

# Positive-First Most Ambiguous: A Simple Active Learning Criterion for Interactive Retrieval of Rare Categories

## Supplementary Material

### 7. Performance of Representativeness-based AL Strategies Across Iterations

Tab. 5 shows the performance of representative-based AL methods for different iterations. The strong performance of the uncertainty-based sampling strategy is consistent from the early retrieval stage.

Table 5. Performance comparison of representativeness-based AL methods vs. uncertainty-based selection on CIFAR100-LT at iterations 5, 15 and 25.

metric	<i>Rand</i>	<i>DAL</i>	<i>CoreSet</i>	<i>MA</i>
<i>cov</i> <sub>5</sub>	0.05	0.039	0.039	<b>0.504</b>
<i>f1</i> <sub>5</sub>	0.182	0.102	0.07	<b>0.828</b>
<i>pos</i> <sub>5</sub>	0.022	0.019	0.02	<b>0.3</b>
<i>cov</i> <sub>15</sub>	0.081	0.066	0.058	<b>0.839</b>
<i>f1</i> <sub>15</sub>	0.294	0.215	0.099	<b>0.857</b>
<i>pos</i> <sub>15</sub>	0.029	0.024	0.026	<b>0.537</b>
<i>cov</i> <sub>25</sub>	0.111	0.096	0.087	<b>0.899</b>
<i>f1</i> <sub>25</sub>	0.381	0.321	0.164	<b>0.857</b>
<i>pos</i> <sub>25</sub>	0.036	0.032	0.035	<b>0.61</b>

### 8. Influence of Coverage Granularity over Iterations

Tab. 6 extends the results of the effect of coverage granularity on each method’s performance to earlier iterations. The results are obtained on ImageNet-LT, using DINOv2 features. Our methods consistently shows strong and stable results since the early retrieval stages, rendering it robust to the choice of  $K$ . In contrast, other methods are sensitive to this value.

### 9. Performance Across Multiple Datasets and Feature Descriptors

#### 9.1. Class coverage scores

We report the class coverage scores of the different methods, datasets and feature extractors in Tab. 7, across different iterations. The results highlight consistently strong performance of *PF-MA* since the early stages of retrieval. Other methods do not keep consistent rankings. This further guarantees that our method promotes very early user satisfaction.

#### 9.2. Classifier performance

Tab. 8 reports the f1-score of the classifier across different methods, datasets and feature extractors. *PF-MA* is always

Table 6. Class coverage scores at iterations 5, 15 and 25 for different  $K$  values. Second best is underlined.

metric	method	$K$		
		16	32	64
<i>cov</i> <sub>5</sub>	<i>MA</i>	<u>0.55</u>	<u>0.423</u>	0.327
	<i>MP</i>	0.413	0.396	<u>0.399</u>
	<i>ALAMP</i>	0.49	0.387	0.315
	<i>MA-S</i>	0.51	0.377	0.28
	<i>MA-D</i>	0.513	0.355	0.246
	<i>MP-S</i>	0.427	0.381	0.343
	<i>MP-D</i>	0.447	0.378	0.297
	<i>PF-MA</i>	<b>0.586</b>	<b>0.493</b>	<b>0.4</b>
<i>cov</i> <sub>15</sub>	<i>MA</i>	<u>0.807</u>	<u>0.736</u>	0.644
	<i>MP</i>	0.69	0.688	<u>0.679</u>
	<i>ALAMP</i>	0.757	0.7	0.622
	<i>MA-S</i>	0.799	0.722	0.62
	<i>MA-D</i>	0.774	0.653	0.527
	<i>MP-S</i>	0.688	0.681	0.666
	<i>MP-D</i>	0.68	0.647	0.582
	<i>PF-MA</i>	<b>0.812</b>	<b>0.79</b>	<b>0.778</b>
<i>cov</i> <sub>25</sub>	<i>MA</i>	<u>0.848</u>	<u>0.804</u>	0.738
	<i>MP</i>	0.776	0.776	0.764
	<i>ALAMP</i>	0.808	0.771	0.712
	<i>MA-S</i>	<u>0.848</u>	0.802	0.731
	<i>MA-D</i>	0.833	0.751	0.648
	<i>MP-S</i>	0.779	0.777	<u>0.764</u>
	<i>MP-D</i>	0.773	0.76	0.723
	<i>PF-MA</i>	<b>0.849</b>	<b>0.844</b>	<b>0.843</b>

among the top-3 methods in terms of classifier generalization, with minimal degradation, although our main objective is retrieval. Other competing methods have much lower coverage results as discussed before. Thus, *PF-MA* is able to maintain a fairly generalizable classifier, while achieving better coverage results.

