

1 Minimizing Calibration Time for Brain Reading 1

2 Anonymous DAGM submission 2

3 Paper ID 006 3

4 **Abstract.** Machine learning is increasingly used to autonomously adapt 4
5 brain-machine interfaces to user-specific brain patterns. In order to min- 5
6 imize the preparation time of the system, it is highly desirable to reduce 6
7 the length of the calibration procedure, during which training data is 7
8 acquired from the user, to a minimum. One recently proposed approach 8
9 is to reuse models that have been trained in historic usage sessions of 9
10 the same or other users by utilizing an ensemble-based approach. In 10
11 this work, we propose two extensions of this approach which are based 11
12 on the idea to combine predictions made by the historic ensemble with 12
13 session-specific predictions that become available once a small amount of 13
14 training data has been collected. These extensions are particularly use- 14
15 ful for *Brain Reading Interfaces* (BRIs), a specific kind of brain-machine 15
16 interfaces. BRIs do not require that user feedback is given and thus, 16
17 additional training data may be acquired concurrently to the usage ses- 17
18 sion. Accordingly, BRIs should initially perform well when only a small 18
19 amount of training data acquired in a short calibration procedure is 19
20 available and allow an increased performance when more training data 20
21 becomes available during the usage session. An empirical offline-study in 21
22 a testbed for the use of BRIs to support robotic telemanipulation shows 22
23 that the proposed extensions allow to achieve this kind of behavior. 23

24 1 Introduction 24

25 Brain Reading Interfaces (BRIs) are one particular kind of brain-machine 25
26 interface (BMI) that allow to provide the machine with information about the 26
27 current mental state and intent of its user such that the machine can optimize 27
28 its behavior accordingly. In contrast to active Brain-Computer Interfaces (BCIs, 28
29 see [3, 13] for a review of works), BRIs estimate the user’s mental state and intent 29
30 based on passive, external observation of brain activity without requiring any 30
31 active participation of the user. This observation can, e.g., be based on electroen- 31
32 cephalography (EEG). Since no active participation of the user is required, BRIs 32
33 are well-suited for scenarios like robotic telemanipulation where a sophisticated 33
34 BMI is expedient but the user needs to be fully immersed in his task. 34

35 Like active BCIs, BRIs must be adapted to the current brain patterns of the 35
36 user since these characteristic patterns vary between different subjects and even 36
37 change over time within the same subject. This can be achieved by using machine 37
38 learning (ML) techniques (see, e.g., Blankertz et al. [4] for an example in an active 38
39 BCI). The common approach for using ML in BCIs is to record labeled training 39
40 data during a so-called calibration procedure that must be conducted prior to 40

41 each usage session. In this calibration procedure, the user acts in a controlled 41
 42 and supervised scenario. The labeled data acquired is then used to adapt the 42
 43 ML-based BCI system to the user’s current brain patterns. The drawback of this 43
 44 approach is that the user has to conduct this calibration procedure each time 44
 45 he wants to use the system. Thus, it is highly desirable to keep this calibration 45
 46 procedure as short as possible (or remove its necessity altogether). 46

47 Different approaches for reducing the calibration time have been proposed: 47
 48 Krauledat et al. [9] proposed an algorithm targeted at long-term BCI users that 48
 49 allows to skip the calibration procedure. This is accomplished by inferring spatial 49
 50 filters and classifiers that generalize well across sessions based on reusing 50
 51 training data from historic sessions of the same user and clustering of historic 51
 52 spatial filters. Fazli et al. [6] proposed a method that allows to skip the cali- 52
 53 bration procedure for both long-term and novel users. Their approach is based 53
 54 on an ensemble of historic spatial-filter/classifier combinations that are trans- 54
 55 ferred to the current session and whose individual predictions are combined into 55
 56 a joint prediction by means of a gating function. Both approaches require that 56
 57 a large number of historic sessions be available. Further approaches for reducing 57
 58 calibration time are multi-task learning [2], semi-supervised learning [10], and a 58
 59 hybrid approach that mixes historic data with session-specific data [11]. 59

