# Adapting Tensor Kernel Machines to Enable Efficient Transfer Learning for Seizure Detection

S. J. S. de Rooij, *Student Member, IEEE*, and B. Hunyadi, *Senior Member, IEEE*

arXiv:2512.02626v1 [cs.LG] 2 Dec 2025

*Abstract*—**Transfer learning aims to optimize performance in a target task by learning from a related source problem. In this work, we propose an efficient transfer learning method using a tensor kernel machine. Our method takes inspiration from the adaptive SVM and hence transfers 'knowledge' from the source to the 'adapted' model via regularization. The main advantage of using tensor kernel machines is that they leverage low-rank tensor networks to learn a compact non-linear model in the primal domain. This allows for a more efficient adaptation without adding more parameters to the model. To demonstrate the effectiveness of our approach, we apply the adaptive tensor kernel machine (Adapt-TKM) to seizure detection on behind-the-ear EEG. By personalizing patient-independent models with a small amount of patient-specific data, the patient-adapted model (which utilizes the Adapt-TKM), achieves better performance compared to the patient-independent and fully patient-specific models. Notably, it is able to do so while requiring around 100 times fewer parameters than the adaptive SVM model, leading to a correspondingly faster inference speed. This makes the Adapt-TKM especially useful for resource-constrained wearable devices.**

*Index Terms*—**transfer learning, tensors, kernel machines, seizure detection**

## I. INTRODUCTION

INSPIRED by how humans use prior information to learn to do new tasks, transfer learning aims to improve performance in a *target* task by leveraging information from a related, or *source*, problem [1]–[3]. This way, the performance of machine learning models can be improved for target tasks for which there is little to no 'information' (i.e. *(supervised) data*) available.

Such situations arise especially in healthcare applications, where both data collection and labelling, which needs to be done by medical experts, are costly [4]. For this reason, generalized (or *patient-independent*) models are easier to implement, even though there is significant inter-subject variability with respect to disease processes [5]. Transfer learning can improve the performance of these generalized models through *personalization*, using only a small amount of patient-specific data to adapt the patient-independent model [6], [7].

S.J.S. de Rooij and B. Hunyadi are with the Signal Processing Systems group in the Microelectronics department at the Faculty of Electrical Engineering, Mathematics and Computer Science (EEMCS), Delft University of Technology, The Netherlands.

S.J.S. de Rooij, and thereby this work, is supported by the TU Delft AI initiative.

In this work, we focus especially on the problem of electroencephalogram (EEG)-based epileptic seizure detection. Epilepsy is a highly heterogeneous disease, where the seizure patterns observed on the EEG vary substantially between patients, depending on the subtype and origin of seizures in the brain [8]. Transfer learning (TL) has been applied previously to improve seizure detection and (the related) epilepsy detection [9].

Most of the work using TL for seizure and epilepsy detection considers the semi-supervised or *transductive* TL case. These models consider there to be no labels available for the target dataset and only use *supervision* for the source data. Thus, they typically aim to minimize the distance between the distribution of the source and target data.

In this manner, in [10], the large-margin-projected transductive SVM (LMPROJ) is used to improve epilepsy detection. The LMPROJ model minimizes the maximum mean discrepancy (MMD) between source and target data during training. Inspired by this work, the authors of [11] developed a semi-supervised learning algorithm using the TSK fuzzy system for model interpretability. Instead of minimizing the distance during the training step, the authors of [12] use transfer component analysis to learn a feature latent subspace where distances between patients are minimized. The classifier is then trained in this feature subspace.

A concern for the implementability of these models may be that the data from the source patients needs to be available for training the TL model, which could raise privacy concerns. Therefore, in [13], the authors develop a privacy-preserving seizure subtype classification model which uses source-free domain adaptation. This way, data from the source patients do not need to be saved.

The downside of these transductive TL (or semi-supervised) models is that they are only able to minimize the distance between the marginal distributions of the data. However, there can also be conditional distribution differences in epileptic EEG signals, i.e. differences in the distributions conditioned on the labels [9]. Methods to account for these differences require labels in the target domain (i.e. *inductive* TL).

In [14], they implemented a max-pooling CNN (MPCNN) to perform end-to-end learning for seizure detection. They found that fine-tuning the patient-independent model with patient-specific data outperformed both the patient-independent (PI) and the patient-specific (PS) model. Their approach is similar to our previous work, where we used fine-tuning on tensor kernel machines to update a PI model with PS data [15].

However, in that case, the fine-tuned model did not outperform the PS model.

Alternatively, in [7] the authors use the adaptive SVM model to personalize a heart-rate based seizure detection algorithm. This adaptive SVM model [16] uses the distance to the source model (in this case, the PI model) as a regularizer to train a new model for the target. This way, an *adapted* model is learned that inherits the support vectors from the *source* and adds new ones from the *target* data. The downside of this is that it increases the model size, which can be especially problematic if the PI model is already large (in the case of a large PI training set).

Moreover, recently, there has been an increase in interest in the development of efficient machine learning models [17]. Both from a sustainability perspective, as well as for application in edge devices [18] such as wearables for health monitoring. To this end, low-rank tensor networks have been proposed as a valuable asset to enable this 'Green AI' [19]. These tensor networks allow for a parameter-efficient expression of the models' parameters by using low-rank approximations.

Therefore, in this paper, we propose an *efficient* transfer learning method using tensor kernel machines [20], [21] based on the adaptive SVM model [16]. To show the effectiveness of the method, we apply it to the detection of epileptic seizures on behind-the-ear EEG [22]. Wearable behind-the-ear EEG devices are being developed to provide continuous, noninvasive monitoring for people with epilepsy [23]–[25]. They aim to improve quality of life by enabling better tracking of seizures and providing a warning of seizure onset [26]–[28].

The adaptability of our model makes it possible to improve detection performance via personalization, even when limited patient-specific data is available. Furthermore, the model's efficiency makes it especially useful for these resource-constrained devices [26]. It may allow for on-device inference as well as adaptation, which could reduce (or even eliminate) the need for data broadcasting.

The rest of this paper is structured as follows. First, in Section II background information is provided on kernel machines and tensor networks. Then, in the following section (III), we introduce the adaptive tensor kernel machine (Adapt-TKM). After which, we describe the seizure detection pipeline (Section IV). Results for the Adapt-TKM are provided in Section V. These results are discussed in Section VI, and some suggestions are made for future work. Finally, Section VII concludes this paper.

## II. KERNEL MACHINES AND TENSORS

### A. Notation and Preliminaries

Throughout the paper, the following notation conventions are used. Vector and matrices are denoted by boldface lowercase ($a \in \mathbb{R}^I$) and boldfase uppercase ($A \in \mathbb{R}^{I \times J}$) letters, respectively. Tensors, i.e. *multidimensional arrays* [29], are denoted by boldface calligraphic letters ($\mathcal{A} \in \mathbb{R}^{I_1 \times I_2 \times \cdots \times I_D}$). The symbols $\otimes$, $\odot$ and $\circledast$ are used to denote the tensor outer product, the (column-wise) Khatri-Rao product and the Hadamard product, respectively.

The Frobenius inner product between two tensors is defined as, $\langle \mathcal{A}, \mathcal{B} \rangle_F := \text{vec}(\mathcal{A})^T \text{vec}(\mathcal{B})$, where $\text{vec}(\mathcal{A})_i = a_{i_1 i_2 \ldots i_D}$

is the vectorization of $\mathcal{A} \in \mathbb{R}^{I_1 \times \cdots \times I_D}$ to $a \in \mathbb{R}^I$ where $I = \prod_{d=1}^{D} I_d$. When it comes to vectorization, we use the little-endian ordering convention, as it is most frequently used in tensor network literature [30].

A $D$-dimensional tensor $\mathcal{A} \in \mathbb{R}^{I_1 \times I_2 \times \cdots \times I_D}$ is considered to be *rank-one* if it can be written as the outer product of $D$ vectors [29],

$$\mathcal{A} = a^{(1)} \otimes a^{(2)} \otimes \cdots \otimes a^{(D)} = \bigotimes_{d=1}^{D} a^{(d)}.$$

Because the number of elements in a tensor grows exponentially with the dimension, low-rank tensor networks are typically used to represent the tensor with a reduced number of parameters [30]. The two tensor networks most often used for this purpose are the tensor-train decomposition (TT) [31] and the canonical polyadic decomposition (CPD) [32], [33].

A rank-$R$ CPD decomposes a tensor $\mathcal{A} \in \mathbb{R}^{I_1 \times I_2 \times \cdots \times I_D}$ into a sum of $R$ rank-one tensors,

$$\mathcal{A} := \sum_{r=1}^{R} \gamma_r \bigotimes_{d=1}^{D} a_r^{(d)}, \tag{1}$$

where the vectors $a_r^{(d)}$ are normalized to unit norm and the scaling is absorbed by $\gamma_r$.

The CPD can also be expressed as the Khatri-Rao product of *factor matrices*,

$$\text{vec}(\mathcal{A}) := \left( A^{(D)} \odot A^{(D-1)} \odot \cdots \odot A^{(1)} \right) \gamma .$$

Here, $\gamma \in \mathbb{R}^R$ and the columns of the factor matrices contain the vectors of the rank-one components, i.e. $A^{(d)} := [a_1^{(d)} \ a_2^{(d)} \ \cdots \ a_R^{(d)}] \in \mathbb{R}^{I_d \times R}$.

