

HCIL: HIERARCHICAL CLASS INCREMENTAL LEARNING FOR LONGLINE FISHING VISUAL MONITORING

Jie Mei^{*}, Suzanne Romain[†], Craig Rose[†], Kelsey Magrane[†], Jenq-Neng Hwang^{*}

^{*} Department of ECE, University of Washington, Seattle, WA, USA

[†]Pacific States Marine Fisheries Commission, National Oceanic and Atmospheric Administration, USA

ABSTRACT

The goal of electronic monitoring of longline fishing is to visually monitor the fish catching activities on fishing vessels based on cameras, either for regulatory compliance or catch counting. The previous hierarchical classification method demonstrates efficient fish species identification of catches from longline fishing, where fishes are under severe deformation and self-occlusion during the catching process. Although the hierarchical classification mitigates the laborious efforts of human reviews by providing confidence scores in different hierarchical levels, its performance drops dramatically under the class incremental learning (CIL) scenario. A CIL system should be able to learn about more and more classes over time from a stream of data, i.e., only the training data for a small number of classes have to be present at the beginning and new classes can be added progressively. In this work, we introduce a Hierarchical Class Incremental Learning (HCIL) model, which significantly improves the state-of-the-art hierarchical classification methods under the CIL scenario.

Index Terms— Hierarchical Classification, Class Incremental Learning, Rehearsal Method, Longline Fishing

1 Introduction

Electronic Monitoring (EM) of Fisheries Automated imagery analysis techniques have drawn increasing attention in fisheries science and industry [1, 2, 3, 4, 5, 6, 7, 8] because they are more scalable and deployable than conventional manual survey and monitoring approach.

The goal of EM is to systematically monitor fish captures using cameras on fishing vessels either for catching counting or regulatory compliance. Then fisheries managers can thus assess the count of fish caught by species and size to monitor catch quotas by vessel or fishery. Besides, managers will detect the retention of specific fish species or sizes that are illegal to be kept. Therefore, accurate detection, length measurement, and species identification are critically needed in the EM systems. In this work, our approach focuses on the species identification task for the video-based longline fishing monitoring, where fish are caught on hooks and viewed

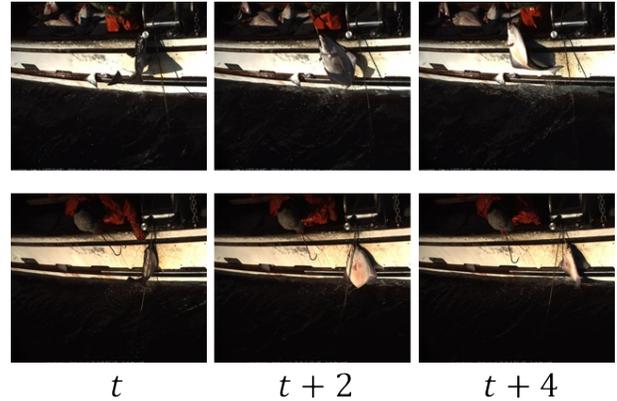


Fig. 1. Longline Fishing: Each row is a temporal sequence of an individual fish caught on a longline hook, as it is being pulled up from the sea and over the rail of the fishing vessel.

as they are pulled up from the sea and over the rail of the fishing vessel as shown in Fig.1.

Hierarchical Classification For the species identification task in the EM systems, a hierarchical classifier has more practical use for the fisheries managers than a flat classifier because it can predict coarse-level groups and fine-level species at the same time. Both levels' predictions are useful for fisheries managers to evaluate the status of fish stocks.

The previous hierarchical fish classification work [9] enforces the hierarchical data structure and introduces an efficient training and inference strategy for video-based fisheries data. With the hierarchical inference, if some input images are predicted with high confidence in one coarse-level group but with low confidence in the corresponding fine-level species, then the hierarchical model allows fisheries personnel to further assign appropriate experts to review those images and get the correct fine-level labels.

Class Incremental Learning Deep neural networks achieve remarkable performance in supervised classification tasks, but only when all the classes to be learned are available at the same time. However, real-world data are constantly acquired through time, leading to ever-changing distributions, i.e., new target fish species are added continuously.

