



# Visual spoofing in content-based spam

- Visual spoofing is a type of threat where the similarities of characters and letters from different languages are used to confuse or trick users.
- The content-based filtering is also known as cognitive filtering that recommends items based on a comparison between the content of the items and a user profile items.

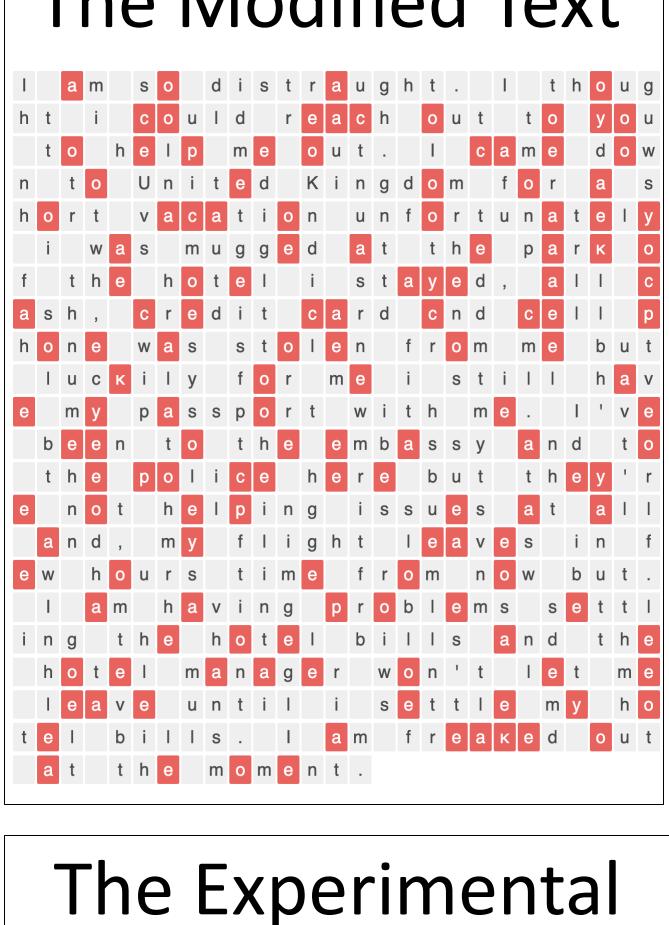
#### Background

- > The Unicode standard provides a unique number for every character, enabling disparate computing systems to exchange and interpret text in the same way.
- > Owing to ease of access, the number of users today continues to grow email rapidly. The use of email extends beyond just messages in text form; it is also used for sharing other types of data – images, videos, archive files, etc.
- Scammers constantly refine their approach of using fraudulent emails to induce individuals to reveal sensitive or private information – a practice known as phishing.
- > 22% of organizations see phishing as their greatest security threat
- > 94% of malware was delivered via email
- 65% of attacker groups used spear phishing as the primary infection vector

