# Robust animal activity recognition using wearable sensors: A correlation distillation-based information recovery method toward data having low sampling rates

A. X. Mao[1], E. D. Huang[2], M. L. Zhu[3] and K. Liu[1,*]

[1]*Department of Infectious Diseases and Public Health, City University of Hong Kong, Hong Kong SAR*
[2]*Department of Computer Science, City University of Hong Kong, Hong Kong SAR*
[3]*Department of Mechanical Engineering, City University of Hong Kong, Hong Kong SAR*
[*]Corresponding author: Kai Liu, kailiu@cityu.edu.hk

## Abstract

Automated animal activity recognition (AAR) has succeeded dramatically due to recent sensing technologies and deep learning advances, enhancing animal health and welfare. Since animals need to be monitored over a long period, factors influencing the energy consumption of sensing devices must be considered carefully. As a critical factor, the sampling rate greatly affects energy usage, battery life, and data storage. To reduce energy costs, existing works often lower sampling rates. However, when the sampling rate falls below a limited threshold, the recognition performance would degrade rapidly due to missing many relevant signals. Therefore, this study proposed a novel correlation distillation-based information recovery (CDIR) method to improve the performance of AAR at low sampling rates. Specifically, we took two convolutional neural networks trained using data having higher and lower sampling rates as the teacher and student models, respectively. The CDIR enabled the student to mimic correlations within teacher features, facilitating missing information recovery. To evaluate its effectiveness, we conducted experiments on a public dataset acquired from six horses using tri-axial accelerometers and gyroscopes with 100 Hz. Data having low sampling rates were obtained by down-sampling the original data at different frequencies (e.g., 50, and 25 Hz). The experimental results demonstrated that our CDIR remarkably boosted the model trained on low sampling rate data (e.g., percentage-point increments in the precision, recall, F1-score, and accuracy of 2.98%, 3.22%, 3.05%, and 1.89%, respectively, for the 12.5-Hz data) while outperforming the existing KD algorithms. This inspired the development of energy-efficient animal monitoring systems.

**Keywords**: behavioral classification, deep learning, resampling, knowledge distillation

## Introduction

Advancements in sensing technologies and smart computing techniques are driving the rapid development of automated and precision livestock management (Kleanthous et al., 2022). Automated animal activity recognition (AAR) is one sector that benefits considerably from using such technologies. Particularly, wearable sensors have been widely used as part of an animal activity monitoring system and in conjunction with deep learning to infer the daily behaviors of individual animals (Smith et al., 2015). Monitoring these activities in real-time could allow early identification of time-sensitive health issues and promote timely human interventions when necessary, effectively improving animal health and welfare (Eerdekens et al., 2021; Mao et al., 2022a, 2022b).

In practical automated AAR systems using wearable devices, as the monitoring of animals needs to be carried out over a long period (e.g., a few weeks up to several months), factors influencing the energy consumption of sensing devices must be considered carefully. As one of the most critical factors, the sampling rate largely affects power usage, battery life, and data storage (Eerdekens et al., 2021; Walton et al., 2018). Walton et al. (2018) stated that increasing the sampling frequency from 16 to 32 Hz reduced the battery life by up to half

(1.81 years vs 3.08 years). Considering practical benefits, existing works often lower sampling rates to reduce energy costs (Riaboff et al., 2022; Walton et al., 2018). Nevertheless, when the sampling rate falls below a limited threshold, the AAR performance would degrade rapidly due to missing many relevant signals. Hence, finding a good way to recover these missing details is necessary for an AAR system.

Knowledge distillation (KD), pioneered by Hinton et al. (2015), transfers the knowledge from a cumbersome teacher model to a small one that is more suitable for deployment, so that the small model's performance can be significantly boosted compared to training the small model alone (Hinton et al., 2015). Particularly, KD has been verified to make it possible to improve the performance of any machine learning algorithms while minimizing deployment and computational resources (Lin et al., 2022). Inspired by the knowledge transfer mechanism in KD, we explored the possibility of transferring knowledge obtained from high sampling rate data to the model trained on low sampling rate data, aiming to alleviate the information loss resulting from the low sampling rate.

