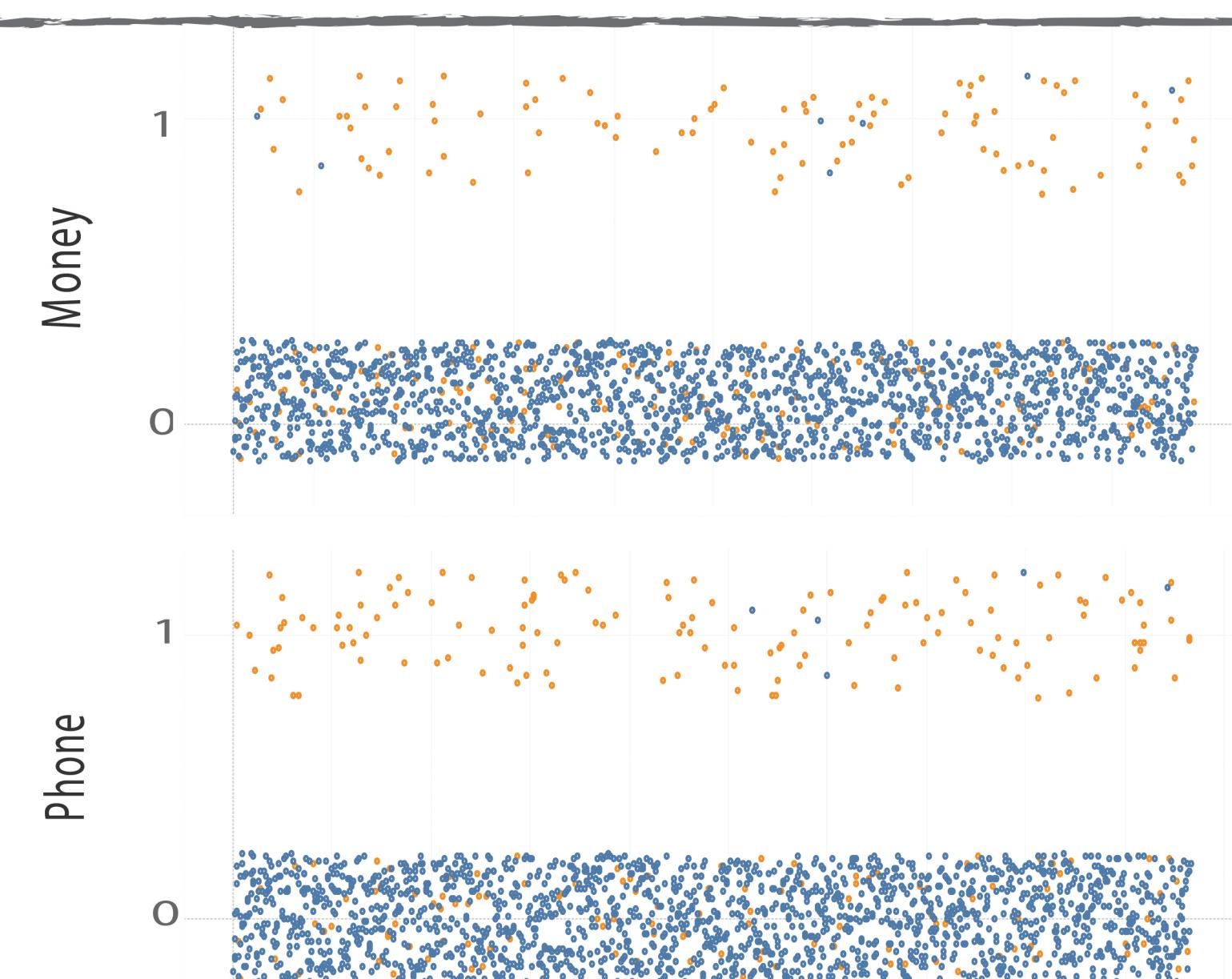


Wide use of SMS results increasing amount of spam messages. Classification of SMS message for spam detection helps manage our messages efficiently and reduce company's cost on message service.

