

Zpravodaj Československého sdružení uživatelů TeXu

Ondřej Sojka; Petr Sojka; Jakub Máca
A Roadmap for Universal Syllabic Segmentation

Zpravodaj Československého sdružení uživatelů TeXu, Vol. 33 (2023), No. 3-4, 125–138

Persistent URL: <http://dml.cz/dmlcz/151995>

Terms of use:

© Československé sdružení uživatelů TeXu, 2023

Institute of Mathematics of the Czech Academy of Sciences provides access to digitized documents strictly for personal use. Each copy of any part of this document must contain these *Terms of use*.



This document has been digitized, optimized for electronic delivery and stamped with digital signature within the project *DML-CZ: The Czech Digital Mathematics Library* <http://dml.cz>

A Roadmap for Universal Syllabic Segmentation

ONDŘEJ SOJKA, PETR SOJKA, JAKUB MÁČA

Space- and time-effective segmentation (word hyphenation) of natural languages remains at the core of every document rendering system, be it $\text{T}_{\text{E}}\text{X}$, web browser, or mobile operating system. In most languages, segmentation mimicking syllabic pronunciation is a pragmatic preference today.

As language switching is often not marked in rendered texts, the typesetting engine needs *universal syllabic segmentation*. In this article, we show the feasibility of this idea by offering a prototype solution to two main problems:

- A) Using Patgen to generate patterns for several languages at once; and
- B) lack of Unicode support in tools like Patgen or $\text{T}_{\text{E}}\text{X}$ (patterns in UTF-16 encoding) is missing.

For A), we have applied it to generating universal syllabic patterns from wordlists of nine syllabic, as opposed to etymology-based, languages (namely, Czech, Slovak, Georgian, Greek, Polish, Russian, Turkish, Turkmen, and Ukrainian). For B), we have created a version of Patgen that uses the Judy array data structure and compared its effectiveness with the trie implementation.

With the data from these nine languages, we show that:

- A) developing universal, up-to-date, high-coverage, and highly generalized universal syllabic segmentation patterns is possible, with a high impact on virtually all typesetting engines, including web page renderers; and
- B) bringing wide character support into the hyphenation part of the $\text{T}_{\text{E}}\text{X}$ suite of programs is possible by using Judy arrays.

Keywords: syllabification, hyphenation, universal syllabic patterns preparation

1. Motivation

Justified alignment achieved with a quality hyphenation algorithm is both optically pleasing and saves time to read, in addition to saving trees. Only *quality hyphenation* allows interword spaces to be as uniform as possible, close to Gutenberg's ideal of spaces of fixed width. A high coverage, space- and time-effective hyphenation (segmentation) algorithm of all natural languages is badly needed¹

This is an updated and enriched version of the paper [1] published in the TUG 2023 proceedings issue of the TUGboat journal.

¹https://bugzilla.mozilla.org/show_bug.cgi?id=672320

as it remains at the core of every document rendering system, be it \TeX , web browsers supporting HTML with CSS 3, or an operating system providing text rendering for mobile applications.

In most languages, segmentation mimicking syllabic pronunciation is pragmatically preferred today. As language switching is often not marked in texts, and cannot be safely guessed from the words themselves, language-agnostic orthographic syllabification, is needed. We call this task *universal syllabic segmentation*, or in short, the syllabification problem.

The syllabification problem has been tackled by several finite state [2] or, more recently, machine learning techniques [3, 4, 5, 6]. Bartlett et al. [3] uses structured support vector machines (SVM) to solve syllabification as a tagging problem. Krantz et al. [7] leverage modern neural network techniques with long short-term memory (LSTM) cells, a convolutional component, and a conditional random field (CRF) output layer, and demonstrated cross-linguistic generalizability, syllabifying English, Dutch, Italian, French, Manipuri, and Basque datasets together.

From an orthographic viewpoint (hyphenation), universal language solutions today should reflect the Unicode standard [8]. Internal support for full (multibyte) Unicode, a must in today’s operating systems and applications, *is missing* in the \TeX family of programs, e.g. in Patgen and \TeX itself. The internal processing is thus limited by the internal one-byte representation of language characters and is hardwired into the optimized code of these programs. Therefore, processing CJK (Chinese, Japanese, Korean) languages or sets of languages whose character representations need multibyte handling is close to impossible. Special “hacks” are needed for character and font encodings both on the input side (package `inputenc`) and output side (packages `fontenc` or `fontspec`) are not backed by internal wide character support.

