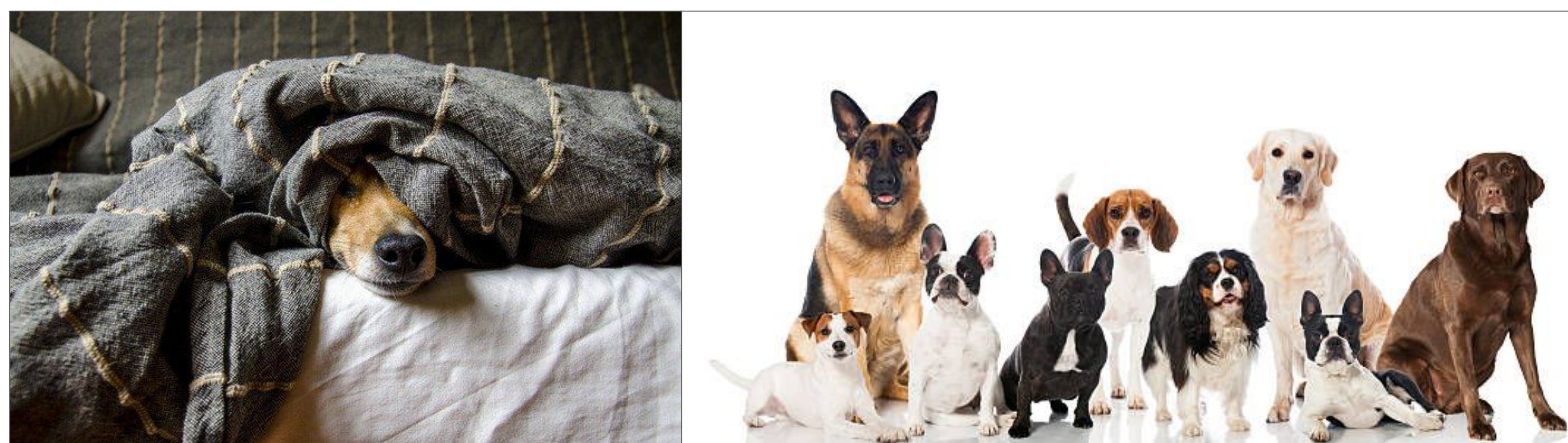




## Abstract

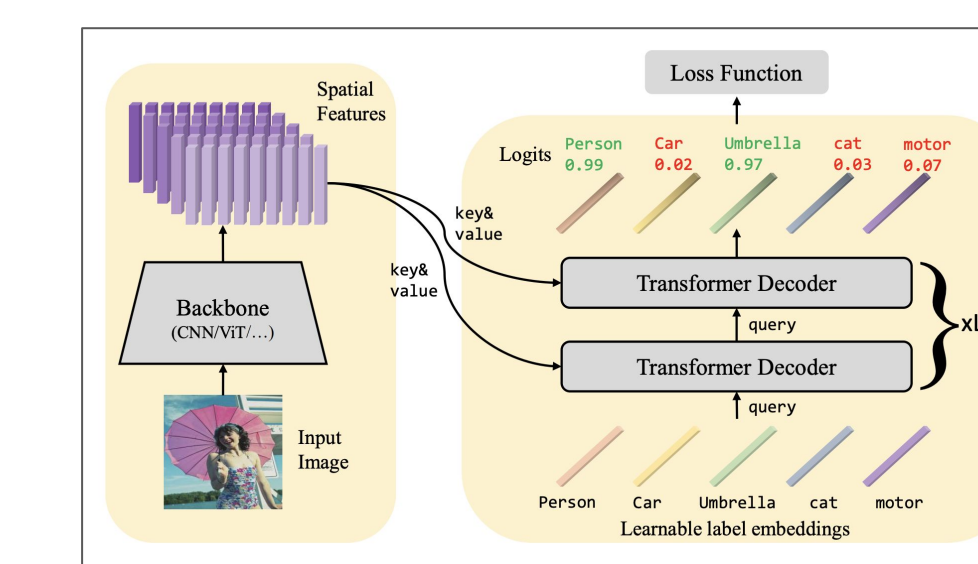
Motivation: Multi-label image recognition is challenging due to large variations in the size and spatial location of objects, even harder when objects are occluded and small.



