

2021 BEETL Competition: Advancing Transfer Learning for Subject Independence & Heterogenous EEG Data Sets

Xiaoxi Wei^a

A. Aldo Faisal^{a,b}

Moritz Grosse-Wentrup^c

Alexandre Gramfort^d

Sylvain Chevallier^e

Vinay Jayaram^f

Camille Jeunet^g

XIAOXI.WEI18@IMPERIAL.AC.UK

ALDO.FAISAL@IMPERIAL.AC.UK

MORITZ.GROSSE-WENTRUP@UNIVIE.AC.AT

ALEXANDRE.GRAMFORT@INRIA.FR

SYLVAIN.CHEVALLIER@UVSQ.FR

VINAYJAYARAM@FB.COM

CAMILLE.JEUNET@U-BORDEAUX.FR

^a*Brain & Behaviour Lab, Imperial College London, United Kingdom*

^b*Institute of Artificial & Human Intelligence, University of Bayreuth, Germany*

^c*Faculty of Computer Science, CogSciHub, Data Science@Univie, University of Vienna, Austria*

^d*Universite Paris-Saclay, Inria, CEA, Palaiseau, France*

^e*LISV, UVSQ, Université Paris-Saclay, France*

^f*Reality Labs, USA*

^g*University of Bordeaux, France*

Stylianos Bakas^{1,2,3}

Siegfried Ludwig^{1,2}

Konstantinos Barmpas^{1,2}

Mehdi Bahri^{1,2}

Yannis Panagakis^{1,2,4}

Nikolaos Laskaris^{1,2,3}

Dimitrios A. Adamos^{1,2,3}

Stefanos Zafeiriou^{1,2}

William C. Duong^{5,6}

Stephen M. Gordon^{5,6}

Vernon J. Lawhern⁶

Maciej Śliwowski^{7,8,9}

Vincent Rouanne⁷

Piotr Tempczyk^{9,10}

STELIOS@COGITAT.IO

SIEGFRIED@COGITAT.IO

NTINOS@COGITAT.IO

MEHDI@COGITAT.IO

YANNIS@COGITAT.IO

NIKOS@COGITAT.IO

DIMITRIOS@COGITAT.IO

STEFANOS@COGITAT.IO

WDUONG@DCSCORP.COM

SGORDON@DCSCORP.COM

VERNON.J.LAWHERN.CIV@ARMY.MIL

MACIEJ.SLIWOWSKI@OPIUM.SH

VINCENT.ROUANNE@GMAIL.COM

PIOTR.TEMPCZYK@OPIUM.SH

¹*Cogitat Ltd., United Kingdom*

²*Intelligent Behaviour Understanding Group, Imperial College London, United Kingdom*

³*Aristotle University of Thessaloniki, Greece*

⁴*National and Kapodistrian University of Athens, Greece*

⁵*DCS Corporation, Alexandria, VA, USA*

⁶*Human Research and Engineering Directorate, DEVCOM Army Research Laboratory, Aberdeen Proving Ground, MD, USA*

⁷*Univ. Grenoble Alpes, CEA, LETI, Clinatec, F-38000 Grenoble, France*

⁸*Université Paris-Saclay, CEA, List, F-91120, Palaiseau, France*

⁹*Polish National Institute for Machine Learning (OPIUM), Warsaw, Poland*

¹⁰*deeptale.ai, Poland*

Editor: Editor's name

Abstract

Transfer learning and meta-learning offer some of the most promising avenues to unlock the scalability of healthcare and consumer technologies driven by biosignal data. This is because current methods cannot generalise well across human subjects’ data and handle learning from different heterogeneously collected data sets, thus limiting the scale of training data. On the other side, developments in transfer learning would benefit significantly from a real-world benchmark with immediate practical application. Therefore, we pick electroencephalography (EEG) as an exemplar for what makes biosignal machine learning hard. We design two transfer learning challenges around diagnostics and Brain-Computer-Interfacing (BCI), that have to be solved in the face of low signal-to-noise ratios, major variability among subjects, differences in the data recording sessions and techniques, and even between the specific BCI tasks recorded in the dataset. Task 1 is centred on the field of medical diagnostics, addressing automatic sleep stage annotation across subjects. Task 2 is centred on Brain-Computer Interfacing (BCI), addressing motor imagery decoding across both subjects *and* data sets. The BEETL competition with its over 30 competing teams and its 3 winning entries brought attention to the potential of deep transfer learning and combinations of set theory and conventional machine learning techniques to overcome the challenges. The results set a new state-of-the-art for the real-world BEETL benchmark.

Keywords: machine learning, transfer learning, domain adaptation, sleep diagnostics, Brain-Computer-Interfaces (BCI), EEG, neuroscience, NeurIPS2021

1. Introduction

The maturing of machine learning methods and their progressive deployments into the real-world brought to the forefront the need for combining similar data sets. Transfer learning has become a promising strategy to align different distributions [Pan and Yang \(2009\)](#). Transfer learning encompasses algorithms to transfer the representations and knowledge from source domains to a target domain. In the machine learning field, fine-tuning is often used to transfer model representations between tasks or data sets [Yosinski et al. \(2014\)](#). Another strategy is model splitting [Huang et al. \(2013a\)](#), which uses multiple sets of parameters for different data sets. Deep domain adaptation [Long et al. \(2015\)](#) is a technique which directly projects features of data sets into a common space with neural networks.

