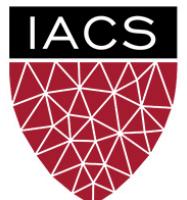


News Sentiment Analysis for Stock Return Prediction

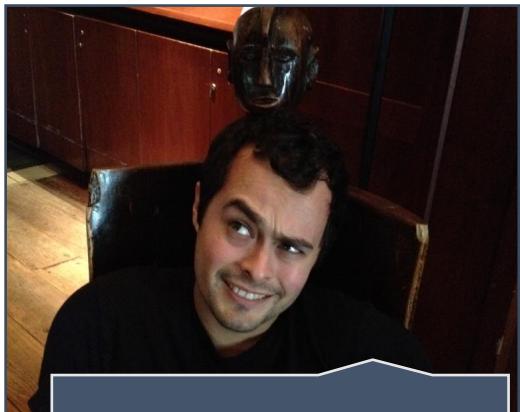
αβnormal Distribution

AC295/CS115

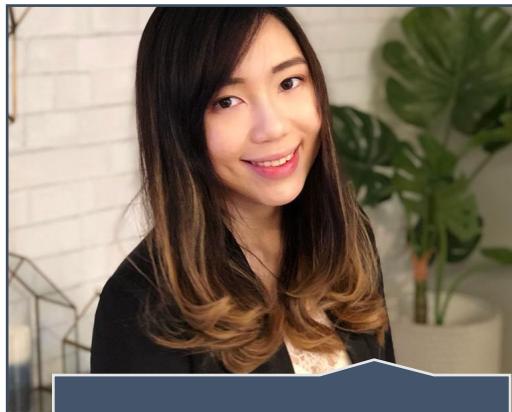
Eduardo, Jessica, Rohit & Stuart



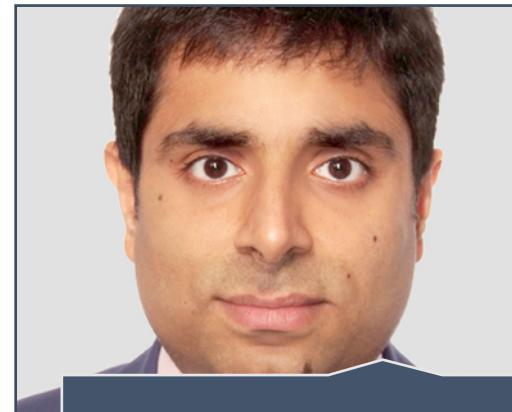
$\alpha\beta$ normal Distribution



Eduardo Peynetti



Jessica Wijaya

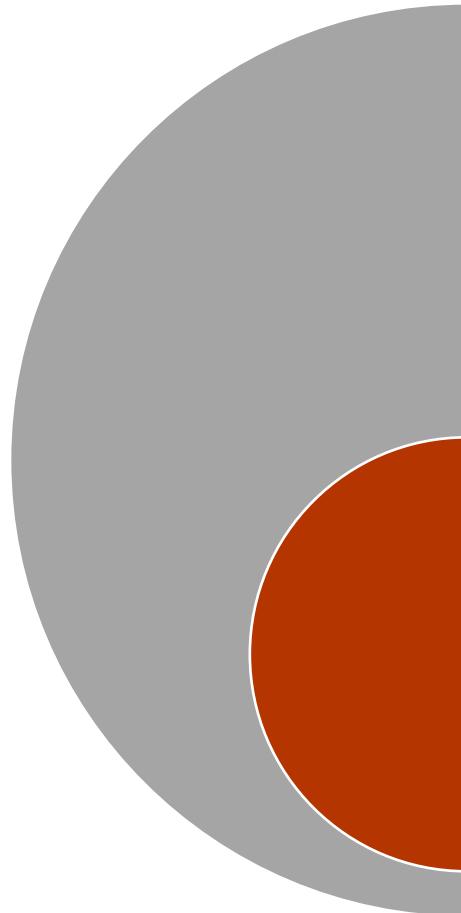


Rohit Beri



Stuart Neilson

Problem Statement



Can we use Transfer Learning and NLP
to extract sentiment from news?

Can news sentiment help to predict
future stock returns?