60 The main contribution of this paper is to propose two extensions of the “pure” 60
 61 ensemble-based approach of Fazli et al. and to present an empirical comparison 61
 62 of these approaches in a testbed for the use of BRIs to support robotic telema- 62
 63 nipulation. The two extensions we propose—Classification Augmentation and 63
 64 Feature Augmentation—are based on the idea of combining the predictions made 64
 65 by the historic ensemble with session-specific predictions that become available 65
 66 once some amount of training data has been collected. We show that these ex- 66
 67 tensions achieve good performance when only a small amount of training data 67
 68 is available and—in contrast to the “pure” ensemble approach—also become 68
 69 increasingly better for more training data. This is particularly important for BRIs, 69
 70 since BRIs allow to interweave the acquisition of training data with the actual 70
 71 usage session. Thus, the system should initially perform well based on a small 71
 72 amount of training data acquired in a short calibration procedure but should also 72
 73 be able to improve performance when increasingly more training data is gathered 73
 74 during the usage session. Furthermore, in contrast to related approaches like [6] 74
 75 and [9], the proposed extensions perform well also when only a small number 75
 76 of historic sessions is available. The paper is structured as follows: In Section 2, 76
 77 a testbed for BRIs in robotic telemanipulation is presented. Subsequently, the 77
 78 baseline BRI as well as different ensemble-based extensions are proposed in Sec- 78
 79 tion 3. In Section 4, the experimental setup and a discussion of our results are 79
 80 given and a conclusion is drawn in Section 5. 80

81 2 Scenario 81

82 *Labyrinth Oddball* The empirical evaluation was conducted on an EEG dataset 82
 83 recorded in the Labyrinth Oddball scenario (see Figure 1), a testbed for the 83

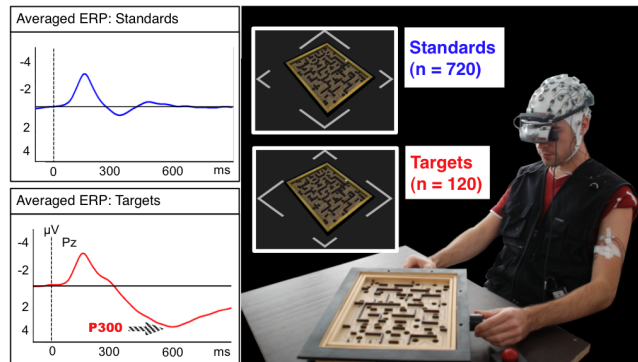


Fig. 1. Labyrinth Oddball: The subject plays a physical simulation of the BRIO[®] labyrinth and has to respond to rare 'target' stimuli by pressing a buzzer. Event-related potentials (ERPs) evoked by 'target' and more frequent 'standard' stimuli are depicted.

84 use of BRIs in robotic telemanipulation. In this testbed, the operator has to 84
 85 simultaneously execute a manipulation task (playing the Labyrinth game) and 85
 86 to distinguish two different kinds of stimuli presented to him while playing the 86
 87 game. The BRI only needs to passively monitor whether the operator of the 87
 88 Labyrinth game correctly recognized and distinguished these stimuli. Since no 88
 89 user feedback is given, the testbed is well suited for evaluation of BRIs (for more 89
 90 details we refer to [1] and the video in the supplementary material). The task 90
 91 for the BRI is to discriminate between the EEG patterns evoked by recognizing 91
 92 so-called 'standard' and 'target' stimuli¹. While 'standard' stimuli are frequent 92
 93 (720 presentations per run) but irrelevant for the user, 'target' stimuli are rare 93
 94 (120 presentations per run) and require him to press a buzzer. Such a scenario is 94
 95 called "oddball discrimination paradigm" and the successful recognition of the 95
 96 rare 'target' stimuli is known to elicit an event-related potential (ERP) called 96
 97 P300 [12]. In contrast to many active BCIs (e.g., [13]), the classification has to 97
 98 be made based on the individual instance and not on an average over several 98
 99 repetitions of the same condition. To avoid differences in early visual brain 99
 100 activity and to make sure that differences in the EEG recorded and classified after 100
 101 the presentation of both stimuli types are actually due to higher cognitive pro- 101
 102 cessing, the visual presentation (shape and color) of standard and target stimuli 102
 103 was kept very similar. Note that neither during the calibration procedure nor 103
 104 during evaluation runs feedback was given to the subject. 104

105 *Data Acquisition* EEG data was acquired in 12 sessions from 6 male subjects; 105
 106 each subject performed 2 sessions. Sessions were recorded on different days; 106
 107 accordingly, the EEG cap was fitted onto the subject's head for each session 107

¹ This is a kind of proxy-task for the actual task of distinguishing between recognized and missed target stimuli (see [1] for a discussion).

