

Simulation of pedestrian dynamics based withemantic trajectory segmentation

*Mikhail Rahmanov*¹, *Andrey Shishkin*^{1*}, *Vladimir Komkov*¹, *Irina Alpackaya*²

¹Moscow Aviation Institute (National Research University), Moscow, Russia

²Moscow State University of Civil Engineering, 26, Yaroslavskoe sh., Moscow, 129337, Russia

Abstract. The article analyzes the existing methods of information processing necessary for the functioning of the system of intelligent control over unregulated pedestrian crossings based on aggregation and data processing by means of IOT. The state space model of the switching Kalman filter is considered, the development of mathematical software for the analysis and processing of information based on the results of intelligent control over unregulated pedestrian crossings, in particular with semantic segmentation of trajectories using agent-based models, is carried out. An MDA (Markov Decision Process) state space model is presented, a Hidden Markov Model (HMM) which has discrete hidden variables. The developments for the development of the following subsystems are presented: activity detector subsystem. Receives video frames as input, supports the static object model (background model) and returns the hotspot mask for the current frame; subsystems for detecting and tracking objects (pedestrians and cars). Based on the video frame and hotspot mask, it detects and accompanies objects of a given class, returning their coordinates; trajectory analysis subsystem. Analyzing the history of movement of pedestrians and cars, returns the facts of traffic violations.

1 Introduction

Currently, there is a rapid increase in the use of video surveillance systems, which is explained by a wide range of tasks solved by such systems and the ever-increasing availability of surveillance and communication tools. A significant proportion of intelligent video surveillance systems is roadsurveillance, which is engaged, among other things, in monitoring compliance with traffic rules, monitoring traffic congestion, and detecting traffic accidents. One of the main violations requiring automatic recognition is the identification of failure to provide an advantage to a pedestrian at unregulated pedestrian crossings, since such accidents, according to statistical data, most often end in death or serious harm to the health of a pedestrian. It should be emphasized that when solving these problems, artificial intelligence systems are widely used, for the training of which a sufficiently large number of experimental data, the collection of which in this particular

* Corresponding author: 17andrew07@gmail.com

case is almost impossible. That is why, the training of systems for recognizing traffic violations should be carried out on the basis of data obtained as a result of modeling the mutual movement of pedestrians and vehicles in the pedestrian crossing zone. The latter makes it urgent to develop mathematical models of pedestrian behavior when crossing the roadway.

Currently, there are several approaches to pedestrian modeling that focus on controlling people (microscopic modeling) to create both individual and group pedestrian behavior. Microscopic models of pedestrians take into account individual interactions and try to simulate the position and speed of each pedestrian in time. Among the most representative microscopic seed models of pedestrians are models of cellular automata, behavioral models based on rules, cognitive models, Helbing's model of social forces, and psychological models [1]. In the microscopic simulator, individuals are modeled as independent entities interacting with others and with the environment, making decisions about changing its dynamic state (including calculating the sum of the set of forces as a kind of decision).

The decision-making process in microscopic simulators is carried out according to a hierarchical scheme: strategic, tactical and operational. Destinations and path planning are selected at the strategic level, route selection is made at the tactical level, and instant decisions to change the kinematic state are made at the operational level. Several microscopic simulators focused on the reproduction of local interactions work only at the operational level [6,7].

A common problem in microscopic models is the relationship between individual behavior and group behavior. Traditionally, rule-based systems are most popular for modeling local interactions in this area [8, 4]. However, due to the complexity of preventing collisions of multiple agents, it is difficult to create a realistic group movement that follows local rules [9]. Most agent-based models separate local interactions from the required global path planning. There are two main approaches to this. One of them is to pre-calculate or edit by the user a path planning map, which is represented as a hover field or a potential and velocity field [10]. Another is the separation of local and global navigation tasks in a multi-level model [11]. The advantage of such a division within the agent model is that intellectual or psychological properties can be introduced into the behavior of agents [5, 12]. One indicator that this relationship is being resolved correctly is that certain collective patterns manifest when groups of pedestrians are in certain situations, as they do in the real world. Several types of collective behavior manifested in specific group situations have been described, such as the formation of lanes in corridors, the "faster - slower" effect and arc-like blockages in narrow places [15, 13]. Social forces and their variants, agent-based models, and animal-based approaches are microscopic models that, through a variety of approaches, successfully shape pedestrian behavior. In pedestrian modeling, the ability to reproduce these phenomena, phenomena of collective behavior or self-organization is an indicator of the quality of the model.

