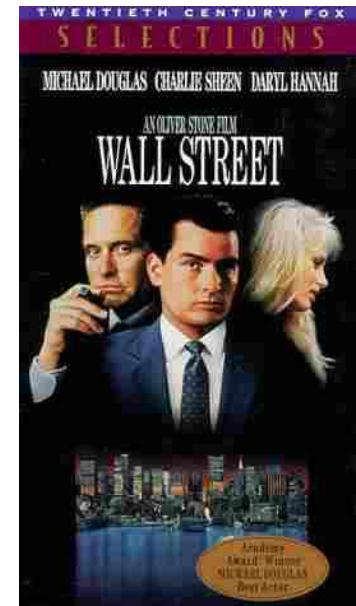# Greedy Algorithms

# Outline

# An Informal Definition

- An algorithm is greedy if it builds up a solution in small steps, choosing a decision at each step myopically to optimize some underlying criterion.

- Many different greedy algorithms for the same problem, each one locally, incrementally optimizing some different measure on its way to a solution.

> Greed is good. Greed is right. Greed works. Greed clarifies, cuts through, and captures the essence of the evolutionary spirit.
>     - Gordon Gecko (Michael Douglas)



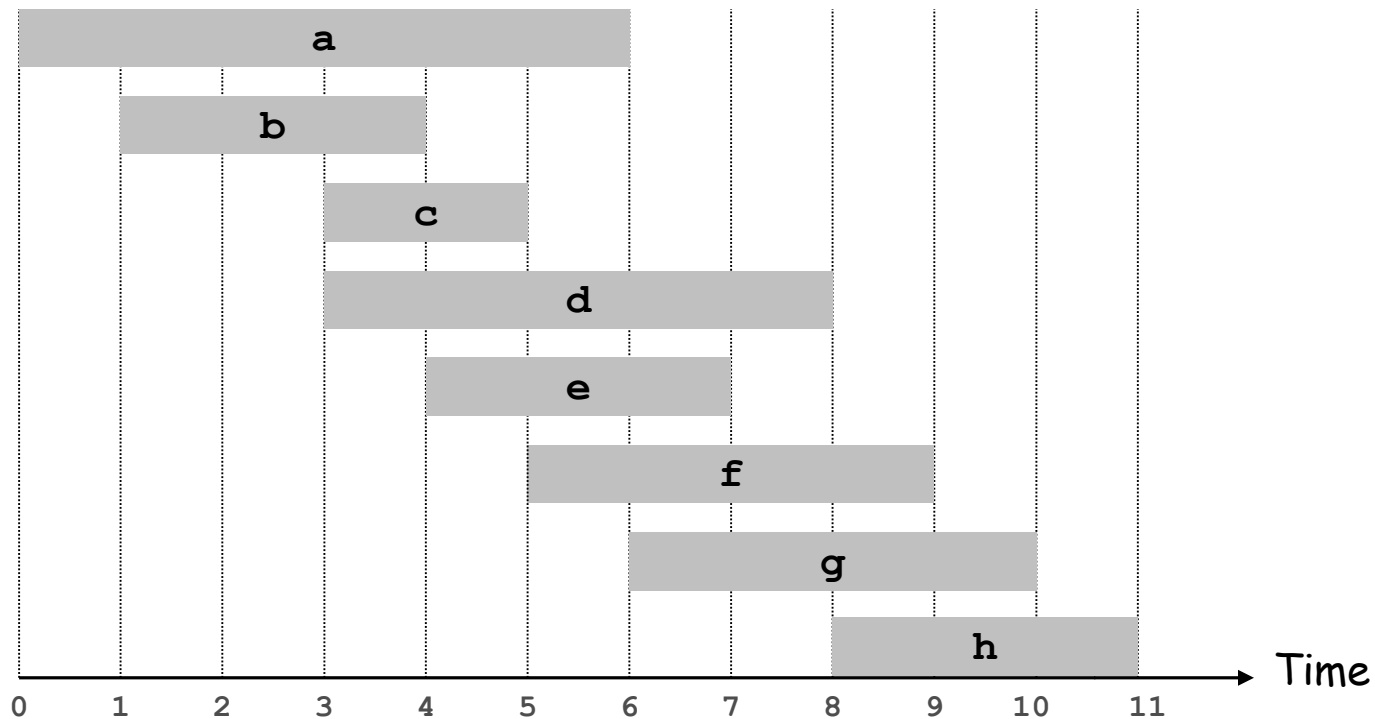Q. for short-sighted greed in the design of algorithms:
- Is greed good?
- Does greed work?

# Interval Scheduling

Interval scheduling.

- Job j starts at $s_j$ and finishes at $f_j$.
- Two jobs compatible if they don't overlap.
- Goal: find maximum subset of mutually compatible jobs.

# Interval Scheduling:  Greedy Algorithms

Greedy template.  Consider jobs in some natural order.
Take each job provided it's compatible with the ones already taken.

- [Earliest start time]  Consider jobs in ascending order of $s_j$.

- [Earliest finish time]  Consider jobs in ascending order of $f_j$.

- [Shortest interval]  Consider jobs in ascending order of $f_j - s_j$.

- [Fewest conflicts]  For each job j, count the number of conflicting jobs $c_j$. Schedule in ascending order of $c_j$.

# Interval Scheduling: Greedy Algorithms

Greedy template. Consider jobs in some natural order.
Take each job provided it's compatible with the ones already taken.

counterexample for earliest start time

counterexample for shortest interval

counterexample for fewest conflicts

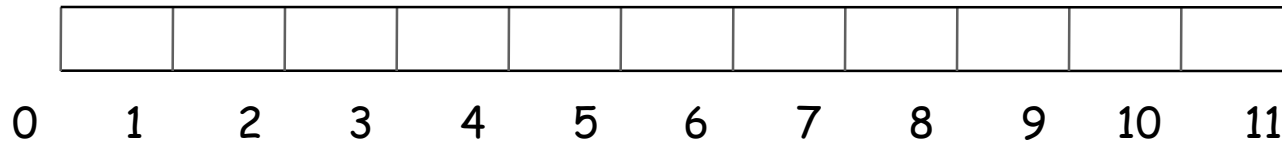# Interval Scheduling:  Greedy Algorithm

Greedy algorithm.  Consider jobs in increasing order of finish time.
Take each job provided it's compatible with the ones already taken.

```
Sort jobs by finish times so that f₁ ≤ f₂ ≤ ... ≤ fₙ.

          ╱  set of jobs selected
         ↙
A ← φ
for j = 1 to n {
    if (job j compatible with A)
        A ← A ∪ {j}
}
return A
```
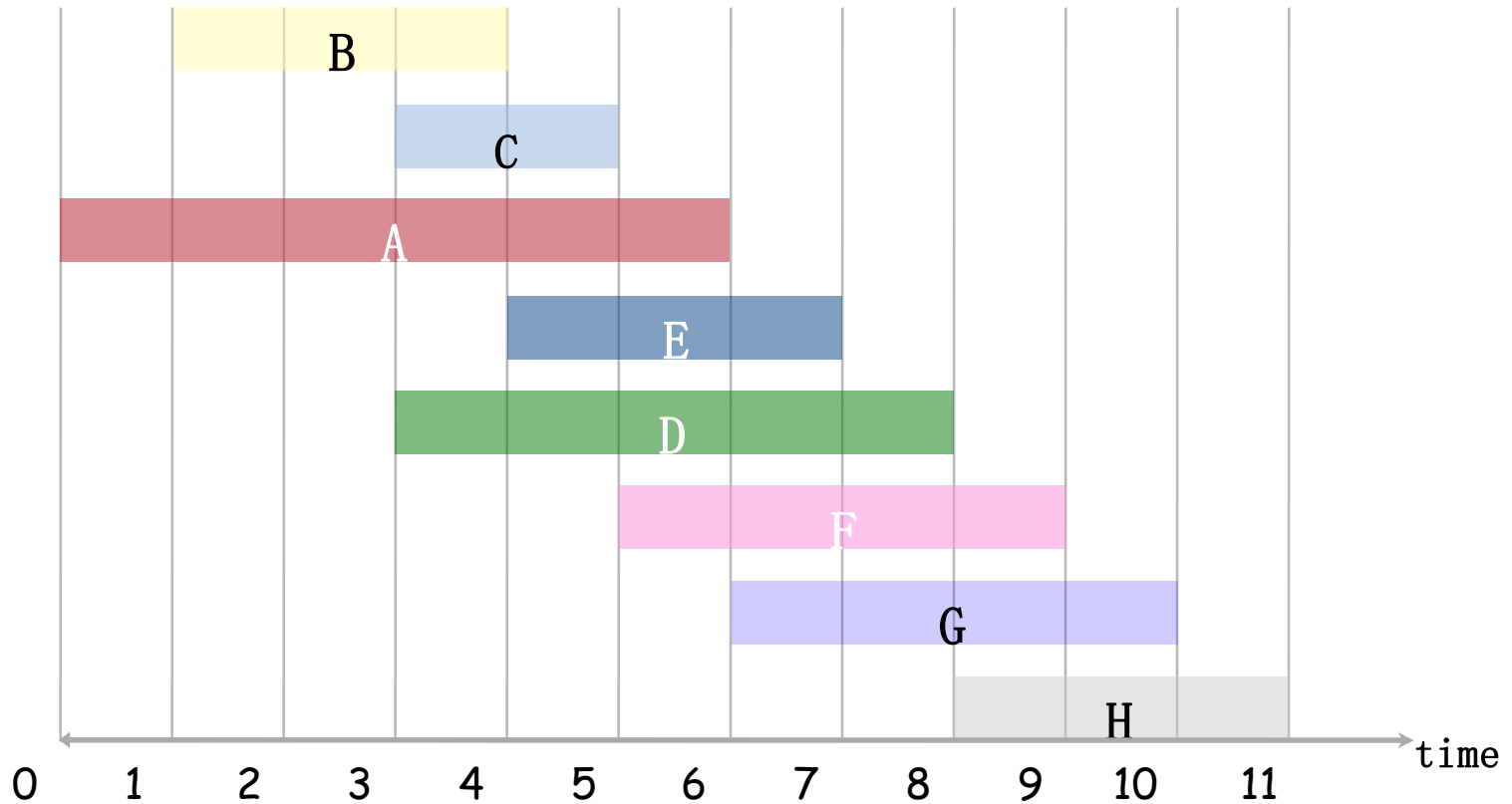
Proposition.  Can implement earliest-finish-time first in $O(n \log n)$ time.
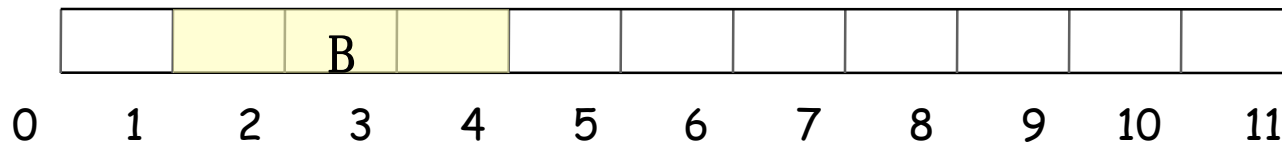- Keep track of job $j*$ that was added last to $S$.
- Job $j$ is compatible with $S$ iff $s_j \geq f_{j*}$.
- Sorting by finish times takes $O(n \log n)$ time.

