





CrudeBERT: Applying Economic Theory Towards Fine-Tuning Transformer-Based Sentiment Analysis Models to the Crude Oil Market

Himmet Kaplan¹ ^a, Ralf-Peter Mundani² ^b, Heiko Rölke² ^c and Albert Weichselbraun² ^d

¹Zurich University of Applied Sciences, Winterthur, Switzerland

²University of Applied Sciences of the Grisons, Chur, Switzerland

Keywords: Natural Language Processing, Sentiment Analysis, Transformers, FinBERT, Crude Oil Market, Fine-Tuning.


Abstract: Predicting market movements based on the sentiment of news media has a long tradition in data analysis. With advances in natural language processing, transformer architectures have emerged that enable contextually aware sentiment classification. Nevertheless, current methods built for the general financial market such as FinBERT cannot distinguish asset-specific value-driving factors. This paper addresses this shortcoming by presenting a method that identifies and classifies events that impact supply and demand in the crude oil markets within a large corpus of relevant news headlines. We then introduce CrudeBERT, a new sentiment analysis model that draws upon these events to contextualize and fine-tune FinBERT, thereby yielding improved sentiment classifications for headlines related to the crude oil futures market. An extensive evaluation demonstrates that CrudeBERT outperforms proprietary and open-source solutions in the domain of crude oil.


1 INTRODUCTION


Crude oil is one of our primary energy sources and also one of the most influential raw materials. Thus, it is of utmost importance for the global economy and even serves as an indicator of economic boom or recession. Since crude oil is a limited natural resource, its price is expected to be determined by supply and demand. Yet, according to literature, crude oil is one of the most volatile markets in the world since its demand is primarily affected by economic activity (business cycle) and exogenous events such as armed conflicts and natural disasters (Buyuksahin and Harris, 2011). Traditionally, analysts draw upon technical analysis which utilizes historical data for prediction. However, historical data rarely provides high-confidence insights (McCarthy et al., 2019). Complementing technical analysis with additional contemporary and relevant information such as news articles could be a promising strategy for achieving more reliable results. Several empirical studies demonstrated that considering news media significantly improved forecasts of large market movements (i.e., higher than 50%) of publicly listed assets (Qian and Rasheed, 2007). Therefore, many researchers studied the ben-


efits of incorporating news data into their prediction models (Baboshkin and Uandykova, 2021). One option for applying news to prediction tasks comes with sentiment analysis which quantifies the impact of news on a certain asset as positive, neutral, or negative. With the latest advancement in computer hardware and software, particularly the development of transformer architectures (Devlin et al., 2018), modern natural language processing (NLP) algorithms emerged that are capable of evaluating text in a contextually aware manner for strategic forecasting. The observations of Jiang et al. (2020) indicate that current transformer-based sentiment classifiers can achieve remarkable accuracies of up to 97.5%. While these sentiment analysis methods gained great traction in prediction tasks for the stock market and cryptocurrencies, they still play only a minor role in forecasting crude oil prices. Thus, the benefits of considering the sentiment of news headlines for crude oil price predictions seem to be evident.

The presented research draws upon FinBERT, a state-of-the-art transformer-based sentiment analysis model that has been pre-trained for the general financial market. However, analyzing over a decade of news headlines relevant to the crude oil market revealed that FinBERT's sentiment classification does not deliver any apparent insights into the contemporary development of oil prices. Therefore, we developed the publicly available CrudeBERT sentiment analysis model that has been optimized for the crude

^a  <https://orcid.org/0000-0002-1115-8669>

^b  <https://orcid.org/0000-0001-6248-714X>

^c  <https://orcid.org/0000-0002-9141-0886>

^d  <https://orcid.org/0000-0001-6399-045X>

oil domain. CrudeBERT extends FinBERT by considering the economic theory of supply and demand. In our experiments, CrudeBERT outperforms FinBERT and provides a promising tool for improving crude oil price predictions by incorporating information on the sentiment conveyed in news headlines.

The main contributions of this paper can be summarized as (i) developing a method that provides transformer models with means for identifying the major supply and demand factors that drive crude oil futures markets, (ii) fine-tuning general transformer-based sentiment analysis methods by incorporating the economic model of supply and demand into these models, and (iii) conducting extensive experiments that draw upon multiple prediction settings to benchmark the developed method against a baseline (random binary classification) and two state-of-the-art (lexicon- and transformer-based) sentiment analysis frameworks.

The remainder of this paper is organized as follows: Chapter 2 discusses related literature that led to the modern transformer-era sentiment analysis applications. Afterward, chapter 3 introduces the FinBERT domain-specific affective model for the domain of crude oil markets and its use in predicting market movements. Chapter 4 describes the evaluation setup, performs a comprehensive evaluation of the CrudeBERT model, and discusses the obtained results. Chapter 5 concludes the paper with a summary and an outlook on future improvements.

2 RELATED WORK

From an industrial standpoint, crude oil is critical to the world's economy. Consequently, many research articles focus on predicting its price with studies varying from technical to fundamental analysis. This literature review focuses on articles aimed at forecasting crude oil prices by including sentiment features.

2.1 Efficient Market Hypothesis

The efficient market hypothesis (EMH) questions whether information retrieved from news articles does contain any predictive value at all since it claims that the price of an asset already considers all publicly available information. Eugene Fama distinguishes between the weak, semi-strong, and strong forms of EMH (Fama, 1970). The weak form claims that the price results only from its historical price history, thus making all available information outside the historical price relevant for forecasting the future price of an asset. The EMH's semi-strong variant on the other hand

considers that the current pricing reflects the historical prices and publicly available information. Therefore, confidential information such as insider knowledge can add value to a prediction, given it hasn't yet altered the current price (Malkiel, 1989). The strong form of EMH assumes that prices reflect historical prices, and publicly available, and confidential information (Fama, 1970). Hence, the strong form assumes that applying fundamental analysis based on any available information cannot lead to abnormal economic returns. This form of the EMH is further supported by various studies that emphasize the notorious difficulty of forecasting crude oil prices, such as the works of Hamilton, which concludes that the oil price appears to be influenced by a random walk with drift (Hamilton, 2008).

