





# Developing a Russian HPSG based on the Russian National Corpus

---

DELPH-IN Summit 2009, Barcelona

# The SlaviGraM Project

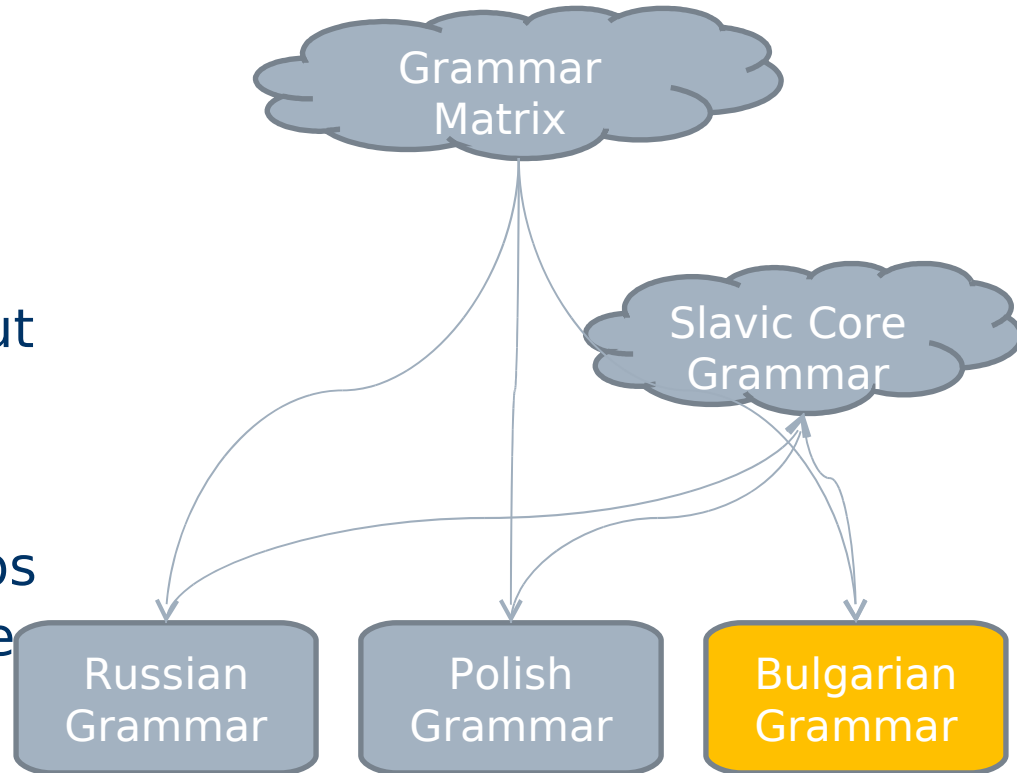
- Proof-of-Concept Implementation
  - strategies for constructing a cross-linguistic resource based on concepts of shared and non-shared grammar in modelling Slavic morphosyntax
- Focus on Russian Resource Grammar
  - construction of Russian HPSG in the DELPH-IN environment combining the Grammar Matrix and a Slavic Core Grammar with corpus-based grammar elaboration; exploiting the Russian National Corpus (RNC) in Russian Resource Grammar (RRG) engineering.
- Cross-Slavic Validation (two showcases)
  - construction of Bulgarian HPSG in the DELPH-IN environment combining the Grammar Matrix and a Slavic Core Grammar with corpus-based grammar elaboration exploiting the BulTreeBank
  - construction of Polish HPSG in the DELPH-IN environment combining the Grammar Matrix and a Slavic Core Grammar with corpus-based grammar elaboration exploiting the IPPI PAN Corpus

# Idea

- construction of Russian HPSG
- in the DELPH-IN environment
- combining the Grammar Matrix
- and a Slavic Core Grammar
- with corpus-based grammar elaboration
  - exploiting the Russian National Corpus (RNC) as structured grammatical knowledge resource

# Top-down vs. Bottom-up

- Use Grammar Matrix to quickly build small grammars for individual languages
- Shared analyses from individual languages are put into the Slavic Core
- When new language is added, the Slavic Core helps to more efficiently build the new grammar, and potentially receives cross-slavic validation

