# NORPPA: NOvel Ringed seal re-identification by Pelage Pattern Aggregation

Ekaterina Nepovinnykh[1*], Ilia Chelak[2], Tuomas Eerola[1] and Heikki Kälviäinen[1]

[1*]Computer Vision and Pattern Recognition Laboratory, Department of Computational Engineering, School of Engineering Science, Lappeenranta-Lahti University of Technology LUT, P.O.Box 20, Lappeenranta, 53851, Finland.
[2]Department of Computer Science, Faculty of Science, University of Helsinki, P.O. Box 4, Helsinki, 00100, Finland.

**Abstract**

We propose a method for Saimaa ringed seal (*Pusa hispida saimensis*) re-identification. Access to large image volumes through camera trapping and crowdsourcing provides novel possibilities for animal monitoring and conservation and calls for automatic methods for analysis, in particular, when re-identifying individual animals from the images. The proposed method NOvel Ringed seal re-identification by Pelage Pattern Aggregation (NORPPA) utilizes the permanent and unique pelage pattern of Saimaa ringed seals and content-based image retrieval techniques. First, the query image is preprocessed, and each seal instance is segmented. Next, the seal's pelage pattern is extracted using a U-net encoder-decoder based method. Then, CNN-based affine invariant features are embedded and aggregated into Fisher Vectors. Finally, the cosine distance between the Fisher Vectors is used to find the best match from a database of known individuals. We perform extensive experiments of various modifications of the method on a new challenging Saimaa ringed seals re-identification dataset. The proposed method is shown to produce the best re-identification accuracy on our dataset in comparisons with alternative approaches.

## 1 Introduction

Animal biometrics, especially image-based individual re-identification, has recently gained extensive attention due to the availability of large volumes of wildlife image data gathered via automatic game cameras and citizen science projects. The benefits of automated re-identification methods are evident as they allow valuable data for conservation efforts to be obtained, for example, accurate population size estimates and novel information about animal migration and behavior patterns (Araujo et al., 2020; McCoy et al., 2018). Compared to traditional methods such as tagging, which may cause stress and change the behavior of the animal, image-based re-identification offers a non-invasive technique for monitoring of endangered species (Norouzzadeh et al., 2018).

The Saimaa ringed seal (*Pusa hispida saimensis*) is an endangered species native to Lake Saimaa, Finland. Seals of this species have a distinct ring-like pelage pattern, which is both permanent and unique for each individual, providing a basis for re-identification. An ongoing

conservation effort (Koivuniemi, Auttila, Niemi, Levänen, & Kunnasranta, 2016; Koivuniemi, Kurkilahti, Niemi, Auttila, & Kunnasranta, 2019; Kunnasranta et al., 2021) uses image-based re-identification to study animal migration and behavior. Currently, however, this re-identification work is carried out manually, which in view of the large number of images is very labour intensive and time consuming. Automated computer vision-based re-identification would clearly be of great benefit when carrying out this task.

A variety of methods for animal re-identification exist that utilize distinct characteristics in fur, feather and skin patterns (T. Berger-Wolf et al., 2015; Crall, Stewart, Berger-Wolf, Rubenstein, & Sundaresan, 2013; Li, Li, Tang, Qian, & Lin, 2020; Moskvyak, Maire, Dayoub, Armstrong, & Baktashmotlagh, 2021), and methods originally developed for human face re-identification have been successfully applied to animals (Agarwal et al., 2019; Crouse et al., 2017; Deb et al., 2018). Visual animal re-identification can be formulated as a task of finding a match for the given query image from a database of known individuals, which is equivalent to a content-based image retrieval (CBIR) problem (Smeulders, Worring, Santini, Gupta, & Jain, 2000) where an image is searched from a database based on the image content. However, despite the clear similarity between CBIR and re-identification tasks, utilizing utilization of CBIR approaches for animal re-identification has remained largely unstudied.

Saimaa ringed seals introduce additional challenges to the re-identification that make the task more difficult compared to many other animals for which re-identification has already been successfully applied. First, the image data is extremely biased. The majority of images are collected using static game cameras producing images with the same viewing angle and background and a limited set of possible seal locations and poses in the frame. At the same time, the high site fidelity (a tendency to return to previously visited locations) and low sociality of Saimaa ringed seals often result in a large portion of images of one individual seal being captured by only one game camera. Machine learning models trained on this kind of data tend to learn features that do not

generalize to new datasets (e.g., data from a different year with different game camera locations). Moreover, as only a small portion of Saimaa ringed seal habitat can be covered with game cameras, datasets for seal identification are usually complemented with DSLR camera images, as well as images obtained via citizen science projects (e.g., mobile phone camera pictures). This image heterogeneity introduces a domain shift and due to the fact that different individuals are often captured with different cameras, it also contributes to the database bias problem. Finally, re-identifying Saimaa ringed seals from images is very challenging per se because of: (i) the large variation in possible poses, which is further exacerbated by the deformable nature of the seals, (ii) the non-uniform pelage patterns, limiting the size of the regions that can be used for the re-identification task, and (iii) the low contrast between the ring pattern and the rest of the pelage, as well as the varying appearance (e.g., wet and dry fur). Re-identification of Saimaa ringed seals is therefore considerably more difficult than, for example, zebra re-identification, where there are clearly visible patterns and limited variation in the pose of the torso.

In this paper, we address the above challenges by proposing the NOvel Ringed seal re-identification by Pelage Pattern Aggregation (NORPPA) method for automatic Saimaa ringed seal re-identification (Fig. 1). The method is inspired by CBIR methods and builds on earlier work (Nepovinnykh, Eerola, & Kälviäinen, 2020) where Siamese networks were utilized to learn a similarity metric for local patches of pelage patterns. We further develop this approach by proposing an improved pattern feature embedding, which is done by utilizing affine invariant local CNN features and aggregating them into a fixed size embedding vector describing global features. The input image is first preprocessed using tone mapping and then segmented to detect and separate the seals from the background. The pelage pattern is further extracted using a U-net encoder-decoder (Ronneberger, Fischer, & Brox, 2015) based method. Affine invariant features are extracted and aggregated into a descriptor. Finally, the re-identification is performed by finding a descriptor with the minimum distance from the database of known individuals.
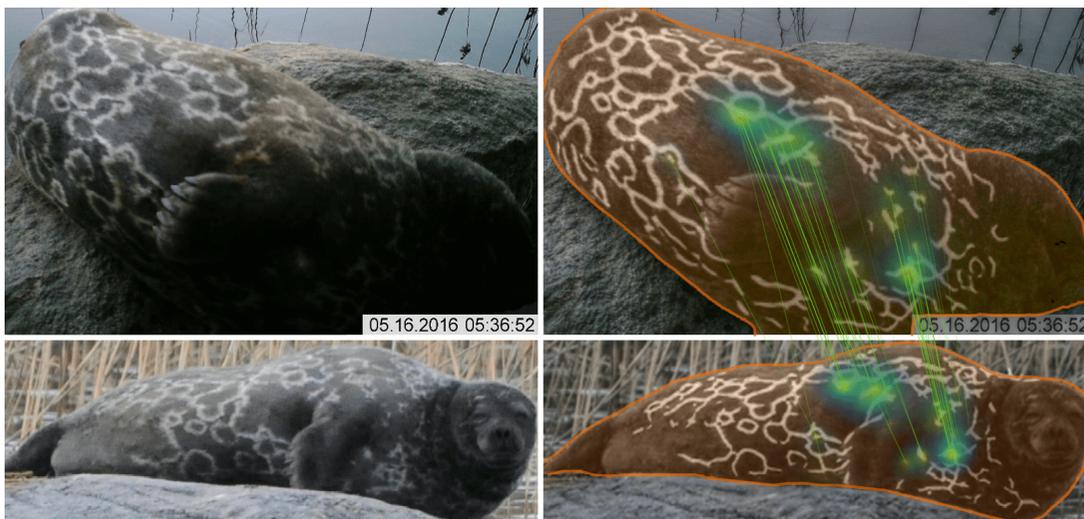
**Fig. 1**: Visualisation of the proposed re-identification method. Input pictures are on the left, the results are on the right. The seal is segmented (orange outline), and matching regions of the pelage pattern are highlighted and connected with lines. The intensity of the highlights corresponds to the similarity of the matched regions.

In the experimental part of the work, we show that the proposed method outperforms previously developed re-identification methods for Saimaa ringed seals as well as HotSpotter (Crall et al., 2013), a popular species agnostic pattern-based re-identification approach, on the challenging task of Saimaa ringed seal re-identification. In addition, different variations of the method are comprehensively evaluated to find the best pattern feature embeddings for the task. The main contribution of this paper can be summarized as follows: (i) a novel Saimaa ringed seal re-identification method (NORPPA) inspired by content based image retrieval methods, (ii) a novel combination of local affine-covariant region learning and CNN-based descriptors and feature aggregation to obtain a single fixed size pattern embedding vector with high discrimination power, and (iii) extensive evaluation of the method and its modifications on a challenging Saimaa ringed seal dataset. While the method was developed for Saimaa ringed seals, it is also possible to apply it to other patterned species as shown by Badreldeen Bdawy Mohamed (2021).

## 2 Related work

### 2.1 Animal re-identification

Animal re-identification is a broad term referring to the process of identifying an individual animal based on its features. The features are based on biological traits, and they can be captured in a number of ways, for example, acoustically (Hartwig, 2005; Pruchova, Jaška, & Linhart, 2017) or visually in the form of images (Vidal, Wolf, Rosenberg, Harris, & Mathis, 2021) or videos (Freytag et al., 2016). Currently, image-based approaches are the most widely utilized approach due to the relative ease of data acquisition and manual analysis (Schneider, Taylor, Linquist, & Kremer, 2019).

