

T-(538|725)-MALV, Málvinnsla Orðflokkar og orðhlutafræði

Hrafn Loftsson¹ Hannes Högni Vilhjálmsson¹

¹Tölvunarfræðideild, Háskólinn í Reykjavík

September 2007

- 1 Orðflokkar
- 2 Orðasafn
- 3 Orðhlutafræði
- 4 Tveggja laga orðhlutagreining

- 1 Orðflokkar
- 2 Orðasafn
- 3 Orðhlutafræði
- 4 Tveggja laga orðhlutagreining

Skilgreining

- Orðflokkur er flokkur orða sem hafa sameiginleg málfræðileg einkenni:
 - Merkingarleg einkenni.
 - Beygingarleg einkenni.
 - Setningarleg einkenni.
 - <http://is.wikibooks.org/wiki/Or%C3%B0flokkagreining>
- e. Part-of-speech, e. lexical categories.

Yfirflokkun

- Lokaður flokkur (e. closed class).
- Opinn flokkur (e. open class).

Skilgreining

- Orðflokkur er flokkur orða sem hafa sameiginleg málfræðileg einkenni:
 - Merkingarleg einkenni.
 - Beygingarleg einkenni.
 - Setningarleg einkenni.
 - `http://is.wikibooks.org/wiki/Or%C3%B0flokkagreining`
- e. Part-of-speech, e. lexical categories.

Yfirflokkun

- Lokaður flokkur (e. closed class).
- Opinn flokkur (e. open class).

Lokaður flokkur

- Ný orð bætast ekki við þennan flokk.
- Kerfisorð.
- Forsetningar, samtengingar, hjálparsagnir, samtengingar, atviksorð (nema háttaratviksorð).

Opinn flokkur

- Getur bætt við sig nýjum orðum.
- Orð geta einnig tapast!
- Inntaksorð: gegna ákveðnu merkingarhlutverki.
- Nafnorð, lýsingarorð, sagnir (nema hjálparsagnir), háttaratviksorð.

Lokaður flokkur

- Ný orð bætast ekki við þennan flokk.
- Kerfisorð.
- Forsetningar, samtengingar, hjálparsagnir, samtengingar, atviksorð (nema háttaratviksorð).

Opinn flokkur

- Getur bætt við sig nýjum orðum.
- Orð geta einnig tapast!
- Inntaksorð: gegna ákveðnu merkingarhlutverki.
- Nafnorð, lýsingarorð, sagnir (nema hjálparsagnir), háttaratviksorð.

11 orðflokkar í íslensku

Fallorð	Sagnorð	Smáorð
Nafnorð (e. nouns)	Sagnorð (e. verbs)	Forsetningar (e. prepositions)
Lýsingarorð (e. adjectives)		Atviksorð (e. adverbs)
Fornöfn (e. pronouns)		Samtengingar (e. conjunctions)
Töluorð (e. numerals)		Upphrópanir (e. interjections)
Greinir (e. article)		Nafnháttarmerki (e. infinitive marker)

Fallorð fallbeygjast. Sagnorð tíðbeygjast. Smáorð beygjast ekki.

Beygingarleg einkenni (e. features) (feitletrun vísar í skammstafanir í íslenska markamenginu)

- Kyn (e. gender)
 - **karlkyn** (e. masculine)
 - **kvenkyn** (e. feminine)
 - **hverugkyn** (e. neuter)
- Tala (e. number)
 - **eintala** (e. singular)
 - **fleirtala** (e. plural)
- Fall (e. case)
 - **nefnifall** (e. nominative)
 - **þolfall** (e. accusative)
 - **þágufall** (e. dative)
 - **eignarfall** (e. genitive)

Beygingarleg einkenni (e. features)