Datasets



27 million+ News Summaries

- 2.7 Million selected



3.3 million+ New Summaries

- 1.0 Million selected



Key Developments for S&P500 Companies

- 900K+ press releases and event announcements



Sharadar Financial Data

- Stocks Returns



10Ks for S&P500 Companies

- 6000+ documents with focus on Risk Factor, MD&A, and Market Risk



Baseline Models

Loughran McDonald

Dictionary based
model that analyzes
word frequencies

TF-IDF scores using
positive and negative
words

BERT for Sentiment Analysis

Pre-trained BERT
model

Trained on SST/2,
baseline sentiment
dataset

Fin-BERT for Sentiment Analysis

Language Model
trained on financial
news

Sentiment Model
trained on pre-labeled
phrases (+/-/-)



Our model

Fine-tuned Fin-BERT

Fine-tuned by using todays returns as a target

-1/+1 for on positive/negative returns

0 when news doesn't have a word in LM dictionary

Use today's returns as a proxy for sentiment

Fine-tuning against Today's Return as Sentiment Score

Noisy estimate of sentiment, easy to obtain for all news

Performance Evaluation

Long-Short Portfolio Construction

Every day:

Estimate average sentiment for each stock which had news

Buy stocks with positive sentiment, sell stocks with negative sentiment

Performance Measure

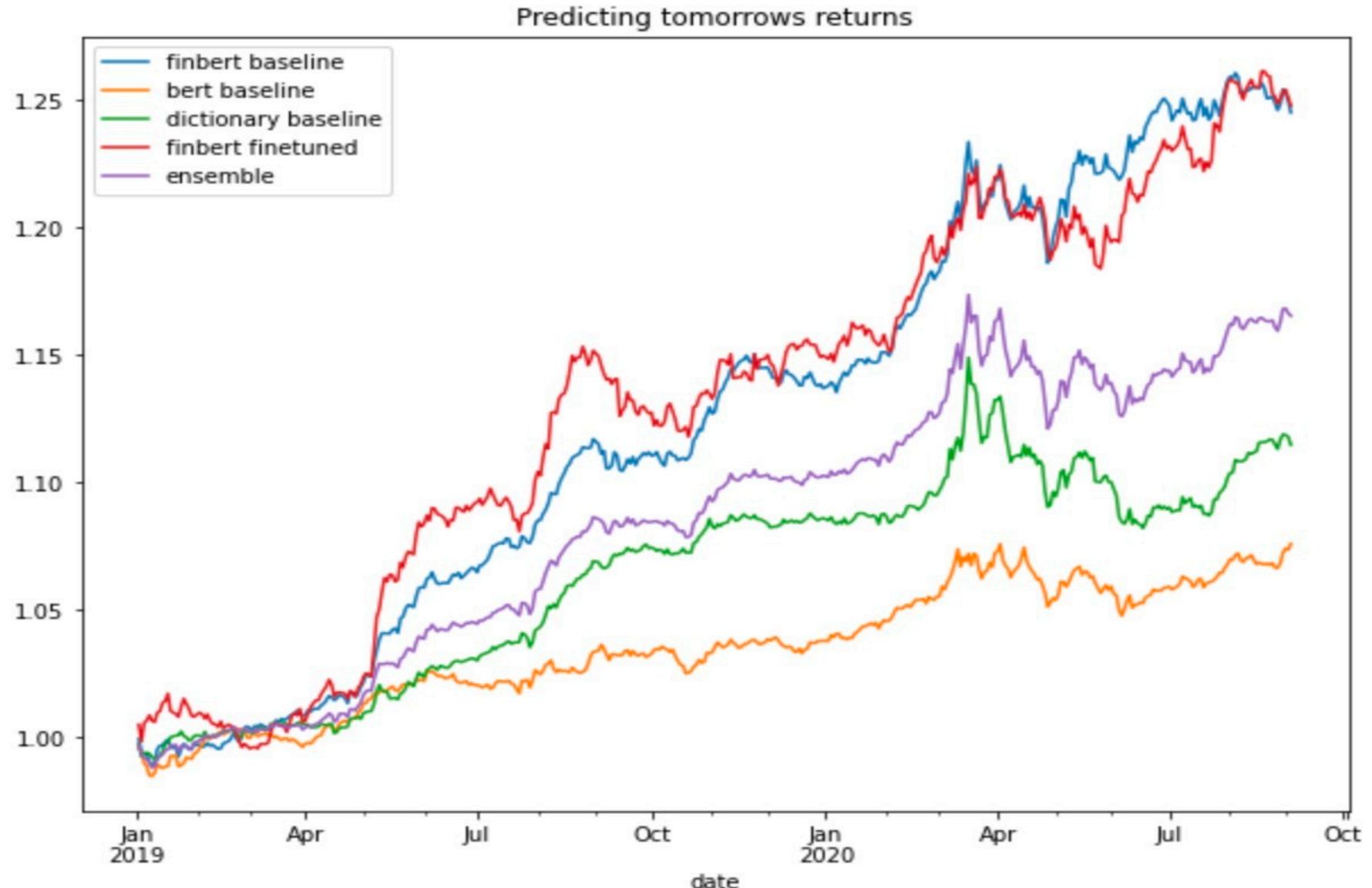
All models have accuracies around 51-52%. Wrong measure!

