



MAX PLANCK INSTITUTE
FOR INFORMATICS

SIC Saarland Informatics
Campus

ICCV  **HONOLULU**
OCT 19-23, 2025 **HAWAII**

VITAL: More Understandable Feature Visualization through Distribution Alignment and Relevant Information Flow



Ada Görgün



Bernt Schiele



Jonas Fischer

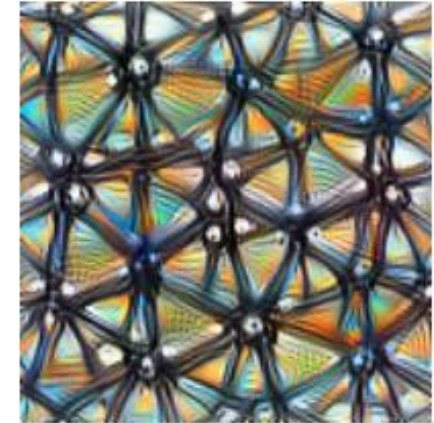
Max Planck Institute for Informatics, Saarland Informatics Campus

Mechanistic Interpretability: Feature Visualization

Attribution Methods vs Feature Visualization



Question: What part of an example is responsible for the network activating a particular way?

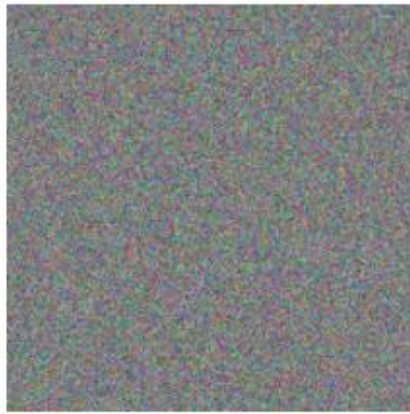


Question: What a network or parts of a network are looking for by generating examples?

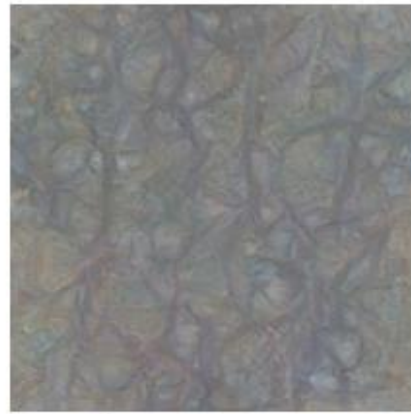
Mechanistic Interpretability: Feature Visualization

- How?

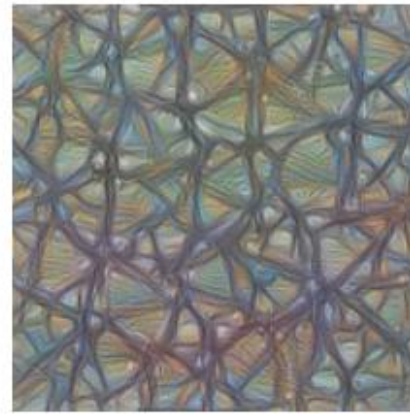
- starting from random noise, optimize an image to activate a particular neuron



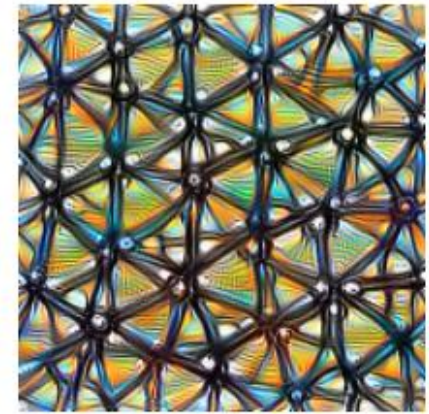
Step 0



Step 4



Step 48

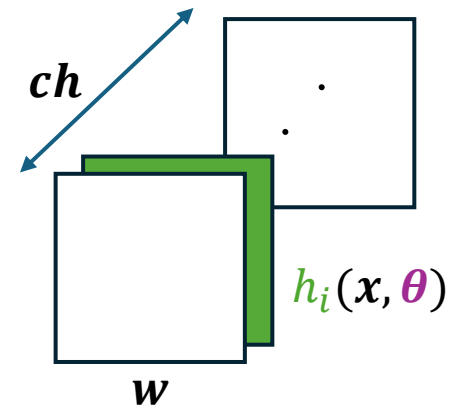


Step 2048

- In a way, this is **the opposite of the learning process**:

- Normally, we have input data and want to learn the right weights θ
- Here, we have the weights and want to generate data samples

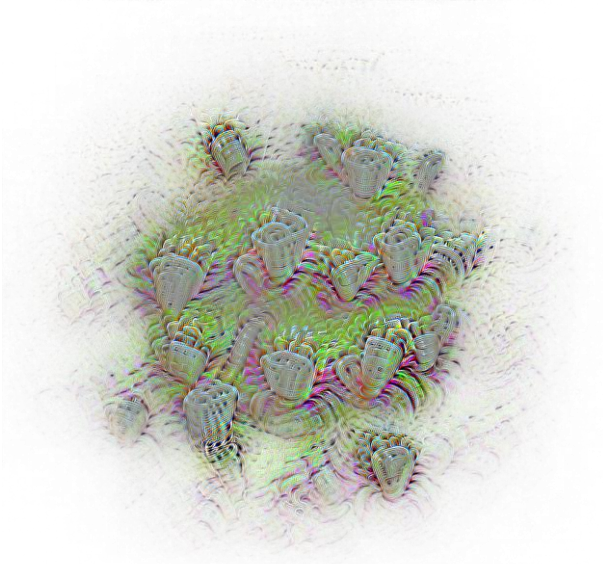
$$x^* = \operatorname{argmax}_x h_i(x, \theta)$$



Feature Visualization: Challenges of Existing Efforts

- Unrealistically Generated Images:

- Visualized features produce abstract, non-realistic images with repetitive patterns [1, 2].
- Results often fail to resemble recognizable patterns or objects [1, 2].
- Fails in larger models [1, 3].
- Leverage statistically learned priors from auxiliary models [3].



[1] Fourier FV, Olah et. al, 2017



[2] MACO, Fel et. al, NeurIPS 2023



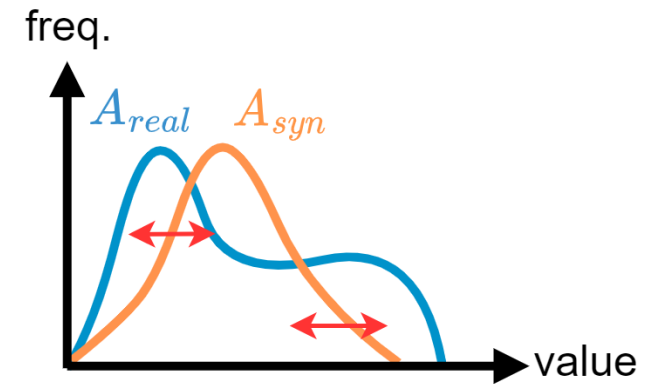
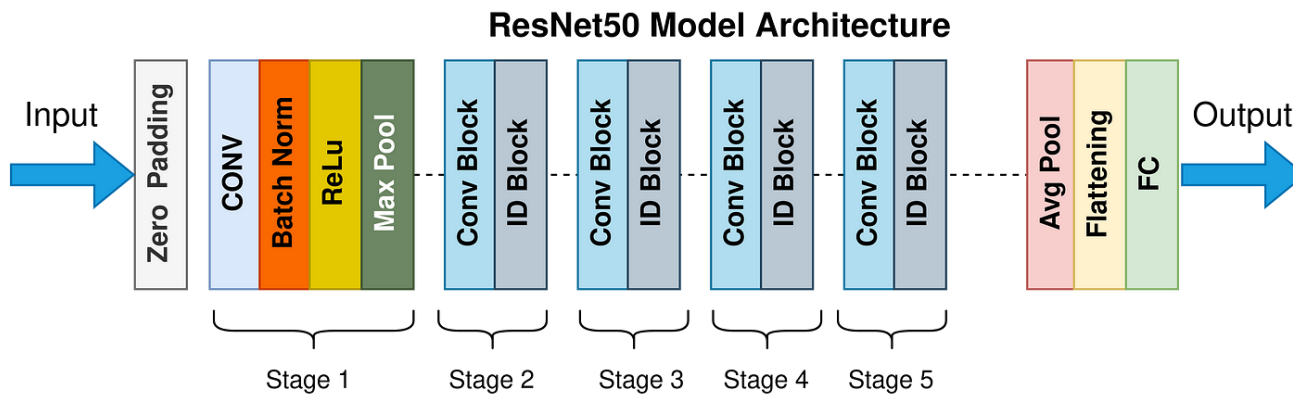
[3] GANs, Nyugen et. al, NeurIPS 2016

VITAL: Motivation and Contributions

- Matching the feature distribution of real images and the synthetic image to:
 - generate a **more interpretable** visualization
 - generate **non-repetitive patterns**
 - capture the intrinsic characteristics of real images
 - **not directly rely on activation maximization**
- Utilization of the method for a more intuitive grasp on:
 - how the **flow of the network** affects the visualized image
- Evaluation of the generated image in terms of:
 - how well it reflects the real images
 - **quantitative** and **qualitative** metrics

VITAL: Proposed Method

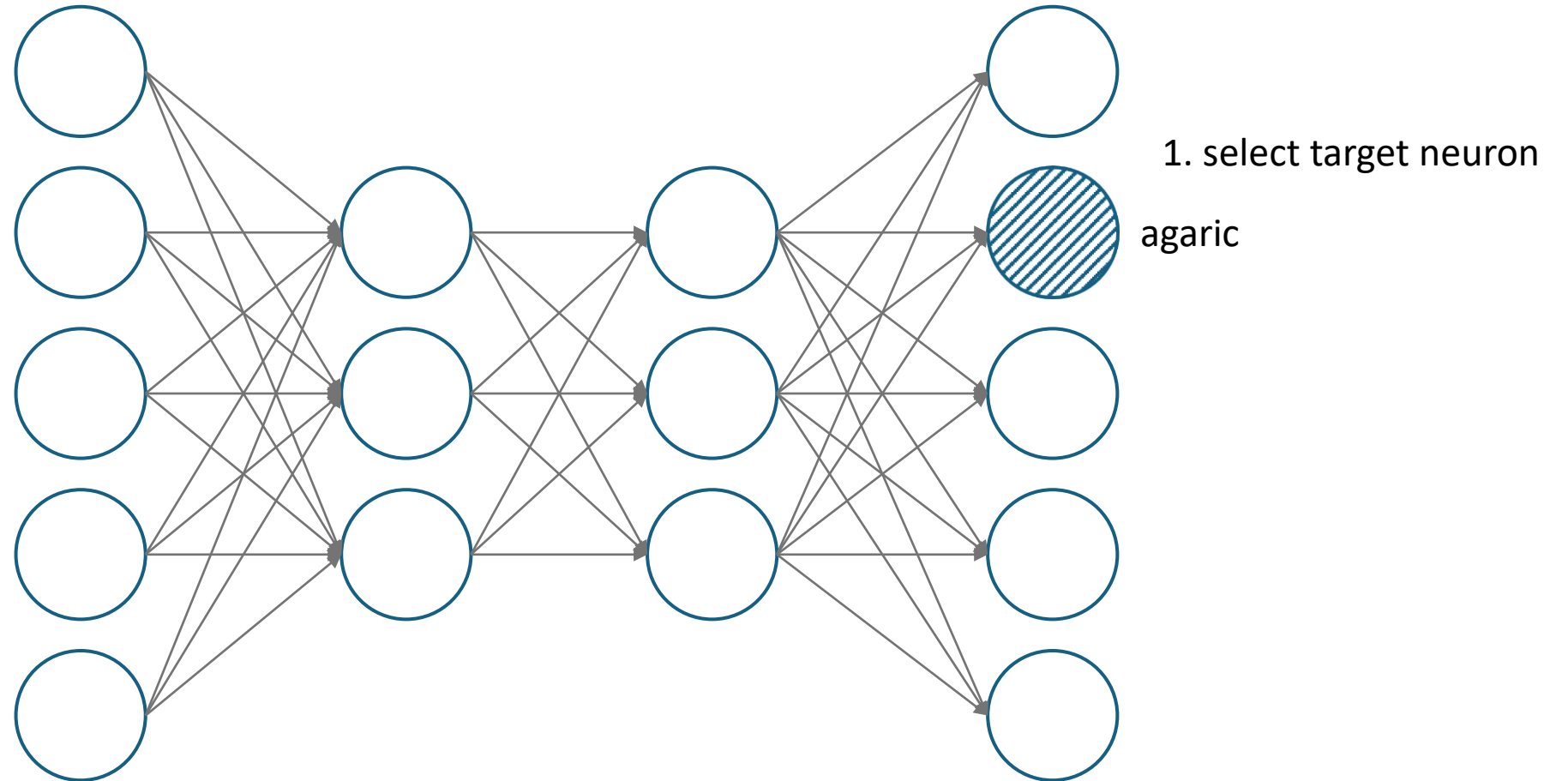
- We want to **match** the feature distribution of real (A_{real}) images and the synthetic (A_{syn}) image.