- Háttur (e. mood) http://is.wikipedia.org/wiki/H%C3%A6ttir_sagna
 - framsöguháttur (e. indicative mood)
 - viðtengingarháttur (e. subjunctive mood)
 - boðháttur (e. imperative mood)
 - nafnháttur (e. infinitive mood)
 - lýsingarháttur nútíðar (e. present participle)
 - lýsingarháttur þátíðar (e. past participle)
- Mynd (e. voice) <http://is.wikipedia.org/wiki/Sagnmyndir>
 - gemynd (e. active)
 - miðmynd (e. middle)
 - þolmynd (e. passive) - mynduð með lýsingarhætti þátíðar
- Persóna (e. person): 1., 2., 3.
- Tala (e. number): eintala, fleirtala
- Tíð (e. tense)
 - nútíð (e. present)
 - þátíð (e. past)

Outline

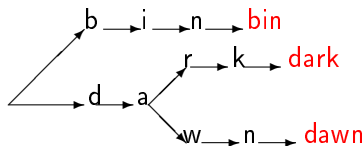
- 1 Orðflokkar
- 2 Orðasafn
- 3 Orðhlutafræði
- 4 Tveggja laga orðhlutagreining

Geymir hvað?

- Listi af (lesmáls)orðum með eða án frekar upplýsinga.
 - Beygingarupplýsingar.
 - Mörk.
 - Setningafræðilegar upplýsingar.
- Oft búið til fyrir tiltekið svið, t.d. tækni, vísindi, fjármál.
- Yfirleitt það fyrsta sem búa þarf til fyrir máltækni kerfi.

Gagnaskipan

- Tætitafla (lykill–gildi): sérhvert orð er lykill og frekari upplýsingar um orðið er gildi þess.
 - Getur tekið mikið minni fyrir stór orðasöfn.
- *Letter trees* (tries) <http://en.wikipedia.org/wiki/Trie>
 - Orð eru geymd sem stafatré (e. tree of characters) sem deila kvíslum eins lengi og einstakir stafir orða eru þeir sömu.



Outline

- 1 Orðflokkar
- 2 Orðasafn
- 3 Orðhlutafræði**
- 4 Tveggja laga orðhlutagreining

Myndan (e. morpheme)

- Myndan (morfem) er minnsta afmörkuð eind í byggingu orðsins.
- Afmarkaður hluti úr orði sem hefur einhverja ákveðna merkingu eða hlutverk.
- Myndön skiptast í stofn (e. stem) og aðskeyti (e. affix) eða beygingarendingu.
- Stofn er sú mynd orðs sem helst eins í öllum beygingarmyndum þess.
- Dæmi: Orðið “orði” samanstendur af tveimur myndönum:
 - “orð” sem er stofn.
 - “i” sem er beygingarending.

Aðskeyti

- Forskeyti (e. prefix)
 - Fer á undan stofni.
 - Dæmi: örfínn
- Viðskeyti (e. suffix)
 - Kemur á eftir stofni.
 - Dæmi: kennari
- Innskeyti (e. infix)
 - Er sett inn í stofn.
 - Dæmi: humingi í tagalog; hingi: að lána, humingi: sá sem fær lánað
- Umskeyti (e. circumfix)
 - Bæði fer á undan og kemur á eftir stofni.
 - Dæmi: gesagt í þýsku (lhpt. af sagen).

- e. morphology
[http://en.wikipedia.org/wiki/Morphology_\(linguistics\)](http://en.wikipedia.org/wiki/Morphology_(linguistics))
- Orðhlutafræði – Orðmyndunarfræði
- Fjallar um hvernig orð eru sett saman úr myndönnum.
- Tvenns konar orðhlutakerfi:
 - Concatenative morphology
 - T.d. Germönsk mál
 - http://en.wikipedia.org/wiki/Germanic_languages
 - Forskeyti Stofn Viðskeyti
 - Non-concatenative (templatic) morphology
http://en.wikipedia.org/wiki/Nonconcatenative_morphology
 - T.d. Arabíska, Hebreska

Orðmyndun í íslensku

Orð eru mynduð á tvennan hátt úr myndönum:

Með beygingu (e. inflection)

- Forskeyti + Stofn + beygingarending
- Nýja orðið tilheyrir sama orðflokki og stofninn.
- Dæmi: hestur, málvinnsla

