

- **Lempel-Ziv-Welchův algoritmus (1978)**
- Kompresní/dekompresní metoda, využití – přenos dat, GIF, TIFF, postscript
- (+) rychlá metoda, (-) horší kompresní poměr (cca o 30% horší než nejlepší met.)
- Slovník – ukládají se opakující se znaky
- Pokud se vyskytne několik stejných posloupností znaků – nahrazení stejným číslem
- Slovník se neukládá, ale je znovu vytvořen ze zakódovaného souboru

- Př. řetězec abcbcabcbcbba (13 znaků), vstupní abeceda (použité znaky): a b c  
přiřazení čísel a = 1, b = 2, c = 3

1. Nalezení nejdelší fráze ze slovníku shodné se vstupem, index na výstup, odebrání
2. Nová fráze = nalezená fráze + jeden další znak

Krok	text vstupu	nalezená fráze	výstup	nová fráze	Index nov. fráze
1	abcbcabcbcba	a	1	ab	4
2	bcabcbcabcbcba	b	2	bc	5
3	cabcbcabcbcba	c	3	ca	6
4	abcbcabcbcba	ab	4	abc	7
5	cabcbcbcba	ca	6	cab	8
6	bcbcbcba	bc	5	bc b	9
7	bcba	bcb	9		
8	a	a	1		

Výstup: řetězec 12346591 (8 čísel), kompresní poměr 61,5%



# Huffmanovo kódování

- **Huffmanovo kódování**

- Algoritmus navržen Davidem Huffmanem (1952)

- Využití prefixového kódu - kód žádného znaku není prefixem jiného znaku, neprefixový kód: Morseova abeceda: A (.-), M(--), J(.---)

- Proměnná délka kódových slov >>> „znaky“, které jsou nejvíce četné mají nejkratší délku a naopak: 111011011, A (0), E(10), G(11)

- **Algoritmus kódování:**

1. Zjištění četnosti jednotlivých „znaků“
2. Vytvoření jednotlivých kódů na základě četností
3. Nahrazení jednotlivými znaky v datovém souboru nalezenými kódy

**(+) Výhody** – rychlá komprese a dekomprese, nenáročné na paměť

**(-) Nevýhody** – nutnost nalezených kódů, menší kompresní poměr

# Huffmanovo kódování

CV06

- **Příklad:** znakový řetězec ABRAKADABRA
- 1) četnosti A (5x – 0,46), R (2x – 0,18), B (2x – 0,18), K (1x – 0,09), D (1x – 0,09)
- 2) vytvoření tabulky dle četností
- 3) poslední dvě četnosti se sečtou a zařadí se do tabulky, sčítá se až do 1

A 0,46		A 0,46		A 0,46		KDBR 0,54	1
R 0,18		R 0,18		KDB 0,36	1	A 0,46	0
B 0,18		B 0,18	1	R 0,18	0		
K 0,09	1	KD 0,18	0				
D 0,09	0						

- 4) Posledním dvěma slovům v každém sloupci tabulce přiřadíme 1 (vyšší četnost) a 0 (nižší četnost)
- 5) Výsledný kód znaku >>> posloupnost 0 a 1 dle toho jak se znak seskupoval s dalšími znaky, např. pro znak K: (1) – KDBR, (1) KDB, (0) KD a (1) K. Výsledné kódy A (0), R (10), B (111), K (1101), D (1100)
- 6) Výsledný řetězec pro ABRAKADABRA: 0 111 10 0 1101 0 1100 0 111 10 0, tj.: 01111001101011000111100
- Kompresní poměr (pokud 1 znak = 8 bit.): 0,26 (26%)



# Aritmetické kódování

- **Aritmetické kódování**
- Huffmanův kód >>> problém při stejné pravděpodobnosti výskytu >>> možné řešení >>> aritmetické kódování
- Pro bezztrátovou kompresi dat, proměnná délka kódových slov jako u Huffmanova kódování, při kódování se však vstupní znak nenahrazuje specifickým kódem, ale výsledek, tj. vstupní datový řetězec se nahradí reálným číslem z intervalu  $<0,1)$ .
- **Algoritmus kódování:**
  1. Zjištění četnosti (pravděpodobnosti výskytu) jednotlivých „znaků“
  2. Dle pravděpodobnosti výskytu znaků se umístí znak v intervalu  $<0,1)$
  3. Celý interval  $<0,1)$  je postupně omezován z obou stran na základě přicházejících znaků.
  4. Každý znak vybere z aktuálního intervalu odpovídající poměrnou část >>> nový základ pro následující symbol.
  5. Po průchodu (načtení) všech znaků dostáváme podinterval z intervalu  $<0,1)$ , výsledkem je pak libovolné reálné číslo z tohoto intervalu.
  6. Na konec kódované zprávy dáme speciální znak, jinak při dekódování není možné určit konec datového toku, nebo uložíme délku původní posloupnosti znaků

