

AWS-EP: A Multi-Task Prediction Approach for MBTI/Big5 Personality Tests

Anonymous ACL submission

Abstract

001 Personality and preferences are essential vari- 042
002 ables in computational sociology and social sci- 043
003 ence. They describe differences between peo- 044
004 ple at both individual and group levels. In re- 045
005 cent years, automated approaches to detect per- 046
006 sonality traits have received much attention due 047
007 to the massive availability of individuals' digi- 048
008 tal footprints. Furthermore, researchers have 049
009 demonstrated a strong link between personality 050
010 traits and various downstream tasks such as per- 051
011 sonalized filtering, profile categorization, and 052
012 profile embedding. Therefore, the detection 053
013 of individuals' personality traits has become a 054
014 critical process for improving the performance 055
015 of different tasks. In this paper, we build on the 056
016 importance of the individual personality and 057
017 propose a novel multitask modeling approach 058
018 that understands and models the user person- 059
019 ality based on its textual posts and comments 060
020 within a multimedia framework. Experiments 061
021 and results demonstrate that our model outper- 062
022 forms state-of-the-art performances across mul- 063
023 tiple famous personality datasets. 064

024 1 Introduction

025 Personality traits highlight the difference among 065
026 the various individuals' characteristic patterns such 066
027 as feeling, thinking, and behaving. Understanding 067
028 people's core personality traits and knowing what 068
029 people are good at can be very important in a wide 069
030 variety of situations. It could ameliorate its social 070
031 relationships, personal development, thinking pat- 071
032 terns, and daily interaction capabilities. People 072
033 are now very familiar with personality test systems 073
034 such as the MBTI (Myers Briggs Type Indicator), 074
035 16 personalities, Big5 (Big five-factor model), and 075
036 other tests. MBTI [Isabel Briggs Myers and Ham- 076](#)
037 [mer \(1987\)](#) and Big5 [Goldberg \(1993\)](#) are the most 077
038 well-known personality test systems. Both are used 078
039 within a large scale of companies and therapy intu- 079
040 itions. 080
041 The MBTI system categorizes a person into 16 081
082

different categories using four main factors (Intro- 042
verted, Intuitive, Thinking, and Perceiving). In 043
this system, a user can only belong to one cate- 044
gory. The Five-Factor (Big5) model measures five 045
key dimensions of people's personalities. It mea- 046
sures its openness' OPN,' its Conscientiousness' 047
CON,' its extraversion' EXT,' its agreeableness' 048
AGR,' and its neuroticism' NEU.' In this personal- 049
ity system, a person belongs to all five categories 050
to a certain degree, unlike the MBTI test, where a 051
person can only be one of 16 categories. Recent 052
research demonstrates that people prefer express- 053
ing their emotions, thoughts, and complaints on 054
social media platforms such as Twitter, Instagram, 055
Facebook, among other ones [Yosephine Susanto 056](#)
[and Cambria. \(2020\)](#). Therefore, in modern times, 057
there is a massive interest in designing automatic 058
learning models that benefit from human digital 059
footprints for different end-goals (example: online 060
posts personality detection). 061

Recent works demonstrate that social media indi- 062
vidual digital footprints are very effective for mea- 063
suring personality traits [Wu Youyou and Stillwell 064](#)
[\(2015\)](#). Despite the serious privacy concerns for 065
individuals [Sandra C Matz and Kosinski \(2020\)](#), 066
this challenging task has gained significant interest 067
from psycholinguistics and natural language pro- 068
cessing researchers due to its extensive downstream 069
applications such as profile categorization and psy- 070
chological treatment. Significant strides in machine 071
learning and deep learning-based personality detec- 072
tion research have taken place in the past few years 073
[Yash Mehta and Eetemadi \(2020\)](#), [Wu Youyou and 074](#)
[Stillwell \(2015\)](#), [Li et al. \(2021\)](#), [Tao Yang \(2021\)](#). 075
Moreover, other psychological research highlights 076
the correlation and the dependency between pair 077
personality test systems [Furnham \(1996\)](#). However, 078
all existing automated approaches focus heavily 079
on personality test systems independently, whether 080
modeling the MBTI or the Big5 system. 081

Motivated by the above discussions, we propose 082

083 the first automated multi-personality test systems
084 modeling approach. We propose a novel multi-
085 task personality prediction model named AWS-EP
086 (All Weight Shared Electra for Personality predic-
087 tion) 3.1. Our proposed model consists of an MLP
088 (Multi-Layer Perceptron) architecture with two pre-
089 diction heads (classification and regression), built
090 on top of a fine-tuned Electra transformer model
091 (see section A.1), to model both MBTI and Big5
092 personality test systems at the same time. We
093 choose to use the Electra model because most re-
094 cent published papers use Bert as their primary
095 model to predict personality traits. No one has in-
096 vestigated the use of the Electra model to predict
097 individuals' personality traits. Therefore, this pa-
098 per aims to explore the benefits of using Electra
099 instead of Bert for the personality trait prediction
100 task by comparing its performance with existing
101 state-of-the-art baselines on different datasets.

102 Moreover, we propose three other baselines, named
103 **OC-EP** (Only Classification Electra for Personal-
104 ity prediction) 3.1, **OR-EP** (Only Regression Elec-
105 tra for Personality prediction)3.1, and **EWS-EP**
106 (Electra Weights Shared for Personality prediction)
107 3.1, to locally evaluate the AWS-EP model and
108 measure its performance compared to local base-
109 lines. Our proposed solution outperforms existing
110 state-of-the-art models in different metrics. To the
111 limit of our knowledge, this is the first automated
112 personality detection approach that models indi-
113 vidual personalities while considering more than
114 one personality test system. Moreover, this is the
115 first work that uses shared weights to predict both
116 the categorical values for the MBTI system and
117 the numerical values for the Big5 system at the
118 same time. Also, it is the first work that tackles the
119 Big5 personality trait prediction as a multi-label re-
120 gression task. Experiments conducted on different
121 benchmark datasets show that our AWS-EP model
122 outperforms state-of-the-art models on different
123 metrics.