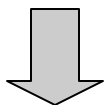


# Integrating External Morphological Analyzers

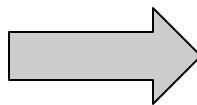


- Mystem (<http://company.yandex.ru/technology/mystem/>)
- Perl-based configurable SPPP wrapper

профессор{профессор=S,муж,од=им,ед}  
читает{читать=V,несов=непрош,ед,изъяв,3-л}  
книгу{книга=S,жен,неод=вин,ед}



S+им := noun-nom-irule.  
S+вин := noun-acc-irule.  
S+ед := noun-sg-irule.  
V+непрош := verb-nonpast-irule.  
V+изъяв := verb-indic-irule.  
V+3-л := verb-3p-irule.



```
<?xml version="1.0" encoding="utf-8"?>
<segment>
  <token form="профессор" from="0" to="9">
    <analysis stem="профессор">
      <rule id="noun-nom-irule" form="профессор"/>
      <rule id="noun-sg-irule" form="профессор"/>
    </analysis>
  </token>
  <token form="читает" from="10" to="16">
    <analysis stem="читать">
      <rule id="verb-nonpast-irule" form="читает"/>
      <rule id="verb-indic-irule" form="читает"/>
      <rule id="verb-3p-irule" form="читает"/>
    </analysis>
  </token>
  <token form="книгу" from="17" to="22">
    <analysis stem="книга">
      <rule id="noun-acc-irule" form="книгу"/>
      <rule id="noun-sg-irule" form="книгу"/>
    </analysis>
  </token>
</segment>
```

# Morphological information in RNC

<b>Part of speech</b> <input type="checkbox"/> noun <input type="checkbox"/> adjective <input type="checkbox"/> numeral <input type="checkbox"/> numeral adjective <input type="checkbox"/> verb <input type="checkbox"/> adverb <input type="checkbox"/> predicative <input type="checkbox"/> parenthesis <input type="checkbox"/> pronoun <input type="checkbox"/> adjective pronoun <input type="checkbox"/> predicative pronoun <input type="checkbox"/> adverbial pronoun <input type="checkbox"/> preposition <input type="checkbox"/> conjunction <input type="checkbox"/> particle <input type="checkbox"/> interjection	<b>Case</b> <input type="checkbox"/> nominative <input type="checkbox"/> vocative* <input type="checkbox"/> genitive <input type="checkbox"/> genitive 2 <input type="checkbox"/> dative <input type="checkbox"/> accusative <input type="checkbox"/> accusative 2* <input type="checkbox"/> instrumental <input type="checkbox"/> locative <input type="checkbox"/> locative 2 <input type="checkbox"/> adnumerative	<b>Mood / Verb form</b> <input type="checkbox"/> indicative <input type="checkbox"/> imperative <input type="checkbox"/> imperative 2 <input type="checkbox"/> infinitive <input type="checkbox"/> participle <input type="checkbox"/> gerund	<b>Degree / Adj. form</b> <input type="checkbox"/> comparative <input type="checkbox"/> comparative 2* <input type="checkbox"/> superlative <input type="checkbox"/> full form <input type="checkbox"/> short form
	<b>Number</b> <input type="checkbox"/> singular <input type="checkbox"/> plural	<b>Tense</b> <input type="checkbox"/> present <input type="checkbox"/> future <input type="checkbox"/> past	<b>Transitivity</b> <input type="checkbox"/> transitive* <input type="checkbox"/> intransitive*
<b>Antroponymic</b> <input type="checkbox"/> family name <input type="checkbox"/> first name <input type="checkbox"/> patronymic	<b>Gender</b> <input type="checkbox"/> masculine <input type="checkbox"/> feminine <input type="checkbox"/> neuter <input type="checkbox"/> common*	<b>Voice</b> <input type="checkbox"/> active <input type="checkbox"/> passive <input type="checkbox"/> middle	
	<b>Animacy</b> <input type="checkbox"/> animate <input type="checkbox"/> inanimate	<b>Aspect</b> <input type="checkbox"/> perfective <input type="checkbox"/> imperfective	

