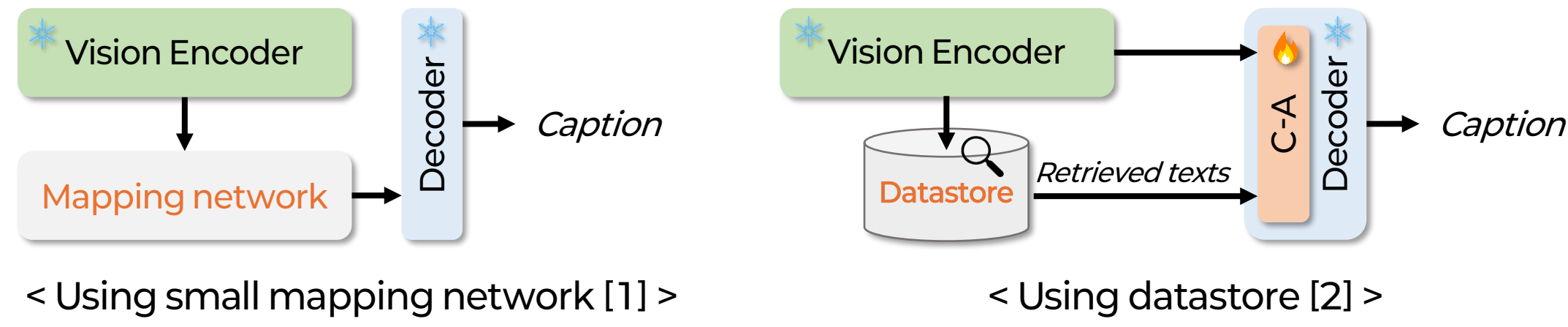


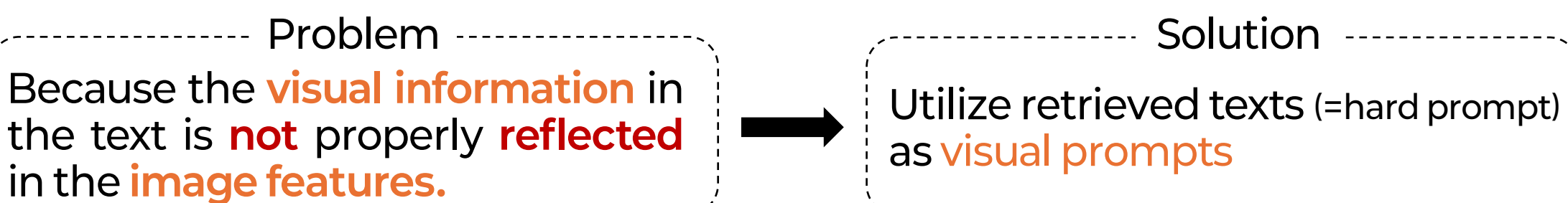
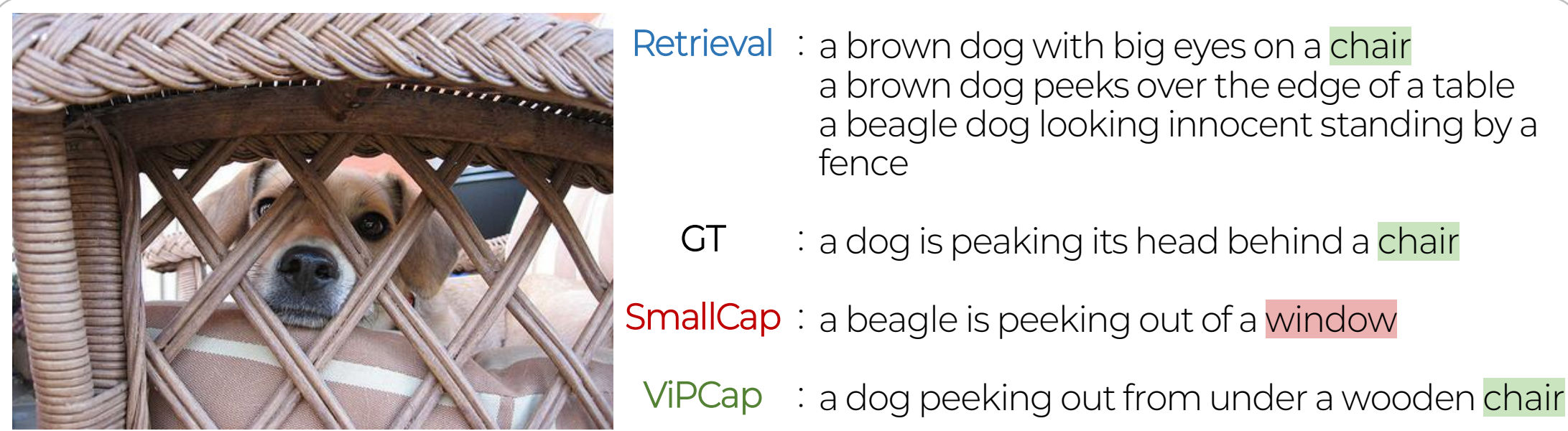
## Introduction