In this study, we proposed a novel method, dubbed correlation distillation-based information recovery (CDIR), to improve the performance of AAR at low sampling rates. Specifically, we took data having higher and lower sampling rates as the teacher and student data, and corresponding trained models as the teacher and student models, respectively. The CDIR could guide the student to mimic correlation within teacher features, further recovering the missing information in student features. To evaluate the effectiveness of our method, we conducted experiments based on a public dataset, which was acquired from six horses (Kamminga et al., 2019) using tri-axial accelerometers and tri-axial gyroscopes with an initial sampling rate of 100 Hz.

## Materials and methods

### Experimental data

The used data was a public dataset created by Kamminga et al. (2019), and it comprised 87,621 2-s labelled samples acquired from six individual horses with neck-attached inertia measurement units. The sampling rate was set at 100 Hz for both the tri-axial accelerometer and tri-axial gyroscope and 12 Hz for the tri-axial magnetometer. Amongst, six activities were included, i.e., eating, galloping, standing, trotting, walking-natural, and walking-rider, and the sample numbers were 16,048, 3,939, 5,113, 25,076, 3,327, and 34,118, respectively. Herein, the tri-axial motion data from the accelerometer and gyroscope were utilized, forming a tensor of $1 \times 6 \times 200$ for each instance. Before being fed into the network, all samples need to be normalized, as Mao et al. (2021).

### Correlation distillation-based information recovery

The proposed correlation distillation-based information recovery (CDIR) method could transfer the knowledge contained in the teacher model to the student model, consequently facilitating the student's missing information recovery. Figure 1 illustrates the overall architecture of the proposed CDIR. Given a teacher instance and a student instance with the same classification label, we first upsampled the student input through the nearest neighbor interpolation (Thévenaz et al., 2000) to keep its shape equal to that of teacher input. Then, the teacher input and upsampled input were separately fed into the teacher model T and student model S to yield the multi-layer teacher feature maps $\emptyset F_{s@\&}^{c'}$ and student feature maps $\emptyset F_{s@\&}^{u'}$.

These feature maps were then fed into the correlation distillation (CD) module at the corresponding layer to calculate loss $\mathcal{L}_{vw}$, which could transfer the time- and axis-aware correlations within feature maps from teacher to student. In addition, the output logit was also obtained from the student model and used to

417

compute the classification loss $\mathcal{L}_{vx}$ combined with the original class label, thereby supervising the student model training. Finally, the overall student network was trained by a joint loss function:

$$\mathcal{L}_{pjpls} = \mathcal{L}_{vx} + \lambda * \mathcal{L}_{vw}, \tag{1}$$

where $\mathcal{L}_{vx}$ denoted the cross-entropy loss and $\lambda$ was the weight factor. Note that the teacher model T was pre-trained using the teacher data in a pure classification way and would be kept fixed during the student model training. The enhanced student classification network can then be directly used in the inference phase in practical scenarios with low sampling rates.



Figure 1: Overall architecture of our correlation distillation-based information recovery.

## Correlation distillation module

Existing literatures have proved that animal behavior movement patterns often show intrinsic periodicity within a particular duration, implying that signals between different temporal positions have particular correlations (Casella et al., 2020; Chakravarty et al., 2019). Meanwhile, specific dependencies also exist in the signals from individual sensing axes (Chen et al., 2018). However, these correlations tend to become weak or even disappear under low sampling rates as some vital information is missing. Hence, we devised a CD module (Figure 1) to enable the student to mimic the time- and axis-aware correlations within teacher feature maps, further facilitating the recovery of the student's missing information.

