

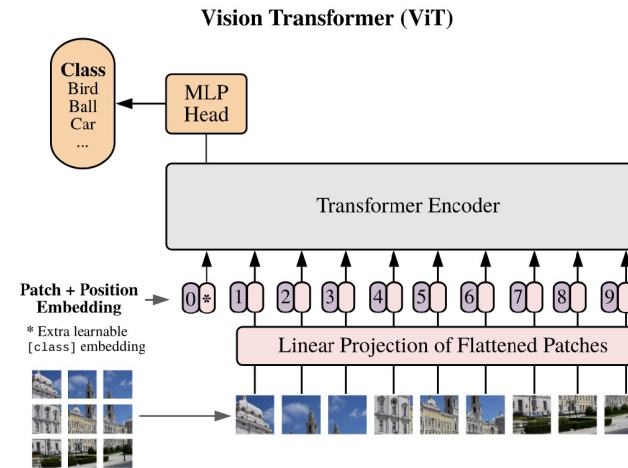
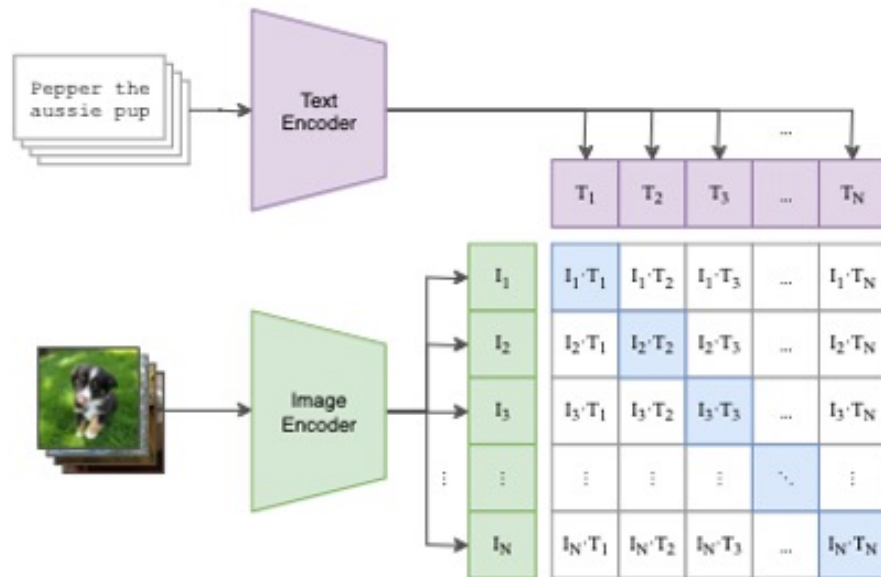
Safety in AI

Nicolas THOME – Prof at Sorbonne University
ISIR Lab, MLIA Team



Context: AI/ML summer

- AI in the last decade: huge performance boost
 - Vision, NLP, multi-modal prediction robotics

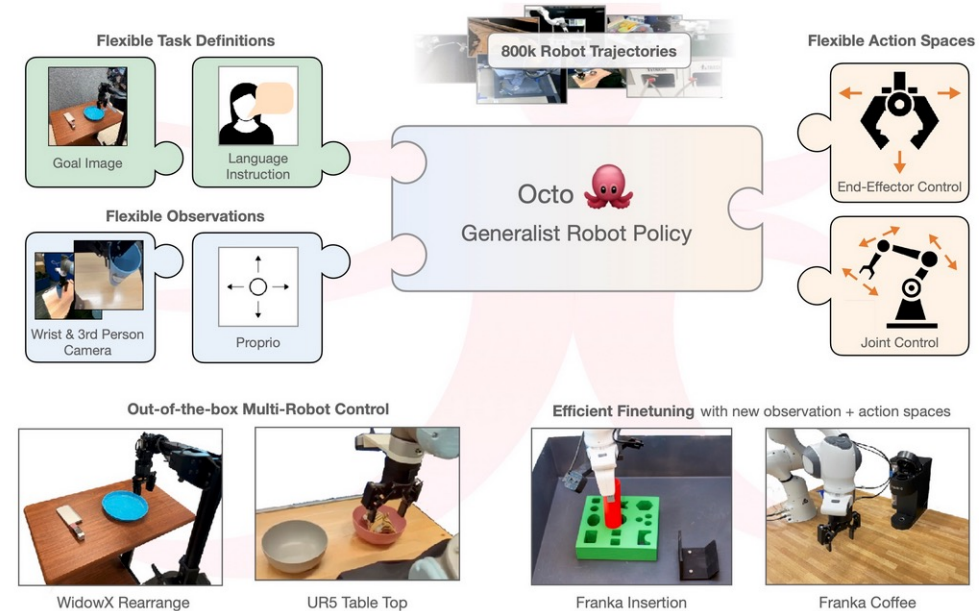


explain quantum theory to a child in the style of snoop dogg

Yo, little homie, let me break it down for you
 Quantum theory, it's a crazy thing, it's true
 See, the world is made up of tiny, tiny particles
 That are always moving, and changing, and actin' wild

These particles, they don't always behave like we think
 Sometimes they're waves, and sometimes they're points, it's a trip
 They can be in two places at once, and they can be in two states
 It's a mind-blowing concept, but that's the quantum fate

So don't be scared, little homie, it's all good
 Quantum theory, it's just the way the world works, it's understood
 It's a wild ride, but it's a cool one too
 Quantum theory, it's the way the world does what it do.



Context: safety in deep learning

Robustness: several brittleness aspects in deep learning

- Explainability, biases & shortcuts, fairness, *etc*

i) Stability: adversarial examples, mistake severity

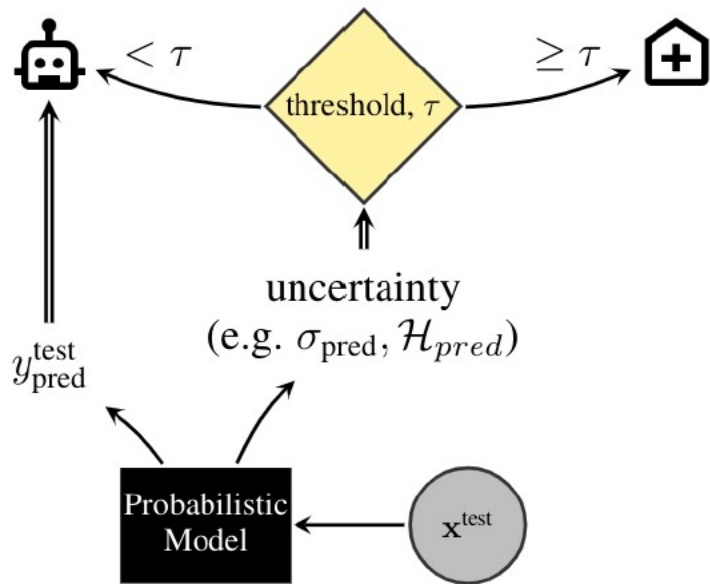
Query image



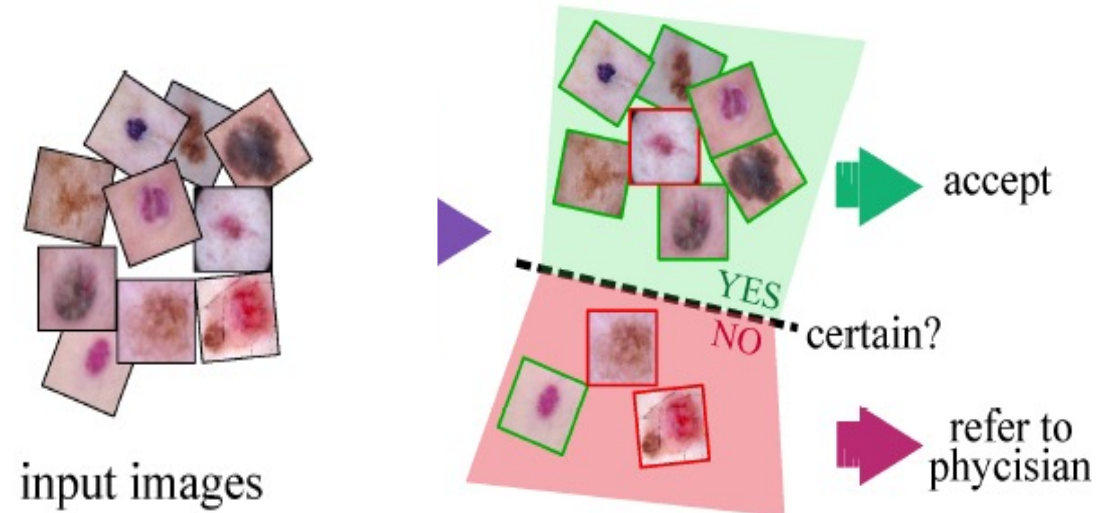
Context: robustness in deep learning

ii) Uncertainty Quantification (UQ)

“Know when you do not know”



Abstain to make a prediction



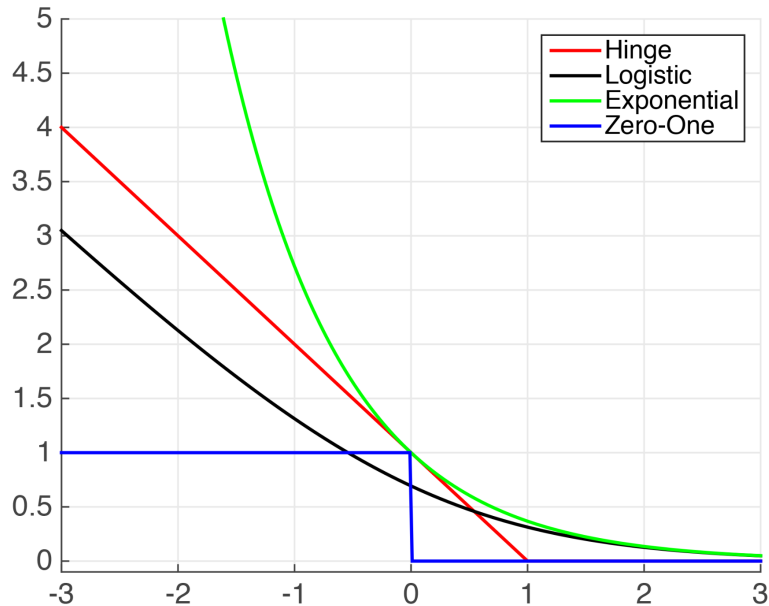
UQ: a challenge in DL

- Which uncertainty score?
- Calibration, ranking correct/incorrect prediction

Context: robustness in deep learning

iii) Training: direct optimization of target metrics, image retrieval

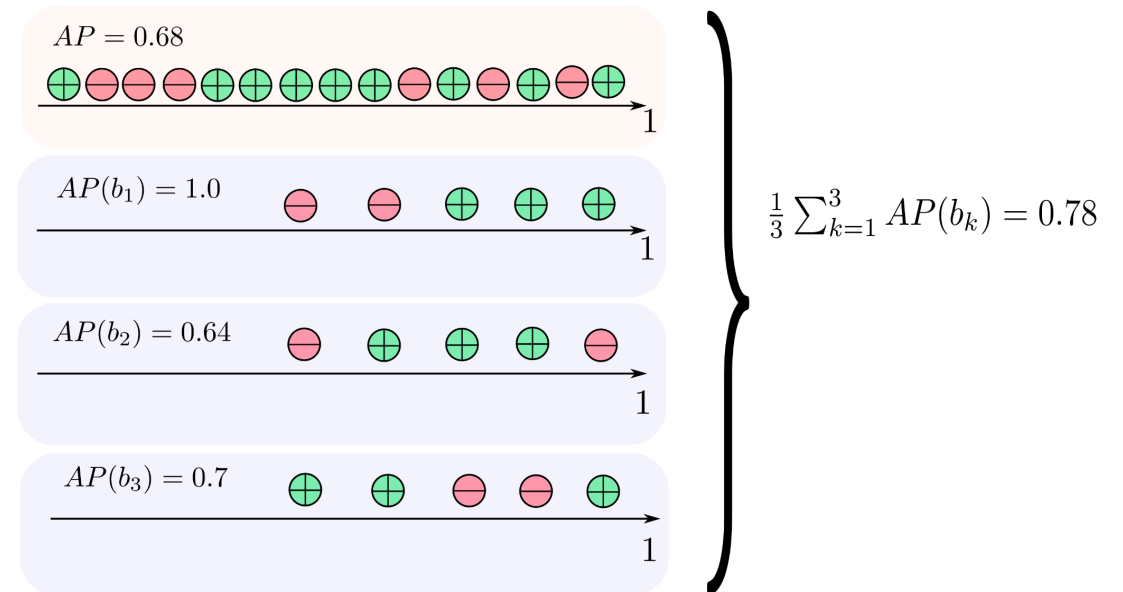
- Non-differentiable losses



“Good” (with guarantees) surrogates for classification, what about other metrics, rank losses?

- Non-decomposable losses

- Rank losses, e.g. Average Precision



- Many other losses (IoU, Dice in segmentation), including global constraints (fairness)

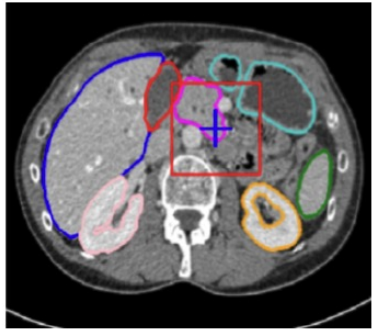
Improving robustness in deep learning

1. Uncertainty quantification
2. Direct optimization of rank losses
3. Controlling mistake severity

Sources of uncertainty

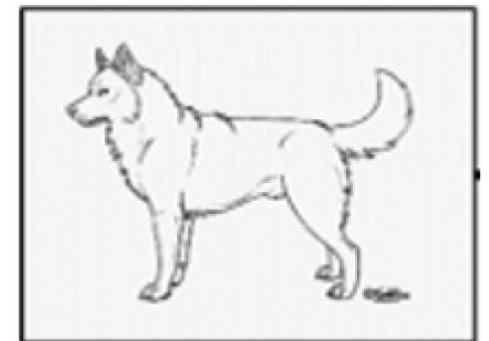
- Aleatoric uncertainty: data

- Class confusion, ambiguous data, sensor noise



- Epistemic uncertainty: model

- Distribution shift in $p(x,y)$, e.g. x (snow, image- \rightarrow cartoon), or y (open set, new classes)