Specific Problem: Existing methods are complex in terms of both model and data; they do not explicitly address problem of small objects and occlusions

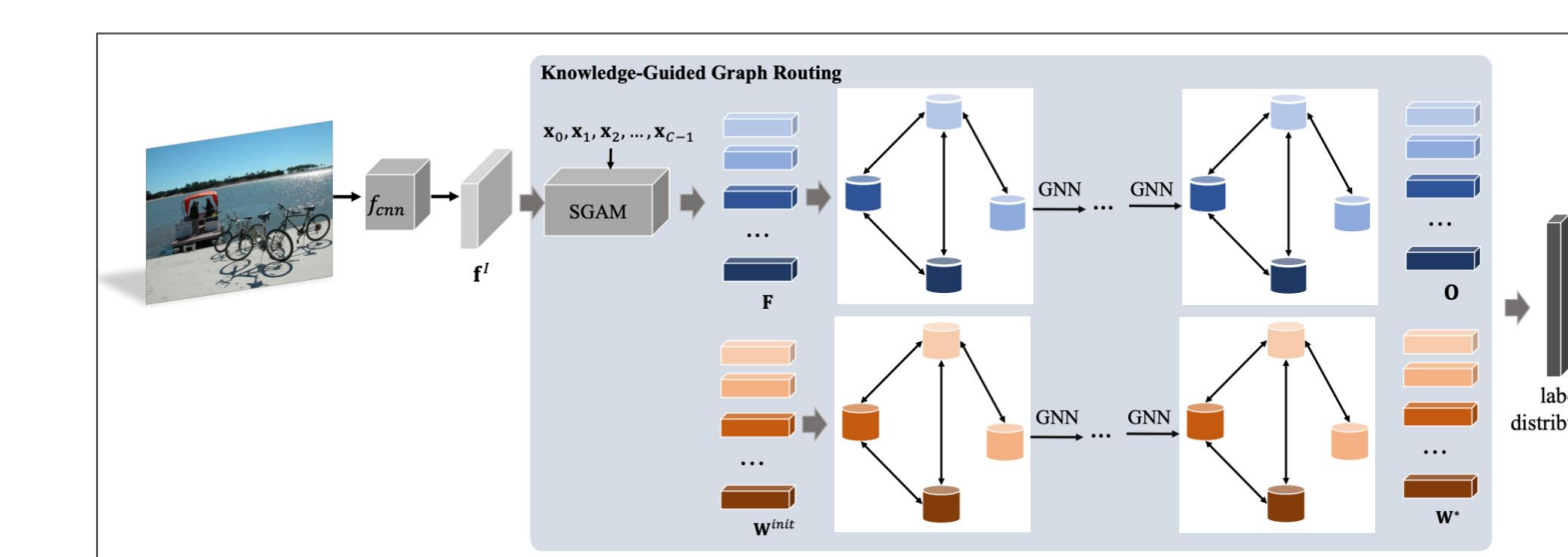
## Model

- Multiple stages of training
- Combination of multiple learnable networks
- Relies on large language models

