# CSE 431/531: Algorithm Analysis and Design (Spring 2020)
# Greedy Algorithms

Lecturer: Shi Li

*Department of Computer Science and Engineering*
*University at Buffalo*

**Def.** In an optimization problem, our goal of is to find a valid solution with the minimum cost (or maximum value).

**Def.** In an optimization problem, our goal of is to find a valid solution with the minimum cost (or maximum value).

## Trivial Algorithm for an Optimization Problem

Enumerate all valid solutions, compare them and output the best one.

**Def.** In an optimization problem, our goal of is to find a valid solution with the minimum cost (or maximum value).

### Trivial Algorithm for an Optimization Problem

Enumerate all valid solutions, compare them and output the best one.

- However, trivial algorithm often runs in exponential time, as the number of potential solutions is often exponentially large.

**Def.** In an optimization problem, our goal of is to find a valid solution with the minimum cost (or maximum value).

**Trivial Algorithm for an Optimization Problem**

Enumerate all valid solutions, compare them and output the best one.

- However, trivial algorithm often runs in exponential time, as the number of potential solutions is often exponentially large.

**Def.** In an optimization problem, our goal of is to find a valid solution with the minimum cost (or maximum value).

Trivial Algorithm for an Optimization Problem

Enumerate all valid solutions, compare them and output the best one.

- However, trivial algorithm often runs in exponential time, as the number of potential solutions is often exponentially large.
- $f(n)$ is a polynomial if $f(n) = O(n^k)$ for some constant $k > 0$.

**Def.** In an optimization problem, our goal of is to find a valid solution with the minimum cost (or maximum value).

## Trivial Algorithm for an Optimization Problem

Enumerate all valid solutions, compare them and output the best one.

- However, trivial algorithm often runs in exponential time, as the number of potential solutions is often exponentially large.
- $f(n)$ is a polynomial if $f(n) = O(n^k)$ for some constant $k > 0$.
- convention: polynomial time = efficient

**Def.** In an optimization problem, our goal of is to find a valid solution with the minimum cost (or maximum value).

## Trivial Algorithm for an Optimization Problem

Enumerate all valid solutions, compare them and output the best one.

- However, trivial algorithm often runs in exponential time, as the number of potential solutions is often exponentially large.
- $f(n)$ is a polynomial if $f(n) = O(n^k)$ for some constant $k > 0$.
- convention: polynomial time = efficient

## Goals of algorithm design

**Def.** In an optimization problem, our goal of is to find a valid solution with the minimum cost (or maximum value).

## Trivial Algorithm for an Optimization Problem

Enumerate all valid solutions, compare them and output the best one.

- However, trivial algorithm often runs in exponential time, as the number of potential solutions is often exponentially large.
- $f(n)$ is a polynomial if $f(n) = O(n^k)$ for some constant $k > 0$.
- convention: polynomial time = efficient

## Goals of algorithm design

1. Design efficient algorithms to solve problems

**Def.** In an optimization problem, our goal of is to find a valid solution with the minimum cost (or maximum value).

## Trivial Algorithm for an Optimization Problem

Enumerate all valid solutions, compare them and output the best one.

- However, trivial algorithm often runs in exponential time, as the number of potential solutions is often exponentially large.
- $f(n)$ is a polynomial if $f(n) = O(n^k)$ for some constant $k > 0$.
- convention: polynomial time = efficient

## Goals of algorithm design

1. Design efficient algorithms to solve problems
2. Design more efficient algorithms to solve problems

# Common Paradigms for Algorithm Design

- Greedy Algorithms
- Divide and Conquer
- Dynamic Programming

# Common Paradigms for Algorithm Design

- Greedy Algorithms
- Divide and Conquer
- Dynamic Programming

- Greedy algorithms are often for optimization problems.

# Common Paradigms for Algorithm Design

- Greedy Algorithms
- Divide and Conquer
- Dynamic Programming

- Greedy algorithms are often for optimization problems.
- They often run in polynomial time due to their simplicity.

## Greedy Algorithm

- Build up the solutions in steps
- At each step, make an irrevocable decision using a "reasonable" strategy

## Greedy Algorithm

- Build up the solutions in steps
- At each step, make an irrevocable decision using a "reasonable" strategy

## Analysis of Greedy Algorithm

- Prove that the reasonable strategy is "safe"
- Show that the remaining task after applying the strategy is to solve a (many) smaller instance(s) of the same problem

## Greedy Algorithm

- Build up the solutions in steps
- At each step, make an irrevocable decision using a "reasonable" strategy

## Analysis of Greedy Algorithm

- Prove that the reasonable strategy is "safe" (key)
- Show that the remaining task after applying the strategy is to solve a (many) smaller instance(s) of the same problem (usually easy)

## Greedy Algorithm

- Build up the solutions in steps
- At each step, make an irrevocable decision using a "reasonable" strategy

## Analysis of Greedy Algorithm

- Prove that the reasonable strategy is "safe" (key)
- Show that the remaining task after applying the strategy is to solve a (many) smaller instance(s) of the same problem (usually easy)

**Def.** A strategy is safe: there is always an optimum solution that agrees with the decision made according to the strategy.

# Outline

1. **Toy Example: Box Packing**

2. Interval Scheduling

3. Offline Caching

4. Data Compression and Huffman Code

5. Summary

## Box Packing

**Input:** $n$ boxes of capacities $c_1, c_2, \cdots, c_n$

$m$ items of sizes $s_1, s_2, \cdots, s_m$

Can put at most 1 item in a box

Item $j$ can be put into box $i$ if $s_j \leq c_i$

**Output:** A way to put as many items as possible in the boxes.

## Box Packing

**Input:** $n$ boxes of capacities $c_1, c_2, \cdots, c_n$

$m$ items of sizes $s_1, s_2, \cdots, s_m$

Can put at most 1 item in a box

Item $j$ can be put into box $i$ if $s_j \leq c_i$

**Output:** A way to put as many items as possible in the boxes.

## Example:

- Box capacities: $60, 40, 25, 15, 12$
- Item sizes:     $45, 42, 20, 19, 16$
- Can put $3$ items in boxes: $45 \rightarrow 60, 20 \rightarrow 40, 19 \rightarrow 25$

**Greedy Algorithm**

- Build up the solutions in steps
- At each step, make an <span style="color:red">irrevocable</span> decision using a "reasonable" strategy

## Greedy Algorithm

- Build up the solutions in steps
- At each step, make an irrevocable decision using a "reasonable" strategy

## Designing a Reasonable Strategy for Box Packing

- Q: Take box 1. Which item should we put in box 1?

## Greedy Algorithm

- Build up the solutions in steps
- At each step, make an irrevocable decision using a "reasonable" strategy

## Designing a Reasonable Strategy for Box Packing

- Q: Take box 1. Which item should we put in box 1?
- A: The item of the largest size that can be put into the box.

## Analysis of Greedy Algorithm

- Prove that the reasonable strategy is "safe"
- Show that the remaining task after applying the strategy is to solve a (many) smaller instance(s) of the same problem

## Analysis of Greedy Algorithm

- Prove that the reasonable strategy is "safe"
- Show that the remaining task after applying the strategy is to solve a (many) smaller instance(s) of the same problem

**Lemma** The strategy that put into box 1 the largest item it can hold is "safe": There is an optimum solution in which box 1 contains the largest item it can hold.

## Analysis of Greedy Algorithm

- Prove that the reasonable strategy is "safe"
- Show that the remaining task after applying the strategy is to solve a (many) smaller instance(s) of the same problem

**Lemma** The strategy that put into box 1 the largest item it can hold is "safe": There is an optimum solution in which box 1 contains the largest item it can hold.

- Intuition: putting the item gives us the easiest residual problem.

**Analysis of Greedy Algorithm**

- Prove that the reasonable strategy is "safe"
- Show that the remaining task after applying the strategy is to solve a (many) smaller instance(s) of the same problem

**Lemma** The strategy that put into box 1 the largest item it can hold is "safe": There is an optimum solution in which box 1 contains the largest item it can hold.

- Intuition: putting the item gives us the easiest residual problem.
- formal proof via exchanging argument:

**Lemma** There is an optimum solution in which box 1 contains the largest item it can hold.

**Lemma** There is an optimum solution in which box 1 contains the largest item it can hold.

**Proof.**

- Let $j =$ largest item that box 1 can hold.

**Lemma** There is an optimum solution in which box 1 contains the largest item it can hold.

Proof.
- Let $j =$ largest item that box 1 can hold.
- Take any optimum solution $S$. If $j$ is put into Box 1 in $S$, done.

**Lemma** There is an optimum solution in which box 1 contains the largest item it can hold.

Proof.

- Let $j =$ largest item that box 1 can hold.
- Take any optimum solution $S$. If $j$ is put into Box 1 in $S$, done.
- Otherwise, assume this is what happens in $S$:

**Lemma** There is an optimum solution in which box 1 contains the largest item it can hold.

- Let $j$ = largest item that box 1 can hold.
- Take any optimum solution $S$. If $j$ is put into Box 1 in $S$, done.
- Otherwise, assume this is what happens in $S$:



- $s_{j'} \leq s_j$, and swapping gives another solution $S'$

**Lemma** There is an optimum solution in which box 1 contains the largest item it can hold.

**Proof.**

- Let $j$ = largest item that box 1 can hold.
- Take any optimum solution $S$. If $j$ is put into Box 1 in $S$, done.
- Otherwise, assume this is what happens in $S$:



- $s_{j'} \leq s_j$, and swapping gives another solution $S'$
- $S'$ is also an optimum solution. In $S'$, $j$ is put into Box 1. $\square$

- Notice that the exchanging operation is only for the sake of analysis; it is not a part of the algorithm.

- Notice that the exchanging operation is only for the sake of analysis; it is not a part of the algorithm.

## Analysis of Greedy Algorithm

- Prove that the reasonable strategy is "safe"
- Show that the remaining task after applying the strategy is to solve a (many) smaller instance(s) of the same problem

- Notice that the exchanging operation is only for the sake of analysis; it is not a part of the algorithm.

Analysis of Greedy Algorithm

- Prove that the reasonable strategy is "safe"
- Show that the remaining task after applying the strategy is to solve a (many) smaller instance(s) of the same problem

- Trivial: we decided to put Item $j$ into Box 1, and the remaining instance is obtained by removing Item $j$ and Box 1.

