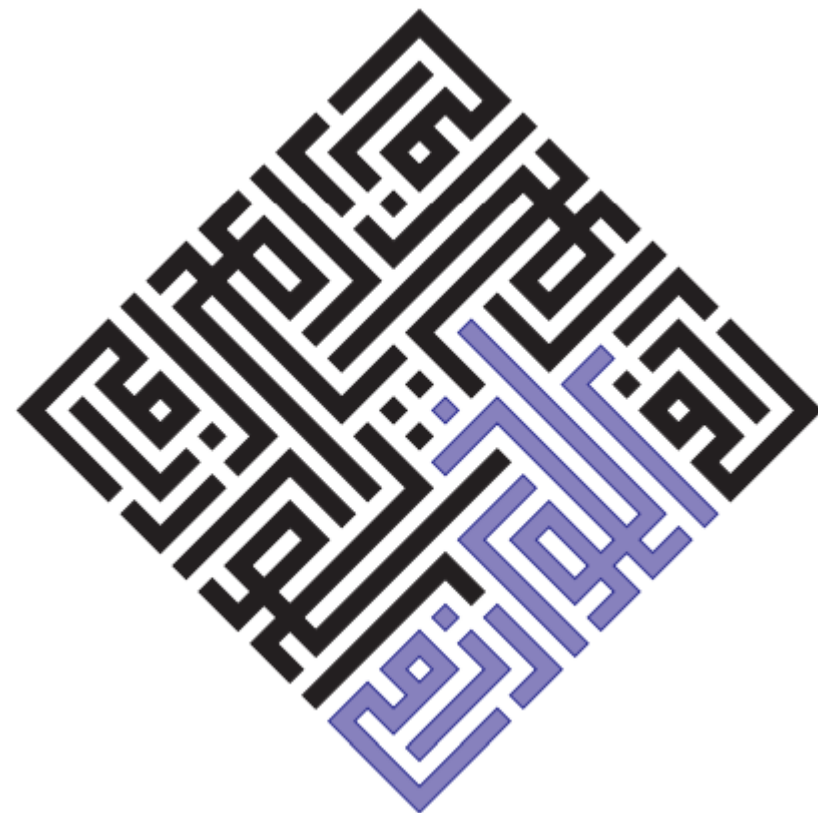


TÖL403G GREINING REIKNIRITA

11. Gráðug reiknirit 2

Hjálmtyr Hafsteinsson
Vor 2022



- Huffman kóðar (*Huffman codes*)
 - Gagnabjöppun með kóðum
 - Forstrengsfríir kóðar (*prefix-free codes*)
- Almennt um gagnabjöppun
 - Notkun á Huffman kóðun
 - Aðrar bjöppunaraðferðir

4.4

- Skrá hefur 100.000 bókstafi af aðeins 6 mismunandi gerðum

Stafur	Tíðni	
<i>a</i>	45%	Ef stafirnir eru geymdir á 8-bita ASCII formi þá þarf $8 \cdot 100.000 = \underline{800.000 \text{ bita}}$ fyrir skrána
<i>b</i>	13%	
<i>c</i>	12%	Við gætum reyndar komist af með 3 bita per staf, þar sem það eru aðeins 6 ólíkir stafir
<i>d</i>	16%	
<i>e</i>	9%	Til dæmis: <i>a</i> : 000, <i>b</i> : 001, o.s.frv.
<i>f</i>	5%	
		Þá þarf aðeins <u>300.000 bita</u> fyrir skrána

Þetta er það besta sem við getum gert með föstum fjölda bita per staf

En hvað með að nota mismunandi marga bita eftir því hversu algengir stafirnir eru?

Notum færri bita fyrir algenga stafi og fleiri bita fyrir sjaldgæfa stafi

Stafur	Kóði	Tíðni*Lengd	Meðallengd
<i>a</i>	0	$0.45 \cdot 1$	0.45
<i>b</i>	101	$0.13 \cdot 3$	0.39
<i>c</i>	100	$0.12 \cdot 3$	0.36
<i>d</i>	111	$0.16 \cdot 3$	0.48
<i>e</i>	1101	$0.09 \cdot 4$	0.36
<i>f</i>	1100	$0.05 \cdot 4$	0.20

Samtals: **2.24** bitar að meðaltali

Hér þurfum við því aðeins 224.000 bita fyrir skrána

Breytileg kóðalengd krefst þess að kóðinn sé forstrengsfrír (*prefix-free*)