**Dataset** 3000 original messages with labels from Kaggle

# Processing      Removing stop word lists and setting mean of document term frequency=5

# Key words in Spams

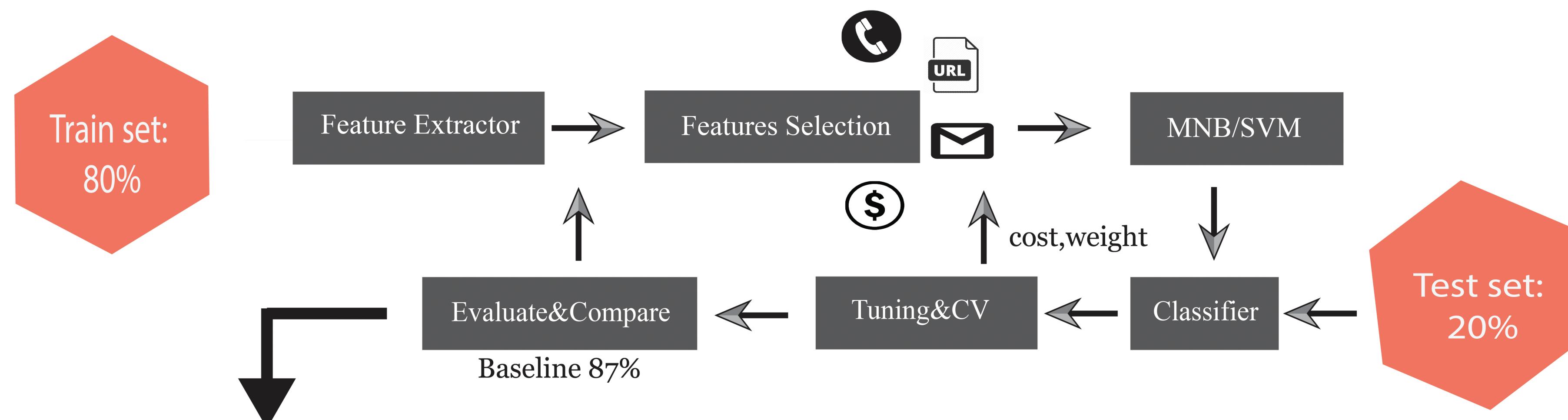


# SMS CLASSIFICATION FOR SPAM DETECTION

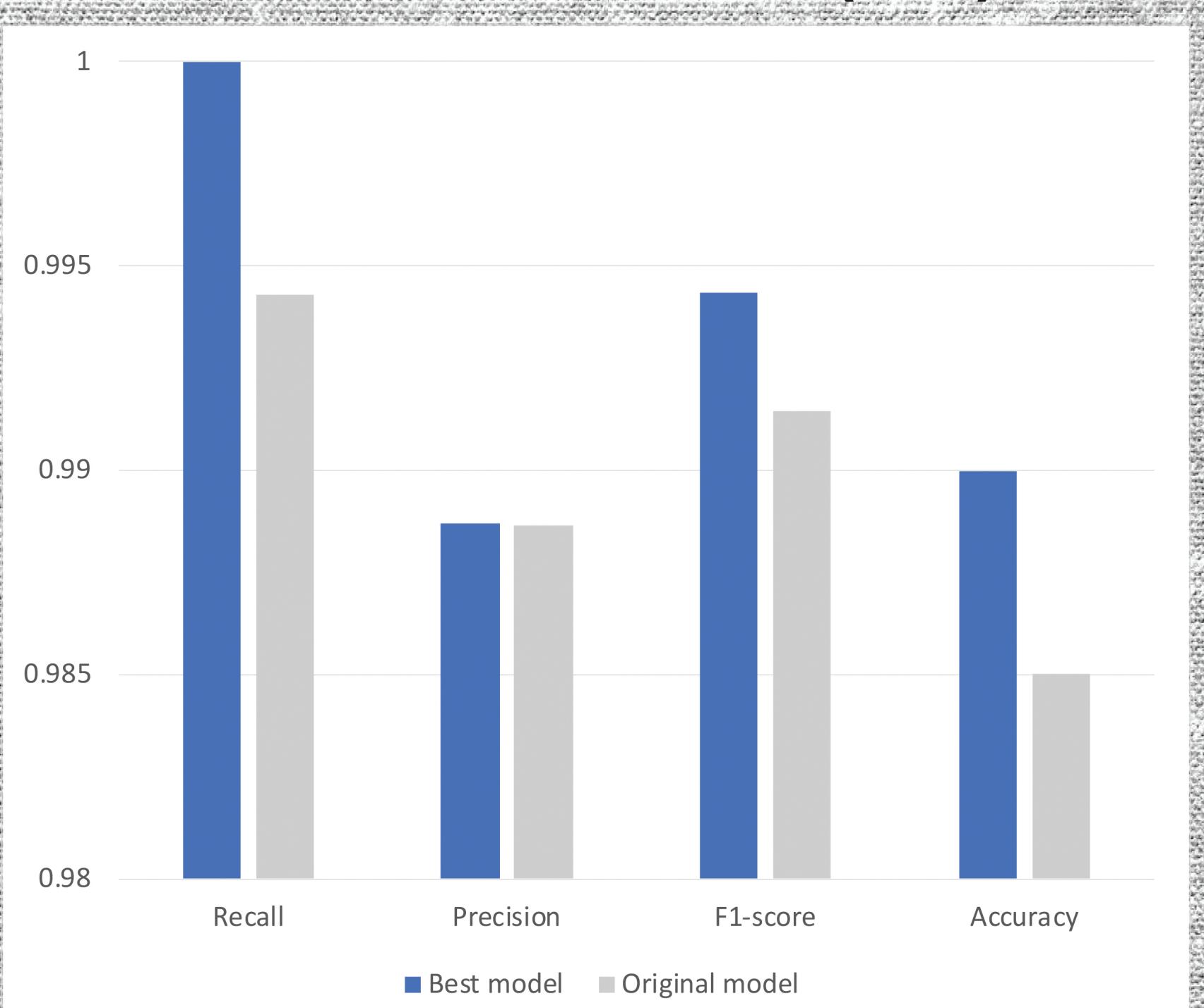
# Yunteng Gu, Nan Sun, Jiarong Yang

# Oops, a spam message!

# Delete



# The best model (M6)



M0: Unigram with TF				
Naive Bayes		SVM		
Actual	Spam	Spam	Not Spam	
Spam	517	8	522	3
Not Spam	3	72	6	69

Recall: 0.984  
Precision: 0.994  
F1-score: 0.989

Recall: 0.994  
Precision: 0.989  
F1-score: 0.991

M1: M0+Phone			
518	7	522	3
2	73	5	70

## M2: M0+URL

Confusion matrix for SVM with TF:

522	3
6	69

Recall: 0.994  
Precision: 0.989  
F1-score: 0.991

M1: M0+Phone			
518	7	522	3
2	73	5	70

## M2: M0+URL

M3: M0+Email

517	8
3	72

NB

M4: M1+Email

522	3
6	69

SVM

517	8
3	7

NB

 improved

**!!Phone number and Money  
Amount can improve perfor-  
mance!!**

# improved

M6: M1+M4+FullRecall

525	0
6	69

Recall: 1.000  
Precision: 0.989  
F1-score: 0.994

M5: M1+M4

	523	2
	2	73

Recall: 0.996  
Precision: 0.996  
F1-score: 0.996



**But recall is more important!**