- We want to achieve visualizations for
 - class neurons
 - intermediate neurons

Matching is done by the sort matching algorithm (Zhang et al. CVPR 2022)

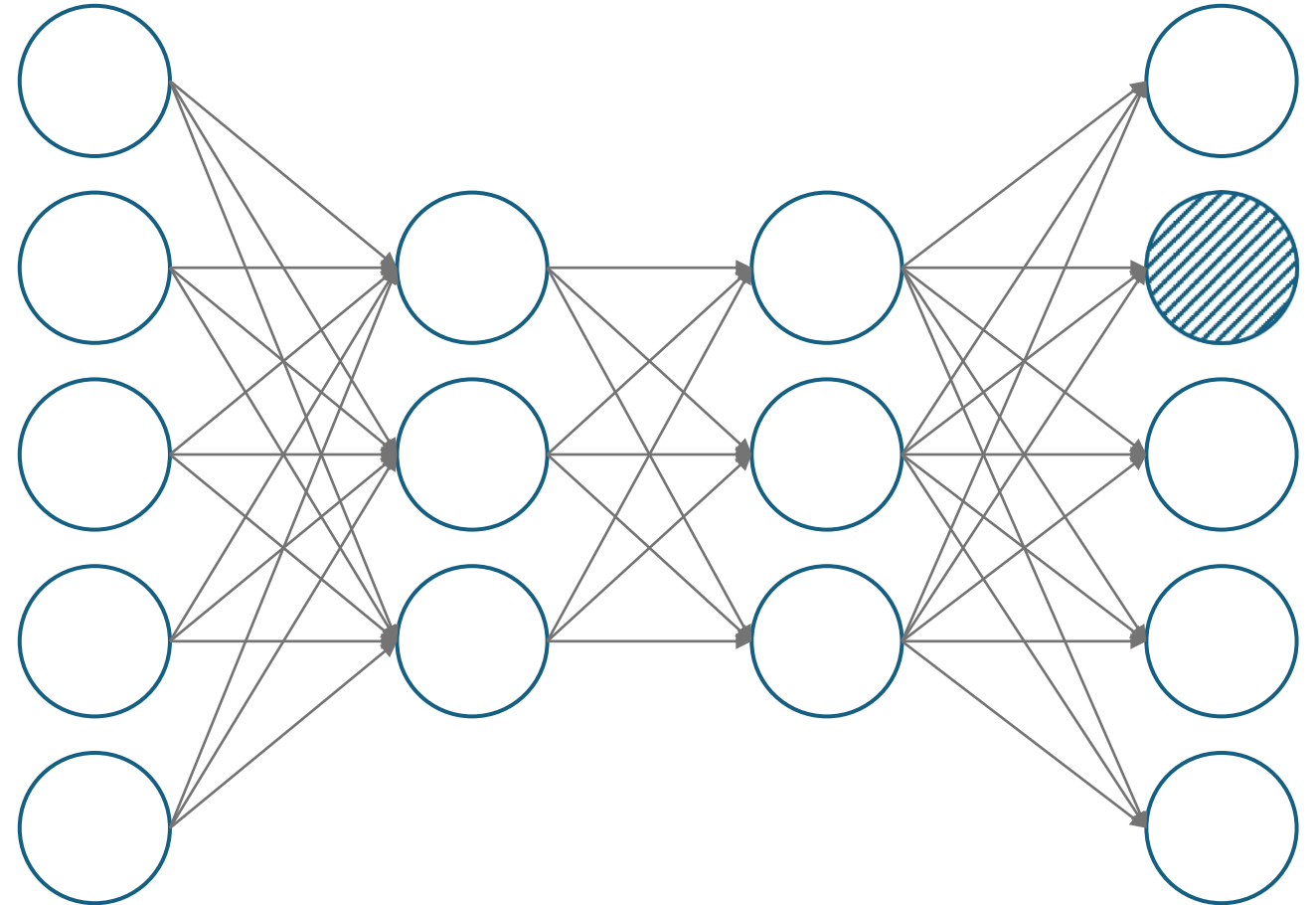
VITAL: Class Neuron Visualization



Pretrained Model Architecture

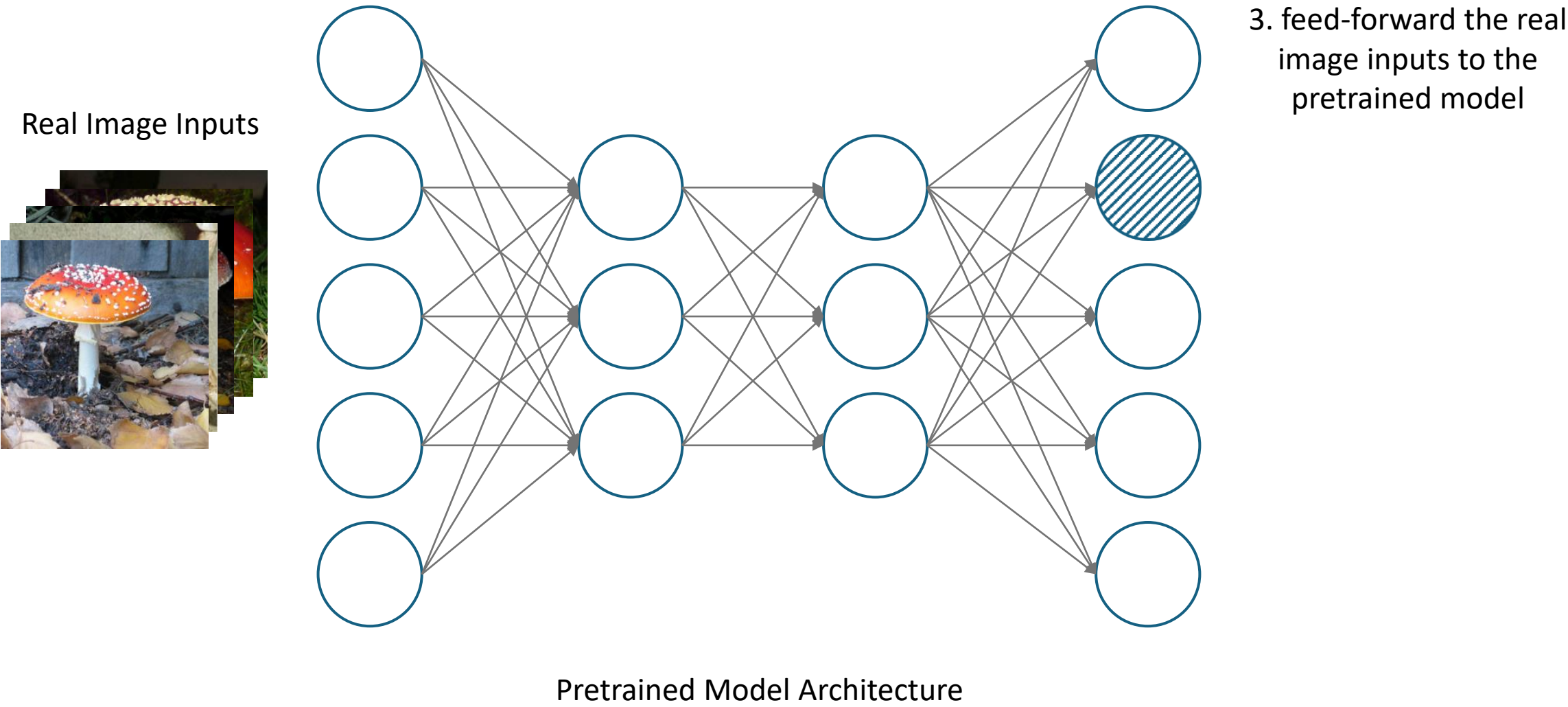
VITAL: Class Neuron Visualization

2. select random N images of the target class neuron



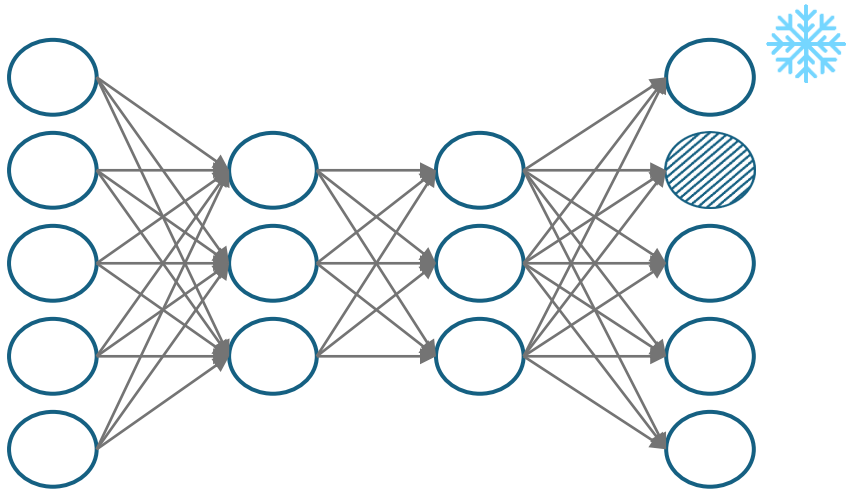
Pretrained Model Architecture

VITAL: Class Neuron Visualization



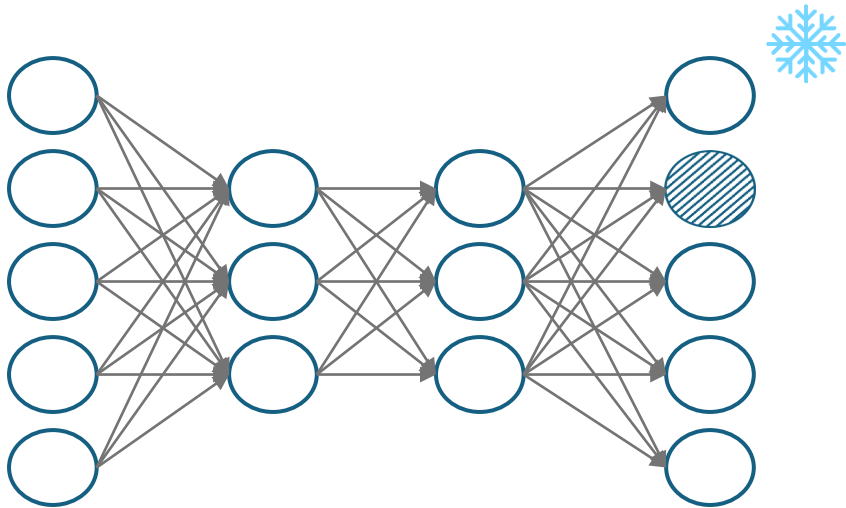
VITAL: Class Neuron Visualization