### Lightweight image captioning



### Problems in prior research

- Computational cost**
  - Despite of mapping network, high computational costs are still required.
- Neglect of rich image descriptions**
  - Prior works only use retrieval data as a text prompt, not a visual prompt.
- Rely on frozen CLIP image encoder**
  - Prior works rely solely on the frozen CLIP encoder.

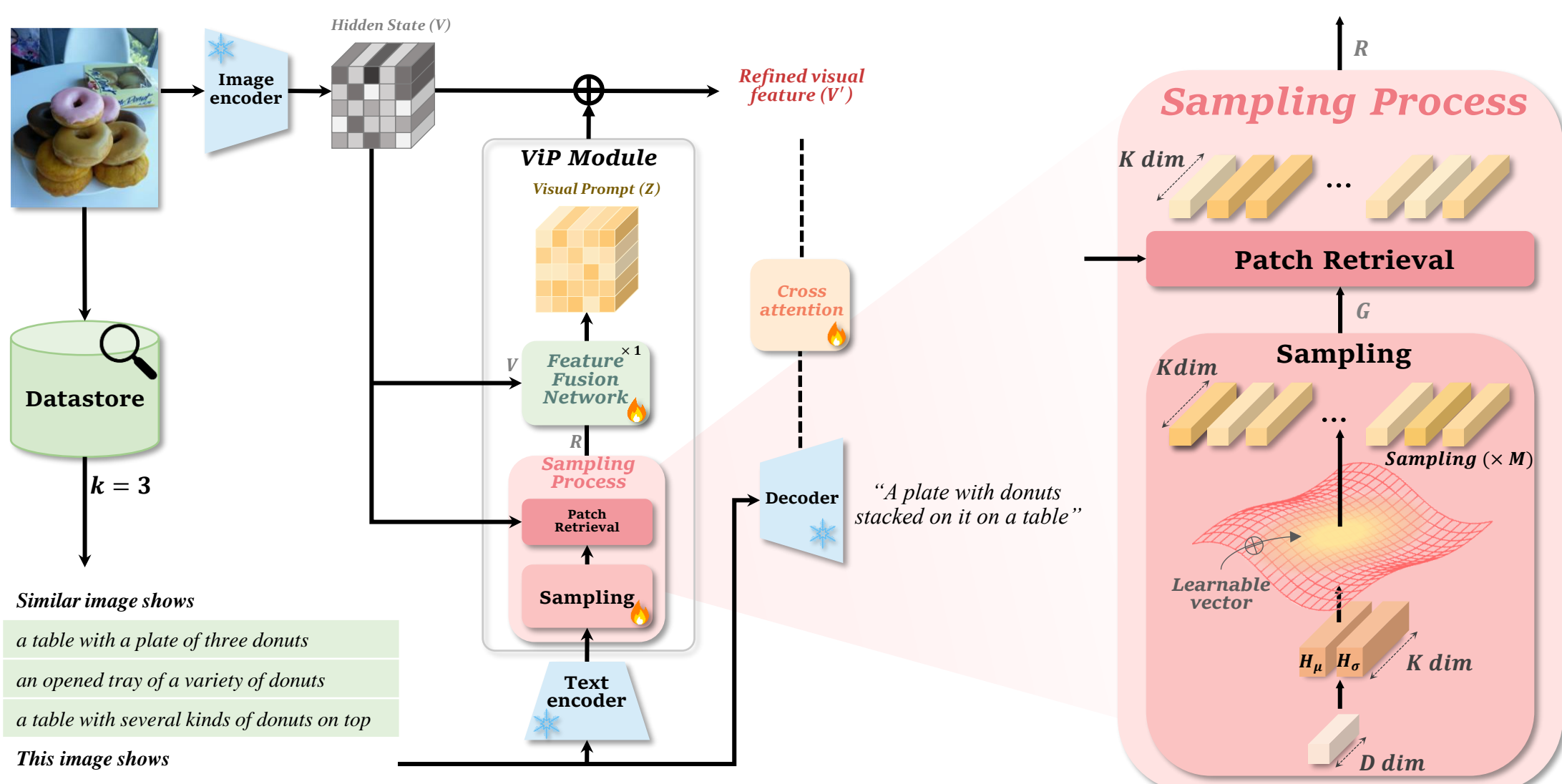


## Proposed Framework

### Contributions

- Generate **visual prompts** from **text prompts**.
