

EXEMPLIFIKACE SUBSTANTIV V LEXIKÁLNÍ DATABÁZI PRALEX (ZPRACOVÁNÍ PŘÍKLADOVÉ ČÁSTI HESLA)

Marta Koutová, Ústav pro jazyk český AV ČR, v. v. i., oddělení současné lexikologie a lexikografie

ÚVOD

Podrobné zpracování exemplifikace je těžištěm naší práce v rámci výzkumného záměru *Vývoje databáze lexikální zásoby českého jazyka počátku 21. století*.

Základní materiálové východisko: Český národní korpus ÚČNK

Používaný korpus: ORIG_SYN (dříve pod názvem SYN – obsahuje korpusy SYN2000, SYN2005 a SYN2006PUB). Zahnuje současně psané texty z oblasti publicistiky, beletrie i odborné literatury.

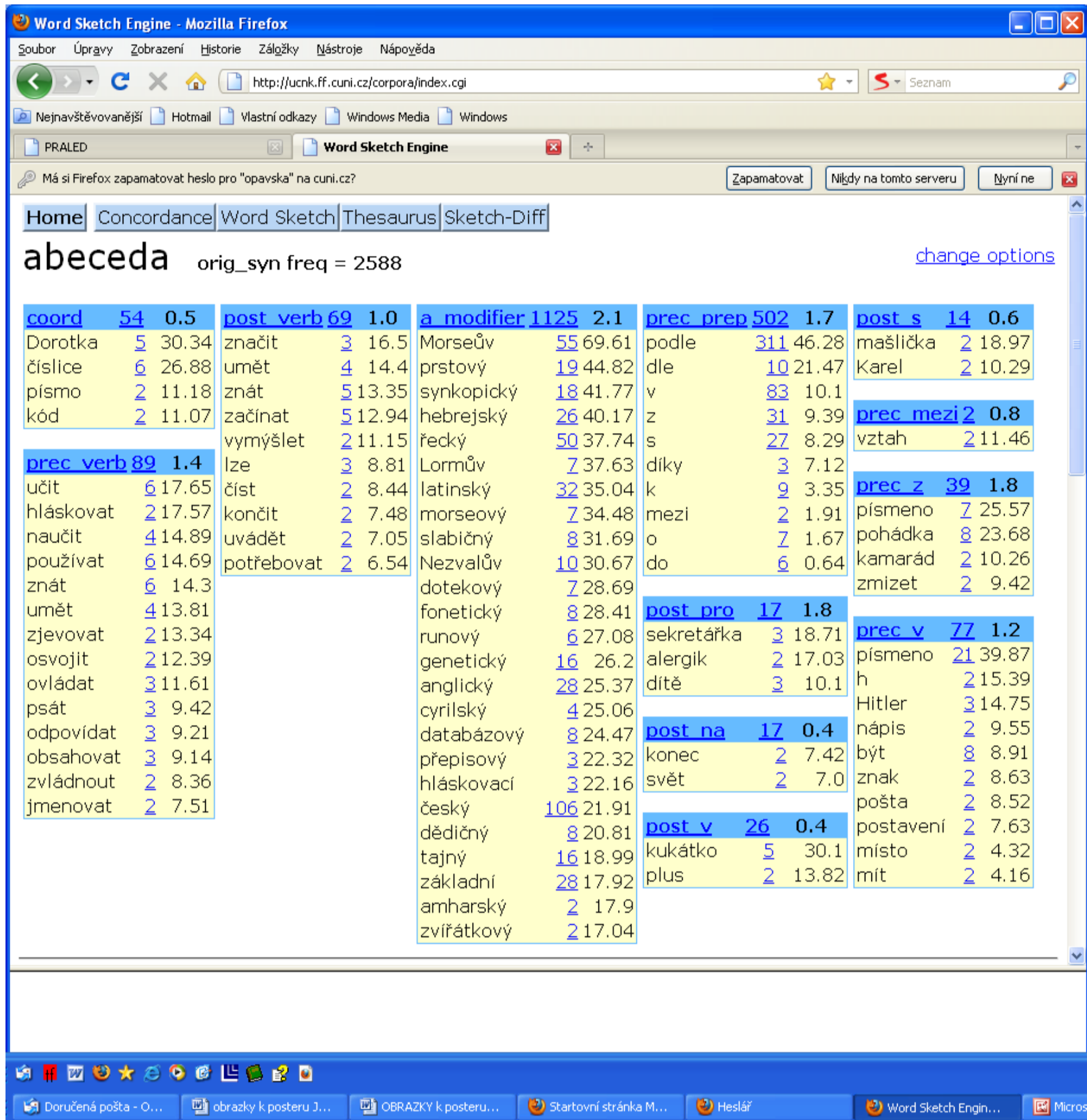
Cíl exemplifikace v LDB PRALEX: zachycení spojitelnosti heslových slov v příkladové části hesla – výběr relevantních dokladů, tj. příkladů na konkrétní užití k předem založeným významům podle SSJC, resp. dalších výkladových slovníků + zachycení dokladů k novým významům.