^{*}e-mail: {jiemei, hwang}@uw.edu

[†]e-mail: {suzanne.romain, craig.rose, kelsey.magrane}@noaa.gov

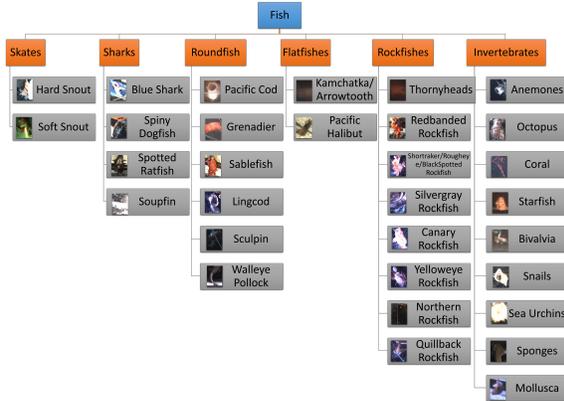


Fig. 2. Hierarchical Data Structure: The dataset, provided by NOAA fisheries scientists, includes frames and corresponding labels which are bounding box location, start and end frames’ IDs of each individual fish, coarse-level group ground truth, and fine-level species ground truth. The sample images shown here are randomly chosen from the dataset.

When a deep neural network model loses access to previous classes data (e.g., for privacy reasons, storage limitations, or data transfer difficulties) and can only be finetuned on new classes data, it could catastrophically forget the old classes, the so-called catastrophic forgetting problem [10, 11, 12, 13].

Although the previous work [9] achieves the state-of-the-art performance on hierarchical species classification task, its performance drops dramatically under the class incremental learning (CIL) scenario. A CIL system should be able to learn more and more classes over time from a stream of data, where only the training data for a small number of classes are present at the beginning and new classes can be added progressively.

Our proposed Hierarchical Class Incremental Learning (HCIL) method can provide coarse-level prediction and fine-level species at the same time, while the system gradually acquires increasing number of training fish classes over the time. More importantly, it significantly improves the state-of-the-art hierarchical classification method under the class incremental learning scenario.

The remaining sections of this paper are organized as follows. In Section 2, overviews of the related works in class incremental learning and hierarchical classification are provided. Section 3 describes details of our proposed method, HCIL. The experimental results are demonstrated and discussed in Section 4. Finally, Section 5 gives conclusions and future work.

2 Related Work

Hierarchical Classification The previous work [9] proposes a hierarchical fish species dataset as shown in Fig. 2. There are 6 coarse-level groups and 31 fine-level species. The total number of frames is more than 186K. We also use this dataset

to perform class incremental learning.

The method proposed in work [9] is an end-to-end training approach with a multi-head CNN-based architecture and two levels’ loss functions. It outperforms the traditional softmax-based flat CNN classifier in the hierarchical classification task. In contrast to it, our proposed method can not only perform the hierarchical classification task but also maintain good performance on both trained classes and newly acquired classes under the class incremental scenario.

Moreover, work [9] utilizes a video-based inference method, i.e., majority vote, for each individual fish to improve the classification accuracy. For a fair comparison, our proposed HCIL method in this paper, also utilizes the same inference method.

Class Incremental Learning Rehearsal-based methods [10, 11, 12, 13] which allow storing a fixed number of data from previously trained classes, have been widely used in the class incremental learning scenarios. More specifically, when training the model on new classes, the model can have access to the stored raw data or feature maps of previously trained classes, which are referred to as memory in the CIL setting. Our proposed HCIL method is also a rehearsal-based method. But contrast to these previous works, which are designed for flat classifiers, our model is a hierarchical classifier that contains a different memory selection module.

Except for the memory, training classes available at the same time are defined as one ’task’. Some incremental learning methods such as [14, 15] require a task identifier at test-time, i.e., need to know which task the test image belongs to. However, our proposed method discards the need for a task identifier by choosing the prediction with the maximum confidence score as the output.

