object by using natural language*



- **Editing** visual data to conform to a desired style, while **preserving the underlying content**
- To propose expressing the desired style through **natural language** (a text prompt), similar to how a commissioned **artist** is provided a verbal or textual description of the desired work.

# Introduction

*Intuitive control the 3D object by using natural language*

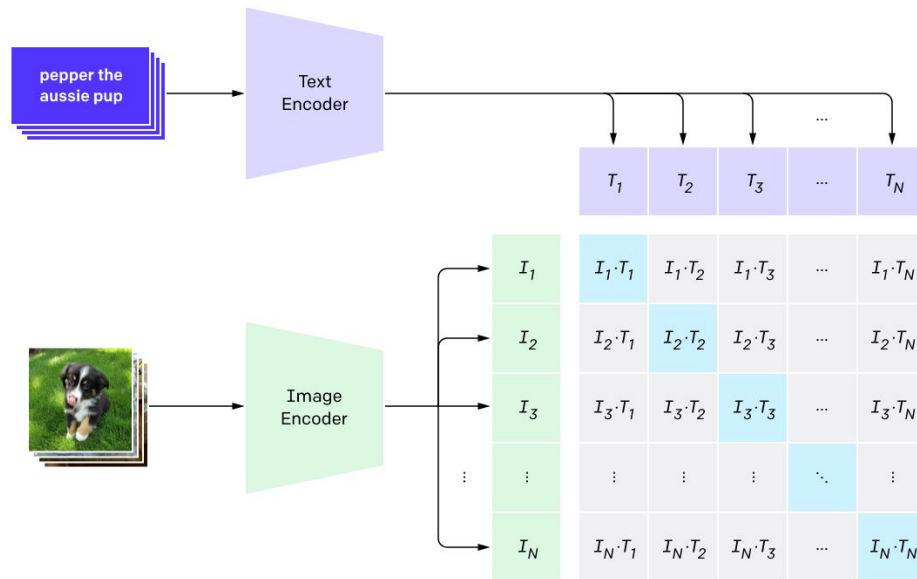


- **Text2Mesh** present a technique for the *semantic manipulation of style* for 3D meshes, harnessing the representational power of **CLIP**.
- This system combines the *advantages* of **explicit mesh surfaces** and the *generality of neural fields* to facilitate intuitive control for stylizing 3D shapes.

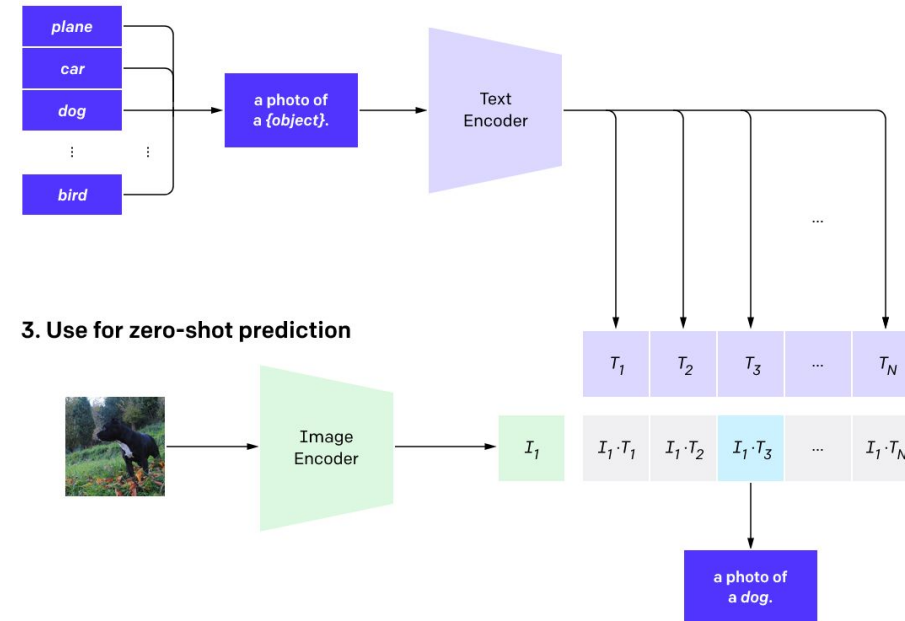
# CLIP (Contrastive Language–Image Pre-training)

*CLIP builds on a large body of work on  
zero-shot transfer, natural language supervision, and multimodal learning*

## 1. Contrastive pre-training



## 2. Create dataset classifier from label text



## 3. Use for zero-shot prediction

CLIP (Contrastive Language–Image Pre-training) learns a joint embedding space for images and text

Radford, A., Kim, J. W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., ... & Sutskever, I. (2021, July). Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning* (pp. 8748-8763). PMLR.

# Necessity of Network

*Ablation on the priors used in Text2Mesh method (full) for a candle mesh  
and target 'Candle made of bark'*



*full*

VS



*-net  
(without network)*

However, a ***straightforward*** optimization of the 3D stylized mesh which maximizes the CLIP similarity score converges to a ***degenerate*** (i.e. noisy) ***solution*** → Employing ***CLIP*** for stylization requires ***careful regularization***

# Related Work

- **Text driven manipulation:**

This work to image manipulation techniques controlled through textual descriptions embedded by CLIP [1](Alex et.al). CLIP learns a joint embedding space for images and text, and other works have incorporate CLIP as means for text-guided image generation.

- **Geometric Style Transfer in 3D:**

Some approaches analyze 3D shapes and identify similarly shaped geometric elements and parts which differ in style [2](Ruizhen et.al). Others transfer geometric style based on content/style separation [3](Xu et.al)

- **Texture Transfer in 3D:**

Aspects of a 3D mesh style can be controlled by texturing a surface through mesh parameterization [4](Mark et.al), but recent work explored a neural representation of texture[5] (Nicholas et.al)

[1] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. arXiv preprint arXiv:2103.00020, 2021.

[2] Ruizhen Hu, Wenchao Li, Oliver Van Kaick, Hui Huang, Melinos Averkiou, Daniel Cohen-Or, and Hao Zhang. Colocating style-defining elements on 3d shapes. ACM Transactions on Graphics (TOG), 36(3):1–15, 2017.

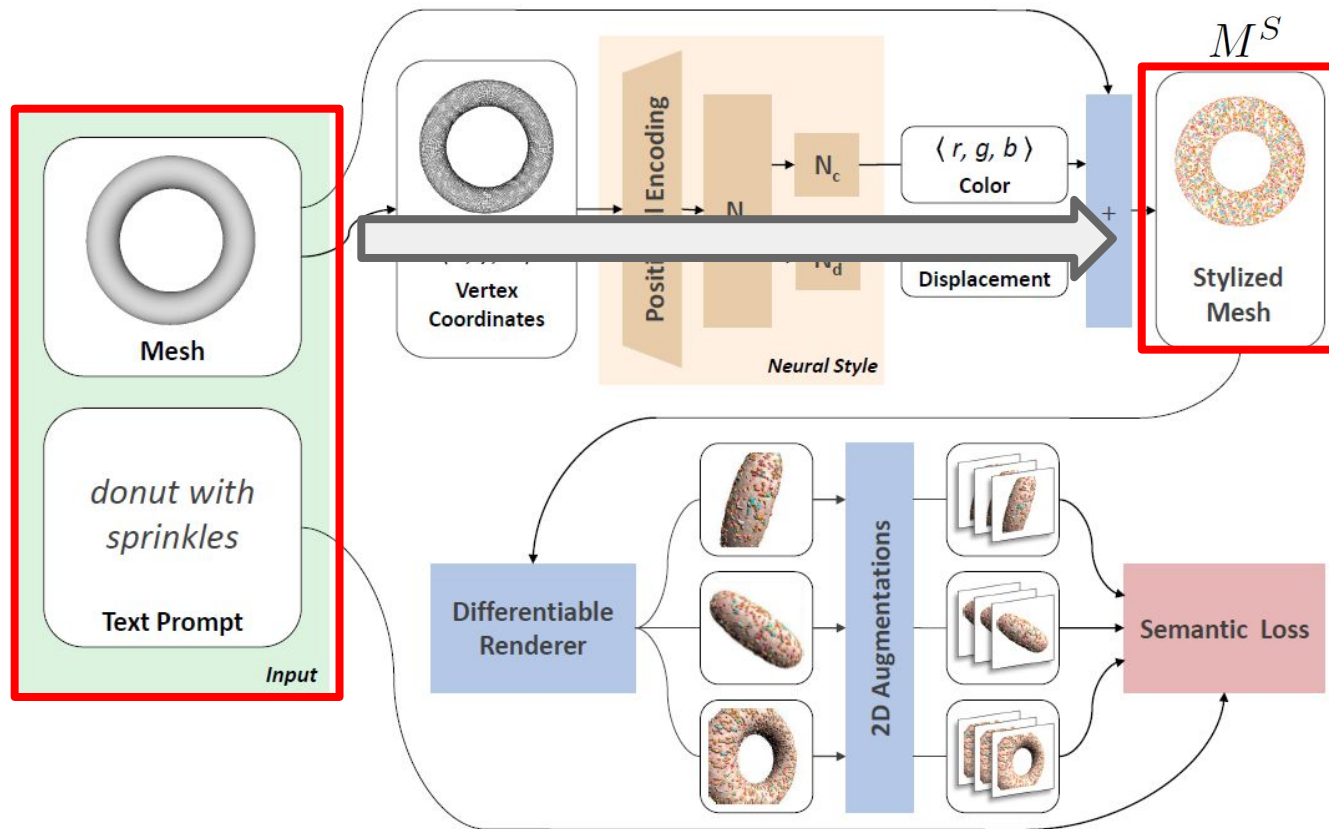
[3] Xu Cao, Weimin Wang, Katashi Nagao, and Ryosuke Nakamura. Psnet: A style transfer network for point cloud stylization on geometry and color. In Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision, pages 3337–3345, 2020.

[4] Mark Gillespie, Boris Springborn, and Keenan Crane. Discrete conformal equivalence of polyhedral surfaces. ACM Transactions on Graphics (TOG), 40(4):1–20, 2021.

[5] Nicholas Sharp. Intrinsic Triangulations in Geometry Processing. PhD thesis, Carnegie Mellon University, August 2021.

# Overall Architecture

Figure 4.

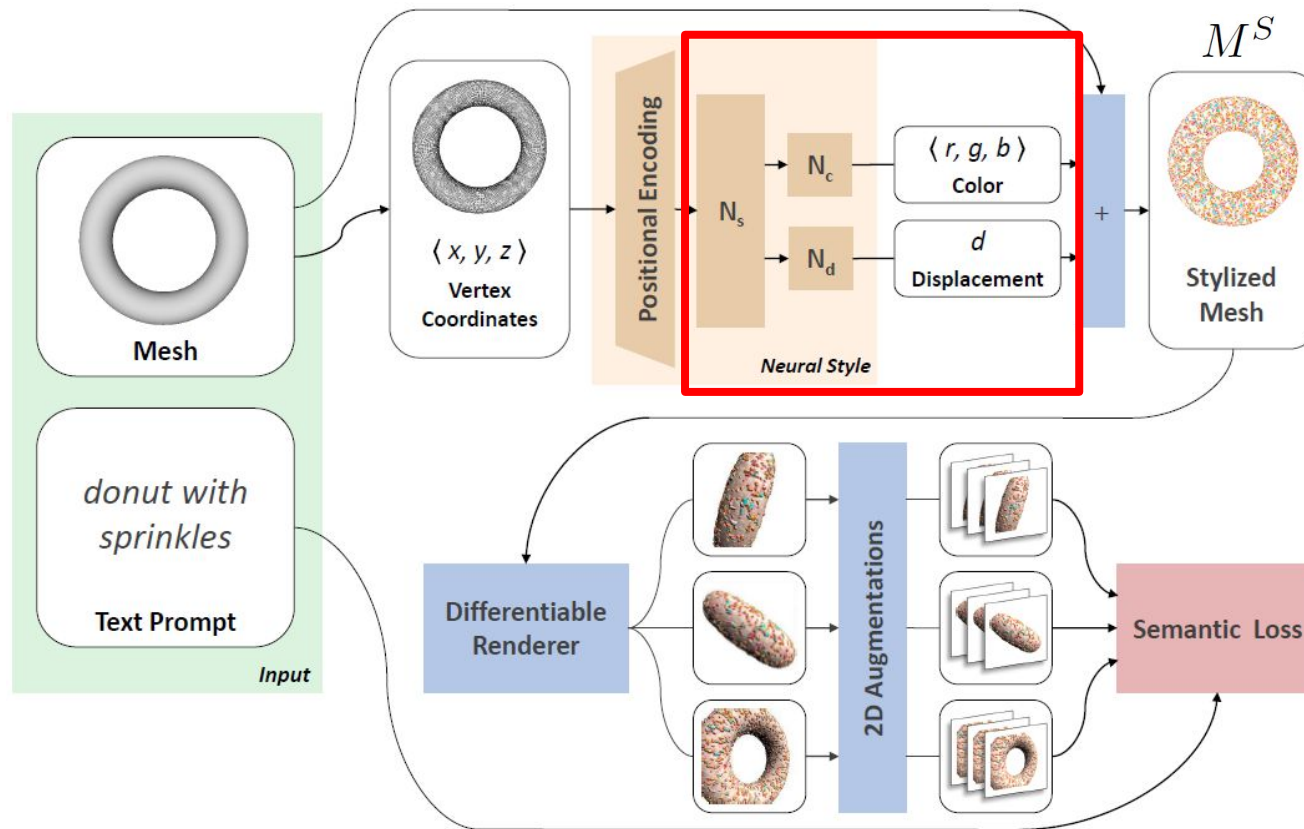


- The object's style (color and local geometry) is **modified** to conform to a **target text prompt**  $t$ , resulting in a **stylized mesh**  $M^S$



# Overall Architecture

Figure 4.