Kóði er forstrengsfrír ef enginn stafakóði er forstrengur fyrir annan stafakóða

Í þessu dæmi þá hefur *a* kóðann 0, það þýðir að enginn annar kóði getur byrjað á 0

Forstrengsfríir (*prefix-free*) kóðar

- Kóðar af fastri lengd eru alltaf forstrengsfríir
 - Það er engin leið fyrir einn stafakóða að vera forstrengur annars, nema þeir séu sami kóðinn!
- Kóðar af breytilegri lengd verða að vera forstrengfríir
 - Annars er ómögulegt að afkóða þá

Stafur	Kóði
<i>a</i>	0
<i>b</i>	101
<i>c</i>	100
<i>d</i>	111
<i>e</i>	1101
<i>f</i>	1100

Til dæmis: Ef stafurinn *b* hefði haft kóðann 01 í dæminu á undan

Fáum bitarununa 0100101

Er fyrsti bitinn (0) afkóðaður sem *a*, eða
eru fyrstu tveir bitarnir (01) afkóðaðir sem *b*?

Eins og kóðinn er í dæminu þá á að afkóða þetta sem *acb* (0 100 101)

Morse kóði

← Fyrst notaður árið 1844, staðlaður 1865

- Gamalt dæmi um kóða af breytilegri lengd
 - Algengir stafir (E, T, I, ...) hafa stuttan kóða
- Hann virðist ekki vera forstrengsfrír
 - Kóðinn fyrir E (●) er forstrengur fyrir marga stafi, A, H, I, ...
- En í raun er Morse kóði forstrengsfrír þríundarkóði
 - Það verður alltaf að vera eyða á milli stafanna
 - Eyðan (□) er því þriðja táknið

Kóðinn fyrir E er því ●□ og er ekki forstrengur fyrir A: ●—□

Til að auðvelda afkóðun eru settar tvær eyður (□□) á milli orða

International Morse Code

1. The length of a dot is one unit.
2. A dash is three units.
3. The space between parts of the same letter is one unit.
4. The space between letters is three units.
5. The space between words is seven units.

A	● —	U	● ● —
B	— ● ● ●	V	● ● — —
C	— — ● —	W	— — ● —
D	— ● ●	X	— ● ● —
E	●	Y	— ● — —
F	● ● — ●	Z	— — ● ●
G	— — ●		
H	● ● ● ●		
I	● ●		
J	● — — —		
K	— ● — —		
L	● — — ●		
M	— —		
N	— ●		
O	— — —		
P	● — — ●		
Q	— — ● —		
R	● — — ●		
S	● ● ●		
T	—		
		1	● — — — —
		2	● ● — — —
		3	● ● ● — —
		4	● ● ● ● —
		5	● ● ● ● ●
		6	— ● ● ● ●
		7	— — ● ● ●
		8	— — — ● ●
		9	— — — — ●
		0	— — — — —

"Greining reiknirita" í [Morse kóða](#):



UTF-8 (*Unicode Transformation Format*)

- Hannaður af Ken Thompson og Rob Pike
 - Tilgangur: Þurfa ekki að nota 2 eða 4 bæti fyrir algeng ensk tákn
 - Getur táknað öll (~1.1 milljón) Unicode táknin með 1 til 4 bætum

Ef kóðinn er aðeins eitt bæti þá er fyrsti bitinn **0**-biti

Annars byrjar fyrsta bætið í kóðanum alltaf á **11**

Næstu bitar á eftir þeim segja til um hversu mörg bætin eru:

- 0**: tvö bæti
- 10**: þrjú bæti
- 110**: fjögur bæti

Framhaldsbætin byrja alltaf á **10**

Þau eru einu bætin sem byrja á **10**

Number of bytes	Bits for code point	First code point	Last code point	Byte 1	Byte 2	Byte 3	Byte 4
1	7	U+0000	U+007F	0xxxxxxx			
2	11	U+0080	U+07FF	110xxxxx	10xxxxxx		
3	16	U+0800	U+FFFF	1110xxxx	10xxxxxx	10xxxxxx	
4	21	U+10000	U+10FFFF	11110xxx	10xxxxxx	10xxxxxx	10xxxxxx

Nokkuð stór hluti bitanna er notaður til að tryggja kóðinn sé forstrengsfrír

Ef 4 bæti: 21 gagnabiti og 11 kóðabitar
þ.e. aðeins ~ 66% bitanna eru gagnabitar

- Hvernig búum við til besta kóðann miðað við gefna stafatíðni?

