

Research Article

News Sentiment and the Risk of a Stock Price Crash Risk: Based on Financial Dictionary Combined BERT-DCA

Shuyi Li  and Junhao Kong 

School of Economics, Zhejiang University, Hangzhou 310030, China

Correspondence should be addressed to Junhao Kong; 12101011@zju.edu.cn

Received 24 March 2022; Revised 15 June 2022; Accepted 29 June 2022; Published 31 July 2022

Academic Editor: Stefan Cristian Gherghina

Copyright © 2022 Shuyi Li and Junhao Kong. This is an open access article distributed under the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

This study combines a financial knowledge dictionary and pretraining method based on BERT (Bidirectional Encoder Representation from Transformers) to construct a deep learning model for identifying stock news sentiments. The study then calculates the sentiment metrics of all stocks and analyzes the impact of news sentiment on the risk of a stock price crash and its heterogeneity. The results show that stocks with more positive sentiment metrics have a higher risk of crash in the following year. We also investigate the information intermediation and investor sentiment channels by which news sentiment affects the risk of a crash. The results show that more net insider sales, lower information transparency, and less analyst coverage amplify the impact of news sentiment on future crash risk, which is consistent with the information intermediation channel. Additionally, more retail investor positions, more active investor sentiment, and divergence between analysts' opinions and news amplify the impact of news sentiment on the risk of a future stock price crash, which is consistent with the investor sentiment channel.

1. Introduction

Stability is a necessary condition for the financial market to efficiently facilitate economic development. In developing countries with weak financial infrastructure and low market efficiency, the risk of firm-specific stock price crash can seriously threaten investors and undermine their confidence in the stock market. The regularly used method of portfolio diversification cannot fully mitigate it, leading researchers to focus on the factors inducing stock price crash. A different strand of literature investigates the effects of media coverage or overall media tone on stock price crash risk, to which the present study belonged. We, however, adopted a slightly different perspective—the asymmetric impact of sentiment heterogeneity (i.e., positive or negative).

Our study aims to examine the asymmetric impact of news sentiment on stock price crash risk attributable to the role of media as information intermediary that influences investors' sentiment. On the one hand, regardless of the bias of media coverage, positive news will inevitably lead to investors' optimism. The media report positive news about

some firms, according to the theory of the spiral of silence [1], pessimistic investors may remain silent and optimistic investors may dominate the market. When short selling is not allowed, pessimism about that firm cannot be hedged normally. Subsequently, a stock price crash will happen when the optimistic sentiment dissipates and pessimistic investors become the marginal buyers [2]. On the other hand, the management has every intention of defending the firm's image and concealing negative information from the public [3, 4], which unintentionally leads to an overvaluation of stock prices, resulting in price bubbles. When the management can no longer delay or conceal negative news and has to release it to the market in a short period, the pessimism of investors will be amplified, inducing a stock price crash [5–8].

We distinguish among the sentiments of the media coverage and study the asymmetric impacts of positive and negative news on stock price crash risk. Thus, we need a precise measurement of the sentiment for objectively evaluating the tone and extent of the news. In this regard, a burgeoning body of finance and accounting literature has used natural language processing (NLP) algorithms to

extract the financial texts' sentiment [9–12]. The emergence of modern deep learning algorithms, such as ELMo, GPT (generative pretraining), and Google BERT, makes NLP move up a gear in the accuracy of sentiment analysis. Related studies have suggested that relative to traditional methods, modern deep learning algorithms that combine prior knowledge in certain fields have better performance. In this study, we customize a state-of-the-art deep learning NLP algorithm (BERT) for financial texts and document its advantages over traditional approaches.

We collect a total of 1,132,856 initial media coverage articles between 2011 and 2020 from the China Stock Market and Accounting Research database [13]. We manually construct our own sentiment dictionary in the financial domain and use it as a corpus for sentiment classification. Moreover, the BERT-based pretraining model is designed to help machines understand the characteristics of human language and extract sentiment information effectively. Based on the model, we classify the related news of each stock. Finally, our sample consisted of 17,267 firm-year observations representing 2,277 individual firms.

To begin with, we examine whether media sentiment is associated with a firm-specific future price crash risk. We measure media sentiment from three dimensions: mixed (on average), positive, and negative sentiments. We define the indicators of the average sentiment and those of the positive and negative sentiment, respectively. We use three proxies of stock price crash risk: the binary variable (CRASH) that equals 1 for a firm-year that experiences one or more crash weeks during the fiscal year and 0 otherwise; the negative coefficient of skewness of firm-specific weekly returns (NCSKEW); and the down-to-up volatility (DUVOL) of firm-specific weekly returns [3, 14, 16]. The results show that firms with more positive media coverage tend to have a higher risk of future stock price crash. Meanwhile, negative media coverage shows a limited effect on the risk of future stock price crash and a significant negative relation with current stock price crash risk, implying that they can lead to a stock price crash in the short term.

We then address the natural question of identifying the channels through which media sentiment affects the stock price crash risk. We hypothesize two possible channels, namely, information intermediation and investor sentiment, and then design various settings to examine them. The stock market is significantly impacted by the media, as an important information intermediary between firm management and market participants. On the one hand, media outfits disseminate value-relevant information on firms' current and future earnings to outside investors, reduce market frictions, improve investor perceptions, and mitigate information asymmetry [14, 17, 18]. On the other hand, media are not the perfect messengers. Media coverage is not always objective and neutral, but rather offers ambiguous, out-of-date, and even exaggerated and biased contents [19, 20]. According to previous findings, more net insider sales, lower information transparency, and less analyst coverage amplify the impact of media sentiment on future crash risk, which is consistent with the "information intermediation" channel.

Investors, owing to their limited attention and overconfidence, can be expected to overreact to catchy, anecdotal, and less relevant information, but underreact to abstract, statistically listed, and relevant information [21]. Furthermore, they may exhibit confirmation bias, which is the tendency to seek and believe information that supports one's beliefs while ignoring later signals that are inconsistent with their prior beliefs after developing a favorable impression of a firm [22]. As such, media coverage of firms (particularly positive news), regardless of whether the content is outdated or not, can easily pique the interest of investors, causing them to overreact or overestimate a firm's prospects and bring about a short- or long-term increase above the fundamental value [23–25]. However, when actual operational problems are revealed or a firm fails to meet expectations, negative sentiment will emerge and the stock price will reverse, resulting in a crash; this process will be reinforced when media coverage is biased and exaggerated. More retail investor positions, more active investor sentiment, and divergence between analysts' opinions and media coverage amplify the impact of news sentiment on future crash risks, which is consistent with the "investor sentiment" channel.

We expect to contribute to the literature in the following ways. First, our work also adds to the growing literature on the determinants of firm-specific price crash risk. Numerous studies have established a link between media sentiment and stock prices [7, 22, 26–28]. Sentiment in news articles contain novel information on stock prices [7, 8], but few studies have paid attention to the relation between the media sentiment's asymmetric effect on the future firm stock price crash risk. We attempt to fill this gap. We provide the formal piece of empirical evidence that positive news sentiment predicts a higher firm-specific future price crash risk, whereas negative sentiment increases the current crash risk. We provide a thorough examination of the impact of media sentiment from both perspectives of inspiration of investor sentiment of media tone and information economics, revealing new evidence that investors' irrational and excessive optimism could be a major cause of stock price bubbles and crashes in China, where retail investors predominate and short selling is restricted.

Second, our research study combines advanced deep learning and dictionary methods, which take full advantage of the performance and intelligence of computer technology and greatly improve the identification accuracy and efficiency of massive sentiment information. The sentiment dictionary approach uses word and syntactic analyses of text to calculate sentiment values as the basis for determining text sentiment tendencies. However, individuals can add necessary semantic words, such as praise words, degree adverbs, and negative words, which play an important role in enhancing or weakening sentiment semantic words [21]. Sentiment dictionary classification methods, which ignore the characteristics of language, such as grammar, context, and subjective construction methods, are likely to have the problem of omission. We attempt to mitigate this problem through integrating a cutting-edge deep learning method. To the best of our knowledge, this study is the first to combine

deep learning and dictionary methods to perform sentiment analysis in the financial field, thereby extending the application of sentiment analysis methods in the financial field.

Third, our study contributes to the literature on text analysis in the economic field [4, 29]. We introduce BERT, a deep learning pretraining model, and combine it with a relatively mature sentiment dictionary. Using the BERT pretraining model, researchers can take full advantage of the contextual information in the news. The vector expression of the same words is different between news and contexts, which was difficult to address in previous studies. In the pretraining model, by pretraining large text corpora as a language model, we create embeddings for the context associations (embedding) of each word in a sentence, which could then be entered into subsequent tasks, thereby enabling a full quantification of the information contained in the text.

2. Theory and Hypothesis

The media, as an important vehicle for information dissemination, play a significant role in the risk of stock price crashes. Sentiments Given the content of a media news, the sentiment contained in the content matters. Pure positive or negative news, which may mask the firms' actual situation, can exacerbate information asymmetry between firms and outside investors. In addition, the problems of irrational sentiment, herding effects, and "chasing the upside and killing the downside" phenomenon are aggravated by biased news. Both channels can increase the risk of a stock price crash. Considering the preceding ideas, we formulate the following competing hypotheses:

Hypothesis H1a: Higher average media sentiment can exacerbate the future stock price crash risk.

Hypothesis H1b: Higher average media sentiment can alleviate the future stock price crash risk.

We also examine the heterogeneous effect of sentiment in news on future stock price crash risks. When the media exaggerate the positive parts of news, they send positive signals on the firms to outside investors. With information asymmetry, investors will overestimate the firms' value, which can lead to abnormal stock price increases. As short selling is not allowed in China, rational and pessimistic investors are unable to engage in the market and stock prices will continue to increase until investors realize that there is an overvaluation component in the news. According to Solomon [28], the media's whitewashing behavior of overusing positive terms to disclose the information of listed firms can lead to a sharp decline in stock prices. Therefore, intense positive news can increase the stock price crash risk.

Alternatively, related research studies reveal the tendency of firm management to conceal negative information from the public [3, 30], which inevitably leads to an overestimate of the firm value, resulting in higher stock prices. Simultaneously, retail investors are more sensitive to negative news [31]. When the management has no choice but to release negative information to the market, retail investors will sell-off stock holdings, which increases the risk of a stock

price crash [32]. Therefore, the coverage of negative news can increase the risk of stock price crash. Considering the preceding ideas, we hypothesize as follows:

Hypothesis H2a: News with positive sentiment can exacerbate the future stock price crash risk.

Hypothesis H2b: News with negative sentiment can exacerbate the future stock price crash risk.

Subsequently, we then proceed with identifying the underlying mechanisms. We hypothesize that the impact may come from the two channels of "information intermediation" and "investor sentiment."

According to Jin and Myers [33], when there is information asymmetry, the agency problem, such as management's rent-seeking and concealment of negative news, can affect the share price crash risk significantly. Insiders have information that is not yet publicly available, which can be used to assess the value of the company and predict future firm performance [34]. Insider sell-off behavior is positively associated with the risk of a stock price crash [35]. The insider's choice to sell stocks sends negative signals to outside investors, thereby raising the probability of a future crash risk. Based on the explanation above, we hypothesize the following.

Hypothesis H3a: The impact of media sentiment on the future stock price crash risk is enhanced when insiders have more net sales of stock.

Information asymmetry can prevent investors from knowing the firm's actual operation and investors may be deceived by false public information. Especially, firms whose information transparency is low and those whose management is more likely to hide bad news are more likely to experience a sharp stock price fall in future [3]. When financial opacity is high, investors cannot fully grasp the true state of a firm through public information. They will rely more on the media coverage to make investment decisions, which will amplify the role of media sentiment. Therefore, we hypothesized as follows.

Hypothesis H3b: The impact of media sentiment on the future stock price crash risk is enhanced when the firm has higher financial opacity.

Analysts serve as both information intermediaries and management monitors [36]. They acquire and process data about firms using public information, field research, and other sources and reduce information asymmetry between firms and investors [37, 38]. According to He et al. [35], analyst coverage reduces stock price crash risk via analysts' role as information intermediaries and monitors. Nonetheless, when analysts cannot fully perform the information intermediary role, it is more difficult for investors to learn the real situation of the firm and they cannot accurately identify the noise in the media coverage; thus, the impact of media sentiment on future stock price risk is more evident through driving the investors' decision-making process. We thus hypothesize as follows:

Hypothesis H3c: The impact of media sentiment on the future stock price crash risk is enhanced when the firm has lower analyst coverage.

Next, we test the existence of the "investor sentiment" channel.

The sentiment of media affects investors differently. Institutional investors have more information and a specialized ability; therefore, it is easier for them to judge the validity of the information contained in the news. Sentiments in news thus have a limited influence on institutional investors. Conversely, retail investors do not have the information and skills possessed by institutional investors; thus, sentiments in news have a greater impact on retail investors. Therefore, we formulate the following hypothesis:

Hypothesis H4a: The impact of media sentiment on the future stock price crash risk is enhanced when the stock has a higher proportion of retail investors.

