# Stock Values and Earnings Call Transcripts: a Dataset Suitable for Sentiment Analysis

Dexter Roozen
Tilburg University
drh.roozen@gmail.com

Francesco Lelli
Tilburg University
f.lelli@tilburguniversity.edu
https://francescolelli.info

#### **Abstract:**

The dataset reports a collection of earnings call transcripts, the related stock prices, and the related sector index. It contains a total of 188 transcripts, 11970 stock prices, and 1196 sector index values. Furthermore, all of these data originated in the period 2016-2020 and are related to the NASDAQ stock market.

The data have been collected using Yahoo Finance and Thomson Reuters Eikon. Specifically, Yahoo Finance offered daily stock prices and traded volume. At the same time, Thomson Reuters Eikon has been used as source for the earnings call transcripts.

The dataset can be used as a benchmark for the evaluation of several NLP techniques as well as machine learning algorithms for understanding their potential for financial applications. Moreover, it is also possible to expand the dataset by extending the period in which the data originated following a similar procedure.

Keywords: dataset, stock, sentiment analysis, nlp, Nasdaq, stock prices

## 1. Introduction:

The dataset contains earning call transcripts and related stock process of 10 popular stocks of the Nasdaq index. These data have been collected for performing a set of "bag of words" analysis in order to evaluate possible correlation between them and stock process.

Several researches try to use yahoo finance Data in order to gather insight in stock prices and volume of trading. This is the case of Rao et all in [2] and authors in [3]. Some of the authors access directly to API. However, in recent times yahoo has made several changes and this practice is becoming more difficult.

Earning call transcript have been collected via Thomson Reuters Eikon. Several authors has use this approach. This is the case of the works presented in [4],[5],[6] and [7] However, the exact procedure for attaining the transcripts is not specified within these papers.

With this preprint we intend to formalize the procedure for collecting the data as well as release an initial version of a dataset that can be extended and reused in order to perform several different analysis.

The table below summarize the data:



Subject	Economics/Linguistics/Computer Science
Specific subject area	The specific subject area of this research is Sentiment Analysis. Sentiment analysis is a natural language processing (NLP) technique to determine the sentiment (positive or negative) behind data.  To elaborate, NLP is a field of research that investigates the ability of computers to understand and manipulate natural languages, such as English.  A crucial step of textual sentiment analysis is to pre-process the text documents. This pre-processing phase consists of multiple 'pre-processing techniques' of which the effects were studied.
Type of data	Table Text
How data were acquired	The stock values and sector index were acquired through Yahoo Finance (website). The earnings call transcripts were acquired through Thomson Reuters Eikon (software).
Data format	Raw
Parameters for data collection	The related companies of the stock values and earnings call transcripts were chosen based on the condition of being NASDAQ listed. Furthermore, the date range for the stock values and earnings call transcripts is 2016-2020.
Description of data collection	The stock values were acquired by using Yahoo Finance. Yahoo Finance provides news, information, commentary, and reports on the subject of finance. This website lets users search for specific companies with its search bar. When entering a company such as "Apple Inc.", the website will direct the user to a summary of general financial information about the company. Besides this 'Summary' tab, there are also other tabs providing different sorts of information. Selecting the 'Historical Data' tab shows the historical stock values of the searched company. Additionally, it also lets users specify the date period and the frequency with which it shows the stock values. Afterward, the presented stock values can be downloaded as a Microsoft Excel Comma Separated Value (CSV) File.  Thomson Reuters Eikon helped acquire the earnings call transcripts. This software provides users with many different sorts of financial

	information. Selecting the 'advanced event search' directs the user to financial information about particular events. Specifying the event type 'Earnings Conference Call' shows information about earnings calls of many different companies. Selecting 'transcript' from the 'Content Type' selector filters out earnings calls without transcripts. Lastly, date preference and company should be specified to find the desired information. With the save batch icon it is possible to download 100 transcripts at a time as text documents.
Data source location	https://doi.org/10.34894/TJE0D0
Data accessibility	Repository name: DataverseNL Data identification number: N.A Direct URL to data: <a href="https://doi.org/10.34894/TJE0D0">https://doi.org/10.34894/TJE0D0</a> Instructions for accessing these data: Data are open access

## 2. Methods for data acquisition

All of the stock values and the sector index were acquired by utilizing the Yahoo Finance search bar. Searching for a company such as Apple Inc. results in a summary of financial information about this company. However, the stock values and sector index within the dataset are presented in the "Historical Data" tab. Selecting this tab and specifying the time period January 1<sup>st</sup>, 2016 – October 1<sup>st</sup>, 2020 and selecting "Apply" will show the data presented in this dataset. Lastly, selecting "Download" provides a CSV file containing all of this data.