#### Research goals

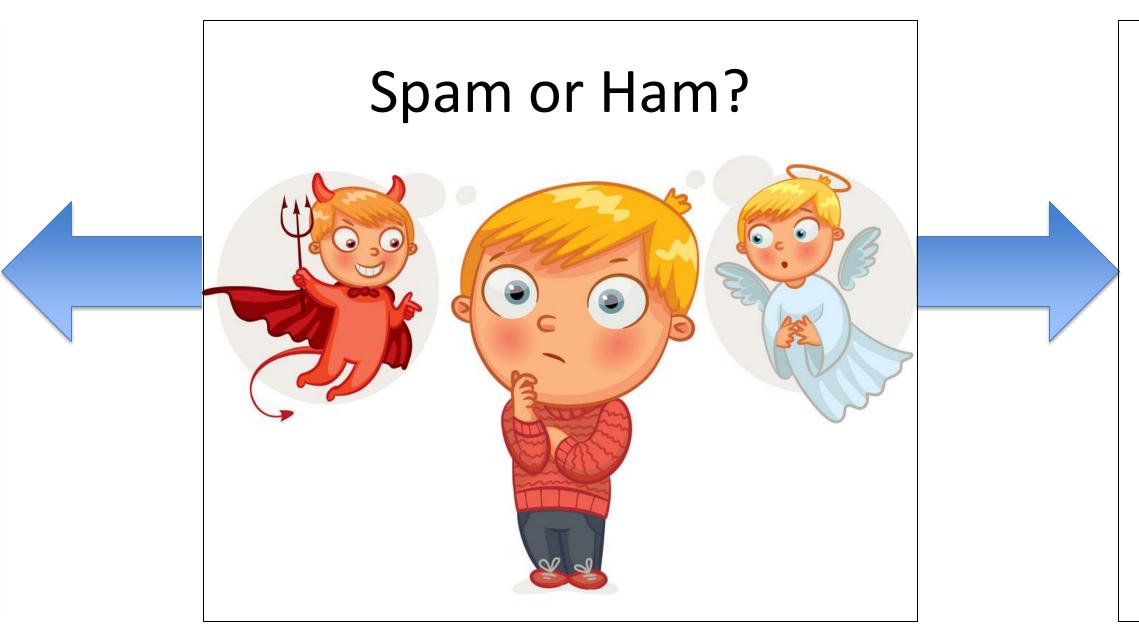
- > Identify new problems in the emails spam filtering.
- Classifying spam emails as ham emails and propose methods to address this threat.
- Protect people from cybercriminals and losing personal information such as bank account and card data, logins and passwords of Internet services, and so on.

I am so distraught. I thought i could reach out to you to help me out. I came down to United Kingdom for a short vacation unfortunately i was mugged at the park of the hotel i stayed, all cash, credit card cnd cell phone was stolen from me but luckily for me i still have my passport with me. I've been to the embassy and to the police here but they're not helping issues at all and, my flight leaves in few hours time from now but. I am having problems settling the hotel bills and the hotel manager won't let me leave until i settle my hotel bills. I am freaked out at the moment.



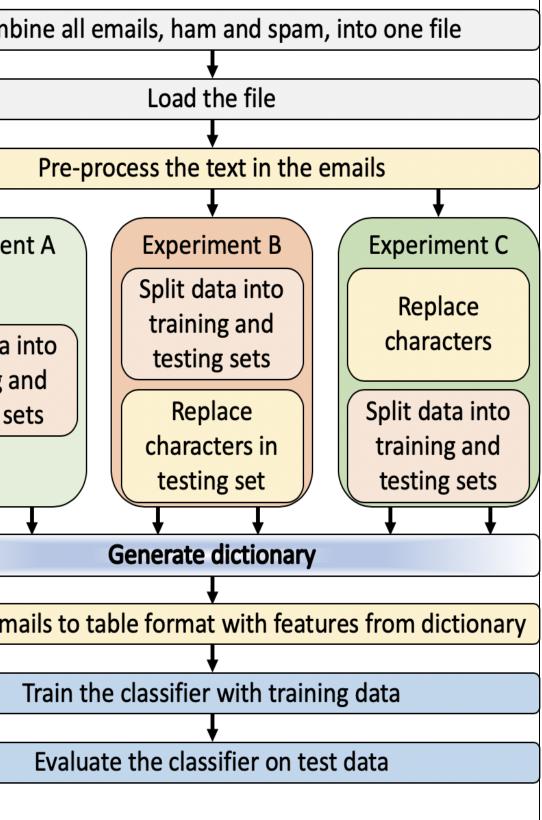
Com
Experime
Split data training testing s
$\overline{}$
Convert en

