

---

# Documentation for The Noisy Ostracods Dataset

---

Jiamian Hu<sup>1</sup>, Yuanyuan Hong<sup>1</sup>, Yihua Chen<sup>1,2</sup>, He Wang<sup>1,3</sup>, Moriaki Yasuhara<sup>1</sup>

<sup>1</sup>The University of Hong Kong

<sup>2</sup>The University of Tokyo

<sup>3</sup>Nanjing Institute of Geology and Palaeontology

{jiamianh, u3001143, yihuaac, hw0701}@connect.hku.hk, yasuhara@hku.hk

## 1 Motivation and Backgrounds

The *Noisy Ostracods* dataset is a real-world taxonomy dataset characterized by various types of noise. It was created out of the need for a clean taxonomy dataset and the challenges we encountered during the cleaning process in our real use case. Our goal was to provide a benchmark for evaluating the performance of robust machine learning methods and label correction algorithms from a practical perspective. The imbalanced and fine-grained nature of the dataset introduces additional challenges to these methods.

The document is made by adapting the most relevant questions from datasheets for datasets[1] according to the property of our datasets.

## 2 Data Collection Detail

The dataset included Ostracods from the Hong Kong marine sediments collected over the past 10 years. The goal for collecting the ostracods is to exploring the quantitative correlation between common ostracod species composition and environmental factors[2, 3]. The details of collection process of sediments are available in the original works[3, 2, 4]. This document primarily focus on the collection process of the photos and the effort we made to ensure the quality of the dataset.

### 2.1 Collection process of Noisy Ostracods 2022

Ostracod samples from the original studies are stored in standard 60-grid microfossil slides. These slides are photographed using a Keyence-VHX-7000 microscope [5]. Images are captured at magnifications ranging from 40X to 80X. A sample slide image is shown in Figure 1. Initially, all images were taken at 50X magnification. In most cases, the photos resemble the surface\_SS6 slide in Table 1. However, when the image resolution is excessively high, the microscope automatically compresses the images. As illustrated in Table 1, for slide HK14TLH1C\_151\_152, the expected resolution at 50X magnification is approximately 23000\*9600. The actual resolution, however, is 11599\*4841, roughly half of the anticipated resolution. We conducted experiments on slide HK14TLH1C\_136\_137 to determine the optimal resolution. Even after compression, the 80X images were still larger in native resolution than the 40X images. However, capturing photos at 80X took four times longer, with no significant improvement in quality. Consequently, we decided to use 40X magnification for slides where 50X images were compressed. Following the preparation sequence shown in Figure 1, the slides are cropped into grids for easier processing. A typical 50X grid is sized around 1600\*1600 pixels. We then annotate the ostracods from the grid images using `labelImg` [6]. Since the original taxonomy record is organized per grid, we simply join the record with the annotation by grid number. Before this joining process, we cleaned the typographical errors in the original identification file. The typos addressed here were different from those present in the final version of the *Noisy Ostracods* dataset. Primarily, we corrected misspellings of genus and species names, such as changing *Xestol-*

