



Structured Inference Networks Using High-Dimensional Sensors for Surveillance Purposes

Vincent Polfliet¹, Nicolas Knudde¹(✉), Baptist Vandersmissen²,
Ivo Couckuyt¹, and Tom Dhaene¹

¹ Department of Information Technology, Ghent University - imec, Ghent, Belgium
`nicolas.knudde@ugent.be`

² Department of Electronics and Information Systems,
Ghent University - imec, Ghent, Belgium

Abstract. Video cameras are arguably the world's most used sensors for surveillance systems. They give a highly detailed representation of a situation that is easily interpreted by both humans and computers. However, these representations can lose part of their representational value when being recorded in less than ideal circumstances. Bad weather conditions, low-light illumination or concealing objects can make the representation more opaque. A radar sensor is a potential solution for these situations, since it is unaffected by the light intensity and can sense through most concealing objects. In this paper, we investigate the performance of a structured inference network on data of a low-power radar device. A structured inference network applies automated feature extraction by creating a latent space out of which the observations can be reconstructed. A classification model can then be trained on this latent space. This methodology allows us to perform experiments for both person identification and action recognition, resulting in competitive error rates ranging from 0% to 6.5% for actions recognition and 10% to 12% for person identification. Furthermore, the possibility of a radar sensor being used as a complement to a camera sensor is investigated.

Keywords: Structured inference network · Person identification
Action recognition · Indoor sensing · Micro-doppler · Sensor fusion

1 Introduction

In recent years, interest in autonomous surveillance systems grew considerably. While these systems improved significantly, the primary sensors remained the same. The dominance of video cameras in autonomous surveillance systems can be explained by their fundamental strengths. They give a detailed high dimensional representation of their environment, which is easily interpretable by humans as well as computers. Moreover, reduction in price and higher resolutions kept driving their success. While the fundamental advantages of video

cameras were a big catalyst for the early development of surveillance systems, their deficiencies are now holding them back. For some of these deficiencies, such as recordings in bad weather conditions or low light environments, workarounds can be found. Others such as concealing clothing are harder to deal with. In contrast, a radar sensor is unaffected by concealing clothing, bad weather conditions, low-light environments and can be placed out of sight, behind a wall.

A radar is an active sensor that transmits an electromagnetic signal, which is reflected by objects in its line of sight. Information about these objects is then extracted out of these signals taking advantage of, e.g. the Doppler effect. Moreover, individual moving parts of a person or object will each reflect their own Doppler signal which are then summarized into a micro-Doppler (MD) signature [3].

These signatures contain information about the movement of the target, providing a promising feature to differentiate between for example cars, bikers, pedestrians or dogs. Another use for these MD signatures is to recognise different actions, ranging from walking to sitting or boxing [10, 11]. However, perhaps the most challenging application is to differentiate individuals based on the way they move, the so called gait-based identification. While there is a noticeable difference between how a dog and a human walks or how a person runs or sits, the difference in the MD signature between two persons walking is more subtle. This subject has been extensively researched, however, previous papers used a high-power radar sensor with relatively simple scenarios. In this paper the data sets are recorded using a low-power frequency modulated continuous wave (FMCW) radar. This radar is a low-cost, power efficient and compact sensor suited for indoor usage. However, the combination of a human’s low radar cross-section and a low-power device poses a significant challenge for this study [4].

Two data sets are used for our experiments. The first uses the IDentification with Radar (IDRad) benchmark, which is an extensive data set where the main objective is to identify individuals moving randomly in a room [19]. An additional data set is recorded where the main objective is to recognise different actions. Previous studies applied either deep convolutional neural networks (DCNN) [11, 19] or clustering methods [9, 20] to MD signatures. Both approaches were successful by exploiting certain properties of the data. The DCNN tries to take advantage of the spatial properties, along the time and velocity axes, of an MD signature. Conversely, the clustering methods are applied on feature vectors of the original noisy data. Hence, a structured inference network (SIN) [14] can potentially exploit both these properties due to its inherent Markovian properties. This model creates a lower dimensional latent space into which each time step is projected without losing their sequential dependencies. The lower dimensional states also implies that the model performs autonomous feature selection on the data. The resulting lower dimensional latent states are then used in a classification model. These properties make the SIN well-suited for high dimensional sequential data, such as radar data.

