

A 3D-ResNet Combined with BRNN: Application in the Auxiliary Diagnosis of ADHD

Abstract

Attention Deficit/Hyperactivity Disorder (ADHD) is a common mental disorder that exhibits a high incidence rate in children and adolescents, and it is also observed in adults. Currently, there is a lack of objective diagnostic methods for ADHD. Therefore, a three-dimensional residual network (3D-ResNet) deep learning method based on feature extraction from rs-fMRI images for assisting in the diagnosis of ADHD based on resting-state functional magnetic resonance imaging (rs-fMRI) and deep learning models was proposed in this paper. Taking into consideration the temporal characteristics of rs-fMRI, we constructed a 3D-ResNet model based on four-dimensional image. The model utilized TimeDistributed to encapsulate residual blocks which allowed the model to extract spatial features from rs-fMRI while preserving its temporal sequence information. We constructed four different hierarchical structures of 3D-ResNet which are subsequently combined with two different bidirectional recurrent neural networks (BRNNs) to extract sequence features. And BRNNs includes bidirectional long short-term memory (Bi-LSTM) and bidirectional gated recurrent unit (Bi-GRU). The proposed method utilized the ADHD-200 Consortium's public dataset for training and was validated by 5-fold cross-validation. The experimental results indicated that the proposed method in this study demonstrated superior performance on the dataset compared to traditional methods (Accuracy: 76.56%, Sensitivity: 80.16%, Specificity: 90.22%). Therefore, adopting this method can further enhance the accuracy of assisting in the diagnosis of ADHD.

1 Introduction

Attention Deficit/Hyperactivity Disorder (ADHD) is a highly prevalent neurodevelopmental disorder¹, the onset of the condition typically occurs before the age of 12², and it is characterized by persistent hyperactivity, excessive impulsivity, or inability to concentrate³. ADHD can be classified into three subtypes, namely Attention Deficit/Hyperactivity Disorder-Combined Type (ADHD-C), Hyperactive-Impulsive Type (ADHD-HI), and Inattentive Type (ADHD-I). To mitigate the challenges associated with image fusion, this study constructed a deep learning model based on rs-fMRI image data. The paper introduced a network architecture named 3D-ResNet, which is employed to extract spatial features from rs-fMRI. Subsequently, it combines Bidirectional Recurrent Neural Networks (BRNNs) to extract temporal features. Traditional Recurrent Neural Networks (RNNs) face challenges such as gradient vanishing and exploding when dealing with sequential data⁷, these problems make it difficult to capture long-range dependencies. BRNNs capture dependency relationships in sequential data by combining information from both the forward and backward directions. Unlike RNNs, which only consider past information, BRNNs simultaneously take into account both past and future information³. This ability helps the model to comprehensively understand the context within the sequence and consequently improve its classification accuracy. The paper combines two different types of BRNN—Bidirectional Long Short-Term Memory (Bi-LSTM) and Bidirectional Gated Recurrent Unit (Bi-GRU)—to find the optimal network composition, Bi-LSTM is an extended form of Long Short-Term Memory (LSTM)⁴, and Bi-GRU is a bidirectional recurrent neural network based on Gated Recurrent Unit (GRU)⁵. Unlike LSTM

85 and GRU, Bi-LSTM and Bi-GRU can
86 simultaneously consider past and future
87 information at each time step, and that enables
88 better capture of long-term dependencies in time
89 series. Compared to traditional methods, the main
90 contributions of our method are as follows:

91 Using scaled rs-fMRI image data as input avoids
92 the cumbersome preprocessing steps associated
93 with multimodal image fusion, and it reduces the
94 need for extensive feature engineering.

95 A deep learning model based on four-dimensional
96 image data was innovatively constructed. This
97 model focuses on extracting spatial features while
98 preserving fMRI time series, and that enhances
99 the correlation of features in both spatial and
100 temporal dimensions.

101 Four different structures of 3D-ResNet were
102 designed, and they were combined with two types
103 of BRNNs. Through ablation experiments, the
104 optimal model combination was identified.