Since both \TeX and Patgen have hardwired 8-bit character representations, to develop practically useable universal syllabic hyphenation, one needs to overcome these constraints.

In this paper we a) constructively show the feasibility of preparation of universal syllabic patterns, b) demonstrate a version of Patgen with wide character support, and c) discuss further steps to do in the \TeX program suite to make language hyphenation Unicode-compliant.

The paper is structured as follows. In Section 2 we define the terminology and describe the language data we have used in our experiments. Section 3 reminds the reader about the principles of the hyphenation algorithm in \TeX and of Patgen-based pattern generation and pattern representation possibilities. Section 4 evaluates the experiments with universal pattern generation. In Section 5 we elaborate on possible routes towards wide character support in the typesetting engines and Patgen. As usual, we sum up and conclude in the final Section 6.

“The concept of the syllable is cross-linguistic, though formal definitions are rarely agreed upon, even within a language. In response, data-driven syllabification methods have been developed to learn from syllabified examples. . . . Syllabification can be considered a sequence labeling task where each label delineates the existence or absence of a syllable boundary.” [7]

2. Syllabification

Human beings convey meaning by pronouncing words as sequences of phonemes. Phonology studies the structure of phonemes we are able to pronounce as syllables [9]. Etymologically, a syllable is an Anglo-Norman variation of Old French *sillabe*, from Latin *syllaba*, from Greek (*syllabē*), “that which is held together; a syllable, several sounds or letters taken together” to make a single sound. [10]

When we delineate boundaries in the orthographic representation of words, we speak about *hyphenation* of words as sequences of *characters*.

2.1. Hyphenation as syllabification

There are subtle differences between syllabification and hyphenation, though. Let us take the Czech word *sestra*. The Czech language authorities [11] allow hyphenations as *se-s-t-ra*, while agreeing that there are only two syllables based on Consonant and Vowel sequencing: either *se-stra* (CV-CCCV), or *ses-tra* (CVC-CCV), or *sest-ra* (CVCC-CV). As with hyphenation, defining segments for syllabification is full of exceptions. The Czech sentence *Strč prst skrz krk* or word *scrnkl* (CCCCCCC) contain consonants-only syllables.

There are also rare cases where word segmentation should differ in different contexts. It may be necessary within one language (different hyphenation *re-cord* and *rec-ord* depending on its part of speech), or between different languages. When developing universal syllabic patterns, these theoretically possible segmentations should not be allowed in the input hyphenated wordlist used for training. But this should not matter, as e.g. Liang’s `hyphen.tex` patterns do not cover more than 10% of positions [12] and few complain about this coverage.

2.2. Data preparation

To show the feasibility of universal pattern generation, we have collected wordlists for a dozen languages, as shown in Table 1. The chosen languages a) have a wide diversity in alphabets and syllables and b) have existing hyphenation patterns as an approximation for syllable segments. The wordlists were collected from public sources or provided for our research as stratified dictionaries from TenTen corpora [15] by Lexical Computing. We used wordlists sorted by frequency and cut at below 5% of word occurrences, to eliminate typos appearing in documents. Each tenth word was taken into a wordlist—a stratified sampling technique inspired by

Table 1: Language resources and patterns used in pattern development experiments. All data was converted to UTF-8 and contains lowercase alphabetic characters only. Alphabet size (# chars) counts characters appearing in the language wordlist collected. Languages were chosen for diversity of size of patterns and syllables.

Language	#words	#chars	#patterns	#syllables	source, alphabet
Czech+Slovak (cz+sk)	606,499	47	8,231	2,288,413	[13], Latin
Georgian (ka)	50,644	33	2,110	224,799	[14], Georgian
Greek (el-monoton)	10,432	48	1,208	37,736	[14], Greek
Panjabi (pa)	892	52	60	2,579	[14], Gurmukhi
Polish (pl)	20,490	34	4,053	65,510	[14], Latin
Russian (ru)	19,698	33	4,808	75,532	[14], Cyrillic
Tamil (ta)	46,526	48	71	209,380	[14], Tamil
Telugu (te)	28,849	66	72	125,508	[14], Telugu
Thai (th)	757	64	4,342	1,185	[14], Thai
Turkish (tr)	24,634	32	597	103,989	[14], Latin
Turkmen (tk)	9,262	30	2,371	33,080	[14], Latin
Ukrainian (ua)	17,007	33	1,990	65,099	[14], Cyrillic

Knuth [16] that was already used successfully in pattern generation [17]. Wordlists were hyphenated by legacy patterns, mostly taken from [14].