2 Related Work

There has been extensive research in the use of radar as a sensor. This section will highlight several relevant studies concerning action recognition and person identification. Afterwards some other recent results will be discussed regarding SINS.

Action recognition and *gait-based identification* are discussed in a wide array of studies. The former is usually defined by the amount of different actions in the data set. In [10, 11], 7 actions are proposed, ranging from walking, walking with a stick, running to even boxing. A wide variety of models have been investigated to differentiate between actions. Kim et al. apply a support vector machine with manual engineered features [10] and an DCNN [11]. In [16], transfer learning is applied to a pretrained CNN. In [5], singular value decomposition with multiple classification models were used for detecting violent intents. The studies [7, 15], investigate autonomous surveillance systems as a tool to monitor elderlies using a wide array of classifiers.

Conversely, mainly data driven models are studied for gait-based identification. In [6] k-means and k-NN clustering is used on thirteen subjects with an accuracy ranging from 92.4% to 100%. The authors of [17] also apply k-NN along with two manual engineered features and Kalgaonkar and Raj obtained an accuracy of 90% by using a Gaussian mixture model (GMM) [9]. Finally, the authors of [19] designed a deep convolutional neural network (DCNN) resulting in an accuracy of 81.61% on lower-power radar data.

Radar data can also be used for non-classification purposes such as person tracking [13].

The structured inference network was proposed in [14]. The authors apply the model to the reconstruction of polyphonic music and the counterfactual prediction of electronic health records of patient data. This model was then also used by the authors of [18] to model human poses. A similar black box variational inference model for state space models is proposed in [2]. While an unsupervised model is proposed in [8], which combines the strengths of a latent graphical variational auto-encoder (VAE) and GMM by using a conditional random field as their inference network. The authors apply their model to a data set of a mouse running in a box, where it successfully clusters different movements of the mouse.

3 Micro-doppler

A large object or body moving through a room at a constant speed induces a constant Doppler Frequency shift. However, smaller moving parts can cause additional micro-motion dynamics, which, in their turn, induce Doppler modulations on the echoed signal. This is referred to as the micro-Doppler effect [3] and causes sidebands around the Doppler frequency, representing the different smaller moving parts. The micro-Doppler map can thus be seen as the power reflected as a function of the speed of the reflector. The radar used in the data

sets is a 77 GHz Frequency Modulated Continuous Wave radar. An FMCW radar has the advantage of being power efficient, but comes at the expense of a low signal to noise ratio, which makes analysing this sensor data more challenging.

4 Structured Inference Network

A structured inference network [14] is a subfield of machine learning where it is assumed that the data conforms to the structure of a *Gaussian state space model* (GSSM). A GSSM assumes that the actual states of a situation are only partly observable and that there exist latent states that fully describe the context of the data without any error. These states are then also assumed to be continuous and only dependent on their previous state. However, in data-oriented problems, the parametric form for a GSSM is usually unknown. A solution for this is a *deep Markov model* (DMM): A GSSM where the emission and transition functions are replaced by multi-layer perceptrons (MLP). The resulting GSSM still has the Markovian structure of an hidden Markov model (HMM) but uses the strength of deep neural networks to help model complex data. An example of a DMM can be seen in Fig. 1.

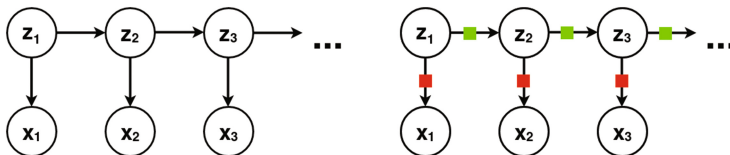


Fig. 1. Generative models of sequential data: **(left)** is a classical HMM. While **(right)** depicts a DMM. The transition (green) and emission (red) functions are both approximated using MLPs. (Color figure online)

The model requires that the latent states are multivariate Gaussian distributions, with a mean and covariance that are functions dependent on the previous latent state. In this paper we also define our observations to be multivariate Gaussian distributions where the parameters are dependent on the latent state. Equation 1 results in a GSSM with model parameters $\theta = \{\alpha, \beta, \kappa, \lambda\}$.

$$\begin{aligned} z_t &\sim \mathcal{N}(G_\alpha(z_{t-1}, \Delta_t), S_\beta(z_{t-1}, \Delta_t)) && (Transition) \\ x_t &\sim \mathcal{N}(G_\kappa(z_t, \Delta_t), S_\lambda(z_t, \Delta_t)) && (Emission) \end{aligned} \quad (1)$$

