# AI for Unstructured Data in Finance

**Machine Learning in Finance Workshop**
**May 17, 2019**

**Amanda Stent**
**NLP Architect, Office of the CTO**

**Bloomberg**
Engineering

TechAtBloomberg.com

# We Are in an "AI Revolution"

BUSINESS | JOURNAL REPORTS: LEADERSHIP

## The Optimistic Promise of Artificial Intelligence

Andrew Ng and Tong Zhang on how AI is going to be like electricity, transforming every industry

Microsoft researchers achieve speech recognition milestone

China's AI Agenda Advances

THE AI REVOLUTION IS ON

Google Translate Has Reached Human-Like Accuracy Thanks To Neural Machine Translation Engine

## Alibaba's AI Outguns Humans in Reading Test

The Sublime and Scary Future of Cameras With A.I. Brains

Stanford team creates computer vision algorithm that can describe photos
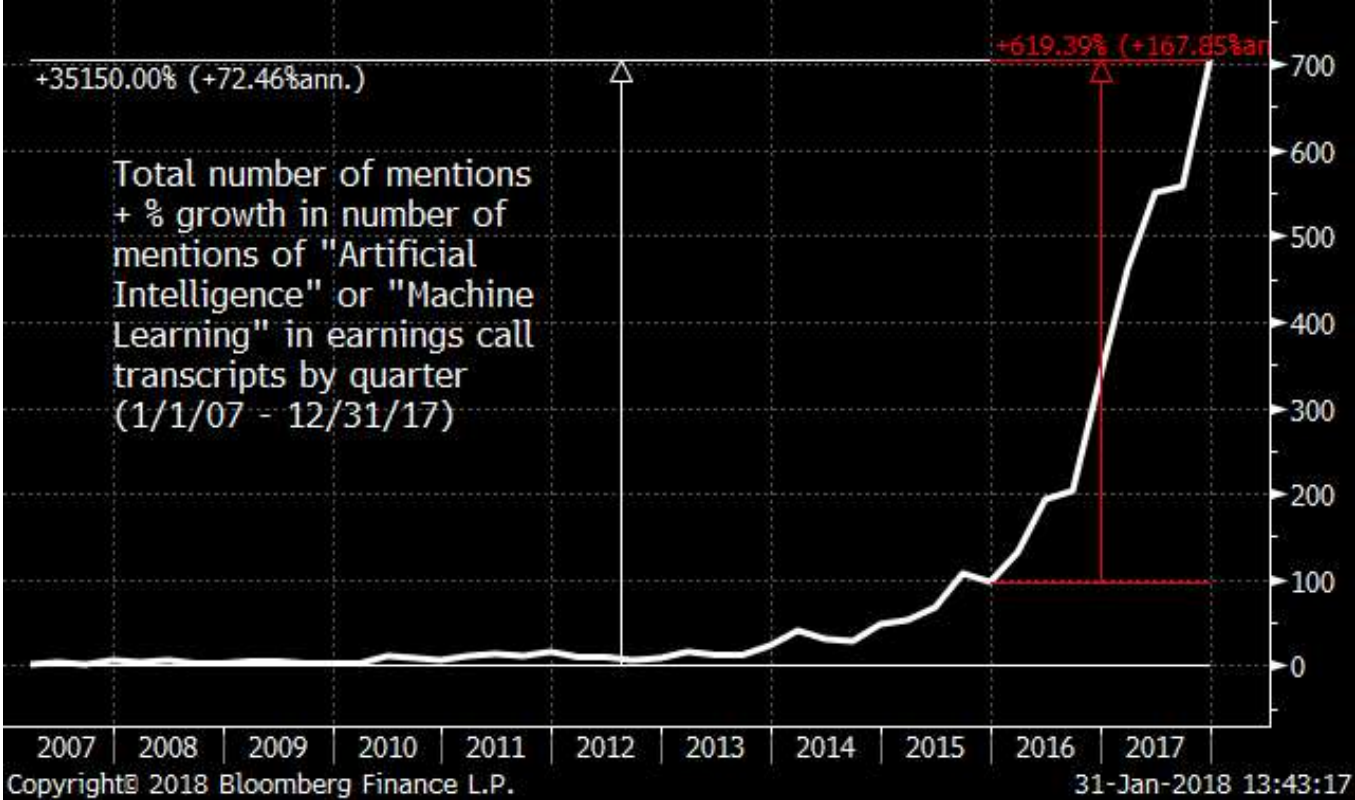
**TechAtBloomberg.com**

**Bloomberg**
Engineering

# The "AI Revolution" is Changing the World



+35150.00% (+72.46%ann.)