10/68

## Generic Greedy Algorithm

1. **while** the instance is non-trivial **do**
2.     make the choice using the greedy strategy
3.     reduce the instance

## Greedy Algorithm for Box Packing

1. $T \leftarrow \{1, 2, 3, \cdots, m\}$
2. for $i \leftarrow 1$ to $n$ do
3.     if some item in $T$ can be put into box $i$, then
4.         $j \leftarrow$ the largest item in $T$ that can be put into box $i$
5.         print("put item $j$ in box $i$")
6.         $T \leftarrow T \setminus \{j\}$

## Generic Greedy Algorithm

1. **while** the instance is non-trivial **do**
2.      make the choice using the greedy strategy
3.      reduce the instance

**Lemma** Generic algorithm is correct if and only if the greedy strategy is safe.

## Generic Greedy Algorithm

1. **while** the instance is non-trivial **do**
2.      make the choice using the greedy strategy
3.      reduce the instance

**Lemma** Generic algorithm is correct if and only if the greedy strategy is safe.

- Greedy strategy is safe: we will not miss the optimum solution

## Generic Greedy Algorithm

1. **while** the instance is non-trivial **do**
2.     make the choice using the greedy strategy
3.     reduce the instance

**Lemma** Generic algorithm is correct if and only if the greedy strategy is safe.

- Greedy strategy is safe: we will not miss the optimum solution
- Greedy stretegy is not safe: we will miss the optimum solution for some instance, since the choices we made are irrevocable.

## Greedy Algorithm

- Build up the solutions in steps
- At each step, make an <span style="color:red">irrevocable</span> decision using a "reasonable" strategy

## Greedy Algorithm

- Build up the solutions in steps
- At each step, make an irrevocable decision using a "reasonable" strategy

## Analysis of Greedy Algorithm

- Prove that the reasonable strategy is "safe"
- Show that the remaining task after applying the strategy is to solve a (many) smaller instance(s) of the same problem

## Greedy Algorithm

- Build up the solutions in steps
- At each step, make an irrevocable decision using a "reasonable" strategy

## Analysis of Greedy Algorithm

- Prove that the reasonable strategy is "safe"
- Show that the remaining task after applying the strategy is to solve a (many) smaller instance(s) of the same problem

**Def.** A strategy is "safe" if there is always an optimum solution that is "consistent" with the decision made according to the strategy.

# Exchange argument: Proof of Safety of a Strategy

- let $S$ be an arbitrary optimum solution.
- if $S$ is consistent with the greedy choice, done.
- otherwise, show that it can be modified to another optimum solution $S'$ that is consistent with the choice.

- let $S$ be an arbitrary optimum solution.
- if $S$ is consistent with the greedy choice, done.
- otherwise, show that it can be modified to another optimum solution $S'$ that is consistent with the choice.

- The procedure is not a part of the algorithm.

# Outline

**Input:** $n$ jobs, job $i$ with start time $s_i$ and finish time $f_i$

$i$ and $j$ are compatible if $[s_i, f_i)$ and $[s_j, f_j)$ are disjoint

**Output:** A maximum-size subset of mutually compatible jobs

## Interval Scheduling

**Input:** $n$ jobs, job $i$ with start time $s_i$ and finish time $f_i$

$i$ and $j$ are compatible if $[s_i, f_i)$ and $[s_j, f_j)$ are disjoint

**Output:** A maximum-size subset of mutually compatible jobs

# Greedy Algorithm for Interval Scheduling

- Which of the following strategies are safe?

# Greedy Algorithm for Interval Scheduling

- Which of the following strategies are safe?
- Schedule the job with the smallest size?

# Greedy Algorithm for Interval Scheduling

- Which of the following strategies are safe?
- Schedule the job with the smallest size? No!

# Greedy Algorithm for Interval Scheduling

- Which of the following strategies are safe?
- Schedule the job with the smallest size? No!

# Greedy Algorithm for Interval Scheduling

- Which of the following strategies are safe?
- Schedule the job with the smallest size? No!

# Greedy Algorithm for Interval Scheduling

- Which of the following strategies are safe?
- Schedule the job with the smallest size? No!

# Greedy Algorithm for Interval Scheduling

- Which of the following strategies are safe?
- Schedule the job with the smallest size? No!

# Greedy Algorithm for Interval Scheduling

- Which of the following strategies are safe?
- Schedule the job with the smallest size? No!
- Schedule the job conflicting with smallest number of other jobs?

# Greedy Algorithm for Interval Scheduling

- Which of the following strategies are safe?
- Schedule the job with the smallest size? No!
- Schedule the job conflicting with smallest number of other jobs? No!

# Greedy Algorithm for Interval Scheduling

- Which of the following strategies are safe?
- Schedule the job with the smallest size? No!
- Schedule the job conflicting with smallest number of other jobs? No!

# Greedy Algorithm for Interval Scheduling

- Which of the following strategies are safe?
- Schedule the job with the smallest size? No!
- Schedule the job conflicting with smallest number of other jobs? No!

# Greedy Algorithm for Interval Scheduling

- Which of the following strategies are safe?
- Schedule the job with the smallest size? No!
- Schedule the job conflicting with smallest number of other jobs? No!

# Greedy Algorithm for Interval Scheduling

- Which of the following strategies are safe?
- Schedule the job with the smallest size? No!
- Schedule the job conflicting with smallest number of other jobs? No!

# Greedy Algorithm for Interval Scheduling

- Which of the following strategies are safe?
- Schedule the job with the smallest size? No!
- Schedule the job conflicting with smallest number of other jobs? No!

# Greedy Algorithm for Interval Scheduling

- Which of the following strategies are safe?
- Schedule the job with the smallest size? No!
- Schedule the job conflicting with smallest number of other jobs? No!
- Schedule the job with the earliest finish time?

# Greedy Algorithm for Interval Scheduling

- Which of the following strategies are safe?
- Schedule the job with the smallest size? No!
- Schedule the job conflicting with smallest number of other jobs? No!
- Schedule the job with the earliest finish time? Yes!

# Greedy Algorithm for Interval Scheduling

- Which of the following strategies are safe?
- Schedule the job with the smallest size? No!
- Schedule the job conflicting with smallest number of other jobs? No!
- Schedule the job with the earliest finish time? Yes!

# Greedy Algorithm for Interval Scheduling

- Which of the following strategies are safe?
- Schedule the job with the smallest size? No!
- Schedule the job conflicting with smallest number of other jobs? No!
- Schedule the job with the earliest finish time? Yes!

# Greedy Algorithm for Interval Scheduling

**Lemma** It is safe to schedule the job $j$ with the earliest finish time: There is an optimum solution where the job $j$ with the earliest finish time is scheduled.

Proof.

# Greedy Algorithm for Interval Scheduling

**Lemma** It is safe to schedule the job $j$ with the earliest finish time: There is an optimum solution where the job $j$ with the earliest finish time is scheduled.

Proof.

- Take an arbitrary optimum solution $S$

$S$:

# Greedy Algorithm for Interval Scheduling

**Lemma** It is safe to schedule the job $j$ with the earliest finish time: There is an optimum solution where the job $j$ with the earliest finish time is scheduled.

Proof.
- Take an arbitrary optimum solution $S$
- If it contains $j$, done

$S$:

# Greedy Algorithm for Interval Scheduling

**Lemma** It is safe to schedule the job $j$ with the earliest finish time: There is an optimum solution where the job $j$ with the earliest finish time is scheduled.

Proof.
- Take an arbitrary optimum solution $S$
- If it contains $j$, done

# Greedy Algorithm for Interval Scheduling

**Lemma** It is safe to schedule the job $j$ with the earliest finish time: There is an optimum solution where the job $j$ with the earliest finish time is scheduled.

Proof.

- Take an arbitrary optimum solution $S$
- If it contains $j$, done
- Otherwise, replace the first job in $S$ with $j$ to obtain another optimum schedule $S'$. □

$S$:

$j$:

# Greedy Algorithm for Interval Scheduling

**Lemma** It is safe to schedule the job $j$ with the earliest finish time: There is an optimum solution where the job $j$ with the earliest finish time is scheduled.

Proof.

- Take an arbitrary optimum solution $S$
- If it contains $j$, done
- Otherwise, replace the first job in $S$ with $j$ to obtain another optimum schedule $S'$. $\square$

# Greedy Algorithm for Interval Scheduling

**Lemma** It is safe to schedule the job $j$ with the earliest finish time: There is an optimum solution where the job $j$ with the earliest finish time is scheduled.

- What is the remaining task after we decided to schedule $j$?
- Is it another instance of interval scheduling problem?

# Greedy Algorithm for Interval Scheduling

**Lemma** It is safe to schedule the job $j$ with the earliest finish time: There is an optimum solution where the job $j$ with the earliest finish time is scheduled.

- What is the remaining task after we decided to schedule $j$?
- Is it another instance of interval scheduling problem? Yes!

# Greedy Algorithm for Interval Scheduling

**Lemma** It is safe to schedule the job $j$ with the earliest finish time: There is an optimum solution where the job $j$ with the earliest finish time is scheduled.

- What is the remaining task after we decided to schedule $j$?
- Is it another instance of interval scheduling problem? Yes!

# Greedy Algorithm for Interval Scheduling

**Lemma** It is safe to schedule the job $j$ with the earliest finish time: There is an optimum solution where the job $j$ with the earliest finish time is scheduled.

- What is the remaining task after we decided to schedule $j$?
- Is it another instance of interval scheduling problem? Yes!

# Greedy Algorithm for Interval Scheduling

Schedule($s, f, n$)

1. $A \leftarrow \{1, 2, \cdots, n\}, S \leftarrow \emptyset$
2. while $A \neq \emptyset$
3.     $j \leftarrow \arg\min_{j' \in A} f_{j'}$
4.     $S \leftarrow S \cup \{j\}; A \leftarrow \{j' \in A : s_{j'} \geq f_j\}$
5. return $S$

# Greedy Algorithm for Interval Scheduling

Schedule($s, f, n$)

1. $A \leftarrow \{1, 2, \cdots, n\}, S \leftarrow \emptyset$
2. while $A \neq \emptyset$
3. $\quad j \leftarrow \arg\min_{j' \in A} f_{j'}$
4. $\quad S \leftarrow S \cup \{j\}; A \leftarrow \{j' \in A : s_{j'} \geq f_j\}$
5. return $S$

# Greedy Algorithm for Interval Scheduling

Schedule($s, f, n$)

1. $A \leftarrow \{1, 2, \cdots, n\}, S \leftarrow \emptyset$
2. while $A \neq \emptyset$
3. $\quad j \leftarrow \arg\min_{j' \in A} f_{j'}$
4. $\quad S \leftarrow S \cup \{j\}; A \leftarrow \{j' \in A : s_{j'} \geq f_j\}$
5. return $S$

# Greedy Algorithm for Interval Scheduling

Schedule($s, f, n$)

1. $A \leftarrow \{1, 2, \cdots, n\}, S \leftarrow \emptyset$
2. while $A \neq \emptyset$
3. $\quad j \leftarrow \arg\min_{j' \in A} f_{j'}$
4. $\quad S \leftarrow S \cup \{j\}; A \leftarrow \{j' \in A : s_{j'} \geq f_j\}$
5. return $S$

# Greedy Algorithm for Interval Scheduling

Schedule($s, f, n$)

1. $A \leftarrow \{1, 2, \cdots, n\}, S \leftarrow \emptyset$
2. while $A \neq \emptyset$
3.     $j \leftarrow \arg\min_{j' \in A} f_{j'}$
4.     $S \leftarrow S \cup \{j\}; A \leftarrow \{j' \in A : s_{j'} \geq f_j\}$
5. return $S$

# Greedy Algorithm for Interval Scheduling

Schedule$(s, f, n)$

1. $A \leftarrow \{1, 2, \cdots, n\}, S \leftarrow \emptyset$
2. while $A \neq \emptyset$
3. $\quad j \leftarrow \arg\min_{j' \in A} f_{j'}$
4. $\quad S \leftarrow S \cup \{j\}; A \leftarrow \{j' \in A : s_{j'} \geq f_j\}$
5. return $S$

# Greedy Algorithm for Interval Scheduling

Schedule($s, f, n$)

1. $A \leftarrow \{1, 2, \cdots, n\}, S \leftarrow \emptyset$
2. while $A \neq \emptyset$
3.     $j \leftarrow \arg\min_{j' \in A} f_{j'}$
4.     $S \leftarrow S \cup \{j\}; A \leftarrow \{j' \in A : s_{j'} \geq f_j\}$
5. return $S$

Running time of algorithm?