108 anew. Each of these sessions consisted of five repetitions (called “runs”) of the 108
 109 Labyrinth Oddball paradigm. After each of the five runs there was a short break 109
 110 of 10 minutes. The EEG was recorded and stored along with information about 110
 111 which stimulus was presented at what time and whether the buzzer was pressed 111
 112 afterwards. EEG was recorded continuously from 64 electrodes (extended 10–20 112
 113 system with reference at electrode FCz), using an actiCap system (Brain Prod- 113
 114 ucts GmbH, Munich, Germany). Two of the 64 channels (replacing the electrodes 114
 115 TP7 and TP8) were used to record electromyography signals of muscles of the 115
 116 lower arm and have been discarded in this study. EEG signals were amplified 116
 117 by two 32 channel BrainAmp DC amplifiers (Brain Products GmbH, Munich, 117
 118 Germany) and were sampled at 1000 Hz. The impedance was kept below 5 k Ω . 118

119 3 Methods 119

120 *Baseline BRI* As a first step of the baseline BRI system used for discrimination of 120
 121 the ‘standard’ and the ‘target’ condition, rectangular time windows starting 0 ms 121
 122 and ending 1000 ms after stimulus presentation are extracted from the continuous 122
 123 signal recorded during the experiment. Thereupon, the extracted time windows 123
 124 are normalized so that the mean value of each channel becomes 0 within this 124
 125 window. Subsequently, the signal is low-pass filtered (cutoff frequency 12 Hz), 125
 126 downsampled from 1000 Hz to 25 Hz, and again low-pass filtered for a cutoff 126
 127 frequency of 4 Hz in order to focus on slow ERPs like the P300. 127

128 After this, the signal is spatially filtered. Spatial filtering denotes a mapping 128
 129 of the original n channels $x(t)$ (that correspond one-to-one to the n electrodes) 129
 130 onto new pseudo-channels $\hat{x}(t) = W^T x(t)$ that are a (linear) mixture of the 130
 131 signals recorded at different electrodes (see Blankertz et al. [5] for a discussion 131
 132 of why spatial filtering is an important step). In this work, we have generated 132
 133 spatial filters based on the common spatial patterns (CSP) algorithm [8]. CSP 133
 134 maps the data onto axes such that the variance for instances of the first class 134
 135 is maximized and the variance for the second class is minimized (or vice versa). 135
 136 This is achieved by a simultaneous diagonalization of the two empirical intra- 136
 137 class covariance matrices $\Sigma_1 = n_1^{-1} \sum_{i=1}^{n_1} (x_i^{(1)} - \mu^{(1)})(x_i^{(1)} - \mu^{(1)})^T$ and $\Sigma_2 =$ 137
 138 $n_2^{-1} \sum_{i=1}^{n_2} (x_i^{(2)} - \mu^{(2)})(x_i^{(2)} - \mu^{(2)})^T$, i.e., by solving $\Sigma_1 W = \Lambda \Sigma_2 W$ where Λ is the 138
 139 vector of generalized eigenvalues and W is the matrix of generalized eigenvectors 139
 140 corresponding to the learned filters. 140

141 The values of the resulting pseudo-channels, i.e., the 26×62 samples of 141
 142 the 62 pseudo-channels that fall into the time window from 0 to 1000 ms, are 142
 143 used as features. Thereupon, each feature dimension is normalized such that its 143
 144 2.5th percentile on the training data is mapped onto 0 and the 97.5th percentile 144
 145 is mapped onto 1. The resulting feature vectors are classified using a support 145
 146 vector machine (SVM) with linear kernel and complexity 0.01. Since the ratio 146
 147 of standard and target class instances in the dataset is highly unbalanced due 147
 148 to the oddball paradigm, the weight for class ‘target’ has been set to 2.0, while 148
 149 the weight of class ‘standard’ was set to 1.0. The feature set and all mentioned 149
 150 parameters have been chosen based on a preliminary investigation conducted on 150

