

Background Paper Presentation

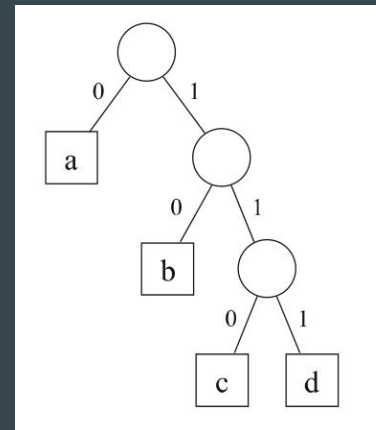
Blanchette, J.C.: *Proof pearl:
Mechanizing the textbook proof of
Huffman's algorithm.*
J. Autom. Reason. 2009

- Huffman Coding
- Huffman's Algorithm
- Proof of Optimality
- Conclusion
- Our Project

Daniel Filipe Nunes Silva & Samuel Chassot

Huffman Coding

- Code, $\mathcal{C} = \{ a \rightarrow 0, b \rightarrow 10, c \rightarrow 110, d \rightarrow 111 \}$
- Full binary tree
- Minimize encoded string length
- Compression



Huffman's Algorithm

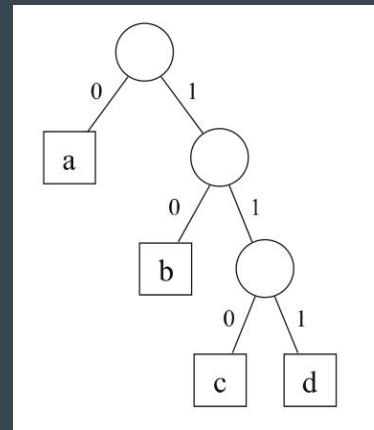
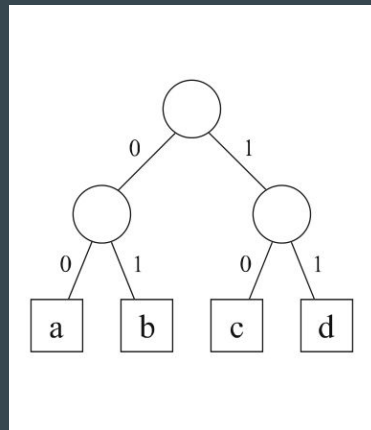
- Fixed-length vs variable-length
- Optimal solution
- Example for string “abacabad”

- **Input**

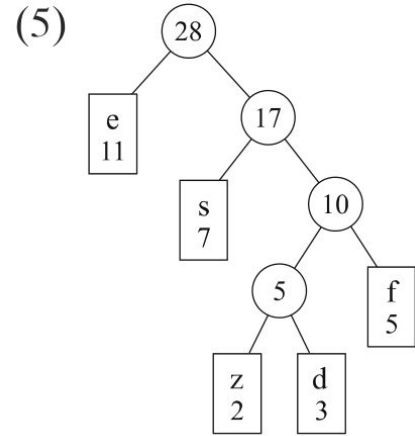
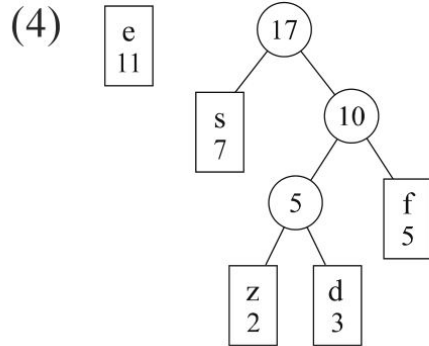
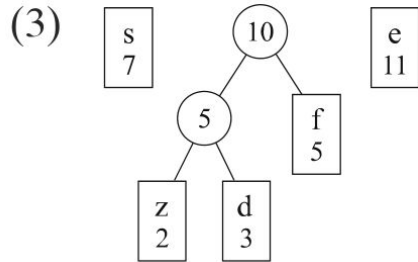
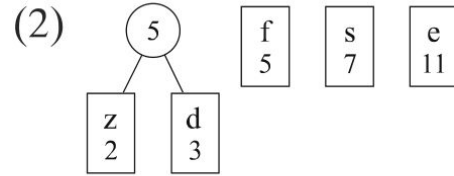
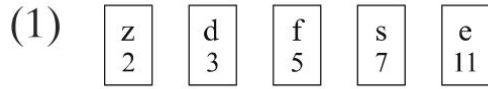
- Alphabet, $A = (a_1, a_2, \dots, a_n)$
- Weights, $W = (w_1, w_2, \dots, w_n)$, where w_i is the weight of a_i