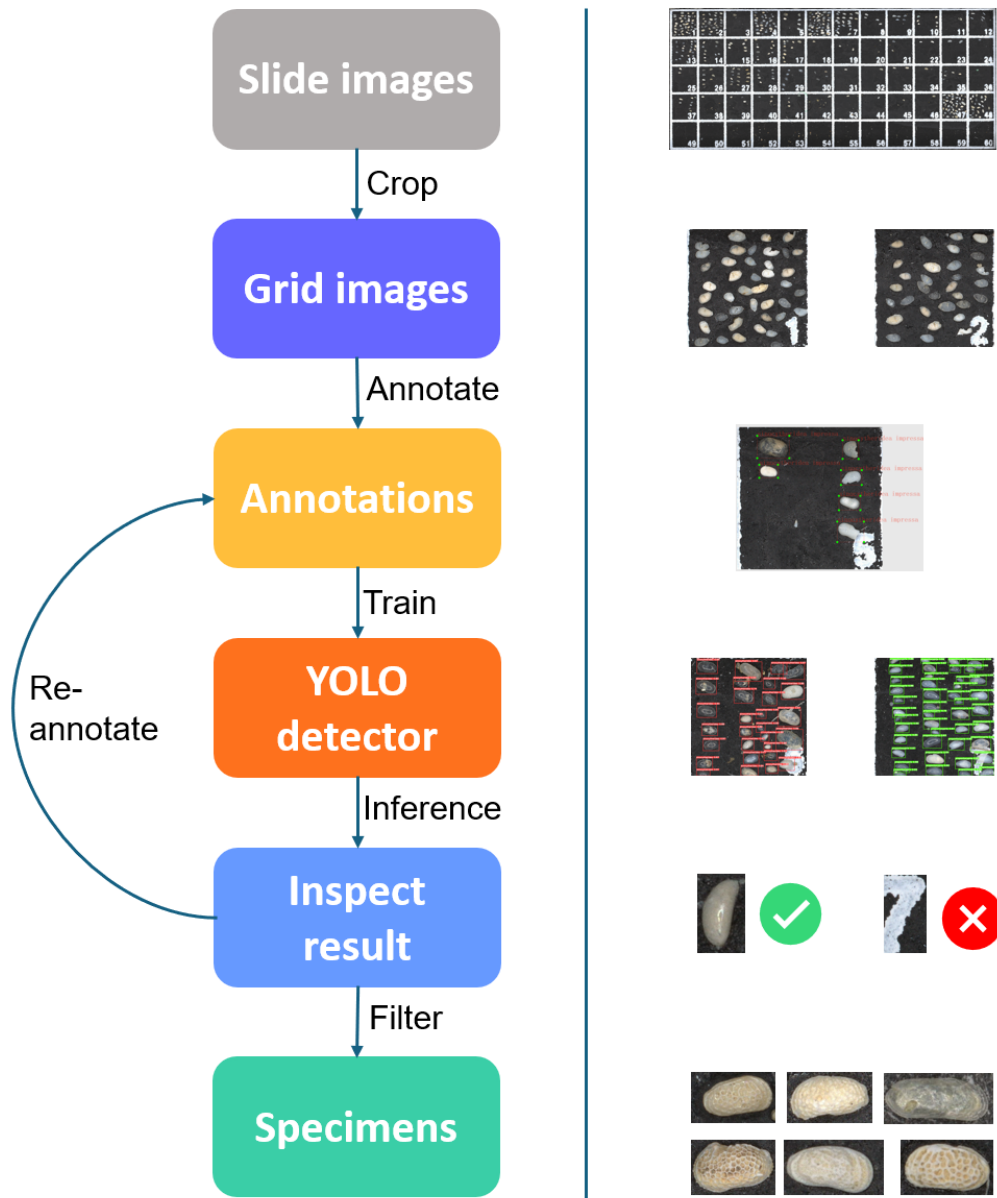


Figure 1: Illustrated collection process of the noisy ostracods dataset. Left column: simplified flow chart. Right column: illustrated examples.

*beris* to *Xestoleberis*. However, semantic typos, such as misnaming *Cytheroia* as *Cyprideis*, could not be detected at this stage because both genera are valid in ostracod taxonomy literature; such errors can only be identified through precise specimen checks.

We employ an iterative approach for annotation: we start by manually annotating around 180 grids and then train the initial detector using YOLO [7]. Using this new YOLO model, we annotate an additional 180 grids, correcting any errors in the bounding boxes. We then train the next model using all 360 annotated grids. This process is repeated until we have a training set of 6000 grids. The model trained on these 6000 grids is then used to annotate subsequent images, which are manually checked to eliminate errors. At this stage, we crop out the specimens to create an ostracod taxonomy dataset. However, the amount of error in the dataset proved to be non-negligible, necessitating multiple revisions. We retained the initial version of the dataset, making it available as *Noisy Ostracods 2022*.

Table 1: Resolution data for slides at different magnifications. The images with red actual resolutions are being compressed.

Slide	Magnification	Expected resolution	Actual resolution
HK14TLH1C_151_152	40X	18548*7764	18548*7764
HK14TLH1C_151_152	50X	23185*9705	11711*4929
HK14TLH1C_136_137	40X	18541*7721	18541*7721
HK14TLH1C_136_137	50X	23176*9651	11599*4841
HK14TLH1C_136_137	80X	37082*15422	18628*7759
surface_SS6	50X	21525*9402	21525*9402

## 2.2 Building the Noisy Dataset

When building the *Noisy Ostracods* dataset, we initially aimed to completely eliminate the errors present in the *Noisy Ostracods 2022* version. We identified several issues in the *Noisy Ostracods 2022* dataset:

- **Including non-ostracods:** Due to failures in the YOLO detector, some slides included non-ostracods. An example of this error, found in the *Xestoleberis* genus, is illustrated in Figure 1.
- **Missing ostracods:** Some ostracods were not detected, also due to failures in the YOLO detector.
- **Bad images:** Some slides photographed at an early stage were too bright, as shown in the left image of Figure 2. This issue was caused by incorrect camera settings.
- **Dual records in grids:** Ideally, each grid should contain only one species of ostracods. However, some grids had more than one species recorded due to practical reasons, such as the mixture of hard-to-distinguish species in the same grid. In creating the 2022 version, we skipped all such records.

Based on the errors identified in the 2022 version, we applied the following fixes:

1. Checked the annotation file: Re-annotated ostracods with incorrect bounding boxes and removed non-ostracods. Also annotated ostracods missed by the YOLO detector.
2. Checked the image files: Retook images that were too bright.
3. Consulted experts: Asked the experts who provided the original identifications to manually annotate the respective species on the images.