1 ZACHYCNÍ SPOJITELNOSTI HESLOVÝCH SLOV

- U substantiv se zaměřujeme na syntaktickou i sémantickou spojitelnost.
- Celou škálu spojitelnosti ilustrujeme primárně pomocí necitátových, tzv. upravených dokladů, tj. příkladových spojení ve formě syntagmat, které představují běžné kolokace zpracovávaných substantiv.
- V menší míře uvádíme i větné doklady (citátové).

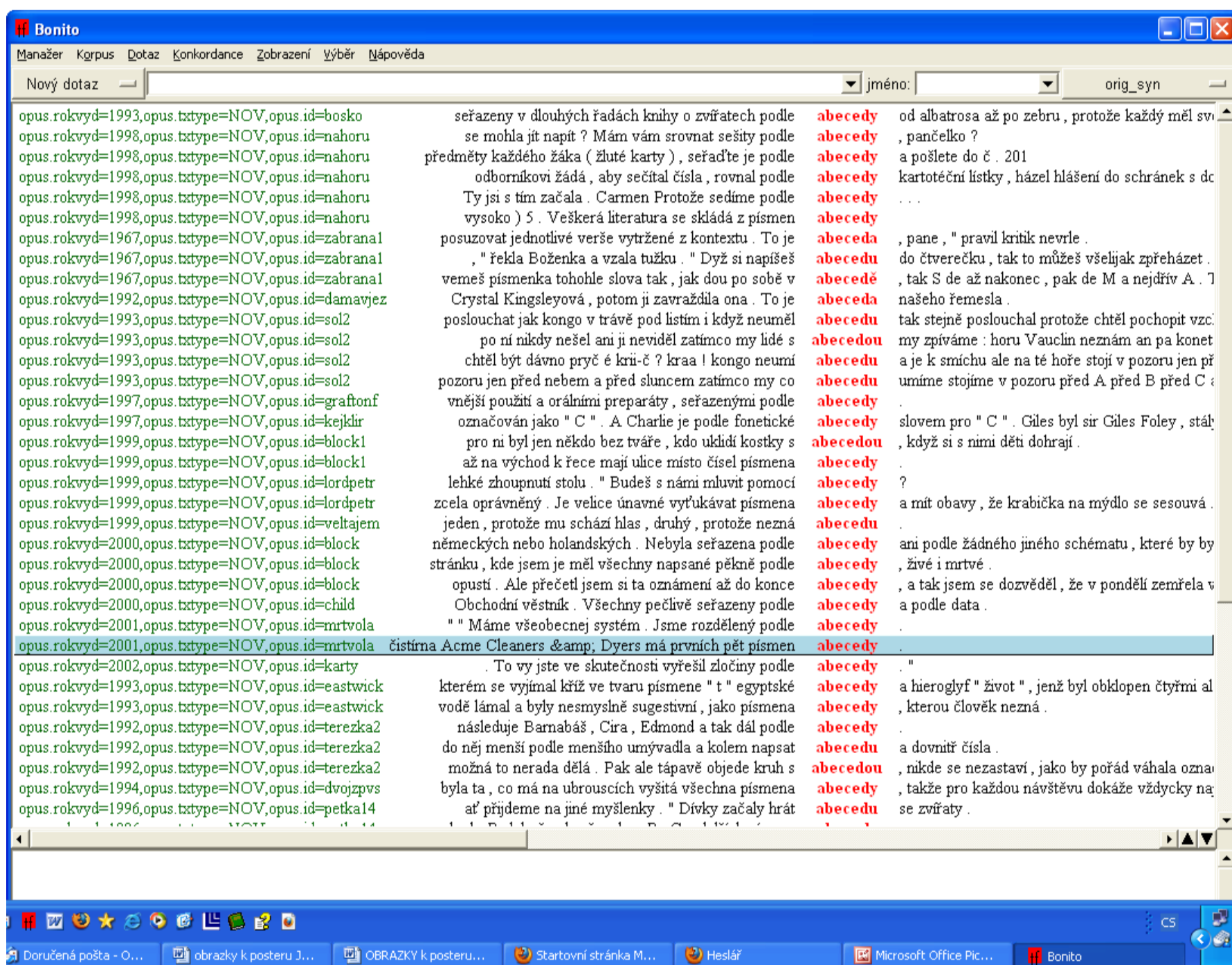
2 PRÁCE S WSE A BONITEM

- Pro třídění dokladů z korpusu využíváme nástroj Word Sketch Engine (WSE) a Bonito.
- WSE používáme pro výběr a podrobné třídění minimálních kontextů k jednotlivým významům.



obr. 1: nástroj WSE

- Z Bonita vybíráme tzv. citátové doklady – v průměru 3–5 větných příkladů s uvedením zdroje.