Various animal species can be re-identified by different types of visually unique biological traits such as fur pattern, face or fin shape. Examples of such traits are presented in Fig. 2. Algorithmically, the methods can be divided into classification and metric-based approaches (Vidal et al., 2021). Classification-based approaches assume that the database of known individuals is known and finite, and the final algorithm can only identify individuals from that database. Metric-based methods, on the other hand, aim to learn a similarity metric between the input images. The re-identification is
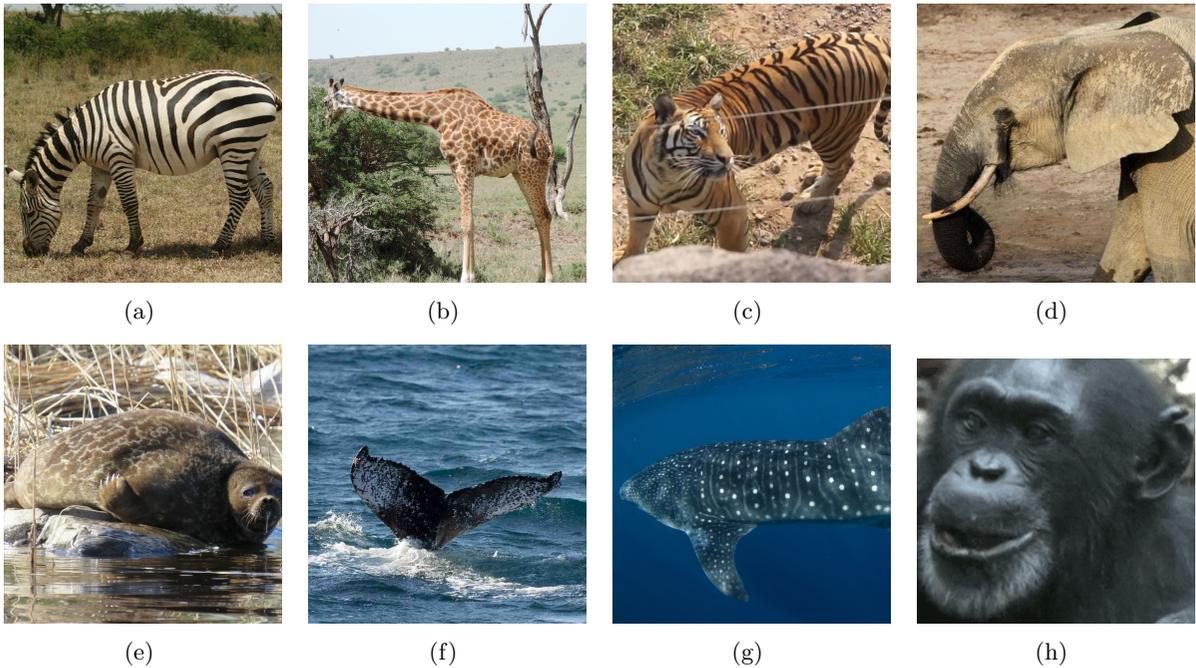
**Fig. 2**: Example images of the main identifiable features from publicly available re-identification data sets: (a) Plains zebra (*Equus quagga*) (Parham et al., 2017): stripe fur pattern; (b) Masai giraffe (*Giraffa tippelskirchi*) (Parham et al., 2017): spot fur pattern; (c) Amur tiger (*Panthera tigris*) (Li et al., 2020): stripe fur pattern; (d) African elephant (*Loxodonta africana*) (Korschens & Denzler, 2019): head shape; (e) Saimaa Ringed seal (*Pusa hispida saimensis*) (Nepovinnykh, Eerola, et al., 2022): ringed fur pattern; (f) Humpback whale (*Megaptera novaeangliae*) (Cheeseman et al., 2017): fluke shape; (g) Whale shark (*Rhincodon typus*) (Holmberg et al., 2009): skin spot pattern; (h) Chimpanzee (*Pan troglodytes*) (Freytag et al., 2016): face

then performed by clustering or matching based on the similarity, which means that metric-based approaches are not limited by the initial database and can be applied to new individuals without retraining. Re-identification algorithms also differ in the feature extraction approaches used, which can be manual or semi-manual, where user input is required to extract salient regions, or automatic, where input images are fully processed by the method. Fully automatic methods are of most interest as they would allow efficient analysis of large data volumes.

One of the largest wildlife re-identification projects, Wildbook (T.Y. Berger-Wolf et al., 2017) uses different kinds of algorithms for edge-based or pattern-based re-identification. The most efficient algorithms for deep learning edge-based re-identification are CurvRank (Weideman et al., 2017), finFindR (Thompson et al., 2019), and

OC/WDTW (Bogucki et al., 2019), which have been applied to marine mammals such as bottlenose dolphins (*Tursiops truncatus*), humpback whales (*Megaptera novaeangliae*), right whales (*Eubalaena glacialis*), and use the unique shape of tail or fins to identify the animals.

Wildbook uses PIE and Hotspotter metric-based algorithms to re-identify animals by pattern. PIE (Moskvyak, Maire, Dayoub, Armstrong, & Baktashmotlagh, 2021) is a deep learning-based method for matching of individuals invariantly to the pose. It receives shape embedding and pose embedding separately and normalizes the shape to match the individual regardless of the specific pose. PIE was originally developed for manta rays (Moskvyak, Maire, Dayoub, Armstrong, & Baktashmotlagh, 2021), but in Wildbook it is also used for humpback whale flukes, orcas, and right whales. HotSpotter (Crall et al., 2013) is a

SIFT-based (D.G. Lowe, 1999) species agnostic algorithm that uses viewpoint invariant descriptors and a scoring mechanism which emphasizes the most distinctive key points, called "hot spots," on an animal pattern. The HotSpotter algorithm has been successfully used for re-identification of zebras (*Equus quagga*) (Crall et al., 2013) and giraffes (*Giraffa tippelskirchi*) (Parham et al., 2017), jaguars (*Panthera onca*) (Crall et al., 2013) and ocelots (*Leopardus pardalis*) (Nipko, Holcombe, & Kelly, 2020).

Most recent methods for animal re-identification utilize deep learning, particularly convolutional neural networks (CNNs) (Schneider, Taylor, & Kremer, 2020; Schneider et al., 2019). CNNs have been successfully applied for re-identification of primate faces (Brust et al., 2017; Deb et al., 2018) and for pattern-based re-identification and recognition of Amur tigers (Panthera tigris)(Li et al., 2020; C. Liu, Zhang, & Guo, 2019; N. Liu, Zhao, Zhang, Cheng, & Zhu, 2019), cattle muzzle (Kumar et al., 2018), zebras (*Equus quagga*) and giraffes (*Giraffa tippelskirchi*) (Badreldeen Bdawy Mohamed, 2021). In order to improve re-identification accuracy, pose estimation and key point alignment have been proposed (Moskvyak, Maire, Dayoub, & Baktashmotlagh, 2021; Yeleshetty, Spreeuwers, & Li, 2020; Yu et al., 2021).

## 2.2 Ringed seal re-identification

A number of methods for the re-identification of Saimaa ringed seals have been developed (Chehrsimin et al., 2018; Chelak, Nepovinnykh, Eerola, Kälviäinen, & Belykh, 2021; Nepovinnykh, Eerola, Kälviäinen, & Radchenko, 2018; Zhelezniakov et al., 2015). Generally, the methods start with preprocessing steps, including seal segmentation, and then proceed to analyzing the unique pelage pattern to generate a descriptor for each individual seal. A seal segmentation method utilizing superpixel classification was proposed in (Zhelezniakov et al., 2015). The re-identification method employs common texture features extracted from the segmented seal and a Bayesian classifier. Additional color normalization and contrast enhancement steps were applied in (Chehrsimin et al., 2018) to make the pattern more visible. The actual re-identification was

performed using the Hotspotter algorithm (Crall et al., 2013).

The first attempt to utilize CNNs for Saimaa ringed seal identification was done in (Nepovinnykh et al., 2018). The individual re-identification was reformulated as a classification problem where each class corresponds to a unique individual, and transfer learning was utilized to train an individual classifier. While the performance is good on a small dataset, the method is only able to reliably perform the re-identification if there is a large set of example images for each individual. Furthermore, the whole system needs to be retrained if a new seal individual is introduced. Finally, it is unclear if the high accuracy is due to the method's ability to learn the necessary features from the pelage pattern, or if it also learns features such as pose, size, or illumination, which separate individuals in the used dataset but do not provide the means to generalize the method to other datasets.

In order to address these issues, a one-shot approach was proposed in (Nepovinnykh et al., 2020). The method starts with CNN-based segmentation of the seal. The pelage pattern is extracted utilizing a Sato tubeness filter-based method. For the re-identification, the whole pattern image is divided into patches, which are then fed into an embedding CNN. The CNN is trained using a triplet loss and essentially provides a metric that measures the visual similarity between the patches. Re-identification is then performed based on this similarity by using topology-preserving projections. The main advantage of using a triplet CNN is the ability to easily add new individuals into the database since no retraining is necessary.

The pattern embedding step is crucial for any re-identification method as distinctive but compact embedding that captures the characteristics of the pattern forms the basis for successful re-identification. The pattern embedding step was considered in more detail in (Chelak et al., 2021), where EDEN, a new pooling layer, was proposed to account for the spatial distribution of pattern features. It was shown that the proposed pooling layer increases the matching accuracy of the pattern patches.

Another version of the re-idenfitication algorithm was proposed and applied to the sister species of Saimaa ringed seals, Ladoga ringed seals (*Pusa hispida ladogensis*) in (Nepovinnykh,

Chelak, et al., 2022). Ladoga ringed seals have a similar pattern to Saimaa ringed seals, which means that the same re-identification algorithm is applicable to both species. Two new steps were introduced into the pipeline: individual grouping and Fisher Vector computation. The individual grouping step focuses on finding multiple instances of the same individual from an image sequence. This rather simple image retrieval-based method was shown to attain high accuracy in matching individuals within an image sequence producing sets of images of each seal to be re-identified. This was shown to be beneficial for the re-identification as it helps to compensate for poor image quality, which often results in an inability to extract patterns from some images, and it allows a larger portion of the pattern to be captured as the seal changes its pose between images. The Fisher Vector (Hutchison et al., 2010; Perronnin & Dance, 2007; Perronnin, Liu, Sánchez, & Poirier, 2010) is used to aggregate patch descriptors from an individual seal into a single image descriptor. The vector further allows the patch descriptors from multiple images to be aggregated, providing a straightforward tool for utilize utilization of the image sets produced by the grouping step. Aggregated image descriptors are used to find a match from the database of known individuals by calculating distances. Promising results were obtained on Ladoga ringed seal re-identification.

## 2.3 Content based image retrieval

The task of visual animal re-identification can be formulated as a task of finding the most similar image from the database to the given query image. This formulation matches the definition of content-based image retrieval (CBIR) (Smeulders et al., 2000) and motivates study of the suitability of CBIR methods for animal re-identification.

