
What's in a Lexicon ?

Cem Bozşahin

Computer Engineering

Middle East Technical University (METU), Ankara

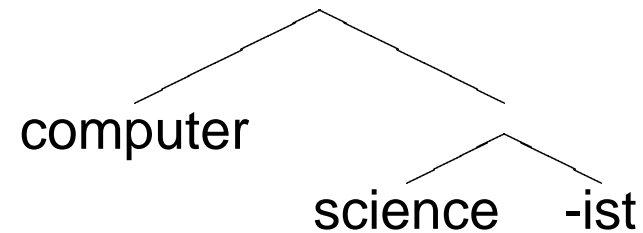
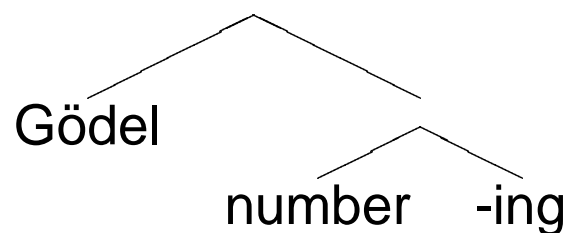
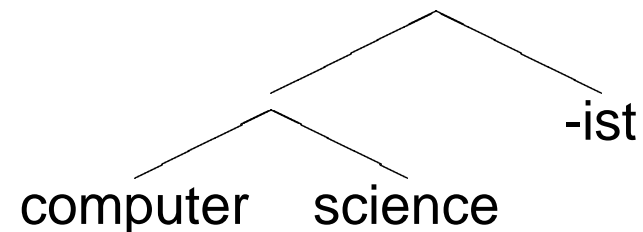
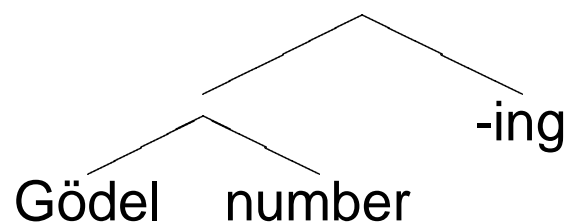
March 17, 2004

Overview

- Examples of bracketing mismatches and phrasal scope of inflections
 - Architectures for morphology-syntax-semantics interface
 - Morphosyntax: with words or morphemes?
 - Morphosyntactic types
 - Lexical representation of free/bound morphemes
 - Sample derivations of the parser (and performance)
-

Derivational morphology

- **bracketing mismatches** were first noted in derivational morphology (Williams, 1981)



Verbal inflection

- The problem arises in inflectional morphology as well
- [West Greenlandic](#) (Fortescue, 1984)

Aatsaat tikeraa-nngi-laq
for.first.time visit-NEG-INDIC/3s

'It is not the first time he has visited.'

- It does not mean 'This is the first time he failed to visit.'
-

Coordination

- German (Müller, 1999)

Wenn **Ihr Lust** und noch nichts anderes **vor-habt**,
if you pleasure and yet nothing else intend

können wir sie ja vom Flughafen abholen
can we them PARTICLE from.the airport pick up

'If you feel like it and have nothing else planned, we can pick them up at the airport.'

- **semantics:** Ihr Lust habt UND noch nichts anderes vorhabt
-

Subordination

- Turkish

Mehmet Ayşe'nin [düzenli uyu]-ma-ma-sı-na kız-ıyor
 M.NOM A.-GEN regularly sleep-NEG-INF-AGR-DAT anger-TENSE
 'Mehmet is angry with Ayşe for not sleeping regularly.'
 not 'Mehmet is constantly angry with Ayşe for not sleeping.'

Mehmet Ayşe'nin kitab-ı oku-ma-sı-nı iste-di
 M.NOM A.-GEN book-ACC read-INF-AGR-ACC want-TENSE
 'Mehmet wanted Ayşe to read the book.'

- semantics: want (read book ayşe) mehmet

Relativisation

- Turkish (Bozsahin, 2002)
- Local and non-local morphosyntactic requirements of rel. noun may be different

Ben Mehmet'in çocuğ-a/*-u ver-diğ-i kitab-ı oku-du-m
 I.NOM M-GEN child-DAT/*ACC give-REL.OP book-ACC read-TENSE-PERS1
 'I read the book that Mehmet gave to the child.'

Ben Mehmet'in kitab-ı ver-diğ-i çocuğ-u/*-a gör-dü-m
 I.NOM M-GEN book-ACC give-REL.OP child-ACC/*DAT see-TENSE-PERS1
 'I saw the child to whom Mehmet gave the book.'