The earnings call transcripts were acquired through Thomson Reuters Eikon. Selecting the "advanced event search" option shows unfiltered financial information about many different sorts of events. Specifying the event type by selecting "Earnings Conference Call" will filter this information by only showing information about earnings calls. Additionally, selecting "Transcript" from the "Content Type" selector will show only earnings calls that can be provided together with a transcript of the earnings call. Lastly, specifying the company and time period will show a list with the earnings call transcripts contained in this dataset. For efficiency purposes, the save batch icon makes it possible to download this whole list of transcripts.

#### 2.1 List of resources used for the collection of the data

The following sources were used in retrieving the historical stock values of the NASDAQ listed companies and the sector index:

- NASDAQ. (2020). Apple Inc. (AAPL). [Historical stock values, 2016-2020]. Retrieved from https://finance.yahoo.com/quote/AAPL/history?p=AAPL
- NASDAQ. (2020). Advanced Micro Devices, Inc. (AMD). [Historical stock values, 2016-2020]. Retrieved from https://finance.yahoo.com/quote/AMD/history?p=AMD
- NASDAQ. (2020). Amazon.com, Inc. (AMZN). [Historical stock values, 2016-2020]. Retrieved from https://finance.yahoo.com/quote/AMZN/history?p=AMZN
- NASDAQ. (2020). ASML Holding N.V. (ASML). [Historical stock values, 2016-2020]. Retrieved from https://finance.yahoo.com/quote/ASML/history?p=ASML
- NASDAQ. (2020). Cisco Systems, Inc. (CSCO). [Historical stock values, 2016-2020]. Retrieved from https://finance.yahoo.com/guote/CSCO/history?p=CSCO
- NASDAQ. (2020). Alphabet Inc. (GOOGL). [Historical stock values, 2016-2020]. Retrieved from https://finance.yahoo.com/quote/GOOGL/history?p=GOOGL
- NASDAQ. (2020). Intel Corporation (INTC). [Historical stock values, 2016-2020]. Retrieved from https://finance.yahoo.com/quote/INTC/history?p=INTC
- NASDAQ. (2020). Microsoft Corporation (MSFT). [Historical stock values, 2016-2020]. Retrieved from https://finance.yahoo.com/quote/MSFT/history?p=MSFT
- NASDAQ. (2020). Micron Technology, Inc. (MU). [Historical stock values, 2016-2020]. Retrieved from https://finance.yahoo.com/quote/MU/history?p=MU
- NASDAQ. (2020). NVIDIA Corporation (NVDA). [Historical stock values, 2016-2020]. Retrieved from https://finance.yahoo.com/quote/NVDA/history?p=NVDA
- NASDAQ. (2020). NASDAQ Composite (^IXIC). [Historical stock values, 2016-2020]. Retrieved from https://finance.yahoo.com/quote/%5EIXIC/history?p=%5EIXIC

The following database provided all of the earnings call transcripts from the selected NASDAQ companies:

• Thomson Reuters Eikon. (2020). [Earnings call transcripts, 2016-2020]. Available at: Thomson Reuters (Accessed: November 19 2020).

# 3. Results: Data Description

The folder named "Stock Values and Sector Index" in the dataset contains all of the CSV files that were acquired through the before mentioned method. These files consist of individual tables for each NASDAQ Company and the NASDAQ sector index. The folder is structured as portrayed in table 1.

Company		Size	Source	
	Type			
Apple Inc. (AAPL)	CSV	1197 rows x 1 column	NASDAQ (2020)	
Advanced Micro Devices, Inc. (AMD)	CSV	1197 rows x 1 column	NASDAQ (2020)	
Amazon.com, Inc. (AMZN)	CSV	1197 rows x 1 column	NASDAQ (2020)	
ASML Holding N.V. (ASML)	CSV	1197 rows x 1 column	NASDAQ (2020)	
Cisco Systems, Inc. (CSCO)	CSV	1197 rows x 1 column	NASDAQ (2020)	
Alphabet Inc. (GOOGL)	CSV	1197 rows x 1 column	NASDAQ (2020)	
Intel Corporation (INTC)	CSV	1197 rows x 1 column	NASDAQ (2020)	
Microsoft Corporation (MSFT)	CSV	1197 rows x 1 column	NASDAQ (2020)	
Micron Technology, Inc. (MU)	CSV	1197 rows x 1 column	NASDAQ (2020)	
NVIDIA Corporation (NVDA)	CSV	1197 rows x 1 column	NASDAQ (2020)	
NASDAQ Composite (^IXIC)	CSV	1196 rows x 1 column	NASDAQ (2020)	

**Table 1**: Folder structure stock values and sector index

Accessing these files can be done through Microsoft Excel. A snippet of what these files look like when opened in Excel is portrayed in figure 1.