Alphabet analysis and statistics are shown in Table 2. The total number of letters appearing in all languages exceeds 245, the maximum number of letters that current Patgen can support. This is why wide-character representation (Unicode UTF-16) support in Patgen (and then in the hyphenator library in a typesetting engine) would be needed to extend our generation to more languages.

3. Pattern development

The idea of squeezing the hyphenated wordlist into the set of patterns was originated in the dissertation of Frank Liang [12], supervised by Donald Knuth. For the automated generation of patterns from a wordlist, Liang wrote the Patgen program. Patgen was one of the very first programs that harnessed the power of data with supervised machine learning. Programmed originally to support English and ASCII, it was later extended to be usable for 8-bit characters and for wordlists that contain at most 245 characters [18]. It is capable of efficient lossy or lossless *compression* of hyphenated dictionaries, with several orders of

magnitude compression ratio. Generated patterns have minimal length, e.g., the shortest context possible, which results in their *generalization* properties.

In general, *exact lossless pattern minimization is non-polynomial* by reduction to the minimum set cover problem [19]. For Czech, *exact lossless pattern generation is feasible* [20], while reaching *100% coverage* and simultaneously *no errors*. Strict pattern minimality (size) is not an issue nowadays.

This idea and its realization is a programming pearl. Motivated by space and time constraints, instead of the classical solution of dictionary problem in the logarithmic time of dictionary size, the word hyphenation is computed from patterns in constant time, where the constant is given by *word length*.

Space needed for patterns in the *packed trie* data structure is typically in tens of kB, which is several orders of magnitude smaller than the wordlist size. With fine-tuned parameters of pattern generation in the so-called *levels*, one can prepare patterns with zero errors and almost full coverage of hyphenation points from the input dictionary.

For practical use, patterns are collected in the repository maintained by the T_EX community [14]. It is no surprise that most if not all leading typesetting engines deploy this “competing pattern engineering technology” [21].

Table 2: Language alphabet overlaps. Cells contain a number of lowercase letters that overlap between languages. In total, 13 languages contain in total 412 different lowercase letters, more than Patgen is capable of digesting.

Language	cz+sk	ka	el	pa	pl	ru	ta	te	th	tr	tk	ua
Czech+Slovak (cz+sk)	47	0	0	0	26	0	0	0	0	25	28	0
Georgian (ka)	0	33	0	0	0	0	0	0	0	0	0	0
Greek (el-monoton)	0	0	48	0	0	0	0	0	0	0	0	0
Panjabi (pa)	0	0	0	52	0	0	0	0	0	0	0	0
Polish (pl)	26	0	0	0	34	0	0	0	0	23	22	0
Russian (ru)	0	0	0	0	0	33	0	0	0	0	0	29
Tamil (ta)	0	0	0	0	0	0	48	0	0	0	0	0
Telugu (te)	0	0	0	0	0	0	0	66	0	0	0	0
Thai (th)	0	0	0	0	0	0	0	0	64	0	0	0
Turkish (tr)	25	0	0	0	23	0	0	0	0	32	25	0
Turkmen (tk)	28	0	0	0	22	0	0	0	0	25	30	0
Ukrainian (ua)	0	0	0	0	0	29	0	0	0	0	0	33

	h y p h e n a t i o n	
p1	1n a	hy-phen-ation → 2 6
p1	1t i o n	... → ...
p2	n2a t	... → ...
p2	2i o	key → data
p2	h e2n	
p3	h y3p h	Solution to the dictionary problem:
p4	h e n a4	For key part (the word) to store
p5	h e n5a t	the data part (its division)
	h0y3p0h0e2n5a4t2i0o0n	

Figure 1: Eight patterns “compete” how to hyphenate *hyphenation*. Winners are patterns **hy3ph** and **hen5at** generated at the highest covering level (odd numbers) generation. The level hierarchy allows for storing exceptions, exceptions to exceptions, exceptions to exceptions to exceptions, . . . , with character contexts as parameters. [12]

3.1. Patterns

The patterns “compete” with each other whether to split the word at a position, given varying characters in both side contexts; see Figure 1.