# Greedy Algorithm for Interval Scheduling

Schedule($s, f, n$)

1. $A \leftarrow \{1, 2, \cdots, n\}, S \leftarrow \emptyset$
2. while $A \neq \emptyset$
3.     $j \leftarrow \arg\min_{j' \in A} f_{j'}$
4.     $S \leftarrow S \cup \{j\}; A \leftarrow \{j' \in A : s_{j'} \geq f_j\}$
5. return $S$

Running time of algorithm?

- Naive implementation: $O(n^2)$ time

# Greedy Algorithm for Interval Scheduling

**Schedule($s, f, n$)**

1. $A \leftarrow \{1, 2, \cdots, n\}, S \leftarrow \emptyset$
2. while $A \neq \emptyset$
3.     $j \leftarrow \arg\min_{j' \in A} f_{j'}$
4.     $S \leftarrow S \cup \{j\}; A \leftarrow \{j' \in A : s_{j'} \geq f_j\}$
5. return $S$

Running time of algorithm?

- Naive implementation: $O(n^2)$ time
- Clever implementation: $O(n \lg n)$ time

# Clever Implementation of Greedy Algorithm

**Schedule($s, f, n$)**

1. sort jobs according to $f$ values
2. $t \leftarrow 0$, $S \leftarrow \emptyset$
3. for every $j \in [n]$ according to non-decreasing order of $f_j$
4.     if $s_j \geq t$ then
5.         $S \leftarrow S \cup \{j\}$
6.         $t \leftarrow f_j$
7. return $S$

# Clever Implementation of Greedy Algorithm

Schedule($s, f, n$)

1. sort jobs according to $f$ values
2. $t \leftarrow 0$, $S \leftarrow \emptyset$
3. for every $j \in [n]$ according to non-decreasing order of $f_j$
4.     if $s_j \geq t$ then
5.         $S \leftarrow S \cup \{j\}$
6.         $t \leftarrow f_j$
7. return $S$

# Clever Implementation of Greedy Algorithm

Schedule($s, f, n$)

1. sort jobs according to $f$ values
2. $t \leftarrow 0$, $S \leftarrow \emptyset$
3. for every $j \in [n]$ according to non-decreasing order of $f_j$
4.    if $s_j \geq t$ then
5.        $S \leftarrow S \cup \{j\}$
6.        $t \leftarrow f_j$
7. return $S$

# Clever Implementation of Greedy Algorithm

Schedule($s, f, n$)

1. sort jobs according to $f$ values
2. $t \leftarrow 0$, $S \leftarrow \emptyset$
3. for every $j \in [n]$ according to non-decreasing order of $f_j$
4.     if $s_j \geq t$ then
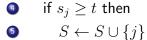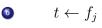5.         $S \leftarrow S \cup \{j\}$
6.         $t \leftarrow f_j$
7. return $S$

# Clever Implementation of Greedy Algorithm

## Schedule($s, f, n$)

1. sort jobs according to $f$ values
2. $t \leftarrow 0$, $S \leftarrow \emptyset$
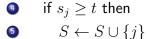3. for every $j \in [n]$ according to non-decreasing order of $f_j$
4.    if $s_j \geq t$ then
5.       $S \leftarrow S \cup \{j\}$
6.       $t \leftarrow f_j$
7. return $S$

# Clever Implementation of Greedy Algorithm

Schedule($s, f, n$)

① sort jobs according to $f$ values

② $t \leftarrow 0$, $S \leftarrow \emptyset$

③ for every $j \in [n]$ according to non-decreasing order of $f_j$

④     if $s_j \geq t$ then
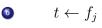
⑤        $S \leftarrow S \cup \{j\}$

⑥        $t \leftarrow f_j$

⑦ return $S$

# Clever Implementation of Greedy Algorithm

Schedule($s, f, n$)

1. sort jobs according to $f$ values
2. $t \leftarrow 0$, $S \leftarrow \emptyset$
3. for every $j \in [n]$ according to non-decreasing order of $f_j$
4.     if $s_j \geq t$ then
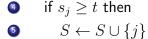5.         $S \leftarrow S \cup \{j\}$
6.         $t \leftarrow f_j$
7. return $S$

# Clever Implementation of Greedy Algorithm

**Schedule($s, f, n$)**

1. sort jobs according to $f$ values
2. $t \leftarrow 0$, $S \leftarrow \emptyset$
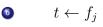3. for every $j \in [n]$ according to non-decreasing order of $f_j$
4.    if $s_j \geq t$ then
5.       $S \leftarrow S \cup \{j\}$
6.       $t \leftarrow f_j$
7. return $S$

# Clever Implementation of Greedy Algorithm

Schedule($s, f, n$)

① sort jobs according to $f$ values

② $t \leftarrow 0$, $S \leftarrow \emptyset$

③ for every $j \in [n]$ according to non-decreasing order of $f_j$

④      if $s_j \geq t$ then
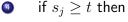
⑤          $S \leftarrow S \cup \{j\}$

⑥          $t \leftarrow f_j$

⑦ return $S$

# Clever Implementation of Greedy Algorithm

**Schedule($s, f, n$)**

1. sort jobs according to $f$ values
2. $t \leftarrow 0$, $S \leftarrow \emptyset$
3. for every $j \in [n]$ according to non-decreasing order of $f_j$
4.     if $s_j \geq t$ then
5.         $S \leftarrow S \cup \{j\}$
6.         $t \leftarrow f_j$
7. return $S$

# Clever Implementation of Greedy Algorithm

Schedule($s, f, n$)

1. sort jobs according to $f$ values
2. $t \leftarrow 0$, $S \leftarrow \emptyset$
3. for every $j \in [n]$ according to non-decreasing order of $f_j$
4.     if $s_j \geq t$ then
5.         $S \leftarrow S \cup \{j\}$
6.         $t \leftarrow f_j$
7. return $S$

# Clever Implementation of Greedy Algorithm

Schedule($s, f, n$)

1. sort jobs according to $f$ values
2. $t \leftarrow 0$, $S \leftarrow \emptyset$
3. for every $j \in [n]$ according to non-decreasing order of $f_j$
4.     if $s_j \geq t$ then
5.         $S \leftarrow S \cup \{j\}$
6.         $t \leftarrow f_j$
7. return $S$

# Clever Implementation of Greedy Algorithm

Schedule$(s, f, n)$

① sort jobs according to $f$ values
② $t \leftarrow 0$, $S \leftarrow \emptyset$
③ for every $j \in [n]$ according to non-decreasing order of $f_j$
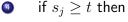④     if $s_j \geq t$ then
⑤       $S \leftarrow S \cup \{j\}$
⑥       $t \leftarrow f_j$
⑦ return $S$

# Clever Implementation of Greedy Algorithm

## Schedule($s, f, n$)

1. sort jobs according to $f$ values
2. $t \leftarrow 0$, $S \leftarrow \emptyset$
3. for every $j \in [n]$ according to non-decreasing order of $f_j$
4.     if $s_j \geq t$ then
5.         $S \leftarrow S \cup \{j\}$
6.         $t \leftarrow f_j$
7. return $S$

# Clever Implementation of Greedy Algorithm

Schedule($s, f, n$)

1. sort jobs according to $f$ values
2. $t \leftarrow 0$, $S \leftarrow \emptyset$
3. for every $j \in [n]$ according to non-decreasing order of $f_j$
4.      if $s_j \geq t$ then
5.          $S \leftarrow S \cup \{j\}$
6.          $t \leftarrow f_j$
7. return $S$

# Clever Implementation of Greedy Algorithm

Schedule($s, f, n$)

1. sort jobs according to $f$ values
2. $t \leftarrow 0$, $S \leftarrow \emptyset$
3. for every $j \in [n]$ according to non-decreasing order of $f_j$
4.    if $s_j \geq t$ then
5.       $S \leftarrow S \cup \{j\}$
6.       $t \leftarrow f_j$
7. return $S$

# Clever Implementation of Greedy Algorithm

Schedule($s, f, n$)

1. sort jobs according to $f$ values
2. $t \leftarrow 0$, $S \leftarrow \emptyset$
3. for every $j \in [n]$ according to non-decreasing order of $f_j$
4.     if $s_j \geq t$ then
5.        $S \leftarrow S \cup \{j\}$
6.        $t \leftarrow f_j$
7. return $S$

# Clever Implementation of Greedy Algorithm

## Schedule($s, f, n$)

1. sort jobs according to $f$ values
2. $t \leftarrow 0$, $S \leftarrow \emptyset$
3. for every $j \in [n]$ according to non-decreasing order of $f_j$
4.     if $s_j \geq t$ then
5.         $S \leftarrow S \cup \{j\}$
6.         $t \leftarrow f_j$
7. return $S$

# Outline

# Offline Caching

- Cache that can store $k$ pages
- Sequence of page requests

# Offline Caching

- Cache that can store $k$ pages
- Sequence of page requests

page sequence

| 1 |
| 5 |
| 4 |
| 2 |
| 5 |
| 3 |
| 2 |
| 1 |

cache

# Offline Caching

- Cache that can store $k$ pages
- Sequence of page requests
- Cache miss happens if requested page not in cache. We need bring the page into cache, and evict some existing page if necessary.

cache

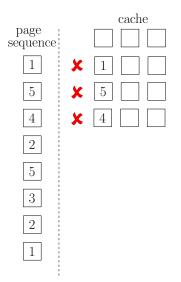page sequence

1

5

4

2

5

3

2

1

# Offline Caching

- Cache that can store $k$ pages
- Sequence of page requests
- Cache miss happens if requested page not in cache. We need bring the page into cache, and evict some existing page if necessary.

