HW 3 due 6/6

# CS3000: Algorithms & Data
# Drew van der Poel
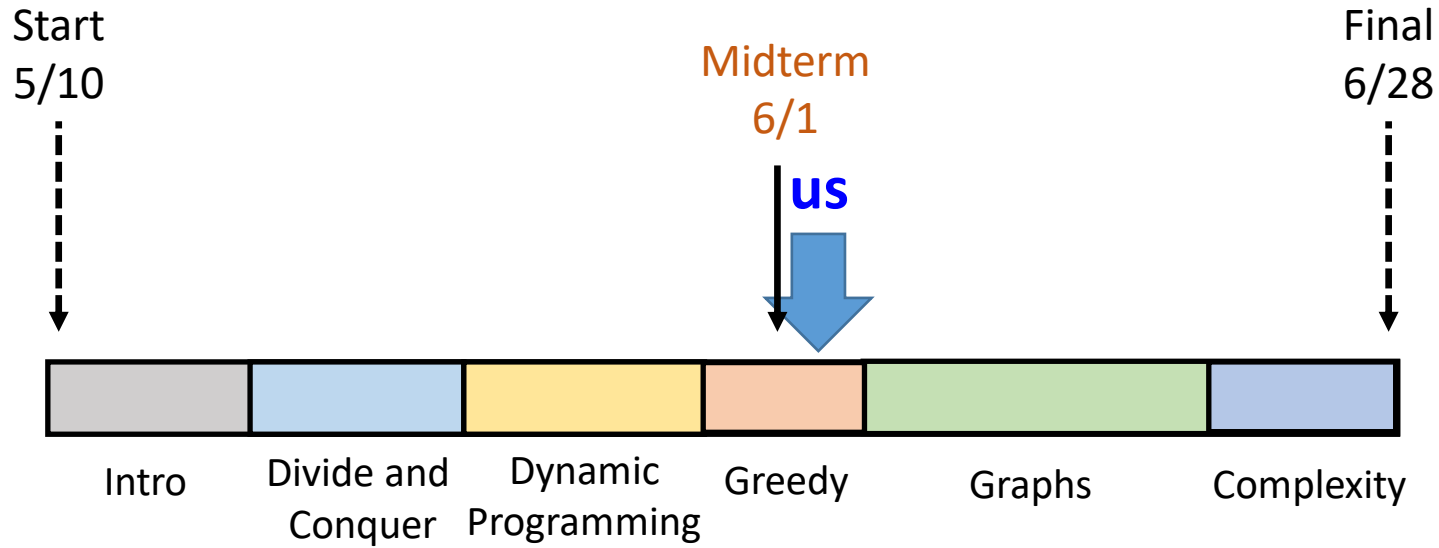
Lecture 13
- Finish Greedy: Exchange Argument
- Huffman Codes

June 2, 2021

# Outline

Start
5/10

Midterm
6/1

**us**

Final
6/28

| Intro | Divide and Conquer | Dynamic Programming | Greedy | Graphs | Complexity |

**Last class:** midterm review

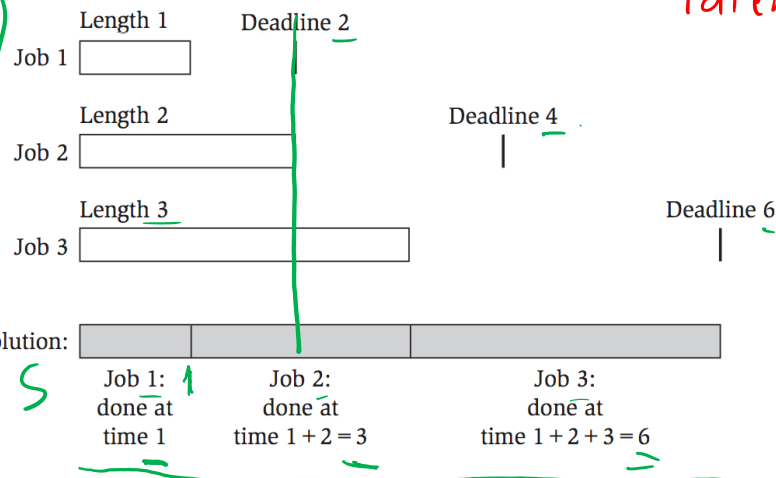**Next class:** graphs – terminology, BFS

