# Learning Geometry-aware Representations by *Sketching*

## THU-PM-259

**Hyundo Lee**

**Inwoo Hwang**

**Hyunsung Go**

**Won-Seok Choi**

**Kibeom Kim**

**Byoung-Tak Zhang**

Seoul National University,

AI Institute, Seoul National University

SEOUL NATIONAL UNIVERSITY
VERI LUX TAS MEA

AIIS
Artificial Intelligence Institute
Seoul National University

# Brief overview

**Concept**


**Image**

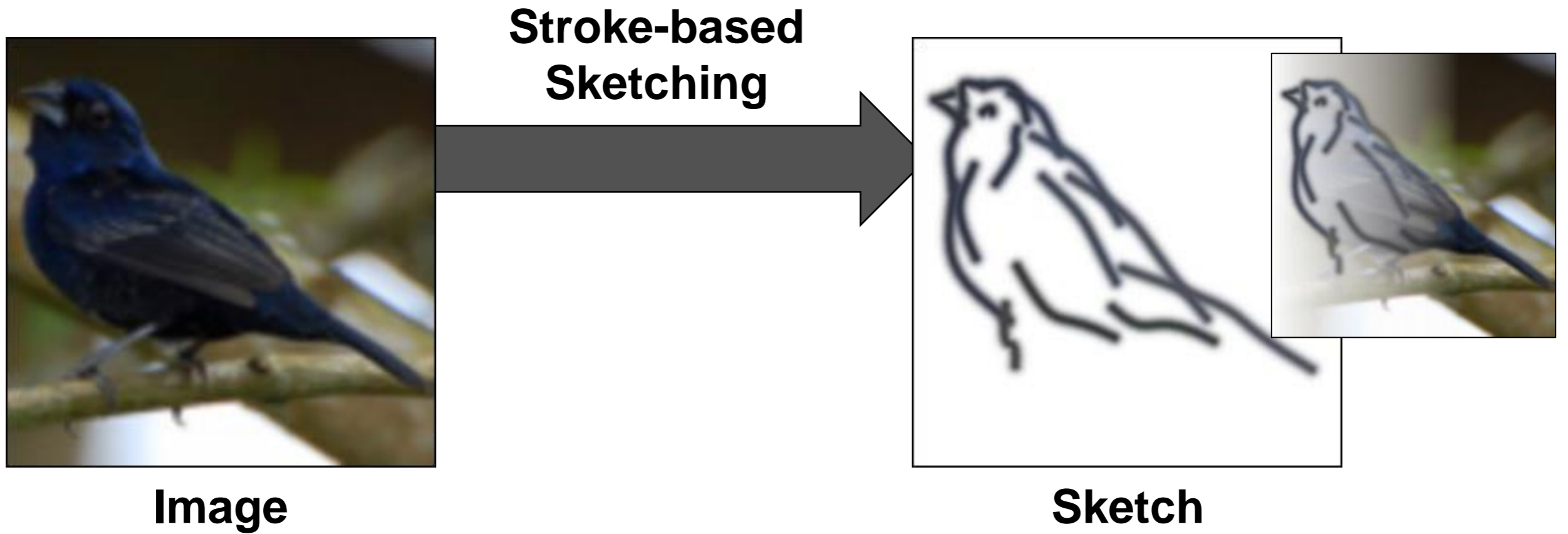**Geometric concepts**
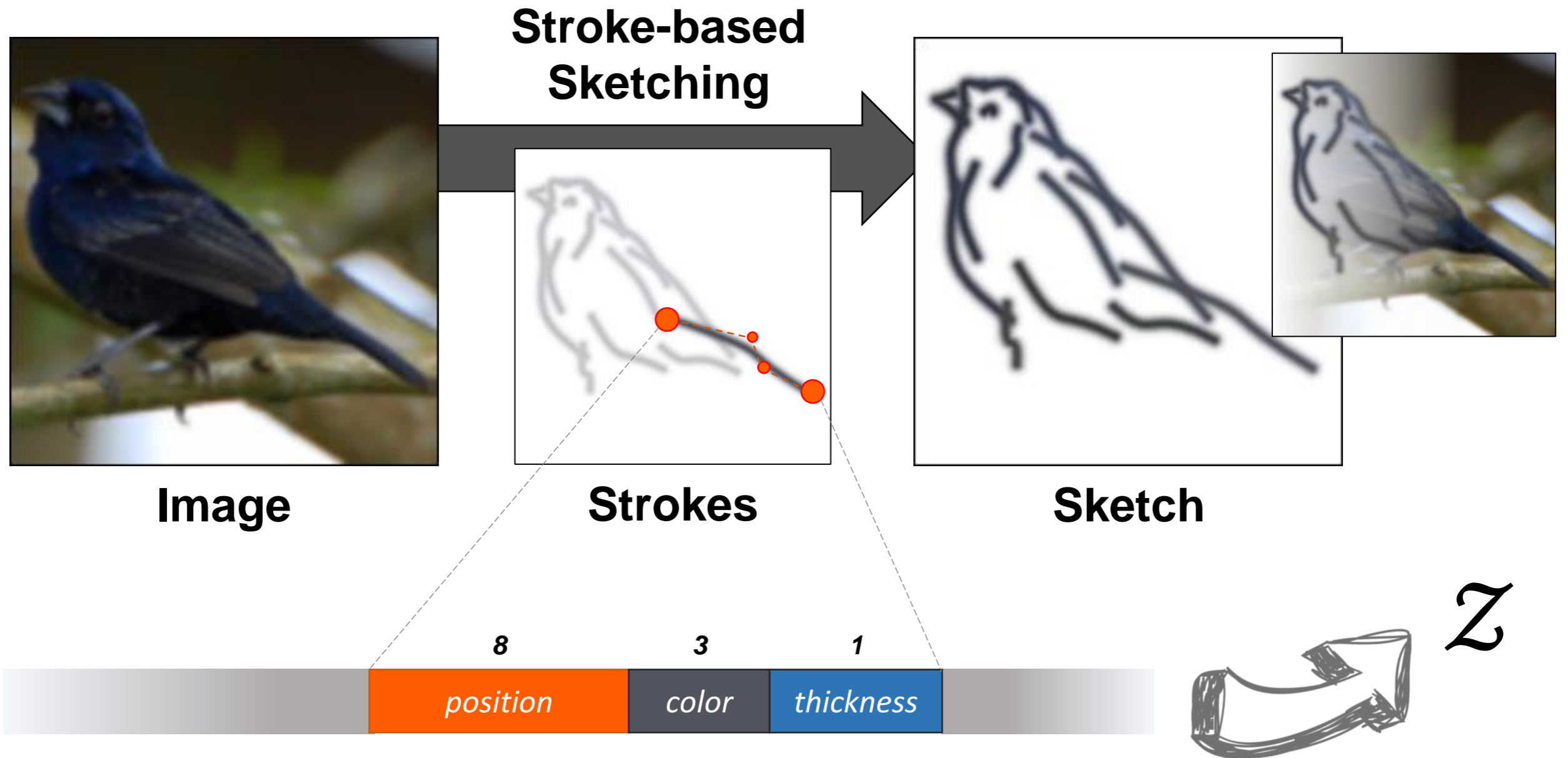ex) position, shape, distance, orientation

