# How to Talk When a Machine is Listening: Corporate Disclosure in the Age of AI

Sean Cao[a]    Wei Jiang[b]    Baozhong Yang[a]    Alan Zhang[c]

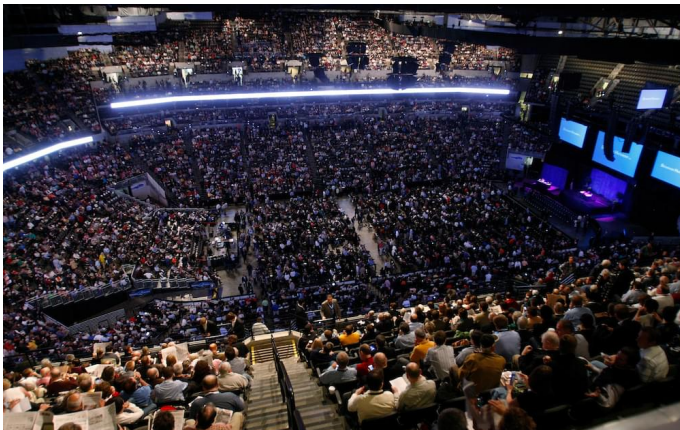[a] J. Mack Robinson College of Business, Georgia State University

[b] Emory University Goizueta Business School, ECGI, and NBER

[c] Florida International University

Prepared for Q Group Research Talks
September 20, 2022

MOTIVATION

- Corporate disclosure communicates financial health, promotes the culture and brand, and engages a full spectrum of stakeholders.
- Warren Buffet's annual letters to shareholders of Berkshire Hathaway showcase Corporate American writing at its best – for human readers.
  - "Be fearful when others are greedy and greedy when others are fearful."
  - "When it's raining gold, reach for a bucket, not a thimble."

## THE CHANGING READERSHIP OF DISCLOSURE



**Artificial Intelligence**

## Robo-surveillance shifts tone of CEO earnings calls

Trading algorithms leave a mark with deeper focus on the spoken word

Hedge funds use natural language processing to scour earnings calls, social media posts and regulatory documents for market-moving clues © FT illustration

**Robin Wigglesworth** in Oslo DECEMBER 5 2020

When Man Group chief executive Luke Ellis discusses his investment company's results with analysts he chooses his words carefully. He knows better than most that the machines are listening.
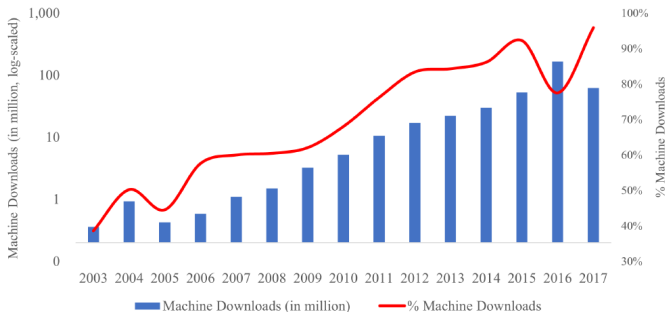
THE CHANGING READERSHIP OF DISCLOSURE

- A substantial amount of buying and selling of shares are triggered by assessment and recommendations made by robots and algorithms.
- Technology makes it feasible: Machine learning and natural language processing kits.
- The sheer volume of regulatory filings makes it inevitable.
    - EDGAR hosts 11 million filings by over 600,000 reporting entities using 478 unique form types. There were 1.5 billion unique requests via SEC.gov in 2016 alone (Bauguess, 2018).
    - The length of 10k increases by five times from 2005 to 2017, and the number of textual changes over previous filings increases by over 12 times (Cohen, Malloy, and Nguyen, 2020).
- The SEC estimates that "as much as 85% of the documents visited are by internet bot" (Bauguess, 2018).
- "It pays to write well" (Hwang and Kim, 2017); but now corporate disclosure needs to resonate with both human and machine readers.

## OBJECTIVES OF THE STUDY

- Research question: Whether and how public companies adjust the way they talk knowing that machines are listening.
  - Quantify and connect expected AI reader base and machine-friendliness of disclosure documents.
  - Identify changes in writing patterns affecting "sentiment" and "tone" after the availability of new algorithms, notably Loughran and McDonald (2011) and BERT (2018).
  - An "out-of-the-sample" test on the machine-assessed voice emotional quality of conference calls.

- Connect and contribute to the growing literature on:
  - Information acquisition and dissemination via downloads of SEC filings.
  - Assessing qualitative information using textual analyses and machine learning.
  - A newly emerging "feedback effect" from machine learning about firm fundamentals to corporate decisions: Encoded rules are at least partially transparent, observable, or reverse-engineerable, agents who are impacted by the decisions thus have the incentive to manipulate the inputs to machine-learning.