Real Image Inputs



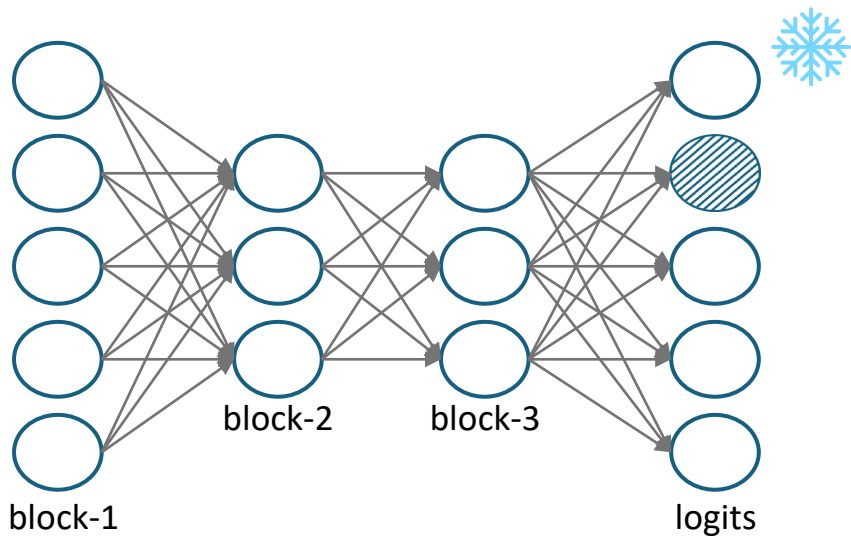
4. generate a noisy input and feed-forward to the same frozen model

Synthetic Input

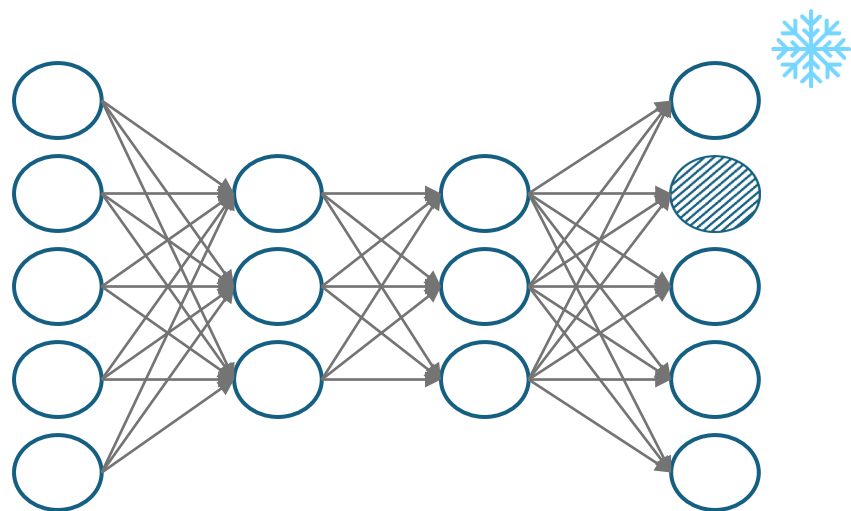


VITAL: Class Neuron Visualization

Real Image Inputs

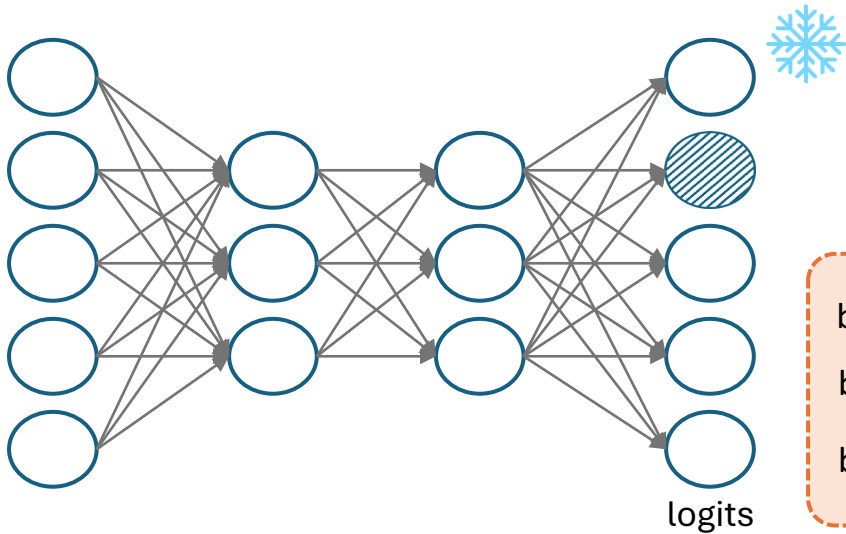


Synthetic Input



VITAL: Class Neuron Visualization

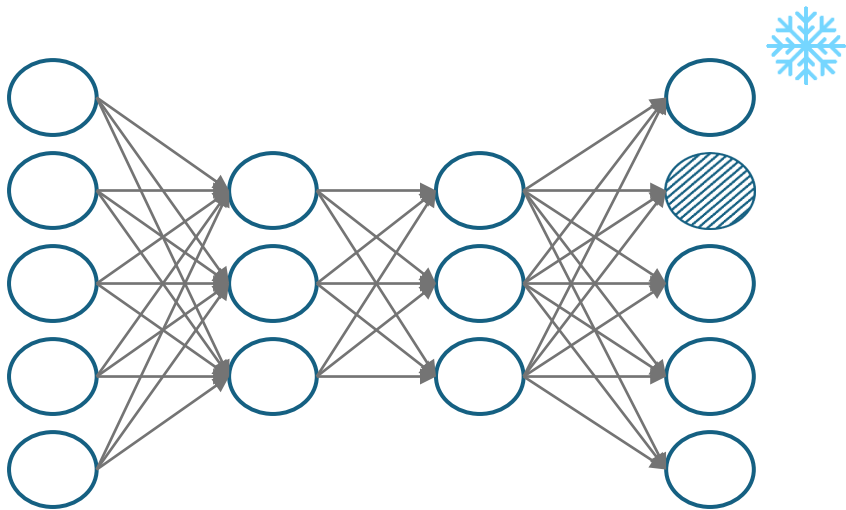
Real Image Inputs



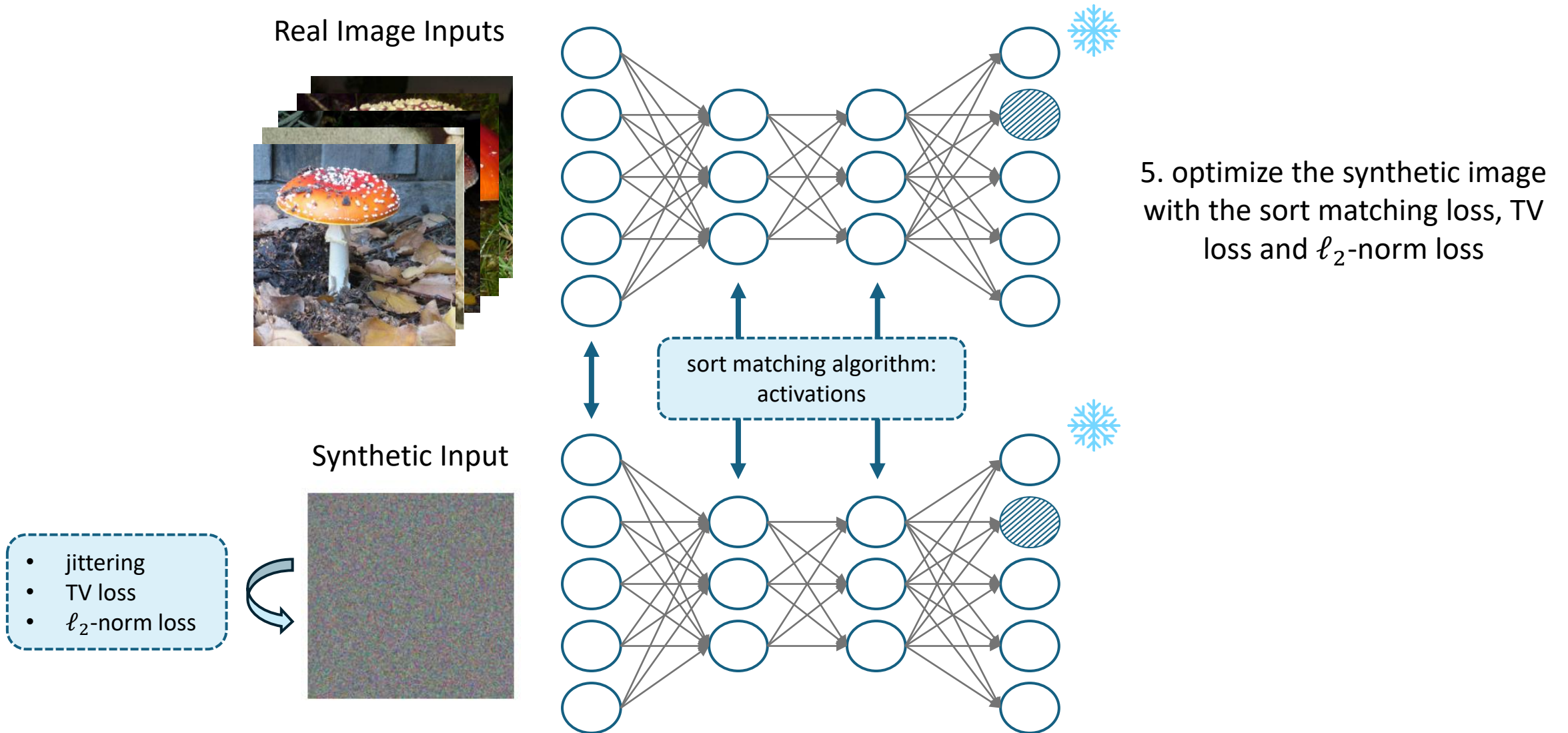
block-1
block-2
block-3

These blocks will be used for aligning the distribution of latent features in each block using the **sort-matching algorithm**.

