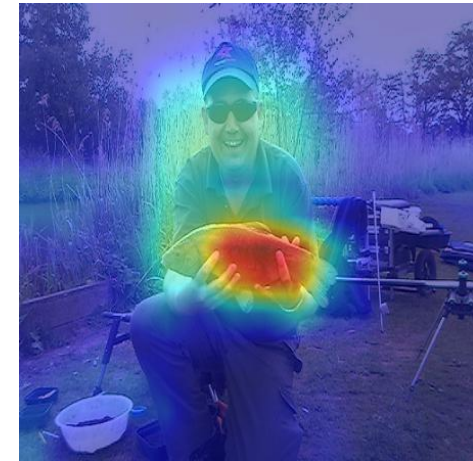


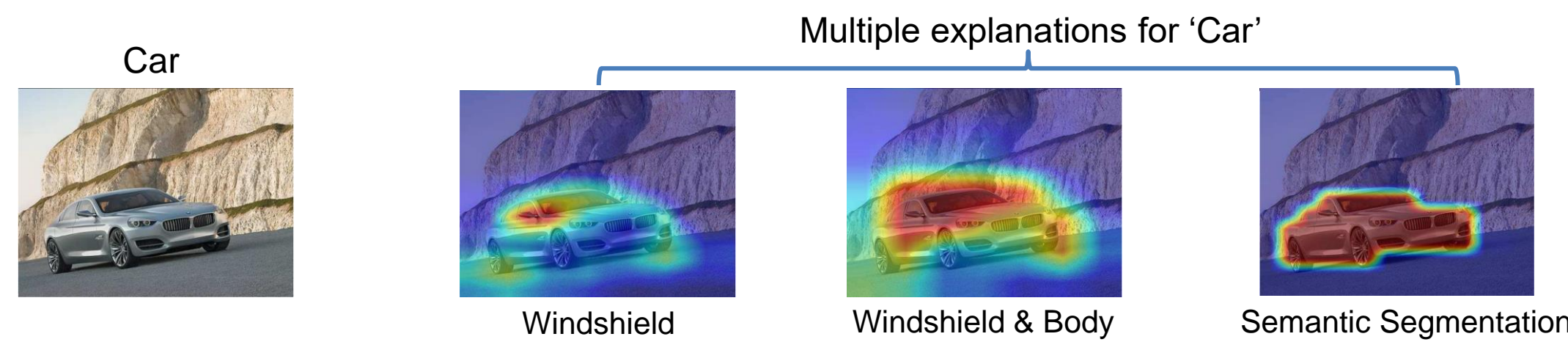
## Motivation

- Our goal is to train models that are more explainable given an explanation tool (e.g., Grad-CAM).

- Image classifiers may learn unwanted contextual biases from data



- Supervising explanation is not feasible since annotating explanations is not a well-defined task.