Use portfolio returns as performance measure



Performance Evaluation

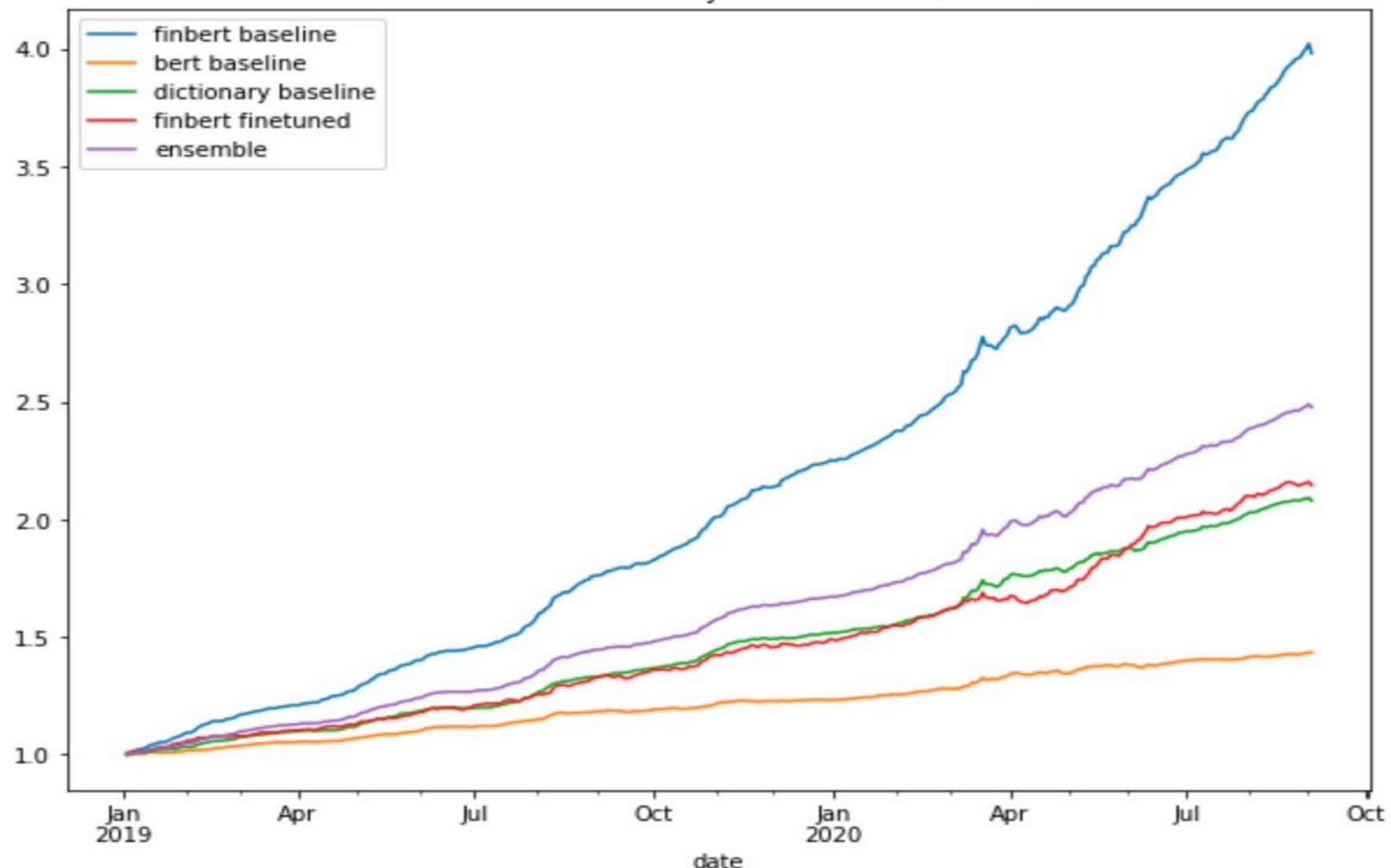
Predicting
next day's
returns



Sentiment Extraction

If we knew
today's news
beforehand?