124 It is important to highlight that our contribution in
125 this work is not creating a novel model architecture
126 for the NLP (Natural Language Processing) field
127 in general. Our contribution is the implementation
128 of different existing NLP mechanisms (pre-trained
129 models, multi-task learning, and weight sharing)
130 to create a novel architecture for the personality
131 trait prediction problem. Moreover, we aim to explore
132 the Electra model performance on the personality
133 trait detection task compared to the existing state-

of-the-art models. 134

2 Related work 135

136 Detecting personality traits can be based on
137 various types of features, such as demographical
138 data (gender, age, followers, etc.), text data (social
139 media content, self-description, etc.) Different
140 research studies have demonstrated that users'
141 online behavior is significantly related to their
142 personality Samuel D Gosling and Gaddis (2011),
143 David John Hughes et al. (2012). Many have
144 successfully applied different learning approaches
145 for a social media-generated content personality
146 trait detection Fabio Celli and Pianesi (2014).
147 Wu Youyou and Stillwell (2015), demonstrate that
148 the digital footprint-based analysis was better at
149 measuring personality traits than close relatives
150 or acquaintances (friends, family, colleagues,
151 etc.) Mayuri Pundlik Kalghatgi and Sidnal (2015)
152 detected the personality trait using an MLP
153 network employing statistical and manual-crafted
154 features. Despite the effectiveness of the manual-
155 crafted features, these types of features are very
156 time-consuming and computationally expensive.
157 That is why researchers have been exploring new
158 data types for personality trait detection.
159 Carducci et al. (2018) were the first to apply
160 textual data for personality detection. They used
161 an SVM (Support Vector Machine) model to do
162 the personality detection on top of textual features
163 instead of the statistical manual-crafted features.
164 Following this work and with the advancement of
165 deep learning approaches, Tommy Tandra (2017)
166 applied personality detection over the text data
167 using LSTM (Long Short-Term Memory) and
168 CNN (Convolutional Neural Network) approaches.
169 Gjurković et al. (2021) used BERT (Bidirectional
170 Encoder Representations from Transformers)
171 Devlin et al. (2019) to set a benchmark for their
172 huge Pandora dataset Gjurković et al. (2021),
173 which include three different personality tests',
174 OCEAN which refers to the Big-Five model
175 categories, MBTI, and Enneagram tests. The
176 authors of this paper developed six regression
177 models to predict age and Big5 traits and eight
178 classification models (The four MBTI features,
179 gender, region, and Enneagram features).
180 Experiments were done using traditional machine
181 learning approaches such as linear/logistic re-
182 gression and deep learning approaches such as
183 MLP. In each model, the comments were encoded

184 using 1024-dimensional vectors derived using
185 BERT, which produced a new benchmark for
186 both regression and classification tasks for this
187 dataset using macro F1- score and P-r-C (Pearson
188 Correlation Coefficient) metrics.

189 Following this work Yang et al. (2021), used both
190 textual and questionnaire answer information to
191 enhance the contextual representation to benefit
192 the personality prediction task.

193 Tao Yang (2021) combined graphical neural net-
194 works with a BERT transformer embedding model
195 to detect personality traits. Their experiments
196 show that their model outperforms the existing
197 state-of-art model by 3.47 and 2.10 points on
198 the average F1-score. To further enhance the
199 effectiveness of personality traits, prediction
200 models Yang Li et al. proposed a new 'Multitask
201 Learning for Emotion and Personality Detection'
202 Li et al. (2021) model. They combined the
203 Bert transformer model and a 3 CNN layers
204 model, allowing information sharing between the
205 different layers to predict user personality and
206 emotion using two different datasets. They also
207 demonstrated that their work surpasses different
208 state-of-the-art models on different metrics such as
209 accuracy, macro-precision, macro-recall, and the
210 macro-F1 metric. The contribution of their work
211 consists of the use of a classification multitask
212 neural network to classify two different tasks
213 (personality and emotion).
214

215 Inspired by all the previous work and the sig-
216 nificant performance improvements that the multi-
217 task learning approach provides, we investigated
218 the effect of using a multi-task MLP approach
219 on top of a fine-tuned Electra transformer model
220 Kevin Clark and Manning (2020), and compared
221 its performance to the already exiting state-of-the-
222 art baselines. We also investigate sharing weights
223 between the MBTI and the Big5 personality tests.
224 Furthermore, we looked at the similarity between
225 these two personality tests.

226 3 Description of Models

227 Throughout this section, we define the different
228 OC-EP (Only Classification Electra for Personality
229 prediction), OR-EP (Only Regression Electra for
230 Personality prediction), EWS-EP (Electra Weights
231 Shared for Personality prediction), and AWS-EP
232 (All Weights Shared Electra for Personality predic-
233 tion) models architecture.

234 The four architectures are built on top of the Elec-
235 tra transformer model. Therefore to understand the
236 proposed architectures, we need first to explain the
237 working mechanism of this model (see Appendix
238 section A.1). Using the pre-trained Electra masked
239 language modeling head, we aim to produce a more
240 contextual representation for each user textual sen-
241 tence to achieve a better text classification perfor-
242 mance.

243 3.1 Models architecture

244 We created four different baseline models to inves-
245 tigate the weight sharing performance for classifi-
246 cation and regression personality prediction tasks:
247 the OC-EP, OR-EP, EWS-EP, and AWS-EP models.
248 We were curious if the independent prediction mod-
249 els would perform better than the weight-shared
250 multi-task models. Figures 1, 2, 3, and 4 describe
251 the main architecture for each baseline.

- 252 • OC-EP:

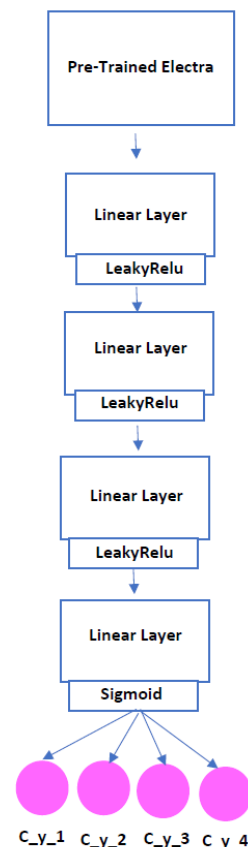


Figure 1: Only Classification Electra for personality prediction Architecture

253 This baseline is designed only for the classi-
254 fication task (predict the MBTI categories),
255 and it is independent of the regression task.

