

# Zpravodaj Československého sdružení uživatelů TeXu

---

Vít Novotný

Nápadovník jmen pro tvůrčí psaní v LuaTeXu

*Zpravodaj Československého sdružení uživatelů TeXu*, Vol. 33 (2023), No. 1-2, 3–38

Persistent URL: <http://dml.cz/dmlcz/151756>

## 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>

Známy výrok informatika Phila Karltona říká, že na informatice jsou obtížné pouze dvě věci: vyprazdňování cache a přidělování jmen. Svě o tom vědí i spisovatelé, kteří musí kromě příběhu a světa vymyslet jména všech svých příběhových postav. V tomto článku vyvineme jazykový model, který spisovatelům umožní automaticky generovat jména postav při tvůrčím psaní v Lua $\TeX$ U. Kromě pomoci při tvůrčím psaní si představíme i další možná použití jazykových modelů v Lua $\TeX$ U, jako je automatické přepínání vzorů dělení slov podle aktuálního jazyka a generování výplňového textu.  $\TeX$ Nicky zaměřeným čtenářům článek poslouží jako prvotní seznámení s programovacími jazyky Lua a `expl3` a s  $\LaTeX$ ovým balíčkem `xparse` pro přípravu uživatelských rozhraní.

**Klíčová slova:** tvůrčí psaní, trie, jazykové modely, Lua $\TeX$ X, Lua, `expl3`, `xparse`

## Úvod

Se vzestupem aplikací postavených na velkých jazykových modelech od firmy OpenAI stoupá zájem laické veřejnosti o jazykové modelování. Nástroj ChatGPT na požádání napíše báseň, esej i dopis, připraví kód v programovacím jazyce nebo zkontroluje gramatické a ortografické chyby v článku. Vyhledávač Bing provede automatické srovnání cen produktů na základě výsledků vyhledávání, zodpoví faktické dotazy včetně citace zdrojů a naplňuje dovolenou. Velké jazykové modely lze ale snadno zneužít pro hromadné generování dezinformací a generování kódu pro napadení informačních systémů, což přitahuje pozornost regulátorů. Stejně tak jsou velké jazykové modely výpočetně náročné a dostupné pouze přes placené webové rozhraní, což omezuje jejich využitelnost.

V tomto článku se svezeme na vlně zájmu o jazykové modely a vyvineme v jazyce Lua malý jazykový model, který můžeme snadno využít v dokumentech sázených pomocí Lua $\TeX$ U. Náš jazykový model bude sloužit spisovatelům pro generování kreativních jmen postav při tvůrčím psaní, ale uvedeme i další potenciální využití jazykových modelů v  $\TeX$ U, jako je automatické přepínání vzorů dělení slov podle aktuálního jazyka a generování výplňového textu pro přípravu maket.

Nejprve v sekci 1 představíme náš jazykový model pomocí obrázků, příkladů a matematických definic. V sekci 2 na straně 8 popsany model implementujeme v programovacím jazyce Lua a v sekci 3 na straně 16 mu v programovacím jazyce `expl3` a pomocí  $\LaTeX$ ového balíčku `xparse` připravíme uživatelské rozhraní pro snadné využití v dokumentech.  $\TeX$ Nicky zaměřeným čtenářům tato sekce poslouží

jako prvotní seznámení s programovacími jazyky Lua a `expl3` a s  $\text{\LaTeX}$ ovým balíčkem `xparse`. Následně v sekci 4 na straně 26 natrénujeme vyvinutý model na několika příkladových databázích jmen a ukážeme, jak bychom vygenerovali jména postav a zahrnuli je do textu povídky. Čtenáři, kteří se nezajímají o detaily jazykového modelu, mohou začít touto sekci. Nakonec si v sekci 5 na straně 31 uvedeme možná vylepšení jazykového modelu a další aplikace jazykových modelů v  $\text{\TeX}$ u, jako je automatické přepínání vzorů dělení slov podle aktuálního jazyka a generování výplňového textu. V sekci 6 na straně 36 článek uzavřeme shrnutím výsledků článku a jejich praktického přínosu pro čtenáře.

## 1. Jazykový model pro generování jmen postav

### 1.1. Trie (prefixové stromy)

Trie [1] je stromová datová struktura, pomocí které můžeme implementovat asociativní pole s textovými klíči. Na rozdíl od hašové tabulky nám trie umožňuje mj. efektivně vyhledávat všechny klíče buď s danou předponou, nebo příponou,<sup>1</sup> vizte např. Obrázek 1 s triemi, které asociují texty „pes filipes“, „pejsek“ a „maxipes fik“ s obrázky pohádkových psů. Pokud chceme získat jména a obrázky všech pohádkových psů s předponou `pe-`, sestoupíme nejprve do vrcholu  $v_2$  a následně projdeme celý podstrom. Takto získáme jména „pes filipes“ a „pejsek“ spolu s odpovídajícími obrázky. Pro jména a obrázky s příponou `-k` sestoupíme nejprve do vrcholu  $v_9$  a opět projdeme celý podstrom.

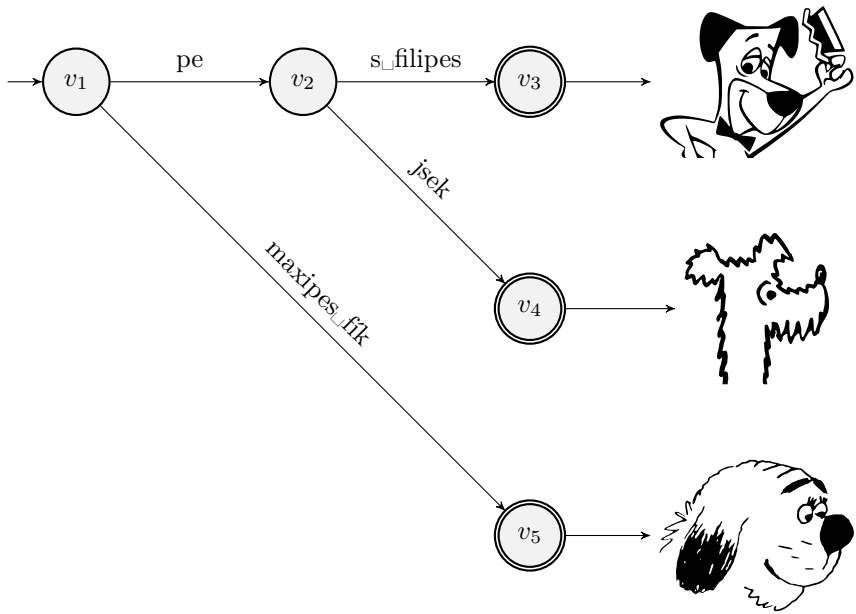
### 1.2. Základní jazykový model

Kromě asociativních polí můžeme trie využít pro implementaci textových množin s efektivním hledáním buď podle předpon, nebo přípon. Příkladem jsou klávesnice v mobilních telefonech, které uživateli našeptávají slova ze slovníku podle rozepsaného slova. Pokud si k hranám trie navíc poznačíme váhu  $c$ , která nám udává, kolikrát jsme hrany během sestavování trie navštívili, poslouží nám trie i jako jazykový model. Na Obrázku 2 vidíme trii, pomocí které můžeme spočítat pravděpodobnost výskytu slova  $S$  jako součin podmíněných pravděpodobností výskytů jednotlivých znaků slova a konce slova:

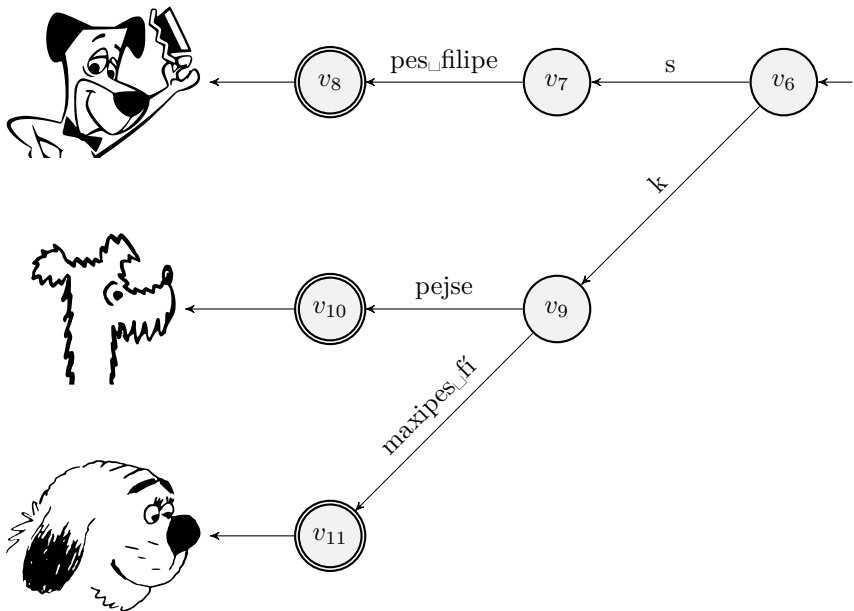
$$P(S) = \left( \prod_{i=1, \dots, |S|} P(S[i] \mid S[1, \dots, i-1]) \right) \cdot P(\epsilon \mid S), \text{ kde } \epsilon \text{ je prázdné slovo.}$$

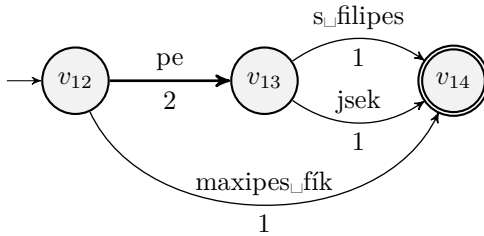
---

<sup>1</sup>Předponami a příponami zde máme na mysli libovolné posloupnosti znaků na začátku a na konci textového klíče. Nejedná se tedy o předpony a přípony ve smyslu nauky o stavbě slov. Kromě vyhledávání textových klíčů buď podle předpon, nebo podle přípon lze v triích vyhledávat i podle kombinace předpony a přípony nebo podle libovolného podřetězce, a to pomocí tzv. *permutermínů*, vizte sekci 5.2.1 na straně 32.



Obrázek 1: Trie, které mapují jména pohádkových psů na jejich obrázky [2, 3, 4]. Trie nahoře umožňuje hledání podle předpon a trie dole hledání podle přípon.





Slovo $S$	$P(S)$
pes filipes	$33,3\%$
pejsek	$33,3\%$
maxipes ffk	$33,3\%$

Obrázek 2: Jazykový model, který nám umožňuje generovat jména pohádkových psů náhodnými procházkami v trii. Tabulka vpravo uvádí všechna jména, která můžeme vygenerovat, a jejich pravděpodobnost.

