

# Top-Down Influence? Predicting CEO Personality and Risk Impact from Speech Transcripts

Anonymous ACL submission

## Abstract

How much does a CEO’s personality influence the performance of their company? Past literature has contested the possibility of predicting the Myers–Briggs Type Indicator (MBTI) from purely textual data. However, we use Transformers to create the first supervised model to regress the MBTI personality of CEOs. We show that moderate to strong predictions can be obtained for three out of four MBTI dimensions. Finally, providing empirical evidence for the *upper echelons theory*, we demonstrate that the predicted CEO personalities have explanatory power of financial risk.

## 1 Introduction

How much influence does the personality of a chief executive officer (CEO) have on the performance of their company? The personal news and antics of famous CEOs like Elon Musk, Jeff Bezos, or Bill Gates make headlines, and their personalities sometimes generate a cult-like following. But what effect do they really have? The *upper echelons theory* (Hambrick and Mason, 1984) suggests that the personalities of CEOs are also reflected in the organizational outcomes of their companies. However, presumably due to the lack of labeled data, no supervised models exist to detect CEOs’s personalities from text and infer their effect on the financial performance of companies.

Psychometric natural language processing (NLP), which is concerned with measuring personality and other psychological traits from text, enables downstream tasks such as political sentiment analysis (Golbeck and Hansen, 2011), deception detection (Enos et al., 2006), and mental health counseling (Calvo et al., 2017). The advent of deep learning has provided alternatives to traditional approaches which are usually based on psycholinguistic dictionaries. Here, we explore Transformer-based models (Devlin et al., 2019; Liu

et al., 2019) to predict personality from the transcripts of earnings calls in a regression task.

**Contributions** We present the first supervised model to predict the Myers–Briggs Type Indicator (MBTI) personality of CEOs. Leveraging crowd-sourced personality votes, we propose an alternative operationalization of the MBTI as a continuum rather than a binary label. We show that this MBTI representation correlates with the Big 5 and can be predicted in a text regression task. Finally, we demonstrate that the predicted MBTI personalities of CEOs have explanatory power of financial risk.

## 2 Background and Related Work

Various personality measures exist in the literature. This section describes the personality model we explore (MBTI), the de-facto standard model (Big 5), and approaches to predict both representations of personality from text.

### 2.1 MBTI

The MBTI is named after Katherine Cook Briggs and Isabel Briggs Myers. They developed it based on the work of the analytical psychologist Carl Jung (Briggs-Myers and Myers, 1995). According to the MBTI, personalities can be classified binarily along the following axes:

- *extraversion vs. introversion* (E–I): describing an out- or inward-oriented social attention;
- *sensing vs. intuition* (S–N): information processing based on perceivable/known facts or conceptualization and imagination;
- *thinking vs. feeling* (T–F): decision-making based on logic and rationality or emotions and empathy;
- *judging vs. perceiving* (J–P): quick judgement and organized action or observation and improvisation on-the-go.

076 Combined, the four labels form one of 16 personal-  
077 ity types (e.g., “ENTJ”). The MBTI is widely used  
078 in human resources management and by laypeople  
079 as a tool for self-exploration.

080 Psychological literature, however, has called as-  
081 sumptions of the MBTI into question. For example,  
082 McCrae and Costa (1989) find no evidence that  
083 personality can be binarized or distinguished into  
084 16 different types. In addition, they find moderate  
085 to strong correlations between MBTI and Big 5  
086 (McCrae et al., 2010), which is described in greater  
087 detail below (cf. Section 2.2). We re-assess these  
088 correlations in our dataset and explore a continuous  
089 representation of the MBTI in line with Big 5.

090 **MBTI Prediction from Text** In a literature  
091 study on text-based personality detection and a  
092 subsequent annotation study, Štajner and Yenikent  
093 (2020, 2021) conclude that predicting the MBTI  
094 from textual data is a difficult task. They hypothe-  
095 size that this is due to the theoretical and qualitative  
096 origin of the index, which distinguishes it from the  
097 empirical and quantitative Big 5. In particular, the  
098 dimensions *sensing* vs. *intuition* (S–N) and *judging*  
099 vs. *perceiving* (J–P) depend on behavioral rather  
100 than linguistic signals (Štajner and Yenikent, 2020,  
101 p. 6291).

102 In a field survey of project managers, Cohen et al.  
103 (2013) show that managers are significantly more  
104 often of the *intuitive* (N) and *thinking* (T) type than  
105 the general population. We observe a similar pat-  
106 tern in our dataset (cf. Section 3.1, Figure 2). Clas-  
107 sifying the MBTI of Twitter users based on count-  
108 based features, gender, and tweet *n*-grams, Plank  
109 and Hovy (2015) outperform a majority class base-  
110 line for the E–I and the T–F dimensions. Gjurković  
111 and Šnajder (2018) predict the self-reported MBTI  
112 of Redditors with support vector machine (SVM)  
113 and multilayer perceptron (MLP) models based on  
114 linguistic and activity-level features. Their model  
115 outperforms a majority class baseline across all di-  
116 mensions with the best results for E–I, followed by  
117 S–N, J–P, and T–F.

118 We compare the best-performing approaches  
119 identified by prior MBTI prediction studies (*n*-  
120 grams and Linguistic Inquiry and Word Counts  
121 (LIWC) dictionaries with SVMs and MLPs) to  
122 novel Transformer architectures. Furthermore, we  
123 consider a different domain (spoken financial dis-  
124 closures) and perform a regression instead of a  
125 classification.

## 2.2 Big 5 126

127 The Big 5 are the established psychometric model.  
128 Here, personality is represented as a continuum  
129 along the five axes *openness*, *conscientiousness*,  
130 *extraversion*, *agreeableness*, and *neuroticism* (Mc-  
131 Crae and John, 1992).

