

# Syntaktisk informert frasesamanstilling

Kevin Brubeck Unhammer

09/08, 2010

# Innhald

<b>1</b>	<b>Innleiing</b>	<b>5</b>
1.1	Vegkart . . . . .	6
1.2	<b>SKRIV</b> Frasesamanstilling frå f-struktur <b>:ROTETE:</b> . . . . .	7
<b>2</b>	<b>Bakgrunn og omgrepsavklaring</b>	<b>9</b>
2.1	<b>SKRIV</b> Relaterte metodar . . . . .	9
2.2	<b>SKRIV</b> Eit kort oversyn over leksikalsk-funksjonell grammatikk og terminologi . . . . .	10
<b>3</b>	<b>Krav til frasesamanstilling</b>	<b>12</b>
3.1	Innleiing . . . . .	12
3.2	Formål med frasesamanstilling . . . . .	12
3.3	Frasesamanstilling i ein LFG-trebank . . . . .	14
3.4	Kva kan lenkjast? . . . . .	15
3.4.1	<b>TOGROK</b> finst det tilfelle der ordlenkjer ikkje impliserer PRED-lenkjer? . . . . .	17
3.5	<b>SKRIV</b> Krav på ordnivå . . . . .	17
3.5.1	<b>SKRIV</b> LPT <b>:ROTETE:</b> . . . . .	17
3.5.2	Lik ordklasse? . . . . .	18
3.6	Funksjonsord . . . . .	18
3.6.1	<b>TOGROK</b> cvimda<PRO> men regne<>expletive – lenkje? <b>:ROTETE:</b> . . . . .	18
3.6.2	<b>TOGROK</b> kva med ekspletivar? ingen PRED men heller ikkje C/F/I <b>:ROTETE:</b> . . . . .	18
3.7	Lenkjing av underordna c-strukturknodar . . . . .	19
3.7.1	<b>SKRIV</b> døme! <b>:ROTETE:</b> . . . . .	21
3.7.2	<b>TOGROK</b> me <sub>OBJ</sub> gusta X <sub>SUBJ</sub> // I <sub>SUBJ</sub> like X <sub>OBJ</sub> ?? <b>:ROTETE:</b> . . . . .	21
3.7.3	<b>TOGROK</b> korleis finn me <i>there is</i> -lenkjer då? <b>:ROTETE:</b> . . . . .	21
3.8	<b>TOGROK</b> mange-til-mange-lenkjing i f-strukturane? <b>:ROTETE:</b> . . . . .	21
3.8.1	<b>SKRIV</b> Kva inneber ei mange-til-mange-lenkjing? <b>:ROTETE:</b> . . . . .	21
3.9	<b>SKRIV</b> Mangel på samsvar i syntaks og semantikk <b>:ROTETE:</b> . . . . .	21

3.10	<b>TOGROK</b> Diskontinuerlege einingar <b>:ROTETE:</b> . . . . .	22
3.10.1	<b>TODO</b> døme på diskontinuerlege konstituentar som er lenkja <b>:ROTETE:</b> . . . . .	22
3.11	<b>TOGROK</b> Er «compounds» frasar? <b>:ROTETE:</b> . . . . .	22
3.12	Krav om lik argumentstruktur . . . . .	22
3.12.1	forsvare «tilsvarande» <b>:ROTETE:</b> . . . . .	24
3.12.2	<b>TODO</b> Sitere eigen korpusundersøking av variasjon i arg-str? . . . . .	24
3.12.3	<b>SKRIV</b> kvifor lik arg-str er bra, så kvifor det er eit problem <b>:ROTETE:</b> . . . . .	24
3.12.4	<b>TODO</b> Ulik følge i argumentstruktur . . . . .	24
	c- og f-strukturar for dømet over <b>:ROTETE:</b> . . . . .	25
3.12.5	<b>SKRIV</b> døme med wager/3 og vedde/4 og gewettet/3 <b>:ROTETE:</b> . . . . .	26
3.12.6	<b>SKRIV</b> (reinskriv) <b>:ROTETE:</b> . . . . .	26
3.12.7	<b>SKRIV</b> True Arguments vs True Adjuncts, Pustejovsky <b>:ROTETE:</b> . . . . .	27
3.13	<b>SKRIV</b> Kan adjunkt lenkjast til nodar <u>undermor</u> -lenkja? . . . . .	27
3.13.1	1. Kausativar og inkorporering . . . . .	28
	<b>TOGROK</b> adjunkt bør ikkje samanføyast? eller? . . . . .	29
3.13.2	2. Adposisjonsobjekt . . . . .	29
3.14	<b>TOGROK</b> kva var poenget med dette? <b>:ROTETE:</b> . . . . .	29
3.15	<b>ULEST</b> Cyrus, FuSe-prosjektet <b>:ROTETE:</b> . . . . .	30
3.16	<b>TODO</b> Konstruksjonar og komposisjonell inekvivalens . . . . .	30
3.17	<b>SKRIV</b> definer sitering frå MRS-suiten <b>:ROTETE:</b> . . . . .	30
3.18	<b>SKRIV</b> setning 7 i MRS-suiten <b>:ROTETE:</b> . . . . .	30
3.19	<b>TOGROK</b> og så finst jo større forskjellar, stilistiske osv... <b>:ROTETE:</b> . . . . .	31
3.20	<b>TOGROK</b> prosessering, kognitive modellar? <b>:ROTETE:</b> . . . . .	31
3.21	<b>TOGROK</b> Retningslinjer for samanstilling <b>:ROTETE:</b> . . . . .	32
3.22	<b>SKRIV</b> Rangering . . . . .	32
<b>4</b>	<b>Implementasjonen av lfgalign</b> . . . . .	<b>33</b>
4.1	Lenkjer mellom f-strukturar . . . . .	34
4.1.1	Overflødige adverbial . . . . .	37
4.1.2	<b>SKRIV</b> Når f-lenkjene ikkje er 1-1 . . . . .	37
	kausativ . . . . .	37
	preposisjonsobjekt . . . . .	37
	notat <b>:ROTETE:</b> . . . . .	38
4.1.3	Kan me gjere f-struktursamanstillinga bottom-up? . . . . .	38
4.2	<b>SKRIV</b> Rangering . . . . .	39
4.2.1	rekursivt lenkja > ulenkja, men LPT . . . . .	39
4.2.2	argument-argument > argument-adjunkt . . . . .	39
4.2.3	arg1-arg1 arg2-arg2 > arg1-arg2 arg2-arg1 (følge) . . . . .	39
4.2.4	flest lenkja adjunkt . . . . .	39

<i>INNHALD</i>	3
4.2.5 Prioritet på rangeringskriterium . . . . .	39
4.3 Lenkjing av c-strukturnodar . . . . .	40
4.3.1 <b>SKRIV</b> viss me har LPT, men ikkje rekursiv f-lenkje . . .	41
<b>5 Diskusjon, resultat av å automatisk samanstille norske og georgiske setningar</b>	<b>42</b>
5.1 Oppdage argumentstrukturnodasjøn . . . . .	42
5.2 Samanlikning med tremetodar og n-grammetodar . . . . .	42
5.2.1 c->f er mange-til-ein . . . . .	43
5.2.2 <b>TOGROK</b> men korleis gjenfinne there is/es gibr? <b>:ROTE-TE:</b> . . . . .	43
<b>6 Avslutning</b>	<b>44</b>

# List of Corrections

TODO: abstract/samandrag . . . . .	5
fortelje om georgisk i kap.5 heller . . . . .	7
limt inn frå prosjektskildringa, dette må omskrivast totalt/fjernast. Kan kanskje ta eit par setningar inn i kap.3, men innleiinga bør iallfall ha eit kort oversyn over metoden . . . . .	7
TODO (innleiing?):	
- reine N-gram-samanstillingar, dependensbaserte	
- ulike formål for samanstilling gir ulike metodar	
- kort oversyn over LFG-termar	
. . . . .	9
«by på fleire problem» – weasel wording, TODO omskriv. . . . .	9
fleire slike? meir om dette, algoritmen . . . . .	9
Dette er motsett retning av det mitt program gjer, nemne seinare? . . . .	10
todo: referere til den faktiske parsaren? det var Bikel kanskje? . . . . .	13
backe det med eksemplar i trebank; kople til adj-arg-lenkje . . . . .	16
der ADJUNKT ikkje er realisert, lenkjer me ikkje PRED. skal me då ikkje lenkje ord heller? . . . . .	17
PRED->ord :: iallfall	
PRED<-ord :: ?	
PRED<->ord	
PRED, ord . . . . .	17
avsnittet over er litt rotete TODO . . . . .	17
LCS, dorr . . . . .	24
intro TODO, kanskje noko om kva eg faktisk har fått ut av implementasjonen . . . . .	33
treng eg ein eigen del om LPT i dette kapitlet? Implementasjonen er jo veldig enkel iallfall. . . . .	33
TODO: nemne føresetnaden om uavhengnad i kapittel 3 . . . . .	35
TODO: forskjellen mellom LPT-krav og rekursjonskrav på argument må inn i kapittel 3 . . . . .	35
TODO: implementere :-> . . . . .	37
og så er det spørsmålet om me kan lenkje adjunkt på ulike nivå i f-strukturane . . . . .	37
Dette må 1. spesifiserast (kap.3), og 2. implementerast... . . . .	37

Rangering er ikkje implementert enno (akkurat no gir rank(f-alignments) berre ut første samanstilling.) . . . . .	39
To problem (kva vil me ha med?)	
1. me får <i>ikkje</i> med LPT-korrespondansar som er OK, men ikkje med i <i>f – alignment</i> ;	
2. me får med LPT-korrespondansar som er med i <i>f – alignment</i> men ikkje <i>aligntable</i> (ikkje er rekursivt lenkja). . . . .	40
til diskusjonsdel: <i>Det er ikkje berre ei N-gramsamanstilling; sidan len- kjene er mellom c-strukturnodar kor kvar node dominerer ein kon- stituent, kunne me kalt det ei konstituentsamanstilling.</i> . . . . .	41
Kan me <b>fjerne</b> visse f-samanstillingar pga. c-strukturinfo? dvs. disambi- guere... (dette er vel heller stoff for diskusjonsdelen?) . . . . .	41
<b>ref</b> . . . . .	42

# Kapittel 1

## Innleiing

Denne masteroppgåva utforskar kva det vil seie at to uttrykk er omsetjingar av kvarandre, og korleis me automatisk kan generere og evaluere samanstilling (*alignment*) av uttrykk som står i eit slikt omsetjingsforhold.

TODO: abstract/samandrag

Omsetjingsforhold finn me mellom setningar i kontekst på ulike språk, men me kan au finne ulike typar ekvivalensforhold (samanstillingar) mellom frasar innanfor setningane, og mellom lingvistiske skildringar av setningane. I samanheng med XPar-prosjektet (XPar, 2008) har eg sett på metodar for automatisk frasesamanstilling – å finne omsetjingsforhold mellom grupper av fleire ord. Resultatet blir ein *annotasjon*, endå ei lingvistisk skildring av tekstene.

Det at me kan omsetje mellom lingvistiske skildringar (t.d. trekkstrukturane til dei grammatiske rammeverka HPSG eller LFG) gjer det tydeleg at me arbeider med ein *modell* av språket; ulike skildringar kan vere sanne innanfor modellen, utan at modellen er lik språket. Sjølve omsetjingsforholdet er au ein teoretisk storleik, og me kan leggje ulike kriterium til grunn for å kalle to uttrykk omsetjingar av kvarandre.

Kriteria avheng av formålet. Samanstillingsannotasjon kan t.d. nyttast som grunnlag for statistisk eller eksempelbasert maskinomsetjing, i tillegg til oppbygging av parallelle korpora for meir teoretiske språkstudie. For statistisk maskinomsetjing vil alle uttrykk vere omsetjingar av kvarandre med eit visst sannsyn (kanskje null), ein har vanlegvis ikkje kriterium som krev lingvistisk analyse. Når samanstillinga skal nyttast i parallelle korpora for lingvistiske undersøkingar vil ein kanskje ha krav om at uttrykk som skal lenkjast er «like» på eit eller anna mål, utover at dei har opptredt saman ofte; i den manuelle samanstillinga i Samuelsson & Volk (2006) har dei t.d. ein del reint semantiske kriterium for å opprette fraselenkjer i ein parallell trebank, men dei har ikkje krav om syntaktisk likskap.

Xpar-prosjektet, som denne masteroppgåva er ein del av, har mellom anna som mål å oppdage relasjonar mellom grammatiske funksjonar, tematiske roller og kasusmarkering, ved hjelp av parallelle trebankar annotert med djupe grammatiske analysar. Samanstillinga planlagt der krev i større grad syntaktisk likskap for å kalle to uttrykk omsetjingar, men kan difor au tene på det at dei grammatiske ana-

lysane er utvikla med tanke på å vere så parallelle som mogleg. Dei grammatiske analysane er gjort i leksikalsk-funksjonell grammatikk, LFG (Bresnan, 2001). Ei grammatisk analyse i LFG involverer både konstituentstruktur (c-struktur) og funksjonell struktur (f-struktur). Konstituentstrukturen liknar på frasestrukturtrea frå andre grammatiske tradisjonar. Dei funksjonelle strukturane er trekkstrukturar, som mellom anna representerer avhengnadsforhold mellom syntaktiske funksjonar som predikat, subjekt og objekt, i tillegg til å halde informasjon om grammatiske trekk som genus, tal eller kasus. Nodar i c-strukturen kan spesifisere informasjon på ulike stader i f-strukturen<sup>1</sup>.

I XPar-prosjektet vil ein finne ut om metodar for frasesamanstilling kan tene på det at LFG-grammatikkane for dei ulike språka er skrivne med same prinsipp lagt til grunn; to parallellstilte setningar bør ha f-strukturar som er like nok til at me kan samanstille frasar ved hjelp av likskapen mellom f-strukturane. I Dyvik et al. (2009, s. 72) finn me følgjande hypotese:

On the basis of monolingual treebanks constructed from a parallel corpus by means of parallel grammars it will be possible to achieve automatic word and phrase alignment with significantly higher precision and recall than hitherto achieved through other means.

kor «parallel grammars» her krev parallellisme i både f-struktur og c-struktur.

Men i tillegg til at ein kanskje kan få betre skåre på desse kvantitative måla, vil lenkjer mellom f-strukturar gi informasjon som er kvalitativt forskjellig frå det ein kan få med å berre sjå på lenkjer mellom ord, N-gram eller konstituentar.

I denne masteroppgåva spesifiserer eg kva for lenkjer mellom f-strukturar og c-strukturknodar me ønskjer, implementerer eit program `lfgalign` som automatisk finn samanstillingar med slike lenkjer, evaluerer resultatet av å køyre programmet mitt, og samanliknar dette med kva me kan få frå andre metodar.

Programmet `lfgalign` opprettar frasesamanstillingar med hjelp av f-strukturinformasjonen gitt av dei parallelle grammatikkane, og bottom-up-informasjon om kva for ordsamanstillingar som er moglege. F-strukturane avgrensar igjen kva for ordsamanstillingar som er moglege, og kva for c-strukturknodar (syntaktiske frasar) som kan lenkjast.