Pro výpočet podmíněné pravděpodobnosti znaku sestoupíme nejprve z počátečního vrcholu po znacích  $S[1, \dots, i - 1]$  do vrcholu  $v$ . Následně zvolíme hranu  $h$  z vrcholu  $v$  se znakem  $S[i]$ . Pravděpodobnost spočítáme jako podíl váhy  $c$  hrany  $h$  vůči součtu vah všech hran  $H$  vycházejících z vrcholu  $v$ :

$$P(S[i] \mid S[1, \dots, i - 1]) = \frac{c(e)}{\sum_{e' \in E} c(e')}.$$

Např. pravděpodobnost slova „pejsek“ spočítáme následovně:

$$P(\text{pejsek}) = P(\text{pe}|\epsilon) \cdot P(\text{jsek}|\text{pe}) \cdot P(\epsilon|\text{pejsek}) = \frac{2}{3} \cdot \frac{1}{2} \cdot 1 = \frac{1}{3} = 33,3\%.$$

Slova můžeme také vybírat náhodnými procházkami v trii. Takto vypadá náhodná procházka po vrcholech  $v_{12}, v_{13}$  a  $v_{14}$ , která vybere jméno „pes filipes“:

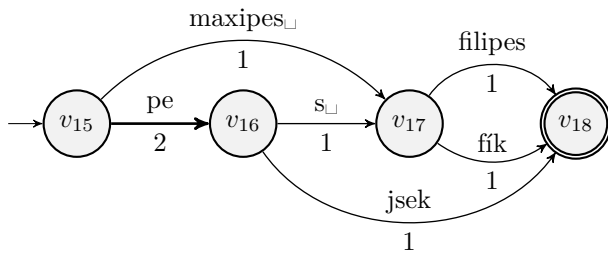
1. Začneme v počátečním vrcholu  $v_{12}$ .
2. S  $66,6\%$  pravděpodobností přejdeme po hraně „pe“ do vrcholu  $v_{13}$ .
3. S  $50\%$  pravděpodobností přejdeme po hraně „s\_filipes“ do vrcholu  $v_{14}$ .
4. Ukončíme náhodnou procházku, protože z vrcholu  $v_{14}$  nevede žádná hrana.

### 1.3. Zapomnětlivý jazykový model

Náš jazykový model sice dokáže vygenerovat slova, pomocí kterých jsme ho sestavili, ale nedovede vymyslet nová slova. Abychom zvýšili kreativitu modelu, upravíme jeho pravděpodobnostní funkci tak, aby místo všech předchozích znaků brala v potaz pouze *kontext* posledních  $n$  znaků a na další znaky „zapomněla“:

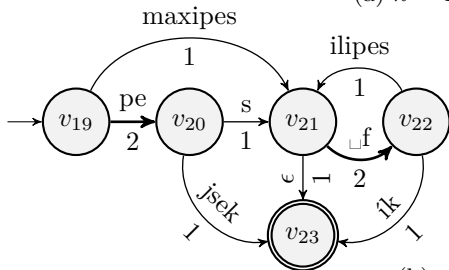
$$P'(S) = \prod_{i=1, \dots, |S|} P(S[i] \mid S[i - n, \dots, i - 1]) \cdot P(\epsilon \mid S[|S| - n, \dots, |S|]).$$

Pro délky kontextu  $n \geq 6$  obdržíme stejný jazykový model jako na Obrázku 2. Jazykové modely pro délky kontextu  $0 < n \leq 5$  najdeme na Obrázku 3. Při  $n = 4$



Slovo $S$	$P(S)$
pejsek	$33,3\bar{3}\%$
pes fík	$16,6\bar{6}\%$
pes filipes	$16,6\bar{6}\%$
maxipes fík	$16,6\bar{6}\%$
maxipes filipes	$16,6\bar{6}\%$

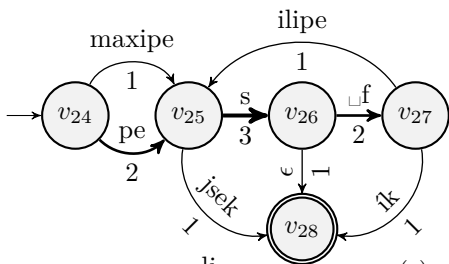
(a)  $n = 4$  a  $n = 5$



Slovo $S$	$P(S)$
pejsek	$33,3\bar{3}\%$
$(\text{maxi})^i \text{pes}(\_ \text{filipes})^j (\_ \text{fík})^k$	$0,3^j \cdot 11,1\bar{1}\%$

$$i \in \{0, 1\}, j \geq 0, k \in \{0, 1\}$$

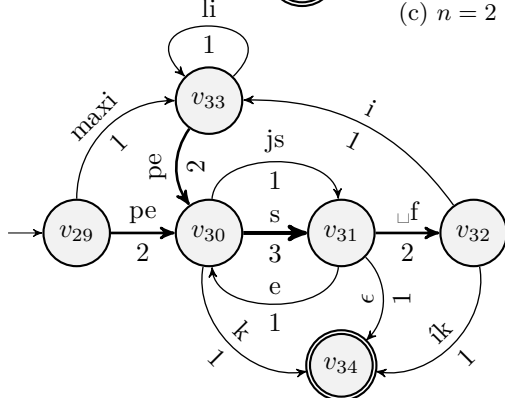
(b)  $n = 3$



Slovo $S$	$P(S)$
$(\text{maxi})^i \text{pe}(\text{s\_filipe})^j \text{jsek}$	$0,5^i \cdot 0,25^j \cdot 16,6\bar{6}\%$
$(\text{maxi})^i \text{pes}(\_ \text{filipes})^j (\_ \text{fík})^k$	$0,25^j \cdot 8,3\bar{3}\%$

$$i \in \{0, 1\}, j \geq 0, k \in \{0, 1\}$$

(c)  $n = 2$



Slovo $S$	$P(S)$
$(\text{maxi}(\text{li})^i)^j \text{p}((\text{ej}^k \text{s})^l \_ \text{fi}(\text{li})^u \text{p})^v (\text{ej}^w \text{s})^x (\_ \text{fík})^y$	$0,3^{i+j+klv+uv+wx} \cdot 0,6^{lv+wx} \cdot 0,16^v \cdot 16,6\bar{6}\%$
$(\text{maxi}(\text{li})^i)^j \text{p}((\text{ej}^k \text{s})^l \_ \text{fi}(\text{li})^u \text{p})^v (\text{ej}^w \text{s})^x \text{ek}$	$0,3^{ij+klv+uv+wx} \cdot 0,6^{lv+wx} \cdot 0,16^v \cdot 3,3\bar{3}\%$

$$i \geq 0, j \in \{0, 1\}, k \in \{0, 1\}, l \geq 0, u \geq 0, v \geq 0, w \in \{0, 1\}, x \geq 1, y \in \{0, 1\}$$

(d)  $n = 1$

Obrázek 3: Zapomnětlivé jazykové modely pro různé velikosti paměti  $n$ , které nám umožňují generovat nová slova a slovní spojení podobná jménům pohádkových psů. Tabulka vpravo uvádí všechny možné výstupy modelu a jejich pravděpodobnost.

a  $n = 5$  dokáže jazykový model vymyslet nová slovní spojení „pes fík“ a „maxipes filipes“. Při  $n = 3$  dostáváme nové slovo „pes“ a v grafu nám vzniká cyklus mezi vrcholy  $v_{21}$  a  $v_{22}$ , který s nízkou pravděpodobností generuje slovní spojení libovolné délky jako „maxipes filipes filipes fík“ s pravděpodobností přibližně 1,23 %. Při  $n = 2$  dokážeme vygenerovat jména jako „maxipes filipejsek“, která obsahují nová odvozená slova. Při  $n = 1$  vznikají v modelu krátké cykly nad vrcholem  $v_{33}$  a mezi vrcholy  $v_{30}$  a  $v_{31}$ , díky kterým dokážeme vygenerovat nová jména s opakovanými slabikami jako „maxilipes fililipes fík“ a „peses filipejsejsek“.

#### 1.4. Katzův couvající jazykový model

Při generování nových jmen můžeme náhodné procházky uvozovat buď požadovanou předponou, nebo příponou. Při dlouhých kontextech ale bude mít většina předpon a přípon nulovou pravděpodobnost, zatímco při krátkých kontextech budou generovaná jména vypadat nahodile. Pokud chceme např. vygenerovat jméno pohádkového psa s předponou *i-*, selžeme pro všechny hodnoty  $n \geq 2$  a při  $n = 1$  získáme nepravděpodobná jména jako „ipek“ a „ililipejsej“.

Možným řešením popsaného problému je Katzův jazykový model [5], který během náhodné procházky „couvá“ s délkou kontextu  $n$  na nejvyšší takovou hodnotu, při které má aktuální předpona nebo přípona nenulovou pravděpodobnost. Náhodná procházka s délkou kontextu  $n = 4$  a předponou *i-* v Katzově modelu neselže a vygeneruje pouze strážlivá jména jako „ilipes“, „ipes fík“ a „ipes filipes“.

## 2. Implementace jazykového modelu

V této sekci rozpracujeme soubor `randomnames.lua`, který bude obsahovat implementaci našeho jazykového modelu v jazyce Lua. Nejprve se v sekci 2.1 zaměříme na sestavení modelu. Následně v sekci 2.2 na straně 12 uděláme malou odbočku a naprogramujeme vlastní generátor pseudonáhodných čísel, abychom mohli nakonec v sekci 2.3 na straně 14 implementovat (pseudo)náhodné procházky v našem modelu. Hotový soubor `randomnames.lua` můžeme stáhnout online [6].

Příklady z této sekce uvozené dosavadním obsahem souboru `randomnames.lua` si můžeme vždy uložit do textového souboru `priklad.lua` a přeložit příkazem `texlua prikklad.lua`.

### 2.1. Sestavení jazykového modelu

Pro sestavení zapomnětlivého jazykového modelu potřebujeme znát pouze délku kontextu  $n$  a množinu trénovacích slov.

```
1 local ForgetfulModel = {}
```

```
2
```

```

3  function ForgetfulModel.new(context_size)
4      local model = {}                -- Vytvoř prázdný objekt.
5      model.context_size = context_size -- Ulož si délku kontextu
6      model.forward_graph = {}        -- a prázdný dopředný
7      model.reverse_graph = {}        -- a zpětný jazykový model.
8      local mt = {                    -- Zděď metody2
9          __index = ForgetfulModel }   -- třídy ForgetfulModel.
10     setmetatable(model, mt)
11     return model
12 end
13
14 local function trim_context(context,      -- Ořízni kontext
15                          max_context_size) -- na danou délku.3
16     local context_size = utf8.len(context)
17     if context_size > max_context_size then
18         local character_offset = context_size - max_context_size + 1
19         local byte_offset = utf8.offset(context, character_offset)
20         context = context:sub(byte_offset, #context)
21     end
22     return context
23 end
24
25 local function add_name(graph, name, max_context_size)
26     name = name                      -- Konec jména kóduj
27         .. "\n"                      -- jako konec řádku.
28     local context = ""
29     for _, code in utf8.codes(name) do -- Procházej znaky jména.
30         if graph[context] == nil then -- Pokud je třeba,
31             graph[context] = {}       -- přidej nový vrchol
32         end                             -- pro současný kontext.
33         local vertex = graph[context]
34         local character = utf8.char(code)
35         if vertex[character] == nil then -- Pokud je třeba,
36             vertex[character] = 0     -- přidej novou hranu
37         end                             -- pro současný znak.

```

---

<sup>2</sup>Jazyk Lua používá tzv. *prototypální dědičnost* pomocí *metatabulek* [7, sekce 20.4], která je obecnější než třídní dědičnost, kterou známe z objektově orientovaných jazyků, jako jsou C++ a Java. Zde jsme vytvořili prázdnou hašovou tabulku `model` s metatabulkou `{ __index = ForgetfulModel }`, která říká, že při přístupu k `modelu` se mají neznámé klíče hledat v hašové tabulce `ForgetfulModel`. To nám umožňuje volat metody jako `model.add_name(model, name)` nebo krátce `model:add_name(name)`.