	Α	В	С	D	Е	F	G	Н
1	Date,Open,H	igh,Low,	Close,Adj C	lose,Volum	е			
2	2016-01-04,2	25.65250	0,26.34250	1,25.50000	00,26.33750	0,24.44303	37,2705976	00
3	2016-01-05,2	6.43750	0,26.46250	0,25.60250	1,25.67750	0,23.83051	13,2231640	00
4	2016-01-06,2	25.13999	9,25.59250	1,24.96750	) <b>1,25.17</b> 499	9,23.36416	51,2738296	00
5	2016-01-07,2	4.67000	0,25.03249	9,24.10750	00,24.11249	9,22.37808	32,3243776	00
6	2016-01-08,2	4.63750	1,24.77750	00,24.19000	1,24.24000	0,22.49640	08,2831920	00
7	2016-01-11,2	4.74250	0,24.76499	9,24.33499	99,24.63250	0,22.86067	76,1989576	00
8	2016-01-12,2	25.13750	1,25.17250	1,24.70999	99,24.99000	0,23.19246	55,1966168	00
9	2016-01-13,2	25.08000	0,25.29750	1,24.32500	1,24.34750	0,22.59617	78,2497584	00
10	2016-01-14,2	4.49000	0,25.12000	1,23.93499	99,24.87999	9,23.09037	70,2526804	00
11	2016-01-15,2	4.04999	9,24.42750	00,23.84000	00,24.28249	9,22.53585	58,3193356	00
12	2016-01-19,2	4.60250	1,24.66250	0,23.87500	00,24.16500	1,22.42680	09,2123508	00
13	2016-01-20,2	23.77500	0,24.54750	1,23.35500	00,24.19750	0,22.45697	70,2893376	00
14	2016-01-21,2	4.26499	9,24.46999	9,23.73500	1,24.07500	1,22.34328	31,2086460	00
15	2016-01-22,2	4.65749	9,25.36500	0,24.59250	1,25.35500	0,23.53121	10,2632020	00
16	2016-01-25,2	25.37999	9,25.38250	0,24.80250	00,24.86000	1,23.07181	17,2071780	00
17	2016-01-26,2	4.98250	0,25.21999	9,24.51750	00,24.99749	9,23.19942	27,3003080	00
18	2016-01-27,2	4.01000	0,24.15749	9,23.33499	99,23.35500	0,21.67507	74,5334788	00
19	2016-01-28,2	23.44750	0,23.62999	9,23.09750	00,23.52249	9,21.83051	19,2227152	00
20	2016-01-29,2	23.69750	0,24.33499	9,23.58750	00,24.33499	9,22.58457	78,2576660	00
21	2016-02-01,2	4.11750	0,24.17750	00,23.85000	00,24.10750	0,22.37344	12,1637740	00
22	2016-02-02,2	23.85500	0,24.01000	0,23.57000	00,23.62000	1,21.92101	11,1494288	00
23	2016-02-03,2	23.75000	0,24.20999	9,23.52000	00,24.08750	0,22.35488	37,1838572	00
24	2016-02-04,2	23.96500	0,24.33250	0,23.79750	1,24.15000	0,22.53450	04,1858868	00
25	2016-02-05,2	4.12999	9,24.23000	0,23.42250	1,23.50499	9,21.93265	50,1856724	00
26	2016-02-08,2	23.28249	9,23.92499	9,23.26000	0,23.75250	1,22.16359	97,2160856	00
27	2016-02-09,2	23.57250	0,23.98500	1,23.48250	0,23.74749	9,22.15893	32,1773248	00
28	2016-02-10,2	23.98000	0,24.08750	0,23.52500	0,23.56749	9,21.99096	57,1693744	00
29	2016-02-11,2	23.44750	0,23.68000	0,23.14749	99,23.42499	9,21.85800	00,2002988	00

Figure 1: Stock values in Excel

Furthermore, the folder named "Transcripts" in the dataset contains multiple folders named after each company of which the earnings call transcripts were acquired. The structure of these folders is portrayed in table 2.

Company Folder	Contents	Source
Apple Inc. (AAPL)	19 Text Documents	Thomson Reuters Eikon (2020)
Advanced Micro Devices, Inc. (AMD)	19 Text Documents	Thomson Reuters Eikon (2020)
Amazon.com, Inc. (AMZN)	19 Text Documents	Thomson Reuters Eikon (2020)
ASML Holding N.V. (ASML)	19 Text Documents	Thomson Reuters Eikon (2020)
Cisco Systems, Inc. (CSCO)	19 Text Documents	Thomson Reuters Eikon (2020)
Alphabet Inc. (GOOGL)	19 Text Documents	Thomson Reuters Eikon (2020)
Intel Corporation (INTC)	19 Text Documents	Thomson Reuters Eikon (2020)
Microsoft Corporation (MSFT)	19 Text Documents	Thomson Reuters Eikon (2020)
Micron Technology, Inc. (MU)	17 Text Documents	Thomson Reuters Eikon (2020)
NVIDIA Corporation (NVDA)	19 Text Documents	Thomson Reuters Eikon (2020)

 Table 2: Folder structure transcripts

The earnings call transcripts are text documents that are all structured in the same manner. Date, time, participants, and the words that are spoken during the earnings call are all registered within these documents. Accessing these data can be done through standard software such as Notepad on Windows devices. A snippet of what these documents look like in Notepad is portrayed in figure 2.

