

Diakritikovač slovenského textu

Autor: Jakub Gederá

Vedúci práce: Ing. Marián Šimko, PhD.

18.6.2015

Motivácia

- Internet, svet bez diakritiky
- Mnoho textu na internete neobsahuje diakritiku
 - Emaily
 - Diskusné príspevky
 - Blogy
- Strata informácie pri ďalšom spracovaní textu

Motivácia

- Problémy so spracovaním takýchto textov:
 - Kategorizácia textu,
 - Extrakcia metadát, kľúčových slov
- Umožniť používateľovi Webu automaticky rekonštruovať diakritiku
- Obsah s diakritikou vyzerá profesionálnejšie
- Takýto text sa lepšie číta

Identifikovaný problém

- Veta: Nova pracka perie dobre.
- Nova:
 - Nova (eruptívna hviezda)
 - Nová (prídavné meno)
- Pracka:
 - Práčka (prístroj na pranie)
 - Pračka (bitka, ruvačka)
 - Pracka (predmet slúžiaci na spínanie niečoho)
- Perie:
 - Perie (sloveso)
 - Perie (podstatné meno)
- Dobre:
 - Dobre (príslovka)
 - Dobré (prídavné meno)

Naiivné riešenia

- Nájsť v slovníku možnosti doplnenia diakritiky pre jednotlivé slová a použiť náhodný výber
- Vybrať slovo s diakritikou, ktoré sa pri vytváraní slovníka vyskytlo najčastejšie
- Slabé výsledky

Využime kontext

- Textová spojitost
- V bare sa uskutočnila ___?___.
 - a) Práčka
 - b) Pračka
 - c) Pracka
- **Ako využiť kontext?**
- Použiť jazykový model slovenského jazyka

Jazykový model

- Pravdepodobnostný a štatistický model jazyka a postupnosti slov [Jurafsky, Martin, 2000]
- Trénovaný na veľkom korpuse textov
- Využívaný aj pri rozpoznávaní reči
- Analogicky problém k doplneniu diakritiky
- Identifikovať najpravdepodobnejšieho slova na zázname
- Pravdepodobnostný odhad nasledujúceho slova
- Nájdenie najpravdepodobnejšej sekvencie slov
-

Použitý jazykový model

- Využívame model, ktorý bol vytvorený JÚLŠ SAV zo Slovenského Národného Korpusu textov
- Je dôležité, na akých typoch textoch je model natrénovaný
- Model použitý v metóde bol natrénovaný na publicistických textoch

Náš algoritmus

1. Tokenizuj vstup a tokeny preved' na malé písmena

Nova pracka perie dobre.

nova pracka perie dobre .

2. Pre každý token nájdí v slovníku jeho možnosti diakritikovania

nova práčka perie dobre .

nová pračka dobré

pracka

Náš algoritmus

3. Prechádzaj text po n-ticiach tokenov a vyhodnoň všetky možné n-gramy jazykovým modelom
 - nova pracka perie dobre. -> -24.7071705
 - nová pračka perie dobre. -> -23.2947597
 - ...
 - nová práčka perie dobré. -> -18.7531738
 - nová práčka perie dobre. -> -17.3096389
- (Pre prácu s jazykovým modelom bola použitá knižnica *kenlm* a jej metóda `score()`, ktorá vráti logaritmus pravdepodobnosti výskytu vyhodnocovaného n-gramu)

Náš algoritmus

4. Vyber n-gram, ktorý bol ohodnotený najlepšie
 - nová práčka perie dobre
5. Rekonštruuj veľkosť písmen
 - Nová práčka perie dobre

Experiment

1. Bez používateľa

- Sledujeme úspešnosť doplnenia diakritiky
- Ako dáta boli použité typy textov, ktoré sa na Webe vyskytujú najčastejšie

2. S používateľom

- Účastníci experimentu vytvorili texty, ktoré sme následne rekonštruovali
- Účastníci ohodnotili svoju spokojnosť s rekonštruovaným textom

Experiment bez používateľa

- Ako testovacia vzorka boli použité 3 typy textov, aby sme ohodnotili vhodnosť použitého jazykového modelu.
- Typy testovaných textov:
 - Publicistické texty
 - Príspevky na diskusných fórach
 - Umelecké texty