# HOW DO WE MEASURE MACHINE READERSHIP: *Machine Downloads*

- The frequency of *Machine Downloads* of corporate filings as an upper bound as well as a proxy for the presence of "machine readers."
- Identify an IP address downloading more than 50 unique firms' filings on a given day, and requests that are attributed to web crawlers in the SEC Log File Data, as a machine (i.e., robot) visitor (Lee, Ma, and Wang, 2015). All remaining requests are labeled as "Other" requests.
- Lag downloads tracking to avoid reserve causality.

## *Machine Downloads*: WHO'S WHO

#### The top machine downloaders by IP address

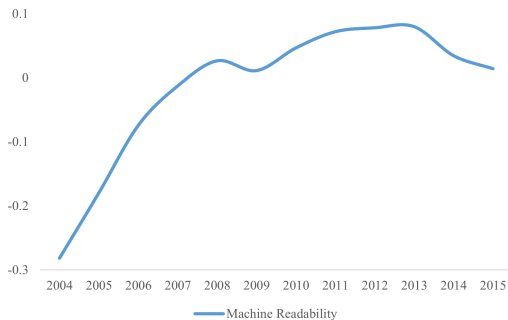| Rank | Name of institution | #MD | Type of institution |
|------|---------------------|-----|---------------------|
| 1 | Renaissance Technologies | 536,753 | Quantitative hedge fund |
| 2 | Two Sigma Investments | 515,255 | Quantitative hedge fund |
| 3 | Barclays Capital | 377,280 | Financial conglomerate with asset management |
| 4 | JPMorgan Chase | 154,475 | Financial conglomerate with asset management |
| 5 | Point72 Asset Management | 104,337 | Quantitative hedge fund |
| 6 | Wells Fargo | 94,261 | Financial conglomerate with asset management |
| 7 | Morgan Stanley | 91,522 | Investment bank with asset management |
| 8 | Citadel LLC | 82,375 | Quantitative hedge fund |
| 9 | RBC Capital Markets | 79,469 | Financial conglomerate with asset management |
| 10 | D. E. Shaw CO. | 67,838 | Quantitative hedge fund |
| 11 | UBS AG | 64,029 | Financial conglomerate with asset management |
| 12 | Deutsche Bank AG | 55,825 | Investment bank with asset management |
| 13 | Union Bank of California | 50,938 | Full service bank with private wealth management |
| 14 | Squarepoint Ops | 48,678 | Quantitative hedge fund |
| 15 | Jefferies Group | 47,926 | Investment bank with asset management |
| 16 | Stifel, Nicolaus Company | 24,759 | Investment bank with asset management |
| 17 | Piper Jaffray | 18,604 | Investment bank with asset management |
| 18 | Lazard | 18,290 | Investment bank with asset management |
| 19 | Oppenheimer Co. | 15,203 | Investment bank with asset management |
| 20 | Northern Trust Corporation | 11,916 | Financial conglomerate with asset management |

ALTERNATIVE MACHINE-READERSHIP MEASURES

- *AI ownership*
  - The percentage of shares outstanding held by investment companies with AI capabilities.
  - Identify AI-equipped investment company if it has AI-related job postings in the past five years according to data from Burning Glass.
- *AI talent supply*
  - Approximate AI talent supply available to institutional investors based on state-year level proportion of working-age population with IT degrees where the investors are headquartered.
  - Headquarters were mostly chosen before the AI era and hence unlikely to be affected by omitted variables or reverse causality.
- Both measures aggregate investor ownership using 13F.

# How do we measure *Machine Readability*?

- Measures the ease at which a filing can be "understood," i.e., processed and parsed, by an automated program.
- Following Allee et al. (2018): The ease of (i) separating tables from text; (ii) extracting numbers from text; (iii) identifying the information contained in the table; (iv) inclusion of all needed information without relying on external exhibits; and (iv) proportion of characters that are standard ASCII characters.
- The average of the five standardized component statistics.