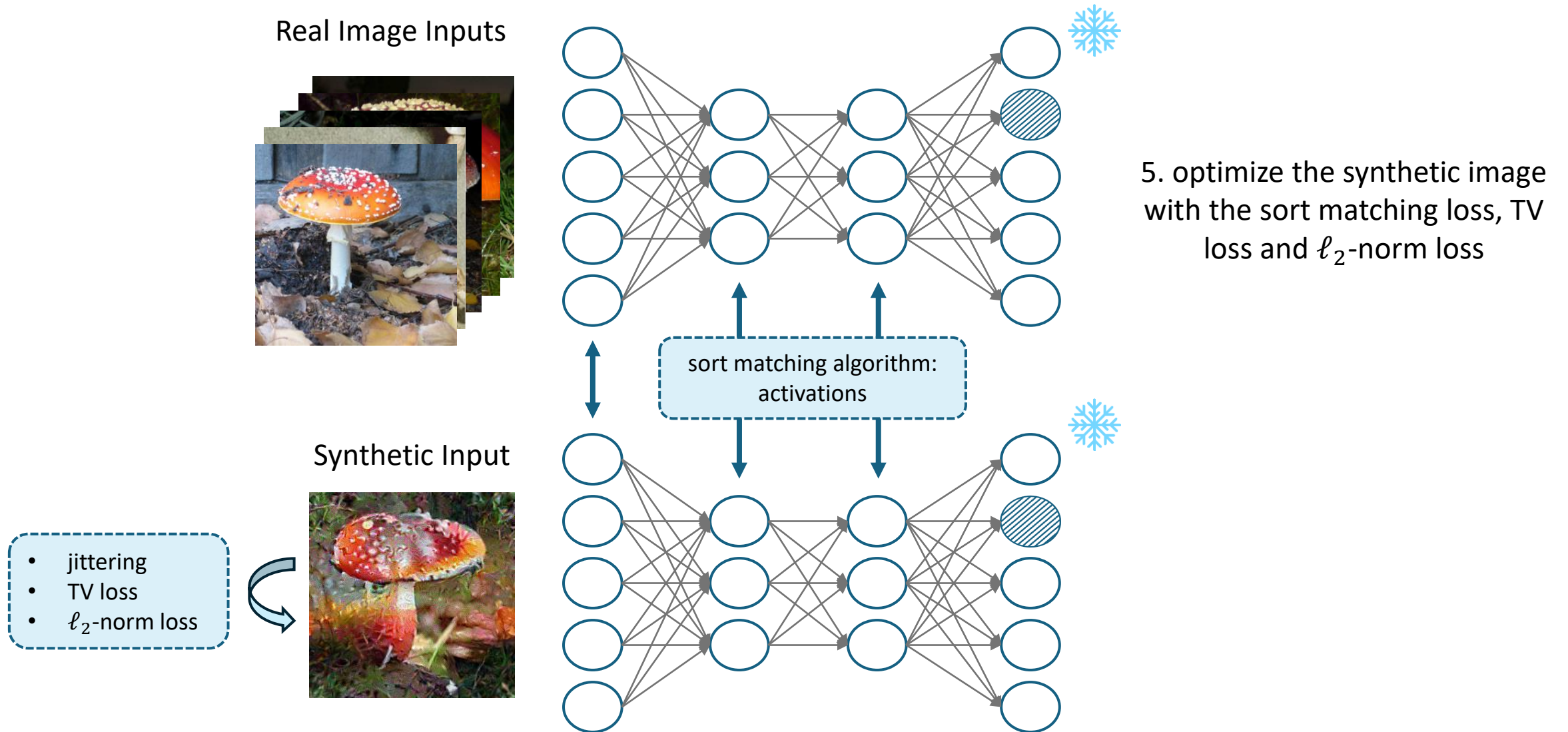
Synthetic Input



VITAL: Class Neuron Visualization



VITAL: Class Neuron Visualization



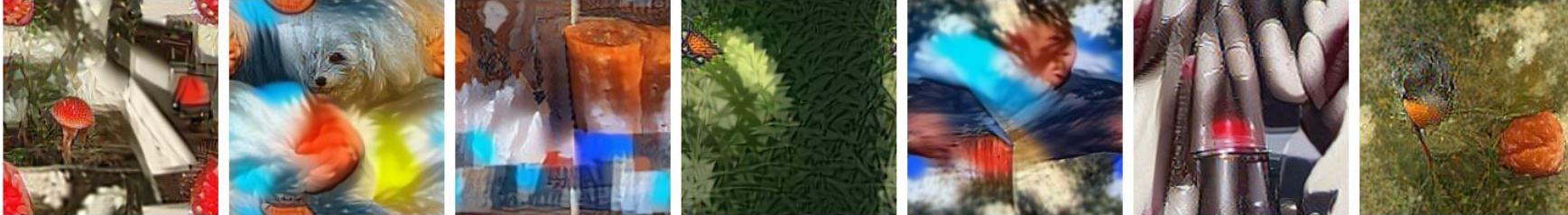
VITAL: Class Neuron Visualization

- Model tested: ResNet50

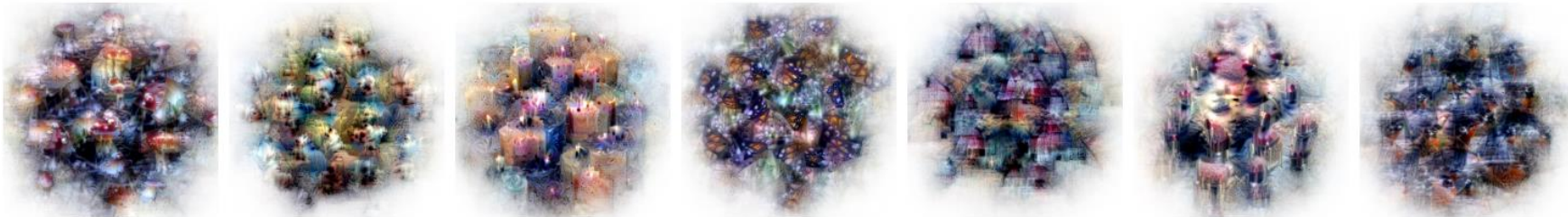
VITAL



DeepInversion



MACO



agaric

maltese dog

candle

butterfly

barn

lipstick

robin

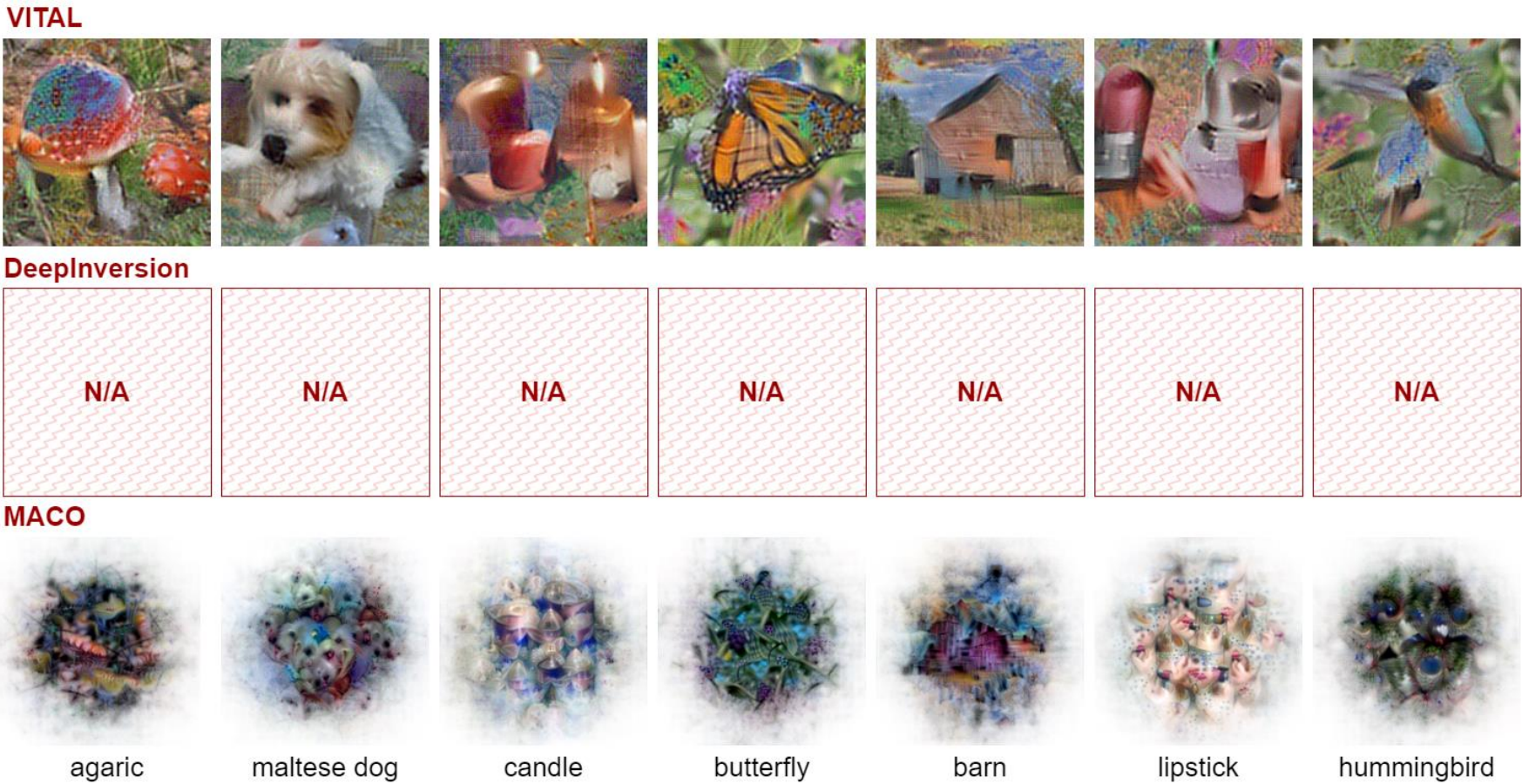
VITAL: Class Neuron Visualization

- Model tested: ConvNext-base



VITAL: Class Neuron Visualization

- Model tested: ViT-L-32



VITAL: Class Neuron Visualization

- Quantitative Results

- ResNet50, ConvNeXt, DenseNet121, ViT-L-16, ViT-L-32
- Compared Methods: Fourier, MACO, DeepInversion