Yet, numerous experts have questioned the hypothesis of the EMH's strong and semi-strong forms, claiming that once a news message is published, the available information changes, and, therefore, the price is expected to adapt. In the experiments of Qian and Rasheed, they concluded that the predictions based on news can correctly forecast price fluctuations with greater than 50% accuracy (Qian and Rasheed, 2007). Furthermore, according to Buyuksahin and Harris, crude oil demand is primarily driven by exogenous events, such as armed conflicts and natural disasters as well as the presence of speculators such as noise traders. They assume that these events considerably contribute towards making crude oil one of the most volatile markets in the world. They also observed a substantial relationship between crude oil price changes and the behavior of politically and economically unstable nations, which often trigger such exogenous events (Buyuksahin and Harris, 2011). This observation is confirmed by Brandt and Gao's more recent study, which shows that macroeconomic and geopolitical news has a strong influence on crude oil, with varying impacts. For example, macroeconomic news influences short-term price movements and also helps to forecast long-term oil prices. On the other hand, the influence of geopolitical news yields typically a robust and instantaneous impact that results in increased volume in trade. However, geopolitical news delivers no conclusive insights in terms of forecasting (Brandt and Gao, 2019). Wex et al. state that forecasts based on the sentiment scores of news articles that cover exogenous events are statistically significant (Wex et al., 2013).

2.2 Sentiment Analysis

Sentiment analysis is considered a prevalent classification task in NLP, which categorizes affective and

subjective information within entire documents, paragraphs, and sentences. It has gained increasing popularity due to its vast potential for a variety of applications such as economics, finance, marketing, political science, psychology, and human-computer interaction (Mohammad, 2021). Sentiment in the context of sentiment analysis, which is also known as opinion mining, refers to the quantification of natural language within pre-defined affective dimensions (Weichselbraun et al., 2022) such as sentiment polarity which distinguishes between positive, neutral, and negative media coverage. Still, there is little discussion about what sentiment in the context of NLP truly represents (Hovy, 2015). Generally, researchers assume that authors always express some sentiment while producing natural language, since emotions, opinions, and expressions in language are fundamental human traits (Taboada, 2016). Therefore, sentiment analysis can also cover complex emotions such as the ones introduced in Plutchik's Wheel of Emotions (Plutchik, 1982) and the Hourglass of Emotions (Susanto et al., 2020). However, most literature in finance tends to break sentiment down into attitudes using binary polarities such as positive and negative, also known as financial sentiment analysis (FSA). Hence a binary classification is more suitable for directly assessing the up or down movements of publicly traded assets (Li et al., 2014).

2.3 Early NLP Methods for FSA

Sentiment analysis in the financial domain has been introduced in the 1980s. One of the first approaches to classifying sentiment in text documents was the Bag-of-Words (BOW) methodology, often referred to as the lexicon-based technique (Liew, 2016). Since a text consists of several words (tokens), BOW simply accumulates the sentiment scores of positive and negative words to compute the overall sentiment classification. As the name suggests, these BOW methods utilize a lexicon consisting of words and their sentimental value, predetermined by, ideally multiple, human annotators. One of the most well-known lexicons for FSA was developed by Loughran and McDonald and aimed at interpreting liabilities concerning 10-K filing returns (Loughran and McDonald, 2011). In their later works (Loughran and McDonald, 2016), they published a survey on the use of text analysis with a focus on accounting and finance. However, creating lexicons that include all possible keywords including negates, or word combinations is very challenging, since a term's sentiment often also depends on the context expressed by surrounding terms or paragraphs.

2.4 Machine Learning-Based Sentiment Analysis

With advancements in computer hard- and software, modern FSA approaches started to heavily rely on machine learning-based approaches. These approaches mostly focused on supervised learning, in which learning is accomplished by training on annotated datasets containing pairs of inputs and matching solutions. As a result, by correcting and optimizing themselves based on the, mostly human-curated, training dataset the machine learning-based approaches identify rules and patterns and attempt to derive a potentially meaningful generalization. These approaches, which require large amounts of annotated training data, are known as supervised learning and are usually used for classification and regression (Chollet, 2018) tasks. For instance, Recurrent Neural Networks (RNN) and Long Short-Term Memory (LSTM) are particularly well suited for sequential data, such as text (Tang et al., 2016). However, the required training dataset is one disadvantage of supervised learning, particularly for classification problems such as sentiment analysis. Since a bigger training dataset usually yields better results using large training datasets is typically expensive (both computationally and financially). Furthermore, RNNs suffer from vanishing and exploding gradients, making them unsuitable for lengthy texts, and are slower to train since their sequential flow is incompatible with parallel processing (Chollet, 2018). One approach to addressing this problem is combining supervised approaches with unsupervised machine learning models. For instance, word embeddings (also known as word vector models), map words into a vector space that aligns semantically related words close to each other in an unsupervised manner. Among the most popular word vector models are word2vec (Mikolov et al., 2013) and Glove (Pennington et al., 2014) which can be trained on large text corpora, therefore, capturing the word semantics within the corpus. Nevertheless, word embeddings still lack the power to fully consider a term's context – i.e., once a model has been trained, words with the same spelling always receive the same word vector, independent of their context – i.e., orange either as a fruit or as a color, or a mixture of both concepts.

2.5 Transformer-Based Sentiment Analysis

More recent language models draw upon the attention mechanism (Bahdanau et al., 2016) which led to significant gains in a range of NLP tasks including senti-

ment analysis. The attention mechanism allows neural networks to resemble the human cognitive function by selectively focusing on particularly relevant information while dismissing other less relevant information. This approach encourages the neural network to spend more computational resources on small but relevant elements of the data (Bahdanau et al., 2016) yielding improvements in terms of speed and accuracy. Further enhancements from Vaswani et al. made use of the attention mechanism to develop the transformer architecture, which allows parallel training and makes it more efficient than RNNs. Initially, the transformer architecture was proposed for neural machine translation, thus, it contains an encoder and decoder. The encoder is a fully connected feed-forward network made out of multiple identical multi-headed attention layers which allows the sequence to be evaluated from contextually varying perspectives (Figure 1). The capability to consider a term's context together with the option to draw upon and customize large pre-trained models has been key to the success of transformer-based language models for sentiment analysis (Vaswani et al., 2017).