### B. Supervised Learning and Kernel Machines

In supervised learning the aim is to find a function $f(\cdot)$ that maps the input data $x \in \mathcal{X}$ to its output $y \in \mathcal{Y}$, given a dataset of $N$ observations of input-output pairs $\{x_n, y_n\}_{n=1}^{N}$. Kernel machines perform nonlinear regression ($y \in \mathbb{R}$) or classification ($y \in \{-1, 1\}$) by mapping the input data into a higher dimensional feature space via a feature map $\Phi(\cdot) : \mathcal{X} \to \mathcal{H}$ [34]. In this feature space, they search for a linear relationship between the features. Kernel machines thus have the following model form[1]:

$$f(x) = \langle \Phi(x), w \rangle \tag{2}$$

The weights of the model, $w$, can be obtained by minimizing the (regularized) empirical risk,

$$w = \underset{w}{\text{argmin}} \left[ \frac{1}{N} \sum_{n=1}^{N} \ell \left( \langle \Phi(x_n), w \rangle, y_n \right) + \lambda \|w\|_p^2 \right] \tag{3}$$

where $\ell(\cdot) : \mathcal{X} \times \mathcal{Y} \to \mathbb{R}^+$ is a convex loss function. Different loss functions lead to the primal formulation of different kernel machines. The perhaps most well-known kernel machine, the

---

[1]Note that we omit a bias term here. This is done to simplify notation and because a bias term is not necessarily needed when the data is properly scaled [35]. However, a bias term can always be added if needed.

SVM [36], uses the hinge loss function in its primal. Kernel ridge regression (or LS-SVMs) [37], on the other hand, use the squared loss function.

Due to the high dimensionality induced by the feature maps, most kernel machines are optimized in the dual (which can be derived using Lagrange multipliers) where the feature maps only appear in the model via their inner products. These inner products can be computed with kernel functions that operate in a Hilbert space: $\kappa(\boldsymbol{x}_i, \boldsymbol{x}_j) = \langle \Phi(\boldsymbol{x}_i), \Phi(\boldsymbol{x}_j) \rangle$ [34]. Thus, eliminating the need for explicit high-dimensional feature maps.

The downside of this dual formulation is that the size of the kernel matrix ($\boldsymbol{K}_{i,j} = \kappa(\boldsymbol{x}_i, \boldsymbol{x}_j)$), which contains the output of the kernel function for each combination of input samples, scales exponentially which the number of input samples ($\mathcal{O}\left(N^2\right)$). This means that to solve the dual and find the $\mathcal{O}\left(N\right)$ dual coefficients ($\hat{\alpha}_n$) an optimization procedure is needed that is at least $\mathcal{O}\left(N^2\right)$ [35]. Furthermore, for model inference (4), it is necessary to store and use (most of) the training set ($\mathcal{O}\left(ND\right)$) as the *support vectors*. All of this could be prohibitively expensive for large enough $N$.

$$f(x) = \sum_{n=1}^{N} \hat{\alpha}_n y_n \kappa(\boldsymbol{x}, \boldsymbol{x}_n) \qquad (4)$$

Therefore, the tensor kernel machines we consider in this paper[2] [20], [21] are optimized in the primal, and the challenges that arise from the high dimensionality of the feature space are tackled using low-rank tensor networks. Thus, we consider feature maps that are the outer product of $D$ local feature maps, $\phi^{(d)}(\cdot) : \mathbb{R} \to \mathbb{R}^{M_d}$,

$$\Phi(\boldsymbol{x}) = \bigotimes_{d=1}^{D} \phi^{(d)}(x^{(d)}), \qquad (5)$$

which results in a rank-1 tensor $\Phi(\boldsymbol{x}) \in \mathbb{R}^{M_1 \times M_2 \times \cdots \times M_D}$ for input data $\boldsymbol{x} \in \mathbb{R}^D$. These tensor feature maps can approximate a whole range of functions by using pure power polynomials or Fourier features as the local feature maps [39], [40]. One such function is the frequently used RBF kernel (Section V-A) [20], [41].

Using these tensor feature maps, the weights of the kernel machine can be expressed as a tensor, $\mathcal{W} \in \mathbb{R}^{M_1 \times M_1 \times \cdots \times M_D}$, leading to the primal of the tensor kernel machine (TKM) with rank-1 feature maps [20], [21],

$$\min_{\mathcal{W}} \frac{1}{N} \sum_{n=1}^{N} \ell\left(\langle \Phi(\boldsymbol{x}_n), \mathcal{W} \rangle_{\mathrm{F}}, y_n\right) + \lambda \left\| \mathcal{W} \right\|_{\mathrm{F}}^2$$
$$\text{s.t. } \Phi(\boldsymbol{x}) = \bigotimes_{d=1}^{D} \phi^{(d)}(x^{(d)}) . \qquad (6)$$

This weight tensor scales exponentially with $D$ (the input data size). Therefore, in TKMs the weight tensor is typically approximated using a low-rank tensor network, such as a CPD or tensor-train. These tensor networks allow for linear [29],

[31], or even log-linear [39], [42], scaling with the input data size.

Throughout, this paper we will use the CPD for the weights, since this decomposition has fewer parameters than the tensor-train decomposition and is thus more efficient to learn. Furthermore, we will use the squared error as loss function, this leads to tensor kernel ridge regression (TKRR) [20],

$$\min_{\mathcal{W}} \frac{1}{N} \sum_{n=1}^{N} \left(\langle \Phi(\boldsymbol{x}_n), \mathcal{W} \rangle_{\mathrm{F}} - y_n\right)^2 + \lambda \left\| \mathcal{W} \right\|_{\mathrm{F}}^2$$
$$\text{s.t. } \Phi(\boldsymbol{x}) = \bigotimes_{d=1}^{D} \phi^{(d)}(x^{(d)}), \ \mathcal{W} = \sum_{r=1}^{R} \gamma_r \bigotimes_{d=1}^{D} \boldsymbol{w}_r^{(d)} \qquad (7)$$

This can be efficiently solved using an block coordinate descent procedure, where the factor matrices of the CPD weight tensor are updated alternatingly. The update of this block coordinate descent is done by solving the following least-squares problem,

$$\min_{\mathrm{vec}(\boldsymbol{w}^{(d)})} \frac{1}{N} \sum_{n=1}^{N} \left( \left\langle \boldsymbol{g}^{(d)}(\boldsymbol{x}_n), \mathrm{vec}(\boldsymbol{W}^{(d)}) \right\rangle_{\mathrm{F}} - y_n \right)^2$$
$$+ \lambda \left\langle \mathrm{vec}\left(\boldsymbol{W}^{(d)^T} \boldsymbol{W}^{(d)}\right), \mathrm{vec}\left(\boldsymbol{H}^{(d)}\right) \right\rangle_{\mathrm{F}}, \qquad (8)$$

where $\mathrm{vec}\left(\boldsymbol{W}^{(d)}\right)$ is the vectorization of the $d$-th factor matrix and,

$$\boldsymbol{g^{(d)}}(\boldsymbol{x}) = \mathrm{vec}\left( \phi^{(d)}(\boldsymbol{x}) \left( \underset{i=1, i \neq d}{\overset{D}{\circledast}} \phi^{(i)}(\boldsymbol{x})^{\mathrm{T}} \boldsymbol{W}^{(i)} \right)^{T} \right) \qquad (9)$$

$$\boldsymbol{H}^{(d)} = \left( \underset{i=1, i \neq d}{\overset{D}{\circledast}} \boldsymbol{W}^{(i)^T} \boldsymbol{W}^{(i)} \right)^{T} . \qquad (10)$$

## C. Transfer Learning and the Adaptive SVM

In transfer learning knowledge from a related domain (*source domain*) is used to improve the performance in a target domain [1]. There are many different categories of transfer learning [43], in this paper we focus on the case of inductive transfer learning. Inductive refers to the availability of labels in both the source and target domain.

In [16] the authors developed an inductive transfer learning method to adapt an existing SVM model (the *source*) to new data (*target*) called Adapt-SVM. In the Adapt-SVM model, the regularization of the weights for the model trained on target data, $\|\boldsymbol{w}\|_{\mathrm{F}}^2$, is replaced by the squared distance between the source and target (or *adapted*) weights, $\|\boldsymbol{w}^s - \boldsymbol{w}\|_{\mathrm{F}}^2$. This way, the source model $f^s(x)$ can be considered as a 'prior' for the adapted model $f^a(x)$. The resulting model will have additional *support vectors* from the target set,

$$f^a(x) = f^s(x) + \sum_{n=1}^{N^t} \hat{\alpha}_n y_n \kappa(x, x_n).$$

Thus, the model size of the adapted model is even larger than the source model ($\mathcal{O}\left((N^s + N^t)D\right)$ instead of $\mathcal{O}\left(N^sD\right)$). This Adapt-SVM method has been used for various applications, including video concept detection [44] and seizure detection using heart rate measurements [7].

---

[2]There also exist support tensor machines [38], which take tensorial data as input. Note that these are different from the tensor kernel machines we discuss in this paper.

## III. ADAPTIVE TENSOR KERNEL MACHINE

In this paper, we propose an adaptive tensor kernel machine (Adapt-TKM) that takes inspiration from the Adapt-SVM model [16]. Thus, we change the regularization of the TKM (6) to the squared Frobenius distance between the source and the *adapted* weights,

$$\min_{\mathcal{W}} \ \frac{1}{N^t} \sum_{n=1}^{N^t} \ell \left( \langle \boldsymbol{\Phi}(\boldsymbol{x}_n^t), \mathcal{W} \rangle_{\mathrm{F}}, y_n^t \right) + \mu \left\| \mathcal{W} - \mathcal{W}^s \right\|_F^2, \tag{11}$$

where $\mathcal{W}^s$ are the weights of the source model, $\mathcal{W}$ the weights of the adapted model, and with $\{\boldsymbol{x}_n^t, y_n^t\}_{n=1}^{N^t}$ as the target dataset. The regularization parameter $\mu$ controls the distance between the source and adapted model weights.

