

# 20SK – Signály a kódy

---

## Přednáška 11 – Slovníkové metody komprese dat (26.11.2018)

Probíraná témata:

- RLE (Run-Length Encoding)
- Slovníkové metody komprese – princip, teoretické vlastnosti
- Algoritmus LZ77
- Algoritmus LZW

Relevantní literatura je [1,2] pro LZ77 a [3] pro LZW. Širší úvod do metod zdrojového kódování (komprese) obsahují knihy Davida Salomona [4,5]. Některé relevantní informace lze nalézt i na Wikipedii nebo jiných online zdrojích.

### Seznam literatury

- [1] Shor, P.: *Lempel-Zip notes* [online]. MIT, 2005. Available from: [http://www-math.mit.edu/~shor/PAM/lempel\\_ziv\\_notes.pdf](http://www-math.mit.edu/~shor/PAM/lempel_ziv_notes.pdf)
- [2] Ziv, J., Lempel, A.: *A Universal Algorithm for Sequential Data Compression*. IEEE Transactions on Information Theory 23 (3), 1977, pp. 337–343. doi:10.1109/TIT.1977.1055714.
- [3] Welch, T.: *A Technique for High-Performance Data Compression*. Computer 17 (6), 1984, pp. 8–19. doi: 10.1109/MC.1984.1659158.
- [4] Salomon, D.: *Concise Introduction to Data Compression*. London: Springer Verlag, 2008, 311 pp.
- [5] Salomon D., Motta G.: *Handbook of Data Compression*. London: Springer Verlag, 2010, 1360 pp. doi: 10.1007/10.1007/978-1-84882-903-9



# Slavnitars' metody komprese

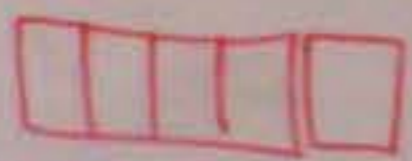
Princip: Kodujeme "fráze" z nějaké tabulky na kódové symboly  
 (tabulka  $\equiv$  slovník, vzniká dynamicky)

Výstup: kódová slova / symboly se stejnou pravděpodobností výskytu

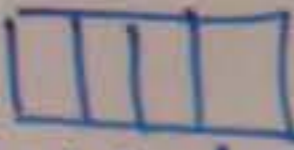
Předchůdce je **RLE** (Run-Length Encoding)

Př: 0000001111100000011000000001111

6, 0, 5, 1, 6, 0, 2, 1, 8, 0, 4, 1



1 bit na hodnotu



4 bity na čítec

⊕ pravidlo: začínáme v 0

6, 5, 6, 2, 8, 4, 1, ...

4 bity na čítec

Ex: 1 000...00111

16x  
 0, 1, 15, 0, 1, 1, 3

0x nula na zač.

16x nula → 15x nula  
 0x jedna  
 1x nula

musí jít o binární formát  
 → třeba 16 barev (nebo 256)  
 → PCX

→ nevhodné pro kratší sekvence!



# LZ77 (Lempel-Ziv)

- textová komprese
- používají posuvné okno jako slovník

Ex: abbcbabcdbbbab...

posuvné okno délky 8 znaků

pozice 0: abbcba**bc** 8 znaků  
 pozice 4: babc**db**ba 8 znaků

Princip: text na výstupu nahrazují odkazy na již zpracovaný text v posuvném okně

Příklad: Kdyby, byly, ryby, nebyly, by, chyby.

$(0,0,k)$   $(0,0,d)$   $(0,0,y)$   $(0,0,b)$   $(2,1,u)$   $(3,2,l)$   $(5,2,r)$   $(10,4,n)$

↑ počet znaků na vyhledávání  
 posuv zřet 0 tolik znaků

$(0,0,e)$   $(12,5,b)$   $(3,2,c)$   $(0,0,h)$   $(16,3,y)$

→ není žádný odkaz na předchozí výstyt

32 znaků  $\rightsquigarrow$  13 trojic, které "nějak" reprezentují

odkaz  $(8 \text{ bitů}, 8 \text{ bitů}, 8 \text{ bitů}) \Rightarrow 13 \times 24 \text{ bitů} = 312$   
 $(5 \text{ bitů}, 3 \text{ bitů}, 7 \text{ bitů}) \Rightarrow 13 \times 15 \text{ bitů} = 195$

$32 \times 8 \text{ bitů} = 256$   
 $32 \times 7 \text{ bitů} = 224$

Nedostatek: je třeba dlouhé okno  $\Rightarrow$  pomalejší  
 dlouhé okno  $\Rightarrow$  delší reprezentace trojic



# LZW (Lempel-Ziv-Welch)

- používá adaptivní slovník, který rozšiřuje  
nějaký předdefinovaný slovník základní  
(slovník lze po komprimaci zakodit a při  
rozbalování zprávy vytvořit znovu)

Princip: 1. začni se základním slovníkem

2. opakuj:

- přidej zakódovanou frázi na výstup
- rozšíř slovník o frázi, zahrnující  
dalsí písmeno

K na výstup a 'Kd' jako novou  
frázi

Příklad:

Kdyby, byly, ryby, nebyly, by, chyby.

Slovník:

K → 0  
d → 1  
y → 2  
b → 3  
l → 4  
r → 5  
n → 6  
e → 7  
c → 8  
h → 9  
· → 10  
→ 11

Kd → 12

dy → 13

yb → 14

by → 15

yl → 16

lb → 17

bl → 18

ly → 19

ylr → 20

ry → 21

yby → 22

yn → 23

nr → 24

eb → 25

byly → 26

yub → 27

byl → 28

lc → 29

ch → 30

by → 31

yby. → 32

K	d	y	b	y	l	by	l	yl	r	yb
0	1	2	3	2	4	15	5	16	6	14
yl	n	e	byl	yl	by	l	c	b		
16	7	8	18	16	15	4	9	10		
		yby	.	EOT						
		22	11							

převodně: 32 x 8 bitů ... 256  
→ bitů ... 224

LZW: 22 x 6 bitů ... 132

- K → 0  
add('Kd', 12)
- d → 1  
add('dy', 13)
- ...

základní