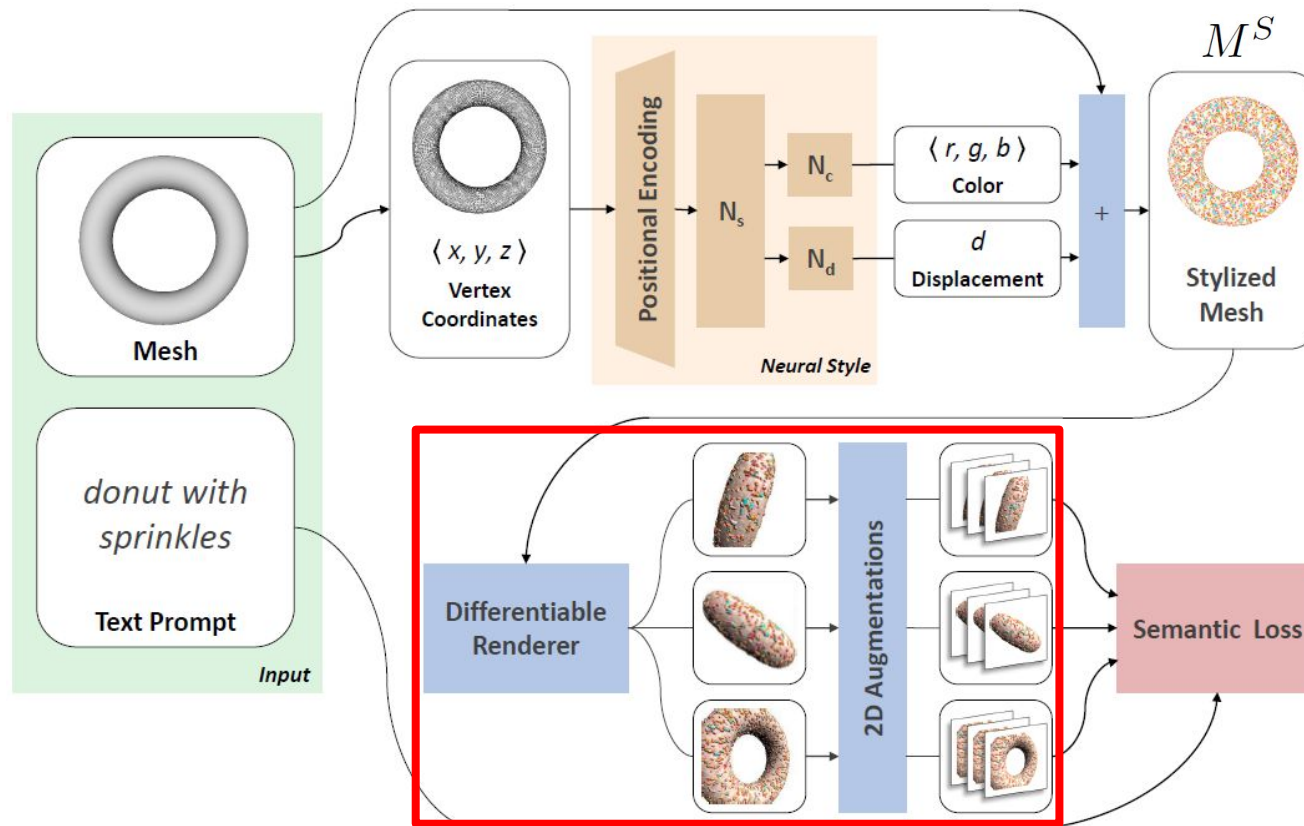


- The **NSF** (Neural Style Field) Network *learns to map points* on the mesh surface to an **RGB color** and **displacement** along the normal direction.



# Overall Architecture

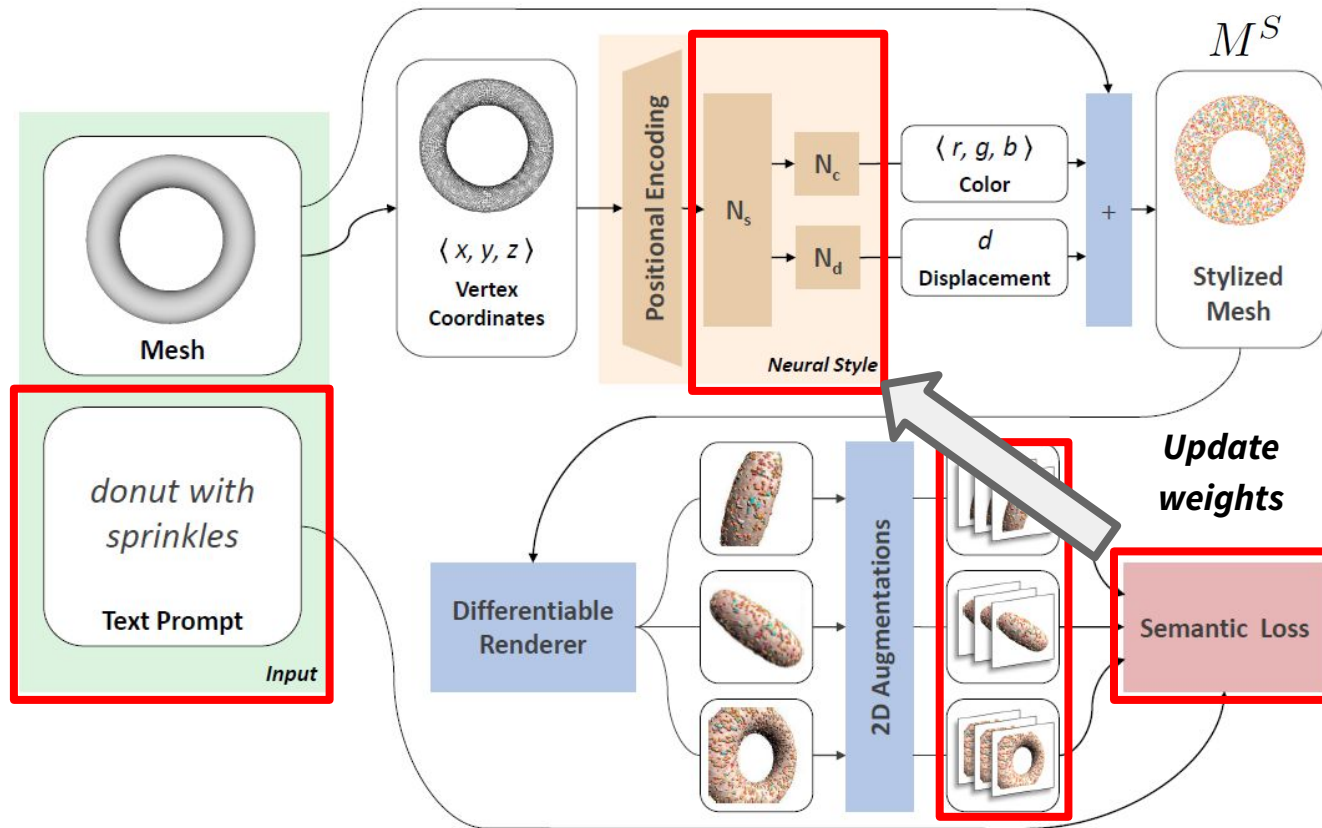
Figure 4.



- Render  $M^S$  from multiple views and apply 2D augmentations that are embedded using **CLIP**.

# Overall Architecture

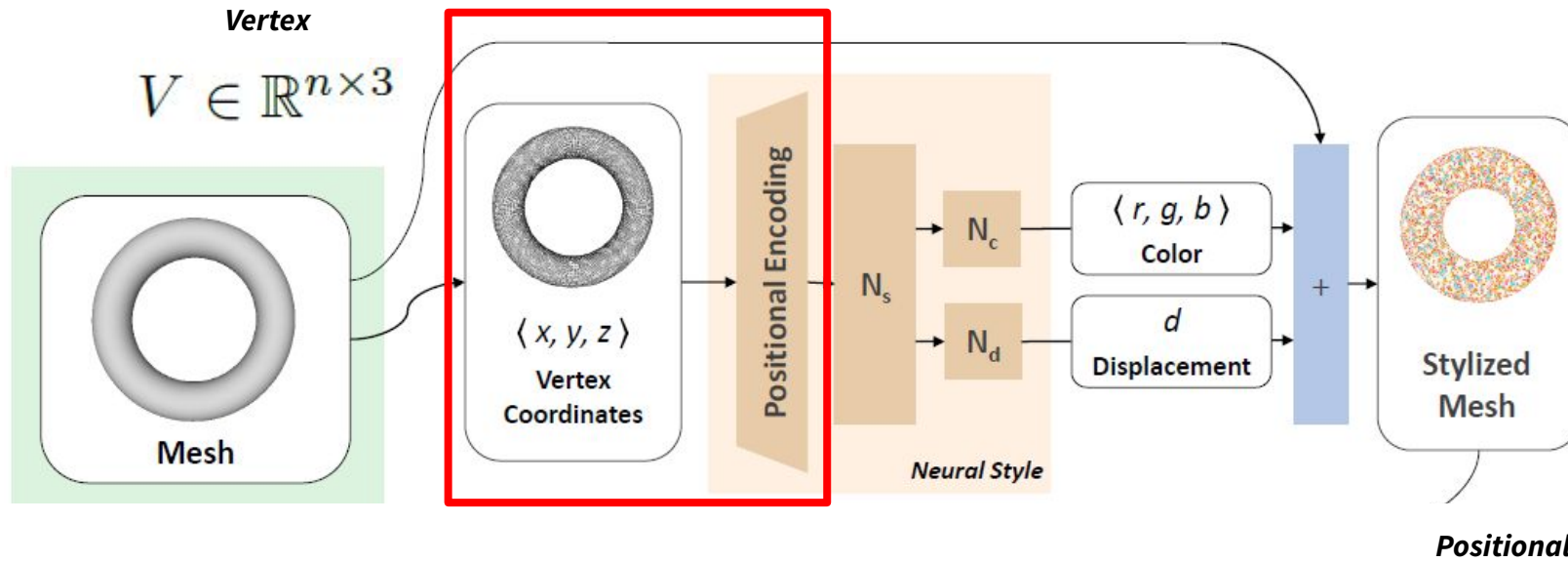
Figure 4.



- The **CLIP similarity** between the rendered and augmented **images** and the **target text** is used as a signal to **update** the neural network **weights**.

# Architecture

Figure 4.