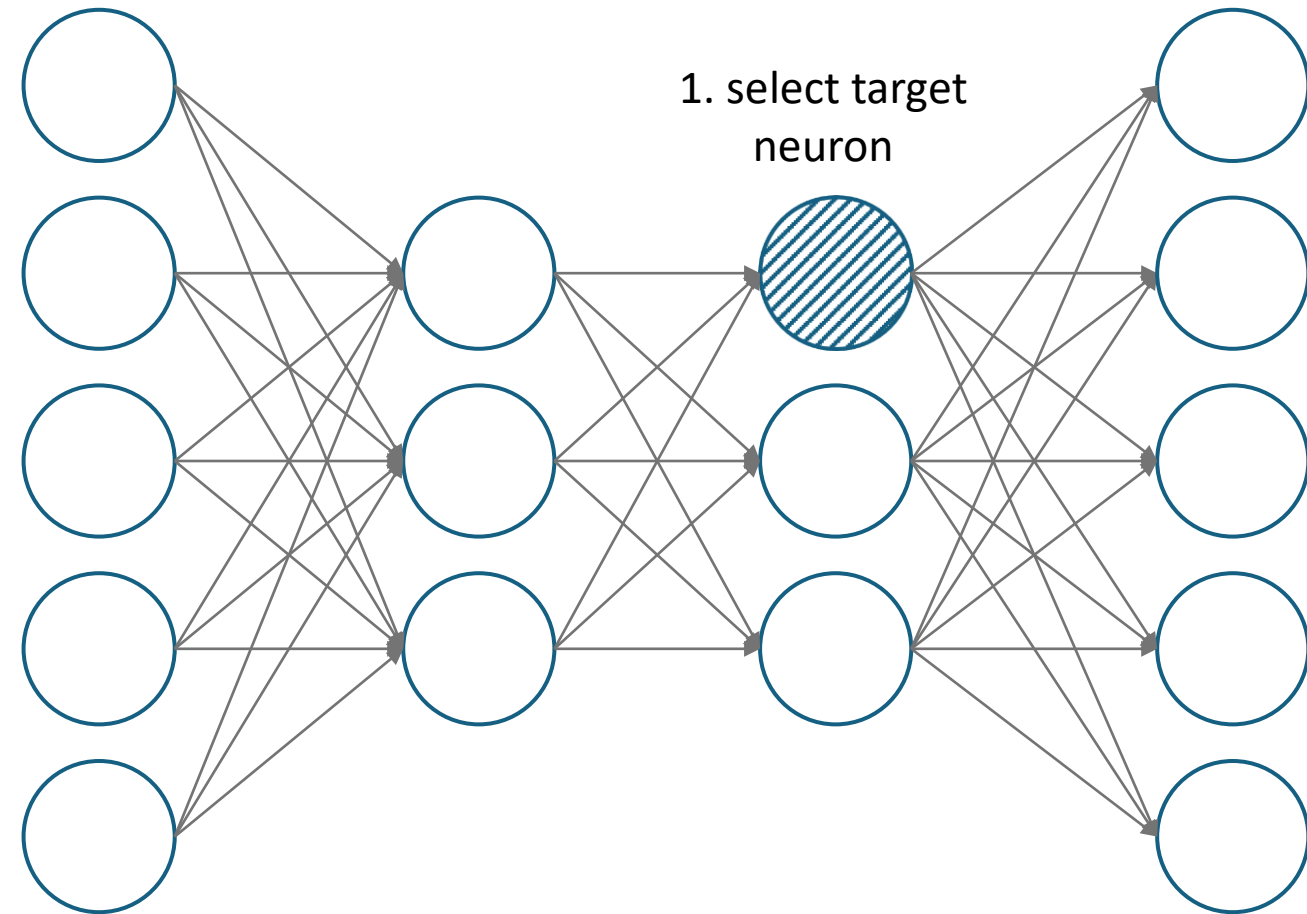
- Metrics:

- Top1 Accuracy** of the generated images
- FID score** to quantify visualization realism
- CLIP Zero-Shot Prediction** to assess general perceptual quality

	Method	Setup	Acc.	FID (↓)	Zero-Shot Prediction	
			Top1 (↑)		Top1 (↑)	Top5 (↑)
ResNet50	<i>ImageNet</i>		-	-	69.11	92.23
	MACO	r: 224	29.43	360.74	12.87	29.73
	Fourier	r: 224	21.30	422.44	6.73	18.27
	DeepInv	bs: 64	100.00	35.76	29.90	55.20
	VITAL		<u>99.90</u>	<u>58.79</u>	66.62	92.56
ConvNeXt	<i>ImageNet</i>		-	-	65.66	89.80
	MACO	r: 224	66.07	62.55	<u>7.20</u>	<u>19.77</u>
	Fourier	r: 224	60.07	<u>59.60</u>	<u>2.77</u>	<u>8.30</u>
	VITAL		99.97	3.92	63.53	90.30
	DenseNet121	<i>ImageNet</i>		-	-	70.64
MACO		r: 224	9.20	1.80	9.33	23.20
Fourier		r: 224	15.53	1.63	4.87	12.17
DeepInv		bs: 64	100.00	0.20	10.00	25.47
VITAL			<u>99.93</u>	<u>0.27</u>	58.70	86.93
ViT-L-16	<i>ImageNet</i>		-	-	64.78	89.31
	MACO	r: 224	44.33	<u>946.96</u>	<u>3.93</u>	<u>10.57</u>
	Fourier	r: 224	25.30	<u>990.51</u>	1.67	5.13
	VITAL		99.80	126.29	68.17	92.80
	ViT-L-32	<i>ImageNet</i>		-	-	65.83
MACO		r: 224	<u>24.87</u>	2318.90	<u>17.53</u>	<u>37.23</u>
Fourier		r: 224	17.03	<u>1983.09</u>	10.30	28.10
VITAL			89.60	147.33	55.97	85.47



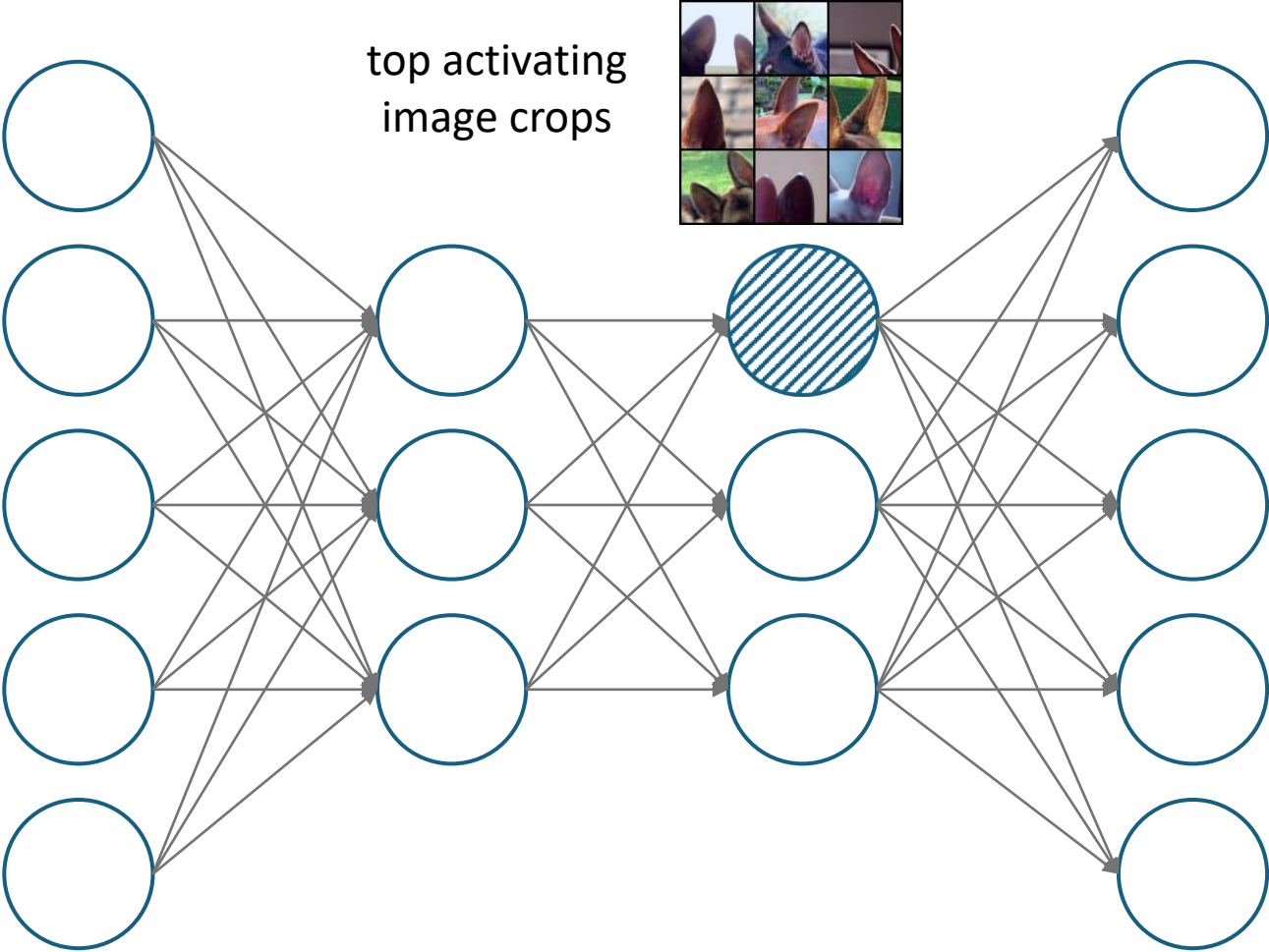
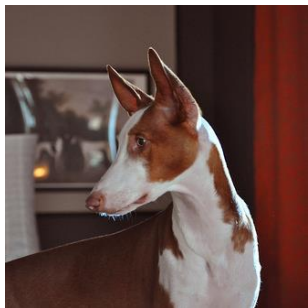
VITAL: Intermediate Neuron Visualization



Pretrained Model Architecture

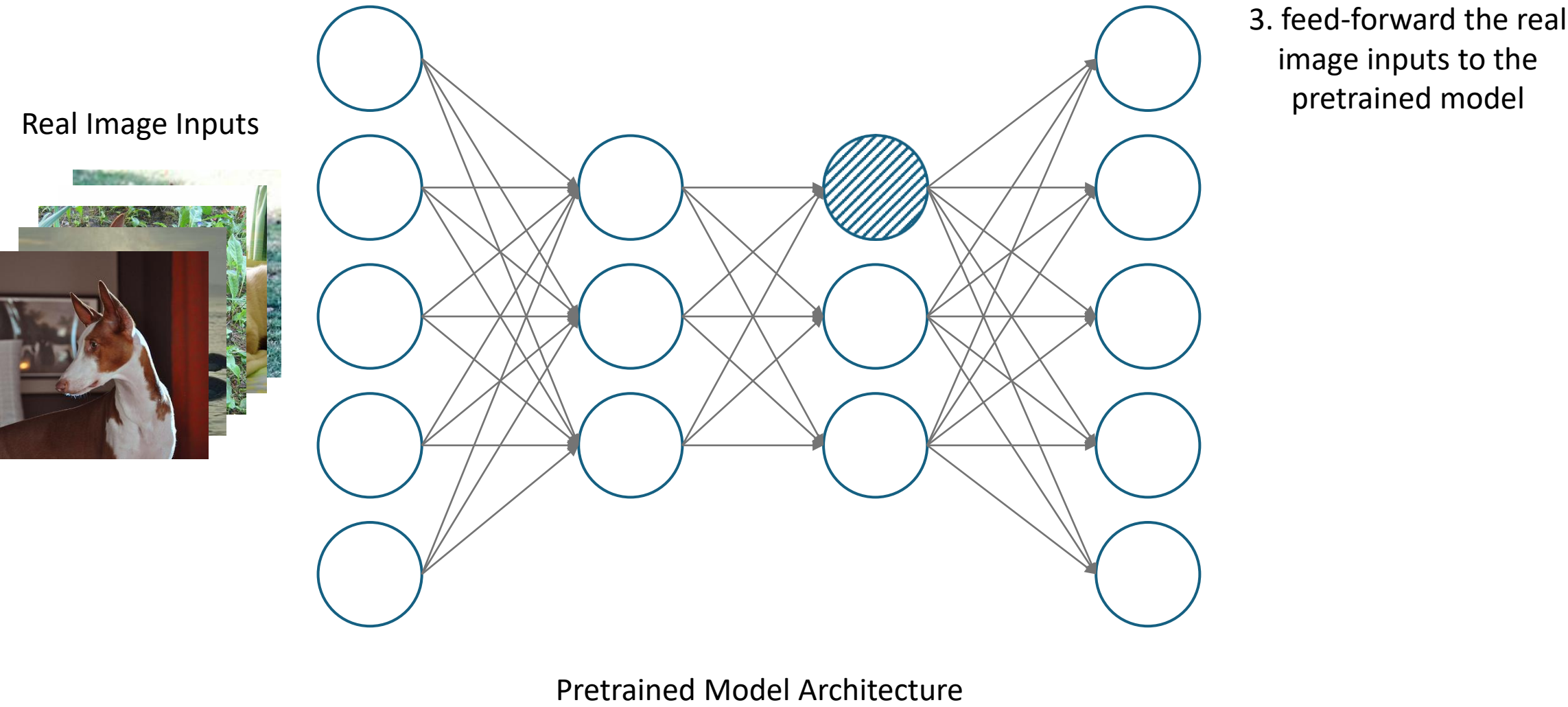
VITAL: Intermediate Neuron Visualization

