

# Araneum Bohemicum Maximum: Velký český webový korpus

**Datum konání:** 28. 4. 2017

**Místo konání:** Filozofická fakulta, budova G, učebna G13

**Název přednášky:** Araneum Bohemicum Maximum: Velký český webový korpus

**Přednášející:** Ing. Vladimír Benko, PhD.

**Počet účastníků:** 18

**Zpracoval:** Radek Kopečný

Dne 28. 4. 2017 v počítačové učebně G13 FF MU jsme se zúčastnili již třetí přednášky v rámci jarního přednáškového cyklu externích odborníků. Nyní byl hostem Ing. Vladimír Benko, PhD., který nám nejdříve představil korpus Araneum Bohemicum Maximum, a v druhé části přednášky jsme si mohli sami vyzkoušet práci s korpusy ve Sketch Engine.

## Ing. Vladimír Benko, PhD.



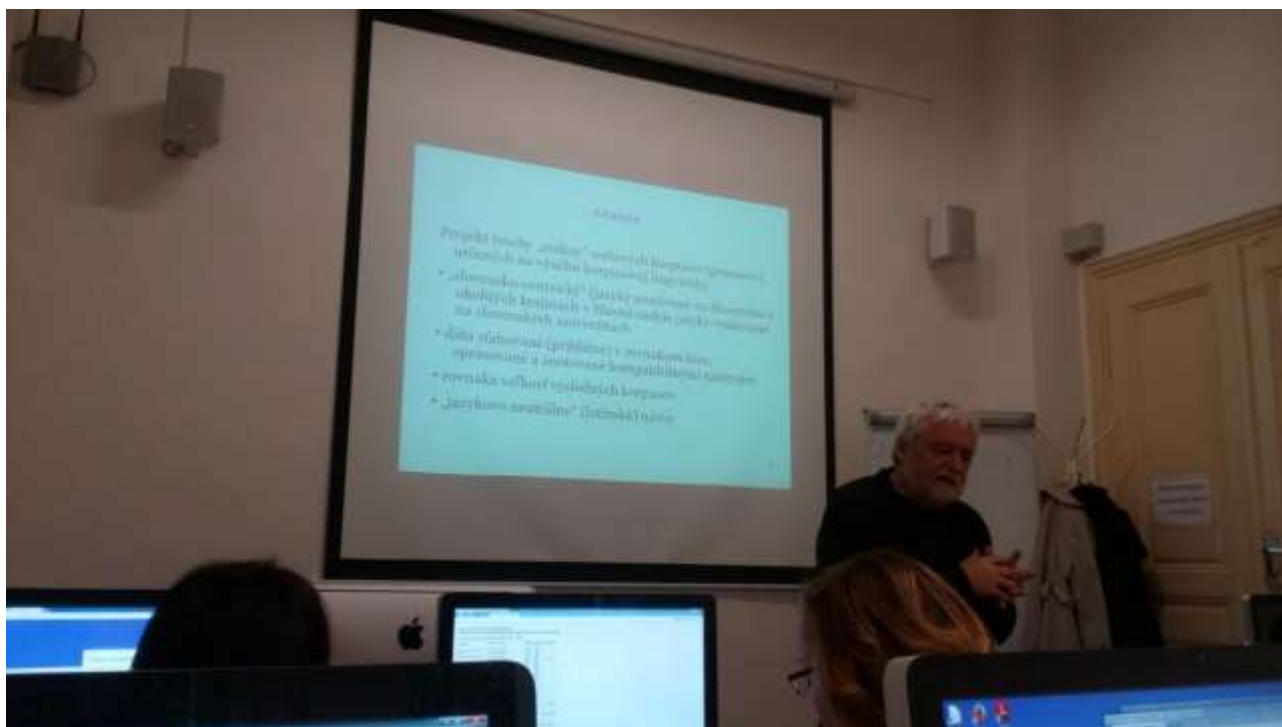
Ing. Vladimír Benko, PhD., působí jako odborný pracovník v Jazykovědném ústavu Ludovíta Štúra na Slovenské akademii věd v Bratislavě. Mezi jeho projekty patří *Slovník súčasného slovenského jazyka – 5. etapa (Koncipovanie a redigovanie slovníkových hesiel a s tým spojený lexikologicko-lexikografický výskum)* a korpus Aranea, na kterém, jak jsme se dozvěděli, pracuje sám. Mezi jeho oblasti výzkumu patří také morfosyntaktická anotace, WSG (Word Sketch Grammar) a počítačová lexikologie a lexikografie.

## První část přednášky

Ing. Vladimír Benko, PhD., nám nejdříve představil historii korpusů. Vůbec první korpus

na světě – Brown Corpus – vznikl v 60. letech minulého století a obsahoval milion slov, zachycuje americkou angličtinu. Dalším korpusem byl britský LOB s podobnou strukturou jako Brown Corpus, jen se zaměřením na britskou angličtinu. Z českých nám byly představeny například korpusy řady SYN nebo TenTen.

