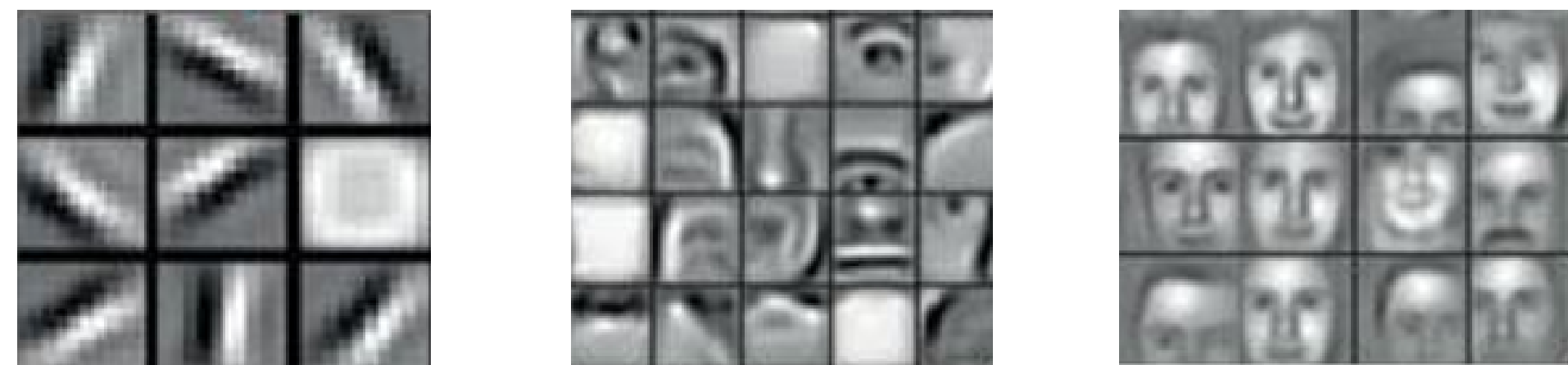


CONTRIBUTIONS

- Reformulating convolution block based local feature embedding as feature assignment through best matching kernel
- Repurposing *soft-max* as a batch-statistics-free replacement of BN-ReLU
- Exploiting mixture of class labels to shape the intermediate features of CNN