The white boxes represent the different layers in the OC-EP architecture. The output of the sigmoid layer defines the probabilities of each category in the MBTI system, where C_{y1} , C_{y2} , C_{y3} , and C_{y4} , define the introverted, intuitive, thinking, and perceiving MBTI axis. The reason behind using the sigmoid function instead of the softmax function is that the softmax function is generally used when we have a multi-classification task (for example, from the five classes, we need to choose only 1 class). However, in our work, we have a multi-label task (from 5 classes, we can choose 1, 2,3, or even all five classes). This baseline is trained using the BCE (Binary Cross Entropy) loss function applied for each class (equation 1).

$$\begin{aligned}
 LOSS_{class} = & -\frac{1}{NM} \sum_{i=1}^N \sum_{j=1}^M C_{y_{ij}} \cdot \log(\hat{C}_{y_{ij}}) \\
 & + (1 - C_{y_{ij}}) \cdot (1 - \log(\hat{C}_{y_{ij}}))
 \end{aligned} \quad (1)$$

where $N \{1..n\}$ defines the data size, M defines the different classes $\{1..4\}$, $C_{y_{ij}}$ defines the i^{eth} row and j^{eth} class original value $\{0, 1\}$, and $\hat{C}_{y_{ij}}$ defines the i^{eth} row and j^{eth} class predicted value $\{0, 1\}$

This model is trained under the objective of minimizing the $LOSS_{class}$ (equation 2) where X defines the training data and θ_{class} defines the OC-EP model learning parameters.

$$\min_{\theta_{class}} \sum_{x \in X} LOSS_{class}(x, \theta_{class}) \quad (2)$$

• OR-EP:

This baseline is designed only for the regression task (predict the Big5 categories). It is independent of the classification task. The output of the last linear layer defines the numerical values (from 0 to 100) of each factor in the Big5 system, where R_{y1} , R_{y2} , R_{y3} , R_{y4} , and R_{y5} define the agreeableness, openness, conscientiousness, extraversion, neuroticism Big5 factors. This baseline is trained using the MSE (Mean Squared Error) loss function

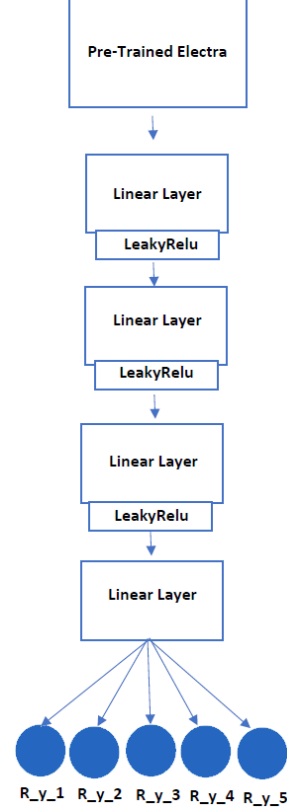


Figure 2: Only Regression Electra for Personality prediction Architecture

for each category (equation 3).

$$LOSS_{reg} = \frac{1}{NM} \sum_{i=1}^N \sum_{j=1}^M (R_{y_{ij}} - \hat{R}_{y_{ij}})^2 \quad (3)$$

where $N \{1..n\}$ defines the data size, M defines the different labels $\{1..5\}$, $R_{y_{ij}}$ defines the i^{eth} row and j^{eth} label original value $\{0..100\}$, and $\hat{R}_{y_{ij}}$ defines the i^{eth} row and j^{eth} predicted value $\{0..100\}$

This model is trained under the objective of minimizing the $LOSS_{reg}$ (equation 4) where X defines the training data and θ_{reg} defines the OR-EP model learning parameters.

$$\min_{\theta_{reg}} \sum_{x \in X} LOSS_{reg}(x, \theta_{reg}) \quad (4)$$

• EWS-EP:

This model is designed to predict both classification (MBTI) and regression (Big5) tasks by sharing only the pre-trained Electra weights $h1$. Regression and classification heads are partially dependent as they share only the $h1$

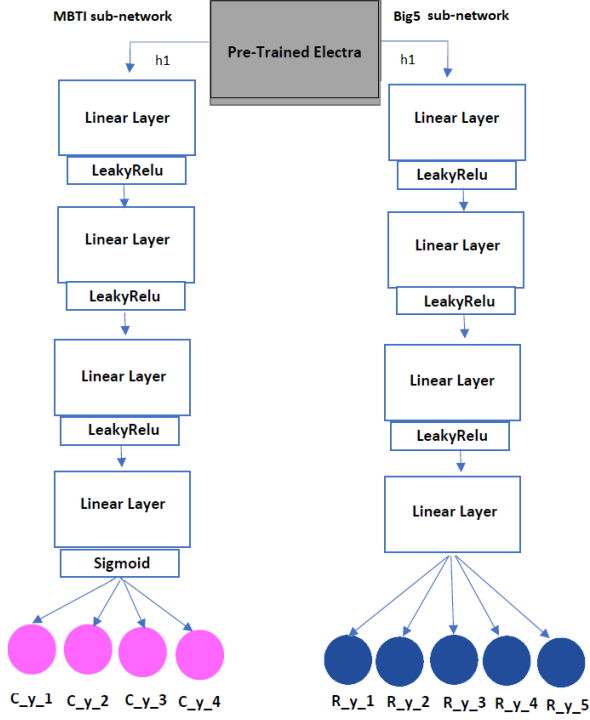


Figure 3: Electra Weights Shared for Personality prediction Architecture

• AWS-EP:

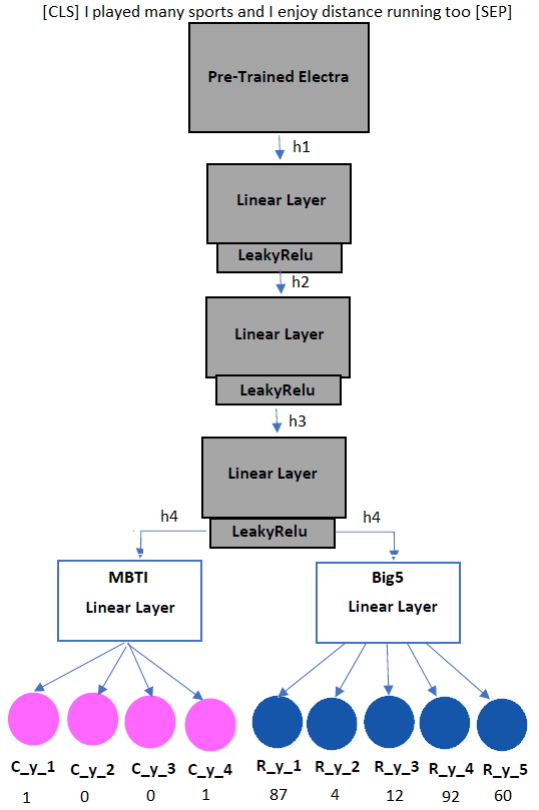


Figure 4: ALL Weights Shared Electra for Personality prediction Architecture

Electra pre-trained model weights. The white boxes represent the independent layers for each sub-architecture. The gray box represents the shared layer 'weights' between the classification MBTI sub-network and the regression Big5 sub-network. The output of the last linear layer defines the numerical values (from 0 to 100) of each Big5 system personality factor (R_{y1} , R_{y2} , R_{y3} , R_{y4}). The output of the last sigmoid layer defines the probabilities of each category in the MBTI system, where C_{y1} , C_{y2} , C_{y3} , and C_{y4} , define the four MBTI personality factors. This model is trained using a combination of the MSE and BCE loss functions.

This model is trained to minimize both the $LOSS_{class}$ and the $LOSS_{reg}$ (equation 5) where X defines the training data, θ_{class} defines the classification sub-model learning parameters, and θ_{reg} defines the regression sub-model learning parameters.

$$\min_{\theta_{class}, \theta_{reg}} \sum_{x \in X} LOSS_{class}(x, \theta_{class}) + LOSS_{reg}(x, \theta_{reg}) \quad (5)$$

Similar to the previous approach, this model is designed to predict both classification (MBTI) and regression (Big5) tasks. However, instead of only sharing the Electra weights $h1$, this approach shares all the network weights (the pre-trained Electra weights $h1$ and the MLP network weights $h2$, $h3$, and $h4$) between the regression and classification heads. The two prediction heads are strongly dependent on each other as they both share the same weights except for the last layer weights. The white boxes represent the independent layers for each sub-AWS-EP architecture. The gray boxes represent the shared layers 'weights' between the classification sub-network and the regression sub-network. Similar to the previous EWS-EP, this model is trained using the same loss function and configuration. The only difference is that all the layers are shared except for the last two heads.

This model is trained under the objective of minimizing both the $LOSS_{class}$ and the $LOSS_{reg}$ (equation 6) where X defines the training data and $\theta_{class,reg}$ defines the model

learning parameters. Unlike the previous model, which had two different model parameters, θ_{reg} and θ_{class} in AWS-EP we have only one model's parameters that combine both regression and classification weights. In AWS-EP, we have only one model's parameters that combine regression and classification weights. Combining the two losses into one loss helps the model focus on both tasks during the training phase.

$$\min_{\theta_{class,reg}} \sum_{x \in X} LOSS_{class}(x, \theta_{class,reg}) + LOSS_{class}(x, \theta_{class,reg}) \quad (6)$$

All proposed baselines in section 3.1 are trained using the same hyper-parameters and share the same work pipeline. First, a sequence of words defined by a sentence start [CLS] and a sentence end [SEP] tokens is given as an input to the Electra base encoder (see the AWS-EP figure 3.1). The encoder will create a contextual vector representation for the sequence of words h_1 . Then the contextual embedding is passed to the MLP network, where we have different linear layers and normalization approaches. The shared MLP linear layers weights (h_2 , h_3 , and h_4) are used to learn the optimal weights that effectively predict both personality factor tests. The last layers (Big5 and MBTI layers) are used as prediction heads. Given the h_4 vector representation, both MBTI and Big5 linear layers try to predict the convenient values for each personality trait factor (example: [1,0,0,1] for the MBTI personality test and [87,4,12,92,60] for the Big5 personality test).

4 Experiments and results

Experiments and results are done using three different datasets. To investigate the performance of our proposed models, we used the Pandora dataset. This dataset combines both Big5 and MBTI features. To evaluate the different model's generalization performances, the MyPersonality Celli et al. (2013) and the MBTI datasets are used. The MyPersonality dataset is used for the Big5 features validation, and the Myers-Briggs Personality Type dataset is used to validate the MBTI features.

4.1 Datasets

- **Pandora dataset Gjurković et al. (2021)** :

Pandora is the largest and the first dataset in the research field that contains more than 17 million Reddit comments written by more than 10k users annotated with both MBTI and Big5 factors with users' demographical features (age, gender, and location). 1.6k users are labeled with the Big5 personality model with more than 3M comments. It also comprises 9k users' annotations with the MBTI personality traits. Due to its massive amount of textual data, throughout this work, Pandora is used as the main dataset to train our baseline models.

It is important to highlight that Pandora is a private dataset and the authors employ different terms of use to protect the users within this dataset Irina Masnikosa and Bakić (2020). Some of the terms consist of not transferring or reproducing any part of the dataset, attempting to identify any user in the dataset, contacting any user in the dataset, displaying users' names and sensitive messages publicly, reporting findings publicly unless it is at an aggregate level. The following two datasets are used as unseen data to evaluate the generalization of our proposed solutions.