105 2 Methods

106 The model proposed in this paper takes rs-fMRI
107 image with a temporal dimension as input. It
108 extracts spatial and temporal features through
109 different modules. The model consists of three
110 modules: spatial feature extraction network,
111 temporal feature extraction network, and classifier.
112 Due to rs-fMRI being a four-dimensional image
113 with a temporal dimension, to preserve its
114 inherent temporal characteristics when extracting
115 the spatial features of rs-fMRI images, this paper
116 innovatively builds a three-dimensional residual
117 network (3D-ResNet) based on three-dimensional
118 convolutional neural networks (3D-CNNs). This
119 is a sequential combination network that
120 simultaneously considers spatial and temporal
121 features. The 3D-ResNet network is utilized to
122 extract spatial features while preserving the
123 temporal features of the images. Subsequently, the
124 extracted feature sequence is used as input for the
125 second module, where a temporal feature
126 extraction network processes the sequence. Finally,
127 the obtained features are fed into a classifier. The
128 following three sections will provide detailed
129 compositions of each module.

130 3 Experiment

131 3.1 Dataset and preprocessing

132 This paper trains on rs-fMRI data that has been
133 preprocessed using the Athena pipeline. The

4 samples are sourced from the ADHD-200 Global
5 Competition dataset. The Athena pipeline
6 provides information on 973 preprocessed
7 subjects, including rs-fMRI scans, T1-weighted
8 structural scans, and preprocessed script files. The
9 preprocessing steps primarily include operations
0 such as slice timing correction, head motion
1 correction, smoothing, and filtering. To mitigate
2 the impact of age differences and the imbalanced
3 distribution of positive and negative samples on
4 model training, after stage exclusion, the
5 remaining rs-fMRI data from 430 subjects is used
6 as input. ADHD subtypes are ignored, and all
7 subclasses are labeled as 1. The average age of
8 participants is 12.62, with an equal proportion of
9 ADHD to Typically Developing Control (TDC)
0 subjects at a ratio of 1:1. For a detailed
1 composition of the dataset, refer to Table 1.

2 Table 1: The detailed composition of multi-site
3 samples.

	Pittsburgh	Peking	Total
ADHD	0	78	215
TDC	49	88	215
Total	49	166	430

4 After preprocessing, the spatial dimensions of the
5 rs-fMRI data are 49x58x47. However, due to
6 variations across different sites, the length of the
7 time series is not uniform. For example, the fMRI
8 data from the NYU site has a scan time of 172,
9 while the image time series length from the
0 NeuroIMAGE site is 257. To mitigate the
1 potential impact of differences in scanners and
2 parameters across different sites on experimental
3 results, the original images are cropped and
4 resized to 20x34x34x34x1 before training. This
5 size is then used as the input dimension for the
6 model, where 20 represents the length of the rs-
7 fMRI time series, and 34 represents the spatial
8 dimensions in terms of length (h), width (w), and
9 depth (d). This approach not only standardizes the
0 image sizes across different sites but also
1 substantially reduces the model's parameter count,
2 simultaneously it can avoid overfitting risks and
3 reduce memory overhead⁶⁻⁸.

5 3.2 Model training

6 In this study, all models are trained using binary
7 cross-entropy loss function with the Adam
8 optimizer. The learning rate is set to 1×10^{-4} .
9 Given that fMRI images are four-dimensional

180 with relatively large dimensions, the images are
 181 scaled proportionally before training. The
 182 standardized size after scaling is 20x34x34x34x1.
 183 Subsequently, the samples are fed into the model
 184 for training. The batch size is set to 32, and the
 185 number of epochs is set to 100. The model's
 186 training progress is evaluated using the accuracy
 187 metric. After multiple parameter adjustments, the
 188 training proceeds successfully.

189 To avoid overfitting during the training process,
 190 this study employs the early stop technique with a
 191 tolerance set to 5. In other words, if the loss does
 192 not show a decreasing trend for five consecutive
 193 iterations, the training is terminated, and the
 194 model parameters from five iterations ago are
 195 saved. Since this paper proposes multiple
 196 networks with different structures, to
 197 comprehensively evaluate the model performance,
 198 multiple metrics are incorporated in the model
 199 evaluation. Additionally, a five-fold cross-
 200 validation is used to obtain more accurate
 201 classification performance.