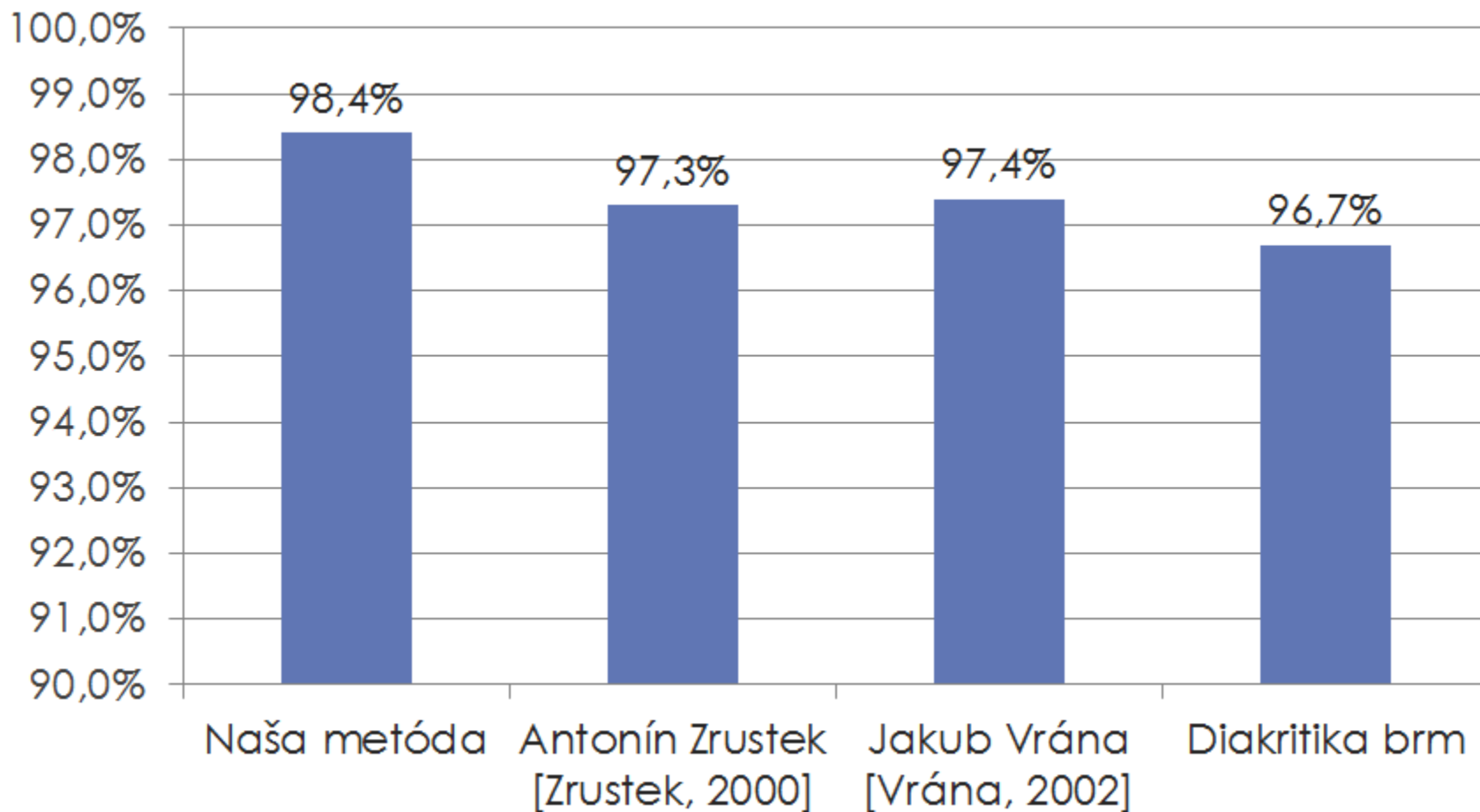
Experiment bez používatele

Typ textu	Počet slov	Počet chýb	Dosiahnutá úspešnosť
Diskusné fórum	39 123	627	98,4 %
Publicistický štýl	45 321	1 108	97,6 %
Umelecké texty	32 110	3 443	89,3 %

Chyby v diakritikovaní

- Dva typy chýb
 1. Pre slovo sa nenašiel správny variant diakritikovania v slovníku (rozšíriť slovník)
 2. Správny variant bol v slovníku, ale vyhodnotenie modelom bolo nesprávne (obohatiť, vytvoriť nový, kvalitnejší model)
- Sú v pomere približne **1:1**

Porovnanie s inými riešeniami



Experiment s používateľmi

- 11 používateľov testovalo spokojnosť s rekonštruovaným textom
- Vytvárali texty vyskytujúce sa bežne na Webe, na diskusných fórach, aplikácia rekonštruovala diakritiku