We pick electroencephalography (EEG), which is broadly considered one of the most promising ways to non-invasively read out of the human brain for diagnostic and human interfacing purposes. EEG reflects the features that make biosignal understanding a hard problem, as it consists of multi-dimensional time-series data that suffers from signal non-stationarities and poor signal-to-noise ratio, variability between users and sessions, different channel numbers and locations, as well as differences in task definitions between data sets. While EEG hardware has steadily evolved, trustworthy and data-efficient decoding methods are still missing. Currently, there are a few reviews [Wan et al. \(2021\)](#); [Wu et al. \(2020\)](#); [Jayaram et al. \(2016a\)](#) on different EEG transfer learning algorithms; and some recent inter-subject studies based on deep transfer learning approaches in Brain-Computer Interfaces (BCI) (e.g. [Wei et al. \(2021\)](#); [Li et al. \(2021a\)](#)) and in sleep studies (e.g. [Chambon et al. \(2018\)](#); [Andreotti et al. \(2018\)](#)). However, the field still lacks a systematic benchmark since these algorithms are tested on different data sets or with different pre-processing and setups in their studies. More importantly, very limited work focuses on cross-dataset transfer learning, which limits the use of big data in EEG decoding.

Two challenging tasks are designed to stimulate such algorithmic innovation. Task 1 is centred on transfer learning in the field of medical diagnostics, addressing automatic sleep stage annotation. The challenge lies in transferring from a large control population data set to clinically relevant cohorts with very little training data (transfer across subjects). Task 2 is centred on transfer learning for BCI, addressing motor imagery decoding. The challenge lies in transferring from multiple data sets, which use different EEG setups comprising hundreds of users, to a set of new users that need to be up and running with only minutes worth of calibration data (transfer across subjects *and* data sets). Currently, most studies focus on transfer learning within a single data set. The Benchmarks for EEG Transfer Learning (BEETL) competition provides a referencing platform for transfer learning strategies to two of the most common EEG applications and a guide on combining and utilising data sets from different sources, which could promote the EEG field towards the use of big data.

2. Task Description

Task 1 is in the field of medical diagnostics and specifically has the goal of automatic sleep stage annotation from sleep EEG data. We provide a data set with adult users (40 users, age 22-65) with 6 label categories for model training, based on which sleep stage annotation has to be transferred to two different age groups (65-80 and 80+) for each of which 5 subjects worth of data are provided. Task 1 is an essential use case for the development of ready-to-use medical diagnostics developed on a standard, large user base that has to be then transferred to many different clinically relevant subpopulations, for which respectively only a few subjects are worth of data can be collected. Beyond requiring subject-independence, the transfer has to work on different user groups (elderly and very elderly subjects) with well documented systematic EEG differences during sleep Landolt et al. (1996); Boselli et al. (1998); Landolt and Borbély (2001); Purdon et al. (2015).

Task 2 is a 3-way motor imagery classification challenge (left-hand, right-hand motor imagery and 'reject') that gets at the heart of the problem of current BCI systems: motor imagery data is exhausting for subjects to record, and historically has been difficult to use in a cross-subject and cross-dataset manner. Currently, there is limited work on cross-dataset transfer learning, and existing methods lack a systematic experimental comparison in the literature. This task provides a platform to compare the performance of current transfer learning algorithms across both subjects and data sets. In the past we organized several motor imagery data sets for BCI challenges in the MOABB (Mother of All BCI Benchmarks, <https://github.com/NeuroTechX/moabb>) database to test the performance of algorithms in terms of their generalisation performance on new data sets. Three source data sets are provided as training data. The algorithms are evaluated on new data sets with different setups, including differences in electrode channels, task definitions, and subjects. The test set contains an unpublished data set that is collected for this purpose and will be added to the MOABB database post-competition. Demonstration figures and more details of tasks can be found on the BEETL website (<https://beetl.ai/challenge>).

3. Data

Sleep Task data set. For Task 1, the sleep stage decoding task, the Physionet sleep data set [Kemp et al. \(2000\)](#); [Goldberger et al. \(2000\)](#) is one of the ideal data sets. The Sleep-EDF is a public database (<https://physionet.org/content/sleep-edfx/1.0.0/>) that contains 197 whole-night sleep recordings with event markers annotated by experts. Sleep patterns consist of sleep stages W, R, 1, 2, 3, 4, M (Movement time) and '?' (not scored). This data set has a clustered distribution of participants of different ages. The number of subjects is large enough for transfer learning algorithms to learn the diversity of distributions. We selected and randomized the subjects and trials in the competition to avoid cheating. Processed data could be found in the competition start kits.

Motor Imagery Task data set. For Task 2, the motor imagery decoding task, we selected three public data sets (Cho2017 [Cho et al. \(2017\)](#), BNCI2014 [Tangemann et al. \(2012\)](#) and PhysionetMI [Goldberger et al. \(2000\)](#); [Schalk et al. \(2004\)](#)) from the MOABB database as sources (<http://moabb.neurotechx.com/docs/datasets.html>). MOABB is a framework for evaluating BCI classification algorithms on publicly available data sets. We have collected a data set in an online racing game format in Cybathlon2020IC [Wei et al. \(2021\)](#) for testing purposes. Some offline subjects from the Weibo2014 data set [Yi et al. \(2014\)](#) and some online collected subjects from the Cybathlon2020IC data set are used as test samples. Data set information can be found in the table 1. Detailed description of the data could be found on the MOABB and BEETL websites. (<https://beetl.ai/data>).

