

1 SUPPLEMENTARY MATERIALS

2 *Test 1 Methods – Novel image matching for known individuals*

3 The first test conducted involved evaluating the algorithm’s predictive accuracy with 100
4 images randomly selected of known individuals. These images were of animals that had a known
5 ID in the historical Risso’s dolphin catalog but were not the same image in the historical catalog
6 that had been provided to the dorsal fin matching algorithm. This approach was to ensure the
7 algorithm was tested on previously unseen images of known individuals. The purpose of this test
8 was to evaluate how the algorithm matched novel images of known individuals, determine where
9 the correct matches appeared in the output queue, assess the similarity scores, and identify
10 potential missed matches. A missed match was defined as one unique ID in the historical catalog
11 matching another unique ID in the historical catalog, where one individual was given two
12 different IDs and treated as two different individuals.

13 Of the 100 random images pulled, 95 were cropped into a square shape centered around
14 the dorsal fin and edited to enhance contrast and bring details in pigmentation and dorsal ridging.
15 The remaining five images were left uncropped and unedited in order to evaluate whether image
16 formatting and post-processing influenced the algorithm’s performance. All images were
17 organized into a single folder for batch upload to the matching platform.

18 One-way ANOVAs and linear regressions were used to examine whether image quality
19 or distinctiveness metrics were associated with the algorithm’s top similarity score for these
20 images, and the MAP score was calculated. Residual diagnostics on ANOVA indicated that
21 assumptions were reasonably met (Shapiro-Wilk $p = 0.7532$; Levene’s test $p = 0.177$), and non-
22 parametric alternatives produced consistent results.

23 *Test 2a Methods – Full historical catalog*

24 Following the results of Test 1, which revealed several missed matches, every image in
25 the historical catalog (n = 2,421 images) was compared against itself in a many-to-many match
26 run. The purpose of this test was to identify missed matches across the entire historical dataset,
27 and assess the typical queue position in which the correct matches appeared, using a larger and
28 more diverse sample set. These were all a “same image” match (i.e., the exact image that was
29 already in the historical catalog was matched to itself, with the expectation that the exact same
30 image should appear first in the queued match options). For manual review, only the top five
31 candidates returned by the algorithm were evaluated as possible matches, based on the results
32 from Test 1 showing that most matches appeared within the first two results.

33 This test was meant to validate the consistency of the algorithm’s performance with
34 identical images and document the typical queue position of missed matches, so no statistical
35 comparisons were made for image quality or distinctiveness within Test 2a. These image quality
36 and distinctiveness characteristics are outlined in Test 2b, which focuses specifically on the
37 missed matches identified during Test 2a.

38 *Test 2b Methods – Missed match rank analysis*

39 While the 2,421 images taken through the algorithm revealed a number of missed
40 matches, every instance of a missed match was not necessarily found with every photo of that
41 individual ID in the first five options presented. Each individual ID could have more than one
42 image in the historical catalog. There could be up to 12 images if it was a Risso’s dolphin that
43 had a variety of photos from ahead, laterally and behind, was a tagged animal, or had mostly
44 poor quality images of the dorsal for measurement but a few good quality images showing other
45 parts of the body for eye-based matching.

46 A detailed investigation was conducted into the missed matches found in Test 2a, by
47 running each of the algorithm-found missed matches through the algorithm again, this time
48 through the entirety of presented possible matches, to determine whether, or where, the known
49 missed match appeared. This allowed further evaluation of how far users may need to search
50 through the algorithm's queue to find correct missed matches.

51 A total of 203 images were re-tested in this step. One-way ANOVAs and linear
52 regressions were used to examine whether image quality or distinctiveness metrics were
53 associated with the algorithm's top similarity score for these images, and the MAP score was
54 calculated. Residuals were approximately normal (Shapiro-Wilk $p \approx 0.015$), and non-parametric
55 Kruskal-Wallis tests for each factor yielded consistent significance patterns.