What if we knew todays news ahead of time?



Sharpe Ratio

$$\frac{\bar{\mu}_{\text{Returns}}}{\sigma_{\text{Returns}}}$$

A way to standardize returns across stocks with different volatilities

Units of reward expected per unit of risk

LM baseline	BERT base	FinBert base	FinBert Fine-Tuned
1.75	1.6	3.05	2.45

Decay in return predictability

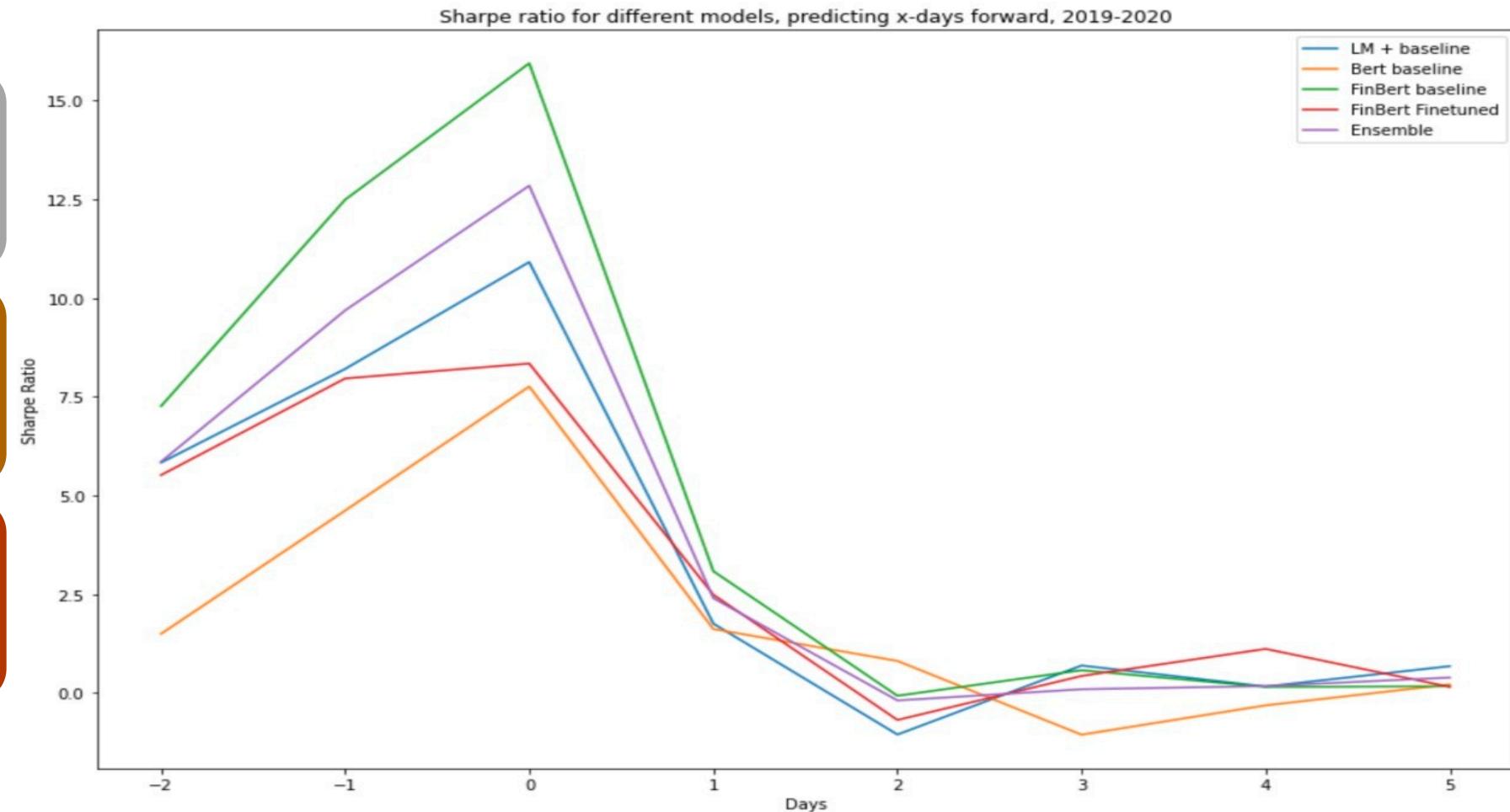
Sharpe Ratio for different models

Predicting x days forward

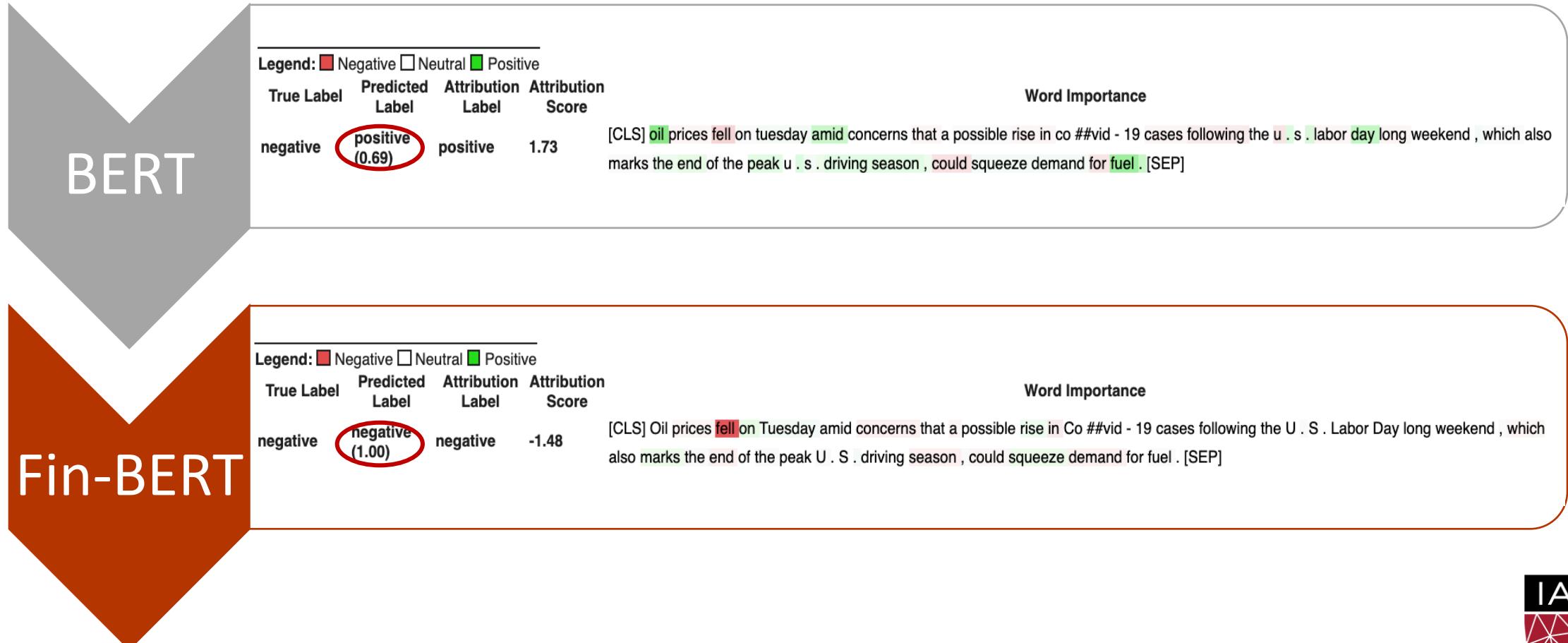
Todays news are predictive of yesterday's returns!