- **MyPersonality dataset Celli et al. (2013)** :

This dataset was collected in 2013 by Celli et al. It contains more than 250 different users with 10000 labeled Facebook statuses in total with the Big5 personality traits. It also combines network properties such as network size, density, transitivity, etc.

- **MBTI Personality Type dataset J (2017)** :

This data was collected using the Personality-Cafe forum, as it displays a large selection of people and their MBTI personality type and what they have written. It contains more than 8k rows of data, where each row represents a different person. For each person, we have the last 50 texts they posted.

4.2 Training properties

The Pandora dataset is randomly partitioned into three parts during the training phase: training, validation, and test subsets. 20% of the data were considered a test set, and 80% were considered a training set. Then to create the validation data, we

split the training set into two sub-parts. 20% were considered validation data, and the rest were kept to train the model, which is done using the scikit-learn library. The same data splitting process was done for all the different experiments using a seed value of zero. The sentence words were embedded into a 256-length token vector, using the pre-trained Electra-small model tokenizer from the pytorch hugging face framework. The pre-trained model was fine-tuned on the Pandora training sub-set, and all models (OC-EP, OR-EP, EWS-EP, AWS-EP) were trained for 10 epochs. We also compared the current validation results with the least validation loss for each epoch and stored the model that gave us the least generalization loss. In our experiments, we reported the performance of a single run (10 epochs) for each model. The hyper-parameters we used during our experiments are defined in table 4 in the appendix.

Different experiments were done to investigate the different generalization performance of the proposed baselines (OC-EP, OR-EP, EWS-EP, AWS-EP), and their performance was compared with state-of-the-art models on different datasets. We used the google collaboratory pro version as our computing infrastructure (166.83 Gb hard drive capacity, 25.46 GB memory capacity, and a 1 Tesla P100-PCIE GPU), which allows us to use a 20h window session of these computational resources.

4.3 Training Results

Training the four baselines using the same hyper-parameters led to different performance results on the training and the validation sets. Figure 6 in the appendix section A highlights the different training and validation performance for each baseline.

Training results show that EWS-EP and AWS-EP models (Multi-task models) have the highest trusted results in terms of generalization performance for the MBTI and Big5 traits. We can see that both are trying to reduce at the same time the training and validation loss during each epoch. This highlights the importance and the good performance of the multi-task learning approach compared to the single-task learning approach results. We can see that the validation error is almost constant along the training epochs for the OC-EP. So during the learning phase, this model is trying to decrease the training loss while keeping the validation loss almost the same. Therefore EWS-EP, and AWS-EP are better than the OC-EP on the

classification generalization task.

4.4 Generalization Results

During the experiment phase, we were more curious about having effective results for predicting the Big5 and MBTI personality traits and investigating to which degree these two tests are similar. We also were curious to know the effect of weight sharing on the model predictions. Tables 1,2, and 3 highlight the generalization performance of the OC-EP, OR- EP, EWS-EP, and AWS-EP models. Table 1 highlights the performance of the proposed baselines (OC-EP, OR-EP, EWS-EP, AWS-EP) on the unseen Pandora test subset. The OC-EP model provides good performance for accuracy and F1-score metrics with 0.738 and 0.844, respectively. Results show that the OR-EP model provides an inferior performance in MSE, r2-score. By introducing a low level of weight sharing in the EWS-EP baseline, both classification and regression results improved. Moreover, allowing for more weight sharing between the MBTI and Big5 prediction tasks in the AWS-EP model significantly improved regression and classification results. It is also clear that the regression head is the one that benefits the most from the weight sharing with a more than 100% increase in terms of the Pearson r correlation metric compared to the OR-EP model. We also report a five-fold decrease in MSE compared to the OR-EP model.

The experiments demonstrate that the more we allow the Big5 prediction head to know and share weights with the MBTI model, the better results the head provides. This demonstrates the high correlation between the Big5 and MBTI personality test systems. The results provided in table 1 show that the most effective model from the 4 baselines is the AWS-EP model. For this reason, we aim to investigate the performance of this model further and evaluate its generalization performance. Tables 2 provide more details for the AWS-EP model performance for each trait factor compared to the Pandora baseline Gjurković et al. (2021) and PQ-Net Yang et al. (2021) baselines.

Our AWS-EP model outperformed the state-of-the-art benchmark of the Pandora paper for both MBTI and Big5 prediction tasks. For the MBTI classification task, we achieved a 0.1461 F1-score increase for the Introverted factor compared to the PQ-Net baseline. Also, we achieved a 0.278 increase for the Intuitive factor, a 0.0882 increase

Models	Classification				Regression		
	Accuracy	Precision	Recall	F1-score	MSE	R2_score	P_r_C
OC-EP	0.738	0.738	1.0	0.844	-	-	-
OR-EP	-	-	-	-	2910.39	-2.82	0.32
EWS-EP	0.739	0.739	1.0	0.845	839.03	0.05	0.47
AWS-EP	0.788	0.792	0.94	0.860	564.12	0.35	0.66

Table 1: The baselines performance on different metrics

Classification performance						
	Pandora	PQ-Net	AWS-EP (Ours)			
MBTI factors	f1	f1	accuracy	precision	recall	f1
Introverted	0.654	0.6894	0.7583	0.7629	0.9233	0.8355
Intuitive	0.606	0.6765	0.9131	0.9131	1.0	0.9545
Thinking	0.739	0.7912	0.7889	0.7939	0.9797	0.8771
Perceiving	0.642	0.6957	0.6916	0.7014	0.8625	0.7736
Average	0.6602	0.7132	0.7880	0.7928	0.9414	0.8602
Regression performance						
Metric	Pandora	AWS-EP				
P-r-C metric	0.2629	0.66				

Table 2: The AWS-EP detailed performance compared to the Pandora paper and the PQ-Net state-of-the-art models on the Pandora benchmark dataset

for the Thinking factor, and a 0.0847 increase for the Perceiving factor. Overall, we achieved a 0.147 F1 score average increase for all the MBTI factors compared to the PQ-Net state-of-the-art model. For the Big5 classification task, we achieved a 0.3971 increase in the Pearson correlation metric. Table 2 show that our AWS-EP model outperforms the state-of-the-art models on the Pandora benchmark dataset.