Investors may place great faith in catchy, anecdotal, and low-relevance information and overreact to it as a result of limited attention and overconfidence [21]. They may also exhibit a confirmation bias, which is the tendency to seek and believe information that supports one's beliefs while ignoring later signals that are inconsistent with prior beliefs after developing a favorable impression of the firm [22]. Thus, positive media coverage of firms attracts investors' attention, thereby causing them to overreact or form over-expectations about the firm's prospects, resulting in a stock price that is briefly or chronically above the underlying value. Therefore, we proposed the following hypothesis:

Hypothesis H4b: The impact of media sentiment on the future stock price crash risk is enhanced when investors are more optimistic.

Heterogeneous beliefs among investors increase when analysts' opinions differ from the sentiment of media coverage. In the absence of a short selling mechanism, more optimistic investors are expected in the market in the short term. However, over time, pessimistic beliefs will eventually emerge, which will increase the likelihood of a future stock price crash. Therefore, we hypothesize as follows:

Hypothesis H4c: The future stock price crash risk increases when analysts' points disagree with the sentiment of media.

3. Sentiment Extraction Model Design

Our study obtains financial news data for the period January 2017 to December 2020 from the China Stock Market and Accounting Research database. We preprocess the data by deleting special symbols and irrelevant information. We select 3,305 news items from the 1,132,856 news texts to label the news regarding the financial entities as positive or negative sentiments and added them to the financial sentiment dictionary manually. Using the BERT pretraining model and financial sentiment dictionary-based attention mechanism, we classify the sentiment of 1,132,856 news items. We then derive the average sentiment index of each stock in each year using the weighted average of all related news sentiment.

3.1. Data Preprocessing. We construct a microsentiment corpus in the financial field. To avoid interference with the hard data contained in the news, we eliminate the company announcements and then preprocess the text by removing

special symbols and using regular matching to remove irrelevant information.

3.1.1. Data Clean-Up. We label the financial entities: we identify the company, person, and brand names in the text. Entity names are marked based on the principle of long matching. We also identify the company and brand names with the help of "Tianyancha." For example, in the following text, "Runtu Shares: Ruan Jiachun (Chairman) plans to reduce no more than the total share capital of 1.28," "Runtu Shares" and "LeTV" are marked as entities.

3.1.2. Label News Sentiment. The sentiment polarity of financial entities is grouped into three categories—neutral, negative, and positive. Each category is defined as follows and Table 1 shows the distribution.

Positive sentiment: The text is marked as positive if the fact favors the operation of the company and there are some artificial positive comments. For example, "Southeast network frame won the bid of 357-million-yuan project."

Negative sentiments: If the information in the text is bad for the company's operation because it includes some facts that are bad for the operation of the company and artificial negative comments. For example, "Tianmaotui will be delisted from the Shenzhen Stock Exchange on July 20."

Neutral sentiment: Unlike positive and negative sentiments, the labeling of neutral sentiments is relatively complex. The text information is related to the operation of the company but cannot be judged as favorable or unfavorable, or it has both favorable and unfavorable facts. For example, "e-commerce is the direction of future development, all enterprises are making efforts. So does Huawei, but at present, the effectiveness needs to be tested."

To construct a microsentiment analysis dataset in the financial field, we select 4,516 samples from the 1,132,856 news texts obtained for annotation. After the independent annotation, all annotators discussed the additional annotations noted in the case of objectionable or uncertain results until consensus was reached. The annotation data were artificially modified and the annotation was completed.

Finally, 3,644 financial entities are sorted out. Each financial entity corresponds to one or more sentences. Each article has a total of 10,112 sentiment sentences. Based on prior financial knowledge, we construct a sentiment dictionary in the financial domain, which contains 2,842 positive words, 1,230 neutral words, and 2,043 negative words (Table 2).

3.2. BERT Pretraining Model. In 2018, Google proposed the natural language pretraining model, BERT, in the article "BERT: Pretraining of Deep Bidirectional Transformers for Language Understanding." Training of the BERT model mainly includes the following two steps.

TABLE 1: Financial entities.

	Positive	Neutral	Negative
Nums.	4,179	3,202	1,627
Pct	46.39%	35.55%	18.06%

The sentiment polarity of financial entities is grouped into three categories—neutral, negative, and positive.

TABLE 2: Financial dictionary.

	Positive	Neutral	Negative
Nums.	2,842	1,230	2,043
Pct	46.48%	20.11%	33.41%

The authors construct a sentiment dictionary in the financial domain, which contains 2,842 positive words, 1,230 neutral words, and 2,043 negative words.

BERT pretraining: The pretraining of BERT helps it learn the characteristics of a character, a word, statement levels, and understatement relationships among massive text data through simultaneous two pre-training tasks—masked language model and next sentence prediction. During pretraining, the same corpus is inputted into the model multiple times, but each input is preprocessed in different forms, allowing the same corpus to be fully utilized. For users, the pretrained models and parameters can be downloaded from the Internet and can be directly fine-tuned without having to do pretraining themselves, which reflects the convenience of BERT.

Fine-tuning: On the basis of the pretrained model, an output layer is customized and added to specific downstream tasks, such as text sentiment classification and sequence annotation. Then, the data from downstream tasks are used to fine-tune the model to generate models with higher prediction accuracy for various NLP tasks.

BERT uses a more powerful bidirectional transformer encoder (Figure 1) along with the masked language model and next sentence prediction (NSP) as an unsupervised goal, to enable the vector representation of each word and word output by the model to describe the overall information of the input text as comprehensively and accurately as possible. Thus, BERT provides better initial values of model parameters for subsequent fine-tuning tasks. Its input embedding is constructed by summing the token, segment, and position embedding of the corresponding word. It also contains more parameters, which gives it a stronger word vector embedding ability.

3.3. Construction of BERT-DCA Model. We construct a BERT-DCA model (Figure 2) that combines the financial sentiment dictionary and attention for sentiment analysis. Two information processing channels—left semantic information attention channel (SAC) and right sentiment information attention channel (EAC)—are adopted in the structure. The SAC extracted semantic information, whereas the EAC allowed the model to pay attention to the

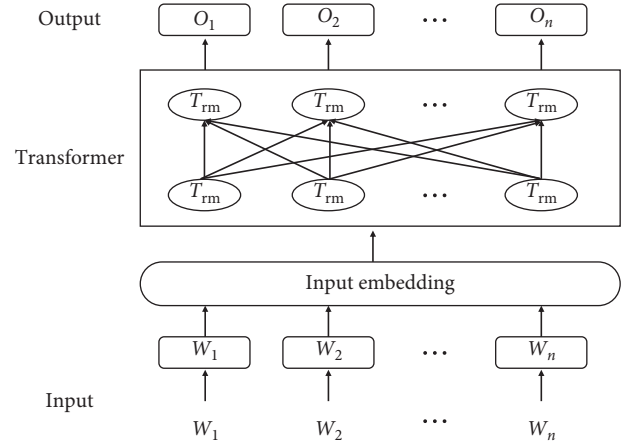


FIGURE 1: BERT model structure.

particularity of different types of words to supplement the weights and obtain more information as a supplement to word-level information.

3.3.1. Input Layer. For the text sentence sequence, after word partitioning, the word sequence $\{W_1, W_2, \dots, W_n\}$ is used as the input for SAC. Based on the domain sentiment dictionary and financial entities, words are classified into the following five categories: Pos, Neg, Neu, Entity, and Other, which denote positive, negative, neutral, financial entities, and others, respectively. They are from the sentimental dictionary discussed above. Then, we derive the sentimental information word collection $\{E_1, E_2, \dots, E_m\}$ as the input of EAC. We then use the pretraining model, BERT, to provide the word vector, which can achieve the dynamic adjustment of the word vector with the context, and train the real sentiment semantic embedding model to obtain the semantic information word vector matrix R_x and the sentiment information word vector matrix R_e .

$$R_x = x_1 \oplus x_2 \oplus \dots \oplus x_n, \quad (1)$$

$$R_e = e_1 \oplus e_2 \oplus \dots \oplus e_m, \quad (2)$$

where \oplus is the row vector connection operator and the dimensions of R_x and R_e denote the number of words in the news and of annotated financial sentiment entities, respectively.

3.3.2. Feature Extraction. For semantic information texts, we used the BiGRU neural network (model structure shown in Figure 3) to handle both forward and reverse text sequences. We extracted the deep text information and then used the financial dictionary to guide attention mechanisms to assign corresponding weights to the extracted feature information. For sentiment information sets, affective information words were encoded using a fully connected network combined with attention mechanisms to obtain affective signals.

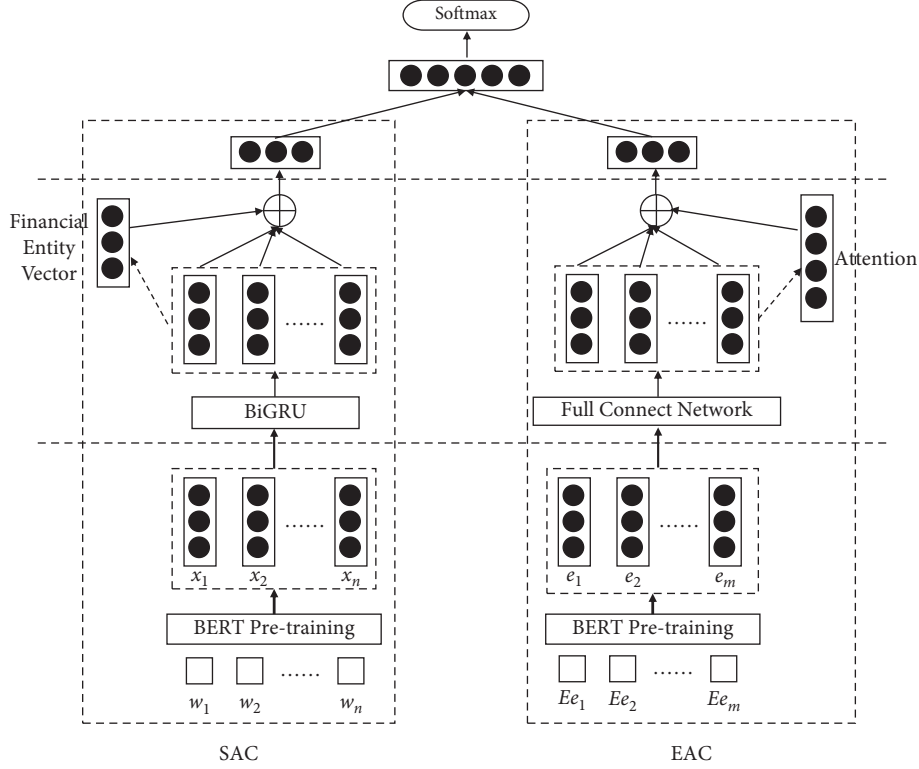


FIGURE 2: A model framework that combines sentiment dictionaries and attention.

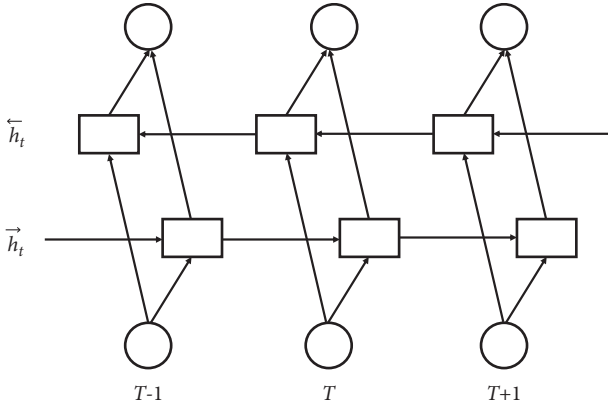


FIGURE 3: Structure of BiGRU.

The output of the BiGRU information extraction model at time t is composed of the output of the forward and reverse GRU and calculated as follows:

$$\begin{aligned}
 x_t &= W_e w_t, \quad t \in [1, T], \\
 \vec{h}_t &= \overrightarrow{GRU}(x_t), \quad t \in [1, T], \\
 \overleftarrow{h}_t &= \overleftarrow{GRU}(x_t), \quad t \in [1, T], \\
 s_t &= \left[\vec{h}_t, \overleftarrow{h}_t \right], \quad t \in [1, T].
 \end{aligned} \tag{3}$$

By combining \vec{h}_t and \overleftarrow{h}_t to obtain the semantic representation s_t , the forward and reverse semantic information elements are considered in the same position.