# Aritmetické kódování

- **Aritmetické kódování** – př.
- Datový řetězec CBAABCADAC (10 znaků)
- 1) Pravděpodobnosti výskytu A – 0.4 (P1), B – 0.2 (P2), C – 0.3 (P3), D – 0.1 (P4)
- 2) rozdělení v intervalu  $<0, 1)$ :  
 $<0, P1)$ ,  $<P1, P1 + P2)$ ,  $<P1 + P2, P1 + P2 + P3)$ ,  $<P1 + P2 + P3, P1 + P2 + P3 + P4)$

Kumulativní pravděpodobnosti:

$$Q_0 = 0, Q_1 = P_1, Q_2 = P_1 + P_2, \dots, Q_N = P_1 + P_2 + \dots + P_N = 1$$

$<0, Q_1)$ ,  $<Q_1, Q_2)$ ,  $<Q_2, Q_3)$ ,  $<Q_3, Q_4)$

$<0, 0.4)$ ,  $<0.4, 0.6)$ ,  $<0.6, 0.9)$ ,  $<0.9, 1)$

**A**

**B**

**C**

**D**

# Aritmetické kódování

CV07

- **Aritmetické kódování** – př.
- Datový řetězec CBAABCADAC (10 znaků)
- Rozdělení intervalu

$\langle 0, 0.4 \rangle$ ,  $\langle 0.4, 0.6 \rangle$ ,  $\langle 0.6, 0.9 \rangle$ ,  $\langle 0.9, 1 \rangle$   
**A**                      **B**                      **C**                      **D**

- 3) Kódování >>> postupné omezování intervalu  $I = \langle 0, 1 \rangle$ ,  $I = \langle L, H \rangle$

>>> Postupně jsou brány znaky z datového řetězce, k nim známe  $I_Z = \langle Z_L, Z_H \rangle$

>>> Nová hodnota intervalu  $I_N = \langle L + Z_L \cdot (H - L), L + Z_H \cdot (H - L) \rangle$

C >>>  $I = \langle 0, 1 \rangle$ ,  $I_N = \langle 0 + 0.6 \cdot (1 - 0), 0 + 0.9 \cdot (1 - 0) \rangle = \langle 0.6, 0.9 \rangle$

B >>>  $I = \langle 0.6, 0.9 \rangle$ ,  $I_N = \langle 0.6 + 0.4 \cdot (0.9 - 0.6), 0.6 + 0.6 \cdot (0.9 - 0.6) \rangle = \langle 0.72, 0.78 \rangle$

A >>>  $I = \langle 0.72, 0.78 \rangle$ ,  $I_N = \langle 0.72 + 0 \cdot (0.78 - 0.72), 0.72 + 0.4 \cdot (0.78 - 0.72) \rangle = \langle 0.72, 0.744 \rangle$

A >>>  $I = \langle 0.72, 0.744 \rangle$ ,  $I_N = \langle 0.72 + 0 \cdot (0.744 - 0.72), 0.72 + 0.4 \cdot (0.744 - 0.72) \rangle = \langle 0.72, 0.7296 \rangle$

B >>>  $I = \langle 0.72, 0.7296 \rangle$ ,  $I_N = \langle 0.72 + 0.4 \cdot (0.7296 - 0.72), 0.72 + 0.6 \cdot (0.7296 - 0.72) \rangle = \langle 0.72384, 0.72576 \rangle$

....

....