$\mathcal{Z}$

**Feature (vector)**

# Brief overview

**Concept**



**Image**

**Stroke-based Sketching**

**Sketch**

# Brief overview

**Concept**



| Image | Strokes | Sketch |

**Stroke-based Sketching**

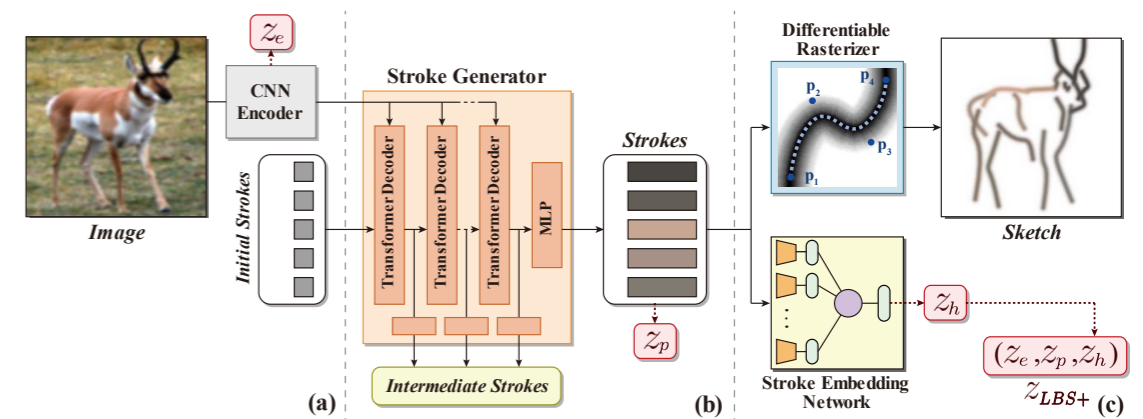| 8 | 3 | 1 |
|---|---|---|
| *position* | *color* | *thickness* |