# Visual Spoofing in Content-Based Spam Detection



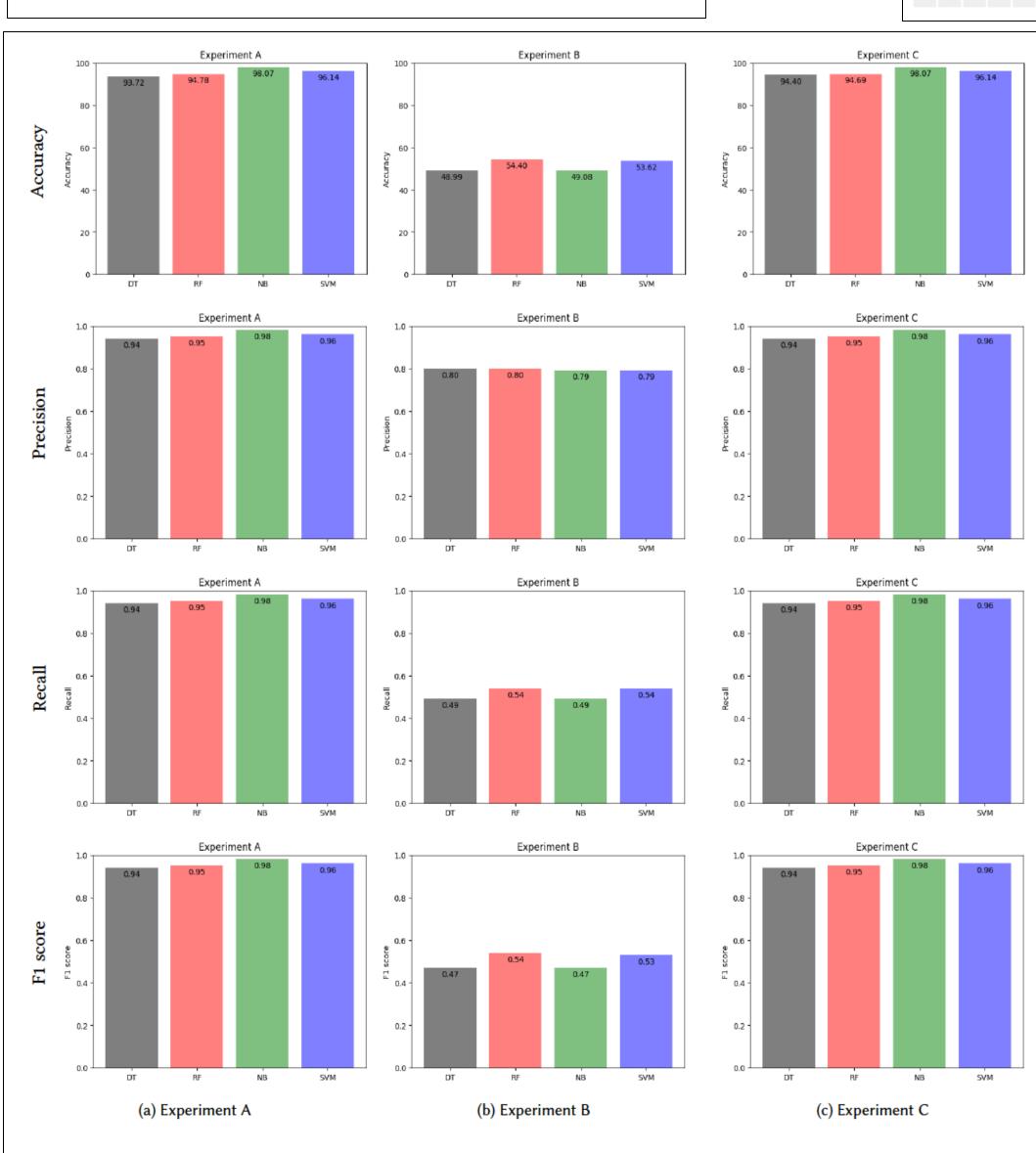
# The Modified Text

# Design



# The **Characters** Replace and their Replace

Character	Replac
a	U+(
e	U+(
k	U+0
ο	U+(
р	U+(
с	U+(
у	U+(
	1



# Mark Sokolov, MS in Software Engineering Advisor: Dr. Nic Herndon

I am so distraught. I thought i could reach out to you to help me out. I came down to United Kingdom for a short vacation unfortunately i was mugged at the park of the hotel i stayed, all cash, credit card and cell phone was stolen from me but luckily for me i still have my passport with me. I've been to the embassy and to the police here but they're not helping issues at all and, my flight leaves in few hours time from now but. I am having problems settling the hotel bills and the hotel manager won't let me leave until i settle my hotel bills. I am freaked out at the moment.