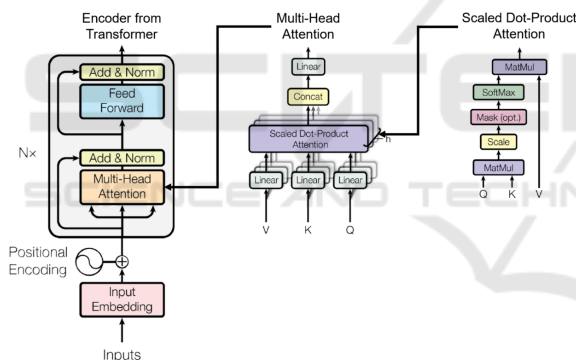


Figure 1: Components of the Multi-Head Attention Design. (Vaswani et al., 2017).

2.6 FinBERT for Financial Sentiment Analysis

Shortly after the release of the transformer architecture Devlin et al. observed that the encoder, when layered, can also serve as a strong representation learning model and for this matter, they developed the Bidirectional Encoder Representations from Transformers (BERT) (Devlin et al., 2018). One noteworthy feature of BERT was its simple customization for a wide range of NLP tasks with the important added capability of contextual perception of words (Yenicelik, 2020). Initially, it was pre-trained on English *Wikipedia* and the *BookCorpus* to give the model a general comprehension of natural language (Zhu

et al., 2015). This model served as a foundation for further adaptations to specific NLP applications and its domain, such as FinBERT (Araci, 2019) which focuses on sentiment analysis of financial news. To achieve this, Araci et al. employed a subset of the *Thomson Reuters Text Research Collection (TRC2)* to adapt the model to the domain of financial news, where occurrences of slang and spelling errors are minimal. For the task-specific fine-tuning process, the training dataset *Financial Phrase Bank* from Malo et al. (2014) was utilized (Figure 2).

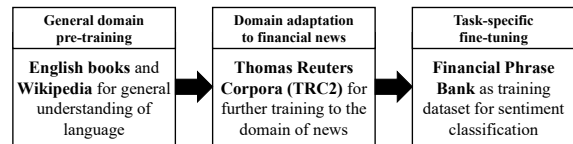


Figure 2: Process of generating FinBERT.

Compared to the number of papers that used FinBERT as a backend for classification, the proportion of papers that use it for classifying sentiments towards crude oil is relatively small.

2.7 RavenPack Event Sentiment Score

To capture the overall sentiment of the market, *RavenPack* developed a lexicon-based news sentiment index namely the *Event Sentiment Score (ESS)*, which is a granular score between -1 (negative sentiment) and 1 (positive sentiment). This score is determined by systematically comparing stories that are often classified as having a good or negative financial or economic impact through manual assessments by experts. Based on this human-curated lexicon the ESS algorithm can examine a wide range of sentiment proxies that are frequently mentioned in financial news allowing it to classify the sentiment from earnings reports to natural disasters. (Hafez et al., 2020)

3 METHOD

This section introduces the CrudeBERT model, which extends FinBERT by incorporating knowledge of an event's expected impact on crude oil supply and demand. Section 3.1 presents an overview of the used news headlines and crude oil price data sets which is followed by a discussion of the data pre-processing steps. Afterward, we analyse the shortcomings and flaws of FinBERT and addressed them by developing the CrudeBERT model.

3.1 Datasets

3.1.1 News Data

The dataset containing news information consists of around 46,000 headlines published between 1 January 2000, and 1 April 2021, with high relevance to the topic of crude oil and obtained through the *RavenPack Realtime news Discovery* platform. Similar to Li et al., we limit our analysis to news headlines, since they are more easily accessible, and have lower requirements in terms of pre-processing, storage, and computational power (Li et al., 2019). The headlines used in this paper originate from 1034 unique news sources, of which the majority has been published on the *Dow Jones newswires* (approx. 21,200), followed by *Reuters* (approx. 3,000), *Bloomberg* (approx. 1,100), and *Platts* (approx. 870). There are also around 400 news sources present that only delivered a single headline. It should be noted that RavenPack has added new publishers over the years, which led to a steady increase in the number of available sources over time, especially till 2012. To ensure rich news coverage with diverse sources, we limit our evaluations to the period after 2012.

3.1.2 Price Data

The oil market is dominated by the two most prevalent grades Brent Crude and Western Texas Intermediate (WTI), which dictate the price of crude oil (EIA, 2021). Brent crude is the benchmark for crude oil in Africa, Europe, and the Middle East, accounting for almost two-thirds of the global supply. WTI, on the other hand, is the favored benchmark used by the United States of America. Since all of the headlines in our dataset are in English, the WTI futures prices were regarded as potentially more relevant for our research. The historical values of WTI have been acquired from the financial market platform *investing.com* for the same period as the headlines.

3.2 Data Pre-Processing

3.2.1 Sentiment Data Normalization

We normalized the sentiment values of headlines, computed by the sentiment classifiers, by using z-statistics with the aim to integrate the market's relative mood into the classification model. By normalizing sentiment data over a sliding window we account for the perfect market theory (i.e., the market price reflects all publicly available information) by assuming that only new information that either disappoints or

excels stakeholder expectations results in significant price changes.

Equation 1 outlines the normalization of the sentiment value at time point t based on a weekly sliding window of size $w = 5$ with

$$sent_{norm,t} = \frac{sent_t - \overline{sent}_{t,w}}{\sigma_{t,w}} \quad (1)$$

where $\overline{sent}_{t,w}$ indicates the average sentiment at time points $t, t-1, \dots, t-w$ within the sliding window, and $\sigma_{t,w}$ the corresponding standard deviation.