# Offline Caching

- Cache that can store $k$ pages
- Sequence of page requests
- Cache miss happens if requested page not in cache. We need bring the page into cache, and evict some existing page if necessary.

# Offline Caching

- Cache that can store $k$ pages
- Sequence of page requests
- Cache miss happens if requested page not in cache. We need bring the page into cache, and evict some existing page if necessary.
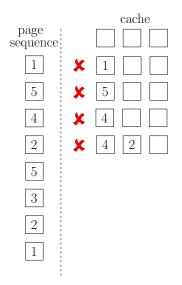
# Offline Caching

- Cache that can store $k$ pages
- Sequence of page requests
- Cache miss happens if requested page not in cache. We need bring the page into cache, and evict some existing page if necessary.

# Offline Caching

- Cache that can store $k$ pages
- Sequence of page requests
- Cache miss happens if requested page not in cache. We need bring the page into cache, and evict some existing page if necessary.

# Offline Caching

- Cache that can store $k$ pages
- Sequence of page requests
- Cache miss happens if requested page not in cache. We need bring the page into cache, and evict some existing page if necessary.
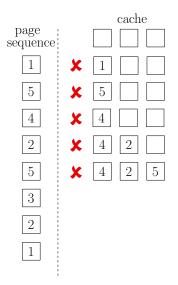
# Offline Caching

- Cache that can store $k$ pages
- Sequence of page requests
- Cache miss happens if requested page not in cache. We need bring the page into cache, and evict some existing page if necessary.

# Offline Caching

- Cache that can store $k$ pages
- Sequence of page requests
- Cache miss happens if requested page not in cache. We need bring the page into cache, and evict some existing page if necessary.
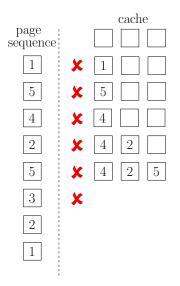
# Offline Caching

- Cache that can store $k$ pages
- Sequence of page requests
- Cache miss happens if requested page not in cache. We need bring the page into cache, and evict some existing page if necessary.

page sequence

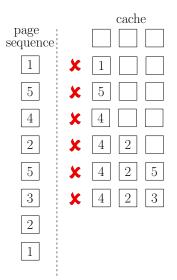| | | cache | | |
|---|---|---|---|---|
| | | | | |
| 1 | ✗ | 1 | | |
| 5 | ✗ | 5 | | |
| 4 | ✗ | 4 | | |
| 2 | ✗ | 4 | 2 | |
| 5 | ✗ | | | |
| 3 | | | | |
| 2 | | | | |
| 1 | | | | |

# Offline Caching

- Cache that can store $k$ pages
- Sequence of page requests
- Cache miss happens if requested page not in cache. We need bring the page into cache, and evict some existing page if necessary.
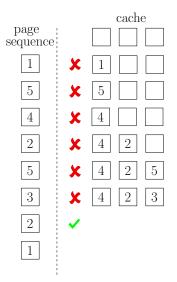
- Cache that can store $k$ pages
- Sequence of page requests
- Cache miss happens if requested page not in cache. We need bring the page into cache, and evict some existing page if necessary.

# Offline Caching

- Cache that can store $k$ pages
- Sequence of page requests
- Cache miss happens if requested page not in cache. We need bring the page into cache, and evict some existing page if necessary.
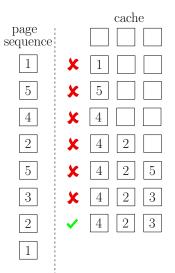
# Offline Caching

- Cache that can store $k$ pages
- Sequence of page requests
- Cache miss happens if requested page not in cache. We need bring the page into cache, and evict some existing page if necessary.
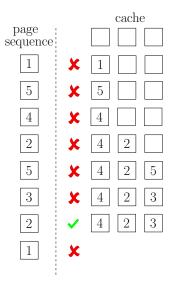- Cache hit happens if requested page already in cache.

# Offline Caching

- Cache that can store $k$ pages
- Sequence of page requests
- Cache miss happens if requested page not in cache. We need bring the page into cache, and evict some existing page if necessary.
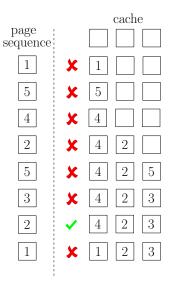- Cache hit happens if requested page already in cache.

# Offline Caching

- Cache that can store $k$ pages
- Sequence of page requests
- Cache miss happens if requested page not in cache. We need bring the page into cache, and evict some existing page if necessary.
- Cache hit happens if requested page already in cache.

# Offline Caching

- Cache that can store $k$ pages
- Sequence of page requests
- Cache miss happens if requested page not in cache. We need bring the page into cache, and evict some existing page if necessary.
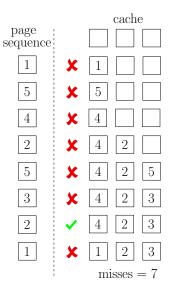- Cache hit happens if requested page already in cache.

# Offline Caching

- Cache that can store $k$ pages
- Sequence of page requests
- Cache miss happens if requested page not in cache. We need bring the page into cache, and evict some existing page if necessary.
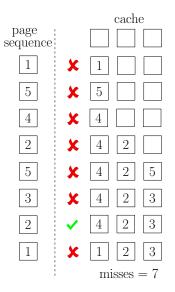- Cache hit happens if requested page already in cache.

# Offline Caching

- Cache that can store $k$ pages
- Sequence of page requests
- Cache miss happens if requested page not in cache. We need bring the page into cache, and evict some existing page if necessary.
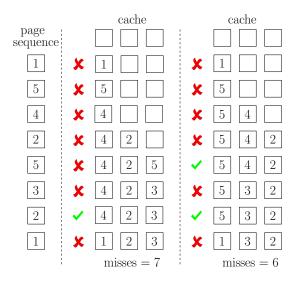- Cache hit happens if requested page already in cache.
- Goal: minimize the number of cache misses.

# A Better Solution for Example

## Offline Caching Problem

**Input:** $k$ : the size of cache

$n$ : number of pages

We use $[n]$ for $\{1, 2, 3, \cdots, n\}$.

$\rho_1, \rho_2, \rho_3, \cdots, \rho_T \in [n]$: sequence of requests

**Output:** $i_1, i_2, i_3, \cdots, i_T \in \{\text{hit}, \text{empty}\} \cup [n]$: indices of pages to evict ("hit" means evicting no page, "empty" means evicting empty page)

## Offline Caching Problem

**Input:** $k$ : the size of cache

$n$ : number of pages

We use $[n]$ for $\{1, 2, 3, \cdots, n\}$.

$\rho_1, \rho_2, \rho_3, \cdots, \rho_T \in [n]$: sequence of requests

**Output:** $i_1, i_2, i_3, \cdots, i_T \in \{\text{hit}, \text{empty}\} \cup [n]$: indices of pages to evict ("hit" means evicting no page, "empty" means evicting empty page)

- Offline Caching: we know the whole sequence ahead of time.
- Online Caching: we have to make decisions on the fly, before seeing future requests.

## Offline Caching Problem

**Input:** $k$ : the size of cache

$n$ : number of pages

We use $[n]$ for $\{1, 2, 3, \cdots, n\}$.

$\rho_1, \rho_2, \rho_3, \cdots, \rho_T \in [n]$: sequence of requests

**Output:** $i_1, i_2, i_3, \cdots, i_T \in \{\text{hit}, \text{empty}\} \cup [n]$: indices of pages to evict ("hit" means evicting no page, "empty" means evicting empty page)

- Offline Caching: we know the whole sequence ahead of time.
- Online Caching: we have to make decisions on the fly, before seeing future requests.

**Q:** Which one is more realistic?

## Offline Caching Problem

**Input:** $k$ : the size of cache

$n$ : number of pages

We use $[n]$ for $\{1, 2, 3, \cdots, n\}$.

$\rho_1, \rho_2, \rho_3, \cdots, \rho_T \in [n]$: sequence of requests

**Output:** $i_1, i_2, i_3, \cdots, i_T \in \{\text{hit}, \text{empty}\} \cup [n]$: indices of pages to evict ("hit" means evicting no page, "empty" means evicting empty page)

- Offline Caching: we know the whole sequence ahead of time.
- Online Caching: we have to make decisions on the fly, before seeing future requests.

**Q:** Which one is more realistic?
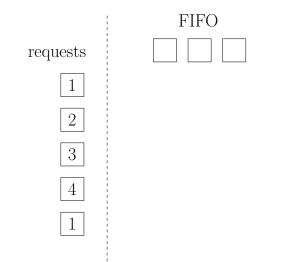
**A:** Online caching

- Offline Caching: we know the whole sequence ahead of time.
- Online Caching: we have to make decisions on the fly, before seeing future requests.

**Q:** Which one is more realistic?

**A:** Online caching

**Q:** Why do we study the offline caching problem?

- Offline Caching: we know the whole sequence ahead of time.
- Online Caching: we have to make decisions on the fly, before seeing future requests.

