

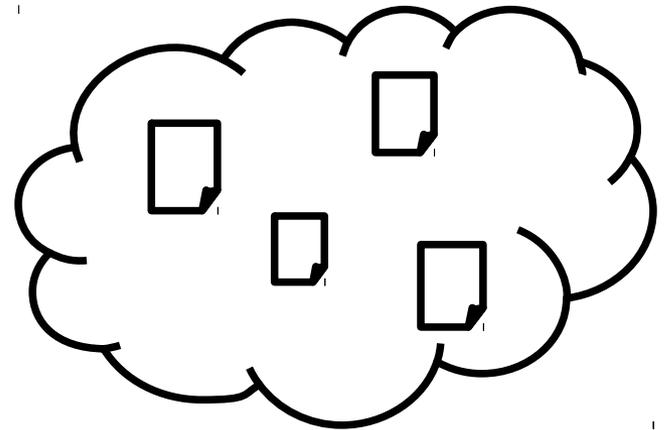
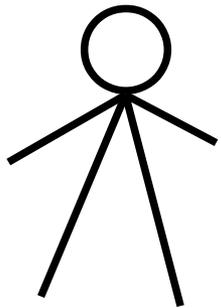
Information Extraction and Natural Language Processing

Marián Šimko



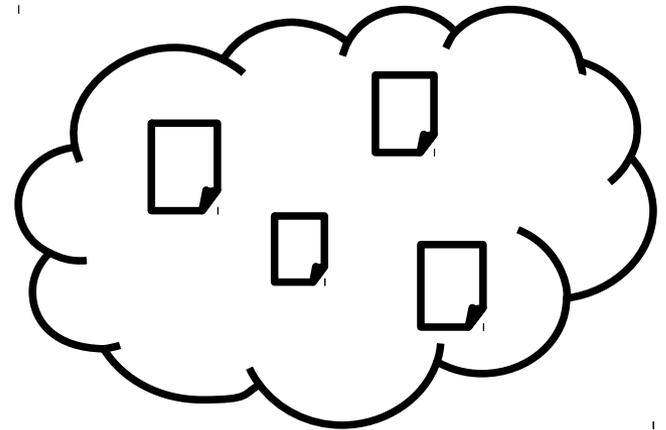
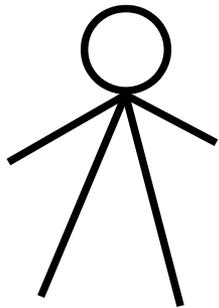
25. 11. 2015

Prístup k informáciám



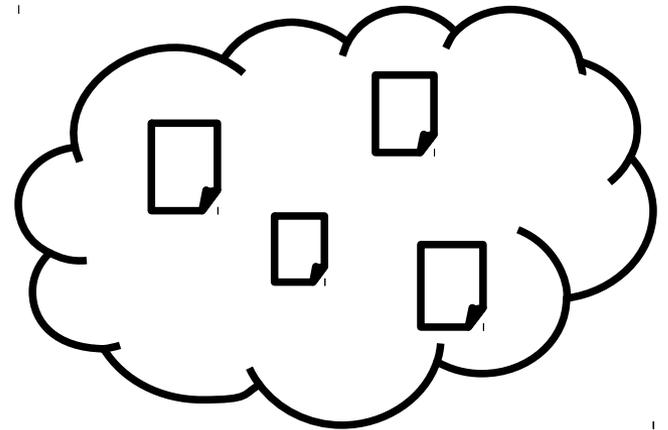
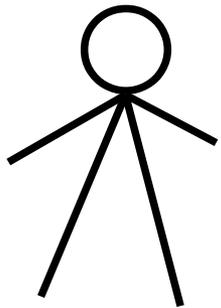
Prístup k informáciám

November 2015:
55 miliárd webových stránok

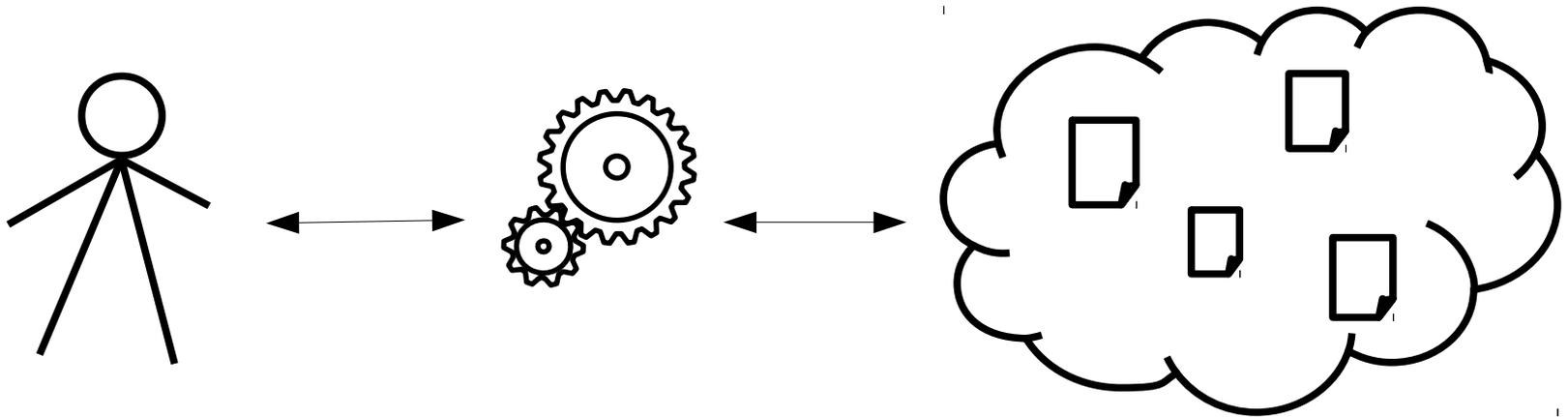




Prístup k informáciám



Prístup k informáciám



Naše informačné potreby

- Information retrieval
- Information extraction
- Document classification/categorisation
- Reasoning
- ...
- Artificial intelligence

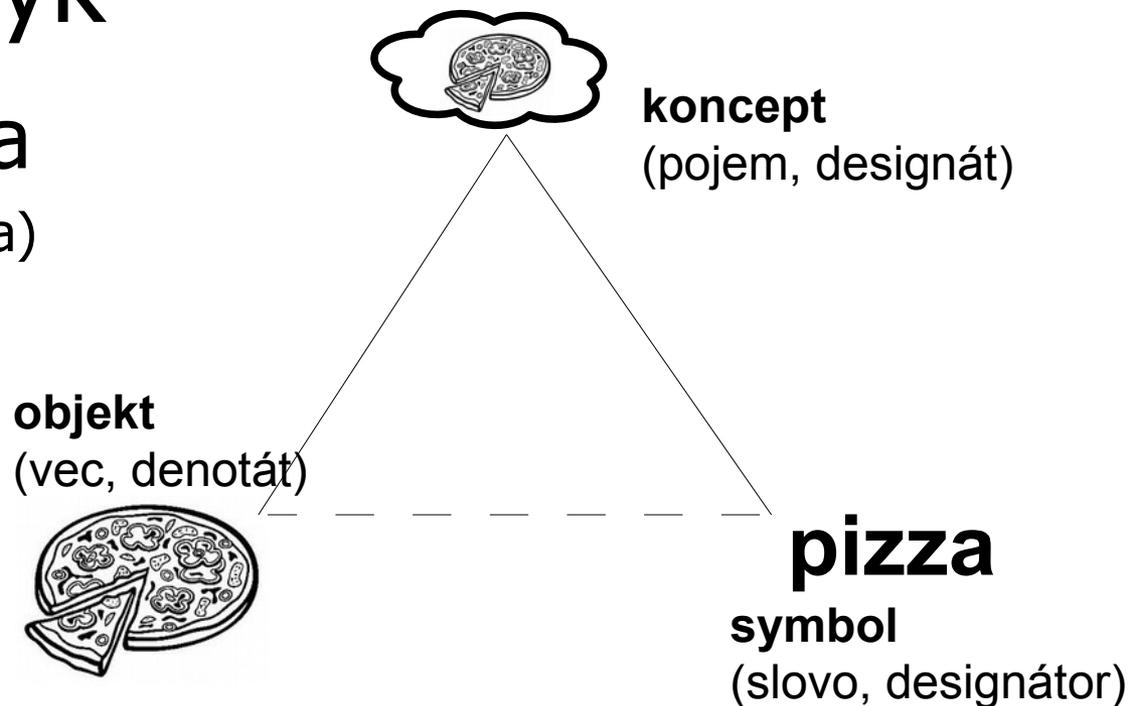
Informácie – kódované v prirodzenom jazyku

Informácie sú v dokumentoch.

Dokumenty =

Prirodzený jazyk

Písomná podoba
(symbolická reprezentácia)



Prirodzený jazyk =
komunikačný prostriedok

Komunikačný postriedok?

Umelá komunikácia

- Vysielateľ – prijímateľ
- Dáta
- Kompresia
- Samoopravný kód
(napr. Hammingov kód)

Komunikačný prostriedok?

