# 10-K / 10-Q Text Analysis

## Eric He – eh1885 – Group 21

# Project Objective

1. Test pipeline and data processing capabilities by systematically downloading years of 10-K data from the SEC EDGAR database.
   a. Bonus: also download 10-Q data and special filings.
2. Can compare changes in text filings Y/Y using text distance and sentiment analysis
   a. Bonus: plug into ChatGPT to highlight and summarize differences

# Concepts

1. Systemic data download - can we parallel process or stream process?
2. Heavy text processing - can we Cythonize?
3. Text data analysis

# Background - EDGAR database

1.  EDGAR is SEC's database for storing corporate filings. The data can be downloaded for free by anyone in the world (but if you do it wrong you can be throttled).
2.  The raw data consists of raw HTML code, and requires a fairly intricate cleaning process to split 10-K or 10-Q into component sections.
3.  Most of the financial statement is boilerplate, and large sections of the statement may not change Q/Q or even Y/Y. However, other times important financial data is presented in tables which are also challenging to process.
4.  Downloading all the data, uncompressed, is several tens of gigabytes

# Example of a 10-K filing

```
<SEC-DOCUMENT>0001628280-16-020309.txt : 20161026
<SEC-HEADER>0001628280-16-020309.hdr.sgml : 20161026
<ACCEPTANCE-DATETIME>20161026164216
ACCESSION NUMBER:            0001628280-16-020309
CONFORMED SUBMISSION TYPE:   10-K
PUBLIC DOCUMENT COUNT:       96
CONFORMED PERIOD OF REPORT:  20160924
FILED AS OF DATE:            20161026
DATE AS OF CHANGE:           20161026

FILER:

        COMPANY DATA:
                COMPANY CONFORMED NAME:           APPLE INC
                CENTRAL INDEX KEY:                0000320193
                STANDARD INDUSTRIAL CLASSIFICATION:  ELECTRONIC COMPUTERS [3571]
                IRS NUMBER:                       942404110
                STATE OF INCORPORATION:           CA
                FISCAL YEAR END:                  0924

        FILING VALUES:
                FORM TYPE:        10-K
                SEC ACT:          1934 Act
                SEC FILE NUMBER:  001-36743
                FILM NUMBER:      161953070

        BUSINESS ADDRESS:
                STREET 1:         ONE INFINITE LOOP
                CITY:             CUPERTINO
                STATE:            CA
                ZIP:              95014
                BUSINESS PHONE:   (408) 996-1010

        MAIL ADDRESS:
                STREET 1:         ONE INFINITE LOOP
                CITY:             CUPERTINO
                STATE:            CA
                ZIP:              95014

        FORMER COMPANY:
                FORMER CONFORMED NAME:  APPLE COMPUTER INC
                DATE OF NAME CHANGE:    19970808
</SEC-HEADER>
<DOCUMENT>
<TYPE>10-K
<SEQUENCE>1
<FILENAME>a201610-k9242016.htm
<DESCRIPTION>10-K
<TEXT>
<!DOCTYPE html PUBLIC "-//W3C//DTD HTML 4.01 Transitional//EN" "http://www.w3.org/TR/html4/loose.dtd"
<html>
        <head>
                <!-- Document created using Wdesk 1 -->
                <!-- Copyright 2016 Workiva -->
                <title>Document</title>
```

# What can a corporate filing tell us?

1. Management will highlight headwinds or other important updates in filings.
2. One might wonder if markets are truly reacting to the information in the filing "instantaneously" or if this market information is incorporated over a longer period of time, e.g. days, weeks or even months. Potential signals:
   a. Sentiment analysis
   b. Text vs. Numeric data
   c. Distance analysis
3. Specific language analysis may be relatively uncorrelated to other common information sources such as correlated company returns or economic data

I am not super interested in comparing to stock returns, I would rather work directly with the text data itself and treat this as a data processing exercise. 6 years ago I downloaded this data with a Google Chrome extension over the course of 4 weeks and felt I could do better.

# Project Roadmap

1. [03/15/23] Proof of Concept Download database
2. [03/22/23] Proof of Concept process filings
3. [03/29/23] Proof of Concept Cosine Distance + Sentiment Analysis
4. [04/15/23] Parallelize data download with multiprocessing, compress, load into text database
5. [04/22/23] Cythonized compute of text features and condensed text summarization algos
6. [04/29/23] Summarize performance of optimization strategies
7. [Extra] Integrate with ChatGPT API to return opinions on textual differences or sentiment
8. [Extra] Do it again but for 10-Qs

# Thank You

## Questions?