2 Model and method

The analysis of pedestrian behavior and trajectories recorded by video cameras is one of the important topics of computer vision, widely studied for decades. In such studies, trajectory segmentation is often performed to reduce the cost of calculations and extract local information. There are three typical approaches [1]:

- Temporal segmentation: separation of the trajectory at points where two observed locations are temporarily separated from each other.
- Shape-based segmentation: Separation at points of greater curvature indicating that the target can change its direction at that point. This is used to simplify the shape of trajectories, and the Ramer Douglas-Pecker algorithm is a well-known approach [3].

- Semantic segmentation: dividing the entire trajectory into semantically significant segments, and a variety of methods have been proposed for different tasks [4,5,6,7,8].

We will focus on the third type, semantic segmentation, of trajectories based on patterns of human behavior (or agents). It would be very useful if we could have segments associated with behavior, however, no methods for segmenting trajectories have been proposed for the task of analyzing human behavior. Our proposed method first evaluates agent models using a mixed dynamic pedestrian agent (MDA) model, and then segments the trajectories using the studied agent models using hidden Markov models (HMM) [10], [11].

In this section, we will briefly describe the mixed model of dynamic pedestrian agents (MDA) proposed by Zhou et al. to study patterns of behavior or agents. MDA is a hierarchical Bayesian model that represents pedestrian trajectories using a mixed model of dynamics and beliefs. Using a modified Kalman filter that processes missing observations in the trajectories, and an iterative EM algorithm, the dynamics and confidence parameters of each agent are evaluated. Finally, the trajectories are grouped based on the evaluated agents.

In our proposed method, we use scoring agents to segment trajectories instead of clustering.

Let $y_t \in R^2$ — two-dimensional coordinates of the pedestrian at time t , and $x_t \in R^2$ — the corresponding state of the next linear dynamical system:

$$x_t \sim P(x_t | x_{t-1}) = N(x_t | Ax_{t-1} + b, Q) \quad (1)$$

$$y_t \sim P(y_t | x_t) = N(y_t | x_t, R) \quad (2)$$

Where is $N(\square)$ is a normal distribution with covariance matrices $Q, R \in R^{2 \times 2}$ and $A \in R^{2 \times 2}$ — is a state transition matrix and $b \in R^2$, assuming that the state transition is a similar transformation. In this article, we explicitly use a translation vector for such a conversion, while Zhou et la. In the work [9], homogeneous coordinates were used for their formulation.

MDA represents pedestrian trajectories as a mixture of dynamics D and belief B. Here's the dynamic $D = (A, b, Q, R)$ describes the dynamics of human movement in a two-dimensional scene. Belief B describes the starting point x_s and endpoint x_e trajectories, each of which is represented by normal distributions as follows:

$$x_s \sim p(x_s) = N(x_s | \mu_s, \Phi_s) \quad (3)$$

$$x_e \sim p(x_e) = N(x_e | \mu_e, \Phi_e) \quad (4)$$

That is, a belief is represented as $B = (\mu_s, \Phi_s, \mu_e, \Phi_e)$, describing where it begins and where it goes. Beca mixtures are recorded as $\pi_m = p(z = m)$, where the hidden variable z represents that the trajectory is generated by agent m . In addition, it is assumed that the observation $y = \{y_1, y_2, \dots, y_\tau\}$ length τ does not begin or end at the exact start and end points x_s and x_e agent, that is, trajectory statistics exist before and after the observed points of the trajectory:

$$x = \{x_s = x_{-t_s}, x_{-t_s+1}, \dots, x_0, x_1, x_2, \dots, x_\tau, x_{\tau+1}, \dots, x_{\tau+t_e} = x_e\} \quad (5)$$

Further $x_{1:T}$ denotes a sequence of x states except x_s and x_e . Figure 1 shows the MDA state space model.