We have shown how effective and powerful the technique is, and that its power depends on the *parameters* of pattern generation [20]. The key is the proper setting of Patgen parameters for pattern generation. The idea of universal segmentation with Patgen has been proposed already in [22]. There, we demonstrated the techniques for the development of two languages together, Czech and Slovak, and developed a joint wordlist and patterns [13].

We wanted to extend the technique to other Slavic and syllabic languages. The bottleneck for adding new languages was Patgen and T_EX’s constraint of one-byte character support only for storing patterns in tries. We thought of using a modern data structure that would allow wide character trie representation. That was the task for a bachelor’s thesis: use a Judy array [23].

3.2. Judy arrays

The Judy array, also known as simply Judy, is a data structure that implements a sparse dynamic array, allowing for versatile applications such as dynamically sized arrays and associative arrays. Judy is internally implemented as a tree structure, where every internal node has 256 ancestor nodes. The most interesting thing about this structure is that it tries to be as memory-efficient as possible by effectively using the available cache, avoiding unnecessary access to the main memory. As a result, Judy is both fast and memory-efficient.

The feasibility of utilizing the Judy structure for storing hyphenation patterns is demonstrated in the thesis [23]. In Chapter 4, it is shown that Judy has the potential to be faster and more memory-efficient compared to the original trie when working with patterns. Further, Chapter 5 explores the potential integration of Judy into Patgen and the consequent impact on Patgen’s generation process. The results from this chapter indicate that rewriting Patgen with Judy is possible but would require an almost complete overhaul of Patgen’s code and algorithms. This redevelopment would yield a Patgen version capable of handling input of any kind, enabling the generation of patterns composed of arbitrary alphabets. However, it is important to note that the generation process would be approximately four times slower than the current implementation. This is due to the hiding of access to the inner nodes of stored tries in Judy. As this access is not needed in \TeX for the hyphenation of individual words, using some variant of Judy in a \TeX successor would make hyphenation faster.

Table 3: Different word hyphenation overlaps. Cells contain a number of the same words that are segmented differently between languages. Differences are caused typically by suboptimal coverage patterns used to hyphenate the wordlist (*vibram* vs. *vib-ram*, *up-gra-de* vs. *upg-ra-de*). We remove the differently hyphenated words when joining wordlists for the final syllabic generation.

Language	cz+sk	ka	el	pa	pl	ru	ta	te	th	tr	tk	ua
Czech+Slovak (cz+sk)	9	0	0	0	388	0	0	0	0	640	69	0
Georgian (ka)	0	0	0	0	0	0	0	0	0	0	0	0
Greek (el-monoton)	0	0	0	0	0	0	0	0	0	0	0	0
Panjabi (pa)	0	0	0	0	0	0	0	0	0	0	0	0
Polish (pl)	388	0	0	0	0	0	0	0	0	187	9	0
Russian (ru)	0	0	0	0	0	0	0	0	0	0	0	125
Tamil (ta)	0	0	0	0	0	0	0	0	0	0	0	0
Telugu (te)	0	0	0	0	0	0	0	0	0	0	0	0
Thai (th)	0	0	0	0	0	0	0	0	0	0	0	0
Turkish (tr)	640	0	0	0	187	0	0	0	0	0	80	0
Turkmen (tk)	69	0	0	0	9	0	0	0	0	80	0	0
Ukrainian (ua)	0	0	0	0	0	125	0	0	0	0	0	0

3.3. Universal pattern generation

To pursue the idea of universal syllabic pattern generation, we have checked whether the legacy patterns hyphenate the same valid word in different languages differently. The result with a short discussion is in Table 3. The expectation

Table 4: Statistics from the generation of universal patterns for cz+sk, ka, el, pl, ru, tr, tk, ua with *custom* parameters and `\lefthyphenmin=2, \righthyphenmin=2`. Generation took 33.23 seconds, 11,238 patterns, 77 kB.

Level	Patterns	Good	Bad	Missed	Lengths	Params
1	2,407	2,066,410	280,020	70,588	1 3	1 3 12
2	2,375	2,025,245	8,866	111,753	2 4	1 1 5
3	4,626	2,118,063	19,213	18,935	3 6	1 2 4
4	2,993	2,117,739	5,920	19,259	3 7	1 4 2

Table 5: Statistics from the generation of universal patterns for cz+sk, ka, el, pl, ru, tr, tk, ua with *correct optimized* parameters and `\lefthyphenmin=2, \righthyphenmin=2`. Generation took 35.43 seconds, 29,742 patterns, 219 kB.