CBIR methods usually consist of two main steps: feature extraction and feature aggregation. The feature extraction problem can be solved using standard hand-crafted features, such as Scale Invariant Feature Transform (SIFT) (Arandjelović & Zisserman, 2012; D. Lowe, 2004), or extraction by convolutional neural networks (see e.g. (Mishchuk, Mishkin, Radenovic, & Matas, 2017)). Then, feature aggregation creates a descriptor for each image that can be used to find the most similar image

from the database. Traditional methods such as Bag of Words (BOW) (Sivic & Zisserman, 2003), Vector of Locally Aggregated Descriptors (VLAD) (Jégou, Douze, Schmid, & Pérez, 2010) and the Fisher Vector (Hutchison et al., 2010; Perronnin & Dance, 2007; Perronnin et al., 2010) do the aggregation using a specially constructed vocabulary. The vocabulary is usually created by an unsupervised clustering algorithm. For example, k-means (MacQueen et al., 1967) is used for VLAD and a Gaussian Mixture Model (GMM) (McLachlan & Basford, 1988) is used for the Fisher Vector. Finally, fixed-size descriptors are created for each image based on the vocabulary and extracted features. The distance between these descriptors is inversely proportional to the visual similarity.

Due to the availability of data and the convenience of end-to-end approaches, deep learning-based methods for CBIR are becoming increasingly more popular. The advantage of CNN-based methods is that the two main steps, feature extraction and feature aggregation, are naturally implemented as a part of the network architecture, with the first part of the network being a feature extractor and a final specialized layer doing the feature aggregation. For example, there have been several attempts to create deep analogues of traditional methods such as NetVLAD (Arandjelovic, Gronat, Torii, Pajdla, & Sivic, 2016) where a generalized VLAD layer is used to aggregate CNN-extracted features.

In (Babenko, Slesarev, Chigorin, & Lempitsky, 2014; Gong, Wang, Guo, & Lazebnik, 2014), fully-connected layers are used to generate the final descriptor, which is a standard approach for CNNs. In (Razavian, Sullivan, Carlsson, & Maki, 2016), a global max pooling approach is introduced which produces the final descriptor from the activation maps by taking a maximum value from each filter activation, resulting in a descriptor of the same size as the number of filters. Different variants of global pooling operations have also been studied. These include integral max-pooling (Tolias, Sicre, & Jégou, 2016), sum pooling (Babenko & Lempitsky, 2015) and generalized mean pooling (Radenović, Tolias, & Chum, 2018). Integral max-pooling (Tolias et al., 2016) is particularly interesting since it creates the final

descriptor by applying max-pooling to the overlapping image regions, which also allows spatial information to be encoded.

# 3 Method

The proposed NORPPA method consists of six steps: 1) image prepossessing, 2) seal instance segmentation, 3) pelage pattern extraction, 4) feature extraction, 5) feature aggregation and 6) individual re-identification (see Fig. 3).

## 3.1 Image preprocessing

Depending on illumination conditions variation in the contrast of the images can be rather high. This could lead to a loss of detail in the region of interest, i.e. the seal and its pelage pattern. In order to rectify this issue, we employ the tone-mapping approach to equalize the contrast in dark and bright image regions. The algorithm proposed by (Mantiuk, Myszkowski, & Seidel, 2006) is used due to its ability to produce realistic tone-mapped images without introducing visual artifacts. This method considers contrast on multiple spatial frequencies while using gradient methods with some additional extensions to ensure that the global brightness levels are not reversed and low-frequency details are properly reconstructed. Examples of images before and after prepossessing are presented in Fig. 4.

## 3.2 Seal instance segmentation

Seal instance segmentation step is important since most of the images are obtained using static camera traps. This together with the fact that seal individuals tend to use same sites or areas inter-annually cause one seal individual to be very often captured with the same camera (same background). This increases the risk that the supervised identification algorithm learns to identify the background instead of the actual seal if the full image or the bounding box around the seal is used. Consequently, algorithm behavior may result in a system that is unable to identify the seal in a new environment.

Instance segmentation is performed using Mask R-CNN (He, Gkioxari, Dollár, & Girshick, 2017). A segmentation model trained for Ladoga ringed seals from (Nepovinnykh, Chelak, et al.,

2022) is utilised. This is possible due to the two species being visually almost indistinguishable. Ladoga ringed seals are more numerous than Saimaa ringed seals and they are often captures in large groups which makes it easier to collect and annotate large training data for the segmentation. For more details about the instance segmentation model and training procedure see (Nepovinnykh, Chelak, et al., 2022).
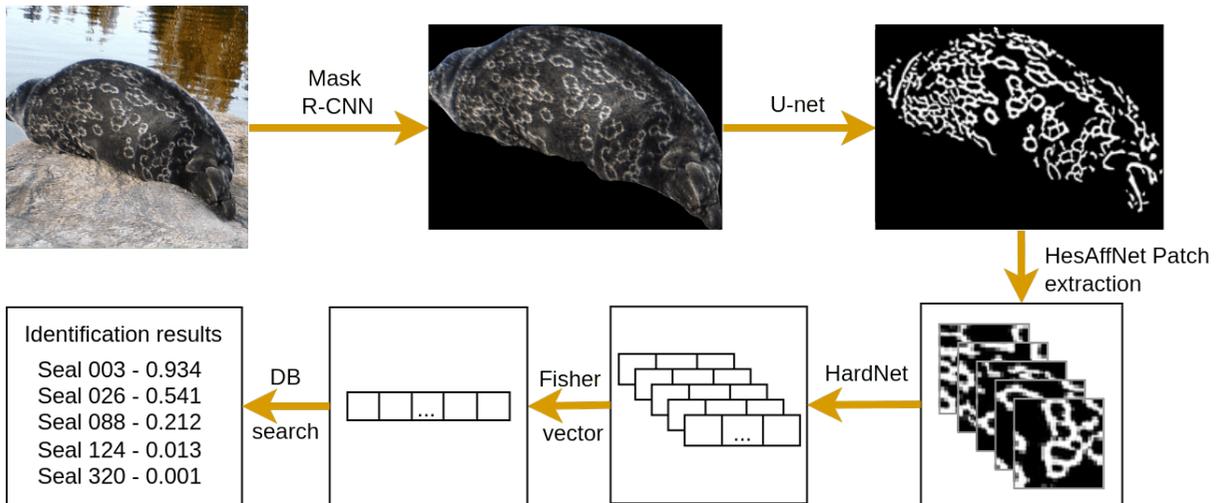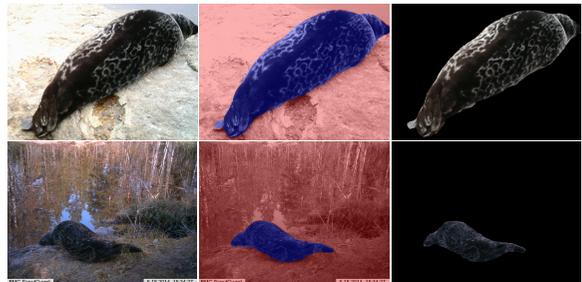
After the segmentation masks are obtained, additional morphological operations are applied to close the holes and smooth the borders by using morphological closing and opening. The examples of segmentation results are presented in Fig. 5.

## 3.3 Pelage pattern extraction

The main distinguishing feature of a seal is its pelage pattern, which is both permanent and unique to each seal allowing the identification of individuals over their whole lifetime. The pelage pattern forms the basis for the proposed re-identification method. In order to focus the attention on the pattern and discard irrelevant information causing database bias such as illumination and other visual factors (e.g., wet fur looks different from the dry fur), the pattern is segmented. This is done using CNN based method utilizing U-net encoder-decoder architecture (Ronneberger et al., 2015). The output of the method is a binarized image of the pelage pattern (see Fig. 6). The pattern image is further post-processed to remove small noise by using unsharp masking and morphological opening. All images are then resized in such way that the mean width of the pattern lines is the same for all images, bringing them into the same scale. This is necessary because the images are obtained from various sources and the image resolution has a large variation. For more detailed explanation of the pattern extraction step, as well as the comparison to other methods, see (Zavialkin, 2020).
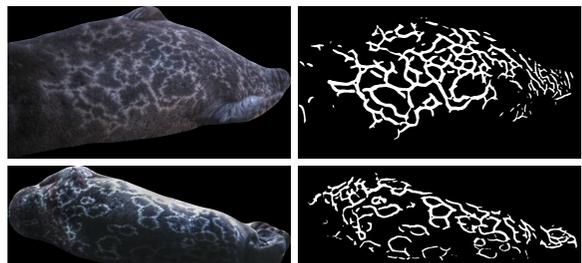
## 3.4 Feature extraction

Seals can be found in a variety of poses. The deformable nature of seals body results in distorted and warped patterns on images. While the pattern as a whole is transformed in a non-linear way, it can be argued that small local regions experience close to affine transformations,

**Fig. 3**: NORPPA re-identification pipeline.



**Fig. 4**: Examples of the image processing of camera trap images. Images on the left are the originals. The right column demonstrates the result of the tone-mapping.



**Fig. 5**: Examples of the segmentation masks. Images on the left are the originals. The mask is highlighted in blue and the background is highlighted in red on the middle images. The last column shows the result of the segmentation.



**Fig. 6**: Example of pattern extraction output.

making an affine invariant feature extractor suitable for the task. For this purpose a combintaion HesAffNet (Mishkin, Radenović, & Matas, 2018) detector and HardNet (Mishchuk et al., 2017) descriptor is used.

The combination of a Hessian-Affine detector (Mikolajczyk & Schmid, 2004) with Root-SIFT (Arandjelović & Zisserman, 2012) used to be considered a gold standard for local feature extraction and description. However, with the increasing size of available datasets and rapidly developing field of deep learning, CNN-based methods are

now able to outperform previous handcrafted features. The combination of HesAffNet (Mishkin et al., 2018) and HardNet (Mishchuk et al., 2017) is able to provide state-of-the-art results in image

retrieval tasks, which makes those methods particularly useful for animal re-identification as well.

HesAffNet is a modification of the classical Hessian Affine Region detector (Mikolajczyk & Schmid, 2002, 2004), where the shape estimation step is done by the AffNet CNN. The detector is based on the Harris cornerness measure (Harris & Stephens, 1988), which uses a second moments matrix to find regions of interest by estimating the most prominent gradient directions. This method is combined with the multiscale approach from (Lindeberg, 1998) which uses Laplacian of Gaussian to find extrema in the scale space. The same concept can be further extended to all affine transformations, not just the scale. However, the degree of freedom is much higher for affine transformations, which complicates the process and requires a special shape adaptation algorithm. The original Hessian Affine detector used Baumberg iteration (Baumberg, 2000), which is replaced by an AffNet CNN in HesAffNet.