Lexemic vs. morphemic lexicons

ver-diğ-i :=

LOCAL	CAT	HEAD <table style="border-collapse: collapse; margin-left: 10px;"> <tr> <td style="padding-right: 5px;">AGR</td> <td style="border-left: 1px solid black; border-right: 1px solid black; padding: 2px 5px;">PERSON <i>third</i></td> </tr> <tr> <td style="padding-right: 5px;">CASE</td> <td style="border-left: 1px solid black; border-right: 1px solid black; padding: 2px 5px;">NUMBER <i>sing</i></td> </tr> </table>	AGR	PERSON <i>third</i>	CASE	NUMBER <i>sing</i>
AGR	PERSON <i>third</i>					
CASE	NUMBER <i>sing</i>					
CONTENT	SUBCAT < 3 NP[gen], 2 NP[acc], 1 NP[dat] > MOD MODSYN LOCAL CONT INDEX 1	RELN <i>give</i> GIVER 3 GIVEE 1 GIFT 2				
NONLOCAL TO-BIND SLASH { 1 }						

-diġ-i :=

$$\left[\begin{array}{l} \text{LOCAL} \left[\begin{array}{l} \text{CAT} \left[\begin{array}{l} \text{HEAD } \textit{noun} [\text{acc or dat}] \\ \text{SUBCAT } \langle \rangle \end{array} \right] \\ \text{CONTENT } \textit{npro} [\text{INDEX } \boxed{1}] \end{array} \right] \\ \text{NONLOCAL | INHER | SLASH } \{ \boxed{1} \} \end{array} \right]$$

Nominal inflection

- Morphological richness of the language does not seem to be the issue
- English (Carpenter, 1997)

four truck-s

semantics: four (plu truck)

alleged thiev-es

semantics: plu (alleged thief), not alleged (plu thief)

Ki-relativisation

- a. araba-da-ki
 car-LOC-REL
 'the one in the car'
- b. [[çocuğ-un ev-i-nde]-ki-ler-in]-ki
 child-GEN house-POSS-LOC2-REL-PLU-GEN-REL
 lit. 'The one that belongs to the ones that are in the child's house'
- c. Ben ev-de-ki-ni hiç kullan-ma-dı -m
 I.NOM house-LOC-REL-ACC2 never use-NEG-TENSE-PERS.1s
 'I never used the one at home.'
- d. [k'üç 'ük ev-de]-ki hediye
 house-LOC-REL present
 'the present_i, the one_i at the little house'
 not 'the little present at the house.'

reduplication

- kırmızı kıpkırmızı yeşil yemyeşil
 - yeni yepyeni/yesyeni açık apaçık
 - Semantics is uniform, yet there seems to be no morphemic representation.
 - It is actually phonology at work on a lexical item (similar processes apply in Tagalog)
 - Lexical rules are *unary* rules; they can create new lexical items, and semantics of that can be associated with the lexical rule (specified only once)
-

-
- NB: lexical rules refer to *substantive* categories (adjective, adverb etc.), whereas combinatory rules use *formal* categories (*X*, *Y* etc.)
-

Resolving the mismatch

- **semantics** may require affixes to have scope larger than the inflected word
 - Alternatives for the morphology-syntax-semantics interface
 - Autonomous levels of morphology, syntax, and semantics (e.g. Sadock, 1991)
 - Morphosyntax-driven semantics (Heylen, 1997; Bozsahin, 2002)
 - The lexicon can be **morphemic** in either case, but it is a **combinatory morphemic lexicon** in a more lexicalist approach
-

Inflectional morphology & linguistic theory

- GB (Anderson, 1982) and LFG (Bresnan, 1995) consider inflectional morphology to be part of syntax, (in GB, it is not part of **combinatory** aspects of grammar)
 - MP (Chomsky, 1995) assumes words enter syntax fully inflected (numeration)
 - HPSG (Pollard & Sag, 1994) keeps it in the lexicon (lexical rules, or lexical inheritance hierarchy)
 - CG work in general (Hoffman, 1995; Heylen, 1997; and others) assumes word-based lexicons, although this is not a theoretical commitment
-