Level	Patterns	Good	Bad	Missed	Lengths	Params
1	7,188	2,049,375	164,224	87,623	1 3	1 5 1
2	4,108	2,042,249	14,094	94,749	1 3	1 5 1
3	15,010	2,134,692	20,544	2,306	2 6	1 3 1
4	6,920	2,133,458	815	3,540	2 7	1 3 1

Table 6: Statistics from the generation of universal patterns for cz+sk, ka, el, pl, ru, tr, tk, ua with *size optimized* parameters and `\lefthyphenmin=2, \righthyphenmin=2`. Generation took 29.75 seconds, 14,321 patterns, 101 kB.

Level	Patterns	Good	Bad	Missed	Lengths	Params
1	1,201	2,092,928	598,321	44,070	1 3	1 2 20
2	2,695	1,736,372	5,274	400,626	2 4	2 1 8
3	4,835	2,102,803	20,094	34,195	3 5	1 4 7
4	6,508	2,099,607	210	37,391	4 7	3 2 1

that syllable-forming principles are universal, as phonology theory suggests, is confirmed. The errors we have found were due to the difference between hyphenation and syllabification caused by inconsistent markup rather than a principled difference in word morphology, e.g. a compound word segmented in one language, and given as a single word in the other.²

²Compound words can evolve in perception into single words even within one language. Examples are the evolution of *e-mail* into *email* or *roz-um* into syllabic *ro-zum* in Czech.

Table 7: Comparison of the efficiency of different approaches to pattern generation of Czechoslovak and of universal patterns. Note that the size of universal patterns grows sublinearly with the number of languages. The generalization ability of universal patterns is only slightly worse than that of Czechoslovak ones. The experience from the development of Czechoslovak patterns shows that performance could be improved by consistent markup of wordlist data.

Wordlist	Parameters	Good	Bad	Missed	Size	Patterns
Czechoslovak	custom	99.87%	0.03%	0.13%	32 kB	5,907
Czechoslovak	correctopt	99.99%	0.00%	0.01%	45 kB	8,231
Czechoslovak	sizeopt	99.67%	0.00%	0.33%	40 kB	7,417
Universal	custom	99.10%	0.28%	0.90%	77 kB	11,238
Universal	correctopt	99.83%	0.04%	0.17%	219 kB	29,742
Universal	sizeopt	98.25%	0.01%	1.75%	101 kB	14,321

Table 8: Results of 10-fold cross-validation (learning on 90%, and testing on remaining 10%). Generalization properties (performance on words not seen during training) are compared with Czechoslovak patterns. By adding 7 languages, the generalization abilities of universal patterns are only slightly worse.

Wordlist	Parameters	Good	Bad	Missed
Czechoslovak	custom	99.85%	0.22%	0.15%
Czechoslovak	correctopt	99.95%	0.15%	0.05%
Czechoslovak	sizeopt	99.58%	0.18%	0.42%
Universal	custom	99.04%	1.07%	0.96%
Universal	correctopt	99.37%	1.30%	0.63%
Universal	sizeopt	98.42%	0.95%	1.56%

We removed all colliding words when joining wordlists into the wordlist universal pattern generation. As mentioned earlier, we collected words for nine languages (cz, sk, ka, el, pl, ru, tr, tk, ua).

We generated universal patterns with the same three sets of Patgen parameters (custom, correct optimized, and size optimized) as when generating Czechoslovak patterns. The results are shown in Tables 4 (custom), 5 (correct optimized) and 6 (size optimized). The results are comparable with generation for two languages and confirm the feasibility of universal pattern development.

We did not pursue 100% coverage at all costs because the source data is noisy, and we do not want the patterns to learn all the typos and inconsistencies. Also, the size of the new languages was rather small, compared to Czechoslovak.

4. Evaluation

We evaluated the quality of developed patterns by two metrics. *Coverage* of hyphenation points in the training wordlist tells how the patterns correctly predicted hyphenation points used in training. *Generalization* means how the patterns behave on unseen data, on words not available in the data used during Patgen training. The methodology is the same as we used in the development of Czechoslovak patterns [13].