As shown in Figure 1, let $F_s^c \in R^{n \times l \times p}$ and $F_s^u \in R^{n \times l \times p}$ represent the feature maps of teacher and student on the same layer l, respectively. Here, c denoted the feature map's channel number, and a and t implicitly referred to the dimensions related to the axis and time, respectively. Specifically, both $F_s^c$ and $F_s^u$ were taken as inputs to two branches of the CD module, i.e., the time-aware correlation distillation (TCD) branch and the axis-aware correlation distillation (ACD) branch. In the TCD branch, we first reshaped $F_s^c$ and $F_s^u$ into two 2D feature maps with a size of $t * ac$, where each vector with length ac was regarded as the feature along the time dimension. Then, we conducted matrix multiplication on each 2D feature map and its corresponding transposition, generating an inter-time correlation matrix $\hat{M}_s^c$ and $\hat{M}_s^u$ for the teacher and student, respectively:

$$\ddot{M}_s^c = \text{Res}(F_s^c) \cdot \text{Res}(F_s^c)^y, \tag{2}$$

$$\ddot{M}_s^u = \text{Res}(F_s^u) \cdot \text{Res}(F_s^u)^y, \tag{3}$$

where $\text{Res}$ and $(\cdot)^y$ represented the reshaping and transposing operations mentioned above, respectively. Afterwards, we performed the $L_2$ normalization of these two matrices over the horizontal dimension. To recover the student's missing information, we encouraged the student to mimic similar time-aware correlations to the teacher by penalizing the negative cosine similarity between the normalized vectors within the teacher's and student's correlation matrices:

$$\mathcal{L}_s^{\ddot{E}} = -\cos\_sim(\text{norm}(\ddot{M}_s^c), \text{norm}(\ddot{M}_s^u))/t, \tag{4}$$

where $\mathcal{L}_s^{\ddot{E}}$ denoted the loss function of the TCD branch, $\text{norm}(\cdot)$ denoted the $L_2$ normalization, and $\cos\_sim(\cdot,\cdot)$ denoted the cosine similarity function.

Meanwhile, similar operations were conducted in the ACD branch to transfer the axis-aware correlations within feature maps from the teacher to the student (Figure 1). Briefly, $F_s^c$ and $F_s^u$ were first reshaped into two 2D feature maps with a size of $a * tc$, which were used to generate two inter-axis correlation matrices $M_s^c$ and $M_s^u$ by carrying out the matrix multiplication. Then, these two inter-axis correlation matrices were normalized using the $L_2$ normalization over the horizontal dimension, and further used to construct the negative cosine similarity loss. The details are formulated as follows:

$$M_s^c = \text{Res}(F_s^c) \cdot \text{Res}(F_s^c)^y, \tag{5}$$

$$M_s^u = \text{Res}(F_s^u) \cdot \text{Res}(F_s^u)^y, \tag{6}$$

$$\mathcal{L}_s^a = -\cos\_sim(\text{norm}(M_s^c), \text{norm}(M_s^u))/a, \tag{7}$$

where $\text{Res}$ and $\mathcal{L}_s^a$ represented the reshaping operation and loss function in the ACD branch. Finally, the total loss $\mathcal{L}_{vw\{,s}$ of the CD module in the l-th layer can be formulated as the loss combination of TCD branch and ACD branch:

$$\mathcal{L}_{vw,s} = \beta_\& * \mathcal{L}_s^{\ddot{E}} + \beta_* * \mathcal{L}_s^a, \tag{8}$$

where $\beta_\&$ and $\beta_*$ were the weight factors and both set to 1/2 in default.