B >>>  $I = \langle 0.72519936, 0.725208576 \rangle$ ,  $I_N = \langle 0.7252048896, 0.7252076544 \rangle$ , C = 0.725205



# Aritmetické kódování

- **Aritmetické dekódování** – př.
- Datový řetězec CBAABCADAC (10 znaků)
- Rozdělení intervalu

$\langle 0, 0.4 \rangle$ ,  $\langle 0.4, 0.6 \rangle$ ,  $\langle 0.6, 0.9 \rangle$ ,  $\langle 0.9, 1 \rangle$   
**A**                      **B**                      **C**                      **D**

- Výsledek kódování >>>  $C = 0.725205$

- 1) Počáteční hodnota intervalu dekódování  $I = \langle 0, 1 \rangle$
- 2) dekódování znaku:  $K = ((C - L) / (H - L))$ ;  $ZL \leq K < ZH$  >>> nalezneme odpovídající znak
- 3) počítáme nový interval  $IN = \langle L + ZL \cdot (H - L), L + ZH \cdot (H - L) \rangle$   
 $I = \langle 0, 1 \rangle$ ,  $K = 0.725205$ , znak = C,  $IN = \langle 0 + 0.6 \cdot (1 - 0), 0 + 0.9 \cdot (1 - 0) \rangle = \langle 0.6, 0.9 \rangle$   
 $I = \langle 0.6, 0.9 \rangle$ ,  $K = 0.41735$ , znak = B,  $IN = \langle 0.6 + 0.4 \cdot (0.9 - 0.6), 0.6 + 0.6 \cdot (0.9 - 0.6) \rangle = \langle 0.72, 0.78 \rangle$   
 $I = \langle 0.72, 0.78 \rangle$ ,  $K = 0.08675$ , znak = A,  $IN = \langle 0.72, 0.744 \rangle$   
 $I = \langle 0.72, 0.744 \rangle$ ,  $K = 0.216875$ , znak = A,  $IN = \langle 0.72, 0.7296 \rangle$   
 ....

# Filtr Move-To-Front (MTF)

- 1980 (B. Ryabko), 1986 (J.K. Bentley) Přesuň na začátek, po provedení Burrowsovy-Wheelerovy transformace
- Dlouhé sekvence symbolů nahrazeny nulami, málo se vyskytuje: velké číslo
- Nahrazení symbolů ze vstupního řetězce jejich indexy ze zásobníku
- zásobník s čísli (pořadí, ASCII...), *i*-tá pozice – kód *i* z pořadí, ASCII
- Aktuálně kódovaný znak přesunut na počátek (dopředu)
- A až Z (26 čísel v zásobníku, + české znaky 42 písmen)
- ANANAS bude zakódován jako 1, 14, 2, 2, 2, 19
- A – 1 (poté: ABCDEFGHIJKLMNOP...)
- N – 14 (poté: NABCDEFGHIJKLMO...)
- A – 2 (poté: ANBCDEFGHIJKLMO...)
- N – 2 (poté: NABCDEFGHIJKLMO...)
- A – 2 (poté: ANBCDEFGHIJKLMO...)
- S – 19 (poté: SANBCDEFGHIJKLMO...)

1	A	14	N
2	B	15	O
3	C	16	P
4	D	17	Q
5	E	18	R
6	F	19	S
7	G	20	T
8	H	21	U
9	I	22	V
10	J	23	W
11	K	24	X
12	L	25	Y
13	M	26	Z

# Filtr Move-To-Front (MTF)

CV11

## ● ANANAS – ASCII

- A – 65, (poté A = 0)
- N – 78, (poté N = 0, A = 1)
- A – 1 (poté A = 0, N = 1)
- N – 1 (poté N = 0, A = 1)
- A – 1 (poté A = 0, N = 1)
- S – 83 (poté S = 0, A = 1...)

