

# Morphological Analysis Unification Grammars

Daniel Zeman

<http://ufal.mff.cuni.cz/course/npf1094>

# Unification Grammars

- Based on
  - context-free grammars
  - **feature structures**
  - their **unifiability**
- Feature structure
  - Sort of database record, or a variable of a structured type: **record** in Pascal or **struct** in C. Description of an object, list of features.
  - features (attributes) ... names of fields
  - values
  - Examples of attribute-value pairs: [number: plural], [case: nominative].

# Feature Structure

<b>entity</b>	
<b>NAME</b>	<b>FF UK</b>
<b>PHONE</b>	<b>258562</b>

<b>POS</b>	<b>noun</b>
<b>GEN</b>	<b>masculine</b>
<b>NUM</b>	<b>singular</b>
<b>CASE</b>	<b>dative</b>

<b>entity</b>	
<b>NAME</b>	<b>Dan</b>
<b>PHONE</b>	<b>221914225</b>

<b>POS</b>	<b>adjective</b>
<b>GEN</b>	<b>masculine</b>
<b>NUM</b>	<b>plural</b>
<b>CASE</b>	<b>accusative</b>
<b>DEG</b>	<b>comparative</b>
<b>NEG</b>	<b>affirmative</b>

<b>faculty</b>	
<b>NAME</b>	<b>MFF UK</b>
<b>DEAN</b>	<b>Netuka</b>
<b>PHONE</b>	<b>221911111</b>

# Feature Structure

- In general: a partial function mapping the set of features to the set of values.

type	
FEATURE <sub>1</sub>	VALUE <sub>1</sub>
FEATURE <sub>2</sub>	VALUE <sub>2</sub>
FEATURE <sub>3</sub>	VALUE <sub>3</sub>