3 Proposed Method

Our hierarchical class incremental learning (HCIL) method consists of a fixed pre-trained feature extraction backbone, a CIL memory selection module, and dynamic support vector machines (SVMs), i.e., SVMs are continually added with appearing of newly acquired classes as shown in Fig. 3.

Pre-trained Feature Extraction Backbone Due to the easy access to large public datasets such as ImageNet-1k [17] and COCO [18], a pre-trained backbone such as ResNet [19] or Swin Transformer [20] can be used to extract discriminative features from images even for new datasets or new classes. These public pre-trained backbones can certainly be utilized in our HCIL model.

Despite the remarkable ability of discriminative features extraction of the pre-trained backbone, inevitably it may generate non-discriminative features for some new classes beyond the dataset the backbone pre-trained on. As a result, we propose a CIL memory selection module, to be discussed later, to select those hard cases for the pre-trained backbone. Our method utilizes selected CIL memory to adjust SVM classifiers’ boundary, instead of updating the backbone, un-

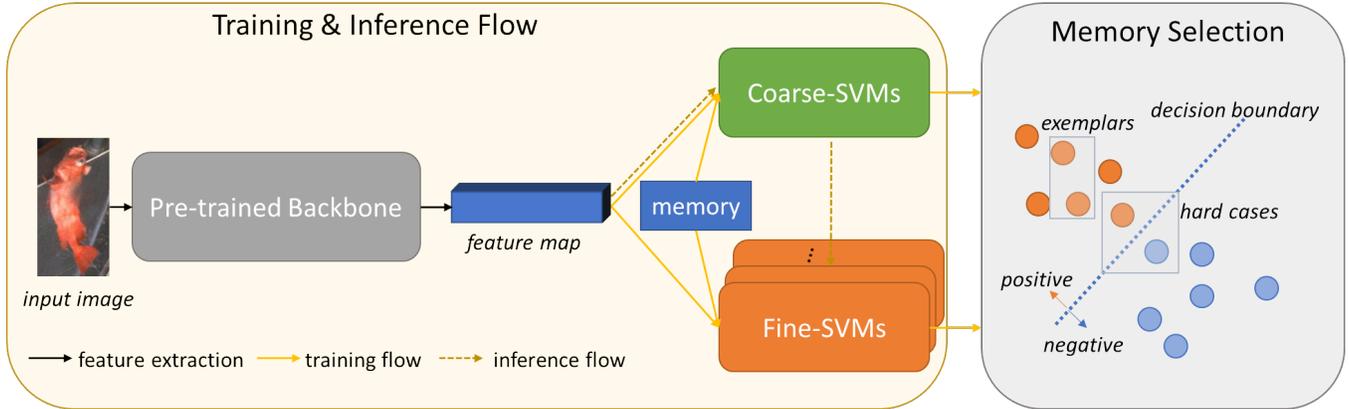


Fig. 3. HCIL: New classes’ feature maps from the fixed pre-trained backbone and old classes’ feature maps from CIL memory are used to train coarse-level group SVMs and corresponding fine-level SVMs. After SVMs training, based on the distance between SVMs’ decision boundary and each training data, CIL memory keeps adding new classes’ features, i.e., hard cases, and adding positive exemplars from new classes based on herding [16]. During inference, the extracted feature map first goes into Coarse-SVMs. And based on the coarse-level group prediction, the extracted feature map goes into corresponding Fine-SVMs for species classification.

der the hierarchical class incremental learning scenario.

Dynamic SVMs Expansion During the training of a class incremental learning scenario, the total number of seen classes is increasing. Only the training data for a small number of classes have to be present at the beginning and new classes are added progressively. Thus, the growing-class classifier can be used in the incremental learning models.