- **Output** : Code, $C = (c_1, c_2, \dots, c_n)$, where c_i is the codeword for a_i

- Minimal sum of $w_i \cdot \text{length}(c_i)$

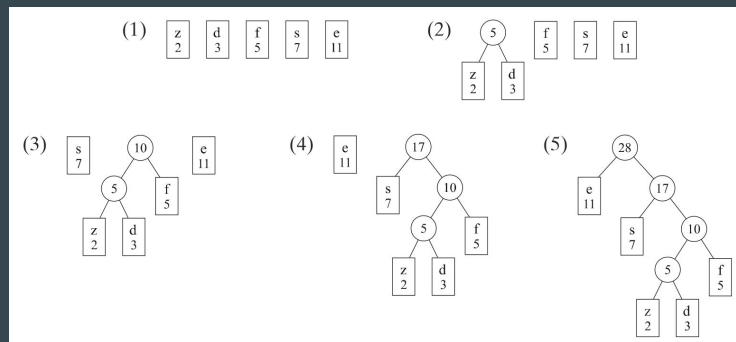


Huffman's Algorithm Example



Functional Implementation

- Datatypes
 - Leaf a w
 - InnerNode w t₁ t₂
- Functions
 - cachedWeight _
 - uniteTrees t₁ t₂
 - insertTree t f
- Huffman's algorithm
 - `huffman (t1 t2 ts) =`
`huffman(insertTree(uniteTrees t1 t2) ts)`
 - `huffman [t] = t`



Basic auxiliary functions

- functions defined on trees
- mostly recursive
- goal \rightarrow define *optimality* and lemmas to prove it

- *alphabet*:

$alphabet(Leaf\ w\ a) = \{a\}; \quad alphabet(InnerNode\ w\ t_1\ t_2) = alphabet\ t_1 \cup alphabet\ t_2.$

- *consistent*:

$consistent(Leaf\ w\ a) = True$

$consistent(InnerNode\ w\ t_1\ t_2) = (consistent\ t_1 \wedge consistent\ t_2$
 $\wedge alphabet\ t_1 \cap alphabet\ t_2 = \emptyset).$

- *depth*:

$$\text{depth}(\text{Leaf } w \ b) \ a = 0$$

$$\text{depth}(\text{InnerNode } w \ t_1 \ t_2) \ a = \begin{cases} \text{if } a \in \text{alphabet } t_1 \text{ then } \text{depth } t_1 \ a + 1 \\ \text{else if } a \in \text{alphabet } t_2 \text{ then } \text{depth } t_2 \ a + 1 \\ \text{else } 0. \end{cases}$$

- *height*:

$$\text{height}(\text{Leaf } w \ a) = 0$$

$$\text{height}(\text{InnerNode } w \ t_1 \ t_2) = \max(\text{height } t_1) (\text{height } t_2) + 1.$$

- *frequency*:

$$\begin{aligned} \text{freq} (\text{Leaf } w \ b) \ a &= (\text{if } a = b \text{ then } w \text{ else } 0) \\ \text{freq} (\text{InnerNode } w \ t_1 \ t_2) \ a &= \text{freq } t_1 \ a + \text{freq } t_2 \ a. \end{aligned}$$

- *weight*:

$$\begin{aligned} \text{weight} (\text{Leaf } w \ a) &= w \\ \text{weight} (\text{InnerNode } w \ t_1 \ t_2) &= \text{weight } t_1 + \text{weight } t_2. \end{aligned}$$

- *cost*:

$$\text{cost}(\text{Leaf } w \ a) = 0$$

$$\text{cost}(\text{InnerNode } w \ t_1 \ t_2) = \text{weight } t_1 + \text{cost } t_1 + \text{weight } t_2 + \text{cost } t_2.$$

- *optimum*:

$$\text{optimum } t = (\forall u. \text{consistent } u \longrightarrow \text{alphabet } t = \text{alphabet } u \longrightarrow \text{freq } t = \text{freq } u \longrightarrow \text{cost } t \leq \text{cost } u).$$

- *swapFourSyms*:

$$\begin{aligned} \text{swapFourSyms } t \ a \ b \ c \ d \ = & \ (\text{if } a = d \ \text{then } \text{swapSyms } t \ b \ c \\ & \ \text{else if } b = c \ \text{then } \text{swapSyms } t \ a \ d \\ & \ \text{else } \text{swapSyms } (\text{swapSyms } t \ a \ c) \ b \ d). \end{aligned}$$