# Brief overview

## Contributions

- Theoretically show that **strokes can convey geometric information**

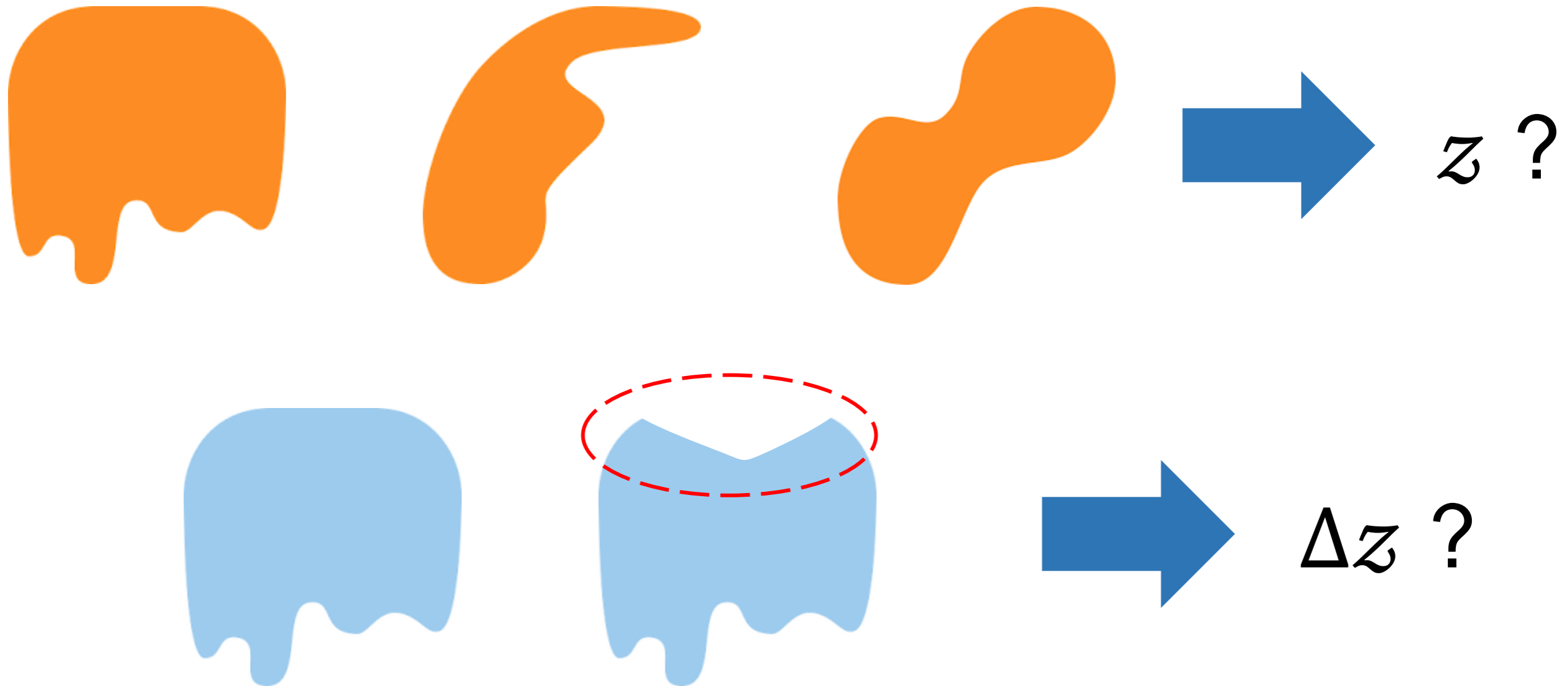- Design a sketch-based visual representation model, coined **LBS (Learning By Sketching)**