Our Coarse-SVMs consist of six SVMs for group-level prediction because there are six groups in the fish dataset as shown in Fig. 2. We use one-vs-all training strategy, where for each SVM only one group data are treated as positive and all rest of groups data are treated as negative. Each Coarse-SVM has its own corresponding Fine-SVMs. For example, the Coarse-SVM for the ‘Sharks’ group has four corresponding Fine-SVMs. And these four Fine-SVMs are added into the model sequentially with the availability of newly acquired shark species, which is called ‘dynamic SVMs expansion’ in this paper. During the inference, based on the Coarse-SVMs prediction, the feature map goes to corresponding Fine-SVMs for species classification so that the model can provide hierarchical predictions.

Algorithm 1 Positive Exemplars Selection

input feature set $F = \{f_1, \dots, f_n\}$ of new class y
input new class mean, $\mu \leftarrow \frac{1}{n} \sum_{f \in F} f$
input target number of exemplars, m_y
for $k = 1, \dots, m_y$ **do**
 $p_k \leftarrow \operatorname{argmin}_{f \in F} \left\| \mu - \frac{1}{k} \left[f + \sum_{j=1}^{k-1} p_j \right] \right\|$
end for
output exemplar set $P_y \leftarrow (p_1, \dots, p_{m_y})$

CIL Memory Selection This module selects some hard cases from newly arrived classes based on their feature maps from the pre-trained backbone. Hard cases is defined by the distances between the feature maps and decision boundary, i.e., the SVM’s outputs, as shown in Fig. 3. These distances can be represented as confidence scores, between 0 and 1, from SVM outputs via logistic transformation. The low confidence of a feature map represents a hard case. As incremental learning goes, more and more hard cases are added. In order to fix the total target number n of hard cases, we sort hard cases by their confidence scores and only keep first n'_i low confidence hard cases for each SVM_i where $n = \sum_i n'_i$.

Besides hard cases selection, the CIL memory module also selects positive exemplars from new classes based on herding [16]. More specifically, positive exemplars are selected based on Algorithm 1, where for each new class, exemplars p_1, \dots, p_{m_y} are selected iteratively until reaching the target number, m_y . Within each iteration, one more sample of the new class is selected to the exemplar set, namely the one that causes the mean feature vector over all current exemplars to best approximate the mean feature vector over all examples of this new class. Thus, the exemplar set P_y , is a prioritized list, i.e., the order of exemplars matters and exemplars earlier in the list are more important.

As incremental learning goes, more and more positive exemplars for new classes are added. In order to fix the total number of positive exemplars from all classes, m , we lower m_y to m'_y for each class y and do the same thing as done for hard cases, i.e., only keeping the first m'_y exemplars in the exemplar set, P_y , for each class y .

4 Experimental Results

We compare our method with state-of-the-art work [9] on NOAA’s dataset under both ‘hierarchical classification setting’ and ‘hierarchical class incremental learning setting’.

Hierarchical Classification Setting This setting serves as a baseline for the next hierarchical class incremental learning setting. More specifically, in this setting, all training data for all classes are available at the same time, i.e., no incremental learning scenario is assumed. Training and testing data split is the same as [9], where each individual fish has its own video data. Training and testing fish are totally different individual fish.

In this setting, our proposed **HCIL** model directly uses the fixed ResNet-101 backbone pre-trained on ImageNet-1k as the feature extractor and only trains our SVMs on NOAA’s dataset. No CIL memory is involved, thus noted as ‘**HCIL w/o m**’ in Table 1. As mentioned in Section 3, during inference, based on the Coarse-SVMs prediction, the feature map goes to corresponding Fine-SVMs for species classification so that the model can provide hierarchical predictions.

For a fair comparison, we allow work [9] to utilize the same pre-trained ResNet-101 as the initial backbone, which is also further finetuned along with its classifier heads on NOAA’s hierarchical fish dataset. In Table 1, we report image-based accuracy and video-based accuracy, noted as *img* and *video* respectively, on both coarse (group) level and fine (species) level, noted as subscript *C* and *F* respectively.

Results are in Table. 1. Even though our proposed ‘**HCIL w/o m**’ method fixes the pre-trained backbone, which is not finetuned by the NOAA’s fish dataset, the performance is still comparable with work [9], that allows updating the backbone and calculating cross-entropy loss functions on both levels. This shows the backbone pre-trained on large public datasets without finetuning can indeed possess the strong ability to extract discriminative features even on new classes or datasets.