OK Clear Cancel

\* - only in the corpus with resolved homonymy

\*\* - only in the corpus with unresolved homonymy

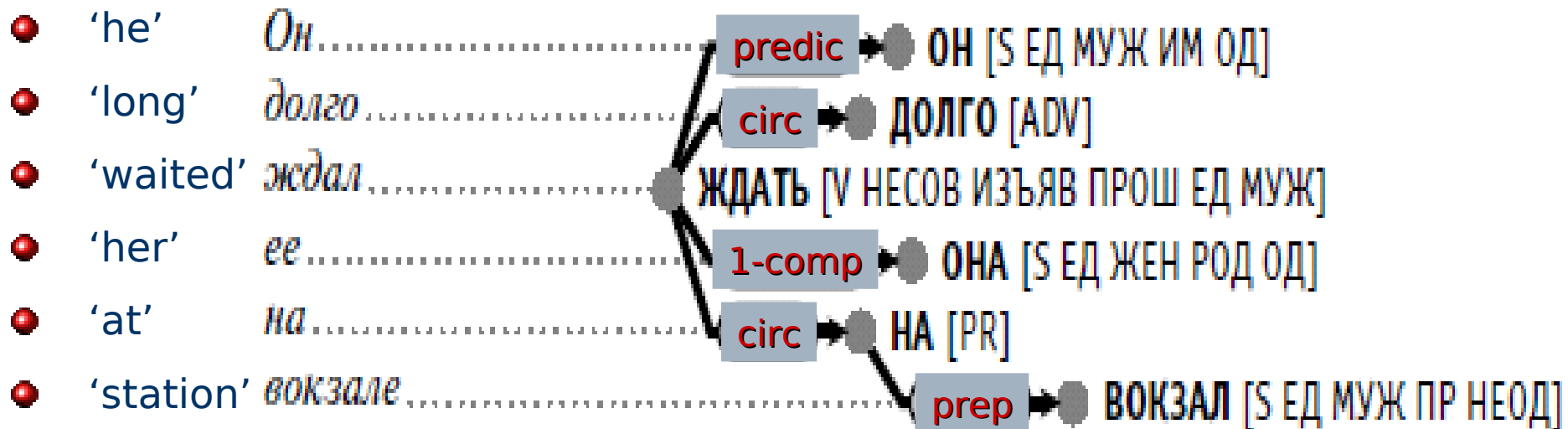
# Making the most of the syntactic Information in RNC



- In the deeply annotated sub-corpus of the RNC, every sentence is marked up with a dependency syntactic structure, with nodes corresponding to the words of the sentence, and labeled edges encoding the syntactic relations.
- The syntactic formalism originates in the Meaning-Text Theory (Mel'cuk 1995), but the inventory of syntactic relations has been extended for the purposes of corpus annotation, incorporating a number of specific linguistic decisions (Boguslavsky et al. 2002; Apresjan et al. 2006).
- Unlike the morphologically annotated portion of the RNC, the deeply annotated sub-corpus only contains fully disambiguated annotations (i.e. both morphological and syntactic ambiguity is resolved).



# A sample structure



# Grammatical features used in the syntactic sub-corpus

<b>Part of speech</b> <input type="checkbox"/> nominal <input type="checkbox"/> adjective <input type="checkbox"/> numeral <input type="checkbox"/> verb <input type="checkbox"/> adverb <input type="checkbox"/> preposition <input type="checkbox"/> conjunction <input type="checkbox"/> particle <input type="checkbox"/> interjection <input type="checkbox"/> compound word <input type="checkbox"/> word-sentence <input type="checkbox"/> foreign word, non-lexical formula	<b>Case</b> <input type="checkbox"/> nominative <input type="checkbox"/> genitive <input type="checkbox"/> partitive <input type="checkbox"/> dative <input type="checkbox"/> accusative <input type="checkbox"/> instrumental <input type="checkbox"/> prepositive <input type="checkbox"/> locative <input type="checkbox"/> vocative	<b>Aspect</b> <input type="checkbox"/> perfective <input type="checkbox"/> imperfective
	<b>Grade</b> <input type="checkbox"/> comparative <input type="checkbox"/> comparative 2 <input type="checkbox"/> superlative	<b>Tense</b> <input type="checkbox"/> present <input type="checkbox"/> non-past <input type="checkbox"/> past
<b>Animacy</b> <input type="checkbox"/> animate <input type="checkbox"/> inanimate	<b>Form</b> <input type="checkbox"/> short form	
<b>Gender</b> <input type="checkbox"/> masculine <input type="checkbox"/> feminine <input type="checkbox"/> neuter	<b>Representation</b> <input type="checkbox"/> finite verb <input type="checkbox"/> infinitive <input type="checkbox"/> participle <input type="checkbox"/> gerund	<b>Voice</b> <input type="checkbox"/> passive
<b>Number</b> <input type="checkbox"/> singular <input type="checkbox"/> plural	<b>Mood</b> <input type="checkbox"/> indicative <input type="checkbox"/> imperative	<b>Other</b> <input type="checkbox"/> part of a compound word