202 4 Results

203 This section primarily presents the experimental
 204 results of training four different residual networks
 205 combined with two types of BRNNs, and
 206 compares them with existing models. In certain
 207 research reports focusing on classification tasks,
 208 the majority often rely solely on the accuracy
 209 metric to assess their methods. However, this
 210 alone is inadequate to substantiate the feasibility
 211 of their approaches, that's because high
 212 classification accuracy may be a result of
 213 imbalanced sample distribution, and it will lead
 214 the model to exhibit bias towards predicting a
 215 specific class in extreme cases. For instance, in a
 216 binary classification task where positive samples
 217 constitute only 10% of the entire dataset, if the
 218 model predicts all samples as negative, the
 219 accuracy can reach 90%. However, for the 10%
 220 positive samples, the model's ability to accurately
 221 predict is uncertain. In this case, the high accuracy
 222 is superficial and lacks practical significance.
 223 Therefore, to accurately assess the model
 224 performance, this paper introduces specificity and
 225 sensitivity. Specificity represents the false positive
 226 rate, and a high specificity indicates a low number
 227 of misdiagnosed samples. In simple terms, it
 228 reflects the model's ability to correctly identify

9 TDC. It can be calculated as follows:

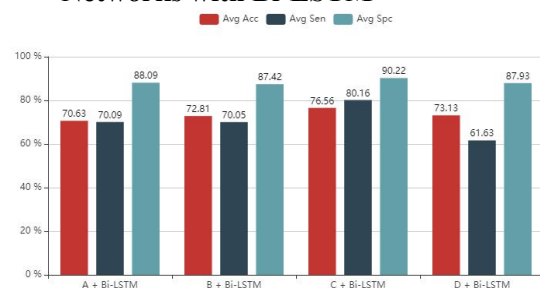
$$specificity = \frac{TN}{TN + FP}$$

1 Sensitivity means the proportion of samples
 2 not missed in the model's prediction results
 3 relative to the total number of samples. It
 4 represents the model's ability to correctly identify
 5 ADHD cases. The calculation for sensitivity is as
 6 follows:

$$sensitivity = \frac{TP}{TP + FN}$$

8 Where True Negative (TN) represents the number
 9 of samples that are negative and predicted as
 0 negative, True Positive (TP) represents the
 1 number of samples that are positive and predicted
 2 as positive, False Negative (FN) represents the
 3 number of samples that are positive but predicted
 4 as negative, and False Positive (FP) represents the
 5 number of samples that are negative but predicted
 6 as positive⁹.

7 4.1 Results of combining Residual 8 Networks with Bi-LSTM

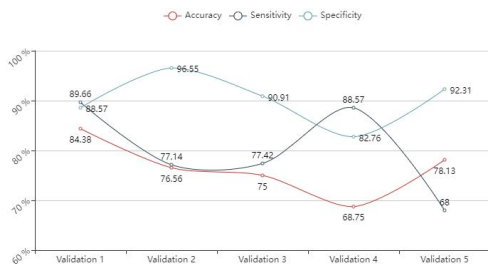


9 Figure 1: The experimental results of the combination
 0 of 3D-ResNet and Bidirectional LSTM model

1 In Figure 1, the experimental results of training
 2 with the combination of four Residual Networks
 3 and Bi-LSTM are presented. The data in the table
 4 represents the average results as the paper utilized
 5 five-fold cross-entropy validation. From the
 6 experimental results, it can be observed that the
 7 performance of these four combined models is
 8 quite close, with the main differences manifesting
 9 in terms of accuracy and sensitivity.

1 In the case of combining with Bi-LSTM, the
 2 accuracy of Residual Networks A and B is inferior
 3 to that of Residual Networks C and D. The
 4 combined model of Residual Network C achieved
 5 the highest sensitivity of 80.16% and the highest
 6 accuracy of 76.56%. During the training process,

267 the models of Residual Networks C and D exhibit
 268 a faster convergence rate compared to those of
 269 Residual Networks A and B. Figure 5 illustrates
 270 the performance of this model using 5-fold cross-
 271 entropy validation. As shown in the figure, except
 272 for validation set 4, the accuracy of other
 273 validation sets is greater than 75%. The
 274 comprehensive performance of sensitivity and
 275 specificity indicates that this combination has
 276 good adaptability and can fit the model's
 277 classification curve well.