- Align** sampled semantic text features with visual representations.
- Flexible integration with **various models and prompts (Plug-and-play)**.

### Overall framework



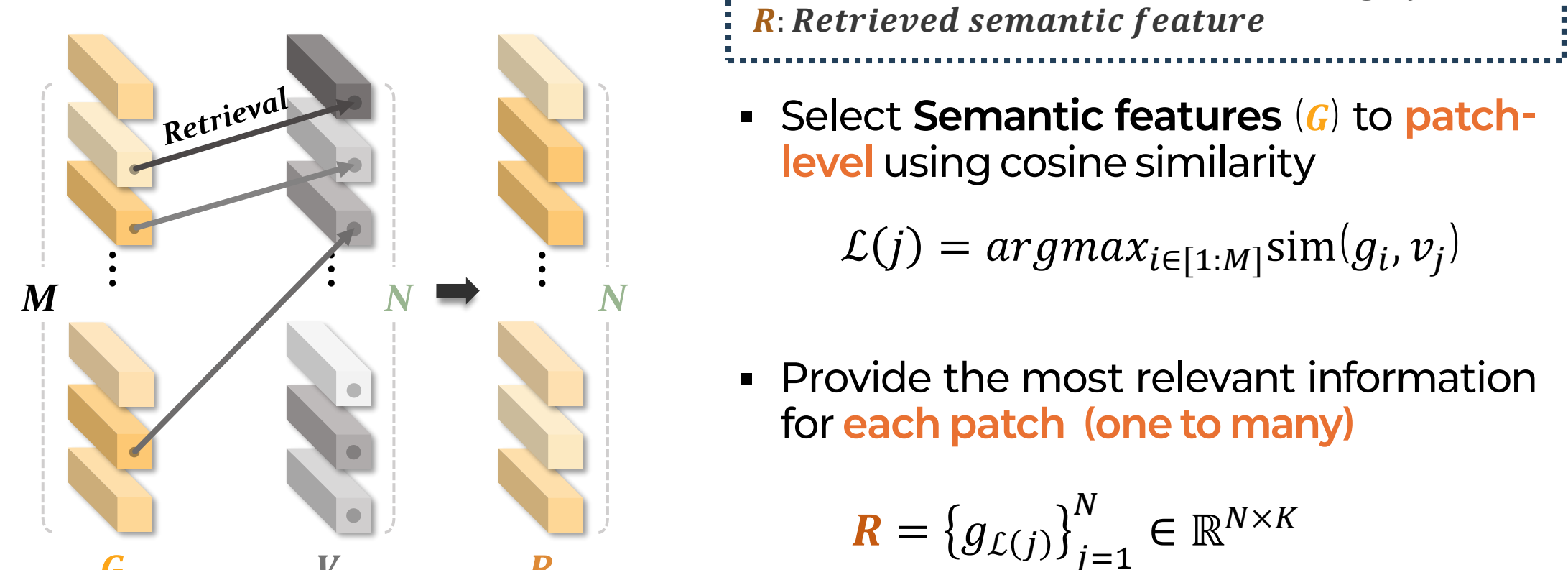
### Method

- Converts text prompts into semantic features using **learnable distribution**.
- Retrieves **semantic features** aligned with patch-level image features.
- Fuse these features with FFN to generate **visual prompts**.