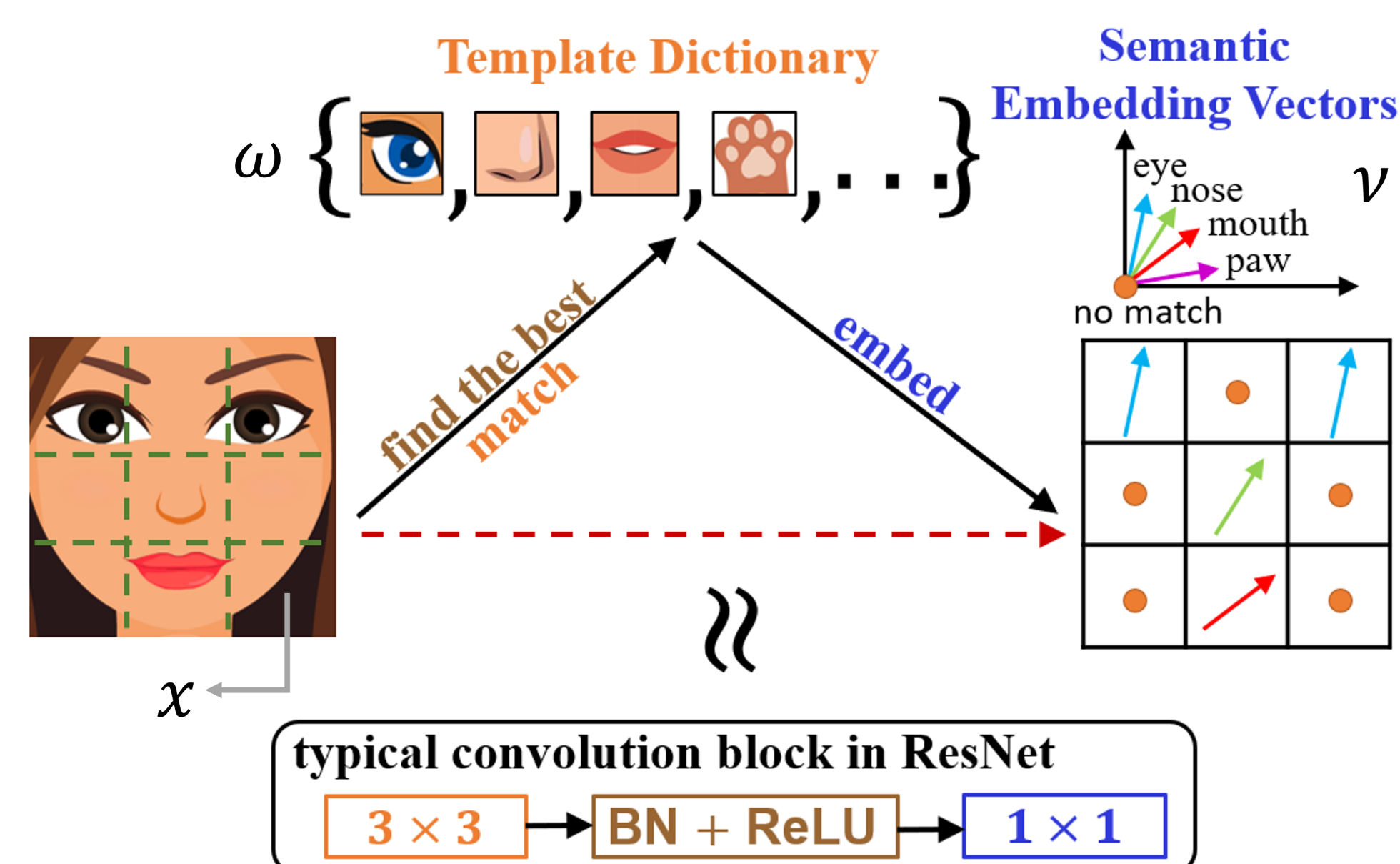
APPROACH

CNNs: Extracting hierarchical features through stacked convolution blocks whose parameters are learned in **top-down** manner via the feedback from a class-supervised loss function