# Minimum Lateness Scheduling

- Input: $n$ jobs with length $t_i$ and deadline $d_i$
  - Simplifying assumption: all deadlines are distinct
- Output: a minimum--lateness schedule for the jobs
  - Can only do one job at a time, no overlap
  - The lateness of job $i$ is $\max\{f_i - d_i, 0\}$
  - The lateness of a schedule is $\max_i\{\max\{f_i - d_i, 0\}\}$

$Lateness(1) = \max(1-2, 0)$
$= 0$
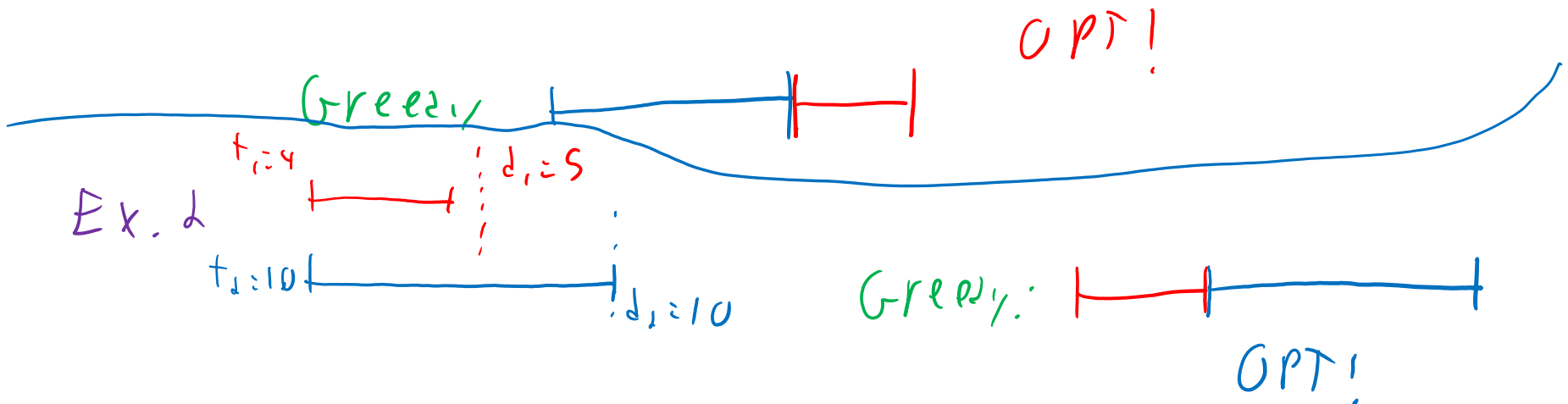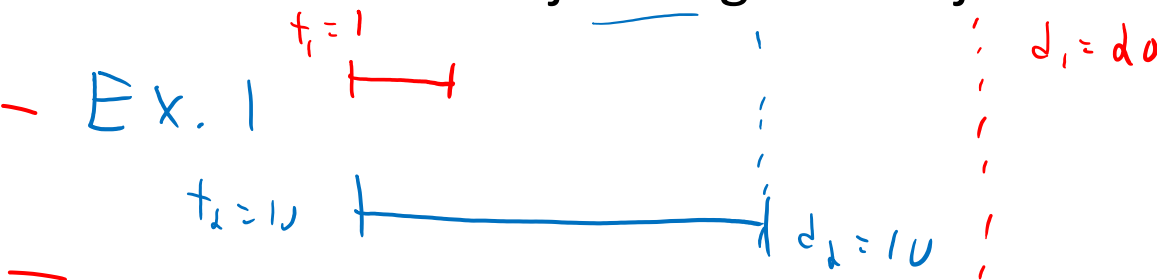
$Lateness(2) = \max(3-4, 0)$
$= 0$

$Lateness(3) = \max(6-6, 0)$
$= 0$

lateness of $i$

$lateness(S) = 0$

| | Length 1 | Deadline 2 |
|---|---|---|
| Job 1 | | |

Length 2    Deadline 4

Job 2

Length 3    Deadline 6

Job 3

Solution:

Job 1: done at time 1    Job 2: done at time $1+2=3$    Job 3: done at time $1+2+3=6$