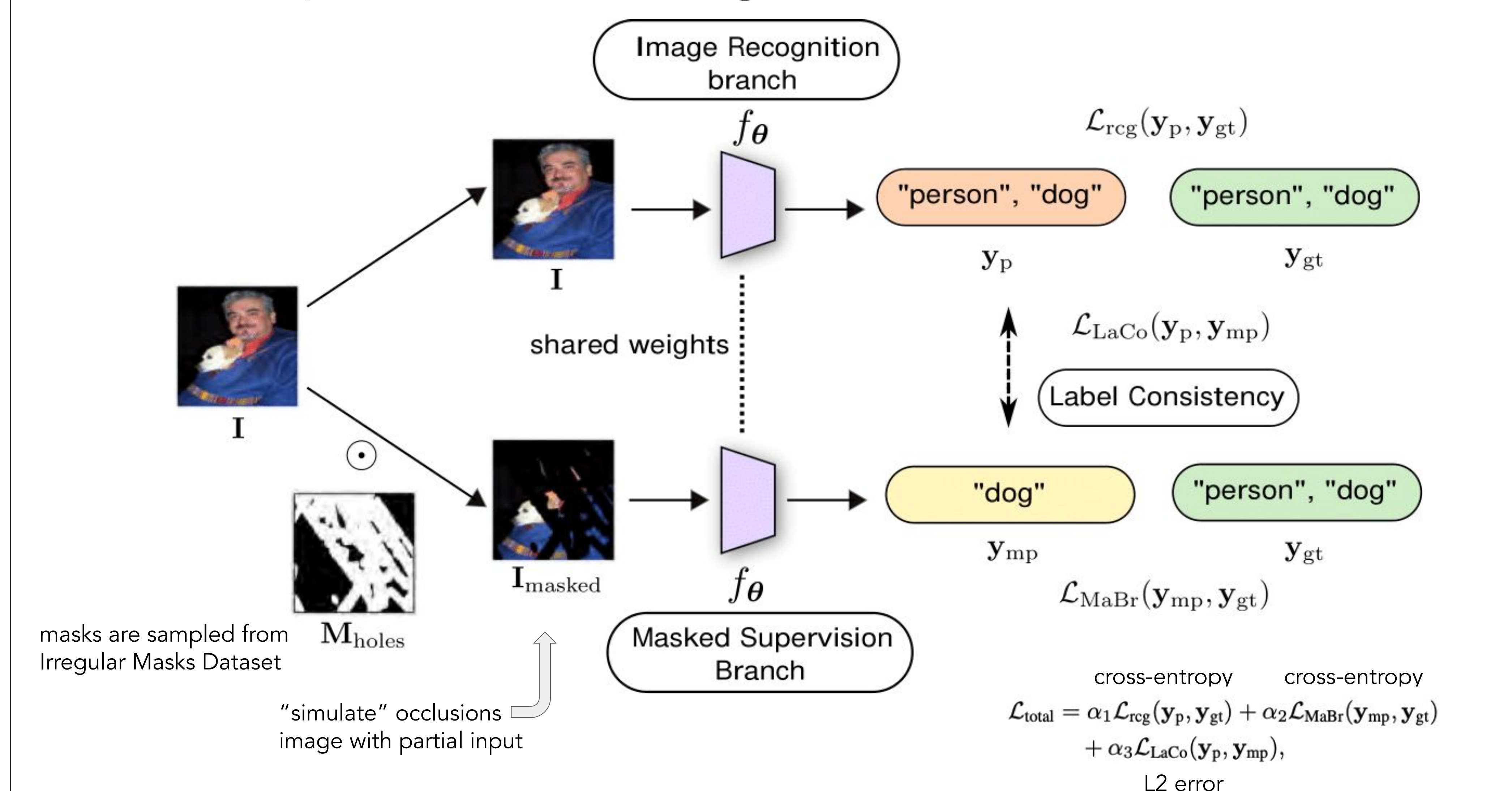


## Data:

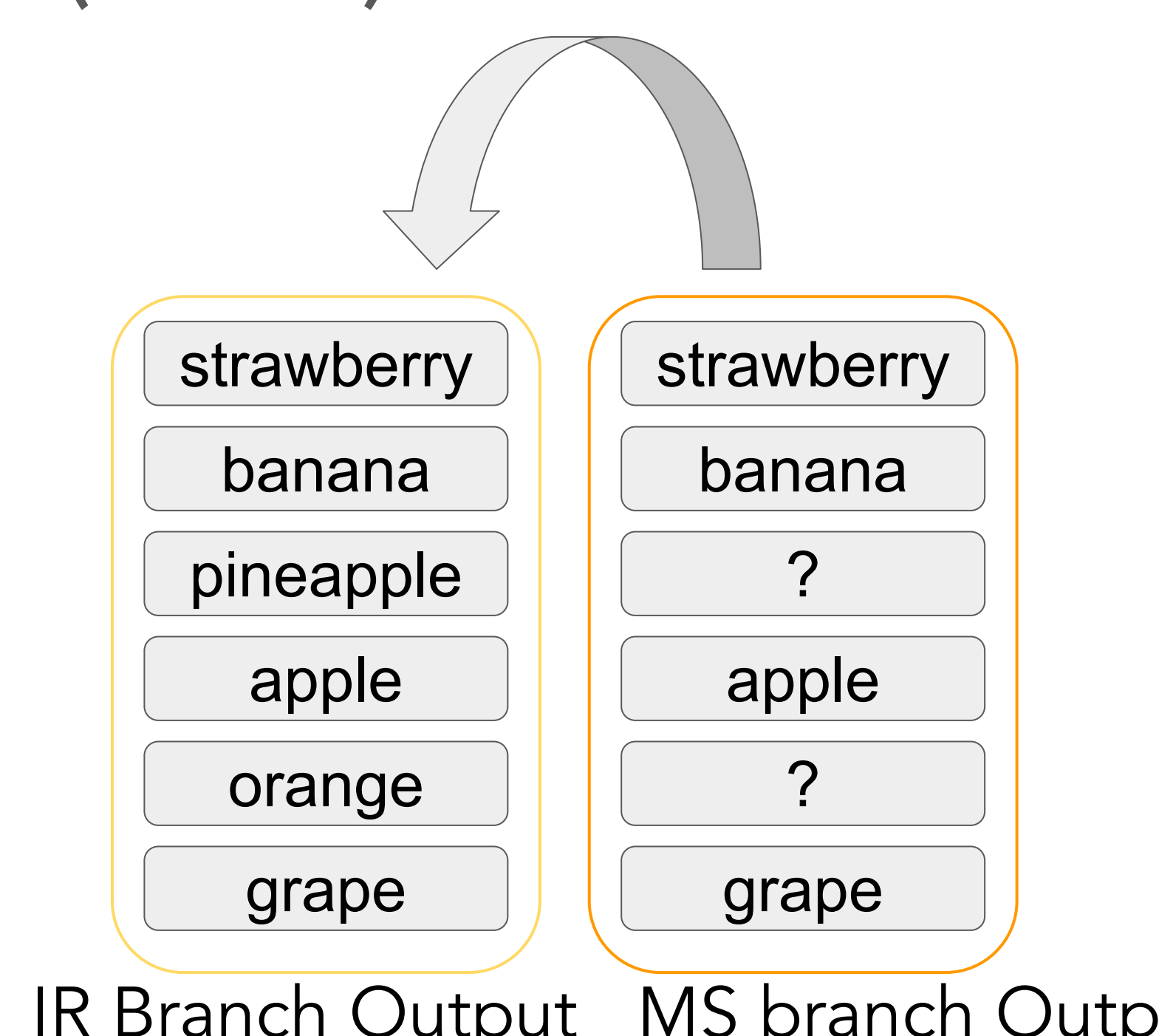
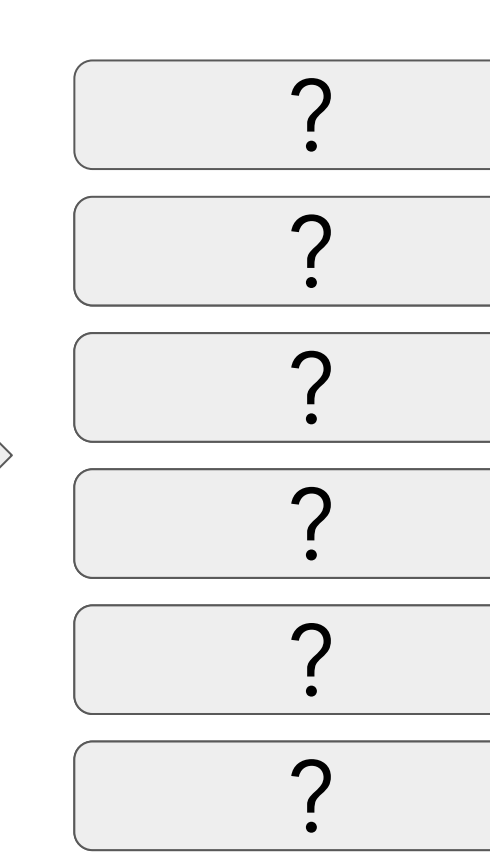
- High input resolution
- Complex data augmentation
- Additional data



## Masked Supervised Learning (MSL)



Masked Branch (MaBr) aims to learn context based representation; Label Consistency (LaCo) models label co-occurrence.



