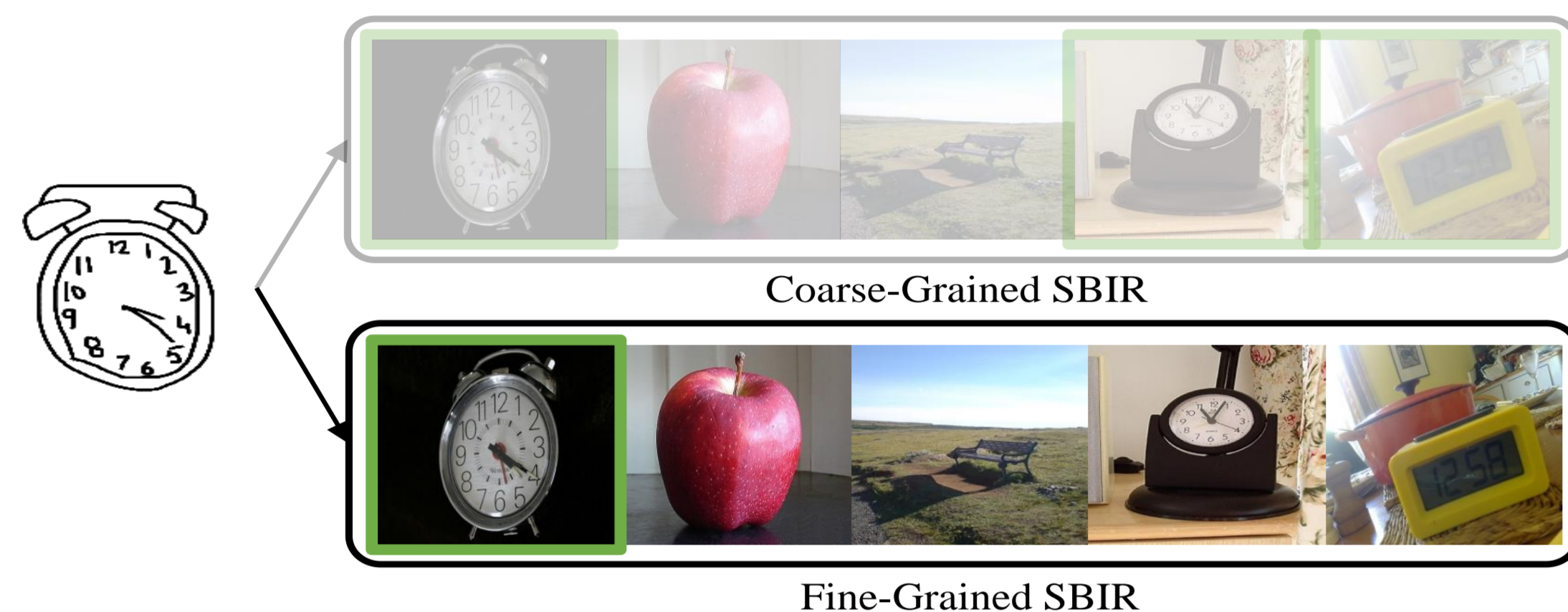


## TL;DR

- **Observation:** Modality-specific features can also be instance-discriminative.
- **Action:** Capture them via cross-attention. Bypass test-time overhead via knowledge distillation.
- **Results:** SOTA on benchmark FG-SBIR datasets.

