

# BEYOND WORD-LEVEL TO SENTENCE-LEVEL SENTIMENT ANALYSIS FOR FINANCIAL REPORTS

Chi-Han Du,<sup>†</sup> Ming-Feng Tsai,<sup>‡</sup> Chuan-Ju Wang<sup>†</sup>

<sup>†</sup>Academia Sinica, Taiwan  
<sup>‡</sup>National Chengchi University, Taiwan



## What is sentiment analysis for financial reports?

Labeled in word-level by financial sentiment word lexicon (Loughran, 2011)

In addition, revenues increased due to fee income on growing **variable** COLI account values, partially **offset** by **declines** in fees on leveraged COLI as that block of business continues to **decline** due to the HIPA Act of 1996. Benefits, **claims** and expenses increased \$593, or 63%, to \$1.5 billion in 1998 from \$938 in 1997 due primarily to the MBL Recapture discussed previously.

Labeled in sentence-level by multiple financial experts (high risk)

In addition, revenues increased due to fee income on growing variable COLI account values, partially offset by declines in fees on leveraged COLI as that block of business continues to decline due to the HIPA Act of 1996. Benefits, claims and expenses increased \$593, or 63%, to \$1.5 billion in 1998 from \$938 in 1997 due primarily to the MBL Recapture discussed previously.

## Motivation

→ Use **existing knowledge (financial sentiment lexicon)** to improve **sentence-level classification performance** of deep learning models.

→ Extend **boundary of financial sentiment out of word range by semantics**, for each sentiment (**positive**, **negative**, **litigious**, and **uncertain**) shown in sentence.

In addition, revenues increased due to **fee income on growing variable COLI account values**, partially **offset by declines** in fees on leveraged COLI as that **block of business** continues to **decline** due to the HIPA Act of 1996. **Benefits, claims and expenses** increased \$593, or 63%, to \$1.5 billion in 1998 from \$938 in 1997 due primarily to the MBL Recapture discussed previously.

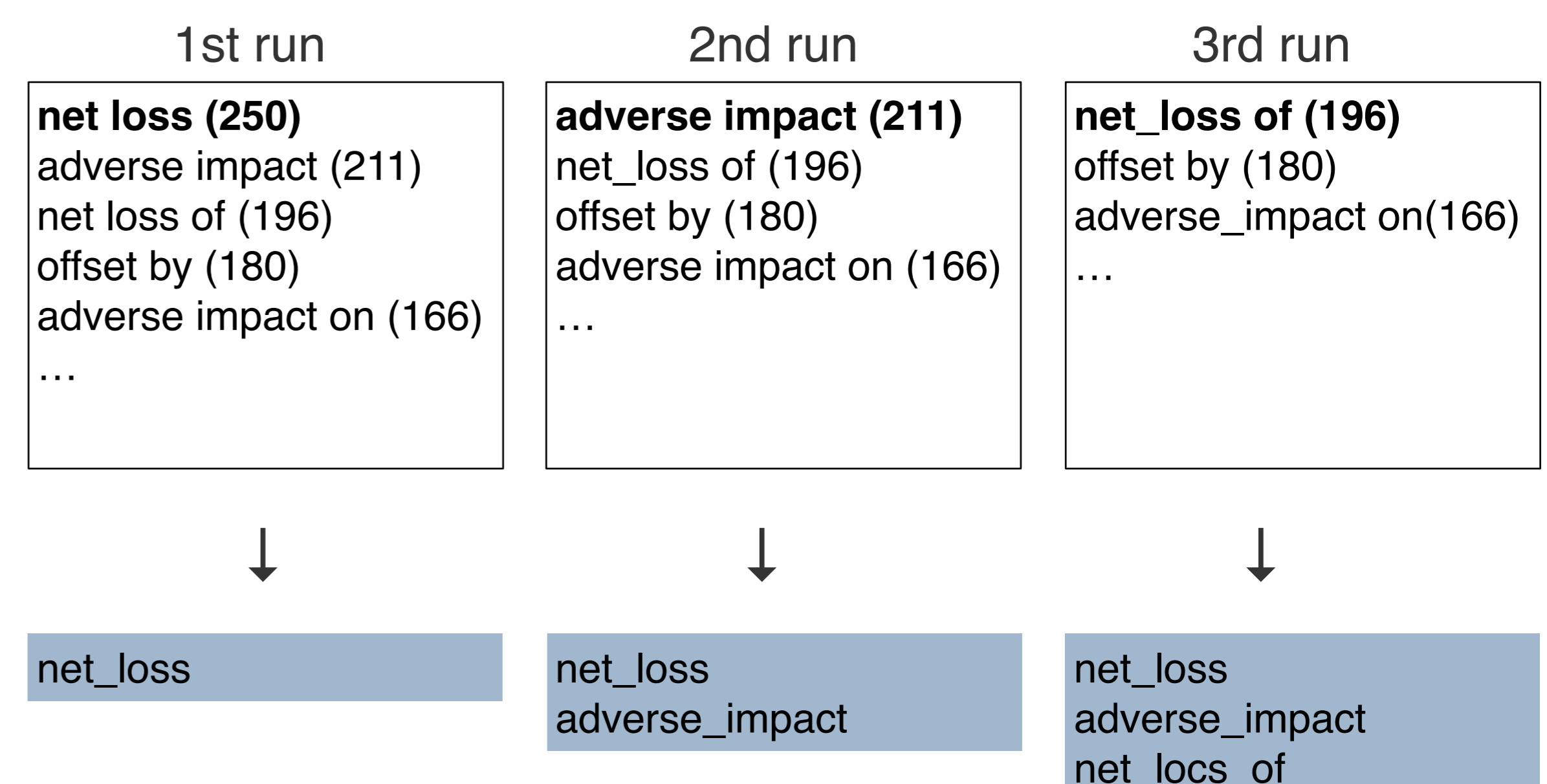
→ Examine **applicability of the proposed approach across models**, including traditional method, naive DL models, and more complicated models.