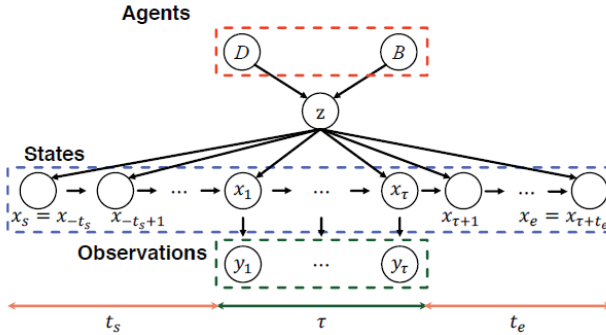


Fig.1. MDA State Space Model

For given K trajectories $Y = \{y^k\}$ MDA Evaluates M Agents $\Theta = \{(D_m, B_m, \pi_m)\}$, maximizing the following logarithmic probability:

$$L = \sum_k \log p(y^k | x^k, z^k, t_s^k, t_e^k, \Theta) \quad (6)$$

this can be rewritten by replacing the hidden variables $Z = \{z^k\}, T = \{(t_s^k, t_e^k)\}$ c $H = \{Z, T\}, h^k = \{z^k, t_s^k, t_e^k\}$ as follows:

$$L = \sum_k \log p(y^k | x^k, h^k, \Theta) \quad (7)$$

The EM algorithm evaluates iteratively because H is not observed.

$$Q(\Theta, \hat{\Theta}) = E_{x, h | y, \hat{\Theta}} [L] \quad (8)$$

Further $E_{x^k | y^k, h^k} [x^k] = \hat{x}^k$ denoted as \hat{x}^k , which is calculated using a modified Kalman filter [12]. The weights are listed as follows:

$$\gamma^k = \frac{p(h^k | \hat{x}_{1:T}^k, \hat{\Theta}) p(y^k | h^k, \hat{x}_{1:T}^k, \hat{\Theta}) p(\hat{x}_s^k, \hat{x}_e^k | \hat{\Theta})}{p(y^k | \hat{x}_{1:T}^k, \hat{\Theta}) p(\hat{x}_s^k, \hat{x}_e^k | \hat{\Theta})} \quad (9)$$

Note that we assume conditional independence between y^k and $\hat{x}_s^k; \hat{x}_e^k$, and between $\hat{x}_s^k, \hat{x}_e^k, h^k, \hat{x}_{1:T}^k$. Further assuming independence between hidden variables z, t_s, t_e and conditional independence between x and $\hat{\Theta}$, we have

$$p\left(h^k \mid x^k, \hat{\Theta}\right) = p\left(h^k\right) = p\left(z^k, t_s^k, t_e^k\right) = p\left(z^k\right) p\left(t_s^k\right) p\left(t_e^k\right) \quad (10)$$

Removing t_s, t_e assuming they are homogeneous, we have

$$\gamma^k = \frac{p\left(z^k\right) p\left(y^k \mid h^k, x_{1:T}^k, \hat{\Theta}\right) p\left(x_s\right) p\left(x_e\right)}{\sum_{h^k} p\left(z^k\right) p\left(y^k \mid h^k, x_{1:T}^k, \hat{\Theta}\right) p\left(x_s\right) p\left(x_e\right)} \quad (11)$$

Where is $p\left(y^k \mid h^k, x_{1:T}^k, \hat{\Theta}\right)$ is calculated using a modified Kalman filter that takes into account unobservable states $\left\{x_{-t_s}, x_{-t_s+1}, \dots, x_0, x_{t+1}, \dots, x_{t+t_e}\right\}$.

Next, we find $\hat{\Theta} = \arg \max_{\Theta} Q\left(\Theta, \hat{\Theta}\right)$ by solving a system of equations obtained by differentiating Q with respect to, which leads to the following analytical solutions: Θ

- 1) for each trajectory y^k , for all $h^k = \left(z^k, t_s^k, t_e^k\right)$ modified Kalman filter is used for evaluation $\left\{x^k\right\}$ and γ^k
- 2) renewal Θ .