In Table 7, we compare the efficiency of different approaches to hyphenating 2 languages and 9 languages from one pattern set. We see that the performance of universal patterns is comparable in size and quality to double- or single-language ones—there is only a negligible difference. Table 8 shows that generalization qualities, given the small input size wordlists, are very good, and comparable to the fine-tuned Czechoslovak results. Investing in the purification and consistency of input wordlists (as we did for Czech and Slovak) would result in near-perfect syllabic patterns with almost 100% coverage and no errors.

5. Future work

A natural further step is to merge further languages where the syllabic principle is used for hyphenation. For that, one would need a version of Patgen we provisionally call UniPatgen. This version would support Unicode not only in I/O but also internally as a wide character (UTF-16) character encoded in the pattern representation in either a packed trie or Judy array. This would allow merging more languages *without* increasing the computational complexity of hyphenation, and only a sublinear increase of pattern size. We believe that coverage may differ from 100% only by words that should be hyphenated differently in different languages—our estimate is in small, single-digit percents, while, as mentioned above, the widely-used `hyphen.tex` patterns do not cover 10+%!

Another possible extension in pattern development is the support of a specific hyphenation penalty for compound word borders. This extension, discussed already 30 years ago [17], would generate patterns first for compound words, and only after fixing them continue with pattern generation for all other hyphenation points. The \TeX engine would then set the hyphenation penalties depending on level ranges in patterns found for the hyphenated word. This extension is orthogonal with support for universal patterns but might require increasing the maximal number of levels allowed in patterns to two digits.

There are several open questions for the \TeX development community:

1. Should the universal syllabic patterns ever be developed?

2. If so, should the needed *internal* wide character representations be added to the T_EX suite of programs? That is, to T_EX-based engines not yet supporting it³ and Patgen or UniPatgen.
3. If not, should it be handled by external segmenters on T_EX's input, based on Patgen's proposed successor, UniPatgen?
4. If UniPatgen was developed, should it be added to the distribution, together with Unicode patterns included and supported in repositories like [14]?
5. Should UniPatgen, and LuaT_EX, add a dependency on a Judy library, or should a more conservative solution be sought and implemented? With a conservative solution, which data structure to use for storing patterns? Should the memory be allocated dynamically, to overcome the abundant explosion of format size that stores the patterns, as output by iniT_EX?
6. Should UniPatgen (and T_EX engines) additionally and orthogonally support patterns and a different hyphenation penalty for compound word borders, currently available in e.g. the German wordlist [24]?

We would appreciate qualified opinions on these decisions being sent to authors.

“All we are saying, give patterns a chance.”

Our paraphrase of John Lennon's protest song refrain

6. Conclusion

Preparation of language-agnostic, i.e. universal, syllabic segmentation patterns could be done! We have demonstrated this possibility by generating patterns based on the wordlists of nine languages with current Patgen. They have superb generalization qualities, high coverage of hyphenation points (more than most legacy patterns), and virtually no errors. Their use could have a high impact on virtually all typesetting engines including web page renderers.

Supporting wide characters in Patgen is a critical requirement for adding more languages. We have shown that bringing wide character support into the hyphenation part of the T_EX suite of programs is possible by using the Judy array. It will allow generating and deploying patterns for the whole Unicode character set. We have discussed a possible roadmap to make this a reality in typesetting engines including T_EX successors.

Acknowledgments

We are indebted to Don Knuth for questioning the common properties of Czech and Slovak hyphenation during our presentation of [20] at TUG 2019, which has led us in this research direction. We also thank everyone on whose shoulders we build our work, e.g. for wordlists by Lexical Computing, and to all who commented on our work at TUG 2021 [13] and TUG 2023. We thank ζ TUG for TUG travel support.

