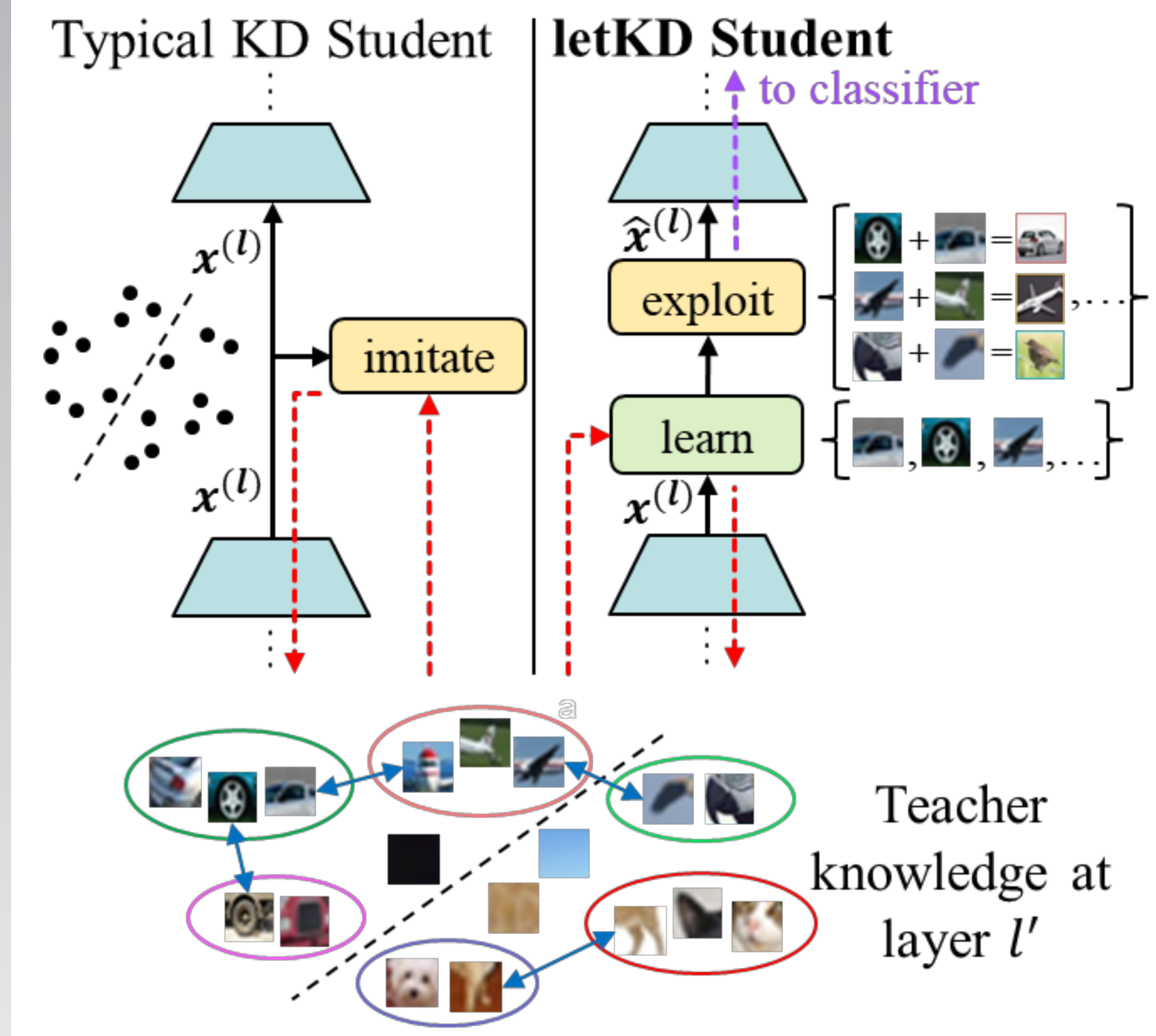


PROBLEM

Typical KD methods use regularization while pushing the student to **imitate** the feature geometry of the teacher.



