

Tento text vznikl jako pomocný text k přípravě na zkoušku z předmětu Úvod do počítačové lingvistiky v roce 2011. Může obsahovat mnoho chyb, dohadů a nepřesností. Tedy jeho autoři se vzdávají zodpovědnosti za důsledky z toho vyplývající... Za každou větší kapitolou následuje seznam možných otázek, jejichž zdrojem jsou wiki matfyz a matfyz fórum.

Markéta Popelová a Jakub Tomek.

MORFOLOGIE

- **Morfologie** studuje vztahy mezi jednotlivými částmi slov, vnitřní struktury slov. Zabývá se tvořením tvarů slov a jejich významem, dále i tvořením nových slov. Studuje způsoby skloňování (**deklinace**) a časování (**konjugace**).
- S morfologií se pojí pojmy:
 - **lexikologie** – slova jsou studována jako jednotky slovní zásoby,
 - **lexikografie** – sestavování slovníků,
 - **tvaroslovné dublety** – stejné slovní tvary odvozené od více slovních základů (žena, tři, hnát, stát, atd.), neboli slova víceznačná, která mají různé slovní druhy,
 - **alternace** – změna hlásek uvnitř kmene (vůz → vozu, švec → ševce, prkno → prken),
 - **alomorfy** – varianty kmene odvozené od stejného slovního základu (nejvíce změn má matka – matce – matek – matčin).
- Slova se dělí na:
 - **autosémantická** = plnovýznamová,
 - **syntémantická** = pomocná.
- Základní jednotkou je **morfém** = nejmenší jednotka nesoucí význam. Existují dva typy morfémů:
 - **lexikální morfém** = kmen slova – nese význam slova,
 - **gramatický morfém** – určuje gramatickou roli slovního tvaru.
- Např. *za-hrad-ou* má předponu „za“, lexikální morfém „hrad“ a gramatický morfém „ou“, který určuje 3 elementární jednotky (sémata): pád, číslo a rod.
- Morfologická typologie jazyků dělí jazyky na:
 - **analytické** (slovo = morfém) → **izolační** – každé slovo je morfém, bez předpon/přípon: vietnamština, čínština
 - **syntetické** (slovo > morfém)
 - **flektivní** – mají předpony, přípony, koncovky (ale míra řetězení je nějak omezená), daný tvar morfému nese více významů (koncovka určuje pár, rod, ...): latina, stará řečtina, slovanské jazyky
 - **aglutinační** – také různé předpony apod., ale jeden morfém nese jeden význam (tedy např. koncovka přidá jeden význam): maďarština, japonština
 - **polysyntetické** (slovo = věta): eskymácké a indiánské jazyky.
- Jak zpracovávat morfologii? Podle toho, na čem je morfologie založená:
 - Na **morfémech** – vidí slovo jako řetízek morfémů.
 - Na **lexémech** – vidí slovo jako výsledek aplikace pravidel, která slovo mění a tím vytváří nový tvar.
 - Na **slovech** – centrální roli mají **vzory**. Známe-li základní tvar slova a jeho vzor, dokážeme vygenerovat všechny zbylé tvary. Vhodné i pokud jeden morfém reprezentuje více gramatických kategorií (např. 3.os, sg., r.ž.), kde předchozí přístupy selhávají.
- **Two-Level Morphology**
 - Systém zpracování morfologie vyvinutý Lauri Karttunenem a Kimmo Koskeniemmin na začátku 80's. Pro češtinu nevhodné.
 - První obecný model zpracování morfologie přirozeného jazyka. Založen a na konečných stavových automatech a na nich definovaných oboustranných přechodech. Mechanismus morfologie byl pro všechny jazyky společný (to byl požadavek, aby to bylo nezávislé na jazyku), ale pro každý jazyk se musel vytvořit slovník a pravidla (přechody mezi stavy).
 - První úroveň lexikální, druhá povrchová.
 - Pravidla se aplikují paralelně. Podmínky se mohou vztahovat k oběma úrovním zároveň či k jedné z nich. Zároveň se prohledává slovník a uplatňují pravidla.
- **Česká morfologie**
 - Činnosti využívající morfologii:
 - **Morfologická analýza** – výsledkem je seznam lemmat a značek. Značky určují kombinace gramatických kategorií. Značky a lemma určují vstupní slovní tvar.

- Vytvářena od r. 1989 prof. Hajičem a kol. Využívá poziční **značky**. Tedy každé slovo má 15-ti písmennou značku, která ji určuje. Nicméně z těchto 15 písmen se využívá jen 13 (2 jsou rezervní), každá kategorii má své pořadí ve výsledné značce. Některé značky se vzájemně vylučují (např. příslovce nemá osobu), pak se na dané pozici píše pomlčka. Kromě značky se každému slovu ještě přiřadí jednoznačné **lemma**, což je základní tvar slova. Značka a lemma dohromady jednoznačně určují slovo a jeho tvar se všim všudy.
 - K některým tvarům daného slova pasuje více značek, pak se tam napíší oboje. Např. slovo „funkci“:
 - <f>funkci<MMl>funkce<MMt>NNFS3-----A----<MMt>NNFS4-----A----<MMt>NNFS6-----A----
 - Lemma je „funkci“ a značky jsou 3, neboť pád může být 3., 4. i 6.
 - Proto u jednoho slova bývá více značek (průměrně 4, nejhorší je to však u měkkých přídavných jmen jako „jarní“, která mají až 27-násobnou víceznačnost).
 - **Morfologické značkování (tagging)** – proces výběru jediné správné značky v daném kontextu. Tedy máme-li slovo v zadaném tvaru, můžeme mu přiřadit (několik?) lemmat a kombinací značek. Značkování pak dostane dané slovo větě a vybere jediné lemma a značku.
 - To je náročné, jsou různé přístupy (algoritmická pravidla/statisticky/kombinace), nejlépe to umí čistě statistické metody → úspěšnost až 96%.
 - **Česká morfologická desambiguace založená na pravidlech** – Oliva a Petkevič vytvořili pravidla, která platí bez výjimek. Pomocí těchto pravidel se vyškrtají nevhodné značky. Všechny ostatní nechá (asi 40%). Náhodou se takto ovšem přišlo na jinou zajímavou aplikaci – pokud se u některého slovního tvaru vyškrtá vše, znamená to, že ve větě musí být gramatická chyba. Tedy jednou z aplikací je **kontrola gramatických chyb**.
 - **Lemmatizace** – proces výběru správného základního tvaru (lemmatu), ze kterého byl odvozen daný vstupní tvar. Klíčová operace pro vyhledávání v textech. U nás má úspěšnost 99,9%.
 - **Generování** – proces výběru správného slovního tvaru ze zadaného lemmatu a množiny značek.
- **Kontrola překlepů jako aplikace morfologie**
 - Požadavky: najít a opravit všechny překlepy, přezkoušet kontextové podmínky korigované verze, neznámá slova se nemají hlásit jako chybná, nedávat falešná chybová hlášení, maximálně automatická korektura, krátký čas zpracování.
 - S tím souvisí následující pojmy:
 - **precision** = počet nahlášených chyb / počet nahlášených slov,
 - **recall** = počet nahlášených chyb / počet všech chyb.
 - Lepší je odhalit méně chyb, ale jen to, co jsou opravdu chyby. To znamená více se snažit maximalizovat precision (aby co nejvíce nahlášených slov byly opravdu chyby) a tolik neřešit recall (aby co nejvíce z chyb bylo odhaleno).
 - Používají se 2 základní metody:
 - **Porovnávání řetězců se slovy ve slovníku**. Buď se slovníkem všech možných slovních tvarů daného jazyka (**wordlist**) nebo se slovníkem lemmat a provádíme morfologickou analýzu.
 - Je to spolehlivé a jednoduché, ale pomalé, náročné na kvalitu slovníku, místo, nerozezná to chybná slova od neznámých a každé zlepšení musí zařídít autor či uživatel.
 - **Porovnávání skupin znaků (dvojic, trojic)** a hledání nedovolených kombinací znaků.
 - Je to nezávislé na slovníku a rychlé. Ale také neúplné a neodhalí překlepy ve slovech, která se skládají jen z vhodných kombinací znaků.
 - Možná vylepšení: vzít v úvahu okolnosti chyb (např. blízké klávesy), zohlednit statistiku chyb, časté pravopisné chyby (mně x mě, jsem x jsme), různé heuristiky na oddělení chyb a neznámých slov, zapojení syntaxe a sémantiky, pracovat s kontextem (např. porovnávat korpusy).
 - **Systém AZIMUT (Automatická Selektce Informací Metodou Úplného Textu)**
 - Vznikl 1990 Králíková, Panevová. Měl dva moduly:
 - **Vyhledávací modul**: Sloužil pro automatické vyhledávání ohýbaných slov v textu na základě parametrů. Na vstupu byla podstatná a přídavná jména v základním tvaru a k nim různé operátory (! = vyskočovat, -číslo- = obě slova musí být bezprostředně za sebou/ob slovo/ve stejné větě/odstavci, apod.). Takto se mohlo nechat v textu vyhledat: „vzdálenost!, rodinný! -1- domek!“, což znamená všechny výrazy, které obsahují buď slovo „vzdálenost“ v nějakém tvaru, či souloví „rodinný domek“, opět v libovolném tvaru (neřeší, jestli gramaticky správném). Předpokládá členění textu na slova, věty, odstavce. Tento modul není tak podstatný.
 - **Jazykový modul**: Modul pro automatické skloňování českých slov. Pro dané vstupní slovo vrátí všechny jeho možné tvary.
 - Využívá **retrográdní slovník** dr. Slavičkové (1975), tedy slovník seřazený dle písmen slov odzadu