TLG and inflectional morphology

- Heylen's (1997, 1999) unary modalities. $\text{Frau} := \square_{\text{case}} \square_{\text{fem}} \square_{\text{sg}} \square_{\text{3p}} \square_{\text{decl}} \mathcal{N}$
- Morphosyntactic type assignment is to inflected forms
- Structural rules regulate scope of inflections, e.g. $\square_{\text{sg}} \square_{\text{case}} \mathcal{N}$ can be turned into $\square_{\text{case}} \square_{\text{sg}} \mathcal{N}$ by a structural rule
- some iterative morphological processes challenge the lexical rules for word-based type assignment (e.g. **-ki** in Turkish)
- A more lexical solution is to have morphemic lexicons and morphosyntactic calculus (i.e. **-ki** as lexical item)

Lexical Syntactic types

- syntactic categories and features

N, NP, S

feature-decorations, *NP_{acc}, S_{fin}*

- But features as such are not part of combinatorics,

unlike e.g. *NP_{case} → Det N Case*

Syntactic calculus

Application (<): $Y:a \quad X \backslash Y:f \Rightarrow X:fa$

Composition (>B): $X/Y:f \quad Y/Z:g \Rightarrow X/Z:\lambda x.f(gx)$

Type Raising (>T): $X:a \Rightarrow T/(T \backslash X):\lambda f.f[a]$

Leftward Contraposition (<XP): $X:a \Rightarrow$
 $S_{+t}/(S/X):\lambda f.f[a]$
 $S_{+t}/(S_{+t}/X):\lambda f.f[a]$

Rightward Contraposition (>XP): $X:a \Rightarrow$
 $S_{-t} \backslash (S \backslash X):\lambda f.f[a]$
 $S_{-t} \backslash (S_{-t} \backslash X):\lambda f.f[a]$

Lexical Morphosyntactic types

- Two kinds of unary modalities on syntactic types

$\triangle^a X$ (flexible morphosyntactic domain for X : “up to certain inflectional type”)

$\boxtimes^a X$ (strict domain : “require certain inflectional type”)

- if inflectional paradigm is Stem-Number-Case,

$\boxtimes^c N$ stands for case-marked nouns

$\triangle^c N$ stands for noun stems , number-marked, and case-marked nouns

- Lattice $L = (\mathcal{D}, \leq, =)$
- The set of basic morphosyntactic types: \mathcal{A}_{ms}
 - $\triangleleft^i X \in \mathcal{A}_{ms}$ and $\bowtie^i X \in \mathcal{A}_{ms}$ if $i \in \mathcal{D}$ and $X \in \mathcal{A}_s$ (\mathcal{A}_s : syntactic types)
- The set of complex morphosyntactic types: \mathcal{B}_{ms}

$$\mathcal{A}_{ms} \subseteq \mathcal{B}_{ms}$$

If $X \in \mathcal{B}_{ms}$ and $Y \in \mathcal{B}_{ms}$, then $X \setminus Y$ and $X / Y \in \mathcal{B}_{ms}$