# Greedy Algorithm: Earliest Deadline First

- Sort jobs so that $d_1 \leq d_2 \leq \cdots \leq d_n$
- For $i = 1, \ldots, n$:
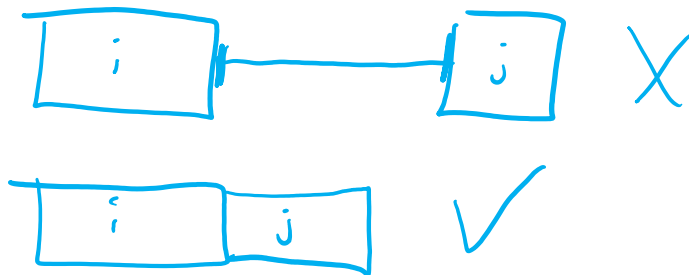  - Schedule job $i$ right after job $i-1$ finishes

# Exchange Argument

- $G$ = greedy schedule, $O$ = (supposedly) optimal schedule

- Exchange Argument:
  - We can transform $O$ to $G$ by exchanging pairs of jobs
  - Each exchange only reduces the lateness of $O$
  - Therefore the lateness of $G$ is at most that of $O$

$$O: \quad j_0, \quad j_1, \quad j_2, \quad \ldots, \quad j_n$$

swap

$$O \longrightarrow swap\ 1 \longrightarrow swap\ 2 \longrightarrow \ldots \longrightarrow G$$

# Exchange Argument

- $G$ = greedy schedule, $O$ = (supposedly) optimal schedule

- Observation: ~~the~~ any optimal schedule has no gaps
  - A schedule is just an ordering of the jobs, with jobs scheduled back–to–back

# Exchange Argument

*any other schedule*

- $G$ = greedy schedule, $O$ = (supposedly) optimal schedule

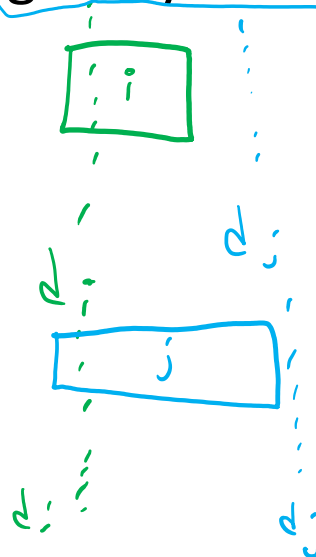- We say that two jobs $i$, $j$ are **inverted** in $O$ if $d_i < d_j$ but $j$ comes before $i$ in the schedule
  - Observation: greedy has no inversions

| j | | i |

Inversion

| i | | j |

$d_i$ $d_j$

No inversion

# Exchange Argument

- We say that two jobs $i, j$ are <span style="color:red">inverted</span> in $O$ if $d_i < d_j$ but $j$ comes before $i$

- <span style="color:#c60">Claim: an optimal schedule has no inversions</span>

  $(a+b)$

  - Step 1: suppose $O$ has an inversion, then it has an inversion $i, j$ where $i, j$ are <span style="color:red">consecutive</span>

  

# Exchange Argument

- We say that two jobs $i, j$ are <span style="color:red">inverted</span> in $O$ if $d_i < d_j$ but $j$ comes before $i$