Počet slov	Počet chýb	Dosiahnutá úspešnosť
1 460	36	97,5 %

- V priemere jedno slovo zo 40 chybné rekonštruované

Experiment s používateľom

- Hodnotili na stupnici
 - Veľmi nespokojný
 - Nespokojný
 - Spokojný
 - Veľmi spokojný

	Veľmi nespokojný	nespokojný	spokojný	Veľmi spokojný
Počet používateľov	0	0	6	5



Diakritikovač

od autora jakub.gedera

★★★★★ (1)

[Produktivita](#)

Počet používateľov: 2

PRIDANÉ DO CHROMU

PREHĽAD

RECENZIE

PODPORA

PODOBNÉ

g+1 0

Doplňovač

diakritiky

Kompatibilné s vaším zariadením

Tento doplnok doplna diakritiku do označeného textu

Aplikácia doplňujúca diakritiku. Označte text a kliknite na ikonu rozšírenia. Označený text je nahradený textom s doplnenou diakritikou. Užitočná vec pri písaní, resp. čítaní emailov alebo iného webového obsahu. Rekonštrukcia s úspešnosťou nad 98%. V priemere jedno slovo z 50 je rekonštruované chybné.

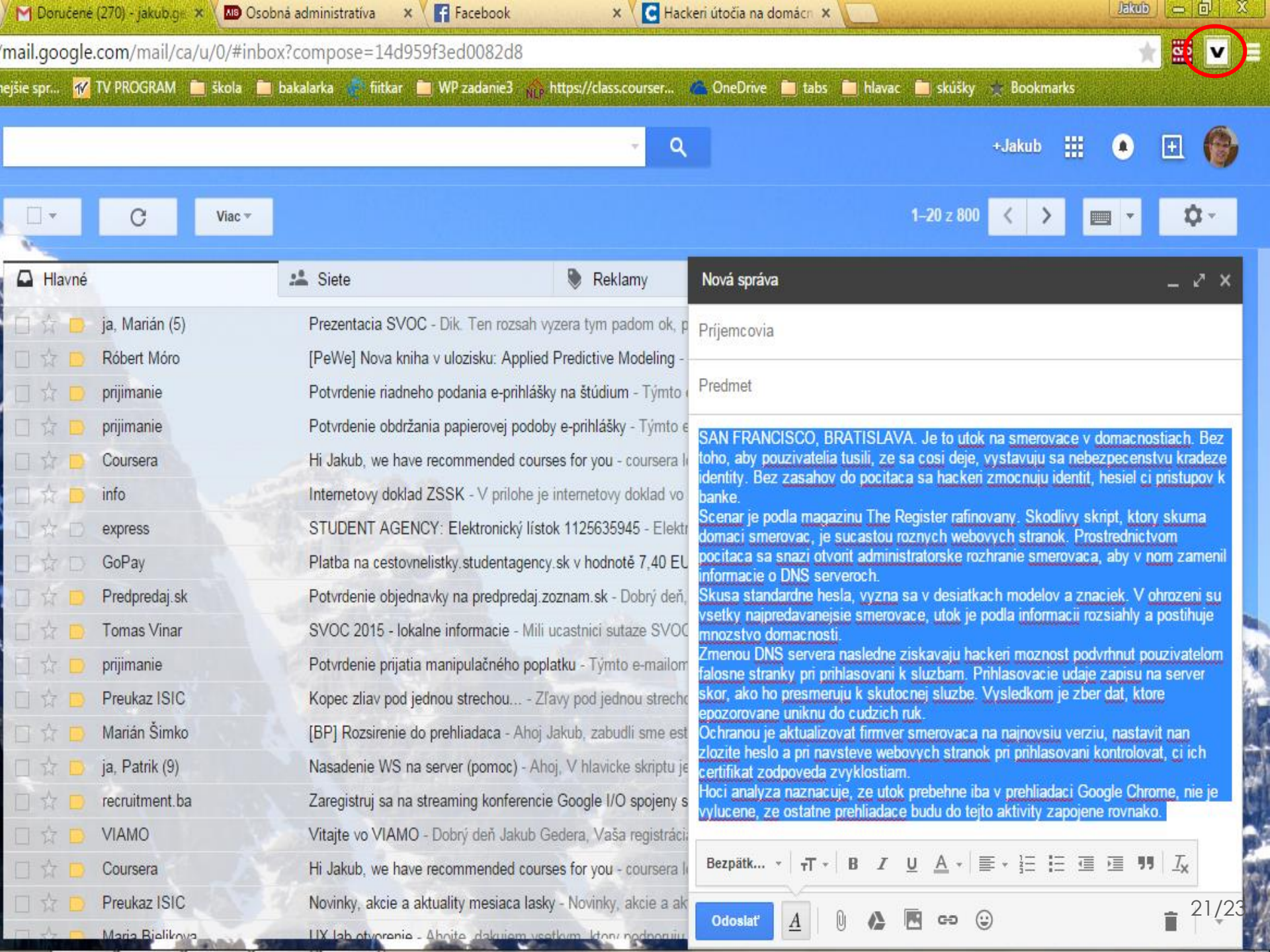
Nahlásiť zneužitie

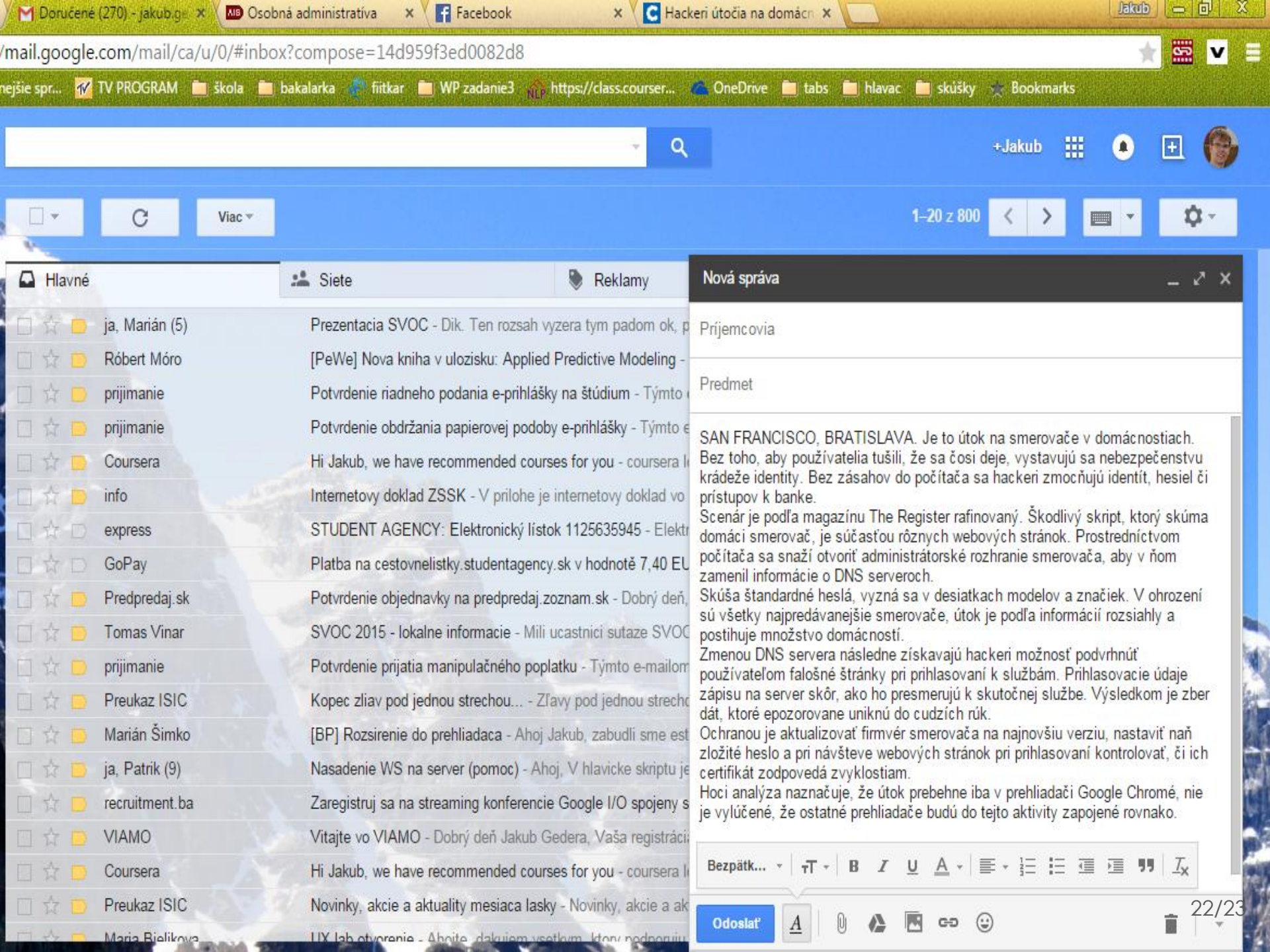
Verzia: 1.0

Aktualizované: 4. júna 2015

Veľkosť: 3.76KB

Jazyk: slovenský





Hlavné Siete Reklamy

- ja, Marián (5) - Prezencia SVOC - Dik. Ten rozsah vyzerá tým padom ok, p
- Róbert Móro - [PeWe] Nova kniha v ulozisku: Applied Predictive Modeling -
- prijimanie - Potvrdenie riadneho podania e-prihlášky na štúdium - Týmto
- prijimanie - Potvrdenie obdržania papierovej podoby e-prihlášky - Týmto e
- Coursera - Hi Jakub, we have recommended courses for you - coursera l
- info - Internetovy doklad ZSSK - V prílohe je internetovy doklad vo