The hidden Markov model (HMM) shown in Figure 2 has discrete hidden variables z .

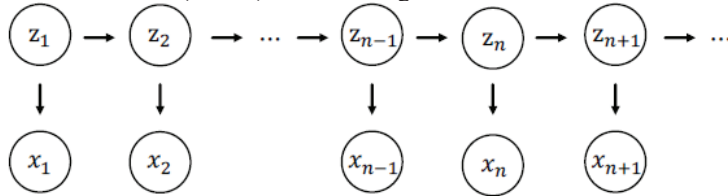


Fig. 2. HMM State Space Model

Using the Baum-Welch algorithm, HMM learns parameters from the training data, then outputs unobservable states. $Z = \left\{z_n\right\}_{n=1}^N$ from observations using the Viterbi algorithm [11]. A possible expansion of MDA to segmentation is shown in Figure 3, where it is assigned to each state using, where $X = \left\{x_n\right\}_{n=1}^N$ z_n x_n x_{n-1} x_n is the state of the Kalman filter and z_n is a hidden variable indicating which agent generates the observation. Both x_n and z_n depend on the previous variables x_{n-1} and z_{n-1} , as shown in Figure 3, which is known as Kalman filter switching [13], [14].

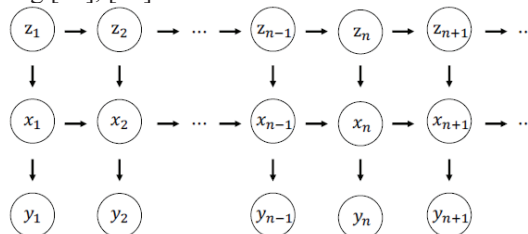


Fig. 3. Model of the state space of the Kalman switching filter

3 Research and results

A Kalman switching filter is a dynamic model of a system whose parameters depend on hidden variables. The state and observation at time x_n, y_n are specified by the formula:

$$x_n \sim P(x_n | x_{n-1}) = N(x_n | A_n x_{n-1}, Q_n) \tag{12}$$

$$y_n \sim P(y_n | x_n) = N(y_n | C_n x_n, R_n) \tag{13}$$

where A_n, C_n, Q_n, R_n are parameters that are switched by the value of the hidden variable. The Kalman switching filter is a useful model, but requires specifying the probabilities of the state transition, so it does not apply to the task presented here. Instead, we propose separating the MDA agent assessment from the HMM output to make the entire procedure work

$$A_n = A[z_n], C_n = C[z_n], Q_n = Q[z_n], R_n = R[z_n], z_n \tag{14}$$

Figure 4 shows an overview of the proposed method. First, we examine multi-agent trajectory models from video using MDA. Then we segment the trajectories using HMM based on the agents studied.

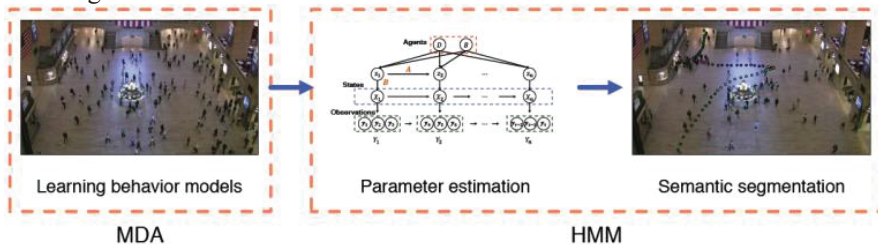


Fig. 4. Overview of the proposed method

Let M agents be, and $D_m = (A_m, b_m, Q_m, R_m)$ $B_m = (\mu_{s,m}, \Phi_{s,m}, \mu_{e,m}, \Phi_{e,m})$ M . Then all agents are denoted as

$$\Theta = \left\{ (D_m, B_m, \pi_m) \right\}_{m=1}^M = \left\{ W_m \right\}_{m=1}^M \tag{14}$$

Figure 5 shows the model of the proposed method.