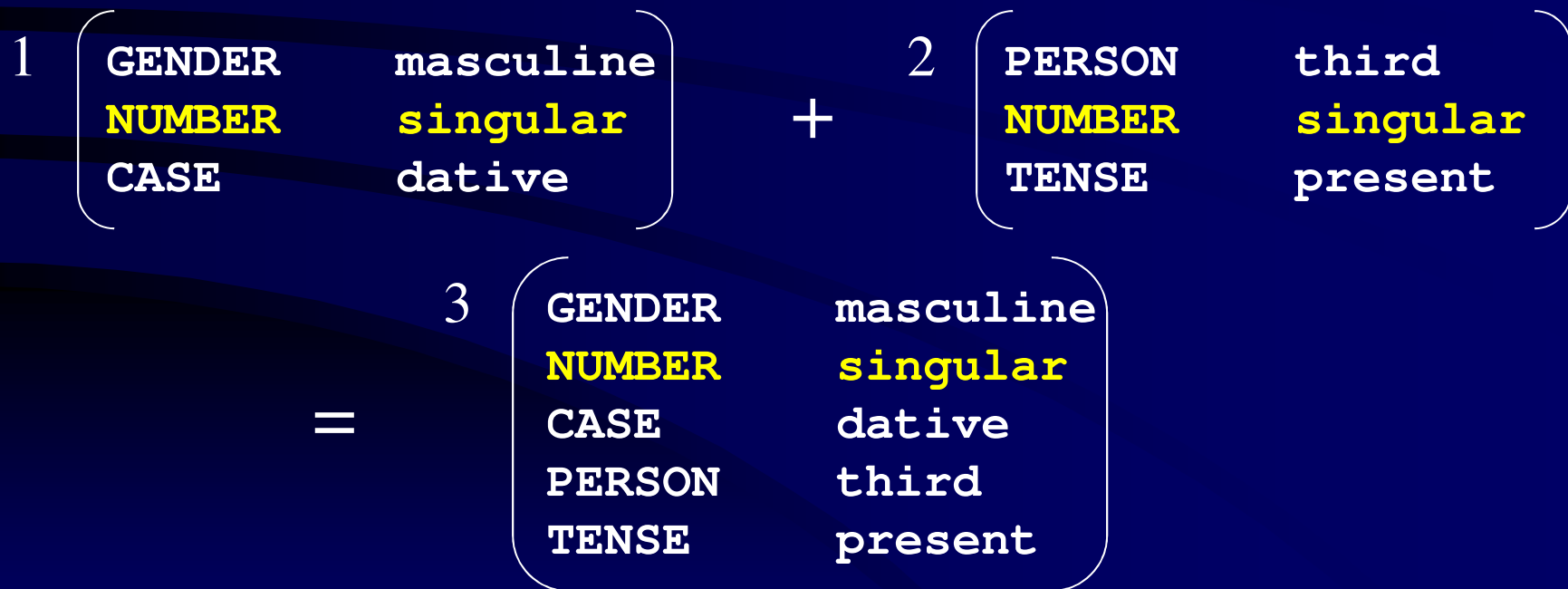
# Unifiability

- Two feature structures are **unifiable** if their values of the features they share are identical.
- Example: structures 1 and 2 are unifiable, so are 2 and 3, while 1 and 3 are not.

1	<table><tbody><tr><td>GENDER</td><td>masculine</td></tr><tr><td>NUMBER</td><td>singular</td></tr><tr><td>CASE</td><td>dative</td></tr></tbody></table>	GENDER	masculine	NUMBER	singular	CASE	dative	2	<table><tbody><tr><td>POS</td><td>verb</td></tr><tr><td>NUMBER</td><td>singular</td></tr><tr><td>TENSE</td><td>present</td></tr></tbody></table>	POS	verb	NUMBER	singular	TENSE	present
GENDER	masculine														
NUMBER	singular														
CASE	dative														
POS	verb														
NUMBER	singular														
TENSE	present														
	3		<table><tbody><tr><td>GENDER</td><td>masculine</td></tr><tr><td>NUMBER</td><td>singular</td></tr><tr><td>CASE</td><td>instrumental</td></tr></tbody></table>	GENDER	masculine	NUMBER	singular	CASE	instrumental						
GENDER	masculine														
NUMBER	singular														
CASE	instrumental														

# Unification

- **Unification** is an operation over two unifiable feature structures. It results in a new feature structure.



# Unification as a Tool for Morphological Generation?

- Input: feature structures “lemma” and “tag”.
- Search lexicon for all structures “entry” that are unifiable with the “lemma” structure.
- For each “entry” structure found look up a “paradigm” structure that is unifiable with both the “entry” and the “tag” structures.
- Unify the corresponding structures “entry”, “paradigm” and “tag”. The resulting structure is called “form”.
- Output: for each “form” structure, concatenate its values of features “paradigm” and “suffix”.

# Unification as a Tool for Morphological Generation?

- Input: feature structures “lemma” and “tag”.

lemma	
LEMMA	háček

tag	
NUMBER	plural
CASE	nominative

Czech noun *háček* has two meanings and belongs to two inflection classes.

“entry” structures unifiable with

entry	
LEMMA	háček
PARADIGM	hrad

entry	
LEMMA	háček
PARADIGM	pán



# Unification as a Tool for Morphological Generation?

- For each “entry” structure find a “paradigm” structure unifiable with both “entry” and “tag”.

<b>entry</b>	
<b>LEMMA</b>	<b>háček</b>
<b>PARADIGM</b>	<b>hrad</b>

<b>entry</b>	
<b>LEMMA</b>	<b>háček</b>
<b>PARADIGM</b>	<b>pán</b>

<b>paradigm</b>	
<b>PARADIGM</b>	<b>hrad</b>
<b>NUMBER</b>	<b>plural</b>
<b>CASE</b>	<b>nominative</b>
<b>SUFFIX</b>	<b>y</b>

<b>paradigm</b>	
<b>PARADIGM</b>	<b>pán</b>
<b>NUMBER</b>	<b>plural</b>
<b>CASE</b>	<b>nominative</b>
<b>SUFFIX</b>	<b>i   ové</b>

# Unification as a Tool for Morphological Generation?

- Unify the corresponding structures “entry”, “paradigm” and “tag”. Call the resulting structure “form”.

<b>form</b>	
<b>LEMMA</b>	<b>háček</b>
<b>PARADIGM</b>	<b>hrad   pán</b>
<b>NUMBER</b>	<b>plural</b>
<b>CASE</b>	<b>nominative</b>
<b>SUFFIX</b>	<b>y   i   ové</b>