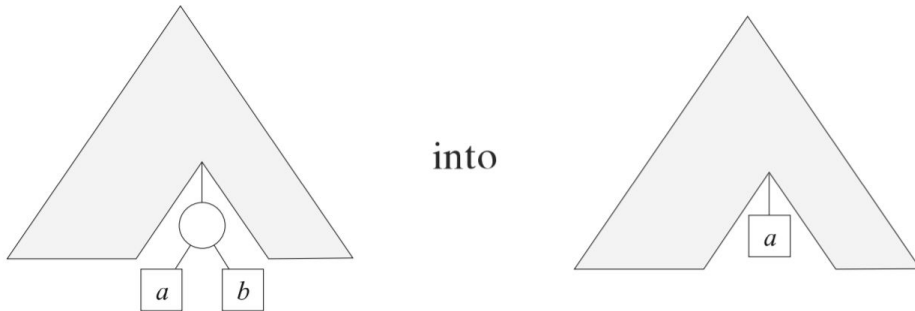
Exchange symbols so that a & b occupies the positions of c & d

$$\begin{aligned} \text{swapLeaves } (\text{Leaf } w_c \ c) \ w_a \ a \ w_b \ b \ = & \ (\text{if } c = a \ \text{then } \text{Leaf } w_b \ b \\ & \ \text{else if } c = b \ \text{then } \text{Leaf } w_a \ a \\ & \ \text{else } \text{Leaf } w_c \ c) \\ \text{swapLeaves } (\text{InnerNode } w \ t_1 \ t_2) \ w_a \ a \ w_b \ b \ = & \ \text{InnerNode } w \ (\text{swapLeaves } t_1 \ w_a \ a \ w_b \ b) \\ & \ (\text{swapLeaves } t_2 \ w_a \ a \ w_b \ b) \\ \text{swapSyms } t \ a \ b \ = & \ \text{swapLeaves } t \ (\text{freq } t \ a) \ a \ (\text{freq } t \ b) \ b \end{aligned}$$

- mergeSiblings:*

$$\text{mergeSibling} (\text{Leaf } w_b \ b) \ a = \text{Leaf } w_b \ b$$

$$\text{mergeSibling} (\text{InnerNode } w \ (\text{Leaf } w_b \ b) \ (\text{Leaf } w_c \ c)) \ a = \begin{cases} \text{Leaf } (w_b + w_c) \ a & \text{if } a = b \vee a = c \\ \text{InnerNode } w \ (\text{Leaf } w_b \ b) \ (\text{Leaf } w_c \ c) & \text{else} \end{cases}$$

$$\text{mergeSibling} (\text{InnerNode } w \ t_1 \ t_2) \ a = \text{InnerNode } w \ (\text{mergeSibling } t_1 \ a) \ (\text{mergeSibling } t_2 \ a).$$


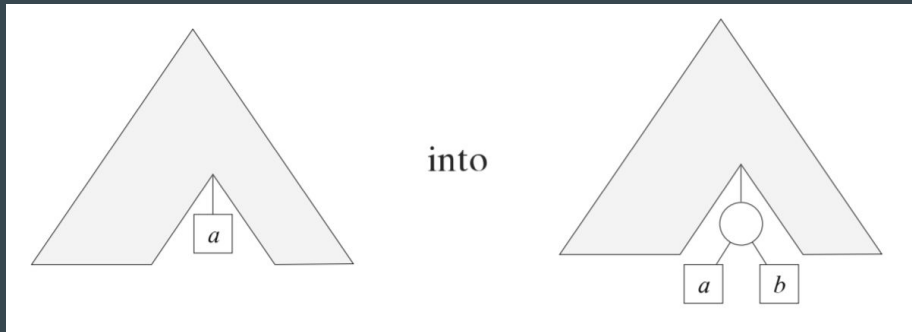
- *splitLeaf*:

$$\text{splitLeaf}(\text{Leaf } w_c \ c) \ w_a \ a \ w_b \ b = \begin{cases} \text{if } c = a \text{ then } \text{InnerNode } w_c \ (\text{Leaf } w_a \ a) \\ \hspace{15em} (\text{Leaf } w_b \ b) \\ \text{else } \text{Leaf } w_c \ c \end{cases}$$

else *Leaf* $w_c \ c$)

$$\text{splitLeaf}(\text{InnerNode } w \ t_1 \ t_2) \ w_a \ a \ w_b \ b = \text{InnerNode } w \ (\text{splitLeaf } t_1 \ w_a \ a \ w_b \ b) \\ (\text{splitLeaf } t_2 \ w_a \ a \ w_b \ b).$$