132 **Big 5 Prediction from Text** As part of the  
133 *myPersonality* project, Kosinski et al. (2015) find  
134 that the Big 5, IQ, and other properties of Face-  
135 book users can be predicted from their liked pages  
136 to varying degrees. Mairesse et al. (2007) create a  
137 text-based Big 5 prediction tool based on student  
138 essays and speech recordings.

139 Benischke et al. (2019) show that CEOs’ Big  
140 5 personalities moderate the relationship between  
141 CEO compensation and risk-taking. Hrazdil et al.  
142 (2020) use *IBM Watson Personality Insight* to pre-  
143 dict the Big 5 of C-level executives in earnings calls  
144 and find that an executive’s personality is associ-  
145 ated with their risk tolerance and company audit  
146 fees. Harrison et al. (2020) find that CEO Big 5  
147 are related to perceived firm risk and shareholder  
148 value. Another finding is that CEO *conscientious-*  
149 *ness* moderates the effect of financial risk on returns  
150 positively, while the opposite holds for *extroversion*  
151 and *neuroticism*.

152 Different to these approaches, we focus on the  
153 MBTI rather than the Big 5. We create the first  
154 supervised model to predict CEOs’ MBTI person-  
155 ality from text by collecting a new dataset of crowd-  
156 annotated MBTI profiles. This sets us apart from  
157 prior work using unsupervised approaches trained  
158 on out-of-domain corpora.

## 3 Personality Prediction 159

160 Using transcribed speech data as an input, we pre-  
161 dict the MBTI personality of CEOs with a text  
162 regression task. The following sheds light on the  
163 dataset collection and validation, methodology, and  
164 results.

### 3.1 Dataset Curation 165

166 For this task, we collect data from two sources: (1)  
167 text data and (2) crowd-sourced personality data.

168 **Text Data** We obtain 88K earnings call tran-  
169 scripts spanning the years 2002–2020 from Re-  
170 finitiv Eikon<sup>1</sup>. Earnings calls are quarterly telecon-  
171 ferences consisting of a scripted presentation and  
172 a spontaneous questions-and-answers (Q&A) ses-  
173 sion, in which company CEOs such as Elon Musk

<sup>1</sup><https://eikon.thomsonreuters.com/index.html>

MBTI	CEO Examples
Extraversion	Steve Jobs (Apple), Lisa Su (AMD), Mary Barra (General Motors)
Introversion	Rupert Murdoch (Fox), Mark Zuckerberg (Facebook), Sheldon Adelson (Las Vegas Sands)
Sensing	Jack Dorsey (Twitter), John Schnatter (Papa John’s), Marcus Lemonis (Camping World)
Intuition	Marissa Mayer (Yahoo), Bob Iger (Disney), Evan Spiegel (Snap)
Thinking	Elon Musk (Tesla), Tim Cook (Apple), Steve Ballmer (Microsoft)
Feeling	Sundar Pichai (Google), Howard Schultz (Starbucks), Naveen Jain (Infospace)
Judging	Jeff Bezos (Amazon), Larry Ellison (Oracle), Martha Stewart (Martha Stewart Living)
Perceiving	Larry Page (Alphabet), Martin Shkreli (Retrophin), Donald Trump (Trump Entertainment)

Table 1: CEO examples for each MBTI dimension from our dataset.

ELON MUSK (CEO): Thank you. So Q1 ended up being a strong quarter despite many challenges in the final few weeks. This is the first time we have achieved positive GAAP net income in a seasonally weak first quarter. Even with all the challenges, we achieved a 20% automotive gross margin, excluding regulatory credits, while ramping 2 major products. What we’ve learned from this is that—we’ve obviously learned a lot here.

Figure 1: Excerpt of Tesla’s Q1 2020 earnings call.

Unit	$\Sigma_x$	$\bar{x}$	$\min_x$	$\max_x$
utterances	13,183	17.91	2	124
sentences	111,781	151.88	2	563
tokens	2,526,473	3432.71	22	9968

Table 2: Statistics of the CEO–call data considered for the personality prediction. Sums ( $\Sigma_x$ ), averages ( $\bar{x}$ ), minima ( $\min_x$ ), and maxima ( $\max_x$ ) are computed across all earnings calls ( $n = 736$ ).

answer questions of banking analysts. This part is characterized by personal style, making it more authentic than written disclosures. Figure 1 shows an excerpt of Tesla’s first-quarter earnings call in 2020.

Given the dialogue nature of the calls, we need to map utterances to individual CEOs as we are not interested in the personality of the analysts. We identify CEO names with regular expressions and minimal preprocessing (e.g., stripping middle name initials or titles). For each CEO, we retrieve all utterances in the presentation and the Q&A session of the calls. Thus, we obtain a mapping of earnings calls ( $n = 88\text{K}$ ), CEOs ( $n = 12.4\text{K}$ ), and utterances ( $n = 157\text{K}$ ).

**Personality Data** We obtain MBTI personality labels for the CEOs from *Personality Database*<sup>2</sup>, which provides crowd-sourced personality profiles for celebrities, managers, and other noteworthy people. While each profile features vote results for the four dimensions of the MBTI, a minority also contains results for the Big 5. We find that 32 CEOs (e.g., Elon Musk and Steve Jobs) from our earnings call sample have at least three MBTI votes available. These CEOs participate in a total of 736 earnings calls. Table 2 gives the descriptive statistics of the merged text–personality data and Table 1 contains example CEOs from our dataset

across the MBTI.