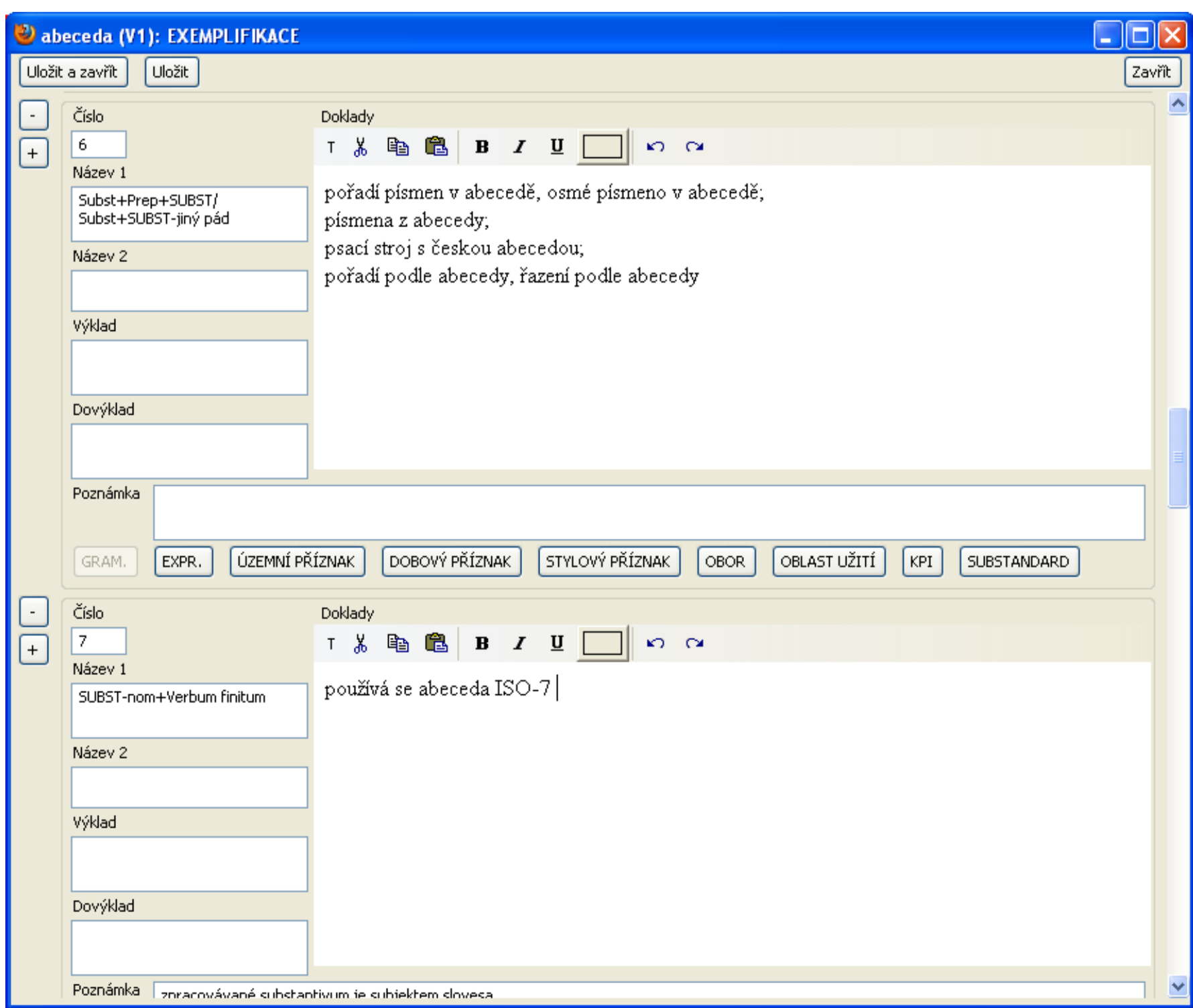
obr. 2: Bonito

3 PRÁCE S KORPUSOVÝMI DOKLADY

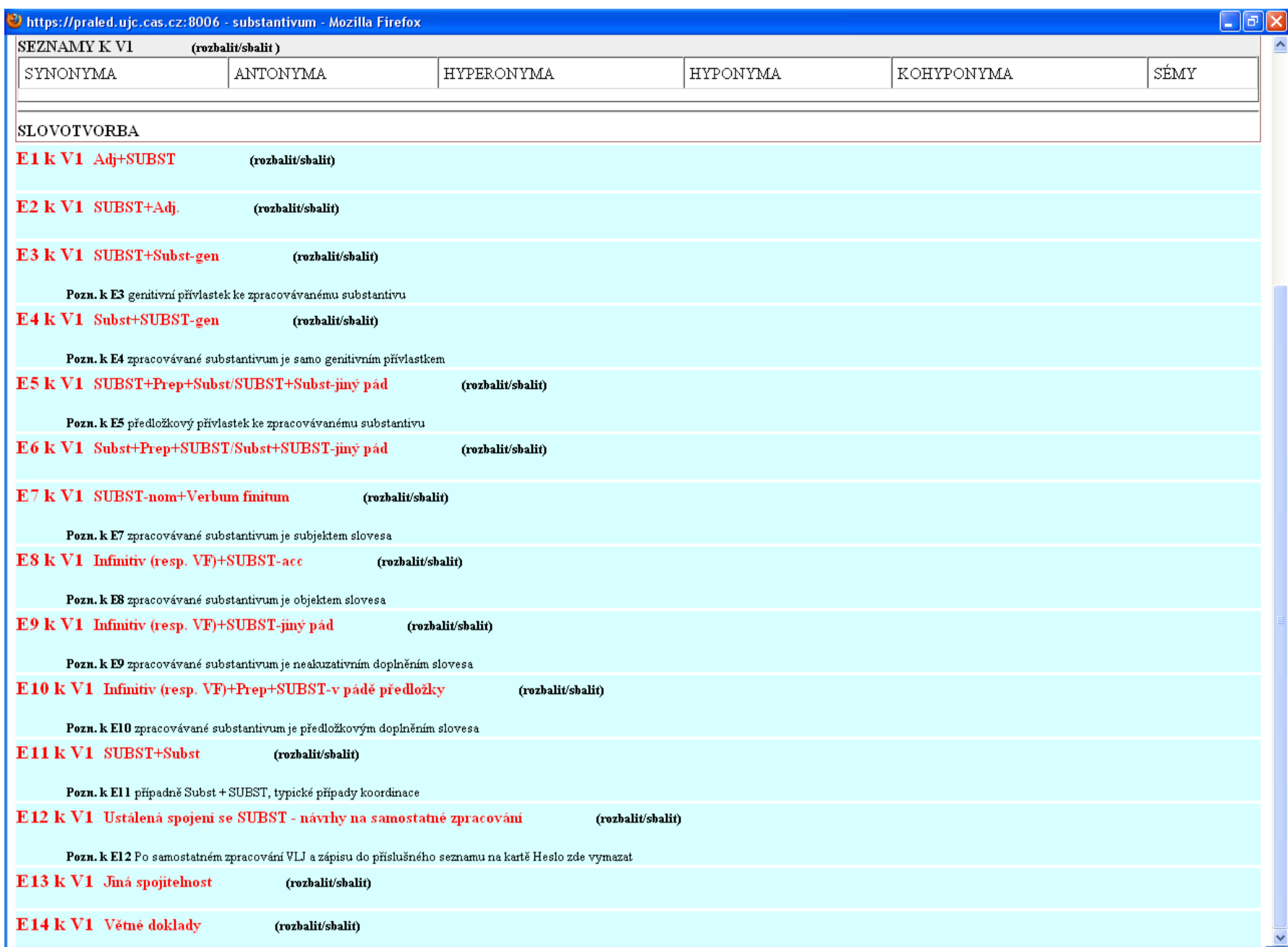
- Vybíráme nejen spojení zcela typická, ale abychom dostali širší obraz sémantické spojitelnosti daného hesla, vybíráme i spojení běžná, v úzu hojně zastoupená.
- V rámci příslušného bloku postupujeme částečně na základě frekvence, částečně na základě sémantiky tak, aby byly vešle uvedeny kolokace z řady lexikálně-sémantických skupin. Pokud nenajdeme žádné relevantní doklady, příslušný blok zůstane prázdný.
- Rozvíjíme minimální syntagmat; v některých případech nestačí uvádět minimální syntagmata, ale je nutné uvádět rozvířejší doklady, aby příklad dával smysl, popř. aby se přesně vymezilo, ko kterému významu daný příklad patří, srov. např.:
 - autodlna: autodlna v Mariánských úlicích (nikoli pouze „autodlna v ulici“); houslista: houslista orchestru Národního divadla (nikoli pouze „houslista orchestru“); kachna: bránit se tiskovým kachnám (k V3: nepravdivá, vymyšlená, senzační zpráva).*
- Někdy sice není nezbytné nutně rozvíjet uvádět, avšak pokud vhodné dokresluje dané spojení a v úzu se běžně vyskytuje, můžeme uvést i rozvířejší doklady tohoto typu:
 - kontrola: jít na kontrolu k lékaři; auto: slabost pro rychlá auta; sřít se protijedoucím autem*