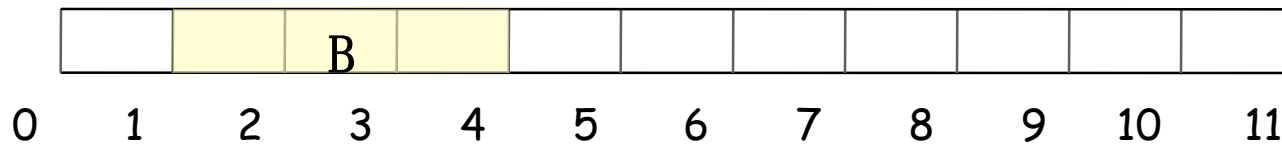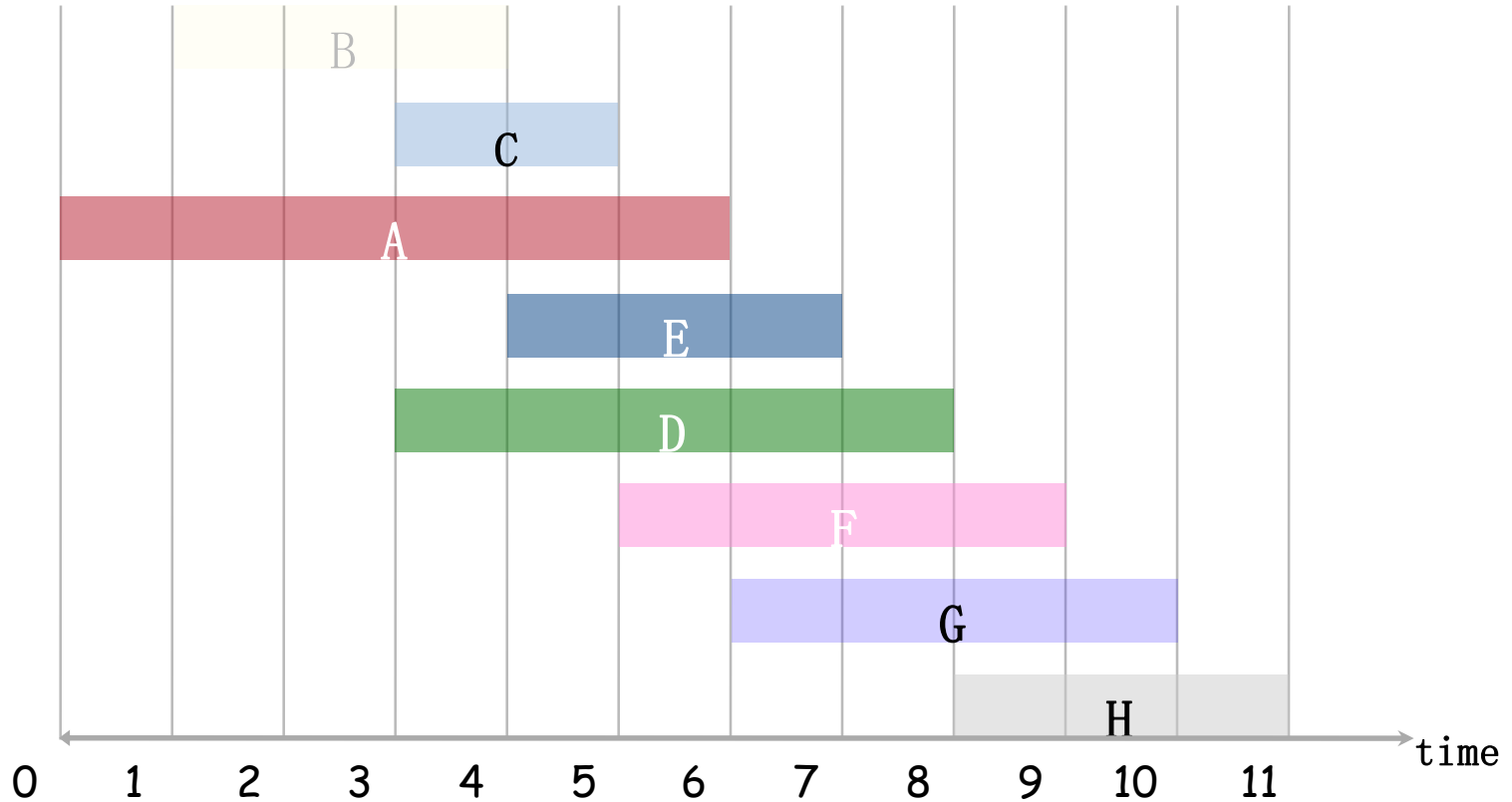# Earliest-finish-time-first algorithm demo
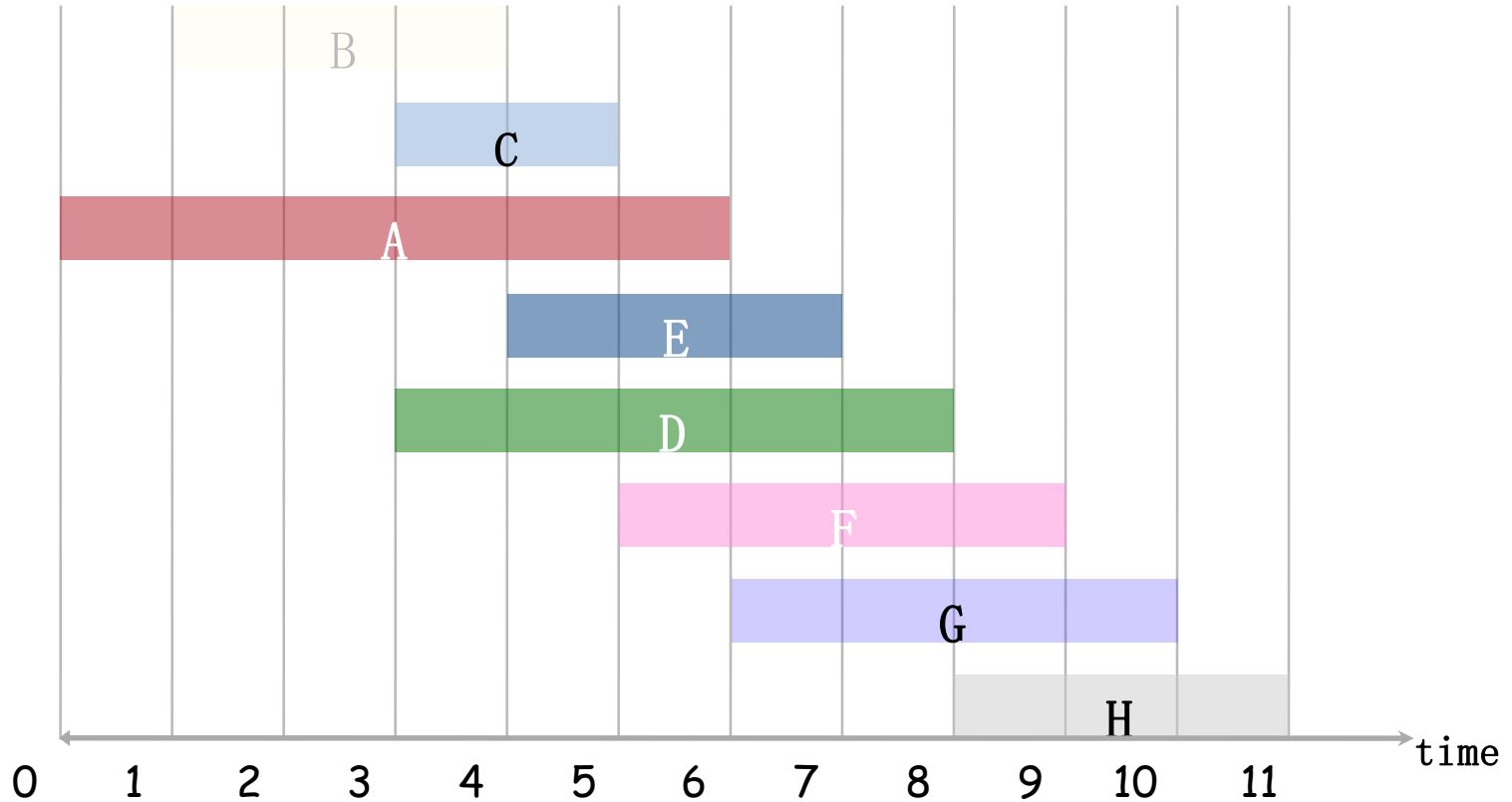
# Earliest-finish-time-first algorithm demo



job B is compatible (add to schedule)

Earliest-finish-time-first algorithm demo

# Earliest-finish-time-first algorithm demo



job C is incompatible (do not add to schedule)

# Earliest-finish-time-first algorithm demo

# Earliest-finish-time-first algorithm demo



job D is incompatible (do not add to schedule)

Earliest-finish-time-first algorithm demo

Earliest-finish-time-first algorithm demo

job F is incompatible (do not add to schedule)

# Earliest-finish-time-first algorithm demo

# Earliest-finish-time-first algorithm demo



job G is incompatible (do not add to schedule)

# Earliest-finish-time-first algorithm demo



job G is incompatible (do not add to schedule)

# Earliest-finish-time-first algorithm demo

Earliest-finish-time-first algorithm demo
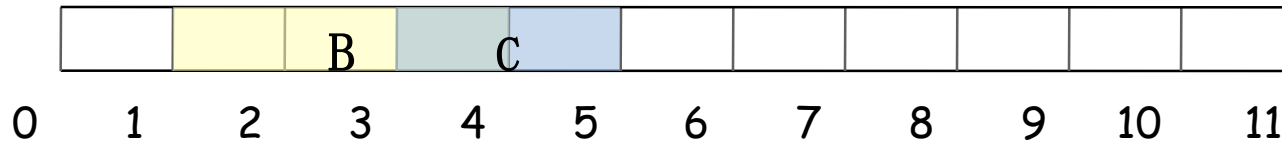
job H is compatible (add to schedule)

# Earliest-finish-time-first algorithm demo

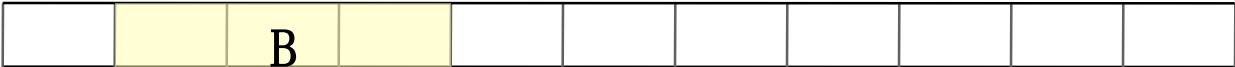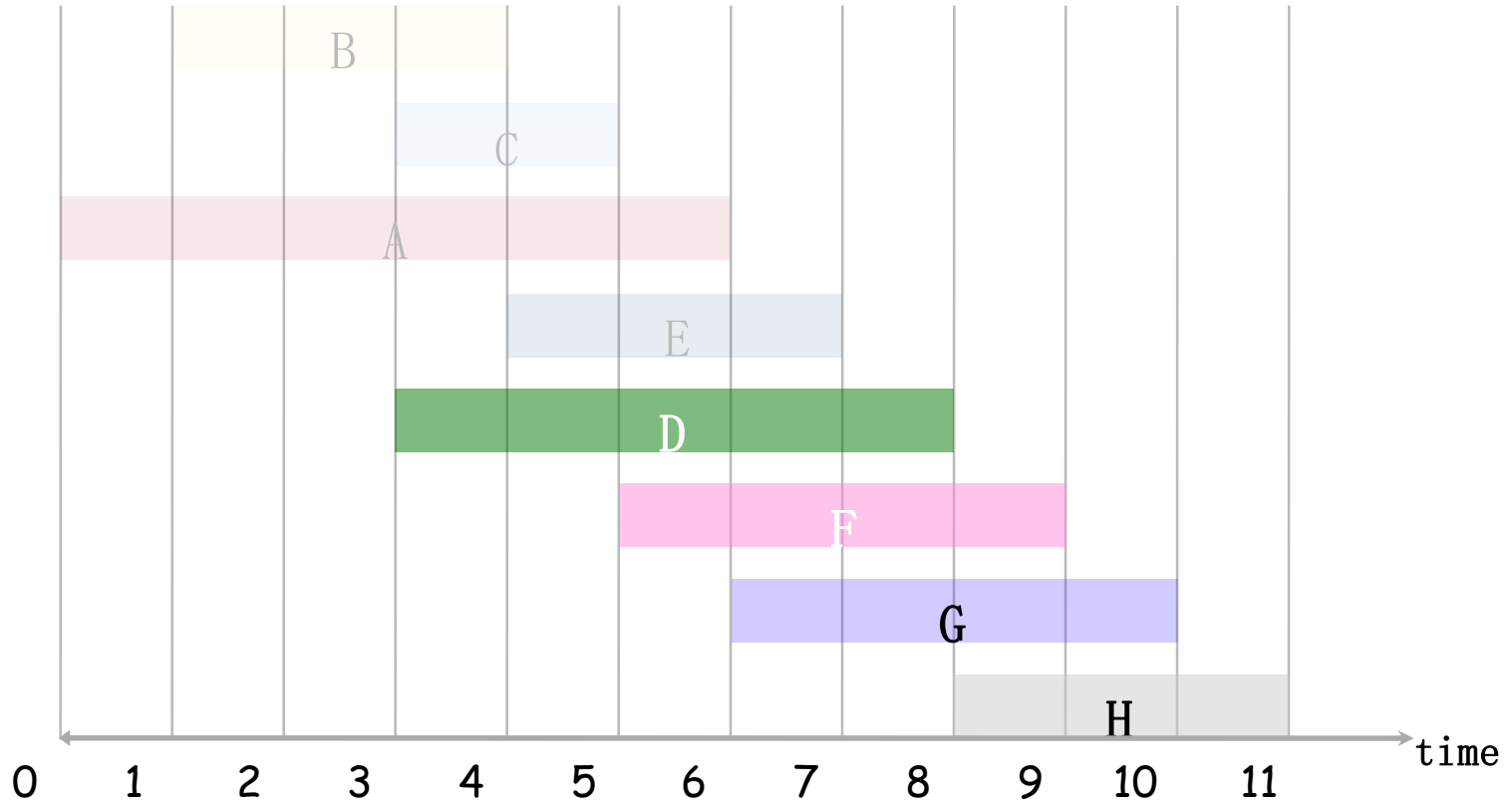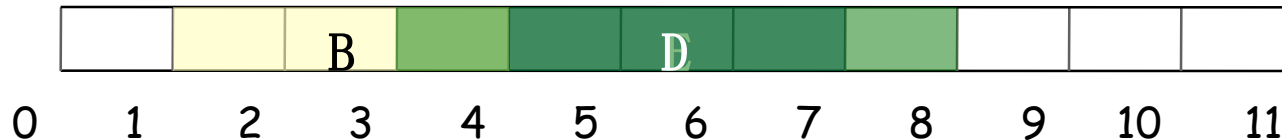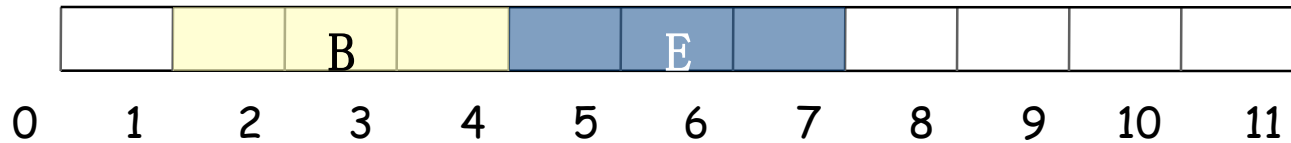| | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| B | | | | | | | | | | | |
| | | | C | | | | | | | | |
| A | | | | | | | | | | | |
| | | | | E | | | | | | | |
| | | D | | | | | | | | | |
| | | | | | F | | | | | | |
| | | | | | | G | | | | | |
| | | | | | | | | H | | | |