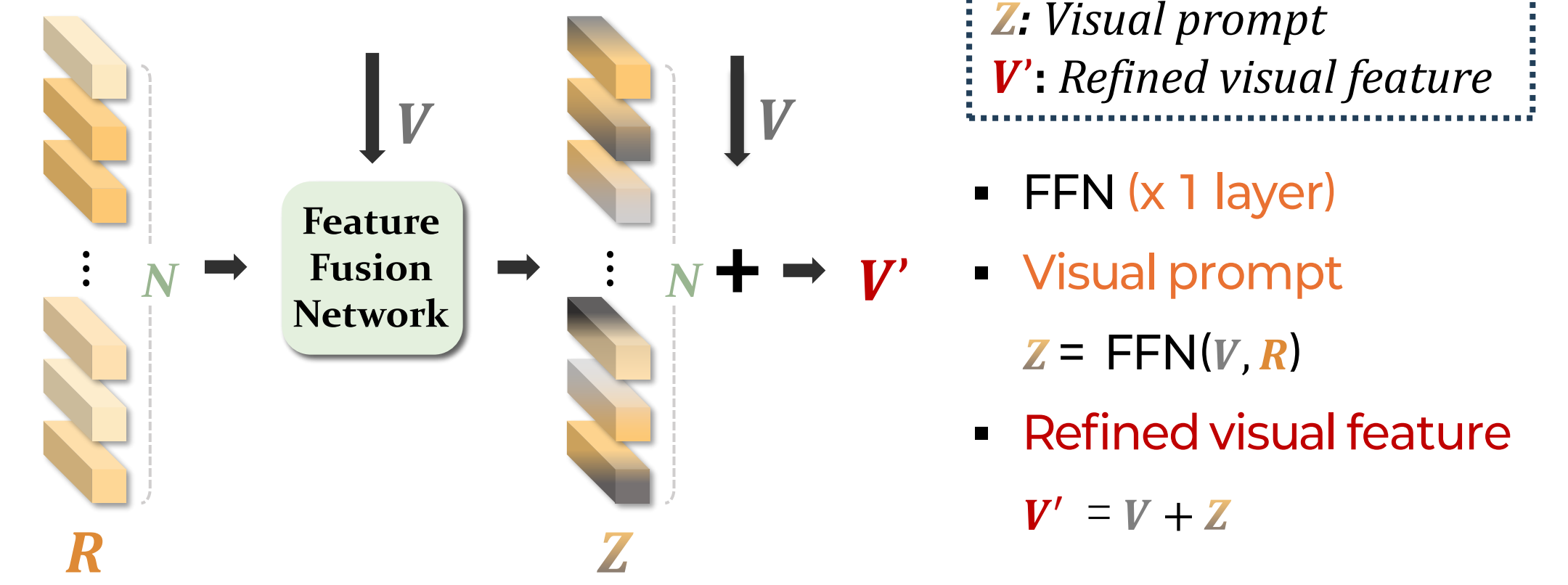
### Sampling process

- Sampling from Gaussian distribution
    - Capture for fine-grained visual information
  - Add learnable vectors
- $$\times M \rightarrow \begin{cases} \mu = H_\mu(T) + \alpha \cdot \omega, \\ \sigma = H_\sigma(T) \end{cases}$$