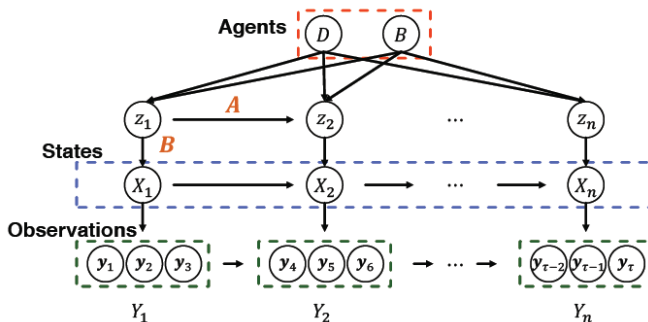


Fig. 5. Model of the state space of the proposed method

Agents pass from each other according to the state transition matrix A , and state X is generated based on the output probability matrix B . We use the Baum-Welch algorithm to estimate the initial probability distributions of the agents studied by M , as well as the matrix A and B . We assume that the agent can switch to another agent at each step, and the observation of one step consists of consecutive three coordinates along the trajectory, as shown $Y_t \in R^6$ $y_{t_1}, y_{t_2}, y_{t_3} \in R^2$ in Figure 5. The state is considered to be the generated agent specified by, associated with step $X_t z_t t$.

Therefore, the trajectory is represented by the hidden variables of, state and observation $Z = \{z_t\}_{t=1}^n$ $X = \{X_t\}_{t=1}^n$ $Y = \{Y_t\}_{t=1}^n$.

Here let be ρ an M -dimensional vector whose m -th element represents the initial distribution of the agent, and A through the MM matrix, the element of which is the probability of transition from agent to agent. $\rho_m \omega_m \times (i, j) a(i, j) \omega_j$.

It is assumed that the distribution of the output signal of the agent m is normal, and let $N(\mu_m, \sum_m) B$ be a vector whose m -th element is the output probability of the agent. $\omega_m \omega_j X_t$

$$b(j, t) \sim p(X_t | \omega_j) = N(X_t | \mu_j, \sum_j) \quad (15)$$

Denoting the parameters of the HMM to be estimated using, we maximize the following logarithmic probability for estimating the given K trajectories $\rho, A, B \Theta = (\rho, A, B)$:

$$Q(\Theta, \Theta^{old}) = \sum_K \sum_Z p(Z | X, \Theta^{old}) \ln p(X, Z, \Theta) \quad (16)$$

using the EM algorithm.

The trajectory is segmented using the Viterbi algorithm with the HMM parameters studied, that is, sequences of hidden variables and agents $\Theta Z^* \Omega^*$.

$$Z^* = \{i_1, i_2, \dots, i_n\} \quad (17)$$

$$\Omega^* = \{\omega_{i_1}, \omega_{i_2}, \dots, \omega_{i_n}\} \quad (18)$$

We compare the proposed method, denoted MDA +HMM in the following, with the Ramer-Douglas Peucker Algorithm (RDP) in terms of segmentation accuracy. Trajectories on the Footpath. For these experiments, a dataset is used [15]. This dataset contains a large number of pedestrian trajectories in videos of different sizes of 1920,1080 pixels. First, we evaluate methods with synthetic trajectories generated from the dataset to compare performance, then with the actual trajectories of the dataset.

The estimation metrics used in these experiments are positional error and step error defined in Algorithm 1.

Note that N_{est} and N_{gt} are the numbers of calculated and actual segmentation points in the trajectory.

To compare methods with a large number of trajectories, we generate 20,000 trajectories from MDA agent models extracted from the dataset. Assuming that the probabilities of the transition are the same, these trajectories are selected from a linear system of equations. (1) and (2). In the future, we use 10,000 trajectories for training HMM (parameter evaluation), and the remaining 10,000 trajectories for Logical inference

(segmentation). Segmentation points are those at which agent models switch from each other.