→ : teacher knowledge propagation
↔ : inter-class relationships captured by the teacher

Considering the *architectural* differences in between, **forcing** the student to imitate the teacher's responses would be demanding, especially for the intermediate layers.

METHOD: FORMULATION

Key Idea: Learn the semantic entities that the teacher finds useful and **exploit** them in feature transform, enabling us to feed forward the knowledge during inference as well.

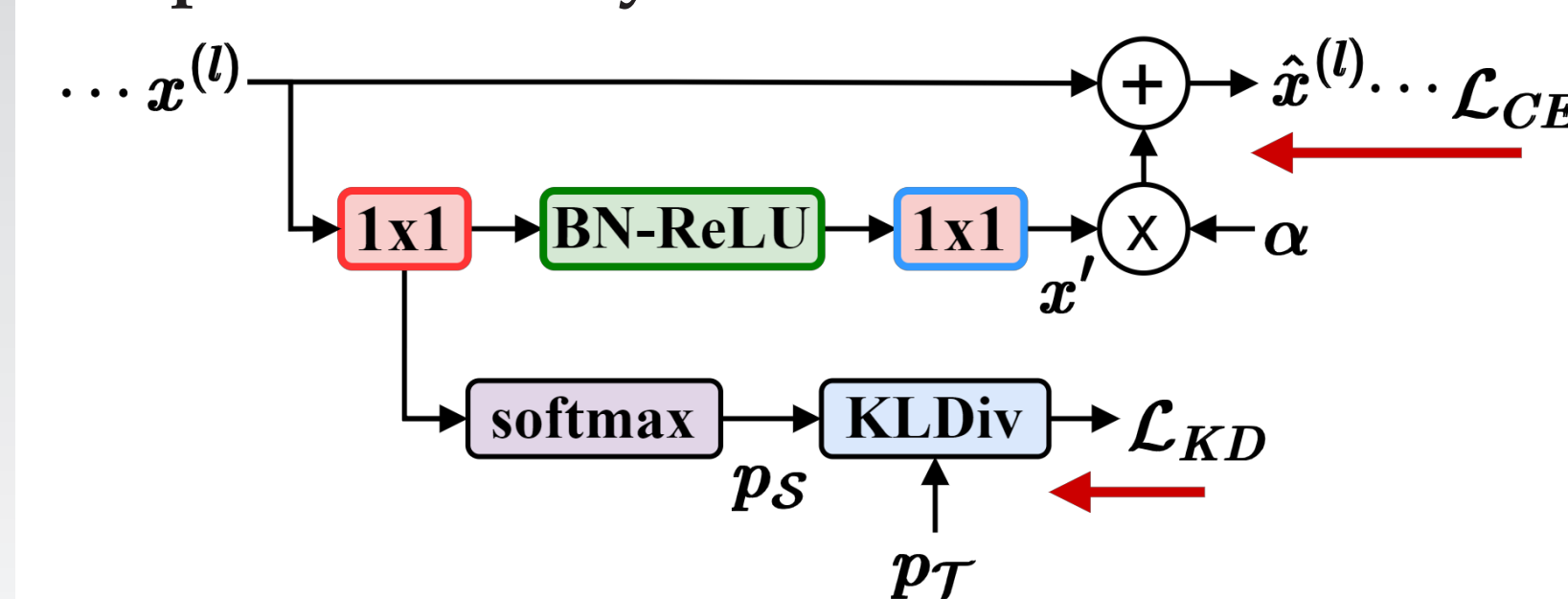
Given a set of matching kernels ω_k and features x_i at spatial location i , **we define** feature embedding by template matching as:

$$p_i = \arg \max_{p, q \geq 0} q \mu + \sum_k p_k \omega_k^T x_i$$

learned through the teacher
 $x'_i = \sum_k p_{ki} \nu_k$

Solver (see paper for details): **1x1-BN-ReLU-1x1** is equivalent to feature embedding by template matching.