### Patch-retrieval module



### Feature Fusion Network (FFN)



## Experimental Results

### Evaluation on COCO, Flickr30k, Nocaps.

Method	Training Param $\theta$	COCO Test				Flickr30k Test		NoCaps Val			
		B@4	M	C	S	C	S	In	Near	Out	Entire
<b>Large scale training models</b>											
OSCAR <sub>Large</sub> (2020)	338M	37.4	30.7	127.8	23.5	-	-	78.8	78.9	77.4	78.6
LEMON <sub>Huge</sub> (2022)	675M	41.5	30.8	139.1	24.1	-	-	118.0	116.3	120.2	117.3
SimVLM <sub>Huge</sub> (2022a)	632M	40.6	33.7	143.3	25.4	-	-	113.7	110.9	115.2	112.2
BLIP2 <sub>ViT-g OPT<sub>T=7B</sub></sub> (2023a)	1.1B	43.7	-	145.8	-	-	-	123.0	117.8	123.4	119.7
CogVLM (2024)	1.5B	-	-	148.7	-	94.9	-	-	-	132.6	128.3
PaLM <sub>TS-XXL</sub> (2023)	1.6B	-	-	149.1	-	-	-	-	-	-	127.0
<b>Lightweight models</b>											
CaMEL (2022)	76M	<b>39.1</b>	<b>29.4</b>	<b>125.7</b>	<b>22.2</b>	-	-	-	-	-	-
I-Tuning <sub>Medium</sub> (2023)	44M	35.5	28.8	120.0	22.0	<b>72.3</b>	<b>19.0</b>	89.6	77.4	58.8	75.4
ClipCap (2021)	43M	33.5	27.5	113.1	21.1	-	-	84.9	66.8	49.1	65.8
I-Tuning <sub>Base</sub> (2023)	14M	34.8	28.3	116.7	21.8	61.5	16.9	83.9	70.3	48.1	67.8
SmallCap (2023)	7M	37.0	27.9	119.7	21.3	60.6	-	87.6	<u>78.6</u>	<u>68.9</u>	<u>77.9</u>
SmallCap <sub>d=16, Large</sub> (2023)	47M	37.2	28.3	121.8	21.5	-	-	-	-	-	-
ViPCap (Ours)	14M	<u>37.7</u>	28.6	<u>122.9</u>	21.9	<u>66.8</u>	<u>17.2</u>	<b>93.8</b>	<b>81.6</b>	<b>71.5</b>	<b>81.3</b>

→ ViPCap demonstrates **competitive performance**

### Evaluation on text-only captioning models.

Method	In-Domain								Cross-Domain											
	COCO				Flickr30k				COCO ⇒ Flickr30k				Flickr30k ⇒ COCO				COCO ⇒ NoCaps			
	B@4	M	C	S	B@4	M	C	S	B@4	M	C	S	B@4	M	C	S	In	Near	Out	Entire
CapDec (2022)	26.4	25.1	91.8	-	17.7	20.0	39.1	-	17.3	18.6	35.7	-	9.2	16.3	27.3	-	60.1	50.2	28.7	45.9
CapDec+ViP	27.0	25.6	94.2	18.8	18.6	20.1	44.4	14.4	15.7	18.0	35.8	11.8	9.5	16.3	30.7	9.2	60.2	50.9	33.7	47.8
$\Delta$	<b>0.6</b>	<b>0.5</b>	<b>2.4</b>	-	<b>0.9</b>	<b>0.1</b>	<b>5.3</b>	-	<b>-1.6</b>	<b>-0.6</b>	<b>0.1</b>	-	<b>0.3</b>	-	<b>3.4</b>	-	<b>0.1</b>	<b>0.7</b>	<b>5.0</b>	<b>1.9</b>
ViECap (2023)	27.2	24.8	92.9	18.2	21.4	20.1	47.9	13.6	17.4	18.0	38.4	11.2	12.6	19.3	54.2	12.5	61.1	64.3	65.0	66.2
ViECap+ViP	27.3	25.1	93.6	18.4	21.2	20.2	48.8	13.9	17.4	18.1	40.2	11.1	13.6	19.3	55.2	12.7	62.2	64.9	67.1	67.2
$\Delta$	<b>0.1</b>	<b>0.3</b>	<b>0.7</b>	<b>0.2</b>	<b>-0.2</b>	<b>0.1</b>	<b>0.9</b>	<b>0.3</b>	-	<b>0.1</b>	<b>1.8</b>	<b>-0.1</b>	<b>1.0</b>	-	<b>1.0</b>	<b>0.2</b>	<b>1.1</b>	<b>0.6</b>	<b>2.1</b>	<b>1.0</b>