Umelá komunikácia

- Vysielateľ – prijímateľ
- Dáta
- Kompresia
- Samoopravný kód
(napr. Hammingov kód)

Prirodzený jazyk

- Autor – čitateľ
- Text
- Stratová kompresia
- Samoopravný kód
(napr. gramatika)

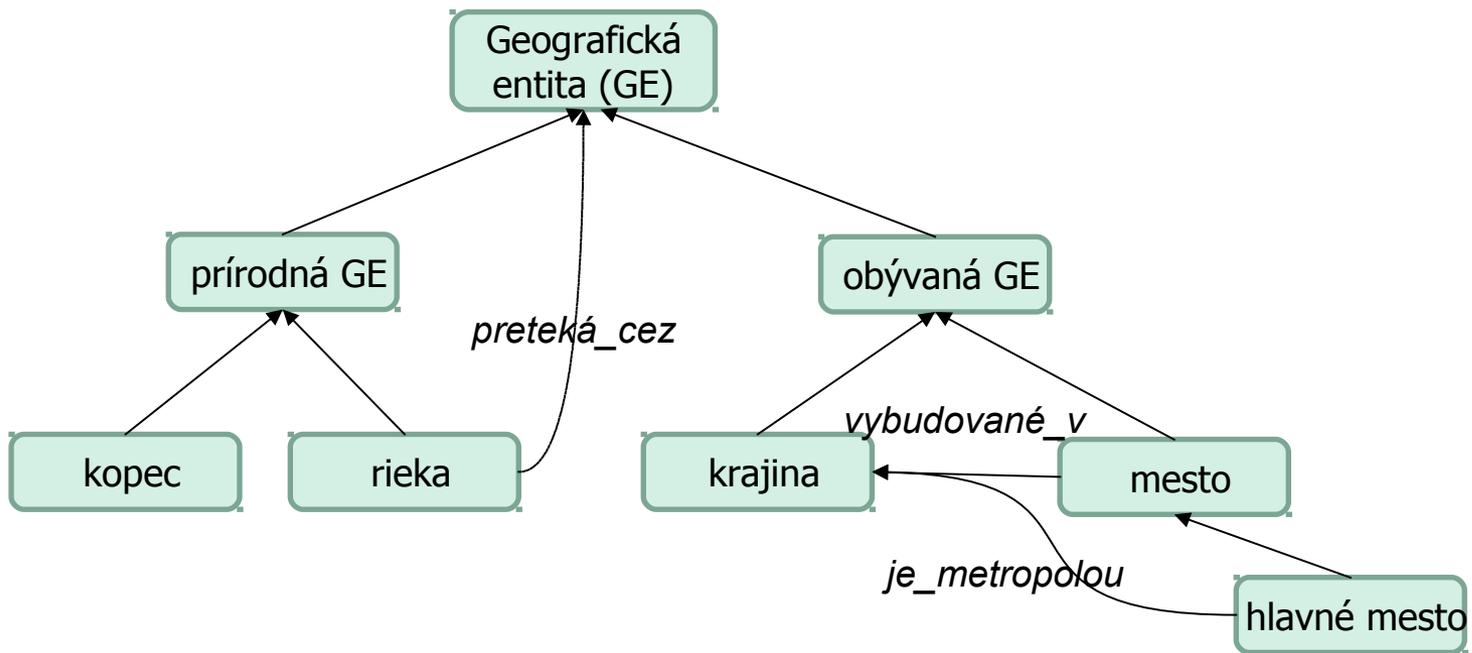
Natural language processing
= spracovanie prirodzeného jazyka

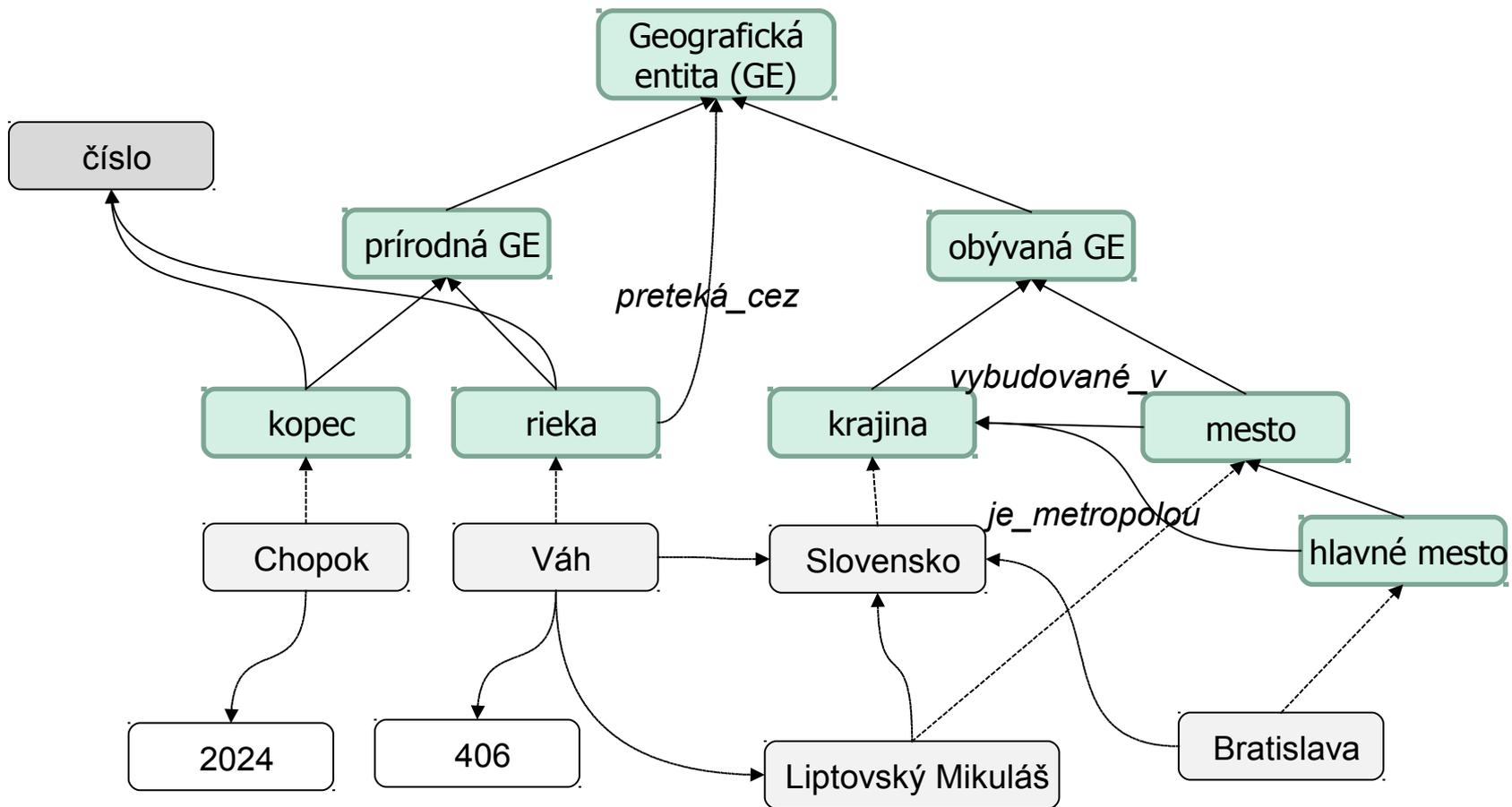
Nachádzanie informácií v komunikácii
v prirodzenom jazyku

Príklad:

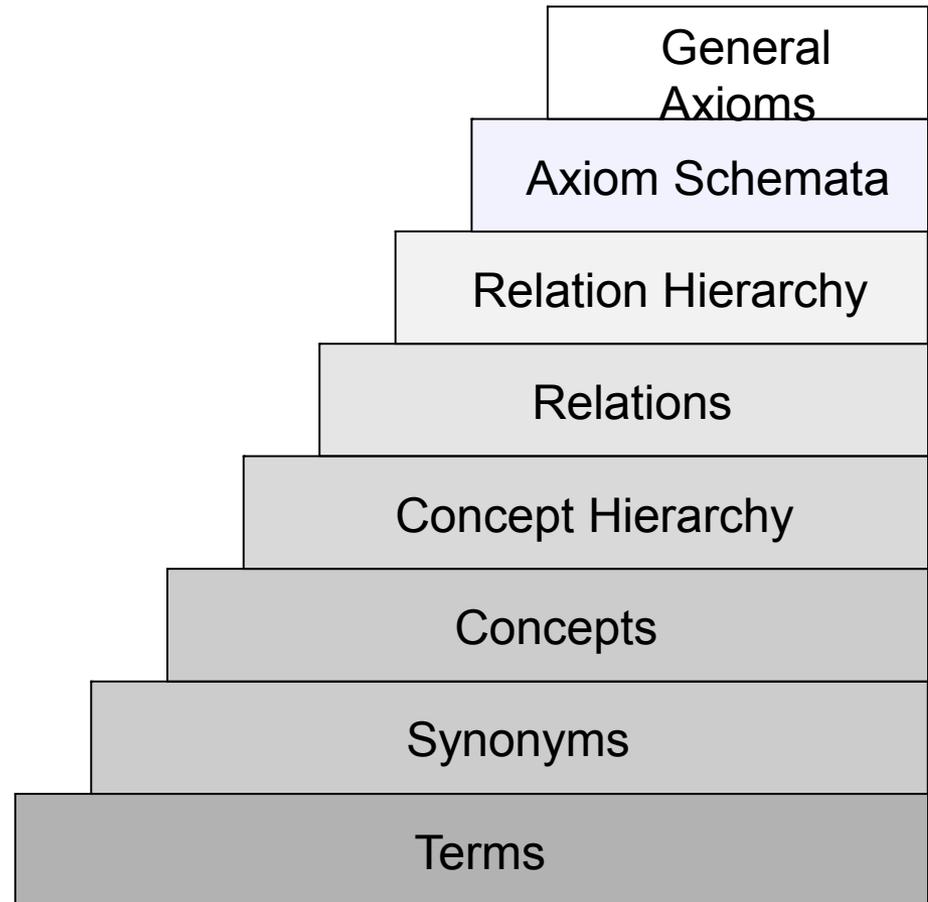
korpus dokumentov z geografie

Chceme opísať danú oblasť



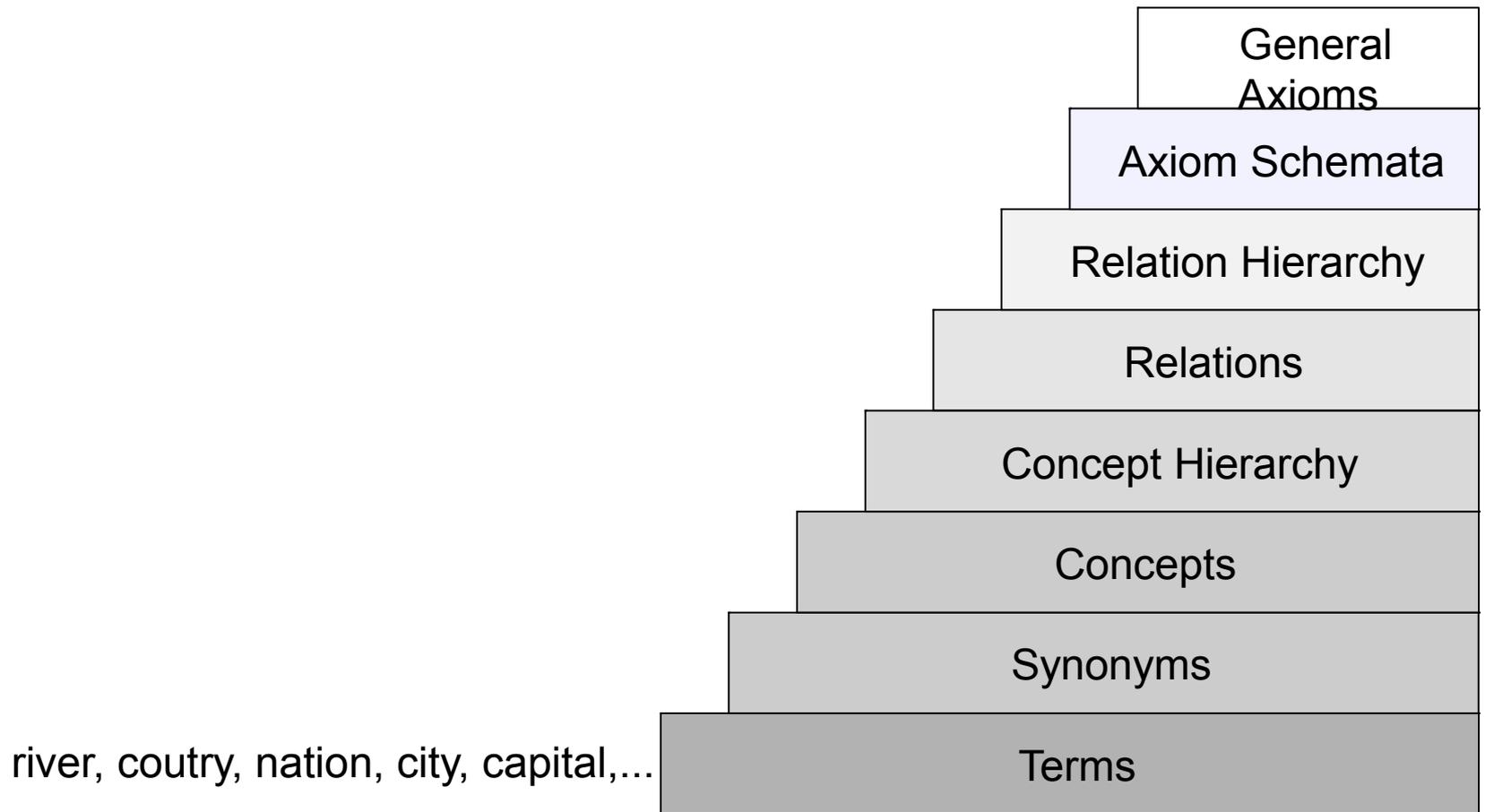


Ontol\u00f3gia



(Cimiano, 2006)

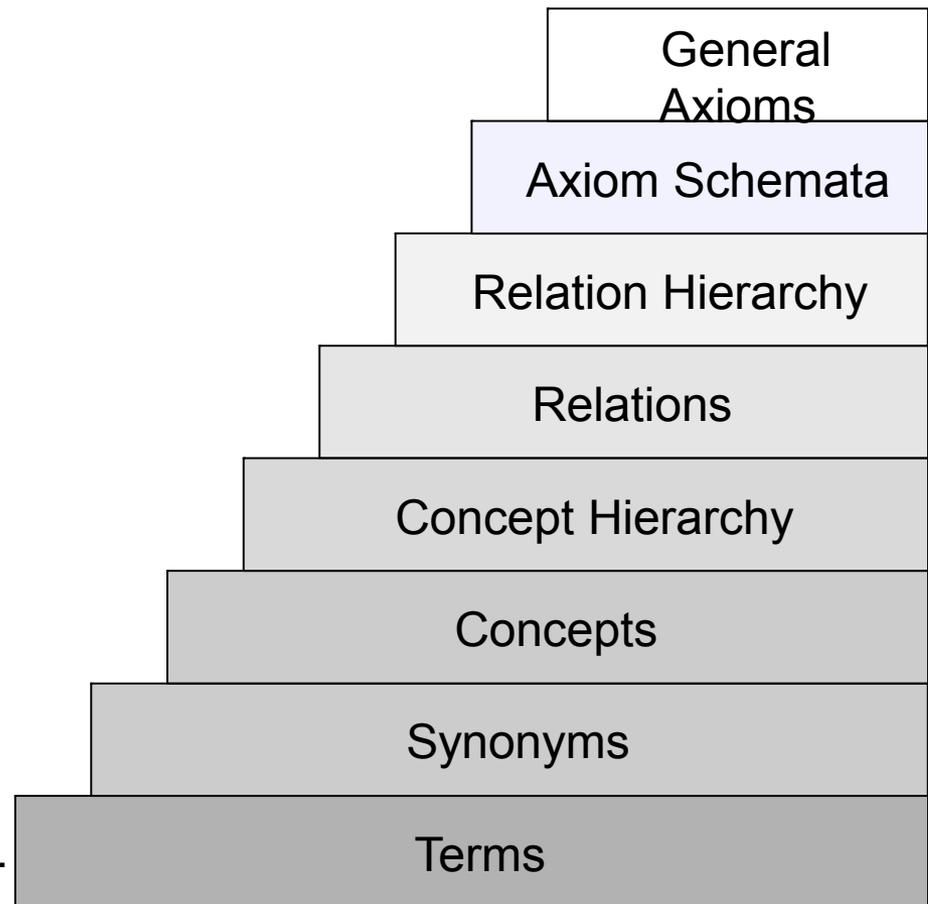
Ontol\u00f3gia



Ontol\u00f3gia

{country, nation}

river, coutry, nation, city, capital,...

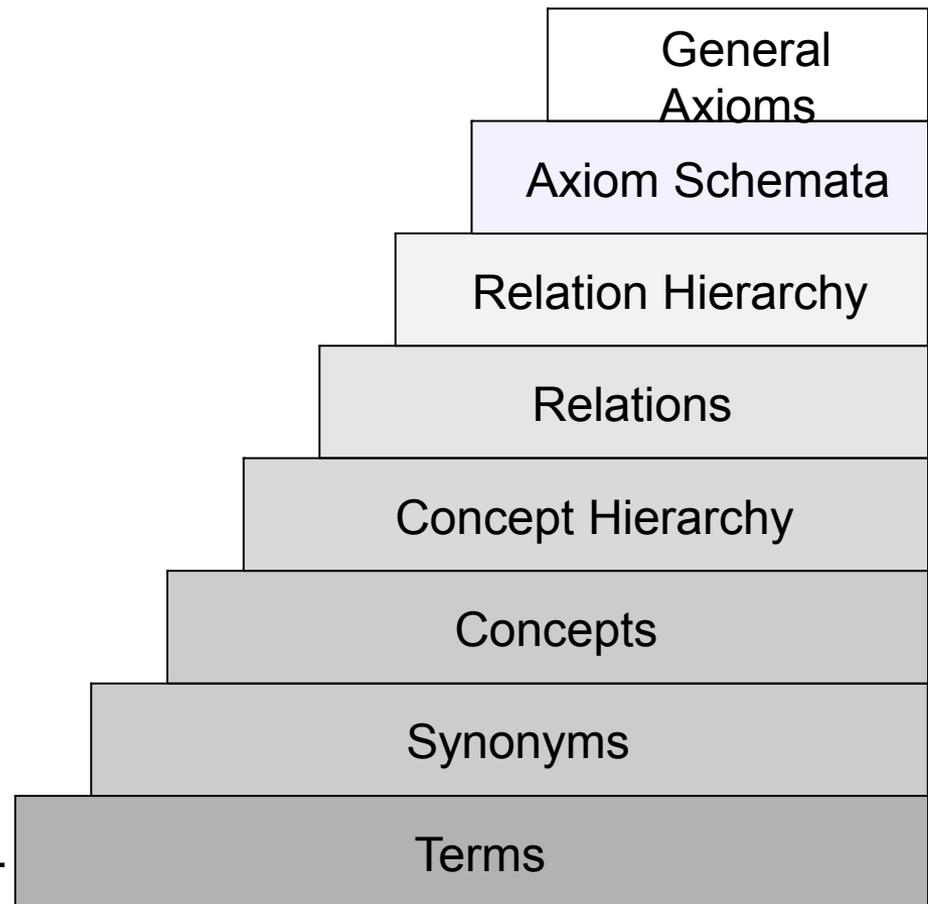


Ontol3gia

$c := \text{country} := \langle i(c), ||c||, \text{Ref}_c(c) \rangle$

{country, nation}

river, coutry, nation, city, capital,...