- (háčky a čárky bere jako samostatná písmena). Vychází z myšlenky, že slova se stejnou koncovkou se skloňují stejně. A výjimky z tohoto pravidla je možno uložit do zvláštního slovníku (jsou jich stovky, maximálně tisíce, což není mnoho).
- Retrográdní slovník nevyužívá přímo, ale na jeho základě byl vytvořen klíč pro určování vzorů slov (seznam pravidel dle konců slov).
 - Porovnává písmena vstupního slova (musí být v základním tvaru!, tedy 1. pád sg.) odzadu, dokud nenajde jednoznačný vzor pro skloňování. Pak slovo vyskloňuje dle vzoru. (Umí i základní alternace.)
 - Má různé problémy: ne vždy lze jednoznačně určit vzor (právnick i trávnick mají stejnou koncovku, ale liší se v životnosti), problém **přegenerování** (systém vygeneruje i neexistující tvary), malý rozsah retrográdního slovníku (je tedy nutno přidávat výjimky), pro slovesa už nefunguje spolehlivě.
- Další pojmy:
 - **Negativní slovník** obsahuje ta slova, která nejsou při vyhledávání v textu důležitá (spojky, citoslovce), proto jsou odstraněna z textu ještě před vyhledáváním.
 - **Konkordance**. Všem slovům mimo negativní slovník byla přiřazena adresa a frekvence výskytu v textu. Slova z negativního slovníku dostala jen adresu (kvůli počítání vzdáleností). Samotné vyhledávání pak neprobíhalo na původním textu, ale na této konkordanci.
 - **Systém MOZAIC (Morphologically/Morphemic Oriented System of Automatic Indexing and Condensation)**
 - 70's MFF Kirschner a kol. Systém pro indexaci dokumentů, tvoření souhrnů, seznamů klíčových slov. Podobně jako ASIMUT nepoužívá rozsáhlé slovníky, ale lingvistické poznatky.
 - Využívá toho, že řada přípon a koncovek nese význam (Aj: -er/-or konatel děje, -tion činnost, -ity/-ness vlastnosti; Čj: -ič/-ač/-čka/-ér/-or/-dlo/-metr/-graf/-fon/-skop přístroje a nástroje, -ací/-ecí účel, atd.)
 - Algoritmus indexování textu:
 - Na vstupu je čistý text se zachovanými typografickým členěním.
 - Lematizace a morfologická analýza → získáme lemata a morfologické značky.
 - Jsou odstraněna lemata, jejichž kmen nemá vztah k dané tematické oblasti (k tomu se využívá malý negativní slovník), či jsou příliš krátká nebo obsahují nevhodné kombinace hlásek.
 - Syntaktická analýza jmenných skupin. Využívá jednoduchou gramatiku v jazyce Systémů Q. To pomůže odhalit tematicky významné několikaslavné termíny (operační zesilovač TESLA KC 415).
 - Vážené ohodnocení termínů podle důležitosti. Záleží na tom, v jak důležité části textu jsou (nadpis, první/poslední odstavec, první/poslední věta). Váhy jsou exponenciální.
 - Normalizace vah vzhledem k délce dokumentu. (Nejčastější termín získá 100 bodů, zbytek poměrně.) Umožňuje porovnávat různě dlouhé dokumenty.
 - Výstupem je 10 nejvýznamnějších termínů, seřazených podle četnosti výskytu.
 - Výhody: Není nutno vytvářet specializované slovníky odborných termínů. Stačí množiny relevantních koncovek a přípon, daný negativní slovník a několik pravidel. Lokální syntaktická analýza umožňuje větší flexibilitu při hledání termínů.
 - Problémy: Neřeší nevyjádřené podměty, odkazování v textu pomocí zájmen apod. Pracné vymýšlení negativního slovníku, pravidel a koncovek.

Úvod do počítačové lingvistiky – příprava na zkoušku

Otázky:

1. Alomorfy.
2. Co je to morfém a jak ho klasifikujeme?
3. Lemmatizace - co to je a kde se používá.
4. Two-level morphology.
5. Morfologie a kontrola pravopisu.
6. **Kontrola překlepů.** Podrobně popište metod používané při automatické kontrole překlepů (i volby dialogu s uživatelem).
7. **ASIMUT** (Co to je, jak funguje, na čem je založen jeho jazykový modul.)
8. **MOZAIKA** (Používá se u systémů MOSAIC syntaktická analýza? Proč ano/proč ne.)