Another technique needed for this model is *variational learning* [12]. Assume that $p(\mathbf{x}, \mathbf{z}) = p_\theta(\mathbf{z})p_\theta(\mathbf{x}|\mathbf{z})$ is a generative model, where \mathbf{x} is the observation and \mathbf{z} the latent variable. The posterior distribution for this generative model is usually intractable. The variational principle then states that there should be

an approximation of the posterior distribution $q_\phi(\mathbf{z}|\mathbf{x})$. Using this principle, a lower bound of the marginal likelihood is found, which is parameterized by a neural network.

$$\log p_\theta(\mathbf{x}) \geq \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} [\log p_\theta(\mathbf{x}|\mathbf{z})] - \text{KL}(q_\phi(\mathbf{z}|\mathbf{x})||p_\theta(\mathbf{z})) \quad (2)$$

Using variational learning a lower bound is found that approximates the posterior distribution of the GSSM [14].

$$\begin{aligned} \mathcal{L}(\mathbf{X}; (\theta, \phi)) = & \sum_{t=1}^T \mathbb{E}_{q_\phi(\mathbf{z}_t|\mathbf{X})} [\log p_\theta(\mathbf{x}|\mathbf{z})] - \text{KL}(q_\phi(\mathbf{z}_1|\mathbf{X})||p_\theta(\mathbf{z}_1)) \\ & - \sum_{t=2}^T \mathbb{E}_{q_\phi(\mathbf{z}_{t-1}|\mathbf{X})} [\text{KL}(q_\phi(\mathbf{z}_t|\mathbf{z}_{t-1}, \mathbf{X})||p_\theta(\mathbf{z}_t|\mathbf{z}_{t-1}))] \end{aligned} \quad (3)$$

Since the latent states of the generative model will be used for classification purposes, we propose an additional modification. By using a different prior for each classification target, we can encourage the latent state to be more accommodating regarding the classification.

5 Methodology

The main objective of this paper is to investigate the efficiency of a SIN applied to MD signatures for two use cases: gait-based person identification and action recognition. Both data sets were recorded using the same low-power FMCW radar, produced by INRAS [1], in an empty indoor environment. The action recognition data set was recorded to study the performance of radar sensors versus camera sensors.

5.1 Preprocessing

Radar: The MD signature is first achieved by calculating a two-dimensional Fourier transform on the range-Doppler map. Afterwards the absolute values are converted to decibels and are summed over the range dimensions. The raw MD signature contains 256 Doppler channels per time step (with 15 fps). Each of these channels represents a speed ranging from -3.8 m/s to 3.8 m/s. The static channels representing the highest and lowest speeds are removed, without any loss of relevant information. Subsequently, the resulting sequence is thresholded by fixing every point under a certain value. After thresholding, a logarithmic scaling step is applied to compress high activated values, which results in a lower variance. Finally, each Doppler channel will be normalized separately for each sequence. Figure 2 displays the different results of the preprocessing steps to transform a raw MD signature to the fully preprocessed MD signature.

Camera: As the video camera data is only used for basic action recognition, there is no need for highly detailed images. Taking this in consideration, the images were first converted to gray scale and then rescaled from 640×480 pixels to 30×20 pixels. The resulting images are then normalized using the mean pixel values. Finally, the camera images are processed by a small convolutional network, as shown in Fig. 4. A partial copy of the camera data set was also created with half of the image occluded (left side). This area will serve as an artificial screen to check the performance between a camera sensor and a radar sensor in less than ideal circumstances. The intermediate results of the preprocessing and an example of an occluded image are shown in Fig. 3.