- Experiments for investigating the effectiveness of our approach in various domains and downstream tasks
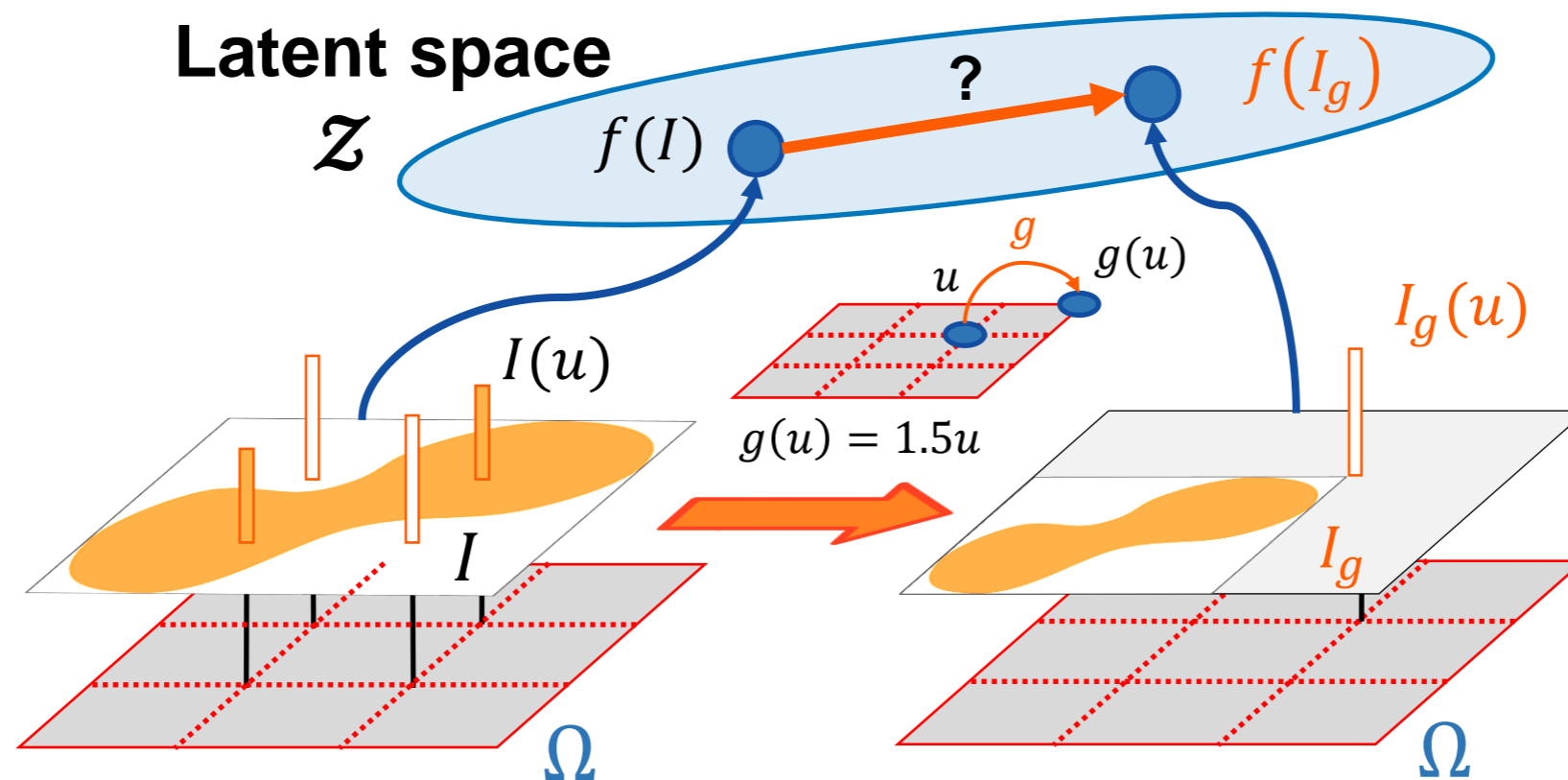
# Question

**How do machines represent shape?**

- Are they interpretable, well-aligned for downstream tasks?

# Question

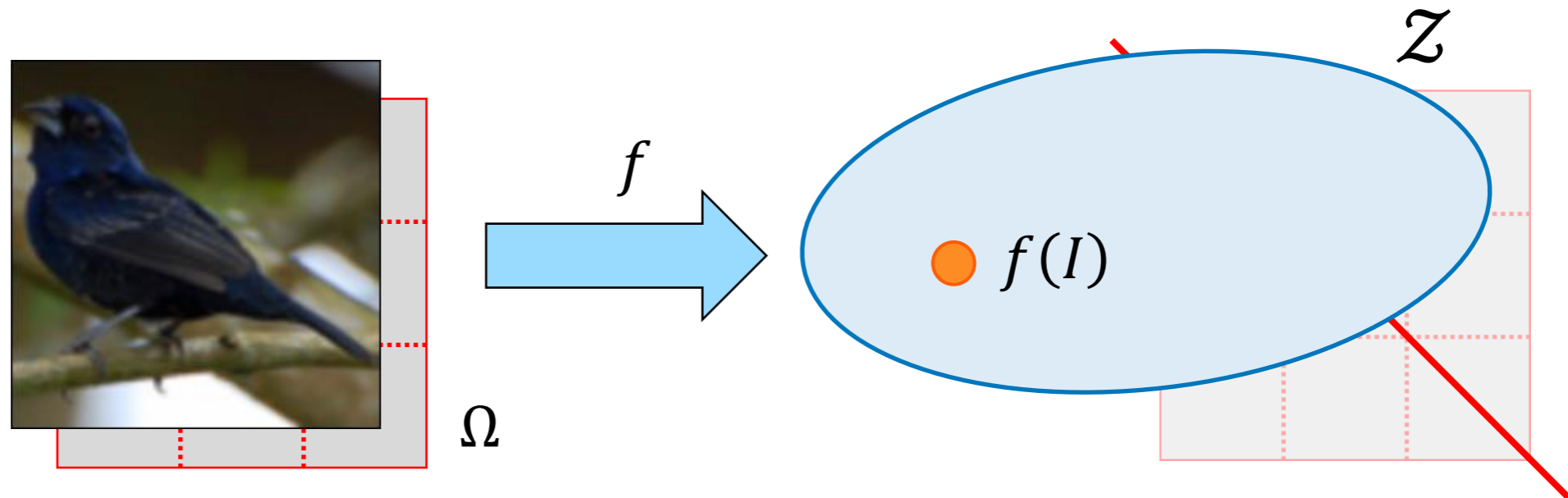## How do machines represent geometric concepts?