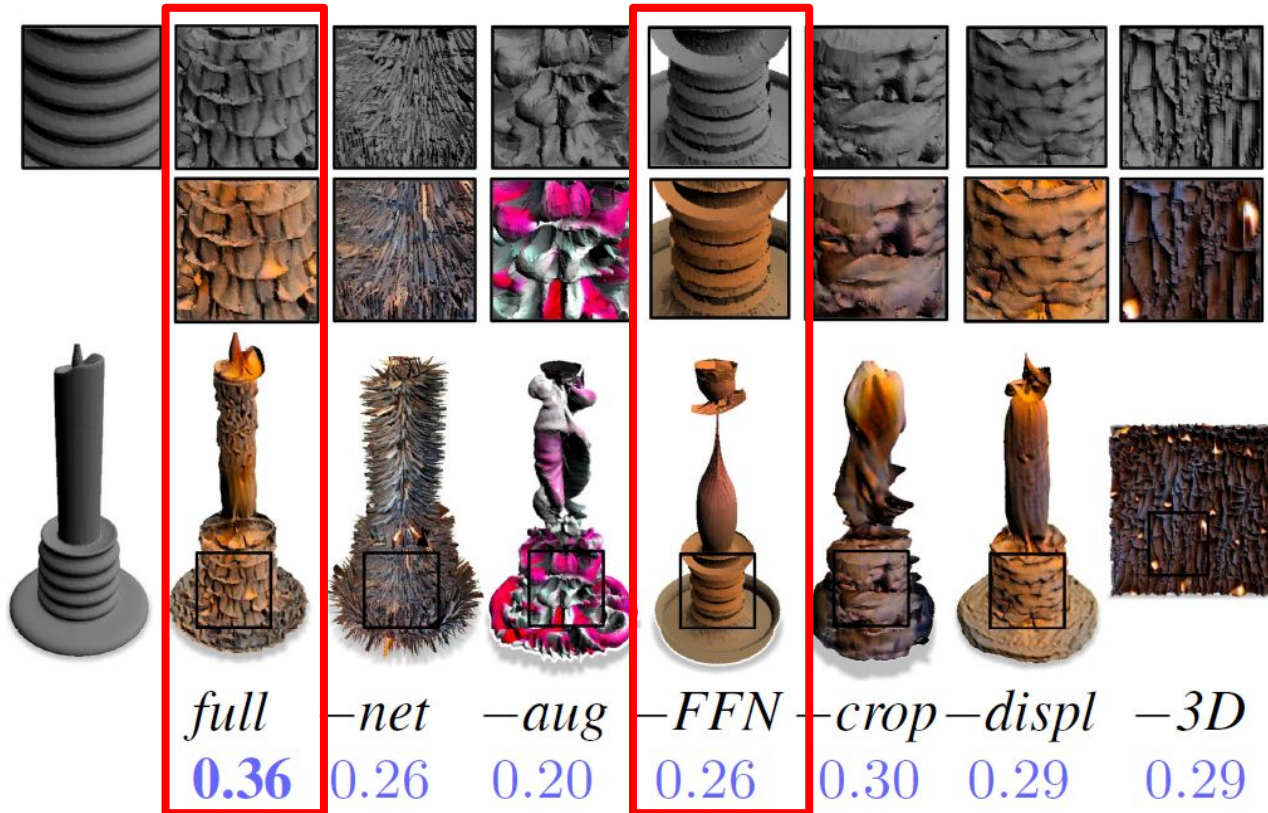
$$\gamma(p) = [\cos(2\pi \mathbf{B}p), \sin(2\pi \mathbf{B}p)]^T$$

- Normalize the coordinates and map a vertex to a 256-dimensional Fourier feature
- Per-vertex **positional encoding** features  $\gamma(p)$  are passed as input to an MLP  $N_s$ , which then branches out to MLPs  $N_d$  and  $N_c$

# Positional Encoding

Figure 5.

Ablation on the priors used in our method (full) for a candle mesh and target 'Candle made of bark'



*full*: our method

*-net*: without our style field network

*-aug*: without 2D augmentations

**-FFN: without positional encoding**

*-crop*: without crop augmentations for  $\psi_{local}$

*-displ*: without the geometry-only component of  $L_{sim}$

*-3D*: learning over a 2D plane in 3D space

CLIP score:  $\text{sim}(\hat{S}^{\text{full}}, \phi_{\text{target}})$

Semantic loss:  $\mathcal{L}_{\text{sim}} = \sum_{\hat{S}} \text{sim}(\hat{S}, \phi_{\text{target}})$

Utilize the **positional encodings** using fast **fourier feature networks** what enables us to obtain the fine grained results to solve the **spectral bias** problem

# Positional Encoding



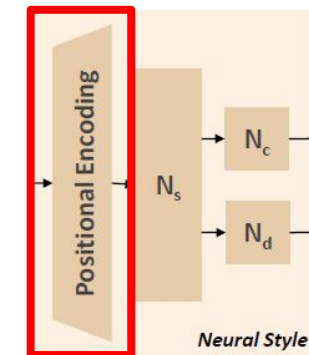
*full*

VS



*-FFN*

**Figure 4.**



Utilize the ***positional encodings*** using fast ***fourier feature networks*** what enables us to obtain the fine grained results to solve the ***spectral bias*** problem



# Positional Encoding

$\sigma$ : The amount of frequencies that are going into the positional encoding

‘Stained glass donut’

For every point  $p$  its positional encoding  $\gamma(p)$  is given by:

$$\gamma(p) = [\cos(2\pi \mathbf{B}p), \sin(2\pi \mathbf{B}p)]^T$$

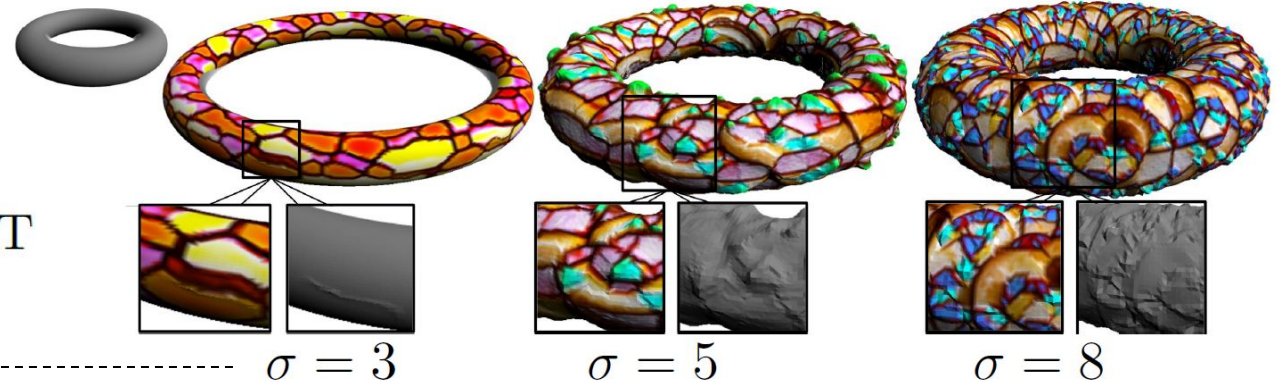
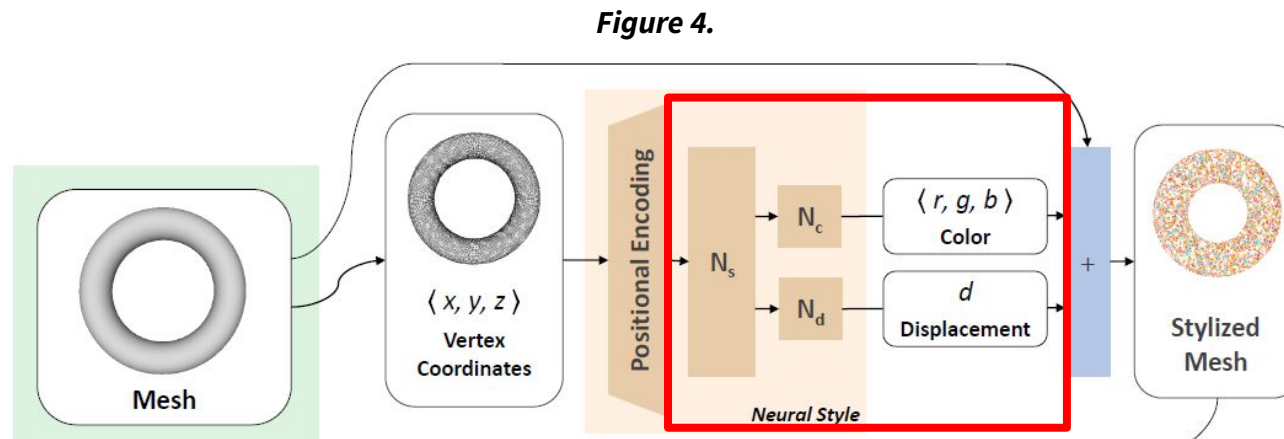


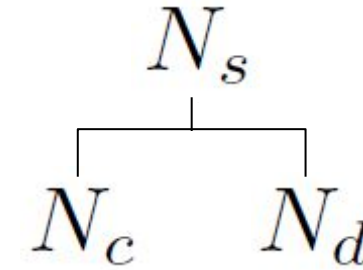
Figure 7. Increasing the range of input frequencies in the positional encoding using increasing SD  $\sigma$  for matrix  $\mathbf{B}$  in Eq. (1).

- The network leverages a **positional encoding** where the range of **frequencies** can be **directly controlled** by the standard deviation  $\sigma$  of the  $\mathbf{B}$  matrix
- Increasing the frequency value increases the frequency of style details on the mesh and produces **sharper** and **more frequent displacements** along the normal direction

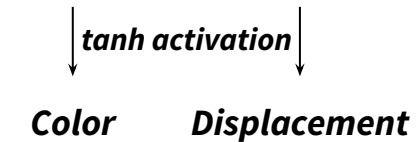
# Architecture



*Four 256-dimensional linear layers  
& ReLU activation*



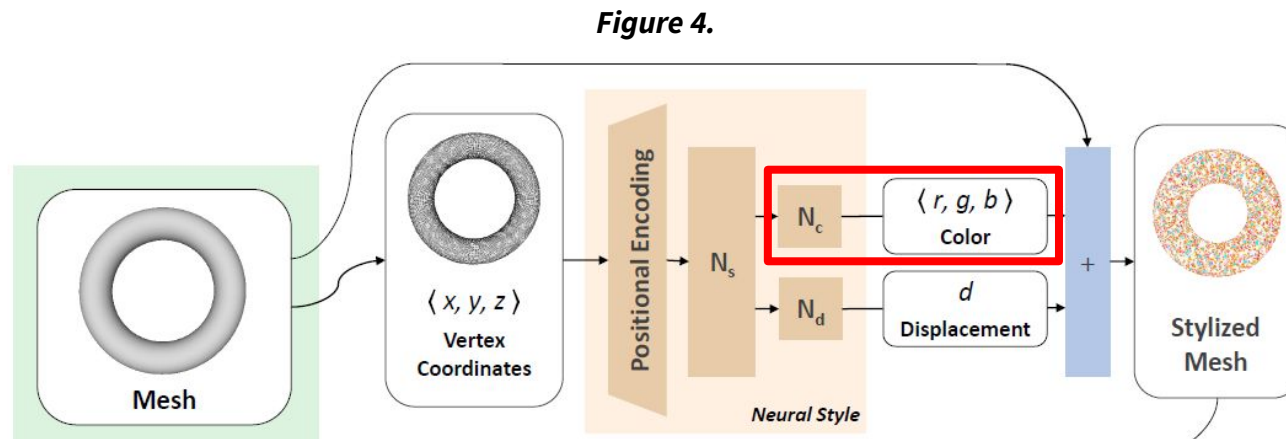
*Each of them have two 256-dimensional linear layers  
& ReLU activation*



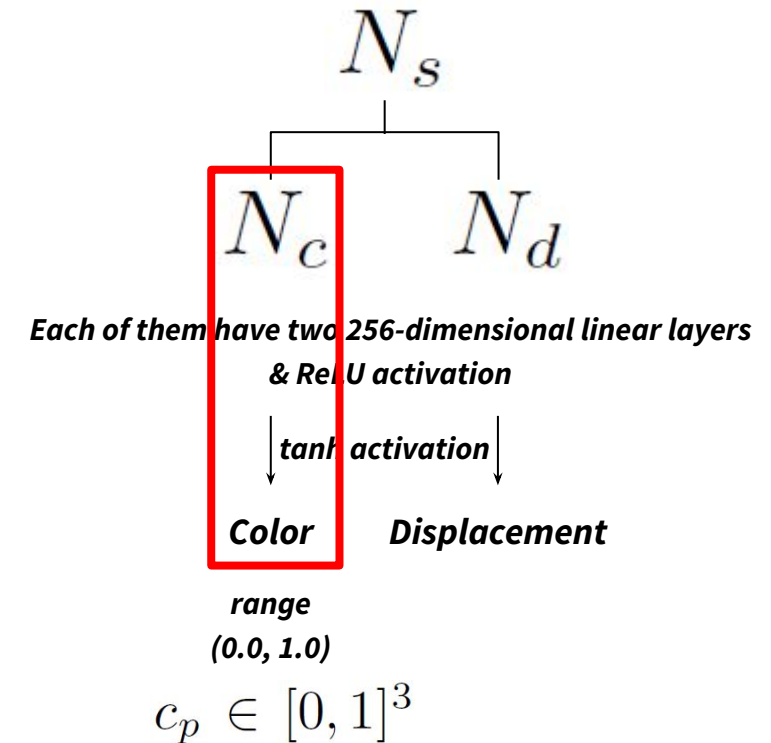
- The shared MLP layers  $N_s$  consist of four 256-dimensional linear layers with ReLU activation.
- The branched layers,  $N_d$  and  $N_c$ , each consist of two 256-dimensional linear layers with ReLU activation.
- After the final linear layer, a tanh activation is applied to each branch.