Despite the promising performance of our AWS-EP model on the Pandora dataset, we were curious to measure its generalization performance using different unseen personality datasets. The datasets we use in this experiment are the MBTI personality dataset from Kaggle J (2017), and the MyPersonality dataset Celli et al. (2013). Table 3 demonstrates the generalization performance of the AWS-EP model on different datasets that it has not been trained on.

As shown in table 3, the AWS-EP model performs exceptionally well on different unseen data. Although it was only trained on the Pandora dataset, this model outperforms state-of-the-art MBTI Kaggle datasets baselines. Without any tuning, our model outperformed state-of-the-art models on different datasets.

Moreover, this model also provides good Pearson r

	MBTI Kaggle	MyPersonality	
Metric	F1	accuracy	P_r_C
TrigNet	0.7086	-	-
PQ-Net	0.7132	-	-
BERT	-	0.7210	-
AWS-EP (Ours)	0.8276	0.487	0.6624

Table 3: AWS-EP generalization performance on different unseen datasets

correlation (P-r-c) results for predicting the regression Big5 trait values on the unseen MyPersonality dataset with a 0.66 correlation value. However, for the Big5 classification task, our model provides a very poor performance compared to the state-of-the-art baseline on the same dataset. While The MyPersonality dataset baseline is a multi-label classification model, and it is trained to classify the Big5 traits categories, our Bi5 sub-model is a regression model. Hence, we cannot compare both model results because they operate on two different tasks. However, despite the good results of our AWS-EP model on the regression task, we were curious to know its performance on the classification task as zero-shot learning. To evaluate its Big5 classification performance, we took the predicted regression values and transformed them

597 into classification values (0 or 1) by applying a
598 50% threshold. This transformation did not surpass
599 the MyPersonality state-of-the-art baseline trained
600 on a classification objective. However, as a zero-
601 shot prediction, the AWS-EP results are promising.
602 Also, this highlights the need to add a new Big5
603 classification head to the AWS-EP model.

604 5 Ethical impact of our work

605 Despite the vast benefits of knowing the user’s per-
606 sonality on his/her daily life services, having the
607 individual personality traits without his/her permis-
608 sion or explicitly indicating his/her personality to
609 us can be unacceptable. We believe that attempt-
610 ing to detect the individual personality can be a
611 personality intrusion. Knowing the individual’s
612 personality can help us know his/her preference,
613 his/her behavior and his/her social relationship with
614 others, etc. If the user did not consent to us know-
615 ing all stated information, then knowing them is
616 simply a privacy intrusion. Moreover, acquiring
617 such information about the users can lead to mental
618 and physical harm. Knowing what the user likes or
619 dislikes can easily affect him/her and can be detri-
620 mental either mentally or physically (for example,
621 manipulating the user to do something dangerous).
622 These are the main reasons why the Pandora dataset
623 (Gjurković et al. (2021)) is not a public dataset, and
624 to use it, you need to submit a request explaining
625 why you are seeking the use of this dataset. Also,
626 the authors of this dataset employ rigorous terms
627 of use (Irina Masnikosa and Bakić (2020)) to pro-
628 tect the users within the dataset. For example, one
629 cannot transfer or reproduce any part of the dataset
630 and attempt to identify or contact any user in the
631 dataset. One cannot publicly display users’ names
632 and sensitive information and messages. Also, one
633 can report findings publicly only on an aggregate
634 level. We believe that the user has the right to keep
635 his/her personality private. Whether personality
636 is consciously or unconsciously revealed in any
637 way, it is the other person’s responsibility to act
638 diligently and protect the shared information to pre-
639 vent from putting anybody in harm’s way. There-
640 fore, our work does not expose any users’ private
641 information, and we do not take users’ unique iden-
642 tifiers or demographic information to predict their
643 personalities. Our predictive model only focuses on
644 the posted users’ social media textual contents. In
645 other words, we do not focus on " who" posted the
646 content but rather on the content itself. Using only

647 the textual content to predict the individual’s per-
648 sonality helps us effectively reduce privacy intru-
649 sion risks. Our work is extremely valuable and can
650 improve many service providers. Only using the
651 content of the users’ posted texts without employ-
652 ing specific users’ information helped us reduce
653 the privacy intrusion issues. However, we think
654 that our model is limited in providing compelling
655 encrypted personality predictions. For now, our
656 model only predicts the personality traits in their
657 original forms. However, it would be more secure
658 in predicting them in an encrypted way. Therefore,
659 we aim to enhance the capability of our model by
660 introducing an encryption mechanism for the pre-
661 dicted results. We believe that it is essential for our
662 personality predictive model to be used in the right,
663 protected, and secured environment that does not
664 harm the users or reveal their personalities in any
665 way.

666 6 Conclusion

667 This work highlighted the effectiveness of using a
668 multi-task learning approach on top of a pre-trained
669 Electra model for the personality prediction task.
670 We also highlighted the effect of sharing weights
671 between the two popular personality trait tests.
672 Empirical results demonstrate that using shared
673 weights between MBTI and Big5 personality
674 tests outperforms state-of-the-art results for both
675 systems on different metrics. Our results show
676 that both personality systems are correlated. Also,
677 we found that despite the good Big5 regression
678 results of our solution, it seems like our model is
679 incapable of effectively classifying the Big5 traits
680 from the regression values. More weight sharing,
681 contextual information, and prediction heads will
682 be considered in future work as we are curious
683 to know the effect of demographical information
684 such as age, gender, and country on personality
685 detection.
686

687 References

- 688 Giulio Carducci, Giuseppe Rizzo, Diego Monti, Enrico
689 Palumbo, and Maurizio Morisio. 2018. *Twitperson-*
690 *ality: Computing personality traits from tweets using*
691 *word embeddings and supervised learning*. *Informa-*
692 *tion (Switzerland)*, 9.
- 693 Fabio Celli, Fabio Pianesi, David Stillwell, and Michal
694 Kosinski. 2013. Workshop on computational person-
695 ality recognition: Shared task.