000	00		043	2B	+	086	56	U	129	81	ü	172	AC	¼	215	D7	
001	01	☐	044	2C	,	087	57	W	130	82	é	173	AD	↓	216	D8	½
002	02	☐	045	2D	-	088	58	X	131	83	â	174	AE	«	217	D9	¾
003	03	♥	046	2E	.	089	59	Y	132	84	ä	175	AF	»	218	DA	█
004	04	♦	047	2F	/	090	5A	Z	133	85	à	176	B0	▨	219	DB	▩
005	05	♠	048	30	0	091	5B	[	134	86	á	177	B1	▩	220	DC	▩
006	06	♣	049	31	1	092	5C	\	135	87	ç	178	B2	▩	221	DD	▩
007	07	•	050	32	2	093	5D	]	136	88	ê	179	B3	▩	222	DE	▩
008	08	■	051	33	3	094	5E	^	137	89	ë	180	B4	▩	223	DF	▩
009	09		052	34	4	095	5F	`	138	8A	è	181	B5	▩	224	E0	α
010	0A		053	35	5	096	60		139	8B	í	182	B6	▩	225	E1	β
011	0B	ø	054	36	6	097	61	a	140	8C	î	183	B7	▩	226	E2	γ
012	0C	♀	055	37	7	098	62	b	141	8D	ï	184	B8	▩	227	E3	π
013	0D		056	38	8	099	63	c	142	8E	ä	185	B9	▩	228	E4	Σ
014	0E	∏	057	39	9	100	64	d	143	8F	å	186	BA	▩	229	E5	σ
015	0F	*	058	3A	:	101	65	e	144	90	é	187	BB	▩	230	E6	μ
016	10	↳	059	3B	;	102	66	f	145	91	æ	188	BC	▩	231	E7	τ
017	11	◀	060	3C	<	103	67	g	146	92	æ	189	BD	▩	232	E8	θ
018	12	±	061	3D	=	104	68	h	147	93	ô	190	BE	▩	233	E9	θ
019	13	!!	062	3E	>	105	69	i	148	94	ö	191	BF	▩	234	EA	Ω
020	14	¶	063	3F	?	106	6A	j	149	95	ò	192	C0	▩	235	EB	δ
021	15	§	064	40	@	107	6B	k	150	96	û	193	C1	▩	236	EC	ω
022	16	—	065	41	A	108	6C	l	151	97	ù	194	C2	▩	237	ED	ø
023	17	‡	066	42	B	109	6D	m	152	98	ÿ	195	C3	▩	238	EE	€
024	18	†	067	43	C	110	6E	n	153	99	ü	196	C4	▩	239	EF	π
025	19	↓	068	44	D	111	6F	o	154	9A	ü	197	C5	▩	240	F0	≡
026	1A		069	45	E	112	70	p	155	9B	ç	198	C6	▩	241	F1	±
027	1B	←	070	46	F	113	71	q	156	9C	£	199	C7	▩	242	F2	≤
028	1C	↳	071	47	G	114	72	r	157	9D	¥	200	C8	▩	243	F3	≤
029	1D	↗	072	48	H	115	73	s	158	9E	ℓ	201	C9	▩	244	F4	∫
030	1E	▲	073	49	I	116	74	t	159	9F	f	202	CA	▩	245	F5	∫
031	1F	▼	074	4A	J	117	75	u	160	A0	á	203	CB	▩	246	F6	÷
032	20		075	4B	K	118	76	v	161	A1	í	204	CC	▩	247	F7	÷
033	21	!	076	4C	L	119	77	w	162	A2	ó	205	CD	▩	248	F8	°
034	22	"	077	4D	M	120	78	x	163	A3	ú	206	CE	▩	249	F9	°
035	23	#	078	4E	N	121	79	y	164	A4	ñ	207	CF	▩	250	FA	·
036	24	\$	079	4F	O	122	7A	z	165	A5	ñ	208	D0	▩	251	FB	√
037	25	%	080	50	P	123	7B	<	166	A6	ª	209	D1	▩	252	FC	²
038	26	&	081	51	Q	124	7C	!	167	A7	ª	210	D2	▩	253	FD	²
039	27	'	082	52	R	125	7D	>	168	A8	¿	211	D3	▩	254	FE	²
040	28	<	083	53	S	126	7E	~	169	A9	¿	212	D4	▩	255	FF	²
041	29	>	084	54	T	127	7F	Δ	170	AA	¿	213	D5	▩			
042	2A	*	085	55	U	128	80	Ç	171	AB	½	214	D6	▩			

# Filtr Move-To-Front (MTF)

## ● Dekódování

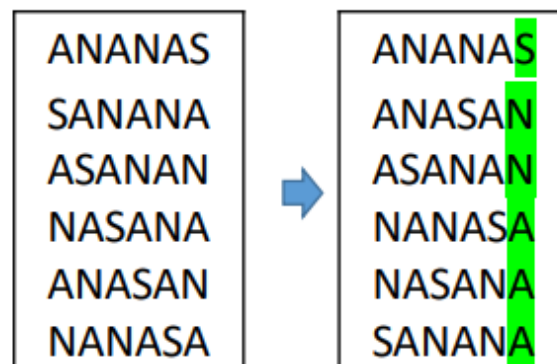
- 1 - A (poté: ABCDEFGHIJKLMNOP...)
- 14 - N (poté: NABCDEFGHIJKLMO...)
- 2 - A (poté: ANABCDEFGHIJKLMO...)
- 2 - N (poté: NABCDEFGHIJKLMO...)
- 2 - A (poté: ANABCDEFGHIJKLMO...)
- 19 - S (poté: SNABCDEFGHIJKLMO...)