Table 1: MI data sets in BEETL

MI Data set	Subjects	Channels	Tasks
Cho2017	52	64	Left/Right hand
BNCI2014	9	22	Left/Right hand/Feet/Tongue
PhysionetMI	109	64	Left/Right hand/Feet/Both hands/Rest
Weibo2014	10	60	Left/Right hand/Feet/Rest
Cybathlon2020IC	5	63	Left/Right hand/Feet/Rest

Training, validation and test set. For the Physionet sleep data set, we provide 80 sessions from 40 subjects (aged from 25-64) with full labels as source data and 5 subjects aged from 65 to 79 as examples of this age group; the performance of the algorithm is tested on more subjects aged from 65 to 79. Similarly, we provide 5 subjects aged from 80 to 95 with labels, while accuracies are reported on other subjects aged from 80 to 95. For the Physionet MI, Cho2017 and BNCI data sets, we provide full data sets with labels as sources. In both the Weibo2014 and Cybathlon2020IC data set, we provide some data with labels per test subject. During the competition, 32 channels around the motor cortex are selected from the Weibo2014 data set. The data set name was not provided during competition to avoid cheating. For the validation data (phase 1, leaderboard phase in the competition), two subjects from the Cybathlon2020IC data set and subjects 3, 4, and 5 of the Weibo2014 data set are used. For the test data (phase 2, final ranking phase in the competition), three subjects from the Cybathlon2020IC and subjects 7 and 9 from Weibo2014 are used as testing samples. Cybathlon2020IC will be added to MOABB post-competition.

Metrics. As both tasks in the BEETL challenge are classification problems, classification accuracy on the test data is the standard for ranking different solutions. To account for class imbalances in the data, balanced accuracy is computed by giving higher weight to classes with less samples Brodersen et al. (2010). The final competition score is the sum of the respective scores on both tasks.

4. Competition Results and Solutions

In task 1, winning teams used a single neural network for all subjects during training and testing based on DeepSleep Chambon et al. (2018) or EEGInception Santamaría-Vázquez et al. (2020). In task 2, the variability across different data sets, tasks, channel locations, and sampling rates make transfer more challenging (table 1). All three teams selected the common channels and aligned the sampling rate of different data sets. Team Cogitat used the Deep Sets Zaheer et al. (2017) framework with EEGInception to align latent distributions across different data sets and subjects. Team Wduong used Label Alignment He and Wu (2020) and Euclidean Alignment He and Wu (2019) to narrow the gap among tasks and data sets. After the alignment, team Wduong used an EEGNet Lawhern et al. (2018) backbone combined with a multi-task learning setup Ruder (2017) and the Maximum Classifier Discrepancy (MCD) Saito et al. (2018) method to perform classification. Team ms01 performed deep transfer learning with an architecture based on Riemannian geometry approaches Huang and Gool (2017). The winners and respective decoding scores on the two tasks can be found in Table 2 (more details available at <https://beetl.ai/prizes>).

Table 2: Method Accuracies of Top 5 Teams

	Task 1	Task 2	Overall
Cogitat	65.55	76.33	141.88
wduong	68.66	71.33	139.99
ms01	65.57	59.87	125.44
nik-sm	69.23	54.47	123.7
michaln	66.78	56.47	123.25

Baseline and Start kits A naive baseline without transfer learning was evaluated with the Shallow ConvNet from Schirrneister et al. (2017). In task 1, training data of the baseline method contains the source subjects and the example subjects from the test population. In task 2, evaluation was done on each subject with their own example data to avoid the negative transfer problem reported in Wei et al. (2021). 17 common channels of all data sets are used in the baseline of task 2. The baseline score in task 1 and task 2 is 57.6% and 49.9% respectively. Start kits for loading data, training models and generating labels are provided in Python (<https://beetl.ai/code>).

4.1. First Place Solution: Latent Subject Alignment

Team Cogitat: Stylianos Bakas, Siegfried Ludwig, Konstantinos Barmpas, Mehdi Bahri, Yannis Panagakis, Nikolaos Laskaris, Dimitrios A. Adamos, Stefanos Zafeiriou

Deep learning models based on the EEGInception architecture were used for both tasks [Santamaría-Vázquez et al. \(2020\)](#). The architecture applies some modifications to the widely used EEGNet architecture [Lawhern et al. \(2018\)](#), including temporal filterbanks with different kernel sizes, as well as additional deeper processing layers. The convolutional layers are followed by a linear classifier. All models are trained with the Adam gradient descent optimizer [Kingma and Ba \(2014\)](#), using default hyperparameters.

As the sleep classification task consists of subject-independent classification on a single data set, a straight-forward approach of training a model on the combined source and calibration data is performed. Latent alignment between subjects is used as described in later paragraphs.

Beyond aligning different subjects, the motor imagery task requires the transfer of trained models from one or more source data sets to perform classification on two target data sets. To make data set transfer simple, only the Physionet MI source data set was used and the 30 common electrodes between the source and the two target data sets were selected. The models are trained on the source and calibration data simultaneously, using a shared feature extractor and separate linear classifier heads per data set [Huang et al. \(2013b\)](#); [Wei et al. \(2021\)](#).

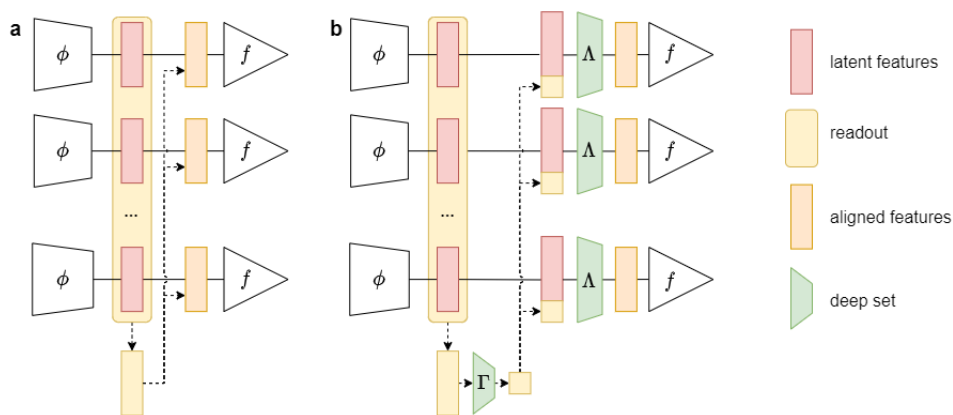


Figure 1: Latent alignment is performed following feature extractor ϕ and the classification function f is applied on the aligned features. a) Statistical alignment standardizes latent distributions of each subject. b) Deep Set alignment uses a trainable function Γ to obtain a distribution embedding and update the features of each trial with a trainable function Λ . Source: [Bakas et al. \(2022\)](#). Solution proposed by team Cogitat.