Proposed KD Layer:

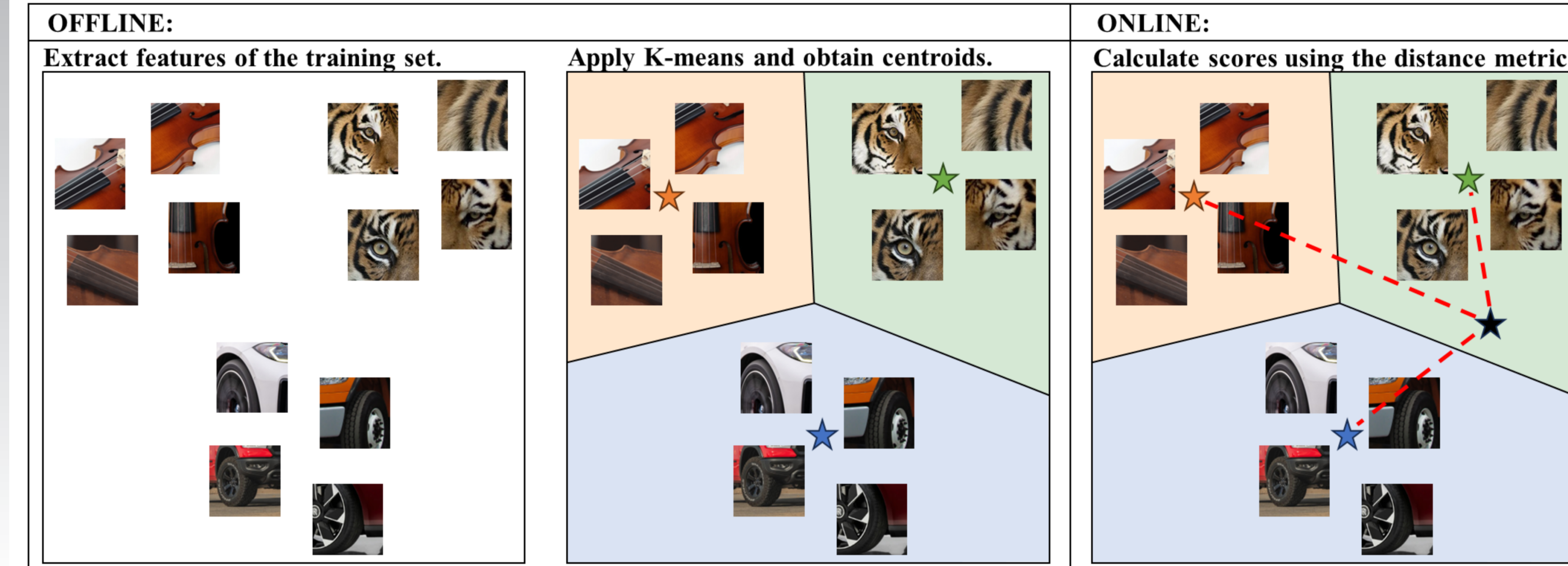


conv. → : learning feedback
 p_S / p_T : student's/teacher's soft predictions

PENULTIMATE LAYER TEACHER SUPERVISION

Viewing each pixel in a feature map as a *semantic entity*, **we propose** to employ *K-means* to the teacher's features to obtain fine-grained labels for these entities.

$$p_{\mathcal{T}}(\star) = \text{softmax} \left(- \frac{\|c - \star\|_2^2}{\tau} \right)_{c \in \{\star, \star, \star\}}$$