Earliest-finish-time-first algorithm demo

done (an optimal set of jobs)

# Interval Scheduling:  Analysis

Theorem.  Greedy algorithm is optimal.

Pf.  (by contradiction)

- Assume greedy is not optimal, and let's see what happens.
- Let $i_1, i_2, \ldots i_k$ denote set of jobs selected by greedy.
- Let $j_1, j_2, \ldots j_m$ denote set of jobs in the optimal solution with $i_1 = j_1, i_2 = j_2, \ldots, i_r = j_r$ for the largest possible value of $r$.

job $i_{r+1}$ finishes before $j_{r+1}$

↓

Greedy:

| $i_1$ | $i_2$ | $i_r$ | $i_{r+1}$ |

OPT:

| $j_1$ | $j_2$ | $j_r$ | $j_{r+1}$ | . . . |

↑

why not replace job $j_{r+1}$
with job $i_{r+1}$?

# Interval Scheduling: Analysis

Theorem. Greedy algorithm is optimal.

Pf. (by contradiction)

- Assume greedy is not optimal, and let's see what happens.
- Let $i_1, i_2, \ldots i_k$ denote set of jobs selected by greedy.
- Let $j_1, j_2, \ldots j_m$ denote set of jobs in the optimal solution with $i_1 = j_1, i_2 = j_2, \ldots, i_r = j_r$ for the largest possible value of r.

job $i_{r+1}$ finishes before $j_{r+1}$

Greedy:

| $i_1$ | $i_2$ | $i_r$ | $i_{r+1}$ |

OPT:

| $j_1$ | $j_2$ | $j_r$ | $i_{r+1}$ | . . . |

solution still feasible and optimal, but contradicts maximality of r.

# Outline

- Interval scheduling

- Interval partitioning

- Scheduling to minimize lateness

- Basic elements of greedy algorithms

- Single-source shortest path: Dijkstra's algorithm

- Minimal spanning trees: Prim, Kruskal

- single-link clustering

# Interval Partitioning

Interval partitioning.

- Lecture j starts at $s_j$ and finishes at $f_j$.
- Goal: find minimum number of classrooms to schedule all lectures so that no two occur at the same time in the same room.

Ex: This schedule uses 4 classrooms to schedule 10 lectures.

# Interval Partitioning

Interval partitioning.

- Lecture j starts at $s_j$ and finishes at $f_j$.
- Goal: find minimum number of classrooms to schedule all lectures so that no two occur at the same time in the same room.

Ex: This schedule uses only 3.

# Interval Partitioning:
# Lower Bound on Optimal Solution

Def.  The depth of a set of open intervals is the maximum number that contain any given time.

Key observation.  Number of classrooms needed ≥ depth.

Ex:  Depth of schedule below = 3 ⇒ schedule below is optimal.

a, b, c all contain 9:30

Q.  Does there always exist a schedule equal to depth of intervals?

# Interval Partitioning:  Greedy Algorithm

Greedy algorithm.  Consider lectures in increasing order of start time: assign lecture to any compatible classroom.

```
Sort intervals by starting time so that s₁ ≤ s₂ ≤ ... ≤ sₙ.
d ← 0  ⟵  number of allocated classrooms

for j = 1 to n {
    if (lecture j is compatible with some classroom k)
        schedule lecture j in classroom k
    else
        allocate a new classroom d + 1
        schedule lecture j in classroom d + 1
        d ← d + 1
}
```

Implementation.  O(n log n).
- For each classroom k, maintain the finish time of the last job added.
- Keep the classrooms in a priority queue.

# Priority Queue

Data structures for dynamically changing set $S$ where each element $v \in S$ has an associated value key$(v)$ that denotes the priority of element $v$; smaller keys represent higher priorities.

A min-heap makes an excellent priority queue.

• Insert$(H, v)$ inserts the item $v$ into heap $H$. If the heap currently has $n$ elements, this takes $O(\log n)$ time.

• FindMin$(H)$ identifies the minimum element in the heap $H$ but does not remove it. This takes $O(1)$ time.

• Delete$(H, i)$ deletes the element in heap position $i$. This is implemented in $O(\log n)$ time for heaps that have $n$ elements.

• ExtractMin$(H)$ identifies and deletes an element with minimum key value from a heap. This is a combination of the preceding two operations, and so it takes $O(\log n)$ time.

# Interval Partitioning:  Greedy Analysis

Observation.  Greedy algorithm never schedules two incompatible lectures in the same classroom.

Theorem.  Greedy algorithm is optimal.
Pf.
- Let d = number of classrooms that the greedy algorithm allocates.
- Classroom d is opened because we needed to schedule a job, say j, that is incompatible with all d-1 other classrooms.
- These d jobs each end after $s_j$.
- Since we sorted by start time, all these incompatibilities are caused by lectures that start no later than $s_j$.
- Thus, we have d lectures overlapping at time $s_j + \varepsilon$.
- Key observation $\Rightarrow$ all schedules use $\geq$ d classrooms.  ▪

# Outline

- Interval scheduling

- Interval partitioning

- Scheduling to minimize lateness

- Basic elements of greedy algorithms

- Single-source shortest path: Dijkstra's algorithm

- Minimal spanning trees: Prim, Kruskal

- single-link clustering

# Scheduling to Minimizing Lateness

Minimizing lateness problem.

- Single resource processes one job at a time.
- Job j requires $t_j$ units of processing time and is due at time $d_j$.
- If j starts at time $s_j$, it finishes at time $f_j = s_j + t_j$.
- Lateness: $\ell_j = \max \{ 0, \ f_j - d_j \}$.
- Goal: schedule all jobs to minimize maximum lateness $L = \max \ell_j$.

Ex:

| | 1 | 2 | 3 | 4 | 5 | 6 |
|---|---|---|---|---|---|---|
| $t_j$ | 3 | 2 | 1 | 4 | 3 | 2 |
| $d_j$ | 6 | 8 | 9 | 9 | 14 | 15 |

lateness = 2          lateness = 0          max lateness = 6

| $d_3 = 9$ | $d_2 = 8$ | $d_6 = 15$ | $d_1 = 6$ | $d_5 = 14$ | $d_4 = 9$ |
|---|---|---|---|---|---|

0    1    2    3    4    5    6    7    8    9    10    11    12    13    14    15

# Minimizing Lateness: Greedy Algorithms

Greedy template. Consider jobs in some order.

- [Shortest processing time first] Consider jobs in ascending order of processing time $t_j$.


- [Earliest deadline first] Consider jobs in ascending order of deadline $d_j$.


- [Smallest slack] Consider jobs in ascending order of slack $d_j - t_j$.

# Minimizing Lateness:  Greedy Algorithms

Greedy template.  Consider jobs in some order.

- [Shortest processing time first]  Consider jobs in ascending order
  of processing time $t_j$.

| | 1 | 2 |
|---|---|---|
| $t_j$ | 1 | 10 |
| $d_j$ | 100 | 10 |

counterexample

- [Smallest slack]  Consider jobs in ascending order of slack $d_j - t_j$.

| | 1 | 2 |
|---|---|---|
| $t_j$ | 5 | 2 |
| $d_j$ | 6 | 4 |

counterexample

# Minimizing Lateness:  Greedy Algorithm

Greedy algorithm.  Earliest deadline first.

```
Sort n jobs by deadline so that d₁ ≤ d₂ ≤ … ≤ dₙ

t ← 0
for j = 1 to n
    Assign job j to interval [t, t + tⱼ]
    sⱼ ← t, fⱼ ← t + tⱼ
    t ← t + tⱼ
output intervals [sⱼ, fⱼ]
```

max lateness = 1

| $d_1 = 6$ | | | $d_2 = 8$ | $d_3 = 9$ | $d_4 = 9$ | | | $d_5 = 14$ | | $d_6 = 15$ |

0   1   2   3   4   5   6   7   8   9   10   11   12   13   14   15

# Minimizing Lateness: No Idle Time

Observation 1. There exists an optimal schedule with no idle time.

an optimal schedule

| | d = 4 | | | d = 6 | | | | | d = 12 | | |
|---|---|---|---|---|---|---|---|---|---|---|---|

0    1    2    3    4    5    6    7    8    9    10    11

an optimal schedule
with no idle time

| d = 4 | | d = 6 | | d = 12 | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|

0    1    2    3    4    5    6    7    8    9    10    11

Observation 2. The earliest-deadline-first schedule has no idle time.
(by the algorithm construction)

# Minimizing Lateness: Inversions

**Def.** Given a schedule S, an inversion is a pair of jobs i and j such that:
i < j but j scheduled before i.

inversion if $i < j$

$f_i$

a schedule with
an inversion

| | | j | i | | | |

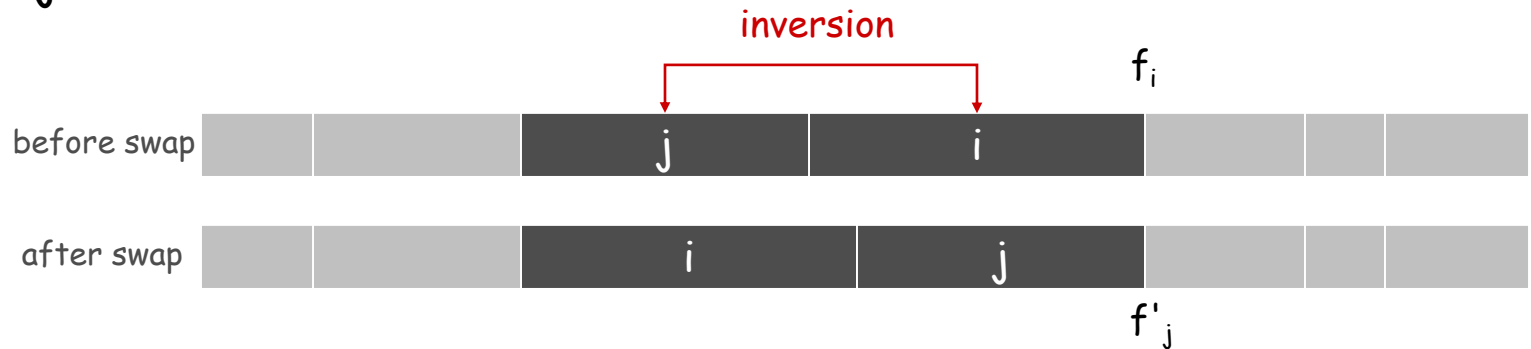[ as before, we assume jobs are numbered so that $d_1 \leq d_2 \leq \dots \leq d_n$ ]

Observation 3. The earliest-deadline-first schedule is the unique idle-free schedule with no inversions.

