

Část I

**Teoretické pojednání**

# Kapitola 1

## Wordnet

### 1.1 Motivace vzniku

Od počátků snah o zpracování přirozeného jazyka (NLP, natural language processing) bylo nutné poskytnout programu data o lexiku ve zpracovávaném textu, ať už ona data byla jakákoliv. Kupříkladu pro překlad se mělo za to, že stačí ekvivalentní dvojice ve zdrojovém a cílovém jazyce, později se přidal kontext v případě statistického strojového překladu spolu s dalšími informacemi, jako je například slovní druh. Tradičně se lexikální materiál ukládá způsobem nikoliv diametrálně odlišným od papírových slovníků určených pro lidské uživatele. Ty obvykle obsahují abecedně (či podle jiného indexu) seřazené jednotlivé záznamy s potřebnými informacemi o slovech, z nichž pak program může čerpat při zpracování textu.

cit?

Jak uvádí Pala; Ševeček, (2013), uspořádání lexikálního materiálu v takovéto formě je sice vhodné pro člověka, ale nikoliv pro strojové zpracování, a to z několika důvodů. Kromě toho, že vyhledávání v abecedním seznamu je relativně pomalé, struktura tradičního slovníku kvůli onomu abecednímu řazení inherentně vzdaluje slova, jež člověk chápe jako nějakým způsobem blízká. Tato blízkost může vyplývat ze vztahu volné synonymie, antonymie, podřazenosti, nadřazenosti, etc. Pokud si tedy například uživatel výkladového slovníku nepříliš obeznámený s daným jazykem vyhledá určité heslo, dozví se sice pravděpodobně jeho význam, ale nebude schopen své znalosti prohlubovat dále zjištěním, kupříkladu jaké je slovo odpovídá opačnému významu.

nejaka citace?, dohledat neco, jak takovy slovníky byly uloženy...

Dalším všeobecným problémem při využití tradičních slovníků k počítačovému zpracování jazyka je fakt, že lexikografové předpokládají u uživatele slovníku značné encyklopedické znalosti. Zařazují tak do slovníku jen informace dle jejich názoru důležité pro rozlišení (differentia specifica) a zařazující do kontextu či přiřazující k určité nadřazené třídě objektů (genus proximum). Vyhledá-li si tedy člověk ve Slovníku spisovného jazyka českého heslo vlk, zjistí následující:

citace

**vlk:** psovitá šelma šedě (n. šedožlutě) zbarvená, žijící v Evropě, Asii a v Sev. Americe

Definice a priori předpokládá, že uživatel je obeznámen s tím, co je šelma a co je pes. Pokud takovou znalostí neslyne (což je vcelku představitelné například u cizince), je nucen si tato slova ve slovníku najít a podívat se na jejich definice (pomiňme nyní netriviální úkol převést slovo psovitá na základní tvar pes). Pokud nerozumí definicím ani nadřazených slov, musí pokračovat v hierarchii dále a dále.

Z uvedeného případu plyne, že jakkoliv je možné správným vyhledáváním hyperonym<sup>1</sup> dospět k tomu, že vlk je konkrétní entita našeho vesmíru, živá bytost o čtyřech končetinách, savec nějakým způsobem příbuzný se psovi, má šedou srst etc., je takový proces dosti komplikovaný. Příklad s cizincem se sice nemusí zdát zcela relevantní, protože se dá předpokladat, že daný člověk má, byť v jiném jazyce, stejné základní znalosti předpokládané lexikografy jako člověk český. Situace je však dramaticky jiná u počítače. Na rozdíl od člověka totiž počítač nemá žádné předchozí znalosti, tudíž musí projít celým procesem popsaným výše, aby byl schopen kupříkladu určit, že vlk může umřít (ježto je živá bytost). Protože však tradiční slovníky typu SSJČ byly vytvářené pro papírové médium, neobsahují žádné propojení ve stylu *toto je odkaz na hyperonymum*, a počítač tudíž jen těžko může zjišťovat, na které vlastně slovo se to má podívat, aby se dobral podstaty pojmu vlk.

### 1.1.1 Strojově čitelné slovníky

V zájmu automatizace vyhledávání ve slovníku vznikaly tzv. strojově čitelné slovníky<sup>2</sup>, což je pojem souhrnně označující lexikální databáze. Podle množství informací, které taková databáze obsahuje, pak lze tyto dělit na slovníky, taxonomie a ontologie. Je evidentní, že obyčejný slovník neobsahuje oproti tradičnímu papírovému slovníku navíc žádné metainformace, takže je počítač při jeho užívání v podstatě omezen na elektronický listovač (Miller et al., 1990).

