

# Robust image retrieval with deep learning.

Scientific seminar @ Orasis 2025 – Elias Ramzi June 13, 2025



### Summary.

Robust image retrieval:

- ROADMAP: optimization of rank losses.
- HAPPIER: Hierarchical Image Retrieval for Robust Ranking.

Few shot robustness in the era of foundation models:

- HEAT: Post-hoc out-of-distribution detection.
- GalLoP: Few shot adaptation of vision-language models.

### Summary and contributions.



Ramzi, Elias, et al. "Robust and Decomposable Average Precision for Image Retrieval." NeurIPS, 2021.
 Ramzi, Elias, et al. "Hierarchical Average Precision Training for Pertinent Image Retrieval." ECCV, 2022.
 Ramzi, Elias, et al. "Optimization of Rank Losses for Image Retrieval." TPAMI, 2025.

### Summary and contributions.



Ramzi, Elias, et al. "Robust and Decomposable Average Precision for Image Retrieval." NeurIPS, 2021.
 Ramzi, Elias, et al. "Optimization of Rank Losses for Image Retrieval." TPAMI, 2025.

#### Image retrieval.



Image retrieval  $\rightarrow$  retrieve similar images to the query image.

### Cham

### **Optimizing rank losses.**







Evaluation metrics are not differentiable  $\rightarrow$  can not be used directly for SGD.

Evaluation metrics are non-decomposable  $\rightarrow$  shortcomings for mini-batch training.

### Ranking in stochastic gradient optimization.

Evaluation metrics in image retrieval  $\rightarrow$  rank-based:

 $egin{aligned} \mathbf{AP} &= rac{1}{|\Omega^+|} \sum\limits_{k \in \Omega^+} rac{\mathrm{rank}^+(k)}{\mathrm{rank}(k)} \ & ext{Recall} @n = rac{\# ext{ number of positive before n}}{|\Omega^+|} \ &= rac{1}{\min(|\Omega^+|,n)} \sum\limits_{k \in \Omega^+} H\left(n - \mathrm{rank}(k)
ight) \end{aligned}$ 

Definition of the rank with Heaviside functions:



 $s_k$  : cosine similarity between

image k and the query.

 $\mathrm{rank}(k) = \# ext{ number of similarities greater than } s_k \ = 1 + \sum_{j \in \Omega} H(s_j - s_k)$ 

7

🐹 ROADMAP

### Surrogate losses: approximation of the metric.<sup>KE</sup> ROADMAP

- Coarse upper-bound of the metrics.
- Not well-aligned with evaluation metrics: supports bottom vs. top of the ranking



Surrogate losses:

Triplet loss – Wu, Chao-Yuan, et al. *ICCV*, 2017. NSM: Zhai, Andrew, et al. *BMVC*, 2018. Multi-similarity loss – Wang, Xun, et al. *CVPR*, 2019.

### Surrogate losses: approximation of the rank.

19

26



0.4

Recall

Query

Precision-Recall curves

Precision

0.0

|P| = 141

+ Smooth AP

Baseline Network (AP = 0.09)

+ Smooth-AP (AP = 0.58)

Long line of work:
Learned Ranking - SoDeep: Engilberge, Martin, et al. CVPR, 2019.
Histogram binning - SoftBin: Revaud, Jerome, et al. ICCV, 2019.
Blackbox optimization of AP: Rolínek, Michal, et al. CVPR, 2020.
Smooth approximations of AP: Brown, Andrew, et al. ECCV, 2020.

No strong theoretical guarantees: not an upper bound.

😹 ROADMAP

 Ill-behaved gradient flow.



### Non-decomposability.



Ranking-based metrics are not decomposable:

 $\mathcal{M}(\Omega) 
eq rac{1}{|\mathcal{B}|} \sum\limits_{b \in \mathcal{B}} \mathcal{M}(b)$ 

- Loss on batches are biased.
- Training can stop before global loss is zero.



### Addressing non-decomposability.



- → Fewer works: brute force approaches
  - Hard negative mining.
  - Doubling the number of forward passes.
  - Storing the dataset.
- $\rightarrow$  Overhead in training time.