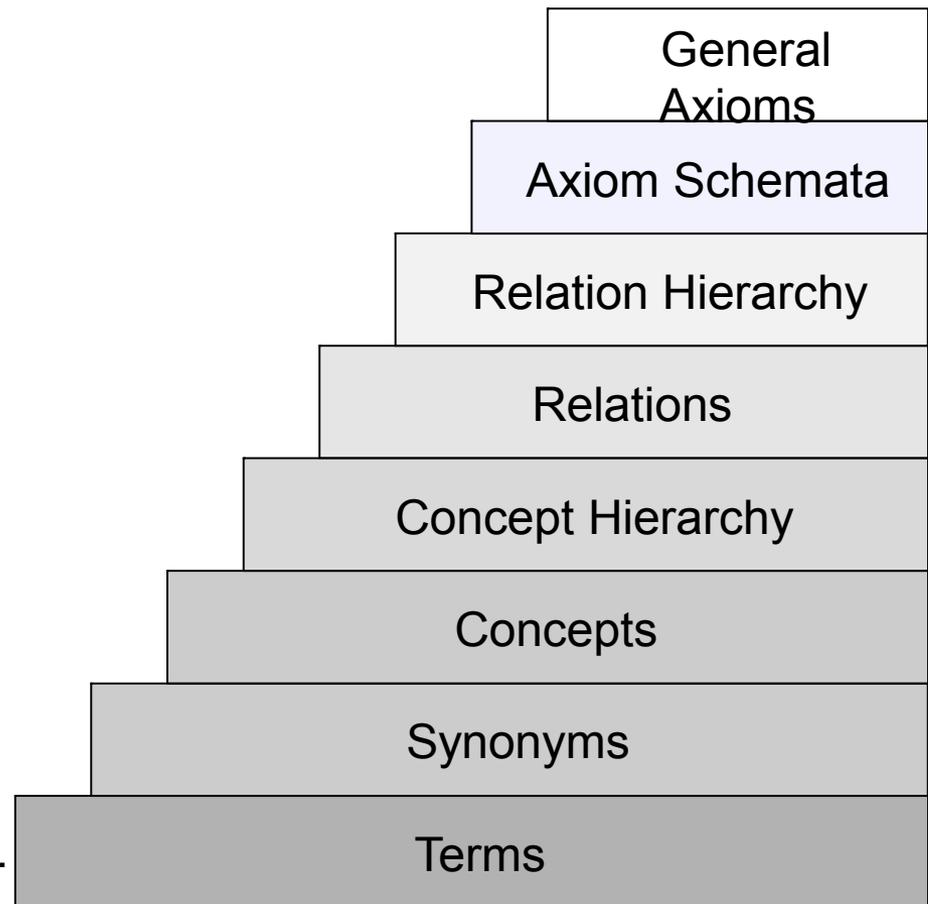
Ontol3gia

capital \leq_c city, city \leq_c GE

c:= country := <i(c), ||c||, Ref_c(c)>

{country, nation}

river, coutry, nation, city, capital,...



Ontol3gia

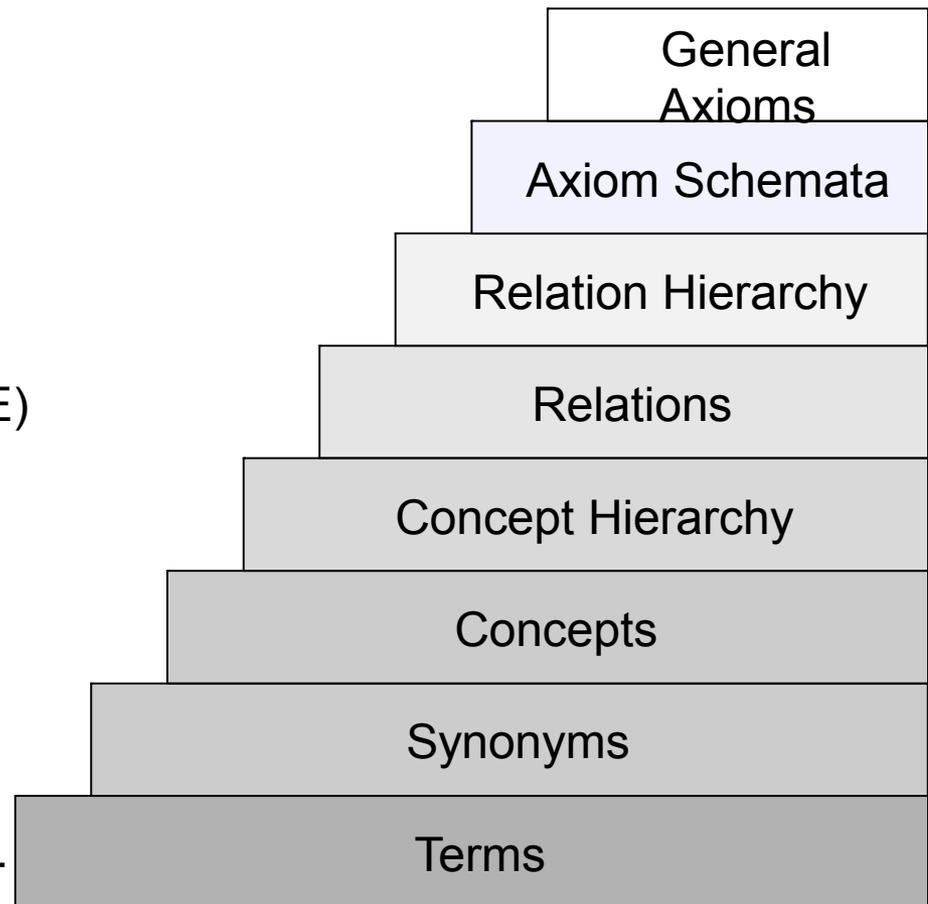
flow_through(dom:river, range:GE)

capital \leq_c city, city \leq_c GE

c:= country := <i(c), ||c||, Ref_c(c)>

{country, nation}

river, coutry, nation, city, capital,...



Ontol3gia

capital_of \leq_R located_in

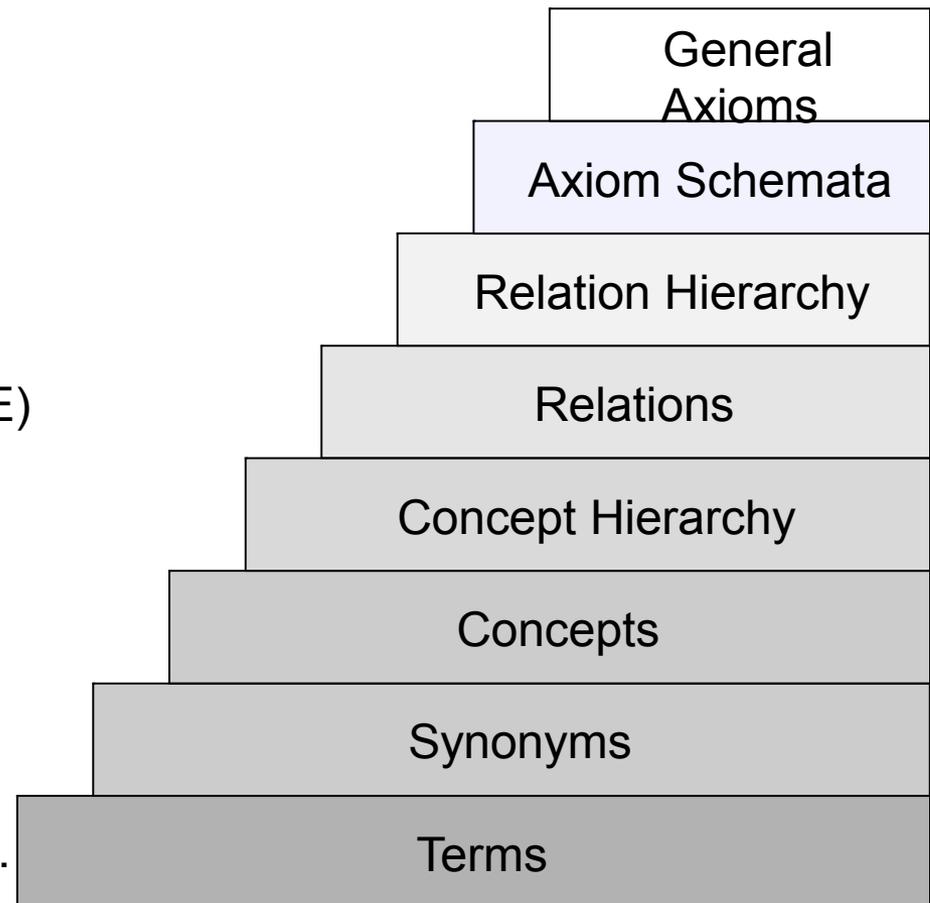
flow_through(dom:river, range:GE)

capital \leq_C city, city \leq_C GE

c:= country := <i(c), ||c||, Ref_C(c)>

{country, nation}

river, coutry, nation, city, capital,...