## 1.1 Vegkart

I neste kapittel ser eg på andre metodar for frasesamanstilling.

I kapittel 3 går eg gjennom kva me ønskjer av ei frasesamanstilling når formålet m.a. er å oppdage relasjonane mellom syntaktiske funksjonar, kasusmarkering og tematiske roller med hjelp av ein parallell trebank. Dette ender opp i ei liste med «krav» som samanstillingane må fylle for å vere lovlege, og som implementasjonen av den automatiske frasesamanstillinga (kapittel 4) må følge.

<sup>1</sup>Ved c-struktur-f-strukturavbildinga  $\phi$ , ein funksjon som tek ein c-strukturnode og returnerer ein (delvis) f-struktur.



Eg evaluerer samanstillingane som kjem ut av denne metoden i kapittel 5, og samanliknar dei med det som er mogleg der me berre har konstituentstruktur (syntaktiske tre) i tillegg til ordsamanstilling.

Eg nyttar språka georgisk og norsk i evalueringa, hovudsakleg fordi dei er svært ulike syntaktisk og morfologisk; Georgisk har t.d. mykje friare ordfølgje og rikare morfologi.

Sidan eg ikkje har tilgang på ferdig setningssamanstilt georgisk-norsk parallelltekst, blir det vanskeleg å køyre den statistiske ordsamanstillinga som er vanleg som første steg i N-grambaserte metodar (utan ein god del forarbeid). Difor konsentrerer eg meg i evalueringa om eit testkorpus kor eg manuelt gjer ordsamanstillinga. Eg veit heller ikkje enno om nokon statistisk parsar av høg kvalitet for georgisk, men testkorpuset er ferdig parsa med LFG-parsaren frå Meurer (2008), c-strukturnodane avgrensar då kva som er ein syntaktisk konstituent.

fortelje om  
georgisk i  
kap.5 heller

## 1.2 SKRIV Frasesamanstilling frå f-struktur :ROTETE:

Om me har f-strukturane til to omsette setningar, burde det kanskje vere mogleg å finne ei f-struktursamanstilling først og så finne ordsamanstillinga ut frå denne. Tanken er at me frå to f-strukturar som skildrar omsette setningar, kan

1. lage ei samanstilling mellom relevante deler av f-strukturane,
2. nytte denne funksjonelle samanstillinga til å finne ei frasesamanstilling, ved å følgje avbildinga frå f-struktur til c-struktur ( $\phi^{-1}$ ).

Eitt problem som byr seg er hypotesemangfaldet: kva for «deler av f-strukturane»? Korleis kan me avgrense søkjerommet? I det minste må me kunne kople det opp mot c-strukturnodar; så PRED-element bør i det minste ha lenkjer, medan t.d. tempus og aspektuell informasjon kanskje er mindre viktig. Men kva kan ignoreras? Vil det oppstå tilfelle då me bør vekte visse element? (Dvs., må me nokon gong disambiguere med slike andre element?)

Vidare må me vite *korleis* me samanstiller desse delene. Me kan t.d. byrje med å kople ytterste PRED frå kvart språk, og så rekursivt kople PRED i dei relevante substrukturane<sup>2</sup>. Gitt ein funksjon  $i$  som returnerer indeksen til ein f-(sub)struktur, kan eit førsteutkast til ei *f-samanstilling*, samanstilling på f-strukturnivå, sjå slik ut:

$$falign(f_1, f_2) = \{(i(f_1(PRED)), i(f_2(PRED)))\} \cup \bigcup_{g_1, g_2 \in fpairs(f_1, f_2)} falign(g_1, g_2)$$

$falign$  vil gi ei mengd av par av indeksar, kor kvart par altså er samanstilt. Ein føresetnad her er at me i tillegg veit kva for par av substrukturar som er «relevante» ( $fpairs(f_1, f_2)$ ).

<sup>2</sup>Dette krev sjølvsagt at ytre PRED faktisk korresponderer i samanstilte setningar, ein ikkje-triviell påstand.

limt inn frå  
prosjektskild-  
ringa, dette  
må omskrivast  
totalt/fjernast.  
Kan kanskje  
ta eit par  
setningar inn i  
kap.3, men  
innleiinga bør  
iallfall ha eit  
kort oversyn  
over metoden

Sjølv om f-strukturar abstraherer frå skilnadene i korleis ulike språk nyttar ordgruppering og ordform til å kode syntaktiske forhold (Bresnan, 2001, s. 14), vil det likevel oppstå forskjellar i f-strukturane til to parallelstilte setningar i eit korpus; både pga. «omsetjarfridom» og det at ulike språk nyttar ulike syntaktiske funksjonar til å uttrykkje det same konseptet. I f-struktursamanstillinga til Riezler & Maxwell (2006, s. 40) får dei t.d. ei lenkje frå ein XCOMP på tysk til eit OBJ på engelsk. Skal ein algoritme gå frå f-strukturar til frasesamanstilling må han i det minste vere robust nok til å takle slik mangel på samsvar. Til å byrje med kan me tenkje oss at *fpairs* gir alle par av GF-ar som har same plass i argumentstrukturen<sup>3</sup> til predikatet, så viss 'sein(SUBJ,XCOMP)' står i  $f_1$  og 'have(SUBJ,OBJ)' i  $f_2$ , vil *fpairs* i det minste returnere  $\{(f_1(\text{SUBJ}), f_2(\text{SUBJ})), (f_1(\text{XCOMP}), f_2(\text{OBJ})), \dots\}$ . Men om me ikkje har slikt samsvar i argumentstruktur, vil *fpairs* ha ein vanskelegare jobb.

Eit større problem er nok adverbial (elementa i  $\text{ADJUNCT}_i$ ), kor f-strukturane ikkje gir like greie hint om kva for substrukturar som høyrer saman<sup>4</sup>. Ein del av masteroppgåva vil altså vere å komme med forslag til funksjonen *fpairs*.

f-samanstillinga kan nyttast til å gi ein samanstilling av frasane dei representerer.  $\phi^{-1}$  gir no ei samanstilling mellom funksjonelle domene i c-strukturane, me har t.d. ei lenkje mellom domenet  $d_1 = \{X, Y, Z\}$  på språk 1 og  $d_2 = \{U, V, W\}$  på språk 2. Kvar node frå  $d_1$  vil kunne (symmetrisk) samanstillast med ein (eller ingen) frå  $d_2$ .

Her kan me utnytte det at frasestrukturane i dei ulike grammatikkane er tufta på same X-bar-prinsipp. Ein  $XP \in d_1$  skal sannsynlegvis samanstillast med ein  $YP \in d_2$  (der  $X$  og  $Y$  gjerne er same symbol, men au kan vere t.d.  $V$  og  $I$ ). I tillegg skal høge nodar sannsynlegvis samanstillast med andre høge nodar, der alt anna er likt, medan mangel på samsvar i samanstillinga til døtre kan føre til at mornodar ikkje skal samanstillast; ein formalisering dette steget, med diskusjon rundt problema, vil au inngå i masteroppgåva.

<sup>3</sup>Ved å nytte argumentplass kan me enkelt få til lenkjer mellom GF-ar med ulike namn, som vist i dømet.

<sup>4</sup>Det er mogleg at f-samanstillinga av adverbial kan tene på informasjon frå (og difor bør skje etter) samanstillinga av frasane som projiserer argumentfunksjonane.

## Kapittel 2

# Bakgrunn og omgrepsavklaring

### 2.1 SKRIV Relaterte metodar

Automatisk frasesamanstilling er eit nytt felt. Det finst allereie veldig gode system for automatisk setningssamanstilling, og automatisk samanstilling av ord har komme langt, men nivåa mellom ord og setning ser ut til å by på fleire problem. «by på fleire problem» – weasel wording, TODO omskriv. Dei ulike tilnærmingane som finst er prega av formåla til utviklarane. Det er verdt å merkje seg at ordet «frase» ofte blir nytta i litteraturen om strenger av ord (N-gram) som ikkje treng vere syntaktiske konstituentar, igjen avhengig av formålet med metoden.

Innanfor korpuslingvistikken har t.d. Piao & McEnery (2001) nytta enkel kolokasjonsinformasjon for å først finne sannsynlege nominale frasar på engelsk og kinesisk (dvs. «chunking»), og så samanstillе desse; her er evalueringsgrunnlaget rett og slett ein manuell gjennomgang av dei mest sannsynlege omsetjingane dei får.

Den manuelle frasesamanstillinga i Samuelsson & Volk (2006), nemnt over, blei nytta som evalueringsstandard for den automatiske metoden i Samuelsson & Volk (2007). Her kjem frasesamanstillinga frå ei ordsamanstilling, der berre N-gram som svarer til ein syntaktisk node blir lenkja som frasar (meir om denne metoden nedanfor). Formålet er å lage ein parallell trebank, kor det altså er unyttig å lenkje «frasar» som *ikkje* er konstituentar.

Sjølv om fraselenkjer kan vere nyttige i korpuslingvistikken er det hovudsakleg innanfor statistisk maskinomsetjing at ein har forska på samanstilling av frasar. Koehn et al. (2003) gir ei grundig evaluering av ulike statistiske metodar for frasesamanstilling til bruk i stokastisk maskinomsetjing. Dei nyttar BLEU-skåren til å rangere resultata (Papineni et al., 2001, i Koehn et al., 2003, s. 51), som gir ei rangering ved (N-grambasert) samanlikning med ferdig omsett tekst.

Den første metoden, AP, er reint N-grambasert. Dei nyttar verktøyet Giza++ (Och og Ney, 2000, i Koehn et al., 2003, s. 50) til å indusere ordsamanstilling frå eit setningssamanstilt korpus (vha. «modell 4» for ordsamanstilling, utvikla ved IBM av Brown et al. (1993)). Denne samanstillinga er 1-til-n (t.d. eitt engelsk ord til to

TODO  
(innleiing?):  
- reine  
N-gram-  
samanstillingar. ■  
dependensba-  
serte  
- ulike formål  
for  
samanstilling  
gir ulike  
metodar  
- kort oversyn  
over  
LFG-termar

fleire slike?  
meir om dette,  
algoritmen

franske), så dei finn ordsamanstilling for begge retningar og tek så snittet av alle moglege N-gramsamanstillingar som ikkje er i konflikt med ordsamanstillingane. Dei føyer så på ord frå unionen av desse vha. nokre enkle heuristikkar.

Den andre metoden, *Syn*, tek berre med dei frasane som står under syntaktiske nodar i eit parsar korpus; frasesamanstillinga til *Syn* er ein delmengd av den i *AP*. Denne syntaktisk informerte modellen gav ein mykje dårlegare BLEU-skåre enn den reint N-grambaserte modellen (faktisk dårlegare enn omsetjingane frå den opphavlege modell 4 for ordsamanstilling, utan frasesamanstilling). Dei forklarar dette med den store mengda uttrykk som ikkje utgjer syntaktiske konstituentar i følge parsaren deira, men likevel konsekvent blir omsett til visse uttrykk på det andre språket (t.d. «es gibt» på tysk til «there is» på engelsk).

Seinare resultat har vist at ein *kombinasjon* av syntaktisk informerte metodar med reint N-grambaserte modellar (dvs. i motsetning til å berre fjerne samanstillingar mellom ikkje-konstituentar) kan auke skåren i ein maskinomsetjingsevaluering, både om ein som i *Syn*-modellen nyttar frasestrukturinformasjon, men i endå større grad om ein nyttar dependensinformasjon (Tinsley et al., 2007; Hearne et al., 2008). Dette er interessant med tanke på at LFG-analysane gir begge typar informasjon.

Riezler & Maxwell (2006) utvikla ein metode for å kombinere frasebasert statistisk maskinomsetjing med LFG-basert setningsgenerering. Dei finn ei n-til-m-ordsamanstilling med Giza++ som i metodane over, men parsar i tillegg setningane i LFG. Dei to moglege f-strukturane som liknar mest blir valt ut, og frå ordsamanstillinga finn dei mange-til-mange-korrespondansar mellom substrukturane i f-strukturane. Ved å leggje til LFG-basert generering fekk det kombinerte systemet betre resultat på langdistanseavhengnader og generalisering til nye uttrykk med strukturell likskap til tidlegare observerte uttrykk.

Så langt har eg ikkje kome over metodar som går i motsett retning, altså prøver å finne eller betre på frase- og ordsamanstilling ut frå ein LFG-parse – det er dette som er strategien til programmet *lfgalign* i kapittel 4 – men det er stor overlapp mellom krava som kjem i kapittel 3 og dei gitt i den første publiseringa i XPar-prosjektet, Dyvik et al. (2009).

Dette er motsett retning av det mitt program gjer, nemne seinare?

## 2.2 SKRIV Eit kort oversyn over leksikalsk-funksjonell grammatikk og terminologi

I dei følgjande kapitla nyttar eg ein del LFG-terminologi (i tillegg til eit par eigne termar). Difor gir eg her eit kort oversyn over det som kan vere nytt for dei som er meir vand med andre grammatiske rammeverk.

**modellteoretisk (vs derivasjonelt)** LFG er eit modellteoretisk, ikkje-derivasjonelt, rammeverk for grammatikk. Pullum & Scholz (2001) gir ein god gjennomgang av forskjellen mellom derivasjonelle (enumerative) grammatikkar og modellteoretiske grammatikkar, kor førstnemnde definerer *mengder av ut-*

*trykk* ved avleiing frå startsymbol, medan sistnemnde gir skildringar av *enkeltuttrykk*. Ein modellteoretisk grammatikk kan i tillegg skildre strukturen (eller dei moglege strukturane) til *fragment* av setningar, og denne strukturen er lik det bidraget som fragmentet tilfører skildringa av heile setninga. Det tilsvarande er ikkje mogleg å gjere derivasjonelt. Pullum & Scholz (2001, s. 32–33) gir t.d. eit fragment som kjem midt i eit høgreforgreina tre; ein derivasjonell skildring ville måtte skildre treet over eller under, men utan informasjon om kva som kjem til høgre eller venstre kan me ikkje (på ein ikkje-vilkårleg måte) skildre subtreet utanfor fragmentet heilt fram til terminal- eller startsymbol.

**modellteoretisk (vs derivasjonelt)** f-struktur

**modellteoretisk (vs derivasjonelt)** c-struktur

**modellteoretisk (vs derivasjonelt)** endosentrisitetsprinsippa

**modellteoretisk (vs derivasjonelt)**  $\bar{X}$

**modellteoretisk (vs derivasjonelt)** diskontinuerlege konstituentar

**modellteoretisk (vs derivasjonelt)**  $\phi, \phi^{-1}$

**fraselenkjer vs frasesamanstilling** Eg nyttar her termene *lenkjing* og *samanstilling* i omtrent same tyding som dei engelske termene *link* og *alignment*, kor ei samanstilling er ei mengd lenkjer. Dette er ekvivalensforhold som me kan finne mellom lingvistiske *representasjonar* (f-struktur, c-struktur) eller *uttrykk* (ord, setningar); lenkjing mellom dei siste altså er meir ateoretisk / datanært.