# Unification as a Tool for Morphological Generation?

- Unification resembles database operations.
- It does not tell how the “form” structure is to be interpreted.
- Rule: **output = form.lemma + form.suffix**
- The rule does not solve phonological changes (and unification cannot help us with this):

instead of

*\*háčekky, \*háčekki, \*háčekkové*

we want

*háčky, hácci, háčkové*

# Unification as a Tool for Morphological Analysis???

- **Non-unification part:** find all possible affixes recognizable in the word  $\Rightarrow$  set of “form” structures.
- The “paradigm” structures tell us what is the set of known affixes.
- **Somehow solve** phonological changes (stem-final palatalization, stem-internal ablaut etc.)
- Then take the dual procedure to the generation: unify form with paradigm, and the result with the lexicon. Entries found in the lexicon are the possible analyses.
  - e.g. *běžím* (“I am running”) = *běžet*(verb:*trpět*)+person(1<sup>st</sup>),  
≠ *běží*(noun:*stavení*)+case(7)

# Unification as a Tool for Morphological Analysis???

- Non-unification part: find all possible affixes recognizable in the word  $\Rightarrow$  set of “form” structures.
- The “paradigm” structures tell us what is the set of known affixes.
- **Somehow** solve phonological changes (stem-final palatalization, stem-internal ablaut etc.)
- Then use unification...
  - In fact, this is what PC Kimmo v.2 does:
  - It combines two-level morphology with a unification grammar.

# Unification Morphology Grammar (UMG)

- Jan Hajič: *Unification Morphology Grammar (PhD thesis)*. Univerzita Karlova, Praha, 1994
- Stuart Shieber: *An Introduction to Unification-based Approaches to Grammar*. CSLI Lecture Notes No. 4, Stanford, California, USA, 1986
- Based on a **context-free grammar**.
- A feature structure is attached to each constituent (label + span).
- Rule: left-hand side (LHS)  $\rightarrow$  right-hand side (RHS) := operation over feature structures.
- Operations can block a rule by requiring unifiability.
- Unification-based chart parser, PATR-II (Shieber).
- Similarly to CFGs, unification grammars were originally designed for sentence syntax analysis and subsequently applied to word analysis as well.

# UMG Syntax

- LHS  $\rightarrow$  RHS := operation over feature structures
  - grammar rule
- $\langle X \rangle$ 
  - non-terminal symbol X. Terminals are written without angle brackets
- #
  - unification operator (it also places requirement on unifiability)
- ^
  - reference operator (it delimits non-terminals / parts of paths to the feature structure we are referencing)
- +
  - concatenation operator
- |
  - disjunction operator. A disjunction of feature structures contains all structures that fulfill the constraints (are unifiable). A disjunction can represent alternate analyses of the same string.

# Example of UMG Rule

$\langle N \rangle \rightarrow \langle L \rangle := [1 = \langle L \rangle^1, \text{umlaut} = \langle L \rangle^{\text{umlaut}} \# \text{no}]$

- Interpretation:
  - If:
    - we recognized constituent  $\langle L \rangle$  and
    - value of the **umlaut** attribute in the feature structure attached to this constituent is “no”
  - Then:
    - we also recognized constituent  $\langle N \rangle$  with the same span
    - we must copy the attributes **1** and **umlaut** from the feature structure of  $\langle L \rangle$  to the feature structure of  $\langle N \rangle$



# Theoretical View of the Lexicon

- A rule that generates the empty string but it provides a gigantic feature structure with the entire lexicon in it.

```
- <LEX> → "" :=  
  [stem=mat, hw=matka, pos=N, x=zn6e] |  
  [stem=atom, hw=atom, pos=N, x=hd1] |  
  [stem=nov, hw=nový, pos=A, x=reg] |  
  [stem=prac, hw=pracovat, pos=V,  
  x=ovatn] |  
  ...;
```

# Theoretical View of the Lexicon

- How the lexicon is bound to the rest of the grammar:
  - $\langle \mathbf{R} \rangle \rightarrow \langle \mathbf{S} \rangle \mathbf{u} \langle \mathbf{LEX} \rangle := \langle \mathbf{LEX} \rangle \# [\mathbf{x}=\mathbf{hd1}, \text{stem}=\langle \mathbf{S} \rangle, \text{case}=\text{gen}|\text{dat}|\text{loc}, \text{num}=\text{sg}]$ 
    - The rule describes formation of singular noun genitive, dative and locative according to the Czech masculine paradigm **hd1** (*hrad*).
    - **R** represents a word unified with the lexicon.
    - **S** is the part of the input that corresponds to the stem of the word. The suffix is shown literally, the **LEX** that follows it corresponds to the empty string.
    - The operation after **:=** says that we are interested in those structures from **LEX** whose stem corresponds to **S** and which encode genitive, dative or locative according to the paradigm **hd1**.
    - Those lexicon entries that pass this filter will form the set of feature structures bound to the non-terminal **R**. In addition these feature structures will contain information on number and case.