SYNTAX

- **Syntax** (skladba) se zabývá vztahy mezi slovy ve větě, tvořením větných konstrukcí, slovosledem.
- Jsou dva typické zápisy syntaxe věty:
 - **Závislostní strom**
 - Velmi se podobá větnému rozboru ze základní školy. Nicméně kořen je jediný a obsahuje přísudek. (**Co když má věta více sloves?**). Uzly jsou právě slova věty. Každé slovo závisí na jiném, závislost je popsána orientovanou hranou mezi těmito slovy. Kromě uzlů, hran, ohodnocení uzlů (např.) a ohodnocení hran, si ještě závislostní strom pamatuje původní slovosled věty jako úplné uspořádání na slovech. Strom se pak může vykreslovat klasicky – či tak, aby byl slovosled zachován.
 - Závislosti zachycuje jednoduše a přehledně. Nicméně ne vždy je jednoduché určit, jak (a zda) jsou slova závislá. Proto mají závislostní stromy další možnosti pro zaznamenávání následujících jevů:
 - **Koordinace** – různé větné členy se stejnou sémantickou rolí. Např. *Jan a Marie; černý nebo bílý.*
 - **Apozice** – různé větné členy se stejnou syntaktickou rolí, shodnou gramatickou kategorií (tzn. gramaticky kongruentní). Např. *Matematicko-fyzikální fakulta (MFF); Ivo Truchlivý, učitel matematiky.*
 - **Parenze** (vsuvka) – věta či větný člen, který syntakticky nesouvisí s okolím, ale upřesňuje, o čem se v okolí mluví. Např. *Mohl bych, prosím, zavřít okno?*
 - Jak s předložkami? Jejich funkcí je vlastně pád. Má být závislá předložka na podstatném jméne nebo naopak?
 - **Složkový (derivační) strom**
 - Odpovídá derivačnímu stromu bezkontextové gramatiky. Tedy věta se rozdělí do částí, které se zase rozdělí do částí, atd. Tedy slova věty (tokeny) odpovídají listům stromu.
 - Je méně přehledný, má větší množství uzlů a má jeden hlavní problém – přirozené jazyky nebývají bezkontextové. Problém mu tedy činí právě neprojektivní konstrukce.
 - Dá se znázornit pouhým uzavorkováním věty, kde uvnitř závorok jsou vždy právě dva prvky, kde prvek je buď jiný uzavorkovaný výraz, či samotné slovo.
 - **Neprojektivní konstrukce**
 - Neprojektivní konstrukce je závislost mezi dvěma slovy ve větě oddělenými slovem třetím, které (ani nepřímo) nezávisí na žádném z nich. Např. Soubor se nepodařilo otevřít. Závislostní strom s tím nemá problém (jen hrany se v něm jakoby kříží). Složkový ano.
 - V češtině jsou běžné, ale jsou i v jiných jazycích.
 - **Transformační gramatika**
 - Historie. Předválečná americká lingvistika se snažila explicitně popsat jazyková pravidla. Tyto směry by se daly považovat za předchůdce transformační gramatiky:
 - **Deskriptivismus** (1933 Bloomfield). Jazyková fakta popisuje, klasifikuje a registruje, ale nevysvětluje. Zabývá se spíše povrchovou větnou strukturou.
 - **Analytická syntax** (1937 Jespersen).
 - **Logický přístup** (1935 Ajdukiewicz) – kategoriální gramatika.
 - Obecně se zavádí **koncept povrchové a hloubkové syntaktické struktury** (surface & deep structure). Povrchová struktura řeší spíše zápis, hloubková význam. Pak je běžné, že jedné povrchové reprezentaci může odpovídat více hloubkových (jedna věta je významově víceznačná), stejně jako naopak (jeden význam se dá vyjádřit různě).
 - Jazyk však nebyl dosud popsán formální matematickou strukturou. Spíše se popisovalo, než že by se vysvětlovalo. Míchala se syntax a sémantika. Syntaktické jevy se popisovaly pomocí sémantiky apod.
 - Noam Chomsky 1957 *Syntactic Structures*. V této knize popsal 3 základní komponenty:
 - **Báze**. Soubor bezkontextových pravidel. Tato pravidla generují složkové stromy, tzv. **frázové ukazatele** (phrase makers). $S \rightarrow NP VP$, kde *NP* je noun phrase, *VP* je verb phrase.
 - **Transformační komponenta**. Soubor transformačních pravidel nad frázovými ukazateli. Z nich vytváří povrchovou strukturu věty. Dva typy transformačních pravidel:
 - **Obligatorní** – transformace musí být provedena (pokud je to možné).
 - **Fakultativní** – transformace je volitelná.
 - **Fonologická komponenta**. Soubor regulárních přepisovacích pravidel. Řetězcům morfémů přidělují fonetickou interpretaci a význam. (*Fonetika* i *fonologie* zkoumají zvukovou stránku jazyka. Fonetika zkoumá, jak se hlásky v těle tvoří a vnímají, zatímco fonologie zkoumá funkci hlásek a zvukové rozdíly, které mají v jazyce nějakou funkci.)
 - Cílem je tvořit věty. Jsou vytvářeny **generativní procedurou**, která používá různá přepisovací pravidla (někdy kontextová, jindy bezkontextová). Není ale schopná zachytit vztahy mezi variantami vět, např.

mezi větou tázací a oznamovací.

- **Transformace** (v transformační komponentě) jsou definovány **strukturním indexem** řetězců (řez stromem, výraz se matchuje na množinu vrcholů) a **strukturní změnou** (co se má s namatchovanými vrcholy provést).
- Pravidla mohou být bezkontextová. Pak má tato složka sílu Turingova stroje, což je moc. V dalších verzích byla tato složka oslabena.
- Vývoj transformační gramatiky.
 - 1965 Aspects of the Theory of Syntax (N. Chomsky) = *Standard Theory*.
 - 1968 *Extended Standard Theory*
 - 1980's *Government-binding Theory* (GB) – založená na obecných principech univerzální gramatiky a parametrech platných pro daný jazyk.
 - 1990's *Teorie minimalismu* – obsahuje jen dvě roviny – opět Chomsky:
 - **Rovina logické formy** (LF) – reprezentace jazyka a významu.
 - **Fonetická rovina** (PF) – zvuková stránka jazyka.
- **Tree Adjoining Grammars** (TAG)
 - Elementární strukturou jsou stromy. Formalismus pro popis gramatik, ale nefunguje tak, že by se řetězec přepisoval řetězcem, ale v gramatickém stromu se uzly nahrazují jinými stromy (např. X miluje Y=Emil, tak za X se dá dosadit „Milan“, ale i strom „Milan, Ferda a Dežo“), ale jen tehdy, když se uzel a kořen stromu shodují. Uzly, které je možno substituovat jsou označeny šipkou. Proces končí, když už nelze žádný uzel nahradit. Svou generativní silou mohou dosahovat až kontextových gramatik (po drobných modifikacích).
 - Typy základních stromů:
 - **základní** (initial) strom: Udává valenční vztahy a strukturu věty
 - **pomocný** (auxiliary) strom: Pomocí těchto se tvoří rekurze ve stromu
 - Typy změn:
 - **Substitute** – list stromu je nahrazen pomocným stromem, jehož kořen je značený stejně, jako list původního stromu.
 - **Adjungace** – vnitřní uzel je nahrazen pomocným strom, kořen opět značen stejně jako list původního stromu.
- **Lexical-Functional Grammar** (LFG)
 - Využívá dva typy struktur:
 - **c-struktura** (constituent). Spojuje slova do frází. Datový typ je složkový strom.
 - **f-struktura** (functional). Reprezentuje funkční vztahy ve větě (např. vazby sloves). Používá datový typ matice atribut-hodnota (že by mapa).
 - Každá c-struktura se spojuje s jednou f-strukturou. Opačně jich může být i více.
- **Unifikační gramatiky**
 - Každý objekt je reprezentován množinou vlastností, tzv. rysů. Stylem <název_vlastnosti>: <hodnota_vlastnosti>. Této množině říkáme **sestava rysů**. Vlastnosti mohou být např. grafémický zápis, slovní druh, rod, číslo, pád, atd. Jejich hodnotami mohou být i další sestavy rysů či proměnné.
 - Máme-li dvě sestavy rysů popisující objekt, můžeme je unifikovat, jestliže nejsou v konfliktu v žádné vlastnosti.
 - Problém je, že lze unifikovat i vlastnosti, které spolu nesouvisí, třeba pád podmětu a způsob přísudku.
 - **Typové sestavy rysů** využívají toho, že některé typy objektů mají společné vlastnosti. Např. u sloves určujeme osobu, číslo, čas, způsob atd. Typy se pak většinou řadí hierarchicky. Slova se dělí na ohebné a neohebné druhy. Ohebné zase na časované a skloňované. Atd.
 - Příklady:
 - **FUG** (Funkční unifikační gramatika). Martin Kay.
 - **HPSG** (Head Driven Phrase Structure Grammar).
 - Zahnuje principy, gramatická pravidla a slovníkové položky (tříděné, dle různých kategorií). Slovo má dva základní rysy – **phon** (zvuk, fonetická forma) a **synsem** (syntaktické, sémantické informace) – tyto rysy dále děleny. Základním typem je **znak**, který se dělí na **slova** a **fráze**.
 - **GPSG** (Generalized Phrase Structure Grammar). 1985.
 - Dnes to ale převládaly statistické metody.
- **Kategoriální gramatiky**
 - Každému slovu přiřazena **kategorie** – množina syntaktických vlastností daného slova.
 - Zápis kategorií X/Y, či X\Y, podle toho, zda argument je vlevo, či vpravo.
 - Dvě základní pravidla: X/Y Y → X; Y X\Y → X.