## Kapittel 3

# Krav til frasesamanstilling

### 3.1 Innleiing

I denne delen prøver eg å finne fram til kva som er den best moglege frasesamanstillinga. Eg argumenterer for at «best» her må tolkast i forhold til eit formål, her å finne samsvar mellom kasusmarkering og semantisk rolletildeling. Som utgangspunkt har eg visse krav for ordsamanstilling gitt i Thunes (2003), saman med krava for frasesamanstilling i Dyvik et al. (2009). Eg viser kvifor ein, for våre formål, må revidere kravet til Thunes om likskap i argumentstruktur. Eg gir nokre døme for å grunngje krava i Dyvik et al. (2009), i tillegg til å utdjupe dei for å gjere dei enklare å implementere i kapittel 4. Dette involverer au å omformulere krava for c-struktursamanstilling slik at dei ikkje refererer til ordlenkjer, berre f-strukturlenkjer. Sidan eit av måla med Xpar-prosjektet er å finne ut kor mykje frasesamanstillingsinformasjon me kan få ut av parallellismen i f-strukturane (eller, sett frå den andre sida, kor uavhengig ein kan gjere seg av den bottom-up-informasjonen ei ordlenkje gir), blir det eit avleidd mål å formulere frasesamanstillingskrava med referanse til f-strukturane der det går an.

### 3.2 Formål med frasesamanstilling

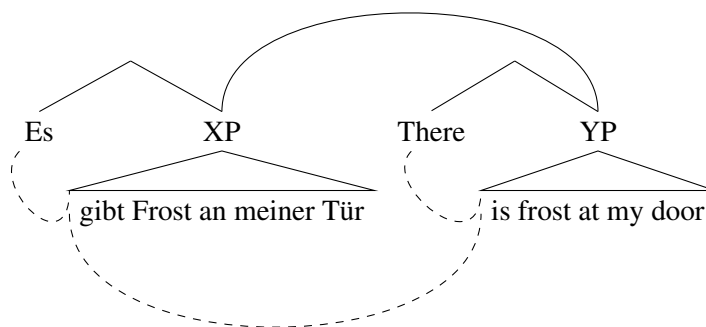
Ei frasesamanstilling er ein slag annotasjon av eit korpus. På same måte som oppbygginga av eit korpus avheng av formålet til korpuset, kan ein ikkje definere den ideelle annotasjonen av eit korpus utan å ta høgd for kva ein skal nytte annotasjonen til.

Me kan illustrere dette med eit enkelt, praktisk døme: ved automatisk ordklassetagging må ein gjerne avvege mellom dekning (å finne flest moglege analysar for flest mogleg ord) og presisjon (å berre ende opp med korrekte analysar). Viss formålet er å annotere ein leksikografisk ressurs, vil det vere viktigare med høg dekning på bekostning av presisjon, sidan leksikografen gjerne leiter etter nye/kreative bruksområde av ord. Skal taggaren nyttast til maskinomsetjing i staden, kan ein ikkje nytte meir enn éin analyse til slutt, så her er presisjon viktigast.

Sjølvsagt kan ein her seie at den *ideelle* annotasjonen vil vere å berre ha korrekte analysar, men sjølv ved ideelle krav er formålet viktig: er ein ute etter å finne N-gram som ofte blir omsett med kvarande, men som *ikkje* er syntaktiske konstituentar, er det klart at retningslinjene nedanfor ikkje er så nyttige.

Sidan utviklinga av automatisk frasesamanstilling hovudsakleg har skjedd innanfor frasebasert statistisk maskinomsetjing (PBSMT), kjem me ikkje utanom ei samanlikning her. I PBSMT er formålet med ei fraselenkje å betre maskinomsetjing på eitt eller anna mål, t.d. BLEU-skåren. BLEU-skåren samanliknar ferdig omsett tekst (ein gullstandard) med det automatisk omsette, ved å sjekke kor mykje N-gram-overlapp det er mellom tekstene. Ei fraselenkje mellom N-grammet *es gibt* og *there is* (dvs. eit auka sannsyn for å nytte slike par i omsetjinga) kan gi ein høgare endeleg skåre i BLEU. Som vist i Koehn et al. (2003) fekk dei ein lågare BLEU-skåre når dei fjerna lenkjer mellom nodar som, i følgje ein robust statistisk PCFG-parsar, ikkje var syntaktiske frasar (konstituentar). Dvs. at i figur 3.1 vil lenkja vist ved den prikkete linja bli fjerna frå mengda over moglege lenkjer om ein berre held seg til syntaktiske konstituentar, og  $p(es\ gibt, there\ is)$  vil ikkje bli tilsvarande auka i den statistiske omsetjingsmodellen. Sidan PBSMT, som skildra i Koehn et al. (2003), er agnostisk til syntaktiske høve i omsetjingssteget<sup>1</sup> er det for dei ingen grunn til å berre halde seg til samanstilling mellom syntaktiske konstituentar; dei har i utgangspunktet meir nytte av kollokasjonsinformasjon.

todo: referere til den faktiske parsaren? det var Bikel kanskje?



Figur 3.1: N-gram-samanstilling versus syntaktiske frasar

Men sett no at me ikkje har som formål å nytte frasesamanstillinga til reint N-grambasert omsetjing. Kva for *lingvistiske* krav kan me stille til å kalle to frasar samanstilte? Me må i alle fall tillate ein del skilnad. I alle større parallelltekster vil parallellstilte setningar ha visse syntaktiske og semantiske<sup>2</sup> omsetjingsskifte, t.d.

<sup>1</sup>Både omsetjingsmodellen og språkmodellane er reint N-grambaserte her, og har difor ikkje nytte av syntaktisk informasjon (i motsetning til syntaktisk informert generering slik Riezler & Maxwell (2006) implementerer).

<sup>2</sup>Sidan eg føreset setningssamanstilte data, kjem eg ikkje inn på diskurs-/pragmatiske verkningar, med mindre dette fører til forskjellar innanfor setningane (sjå t.d. del 3.5 om lenkjer mellom koreferente substantiv og pronomen).

leksikalisering av syntaktiske konstruksjonar eller omvendt, endring av ordklasse, presisering/depresisering, endringar i leksikale trekk (t.d. telleleg/utelleleg), osb. (Munday, 2001, s. 56–62), slik at den einaste fullstendige, «perfekte» samanstillinga vil vere identitetsfunksjonen. Kor mykje mangel på samsvar me godtek blir då avgjort av formålet med samanstillinga.

Eitt av formåla med samanstillinga i denne oppgåva er å kunne oppdage korleis ulike språk realiserer semantiske roller syntaktisk; då spesielt i forhold til hypotesane gitt i XPar (2008, s. 7), t.d. at «case marking might be useful to further determine a given argument's semantic role». Skal me finne det siste, må me altså kunne lenkje frasar med ulik kasusmarkering, men ha krav om lik tildeling av semantiske roller; samtidig skal me sjå at me ikkje kan ha krav om lik syntaktisk funksjon. I tillegg vil me sjølvsagt ikkje lenkje på tvers av konstituentgrenser, sidan det er fullstendige konstituentar<sup>3</sup> som fyller dei semantiske rollene.

Eit anna mogleg formål er å nytte desse frasesamanstillingane til maskinomsetjing. Riezler & Maxwell (2006) nyttar ein stokastisk frasesamanstilling til å oppdage transfer-reglar for bruk i LFG-basert generering i maskinomsetjing. Dette er reglar som omsett fragment av ein f-struktur på kjeldespråket til f-strukturfragment på målspråket. (Eit krav på utforminga av moglege transfer-reglar hindrar at ein får reglar som lenkjar ikkje-konstituentar, eg kjem tilbake til dette nedanfor.) Samanstillinga utvikla her burde au kunne nyttast til å finne slike transfer-reglar, men dette er ikkje noko eg har lagt vekt på.

Nedanfor gir eg eit forslag til krav for frasesamanstilling, med desse formåla i tankane. Om alle krava er moglege å implementere, er eit separat problem.

### 3.3 Frasesamanstilling i ein LFG-trebank

Samanstilte frasar bør ha nok semantisk likskap til å kunne opptre som omsetjingar i liknande omgivnader (Dyvik et al., 2009, s. 74). Thunes (2003) gir nokre prinsipp – som er passande å ha som utgangspunkt – for å fastslå det som kan kallast *omsetjingsmessig korrespondanse* (her for ordsamanstilling). Dette er prinsipp som skal gjelde for eit litt forskjellig formål, men som au «ligger nær opp til det vi intuitivt mener er riktig» (Thunes, 2003, s. 2). Prinsippa blir nytta til å lage ein gullstandard for ordsamanstilling<sup>4</sup>, hovudsakleg for dei opne klassene, og er definert ved å vise til kva for rolle eit argumentord spelar, eller kva for rolletildeling eit predikat eller modifierande ord gir. Så for å t.d. samanstill to verb må dei ha like mange semantiske argument (men argumenta treng ikkje alle realiserast syntaktisk) og dei

<sup>3</sup>LFG tillèt som nemnt diskontinuerlege konstituentar, men dette er ikkje det same som ikkje-konstituentar av typen «es gibt» / «there is».

<sup>4</sup>(Thunes, 2003, s. 2): «Våre prinsipper er satt opp for å tjene et bestemt formål, nemlig å samle inn data som metoden i Semantic Mirrors skal anvendes på», ein metode for å automatisk finne WordNet-liknande relasjonar frå parallelltekst. I denne metoden vil det vere naturleg med høge krav til presisjon, men kanskje lågare krav til dekning: speilmetoden skal finne leksikale semantiske forhold som held på *typenivå*, medan for trebanken er det viktigare korleis me kan annotere eit *token* av t.d. eit verb i ein viss VP i ei gitt korpussetning.



må *tildele same roller*; medan argumenta må *spele same rolle*, og både argument og adjunkt må vere *koreferente*. Lenkja ord må vere del av frasar som spelar same rolle i «det som er felles i interpretasjonene av [dei to setningane]» (Thunes, 2003, s. 3).

Viss me tek utgangspunkt i det siste, vil det vere naturleg å i tillegg lenkje desse frasane som spelar same rolle i «det som er felles i interpretasjonene».

Krava for ordsamanstillinga må au vere fylt for at desse frasane kan samanstillast. Ei ordsamanstilling er altså naudsynt for ein frasesamanstilling, og omvendt. Dette er berre problematisk om me føreset at det eine er derivert av det andre; men dette har me ingen *a priori* grunn til å gjere. Krava eg her utviklar bør i staden sjåast på som *skrankar* på moglege samanstillingar i modellen (jamfør 2.2 om modellteoretiske grammatikkar), heller enn derivasjonelle forhold. Samtidig er det som nemnt eit mål å finne ut kor uavhengig me kan gjere oss av ordlenkjingsinformasjonen (dette er au nyttig for implementasjonen), utan at det treng å gi krava ei *retning*.

Ei frasesamanstilling er ei skildring av forhold mellom *fragment* av setningar, dette er endå ein grunn til at det er naturleg å skildre dei ønskelege forholda som skrankar på moglege samanstillingar. Me kan setje skrankar på f-struktur-, konstituent- og ordsamanstilling samtidig, utan å måtte ha krav om at den eine samanstillinga er fullstendig (eller delvis) avleia av den andre, før me veit om eit slikt avleiingsforhold er empirisk fundert. Me kan i tillegg ha ufullstendige samanstillingar i dei tilfella der det er ufullstendig samsvar mellom setningane (der ei fullstendig samanstilling ville brutt visse krav).

Sidan metoden er mynta på bruk i ein LFG-parsa trebank, og delvis vil nytte denne annotasjonen som datagrunnlag, er det naturleg å nytte same konsept som blir nytta i LFG<sup>5</sup> (f-struktur, c-struktur, endosentrisitetsprinsipp,  $\bar{X}$ -tre, osv.) au i desse krava til den «beste» frasesamanstillinga; i den grad LFG gir ein generaliserbar skildring av syntaks, bør desse krava vere generaliserbare til andre teoriar.

Utan skrankar i det heile vil alt kunne lenkjast til alt (noko som er like unyttig som å ikkje lenkje noko); i del 3.4 ser eg på kva for typar element i dei lingvistiske analysane (ord, grammatiske trekk, konstituentar, ...) det er fornuftig å tillate lenkjer mellom. I avsnitta nedanfor spesifiserer eg kva som må til for at me skal lenkje element av desse typane.

### 3.4 Kva kan lenkjast?

Viss to uttrykk er samanstilt på setningsnivå (slik at me dimed kan gå ut frå at dei er omsetjingar av kvarandre), og begge har ein LFG-analyse, så har me iallfall tre ulike nivå kor me kan finne ekvivalensforhold under setningsnivå:

<sup>5</sup>I tillegg finst andre positive biverknader av ein LFG-basert frasesamanstilling for bruk i denne samanhengen, som at ein kan studere kor parallelle dei parallelle grammatikkane i ParGram-prosjektet (Butt et al., 2002) faktisk er, på ulike nivå (leksikon og argumentstruktur, c-struktur, f-struktur).

1. mellom ord i setningane,
2. mellom f-strukturar,
3. mellom c-strukturnodar.

På begge språk har me alle nivå – det er ingen grunn til å lenkje på tvers av nivå sidan forhold mellom desse nivåa er implisitt i LFG-analysen.

Alle ord i setninga er *kandidatar* for samanstilling med ord i omsetjinga, men det kan godt hende at eit ord *ikkje* har ei lenkje, og me kan heller ikkje utelukke at det finst mange-til-mange-lenkjer som ikkje kan «delast opp». Dette gjeld au nodane i c-strukturen.

Når det gjeld f-strukturane er det ganske mange element me teoretisk sett kunne ha lenkja, t.d. enkelttrekk som kasus eller dei uordna mengdene med adjunkt, men det som er mest *nyttig* er nok å berre lenkje der det er ei nær kopling til orda i setninga. Sidan alle PRED-element i ein f-struktur unikt står for predikerande ord, kan me – gitt to samanstilte setningar – la *kandidatane for samanstilling på f-strukturnivå* inkludere<sup>6</sup> alle desse PRED-elementa i f-strukturane til setningane. PRED-element representerer semantiske bidrag som oftare er naudsyne på begge språk i omsetjingar, medan andre f-strukturtrekk gjerne er valfrie på det eine av språka; det er ikkje alle språk som har t.d. obligatorisk kasusmarkering, og ein vil kanskje nytte trebanken til å oppdage nettopp slik variasjon. PRED-elementa er i tillegg gjerne enklare å knyte direkte opp mot konkrete tekststrengen, medan t.d. aspekt kanskje er umogleg å skilje frå tempus i affikset.

Eg føreslår følgjande føringar:

- (1) Ei samanstilling av to PRED-element i f-strukturane tilseier at:
  - a. f-strukturane til desse er lenkja,
  - b. orda i setningane som projiserer PRED-elementa tek del i ei samanstilling med kvarandre (kor andre ord kan vere involvert), og at
  - c. iallfall dei øvste nodane i det funksjonelle domenet<sup>7</sup> til f-strukturen er samanstilt.