3.2.2 Price Data Normalization

Due to market volatility, commodity and stock prices show random fluctuations that overlap short-term and long-term trends within the market. Therefore, we also normalize price data using z-statistics to better distinguish between significant market movements and random fluctuations. As with the $sent_{norm,t}$ we normalize price data for a weekly sliding window of $w = 5$ (due to the market being closed over the weekends) as outlined in the equation below:

$$price_{norm,t} = \frac{price_t - \overline{price}_{t,w}}{\sigma_{t,w}} \quad (2)$$

with $\overline{price}_{t,w}$ and $\sigma_{t,w}$ indicating the average price and standard deviations within the chosen sliding window.

3.2.3 Handling Multiple Daily Sentiment Scores

Days covered by multiple news headlines yield multiple sentiment scores, which need to be merged for that given day. Prior work by Hafez et al. concluded that using the sum rather than the mean can examine both the sentiment score and the sentiment volume at the same time, even though, the normalization of scores would severely shrink the relative impact of days with a lower news volume. According to their research, this strategy resulted in superior outcomes in their experiments (Hafez et al., 2018).

3.2.4 Handling Gaps in the Dataset

Rows containing gaps caused either by missing headlines or missing prices (due to closings of the market) were dropped entirely as a row. Furthermore, all the values have been scaled between -1 and 1 . The final dataset covers the period from 1 January 2012 to 1 April 2021 and yields 3376 rows of data.

The dataset aligns the summarized daily sentiment scores with the price change of the following day ($Return_{t+1}$), i.e., assumes that markets will adapt to new information by the next day at the latest. Sentiment scores vary between positive (1) and negative

(−1) values. The price, in contrast, always remains positive with the notable exception of 20 April 2020 when prices became negative for a short period. Price changes (i.e., *Returns*) are, therefore, better suited for indicating the market’s reaction to news coverage. We compute the daily *Returns* of WTI crude oil futures as outlined in Equation 3 and compare them to the sentiment scores.

$$Return = \frac{Price_t - Price_{t-1}}{Price_{t-1}} \quad (3)$$

Interpreting the oil price as the result of cumulative returns allows a comparison to the cumulative sentiment scores, as illustrated in Figure 7.

3.3 Shortcomings of the FinBERT

Figure 7 performs a visual comparison of FinBERT’s cumulative sentiment scores (red) and the price (blue) history of WTI crude oil futures to assess FinBERT’s forecasting potential. The plot does not show any apparent relationship or trends and outlines the need for an additional inquiry into the underlying causes of this poor relationship and potential ways for correcting it.

Adam Smith’s (1776) price theory advocates that the price of limited resources such as crude oil is determined by supply and demand. In this context, *supply* refers to the amount of a product or service that a provider will sell at a given price during a specific period. *Demand* denotes the amount of a product or service that a buyer is willing to acquire during the same period for a given price. The interaction between suppliers and customers yields a competitive market in which the price of products and services is determined by the equilibrium between supply and demand (Smith, 1776). For example, if demand remains constant but supply falls, the resulting shortage will cause prices to rise. A shortage can also occur if the supply remains constant but the demand rises.

In contrast, increased supply with constant demand will result in a surplus and consequently a decrease in prices. A surplus can also emerge if supply remains constant but demand falls. This logic behind supply and demand can be summed up as follows:

Less supply → shortage → higher price
 More supply → surplus → lower price
 Less demand → surplus → lower price
 More demand → shortage → higher price

A drill-down analysis that compared news headlines to FinBERT sentiment scores revealed that FinBERT tended to produce dubious outcomes. Given that crude oil is a publicly-traded asset and FinBERT has been trained on general financial market news this

result seems arguably surprising. Having said that, according to Xing et al. such behavior is expected when utilizing general sentiment analysis methods for a specific domain and is known as the domain adaptation problem (Xing et al., 2020). Weichselbraun et al. (2022) also emphasize the need for domain-specific affective models and present methods for creating such models.

Interpreting the FinBERT scores of news headlines listed in Table 1 based on the impact of supply and demand on prices reveals some serious issues with FinBERT’s assessment of strongly positive (+1), highly negative (-1), and neutral (0) events. Headlines suggesting a drop in supply (e.g., due to accidents at oil refineries and oil platforms), tend to ensue negative FinBERT scores although the corresponding events likely lead to higher crude oil prices. The FinBERT model probably returns these negative scores since accidents are rarely good news in finance and due to moral assessments derived from the human-made annotations within the *Financial Phrase Bank*.

Headlines implying a rise in supply (e.g., due to oil discoveries and increasing exports), in contrast, frequently yield neutral FinBERT scores.

When it comes to a decline in demand (e.g., induced by decreased imports), the resulting surplus should lead to a price decrease. This assessment is also confirmed by two of the three FinBERT scores for the analyzed headlines (row *supply surplus* due to *decreasing demand*) in Table 1. The first headline, in contrast, yields a positive FinBERT score, since FinBERT is not able to correctly interpret the fall in imports indicated by negative numbers such as -16.0%. This limitation should be taken under consideration when utilizing FinBERT since a substantial number of headlines do contain such values. Lastly, headlines that indicate increasing demand should result in higher oil prices. The experiments with FinBERT confirm that it considers headlines conveying increased demand mostly as positive.

3.4 Extending FinBERT to CrudeBERT

In the next step, we extended FinBERT to CrudeBERT to consider the economic law of supply and demand in the model’s assessment.

3.4.1 Training Dataset Generation

To generate a domain-specific labeled silver standard for CrudeBERT, we analyzed several hundred headlines to determine frequently recurring topics, keywords indicating these topics, and their likely impact on the supply and demand of crude oil. This process identified the following major topics: *accidents, oil*

Table 1: Sample of headlines and output of FinBERT.

