



Title	Multi-Modal Representation Learning for Mapping Between Body Motion and Visual Imagery
Author(s)	Jülg, Tobias; Walter, Florian; Kim, Dongmin et al.
Citation	The 11th International Symposium on Adaptive Motion of Animals and Machines (AMAM2023). 2023, p. 150-151
Version Type	VoR
URL	<a href="https://doi.org/10.18910/92312">https://doi.org/10.18910/92312</a>
rights	
Note	

*The University of Osaka Institutional Knowledge Archive : OUKA*

<https://ir.library.osaka-u.ac.jp/>

The University of Osaka

# Multi-Modal Representation Learning for Mapping Between Body Motion and Visual Imagery

Tobias Jülg<sup>1,2</sup>, Florian Walter<sup>2</sup>, Dongmin Kim<sup>1</sup>, Hoshinori Kanazawa<sup>1</sup>, Alois Knoll<sup>2</sup>, Yasuo Kuniyoshi<sup>1</sup>  
<sup>1</sup>Graduate School of Information Science and Technology, The University of Tokyo, Japan  
<sup>2</sup>School of Computation, Information and Technology, Technical University of Munich, Germany  
 {juelg, walter, knoll}@cit.tum.de, {d-kim, kanazawa, kuniyosh}@isi.imi.i.u-tokyo.ac.jp

## 1 Introduction

Human infants need to learn autonomously to control their body in order to achieve higher level tasks. For this, it is crucial to have an internal representation of the body and environment states. An important requirement is that it should be multi-modal, to learn a bidirectional mapping between body states, represented as mental images, and related joint angles.

In this work, we model this mapping in a novel way by combining two basic tools from representation learning: We use a multi-modal self-organizing map (SOM) [1] to learn an internal representation for vision and proprioception of a simulated infant. As we will show, adding an autoencoder (AE) represents the visual stimuli in a low-dimensional latent space, reduces the computational complexity and increases the quality of the learned representation.

## 2 Method

### 2.1 Dataset

We have created a dataset with  $N = 20000$  samples generated with the infant simulation from [2]. It contains data from two modalities for training the AE and the multi-modal SOM. Each sample in the dataset corresponds to an arm pose of the simulated infant and contains data from vision (modality  $m_1$ ) and proprioception (modality  $m_2$ ):

$$\mathcal{D} = \{(d_i^{(m_1)}, d_i^{(m_2)})\}_{i=1}^N \quad (1)$$

The experiment setup and a sample can be seen in Figure 1. For the vision modality, we use a camera which is attached to the position of the infant’s right eye. It can take wide angle RGB images of 300 by 500 pixels. For the proprioception modality the five joint angles of the right arm are obtained.

The poses used in the dataset were sampled randomly in joint space. Samples where the arm was not visible were discarded. For training we used a 60/20/20 split.

### 2.2 Multi-Modal SOM

The multi-modal SOM for learning a representation of the infants’s arm pose builds upon the Multi-Modal Convergence Maps framework from [1].

We use  $x$  and  $y$  to refer to the map dimensions and  $m$  to refer to the modality dimension. We further define the weights of the map  $w_{x,y}^m$  for the neuron  $x, y \in \mathbb{N}$  and modality  $m \in \mathcal{M}$ .

In order to train the map for a given training sample  $d_i \in \mathcal{D}$ ,

The 11th International Symposium

on Adaptive Motion of Animals and Machines (AMAM2023)

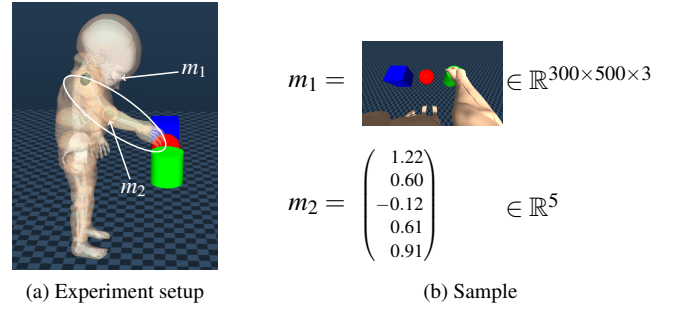


Figure 1: (a) shows the experiment setup in the MuJoCo simulation with the infant performing a valid pose and (b) shows the corresponding dataset sample.