# Minimizing Lateness: Inversions

Def. Given a schedule S, an inversion is a pair of jobs i and j such that:
i < j but j scheduled before i.

inversion

$f_i$

before swap | j | i

after swap | i | j

$f'_j$

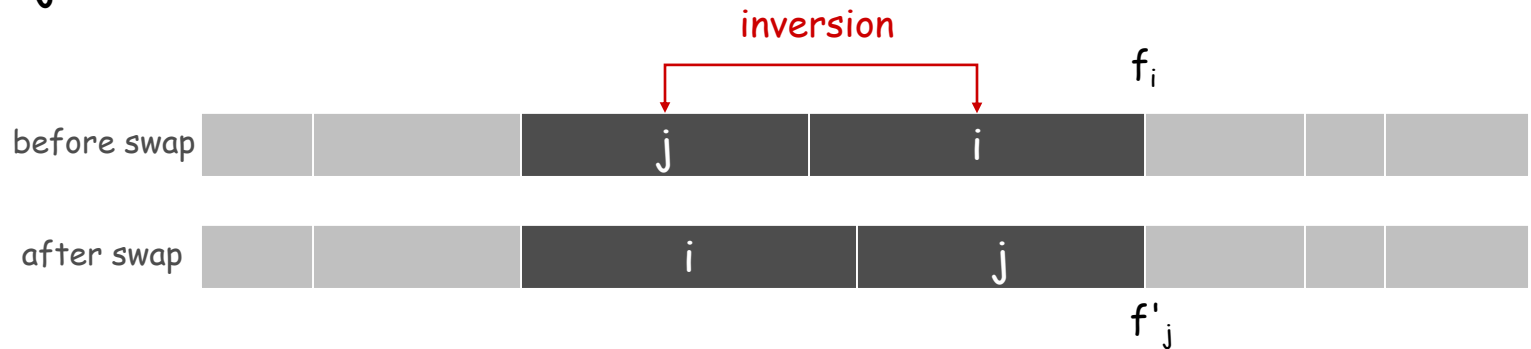Observation 4. If an idle-free schedule has an inversion, then it has an adjacent inversion.

Pf.
- Let $i$–$j$ be a closest inversion.
- Let $k$ be element immediately to the right of $j$.
- Case 1. $[\,j > k\,]$ Then $j$–$k$ is an adjacent inversion.
- Case 2. $[\,j < k\,]$ Then $i$–$k$ is a closer inversion since $i < j < k$.

| | $j$ | $k$ | | | $i$ | | |
|---|---|---|---|---|---|---|---|

# Minimizing Lateness: Inversions

Def. Given a schedule S, an inversion is a pair of jobs i and j such that:
i < j but j scheduled before i.

inversion

$f_i$

before swap

| | | j | i | | | |

after swap

| | | i | j | | | |

$f'_j$

Claim. Swapping two consecutive, inverted jobs reduces the number of inversions by one and does not increase the max lateness.

Pf. Let $\ell$ be the lateness before the swap, and let $\ell$ ' be it afterwards.

- $\ell'_k = \ell_k$ for all $k \neq i, j$
- $\ell'_i \leq \ell_i$
- If job j is late:

$$
\begin{aligned}
\ell'_j &= f'_j - d_j && \text{(definition)} \\
&= f_i - d_j && (j \text{ finishes at time} f_i) \\
&\leq f_i - d_i && (i < j) \\
&\leq \ell_i && \text{(definition)}
\end{aligned}
$$

# Minimizing Lateness:
# Analysis of Greedy Algorithm

**Theorem.** The earliest-deadline-first schedule S is optimal.

**Pf.** [by contradiction]

Define $S*$ to be an optimal schedule with the fewest inversions.

- Can assume $S*$ has no idle time. &larr; Observation 1
- Case 1. [ $S*$ has no inversions ] Then $S = S*$. &larr; Observation 3
- Case 2. [ $S*$ has an inversion ]
  - let $i$–$j$ be an adjacent inversion &larr; Observation 4
  - exchanging jobs $i$ and $j$ decreases the number of inversions by 1 without increasing the max lateness &larr; key claim
  - contradicts "fewest inversions" part of the definition of $S*$

# Outline

- Interval scheduling

- Interval partitioning

- Scheduling to minimize lateness

- Basic elements of greedy algorithms

- Single-source shortest path: Dijkstra's algorithm

- Minimal spanning trees: Prim, Kruskal

- single-link clustering

# Greedy Analysis Strategies

**Greedy algorithm stays ahead.**  Show that after each step of the greedy algorithm, its solution is at least as good as any other algorithm's.

- Interval scheduling

**Structural.**  Discover a simple "structural" bound asserting that every possible solution must have a certain value. Then show that your algorithm always achieves this bound.

- Interval partitioning

**Exchange argument.**  Gradually transform any solution to the one found by the greedy algorithm without hurting its quality.

- Minimizing lateness

**Remark.** Greedy algorithms do not always give an optimal solution but can produce a solution that is guaranteed to be close to optimal.

# Basic Elements of Greedy Algorithms

Usually, if the given problem has the following ingredients, then it is more likely to develop a greedy algorithm for it.

Greedy-choice property. A locally optimal choice is globally optimal.

We can assemble a globally optimal solution by making locally optimal (greedy) choices. That is, when we make the choice that looks best in the current problem, without considering results from subproblems.

Optimal substructure. A problem exhibits optimal substructure if an optimal solution to the problem contains within it optimal solutions to subproblems.

# 0-1 Knapsack Problem

Given $n$ items and a "knapsack."

▪ Item $i$ weighs $w_i > 0$ and has value $v_i > 0$.

▪ Knapsack has weight capacity of $W$.

▪ Goal: pack knapsack so as to maximize total value.

Ex. $\{1, 2, 5\}$ has value 35 (and weight 10).

Ex. $\{3, 4\}$ has value 40 (and weight 11).

Ex. $\{3, 5\}$ has value 46 (but exceeds weight limit).

| $i$ | $v_i$ | $w_i$ |
|---|---|---|
| 1 | 1 | 1 |
| 2 | 6 | 2 |
| 3 | 18 | 5 |
| 4 | 22 | 6 |
| 5 | 28 | 7 |

knapsack instance
(weight limit W = 11)

Greedy by value. Repeatedly add item with maximum $v_i$.

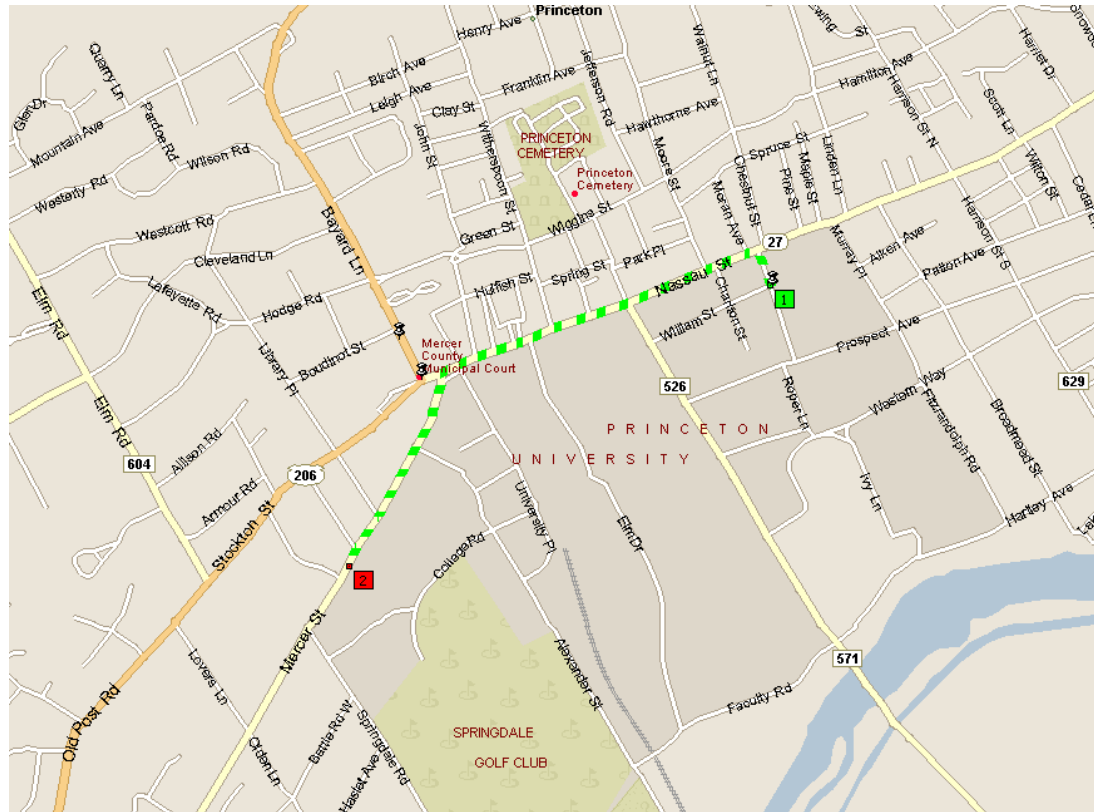Greedy by weight. Repeatedly add item with minimum $w_i$.

Greedy by ratio. Repeatedly add item with maximum ratio $v_i / w_i$.

Observation. None of greedy algorithms is optimal even if it exhibits optimal substructure.

# Outline

- Interval scheduling

- Interval partitioning

- Scheduling to minimize lateness

- Basic elements of greedy algorithms

- Single-source shortest path: Dijkstra's algorithm

- Minimal spanning trees: Prim, Kruskal

- single-link clustering

# Shortest Paths in a Graph



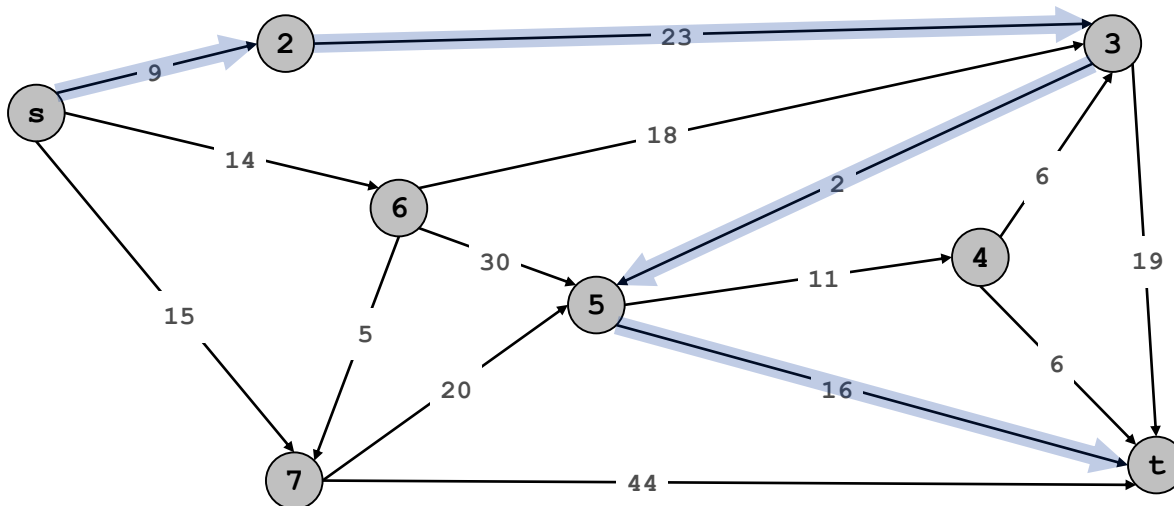shortest path from Princeton CS department to Einstein's house

# Shortest Path Problem

Shortest path network.

- Directed graph G = (V, E).
- Source s, destination t.
- Length $\ell_e$ = length of edge e.