Instead of representing each personality as one of 16 types, we represent each personality profile as a vector of 4 continuous variables ranging from 0 to 1, based on the crowd-sourced votes. We normalize the votes for the right-hand side of a scale  $s$  by the total votes:

$$\text{personality}_s = \frac{\text{votes}_{1,s}}{\text{votes}_{0,s} + \text{votes}_{1,s}}. \quad (1)$$

For example, for the E–I scale, we divide the votes for introversion (I) by the total number of votes for E and I. The resulting number is thus the likelihood of the CEO being intro- or extroverted. This representation is similar to the Big 5 model (excluding the *neuroticism* dimension) and allows for a more granular representation of personality than the usual operationalization of the MBTI. Figure 2 shows the distributions of the such obtained continuous labels. Most CEOs in our sample are rather *extroverted*, *intuitive*, *thinking*, and *judging* (cf. Figure 2), which corresponds to the ENTJ “Decisive Strategist” MBTI type.<sup>3</sup>

**Internal Validation** To assess the validity of the crowd-sourced votes, we analyze the inter-annotator agreement for all MBTI raters (cf. Table

<sup>2</sup><https://www.personality-database.com/>

<sup>3</sup><https://eu.themyersbriggs.com/en/tools/MBTI/MBTI-personality-Types/ENTJ>

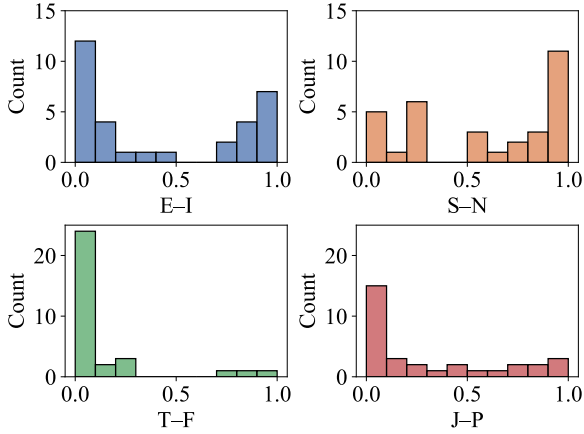


Figure 2: Label distributions for all CEOs considered in the personality prediction ( $n = 32$ ) across the MBTI dimensions *extraversion–introversion* (E–I), *sensing–intuition* (S–N), *thinking–feeling* (T–F), and *judging–perceiving* (J–P).

MBTI	$p_a$	$\alpha$	$\kappa_{bp}$	$\gamma$
E–I	87.45	0.40	0.75	0.76
S–N	80.20	0.43	0.60	0.62
T–F	83.33	0.14	0.67	0.71
J–P	90.62	0.17	0.81	0.88

Table 3: IAA per MBTI dimension in terms of percentage agreement ( $p_a$ ), Krippendorff’s  $\alpha$ , Brennan–Prediger coefficient ( $\kappa_{bp}$ ), and Gwet’s  $\gamma$ .

3). While  $p_a$  is high with values ranging between ca. 80 and 90%, Krippendorff’s  $\alpha$  (Krippendorff, 2013) yields only slight to moderate values between 0.14 and 0.43. Quarfoot and Levine (2016) call this phenomenon the “frequency distribution paradox,” where highly skewed label distributions combined with high percentage agreements can lead to low values of  $\alpha$ . As alternative measures robust to this undesirable property, they propose the Brennan–Prediger coefficient  $\kappa_{bp}$  (Brennan and Prediger, 1981) and Gwet’s  $\gamma$  (Gwet, 2008), which in our case yield a high IAA between 0.60 to 0.88.

**External Validation** To get a notion of external validity, we construct a correlation matrix between the crowd-based MBTI and Big 5 votes of all 2.2K profiles with more than three votes available on *Personality Database* (cf. Figure 3). According to McCrae and Costa (1989) and subsequent work (Furnham, 1996; Furnham et al., 2003), strong correlations should exist between MBTI *introversion* and Big 5 *extraversion* ( $r = -0.74$ ) as well as between MBTI *intuition* and Big 5 *openness* ( $r = 0.72$ ).

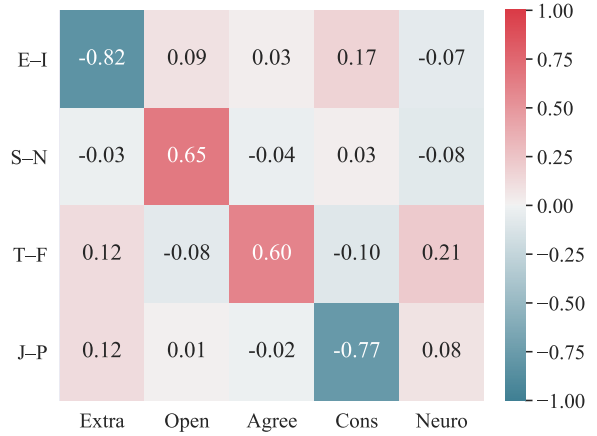


Figure 3: Correlation of MBTI (y-axis) and Big 5 (x-axis) scales for all profiles on the *Personality Database* with at least three votes ( $n = 2.2K$ ).

Furthermore, moderate correlations should exist between MBTI *feeling* and Big 5 *agreeableness* ( $r = 0.44$ ) and between MBTI *perceiving* and Big 5 *conscientiousness* ( $r = -0.49$ ). Our results confirm the findings of McCrae and Costa (1989) with similar correlations in the first two rows and stronger correlations in the third and fourth rows. This is most likely due to our increased sample size ( $n = 2.2K$  vs.  $n = 267$ ).

### 3.2 Methodology

For each of the 32 CEOs appearing in 736 CEO–call instances, we compare sparse approaches suggested by past literature to novel Transformer architectures for a regression of MBTI personality.<sup>4</sup>


**Data Split** We use an 80:10:10 to split our data into separate training ( $n = 568$ ), validation ( $n = 84$ ), and test sets ( $n = 84$ ). To avoid overfitting, we use sklearn’s `GroupShuffleSplit` with the CEO names as group splitting criterion, i.e., we split the data such that no CEO present in the training data appears in the validation or test data.