After applying these fixes, the majority of errors caused by the YOLO detector and camera settings were resolved. The number of YOLO detector failure found in current cleaning process on the *Noisy Ostracods* dataset is less than 20. However, unlike the thorough cleaning described in the main article, this round of inspection was conducted by non-experts. This led to some remaining errors, such as fragmentation errors and semantic typos, as detailed in the main article. Nevertheless, we believe that the errors in the current version of the *Noisy Ostracods* dataset have been minimized from a procedural perspective. The remaining noise primarily consists of hard-to-fix errors, such as expert misidentifications. Consequently, this version can serve as a valuable resource for researchers studying real-world noise and its effects on the performance of machine learning models, particularly in the context of trustworthy machine learning.

## 3 Dataset composition

### 3.1 Images

The *Noisy Ostracods* dataset comprises 71,466 images, while the *Noisy Ostracods 2022* dataset contains 68,458 images. The difference in image count is primarily due to the addition of specimens from grids with mixed records. The image files are organized by their *annotated* species. If a



Figure 2: Illustration of camera parameter failure error. Left: a grid image that is too bright. Right: the re-took image.

specimen’s species is unconfirmed in the original identification, it is moved to the genus class *genus unidentified*. A concrete example of the file structure of the dataset is shown in the file tree below.

```

Root Folder
├── Noisy_ostracods
│   ├── sinocythere sinensis
│   │   └── images
│   ├── sinocythere unidentified
│   │   └── images
│   └── other species
│       └── images
  
```

In the file tree, the *Sinocythere unidentified* included the specimens of *Sinocythere* with unconfirmed species. All versions of the dataset adhere to this structure. To preserve the best quality of the images, all RGB images are stored in uncompressed .tif format. In addition to ostracod images, we included 9,117 negative samples by randomly cropping backgrounds that do not belong to ostracod bounding boxes and by adding random photos from the Endless Forams dataset [8] as a negative class. We also have 876 images without labels, as the original records for these grids marked their genus and species with a ? (question mark). Some of these can be identified to existing genera and species. In total, the *Noisy Ostracods* dataset contains 81,459 labeled images.

### 3.2 Labels

The labels are stored in .csv files to enable flexibility. Each record in the label file includes two entries: the image path and the label number. For example, a row in the label file might look like `alocopocythere goujoni/HKUV12_465_467_38_ind6.tif,2`. This line indicates that the file `HKUV12_465_467_38_ind6.tif` located in the `alocopocythere goujoni` folder has an image label of 2. Each label file is accompanied by a corresponding guidance file that maps the image label number to its genus/species name in string format. These guidance files are .txt files that list the genus/species names, with the line number corresponding to the label number in the label files.

This structure provides the flexibility to update the genus/species of the images by simply changing the image label number in the corresponding label file whenever an error is detected. If a typographical or semantic error is found, we can correct it by updating the corresponding guidance and label files. Such a structure ensures that any necessary adjustments can be made efficiently, maintaining the accuracy and integrity of the dataset.

The provided label files and label guidance files are:

- **Noisy label files:** ostracods\_genus\_final\_train.csv, ostracods\_genus\_final\_val.csv, ostracods\_genus\_final\_test.csv, ostracods\_species\_final\_train.csv, ostracods\_species\_final\_val.csv, and ostracods\_species\_final\_test.csv. These files contain the training, validation, and test splits for the *Noisy Ostracods* dataset at the genus and species levels.
- **Clean label files:** ostracods\_genus\_clean\_test.csv and ostracods\_genus\_clean\_val.csv. These files contain the cleaned test and validation splits at the genus level.
- **Guidance files:** ostracods\_genus\_final\_guide.txt and ostracods\_species\_final\_guide.txt. The guidance files include 139 rows for species and 79 rows for genus. The same guidance is used for both noisy and clean data.

The clean label files are slightly smaller than their noisy counterparts because some noisy files have been deleted. At the current stage, we are beginning the comprehensive cleaning of the dataset at the genus level and have inspected some problematic species. The known issues are listed in Table 2.

As shown in the table, the inconsistent usage of *Confer* (cf.) and *Affinis* (aff.) across different projects has introduced many pseudo-classes. Additionally, a typo of *Spinileberis quadriaculeata* was introduced during the re-annotation of grids containing multiple species. Resolving issues with some challenging species, such as *Pistocythereis subovata* and *Pistocythereis bradyi*, may result in multi-class labels, as it is nearly impossible to distinguish these species from a single image.

Table 2: Full list of genus and species in the *Noisy Ostracods* dataset. Known problems are listed.