# Architecture

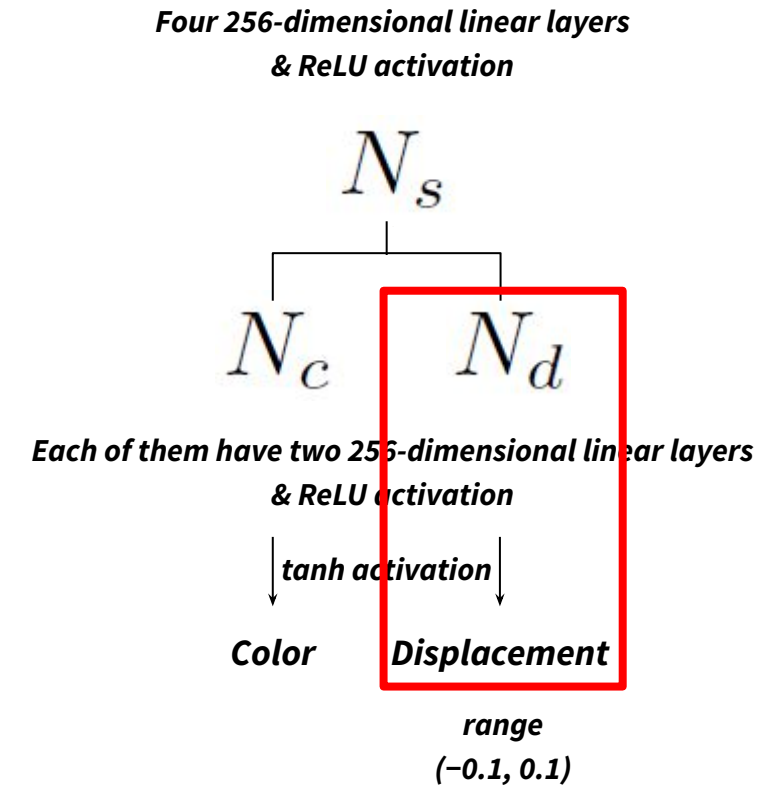
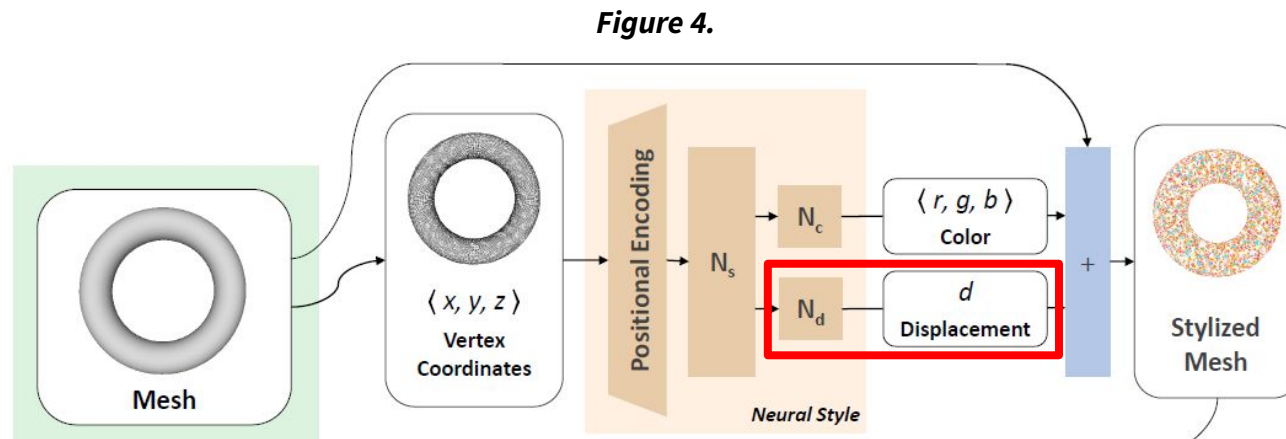


*Four 256-dimensional linear layers  
& ReLU activation*



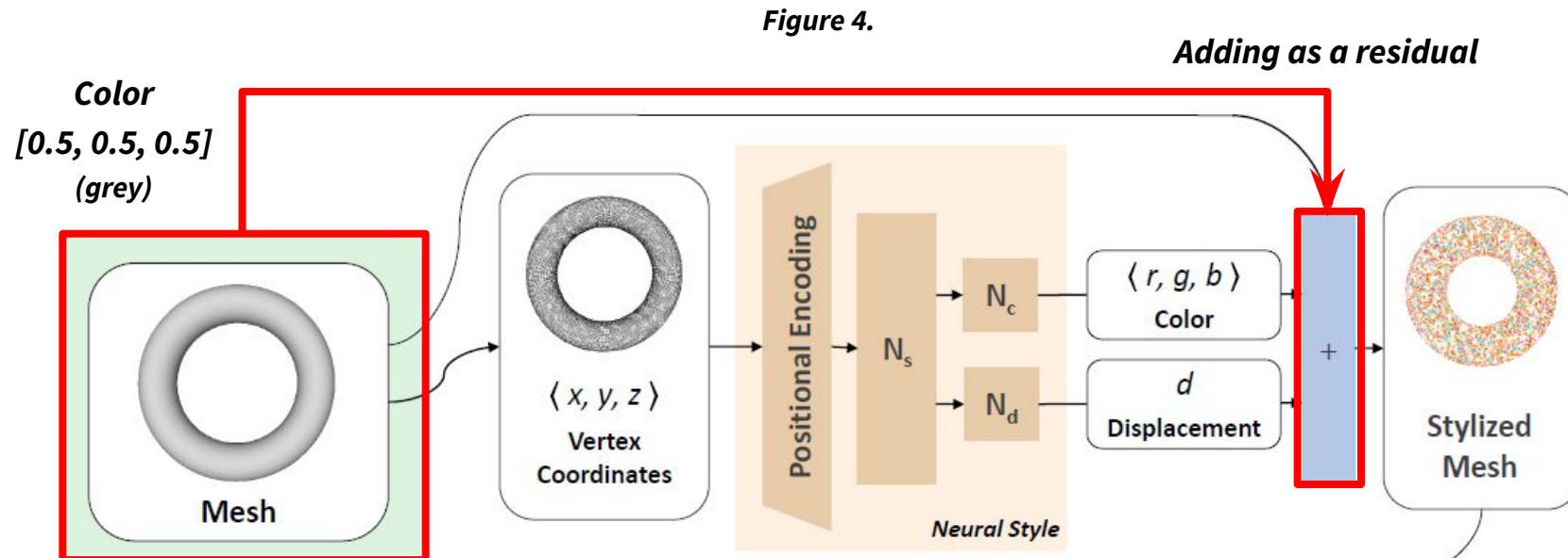
- Divide the output of  $N_c$  by 2 and add it to  $[0.5, 0.5, 0.5]$
- This enforces the final **color prediction**  $c_p$  to be in range  $(0.0, 1.0)$

# Architecture



- For the branch  **$N_d$** , multiply the final tanh layer by 0.1 to get displacements in the range  **$(-0.1, 0.1)$**
- Constrain  $d_p$  to be in the range  $(-0.1, 0.1)$  to **prevent content-altering displacements**

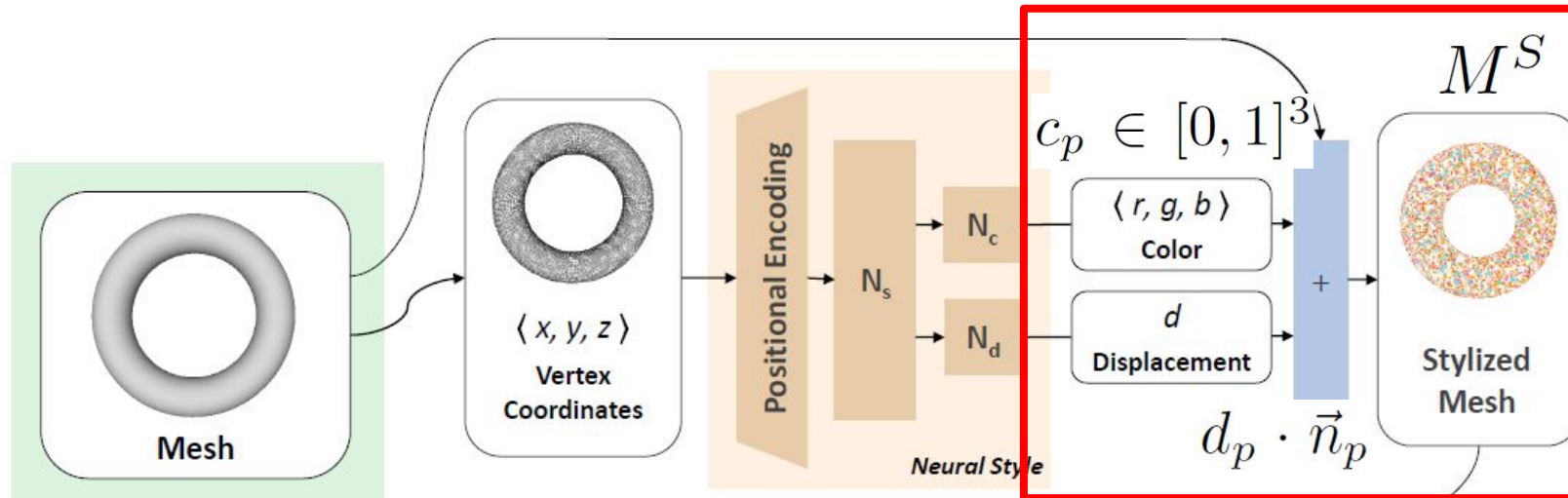
# Architecture



Author find that initializing the *mesh color* to  $[0.5, 0.5, 0.5]$  (grey) and *adding the network output as a residual* helps *prevent undesirable solutions* in the early iterations of training.

# Architecture

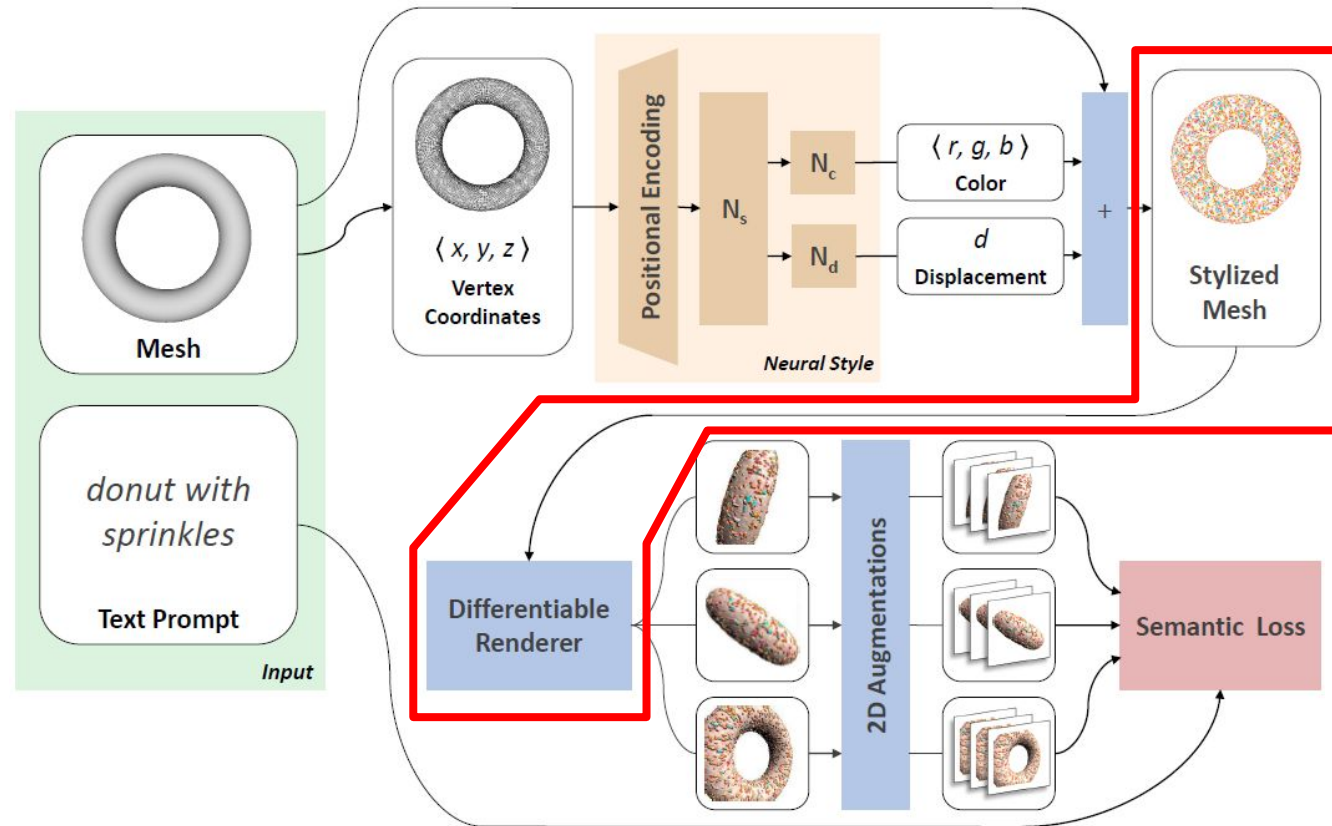
Figure 4.