→ Demonstrate its potential as **image feature**

### Ablation Studies

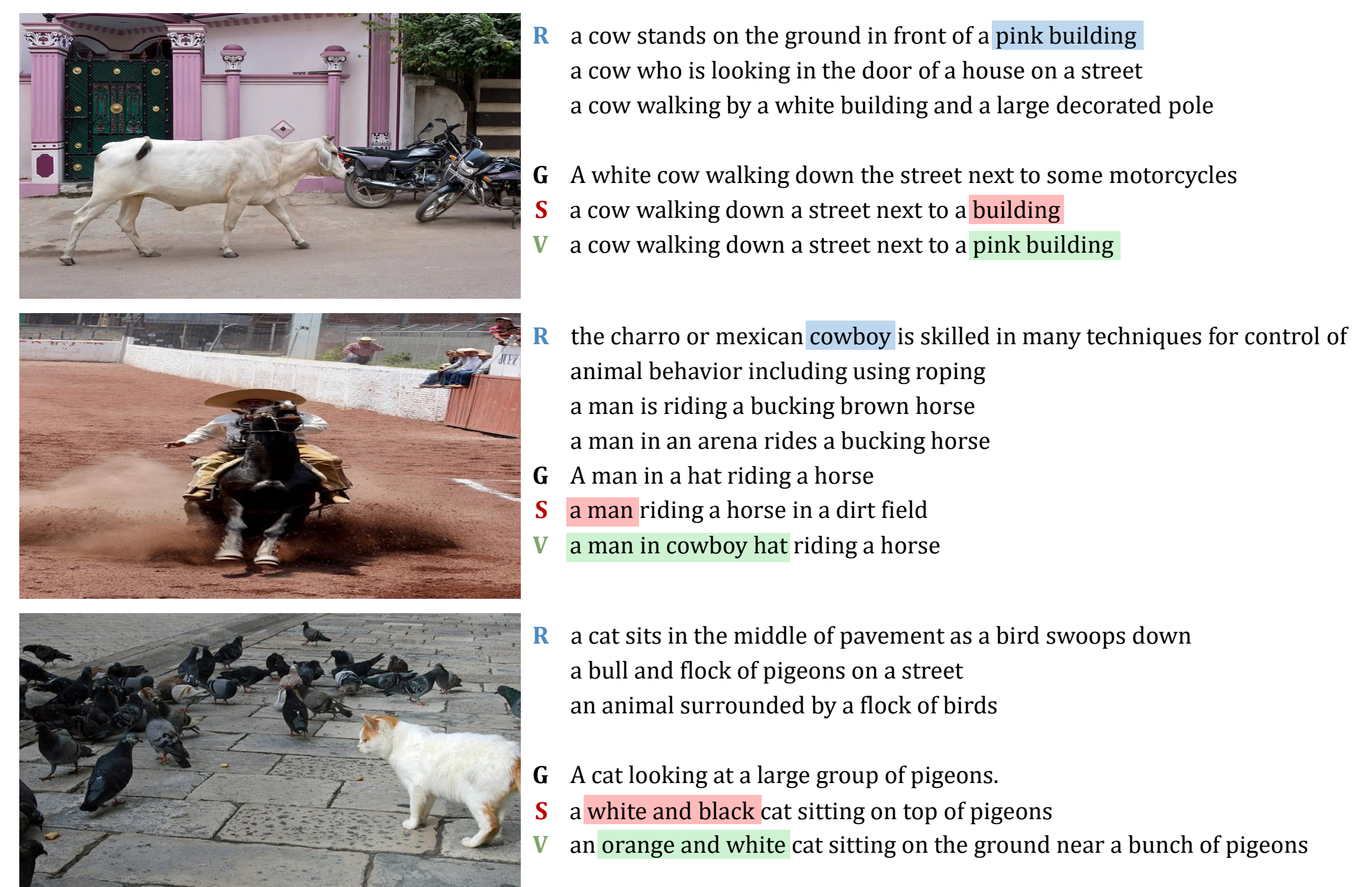
Prompt	ViP	
	×	✓
"This image shows"	111.1	116.0 ( <b>4.9</b> ↑)
Retrieval prompt	117.3	119.9 ( <b>2.6</b> ↑)

→ Achieves notable performance **without using retrieval prompts**

Method	Enc.	Dec.	ViP	Ret	CIDEr
ViPCap	ViT	OPT	-125M	✓	122.0
(Ours)	-B/32	XGLM		✓	122.5 ( <b>0.5</b> ↑)
				✓	116.8
				✓	121.2 ( <b>4.4</b> ↑)
EVCap	EVA-CLIP-g	Vicuna	-13B	×	140.1
				×	141.3 ( <b>1.2</b> ↑)

→ Outperformance based on **various models**

### Qualitative Results



## References

- Junnan Li, Dongxu Li et al., "BLIP-2: Bootstrapping Language-Image Pre-training with Frozen Image Encoders and Large Language Models", ICML 2023
- Ramos, Rita and Martins, Bruno and Elliott et al., "SmallCap: Lightweight Image Captioning Prompted With Retrieval Augmentation", CVPR 2023
- Li, Jiaxuan and Vo et al., "EVCap: Retrieval-Augmented Image Captioning with External Visual-Name Memory for Open-World Comprehension", CVPR 2024