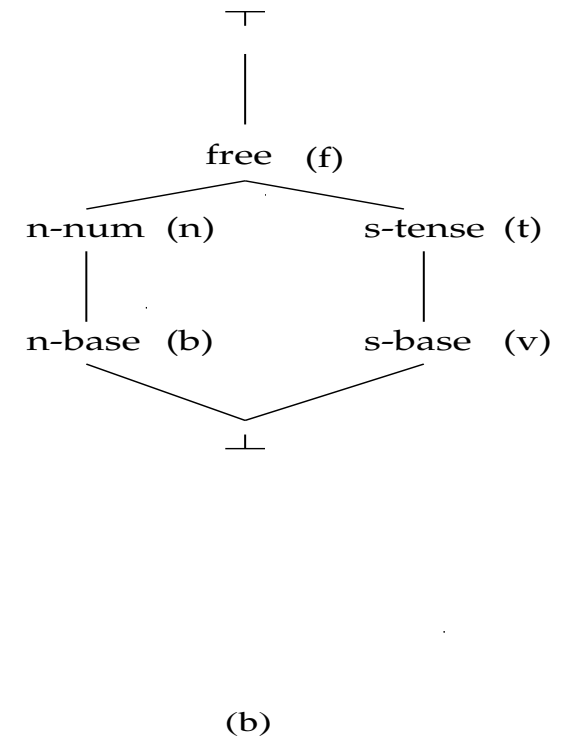
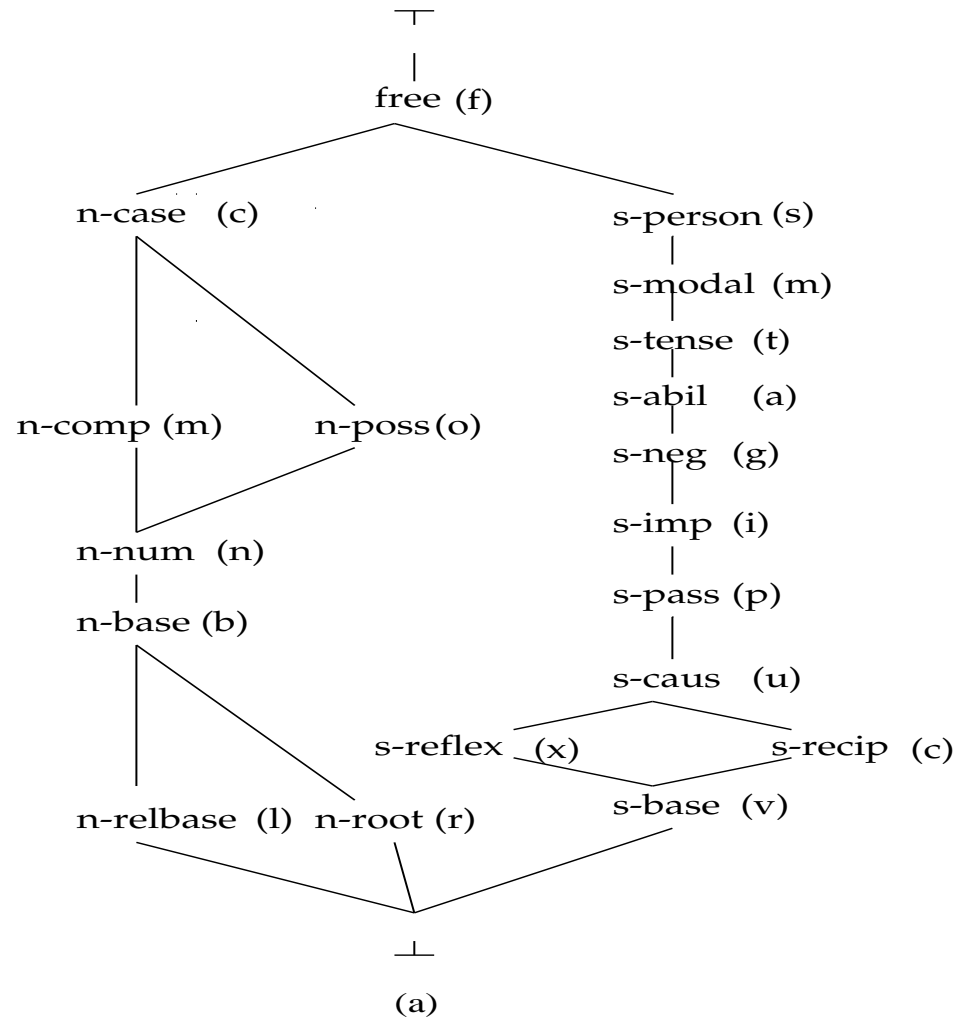
Lattice of diacritics (inflectional types)

- Inclusion of domains is specified in a language-particular lattice

This comes in handy for specifying morphotactics as well

- More importantly, it allows morphosyntactic types to pick semantic domains independent of surface attachment
- All of this is specified in the lexical entry

attachment type, morphosyntactic type, diacritic, semantic type



Morphosyntactic lexicon & grammar

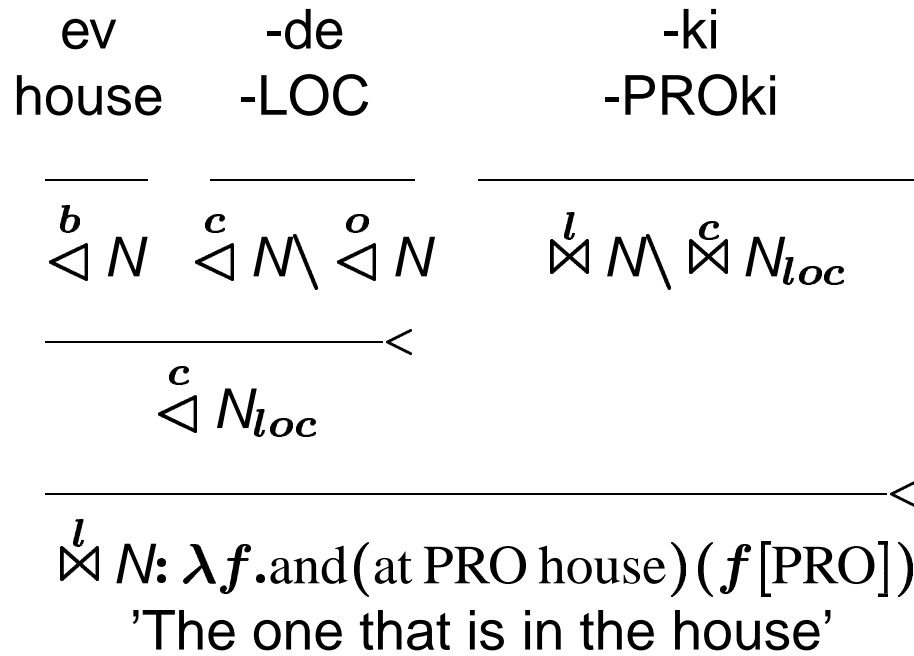
• -PLU := $\overset{a}{\circ} s - \overset{n}{\triangleleft} M \setminus \overset{b}{\triangleleft} N: \lambda x. \text{plu } x$