Old news / market knew the news before-hand

Small but persistent effect of news into the future



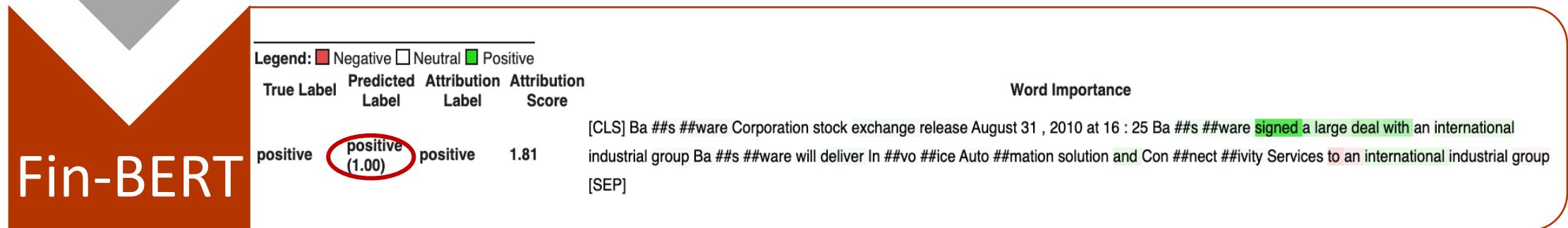
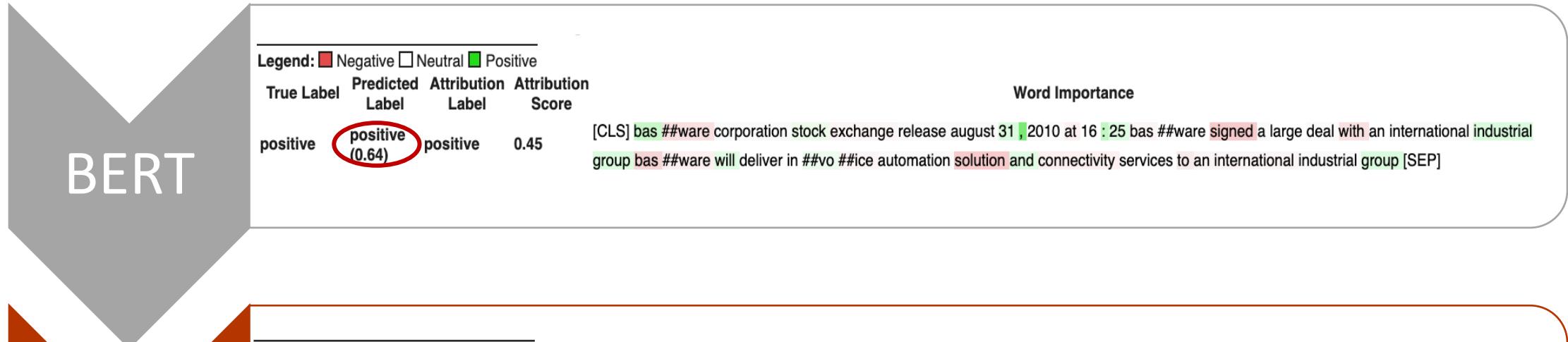
BERT vs. Fin-BERT – Negative (“oil prices fell”)



BERT vs. Fin-BERT – Positive (shares “doubled”)



BERT vs. Fin-BERT – Positive (“signed a large deal”)



Challenges

Financial Data is noisy!

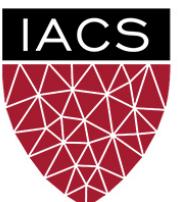
- Wrong timestamps mean fake performance
- Repeated news across datasets

Long training and evaluation time!