This objective (11) can be solved using block coordinate descent if $\mathcal{W}$ is a low-rank tensor-network. To derive the update step, we first expand the norm,

$$\|\mathcal{W} - \mathcal{W}^s\|_{\mathrm{F}}^2 = \langle \mathcal{W} - \mathcal{W}^s, \mathcal{W} - \mathcal{W}^s \rangle_{\mathrm{F}}$$
$$= \langle \mathcal{W}, \mathcal{W} \rangle_{\mathrm{F}} - 2\langle \mathcal{W}, \mathcal{W}^s \rangle_{\mathrm{F}} + \langle \mathcal{W}^s, \mathcal{W}^s \rangle_{\mathrm{F}}. \tag{12}$$

From this expansion, we see that $\langle \mathcal{W}, \mathcal{W} \rangle_{\mathrm{F}}$ is the same as the regularization term from the 'regular' TKRR (7), and $\langle \mathcal{W}^s, \mathcal{W}^s \rangle_{\mathrm{F}}$ is constant with respect to $\mathcal{W}$ and therefore has no influence on the optimal value of $\mathcal{W}$. Therefore, the only 'new' term in the optimization procedure is $\langle \mathcal{W}, \mathcal{W}^s \rangle_{\mathrm{F}}$.

If both $\mathcal{W}^s$ and $\mathcal{W}$ are CPD tensors, we can separate out the $d$-th factor matrix of $\mathcal{W}$ for the update step. First, to simplify notation we define $\mathcal{S} = \mathcal{W}^s$ with CPD-rank $P$. The dot product $\langle \mathcal{W}, \mathcal{S} \rangle_{\mathrm{F}}$ can then be rewritten as follows,

$$\langle \mathcal{W}, \mathcal{S} \rangle_{\mathrm{F}} = \left\langle \sum_{r=1}^{R} \bigotimes_{d=1}^{D} \boldsymbol{w}_r^{(d)}, \sum_{p=1}^{P} \bigotimes_{d=1}^{D} \boldsymbol{s}_r^{(d)} \right\rangle_{\mathrm{F}}$$

$$= \sum_{m_1=1}^{\hat{M}} \cdots \sum_{m_D=1}^{\hat{M}} \sum_{r=1}^{R} \sum_{p=1}^{P} w_{m_1 r}^{(1)} s_{m_1,p}^{(1)} \cdots w_{m_D,r}^{(D)} s_{m_D,p}^{(D)}$$

$$= \mathrm{vec}\left(\boldsymbol{W}^{(d)}\right)^{\mathrm{T}} \mathrm{vec}\left( \boldsymbol{S}^{(d)} \left( \underset{i=1, i \neq d}{\overset{D}{\circledast}} \boldsymbol{W}^{(i)\mathrm{T}} \boldsymbol{S}^{(i)} \right)^{\mathrm{T}} \right)$$

$$= \left\langle \mathrm{vec}\left(\boldsymbol{W}^{(d)}\right), \mathrm{vec}\left(\boldsymbol{Q}^{(d)}\right) \right\rangle_{\mathrm{F}}. \tag{13}$$

Combining this into the block coordinate descent step of the TKRR model (14) we obtain the following update step for the Adapt-TKRR model,

$$\min_{\mathrm{vec}(\boldsymbol{W}^{(d)})} \ \frac{1}{N} \sum_{n=1}^{N} \left( \left\langle \boldsymbol{g}^{(d)}(\boldsymbol{x}_n), \mathrm{vec}(\boldsymbol{W}^{(d)}) \right\rangle_{\mathrm{F}} - y_n \right)^2$$
$$+ \mu \left\langle \mathrm{vec}\left(\boldsymbol{W}^{(d)\mathrm{T}} \boldsymbol{W}^{(d)}\right), \mathrm{vec}\left(\boldsymbol{H}^{(d)}\right) \right\rangle_{\mathrm{F}} \tag{14}$$
$$- 2\mu \left\langle \mathrm{vec}\left(\boldsymbol{W}^{(d)}\right), \mathrm{vec}\left(\boldsymbol{Q}^{(d)}\right) \right\rangle_{\mathrm{F}},$$

where $\boldsymbol{g}^{(d)}(\boldsymbol{x})$, $\boldsymbol{H}^{(d)}$ and $\boldsymbol{Q}^{(d)}$ are defined as in (9), (10) and (13). Algorithm 1 shows the full Adapt-TKRR algorithm.

---

**Algorithm 1** Adapt-TKM

**Input:** Training data: $\{\boldsymbol{x}_n, y_n\}_{n=1}^{N}$.
Regularization parameter, $\mu$, maximum number of iterations, $n_{\max}$, feature map, $\boldsymbol{\Phi}$, weights of the source model, $\{\boldsymbol{S}^{(d)}\}_{d=1}^{D}$.

1: **Initialize:** Randomly initialize of the CPD weights, $\{\boldsymbol{W}_0^{(d)}\}_{d=1}^{D}$, or set it equal to the source weights $\{\boldsymbol{W}_0^{(d)}\}_{d=1}^{D} = \{\boldsymbol{S}^{(d)}\}_{d=1}^{D}$
2: **repeat**
3:     **for** $d = 1 : D$ **do**
4:         Compute $\forall n \ \boldsymbol{g}^{(d)}(\boldsymbol{x}_n)$ as defined in (9).
5:         Set the regularization matrices $\boldsymbol{H}^{(d)}$ and $\boldsymbol{Q}^{(d)}$ as in (10) and (13).
6:         Set $\boldsymbol{W}^{(d)}$ to solution of (14).
7:         Normalize $\boldsymbol{W}^{(d)}$ and set the scaling factors to $\boldsymbol{\gamma}$.
8:         Update $n_{\mathrm{iter}} += 1$
9:     **end for**
10: **until** $n_{\mathrm{iter}} = n_{\max}$
11: **return** $[\boldsymbol{\gamma}; \boldsymbol{W}^{(1)}, \boldsymbol{W}^{(2)}, \ldots, \boldsymbol{W}^{(D)}]$

---

## IV. SEIZURE DETECTION

In this section we describe the seizure detection pipeline. A schematic of this pipeline can be found in Figure 1. We aimed to keep the pipeline close to that of [28], a previous work on seizure detection with behind-the-ear EEG using an SVM classifier.

### A. Preprocessing

The EEG recordings are first band-pass filtered between 1-25Hz. Then, the EEG data is divided into 2s segments. For testing 50% overlap was used between the segments. However, for the training data different to reduce the data imbalance different amounts of overlap were used for seizure and non-seizure segments. For seizure data 90% overlap was used and for the non-seizure data no overlap was used. Segments with an RMS amplitude less than $13\mu V$ or larger than $150\mu V$ were removed, as they contained mainly background EEG or high-amplitude artifacts [28].

TABLE I: Extracted features [28]

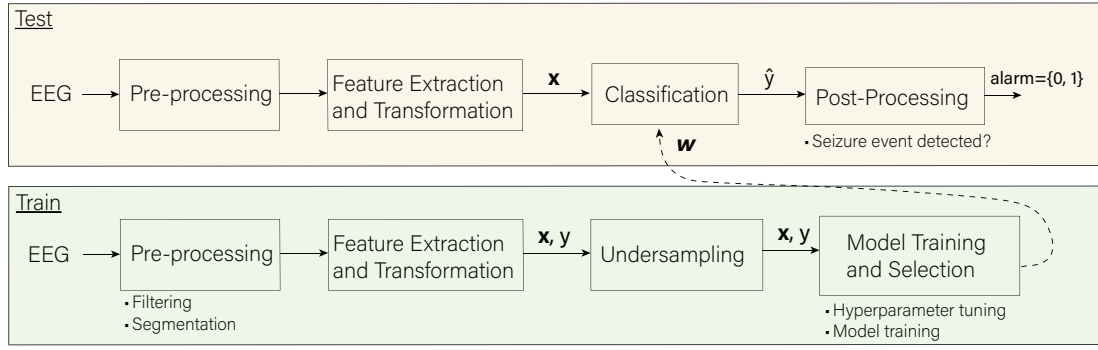| | |
|---|---|
| Time domain | 1-3. Number of zero crossings, maxima and minima<br>4. Skewness<br>5. Kurtosis<br>6. RMS amplitude |
| Frequency domain | 7. Total power<br>8. Peak frequency<br>9-18. Mean and normalized power in frequency bands: $\delta$ (1-3 Hz), $\theta$ (4-8 Hz), $\alpha$ (9-13 Hz), $\beta$ (14-20 Hz), HF (40-80 Hz). |
| Entropy | 19. Spectral entropy<br>20. Sample entropy<br>21. Shannon entropy |

Fig. 1: Seizure detection pipeline.

## B. Feature Extraction and Transformation

To reduce the data size and facilitate classification, we extract features for each segment. These features are a combination of frequently used time domain, frequency domain and non-linear entropy-based features [28]. The features are extracted per EEG channel, thus we end up with a total of $21 \times n_{ch}$ features.

Similar to what was done in [28] we use a power transformation, in our case the Yeo-Johnson transformation [45], to improve the 'normality' of certain features: [5-7, 9-13, 15,16,18,21]. Lastly, all features were scaled to lie in the range of $[-0.5, 0.5]$.

## C. Class Weighting and Undersampling

Even though we use more overlap for the seizure segments than the non-seizure segments, the data is still highly imbalanced. Learning with imbalanced data is a challenging topic as it 'destroys' the i.i.d assumption [46]. Thus, we use a combination of undersampling and class weighting (i.e. cost sensitivity) to deal with this problem.