Next, we use the financial dictionary to guide attention mechanisms. To improve the accuracy of the sentiment analysis of financial text, we model the relation between sentiment and each word, assigning different weights to the semantic characteristics of the clause using the attention mechanism. In this way, more important words get more attention. Based on the financial entities and BiGRU layer output $H^s = \{h_1^s, h_2^s, \dots, h_t^s\}$, we obtained vectorized representations as word-level attention. The weight is calculated as follows:

$$\alpha_{st} = \frac{\exp(\gamma(h_c^s, e^E))}{\sum_c \exp(\gamma(h_c^s, e^E))}, \tag{4}$$

$$\gamma(h_c^s, e^E) = \tanh(h_c^s \cdot \omega_m^T \cdot e^{EF} + b_a).$$

The output after BiGRU processing is expressed as $[h_1^s, h_2^s, \dots, h_c^s]$, where ω_m^T is the weight matrix, b_a is the offset, e^E is the word vector of the financial entity, and α_{st} is the attention weight of the word w_{st} relative to the financial entity e^E . The text features with attention-weighted sentences are represented as follows:

$$o_s = \sum_t \alpha_{st} h_{st}, \tag{5}$$

where o_s is a semantic representation weighted by attention.

3.3.3. Feature Fusion Layer. The main task of the feature fusion layer is to combine the feature vector O^s generated in

SAC and feature vector O^e generated in EAC, to construct the overall sentimental feature vector. To simplify the calculation of the model, we perform feature fusion by row connection, constructing a matrix $O^* = (r_s + r_e) * c$ to generate the feature vector, where r_s and r_e are the number of rows of O^s and O^e , respectively, and c gives the column numbers for O^s and O^e .

3.3.4. Output Layer. We input the sentiment feature vector O^* generated by the feature fusion layer into the SoftMax classifier to obtain the final sentiment classification result predicted by the model as follows:

$$p = \text{softmax}(w_o O^* + b_o), \quad (6)$$

where w_o is the weight coefficient matrix, b_o is the bias matrix, and p is the predicted sentiment label.

3.3.5. Model Training. To use the constructed financial sentiment dictionary that could correspond to the input sentence, we need to construct a sentiment word vector of the same length as the term after the segmentation: Vec_{Att} , initialized as 0. After traversing the words in the input financial text, we set the corresponding position to 1 in the sentiment word vector if they appear in the financial sentiment dictionary. For example, assuming that the input financial text is “Langma cloud business development uncertainty,” we first initialize the sentiment word vector $[0, 0, 0, 0, 0]$. After the input sentence, the word “uncertainty” appears in the financial sentiment dictionary, and it is a negative word. Then, we set the word “uncertainty” in the corresponding position of the sentiment word vector to 1, after which the sentiment word vector of the sentence becomes $[0, 0, 0, 0, 1]$.

To employ the financial dictionary shown above as a guiding attention mechanism, we modify the loss function and add $\lambda(\alpha - Vec_{Att})^2$ after the cross-entropy loss. Here, λ is the hyperparameter that determines the importance of sentiment dictionary loss, α is the score of the attention mechanism, and Vec_{Att} is the sentiment dictionary vector. Thus, the attention mechanism score α can fit the financial sentiment word vector for the model to pay more attention to the input financial text—the financial sentiment words. The loss function is as follows:

$$L = - \sum_{i \in D} y_i \log p_i + \lambda(\alpha_{norm} - Vec_{Att}), \quad (7)$$

where D is the collection of samples, y_i is the true label, and p_i is the prediction result of the model. λ is the hyperparameter that determines the importance of affective dictionary loss and α_{norm} is the average attention score.

We thus use the predicted labels of 1,132,856 news items as the sentiment score in the empirical analysis.

4. Empirical Model

4.1. Sample Selection and Data Sources. Our sample covers all A-share listed firms from 2011 to 2020. Owing to the lag

phase of the study, the data time span is nine years (from 2012 to 2020).

We obtain the financial stock trading data from the WIND database. Among them, stock yield is given as weekly data, and the rest, as annual data. Following prior studies [38, 39], we process the original sample as follows: (1) We exclude financial and insurance listed firms; (2) we exclude listed firms with ST or * ST (ST: the company has suffered losses for two consecutive years and is specially treated, ST*: the company has suffered losses for three consecutive years and warned with delisting.); (3) we exclude listed firms with missing or abnormal data; and (4) we exclude listed firms with less than 15 weekly yield data.

We obtain data regarding Internet media news from the GuoTai'an (CASMAR) database. We perform positive and negative analyses of each report using sentiment analysis technology and assign sentiment scores. We then calculate the sum of the number of relevant news reports during the research period, average level of sentiment scores, and sentiment scores weighted by the number of news reports.

4.2. Econometric Model and Variables. To study the relation between diversification and future stock price crash risk, we construct a multiple regression model as follows:

$$\begin{aligned} \text{CrashRisk}_{i,t} = & \beta_0 \text{Cons} + \beta_1 \text{NewsSentiment}_{i,t-1} + \beta_2 \text{Size}_{i,t-1} \\ & + \beta_3 \text{Level}_{i,t-1} + \beta_4 \text{ROA}_{i,t-1} + \beta_5 \text{Ret}_{i,t-1} \\ & + \beta_6 \text{Sigma}_{i,t-1} + \sum \text{Year} + \sum \text{Firm} + \varepsilon_{i,t}. \end{aligned} \quad (8)$$

In the model, the explained variable CrashRisk represents the risk of a crash in individual stocks. NewsSentiment is the core explanatory variable, indicating the calculated sentimental score for the individual stock news report. Size denotes the size of the enterprise, Level represents the financial leverage of the company, ROA is the return on equity of the company, Ret is the average of the enterprise-specific weekly rate, and Sigma is the standard deviation of the enterprise-specific weekly rate. Year and Firm are time and firm fixed effects, respectively. Here, we focus on the coefficient β_1 . If β_1 is significantly positive, then there is a positive relation between the news reporting sentiment and risk of a stock price crash. Conversely, if β_1 is significantly negative, then there is a negative relation between the news reporting sentiment and risk of a stock price crash. The variables are introduced in Table 3 and elaborated as follows:

4.2.1. Explained Variables. Based on the methods of Jin and Myers [33] and Xu et al. [40], our study employs three approaches to measure the risk of a stock price crash. The specific algorithm is as follows:

The unexplained weekly yield of individual stocks in the market is calculated using the following model:

$$R_{i,t} = R_{m,t-2} + R_{m,t-1} + R_{m,t} + R_{m,t+1} + R_{m,t+2} + \varepsilon_{i,t}, \quad (9)$$

where $R_{i,t}$ represents the weekly yield of stock i in week t . $R_{m,t}$ is the weighted average of the weekly yield of week t . $\varepsilon_{i,t}$

TABLE 3: Variables and definitions.

	Indicator	Definition
Dependent variables	CRASH	Used to measure the risk of a crash: 1, 0 indicator
	NCSKEW	Used to measure the risk of a crash: negative return bias coefficient
	DUVOL	Used to measure the risk of a crash: the earnings fluctuation ratio
Independent variables	newSenti newPos newNeg	News coverage' sentiment
	anaSenti	Analysts' sentiment
Control variables	Size	The natural logarithm of the firm's market value
	Level	Firm's leverage, equal to the ratio of the firm's total liabilities to total assets
	ROA	The ratio of net profit to total assets
	Ret	The average value of the firm's annual weekly rate of return
	Sigma	The standard deviation of the firm's annual weekly rate of return
	Year	Yearly dummy variable
	Firm	Firm dummy variable

is the residual in equation (2), which represents the weekly return of stocks not explained in the market. Because $\epsilon_{i,t}$ is highly biased, we use $W_{i,t} = \ln(1 + \epsilon_{i,t})$ to represent stock-specific weekly yields. Based on $W_{i,t}$, we measure the risk of a stock price crash using three indicators—(CRASH), a negative return bias coefficient (NCSKEW), and the earnings fluctuation ratio (DUVOL).

CRASH is calculated as follows:

$$\text{CRASH}_{i,t} = 1 \left[\exists t, W_{i,t} \leq \text{Average}(W_{i,t}) - 3.09\sigma_{i,t} \right]. \quad (10)$$

$\text{CRASH}_{i,t}$ equals 1 if a firm experiences one or more firm-specific weekly returns $W_{i,t}$ falling 3.09 standard deviations below the mean firm-specific weekly return, and 0 otherwise.

NCSKEW is calculated as follows:

$$\text{NCSKEW}_{i,t} = \frac{-[n(n-1)^{3/2} \sum W_{i,t}^3]}{[(n-1)(n-2)(\sum W_{i,t}^2)^{3/2}]}, \quad (11)$$

where n represents the number of stock i in year t . The coefficient of the negative return bias is a positive measure of the risk of a stock price crash. Thus, the greater the coefficient is, the higher the possibility of a stock price crash.

DUVOL is calculated as follows:

$$\text{DUVOL}_{i,t} = \log \left\{ \frac{[(n_u - 1) \sum_{\text{DOWN}} W_{i,t}^2]}{[(n_d - 1) \sum_{\text{UP}} W_{i,t}^2]} \right\}. \quad (12)$$

The core explanatory variable in our study is a quantitatively weighted news reporting sentiment propensity, which is calculated as follows:

4.2.2. Explanatory Variables. The core explanatory variable in this study is a quantitatively weighted news reporting sentimental propensity, which is calculated as follows:

$$\text{newSenti}_{i,T} = \frac{\text{NewsCount}_{i,t} * \text{SentimentScore}_{i,t}}{\sum \text{NewsCount}_{i,t}}, \quad (13)$$

where $\text{newSenti}_{i,T}$ represents the media report sentiment tone of stock in year T . $\text{NewsCount}_{i,t}$ represents the number of news items regarding stock i in year T . $\text{SentimentScore}_{i,t}$ represents the average media reporting sentiment scores of

trading day t of stock i in year T , each calculated by our BERT-DCA model. Regarding the number of news, the higher it is, the more likely the investors will read the news. Thus, the probability that the reported sentiments are transmitted to investors is also higher. To examine how different types of Internet news sentiment work, we construct both positive and negative news coverage sentiment indicators.

$$\begin{aligned} \text{newPos}_{i,T} &= \frac{\text{PosNewsCount}_{i,t} * \text{PosSentimentScore}_{i,t}}{\sum \text{PosNewsCount}_{i,t}}, \\ \text{newNeg}_{i,T} &= \frac{\text{NegNewsCount}_{i,t} * \text{NegSentimentScore}_{i,t}}{\sum \text{NegNewsCount}_{i,t}}. \end{aligned} \quad (14)$$

To study the impact of market differences on the risk of a stock price crash, we use analysts' rating data to calculate their sentiments. First, we grade analysts on five points: +2, +1, 0, -1, and -2, indicating buy, overweight, neutral, reduction, and sell, respectively. We then calculate the total score of the year, divide it by the rating number, and finally standardize it using $\sim N(0,1)$, given as *Ana_senti*.

4.2.3. Control Variables. Following Jin and Myers [33], the control variables are defined as follows.

Enterprise size (Size) is expressed as the natural logarithm of the enterprise market value.

Operating leverage (Level) is the enterprise asset-liability ratio.

Compensation rate of corporate total assets (ROA) is an indicator used to measure corporate profitability.

Previous value of negative return bias coefficient/fluctuation ratio is used to control the impact of the lag phase of the risk of a stock price crash.

Stock-specific weekly earnings annual average (Ret) reflects the average level of stock yield.

Weekly earnings volatility (Sigma) reflects the volatility levels of stock-specific weekly earnings.

TABLE 4: Baseline regression.

	(1) CRASH	(2) CRASH	(3) NCSKEW	(4) NCSKEW	(5) DUVOL	(6) DUVOL
L.newSenti	0.087*** (0.031)	0.064** (0.032)	0.395*** (0.092)	0.190** (0.094)	0.203*** (0.062)	0.103* (0.060)
L.CRASH		-0.147*** (0.009)				
L.NCSKEW				-0.096*** (0.009)		
L.DUVOL						-0.106*** (0.008)
L.ret		0.957** (0.410)		8.770*** (0.842)		6.232*** (0.584)
L.sigma		0.594*** (0.195)		1.039** (0.441)		1.192*** (0.299)
L.roa		0.360 (0.224)		1.339*** (0.459)		0.860*** (0.307)
L.level		-0.012 (0.033)		-0.130** (0.059)		-0.113*** (0.039)
L.size		0.000 (0.000)		0.000* (0.000)		0.000 (0.000)
<i>Fixed effects:</i>						
Year dummy	Yes	Yes	Yes	Yes	Yes	Yes
Firm dummy	Yes	Yes	Yes	Yes	Yes	Yes
N	18,775	17,267	18,775	17,267	18,775	17,267
R ²	0.000	0.035	0.001	0.060	0.001	0.065

This table reports baseline regression estimates of stock crash risk on the quantitatively weighted news coverage sentiment. The sample of CRASH is used for columns (1)-(2), the sample of NCSKEW is used for columns (3)-(4), and that of DUVOL is used columns (5)-(6). Columns (1), (3), and (5) only include the news coverage sentiment; columns (2), (4), and (6) control for lagged crash risk and firm characteristic; year and firm dummy variable are included; all control variables are included as lagged one year. Robust standard errors are reported in parentheses. The labels ***, **, and * indicate 1%, 5%, and 10% levels of significance, respectively.