Considering that the time- and axis-aware correlations are included in intermediate features across various neural network layers, we applied the CD module to all layers. Suppose a total of L pairwise (teacher-student) feature maps from the model, the final loss function $\mathcal{L}_{vw}$ can be defined as the average value of losses across L layers:

$$\mathcal{L}_{vw} = \frac{1}{L} \sum_{s@\&}^{t} \mathcal{L}_{vw,s}. \tag{9}$$

With the help of loss $\mathcal{L}_{vw}$, the CD module could compel the student to imitate the inter-time and inter-axis correlations within the teacher feature map across hierarchical layers. Obviously, these correlations heavily rely on ample information within features. Therefore, minimizing the loss $\mathcal{L}_{vw}$ in Equation 9 can facilitate the recovery of the student's missing information.

<u>Implementation</u>

This study employed the cross-modality interaction network (CMI-Net), which has been previously validated in improving horse activity recognition performance (Mao et al., 2021), as the classification network architecture of the teacher model and student model. Four common evaluation metrics were measured to indicate the comprehensive classification performance, i.e., precision, recall, F1-score, and accuracy. The performances of CMI-Net trained at different sampling rates were used as the baseline values. Herein, we down-sampled the original data sampled at 100 Hz to sampling rates of 50, 25, 12.5, 10, 5, and 2 Hz. Amongst, the sampling rate at which CMI-Net performed best was selected as the sampling rate of the teacher data, and the corresponding trained model served as the teacher model. The student model was then trained using the data having a lower sampling rate than the teacher sampling rate while being guided by the knowledge obtained from the trained teacher model.

During training, we adopted the same configurations as Mao et al. (2021), and we used the grid search method to set exact values for $\lambda$ in Equation 1. The best model with the highest validation accuracy was saved and further verified using test data. To evaluate the effectiveness of our method in improving the student model's performance, we compared it against the baseline method (i.e., training the student model alone) and various existing KD methods (i.e., basic KD (Hinton et al., 2015), AT (Zagoruyko and Komodakis, 2017), ICKD (Liu et al., 2021)). All experiments were executed using the PyTorch framework on a single NVIDIA GeForce RTX 3090 GPU.

## Results and discussion

Overall, the baseline model (CMI-Net) arrived at the best performance when using data having a sampling rate of 25 Hz. Comparison experiments demonstrated that our proposed CDIR remarkably boosted the model trained on data having low sampling rates while exhibiting superior performance to the existing KD methods.

<u>Baseline performance</u>

Figure 2 illustrates the performance of the baseline model under different sampling rates, including 100, 50, 25, 12.5, 10, 5, and 2 Hz. We can observe that as the sampling rate gradually dropped from 100 Hz, the performance increased slightly until it peaked at 25 Hz; that is, a precision, recall, and F1-score of 82.92%, 85.12%, and 83.50%, respectively. This result was consistent with the results of several previous studies on horse behaviour recognition using wearable sensors (Eerdekens et al., 2020a, 2020b, 2021), which found that reducing the sampling rate within a certain range led to higher performance. This phenomenon might be explained by the fact that too high sampling rates raised irrelevant noise in the data, thereby misleading the final behaviour classification. Subsequently, the performance, as expected, started to degrade constantly as the sampling rate continued to drop, which was because that less information was included in data having a low sampling rate. According to the selection criteria described in the implementation details, we chose 25 Hz as the teacher sampling rate.

Figure 2: Classification performance of baseline method for data at different sampling rates (i.e., 100, 50, 25, 12.5, 10, 5, and 2 Hz).

## Comparisons with existing methods

Table 1 presents the experimental results when using the data with a sampling rate of 12.5 Hz ($\lambda$=1.1) and 5 Hz ($\lambda$=1) as the student input. For comparison, we also report the performance of the teacher model that was trained on data obtained at a sampling rate of 25 Hz. The results revealed that all KD approaches outperformed the baseline method in terms of all evaluation metrics, which demonstrated the promising capabilities of the knowledge transfer mechanism in alleviating the performance degradation resulting from the use of data obtained at low sampling rates. Our proposed method exhibited superior performance to existing KD methods with a precision, recall, F1-score, and accuracy of 83.51%, 84.70%, 83.93%, and 91.20%, respectively, for the 12.5-Hz data and 76.77%, 76.36%, 76.45%, and 87.05%, respectively, for the 5-Hz data. This can be ascribed to the ability of our architecture to effectively recover the student's missing information by sufficiently exploiting the teacher's knowledge. Moreover, it was also worth noting that our method even exceeded the teacher model performance with relative percentage-point gains of 0.59%, 0.43%, and 1.13% in precision, F1-score, and accuracy, respectively, for the 12.5-Hz data. This finding was consistent with several previous works (Park et al., 2019; Romero and Kahou, 2015). This was also greatly in line with practical requirements that aimed to relieve the burden of energy costs while maintaining desirable performance.