<i>Genus</i>	<i>Species</i>	<b>Count</b>	<b>Known problems</b>
<i>aglaiocypris</i>	-	426	-
<i>alataconcha</i>	<i>alataconcha cf. pterogona</i>	8	-
<i>alataconcha</i>	-	1	-
<i>alocopocythere</i>	<i>alocopocythere goujoni</i>	285	-
<i>alocopocythere</i>	<i>alocopocythere kendengensis</i>	67	-
<i>alocopocythere</i>	<i>alocopocythere profusa</i>	108	-
<i>alocopocythere</i>	-	5	-
<i>argilloecia</i>	<i>argilloecia lunata</i>	23	-
<i>argilloecia</i>	-	29	-
<i>atjehella</i>	<i>atjehella cf. semiplicata</i>	68	Possible duplicate with <i>atjehella semiplicata</i>
<i>atjehella</i>	<i>atjehella kingmai</i>	4	-
<i>atjehella</i>	<i>atjehella semiplicata</i>	2	Possible duplicate with <i>atjehella cf. semiplicata</i>
<i>aurila</i>	<i>aurila aff. corniculata</i>	11	Possible duplicate with <i>aurila cf. corniculata</i>
<i>aurila</i>	<i>aurila cf. corniculata</i>	4	Possible duplicate with <i>aurila aff. corniculata</i>
<i>aurila</i>	<i>aurila cf. disparata</i>	245	Possible merge with <i>aurila disparata</i>
<i>aurila</i>	<i>aurila cf. hataii</i>	83	-
<i>aurila</i>	<i>aurila cf. uranouchiensis</i>	5	-
<i>aurila</i>	<i>aurila disparata</i>	30	Possible merge with <i>aurila cf. disparata</i>
<i>aurila</i>	-	314	-
<i>bicornucythere</i>	<i>bicornucythere bisanensis</i>	3177	-
<i>bicornucythere</i>	-	1633	-
<i>bythoceratina</i>	<i>bythoceratina callidictya</i>	1	Possible merge with <i>bythoceratina cf. callidictya</i>
<i>bythoceratina</i>	<i>bythoceratina cassidoidea</i>	13	-
<i>bythoceratina</i>	<i>bythoceratina cf. angulata</i>	1	-
<i>bythoceratina</i>	<i>bythoceratina cf. callidictya</i>	21	Possible merge with <i>bythoceratina callidictya</i>

<b>Genus</b>	<b>Species</b>	<b>Count</b>	<b>Known problems</b>
<i>bythoceratina</i>	<i>bythoceratina cf. orientalis</i>	11	-
<i>bythoceratina</i>	<i>bythoceratina cf. robusta</i>	2	Possibile merge with <i>bythoceratina robusta</i>
<i>bythoceratina</i>	<i>bythoceratina orientalis</i>	4	-
<i>bythoceratina</i>	<i>bythoceratina robusta</i>	9	Possibile merge with <i>bythoceratina cf. robusta</i>
<i>bythoceratina</i>	<i>bythoceratina sheyangensis</i>	32	-
<i>bythoceratina</i>	-	266	-
<i>bythocypris</i>	-	7	-
<i>bythocythere</i>	-	20	-
<i>callistocythere</i>	<i>callistocythere aff. reticulata</i>	9	Possibile duplicate with <i>callistocythere cf. reticulata</i>
<i>callistocythere</i>	<i>callistocythere aff. undulatifacialis</i>	263	Possibile merge with <i>callistocythere undulatifacialis</i>
<i>callistocythere</i>	<i>callistocythere asiatica</i>	16	-
<i>callistocythere</i>	<i>callistocythere cf. multirugosa</i>	19	-
<i>callistocythere</i>	<i>callistocythere cf. nipponica</i>	7	-
<i>callistocythere</i>	<i>callistocythere cf. reticulata</i>	2	Possibile duplicate with <i>callistocythere aff. reticulata</i>
<i>callistocythere</i>	<i>callistocythere undata</i>	5	-
<i>callistocythere</i>	<i>callistocythere undulatifacialis</i>	10	Possibile merge with <i>callistocythere aff. undulatifacialis</i>
<i>callistocythere</i>	-	45	-
<i>cathacythere</i>	<i>cathacythere reticulata</i>	6	-
<i>caudites</i>	-	39	-
<i>cletocythereis</i>	-	3	-
<i>copytus</i>	<i>copytus posterosulcus</i>	708	-
<i>coquimba</i>	<i>coquimba cf. ishizakii</i>	12	-
<i>coquimba</i>	-	10	-
<i>cornucoquimba</i>	<i>cornucoquimba cf. gibboidea</i>	491	-
<i>cornucoquimba</i>	<i>cornucoquimba leizhouensis</i>	77	-
<i>cornucoquimba</i>	<i>cornucoquimba pustulata</i>	5	-
<i>cornucoquimba</i>	-	69	-
<i>cyprideis</i>	-	12	According to the photos, should be <i>cytheroideis</i>
<i>cythere</i>	<i>cythere omotenipponica</i>	67	-
<i>cythere</i>	-	48	-
<i>cytherelloidea</i>	<i>cytherelloidea cingulata</i>	2	-
<i>cytherelloidea</i>	-	2	-
<i>cytherelloidea</i>	<i>cytherelloidea yingliensis</i>	3	-
<i>cytheroideis</i>	<i>cytheroideis leizhouensis</i>	199	-
<i>cytheroideis</i>	-	553	-
<i>cytheropteron</i>	<i>cytheropteron cf. ignobilis</i>	122	-
<i>cytheropteron</i>	<i>cytheropteron higashikawai</i>	1	-
<i>cytheropteron</i>	<i>cytheropteron miurense</i>	349	-
<i>cytheropteron</i>	-	114	-
<i>cytherura</i>	-	5	-
<i>darwinula</i>	-	2	-
<i>eucythere</i>	-	4	-
<i>hanaiborchella</i>	<i>hanaiborchella cf. miurense</i>	77	-
<i>hanaiborchella</i>	-	7	-
<i>haplocythereidea</i>	<i>haplocythereidea agilis</i>	5	According to the photos, should be <i>neocyprideis</i>
<i>haplocythereidea</i>	<i>haplocythereidea cf. agilis</i>	2	According to the photos, should be <i>neocyprideis</i>