Ontol3gia

disjoint(river, mountain)

capital_of \leq_R located_in

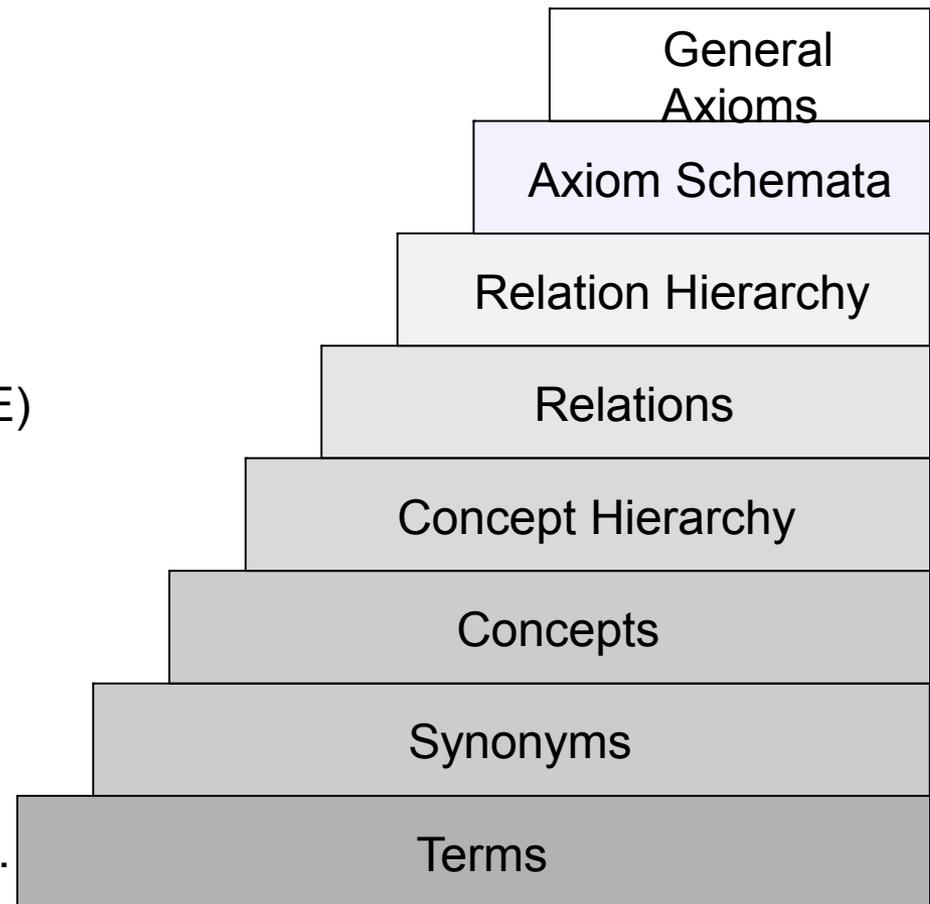
flow_through(dom:river, range:GE)

capital \leq_C city, city \leq_C GE

c:= country := <i(c), ||c||, Ref_C(c)>

{country, nation}

river, coutry, nation, city, capital,...



Ontol3gia

$\forall x(\text{country}(x) \rightarrow \exists y \text{ capital_of}(y,x) \wedge \forall z(\text{capital_of}(z,x) \rightarrow y=z))$

$\text{disjoint}(\text{river}, \text{mountain})$

$\text{capital_of} \leq_R \text{located_in}$

$\text{flow_through}(\text{dom:river}, \text{range:GE})$

$\text{capital} \leq_C \text{city}, \text{city} \leq_C \text{GE}$

$c := \text{country} := \langle i(c), ||c||, \text{Ref}_C(c) \rangle$

$\{\text{country}, \text{nation}\}$

$\text{river}, \text{country}, \text{nation}, \text{city}, \text{capital}, \dots$

General
Axioms

Axiom Schemata

Relation Hierarchy

Relations

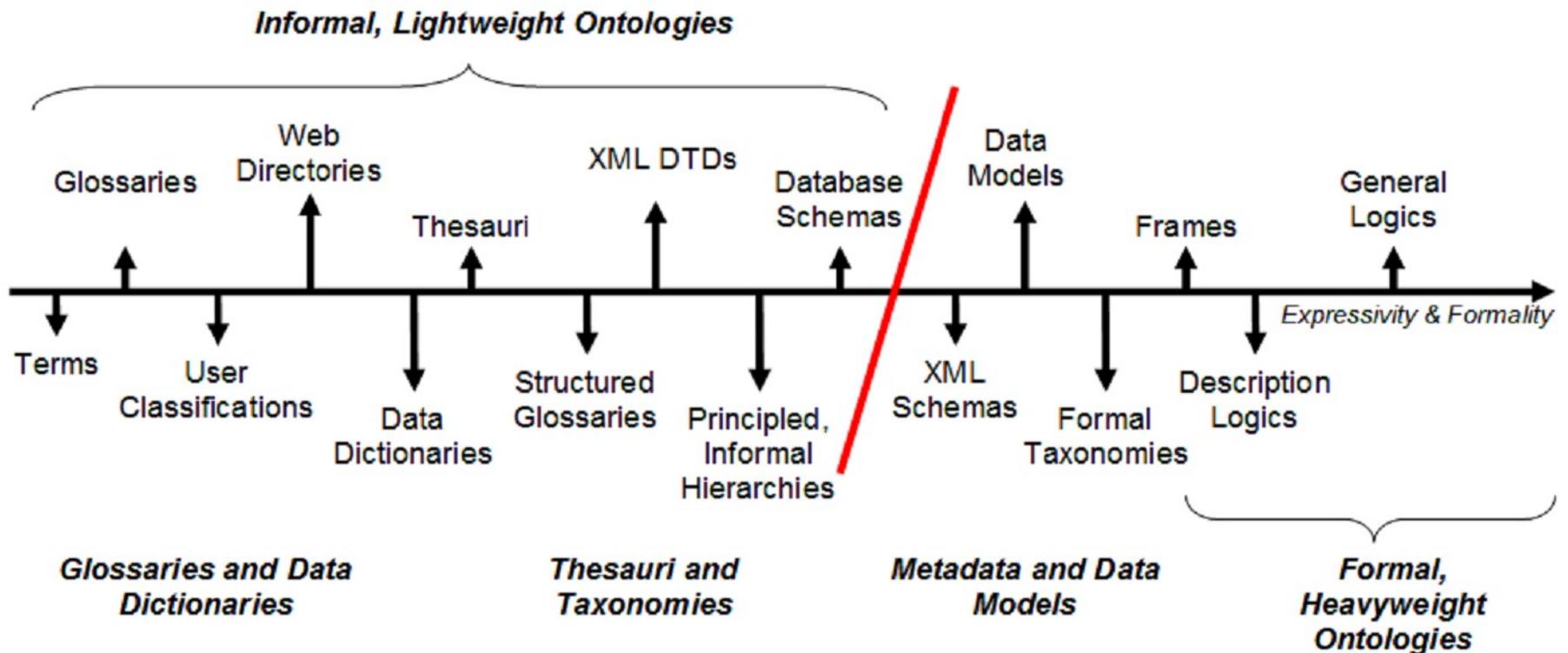
Concept Hierarchy

Concepts

Synonyms

Terms

Lightweight vs. Heavyweight ontologie

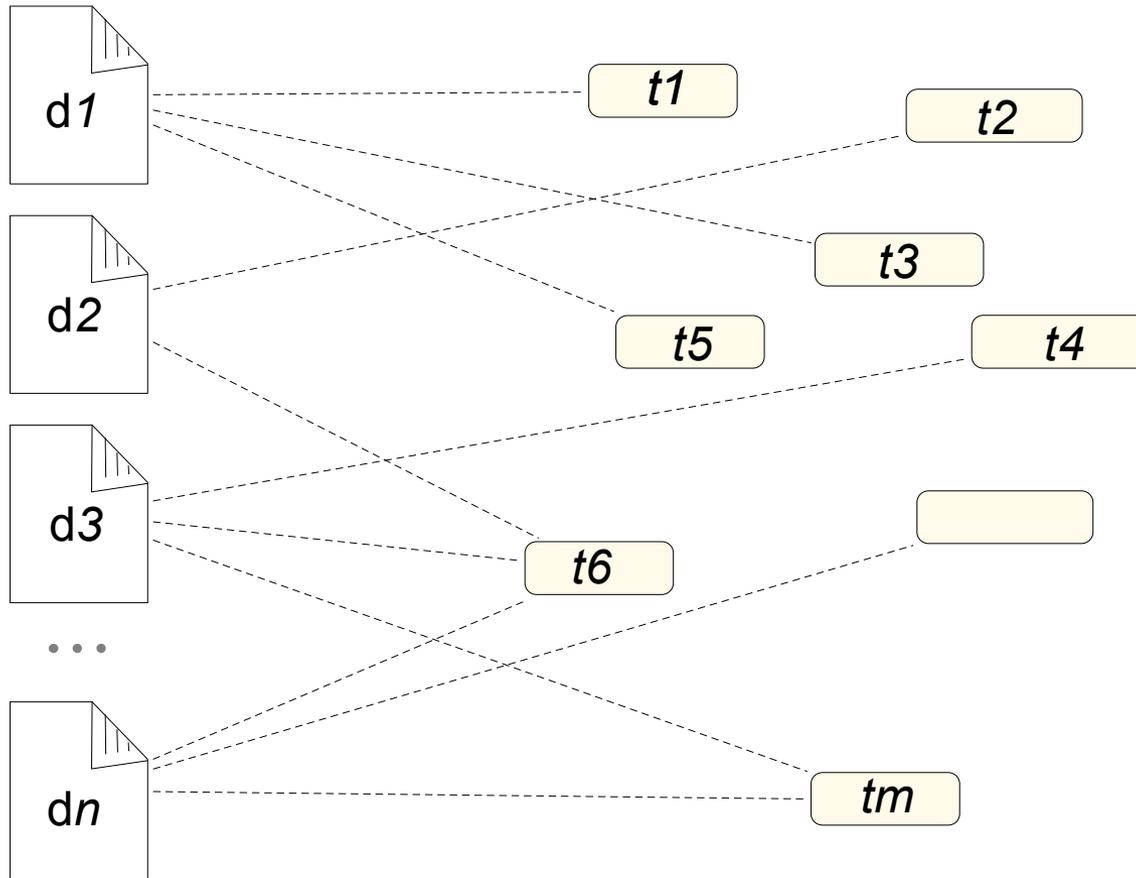


(Wong et al., 2012)