4 CELKOVÁ STRUKTURA EXEMPLIFIKACE U SUBSTANTIV

- Příkladovou část hesla rozdělujeme do 14–15 exemplifikačních bloků. V členění do bloků se projevuje syntaktická spojitelnost daného substantiva, uvnitř jednotlivých bloků je zachycena spojitelnost sémantická.



obr. 3: Nástroj Exemplifikace v LDB PRALEX (část formuláře pro zápis dat)



obr. 4: Náhled Exemplifikace v LDB PRALEX (před vyplněním)

- Do exemplifikačních bloků E1–E13 uvádíme výhradně tzv. upravené doklady, tj. minimální nebo někdy i rozvířejší syntagmata, která získávají především pomocí WSE. Pro větné doklady slouží blok E14. Fakultativní blok E15 slouží pro doklady na lemma užité v názvech a přidáváme ho jen v případě potřeby.



obr. 5: Exemplifikace (příkladová část) hesla *abeceda* v LDB PRALEX

5 VNITŘNÍ ČLENĚNÍ V RAMCI JEDNOHO EXEMPLIFIKAČNÍHO BLOKU

- V každém exemplifikačním bloku uvádíme nejprve doklady na běžné, nepřízvučné užití. Pokud u nějakého spojení převládá užití v plurálu, uvádíme doklad v této podobě, srov. např.:
 - automobil: majitel automobilu, řidič automobilu*
 - obyvatel: místní obyvatelé, zdejší obyvatelé, tamní obyvatelé*
 - podpatek: vysoký podpatek i vysoké podpatky*
- Pokud v korpusu najdeme i doklady na užití odborné, přenesené apod., uvádíme je až po „neutrálních“ dokladech.
- Doklady na odborné užití signalizujeme titulkem „odb.“ (zatím nejde o oborové kvalifikátory, ale jen o vymezení jasnějších případů zvlášť), srov. např. u hesla *mezera*:
 - vzdálová mezera;*
 - větrná mezera;*
 - vohná mezera, zejčí mezera;*
 - úzká mezera, široká mezera;*
 - dvoucíselná mezera;*
 - odb. dilatační mezera*
- Ve výjimečných případech používáme v titulku i konkrétnější označení dané oblasti, např. „sport.“ aj.
- Doklady na přenesené užití registrujeme pod titulkem „přen.“, srov. např. u hesla *čtěník*:
 - přen. role čtěníka Balkanu*
- Podobně u hesla *kometka*:
 - přen. zářivá kometka zápasu; vyhoopnout se mezi komety letošního roka.*
- Někdy může být určité spojení užitév v přímém i přeneseném významu. V těchto případech dáváme za takový doklad do závorky komentář „(i přen.)“.
- Je-li spojení se substantivem užíváno jako neustálené (v přímém významu) i jako frazém, případně soulovi, pak uvádíme v závorce komentář „(i fraz.)“, případně „(i sousl.)“:
 - podpatek: souzří podpatky (i fraz.); odskočnout podpatky (i fraz.); svuknout podpatky (i fraz.); skáka: býti káň (i fraz.); čerň káň (i fraz.); křák: černá skřípka (i sousl.).*