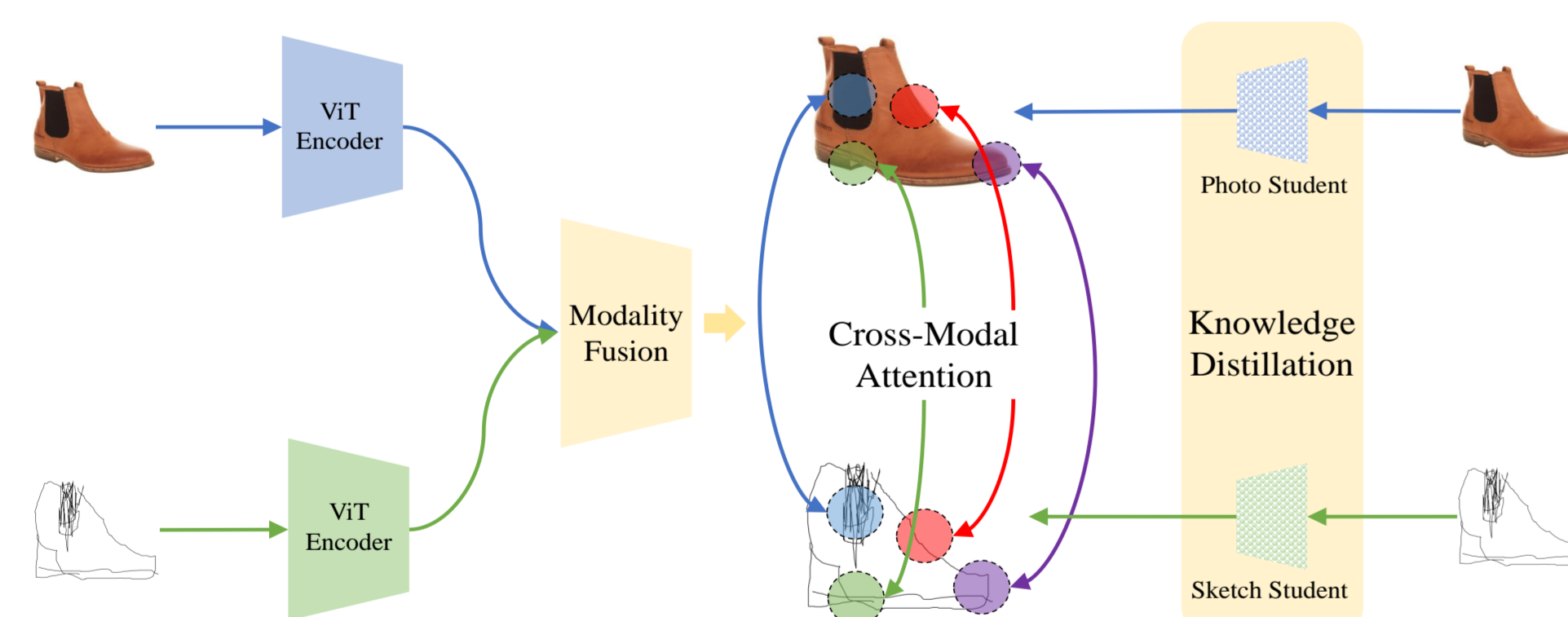
## Coarse-Grained SBIR vs. Fine-Grained SBIR



Fine-grained SBIR aims to retrieve a *specific* photo corresponding to a query sketch, unlike the coarse-grained setting, where all photos of the corresponding class are returned.

## Types of Cross-Modal Features

1. **Modality-shared** (useful for retrieval).
2. **Modality-specific**, that represents a *shared underlying concept*, but manifests differently in the two modalities (useful for retrieval).
3. **Modality-specific**, that does not represent a shared concept (**not useful for retrieval**).