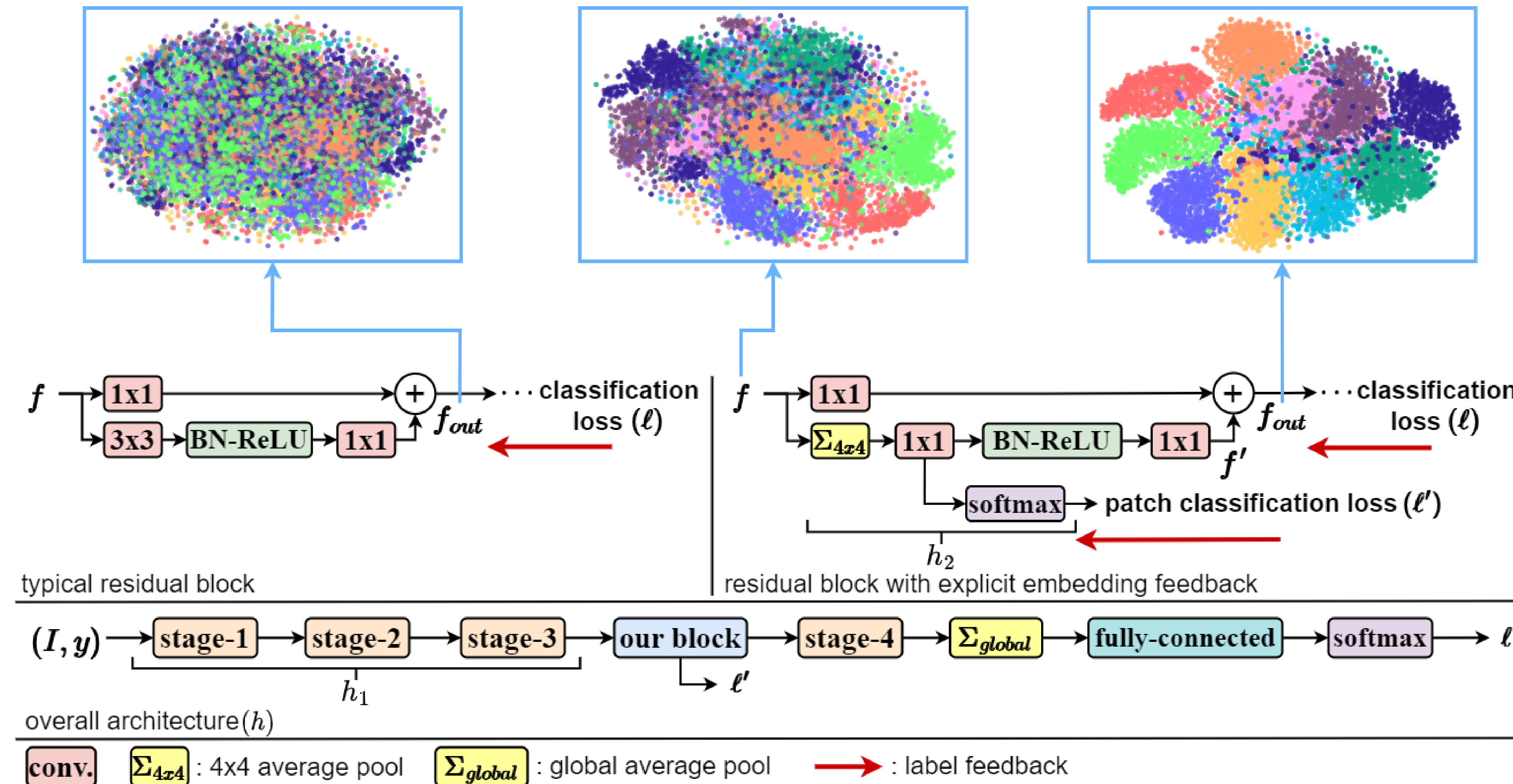
Low level features Mid level features High level features

- Differently, we want to use the **mixture of class labels** to synthesize new labels that can **explicitly supervise lower level feature extraction** (e.g. *plane* + *bird* \approx *wing*).
- We first formulate an optimization problem to analytically express selecting the best matching kernel and assigning its semantic vector.



- We then show its resemblance to a typical ResNet block.

METHOD OVERVIEW



- 1 We augment the training loss with an auxiliary per patch (x_{\square}) classification loss and enable label supervision in lower layers:

$$\mathcal{L}(\mathcal{D}) = \frac{1}{|\mathcal{D}|} \sum_{(I,y) \in \mathcal{D}} [(1-\lambda)\ell(h(I), y) + \frac{1}{wh} \sum_{x \in h_1(I)} \lambda \ell(h_2(x_{\square}), y)]$$

- 2 With our new block, we yield novel semantic entities as the convex combination of the class features and explicitly shape the feature embeddings according to these semantics.

THEORY TO PRACTICE

- Given a set of matching kernels ω_k and 3×3 features as $x_{3 \times 3}$, we define the problem as:

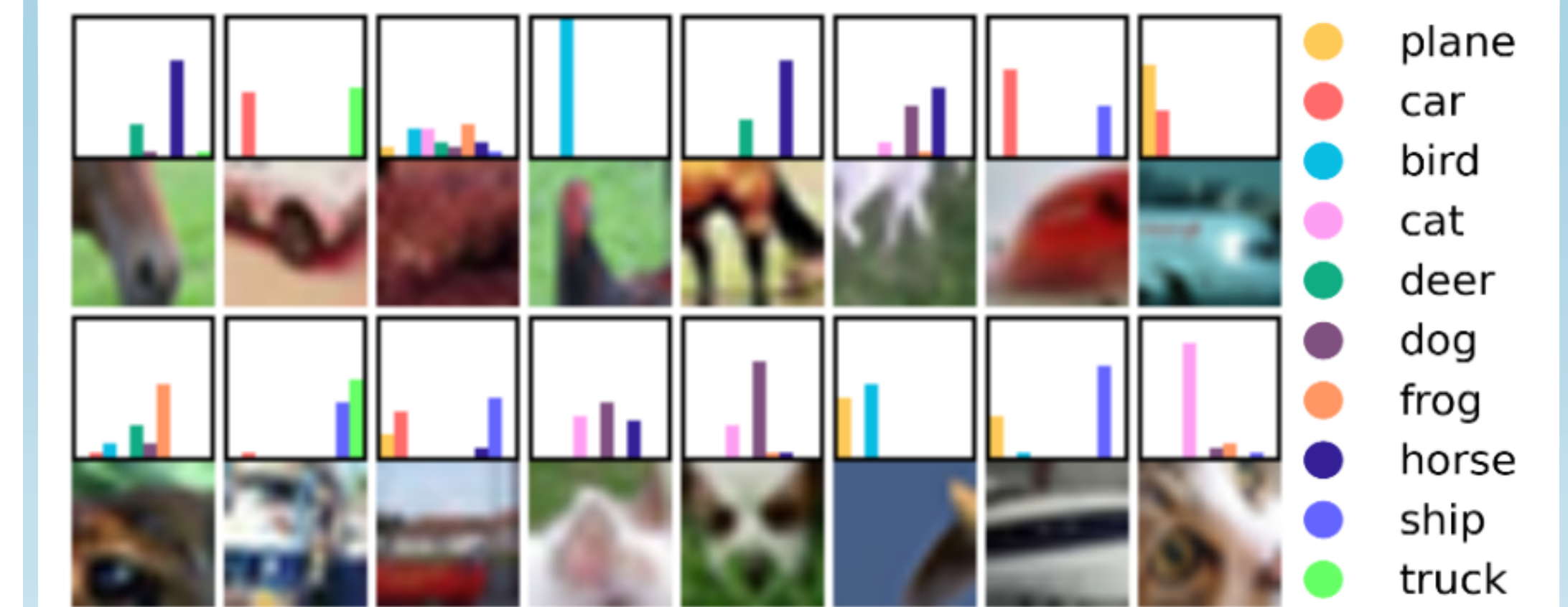
$$p^* = \arg \max_{p, q \geq 0, q + \sum_k p_k = 1} q \mu + \sum_k p_k \omega_k^T x_{3 \times 3}$$

- We make the above operation differentiable by smoothing the objective with entropy.
- We represent 3×3 pattern with $x' = \sum_k p_k^* \nu_k$.
- Relating BN-ReLU to *soft-max*, we show that 3×3 -BN-ReLU- 1×1 inherently perform these operations.