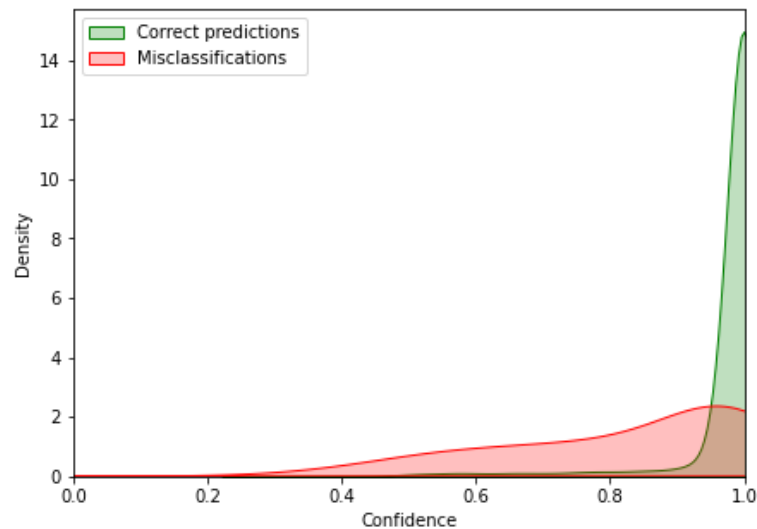
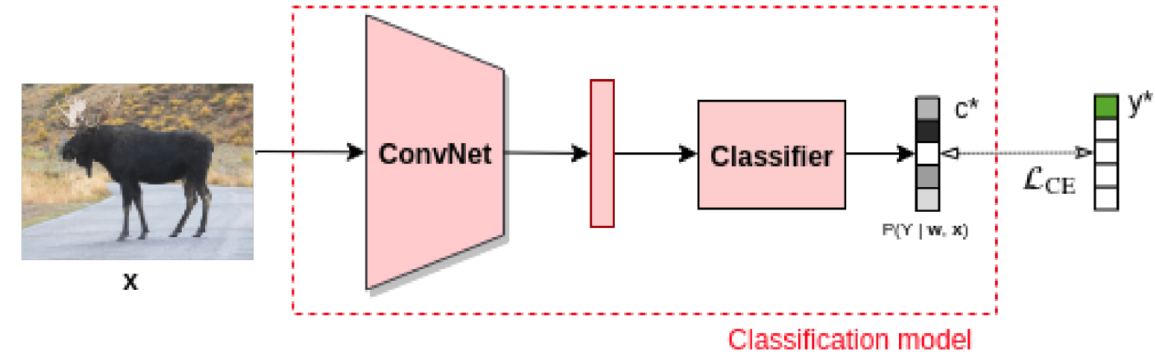


Uncertainty quantification in deep learning

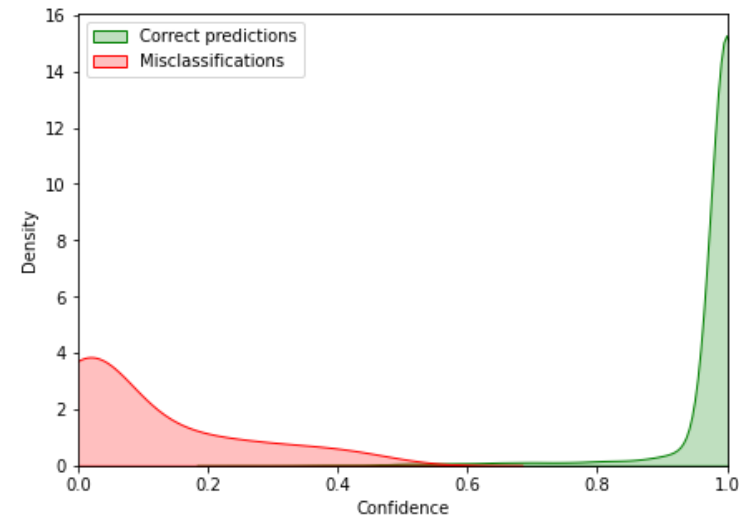
- **Uncertainty for failure prediction [CBT+19]:** correct vs incorrect predictions

- **Our proposal: True Class probability (TCP)** vs Maximum Class Probability (MCP)

- TCP better than MCP for failure prediction



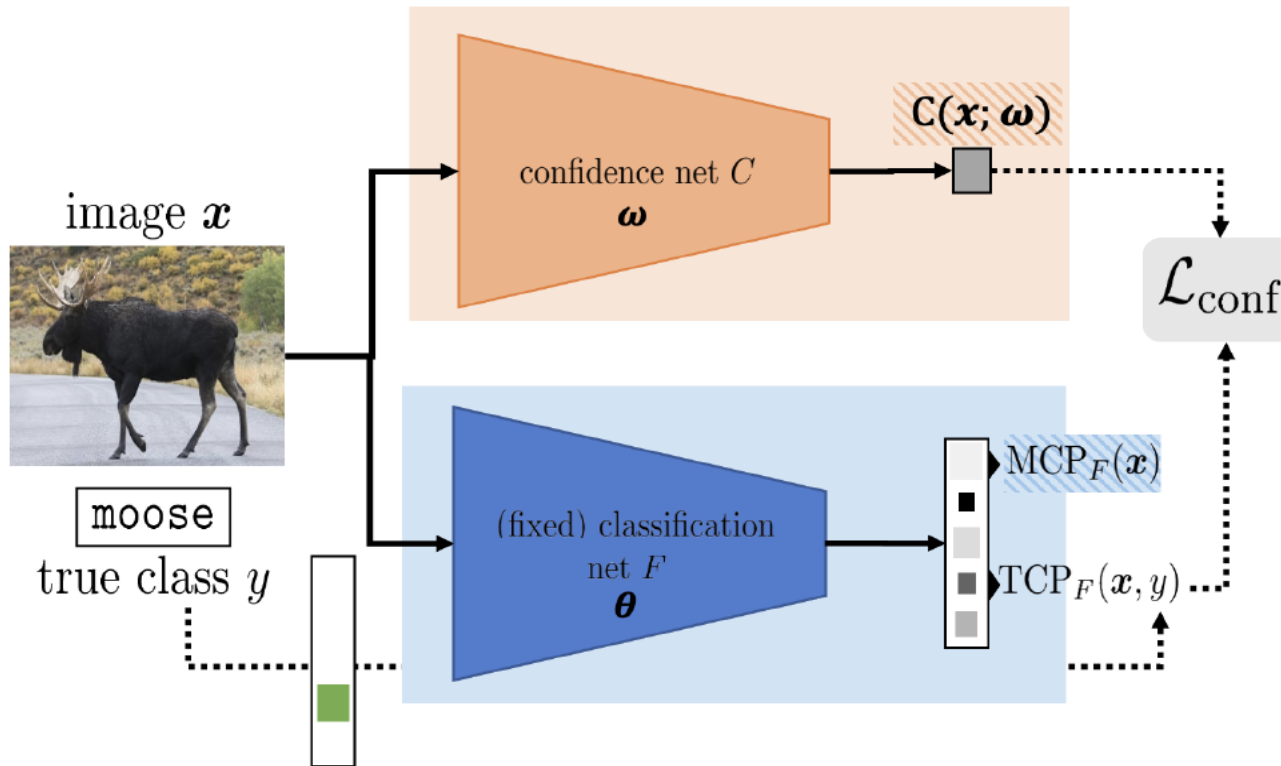
MCP



TCP

Uncertainty quantification in deep learning

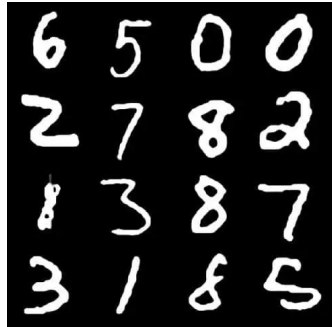
TCP unknown at test time: learning it! => ConfidNet



- Pre-trained prediction model (blue)
- Learning to regress TCP with an auxiliary model (orange)

$$\mathcal{L}_{conf}(\theta; \mathcal{D}) = \frac{1}{N} \sum_{i=1}^N (\hat{c}(\mathbf{x}_i, \theta) - c^*(\mathbf{x}_i, y_i^*))^2$$

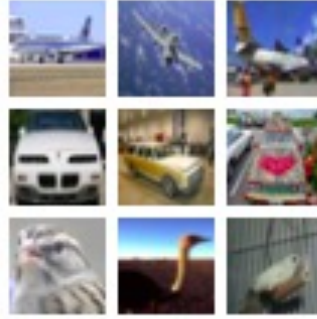
Results



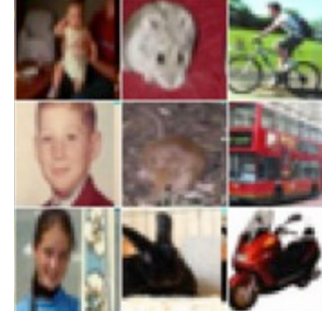
MNIST



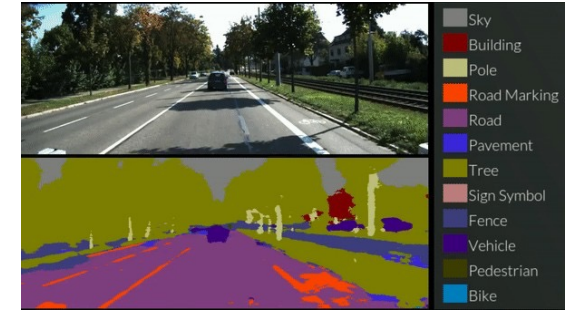
SVHN



CIFAR-10



CIFAR-100



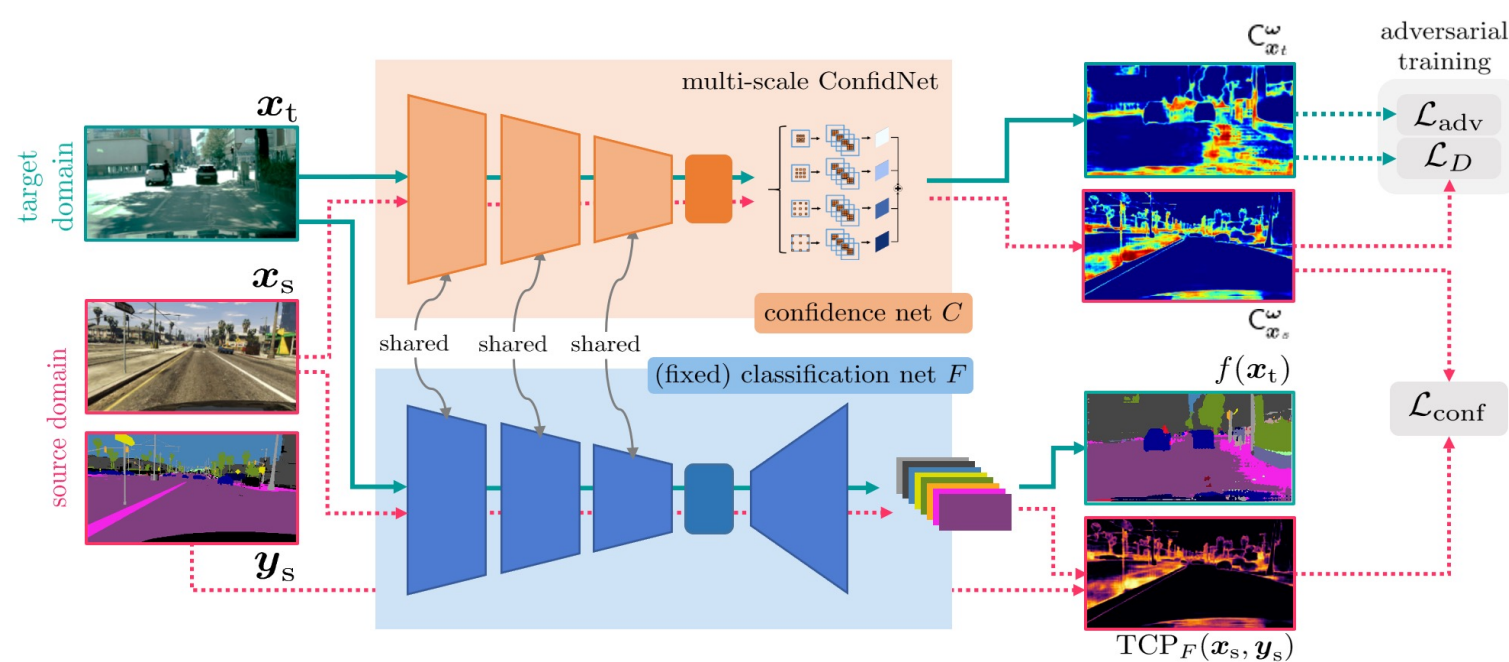
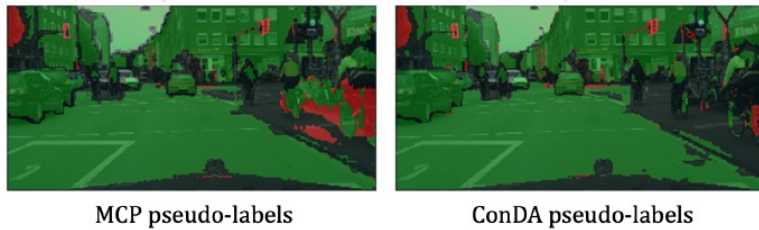
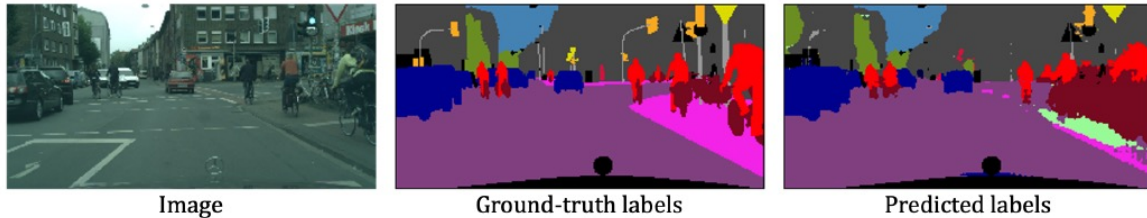
CamVid

| | MNIST | | SVHN | CIFAR-10 | CIFAR-100 | CamVid |
|--|-------------------|-------------------|-------------------|-------------------|-------------------|-------------------|
| | <i>MLP</i> | <i>LeNet-5</i> | <i>LeNet-5</i> | <i>VGG-16</i> | <i>VGG-16</i> | <i>SegNet</i> |
| MCP [Hendrycks & Gimpel, 2017] | 47.3 ± 1.7 | 36.1 ± 3.6 | 46.2 ± 0.5 | 48.4 ± 0.7 | 71.3 ± 0.4 | 48.5 ± 0.3 |
| MC Dropout [Gal et Ghahramani, 2015] | 41.0 ± 1.2 | 42.1 ± 5.5 | 45.2 ± 1.3 | 48.1 ± 1.0 | 71.9 ± 0.7 | 49.4 ± 0.3 |
| TrustScore [Jiang et al., 2019] | 52.1 ± 1.8 | 33.5 ± 3.8 | 44.8 ± 1.3 | 41.8 ± 2.0 | 66.8 ± 0.5 | 20.4 ± 1.0 |
| ConfidNet | 59.7 ± 1.9 | 45.5 ± 3.8 | 48.6 ± 1.0 | 53.7 ± 0.6 | 73.6 ± 0.6 | 50.5 ± 0.3 |

AP errors (%)