		Headlines	Sentiment Score Expected	Sentiment Score FinBERT
Shortage	Supply Decrease	Major Explosion, Fire at Oil Refinery in Southeast Philadelphia	Positive	-0.886292
		PETROLEOS confirms Gulf of Mexico oil platform accident	Positive	-0.507213
		CASUALTIES FEARED AT OIL ACCIDENT NEAR IRANS BORDER	Positive	-0.901763
	Demand Increase	EIA Chief expects Global Oil Demand Growth 1 M B/D to 2011	Positive	0.930822
		Turkey Jan-Oct Crude Imports +98.5% To 57.9M MT	Positive	0.866315
		China's crude oil imports up 78.30% in February 2019	Positive	0.922963
Surplus	Demand Decrease	China February Crude Imports -16.0% On Year	Negative	0.540711
		Turkey May Crude Imports down 11.0% On Year	Negative	-0.965965
		Japan June Crude Oil Imports decrease 10.9% On Yr	Negative	-0.955271
	Supply Increase	Iran's Feb Oil Exports +20.9% On Mo at 1.56M B/D - Official	Negative	0.139093
		Apache announces large petroleum discovery in Philadelphia	Negative	0.089624
		Turkey finds oil near Syria, Iraq border	Negative	0.076210

discoveries, changes in exports, changes in imports, changes in demand, pricing, supply, pipeline limitations, drilling, and spillage.

We then queried the RavenPack repository for headlines containing the identified keywords and assigned them to the corresponding topic. The headline in Figure 3, for instance, was assigned the topic *changes in imports* due to the occurrence of the word *import* in the headline. Afterward, we determined the direction of the change by classifying the headline's polarity based on the presence of terms that indicate an increase, a decrease, or constant levels.



Figure 3: Example of Detected Topic and Polarity.

The described approach enabled us to provide topic and direction labels for around 30,000 headlines. Table 2 summarizes the ten detected frequently reoccurring topics and their corresponding frequencies (we do not report the number of headlines with overlapping topics since it has been negligibly small):

Table 2: Frequency of Reoccurring Topics in the Domain of Crude Oil.

Supply change			Demand change		
Increase ca. 5900	No change ca. 350	Decrease ca. 5700	Increase ca. 1300	No change ca. 50	Decrease ca. 800
Export change			Import change		
Increase ca. 2000	No change ca. 150	Decrease ca. 1500	Increase ca. 2800	No change ca. 50	Decrease ca. 2300
Price change			Spill	Discovery	Drilling
Increase ca. 1600	Decrease ca. 1300		ca. 2300	ca. 1600	ca. 100
			Accident	Pipeline issue	
			ca. 400	ca. 100	

Assessing these labels based on the price theory of supply and demand, allowed the creation of a domain-

specific silver standard that classifies the headlines into positive (i.e., indicating increasing crude oil prices), negative (i.e., likely to cause decreasing crude oil prices), and neutral (i.e., should not affect the crude oil price), as outlined below (Figure 4):

- *Lower Prices (score: -1)*: Headlines covering events such as drilling, discovery, increased exports, or simply a rise in oil production are likely to cause an increase in supply. Similarly, headlines stating that oil imports or consumptions are decreasing should, in principle, result in a surplus of oil and, therefore, lower prices.
- *Higher Prices (score: +1)*: Headlines announcing accidents, pipeline constraints, oil spills, or a direct decline in oil supply, in turn, indicate a possible oil shortage due to the negative impact of these events on supply. Likely shortages can also be inferred from news indicating a rise in demand, an increase in imports, or a drop in exports. Generally, news that signals a scarcity of oil or a price increase should have a positive impact on the price.
- *No Price Changes (score: 0)*: A neutral score has been assigned to the relatively small number of headlines that report no signs of supply, demand, imports, or exports.

3.4.2 Model Fine-Tuning

The labeled headlines with the corresponding domain-specific sentiment scores yielded the S&D-Dataset which contains approximately 14,000 negative, 500 neutral, and 15,000 positive headlines.

OIL PRICE DECREASE	OIL PRICE SAME	OIL PRICE INCREASE
«Price decrease»	«Supply steady»	«Price increase»
«Supply increase»	«Demand steady»	«Supply decrease»
«Demand decrease»	«Exports steady»	«Demand increase»
«Exports increase»	«Imports steady»	«Exports decrease»
«Imports decrease»		«Imports increase»
«Oil discovery»		«Spills»
«Drilling»		«Pipeline constraint»
		«Accident»
Number of headlines: ca. 14'000	Number of headlines: ca. 500	Number of headlines: ca. 15'000

Figure 4: Assignment of the Labeled Topics.

We split the dataset into training (60%), test (20%), and validation (20%) partitions (keeping the distribution across classes), and used the test dataset for fine-tuning FinBERT resulting in the CrudeBERT classifier (Figure 5):

Despite the relatively low number of neutral headlines, we included them in training to provide the neural network with examples of lower domain-specific sentiment scores that have not been assigned to one

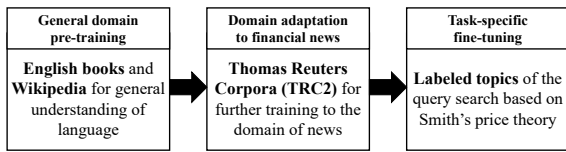


Figure 5: Process of Fine-tuning FinBERT to CrudeBERT.

of the two extremes (i.e., +1 for positive and -1 for negative news).

A preliminary evaluation of the CrudeBERT classifier on the silver standard dataset yielded, despite the class imbalance, a macro F1 score of 0.97, a macro accuracy of 0.98, and a macro recall of 0.97 (Figure 6). On the other hand, the same evaluation with the FinBERT classifier yielded a macro F1 score of 0.29, a macro accuracy of 0.59, and a macro recall of 0.32 on the silver test dataset (Figure 6). Given the substantial amount of headlines used for fine-tuning and their relatively short length (on average 10.4 words per headline), these improvements are not surprising.