# Access to the syntactic sub-corpus



[Main corpus](#) [Syntactic corpus](#) [Spoken corpus](#)

[русская версия](#)

## Search by exact form <sup>?</sup> A B B

Word or phrase

search

clear

## Lexico-grammatical search <sup>?</sup>

Word <sup>?</sup> A B B

Gramm. features <sup>?</sup> [select](#)



Distance from parent: from  to  <sup>?</sup>

Syntactic relationship

<sup>?</sup> [select](#)

Word <sup>?</sup> A B B

Gramm. features <sup>?</sup> [select](#)



search

clear

# Syntactic relations in RNC

<p><b>Actantial relationships</b></p> <ul style="list-style-type: none"> <li><input type="checkbox"/> predicative</li> <li><input type="checkbox"/> dative subjective</li> <li><input type="checkbox"/> agentive</li> <li><input type="checkbox"/> quasi-agentive</li> <li><input type="checkbox"/> non-intrinsic agentive</li> <li><input type="checkbox"/> I completive</li> <li><input type="checkbox"/> II completive</li> <li><input type="checkbox"/> III completive</li> <li><input type="checkbox"/> IV completive</li> <li><input type="checkbox"/> V completive</li> <li><input type="checkbox"/> copula</li> <li><input type="checkbox"/> I non-intrinsic completive</li> <li><input type="checkbox"/> II non-intrinsic completive</li> <li><input type="checkbox"/> III non-intrinsic completive</li> <li><input type="checkbox"/> non-actantial completive</li> <li><input type="checkbox"/> completive appositive</li> <li><input type="checkbox"/> prepositional</li> <li><input type="checkbox"/> subordinating conjunctive</li> <li><input type="checkbox"/> comparative</li> <li><input type="checkbox"/> comparative conjunctive</li> <li><input type="checkbox"/> elective</li> </ul>	<p><b>Attributive</b></p> <p><b>determinative</b></p> <ul style="list-style-type: none"> <li><input type="checkbox"/> (proper) determinative</li> <li><input type="checkbox"/> descriptive determinative</li> <li><input type="checkbox"/> approximative ordinal</li> <li><input type="checkbox"/> relative</li> </ul> <p><b>General attributive</b></p> <ul style="list-style-type: none"> <li><input type="checkbox"/> (proper) attributive</li> <li><input type="checkbox"/> compound</li> </ul> <p><b>appositive</b></p> <ul style="list-style-type: none"> <li><input type="checkbox"/> (proper) appositive</li> <li><input type="checkbox"/> dangling appositive</li> <li><input type="checkbox"/> nominative appositive</li> <li><input type="checkbox"/> numerative appositive</li> </ul> <p><b>quantitative</b></p> <ul style="list-style-type: none"> <li><input type="checkbox"/> (proper) quantitative</li> <li><input type="checkbox"/> approximative quantitative</li> <li><input type="checkbox"/> approximative co-predicative</li> <li><input type="checkbox"/> approximative delimitative</li> <li><input type="checkbox"/> distributive</li> <li><input type="checkbox"/> additive</li> </ul> <p><b>circumstantial</b></p> <ul style="list-style-type: none"> <li><input type="checkbox"/> (proper) circumstantial</li> <li><input type="checkbox"/> durative</li> <li><input type="checkbox"/> multiple durative</li> <li><input type="checkbox"/> distantional</li> <li><input type="checkbox"/> circumstantial tautological</li> <li><input type="checkbox"/> subjective circumstantial</li> <li><input type="checkbox"/> objective circumstantial</li> <li><input type="checkbox"/> subjective co-predicative</li> <li><input type="checkbox"/> objective co-predicative</li> <li><input type="checkbox"/> delimitative</li> <li><input type="checkbox"/> parenthetical</li> <li><input type="checkbox"/> complement clause</li> <li><input type="checkbox"/> expository</li> <li><input type="checkbox"/> adjunctive</li> <li><input type="checkbox"/> precisising</li> </ul>	<p><b>Coordinative</b></p> <ul style="list-style-type: none"> <li><input type="checkbox"/> coordinative</li> <li><input type="checkbox"/> sentential coordinative</li> <li><input type="checkbox"/> conjunctive coordinative</li> <li><input type="checkbox"/> communicative coordinative</li> <li><input type="checkbox"/> multiple</li> </ul> <p><b>Syncategorematic</b></p> <ul style="list-style-type: none"> <li><input type="checkbox"/> analytical</li> <li><input type="checkbox"/> passive analytical</li> <li><input type="checkbox"/> auxiliary</li> <li><input type="checkbox"/> quantitative auxiliary</li> <li><input type="checkbox"/> correlative</li> <li><input type="checkbox"/> expletive</li> <li><input type="checkbox"/> proleptic</li> <li><input type="checkbox"/> elliptic</li> </ul>
---	---	---