Poté se Ing. Vladimír Benko, PhD., zaměřil na představení toho, jak vzniká webový korpus, jenž obsahuje články z internetu. Nejprve dojde ke stažení dat z internetu a provede se konverze těchto dat na text (odstranění značek HTML a skriptů). Následně dochází k odstranění „orámování“ textu (boilerplate) a identifikaci a normalizaci kódování. Poté je identifikován samotný jazyk a cizojazyčné texty jsou vyfiltrovány. Následnou „deduplikací“ dojde k odstranění duplicitních textů a také dochází k filtraci defektních textů (absence diakritiky, nesprávně interpretované kódování, ...).



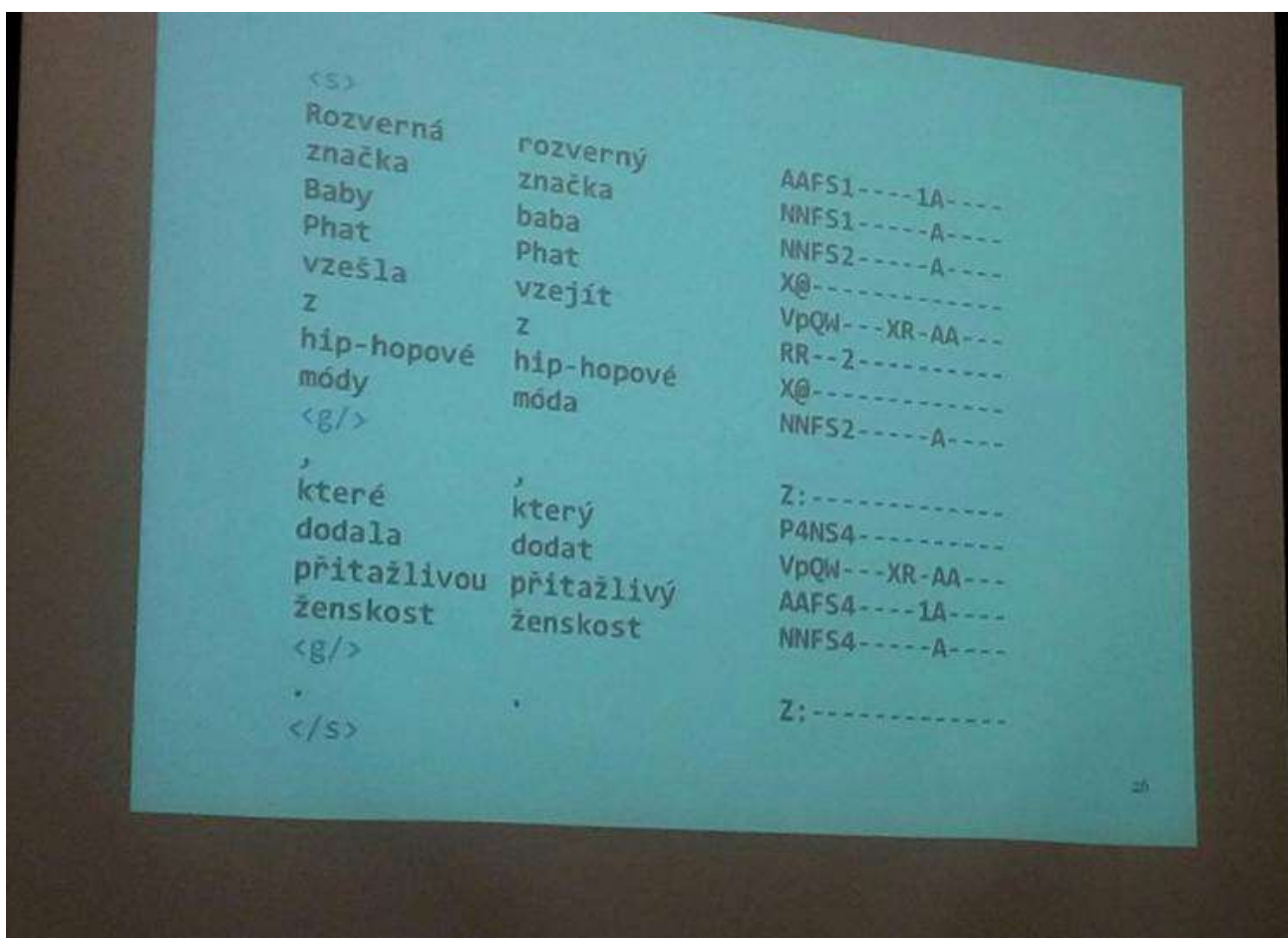
Po uvedení webových korpusů obecně jsme se dostali k samotnému projektu Aranea, jenž je projektem tvorby „rodiny“ webových korpusů primárně určených k výuce počítačové lingvistiky. Aranea je tzv. „slovensko-centrický“, zaměřený hlavně na jazyky používané na Slovensku a v okolních zemích + cizí jazyky, jež jsou vyučovány na slovenských univerzitách. Stahovaná data pro jednotlivé jazyky probíhají zhruba ve stejném období, jsou zpracovávána a anotována kompatibilními nástroji a je zde snaha o stejnou velikost výsledných korpusů. Zajímavostí je používání latinských jazykově neutrálních názvů.