+619.39% (+167.85%ar)

Total number of mentions
+ % growth in number of
mentions of "Artificial
Intelligence" or "Machine
Learning" in earnings call
transcripts by quarter
(1/1/07 - 12/31/17)

2007  2008  2009  2010  2011  2012  2013  2014  2015  2016  2017

Copyright© 2018 Bloomberg Finance L.P.

31-Jan-2018 13:43:17

**TechAtBloomberg.com**

**Bloomberg**

Engineering

# The "AI Revolution" is Data Driven

- Deep-learning based AI depends on:
  - Lots of specialized (tailored, expensive) compute
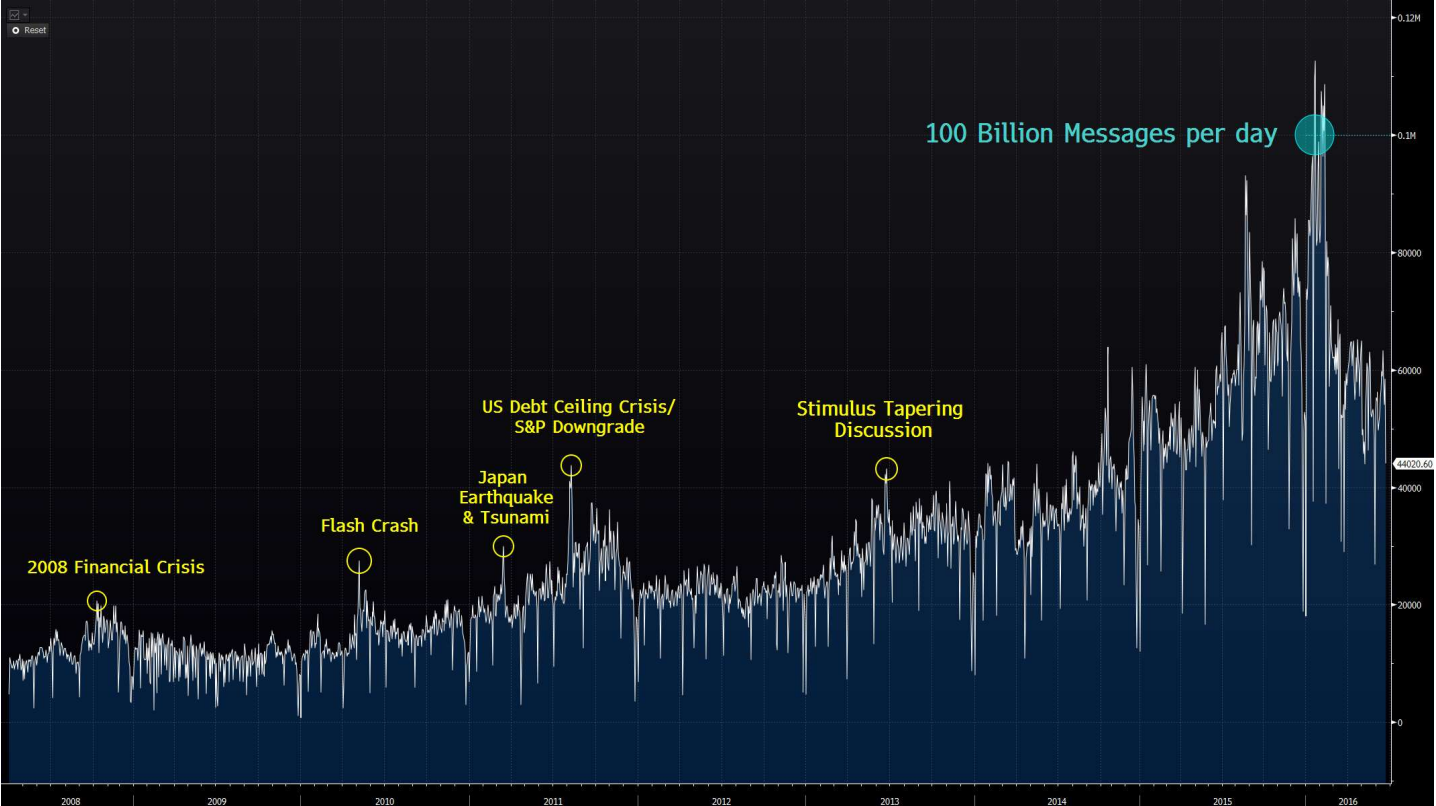  - Lots of labeled data (supplied by you and me!)