**Normalization** Given the highly skewed distributions, after the train–validation–test split, we apply a Box-Cox transformation (Box and Cox, 1964) to  $y$  with the following formula:

$$y(\lambda) = \begin{cases} \frac{y^\lambda - 1}{\lambda} & \text{for } \lambda \neq 0, \\ \ln(y) & \text{for } \lambda = 0. \end{cases} \quad (2)$$

<sup>4</sup>The supplementary material contains our implementation. Due to intellectual property restrictions of the earnings call data, we are in the process of clarifying how to make the dataset available.

We obtain  $\lambda$  via maximum-likelihood estimation. The resulting transformation makes the four label distributions more Gaussian-like by stabilizing variance.

**Transformers** We explore cased-vocabulary BERT<sub>base</sub> (12-layer, 768-hidden, 12-heads, 109M parameters) (Devlin et al., 2019) and RoBERTa<sub>base</sub> (12-layer, 768-hidden, 12-heads, 125M parameters) (Liu et al., 2019) models with a linear regression head. The models are trained with a maximum sequence length of 512 and a sliding window approach. We determine the training batch size and learning rate by running a Bayesian optimization over the grid of batch sizes  $b \in \{32, 64, 128, 256\}$  and learning rates  $l \in [0, 5 \times 10^{-5}]$ .<sup>5</sup> We train a model for up to 10 epochs and early stopping with a patience of one epoch. For each of the four MBTI dimensions, we evaluate 40 combinations of hyperparameters and select the model with minimal loss on the validation set. Different to the mean-squared error (MSE) loss, which is implemented per default in the  Transformers (Wolf et al., 2020) regressors, we minimize the L1 or alternatively called mean absolute error (MAE) loss, which is less sensitive to outliers.

**Sparse Methods** We also explore the sparse representations suggested by Plank and Hovy (2015) and Gjurković and Šnajder (2018). These include term frequency-inverse document frequency (tf-idf) vectors with  $n$ -grams of length  $n \in \{1, 2, 3\}$  and dictionary features across all dimensions of LIWC 2015 (Pennebaker et al., 2015) fed into SVM and three-layer MLP regressors. We compare all possible feature-algorithm combinations with respect to their average MAE on the validation set and select the combination with the lowest error. This was achieved with trigram tf-idf vectors fed into an SVM.

**Evaluation** The final model performance is evaluated on the test set in terms of the correlation coefficient Pearson’s  $r$  and the rank coefficients Spearman’s  $\rho$ , and Kendall’s  $\tau$ . While  $r$  measures linear relationships,  $\rho$  and  $\tau$  measure monotonic relationships and are more robust to outliers. In addition, we consider the error measure MAE, which is the minimized loss function of the transformers. In case of a tie, we give precedence to  $\tau$ , as this measure is least sensitive to outliers and particu-

larly suited for small sample sizes.

### 3.3 Results and Discussion

The results of the personality prediction task are depicted in Table 4. An SVM performs competitive, especially for the dimensions E-I ( $\tau = 0.44$ ) and S-N ( $\tau = 0.20$ ). While the SVM outperforms BERT for all dimensions except for J-P, RoBERTa achieves the best results in most cases.

The largest correlations across all models are achieved for the *extraversion-introversion* (E-I) scale with strong linear and rank correlations for the RoBERTa regressor ( $r = 0.70$ ,  $\rho = 0.66$ ). This result is not surprising, as distinguishing between *extra-* and *introverted* CEOs should be comparably easy to make based on linguistic style. This is followed by the *sensing-intuition* (S-N) scale with moderate to strong correlations ( $r = 0.45$ ,  $\rho = 0.53$ ) and the *judging-perceiving* (J-P) scale with weak to moderate correlations ( $r = 0.40$ ,  $\rho = 0.36$ ). The worst results are obtained for the *thinking-feeling* (T-F) scale, with the SVM and RoBERTa obtaining correlations of around zero and BERT even obtaining weak to moderate negative correlations. There are several possible explanations for this: Conceptually, it could be the case that this dimension simply can not be captured by analyzing linguistic data. Furthermore, the predictive power could be low due to the comparably small sample size. Lastly, we hypothesize that the skewness of the label distribution, which was the highest across all MBTI dimensions for the T-F scale (cf. Figure 2), has contributed to the weak performance. This warrants further research to explore whether our findings can be confirmed for larger datasets with less skewed label distributions.

Štajner and Yenikent (2020) hypothesize that the S-N and J-P dimensions should theoretically make for the worst candidates in a text-based personality prediction task since they capture behavioral rather than linguistic dimensions of personality. Although our regressors perform worse on these dimensions than for the *extraversion-introversion* dimension, they still achieve moderate to strong correlations, showing that even the more latent dimensions of personality can be uncovered with NLP.

**Qualitative Analysis** As a brief qualitative analysis, we use Shapley Additive Explanations (SHAP) developed by Lundberg and Lee (2017) to visualize the personality predictions for an exemplary text snippet across the four MBTI dimensions

<sup>5</sup>Final hyperparameter choices and results on our validation set can be found in Appendices A and B.

Thank you. So Q1 ended up being a strong quarter despite many challenges in the final few weeks. This is the first time we have achieved positive GAAP net income in a seasonally weak first quarter. Even with all the challenges, we achieved a 20% automotive gross margin, excluding regulatory credits, while ramping 2 major products. What we've learned from this is that -- we've obviously learned a lot here.

(a) Result of the E-I regressor.

Thank you. So Q1 ended up being a strong quarter despite many challenges in the final few weeks. This is the first time we have achieved positive GAAP net income in a seasonally weak first quarter. Even with all the challenges, we achieved a 20% automotive gross margin, excluding regulatory credits, while ramping 2 major products. What we've learned from this is that -- we've obviously learned a lot here.

(b) Result of the S-N regressor.