AffNet and HardNet are closely related, sharing the architecture and similar training procedure. During the training of HardNet, batches of matching patch pairs are chosen, each containing an anchor $a_i$ and positive match $p_i$. Each pair correspond to a different location, i.e. there are no other matches except for the ones in each pair. Each patch is encoded by the network, and a matrix of pair-wise distances between all anchors and positive matches are computed. For each pair, a closest non-matching descriptor from the batch is chosen, and a final hard negative margin loss is computed as

$$L = \frac{1}{n} \sum_{i=1}^{n} \max(0, 1 + d(a_i, p_i)) \\ - \min(d(a_i, p_{j\,\min}), d(a_{j\,\min}, p_i)), \tag{1}$$

where $p_{j\,\min}$ is the closest non-matching positive to $a_i$, and $a_{j\,\min}$ is the closest non-matching anchor to $p_i$.

AffNet utilizes a slightly different training procedure, the main difference being that the derivative for the negative term in the loss is set to 0. This loss is called hard negative-constant and helps avoid the situations where positive samples cannot be moved closer together because of a negative sample lying between them in the metric space. The training procedure for AffNet is also

more complicated, since it is learning affine shapes and not just a distance metric. Therefore, spatial transformers are used to transform input patches according to the predicted shape, which are then fed into a descriptor network, e.g. HardNet, and only then is the loss calculated and backpropagated through both networks. The example of HesaffNet application to a preprocessed image is visualised in Figure 7.
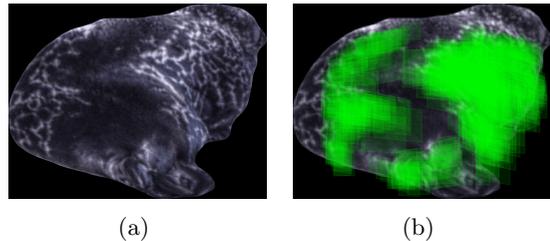


**Fig. 7**: Visualisation of Hessian Affine patch extraction: (a) segmented image; (b) HesAffNet-based patch extraction. Extracted regions are highlighted in green.

## 3.5 Feature aggregation

Features are aggregated using Fisher Vector (Hutchison et al., 2010; Perronnin & Dance, 2007; Perronnin et al., 2010). First, Principal Component Analysis (PCA) is applied to the resulting the feature embeddings to decorrelate the features and reduce the dimensionality. This is an important for Fisher Vectors, which are known to produce large descriptors. The images in the database of known individuals are used to learn principal components. Next, a visual vocabulary is constructed by applying Gaussian Mixture Model (GMM) to the features from the database. Then, Fisher Vectors are created for each image by computing the partial derivatives of the log-likelihood function with respect to the GMM parameters and concatenating them. Kernel PCA (Schölkopf, Smola, & Müller, 1998) is applied to further reduce the dimensionality of the resulting image descriptors which helps to reduce the storage requirements for the database, as well as speed up the database search for the re-identification.

### 3.5.1 Fisher Vector

Let $X = \{x_t, t = 1, \ldots, T\}$ be a sample of $T$ observations and $u_\lambda$ be a probability density function modelling the distribution of the data, where $\lambda$ is a vector of its parameters. The score is defined as the gradient of the log-likelihood of the data on the model:

$$G_\lambda^X = \nabla_\lambda \log u_\lambda(X). \qquad (2)$$

This score function can be used to define the Fisher Information Matrix (FIM) (Amari & Nagaoka, 2000):

$$F_\lambda = E_{x \sim u_\lambda}[G_\lambda^X G_\lambda^{X'}], \qquad (3)$$

which acts as a local metric for a parametric family of distributions. This metric can also be used to measure the similarity between 2 samples using the Fisher Kernel (FK) (Jaakkola & Haussler, 1999):

$$
\begin{aligned}
K_{FK}(X,Y) &= G_\lambda^{X'} F_\lambda^{-1} G_\lambda^Y \\
&= G_\lambda^{X'} L_\lambda' L_\lambda G_\lambda^Y \qquad (4) \\
&= \mathscr{G}_\lambda^{X'} \mathscr{G}_\lambda^Y,
\end{aligned}
$$

where $L_\lambda' L_\lambda$ is the Cholesky decomposition of $F_\lambda^{-1}$, $G_\lambda^X$ and $G_\lambda^Y$ are the Fisher Vectors of samples $X$ and $Y$ respectively. By using Fisher Vectors, it is possible to calculate the kernel as simple dot product, which can efficiently be utilized by linear classifiers. When constructing a Fisher Vector for an image, a set of local features is assumed to be independent, meaning that the final descriptor can be constructed as a sum of Fisher Vectors for each local feature, i.e.

$$G_\lambda^X = \sum_{t=1}^{T} L_\lambda \nabla_\lambda \log u_\lambda(X). \qquad (5)$$

Usually, Gaussian Mixture Model (GMM) is used as $u_\lambda$, since it can be used to approximate any continuous distribution with arbitrary precision (Titterington, Afm, Smith, Makov, et al., 1985). Then, the vector of parameters $\lambda$ contains mixture weights $w_k$, mean vectors $\mu_k$ and covariance matrices $\Sigma_k$ for each Gaussian $u_k, k = 1, \ldots, K$. Using the assumption that the assignment of each feature to mixture components is almost hard, i.e.

each feature is assigned to only one cluster, it could be inferred (Sánchez, Perronnin, Mensink, & Verbeek, 2013) that the FIM is diagonal, which means that $L_\lambda$ is just a coordinate-wise normalization of the gradient vectors. The final normalized gradients are then defined as follows

$$\mathscr{G}_{\mu_k}^X = \frac{1}{\sqrt{w_k}} \sum_{t=1}^{T} \gamma_t(k) (\frac{x_t - \mu_k}{\sigma_k}), \qquad (6)$$

$$\mathscr{G}_{\sigma_k}^X = \frac{1}{\sqrt{w_k}} \sum_{t=1}^{T} \gamma_t(k) \frac{1}{\sqrt{2}} \left[ \frac{(x_t - \mu_k)^2}{\sigma_k^2} - 1 \right], \quad (7)$$

where $\gamma_t$ is the soft assignment function

$$\gamma_t(k) = \frac{w_k u_k(x_t)}{\sum_{j=1}^{K} w_j u_j(x_t)}. \qquad (8)$$

It should be noted that the gradients for the weight parameters $w_k$ are usually omitted, since they do not provide much additional information (Hutchison et al., 2010). Those gradients are concatenated into a vector of size $2DK$, where $D$ is the dimensionality of samples and $K$ is the number of components in GMM. It has been shown (Hutchison et al., 2010) that $L2$ and Power normalization generally improve the performance of the method. Therefore, it is common to apply Power and $L2$ normalization to the Fisher Vector to get the final descriptor.

## 3.6 Individual re-identification

Re-identification is done by calculating the cosine distance from the query image descriptor to each image descriptor in the database of known individuals and selecting the individual ID with the lowest distance. To visualize the re-identification and to provide semi-automatic tool for experts, heatmaps highlighting the similar areas in patterns of the query image and database images are computed. This is done using the following method. First, features from a query are paired with the closest database features. Then, pairs with distance larger than 10th percentile of distances are discarded. The remaining pairs are used to find the homography using Direct Linear Transform (DLT) (Hartley & Zisserman, 2004) and Random Sample Consensus (RANSAC) (Fischler & Bolles, 1981). The inliers of the final homography are highlighted with ellipses aligned

and transformed according to the extracted affine regions. The intensity of each ellipse is inversely proportional to the distance between the local features in the corresponding pair, i.e. directly proportional to their similarity.

# 4 Experiments and results

## 4.1 Data

The dataset consists of 57 individual seals with a total of 2080 images. The dataset is divided into two subsets: database and query. The database subset contains a minimal number of high-quality unique images that are enough to cover the full body pattern of each seal. The query subset contains the remaining images and contains the same individuals as in the database. It should be noted that the high-quality images were prioritized when constructing the database and, therefore, images in the query subset often have lower quality. Examples of images from both subsets are presented in Fig. 8. The dataset has been made publicly available. For further description of the dataset, see (Nepovinnykh, Eerola, et al., 2022).
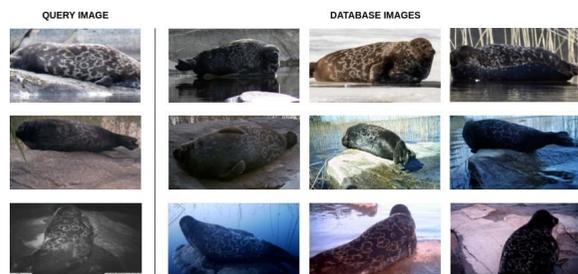


**Fig. 8**: Examples from the database and query datasets. Every row contains images of an individual seal. For every image from the query dataset (left) there is a corresponding subset of images from the database (right).

To train and evaluate the patch embedding (feature extraction) and matching (finding the corresponding patch in other images) a separate dataset of pattern image patches (see Fig.9) was constructed (Chelak et al., 2021). The dataset contains, in total, 4599 images (patches of the size $256 \times 256$ pixels). The data is divided into training and testing subsets. The training subset contains 3016 images and 16 classes. The testing subset
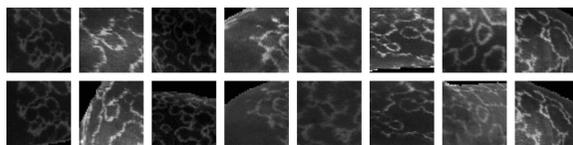


**Fig. 9**: Examples of pattern image patches. The patches in the second row match the patches in the first row.

contains 1583 images and 26 classes that are different from the training classes in the training set. Each class corresponds to one manually selected location in the pelage pattern of one individual seal. Each sample from one class was extracted from different images of the same seal. For estimation of the accuracy of the method, the testing set was divided into the database and query subset with a ratio of 1 to 2. The images that were used to construct the dataset of pattern image patches are not included in the database and query subsets of the re-identification dataset.