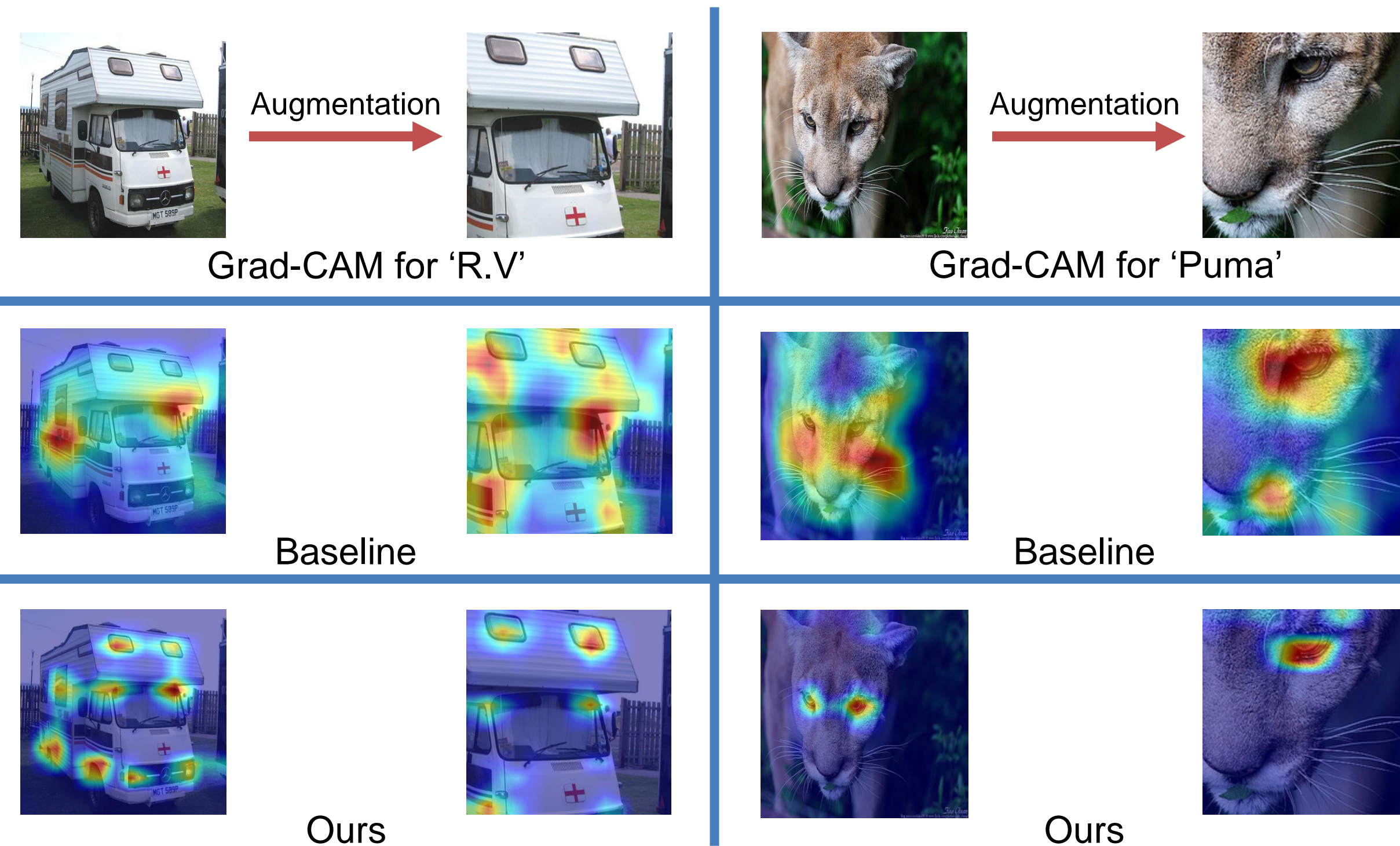


- We adopt ideas from self-supervised learning to improve explanation without annotating explanation.

- Our method can also leverage unlabeled data.