The approach for improving subject-independence is to align the latent feature distributions of the deep learning models between different subjects and sessions [Bakas et al. \(2022\)](#). This is similar to the Euclidean Alignment method, which performs spatial whitening with the average covariance matrix per subject [He and Wu \(2019\)](#), although the methods are not constrained to be applied on the model input. Performing distribution alignment in

the classifier latent space is related to Riemannian manifold methods [Zanini et al. \(2017\)](#), which have gained traction in recent years, but integration with deep learning models has not yet matured.

For the first implementation of latent alignment, team Cogitat developed a novel statistical alignment method, which standardizes latent feature distributions per subject and can be seen as a subject-wise batch normalization [Li et al. \(2021b\)](#), although a separate module per subject is not used (figure 1a). This method is very computationally efficient, straightforward to implement and can be used in the place of the batch normalization layers in a deep learning model. During training, multiple trials need to be present for each subject in the batch, in order to estimate feature distributions. During inference, a statistical estimation obtained on unlabelled trials of the target subject can be used [Xu et al. \(2021\)](#). This worked well on the sleep classification task, where no subject-wise labelled calibration data is available.

Taking the perspective of classification on a set of EEG trials from a given subject, the application of Deep Set architectures follows naturally as a second implementation of latent alignment [Zaheer et al. \(2017\)](#). Team Cogitat developed a deep learning layer that uses the statistical mean to obtain a distribution embedding of the given subject to update trial-wise features (figure 1b). This worked particularly well on the MI task, which provides subject-wise calibration data.

4.2. Second Place Solution: Multi-Source EEGNet with Domain and Label Adaptation

Team `wduong`: William C. Duong, Stephen M. Gordon and Vernon J. Lawhern

Team Wduong used a modified version of DeepSleep in Task 1 [Chambon et al. \(2018\)](#) with some minor changes, including adjusting the convolution kernel size, adding batch normalization after each layer, and setting up a pre-training scheme. For the motor imagery decoding (Task 2), team Wduong used a combination of data alignment, multi-task EEGNet model, and maximum classifier discrepancy domain adaptation to train a robust model. The code is available at https://github.com/mcd4874/NeurIPS_competition.

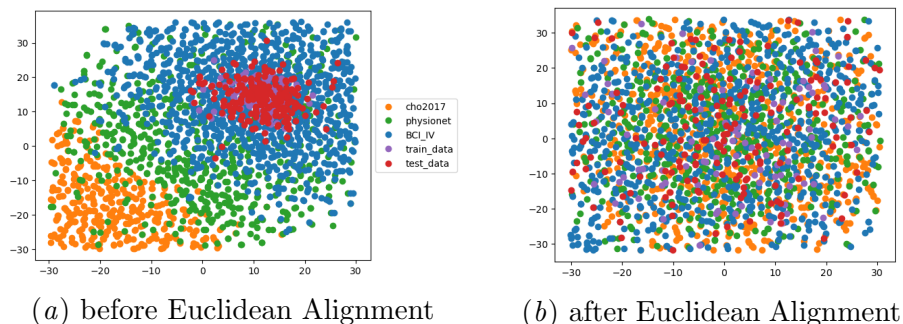


Figure 2: Data sets alignment with Euclidean Alignment. Solution proposed by team Wduong.

Due to limited available labeled trials, the motor imagery decoding task needs transfer learning strategies to train a robust model with good performance. To create a generalized framework to handle multi-source data set transfer learning, all three provided data sets, included Cho2017, Physionet, and BNCI2014, are used in the training procedure. 17 common channels among three source data sets and two target data sets are utilised. To demonstrate the generalization of the approach, the tongue category in the BNCI2014 data sets is kept, which is not available in two target data sets. Two main problems exist among source data sets and target data sets. First, there is a category gap problem due to different label categories between source and target data sets. Label Alignment is applied to solve this problem [He and Wu \(2020\)](#). For example, the tongue motor imagery category can be projected to the rest category to align corresponding distributions; therefore, the training phase can use source data with different labels compared to target data. In addition to the label gap problem, other data sets with different experimental recording setups led to a domain gap between subjects from source data sets and target data sets. Euclidean Alignment is used to close the gap between subjects in source domains and target domain [He and Wu \(2019\)](#). All trials of each subject are aligned such that the mean covariance matrices of all subjects are equal to the identity matrix after alignment. After alignment, data set distribution in both source domains and target domain are more similar as seen in figure 2.

Finally, a multi-task EEGNet combined with Maximum Classifier Discrepancy (MCD) [Saito et al. \(2018\)](#) is developed to solve the motor-imagery decoding problem. A multi-task learning (MTL) setup can learn a common EEGNet backbone to increase feature representation generalization among data sets, where MTL treats each data set from source domains and target domain as an individual task [Ruder \(2017\)](#). Combining losses with a weighted sum can optimize the shared EEGNet backbone and all the classifiers jointly [Lawhern et al. \(2018\)](#). An adversarial training procedure in the MCD method is conducted between the target classifier and EEGNet backbone to both increase the discriminative ability of target classifiers and enhance the feature extraction of the backbone [Saito et al. \(2018\)](#). The model architecture is in Figure 3.