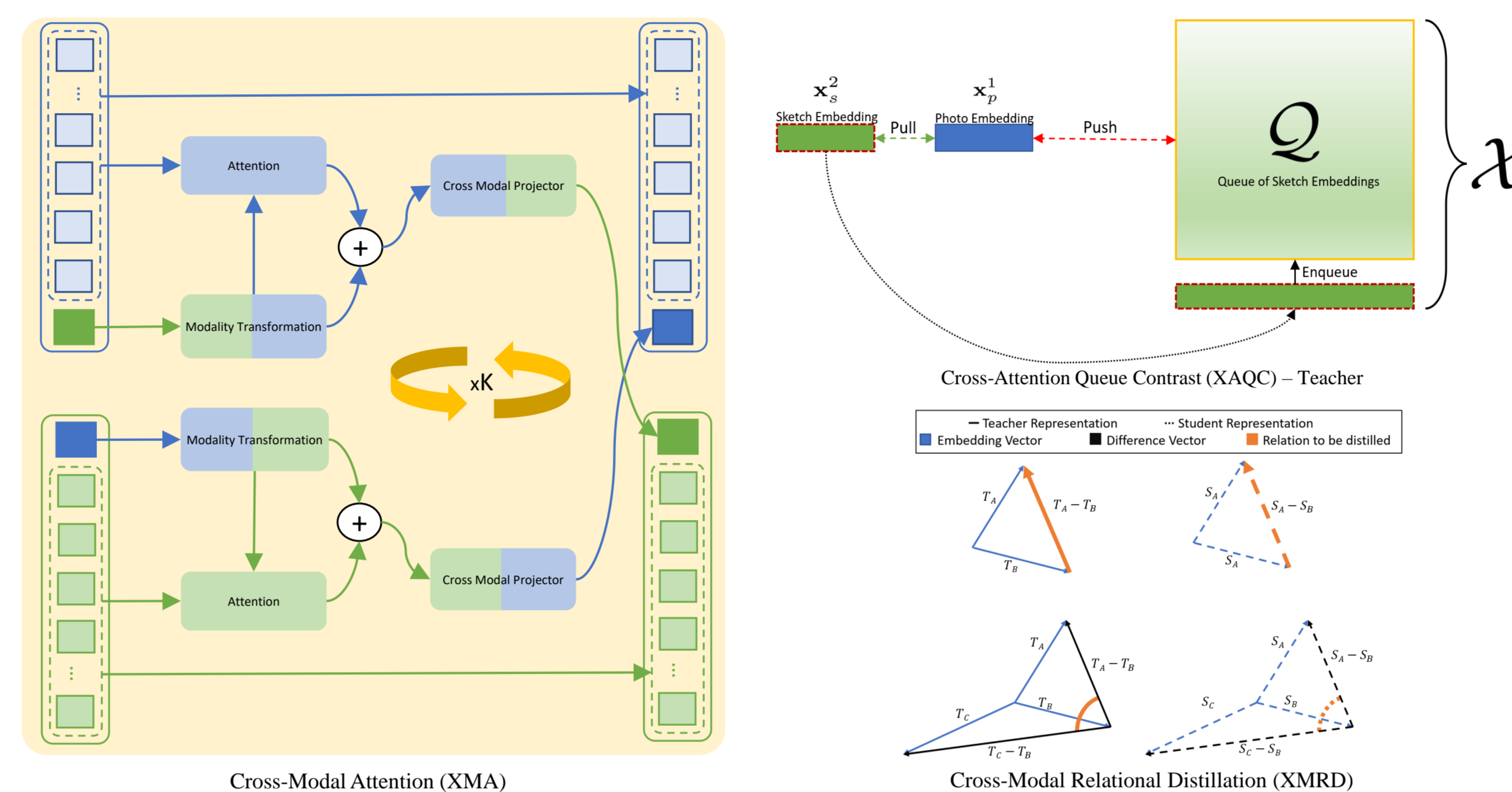
Color contrast may be depicted in sketches via scribbles (area under blue circles).

We simultaneously preserve (1) and (2) while discarding (3) via **Cross-Modal Attention (XMA)**. We bypass the test-time computational overhead of XMA via **Cross-Modal Knowledge Distillation**.