## 4.2 Feature extraction

The feature extraction step contains two differences compared to the previous version of the Saimaa ringed seal re-identification algorithm (Nepovinnykh et al., 2020). The first difference is that the region of interest detection approach uses the affine invariant regions (HesAffNet) instead of dense patches. The second difference is a switch to HardNet network to compute patch embedding. To assess the necessity of each of these changes both modifications were evaluated separately. Hyperparameters for all versions of the algorithm were chosen using the Tree Parzen Estimator (Bergstra, Bardenet, Bengio, & Kégl, 2011) algorithm. The results of the experiments are presented in Table 1.

As can be seen, both HesAffNet for region of interest detection and HardNet for patch embedding computation improve the accuracy noticeably. This finding leads to the conclusion that the dense patches approach cannot handle more general cases, whereas fine invariant features provide much needed robustness to various imaging conditions.

In order to evaluate the effect of the pelage pattern extraction on the algorithm's accuracy, an ablation study has been performed. The results with and without the pattern extraction step are

**Table 1**: Re-identification accuracy for different variants of the algorithm.

| Patch extraction | Patch embedding | Top-1 | Top-2 | Top-3 | Top-4 | Top-5 |
|---|---|---|---|---|---|---|
| Dense | Triplet | 52.06% | 56.91% | 60.36% | 63.58% | 65.70% |
| | ArcFace | 39.94% | 45.15% | 50.06% | 53.58% | 56.67% |
| | HardNet | 52.18% | 57.88% | 61.70% | 64.61% | 67.27% |
| HessAffNet | Triplet | 60.42% | 65.03% | 69.27% | 71.64% | 73.52% |
| | ArcFace | 47.03% | 51.64% | 55.58% | 58.36% | 60.55% |
| | HardNet | **77.64%** | **80.91%** | **82.97%** | **84.18%** | **85.27%** |

presented in Table 2. It is clear that the pelage feature extraction significantly increases the accuracy of the algorithm.

### 4.3 Patch embedding network

The following experiments were conducted in order to further improve the method:

1. Training and fine-tuning of HardNet on different datasets,
2. Various architecture modifications to the HardNet model.

#### 4.3.1 Training and fine-tuning

The original HardNet was trained on the union of HPatches (Balntas, Lenc, Vedaldi, & Mikolajczyk, 2017) and Brown (Brown & Lowe, 2007) datasets. Typically, fine-tuning a machine learning model on domain-specific training data improves the method performance in a new domain. To test this on Saimaa ringed seal re-identification, we fine-tuned the HardNet model on patches of pelage pattern images. Fine-tuned models were compared to the pretrained model, a model trained from scratch on the pattern patches, and a model trained on the union of all datasets.

The results are presented in Table 3. For the training, all hyperparameters and random seeds were taken from the original implementation of HardNet (Mishchuk et al., 2017).

While fine-tuning on the patches dataset improved the accuracy of the patch matching, the overall accuracy of the full-image matching dropped significantly. One possible reason is that the patches dataset was created using patches of the same scale, while the patches extracted by

HesAffNet during the full re-identification algorithm vary in scale, leading to a different level of detail.

Training on the union of all datasets showed no considerable improvements. This result can be explained by the size of the pelage pattern patches dataset in comparison to the combined sizes of the Brown and HPatches datasets. In other words, since HardNet utilizes triplet sampling during the training stage, the probability of an image from the pelage pattern dataset appearing in the triplet is extremely small.

#### 4.3.2 Architecture modifications

Several further modifications to the HardNet architecture were also considered. First, a Self-Organized Operational Neural Network (Self-ONN) (Malik, Kiranyaz, & Gabbouj, 2021) was incorporated into the HardNet model. Self-ONNs are networks consisting of layers that are the generalizations of convolutional layers. Simply put, each value in a convolutional kernel can be seen as a linear function, and this function can be generalized through Taylor series approximation with coefficients learned by the network. Such an approach leads to great nonlinearity even with shallow networks. Other modifications include the use of an EDEN pooling layer (Chelak et al., 2021), as well as changes to the number of channels and the output vector size.

The following models were evaluated:

- **HardNetONN**. This model has the same architecture as HardNet in terms of layers and number of channels in each layer. The only difference is that each convolution layer is replaced by a self ONN layer with a Taylor series degree

**Table 2**: Comparison of re-identification results by the NORPPA method on the SealID dataset with and without the pattern extraction step.

| Input data | TOP-1 | TOP-5 | TOP-10 | TOP-20 |
|---|---|---|---|---|
| Original images | 55.03% | 68.48% | 76.36% | 84.73% |
| Pattern images | 77.64% | 85.27% | 89.09% | 92.18% |

**Table 3**: Comparison of results for HardNet trained and fine-tuned on various datasets. We report mean with standard deviation.

| Training | Fine-tuning | Patches top-1 | Full top-1 | Full top-5 |
|---|---|---|---|---|
| Pattern patches | - | $86.44 \pm 0.41\%$ | $59.86 \pm 1.36\%$ | $71.39 \pm 1.31\%$ |
| Brown+HPatches | - | $93.02 \pm 0.51\%$ | $\mathbf{77.19 \pm 0.93}\%$ | $\mathbf{85.11 \pm 0.77}\%$ |
| Brown+HPatches | Pattern patches | $\mathbf{93.76 \pm 0.12}\%$ | $70.69 \pm 0.41\%$ | $80.46 \pm 0.42\%$ |
| Brown+HPatches+ Pattern patches | - | $92.48 \pm 0.94\%$ | $76.99 \pm 1.19\%$ | $84.85 \pm 1.03\%$ |

equal to 3 for all layers, which leads to three times as many parameters.

- **HardNetONNDrop**. This model has the same architecture as HardNetONN but the last layer has a dropout with a probability of 0.3 similarly to the original HardNet.
- **HardNetONN + EDEN**. This model has the same architecture as HardNetONN, albeit that the last convolutional layer kernel is downsampled from $8 \times 8$ to $3 \times 3$ so that it would be possible to apply pooling. After the pooling a vector of size 128 is fed into a fully connected layer with the output size of 128, resulting in a compact embedding.
- **HardNetONNSmall**. This model has the same number of parameters as HardNet. All the layers were shrunk by half and the Taylor series degree was set to 4 for all layers. Consequently, the final vector has a size of 64 instead of 128.
- **HardNet3_384**. This model is an original HardNet with 3 times as many channels in all of the layers. Therefore, it has an output vector of size 384 instead of 128.
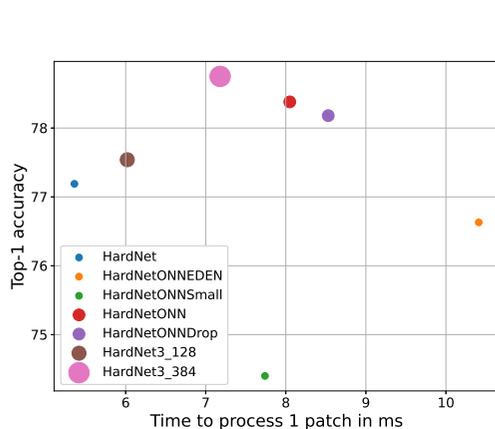- **HardNet3_128**. This model is the same as HardNet3_384 but with an output vector size of 128.

A comparison of the models is presented in Table 4.

The HardNetONN and HardNet3_384 models show higher accuracy on both patch matching and full re-identification tasks than other versions of HardNet. Moreover, although HardNet3_384 has 12 million parameters and a vector size of 384, the difference of scores with HardNetONN is small, with the TOP-5 full re-identification score difference being negligible. A comparison of the processing speed of the models is presented in Fig. 10. Overall, the improvements over the baseline HardNet are rather small while result in a noticeable increase in computer time, limiting their usability in practice.

The accuracy of HardNetONNSmall is worse compared to HardNet, although it has the same number of parameters. This can be explained by the fact that the embedding vector is cut in half for HardNetONNSmall and may not contain enough information to learn a good metric. Additionally, HardNetONN + EDEN also scored lower than the original HardNet, although higher than HardNetONNSmall. The reason could lie in the redundancy of inductive bias provided by the pooling, as well as the worse convergence of the model.

**Table 4**: Comparison of the proposed models.

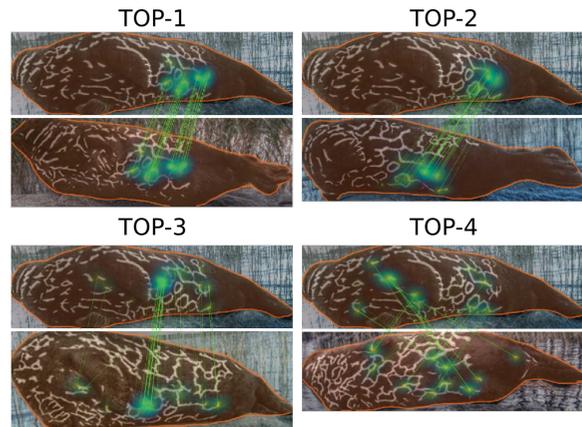| Method | Parameters (M) | Patches TOP-1 | Full TOP-1 | Full TOP-5 |
|---|---|---|---|---|
| HardNet | 1.3 | $93.02 \pm 0.51\%$ | $77.19 \pm 0.93\%$ | $85.11 \pm 0.77\%$ |
| HardNetONN + EDEN | 1.3 | $91.73 \pm 0.48\%$ | $76.63 \pm 0.81\%$ | $84.5 \pm 0.47\%$ |
| HardNetONNSmall | 1.3 | $90.48 \pm 0.68\%$ | $74.40 \pm 1.03\%$ | $82.06 \pm 0.91\%$ |
| HardNetONN | 4 | $93.41 \pm 0.43\%$ | $78.38 \pm 0.35\%$ | $86.26 \pm 0.44\%$ |
| HardNetONNDrop | 4 | $92.55 \pm 0.31\%$ | $78.18 \pm 0.98\%$ | $85.37 \pm 0.63\%$ |
| HardNet3_128 | 5.7 | $93.50 \pm 0.54\%$ | $77.54 \pm 0.31\%$ | $86.18 \pm 0.44\%$ |
| HardNet3_384 | 12 | $\mathbf{94.30 \pm 0.56\%}$ | $\mathbf{78.75 \pm 0.33\%}$ | $\mathbf{86.28 \pm 0.44\%}$ |



**Fig. 10**: Plot of TOP-1 re-identification accuracy versus processing speed for all models. The circle size corresponds to the number of parameters in the model.



**Fig. 11**: TOP-4 examples for the NORPPA method. First line: query image. Second line: four best matches in a decreasing order of similarity from left to right. Matched hotspots are highlighted in green. TOP-1–TOP-3 matches are correct. TOP-4 is incorrect.

