

# Pattern Generation Revisited

David Antoš, Petr Sojka

Faculty of Informatics  
Masaryk University Brno

Czech Republic

# What is this lecture about

- How T<sub>E</sub>X hyphenates words
- Generation of patterns
- Using PATGEN
- Why PATGEN doesn't suit us today
- PATLIB and OPATGEN, system architecture
- Pattern recognition in typesetting
- Future applications

# How T<sub>E</sub>X hyphenates words

- Before 2<sup>nd</sup> pass of paragraph breaking, all words
- T<sub>E</sub>X finds all *patterns* of the word, all hyphens
- Pattern is a subword, with hyphenation information between symbols
  - n2at hen5at
- Hyphenation information—numbers
  - *odd*—allow breaking—*covering*
  - *even*—disallow breaking—*inhibiting*
- Patterns *compete*—higher value wins

## Example of application of patterns

h y p h e n a t i o n

1 n a

1 t i o

n 2 a t

2 i o

h e 2 n

. h y 3 p h

h e n a 4

h e n 5 a t

. h 0 y 3 p 0 h 0 e 2 n 5 a 4 t 2 i 0 o 0 n .

h y - p h e n - a t i o n

- Patterns hold context of hyphenation point occurrence

- Pattern recognition—a kind of intelligence (magic?)

# Pattern generating process

- We want *minimal* set of patterns
- Generating completely covering minimal patterns is NP-complete
- Good iterative methods for “near optimal” solution
- One way to do it:
  - dic-tio-nary with marked hy-phen points
  - we repeat going through the dictionary making *levels*
    - ★ *odd*—covering
    - ★ *even*—inhibiting

- *candidate choosing rule*—we take subwords (k-character subword)
  - ★ and count the number of their good and bad work
  - ★  $hen(35, 4)a$
- not all the candidates are good
  - ★ *pattern choosing rule*—linear function over the number of good and bad work of the pattern compared to a threshold:
 
$$good\_count * good\_wt - bad\_count * bad\_wt \geq threshold$$
- we put good candidates into pattern set with current level number
  - ★  $hen5a$
- they still make errors, next level will correct errors of patterns selected so far
- exception list at the end of the whole process, to correct remaining errors

# How good the generation process is

- Input data  $\approx$  MB
  - patterns  $\approx$  tens of KB
  - covering  $> 98\%$
  - with error  $< 0.1\%$
  - usual number of levels  $\leq 5$
- Results depend on parameter setting
- Nobody knows how to set the parameters

# Generating with PATGEN

- PATGEN
  - Frank Liang, 1982, 7-bit ASCII
  - later modified for accents, 8-bit ASCII, ... (Peter Breitenlohner, Yannis Haralambous, Karl Berry, ...)
- Patterns for tens of languages exist
  - from a dictionary
  - first levels by hand



# Why PATGEN does not suit us today

- Monolithic structured code, difficult to maintain
- Only ASCII, approx. 240 symbols
- $\Omega$  handles UNICODE
- Static data structures
- Difficult to use for purposes other than hyphenation

# PATLIB and OPATGEN

- PATtern manipulating LIBrary—from scratch
- O (Ω?) PATtern GENerator—handles UTF-8 UNICODE encoding
- Generalisation of PATGEN
- Implementation: C++ in CWEB
- <http://www.fi.muni.cz/~xantos/patlib>
- GPL

# PATLIB architecture

- PATLIB
  - Finite language store (pattern manipulator) implemented using packed trie (see the Proceedings)
    - ★ insert pattern
    - ★ delete pattern
    - ★ ...low-level operations
  - Generator
    - ★ handles creating patterns using the same strategy as PATGEN
    - ★ easy to change
- OPATGEN is in fact “only” the word list input/output interface

# Pros and cons

- What do we get
  - implementation using logical operations
  - any strategy of generating is easier to implement
  - generality and flexibility
  - user manual :-)
- What do we pay for it
  - performance loss

# Generalisation, abstraction

- Hyphenating point → *point of interest*
- Patterns *capture* information related to points of interest
- Patterns *compress* information on points of interest
- General, not tight to hyphenation
- Using of patterns is effective, information retrieval is *linear* wrt. the “word length”
- May be implemented as Pattern Translation Process (PTP)
- Some can be done as  $\Omega$ TP in Omega

## Examples of usage

- Hyphenation of compound words
  - we want to prefer breaking on word boundaries
  - needed for German (Wort-silben-trennung)
- Context dependent ligatures
  - ligatures over compound word boundaries are wrong
  - English: shelfful vs. shelfful
  - Czech: šéflékař vs. šéflékař (doctor in chief)
- Fraktur long s versus short s
  - morphology dependent

- End of sentence
  - different space width
  - $\backslash@$  and  $\backslash_$  in L<sup>A</sup>T<sub>E</sub>X (sorry, the slides are in ConT<sub>E</sub>Xt)
- Thai segmentation
  - word/sentence boundaries are not present in the input Thai transcription, needed for line-breaking
- Arabic letter hamza
  - five hamza variants, depending on context
- Adding Greek accents
- Adding Czech accents to 7-bit ASCII texts
  - e-mail,...





# Performance loss

- 10–15 times to PATGEN
- Caused by
  - breaking the program into logical layers
  - using dynamic memory
  - alphabet symbol is an object, not mere number
  - not (yet, if ever) optimised