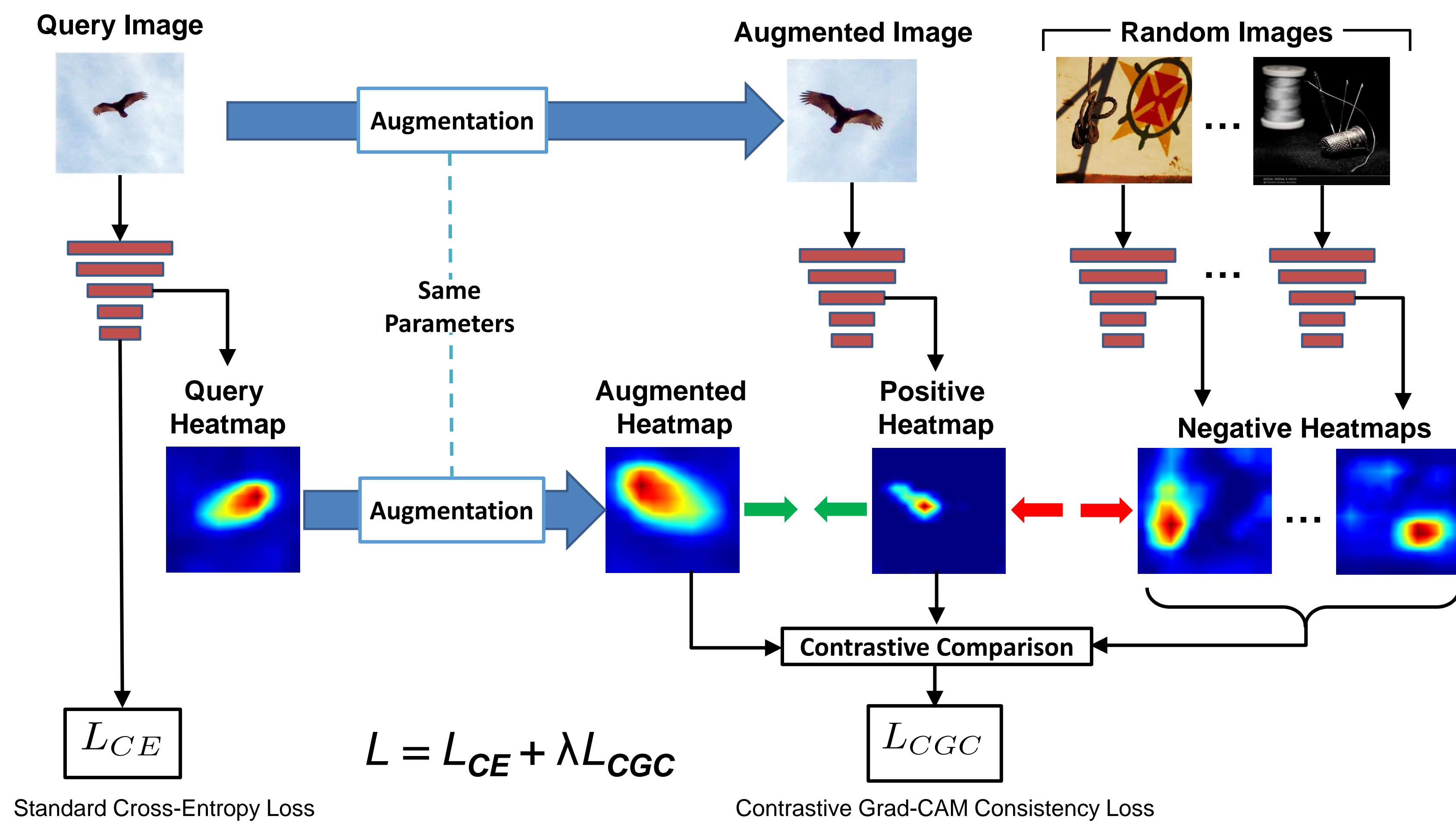
## Key Idea

- Consistent explanations w.r.t. spatial transformations

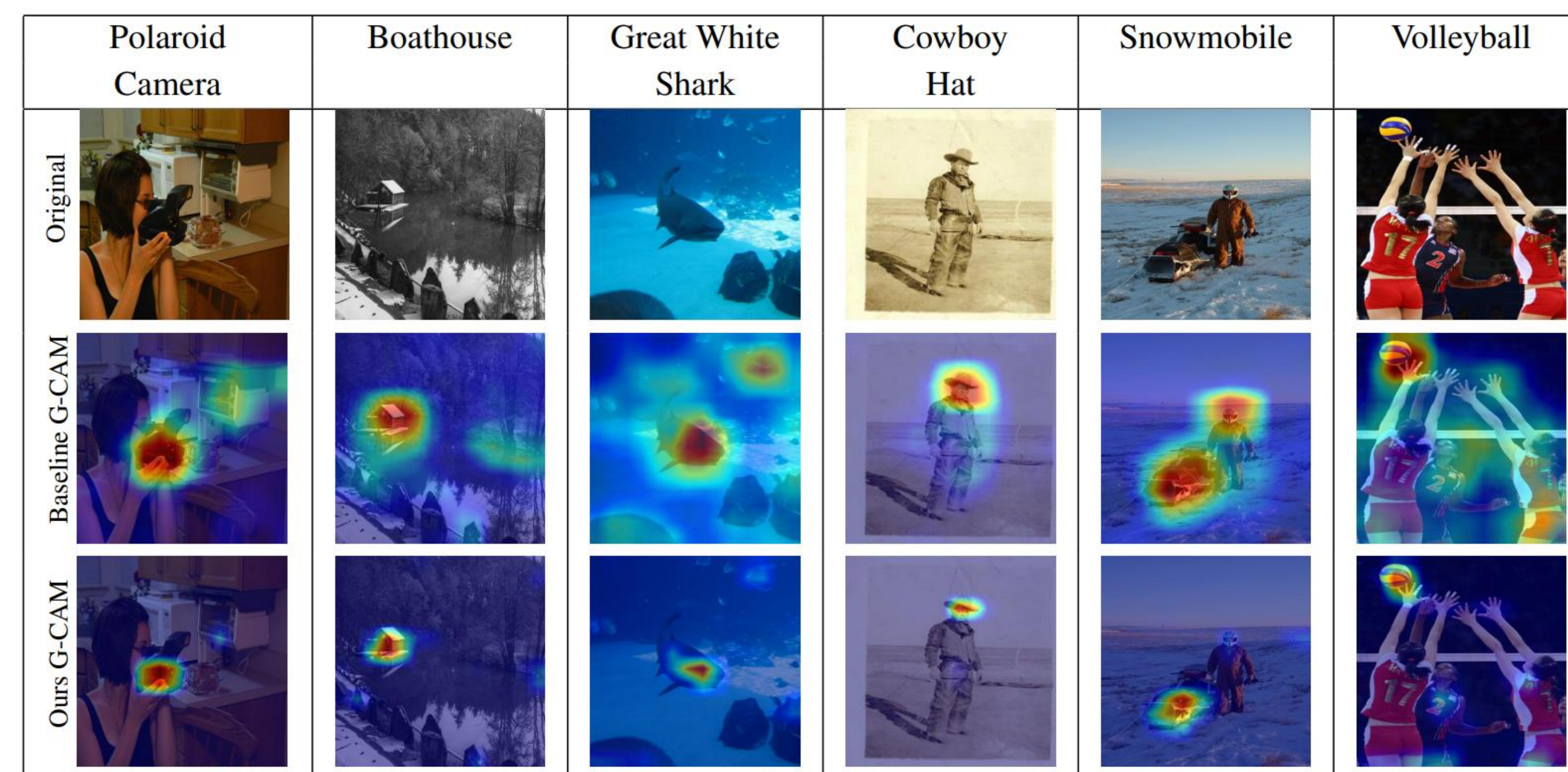


# Consistent Explanations by Contrastive Learning

## Method: Contrastive learning on the explanation space



## Results



Our model focuses on the most discriminative regions of the object instead of the background pixels.

**Content Heatmap (CH)** <sup>[1]</sup>: Percentage of heatmap strictly within object bounding box.

## ImageNet

| Architecture | Method        | Top-1 Acc (%) | CH (%)       | CGC Loss    |
|--------------|---------------|---------------|--------------|-------------|
| ResNet18     | Cross-Entropy | <b>69.76</b>  | 54.47        | 3.19        |
|              | GCC [1]       | 67.74         | 57.73        | 3.14        |
|              | Ours (CGC)    | 66.37         | <b>65.83</b> | <b>2.59</b> |
| ResNet50     | Cross-Entropy | <b>76.13</b>  | <b>54.77</b> | 3.15        |
|              | GCC [1]       | 74.40         | 59.42        | 3.09        |
|              | Ours (CGC)    | <b>74.60</b>  | <b>71.75</b> | <b>2.64</b> |

Our method results in models with improved explanation with marginal drop in classification accuracy.

## Fine-grained Classification

| Method        | CUB-200             | FGVC-Aircraft       | Cars-196            | VGG Flowers-102     |
|---------------|---------------------|---------------------|---------------------|---------------------|