The Initial Text														t											
1		а	m		s	0	•	d	i	s	t	r	а	u	g	h	t	•				t	h	0	u
h	t		i		с	0	u	I	d		r	е	a	с	h		0	u	t		t	0		у	0
	t	0		h	е	Ι	р		m	е		0	u	t			Т		С	а	m	е		d	0
n		t	0		U	n	i	t	е	d		K	i	n	g	d	0	m		f	0	r		а	
h	0	r	t		۷	а	С	а	t	i	0	n		u	n	f	0	r	t	u	n	а	t	е	Ι
	i		w.	а	S		m	u	g	g	е	d					t	h	e		р	а	r	k	
Ť		t	h	е		h	0	t	e	1						а	,			,			1	1	
a	s	n	,		C	r	e	a	1	τ +	0	C	a	r	a	f	a				c	e	1	۱ h	
	1		6	k	i	a	5 V		5 f	ι 0	r	'	m			' i	1	0	+	i	1	e I		h	a
e	Ľ	m	v	K	' p	' a	y	s	' p	0	r	t		w	i	' t	h	0	'n	'e			T		v
Ū	b	e	e,	n	9	t	0	Ū	t	h	e		е	m	b	a	s	s	y	Ū	a	n	d		t
	t	h	е		р	0	I	i	с	е		h	е	r	е		b	u	t		t	h	е	у	١
е		n	0	t		h	е	I	р	i	n	g		i	s	s	u	е	s		а	t		а	I
	а	n	d	,		m	у		f	I	i	g	h	t		T	е	а	v	е	s		i	n	
е	W		h	0										f	r	0	m		n	0	W		b	u	t
													-										е		
i		-																		а			+		
	1										-	е								e					
t	e		u			1																	-		
						е			0																
	h e e	e w I i n f h t	htII	ashhonellucllucelmylbeeltheeundewlhlluainguihotlleaileaileailiailiailiailia	Iiihtiihoihniioiniiiinii <th>II<t< th=""><th>IIamISOhtiiiiiintoiiiiiiinitoiiiiiiiinii</th></t<><th>Iiii<t< th=""><th>IIamSoIdhtiIC0IIntoIMeIpIntoIIIIIIInitoIIIIIIIniINIIIIIIIIniINIIIIIIIIIniINII&lt;</th><th>I   I   A   M   S   O   I   I     h   t   I   I   C   O   I   I     h   t   O   I   A   E   I   P   I   M     n   t   O   I   P   I   P   I   M     n   t   O   I   P   I   P   I   M     n   T   T   N   I</th><th>IIAmSOJJJhtIICOUIJJITOIPIMPIMPITTOIIIIIIIIITTIII</th><th>IIamISOIIISIhtOICAIPIAIIITOIPIPIMIIIITOIPIPIMIIIITOIPIIPIMIIIITOIIIIIIIIIIIIIIIOIII&lt;</th><th>I&lt;</th><th>I a m s o u i s t r a   h t o i c o u l d r a   i t o i f o i i p i m e i a   i t o i i p i m i<i< th=""><th>I a m s o d i s t r a u   h t o i c o u l d r e a c   t o i h e l p i d i s i a c   n t o i u n i t a c a i<i< th=""><th>I a m s s a u g   h t i i c o u i d i s i<i i=""> i i i i i i i i i i i i i i i i<iii< td=""> i<i i=""> i i<i< th=""><th>I a m s s o d i s t r a u g h   h t i i c o u i d r e a c h i   i t o u h e i t i<th>I a m s s o d i s t r a u g h t   h t i i c o u l d r e a c h o o   t o h e l p m e d u t o<th>I a m s o u i s t r a u g h t .   h t i i c o u l d r e a c h u</th><th>I a m s o d i s t r a u g h t .   h t o h e l p m o u t d u r e a c h u t d u t d u t d u t d u t d u t d u t d u t d u t t t t t t u t t u t</th><th>I a m s o d i s t r a u g h t . I   h t i c o u l d v r e a c h u t .</th><th>I a m s 0 d i s t r a u g h t . I   h t i c 0 u I d r e a c h u t</th><th>I a m s o d i s t r a u g h t . I t t   h t i c o u I d r e a c h o u t t o u t t o u t o u t o u t t o u t</th><th>h t i c o u l r e a c h o u t t o   t o h e l p m e o u t . I c a m e   n t o U n i t e d K i n g d o m f o n a n<th>I a m s 0 u i s t r a u g h t . I u v v a u g h t . I u v r a u g h t . I u y</th></th></th></th></i<></i></iii<></i></th></i<></th></i<></th></t<></th></th>	II <t< th=""><th>IIamISOhtiiiiiintoiiiiiiinitoiiiiiiiinii</th></t<> <th>Iiii<t< th=""><th>IIamSoIdhtiIC0IIntoIMeIpIntoIIIIIIInitoIIIIIIIniINIIIIIIIIniINIIIIIIIIIniINII&lt;</th><th>I   I   A   M   S   O   I   I     h   t   I   I   C   O   I   I     h   t   O   I   A   E   I   P   I   M     n   t   O   I   P   I   P   I   