Cham

Negative mining – Wu, Chao-Yuan, et al. *ICCV*, 2017. Large batches – Revaud, Jerome, et al. *ICCV*, 2019. cross batch memory – Wang, Xun, et al. *CVPR*, 2020. 🔀 ROADMAP

### Robust and decomposable rank losses.





## le c**nam**

### SupRank: smooth approximation of the rank.



😹 ROADMAP

## le c**nam**

### SupRank: smooth approximation of the rank.



$$\mathrm{AP} = rac{1}{|\Omega^+|} \sum_{k \in \Omega^+} rac{\mathrm{rank}^+(k)}{\mathrm{rank}^+(k) + \mathrm{rank}^-(k)}$$

$$\mathrm{rank}^-_s(k) = \sum_{j\in\Omega^-_k} {H^-(s_j-s_k)}$$

- Optimizing rank<sup>-</sup> → smooth approximation + upper bound.
- Not optimizing rank<sup>+</sup> → well-behaved gradients.

😹 ROADMAP

### The decomposability gap.

$$\mathrm{DG}(\Omega) = rac{1}{|\mathcal{B}|} \sum\limits_{b \in \mathcal{B}} \mathcal{M}(b) - \mathcal{M}(\Omega)$$

→ Difference between the average ranking-based metrics on batch and its value on the whole dataset.





ROADMAP





$$\mathcal{L}_{DG}(oldsymbol{ heta}) = rac{1}{|\Omega^+|}\sum_{oldsymbol{x}_{oldsymbol{j}}\in\Omega^+} [lpha-oldsymbol{s}_j]_+ + rac{1}{|\Omega^-|}\sum_{oldsymbol{x}_{oldsymbol{j}}\in\Omega^-} [oldsymbol{s}_j-eta]_+$$

→ Calibrates scores across batches:
♦ positive scores greater than *Q*.
♦ negative scores lower than *β*.
→ Explicitly optimize an upper-bound on the decomposability gap.



Framework applicable with ranking-based loss: *e.g.* AP, Recall, NDCG.

Application to Average Precision:

$$egin{aligned} \mathcal{L}_{ ext{SupAP}} &= 1 - rac{1}{|\Omega^+|} \sum\limits_{k \in \Omega^+} rac{ ext{rank}^+(k)}{ ext{rank}^+(k) + ext{rank}^-_s(k)} \ \mathcal{L}_{ ext{ROADMAP}}(oldsymbol{ heta}) &= (1 - \lambda) \cdot \mathcal{L}_{ ext{SupAP}}(oldsymbol{ heta}) + \lambda \cdot \mathcal{L}_{ ext{DG}}(oldsymbol{ heta}) \end{aligned}$$

### Experimental validation.





SOP: retail products.

iNaturalist: wildlife images.

Song, Hyun Oh, et al. "Deep Metric Learning via Lifted Structured Feature Embedding." *CVPR*, 2016. Van Horn, Grant, et al. "The iNaturalist Species Classification and Detection Dataset." *CVPR*, 2018.

### Comparison to AP methods.





- → Fair comparison under the same setting (backbone, batch size, data, optimization...)
- → Significant gains when **changing the loss**.

#### Ablation studies.

🌃 ROADMAP



## Impact of the decomposability loss.



 $\rightarrow$  Biggest **relative** increase on small batches.

 $\rightarrow$  ROADMAP works with small batches.

🔀 ROADMAP

### Conclusion on ROADMAP.



Pros:

- **Smooth** and **decomposable** surrogate for rank losses. Consistent and **significant improvement**s.  $\rightarrow$
- $\rightarrow$

Limitation:

Only defined for binary labels.  $\rightarrow$ 



### Summary and contributions.



Ramzi, Elias, et al. "Hierarchical Average Precision Training for Pertinent Image Retrieval." ECCV,2022.
 Ramzi, Elias, et al. "Optimization of Rank Losses for Image Retrieval." TPAMI, 2025.

### Hierarchical Image Retrieval for Robust Ranking. HAPPIER



- Image retrieval is binary  $\rightarrow$  do not take into account mistake severity.
- Extend AP to graded setting to take importance of errors into account.