## Sub-phrase Algorithm

```

1 function Sub-Phrase ( $T_M, k, \ell$ );
   Input : A frequency table  $T_M$  including the top  $k$  most frequent sentiment  $n$ -grams and their frequencies, for  $n = 2, \dots, M$ ; the number of iterations,  $\ell$ 
   Output: A reference table,  $W$ 
2  $W \leftarrow \{\}$ ;
3 for  $e \leftarrow 1$  to  $\ell$  do
4   Find the most frequent word pair  $w_i$  and  $w_j$  in  $T_M$ ;
5   Find all  $n$ -grams containing  $w_i$  and  $w_j$  within  $T_M$ ;
6   Merge these two words into a new "word";
7   Add the merged new "word"  $w_i-w_j$  to the reference table  $W$ ;
8   Delete the most frequent word pair  $w_i$  and  $w_j$  in  $T_M$ ;
9   Update the frequency table  $T_M$  by replacing  $(w_i, w_j)$  as  $(w_i-w_j)$ ;
10 end
11 return  $W$ ;

```



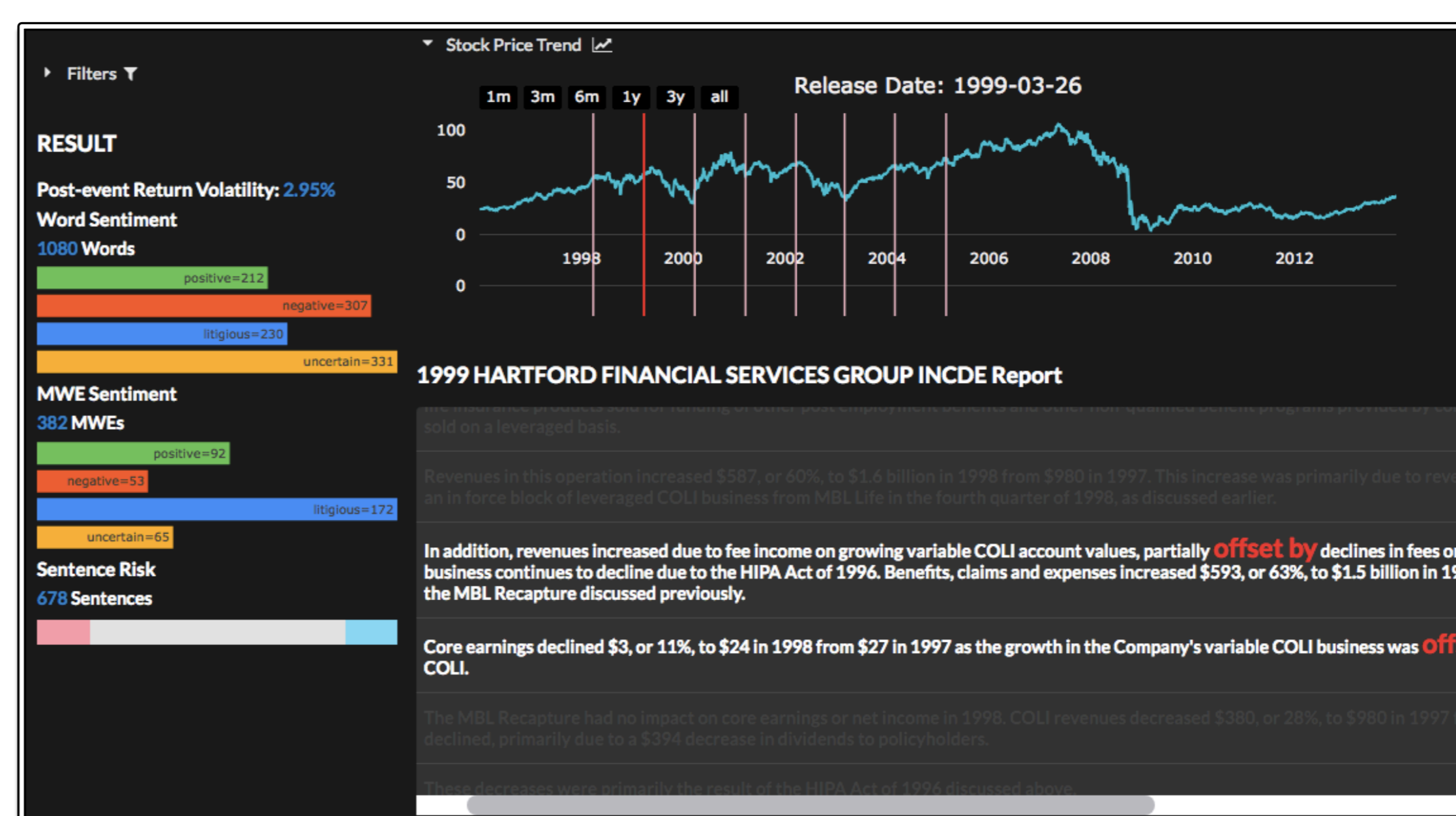
→ **The time spent is primarily proportional to  $k$** , which means that it is fast to implement.

## Main Results

	Accuracy	F1 score	
		High-risk	Neutral
tf-idf	<b>88.27</b>	<b>0.889</b>	0.876
tf-idf+senti-phrases	87.15	0.883	<b>0.880</b>
LSTM [8]	86.96	<b>0.893</b>	0.851
LSTM+senti-phrases	<b>87.14</b>	0.889	<b>0.857</b>
CNN [9]	86.33	0.852	0.891
CNN+senti-phrases	<b>86.35</b>	<b>0.861</b>	<b>0.915</b>
fastText [10]	87.76	0.858	0.895
fastText+senti-phrases	<b>88.03</b>	<b>0.922</b>	<b>0.901</b>
SiameseCBOW [11]	87.92	0.890	<b>0.902</b>
SiameseCBOW+senti-phrases	<b>88.79</b>	<b>0.927</b>	0.888

→ Combining words to generate senti-phrases is not beneficial to the traditional bag-of-words model.  
 → **Complicated DL models achieve better performance than naive models**, but all DL models perform better when using senti-phrases.

## An application



→ Our new developed tool: **Financial Risk Information Detecting and analyzing System (FRIDAYS) (AAAI'19)**



<https://cfda.csie.org/FRIDAYS/>

→ The proposed algorithm is fast to compress data and even improve the semantics of NLP models for financial texts.  
 → As a result, in the future it could be applied for **summarization of financial corpus**, or even **automatic generation (NLU) for financial reports**.