*Source: Michael Felsberg*

**Bloomberg**

Engineering

# Bloomberg is the Home of Market Data



100 Billion Messages per day

US Debt Ceiling Crisis/
S&P Downgrade

Stimulus Tapering
Discussion

Japan
Earthquake
& Tsunami

Flash Crash

2008 Financial Crisis

**TechAtBloomberg.com**

**Bloomberg**

Engineering

# Bloomberg Is the Home of Finance News

## 325K+
subscribers

## News search
### 16M queries
per day
Stories available for search in ~100ms
Average query response time <200ms

## Volume
### 2M
stories per day

### Index of
### 650M
stories

### 500
documents ingested
per second

## News alerting
### 1.5M
subscriptions

### 500 stories
matched per second
Alerts delivered in <100ms

**TechAtBloomberg.com**

**Bloomberg**

Engineering

# News Moves Markets (Fast)



First Bloomberg Headline

New York Times Story

SEC Announcement

12%

<20 minutes

Day Session
Last 166.15 -18.12
High on 04/16 10:40 186.41
Average 172.2602
Low on 04/16 10:56 163.55
Prev Close ----- 184.27

10:35  10:40  10:45  10:50  10:55
16 Apr 2010

# Social Media Moves Markets



**Over $136 <u>billion</u> was wiped out in minutes**

# Our Key Challenges

- The need for speed

- The need for precision

- The need for human in the loop

**Bloomberg**

Engineering

# Text Analytics vs. NLP

## Natural Language Processing (computational linguistics)

"the scientific and engineering discipline concerned with understanding written and spoken language from a computational perspective, and building artifacts that usefully process and produce language" (*Stanford Encyclopedia of Philosophy*)

## Text Analytics

"the process of deriving information from text sources" (*Gartner IT Glossary*)

**Bloomberg**

Engineering

# Why Do Text Analytics?

**Most companies that use text analytics have one goal**

- Assess customer feedback via voice, SMS, email… (e.g., to assess NPS)
- NLP tasks: entity recognition, sentiment analysis, key terms, visualization

**A few have the goal of ranking/serving the text**

- Treat text as a "first class object" (i.e., text is both input and output)
- NLP tasks: entity recognition and linking, topic detection, sentiment, clustering

**A few have a different kind of goal**

- Assess market signals (financial reports, news articles, tweets…) in order to predict some (financial) outcome (e.g., movie grosses, stock prices)
- NLP tasks: entity recognition and linking, sentiment analysis, relation/event extraction, signal extraction, classification

**TechAtBloomberg.com**

**Bloomberg**

Engineering

# How Text Analytics Is Done

**Generation 0: write a bunch of rules ("templates", "grammars")**

- High-precision
- Slow, manual, difficult to maintain or update

**Generation 1: train a statistical classifier**

- For sequence tagging: conditional random fields; for document classification: linear/logistic regression, SVMs, decision trees/random forests
- Need labeled data

**Generation 2: deep learning and human in the loop**

- Need a lot of labeled data, or distant supervision
- May be slower

**We still need separate "models" for each language/genre pair**

English news, English social, English generic; Chinese news, Chinese social…

**TechAtBloomberg.com**

**Bloomberg**

Engineering

# The Life of a Document

Ingestion

Content Extraction — Text, tables, figures

Information Extraction — Topics, entities, quantities, events, sentiment

Signal Extraction — Counts of topics, entities, events… over time

Key Insights, News Importance…

Ranking and Recommendation

Analytics

News Trends, Market Moving News…

# Topics are Central to AI for Finance

- Human curated in-house taxonomy of topic codes
- Exposed as an important metadata of text documents