Thank you. So Q1 ended up being a strong quarter despite many challenges in the final few weeks. This is the first time we have achieved positive GAAP net income in a seasonally weak first quarter. Even with all the challenges, we achieved a 20% automotive gross margin, excluding regulatory credits, while ramping 2 major products. What we've learned from this is that -- we've obviously learned a lot here.

(c) Result of the T-F regressor.

Thank you. So Q1 ended up being a strong quarter despite many challenges in the final few weeks. This is the first time we have achieved positive GAAP net income in a seasonally weak first quarter. Even with all the challenges, we achieved a 20% automotive gross margin, excluding regulatory credits, while ramping 2 major products. What we've learned from this is that -- we've obviously learned a lot here.

(d) Result of the J-P regressor.

Figure 4: Example snippet from our dataset (uttered by Elon Musk in Tesla’s Q1 2020 earnings call) with SHAP heatmap across the MBTI. Red indicates a positive and blue a negative influence on the prediction.

MBTI	Model	$r$	$\rho$	$\tau$	MAE
E-I	SVM	0.57	0.58	0.44	0.38
	BERT	0.39	0.35	0.22	0.59
	RoBERTa	<b>0.70</b>	<b>0.66</b>	<b>0.52</b>	<b>0.34</b>
S-N	SVM	0.32	0.36	0.20	0.30
	BERT	0.08	0.23	0.16	0.46
	RoBERTa	<b>0.45</b>	<b>0.53</b>	<b>0.38</b>	<b>0.28</b>
T-F	SVM	<b>0.03</b>	-0.12	-0.08	<b>0.37</b>
	BERT	-0.47	-0.41	-0.27	0.41
	RoBERTa	0.01	<b>-0.10</b>	<b>-0.07</b>	0.39
J-P	SVM	-0.05	0.04	0.02	<b>0.35</b>
	BERT	0.39	<b>0.38</b>	<b>0.25</b>	0.52
	RoBERTa	<b>0.40</b>	0.36	0.21	0.36

Table 4: Correlation results of the personality regression task. CEO personality is predicted across the MBTI dimensions *extraversion-introversion* (E-I), *sensing-intuition* (S-I), *thinking-feeling* (T-F), and *judging-perceiving* (J-P). SVM is trained on trigram tf-idf vectors, BERT<sub>base</sub>, and RoBERTa<sub>base</sub> on text. Best results in bold.

with heatmaps (cf. Figure 4). The analyzed personality is Elon Musk, who, according to the crowd votes, scores high on E-I (*introversion*) and on S-N (*intuitive*), low on T-F (*thinking*), and medium on J-P (*judging/perceiving*). Particularly, the results for T-F (cf. Figure 4c) are interesting, where statements related to factual content are related to increased T, and interpretative statements (e.g.,

“[e]ven with all the challenges”) to increased F.

## 4 Risk Regression

According to *upper echelons theory* (Hambrick and Mason, 1984), strategic choices and performance measures of organizations can be predicted by characteristics of their top management. As a use case for our personality prediction task, we explore whether we can find empirical support for this theory. We hypothesize that having a different personality to most CEOs (i.e., ENTJ, cf. Figure 2 and Cohen et al. (2013)) should translate into increased financial risk.

### 4.1 Dataset Curation

As a basis for the risk regression task, we take the sample of 88K earnings calls and merge it with financial data obtained from the databases CRSP and IBES.<sup>6</sup> To measure risk, we calculate the stock return volatility in the business week following each call as a label. We use the sample standard deviation of logarithmic stock returns for more robust measures. As features, we incorporate a comprehensive set of risk proxies (cf. Table 5) suggested by Price et al. (2012) and Theil et al. (2019).<sup>7</sup>

<sup>6</sup><https://wrds-www.wharton.upenn.edu>

<sup>7</sup>We initially also considered including a market volatility index (VIX), but decided against it as its low explanatory power and high variation inflation factor (VIF) indicated redundancy of this variable (Johnston et al., 2018).

Feature	Definition
Past Vola	Standard deviation of logarithmic returns in the business quarter before the call
Size	Market value of the firm, i.e., the number of outstanding shares times stock price one day before the call
BTM	Book-to-Market = book value of the firm divided by market value
SUE	Mean absolute deviation of analysts' earnings-per-share forecasts from the actual value in the preceding quarter
Spread	Difference between the stock's bid and ask price on the call date
Leverage	Total liabilities divided by assets
ROA	Return on Assets, i.e., net income divided by assets
Volume	Stock trading volume on the call date
Industry	Fama-French 12 industry dummies (e.g., <i>health</i> or <i>finance</i> ). <sup>8</sup>
Time	Year-quarter dummies

Table 5: Financial features used for the risk regression task. BTM is calculated following (Fama and French, 2001) and firms with a negative value are removed. Size, BTM, and volume are  $\log 1p$ -transformed.<sup>9</sup>

## 4.2 Methodology

We use the best-performing personality prediction model (RoBERTa) to infer the personality of the 12.4K unlabelled CEOs present in the 88K calls. Together with the financial covariates (see above), the predicted CEO MBTI is then used to explain short-term stock return volatility following the calls with multiple linear regression.<sup>10</sup> Volatility is the most common financial risk measure, and its prediction is an important task for firm valuation and financial decision-making. Importantly, “risk” is a purely descriptive concept in finance, as it measures the fluctuation of stock returns.

## 4.3 Results and Discussion

The results of this risk regression task are shown in Table 6 and Figure 5. We find that all MBTI dimensions are significantly associated with risk following the call date. This significance is high ( $p \leq 0.001$ ) for all scales except J–P. The direction of this association behaves as expected for all MBTI scales: a CEO communicating in an *introverted*, *feeling* and *perceptive* manner is associated with increased risk ( $\beta_i = 0.03$ ,  $\beta_f = 0.11$ ,  $\beta_p = 0.01$ ), while an *intuitive* communication is associated with decreased risk ( $\beta_s = -0.02$ ). We compare these results with the predominant personality type of CEOs in Cohen et al. (2013) and our

<sup>10</sup>The supplementary material contains our dataset and implementation.