278

279 Figure 2: The performance of the combination of
 280 Residual Network C and Bi-LSTM using 5-fold cross-
 281 entropy validation

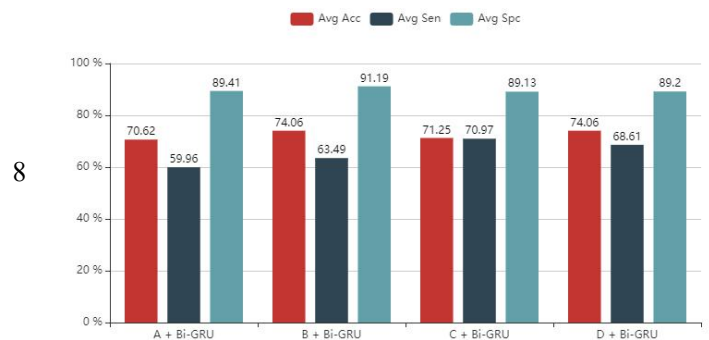
282 4.2 Results of combining Residual 283 Networks with Bi-GRU

284 Following the same method as the previous
 285 section, four Residual Networks were individually
 286 combined with Bi-GRU for training and
 287 validation. Figure 3 illustrates the performance of
 288 the models after 5-fold cross-entropy validation.
 289 The main difference in the current combination
 290 method is reflected in sensitivity, while the four
 291 Residual Networks show similar performance in
 292 accuracy and specificity. From the figure, it can be
 293 observed that the combined model of Residual
 294 Network C has the best overall performance
 295 across these three metrics, and its accuracy,
 296 sensitivity, and specificity are 71.25%, 70.97%,
 297 and 89.13%. When the model's accuracy and
 298 specificity exhibit similar performance, sensitivity
 299 becomes a key indicator representing the
 300 performance differences among the four models.
 301 The performance of the model on the validation
 302 set at this time is shown in Figure 4.

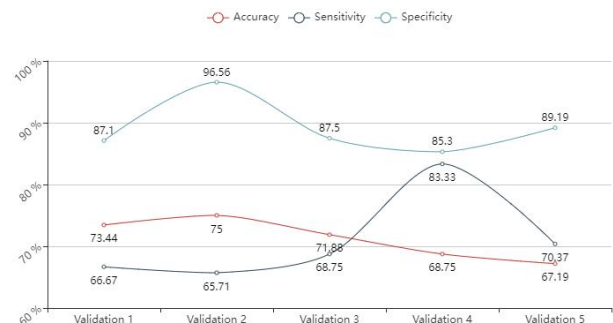
303 The data in the figure indicates that except for
 304 validation set 2, the accuracy of the model
 305 combining Residual Network C with Bi-GRU is
 306 less than 75%. The accuracy on validation sets 4

337 Table 2: Comparison of performance metrics across models.

17 and 5 is even less than 70%. Compared to the
 18 model combined with Bi-LSTM, the combination
 19 of Residual Network C and Bi-GRU is not stable
 20 enough, especially in terms of sensitivity. It can't
 21 fit the classification curve well and only performs
 22 well in certain specific intervals, and the results
 23 lack generality. Therefore, by comparing the data,
 24 it can be concluded that the model combining
 25 Residual Network C with Bi-LSTM exhibits the
 26 best performance in the task of ADHD
 27 classification recognition.



9 Figure 3: The experimental results of the combination
 10 of 3D-ResNet and Bidirectional GRU model



11 Figure 4: The performance of the combination of
 12 Residual Network C and Bi-GRU using 5-fold cross-
 13 entropy validation

15 4.3 Compare with existing models

16 Table 2 compares various existing methods with
 17 the approach proposed in this paper. Firstly, all of
 18 them are based on the ADHD-200 dataset, it can
 19 be observed that the proposed method in this
 20 paper shows a significant improvement compared
 21 to existing methods, both in terms of accuracy,
 22 sensitivity, and specificity. However, due to
 23 variations in data partitioning across different
 24 studies, and the fact that some studies do not
 25 report specificity and sensitivity, so making direct
 26 comparisons is challenging.