M     n   t   O   I   P   I   P   I   M     n   T   T   N   I</th><th>IIAmSOJJJhtIICOUIJJITOIPIMPIMPITTOIIIIIIIIITTIII</th><th>IIamISOIIISIhtOICAIPIAIIITOIPIPIMIIIITOIPIPIMIIIITOIPIIPIMIIIITOIIIIIIIIIIIIIIIOIII&lt;</th><th>I&lt;</th><th>I a m s o u i s t r a   h t o i c o u l d r a   i t o i f o i i p i m e i a   i t o i i p i m i<i< th=""><th>I a m s o d i s t r a u   h t o i c o u l d r e a c   t o i h e l p i d i s i a c   n t o i u n i t a c a i<i< th=""><th>I a m s s a u g   h t i i c o u i d i s i<i i=""> i i i i i i i i i i i i i i i i<iii< td=""> i<i i=""> i i<i< th=""><th>I a m s s o d i s t r a u g h   h t i i c o u i d r e a c h i   i t o u h e i t i<th>I a m s s o d i s t r a u g h t   h t i i c o u l d r e a c h o o   t o h e l p m e d u t o<th>I a m s o u i s t r a u g h t .   h t i i c o u l d r e a c h u</th><th>I a m s o d i s t r a u g h t .   h t o h e l p m o u t d u r e a c h u t d u t d u t d u t d u t d u t d u t d u t d u t t t t t t u t t u t</th><th>I a m s o d i s t r a u g h t . I   h t i c o u l d v r e a c h u t .</th><th>I a m s 0 d i s t r a u g h t . I   h t i c 0 u I d r e a c h u t</th><th>I a m s o d i s t r a u g h t . I t t   h t i c o u I d r e a c h o u t t o u t t o u t o u t o u t t o u t</th><th>h t i c o u l r e a c h o u t t o   t o h e l p m e o u t . I c a m e   n t o U n i t e d K i n g d o m f o n a n<th>I a m s 0 u i s t r a u g h t . I u v v a u g h t . I u v r a u g h t . I u y</th></th></th></th></i<></i></iii<></i></th></i<></th></i<></th></t<></th>	IIamISOhtiiiiiintoiiiiiiinitoiiiiiiiinii	Iiii <t< th=""><th>IIamSoIdhtiIC0IIntoIMeIpIntoIIIIIIInitoIIIIIIIniINIIIIIIIIniINIIIIIIIIIniINII&lt;</th><th>I   I   A   M   S   O   I   I     h   t   I   I   C   O   I   I     h   t   O   I   A   E   I   P   I   M     n   t   O   I   P   I   P   I   M     n   t   O   I   P   I   P   I   M     n   T   T   N   I</th><th>IIAmSOJJJhtIICOUIJJITOIPIMPIMPITTOIIIIIIIIITTIII</th><th>IIamISOIIISIhtOICAIPIAIIITOIPIPIMIIIITOIPIPIMIIIITOIPIIPIMIIIITOIIIIIIIIIIIIIIIOIII&lt;</th><th>I&lt;</th><th>I a m s o u i s t r a   h t o i c o u l d r a   i t o i f o i i p i m e i a   i t o i i p i m i<i< th=""><th>I a m s o d i s t r a u   h t o i c o u l d r e a c   t o i h e l p i d i s i a c   n t o i u n i t a c a i<i< th=""><th>I a m s s a u g   h t i i c o u i d i s i<i i=""> i i i i i i i i i i i i i i i i<iii< td=""> i<i i=""> i i<i< th=""><th>I a m s s o d i s t r a u g h   h t i i c o u i d r e a c h i   i t o u h e i t i<th>I a m s s o d i s t r a u g h t   h t i i c o u l d r e a c h o o   t o h e l p m e d u t o<th>I a m s o u i s t r a u g h t .   h t i i c o u l d r e a c h u</th><th>I a m s o d i s t r a u g h t .   h t o h e l p m o u t d u r e a c h u t d u t d u t d u t d u t d u t d u t d u t d u t t t t t t u t t u t</th><th>I a m s o d i s t r a u g h t . I   h t i c o u l d v r e a c h u t .