5. Empirical Analysis

5.1. News Sentiment and the Risk of a Stock Price Crash. Table 4 shows the analysis results for the relation between the media sentiment and future stock price crash risk. Columns 1, 3, and 5 in Table 4 are the effects of the sentiment weighted by the number of news items (*newSenti*) on the risk of a future stock price crash where control variables are not included. We found a significant positive effect of CRASH (0.087), NCSKEW (0.395), and DUVOL (0.203). Columns 2, 4, and 6 show the effect of sentiment (*newSenti*) on the future stock price crash risk after adding all control variables. Similarly, we find a significant positive effect on CRASH (0.064), NCSKEW (0.190), and DUVOL (0.103). Although the coefficient decreases after including control variables, they are still economically and statistically significant. Thus, the more positive the average media sentiment is, the higher the future stock price crash risk, or the more negative the current average news sentiment is, the lower the risk of a future stock price crash. These findings support *H1a* but not *H1b*.

Regarding the control variables, we observe a negative effect of firm financial leverage on the risk of a future stock price crash, implying that the latter risk is higher in firms with lower financial leverage—the smaller the size of the firm, the higher the risk. In addition, we find that a firm's

stock return (ROA) is significantly and negatively related to the risk of a future stock price crash, implying that the better the performance of a firm, the less likely it is to have a stock price crash in future. The effects of the other control variables on future stock price crash risk are not robust.

According to the results of the baseline regression, media sentiment is positively related to the future stock price crash risk. We replace average media sentiment in the baseline regression with media coverage positive and negative sentiment indicators to examine hypothesis *H2*. The results are shown in Table 5.

Columns 1, 2, and 3 in Panel A indicate that the effect of positive sentiment is significant at the 1% level, whereas Columns 4, 5, and 6 indicate that the effect of negative sentiment is insignificant. The regression results suggest that positive media sentiment plays a dominant role in China; the more positive the media sentiment, the higher the future stock price crash risk. Under information asymmetry, uninformed investors receive more positive information and irrational investors develop an overvaluation of stock prices [2, 14]. The negative sentiment appears to curb future crash risk, but the effect is insignificant.

Indeed, investors react more strongly to negative news [31]. Negative news causes retail investors to sell and increases the risk of a future stock price crash [32]. However, our findings indicate that negative news in the previous one

TABLE 5: Positive sentiment and negative sentiments of news coverage.

<i>Panel A: one lag period</i>						
	(1) CRASH	(2) NCSKEW	(3) DUVOL	(4) CRASH	(5) NCSKEW	(6) DUVOL
L.newPos	0.135*** (0.045)	0.516*** (0.187)	0.323*** (0.125)			
L.newNeg				0.033 (0.073)	0.087 (0.162)	0.015 (0.110)
L.CRASH	-0.152*** (0.009)			-0.151*** (0.009)		
L.NCSKEW		-0.102*** (0.009)			-0.107*** (0.009)	
L.DUVOL			-0.111*** (0.009)			-0.116*** (0.009)
L.ret	1.060** (0.426)	9.014*** (0.869)	6.415*** (0.606)	1.016** (0.431)	9.305*** (0.888)	6.560*** (0.614)
L.sigma	0.618*** (0.204)	1.204*** (0.455)	1.261*** (0.310)	0.545*** (0.203)	1.059** (0.459)	1.183*** (0.312)
L.roa	0.371 (0.236)	1.388*** (0.489)	0.877*** (0.321)	0.425* (0.235)	1.427*** (0.485)	0.812** (0.326)
L.level	-0.016 (0.035)	-0.111* (0.059)	-0.108*** (0.040)	-0.004 (0.037)	-0.147** (0.063)	-0.142*** (0.042)
L.size	0.000 (0.000)	0.000* (0.000)	0.000 (0.000)	0.000 (0.000)	0.000** (0.000)	0.000* (0.000)
<i>Fix effects:</i>						
Year dummy	Yes	Yes	Yes	Yes	Yes	Yes
Firm dummy	Yes	Yes	Yes	Yes	Yes	Yes
N	16,267	16,267	16,267	15,966	15,966	15,966
R ²	0.036	0.060	0.067	0.036	0.062	0.068
<i>Panel B: same period</i>						
newPos	-0.068 (0.088)	-0.023 (0.187)	-0.010 (0.129)			
newNeg				-0.185** (0.077)	-0.427*** (0.159)	-0.303*** (0.111)
L.CRASH	-0.161*** (0.009)			-0.161*** (0.009)		
L.NCSKEW		-0.109*** (0.009)			-0.110*** (0.009)	
L.DUVOL			-0.113*** (0.009)			-0.116*** (0.009)
L.ret	0.662 (0.417)	5.425*** (0.888)	3.671*** (0.630)	0.570 (0.425)	5.596*** (0.908)	3.864*** (0.640)
L.sigma	0.566*** (0.202)	0.138 (0.466)	0.219 (0.322)	0.556*** (0.201)	0.105 (0.475)	0.311 (0.326)
L.roa	0.647*** (0.221)	1.765*** (0.461)	1.094*** (0.316)	0.489** (0.214)	1.272*** (0.467)	0.822** (0.321)
L.level	-0.048 (0.036)	-0.112* (0.066)	-0.108** (0.044)	-0.049 (0.036)	-0.116* (0.065)	-0.115*** (0.044)
L.size	-0.000 (0.000)	0.000 (0.000)	-0.000 (0.000)	-0.000 (0.000)	0.000 (0.000)	0.000 (0.000)
<i>Fix effects:</i>						
Year dummy	Yes	Yes	Yes	Yes	Yes	Yes
Firm dummy	Yes	Yes	Yes	Yes	Yes	Yes
N	16,267	16,267	16,267	15,966	15,966	15,966
R ²	0.041	0.066	0.072	0.037	0.064	0.070

This table reports regression estimates of stock crash risk on the news coverage positive and negative sentiment indicators. Panel A is one year lag; panel B is the same year; year and firm dummy variable are included; all control variables are included as lagged one year. Robust standard errors are reported in parentheses. The labels ***, **, and * indicate 1%, 5%, and 10% levels of significance, respectively.

TABLE 6: Quarterly crash risk for different future periods.

	(1) CRASH	(2) CRASH	(3) CRASH	(4) CRASH	(5) CRASH	(6) CRASH	(7) CRASH	(8) CRASH
newPos	-0.035 (0.039)							
L.newPos		0.116*** (0.039)						
L2.newPos			0.082** (0.039)					
L3.newPos				0.047 (0.039)				
newNeg					-0.164*** (0.036)			
L.newNeg						0.046 (0.036)		
L2.newNeg							0.042 (0.036)	
L3.newNeg								-0.048 (0.036)
L.CRASH	-0.052*** (0.005)	-0.053*** (0.005)	-0.055*** (0.005)	-0.052*** (0.005)	-0.052*** (0.005)	-0.053*** (0.005)	-0.052*** (0.005)	-0.053*** (0.005)
L.ret	1.523*** (0.479)	1.824*** (0.484)	2.339*** (0.540)	2.999*** (0.565)	1.709*** (0.476)	1.917*** (0.478)	2.418*** (0.533)	2.833*** (0.540)
L.sigma	1.734*** (0.253)	1.633*** (0.260)	1.807*** (0.282)	2.052*** (0.283)	1.698*** (0.260)	1.617*** (0.256)	2.066*** (0.276)	1.916*** (0.280)
L.roa	-0.070 (0.052)	-0.132*** (0.051)	-0.194*** (0.053)	-0.141*** (0.052)	-0.095* (0.052)	-0.143*** (0.051)	-0.134*** (0.052)	-0.155*** (0.053)
L.level	-0.038** (0.019)	-0.037** (0.018)	-0.038* (0.020)	-0.034* (0.019)	-0.029 (0.020)	-0.030 (0.020)	-0.038* (0.020)	-0.009 (0.020)
L.size	0.000 (0.000)	-0.000 (0.000)	0.000 (0.000)	0.000 (0.000)	0.000 (0.000)	0.000 (0.000)	0.000 (0.000)	0.000 (0.000)
<i>Fixed effects</i>								
Quarter dummy	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Firm dummy	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
N	51,110	51,110	51,110	51,110	51,919	51,919	51,919	51,919
R ²	0.040	0.039	0.039	0.038	0.037	0.038	0.039	0.038

This table reports regression estimates of quarterly stock crash risk on the news coverage positive and negative sentiment indicators. Columns (1) to (4) are positive indicators. Columns (5) to (8) are negative indicators; quarter and firm dummy variable are included; all control variables are included as lagged one year. Robust standard errors are reported in parentheses. The labels ***, **, and * indicate 1%, 5%, and 10% levels of significance, respectively.

TABLE 7: Two stages OLS regression.

	First stage (1) Newsenti	(2) CRASH	Second stage (3) NCSKEW	(4) DUVOL
L.newSenti		0.074*** (0.028)	0.510** (0.253)	0.313** (0.157)
L.newSentiInd	0.873*** (0.022)			
L.newSentiPro	0.747*** (0.041)			
Controls	Yes	Yes	Yes	Yes
Cragg-Donald Wald F	185.966***			
Sargan chi (p)		0.257 (0.612)	0.571 (0.450)	1.466 (0.226)
N	16,845	16,845	16,845	16,845
R ²	0.212	0.035	0.059	0.064

This table reports regression estimates of stock crash risk on the news coverage sentiment indicators two stage OLS. Columns (1) is first stage. Columns (2) to (4) are second stage; year and firm dummy variable are included; all control variables are included as lagged one year. Robust standard errors are reported in parentheses. The labels ***, **, and * indicate 1%, 5%, and 10% levels of significance, respectively.

TABLE 8: Regression results for corporate insider trading.

	Low			High		
	(1)	(2)	(3)	(4)	(5)	(6)
	CRASH	NCSKEW	DUVOL	CRASH	NCSKEW	DUVOL
L.newSenti	0.030 (0.168)	-0.394 (0.333)	-0.154 (0.224)	0.103*** (0.037)	0.336*** (0.114)	0.119** (0.059)
L.CRASH	-0.144*** (0.010)			-0.185*** (0.033)		
L.NCSKEW		-0.088*** (0.034)			-0.096*** (0.009)	
L.DUVOL			-0.072** (0.032)			-0.109*** (0.009)
L.ret	1.570 (1.382)	10.925*** (2.656)	7.660*** (1.897)	0.933** (0.454)	8.500*** (0.954)	6.077*** (0.661)
L.sigma	1.427* (0.744)	3.266** (1.595)	2.124* (1.135)	0.541** (0.214)	0.707 (0.485)	1.049*** (0.326)
L.roa	-0.717 (0.880)	0.484 (1.754)	0.712 (1.204)	0.453* (0.235)	1.598*** (0.498)	1.016*** (0.336)
L.level	-0.036 (0.112)	0.014 (0.214)	-0.035 (0.150)	0.020 (0.036)	-0.120* (0.064)	-0.104** (0.041)
L.size	-0.000 (0.000)	-0.000 (0.000)	-0.000* (0.000)	0.000 (0.000)	0.000** (0.000)	0.000*** (0.000)
<i>Fixed effects:</i>		Yes	Yes	Yes	Yes	Yes
Year dummy	Yes	Yes	Yes	Yes	Yes	Yes
Firm dummy						
N	5,488	5,488	5,488	5,884	5,884	5,884
R ²	0.035	0.064	0.069	0.043	0.036	0.055

This table reports panel estimates of stock crash risk on the news coverage sentiment by using the subsamples of high insider trading groups and low insider trading groups. The total sample is divided into two subsamples: lower and higher groups, based on 30% and 70% quartiles of insiders' net stock sales, respectively. Columns (1) to (3) are low subgroup, the dependent variables are CRASH, NCSKEW, and DUVOL, respectively. Columns (4) to (6) are high subgroup, the dependent variables are CRASH, NCSKEW, and DUVOL, respectively. All control variables are included as lagged one year. Robust standard errors are reported in parentheses. The labels ***, **, and * indicate 1%, 5%, and 10% levels of significance, respectively.

TABLE 9: Regression results for information disclosure quality.