<b>Genus</b>	<b>Species</b>	<b>Count</b>	<b>Known problems</b>
<i>hemicytheridea</i>	<i>hemicytheridea cancellata</i>	12	According to the photos, majority are <i>bicornucythere bisanensis</i>
<i>hemicytheridea</i>	<i>hemicytheridea reticulata</i>	387	-
<i>hemicytheridea</i>	-	41	-
<i>hemicytherura</i>	<i>hemicytherura cf. cuneata</i>	121	Possibile merge with <i>hemicytherura cuneata</i>
<i>hemicytherura</i>	<i>hemicytherura cf. kajiyamai</i>	10	Possibile merge with <i>hemicytherura kajiyamai</i>
<i>hemicytherura</i>	<i>hemicytherura cuneata</i>	57	Possibile merge with <i>hemicytherura cf. cuneata</i>
<i>hemicytherura</i>	<i>hemicytherura kajiyamai</i>	14	Possibile merge with <i>hemicytherura cf. kajiyamai</i>
<i>hemicytherura</i>	-	131	-
<i>hemikrithe</i>	<i>hemikrithe orientalis</i>	158	-
<i>hemikrithe</i>	<i>hemikrithe reticulata</i>	2	Typo of <i>hemicytheridea reticulata</i>
<i>hemikrithe</i>	-	15	-
<i>hermanites</i>	<i>hermanites bicostata</i>	2	-
<i>javanella</i>	<i>javanella kendengensis</i>	8	-
<i>javanella</i>	-	1	-
<i>keijella</i>	<i>keijella apta</i>	1	Possibile merge with <i>keijella cf. apta</i>
<i>keijella</i>	<i>keijella cf. apta</i>	10	Possibile merge with <i>keijella apta</i>
<i>keijella</i>	<i>keijella demissa</i>	108	Typo of <i>keijia demissa</i>
<i>keijella</i>	<i>keijella kloempritensis</i>	2915	-
<i>keijella</i>	-	12	-
<i>keijia</i>	<i>keijia demissa</i>	9	-
<i>kotoracythere</i>	<i>kotoracythere doratus</i>	2	-
<i>krithe</i>	<i>krithe japonica</i>	28	-
<i>krithe</i>	-	1	-
<i>leptocythere</i>	<i>leptocythere pulchra</i>	6	-
<i>leptocythere</i>	-	3	-
<i>loxoconcha</i>	<i>loxoconcha aff. hattorii</i>	2	Possibile merge with <i>loxoconcha hattorii</i>
<i>loxoconcha</i>	<i>loxoconcha aff. uranouchiensis</i>	2	Possibile duplicate with <i>loxoconcha cf. uranouchiensis</i>
<i>loxoconcha</i>	<i>loxoconcha aff. viva</i>	7	Possibile duplicate with <i>loxoconcha cf. viva</i>
<i>loxoconcha</i>	<i>loxoconcha cf. kattoi</i>	82	Possibile merge with <i>loxoconcha kattoi</i>
<i>loxoconcha</i>	<i>loxoconcha cf. kosugi</i>	5	-
<i>loxoconcha</i>	<i>loxoconcha cf. uranouchiensis</i>	6	Possibile duplicate with <i>loxoconcha aff. uranouchiensis</i>
<i>loxoconcha</i>	<i>loxoconcha cf. viva</i>	1	Possibile duplicate with <i>loxoconcha aff. viva</i>
<i>loxoconcha</i>	<i>loxoconcha epeterseni</i>	149	-
<i>loxoconcha</i>	<i>loxoconcha hattorii</i>	11	Possibile merge with <i>loxoconcha cf. hattorii</i>
<i>loxoconcha</i>	<i>loxoconcha japonica</i>	221	-
<i>loxoconcha</i>	<i>loxoconcha kattoi</i>	145	Possibile merge with <i>loxoconcha cf. kattoi</i>
<i>loxoconcha</i>	<i>loxoconcha malayensis</i>	1177	-
<i>loxoconcha</i>	<i>loxoconcha ocellata</i>	1	-
<i>loxoconcha</i>	<i>loxoconcha pulchra</i>	16	-
<i>loxoconcha</i>	-	622	-