## *Machine Readability* & *Machine Downloads* ARE POSITIVELY RELATED: IN LEVELS AND "UPGRADES"

Sample: All annual and quarterly filings by publicly listed firms from 2003-2016.
"Upgrades:" A one-standard deviation increase in machine readability over the previous year.

| Dependent Variable | (1) | (2) | (3) | (4) | (5) | (6) |
|---|---|---|---|---|---|---|
| | | *Machine Readability* | | | *MR Upgrade* | |
| | | | | | | |
| *Machine Downloads* | 0.076*** | 0.075*** | 0.060*** | 0.078*** | | |
| | (13.89) | (17.45) | (10.33) | (15.93) | | |
| Δ*Machine Downloads* | | | | | 0.005*** | 0.006*** |
| | | | | | (2.90) | (3.40) |
| *Other Downloads* | 0.005 | 0.002 | -0.007 | -0.006 | 0.000 | -0.001 |
| | (1.15) | (0.47) | (-1.44) | (-1.33) | (0.20) | (-0.44) |
| | | | | | | |
| Observations | 198,358 | 199,241 | 150,425 | 150,346 | 135,146 | 135,068 |
| R-squared | 0.082 | 0.363 | 0.084 | 0.357 | 0.025 | 0.144 |
| Control Variables | Yes | Yes | Yes | Yes | Yes | Yes |
| Firm FE | No | Yes | No | Yes | No | Yes |
| Industry FE | Yes | No | Yes | No | Yes | No |
| Year FE | Yes | Yes | Yes | Yes | Yes | Yes |

## SAME RELATION HOLDS WITH INVESTOR AI CAPABILITY MEASURES

| Dependent Variable | (1) | (2) | (3) | (4) |
|---|---|---|---|---|
| | | *Machine Readability* | | |
| | | | | |
| *AI Ownership* | 0.515*** | 0.356*** | | |
| | (8.06) | (8.29) | | |
| *AI Talent Supply* | | | 0.160*** | 0.192** |
| | | | (3.09) | (2.29) |
| | | | | |
| Observations | 50,747 | 50,608 | 70,969 | 70,912 |
| R-squared | 0.093 | 0.373 | 0.088 | 0.361 |
| Control Variables | Yes | Yes | Yes | Yes |
| Firm FE | No | Yes | No | Yes |
| Industry FE | Yes | No | Yes | No |
| Year FE | Yes | Yes | Yes | Yes |

## Machines speed up information dissemination

- Resort to TAQ high-frequency data to track down "time to trade" and "time to directional trade" from filing posting.
- A one standard deviation increase in machine downloads is associated with 7-12 seconds faster in the first trade, or 10-15 seconds faster in the first directional trade.
- The above effect is significantly stronger when interacted with machine readability.
- When machine downloads are high, bid-ask spread enlarges more right after the posting of filings: Machines are creating new information asymmetry based on *public* information.

| Dependent Variable | *Bid-Ask Spread* | | |
|---|---|---|---|
| *Machine Downloads* $\times$ *After* | 0.028*** | 0.063*** | 0.055*** |
| | (3.11) | (7.24) | (8.46) |
| *Machine Downloads* | 0.993*** | 0.877*** | |
| | (49.59) | (36.07) | |
| | | | |
| Observations | 2,328,247 | 2,328,190 | 2,673,992 |
| R-squared | 0.116 | 0.232 | 0.720 |
| Control Variables | Yes | Yes | Subsumed |
| Company FE | No | Yes | No |
| Filing FE | No | No | Yes |
| Minute FE | Yes | Yes | Yes |

## WRITING FOR MACHINES

- *Machine readability* is about format for easy processing. How about communicating the content?
- Will a machine understand:
  "The period for a new election of a citizen to administer the executive government of the United States being not far distant, and the time actually arrived when your thoughts must be employed in designating the person who is to be clothed with that important trust, it appears to me proper, especially as it may conduce to a more distinct expression of the public voice, that I should now apprise you of the resolution I have formed, to decline being considered among the number of those out of whom a choice is to be made." (George Washington's 1796 Farewell Address)
- Will a machine misinterpret "The concern about a possible delay is unwarranted."
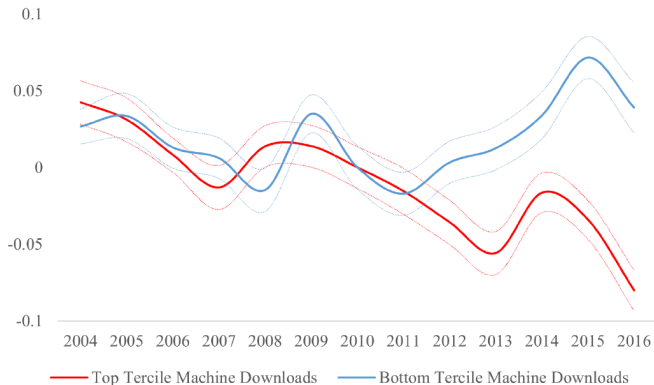- Be mindful when job seekers compose resumes!