Shortest path problem:  find shortest directed path from s to t.

cost of path = sum of edge costs in path



Cost of path s-2-3-5-t
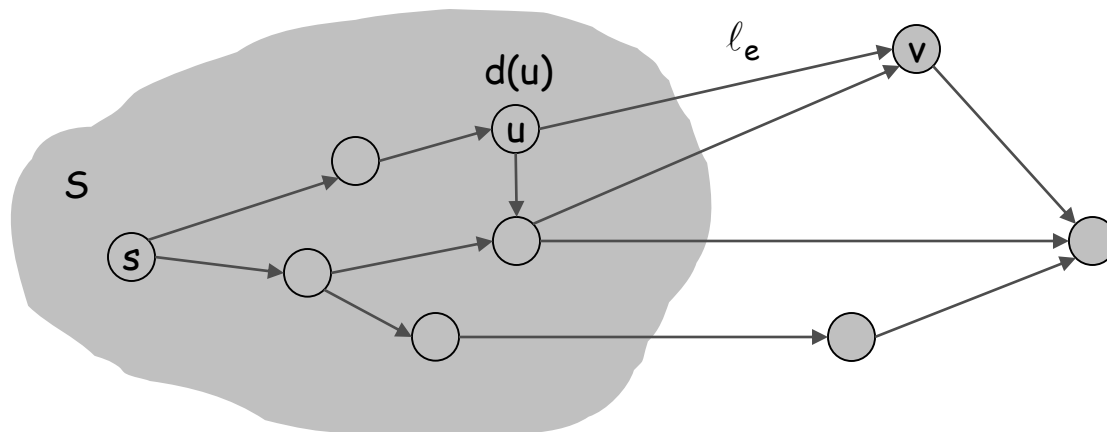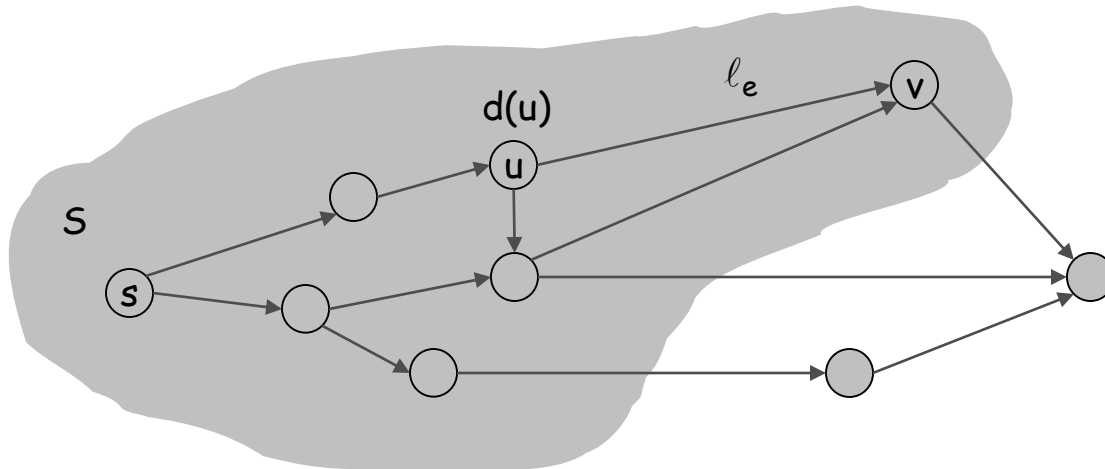    =  9 + 23 + 2 + 16
    = 50.

# Dijkstra's Algorithm

Dijkstra's algorithm.

- Maintain a set of explored nodes S for which we have determined the shortest path distance d(u) from s to u in S.
- Initialize S = { s }, d(s) = 0.
- Repeatedly choose unexplored node v which minimizes

$$\pi(v) = \min_{e = (u,v)\,:\,u \in S} d(u) + \ell_e\,,$$

add v to S, and set d(v) = π(v).

shortest path to some u in explored
part, followed by a single edge (u, v)

# Dijkstra's Algorithm

Dijkstra's algorithm.

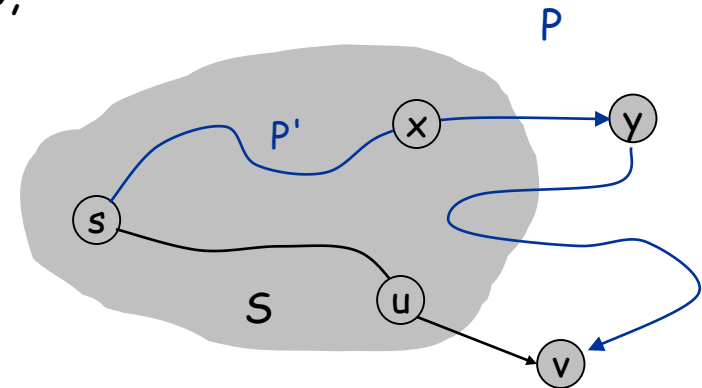- Maintain a set of explored nodes S for which we have determined the shortest path distance d(u) from s to u.
- Initialize S = { s }, d(s) = 0.
- Repeatedly choose unexplored node v which minimizes

$$\pi(v) = \min_{e = (u,v)\,:\,u \in S} d(u) + \ell_e,$$

add v to S, and set d(v) = π(v).

shortest path to some u in explored part, followed by a single edge (u, v)

# Dijkstra's Algorithm:  Proof of Correctness

Invariant.  For each node u $\in$ S, d(u) is the length of the shortest s-u path.
Pf.  (by induction on |S|)
Base case:  |S| = 1 is trivial.
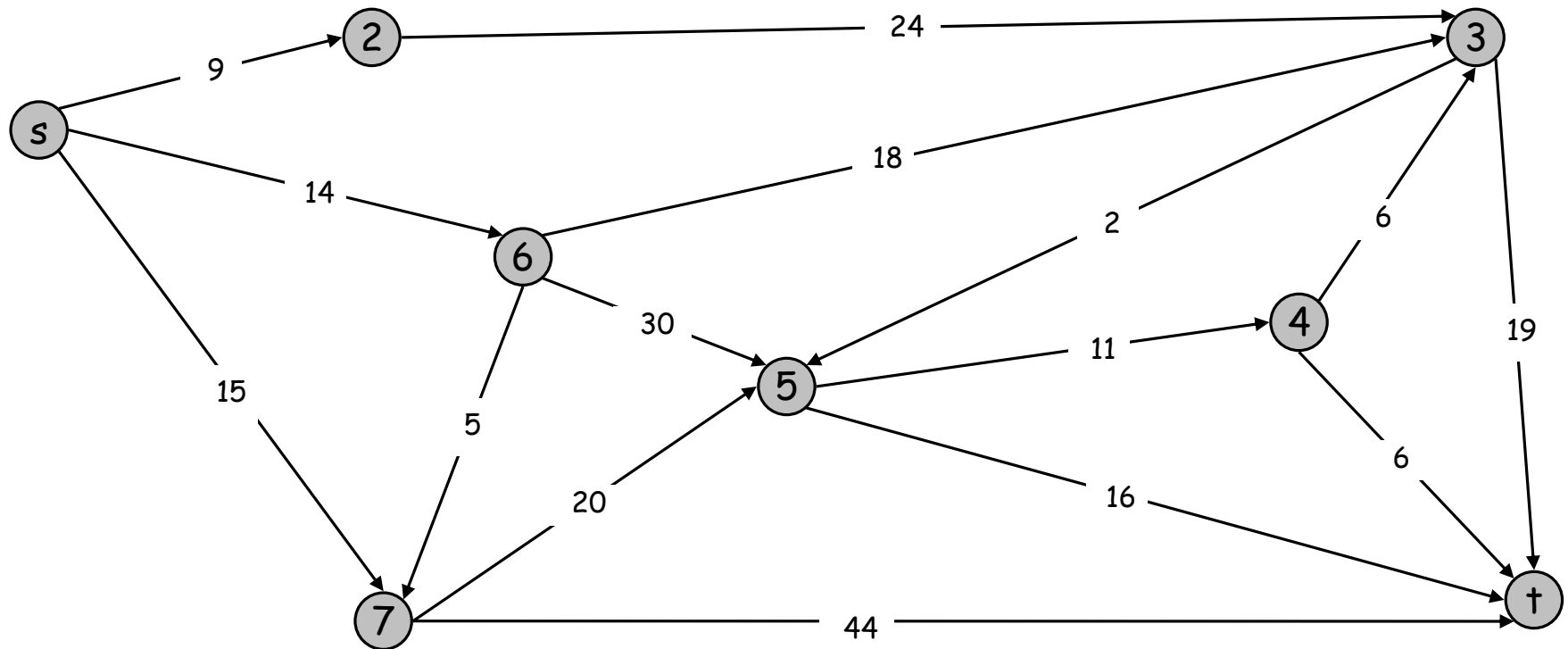Inductive hypothesis:  Assume true for |S| = k $\geq$ 1.

- Let v be next node added to S, and let u-v be the chosen edge.
- The shortest s-u path plus (u, v) is an s-v path of length $\pi(v)$.
- Consider any s-v path P. We'll see that it's no shorter than $\pi(v)$.
- Let x-y be the first edge in P that leaves S, and let P' be the subpath to x.
- P is already too long as soon as it leaves S.



$$\ell (P) \geq \ell (P') + \ell (x,y) \geq d(x) + \ell (x, y) \geq \pi(y) \geq \pi(v)$$

nonnegative weights        inductive hypothesis        defn of $\pi(y)$        Dijkstra chose v instead of y
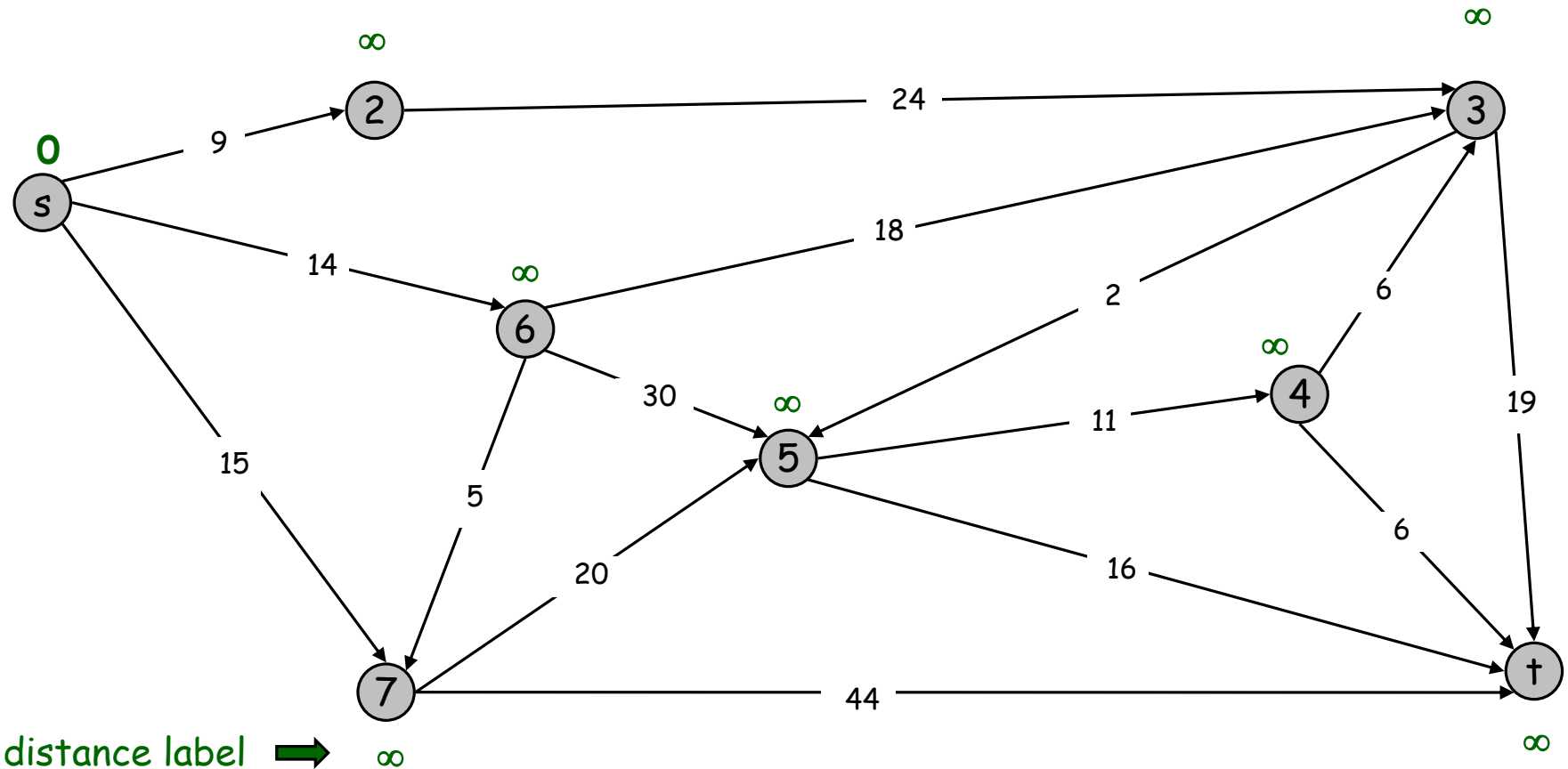
# Dijkstra's Shortest Path Algorithm

Find shortest path from s to t.