	Low			High		
	(1)	(2)	(3)	(4)	(5)	(6)
	CRASH	NCSKEW	DUVOL	CRASH	NCSKEW	DUVOL
L.newSenti	0.051 (0.106)	0.026 (0.222)	-0.093 (0.144)	0.114*** (0.043)	0.384** (0.179)	0.196** (0.098)
L.CRASH	-0.114*** (0.019)			-0.162*** (0.019)		
L.NCSKEW		-0.098*** (0.021)			-0.111*** (0.015)	
L.DUVOL			-0.101*** (0.020)			-0.124*** (0.016)
L.ret	2.489*** (0.939)	11.948*** (2.231)	8.111*** (1.424)	0.494 (0.853)	5.299*** (1.451)	4.534*** (1.117)
L.sigma	1.725*** (0.411)	1.739* (1.045)	1.502** (0.662)	0.239 (0.437)	1.582** (0.744)	1.326** (0.579)
L.roa	-0.430 (0.544)	-0.536 (1.349)	0.197 (0.864)	0.768* (0.405)	1.134* (0.665)	0.586 (0.487)
L.level	-0.141*** (0.042)	0.037 (0.112)	0.040 (0.069)	-0.025 (0.059)	-0.218** (0.102)	-0.205*** (0.076)
L.size	0.000 (0.000)	0.000** (0.000)	0.000 (0.000)	0.000 (0.000)	0.000 (0.000)	0.000 (0.000)
<i>Fixed effects:</i>		Yes	Yes	Yes	Yes	Yes
Year dummy	Yes	Yes	Yes	Yes	Yes	Yes
Firm dummy						
N	5,488	5,488	5,488	5,884	5,884	5,884
R ²	0.038	0.064	0.070	0.044	0.080	0.079

This table reports panel estimates of stock crash risk on the news coverage sentiment by using the subsamples of high information disclosure quality groups and low information disclosure quality groups. The total sample is divided into two subsamples: lower and higher groups, based on 30% and 70% quartiles of KV index, respectively. Columns (1) to (3) are low subgroup, the dependent variables are CRASH, NCSKEW, and DUVOL, respectively. Columns (4) to (6) are high subgroup, the dependent variables are CRASH, NCSKEW, and DUVOL, respectively. All control variables are included as lagged one year. Robust standard errors are reported in parentheses. The labels ***, **, and * indicate 1%, 5%, and 10% levels of significance, respectively.

TABLE 10: Regression results for analysts' coverage.

	Low			High		
	(1)	(2)	(3)	(4)	(5)	(6)
	CRASH	NCSKEW	DUVOL	CRASH	NCSKEW	DUVOL
L.newSenti	0.129** (0.058)	0.543** (0.221)	0.413** (0.167)	0.029 (0.078)	0.132 (0.220)	0.147 (0.146)
L.CRASH	-0.200*** (0.029)			-0.179*** (0.012)		
L.NCSKEW		-0.163*** (0.018)			-0.102*** (0.022)	
L.DUVOL			-0.156*** (0.019)			-0.121*** (0.021)
L.ret	2.276 (1.427)	6.842*** (1.432)	6.149*** (1.104)	0.841 (0.531)	3.826* (2.198)	2.842* (1.504)
L.sigma	1.788*** (0.600)	0.335 (0.839)	0.675 (0.641)	0.417 (0.274)	0.063 (1.060)	0.920 (0.721)
L.roa	0.159 (0.783)	2.087*** (0.760)	0.600 (0.538)	0.480* (0.284)	0.598 (1.204)	0.891 (0.798)
L.level	0.029 (0.093)	-0.158 (0.134)	-0.131 (0.103)	-0.030 (0.041)	-0.051 (0.088)	-0.031 (0.060)
L.size	0.000* (0.000)	0.000 (0.000)	0.000** (0.000)	0.000 (0.000)	0.000 (0.000)	-0.000 (0.000)
<i>Fixed effects:</i>						
Year dummy	Yes	Yes	Yes	Yes	Yes	Yes
Firm dummy	Yes	Yes	Yes	Yes	Yes	Yes
N	4618	4618	4618	4569	4569	4569
R ²	0.066	0.073	0.069	0.078	0.081	0.086

This table reports panel estimates of stock crash risk on the news coverage sentiment by using the subsamples of low analysts' attention groups and high analysts' attention groups. We used the total number of analysts coverage of firms to measure analysts' attention. We used the 30% and 70% quartiles of analysts' attention as the cut-off, and the firms were divided into low and high attention groups. Columns (1) to (3) are low subgroup, the dependent variables are CRASH, NCSKEW and DUVOL, respectively; column (4) to (6) are high subgroup, the dependent variables are CRASH, NCSKEW and DUVOL, respectively. All control variables are included as lagged one year. Robust standard errors are reported in parentheses. The labels ***, **, and * indicate 1%, 5%, and 10% levels of significance, respectively.

TABLE 11: Regression results for different institutional holding groups.

	Low			High		
	(1)	(2)	(3)	(4)	(5)	(6)
	CRASH	NCSKEW	DUVOL	CRASH	NCSKEW	DUVOL
L.newSenti	0.092** (0.044)	0.481*** (0.176)	0.358*** (0.124)	0.018 (0.094)	0.037 (0.220)	-0.027 (0.144)
L.CRASH	-0.172*** (0.016)			-0.204*** (0.022)		
L.NCSKEW		-0.131*** (0.015)			-0.161*** (0.020)	
L.DUVOL			-0.170*** (0.014)			-0.182*** (0.019)
L.ret	0.933 (0.720)	6.633*** (1.339)	5.416*** (0.959)	0.316 (1.001)	2.351 (2.192)	2.020 (1.487)
L.sigma	0.697** (0.343)	0.236 (0.726)	0.744 (0.527)	0.632 (0.447)	0.272 (1.023)	0.755 (0.678)
L.roa	0.183 (0.376)	0.638 (0.706)	0.310 (0.537)	0.423 (0.538)	0.861 (1.256)	0.850 (0.784)
L.level	-0.062 (0.049)	-0.072 (0.070)	-0.087* (0.049)	-0.100 (0.070)	-0.346*** (0.171)	-0.181* (0.109)
L.size	0.000 (0.000)	0.000*** (0.000)	0.000*** (0.000)	-0.000 (0.000)	0.000 (0.000)	0.000 (0.000)
<i>Fixed effects:</i>						
Year dummy	Yes	Yes	Yes	Yes	Yes	Yes
Firm dummy	Yes	Yes	Yes	Yes	Yes	Yes
N	4,734	4,734	4,734	6,513	6,513	6,513
R ²	0.046	0.071	0.087	0.057	0.080	0.092

This table reports panel estimates of stock crash risk on the news coverage sentiment by using the subsamples of low institutional shareholding groups and high institutional shareholding groups. We used the 30% and 70% quartiles of the institutional shareholding as the cut-off point, and the firms were divided into low and high shareholding groups. Columns (1) to (3) are low subgroup, the dependent variables are CRASH, NCSKEW, and DUVOL, respectively. Columns (4) to (6) are high subgroup, the dependent variables are CRASH, NCSKEW, and DUVOL, respectively. All control variables are included as lagged one year. Robust standard errors are reported in parentheses. The labels ***, **, and * indicate 1%, 5%, and 10% levels of significance, respectively.

TABLE 12: Regression results for different investor sentiment.

	Low			High		
	(1)	(2)	(3)	(4)	(5)	(6)
	CRASH	NCSKEW	DUVOL	CRASH	NCSKEW	DUVOL
L.newSenti	-0.011 (0.109)	-0.063 (0.231)	-0.097 (0.151)	0.176** (0.073)	0.441** (0.196)	0.254** (0.128)
L.CRASH	-0.178*** (0.019)			-0.173*** (0.016)		
L.NCSKEW		-0.144*** (0.019)			-0.104*** (0.018)	
L.DUVOL			-0.133*** (0.019)			-0.126*** (0.017)
L.ret	2.540*** (0.677)	7.085*** (1.893)	4.014*** (1.318)	2.426*** (0.759)	6.721*** (1.778)	4.536*** (1.205)
L.sigma	0.665 (0.437)	1.373 (0.929)	0.977 (0.606)	0.789** (0.371)	0.826 (0.925)	0.019 (0.634)
L.roa	0.755 (0.605)	1.529 (0.981)	1.205* (0.617)	0.554 (0.409)	2.048** (0.995)	1.666** (0.699)
L.level	0.036 (0.067)	-0.196* (0.106)	-0.102 (0.066)	-0.003 (0.055)	-0.089 (0.128)	-0.087 (0.088)
L.size	0.000 (0.000)	0.000* (0.000)	0.000 (0.000)	0.000 (0.000)	0.000 (0.000)	0.000 (0.000)
<i>Fixed effects:</i>	Yes	Yes	Yes	Yes	Yes	Yes
Year dummy	Yes	Yes	Yes	Yes	Yes	Yes
Firm dummy						
N	4,881	4,881	4,881	5,863	5,863	5,863
R ²	0.047	0.078	0.085	0.050	0.071	0.083

This table reports panel estimates of stock crash risk on the news coverage sentiment by using the subsamples of pessimistic investor groups and optimistic investor groups. We used the 30% and 70% quartiles of the investor sentiment as the cut-off point, and the firms were divided into low and high shareholding groups. Columns (1) to (3) are low subgroup, the dependent variables are CRASH, NCSKEW, and DUVOL, respectively. Columns (4) to (6) are high subgroup, the dependent variables are CRASH, NCSKEW, and DUVOL, respectively. All control variables are included as lagged one year. Robust standard errors are reported in parentheses. The labels ***, **, and * indicate 1%, 5%, and 10% levels of significance, respectively.

TABLE 13: Regression of adding analyst and media coverage sentiment crossterms.

	(1)	(2)	(3)
	CRASH	NCSKEW	DUVOL
L.newSenti	0.186* (0.109)	0.360** (0.171)	0.266 (0.158)
L.anaSenti	0.039*** (0.011)	0.057** (0.024)	0.049*** (0.016)
L.newSenti#L.anaSenti	-0.247* (0.145)	-0.172* (0.100)	-0.156* (0.092)
L.CRASH	-0.151*** (0.009)		
L.NCSKEW		-0.102*** (0.009)	
L.DUVOL			-0.111*** (0.009)
L.ret	0.942** (0.426)	8.677*** (0.867)	6.188*** (0.603)
L.sigma	0.524** (0.206)	0.874* (0.460)	1.120*** (0.315)
L.roa	0.285 (0.231)	1.286*** (0.484)	0.813** (0.325)
L.level	0.006 (0.036)	-0.141** (0.062)	-0.109*** (0.041)
L.size	0.000 (0.000)	0.000* (0.000)	0.000 (0.000)
<i>Fixed effects:</i>			
Year dummy	Yes	Yes	Yes
Firm dummy	Yes	Yes	Yes
N	15,753	15,753	15,753
R ²	0.037	0.064	0.069

This table reports panel estimates of stock crash risk on the interaction between the news coverage sentiment and the analysis report sentiment. The dependent variables in Columns (1), (2), and (3) are CRASH, NCSKEW, and DUVOL, respectively. All control variables are included as lagged one year. Robust standard errors are reported in parentheses. The labels ***, **, and * indicate 1%, 5%, and 10% levels of significance, respectively.

year does not increase the risk of a future stock price crash. A reason may be that investors are overly sensitive to negative news, particularly in markets with a high proportion of irrational investors, such as China. When negative news breaks, investors panic and quickly sell their stocks, resulting in a significant drop in the stock price below its fundamental value, which means the occurrence of the price crash in the current term rather than the future. Thus, negative news coverage can be hypothesized to increase the current stock crash risk, whereas positive media coverage decreases such a risk. The results are represented in Table 5, Panel B. It shows that the effects of negative news in the current period is significant, confirming our analysis that the more negative the news, the greater the risk of a crash in the current period. Meanwhile, the effects of positive media coverage are insignificant, indicating that positive news suppresses the risk of a crash in the current period. However, because investors are less sensitive to positive than negative news, their impact is not significant.

To investigate the robustness of the impact of positive and negative news, we considered quarterly level crash risk regressions that included current, prior one period to prior three period indicators for positive and negative news, respectively. The results are shown in Table 6. Overall, the quarterly regression results are consistent with the annual regression results. Positive news reduces the risk of a stock price crash in the current period (-0.035); however, the effect is not statistically significant. Columns 2–4 show that positive news significantly increases the risk of a stock price crash in the future period, consistent with the annual regression. Negative media coverage increases the current period crash risk (-0.164) and had a significant effect. Columns 6–8 reveal that negative media coverage reduces the future stock price crash risk, but the suppression effect is insignificant except for the previous period. These results confirm *H2a* but not *H2b*.

5.2. Endogeneity. The relation between media coverage and stock returns is endogenous. The reverse causality, as media coverage is more likely to focus on stocks with higher returns; also, there may be control variables that we are unaware of, resulting in the omitted variables problem.

Following Xu et al. [41] and Ertugrul et al. [42], we select industry-level news sentiment means (*newSentiInd*) but exclude the company and province levels (*newSentiPro*) as the instrument variables for firms' media sentiment. We presume that other publicly traded firms in the same industry or province would face similar industry characteristics and external environments; thus, their media coverage may have a certain correlation. Furthermore, there is no evidence that media coverage of other publicly traded firms in the same industry or province will influence a firm's stock trading behavior, which satisfies the exclusion restriction to some extent.

Table 7 shows the regression results. The coefficients of the *newSentiInd* and *newSentiPro* variables are significantly positive in Column 1, indicating that the higher the media sentiment of listed firms in the industry and province, the higher the mean value of the sentiment of the listed firms. The Cragg–Donald F statistic equals 185.966, which is much larger than the critical value, and this statistic rejects the hypothesis that the instrumental variables are weak at the 1% level. The results of the second stage regression in Columns 2, 3, and 4 show that none of the values of the Sargan statistic reject the original hypothesis of instrumental variable exogeneity. The results of *newSenti* continue to be significantly positive, which is in line with the results of the main regression.