INTERMEDIATE LAYER TEACHER SUPERVISION

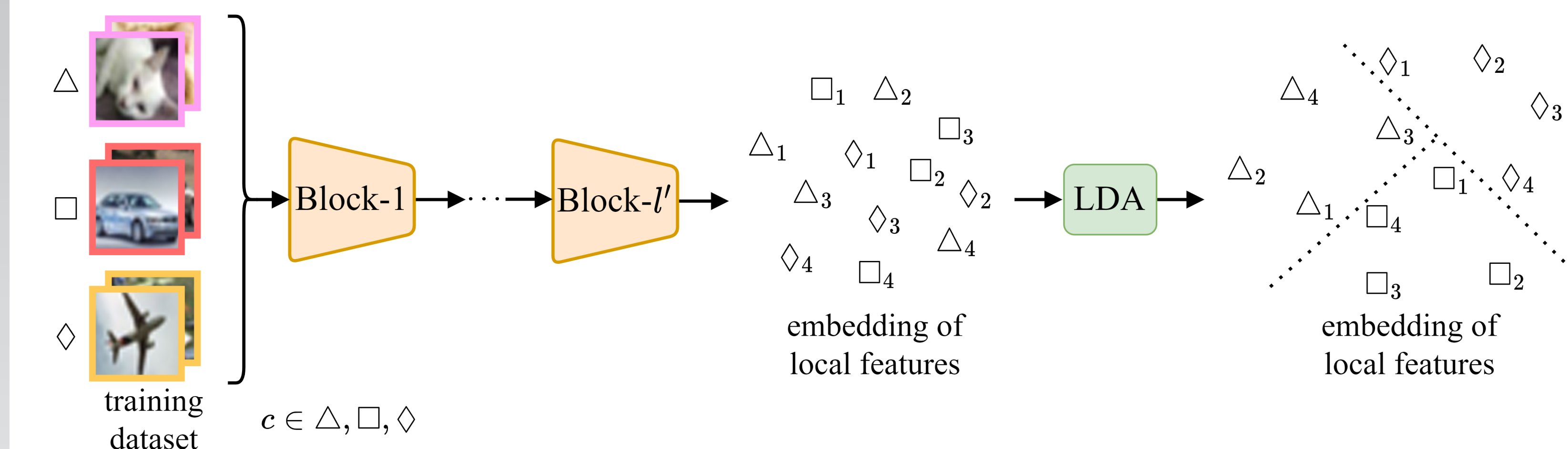
Data-driven mechanisms (e.g. *K-means*) lead to imitation of the teacher's geometry, hurting the performance of the student due to capacity differences.

Key Idea: Exploit the relationship between different classes (e.g. *leg* for *horse* and *deer*) that are learned by the teacher.

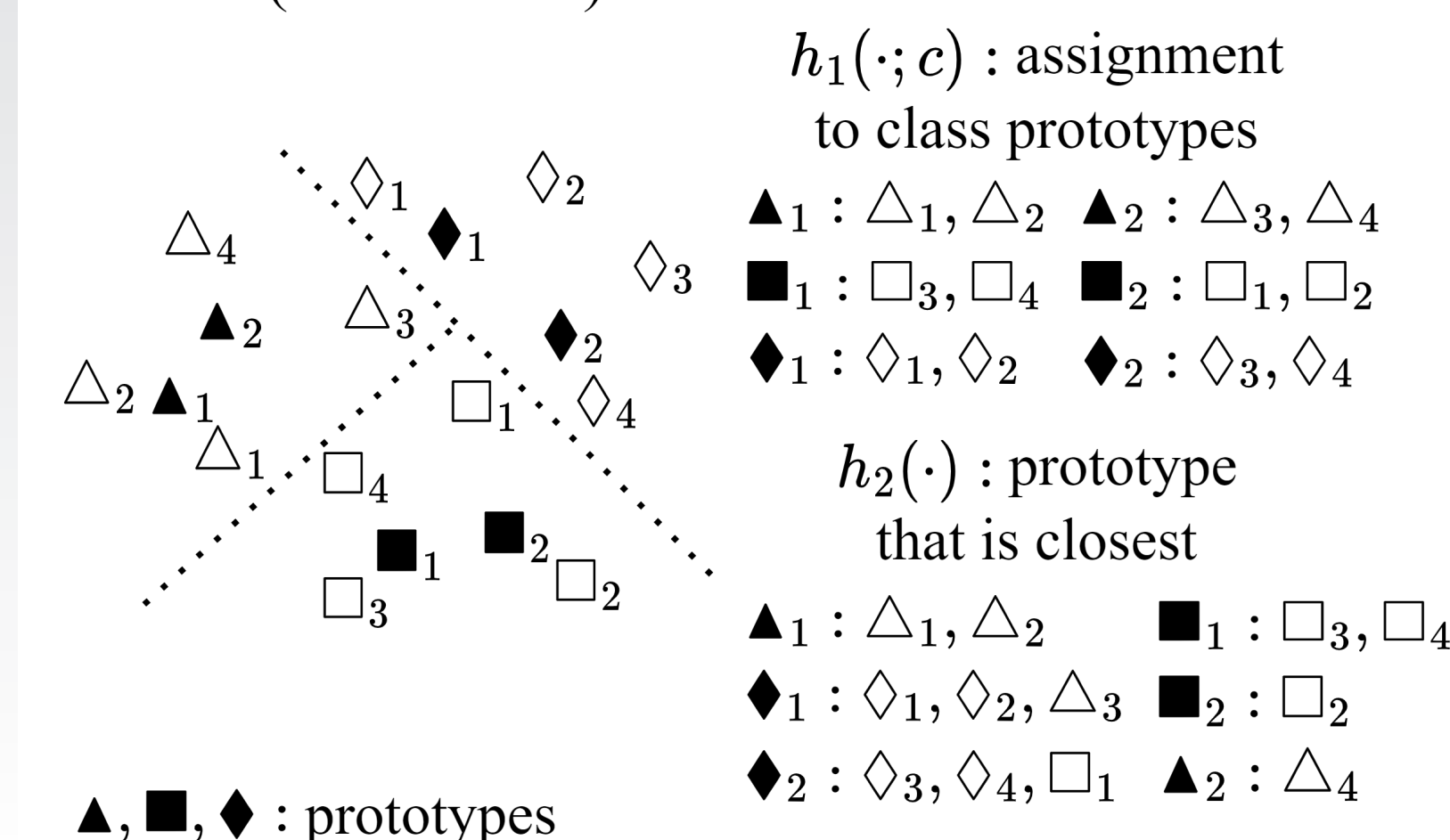
$p_{\mathcal{T}}(i)$: probability of classifying feature i as a particular prototype given i is assigned to its corresponding prototype.