(Underordna nodar i det funksjonelle domenet kan berre lenkjast om visse krav, gitt nedanfor, er oppfylt. Me kan altså gjerne ha c-strukturnodar som ikkje er lenkja til andre nodar.)

Påstandane over må forsvarast. Punkt (1-a) og (1-c) over seier at viss PRED-elementa projisert av t.d. to verb i verbfrasar er lenkja, vil *heile* VP-ane vere lenkja (både VP-nodane som dominerer dei lenkja funksjonelle domena og f-strukturane

backe det med  
eksemplar i  
trebank; kople  
til  
adj-arg-lenkje

<sup>6</sup>I del 3.6 kjem eg tilbake til spørsmålet om me vil inkludere visse f-strukturar utan PRED-element i kandidatane for samanstilling.

<sup>7</sup>Det funksjonelle domenet til ein f-struktur er gitt ved  $\phi^{-1}$ , inversen av c-til-f-strukturavbildinga, og tilsvarende dei nodane i c-strukturen som projiserer denne f-strukturen, t.d. ein VP-node med dominerande IP og CP (Bresnan, 2001, s. 126). Sidan dette er inversen av ein funksjon, kan me ha diskontinuerlege konstituentar i same funksjonelle domene (fleire funksjonsargument som gir same verdi).

frå ytre PRED til verba), det er dette som gjer det til ei fraselenkje; medan i følgje punkt (1-b) vil denne fraselenkja leie til at sjølve verba au er lenkja, ein sterkare påstand sidan dette tilseier at *PRED-samanstilling impliserer ordsamanstilling*. I visse tilfelle er dette heilt uproblematisk, t.d. viss *I slept down by the river* skal lenkjast med *Eg sov nede med elva* vil me uansett lenkje *slept* og *sov*; dette kan gjelde transitive verb au:

- (2) a. The locusts have no king, just noise and hard language  
 ↔  
 b. Grashoppene har ingen konge, berre støy og krasse ord

*have/har* tek del i VP-samanstillinga *have no king.../har ingen konge....*

Som nemnd over; ordsamanstillinga treng ikkje vere ein-til-ein, det punkt (1-b) seier er at desse orda iallfall er ein del av ein samanstilling med kvarandre (i (2) altså VP-samanstillinga). Kanskje er dette ei mange-til-mange-lenkjing som ikkje *kan* reduserast til ein-til-ein-lenkjingar; eller kanskje er det som i (2) mogleg å skilje ut delsamanningar, som *have/har*. Eg kjem tilbake til dette i del 3.12 om argumentstruktur og adjunkt.

Alle nodar i c-strukturen (alle syntaktiske *frasar/konstituentar* i setninga) som kan koplast til PRED-haldande f-strukturar, vil altså vere kandidatar for samanstilling på c-strukturnivå (dette inkluderer diskontinuerlege konstituentar), men ikkje alle vil bli samanstilt.

der  
 ADJUNKT  
 ikkje er  
 realisert,  
 lenkjer me  
 ikkje PRED.  
 skal me då  
 ikkje lenkje  
 ord heller?  
 PRED->ord ::  
 iallfall  
 PRED<-ord ::  
 ?  
 PRED<->ord  
 PRED, ord  
 vsnittet over  
 er litt rotete  
 TODO

### 3.4.1 TOGROK finst det tilfelle der ordlenkjer ikkje impliserer PRED-lenkjer?

hypotese: det er alltid slik at  
 ordlenkjing av predikerande ord => PRED-lenkje

## 3.5 SKRIV Krav på ordnivå

### 3.5.1 SKRIV LPT :ROTETE:

«a source word WS and a target word WT are taken to correspond translationally only if (i) WT can in general (out of context) be taken to be among the semantically plausible translations of WS, i.e., WT belongs to the set of ‘linguistically predictable translations (LPT)’ of WS, and (ii) WS and WT occupy corresponding positions within corresponding argument structures.»

«a source phrase PHS and target phrase PHT are taken to correspond if (i) they contain corresponding words, (ii) PHS contains no word or phrase corresponding to a target word or phrase outside PHT, and similarly (iii) PHT contains no word or phrase corresponding to a source word or phrase outside PH.»

### 3.5.2 Lik ordklasse?

Ulike språk leksikaliserer same konsept på ulike måtar. Cheung et al. (2002, s. 3) skriv at det engelske ordet *fulfilment* meir naturleg blir omsett til eit verb på kinesisk. Det same gjeld t.d. *solitude* omsett til norsk. Eit georgisk verbalsubstantiv (*masdar*) kan bli omsett til eit verb i infinitiv på norsk<sup>8</sup>. Slike skifte mellom ordklassar er svært vanlege i omsetjing<sup>9</sup>.

Me kan opne for ordklasseoverskridande lenkjer der det er samsvar mellom visse *trekk*, t.d. kan to predikerande ord lenkjast, eller to «nominale» ord. Ein annan måte å gjere dette på er rett og slett å krevje ein viss likskap i argumentstruktur.

## 3.6 Funksjonsord

I tillegg kan me ha ord i setninga som ikkje tilsvarer PRED-element i f-strukturen, typisk funksjonsord (t.d. *som*, *at*). Ved endosentrisitetsprinsippa til Bresnan (2001) er komplementet til funksjonelle kategoriar (C, I, P) ein funksjonell ko-kjerne.

- (3) Skal nodar for ord som ikkje projiserer PRED-element<sup>10</sup> samanstillast, må følgjande krav vere oppfylt:
- a. det funksjonelle domenet (gitt ved komplementet) må vere samanstilt, og
  - b. dei er begge c-strukturhovud.

Om (3-a og -b) er oppfylt, kan me få samanstillinga vist i figur 3.2, og i dette tilfellet er (3-b) oppfylt og (3-a) vil vere oppfylt om me kan samanstillе *cvimda* med *det regnet*.

### 3.6.1 TOGROK *cvimda*<PRO> men regne<>expletive – lenkje? :RO-TETE:

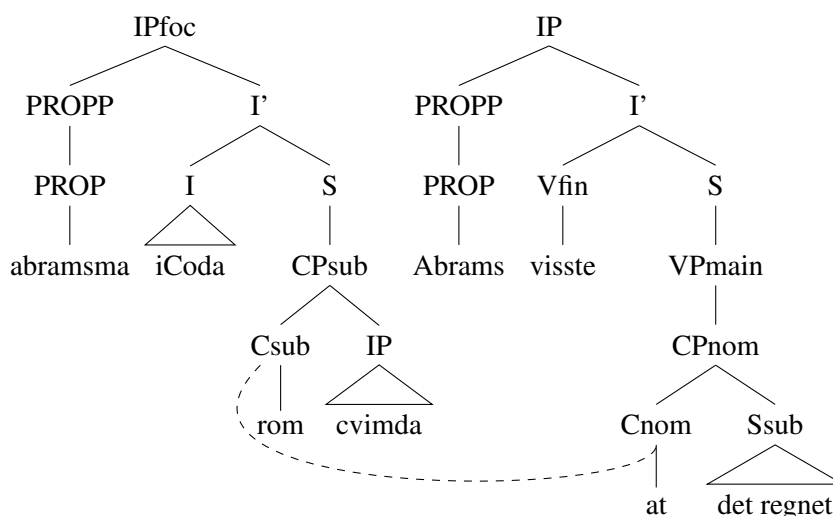
### 3.6.2 TOGROK kva med ekspletivar? ingen PRED men heller ikkje C/F/I :ROTETE:

Kandidatane på f-strukturnivå må jo inkludere desse au. . .

<sup>8</sup>Det georgiske verbalsubstantivet (*masdar*) er i følge Aronson (1990, kap. 2.5) ein *nominal* form, det kan i motsetning til norske verbalsubstantiv og engelske gerundium ikkje ta objekt, men kan ha modifierande substantiv i genitiv.

<sup>9</sup>Munday (Catford (1965), i 2001, s. 61) gir ein gjennomgang av slike *klaseskifte*, og andre typar omsetjingsskifte.

<sup>10</sup>Skal ein lenkje ordet *som* (utan PRED) med ordet *which* (med PRED)? Viss begge står under C i treet, kan det kanskje vere informativt med ein type «defekt» lenkje, sjølv om berre det eine ordet blir rekna for å vere eit innhaldsord. Frasane til deira funksjonelle domene vil uansett vere samanstilt via toppnodane (t.d. CP).

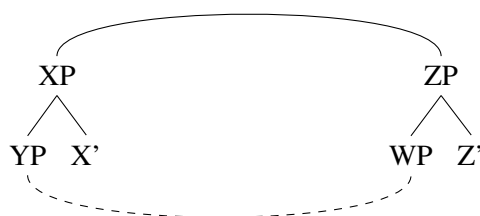


Figur 3.2: Mogleg samanstilling av funksjonsord mellom georgisk og norsk (bok-mål)

### 3.7 Lenkjing av underordna c-strukturnodar

Toppnodane i eit lenkja funksjonelt domene i c-struktur (XP på språk 1, ZP på språk 2) vil ha ein informasjonsmessig korrespondanse, og kan samanstillast. Men det er mogleg å samanstille to toppnoder i funksjonelle domene i c-strukturen utan at nodane under ( $X'$ ,  $Z'$ ) er samanstilt. Ein grunn til å ikkje samanstille desse underordna nodane, vil vere viss spesifikator til X ikkje spelar same rolle i tolkinga som spesifikator til Z, dvs. viss YP og WP i figur 3.3 ikkje er lenkja.

Me kan utelukke lenkjing av ikkje-konstituentar som *there is* ved å krevje at ei fullstendig samanstilling mellom to frasar må vere slik at heile substrukturen au er samanstilt. *There is* og *Es gibt* i figur 3.1 kan då ikkje samanstillast åleine, men berre som del av ei ytre frasesamanstilling. Så når *kan* me samanstille nodane som står under øvste node i f-domenet?



Figur 3.3: Lenkjing av underordna c-strukturnodar

I figur 3.3 der XP og ZP er lenkja, vil YP og WP – i kraft av å vere toppnoder i sine domene – måtte ha ei lenkje i f-strukturen for at c-strukturnodane kan lenkjast

(det kunne jo t.d. hende at f-strukturen projisert av YP samsvarte med den projisert av Z', eller ein struktur under Z').

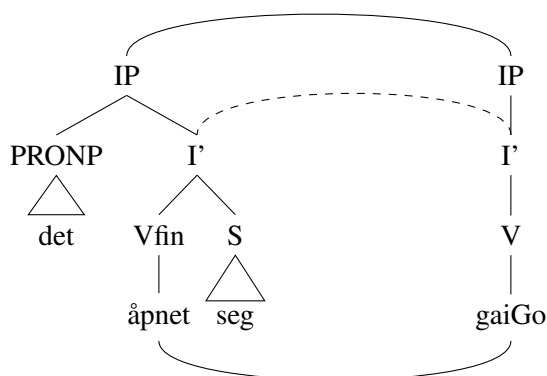
Om me skal lenkje Z' og X' i figuren over må dei respektive spesifikatornodane vere lenkja. Me får då følgjande krav:

- (4) Krav for lenkjing av underordna c-strukturnodar:
- c-strukturnodar som ligg under øvste node i to funksjonelle domena kan berre samanstillast med nodar som ligg innanfor desse domena,
  - c-strukturnodar kan berre samanstillast om deira funksjonelle domene er lenkja på f-strukturnivå,
  - om ein c-strukturnode X' som ikkje er toppnode i det funksjonelle domenet har ein søsternode YP, må YP vere samanstilt med ein søsternode til Z' for å samanstille X' og Z'

(4-a) seier at om XP og ZP er samanstilt, der XP er t.d. OBJ til IP, kan ikkje Z' samanstillast med SUBJ til IP osv., men berre til nodar innanfor OBJ-domenet. (4-c) påført figur 3.3 seier altså at spesifikatornodane må vere lenkja for at X' og Z' skal lenkjast (manglande søsternode på den eine sida vil au hindre samanstilling).

I figur 3.2 er alle nodane under S vist i dei to trea i same funksjonelle domene (kvar node under S er annotert med  $\uparrow=\downarrow$ ), så om dei funksjonelle domena er samanstilt (som krev at *rom cvimda* og *at det regner* er samanstilt), vil (4-a og -b) vere oppfylt kva gjeld CP-komplementa – lenkjinga går ikkje ut over dei funksjonelle domena. Sidan Csub og Cnom er funksjonelle kategoriar er dei au samanstilt via samanstillinga av S-nodane og føringane i (3), og (4-c) er då oppfylt. (4) står altså ikkje i vegen for å samanstille IP-en over *cvimda* og Ssub.

I figur 3.4 derimot (?), kan me ikkje samanstille I'-nodane. PRONP-noden, spesifikator på den norske sida, er ikkje lenkja med nokon spesifikator på den georgiske sida. Den informasjonen (her reint syntaktisk) som ordet *det* tilfører IP, ligg under I' på georgisk. Om me skulle lenkja I', måtte me altså hatt ein georgisk spesifikator som var lenkja til den norske PRONP.



Figur 3.4: Umogleg samanstilling av funksjonsord mellom bokmål og georgisk

**3.7.1 SKRIV døme! :ROTETE:****3.7.2 TOGROK me<sub>OBJ</sub> gusta X<sub>SUBJ</sub> // I<sub>SUBJ</sub> like X<sub>OBJ</sub> ?? :ROTETE:****3.7.3 TOGROK korleis finn me *there is*-lenkjer då? :ROTETE:**

(og kva skal me med dei?)

«Til gjengjeld vil me få lenkjer sjølv om me har mellomståande ord (*There never is*) som opptre utanfor N-grammet på det andre språket.»

**3.8 TOGROK mange-til-mange-lenkjing i f-strukturane? :ROTETE:**

Eg er litt usikker på om me skal ha slike mange-til-mange-korrespondansar i f-strukturane; eg har rekna med at ei f-strukturlenkje *impliserer* ei slag lenkjing mellom det som er innanfor f-strukturane; men i Riezler & Maxwell (2006) er det i staden berre eit krav om at desse f-strukturane er lenkja i same transfer-regel.

Riezler & Maxwell (2006, s. 40–41) tillet mange-til-mange-lenkjing mellom f-strukturar, så lenge alle f-strukturane som blir lenkja til slutt opptre i same transfer-regel. Frå følgjande setningspar:

- (5) Dafür bin ich zutiefst dankbar  
I have a deep appreciation for that

lenkjar dei {*zutiefst*} med { *a, deep, appreciation* }, men sidan {*appreciation*} er samanstilt med {*dankbar*}, må transfer-regelen inkludere { *zutiefst, dankbar* } på den eine sida og { *a, deep, appreciation* } på den andre.

**3.8.1 SKRIV Kva inneber ei mange-til-mange-lenkjing? :ROTETE:****3.9 SKRIV Mangel på samsvar i syntaks og semantikk :ROTETE:**

(Kruijff-Korbayova et al., 2006, s. 5) gir følgjande døme:

- (6) nikdy nebyl  
never was.not  
'has never been'

*nebyl* blir «svakt» samanstilt med *never*, men «sterkt» samanstilt med *has ... been* i deira system. I tillegg er det ein sterk samanstilling mellom *never* og *nikby*.