# UMG Example

`<L> → a := [l=a];`

`<L> → b := [l=b];`

...

`<N> → <L> := [l=<L>^1];`

`<N> → <L> <N> := [l=<L>^1+<N>^1];`

`<S> → <N> := <N>;`

`<R> → <S> := <LEX> # [stem=<S>^1, x=hd1, num=sg,  
case=nom|acc, ...];`

`<R> → <S>u := <LEX> # [stem=<S>^1, x=hd1, num=sg,  
case=gen, ...];`

`<LEX> → "" := ... | [stem=hrad, x=hd1, ...] | ...`

# UMC Example

<L> is a letter

<L> → a := [l=a];

<L> → b := [l=b];

...

<N> → <L> := [l=<L>^1];

<N> → <L> <N> := [l=<L>^1+<N>^1];

<S> → <N> := <N>;

<R> → <S> := <LEX> # [stem=<S>^1, x=hd1, num=sg, case=nom, acc, ...];

<R> → <S>u := <LEX> # [stem=<S>^1, x=hd1, num=sg, case=gen, ...];

<LEX> → "" := ... | [s

<N> is a string

<S> is a potential word stem

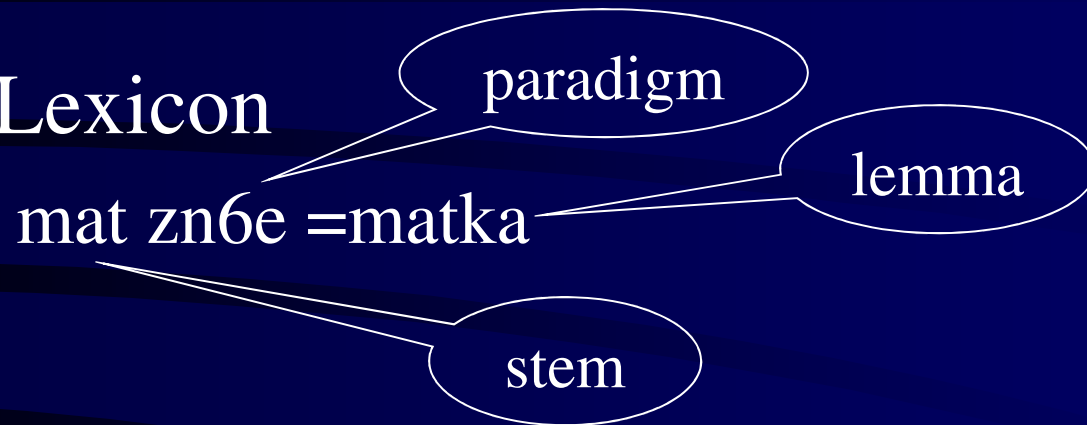
<R> is a recognized word form checked against lexicon

# The Lexicon in Practice

- It is not efficient to treat the lexicon as part of grammar.
- Real implementation looks different:
  - Store the lexicon in a separate data structure with fast search access.
  - Cover rules containing **<LEX>** by the algorithm accessing the data structure.
  - Use the unifying chart parser to process the rest of the grammar.

# UMG Example

- Lexicon



Typical system with many paradigms, e.g. 44 distinct paradigms correspond to the “school” paradigm *žena* (yet the paradigms don’t solve the shortening of the stem-internal vowel).

# UMG Example

```
{ paradigm = stavení; omitting LHS, always the same }  
<_><í>$ := [key=<_>í, x=(st|rž), cat=[pos=n],  
  morf=[infl=[pf=( [gnd=n, num=sg,  
  case=(nom|gen|dat|acc|voc|loc) ] | [gnd=n,  
  num=pl, case=(nom|gen|acc|voc) ] ) ] ] ] ;  
<_><í><m>$ := [key=<_>í, x=(st|rž), cat=[pos=n],  
  morf=[infl=[pf=( [gnd=n, num=sg, case=ins] |  
  [gnd=n, num=pl, case=dat] ) ] ] ] ;  
<_><í><c><h>$ := [key=<_>í, x=(st|rž),  
  cat=[pos=n], morf=[infl=[pf=[gnd=n, num=pl,  
  case=loc] ] ] ] ;  
<_><í><m><i>$ := [key=<_>í, x=(st|rž),  
  cat=[pos=n], morf=[infl=[pf=[gnd=n, num=pl,  
  case=ins] ] ] ] ;
```