Með gráðugu reikniriti!

David Huffman fann þetta reiknirit árið 1951

Hann var nemendi í námskeiði í MIT og þetta var verkefni í námskeiðinu

Gróf lýsing reikniritsins: Sameina tvö sjaldgæfustu táknin í eitt tákn og halda áfram endurkvæmt þar til aðeins eitt tákn eftir

Getum táknað forstrengsfría kóða sem tvíundartré:

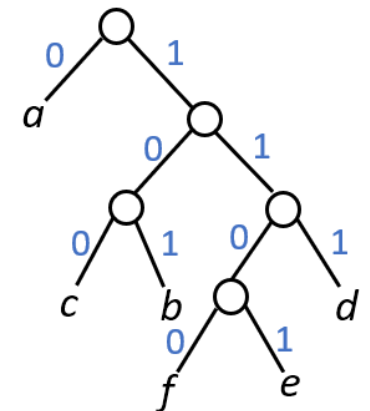
Stafir eru í laufunum

Hver leggur er merktur með 0 eða 1

Vegurinn frá rót í lauf er þá kóði þess stafs

6-stafa kóðinn
frá fyrri glæru:

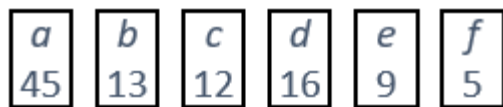
Engir stafir í innri
hnútum, annars væri
kóðinn ekki forstrengsfrír



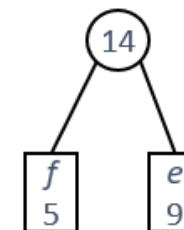
Einfalt sýnidæmi

Í upphafi er hver stafur sjálfstætt tré, sameinum svo eitt og eitt tré í einu

6 tré

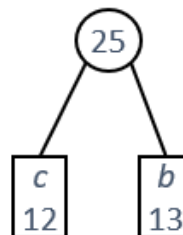


Sameinum tré með lægstu tíðnir, *f* og *e*:

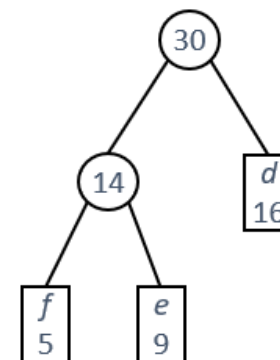


Nú eru 5 tré, með tíðnir 45, 13, 12, 16 og 14

Sameinum tré með lægstu tíðnir, *c* og *b*:

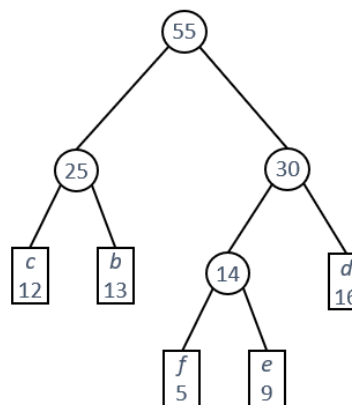


Sameinum næst *d* og tréð með tíðnina 14:

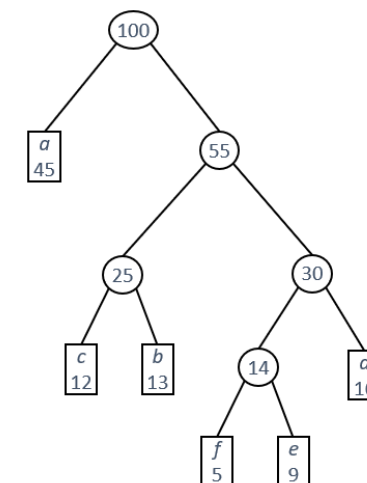


Nú eru 3 tré eftir, með tíðnir 45, 30 og 25

Sameinum tré með tíðnir 25 og 30:



Endum svo með:



- Höfum 3 tákn: a , b , c og þrjár ólíkar kóðanir fyrir þau. Eru þær Huffman kóðanir? Ef svo, teiknið kóðunartréð. Annars, af hverju ekki?
- a : 0, b : 10, c : 11
- a : 0, b : 1, c : 00
- a : 10, b : 01, c : 00

Stærra sýnidæmi

Byrjum með textann:

THISSENTENCECONTAINSTHREEASTHREECSTWODSTWENTYSIXESFIVEFST
HREEGSEIGHTHSTHIRTEENISTWOLSSIXTEENNSNINEOSSIXRSTWENTYSEV
ENSSTWENTYTWOTSTWOUSFIVEVSEIGHTWSFOURXS FIVEYSANDONLYONEZ⁶