### Hierarchical learning.





- Hierarchy is a good proxy for human perception of mistake severity.
- Several public datasets include hierarchical annotations.

### Hierarchical image retrieval in the literature.





- → Introduction of the "Dynamic Metric learning" datasets.
- → CSL = triplet + proxy loss for hierarchical setting.
  - Limitation: **not aligned** with evaluation metrics.

Sun, Yifan, et al. "Dynamic metric learning: Towards a scalable metric space to accommodate multiple semantic scales." *CVPR*, 2021.

### HAPPIER.





### Relevance function: graded similarities.









- Leverage a hierarchical tree to design a graded similarity, or "relevance" between categories.
- **Decreasing function of the distance** in the hierarchical tree.

 $\mathrm{rel}(k) = rac{l/L}{|\Omega^{(l)}|}$ 

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

### Hierarchical rank.



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



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



### $\mathcal{H} ext{-}\mathrm{AP} = rac{1}{\sum\limits_{k\in\Omega^+}\mathrm{rel}(k)}\sum\limits_{k\in\Omega^+}rac{\mathcal{H} ext{-}\mathrm{rank}(k)}{\mathrm{rank}(k)}$

 $\mathcal{H}$ -AP properties:

- Consistent generalization of AP.
- Keeps desirable properties of AP.
  - Penalize wrong rankings.
  - Emphasis for the top of the ranking.
  - Relevant in imbalanced settings.
- Is flexible wrt. the relevance.



### Optimizing $\mathcal{H}$ -AP.

😸 HAPPIER



$$egin{aligned} \mathcal{L}_{ ext{Sup-}\mathcal{H} ext{-} ext{AP}}(m{ heta}) &= 1 - rac{1}{\sum\limits_{k\in\Omega^+} ext{rel}(k)}\sum\limits_{k\in\Omega^+}rac{\mathcal{H} ext{-} ext{rank}(k)}{ ext{rank}^+(k) + ext{rank}^-_s(k)} \ \mathcal{L}_{ ext{HAPPIER}}(m{ heta}) &= (1-\lambda)\cdot\mathcal{L}_{ ext{Sup-}\mathcal{H} ext{-} ext{AP}}(m{ heta}) + \lambda\cdot\mathcal{L}_{ ext{DG}}^*(m{ heta}) \end{aligned}$$

### Comparison to hierarchical methods.





→ HAPPIER outperforms the state-of-the-art hierarchical loss CSL on fine-grained and hierarchical retrieval.

iNat-base

iNat-full

TL: Wu, Chao-Yuan, et al. *ICCV*, 2017.
 NSM: Zhai, Andrew, et al. *BMVC*, 2018.
 CSL: Sun, Yifan, et al. *CVPR*, 2021.

SOP

### Comparison to fine-grained methods.



- On par for fine-grained retrieval ("Species").  $\rightarrow$
- Large gains on other hierarchical levels from "Family".  $\rightarrow$

TL: Wu, Chao-Yuan, et al. ICCV, 2017. NSM: Zhai, Andrew, et al. BMVC, 2018. HAPPIER

### Qualitative results.

HAPPIER



- High quality of clusters.
  - Better mistakes, even in failure cases.

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

Ouer



Relevant index images:



 $GLDv2 \rightarrow large$ scale landmarks retrieval dataset.

HAPPIER

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

Weyand, Tobias, et al. "Google landmarks dataset v2 - a large-scale benchmark for instance-level recognition and retrieval." *CVPR*, 2020.