151 a hold-out dataset. The implementation of the data processing system is based 151
 152 on the “Modular toolkit for Data Processing” [14]. 152

153 *Ensemble approach* The baseline BRI outlined above adapts to the specific user 153
 154 by supervised training of subject- (and session)-specific spatial filters, feature 154
 155 normalization, and classifiers. Once trained, these three components form a 155
 156 subject- and session-specific classification system c_s (subsequently called a *clas-* 156
 157 *sification flow*) that maps preprocessed time series x onto the scalar classifier 157
 158 prediction $c_s(x) \in \mathbb{R}$. Unfortunately, training of a classification flow requires a 158
 159 large training dataset that needs to be recorded at the start of each session. In 159
 160 order to reduce the required amount of training data (possibly even to zero), Fa- 160
 161 zli et al. [6] proposed to reuse classification flows trained on N historic sessions 161
 162 from the same and other subjects; such a set $h = (c_{h_1}, \dots, c_{h_N})$ of historical 162
 163 classification flows c_{h_i} is called an *ensemble*. An ensemble can be used to gener- 163
 164 ate a vector of class predictions $h(x) = (c_{h_1}(x), \dots, c_{h_N}(x)) \in \mathbb{R}^N$ for a given 164
 165 time series x . 165

166 Thereupon, a so-called *gating function* g combines the ensemble’s predic- 166
 167 tions $h(x) \in \mathbb{R}^N$ into a joint prediction $g(h(x)) \in \mathbb{R}$ (in the linear case $g(x) =$ 167
 168 $\sum_{i=1}^N w_i c_{h_i}(x)$). A gating function can be defined without requiring session- 168
 169 specific training data by, e.g., training it on historic data (compare Fazli et 169
 170 al. [6]) or, alternatively, without any training by predicting according to the 170
 171 equally-weighted mean of the ensemble’s predictions ($w_i = 1/N$). Furthermore, 171
 172 in situations where a small amount of session-specific training data is available, 172
 173 it is possible to train a gating function such that higher weights w_i are assigned 173
 174 to historic flows c_{h_i} that have high predictive performance for the current ses- 174
 175 sion. In this paper, we focus on the latter approach since it can be combined 175
 176 naturally with the proposed augmentation approaches (see below). We use an 176
 177 SVM with linear kernel for learning the gating function’s parameters w_i since 177
 178 this SVM-based gating function achieved superior performance on hold-out test 178
 179 data of the given scenario compared to other common methods for learning gat- 179
 180 ing functions. The outlined “pure” ensemble approach is depicted as the middle 180
 181 layer in Figure 2. 181

182 *Augmentation approaches* While ensemble approaches have been successful in 182
 183 achieving good performance when only a limited amount (or even no) training 183
 184 data from the current session is available (see, e.g., [6]), it is unlikely that they 184
 185 can achieve competitive results when more session-specific training data becomes 185
 186 available since they can not exploit novel patterns or shifts present in the current 186
 187 session that have not been observed in any of the historic sessions. We propose 187
 188 to use the ensemble approach presented above not instead but in addition to 188
 189 the training of a session-specific flow c_s , i.e., to *augment* the session-specific 189
 190 flow c_s by the predictions of the ensemble h . In this approach, the available 190
 191 training data is used for two purposes: training of a session-specific flow c_s and 191
 192 training of the gating function g which determines the final classification based on 192
 193 the ensemble’s predictions and the session-specific information. We propose and 193

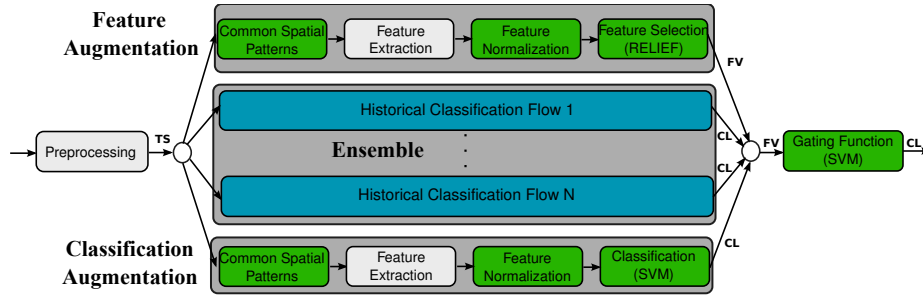


Fig. 2. Different ensemble and augmentation approaches. Feature Augmentation and Classification Augmentation are two alternative approaches for augmenting the ensemble’s predictions by session-specific information. TS denotes a time-series, FV a feature vector, and CL a scalar classifier prediction.