Með afleiðslu (e. derivation)

- Forskeyti + Stofn + viðskeyti
- Nýja orðið getur tilheyrt öðrum orðflokki en stofninn.
- Dæmi: óljós, klofningur, kennari

Orðmyndun í íslensku

Orð eru mynduð á tvennan hátt úr myndönum:

Með beygingu (e. inflection)

- Forskeyti + Stofn + beygingarending
- Nýja orðið tilheyrir sama orðflokki og stofninn.
- Dæmi: hestur, málvinnsla

Með afleiðslu (e. derivation)

- Forskeyti + Stofn + viðskeyti
- Nýja orðið getur tilheyrt öðrum orðflokki en stofninn.
- Dæmi: óljós, klofningur, kennari

Samsett orð (e. compound words)

- Tvö (eða fleiri) orð sett saman til að mynda nýtt orð.
- Algengt að setja saman tvö nafnorð.
- http://rettritun.is/?id=namsefni_k9_r1

Orðhlutafræðileg greining/þáttun

- e. morphological analysis/parsing
- Sú aðgerð að brjóta orð upp í einstök myndön.
- Nauðsynleg í mörkum máltækniakerfum:
 - “Lemmatization”. Það að finna uppflettimynd (lemmu, flettu) orðsins (e. dictionary form).
 - “Stemming”. Það að finna stofn, notað t.d. í upplýsingaheimt.
 - Orð sem ekki finnast í orðasafni þurfa greiningu (t.d. í mörkun).

Orðhlutafræðileg myndun

- Að mynda orð að gefnum myndönnum.

Orðhlutafræðileg greining/þáttun

- e. morphological analysis/parsing
- Sú aðgerð að brjóta orð upp í einstök myndön.
- Nauðsynleg í mörkum máltækniakerfum:
 - “Lemmatization”. Það að finna uppflettimynd (lemmu, flettu) orðsins (e. dictionary form).
 - “Stemming”. Það að finna stofn, notað t.d. í upplýsingaheimt.
 - Orð sem ekki finnast í orðasafni þurfa greiningu (t.d. í mörkun).

Orðhlutafræðileg myndun

- Að mynda orð að gefnum myndönnum.

Lemmatization vs. Stemming

- Dæmi: “hesti”. Lemma = “hestur”. Stofn = “hest”

Margræðni í lemmun

- 1 A **run** in the forest. Lemma: **run**, nafnorð í eintölu.
- 2 The sportsmen **ran** everyday. Lemma: **run**, sögn í þátíð, þriðju persónu, fleirtölu.
- 3 **Rétturinn** var fullskipaður. Lemma: **réttur**, nafnorð í karlkyni, eintölu, nefnifalli, með greini.
- 4 Dæmið var **rétt**. Lemma: **réttur**, lýsingarorð í hvorugkyni, eintölu, nefnifalli, sterk beyging, frumstig.

Lemmatization vs. Stemming

- Dæmi: “hesti”. Lemma = “hestur”. Stofn = “hest”

Margræðni í lemmun

- 1 A **run** in the forest. Lemma: **run**, nafnorð í eintölu.
- 2 The sportsmen **ran** everyday. Lemma: **run**, sögn í þátíð, þriðju persónu, fleirtölu.
- 3 **Rétturinn** var fullskipaður. Lemma: **réttur**, nafnorð í karlkyni, eintölu, nefnifalli, með greini.
- 4 Dæmið var **rétt**. Lemma: **réttur**, lýsingarorð í hvorugkyni, eintölu, nefnifalli, sterk beyging, frumstig.

Outline

- 1 Orðflokkar
- 2 Orðasafn
- 3 Orðhlutafræði
- 4 Tveggja laga orðhlutagreining**

Hvað á að vera í orðasafninu?

- Því ekki hafa allar orðmyndir í orðasafni?
 - T.d. fall- og töluendingar?
- Mörg ferli eru fyrirsegjanleg (e. predictable)
 - Þess vegna er óhagkvæmt að telja upp myndir.
- Mörg ferli eru frjó (e. productive)
 - Þess vegna útilokað að telja upp myndir.
- Þess vegna er betra að hafa aðskilin söfn og reglur sem raða orðhlutum saman.