<b>Genus</b>	<b>Species</b>	<b>Count</b>	<b>Known problems</b>
<i>loxoconcha</i>	<i>loxoconcha zhejiangensis</i>	231	-
<i>macrocypris</i>	-	1	-
<i>microcythere</i>	-	8	-
<i>morkhovenia</i>	<i>morkhovenia inconspicua</i>	8	-
<i>munseyella</i>	<i>munseyella japonica</i>	341	-
<i>munseyella</i>	<i>munseyella oblonga</i>	3	-
<i>munseyella</i>	-	304	-
<i>neocyprideis</i>	-	13	-
<i>neocytheretta</i>	<i>neocytheretta elongata</i>	2	Typo of <i>neosinocythere elongata</i>
<i>neocytheretta</i>	<i>neocytheretta faceta</i>	160	-
<i>neocytheretta</i>	<i>neocytheretta snellii</i>	32	-
<i>neocytheretta</i>	-	408	-
<i>neocytheromorpha</i>	<i>neocytheromorpha regalis</i>	26	-
<i>neocytheromorpha</i>	-	34	-
<i>neomonoceratina</i>	<i>neomonoceratina delicata</i>	9650	-
<i>neomonoceratina</i>	<i>neomonoceratina elongata</i>	1	Typo of <i>neosinocythere elongata</i> . However, according to photo, should be <i>spinileberis</i> .
<i>neonesidea</i>	<i>neonesidea elegans</i>	210	-
<i>neonesidea</i>	<i>neonesidea oligodentata</i>	129	-
<i>neonesidea</i>	-	156	-
<i>neopellucistoma</i>	<i>neopellucistoma inflatum</i>	17	-
<i>neopellucistoma</i>	-	7	-
<i>neosinocythere</i>	<i>neosinocythere elongata</i>	1666	-
<i>neosinocythere</i>	-	799	-
<i>nipponocythere</i>	<i>nipponocythere bicarinata</i>	256	-
<i>nipponocythere</i>	<i>nipponocythere delicata</i>	632	-
<i>nipponocythere</i>	-	215	-
<i>orionina</i>	<i>orionina yongleensis</i>	1	Species may be wrong
<i>pacambocythere</i>	-	1	-
<i>palmenella</i>	-	8	-
<i>paracathaycythere</i>	<i>paracathaycythere cf. costaereticulata</i>	1	-
<i>paracypris</i>	-	71	-
<i>paracytheridea</i>	<i>paracytheridea tschoppi</i>	3	-
<i>paracytherois</i>	<i>paracytherois cf. acuminata</i>	45	-
<i>paracytherois</i>	<i>paracytherois cf. tosaensis</i>	2	-
<i>paracytherois</i>	-	51	-
<i>paradoxostomatid</i>	-	1257	-
<i>parakrithe</i>	<i>parakrithe cf. elongata</i>	1	Could be <i>parakrithella</i>
<i>parakrithe</i>	<i>parakrithe japonica</i>	1	-
<i>parakrithella</i>	<i>parakrithella cf. pseudadonta</i>	26	Possibile merge with <i>parakrithella pseudadonta</i>
<i>parakrithella</i>	<i>parakrithella pseudadonta</i>	2	Possibile merge with <i>parakrithella cf. pseudadonta</i>
<i>parakrithella</i>	-	8	-
<i>phlyctocythere</i>	<i>phlyctocythere japonica</i>	328	-
<i>phlyctocythere</i>	-	142	-
<i>pistocythereis</i>	<i>pistocythereis aff. miaoliensis</i>	12	-
<i>pistocythereis</i>	<i>pistocythereis bradyformis</i>	1015	-
<i>pistocythereis</i>	<i>pistocythereis bradyi</i>	1758	Possibile merge with <i>pistocythereis cf. bradyi</i>
<i>pistocythereis</i>	<i>pistocythereis cf. bradyi</i>	41	Possibile merge with <i>pistocythereis bradyi</i>