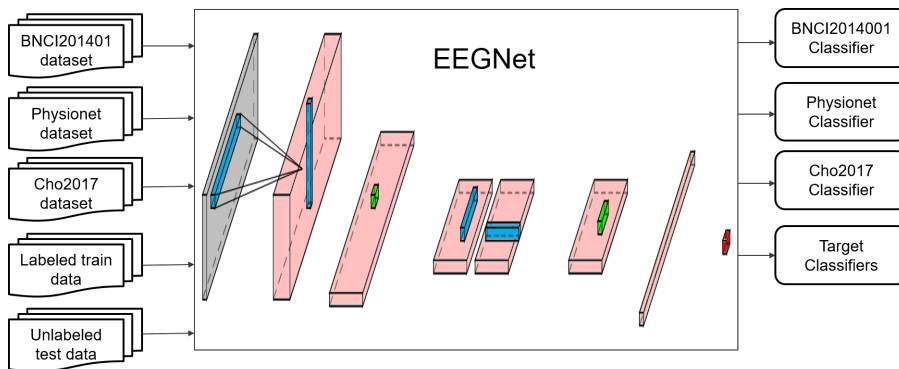


Figure 3: Multi-task EEGNet. Solution proposed by team Wduong.

Each targeted subject is treated as a separate target domain for the training strategy due to the highly subject-specific nature of motor imagery EEG signals. Therefore, five

model groups known as A_0, A_1, A_2, B_0 , and B_1 were created. Five-fold block-wise cross-validation is applied for each model group, such that 80% of the target domain data is used to train a model and 20% is used for validation to pick the best model. An ensemble learning strategy combines five-fold models via majority vote to learn the final prediction. During the training phase, the target data set samples a batch size of 16, while each source data set samples a batch size of 64 for each mini-batch to train the model. An Adam optimizer with the learning of 0.001 is used to train the model for 20 epochs.

4.3. Third place solution: Classification of covariance matrices with SPDNet

Team ms01: Maciej Śliwowski, Vincent Rouanne, Piotr Tempczyk

Team ms01 performed different methods for task 1 and task 2 because of the difference in the available number of EEG channels. In the sleep decoding task (task 1), a DeepSleep model [Chambon et al. \(2018\)](#) was trained and evaluated without major modifications. For the motor imagery decoding task (task 2), methods analyzing spatial covariance matrices of the signal was studied with a focus on Riemannian geometry approaches.

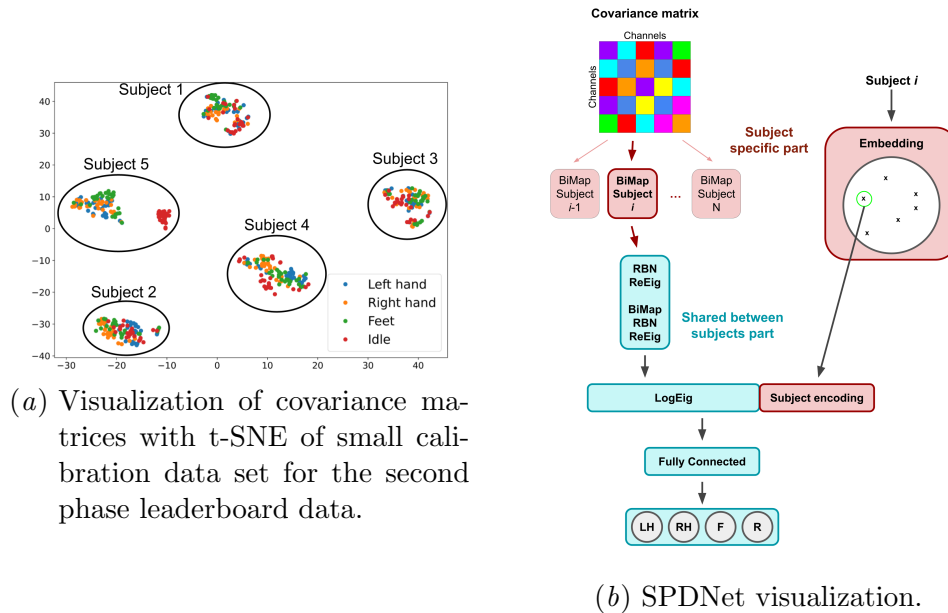


Figure 4: Solution proposed by team ms01.

The approach performs classification of the spatial covariance matrices for different motor imagery (MI) tasks with Minimum Distance to Riemannian Mean (MDRM) classifiers [Barachant et al. \(2012\)](#). Covariance matrices lie on a manifold of positive semi-definite matrices in which the distance between data points can be computed using Riemannian geometry. The Riemannian mean of covariance matrices is computed separately for each motor imagery task. New observations are classified using the minimal distance between the observation’s covariance matrix and the tasks mean covariance matrices. To visualize the

data distribution and interpret the results obtained with the MDRM classifier, t-Stochastic Neighbor Embedding (t-SNE) [van der Maaten and Hinton \(2008\)](#) with Riemannian distance are used (Figure 4a).

Based on the MDRM results and the t-SNE visualization (Figure 4a), the solution is proposed using nonlinear methods designed specifically for the analysis of symmetric positive definite matrices (SPD). SPD networks, similar to standard deep learning, are able to learn an efficient representation of the data that is not linearly separable as in our case. Covariance matrices are classified using SPDNet [Huang and Gool \(2017\)](#), a type of Symmetric Positive Definite manifold networks, together with Riemannian batch normalization (RBN) [Brooks et al. \(2019\)](#) (code from [Brooks et al. \(2019\)](#)). All subjects from all sources and leaderboard labelled data sets are combined to form the training and validation data sets (30% of the leaderboard labelled data set was used as validation).