### Scraping Wikimedia Commons.

Upload media

Wikipedia

Japan

Height 350 m

above sea 1,260 m level

Elevation



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 calciferous sandstone

material Location Dean, Allerdale, Cumbria, North West

England, England, UK Architectural English Gothic style architecture

Norman architecture

**Diocese** Diocese of Carlisle Heritage Grade I listed building designation (1986-)

Inception 12th century Religion or Anglicanism worldview

official website 2 E 7 + 1.4 -Dean 200 m Wikimedia mansi 2 I Man data D

Q 54° 36' 53.64" N, 3° 26' 25.08" WC

Buena Vista Lagoon [Collapse] Shōmvō Falls [Collapse] lake in United States of America Upload mediar? Wikipedia Instance of lake Location California Instance of waterfall Named after Shōmyō Located in Chübu-Sangaku protected area National Park Wikimedia mapsi2" | Map data @. a3° 10' 23.16" N, 117° 21' 00" W& Authority control Collapse Located in or Q4985455 next to body Shomyo River Reasonator @ . Scholia @ . PetScan @ . of water statistics 2 · WikiMap 2 · Locator tool 2 · KML file (? • WikiShootMe (? • OpenStreetMap (? Search depicted designation of Japan

building in Senegal waterfall in Toyama Prefecture, Japan Upload media (1) Wikipedia Instance of mosque Location Dakar, Dakar Department, Dakar, Senegal Start time 1997 Part of Japan's Top 100 **Religion or** worldview Waterfalls (38) **F 7** 1.1 Location Ashikuraji, Tatevama Nakaniikawa district Toyama Prefecture 300 m | Wikimedia maps (≥ I Map data © O. Heritage Place of Scenic Beauty Ca 14° 42' 56.5" N, 17° 29' 25.8" W& Authority control [Collaps natural monument O332492 Reasonator 12 · Scholia 12 · PetScan 12 · 126 m (Q46868895) statistics @ • WikiMap @ • Locator tool @ • KML file & • WikiShootMe & • OpenStreetMap & Search depicted



Mosquée de la Divinité [Collapse]

Wikimedia Commons  $\rightarrow$  the largest open database of landmarks.

#### Scraped labels include:

- 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.
## e cnam

#### Final super categories.







Bridge.



Waterfall.



Castle.



37

Cham  $\mathcal{H}$ -GLDv2 results.





# e cnam

#### Results on Coexya data.







→ HAPPIER can be extended to multi-label setting.

→ Optimizing HAPPIER works best on both fine-grained and hierarchical metrics.

### Qualitative results on Coexya data.





HAPPIER  $\rightarrow$  retrieves images with **both labels**. 



- Optimizing hierarchical metrics:
  - Better robustness wrt. to mistake severity.
  - Performances on par on fine-grained metrics.
- Hierarchical annotations  $\rightarrow$  **not too costly** + boost models' robustness.
- HAPPIER  $\rightarrow$  production in Acsepto



#### Short term extensions.

#### 🔀 ROADMAP:

- Optimize other metrics: Recall, NDCG.
- Use in GNN recommender systems.



😸 HAPPIER:

- Use in multi-label setting + business driven insights.
- Aerial image time series ranking.



2019-06-01



2019-08-30



2019-09-19

Few shot robustness in the era of foundation models.

### Cedric EA4629COEXYA<br/>Connect skills, create more<br/>Valeo.ai

### Summary and contributions.



Lafon, Marc, et al. "Hybrid Energy Based Model in the Feature Space for Out-of-Distribution Detection." ICML, 2023.
 Lafon, Marc, et al. "GalLoP: Learning global and local prompts for vision-language models", ECCV, 2024

### Summary and contributions.



Lafon, Marc, et al. "Hybrid Energy Based Model in the Feature Space for Out-of-Distribution Detection." *ICML*, 2023.

#### Out-of-distribution detection.





Al confuses a horse-drawn carriage (OOD) for a truck.

- Detection of in-distribution (ID) = similar to train dataset, vs. out-of-distribution (OOD) = dissimilar to train dataset.
- Important task for real world deployment.

 $OOD \rightarrow binary problem:$ 

$$G_\lambda(x) = egin{cases} ext{ID, if } E(x) \leq \lambda_{95\%} \ ext{OOD, if } E(x) > \lambda_{95\%} \end{cases}$$

#### Post-hoc out-of-distribution detection.





• Enjoy strong predictive performances from off-the-shelf neural networks.

 Use any backbones → ConvNets, transformers.

• Limited compute overhead at inference time.