<sup>3</sup>Unární operátor mřížky (`#`) v jazyce Lua vrací délku pole nebo textového řetězce.



```

38     end
39     vertex[character] =                               -- Navyš hodnotu hrany
40         vertex[character] + 1                         -- pro současný znak.
41     context = trim_context(                           -- Aktualizuj kontext
42         context .. character,                         -- a seřizni ho na
43         max_context_size)                            -- požadovanou délku.
44     end
45 end
46
47 local function reverse(name)                          -- Ulož znaky jména
48     local index = utf8.len(name)                     -- v obráceném pořadí.
49     local result = {}
50     for _, code in utf8.codes(name) do
51         local character = utf8.char(code)
52         result[index] = character
53         index = index - 1
54     end
55     local reversed_name = table.concat(result)
56     return reversed_name
57 end
58
59 function ForgetfulModel.add_name(model, name)
60     add_name(model.forward_graph,                    -- Zaznamenej jméno
61         name,                                       -- zepředu i pozpátku
62         model.context_size)                         -- pro náhodně procházky
63     add_name(model.reverse_graph,                   -- vvozené buď
64         reverse(name),                             -- požadovanou předponou,
65         model.context_size)                         -- nebo příponou.
66 end
67
68 function ForgetfulModel.input_names(model, filename)
69     local file = io.open(filename, "r")             -- Procházej řádky
70     assert(file)                                    -- zadaného textového
71     for name in file:lines() do                     -- souboru se jmény.
72         if #name:match("~%s*(.)%s*$") == 0        -- Přeskoč prázdné řádky.
73             then
74                 goto continue
75             end
76         model:add_name(name)                         -- Zanes jméno do modelu.
77         ::continue::
78     end
79 end

```

Katzůvouvající jazykový model sestavíme ze zapomnětlivých modelů pro délky kontextu  $n = i, i + 1, \dots, j - 1, j$ , kde hraniční hodnoty  $i, j$  zadává uživatel.

```
80 local BackoffModel = {}
81
82 function BackoffModel.new(min_context_size, max_context_size)
83     local model = {}           -- Vytvoř prázdný objekt.
84     model.min_context_size =   -- Ulož si hraniční hodnoty
85         min_context_size       -- délky kontextu.
86     model.max_context_size =
87         max_context_size
88     model.submodels = {}       -- Vytvoř zapomnětlivé modely.
89     for context_size = min_context_size, max_context_size do
90         model.submodels[context_size]
91             = ForgetfulModel.new(context_size)
92     end
93     local mt = {               -- Zděd metody
94         __index = BackoffModel } -- třídy BackoffModel.
95     setmetatable(model, mt)
96     return model
97 end
98
99 function BackoffModel.for_each_submodel(model, func,
100                                     min_context_size,
101                                     max_context_size)
102     assert(min_context_size >= model.min_context_size)
103     assert(max_context_size <= model.max_context_size)
104     local i, j = min_context_size, max_context_size
105     for context_size = j, i, -1 do -- V pořadí klesajících
106         local submodel           -- délky kontextu uplatni
107             = model.submodels[   -- funkci nad všemi
108                 context_size]    -- zapomnětlivými modely.
109         local result = func(submodel) -- Pokud funkce vrátí hodnotu
110         if result == false then     -- false, přeruš zpracování
111             break                 -- předčasně.
112         end
113     end
114 end
115
116 function BackoffModel.add_name(model, name)
117     model:for_each_submodel(
118         function(submodel)
```

```

119     submodel:add_name(name)           -- Zanes jméno do všech
120 end,                                  -- zapomnětlivých modelů.
121 model.min_context_size,
122 model.max_context_size)
123 end
124
125 function BackoffModel.input_names(model, filename)
126     model:for_each_submodel(
127         function(submodel)
128             submodel:input_names(     -- Zanes jména do všech
129                 filename)           -- zapomnětlivých modelů.
130         end,
131         model.min_context_size,
132         model.max_context_size)
133 end

```

Katzův couvající jazykový model z Obrázku 3 pro délky kontextu  $n = 1, 2, \dots, 5$  sestavíme v Lua<sub>T<sub>E</sub>X</sub>u následně:

```

model = BackoffModel.new(1, 6)         -- Vytvoř prázdný model.
model:add_name("pes filipes")         -- Zanes do modelu jména
model:add_name("pejsek")             -- pohádkových psů.
model:add_name("maxipes fík")

```

## 2.2. Generátor pseudonáhodných čísel

Pro náhodné procházky v jazykovém modelu potřebujeme hrací kostku, která nám řekne, po které hraně grafu se máme vydat. Jazyk Lua poskytuje vestavěný generátor pseudonáhodných čísel, avšak ten má několik slabín: Pro generování čísel slouží funkce `math.random()`, která generuje různá čísla na různých platformách. My ale chceme, aby náš model generoval stejná náhodná jména na všech aktuálních instalacích <sub>T<sub>E</sub>X</sub>u. Pro resetování vestavěného generátoru pomocí náhodného semínka slouží funkce `math.randomseed()`. Naším záměrem však je současně používat několik jazykových modelů se vzájemně nezávislými generátory. Uvedené problémy vyřešíme implementací vlastního generátoru pseudonáhodných čísel.

Jazyk Lua 5.3 používá alespoň 32bitová celá čísla se znaménkem a neurčitou endianitou [8, sekce 2.1]. Pro nezávislost na platformě tedy musíme zvolit algoritmus, který nepoužívá bitové operace ani jinak nezávisí na endianitě a který pracuje pouze s čísly v rozmezí od 0 do  $2^{31} - 1$ . Takovým algoritmem je například multiplikativní lineární kongruenční generátor [9, sekce 7.1]:

$$x_n = a \cdot x_{n-1} \pmod{m},$$

kde  $x_1 \neq 0 \pmod{m}$  je náhodné semínko,  $x_n$  je  $n$ -té pseudonáhodné číslo a  $a, m$

jsou tabulkové hodnoty. Pro dosažení nezávislosti na platformě požadujeme, aby mezivýsledky nepřesáhly hodnotu  $2^{31} - 1$ . Toho dosáhneme vhodným výběrem hodnot  $a = 771\,645\,345$ ,  $m = 2^{30} - 35 = 1\,073\,741\,789$  [10, tabulka 2], omezením semínka na hodnoty od 1 do  $m - 1$  a použitím algoritmu pro modulární násobení.

```

138 local Random = { A = 771645345, M = 1073741789 }
139
140 function Random.new(seed)
141     assert(seed >= 0 and                    -- Povol semínko z rozsahu
142            seed <= 2147483647)             -- od 0 do  $2^{31} - 1$ .
143     seed = seed                               -- Převed' semínko do rozsahu
144         % (Random.M - 1) + 1                 -- od 1 do  $m - 1$ .
145     local random = { x = seed }              -- Vytvoř objekt se semínkem.
146     local mt = { __index = Random }         -- Zděd' metody třídy Random.
147     setmetatable(random, mt)
148     return random
149 end
150
151 local function multiply_modulo(a, x, m)
152     local result = 0                          -- Algoritmus pro modulární
153     a = a % m                                  -- násobení zajistí, že
154     while x > 0 do                             -- mezivýsledky nepřesáhnou
155         if x % 2 == 1 then                     -- hodnotu  $2 * m$ .
156             result = (result + a) % m
157         end
158         a = (a * 2) % m
159         x = x // 2
160     end
161     return result
162 end
163
164 function Random.get_next_number(random, from, to)
165     random.x = multiply_modulo(Random.A, random.x, Random.M)
166     local result = from                        -- Navrať celé číslo
167         + (random.x                            -- v zadaném rozsahu.
168            % (to - from + 1))
169     return result
170 end
171

```

Generátor můžeme v LuaTeXu použít následně:

```

random = Random.new(42)      -- Vytvoř generátor se semínkem 42.
for i = 1, 30 do            -- Vypiš 30 hodů 6stěnnou kostkou.

```

```

    tex.sprint(random:get_next_number(1, 6))
    if i < 30 then tex.sprint(", ") end
end
4, 2, 1, 2, 3, 1, 1, 4, 4, 2, 4, 5, 3, 4, 6, 3, 1, 5, 4, 2, 1, 3, 3, 4, 6, 5, 5, 2, 6, 6

```

### 2.3. Náhodné procházky

Náhodné procházky Katzova couvajícího modelu sestávají z elementárních náhodných kroků jednotlivých zapomnětlivých modelů.

```

177 local function take_random_step(graph, context, random)
178     if graph[context] == nil then -- Pokud neexistuje vrchol pro
179         return nil -- současný kontext, selži.
180     end
181     local vertex = graph[context] -- Vyber vrchol pro současný
182     assert(#graph[context] > 0) -- kontext.
183     local total_weight = 0 -- Sečti hodnoty odchozích hran.
184     local weight
185     for _, character in ipairs(vertex) do
186         weight = vertex[character]
187         total_weight = total_weight + weight
188     end
189     local random_number = -- Vyber náhodnou hranu.
190         random:get_next_number(0, total_weight)
191     local weight_accumulator = 0
192     for _, character in ipairs(vertex) do
193         weight = vertex[character]
194         weight_accumulator = weight_accumulator + weight
195         if weight_accumulator >= random_number then
196             return character -- Navrať znak na hraně.
197         end
198     end
199     assert(false) -- Sem bychom se neměli dostat.
200 end
201
202 function ForgetfulModel.take_random_step(model, graph_type,
203     context, random)
204     local graph = model[ -- Umožni náhodné procházky
205         graph_type .. "_graph"] -- zepředu i pozpátku.
206     assert(graph ~= nil)
207     context = trim_context( -- Seřídni kontext na
208         context, -- požadovanou délku.
209         model.context_size)

```

```

210     return take_random_step(graph, context, random)
211 end
212
213 function BackoffModel.take_random_walk(model, graph_type,
214                                         affix, random,
215                                         min_context_size,
216                                         max_context_size)
217     local name = affix -- Umožni zadat buď předponu,
218     if graph_type == "reverse" -- nebo příponu jména.
219         then name = reverse(name)
220     end
221     while true do
222         local character
223         model:for_each_submodel(
224             function(submodel)
225                 character = submodel:take_random_step(
226                     graph_type, name, random)
227                 if character ~= nil then -- Najdi zapomnětlivý
228                     return false -- model s nejvyšší
229                     end -- délkou kontextu, ve
230                     -- kterém existuje vrchol
231                     min_context_size, -- pro současný kontext.
232                     max_context_size)
233                 if character == nil then -- Pokud takový model
234                     return nil -- neexistuje, selží.
235                 end
236                 if character == "\n" then -- Pokud jsme vygenerovali
237                     break -- znak konce řádku,
238                     end -- ukončíme procházku.
239                 name = name .. character
240             end
241             if graph_type == "reverse" -- Při náhodné procházce
242                 then name = reverse(name) -- pozpátku otoč
243             end -- vygenerované jméno.
244             return name -- Navrať vygenerované jméno.
245         end

```

Náhodné procházky provedeme v LuaTeXu následně:

```

for i = 1, 5 do -- Vypiš výsledky 5 procházek s předponou pe- se
    tex.sprint( -- zapomnětlivým modelem s délkou kontextu 6.
        model:take_random_walk("forward", "pe", random, 6, 6))
    if i < 5 then tex.sprint(", ") end