- <span style="color:orange">Claim: an optimal schedule has no inversions</span>
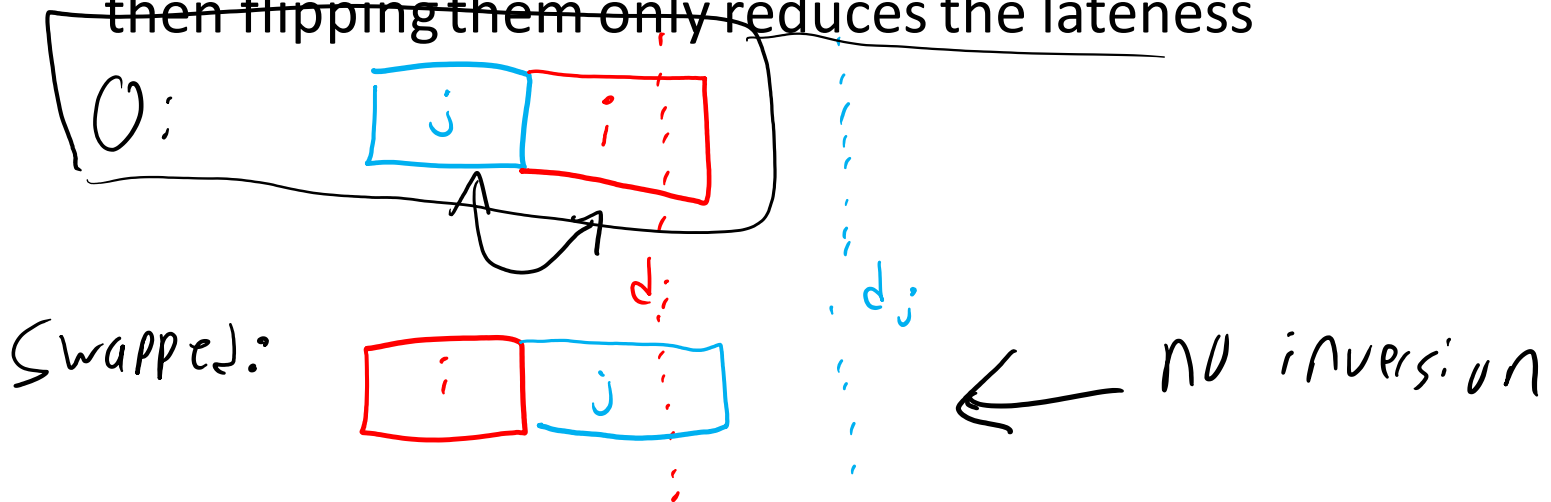  - Step 1: suppose $O$ has an inversion, then it has an inversion $i, j$ where $i, j$ are <span style="color:red">consecutive</span>
  - Step 2: if $i, j$ are a consecutive jobs that are inverted then flipping them only reduces the lateness

$O$:  [ j ][ i ]   $d_i$   $d_j$

Swapped:  [ i ][ j ]   ← No inversion

# Exchange Argument

- If $i, j$ are a consecutive jobs that are inverted then flipping them only reduces the lateness

Inversion

$d_i < d_j$ ✓

No Inversion

$j$     $i$

$i$     $j$

$i$     $j$

$s$        $f$      $d_j$

$s$        $f$

$d_i$

Case

$d_i > f = s + t_i + t_j$

max.
lateness($i$)
$0$

max.
lateness($j$)
$0$

$j$     $i$

max.
lateness($j$)
$0$

max.
lateness($i$)
$0$

max.
lateness $\downarrow$ $0$

$s + t_j - d_j$

$\boxed{s + t_i + t_j - d_i}$   $d_i < s + t_i + t_j$

① $s + t_i + t_j - d_j$

② $s + t_i - d_i$

$s + t_j - d_j < s + t_i + t_j - d_i$

$d_i < d_j + t_i$

① $s + t_i + t_j - d_j > s + t_i + t_j - d_i$

$d_j > d_i$ ✓

② $s + t_i + t_j - d_i > s + t_i - d_i$
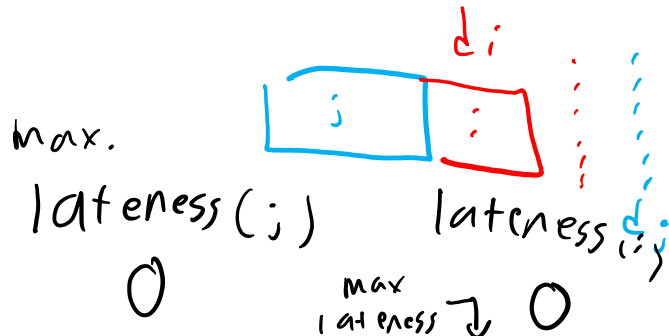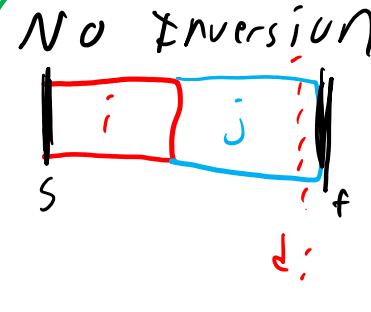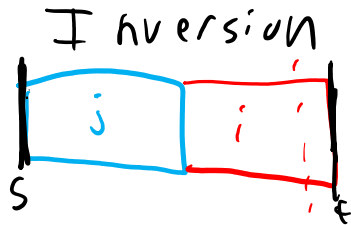
$t_j > 0$ ✓