For the RDP method, we segment the trajectories by changing the parameter values, and then choose the best one. In this case, $\zeta = 69$ and $\zeta = 80$ minimize each error. For the proposed method, we select a different number of agents for segmentation (from 5 to 10) to evaluate the HMM parameters and output. Since MDA has studied 10 agents, we perform the same procedure 10 times (except when using all 10 agents) and then report average results. $\zeta = 5$ Table 1 shows the results of the comparison. The proposed method works better when the number of agents used is greater than eight.

Table 1. Results of experiments using synthetic trajectories

Method	Agent	Error Positioning	Step Error
MDA+HMM	5	44.99 \pm 6.25	2.35 \pm 0.17
	6	38.81 \pm 2.70	2.10 \pm 0.18
	7	33.76 \pm 3.01	1.81 \pm 0.13
	8	30.91 \pm 4.09	1.57 \pm 0.13
	9	25.59 \pm 3.34	1.32 \pm 0.10
	10	20.88	1.09
RDP	ζ		
	69	33.69	1.84
	80	34.17	1.82

To evaluate the methods using a real dataset, we selected and manually annotated 104 trajectories so that the trajectories were segmented at the point where pedestrians turn in their walking direction. For the RDP method, we segment all the trajectories by changing the parameter values, and then select the best one. In this case, $\zeta = 38$ and $\zeta = 29$ minimize each error. For the proposed method, we select a different number of agents for segmentation (from 5 to 10), as we did in the previous section.

To separate the dataset to evaluate the HMM parameters and the output, we perform a four-fold cross-validation. The results are shown in Table 2.

Table 2. Results

Method	Agent	Error Positioning	Step Error
MDA+HMM	5	62.12 \pm 5.95	2.92 \pm 0.27
	6	57.69 \pm 5.22	2.70 \pm 0.20
	7	56.41 \pm 5.73	2.61 \pm 0.25
	8	52.32 \pm 5.64	2.39 \pm 0.22
	9	53.69 \pm 3.53	2.46 \pm 0.15
	10	53.41	2.44
RDP	ζ		
	29	27.85	1.21
	38	26.61	1.24

4 Conclusion

As can be seen from the results obtained, the RDP is smaller than in the proposed method, but it does not provide any semantic information about segmentation. In contrast, the proposed method divides the trajectories into semantically significant segments with corresponding agent-based models, which helps to understand the behavior of pedestrians in real conditions.