```

end

*pejsek, pes filipes, pes filipes, pejsek, pes filipes*

```
for i = 1, 20 do -- Vypiš výsledky 20 procházek s příponou -k se
  tex.sprint( -- zapomnětlivým modelem s délkou kontextu 6.
    model:take_random_walk("reverse", "k", random, 6, 6))
  if i < 20 then tex.sprint(", ") end
end
```

end

*pejsek, pejsek, pejsek, pejsek, pejsek, pejsek, pejsek, pejsek, pejsek, pejsek, pejsek,
pejsek, maxipes fík, pejsek, pejsek, pejsek, pejsek, pejsek, maxipes fík, pejsek*

```
for i = 1, 5 do -- Vypiš výsledky 5 procházek bez afixu se
  tex.sprint( -- zapomnětlivým modelem s délkou kontextu 3.
    model:take_random_walk("forward", "", random, 3, 3))
  if i < 5 then tex.sprint(", ") end
end
```

end

*pes filipes, pes filipes, maxipes fík, pes filipes filipes filipes, pes*

```
for i = 1, 10 do -- Vypiš výsledky 10 procházek bez afixu se
  tex.sprint( -- zapomnětlivým modelem s délkou kontextu 2.
    model:take_random_walk("forward", "", random, 2, 2))
  if i < 10 then tex.sprint(", ") end
end
```

end

*pes filipes filipes, pes filipes, maxipejsek, maxipes filipes filipes fík, maxipes fík,
maxipes filipes filipejsek, pes, pes fík, pes filipejsek, pes*

```
for i = 1, 10 do -- Vypiš výsledky 10 procházek bez afixu se
  tex.sprint( -- zapomnětlivým modelem s délkou kontextu 1.
    model:take_random_walk("forward", "", random, 1, 1))
  if i < 10 then tex.sprint(", ") end
end
```

end

*pes, pes fipes fililipes, maxipejs fík, pes fipes fipes fipek, pek, pes fipejs, pes filililipes
filipek, maxilipes filipesejs fililililipejs fipek, maxipes fipejseses, maxipes*

```
for i = 1, 5 do -- Vypiš výsledky 5 procházek s předponou i- nad
  tex.sprint( -- Katzovým modelem s délkami kontextu 1 až 4.
    model:take_random_walk("forward", "i", random, 1, 4))
  if i < 5 then tex.sprint(", ") end
end
```

end

*ipes filipes filipes filipes, ipes filipes fík, ilipes filipes, ilipes, ilipes*

### 3. Uživatelské rozhraní pro vkládání jmen do dokumentů

V této sekci ukončíme soubor `randomnames.lua` a vytvoříme dva nové soubory `randomnames.tex` a `randomnames.sty`, které budou obsahovat uživatelské roz-

hraní pro snadné využití našeho jazykového modelu v  $\text{\LaTeX}$ ových dokumentech. Hotové soubory `randomnames.lua`, `randomnames.tex` a `randomnames.sty` můžeme stáhnout online [6].

Nejprve v sekci 3.1 doplníme do souboru `randomnames.lua` fasádu, která poslouží jako programátorské rozhraní pro naši implementaci v jazyce Lua. Následně v sekci 3.2 na straně 19 vytvoříme soubor `randomnames.tex` s přemostěním mezi Luou a  $\text{\TeX}$ em pomocí programovacího jazyka `expl3`. Nakonec v sekci 3.3 na straně 25 vytvoříme soubor `randomnames.sty` s uživatelským rozhraním pro snadné využití našeho jazykového modelu v  $\text{\LaTeX}$ ových dokumentech.

### 3.1. Fasáda pro vytváření jazykových modelů v jazyce Lua

Dosud jsme náš jazykový model a generátor pseudonáhodných čísel ukládali do proměnných `model` a `random`. Při práci s mnoha modely a generátory je ale použití proměnných nepřehledné. V této sekci vytvoříme fasádu, která nám umožní vytvářet pojmenované modely s přidruženými generátory a snadno k nim přistupovat bez použití proměnných.

```
276 local randomnames = {}      -- Veřejné rozhraní našeho Lua modulu.
277
278 local models = {}          -- Vytvoř prázdnou databázi modelů
279 local randoms = {}        -- a generátorů pseudonáhodných čísel.
280
281 local function hash(name)  -- Vypočti ze jména modelu počáteční
282     local result = 0        -- semínko s hašovací funkcí sdbm.
283     local modulus = Random.M - 1
284     for _, code in utf8.codes(name) do
285         result = multiply_modulo(result, 65599, modulus)
286         result = (result + code) % modulus
287     end
288     return result
289 end
290
291 function randomnames.new_model(name, min_context_size,
292                               max_context_size, seed)
293     assert(models[name] == nil,
294            [[Model named "]] .. name .. [[ already exists]])
295     min_context_size =      -- Pokud uživatel nedodal minimální
296         min_context_size or 1 -- nebo maximální délku kontextu,
297     max_context_size =      -- použij výchozí hodnoty 1 a 3.
298         max_context_size or 3
299     local model =          -- Vytvoř nový jazykový model.
300         BackoffModel.new(min_context_size, max_context_size)
```



```

301     if seed == nil then           -- Pokud uživatel nedodal náhodné
302         seed = hash(name)       -- semínko, spočítej ho jako haš
303     end                           -- jména jazykového modelu.
304     local random =              -- Vytvoř nový generátor
305         Random.new(seed)        -- pseudonáhodných čísel.
306     models[name] = model        -- Ulož model a generátor do
307     randoms[name] = random      -- databáze pod zadaným jménem.
308 end
309
310 local function get_model(name)
311     local model
312     model = models[name]        -- Zkus načíst z databáze model.
313     assert(model ~= nil,       -- Pokud neexistuje, vypiš chybu.
314         [[Model named "]] .. name .. [[ does not exist]])
315     return model
316 end
317
318 function randomnames.add_name(model_name, name)
319     local model =               -- Načti z databáze model
320         get_model(model_name)  -- a zanes do něj jméno.
321     model:add_name(name)
322 end
323
324 function randomnames.input_names(model_name, filename)
325     local model =               -- Načti z databáze model
326         get_model(model_name)  -- a zanes do něj jména.
327     model:input_names(filename)
328 end
329
330 function randomnames.take_random_walk(name, graph_type,
331                                     affix, seed,
332                                     min_context_size,
333                                     max_context_size)
334     local model =               -- Načti z databáze model.
335         get_model(name)
336     graph_type = graph_type     -- Pokud uživatel nedodal směr
337         or "forward"           -- procházky, jdi zepředu.
338     affix = affix              -- Pokud uživatel nedodal předponu
339         or ""                  -- ani příponu, použij prázdnou.
340     local random                -- Pokud uživatel nedodal náhodné
341     if seed == nil then        -- semínko, načti z databáze
342         random = randoms[name] -- generátor náhodných čísel

```

```

343 else                                     -- přidružený k modelu.
344     random =                             -- Jinak vytvoř nový generátor
345     Random.new(seed)                     -- pomocí náhodného semínka.
346 end
347 if min_context_size                       -- Pokud uživatel nedodal
348     == nil then                           -- minimální nebo maximální
349     min_context_size =                     -- délku kontextu, použij
350     model.                                 -- hodnoty z modelu.
351     min_context_size
352 end
353 if max_context_size == nil then
354     max_context_size = model.max_context_size
355 end
356 local result =                            -- Proveď náhodnou procházku.
357     model:take_random_walk(graph_type, affix, random,
358                             min_context_size,
359                             max_context_size)
360 return result
361 end

```

Pomocí naší fasády vytvoříme Katzův couvající jazykový model a provedeme v něm náhodnou procházku v LuaTeXu následně:

```

randomnames.new_model("pes", 1, 3, 22)      -- Vytvoř model.
randomnames.add_name("pes", "pes filipes")  -- Zanes do modelu
randomnames.add_name("pes", "pejsek")      -- jména
randomnames.add_name("pes", "maxipes fík") -- pohádkových psů.
tex.print(randomnames.take_random_walk("pes")) -- Vygeneruj jméno.
maxipes filipes fík

```

Na konci souboru `randomnames.lua` navrátíme naši fasádu:

```

367 return randomnames

```

Poté můžeme v LuaTeXu načíst fasádu odkudkoliv příkazem `require()`:

```

local randomnames = -- Načti fasádu ze souboru randomnames.lua.
require("randomnames")

```

### 3.2. Přemostění mezi jazyky Lua a TeX pomocí jazyka expl3

Abychom mohli jazykové modely používat přímo z TeXu, budeme nyní chtít napojit naši fasádu na TeXová makra. Na rozdíl od jazyka Lua, ve kterém můžeme použít speciální prázdnou hodnotu `nil` pro nepovinné parametry, nemáme v jazyce TeX přímočarý způsob, jak část parametrů funkce vynechat. Využijeme proto datový typ `key-value` programovacího jazyka `expl3` [11, kapitola 26], která nám při volání funkcí umožní uvést pouze část nepovinných parametrů.

```

1 \ifx \ExplSyntaxOn \undefined           % Načti programovací jazyk
2   \input expl3-generic                   % expl3 do plain TeXu.
3 \fi
4 \ExplSyntaxOn
5
6 \tl_new:N                               % Vytvoř lokální proměnné
7   \l_randomnames_min_context_size_tl    % pro nepovinné parametry
8 \tl_new:N                               % funkce randomnames.
9   \l_randomnames_max_context_size_tl    % new_model().
10 \tl_new:N \l_randomnames_seed_tl
11 \tl_gset:Nn \l_randomnames_min_context_size_tl { nil }
12 \tl_gset:Nn \l_randomnames_max_context_size_tl { nil }
13 \tl_gset:Nn \l_randomnames_seed_tl { nil }
14
15 \keys_define:nn                         % Vytvoř sběrnici,4 která
29   { randomnames / new_model }          % naplní lokální proměnné
30   {                                     % zadanými hodnotami
31     min_context_size .code:n =         % nepovinných parametrů
32     {                                   % ve formátu požadovaném
33       \tl_set:Nn                         % jazykem Lua.
34         \l_randomnames_min_context_size_tl
35         { tonumber(" \lua_escape:e { #1 } ") }
36     },
37     max_context_size .code:n =
38     {

```

---

<sup>4</sup> V programovacím jazyku expl3 můžeme s datovým typem key–value nakládat dvěma způsoby. Mějme např. key–value { min\_context\_size = 3, seed = 42 }. Tento key–value můžeme uložit do hašové tabulky [11, kapitola 24] a následně přistupovat k hodnotám jednotlivých klíčů:

```

\prop_set_from_keyval:Nn                 % Vytvoř v pomocné proměnné
  \l_tmpa_prop                           % \l_tmpa_prop hašovou tabulku
  { min_context_size = 3, seed = 42 }     % a naplní ji naším key–value.
Náhodné-semínko-je:-

```

```

\prop_item:Nn                             % Získej hodnotu klíče seed
  \l_tmpa_prop                             % z naší hašové tabulky.
  { seed } .