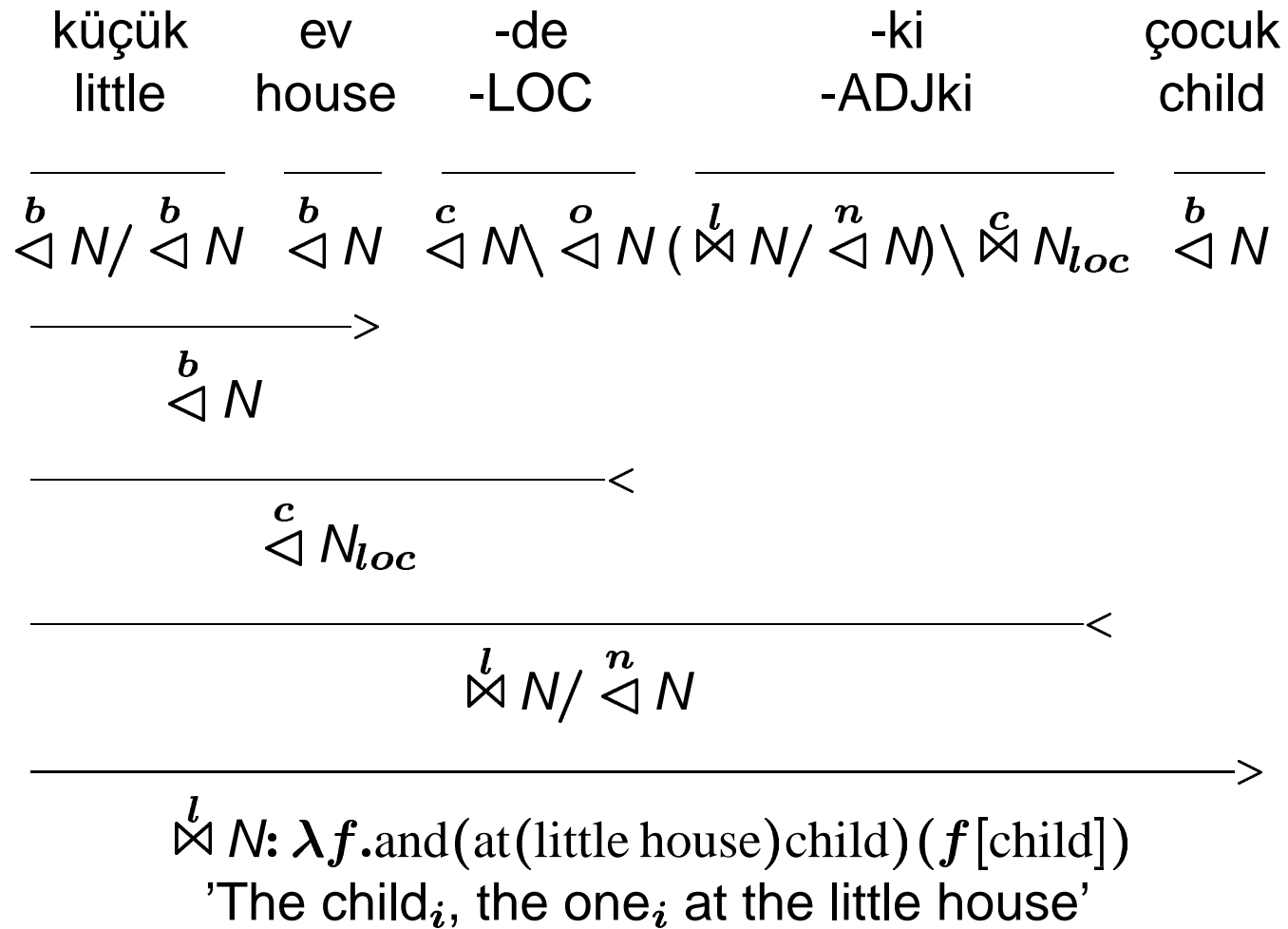
• Forward Application ($>$):