Glæra fengin að láni úr námskeiðinu: Inngangur að tungutækni,
Eiríkur Rögnvaldsson, HÍ, 2002.

Tveggja laga orðhlutagreining

Tilgangur

- e. two-level morphology (Kimmo Koskeniemi 1983).
- Tengir saman “surface form” orðs – orðið eins og það kemur fyrir í texta – við “lexical form” (underlying form) – þ.e. röð myndana orðsins.
- Vörpunin á milli “surface form” og “lexical form” á sér stað með *stöðuferjali* í báðar áttir.

Dæmi (0 merkir tómi strengurinn)

Lexical form: hest+ur (myndan notað)

Surface form: hest0ur

Lexical form: hest +n +k +e +n (beygingarleg einkenni notuð)

Surface form: hest 0 0 u r



Tveggja laga orðhlutagreining

Tilgangur

- e. two-level morphology (Kimmo Koskeniemi 1983).
- Tengir saman “surface form” orðs – orðið eins og það kemur fyrir í texta – við “lexical form” (underlying form) – þ.e. röð myndana orðsins.
- Vörpunin á milli “surface form” og “lexical form” á sér stað með *stöðuferjali* í báðar áttir.

Dæmi (0 merkir tómi strengurinn)

Lexical form: hest+ur (myndan notað)

Surface form: hest0ur

Lexical form: hest +n +k +e +n (beygingarleg einkenni notuð)

Surface form: hest 0 0 u r



Stöðuferjald

- Endanleg stöðuvél sem ber kennsl á eða myndar **par af strengjum**.
- Leggir merktir með tveimur stöfum, sá fyrsti er inntaksstafur, sá seinni er úttaksstafur.
- Vélin breytir (e. transduces) inntaksstaf í úttaksstaf þegar færsla á sér stað eftir legg.
- Sjá Fig. 5.6. bls. 132 í kennslubók.

Stöðuferjald (e. finite-state transducer)

Stærðfræðileg skilgreining

Stöðuferjald samanstendur af fimm hlutum $(Q, \Sigma, q_0, F, \delta)$:

- 1 Q er mengi af endanlegum stöðum, $q_0, q_1 \dots q_n$.
- 2 Σ er endanlegt mengi af inntakspörum $i : o$, i er úr inntaksstafrófi, o úr úttaksstafrófi.
- 3 q_0 er upphafsstaða, $q_0 \in Q$.
- 4 F er mengi lokastaða, $F \subseteq Q$.
- 5 δ er breytingafall (e. transition function) $Q \times \Sigma \rightarrow Q$, $\delta(q, i, o)$ skilar þeirri stöðu sem vélin fer í úr stöðu q miðað við inntaksparið $i : o$.

Stöðuferjald: Dæmi

hestur

Eintala

Lexical: hest +n +k +e +n

Surface: hest 0 0 u r

Lexical: hest +n +k +e +o

Surface: hest 0 0 0 0

Lexical: hest +n +k +e +þ

Surface: hest 0 0 0 i

Lexical: hest +n +k +e +e

Surface: hest 0 0 0 s

Fleirtala

hest +n +k +f +n

hest 0 0 a r

hest +n +k +f +o

hest 0 0 0 a

hest +n +k +f +þ

hest 0 0 u m

hest +n +k +f +e

hest 0 0 0 a

- Prolog: Sjá Appendix A í kennslubók.
- SWI-Prolog: <http://www.swi-prolog.org/>
- Útfærslan á stöðuferjaldinu (Prolog-kóði): Sjá undir “Annað efni”.
- Lesa forritskóða inn í Prolog: `consult('transduce.pro')`.
- Prófa ferjaldið: `transduce(1, Final, Lexical, [h,e,s,t,u,r])`.
- Prófa ferjaldið: `transduce(1, Final, [h,e,s,t,+n,+k,+e,+n], Surface)`.