- Nástroje na syntaktickou analýzu.
 - **Augmented Transition Networks** (Woods, 1970). Neboli **rozšířené přechodové systémy**.
 - Má tři typy hran: CAT (přechod do stavu, nalezne-li příslušnou kategorii), JUMP (přechod do stavu bez hledání kategorie), SEEK (přechod k podsíti).
 - **Q-systémy** (Colmerauer – otec prologu, 1969).
 - Formalismus pro transformaci grafů. Grafy jsou linearizovány, např. S(NP,VP(V,NP)).
 - Nějak to vytvoří takový zvláštní graf s uzly přímo těmi stromy (resp. jejich kategoriemi) a mezi nimi zakresluje lineární hrany. Z toho to pak vytváří pravidla. Nějak to pracuje s proměnnými, které jsou značeny *.
 - Používal je např. RUSLAN. Ale jinak byly oblíbené téměř jen ve fancouzštině. Jinde byly oblíbenější rozšířené přechodové systémy.
- **Funkční generativní popis** (Sgall, 1967)
 - **Stratifikační teorie** – řeší popis na více vrstvách:
 - fonetická
 - fonologická
 - morfématická
 - povrchová
 - tektogramatická.
 - Princip **forem** a **funkcí** – jednotka na vyšší rovině reprezentuje funkci jednotky na nižší rovině.
 - Na vyšších úrovních (povrchová a tektogramatická) se jazyk popisuje **závislostní reprezentací**, typicky závislostními stromy. Teorie valence, viz níže.
- **Valence**
 - Schopnost některých slov (především sloves) „vyžadovat“ jiné větné členy a tvořit s nimi věty.
 - V dané teorii se rozlišují na tektogramatické rovině dva typy základních členů:
 - **aktant** – může být ve větě pouze jednou
 - **volné doplnění** – může být vícekrát
 - Dále vazby se v TG rovině dělí na **obligatorní** (povinné) a **fakultativní** (nepovinné). Obligatorní nesmí na tektogramatické rovině chybět. Nicméně smí chybět na povrchové rovině, známe-li ho např. z kontextu.
 - **Valenční rámec** – množina aktantů a obligatorních volných doplnění ve větě.
 - 2007 sestavili Lopatková a Žabokrtský **Vallex** – valenční slovník.
- Kontrola gramatické správnosti:
 - Co lze kontrolovat?
 - shoda podmětu s přísudkem
 - interpunkce
 - neprojektivní konstrukce
 - zájmena (mě/mně)
 - Jak kontrolovat?
 - **Chybové vzorky**. Vhodné hlavně pro jazyky s pevným slovosledem), kde chybné konstrukce zůstávají v lokálním kontextu a nerozlézají se daleko po větě.
 - **Gramatika** Nelze ale rozeznat, zda je věta špatně, nebo zda je správně vzhledem k neúplné gramatice. To se snaží řešit RFODG (Robust Free-Order Dependency Grammar).
 - RFODG: Výpočet probíhá ve fázích. Interpret gramatiky rozhoduje, jak se bude stejné gramatické pravidlo používat. Je snaha o co nejplynulejší fázování výpočtu, což 2001 zlepšil např. i Holan.
 - LanGr: 2003 Pavel Květoň. To, co používá MS Word. Pravidla psána ručně na základě korpusu, pracují v cyklech. Snaha o vysokou precision (85%) oproti recall (30%), uživatele „otravuje“ tak jednou za 3-4 stránky. Každé pravidlo má 4 části: kontext, desambiguace, report a akce. Pravidla byla tvořena speciálně pro češtinu. Pro jiný jazyk by byla téměř úplně jiná.
 - Používá desambiguaci na to, že když se odstraní všechny tagy, tak víme, že je něco špatně. V tu chvíli se ale musí určit, co je špatně a jak to opravit. (A toto samozřejmě neopraví všechny chyby.)
 - Neřeší to ten problém, že věta může být správně, ale až v dalekém kontextu přes jiné věty. (*Tatínek šel do práce.*)
 - Je problém, jak hodnotit kvalitu grammer-checkeru. Je nutno to dělat ručně.
 - Obecně používá tuto přípravu (klasický postup při zpracování psaného textu): segmentace (rozseká na věty) → tokenizace (rozseká na slova) → morfologická analýza (každému tokenu dá seznam dvojic lemma – tag) → morfologická desambiguace (každému tokenu vybere ideálně jeden token) → syntaktická analýza (větný rozbor) → sémantická analýza (rozbor významu věty).

Úvod do počítačové lingvistiky – příprava na zkoušku

Otázky:

1. Syntaktická analýza.
2. Závislostní (D-tree) a složkový (C-tree) strom pro větu "Ve včerejším závodě startovali výborní skokani."
3. Převedte složkový strom na závislostní.
4. Chomského transformační gramatika (v jiné otázce jako teorie popsaná v knize Syntactic structures).
5. Na co slouží strukturní index u Chomského gramatiky?
6. Sestavy rysů a jejich použití.
7. Co jsou to unifikační gramatiky, jejich výhody, nevýhody.
8. Hlubková a povrchní syntaxe, vztahy mezi nimi.
9. Valence.
10. Co je to HPSG, LFG, FGD.
11. Co znamená zkratka TAG, stručně vysvětlit princip.
12. Na čem je založena teorie funkčního generativního popisu.
13. Teorie minimalismu - autor a na které teorii navazuje.