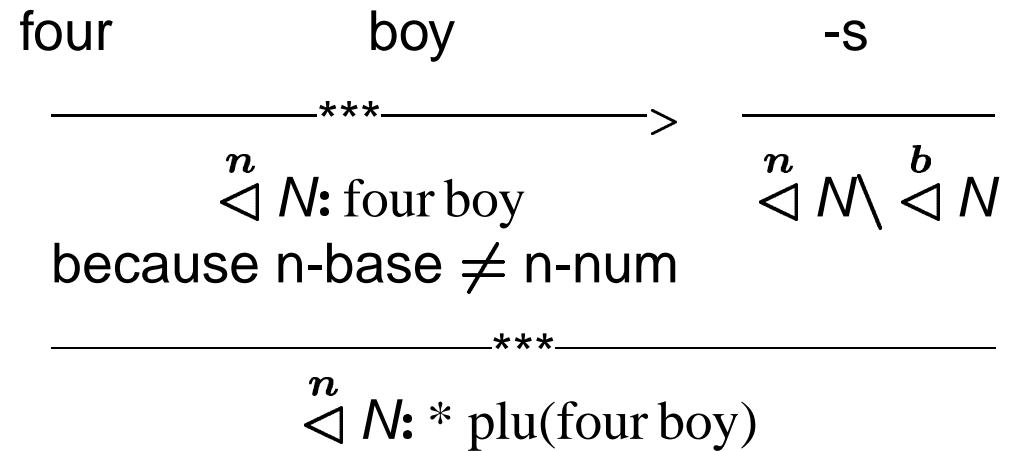
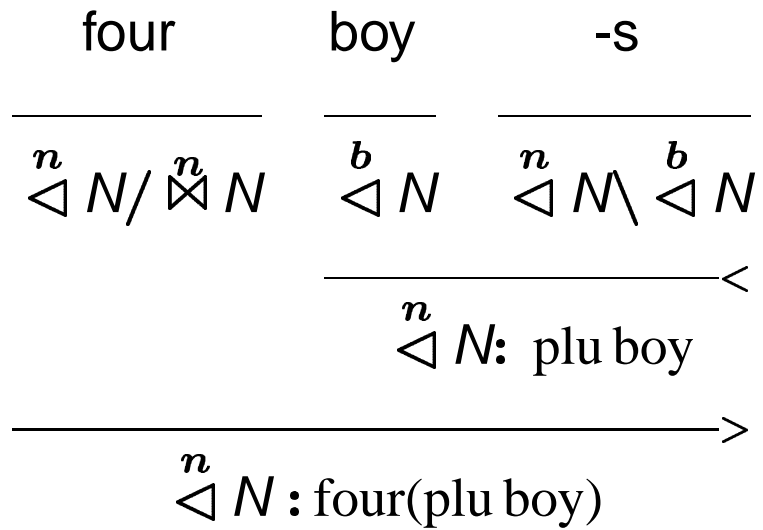
$$\frac{\overset{i}{\circ} s_1 - X / \overset{\alpha_1}{\square_1} Y: f \quad \overset{j}{\circ} s_2 - \overset{\alpha_2}{\square_2} Y: a}{\overset{k}{\circ} (s_1 \bullet s_2) - X: fa} \rightarrow$$

if $\alpha_2 \square_1 \alpha_1$ in lattice L , for:

$\square_1, \square_2 \in \{\bowtie, \triangleleft\}$,
 $\alpha_1, \alpha_2 \in \mathcal{D}$ in L ,
 $i, j, k \in \{a, s, c\}$,
 $\overset{i}{\circ} \overset{j}{\circ} \vdash_a \overset{k}{\circ}$







toy gun -s

$\overset{b}{\triangleleft} N / \overset{b}{\triangleleft} N$ $\overset{n}{\triangleleft} N : \text{plu gun}$ <

$\overset{n}{\triangleleft} N : * \text{toy(plu gun)}$
 because n-num $\not\leq$ n-base

toy gun -s

$\overset{b}{\triangleleft} N / \overset{b}{\triangleleft} N$ $\overset{b}{\triangleleft} N$ $\overset{n}{\triangleleft} N \setminus \overset{b}{\triangleleft} N$