194 compare two alternative approaches: *Classification Augmentation* and *Feature* 194
 195 *Augmentation* (see Figure 2). 195

196 In the classification augmentation approach, the prediction of the session- 196
 197 specific classification flow $c_s(x)$ is treated like any of the ensemble flow’s pre- 197
 198 dictions $c_{h_i}(x)$: An augmented ensemble $\tilde{h} = (c_{h_1}, \dots, c_{h_N}, c_s)$ is generated 198
 199 and the gating function g chooses the joint prediction $g(\tilde{h}(x))$ based on \tilde{h} ’s output 199
 200 ($\tilde{h}(x) \in \mathbb{R}^{N+1}$). Both c_s and g need to be trained based on data acquired in the 200
 201 current session; using the same data for both tasks, however, would result in a 201
 202 too strong reliance of the gating function on c_s since the predictive performance 202
 203 of c_s would be evaluated on its own training data. Thus, the available training 203
 204 data needs to be split into two parts. Empirically, we have found that using 2/3 204
 205 for training of c_s and 1/3 for training of g is a good compromise. 205

206 In contrast, in the feature augmentation approach, the session-specific in- 206
 207 formation added to the ensemble’s predictions is not the classifier’s prediction 207
 208 $c_s(x)$ but the values of the n most informative features $f_1(x), \dots, f_n(x)$, i.e., 208
 209 $\tilde{h}(x) = (c_{h_1}(x), \dots, c_{h_N}(x), f_1(x), \dots, f_n(x)) \in \mathbb{R}^{N+n}$. Thus, $\tilde{h}(x)$ consists of 209
 210 two very different kinds of values: classifier predictions and CSP-pseudo-channel 210
 211 values (the selected features). However, this does not impose a problem and has 211
 212 the advantage that the available training data can be used more efficiently than 212
 213 in classification augmentation (note that while in principle feature selection and 213
 214 training of the gating function should be done on disjoint training sets, we have 214
 215 found empirically that it is favorable to train both on the same data). The choice 215
 216 of n is one additional parameter of this approach. The determination of the most 216
 217 informative features is made using the RELIEF feature selection algorithm [7]. 217

218 4 Evaluation 218

219 *Experimental Setup* One historic classification flow has been trained for each 219
 220 historic session, resulting in 12 historic classification flows. Each of the 12 sessions 220

has been used once as evaluation session with the remaining 11 sessions being considered accordingly as historic sessions. Two different settings have been compared: In the “LeaveOneSessionOut” setting, the classification flows belonging to all but the current evaluation session have been used in the ensemble (resulting in ensembles of $N = 11$ flows), while in the “LeaveOneSubjectOut” setting, all classification flows that have not been generated from usage sessions of the current subject are used in the ensemble (resulting in ensembles of $N = 10$ flows). For each evaluation session, the data recorded in the first run has been used as training data and each of the remaining four runs has been used once as test dataset (intra-session setup), resulting in $4 * 12 = 48$ performance samples per method. Training datasets of six different sizes $t \in \{42, 84, 168, 252, 420, 840\}$ have been randomly sampled from the 840 labeled instances of the first run, where $t = 840$ corresponds to a calibration time of approximately 16 minutes. We refer to “experimental_design.pdf” in the supplementary material for more details.

Parameters of the SVM gating function have been selected using 5-fold internal cross-validation on the training data (complexity $C \in \{0.001, 0.01, 0.1, 1.0\}$ and target class weight $w_t \in \{1, 2, 5, 10\}$ for standard class weight 1). The parameter n of the feature-augmentation approach has been linearly increased from $n = 2$ for $t = 42$ to $n = 50$ for $t = 840$ to account for a stronger influence of the session-specific information when more training data becomes available. The scalar output of the gating function g is mapped onto the binary classes by choosing a threshold that maximizes the performance on the training data. For comparison, the results of the “zero-training” gating function that predicts according to the equally-weighted ensemble mean are given for the pure ensemble for $t = 0$. The performance of the session-specific flow c_s is given as “baseline”. No value for classification augmentation is given for $t = 42$ since not enough target class training examples were available for the two-stage training procedure.