```

Alternativně můžeme vytvořit datovou sběrnici [11, kapitola 26], u které nejprve nastavíme způsob zpracování hodnot jednotlivých klíčů a následně do ní náš key–value vpustíme:

```

\keys_define:nn                           % Vytvoř datovou sběrnici, která
  { randomnames / process_key_value }     % při zpracování klíče seed
  { seed .code:n = { Náhodné-semínko-je:~#1. } } % vypíše jeho hodnotu.
\keys_set:nn                               % Vpusť do datové sběrnice náš
  { randomnames / process_key_value }     % key–value.
  { min_context_size = 3, seed = 42 }

```

Oba příklady vypíší text „Náhodné semínko je 42.“. V našem přemostění používáme druhý popsáný způsob s datovou sběrnici.

```

39     \tl_set:Nn
40     \l_randomnames_max_context_size_tl
41     { tonumber(" \lua_escape:e { #1 } ") }
42   },
43   context_size .meta:n =           % Umožni snadno vytvářet
44   {                                 % zapomnětlivé modely
45     min_context_size = { #1 },     % s pevnou délkou kontextu,
46     max_context_size = { #1 },     % které nemohou couvat.
47   },
48   seed .code:n =
49   {
50     \tl_set:Nn
51     \l_randomnames_seed_tl
52     { tonumber(" \lua_escape:e { #1 } ") }
53   },
54 }
55
56 \cs_new:Nn
57 \randomnames_new_model:nn
58 {
59   \group_begin:
60   \keys_set:nn                     % Naplň lokální proměnné
61   { randomnames / new_model }     % zadanými hodnotami
62   { #1 }                           % nepovinných parametrů.
63   \lua_now:e                       % Zavolej funkci
64   {                                 % randomnames.new_model()
65     local-randomnames =           % v jazyce Lua.5
66     require("randomnames")
67     randomnames.new_model(
68     " \lua_escape:e { #2 } ",
69     \l_randomnames_min_context_size_tl,
70     \l_randomnames_max_context_size_tl,
71     \l_randomnames_seed_tl
72   )
73 }
74 \group_end:                       % Resetuj hodnoty lokálních
75 }                                 % proměnných zpět na nil.
76

```

---

<sup>5</sup>Znak vlnovky (~) v jazyce expl3 vkládá mezeru, která nám v kódu v jazyce Lua odděluje klíčové slovo `local` od názvu proměnné `randomnames`. Pokud bychom vlnovku neuvdli, definovali bychom místo lokální proměnné `randomnames` globální proměnnou `localrandomnames`.

```

77 \cs_new:Nn
78   \randomnames_add_name:nn
79   {
80     \lua_now:e                               % Zavolej funkci
81     {                                         % randomnames.add_name()
82       local~randomnames =                   % v jazyce Lua.
83         require("randomnames")
84       randomnames.add_name(
85         " \lua_escape:e { #1 } ",
86         " \lua_escape:e { #2 } "
87       )
88     }
89   }
90
91 \cs_new:Nn
92   \randomnames_input_names:nn
93   {
94     \lua_now:e                               % Zavolej funkci
95     {                                         % randomnames.
96       local~randomnames =                   % input_names() v jazyce
97         require("randomnames")              % Lua.
98       randomnames.input_names(
99         " \lua_escape:e { #1 } ",
100        " \lua_escape:e { #2 } "
101      )
102    }
103  }
104
105 \tl_new:N                               % Vytvoř lokální proměnné
106   \l_randomnames_graph_type_tl          % pro nepovinné parametry
107 \tl_new:N                               % funkce randomnames.
108   \l_randomnames_affix_tl              % take_random_walk().
109 \tl_gset:Nn \l_randomnames_graph_type_tl { nil }
110 \tl_gset:Nn \l_randomnames_affix_tl { nil }
111
112 \bool_new:N                             % Vytvoř proměnné, pomocí
113   \l_randomnames_save_bool             % kterých bude uživatel moci
114 \bool_gset_false:N                      % ukládat jména, jež využívá
115   \l_randomnames_save_bool             % opakovaně.
116 \tl_new:N \l_randomnames_save_tl
117
118 \keys_define:nn                         % Vytvoř sběrnici, která

```

```

119 { randomnames / take_random_walk } % naplní lokální proměnné
120 { % zadanými hodnotami
121 graph_type .code:n = % nepovinných parametrů
122 { % ve formátu požadovaném
123 \tl_set:Nn % jazykem Lua.
124 \l_randomnames_graph_type_tl
125 { " \lua_escape:e { #1 } " }
126 },
127 affix .code:n =
128 {
129 \tl_set:Nn
130 \l_randomnames_affix_tl
131 { " \lua_escape:e { #1 } " }
132 },
133 save .code:n = % Umožni ukládat
134 { % vygenerovaná jména
135 \bool_set_true:N % pro opakované použití.
136 \l_randomnames_save_bool
137 \tl_set:Nn \l_randomnames_save_tl { #1 }
138 },
139 prefix .meta:n = % Umožni zadávat předpony
140 { % a přípony náhodných
141 graph_type = { forward }, % procházek tak, aby se
142 affix = { #1 }, % zároveň automaticky
143 }, % nastavil správný směr
144 suffix .meta:n = % náhodně procházky.
145 {
146 graph_type = { reverse },
147 affix = { #1 },
148 },
149 }
150 \keys_define:nn % Poděť zpracování
151 { randomnames } % nepovinných parametrů
152 { % funkce randomnames.
153 take_random_walk .inherit:n = % new_model(), abychom
154 { randomnames / new_model }, % se vyhnuli duplikaci
155 } % kódu.
156
157 \cs_new:Nn
158 \randomnames_take_random_walk:nn
159 {
160 \group_begin:

```

```

161 \keys_set:nn                               % Naplň lokální proměnné
162 {                                           % zadanými hodnotami
163     randomnames /                           % nepovinných parametrů.
164     take_random_walk
165 }
166 { #1 }
167 \lua_now:e                                 % Zavolej funkci randomnames.
168 {                                           % take_random_walk()
169     local~randomnames =                     % v jazyce Lua.
170     require("randomnames")
171     local~result = randomnames.take_random_walk(
172     " \lua_escape:e { #2 } ",
173     \l_randomnames_graph_type_tl,
174     \l_randomnames_affix_tl,
175     \l_randomnames_seed_tl,
176     \l_randomnames_min_context_size_tl,
177     \l_randomnames_max_context_size_tl
178     )
179     token.set_macro(                         % Ulož vygenerované jméno
180     "l_tmpa_tl", result)                    % do proměnné \l_tmpa_tl
181 }                                           % jazyka expl3.
182 \bool_if:NT                                % Pokud uživatel uvedl
183 \l_randomnames_save_bool                   % parametr save, ulož
184 {                                           % vygenerované jméno
185     \cs_gset:cpx                             % do zadaného TeXového
186     \l_randomnames_save_tl                 % příkazu.
187     { \l_tmpa_tl }
188 }
189 \tl_use:N \l_tmpa_tl                       % Vlož jméno na výstup.
190 \group_end:                                % Resetuj hodnoty lokálních
191 }                                           % proměnných zpět na nil.
192
193 \ExplSyntaxOff

```

Pomocí našeho přemostění vytvoříme Katzův couvající jazykový model a provedeme v něm náhodnou procházku v plain LuaTeXu takto:

```

\input randomnames                           % Načti naše
\ExplSyntaxOn                                % přemostění.
\randomnames_new_model:nn                    % Vytvoř model.
{
    min_context_size = 1,
    max_context_size = 3,

```

```

    seed = 22,
  }
  { pes }
\randomnames_add_name:nn { pes } { pes~filipes } % Zanes do modelu
\randomnames_add_name:nn { pes } { pejsek }      % jména
\randomnames_add_name:nn { pes } { maxipes~fik } % pohádkových psů.
\randomnames_take_random_walk:nn { } { pes }     % Vygeneruj jméno.
\ExplSyntaxOff \bye
maxipes filipes fik

```

### 3.3. Uživatelské rozhraní pomocí L<sup>A</sup>T<sub>E</sub>Xového balíčku xparse

Programovací jazyk expl3 má pro uživatele T<sub>E</sub>Xu poměrně neobvyklou syntax. Pomocí L<sup>A</sup>T<sub>E</sub>Xového balíčku xparse proto vytvoříme uživatelské rozhraní, které bude pro uživatele přirozenější.

```

1 \input randomnames\relax           % Načti naše přemostění.
2
3 \ProvidesExplPackage
4   {randomnames}%
5   {2023-06-25}%
6   {1.0.2}%
7   {Character name generators for creative writing in LuaLaTeX}
8
9 \RequirePackage{xparse}           % Načti balíček xparse.
10
11 \NewDocumentCommand              % Vytvoř LaTeXový příkaz
12   { \newmodel }                 % \newmodel s jedním
13   { 0{ } m }                   % nepovinným a jedním
14   {                             % povinným argumentem,
15     \randomnames_new_model:nn   % který zavolá funkci
16     { #1 }                      % randomnames.new_model()
17     { #2 }                      % v jazyce Lua.
18   }
19
20 \NewDocumentCommand              % Vytvoř LaTeXový příkaz
21   { \addname }                 % \addname se dvěma
22   { m m }                      % povinnými argumenty,
23   {                             % který zavolá funkci
24     \randomnames_add_name:nn   % randomnames.add_name()
25     { #1 }                    % v jazyce Lua.
26     { #2 }
27   }

```



```

28 \NewDocumentCommand                               % Vytvoř LaTeXový příkaz
29   { \inputnames }                                % \inputnames se dvěma
30   { m m }                                         % povinnými argumenty,
31   {                                               % který zavolá funkci
32     \randomnames_input_names:nn                  % randomnames.input_names()
33     { #1 }                                        % v jazyce Lua.
34     { #2 }
35   }
36
37 \NewDocumentCommand                               % Vytvoř LaTeXový příkaz
38   { \randomname }                                 % \randomname s jedním
39   { 0{ } m }                                     % nepovinným a jedním
40   {                                               % povinným argumentem,
41     \randomnames_take_random_walk:nn           % který zavolá funkci
42     { #1 }                                        % randomnames.
43     { #2 }                                        % take_random_walk()
44   }                                               % v jazyce Lua.

```

Pomocí našeho uživatelského rozhraní příkazů vytvoříme Katzův couvající jazykový model a provedeme v něm náhodnou procházku v Lua<sub>L</sub>AT<sub>E</sub>Xu následně:

```

\documentclass{article}
\usepackage[czech]{babel}
\usepackage{randomnames}
\newmodel[min_context_size=1,                       % Vytvoř model.
          max_context_size=3,
          seed=22]{pes}
\addname{pes}{pes filipes}                         % Zanes do modelu
\addname{pes}{pejsek}                              % jména
\addname{pes}{maxipes fík}                         % pohádkových psů.
\begin{document}
\randomname{pes}                                    % Vygeneruj jméno.
\end{document}
maxipes filipes fík

```

## 4. Příklady užití

V této sekci stáhneme několik příkladových databází jmen, natrénujeme na nich náš model a ukážeme si, jak bychom vygenerovali jména postav a zahrnuli je do textu povídky.

Při stahování příkladových databází používáme terminál UNIXového systému s příkazovým procesorem `/bin/bash` a s nainstalovaným programem `xmllint` z knihovny `libxml2`. Příkladové dokumenty si můžeme uložit do textového souboru

příklad.tex a přeložit příkazem `lualatex příklad.tex`. Pro úspěšný překlad stáhneme do pracovního adresáře soubory `randomnames.lua`, `randomnames.tex` a `randomnames.sty` [6], které jsme vytvořili v sekcích 2 a 3.