>

$\overset{b}{\triangleleft} N : \text{toy gun}$

<

$\overset{n}{\triangleleft} N : \text{plu(toy gun)}$

Aatsaat
for.the.first.time

tikeraa
visit

-nngi
-NEG

-laq
-INDIC

$(\overset{i}{\triangleleft} S \setminus \overset{f}{\triangleleft} NP) / (\overset{i}{\triangleleft} \setminus \overset{f}{\triangleleft} NP) \quad \overset{v}{\triangleleft} S \setminus \overset{f}{\triangleleft} NP \quad (\overset{n}{\triangleleft} S \setminus \overset{f}{\triangleleft} NP) \setminus (\overset{i}{\triangleleft} S \setminus \overset{f}{\triangleleft} NP)$

$\overset{i}{\triangleleft} S \setminus \overset{f}{\triangleleft} NP$

$\overset{n}{\triangleleft} S \setminus \overset{f}{\triangleleft} NP$

'This is not the first time he visited.'

Aatsaat
for.the.first.time

tikeraa
visit

-nngi
-NEG

$$(\overset{i}{\triangleleft} S \setminus \overset{f}{\triangleleft} NP) / (\overset{i}{\triangleleft} S \setminus \overset{f}{\triangleleft} NP) \quad \overset{v}{\triangleleft} S \setminus \overset{f}{\triangleleft} NP \quad (\overset{n}{\triangleleft} S \setminus \overset{f}{\triangleleft} NP) \setminus (\overset{i}{\triangleleft} S \setminus \overset{f}{\triangleleft} NP)$$

$$\overset{n}{\triangleleft} S \setminus \overset{f}{\triangleleft} NP$$

because $n \not\leq i$

Çocuk child.NOM	kız-a girl-DAT	kalem-i pen-ACC	ver give	-me -SUB1i	-yi -ACC	unut-tu forgot
—>T	—>T	—>T	—————	—————	—————	—<B
↑ N_{nom}	↑ N_{dat}	↑ N_{acc}	DV	$b \triangleleft N \setminus (\triangleleft S \setminus \triangleleft NP_{nom})$	$c \triangleleft N_{acc} \setminus \triangleleft N$	TV
: $\lambda f.f$ [child]	: $\lambda g.g$ [girl]	: $\lambda h.h$ [pen]	: $\lambda x.\lambda y.\lambda z.$ give yxz	: $\lambda f.f$: $\lambda f.f$: $\lambda f.\lambda x.$ forget (f [ana x]) x

—————>
 $v \triangleleft S \setminus \triangleleft NP_{nom} \setminus \triangleleft NP_{dat}$

—————>
 $v \triangleleft S \setminus \triangleleft NP_{nom}$

—————<
 $b \triangleleft N$

—————<
 $c \triangleleft N_{acc}$

—————>T
 $(S \setminus NP) / (S \setminus NP \setminus \triangleleft NP_{acc})$

—————>
 $t \triangleleft S \setminus \triangleleft NP_{nom}$

—————>
 $t \triangleleft S$: forget(give girl pen(ana child))child
 'The child forgot to give the pen to the girl.'

çocuğ-un kitab-ı ver -diği adam uyu-du
 child-GEN book-ACC give -OP.AGR man sleep-TENSE

$\overline{\text{NP}}_{agr} <$ $\overline{\text{TV/DV}} >^T$ $\overline{\text{DV}}$ $\overline{(N^\uparrow / N) \setminus IV_{agr}}$ \overline{N} $\overline{IV} <^B$

$\overline{\text{TV}} >$

$\overline{IV_{agr}} <$

$\overline{N^\uparrow / N} <$

$\overline{N^\uparrow = S \setminus IV} >$

$\overline{\hspace{15em}} >$

S: and(sleep man) (give man book child)

'The man to whom the child gave the book slept.'

completely

destroy

-ed

$$\overline{(\overset{v}{\triangleleft} S \setminus \overset{f}{\triangleleft} NP) / (\overset{v}{\triangleleft} S \setminus \overset{f}{\triangleleft} NP)} \quad \overline{(\overset{v}{\triangleleft} S \setminus \overset{f}{\triangleleft} NP) / \overset{f}{\triangleleft} NP} \quad \overline{(\overset{t}{\triangleleft} S \setminus \overset{f}{\triangleleft} NP) / (\overset{v}{\triangleleft} S \setminus \overset{f}{\triangleleft} NP)}$$

