

The Same is Not The Same

Postcorrection of Alphabet Confusion Errors
in Mixed-Alphabet OCR Recognition

by Jonas Hempel

Bulgerian-German to English

Иван ора нивата

In English: Ivan plowed the field.

'opa' is German word for 'grandfather'

Alphabet Similarities (1)

<i>Latin</i>	<i>A</i>	<i>B</i>	<i>C</i>	<i>E</i>	<i>H</i>	<i>K</i>	<i>M</i>	<i>O</i>	<i>P</i>	<i>T</i>	<i>X</i>	<i>Y</i>	<i>a</i>	<i>c</i>	<i>e</i>	<i>g</i>	<i>k</i>	<i>m</i>	<i>n</i>	<i>o</i>	<i>p</i>	<i>u</i>	<i>x</i>	<i>y</i>
<i>Cyrillic</i>	<i>A</i>	<i>B</i>	<i>C</i>	<i>E</i>	<i>H</i>	<i>K</i>	<i>M</i>	<i>O</i>	<i>P</i>	<i>T</i>	<i>X</i>	<i>Y</i>	<i>a</i>	<i>c</i>	<i>e</i>	-	-	<i>m</i>	<i>n</i>	<i>o</i>	<i>p</i>	<i>u</i>	<i>x</i>	<i>y</i>
<i>Latin</i>	A	B	C	E	H	K	M	O	P	T	X	Y	a	c	e	g	k	m	n	o	p	u	x	y
<i>Cyrillic</i>	A	B	C	E	H	K	M	O	P	T	X	Y	a	c	e	g	k	m	n	o	p	u	x	y

- Latin-Cyrillic transition table
- Upper font is Times New Roman
- Lower font is Universum

Alphabet Similarities (2)

<i>Latin</i>	<i>A</i>	<i>B</i>	<i>E</i>	<i>Z</i>	<i>H</i>	<i>I</i>	<i>K</i>	<i>M</i>	<i>N</i>	<i>O</i>	<i>P</i>	<i>T</i>	<i>Y</i>	<i>a</i>	<i>y</i>	<i>n</i>	<i>i</i>	<i>o</i>	<i>p</i>	<i>u</i>	<i>w</i>
<i>Greek</i>	<i>A</i>	<i>B</i>	<i>E</i>	<i>Z</i>	<i>H</i>	<i>I</i>	<i>K</i>	<i>M</i>	<i>N</i>	<i>O</i>	<i>P</i>	<i>T</i>	<i>Y</i>	<i>α</i>	<i>γ</i>	<i>η</i>	<i>ι</i>	<i>ο</i>	<i>ρ</i>	<i>υ</i>	<i>ω</i>
<i>Latin</i>	<i>A</i>	<i>B</i>	<i>E</i>	<i>Z</i>	<i>H</i>	<i>I</i>	<i>K</i>	<i>M</i>	<i>N</i>	<i>O</i>	<i>P</i>	<i>T</i>	<i>Y</i>	<i>a</i>	<i>y</i>	<i>n</i>	<i>i</i>	<i>o</i>	<i>p</i>	<i>u</i>	<i>w</i>
<i>Greek</i>	<i>A</i>	<i>B</i>	<i>E</i>	<i>Z</i>	<i>H</i>	<i>I</i>	<i>K</i>	<i>M</i>	<i>N</i>	<i>O</i>	<i>P</i>	<i>T</i>	<i>Y</i>	<i>α</i>	<i>γ</i>	<i>η</i>	<i>ι</i>	<i>ο</i>	<i>ρ</i>	<i>υ</i>	<i>ω</i>

- Latin-Greek transition table
- Upper font is Times New Roman
- Lower font is Verdana Cursive

Training and Test corpora

- Sophia-Munich corpus
- Bulgarian EC corpus
- Greek-Latin corpus

Algorithm

- Levenshtein distance $d_0(w_i, v)$
- Normalized similarity value $s(v, w_i)$
- collocation frequency value $f(v, w_{i-1}, w_{i+1})$
→ $\text{score}(v) = \alpha * s(v, w_i) + (1-\alpha) * f(v)$
- α balance parameter
- τ threshold parameter

Evaluation Results (1)

Corpus/OCR	tokens	error rate OCR → pc	ac-error rate OCR → pc
SM-OCR1-Tr	8110	11.22 → 6.57%	5.42 → 1.25%
SM-OCR1-Te	7923	10.59 → 6.25%	5.44 → 1.26%
SM-OCR2-Tr	5099	37.24 → 15.87%	9.96 → 1.75%
SM-OCR2-Te	5115	43.63 → 16.81%	10.28 → 2.01%
EC-OCR1-Tr	6571	15.05 → 5.68%	10.94 → 1.81%
EC-OCR1-Te	6230	16.44 → 7.03%	13.00 → 3.74%
EC-OCR2-Tr	6571	48.52 → 9.0%	27.71 → 4.14%
EC-OCR2-Te	6230	48.81 → 11.35%	27.50 → 3.23%

- Bulgarian Sophia-Munich and Bulgarian EC corpus
- Error rate for plain OCR recognition and postcorrection
- Training (Tr) and Test (Te) data
- ac-error: alphabet confusion error

Evaluation Results (2)

Corpus/OCR	err.r.	ac-errors	ac-err. r.
GR-OCR1-Ti	2.37%	62	0.74%
GR-OCR1-Vd	2.12%	49	0.58%
GR-OCR2-Ti	12.00%	550	6.54%
GR-OCR2-Vd	10.87%	457	5.43%

- Greek newspaper corpus
- Cursive Times (Ti) and cursive Verdana (Vd) font
- ac-error: alphabet confusion error