# Comparison of UMG and CFG

- The feature structure contains the required output (tag)  $\Rightarrow$  no need for supplementary non-terminal naming convention
- At the same time, the features and their unifiability constrains rule application  $\Rightarrow$  no need to split non-terminals
- Disjunction of structures represents homonymous analyses
- Phonology is still an issue. Either combinatorial explosion of paradigms (UMG) or use in tandem with two-level rules (see below)



# PC-Kimmo Word Grammar

- Unification grammar by Stuart Shieber. Rule syntax somewhat different from UMG, application is similar.
- lexicon
  - recognize possible morphemes in the word
- rules
  - phonological changes, esp. on morpheme boundary
- grammar
  - analysis of inter-morpheme relations
  - derivation of word features from morpheme features
  - constraints on morphotactics (what morphemes can combine and in what order)

# PC-Kimmo Word Grammar

```

en   +`large  +ment  +s
VR1a +`large  +NR25  +PL

          Word
        _____|_____
          Stem          INFL
          _____|_____
          Stem          SUFFIX  +s
          _____|_____
          PREFIX Stem  +ment   +PL
          _____|_____
          ent+         |
          VR1a+       ROOT
                    `large
                    `large
    
```

```

Word:
[ cat:      Word
  head:     [ agr:
              [ 3sg: - ]
              number:PL
              pos:   N ]
  root:     `large
  root_pos:AJ
  clitic:-
  drvstem:- ]
    
```

# PC-Kimmo Word Grammar

- First the old part of PC-Kimmo segments the word into morphemes.
- Then the new part parses the sequence of morphemes according to the grammar.
  - Grammar can reject some morpheme sequences.
  - It assigns an interpretation (a feature structure) to the sequences it accepts. The old PC-Kimmo was able to gloss morphemes but it could not tell how to combine the morpheme glosses into the interpretation of the whole word (e.g. that the suffix *-able* changes a verb to an adjective).
- A grammar rule looks like this:
  - Word → Stem INFL
    - <Stem head pos> = <INFL from\_pos>
    - <Word head> = <INFL head>

# Grammar Rule

- **Word**  $\rightarrow$  **Stem** **INFL**  
    <Stem head pos> = <INFL from\_pos>  
    <Word head> = <INFL head>
- The rule cannot be used if the feature **pos** of the substructure **head** of the morpheme **Stem** is not equal to the feature **from\_pos** of the morpheme **INFL**.
  - The morpheme symbols are pre-terminals and they correspond to the names of sublexicons where the morphemes were found.
- If the rule is used it shall copy the value of the **head** feature from the **INFL** constituent to the **head** feature of the **Word** constituent.