Table 1: Comparison of our method against the baseline method and existing knowledge distillation methods with low sampling rates.

| Methods[#] | 12.5 Hz | | | | 5 Hz | | | |
|---|---|---|---|---|---|---|---|---|
| | Prec (%) | Rec (%) | F1 (%) | Acc (%) | Prec (%) | Rec (%) | F1 (%) | Acc (%) |
| Teacher | 82.92 | 85.12 | 83.50 | 90.07 | 82.92 | 85.12 | 83.50 | 90.07 |
| Baseline | 80.18 | 81.12 | 80.48 | 89.01 | 74.53 | 74.94 | 74.67 | 85.44 |
| KD (Hinton et al., 2015) | 81.54 | 82.12 | 81.68 | 89.99 | 76.01 | 75.05 | 75.25 | **87.20** |
| AT (Zagoruyko et al., 2017) | 82.45 | 82.83 | 82.61 | **91.17** | 75.07 | **75.76** | 75.17 | 86.47 |
| ICKD (Liu et al., 2021) | 81.36 | 82.43 | 81.77 | 90.15 | 75.15 | 74.59 | 74.39 | 86.02 |
| Our CDIR | **83.16** | **84.34** | **83.53** | 90.90 | **76.45** | 74.99 | **75.57** | 86.45 |

[#] Prec: precision; Rec: recall; F1: F1-score; Acc: accuracy. The best result on student datasets for each metric is highlighted in bold.

## Conclusions

This study proposed a novel method called CDIR to improve the performance of AAR at low sampling rates. The CDIR enabled the student to mimic the time- and axis-aware correlations within teacher feature maps, further facilitating the recovery of the student's missing information. The experimental results revealed that our approach remarkably boosted the performance of a model trained on data obtained at low sampling rates while outperforming existing KD algorithms. This provided considerable inspiration for the practical application of AAR, especially in scenarios with limited energy sources.

## Acknowledgments

## References

Casella, E., Khamesi, A.R., and Silvestri, S. (2020) A framework for the recognition of horse gaits through wearable devices. *Pervasive and Mobile Computing* 67, 101213.

Chakravarty, P., Cozzi, G., Ozgul, A., and Aminian, K. (2019) A novel biomechanical approach for animal behaviour recognition using accelerometers. *Methods in Ecology and Evolution* 10(6), 802–814.

Chen, K., Zhang, D., Yao, L., Guo, B., Yu, Z., and Liu, Y. (2018) Deep learning for sensor-based human activity recognition: overview, challenges and opportunities. *ACM Computing Surveys* 54(4), 1-40.

Eerdekens, A., Deruyck, M., Fontaine, J., Martens, L., De Poorter, E., and Joseph, W. (2020) Automatic equine activity detection by convolutional neural networks using accelerometer data. *Computers and Electronics in Agriculture* 168, 105139.

Eerdekens, A., Deruyck, M., Fontaine, J., Martens, L., de Poorter, E., Plets, D., and Joseph, W. (2020) Resampling and data augmentation for equines' behaviour classification based on wearable sensor accelerometer data using a convolutional neural network. In: Vassos Soteriou *2020 International Conference on Omni-Layer Intelligent Systems.* Barcelona, Spain, pp. 1-6.