# Exchange Argument

- We say that two jobs $i, j$ are inverted in $O$ if $d_i < d_j$ but $j$ comes before $i$

- Claim: an optimal schedule has no inversions
  - Step 1: suppose $O$ has an inversion, then it has an inversion $i, j$ where $i, j$ are consecutive
  - Step 2: if $i, j$ are consecutive jobs that are inverted then flipping them only reduces the lateness

    $$O \longrightarrow \text{Swap 1} \longrightarrow \dots \longrightarrow G$$

- $G$ is the unique schedule with no inversions, lateness(G) $\leq$ lateness(O)

- Problems: counting students, stable matching, sorting, n-digit multiplication, array searching, selection, weighted interval scheduling, segmented least squares, knapsack

- Alg. techniques: divide & conquer, dynamic programming, greedy

- Analysis: asymptotic analysis, recursion trees, Master Thm.

- Proof techniques: (strong) induction, contradiction, greedy stays ahead, **exchange argument**

# Data Compression

- How do we store strings of text compactly?

- A binary code is a mapping from alphabet $\Sigma \to \{0,1\}*$

  - Simplest code: assign numbers $0, 1, \ldots, |\Sigma|-1$ to each symbol, map to binary numbers of $\lceil \log_2 |\Sigma| \rceil$ bits (**fixed-length**)

  - Morse Code: (variable length)

    | | | | | | |
    |---|---|---|---|---|---|
    | A | ●━ | J | ●━━━ | S | ●●● |
    | B | ━●●● | K | ━●━ | T | ━ |
    | C | ━●━● | L | ●━●● | U | ●●━ |
    | D | ━●● | M | ━━ | V | ●●●━ |
    | E | ● | N | ━● | W | ●━━ |
    | F | ●●━● | O | ━━━ | X | ━●●━ |
    | G | ━━● | P | ●━━● | Y | ━●━━ |
    | H | ●●●● | Q | ━━●━ | Z | ━━●● |
    | I | ●● | R | ●━● | | |

$\Sigma = \{a, b, c, d\}$

$|\Sigma| = 4$      2-bits

$a \to 0$      $00$

$b \to 1$      $01$

$c \to 2$      $10$

$d \to 3$      $11$

# Data Compression

- Letters have uneven frequencies!
  - Want to use short encodings for frequent letters, long encodings for infrequent letters -> smaller files/more compression
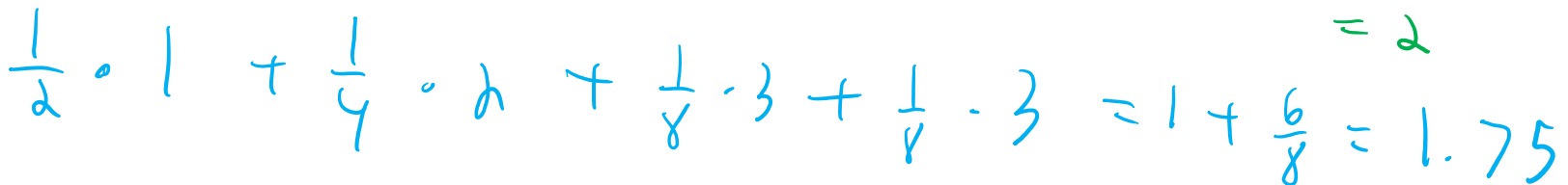
|  | a | b | c | d | exp. enc. len. |
|---|---|---|---|---|---|
| Frequency | 1/2 | 1/4 | 1/8 | 1/8 | |
| Encoding 1 | 00 | 01 | 10 | 11 | 2.0 |
| Encoding 2 | 0 | 10 | 110 | 111 | 1.75 |

$$\frac{1}{2} \cdot 2 + \frac{1}{4} \cdot 2 + \frac{1}{8} \cdot 2 + \frac{1}{8} \cdot 2 = 1 + \frac{1}{2} + \frac{1}{4} + \frac{1}{4}$$

$$= 2$$

$$\frac{1}{2} \cdot 1 + \frac{1}{4} \cdot 2 + \frac{1}{8} \cdot 3 + \frac{1}{8} \cdot 3 = 1 + \frac{6}{8} = 1.75$$

# Data Compression

- What properties would a good code of A-Z have?