→B

$$(\overset{v}{\triangleleft} S \setminus \overset{f}{\triangleleft} NP) / \overset{f}{\triangleleft} NP$$
←B_x

$$(\overset{t}{\triangleleft} S \setminus \overset{f}{\triangleleft} NP) / \overset{f}{\triangleleft} NP$$

*completely

did

destroy

$$\overline{(\overset{v}{\triangleleft} S \setminus \overset{f}{\triangleleft} NP) / (\overset{v}{\triangleleft} S \setminus \overset{f}{\triangleleft} NP)} \quad \overline{(\overset{t}{\triangleleft} S \setminus \overset{f}{\triangleleft} NP) / (\overset{v}{\triangleleft} S \setminus \overset{f}{\triangleleft} NP)} \quad \overline{(\overset{v}{\triangleleft} S \setminus \overset{f}{\triangleleft} NP) / \overset{f}{\triangleleft} NP}$$

→B

$$(\overset{t}{\triangleleft} S \setminus \overset{f}{\triangleleft} NP) / \overset{f}{\triangleleft} NP$$

because $t \not\leq v$

did

destroy

completely

$$(\overset{t}{\triangleleft} S \setminus \overset{f}{\triangleleft} NP) / (\overset{v}{\triangleleft} S \setminus \overset{f}{\triangleleft} NP) \quad (\overset{v}{\triangleleft} S \setminus \overset{f}{\triangleleft} NP) / \overset{f}{\triangleleft} NP \quad (\overset{v}{\triangleleft} S \setminus \overset{f}{\triangleleft} NP) \setminus (\overset{v}{\triangleleft} S \setminus \overset{f}{\triangleleft} NP)$$

$$(\overset{v}{\triangleleft} S \setminus \overset{f}{\triangleleft} NP) / \overset{f}{\triangleleft} NP$$
 $\leftarrow B_x$ $\rightarrow B$

$$(\overset{t}{\triangleleft} S \setminus \overset{f}{\triangleleft} NP) / \overset{f}{\triangleleft} NP$$

destroy

-ed

completely

$$(\overset{v}{\triangleleft} S \setminus \overset{f}{\triangleleft} NP) / \overset{f}{\triangleleft} NP \quad (\overset{t}{\triangleleft} S \setminus \overset{f}{\triangleleft} NP) / (\overset{v}{\triangleleft} S \setminus \overset{f}{\triangleleft} NP) \quad (\overset{v}{\triangleleft} S \setminus \overset{f}{\triangleleft} NP) \setminus (\overset{v}{\triangleleft} S \setminus \overset{f}{\triangleleft} NP)$$

$$(\overset{t}{\triangleleft} S \setminus \overset{f}{\triangleleft} NP) / \overset{f}{\triangleleft} NP$$
 $\leftarrow B_x$ $\leftarrow B_x$

???

Experiments with the CKY parser

- a 21-morpheme sentence (12 words) parsed in 2.9 seconds

37-morphemes (20 words) in 40 seconds

- Güngördü & Oflazer's LFG parser takes 10 seconds/sentence with 24,000 word lexicon
 - separate morphological analyzers deliver 2 to 5 analyses/second (Oflazer, 1996; Komagata, 1997)
 - 2.8 morphemes/word on the average including derivations (Turkish)
less than 2 inflections/word (Oflazer et. al, 2001)
-

Sample text type	Number of items in text			Avg. number of parses/gram. input		Avg. CPU time per test (milliseconds)		
	tests	words	morphs	PAS	NF	Unrestr.	PAS	NF
				check	parse		check	parse
Word order & case	58	216	384	1.26	3.68	39	39	30
Subordination	14	70	137	3.00	5.09	267	270	180
Relativisation	23	130	232	2.04	2.32	796	783	266
Control verbs	33	147	291	1.42	3.34	166	163	137
Possessives & compounds	26	109	200	1.23	2.47	137	135	98
Adjuncts	14	57	100	1.12	4.87	89	88	72
-ki relatives	24	66	179	1.07	1.54	36	36	35

Conclusion

- The key to integration of inflectional morphology and syntax is granting representational status to morphemes
 - Morphosyntactic mismatches do not necessitate multi-tiered grammars
 - Lexical items can be smaller or larger than words, and project their own semantic domains and attachment characteristics
 - Loss of efficiency is tolerable up to medium-length sentences
 - Modular grammar-lexicon (in fact, just the lexicon!)
-