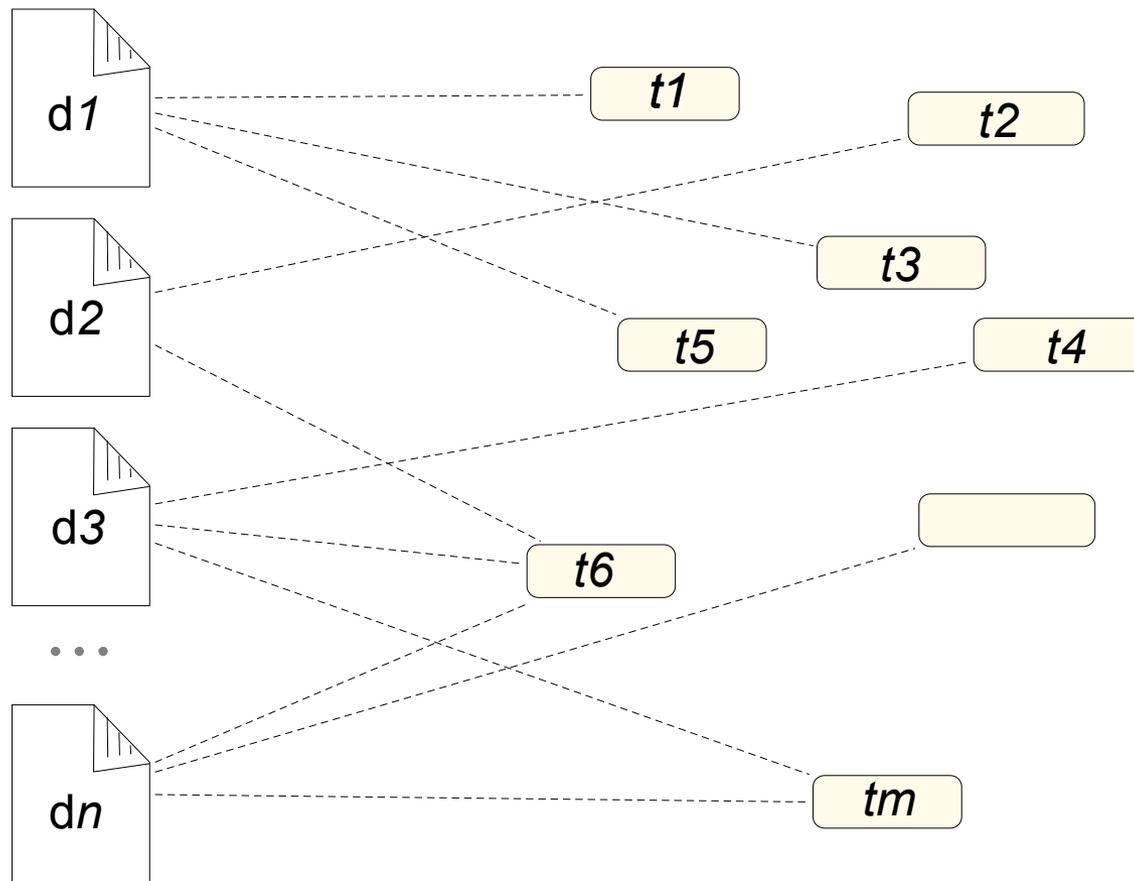
Lightweight ontologie

- Termy
 - jedno a viac-slovné výrazy
 - *pizza*
- Slovníky
 - Termy s výkladom, vysvetlením
 - *pizza* - A baked Italian dish of a thinly rolled bread dough crust typically topped before baking with tomato sauce, cheese and other ingredients such as meat, vegetables or fruit

Termy



Slovníky



$t1$ = a sauce
made of tomatos

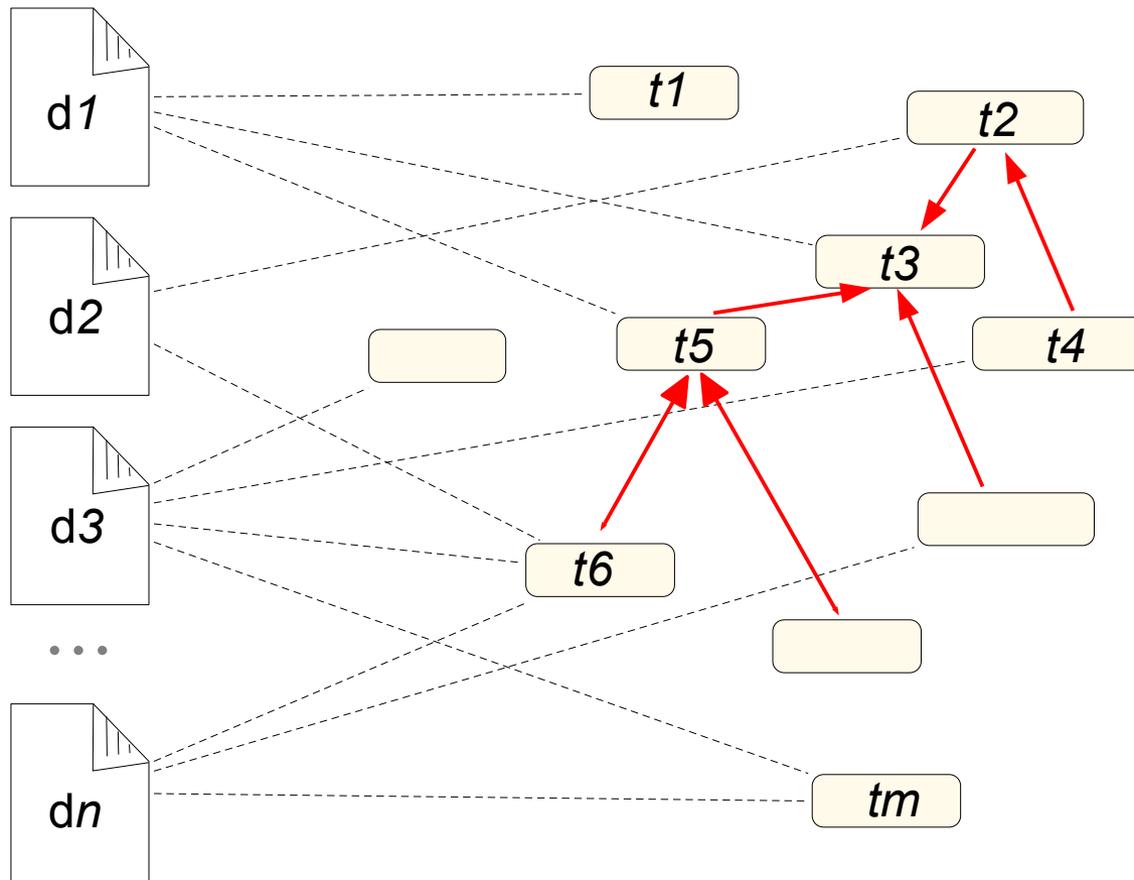
$t2$ = a baked
italian dish ...

$t3$ = ...

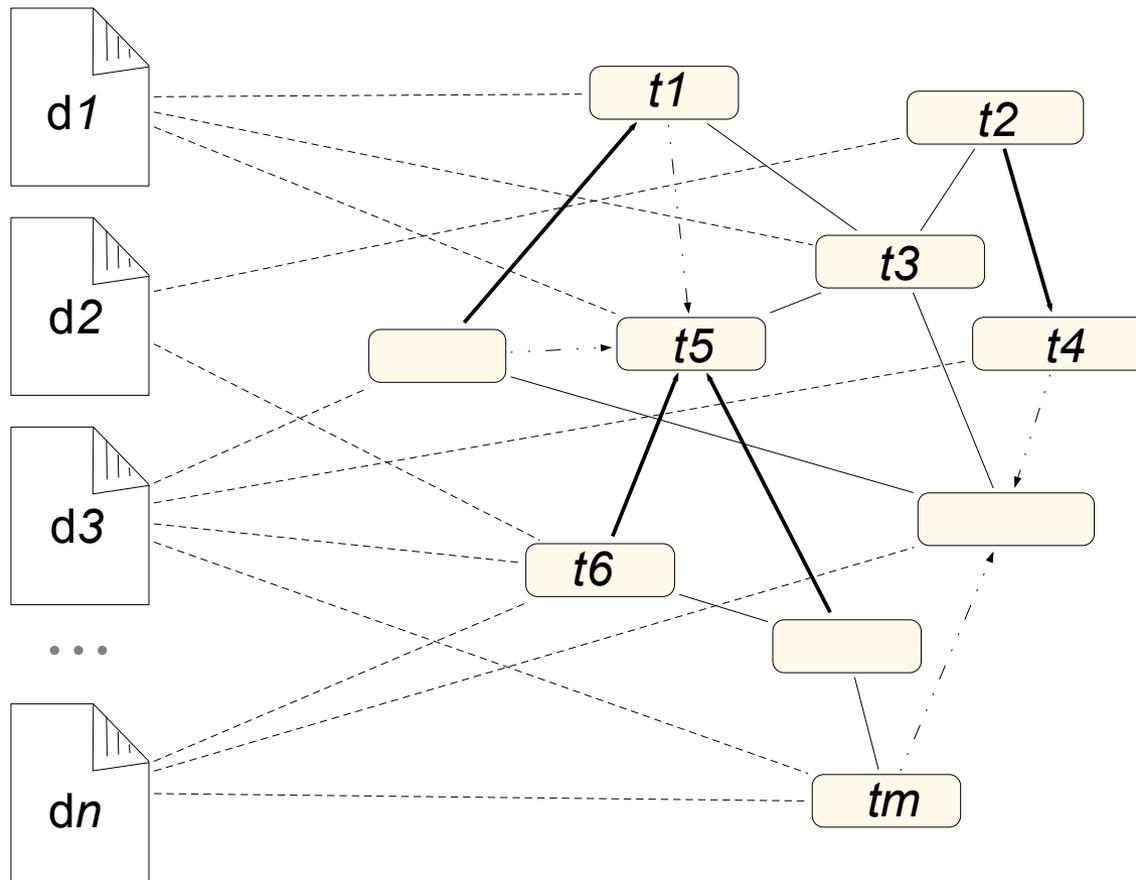
Lightweight ontológia

- Taxonómie
 - Hierarchické vzťahy typu is-a
 - *food, pizza, margherita*
- Tezaury
 - Slovníky, navyše s:
 - Hierarchickými vzťahmi: všeobecnosť (hypero-/hyponymá),
časť celku (mero-/holonymá)
 - Vzťahmi ekvivalencie
 - Vzťahmi paradigmatickej podobnosti (relatedness)