  - The encoding is short on average
    
    $\leq 4$ bits per letter

  - Easy to encode a string
    
    Encode(ALGO) = ●— — ●—●● — — ●— — —

  - Easy to decode a string? *NO*
    
    *N* *G* *I*
    
    Decode(— ● — — ● ● ●) =
    
    *Y* *S*
    
    *K* *B*

```
A ●—          J ●— — —      S ● ● ●
B —● ● ●       K —●—         T —
C —●—●         L ●—● ●       U ● ●—
D —● ●         M — —         V ● ● ●—
E ●           N —●          W ●— —
F ● ●—●        O — — —       X —● ●—
G — —●         P ●— —●       Y —●— —
H ● ● ● ●       Q — —●—       Z — —● ●
I ● ●          R ●—●
```

# Prefix Free Codes

- Cannot decode if there are ambiguities
    - e.g. $enc(\text{"}E\text{"})$ is a prefix of $enc(\text{"}S\text{"})$ in Morse code

$\bullet$            $\bullet \; \circ \; \circ$

$\uparrow$

$E$

- Prefix-Free Code:
    - A binary $enc: \Sigma \to \{0,1\}^*$ such that
      for every $x \neq y \in \Sigma$, $enc(x)$ is not a prefix of $enc(y)$
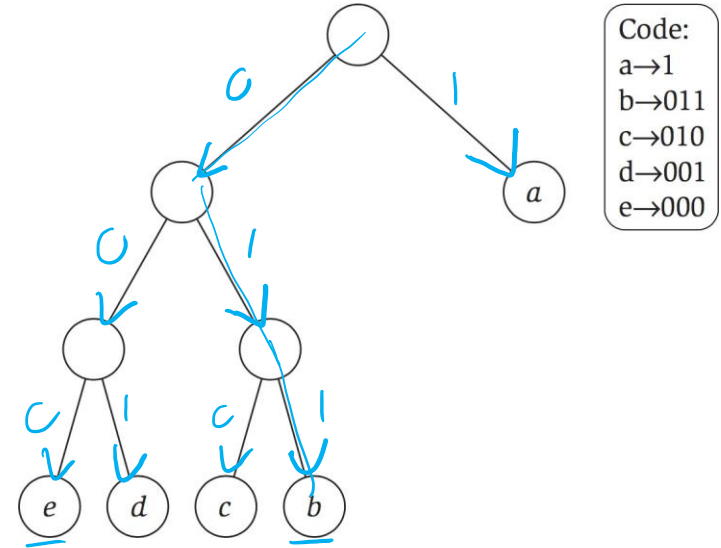
    - Any fixed-length code is prefix-free

    - Are all prefix-free codes fixed-length?  $No$

# Prefix Free Codes

- Can represent a prefix-free code as a binary tree
  - Each leaf is a character from the alphabet

Code:
a→1
b→011
c→010
d→001
e→000

- Encode by going up the tree (or using a table)
  - b e a d → 011 000 1 001

- Decode by going down the tree
  - 0 1 1 0 0 0 1 0 0 1 0 1 0 1 0 1 1
  
    b e a d c a b

# Huffman Codes

- (An algorithm to find) an **optimal** prefix-free code

$\downarrow \ len(T)$

| | **a** | **b** | **c** | **d** | **exp. enc. len.** |
|---|---|---|---|---|---|
| **Frequency** | 1/2 | 1/4 | 1/8 | 1/8 | |
| **Encoding 1** | 00 | 01 | 10 | 11 | 2.0 |
| **Encoding 2** | 0 | 10 | 110 | 111 | 1.75 |

$$len(E2) = f_a \cdot len_{E2}(a) + f_b \cdot len_{E2}(b) + f_c \cdot len_{E2}(c) + f_d \cdot len_{E2}(d)$$

$$= \frac{1}{2} \cdot 1 + \frac{1}{4} \cdot 2 + \frac{1}{8} \cdot 3 + \frac{1}{8} \cdot 3$$

# of bits
for i in T

- **optimal** $= \min\limits_{\text{prefix-free T}} \quad len(T) = \Sigma_{i \in \Sigma} \left( f_i * len_T(i) \right)$

  - Note, optimality depends on what you're compressing
  - H is the 8th most frequent letter in English (6.094%) but the 20th most frequent in Italian (0.636%)

- Problems: counting students, stable matching, sorting, n-digit mulitiplication, array searching, selection, weighted interval scheduling, segmented least squares, knapsack, **prefix-free encoding**

- Alg. techniques: divide & conquer, dynamic programming, greedy

- Analysis: asymptotic analysis, recursion trees, Master Thm.