### 10. Performance Across the Search Iterations and Per range of Class Size

#### 10.1. Class coverage scores

Fig. 7 and Fig. 8 show the strong performance of *PF-MA* across iterations and per range of class size, for different datasets. The results are less pronounced on PlantNet300K because of its specific nature. *MP* and variants perform well on smaller classes, but they struggle for larger ones, where *MA* perform better thanks to the less imbalanced distribution of the negative vs. positive class problem. Our method yields consistently higher coverage across iterations, and

Table 7. Class coverage scores at iterations 5, 15 and 25 for different AL methods, datasets and descriptors. Second best is underlined.

metric	method	Cifar100-LT		ImageNet-LT		PlantNet300K	
		CLIP	DINOv2	CLIP	DINOv2	CLIP	DINOv2
<i>cov</i> <sub>5</sub>	<i>MA</i>	0.402	<u>0.502</u>	0.351	<u>0.423</u>	0.203	0.282
	<i>MP</i>	<u>0.41</u>	0.454	<b>0.368</b>	0.396	<b>0.212</b>	0.29
	<i>ALAMP</i>	0.349	0.471	0.289	0.387	0.148	0.235
	<i>MA-S</i>	0.363	0.442	0.321	0.377	0.188	0.257
	<i>MA-D</i>	0.388	0.407	0.327	0.355	0.193	0.284
	<i>MP-S</i>	0.381	0.435	0.34	0.381	0.193	0.27
	<i>MP-D</i>	0.409	0.429	0.347	0.378	0.198	<u>0.294</u>
	<i>PF-MA</i>	<b>0.411</b>	<b>0.56</b>	<u>0.36</u>	<b>0.493</b>	<u>0.204</u>	<b>0.298</b>
<i>cov</i> <sub>15</sub>	<i>MA</i>	0.801	<u>0.838</u>	<u>0.76</u>	<u>0.736</u>	<u>0.487</u>	<u>0.59</u>
	<i>MP</i>	0.773	0.785	0.713	0.688	0.456	0.537
	<i>ALAMP</i>	0.704	0.83	0.638	0.7	0.318	0.514
	<i>MA-S</i>	0.775	0.821	0.735	0.722	0.465	0.57
	<i>MA-D</i>	<u>0.815</u>	0.74	0.734	0.653	0.48	0.588
	<i>MP-S</i>	0.756	0.785	0.694	0.681	0.435	0.523
	<i>MP-D</i>	0.774	0.736	0.697	0.647	0.447	0.542
	<i>PF-MA</i>	<b>0.824</b>	<b>0.915</b>	<b>0.773</b>	<b>0.79</b>	<b>0.488</b>	<b>0.604</b>
<i>cov</i> <sub>25</sub>	<i>MA</i>	0.89	<u>0.899</u>	<u>0.852</u>	<u>0.804</u>	<b>0.596</b>	0.678
	<i>MP</i>	0.875	0.877	0.821	0.776	0.565	0.628
	<i>ALAMP</i>	0.811	0.892	0.744	0.771	0.406	0.605
	<i>MA-S</i>	0.88	0.895	0.844	0.802	0.583	0.671
	<i>MA-D</i>	<u>0.903</u>	0.84	0.842	0.751	<u>0.595</u>	<u>0.679</u>
	<i>MP-S</i>	0.869	0.882	0.815	0.777	0.552	0.625
	<i>MP-D</i>	0.881	0.86	0.816	0.76	0.563	0.639
	<i>PF-MA</i>	<b>0.908</b>	<b>0.954</b>	<b>0.861</b>	<b>0.844</b>	<b>0.596</b>	<b>0.684</b>

presents the best compromise between smaller and larger classes. This highlights its practicality at test time, where early user satisfaction is desired, and the actual size of the user-defined concept is unknown.

## 10.2. Classifier performance

Both Fig. 9 and Fig. 10 show the ability of *PF-MA* to achieve competitive classifier generalization performance. This performance is steady across iterations, and across class sizes. *MP* on the other hand, has very poor classifier performance overall, always with a significant drop on larger class sizes.