696	Mark Batey David John Hughes, Moss Rowe, Giuseppe Riccardi Andrew Lee, and Fabio Pianesi. 2012. A tale of two sites: Twitter vs. facebook and the personality predictors of social media usage. <i>Computers in Human Behavior</i> , pages 561–569.	751
697		752
698		753
699		754
700		755
701	Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. <i>ArXiv</i> , abs/1810.04805.	756
702		757
703		758
704		759
705	Joan-Isaac Biel Daniel Gatica-Perez Giuseppe Riccardi Fabio Celli, Bruno Lepri and Fabio Pianesi. 2014. The workshop on computational personality recognition. <i>ACM</i> , pages 1245–1246.	760
706		761
707		762
708		763
709	Adrian Furnham. 1996. The big five versus the big four: the relationship between the myers-briggs type indicator (mbti) and neo-pi five factor model of personality. <i>Personality and Individual Differences</i> , pages 303–307.	764
710		765
711		766
712		767
713		768
714	Matej Gjurković, Mladen Karan, Iva Vukojević, Mihaela Bošnjak, and Jan Snajder. 2021. PANDORA talks: Personality and demographics on Reddit. In <i>Proceedings of the Ninth International Workshop on Natural Language Processing for Social Media</i> , pages 138–152, Online. Association for Computational Linguistics.	769
715		770
716		771
717		772
718		773
719		774
720		775
721	Lewis R Goldberg. 1993. The structure of phenotypic personality traits. <i>American psychologist</i> , pages 26–34.	776
722		777
723		778
724	Ivan Crnomarković Jan Šnajder Josip Jukić Matej Gjurković Mihaela Bošnjak Mladen Karan Irina Masnikosa, Iva Vukojević and Sara Bakić. 2020. Pandora. Preprint at https://psy.takefab.hr/datasets/all/pandora/ Last visited 09-01-2022.	779
725		780
726		781
727		782
728		783
729		784
730	Mary H. Mc Caulley Isabel Briggs Myers and Allen L. Hammer, editors. 1987. <i>Introduction to Type: A description of the theory and applications of the Myers-Briggs type indicator</i> . Consulting Psychologists Press.	785
731		786
732		787
733		788
734		789
735	Mitchell J. 2017. (mbti) myers-briggs personality type dataset. Preprint at https://www.kaggle.com/datasnaek/mbti-type Last visited 1-06-2022.	790
736		791
737		792
738		793
739	Quoc V. Le Kevin Clark, Minh-Thang Luong and Christopher D. Manning. 2020. Electra: Pre-training text encoders as discriminators rather than generators. <i>ArXiv</i> , abs/2003.10555.	794
740		795
741		796
742		797
743	Yang Li, Amirmohammad Kazameini, Yash Mehta, and E. Cambria. 2021. Multitask learning for emotion and personality detection. <i>ArXiv</i> , abs/2101.02346.	798
744		799
745		800
746	Manjula Ramannavar Mayuri Pundlik Kalghatgi and Nandini S Sidnal. 2015. A neural network approach to personality prediction based on the big-five model. <i>International Journal of Innovative Research in Advanced Engineering (IJIRAE)</i> , pages 56–63.	801
747		802
748		803
749		804
750		805
	Simine Vazire Nicholas Holtzman Samuel D Gosling, Adam A Augustine and Sam Gaddis. 2011. Manifestations of personality in online social networks: Self-reported facebook-related behaviors and observable profile information. <i>Cyberpsychology, Behavior, and Social Networking</i> , pages 483–488.	806
		807
	Ruth E Appel Sandra C Matz and Michal Kosinski. 2020. Privacy in the age of psychological targeting. <i>Current opinion in psychology</i> .	808
		809
	Haolan Ouyang Xiaojun Quan Tao Yang, Feifan Yang. 2021. Psycholinguistic tripartite graph network for personality detection. <i>The Joint Conference of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing</i> .	810
		811
	Rini Wongso Yen Lina Prasetyo et al Tommy Tandera, Derwin Suhartono. 2017. Personality prediction system from facebook users. <i>Procedia computer science</i> , page 604–611.	812
		813
	Michal Kosinski Wu Youyou and David Stillwell. 2015. Computer-based personality judgments are more accurate than those made by humans. <i>Proceedings of the national academy of sciences</i> , pages 1036–1040.	814
		815
	Feifan Yang, Tao Yang, Xiaojun Quan, and Qinliang Su. 2021. Learning to answer psychological questionnaire for personality detection. In <i>Findings of the Association for Computational Linguistics: EMNLP 2021</i> , pages 1131–1142.	816
		817
	Amir mohammad Kazameini Clemens Stachl Erik Cambria Yash Mehta, Samin Fatehi and Sauleh Eetemadi. 2020. bottom-up and top-down: Predicting personality with psycholinguistic and language model features. <i>2020 IEEE International Conference on Data Mining (ICDM)</i> , pages 1184–1189.	818
		819
	Bee Chin Ng Yosephine Susanto, Andrew Livingstone and Erik Cambria. 2020. The hourglass model revisited. in <i>IEEE Intelligent Systems</i> , pages 96–102.	820
		821
		822
		823
		824
		825
		826
		827
		828
		829
		830
		831
		832
		833
		834
		835
		836
		837
		838
		839
		840
		841
		842
		843
		844
		845
		846
		847
		848
		849
		850

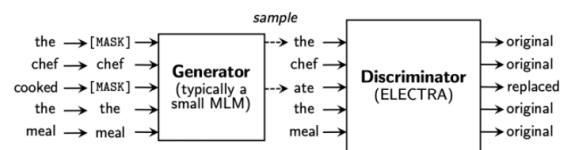


Figure 5: Electra Architecture Kevin Clark and Manning (2020)