Normally, $w_a + w_b = \text{freq}(t, a)$

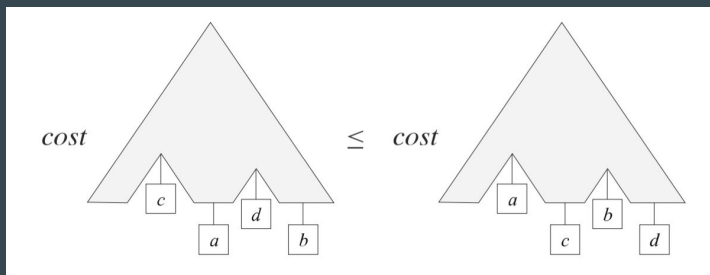


Lemma 8.5 - Leaf Split Optimality

If consistent t , optimum t , $a \in \text{alphabet } t$, $b \notin \text{alphabet } t$, $\text{freq } t a = w_a + w_b$, $\forall c \in \text{alphabet } t. w_a \leq \text{freq } t c \wedge w_b \leq \text{freq } t c$, and $w_a \leq w_b$, then optimum (splitLeaf $t w_a a w_b b$).

Intermediary lemmas - intuition

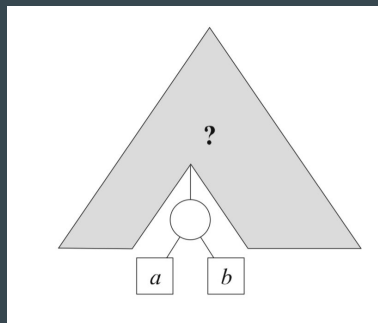
- **Lemma 8.4:** If a and b are the minima of the tree:



- **Lemma 7.2:** merging two siblings a and b decreases the *cost* by $w_a + w_b$

Lemma 8.5 - Leaf Split Optimality - Proof

- assume $\forall u, \text{cost}(t) \leq \text{cost}(u)$
- $\text{height}(t) > 0$: exists c and d , two sibling symbols at the bottom of u .
- swap c and d with a and b , minima of u .
 - We obtain a new tree v .

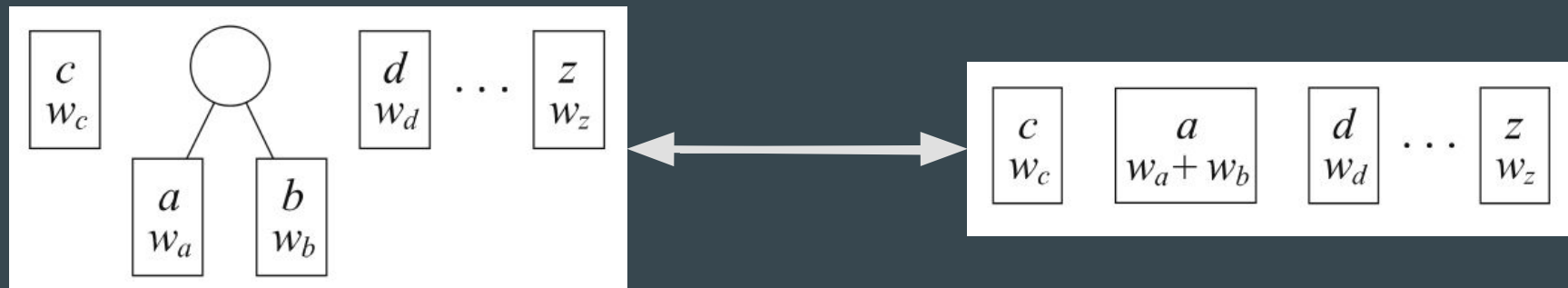


$$\begin{aligned}
 & \text{cost}(\text{splitLeaf } t \ a \ w_a \ b \ w_b) \\
 = & \text{cost } t + w_a + w_b && \text{by Lemma 7.3} \\
 \leq & \text{cost}(\overbrace{\text{mergeSibling}(\text{swapFourSyms } u \ a \ b \ c \ d) \ a}^v) + w_a + w_b && \text{by assumption} \\
 = & \text{cost}(\text{swapFourSyms } u \ a \ b \ c \ d) && \text{by Lemma 7.2} \\
 \leq & \text{cost } u. && \text{by Lemma 8.4}
 \end{aligned}$$