- Ex) position, shape, distance, orientation

- An image $I$ belongs to a physical domain $\Omega$: a 2D grid space

- For given geometric transformation $g: \Omega \to \Omega$,

  ▸ The transformed image $I_g$, $I_g(u) = I\big(g^{-1}(u)\big)$

- With given representations $f(I)$ and $f(I_g)$, can we predict the transformation $g$?
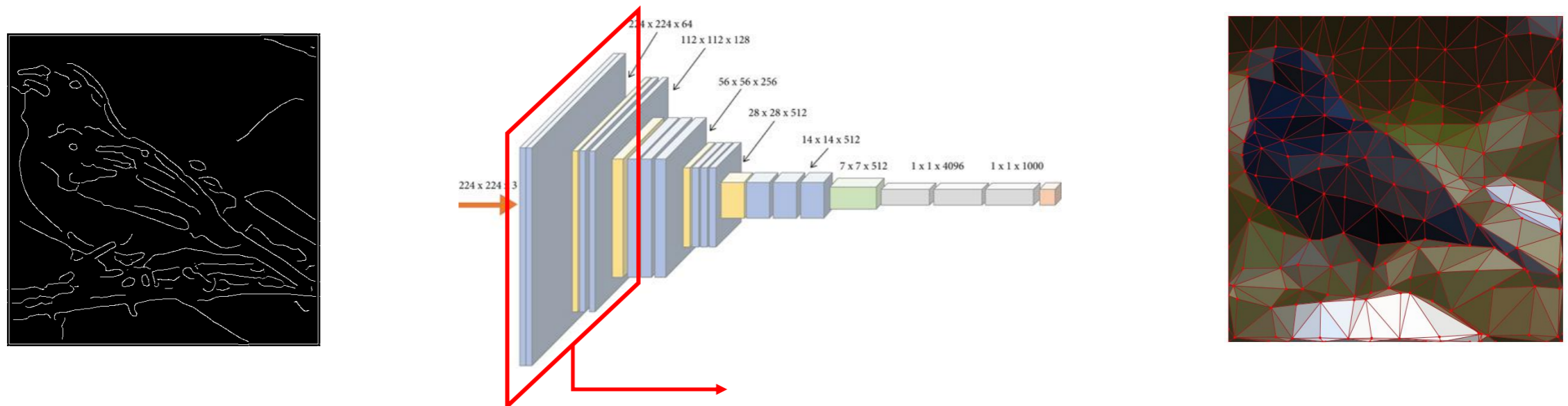
# Background

## Previous studies

- Representations on semantic latent space
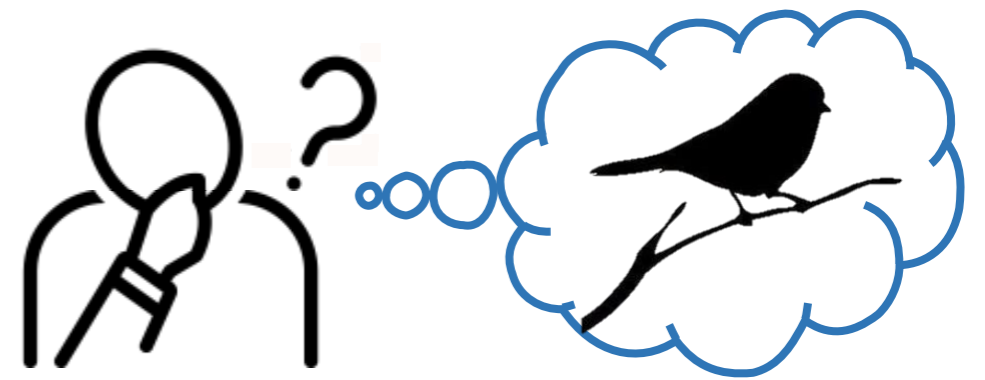


- Representations with same domain $\Omega$
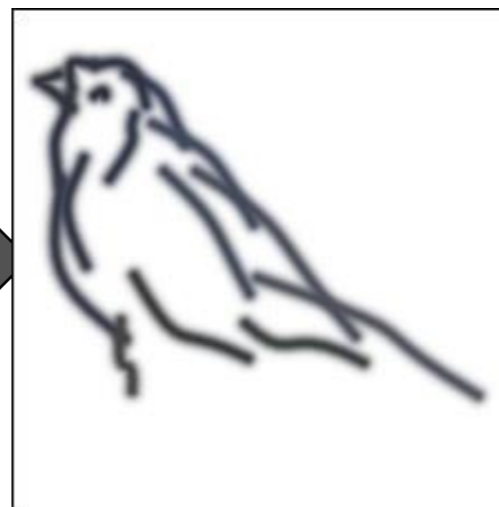
# Motivation

## How do human represent image?

- Language



- Drawing/sketching



**Image**                    **Sketch**

# Motivation

## Core properties

- Sketching

  ▸ salient features of an image $\Rightarrow$ abstract image based on set of colored strokes

  ▸ **Preserving geometry**: shares the same physical domain $\Omega$

  ▸ **Abstraction**: representation with limited # of strokes

  ▸ **Compactness**: can be represented as a set of parametric curves



| | 8 | 3 | 1 |
| --- | --- | --- | --- |
| | *position* | *color* | *thickness* |

**Stroke**                                    **Parameterization**

$\Rightarrow$ Based on these properties,
   we use **strokes** as a **geometry-aware representation** for various downstream tasks

# Mathematical framework

## Does sketches & strokes represent geometric information?

- With given representations $f(I)$ and $f(I_g)$, can we predict the transformation $g$?