**Bloomberg**

Engineering

# Topics & Entities Convert Unstructured Content to Structured Data

- Over 1B tweets processed in the past 12 months
- Over 700 topic models deployed (still growing)
- Over 1M tweets labelled in-house for model training
- Peak-rate capacity of the real-time classification system



2016 U.S. Presidential Election
- **~100M** tweets in 4 days
- peak rate at **4K+** tweets/sec

**Bloomberg**

Engineering

# Topics & Entities Link Unstructured Content to the World

**Bloomberg**

Engineering

# Topics + Entities + Interaction = Insight

# (Topics + Entities + Market Fundamentals) / Time = Insight



TREN  News Trends ☆

Set Alert | Set View as Default

| Companies | Security List | **Topics** | | Trend Period **8 Hours** ▾ |
|---|---|---|---|---|

| News Reader Activity | News Sentiment | Twitter Sentiment | News Volume | **Twitter Volume** |
|---|---|---|---|---|

**Largest Increase** | Largest Total

| News Topic | Δ Pub. ↓ | NT | Representative Tweet |
|---|---|---|---|
| 1) Agricultural Chemicals | | | A. Smith: German Beer Industry in Shock Over Monsanto Glyphosate Contamination (deadly mutagenic Roundup) … |
| 2) Fertilizer | | | Livemint: RT @livemint_m2m: Is India turning the corner on usage of fertilizers?Urea sales till 20 March stood at… |
| 3) Iceland | | | Skift: In parts of Iceland, you're more likely to meet a tourist than you are a local. The country doesn't want to … |
| 4) Health Care Facilities & Svcs | | | Richard James: These Women Suffering Complications From Mesh Implants Were Let Down By The Health System, … |
| 5) Hospitals, Clinics | | | Manhattan Institute: .@CPopeHC: What can the NHS learn from the American healthcare system, and vice versa? h… |
| 6) Space Exploration | | | Hindustan Times: Astronomers are puzzled by a 'see-through' galaxy that doesn't have dark matter https://t.co/… |
| 7) Env. Protection Agency | | | NorthJersey.com: EPA and Honeywell on Edgewater Superfund site https://t.co/Z43drJHzPr https://t.co/Dmp6S2L… |
| 8) Pakistan Government | | | Govt of Pakistan: #AlbertEinstein #Quoteoftheday #QOTD #quote #Pakistan https://t.co/u3C72mVCq9 |
| 9) Pakistan Prime Minister | | | RT: Pakistani PM forced to undergo security check at US airport as Trump reportedly mulls sanctions against #Isl… |
| 10) Chile Government | | | NYT National News: The Chilean government said it would start an inquiry into whether the remains of a tiny bab… |
| 11) Korea Economy | | | Bloomberg Markets: #5Things- North Korea and…Japan?- Tech drama- Indonesia's new central bank chief- Upwa… |
| 12) Copyrights | | | ITWeb News Publication: Appeals court sends multibillion-dollar Java copyright case between Oracle and Google b… |
| 13) Hybrid Bonds | | | Christopher Mahoney: "Tesla's ratings reflect significant shortfall in the production rate of the Model 3. Tesla fac… |
| 14) Canada | | | Calgary Herald: Egg Week: The absolutely official, completely objective, not biased at all, definitive ranking of w… |
| 15) East African Community | | | The African Voice: US President Donald Trump nominates Illinois Senator Kyle McCarter to serve as Ambassador o… |

**TechAtBloomberg.com**

**Bloomberg**

Engineering

# (Extracted Signals + Market Fundamentals) / Time = Insight

# (Extracted Signals + Market Fundamentals) / Time = Insight

**Bloomberg**

Engineering

# Bloomberg in a Nutshell



Ingestion + Creation

Analyses

Discovery

# The Need for Precision



FORTUNE | How Bill O'Reilly Leaving Fox Fired Up O'Reilly Auto Parts Stock

accessories and tools, has no relation to the host of Fox's now-canceled *O'Reilly Factor*. But without any other news that day to explain the move in the retailer's stock price, some investors had a different theory: Computerized algorithms that trade stocks based on Twitter and social media alone had picked up on a surge in posts about Bill O'Reilly, and interpreted it as a signal to buy O'Reilly stock.