</th><th>I a m s 0 d i s t r a u g h t . I   h t i c 0 u I d r e a c h u t</th><th>I a m s o d i s t r a u g h t . I t t   h t i c o u I d r e a c h o u t t o u t t o u t o u t o u t t o u t</th><th>h t i c o u l r e a c h o u t t o   t o h e l p m e o u t . I c a m e   n t o U n i t e d K i n g d o m f o n a n<th>I a m s 0 u i s t r a u g h t . I u v v a u g h t . I u v r a u g h t . I u y</th></th></th></th></i<></i></iii<></i></th></i<></th></i<></th></t<>	IIamSoIdhtiIC0IIntoIMeIpIntoIIIIIIInitoIIIIIIIniINIIIIIIIIniINIIIIIIIIIniINII<	I   I   A   M   S   O   I   I     h   t   I   I   C   O   I   I     h   t   O   I   A   E   I   P   I   M     n   t   O   I   P   I   P   I   M     n   t   O   I   P   I   P   I   M     n   T   T   N   I	IIAmSOJJJhtIICOUIJJITOIPIMPIMPITTOIIIIIIIIITTIII	IIamISOIIISIhtOICAIPIAIIITOIPIPIMIIIITOIPIPIMIIIITOIPIIPIMIIIITOIIIIIIIIIIIIIIIOIII<	I<	I a m s o u i s t r a   h t o i c o u l d r a   i t o i f o i i p i m e i a   i t o i i p i m i <i< th=""><th>I a m s o d i s t r a u   h t o i c o u l d r e a c   t o i h e l p i d i s i a c   n t o i u n i t a c a i<i< th=""><th>I a m s s a u g   h t i i c o u i d i s i<i i=""> i i i i i i i i i i i i i i i i<iii< td=""> i<i i=""> i i<i< th=""><th>I a m s s o d i s t r a u g h   h t i i c o u i d r e a c h i   i t o u h e i t i<th>I a m s s o d i s t r a u g h t   h t i i c o u l d r e a c h o o   t o h e l p m e d u t o<th>I a m s o u i s t r a u g h t .   h t i i c o u l d r e a c h u</th><th>I a m s o d i s t r a u g h t .   h t o h e l p m o u t d u r e a c h u t d u t d u t d u t d u t d u t d u t d u t d u t t t t t t u t t u t</th><th>I a m s o d i s t r a u g h t . I   h t i c o u l d v r e a c h u t .</th><th>I a m s 0 d i s t r a u g h t . I   h t i c 0 u I d r e a c h u t</th><th>I a m s o d i s t r a u g h t . I t t   h t i c o u I d r e a c h o u t t o u t t o u t o u t o u t t o u t</th><th>h t i c o u l r e a c h o u t t o   t o h e l p m e o u t . I c a m e   n t o U n i t e d K i n g d o m f o n a n<th>I a m s 0 u i s t r a u g h t . I u v v a u g h t . I u v r a u g h t . I u y</th></th></th></th></i<></i></iii<></i></th></i<></th></i<>	I a m s o d i s t r a u   h t o i c o u l d r e a c   t o i h e l p i d i s i a c   n t o i u n i t a c a i <i< th=""><th>I a m s s a u g   h t i i c o u i d i s i<i i=""> i i i i i i i i i i i i i i i i<iii< td=""> i<i i=""> i i<i< th=""><th>I a m s s o d i s t r a u g h   h t i i c o u i d r e a c h i   i t o u h e i t i<th>I a m s s o d i s t r a u g h t   h t i i c o u l d r e a c h o o   t o h e l p m e d u t o<th>I a m s o u i s t r a u g h t .   h t i i c o u l d r e a c h u</th><th>I a m s o d i s t r a u g h t .   h t o h e l p m o u t d u r e a c h u t d u t d u t d u t d u t d u t d u t d u t d u t t t t t t u t t u t</th><th>I a m s o d i s t r a u g h t . I   h t i c o u l d v r e a c h u t .</th><th>I a m s 0 d i s t r a u g h t . I   h t i c 0 u I d r e a c h u t</th><th>I a m s o d i s t r a u g h t . I t t   h t i c o u I d r e a c h o u t t o u t t o u t o u t o u t t o u t</th><th>h t i c o u l r e a c h o u t t o   t o h e l p m e o u t . I c a m e   n t o U n i t e d K i n g d o m f o n a n<th>I a m s 0 u i s t r a u g h t . I u v v a u g h t . I u v r a u g h t . I u y</th></th></th></th></i<></i></iii<></i></th></i<>	I a m s s a u g   h t i i c o u i d i s i <i i=""> i i i i i i i i i i i i i i i i<iii< td=""> i<i i=""> i i<i< th=""><th>I a m s s o d i s t r a u g h   h t i i c o u i d r e a c h i   i t o u h e i t i<th>I a m s s o d i s t r a u g h t   h t i i c o u l d r e a c h o o   t o h e l p m e d u t o<th>I a m s o u i s t r a u g h t .   