Araneum Bohemicum Maximum patří k největším korpusům v projektu. Stahování dat

probíhalo v pěti obdobích – 2013, 2x 2015 a 2x 2016 pomocí nástroje SpiderLing.

Program Unitok následně provádí tokenizaci (rozčlenění textu složeného z písmen, mezer a interpunkčních znamének na tokeny). Morfologickou analýzu textů provádí program MorphoDita. Závěrem se texty segmentují na věty, filtrují, provádí se konverze tagsetu Evy Hajičové na tagset AUT (Araneum Universal Tagset) a k tomu používá Ing. Vladimír Benko, PhD., vlastní nástroje.

Jak taková analýza probíhá, jsme si ukázali na větě: „Rozverná značka Baby Phat vzešla z hip-hopové módy, které dodala přitažlivou ženskost.“



## Druhá část – workshop

Druhá část přednášky měla charakter workshopu. Vyzkoušeli jsme si sami, jak pracovat s korpusem Aranea a poté i Sketch Engine. Korpusy Aranea <http://aranea.juls.savba.sk/> jsou veřejně přístupné pod účtem hosta s jistými omezeními, ale dostali jsme možnost se registrovat pro plný přístup a Ing. Vladimír Benko, PhD., nám naši registraci potvrdil.

## Collocation candidates

Page 1  Go [Next >](#)

	<u>Cooccurrence</u> <u>count</u>	<u>Candidate</u> <u>count</u>	<u>T-score</u>	<u>MI</u>	<u>logDice</u>
<a href="#">P</a>   <a href="#">N</a> šálek	194	890	13.925	11.913	9.694
<a href="#">P</a>   <a href="#">N</a> zrnkový	77	94	8.774	13.823	8.519
<a href="#">P</a>   <a href="#">N</a> pítí	78	3,413	8.810	8.659	7.970
<a href="#">P</a>   <a href="#">N</a> ranní	75	4,145	8.633	8.322	7.813
<a href="#">P</a>   <a href="#">N</a> pražený	45	283	6.706	11.458	7.705
<a href="#">P</a>   <a href="#">N</a> pít	72	5,736	8.447	7.795	7.558
<a href="#">P</a>   <a href="#">N</a> zelený	112	14,256	10.507	7.119	7.447
<a href="#">P</a>   <a href="#">N</a> mletý	39	951	6.236	9.503	7.368
<a href="#">P</a>   <a href="#">N</a> plantážní	30	36	5.477	13.848	7.172
<a href="#">P</a>   <a href="#">N</a> milovník	47	4,185	6.821	7.634	7.134
<a href="#">P</a>   <a href="#">N</a> instantní	30	343	5.474	10.596	7.108
<a href="#">P</a>   <a href="#">N</a> odpolední	36	3,040	5.971	7.711	6.908
<a href="#">P</a>   <a href="#">N</a> cibetkový	23	31	4.795	13.680	6.789
<a href="#">P</a>   <a href="#">N</a> ledový	30	2,761	5.449	7.587	6.687
<a href="#">P</a>   <a href="#">N</a> příprava	84	20,455	9.039	6.183	6.659
<a href="#">P</a>   <a href="#">N</a> italský	37	6,379	6.023	6.681	6.525
<a href="#">P</a>   <a href="#">N</a> výborný	50	11,726	6.977	6.237	6.468
<a href="#">P</a>   <a href="#">N</a> lahodný	21	1,042	4.570	8.478	6.458
<a href="#">P</a>   <a href="#">N</a> čerstvý	38	7,396	6.097	6.506	6.456
<a href="#">P</a>   <a href="#">N</a> uvařit	22	1,790	4.669	7.764	6.394
<a href="#">P</a>   <a href="#">N</a> silný	77	24,604	8.616	5.791	6.329

Na lemmatu káva jsme si ukázali jeden ze způsobů využití korpusu – hledání kolokací.