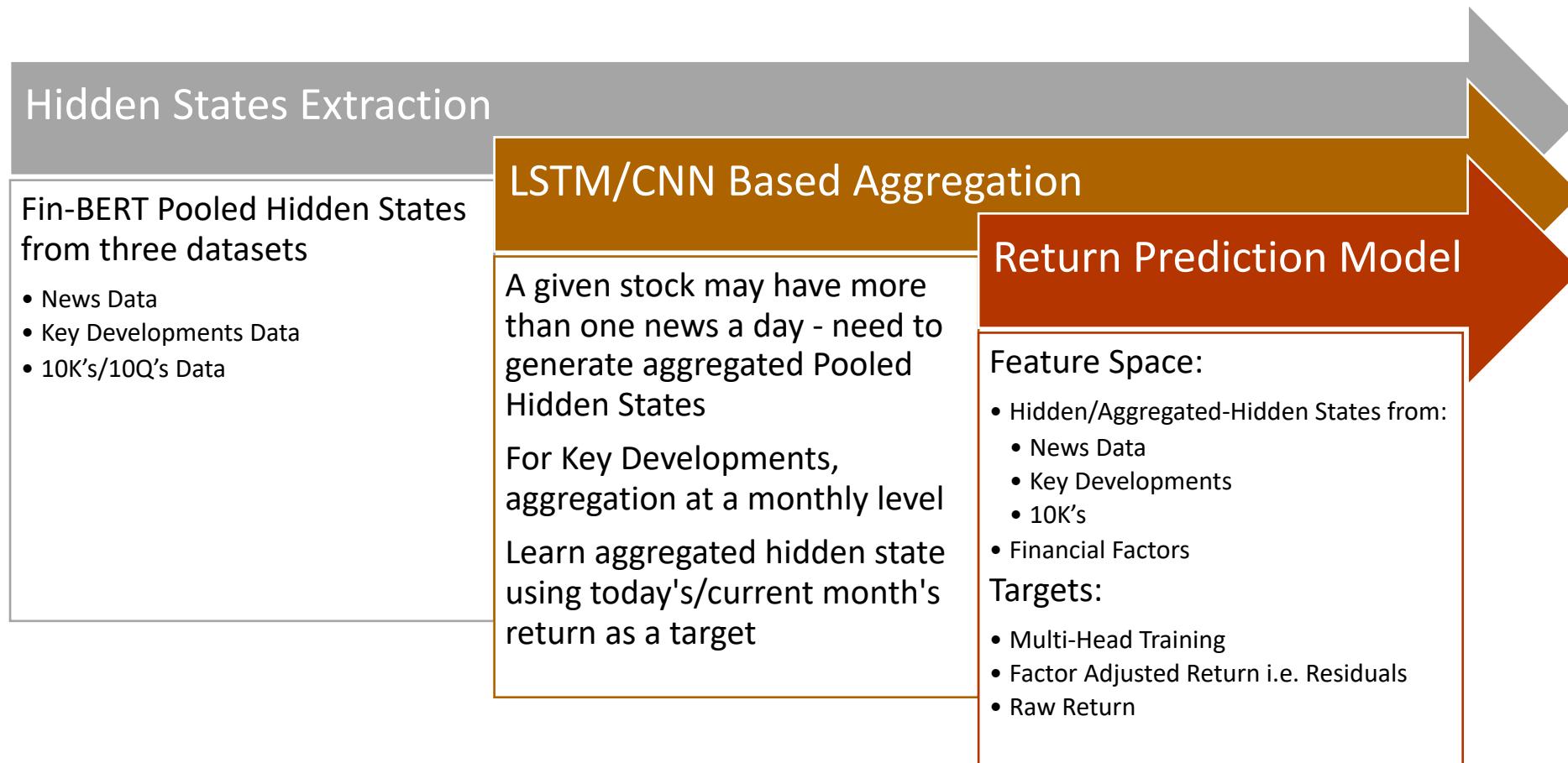
- Need prediction for all news in test set, over 2 million forward passes in BERT
- One lonely news in a small stock might dictate performance

Targets are noisy!

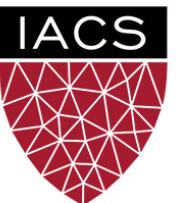
- Returns are just a soft measure of sentiment
- Easy to see model “forget” while training