Feature	FIN	FIN + MBTI
E–I		0.03*** (8.72)
S–N		-0.02*** (-6.19)
T–F		0.11*** (31.47)
J–P		0.01* (2.35)
Past Vola	0.38*** (87.71)	0.37*** (83.98)
Size	-0.21*** (-41.19)	-0.22*** (-42.21)
Volume	0.11*** (23.70)	0.10*** (22.84)
$n$	87,826	87,826
Adj. $R^2$	32.50%	33.30%
AIC	214,900	213,800
BIC	215,700	214,700

\* $p \leq 0.05$ , \*\* $p \leq 0.01$ , \*\*\* $p \leq 0.001$

Table 6: Results of the risk regression with  $z$ -standardized coefficients and  $t$ -statistics in parentheses. The sample consists of 88K earnings calls spanning 5K firms and years 2002–2020. Regressions include industry- and time-fixed effects. FIN is a model with just the financial features (defined in Section 4.1) and FIN + MBTI is a joint model including the MBTI (E–I, S–N, T–F, and J–P). For brevity, only the three financials with the largest coefficients are included; graphical results for all features are shown in Figure 5.

study (cf. Figure 2). It seems that sharing the modal personality (being *extroverted*, *intuitive*, *thinking*, and *judging*) correlates with decreased risk.

Less surprisingly, a larger past volatility and a smaller firm size are strongly associated with increased levels of risk following an earnings call (see Figure 5). Notably, T–F has the third-largest association with future risk ( $\beta_f = 0.11$ ). Though only weakly correlated with the ground truth (cf. Table 4), the results suggest that the predictions for this scale contain strong economic signal for risk regression. Lastly, incorporating the MBTI leads to a slight increase in adjusted  $R^2$  and decreases in Akaike information criterion (AIC) and Bayesian information criterion (BIC), indicating that including personality as a feature helps to explain variation of risk.

In sum, these results provide new empirical evidence to support the *upper echelons theory*. We show that situational aspects of CEO personality, predicted with our MBTI regressor, also reflect firm performance measured by stock return volatility, the most common financial risk measure.

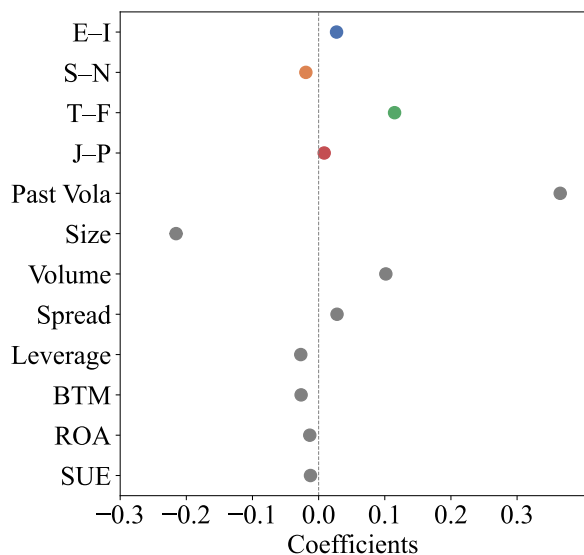


Figure 5:  $z$ -standardized feature coefficients of the risk regression task ( $n = 88\text{K}$ ) with 95% confidence intervals (error bars are smaller than the symbol size). Regressions include industry- and time-fixed effects. Financial features are defined in Section 4.1 and numerical results are presented in Table 6.

## 5 Ethical Considerations

In the following, we discuss possible biases and environmental considerations related to a personality prediction task from text.

**Social Desirability Bias** Past literature has uncovered a trait-related social desirability bias for the Big 5: This bias is most pronounced for the *neuroticism* trait (which is omitted in the MBTI), followed by *conscientiousness* and *agreeableness*, and to a lesser extent for *extraversion* and *openness* (Ones et al., 1996, Table 2). For the MBTI, in contrast, there exist no “bad” personality traits. As shown in Section 3.1, however, the Big 5 and the MBTI correlate strongly. Therefore, the points raised about social desirability, albeit to a lesser extent, should apply to the MBTI, too.

**Sample Biases** Critically, our dataset comprises a limited sample of 32 CEOs of large American (mostly tech) companies. While these companies (Alphabet, Facebook, Apple, etc.) constitute a large share of the American market, this renders the personality prediction model less applicable to non-American, small, or non-tech companies. Out of the 32 CEOs present in the dataset, only four (i.e., 12.5%) are female. While this gender ratio is twice as high as that of the S&P 500 (Catalyst, 2021), this highlights that the findings of this study might

generalize poorly to non-male CEOs. In addition, as shown in Section 3.1, Figure 2, our findings coincide with those of Cohen et al. (2013) in the sense that CEOs as a social cohort share a distinct distribution of personality traits, which is why we argue that the MBTI regressors should only be applied with caution, if at all, to non-CEO samples.

**Energy Consumption** Training neural models can have substantial financial and environmental costs (Strubell et al., 2019), which motivates us to discuss the computational efficiency of the proposed models. Using an NVIDIA Tesla P100 GPU, we run a hyperparameter optimization over 40 configurations per MBTI dimension for both BERT and RoBERTa. The average power consumption is 200W and the optimization takes ca. 16 hours, i.e., 3.2 kilowatt hours (kWh) with an electricity cost of 42 cents per model.<sup>11</sup> Labeling the 88K earnings call instances with no available ground truth takes ca. 18 hours and 140W, i.e., 2.52 kWh of GPU time and 33 cents, respectively. Training time of the trigram tf-idf representation with an SVM algorithm is negligible with training taking ca. 2 minutes on a quad-core processor with 8GB RAM. Whether the performance increases of the Transformers over a sparse method justify the added computational costs should be considered carefully on a case-by-case basis.

