

Někdo rozumí počítačům. Počítače nerozumí nikomu. Naučíme je to?

Diacritics Restoration

Otázka: Nevadi, když píšou bez háček a carek?

Odpověď: Ne.

Lidé často píšou bez diakritiky. Pokud chceme takovým textům rozumět a třeba v nich vyhledávat, musíme si diakritiku domyslet. Někdy je to snadné, protože slovo bez diakritiky ve slovníku nenajdeme (např. *pocitac*). Jindy je třeba využít znalostí o frekvencích výskytu a spoluvýskytu slov, které získáme z korpusu. Z nich program spočítá pravděpodobnost oháčkovaného slova (např. *měj se místo mej se*).



Jeste v patek se zdalo, ze to bude dobre, dnes je to podstatne horsi. V kazdem pripade se neco musi stat, abychom mohli se cti dostat svym zavazkum.

Ještě v pátek se zdálo, že to bude dobře, dnes je to podstatně horší. V každém případě se něco musí stát, abychom mohli se ctí dostat svým závazkům.

Kopírovat

Problémy s chybějící diakritikou jsou běžné i v mnoha dalších jazycích, např. ve slovenštině, galštině, španělštině, rumunštině či haitské kreolštině.

Centrum zpracování přirozeného jazyka
Fakulta informatiky
Masarykova univerzita