- Every point  $p$  is displaced by  $d_p \cdot \vec{n}_p$  and colored by  $c_p$  to obtain stylized mesh prediction  $M^S$
- **Vertex colors propagate over the entire mesh surface** using an interpolation-based differentiable renderer

# Architecture

Figure 4.



**Sample  $n\theta$  ( $= 5$ ) views** around a predefined **anchor view** and render them using a **differentiable renderer**, for given the stylized mesh  $M^S$  and the displaced mesh  $M_{\text{displ}}^S$

# Anchor View Choice

Views with Highest and lowest CLIP Score

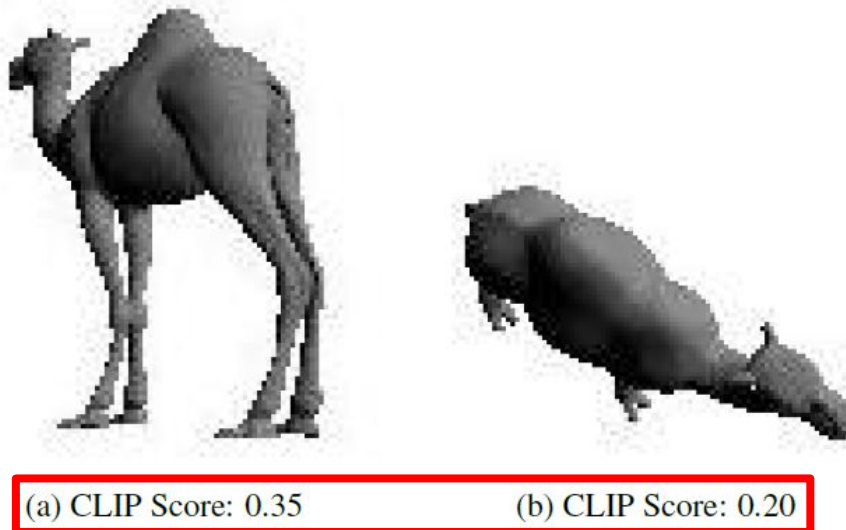


Figure 20. Example views with CLIP similarities.

CLIP score of the view  
that passes from the vertex to the center of the mesh

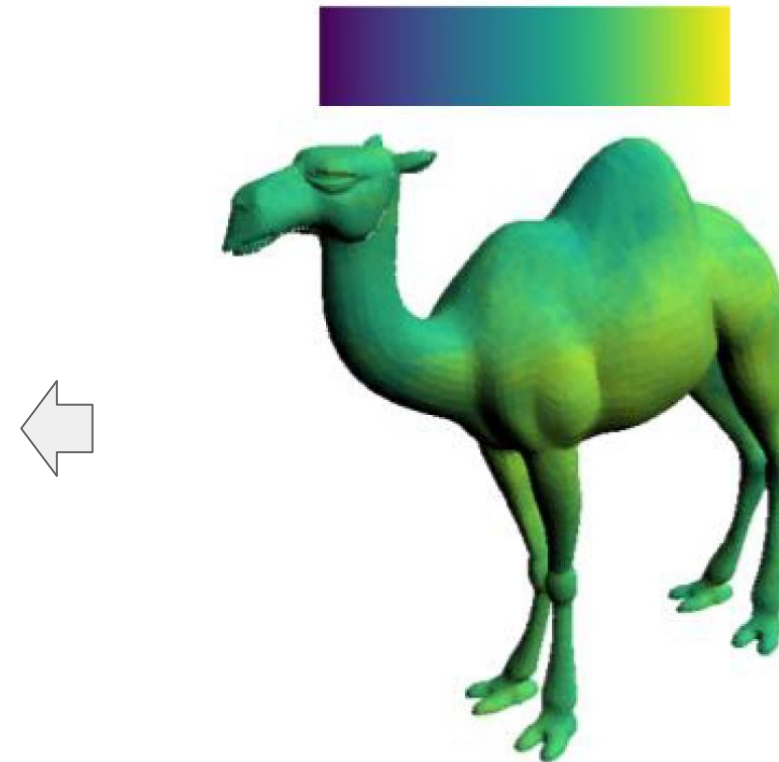


Figure 19. CLIP scores for each vertex view.

- Select the view with the **highest** (i.e. best) **CLIP similarity** to the content **as the anchor** which will allow a high-quality stylization
- This **metric is limited in expressiveness**, however, as demonstrated by the **constrained range** that the scores fall within for all the views around the mesh.



# Anchor View Choice

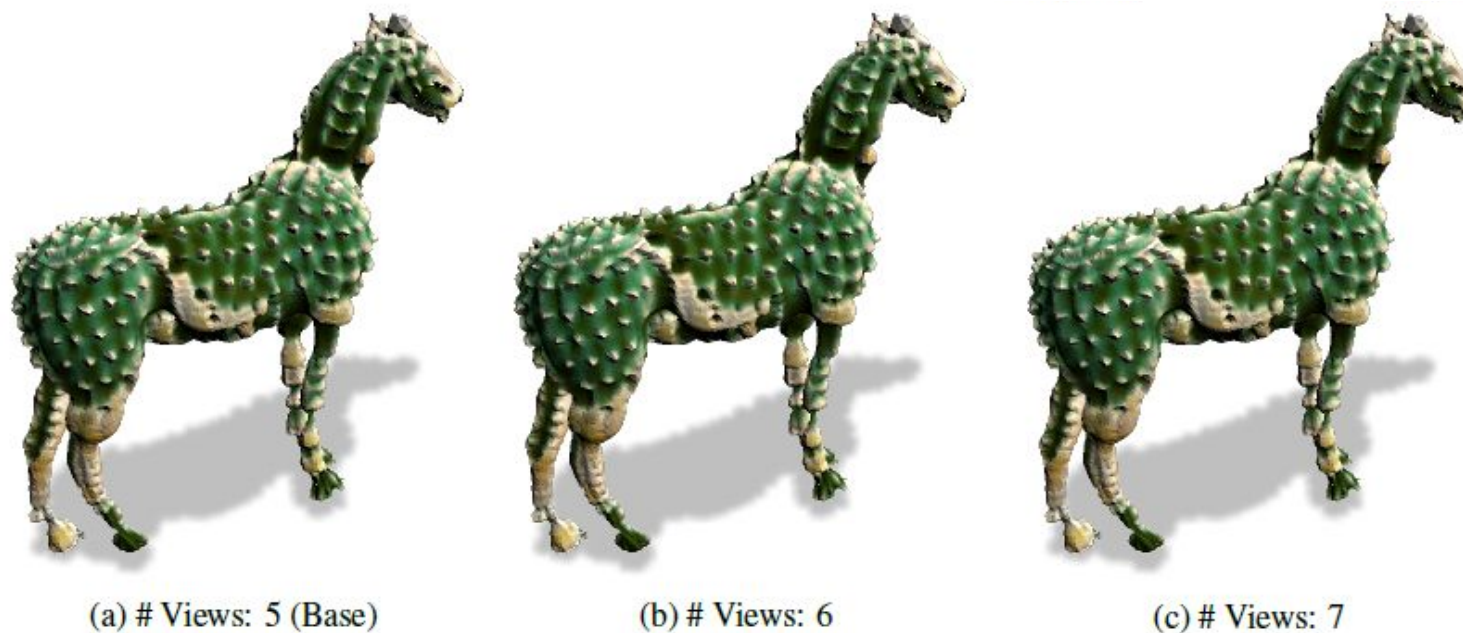


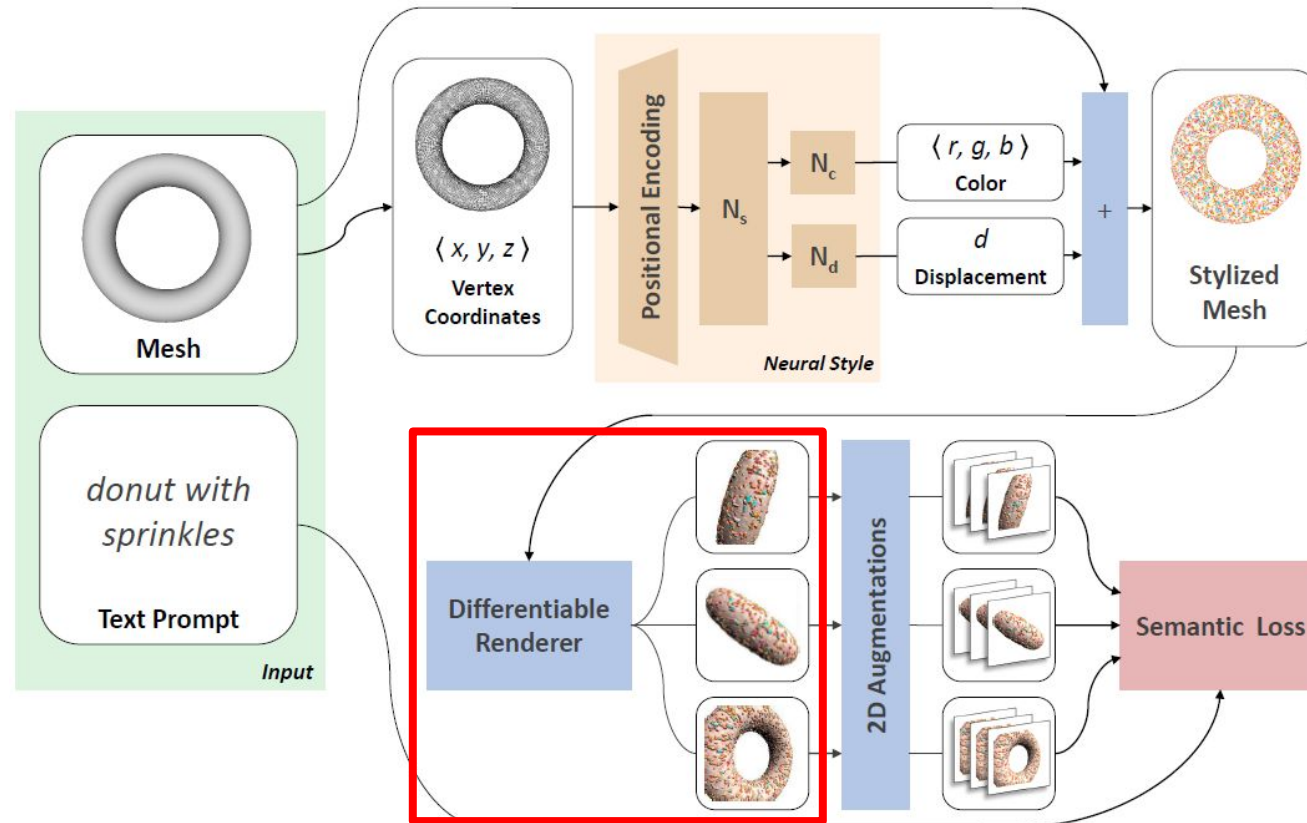
Figure 21. Style outputs sampling different # views. Prompt: 'A horse made of cactus'

***Increasing*** the number of views ***beyond 5*** does ***little to change the quality*** of the output stylization



# Architecture

Figure 4.