IMPLICATIONS

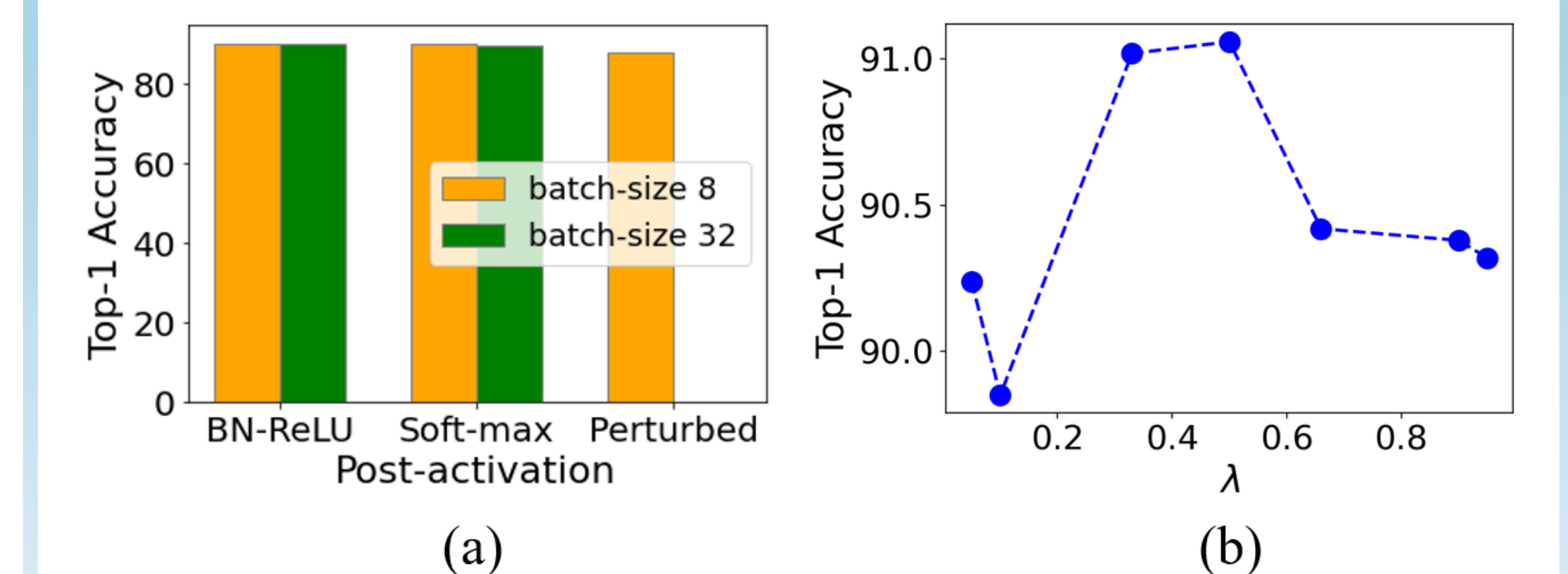
- We observe better clustering effects in which a clear distinction between *animals* and *vehicles* exists.
- Our method enables mixture of class labels and creates fine labels by using existing classes with the help of our embedding vectors which semantically reshape the overall geometry.
- Our formulation mimicks *cross-attention* between $x_{3 \times 3}$ (*queries*) and ω_k (*keys*) in which the final representation is given as the convex combination of ν_k (*values*).

ANALYSIS OF BEHAVIOUR



- We generate new semantic entities from the combination of class vectors such as *wing* from *bird* and *plane*, *tire* from *car* and *truck*.
- We observe discriminative patches specific to certain classes.
- We observe generic entities as the mixture of many classes such as *fur* for *animal* classes.

ABLATION STUDIES



(a) BN-ReLU \approx *soft-maximizer*

(b) The inclusion of our mechanism through auxiliary loss boosts the performance.

QUANTITATIVE RESULTS

Dataset \rightarrow		Cifar10	Cifar100	Mini-ImageNet
Architecture \downarrow	Params	top-1 acc.	top-1 acc.	top-1 acc.
RN26	0.96M+257C	89.52	65.94	60.43
RN26-aux.	0.96M+386C	90.57	66.21	60.70
RN26-Ours	0.98M+516C	91.06	66.78	61.23
RN38	1.42M+257C	90.78	68.15	60.72
RN38-Ours	1.44M+516C	91.36	69.01	63.83
WRN16	1.28M+129C	90.52	67.11	60.73
WRN16-Ours	1.30M+388C	91.10	67.36	62.92
DN100	1.20M+535C	92.62	71.65	65.03
DN100-Ours	1.32M+1222C	92.92	71.25	68.86
DN100-Ours-C	1.36M+1264C	92.71	72.14	68.93