**Q:** Which one is more realistic?

**A:** Online caching

**Q:** Why do we study the offline caching problem?

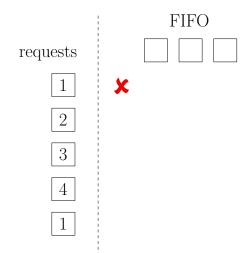**A:** Use the offline solution as a benchmark to measure the "competitive ratio" of online algorithms

# Offline Caching: Potential Greedy Algorithms

- FIFO(First-In-First-Out): always evict the first page in cache

# Offline Caching: Potential Greedy Algorithms

- FIFO(First-In-First-Out): always evict the first page in cache
- LRU(Least-Recently-Used): Evict page whose most recent access was earliest
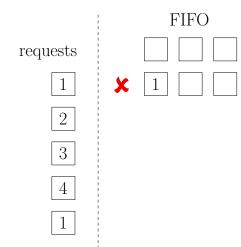
# Offline Caching: Potential Greedy Algorithms

- FIFO(First-In-First-Out): always evict the first page in cache
- LRU(Least-Recently-Used): Evict page whose most recent access was earliest
- LFU(Least-Frequently-Used): Evict page that was least frequently requested
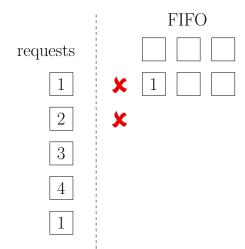
# Offline Caching: Potential Greedy Algorithms

- FIFO(First-In-First-Out): always evict the first page in cache
- LRU(Least-Recently-Used): Evict page whose most recent access was earliest
- LFU(Least-Frequently-Used): Evict page that was least frequently requested

- All the above algorithms are not optimum!
- Indeed all the algorithms are "online", i.e, the decisions can be made without knowing future requests. Online algorithms can not be optimum.
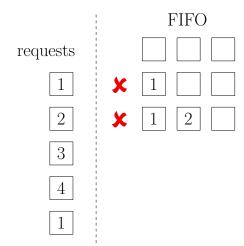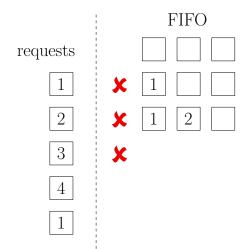
FIFO

requests

# FIFO is not optimum
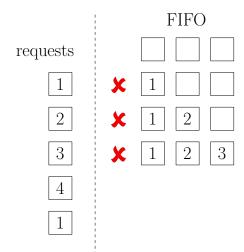


FIFO

requests

| 1 |

| 2 |

| 3 |

| 4 |

| 1 |

# FIFO is not optimum

# FIFO is not optimum

# FIFO is not optimum

# FIFO is not optimum

# FIFO is not optimum

# FIFO is not optimum

# FIFO is not optimum

# FIFO is not optimum

# FIFO is not optimum

# FIFO is not optimum

# Optimum Offline Caching

## Furthest-in-Future (FF)

- Algorithm: every time, evict the item that is not requested until furthest in the future, if we need to evict one.
- The algorithm is not an online algorithm, since the decision at a step depends on the request sequence in the future.

# Furthest-in-Future (FF)

# Example

requests

| 1 | 5 | 4 | 2 | 5 | 3 | 2 | 4 | 3 | 1 | 5 | 3 |

requests

| 1 | 5 | 4 | 2 | 5 | 3 | 2 | 4 | 3 | 1 | 5 | 3 |

✘ ✘ ✘

| | 1 | 1 | 1 |
| | | 5 | 5 |
| | | | 4 |

# Example

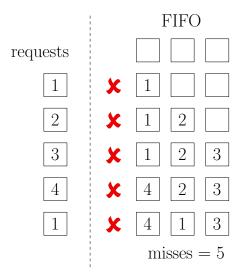requests

| 1 | 5 | 4 | 2 | (5) | 3 | 2 | (4) | 3 | (1) | 5 | 3 |

✘ ✘ ✘

| | 1 | 1 | 1 |
| | | 5 | 5 |
| | | | 4 |

# Example

# Example

requests

| 1 | 5 | 4 | 2 | 5 | 3 | 2 | 4 | 3 | 1 | 5 | 3 |

✗ ✗ ✗ ✗

| | 1 | 1 | 1 | 2 |

| | | 5 | 5 | 5 |

| | | | 4 | 4 |

# Example

requests

| 1 | 5 | 4 | 2 | 5 | 3 | 2 | 4 | 3 | 1 | 5 | 3 |

✘ ✘ ✘ ✘ ✔

| | 1 | 1 | 1 | 2 | 2 |
| | | 5 | 5 | 5 | 5 |
| | | | 4 | 4 | 4 |

# Example

# Example

# Example

# Example

# Example

# Example

requests

| 1 | 5 | 4 | 2 | 5 | 3 | 2 | 4 | 3 | 1 | 5 | ③ |
|---|---|---|---|---|---|---|---|---|---|---|---|

| ✘ | ✘ | ✘ | ✘ | ✔ | ✘ | ✔ | ✔ | ✔ |
|---|---|---|---|---|---|---|---|---|

|   | 1 | 1 | 1 | 2 | 2 | 2 | 2 | 2 |
|---|---|---|---|---|---|---|---|---|
|   |   | 5 | 5 | 5 | 5 | 3 | 3 | 3 | 3 |
|   |   |   | 4 | 4 | 4 | 4 | 4 | 4 | 4 |

# Example

# Example

requests

| 1 | 5 | 4 | 2 | 5 | 3 | 2 | 4 | 3 | 1 | 5 | (3) |
|---|---|---|---|---|---|---|---|---|---|---|-----|

| ✖ | ✖ | ✖ | ✖ | ✔ | ✖ | ✔ | ✔ | ✔ | ✖ | ✖ |
|---|---|---|---|---|---|---|---|---|---|---|

| | 1 | 1 | 1 | 2 | 2 | 2 | 2 | 2 | 2 | 1 | 5 |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | 5 | 5 | 5 | 5 | 3 | 3 | 3 | 3 | 3 | 3 |
| | | | 4 | 4 | 4 | 4 | 4 | 4 | 4 | 4 | 4 |

# Example

# Recall: Designing and Analyzing Greedy Algorithms

## Greedy Algorithm
- Build up the solutions in steps
- At each step, make an <span style="color:red">irrevocable</span> decision using a "reasonable" strategy

## Analysis of Greedy Algorithm
- Prove that the reasonable strategy is "safe" (key)
- Show that the remaining task after applying the strategy is to solve a (many) smaller instance(s) of the same problem (usually easy)

# Recall: Designing and Analyzing Greedy Algorithms

## Greedy Algorithm
- Build up the solutions in steps
- At each step, make an irrevocable decision using a "reasonable" strategy

## Analysis of Greedy Algorithm
- Prove that the reasonable strategy is "safe" (key)
- Show that the remaining task after applying the strategy is to solve a (many) smaller instance(s) of the same problem (usually easy)

## Offline Caching Problem

**Input:** $k$ : the size of cache
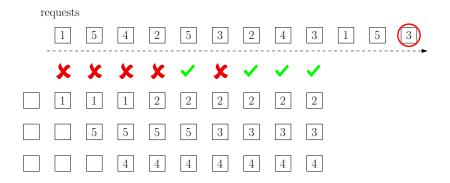
$n$ : number of pages

$\rho_1, \rho_2, \rho_3, \cdots, \rho_T \in [n]$: sequence of requests

**Output:** $i_1, i_2, i_3, \cdots, i_t \in \{\text{hit}, \text{empty}\} \cup [n]$

- empty stands for an empty page
- "hit" means evicting no pages

## Offline Caching Problem

**Input:** $k$ : the size of cache

$n$ : number of pages

$\rho_1, \rho_2, \rho_3, \cdots, \rho_T \in [n]$: sequence of requests

$p_1, p_2, \cdots, p_k \in \{\text{empty}\} \cup [n]$: initial set of pages in cache

**Output:** $i_1, i_2, i_3, \cdots, i_t \in \{\text{hit}, \text{empty}\} \cup [n]$

- empty stands for an empty page
- "hit" means evicting no pages

## Analysis of Greedy Algorithm

- Prove that the reasonable strategy is "safe" (key)
- Show that the remaining task after applying the strategy is to solve a (many) smaller instance(s) of the same problem (usually easy)

## Analysis of Greedy Algorithm

- Prove that the reasonable strategy is "safe" (key)
- Show that the remaining task after applying the strategy is to solve a (many) smaller instance(s) of the same problem (usually easy)

**Lemma** Assume at time 1 a page fault happens and there are no empty pages in the cache. Let $p^*$ be the page in cache that is not requested until furthest in the future. It is safe to evict $p^*$ at time 1.

## Analysis of Greedy Algorithm

- Prove that the reasonable strategy is "safe" (key)
- Show that the remaining task after applying the strategy is to solve a (many) smaller instance(s) of the same problem (usually easy)

**Lemma** Assume at time 1 a page fault happens and there are no empty pages in the cache. Let $p^*$ be the page in cache that is not requested until furthest in the future. There is an optimum solution in which $p^*$ is evicted at time 1.

| 4 | ? | ? | ? | 1 | 2 | ? | 3 |

$S:$

| 1 |
|---|
| 2 |
| 3 |

**Proof.**

1. $S$: any optimum solution
2. $p^*$: page in cache not requested until furthest in the future.
   - In the example, $p^* = 3$.

$S:$

| 4 | ? | ? | ? | 1 | 2 | ? | 3 |

| 1 | 1 |
| 2 | 4 |
| 3 | 3 |

**Proof.**

1. $S$: any optimum solution
2. $p^*$: page in cache not requested until furthest in the future.
   - In the example, $p^* = 3$.
3. Assume $S$ evicts some $p' \neq p^*$ at time 1; otherwise done.
   - In the example, $p' = 2$.

## Proof.

1. $S$: any optimum solution
2. $p^*$: page in cache not requested until furthest in the future.
   - In the example, $p^* = 3$.
3. Assume $S$ evicts some $p' \neq p^*$ at time 1; otherwise done.
   - In the example, $p' = 2$.

## Proof.

4. Create $S'$. $S'$ evicts $p^*(=3)$ instead of $p'(=2)$ at time 1.

---

**Proof.**

4. Create $S'$. $S'$ evicts $p^*(=3)$ instead of $p'(=2)$ at time 1.

## Proof.

4. Create $S'$. $S'$ evicts $p^*(=3)$ instead of $p'(=2)$ at time 1.
5. After time 1, cache status of $S$ and that of $S'$ differ by only 1 page. $S$ contains $p'(=2)$ and $S$ contains $p^*(=3)$.

## Proof.

4. Create $S'$. $S'$ evicts $p^*(=3)$ instead of $p'(=2)$ at time 1.

5. After time 1, cache status of $S$ and that of $S'$ differ by only 1 page. $S$ contains $p'(=2)$ and $S$ contains $p^*(=3)$.

6. From now on, $S'$ will "copy" $S$.

S :

S' :

### Proof.

4. Create $S'$. $S'$ evicts $p^*(=3)$ instead of $p'(=2)$ at time 1.

5. After time 1, cache status of $S$ and that of $S'$ differ by only 1 page. $S$ contains $p'(=2)$ and $S$ contains $p^*(=3)$.

6. From now on, $S'$ will "copy" $S$.

### Proof.

4. Create $S'$. $S'$ evicts $p^*(=3)$ instead of $p'(=2)$ at time 1.
5. After time 1, cache status of $S$ and that of $S'$ differ by only 1 page. $S$ contains $p'(=2)$ and $S$ contains $p^*(=3)$.
6. From now on, $S'$ will "copy" $S$.

**Proof.**

4. Create $S'$. $S'$ evicts $p^*(=3)$ instead of $p'(=2)$ at time 1.

5. After time 1, cache status of $S$ and that of $S'$ differ by only 1 page. $S$ contains $p'(=2)$ and $S$ contains $p^*(=3)$.

6. From now on, $S'$ will "copy" $S$.

**Proof.**

④ Create $S'$. $S'$ evicts $p^*(=3)$ instead of $p'(=2)$ at time 1.