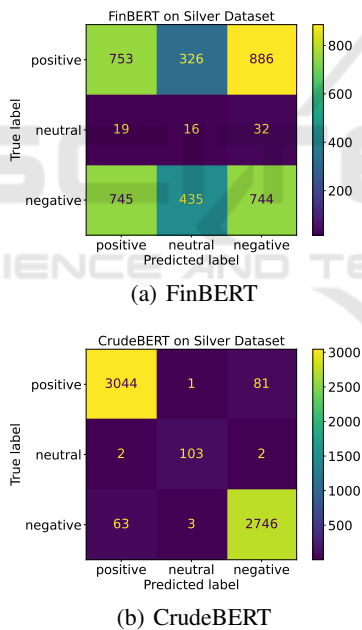


Figure 6: Confusion Matrices of the Two Transformer-based Financial Sentiment Classifiers on the Silver Dataset.

The qualitative comparison in Figure 7 further supports our initial intuition that FinBERT’s lack of knowledge of an event’s impact on supply and demand seriously limits its suitability for prediction tasks. Consequently, it fails to track historical price movements compared to the fine-tuned CrudeBERT model and the commercial classifier of RavenPack.

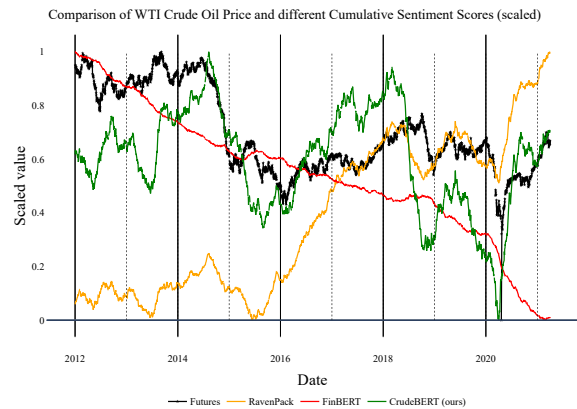


Figure 7: Comparison of WTI Crude Oil Futures Prices and Different Cumulative Sentiment Scores (Scaled).

4 EVALUATION

The following experiments leverage three different sentiment classifiers (FinBERT, CrudeBERT, and RavenPack ESS) to assess the potential of analyzing headlines for predicting the direction of the next day’s ($Return_{t+1}$) change in crude oil futures prices, using a two-class higher/lower price classification schema.

The evaluation considers the period between 1 January 2012 and 1 April 2021 consisting of 3376 days’ worth of data. We use precision, recall, and the F1 metric to assess the predictive potential of the evaluated classifiers.

Table 3 summarizes the evaluation results. On average CrudeBERT outperforms FinBERT, RavenPack, and a random baseline for binary classification. Applying FinBERT without any customizations to the prediction task seems to be contra-productive since it yields worse results than the random baseline. Fine-tuning FinBERT with the presented domain adaptation method considerably improves the method’s performance. CrudeBERT’s overall predictions also surpass the results from RavenPack’s proprietary sentiment classifier, although these differences are less pronounced. CrudeBERT performs slightly worse for price-up predictions but considerably better at predicting pre-down movements.

Figure 8 presents a confusion matrix that compares the predicted label for each classifier with the following day’s price changes of WTI crude oil futures ($Return_{t+1}$).

We, therefore, drew upon the SciPy¹ stats package to perform Pearson’s chi-square test to determine whether the improvements provided by CrudeBERT are statistically significant. When compared to ei-

¹<https://scipy.org>

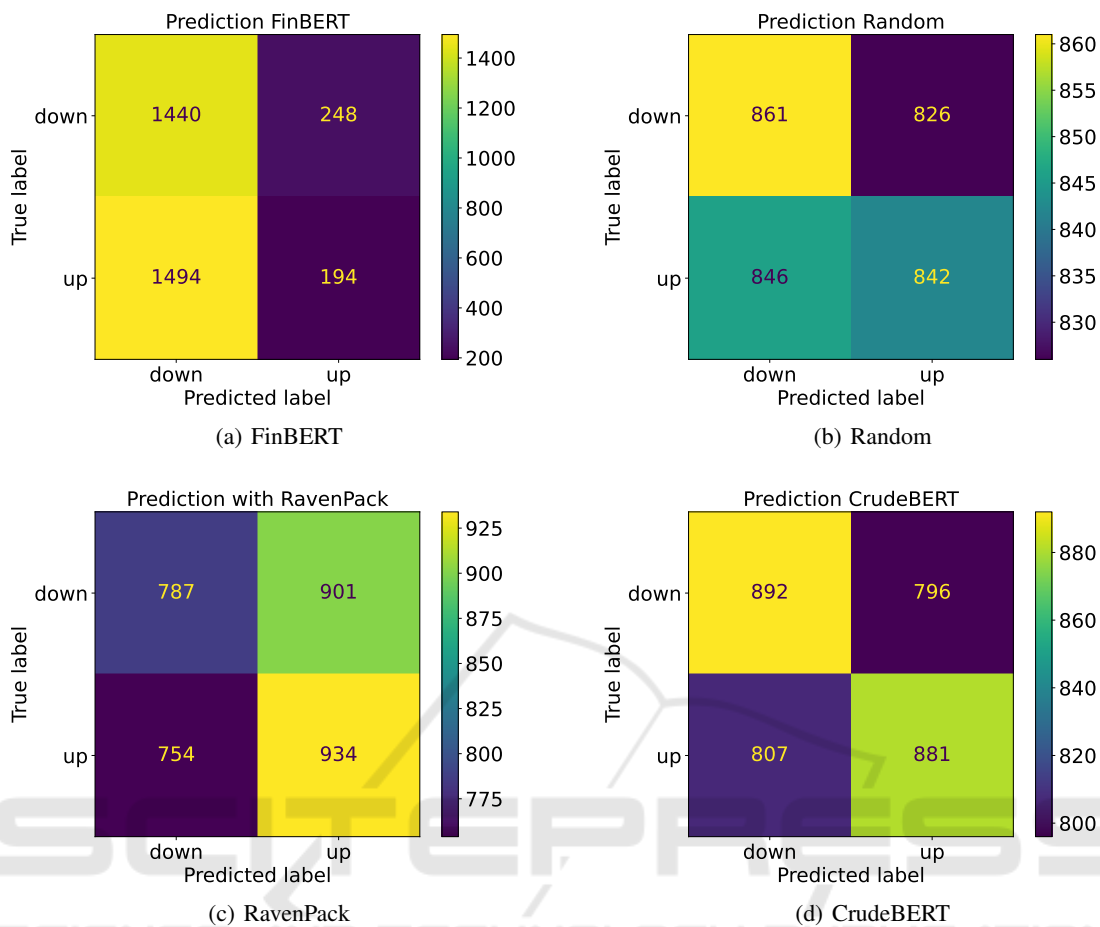


Figure 8: Confusion Matrices Comparing the Price Changes Predicted by RavenPack, FinBERT, and CrudeBERT with the Recorded Price Changes of the Following Day WTI Crude Oil Futures.