Offline Teacher Supervision:

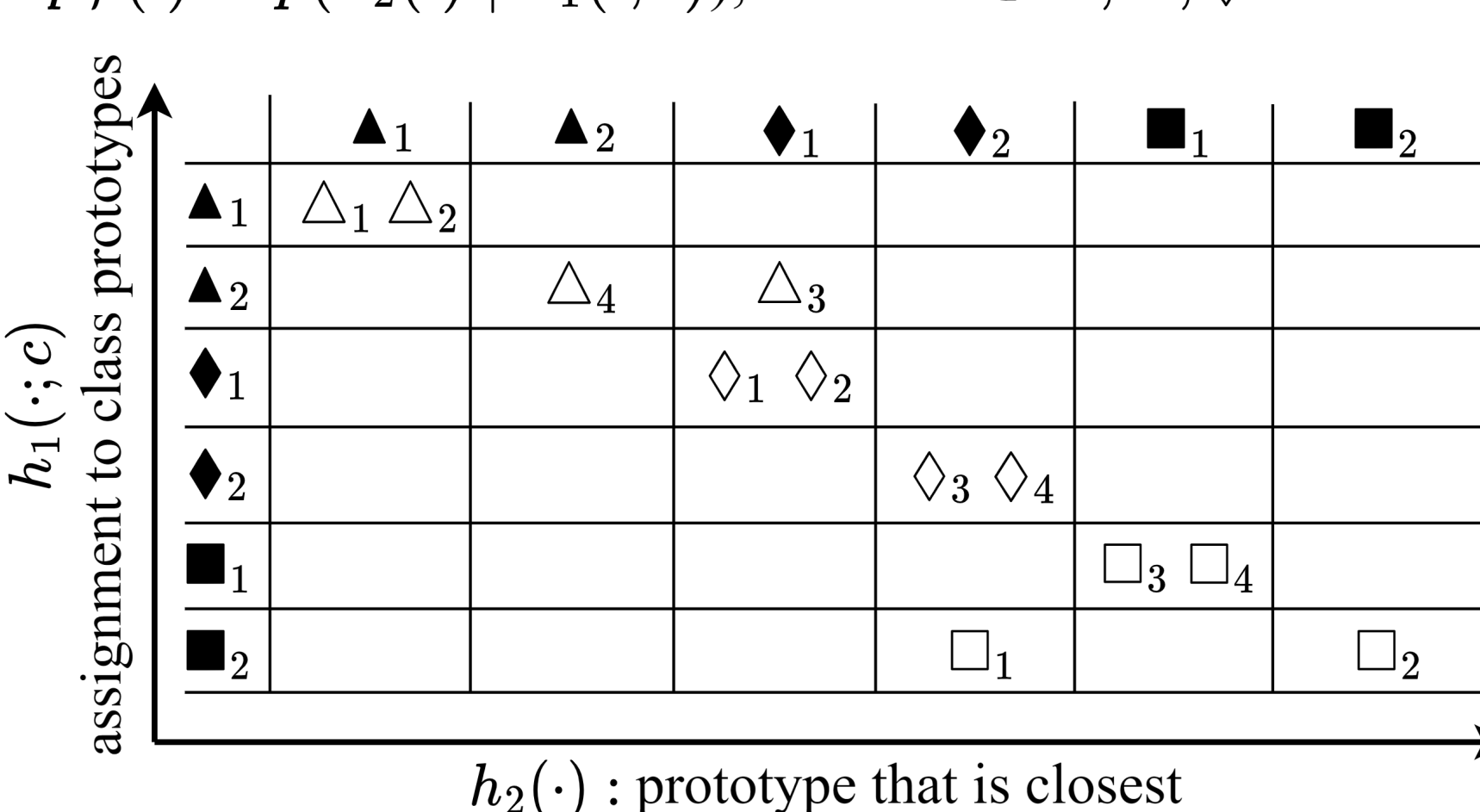
Step-1: Extract features of the training set and apply LDA.