Míra, jakou se strojově čitelný slovník odlišuje od pouhé zdigitalizované formy papírového slovníku a přiblíží se k pokročilé lexikální databázi, lze vyjádřit v několika stupních. V případě, že slovník má jednotlivé koncepty uspořádány v hierarchii dle nadřazenosti–podřazenosti, lze jej označit za taxonomii, tedy systém s hlubší strukturou než pouze abecedním řazením hesel.

Dalším stupněm už je skutečná lexikální databáze, která má jednotlivé koncepty propojeny rozličnými vztahy, počínaje onou základní hyperonymií a hyponymií a pokračuje kupříkladu vztahy meronymie<sup>3</sup> či antony-

neja-  
kej link,  
kde bu-  
dou kon-  
cepty/sen-  
ses vysvet-  
leny

<sup>1</sup>nadřazené slovo

<sup>2</sup>machine readable dictionary

<sup>3</sup>vztah je částí, tedy např. dveře je meronymem trolejbusu

mie<sup>4</sup>. Kromě vztahů mezi koncepty bude taková lexikální databáze obsahovat zřejmě i další informace, tedy nějaké kategorie slov, jejich popis, etc. Databáze tak popsaných konceptů propojených sémantickými vztahy může být nazývána ontologií. (GarshoI, n.d.)

### 1.1.2 Od slovníků k Wordnetu

Výše uvedená opozice papírového slovníku a ontologie ilustruje rozdíly tradičního slovníku a počítačově zpracovatelné lexikální databáze. Už ze samotného konceptu takové databáze je evidentní jeden klíčový rozdíl – tradiční slovníky, jsouce řazené abecedně, od sebe oddalují některá hesla, jež by bylo vhodné mít pohromadě (Pala; Ševeček, 2013). Příkladem budiž *kostra* a její části, např. *lebka*. V SSČ<sup>5</sup> i SSJČ se u *lebky* uvádí, že jde o *kostru hlavy*. Lze tedy s jistou rezervou tvrdit, že heslo obsahuje své holonymum<sup>6</sup>, opačně to však již nefunguje. Z celkem evidentních důvodů nejsou u hesla *kostra* uvedeny všechny její části. Tento příklad příhodně ukazuje i jistou nesystémovost tradičních slovníků, která je pro počítačové zpracování fatální.

Naznačeny tedy byly vlastnosti, jež by databáze významů měla oproti tradičnímu slovníku mít, aby byla použitelná pro počítačové zpracování přirozeného jazyka. Především jde o systémovost vztahů. Hypero/hyponymie je vztah oboustranný, tudíž by mělo být možné se stejnou cestou dostat od nadřazeného slova k podřazenému a naopak. Dále jsou podstatné pojmenované sémantické vztahy mezi slovy. Díky nim je totiž možno jednoznačně určit, které slovo (či slova) je v takové databázi konkrétnímu slovu nadřazené, které je jeho specifikací, označením jeho částí, etc.

S touto myšlenkou tedy vznikl Wordnet - lexikální síť provázaná sémantickými vztahy, která by dle poznatků psycholingvistiky odrážela uspořádání lexikálního materiálu v lidském mozku (o tom v dalších kapitolách ). (Pala; Ševeček, 2013) Zde by bylo na místě poznamenat, že ačkoliv se tak z odstavců výše může čtenáři jevit a i všeobecně je to často tvrzeno, Wordnet není ontologií v pravém slova smyslu, protože něco něco.. [https://en.wikipedia.org/wiki/WordNet#WordNet\\_as\\_a\\_lexical\\_ontology](https://en.wikipedia.org/wiki/WordNet#WordNet_as_a_lexical_ontology)

ref

a tedy tomu vůbec nerozumím, ale přijde mi to relevantní

## 1.2 Historie Wordnetu

<sup>4</sup>protikladu

<sup>5</sup>Slovník spisovné češtiny

<sup>6</sup>vztah opačný k meronymii; tedy např. *dům* je holonymem pro *okno*, *dveře*, *práh* etc.

# Seznam literatury

- Garshol, Lars Marius. Metadata? Thesauri? Taxonomies? Topic Maps! 2004—10—26)[2010-6—19] <http://www.ontopia.net/topicmaps/materials/tm-VS—thesauri.htm1>. Dostupné také z: <http://www.ontopia.net/topicmaps/materials/tm-vs-thesauri.html#N773>.
- Miller, George A et al., 1990. Introduction to WordNet: An on-line lexical database. *International journal of lexicography*, 3(4): s. 235–244.
- Pala, Karel; Ševeček, Pavel, 2013. Česká lexikální databáze typu WordNet. *Feedback*, 48(A47).