Bloomberg

Engineering

# Human in the Loop

- In industry, no model is static
  - New entities
  - New contexts
  - New events and relationships

- Humans in the loop can:
  - Ameliorate lack of recall, precision
  - Provide important training data

- To enable hybrid workflow, more work is needed on:
  - Models "knowing what they know"
  - Active learning
  - Tools for humans to work with ML algorithms

**Bloomberg**

Engineering

# Going Beyond Topics, Entities and Sentiment



Before → Call → After

Before: {Bearish, neutral, bullish} / Price target

After: {Bearish, neutral, bullish} / Price target

**Bloomberg**

Engineering

# Going Beyond Topics, Entities and Sentiment

**Brian Nowak, Analyst:** Thanks for taking my questions. One on YouTube, I guess. Could you just talk to some of the qualitative drivers that are really bringing more advertising dollars on to YouTube? And then I think last quarter you had mentioned the top 100 advertiser spending was up 60% year-on-year on YouTube, wondering, if you could update us on that? And the second one on search, it sounds like mobile is accelerating. Where are you now in the mobile versus desktop monetization gap? And, Sundar, how do you think about that long-term? Do you see mobile being higher, reaching equilibrium? How do you see that trending?

**Sundar Pichai, CEO:** On the YouTube one. Look, I mean, the shift to video is a profound medium shift and especially in the context of mobile, you know and obviously users are following that. You're seeing it in YouTube as well as elsewhere in mobile. And so, advertisers are being increasingly conscious. They're being very, very responsive. So, we're seeing great traction there and we'll continue to see that. They are moving more off their traditional budgets to YouTube and that's where we are getting traction. On mobile search, to me, increasingly we see we already announced that over 50% of our searches are on mobile. Mobile gives us very unique opportunities in terms of better understanding users and over time, as we use things like machine learning, I think we can make great strides. So, my long-term view on this is, it is as-compelling or in fact even better than desktop, but it will take us time to get there. We're going to be focused till we get there.

Table 1: Earnings calls are extremely complex examples of naturally-occurring discourse. In this example question-answer pair from a Google earnings call on October 27, 2016, the analyst asks **six distinct questions** in a single turn. Because the interaction originates as speech, there are **discourse markers and hedging.** The analyst and executive discuss **concrete entities and performance statistics** and **past, present and future** performance.

Bloomberg

Engineering

# Going Beyond Topics, Entities and Sentiment

- Bullish analysts are called on earlier in earnings calls

- Bullish analysts ask more positive questions; neutral/bearish analysts ask more negative questions

- Bullish analysts ask more concrete questions

- Neutral/bearish analysts ask about the past more

- Earnings calls are moderately predictive of changes in analysts' forecasts

# Some Recent Work

- Meerkamp, P. & Zhou, Z. (2017) "Information extraction with character-level neural networks and free noisy supervision." arXiv 1612.04118.

- Yang, Y., Irsoy, O., & Rahman, S. (2018) *"Collective entity disambiguation with structured gradient tree boosting. In Proc. NAACL.*

- Keith, K. & Stent, A. (2019) "Modeling financial analysts' decision making via the pragmatics and semantics of earnings calls." In *Proc. ACL.*

**TechAtBloomberg.com**

# Data Science at Bloomberg

- **Our data:** We have the most text labeled with financial (*and legal!*) entities, instruments and topics anywhere

- **Our databases:** We have the most information about financial entities, instruments and topics anywhere

- **Our communications:** We connect the most extraordinary individuals across global finance

- **Our challenges:** Our business involves *every* AI challenge

**TechAtBloomberg.com**

**Bloomberg**

Engineering

# Where to Learn More

- **TechAtBloomberg.com**
  —Papers published in 2018 at EMNLP, SIGIR, NAACL HLT, WWW, and more…
  —Data for Good Exchange, Machines + Media, and other exciting events
  —Our work with JupyterHub, Solr, and other open source frameworks
  —Foundations of Machine Learning (FOML), David Rosenberg's ML class

- **GitHub.com/Bloomberg**
  —bqplot: fast interactive plotting
  —sgtb: structured gradient tree boosting
  —foml: Foundations of Machine Learning
  —cnn-rnf: Convolutional Neural Networks with Recurrent Neural Filters
  —scatteract: extraction of labeled quantities from scatter plots

**TechAtBloomberg.com**

**Bloomberg**
Engineering

# Thank You!

## Questions?

**TechAtBloomberg.com**

**Bloomberg**
Engineering