#### Prior out-of-distribution scorers.







Classification based methods:  $\rightarrow$  Better at distinguishing near-OOD.  $\rightarrow$  Strong guarantees for far-OOD.  $\rightarrow$  No control for far-OOD.

Energy logits (EL): Liu, Weitang, et al. NeurIPS, 2020. DICE: Sun, Yiyou, et al. ECCV, 2022.

Parametric methods:  $\rightarrow$  Lower performances on near-OOD.

SSD (GMM): Sehwag, Vikash, et al. ICLR, 2021. KNN: Sun, Yiyou, et al. ICML, 2022.

#### Motivations.



\_



### HEAT for out-of-distribution detection.



- $EBM \rightarrow data driven residual learning.$
- Composition  $\rightarrow$  leverage  $\neq$  modeling biases.

HEAT

### Hybrid energy-based density estimation.

Hybrid density:

$$p^h_{ heta_k}(z) = rac{1}{Z( heta_k)} p^r_{ heta_k}(z) q_k(z)$$

Hybrid energy:

 $E^{\,h}_{ heta_k}(z)=E_{q_k}(z)+E_{ heta_k}(z)$ 



HEAT

#### Composition of refined prior density estimators.

Class prior:

$$E^h_{ heta_l}(z) = -\log \sum_c e^{f(z)[c]} + E^r_{ heta_l}(z)$$
 .

Feature & style prior:

$$E^{\,h}_{ heta_g}(z) = -\log \sum_c e^{-rac{1}{2}(z-\mu_c)^T \Sigma^{-1}(z-\mu_c)} + E^{\,r}_{ heta_g}(z)$$

Energy function composition:

$$E^{eta}_{ ext{HEAT}} = rac{1}{eta} {\log \sum_{k=1}^{K} e^{eta E^{h}_{ heta_{k}}}}$$



EL – Liu, Weitang, et al. NeurIPS, 2020. GMM – Lee, Kimin, et al. NeurIPS, 2018. Gram matrix – Sastry, Chandramouli Shama, et al. ICML, 2020.



#### Residual learning.





• Residual learning → better performances

### Composition.





• Results better than individual scorers.

• Composition  $\rightarrow$  performance boost.

#### Comparison to state-of-the-art methods.



• HEAT  $\rightarrow$  competitive performances on 3 datasets

HEAT

#### Conclusion on HEAT.



- Post-hoc OOD detection  $\rightarrow$  no assumption on models.
- EBM models boost prior OOD scorers.
- Composing refined OOD scorers  $\rightarrow$  boost performances.



#### Summary and contributions.



Lafon, Marc, et al. "GalLoP: Learning global and local prompts for vision-language models", ECCV, 2024

#### Vision-language models.



(1) Contrastive pre-training



(2) Create dataset classifier from label text

#### CLIP:

- Vision-language pre-training.
- Zero-shot classification on downstream datasets.

CLIP – Radford, Alec, et al. ICML, 2021. ALIGN – Jia, Chao, et al. ICML, 2021.

#### Few shot adaptation with prompt learning.



- Few shot adaptation of VLMs to downstream dataset.
- No need for "prompt engineering".

📕 CoOp – Zhou, Kaiyang, et al. IJCV, 2022.

Cham

 $\bigcirc$ 

🐎 GalLoP

#### Motivations.





- Prompt learning trade off top-1 vs. robustness.
- Ensembling and local features  $\rightarrow$  robustness.

LoCoOp – Miyai, Atsuyuki, et al. NeurIPS, 2023. PromptSRC – Khattak, Muhammad Uzair, et al. ICCV, 2023.

### Prompt learning with local features.





🐎 GalLoP

- Optimal transport for assignment.
- All local features → including irrelevant features.

PLOT – Chen, Guangyi, et al. "Plot: Prompt learning with optimal transport for vision-language models." ICLR, 2023.

- Local features for regularization.
- Lower top-1 accuracy

LoCoOp – Miyai, Atsuyuki, et al. "Locoop: Few-shot out-of-distribution detection via prompt learning." NeurIPS, 2023.

### GalLoP.





- Multiple global prompts w/ "prompt dropout".
- Local alignment w/ sparsity and linear projection.
- Multiple local prompts w/ multiscale loss.