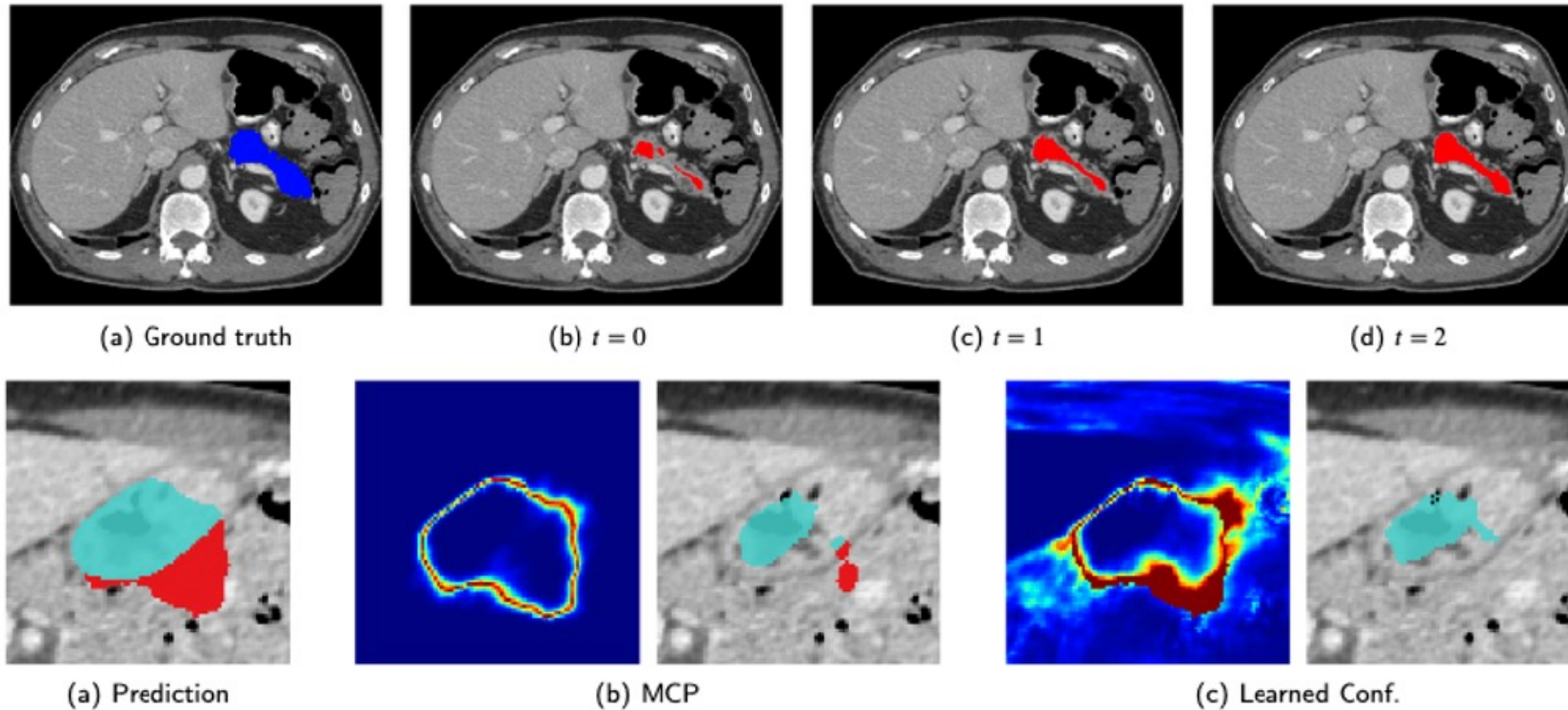
Learning confidence for self-labelling

- Extension for domain adaptation [CTS+21]



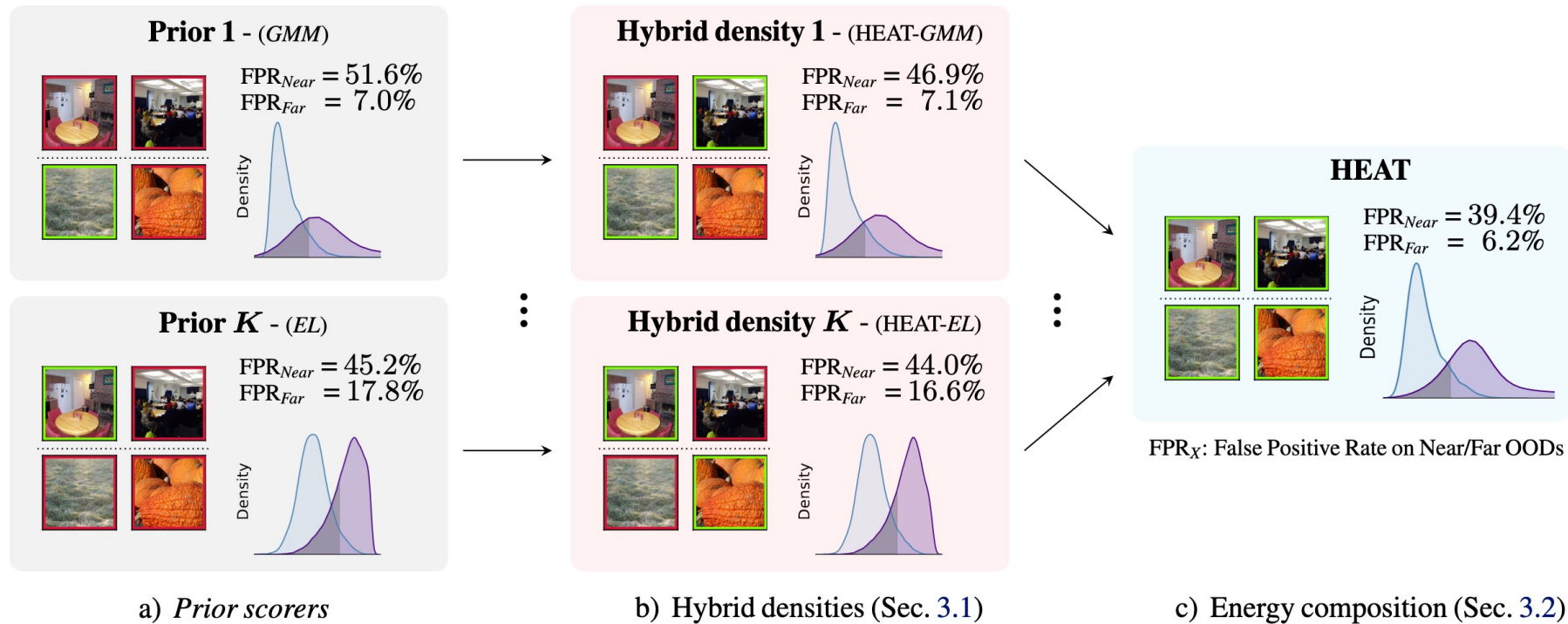
Learning confidence for self-labelling

- Extension for Medical image segmentation [PTS21]



Out-Of-Distribution (OOD) detection

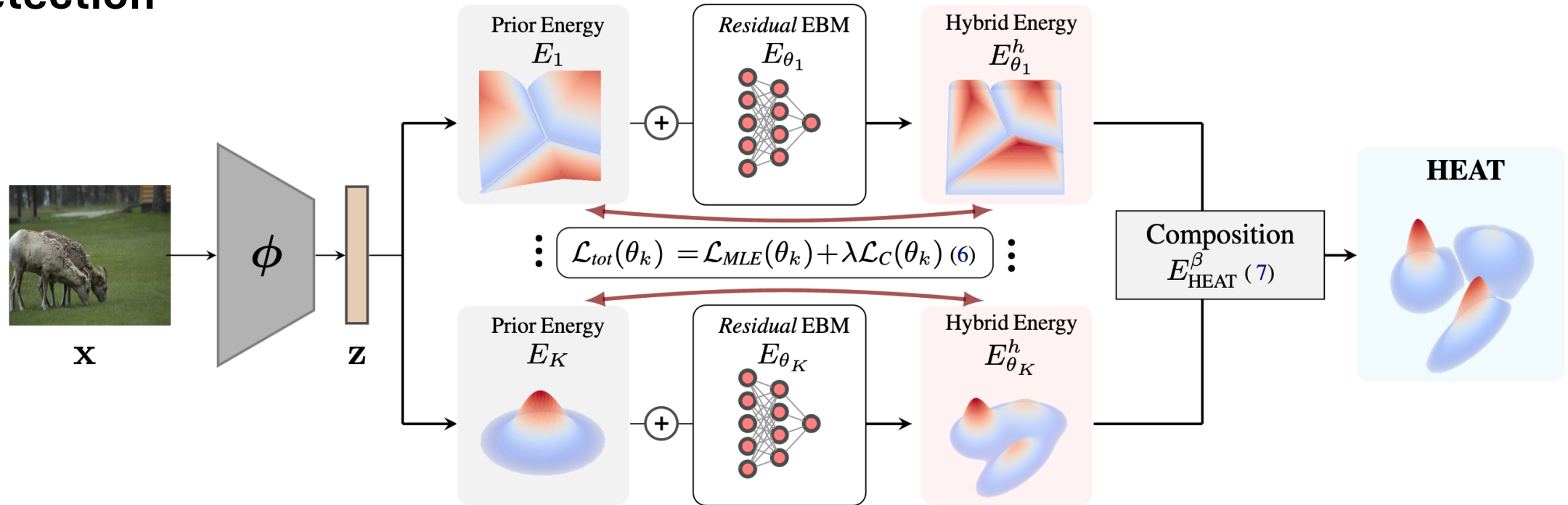
- Post-hoc OOD detection: leveraging any state-of-the-art prediction model
- Accurate OOD detection \Leftrightarrow accurate in-distribution (ID) density estimation
 - State-of-the-art ID density estimation: prior densities, e.g., GMM, Energy Logits (EL)



- Prior density: not accurate \Rightarrow **Energy correction**
- GMM good for far-OOD, EL for near-OOD \Rightarrow **Energy composition**

OOD detection

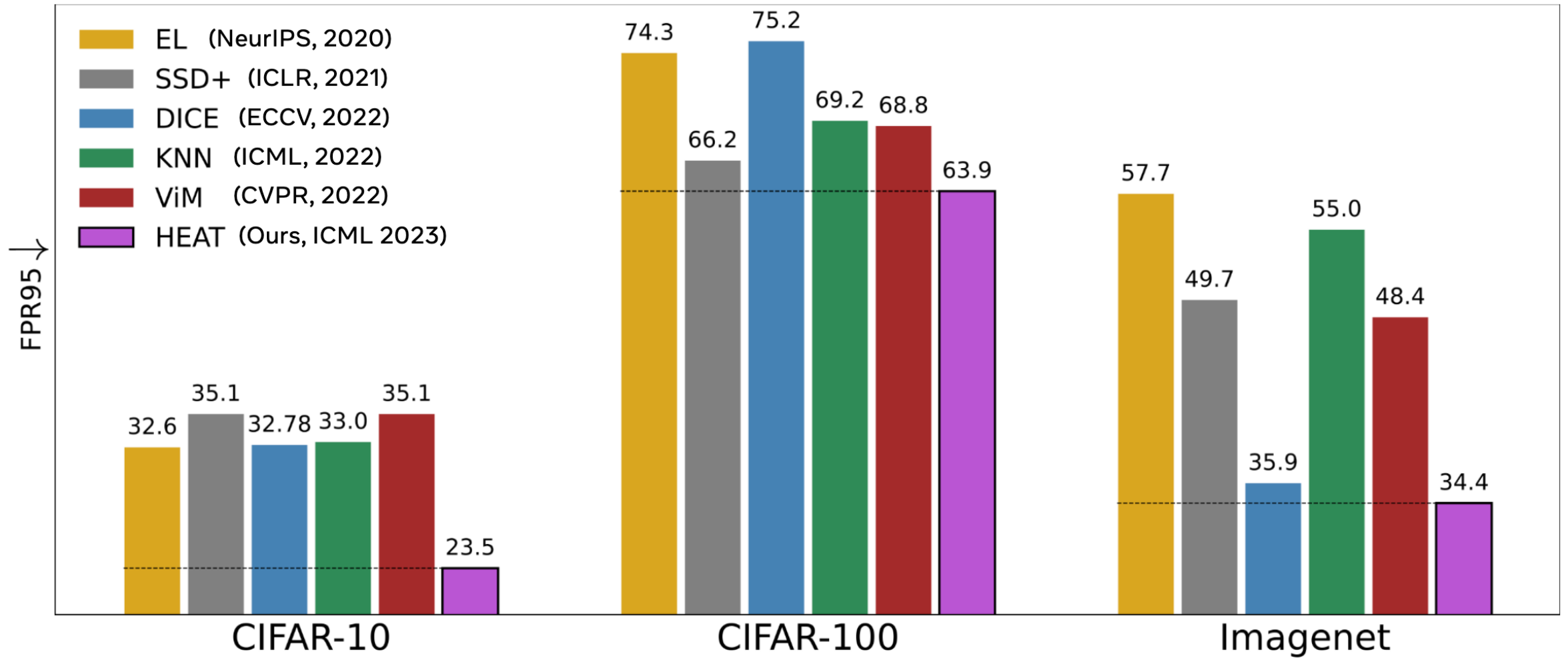
- **HEAT [LRR+23]: Hybrid Energy Based Model (EBM) in the feature space for OOD detection**



- **Energy-based correction** of prior energy terms, e.g. Gaussians

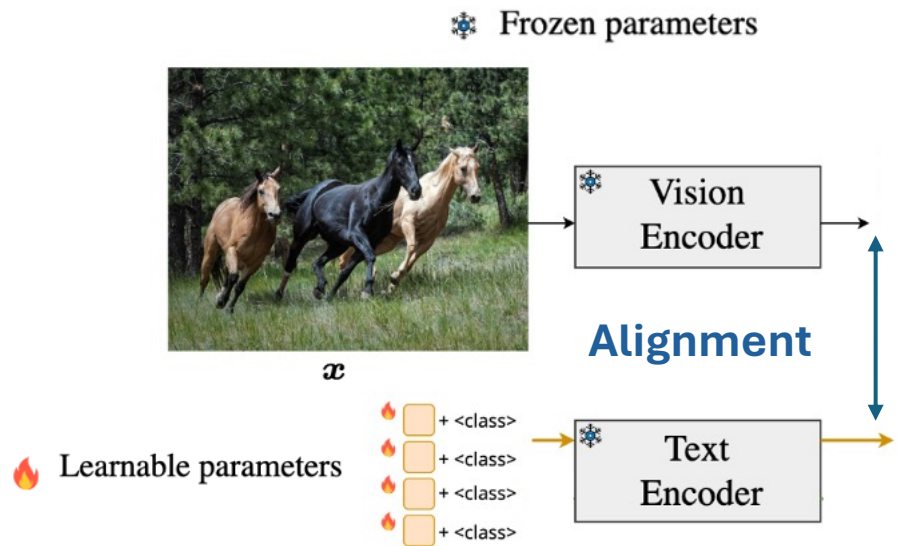
- **Energy composition** of several terms (Gaussian, Energy Logits, std for style)

Results



Robustness in prompt learning

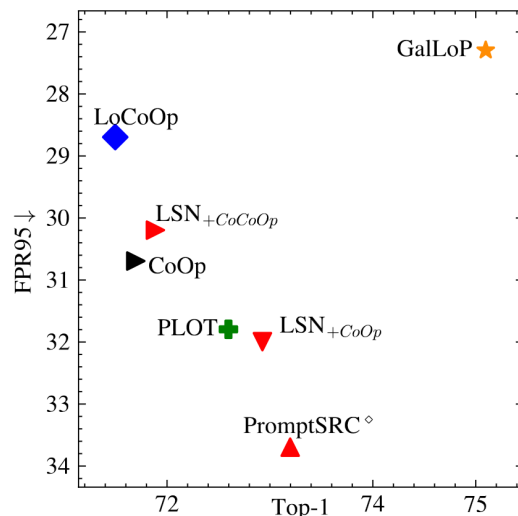
- Learning prompts from frozen vision-language models (VLMs), e.g., CLIP



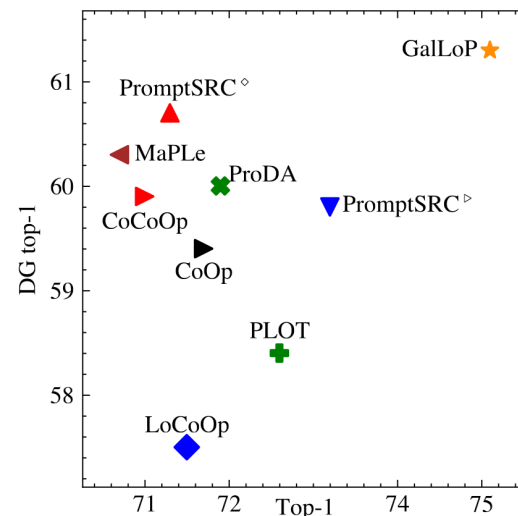
State-of-the-art methods' shortcomings:

- Local prompt: aligning local features
- Accuracy and robustness, e.g., OOD detection, domain generalization

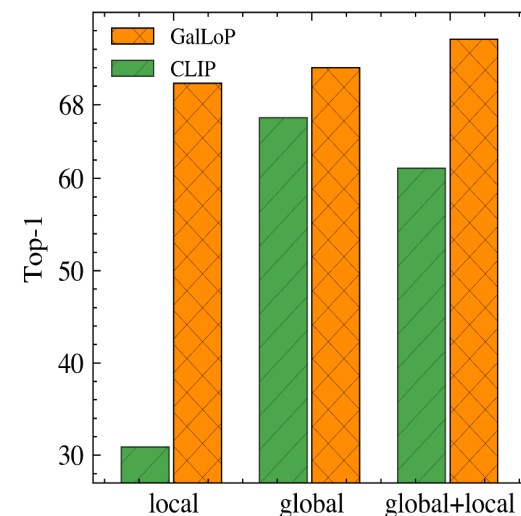
Our method: Global and Local Prompts for VLMs GalLoP



(a) OOD detection

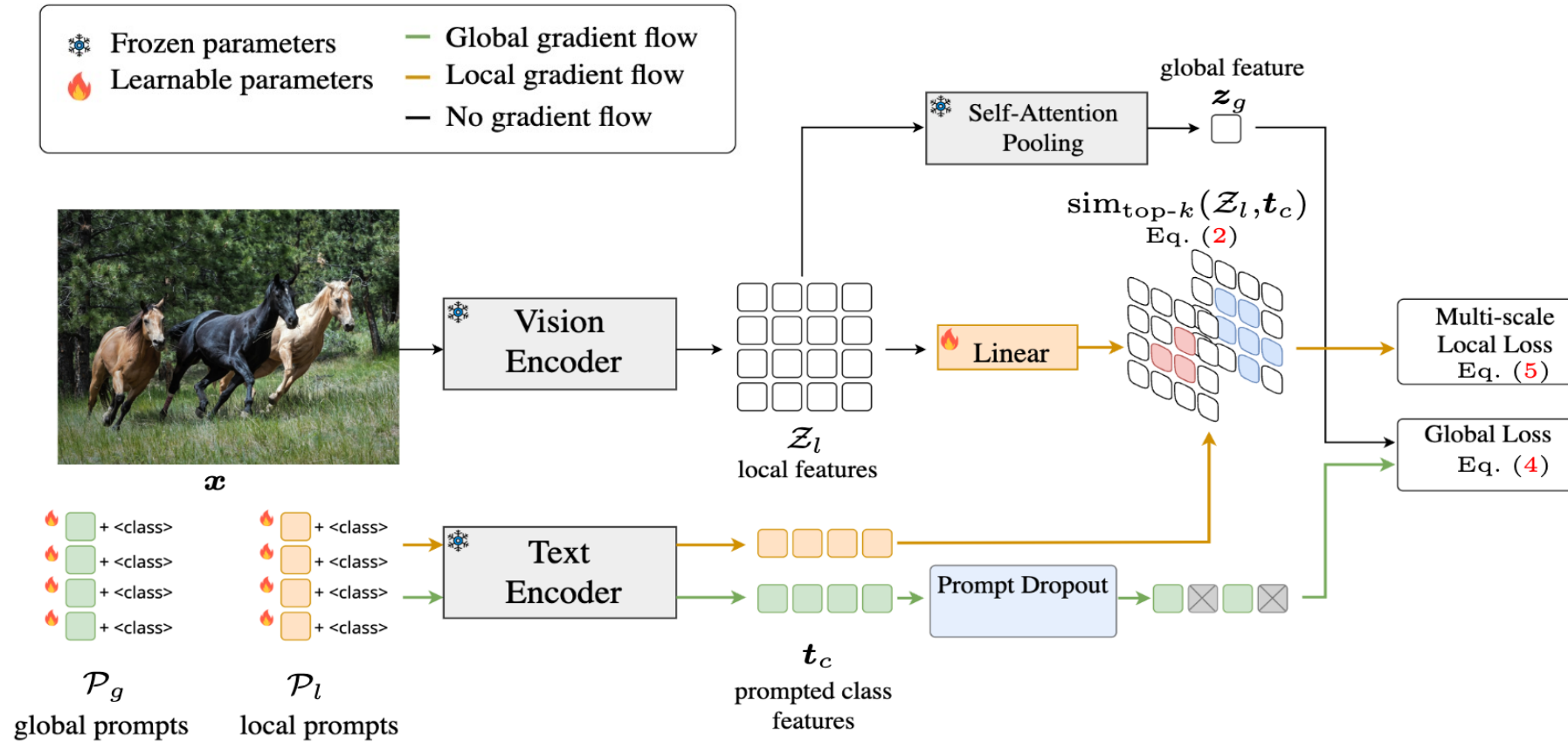


(b) Domain Generalization



(c) Global-local

Learning Global and Local Prompts for VLMs (GalLoP)



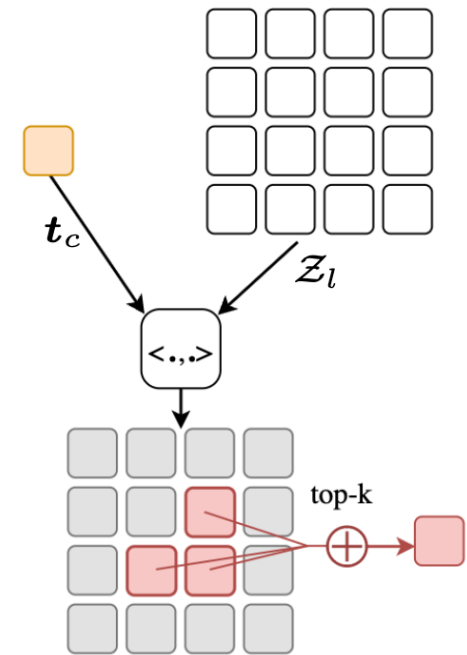
- **Local prompts:** sparse local matching + linear alignment
- **Global prompts diversity:** prompt dropout = multiscale

Local prompts with GalLoP

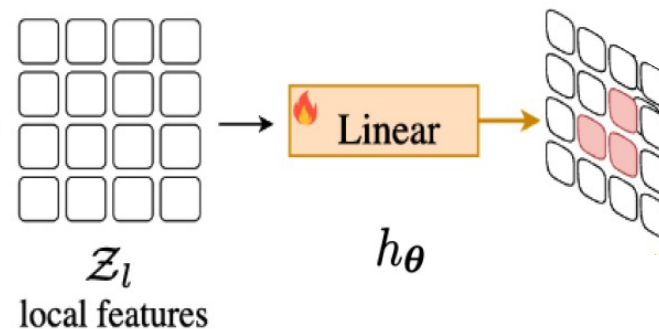
- **Sparse local alignment**

$$\text{sim}_{\text{top-}k}(\mathcal{Z}_l, \mathbf{t}_c) := \frac{1}{k} \sum_{i=1}^L \mathbb{1}_{\text{top-}k}(i) \cdot \langle \mathbf{z}_i^l, \mathbf{t}_c \rangle$$

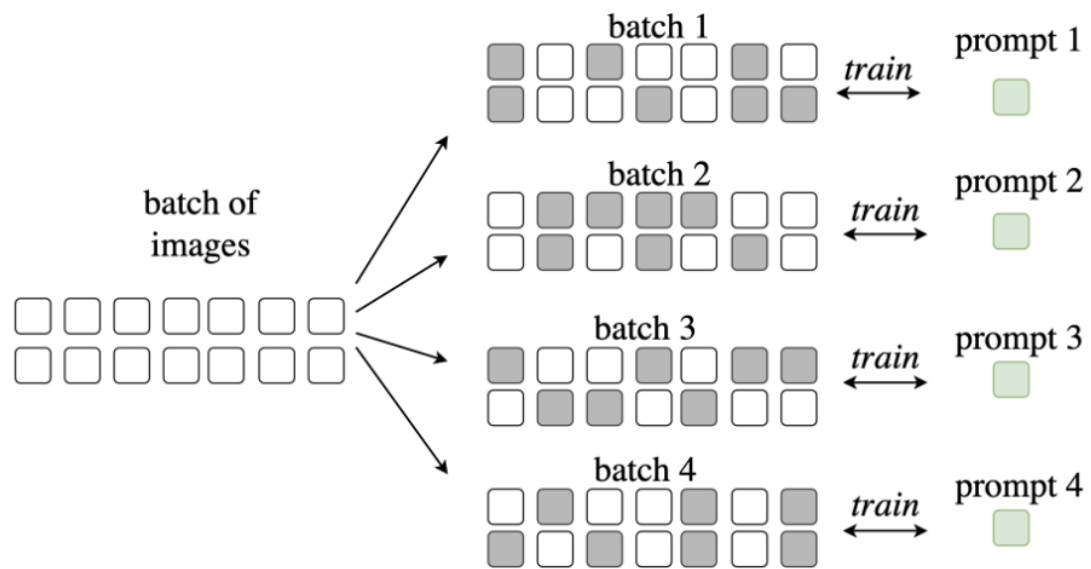
$$\text{where } \mathbb{1}_{\text{top-}k}(i) = \begin{cases} 1 & \text{if } \text{rank}_i(\langle \mathbf{z}_i^l, \mathbf{t}_c \rangle) \leq k, \\ 0 & \text{otherwise.} \end{cases}$$



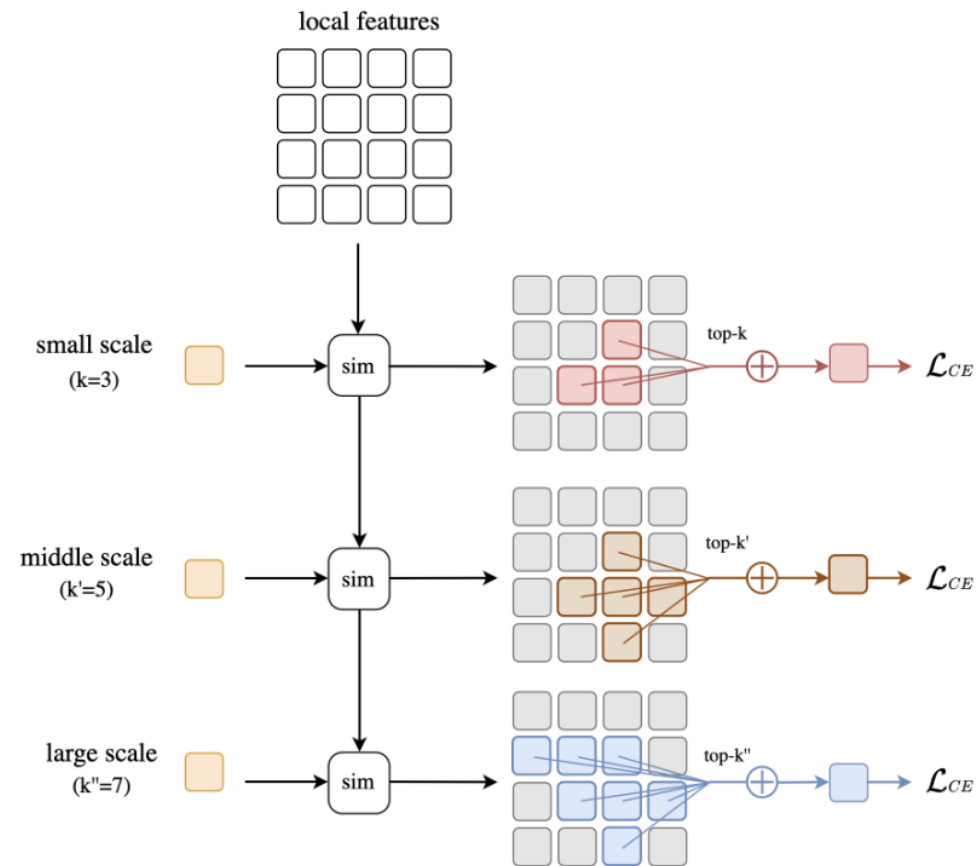
- **Learned alignment**



Prompts diversity with GalLoP

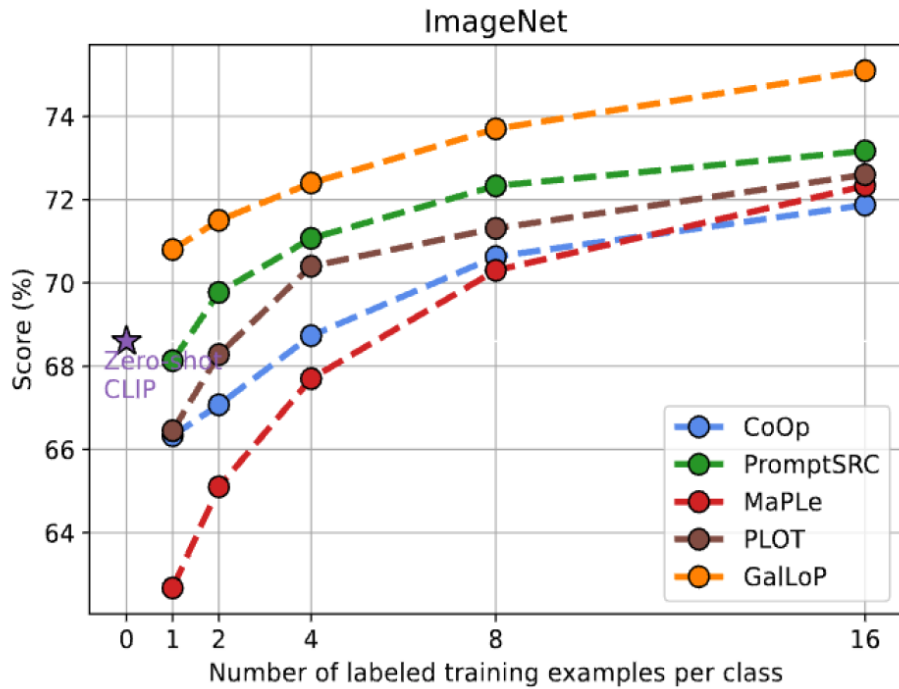


(a) Prompt dropout



(b) Multiscale loss

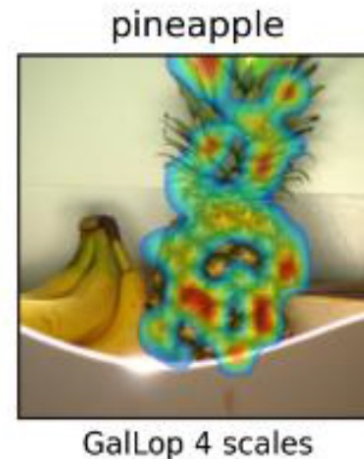
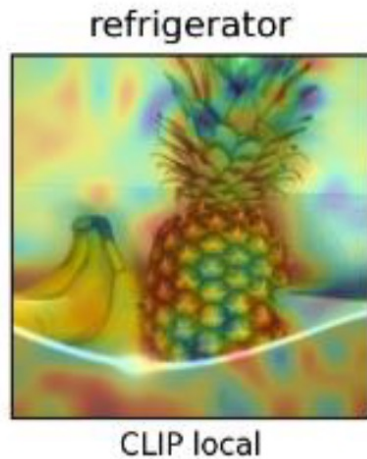
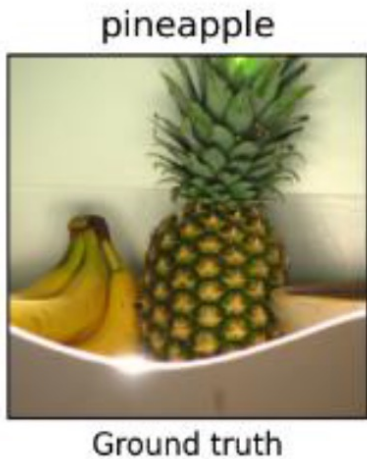
GalLoP Results



| | Top-1 | DG | FPR95↓ | AUC |
|--------------------------|-------------|-------------|-------------|-------------|
| CLIP _{Global} | 66.6 | 57.2 | 42.8 | 90.8 |
| CLIP _{Local} | 12.5 | 9.49 | 73.3 | 73.7 |
| CLIP _{GL} | 61.1 | 49.3 | 35.5 | 90.8 |
| CoOp _{Global} | 71.4 | 59.2 | 39.1 | 91.1 |
| CoOp _{Local} | 41.2 | 30.1 | 65.2 | 78.3 |
| CoOp _{GL} | 69.5 | 55.6 | 33.7 | 90.5 |
| GalLoP _{Global} | 72.0 | 60.4 | 37.0 | 91.7 |
| GalLoP _{Local} | 70.9 | 54.1 | 36.0 | 90.1 |
| GalLoP | 75.1 | 61.3 | 27.3 | 93.2 |