### 3.10 TOGROK Diskontinuerlege einingar :ROTETE:

- diskontinuerlege einingar (Cheung et al., 2002, s. 4) @books.google – skal dei eigentleg samanstillast? Kva for problem gir dei i forhold til c-strukturnivåsamanning?■

#### 3.10.1 TODO døme på diskontinuerlege konstituentar som er lenkja :ROTETE:

### 3.11 TOGROK Er «compounds» frasar? :ROTETE:

(Giegerich, 2006, p. 1)

### 3.12 Krav om lik argumentstruktur

Thunes (2003) gir som nemnd eit krav om at *predikat må ha tilsvarande semantiske argument* for å samanstillast.

Om det alltid er slik at to predikat har like mange argument, som kjem i same rekkjefølgje i argumentstrukturen, vil det gjere den praktiske oppgåva med å samanstille predikata, og argument med argument, mykje enklare. Men kan me stille så sterke krav?

Sett at ein setning på språk 1 har ei *at*-setning som adjunkt, medan denne setninga på språk 2 er eit argument, og at desse setningane ville vore samanstilte om dei opptredde åleine. Om dei uttrykkjer same proposisjon og *speler same rolle i verbsituasjonen*, synest det naturleg å lenkje desse.

Omsetjingsrelasjonar gir data for verbsituasjon, på eit meir generelt grunnlag enn det me kan få frå einspråklege analysar åleine. Om me har gode semantiske grunnar for å kalle ein deltakar i ein verbsituasjon eit argument på eitt språk, vil dei same grunnane gjelde for omsetjingsmessig korresponderande verb på andre språk. Ein kan då nytte unionen over alle argument til korresponderande verb til å karakterisere kva ein meiner med *deltakarane i verbsituasjonen*. Syntaktiske forhold i språket kan sjølvstendig gi grunnar til å *ikkje* kalle dette eit argument (om det er mogleg å finne akseptable syntaktiske grunnar for å kalle noko ein adjunkt heller enn eit argument).

For å gjere dette konkret kan me sjå på setning 7 i MRS-suiten (?)<sup>11</sup>:

- (7)    abramsi            brouns            daenajleva sigaretze,    rom cvimda  
          Abrams.NOM Brown.DAT vedde.3SG sigarett.om, at    regne.3SG.IMP  
          ‘Abrams veddet en sigarett med Brown på at det regnet’

I følgje LFG-parsen til desse setningane har hovudpredikata svært ulike argument-

<sup>11</sup>Setningane i første og tredje linje i døma er direkte henta frå MRS-suiten, med mindre anna er opplyst.



struktur<sup>12</sup>. Det norske *vedde* har fire argument, medan *da-najleveba* har to (Abrams og Browne), kor at-setninga på norsk og *rom cvimda* uttrykkjer same proposisjon og spelar same rolle i verbsituasjonen. Den engelske LFG-parsen av den tilsvarende setninga (mine omsetjingar) gir tre argument, *with* blir her adjunkt, medan den tyske grammatikken, som au har tre argument, gjer *at*-setninga til adjunkt. I (8) nedanfor har eg representert dei omsetjingsmessig korresponderande frasane i f-strukturane med dei norske omsetjingane for å illustrere dette:

- (8) a. Adams veddet en sigarett med Browne (norsk bokmål)  
på at det regnet.

$$\left[ \begin{array}{ll} \text{PRED} & \text{'vedde<Abrams, sigarett, Browne, regne>'} \\ \text{ADJUNCT} & \{\} \end{array} \right]$$

- b. abramsi brouns daenajleva sigaretze, rom cvimda. (georgisk)

$$\left[ \begin{array}{ll} \text{PRED} & \text{'da-najleveba<Abrams, Browne, regne>'} \\ \text{ADJUNCT} & \{\text{sigarett}\} \end{array} \right]$$

- c. Abrams hat mit Browne um eine Zigarette gewettet, (tysk)  
daß es regnet.

$$\left[ \begin{array}{ll} \text{PRED} & \text{'wetten<Abrams, sigarett>'} \\ \text{ADJUNCT} & \{\text{Browne, sigarett}\} \end{array} \right]$$

- d. Abrams bet a cigarette with Brown that it was raining. (engelsk)

$$\left[ \begin{array}{ll} \text{PRED} & \text{'bet<Abrams, sigarett, regne>'} \\ \text{ADJUNCT} & \{\text{Browne}\} \end{array} \right]$$

Om ein skal ha grammatikkane som datagrunnlag er det altså eit reellt problem kva ein skal gjere med mangel på samsvar i argumentstruktur. Om det alltid var fullstendig samsvar i argumentstruktur, ville det vore trivielt å lenkje argument: viss to korresponderande verb hadde tre argument, ville me lenkja det første med det første, det andre med det andre og det tredje med det tredje. Men om me har analysar som dei over, ser det ut til at me treng bottom-up-informasjon om kva for adjunkt og argument som samsvarer.

Det same gjeld forøvrig lenkjing av adjunkt til adjunkt. Adjunkt plukker ut si eiga rolle der argument får rolla tildelt frå verbet, og f-strukturane har ingen hierarkisk inndeling av desse slik me har for verb og argument, dei er i staden representert som *uordna mengder*.

<sup>12</sup>Analysane er henta 18. mai, 2009, frå <http://decentius.aksis.uib.no/logon/xle.xml>, som implementerer LFG-grammatikkane frå ParGram-prosjektet (Butt et al., 2002).

**3.12.1 forsvare «tilsvarande» :ROTETE:**

Tilsvarende på engelsk: <sup>13</sup>

**3.12.2 TODO Sitere eigen korpusundersøking av variasjon i arg-str?**

Ei undersøking av den frasesamanstilte trebanken SMULTRON (Samuelsson & Volk, 2006) mot LFG-grammatikkane for engelsk og tysk fann at 2 av 15 korresponderande verbtoken<sup>14</sup> for høgfrekvente innhaldsverb fekk analysar kor argument korresponderte med adjunkt (?).

LCS, dorr

**3.12.3 SKRIV kvifor lik arg-str er bra, så kvifor det er eit problem :ROTETE:****3.12.4 TODO Ulik følgje i argumentstruktur**

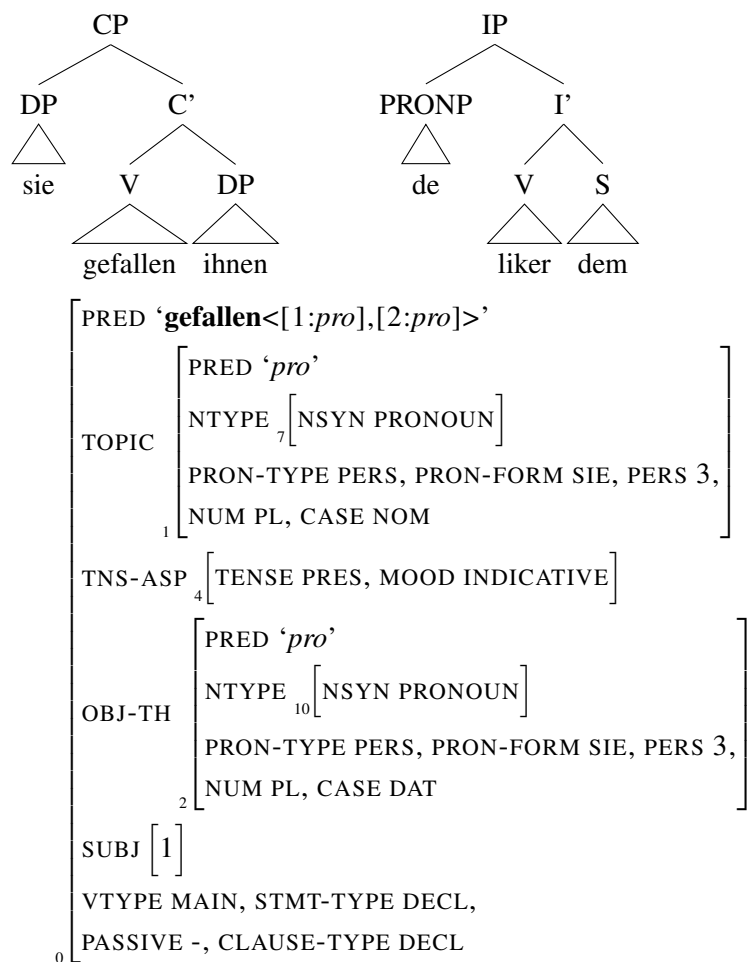
I tillegg til at argument kan lenkjast til adjunkt, kan koreferente argument ha ulik følgje i argumentstrukturen. Det er klart at me vil lenkje objektet til *gefallen* (eller bokmål: *behage*) med subjektet til *like*, og omvendt. Men rekkjefølgje i argumentstrukturane i ParGram-prosjektet er ofte basert på syntaktisk funksjon heller enn rolle, slik at eit verb som har opplevar som objekt og tema som subjekt vil ha opplevar nedanfor tema i argumentstrukturen, medan ei omsetjing av dette verbet kan ha tema nedanfor:

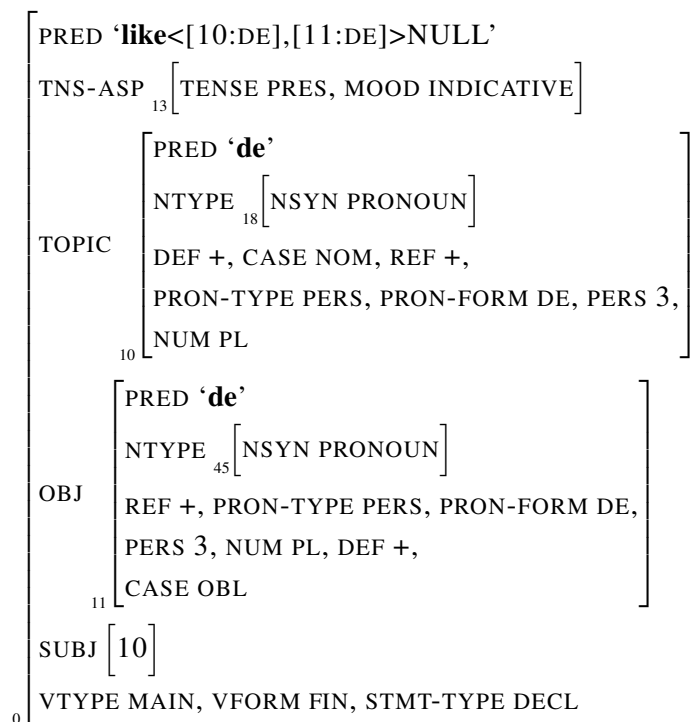
- (9) a. sie<sub>j</sub> gefallen ihnen<sub>i</sub>  
 $\left[ \text{PRED} \quad \text{'gefallen'} \langle \text{de}_j, \text{de}_i \rangle \right]$
- ↔
- b. de<sub>i</sub> liker dem<sub>j</sub>  
 $\left[ \text{PRED} \quad \text{'like'} \langle \text{de}_i, \text{de}_j \rangle \right]$

Argumentstrukturane i (9) har omvendt intern følgje, og som vist ved dette dømet er det heller ikkje noko f-strukturinformasjon me kunne nytta til å sikre lenkjinga *sie/dem* og *ihnen/de*. Igjen ser det ut til at bottom-up-informasjon trengst.

<sup>13</sup>„wagered \* with \* that \*” på Google gir 215 treff, kor 9 av dei første 10 følgjer det intenderte mønsteret.

<sup>14</sup>25 om ein inkluderer analysar kor minst eitt av argumenta ikkje hadde korrekt analyse (t.d. eit PRO der grammatikken burde funne eit substantiv).

**c- og f-strukturar for dømet over :ROTETE:**



### 3.12.5 SKRIV d me med w ger/3 og vedde/4 og gewettet/3 :ROTE-TE:

### 3.12.6 SKRIV (reinskriv) :ROTETE:

Same globale tyding krev i det minste at, i situasjonen verbet denoterer, speler deltakarane same rolle. Men dette er end  meir abstrakt/semantisk enn (semantisk) argumentstruktur. . .

Problem: ikkje-komposisjonell omsetjing. Same globale tyding. Det treng ikkje vere berre pragmatisk forskjell type *kan du lukke d ra* vs *lukk d ra*, kor situasjon gjer setningane like sidan me kan ha konvensjonaliserte konstruksjoner p  L1 kor heile tilsvarer enkeltord p  L2, a la japansk *viss eg ikkje g r p  skulen s  kan det ikkje vere* ~ = *eg m  g  p  skulen*.

Ein f resetnad eg har, er at setningar som er samanstilte faktisk har ein omsetjingsmessig korrespondanse (dette er min data). S  om eit par av ytre predikat ikkje korresponderer er det au ein type data; nemleg at me har ein omsetjingsmessig korrespondanse der det var ein mismatch i ytre argumentstruktur. (Algoritmen b r d  lagre slike mismatches eksplisitt, ikkje berre la vere   lenkje, for det kan vere andre grunnar til at det ikkje kom ei lenkjing. A la ekspertsystem: forklare resonnementet.)

Alternativt ein konstruksjonslenkjing. . .

Kan au ha eit krav om at argstr til  $PRED_{L1}$  er ein slag delmengd av argstr til  $PRED_{L2}$ .

### 3.12.7 SKRIV True Arguments vs True Adjuncts, Pustejovsky :RO-TETE:

- Treng d me f rst...
- Er «with Browne» eit Default Argument for «wager»?
  - D-ARG: he built a house out of bricks
- Adjunkt plukker ut sine egne roller, per definisjon, ved vedde/4 og wager/3 har me ein slik situasjon:

vedde <-----wager >-----<-----wetten  
 \\_\_\_\_with\_/ \\_\_dass

Bottom-up-informasjon vil au vere naudsynt for dei 3 rollene som *er* argument, sidan me kan ha vedde<1,2,3,4> og wager<a,b,c>with<d>, kor det er umogleg   seie om d skal p  plass 1,2,3 eller 4 (dvs. me kan ha vedde<a,b,c,d>, vedde<a,b,d,c>, vedde<a,d,b,c> og vedde<d,a,b,c> – men sannsynlegvis er alt  a,b,c i same rekkjef lgje uansett...)

### 3.13 SKRIV Kan adjunkt lenkjast til nodar under morlenkja?

Krav (vi) i Dyvik et al. (2009, s. 75) krev at viss  $F_s$  og  $F_t$  er lenkja, s  kan ingen adjunkt  $D_s$  til  $F_s$  vere lenkja til nodar utanfor  $F_t$ . Men kan ein  $D_s$  lenkjast til ei dotternode av argument eller adjunkt til  $F_t$ ?

$R_t$  er dotter til  $F_t$ , og m  d  vere lenkja til ei dotter av  $F_s$ ,  $A_s$ . D  m  au alle argument til  $R_t$  vere lenkja til d tre av  $A_s$ , s   $D_s$  kan ikkje lenkjast til argument av dotternodar til  $F_t$ . Kva med adjunkt? Om me finn eit ulenkja adjunkt til  $R_t$  kan me heller ikkje lenkje dette til  $D_s$  ved krav (vi) igjen, sidan  $D_s$  st r utanfor  $A_s$ .

