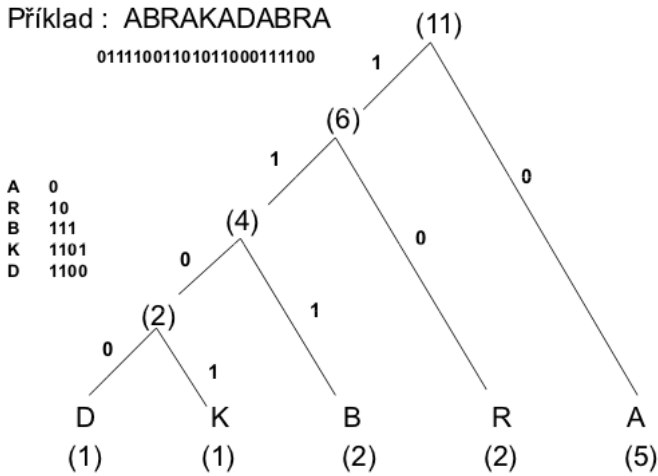




- Vybereme z S dva prvky M,N s nejnižšími četnostmi m, n,  $m < n$
- Vytvoříme uzel p, kde vlevo je M a vpravo je N, a jeho četnost bude  $m+n$
- Vložíme p do S a zpět na bod 3 dokud v S není jen jeden uzel

Strom samozřejmě musíme přenášet spolu s komprimovanými daty. Strom můžeme reprezentovat řetězcem tak, že jej projdeme DFS algoritmem a pro každý vrchol uložíme 0 a pro každý list 1 a znak tohoto listu. Při načítání čteme nuly a vytváříme uzly (směrem doleva). Pokud narazíme na 1, zapíšeme písmeno a vrátíme se do nejbližšího pravého uzlu (pravý uzel rodiče, případně nejvíce levý prvek podstromu - u rodiče vpravo)

Při dekompresi procházíme strom a rekonstruujeme zprávu.



## Aritmetické kódování

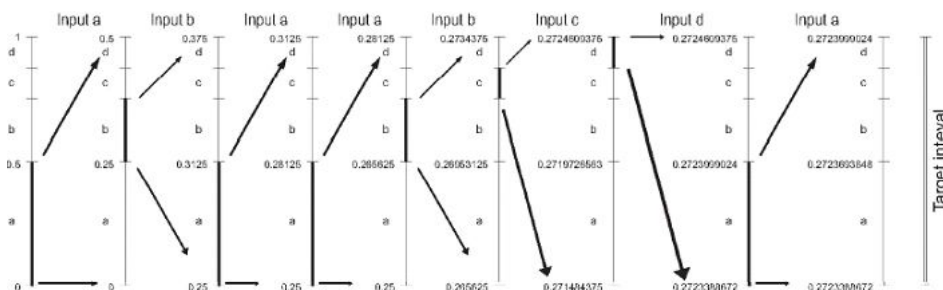
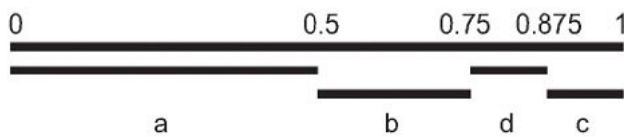
Také statistická metoda. Aritmetické kódování na rozdíl od jiných metod nepracuje na principu nahrazování vstupního znaku specifickým kódem. Místo toho se kódovaný vstupní proud znaků nahradí jediným reálným číslem z intervalu  $(0,1)$ .

Na základě pravděpodobnosti výskytu jednotlivých symbolů vstupní abecedy je každému symbolu přiřazena odpovídající poměrná část intervalu  $(0,1)$ . Při kódování je pak celý interval  $(0,1)$  postupně omezován z obou stran na základě postupně přicházejících symbolů (Interval ořízneme jen na tu část, které odpovídá pravděpodobnost písmene v textu, a tento interval opět rozdělíme podle pravděpodobností pro další písmeno).

Kódovaná hodnota se reprezentuje libovolným reálným číslem, které leží ve výsledném intervalu získaném po přičtení všech vstupních symbolů. Vzhledem k tomu, že z takto reprezentované hodnoty nelze při dekódování určit konec zprávy, je třeba navíc ke zprávě přidat speciální znak označující konec, případně musí být uložena i délka původní posloupnosti.

Příklad: kódujeme zprávu 'abaabca'

$$P(a)=0.5, P(b)=0.25, P(c)=0.125, P(d)=0.125$$



Co z toho číst?