# Mapping the components of a RNC dependency relation to HPSG categories

- Observe the following straightforward convention:
  - in a given syntactic dependency relation,
  - the “governor X” corresponds to the lexical head of the head daughter, while
  - the “dependent Y” corresponds to the lexical head of the non-head daughter.

# E.g. predicative syntactic relation predicate(X) – first argument(Y)



## Head-Subject Schema (trivial cases)

Профессор[*Y*] читает[*X*].  
professor·NOMINATIVE read·FINITE.ACTIVE

'The professor reads'

Заявка[*Y*] изучается[*X*].  
proposal·NOMINATIVE study·FINITE.REFLEXIVE.PASSIVE

'The proposal is being studied.'

vs.

Комитет[*Y*] изучает[*X*] заявку  
committee·NOMINATIVE study·FINITE.ACTIVE proposal·ACCUSATIVE

'The committee studies the proposal.'

# E.g. predicative syntactic relation predicate(X) – first argument(Y)



## Head-Subject Schema (non-verbal predication)

● noun(X): Москва[Y] — столица[X] России.  
Moscow·NOMINATIVE capital·NOMINATIVE Russia·GENITIVE  
'Moscow is the capital of Russia.'

● adjective(X): Профессор[Y] какой-то странный[X].  
professor·NOMINATIVE somewhat strange·FULL-ADJECTIVE.NOMINATIVE  
'The professor is kind of strange.'

Профессор[Y] очень добр[X].  
professor·NOMINATIVE very kind·SHORT-ADJECTIVE  
'The professor is very kind.'

Профессор[Y] должен[X] уходить.  
professor·NOMINATIVE obliged·PREDICATIVE-ADJECTIVE leave·INFINITIVE  
'The professor must leave.'

# 1.1. predicative syntactic relation

## predicate(X) – first argument(Y)



### Head-Subject Schema (non-nominative subject)

- Y is a noun in genitive or partitive (also known as “second genitive”)

Хлеба[Y]      не      осталось[X]  
bread.<sub>GENITIVE</sub>    NEG    remain.<sub>IMPERSONAL.REFLEXIVE.PAST</sub>

‘No bread left.’

Сахару[Y]      хватит[X]      на всех  
sugar.<sub>PARTITIVE</sub>    suffice.<sub>IMPERSONAL.NON-PAST</sub>    on all

‘There will be enough sugar for everybody.’

- Y is a distributive or approximative PP

Пришло[X]      до[Y]    десяти    детей.  
come.<sub>IMPERSONAL.PAST</sub>    to    ten.<sub>GENITIVE</sub>    children

‘Up to ten children came.’

Нам      досталось[X]      по [Y]    груше.  
we.<sub>DATIVE</sub>    recieve.<sub>IMPERSONAL.REFLEXIVE.PAST</sub>    of    pear.<sub>DATIVE</sub>

‘We’ve got a pear each.’