STROJOVÝ PŘEKLAD

- Problémy strojového překladu:
 - Cílem překladu není převést slovo na slovo, ani větu na větu, ale převést sdělení z jednoho jazyka do druhého, aby se dalo pochopit, co chtělo říci.
 - Nestačí tedy překládat doslovně, dokonce nestačí ani překládat se znalostí morfologie, ustálených slovních spojení, či syntaxe, ale je třeba znát i kontext, který v jazyce obsažen vůbec není.
 - Různé počty výrazů pro určité slovo, některé jazyky jsou „jemnější“ (např. eskymáci budou asi mít širší rybářskou slovní zásobu nežli Němci). Dále podobné významy pro danou oblast nemusejí jít přímo namapovat na obdobné výrazy v jiném jazyce (např. vaření v Aj a Japonštině).
 - Význam slova závisí na kontextu – „otevřít“ může být různé pro program, okno, plechovku...
- Historie ve světě:
 - Automatickým překladem se zabývali již od 1933, postupně zkoušeli překlady slovo od slova → slovníky pro předpony, kmeny, přípony a koncovky zvláště → zaveden **preediting** (ručně se text oseká o různé zvláštnosti, zkrátí se věty, víceznačná slova se nahradí jednoznačnými – a stroj pak tyto snadné věty už přeloží) a **postediting** (stroj přeloží, co umí – a zbytek přeloží člověk – pro lidi to je spíše otrava) → zaveden **pivotní jazyk** (ale politický problém, který jazyk se má zvolit + kumulují se zbytečné chyby, např. pro překlad příbuzných jazyků např. Čs-Aj-Sk) → experiment překladu jednoduchých vět mezi Aj a R (úspěšné) → v roce 1966 **Alpac** – zpráva, která v USA utlumila výzkum, tedy objevy po tomto roce především odjinud.
 - **TAUM METEO** (1976) – z Montrealu, překlad meteorologických zpráv z angličtiny do francouzštiny – byla vymezena a rozumně omezena podmnožina syntaxe a sémantiky. Díky vhodné implementaci (Q-systém) šlo rozpoznat, že se neví, jak text přeložit. Což se pak udělalo ručně.
 - **Systran**. Překlad dokumentů EU, přímo (každý pár jazyků odděleně, uspokojivě jen několik málo prvních A-F-N). Data oddělena od programu.
 - **EUROTRA**. Mělo nahradit systran, ale každý s každým se měl dohodnout (opět přímý překlad) na analýze, dohodnout rozhraní atd. – nezládlo se. Příliš megalomanské (72 jazykových párů).
 - **VERBMOBIL** – Německý nástupce Eurotry, v universitním prostředí, překlad mluvené řeči, domluva obchodníků na další schůzce. V současnosti asi nežije.
- Historie v ČR:
 - **APAČ** (Kirschner, 1980's): Z angličtiny do češtiny – překlad z oblasti vodních pump.
 - Využíval **transdukční slovník** – překladač koncovek (-ation → -ace; -ic → -ický) + slovník s asi 1500 výrazy.
 - **RUSLAN**: Překlad manuálů k OS sálových počítačů – pomalý (1 věta trvala asi 4 minuty).
 - Slovník o velikosti cca 8500 slov + transdukční slovník.
 - Využit transfer.
 - Gramatika zapsána pomocí Q-systémů (TAUM)
 - Navíc záchranná pravidla pro případ problémů při analýze.
 - **Česilko** – lokalizace velkých SW systémů.
 - Snaha minimalizovat podíl na překladu.
 - Místo lokalizace z jazyka typově odlišného, bylo myšlenkou překládat z jazyka blízkého.
 - Typ překladu – FAHQ (plně automatický, vysoce kvalitní), blízké jazyky, plně morfologické slovníky, statistická analýza češtiny.
 - Blízké jazyky typu Čj-Sk mívají shodnou syntaxi a slovosled, jiné slovníky, avšak ne úplně, odlišné tvarosloví. Např. tedy funguje doslovný překlad slovo od slova.
 - Myšlenkou bylo, pomocí lidí přeložit ze zdrojového jazyka test do češtiny, a z češtiny strojově do polštiny, slovenštiny, ruštiny...
- Dnešní situace: Neexistují obecně použitelné systémy, přitom překlad je potřebný, určitá automatizace se hodí; je třeba spojit síly člověka a počítače.
 - **Překlad podpořený počítačem** (Computer Assisted Translation) – není strojový překlad:
 - Pracuje s **překladovou pamětí** – text je dělen na segmenty (věty, či polygramy), které lidský překladatel překládá a systém ukládá dvojici: text v původním jazyce a překlad. Pokud se v textu časem objeví podobná, či stejná fráze, bude nabídnut uložený překlad. Je vhodné pro techničtější texty, pro beletrii ne tolik.
 - Taktéž využívá **terminologickou databázi** – opět plněna ručně, sestává z konkrétních termínů zdrojového i cílového jazyka – např. carbon dioxide – oxid uhličitý (ne dioxid uhlíku apod.).
 - Dnes se kombinuje se statistickým překladem.
 - **Strojový překlad** (v základu bez intervence člověka)
 - **Pravidlový** – systémy se snaží překládat slova a pomocí pravidel je k sobě skládat.
 - **Statistický**
 - Především se využívají korpusová data (dvojjazyčná).

- Základním přístupem je generování určitého počtu možných překladů, kterým je pak přisouzena pravděpodobnost, že se jedná o překlad správný.
- Pokud jsou použity specializované korpusy, dobré výsledky jsou především v dané oblasti, není to úplně zobecnitelné.
- Jednoduchý pravděpodobnostní model – uváží se frekvence slova v trénovacích datech, což je výsledná pravděpodobnost slova.
- **Model zašuměného kanálu (Noisy Channel Model):**
 - Idea byla překlad z francouzštiny do angličtiny.
 - Uvažoval se model $P(e|f)$ – udává pravděpodobnost anglické věty e , za předpokladu francouzské věty f . Model zašuměného kanálu má dvě složky:
 - $P(e)$ – jazykový model – např. trigramový model (z libovolných dat, ne nutně paralelní korpus).
 - $P(f|e)$ – překladový model – trénovaný z paralelního korpusu francouzsko-anglického.
 - Pak $P(e|f) = P(e, f) / P(f)$
 - Na základě překladového modelu byli vytvořeni určití kandidáti, pomocí $P(e)$ – jazykového modelu, vybráno mezi nimi.
- K hodnocení používán tzv. **Bleu index** – porovnává kvalitu automatického překladu vůči lidskému překladu. Bohužel nebere v úvahu morfologii – drobnosti typu chybné koncovky, které srozumitelnosti nebrání, jsou brány stejně jako zcela špatný překlad.
- Jiné dělení překladů:
 - **Přímý překlad** – z jazyka do jazyka. Problém je, že je pro n států třeba hodně párů překladačů.
 - **Přes mezijazyk:**
 - **Pivotní jazyk** – pokud je třeba překládat mezi více jazyky, místo stylu „každý z každým“ se každý naučí převést text ze svého jazyka do jazyka pivotního, typicky přirozeného jazyka.
 - **Interlingua** – hypotetický formální logický zápis sdělení. Konstrukce obecné interlingui zatím moc neexistuje, neboť význam se těžko zapisuje (viz kapitola sémantická). Když se interlingua používá, bývá to umělý mezijazyk, volně založený na románských jazycích.
 - S Interlinguou se pojí tzv. **Vauquoisův trojúhelník** – čím vyšší patro, tím obtížnější provedení:

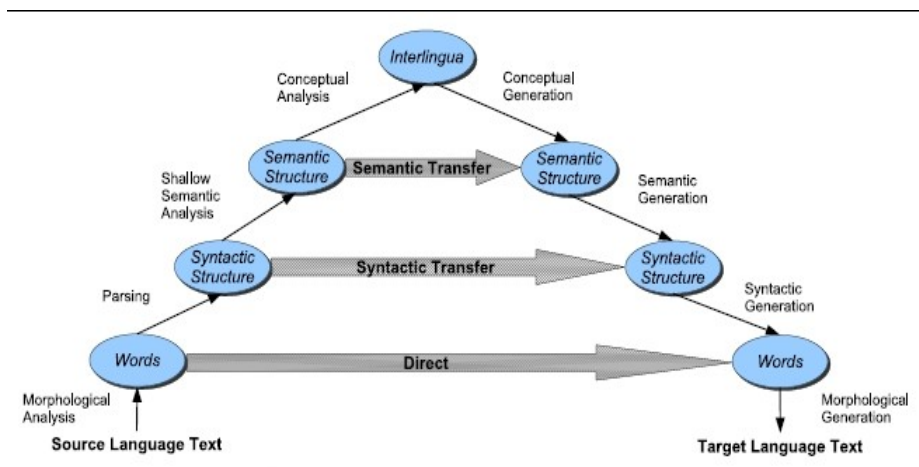
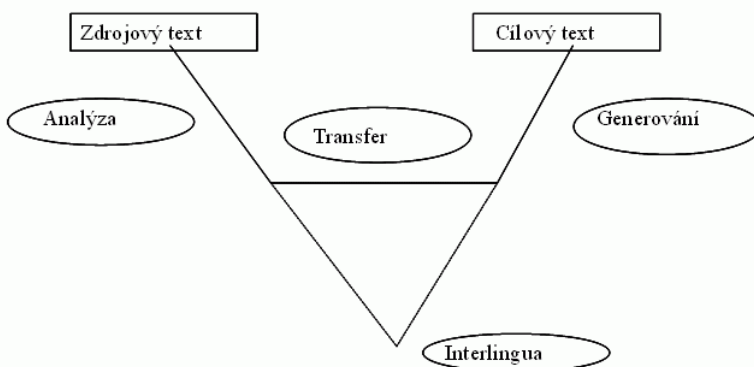


Figure 24.3 The Vauquois triangle.