Improved local prompts => effective combination with global prompts

Better accuracy and robustness (OOD detection and domain generalization)



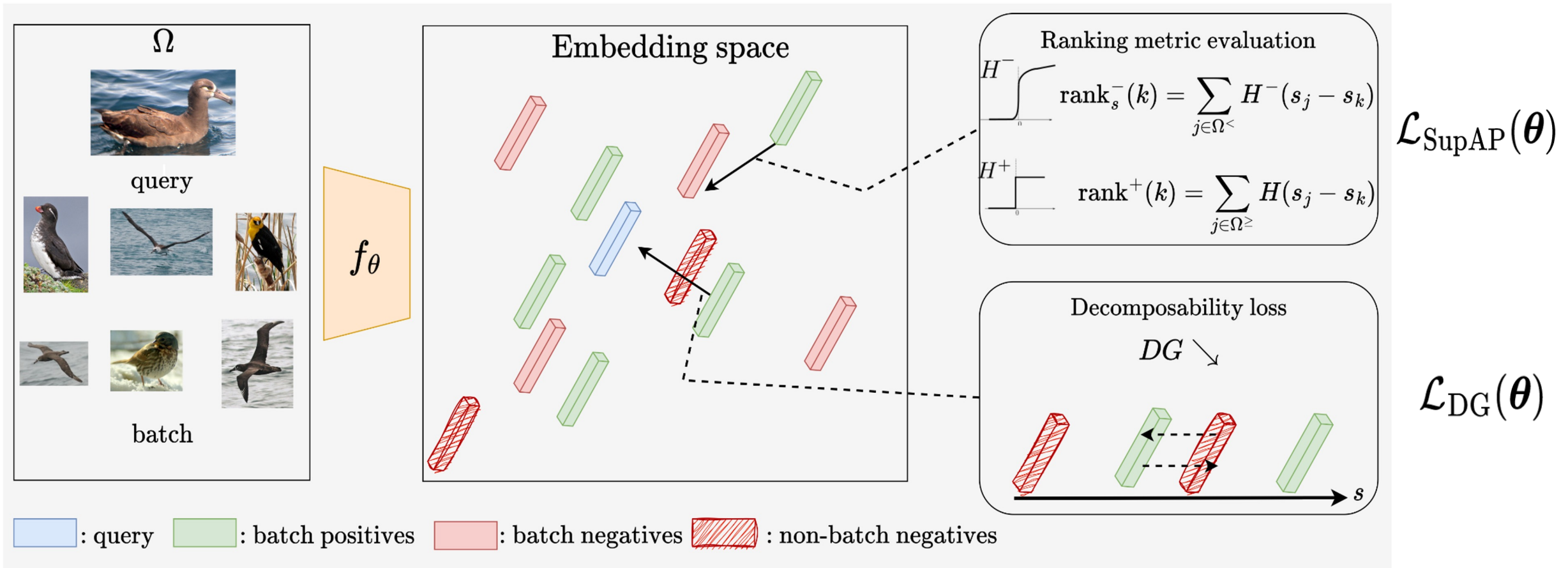
CLIP vs GalLoP with local features

Robustness: recent contributions

1. Uncertainty quantification
2. Direct optimization of rank losses
3. Robustness

Direct optimization of rank losses for image retrieval

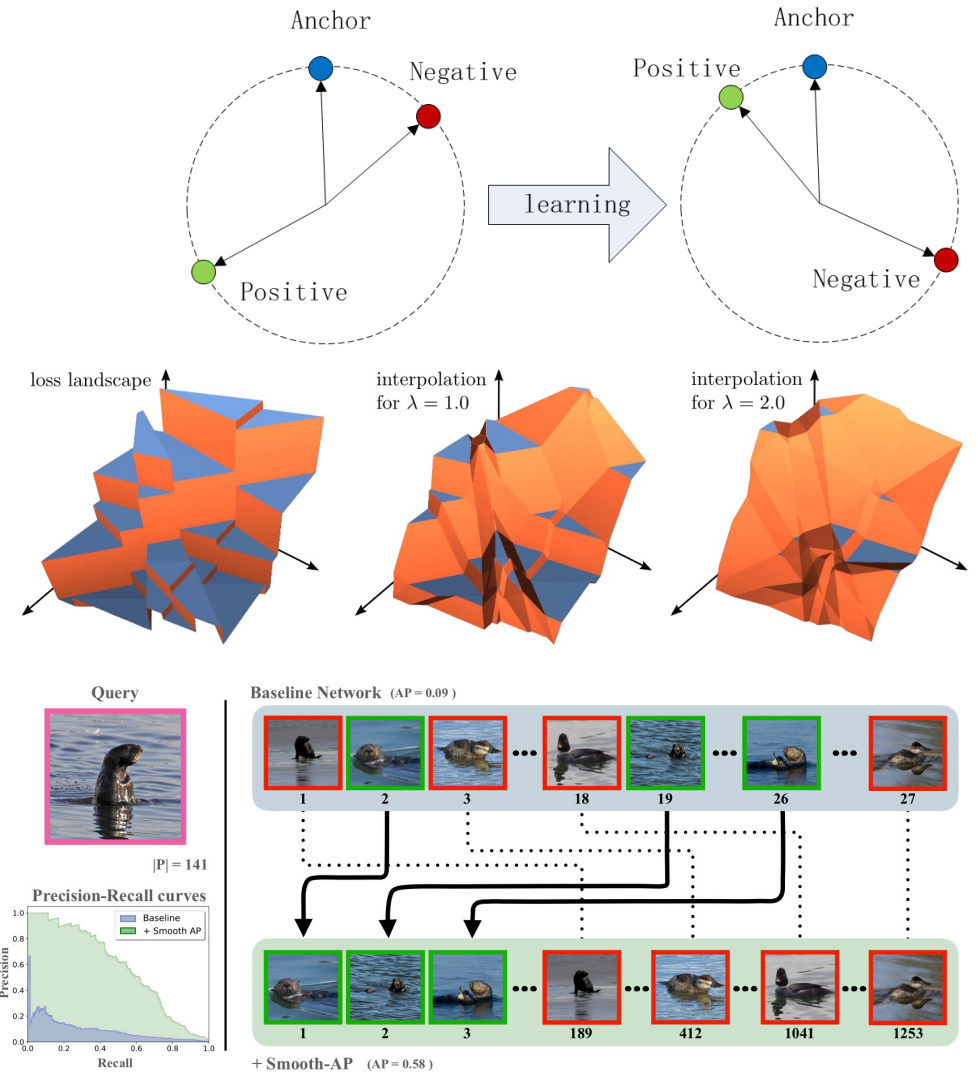
1. Theoretically sound surrogates for non-differentiable rank losses, e.g., Average Precision (AP)
2. Reducing the decomposability gap



$$\mathcal{L}_{\text{ROADMAP}}(\theta) = (1 - \lambda) \cdot \mathcal{L}_{\text{SupAP}}(\theta) + \lambda \cdot \mathcal{L}_{\text{DG}}(\theta)$$

Image retrieval: non-smooth metrics & losses

- Standard losses, e.g., : triplet loss, NSM [A]
 - ⊖ Coarse upper-bounds, not well-aligned with metrics: supports bottom vs. top of the ranking
- Upper bounds: structural SVMs, Blackbox optim [B]
 - ⊕ General methods, theoretical guarantees
 - ⊖ Coarse upper bounds
- Rank approximation: binning approaches, smoothAP [C,D]
 - ⊕ Tighter approximations
 - ⊖ No theoretical guarantees



[A] A. Zhai, and H.Y Wu. Classification is a strong baseline for deep metric learning. BMVC 2018

[B] M. Rolínek, V. Musil, A. Paulus, M. Vlastelica, C. Michaelis, G. Martius. Optimizing rank-based metrics with blackbox differentiation. CVPR 2020

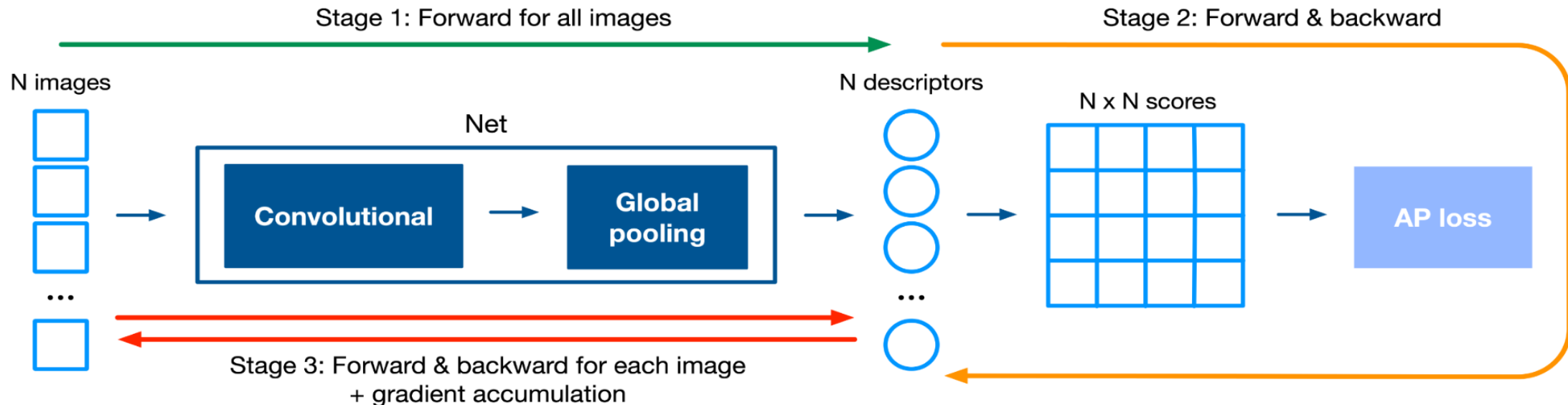
[C] A. Brown, W. Xie, V. Kalogeiton, A. Zisserman. Smooth-ap: Smoothing the path towards large-scale image retrieval. ECCV 2020

[D] Y. Patel, G. Tolas, and J. Matas, "Recall@ k surrogate loss with large batches and similarity mixup," in CVPR, 2022

Image retrieval: addressing non-decomposability

Fewer works, brute-force approaches

- Sampling informative batches or constraints in batch
- Storing the datasets, e.g., x-batch memory [D]: **increased in memory**
- Large batches + 2-step approach for AP and back-prop [E] : **increase in training time**

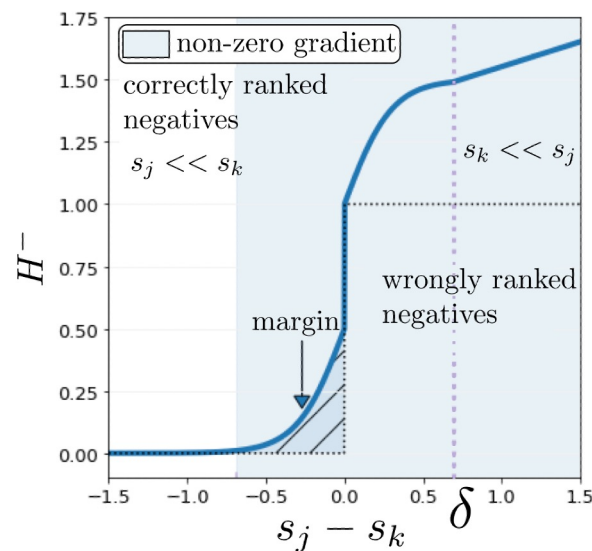
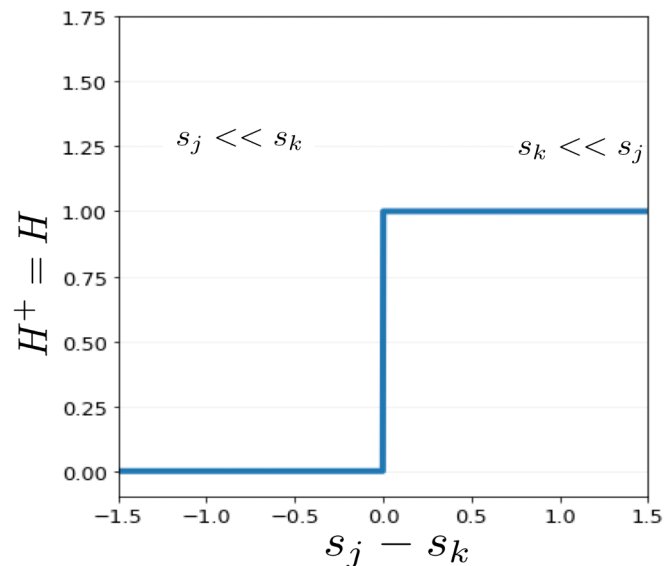


[D] X. Wang, H. Zhang, W. Huang, and M. R. Scott, "Cross-batch memory for embedding learning," in CVPR, 2020

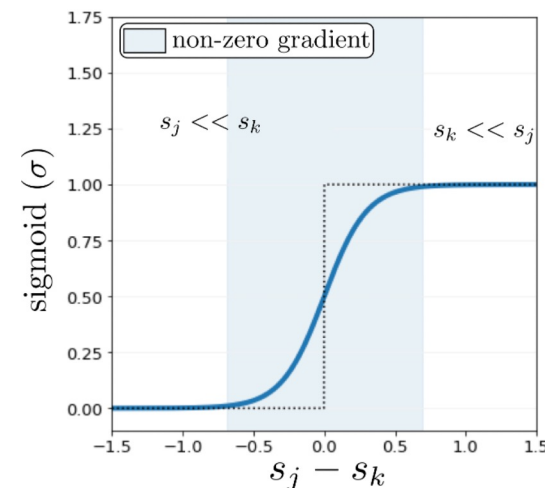
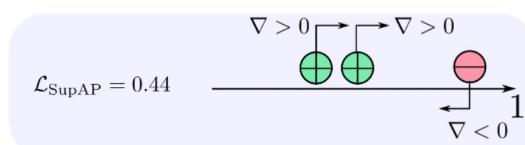
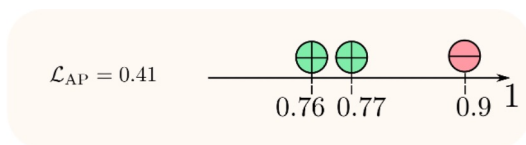
[E] J. Revaud, J. Almazan, R. S. Rezende, and C. R. d. Souza, "Learning with average precision: Training image retrieval with a listwise loss," in ICCV, 2019.

