

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

# Corpus-Based Knowledge

*Otázka:* **Kdo ví víc?**

*Odověď:* **Bůh ví  $4,2578537320 \times$  častěji než čert.**

Vytváříme obrovské soubory textů (jazykové korpusy), ze kterých se lingvisté o jazyku leccos dozvídají: která slova se používají často, málo často, často spolu, která slova jsou nová a co znamenají. Vývoj software, který umožňuje indexovat miliardy slov a pomocí regulárních výrazů v nich rychle vyhledávat, je informatická výzva.

## vědět

czTenTen12 [Majka] frekvence = 5931659 (1091.0 v miliónu)

<u>has_subj</u>	<u>440718</u>	<u>-8.1</u>	<u>post_prep</u>	<u>265921</u>	<u>-4.5</u>	<u>coord</u>	<u>91437</u>	<u>-1.7</u>
fakt	<u>19690</u>	9.3	o	<u>165643</u>	8.1	znát	<u>6088</u>	7.24
bůh	<u>16942</u>	8.83	od	<u>8189</u>	4.96	tušit	<u>830</u>	6.81
moc	<u>29095</u>	8.68	z	<u>18031</u>	4.65	tápat	<u>332</u>	6.67
čert	<u>3979</u>	8.03	díky	<u>643</u>	4.36	chtít	<u>2209</u>	6.41
dávno	<u>3580</u>	7.8	podle	<u>1554</u>	3.79	neznat	<u>904</u>	6.41
rada	<u>9698</u>	7.29	prostřednictvím	<u>231</u>	3.7	vědět	<u>4488</u>	6.23
člověk	<u>31112</u>	7.23	kvůli	<u>416</u>	3.67	umět	<u>1462</u>	6.1
málo	<u>3868</u>	7.15	za	<u>4609</u>	3.53	rozumět	<u>936</u>	6.08
prd	<u>1625</u>	6.83	s	<u>9095</u>	3.5	chápat	<u>1053</u>	5.88

Centrum zpracování přirozeného jazyka  
Fakulta informatiky  
Masarykova univerzita