6 SPECIFIKA A PŘÍKLADY ZPRACOVÁNÍ JEDNOTLIVÝCH BLOKŮ EXEMPLIFIKACE U SUBSTANTIVNÍCH HESEL

Blok E1 (Adj+SUBST)

- První blok zachycuje spojení zpracovávaného substantiva s anteposovaným adjektivním přívlaskem. V tomto bloku zachycujeme oddělené doklady na singulárové a plurárové užití zpracovávaného substantiva podle toho, co u daných spojení převládá. Jako příklad uvedeme heslo *čeština*:
 - spisovná čeština, obecná čeština, hovorová čeština;*
 - stará čeština, archaická čeština, obrozenecká čeština, velevslavská čeština, středověká čeština;*
 - současná čeština, dnešní čeština, novodobá čeština;*
 - lámaná čeština, špatná čeština, dobrá čeština, plynná čeština, perfektní čeština, dokonalá čeština;*
 - dobře srozumitelná čeština, slušná čeština, bezchybná čeština; krásná čeština, kultivovaná čeština; jadrná čeština;*
 - mluvněná čeština, psaná čeština;*
 - auto: česká čeština*

Blok E2 (SUBST+Adj)

- Do tohoto bloku dáváme výhradně odborná spojení zpracovávaného substantiva s adjektivem v postpozici typu *komalínka vonná, kyselina sírová, atribut verbální*. Nepatří sem případy, kdy je postpozice ovlivněna stylisticky nebo aktuálním větným členěním, srov. např. u hesla *ozon*: *troupičkový ozon má vůči UV-záření stejné chování jako ozon stratosférický*.

Blok E3 (SUBST+Subst-gen)

- Do bloku E3 vybíráme příklady na spojení daného substantiva s jiným substantivem v genitivu. Snažíme se přitom o oddělení objektových a subjektových genitivů:
 - autocenzura textů x autocenzura redaktorů;*
 - makrofotografie hmyzu x makrofotografie Jiřího Štreita*

Příklad zpracování – heslo *čeština*:
čeština emigrantů, čeština cizinců;
čeština Karla Hynka Máchy

Blok E4 (Subst+SUBST-gen)

- v bloku E4 je zpracovávané substantivum samo genitivním přívlaskem.

Příklad zpracování – heslo *čeština*:
znalost češtiny; vjuka češtiny, kurs češtiny, studium češtiny, hodina češtiny;
učitel češtiny, profesor češtiny, lektor češtiny pro cizince, uživatel češtiny;
učivání češtiny, výslovnost češtiny;
slovník češtiny, mluvnice češtiny, učebnice češtiny

Blok E5 (SUBST+Prep+Subst / SUBST+Subst-jiný pád)

- Blok E5 zachycuje příklady ke zpracovávanému substantivu. Do tohoto bloku patří rovněž doklady na spojení zpracovávaného substantiva se substantivem v jiném pádě, srov. např. u hesla *dárek*:
 - dárek dětem, dárek dětem.*
- Vazby se stejnou předložkou a stejným pádem dáváme k sobě, převládaje zde tedy formální kritérium nad sémantickým.
- Pokud se však u některých dokladů předložka pojí se stejným pádem, ale sémantika je odlišná, pak rozdělujeme dané doklady na více řádků. Srov. např. u hesla *automa*:
 - automa na nápoje, automa na kávu, automa na jízdenky;*
 - automa na mince, automa na karty*

Příklad zpracování – heslo *čeština*:
čeština pro cizince, čeština pro samouky;
čeština na univerzitě, čeština na gymnáziu;
čeština ve škole;
čeština se silným přízvukem