## 4.4 Qualitative evaluation

Visual examples of the re-identification results for the proposed NORPPA method are presented in Fig. 11. For the final version we use HardNet trained on Brown and HPatches datasets. Upon inspecting the results with highlighted areas, it is evident that the proposed method learns to perform the matching between query and database images based on the characteristics of the pelage pattern. Furthermore, it can be seen that the method is able to find the corresponding regions in the patterns in very challenging cases (Fig. 12).

## 4.5 Quantitative evaluation

SaimaaReID (Nepovinnykh et al., 2020), Ladoga-ReID (Nepovinnykh, Chelak, et al., 2022) without grouping step and NORPPA seal re-identification

methods have been compared to HotSpotter(Crall et al., 2013), which is another method developed for patterned animal re-identification. HotSpotter is species-agnostic, and as such can be applied to Saimaa ringed seals as well. The results of NORPPA and HotSpotter for the Saimaa ringed seal dataset are presented in Table 5. It can be seen that the proposed method clearly outperforms HotSpotter based on TOP-1 accuracy. The difference is even more clear on TOP-5 accuracy, implying that even when NORPPA fails to correctly re-identify the seal, it is often able to provide a high rank for the correct match in the database. This is especially useful when the method is applied in a semi-supervised manner
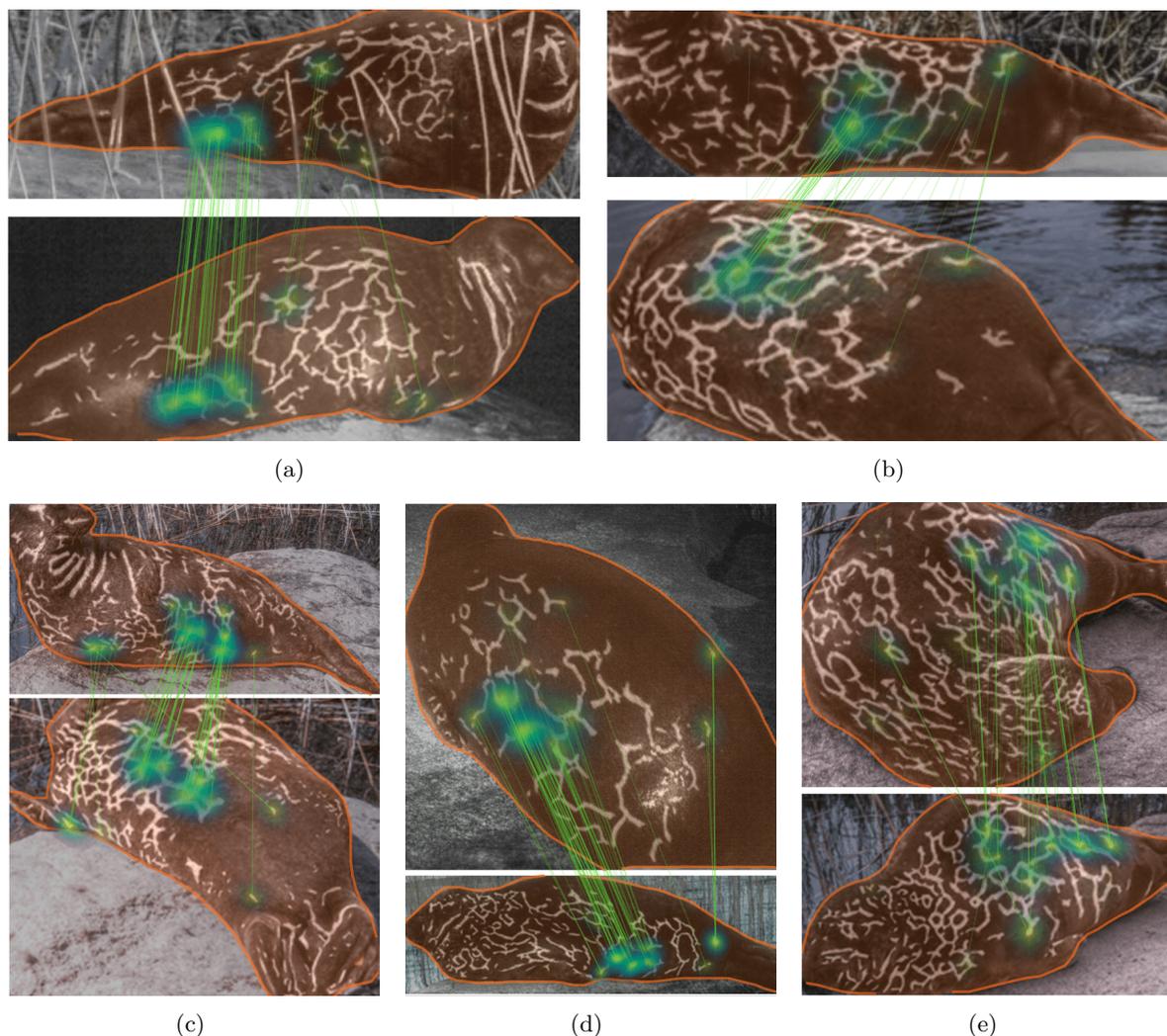
**Fig. 12**: Examples of some challenging cases. Top images are matched to the bottom images. The seal segmentation is shown in orange. The matching regions are highlighted and connected with green lines, the intensity corresponds to the similarity of a matched pair. The algorithm is able to match patterns even when the pose and point of view change significantly. (a) shows a successful match in the presence of some foreground obstacles.

**Table 5**: Comparison of re-identification results between HotSpotter, NORPPA and previous iterations of the algorithm: SaimaaReID (Nepovinnykh et al., 2020) and LadogaReID (Nepovinnykh, Chelak, et al., 2022).

| Method | TOP-1 | TOP-2 | TOP-3 | TOP-4 | TOP-5 |
|---|---|---|---|---|---|
| SaimaaReID (Nepovinnykh et al., 2020) | 35.23% | 41.45% | 44.61% | 47.92% | 60.39% |
| LadogaReID (Nepovinnykh, Chelak, et al., 2022) | 39.94% | 45.15% | 50.06% | 53.58% | 56.67% |
| HotSpotter (Crall et al., 2013) | 61.87% | 63.09% | 63.63% | 63.93% | 64.42% |
| NORPPA (ours) | **77.64%** | **80.91%** | **82.97%** | **84.18%** | **85.27%** |

where the algorithm provides a set of possible matches for the expert to verify.

By considering a larger number of top matches, it is possible to further increase the chances of finding a correct individual. The plot of the top-$k$ accuracy relative to the $k$ value is presented in Fig. 13. The relationship for the NORPPA, SaimaaReID and LadogaReID methods is logarithmic in nature with fast growth for small $k$ values, which slows down significantly with higher values. HotSpotter, on the other hand, exhibits almost no improvement after TOP-2 accuracy, with the difference between TOP-1 and TOP-5 accuracy being only about 2%, while the difference for NORPPA is almost 10%. The improvement in accuracy is a desirable property for a semi-automatic approach, offering a considerable accuracy improvement in exchange for a relatively small increase in the manual work required (as compared to a fully manual approach). Depending on the final application and available data, the relationship between the top-$k$ accuracy and $k$ can be used to determine the optimal number of matches to be returned by the algorithm.
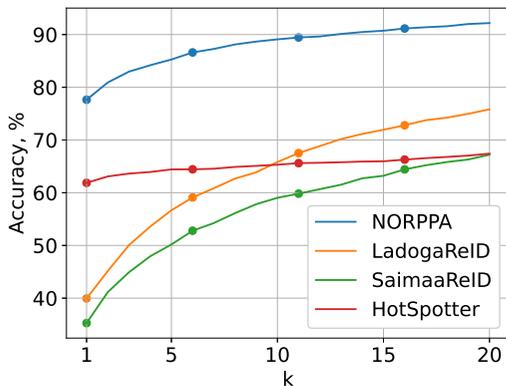


**Fig. 13**: Plot of top-$k$ re-identification accuracy for the proposed NORPPA method relative to $k$.

## 5  Conclusion

A novel method for Saimaa ringed seal re-identification called NOvel Ringed seal re-identification by Pelage Pattern Aggregation (NORPPA) was proposed in this paper. The method utilizes pelage pattern extraction and feature aggregation inspired by content-based image retrieval techniques. The re-identification pipeline consists of image enhancement, seal instance segmentation by Mask R-CNN, U-net based pelage pattern extraction, pattern feature extraction, feature aggregation, and individual re-identification by database search. Improved pattern feature embeddings were proposed by employing affine-invariant region of interest detection, CNN based feature descriptors, and Fisher Vector feature aggregation to obtain fixed size embedding vectors with high discriminative power. The proposed method was applied to a novel and challenging Saimaa ringed seal dataset and showed superior performance compared to HotSpotter and earlier versions of the Saimaa ringed seal re-identification method by the authors. One additional benefit of the proposed method is that it allows features to be aggregated over multiple images. This opens interesting possibilities for further research as sequences of game camera images can be utilized to create a single descriptor for a larger portion of a pelage pattern by filling in the gaps created by obstructions and viewpoints. While the method was developed for Saimaa ringed seals, it is also possible to apply it on other patterned animal species.

## Declarations

### Funding

## Conflict of interest

We declare no competing interests.

## Ethics approval

Data collection was done under permits by the Finnish environmental authorities ELY centre (ESAELY/1290/2015, POKELY/1232/2015, KASELY/2014/2015 and POSELY/313/07.01/2012) and Metsähallitus (MH 5813/2013 and MH 6377/2018/05.04.01).

## Consent for publication

All authors consent that the publisher has the author's permission to publish research findings. All authors guarantee that the research findings have not been previously published.

## Availability of data and materials

All data and materials are publicly available at https://doi.org/10.23729/0f4a3296-3b10 -40c8-9ad3-0cf00a5a4a53

## Code availability

The codes for the described experiments are available at https://github.com/kwadraterry/Norppa

## Authors' contributions

T. Eerola and H. Kälviäinen were responsible for the supervision of the research, designing methodology, and project administration; E.Nepovinnykh and I.Chelak implemented the algorithm. E.Nepovinnykh, I.Chelak, T.Eerola, and H. Kälviäinen prepared the original draft of the manuscript. All the authors gave the final approval for publication.