one first needs to calculate the neuron activity per modality as the weighted MAE between the current learned neuron representation and the sample. This is done for all modalities separately:

$$a_{x,y}^{(m)} = -\frac{c^{(m)}}{n^{(m)}} \cdot \sum_{j=1}^{n^{(m)}} |w_{x,y}^{(m)} - d_j^{(m)}| \quad \forall x, y \forall m \in \mathcal{M} \quad (2)$$

$c^{(m)}$  is a factor which attributes each modality a certain influence during training and  $n^{(m)}$  is the flattened dimensionality of modality  $m$ . Afterwards, one calculates the combined neuron activity over all modalities:

$$A_{x,y} = \frac{\sum_{m \in \mathcal{M}} a_{x,y}^{(m)}}{\sum_{m \in \mathcal{M}} c^{(m)}} \quad (3)$$

The neuron with the highest activity, i.e. the one that most closely represents the training sample, is picked as the winning neuron:  $x_{win}, y_{win}$ . The weight update is then calculated by the error, a learning rate  $\lambda$  and an exponentially decaying neighborhood function  $\mu : \mathbb{N}^4 \rightarrow \mathbb{R}_{\geq 0}$  which has the highest value at the winning neuron:

$$dw_{x,y}^{(m)} = (w_{x,y}^{(m)} - d_i^{(m)}) \cdot \lambda \cdot \mu(x, y, x_{win}, y_{win}) \quad (4)$$

We extended the weight update proposed in [1] to support mini batches by averaging over individual updates:

$$dw_{x,y}^{(m)} = \frac{1}{|B_i|} \sum_{s \in B_i} dw_{x,y,s}^{(m)} \quad (5)$$

where  $B_i$  is the  $i$ th batch and  $dw_{x,y,s}^{(m)}$  is the update calculated by the  $s$ th item of this batch for modality  $m$ .

After training a multi-modal SOM, it is possible to map between its modalities. Given, for example, a sample  $s^{(m_2)}$  where only  $m_2$  is known, then  $\hat{s}^{(m_1)}$  can be predicted by simply calculating the activity for  $s^{(m_2)}$  and returning the  $m_1$  component of the winning neuron:  $w_{x_{win},y_{win}}^{(m_1)}$ .

### 2.3 Architecture

Our SOM with pre-trained AE (SOM-AE) architecture is summarized in Figure 2. Visual input is represented in a latent space of dimension 8 by an AE that was trained in advance. It consists out of an encoder  $e : \mathbb{R}^{64 \times 64} \rightarrow \mathbb{R}^8$  (dimensions result from pre-processing) and a decoder  $d : \mathbb{R}^8 \rightarrow \mathbb{R}^{64 \times 64}$  which are trained on the MSE reconstruction loss. Thus, the training input of the SOM for a sample  $d_i^{(m_1)}$  becomes  $l_i^{(m_1)} = e(d_i^{(m_1)})$ . This reduces the map size from 10.5GB to only 33.3MB. Note, that it is possible to cache the encodings offline, thus avoiding a new bottleneck.

We argue that this combination of SOMs and AEs is biologically plausible, as it can be compared to the hierarchical structure of visual processing in the human brain. The AE uses a convolutional neural network architecture that has been argued to share similarities with the visual cortex [3]. Thus, the latent space can be compared to the highest level of abstraction which is fed to multisensory areas which are modeled by the multi-modal SOM.

Combining SOMs and AEs has been studied in a few ways before. The approach which might come closest to ours is Deep Neural Maps [4] which learn both AE and SOM simultaneously, resulting in improved embeddings on unimodal image focused benchmark datasets. In contrast, we first train the AE separately and use the learned latent space to train a biologically plausible multi-modal SOM [1]. We also use a complex dataset composed of different modalities of arm postures from a simulated embodied infant.