³https://cs.overleaf.com/learn/latex/TeX_primitives_listed_by_TeX_engine

References

1. SOJKA, Ondřej; SOJKA, Petr; MÁČA, Jakub. A roadmap for universal syllabic segmentation. *TUGboat*. 2023, vol. 44, no. 2. ISSN 0896-3207. Dostupné také z: <https://doi.org/10.47397/tb/44-2/tb137sojka-syllabic>.
2. HARALAMBOUS, Yannis. New hyphenation techniques in Ω_2 . *TUGboat*. 2006, vol. 27, no. 1, s. 98–103. Dostupné také z: <https://tug.org/TUGboat/tb27-1/tb86haralambous-hyph.pdf>.
3. BARTLETT, Susan; KONDRÁK, Grzegorz; CHERRY, Colin. Automatic Syllabification with Structured SVMs for Letter-to-Phoneme Conversion. In: *Proceedings of ACL-08: HLT*. Columbus, Ohio: Assoc. for Computational Linguistics, 2008, s. 568–576. Dostupné také z: <https://aclweb.org/anthology/P08-1065>.
4. MARCHAND, Yannick; ADSETT, Connie R.; DAMPER, Robert I. Automatic Syllabification in English: A Comparison of Different Algorithms. *Language and Speech*. 2009, vol. 52, no. 1, s. 1–27. Dostupné z DOI: 10.1177/0023830908099881.
5. SHAO, Yan; HARDMEIER, Christian; NIVRE, Joakim. Universal Word Segmentation: Implementation and Interpretation. *Transactions of the Association for Computational Linguistics*. 2018, vol. 6, s. 421–435. Dostupné z DOI: 10.1162/tacl_a_00033.
6. TROGKANIS, Nikolaos; ELKAN, Charles. Conditional Random Fields for Word Hyphenation. In: *Proceedings of the 48th Annual Meeting of the ACL*. Uppsala, Sweden: ACL, 2010, s. 366–374. Dostupné také z: <https://aclweb.org/anthology/P10-1038>.
7. KRANTZ, Jacob; DULIN, Maxwell; PALMA, Paul De. Language-Agnostic Syllabification with Neural Sequence Labeling. *CoRR*. 2019, vol. abs/1909.13362. Dostupné také z: <https://arxiv.org/abs/1909.13362>.
8. THE UNICODE CONSORTIUM. *The Unicode Standard: Worldwide Character Encoding. Version 15.1*. Mountain View, CA, USA: Unicode, Inc., 2023. ISBN 978-1-936213-32-0. Dostupné také z: <https://unicode.org/versions/Unicode15.1.0>.
9. MADDIESON, Ian. Syllable Structure. In: DRYER, Matthew S.; HASPELMATH, Martin (eds.). *The World Atlas of Language Structures Online*. Leipzig: Max Planck Institute for Evolutionary Anthropology, 2013. Dostupné také z: <https://wals.info/chapter/12>.
10. *Online Etymology Dictionary*. "syllable" [online]. [N.d.]. [cit. 2023-07-23]. Dostupné z: <https://www.etymonline.com/word/syllable>.
11. *Internetová jazyková příručka (Internet Language Reference Book)* [online]. 2023. [cit. 2023-07-06]. Dostupné z: <https://prirucka.ujc.cas.cz/?id=135>.

12. LIANG, Franklin M. *Word Hyphenation by Computer*. 1983. Dostupné také z: <https://tug.org/docs/liang/liang-thesis.pdf>. Dis. pr. Stanford University.
13. SOJKA, Petr; SOJKA, Ondřej. New Czechoslovak Hyphenation Patterns, Word Lists, and Workflow. *TUGboat*. 2021, vol. 42, no. 2. ISSN 0896-3207. Dostupné také z: <https://doi.org/10.47397/tb/42-2/tb131sojka-czech>.
14. ROSENDAHL, Arthur; MIKLAVEC, Mojca. *TeX hyphenation patterns*. TUG, 2023. Dostupné také z: <http://hyphenation.org/tex>. Accessed 2023-07-05.
15. JAKUBÍČEK, Miloš; KILGARRIFF, Adam; KOVÁŘ, Vojtěch; RYCHLÝ, Pavel; SUCHOMEL, Vít. The TenTen Corpus Family. In: *Proc. of the 7th International Corpus Linguistics Conference (CL)*. Lancaster, 2013, s. 125–127.
16. KNUTH, Donald E. *3:16 Bible Texts Illuminated*. A-R Editions, Inc., 1991. ISBN 0-89579-252-4.
17. SOJKA, Petr; ŠEVEČEK, Pavel. Hyphenation in TeX—Quo Vadis? *TUGboat*. 1995, vol. 16, no. 3, 280–289. Dostupné také z: <https://tug.org/TUGboat/tb16-3/tb48soj1.pdf>.
18. LIANG, Franklin M.; BREITENLOHNER, Peter. *PATtern GENeration Program for the TeX82 Hyphenator* [Electronic documentation of PATGEN program version 2.4 on CTAN. <https://ctan.org/pkg/patgen>]. 1999.
19. SOJKA, Petr. *Competing Patterns in Language Engineering and Computer Typesetting*. 2005. Dostupné také z: https://researchgate.net/publication/265246931_Competing_Patterns_in_Language_Engineering_and_Computer_Typesetting/. Dis. pr. Faculty of Informatics.
20. SOJKA, Petr; SOJKA, Ondřej. The Unreasonable Effectiveness of Pattern Generation. *TUGboat*. 2019, vol. 40, no. 2, s. 187–193. Dostupné také z: <https://tug.org/TUGboat/tb40-2/tb125sojka-patgen.pdf>.
21. SOJKA, Petr. Competing Patterns for Language Engineering. In: SOJKA, Petr; KOPEČEK, Ivan; PALA, Karel (eds.). *Proceedings of the Third International Workshop on Text, Speech and Dialogue—TSD 2000*. Brno, Czech Republic: Springer-Verlag, 2000, s. 157–162. LNAI 1902. Dostupné z DOI: 10.1007/3-540-45323-7_27.
22. SOJKA, Petr; SOJKA, Ondřej. Towards Universal Hyphenation Patterns. In: HORÁK, Aleš; RYCHLÝ, Pavel; RAMBOUSEK, Adam (eds.). *Proceedings of Recent Advances in Slavonic Natural Language Processing—RASLAN 2019*. Karlova Studánka, Czech Republic: Tribun EU, 2019, 63–68. Dostupné také z: <https://nlp.fi.muni.cz/raslan/2019/paper13-sojka.pdf>. <https://is.muni.cz/publication/1585259/?lang=en>.