MaBr: uses nearby non-masked regions around objects to make predictions for partly visible/masked objects

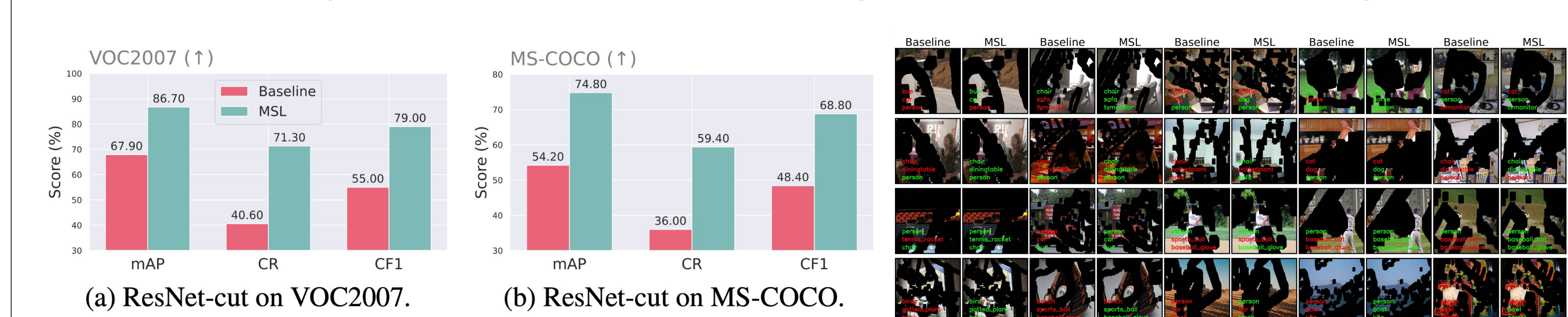
LaCo: learn a distribution of association across different classes  
 "use context to infer presence of a class"

## Results

Outperforms SOTA methods on VOC2007, MS-COCO and WIDER-Attribute.