- Proof techniques: (strong) induction, contradiction, greedy stays ahead, exchange argument
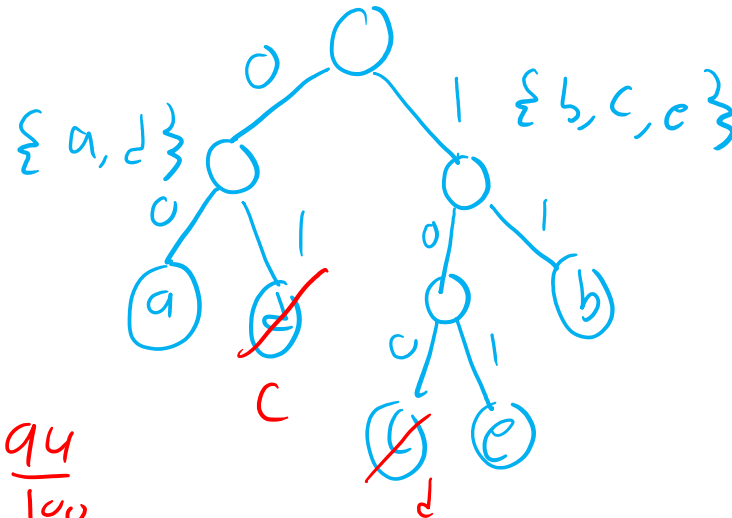
# Huffman Codes

$$\Sigma_{i \in \Sigma} \left( f_i * \underline{\text{len}}_T(i) \right)$$

- **Idea:** Balanced binary trees should have low depth -> small lengths
- **First Try:** Split letters into two sets of roughly equal frequency and repeat (greedy!!)

| a | b | c | d | e |
|---|---|---|---|---|
| .32 | .25 | .20 | .18 | .05 |

$f(i)$

after swap:

$\text{len}(c) = 2, \quad f_c = 0.2$

$\text{len}(d) = 3, \quad f_d = 0.18$

$2 \cdot \frac{20}{100} + 3 \cdot \frac{18}{100}$

$= \frac{40 + 54}{100} = \frac{94}{100}$

$\{a, d\}$ $\qquad$ $\{b, c, e\}$

0 $\quad$ 1

a $\quad$ d

C

d $\quad$ e

Originally: $\text{len}(d) = 2, \quad f_d = .18$

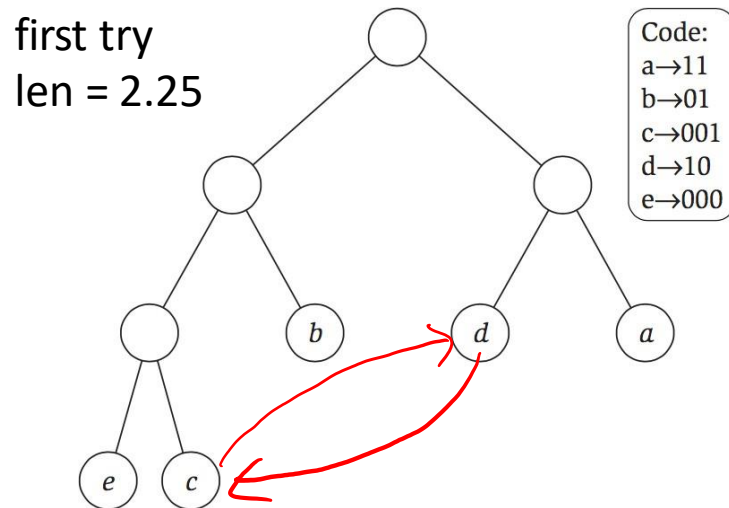$\text{len}(c) = 3, \quad f_c = .2$

$2 \cdot \frac{18}{100} + 3 \cdot \frac{20}{100} = \frac{36 + 60}{100}$

$= \frac{96}{100}$

# Huffman Codes

- First Try: Split letters into two sets of roughly equal frequency and repeat

| a | b | c | d | e |
|---|---|---|---|---|
| .32 | .25 | .20 | .18 | .05 |

first try
len = 2.25

Code:
a→11
b→01
c→001
d→10
e→000

optimal
len = 2.23

# Huffman Codes

| a | b | c | d | e |
|---|---|---|---|---|
| .32 | .25 | .20 | .18 | .05 |

optimal
len = 2.23

Code:
a→11
b→10
c→01
d→001
e→000

# Huffman Codes

- Huffman's Algorithm: pair up the two letters with the lowest frequency and repeat