Table 8. F1 scores at iterations 5, 15 and 25 for different AL methods, datasets and descriptors. Top-3 results are underlined. A star (\*) is used only when our method achieves the top-1 result.

metric	method	Cifar100-LT		ImageNet-LT		PlantNet300K	
		CLIP	DINOv2	CLIP	DINOv2	CLIP	DINOv2
$f_{15}$	<i>MA</i>	<u>0.604</u>	<u>0.83</u>	<u>0.504</u>	<u>0.675</u>	<u>0.191</u>	<u>0.413</u>
	<i>MP</i>	0.258	0.597	0.187	0.573	0.077	0.221
	<i>ALAMP</i>	0.56	0.806	0.45	0.646	0.147	0.373
	<i>MA-S</i>	<u>0.57</u>	<u>0.816</u>	<u>0.476</u>	<u>0.667</u>	<u>0.176</u>	<u>0.398</u>
	<i>MA-D</i>	0.517	0.803	0.437	0.661	<u>0.159</u>	<u>0.398</u>
	<i>MP-S</i>	0.325	0.666	0.241	0.59	0.091	0.279
	<i>MP-D</i>	0.358	0.637	0.245	0.579	0.093	0.296
	<i>PF-MA</i>	<u>0.601</u>	<u>0.814</u>	<u>0.505*</u>	<u>0.67</u>	<u>0.191*</u>	<u>0.41</u>
$f_{15}$	<i>MA</i>	<u>0.659</u>	<u>0.86</u>	<u>0.59</u>	<u>0.727</u>	<u>0.296</u>	<u>0.498</u>
	<i>MP</i>	0.431	0.71	0.378	0.63	0.174	0.388
	<i>ALAMP</i>	0.644	<u>0.858</u>	0.56	0.709	0.237	0.476
	<i>MA-S</i>	<u>0.656</u>	<u>0.858</u>	<u>0.584</u>	<u>0.725</u>	0.283	<u>0.493</u>
	<i>MA-D</i>	0.637	0.85	0.571	<u>0.723</u>	<u>0.268</u>	0.492
	<i>MP-S</i>	0.445	0.727	0.387	0.631	0.174	0.39
	<i>MP-D</i>	0.438	0.725	0.383	0.624	0.17	0.396
	<i>PF-MA</i>	<u>0.665*</u>	<u>0.861*</u>	<u>0.587</u>	<u>0.723</u>	<u>0.296*</u>	<u>0.495</u>
$f_{25}$	<i>MA</i>	0.648	<u>0.859</u>	<u>0.593</u>	<u>0.735</u>	<u>0.298</u>	<u>0.508</u>
	<i>MP</i>	0.456	0.73	0.416	0.63	0.196	0.421
	<i>ALAMP</i>	<u>0.647</u>	<u>0.864</u>	0.58	0.722	0.266	0.493
	<i>MA-S</i>	<u>0.649</u>	<u>0.862</u>	<u>0.593</u>	<u>0.735</u>	<u>0.295</u>	<u>0.506</u>
	<i>MA-D</i>	<u>0.648</u>	0.856	<u>0.589</u>	<u>0.734</u>	0.29	<u>0.508</u>
	<i>MP-S</i>	0.472	0.746	0.427	0.632	0.202	0.429
	<i>MP-D</i>	0.49	0.742	0.434	0.632	0.207	0.434
	<i>PF-MA</i>	<u>0.649*</u>	<u>0.859</u>	<u>0.583</u>	<u>0.73</u>	<u>0.299*</u>	<u>0.507</u>

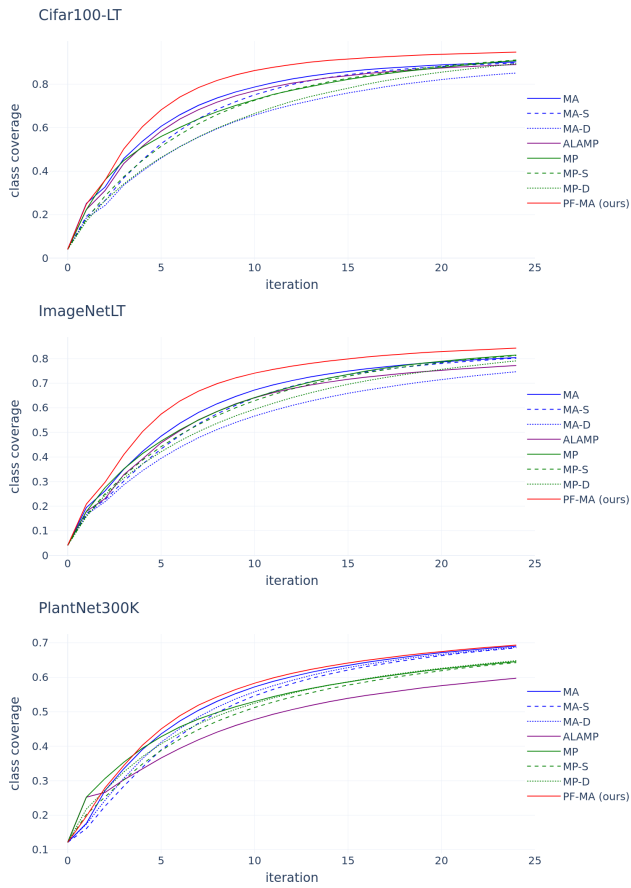


Figure 7. Class coverage results per iteration.