Robust and decomposable AP (ROADMAP)

$$\text{AP} = \frac{1}{|\Omega^+|} \sum_{k \in \Omega^+} \frac{\text{rank}^+(k)}{\text{rank}(k)} = \frac{1}{|\Omega^+|} \sum_{k \in \Omega^+} \frac{\text{rank}^+(k)}{\text{rank}^+(k) + \text{rank}^-(k)} \quad \text{rank}(k) = 1 + \sum_{j \in \Omega} H(s_j - s_k)$$



- Optimizing rank⁻ → smooth approximation & **upper bound of AP**
- Not optimizing rank⁺ → well-behaved gradients.



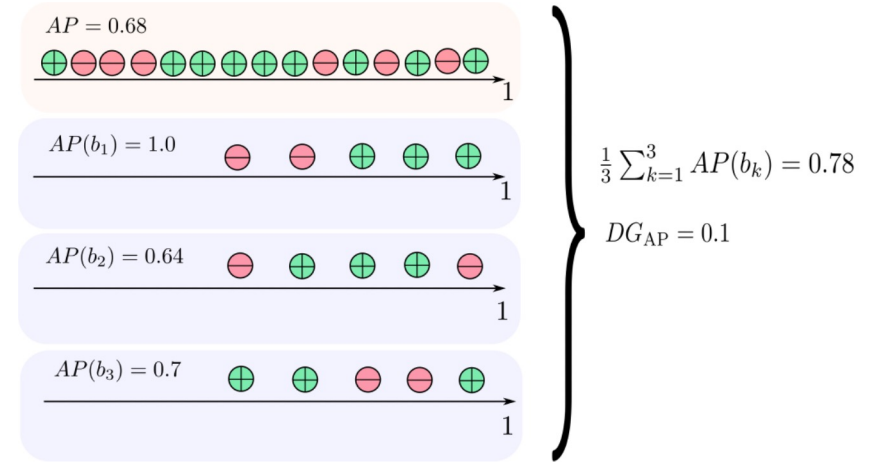
Better than

$$\mathcal{L}_{\text{SupAP}} = 1 - \frac{1}{|\Omega^+|} \sum_{k \in \Omega^+} \frac{\text{rank}^+(k)}{\text{rank}^+(k) + \text{rank}_s^-(k)}$$

Improving decomposability

Decomposability gap:

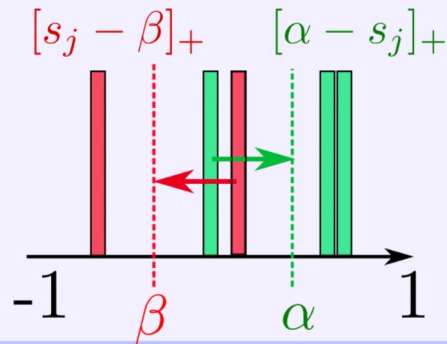
$$DG(\Omega) = \frac{1}{|\mathcal{B}|} \sum_{b \in \mathcal{B}} \mathcal{M}(b) - \mathcal{M}(\Omega)$$



Decomposability

$DG \rightarrow 0$

\mathcal{L}_{DG}

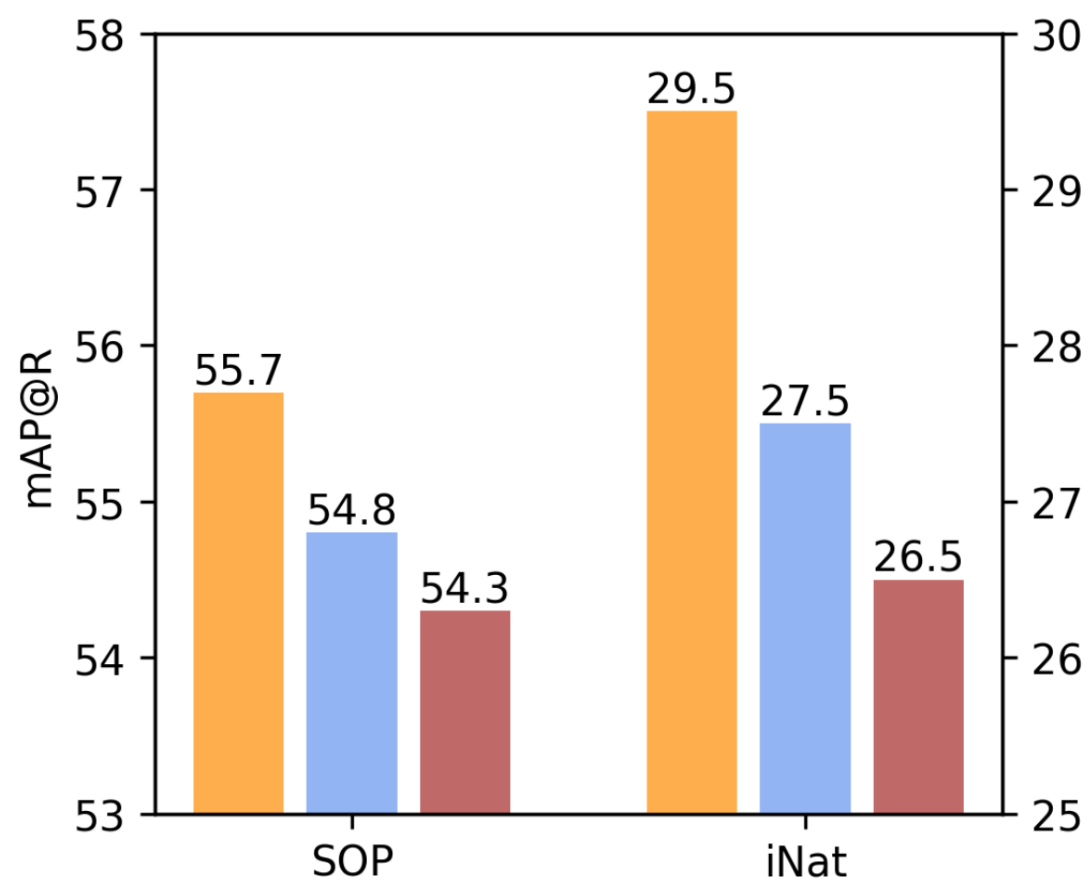
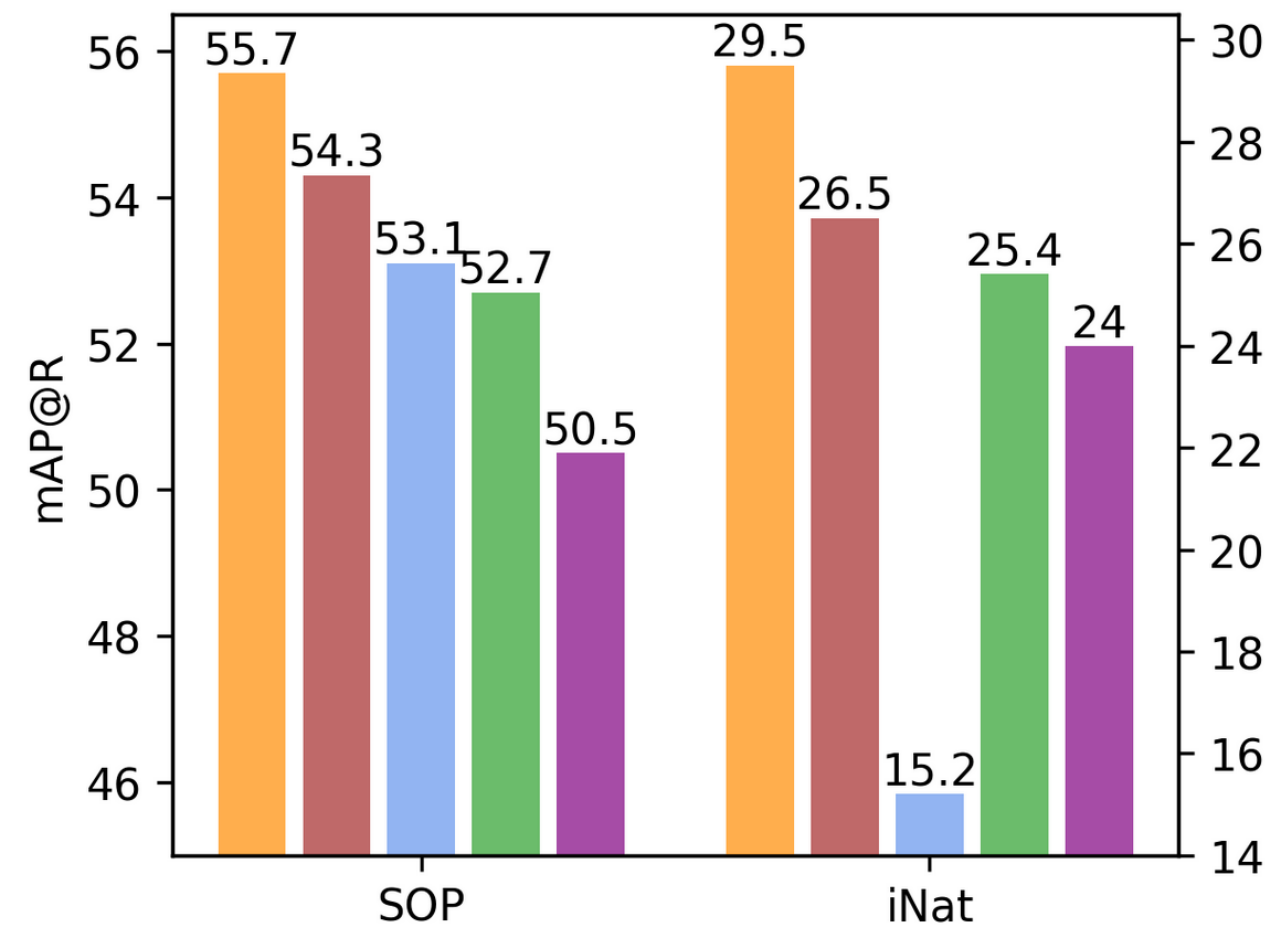
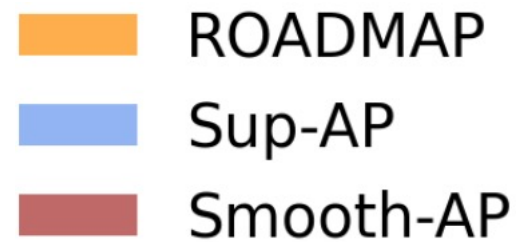


$$\mathcal{L}_{DG}(\theta) = \frac{1}{|\Omega^+|} \sum_{x_j \in \Omega^+} [\alpha - s_j]_+ + \frac{1}{|\Omega^-|} \sum_{x_j \in \Omega^-} [s_j - \beta]_+$$

- Calibrates scores across batches
 - Positive scores $\geq \alpha$
 - Negative scores $< \beta$
- **Proof:** \mathcal{L}_{DG} reduces decomposability gap

$$\mathcal{L}_{ROADMAP}(\theta) = (1 - \lambda) \cdot \mathcal{L}_{SupAP}(\theta) + \lambda \cdot \mathcal{L}_{DG}(\theta)$$

Results



Robustness: recent contributions

1. Uncertainty quantification
2. Direct optimization of rank losses
3. Controlling mistake severity

Hierarchical Image Retrieval for Robust Ranking

- Binary image retrieval → do not take into account mistake severity
- HAPPIER: Hierarchical Average Precision training for Pertinent Image Retrieval
Extending AP to graded setting to take importance of errors into account

HAPPIER

rank 1

rank 2

rank 3

rank 4

rank 5

rank 6

\mathcal{H} -AP AP

0.94 0.9



Query image



Baseline

\mathcal{H} -AP AP

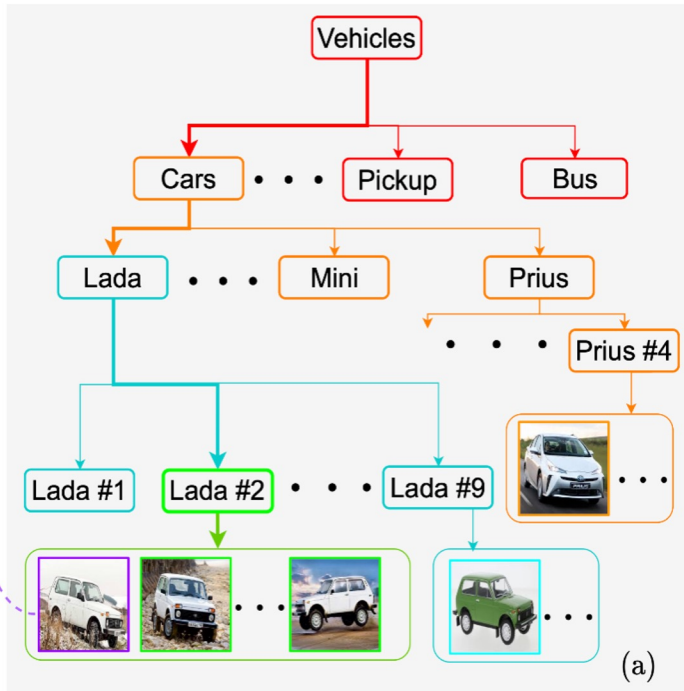
0.68 0.9



Relevance function: graded similarities



Query Image: Lada #2



- **Relations between categories** → proxy for mistake severity.
- **Decreasing function of the distance** in the hierarchical tree.

$$\text{rel}(k) = \frac{l/L}{|\Omega^{(l)}|}$$

- l : level of the closest ancestor in the tree.
- L total number of levels.