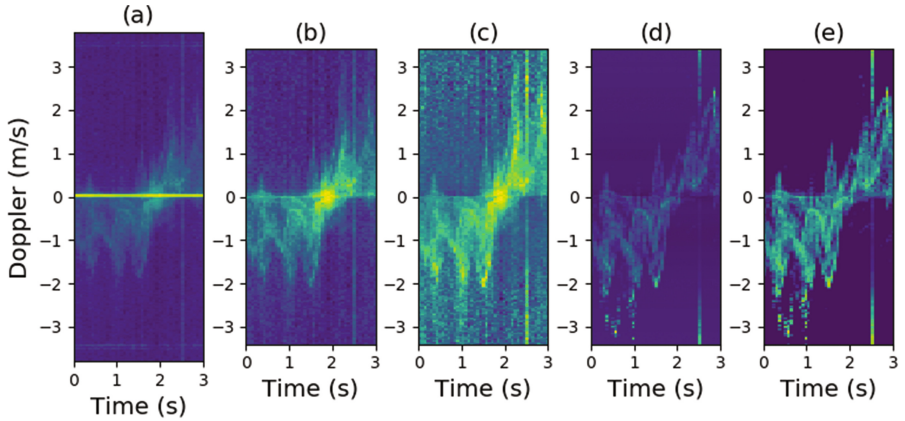


Fig. 2. A 3s MD signature, each figure shows the results of a preprocessing step, with first (a) showing the raw signature. (b) is then obtained by removing the static channels. (c) is the normalized MD signature of (b) and still displays a lot of noise. This is then solved by applying thresholding (d) and finally the variance in the high activated areas is reduced by log scaling (e).



Fig. 3. Camera images from the Actions data set: From left to right we have the raw image (a), conversion to gray scale with rescaling (b), normalized image (c) and the occluded version of the image (d).

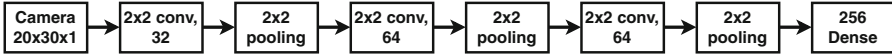


Fig. 4. Convolutional neural network to compress the camera images to lower dimensional vectors.

Sensor Fusion: The high-dimensional radar and camera data are represented by vectors after their respective preprocessing steps. A straightforward form of sensor fusion is obtained by concatenating them. However, both vectors might contain duplicate information. This is filtered out by sending the concatenated vectors through a dense layer. The resulting vector can then be mapped by the SIN to obtain a latent space containing the information of both sensors.

5.2 Model

We implemented the SIN, using the theory mentioned in Sect. 4, in Tensorflow. An outline of the model is shown in Fig. 5. The data will be fed into the Recurrent Neural Network (RNN), which is used as a generative model to create the latent space. Afterwards these states will go through the emission and transition MLPs to find a prediction for respectively the observations and the next latent state. These three predictions and the actual data are then used to calculate the likelihood. Once the SIN is trained, a classification model is applied on the latent states from the generative model. Three different classification models were tested and can be seen in Fig. 6.

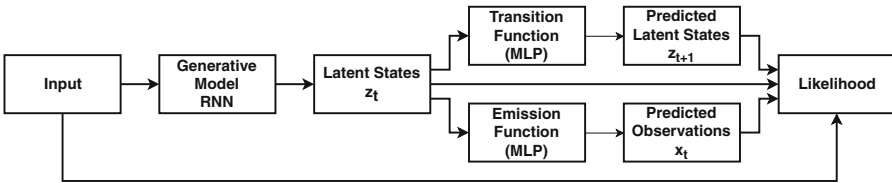


Fig. 5. An outline of the SIN: A generative model is coupled with two MLPs that represent the transition and emission functions.

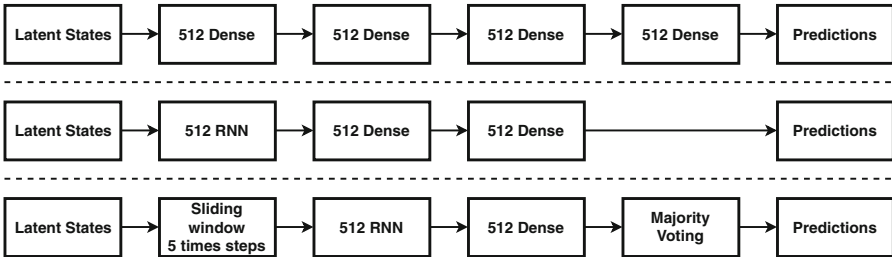


Fig. 6. Three different possible classification models from top to bottom: a MLP, a RNN, and a RNN with majority voting.

6 Experiments

First the efficiency of the model for gait based identification is investigated using the IDRad data set. Afterwards the results of the action recognition data set will be discussed, comparing both camera and radar sensors.

6.1 Person Identification

The IDRad contains recordings of 5 people. Each test person was required to walk for 20 min in random directions with abrupt stops and turns in 2 empty rooms. Each model is trained using sequences of 3 s, which allows us to compare our results with [19].