Blok E6 (Subst+Prep+SUBST / Subst+SUBST-jiný pád)

- Do bloku E6 dáváme případy, kde se zpracovávané substantivum vyskytuje jako předložkový přívlasek k jiným substantivum nebo se s jinými substantivy pojí jako přívlasek dativní či instrumentální (např. *arterie: přávek, arterie*).
- Příklad zpracování – heslo *čeština*:
 - překladatel z češtiny; překlad z češtiny do němčiny;*
 - zkouška z češtiny, maturita z češtiny, píseňka z češtiny, díkání z češtiny, známka z češtiny;*
 - učitelka na češtinu;*
 - návod v češtině, text v češtině, nápis v češtině, vysílání v češtině;*
 - vztah k češtině, láska k češtině;*
 - prohřešek proti češtině;*
 - zájem o češtinu;*
 - učí pro češtinu;*
 - problémy s češtinou, zachzení s češtinou*

Blok E7 (SUBST-nom+Verbum finitum)

- V bloku E7 je zpracovávané substantivum subjektem slovesa.

Příklad zpracování – heslo *čeština*:
ze všech stran zněla čeština;
úředním jazykem je čeština

Blok E8 (Infinitiv (resp. VF)+SUBST-acc)

- V bloku E8 je dané substantivum objektem slovesa.

Příklad zpracování – heslo *čeština*:
učit češtinu;
studovat češtinu, učit se češtinu;
ovládat češtinu, užívat češtinu;
šlyšet češtinu;
komolít češtinu, lámat češtinu, prznit češtinu

Blok E9 (Infinitiv (resp. VF)+SUBST-jiný pád)

- Blok E9 zachycuje případy, kdy zpracovávané substantivum je neakuzativním doplňním slovesa, tj. vyskytuje se v genitivu, dativu či instrumentálu.

Příklad zpracování – heslo *čeština*:
učit se češtině, rozumět češtině;
mluvit spisovnou češtinou, hovořit lánanou češtinou, vládnout krásnou češtinou;
knihu je psána krásnou češtinou

Blok E10 (Infinitiv (resp. VF)+Prep+SUBST-v pádě předložky)

- V bloku E10 je zpracovávané substantivum předložkovým doplňním slovesa.
- Slovesa v blocích E8–E10 uvádíme většinou v infinitivním tvaru. Pokud je ale důležitý původce děje, je nutné uvádět daná slovesa ve finitivním tvaru, srov. např.:
 - veřej: vrata vzrvala ve veřejích, dveře skřípou ve veřejích (tj. nikoli pouze skřípá ve veřejích, vzrval ve veřejích).*

Příklad zpracování – heslo *čeština*:
překládat knihu do češtiny, překládat do češtiny, tlumočit do češtiny, převést text do češtiny;
překládat z češtiny, přecházet z češtiny do angličtiny, maturovat z češtiny;
knihu vyšla v češtině;
zápasit s češtinou, válčit s češtinou, být s češtinou na štitu

Blok E11 (Subst+Subst, případně Subst + SUBST)

- V bloku E11 je našim díkolem uvádíme typické příklady koordinace. Zapisujeme sem především koordinace se spojkami *a* nebo (protože jím nám WS nezobrazí), popř. též se spojkou *i* nebo s předložkou *s*, jako např. u lemmatu *otec*: *otec s matkou, matka s otcem*.

Příklad zpracování – heslo *čeština*:
čeština a angličtina, čeština a slovenština; směsice češtiny a němčiny

Blok E12 (Ustálené spojení se SUBST – návrh na samostatné zpracování)

- Tento blok není na rozdíl od bloků E1–E11 a E13 vymezen syntakticky. Je to blok, do kterého zpracovatelé substantiv zatím pracovně zapisují ta ustálená spojení se zpracovávaným substantivem, která identifikují v korpusu a na základě lexikografického kritéria potřebnosti výkladu navrhuji k samostatnému zpracování.