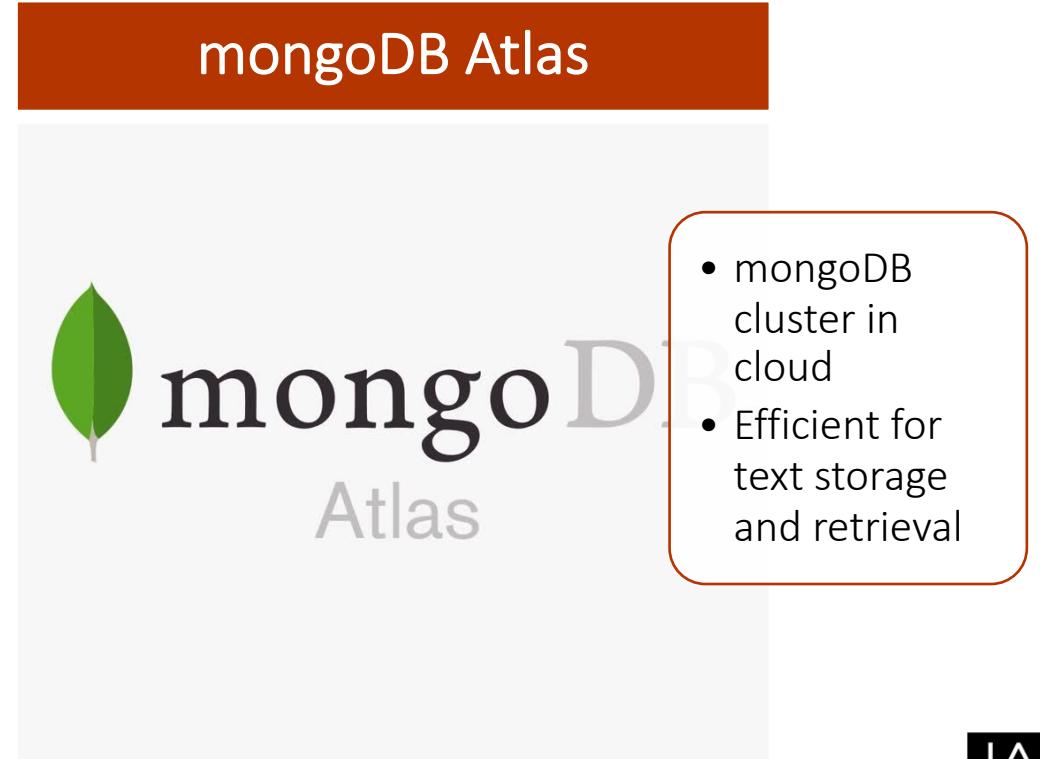
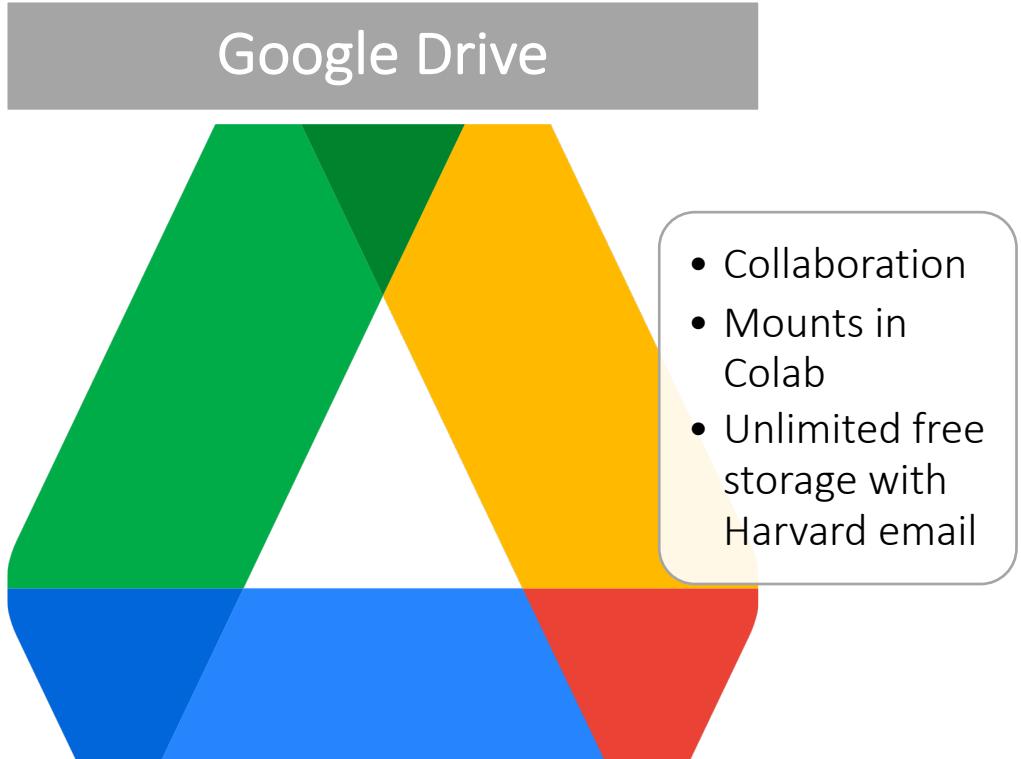
Work in Progress



Thanks!



Data Storage



Data Management



Raw text data

- Meta data
- Processed feature vectors
- **Cluster level pre-processing**

Financial and returns data

- Efficient integration with zipline – financial modelling library by Quantopian

TensorFlow model pipelines

- **Storing preprocessed data for efficient training**

Intermediate storage

- JSON
- CSV
- PKL

Data Preprocessing

News/Text Data

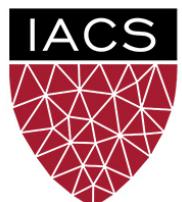
- Sentiment Score from BERT
- Pooled hidden layer representation from BERT Model

10K's

- Extracting relevant items from 10K's – Risk Factors, MD&A, Market Risk
- Summary for each item using LSA due to long length – up to 150K words per item

Financial Data

- Financial data from Quandl ingested into Zipline for performance analysis



BERT vs. Fin-BERT – Positive

(“announced a multi-year agreement”)