## 3 Results and Conclusion

The performance comparison between the plain SOM and SOM-AE can be seen in Table 1. MAE losses and standard deviations for four mapping scenarios over the test set are provided in the lower part. In the case of SOM-AE, for the vision modality, the decoded output  $d(w_{x_{win},y_{win}}^{(m_1)})$  is used for the er-

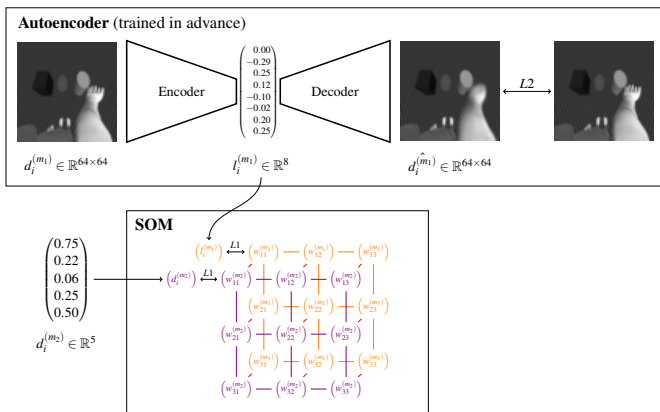


Figure 2: SOM-AE architecture, the visual modality  $m_1$  is first used to train an AE with a latent space of size 8. Afterwards the latent space is used for the multi-modal SOM training.

ror calculation with the sample’s image  $d_i^{(m_1)}$ . SOM-AE uses much less computational resources for training, even when the training time of the AE is also included. The performance does not significantly reduce and even increases in some cases. SOM-AE also performs better in terms of learned map quality, as the randomly chosen validation samples in Table 2 suggest.

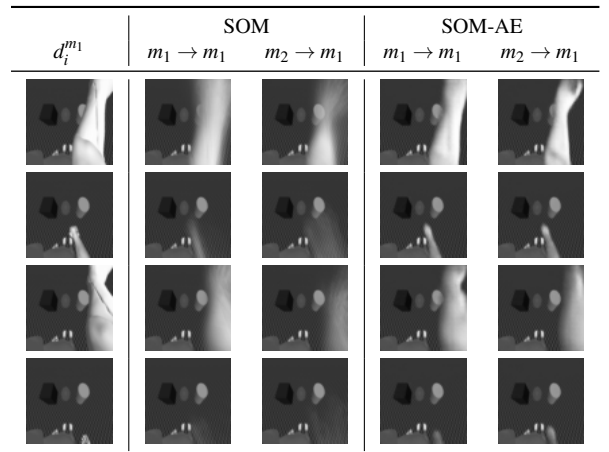
Given the trained SOM-AE, it is possible to let an agent, which embodies the infant, “imagine” a situation and let it “remember” how to get to it by mapping between modalities, as described earlier.

In conclusion, SOM-AE not only significantly reduces the computational requirements to train a multi-modal SOM model, but also improves the qualitative mapping results while keeping the quantitative loss at the same level. Therefore, this framework is a step towards better multi-modal representations and bidirectional modality mapping. As argued before, it can partly be mapped to cognitive perception processes in the human brain and it, thus, models the emergence of interconnection between body motion and visual imagery in an infant’s developmental process.

Table 1: Quantitative performance comparison

	SOM	SOM-AE
Batch size	1	<b>10</b>
Epochs	200	200
SOM training time	61.5 h	<b>15.5 min</b>
AE training time	<b>0</b>	9.5 h
GPU	A100 SXM4	<b>RTX 2080</b>
VRAM used	42GB	<b>2.6GB</b>
Overall test loss	0.068±0.077	<b>0.058±0.065</b>
Test loss: $m_1 \rightarrow m_1$	0.025±0.020	<b>0.021±0.020</b>
Test loss: $m_1 \rightarrow m_2$	0.163±0.096	<b>0.127±0.087</b>
Test loss: $m_2 \rightarrow m_1$	0.046±0.038	<b>0.041±0.041</b>
Test loss: $m_2 \rightarrow m_2$	<b>0.039±0.021</b>	0.042±0.023

Table 2: Qualitative inference comparison



## References

- [1] S. Lallee and P. F. Dominey, “Multi-modal convergence maps: from body schema and self-representation to mental imagery,” *Adaptive Behavior*, vol. 21, no. 4, pp. 274–285, Aug. 2013.
- [2] D. Kim, H. Kanazawa, and Y. Kuniyoshi, “Simulating a human fetus in soft uterus,” in *2022 IEEE International Conference on Development and Learning (ICDL)*, 2022, pp. 135–141.
- [3] D. Hassabis, D. Kumaran, C. Summerfield, and M. Botvinick, “Neuroscience-inspired artificial intelligence,” *Neuron*, vol. 95, no. 2, pp. 245–258, 2017.
- [4] M. Pesteie, P. Abolmaesumi, and R. Rohling, “Deep neural maps,” 2018. [Online]. Available: <https://openreview.net/forum?id=HyG76D1wf>