Class weighting is a popular approach to deal with a class imbalance [47]. This way more weight is put on the misclassification of samples from one class (typically the *minority* class) versus the other. Adding class weighting to the data fitting term of (8) and (14) we obtain,

$$\min_{\boldsymbol{\mathcal{W}}} \quad \frac{1}{N} \left[ C^+ \sum_{n^+=1}^{N^+} \left( \langle \boldsymbol{\Phi}(\boldsymbol{x}_n^+), \boldsymbol{\mathcal{W}} \rangle_{\mathrm{F}} - y_n^+ \right)^2 \right.$$
$$\left. + C^- \sum_{n^-=1}^{N^-} \left( \langle \boldsymbol{\Phi}(\boldsymbol{x}_n^-), \boldsymbol{\mathcal{W}} \rangle_{\mathrm{F}} - y_n^- \right)^2 \right] \quad (15)$$
$$+ \text{ regularization}$$

A popular choice for the class weights, and the one used in this paper, is $C^+ = \frac{N}{2N^+}$ and $C^- = \frac{N}{2N^-}$, which reduces to $C^+ = C^- = 1$ if $N^+ = N^- = \frac{1}{2}N$.

Besides using class-weighting, we also undersample the the non-seizure data. This is done until a ratio of 1:10 (seizure to non-seizure data) is obtained. For the patient independent model the same stategy is used an in [28]. Thus, for every 15 minutes in the first 24 hour period, 10 non-seizure segments are 'selected'. These selected samples are then randomly

sampled to obtain the desired ratio. This way, we aim to obtain data from all the different brain stages that occur throughout the day. For the patient-specific or patient-adapted case, the undersampling is done at random per crossvalidation fold.

## D. Post-Processing and Model Evaluation

Seizure detection algorithms typically suffer from many false positives, which can create a barrier to the adoption of a wearable [48]. Post-processing can help to reduce these false positive by filtering out 'sporadic' positive predictions. These strategies are effective because seizures typically are of longer duration than the chosen segment length. A minimum duration of 10 seconds is often considered for automatic seizure detection algorithms. In this paper, we only consider positive detections when at least 8 out of 10 consecutive segments are classified as a seizure [28].

Because a seizure is an *event*, it is typically more relevant to consider event-based scores rather than segment-based scores. In this work, we use the any-overlap method to evaluate the event-based performance [49]. The any-overlap method considers a prediction correct if it overlaps with any portion of the true seizure event. The event-based metrics we use are the same as proposed in the SzCORE framework [50]: sensitivity, precision, F1-score and false alarms per day (FA/24hr).

## E. Dataset

For the seizure detection experiments, we use the SeizIT1 behind-the-ear EEG dataset [22]. The SeizeIT1 project was approved by the UZ Leuven ethics committee and informed consent was obtained from every participant [28]. This committee also approved the use of the dataset for the current study (S67350 - amendment 3). Furthermore, approval of the current study was given by the TU Delft Human Research Ethics Committee (HREC) (application ID: 5759).

The SeizeIT1 data was acquired in the hospital during presurgical evaluation, where patients were monitored with vEEG for multiple days. The data includes behind-the-ear EEG (two electrodes per side), 10-20 full scalp EEG and single-lead ECG. The median duration of the seizures was 50 seconds, with seizures varying from 11 to 695 seconds in length. The seizures consisted mainly of Focal Impared Awareness (FIA) seizures (89%), with most of the seizures originating from the (fronto-)temporal lobe (91%).

Each EEG file is accompanied by two different annotation files. One file ('a1') contains the annotations from the neurologist, which did not always contain an end time for each seizure. The neurologist also noted whether the seizure was visible on the behind-the-ear channels and (if visible) from which hemisphere the seizure originated. The other file ('a2') contains annotations adapted by an engineer [28] where each seizure has the same length (10 s). Because limiting the seizure duration to only 10 seconds increases the data imbalance even further, we combine the neurologists' annotations with those of the engineer. Thus, in this 'combined annotation', each seizure has an end time and a duration of at least 10 seconds.

For the seizure detection algorithm, we use a bipolar EEG montage that includes a cross-head channel (between the 'top' behind-the-ear (bhe) electrodes on each side) and a 'lateral' channel. The lateral channel is the difference between the two bhe electrodes on one side of the head. This side is different per patient and depends on the lateralization of the patient's seizures.

## V. EXPERIMENTS AND RESULTS

The code to obtain the experiments and results presented below can be found on Github[3]. For the event-based scores, we used the `timescoring` library from the SzCORE framework with its default settings [50].

### A. Feature Map

As in [15], [20] we approximate the RBF kernel with sinusoidal basis functions [40]. This means the local feature maps of (5) are defined as,

$$
\left( \phi^{(d)}\left(x^{(d)}\right) \right)_{i_d} = \frac{1}{\sqrt{U_d}} p\left(\frac{\pi i_d}{2U_d}\right) \sin\left(\frac{\pi i_d \left(x^{(d)} + U_d\right)}{2U_d}\right),
$$
(16)

for $i_d = 1, \ldots, M_d$, with input data, $x^{(d)} \in [-U_d, U_d]$, and where $p(\cdot)$ is the spectral density of the RBF kernel, $p(z) = \sqrt{2\pi\sigma^2} \exp\left(-2\pi^2\sigma^2 z^2\right)$ [41], with $\sigma$ as the lengthscale of the RBF kernel.

To simplify things we choose $\forall d, \ M_d = M$ and $U_d = U$. This leads to three hyperparameters for the feature map, $U$, $M$ and $\sigma$. How well the feature map approximates the RBF kernel for a given lengthscale depends on the number of basis functions $M$. Generally, for larger $\sigma$, more basis functions are needed because the approximation becomes less accurate close to the boundary. Another way to deal with these boundary effects is to increase $U$, while keeping the scaling of the input data the same. This ensures that the input data lies further away from the boundary of the 'approximation' interval.

Below we describe a practical heuristic approach to obtain 'good' values for $M$ and $U$ given some $\sigma$:

1) In case of a large training set, collect a *small* sample of the training dataset[4], and scale the data to a certain interval (e.g. $[-0.5, 0.5]$)

---

[3]https://github.com/sderooij/Adapt-TKRR-experiments
[4]We used a sample of 100 datapoints with a 1:1 ratio between positive and negative samples.



(a) Source model

(b) Target-only model

(c) Adapted model $\mu = 10^{-6}$

(d) Adapted model $\mu = 10^{-4}$

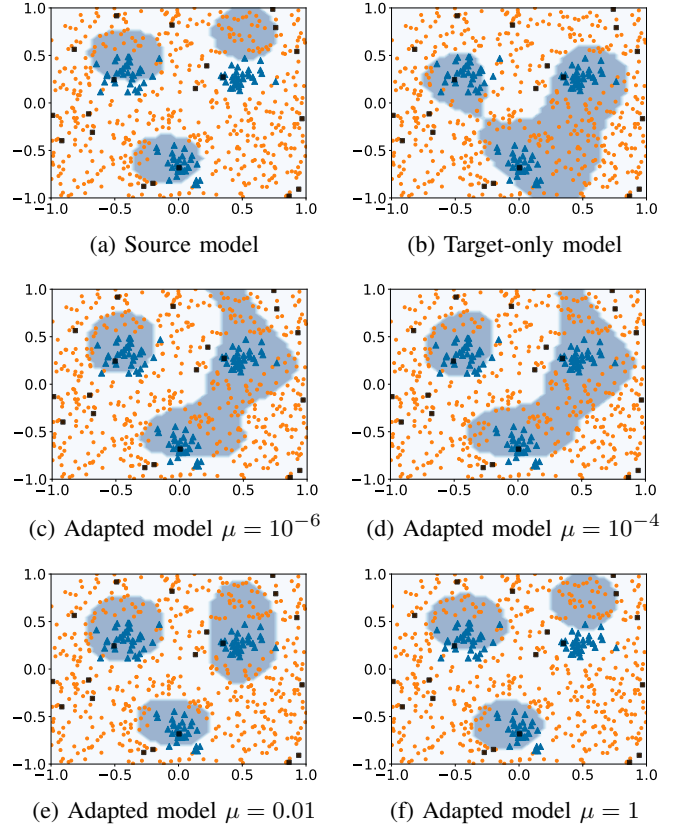(e) Adapted model $\mu = 0.01$

(f) Adapted model $\mu = 1$

Fig. 2: Decision boundaries of different TKRR models on the synthetic target data. The orange dots show the negative samples, the blue triangles the positive samples and the samples used for training are highlighted by the black squares. The shaded blue area shows where the classifier labels samples as *positive*.

2) Set $\sigma$ and make a grid for possible values of $M$ and $U$.
3) For each combination of $M$ and $U$, compute the kernel matrix for the feature map, $\boldsymbol{K}_\Phi(i,j) = \langle \Phi(\boldsymbol{x}_i), \Phi(\boldsymbol{x}_j) \rangle_F$.
4) Compare the $\boldsymbol{K}_\Phi$ matrices with the true RBF kernel matrix, $\|\boldsymbol{K}_\Phi - \boldsymbol{K}_{\mathrm{RBF}}\|_F / \|\boldsymbol{K}_{\mathrm{RBF}}\|_F$.
5) Select $M$ and $U$ based on the relative error (and desired model size).

It should be noted that the choice for $M$ changes the size and complexity of the model. Therefore, there might be a trade-off between the 'accuracy' of the feature map and the desired model size.

