

# PLIN021 Sémantická analýza v praxi

OP VK Mezi bohemistikou a informatikou  
[www.projekt-inova.cz](http://www.projekt-inova.cz)

Zuzana Nevěřilová  
[xpopelk@fi.muni.cz](mailto:xpopelk@fi.muni.cz)

Centrum zpracování přirozeného jazyka, B203  
Fakulta informatiky, Masarykova univerzita

11. dubna 2012

Sémantické rámce

Taxonomie

Sémantické sítě

Odvozování

Existující sémantické sítě

## Rámce – použití

Rámce můžeme použít pro desambiguaci slov i celých vět  
[Laparra and Rigau, 2009].

[Bernard Lansky]*STUDENT* studied [the piano]*SUBJECT*  
[with Peter Wallfisch]*TEACHER*.

## Rámce – použití

Rámce můžeme použít pro doplnění implicitní (nezmiňované) znalosti.

*Koupila jsem ojetou felicii. Byly to vyhozené peníze.*

koupit:

- má \_\_činitele člověk/instituce/skupina
- má \_\_benefaktora člověk/instituce/skupina
- má \_\_předmět výrobek/nemovitost/zvíře/rostlina/přírodnina
- má \_\_část činitel dá peníze
- má \_\_část benefaktor dá předmět

## └ Sémantické rámce

## └ Rámce – použití

Rámce můžeme použít pro **hledání** i **implikaci** i **semantickou** i **syntaktickou**.

**Komplexy** jsem objevil kvůli. Byly to vyhledání **práce**.

komplexy:

- má \_číslo (číslo) i (množina) / složit
- má \_konfakto (číslo) i (množina) / složit
- má \_předmět vým (be) / se množin / zvláš / (množina) / předmět
- má \_číslo (číslo) i **práce**
- má \_číslo konfakto i předmět

Rámce se používají hodně. FrameNet nebo VerbaLex jsou zajímavé i svým velkým rozsahem, nejsou to jen nějaké experimenty s uměle vybranými jevy.

zjistit, jak je na tom český FrameNet (PhD práce J. Materny)

# Ontologie

O. je značně nadužívaný pojem, v informatice znamená „formální a explicitní specifikaci sdílené konceptualizace“ [Gruber, 2009]

- formální
- explicitní
- sdílené pojmy

# Ontologie

- slovník (glosář, inventář pojmů ...)
- taxonomie (tezaurus, inventář relací ...)

└ Sémantické rámce

└ Ontologie

- slovesit [glósiť, ievestítʹ sojmá...]
- taxonomie [sozdať, ievestítʹ náci...]

Ontologie se skládá z uvedených součástí a má uvedené vlastnosti, jinak ale může mít libovolný „tvar“. O. můžeme chápat jako nadpojem pro taxonomie, sémantické sítě atd. Bohužel kvůli nadužívání termínu v mnoha oblastech se setkáme s odmítáním zařadit určité projekty pod pojem ontologie. Vždy, když se hovoří o o., je potřeba ujasnit si, co tím myslíme. Nám v tomto kurzu bude stačit tento jednoduchý pohled a budeme se soustředit hlavně na různé „tvary“ a využití o.



## Taxonomie (stromy)

- Aristoteles – kategorie (všech) entit, které mohou lidé vnímat
- Porfyrios – uspořádal kategorie
- Carl Linné – klasifikace (všech) organismů

důležité rysy: uzly jsou **třídy** (organismů, entit ...), třídy jsou **strukturované** do stromu (podtřída, nadtřída), uzly na stejné úrovni se **vzájemně vylučují** (implicitní předpoklad)

## └ Taxonomie

## └ Taxonomie (stromy)

- Aristotelis - kategórie (súčasť etnik, čo je motus kiti v omet)
- Porphyra - aspekt d' kategórie
- Carl Linnaeus - hierarchia (všetky orgány)

• Delenie typy: aký je **rod** (genus), a aké... (špecifický je a  
 aké **rod** do akého (rod, aké, aké), aký je aký a aký  
 aké **rod** (rod) (rod) (rod)

opět – každý z nich má pro své dílo zcela jiné motivy, výsledek je ale docela podobný...

## └ Taxonomie

## └ Taxonomie (stromy)

- Aristotelis - kategórie (dělení) ontiké, toú é mótos káti vónat
- Porphyrius - aspektá é katagoriké
- Carl Linnaeus - de philosophia botanica

É délitá spék: sly jóna **óný** jn gáris má, é mít...]. Wády jóna **strukturovat** éu stromu (podřídka, nadřídka), sly u stajé é novší **se vzájemně vylučují** (implikace) přívěpka é]

Najít dobré ukázky (obrázky k výše uvedeným). U Aristotela je těch kategorií 10: substance, kvantita, kvalita, relace, místo, čas, bytí na pozici, bytí ve stavu, děláni, ovlivnění