Method	mAP	CR	CF1	Method	Input Resolution	mAP	CR	CF1	OP	OF1
ResNet [11]	82.9	-	-	ResNet [11]	448 x 448	79.4	83.4	66.6	74.0	88.8
ViViA-V [20]	93.0	-	-	FLA [17]	256 x 256	-	80.4	68.9	74.2	81.5
Ames-ReliefNet [1]	92.0	-	-	ResNet-cut [15]	448 x 448	82.1	86.2	68.7	76.4	88.9
SCRF [22]	92.5	-	-	ML-GCN [9]	448 x 448	83.0	85.1	72.0	82.8	85.4
SSRL [4]	93.4	-	-	MS-CMA [19]	448 x 448	82.8	82.9	74.4	78.4	84.4
MSL [8]	94.0	-	-	KSRN [10]	448 x 448	82.7	84.6	72.2	87.8	92.5
ADG-GCN [23]	94.0	-	-	MCAR [13]	448 x 448	82.8	82.9	72.1	78.0	88.0
ADG-GCN [23]	94.6	-	-	CSRA [15]	448 x 448	84.3	83.5	74.3	78.6	83.1
BMML (gov) [17]	95.0	-	-	TORIS [11]	448 x 448	84.0	82.0	72.8	78.4	83.3
DA-Ref [21]	94.3	-	-	OSL-RN [22]	448 x 448	84.2	84.1	72.1	76.4	81.2
ANL [1]	94.6	-	-	DVA-RN [11]	448 x 448	82.8	82.8	-	-	-
MCAR [13]	94.8	-	-	SST [1]	448 x 448	84.2	84.1	72.1	76.4	81.2
CSRA [15]	93.7	82.3	88.3	PAGEN [7]	448 x 448	83.2	84.9	72.7	78.3	83.0
KQGR [4]	93.6	-	-	KQGR [4]	448 x 448	84.3	84.6	72.3	76.4	81.2
KQGR (gov) [4]	93.9	-	-	ADG-GCN [23]	256 x 256	83.2	84.7	75.9	80.1	84.9
SST [1]	94.5	-	-	SOGR [6]	256 x 256	83.8	83.9	68.5	76.1	80.8
MSL-V	95.0	84.8	89.5	CTra [16]	256 x 256	83.1	84.3	74.3	79.2	83.7
MSL-C	96.1	92.4	93.6	MCAR [13]	256 x 256	84.5	84.3	73.9	78.2	86.9
				MSL-C	448 x 448	86.4	90.1	76.2	80.4	89.1

Robust to heavily masked inputs; would be good at occlusions naturally.



Better than self-supervised method; better at heavy masking; generic and applicable to any model architecture.

Masking	VOC2007	MS-COCO
MAE [14]	95.3	85.5
MSL	<b>96.1</b>	<b>86.4</b>

Method	Masking	VOC2007	MS-COCO
MSL-V	Low	94.6	77.8
MSL-V	High	<b>95.0</b>	<b>79.0</b>
MSL-C	Low	95.0	85.1
MSL-C	High	<b>96.1</b>	<b>86.4</b>

Architecture	VOC2007, mAP (%)	MS-COCO, mAP (%)
VIT	94.4	76.8
+ MSL	<b>95.0</b>	<b>79.0</b>
ResNet	93.7	84.3
+ MSL	<b>96.1</b>	<b>86.4</b>

Method	VOC2007, mAP (%)
MCAR [13]	94.8
MCAR [13] w/ MSL	<b>95.6</b>
SST [8]	94.5
SST [8] w/ MSL	<b>95.8</b>

## Conclusion

MSL is a simple yet effective single-stage, model-agnostic learning paradigm using masking. Use context on both image- and label-level to infer presence of multiple objects, even in cases where object are small and occluded.

Compare to recent methods, MSL does not require:

- multiple stages of training, combination of multiple networks, reliance on large language models (LLMs)
- high input resolution, complex data augmentation strategies, additional data for pretraining

TLDR: MSL is a simple and effective single-stage model-agnostic learning paradigm for visual learning tasks. It is similar to how us humans use context to perceive the visual world. Do try it!