23. MÁČA, Jakub. *Judy*. Brno, Czech Republic, 2023. Dostupné také z: <https://is.muni.cz/th/kru3j>. Bachelor Thesis supervised by Petr Sojka and defended at Masaryk University, Faculty of Informatics.
24. LEMBERG, Werner. *A database of German words with hyphenation information*. 2023. Dostupné také z: <https://repo.or.cz/wortliste.git>.

Plán pro univerzální slabičnou segmentaci

Prostorově a časově efektivní segmentace (dělení slov) přirozených jazyků zůstává jádrem každého sázecího systému, ať už jde o $\text{T}_{\text{E}}\text{X}$, webový prohlížeč nebo mobilní operační systém. Ve většině jazyků je dnes pragmaticky preferováno slabičné dělení reflektující výslovnost při čtení.

Vzhledem k tomu, že přepínání jazyků často není v textech označeno, renderovací stroj (webový prohlížeč či $\text{T}_{\text{E}}\text{X}$) potřebuje *univerzální* slabikovou segmentaci. V předloženém článku ukazujeme proveditelnost této myšlenky tím, že nabízíme prototypové řešení dvou hlavních problémů:

- A) použití Patgenu ke generování vzorů pro několik jazyků najednou; a
- B) neexistence podpory Unicode v nástrojích jako Patgen nebo $\text{T}_{\text{E}}\text{X}$ (vzory v kódování UTF-16).

Pro A) jsme ke generování univerzálních slabičných vzorů použili seznamy slov devíti slabičných jazyků (čeština, slovenština, gruzínština, řečtina, polština, ruština, turečtina, turkmenština a ukrajinština). Pro B) jsme vytvořili verzi Patgen, která používá datovou strukturu Judy array, a porovnali její efektivitu s implementací trie.

S údaji z těchto devíti jazyků ukazujeme, že:

- A) vyvinutí univerzálních, obecných slabičných vzorů s vysokým pokrytím je možné, a to s velkým dopadem na prakticky všechny sázecí stroje včetně webových prohlížečů; a
- B) podpora Unicode znaků ve vzorech dělení slov v programech $\text{T}_{\text{E}}\text{X}$ a Patgen je možná pomocí Judy array.

Klíčová slova: slabikování, dělení slov, příprava univerzálních slabičných vzorů

*Ondřej Sojka, Fakulta informatiky Masarykovy univerzity, Brno, ČR
454904 (at) mail dot muni dot cz, ORCID 0000-0003-2048-9977*

*Petr Sojka, Fakulta informatiky Masarykovy univerzity, Brno, ČR
sojka (at) fi dot muni dot cz, ORCID 0000-0002-5768-4007*

*Jakub Máca, Fakulta informatiky Masarykovy univerzity, Brno, ČR
514024 (at) mail dot muni dot cz, ORCID 0009-0008-1583-3183*