Table 3: Classification Report of Different Sentiment Classifiers for Predicting Following Day WTI Oil Futures.

Metric	Category	FinBERT	Random	RavenPack	CrudeBERT
Precision	Price down	0.49	0.51	0.51	0.53
	Price up	0.44	0.50	0.51	0.53
	Macro	0.46	0.51	0.51	0.53
Recall	Price down	0.85	0.51	0.47	0.53
	Price up	0.11	0.50	0.55	0.52
	Macro	0.48	0.51	0.51	0.53
F1-Score	Price down	0.62	0.51	0.49	0.53
	Price up	0.18	0.50	0.53	0.52
	Macro	0.40	0.51	0.51	0.53

ther FinBERT or the random baseline, both CrudeBERT and RavenPack yield significantly better results at the 0.05 significance level. The improvements from RavenPack to CrudeBERT (1721 versus 1773 correct predictions) are less substantial and have only been judged significant at the 0.10 significance level.

5 OUTLOOK AND CONCLUSION

Predicting market movements based on news headlines is still a very challenging task, as outlined in the experiments conducted in Section 4. Even FinBERT, a state-of-the-art sentiment classifier that contains knowledge about the general financial domain, is unable to offer helpful insights into the future price fluctuations of commodities like crude oil when used without any domain adaptations.

The presented paper, therefore, introduces a method for fine-tuning FinBERT based on news headlines. Our approach selects frequently reoccurring topics that cover events illustrating fundamental market dynamics such as the interplay between supply and demand. A frequency analysis identifies these topics which are then used as keywords in search queries for collecting additional suitable headlines. Classifying the retrieved headlines based on Adam

Smith's price theory allows the creation of a silver standard dataset, which serves as a practical and cost-effective alternative to human-curated training datasets. Applying this method to the domain of crude oil led to the creation of a silver standard that has then been used for fine-tuning FinBERT to create CrudeBERT, a domain-specific affective model that provides significantly better results than the original transformer model. In our experiments, which cover crude oil futures price movements over a nine-year period, CrudeBERT outperforms FinBERT and a random baseline on a significance level of 0.05. CrudeBERT even yields better results than RavenPack's proprietary sentiment analysis model which has been optimized in years of development, although the observed improvements are only significant on the 0.10 significance level.

Future research on evaluation methods and metrics will help to better understand the relationship between the model's predictions and future crude oil prices. The presented experiments only shed light upon its short-term prediction performance (i.e., $Return_{t+1}$ which covers the next business day). Thus, further research is required to investigate CrudeBERT's suitability for long-term strategies and in different economic environments (e.g., during times of economic boom or recession).

It is also noteworthy that news headlines alone rather than the whole article seem to be sufficient for providing insights into the likely direction of price changes. Despite the presented improvements, CrudeBERT still has limitations and will be subject to further developments. We also intend to provide CrudeBERT with the ability to distinguish named entities (e.g. countries and oil companies) and numerical clues (e.g. *increased by 10%* and *increased by 1%*) to obtain a more fine-grained sentiment score. This improved indicator should no longer be limited to providing information on the direction of price movements but also express their valence. Considering news volume seems to be another strategy for assessing an event's impact on the market.

Future research will also address the silver standard generation process. The current process, for instance, does not contain any additional logic for handling headlines with contradictory information on future supply and demand (e.g., "Ivory Coast Jan Crude Oil Exports -1 % On Yr, Imports -3 %"). We, therefore, plan to develop strategies for identifying and processing such mixed-signal news headlines.

Furthermore, we aim to assess the feasibility of extending the presented method to other commodities such as perishable (e.g., coffee beans), non-perishable (e.g., natural gas), precious (e.g., gold), and non-

precious (e.g., iron ore) commodities, where pricing may be influenced by similar factors.

ACKNOWLEDGMENT

We would like to extend our gratitude to Prof Dr Hans Wernher van de Venn and the Institute of Mechatronic Systems at Zurich University of Applied Sciences for their generous support of this research. In addition, we would like to thank Dr Adrian M.P. Braşoveanu for his valuable inputs on suitable evaluations for the CrudeBERT model. We would also like to thank Dr Sahand Haji Ali Ahmad for his assessment of the relevance of the news categories used in developing the silver standard dataset.

REFERENCES

- Araci, D. (2019). FinBERT: Financial Sentiment Analysis with Pre-trained Language Models. *arXiv:1908.10063 [cs]*. arXiv: 1908.10063.
- Baboshkin, P. and Uandykova, M. (2021). Multi-source Model of Heterogeneous Data Analysis for Oil Price Forecasting. *International Journal of Energy Economics and Policy*, 11(2):384–391.
- Bahdanau, D., Cho, K., and Bengio, Y. (2016). Neural Machine Translation by Jointly Learning to Align and Translate. *arXiv:1409.0473 [cs, stat]*. arXiv: 1409.0473.
- Brandt, M. W. and Gao, L. (2019). Macro fundamentals or geopolitical events? A textual analysis of news events for crude oil. *Journal of Empirical Finance*, 51:64–94.
- Buyuksahin, B. and Harris, J. (2011). Do Speculators Drive Crude Oil Futures Prices? *The Energy Journal*, Volume 32(Number 2):167–202.
- Chollet, F. (2018). *Deep learning with Python*. Manning Publications Co, Shelter Island, New York. OCLC: ocn982650571.
- Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K. (2018). BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. *arXiv:1810.04805 [cs]*. arXiv: 1810.04805.
- EIA, U. E. I. A. (2021). Table Definitions, Sources, and Explanatory Notes.
- Fama, E. F. (1970). Efficient Capital Markets: A Review of Theory and Empirical Work. *The Journal of Finance*, 25(2):383.
- Hafez, P., Matas, R., Grinis, I., Gomez, F., Kangrga, M., and Liu, A. (2020). Factor Investing With Sentiment: A Look at Asia-Pacific Markets. White Paper.
- Hafez, P., Matas, R., Lautizi, F., A. Guerrero-Colón, J., Gómez, M., and Gómez, F. (2018). Effects of Event Sentiment Aggregation: Sum vs. Mean. White Paper, RavenPack.