Table 1. Hierarchical Classification Setting

Method	<i>img_C</i>	<i>img_F</i>	<i>video_C</i>	<i>video_F</i>
[9]	92.0	82.9	96.5	91.2
HCIL w/o m	91.8	81.2	95.9	90.7

Table 2. Hierarchical Class Incremental Learning Setting

Method	<i>img_C</i>	<i>img_F</i>	<i>video_C</i>	<i>video_F</i>
[9] w/ m	80.5	65.7	81.1	70.4
HCIL w/o m	82.8	69.9	86.2	75.3
HCIL	91.0	80.4	92.1	82.8
HCIL w/ Swin	92.6	83.5	93.2	84.3

Hierarchical Class Incremental learning Setting In this setting, both our method, HCIL, and work [9] still utilize

ResNet-101 [19] backbone pre-trained on ImageNet-1k. Testing data are still the same as the previous setting and cover all species. However, training data are divided into three **tasks**. The first task includes one-third of the species within each group. The second task includes another one-third of the species within each group. The third task includes the rest species’ data. There are no overlapping classes between the three tasks. However, the ‘Skates’ and ‘Flatfish’ groups have only two species each so one Fine-SVM for each group is enough. As a result, all data from these two groups are included in the first task.

When training our HCIL model, on the first task data, the feature maps from the fixed pre-trained backbone are used to train six Coarse-SVMs and corresponding Fine-SVMs. Based on SVMs confidence scores and the herding method introduced in Section 3, HCIL constructs the memory. Next, when training on the following tasks, the memory’s feature maps are also used to train Coarse-SVMs and newly added Fine-SVMs. We set the total number of hard cases *n* to 200, and the total number of positive exemplars *m* to 1800 so that the memory size won’t increase as the incremental learning goes.

Work [9] does not have memory selection or classifiers expansion. For a fair comparison, when finetuning work [9]’s both pre-trained backbone and classifiers on later two tasks, we allow it to use the same size memory but randomly sampled from previously trained classes, denoted as ‘**[9] w/ m**’ in Table 2. When training on each task, its classifiers always output predictions over 31 species and calculate loss functions. We evaluate the final trained models on all testing data.

Results are in Table. 2. Compared with work [9] which is not designed for incremental learning, our HCIL model achieves significantly better performance. For the CIL memory ablation study, we remove CIL memory, noted as ‘**HCIL w/o m**’ in Table 2, and the performance drops dramatically but is still better than work [9] with randomly sampled memory. This ablation study shows the benefits of CIL memory selection module. It also tells that under incremental learning setting, updating the backbone and classifiers even with some randomly sampled memory, makes the deep model suffer from catastrophic forgetting. For the backbone ablation study, we replace ResNet-101 [19] with Swin-Transformer [20], noted as ‘**HCIL w/ Swin**’, and get the best performance.

5 Conclusions and Future

Our proposed HCIL model combines the advantages of both hierarchical classification and incremental learning. It consists of a fixed backbone pre-trained on large public datasets, a CIL memory selection module, and dynamic SVMs expansion. Our HCIL is also a backbone-agnostic approach. Future experiments may include dividing training data into more tasks to form longer incremental learning scenarios.