<b>Genus</b>	<b>Species</b>	<b>Count</b>	<b>Known problems</b>
<i>pistocythereis</i>	<i>pistocythereis cf. subovata</i>	11	Possibile merge with <i>pistocythereis subovata</i>
<i>pistocythereis</i>	<i>pistocythereis euplectella</i>	74	-
<i>pistocythereis</i>	<i>pistocythereis subovata</i>	56	Possibile merge with <i>pistocythereis cf. subovata</i> , very hard to be distinguished from <i>pistocythereis bradyi</i> visually.
<i>pistocythereis</i>	-	1884	-
<i>pontocythere</i>	<i>pontocythere cf. subjaponica</i>	44	-
<i>pontocythere</i>	-	96	-
<i>propontocypris</i>	<i>propontocypris clara</i>	1	-
<i>propontocypris</i>	-	5021	-
<i>pseudocythere</i>	<i>pseudocythere cf. caudata</i>	1	-
<i>pseudocythere</i>	-	5	-
<i>robustauria</i>	<i>robustauria cf. ishizakii</i>	34	-
<i>robustauria</i>	<i>robustauria salebroza</i>	2	-
<i>robustauria</i>	-	5	-
<i>semicytherura</i>	<i>semicytherura cf. miurensis</i>	12	-
<i>semicytherura</i>	<i>semicytherura cf. undata</i>	1	-
<i>semicytherura</i>	<i>semicytherura cf. wakamurasaki</i>	3	-
<i>semicytherura</i>	<i>semicytherura indonesiana</i>	6	-
<i>semicytherura</i>	-	120	-
<i>sinocythere</i>	<i>sinocythere dongtaiensis</i>	3	-
<i>sinocythere</i>	<i>sinocythere sinensis</i>	87	-
<i>sinocythere</i>	-	183	-
<i>sinocytheridea</i>	<i>sinocytheridea impressa</i>	22429	-
<i>sinoleberis</i>	<i>sinoleberis cf. tosaensis</i>	3	-
<i>spinileberis</i>	<i>spinileberis quadriaculeata</i>	1484	-
<i>spinileberis</i>	<i>spinileberis quadriculeata</i>	1	Typo of <i>spinileberis quadriaculeata</i>
<i>spinileberis</i>	<i>spinileberis rhomboidalis</i>	342	-
<i>spinileberis</i>	-	111	-
<i>stigmatocythere</i>	<i>stigmatocythere aff. roesmani</i>	8	Possibile duplicate with <i>stigmatocythere cf. roesmani</i>
<i>stigmatocythere</i>	<i>stigmatocythere bona</i>	78	-
<i>stigmatocythere</i>	<i>stigmatocythere cf. roesmani</i>	17	Possibile duplicate with <i>stigmatocythere aff. roesmani</i>
<i>stigmatocythere</i>	<i>stigmatocythere costa</i>	58	-
<i>stigmatocythere</i>	<i>stigmatocythere kingmai</i>	64	-
<i>stigmatocythere</i>	<i>stigmatocythere roesmani</i>	767	Possibile merge with <i>stigmatocythere cf. roesmani</i> and <i>stigmatocythere aff. roesmani</i>
<i>stigmatocythere</i>	-	270	-
<i>tanella</i>	<i>tanella gracillis</i>	81	-
<i>tanella</i>	-	29	-
<i>trachyleberididae</i>	-	25	-
<i>trachyleberis</i>	-	2	-
<i>triebelina</i>	<i>triebelina aff. sertata</i>	8	-
<i>xestoleberis</i>	<i>xestoleberis hanaii</i>	222	-
<i>xestoleberis</i>	<i>xestoleberis suetsumuhana</i>	2	-
<i>xestoleberis</i>	-	860	-
<i>xiphichilus</i>	-	87	-
<b>Grand Total</b>		<b>71466</b>	

### 3.3 Version Difference

The images in the *Noisy Ostracods* dataset and its 2022 version do not have a one-to-one correspondence. This means that images with the same name may not be the same across versions. The reason is straightforward: we performed per-image re-labeling when building the *Noisy Ostracods* dataset based on the 2022 version to identify and correct possible false labels and missing labels.

We provide the original train, test, and validation splits from 2022. However, this version has the following issues:

- **Random Split:** The splits for genus and species are not consistent. This means that the test and validation images at the genus and species levels do not match. Some images in the test set for genus identification may appear in the training set for species identification.
- **Minor Class Elimination:** All minor classes with fewer than 10 images were removed rather than being moved to the training set, unlike in the current version.

Please consider these differences when comparing methods using the 2022 version of the dataset.

## 4 Data Availability and Maintenance

We are making the Noisy Ostracods datasets and the Noisy Ostracods 2022 datasets available online.

Noisy Ostracods images: [Noisy Ostracods images: Click to download](#)

Noisy Ostracods 2022 images: [Click to download](#)

Croissant[9] metadata: [Click to download](#)

Code: [https://github.com/H-Jamieu/Noisy\\_ostracods](https://github.com/H-Jamieu/Noisy_ostracods)

We are actively following the taxonomy changes in ostracods and will revise the taxonomy accordingly. New samples from our studies will be annotated and published after verification. Any future updates on the dataset will be included in the code repository. We are currently cleaning the entire dataset and discussing potential changes in some ambiguous species. The link to download the fully cleaned dataset will be released on the code repository. One copy of the dataset will be hosted in the Data Repository of The University of Hong Kong to ensure long term preservation.

## 5 Dataset Licence and Author Responsibility

The authors hereby declare that they bear all responsibility in case of violation of rights, including but not limited to intellectual property rights, data privacy rights, and any other applicable laws. They confirm that the data provided in this work is licensed under the Creative Commons Attribution 4.0 International (CC BY 4.0) license[10].

## 6 Additional Contents

### 6.1 Error Correction

During the manual cleaning of the dataset, we initially reported the discovery of two new genera, including *Pseudocythere* and another unnamed genus. Upon further inspection, we realized that *Pseudocythere* was already present in our dataset. This oversight occurred due to an error in our initial data review process. We deeply apologize for this mistake and any confusion it may have caused.

Furthermore, we are not certain if the specimen is indeed *Pseudocythere* since the key identification parts of the specimen are broken. We are still checking the relevant taxonomic materials to address this issue at the time of writing. Meanwhile, the other specimen is confirmed to be a new genus.