Men om  $D_t$  er ei ulenkja *adjunktdotter* av  $F_t$ , s  vil d tre av  $D_t$  kunne lenkjast til  $D_s$ , s  lenge  $D_t$  forblir ulenkja. Me kan alt  sj  ned i adjunktd tre av  $F_t$  for   lenkje  $D_s$ .

P  same m te b r ein kunne rekursivt sj  ned i ulenkja adjunktd tre av  $R_t$ , men ein b r kanskje ikkje kunne lenkje s  djupt uansett? Ikkje automatisk, uansett.

Programmet mitt vil, gitt to initielle f-strukturar med LPT-korrespondanse, finne alle moglege kombinasjonar av lenkjer som inneheld alle argument og kanskje adjunkt, dvs. om me har

$F_s$  [ PRED p<1,2> ADJUNCT { 3 } ]

$F_t$  [ PRED p<4> ADJUNCT { 5,6 } ]

vil dette vere logisk moglege samanstillingar av «f-strukturd tre»:

(( (1 . 4) (2 . 5)) ((1 . 4) (2 . 6)) ((1 . 5) (2 . 4))  
 ((1 . 5) (2 . 6) (3 . 4)) ((1 . 6) (2 . 4)) ((1 . 6) (2 . 5) (3 . 4)))■

Me luker ut kombinasjonar som bryt med LPT-korrespondanse. Med full informasjon bør me sjølvsagt berre ende opp med éin kombinasjon, t.d. ((1 . 4) (2 . 5)).

Så langt bør altså krav (i-iv) frå Dyvik et al. (2009) vere dekkja.

Me kan krevje at f strukturane-til f strukturdøtre-kan lenkjast rekursivt for at  $F_s$  og  $F_t$  skal lenkjast, t.d. både (1 . 4) og (2 . 5). Men her kjem det (iallfall) to problem.

### 3.13.1 1. Kausativar og inkorporering

Om me har

$F_s$  [ PRED p<SUBJ, 1, 2> XCOMP 2[ PRED q<1> ] ]

$F_t$  [ PRED pq<SUBJ, OBJ> ]

kor pq er t.d. ein kausativ som tilsvare p<..., q>, så vil me ikkje kunne lenkje  $F_s$  og  $F_t$  sidan det bryt med krav (iii),  $F_s$  har eit argument for mykje. Men her vil det kanskje vere naturleg å ha ei ein-mange-lenkje:

(( $F_s$  2) .  $F_t$ )

No kan me sjå på unionen av argument av  $F_s$  (minus XCOMP) og argument av XCOMP, alle argument i denne unionen må då ha LPT-korrespondanse med argument/adjunkt av  $F_t$ , og alle argument av  $F_t$  må ha LPT-korrespondanse med argument/adjunkt av unionen.

Det same bør kanskje skje ved vanleg inkorporering av substantiv, då må det altså vere mogleg å føye saman t.d. verb og objekt; ein kombinasjon av dette og kausativ bør vel vere mogleg, t.d.

$F_s$  [ PRED la<SUBJ, 1> XCOMP 2[ få<1, 3:pengar> ] ]

$F_t$  [ PRED belønn<SUBJ, 1> ]

Igjen ser me på argument frå unionen av ( $F_s$  2 3) minus 2 og 3, og om det er mogleg å lenkje dei til argument/adjunkt av  $F_t$ , og omvendt.

Men det bør kanskje vere grenser for kor langt samanføyning kan gå... eg kan ikkje tenkje meg at me vil lenkje (( $F_s$  2) .  $F_t$ ) eller (( $F_s$  1 2) .  $F_t$ ) her:

$F_s$  [ PRED p<..., 1> XCOMP 1[ PRED q<..., 2> XCOMP 2[ PRED r<...> ] ] ]■

$F_t$  [ PRED pr<...> ]

... men det kan jo hende det finst situasjonar der dette au vil vere rett. Problemet er altså kor me skal setje grensene i implementasjonen. Om me skal prøve å samanføye på alle moglege måtar (altså, der me ikkje har informasjon om LPT), i tillegg til «vanlege» lenkjer, blir det fort komputasjonelt vanskeleg. Me kan sjølvstundt snu på LPT-kravet her, og seie at dette er berre lov der me har positiv informasjon om LPT-korrespondanse, i staden for at det ikkje er lov om me har motstridande LPT-informasjon, det vil nok hjelpe, men det er vanskeleg å finne prinsipelle avgrensingar her.

### TOGROK adjunkt bør ikkje samanføyast? eller?

Det einaste eg kan tenkje meg er at adjunkt ikkje bør vere kandidat for samanføyning (i såfall burde dei vel heller vore analysert som argument?).

#### 3.13.2 2. Adposisjonsobjekt

I følgjande setningspar har me eit objekt «sigarett» som svarer til PP-en «sigaretze» («sigareti» + «ze»), eit adjunkt:

Abrams veddet en sigarett med Browne på at det regnet.  
 abramsi brouns daenajleva sigaretze, rom cvimda.

$F_s$  [ PRED sigarett ]

$F_t$  [ PRED ze<1> 1[ PRED sigareti ] ]

$F_s$  og  $F_t$  er døtre av dei ytre predikata i kvar setning, krav (iii) seier at det må vere LPT-korrespondanse mellom desse for at me skal kunne lenkje «veddet» og «daenajleva». Her synest det feil å føye saman «sigareti» og «ze», ( $F_s . (F_t 1)$ ), sidan «sigarett» ikkje inneheld informasjonen gitt av «ze».

Det finst då to løysingar. Me kan slakke på LMT-kravet ved å la  $L'(F_t) = \{\text{sigaretze}, \text{ze}\}$  (evt.  $\{\text{sigaret}, \text{ze}\}$ ), då kan me lenkje ( $F_s . F_t$ ), medan 1 er ulenkja.

Eller me kan lenkje ( $F_s . 1$ ), kor me har skikkeleg LMT-korrespondanse, men då må me slakke på (iii) og (iv), og altså ha lov til å «hoppe over» ein f-struktur for å lenkje «veddet» og «daenajleva».  $F_t$  er då ulenkja. Det er løysinga valt i Dyvik et al. (2009, s. 75, fotnote 3), og den løysinga eg følgjer nedanfor.

### 3.14 TOGROK kva var poenget med dette? :ROTETE:

«etter og uten er dei einaste prep som tek setn utan å vere arg»

### 3.15 ULEST Cyrus, FuSe-prosjektet :ROTETE:

Cyrus et al. (2004) «Abstract: We report on a recently initiated project which aims at building a multi-layered parallel treebank of English and German. Particular attention is devoted to a dedicated predicate-argument layer which is used for aligning translationally equivalent sentences of the two languages. We describe both our conceptual decisions and aspects of their technical realisation. We discuss some selected problems and conclude with a few remarks on how this project relates to similar projects in the field.»

### 3.16 TODO Konstruksjonar og komposisjonell inekvivallens

$\bar{X}$ -teori føreset at det finst éi dotter i kvart ledd som kan reknast som predikatet for dette leddet. Ei utfordring for  $\bar{X}$ -baserte teoriar er då handsaming av *komplekse predikat*. Desse har fleire grammatiske element innanfor same ledd som alle bidrar med «a non-trivial part of the information of the complex predicate» (Alsina et al., 1997). I LFG er det ein føresetnad at me berre har éin PRED ytterst i kvar f-struktur; ulike mekanismar har blitt føreslått for å handsame dette fenomenet.

I omsette tekster kan me få eit analogt problem:

- (10) It can't be done  
Det lar seg ikke gjøre

Her vil ytre predikat i f-strukturen på norsk vere 'la<det<sub>1</sub>,XCOMP>PRO', kor XCOMP[PRED 'gjøre<NULL, det<sub>1</sub>>NULL'].

På engelsk får me 'can<XCOMP, it<sub>2</sub>>', kor XCOMP[PRED 'do<NULL, it<sub>2</sub>>'].

Skal me lenkje orda *can* og *la*? På *heile konstruksjonen* finn me iallfall eit omsetjingsforhold:

It can't be done	Det lar seg ikke gjøre	
can't be done	lar seg ikke gjøre	
be done	gjøre	s?
_ can't be VPASS	_ lar seg ikke VPASS	??
_1 can _2 be VPASS <sub>3</sub>	_1 lar seg _2 VPASS <sub>3</sub>	??

(kan me få den siste generaliseringa frå trebanken?)

### 3.17 SKRIV definer sitering frå MRS-suiten :ROTETE:

### 3.18 SKRIV setning 7 i MRS-suiten :ROTETE:

Ein samanstilling bør i det minste gi følgjande:



abramsi brouns daenajleva sigaretze, rom cvimda	Abrams veddet en sigarett med Brown på at det regnet
abramsi brouns daenajleva sigaretze	Abrams veddet en sigarett med Brown
brouns daenajleva sigaretze	veddet en sigarett med Brown
daenajleva sigaretze	veddet (en) sigarett (på)
daenajleva	veddet
sigaretze	(en) sigarett (på)
rom cvimda	at det regnet
cvimda	(det) regnet
abramsi	Abrams
brouns	Brown

### 3.19 TOGROK og så finst jo større forskjellar, stilistiske osv... :ROTETE:

### 3.20 TOGROK prosessering, kognitive modellar? :ROTE-TE:

finne empiri frå korleis menneske samanstillar? (dvs., korleis skjer omsetjing)

- Maier (2009), <http://linguistlist.org/issues/20/20-1786.html>

«cross-linguistic structural phenomena in the language production of bilinguals in the specific context of translation.»

- <http://www.linguistlist.org/pubs/diss/browse-diss-action.cfm?DissID=143>
- books.google bialystok????lpb: «Translation has been called “interlanguage paraphrase”», «a metalinguistic skill». «Paraphrasing consists in finding the meaning of two compared sequences and showing its equivalence, and this identification constitutes a judgment on the sequences»[s.~151]
- books.google house????iic: «The process of translation, particularly if successful, necessitates a complex text and discourse processing. The process of interpretation performed by the translator on the source text might lead to a TL text which is more redundant than the SL text. This argument may be stated as “the explicitation hypothesis”, [...] especially marked in the work of “non professional” translators» [s.~19–20]
- Hutchinson: «What is a grammatical sentence?» (vanskeleg å unngå *talaren* i akseptabilitetvurderingar); kva er ei frasesamanstilling, sånn ute i naturen?

### **3.21 TOGROK Retningslinjer for samanstilling :ROTE-TE:**

Ved korpusbygging er det vanleg at retningslinjer for samanstilling blir utvikla *etter kvart som ein finn problem...* (det er vanskeleg å seie noko *a priori* om kva for vanskar ein kan finne).

### **3.22 SKRIV Rangering**

## Kapittel 4

# Implementasjonen av `lfgalign`

For å finne ut av kor godt krava i forrige kapittel fungerer til å avgrense kva for lenkjer som er moglege, har eg implementert dei etter beste evne i eit Lisp<sup>1</sup>-program.

Ei implementering gjer det svært synleg om det finst manglar i eit formelt krav, eller om noko ikkje er godt nok spesifisert.

Programmet `lfgalign`<sup>2</sup> tek inn LFG-analysane av to setningar som me av uavhengige grunnar trur er omsetjingar av kvarandre. LFG-analysane må vere disambiguerte og i Prolog-formatet frå XLE<sup>3</sup>. Programmet les inn dei to filene og opprettar ein intern representasjon av LFG-analysen.

Me kan i tillegg gi programmet informasjon om kva for ord-omsetjingar me ser på som lingvistisk prediktable. Intensjonen er at dette kan vere informert av omsetjingstabellen frå eit automatisk ordsamanstillingsprogram, eller av handskrivne omsetjingsordbøker.

Programmet byrjar lenkjinga med f-strukturane. Ei f-struktursamanstilling er ei mengd med *lenkjer* mellom individuelle f-strukturar. Resultatet av lenkjinga på dette nivået kan vere tvitydig: sidan det ofte finst fleire måtar å lenkje argument og adjunkt på, får me i første omgang mange samanstillingar mellom kjelde- og mål-f-strukturar.

Difor rangerer me f-struktursamanstillingane, og den beste sender me vidare til c-struktursamanstillinga. Denne delen av programmet gir ut éi, utvitydig mengd med mange-til-mange-lenkjer mellom c-strukturane (her treng me ingen rangering). Nodane i kvar av desse mange-til-mange-lenkjene definerer no den endelege frasesamanstillinga.

<sup>1</sup>Dette språkvalet kan gjere eventuell integrering med andre LFG-system lettare (Common Lisp er m.a. nytta i LFG Parsebanker (Rosén et al., 2009)).

<sup>2</sup>Tilgjengeleg frå <http://github.com/unhammer/lfgalign> som fri og open programvare under GNU General Public License.

<sup>3</sup>Formatet er dokumentert på <http://www2.parc.com/isl/groups/nlt/xle/doc/xle.html>. Importeringa til Lisp-strukturar handterer «pakka representasjonar» og kjenner igjen ekvivalensforhold (t.d. der fleire  $\phi$ -variablar refererer til same f-struktur, eller fleire Prolog-variablar refererer til same analyseval); men filene eg har testa utnyttar ikkje det fulle spennet til formatet, så det finst ganske sikkert feil.

intro TODO,  
kanskje noko  
om kva eg  
faktisk har fått  
ut av imple-  
mentasjonen

treng eg ein  
eigen del om  
LPT i dette  
kapittelet?  
Implementa-  
sjonen er jo  
veldig enkel  
iallfall.

Nedanfor går eg gjennom detaljane rundt dei relevante delene av programmet.

## 4.1 Lenkjer mellom f-strukturar

Hovudalgoritmen for lenkjing mellom f-strukturar er vist i kodefigur 1. Funksjonen *f-align* returnerer ei mengd med moglege samanstillingar. Kvar samanstilling er ei mengd med par av f-strukturar<sup>4</sup>. Eit par  $(F_s, F_t)$  representerer ei lenkje frå ein f-struktur på kjeldespråket, til ein f-struktur på målspråket. Me føreset at dette paret har LPT-korrespondanse<sup>5</sup>, dette blir sjekka før alle kall på *f-align*. Der me ikkje har informasjon om LPT-korrespondanse mellom to ord (orda er ukjende), er lenkjing lov. Pro-element og substantiv kan alltid lenkjast med kvarandre.

Hjelpfunksjonen *argalign* (som igjen kallar *argalign-p*, vist i kodefigur 2) gir alle moglege «argumentpermutasjonar», dvs. moglege kombinasjonar av lenkjer mellom argumenta til  $F_s$  og  $F_t$  som tilfredsstiller kravet om LPT-korrespondanse, men utan å sjekke at desse argumenta igjen kan samanstillast. Funksjonen prøver å lenkje kvart argument til eit argument eller eit adjunkt, men gir ingen lenkjer mellom to adjunkt (sjå del 4.1.1 nedanfor om dette). Funksjonen gir heller ikkje kombinasjonar der minst eitt argument ikkje er lenkja – alle kombinasjonane må inkludere alle argument frå  $F_s$  og  $F_t$ , jf. krav (iii) og (iv) i Dyvik et al. (2009, s. 75). Elles er krav (i) er tautologisk oppfylt, medan me som nemnt føreset at krav (ii) er oppfylt før alle kall på *f-align*.