Analysis of the Generative Model: The classification models are trained on the latent space created by the SIN. However, this model is trained on the likelihood of the reconstruction of the data and is thus independent of the targets of the data. This means that the performance of the classification depends on how well the SIN generalizes the latent space regarding the classification, making the training time of the SIN a hyperparameter. Figure 7 shows the impact of the training time of the SIN on both the classification loss as well as the reconstruction likelihood. While the structured inference network keeps improving over time, the classification model reaches its peak performance in the 100 to 200 epochs interval.

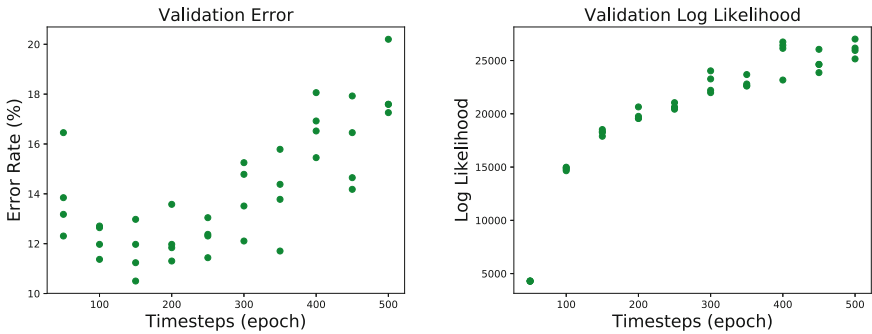


Fig. 7. These figures display the impact of the training time of the SIN on the validation error rate of the classification and on the validation log-likelihood of the SIN itself. It can be seen that while the log-likelihood keeps improving over time, this is not the case for the classification error. The best performing classification models are thus trained on the latent states of SINs with a training time between 100 to 200 epochs.

Results: The structured inference network was trained for 150 epochs and repeated between 10 and 20 times for each experiment.

Table 1 illustrates the impact of the preprocessing. We can see that the results improve when removing the static channels. The results do not improve when adding either the thresholding or log scaling preprocessing step. However when combining both these preprocessing steps, we are able to obtain the best performing models. This is due to the variance of the reconstruction being lowered in the low activated areas by thresholding and in the high activated areas by log scaling. The results thus only improve when both variances are lowered.

Table 2 shows the results of the different classification models on the latent state. Each classification model was tested on the same latent space created by a SIN. Here can be seen that the RNN model outperforms both other models.

Table 1. The impact of adding or removing a subsequent preprocessing step on the error rate. The error rate displays the mean and standard deviation of results over the 5 runs.

Preprocessing	Validation	Test
Raw	$40.90 \pm 3.24\%$	$32.72 \pm 2.21\%$
Remove Static	$21.24 \pm 1.25\%$	$29.59 \pm 1.57\%$
Remove Static + Thresholded	$22.11 \pm 3.13\%$	$32.86 \pm 3.53\%$
Remove Static + Log Scaled	$32.37 \pm 6.54\%$	$38.78 \pm 3.26\%$
Full Preprocessing	$11.92 \pm 0.92\%$	$10.44 \pm 0.76\%$

Table 2. The performance of the different classification models by their error rate.

Classification Model	Validation	Test
MLP	23.14%	18.39%
RNN	11.64%	10.34%
RNN MV	11.78%	10.17%

Table 3. The performance of the two types of structured inference networks and the results of the DCNN as stated in [19]. The error rate displays the mean and standard deviation of results over the 5 runs.

Model	Validation	Test
PCA + RF [19]	48.86%	38.59%
DCNN [19]	24.70%	21.54%
RNN	$15.26 \pm 1.62\%$	$12.20 \pm 0.26\%$
SIN + RNN	$12.24 \pm 1.49\%$	$10.66 \pm 0.74\%$
SIN multiple priors + RNN	$11.92 \pm 0.92\%$	$10.44 \pm 0.76\%$

Finally, Table 3 compares the results found in [19] using a DCNN and principle component analysis with an SVM versus a basic RNN, a SIN and a SIN with different priors. It can be seen that the previous benchmark is improved by up to 12% on the validation set and 11% on the test set using the extra log scaling preprocessing step and a SIN.