5.3. Channels between News Coverage and Crash Risk. Media sentiment affects the risk of a future stock price crash via two mechanisms. The first is through the investor sentiment, which, in turn, affects stock crash risk. The second is that media coverage serves as an information intermediary, conveying true or false information on listed firms; investors influence the crash risk by interpreting the information they receive.

TABLE 14: Robust test: baseline regression.

	(1)	(2)	(3)	(4)	(5)	(6)
	CRASH	CRASH	NCSKEW	NCSKEW	DUVOL	DUVOL
L..newSenti2	0.100*** (0.037)	0.078** (0.038)	0.561*** (0.091)	0.221** (0.106)	0.317*** (0.060)	0.151** (0.071)
L.CRASH		-0.145*** (0.008)		-0.092*** (0.008)		
L.NCSKEW						-0.101*** (0.008)
L.DUVOL				8.451*** (0.814)		6.011*** (0.565)
L.ret		0.948** (0.391)		0.794* (0.427)		1.015*** (0.288)
L.sigma		-0.589*** (0.188)		1.303*** (0.431)		0.734* (0.292)
L.roa		0.435** (0.210)		-0.142 (0.057)		-0.114*** (0.037)
L.level		0.019 (0.031)		0.000* (0.000)		0.000 (0.000)
L.size		0.000 (0.000)		0.000* (0.000)		0.000 (0.000)
<i>Fixed effects:</i>						
Year dummy	Yes	Yes	Yes	Yes	Yes	Yes
Firm dummy	Yes	Yes	Yes	Yes	Yes	Yes
N	20,107	18,456	20,107	18,456	20,107	18,456
R ²	0.000	0.037	0.002	0.062	0.002	0.068

This table reports robust baseline regression estimates of stock crash risk on the quantitatively weighted news coverage sentiment. The sample of CRASH is used for columns (1)-(2), the sample of NCSKEW is used for columns (3)-(4), and that of DUVOL is used columns (5)-(6). Columns (1), (3), and (5) only include the news coverage sentiment. Columns (2), (4), and (6) control for lagged crash risk and firm characteristic; year and firm dummy variable are included; all control variables are included as lagged one year. Robust standard errors are reported in parentheses. The labels ***, **, and * indicate 1%, 5%, and 10% levels of significance, respectively.

TABLE 15: Robust test: positive sentiment and negative sentiments of news coverage.

	(1) CRASH	(2) NCSKEW	(3) DUVOL	(4) CRASH	(5) NCSKEW	(6) DUVOL
<i>Panel A: one lag period</i>						
L.newPos2	0.161*** (0.061)	0.331** (0.164)	0.228** (0.113)			
L.newNeg2				0.066 (0.068)	0.180 (0.154)	0.045 (0.105)
L.CRASH	-0.147*** (0.008)			-0.147*** (0.009)		
L.NCSKEW		-0.095*** (0.008)			-0.099*** (0.009)	
L.DUVOL			-0.104*** (0.008)			-0.107*** (0.008)
L.ret	1.058*** (0.400)	8.719*** (0.830)	6.152*** (0.577)	0.896** (0.404)	8.770*** (0.844)	6.304*** (0.585)
L.sigma	-0.627*** (0.194)	0.937** (0.439)	1.058*** (0.298)	-0.550*** (0.192)	0.966** (0.440)	1.160*** (0.298)
L.roa	0.417* (0.220)	1.291*** (0.454)	0.773** (0.301)	0.462** (0.218)	1.487*** (0.452)	0.845*** (0.307)
L.level	0.018 (0.032)	-0.116** (0.057)	-0.103*** (0.038)	0.010 (0.034)	-0.163*** (0.060)	-0.141*** (0.040)
L.size	0.000 (0.000)	0.000* (0.000)	0.000 (0.000)	0.000 (0.000)	0.000** (0.000)	0.000 (0.000)
<i>Fix effects:</i>						
Year dummy	Yes	Yes	Yes	Yes	Yes	Yes
Firm dummy	Yes	Yes	Yes	Yes	Yes	Yes
N	17,363	17,363	17,363	16,276	16,276	16,276
R ²	0.036	0.060	0.067	0.036	0.062	0.068
<i>Panel B: same period</i>						
newPos2	-0.065 (0.085)	0.006 (0.180)	0.042 (0.123)			
newNeg2				-0.303*** (0.080)	-0.739*** (0.154)	-0.531*** (0.108)
L.CRASH	-0.153*** (0.009)			-0.154*** (0.009)		
L.NCSKEW		-0.102*** (0.009)			-0.105*** (0.009)	
L.DUVOL			-0.108*** (0.009)			-0.108*** (0.009)
L.ret	0.363 (0.395)	4.913*** (0.844)	3.543*** (0.594)	0.423 (0.407)	5.466*** (0.863)	3.870*** (0.610)
L.sigma	-0.508*** (0.190)	0.146 (0.446)	0.299 (0.307)	-0.519*** (0.192)	0.104 (0.448)	0.297 (0.310)
L.roa	0.661*** (0.207)	1.581*** (0.428)	0.894*** (0.296)	0.524** (0.207)	1.468*** (0.439)	0.919*** (0.304)
L.level	0.053 (0.033)	-0.110* (0.060)	-0.104** (0.041)	0.070** (0.032)	-0.083 (0.059)	-0.083** (0.040)
L.size	-0.000 (0.000)	0.000 (0.000)	0.000 (0.000)	-0.000 (0.000)	0.000 (0.000)	0.000 (0.000)
<i>Fix effects:</i>						
Year dummy	Yes	Yes	Yes	Yes	Yes	Yes
Firm dummy	Yes	Yes	Yes	Yes	Yes	Yes
N	17,363	17,363	17,363	16,276	16,276	16,276
R ²	0.041	0.066	0.072	0.037	0.064	0.070

This table reports robust regression estimates of stock crash risk on the news coverage positive and negative sentiment indicators. Panel A is one year lag; panel B is the same year; year and firm dummy variable are included; all control variables are included as lagged one year. Robust standard errors are reported in parentheses. The labels ***, **, and * indicate 1%, 5%, and 10% levels of significance, respectively.

TABLE 16: Robust test: quarterly crash risk for different future periods.

	(1) CRASH	(2) CRASH	(3) CRASH	(4) CRASH	(5) CRASH	(6) CRASH	(7) CRASH	(8) CRASH
newPos2	-0.049 (0.038)							
L.newPos2		0.096** (0.038)						
L2.newPos2			0.084** (0.039)					
L3.newPos2				0.025 (0.038)				
newNeg2					-0.248*** (0.035)			
L.newNeg2						0.034 (0.033)		
L2.newNeg2							-0.015 (0.035)	
L3.newNeg2								0.027 (0.033)
L.CRASH	-0.048*** (0.005)	-0.049*** (0.005)	-0.047*** (0.005)	-0.047*** (0.005)	-0.049*** (0.005)	-0.049*** (0.005)	-0.048*** (0.005)	-0.045*** (0.005)
L.ret	1.921*** (0.447)	1.909*** (0.450)	2.445*** (0.502)	3.096*** (0.508)	1.837*** (0.457)	1.773*** (0.458)	2.562*** (0.505)	3.184*** (0.525)
L.sigma	-1.686*** (0.237)	-1.464*** (0.234)	-2.182*** (0.259)	-2.174*** (0.262)	-1.780*** (0.239)	-1.411*** (0.239)	-1.931*** (0.263)	-2.287*** (0.264)
L.roa	-0.091** (0.046)	-0.123*** (0.046)	-0.122*** (0.047)	-0.121*** (0.046)	-0.062 (0.047)	-0.124*** (0.046)	-0.141*** (0.049)	-0.118** (0.047)
L.level	-0.037** (0.016)	-0.034* (0.018)	-0.029* (0.017)	-0.018 (0.019)	-0.036** (0.017)	-0.045*** (0.017)	-0.038** (0.017)	-0.029 (0.018)
L.size	0.000 (0.000)	-0.000 (0.000)	0.000 (0.000)	0.000 (0.000)	0.000 (0.000)	0.000 (0.000)	0.000 (0.000)	0.000 (0.000)
<i>Fixed effects:</i>								
Quarter dummy	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Firm dummy	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
N	58,530	58,530	58,530	58,530	59,748	59,748	59,748	59,748
R ²	0.038	0.037	0.039	0.038	0.040	0.039	0.039	0.038

This table reports robust regression estimates of stock crash risk on the news coverage sentiment indicators two stage OLS. Column (1) is the first stage. Columns (2) to (4) are the second stage; year and firm dummy variable are included; all control variables are included as lagged one year. Robust standard errors are reported in parentheses. The labels ***, **, and * indicate 1%, 5%, and 10% levels of significance, respectively.

TABLE 17: Robust test: two stages OLS regression.

	First stage (1) newsenti	(2) CRASH	Second stage (3) NCSKEW	(4) DUVOL
L.newSenti2		0.092** (0.045)	0.994* (0.456)	0.419* (0.246)
L.newSentiInd	0.494*** (0.020)			
L.newSentiPro	0.476*** (0.037)			
Controls	Yes	Yes	Yes	Yes
Cragg-Donald Wald F	185.966***			
Sargan chi (p)		0.121 (0.788)	0.004 (0.947)	0.254 (0.614)
N	16,845	16,845	16,845	16,845
R ²	0.212	0.035	0.059	0.064

TABLE 18: Robust test: regression results for corporate insider trading.

	Low			High		
	(1)	(2)	(3)	(4)	(5)	(6)
	CRASH	NCSKEW	DUVOL	CRASH	NCSKEW	DUVOL
L.newSenti2	0.011 (0.052)	-0.394 (0.333)	-0.154 (0.224)	0.097** (0.048)	0.336*** (0.114)	0.119** (0.059)
L.CRASH	-0.143*** (0.009)			-0.174*** (0.031)		
L.NCSKEW		-0.091*** (0.032)			-0.094*** (0.009)	
L.DUVOL			-0.084*** (0.030)			-0.104*** (0.009)
L.ret	1.834 (1.346)	10.422*** (2.642)	7.613*** (1.858)	1.020** (0.429)	8.353*** (0.917)	5.892 (0.638)
L.sigma	-1.305* (0.721)	2.975* (1.559)	2.061* (1.090)	-0.586*** (0.205)	0.550 (0.467)	0.891*** (0.313)
L.roa	-0.316 (0.801)	1.063 (1.615)	0.783 (1.111)	0.516** (0.218)	1.509*** (0.465)	0.876*** (0.317)
L.level	-0.027 (0.108)	-0.058 (0.208)	-0.020 (0.144)	0.028 (0.033)	-0.138** (0.062)	-0.108*** (0.040)
L.size	-0.000 (0.000)	-0.000 (0.000)	-0.000** (0.000)	0.000 (0.000)	0.000** (0.000)	0.000*** (0.000)
<i>Fixed effects:</i>	Yes	Yes	Yes	Yes	Yes	Yes
Year dummy	Yes	Yes	Yes	Yes	Yes	Yes
Firm dummy						
N	5,488	5,488	5,488	5,884	5,884	5,884
R ²	0.035	0.065	0.071	0.043	0.047	0.067

This table reports robust panel estimates of stock crash risk on the news coverage sentiment by using the subsamples of high insider trading groups and low insider trading groups. The total sample is divided into two subsamples: lower and higher groups, based on 30% and 70% quartiles of insiders' net stock sales, respectively. Columns (1) to (3) are low subgroup, the dependent variables are CRASH, NCSKEW, and DUVOL, respectively. Columns (4) to (6) are high subgroup, the dependent variables are CRASH, NCSKEW, and DUVOL, respectively. All control variables are included as lagged one year. Robust standard errors are reported in parentheses. The labels ***, **, and * indicate 1%, 5%, and 10% levels of significance, respectively.

TABLE 19: Robust test: regression results for information disclosure quality.