Step-2: Apply *K-means* for each class separately. Obtain prototypes as the cluster centers (sub-classes).

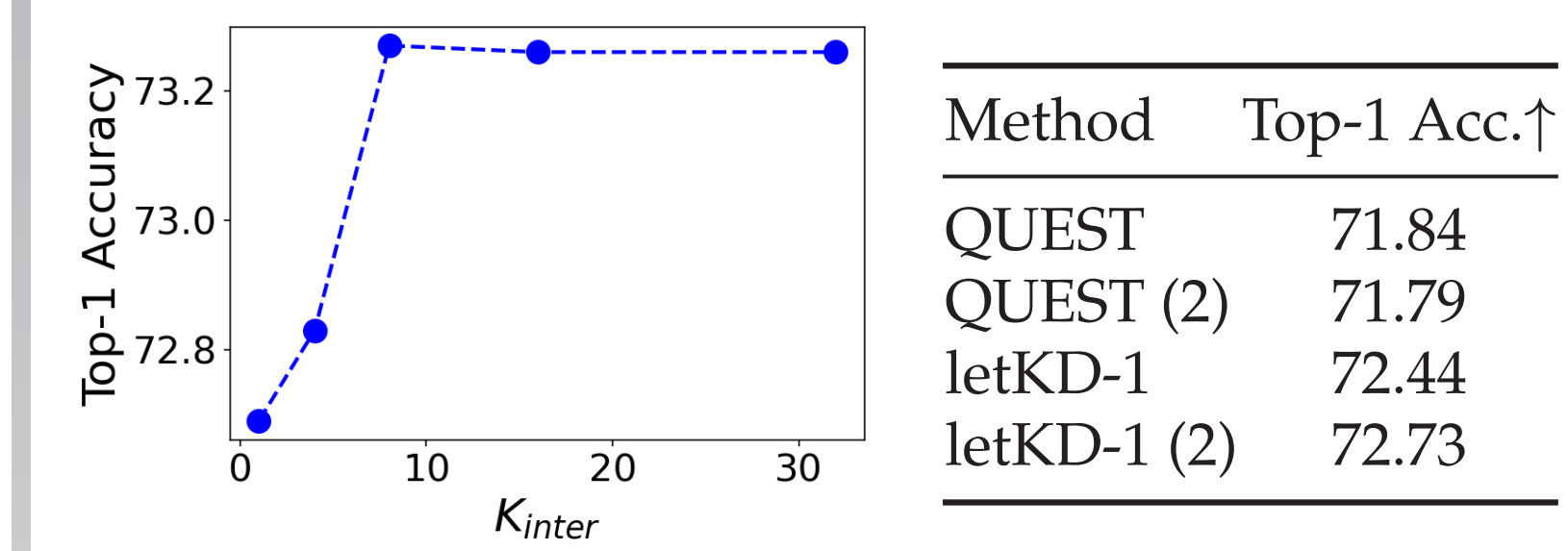


Step-3: Assign the nearest-neighbor prototype to each sample and normalize each row to calculate $p_{\mathcal{T}}(\cdot) = p(h_2(\cdot) | h_1(\cdot; c))$, where $c \in \Delta, \square, \diamond$.



ABLATION STUDIES

Clusters →	3x3 Kernels		K-Means		Teacher Student	
Archs. ↓	QUEST	letKD-1	QUEST	letKD-1	QUEST	letKD-1
RN56-RN20	71.92	72.11	71.84	72.44	72.34	69.06
RN110-RN32	74.31	74.44	74.08	74.40	74.31	71.14
RN83-RN29	72.41	72.61	72.48	73.33	73.84	70.53