## 6 Conclusion and Future Work

We present the first text regression approach for predicting the MBTI personality of CEOs. Although past research has contested the possibility of predicting MBTI from purely textual data, we observe moderate to strong correlations with the ground truth for three out of four dimensions. In a risk regression task, we demonstrate that—consistent with the *upper echelons theory*—the predicted CEO personality is significantly associated with financial risk in the form of stock return volatility. Qualitatively, extroverted, intuitive, thinking, and judging CEOs seem to incur less financial risk.

In the future, we plan to model the personality prediction task as a multi-task learning problem, in which one single regressor is trained to predict all four MBTI dimensions at once. In addition, it would be interesting to incorporate speech signals of executives (e.g., voice modulation, tonality, and silence) into the personality predictions.

<sup>11</sup>Calculations assume the average U.S. electricity rate of 13.19 cents per 6 September 2021: <https://www.electricchoice.com/electricity-prices-by-state>



528  
529  
530  
531  
532  
533  
  
534  
535  
536  
  
537  
538  
539  
540  
  
541  
542  
543  
  
544  
545  
546  
547  
548  
  
549  
  
550  
551  
552  
553  
  
554  
555  
556  
557  
558  
  
559  
560  
561  
562  
563  
564  
  
565  
566  
567  
568  
  
569  
570  
571  
572  
573  
  
574  
575  
576  
577  
578

## References

Mirko H. Benischke, Geoffrey P. Martin, and Lotte Glaser. 2019. [CEO Equity Risk Bearing and Strategic Risk Taking: The Moderating Effect of CEO Personality](#). *Strategic Management Journal*, 40(1):153–177.

George E. P. Box and David R. Cox. 1964. [An Analysis of Transformations](#). *Journal of the American Statistical Association*, 26(2):211–252.

Robert L. Brennan and Dale J. Prediger. 1981. [Coefficient Kappa: Some Uses, Misuses, and Alternatives](#). *Educational and Psychological Measurement*, 41:687–699.

Isabel Briggs-Myers and Peter B. Myers. 1995. [Gifts Differing: Understanding Personality Type](#). Davies-Black.

Rafael A. Calvo, David N. Milne, M. Sazzad Husain, and Helen Christensen. 2017. [Natural Language Processing in Mental Health Applications Using Non-Clinical Texts](#). *Natural Language Engineering*, 23(5):649–685.

Catalyst. 2021. [Women CEOs of the S&P 500](#).

Yuval Cohen, Hana Ornoy, and Baruch Keren. 2013. [MBTI Personality Types of Project Managers and Their Success: A Field Survey](#). *Project Management Journal*, 44(3):78–87.

Jacob Devlin, Ming Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding](#). *Proceedings of NAACL*, pages 4171–4186.

Frank Enos, Stefan Benus, Robin L. Cautin, Martin Graciarena, Julia Hirschberg, and Elizabeth Shriberg. 2006. [Personality Factors in Human Deception Detection: Comparing Human to Machine Performance](#). In *Proceedings of Interspeech*, pages 813–816.

Eugene F. Fama and Kenneth R. French. 2001. [Disappearing Dividends: Changing Firm Characteristics or Lower Propensity to Pay?](#) *Journal of Financial Economics*, 60(1):3–43.

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](#). *Personality and Individual Differences*, 21(2):303–307.

Adrian Furnham, Joanna Moutafi, and John Crump. 2003. [The Relationship Between the Revised Neo-Personality Inventory and the Myers–Briggs Type Indicator](#). *Social Behavior and Personality*, 31(6):577–584.

Matej Gjurković and Jan Šnajder. 2018. [Reddit: A Gold Mine for Personality Prediction](#). In *Proceedings of the ACL Workshop on Computational Modeling of People’s Opinions, Personality, and Emotions in Social Media*, pages 87–97.

Jennifer Golbeck and Derek Hansen. 2011. [Computing Political Preference among Twitter Followers](#). In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, pages 1105–1108.

Kilem Li Gwet. 2008. [Computing inter-rater reliability and its variance in the presence of high agreement](#). *British Journal of Mathematics and Statistical Psychology*, 61:29–48.

Donald C. Hambrick and Phyllis A. Mason. 1984. [Upper Echelons: The Organization as a Reflection of Its Top Managers](#). *Academy of Management Review*, 9(2):193–206.

Joseph S. Harrison, Gary R. Thurgood, Steven Boivie, and Michael D. Pfarrer. 2020. [Perception Is Reality: How CEOs’ Observed Personality Influences Market Perceptions of Firm Risk and Shareholder Returns](#). *Academy of Management Journal*, 63(4):1166–1195.

Karel Hrazdil, Jiri Novak, Rafael Rogo, Christine Wiedman, and Ray Zhang. 2020. [Measuring Executive Personality Using Machine-Learning Algorithms: A New Approach and Audit Fee-Based Validation Tests](#). *Journal of Business Finance and Accounting*, 47(3–4):519–544.

Ron Johnston, Kelvyn Jones, and David Manley. 2018. [Confounding and Collinearity in Regression Analysis: A Cautionary Tale and an Alternative Procedure, Illustrated by Studies of British Voting Behaviour](#). *Quality & Quantity*, 52:1957–1976.

Michal Kosinski, Sandra C. Matz, Samuel D. Gosling, Vesselin Popov, and David Stillwell. 2015. [Facebook as a Research Tool for the Social Sciences: Opportunities, Challenges, Ethical Considerations, and Practical Guidelines](#). *American Psychologist*, 70(6):543–556.

Klaus Krippendorff. 2013. *Content Analysis: An Introduction to Its Methodology*, 3rd edition. Sage, Thousand Oaks (CA), USA.

Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. [RoBERTa: A Robustly Optimized BERT Pretraining Approach](#).