#### 4.1. Soudobá křestní jména

V soudobé fikci se nám může hodit generovat křestní jména postav. Stáhneme si proto z webu `rodina.cz` seznam 629 mužských a 605 ženských křestních jmen:

```
$ wget -O- https://www.rodina.cz/scripts/jmena/default.asp?muz=0 |
  xmllint -html -xpath '//a[contains(@href,"jmeno.asp")]/text()' \
  - > krestni-jmena-zeny.txt
$ wget -O- https://www.rodina.cz/scripts/jmena/default.asp?muz=1 |
  xmllint -html -xpath '//a[contains(@href,"jmeno.asp")]/text()' \
  - > krestni-jmena-muzi.txt
```

Následně sestavíme dva modely, jeden pro mužská křestní jména a druhý pro ženská křestní jména, a použijeme je v textu povídky:

```
\documentclass{article}
\usepackage[czech]{babel}
\usepackage{randomnames}
\newmodel{zena}
\newmodel{muz}
\inputnames{zena}{krestni-jmena-zeny.txt}
\inputnames{muz}{krestni-jmena-muzi.txt}
\begin{document}
Když se setkali na výstavě psů, \randomname[prefix=Ro]{muz}
a \randomname[suffix=lie]{zena} si okamžitě uvědomili, že se už
nikdy nebudou chtít rozloučit.
\end{document}
```

*Když se setkali na výstavě psů, Rodan a Aurélie si okamžitě uvědomili, že se už nikdy nebudou chtít rozloučit.*

V našem příkladu jsme použili výchozí maximální délku kontextu 3, abychom modely pobídli k vyšší kreativitě. Už s délkou kontextu 6 bychom ale dokázali vygenerovat 23 nových přesvědčivých jmen, která nebyla součástí množiny trénovacích jmen z webu `rodina.cz`: Adalbertina, Adalbertýna, Budimíra, Budislava, Erharda, Fridolína, Gvendolín, Jarolína, Jonathanael, Karolím, Karolíma, Klementýn, Mojmíra, Mstislava, Něhoslava, Stanimíra, Stojmíra, Svatomíra, Šebastian, Tichomíra, Velimíra, Vítoslava a Věslav.

#### 4.2. Postavy světa J. R. R. Tolkiena

Ve fanfikci z fantasy světa anglického spisovatele J. R. R. Tolkiena se nám může hodit generovat jména postavám různých ras. Stáhneme proto z webu

behindthename.com seznam 51 jmen trpaslíků, 19 jmen elfek a 316 jmen lidských mužů.

```
$ download_tolkien_names() {  
  wget -O- https://behindthename.com/namesakes/list/tolkien/alpha|  
  xmllint -html -xpath '//div[@id="div_refinepage"]/table/  
  tr[td[2]/text() = "$2" and td[3]/text() = "$1"]/td[1]' - |  
  sed 's/<[\>]*>//g'; }  
$ download_tolkien_names Dwarf m > tolkien-jmena-trpasliku.txt  
$ download_tolkien_names Elf f > tolkien-jmena-elfek.txt  
$ download_tolkien_names Man m > tolkien-jmena-lidskych-muzu.txt
```

Následně sestavíme tři modely, jeden pro jména trpaslíků, druhý pro jména elfek a třetí pro jména lidských mužů, a použijeme je v textu fanfikce:

```
\documentclass{article}  
\usepackage[czech]{babel}  
\usepackage{randomnames}  
\newmodel[context_size=2,seed=35]{trpaslik}  
\newmodel[max_context_size=2,seed=6]{elfka}  
\newmodel[seed=2419]{clovek}  
\inputnames{trpaslik}{tolkien-jmena-trpasliku.txt}  
\inputnames{elfka}{tolkien-jmena-elfek.txt}  
\inputnames{clovek}{tolkien-jmena-lidskych-muzu.txt}  
\begin{document}
```

Tři hrdinové - trpaslík, elf a člověk - se vydali do temného lesa, aby zastavili invazi nepřátelských vrrků. Trpaslík

```
\randomname{trpaslik} byl nastražen jako návnada na zuby těchto  
dravých zvířat, elfka \randomname[suffix=iel]{elfka} připravovala  
luk, zatímco člověk \randomname{clovek} divoce máchal mečem.  
\end{document}
```

*Tři hrdinové - trpaslík, elf a člověk - se vydali do temného lesa, aby zastavili invazi nepřátelských vrrků. Trpaslík Duri byl nastražen jako návnada na zuby těchto dravých zvířat, elfka Amariel připravovala luk, zatímco člověk Túrindor divoce máchal mečem.*

Pokud neuvedeme náhodné semínko, každé použití příkazu `\randomname` vygeneruje nové jméno. Pro zachování vygenerovaných jmen můžeme použít parametr `save`, kterým jméno postavy trvale uložíme:

```
\randomname[save=Žoldněř]{clovek} vyrostl v bohaté obchodnické  
rodině. Když však jeho otec zesnul a rodinné jmění zpustlo,  
\Žoldněř se stal dobrodruhem a nechal se najímat jako žoldněř.  
Během svých dobrodružství si \Žoldněř vydobyl pověst neohroženého  
válečníka a jeho meč se stal jeho nejcennějším jméním.
```

---

Vstupní procesor LuaTeXu zná Unicode a můžeme tedy psát příkazy s diakritikou jako `\Žoldněř`.

*Túrindor vyrostl v bohaté obchodnické rodině. Když však jeho otec zesnul a rodinné jmění zpustlo, Túrindor se stal dobrodruhem a nechal se najímat jako žoldněř. Během svých dobrodružství si Túrindor vydobyl pověst neohroženého válečníka a jeho meč se stal jeho nejcenějším jměním.*

### 4.3. Bytosti v Cthulhu mýtu H. P. Lovecrafta

Ve fanfikci z hororového světa amerického spisovatele H. P. Lovecrafta se nám může hodit generovat jména kosmických božstev, tzv. *prastarých*. Stáhneme proto z anglické wikipedie seznam 157 jmen prastarých:

```
$ wget -O- https://en.wikipedia.org/wiki/List_of_Great_Old_Ones |
  xmllint -html -xpath '//table[2]//td[1]/text()' - | grep . \
  > lovecraft-jmena-prastarych.txt
```

Následně sestavíme model jmen prastarých a použijeme ho v textu fanfikce:

```
\documentclass{article}
\usepackage[czech]{babel}
\usepackage{randomnames}
\newmodel[seed=140]{prastary}
\inputnames{prastary}{lovecraft-jmena-prastarych.txt}
\begin{document}
```

Po západu slunce se skupina šílených kultistů shromáždila kolem oltáře, aby přivolali svého pána. Prastarý

```
\randomname[prefix=C']{prastary}, jehož jméno nebylo možné
vyslovit lidskými ústy, se brzy objevil a jeho oči temně zářily.
\end{document}
```

*Po západu slunce se skupina šílených kultistů shromáždila kolem oltáře, aby přivolali svého pána. Prastarý C'tya-Yg'Nalla, jehož jméno nebylo možné vyslovit lidskými ústy, se brzy objevil a jeho oči temně zářily.*

### 4.4. Kouzla ve světě Harryho Pottera

Ve fanfikci ze světa Harryho Pottera se nám může hodit generovat nová jména kouzel. Stáhneme proto z webu harrypotter.fandom.com seznam 192 kouzel:

```
$ wget -O- https://harrypotter.fandom.com/wiki/List_of_spells |
  xmllint -html -xpath '//h3//i/a/text()' - \
  > harry-potter-kouzla.txt
```

Následně sestavíme model jmen kouzel a použijeme ho v textu fanfikce:

```
\documentclass{article}
\usepackage[czech]{babel}
\usepackage{randomnames}
```

```

\newmodel{kouzlo}
\inputnames{kouzlo}{harry-potter-kouzla.txt}
\begin{document}
Tesák, Hagridův pes, ležel před chatrčí a pozoroval okolní les.
Najednou uslyšel hlasitý křik z nedaleké cesty. Překvapeně se
postavil na nohy a natahoval uši. Když zpozoroval kouzelníka, jak
na něj mává hůlkou a křičí „\randomname[seed=2136]{kouzlo}!“,
okamžitě sebou trhl a rozběhl se k němu s otevřenou tlamou.
\end{document}

```

*Tesák, Hagridův pes, ležel před chatrčí a pozoroval okolní les. Najednou uslyšel hlasitý křik z nedaleké cesty. Překvapeně se postavil na nohy a natahoval uši. Když zpozoroval kouzelníka, jak na něj mává hůlkou a křičí „Furnunculus!“, okamžitě sebou trhl a rozběhl se k němu s otevřenou tlamou.*

#### 4.5. Kapesní příšerky Pokémon

Ve fanfikci ze světa kapesních příšerek Pokémon se nám může hodit generovat jména pokémonů. Stáhneme proto z webu [bulbapedia.bulbagarden.net](http://bulbapedia.bulbagarden.net) seznam 1015 jmen pokémonů:

```

$ wget -O- 'https://bulbapedia.bulbagarden.net/wiki/List_of_' \
  'Pokémon_by_name' | xmllint -html -xpath '//table//img/@alt' - |
  sed -r 's/ alt="(.*")/1/' > jmena-pokemonu.txt

```

Následně sestavíme model jmen pokémonů a použijeme ho v textu fanfikce:

```

\documentclass{article}
\usepackage[czech]{babel}
\usepackage{randomnames}
\newmodel[max_context_size=4,seed=4]{pokemon}
\inputnames{pokemon}{jmena-pokemonu.txt}
\begin{document}
Cestování po světě pokémonů může být plné nečekaných překvapení.
Onoho dne se trenér z Rumělkového města procházel po parku, když
tu náhle na něj z vysokého podrostu vyskočil psí pokémon
\randomname[prefix=Haf]{pokemon} a hlasitě zaštěkal. Od té doby
trenér věděl, že musí být vždy připraven na nečekané situace.
\end{document}

```

*Cestování po světě pokémonů může být plné nečekaných překvapení. Onoho dne se trenér z Rumělkového města procházel po parku, když tu náhle na něj z vysokého podrostu vyskočil psí pokémon Hafeon a hlasitě zaštěkal. Od té doby trenér věděl, že musí být vždy připraven na nečekané situace.*

## 5. Budoucí práce

V této sekci uvedeme možná vylepšení vyvinutého jazykového modelu a další aplikace jazykových modelů v  $\text{T}_{\text{E}}\text{X}$ u, na které by se měla zaměřit budoucí práce.

### 5.1. Automatické skloňování českých a slovenských jmen

Pokud nemáme trénovací data v různých gramatických pádech a číslech, bude náš jazykový model generovat jména pouze v prvním pádě jednotného čísla. V analytickém jazyce, jako je angličtina, to není problém. U flektivních jazyků, jako jsou čeština a slovenština, bychom ale chtěli, aby náš model uměl jména nejen generovat, ale také skloňovat.

Tereza Vrabcová [12] navrhla následující algoritmus skloňování českých jmen:

1. Pro každé jméno  $J$  v trénovací množině:
  - (a) Pokud jméno existuje ve wikislovníku, stáhni jeho skloňování odtud.
  - (b) Jinak stáhni skloňování jména z webu sklonuj.cz.
2. Pro každé automaticky vygenerované jméno  $J'$ :
  - (a) Najdi v trénovací množině jméno  $J$  s nejpodobnější koncovkou:

$$J = \arg \max_J \sum_{i=1, \dots, \min(|J|, |J'|)} \begin{cases} \frac{1}{2^i} & J[|J| - i + 1] = J'[|J'| - i + 1], \\ 0 & \text{jinak.} \end{cases}$$

- (b) Skloňuj jméno  $J'$  stejně jako jméno  $J$ .