- První řádka ukazuje pravděpodobnosti písmen (ty vypočteme jako četnost písmene / písmen celkem)
- Dále vidíme, jak bude interval (a podinterval) rozdělený. 'a' jako nejpravděpodobnější písmeno dostane největší část intervalu atd...
- A dále už kódujeme. Jdeme zleva do prava. Přejde písmeno 'a', vybereme tedy spodní interval. Tento opět rozdělíme na stejné části jako interval původní a přijde 'b'. Opět vybereme odpovídající interval, který rozdělíme atd.
- Zvýrazněná část intervalu se na obrázku vždy zaozobuje ;) (jsou tam šipky).
- Na konci máme jen jeden interval, vybereme tedy číslo, které do něj patří. Zpráva je zakódována v desetinné části tohoto čísla (a ano, celá zpráva je reprezentována jedním číslem)

## LZW

Z přednášek, doplněno: Slovníkový algoritmus. Využívá skupin znaků označených kódy a přenos kódů, které mají menší velikost než původní skupiny. Princip: Jednoprůchodová metoda (nevyžaduje předběžnou analýzu souboru. Ve výsledku se jen posunujeme po souboru jedním ukazatelem.) Vyhledávání opakujících se posloupností znaků, ukládání těchto posloupností do slovníku pro další použití a přiřazení jednoznakového kódu těmto posloupnostem. Kódy musí mít délku (počet bitů) větší než délka původních znaků (např. pro ASCII znaky (8 bitů) se obvykle používá délka kódu 12 bitů popř. větší. Na osmi bitech kódu by se dal uložit pouze slovník o 256 slovech (posloupnostech znaků), což odpovídá pouze počtu písmen, což je k ničemu.) Při průchodu komprimovaným souborem se vytváří slovník (počet položek slovníku odpovídá hodnotě 2(počet bitů kódu), kde prvních 2(počet bitů původních znaků) položek jsou znaky původní abecedy a zbývající položky tvoří posloupnosti znaků obsažené v komprimovaném souboru (používáme tedy i kódy z již komprimovaného souboru k vytváření nových kódů). Při dekompresi procházíme kódy, a podle tabulky je rozebíráme na menší a menší, až se dostaneme na elementární znaky, které vypíšeme.

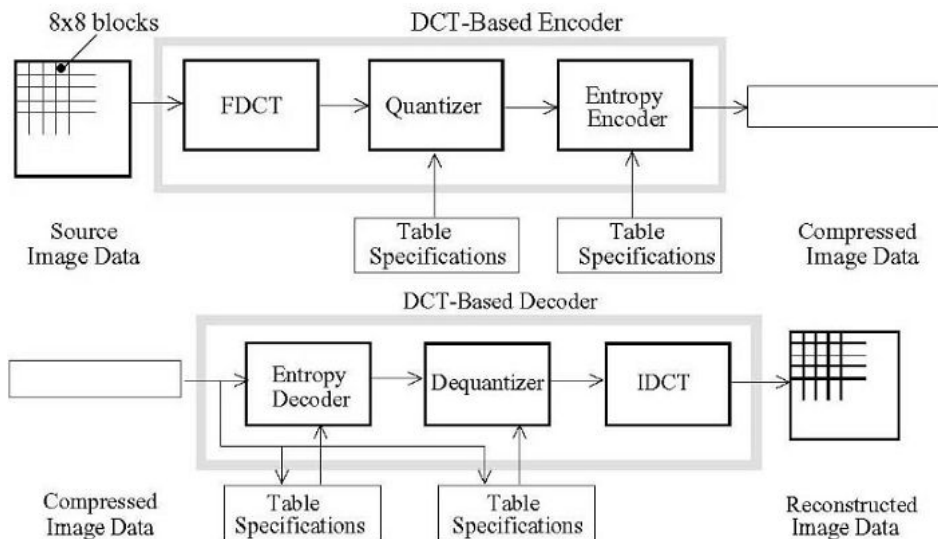
Aby to celé fungovalo, musíme slovník na dekomprimační straně "naučit" slova od menších po větší. Musíme proto na začátku souboru definovat elementární znaky, a z nich poté skládat slovník. Pro ASCII tedy prvních 256 znaků tabulky může odpovídat 1:1 ASCII abecedě, až poté následuje slovník. Pokud při dekompresi narazíme na neznámý kód, víme, že je složen z předchozích 2 kódů. Proto se jej naučíme, a příště už víme, co zapsat (případně jak jej rozložit). Tabulka na obou stranách tedy je stavěna stejně, jen používána z jiné strany (při komprimaci ukládáme kódy, při dekomprimaci jejich znaky)

Pro obrázky viz přednášky z PT, jsou celkem snadno pochopitelné, i když místo "kód" občas autor používá "znak" či "nový znak", což je trochu matoucí.

## JPG

V současné době patří mezi nejvíce používané komprese u obrázků. Je vhodná pro komprimaci fotek, nevhodná pro např. technické výkresy (čárové výkresy) - dochází k viditelnému rozmazání (na webové aplikace pomalu převládá png (zejména kvůli neztrátovosti a průhlednosti))

Princip: části obrazu se transformují do frekvenční oblasti (výsledkem je matice „frekvenčních“ koeficientů, bezztrátový převod) z matice koeficientů se odstraní koeficienty odpovídající vyšším frekvencím (rychlejší změny jasu - např. hrany v obraze, zde jsou ztáty) zbývající koeficienty se vhodným způsobem zkomprimují (huffman, aritmetický)



## Fraktálová komprese

Fraktál je obraz, který je matematicky popsateľný a tím jej lze libovolně přibližovat a oddalovat bez stráty informace. Také lze fraktál popsat jako soubor matematicky popsáných sebepodobných částí, jejichž spojováním dostáváme také sebepodobný obraz ve větším detailu (jako když děláme koláž z fotek, které dohromady a z dálky vypadají jako fotka)

Toho lze využít pro kompresi. Pokud se podíváme na libovolný obrázek, často zjistíme, že některé jeho části jsou si podobné