# Porfyriův strom (John. F. Sowa)

*Supreme genus:*

**Substance**

*Differentiae:*

**material**

**immaterial**

*Subordinate genera:*

**Body**

**Spirit**

*Differentiae:*

**animate**

**inanimate**

*Subordinate genera:*

**Living**

**Mineral**

*Differentiae:*

**sensitive**

**insensitive**

*Proximate genera:*

**Animal**

**Plant**

*Differentiae:*

**rational**

**irrational**

*Species:*

**Human**

**Beast**

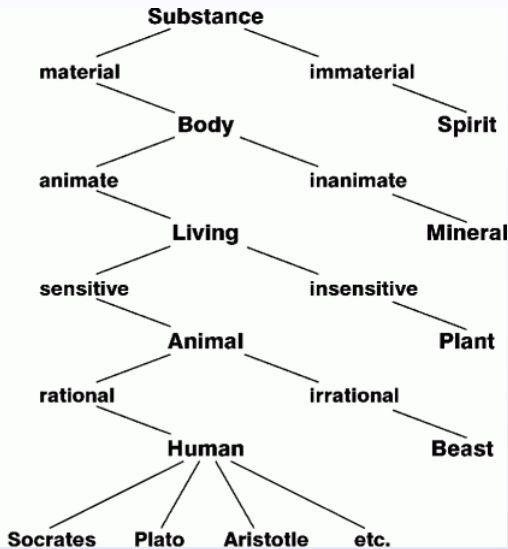
*Individuals:*

**Socrates**

**Plato**

**Aristotle**

**etc.**



# Taxonomie (stromy)

relace is a

relace member of

└ Taxonomie

└ Taxonomie (stromy)

Připomenout rozdíl mezi třídou a instancí. Jak je v přirozeném jazyce rozeznáváme? Například podle jména:

Pes je masožravec, nepohrdne však ani ovocem.

Alík má rád švestky.

V mnoha případech je to složitější:

Americký prezident je zároveň předsedou vlády.

Americký prezident má babičku v ohrožení.

## Taxonomie (stromy) a slovníkové definice

klasická definice = **genus proximum** + **differentia specifica**

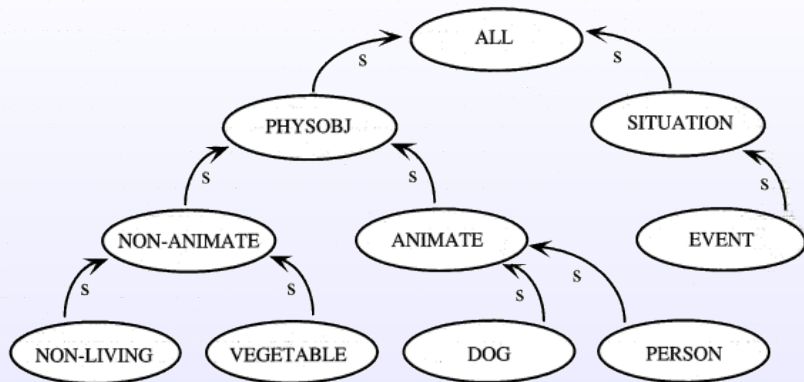
Počítač je v informatice **elektronické zařízení**, které **zpracovává data pomocí** předem vytvořeného **programu**.

Elektronické zařízení je **zařízení**, jehož **funkce závisí na elektrickém proudu** nebo na elektromagnetickém poli.





# Sémantické sítě



sémantická síť = reprezentace lexikálních znalostí  
[Collins and Quillian, 1969]

uzly = entity (třídy nebo instance), jednomu konceptu odpovídá jeden uzel

# Sémantické síť

- nadtyp–podtyp, is a, is-a, isa (hypo/hyperonymie)
- instance třídy, member of
- část–celek, has a (holo/meronymie)
- upřesnění akce (troponymie)
- příčina–následek
- ...

└ Sémantické sítě

└ Sémantické sítě

- nadtyp–podtyp, k a, is-a, isa (typy/typy-nymik)
- část–celek, member of
- část–část, k a a (kolo/nem-nymik)
- přísluší a lež (troupy-nymik)
- přísluší–vzájemně
- ...

Jednotlivé podsítě, kde uzly spojují relace jednoho druhu, jsou stromy (tj. taxonomie). Je to docela logické – v každém druhu relace máme nějaké uspořádání „od nejmenšího po největší“.

Např. nadtyp–podtyp je klasická taxonomie. Část–celek taky, protože objekt  $x$  se skládá z částí  $a$  a  $b$ , část  $a$  se skládá z částí  $m$  a  $n$  (které jsou patrně menší než  $a$  i menší než  $x$ ).

└ Sémantické sítě

└ Sémantické sítě

- *rod* = typ = rod by p, k, a, i, se, ísa (typy) by p, se, mynk
- *číslo* = a třídy = množstvo
- *číslo* = k, ísa a (kolik) se, mynk
- *průběh* = a ísa (troupa mynk)
- *průběh* = řádění k
- ...

Důležitou vlastností taxonomií (tj. i těch částí sém. sítě, které tvoří taxonomii) je tranzitivita. Využíváme ji v odvozování (viz dál).