Pro vyzkoušení algoritmu stáhneme do pracovního adresáře soubory `get_data.py`, `requirements.txt`, `declension.lua`, `declension.tex` a `declension.sty` [12]. Následně si stáhneme z webu `rodina.cz` česká křestní jména podle pokynů v sekci 4.1 na straně 27. Dále si pomocí skriptu `get_data.py` stáhneme skloňování jmen:

```
$ python3 -m pip install -r requirements.txt
$ python3 ./get_data.py krestni-jmena-zeny.txt krestni-jmena-zeny
$ python3 ./get_data.py krestni-jmena-muzi.txt krestni-jmena-muzi
```

Nakonec sestavíme dva modely, jeden pro mužská křestní jména a druhý pro ženská křestní jména, a použijeme je v textu povídky:

```
\documentclass{article}
\usepackage[czech]{babel}
\usepackage{randomnames}
\newmodel{zena}
\newmodel{muz}
\inputnames{zena}{krestni-jmena-zeny.txt}
\inputnames{muz}{krestni-jmena-muzi.txt}
\usepackage{declension}
\newdeclmodel{zena}
```

```

\newdeclmodel{muz}
\loaddata{zena}{krestni-jmena-zeny.decl}{krestni-jmena-zeny.names}
\loaddata{muz}{krestni-jmena-muzi.decl}{krestni-jmena-muzi.names}
\begin{document}
Když se setkali na výstavě psů, \randomname[prefix=Ro,save=0n]
{muz} a \randomname[suffix=lie,save=0na]{zena} si okamžitě
uvědomili, že se už nikdy nebudou chtít rozloučit. \decline{zena}
{\0na}{1}{s} si uvědomila, že bez \decline{muz}{\0n}{2}{s} nechce
žít. \decline{muz}{\0n}{1}{s} si ve stejný okamžik uvědomil, že
s \decline{zena}{\0na}{7}{s} chce žít navždy.
\end{document}

```

*Když se setkali na výstavě psů, Rodan a Aurélie si okamžitě uvědomili, že se už nikdy nebudou chtít rozloučit. Aurélie si uvědomila, že bez Rodana nechce žít. Rodan si ve stejný okamžik uvědomil, že s Aurélií chce žít navždy.*

Popsaný algoritmus funguje pouze pro česká jména a hledání jména s nejpodobnější koncovkou má lineární časovou složitost ve velikosti trénovací množiny. Budoucí práce by se měla zaměřit na rozšíření algoritmu o další flektivní jazyky jako slovenština a použít datovou strukturu jako trie pro efektivní hledání trénovacích jmen s nejdelsí společnou koncovkou.

## 5.2. Omezující podmínky pro generovaná jména

Náš jazykový model umožňuje uživateli zadat buď předponu, nebo příponu vygenerovaného jména, ale ne obojí zároveň. Dále uživatel nemůže ovlivnit střed vygenerovaného jména ani omezit jeho délku. Ve chvíli, kdy uživatel není s vygenerovaným jménem spokojený, musí jméno generovat opakovaně s různými náhodnými semínky. Tento zdlouhavý proces bychom rádi urychlili tak, že uživateli umožníme zadat dodatečné omezující podmínky.

### 5.2.1. Permutermíny

Pro rozšířené vyhledávání v triích se využívají tzv. *permutermíny* [13, sekce 3.2.1]. Při sestavování trie použijeme místo slova „pejsek“ veškeré jeho rotace (permutermíny) doplněné o metasymbol \$, který udává konec slova: „pejsek\$“, „ejsek\$p“, „jsek\$pe“, „sek\$pej“, „ek\$pejs“, „k\$pejse“ a „\$pejsek“. Následně můžeme hledat slova a slovní spojení s předponou pe- a příponou -ek hledáním permutermínů s předponou ek\$pe-. Dále můžeme hledat slova a slovní spojení obsahující text „ejs“ hledáním permutermínů s předponou ejs-.

Vzhledem k tomu, že jsme náš zapomnětlivý jazykový model odvodili od trií, mohlo by se zdát, že permutermíny budeme moci využít i v něm. Pokud ale do modelu vložíme místo jmen permutermíny, všimneme si několika nežádoucích jevů: ačkoliv náš jazykový model začne generovat permutermín s předponou ek\$pe-,

nic mu nebrání v tom, aby v permutermínu vygeneroval několik metasymbolů \$ a přestal generovat v bodě, který není pro požadovanou koncovku -ek vhodný:

```

local function add_permuterminals(model_name, name)
  randomnames.add_name(model_name, name .. "$")
  for position, code in utf8.codes(name) do
    local character = utf8.char(code)
    local prefix = name:sub(1, position + #character - 1)
    local suffix = name:sub(position + #character, -1)
    randomnames.add_name(model_name, suffix .. "$" .. prefix)
  end
end

randomnames.new_model("perm", 1, 3, 17) -- Vytvoř model.
add_permuterminals("perm", "pes filipes") -- Zanes do modelu
add_permuterminals("perm", "pejsek") -- permutermíny
add_permuterminals("perm", "maxipes fík") -- pohádkových psů.
local name = randomnames. -- Vygeneruj permutermín.
  take_random_walk("perm", "forward", "ek$pe")
tex.print(name)

ek$pes fík$maxi

```

Popsané problémy souvisí se zapomnětlivostí našeho jazykového modelu a nemají přímočaré řešení, tj. bez úprav našeho modelu nelze permutermíny použít. Budoucí práce by se měla zaměřit na popis a implementaci potřebných úprav.

### 5.2.2. Regulární výrazy

Dalším způsobem, jak omezit vygenerované jméno, jsou regulární výrazy:

```

local name = nil
while name == nil or not name:match( -- Vygeneruj jméno s danou
  "maxipes " .. (""):rep(60) .. -- předponou, příponou a
  ("?.?"):rep(10) .. " fík") do -- 60 až 70 znaky uprostřed.
  name = randomnames.take_random_walk("pes")
end
tex.print(name)

maxipes filipes filipes filipes filipes filipes filipes filipes filipes fík

```

Ve výše uvedeném příkladu jsme opakovaně generovali jména, dokud jsme nenašli takové, které by odpovídalo zadanému regulárnímu výrazu. Takovéto hledání je implementačně přímočaré, ale může být výpočetně náročné, pokud hledáme nepravděpodobné jméno.

Pokud sestavíme pro zadaný regulární výraz nedeterministický konečný automat, můžeme již během náhodné procházky navštívit pouze hrany, které automat



nedostanou do neakceptujícího stavu, a přestat generovat, pouze když je automat v akceptujícím stavu. Tento přístup je efektivnější, ale ani zde se nevyhne opakovanému generování jmen: snadno se totiž dostaneme do slepé uličky, kdy náš model není schopný vygenerovat odpovídající jméno. Pokud bychom například, podobně jako ve výše uvedeném příkladu, chtěli vygenerovat jméno pohádkového psa, ale omezili bychom ho regulárním výrazem  $maxipes_{\perp}.*_{\perp}pejsek$ , náš model bude generovat nekonečný textový řetězec „maxipes filipes filipes filipes ...“. Automat modelu nedovolí přestat generovat, dokud nezakončí jméno slovem „pejsek“, což ale nemůže nastat, vizte Obrázek 3b na straně 7. Jelikož je náš jazykový model také (pravděpodobnostní) konečný automat, který generuje regulární jazyky, mohli bychom podobné slepé uličky detekovat algoritmem, který zkontroluje, že regulární jazyk generovaný modelem a regulární jazyk akceptovaný automatem s aktuálním stavem nastaveným jako počátečním mají neprázdný průnik. Nejefektivnější známý algoritmus má ale exponenciální časovou složitost [14], což ho pro naše potřeby činí nepraktickým. Jako praktické řešení se jeví shora omezit počet kroků modelu během náhodné procházky a počet opakovaných pokusů o náhodnou procházku.

### 5.3. Vyhlažování pravděpodobnosti

V sekci 4 jsme pracovali s trénovacími množinami, které často obsahovaly jen malé desítky jmen. Jazykové modely trénované na těchto množinách jmen mají omezené schopnosti generalizace a nedokážou např. vygenerovat jména se znaky, které se neobjevily v trénovací množině. U mnoha jmen však známe dodatečné informace, kterými bychom mohli jazykové modely obohatit: jména prastarých v Cthulhu mýtu H. P. Lovecrafta (vizte sekci 4.3) jsou psána v latinkové transkripci r'lyehštiny, která byla pravděpodobně inspirována psanou velštinou [15] podobně jako jména v elfské sindarinštině [16] (vizte sekci 4.2). Anglická jména kapesních přišerek Pokémon (vizte sekci 4.5) jsou často anglické složeniny, např. bulbasaur = bulb (žárovka) + dinosaur a beedrill = bee (včela) + drill (vrták) [17].

Jedním způsobem, jak zlepšit generalizační schopnosti našich modelů, je pomocí tzv. *vyhlažování pravděpodobnosti* smíšením několika modelů. Pokud máme např. model sindarinštiny  $M_s$  a model velštiny  $M_v$ , můžeme pro výpočet pravděpodobnosti dalšího znaku  $S[i]$  použít lineární vyhlažování:

$$P(S[i] | k) = \lambda \cdot P(S[i] | k; M_s) + (1 - \lambda) \cdot P(S[i] | k; M_v),$$

kde  $k = S[i-n, \dots, i-1]$  je kontext posledních  $n$  znaků a  $\lambda$  je relativní váha sindarinštiny oproti velštině. Díky velké trénovací množině pro velštinu bude výsledný model schopný lépe generalizovat mimo trénovací množinu pro sindarinštinu a generovat kreativní nová jména.

## 5.4. Možné aplikace jazykových modelů v $\text{T}_{\text{E}}\text{X}$ u

V tomto článku jsme se zaměřili na využití jazykových modelů pro generování jmen postav, což je zajímavá, ale okrajová aplikace jazykových modelů. Vzhledem k tematickému zaměření *Zpravodaje  $\mathcal{C}\mathcal{S}\mathcal{T}\mathcal{U}\mathcal{G}\mathcal{U}$*  by mohlo čtenáře zajímat, jakým způsobem mohou jazykové modely pomoci nejen autorovi textu, ale také sazeči.

### 5.4.1. Detekce jazyka pro účely dělení slov

V naší implementaci jsme se zaměřili na generativní využití jazykových modelů, kdy pomocí náhodných procházek vytváříme nová jména postav. Pokud bychom ale naše modely místo jmen trénovali na větách a pokud bychom rozšířili naši implementaci o výpočet pravděpodobnosti  $P(V; M)$  vět  $V$  v daném jazykovém modelu  $M$ , pak můžeme jazykové modely využít pro detekci jazyka.