### B. An Example with Synthetic Data

To showcase the effectiveness of the Adapt-TKM algorithm, we first apply it to some synthetically generated data. We use the same strategy as in the Adapt-SVM paper [16]. Thus, we generate the positive data for the source and target dataset from a Gaussian mixture with 3 components, and the negative data is sampled from a uniform distribution outside of the area of the positive data. In the source dataset, the Gaussian components are centered at $(-0.4, 0.5)$, $(0.5, 0.7)$ and $(-0.1, -0.6)$, while in the target dataset, the means are

| $M$ \ $U$ | 1 | 1.25 | 1.5 | 1.75 | 2 | 2.25 |
|---|---|---|---|---|---|---|
| 10 | 0.42 | 0.051 | $1.3 \cdot 10^{-3}$ | $5.6 \cdot 10^{-3}$ | $1.6 \cdot 10^{-2}$ | $3.7 \cdot 10^{-2}$ |
| 11 | 0.42 | 0.051 | $8.1 \cdot 10^{-4}$ | $2.2 \cdot 10^{-3}$ | $8.0 \cdot 10^{-3}$ | $2.0 \cdot 10^{-2}$ |
| 12 | 0.42 | 0.051 | $7.4 \cdot 10^{-4}$ | $7.5 \cdot 10^{-4}$ | $3.7 \cdot 10^{-3}$ | $1.0 \cdot 10^{-2}$ |
| 13 | 0.42 | 0.051 | $7.3 \cdot 10^{-4}$ | $2.5 \cdot 10^{-4}$ | $1.5 \cdot 10^{-3}$ | $5.2 \cdot 10^{-3}$ |
| 14 | 0.42 | 0.051 | $7.3 \cdot 10^{-4}$ | $8.0 \cdot 10^{-5}$ | $6.0 \cdot 10^{-4}$ | $2.6 \cdot 10^{-3}$ |
| 15 | 0.42 | 0.051 | $7.2 \cdot 10^{-4}$ | $2.3 \cdot 10^{-5}$ | $2.3 \cdot 10^{-4}$ | $1.2 \cdot 10^{-3}$ |
| 16 | 0.42 | 0.051 | $7.2 \cdot 10^{-4}$ | $5.8 \cdot 10^{-6}$ | $8.5 \cdot 10^{-5}$ | $5.1 \cdot 10^{-4}$ |
| 17 | 0.42 | 0.051 | $7.2 \cdot 10^{-4}$ | $1.7 \cdot 10^{-6}$ | $2.8 \cdot 10^{-5}$ | $2.1 \cdot 10^{-4}$ |
| 18 | 0.42 | 0.051 | $7.2 \cdot 10^{-4}$ | $1.2 \cdot 10^{-6}$ | $8.4 \cdot 10^{-6}$ | $8.8 \cdot 10^{-5}$ |
| 19 | 0.42 | 0.051 | $7.2 \cdot 10^{-4}$ | $1.0 \cdot 10^{-6}$ | $2.6 \cdot 10^{-6}$ | $3.3 \cdot 10^{-5}$ |
| 20 | 0.42 | 0.051 | $7.2 \cdot 10^{-4}$ | $9.9 \cdot 10^{-7}$ | $7.4 \cdot 10^{-7}$ | $1.2 \cdot 10^{-5}$ |

TABLE II: Relative error, $\|\boldsymbol{K}_\Phi - \boldsymbol{K}_{\mathrm{RBF}}\|_{\mathrm{F}} / \|\boldsymbol{K}_{\mathrm{RBF}}\|_{\mathrm{F}}$, at different values for $M$ and $U$ on the synthetic source dataset.

| Model | Sensitivity | Precision | F1-score | FA/24hr |
|---|---|---|---|---|
| SVM-PI | 0.620 (0.22) | 0.548 (0.42) | 0.458 (0.31) | 9.06 (14) |
| TKRR-PI | 0.630 (0.18) | 0.575 (0.44) | 0.486 (0.33) | 12.0 (23) |
| SVM-PS | 0.623 (0.12) | 0.508 (0.31) | 0.443 (0.21) | 16.6 (22) |
| TKRR-PS | 0.625 (0.14) | 0.414 (0.35) | 0.347 (0.25) | 26.2 (30) |
| SVM-PA | 0.621 (0.15) | 0.555 (0.33) | 0.484 (0.22) | 6.5 (7.1) |
| TKRR-PA | 0.632 (0.17) | 0.601 (0.37) | 0.523 (0.28) | 5.19 (7.3) |

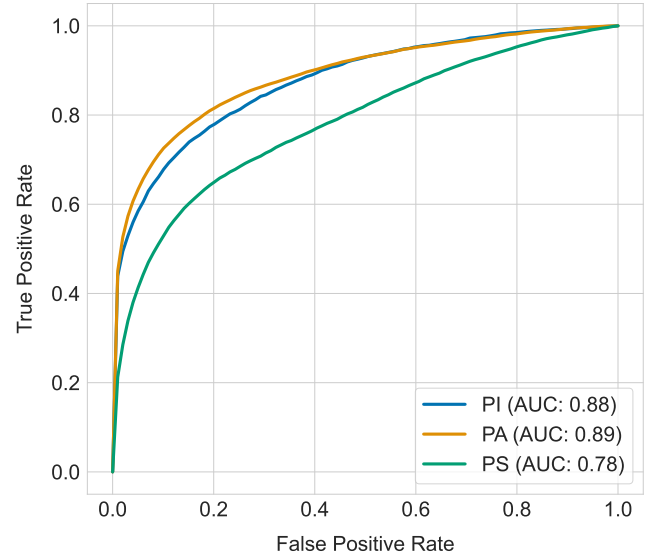TABLE III: Average performance across the patients, standard deviation between brackets.



Fig. 3: ROC curves of the TKRR models

at $(-0.4, 0.3)$, $(0.5, 0.3)$ and $(0, -0.65)$. Each dataset has 500 negative and 100 positive samples.

First, we train a source model using the TKRR classifier with the aforementioned feature map (16). Due to the data imbalance, class weighting was used as described in Section IV-C. The feature map parameter $\sigma$ was chosen to coincide with the RBF parameter used in [16][5]. Then, the other feature map parameters ($M$ and $U$) are determined using the procedure described in the previous section. In Table II, we show the relative error between the true RBF kernel matrix, $K_{\mathrm{RBF}}$ and the one calculated with the feature map, $K_{\boldsymbol{\Phi}}$. Based on these values, $M = 14$ and $U = 1.75$ were selected, as the error for these values was deemed sufficiently low. The rank of the CPD weight tensor was set to $R = 4$ and the regularization parameter to $\lambda = 0.001$. Figure 2a shows the decision boundary of the source model applied to the target dataset.

The same model parameters were used to train a target-only TKRR model (Figure 2b) using a small subset of samples from the target dataset: 17 negative and 3 positive samples. This same subset was used to train the Adapt-TKRR model for varying values of $\mu$ (Figures 2c-2f).

From Figure 2, it is clear that the adapted model becomes more similar to the source model for high $\mu$ and more similar to the target model for low $\mu$. The best performance is obtained at $\mu = 0.01$. In this case, the F1-score is $0.67$, which is better than both the source ($0.48$) and target-only model's ($0.61$) F1-scores.

### C. Performance of Adapt-TKRR on Seizure Detection

In this section, we test the effectiveness of the Adapt-TKRR model in personalizing seizure detection models. Thus, we first implement a patient-independent (i.e. *source*) model (TKRR-PI), which is then *adapted* using a small amount of patient-specific (or *target*) data. The performance of these patient-*adapted* models (TKRR-PA) is also compared to patient-specific models (TKRR-PS), which are trained only on the

patient-specific data. As a reference, we compare the performance of the tensor kernel machines to the performance of the SVM and Adapt-SVM classifiers.

The patient-independent (PI) models are validated using a leave-one-patient-out cross-validation procedure. Furthermore, we only use the seizures visible on the behind-the-ear channels to train the PI models, as we found that this led to the best results. For testing, we did not adhere to these restrictions, so the models were tested on all the seizures. Hyperparameter tuning for the PI models was done using five-fold cross-validation on the training sets, maximizing the area under the ROC curve (AUROC).

The patient-adapted (PA) and patient-specific (PS) models are validated using a leave-one-seizure-*in* cross-validation procedure. The splits of the cross-validation folds are set in the middle between two seizures. Thus, every fold contains one seizure with varying amounts of non-seizure data. As stated in Section IV-C, this non-seizure data is randomly sampled to obtain the desired ratio for the training set. Furthermore, we only consider patients who experienced five or more seizures in order to have sufficient data for testing.

For the PS and PA models, the training datasets were too small to obtain 'good' hyperparameters through cross-validation on the training data. Thus, we used the same

---

[5]In [16] $\kappa(x_i, x_j) = e^{-\rho\|x_i - x_j\|^2}$ with $\rho = 5$, so $\sigma = \sqrt{\frac{1}{2 \cdot \rho}} = \sqrt{\frac{1}{10}}$.
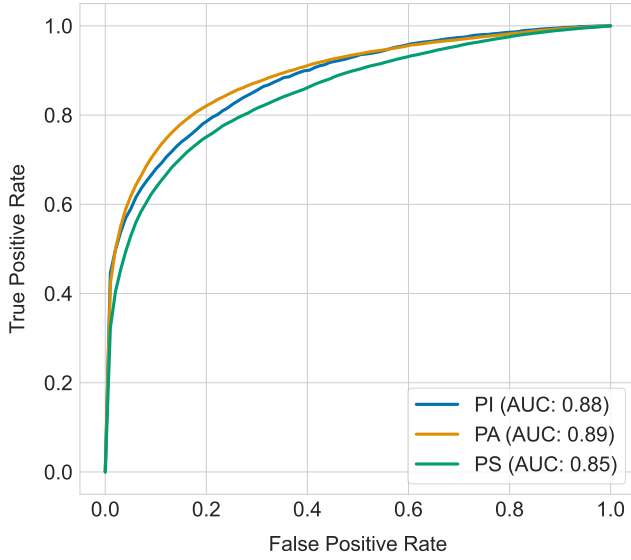
Fig. 4: ROC curves of the SVM models

| $\mu$ | Sensitivity | Precision | F1-score | FA/24hr |
|---|---|---|---|---|
| $1 \times 10^{-4}$ | 0.628 (0.15) | 0.508 (0.33) | 0.456 (0.24) | 11.1 (12) |
| $1 \times 10^{-3}$ | 0.636 (0.15) | 0.552 (0.34) | 0.490 (0.25) | 5.54 (5.5) |
| 0.01 | 0.632 (0.17) | 0.601 (0.37) | 0.523 (0.28) | 5.19 (7.3) |
| 0.1 | 0.635 (0.19) | 0.608 (0.41) | 0.512 (0.31) | 8.04 (13) |
| 'Optimal' | 0.648 (0.16) | 0.671 (0.35) | 0.582 (0.25) | 4.19 (7.1) |

TABLE IV: Average performance (and standard deviation) across the patients of the TKRR-PA model for different $\mu$ values. 'Optimal' is the case when using the optimal value for each patient based on the F1-score.

hyperparameters as for the corresponding PI model. Because of the difference in regularization, the regularization parameter of the adaptive models could not be directly inherited from the PI model. Thus, we set them to a fixed value: $\mu = 0.01$ (for Adapt-TKRR, same as the optimum for the synthetic data V-B) and $C = 10$ (for Adapt-SVM, same as in [16]). In the next section (V-D), we will discuss the influence of the $\mu$ hyperparameter for the TKRR-PA models more in depth.