	Low			High		
	(1)	(2)	(3)	(4)	(5)	(6)
	CRASH	NCSKEW	DUVOL	CRASH	NCSKEW	DUVOL
L.newSenti2	0.011 (0.098)	0.026 (0.222)	-0.093 (0.144)	0.122*** (0.039)	0.384** (0.179)	0.196** (0.098)
L.CRASH	-0.120*** (0.017)			-0.157*** (0.018)		
L.NCSKEW		-0.094*** (0.020)			-0.101*** (0.014)	
L.DUVOL			-0.100*** (0.018)			-0.117*** (0.015)
L.ret	2.396*** (0.884)	10.914*** (2.111)	7.517*** (1.351)	0.510 (0.787)	5.451*** (1.366)	4.603*** (1.061)
L.sigma	-1.702*** (0.384)	1.233 (0.995)	1.136* (0.631)	-0.165 (0.410)	1.549** (0.712)	1.347** (0.548)
L.roa	-0.277 (0.495)	-0.455 (1.223)	-0.010 (0.789)	0.665* (0.374)	1.002 (0.632)	0.514 (0.468)
L.level	0.140*** (0.040)	0.007 (0.112)	0.033 (0.067)	-0.007 (0.055)	-0.211** (0.097)	-0.205*** (0.072)
L.size	0.000 (0.000)	0.000** (0.000)	0.000 (0.000)	0.000 (0.000)	0.000 (0.000)	0.000 (0.000)
<i>Fixed effects:</i>	Yes	Yes	Yes	Yes	Yes	Yes
Year dummy	Yes	Yes	Yes	Yes	Yes	Yes
Firm dummy						
N	5,906	5,906	5,906	6,243	6,243	6,243
R ²	0.040	0.064	0.071	0.047	0.080	0.080

This table reports robust panel estimates of stock crash risk on the news coverage sentiment by using the subsamples of high information disclosure quality groups and low information disclosure quality groups. The total sample is divided into two subsamples: lower and higher groups, based on 30% and 70% quartiles of KV index respectively. Columns (1) to (3) are low subgroup, the dependent variables are CRASH, NCSKEW, and DUVOL, respectively. Columns (4) to (6) are high subgroup, the dependent variables are CRASH, NCSKEW, and DUVOL, respectively. All control variables are included as lagged one year. Robust standard errors are reported in parentheses. The labels ***, **, and * indicate 1%, 5%, and 10% levels of significance, respectively.

TABLE 20: Robust test: regression results for different analysts coverage.

	Low			High		
	(1)	(2)	(3)	(4)	(5)	(6)
	CRASH	NCSKEW	DUVOL	CRASH	NCSKEW	DUVOL
L.newSenti2	0.074 (0.076)	0.515*** (0.154)	0.345*** (0.107)	0.029 (0.113)	-0.027 (0.267)	-0.073 (0.190)
L.CRASH	-0.178*** (0.012)			-0.203*** (0.027)		
L.NCSKEW		-0.136*** (0.012)			-0.182*** (0.025)	
L.DUVOL			-0.139*** (0.012)			-0.192*** (0.024)
L.ret	0.713 (0.507)	7.801*** (1.019)	5.910*** (0.729)	2.251* (1.342)	6.580** (3.080)	4.741** (2.039)
L.sigma	-0.351 (0.265)	-0.211 (0.594)	-0.666 (0.421)	-1.794*** (0.566)	-1.440 (1.363)	-1.033 (0.868)
L.roa	0.569** (0.277)	1.054* (0.556)	0.414 (0.380)	0.179 (0.742)	0.860 (1.585)	0.371 (1.046)
L.level	-0.010 (0.040)	-0.106 (0.087)	-0.110* (0.062)	0.020 (0.086)	-0.181 (0.206)	-0.090 (0.128)
L.size	0.000 (0.000)	0.000** (0.000)	0.000*** (0.000)	0.000 (0.000)	0.000* (0.000)	-0.000 (0.000)
<i>Fixed effects:</i>	Yes	Yes	Yes	Yes	Yes	Yes
Year dummy	Yes	Yes	Yes	Yes	Yes	Yes
Firm dummy						
N	5,028	5,028	5,028	4,809	4,809	4,809
R ²	0.067	0.069	0.068	0.071	0.118	0.119

This table reports robust panel estimates of stock crash risk on the news coverage sentiment by using the subsamples of low analysts' attention groups and high analysts' attention groups. We used the total number of analysts coverage of firms to measure analysts' attention. We used the 30% and 70% quartiles of analysts' attention as the cut-off, and the firms were divided into low and high attention groups. Columns (1) to (3) are low subgroup, the dependent variables are CRASH, NCSKEW, and DUVOL, respectively. Columns (4) to (6) are high subgroup, the dependent variables are CRASH, NCSKEW, and DUVOL, respectively. All control variables are included as lagged one year. Robust standard errors are reported in parentheses. The labels ***, **, and * indicate 1%, 5%, and 10% levels of significance, respectively.

TABLE 21: Robust test: regression results for different institutional holding group.

	Low			High		
	(1)	(2)	(3)	(4)	(5)	(6)
	CRASH	NCSKEW	DUVOL	CRASH	NCSKEW	DUVOL
L.newSenti2	0.102** (0.044)	0.481*** (0.176)	0.358*** (0.124)	0.028 (0.106)	0.037 (0.220)	-0.027 (0.144)
L.CRASH	-0.198*** (0.020)			-0.174*** (0.015)		
L.NCSKEW		-0.131*** (0.015)			-0.161*** (0.020)	
L.DUVOL			-0.170*** (0.014)			-0.182*** (0.019)
L.ret	0.698 (0.949)	6.795*** (1.287)	5.479*** (0.932)	0.933 (0.670)	2.429 (2.091)	1.849 (1.426)
L.sigma	-0.777* (0.423)	0.043 (0.696)	0.806 (0.503)	-0.621* (0.335)	0.018 (0.970)	0.540 (0.642)
L.roa	0.464 (0.509)	0.860 (0.655)	0.321 (0.505)	0.277 (0.362)	0.723 (1.178)	0.699 (0.737)
L.level	-0.094 (0.065)	-0.074 (0.068)	-0.085* (0.048)	0.064 (0.048)	-0.367*** (0.159)	-0.156 (0.104)
L.size	-0.000 (0.000)	0.000*** (0.000)	0.000*** (0.000)	0.000 (0.000)	0.000 (0.000)	0.000 (0.000)
<i>Fixed effects:</i>	Yes	Yes	Yes	Yes	Yes	Yes
Year dummy	Yes	Yes	Yes	Yes	Yes	Yes
Firm dummy						
N	5,120	5,120	5,120	6,959	6,959	6,959
R ²	0.046	0.071	0.086	0.067	0.084	0.095

This table reports robust panel estimates of stock crash risk on the news coverage sentiment by using the subsamples of low institutional shareholding groups and high institutional shareholding groups. We used the 30% and 70% quartiles of the institutional shareholding as the cut-off point, and the firms were divided into low and high shareholding groups. Columns (1) to (3) are low subgroup, the dependent variables are CRASH, NCSKEW, and DUVOL, respectively. Columns (4) to (6) are high subgroup, the dependent variables are CRASH, NCSKEW, and DUVOL, respectively. All control variables are included as lagged one year. Robust standard errors are reported in parentheses. The labels ***, **, and * indicate 1%, 5%, and 10% levels of significance, respectively.

TABLE 22: Robust test: regression results for different investor sentiment.

	Low			High		
	(1)	(2)	(3)	(4)	(5)	(6)
	CRASH	NCSKEW	DUVOL	CRASH	NCSKEW	DUVOL
L.newSenti2	-0.039 (0.101)	0.022 (0.216)	-0.017 (0.147)	0.173*** (0.063)	0.544*** (0.210)	0.348** (0.137)
L.CRASH	-0.177*** (0.018)			-0.165*** (0.015)		
L.NCSKEW		-0.141*** (0.018)			-0.103*** (0.017)	
L.DUVOL			-0.132*** (0.018)			-0.125*** (0.016)
L.ret	2.871*** (0.947)	6.671*** (1.836)	3.952*** (1.283)	2.329*** (0.710)	6.885*** (1.703)	4.676*** (1.149)
L.sigma	-0.537 (0.422)	0.860 (0.901)	0.689 (0.592)	0.784** (0.359)	0.694 (0.865)	0.072 (0.590)
L.roa	0.751 (0.574)	1.287 (0.945)	0.873 (0.612)	0.676* (0.390)	2.094** (0.875)	1.581** (0.616)
L.level	0.035 (0.064)	-0.193* (0.101)	-0.109* (0.064)	0.023 (0.051)	-0.084 (0.120)	-0.078 (0.082)
L.size	0.000 (0.000)	0.000 (0.000)	0.000 (0.000)	0.000 (0.000)	0.000 (0.000)	0.000 (0.000)
<i>Fixed effects:</i>	Yes	Yes	Yes	Yes	Yes	Yes
Year dummy	Yes	Yes	Yes	Yes	Yes	Yes
Firm dummy						
N	5,213	5,213	5,213	6,258	6,258	6,258
R ²	0.047	0.077	0.085	0.051	0.082	0.085

This table reports robust panel estimates of stock crash risk on the news coverage sentiment by using the subsamples of pessimistic investor groups and optimistic investor groups. We used the 30% and 70% quartiles of the investor sentiment as the cut-off point, and the firms were divided into low and high shareholding groups. Columns (1) to (3) are low subgroup, the dependent variables are CRASH, NCSKEW, and DUVOL, respectively. Columns (4) to (6) are high subgroup, the dependent variables are CRASH, NCSKEW, and DUVOL, respectively. All control variables are included as lagged one year. Robust standard errors are reported in parentheses. The labels ***, **, and * indicate 1%, 5%, and 10% levels of significance, respectively.

TABLE 23: Robust test: regression of adding analyst and media coverage sentiment crossterms.

	(1)	(2)	(3)
	CRASH	NCSKEW	DUVOL
L.newSenti2	0.238** (0.117)	0.618** (0.269)	0.336* (0.185)
L.anaSenti	0.035*** (0.011)	0.061*** (0.023)	0.050*** (0.016)
L.newSenti2#L.anaSenti	-0.121** (0.057)	-0.424* (0.235)	-0.186* (0.105)
L.CRASH	-0.150*** (0.009)		
L.NCSKEW		-0.099*** (0.009)	
L.DUVOL			-0.108*** (0.009)
L.ret	0.933** (0.406)	8.405*** (0.839)	5.992*** (0.584)
L.sigma	0.528*** (0.198)	0.664*** (0.144)	0.953*** (0.303)
L.roa	0.377* (0.217)	1.260*** (0.453)	0.688** (0.308)
L.level	0.014 (0.034)	-0.146** (0.060)	-0.112*** (0.039)
L.size	0.000 (0.000)	0.000** (0.000)	0.000 (0.000)
<i>Fixed effects:</i>			
Year dummy	Yes	Yes	Yes
Firm dummy	Yes	Yes	Yes
N	16,829	16,829	16,829
R ²	0.039	0.066	0.072

This table reports robust panel estimates of stock crash risk on the interaction between the news coverage sentiment and the analysis report sentiment. The dependent variables in columns (1), (2), and (3) are CRASH, NCSKEW, and DUVOL, respectively. All control variables are included as lagged one year. Robust standard errors are reported in parentheses. The labels ***, **, and * indicate 1%, 5%, and 10% levels of significance, respectively.

5.3.1. Information Intermediation. Insiders have information that is not yet publicly available, which can be used to judge the value of the firm and predict future firm performance [34]. Insider sell-off behavior is positively associated with stock price crash risk [35]. The insider's choice to sell stocks sends negative signals to outside investors, thereby raising the probability of a future crash risk. Furthermore, the more overpriced a stock is, the greater the chance of a crash.

We thus divide the total sample into two subsamples, lower and higher groups, based on 30% and 70% quartiles of insiders' net stock sales, respectively. The regression results for various groups are shown in Table 8. The coefficients of the high net selling subgroups are significant, whereas those coefficients of the low net selling subgroups are insignificant. Thus, insiders sell more stocks, thereby amplifying the impact of media sentiment on the risk of a future stock price crash, confirming *H3a*.

We examine the quality of firm disclosure to determine if it would mitigate the bubble created by media and reduce the

risk of a stock price crash. We followed the method of Kim and Verrecchia [43] (KV index) to measure the quality of information disclosure. The higher the KV index, the lower the quality of the information disclosure of listed firms.

We again divide the total sample into two subsamples based on the 30% and 70% quartiles of the KV index into lower and higher groups. The regression results for the different groups are presented in Table 9. The coefficients of the lower disclosure group are significant, whereas those coefficients of the higher disclosure group are insignificant. Thus, the effect of news coverage on the risk of a future stock price crash is significantly enhanced in firms with poor disclosure. Meanwhile, the effect of media sentiment was significantly weaker when the quality of firm disclosure is high. Thus, *H3b* is supported.

Next, we examine the Hypothesis *H3c*. Follow He et al. [35], we calculate the analysts coverage as the number of analysts forecast over the past three years. Then, we divide the total sample into two subsamples: lower coverage and higher coverage groups, according to 30% and 70% quartiles of analysts' coverage to firms. Table 10 presents the regression results.

The coefficients of the low analyst coverage groups are significant, indicating that positive media coverage in the previous year increases the future stock price crash risk of firms with lower analyst coverage. The coefficients of the high analyst coverage groups are insignificant, and the impact of news coverage on stock price crash risk is attenuated. Thus, news coverage sentiment has a stronger impact on stock price crash risk when analysts pay less attention to a firm, supporting *H3c*.

5.3.2. Investor Sentiment. In exploring the investor sentiment channel, we investigate whether an increase in number of retail investors could increase the emotional impact of the media. Table 11 shows the regression results for different groups, divided into low and high groups based on 30% and 70% quartiles of institutional holding.