⑤ After time 1, cache status of $S$ and that of $S'$ differ by only 1 page. $S$ contains $p'(=2)$ and $S$ contains $p^*(=3)$.

⑥ From now on, $S'$ will "copy" $S$.

**Proof.**

4. Create $S'$. $S'$ evicts $p^*(=3)$ instead of $p'(=2)$ at time 1.

5. After time 1, cache status of $S$ and that of $S'$ differ by only 1 page. $S$ contains $p'(=2)$ and $S$ contains $p^*(=3)$.

6. From now on, $S'$ will "copy" $S$.
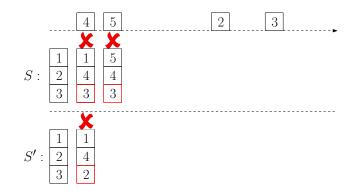
## Proof.

4. Create $S'$. $S'$ evicts $p^*(=3)$ instead of $p'(=2)$ at time 1.

5. After time 1, cache status of $S$ and that of $S'$ differ by only 1 page. $S$ contains $p'(=2)$ and $S$ contains $p^*(=3)$.

6. From now on, $S'$ will "copy" $S$.
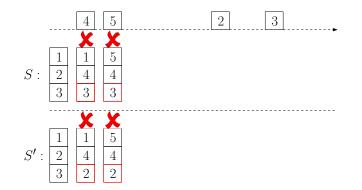
## Proof.

④ Create $S'$. $S'$ evicts $p^*(=3)$ instead of $p'(=2)$ at time 1.

⑤ After time 1, cache status of $S$ and that of $S'$ differ by only 1 page. $S$ contains $p'(=2)$ and $S$ contains $p^*(=3)$.

⑥ From now on, $S'$ will "copy" $S$.

## Proof.

4. Create $S'$. $S'$ evicts $p^*(=3)$ instead of $p'(=2)$ at time 1.

5. After time 1, cache status of $S$ and that of $S'$ differ by only 1 page. $S$ contains $p'(=2)$ and $S$ contains $p^*(=3)$.

6. From now on, $S'$ will "copy" $S$.
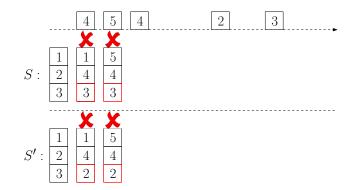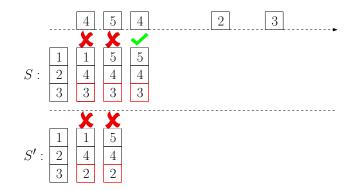
Proof.

4. Create $S'$. $S'$ evicts $p^*(=3)$ instead of $p'(=2)$ at time 1.

5. After time 1, cache status of $S$ and that of $S'$ differ by only 1 page. $S$ contains $p'(=2)$ and $S$ contains $p^*(=3)$.
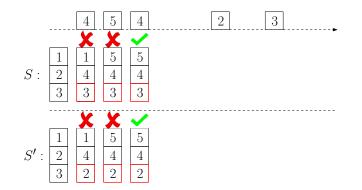
6. From now on, $S'$ will "copy" $S$.

## Proof.

7. If $S$ evicted the page $p'$, $S'$ will evict the page $p^*$. Then, the cache status of $S$ and that of $S'$ will be the same. $S$ and $S'$ will be exactly the same from now on.

### Proof.

7. If $S$ evicted the page $p'$, $S'$ will evict the page $p^*$. Then, the cache status of $S$ and that of $S'$ will be the same. $S$ and $S'$ will be exactly the same from now on.

8. Assume $S$ did not evict $p'(=2)$ before we see $p'(=2)$.

## Proof.

7. If $S$ evicted the page $p'$, $S'$ will evict the page $p^*$. Then, the cache status of $S$ and that of $S'$ will be the same. $S$ and $S'$ will be exactly the same from now on.

8. Assume $S$ did not evict $p'(=2)$ before we see $p'(=2)$.

Proof.

**Proof.**

## Proof.

**Proof.**

9. If $S$ evicts $p^*(=3)$ for $p'(=2)$, then $S$ won't be optimum. Assume otherwise.

## Proof.

9. If $S$ evicts $p^*(=3)$ for $p'(=2)$, then $S$ won't be optimum.
   Assume otherwise.

## Proof.

9. If $S$ evicts $p^*(=3)$ for $p'(=2)$, then $S$ won't be optimum. Assume otherwise.

## Proof.

9. If $S$ evicts $p^*(=3)$ for $p'(=2)$, then $S$ won't be optimum. Assume otherwise.

## Proof.

9. If $S$ evicts $p^*(=3)$ for $p'(=2)$, then $S$ won't be optimum. Assume otherwise.

10. So far, $S'$ has 1 less page-miss than $S$ does.

### Proof.

- **9** If $S$ evicts $p^*(=3)$ for $p'(=2)$, then $S$ won't be optimum. Assume otherwise.
- **10** So far, $S'$ has 1 less page-miss than $S$ does.
- **11** The status of $S'$ and that of $S$ only differ by 1 page.

**Proof.**

**Proof.**

12. We can then guarantee that $S'$ make at most the same number of page-misses as $S$ does.

## Proof.

12. We can then guarantee that $S'$ make at most the same number of page-misses as $S$ does.
    - Idea: if $S$ has a page-hit and $S'$ has a page-miss, we use the opportunity to make the status of $S'$ the same as that of $S$. □

- Thus, we have shown how to create another solution $S'$ with the same number of page-misses as that of the optimum solution $S$. Thus, we proved

**Lemma** Assume at time 1 a page fault happens and there are no empty pages in the cache. Let $p^*$ be the page in cache that is not requested until furthest in the future. There is an optimum solution in which $p^*$ is evicted at time 1.

- Thus, we have shown how to create another solution $S'$ with the same number of page-misses as that of the optimum solution $S$. Thus, we proved

**Lemma** Assume at time 1 a page fault happens and there are no empty pages in the cache. Let $p^*$ be the page in cache that is not requested until furthest in the future. It is safe to evict $p^*$ at time 1.

- Thus, we have shown how to create another solution $S'$ with the same number of page-misses as that of the optimum solution $S$. Thus, we proved

**Lemma** Assume at time 1 a page fault happens and there are no empty pages in the cache. Let $p^*$ be the page in cache that is not requested until furthest in the future. It is safe to evict $p^*$ at time 1.

**Theorem** The furthest-in-future strategy is optimum.

1. **for** $t \leftarrow 1$ to $T$ **do**
2.     **if** $\rho_t$ is in cache, **then** do nothing
3.     **else if** there is an empty page in cache, **then**
4.         evict the empty page and load $\rho_t$ in cache
5.     **else**
6.         $p^* \leftarrow$ the page in cache that is not used furthest in the future
7.         evict $p^*$ and load $\rho_t$ in cache

**Q:** How can we make the algorithm as fast as possible?

**A:**

**Q:** How can we make the algorithm as fast as possible?

**A:**
- The running time can be made to be $O(n + T \log k)$.

**Q:** How can we make the algorithm as fast as possible?

**A:**

- The running time can be made to be $O(n + T \log k)$.
- For each page $p$, use a linked list to store the time steps in which $p$ is requested.

**Q:** How can we make the algorithm as fast as possible?

**A:**

- The running time can be made to be $O(n + T \log k)$.
- For each page $p$, use a linked list to store the time steps in which $p$ is requested.
  - We can find the next time a page is requested easily.

**Q:** How can we make the algorithm as fast as possible?

**A:**

- The running time can be made to be $O(n + T \log k)$.
- For each page $p$, use a linked list to store the time steps in which $p$ is requested.
  - We can find the next time a page is requested easily.
- Use a priority queue data structure to hold all the pages in cache, so that we can easily find the page that is requested furthest in the future.

1. **for** every $p \leftarrow 1$ to $n$ **do**
2.     $lists[p] \leftarrow$ linked list of times in which $p$ is requested, in increasing order \\\\put $\infty$ at the end of the list
3.     $pointer[p] \leftarrow$ head of $lists[p]$
4.     $nexttime[p] \leftarrow$ value pointed by $pointer[p]$
5. $Q \leftarrow$ empty priority queue
6. **for** every $t \leftarrow 1$ to $T$ **do**
7.     move $pointer[\rho_t]$ to right by one position
8.     $nexttime[\rho_t] \leftarrow$ value pointed by $pointer[\rho_t]$
9.     **if** $\rho_t \in Q$ **then** $Q$.update-priority($\rho_t, nexttime[\rho_t]$), continue
10.    **if** $Q$ has size $k$ **then** $p \leftarrow Q$.extract-max() and evict $p$
11.    load $\rho_t$
12.    add $\rho_t$ to $Q$ with priority value $nexttime[\rho_t]$

# Outline

# Encoding Letters Using Bits

- 8 letters $a, b, c, d, e, f, g, h$ in a language
- need to encode a message using bits
- idea: use 3 bits per letter

| $a$ | $b$ | $c$ | $d$ | $e$ | $f$ | $g$ | $h$ |
|-----|-----|-----|-----|-----|-----|-----|-----|
| 000 | 001 | 010 | 011 | 100 | 101 | 110 | 111 |

$$deacfg \rightarrow 011100000010101110$$

**Q:** Can we have a better encoding scheme?

- Seems unlikely: must use 3 bits per letter

**Q:** What if some letters appear more frequently than the others?

**Q:** If some letters appear more frequently than the others, can we have a better encoding scheme?

**A:** Using variable-length encoding scheme might be more efficient.

**Idea**
- using fewer bits for letters that are more frequently used, and more bits for letters that are less frequently used.

**Q:** What is the issue with the following encoding scheme?

- $a$: 0      $b$: 1      $c$: 00

**Q:** What is the issue with the following encoding scheme?
- $a$: 0      $b$: 1      $c$: 00

**A:** Can not guarantee a unique decoding. For example, $00$ can be decoded to $aa$ or $c$.

**Q:** What is the issue with the following encoding scheme?

- $a$: 0    $b$: 1    $c$: 00

**A:** Can not guarantee a unique decoding. For example, $00$ can be decoded to $aa$ or $c$.

**Solution**

Use prefix codes to guarantee a unique decoding.

# Prefix Codes

**Def.** A prefix code for a set $S$ of letters is a function $\gamma : S \to \{0, 1\}^*$ such that for two distinct $x, y \in S$, $\gamma(x)$ is not a prefix of $\gamma(y)$.