## References

Agarwal, M., Sinha, S., Singh, M., Nagpal, S., Singh, R., Vatsa, M. (2019). Triplet transform learning for automated primate face recognition. *International Conference on Image Processing (ICIP).* https://doi.org/ 10.1109/ICIP.2019.8803501

Amari, S., & Nagaoka, H. (2000). *Methods of information geometry.* American Mathematical Society.

Arandjelovic, R., Gronat, P., Torii, A., Pajdla, T., Sivic, J. (2016). NetVLAD: CNN Architecture for Weakly Supervised Place Recognition. *Conference on Computer Vision and Pattern Recognition (CVPR).* https://doi .org/10.1109/CVPR.2016.572

Arandjelović, R., & Zisserman, A. (2012). Three things everyone should know to improve object retrieval. *Conference on Computer Vision and Pattern Recognition (CVPR).* https://doi.org/10.1109/CVPR .2012.6248018

Araujo, G., Ismail, A., McCann, C., McCann, D., Legaspi, C., Snow, S., ... Ponzo, A. (2020). Getting the most out of citizen science for endangered species such as Whale Shark. *Journal of Fish Biology*, *96*, 864–867. https://doi.org/10.1111/jfb.14254

Babenko, A., & Lempitsky, V. (2015). Aggregating Deep Convolutional Features for Image Retrieval. *International Conference on Computer Vision (ICCV).* https://doi.org/ 10.1109/iccv.2015.150

Babenko, A., Slesarev, A., Chigorin, A., Lempitsky, V. (2014). Neural Codes for Image Retrieval. *European Conference on Computer Vision (ECCV).* https://doi.org/10 .1007/978-3-319-10590-1_38

Badreldeen Bdawy Mohamed, O. (2021). *Metric learning based pattern matching for species agnostic animal re-identification.* Master's thesis, Lappeenranta-Lahti University of Technology LUT, Finland.

Balntas, V., Lenc, K., Vedaldi, A., Mikolajczyk, K. (2017). Hpatches: A benchmark and evaluation of handcrafted and learned local descriptors. *Conference on Computer Vision and Pattern Recognition (CVPR).* https:// doi.org/10.1109/TPAMI.2019.2915233

Baumberg, A. (2000). Reliable feature matching across widely separated views. *Conference on Computer Vision and Pattern Recognition (CVPR).* https://doi.org/10.1109/ CVPR.2000.855899

Berger-Wolf, T., Rubenstein, D., Stewart, C., Holmberg, J., Parham, J., Crall, J. (2015). Ibeis: Image-based ecological information system: From pixels to science and conservation. *Bloomberg Data for Good Exchange Conference.*

Berger-Wolf, T.Y., Rubenstein, D.I., Stewart, C.V., Holmberg, J.A., Parham, J., Menon, S., ... Joppa, L. (2017). Wildbook: Crowdsourcing, computer vision, and data science for conservation. *arXiv preprint arXiv:1710.08880.*

Bergstra, J., Bardenet, R., Bengio, Y., Kégl, B. (2011). Algorithms for Hyper-Parameter Optimization. *Conference on Neural Information Processing Systems (NeurIPS).*

Bogucki, R., Cygan, M., Khan, C.B., Klimek, M., Milczek, J.K., Mucha, M. (2019). Applying deep learning to right whale photo identification. *Conservation Biology*, *33*, 676–684. https://doi.org/10.1111/cobi.13226

Brown, M., & Lowe, D.G. (2007). Automatic Panoramic Image Stitching using Invariant Features. *International Journal of Computer Vision*, *74*, 59–73. https://doi.org/10.1007/s11263-006-0002-3

Brust, C.-A., Burghardt, T., Groenenberg, M., Kading, C., Kuhl, H.S., Manguette, M.L., Denzler, J. (2017). Towards automated visual monitoring of individual gorillas in the wild. *International Conference on Computer Vision Workshop (ICCVW).* https://doi.org/10.1109/iccvw.2017.333

Cheeseman, T., Johnson, T., Muldavin, N. (2017). *Happywhale: Globalizing marine mammal photo identification via a citizen science web platform.* Paper SC/67A/PH/02 presented to the Scientific Committee of the Report to the International Whaling Commission.

Chehrsimin, T., Eerola, T., Koivuniemi, M., Auttila, M., Levänen, R., Niemi, M., ... Kälviäinen, H. (2018). Automatic individual identification of Saimaa ringed seals. *IET Computer Vision*, *12*, 146–152. https://doi.org/10.1049/iet-cvi.2017.0082

Chelak, I., Nepovinnykh, E., Eerola, T., Kälviäinen, H., Belykh, I. (2021). EDEN: Deep Feature Distribution Pooling for Saimaa Ringed Seals Pattern Matching. *arXiv preprint arXiv:2105.13979.*

Crall, J., Stewart, C., Berger-Wolf, T., Rubenstein, D., Sundaresan, S. (2013). Hotspotter - patterned species instance recognition. *Winter Conference on Applications of Computer Vision (WACV).* https://doi.org/10.1109/2013.6475023

Crouse, D., Jacobs, R., Richardson, Z., Klum, S., Jain, A., Baden, A., Tecot, S. (2017). Lemurfaceid: A face recognition system to facilitate individual identification of lemurs. *BMC Zoology*, *2*, 1–14. https://doi.org/10.1186/s40850-016-0011-9

Deb, D., Wiper, S., Gong, S., Shi, Y., Tymoszek, C., Fletcher, A., Jain, A.K. (2018). Face recognition: Primates in the wild. *International Conference on Biometrics Theory, Applications and Systems (BTAS).* https://doi.org/10.1109/btas.2018.8698538

Fischler, M.A., & Bolles, R.C. (1981). Random sample consensus: a paradigm for model fitting with applications to image analysis and automated cartography. *Communications of the ACM*, *24*, 381–395. https://doi.org/10.1145/358669.358692

Freytag, A., Rodner, E., Simon, M., Loos, A., Kühl, H., Denzler, J. (2016). Chimpanzee Faces in the Wild: Log-Euclidean CNNs for Predicting Identities and Attributes of Primates. *German Conference on Pattern Recognition (GCPR).* https://doi.org/10.1007/978-3-319-45886-1_5

Gong, Y., Wang, L., Guo, R., Lazebnik, S. (2014). Multi-scale Orderless Pooling of Deep Convolutional Activation Features. *European Conference on Computer Vision (ECCV).* https://doi.org/10.1007/978-3-319-10584-0_26

Harris, C.G., & Stephens, M.J. (1988). A Combined Corner and Edge Detector. *Alvey Vision Conference.* https://doi.org/10

.5244/c.2.23

Hartley, R., & Zisserman, A. (2004). *Multiple View Geometry in Computer Vision*. Cambridge University Press. https://doi.org/10.1017/CBO9780511811685

Hartwig, S. (2005). Individual acoustic identification as a non-invasive conservation tool: An approach to the conservation of the African wild dog Lycaon pictus (Temminck, 1820). *Bioacoustics The International Journal of Animal Sound and its Recording*, *15*, 35–50. https://doi.org/10.1080/09524622.2005.9753537

He, K., Gkioxari, G., Dollár, P., Girshick, R. (2017). Mask R-CNN. *International Conference on Computer Vision (ICCV)*. https://doi.org/10.1109/iccv.2017.322

Holmberg, J., Norman, B., Arzoumanian, Z. (2009). Estimating population size, structure, and residency time for whale sharks Rhincodon typus through collaborative photo-identification. *Endangered Species Research*, *7*, 39–53. https://doi.org/10.3354/esr00186

Hutchison, D., Kanade, T., Kittler, J., Kleinberg, J.M., Mattern, F., Mitchell, J.C., . . . Mensink, T. (2010). Improving the Fisher Kernel for Large-Scale Image Classification. *European Conference on Computer Vision (ECCV)*. https://doi.org/10.1007/978-3-642-15561-1_11

Jaakkola, T., & Haussler, D. (1999). Exploiting Generative Models in Discriminative Classifiers. *Conference on Neural Information Processing Systems (NeurIPS)*.

Jégou, H., Douze, M., Schmid, C., Pérez, P. (2010). Aggregating local descriptors into a compact image representation. *Conference on Computer Vision and Pattern Recognition (CVPR)*. https://doi.org/10.1109/CVPR.2010.5540039

Koivuniemi, M., Auttila, M., Niemi, M., Levänen, R., Kunnasranta, M. (2016). Photo-ID as a tool for studying and monitoring the endangered Saimaa ringed seal. *Endangered Species Research*, *30*, 29–36. https://doi.org/10.3354/esr00723

Koivuniemi, M., Kurkilahti, M., Niemi, M., Auttila, M., Kunnasranta, M. (2019). A mark–recapture approach for estimating population size of the endangered ringed seal (Phoca hispida saimensis). *PLOS ONE*, *14*, 214–269. https://doi.org/10.1371/journal.pone.0214269

Korschens, M., & Denzler, J. (2019). ELPephants: A Fine-Grained Dataset for Elephant Re-Identification. *International Conference on Computer Vision Workshop (ICCVW)*. https://doi.org/10.1109/iccvw.2019.00035

Kumar, S., Pandey, A., Sai Ram Satwik, K., Kumar, S., Singh, S.K., Singh, A.K., Mohan, A. (2018). Deep learning framework for recognition of cattle using muzzle point image pattern. *Measurement*, *116*, 1–17. https://doi.org/10.1016/j.measurement.2017.10.064

Kunnasranta, M., Niemi, M., Auttila, M., Valtonen, M., Kammonen, J., Nyman, T. (2021). Sealed in a lake – Biology and conservation of the endangered Saimaa ringed seal: A review. *Biological Conservation*, *253*, 108908. https://doi.org/10.1016/j.biocon.2020.108908

Li, S., Li, J., Tang, H., Qian, R., Lin, W. (2020). ATRW: A Benchmark for Amur Tiger Re-identification in the Wild. *ACM International Conference on Multimedia*. https://doi.org/10.1145/3394171.3413569

Lindeberg, T. (1998). Feature Detection with Automatic Scale Selection. *International Journal of Computer Vision*, *30*, 77–116. https://doi.org/10.1023/A:1008045108935

Liu, C., Zhang, R., Guo, L. (2019). Part-Pose Guided Amur Tiger Re-Identification. *International Conference on Computer Vision Workshop (ICCVW)*. https://doi.org/10.1109/ICCVW.2019.00042

Liu, N., Zhao, Q., Zhang, N., Cheng, X.,

Zhu, J. (2019). Pose-Guided Complementary Features Learning for Amur Tiger Re-Identification. *International Conference on Computer Vision Workshop (ICCVW)*. https://doi.org/10.1109/ICCVW.2019.00038

Lowe, D. (2004). Distinctive Image Features from Scale-Invariant Keypoints. *International Journal of Computer Vision*, *60*, 91–110. https://doi.org/10.1023/B:VISI.0000029664.99615.94

Lowe, D.G. (1999). Object Recognition from Local Scale-Invariant Features. *International Conference on Computer Vision (ICCV)*. https://doi.org/10.5555/850924.851523

MacQueen, J., et al. (1967). Some methods for classification and analysis of multivariate observations. *Berkeley Symposium on Mathematical Statistics and Probability*.