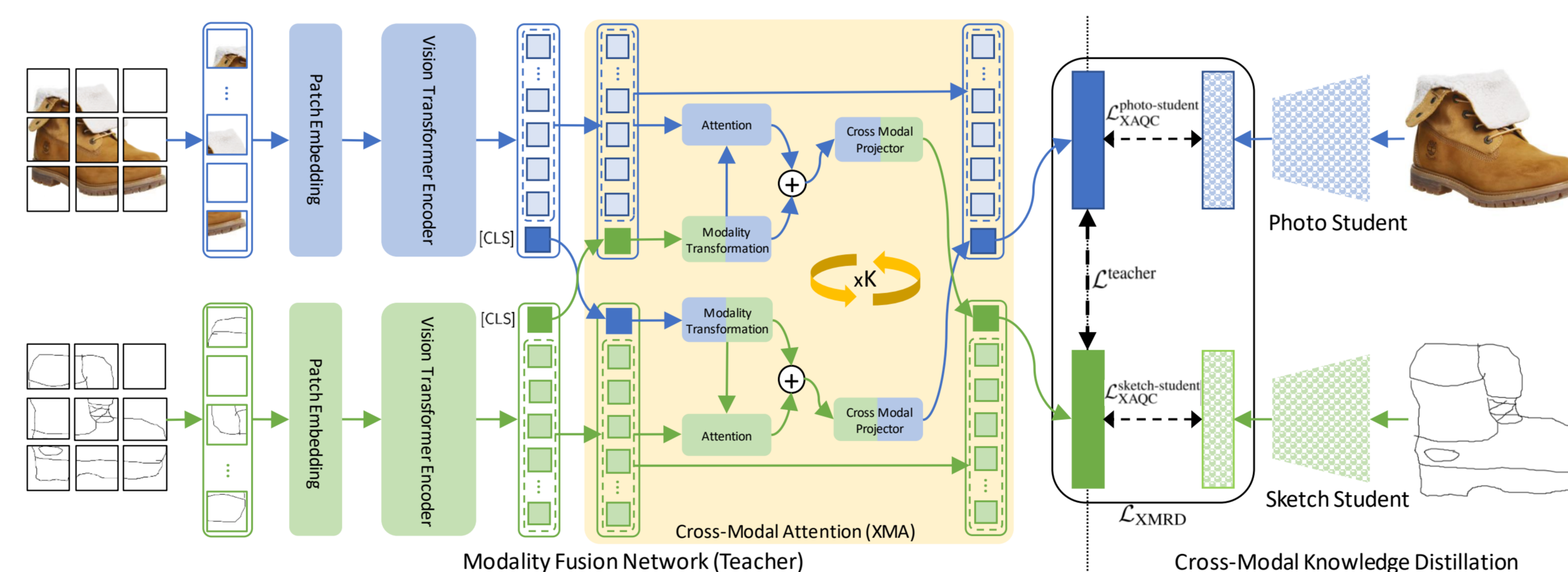
## Contributions

- Leveraging **modality-specific, instance-discriminative** features for FG-SBIR.
- Cross-Modal Knowledge Distillation for **fast retrieval without test-time cross-attention** [1].
- **State-of-the-art results** on benchmark FG-SBIR datasets.

## The XModalViT Framework



XMA → Modality-Fusion | XAQC → Ranking, XModal KD | XMRD → Relational Invariants



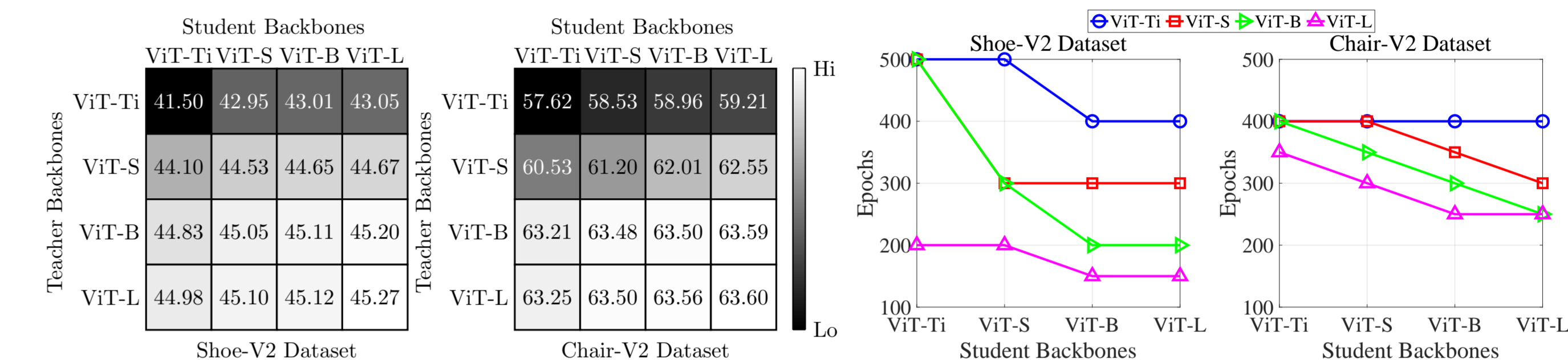
XModalViT: Modality Fusion (XMA) followed by disentanglement (XMKD) for fast retrieval.

## Comparison with SOTA

Method	Shoe-V2		Chair-V2		Sketchy	
	acc@1	acc@10	acc@1	acc@10	acc@1	acc@10
Yu <i>et al.</i> , CVPR '16	28.71	71.56	47.65	84.24	54.27	-
Yang <i>et al.</i> , ICCV '21	32.33	79.63	52.89	94.88	25.87	-
Sain <i>et al.</i> , CVPR '21	36.47	81.83	62.86	91.14	37.10	-
Bhunia <i>et al.</i> , CVPR '21	39.10	87.50	62.20	90.80	50.14	-
Chowdhury <i>et al.</i> , CVPR '22	39.90	82.90	-	-	40.16	92.00
Bhunia <i>et al.</i> , CVPR '22	43.70	-	<b>64.80</b>	-	46.20	96.49
<b>Ours (XModalViT)</b>	<b>45.05</b>	<b>90.23</b>	<b>63.48</b>	<b>95.02</b>	<b>56.15</b>	<b>96.86</b>