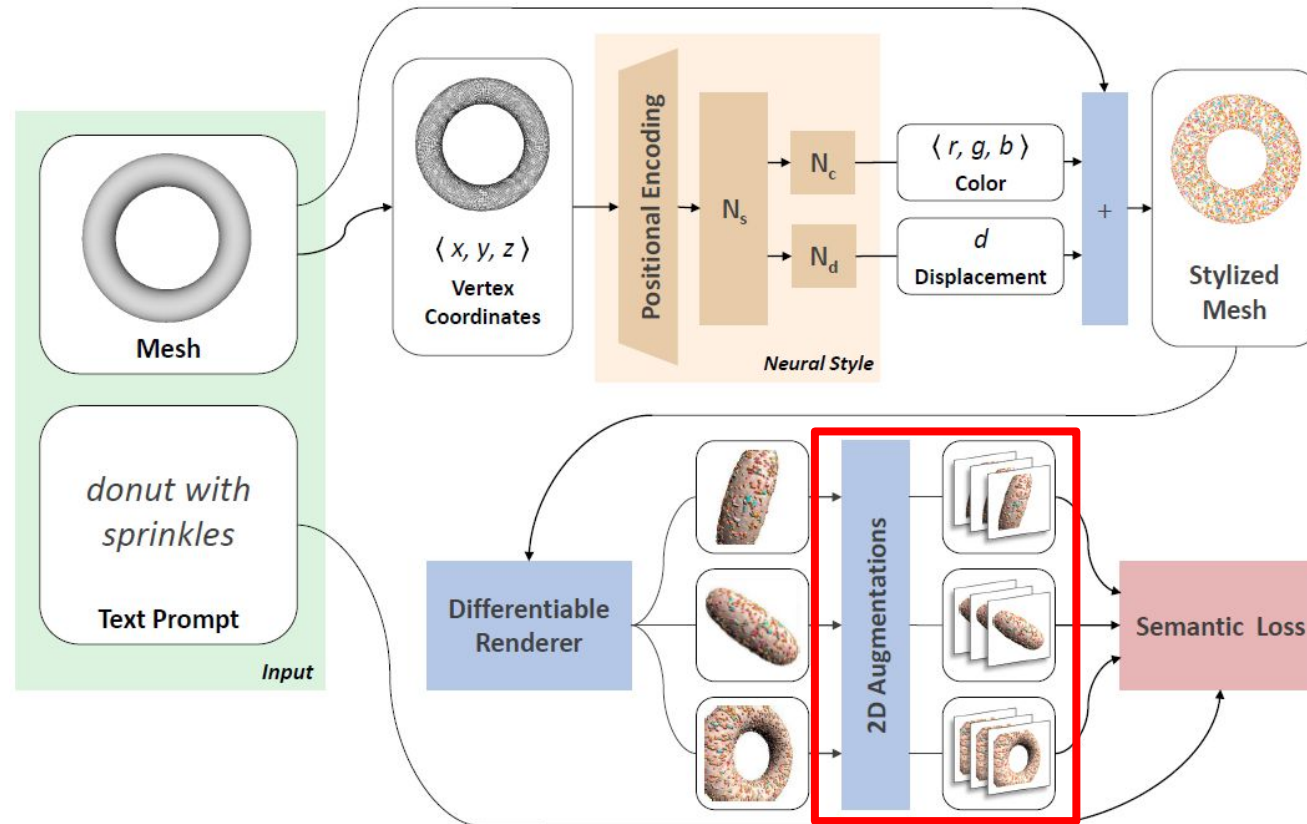


For each view,  $\theta$ , render two 2D projections of the surface  $I_{\theta}^{\text{full}}$  for  $M^S$  and  $I_{\theta}^{\text{displ}}$  for  $M_{\text{displ}}^S$

$M^S$  : Stylized mesh  
 $M_{\text{displ}}^S$  : Displaced mesh

# Architecture

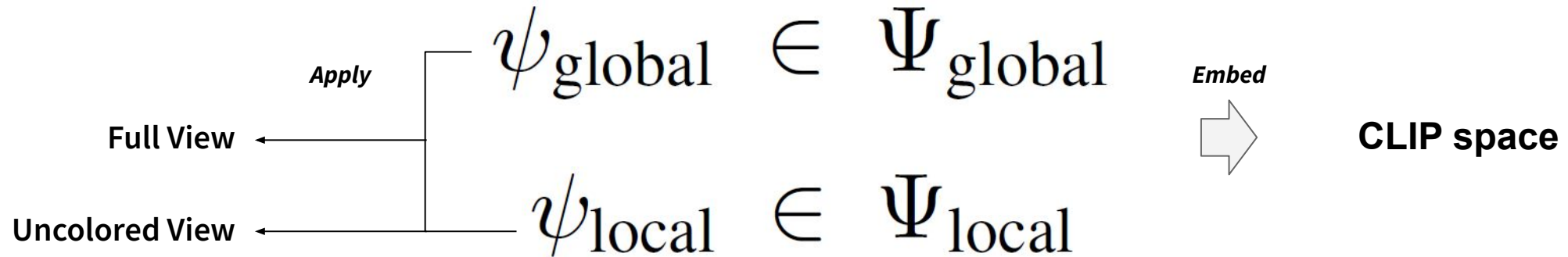
Figure 4.



- Draw a 2D augmentation  $\psi_{global} \in \Psi_{global}$  and  $\psi_{local} \in \Psi_{local}$
- Apply  $\psi_{global}$ ,  $\psi_{local}$  to the full view and  $\psi_{local}$  to the uncolored view, and embed them into CLIP space

# 2D augmentations ( $\psi_{global}$ , $\psi_{local}$ and cropping)

$\psi_{global}$ : Involves a **random perspective** transformation



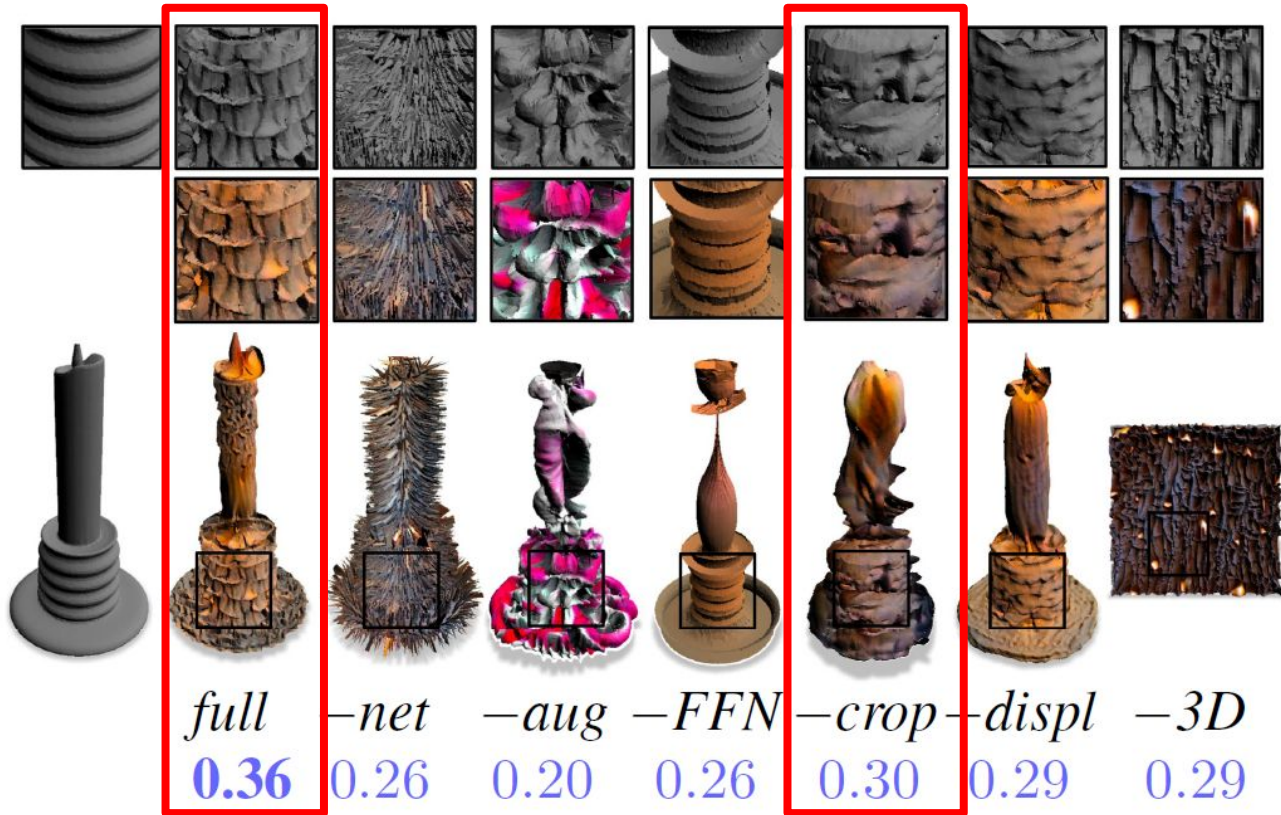
$\psi_{local}$ : Generates both a **random perspective** and a **random crop** that is 10% of the original image

The 2D augmentations generated using  $\psi_{global}$  and  $\psi_{local}$  are critical for method to **avoid degenerate solutions**

# 2D augmentations ( $\psi_{global}$ , $\psi_{local}$ and cropping)

Figure 5.

Ablation on the priors used in our method (full) for a candle mesh and target ‘Candle made of bark’



full: our method

-net: without our style field network

-aug: without 2D augmentations

-FFN: without positional encoding

**-crop: without crop augmentations for  $\psi_{local}$**

-displ: without the geometry-only component of  $L_{sim}$

-3D: learning over a 2D plane in 3D space

CLIP score:  $\text{sim}(\hat{S}^{\text{full}}, \phi_{\text{target}})$

Semantic loss:  $\mathcal{L}_{\text{sim}} = \sum_{\hat{S}} \text{sim}(\hat{S}, \phi_{\text{target}})$

**Cropping** allows the network to focus on localized regions when making *fine grained adjustments* to the surface geometry and color (Check the -crop)

## 2D augmentations ( $\psi_{global}$ , $\psi_{local}$ and cropping)

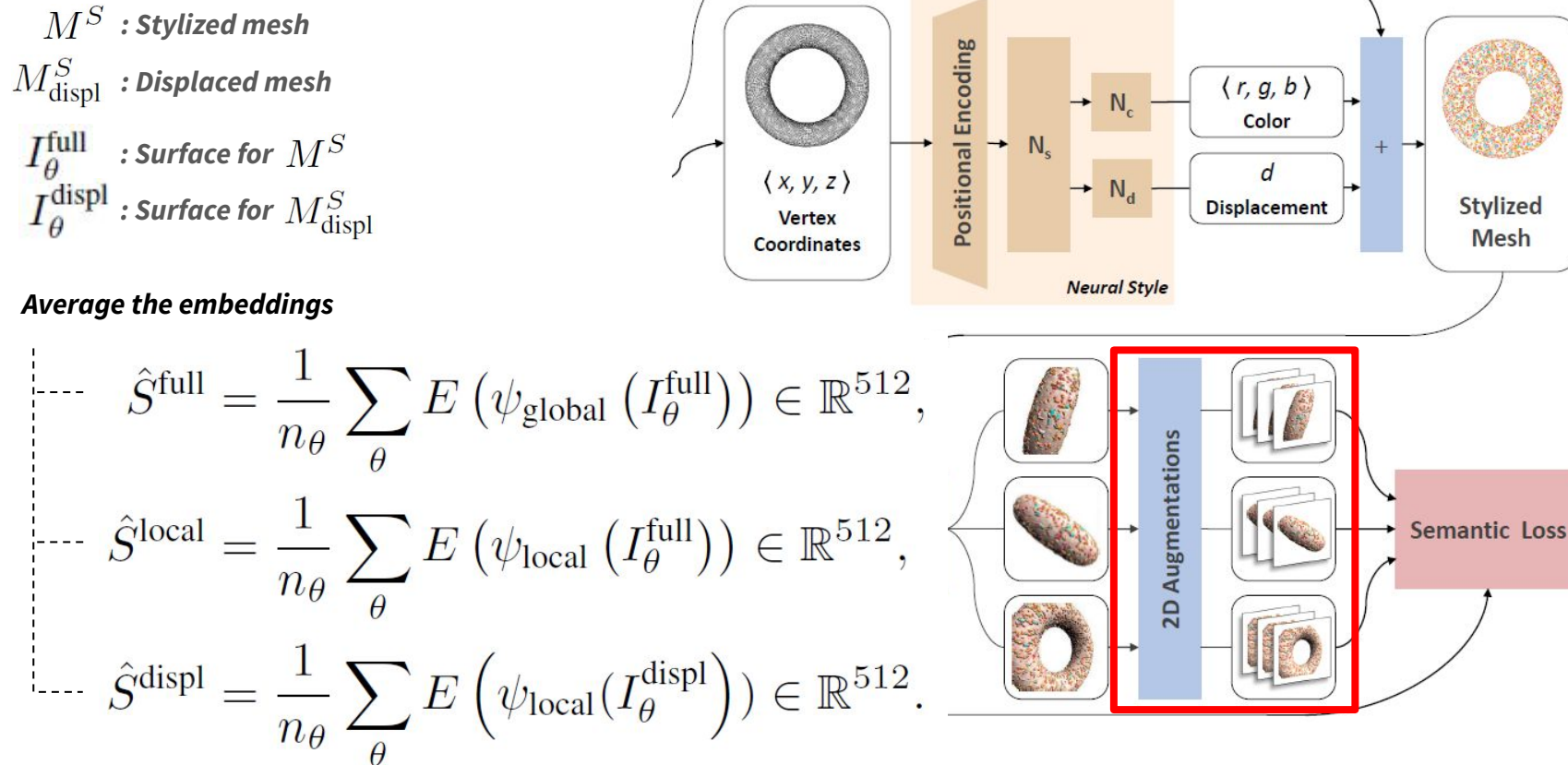


**Cropping** allows the network to focus on localized regions when making *fine grained adjustments* to the surface geometry and color (Check the *-crop*)