- **$\mathcal{G}$-equivariance**: $\exists \rho' \ s.t. \ f(\rho(g) \cdot I) = \rho'(g) \cdot f(I)$

  - Sketching is equivariant to arbitrary geometric transformation $g \in \mathcal{G}$ $(\rho' = \rho)$

  - Converting into stroke is equivariant to affine transformation $a \in \mathcal{A}$

# Architecture

## Learning By Sketching (LBS)

- Abstraction & reflecting geometric information in a short inference time

- Without using explicit sketch dataset.



$$\mathcal{L}_{LBS} = \mathcal{L}_{percept} + \lambda_g \cdot \mathcal{L}_{guide} + \lambda_e \cdot \mathcal{L}_{embed}$$

# Architecture

## Learning By Sketching (LBS)

- $\mathcal{L}_{percept}$: CLIP-based perceptual loss [1]



[1] Vinker, Yael, et al. "Clipasso: Semantically-aware object sketching." *ACM Transactions on Graphics (TOG)* 41.4 (2022): 1-11.

# Architecture

## Learning By Sketching (LBS)

- Guidance loss $\mathcal{L}_{guide}$

  ▸ Optimization-based generation with $\mathcal{L}_{percept}$: $\boldsymbol{p}_{gt}$

  ▸ Predicting $\Delta\boldsymbol{p}_{gt}$ for each layer

# Architecture

## Learning By Sketching (LBS)

- Stroke embedding loss $\mathcal{L}_{embed}$

  - ▸ $z_h$: combines the information of all strokes through a stroke embedding network

# Experiments

## Quantitative results

- Understanding geometric primitives & concepts



Euclidean geometry concepts     10 realizations per class     concept classification