	Methods	Validation	Accuracy	Sensitivity	Specificity
Dai et al. ⁷	MKL	10-fold cross	67.79%	38.29%	84.08%
Zou et al. ¹¹	3D-CNN	hold-out set ¹	69.15%	-	-
Mao et al. ¹³	4D-CNN	hold-out set ¹	71.30%	73.20%	69.70%
Zhang et al. ¹⁴	SC-CNN-Attention	loocv ²	68.60%	-	-
Niu et al. ¹⁶	3D-CNN+GRU	5-fold cross	71.65%	68.00%	73.80%
Proposed	3D-ResNet+Bi-LSTM	5-fold cross	76.56%	80.16%	90.22%

339 5 Discussion

340 In existing research, the majority of studies adopt
341 dimensionality reduction to extract low-order
342 features from fMRI images. This machine
343 learning approach, by employing straightforward
344 feature selection to discard irrelevant features¹³,
345 and it often overlooks the temporal and spatial
346 correlations in fMRI data. Consequently, the
347 features extracted lack representational capacity.
348 Recent studies tend to lean towards employing
349 deep learning methods for feature extraction, for
350 example, Niu et al.¹⁰ utilized 3D-CNN to process
351 the three-dimensional spatial information of rs-
352 fMRI. Through one-dimensional filters of
353 different scales, significant features were
354 extracted. Subsequently, an independent GRU
355 was employed to handle the one-dimensional
356 temporal information. Finally, feature fusion was
357 performed. Hong et al.¹¹ employed
358 TimeDistributed to encapsulate 3D-CNN for
359 processing the spatiotemporal information of rs-
360 fMRI. Through this network, a feature sequence
361 was obtained. Then, it was combined with LSTM
362 to extract sequence features, this method achieved
363 a classification accuracy of 68% for ADHD. The
364 method of using 3D-CNN and GRU to separately
365 process the spatiotemporal information of rs-
366 fMRI takes into account the correlation between
367 spatial and temporal dimensions. However, in
368 reality, fMRI images exhibit dynamic spatial
369 characteristics with temporal features. The
370 approach of using two independent networks to
371 extract features separately may contradict the
372 realistic nature of the data¹². This results in lower
373 spatiotemporal correlation of features. The
374 concatenated network of 3D-CNN combined with
375 LSTM which is the inspiration behind this paper
376 avoided this issue. Due to the use of
377 TimeDistributed encapsulation for 3D-CNN with
378 four-dimensional input data, the model generated

379 a large number of parameters. Setting a high
380 number of convolutional layers may lead to
381 overfitting during the training process. On the
382 other hand, a shallower 3D-CNN may not
383 effectively extract meaningful information, and it
384 will lead to a lack of expressive power in the
385 model and consequently lower accuracy.

386 In this study, a 3D-ResNet encapsulated by
387 TimeDistributed was constructed to process the
388 spatial information of rs-fMRI. It is combined
389 with BRNNs to form a concatenated network, and
390 that effectively addressed the issues of
391 insufficient spatiotemporal correlation of features
392 and overfitting. By analyzing the experimental
393 results in the fourth section, it can be concluded
394 that the overall performance of the model
395 combining 3D-ResNet with Bi-LSTM is better
396 than that with Bi-GRU. Among them, the model
397 combining Residual Network C with Bi-LSTM
398 exhibits the best performance. However, it is
399 essential to note that this model requires a certain
400 amount of memory as a basis. The experimental
401 model was trained on an NVIDIA A100 GPU
402 with 40GB of memory. Despite the increased
403 memory overhead, the noticeable improvement in
404 accuracy justifies the additional memory usage.