Scott M Lundberg and Su-In Lee. 2017. [A unified approach to interpreting model predictions](#). In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.

François Mairesse, Marilyn A. Walker, Matthias R. Mehl, and Roger K. Moore. 2007. [Using Linguistic Cues for the Automatic Recognition of Personality](#)

- 634           in Conversation and Text. *Journal of Artificial Intel-*  
635           *ligence Research*, 30:457–500. 688
- 636 Robert R. McCrae and Paul T. Costa. 1989. Reinter- 689  
637           preting the Myers–Briggs Type Indicator from the 690  
638           Perspective of the Five-Factor Model of Personality. 691  
639           *Journal of Personality*, 57(1):17–40. 692
- 640 Robert R. McCrae, Paul T. Costa, and Thomas A. Mar- 693  
641           tin. 2010. The NEO-PI-3: A More Readable Re-  
642           vised NEO Personality Inventory. *Journal of Per-*  
643           *sonality Assessment*, 84(3):261–270.
- 644 Robert R. McCrae and Oliver P. John. 1992. An In-  
645           troduction to the Five-Factor Model and Its Applica-  
646           tions. *Journal of Personality*, 60(2):175–215.
- 647 Deniz S. Ones, Chockalingam Viswesvaran, and An-  
648           gelika D. Reiss. 1996. Role of Social Desirabil-  
649           ity in Personality Testing for Personnel Selection:  
650           The Red Herring. *Journal of Applied Psychology*,  
651           81(6):660–679.
- 652 James W. Pennebaker, Ryan L. Boyd, Kayla Jordan,  
653           and Kate Blackburn. 2015. The Development and  
654           Psychometric Properties of LIWC2015. White pa-  
655           per, University of Texas at Austin.
- 656 Barbara Plank and Dirk Hovy. 2015. Personality Traits  
657           on Twitter—or—How to Get 1500 Personality Tests  
658           in a Week. In *Proceedings of the 6th Workshop*  
659           *on Computational Approaches to Subjectivity, Sen-*  
660           *timent and Social Media Analysis*, pages 92–98.
- 661 S. McKay Price, James S. Doran, David R. Peterson,  
662           and Barbara A. Bliss. 2012. Earnings Conference  
663           Calls and Stock Returns: The Incremental Informa-  
664           tiveness of Textual Tone. *Journal of Banking and*  
665           *Finance*, 36(4):992–1011.
- 666 David Quarfoot and Richard A. Levine. 2016. How Ro-  
667           bust Are Multirater Interrater Reliability Indices to  
668           Changes in Frequency Distribution? *The American*  
669           *Statistician*, 70:373–384.
- 670 Emma Strubell, Ananya Ganesh, and Andrew McCal-  
671           lum. 2019. Energy and Policy Considerations for  
672           Deep Learning in NLP. In *Proceedings of ACL*,  
673           pages 3645–3650.
- 674 Christoph Kilian Theil, Samuel Broscheit, and Heiner  
675           Stuckenschmidt. 2019. PRoFET: Predicting the  
676           Risk of Firms from Event Transcripts. In *Proceed-*  
677           *ings of IJCAI*, pages 5211–5217.
- 678 Sanja Štajner and Seren Yenikent. 2020. A Survey of  
679           Automatic Personality Detection from Texts. In *Pro-*  
680           *ceedings of COLING*, pages 6284–6295.
- 681 Sanja Štajner and Seren Yenikent. 2021. Why Is MBTI  
682           Personality Detection from Texts a Difficult Task?  
683           In *Proceedings of EACL*, pages 3580–3589.
- 684 Thomas Wolf, Lysandre Debut, Victor Sanh, Julien  
685           Chaumond, Clement Delangue, Anthony Moi, Pier-  
686           ric Cistac, Tim Rault, Rémi Louf, Morgan Fun-  
687           towicz, Joe Davison, Sam Shleifer, Patrick von

694 **A Hyperparameter Configurations**

695 Using a Bayesian hyperparameter optimization as  
 696 specified in Section 3.2, the following configura-  
 697 tions led to minimal loss on the validation set. Ta-  
 698 ble 7a summarizes the optimal configuration for  
 699 BERT and Table 7b the one for RoBERTa.

MBTI	Batch Size	Learning Rate
E-I	128	$4.8 \times 10^{-5}$
S-N	32	$4.9 \times 10^{-5}$
T-F	32	$1.0 \times 10^{-6}$
J-P	256	$8.6 \times 10^{-6}$

(a) Hyperparameters for BERT.

MBTI	Batch Size	Learning Rate
E-I	256	$4.3 \times 10^{-5}$
S-N	32	$4.6 \times 10^{-5}$
T-F	128	$9.4 \times 10^{-8}$
J-P	128	$4.7 \times 10^{-5}$

(b) Hyperparameters for RoBERTa.

Table 7: Final hyperparameter configurations found by the Bayesian optimization searching over 40 configurations per MBTI dimension.

700 **B Results on the Validation Set**

701 The results of the MBTI regressors on the valida-  
 702 tion set are depicted in Table 8.

MBTI	Model	$r$	$\rho$	$\tau$	MAE
E-I	SVM	0.70	0.69	0.55	0.38
	BERT	0.46	0.42	0.28	0.62
	RoBERTa	0.72	0.60	0.48	0.35
S-N	SVM	0.34	0.48	0.30	0.28
	BERT	0.20	0.35	0.24	0.53
	RoBERTa	0.43	0.61	0.43	0.27
T-F	SVM	0.13	-0.05	-0.03	0.33
	BERT	-0.43	-0.32	-0.22	0.38
	RoBERTa	0.11	-0.07	-0.03	0.36
J-P	SVM	-0.05	0.05	0.03	0.35
	BERT	0.32	0.28	0.19	0.53
	RoBERTa	0.25	0.14	0.06	0.40

Table 8: Results of the personality prediction task on the validation set.