# Architecture

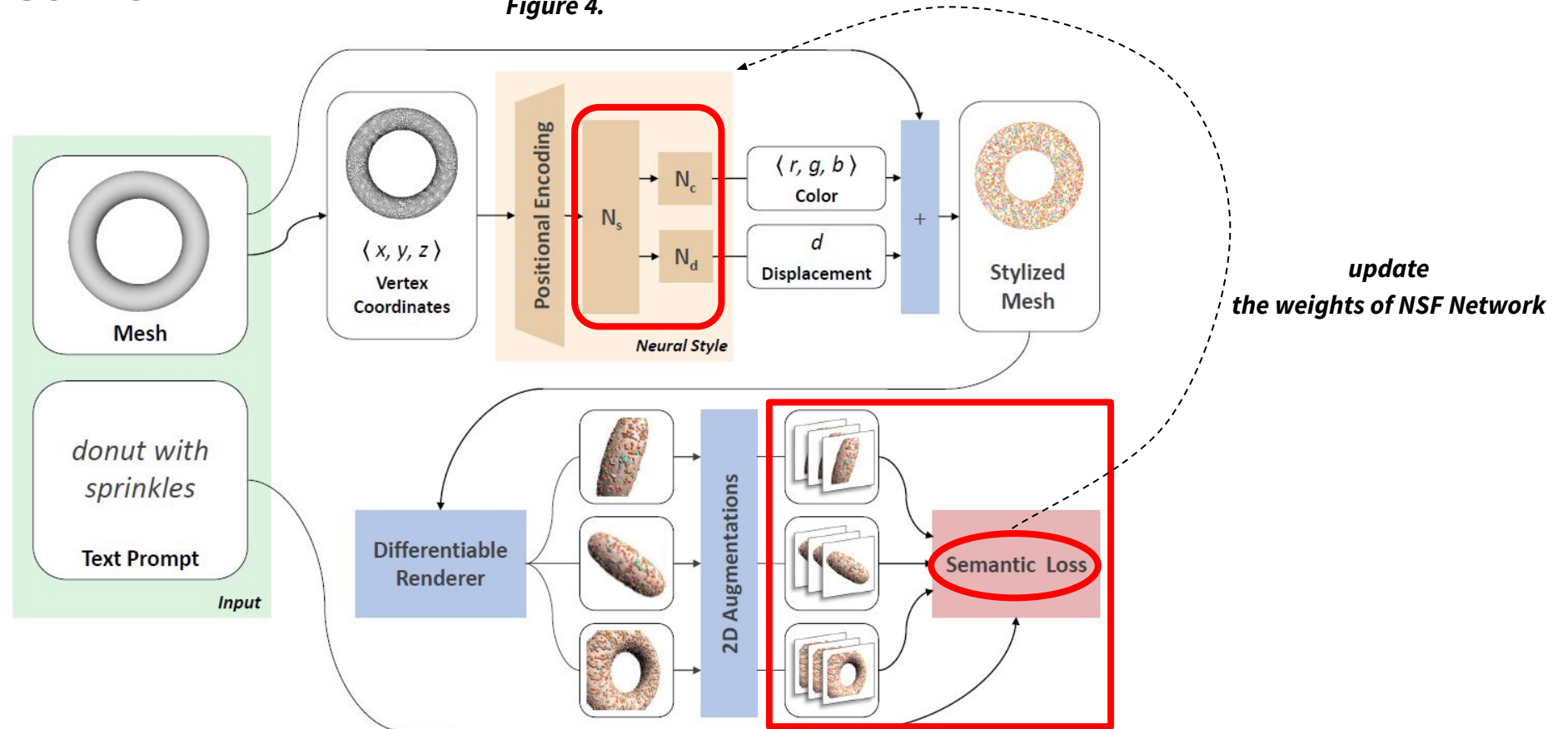
Figure 4.



- Average the embeddings across all views
- Consider an augmented representation of our input mesh as the average of its encoding from multiple augmented views

# Architecture

Figure 4.



- The target  $t$  is embedded through CLIP by  $\phi_{\text{target}} = E(t) \in \mathbb{R}^{512}$

- Semantic Loss:

$$\mathcal{L}_{\text{sim}} = \sum_{\hat{S}} \text{sim}(\hat{S}, \phi_{\text{target}}) \quad \left\{ \begin{array}{l} \hat{S} \in \{\hat{S}^{\text{full}}, \hat{S}^{\text{displ}}, \hat{S}^{\text{local}}\} \\ \text{sim}(a, b) = \frac{a \cdot b}{|a| \cdot |b|} \end{array} \right.$$

$\text{sim}(a, b)$  is the cosine similarity between  $a$  and  $b$



# Updating the Weights of NSF Network

$M^S$  : Stylized mesh       $I_{\theta}^{\text{full}}$  : Surface for  $M^S$   
 $M_{\text{displ}}^S$  : Displaced mesh       $I_{\theta}^{\text{displ}}$  : Surface for  $M_{\text{displ}}^S$

**Average the embeddings**

$$\hat{S}^{\text{full}} = \frac{1}{n_{\theta}} \sum_{\theta} E(\psi_{\text{global}}(I_{\theta}^{\text{full}})) \in \mathbb{R}^{512},$$

$\hat{S}^{\text{full}}$

**Update**



**All**  
 $N_s \quad N_c \quad N_d$

$$\hat{S}^{\text{local}} = \frac{1}{n_{\theta}} \sum_{\theta} E(\psi_{\text{local}}(I_{\theta}^{\text{full}})) \in \mathbb{R}^{512},$$

$\hat{S}^{\text{local}}$

$$\hat{S}^{\text{displ}} = \frac{1}{n_{\theta}} \sum_{\theta} E(\psi_{\text{local}}(I_{\theta}^{\text{displ}})) \in \mathbb{R}^{512}.$$



**Semantic loss**

$$\mathcal{L}_{\text{sim}} = \sum_{\hat{S}} \text{sim}(\hat{S}, \phi_{\text{target}})$$

$\hat{S}^{\text{displ}}$

**Update**



$N_s \quad N_d$

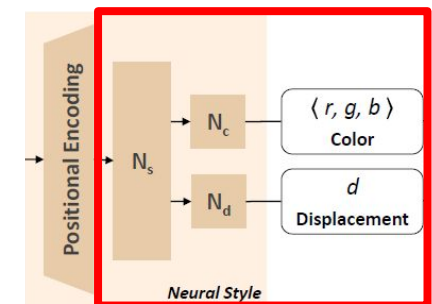
**Unrelated to Color**

$N_c \rightarrow \text{Color } c_p \in [0, 1]^3$

$N_d \rightarrow \text{Displacement along the vertex normal } d_p$

$\hat{S}^{\text{full}}$  and  $\hat{S}^{\text{local}}$  update  $N_s, N_c$  and  $N_d$  while  $\hat{S}^{\text{displ}}$  only updates  $N_s$  and  $N_d$

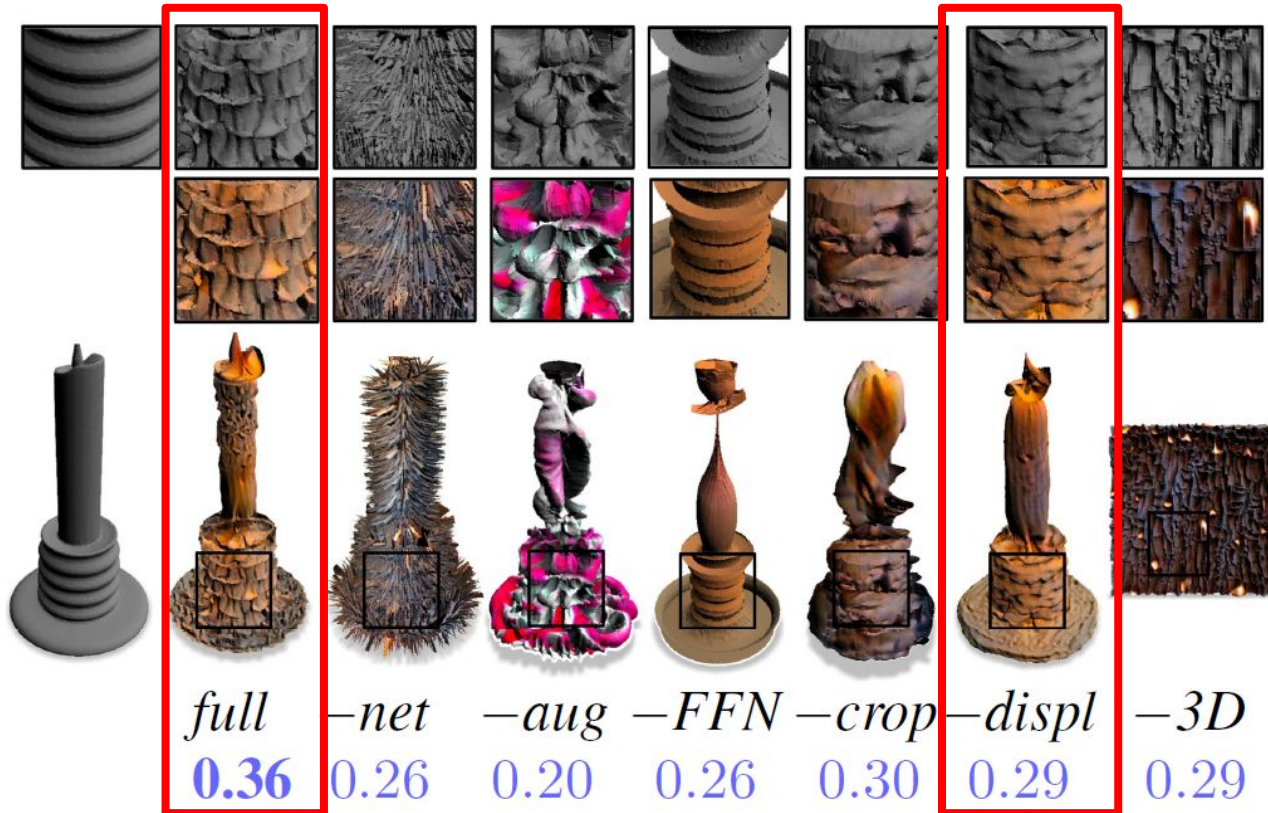
**Figure 4.**



# Separation (geo-only loss | geo-and-color loss)

Figure 5.

Ablation on the priors used in our method (full) for a candle mesh and target 'Candle made of bark'



full: our method

-net: without our style field network

-aug: without 2D augmentations

-FFN: without positional encoding

-crop: without crop augmentations for  $\psi_{local}$

**-displ: without the geometry-only component of  $L_{sim}$**

-3D: learning over a 2D plane in 3D space

CLIP score:  $\text{sim}(\hat{S}^{\text{full}}, \phi_{\text{target}})$

Semantic loss:  $\mathcal{L}_{\text{sim}} = \sum_{\hat{S}} \text{sim}(\hat{S}, \phi_{\text{target}})$

The **separation** into a **geometry-only loss** and **geometry-and-color loss** is an effective tool for **encouraging meaningful changes in geometry** (check the -displ)

# Separation (geo-only loss | geo-and-color loss)



*full*

VS



*-displ*

The ***separation*** into a ***geometry-only loss*** and ***geometry-and-color loss*** is an effective tool for ***encouraging meaningful changes in geometry*** (check the *-displ*)