Tíðni bókstafa í honum:

A	C	D	E	F	G	H	I	L	N	O	R	S	T	U	V	W	X	Y	Z
3	3	2	26	5	3	8	13	2	16	9	6	27	22	2	5	8	4	5	1

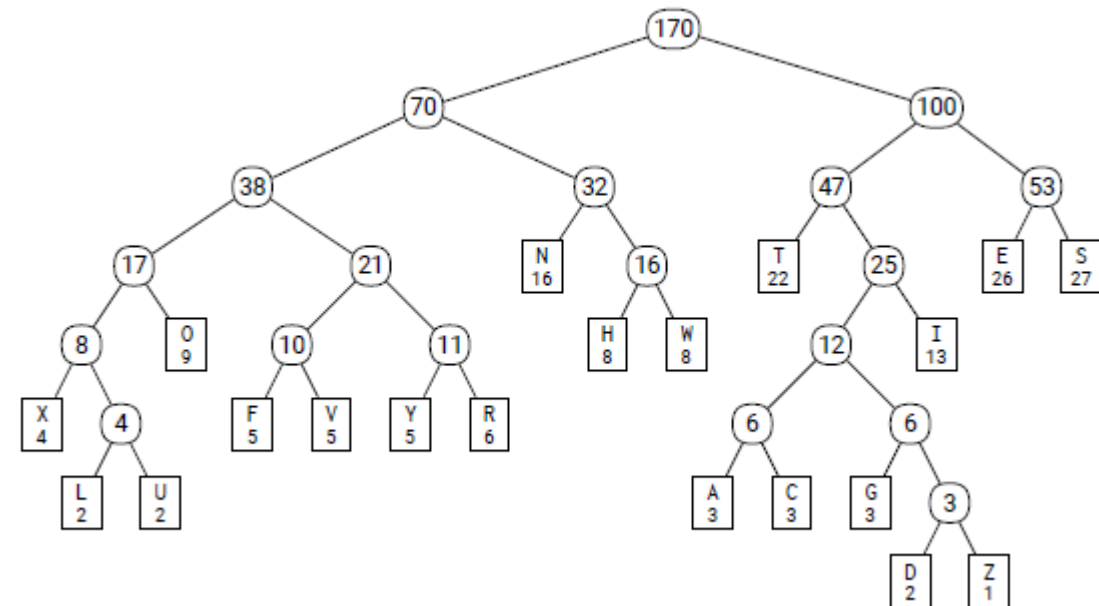
Keyrum reiknirit Huffman og fáum kóðunartréð:

Hér er stysti kóðinn 3 bitar (N, T, E, S)

Lengsti kóðinn er 7 bitar (D, Z)

Textinn hefur 170 bókstafi

Kóðinn er 649 bitar (~3.8 bitar per staf)

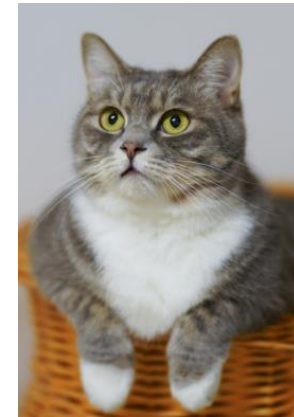


- Geymum táknin (trén) í forgangsbiðröð
 - Finna tvær lægstu tíðnirnar kostar þá $O(\log(n))$
 - Ef við byrjum með tóma forgangsbiðröð, þá
 - þurfum við $2n-1$ Innsetningar (*Insert*)
 - og $2n-2$ EyðaMinnsta (*ExtractMin*)
 - Heildartíminn er þá $O(n \log(n))$
-
- Getum gert kóðun/afkóðun með trénu á $O(m)$ tíma, þar sem m er lengd skilaboðanna

```
BUILDHUFFMAN( $f[1..n]$ ):
  for  $i \leftarrow 1$  to  $n$ 
     $L[i] \leftarrow 0$ ;  $R[i] \leftarrow 0$ 
    INSERT( $i, f[i]$ )
  for  $i \leftarrow n$  to  $2n-1$ 
     $x \leftarrow \text{EXTRACTMIN}()$     ⟨⟨find two rarest symbols⟩⟩
     $y \leftarrow \text{EXTRACTMIN}()$ 
     $f[i] \leftarrow f[x] + f[y]$     ⟨⟨merge into a new symbol⟩⟩
    INSERT( $i, f[i]$ )
     $L[i] \leftarrow x$ ;  $P[x] \leftarrow i$     ⟨⟨update tree pointers⟩⟩
     $R[i] \leftarrow y$ ;  $P[y] \leftarrow i$ 
   $P[2n-1] \leftarrow 0$ 
```