h t i i c o u l d r e a c h u</th><th>I a m s o d i s t r a u g h t .   h t o h e l p m o u t d u r e a c h u t d u t d u t d u t d u t d u t d u t d u t d u t t t t t t u t t u t</th><th>I a m s o d i s t r a u g h t . I   h t i c o u l d v r e a c h u t .</th><th>I a m s 0 d i s t r a u g h t . I   h t i c 0 u I d r e a c h u t</th><th>I a m s o d i s t r a u g h t . I t t   h t i c o u I d r e a c h o u t t o u t t o u t o u t o u t t o u t</th><th>h t i c o u l r e a c h o u t t o   t o h e l p m e o u t . I c a m e   n t o U n i t e d K i n g d o m f o n a n<th>I a m s 0 u i s t r a u g h t . I u v v a u g h t . I u v r a u g h t . I u y</th></th></th></th></i<></i></iii<></i>	I a m s s o d i s t r a u g h   h t i i c o u i d r e a c h i   i t o u h e i t i <th>I a m s s o d i s t r a u g h t   h t i i c o u l d r e a c h o o   t o h e l p m e d u t o<th>I a m s o u i s t r a u g h t .   h t i i c o u l d r e a c h u</th><th>I a m s o d i s t r a u g h t .   h t o h e l p m o u t d u r e a c h u t d u t d u t d u t d u t d u t d u t d u t d u t t t t t t u t t u t</th><th>I a m s o d i s t r a u g h t . I   h t i c o u l d v r e a c h u t .</th><th>I a m s 0 d i s t r a u g h t . I   h t i c 0 u I d r e a c h u t</th><th>I a m s o d i s t r a u g h t . I t t   h t i c o u I d r e a c h o u t t o u t t o u t o u t o u t t o u t</th><th>h t i c o u l r e a c h o u t t o   t o h e l p m e o u t . I c a m e   n t o U n i t e d K i n g d o m f o n a n<th>I a m s 0 u i s t r a u g h t . I u v v a u g h t . I u v r a u g h t . I u y</th></th></th>	I a m s s o d i s t r a u g h t   h t i i c o u l d r e a c h o o   t o h e l p m e d u t o <th>I a m s o u i s t r a u g h t .   h t i i c o u l d r e a c h u</th> <th>I a m s o d i s t r a u g h t .   h t o h e l p m o u t d u r e a c h u t d u t d u t d u t d u t d u t d u t d u t d u t t t t t t u t t u t</th> <th>I a m s o d i s t r a u g h t . I   h t i c o u l d v r e a c h u t .</th> <th>I a m s 0 d i s t r a u g h t . I   h t i c 0 u I d r e a c h u t</th> <th>I a m s o d i s t r a u g h t . I t t   h t i c o u I d r e a c h o u t t o u t t o u t o u t o u t t o u t</th> <th>h t i c o u l r e a c h o u t t o   t o h e l p m e o u t . I c a m e   n t o U n i t e d K i n g d o m f o n a n<th>I a m s 0 u i s t r a u g h t . I u v v a u g h t . I u v r a u g h t . I u y</th></th>	I a m s o u i s t r a u g h t .   h t i i c o u l d r e a c h u	I a m s o d i s t r a u g h t .   h t o h e l p m o u t d u r e a c h u t d u t d u t d u t d u t d u t d u t d u t d u t t t t t t u t t u t	I a m s o d i s t r a u g h t . I   h t i c o u l d v r e a c h u t .	I a m s 0 d i s t r a u g h t . I   h t i c 0 u I d r e a c h u t	I a m s o d i s t r a u g h t . I t t   h t i c o u I d r e a c h o u t t o u t t o u t o u t o u t t o u t	h t i c o u l r e a c h o u t t o   t o h e l p m e o u t . I c a m e   n t o U n i t e d K i n g d o m f o n a n <th>I a m s 0 u i s t r a u g h t . I u v v a u g h t . I u v r a u g h t . I u y</th>	I a m s 0 u i s t r a u g h t . I u v v a u g h t . I u v r a u g h t . I u y