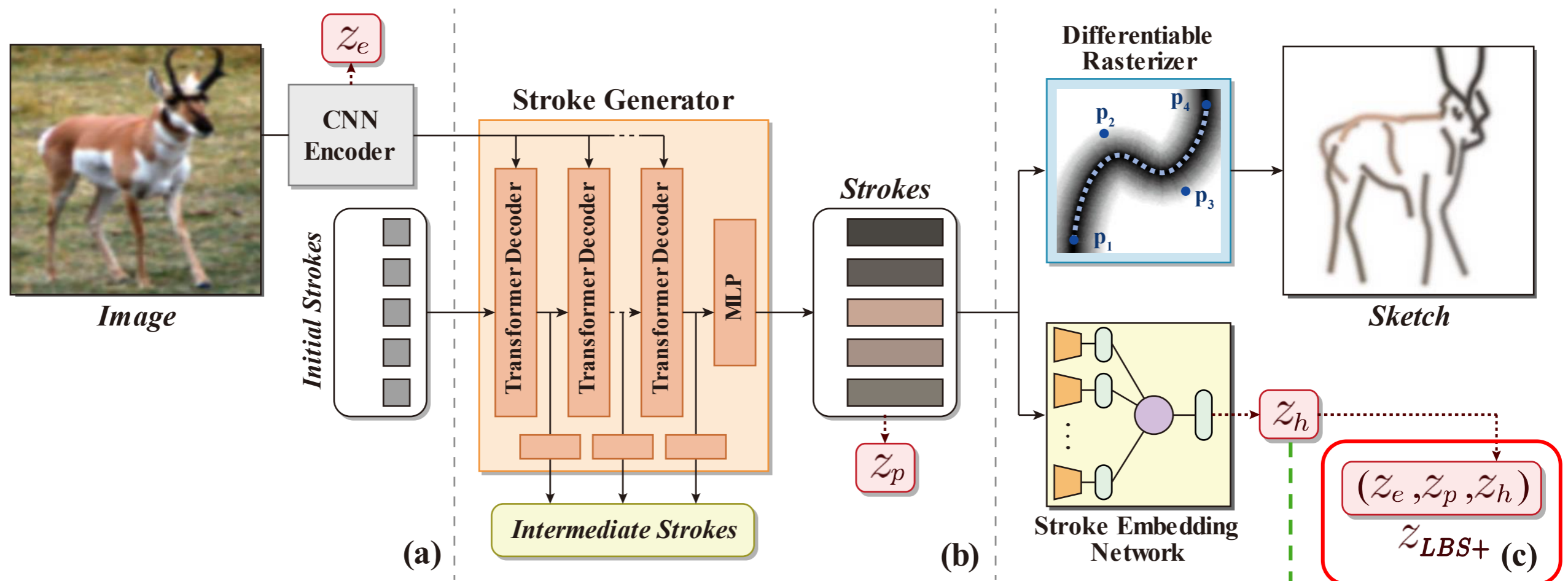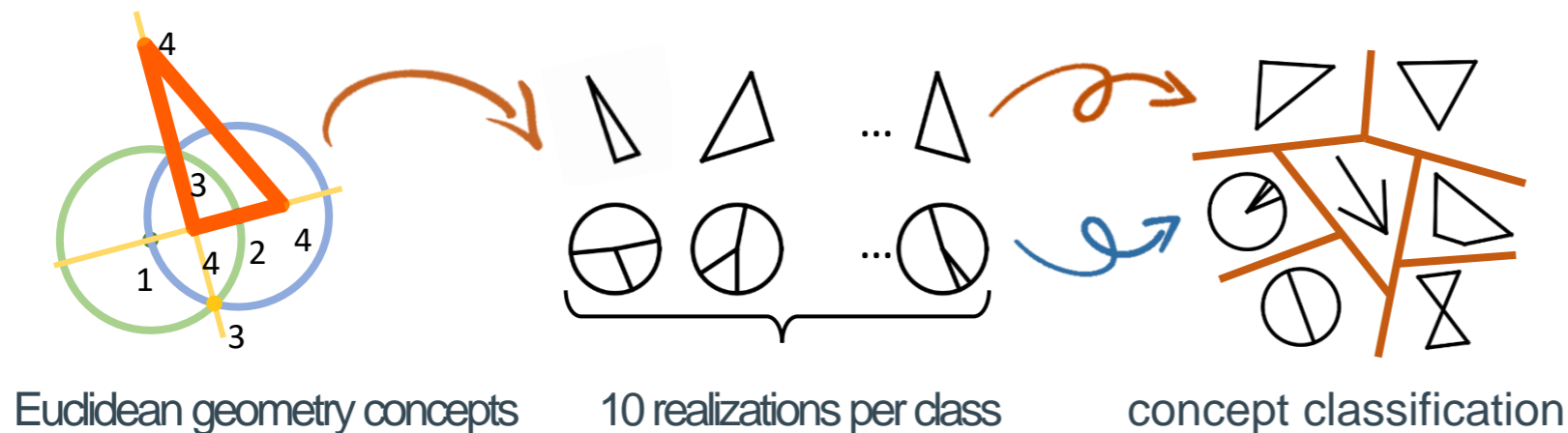| La-bel | Method | Geoclidean | |
|---|---|---|---|
| | | Constraints | Elements |
| | CE | $53.89_{\pm1.58}$ | $70.57_{\pm4.29}$ |
| | SupCon [38] | $42.41_{\pm3.16}$ | $55.83_{\pm4.28}$ |
| ✓ | LtD-diff [54] | $57.26_{\pm2.19}$ | $69.47_{\pm2.11}$ |
| | E(2)-CNN [73] | $\mathbf{71.03}_{\pm1.94}$ | $69.28_{\pm1.46}$ |
| | **LBS (CE)** | $50.01_{\pm1.58}$ | $\mathbf{81.06}_{\pm3.14}$ |
| | SimCLR [9] | $32.04_{\pm0.64}$ | $65.14_{\pm4.11}$ |
| | $\beta$-TCVAE [8] | $17.18_{\pm1.35}$ | $33.82_{\pm1.64}$ |
| ✗ | GeoSSL [55] | $18.66_{\pm3.33}$ | $33.47_{\pm2.80}$ |
| | HoG [14] | $23.82$ | $52.05$ |
| | **LBS** | $\mathbf{47.43}_{\pm1.34}$ | $\mathbf{81.34}_{\pm0.16}$ |

- Local geometric information & Simple spatial reasoning
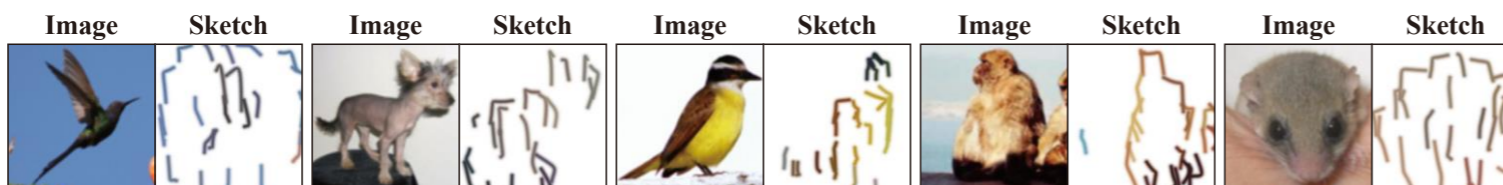


Q) **Color** of the **leftmost** object?