# Prefix Codes

**Def.** A prefix code for a set $S$ of letters is a function $\gamma : S \to \{0, 1\}^*$ such that for two distinct $x, y \in S$, $\gamma(x)$ is not a prefix of $\gamma(y)$.

| $a$ | $b$ | $c$ | $d$ |
|-----|------|------|-----|
| 001 | 0000 | 0001 | 100 |

| $e$ | $f$ | $g$ | $h$ |
|-----|------|------|-----|
| 11 | 1010 | 1011 | 01 |

# Prefix Codes Guarantee Unique Decoding

- Reason: there is only one way to cut the first code.

# Prefix Codes Guarantee Unique Decoding

- Reason: there is only one way to cut the first code.

| $a$ | $b$ | $c$ | $d$ |
|------|------|------|------|
| 001 | 0000 | 0001 | 100 |

| $e$ | $f$ | $g$ | $h$ |
|------|------|------|------|
| 11 | 1010 | 1011 | 01 |

# Prefix Codes Guarantee Unique Decoding

- Reason: there is only one way to cut the first code.

| $a$ | $b$ | $c$ | $d$ |
|-----|------|------|-----|
| 001 | 0000 | 0001 | 100 |

| $e$ | $f$ | $g$ | $h$ |
|-----|------|------|-----|
| 11 | 1010 | 1011 | 01 |



- 00010011000000010111101000001001

# Prefix Codes Guarantee Unique Decoding

- Reason: there is only one way to cut the first code.

| $a$ | $b$ | $c$ | $d$ |
|-----|------|------|-----|
| 001 | 0000 | 0001 | 100 |

| $e$ | $f$ | $g$ | $h$ |
|-----|------|------|-----|
| 11 | 1010 | 1011 | 01 |



- 0001/0011000000010111101000001001

- c

# Prefix Codes Guarantee Unique Decoding

- Reason: there is only one way to cut the first code.

| $a$ | $b$ | $c$ | $d$ |
|------|------|------|------|
| 001 | 0000 | 0001 | 100 |

| $e$ | $f$ | $g$ | $h$ |
|------|------|------|------|
| 11 | 1010 | 1011 | 01 |



- 0001/001/10000001011110100001001
- ca

# Prefix Codes Guarantee Unique Decoding

- Reason: there is only one way to cut the first code.

| $a$ | $b$ | $c$ | $d$ |
|-----|------|------|-----|
| 001 | 0000 | 0001 | 100 |

| $e$ | $f$ | $g$ | $h$ |
|-----|------|------|-----|
| 11 | 1010 | 1011 | 01 |



- 0001/001/100/000001011110100001001
- cad

# Prefix Codes Guarantee Unique Decoding

- Reason: there is only one way to cut the first code.

| $a$ | $b$ | $c$ | $d$ |
|-----|------|------|-----|
| 001 | 0000 | 0001 | 100 |

| $e$ | $f$ | $g$ | $h$ |
|-----|------|------|-----|
| 11 | 1010 | 1011 | 01 |



- 0001/001/100/<span style="color:red">0000</span>/01011110100001001
- cad<span style="color:red">b</span>

# Prefix Codes Guarantee Unique Decoding

- Reason: there is only one way to cut the first code.

| $a$ | $b$ | $c$ | $d$ |
|-----|------|------|-----|
| 001 | 0000 | 0001 | 100 |

| $e$ | $f$ | $g$ | $h$ |
|-----|------|------|-----|
| 11  | 1010 | 1011 | 01  |



- 0001/001/100/0000/01/011110100001001
- cadbh

# Prefix Codes Guarantee Unique Decoding

- Reason: there is only one way to cut the first code.

| $a$ | $b$ | $c$ | $d$ |
|------|------|------|------|
| 001 | 0000 | 0001 | 100 |

| $e$ | $f$ | $g$ | $h$ |
|------|------|------|------|
| 11 | 1010 | 1011 | 01 |



- 0001/001/100/0000/01/01/1110100001001
- cadbhh

# Prefix Codes Guarantee Unique Decoding

- Reason: there is only one way to cut the first code.

| $a$ | $b$ | $c$ | $d$ |
|------|------|------|------|
| 001 | 0000 | 0001 | 100 |

| $e$ | $f$ | $g$ | $h$ |
|------|------|------|------|
| 11 | 1010 | 1011 | 01 |



- 0001/001/100/0000/01/01/<span style="color:red">11</span>/10100001001
- cadbhh<span style="color:red">e</span>

# Prefix Codes Guarantee Unique Decoding

- Reason: there is only one way to cut the first code.

| $a$ | $b$ | $c$ | $d$ |
|-----|------|------|-----|
| 001 | 0000 | 0001 | 100 |

| $e$ | $f$ | $g$ | $h$ |
|-----|------|------|-----|
| 11 | 1010 | 1011 | 01 |



- 0001/001/100/0000/01/01/11/1010/0001001
- cadbhhef

# Prefix Codes Guarantee Unique Decoding

- Reason: there is only one way to cut the first code.

| $a$ | $b$ | $c$ | $d$ |
|-----|------|------|-----|
| 001 | 0000 | 0001 | 100 |

| $e$ | $f$ | $g$ | $h$ |
|-----|------|------|-----|
| 11 | 1010 | 1011 | 01 |



- 0001/001/100/0000/01/01/11/1010/0001/001
- cadbhhefc

- Reason: there is only one way to cut the first code.

| $a$ | $b$ | $c$ | $d$ |
|-----|------|------|-----|
| 001 | 0000 | 0001 | 100 |

| $e$ | $f$ | $g$ | $h$ |
|-----|------|------|-----|
| 11 | 1010 | 1011 | 01 |



- 0001/001/100/0000/01/01/11/1010/0001/001/
- cadbhhefca

Properties of Encoding Tree

## Properties of Encoding Tree

- Rooted binary tree

## Properties of Encoding Tree

- Rooted binary tree
- Left edges labelled 0 and right edges labelled 1

### Properties of Encoding Tree

- Rooted binary tree
- Left edges labelled 0 and right edges labelled 1
- A leaf corresponds to a code for some letter

## Properties of Encoding Tree

- Rooted binary tree
- Left edges labelled 0 and right edges labelled 1
- A leaf corresponds to a code for some letter
- If coding scheme is not wasteful: a non-leaf has exactly two children

## Properties of Encoding Tree

- Rooted binary tree
- Left edges labelled 0 and right edges labelled 1
- A leaf corresponds to a code for some letter
- If coding scheme is not wasteful: a non-leaf has exactly two children

## Best Prefix Codes

**Input:** frequencies of letters in a message

**Output:** prefix coding scheme with the shortest encoding for the message

| letters | $a$ | $b$ | $c$ | $d$ | $e$ | |
|---|---|---|---|---|---|---|
| frequencies | 18 | 3 | 4 | 6 | 10 | |
| | | | | | | |
| | | | | | | |
| | | | | | | |
| | | | | | | |



scheme 1                    scheme 2                    scheme 3

| letters | a | b | c | d | e | |
|---|---|---|---|---|---|---|
| frequencies | 18 | 3 | 4 | 6 | 10 | |
| scheme 1 length | 2 | 3 | 3 | 2 | 2 | total = 89 |
| scheme 2 length | 1 | 3 | 3 | 3 | 3 | total = 87 |
| scheme 3 length | 1 | 4 | 4 | 3 | 2 | total = 84 |



scheme 1      scheme 2      scheme 3

- Example Input: ($a$: 18, $b$: 3, $c$: 4, $d$: 6, $e$: 10)

- Example Input: ($a$: 18, $b$: 3, $c$: 4, $d$: 6, $e$: 10)

**Q:** What types of decisions should we make?

- Example Input: ($a$: 18, $b$: 3, $c$: 4, $d$: 6, $e$: 10)

**Q:** What types of decisions should we make?

- Can we directly give a code for some letter?

- Example Input: ($a$: 18, $b$: 3, $c$: 4, $d$: 6, $e$: 10)

**Q:** What types of decisions should we make?

- Can we directly give a code for some letter?
- Hard to design a strategy; residual problem is complicated.

- Example Input: ($a$: 18, $b$: 3, $c$: 4, $d$: 6, $e$: 10)

**Q:** What types of decisions should we make?

- Can we directly give a code for some letter?
- Hard to design a strategy; residual problem is complicated.

- Can we partition the letters into left and right sub-trees?

- Example Input: ($a$: 18, $b$: 3, $c$: 4, $d$: 6, $e$: 10)

**Q:** What types of decisions should we make?

- Can we directly give a code for some letter?
- Hard to design a strategy; residual problem is complicated.

- Can we partition the letters into left and right sub-trees?
- Not clear how to design the greedy algorithm

- Example Input: ($a$: 18, $b$: 3, $c$: 4, $d$: 6, $e$: 10)

**Q:** What types of decisions should we make?

- Can we directly give a code for some letter?
- Hard to design a strategy; residual problem is complicated.

- Can we partition the letters into left and right sub-trees?
- Not clear how to design the greedy algorithm

**A:** We can choose two letters and make them brothers in the tree.

# Which Two Letters Can Be Safely Put Together As Brothers?

- Focus on the "structure" of the optimum encoding tree

# Which Two Letters Can Be Safely Put Together As Brothers?

- Focus on the "structure" of the optimum encoding tree
- There are two deepest leaves that are brothers

# Which Two Letters Can Be Safely Put Together As Brothers?

- Focus on the "structure" of the optimum encoding tree
- There are two deepest leaves that are brothers



best to put the two least frenquent symbols here!

# Which Two Letters Can Be Safely Put Together As Brothers?

- Focus on the "structure" of the optimum encoding tree
- There are two deepest leaves that are brothers



best to put the two least frenquent symbols here!

**Lemma** It is safe to make the two least frequent letters brothers.

**Lemma** There is an optimum encoding tree, where the two least frequent letters are brothers.

**Lemma** There is an optimum encoding tree, where the two least frequent letters are brothers.

- So we can irrevocably decide to make the two least frequent letters brothers.

**Lemma** There is an optimum encoding tree, where the two least frequent letters are brothers.

- So we can irrevocably decide to make the two least frequent letters brothers.

**Q:** Is the residual problem another instance of the best prefix codes problem?

**Lemma** There is an optimum encoding tree, where the two least frequent letters are brothers.

- So we can irrevocably decide to make the two least frequent letters brothers.