Eit døme: viss  $F_s$  har argumenta SUBJ og OBJ og ingen adjunkt, og  $F_t$  har argumentet SUBJ og eitt adjunkt ADJ, der alle ord-omsetjingar er moglege, vil *argalign* gi dei to samanstillingane  $\{(SUBJ, SUBJ), (OBJ, ADJ)\}$  og  $\{(SUBJ, ADJ), (OBJ, SUBJ)\}$ . Viss adjunktet til  $F_t$  ikkje fantest, eller ikkje hadde LPT-korrespondanse med nokon av argumenta til  $F_s$ , ville me ikkje fått nokon samanstillingar; medan viss paret (SUBJ, SUBJ) ikkje hadde LPT-korrespondanse og alt anna var likt, ville me berre fått den siste samanstillinga.

Funksjonen *f-align* går så gjennom kvar lenkje i kvar argumentpermutasjon, og prøver å kalle *f-align* på alle lenkjene. Sidan lenkjene som *argalign* gir har LPT-korrespondanse, vil alle f-strukturane i dei rekursive kalla i *f-align* ha LPT-korrespondanse. Eit rekursivt kall kan gi nye samanstillingar i dei indre f-strukturane, viss dei relevante krava er oppfylte.

Det er mogleg at ei lenkje frå éi samanstilling kan finnast i andre samanstillingar, me unngår dobbeltarbeid ved å lagre alle delvise samanstillingar i tabellen *alignable*. Dette føreset at *f-align(s, t)* er uavhengig av konteksten rundt; t.d. må mengda av samanstillingar som kjem ved å lenkje subjektet til  $F_s$  mot subjektet

<sup>4</sup>Eigentleg eit slag avgjerdstre; kvart element er eit par, kor første element er lenkja mellom dei yttarste f-strukturane, og andre element er dei moglege samanstillingane for dei indre strukturane. Denne strukturen kan vere nyttig for å rangere samanstillingar, og *f-align* blir mykje meir oversiktleg av å jobbe med eit slikt tre. Ein funksjon *flatten* omformar det ferdige treet til ei enkel liste med samanstillingar, kor kvar samanstilling er ei flat liste med lenkjer mellom f-strukturar.

<sup>5</sup>Når eg her skriv at to f-strukturar har LPT-korrespondanse, meiner eg sjølv sagt at ordformen til PRED-verdien til kvar f-struktur har LPT-korrespondanse.

til  $F_t$  vere uavhengig av om objektet til  $F_s$  er lenkja mot eit objekt eller eit adjunkt osb. av  $F_t$ .

```

alignments  $\leftarrow \emptyset$  ;
forall the argperm in argalign( $F_s$ ,  $F_t$ ) do
     $p \leftarrow \emptyset$  ;
    forall the  $A_s$ ,  $A_t$  in argperm do
        if not(aligntable[ $A_s$ ,  $A_t$ ]) then
            | aligntable[ $A_s$ ,  $A_t$ ]  $\leftarrow$  f-align( $A_s$ ,  $A_t$ );
        if aligntable[ $A_s$ ,  $A_t$ ] then
            | add aligntable[ $A_s$ ,  $A_t$ ] to  $p$ ;
        else
            | add ( $A_s$ ,  $A_t$ ) to  $p$ 
    end
    add  $p$  to alignments ;
    forall the adjperm in adjalign(argperm,  $F_s$ ,  $F_t$ ) do
         $a \leftarrow$  copy-of( $p$ ) ; // optional adjunct links
        forall the  $A_s$ ,  $A_t$  in adjperm do
            if not(aligntable[ $A_s$ ,  $A_t$ ]) then
                | aligntable[ $A_s$ ,  $A_t$ ]  $\leftarrow$  f-align( $A_s$ ,  $A_t$ );
            if aligntable[ $A_s$ ,  $A_t$ ] then
                | add aligntable[ $A_s$ ,  $A_t$ ] to  $a$ ;
            else
                | add ( $A_s$ ,  $A_t$ ) to  $a$ 
        end
        add  $a$  to alignments ;
    end
end
// loop through adjalign if no arguments exist
if alignments =  $\emptyset$  then return  $\emptyset$  ; // Fail
else return (( $F_s$ ,  $F_t$ ), alignments) ;

```

**Funksjon 1: f-align( $F_s$ ,  $F_t$ )**

TODO:  
nemne  
føresetnaden  
om  
uavhengnad i  
kapittel 3

Sjølvs om det er krav om LPT-korrespondanse mellom kvart argument og eit argument/adjunkt for å lenkje  $F_s$  og  $F_t$ , er det ikkje noko krav om at alle para i ein argumentpermutasjon tilfredsstiller alle lenkjingskrava. Viss f-align(OBJ, ADJ) frå dømet over gir null, og ikkje kan lenkjast (t.d. fordi ADJ hadde eitt argument, og OBJ ingen argument/adjunkt), medan f-align(SUBJ, SUBJ) kan lenkjast, vil f-align likevel returnere samanstillinga som inneheld (OBJ, ADJ) og (SUBJ, SUBJ). Me kan sjå i *aligntable* for å finne ut av om kvar av f-strukturane kunne lenkjast; i dette tilfellet vil *aligntable*[OBJ, ADJ] vere tom.

Om me i tillegg krev at substrukturar kan samanstillast kan me utelukke len-

TODO:  
forskjellen  
mellom  
LPT-krav og  
rekursjons-  
krav på  
argument må  
inn i kapittel 3

**usage:** Kalt av argalign slik:

argalign-p(arguments( $F_s$ ), adjuncts( $F_s$ ), arguments( $F_t$ ), adjuncts( $F_t$ ))

$a \leftarrow \emptyset$ ;

**if**  $args_s$  **then**

$s \in args_s$ ;

**forall the**  $t \in args_t$  **where**  $LPT(s,t)$  **do**

**forall the**  $p \in argalign-p(args_s - \{s\}, adj_s, args_t - \{t\}, adj_t)$  **do**

            add  $\{(s,t)\} \cup p$  to  $a$ ;

**end**

**forall the**  $t \in adj_t$  **where**  $LPT(s,t)$  **do**

**forall the**  $p \in argalign-p(args_s - \{s\}, adj_s, args_t, adj_t - \{t\})$  **do**

            add  $\{(s,t)\} \cup p$  to  $a$ ;

**end**

**return**  $a$ ;

**else if**  $args_t$  **then**

**if**  $adj_s$  **then**

$s \in adj_s$ ;

**forall the**  $t \in args_t$  **where**  $LPT(s,t)$  **do**

**forall the**  $p \in argalign-p(args_s, adj_s - \{s\}, args_t - \{t\}, adj_t)$

**do** add  $\{(s,t)\} \cup p$  to  $a$ ;

**end**

**return**  $a$ ;

**else**

**return**  $\emptyset$ ; // Fail

**else**

**return**  $\{\emptyset\}$ ; // End

**Funksjon 2:**  $argalign-p(args_s, adj_s, args_t, adj_t)$

kjing av f-strukturane  $F_s$  og  $F_t$  i (1) under:

- (1) a. 
$$F_s \left[ \begin{array}{l} \text{PRED 'planlegge<eg,[1:gi]>'} \\ \text{XCOMP}_1 \left[ \text{PRED 'gi (opp)'} \right] \end{array} \right]$$
- b. 
$$F_t \left[ \begin{array}{l} \text{PRED 'plan<I,[2:give]>'} \\ \text{XCOMP}_2 \left[ \text{PRED 'give<I,him,it>'} \right] \end{array} \right]$$

Men det kan vere at me ikkje *vil* krevje dette i alle moglege tilfelle. Ei tryggare løysing er å rangere ulike løysingar i etterkant, ved å spørje etter dei argumentsamanstillingane som har flest medlem i *aligntable*, dette kjem eg tilbake til i 4.2 nedanfor.

#### 4.1.1 Overflødige adverbial

Argumentpermutasjonane frå *argalign* prøver som nemnt ikkje reine adjunkt-adjunkt-lenkjer, sidan me ikkje vil forkaste lenkjing av  $F_s$  og  $F_t$  berre på grunn av at ikkje alle adjunkt kunne lenkjast. Men når me har prøvd ein argumentpermutasjon, kan me lage ein kopi av denne som i tillegg inneheld lenkjer mellom «overflødige» adverbial, altså dei adjunkt-adjunkt-lenkjene som *argalign* ikkje prøver. Hjelpefunksjonen *adjoin* (ikkje vist her **TODO: implementere :->**) konstruerer moglege permutasjonar av lenkjer mellom adjunkt som ikkje er inkludert i *argperm*, og *f-align* prøver desse rekursivt på same måte som med argumentlenkjene. Lenkjene blir lagt til ein *kopi* av argumentpermutasjonane, sidan det ikkje er sikkert at me ønskjer å lenkje alle adjunktdøtre. Viss me har to overflødige adjunkt på kvar side, og kravet om LPT-korrespondanse er dekkja for alle fire moglege par, får me seks moglege permutasjonar, sidan me inkluderer dei fire permutasjonane der eitt adjunktpar er ulenkja.

og så er det spørsmålet om me kan lenkje adjunkt på ulike nivå i f-strukturane

Viss  $F_s$  og  $F_t$  ikkje hadde argument i det heile teke, går me au gjennom moglege permutasjonar av adjunktdøtre, på same måte (ikkje vist i kodefigur 1).

#### 4.1.2 SKRIV Når f-lenkjene ikkje er 1-1

**kausativ**

**preposisjonsobjekt**

“sigaretten” og “sigaretze” er ikkje på same nivå i dei respektive f-strukturane nedanfor:

```
0[ PRED vedde<28,29,27,30>
  29[ PRED sigarett<> ] ]

0[ PRED da-najleveba<37,10,46>
  ADJUNCT { 2 }
```

Dette må 1. spesifiserast (kap.3), og 2. implementerast...

```
2[ ze<5>
  OBJ 5[ sigareti ] ] ]
```

Men dette er fordi sigaretze er adjunkt, og i følge fotnote 3 i Dyvik et al. (2009) er P(ADJ) lik den semantiske formen til preposisjonsobjektet til ADJ. (Er dette heilt enkelt? Kan me ha adjunkt med fleire «preposisjonsargument»?)

I motsetning til å endre på P(ADJ) slik at pred-verdien til eit adjunkt hoppar over preposisjonar og inn i preposisjonsobjekt, har eg i implementasjonen min gjort det slik at funksjonen som henter ut *f-strukturen* til adjunkt-døtre av ein *f-struktur* hoppar over preposisjonsobjekt (m.a. fordi ved eit kall P(X) kan ikkje P vite om X er eit adjunkt eller argument).

Sjå au del 3.13.

#### notat :ROTETE:

filene

```
((tab_s (open-and-import "dev/TEST_argadj_s.pl"))
 (tab_t (open-and-import "dev/TEST_argadj_t.pl")))
```

viser at me kan trenge samanføyning av pred på ulike nivå.

#### 4.1.3 Kan me gjere f-struktursamanstillinga bottom-up?

Ein alternativ metode for lenkjing av *f-strukturane* er å byrje med alle logisk moglege permutasjonar av LPT-korrespondansar, og så sile ut dei som ikkje svarer til krava. Prosessen ville nok blitt mykje meir oversiktleg på denne måten, sidan det då berre er snakk om å sjekke krav for kvar enkelt lenkje. Men ein slik metode er vanskeleg i praksis; når avskjeringa skjer så seint, blir det alt for mange moglege kombinasjonar for lengre setningar med mange ukjende ord til at ein vanleg datamaskin kan halde styr på dei.

Me må i alle tilfelle vere klar for ei setning der alle ord er ukjende (me har ingen informasjon om LPT-korrespondanse), slik at kvart kjeldeord kan lenkjast til kvart målord. Viss begge setningane er 4 ord, får me 16 moglege samanstillingar der alle ord er med i nøyaktig éi lenkje ( $2^l$ , kor  $l$  er setningslengd). Men ofte har me null-lenkjer, me må altså i tillegg tillate samanstillingar der minst eitt ord er ulenkja, utan at me treng å vite kva for ord det er; med desse kortare listene inkludert får me endå fleire moglege samanstillingar per setning (4 ord gir 26, 8 ord gir 2186 moglege samanstillingar). Sjølv om me heile tida vel dei samanstillingane som lenkjar flest ord, ville maskinen raskt fått problem. I tillegg har me problemet med 1-mange-lenkjer, som skaper endå fleire moglege samanstillingar.

Ein sideverknad av å byrje med ytre lenkjer og gå innover (prosessen skildra i del 4.1) er at me automatisk unngår å prøve «kryssande» lenkjer, t.d. å lenkje  $F_s$  med XCOMP av  $F_t$ , og XCOMP av  $F_s$  med  $F_t$  (denne kombinasjonen av lenkjer vil jo vere ein del av alle logisk moglege permutasjonar). Me får au prioritert å lenkje



ytre element, som jo er sikrere lenkjer: gitt to f-strukturar for setningar der alt me veit om lenkjinga er at *setningane* er omsetjingar av kvarandre, vil dei to ytre f-strukturane ha størst sjanse for å korrespondere med kvarandre. For kvart steg du går innover må du multiplisere inn sjansen for å trå feil i argumentpermutasjonane.

## 4.2 SKRIV Rangering

Rangering er ikkje implementert enno (akkurat no gir *rank(f-alignments)* berre ut første samanstilling.) Rangering foregår etter ulike kriterium. Her er eit par forslag:

### 4.2.1 rekursivt lenkja > ulenkja, men LPT

*alignable* seier om noko er rekursivt lenkja eller ikkje, plusspoeng viss me har klart å lenkje rekursivt.

### 4.2.2 argument-argument > argument-adjunkt

Plusspoeng for argument-argument-lenkjer, burde vere eit bra kriterium, men me får sjølvsagt problem viss LPT ikkje seier noko i døme (8) med  
 da-najleveba<Abrams,Browne,regne> adjunkt: sigarett  
 bet<Abrams,sigarett,regne> adjunkt: Browne

### 4.2.3 arg1-arg1 arg2-arg2 > arg1-arg2 arg2-arg1 (følge)

Dette kjem til å gi problem når me vil lenkje «behage» og «like», viss me ikkje har motstridande LPT-informasjon (og argumentfølge i leksikon ikkje er basert på semantikk, men syntaks). Men elles er det vel OK.

Enklaste implementasjon: Levenshtein-avstand. Men burde visse argument vek-  
 tast? (T.d. vekte subjekt om alt anna er likt.)

Andre forslag: [http://en.wikipedia.org/wiki/Edit\\_distance](http://en.wikipedia.org/wiki/Edit_distance)

### 4.2.4 flest lenkja adjunkt

Usikker på dette... avheng av om me tillèt lenkjer på tvers av f-strukturar.