Q) **Shape** of the rightmost object?

Q) Results after **shifting** the rightmost object?

…

Label: brown

| Label | Method | RC | BC | Size | Shape | Third | Shift |
|---|---|---|---|---|---|---|---|
| | CE | $98.71_{\pm0.10}$ | $76.02_{\pm0.90}$ | $92.51_{\pm0.40}$ | $49.97_{\pm0.29}$ | $40.66_{\pm0.20}$ | $\mathbf{62.06}_{\pm1.62}$ |
| | SupCon [38] | $\mathbf{98.75}_{\pm0.08}$ | $66.04_{\pm2.65}$ | $91.88_{\pm0.37}$ | $49.15_{\pm0.85}$ | $37.93_{\pm0.27}$ | $56.05_{\pm2.56}$ |
| ✓ | LtD-diff [54] | $62.29_{\pm0.48}$ | $15.84_{\pm0.43}$ | $63.98_{\pm3.38}$ | $43.96_{\pm3.05}$ | $16.47_{\pm0.59}$ | $17.21_{\pm0.29}$ |
| | E(2)-CNN [73] | $98.50_{\pm0.10}$ | $73.51_{\pm2.50}$ | $89.84_{\pm0.46}$ | $45.85_{\pm0.93}$ | $\mathbf{41.95}_{\pm0.18}$ | $59.29_{\pm0.91}$ |
| | **LBS (CE)** | $97.49_{\pm0.22}$ | $\mathbf{84.09}_{\pm0.84}$ | $\mathbf{93.22}_{\pm0.29}$ | $\mathbf{70.03}_{\pm0.68}$ | $38.23_{\pm0.25}$ | $51.56_{\pm0.16}$ |
| | SimCLR [9] | $60.61_{\pm1.24}$ | $63.77_{\pm2.29}$ | $83.35_{\pm0.60}$ | $41.95_{\pm0.33}$ | $33.42_{\pm0.55}$ | $43.05_{\pm0.55}$ |
| | E(2)-CNN [73] | $53.50_{\pm7.30}$ | $55.52_{\pm7.60}$ | $83.52_{\pm1.56}$ | $42.06_{\pm1.12}$ | $30.74_{\pm2.84}$ | $38.03_{\pm4.44}$ |
| | $\beta$-TCVAE [8] | $17.09_{\pm0.20}$ | $20.04_{\pm0.71}$ | $71.27_{\pm0.10}$ | $36.30_{\pm0.10}$ | $15.38_{\pm0.19}$ | $16.35_{\pm0.18}$ |
| ✗ | GeoSSL [55] | $20.16_{\pm0.63}$ | $21.61_{\pm0.67}$ | $73.79_{\pm0.78}$ | $44.08_{\pm1.10}$ | $15.39_{\pm0.16}$ | $16.94_{\pm0.34}$ |
| | DefGrid [21] | $73.81_{\pm0.91}$ | $73.38_{\pm0.80}$ | $81.50_{\pm0.22}$ | $46.34_{\pm0.77}$ | $24.90_{\pm0.27}$ | $36.28_{\pm0.13}$ |
| | **LBS** | $\mathbf{84.31}_{\pm0.08}$ | $\mathbf{83.00}_{\pm0.39}$ | $\mathbf{92.66}_{\pm0.41}$ | $\mathbf{70.01}_{\pm0.53}$ | $\mathbf{37.41}_{\pm0.29}$ | $\mathbf{49.32}_{\pm0.17}$ |
| | CLIP [59] | $37.39$ | $54.98$ | $77.51$ | $66.91$ | $34.75$ | $34.80$ |
| | HoG [14] | $56.83$ | $58.69$ | $81.73$ | $61.14$ | $24.28$ | $33.25$ |

- Domain transfer



(b) CLEVR → STL-10

| Dataset | Labeled | | Unlabeled | |
|---|---|---|---|---|
| | Method | Accuracy | Method | Accuracy |
| | CE | $46.15_{\pm0.12}$ | SimCLR [9] | $41.68_{\pm0.05}$ |
| CLEVR | SupCon [38] | $43.41_{\pm0.36}$ | $\beta$-TCVAE [8] | $27.35_{\pm0.38}$ |
| ↓ | LtD-diff [54] | $50.81_{\pm0.67}$ | GeoSSL [55] | $35.93_{\pm0.96}$ |
| STL-10 | E(2)-CNN [73] | $45.19_{\pm0.84}$ | E(2)-CNN [73] | $38.50_{\pm0.49}$ |
| | **LBS (CE)** | $\mathbf{56.48}_{\pm0.89}$ | DefGrid [21] | $33.13_{\pm0.17}$ |
| | | | **LBS** | $55.35_{\pm0.18}$ |

# Experiments

## Qualitative results



- Progressive optimization process



Initial                                                                 Final

# Thank you