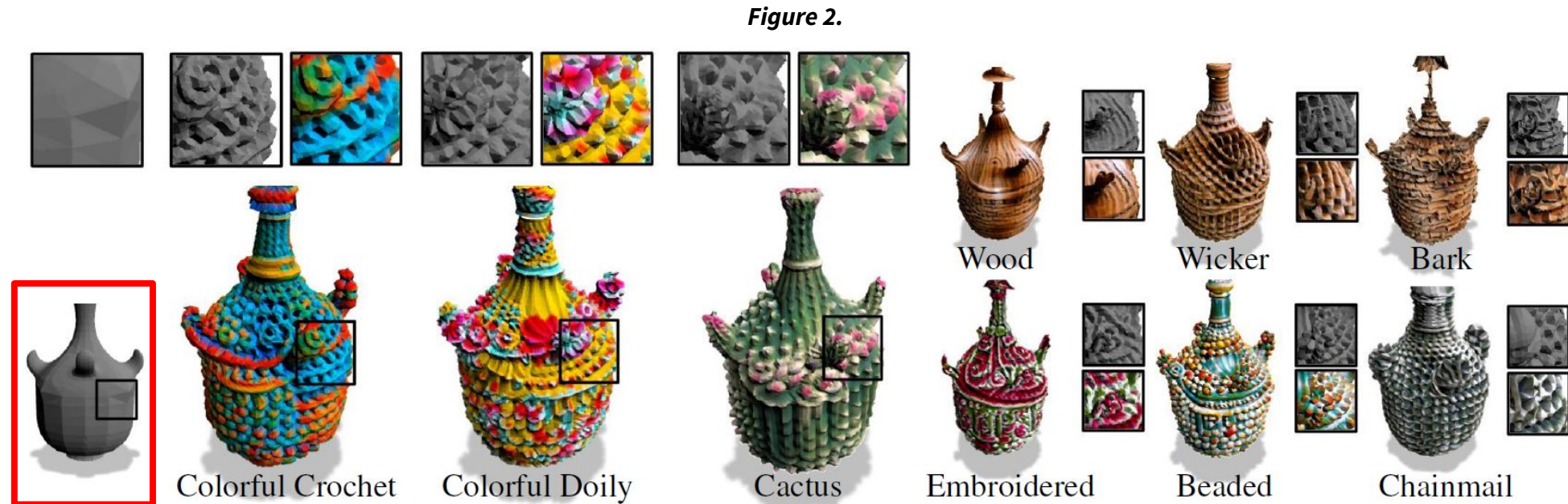
# Experiments and Results

## Input source meshes

**Sources:** COSEG, Thingi10K, Shapenet, Turbo Squid, and ModelNet

**Features:** Average of 79,366 faces, 16% non-manifold edges, 0.2% non-manifold vertices, and 12% boundaries

# Experiments and Results

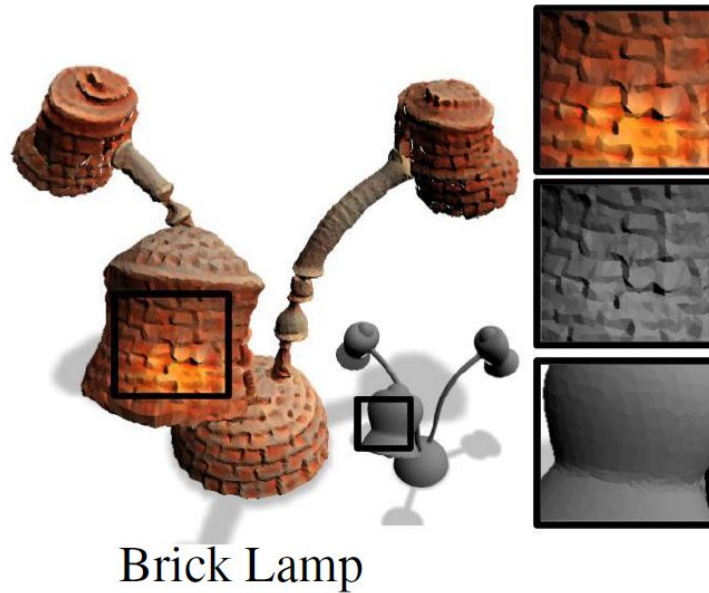


Maintaining global semantics and preserving the underlying content



# Experiments and Results

*Figure 1.*



*Figure 2.*



Generates structured textures which are aligned to sharp curves and features

# Experiments and Results



Figure 6. Our neural texture field stylizes the entire 3D shape.

Styles the entire mesh in a consistent manner that is part-aware and exhibits natural variation in texture



# Experiments and Results

$\sigma$ : The amount of frequencies that are going into the positional encoding

‘Stained glass donut’

For every point  $p$  its positional encoding  $\gamma(p)$  is given by:

$$\gamma(p) = [\cos(2\pi \mathbf{B}p), \sin(2\pi \mathbf{B}p)]^T$$

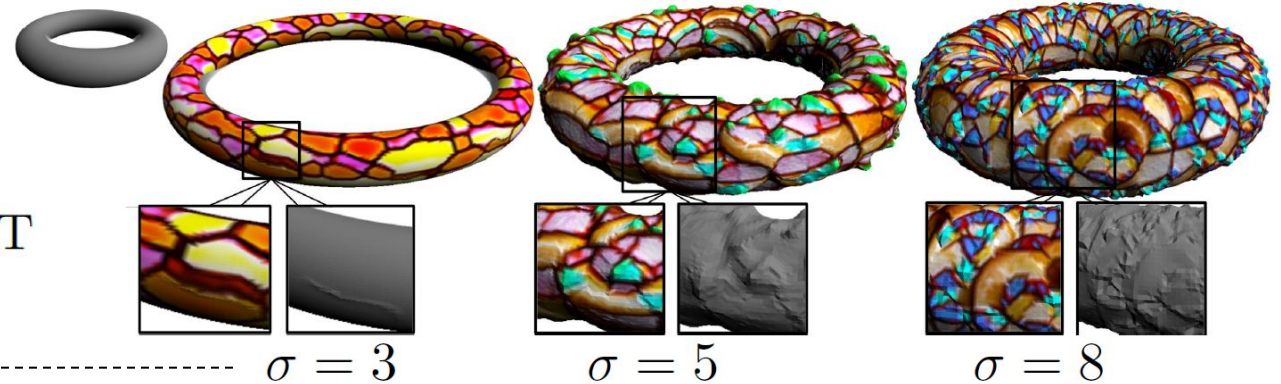
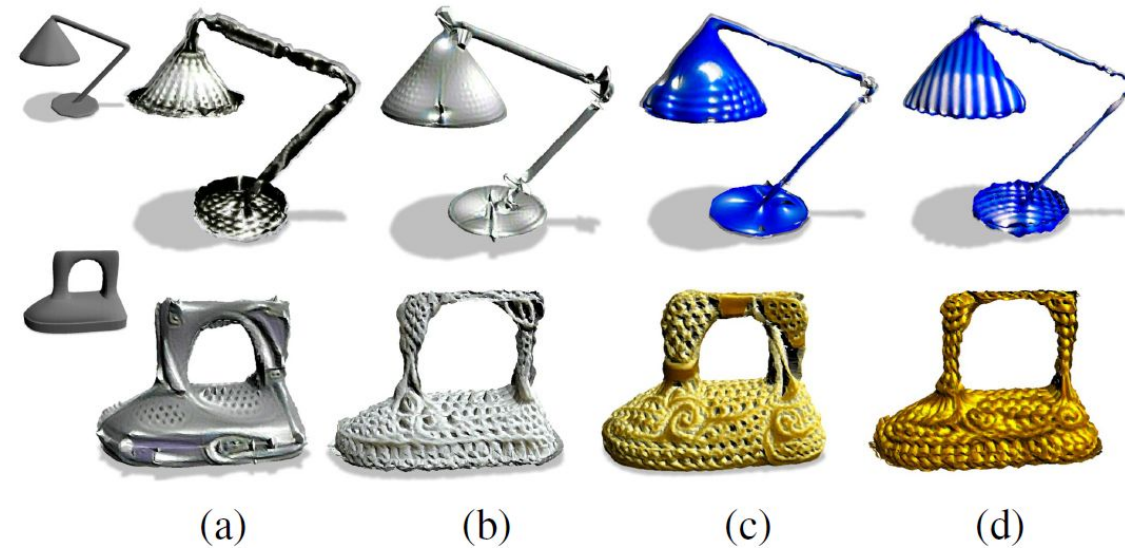


Figure 7. Increasing the range of input frequencies in the positional encoding using increasing SD  $\sigma$  for matrix  $\mathbf{B}$  in Eq. (1).

Increasing the  $\sigma$  (frequency value) increases the frequency of style details on the mesh and produces **sharper** and **more frequent displacements** along the normal direction

# Experiments and Results

**Figure 8.**  
*Increasing the target text prompt granularity for a source mesh of a lamp and iron.*



*Increasing style details*

(a). 'Lamp', (b). 'Luxo lamp', (c). 'Blue steel luxo lamp', (d). 'Blue steel luxo lamp with corrugated metal.'

(a). 'Clothes iron', (b). 'Clothes iron made of crochet', (c). 'Golden clothes iron made of crochet', (d). 'Shiny golden clothes iron made of crochet.'

- Successfully synthesize styles of varying levels of specificity.
- Retention of the style details from each level of target granularity to the next.

# Experiments and Results



Figure 22. Prompt: ‘A shoe made of cactus’

Figure 23. Prompt: ‘A chair made of brick’

## 57 users evaluate 8 samples

- (Q1) “How natural is the output depiction of {content} + {style}?”
- (Q2) “How well does the output match the original {content}?”
- (Q3) “How well does the output match the target {style}?”

Outperforms the VQGAN baseline across all questions, with a difference of 1.07, 0.44, and 1.32 for Q1-Q3

### VQGAN-CLIP:

Synthesizes color inside a binary 2D mask projected from the 3D source shape (without 3D deformations) guided by CLIP

	(Q1): Overall	(Q2): Content	(Q3): Style
VQGAN	2.83 ( $\pm 0.39$ )	3.60 ( $\pm 0.68$ )	2.59 ( $\pm 0.44$ )
Ours	<b>3.90</b> ( $\pm 0.37$ )	<b>4.04</b> ( $\pm 0.53$ )	<b>3.91</b> ( $\pm 0.51$ )

Table 1. Mean opinion scores (1-5) for Q1-Q3 (see Sec. 4.3), for our method and baseline (control score: 1.16).

# Limitation

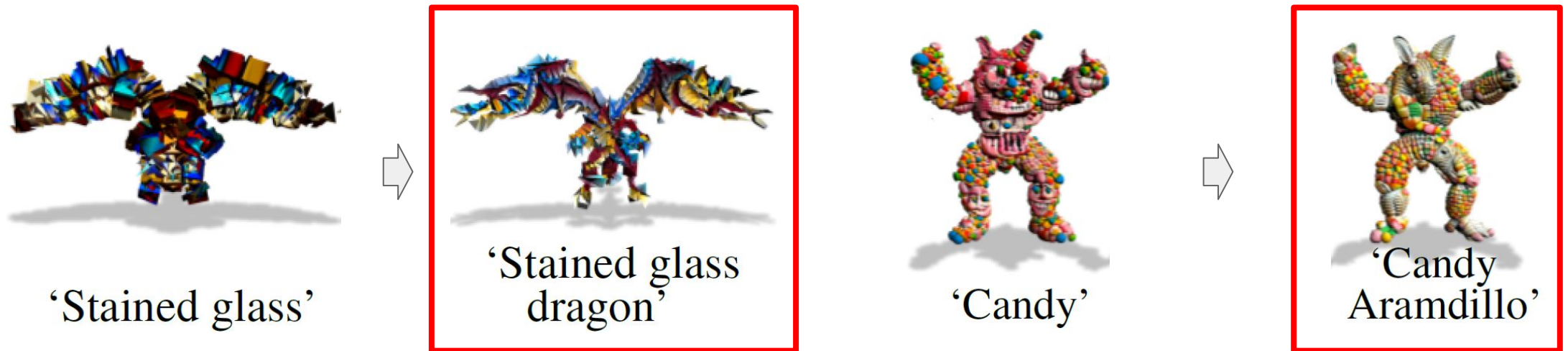
Figure 15.  
Geometric content and target style synergy.



- **Assumes** there *exists* a **synergy** between the *input 3D geometry* and the *target style prompt*.
- However, stylizing a 3D mesh (e.g., dragon) towards an *unrelated/unnatural prompt* (e.g., stained glass) may result in a stylization that *ignores the geometric prior* and *erases the source shape content*.

# Limitation

Figure 15.  
Geometric content and target style synergy.



The author *solve* this by simply *including the object category in the text prompt* (e.g., stained glass dragon) which adds a content preservation constraint into the target.



# Limitation



Figure 25. Our method enables visualizing the biases in the CLIP embedding space. Given a human male input (source in Figure 3), and target prompt: ‘a nurse’, we observe a gender bias in CLIP to favor female shapes.

- ***Societal bias***: The nurse style in Fig. 25 is biased towards adding female features to the input male shape.
- Present in joint image-text embeddings through our stylization framework



# Conclusion

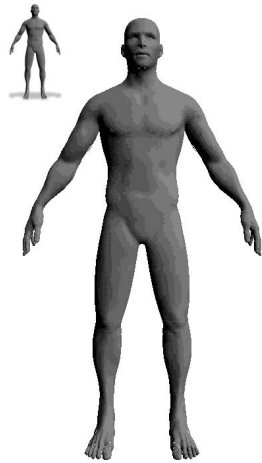
## Mesh stylization through 2D projections

### *Key Points*

1. Intuitive control over 3D shape manipulation
2. Without a directional field or mesh parameterization
3. Without relying on a pre-trained GAN network or a 3D dataset



a vase made of colorful crochet



**Thank you.**