56 *Test 3 Methods – Updating Risso's Dolphin Catalog*

57 During Tagless BRS and opportunistic survey efforts in 2022-2024, Cascadia collected
58 images of Risso's dolphins from four different vessels operating in the vicinity of Catalina
59 Island. Across 18 days of survey effort in the area (comprising of 27 boat days), 108 total
60 sightings of Risso's dolphins were documented. The images collected from these sightings were
61 compiled and assigned image quality and distinctiveness metrics through the Risso's database by
62 a human cataloger. Unlike previous catalog updates, these fins were not measured using ImageJ,
63 as this step was no longer necessary given the algorithm's prior success. These images were then
64 run against the multi-species algorithm for matching with the goal of updating Cascadia's
65 Risso's dolphin catalog. Each image was compared against the first 10 match options presented
66 by the algorithm's queue. To evaluate which factors may influence algorithm performance, linear
67 regression and ANOVA were applied to test the effects of image quality and individual
68 distinctiveness metrics on the algorithm similarity score, and residual diagnostics indicated that

69 model assumptions were reasonably met, with residuals approximately normal (Shapiro-Wilk $p =$
70 0.75) and homogeneity of variance was satisfied (Levene's test $p = 0.18$). The MAP score was
71 calculated to quantify overall accuracy and to account for cases where the algorithm may not
72 identify a known individual, the MAP was recalculated by introducing a dummy "new
73 individual" prediction at a score of 0.5. If this dummy prediction ranked first, meaning the top
74 algorithm score was <0.5 , the precision for that image was 1.0; if it ranked second, meaning the
75 top score was >0.5 , the precision was 0.5. This adjustment allowed all images, including those
76 without confirmed matches, to contribute to the overall MAP score. Unlike earlier catalog
77 updates, where the matching process involved two human reviewers (one conducting the initial
78 search for matches, and a second reviewing any unfound matches to ensure the first person did
79 not miss any), here only a single human matcher reviewed the results. Given the algorithm's
80 demonstrated reliability in prior testing, there was enough confidence in its success to forgo the
81 need for a second person to verify whether the algorithm may have missed any unfound matches.

82 *Test 4 Methods – Poor Quality Images*

83 During the initial image compilation process, photographs of individual Risso's dolphins
84 from each sighting were imported into a database and assigned a set of image quality and
85 distinctiveness metrics, as described previously. These metrics were used to both determine the
86 suitability of each image for inclusion in the catalog and assist in eye-based matching by
87 categorizing the animals. Images were often flagged as poor quality in these early stages due to
88 factors such as suboptimal angle, visibility, exposure or sharpness, and did not move forward in
89 the matching process from that point.

90 For this test, we selected 145 images that had been previously flagged as too "poor
91 quality" for eye-based matching in the compiling stage, with an overall quality score of 2 or

92 higher (range 2-2.5). These images were compared to the historical catalog through the multi-
93 species matching algorithm to evaluate whether the algorithm was able to identify matches that
94 human matchers had previously deemed unsuitable for catalog inclusion. The MAP score was
95 calculated for both found matches and using the dummy “new individual” prediction calculation,
96 and linear regression and ANOVA analysis were conducted to assess whether image quality or
97 distinctiveness metrics may have influenced the algorithm similarity score.

98 *Match Type Assessment Methods*

99 To evaluate comparative matching performance, all confirmed matches from the multiple
100 datasets were grouped together into three categories: matches identified by both the algorithm
101 and the human (Y/Y), matches identified by the algorithm but missed by the human (Y/N), and
102 images not identified by either the algorithm or the human (N/N). These categories were used to
103 assess whether the algorithm performed comparably to, better than, or worse than human
104 matchers. To further explore potential performance patterns, image quality and distinctiveness
105 metrics were proportionally compared across the three match types to identify specific attributes
106 that may influence algorithmic versus human matching success.

107