2. find top-k images that activate the target neuron the most

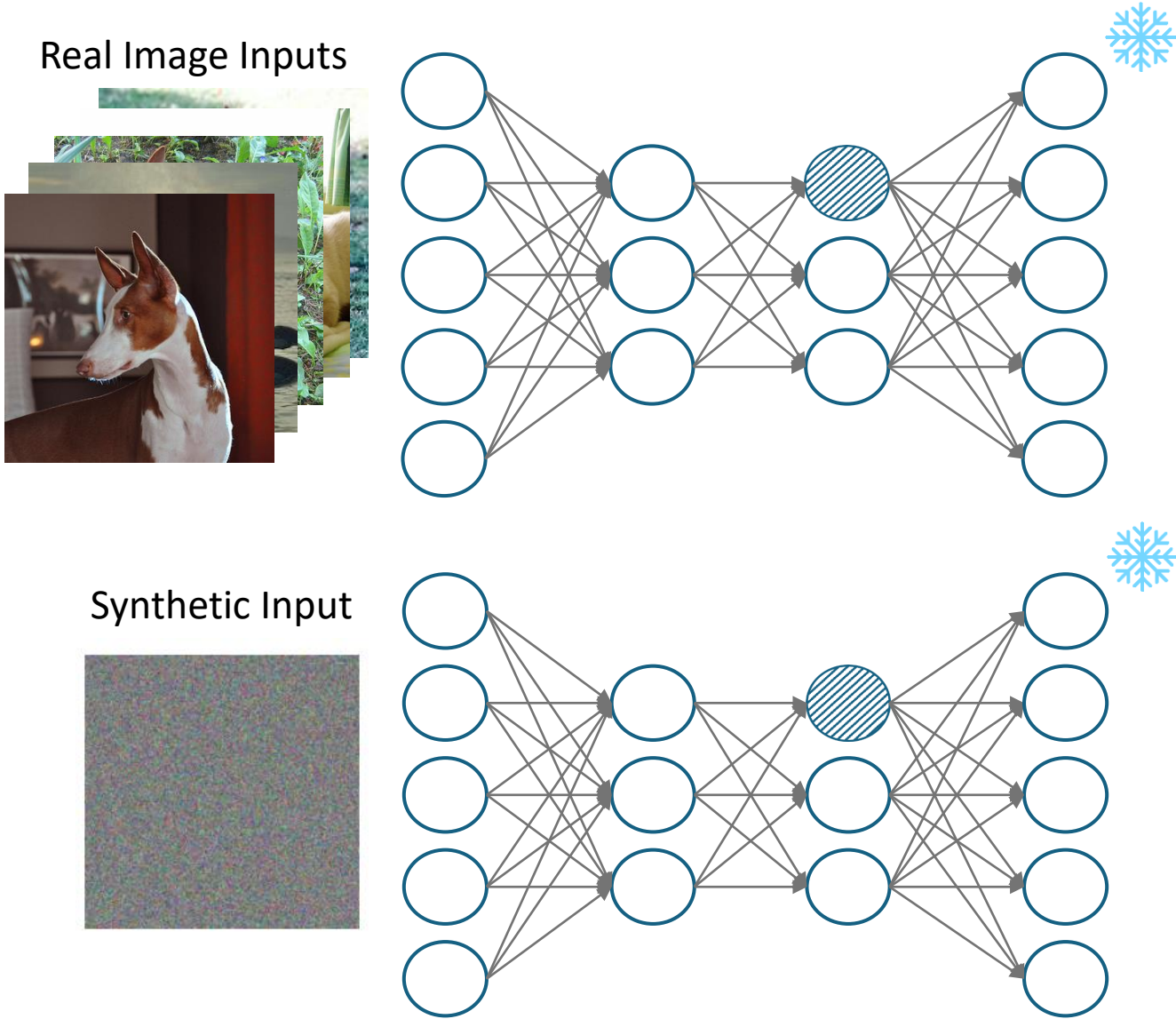


Pretrained Model Architecture

VITAL: Intermediate Neuron Visualization

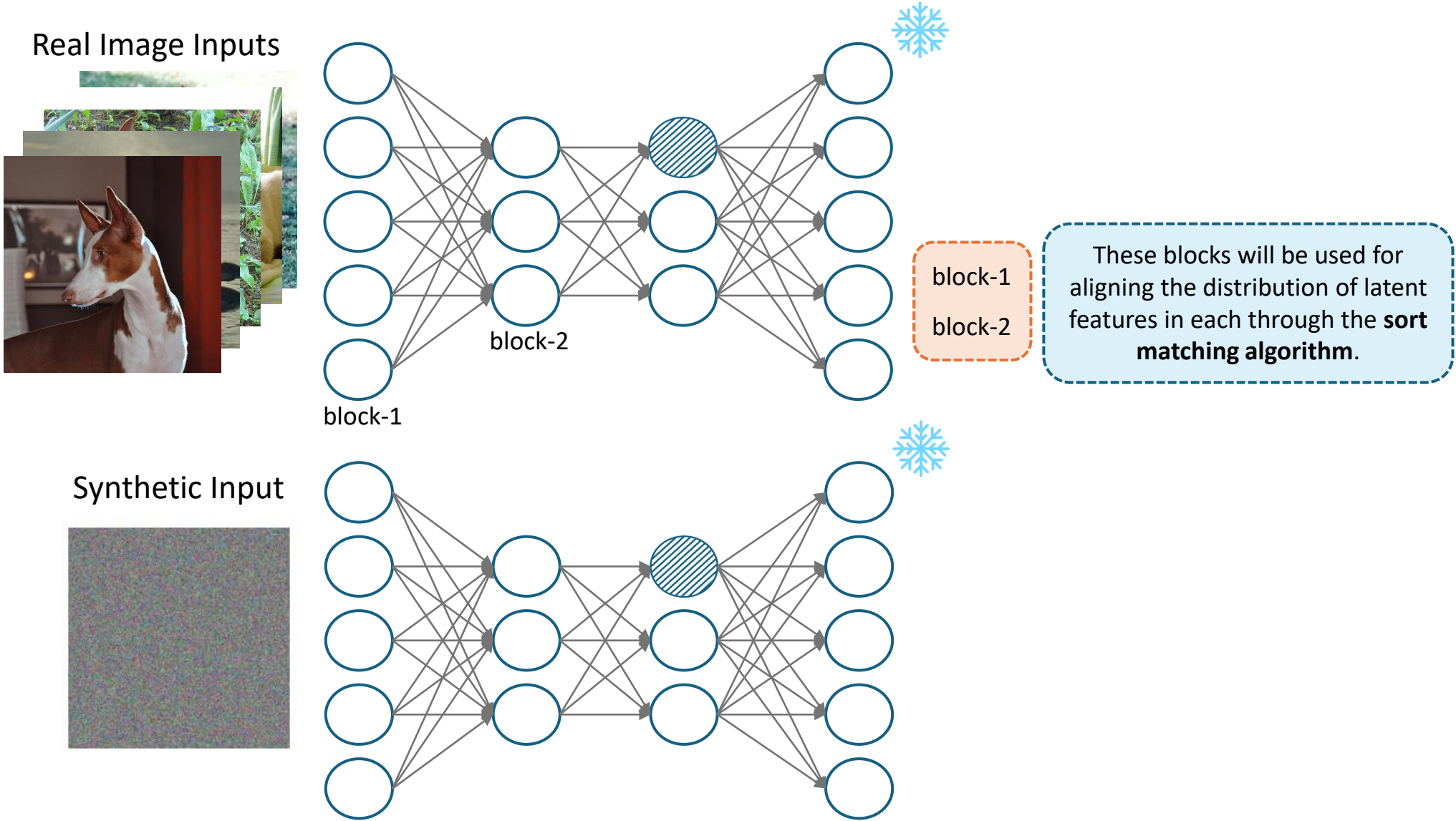


VITAL: Intermediate Neuron Visualization

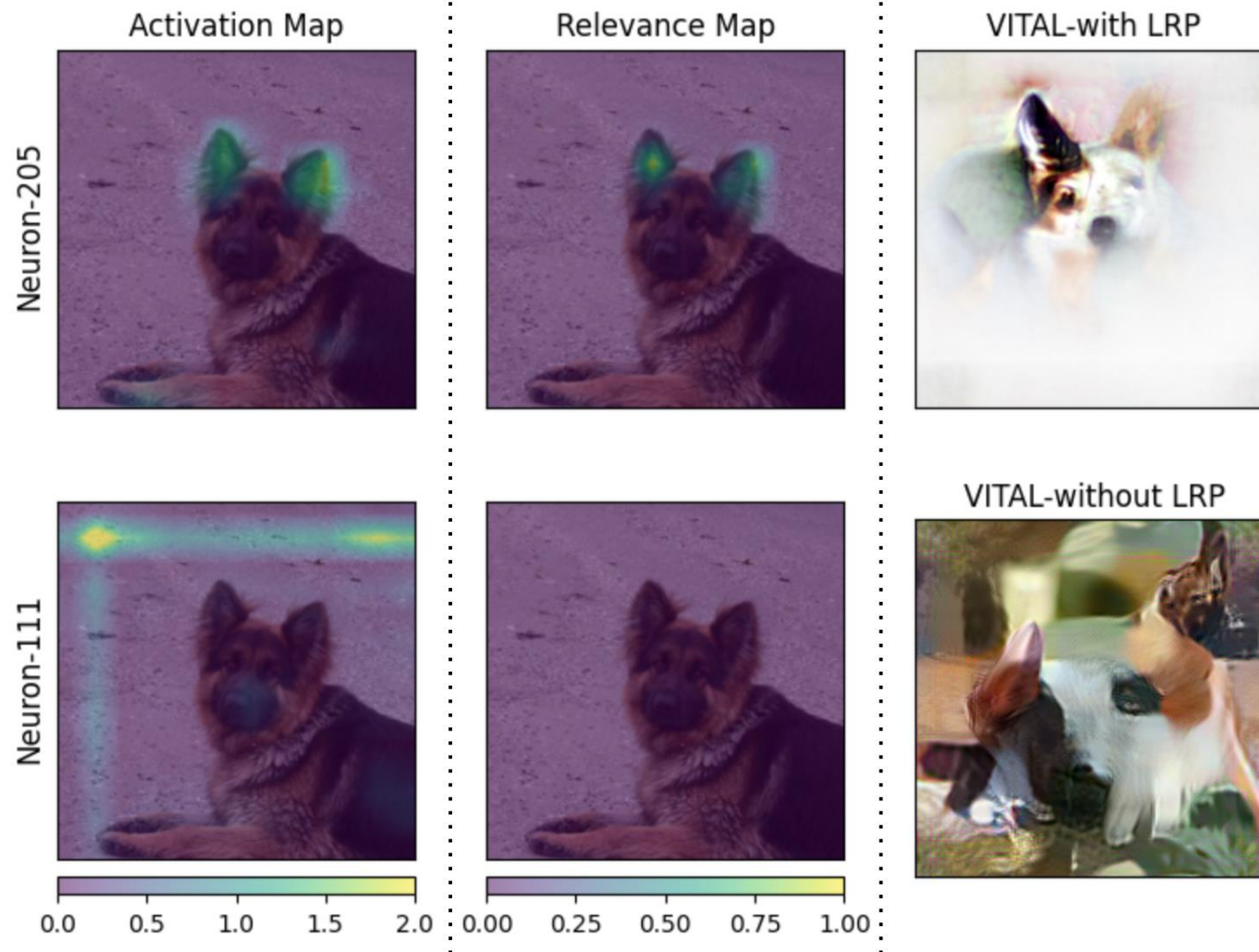


4. generate a noisy input and feed-forward to the same frozen model

VITAL: Intermediate Neuron Visualization