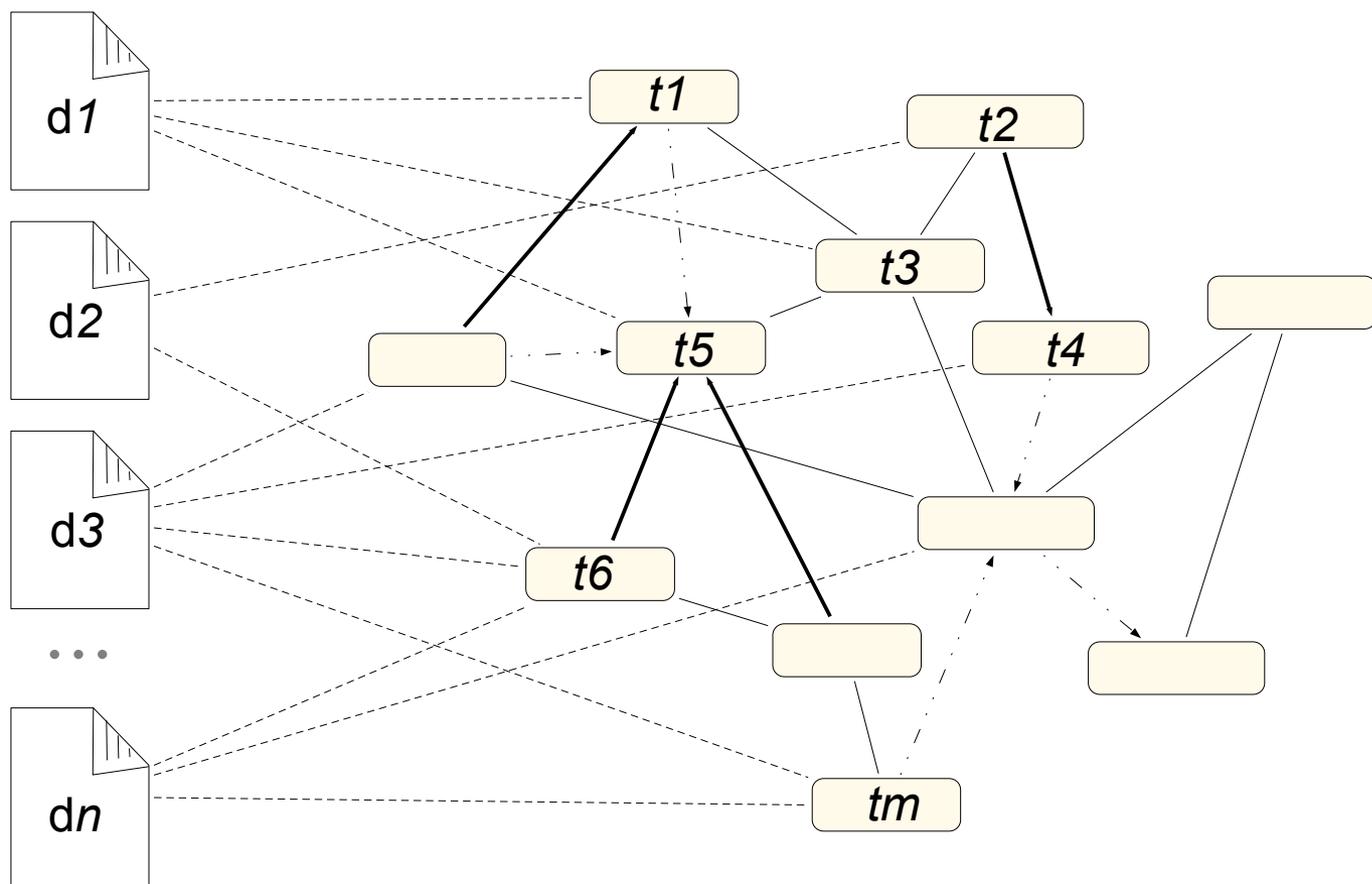
Taxonómie



Tezaury



Lahká ontológia



Heavyweight ontología

$$O = \{C, \leq_C, R, \sigma_R, \leq_R, A, \sigma_A, T\}$$

- four disjoint sets C , R , A and T whose elements are called concept identifiers, relation identifiers, attribute identifiers and data types, respectively,
- a semi-upper lattice \leq_C on C with top element root_C , called concept hierarchy or taxonomy,
- function $\sigma_R: R \rightarrow C^+$ called relation signature,
- a partial order \leq_R on R , called relation hierarchy, where $r_1 \leq_R r_2$ implies $|\sigma_R(r_1)| = |\sigma_R(r_2)|$ and $\pi_i(\sigma_R(r_1)) \leq_C \pi_j(\sigma_R(r_2))$, for each $1 \leq i \leq |\sigma_R(r_1)|$, and
- a function $\sigma_A: A \rightarrow C \times T$, called attribute signature,
- a set T of datatypes such as strings, integers, etc.

(Cimiano, 2006)

Heavyweight ontológia

- Definition 2 (Domain and Range)
- Definition 3 (\mathcal{L} -axiom System)
- Definition 4 (Lexicon)
- Definition 5 (Knowledge Base (KB))
- Definition 6 (Instance Lexicon)
- Definition 7 (Extension)
- Definition 8 (Intension)

(Cimiano, 2006)

Koncept - trojica

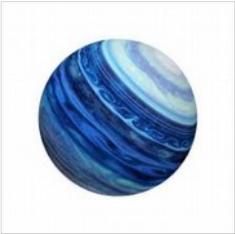
($i(c)$, $e(c)$, $ref(c)$)

- $i(c)$ – intenzionálny opis
- $e(c)$ – extenzionálny opis
- $ref(c)$ – lexikálna realizácia v korpuse

Príklad: Plynný obor



Príklad: Plynný obor



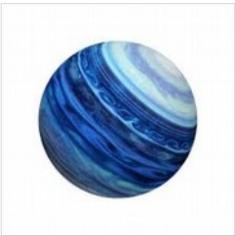
$i(c)$:

Príklad: Plynný obor



i(c): veľká planéta, ktorá nie je zložená prevažne z hornín alebo inej pevnej látky

Príklad: Plynný obor



i(c): veľká planéta, ktorá nie je zložená prevažne z hornín alebo inej pevnej látky

e(c):

Príklad: Plynný obor



i(c): veľká planéta, ktorá nie je zložená prevažne z hornín alebo inej pevnej látky

e(c): Jupiter, Saturn, 47 Ursae Majoris c

Príklad: Plynný obor

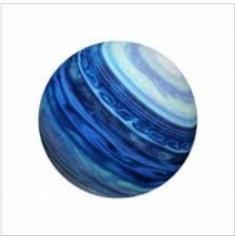


i(c): veľká planéta, ktorá nie je zložená prevažne z hornín alebo inej pevnej látky

e(c): Jupiter, Saturn, 47 Ursae Majoris c

ref(c):

Príklad: Plynný obor



i(c): veľká planéta, ktorá nie je zložená prevažne z hornín alebo inej pevnej látky