Evaluation metrics Experiment for Α characters (no Experiment replaced), (characters replaced in the test set), and Experiment (characters replaced train and test IN the sets) from following machine algorit hms used: decision tree (DT), random forest (RF), naïve Bayes (NB), and support vector machine (SVM).



# Experiment A

Each email text was preserved in its original form – i.e., no characters were replaced. We trained and evaluated our set of classifiers using data from this distribution only.

## Experiment B

We modified data by testing our introducing corresponding Cyrillic letters in place of the Latin letters 'a', 'e', 'k', 'o', 'p', 'c', and 'y'. As desired. This resulted in single and mixed-script confusables in our testing dataset. The intent here was to simulate a visual spoofing effect in the email messages that our testing set consists of. With the original character encoding in our training dataset still preserved, we trained all classifiers using data from the same distribution as in experiment A. However, model evaluation is done using data effectively from a different distribution than the one used in experiment A.

## Experiment C

We modified both our training and testing datasets to replicate visual spoofing in the entire distribution used to develop our model, using the set of confusables introduced in experiment B. Training and evaluation was performed using data from this modified distribution only. As a result, unlike experiments A or B, each of our models would simulate spam filters designed to classify emails that contain visual spoofing.

## Conclusion and future work

- $\succ$  These experiments indicate that using a classifier trained on data using Latin alphabet, to classify a message with a combination of Latin and Cyrillic letters leads to much lower classification accuracy compared to the same classifier used with a message with Latin characters only.
- $\succ$  In future work we plan to evaluate this approach with characters from multiple alphabets. In addition, we would like to investigate the impact of this method with applications for other used text communication.