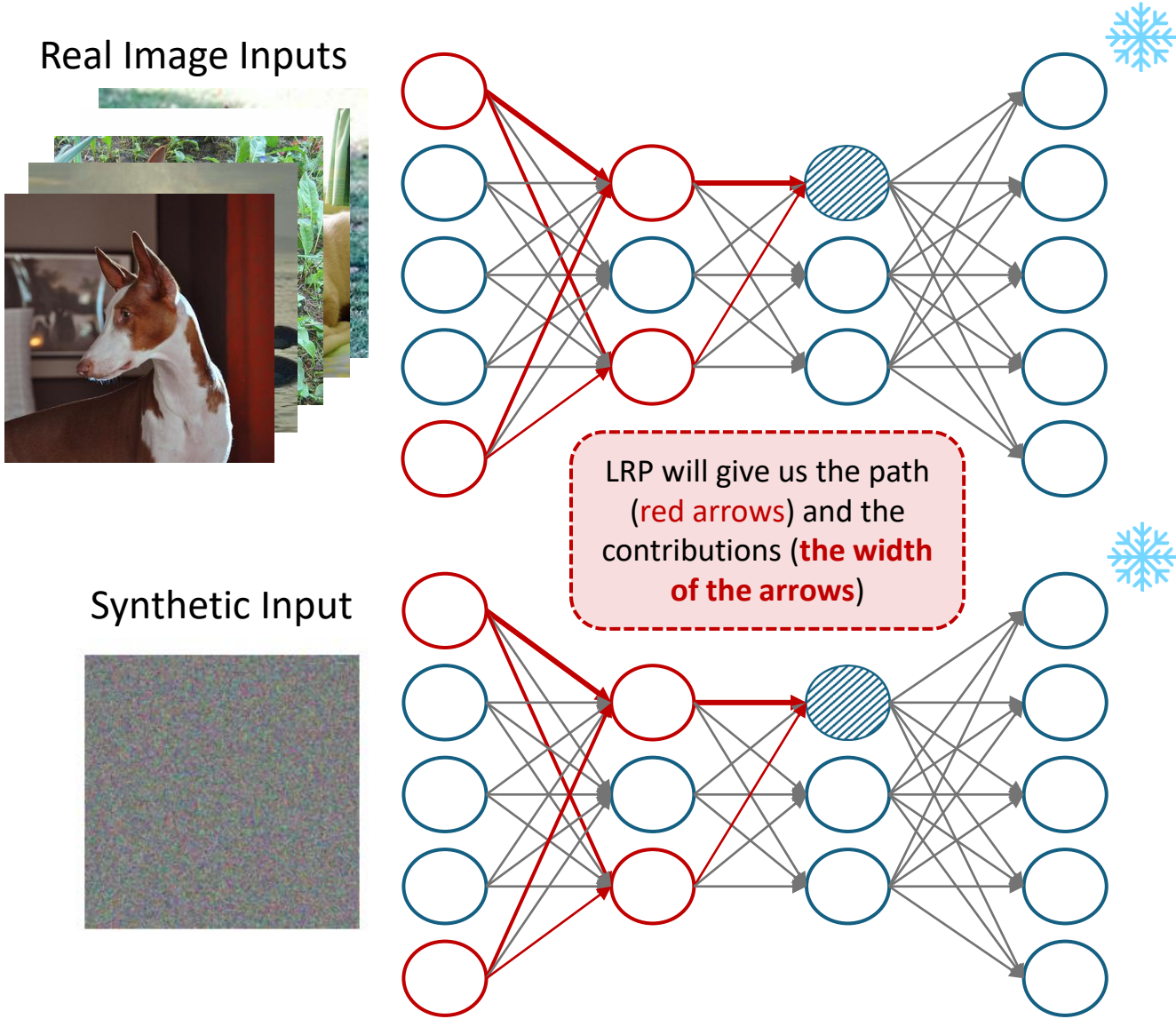
VITAL: Intermediate Neuron Visualization



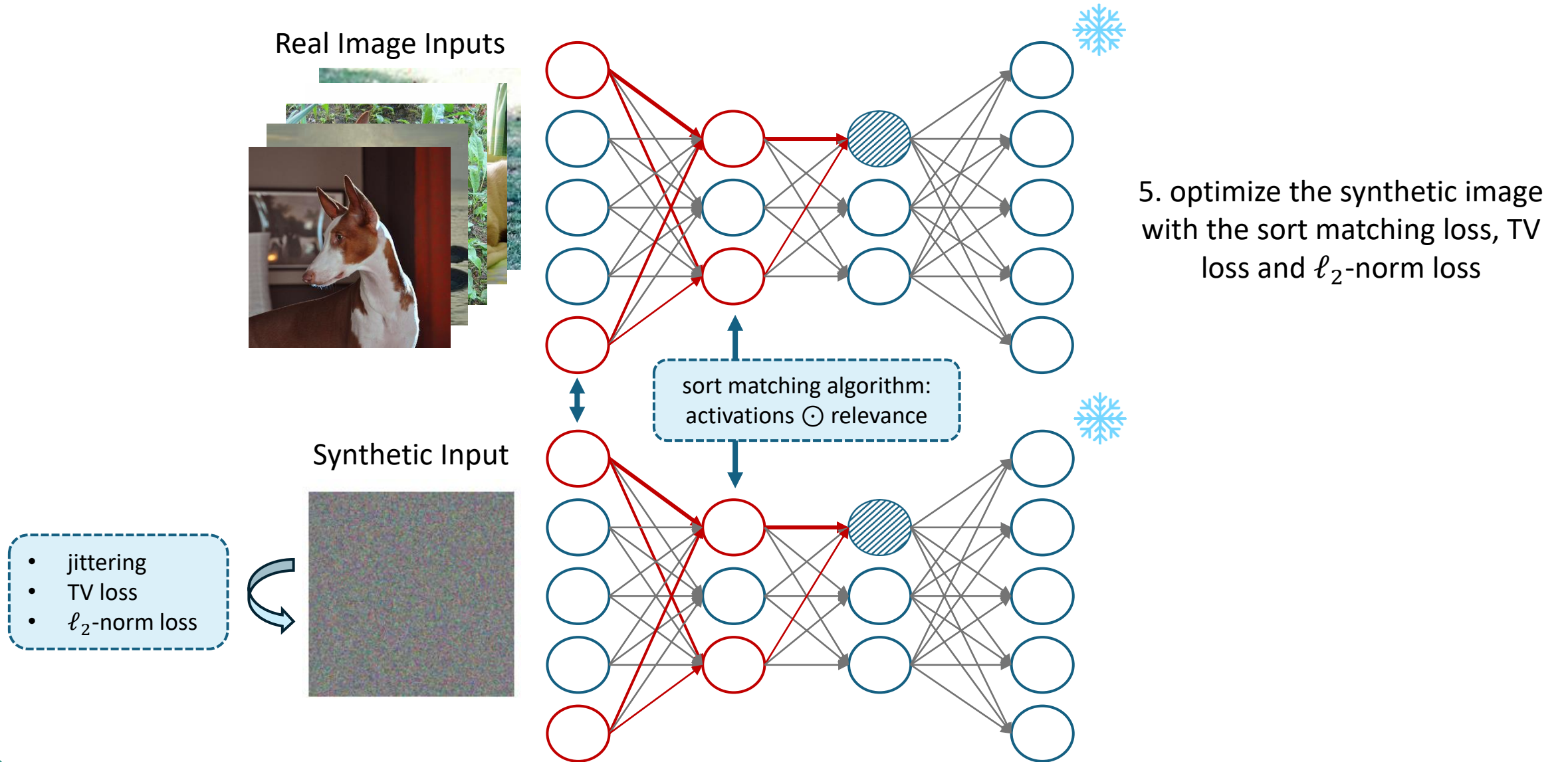
For intermediate neuron visualization, we need explicit supervision from the target neuron to find the regions for matching the distribution.

This supervision can be obtained through an attribution method (e.g., LRP).

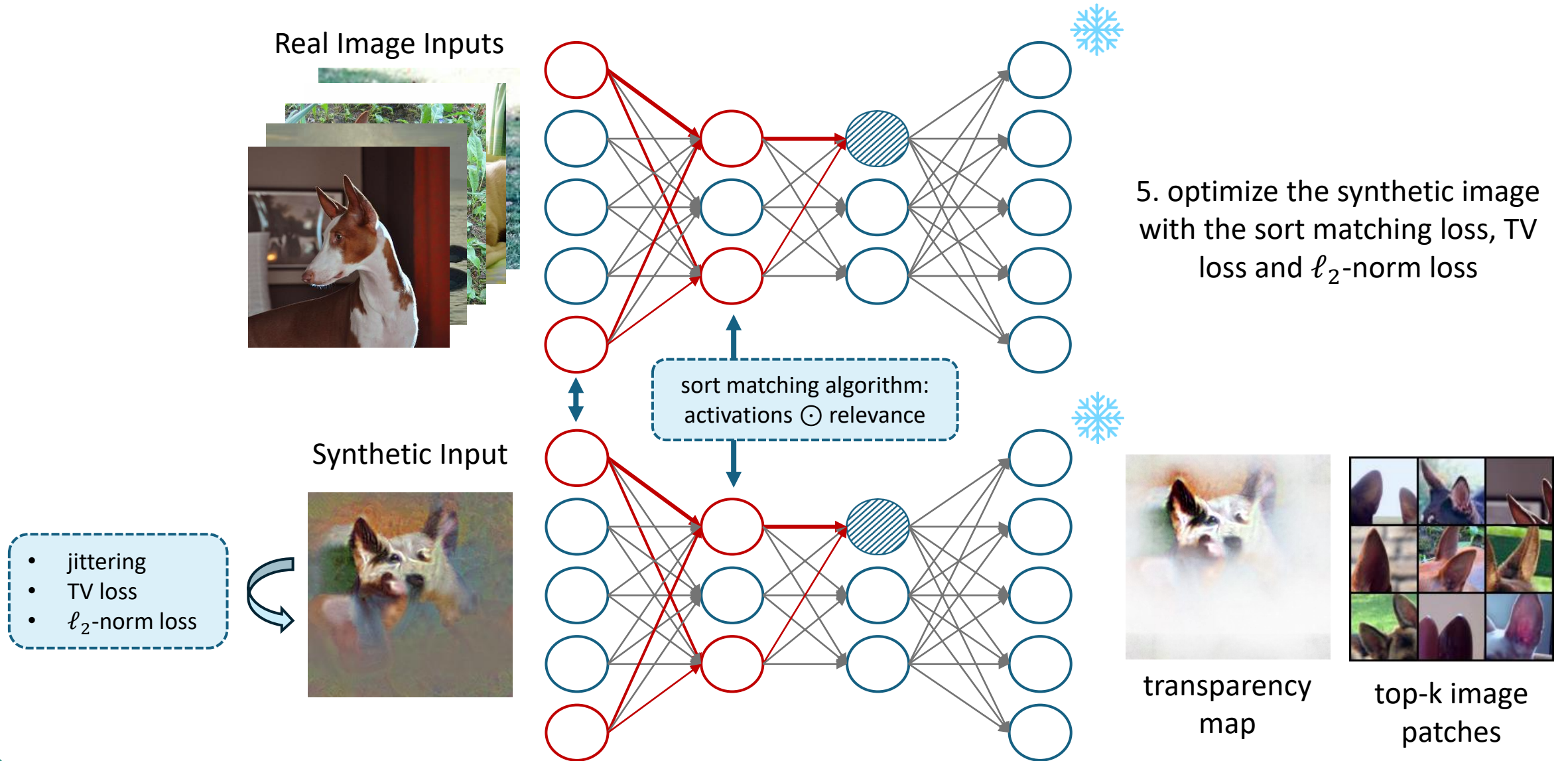
VITAL: Intermediate Neuron Visualization