- express - STUDENT AGENCY: Elektronický lístok 1125635945 - Elekt
- GoPay - Platba na cestovnelistky.studentagency.sk v hodnote 7,40 EU
- Predpredaj.sk - Potvrdenie objednávky na predpredaj.zoznam.sk - Dobry deň,
- Tomas Vinar - SVOC 2015 - lokalne informacie - Mili ucastnici sutaze SVOC
- prijimanie - Potvrdenie prijatia manipulačného poplatku - Týmto e-mailom
- Preukaz ISIC - Kopec zliav pod jednou strechou... - Zľavy pod jednou strecho
- Marián Šimko - [BP] Rozsirenie do prehliadaca - Ahoj Jakub, zabudli sme est
- ja, Patrik (9) - Nasadenie WS na server (pomoc) - Ahoj, v hlavicke skriptu je
- recruitment.ba - Zaregistruj sa na streaming konferencie Google I/O spojeny s
- VIAMO - Vitajte vo VIAMO - Dobry deň Jakub Gedera, Vaša registráci
- Coursera - Hi Jakub, we have recommended courses for you - coursera l
- Preukaz ISIC - Novinky, akcie a aktuality mesiaca lasky - Novinky, akcie a ak
- Maria Bielkova - LIX lab otvorenie - Ahojte, dakujem vsetom, ktory podporuju

Nová správa [close icon]

Prijemcovia

Predmet

SAN FRANCISCO, BRATISLAVA. Je to útok na smerovače v domácnostiach. Bez toho, aby používatelia tušili, že sa čosi deje, vystavujú sa nebezpečenstvu krádeže identity. Bez zásahov do počítača sa hackeri zmocňujú identít, hesiel či prístupov k banke.

Scenár je podľa magazínu The Register rafinovaný. Škodlivý skript, ktorý skúma domáci smerovač, je súčasťou rôznych webových stránok. Prostredníctvom počítača sa snaží otvoriť administrátorské rozhranie smerovača, aby v ňom zamenil informácie o DNS serveroch.