- Hamilton, J. (2008). Understanding Crude Oil Prices. Technical Report w14492, National Bureau of Economic Research, Cambridge, MA.
- Hovy, E. H. (2015). What are Sentiment, Affect, and Emotion? Applying the Methodology of Michael Zock to Sentiment Analysis. In Gala, N., Rapp, R., and Bel-Enguix, G., editors, *Language Production, Cognition, and the Lexicon*, pages 13–24. Springer International Publishing, Cham.
- Jiang, H., He, P., Chen, W., Liu, X., Gao, J., and Zhao, T. (2020). SMART: Robust and Efficient Fine-Tuning for Pre-trained Natural Language Models through Principled Regularized Optimization. *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 2177–2190. arXiv: 1911.03437.
- Li, X., Shang, W., and Wang, S. (2019). Text-based crude oil price forecasting: A deep learning approach. *International Journal of Forecasting*, 35(4):1548–1560.
- Li, X., Xie, H., Chen, L., Wang, J., and Deng, X. (2014). News impact on stock price return via sentiment analysis. *Knowledge-Based Systems*, 69:14–23.
- Liew, J. S. Y. (2016). *Fine-grained Emotion Detection in Microblog Text*. PhD thesis.
- Loughran, T. and McDonald, B. (2011). When Is a Liability Not a Liability? Textual Analysis, Dictionaries, and 10-Ks. *The Journal of Finance*, 66(1):35–65.
- Loughran, T. and McDonald, B. (2016). Textual analysis in accounting and finance: A survey. *Journal of Accounting Research*, 54(4):1187–1230.
- Malkiel, B. G. (1989). Efficient market hypothesis. In *Finance*, pages 127–134. Springer.
- Malo, P., Sinha, A., Korhonen, P., Wallenius, J., and Takala, P. (2014). Good debt or bad debt: Detecting semantic orientations in economic texts: Good Debt or Bad Debt. *Journal of the Association for Information Science and Technology*, 65(4):782–796.
- McCarthy, R. V., McCarthy, M. M., Ceccucci, W., Halawi, L., and SpringerLink (Online service) (2019). *Applying Predictive Analytics Finding Value in Data*. OCLC: 1204071994.
- Mikolov, T., Chen, K., Corrado, G., and Dean, J. (2013). Efficient estimation of word representations in vector space.
- Mohammad, S. M. (2021). Sentiment Analysis: Automatically Detecting Valence, Emotions, and Other Affective States from Text. Number: arXiv:2005.11882 arXiv:2005.11882 [cs].
- Pennington, J., Socher, R., and Manning, C. (2014). Glove: Global Vectors for Word Representation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1532–1543, Doha, Qatar. Association for Computational Linguistics.
- Plutchik, R. (1982). A psychoevolutionary theory of emotions. *Social Science Information*.
- Qian, B. and Rasheed, K. (2007). Stock market prediction with multiple classifiers. *Applied Intelligence*, 26(1):25–33. Publisher: Springer.
- Smith, A. (1776). *An Inquiry into the Nature and Causes of the Wealth of Nations*. McMaster University Archive for the History of Economic Thought.
- Susanto, Y., Livingstone, A., Ng, B. C., and Cambria, E. (2020). The Hourglass model revisited. *IEEE Intelligent Systems*, 35(5).
- Taboada, M. (2016). Sentiment analysis: An overview from linguistics. *Annual Review of Linguistics*, 2(1):325–347.
- Tang, D., Qin, B., Feng, X., and Liu, T. (2016). Effective LSTMs for target-dependent sentiment classification. In *Proceedings of COLING 2016, the 26th International Conference on Computational Linguistics: Technical Papers*, pages 3298–3307, Osaka, Japan. The COLING 2016 Organizing Committee.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L., and Polosukhin, I. (2017). Attention Is All You Need. *arXiv:1706.03762 [cs]*. arXiv: 1706.03762.
- Weichselbraun, A., Steixner, J., Brasoveanu, A. M. P., Scharl, A., Göbel, M., and Nixon, L. J. B. (2022). Automatic Expansion of Domain-Specific Affective Models for Web Intelligence Applications. *Cognitive Computation*, 14(1):228–245.
- Wex, F., Widder, N., Liebmann, M., and Neumann, D. (2013). Early Warning of Impending Oil Crises Using the Predictive Power of Online News Stories. In *2013 46th Hawaii International Conference on System Sciences*, pages 1512–1521, Wailea, HI, USA. IEEE.
- Xing, F., Malandri, L., Zhang, Y., and Cambria, E. (2020). Financial Sentiment Analysis: An Investigation into Common Mistakes and Silver Bullets. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 978–987, Barcelona, Spain (Online). International Committee on Computational Linguistics.
- Yenicelik, K. D. (2020). *Understanding and Exploiting Subspace Organization in Contextual Word Embeddings*. Masterthese, Eidgenössische Technische Hochschule Zürich, Zürich 8006, Schweiz.
- Zhu, Y., Kiros, R., Zemel, R., Salakhutdinov, R., Urtasun, R., Torralba, A., and Fidler, S. (2015). Aligning Books and Movies: Towards Story-Like Visual Explanations by Watching Movies and Reading Books. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*.