e(c): Jupiter, Saturn, 47 Ursae Majoris c

ref(c): plynný obor, Joviálna (Ioviálna) planéta

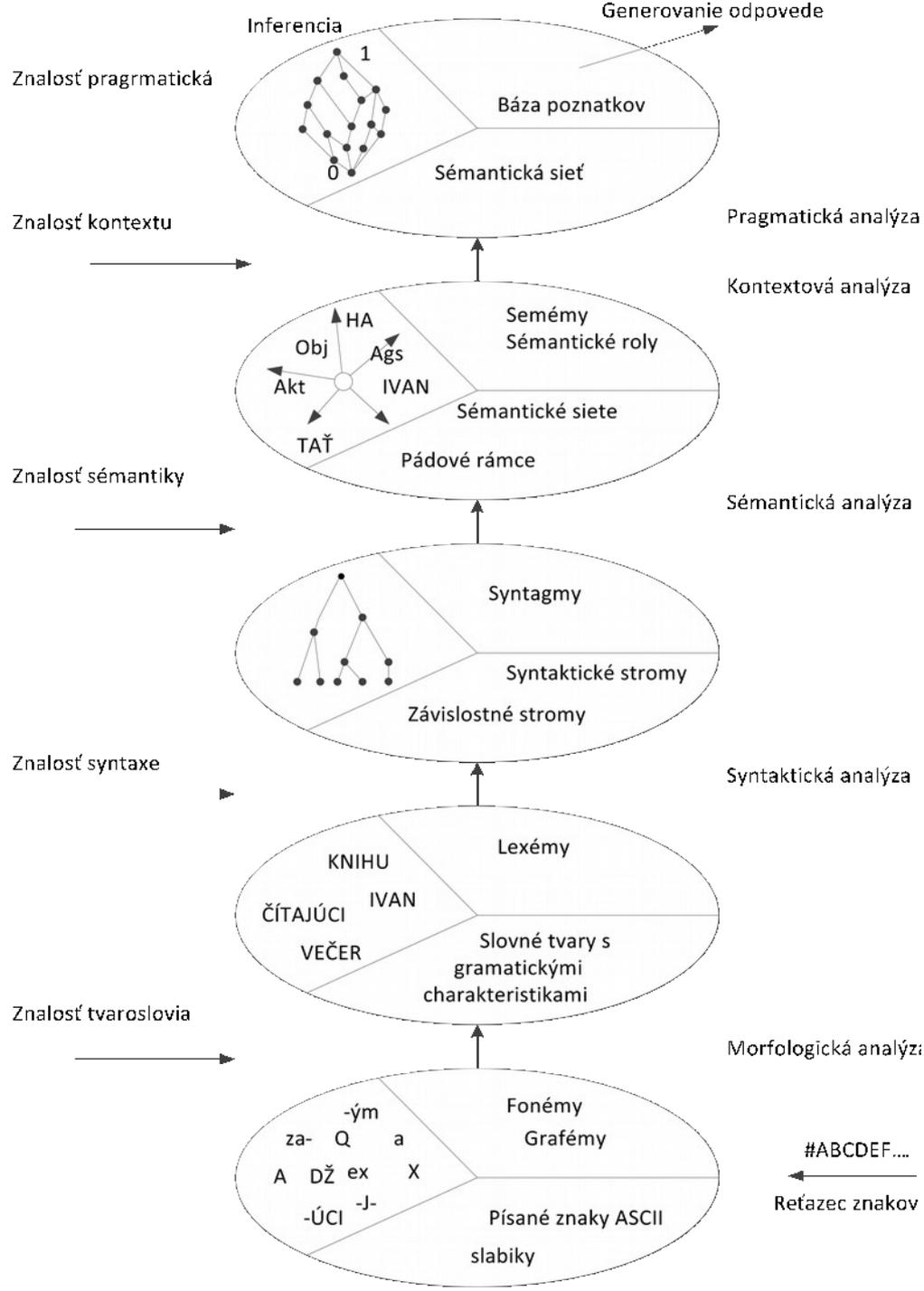
Zápis ontologií/sémantiky

- RDF
 - Resource Description Framework
- RDFS
 - RDF Schema
- OWL
 - Web Ontology Language
 - OWL Lite
 - OWL DL
 - OWL Full

Spracovanie prirodzeného jazyka

Spracovanie prirodzeného jazyka

The title is centered on a dark grey background. Below the text, there is a decorative graphic consisting of a solid teal horizontal bar, followed by a white horizontal bar, and then two thin, parallel teal horizontal lines.



- Zelené žaby žijú v rybníku.
- Červené žaby majú dlhé nosy.

- Zelené žaby žijú v rybníku.

- Červené žaby majú dlhé nosy.

- Červené idey majú dlhé nosy.

pragmatically

- Zelené žaby žijú v rybníku.

- Červené žaby mají dlhé nosy.

pragmaticky

- Červené idey mají dlhé nosy.

sémanticky

- Červenými ideou mať nos dlhý.

- Zelené žaby žijú v rybníku.

- Červené žaby majú dlhé nosy.

pragmatický

- Červené idey majú dlhé nosy.

sémantický

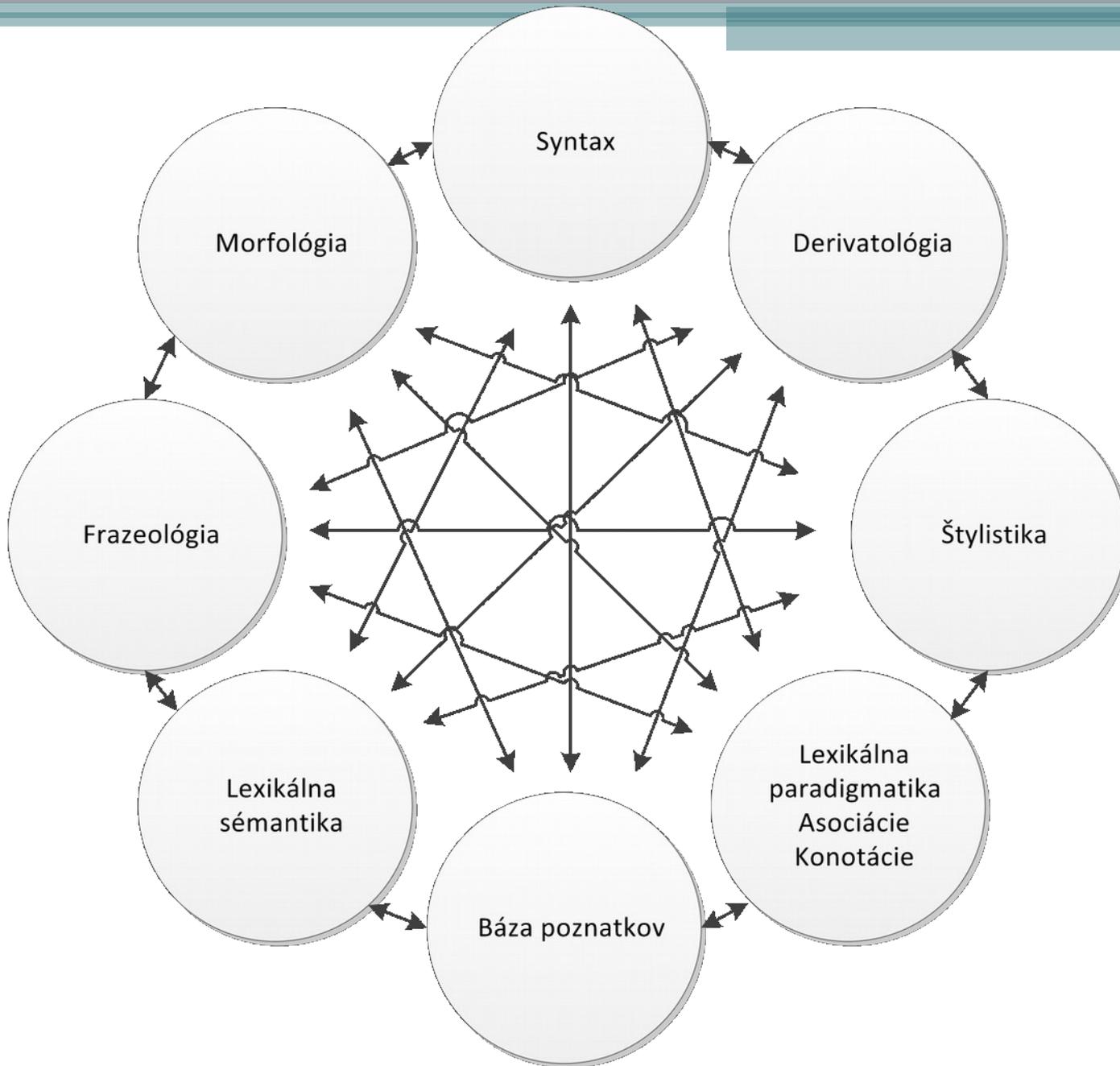
- Červenými ideou mať nos dlhý.

syntaktický

- Červenskami ideová maty nosník prídlžie.

- Zelené žaby žijú v rybníku.
- Červené žaby majú dlhé nosy. # pragmatický
- Červené idey majú dlhé nosy. # sémantický
- Červenými ideou mať nos dlhý. # syntaktický
- Červenskami ideová maty nosník prídlžie. # morfemický
- Cfvskmn idéé math noss d'ľha.

- Zelené žaby žijú v rybníku.
- Červené žaby majú dlhé nosy. # pragmaticky
- Červené idey majú dlhé nosy. # sémanticky
- Červenými ideou mať nos dlhý. # syntaktický
- Červenskami ideová maty nosník prídlžie. # morfematicky
- Cfvskmn idéé math noss d'ľha. # fonologicky



Spracovanie prirodzeného jazyka

- Predspracovanie textu
- Určovanie slovných druhov
- Lematizácia a stemovanie
- Rozpoznávanie vlastných pomenovaní
- Identifikácia viacslovných pojmov
- Rozpoznanie koreferencií
- Syntaktická analýza

Predspracovanie textu

1. Konverzia dokumentov na jednotné kódovanie
2. Extrakcia čistého textu
3. Tokenizácia a segmentácia
4. Normalizácia
5. Odstránenie/označenie stop slov

text.fiit.stuba.sk



STU
FIIT
SLOVAK UNIVERSITY OF
TECHNOLOGY IN BRATISLAVA
FACULTY OF INFORMATICS
AND INFORMATION TECHNOLOGIES



ÚLOHY SPRACOVANIA TEXTU SLUŽBY A NÁSTROJE KONTAKT

Spracovanie textu na FIIT STU

CHCEM VEDIĤ VIAC

Predspracovanie textu: Tokenizácia

Vstup: Joviálna (Ioviálna) planéta alebo plynný obor je planéta, ktorá je svojou veľkosťou a zložením podobná Jupiteru.

Výstup: [Joviálna][][(][Ioviálna][)][][planéta][]
[alebo][][plynný][][obor][][je][]
[planéta][,][][ktorá][][je][][svojou][]
[veľkosťou][][a][][zložením][]
[podobná][][Jupiteru][.]

Určovanie slovných druhov

Vstup: 'rybníkom'

Výstup: 'podstatné meno, mužský rod,
jednotné číslo, inštrumentál'
(a niekedy aj viac)

Rozpoznanie koreferencií

Monika miluje tenis. Dokáže obetovať hodiny, aby ho trénovala.

Rytier zosadol z koňa, aby sa napil.

Syntaktická analýza