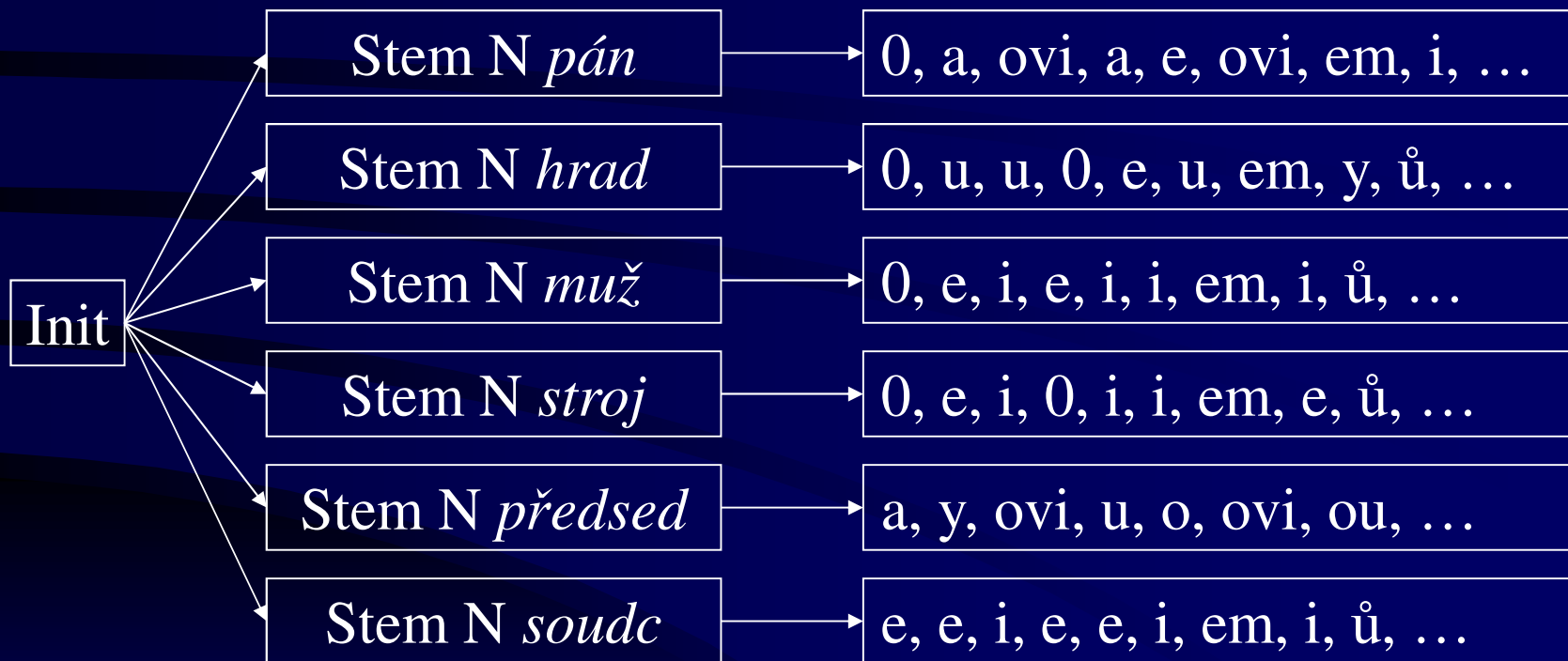
# Grammar Rule

- **RULE** <rule>  
    <rule constraints>
- Left-hand side is separated from the right-hand side by  $\rightarrow$  or  $=$ .
- **RULE Stem\_1 = Stem\_2 SUFFIX**
- **X** represents any terminal or non-terminal.
- Characters **() [] {} <> = : /** are special.
  - Underscore is used only to append an index to a symbol.
- Left-hand side of the first rule is the start symbol of the grammar.
- **N = Nstem {Sing / Plural}**

# Advantages of the Grammar

- Czech examples:
  - Grammar blocks combination of a stem belonging to the *žena* paradigm with a suffix belonging to the *ruže* paradigm.
  - It can even check long-distance dependencies such as:
    - nejchytřejší
  - Take feminine noun *žena* (*woman*). Derive from it the possessive adjective *ženin* (*woman's*). The gender changed from feminine to masculine (the suffix says that the possessed object is masculine). The grammar can store the original gender as the lexical possessor's gender.

# Without Grammar



# With Grammar

Init

Stem N MASC

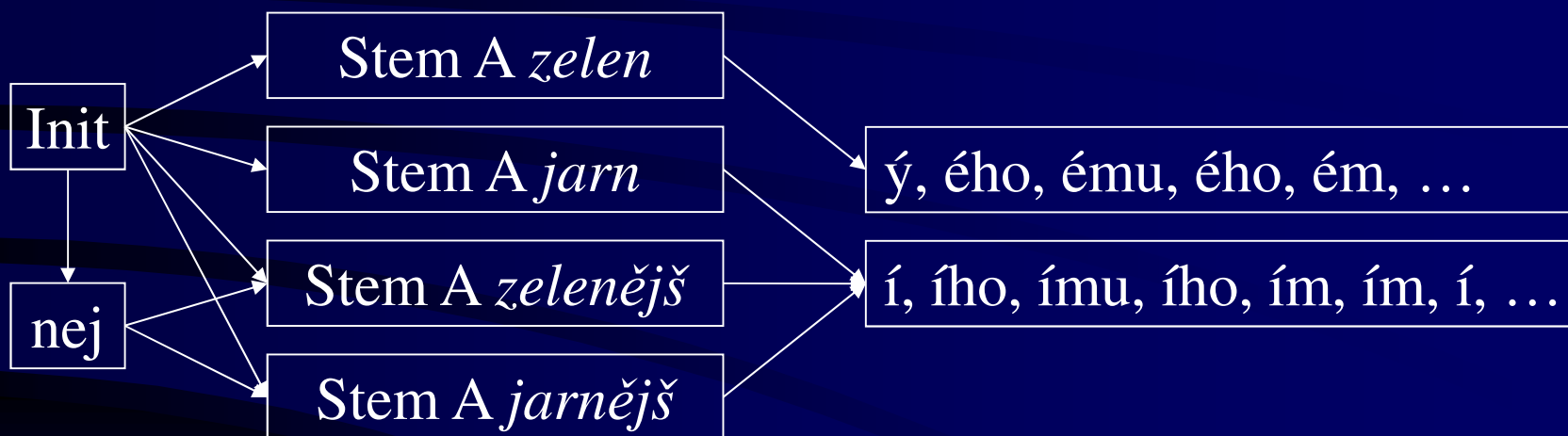
[paradigm: x]