Eerdekens, A., Deruyck, M., Fontaine, J., Martens, L., de Poorter, E., Plets, D., and Joseph, W. (2021) A framework for energy-efficient equine activity recognition with leg accelerometers. *Computers and Electronics in Agriculture* 183, 106020.

Hinton, G., Vinyals, O., and Dean, J. (2015) Distilling the knowledge in a neural network. In *Proceedings of Conference on Neural Information Processing Systems, Deep Learning and Representation Learning Workshop.* Quebec, Canada, pp. 1–9.

Kamminga, J.W., Janßen, L.M., Meratnia, N., and Havinga, P.J.M. (2019) Horsing around—A dataset comprising horse movement. *Data* 4(4), 1–13.

Kleanthous, N.T.L., Hussain, A., Khan, W., Sneddon, J., and Liatsis, P. (2022) Deep transfer learning in sheep activity recognition using accelerometer data. *Expert Systems with Applications* 207, 117925.

Lin, S., Xie, H., Wang, B., Yu, K., Chang, X., Liang, X., and Wang, G. (2022) Knowledge distillation via the target-aware transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition.* New Orleans, Louisiana, pp. 10915-10924.

Liu, L., Huang, Q., Lin, S., Xie, H., Wang, B., Chang, X., and Liang, X. (2021) Exploring inter-channel correlation for diversity-preserved knowledge distillation. In *Proceedings of the IEEE International Conference on Computer Vision.* virtually, pp. 8251–8260.

Mao, A., Giraudet, C., Liu, K., De, I., Nolasco, A., Xie, Z., Xie, Z., Gao, Y., Theobald, J., Bhatta, D., Stewart, R., and Mcelligott, A.G. (2022) Automated identification of chicken distress vocalizations using deep learning models. *Journal of Royal Society Interface* 19(191), 20210921.

Mao, A., Huang, E., and Gan, H. (2022) FedAAR : a novel federated learning framework for animal activity recognition with wearable sensors. *Animals* 12(16), 2142.

Mao, A., Huang, E., Gan, H., Parkes, R.S.V, and Xu, W. (2021) Cross-modality interaction network for equine activity recognition using imbalanced multi-modal data. *Sensors* 21(17), 5818.

Park, W., Kim, D., Lu, Y., and Cho, M. (2019) Relational knowledge distillation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition.* Long Beach, CA, pp. 3962–3971.

Riaboff, L., Shalloo, L., Smeaton, A.F., Couvreur, S., Madouasse, A., and Keane, M.T. (2022) Predicting livestock behaviour using accelerometers: A systematic review of processing techniques for ruminant behaviour prediction from raw accelerometer data. *Computers and Electronics in Agriculture* 192, 106610.

Romero, A., and Kahou, S.E. (2015) FitNets : hints for thin deep nets. In *3rd International Conference on Learning Representation.* San Diego, California, USA.

Smith, D., Dutta, R., Hellicar, A., Bishop-Hurley, G., Rawnsley, R., Henry, D., Hills, J., and Timms, G. (2015) Bag of class posteriors, a new multivariate time series classifier applied to animal behaviour identification. *Expert  Systems with Applications* 42(7), 3774–3784.

Thévenaz, P., Blu, T., and Unser, M. (2000) Interpolation revisited. *IEEE Transactions on Medical Imaging* 19(7), 739–758.

Walton, E., Casey, C., Mitsch, J., Vázquez-Diosdado, J.A., Yan, J., Dottorini, T., Ellis, K.A., Winterlich, A., and Kaler, J. (2018) Evaluation of sampling frequency, window size and sensor position for classification of sheep behaviour. *Royal Society Open Science* 5(2), 171442.

Zagoruyko, S., and Komodakis, N. (2017) Paying more attention to attention: Improving the performance of convolutional neural networks via attention transfer. In *5th International Conference on Learning Representations.* Toulon, France, pp. 1–13.