6.2 Action Recognition

The data used in this experiment contains radar and camera data of actions generated by 3 people. It consists of 540 samples of 3s each. Each sample represents either a person walking, sitting down or falling. For these experiments the same optimal preprocessing was used as described in Sect. 5.1.

Correlation Between Camera and Radar Sequences: The structured inference network can be used to check for correlation between the two sensors.

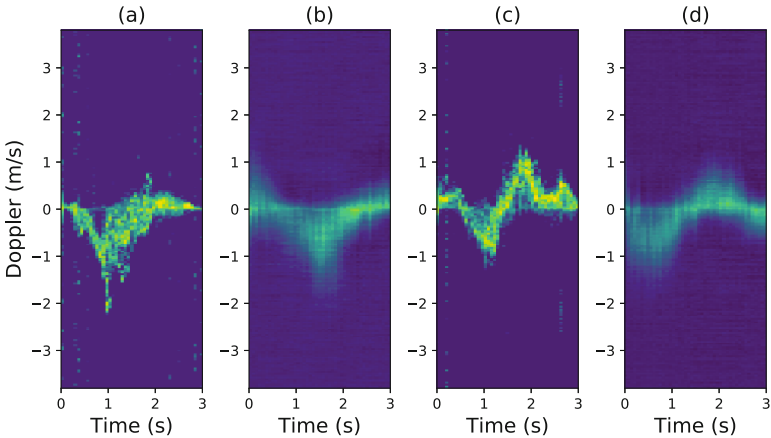


Fig. 8. Artificially generated MD signatures created from camera sequences (b) and (d) versus the original MD signatures (a) and (c).

Table 4. The performance with different sensors by their error rate. SL implies that a screen was artificially inserted on the left side of the camera images, occluding half of the image.

Input	Validation	Test
Radar	$3.78 \pm 1.29\%$	$6.33 \pm 3.13\%$
Camera	$0.11 \pm 0.33\%$	$0.67 \pm 0.74\%$
Camera SL	$4.78 \pm 1.61\%$	$5.56 \pm 2.19\%$
Radar and Camera	$0.11 \pm 0.33\%$	$0.72 \pm 0.73\%$
Radar and Camera SL	$3.33 \pm 1.17\%$	$3.39 \pm 1.43\%$

This is done by training the model on the reconstruction of the first sensor's data and using the second sensor as input. The results of reconstructing MD signatures out of camera sequences can be seen in Fig. 8. While these are not exact reconstructions, the shape of the MD signatures are very similar, confirming the correlation that the log-likelihood of the model suggested.

Results: Table 4 shows difference in results between the radar and camera sensor. The camera data performs better than the radar, with an error rate of 0.67% compared to 6.33%. However, the radar data performs equally well when half of the camera image is occluded. Then by combining the radar and camera data, the problem of the screen was partially alleviated resulting in error rates or 3.33%, which is 2 to 3% lower than the individual sensors.

7 Conclusion and Future Work

We propose to use a classification model on top of the latent space created by a structured inference network and show it outperforms previous methods such as a deep convolutional neural network. This is illustrated on novel use cases of high dimensional camera and radar sequences, where we also show its potential to be used for sensor fusion.

It is noted that the performance of the classification model naturally depends on the amount of trained epochs of the structured inference network, since the latent space is created without consideration of the targets. A possible solution for this could be the unsupervised model mentioned in [8], which combines the strengths of a structured variational auto-encoder with a GMM. Another research point is to apply this model on more challenging radar data, such as walking around with an object or walking in a furnished room.

References

1. Inras gmbh (2017). <http://www.inras.at>
2. Archer, E., Memming Park, I., Buesing, L., Cunningham, J., Paninski, L.: Black box variational inference for state space models. ArXiv e-prints, November 2015
3. Chen, V.C., Li, F., Ho, S.S., Wechsler, H.: Micro-doppler effect in radar: phenomenon, model, and simulation study. *IEEE Trans. Aerosp. Electr. Syst.* **42**(1), 2–21 (2006). <https://doi.org/10.1109/TAES.2006.1603402>
4. Chen, V., Tahmoush, D., Miceli, W.: Radar micro-doppler signatures: processing and applications (2014)
5. Fioranelli, F., Ritchie, M., Griffiths, H.: Classification of unarmed/armed personnel using the netrad multistatic radar for micro-doppler and singular value decomposition features. *IEEE Geosci. Remote Sens. Lett.* **12**(9), 1933–1937 (2015). <https://doi.org/10.1109/LGRS.2015.2439393>
6. Garreau, G., et al.: Gait-based person and gender recognition using micro-doppler signatures. In: 2011 IEEE Biomedical Circuits and Systems Conference (BioCAS), pp. 444–447, November 2011. <https://doi.org/10.1109/BioCAS.2011.6107823>