# Dijkstra's Shortest Path Algorithm

S = { }
PQ = { s, 2, 3, 4, 5, 6, 7, † }
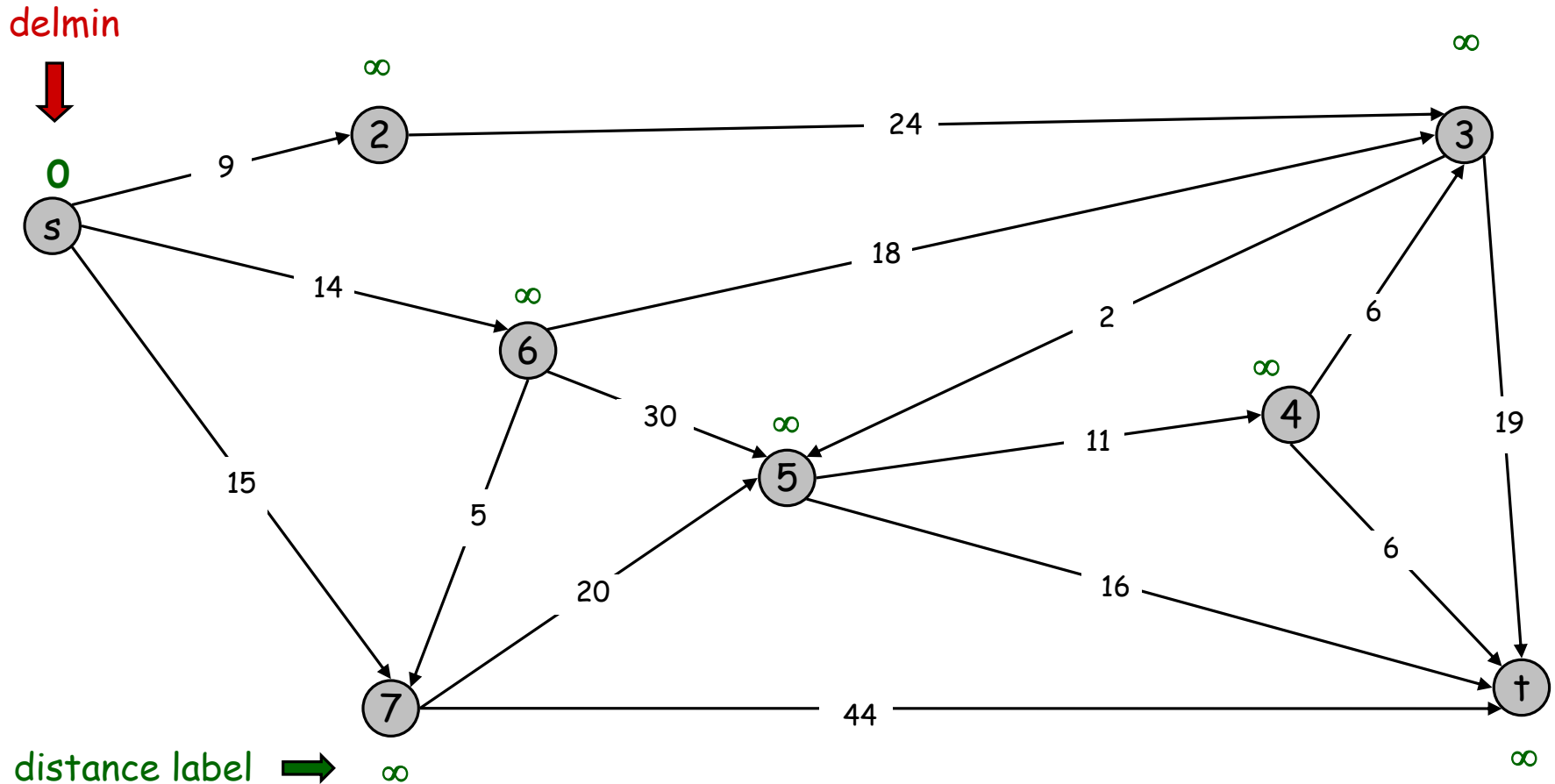


distance label ➡

# Dijkstra's Shortest Path Algorithm

S = { }
PQ = { s, 2, 3, 4, 5, 6, 7, † }



delmin

∞

∞

0

2 — 24 — 3

9

s

18

14

∞

6

2

6

∞

4

30

∞

11

15

5

5

6

20

16

19

7 — 44 — †

distance label ➡ ∞

∞

# Dijkstra's Shortest Path Algorithm

S = { s }
PQ = { 2, 3, 4, 5, 6, 7, † }



decrease key

∞

✗ 9

0

24

9

18

14

✗ 14

2

6

∞

30

∞

11

19

6

5

20

16

6

15

44

distance label ➡ ✗ 15

∞

# Dijkstra's Shortest Path Algorithm

S = { s }
PQ = { 2, 3, 4, 5, 6, 7, † }



delmin

✗ 9

∞

2

0

9

s

18

24

3

14

✗ 14

6

2

6

∞

∞

4

30

11

∞

5

19

15

5

20

6

16

7

44

†

distance label ➡ ✗ 15

∞

# Dijkstra's Shortest Path Algorithm

S = { s, 2 }
PQ = { 3, 4, 5, 6, 7, † }

# Dijkstra's Shortest Path Algorithm

S = { s, 2 }
PQ = { 3, 4, 5, 6, 7, † }

decrease key

✗ 33

✗ 9

2

0

s

24

9

14

✗ 14

6

18

2

6

∞

4

15

30

∞

5

11

19

5

20

16

6

7

44

†

✗ 15

∞

# Dijkstra's Shortest Path Algorithm

S = { s, 2 }
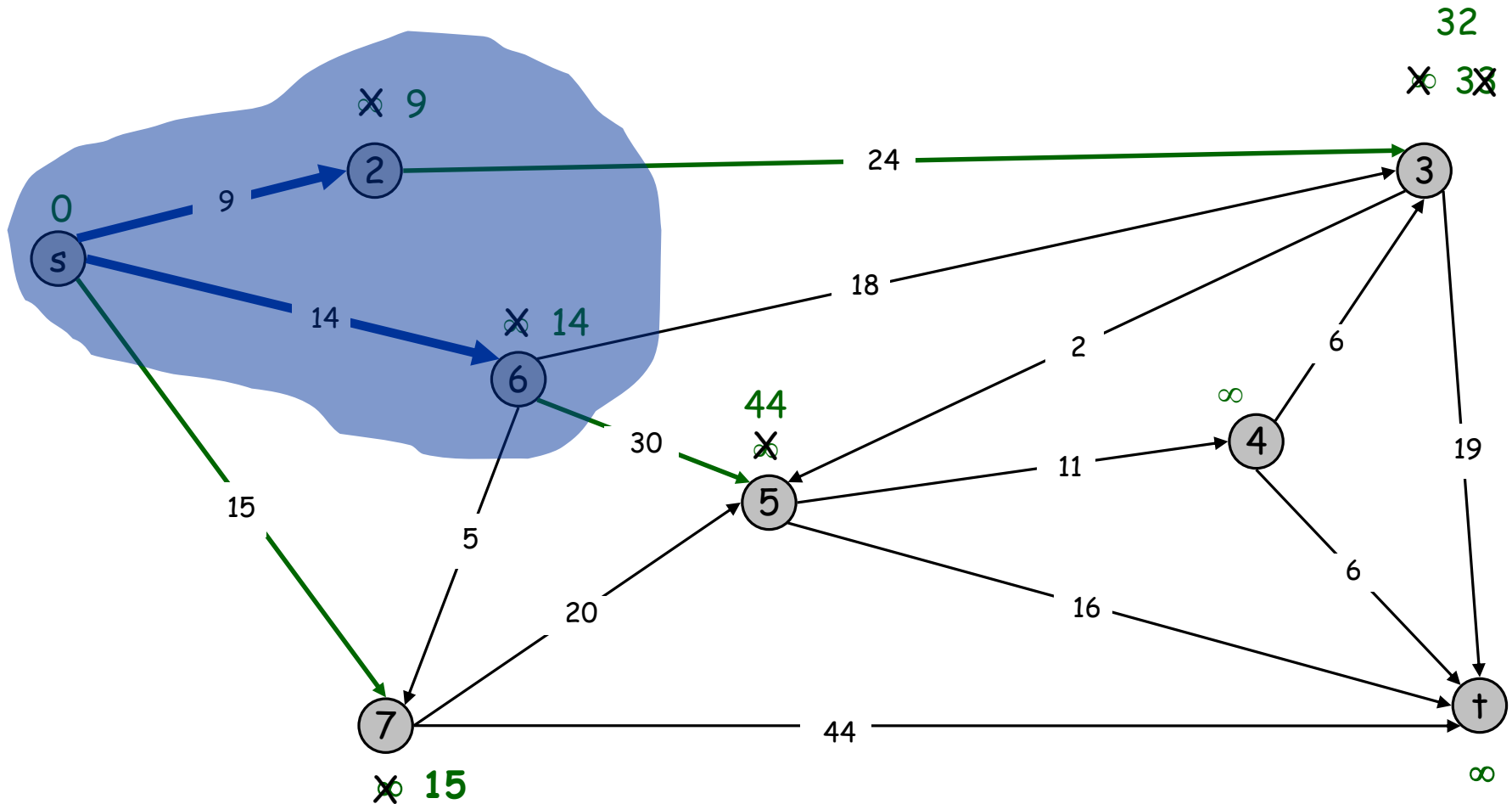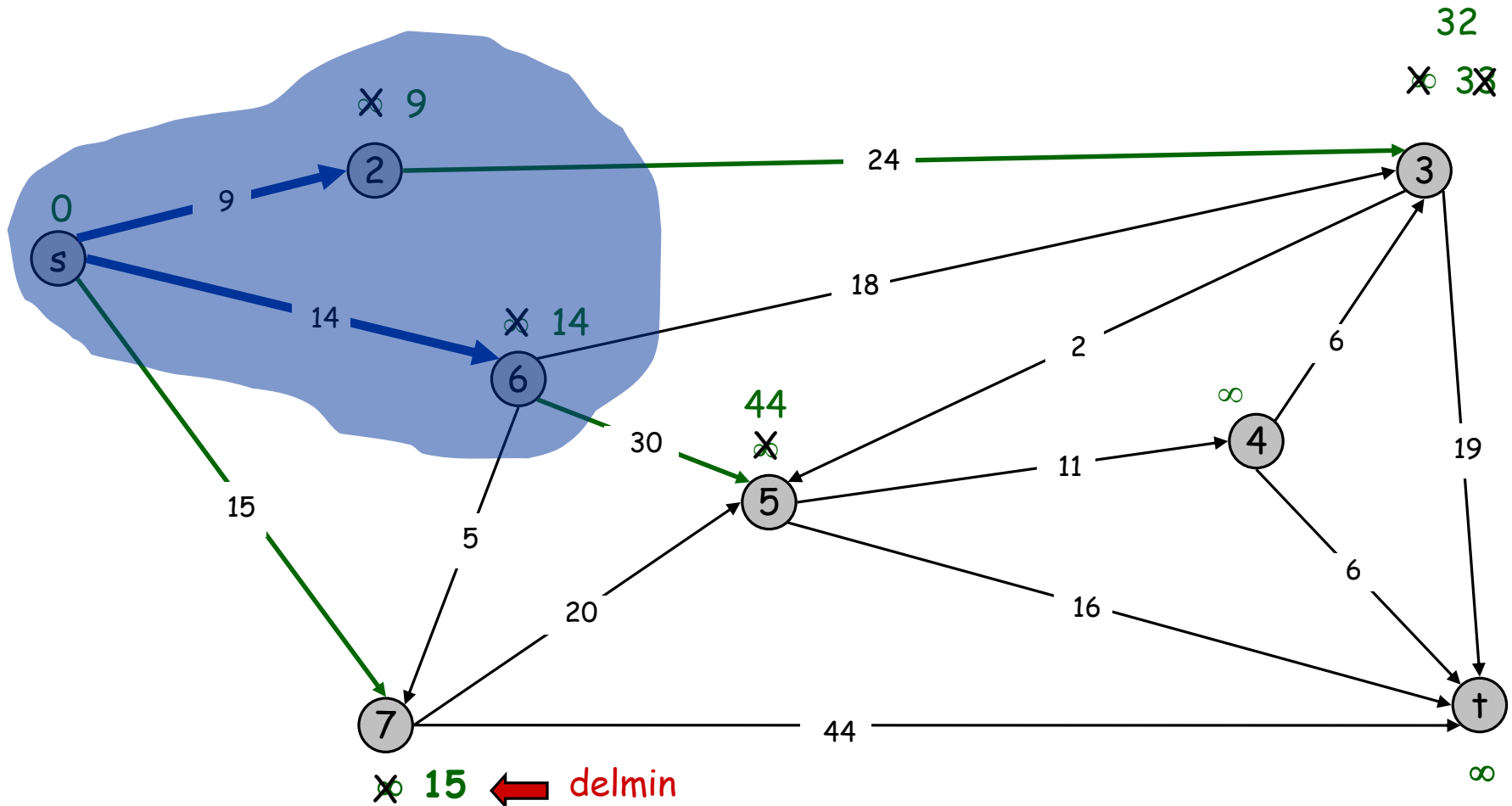PQ = { 3, 4, 5, 6, 7, † }