In the proposed SPDNet architecture for transfer learning (Figure 4b), data for each subject is processed by an individual feature extractor consisting of one bilinear transformation (BiMap layer) [Huang and Gool \(2017\)](#). As further steps in the model are shared between subjects, the first BiMap layer is able to unify the representations and make it subject invariant. The shared part of the network consists of RBN and rectified eigenvalues activation (ReEig) [Huang and Gool \(2017\)](#) followed by another block of BiMap, RBN, and ReEig. Extracted features are transformed from Riemannian space into Euclidean space (LogEig layer). Two fully connected layers (with a dropout layer in-between) are added to predict four MI classes. To catch between-subject variance, an embedding/encoding of patients of size fifteen is created, which is concatenated to the first fully connected euclidean layer. The embedding is trained together with the whole network with backpropagation. A hypothesis is that the subject encoding could ease features standardization across subjects.

5. Discussion

In this section, we will discuss key observations from the BEETL competition. This section includes the reasoning for choices of algorithms, the commonalities and differences among solutions, limitations of the winning methods and task design, and future directions.

Architecture design is a challenge for cross-dataset EEG transfer learning. In light of the relatively large data sets provided, team Cogitat chose the EEGInception architecture for both tasks, which has enough capacity to learn discriminative features across subjects and data sets. Both team Wduong and ms01 used a modified DeepSleep model for task 1. In task 2, team Wduong used a multi-task EEGNet setup. Team ms01 used the SPDNet architecture that analyzes covariance matrices with nonlinear transformations. Despite the diversity of algorithms, all teams used deep learning approaches. The choice of using neural networks might be explained by the need for increased model capacity when attempting to find classifiers that can generalize across many subjects and data sets. Further, deep learning architectures allow for end-to-end training of various processing stages, including representation learning without hand-crafted features and distribution alignment. This provides flexibility in the design of architectures to test out different transfer learning setups.

Another commonality is that all employed a combination of shared network backbones and individualized processing layers for different data sets or subjects (either in shallow layers or classification layers). Team Cogitat and Wduong used a shared feature extractor

with different classifiers for different data sets. Team ms01 had unique layers for each subject at shallow layers of their architecture. The splitting approach allows model backbones to be trained on large amounts of data while allowing for adaptation to the existing differences.

EEG decoding suffers from substantial covariate shifts among different subjects and data sets. All three teams used alignment approaches to reduce differences among data sets, subjects and sessions. Team Cogitat performed statistical alignment of latent features at various stages of the model. Team Wduong performed Euclidean alignment on the EEG trials based on their covariances and an additional label alignment approach. Team ms01 analyzed covariance matrices with individual per-subject layers to unify the representations. Although the exact implementations vary, the winning solutions indicate the utility of alignment methods to perform domain adaptation in EEG transfer learning methods.

How to select and process source data sets is a key question for EEG transfer learning towards the use of big data. In task 2, team Wduong and team ms01 used all three source data sets with 17 common channels to utilise more subjects and tasks. Team Cogitat used one source data set (PhysionetMI) with more electrodes. Here we highlight a trade-off between keeping more common electrode channels or increasing the number of subjects by adding more data sets. One limitation of the winning solutions is that they either abandoned some channels or some subjects. Discriminative information could be lost by this selection approach. Therefore, aiming to apply some form of upsampling approach or building a better architecture to use all channels of data sets would be possible future directions.

We observe that the gaps of decoding accuracies between participants are relatively small in the sleep task. Most teams did not utilise complex domain adaptation strategies. This may be caused by various reasons. Firstly, sleep EEG could be relatively similar across subjects by nature compared to motor imagery. Secondly, the sleep task is designed on only a single data set with two channels, which simplified the scenario for medical diagnostics. Another human factor could be that sleep data is manually labelled by human experts which unified the data in terms of classes. More studies on utilising multiple sleep data sets with more channels and setups are needed in the literature to conclude a precise explanation.

There are various directions that require more studies. In the competition, backbone architectures are based on similar ideas of temporal and spatial filters [Lawhern et al. \(2018\)](#); [Santamaría-Vázquez et al. \(2020\)](#). More architectures in the EEG field should be explored for multi-dataset transfer learning, including conventional transfer learning algorithms [Jayaram et al. \(2016b\)](#). Winning solutions used various alignment approaches, whereas a systematical comparison of these alignment approaches in the same experiment setup is required. In addition, the winning methods either leave out some channels or source tasks in their approaches. Therefore, algorithms that could use the full potential of sources should be developed in the future. Furthermore, in light of the fact that all winning solutions used deep learning, the interpretability [Goebel et al. \(2018\)](#); [Samek et al. \(2019\)](#) of features used for deep transfer learning needs to be further investigated.

To conclude, the BEETL competition brought attention to the challenge of using large-scale biosignal transfer learning on a large amount of subjects and data sets. It provides a common platform to evaluate transfer learning algorithms and some valuable examples in EEG transfer learning with insights. We hope the BEETL competition could open up the opportunity to utilise heterogenous open-source data sets and move the EEG field forward.