1	A	14	N
2	B	15	O
3	C	16	P
4	D	17	Q
5	E	18	R
6	F	19	S
7	G	20	T
8	H	21	U
9	I	22	V
10	J	23	W
11	K	24	X
12	L	25	Y
13	M	26	Z



# Burrows-Wheelerova Transformace (BWT)

- Michael Wheeler (1984) + David Burrows (1994)
- Přeskupí symboly tak, že shodné budou s vysokou pravděpodobností vedle sebe, pak můžeme použít RLE...
- Bzip2 – BWT >>> MTF >>> Huffmanovo kódování
- Vstupní řetězec se zpracovává po blocích o konstantní velikosti, větší bloky – větší efektivita, blok cca desítky tisíc symbolů (Block Sorting Algorithm)
- 1) blok o délce N (např. ANANAS), N cyklických posuvů o jeden symbol doprava >>> matice N \* N, řádky => cyklické posuny
- 2) lexikografické (podle abecedy) seřídění matice
- 3) poslední sloupec je výsledek transformace
- 4) přidáme číslo řádku v seříděné matici, kde se nachází vstupní řetězec: **SNNAAA, 1. ŘÁDEK**





# Burrows-Wheelerova Transformace (BWT)

- Dekódování: SNNAAA, 1. ŘÁDEK

1	sort	2	sort	3	sort	4	sort	5	sort	6	sort
S	A	SA	AN	SAN	ANA	SANA	ANAN	SANAN	ANANA	SANANA	ANANAS
N	A	NA	AN	NAN	ANA	NANA	ANAS	NANAS	ANASA	NANASA	ANASAN
N	A	NA	AS	NAS	ASA	NASA	ASAN	NASAN	ASANA	NASANA	ASANAN
A	N	AN	NA	ANA	NAN	ANAN	NANA	ANANA	NANAS	ANANAS	NANASA
A	N	AN	NA	ANA	NAS	ANAS	NASA	ANASA	NASAN	ANASAN	NASANA
A	S	AS	SA	ASA	SAN	ASAN	SANA	ASANA	SANAN	ASANAN	SANANA

- VÝSTUP: ANANAS

# BWT + MTF

IDK bylo to  
v  
přednášce  
tak jsem to  
sem raději  
dal

## ● BWT Dekódování: SNNAAA, 1. ŘÁDEK

## ● MTF:

S – 83 (poté: S = 0)

N – 78 (poté: N = 0, S = 1)

N – 0

A – 65 (poté: A = 0, N = 1, S = 2)

A – 0

A – 0

pro Huffmanovo kódování:

0 – 3x

65 – 1x

78 – 1x

83 – 1x

U delších řetězců – vyšší počet 0

000	00		043	2B	+	086	56	U	129	81	ü	172	AC	¼	215	D7	
001	01	⊖	044	2C	,	087	57	U	130	82	é	173	AD	↓	216	D8	
002	02	⊖	045	2D	-	088	58	X	131	83	â	174	AE	«	217	D9	
003	03	♥	046	2E	.	089	59	Y	132	84	ã	175	AF	»	218	DA	
004	04	♦	047	2F	/	090	5A	Z	133	85	ä	176	B0		219	DB	
005	05	♠	048	30	0	091	5B	[	134	86	å	177	B1		220	DC	
006	06	♠	049	31	1	092	5C	\	135	87	ç	178	B2		221	DD	
007	07	•	050	32	2	093	5D	]	136	88	ê	179	B3		222	DE	
008	08	■	051	33	3	094	5E	^	137	89	ë	180	B4		223	DF	
009	09		052	34	4	095	5F	`	138	8A	è	181	B5		224	E0	α
010	0A	δ	053	35	5	096	60		139	8B	ï	182	B6		225	E1	β
011	0B	δ	054	36	6	097	61	a	140	8C	î	183	B7		226	E2	Γ
012	0C	♀	055	37	7	098	62	b	141	8D	ì	184	B8		227	E3	Π
013	0D		056	38	8	099	63	c	142	8E	ñ	185	B9		228	E4	Σ
014	0E		057	39	9	100	64	d	143	8F	â	186	BA		229	E5	σ
015	0F	✱	058	3A	:	101	65	e	144	90	é	187	BB		230	E6	μ
016	10	▷	059	3B	;	102	66	f	145	91	æ	188	BC		231	E7	τ
017	11	◄	060	3C	<	103	67	g	146	92	æ	189	BD		232	E8	δ
018	12	±	061	3D	=	104	68	h	147	93	ô	190	BE		233	E9	θ
019	13	!!	062	3E	>	105	69	i	148	94	ö	191	BF		234	EA	Ω
020	14	¶	063	3F	?	106	6A	j	149	95	ò	192	C0		235	EB	δ
021	15	§	064	40	0	107	6B	k	150	96	û	193	C1		236	EC	ω
022	16	—	065	41	A	108	6C	l	151	97	ù	194	C2		237	ED	ø
023	17	±	066	42	B	109	6D	m	152	98	ÿ	195	C3		238	EE	€
024	18	↑	067	43	C	110	6E	n	153	99	ö	196	C4		239	EF	ñ
025	19	↓	068	44	D	111	6F	o	154	9A	ü	197	C5		240	F0	≡
026	1A		069	45	E	112	70	p	155	9B	ç	198	C6		241	F1	±
027	1B	←	070	46	F	113	71	q	156	9C	É	199	C7		242	F2	≥
028	1C	↳	071	47	G	114	72	r	157	9D	¥	200	C8		243	F3	≤
029	1D	↔	072	48	H	115	73	s	158	9E	ℓ	201	C9		244	F4	ƒ
030	1E	▲	073	49	I	116	74	t	159	9F	f	202	CA		245	F5	J
031	1F	▼	074	4A	J	117	75	u	160	A0	á	203	CB		246	F6	÷
032	20		075	4B	K	118	76	v	161	A1	í	204	CC		247	F7	≈
033	21	†	076	4C	L	119	77	w	162	A2	ó	205	CD		248	F8	°
034	22	″	077	4D	M	120	78	x	163	A3	ú	206	CE		249	F9	·
035	23	‡	078	4E	N	121	79	y	164	A4	ñ	207	CF		250	FA	√
036	24	§	079	4F	O	122	7A	z	165	A5	ñ	208	D0		251	FB	√
037	25	‰	080	50	P	123	7B	¿	166	A6	â	209	D1		252	FC	√
038	26	‰	081	51	Q	124	7C	¡	167	A7	é	210	D2		253	FD	√
039	27	’	082	52	R	125	7D	¢	168	A8	ç	211	D3		254	FE	√
040	28	‹	083	53	S	126	7E	£	169	A9	ç	212	D4		255	FF	√
041	29	›	084	54	T	127	7F	¤	170	AA	ç	213	D5				
042	2A	✱	085	55	U	128	80	¥	171	AB	½	214	D6				

# Bezpečnostní kódování

CV12

- **Inverzní kód**

- Vysílaný kód se při vysílání opakuje beze změn, pokud je v něm sudý počet 0, pokud máme lichý počet, tak se opakovaný kód invertuje,  $D = 4$ , lze opravit jednoduchou chybu

sudý poč. 0: 01011 01011	vysláno:	01011 01011	01001 10110	10111 01000
lichý poč. 0: 10111 01000	přijato:	00011 01011	01011 10110	10111 00000
	překódováno:	00011 10100	01011 10110	10111 11111
	oprava:	00011 $\oplus$ 10100 10111 chyba infor. části kódu	01011 $\oplus$ 10110 11101 chyba infor. části kódu	10111 $\oplus$ 11111 01000 chyba zab. části kódu

- **Dekódování**

- Přijatý kód >>> kontrola >>> není správně >>> oprava – dekodování, Není to inverzní postup ke kódování
- Máme kódy 000000, 111111, alespoň 4 znaky jsou správně 100100 >>> 000000