**Q:** Is the residual problem another instance of the best prefix codes problem?

**A:** Yes, though it is not immediate to see why.

- $f_x$: the frequency of the letter $x$ in the support.
- $x_1$ and $x_2$: the two letters we decided to put together.
- $d_x$ the depth of letter $x$ in our output encoding tree.



$$\sum_{x \in S} f_x d_x$$
$$= \sum_{x \in S \setminus \{x_1, x_2\}} f_x d_x + f_{x_1} d_{x_1} + f_{x_2} d_{x_2}$$
$$= \sum_{x \in S \setminus \{x_1, x_2\}} f_x d_x + (f_{x_1} + f_{x_2}) d_{x_1}$$

- $f_x$: the frequency of the letter $x$ in the support.
- $x_1$ and $x_2$: the two letters we decided to put together.
- $d_x$ the depth of letter $x$ in our output encoding tree.



$$\sum_{x \in S} f_x d_x$$

$$= \sum_{x \in S \setminus \{x_1, x_2\}} f_x d_x + f_{x_1} d_{x_1} + f_{x_2} d_{x_2}$$

$$= \sum_{x \in S \setminus \{x_1, x_2\}} f_x d_x + (f_{x_1} + f_{x_2}) d_{x_1}$$

- $f_x$: the frequency of the letter $x$ in the support.
- $x_1$ and $x_2$: the two letters we decided to put together.
- $d_x$ the depth of letter $x$ in our output encoding tree.



$$\sum_{x \in S} f_x d_x$$

$$= \sum_{x \in S \setminus \{x_1, x_2\}} f_x d_x + f_{x_1} d_{x_1} + f_{x_2} d_{x_2}$$

$$= \sum_{x \in S \setminus \{x_1, x_2\}} f_x d_x + (f_{x_1} + f_{x_2}) d_{x_1}$$

Def: $f_{x'} = f_{x_1} + f_{x_2}$

- $f_x$: the frequency of the letter $x$ in the support.
- $x_1$ and $x_2$: the two letters we decided to put together.
- $d_x$ the depth of letter $x$ in our output encoding tree.



$$\sum_{x \in S} f_x d_x$$

$$= \sum_{x \in S \setminus \{x_1, x_2\}} f_x d_x + f_{x_1} d_{x_1} + f_{x_2} d_{x_2}$$

$$= \sum_{x \in S \setminus \{x_1, x_2\}} f_x d_x + (f_{x_1} + f_{x_2}) d_{x_1}$$

$$= \sum_{x \in S \setminus \{x_1, x_2\}} f_x d_x + f_{x'} (d_{x'} + 1)$$

Def: $f_{x'} = f_{x_1} + f_{x_2}$

- $f_x$: the frequency of the letter $x$ in the support.
- $x_1$ and $x_2$: the two letters we decided to put together.
- $d_x$ the depth of letter $x$ in our output encoding tree.



$$\sum_{x \in S} f_x d_x$$

$$= \sum_{x \in S \setminus \{x_1, x_2\}} f_x d_x + f_{x_1} d_{x_1} + f_{x_2} d_{x_2}$$

$$= \sum_{x \in S \setminus \{x_1, x_2\}} f_x d_x + (f_{x_1} + f_{x_2}) d_{x_1}$$

$$= \sum_{x \in S \setminus \{x_1, x_2\}} f_x d_x + f_{x'}(d_{x'} + 1)$$

$$= \sum_{x \in S \setminus \{x_1, x_2\} \cup \{x'\}} f_x d_x + f_{x'}$$

Def: $f_{x'} = f_{x_1} + f_{x_2}$

- $f_x$: the frequency of the letter $x$ in the support.
- $x_1$ and $x_2$: the two letters we decided to put together.
- $d_x$ the depth of letter $x$ in our output encoding tree.



encoding tree for
$S \setminus \{x_1, x_2\} \cup \{x'\}$

Def: $f_{x'} = f_{x_1} + f_{x_2}$

$$
\begin{aligned}
& \sum_{x \in S} f_x d_x \\
=& \sum_{x \in S \setminus \{x_1, x_2\}} f_x d_x + f_{x_1} d_{x_1} + f_{x_2} d_{x_2} \\
=& \sum_{x \in S \setminus \{x_1, x_2\}} f_x d_x + (f_{x_1} + f_{x_2}) d_{x_1} \\
=& \sum_{x \in S \setminus \{x_1, x_2\}} f_x d_x + f_{x'}(d_{x'} + 1) \\
=& \sum_{x \in S \setminus \{x_1, x_2\} \cup \{x'\}} f_x d_x + f_{x'}
\end{aligned}
$$

In order to minimize

$$\sum_{x \in S} f_x d_x,$$

we need to minimize

$$\sum_{x \in S \setminus \{x_1, x_2\} \cup \{x'\}} f_x d_x,$$

subject to that $d$ is the depth function for an encoding tree of $S \setminus \{x_1, x_2\}$.

- This is exactly the best prefix codes problem, with letters $S \setminus \{x_1, x_2\} \cup \{x'\}$ and frequency vector $f$!

# Example



$\textcircled{A}$ $^{27}$  $\textcircled{B}$ $^{15}$  $\textcircled{C}$ $^{11}$  $\textcircled{D}$ $^{9}$  $\textcircled{E}$ $^{8}$  $\textcircled{F}$ $^{5}$

# Example

# Example

# Example

# Example

# Example

# Example

# Example



$A : 00$
$B : 10$
$C : 010$
$D : 011$
$E : 110$
$F : 111$

**Def.** The codes given the greedy algorithm is called the Huffman codes.

**Def.** The codes given the greedy algorithm is called the Huffman codes.

Huffman$(S, f)$

1. **while** $|S| > 1$ **do**
2.     let $x_1, x_2$ be the two letters with the smallest $f$ values
3.     introduce a new letter $x'$ and let $f_{x'} = f_{x_1} + f_{x_2}$
4.     let $x_1$ and $x_2$ be the two children of $x'$
5.     $S \leftarrow S \setminus \{x_1, x_2\} \cup \{x'\}$
6. **return** the tree constructed

# Algorithm using Priority Queue

Huffman$(S, f)$

1. $Q \leftarrow$ build-priority-queue$(S)$
2. **while** $Q$.size $> 1$ **do**
3.     $x_1 \leftarrow Q$.extract-min$()$
4.     $x_2 \leftarrow Q$.extract-min$()$
5.     introduce a new letter $x'$ and let $f_{x'} = f_{x_1} + f_{x_2}$
6.     let $x_1$ and $x_2$ be the two children of $x'$
7.     $Q$.insert$(x')$
8. **return** the tree constructed

# Outline

# Summary for Greedy Algorithms

## Greedy Algorithm

- Build up the solutions in steps
- At each step, make an irrevocable decision using a "reasonable" strategy

# Summary for Greedy Algorithms

## Greedy Algorithm

- Build up the solutions in steps
- At each step, make an <span style="color:red">irrevocable</span> decision using a "reasonable" strategy

- Interval scheduling problem: schedule the job $j^*$ with the earliest deadline

# Summary for Greedy Algorithms

**Greedy Algorithm**

- Build up the solutions in steps
- At each step, make an <span style="color:red">irrevocable</span> decision using a "reasonable" strategy

- Interval scheduling problem: schedule the job $j^*$ with the earliest deadline
- Offline Caching: evict the page that is used furthest in the future

# Summary for Greedy Algorithms

**Greedy Algorithm**
- Build up the solutions in steps
- At each step, make an <span style="color:red">irrevocable</span> decision using a "reasonable" strategy

- Interval scheduling problem: schedule the job $j^*$ with the earliest deadline
- Offline Caching: evict the page that is used furthest in the future
- Huffman codes: make the two least frequent letters brothers

# Summary for Greedy Algorithms

**Analysis of Greedy Algorithm**

- Prove that the reasonable strategy is "safe" (key)
- Show that the remaining task after applying the strategy is to solve a (many) smaller instance(s) of the same problem (usually easy)

# Summary for Greedy Algorithms

**Analysis of Greedy Algorithm**
- Prove that the reasonable strategy is "safe" (key)
- Show that the remaining task after applying the strategy is to solve a (many) smaller instance(s) of the same problem (usually easy)

**Def.** A strategy is "safe" if there is always an optimum solution that "agrees with" the decision made according to the strategy.

# Proving a Strategy is Safe

- Take an arbitrary optimum solution $S$

# Proving a Strategy is Safe

- Take an arbitrary optimum solution $S$
- If $S$ agrees with the decision made according to the strategy, done

# Proving a Strategy is Safe

- Take an arbitrary optimum solution $S$
- If $S$ agrees with the decision made according to the strategy, done
- So assume $S$ does not agree with decision

# Proving a Strategy is Safe

- Take an arbitrary optimum solution $S$
- If $S$ agrees with the decision made according to the strategy, done
- So assume $S$ does not agree with decision
- Change $S$ slightly to another optimum solution $S'$ that agrees with the decision

# Proving a Strategy is Safe

- Take an arbitrary optimum solution $S$
- If $S$ agrees with the decision made according to the strategy, done
- So assume $S$ does not agree with decision
- Change $S$ slightly to another optimum solution $S'$ that agrees with the decision
  - Interval scheduling problem: exchange $j^*$ with the first job in an optimal solution

# Proving a Strategy is Safe

- Take an arbitrary optimum solution $S$
- If $S$ agrees with the decision made according to the strategy, done
- So assume $S$ does not agree with decision
- Change $S$ slightly to another optimum solution $S'$ that agrees with the decision
  - Interval scheduling problem: exchange $j^*$ with the first job in an optimal solution
  - Offline caching: a complicated "copying" algorithm

# Proving a Strategy is Safe

- Take an arbitrary optimum solution $S$
- If $S$ agrees with the decision made according to the strategy, done
- So assume $S$ does not agree with decision
- Change $S$ slightly to another optimum solution $S'$ that agrees with the decision
  - Interval scheduling problem: exchange $j^*$ with the first job in an optimal solution
  - Offline caching: a complicated "copying" algorithm
  - Huffman codes: move the two least frequent letters to the deepest leaves.

# Summary for Greedy Algorithms

## Analysis of Greedy Algorithm

- Prove that the reasonable strategy is "safe" (key)
- Show that the remaining task after applying the strategy is to solve a (many) smaller instance(s) of the same problem (usually easy)

# Summary for Greedy Algorithms

**Analysis of Greedy Algorithm**

- Prove that the reasonable strategy is "safe" (key)
- Show that the remaining task after applying the strategy is to solve a (many) smaller instance(s) of the same problem (usually easy)

- Interval scheduling problem: remove $j^*$ and the jobs it conflicts with

# Summary for Greedy Algorithms

## Analysis of Greedy Algorithm

- Prove that the reasonable strategy is "safe" (key)
- Show that the remaining task after applying the strategy is to solve a (many) smaller instance(s) of the same problem (usually easy)

- Interval scheduling problem: remove $j^*$ and the jobs it conflicts with
- Offline caching: trivial

# Summary for Greedy Algorithms

## Analysis of Greedy Algorithm

- Prove that the reasonable strategy is "safe" (key)
- Show that the remaining task after applying the strategy is to solve a (many) smaller instance(s) of the same problem (usually easy)

- Interval scheduling problem: remove $j^*$ and the jobs it conflicts with
- Offline caching: trivial
- Huffman codes: merge two letters into one