0, a, ovi, e, em, i, ové, ů,  
ŭm, y, ech, u, é, ích, o, ou

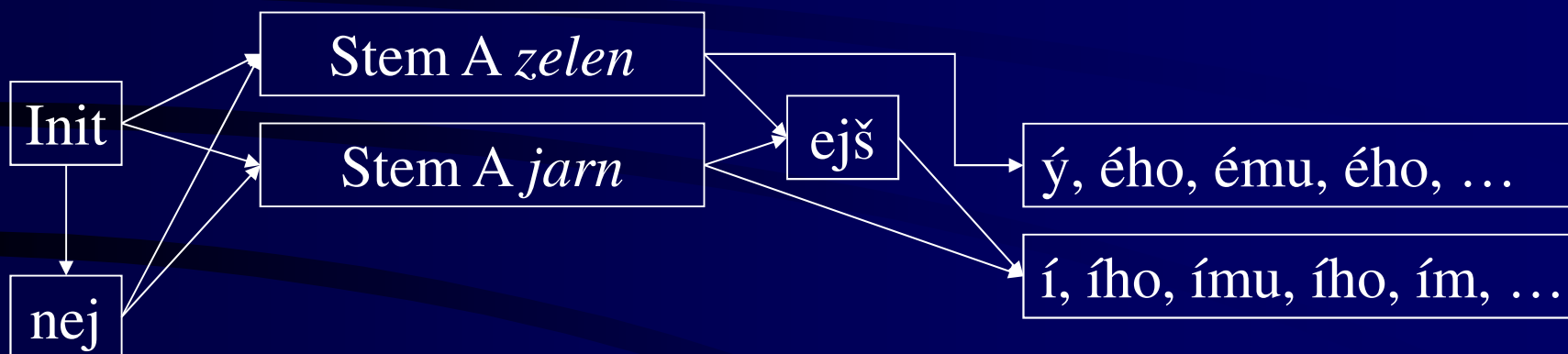
[paradigm: x]