6 References

- [1] Jie Mei, Jingxi Yu, Suzanne Romain, Craig Rose, Kelsey Magrane, Graeme LeeSon, and Jenq-Neng Hwang, “Unsupervised severely deformed mesh reconstruction (dmr) from a single-view image,” 2022. 1
- [2] Shilpi Gupta, Nissreen Abu-Ghannam, Roberto Massini, Yaseen Mottiar, Illimar Altosaar, Micha Peleg, Mark D Normand, and Maria G Corradini, “Trends in application of imaging technologies to inspection of fish and,” . 1
- [3] Darren J White, C Svellingen, and Norval JC Strachan, “Automated measurement of species and length of fish by computer vision,” *Fisheries Research*, vol. 80, no. 2-3, pp. 203–210, 2006. 1
- [4] Jie Mei, Jenq-Neng Hwang, Suzanne Romain, Craig Rose, Braden Moore, and Kelsey Magrane, “Absolute 3d pose estimation and length measurement of severely deformed fish from monocular videos in longline fishing,” in *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2021, pp. 2175–2179. 1
- [5] Tsung-Wei Huang, Jenq-Neng Hwang, and Craig S Rose, “Chute based automated fish length measurement and water drop detection,” in *2016 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2016, pp. 1906–1910. 1
- [6] Kresimir Williams, Nathan Lauffenburger, Meng-Che Chuang, Jenq-Neng Hwang, and Rick Towler, “Automated measurements of fish within a trawl using stereo images from a camera-trawl device (camtrawl),” *Methods in Oceanography*, vol. 17, pp. 138–152, 2016. 1
- [7] Tsung-Wei Huang, Jenq-Neng Hwang, Suzanne Romain, and Farron Wallace, “Live tracking of rail-based fish catching on wild sea surface,” in *2016 ICPR 2nd Workshop on Computer Vision for Analysis of Underwater Imagery (CVAUI)*. IEEE, 2016, pp. 25–30. 1
- [8] Gaoang Wang, Jenq-Neng Hwang, Kresimir Williams, and George Cutter, “Closed-loop tracking-by-detection for rov-based multiple fish tracking,” in *2016 ICPR 2nd Workshop on Computer Vision for Analysis of Underwater Imagery (CVAUI)*. IEEE, 2016, pp. 7–12. 1
- [9] J. Mei, Jenq-Neng Hwang, S. Romain, Craig S. Rose, Braden Moore, and Kelsey Magrane, “Video-based hierarchical species classification for longline fishing monitoring,” in *ICPR Workshops*, 2020. 1, 2, 4, 1, 2, 4
- [10] Arthur Douillard, Alexandre Ramé, Guillaume Couairon, and Matthieu Cord, “Dytox: Transformers for continual learning with dynamic token expansion,” *arXiv preprint arXiv:2111.11326*, 2021. 1, 2
- [11] Sylvestre-Alvise Rebuffi, Alexander Kolesnikov, Georg Sperl, and Christoph H Lampert, “icarl: Incremental classifier and representation learning,” in *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, 2017, pp. 2001–2010. 1, 2
- [12] Saihui Hou, Xinyu Pan, Chen Change Loy, Zilei Wang, and Dahua Lin, “Learning a unified classifier incrementally via rebalancing,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 831–839. 1, 2
- [13] Francisco M Castro, Manuel J Marín-Jiménez, Nicolás Guil, Cordelia Schmid, and Karteek Alahari, “End-to-end incremental learning,” in *Proceedings of the European conference on computer vision (ECCV)*, 2018, pp. 233–248. 1, 2
- [14] Chrisantha Fernando, Dylan Banarse, Charles Blundell, Yori Zwols, David Ha, Andrei A Rusu, Alexander Pritzel, and Daan Wierstra, “Pathnet: Evolution channels gradient descent in super neural networks,” *arXiv preprint arXiv:1701.08734*, 2017. 2
- [15] Joan Serra, Didac Suris, Marius Miron, and Alexandros Karatzoglou, “Overcoming catastrophic forgetting with hard attention to the task,” in *International Conference on Machine Learning*. PMLR, 2018, pp. 4548–4557. 2
- [16] Max Welling, “Herding dynamical weights to learn,” in *Proceedings of the 26th Annual International Conference on Machine Learning*, 2009, pp. 1121–1128. 3, 3
- [17] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al., “Imagenet large scale visual recognition challenge,” *International journal of computer vision*, vol. 115, no. 3, pp. 211–252, 2015. 3
- [18] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick, “Microsoft coco: Common objects in context,” in *European conference on computer vision*. Springer, 2014, pp. 740–755. 3
- [19] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun, “Deep residual learning for image recognition,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778. 3, 4
- [20] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo, “Swin transformer: Hierarchical vision transformer using shifted windows,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 10012–10022. 3, 4