Because of the large class-skew of the classification task, standard measures such as accuracy are not well suited as performance metric. Instead, performance is measured according to the *mutual information* metric $I(T; Y) = H(T) - H(T|Y)$ with $H(T) = -\sum_{i=1}^n p(x_i) \log_2 p(x_i)$ being the Shannon entropy of the class label T and $H(T|Y)$ the conditional entropy of the class label T given the classifier’s prediction Y . The values of the metric correspond to the bits of information about the true class label conveyed by the classifier. The main advantage of this metric is that any kind of random classifier has mutual information 0. Note that the class label’s entropy (and thus $I(T; Y)$) is upper bounded by $H(T) \approx 0.533$ for the given class ratio of 6 : 1. The optimally achieved performance (mutual information of 0.22) corresponds roughly to 94% correct classifications.

Results and Discussion We compared the four different approaches (factor e) for different training set sizes (factor t) by repeated measures ANOVA with t and e as within-subjects factors. This statistic model was separately performed for each setting $s \in \{ \text{“LeaveOneSessionOut”}, \text{“LeaveOneSubjectOut”} \}$ because of the different ensemble sizes N for the two settings. Whenever the results of

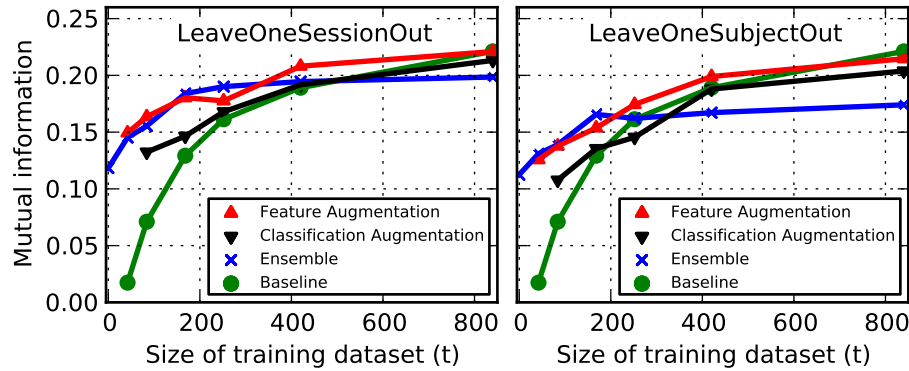


Fig. 3. Effect of training set size. Comparison of baseline, ensemble, and augmentation approaches for maximal N (LeaveOneSessionOut: $N = 11$, LeaveOneSubjectOut: $N = 10$) and for different training set sizes t .

266 the two different settings were compared, the additional factor s was added to 266
 267 the statistic model. Furthermore, in order to avoid that the different values of 267
 268 N for the two settings s affect these comparisons, one randomly selected session 268
 269 of another subject was removed from the “LeaveOneSessionOut” setting such 269
 270 that $N = 10$ in both cases. If needed, the Greenhouse-Geisser correction was 270
 271 applied. For pairwise comparisons, Bonferroni correction was applied. All tests 271
 272 have been performed for a significance level of $p < 0.05$ (see “statistics.pdf” in 272
 273 supplementary material for more detailed results). 273

274 Figure 3 summarizes the results of the study. In the “LeaveOneSessionOut” 274
 275 setting, the ensemble approach is significantly better than the baseline for $t \leq$ 275
 276 252 and worse for $t = 840$. This supports the hypothesis that historic predictors 276
 277 provide good performance when only a small amount of training data is available 277
 278 but are outperformed by session-specific predictors when larger amounts of training 278
 279 data have been acquired. Among the two augmentation approaches, feature 279
 280 augmentation is clearly better with statistical significance for $t \in \{42, 84, 168, 420\}$. 280
 281 This may be attributed to the inefficient usage of training data in the classification 281
 282 augmentation approach where it is necessary to split the training data 282
 283 into two disjoint parts (see Section 3). Furthermore, feature augmentation can 283
 284 be considered to be superior to both the ensemble and the baseline approach 284
 285 since performance is never significantly worse than any of the two, but signif- 285
 286 icantly better than the ensemble for $t \geq 420$ and better than the baseline for 286
 287 $t \in \{42, 84, 168, 420\}$. This indicates that feature augmentation provides an eff- 287
 288 icient way of combining historic and session-specific information by adaptively 288
 289 learning which source of information should be trusted more. 289