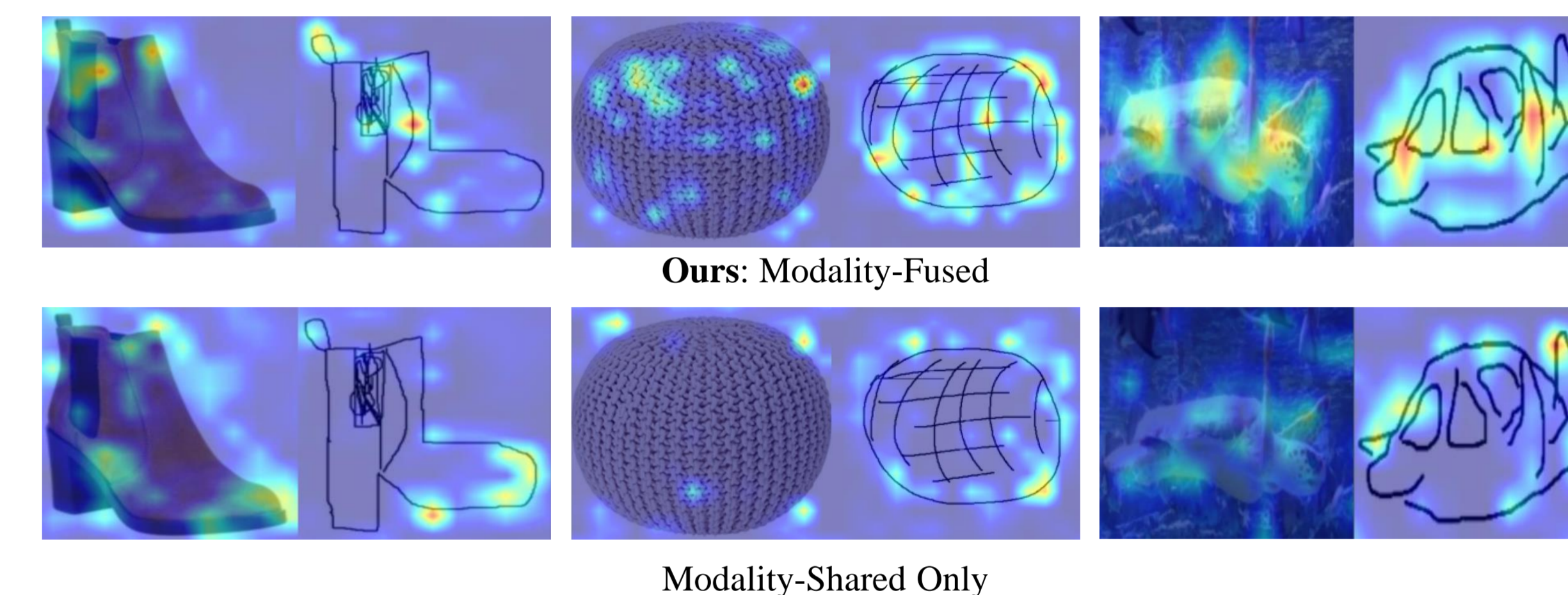
Surpasses **SOTA** on both **single-class** (Shoe-V2 & Chair-V2) & **multi-class** (Sketchy) benchmarks.

## Effect of Varying Backbone Sizes



(Left) Variation in acc@1 and (right) student convergence time (# epochs) with encoder size.

## Qualitative Results



Our method attends to **modality-specific features** like scribbles (Shoe), grids/meshes (Chair), or primitive shapes like square/triangle (Turtle), that are **instance discriminative**.



Sample **top-10** retrieval results in order obtained using our method (green → ground-truth).

## References

- [1] Chun-Fu (Richard) Chen, Quanfu Fan, and Rameswar Panda. CrossViT: Cross-Attention Multi-Scale Vision Transformer for Image Classification. In ICCV, 2021.
- [2] Patsorn Sangkloy, Nathan Burnell, Cusuh Ham, and James Hays. The sketchy database: Learning to retrieve badly drawn bunnies. ACM SIGGRAPH, 2016.