Malik, J., Kiranyaz, S., Gabbouj, M. (2021). Self-organized operational neural networks for severe image restoration problems. *Neural Networks*, *135*, 201–211. https://doi.org/10.1016/j.neunet.2020.12.014

Mantiuk, R., Myszkowski, K., Seidel, H.-P. (2006). A perceptual framework for contrast processing of high dynamic range images. *ACM Transactions on Applied Perception*, *3*, 286–308. https://doi.org/10.1145/1166087.1166095

McCoy, E., Burce, R., David, D., Aca, E., Hardy, J., Labaja, J., . . . Araujo, G. (2018). Long-Term Photo-Identification Reveals the Population Dynamics and Strong Site Fidelity of Adult Whale Sharks to the Coastal Waters of Donsol, Philippines. *Frontiers in Marine Science*, *5*, 271. https://doi.org/10.3389/fmars.2018.00271

McLachlan, G.J., & Basford, K.E. (1988). *Mixture models: Inference and applications to clustering*. M. Dekker New York.

Mikolajczyk, K., & Schmid, C. (2002). An affine invariant interest point detector. *European Conference on Computer Vision (ECCV)*. https://doi.org/10.1007/3-540-47969-4_9

Mikolajczyk, K., & Schmid, C. (2004). Scale & Affine Invariant Interest Point Detectors. *International Journal of Computer Vision*, *60*, 63–86. https://doi.org/10.1023/B:VISI.0000027790.02288.f2

Mishchuk, A., Mishkin, D., Radenovic, F., Matas, J. (2017). Working hard to know your neighbor's margins: Local descriptor learning loss. *Conference on Neural Information Processing Systems (NeurIPS)*.

Mishkin, D., Radenović, F., Matas, J. (2018). Repeatability Is Not Enough: Learning Affine Regions via Discriminability. *European Conference on Computer Vision (ECCV)*. https://doi.org/10.1007/978-3-030-01240-3_18

Moskvyak, O., Maire, F., Dayoub, F., Armstrong, A.O., Baktashmotlagh, M. (2021). Robust re-identification of manta rays from natural markings by learning pose invariant embeddings. *International Conference on Digital Image Computing: Techniques and Applications (DICTA)*. https://doi.org/10.1109/DICTA52665.2021.9647359

Moskvyak, O., Maire, F., Dayoub, F., Baktashmotlagh, M. (2021). Keypoint-Aligned Embeddings for Image Retrieval and Re-identification. *Winter Conference on Applications of Computer Vision (WACV)*. https://doi.org/10.1109/48630.2021.00072

Nepovinnykh, E., Chelak, I., Lushpanov, A., Eerola, T., Kälviäinen, H., Chirkova, O. (2022). Matching individual Ladoga ringed seals across short-term image sequences. *Mammalian Biology*, 1-16. https://doi.org/10.1007/s42991-022-00229-3

Nepovinnykh, E., Eerola, T., Biard, V., Mutka, P., Niemi, M., Kunnasranta, M., Kälviäinen, H. (2022). SealID: Saimaa ringed seal re-identification database. *arXiv preprint arXiv:2206.02260*.

Nepovinnykh, E., Eerola, T., Kälviäinen, H.

(2020). Siamese Network Based Pelage Pattern Matching for Ringed Seal Re-identification. *Winter Conference on Applications of Computer Vision Workshops (WACVW).* https://doi.org/10.1109/wacvw50321.2020.9096935

Nepovinnykh, E., Eerola, T., Kälviäinen, H., Radchenko, G. (2018). Identification of Saimaa Ringed Seal Individuals Using Transfer Learning. *International Conference on Advanced Concepts for Intelligent Vision Systems (ACIVS).* https://doi.org/10.1007/978-3-030-01449-0_18

Nipko, R., Holcombe, B., Kelly, M. (2020). Identifying Individual Jaguars and Ocelots via Pattern-Recognition Software: Comparing HotSpotter and Wild-ID. *Wildlife Society Bulletin*, *44*, 424-433. https://doi.org/10.1002/wsb.1086

Norouzzadeh, M.S., Nguyen, A., Kosmala, M., Swanson, A., Palmer, M.S., Packer, C., Clune, J. (2018). Automatically identifying, counting, and describing wild animals in camera-trap images with deep learning. *Proceedings of the National Academy of Sciences*, *115*, 5716–5725. https://doi.org/10.1073/pnas.1719367115

Parham, J.R., Crall, J., Stewart, C., Berger-Wolf, T., Rubenstein, D. (2017). Animal Population Censusing at Scale with Citizen Science and Photographic Identification. *AAAI Spring Symposium Series.*

Perronnin, F., & Dance, C. (2007). Fisher Kernels on Visual Vocabularies for Image Categorization. *Conference on Computer Vision and Pattern Recognition (CVPR).* https://doi.org/10.1109/CVPR.2007.383266

Perronnin, F., Liu, Y., Sánchez, J., Poirier, H. (2010). Large-scale image retrieval with compressed Fisher vectors. *Conference on Computer Vision and Pattern Recognition (CVPR).* https://doi.org/10.1109/CVPR.2010.5540009

Pruchova, A., Jaška, P., Linhart, P. (2017). Cues to individual identity in songs of songbirds: testing general song characteristics in Chiffchaffs Phylloscopus collybita. *Journal of Ornithology*, *158*, 911–924. https://doi.org/10.1007/s10336-017-1455-6

Radenović, F., Tolias, G., Chum, O. (2018). Fine-tuning CNN image retrieval with no human annotation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, *41*, 1655–1668. https://doi.org/10.1109/tpami.2018.2846566

Razavian, A.S., Sullivan, J., Carlsson, S., Maki, A. (2016). Visual Instance Retrieval with Deep Convolutional Networks. *ITE Transactions on Media Technology and Applications*, *4*, 251–258. https://doi.org/10.3169/mta.4.251

Ronneberger, O., Fischer, P., Brox, T. (2015). U-Net: Convolutional Networks for Biomedical Image Segmentation. *International conference on medical image computing and computer assisted intervention (MICCAI).* https://doi.org/10.1007/978-3-319-24574-4_28

Sánchez, J., Perronnin, F., Mensink, T., Verbeek, J. (2013). Image Classification with the Fisher Vector: Theory and Practice. *International Journal of Computer Vision*, *105*, 222–245. https://doi.org/10.1007/s11263-013-0636-x

Schneider, S., Taylor, G., Kremer, S. (2020). Similarity Learning Networks for Animal Individual Re-Identification - Beyond the Capabilities of a Human Observer. *Winter Applications of Computer Vision Workshops (WACVW).* https://doi.org/10.1109/WACVW50321.2020.9096925

Schneider, S., Taylor, G.W., Linquist, S., Kremer, S.C. (2019). Past, present and future approaches using computer vision for animal re-identification from camera trap data. *Methods in Ecology and Evolution*, *10*, 461–470. https://doi.org/10.1111/2041-210x.13133

Schölkopf, B., Smola, A., Müller, K.-R. (1998). Nonlinear Component Analysis as a Kernel

Eigenvalue Problem. *Neural Computation*, *10*, 1299–1319. https://doi.org/10.1162/089976698300017467

Sivic, & Zisserman. (2003). Video Google: a text retrieval approach to object matching in videos. *International Conference on Computer Vision (ICCV)*. https://doi.org/10.1109/ICCV.2003.1238663

Smeulders, A., Worring, M., Santini, S., Gupta, A., Jain, R. (2000). Content-based image retrieval at the end of the early years. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, *22*, 1349–1380. https://doi.org/10.1109/34.895972

Thompson, J., Zero, V., Schwacke, L., Speakman, T., Quigley, B., Morey, J., McDonald, T. (2019). finFindR: Computer-assisted Recognition and Identification of Bottlenose Dolphin Photos in R. *bioRxiv*, 825661. https://doi.org/10.1101/825661

Titterington, D.M., Afm, S., Smith, A.F., Makov, U., et al. (1985). *Statistical analysis of finite mixture distributions.* John Wiley & Sons Incorporated.

Tolias, G., Sicre, R., Jégou, H. (2016). Particular Object Retrieval With Integral Max-Pooling of CNN Activations. *International conference on learning representations (ICLR).*

Vidal, M., Wolf, N., Rosenberg, B., Harris, B., Mathis, A. (2021). Perspectives on Individual Animal Identification from Biology and Computer Vision. *Integrative and Comparative Biology*, *61*, 900-916. https://doi.org/10.1093/icb/icab107

Weideman, H.J., Jablons, Z.M., Holmberg, J., Flynn, K., Calambokidis, J., Tyson, R.B., . . . others (2017). Integral curvature representation and matching algorithms for identification of dolphins and whales. *International Conference on Computer Vision Workshop (ICCVW)*. https://doi.org/10.1109/iccvw.2017.334

Yeleshetty, D., Spreeuwers, L., Li, Y. (2020).

3D Face Recognition For Cows. *International Conference of the Biometrics Special Interest Group (BIOSIG).*

Yu, H., Xu, Y., Zhang, J., Zhao, W., Guan, Z., Tao, D. (2021). AP-10k: A benchmark for animal pose estimation in the wild. *Conference on Neural Information Processing Systems (NeurIPS) Datasets and Benchmarks Track.*

Zavialkin, D. (2020). *CNN-based ringed seal pelage pattern extraction.* Master's thesis, Lappeenranta-Lahti University of Technology LUT, Finland.

Zhelezniakov, A., Eerola, T., Koivuniemi, M., Auttila, M., Levänen, R., Niemi, M., . . . Kälviäinen, H. (2015). Segmentation of Saimaa ringed seals for identification purposes. *International Symposium on Visual Computing (ISVC).* https://doi.org/10.1007/978-3-319-27863-6_21