Detekce jazyka má přímočaré využití v  $\text{T}_{\text{E}}\text{X}$ ovém algoritmu řádkového zlomu. Můžeme např. připravit  $\text{T}_{\text{E}}\text{X}$ ové makro `\foreignlanguage{věta}`, které rozpozná jazyk věty a vysází ji se správnými vzory dělení slov. Podobný příkaz má internacionalizační  $\text{L}_{\text{A}}\text{T}_{\text{E}}\text{X}$ ový balíček `babel`, u kterého ale musíme kromě věty vždy ručně uvést i jazyk. Ještě vyšší úroveň uživatelského komfortu můžeme dosáhnout, když budeme pomocí callbacku `pre_linebreak_filter` `LuaTEXu` [18, sekce 9.5.4] automaticky přepínat vzory dělení slov podle jazyka aktuálního odstavce.

Pokud máme např. model češtiny  $M_{\text{CZ}}$  a model slovenštiny  $M_{\text{SK}}$ , pak můžeme rozhodnout, je-li věta  $V$  v češtině, nebo ve slovenštině porovnáním pravděpodobností  $P(\text{CZ} | V)$  a  $P(\text{SK} | V)$ . Aplikací Bayesovy věty dostáváme:

$$P(\text{CZ} | V) = \frac{P(V; M_{\text{CZ}}) \cdot P(\text{CZ})}{P(V)}.$$

Pokud předpokládáme, že všechny jazyky jsou apriori stejně pravděpodobné, pak  $P(\text{CZ})/P(V)$  je jen multiplikační konstanta a dostáváme, že věta  $V$  je v češtině právě tehdy, když  $P(V; M_{\text{CZ}}) > P(V; M_{\text{SK}})$ , což dokážeme snadno spočítat pouze s využitím našich dvou modelů.

V praxi rozlišovat češtinu od slovenštiny nepotřebujeme, protože oba jazyky mají v  $\text{T}_{\text{E}}\text{X}$ u společné vzory pro dělení slov [19, 20]. Máme ale mnoho jiných jazyků  $J$ , pro které společné vzory pro dělení slov (zatím) neexistují. Při detekci jazyka bychom vždy vybrali nejpravděpodobnější jazyk  $\text{CZ} = \arg \max_{\text{CZ} \in J} P(V; M_{\text{CZ}})$ .

### 5.4.2. Výplňový text

Při přípravě dokumentových šablon se nám může hodit generovat výplňový text, který do šablon umístíme při přípravě ukázkových maket. Pokud budeme, stejně jako v předchozí sekci, trénovat naše modely místo jmen na větách, můžeme snadno generovat věrohodný výplňový text. Stáhněme si například 114 vět z knihy *Dášeňka čili život štěněte* spisovatele Karla Čapka [21]:

```
$ for NAME in dasenka-cili-zivot-stenete{,-{2..5}}.html; do
  wget -O- https://www.cesky-jazyk.cz/citanka/karel-capek/$NAME |
  xmllint -html -xpath '//div[@id="dbtext"]//p[2]/text()' -
  done | grep '\w' | sed -e 's/€#13;//g' -e 's/\./.\n/g' \
  > karel-capek-dasenka.txt
```

Následně sestavíme model vět o psech a použijeme ho k vygenerování odstavce výplňového textu:

```
\documentclass{article}
\usepackage[czech]{babel}
\usepackage{randomnames}
\newmodel[context_size=9]{veta-o-psu}
\inputnames{veta-o-psu}{karel-capek-dasenka.txt}
\begin{document}
\randomname[seed=417]{veta-o-psu}
\randomname[seed=10204]{veta-o-psu}
\randomname[seed=9335]{veta-o-psu}
\end{document}
```

*Ba, lidi, takové bílé nic, do hrsti se to vešlo; ale anžto to mělo pár černých ušísek a vzadu ocásek, ocáskem se nedá chodit. Nevyhne se žádné lidské náruči; pak dostane pro útěchu špičku člověčího nosu, aby se naučila sedět a chodit a ledacos jiného důležitého pro život. A musí se olíznout jazejkem umývala, česala a hladila, pěstovala ji, líbala a lízala, čistila a jazejkem, a namoutě kutě, ono je to dobré.*

Vygenerovaný text splňuje požadavky kladené na výplňový text: je nonsensový, ale zachovává charakter trénovacích vět do té míry, že je dobře zalomitelný algoritmem řádkového zlomu. Vzhledem k tomu, že u výplňového textu nečekáme přítomnost nových slov, může být vhodnější (a efektivnější) použít jazykový model, který negeneruje věty po znacích, ale po celých slovech.

## 6. Závěr

V tomto článku jsme ukázali, jak můžeme pomocí jazykových modelů automaticky generovat jména postav při tvůrčím psaní v Lua $\TeX$ u. Dále jsme vyvinuli balíček, který je přístupný laickým uživatelům  $\LaTeX$ u, a zveřejnili jsme ho online [6]. Nakonec jsme analyzovali možná vylepšení našeho modelu a další možné aplikace jazykových modelů v Lua $\TeX$ u. Článek názorně ukazuje vývoj  $\LaTeX$ ového balíčku, který kombinuje kód v  $\TeX$ u s kódem v jazyce Lua, a demonstruje možnosti využití jazykových modelů v  $\TeX$ u.

## Odkazy

1. KNUTH, Donald. 6.3: Digital Searching. In: *The Art of Computer Programming Volume 3: Sorting and Searching*. 2. vyd. Addison-Wesley, 1997, s. 492. ISBN 0-201-89685-0.
2. *Check out this transparent Huckleberry Hound hello PNG image* [online]. [cit. 2023-02-22]. Dostupné z: [https://cartoongoodies.com/png\\_images/huckleberry-hound-hello/](https://cartoongoodies.com/png_images/huckleberry-hound-hello/).
3. ČAPEK, Josef. *Povídání o pejskovi a kočičce: jak spolu hospodařili a ještě o všelijakých jiných věcech*. 13. vyd. Albatros, 1972.
4. *GoGEN Maxipes Fík: Česká značka elektroniky pro děti* [online]. [cit. 2023-02-22]. Dostupné z: <https://www.gogen.cz/gogen-maxipes-fik/>.
5. KATZ, S. M. Estimation of probabilities from sparse data for the language model component of a speech recognizer. *IEEE Transactions on Acoustics, Speech, and Signal Processing*. 1987, roč. 35, č. 3, s. 400–401.
6. NOVOTNÝ, Vít. *Nápadovnick jmen postav pro tvůrčí psaní v LuaTeXu: Release The latest version* [online]. GitHub, 2023-05-03 [cit. 2023-05-03]. Dostupné z: <https://github.com/witiko/character-name-generator-for-creative-writing-in-luatex/releases/tag/latest>.
7. IERUSALIMSKY, Roberto. *Programming in Lua*. 4. vyd. Lua.org, 2016. ISBN 978-8590379867.
8. IERUSALIMSKY, Roberto; FIGUEIREDO, Luiz Henrique de; CELES, Waldemar. *Lua 5.3 Reference Manual* [online]. 2020-07-14. [cit. 2023-03-10]. Dostupné z: <https://www.lua.org/manual/5.3/manual.html>.
9. PRESS, W.H.; TEUKOLSKY, S.A.; VETTERLING, W.T.; FLANNERY, B.P. *Numerical Recipes 3rd Edition: The Art of Scientific Computing*. 3. vyd. Cambridge University Press, 2007. ISBN 978-0-521-88068-8.
10. L'ECUYER, Pierre. Tables of linear congruential generators of different sizes and good lattice structure. *Mathematics of Computation*. 1999, roč. 68, č. 225, s. 249–260.
11. THE L<sup>A</sup>T<sub>E</sub>X PROJECT. *The L<sup>A</sup>T<sub>E</sub>X3 Interfaces* [online]. CTAN, 2023-03-30 [cit. 2023-04-09]. Dostupné z: <http://mirrors.ctan.org/macros/latex/contrib/l3kernel/interface3.pdf>.
12. VRABCOVÁ, Tereza. *Model pro automatické skloňování českých jmen v LuaTeXu: Release The latest version* [online]. GitHub, 2023-05-03 [cit. 2023-05-03]. Dostupné z: [https://github.com/xvrabcov/declension\\_names/releases/tag/latest](https://github.com/xvrabcov/declension_names/releases/tag/latest).
13. MANNING, Christopher D.; RAGHAVAN, Prabhakar; SCHÜTZE, Hinrich. *Introduction to Information Retrieval*. Cambridge University Press, 2008. ISBN 0521865719. Dostupné také z: <https://nlp.stanford.edu/IR-book/>.

14. RABIN, Michael O; SCOTT, Dana. Finite automata and their decision problems. *IBM J. Res. Dev.* 1959, roč. 3, č. 2, s. 114–125.
15. NOVOTNÝ, Vít; STARÁ, Marie. Cthulhu Hails from Wales:  $N$ -gram Frequency Analysis of R'lyehian. In: HORÁK, Aleš; RYCHLÝ, Pavel; RAMBOUSEK, Adam (ed.). *Recent Advances in Slavonic Natural Language Processing*. Tribun EU, 2020. ISBN 978-80-263-1600-8. ISSN 2336-4289. Dostupné také z: <https://nlp.fi.muni.cz/raslan/2020/paper12.pdf>.
16. HOOKER, Mark T. *Tolkien and Welsh (Tolkien a Chymraeg): Essays on J. R. R. Tolkien's Use of Welsh in his Legendarium*. 1. vyd. Llyfrau, 2012. ISBN 978-1477667736.
17. CREATIVEBLOGGER. *Pokémon in Translation: Where Do Their Names Come From?* [online]. 2017-08-11. [cit. 2023-04-17]. Dostupné z: <https://creativetranslation.com/blog-pokemon-names-in-translation/>.
18. L<sup>A</sup>T<sub>E</sub>X DEVELOPMENT TEAM. *LuaT<sub>E</sub>X Reference Manual* [online]. CTAN, 2023-04-06 [cit. 2023-04-17]. Dostupné z: <https://ctan.org/pkg/luatex>. Verze 1.16.
19. SOJKA, Petr; SOJKA, Ondřej. The Unreasonable Effectiveness of Pattern Generation. *Zpravodaj ČSTUGu*. 2019, roč. 29, č. 1–4, s. 73–86. ISSN 1211-6661. Dostupné z DOI: 10.5300/2019-1-4/73.
20. SOJKA, Petr; SOJKA, Ondřej. Towards New Czechoslovak Hyphenation Patterns. *Zpravodaj ČSTUGu*. 2020, roč. 30, č. 3–4, s. 118–126. ISSN 1211-6661. Dostupné z DOI: 10.5300/2020-3-4/118.
21. ČAPEK, Karel. *Dášeňka čili život štěněte*. 1. vyd. Fr. Borový, 1933.

## Summary: Character Name Generator for Creative Writing in LuaT<sub>E</sub>X

A famous dictum of the computer scientist Phil Karlton says that there are only two difficult things in computer science: cache invalidation and naming things. This is also true in creative writing, where authors have to come up not just with a story and a setting but also the names of all their fictional characters. In this article, we develop a language model in LuaT<sub>E</sub>X, which allows authors to automatically generate names for their characters. Besides creative writing, we also discuss other uses of language models in LuaT<sub>E</sub>X, namely the automatic switching of hyphenation patterns based on the current language and blind text generation. For the T<sub>E</sub>Xnically-minded users, the article acts as an introduction to the programming languages of Lua and expl3, and also the xparse L<sup>A</sup>T<sub>E</sub>X package for defining document commands in L<sup>A</sup>T<sub>E</sub>X.

**Keywords:** creative writing, trie, language models, LuaT<sub>E</sub>X, Lua, expl3, xparse

*Vít Starý Novotný, witiko@mail.muni.cz*