7. Gurbuz, S.Z., Clemente, C., Balleri, A., Soraghan, J.J.: Micro-doppler-based in-home aided and unaided walking recognition with multiple radar and sonar systems. *IET Radar Sonar Navig.* **11**(1), 107–115 (2017). <https://doi.org/10.1049/iet-rsn.2016.0055>
8. Johnson, M., Duvenaud, D.K., Wiltschko, A., Adams, R.P., Datta, S.R.: Composing graphical models with neural networks for structured representations and fast inference. In: Lee, D.D., Sugiyama, M., Luxburg, U.V., Guyon, I., Garnett, R. (eds.) *Advances in Neural Information Processing Systems* 29, pp. 2946–2954. Curran Associates, Inc. (2016). <http://papers.nips.cc/paper/6379-composing-graphical-models-with-neural-networks-for-structured-representations-and-fast-inference.pdf>
9. Kalgaonkar, K., Raj, B.: Acoustic doppler sonar for gait recognition. In: 2007 IEEE Conference on Advanced Video and Signal Based Surveillance, pp. 27–32, September 2007. <https://doi.org/10.1109/AVSS.2007.4425281>
10. Kim, Y., Ling, H.: Human activity classification based on micro-doppler signatures using a support vector machine. *IEEE Trans. Geosci. Remote Sens.* **47**(5), 1328–1337 (2009). <https://doi.org/10.1109/TGRS.2009.2012849>
11. Kim, Y., Moon, T.: Human detection and activity classification based on micro-doppler signatures using deep convolutional neural networks. *IEEE Geosci. Remote Sens. Lett.* **13**(1), 8–12 (2016). <https://doi.org/10.1109/LGRS.2015.2491329>
12. Kingma, D.P., Welling, M.: Auto-encoding variational bayes. arXiv preprint [arXiv:1312.6114](https://arxiv.org/abs/1312.6114) (2013)
13. Knudde, N., et al.: Indoor tracking of multiple persons with a 77 GHz MIMO FMCW radar. In: 2017 European Radar Conference (EURAD), pp. 61–64, October 2017. <https://doi.org/10.23919/EURAD.2017.8249147>
14. Krishnan, R., Shalit, U., Sontag, D.: Structured inference networks for nonlinear state space models (2017). <https://aaai.org/ocs/index.php/AAAI/AAAI17/paper/view/14215>
15. Liu, L., Popescu, M., Skubic, M., Rantz, M., Yardibi, T., Cuddihy, P.: Automatic fall detection based on doppler radar motion signature. In: 2011 5th International Conference on Pervasive Computing Technologies for Healthcare (PervasiveHealth) and Workshops, pp. 222–225, May 2011. <https://doi.org/10.4108/icst.pervasivehealth.2011.245993>
16. Park, J., Javier, R.J., Moon, T., Kim, Y.: Micro-doppler based classification of human aquatic activities via transfer learning of convolutional neural networks. *Sensors.* **16**(12), 1990 (2016)
17. Tahmouh, D., Silvius, J.: Radar micro-doppler for long range front-view gait recognition. In: 2009 IEEE 3rd International Conference on Biometrics: Theory, Applications, and Systems, pp. 1–6, September 2009. <https://doi.org/10.1109/BTAS.2009.5339049>
18. Toyer, S., Cherian, A., Han, T., Gould, S.: Human pose forecasting via deep markov models. arXiv preprint [arXiv:1707.09240](https://arxiv.org/abs/1707.09240) (2017)
19. Vandersmissen, B., et al.: Indoor person identification using a low-power FMCW radar. *IEEE Trans. Geosci. Remote Sens.* **PP**, 1–12 (2018). <https://doi.org/10.1109/TGRS.2018.2816812>
20. Zhang, Z., Andreou, A.G.: Human identification experiments using acoustic micro-doppler signatures. In: 2008 Argentine School of Micro-Nanoelectronics, Technology and Applications, pp. 81–86, September 2008