■ Tapandi (*lossy*) þjöppun

- Þjappar gögnum með því að einfalda þau og geyma nálgun á þau (*approximation*)
- Getur náð mjög mikilli þjöppun (10-100 föld) án mikillar sjáanlegrar þjögunar
- Ekki hægt að ná aftur upphaflegu gögnunum (*irreversible*)
- Gerð þjöppunar fer eftir tegund gagna: myndir, tónlist, myndskaið, ...



65.784 bæti



1.998 bæti

~33 föld þjöppun

Dæmi: JPEG, MP3, MPEG, ...

■ Taplaus (*lossless*) þjöppun

- Nýtir sér tölfræðilega umfremd (*redundancy*) til að tákna gögn á styttri máta
- Nær yfirleitt ekki mjög mikilli þjöppun (dæmigert 2-3 föld)
- Þjöppunin er viðsnúanleg, við fáum aftur nákvæmlega sömu gögnin
- Oftast frekar almennar aðferðir sem virka á allar gerðir gagna

Dæmi: Huffman, LZ, TIFF, PNG, ...

- Tvö þekktustu taplausu reikniritin eru LZ77 og LZ78
 - Sett fram af [Abraham Lempel](#) og [Jacob Ziv](#) árin 1977 og 1978
- Notuð mjög víða í þjöppunarforritum
 - GIF-skrár, ZIP-skrár, ýmis netbúnaður, ...
- LZW þjöppunin
 - Terry Welch bjó til sniðuga útfærslu á LZ78
 - Fyrirtækið hans, Sperry Corp. sótti um einkaleyfi á LZW árið 1983
 - LZW var notað í GIF myndaforminu og Unisys (sem keypti Sperry) fór að rukka fyrir það!
 - Varð til þess að PNG myndaformið var skilgreint
 - Einkaleyfið féll úr gildi árið 2002 (20 ára gildistími)

← Bæði "orðasafns"-reiknirit
(*dictionary coders*)

LZ77 notar hliðrandi glugga, finnur endurteknaðar gagnablokkir og setur í staðinn tilvísun á fyrra tilvik

LZ78 býr jafnóðum til orðasafn og setur tilvísun í hana í stað endurtekinna gagnablokka

- Galli við Huffman kóðun:

Þurfum að vita tíðnina á öllum stöfunum til að búa til kóðann

- Hentar þá vel á eftir annari þjöppun (teljum fjölda tákna um leið!)

Í `gzip` er Huffman kóðun beitt á eftir [LZ77](#) aðferðinni

Í `bzip2` er Huffman kóðun beitt á eftir [Burrows-Wheeler](#) aðferðinni

Huffman kóðun er lokaprep í JPEG, MP3 og PNG

Nær 20 – 40%
aukaþjöppun

Oftast um 20%
aukaþjöppun

Huffman kóðun notuð til að
"kreista út" alla mögulega
umfremd (*redundancy*)

- Tveir gallar á venjulegri Huffman kóðun:
 - Þurfum tvær umferðir (til að safna tíðnigildum)
 - Þurfum að senda kóðunartréð með þjöppuðu gögnunum til að hægt sé að afkóða þau

Minskar þjöppunina, sérstaklega ef gagnamagnið er lítið

- Lausn: Kvik Huffman kóðun

Búum til kóðunartréð jafnóðum

Á hverjum tíma er tréð Huffman tré fyrir þá stafi sem hafa sést hingað til

Þegar við fáum nýjan staf þá notum við tréð til að kóða stafinn og uppfærum svo tréð út frá nýrri tíðni hans

Oftast ekki eins góð þjöppun og venjuleg Huffman

En ef tíðni stafa tákna er breytileg í stórri skrá þá getur Kvik Huffman kóðun verið betri

1. Ef við höfum 16 mismunandi tákn, sem öll hafa nákvæmlega sömu tíðni, hvaða kóðalengd gefur þá Huffman kóðun fyrir þau?
2. Búið til Huffman kóða fyrir táknin $a:20$, $b:10$, $c:5$, $r:8$ og sýnið kóðann fyrir strenginn "abraca"
3. Við höfum 1000 stafa streng. Hversu mikið er hægt að þjappa honum?