milovat/nenávidět Araneum Bohemicum III Maius (Czech, 17

milovat 6.0 4.0 2.0 0 -2.0 -4.0 -6.0 nenávidět

X/Y , X/Y	13,878	1,781	0.13	0.12	Y X	48,229	6,980	0.45	0.48
zachovávat	<a href="#">59</a>	0	7.0	--	bezmezně	<a href="#">240</a>	0	7.3	--
žít	<a href="#">138</a>	0	6.9	--	svobodně	<a href="#">231</a>	0	7.1	--
meditovat	<a href="#">50</a>	0	6.8	--	nadevše	<a href="#">186</a>	0	7.0	--
odpouštět	<a href="#">39</a>	0	6.4	--	heidi	<a href="#">155</a>	0	6.7	--
toužit	<a href="#">41</a>	0	6.3	--	tajně	<a href="#">108</a>	0	6.0	--
ctít	<a href="#">31</a>	0	6.1	--	doopravdy	<a href="#">111</a>	0	5.8	--
obdivovat	<a href="#">33</a>	0	6.1	--	vroucně	<a href="#">82</a>	0	5.8	--
dávat	<a href="#">63</a>	0	6.1	--	bezpodmínečně	<a href="#">69</a>	0	5.4	--
věřit	<a href="#">70</a>	0	6.0	--	smírně	<a href="#">84</a>	0	5.3	--
cítit	<a href="#">61</a>	0	5.9	--	prostě	<a href="#">696</a>	<a href="#">28</a>	5.9	1.4
respektovat	<a href="#">31</a>	0	5.9	--	opravdu	<a href="#">796</a>	<a href="#">43</a>	5.5	1.4
poznat	<a href="#">53</a>	0	5.9	--	skutečně	<a href="#">387</a>	<a href="#">22</a>	5.5	1.6
radovat	<a href="#">27</a>	0	5.8	--	bůh	<a href="#">685</a>	<a href="#">59</a>	6.8	3.5
přijímat	<a href="#">32</a>	0	5.8	--	tolik	<a href="#">513</a>	<a href="#">72</a>	6.4	3.9
trpět	<a href="#">33</a>	0	5.7	--	hluboce	<a href="#">158</a>	<a href="#">16</a>	6.2	4.2
dělat	<a href="#">235</a>	<a href="#">15</a>	6.5	2.8	přestat	<a href="#">179</a>	<a href="#">33</a>	5.3	3.4
znát	<a href="#">91</a>	<a href="#">10</a>	6.0	3.2	doslova	<a href="#">187</a>	<a href="#">41</a>	5.7	4.2
dokázat	<a href="#">85</a>	<a href="#">9</a>	6.2	3.5	strašně	<a href="#">148</a>	<a href="#">33</a>	5.7	4.4
milovat	<a href="#">1,222</a>	<a href="#">142</a>	10.5	8.2	vášnivě	<a href="#">109</a>	<a href="#">9</a>	6.1	5.0
snášet	<a href="#">58</a>	<a href="#">20</a>	6.8	7.2	navzájem	<a href="#">124</a>	<a href="#">38</a>	5.7	5.1
nenávidět	<a href="#">142</a>	<a href="#">108</a>	8.2	10.0	upřímně	<a href="#">93</a>	<a href="#">104</a>	5.5	7.0
závidět	0	<a href="#">12</a>	--	6.9	duše	<a href="#">40</a>	<a href="#">212</a>	3.5	6.6
odsuzovat	0	<a href="#">12</a>	--	7.0	svorně	0	<a href="#">18</a>	--	6.0
proklínat	0	<a href="#">10</a>	--	7.3	spam	0	<a href="#">45</a>	--	7.0
pohrdat	0	<a href="#">12</a>	--	7.4	bylostně	0	<a href="#">56</a>	--	7.6

Následně jsme si ukazovali a vyzkoušeli nástroj Sketch Engine (dále SE), jenž vyvíjí společnost Lexical Computing Ltd. ve spolupráci s Centrem zpracování přirozeného jazyka Fakulty informatiky Masarykovy univerzity, a proto k němu máme bezplatný přístup. SE slouží k vytváření

korpusů a vyhledávání v nich, má přístup k mnoha korpusům ve více než 80 jazycích a například Aranea mezi ně též patří. Oproti běžným nástrojům má rozšířené možnosti, například Thesaurus (spojení slov na základě jejich podobnosti ve významu – SE ještě vyobrazí jejich frekvenci výskytu v konkrétním korpusu). Právě díky mnoha jazykům je SE velmi užitečný nástroj pro překladatele, neboť umožňuje vyhledávat kolokace paralelně, dělí je dle četnosti použití, dle pozic (za slovem, nebo před slovem) a také dle slovních druhů. Zajímavé je také vyhledávání dvojic slov a možnost podívat se, se kterými slovy se více spojuje první z dvojice, se kterými to druhé a která slova mají společná (viz obrázek).

Na závěr byl prostor pro dotazy na samotný projekt Araneum i korpusy obecně. Přednáška byla velmi zajímavá a díky širokému využití korpusů si přišli na své i studenti jiného oboru než počítačové lingvistiky.