Acknowledgments

Our competition is officially affiliated with the BCI Society. We thank Facebook Reality Labs for sponsoring 5000\$ in competition prizes. Many thanks to NeurIPS2021 Organisers and all researchers who devote themselves into this competition. We thank all organisers and teams in BEETL. Aldo Faisal acknowledges a UKRI Turing AI Fellowship Grant (EP/V025449/1). Alexandre Gramfort acknowledges the support of ANR BrAIN AI chair (ANR-20-CHIA-0016). Maciej Śliwowski was supported by the CEA NUMERICS program, which has received funding from European Union’s Horizon 2020 research and innovation program under the Marie Skłodowska-Curie grant agreement No 800945.

References

- Fernando Andreotti, Huy Phan, Navin Cooray, Christine Lo, Michele TM Hu, and Maarten De Vos. Multichannel sleep stage classification and transfer learning using convolutional neural networks. In *2018 40th Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC)*, pages 171–174. IEEE, 2018.
- Stylios Bakas, Siegfried Ludwig, Konstantinos Barmpas, Mehdi Bahri, Yannis Panagakis, Nikolaos Laskaris, Dimitrios A Adamos, and Stefanos Zafeiriou. Team cogitat at neurips 2021: Benchmarks for eeg transfer learning competition. *arXiv preprint arXiv:2202.03267*, 2022.
- Alexandre Barachant, Stéphane Bonnet, Marco Congedo, and Christian Jutten. Multiclass brain–computer interface classification by riemannian geometry. *IEEE Transactions on Biomedical Engineering*, 59(4):920–928, 2012. doi: 10.1109/TBME.2011.2172210.
- Mirella Boselli, Liborio Parrino, Arianna Smerieri, and Mario Giovanni Terzano. Effect of age on eeg arousals in normal sleep. *Sleep*, 21(4):361–367, 1998.
- Kay Henning Brodersen, Cheng Soon Ong, Klaas Enno Stephan, and Joachim M Buhmann. The balanced accuracy and its posterior distribution. In *2010 20th international conference on pattern recognition*, pages 3121–3124. IEEE, 2010.
- Daniel Brooks, Olivier Schwander, Frederic Barbaresco, Jean-Yves Schneider, and Matthieu Cord. Riemannian batch normalization for spd neural networks. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- Stanislas Chambon, Mathieu N Galtier, Pierrick J Arnal, Gilles Wainrib, and Alexandre Gramfort. A deep learning architecture for temporal sleep stage classification using multivariate and multimodal time series. *IEEE Transactions on Neural Systems and Rehabilitation Engineering*, 26(4):758–769, 2018.
- Hohyun Cho, Minkyu Ahn, Sangtae Ahn, Moonyoung Kwon, and Sung Chan Jun. EEG datasets for motor imagery brain–computer interface. *GigaScience*, 6(7):gix034, 2017.

- Randy Goebel, Ajay Chander, Katharina Holzinger, Freddy Lecue, Zeynep Akata, Simone Stumpf, Peter Kieseberg, and Andreas Holzinger. Explainable ai: the new 42? In *International cross-domain conference for machine learning and knowledge extraction*, pages 295–303. Springer, 2018.
- Ary L Goldberger, Luis AN Amaral, Leon Glass, Jeffrey M Hausdorff, Plamen Ch Ivanov, Roger G Mark, Joseph E Mietus, George B Moody, Chung-Kang Peng, and H Eugene Stanley. Physiobank, physiotoolkit, and physionet: components of a new research resource for complex physiologic signals. *circulation*, 101(23):e215–e220, 2000.
- He He and Dongrui Wu. Transfer learning for brain–computer interfaces: A euclidean space data alignment approach. *IEEE Transactions on Biomedical Engineering*, 67(2):399–410, 2019.
- He He and Dongrui Wu. Different set domain adaptation for brain-computer interfaces: A label alignment approach. *IEEE Transactions on Neural Systems and Rehabilitation Engineering*, 28(5):1091–1108, 2020.
- Jui-Ting Huang, Jinyu Li, Dong Yu, Li Deng, and Yifan Gong. Cross-language knowledge transfer using multilingual deep neural network with shared hidden layers. In *2013 IEEE Intl. Conf. on Acoustics, Speech and Signal Processing*, pages 7304–7308. IEEE, 2013a.
- Jui-Ting Huang, Jinyu Li, Dong Yu, Li Deng, and Yifan Gong. Cross-language knowledge transfer using multilingual deep neural network with shared hidden layers. In *2013 IEEE International Conference on Acoustics, Speech and Signal Processing*, pages 7304–7308. IEEE, 2013b.
- Zhiwu Huang and Luc Van Gool. A riemannian network for spd matrix learning. In *Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence, AAAI’17*, page 2036–2042. AAAI Press, 2017.
- Vinay Jayaram, Morteza Alamgir, Yasemin Altun, Bernhard Scholkopf, and Moritz Grosse-Wentrup. Transfer learning in brain-computer interfaces. *IEEE Computational Intelligence Magazine*, 11(1):20–31, 2016a.
- Vinay Jayaram, Morteza Alamgir, Yasemin Altun, Bernhard Scholkopf, and Moritz Grosse-Wentrup. Transfer learning in brain-computer interfaces. *IEEE Comp. Intel. Magazine*, 11(1):20–31, 2016b.
- Bob Kemp, Aeilko H Zwinderman, Bert Tuk, Hilbert AC Kamphuisen, and Josefien JL Oberye. Analysis of a sleep-dependent neuronal feedback loop: the slow-wave microcontinuity of the EEG. *IEEE Transactions on Biomedical Engineering*, 47(9):1185–1194, 2000.
- Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Hans-Peter Landolt and Alexander A Borbély. Age-dependent changes in sleep eeg topography. *Clinical Neurophysiology*, 112(2):369–377, 2001.