290 Results in the “LeaveOneSubjectOut” setting are qualitatively similar, with 290
 291 the notable difference that the ensemble’s performance is significantly worse than 291
 292 in the “LeaveOneSessionOut” setting for all t . This shows that a historic session 292

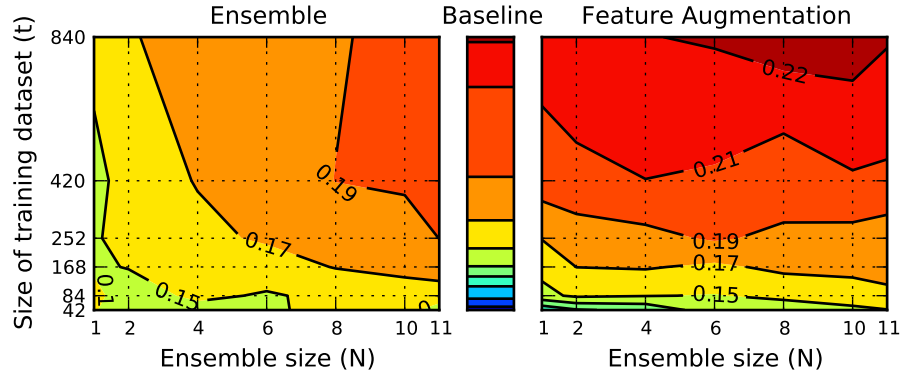


Fig. 4. Effect of the ensemble size. Mutual influence of ensemble size N and the training set size t onto performance (mutual information) in the LeaveOneSessionOut setting. For comparison, the baseline performance is shown for the same values of t .

293 of the same user helps to increase the performance of the ensemble approach. 293
 294 As a result, in the “LeaveOneSubjectOut” setting, the ensemble is significantly 294
 295 better than the baseline only for $t \leq 168$ but worse for $t = 840$. Performance 295
 296 of the feature augmentation approach deteriorates significantly as well in the 296
 297 “LeaveOneSubjectOut” setting for all $t \neq 252$; however, this deterioration is less 297
 298 strong since the session-specific flow compensates partly for the missing historic 298
 299 session of the same user. Accordingly, the feature augmentation approach is still 299
 300 never significantly worse than the baseline but significantly better for $t \leq 168$. 300

301 Figure 4 shows how the size N ($N \in \{1, 2, 4, 6, 8, 10, 11\}$) of the historic 301
 302 ensemble and the size of the training dataset t mutually affect the performance of 302
 303 the pure ensemble and the feature augmentation approach (in the “LeaveOneSes- 303
 304 sionOut” setting). These results have been separately analyzed for each setting 304
 305 by repeated measures ANOVA with the within-subjects factors N , t , and e . The 305
 306 performance of the pure ensemble approach depends strongly on the ensemble’s 306
 307 size: Even for large t , no performance above 0.17 is achieved for $N \leq 2$ and 307
 308 no performance above 0.19 for $N \leq 6$. This dependence on N is even stronger 308
 309 in the “LeaveOneSubjectOut” setting (see “ensemble_size_LOSubjectO.pdf” in 309
 310 supplementary material). On the other hand, the feature augmentation approach 310
 311 depends less strongly on N , outperforming the baseline for small t significantly 311
 312 even when N is very small ($t < 84$ for $N = 1$; $t < 168$ for $N \in \{2, 4\}$) but never 312
 313 being significantly worse than the baseline. 313

314 5 Conclusion 314

315 We have presented two alternative approaches for combining predictions made by 315
 316 an ensemble trained on historic sessions with a flow that has been trained on data 316
 317 acquired in the current usage session. This hybrid approach allows to achieve a 317