VITAL: Intermediate Neuron Visualization

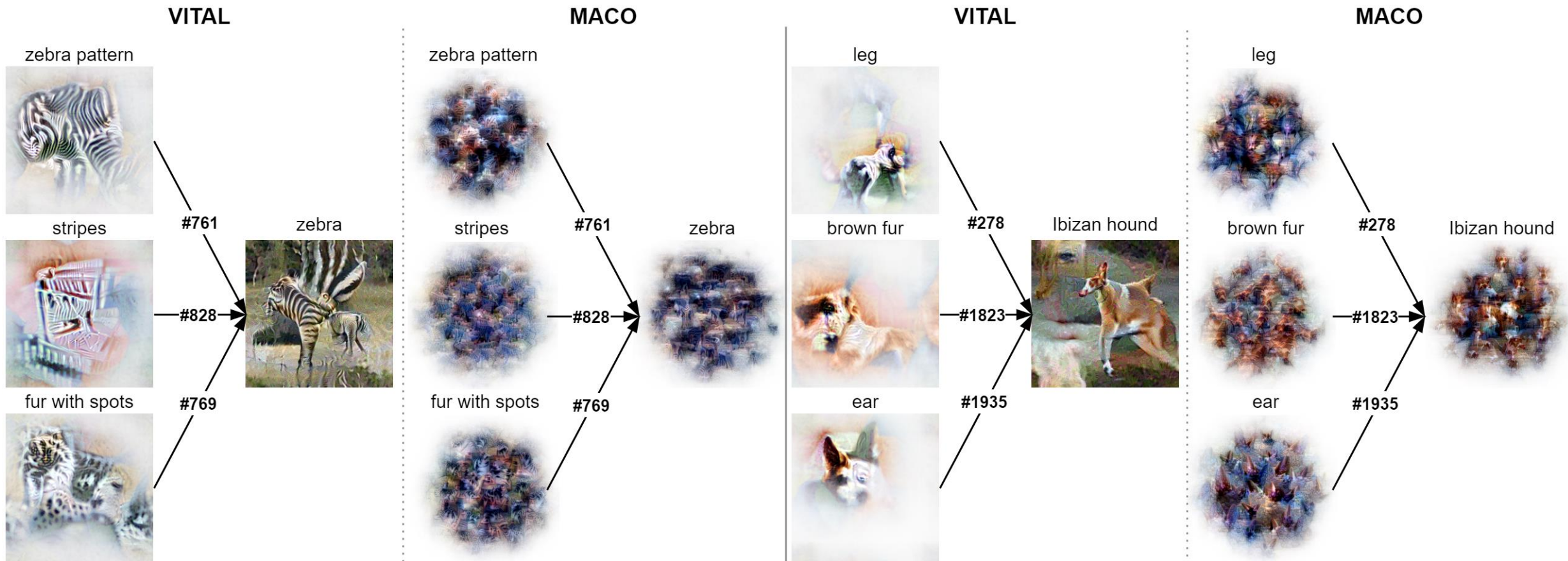


VITAL: Intermediate Neuron Visualization



VITAL: Intermediate Neuron Visualization

- Model tested: ResNet50



VITAL: Human User Study

- The user study aims to evaluate **interpretability / understandability** of the generated feature visualizations (FVs) from a **human perspective**.
 - complements the quantitative metrics (e.g. classification accuracy, FID, CLIP zero-shot)
- The study is organized into **3 tasks**, with increasing difficulty / less prior information.

Task	Description / Input	What participants do	Scoring / Measure
(1) <i>Class-level match</i>	A class name (single word) + a FV	Rate how well the FV matches that class	1 (worst) to 5 (best) scale
(2) <i>Inner neuron match</i>	Reference images that highly activate a hidden neuron + the FV for that neuron	Rate how well the FV corresponds to those images	1 (worst) to 5 (best) scale
(3) <i>Labeling / open description</i>	A FV only (no class name)	Write a word / short description for what the FV shows	Compare description to ground truth via embedding similarity (Universal Sentence Encoder)

VITAL: Human User Study

- The user study aims to evaluate **interpretability / understandability** of the generated feature visualizations (FVs) from a **human perspective**.
 - complements the quantitative metrics (e.g. classification accuracy, FID, CLIP zero-shot)
- The study is organized into **3 tasks**, with increasing difficulty / less prior information.

Task	Description / Input	What participants do	Scoring / Measure
(1) <i>Class-level match</i>	A class name (single word) + a FV	Rate how well the FV matches that class	1 (worst) to 5 (best) scale
(2) <i>Inner neuron match</i>	Reference images that highly activate a hidden neuron + the FV for that neuron	Rate how well the FV corresponds to those images	1 (worst) to 5 (best) scale
(3) <i>Labeling / open description</i>	A FV only (no class name)	Write a word / short description for what the FV shows	Compare description to ground truth via embedding similarity (Universal Sentence Encoder)

VITAL: Human User Study

- The user study aims to evaluate **interpretability / understandability** of the generated feature visualizations (FVs) from a **human perspective**.
 - complements the quantitative metrics (e.g. classification accuracy, FID, CLIP zero-shot)
- The study is organized into **3 tasks**, with increasing difficulty / less prior information.

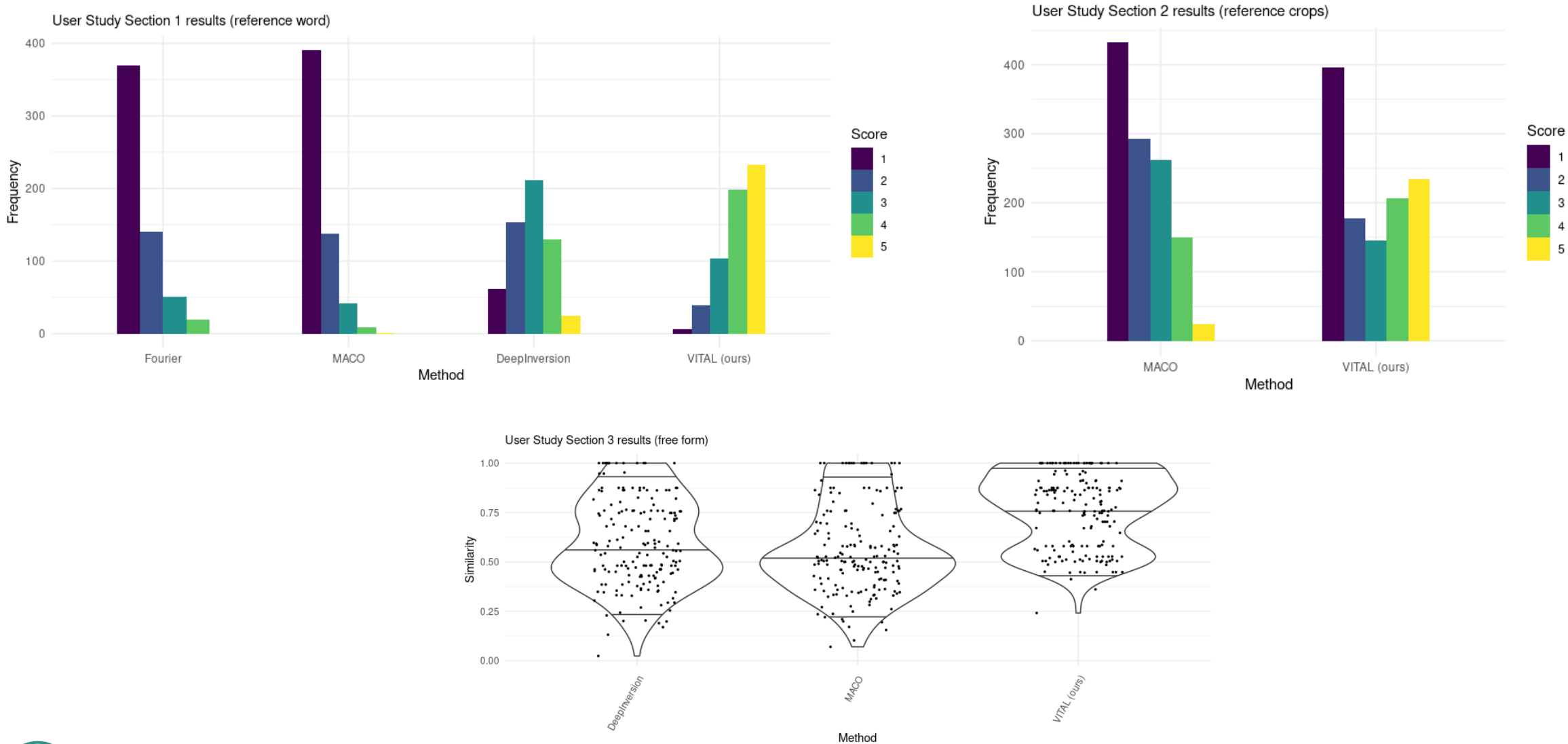
Task	Description / Input	What participants do	Scoring / Measure
(1) <i>Class-level match</i>	A class name (single word) + a FV	Rate how well the FV matches that class	1 (worst) to 5 (best) scale
(2) <i>Inner neuron match</i>	Reference images that highly activate a hidden neuron + the FV for that neuron	Rate how well the FV corresponds to those images	1 (worst) to 5 (best) scale
(3) <i>Labeling / open description</i>	A FV only (no class name)	Write a word / short description for what the FV shows	Compare description to ground truth via embedding similarity (Universal Sentence Encoder)

VITAL: Human User Study

- The user study aims to evaluate **interpretability / understandability** of the generated feature visualizations (FVs) from a **human perspective**.
 - complements the quantitative metrics (e.g. classification accuracy, FID, CLIP zero-shot)
- The study is organized into **3 tasks**, with increasing difficulty / less prior information.

Task	Description / Input	What participants do	Scoring / Measure
(1) <i>Class-level match</i>	A class name (single word) + a FV	Rate how well the FV matches that class	1 (worst) to 5 (best) scale
(2) <i>Inner neuron match</i>	Reference images that highly activate a hidden neuron + the FV for that neuron	Rate how well the FV corresponds to those images	1 (worst) to 5 (best) scale
(3) <i>Labeling / open description</i>	A FV only (no class name)	Write a word / short description for what the FV shows	Compare description to ground truth via embedding similarity (Universal Sentence Encoder)

VITAL: Human User Study



VITAL: Limitations and Future Work

- While our visualizations, like those of other methods, are not photo-realistic, they sometimes lack spatial coherence, particularly for non-rigid objects.
- Exhaustive benchmarks and well-defined evaluation metrics for inner neurons are actively discussed but still lacking.
- Future Works:
 - understanding networks in the medical domain
 - downstream applications in global representation learning
 - studying what neurons encode after knowledge transfer
 - how pruning affects the representations



Thanks for listening!