References

- 1 A. Amer Mohammed Salih, M. Al-Khannaq, K. Hasikin, & N. Ashidi Mat Isa., *Adaptive local exposure based region determination for non-uniform illumination and low contrast images*. Alexandria Engineering Journal, **61**(12), 11185-11195. (2022) doi:10.1016/j.aej.2022.04.023
- 2 T. Li, Z. Zhan, & G. Tan, Accurate visual localization with semantic masking and attention. Eurasip Journal on Advances in Signal Processing, 2022(1) (2022) doi:10.1186/s13634-022-00875-2
- 3 C. Zheng, D. Cao, & C. Hu, *A similarity-guided segmentation model for garbage detection under road scene*. Frontiers of Optoelectronics, **15**(1) (2022) doi:10.1007/s12200-022-00004-9
- 4 W. Shi, Z. Huang, H. Huang, C. Hu, M. Chen, S. Yang, & H. Chen., *LOEN: Lensless optoelectronic neural network empowered machine vision*. Light: Science and Applications, **11**(1) (2022) doi:10.1038/s41377-022-00809-5
- 5 M. Rusanovsky, O. Beeri, & G. Oren, *An end-to-end computer vision methodology for quantitative metallography*. Scientific Reports, **12**(1) (2022) doi:10.1038/s41598-022-08651-w
- 6 Y. Bitjukov, Y. Deniskin, G. Deniskina, & I. Pocebneva, Application of wavelets and conformal reflections to finding optimal scheme of fiber placement at 3d printing constructions from composition materials. Paper presented at the E3S Web of Conferences, **244** (2021) doi:10.1051/e3sconf/202124405004
- 7 R. Ojala, J. Vepsäläinen, & K. Tammi, Motion detection and classification: Ultra-fast road user detection. Journal of Big Data, **9**(1) (2022) doi:10.1186/s40537-022-00581-8
- 8 M. A. Al-Malla, A. Jafar, & N. Ghneim., Image captioning model using attention and object features to mimic human image understanding. Journal of Big Data, **9**(1) (2022) doi:10.1186/s40537-022-00571-w
- 9 W. Yang, W. Liow, S. Chen, J. Yang, P. Chung, & S. Mao, Improved vehicle detection systems with double-layer LSTM modules. Eurasip Journal on Advances in Signal Processing, 2022(1) (2022) doi:10.1186/s13634-022-00839-6
- 10 S. Hao, X. Han, Y. Guo, & M. Wang., Decoupled low-light image enhancement. ACM Transactions on Multimedia Computing, Communications and Applications, **18**(4) (2022) doi:10.1145/3498341
- 11 M. Lupión, A. Polo-Rodríguez, J. Medina-Quero, J. F. Sanjuan, & P. M. Ortigosa, On the limits of conditional generative adversarial neural networks to reconstruct the identification of inhabitants from IoT low-resolution thermal sensors. Expert Systems with Applications, **203** (2022) doi:10.1016/j.eswa.2022.117356
- 12 H. Zhang, S. Zhang, Y. Zhang, J. Liang, & Z. Wang., Machining feature recognition based on a novel multi-task deep learning network. Robotics and Computer-Integrated Manufacturing, **77** (2022) doi:10.1016/j.rcim.2022.102369
- 13 A. R. Deniskina, I. V. Pocebneva, & A. V. Smolyaninov. Multidimensional object management. Paper presented at the Proceedings - 2021 International Russian Automation Conference, RusAutoCon 2021, 17-22. (2021) doi:10.1109/RusAutoCon52004.2021.9537333
- 14 A. Smolyaninov, I. Pocebneva, I. Fateeva, & K. Singur. Software implementation of a virtual laboratory bench for distance learning. Paper presented at the E3S Web of Conferences, **244** (2021) doi:10.1051/e3sconf/202124411009
- 15 M. Guo, T. Xu, J. Liu, Z. Liu, P. Jiang, T. Mu, S. Hu., Attention mechanisms in computer vision: A survey. Computational Visual Media, **8**(3), 331-368. (2022) doi:10.1007/s41095-022-0271-y
- 16 A. Smolyaninov, I. Pocebneva, I. Fateeva, & K. Singur. Software implementation of a virtual laboratory bench for distance learning. Paper presented at the E3S Web of Conferences, **244** (2021) doi:10.1051/e3sconf/202124411009
- 17 B. Li, L. Ye, J. Liang, Y. Wang, & J. Han., Region-of-interest and channel attention-based joint optimization of image compression and computer vision. Neurocomputing, **500**, 13-25. (2022) doi:10.1016/j.neucom.2022.05.047

- 18 B. Li, L. Ye, J. Liang, Y. Wang, & J. Han,. Region-of-interest and channel attention-based joint optimization of image compression and computer vision. *Neurocomputing*, 500, 13-25. (2022) doi:10.1016/j.neucom.2022.05.047
- 19 S. Cha, & Y. Wang, Zero-shot semantic segmentation via spatial and multi-scale aware visual class embedding. *Pattern Recognition Letters*, 158, 87-93. (2022) doi:10.1016/j.patrec.2022.04.011
- 20 M. Schellenberg, K. K. Dreher, N. Holzwarth, F. Isensee, A. Reinke, N. Schreck, J. Gröhl,. Semantic segmentation of multispectral photoacoustic images using deep learning. *Photoacoustics*, 26 (2022) doi:10.1016/j.pacs.2022.100341
- 21 F. Li, Y. Lu, X. Mao, J. Duan, & X. Liu. Multi-task deep learning model based on hierarchical relations of address elements for semantic address matching. *Neural Computing and Applications*, 34(11), 8919-8931. (2022) doi:10.1007/s00521-022-06914-1