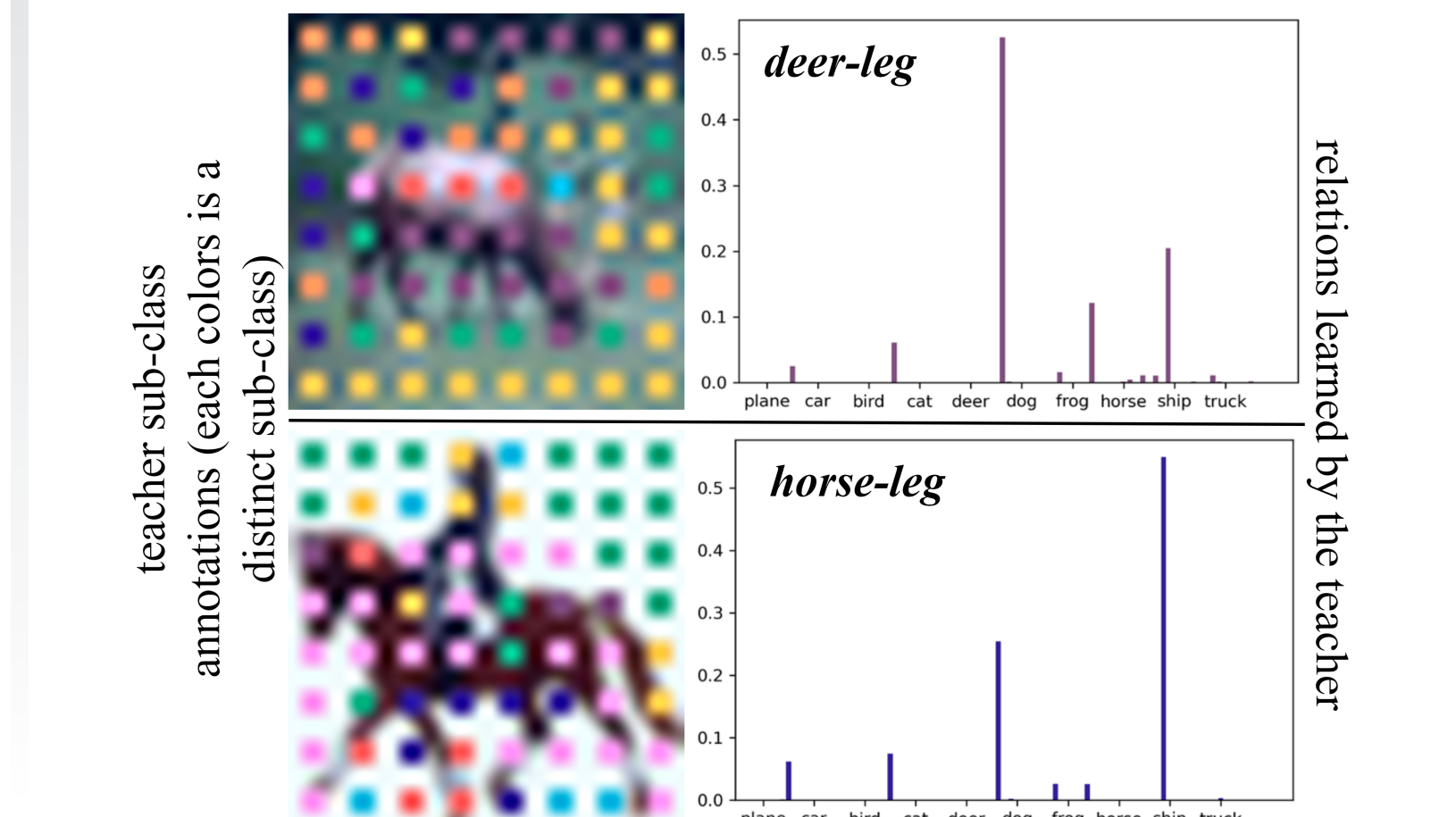


Effect of our KD layer:

Classification capacity in the intermediate layers is improved with our KD layer.

Marginal increase in the computation is not the main source of improvement.

The major contribution to the performance occurs upon combining the KD layer with our supervision.



Without K-means. Kernels of 3x3 of the residual blocks in architectures such as ResNet correspond to learnable templates (i.e., cluster centers) of some semantic entities.

Employing sub-classes (K_{inter}) for the intermediate layer distillation enables better knowledge transfer.

Naive multi-layer distillation (denoted as (2)) without our KD layer hurts the performance, suggesting that a better way should be found.

	letKD-2 ($\alpha_{inter} = 0$)	letKD-2 ($\alpha_{inter} = 1$)
	x	\hat{x}
Intermediate Layer Top-1 Acc.↑	52.71	51.71 56.36

Methods	Top-1 Acc.↑
FitNet	71.59
FitNet+KD layer without supervision	71.80
FitNet+KD layer with supervision	73.36

Inter.	α_{inter}	Penult.	α_{penult}	Top-1 Acc.↑
✓	0	-	0	70.64
✓	1	-	0	70.80
-	0	✓	0	71.84
-	0	✓	1	72.44