□

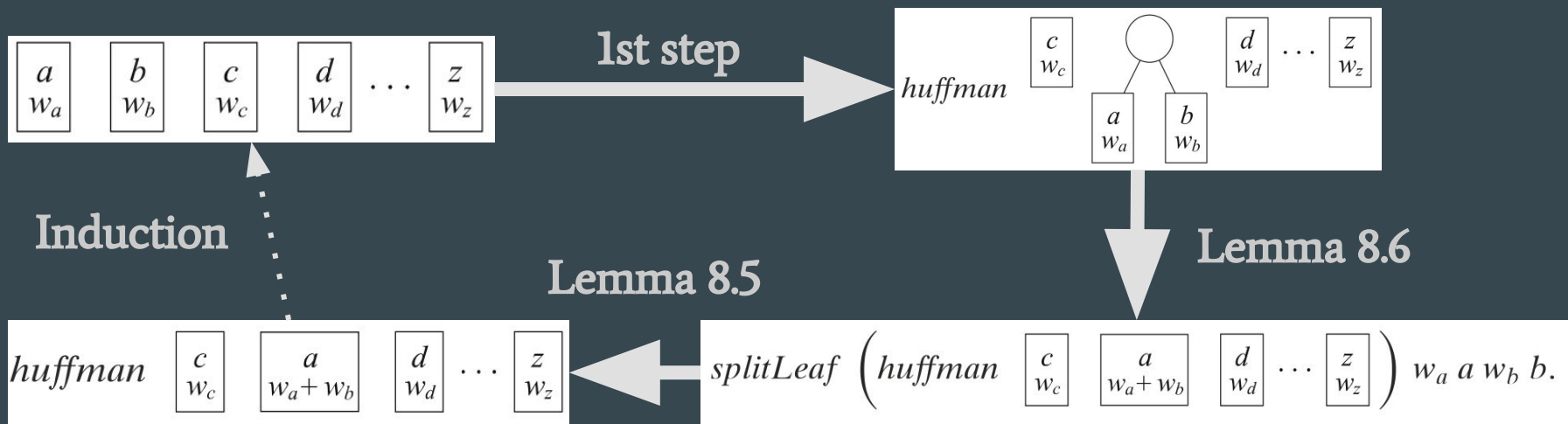
Leaf Split Commutativity

Lemma 8.6 (Leaf Split Commutativity) *If $\text{consistent}_F ts$, $ts \neq []$, and $a \in \text{alphabet}_F ts$, then $\text{splitLeaf}(\text{huffman } ts) w_a a w_b b = \text{huffman}(\text{splitLeaf}_F ts w_a a w_b b)$.*



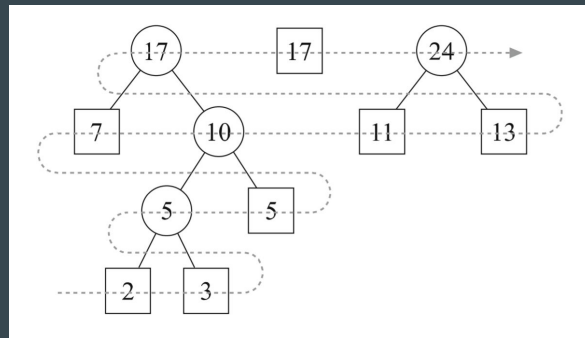
Huffman Optimality

Theorem 8.7 (Huffman Optimality) *If $\text{consistent}_F ts$, $\text{height}_F ts = 0$, $\text{sortedByWeight} ts$, and $ts \neq []$, then $\text{optimum}(\text{huffman } ts)$.*



Conclusion

- They formalized and demonstrated the textbook's proof of the Huffman algorithm
- They found that custom induction rules simplify a lot the proof using Isabelle
- They defined Lemmas that were not seen in the literature
 - e.g. each step of the algorithm preserves this invariant:
 - the nodes of the forest are ordered by weight from left to right, bottom to top:
 - They could not prove the preservation though



Our project

- They proposed to then extend the proof's scope to the applications of the algorithm.
- That is what we will do:

We will implement a pair of *encode/decode* functions using the tree produced by the Huffman's algorithm and prove that they are bijectives.
i.e. that

$$\mathit{decode}(\mathit{encode}(\mathbf{x})) = \mathbf{x}$$