Blok E13 (Jiná spojitelnost)

- Tento blok zatím slouží jako sběrný koš pro případy, které se syntakticky nehodí do žádného jiného bloku. Ukázky různých spojení, která patří do E13, lze nalézt např. u hesla *puchýř*:
 - nudy plní puchýře; puchýře naplněné vodou; látky vyvolávající na kůži puchýře; opar ve stadiu puchýře; mit dlaně samý puchýři*

Blok E14 (Větné doklady)

- Na rozdíl od předchozích bloků zachycuje blok E14 celé větné doklady z korpusu, tj. příklady rozvířejší, často i takové, které díky širšímu kontextu vysvětlují a ilustrují méně jasná spojení z ostatních exemplifikačních bloků.

Příklad zpracování – heslo *čeština*:
opus.rokyid=2003,opus.tytype=SCI,opus.id=ener0302. Konference bude vedena v angličtině a < češtině>.
opus.rokyid=1993,opus.tytype=PIB,opus.id=ln493161 V Tatrách zni letos < čeština> sporadicky.
opus.rokyid=1995,opus.tytype=PIB,opus.id=mf95102 V < češtině> zni vaše jméno podobně jako kafe.

Blok E15 (V názvech)

- Do bloku E15 dáváme příklady na užití daného substantiva v nejrůznějších názvech. Tento blok v základní době exemplifikační struktury substantiv není, přidáváme jej fakultativně – pouze v případě, že nějaké relevantní doklady v korpusu nalezneme.

Příklad zpracování – heslo *čeština*:
Nová slova v češtině; slovník neologizmů;
Přiruček mluvnice češtiny;
Slovník spisovné češtiny pro školu a veřejnost

7 ČLENĚNÍ NA VÝZNAMY NA ZÁKLADĚ PRIMÁRNÍCH A REFERENČNÍCH ZDROJŮ

Vyhodnocení celé sémantické struktury polysémních hesel a nové formulace příslušných definic budou až takolem něho výkladového slovníku pro potřeby LDB Pralex zatím pracovně používáme výklady ze Slovníku spisovného jazyka českého (SSJC) a případně i z dalších výkladových slovníků. Pokud v korpusu najdeme doklady pro významy nezachycené ve slovníčích, nebo zaznamenáme nějaké změny vůči předchozímu zpracování v nich, založíme další význam. Polysémní hesla tedy primárně členíme na významy ve shodě se SSJC, při práci s korpusem a dalšími slovníky někdy doplňujeme další významy.

- Např. u substantiva *delegát* jsme zachytili na základě výskytu v korpusu ORIG_SYN další dva významy:
 - V2: *pracovní cestovní kancelář pečující o rekreaty v místě jejich (zahraničního) pobytu (delegát cestovní kanceláře, delegát Cedokey);*
 - V3: *osoba vyslaná fotbalovým svazem a pověřená dohledem nad regularitou hodnocení utkání (fotbalový delegát; svazový delegát; delegát zápasu, delegát utkání; delegát svazu).*
- Zpracování v SSJC porovnávané se zpracováním v SSC, ASCS aj., a pokud tyto slovníky zachycují oproti SSJC další významy, v relevantních případech na základě těchto nových výkladových slovníků a doplňují je v korpusu významovou strukturou zpracovávaného hesla doplňujeme.
- Další význam někdy vymezujeme na základě dokončeného procesu lexikalizace, která byla v SSJC jen naznačena. Pokud bylo registrováno v SSJC nějaké přenesené užití nebo významový odstín a dnes je tento významový posun v korpusu hojně doložen, může to být signálem, že se ustáluje další samostatný význam. Srov. např. u hesla *buňka* – v SSJC uvedeno u 2. významu jako „zř. přen.“ „malá místnost, jedna z mnoha stejných“; v LDB vyčleněno jako další význam (V5), „malá místnost, jedna z mnoha stejných; dočasná stavba menších rozměrů sloužící konkrétním účelům“:
 - stavební buňka, pracovní buňka, prodejní buňka;*
 - rodinná buňka, obytná buňka;*
 - sanitární buňka, servisní buňka;*
 - sklepní buňka, skladovací buňka;*
 - provizní buňka;*
 - prefabrikovaná buňka, montovaná buňka, smontovaná buňka;*
- Pokud najdeme jen několik dokladů, které signalizují nové významové posuny, na jejichž základě se zatím nemůžeme přesně vymezit, uvádíme je v poli „Doklady pro nové významy“, aby bylo možné je později ověřit dohledáním v dalších zdrojích a vyhodnotit, zda jde opravdu o další význam. Všechny tyto signály budou využití při práci na novém výkladovém slovníku.