Skúša štandardné heslá, vyzná sa v desiatkach modelov a značiek. V ohrození sú všetky najpredávanejšie smerovače, útok je podľa informácií rozsiahly a postihuje množstvo domácností.

Zmenou DNS servera následne získavajú hackeri možnosť podvrhnúť používateľom falošné stránky pri prihlasovaní k službám. Prihlasovacie údaje zapisu na server skôr, ako ho presmerujú k skutočnej službe. Výsledkom je zber dát, ktoré epozorovane uniknú do cudzích rúk.

Ochranou je aktualizovať firmvér smerovača na najnovšiu verziu, nastaviť na ň zložité heslo a pri návšteve webových stránok pri prihlasovaní kontrolovať, či ich certifikát zodpovedá zvyklostiam.

Hoci analýza naznačuje, že útok prebehne iba v prehliadači Google Chromé, nie je vylúčené, že ostatné prehliadače budú do tejto aktivity zapojené rovnako.

Bezpätk... [font size icon] [bold icon] [italic icon] [underline icon] [text color icon] [bullet list icon] [numbered list icon] [link icon] [quote icon] [undo icon]

Odoslať [text color icon] [attach icon] [insert icon] [undo icon] [redo icon] [emoji icon]

22/23

Záver

- Pri dopĺňaní diakritiky sme **využili kontext** pre nájdenie správneho variantu slova s diakritikou
- Priblížili sme aplikáciu diakritikovania bežným používateľom Webu formou **rozšírenia pre webový prehliadač**
- **Overili** sme **úspešnosť** diakritikovania
- Aplikácia bola zo strany používateľov ohodnotená **pozitívne** a vedia si ju predstaviť v **praxi** (emaily, blogy, neformálne dokumenty)
- Nie len pri písaní textov ale aj pri čítaní