## MANAGE SENTIMENT WITH HUMAN AND MACHINE READERS

- Representation of "positive" and especially "negative" words had been based on the Harvard Psychosociological Dictionary which provides predictive information about firm outcomes and stock returns.
- Loughran and McDonald (2011) presented a specialized dictionary of positive/negative and tone words that fits the unique financial text, which has been feeding into algorithms.
- "Sentiment" is defined as the representation of "negative" words in the documents.

| Dependent Variable | (1) | (2) | (3) | (4) | (5) | (6) |
|---|---|---|---|---|---|---|
| | LM – Harvard Sentiment | | LM Sentiment | | Harvard Sentiment | |
| | | | | | | |
| Machine Downloads | -0.072*** | -0.079*** | -0.062*** | -0.050*** | 0.010 | 0.029*** |
| × Post | (-6.95) | (-8.94) | (-4.98) | (-4.99) | (0.76) | (2.65) |
| Machine Downloads | -0.007 | -0.011** | -0.009 | -0.019*** | -0.002 | -0.008 |
| | (-1.17) | (-2.46) | (-1.18) | (-3.72) | (-0.23) | (-1.43) |
| | | | | | | |
| Observations | 158,578 | 158,515 | 158,578 | 158,515 | 158,578 | 158,515 |
| R-squared | 0.217 | 0.568 | 0.241 | 0.632 | 0.208 | 0.590 |
| Control Variables | Yes | Yes | Yes | Yes | Yes | Yes |
| Firm FE | No | Yes | No | Yes | No | Yes |
| Industry FE | Yes | No | Yes | No | Yes | No |
| Year FE | Yes | Yes | Yes | Yes | Yes | Yes |

## Parallel pre-trends of $LM$ - $Harvard$

Firms with high expected machine downloads differentially avoid LM-negative words relative to the Harvard-negative words, only after the publication of LM (2011), the exact timing of which is quasi-random.

## WORDS WITH THE BIGGEST REDUCTION IN FREQUENCY AFTER 2011

| Neg. Word | Change in Frequency (in percentage points) | Neg. Word | % Reduction in Frequency |
|-----------|---------------------------------------------|-----------|--------------------------|
| restructuring | -0.35% | correction | -37.6% |
| termination | -0.34% | destroyed | -34.5% |
| restatement | -0.25% | restatement | -32.5% |
| declined | -0.25% | declined | -20.6% |
| correction | -0.21% | purported | -20.0% |
| misstatement | -0.21% | encumbrance | -19.2% |
| terminated | -0.16% | counterclaim | -18.4% |
| late | -0.16% | misstatement | -18.0% |
| alleged | -0.15% | writeoff | -17.5% |
| omitted | -0.15% | closure | -17.0% |

## OTHER TONES DEVELOPED IN LM (2011)

- *Litigious* words (such as "claimant" and "tort") reflect a litigious environment.
- *Uncertain* words (such as "approximate" and "contingency") capture a general notion of imprecision.
- *Weak Modal* (such as "possibly" and "could") and *Strong Modal* (such as "always" and "must") words convey levels of confidence.
- Measured as the ratio of each category of words to the length of the filing.
- A high level of each of the four tones predicts one or more of negative outcomes: More "material weakness," fraud, and law suits; and is met with lower short-term stock return.
- Do managers avoid these words after the dictionary became publicly known?

## TONE FOR MACHINES

Firms avoid all four categories of tone words significantly more after the public knowledge of their impact.

| Dependent Variable | Litigious | | Uncertainty | | Weak Modal | | Strong Modal | |
|---|---|---|---|---|---|---|---|---|
| Machine Downloads × Post | -0.056*** | -0.057*** | -0.016** | -0.021*** | -0.028*** | -0.034*** | -0.008*** | -0.007*** |
| | (-5.38) | (-6.02) | (-2.01) | (-3.49) | (-4.85) | (-8.86) | (-4.39) | (-4.39) |
| Machine Downloads | 0.011* | 0.007 | -0.006 | -0.009*** | -0.018*** | -0.021*** | -0.003** | -0.004*** |
| | (1.71) | (1.44) | (-1.33) | (-3.05) | (-5.39) | (-10.05) | (-2.19) | (-4.98) |
| | | | | | | | | |
| Observations | 158,578 | 158,515 | 158,578 | 158,515 | 158,578 | 158,515 | 158,578 | 158,515 |
| R-squared | 0.188 | 0.509 | 0.196 | 0.6 | 0.238 | 0.624 | 0.277 | 0.571 |
| Controls included | Yes | Yes | Yes | Yes | Yes | Yes | Yes | Yes |
| Company FE | No | Yes | No | Yes | No | Yes | No | Yes |
| Industry FE | Yes | No | Yes | No | Yes | No | Yes | No |
| Year FE | Yes | Yes | Yes | Yes | Yes | Yes | Yes | Yes |

## Managing sentiment in response to BERT

Sample includes all annual and quarterly filings between 2016 and 2019.
*BERT Sentiment* is defined as the number of "negative" sentences, scaled by the total number of sentences.