#### Prompt learning on local features.

$$egin{aligned} & ext{sim}_{ ext{top-}k}(\mathcal{Z}_l,oldsymbol{t}_c)\!\coloneqq\!rac{1}{k}\!\sum_{i=1}^L \mathbbm{1}_{ ext{top-}k}(i)\!\cdot\!ig\langleoldsymbol{z}_i^l,oldsymbol{t}_cig
angle \ & ext{top-}k(i)\!=\!igg\{egin{aligned} 1 & ext{if} & ext{rank}_i(ig\langleoldsymbol{z}_i^l,oldsymbol{t}_cig
angle)\!\leq\!k, \ & ext{otherwise.} \end{aligned}$$

Exploit local features for prompt learning:

- Sparse local similarity.
- Learnable projection layer.



🐎 GalLoP

# le c**nam**

### Diversity: prompt dropout & multiscale loss.





• Diversity by enforcing different gradient signal.

#### Few shots results.



le c**nam** 



- Strong performances in **low shots** settings.
- Outperforms state-of-the-art prompt learning methods.

# e c**nam**

#### Robustness results.





- Strong domain generalization results (top-1).
- Outperforms dedicated OOD detection methods.

### Combining global and local features





- Vanilla local features  $\rightarrow$  no improvement.
- **Complementarity** of global / local features in GalLop.

### The need for sparsity.



- Prompt learning → boosts performances.
- Sparsity → large gains in three regimes.
- Linear projection → better alignment.



# Learning multiple prompts.







(a) Impact of prompt dropout when learning multiple global prompts.

(b) Impact of multiple scales when learning local prompts, with  $k_1 = 10$ ,  $\Delta_k = 10$ .

- Prompt dropout  $\rightarrow$  gains when using multiple global prompts.
- Multiscale loss  $\rightarrow$  boosts local performances.

#### Qualitative results.

🐎 GalLoP



- Localization of objects.
- Few shot and weakly supervised segmentation.



- Exploit local Features: sparsity & learnable projection.
- Global + local prompt learning  $\rightarrow$  strong top-1 + robustness.

impala (antelope)



Ground truth

recreational vehicle



CLIP local



GalLop 1 scale (k=10)

gazelle



GalLop 4 scales

#### Short term extensions.

#### HEAT:

- Extend other scorers, e.g. kNN.
- Apply HEAT to dense prediction.
- Extension to other modalities, e.g. NLP, audio.



#### 🐎 GalLoP:

- Integrate an adaptive sparsity ratio.
- Learn the linear projection on a bigger vision-language dataset.


# Current research interests.



### valeo.ai

#### World model for autonomous driving



Bartoccioni, Florent, et al. "VaViM and VaVAM: Autonomous Driving through Video Generative Modeling." ArXiv, 2025.

#### Vision langue models

VLM

black-box





Cardiel, Amaia, et al. "LLM-wrapper: Black-Box Semantic-Aware Adaptation of Vision-Language Models for Referring Expression Comprehension." ICLR, 2025.

## Thank you for your attention.

Elias Ramzi, et al. "Robust and Decomposable Average Precision for Image Retrieval." *NeurIPS*, 2021 & *RFIAP*, 2022. https://github.com/elias-ramzi/ROADMAP

EliasRamzi, et al. "Hierarchical Average Precision Training for Pertinent Image Retrieval." *ECCV*,2022. <br/>
<a href="https://github.com/elias-ramzi/HAPPIER">https://github.com/elias-ramzi/HAPPIER</a>

Marc Lafon, et al. "Hybrid Energy Based Model in the Feature Space for Out-of-Distribution Detection." ICML, 2023.
<u>https://github.com/MarcLafon/heatood</u>

Elias Ramzi, et al. "Optimization of Rank Losses for Image Retrieval." *TPAMI*, 2025. <u>https://github.com/cvdfoundation/google-landmark</u>

Marc Lafon, et al. "GalLoP: Learning global and local prompts for vision-language models", ECCV, 2024 <u>https://github.com/MarcLafon/gallop</u>