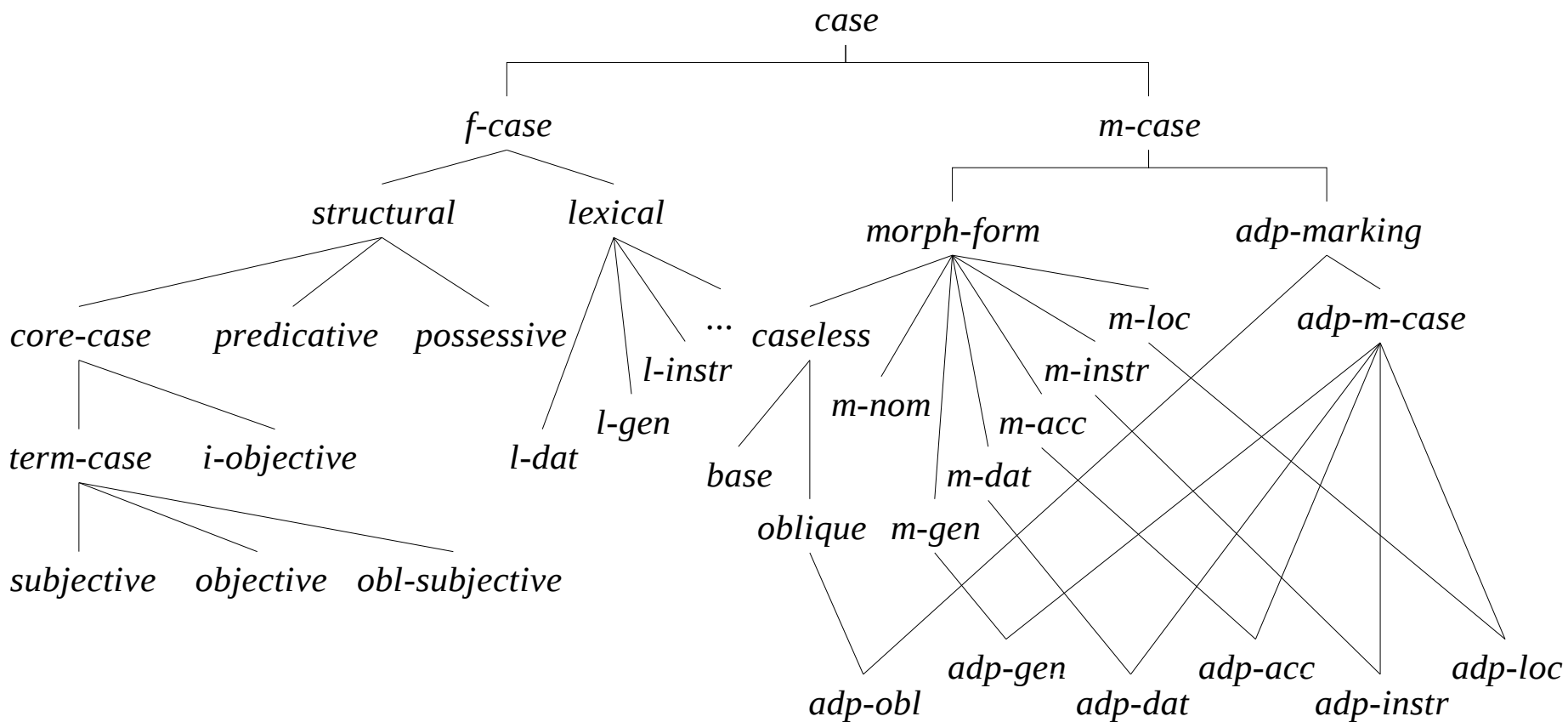
# More Than Just a Russian Grammar

## E.g. Shared Slavic Case Hierarchy



functional dimension

marking dimension

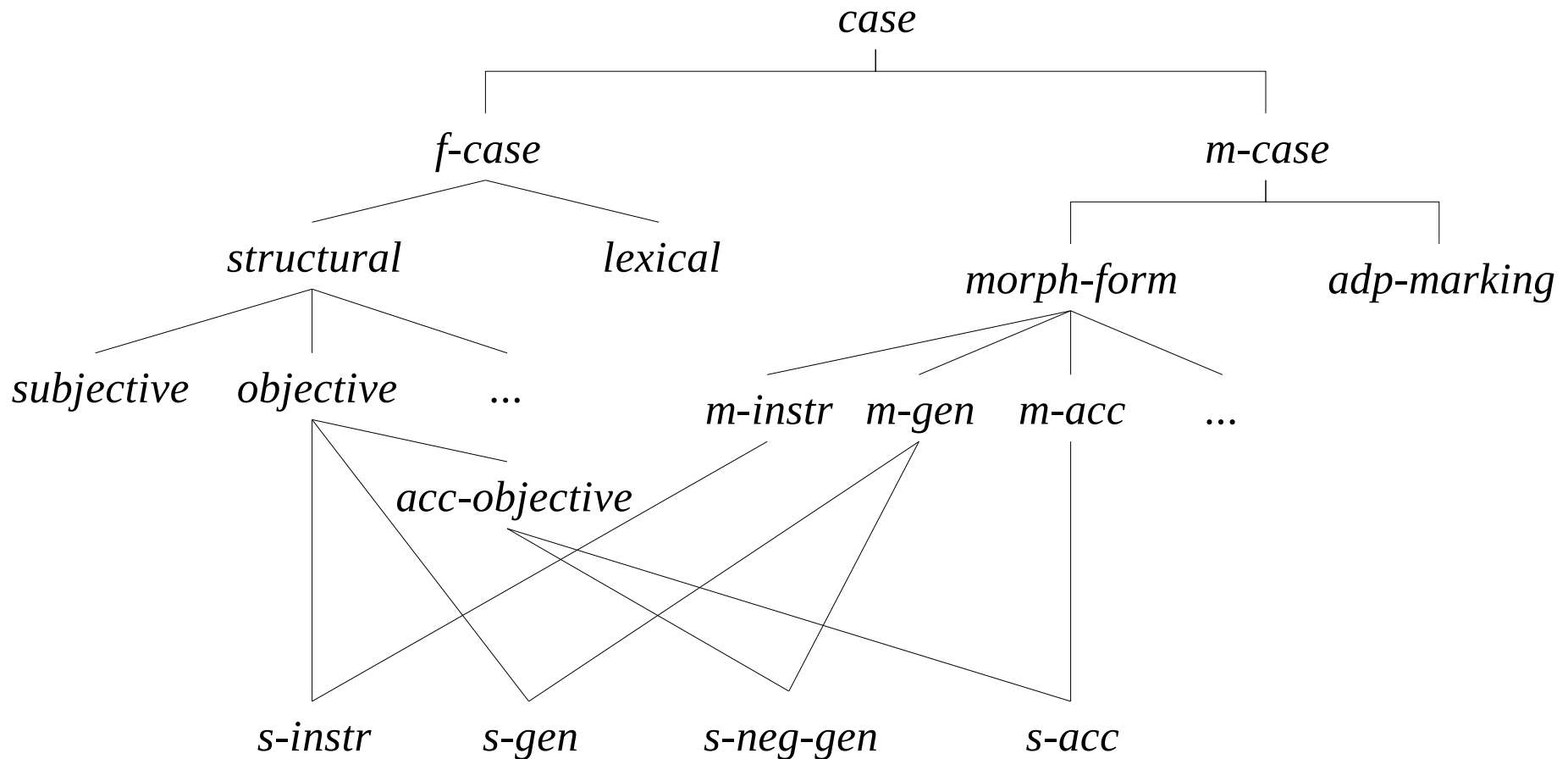


# Parametric constraints on the instantiation of functional case



- Language-specific: the particular encoding of functional case, e.g., the instrumental / prepositional marking of the agentive phrase in passive.
- Idiosyncratic: lexically determined subclasses of structural arguments, e.g. accusative, genitive and instrumental direct objects in Polish.
- Context-sensitive: case alternations triggered by syntactic or semantic context, e.g., genitive of negation.

# Structural case in Polish: the case of direct object NPs



# Future Work

---



- Extend the Russian grammar
- Extend the Slavic Core Grammar
- Start grammars for Polish, Bulgarian, Czech, ... ..