## Odbočka k odvozování

Fakt  $F$  = tvrzení s pravdivostní hodnotou (např. ptáci létají)

Báze znalostí (knowledge base)  $KB$  = (pokud možno konzistentní) soubor faktů (např. ptáci létají, vlaštovka je pták)

Pokud z  $KB$  plyne  $F$  a přidáme další fakt takový, že  $KB$  je stále konzistentní, je  $KB$  **monotónní** reprezentace. [Allen, 1995]

ptáci létají

vlaštovka je pták

---

vlaštovka létá

## Odbočka k odvozování

ptáci létají  
tučňák je pták

---

tučňák létá

kromě tučňáka ptáci létají  
tučňák je pták

---

NOT(tučňák létá)

kromě tučňáka ptáci létají  
pštros je pták

---

pštros létá

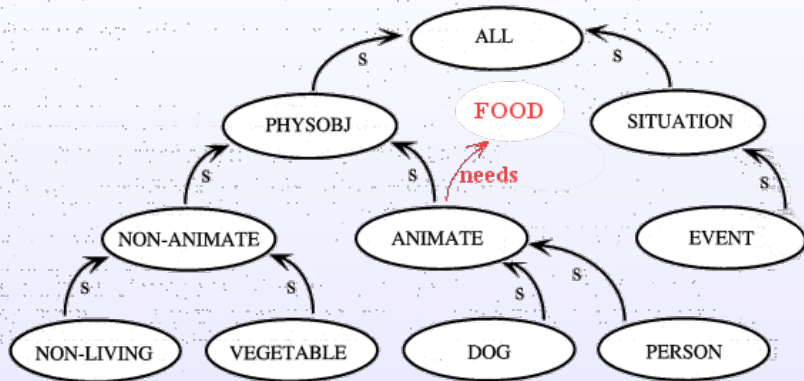
kromě tučňáka, pštroso, mlád'at, mrtvých ptáků, ptáků se zraněnými křídly . . . ptáci létají

## Odbočka k odvozování

Používáme **implicitní pravidlo** (*default rule*), tj. ptáci létají, dokud neřekneme jinak. Uvedeme-li implicitní pravidlo, má přednost před obecným faktem.

Ptáci létají, ale tučňák ne.

# Sémantické sítě – dědičnost



odvozování je **monotónní**



└ Odvozování

└ Sémantické sítě – dědičnost



Praktická ukázka dědičnosti a odvozování:

1. silniční vozidlo má (has part) volant
2. dodávka je (isa) silniční vozidlo
3. dodávka má (has part) volant
4. Mercedes Sprinter je (member of) dodávka
5. Mercedes Sprinter má (has part) volant

Podíváme se, jestli uvedené „lezení po větvích“ platí pokaždé.

# Sémantické sítě

WordNet a EuroWordNet

český WordNet

Sítě z Wikipedie, dbpedia, ArtNet

└ Existující sémantické sítě

└ Sémantické sítě

Nachystat ukázky. Český WordNet je jistě dobré téma pro BP. Podobně česká dbpedia (která dosud neexistuje).

## Asociativní sítě

Někdy totožné se sémantickými sítěmi, jindy u asociativních sítí neplatí předpoklad, že jeden koncept odpovídá jednomu uzlu.

└ Existující sémantické sítě

└ Asociativní sítě





Uvedený předpoklad typicky neplatí u sítí, které jsou generovány automaticky.

- └ Existující sémantické sítě

- └ Asociativní sítě

Sémantické (asociativní) sítě se používají velmi mnoho. Je jich mnoho druhů, je mnoho způsobů jejich zápisu, jsou v módě (protože sémantický web je v módě a protože web tvoří také síť).

V souvislosti se sémantickým webem můžeme zmínit jazyky sémantického webu: RDF, RDFS, OWL. Podrobněji se jim můžeme věnovat, pokud bude čas (a chuť). Trochu jsme nakousli odvozování. Podle experimentů z kognitivní vědy (dohledat) lidé nemají v hlavě všechno, ale odvozují. Oblast odvozování znalostí ze znalostníchází by vydala na zvláštní seminář, přesto se odvozování nemůžeme úplně vyhnout. Logická reprezentace bývá zpravidla chápána mimo jazyk, přesto cítíme, že je neoddělitelnou součástí promluvy.

-  Allen, J. (1995).  
*Natural Language Understanding (2nd ed.)*.  
Benjamin-Cummings Publishing Co., Inc., Redwood City, CA,  
USA.
-  Collins, A. M. and Quillian, M. R. (1969).  
Retrieval time from semantic memory.  
*Journal of Verbal Learning and Verbal Behavior*, 8(2):240–247.
-  Gruber, T. (2009).  
Ontology.  
In Liu, L. and Özsu, M. T., editors, *Encyclopedia of Database  
Systems*, page 1963–1965. Springer Verlag.
-  Laparra, E. and Rigau, G. (2009).  
Integrating wordnet and framenet using a knowledge-based  
word sense disambiguation algorithm.  
In *RANLP*, Borovets, Bulgaria.