Table III shows the average performance across the patients for the different models. The scores presented are the event-based metrics (so after post-processing as described in Section IV-D). To better compare the performance of the different models, we tuned the thresholds to obtain similar sensitivity levels. Thus, we can see that the PA models have a lower false positive rate at comparable sensitivity to the PI models. On the other hand, the TKRR-PS model has the worst performance across all metrics.

This result is supported by the difference in ROC curves, which are plotted in Figure 3 for the TKRR models and in Figure 4 for the SVM models. From these ROC curves, it is also clear that when it comes to the performance of the PI and PA models, the (Adapt-)TKRR and (Adapt-)SVM classifiers lead to fairly similar results. The TKRR-PS model, however, performs much worse than its SVM counterpart. This is likely due to the inheritance of the model parameters from the PI model, which for the TKRR classifier also means inheriting its size.

So far, the results presented have been averaged over the patients. It might, however, be more interesting to see the patient-by-patient results. To this end, we compute the (event-based) F1-scores per patient for the TKRR models. These are shown in Figure 5. Here, we can see that the patient-adapted model performs the best for 9 out of 15 patients. There are, however, also patients for which the PA model decreases the performance of the PI model (patients 33, 65, 78, 82 and 99). These seem to be patients where the PI model

performs relatively well (F1-score $> 0.6$), and the PS model performs significantly worse, suggesting that adaptation may have a negative effect here.

### D. Influence of the Regularization Parameter

As previously seen on the synthetic data in Section V-B, the choice for the regularization parameter, $\mu$, significantly influences the performance of the Adapt-TKRR model. To investigate the influence of these parameters on the performance, we compare the F1-scores for each patient for different $\mu$ values in Figure 6. From this Figure, it is clear that some patients (33, 65, 78, 79, 82, 94 and 99) benefit from high $\mu$ values, while for others, the Adapt-TKRR performs best when $\mu$ is more 'moderate' (16, 36, 70, 71 and 76) or low (13, 61 and 83).

Typically, the patients for whom a high $\mu$ value is optimal are the ones where the performance of the PI model is already 'good'. Whereas the ones that benefit from lower $\mu$ values are typically those for whom the PS model performs (relatively) better. If we choose the optimal value for $\mu$ for each patient (based on Figure 6), we can reach performance that is better than (or at least on par with) the PI and PS models for all but two patients. For these two patients (78 82), the $\mu$ value needs to be even higher than $0.1$ to reach the same performance as the PI model, which essentially means that no transfer occurs and the model stays the 'same' (such as the case for $\mu = 1$ on the synthetic data, Figure 2f).

Besides the patient-by-patient results, we also look at the overall performance of the TKRR-PA models for different $\mu$ values. The event-based scores of these models are found in Table IV. In this table, we also present the performance in the case where the 'optimal' $\mu$ is chosen for each patient (based on Figure 6). As expected, this 'optimal' model has the best performance among all metrics. Interestingly, the next best performing model is the one where $\mu = 0.01$, which corresponds with the synthetic data experiment (and was therefore chosen for the TKRR-PA model). This suggests that selecting $\mu = 0.01$ for the Adapt-TKRR model provides a reasonable trade-off between source and target if no prior information is available.

### E. Initialization of the Adapt-TKM Model

In this section, we investigate whether initialization of the Adapt-TKRR model with the source model weights im-
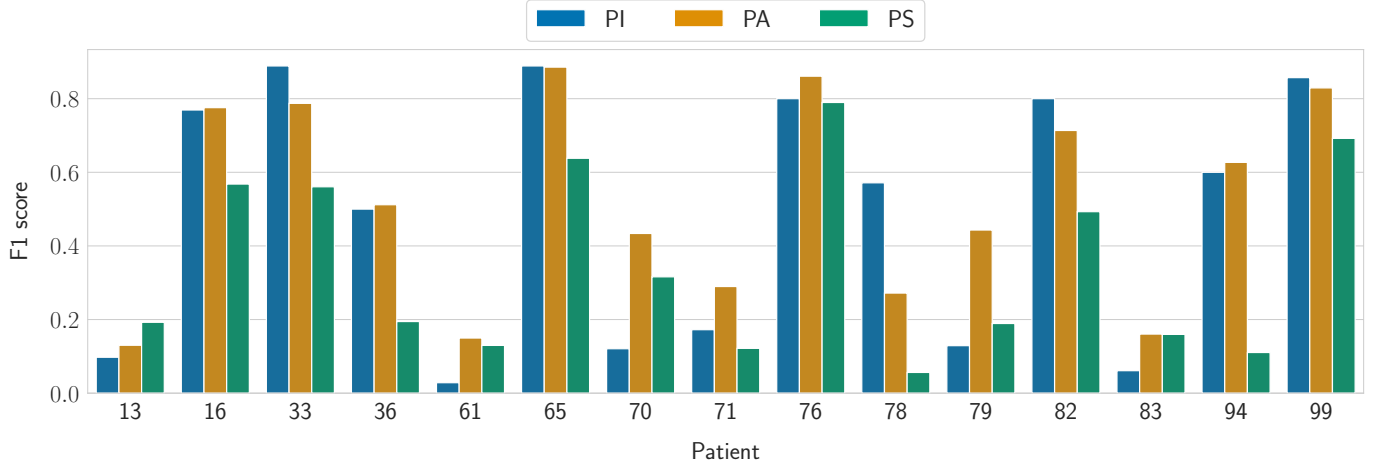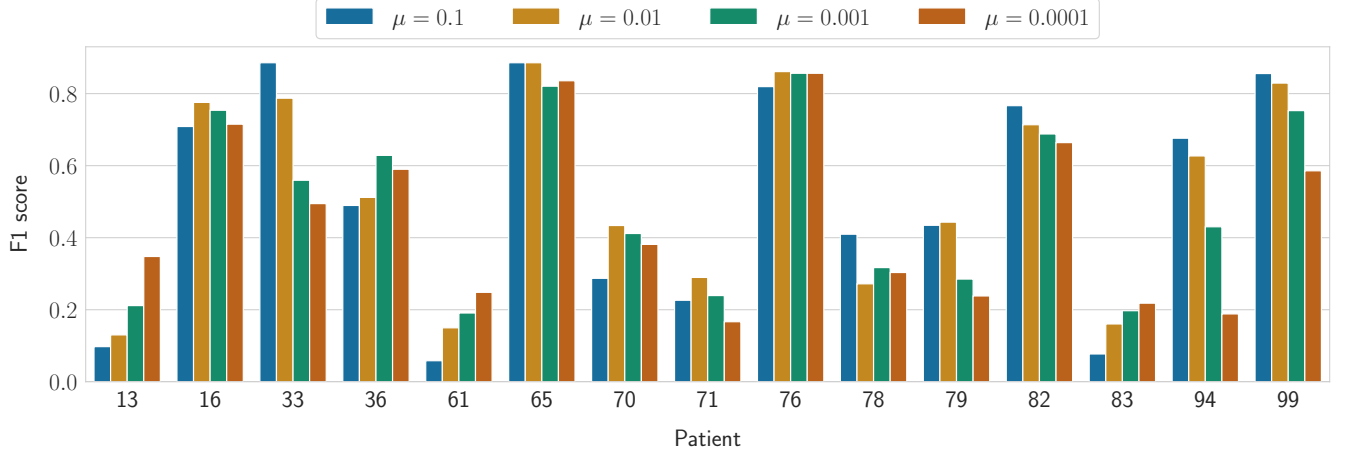
Fig. 5: F1-scores per patient for the TKRR models.



Fig. 6: F1-scores per patient for different values of $\mu$ in the Adapt-TKRR model.

proves convergence. In our previous work [15], we found that initialising a patient-specific TKRR model with the PI weights improved the convergence speed significantly ('patient-*finetuning*'). To test whether this also holds for the Adapt-TKRR models, we plot the normalized training loss of the TKRR-PA models against the number of iterations (Figure 7). The loss is computed by calculating the data fitting term of (15), i.e. $\ell = \frac{1}{N}[C^+ \sum_{n^+=1}^{N^+} (\langle \boldsymbol{\Phi}(\boldsymbol{x}_n^+), \boldsymbol{\mathcal{W}} \rangle_{\mathrm{F}} - y_n^+)^2 + C^- \sum_{n^-=1}^{N^-} (\langle \boldsymbol{\Phi}(\boldsymbol{x}_n^-), \boldsymbol{\mathcal{W}} \rangle_{\mathrm{F}} - y_n^-)^2]$.

The models using random initialization start to converge after around 42 iterations; this equals exactly one full sweep over the factor matrices of the CPD weights ($D = 42$). The optimum for the randomly initialized models is not reached until about 110 iterations are done. The models initialized with the source weights, on the other hand, reach their optimum already after about one sweep (42 iterations). Random initialization does seem to lead to a better optimum than initialization with source weights, though the difference is small: 0.011.