### 4.2.5 Prioritet på rangeringskriterium

Dette bør sjølvsagt testast empirisk, blir kanskje utanfor denne oppgåva (diskusjonsdel?), men kan jo prøve meg litt rundt.

### 4.3 Lenkjing av c-strukturnodar

Samanstilling mellom f-strukturar treng i *lfgalign* ikkje informasjon om c-strukturen, medan lenkjing av c-strukturnodar skjer på grunnlag av f-struktursamanstillinga. Programmet utfører difor samanstilling av c-strukturar sist.

Funksjonen *c-align* har som inndata c-strukturanalysane av kjelde- og målsetninga, og éi f-struktursamanstilling; utdata er ei mengd med lenkjer. Ei lenkje er eit par der første element er ei mengd c-strukturnodar på kjeldespråket, og andre element ei mengd nodar på målspråket. Det er ingen overlapp mellom medlem av lenkjer (ein node er aldri med i meir enn eitt par).

I Dyvik et al. (2009, s. 77) er kravet for å lenkje to c-strukturnodar er at dei dominerer same mengd med ordlenkjer<sup>6</sup>. Ein node *n* dominerer ei mengd lenkjer *l* viss unionen av lenkjene dominert av døtrene til *n* er lik *l*. I *lfgalign* opererer eg ikkje med *ordlenkjer* i seg sjølv; f-struktursamanstillinga er basert på LPT-korrespondansar, som definerer moglege ordlenkjer utan å sjå på kontekst, og f-struktursamanstillinga avgrensar vidare moglege ordlenkjer gitt f-strukturinformasjon. Preterminale nodar er dei mest ordnære nodane som kan ha ei f-strukturlenkje (ved  $\emptyset$ ); når formålet er å lenkje c-strukturnodar kan me nytte f-strukturlenkja til den preterminale noden i staden for ordlenkjer.

```

c-alignments  $\leftarrow \emptyset$  ;
splitss  $\leftarrow$  new table ;
add-links(f-alignment, trees, splitss) ;
splitst  $\leftarrow$  new table ;
add-links(f-alignment, treet, splitst) ;
forall the links being the keys in splitss do
    if (links in splitst) then
        | add (splitss[links], splitst[links]) to c-alignments ;
    end
return c-alignments ;

```

**Funksjon 3:** *c-align*(f-alignment, tree<sub>s</sub>, tree<sub>t</sub>)

Hjelpeprosedyren *add-links* utfører hovudjobben. Inndata er rotnoden til c-strukturtrreet for eitt av språka, og f-samanstillinga. Prosedyren kappar opp treet i nodemengder, kor kvar nodemengd dominerer same lenkjemengd (som definert over). Nodemengdene blir lagra i ein tabell, indeksert på lenkjemengdene. Prosedyren går rekursivt gjennom treet frå rot til lauv; lenkjemengden for kvar node er unionen av lenkjemengdene returnert av *add-links* kalt på kvar av døtrene. Viss ein node dominerer ei lenkjemengd *links*, legg me til denne noden i tabellen *splits*[*links*].

To problem  
(kva vil me ha med?)  
1. me får ikkje med LPT-korrespondansar som er OK, men ikkje med i *f-alignment*;  
2. me får med LPT-korrespondansar som er med i *f-alignment* men ikkje *alignable* (ikkje er rekursivt lenkja).

<sup>6</sup>Dette er ein litt enklare måte å definere kravet på; ei *lenkje* refererer til både kjelde og mål, dimed blir det mogleg å seie at ein node på kjeldespråket kan dominere same mengd som ein node på målspråket.

```

links ← ∅;
if node then
  if preterminal?(node) then
    let link ∈ f-alignment s.t.  $\phi(\text{node}) \in \text{link}$  ;
    if link then links ← {link}
  else
    links ← add-links(f-alignment, left-branch(node)) ∪
    add-links(f-alignment, right-branch(node)) ;
  add node to splits[links] ;
return links ;
Funksjon 4: add-links(f-alignment, node, splits)

```

Sidan c-align kallar add-links for kvar av sidene, får me to tabellar  $splits_s$  og  $splits_t$ . Me hentar så ut alle dei lenkjemengdene som er i begge tabellane (dvs. snittet av oppslagsnøklerne til tabellen); nodane som er lagra med mengd med f-strukturlenkjer skal lenkjast på c-strukturnivå. Alle desse mange-til-mange-lenkjene blir til slutt returnert av c-align.

Proessen er no ferdig, mange-til-mange-lenkjene mellom c-strukturnodar definerer frasesamanstillinga

*til diskusjonsdel: Det er ikkje berre ei N-gramsamanstilling; sidan lenkjene er mellom c-strukturnodar kor kvar node dominerer ein konstituent, kunne me kalt det ei konstituentsamanstilling..*

### 4.3.1 SKRIV viss me har LPT, men ikkje rekursiv f-lenkje

Dette bør kanskje vere valfritt i programmet, for å sjå kva det fører til: vil du ta med LPT-korrespondansar som ikkje har f-lenkjer i add-links?

Og omvendt, finst det LPT-korrespondansar som ikkje kjem med i f-alignment i det heile teke, men som likevel burde ha noko å seie for c-strukturlenkjinga? (Men burde dei ikkje då vere med i f-alignment au?)

Kan me fjerne visse f-samanstillingar pga. c-strukturinfo? dvs. disambiguere... (dette er vel heller stoff for diskusjonsdelen?)

## Kapittel 5

# Diskusjon, resultat av å automatisk samanstille norske og georgiske setningar

- om kjeldematerialet
- manglar med implementasjonen
- samanlikning av lenkjing basert på f-struktur og lenkjing basert på N-gram
- bruksområde for samanstillingar

### 5.1 Oppdage argumentstrukturalternasjon

I døme TODO i kapittel TODO viste eg at f-strukturar og LPT-korrespondanse kan-  
skje ikkje har nok informasjon til å kunne handtere ulik følgje i argumentstruktur. ref  
Programmet mitt vil her gi begge løysingar, ei rangering basert på lik argumentføl-  
gje vil gi feil løysing på topp.

Kanskje kan me nytte data frå fleire førekomstar med andre subjekt og objekt  
til å lære slike argumentstrukturalternasjonar. Om me observerer *sie gefällt mir/jeg  
liker henne* vil me jo ha f-strukturinformasjon som kan nyttast til å informere ar-  
gumentstrukturalternasjon (*sie/henne* er hokjønn, etc.), om det var substantiv der  
ville LPT-korrespondanse kunne informere dette.

### 5.2 Samanlikning med tremetodar og n-grammetodar

I tillegg til at ein kanskje kan få betre skåre på kvantitative mål som presisjon  
og gjenkjenning, vil lenkjer mellom f-strukturar gi informasjon som er kvalitativt  
forskjellig frå det ein kan få med å berre sjå på lenkjer mellom ord, n-gram eller  
konstituentar.

### 5.2.1 c->f er mange-til-ein

Avbildinga frå c-strukturknodar til f-struktur er mange-til-ein, kan me t.d. innanfor eitt tre ha fleire N-gram per f-strukturhovud; ein metode som berre ser på enkle N-gramlenkjer vil ikkje registrere desse relasjonane (t.d. metoden i Samuelsson & Volk (2007)).

### 5.2.2 TOGROK men korleis gjenfinne *there is/es gibt*? :ROTETE:

1. N-gram kjem like ofte som heile konstruksjonen, då kan dette gjenfinnast
  - dvs., *there is NP/es gibt NP*-samanstilling kjem like ofte som *there is* eller *es gibt* førekjem. Eit TigerXML-type søk etter *there is NP/es gibt NP* burde jo vere mogleg, sjekk om dette er delmengd av *there is/es gibt*. \* Avslutning

## **Kapittel 6**

## **Avslutning**

# Referansar

- Alsina, A., Bresnan, J. & Sells, P. (red.). (1997). *Complex predicates*. Stanford, CA, USA: Center for the Study of Language and Information. Paperback.
- Aronson, H. (1990). *Georgian. A Reading Grammar. Corrected Edition*. Columbus, OH: Slavica Publishers.
- Bresnan, J. (2001). *Lexical-Functional Syntax*. Oxford, UK: Blackwell Publishers. Tilgjengeleg frå <http://books.google.com/books?id=7elu0CcxQWkC> (ISBN: 0631209743)
- Brown, P.F., Della Pietra, S.A., Della Pietra, V.J. & Mercer, R.L. (1993). The mathematics of statistical machine translation: Parameter estimation. *Computational Linguistics*, 19(2), 263–311. Tilgjengeleg frå <http://citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.1.13.8919>
- Butt, M., Dyvik, H., King, T., Masuichi, H. & Rohrer, C. (2002). The Parallel Grammar Project. I *COLING-02 on Grammar engineering and evaluation* (vol. 15, s. 1–7). Morristown, NJ: Association for Computational Linguistics. Tilgjengeleg frå <http://portal.acm.org/citation.cfm?id=1118783.1118786>
- Cheung, L., Lai, T., Luk, R., Kwong, O., Sin, K., Tsou, B. et al. (2002). Some Considerations on Guidelines for Bilingual Alignment and Terminology Extraction. , 1–5. Tilgjengeleg frå <http://www.aclweb.org/anthology-new//W/W02/W02-1802.pdf>
- Cyrus, L., Feddes, H. & Schumacher, F. (2004). Annotating predicate-argument structure for a parallel treebank. *LISBON, 2004*, 39. Tilgjengeleg frå <http://arxiv.org/abs/cs/0407002>
- Dyvik, H., Meurer, P., Rosén, V. & Smedt, K.D. (2009). Linguistically motivated parallel parsebanks. I M. Passarotti, A. Przepiórkowski, S. Raynaud & F.V. Eynde (red.), *Proceedings of the eighth international workshop on treebanks and linguistic theories* (s. 71–82). Milan, Italy: EDUCatt. Tilgjengeleg frå [http://tlt8.unicatt.it/allegati/Proceedings\\_TLT8.pdf#page=83](http://tlt8.unicatt.it/allegati/Proceedings_TLT8.pdf#page=83)

- Giegerich, H. (2006). Attribution in English and the distinction between phrases and compounds'. *Englisch in Zeit und Raum-English in Time and Space: Forschungsbericht für Klaus Faiss. Trier: Wissenschaftlicher Verlag Trier*. Tilgjengeleg frå <http://www.englang.ed.ac.uk/people/attributioninenglish.pdf>
- Hearne, M., Ozdowska, S. & Tinsley, J. (2008). Comparing Constituency and Dependency Representations for SMT Phrase-Extraction. I *Actes de la 15e Conférence Annuelle sur le Traitement Automatique des Langues Naturelles (TALN '08)*. Avignon, France. Tilgjengeleg frå <http://www.computing.dcu.ie/~mhearne/publications.html>
- Koehn, P., Och, F. & Marcu, D. (2003). Statistical phrase-based translation. I *NAACL '03: Proceedings of the 2003 Conference of the North American Chapter of the Association for Computational Linguistics on Human Language Technology* (s. 48–54). Morristown, NJ, USA. Tilgjengeleg frå <http://www.iccs.inf.ed.ac.uk/~pkoehn/publications/phrase2003.pdf>
- Kruijff-Korbyova, I., Chvatalova, K. & Postolache, O. (2006). Annotation guidelines for Czech-English word alignment. , 1256–1261. Tilgjengeleg frå <http://www.mt-archive.info/LREC-2006-Kruijff.pdf>
- Maier, R.M. (2009). *Structural Interference from the Source Language: A psycholinguistic investigation of syntactic processes in non-professional translation*. Upublisert akademisk avhandling, University of Edinburgh. Tilgjengeleg frå <http://linguistlist.org/issues/20/20-1786.html>
- Meurer, P. (2008, March). *A Computational Grammar for Georgian*. Tilgjengeleg frå <http://maximos.aksis.uib.no/~paul/articles/Tbilisi2007-LNAI.pdf>
- Munday, J. (2001). *Introducing Translation Studies: Theories and Applications*. London: Routledge.
- Piao, S. & McEnery, T. (2001). Multi-word Unit Alignment in English-Chinese Parallel Corpora. I P. Rayson, A. Wilson, T. McEnery, A. Hardie & S. Khoja (red.), *Proceedings of the Corpus Linguistics 2001 Conference* (s. 466–475). Lancaster, UK. Tilgjengeleg frå [http://personalpages.manchester.ac.uk/staff/scott.piao/research/papers/mwu\\_align4.pdf](http://personalpages.manchester.ac.uk/staff/scott.piao/research/papers/mwu_align4.pdf)
- Pullum, G. & Scholz, B. (2001). On the Distinction between Model-Theoretic and Generative-Enumerative Syntactic Frameworks. *Logical Aspects of Computational Linguistics: 4th International Conference, Lacl 2001, Le Croisic, France, June 27-29, 2001, Proceedings*. Tilgjengeleg frå <http://portal.acm.org/citation.cfm?id=645668.665062>



- Riezler, S. & Maxwell, J. (2006). Grammatical Machine Translation. I M. Butt, M. Dalrymple & T.H. King (red.), *Intelligent Linguistic Architecture: Variations on themes by Ronald M. Kaplan* (s. 35–52). Stanford, CA: CSLI Publications. Tilgjengeleg frå <http://www.parc.com/research/publications/details.php?id=5675>
- Rosén, V., Meurer, P. & Smedt, K. de. (2009). LFG Parsebanker: A Toolkit for Building and Searching a Treebank as a Parsed Corpus. I F.V. Eynde, A. Frank, G. van Noord & K.D. Smedt (red.), *Proceedings of the 7th International Workshop on Treebanks and Linguistic Theories (TLT7)* (s. 127–133). Utrecht: LOT. Tilgjengeleg frå <http://ling.uib.no/~desmedt/papers/tlt7rosen-submitted.pdf>
- Samuelsson, Y. & Volk, M. (2006). Phrase Alignment in Parallel Treebanks. I *Proceedings of Treebanks and Linguistic Theories (TLT '06)*. Prague. Tilgjengeleg frå [http://ling16.ling.su.se:8080/new\\_PubDB/doc\\_repository/229\\_align.pdf](http://ling16.ling.su.se:8080/new_PubDB/doc_repository/229_align.pdf)
- Samuelsson, Y. & Volk, M. (2007). Automatic Phrase Alignment: Using Statistical N-Gram Alignment for Syntactic Phrase Alignment. I *Proceedings of Treebanks and Linguistic Theories (TLT '07)*. Bergen, Norway.
- Thunes, M. (2003). *Ekserpering av leksikalske oversettelsekorrespondanser fra parallelltekst*. Tilgjengeleg frå <http://www.hf.uib.no/i/LiLi/SLF/ans/Dyvik/marthaex.pdf>
- Tinsley, J., Hearne, M. & Way, A. (2007). Exploiting Parallel Treebanks to Improve Phrase-Based Statistical Machine Translation. I *Proceedings of Treebanks and Linguistic Theories (TLT '07)*. Bergen, Norway.
- XPar. (2008). *XPAR: Language diversity and parallel grammars*. (Submitted to the Research Council of Norway.)