| a | b | c | d | e |
|---|---|---|---|---|
| .32 | .25 | .20 | .18 | .05 |

{e,d} .23

.43

.57

a → 10        d → 000

b → 11        e → 001

c → 01

# Now You Try!

- Huffman's Algorithm: pair up the two letters with the lowest frequency and repeat

| a | b | c | d | e |
|---|---|---|---|---|
| .24 | .45 | .11 | .07 | .13 |

# Huffman Codes

- Huffman's Algorithm: pair up the two letters with the lowest frequency and repeat

- Theorem: Huffman's Algorithm produces a prefix-free code of optimal length
  - We can prove the theorem using an exchange argument

# An Experiment

- Take the Dickens novel *A Tale of Two Cities*
  - File size is 799,940 bytes

- Build a Huffman code and compress - what letters have long codes?

| char | frequency | code |
| --- | --- | --- |
| 'A' | 48165 | 1110 |
| 'B' | 8414 | 101000 |
| 'C' | 13896 | 00100 |
| 'D' | 28041 | 0011 |
| 'E' | 74809 | 011 |
| 'F' | 13559 | 111111 |
| 'G' | 12530 | 111110 |
| 'H' | 38961 | 1001 |

| char | frequency | code |
| --- | --- | --- |
| 'I' | 41005 | 1011 |
| 'J' | 710 | 1111011010 |
| 'K' | 4782 | 11110111 |
| 'L' | 22030 | 10101 |
| 'M' | 15298 | 01000 |
| 'N' | 42380 | 1100 |
| 'O' | 46499 | 1101 |
| 'P' | 9957 | 101001 |
| 'Q' | 667 | 1111011001 |

| char | frequency | code |
| --- | --- | --- |
| 'R' | 37187 | 0101 |
| 'S' | 37575 | 1000 |
| 'T' | 54024 | 000 |
| 'U' | 16726 | 01001 |
| 'V' | 5199 | 1111010 |
| 'W' | 14113 | 00101 |
| 'X' | 724 | 1111011011 |
| 'Y' | 12177 | 111100 |
| 'Z' | 215 | 1111011000 |

  - File size is now 439,688 bytes

| | Raw | Huffman |
| --- | --- | --- |
| **Size** | 799,940 | 439,688 |

# But Wait!

- Take the Dickens novel *A Tale of Two Cities*
  - File size is 799,940 bytes

- Build a Huffman code and compress

| char | frequency | code |
|------|-----------|------|
| 'A' | 48165 | 1110 |
| 'B' | 8414 | 101000 |
| 'C' | 13896 | 00100 |
| 'D' | 28041 | 0011 |
| 'E' | 74809 | 011 |
| 'F' | 13559 | 111111 |
| 'G' | 12530 | 111110 |
| 'H' | 38961 | 1001 |

| char | frequency | code |
|------|-----------|------|
| 'I' | 41005 | 1011 |
| 'J' | 710 | 1111011010 |
| 'K' | 4782 | 11110111 |
| 'L' | 22030 | 10101 |
| 'M' | 15298 | 01000 |
| 'N' | 42380 | 1100 |
| 'O' | 46499 | 1101 |
| 'P' | 9957 | 101001 |
| 'Q' | 667 | 1111011001 |

| char | frequency | code |
|------|-----------|------|
| 'R' | 37187 | 0101 |
| 'S' | 37575 | 1000 |
| 'T' | 54024 | 000 |
| 'U' | 16726 | 01001 |
| 'V' | 5199 | 1111010 |
| 'W' | 14113 | 00101 |
| 'X' | 724 | 1111011011 |
| 'Y' | 12177 | 111100 |
| 'Z' | 215 | 1111011000 |

  - File size is now 439,688 bytes

- But we can do better!

| | Raw | Huffman | gzip | bzip2 |
|------|-----|---------|------|-------|
| Size | 799,940 | 439,688 | 301,295 | 220,156 |