405 6 Conclusion

406 A three-dimensional residual network named 3D-
407 ResNet which was combined with BRNNs was
408 introduced in this study. Compared to techniques
409 that involve fusing fMRI with MRI, the method
410 proposed in this paper eliminated the need for
411 complex image preprocessing; And compared to
412 methods that extract low-level features from
413 fMRI, this model retained spatial correlations
414 while extracting features. This paper constructed
415 four different structures of residual networks, and
416 through ablation experiment, it demonstrated that
417 the model combining Residual Network C with
418 Bi-LSTM has the best performance. Under the 5-

419 fold cross-entropy validation method, the average
420 accuracy, sensitivity, and specificity are 76.56%,
421 80.16%, and 90.22%. Compared to existing
422 methods, there is a significant improvement in
423 accuracy when performing classification tasks on
424 the multi-site ADHD-200 dataset. This result
425 indicated that combining 3D-ResNet with
426 BRNNs for assisting in the diagnosis of ADHD is
427 feasible. What is even more promising is that this
428 technology can be applied to the classification
429 and diagnosis of other neurological disorders. It
430 holds considerable prospects in studies based on
431 rs-fMRI.

432 References

433 Zhou, Dingfu et al. "Deep Learning Enabled
434 Diagnosis of Children's ADHD Based on the Big
435 Data of Video Screen Long-Range EEG." *Journal*
436 *of healthcare engineering* vol. 2022 5222136. 4
437 Apr. 2022.

438 Cherkasova, Mariya V et al. "Review: Adult
439 Outcome as Seen Through Controlled Prospective
440 Follow-up Studies of Children With Attention-
441 Deficit/Hyperactivity Disorder Followed Into
442 Adulthood." *Journal of the American Academy of*
443 *Child and Adolescent Psychiatry* vol. 61,3 (2022):
444 378-391.

445 Bellec, Pierre et al. "The Neuro Bureau ADHD-200
446 Preprocessed repository." *NeuroImage* vol. 144,Pt
447 B (2017): 275-286.

448 ADHD-200-Results-Webpage. "Adhd-200 global
449 competition results" 2011. Available online at:
450 [http://fcon_1000.projects.nitrc.org/indi/adhd200/re](http://fcon_1000.projects.nitrc.org/indi/adhd200/results.html)
451 [sults.html](http://fcon_1000.projects.nitrc.org/indi/adhd200/results.html).

452 Dai, Dai et al. "Classification of ADHD children
453 through multimodal magnetic resonance imaging."
454 *Frontiers in Systems Neuroscience* 6 (2012): n.
455 pag.

456 Guo, Xiaojiao and Lianghua He. "ADHD
457 Discrimination Based on Social Network." 2014
458 *International Conference on Cloud Computing and*
459 *Big Data* (2014): 55-61.

460 Kuang, Deping et al. "Discrimination of ADHD
461 Based on fMRI Data with Deep Belief Network."
462 *International Conference on Intelligent Computing*
463 (2014).

464 Hao, A.J.; He, B.L.; Yin, C.H.: 'Discrimination of
465 ADHD children based on Deep Bayesian Network',
466 *IET Conference Proceedings*, 2015, p. 6.-6 .

467 Zou, Liang et al. "3D CNN Based Automatic
468 Diagnosis of Attention Deficit Hyperactivity

59 Disorder Using Functional and Structural MRI."
70 *IEEE Access* 5 (2017): 23626-23636.

71 Vu, Hanh et al. "3D convolutional neural network for
72 feature extraction and classification of fMRI
73 volumes." 2018 *International Workshop on Pattern*
74 *Recognition in Neuroimaging (PRNI)* (2018): 1-4.

75 Mao, Zhenyu et al. "Spatio-temporal deep learning
76 method for ADHD fMRI classification." *Inf. Sci.*
77 499 (2019): 1-11.

78 Zhang, Tao et al. "Separated Channel Attention
79 Convolutional Neural Network (SC-CNN-
80 Attention) to Identify ADHD in Multi-Site Rs-
81 fMRI Dataset." *Entropy (Basel, Switzerland)* vol.
82 22,8 893. 14 Aug. 2020.

83 Tran, Du, et al. "Learning spatiotemporal features
84 with 3d convolutional networks." *Proceedings of*
85 *the IEEE international conference on computer*
86 *vision*. 2015.

87