- Hans-Peter Landolt, Derk-Jan Dijk, Peter Achermann, and Alexander A Borbély. Effect of age on the sleep eeg: slow-wave activity and spindle frequency activity in young and middle-aged men. *Brain research*, 738(2):205–212, 1996.
- Vernon J Lawhern, Amelia J Solon, Nicholas R Waytowich, Stephen M Gordon, Chou P Hung, and Brent J Lance. EEGNet: a compact convolutional neural network for EEG-based brain–computer interfaces. *Journal of neural engineering*, 15(5):056013, 2018.
- Denghao Li, Pablo Ortega, Xiaoxi Wei, and A.A. Faisal. Meta-Learning EEG Motor Imagery Decoder for Brain Computer Interface. *IEEE Neural Engineering (NER)*, 10, 2021a.
- Ruilin Li, Lipo Wang, and Olga Sourina. Subject matching for cross-subject eeg-based recognition of driver states related to situation awareness. *Methods*, 2021b.
- Mingsheng Long, Yue Cao, Jianmin Wang, and Michael I Jordan. Learning transferable features with deep adaptation networks. *arXiv preprint arXiv:1502.02791*, 2015.
- Sinno Jialin Pan and Qiang Yang. A survey on transfer learning. *IEEE Transactions on knowledge and data engineering*, 22(10):1345–1359, 2009.
- PL Purdon, KJ Pavone, O Akeju, AC Smith, AL Sampson, J Lee, DW Zhou, K Solt, and EN Brown. The ageing brain: age-dependent changes in the electroencephalogram during propofol and sevoflurane general anaesthesia. *British journal of anaesthesia*, 115 (suppl.1):i46–i57, 2015.
- Sebastian Ruder. An overview of multi-task learning in deep neural networks, 2017.
- Kuniaki Saito, Kohei Watanabe, Yoshitaka Ushiku, and Tatsuya Harada. Maximum classifier discrepancy for unsupervised domain adaptation. *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.
- Wojciech Samek, Grégoire Montavon, Andrea Vedaldi, Lars Kai Hansen, and Klaus-Robert Müller. *Explainable AI: interpreting, explaining and visualizing deep learning*, volume 11700. Springer Nature, 2019.
- Eduardo Santamaría-Vázquez, Víctor Martínez-Cagigal, Fernando Vaquerizo-Villar, and Roberto Hornero. Eeg-inception: A novel deep convolutional neural network for assistive erp-based brain-computer interfaces. *IEEE Transactions on Neural Systems and Rehabilitation Engineering*, 28(12):2773–2782, 2020.
- Gerwin Schalk, Dennis J McFarland, Thilo Hinterberger, Niels Birbaumer, and Jonathan R Wolpaw. Bci2000: a general-purpose brain-computer interface (bci) system. *IEEE Transactions on biomedical engineering*, 51(6):1034–1043, 2004.
- Robin Tibor Schirrmeister, Jost Tobias Springenberg, Lukas Dominique Josef Fiederer, Martin Glasstetter, Katharina Eggenberger, Michael Tangermann, Frank Hutter, Wolfram Burgard, and Tonio Ball. Deep learning with convolutional neural networks for EEG decoding and visualization. *Human brain mapping*, 38(11):5391–5420, 2017.

- Michael Tangermann, Klaus-Robert Müller, Ad Aertsen, Niels Birbaumer, Christoph Braun, Clemens Brunner, Robert Leeb, Carsten Mehring, Kai J Miller, Gernot Mueller-Putz, et al. Review of the bci competition iv. *Frontiers in neuroscience*, 6:55, 2012.
- Laurens van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of Machine Learning Research*, 9(86):2579–2605, 2008.
- Zitong Wan, Rui Yang, Mengjie Huang, Nianyin Zeng, and Xiaohui Liu. A review on transfer learning in EEG signal analysis. *Neurocomputing*, 421:1–14, 2021.
- Xiaoxi Wei, Pablo Ortega, and A Aldo Faisal. Inter-subject deep transfer learning for motor imagery eeg decoding. In *2021 10th International IEEE/EMBS Conference on Neural Engineering (NER)*, pages 21–24. IEEE, 2021.
- Dongrui Wu, Yifan Xu, and Bao-Liang Lu. Transfer learning for EEG-based brain-computer interfaces: A review of progress made since 2016. *IEEE Transactions on Cognitive and Developmental Systems*, 2020.
- Lichao Xu, Zhen Ma, Jiayuan Meng, Minpeng Xu, Tzyy-Ping Jung, and Dong Ming. Improving transfer performance of deep learning with adaptive batch normalization for brain-computer interfaces. In *2021 43rd Annual International Conference of the IEEE Engineering in Medicine & Biology Society (EMBC)*, pages 5800–5803. IEEE, 2021.
- Weibo Yi, Shuang Qiu, Kun Wang, Hongzhi Qi, Lixin Zhang, Peng Zhou, Feng He, and Dong Ming. Evaluation of eeg oscillatory patterns and cognitive process during simple and compound limb motor imagery. *PloS one*, 9(12):e114853, 2014.
- Jason Yosinski, Jeff Clune, Yoshua Bengio, and Hod Lipson. How transferable are features in deep neural networks? In *Adv. in neural information processing systems (NIPS)*, pages 3320–3328, 2014.
- Manzil Zaheer, Satwik Kottur, Siamak Ravanbakhsh, Barnabas Poczos, Ruslan Salakhutdinov, and Alexander Smola. Deep sets. *arXiv preprint arXiv:1703.06114*, 2017.
- Paolo Zanini, Marco Congedo, Christian Jutten, Salem Said, and Yannick Berthoumieu. Transfer learning: A riemannian geometry framework with applications to brain-computer interfaces. *IEEE Transactions on Biomedical Engineering*, 65(5):1107–1116, 2017.