- Typicky se používá následující schéma: Text v prvním jazyce projde morfologickou, syntaktickou a sémantickou analýzou. Dále proběhne **transfer** do interlingui. A z ní se generuje text v novém jazyce. Některé části dělají lidé ručně (nějaké části syntaktické a hl. asi sémantické analýzy), některé se dějí strojově (např. ten samotný transfer). Opět do souvislého pěkného textu to asi upravuje typicky zase člověk.



Úvod do počítačové lingvistiky – příprava na zkoušku

Otázky:

1. **Strojový překlad** (kategorie, principy u jednotlivých kategorií + uvést příklady, metody využívající člověka).
2. **Noise Channel** v překladu. (Podrobně popište automatický překlad metodou zašuměného kanálu.)
3. Překladová paměť.
4. Interlingua a k čemu se používá, rozdíl mezi interlinguou a pivotním jazykem.
5. Transfer v automatickém překladě.
6. České překladové systémy.
7. Česílko.
8. Ruslan.
9. Popište Vauquoisův trojúhelník. (trojúhelník s interlinguou na vrcholu).

KORPUSOVÁ LINGVISTIKA

- **Korpus** je rozsáhlý soubor textů (v digitální podobě) v daném jazyku, většinou anotovaný (označkováný) na základě přechodí morfologické a někdy i syntaktické analýzy. Je to cenný soubor dat, ale někdy se chybně považuje za reprezentativní vzorek či rovnou celý jazyk.
- **Brown Corpus of Standard American English**
 - První moderní elektronický korpus (1961 W.N.Francis a H.Kučera). Skládal se z textů, které toho roku v Americe vyšly (v novinách, krásné literatuře apod.) Obsahoval milion slov, 15 druhů textu (novinové reportáže, humor, krásná literatura, ...), dohromady 500 textů každý cca 2000 slov. Texty byly vybírány schválně náhodně. Celé to bylo pečlivé, ale milion slov není moc. A texty nebyly anotované.
- **Penn Treebank**
 - První a nejznámější syntakticky anotovaný korpus (1990's Univerzita v Pensylvánii). Také milion slov, asi 2500 článků. Všechno byly ale články z Wall Street Journal za poslední 3 roky, což je dosti omezující. Navíc články byly různě dlouhé.
 - Syntaktická analýza využívala složkové systémy. Tedy anotace pomocí uzávorkování a různých značek.
 - Byla snaha ho přeložit do češtiny (PCEDT), což se podařilo, ale s obtížemi – vyžadovalo to lidi, kteří by uměli dobře česky i anglicky a vyznali se v prostředí burzovních textů z 90's. (Nejen překladatel, ale i ten, kdo to reviduje.) Motivace pro překlad včetně podobného značkování byl takový aby se nějaký statistický program mohl učit rozdily.
- **Český národní korpus (CNC)**
 - Od 1994 společně UK, MU a Ústav pro jazyk český. Morfologicky označkován („otaggovaný“) - ne ručně, ale automatickými nástroji pro morfologickou analýzu. Současně obsahuje 500 miliónů slov a je složený z převážné části z novinových článků, dále z literatury a odborných textů.
 - Morfologická analýza používá ony výše popsané 15-ti poziční značky. V korpusu je rozeznáno 700 000 lemat, 15 miliónů slovních forem a po stochastické desambiguaci zůstane u každého slova průměrně 4,29 tagů. Používá statistické metody. Na učení se využívá ručně označkováný korpus s 1,2 milióny tokenů (slov). K automatickému učení používá kontextová pravidla (asi 11 000 pravidel). Automaticky určuje váhy. Dosahuje rychlosti 200 tokenů za sekundu a výsledná úspěšnost je přes 94%.
- **Pražský závislostní korpus (Prague Dependency Treebank)**
 - Od 1967 P. Sgall, dále i ostatní. Obsahuje 100 000 vět a 1,25 miliónů slov. Je anotovaný (opět automaticky) na několika rovinách (na některých rovinách jsou anotovány jen jeho části):
 - **Slovní rovina** (*w-rovina*). Pouze surový text bez anotace, ovšem včetně členění.
 - **Morfologická rovina** (*m-rovina*).
 - Každému slovu ve větě přiřadí několik atributů (lemma, tag (15-ti poziční značka), jednoznačné id využitě při propojování rovin, odkaz do slovní roviny, atd.).
 - Anotace probíhala dvoufázově: nejdříve anotoval automatický morfologický analyzátor → a pak dva lidští anotátoři na sobě nezávisle vybírali správná lemmata a tagy z výsledků automatického → nakonec třetí lidský anotátor vybral nejlepší možnost z předchozích dvou.
 - **Analytická rovina** (*a-rovina*).
 - Každá věta je reprezentována stromem orientovaným do kořene s ohodnocenými hranami mezi uzly. Uzly jsou právě prvky morfologické roviny, hrany jsou ohodnoceny podle závislostních vztahů uzlů, či určují další jevy (koordinace – s předchozí větou, apozice, interpunkce). Každý uzel si i pamatuje své pořadí ve větě kvůli grafickému znázornění.
 - Byl použit automatický parser na předzpracování textu a dále automatický nástroj, který na základě pravidel určoval ohodnocení hran, ale výstup byl často chybný či neúplný, tedy museli nastoupit ruční anotátoři. Následně byly provedeny automatické kontrolní testy (např. slovesný jmenný predikát závisí na být) a porušení byla ručně opravena.
 - Nakonec byla provedena společná revize morfologické a analytické roviny (např. shoda v pádě, rodu a čísle závislého a nadřazeného uzlu, atd.).
 - **Tektogramatická rovina** (*t-rovina*).
 - Opět je každá věta reprezentována stromem. Nicméně jeho uzly už nemusí být právě prvky morfologické analýzy (některé prvky zde nejsou (např. předložky) a některé uzly tu jsou navíc (např. nevyjádřený podmět)). Zachycuje hloubkovou strukturu věty. K některým uzlům jsou připojeny gramatémy poskytující o uzlu informaci, kterou nelze jinak odvodit. K uzlům reprezentujícím sloveso či některé typy podstatným jmen je přiřazen valenční rámec (odkaz do vallexu). Dále nějaké koreference.

Otázky:

1. **Korpusy**. Charakterizovat korpusy, které jsme probírali (zdroje textu, co je v nich značkováno atd.). K čemu jsou korpusy dobré v teoretickém i aplikovaném výzkumu?
2. Brownův korpus.
3. PennTreeBank.
4. Český národní korpus (složení, velikost, typy značek).
5. Co víte o Pražském závislostním korpuse? (Tady toho chtěl trochu víc - velikost, zdroj, použité značky, ...)