### F. Model Efficiency

To analyze the efficiency of the models, we compare the total number of parameters and the inference speed for each model. The number of model parameters for the SVM models are equal to the size of support vectors ($N_{sv}D$) plus the dual coefficients ($N_{sv}$). The amount of model parameters for the TKRR models is simply determined by the size of the CPD weight tensor ($DMR$).

The second column of Table V presents the average number of parameters for each model. It should be noted that because the model hyperparameters for the TKRR models were kept the same for the PA and PS models as the corresponding PI model, the amount of model hyperparameters is the same across the models. From these values, it is clear that the TKRR models have about 100 times fewer model parameters than the SVM-PI model. Only the SVM-PS models have slightly fewer hyperparameters, though it is still in the same order of magnitude ($10^4$) as the TKRR models.

The inference speeds were determined by randomly selecting 1000 samples (same for each model) and computing the
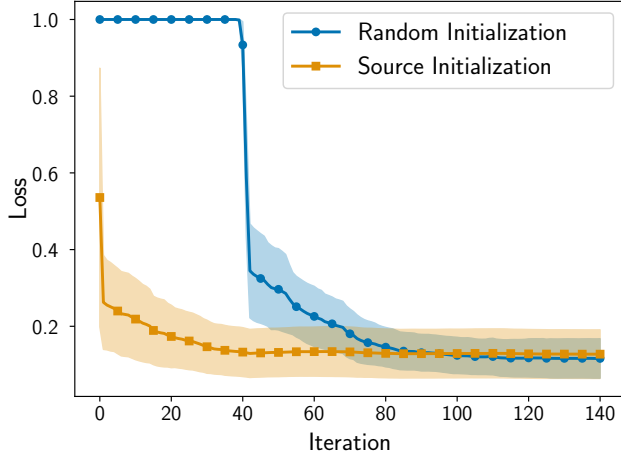
Fig. 7: Mean convergence of the Adapt-TKRR models (PA models) for initialization with the source model (PI model) and random initialization. The y-axis shows the loss (i.e. the data fitting term of (15)), and the x-axis shows the iteration count. The shaded areas show the spread of $\pm 1$ standard deviation. Dots are placed every 5 iterations.

| Model | No. Parameters | Inference Time (s) |
|---|---|---|
| SVM-PI | $1.363 \cdot 10^6$ | $3.3 \cdot 10^{-3}$ |
| TKRR-PI | $1.988 \cdot 10^4$ | $3.9 \cdot 10^{-5}$ |
| SVM-PS | $1.140 \cdot 10^4$ | $3.3 \cdot 10^{-5}$ |
| TKRR-PS | $1.988 \cdot 10^4$ | $3.9 \cdot 10^{-5}$ |
| SVM-PA | $1.369 \cdot 10^6$ | $5.2 \cdot 10^{-3}$ |
| TKRR-PA | $1.988 \cdot 10^4$ | $3.9 \cdot 10^{-5}$ |

TABLE V: Average number of parameters and inference time (in seconds) per model.

average time it takes the model to determine the label for a sample. The computations were performed on the same hardware: a laptop computer with a quad-core Intel i7-1185G7 (3.0 GHz) processor and 32GB of RAM. The third column of Table V presents the average inference speed per model. Again, we see that the SVM-PI and SVM-PA models are the least efficient. They are about 100 times slower than the TKRR and the SVM-PS models.

## VI. Discussion and Future Work

From the results presented in the previous section, it is clear that the adaptive tensor kernel machine can be an effective and efficient transfer learning model. It is able to improve the seizure detection performance (for some patients) by personalizing a patient-independent model. Furthermore, it can do this more efficiently than the Adapt-SVM model and without increasing the model size. Besides the efficiency gain, there is also another advantage to the Adapt-TKM model over the Adapt-SVM. The Adapt-SVM uses the support vectors of the source model, these support vectors are the training data of the source model. In the seizure detection case, this may raise privacy concerns due to 'saving' data from other patients.

The Adapt-TKM model does not 'save' any source data and is thus more 'privacy-preserving' [13].

While we have seen that the Adapt-TKM can be effective in personalizing seizure detection models, there are still limitations to its effectiveness. Most notable is that seizure detection performance did not improve for all patients. Especially for patients who already had good performance from the PI models the adaptation could have a negative effect. Part of this negative effect is because choosing the regularization of the Adapt-TKM model (or the rate of transfer) is not straightforward. While the value we chose ($\mu = 0.01$) seemed to have decent performance overall, this did not hold for all cases (Figure 6). Thus, finding an effective strategy to choose this parameter remains a topic of future research.

Another limitation of the Adapt-TKM model, or actually tensor kernel machines in general, is that it introduces extra hyperparameters compared to 'standard' kernel machines. For the extra hyperparameters of the feature map, $M$ and $U$, we propose a heuristic approach in Section V-A. However, there is no straightforward way to select the rank, $R$, of the CPD weight tensor, and it was done via a grid search on the PI data. Inheriting this hyperparameter from the PI model may have been suboptimal for the PA and PS models and could have led to larger models than necessary. A more informed way to choose (or update) the rank of the CPD weight tensor, may improve performance and reduce the need for an exhaustive grid search.

## VII. Conclusion

In this work, we have presented an efficient transfer learning method using tensor kernel machines. This method, the adaptive tensor kernel machine (Adapt-TKM), can adapt a previously learned source model to new data. It does so by transferring 'knowledge' from the source weights through regularization while learning the adapted model weights. The advantage of this Adapt-TKM model over the conventional adaptive SVM is that it is able to learn a compact non-linear model in the primal, using low-rank tensor networks to deal with the high-dimensionality of the feature space. This way the Adapt-TKM does not increase in size compared to the source, as opposed to the kernelized adaptive SVM.

We showcased the effectiveness of the method by applying it to seizure detection for wearable behind-the-ear EEG. The Adapt-TKM model efficiently personalizes a patient-independent model with minimal patient-specific data. This patient-adapted model outperformed both the patient-independent and patient-specific models. Moreover, the Adapt-TKM model was about a hundred times more efficient than its SVM counterpart, both in terms of the number of parameters and its inference speed.

## REFERENCES

[1] F. Zhuang, Z. Qi, K. Duan, D. Xi, Y. Zhu, H. Zhu, H. Xiong, and Q. He, "A Comprehensive Survey on Transfer Learning," *Proceedings of the IEEE*, vol. 109, no. 1, pp. 43–76, Jan. 2021.

[2] A. Hosna, E. Merry, J. Gyalmo, Z. Alom, Z. Aung, and M. A. Azim, "Transfer learning: A friendly introduction," *Journal of Big Data*, vol. 9, no. 1, p. 102, Oct. 2022.

[3] S. J. Pan and Q. Yang, "A Survey on Transfer Learning," *IEEE Transactions on Knowledge and Data Engineering*, vol. 22, no. 10, pp. 1345–1359, Oct. 2010.

[4] G. Balakrishnan and Z. Syed, "Scalable Personalization of Long-Term Physiological Monitoring: Active Learning Methodologies for Epileptic Seizure Onset Detection," in *Proceedings of the Fifteenth International Conference on Artificial Intelligence and Statistics*. PMLR, Mar. 2012, pp. 73–81.

[5] L. H. Goetz and N. J. Schork, "Personalized medicine: Motivation, challenges, and progress," *Fertility and Sterility*, vol. 109, no. 6, pp. 952–963, Jun. 2018.

[6] Z. Wan, R. Yang, M. Huang, N. Zeng, and X. Liu, "A review on transfer learning in EEG signal analysis," *Neurocomputing*, vol. 421, pp. 1–14, Jan. 2021.

[7] T. De Cooman, K. Vandecasteele, C. Varon, B. Hunyadi, E. Cleeren, W. Van Paesschen, and S. Van Huffel, "Personalizing Heart Rate-Based Seizure Detection Using Supervised SVM Transfer Learning," *Frontiers in Neurology*, vol. 11, 2020.

[8] I. E. Scheffer, S. Berkovic, G. Capovilla, M. B. Connolly, J. French, L. Guilhoto, E. Hirsch, S. Jain, G. W. Mathern, S. L. Moshé, D. R. Nordli, E. Perucca, T. Tomson, S. Wiebe, Y.-H. Zhang, and S. M. Zuberi, "ILAE classification of the epilepsies: Position paper of the ILAE Commission for Classification and Terminology," *Epilepsia*, vol. 58, no. 4, pp. 512–521, Apr. 2017.

[9] X. Cui, J. Cao, T. Jiang, and F. Gao, "Transfer Learning Based Seizure Detection: A Review," in *Cognitive Computation and Systems*, ser. Communications in Computer and Information Science, F. Sun, J. Li, H. Liu, and Z. Chu, Eds. Singapore: Springer Nature, 2023, pp. 160–175.

[10] C. Yang, Z. Deng, K.-S. Choi, Y. Jiang, and S. Wang, "Transductive domain adaptive learning for epileptic electroencephalogram recognition," *Artificial Intelligence in Medicine*, vol. 62, no. 3, pp. 165–177, Nov. 2014.

[11] Y. Jiang, D. Wu, Z. Deng, P. Qian, J. Wang, G. Wang, F.-L. Chung, K.-S. Choi, and S. Wang, "Seizure Classification From EEG Signals Using Transfer Learning, Semi-Supervised Learning and TSK Fuzzy System," *IEEE Transactions on Neural Systems and Rehabilitation Engineering*, vol. 25, no. 12, pp. 2270–2284, Dec. 2017.