The coefficients for the low institutional holding subgroups are significant, whereas those coefficients for the high institutional holding subgroups are insignificant. The findings suggest that as retail investors increase their holdings, they tend to behave more irrationally, thereby amplifying the emotional tendency of media sentiment. Our results support *H4a*.

We consider direct proxy variables for investor sentiment and construct investor sentiment indicators according to Rhodes-Kropf et al. [44], dividing the total sample into pessimistic and optimistic groups based on 30% and 70% quartiles of investor sentiment. The regression results for different groups are shown in Table 12. The coefficients for the pessimistic subgroups are insignificant, whereas those coefficients for the optimistic subgroups are significant. The findings suggest that optimistic investor sentiment, which increases the likelihood that the current stock price is overvalued, increases the impact of news coverage on future crash risk, supporting *H4b*.

To investigate whether disagreement between professional analysts and retail investors also enhances the risk of a stock price crash, we add the cross-term of analysts' sentiment and media sentiment into the regression (Table 13).

The proxy variable used in this study to represent consistency and disagreement is the crossterm of analysts and news sentiment. When there are significant differences between the two opinions, often one is less than 0 and the other is greater than 0, so the cross-term is negative and represents the differences in opinions. Conversely, the two types of views have the same symbol and positive multiplication, implying that they are consistent.

The results of the regression are reported in Table 13. The results demonstrate the impact of analyst-rated sentiment on future stock price crash risk. The coefficient of the cross-term is negative, consistent with our theoretical hypothesis that the future stock price crash risk decreases when media sentiment and analyst sentiments are consistent, and increases otherwise. The regression results show that divergence of opinions in the market could increase media sentiment tendencies, thereby supporting *H4c*.

5.4. Robust Test. *newSenti**newSenti2* To test the robustness of the indicators, we replace the main explanatory variable of the previous regression (*newSenti*), which is the weighted average of news sentiment, with an equally weighted average (*newSenti2*) and repeat the previous regressions. The results are shown in Tables 14--23 and are consistent with those in the previous contents.

This table reports robust regression estimates of quarterly stock crash risk on the news coverage positive and negative sentiment indicators. Columns (1) to (4) are positive indicators. Columns (5) to (8) are negative indicators; quarter and firm dummy variable are included; all control variables are included as lagged one year. Robust standard errors are reported in parentheses. The labels ***, **, and * indicate 1%, 5%, and 10% levels of significance, respectively.

6. Conclusions

We construct a deep learning model of stock news sentiment recognition based on the advanced approach of financial knowledge dictionary and NLP (BERT-based pretraining) technology. We use this model to calculate the sentiment indicators of all stocks from 2011 to 2020. Subsequently, we analyze the impact of media sentiment on future stock price crash risk and its heterogeneity.

We find that average media sentiment exacerbates the risk of future stock price crashes. The heterogeneity results indicate that positive coverage significantly increases future stock price crash risk, whereas negative coverage has a limited effect. However, negative coverage is highly correlated with current stock price crash risk. We also investigate the information intermediation and investor sentiment channels by which media sentiment affects the risk of a crash. The results show that more net insider sales, lower information transparency, and less analyst coverage amplify the impact of media sentiment on future crash risk, which is

consistent with the information intermediation channel. Additionally, more retail investor positions, more active investor sentiment, and divergence between analysts' opinions and news amplify the impact of news sentiment on the risk of a future stock price crash, which is consistent with the investor sentiment channel.

Our finding that positive media sentiment can lead to an extreme outcome in the stock market is useful for both regulators and investors. Our examination of the impact of media sentiment on the future stock price crash risk adopted both behavioral finance and information economics perspectives, and revealed that investors' irrational and excessive optimism could be a major cause of stock price bubbles and crashes in China's stock market, which is dominated by retail investors who are restricted from short selling.

In terms of research methodology, our study combined advanced deep learning and dictionary methods, fully utilizing computer performance and intelligence to significantly improve the recognition accuracy and efficiency of massive amounts of sentiment data. To the best of our knowledge, we are the first to combine deep learning and dictionary methods for sentiment analysis in finance, thereby broadening the scope of sentiment analysis methods in finance.

Data Availability

The data copyright belongs to the GuoTai'an and Wind, disclosing is not allowed.

Conflicts of Interest

The authors declare that there are no conflicts of interest regarding the publication of this paper.

References

- [1] E. Noelle-Neumann, "The spiral of silence a theory of public opinion," *Journal of Communication*, vol. 24, no. 2, pp. 43-51, 1974.
- [2] H. Hong and J. C. Stein, "Differences of opinion, short-sales constraints and market crashes," *Review of Financial Studies*, vol. 16, no. 2, pp. 487-525, 2003.
- [3] A. P. Hutton, A. J. Marcus, and H. Tehrani, "Opaque financial reports, R2, and crash risk," *Journal of Financial Economics*, vol. 94, no. 1, pp. 67-86, 2009.
- [4] D. F. Larcker and A. A. Zakolyukina, "Detecting deceptive discussions in conference calls," *Journal of Accounting Research*, vol. 50, no. 2, pp. 495-540, 2012.
- [5] W. S. Chan, "Stock price reaction to news and no-news: drift and reversal after headlines," *Journal of Financial Economics*, vol. 70, no. 2, pp. 223-260, 2003.
- [6] J. R. Joe, H. Louis, and D. Robinson, "Managers' and investors' responses to media exposure of board ineffectiveness," *Journal of Financial and Quantitative Analysis*, vol. 44, no. 3, pp. 579-605, 2009.
- [7] P. C. Tetlock, "Giving content to investor sentiment: the role of media in the stock market," *The Journal of Finance*, vol. 62, no. 3, pp. 1139-1168, 2007.

- [8] P. C. Tetlock, M. Saar-Tsechansky, and S. Macskassy, "More than words: quantifying language to measure firms' fundamentals," *The Journal of Finance*, vol. 63, no. 3, pp. 1437–1467, 2008.
- [9] S. R. Das, "Text and context: language analytics in finance," *Foundations and Trends® in Finance*, vol. 8, no. 3, pp. 145–261, 2014.
- [10] M. El-Haj, P. Rayson, M. Walker, S. Young, and V. Simaki, "In search of meaning: I," *Journal of Business Finance & Accounting*, vol. 46, no. 3-4, pp. 265–306, 2019.
- [11] F. Li, "The information content of forward-looking statements in corporate filings-A naïve bayesian machine learning approach," *Journal of Accounting Research*, vol. 48, no. 5, pp. 1049–1102, 2010.
- [12] T. I. M. Loughran and B. McDonald, "Textual analysis in accounting and finance: a survey," *Journal of Accounting Research*, vol. 54, no. 4, pp. 1187–1230, 2016.
- [13] J. You, B. Zhang, and L. Zhang, "Who captures the power of the pen?" *Review of Financial Studies*, vol. 31, no. 1, pp. 43–96, 2017.
- [14] J. Chen, H. Hong, and J. C. Stein, "Forecasting crashes: trading volume, past returns, and conditional skewness in stock prices," *Journal of Financial Economics*, vol. 61, no. 3, pp. 345–381, 2001.
- [15] C. R. Harvey and A. Siddique, "Conditional skewness in asset pricing tests," *The Journal of Finance*, vol. 55, no. 3, pp. 1263–1295, 2000.
- [16] J. B. Kim, Y. Li, and L. Zhang, "CFOs versus CEOs: equity incentives and crashes," *Journal of Financial Economics*, vol. 101, no. 3, pp. 713–730, 2011.
- [17] B. J. Bushee, J. E. Core, W. Guay, and S. J. Hamm, "The role of the business press as an information intermediary," *Journal of Accounting Research*, vol. 48, pp. 1–19, 2010.
- [18] Q. Li, J. Wang, and L. Bao, "Media tone, bias, and stock price crash risk: evidence from China," *Asia-Pacific Journal of Accounting & Economics*, vol. 29, pp. 1–35, 2022.
- [19] X. Cui, A. Sensoy, D. K. Nguyen, S. Yao, and Y. Wu, "Positive information shocks, investor behavior and stock price crash risk," *Journal of Economic Behavior & Organization*, vol. 197, pp. 493–518, 2022.
- [20] A. Hillert, H. Jacobs, and S. Müller, "Media makes momentum," *Review of Financial Studies*, vol. 27, no. 12, pp. 3467–3501, 2014.
- [21] T. Odean, "Volume, volatility, price, and profit. When all traders are above average," *The Journal of Finance*, vol. 53, no. 6, pp. 1887–1934, 1998.
- [22] S. Pouget, J. Sauvagnat, and S. Villeneuve, "A mind is a terrible thing to change: confirmatory bias in financial markets," *Review of Financial Studies*, vol. 30, no. 6, pp. 2066–2109, 2017.
- [23] K. R. Ahern and D. Sosyura, "Rumor has it: sensational is min financial media," *Review of Financial Studies*, vol. 28, pp. 2050–2093, 2015.
- [24] K. R. Ahern and D. Sosyura, "Who writes the news? Corporate press releases during merger negotiations," *The Journal of Finance*, vol. 69, pp. 241–291, 2014.
- [25] J. Engelberg, C. Sasseville, and J. Williams, "Madness? The case of mad money," *Management Science*, vol. 58, no. 2, pp. 351–364, 2012.
- [26] S. P. Fraiberger, D. Lee, D. Puy, and R. Ranciere, "Media sentiment and international asset prices," *IMF Working Papers*, vol. 18, no. 274, p. 1, 2018.
- [27] B. Gan, V. Alexeev, R. Bird, and D. Yeung, "Sensitivity to sentiment: news vs social media," *International Review of Financial Analysis*, vol. 67, Article ID 101390, 2020.
- [28] D. H. Solomon, "Selective publicity and stock prices," *The Journal of Finance*, vol. 67, no. 2, pp. 599–638, 2012.
- [29] T. A. Hassan, S. Hollander, L. van Lent, and A. Tahoun, "Firm-level political risk: measurement and effects," *Quarterly Journal of Economics*, vol. 134, no. 4, pp. 2135–2202, 2019.
- [30] S. P. Kothari, S. Shu, and P. D. Wysocki, "Do managers withhold bad news?" *Journal of Accounting Research*, vol. 47, no. 1, pp. 241–276, 2009.
- [31] P. Rozin and E. B. Royzman, "Negativity bias, negativity dominance, and contagion," *Personality and Social Psychology Review*, vol. 5, no. 4, pp. 296–320, 2001.
- [32] L. Zhang, C. Wang, and H. Yao, "Information contagion and stock price crash risk," *Mathematical Problems in Engineering*, vol. 2021, Article ID 8891338, 2021.
- [33] L. Jin and S. Myers, "R2 around the world: new theory and new tests," *Journal of Financial Economics*, vol. 79, no. 2, pp. 257–292, 2006.
- [34] J. D. Piotroski and D. T. Roulstone, "Do insider trades reflect both contrarian beliefs and superior knowledge about future cash flow realizations?" *Journal of Accounting and Economics*, vol. 39, no. 1, pp. 55–81, 2005.
- [35] G. He, H. M. Ren, and R. Taffler, "Do corporate insiders trade of future stock price crash risk?" *Review of Quantitative Finance and Accounting*, vol. 56, no. 4, pp. 1561–1591, 2021.
- [36] K. Schipper, "Analysts' forecasts," *Accounting Horizons*, vol. 5, pp. 105–121, 1991.
- [37] S. C. Andrade, J. Bian, and T. R. Burch, "Analyst coverage, information, and bubbles," *Journal of Financial and Quantitative Analysis*, vol. 48, no. 5, pp. 1573–1605, 2013.
- [38] C. T. Li, F. M. Song, and X. Zhang, "Analyst following and corporate earnings management: evidence from China," *Journal of Financial Research*, vol. 7, pp. 124–139, 2014.
- [39] P. C. Tetlock, "Does public financial news resolve asymmetric information?" *Review of Financial Studies*, vol. 23, no. 9, pp. 3520–3557, 2010.
- [40] N. Xu, X. Jiang, Z. Yi, and X. Xu, "Conflicts of interest, analyst optimism and stock price crash risk," *Economic Research Journal*, vol. 7, pp. 127–140, 2012.
- [41] N. Xu, X. Li, Q. Yuan, and K. C. Chan, "Excess perks and stock price crash risk: evidence from China," *Journal of Corporate Finance*, vol. 25, pp. 419–434, 2014.
- [42] M. Ertugrul, J. Lei, J. Qiu, and C. Wan, "Annual report readability, tone ambiguity, and the cost of borrowing," *Journal of Financial and Quantitative Analysis*, vol. 52, no. 2, pp. 811–836, 2017.
- [43] O. Kim and R. E. Verrecchia, "The relation among disclosure, returns, and trading volume information," *The Accounting Review*, vol. 76, no. 4, pp. 633–654, 2001.
- [44] M. Rhodes-Kropf, D. T. Robinson, and S. Viswanathan, "Valuation waves and merger activity: the empirical evidence," *Journal of Financial Economics*, vol. 77, no. 3, pp. 561–603, 2005.