PRAVDĚPODOBNOSTNÍ A STATISTICKÉ METODY V AUTOMATICKÉM PŘEKLADU

- Motivace: víme, že existují 3 překlady pro dané slovo. Je těžké určit, který je pro danou situaci vhodný. Nicméně mohl by nám k tomu pomoci kontext okolních slov. Statistické překladové metody v podstatě zkoumají, jakou mají různé kombinace slov v daném jazyce pravděpodobnost – a dle toho se rozhodují o překladu.
- Pravděpodobnost výskytu slova w v textu T je $P(w)$ = počet výskytů slova S v textu T / počet slov textu T .
- **Modelování jazyka** je technika, která se snaží předpovídat, co bude následující slovo na základě předchozího kontextu. Necht' jsme před slovem w . Označme h dosavadní historii (text před slovem w). Pak nás zajímá $P(w|h)$. Což z Bayesovy věty spočítáme jako $P(w|h) = P(h|w) * P(w) / P(h)$. Díky větě o úplné pravděpodobnosti pak můžeme počítat pravděpodobnost celé věty W jako:

$$P(W) = P(\langle w_i \rangle_{i=1..n}) = P(w_n | \langle w_i \rangle_{i=1..n-1}) * P(w_{n-1} | \langle w_i \rangle_{i=1..n-2}) * \\ P(w_{n-2} | \langle w_i \rangle_{i=1..n-3}) * \dots * P(w_2 | w_1) * P(w_1)$$

- Jelikož příliš dlouhá historie by byla výpočetně náročná a zároveň by mnohé pravděpodobnosti byly příliš malé (kombinace dlouhých sousloví nejsou příliš pravděpodobná), tak se historie omezuje pouze na 3 slova. Což se nazývá trigramový model:

$$P(W) = P(w_3 | w_2 w_1) * P(w_2 | w_1) * P(w_1)$$

- Termín ***n-gram*** znamená n -tice slov za sebou (lépe by však bylo upřesnit „slovní n -gram“, jindy se n -gramem totiž myslí n -tice písmen).
- **Výhlazování**. Ve velkém slovníku je příliš mnoho nulových pravděpodobností (kombinací trigramů je hodně, ale v daných textech se jich vyskytne jen malá část). To se řeší tak, že nulové pravděpodobnosti se nahradí nějakými malými hodnotami. **Proč? Asi aby nám to vůbec dávalo nějaké výsledky. Abychom nedostávali samé nuly...?**

Úvod do počítačové lingvistiky – příprava na zkoušku

- **Noise Channel viz kapitola strojový překlad.**
- Jak měřit kvalitu překladu? To je obtížná záležitost i ručně, natož automaticky. V roce 2002 vznikla míra Bleu. Pro porovnání dvou překladů vyžaduje mít daný text ještě alespoň jednou kvalitně přeložený. Následně zkoumá, zda se dané slovní n-gramy (slova, dvojice, trojice a čtveřice slov) vyskytují v některých z referenčních překladů. Čím více je překladů, tím lépe nám to řeší problém synonym nebo pouze jinak správně uspořádaného slovosledu.
 - $Bleu = BP * (p_1 * p_2 * p_3 * p_4)^{1/4}$. BP je koeficient kvůli krátkým větám **či co**.
 - Problémy jsou zřejmé – jiný slovosled může způsobit velmi špatné výsledky v této metrice. Nebere v úvahu morfolonii, tedy pouze chybná koncovka (ale správný význam) pokazí skóre stejně jako úplně špatný překlad. Je to hodně náročné na velikost trénovacích dat. Proto se lépe překládá mezi „velkými jazyky“, kde se data lehko shání.

Otázky:

1. Podrobně popište **statistické metody v automatickém překladě**.

2. Co to jsou n-gramy? (Pozor na to, že zde se mluví o slovních n-gramech, ne písmenkových.)

3. Co je vyhlazování?

4. Bleu metoda.

SÉMANTIKA

- Sémantika přirozeného jazyka
 - Pomocí syntaxe můžeme rozlišovat gramaticky správné a nesprávné věty. Nicméně nic to neříká o jejich pravdivosti. Zároveň je nutno rozlišovat mezi **významem** a **pravdivostí** věty. (Naopak sémantika formálních jazyků tyto pojmy často ztotožňuje.) Pravdivost je dána kontextem, není obsažena v jazyce. Jsou k ní potřeba různá pravidla a předpoklady světa, ze kterého vycházíme. Navíc i nepravdivé věty mohou mít svůj význam. U některých vět zase není možno ověřit pravdivost. Obdobně je těžké obecně rozlišit věty se stejným významem. (*Pozorovali ho dobrovolně. X Byl jimi pozorován dobrovolně.*)
 - **Výplývání** – z pravdivé věty často vyplývají různé další skutečnosti (na základě obecných pravidel a zákonitostí), nicméně tyto zákonitosti nejsou stoprocentní, mohou mít výjimky, které nás předem nenapadnou. (*Tučnáci jsou ptáci. => Tučnáci mají křídla a létají.*)
 - **Fregeho princip kompozicionality** (1925 Gottlob Frege). Význam složeného výrazu je jednoznačně určen významy jeho částí a způsobem jejich kombinace. Tedy např. význam textu je určen významy jednotlivých vět a jejich poskládáním. Obdobně význam vět je určen významem jejich slov, atd.
- Lexikální sémantika
 - Pro popis významu slov bychom potřebovali opět nějaký (*meta*)jazyk – buď **formální** (např. vycházející z něčeho již vystavěného, např. predikátové logiky) nebo **přirozený** (ten stejný nebo jiný) + se dá využívat okolní svět (*Toto je křída.*). Přirozený jazyk používají například výkladové slovníky, slovníky synonym, definice slov apod.
 - Problémy lexikální sémantiky: Význam slova závisí na kontextu okolních slov a vět (např. *Střílení poslanců ohrožuje naši demokracii*). Význam slov není jednoznačný (*oko, list, ...*).
 - Jednou z možností popisu významu slov jsou **významové třídy (rasy)**. **Ontologie** je množina objektů, která představuje klasifikaci objektů universa U na různé třídy (např. fyzické objekty, vlastnosti, vztahy, činnosti, živé bytosti apod.), které lze dále dělit. Dané slovo (objekt) pak popíšeme pomocí příznaků ke každé třídě: + (patří do ní), - (nepatří do ní), 0 (nezávisí na ní). Ontologie jsou buď **doménové** (někde jsem našla, že zpracovává jen jednu doménu – obor; jinde že to je množina názvů oborů) či **vrcholové** (*Top Ontology* – prý množina nejzákladnějších výrazů, nezávislých na jazyku).
 - Jinou možností popisu významu slov jsou **sémantické sítě**. Ty umožňují určit různé vztahy a směry vztahů mezi pojmy, tj. nejen hierarchii sémantických tříd, ale i vztahy napříč nimi. Zabývají se vztahy jako **hyponimie** (slovo nadřazené) a **hyponimie** (slovo podřazené), **synonymie** (ekvivalentní význam ale jiná forma) a **antonimie** (slova protikladná), **meronymie** (býti částí) a **holonymie** (obsahovat). Navíc se dobře zpracovávají počítačově.
 - **WordNet** 1993 G. A. Miller. Rozsáhlá lexikální databáze anglických slov (podstatná a přídavná jména, slovesa a příslovce) seskupených do množin synonym, tzv. **synonymických řad** neboli **synsetů**. Každý synset vyjadřuje určitý koncept. (**Co tím chtěl autor říct?** Prostě všechna slova v synsetu mají stejný význam.) Mezi synsety jsou propojení v podobě sémantických a lexikálních relací. Síť lze procházet pomocí počítače. Nicméně vznikala hlavně ručně.
 - Z příkladu mi přijde, že to vypadá jako takový výkladový slovník – k danému slovu to vypisuje určité jeho významy. Pro dané významy to vypíše slovní popis, příklad a seznam synonym. Co z toho dělá sémantickou síť je asi to, že to tvoří vnitřní strukturu (ostatní pojmy fungují jako odkazy), kterou lze procházet a jsou tam zaznamenány i různé relace nadřazenosti/ podřazenosti/ býti částí apod.
 - **EuroWordNet** 1997 Vossen. Rozšíření WordNetu do více jazyků (nejdříve přidány holandština, italština a španělština, později francouzština, němčina, čeština a estonština). Navíc byla zavedena vrcholová ontologie, což byla množina 63 nejzákladnějších výrazů (konceptů), nezávislých na jazyku. Ke každému jazyku pak bylo vybráno 1000 základních jazykově závislých konceptů (*Base Concepts*), tvořících jádra slov. V Aj WordNetu získal každý synset jednoznačný identifikátor, díky kterému vznikl mezi-jazykový index (*Inter-Lingual Index, ILI*). Pak byly na sebe různojazyčné WordNety navázány a vznikly vztahy ekvivalence (*EQ-relations, ale to moc nevím, co bylo.*)
 - Aplikace
 - Překlad. Jednak může pomáhat v počítačem asistovaném překladu (Computer Aided Translation). Překladatel si v něm může hledat významy slov, jejich synonyma, antonyma, příklady použití, slova odvozená apod. Druhá společně s morfologickou a syntaktickou analýzou může díky tomu, že ukládá valenční rámce pro slovesa, sloužit k automatickému strojovému překladu.
 - Extrakce informací. Např. může sloužit při vícejazyčném vyhledávání a kdekoliv kde jsou potřeba sémantické vztahy jako synonymie apod.
 - Určování významů slov (Word Sense Disambiguation).
 - Nějak může pomáhat sémantickému webu.
 - Vyhodnocování kvality překladu – zlepšení automatických metrik typu BLEU.