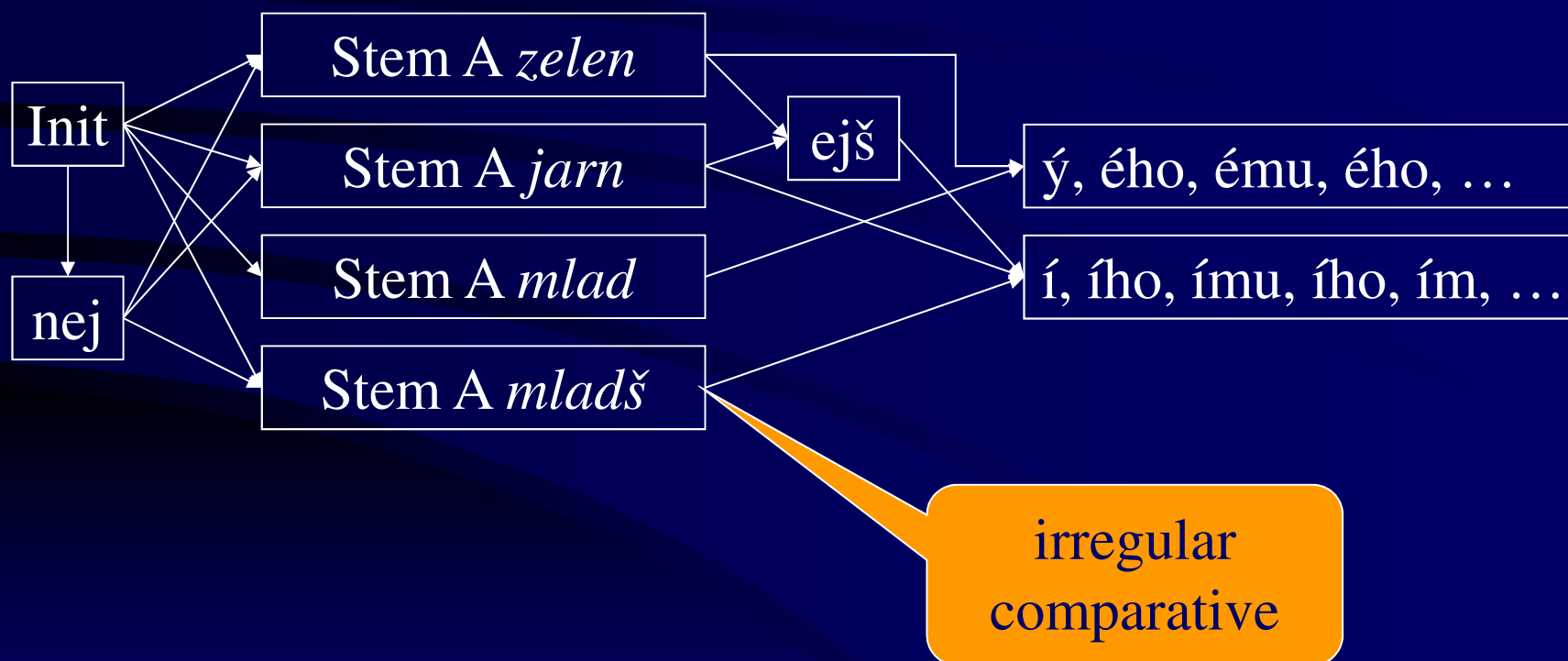
# Without Grammar



# With Grammar



# With Grammar



# Grammar Cannot Interact with Phonology

- Phonological rule of consonant softening in Czech imperative:
  - metěš (“you sweep”) → met’(me,te) (“sweep!”)
  - t:t’ ⇔ \_ +:0 λ:0 or m:m e:e or t:t e:e
- The rule must not apply in genitive plural form of feminine nouns:
  - kóta (“spot elevation”) → \*kót’
- Phonological rules cannot read the feature structures to constrain their application.
- There have been extensions other than PC-Kimmo that combined phonological rules with feature structures, e.g. Trost (1990).

# Automatic Features

- Every lexicon entry automatically receives the following features:
  - cat = name of sublexicon (\lx)
  - lex = morpheme, lexical string (\lf)
  - gloss = gloss from the lexicon entry (\gl)

# Assigning Values to Features

- Abbreviations of feature assignments
  - If we are going to assign a value to thousands of lexicon entries we want it to be as short as possible.
  - **LET** <shortcut | category> **be** <definition>
  - e.g.
  - **Let pl be [number: PL]**
  - **Let pl be <number> = PL**
  - **Let 3sg be [tense: PRES  
agr: 3SG]**
- **Disjunction:**
  - **Let sg/pl be {[number: SG] [number: PL]}**
  - **Let sg/pl be <number> = {SG PL}**
- Default values:
  - **Let N be <number> = !SG**
  - Unless someone explicitly assigns the value of number to a noun, the noun is assumed to be in singular.

# Lexical Rules

- Not shortcuts but systematic transformation of features for groups of lexicon entries. They transform a feature structure to another one.
- **DEFINE** <lexical rule name> **as** <mapping>
- The example in the on-line documentation is invalid.
- When the analysis is done and the feature structure for the whole word is ready, we can apply a lexical rule that will modify the structure.

# Parameter Setting

- **PARAMETER** <name> **is** <value>
  - Parameter Start symbol is **Word**
  - Parameter Attribute order is **cat head root**
    - In which order shall PC-Kimmo display the features?
  - **Category feature** (default: **cat**)
  - **Lexical feature** (default: **lex**)
  - **Gloss feature** (default: **gloss**)
    - What are the names of important features with special meaning?



# PC Kimmo Demonstration

- r ženě
- Synthesis (new in PCK v. 2, but grammar is still not necessary)
- l synthesis-lexicon cs.lex
- s N(žena) +SG+LOC
- If a grammar is available it may block synthesis of invalid combinations
- Unfortunately we cannot generate from a feature structure