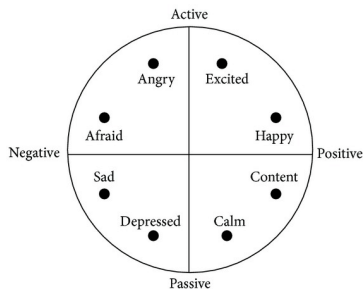| Dependent Variable | (1) | (2) | (3) | (4) |
|---|---|---|---|---|
| | \multicolumn | *BERT Sentiment* | | |
| | NegSent/TotalSent | | NegSent/TotalWords | |
| | | | | |
| *AI Ownership × Post-BERT* | -4.276** | | -0.190** | |
| | (-2.13) | | (-2.37) | |
| *AI Ownership* | 2.025 | | 0.096 | |
| | (1.08) | | (1.27) | |
| *AI Talent Supply × Post-BERT* | | -0.983*** | | -0.041*** |
| | | (-3.61) | | (-3.98) |
| *AI Talent Supply* | | -0.522 | | -0.010 |
| | | (-1.18) | | (-0.65) |
| | | | | |
| Observations | 6,399 | 6,627 | 6,399 | 6,627 |
| R-squared | 0.795 | 0.796 | 0.803 | 0.804 |
| Control Variables | Yes | Yes | Yes | Yes |
| Firm FE | Yes | Yes | Yes | Yes |
| Year FE | Yes | Yes | Yes | Yes |

## Managing audio tones: The real "talk"

- Starting around 2008, voice analytic software (e.g., Layered Voice Analysis (LVA)) has gained popularity among investors looking for an edge in information processing.
- Such software has enabled researchers to study the vocal expressions of managers and their implications on capital markets.
- Is there a feedback effect to how managers talk? Learn from a sample of 43,462 earnings-related conference call speeches from 3,290 unique companies during 2010–2016.
- Two key measures based on the existent literature:
  Emotional *Valence* and *Arousal* correspond to positivity and excitedness of voices.
- Hu and Song (2020) showed that venture capitalists are more likely to invest in start-ups whose founders give pitches that are rated high in either and both.

## More Details of Voice Analytics

- Open-source pre-trained Python machine learning package *pyAudioAnalysis* (Giannakopoulos, 2015) to extract emotional measures from earnings calls.
- *Emotion Valence*: the extent to which an emotion is positive or negative, with a larger value indicating greater positivity.
- *Emotion Arousal*: the intensity or the strength of the associated emotion state, and a greater (lower) value suggests that the speaker is more excited (calmer).
- Both measures are bounded between –1 and 1.

Valence and Arousal in a 2D Cartesian Coordinates system

## HOW TO TALK TO MACHINES

Managers talk with higher valence and, to a lesser degree, higher arousal, when there are more machines expected in the audience. (Control variables include earnings surprise.)

| Dependent Variable | Emotion Valence | | | Emotion Arousal | | |
|---|---|---|---|---|---|---|
| Machine Downloads | 0.043*** | 0.035*** | 0.042*** | 0.004* | 0.003 | 0.005** |
| | (11.40) | (8.13) | (11.14) | (1.79) | (0.94) | (2.28) |
| Other Downloads | -0.017*** | -0.014*** | -0.017*** | -0.006*** | 0.000 | -0.006*** |
| | (-5.74) | (-4.32) | (-5.67) | (-3.65) | (0.19) | (-3.71) |
| Observations | 43,336 | 41,340 | 41,224 | 43,336 | 41,340 | 41,224 |
| R-squared | 0.389 | 0.189 | 0.383 | 0.395 | 0.132 | 0.395 |
| Controls included | No | Yes | Yes | No | Yes | Yes |
| Company FE | Yes | No | Yes | Yes | No | Yes |
| Industry FE | No | Yes | No | No | Yes | No |
| Year FE | Yes | Yes | Yes | Yes | Yes | Yes |

## CONCLUSION

- Corporate disclosure in writing and speaking has been reshaped by machine readership employed by algorithmic traders and quantitative analysts.
- Increasing AI readership motivates firms to prepare filings that are more friendly to machine parsing and processing.
- Firms adapt sentiment and tone management to evolving algorithms.
- The feedback effect from technology calls for more studies to understand the induced behavior by AI and algorithms in financial markets.