[12] X. Jiang, K. Xu, and W. Chen, "Transfer Component Analysis to Reduce Individual Difference of EEG Characteristics for Automated Seizure Detection," in *2019 IEEE Biomedical Circuits and Systems Conference (BioCAS)*, Oct. 2019, pp. 1–4.

[13] C. Zhao, R. Peng, and D. Wu, "Source-Free Domain Adaptation (SFDA) for Privacy-Preserving Seizure Subtype Classification," *IEEE Transactions on Neural Systems and Rehabilitation Engineering*, vol. 31, pp. 2315–2325, 2023.

[14] A. Page, C. Shea, and T. Mohsenin, "Wearable seizure detection using convolutional neural networks with transfer learning," in *2016 IEEE International Symposium on Circuits and Systems (ISCAS)*, May 2016, pp. 1086–1089.

[15] S. De Rooij, F. Wesel, and B. Hunyadi, "Efficient Patient Fine-Tuned Seizure Detection with a Tensor Kernel Machine," in *2024 32nd European Signal Processing Conference (EUSIPCO)*, Aug. 2024, pp. 1372–1376.

[16] J. Yang, R. Yan, and A. G. Hauptmann, "Adapting SVM Classifiers to Data with Shifted Distributions," in *Seventh IEEE International Conference on Data Mining Workshops (ICDMW 2007)*, Oct. 2007, pp. 69–76.

[17] R. Schwartz, J. Dodge, N. A. Smith, and O. Etzioni, "Green AI," *Commun. ACM*, vol. 63, no. 12, pp. 54–63, Nov. 2020.

[18] Y. Mao, X. Yu, K. Huang, Y.-J. Angela Zhang, and J. Zhang, "Green Edge AI: A Contemporary Survey," *Proceedings of the IEEE*, vol. 112, no. 7, pp. 880–911, Jul. 2024.

[19] E. Memmel, C. Menzen, J. Schuurmans, F. Wesel, and K. Batselier, "Position: Tensor networks are a valuable asset for green AI," in *Proceedings of the 41st International Conference on Machine Learning*, ser. ICML'24, vol. 235. Vienna, Austria: JMLR.org, Jul. 2024, pp. 35340–35353.

[20] F. Wesel and K. Batselier, "Large-Scale Learning with Fourier Features and Tensor Decompositions," in *Advances in Neural Information Processing Systems*, vol. 34. Curran Associates, Inc., 2021, pp. 17543–17554.

[21] S. Klus and P. Gelß, "Tensor-based algorithms for image classification," *Algorithms*, vol. 12, no. 11, p. 240, Nov. 2019.

[22] C. Chatzichristos and M. Claro Bhagubai, "SeizeIT1," 2023. [Online]. Available: https://doi.org/10.48804/P5Q0OJ

[23] J. Macea, M. Bhagubai, V. Broux, M. De Vos, and W. Van Paesschen, "In-hospital and home-based long-term monitoring of focal epilepsy with a wearable electroencephalographic device: Diagnostic yield and user experience," *Epilepsia*, vol. 64, no. 4, pp. 937–950, 2023.

[24] J. Lehnen, P. Venkatesh, Z. Yao, A. Aziz, P. V. P. Nguyen, J. Harvey, S. Alick-Lindstrom, A. Doyle, I. Podkorytova, G. Perven, R. Hays, R. Zepeda, R. R. Das, and K. Ding, "Real-Time Seizure Detection Using Behind-the-Ear Wearable System," *Journal of Clinical Neurophysiology*, vol. 42, no. 2, p. 118, Feb. 2025.

[25] M. Bhagubai, C. Chatzichristos, L. Swinnen, J. Macea, J. Zhang, L. Lagae, K. Jansen, A. Schulze-Bonhage, F. Sales, B. Mahler, Y. Weber, W. V. Paesschen, and M. D. Vos, "SeizeIT2: Wearable Dataset Of Patients With Focal Epilepsy," Feb. 2025.

[26] K. Komal, F. Cleary, J. Wells, and L. Bennett, "A systematic review of the literature reporting on remote monitoring epileptic seizure detection devices," *Epilepsy Research*, vol. 201, p. 107334, Mar. 2024.

[27] L. Hadady, P. Klivényi, D. Fabó, and S. Beniczky, "Real-world user experience with seizure detection wearable devices in the home environment," *Epilepsia*, vol. 64, no. S4, pp. S72–S77, 2023.

[28] K. Vandecasteele, T. De Cooman, J. Dan, E. Cleeren, S. Van Huffel, B. Hunyadi, and W. Van Paesschen, "Visual seizure annotation and automated seizure detection using behind-the-ear electroencephalographic channels," *Epilepsia*, vol. 61, no. 4, pp. 766–775, Apr. 2020.

[29] T. G. Kolda and B. W. Bader, "Tensor Decompositions and Applications," *SIAM Review*, vol. 51, no. 3, pp. 455–500, Aug. 2009.

[30] A. Cichocki, N. Lee, I. Oseledets, A.-H. Phan, Q. Zhao, and D. P. Mandic, "Tensor Networks for Dimensionality Reduction and Large-scale Optimization: Part 1 Low-Rank Tensor Decompositions," *Foundations and Trends® in Machine Learning*, vol. 9, no. 4-5, pp. 249–429, 2016.

[31] I. V. Oseledets, "Tensor-Train Decomposition," *SIAM Journal on Scientific Computing*, vol. 33, no. 5, pp. 2295–2317, Jan. 2011.

[32] J. B. Kruskal, "Three-way arrays: Rank and uniqueness of trilinear decompositions, with application to arithmetic complexity and statistics," *Linear Algebra and its Applications*, vol. 18, no. 2, pp. 95–138, Jan. 1977.

[33] F. L. Hitchcock, "The Expression of a Tensor or a Polyadic as a Sum of Products," *Journal of Mathematics and Physics*, vol. 6, no. 1-4, pp. 164–189, 1927.

[34] B. Schölkopf and A. J. Smola, *Learning with Kernels: Support Vector Machines, Regularization, Optimization, and Beyond*, ser. Adaptive Computation and Machine Learning. Cambridge, Mass: MIT Press, 2002.

[35] C.-C. Chang and C.-J. Lin, "LIBSVM: A library for support vector machines," *ACM Transactions on Intelligent Systems and Technology*, vol. 2, no. 3, pp. 27:1–27:27, May 2011.

[36] C. Cortes and V. Vapnik, "Support-vector networks," *Machine Learning*, vol. 20, no. 3, pp. 273–297, Sep. 1995.

[37] J. Suykens and J. Vandewalle, "Least Squares Support Vector Machine Classifiers," *Neural Processing Letters*, vol. 9, no. 3, pp. 293–300, Jun. 1999.

[38] X. Guo, X. Huang, L. Zhang, L. Zhang, A. Plaza, and J. A. Benediktsson, "Support Tensor Machines for Classification of Hyperspectral Remote Sensing Imagery," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 54, no. 6, pp. 3248–3264, Jun. 2016.

[39] F. Wesel and K. Batselier, "Quantized Fourier and Polynomial Features for more Expressive Tensor Network Models," in *Proceedings of The 27th International Conference on Artificial Intelligence and Statistics*. PMLR, Apr. 2024, pp. 1261–1269.

[40] A. Solin and S. Särkkä, "Hilbert space methods for reduced-rank Gaussian process regression," *Statistics and Computing*, vol. 30, no. 2, pp. 419–446, Mar. 2020.

[41] C. E. Rasmussen and C. K. I. Williams, *Gaussian Processes for Machine Learning*. The MIT Press, Nov. 2005.

[42] B. N. Khoromskij, "O(dlog N)-Quantics Approximation of N-d Tensors in High-Dimensional Numerical Modeling," *Constructive Approximation*, vol. 34, no. 2, pp. 257–280, Oct. 2011.

[43] S. Niu, Y. Liu, J. Wang, and H. Song, "A Decade Survey of Transfer Learning (2010–2020)," vol. 1, no. 2, 2020.

[44] J. Yang, R. Yan, and A. G. Hauptmann, "Cross-Domain Video Concept Detection Using Adaptive SVMs," in *Proceedings of the 15th ACM International Conference on Multimedia*. Augsburg Germany: ACM, Sep. 2007, pp. 188–197.

[45] I.-K. Yeo and R. A. Johnson, "A new family of power transformations to improve normality or symmetry," *Biometrika*, vol. 87, no. 4, pp. 954–959, Dec. 2000.

[46] G. Haixiang, L. Yijing, J. Shang, G. Mingyun, H. Yuanyue, and G. Bing, "Learning from class-imbalanced data: Review of methods and applications," *Expert Systems with Applications*, vol. 73, pp. 220–239, May 2017.

[47] R. Akbani, S. Kwek, and N. Japkowicz, "Applying Support Vector Machines to Imbalanced Datasets," in *Machine Learning: ECML 2004*. Berlin, Heidelberg: Springer Berlin Heidelberg, 2004, vol. 3201, pp. 39–50.

[48] M. E. Thompson, J. Langer, and M. Kinfe, "Seizure detection watch improves quality of life for adolescents and their families," *Epilepsy & Behavior*, vol. 98, pp. 188–194, Sep. 2019.

[49] V. Shah, M. Golmohammadi, I. Obeid, and J. Picone, "Objective Evaluation Metrics for Automatic Classification of EEG Events," in *Biomedical Signal Processing*, I. Obeid, I. Selesnick, and J. Picone, Eds. Cham: Springer International Publishing, 2021, pp. 223–255.

[50] J. Dan, U. Pale, A. Amirshahi, W. Cappelletti, T. M. Ingolfsson, X. Wang, A. Cossettini, A. Bernini, L. Benini, S. Beniczky, D. Atienza, and P. Ryvlin, "SzCORE: Seizure Community Open-Source Research Evaluation framework for the validation of electroencephalography-based automated seizure detection algorithms," *Epilepsia*, 2024.