As a result, the images are still being deleted and will not affect the reported results. The two images in question are `surface_VS13_6_ind4.tif` and `rawSample_B1a_29_ind2.tif`, for users of the dataset. The corresponding paragraph has been revised in our latest version of the main article.

## 6.2 Additional Hyper-parameters

In the main article, we forgot to mention the size of the images for training. We used 224\*224 for all the training. For Co-teaching[11], Co-teaching+[12], loss-clip, mixup-cutmix[13, 14], CL[15] and CE training, we used FP16 scaling to accelerate training. As for Divide-mix[16], we follow the official implementation not scaling the training to FP16. For embedding calculated using CLIP-ViT-L-14@336[17], we scaled the images to 336\*336 to preserve the consistency. All other embeddings are calculated using the image size the models trained on for SimiFeat[18].

## References

- [1] Timnit Gebru et al. “Datasheets for Datasets”. In: (2018). URL: <http://arxiv.org/abs/1803.09010>.
- [2] Y. Hong et al. “Baseline for ostracod-based northwestern Pacific and Indo-Pacific shallow-marine paleoenvironmental reconstructions: ecological modeling of species distributions”. In: *Biogeosciences* 16.2 (2019), pp. 585–604. DOI: 10.5194/bg-16-585-2019.
- [3] Yuanyuan Hong. “Hong Kong shallow marine benthic ecosystem history : conservation paleoecology approach based on microfossil ostracods”. 2016. URL: <http://hdl.handle.net/10722/240648>.
- [4] Emma Cieslak-Jones. *Ostracods: Fossil time machines into past and future ecosystems*. NatureVolve digital magazine. 2022. URL: [https://www.su.se/polopoly\\_fs/1.633742.1666965417!/menu/standard/file/Yasuhara%26Hong2022NatureVolve.pdf](https://www.su.se/polopoly_fs/1.633742.1666965417!/menu/standard/file/Yasuhara%26Hong2022NatureVolve.pdf).
- [5] Keyence Corporation. *VHX-7000 Digital Microscope*. Available online: <https://www.keyence.com/products/microscope/digital-microscope/vhx-7000/>. 2022. URL: <https://www.keyence.com/products/microscope/digital-microscope/vhx-7000/>.
- [6] Tzutalin. *LabelImg*. <https://github.com/tzutalin/labelImg>. 2015.
- [7] Glenn Jocher, Ayush Chaurasia, and Jing Qiu. *Ultralytics YOLO*. Version 8.0.0. 2023. URL: <https://github.com/ultralytics/ultralytics>.
- [8] AY Hsiang et al. “Endless Forams: >34,000 modern planktonic foraminiferal images for taxonomic training and automated species recognition using convolutional neural networks”. In: *Paleoceanography & Paleoclimatology* 34 (2019). DOI: 10.1029/2019PA003612.
- [9] Mubashara Akhtar et al. “Croissant: A Metadata Format for ML-Ready Datasets”. In: New York, NY, USA: Association for Computing Machinery, 2024. DOI: 10.1145/3650203.3663326. URL: <https://doi.org/10.1145/3650203.3663326>.
- [10] Creative Commons. *CC BY 4.0 License*. Accessed: 2024-06-12. 2013. URL: <https://creativecommons.org/licenses/by/4.0/>.
- [11] Bo Han et al. “Co-teaching: Robust training of deep neural networks with extremely noisy labels”. In: *NeurIPS*. 2018, pp. 8535–8545.
- [12] Xingrui Yu et al. “How does Disagreement Help Generalization against Label Corruption?” In: *International Conference on Machine Learning*. 2019, pp. 7164–7173.
- [13] Hongyi Zhang et al. “mixup: Beyond Empirical Risk Minimization”. In: *International Conference on Learning Representations*. 2018. URL: <https://openreview.net/forum?id=r1Ddp1-Rb>.
- [14] Sangdoon Yun et al. “CutMix: Regularization Strategy to Train Strong Classifiers With Localizable Features”. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*. 2019.
- [15] Curtis G. Northcutt, Lu Jiang, and Isaac L. Chuang. “Confident Learning: Estimating Uncertainty in Dataset Labels”. In: *Journal of Artificial Intelligence Research (JAIR)* 70 (2021), pp. 1373–1411.
- [16] Junnan Li, Richard Socher, and Steven C.H. Hoi. “DivideMix: Learning with Noisy Labels as Semi-supervised Learning”. In: *International Conference on Learning Representations*. 2020.
- [17] Alec Radford et al. *Learning Transferable Visual Models From Natural Language Supervision*. 2021. arXiv: 2103.00020 [cs.CV].
- [18] Zhaowei Zhu, Zihao Dong, and Yang Liu. “Detecting Corrupted Labels Without Training a Model to Predict”. In: *Proceedings of the International Conference on Machine Learning*. Vol. 139. PMLR, 2022, pp. 12320–12330. arXiv: 2110.06283 [cs.LG]. URL: <http://proceedings.mlr.press/v139/zhu22f.html>.