| Cross-Entropy | 80.09 ± 0.89        | 83.65 ± 0.15        | 89.71 ± 0.14        | <b>96.09 ± 0.23</b> |
| Ours (CGC)    | <b>81.49 ± 0.09</b> | <b>85.72 ± 0.20</b> | <b>90.28 ± 0.08</b> | <b>96.18 ± 0.09</b> |

Our method improves the classification accuracy on fine-grained datasets.

## 1% Labeled ImageNet Subset

| Method        | Top-1 Acc (%) | Top-5 Acc (%) | CH (%)       |
|---------------|---------------|---------------|--------------|
| Cross-Entropy | 54.00         | 78.69         | 46.08        |
| Ours (CGC)    | <b>55.18</b>  | <b>79.12</b>  | <b>46.76</b> |

Our method leverages unlabeled data for  $L_{CGC}$  loss term (both models are initialized from SwaV[2])

## Conclusion

- We introduce a contrastive learning method for training **more explainable** models for a given explanation tool (e.g., Grad-CAM).
- Compared to manual annotation, our method significantly improves explanation consistency on ImageNet, UnRel, and fine-grained datasets.
- Our method acts as a regularizer that focuses more attention on the discriminating aspects of the image, thereby **reducing contextual bias**.
- We further leverage **unlabeled data** to improve the classification accuracy in **limited-label settings**.

## References

- [1] Vipin Pillai and Hamed Pirsiavash. Explainable models with consistent interpretations. AAAI 2021.
- [2] Caron, Misra, Mairal, Goyal, Bojanowski, Joulin. Unsupervised learning of visual features by contrasting cluster assignments. NeurIPS 2020.