- **Problémy**
 - Především je problém, že cizojazyčné WordNety vznikaly především jako překlad toho anglického. Tedy nezachycují typické vlastnosti jazyků. Také tam vzniklo mnoho chyb a nekonzistencí. A navíc pak přestaly být projekty financovány a přestaly se rozvíjet.
 - Jiný problém je, že podobných výsledků (a lepších, rychlejších, s méně usilím) lze dnes dosahovat pomocí statistických metod. A obecně v době Googlu apod. některé projekty jako je WordNet už nemají tak dobrý smysl.
- **Reprezentace významu vět**
 - Pomocí predikátové logiky 1. řádu + se přidávají nové vlastnosti. Vychází z principu kompozicionality. Složkám věty odpovídají části sémantického zápisu. Problémy začíná tvořit modalita, čas, postoj, předpoklady (*presupozice*), neurčitost (*fuzziness*), apod. Jsou tedy vytvořeny nové operátory (např. *possible(F)*, *necessary(F)*). Ale tento přístup nemá dostatečnou sílu. Nastávají problémy při nahrazování částí vět.
 - **Extenze** je souhrn věcí, které pod pojem spadají. **Intenze** je samotný popis (charakteristika, definice) pojmu. Např. intenzí pojmu „čtyřúhelník“ je rovinný mnohoúhelník se čtyřmi vrcholy a čtyřmi stranami. Extenzí stejného pojmu jsou asi pojmy různoběžník, rovnoběžník = kosodélník (tedy obdélník, kosočtverec, čtverec), lichoběžník, deltoid, atd.
- **Základní přístupy k sémantice**
 - **Modelově-teoretická sémantika**. Pracuje s pravdivostními podmínkami vztaženými k určitému modelu. Syntaktické kategorie odpovídají sémantickým typům. Obsahuje základní lexikální výrazy a jejich interpretaci, syntaktická a sémantická pravidla.
 - **Montagueovská gramatika**. Původně *Universal Grammar* 1971. Založena na formální logice, lambda kalkulu, teorii množin, používá pojmy intenzionální logiky a teorie typů. Vychází z předpokladu, že neexistuje zásadní rozdíl mezi sémantikou přirozených a formálních jazyků. Obsahuje syntaktické kategorie s množinami konkrétních slov a syntaktická pravidla pro slova z těchto kategorií.
 - **Kompozicionální sémantika**. Vychází z principu kompozicionality. Používá různé reprezentace (sémantické rysy a jejich skládání, koncepty a převod ze syntaktické reprezentace, logickou reprezentaci a zjišťování pravdivosti).
 - Další příklad **Transparentní intenzionální logika** (TIL). Je založen na typovém lambda kalkulu. Nemá vlastní logické spojky, kvantifikátory apod. Nějak řeší možné světy. Univerzum je množina společná všem možným světům. Používá nějaké nálepky individuí. Či co.
- **Rozpoznávání vztahů v textu**. Pochopení smyslu textu je ještě těžší než smyslu věty. Problém je, že věty v textu na sebe navazují a odkazují se (např. nevyjádřeným podmíněm). Jsou tyto typy vztahů v textu (resp. spíše referencí):
 - **Exoforma** – odkazování mimo text – zájmeno poukazuje k mimotextové situaci či skutečností. *Slyšíš ji? Dej mi, prosím, tyhle tři.*
 - **Endoforma** – odkazování v rámci textu.
 - **Anafora** – odkazování zpětně v textu.
 - **Zájmena a nevyjádřený podmět či jiný větný člen** (tzv. nulové výrazy). *Jakub šel udělat čaj. To neměl dělat. On se u toho vždy opaří.*
 - **Určité jmenné skupiny**. Dvojí sousloví označující to samé. *Budeme první, prohlásil Pavel. Nedvěd. To mělo mužstvo rádo, když jejich kapitán takto promluvil.*
 - **Elipsa**. Vynechání části věty obsahující informaci, která je příjemci známa a bez níž větu dokáže pochopit. *Jak se máš? Dobře. Kolik je? Petr přinesl dva stoly. Dřevěný a kovový.*
 - **Textové spojovací výrazy**: spojky, například, na jedné straně – na druhé straně, atd. *Nejdříve si dáme jídlo, pak siestu.*
 - **Katafora** – odkazování dopředu v textu. *Věřte tomu nebo ne, máme schodkový rozpočet. Dej mu pohlavěk, tomu uličníkovi.*
 - **Aplikace**.
 - Získávání informací z textu.
 - Automatický překlad.
 - Dialogové systémy.
 - **Řešení**.
 - Co můžeme využít? Morfologické značky (shoda v rodě a pádě), syntaktickou strukturu věty (valence pomáhá doplnění elipsy), statistické přístupy, aktuální členění (odkazujeme se jinak na něco, co bylo zmíněno na začátku a uprostřed (resp. v základu a ohnisku)) a pomocné znalosti (ontologie, sémantická síť, atd.).
 - Např. *Stock of Shared Knowledge*. Každému podstatnému jménu přidělí určitý index podle „důležitosti“, resp. pravděpodobnosti, že se na něj pak někdo bude odkazovat.

Úvod do počítačové lingvistiky – příprava na zkoušku

Otázky:

1. Rozdíly mezi významem a pravdivostí věty.
2. Fregova koncepce (Fregeho princip kompozicionality).
3. Ontologie - co to je a jak se používá.
4. EuroWordNet, WordNet.
5. Rozdíl extenze / intenze v sémantice.
6. Rozdíly mezi modelově teoretickou a kompozicionální sémantikou.
7. Uveďte 4 typy anaforických vztahů v textu + příklady.