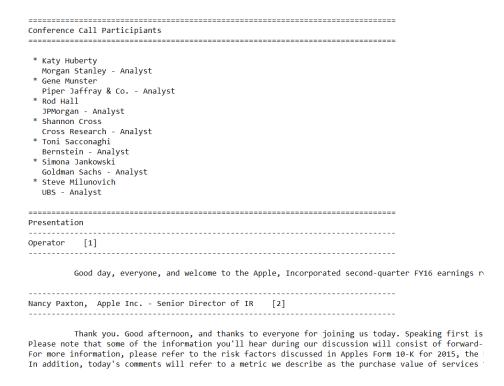


Figure 2: Snippet earnings call transcript

#### 3.1 Determining positive and negative transcripts

Formulas are used to determine whether a transcript is positive or negative. Firstly, the stock ratio formula, which has the following form:

stock ratio = stock value one day after earnings call / stock value n days before earnings call

The stock ratio shows the percentage increase or decrease of the stock value. However, a percentage increase in stock value does not immediately imply that the earnings call is positive as there are other variables to consider. To factor in an additional variable called investor mood, a second formula is defined:

sector ratio = sector value one day after earnings call / sector value n days before earnings call

Sector refers to the NASDAQ composite. The sector ratio is taken into account to consider the mood of the sector index. If the increase in stock ratio turns out to be higher than the increase of the sector ratio, the earnings call can be determined positive. If not, the transcript is deemed negative.

## 4. Discussion

The collected dataset provides the following value:

- These data can prove useful as they may help to further uncover dynamics related to correlational relationships between stock values and earnings call transcripts.
- Furthermore, the data can easily be expanded by i.e. extending the date range. Additionally, the data is easy to use and readable by multiple programming languages.
- Both practitioners at companies as well as scholars can benefit from the use of these data. Every
  company and scholar uses homemade datasets with consequential discrepancies. The adoption
  of a shared dataset for benchmarking analysis will promote a homogeneous evaluation of the
  results.
- The data was used primarily for the application of a limited amount of NLP techniques and machine learning algorithms. Consequently, this dataset offers the possibility to explore different approaches.

## 5. Conclusions

In this preprint we presented a dataset that has been designed for performing sentiment analysis in the stock market. Information regarding daily price and volume has been collected using yahoo finance. At the same, Thomson Reuters has been used for collecting earning transcripts. Details about the procedure has been described and presented in the previous sections. The dataset contains 11970 stock prices, and 1196 sector index values. Furthermore, all of these data originated in the period 2016-2020 and are related to the NASDAQ stock market.

The dataset can be used for developing and benchmarking NLP techniques and machine learning algorithms.

#### **Competing Interests**

The authors declare that they have no known competing financial interests or personal relationships which have or could be perceived to have influenced the work reported in this article.

## 6. References

[1] Roozen, D. R. H. (2021). Correlations between Earnings Call Transcripts and Stock Price: a Sentiment Analysis Approach. Master thesis at Tilburg University

[2] Rao, T., & Srivastava, S. (2012). Analyzing stock market movements using twitter sentiment analysis.

- [3] Kordonis, J., Symeonidis, S., & Arampatzis, A. (2016). Stock price forecasting via sentiment analysis on twitter. In *Proceedings of the 20th Pan-Hellenic Conference on Informatics* (pp. 1-6).
- [4] Jha, V., Blaine, J., & Montague, W. (2015). Finding Value in Earnings Transcripts Data with AlphaSense. Available at extractalpha.com
- [5] Heinrichs, A., Park, J., & Soltes, E. F. (2019). Who consumes firm disclosures? Evidence from earnings conference calls. *The Accounting Review*, *94*(3), 205-231.
- [6] Theil, C. K., Broscheit, S., & Stuckenschmidt, H. (2019). PRoFET: Predicting the Risk of Firms from Event Transcripts. In *IJCAI* (pp. 5211-5217).
- [7] Jenkins, P. (2020, April). Structured Paragraph Embeddings of Financial Earnings Calls. In *Companion Proceedings of the Web Conference 2020* (pp. 264-268).