# Dijkstra's Shortest Path Algorithm

S = { s, 2, 6 }
PQ = { 3, 4, 5, 7, † }

# Dijkstra's Shortest Path Algorithm

S = { s, 2, 6 }
PQ = { 3, 4, 5, 7, † }



32
✗ 3✗

✗ 9

2

0

s

9

24

3

18

14

✗ 14

6

2

6

44
✗

∞

4

30

5

11

19

15

5

20

16

6

7

44

†

✗ 15 ⬅ delmin

∞

# Dijkstra's Shortest Path Algorithm

S = { s, 2, 6, 7 }
PQ = { 3, 4, 5, † }
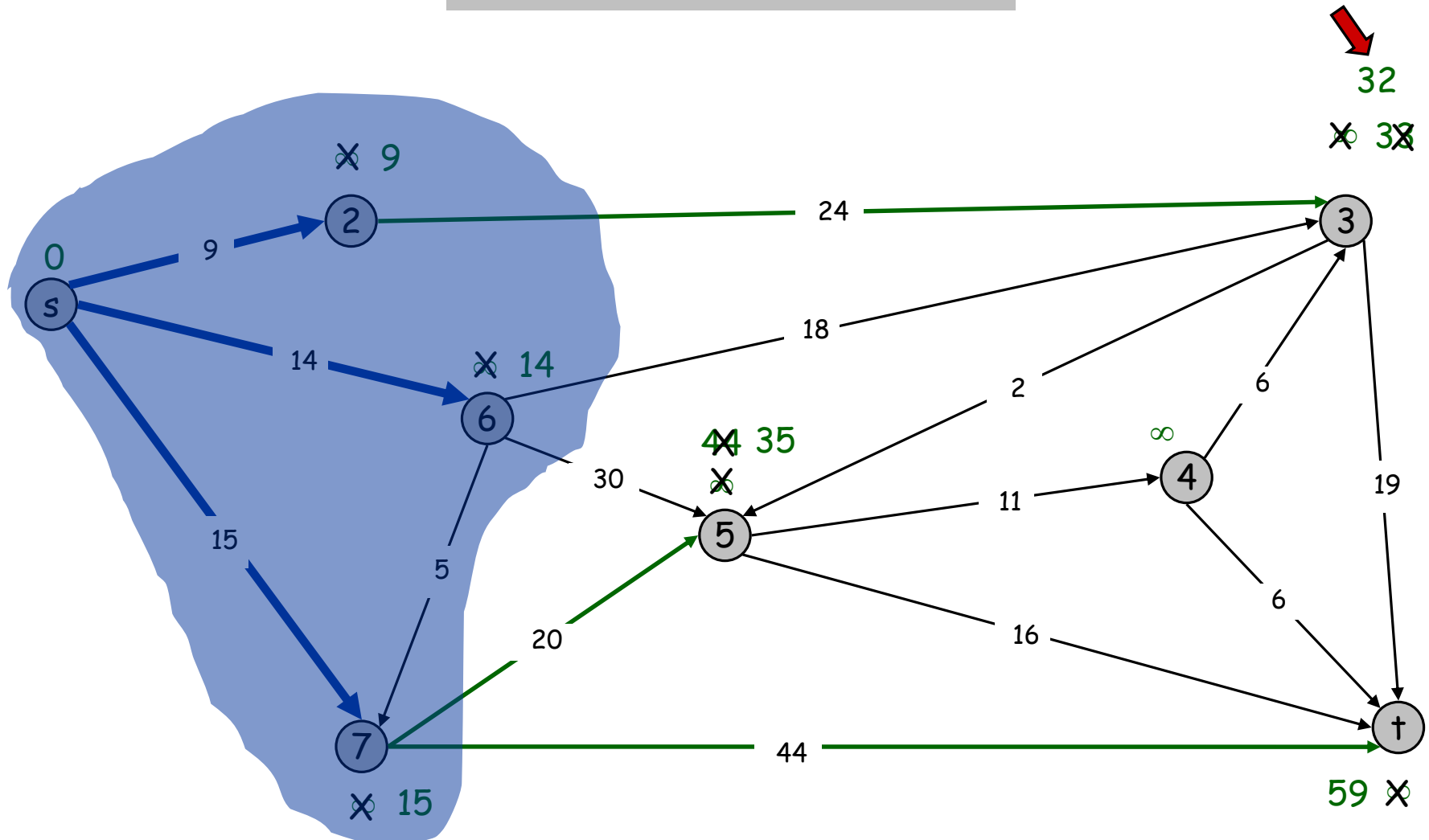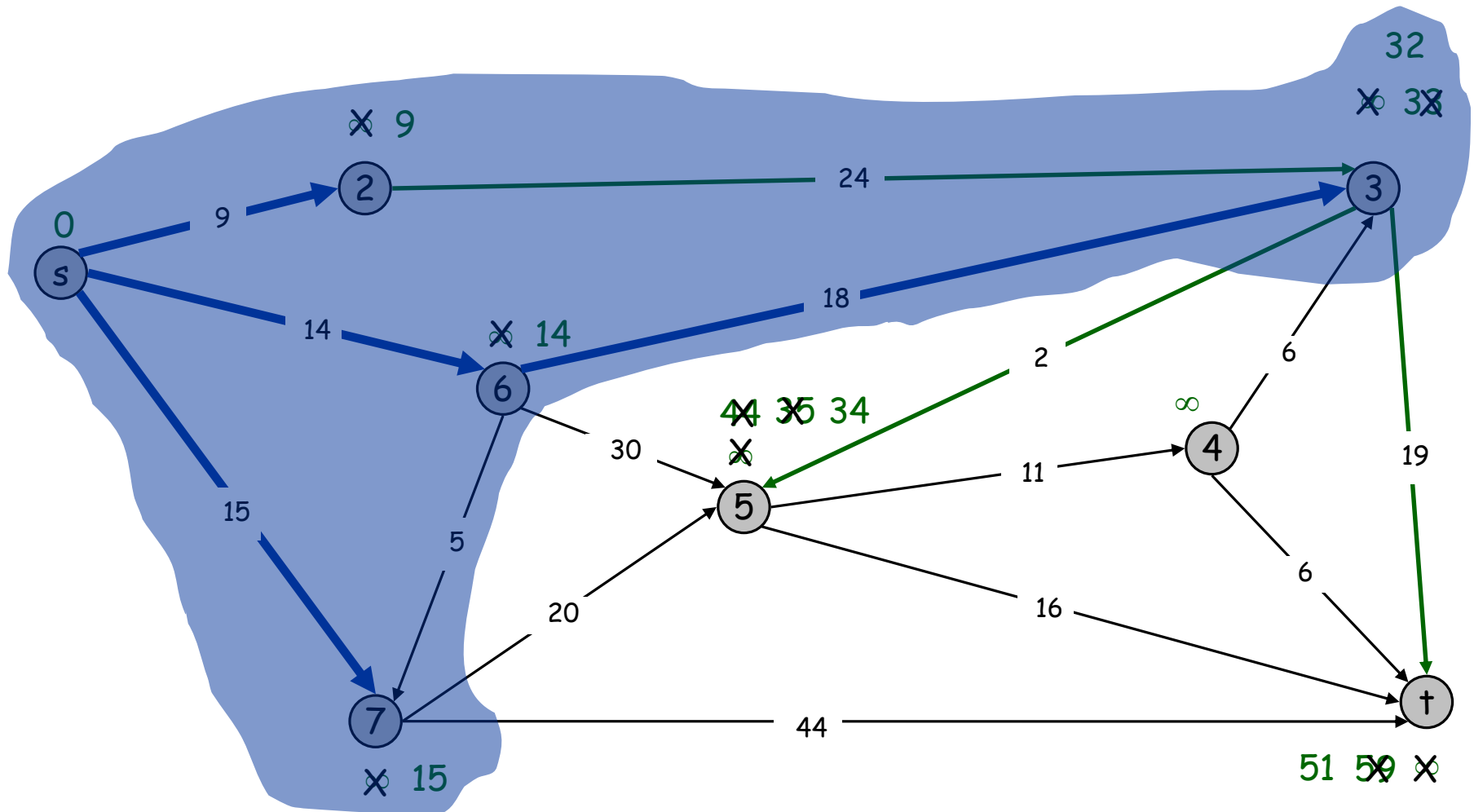
# Dijkstra's Shortest Path Algorithm