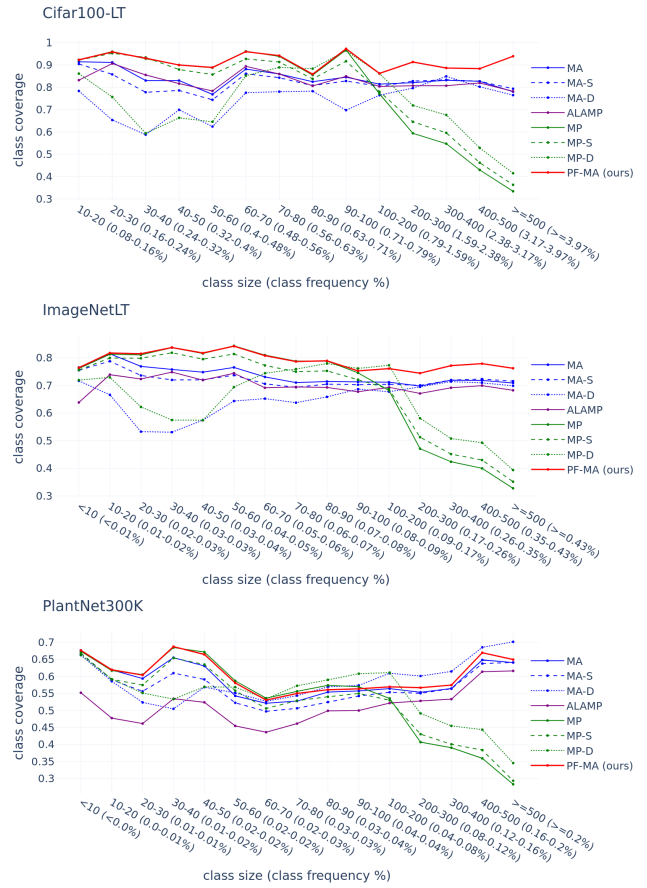


Figure 8. Class coverage results at iteration 15 per range of class size.

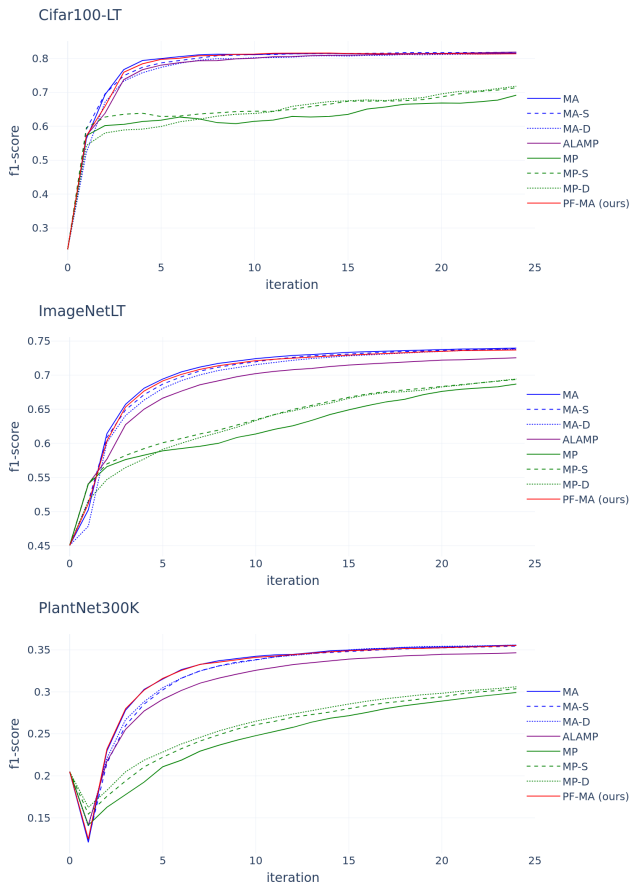


Figure 9. F1 score results per iteration.



Figure 10. F1 score results at iteration 15 per range of class size.