795 Unlike BERT which heavily relies on the MLM
796 approach during the pre-training phase, this model
797 uses a new training approach called the replaced
798 token detection approach. Figure 5 highlights the
799 Electra model architecture. The Electra architec-
800 ture combines both a generator and discriminator
801 components. The generator will be trained using
802 the MLM goal and the discriminator will try to
803 predict for each word provided by the generator
804 whether it has been replaced or not. Therefore in-
805 stead of only knowing the 15% masked words in
806 the sentence as BERT does, this model will have
807 knowledge of all the tokens within the sentence
808 and predict whether it is the original token or the
809 replaced one. Having knowledge about all the
810 words instead of only 15% of them gives the Electra
811 model much more insights about the context within
812 a group of words. Moreover, using the discrimina-
813 tor as a binary classifier to predict whether the word
814 has been replaced or not will help the model gain
815 time during the training phase. As binary classifi-
816 cation is less computationally expensive compared
817 to the word generation task. To effectively train
818 this model the authors propose two losses, one for
819 the generator L_{MLM} (equation 7), and one for the
820 Discriminator L_{Disc} (equation 8).

$$L_{MLM}(x, \theta_G) = E\left(\sum_{i \in m} -\log p_G(x_i/x^{masked})\right) \quad (7)$$

821 θ_G is the generator learning parameters, x_i is the
822 current token input and x^{masked} is the replacement
823 tokens vector.
824

$$L_{Disc}(x, \theta_D) = E\left(\sum_{t=1}^n -1(x_t^{corrupt} = x_t) \log D(x_t^{corrupt}) - 1(x_t^{corrupt} \neq x_t) \log(1 - D(x_t^{corrupt}))\right) \quad (8)$$

826 θ_D defines the discriminator learning parameters, 1
827 defines the indicator function and $x_t^{corrupt}$ defines
828 the replaced token. To train both the generator
829 and the discriminator in an END-2-END process
830 the authors combined both losses into a single loss
831 function (equation 9) with the addition of a new
832 penalty term λ for the discriminator loss.

$$\min_{\theta_G, \theta_D} \sum_{x \in X} L_{MLM}(x, \theta_G) + \lambda L_{Disc}(x, \theta_D) \quad (9)$$

Table 4: The models hyperparameters

Hyper-Parameter	Value
Epochs number	10
Optimizer	Adam optimizer
Learning rate	2e-5
Weight decay	0.01
Activation function	LeakyRelu
Dropout degree	0.4
Classification loss (CL)	BCE with logits loss
Regression loss (RL)	MSE loss
Global loss	(CL+RL)/2
Batch size	15
Trainable parameters	13542167

A.2 The training hyperparameters

Supplementary Material

Datasets supplementary material:

- [MyPersonality dataset](#)
- [MBTI Personality Type dataset](#)
- [Big Five personality traits explanation](#)
- [MBTI personality traits explanation](#)
- [Pendora dataset request platform](#)

Models card supplementary material:

- [OC-EP model card](#)
- [OR-EP model card](#)
- [EWS-EP model card](#)
- [AWS-EP model card](#)
- [AWS-EP model code](#)

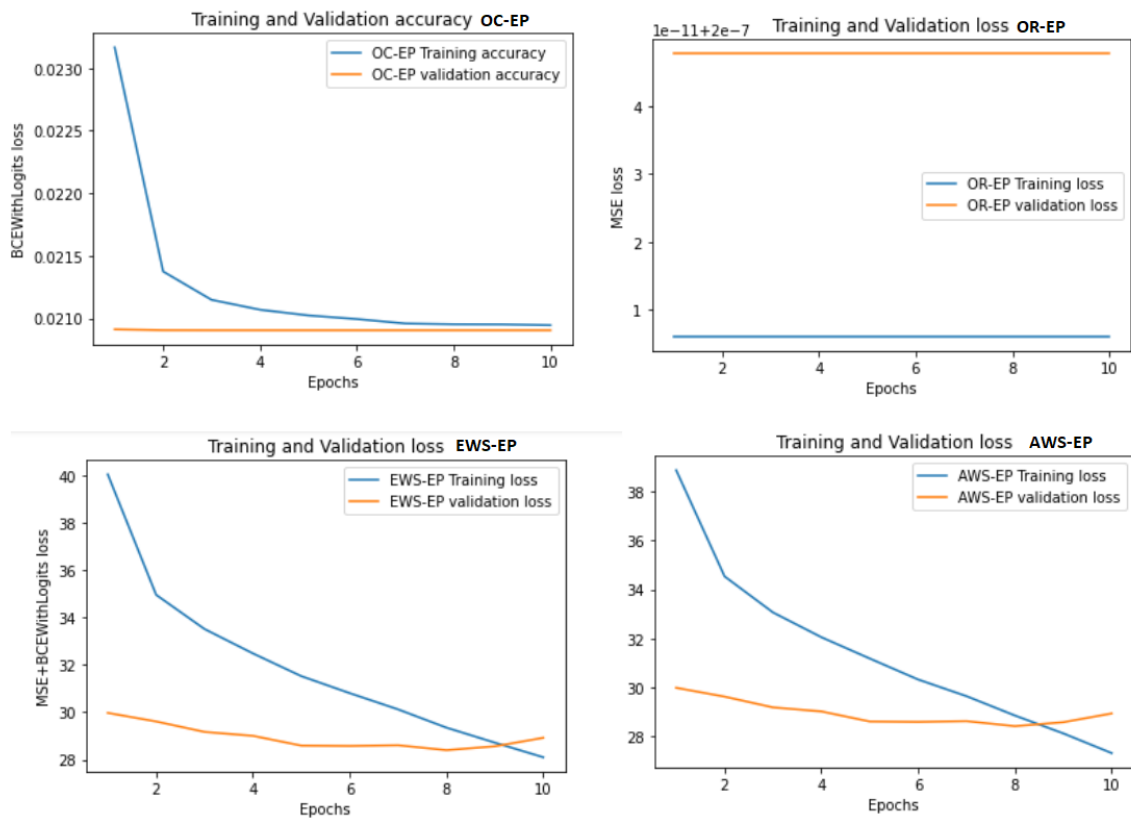


Figure 6: OC-EP, OR-EP, EWS-EP, and AWS-EP training performances