Hierarchical average precision (\mathcal{H} -AP)

$$\mathcal{H}\text{-rank}(k) = \text{rel}(k) + \sum_{j \in \Omega^+} \min(\text{rel}(k), \text{rel}(j)) \cdot H(s_j - s_k)$$



| | | | | | |
|---|---|---|---|---|------------------------|
| 1 | 2 | 3 | 4 | 5 | |
| | | | | | |
| | | | | | Lada #2 rel := 1 |
| | | | | | Lada #9 rel := 2/3 |
| | | | | | Prius #4 rel := 1/3 |
| | | | | | Bus rel := 0 |

$\mathcal{H}\text{-rank}(1) = 1/3 = 1/3$
 $\mathcal{H}\text{-rank}(2) = 1 + (1/3) = 4/3$
 $\mathcal{H}\text{-rank}(4) = 2/3 + (1/3 + 2/3) = 5/3$
 $\mathcal{H}\text{-rank}(5) = 1 + (1/3 + 1 + 2/3) = 3$

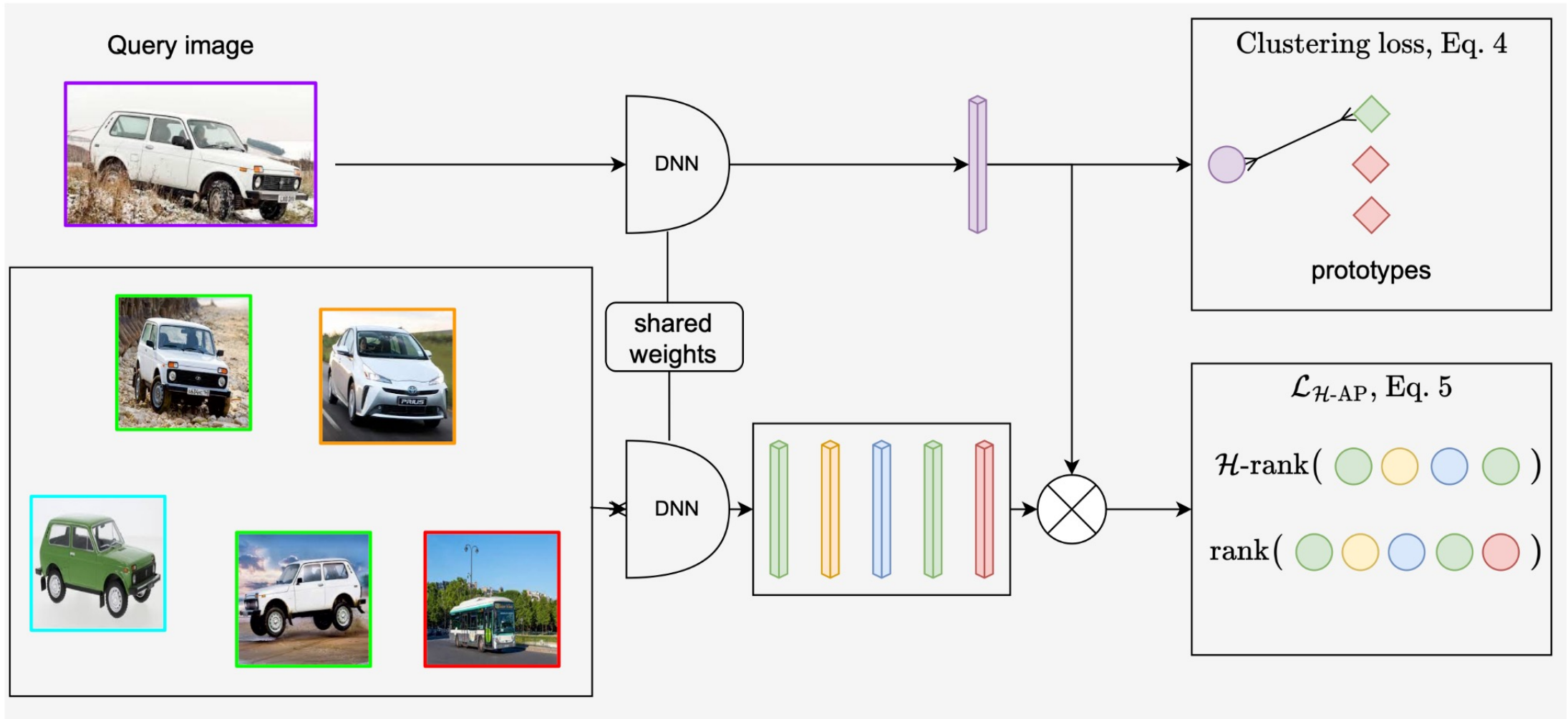
- Errors in ranking \rightarrow **weighted by relevance**
- Correct \mathcal{H} -rank \rightarrow decreasing order of relevance

$$\mathcal{H}\text{-AP} = \frac{1}{\sum_{k \in \Omega^+} \text{rel}(k)} \sum_{k \in \Omega^+} \frac{\mathcal{H}\text{-rank}(k)}{\text{rank}(k)}$$

- Consistent generalization of AP.
- Flexible wrt. the relevance

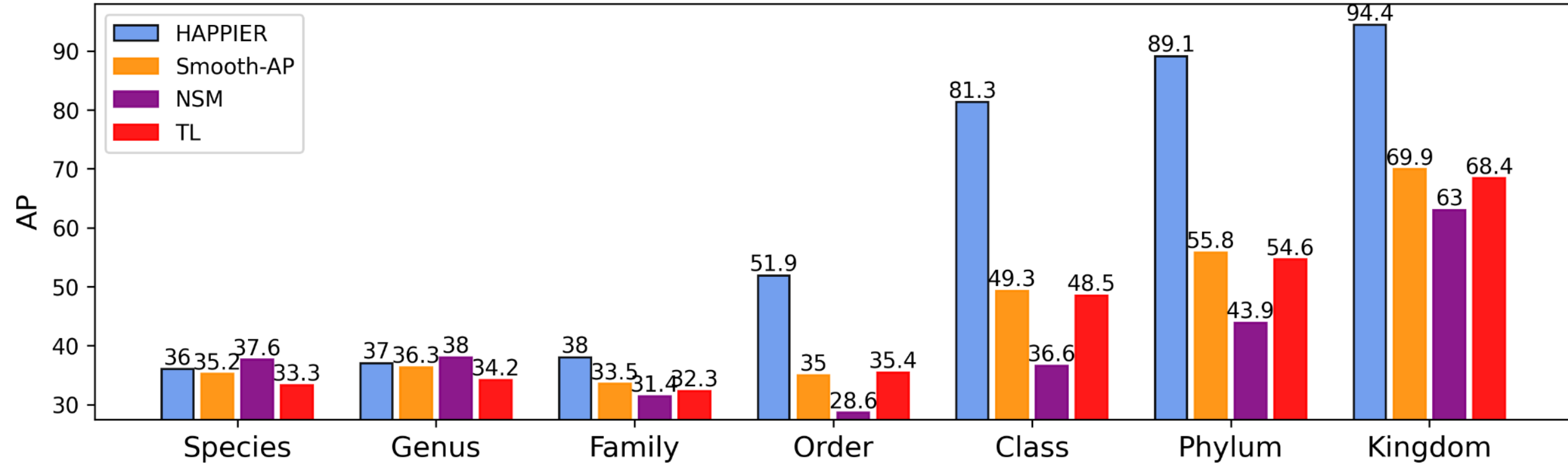


HAPPIER training



$$\mathcal{L}_{\text{Sup-H-AP}}(\theta) = 1 - \frac{1}{\sum_{k \in \Omega^+} \text{rel}(k)} \sum_{k \in \Omega^+} \frac{\mathcal{H}\text{-rank}(k)}{\text{rank}^+(k) + \text{rank}_s^-(k)}$$

Results



- On par for fine-grained retrieval (“Species”)
- **Large gains** on other hierarchical levels from “Family”

\mathcal{H} -GLDv2: a hierarchical landmark dataset

Query:



GLDv2 → large scale
landmarks retrieval
dataset [F]

No hierarchical annotations
→ how difficult is it to create
hierarchical annotations?

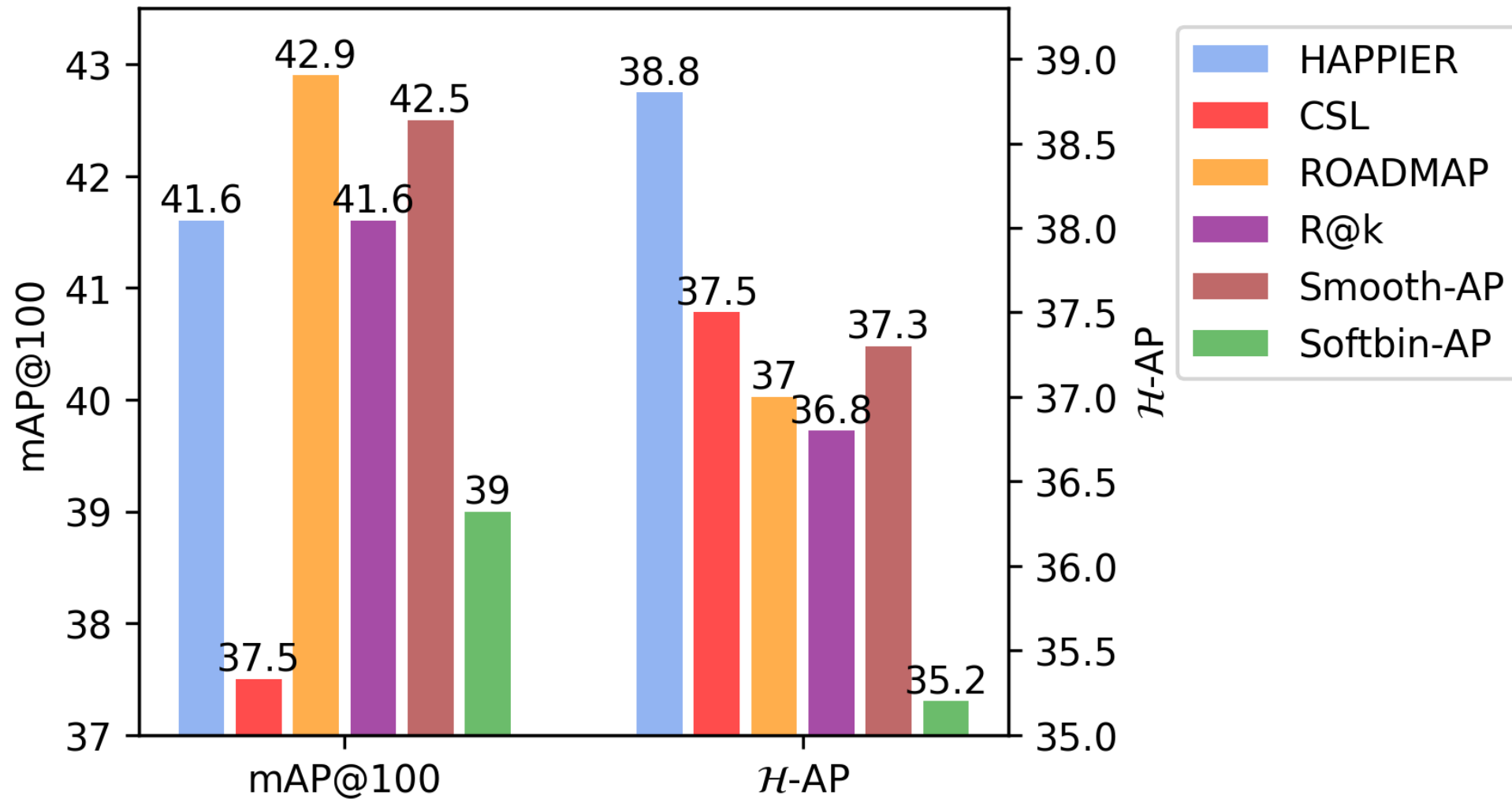
\mathcal{H} -GLDv2

Relevant index images:



1. Scraping Wikimedia Commons
2. Post-processing

Results



Perspectives

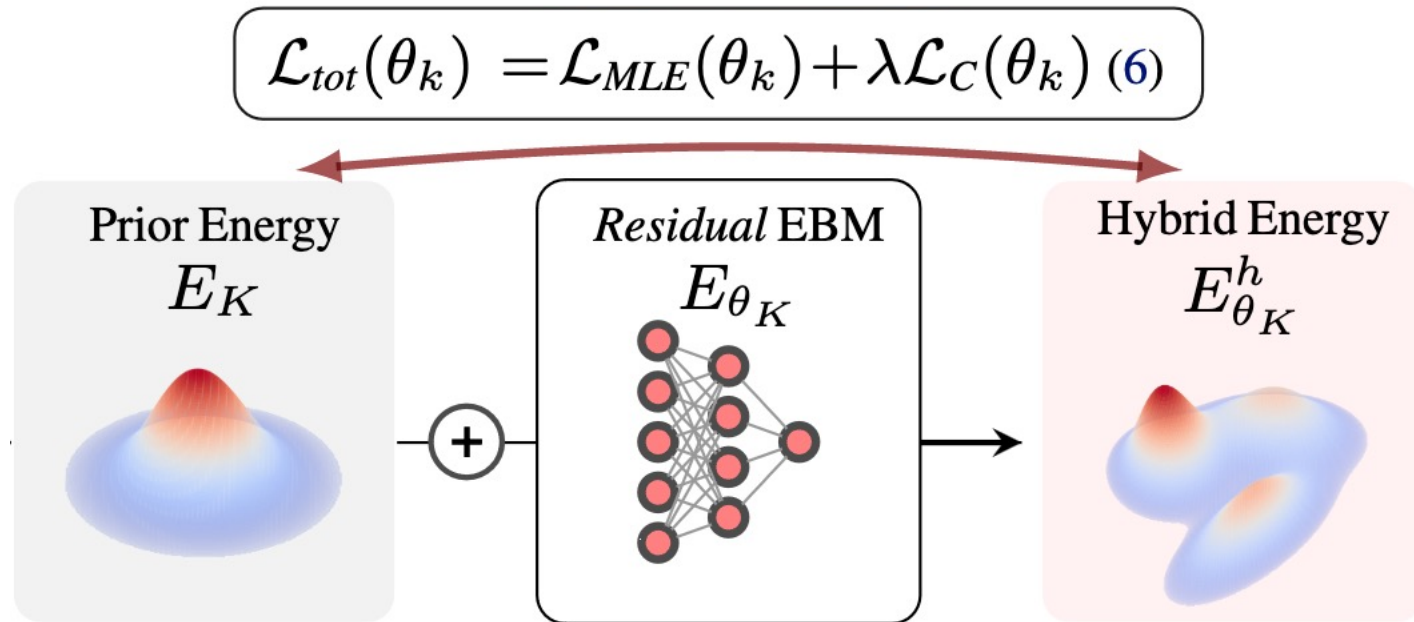
- Uncertainty quantification:
 - Global measure of uncertainty (aleatoric, epistemic)
 - For foundation models, e.g., CLIP
 - Test-time adaptation
- Non-smooth & non-decomposable metrics beyond ranking
- Mistake severity robustness
 - Adaptation to multi-modal models

Thank you for you attention!