318 better performance than the session-specific predictor when only small amounts 318
 319 of training data are available and a better performance than the historic ensemble 319
 320 when more training data becomes available. The proposed approach performs 320
 321 well for subjects for which historic sessions exist but also for novel subjects 321
 322 for which no historic sessions have been conducted. Furthermore, in contrast 322
 323 to related approaches like [6] and [9], the proposed method also achieves good 323
 324 performance when only a small number of historic sessions is available, where it 324
 325 still outperforms the session-specific predictor for small training datasets. Future 325
 326 work is to conduct online studies in which the acquisition of training data is 326
 327 performed concurrently to the usage session. 327

328 References 328

- 329 1. Supplementary Material I. (supplementary_material.I.pdf) 329
- 330 2. Alamgir, M., Grosse-Wentrup, M., Altun, Y.: Multi-task learning for Brain- 330
 331 Computer Interfaces. In: Proceedings of the 13th International Conference on Ar- 331
 332 tificial Intelligence and Statistics. vol. 9 of JMLR: W&CP 9 (2010) 332
- 333 3. Birbaumer, N.: Breaking the silence: Brain-Computer Interfaces (BCI) for commu- 333
 334 nication and motor control. *Psychophysiology* 43(6), 517–532 (Nov 2006) 334
- 335 4. Blankertz, B., Dornhege, G., Lemm, S., Krauledat, M., Curio, G., Müller, K.R.: 335
 336 The Berlin Brain-Computer Interface: Machine learning based detection of user 336
 337 specific brain states. *Journal of Universal Computer Science* 12(6), 581–607 (2006) 337
- 338 5. Blankertz, B., Tomioka, R., Lemm, S., Kawanabe, M., Müller, K.R.: Optimizing 338
 339 spatial filters for robust EEG Single-Trial analysis. *Signal Processing Magazine,* 339
 340 *IEEE* 25(1), 41–56 (2008) 340
- 341 6. Fazli, S., Popescu, F., Danóczy, M., Blankertz, B., Müller, K., Grozea, C.: Subject- 341
 342 independent mental state classification in single trials. *Neural Networks* 22(9), 342
 343 1305–1312 (Nov 2009) 343
- 344 7. Kira, K., Rendell, L.A.: The feature selection problem: Traditional methods and a 344
 345 new algorithm. In: *AAAI*. pp. 129–134 (1992) 345
- 346 8. Koles, Z.J.: The quantitative extraction and topographic mapping of the abnormal 346
 347 components in the clinical EEG. *Electroencephalography and Clinical Neurophysi-* 347
 348 *ology* 79, 440–447 (1991) 348
- 349 9. Krauledat, M., Tangermann, M., Blankertz, B., Müller, K.: Towards zero training 349
 350 for Brain-Computer interfacing. *PLoS ONE* 3(8), e2967 (2008) 350
- 351 10. Li, Y., Guan, C., Li, H., Chin, Z.: A self-training semi-supervised SVM algorithm 351
 352 and its application in an EEG-based brain computer interface speller system. *Pat-* 352
 353 *tern Recognition Letters* 29(9), 1285–1294 (Jul 2008) 353
- 354 11. Lotte, F., Guan, C.: Learning from other subjects helps reducing Brain-Computer 354
 355 interface calibration time. In: *International Conference on Audio Speech and Signal* 355
 356 *Processing (ICASSP)* (2010) 356
- 357 12. Squires, N.K., Squires, K.C., Hillyard, S.A.: Two varieties of long-latency posi- 357
 358 tive waves evoked by unpredictable auditory stimuli. *Electroencephalography and* 358
 359 *Clinical Neurophysiology* 38(4), 387–401 (April 1975) 359
- 360 13. Wolpaw, J.R., Birbaumer, N., McFarland, D.J., Pfurtscheller, G., Vaughan, T.M.: 360
 361 Brain-computer interfaces for communication and control. *Clinical Neurophysiol-* 361
 362 *ogy* 113(6), 767–791 (Jun 2002) 362
- 363 14. Zito, T., Wilbert, N., Wiskott, L., Berkes, P.: Modular toolkit for data processing 363
 364 (MDP): a python data processing framework. *Front. Neuroinform.* 2, 8 (2008) 364