✓	0	✓	0	71.70
✓	0	✓	1	72.13
✓	1	✓	0	72.78
✓	1	✓	1	73.27

IMAGE RECOGNITION EXPERIMENTS: TOP-1 ACC.↑

CIFAR100	Homogeneous						Heterogeneous			
	Teacher Student	WRN-40-2 WRN-16-2	WRN-40-1	RN56 RN20	RN110 RN20	RN110 RN32	RN32x4 RN8x4	WRN-40-2 SNV1	RN32x4 SNV1	RN32x4 SNV2
Methods ↓	75.61	75.61	72.34	74.31	74.31	79.42	75.61	79.42	79.42	79.34
	73.26	71.98	69.06	69.06	71.14	72.50	70.50	70.50	71.82	64.60
SimKD	76.06	74.92	68.95	69.35	72.15	78.08	76.95	77.18	77.78	68.91
TDD	75.01	74.04	71.53	-	-	-	75.60	-	-	68.37
QUEST	76.10	74.58	71.84	71.89	74.08	75.88	76.75	76.28	77.09	69.81
letKD-1	76.29	75.01	72.44	72.68	74.40	76.70	76.93	76.65	77.75	69.97
	± 0.15	± 0.09	± 0.24	± 0.31	± 0.14	± 0.06	± 0.16	± 0.24	± 0.17	± 0.18
letKD-2	76.56	75.19	73.27	73.38	74.62	77.09	77.08	77.30	77.95	70.39
	± 0.22	± 0.13	± 0.16	± 0.14	± 0.20	± 0.18	± 0.12	± 0.12	± 0.06	± 0.23

ImageNet	Teacher	Student	KD	DKD	QUEST	letKD-1	letKD-2
RN34-RN18	Top-1	73.31	69.75	70.66	71.70	72.33	72.38
	Top-5	91.42	89.07	89.88	90.41	91.06	91.15
RN50-MNV2	Top-1	76.13	68.87	68.58	72.05	73.78	73.98
	Top-5	92.86	88.76	88.98	91.05	91.81	92.00