- C. Corbière, N. Thome, A. Bar-Hen, M. Cord, P. Pérez. Addressing Failure Detection by Learning Model Confidence. NeurIPS 2019. <https://github.com/valeoai/ConfidNet>
- C. Corbière, N. Thome, A. Saporta, T-H. Vu, M. Cord, P. Pérez. Confidence Estimation via Auxiliary Models. IEEE T-PAMI, vol. 44, no. 10, pp. 6043-6055, June 2021.
- O. Petit, N. Thome, L. Soler. 3D Spatial Priors for Semi-Supervised Organ Segmentation with Deep Convolutional Neural Networks. IJCARS, Springer Verlag, In press, 2021.
- E. Ramzi, N. Thome, C. Rambour, N. Audebert, X. Bitot. “Robust and Decomposable Average Precision for Image Retrieval.” *NeurIPS* 2021 <https://github.com/elias-ramzi/ROADMAP>
- E. Ramzi, , N. Audebert, C. Rambour, N. Thome, X. Bitot. “Hierarchical Average Precision Training for Pertinent Image Retrieval.” *ECCV*, 2022. <https://github.com/elias-ramzi/HAPPIER>
- Lafon, Marc, E. Ramzi, C. Rambour, N. Thome. “Hybrid Energy Based Model in the Feature Space for Out-of-Distribution Detection.” *ICML*, 2023. <https://github.com/MarcLafon/heatood>
- Lafon, Marc, E. Ramzi, C. Rambour, N. Audebert. N. Thome. “GalLoP: Learning Global and Local Prompts for Vision-Language Models.” *ECCV*, 2024.
- Ramzi, Elias, et al. “Optimization of Rank Losses for Image Retrieval.” *under-review TPAMI*, 2024. <https://github.com/cvdfoundation/google-landmark>

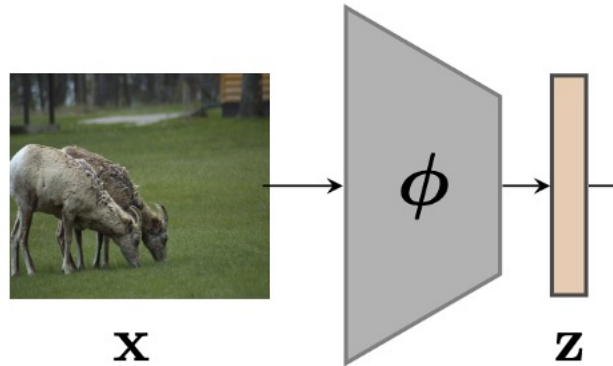
HEAT: Energy correction

- Hybrid energy: $E_{\theta_k}^h(\mathbf{z}) = E_{q_k}(\mathbf{z}) + E_{\theta_k}(\mathbf{z})$ $p_{\theta_k}^h(\mathbf{z}) = \frac{1}{Z(\theta_k)} \exp(-E_{\theta_k}^h(\mathbf{z}))$
- Controlling the residual $\mathcal{L}_C(\theta_k) = \mathbb{E}_{p_{in}, p_{\theta_k}^h} [(E_{\theta_k}^h - E_{q_k})^2]$
 - Correction with mimical norm



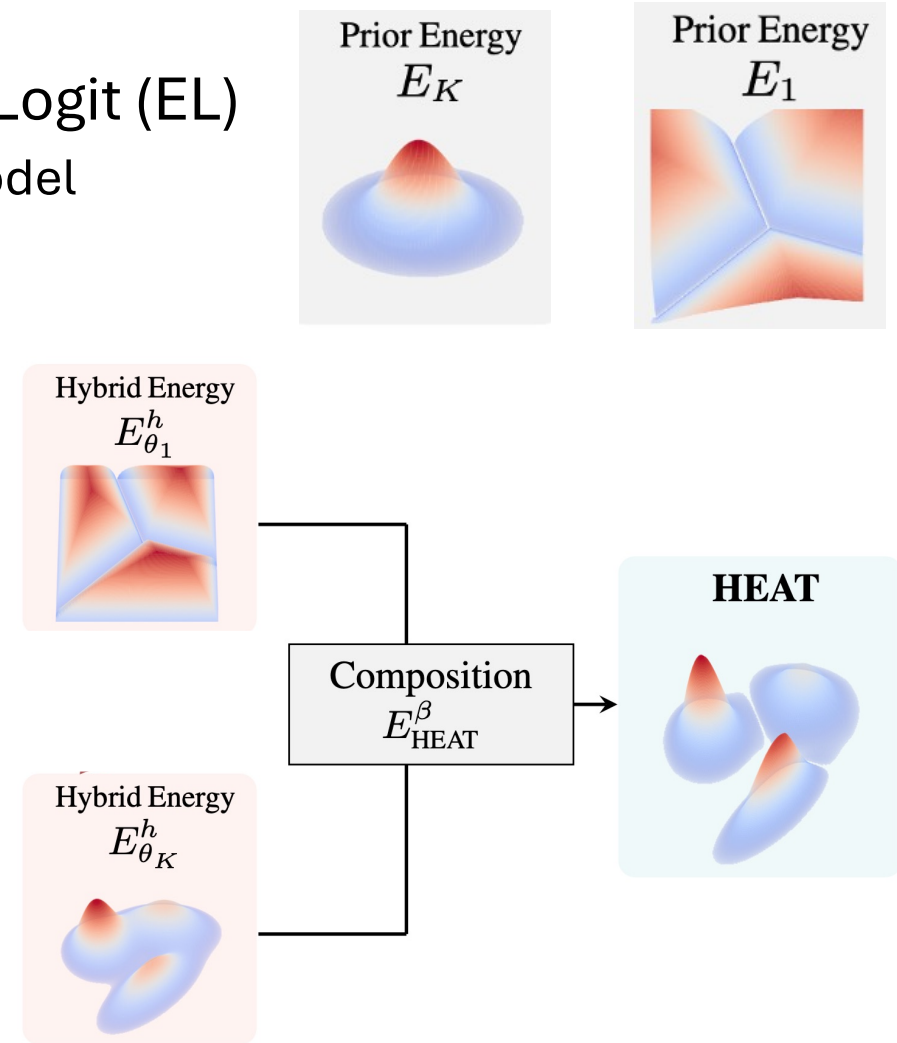
HEAT: Energy composition

- Prior energy function, e.g. Gaussians or Energy Logit (EL)
 - Avg or Std (style) features as inputs for the energy model



- Energy composition:
 - $\beta \rightarrow +/\infty, \max/$
 - $\beta = -1, \text{logsumexp}$

$$E_{\text{HEAT}}^{\beta} = \frac{1}{\beta} \log \sum_{k=1}^K e^{\beta E_{\theta_k}^h}$$



1. Scraping Wikimedia Commons


Wikimedia Commons → largest open database of landmarks.

- GLDv2 sourced from wikimedia commons
- « Instance of » => super-category

Ex of scraped labels:

- Church building.
- Church building (1172–1954)
- Cathedral
- Castle
- Corsican nature reserve
- New Zealand great walks
- Waterfall
- Arch-gravity dam
- Canal
- Association football venue
- Astronomical observatory
- Village

Tell Qasile [Collapse]
archaeological site in Tel Aviv District, Israel

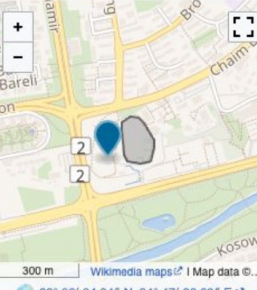


Upload media [↗](#)
Wikipedia

Instance of [archaeological site](#)

Culture Philistines

Location Tel Aviv District, Israel



Authority control [Collapse]
Q2690025
VIAF ID: 243213562
National Library of Israel J9U ID: 987007544134305171
Reasonator · Scholia · PetScan · statistics · WikiMap · Locator tool · KML file · WikiShootMe · OpenStreetMap · Search depicted

St Oswald's Church, Dean [Collapse]
church in Dean, Cumbria, UK



Upload media [↗](#)
Wikipedia

Instance of [church building](#)

Dedicated to Oswald of Worcester

Made from material calciferous sandstone

Location Dean, Allerdale, Cumbria, North West England, England, UK

Architectural style English Gothic architecture
Norman architecture

Diocese Diocese of Carlisle

Heritage designation Grade I listed building (1986–)

Inception 12th century


Religion or worldview Anglicanism

[official website](#)



Authority control [Collapse]
Q4985455
Reasonator · Scholia · PetScan · statistics · WikiMap · Locator tool · KML file · WikiShootMe · OpenStreetMap · Search depicted


Buena Vista Lagoon [Collapse]
lake in United States of America



Upload media [↗](#)
Wikipedia

Instance of [lake](#)

Location California



Authority control [Collapse]
Q4985455
Reasonator · Scholia · PetScan · statistics · WikiMap · Locator tool · KML file · WikiShootMe · OpenStreetMap · Search depicted

Shōmyō Falls [Collapse]
waterfall in Toyama Prefecture, Japan



Upload media [↗](#)
Wikipedia

Instance of [waterfall](#)

Part of Japan's Top 100 Waterfalls (38)

Named after Shōmyō

Located in protected area Chūbu-Sangaku National Park

Location Ashikuraji, Tateyama, Nakanikawa district, Toyama Prefecture, Japan


Located in or next to body of water Shōmyō River

Heritage designation Place of Scenic Beauty of Japan
natural monument

Height 350 m
126 m (Q46868895)

Elevation above sea level 1,260 m

Mosquée de la Divinité [Collapse]
building in Senegal



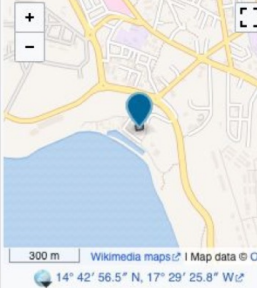
Upload media [↗](#)
Wikipedia

Instance of [mosque](#)

Location Dakar, Dakar Department, Dakar, Senegal

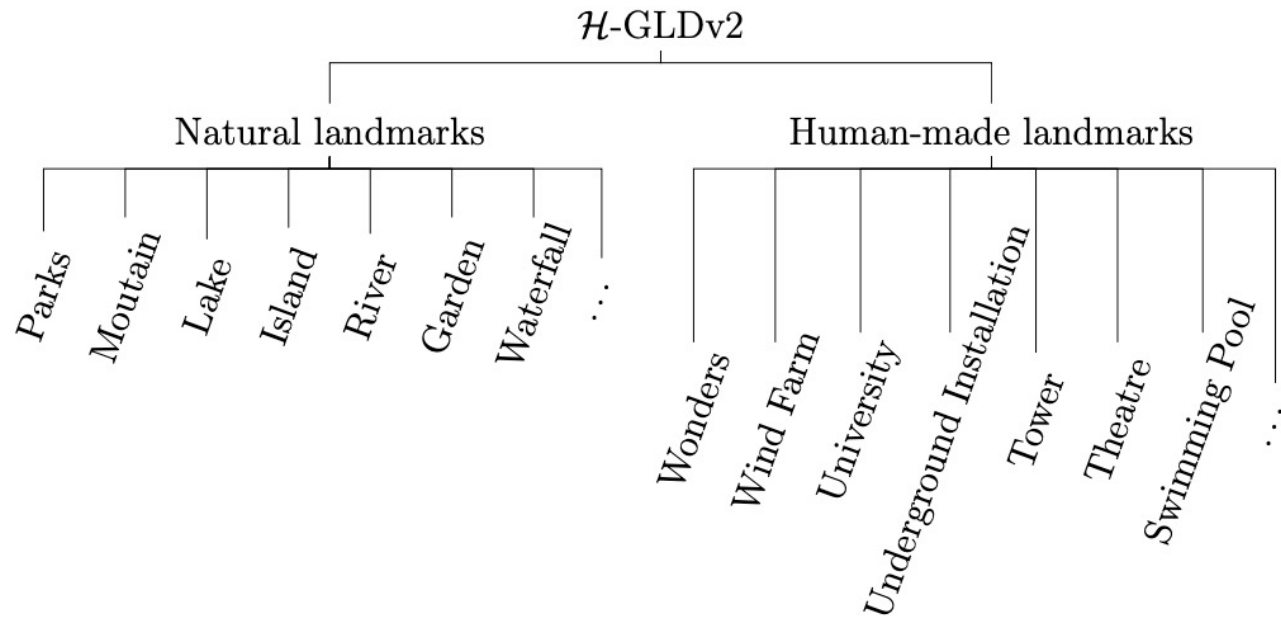
Start time 1997

Religion or worldview Islam



Authority control [Collapse]
Q3324921
Reasonator · Scholia · PetScan · statistics · WikiMap · Locator tool · KML file · WikiShootMe · OpenStreetMap · Search depicted

2. Post-processing (manual + automatic)



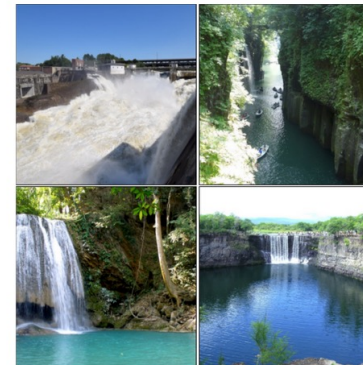
- K-Means clustering from CLIP's textual encoder
- Manual verification + adding natural/man made in hierarchy
- 78 super categories



Bridge.



Castle.

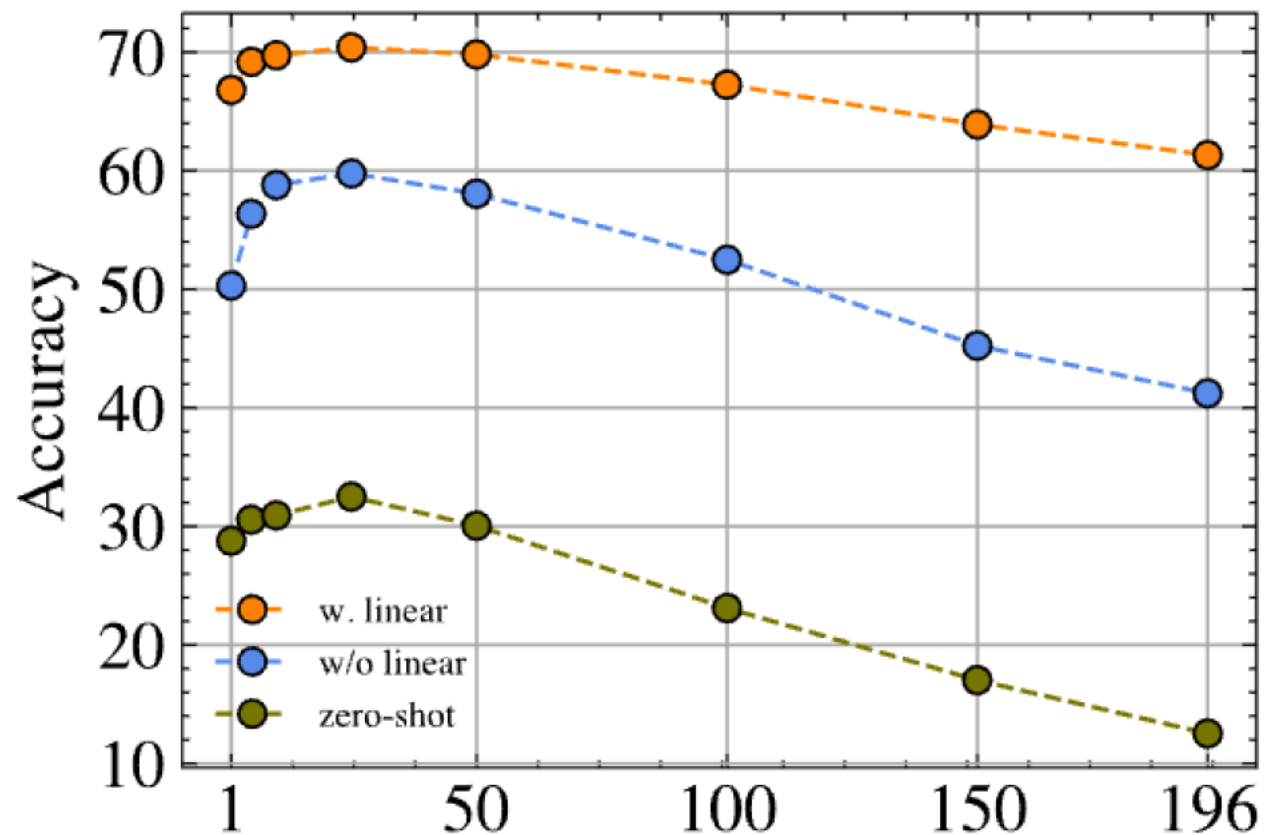


Waterfall.



Volcano.

GalLoP Results



Impact of number of regions k

| Local Method | Sparse | Dense |
|--------------------------------|--------|-------|
| $\text{CLIP}_{\text{Local}}$ | 30.9 | 12.5 |
| Local Align. | 67.9 | 59.4 |
| $\text{Prompt}_{\text{Local}}$ | 59.8 | 41.2 |
| $\text{GalLoP}_{\text{Local}}$ | 69.8 | 61.1 |

Main GalLop's components