S = { s, 2, 6, 7 }
PQ = { 3, 4, 5, † }

delmin

32

✗ 3̶3̶

✗ 9

2 ——— 24 ——— 3

0

9

s

14

18

✗ 14

6

2

6

4̶4̶ 35

∞

30

11

4

5

5

19

15

20

6

16

7

44

†

✗ 15

59 ✗

# Dijkstra's Shortest Path Algorithm

S = { s, 2, 3, 6, 7 }
PQ = { 4, 5, † }



32

⌖ 3⌖

⌖ 9

0

2

24

3

9

s

18

14

6

2

6

⌖ 14

44 ⌖5 34
⌖

∞

19

30

4

5

11

15

5

6

20

16

7

44

†

⌖ 15

51 5⌖ ⌖

# Dijkstra's Shortest Path Algorithm

S = { s, 2, 3, 6, 7 }
PQ = { 4, 5, † }

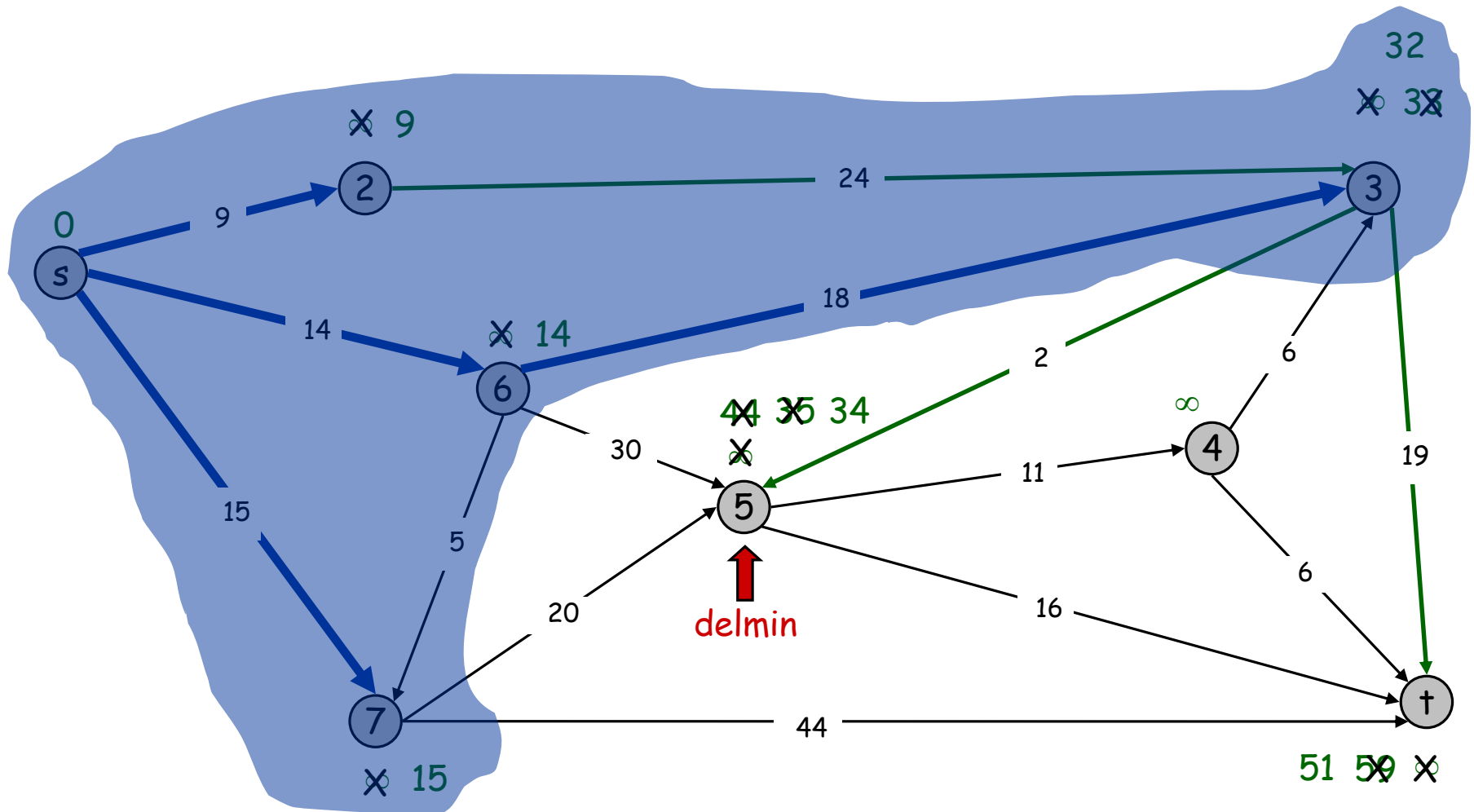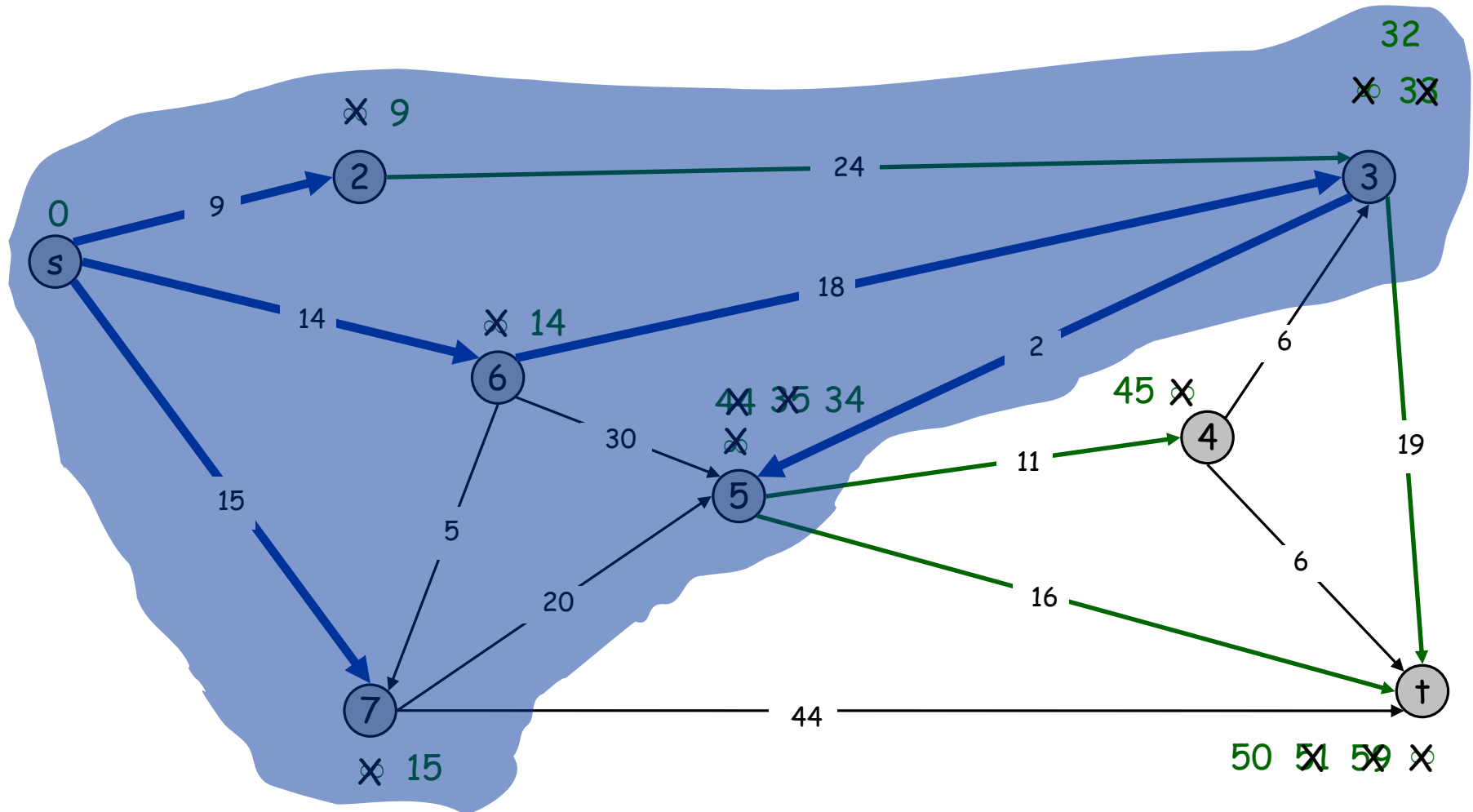# Dijkstra's Shortest Path Algorithm



S = { s, 2, 3, 5, 6, 7 }
PQ = { 4, † }

# Dijkstra's Shortest Path Algorithm

S = { s, 2, 3, 5, 6, 7 }
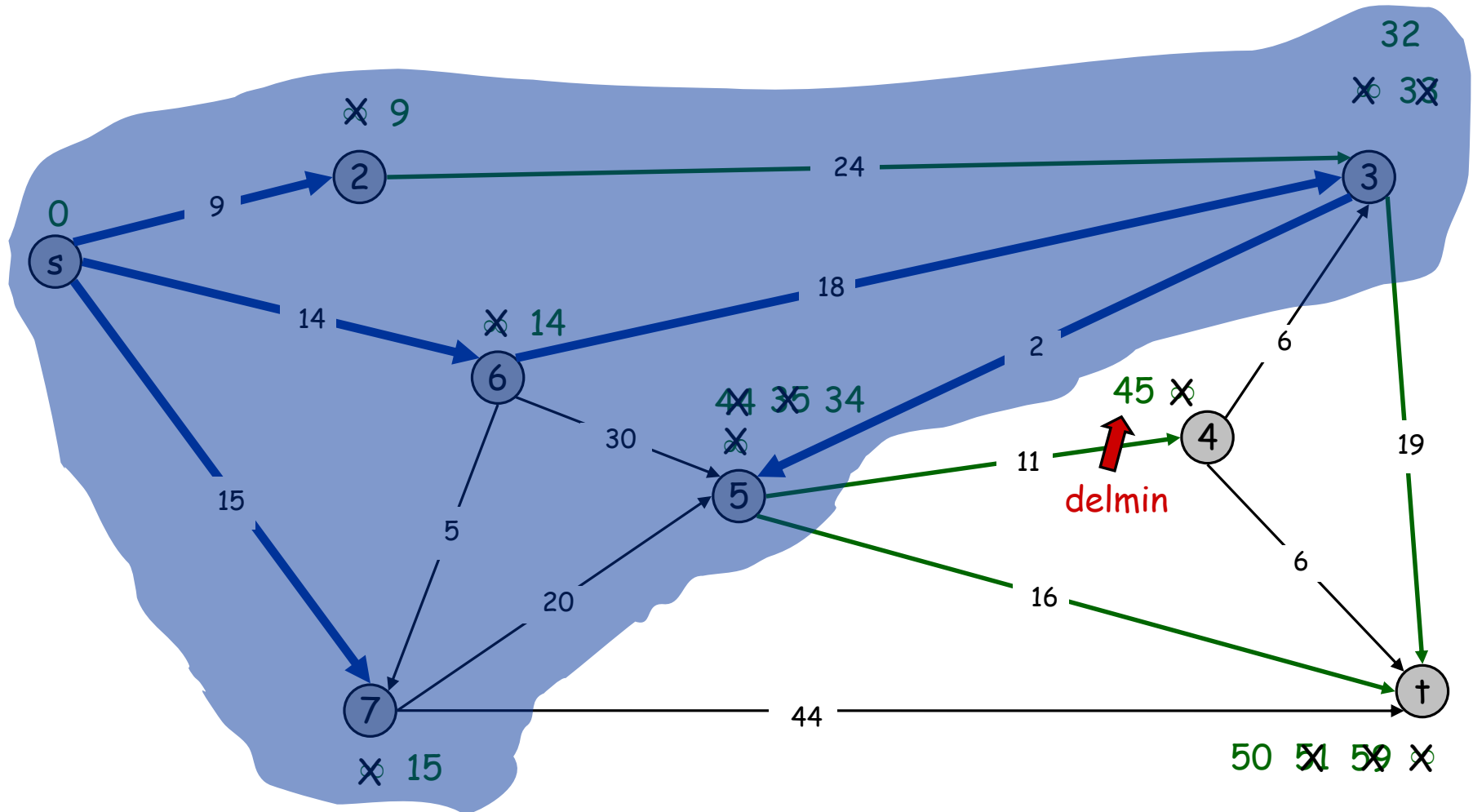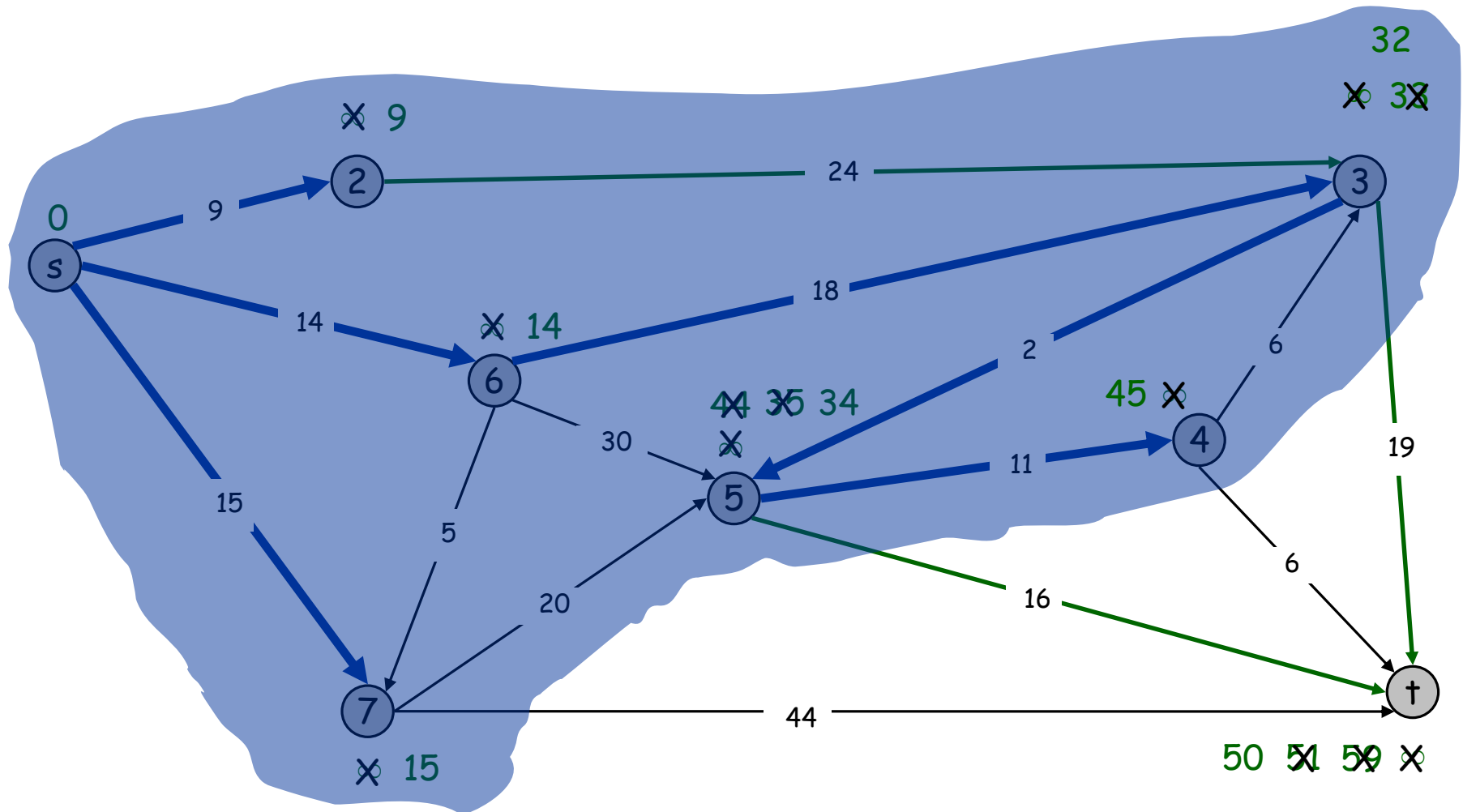PQ = { 4, † }

# Dijkstra's Shortest Path Algorithm

S = { s, 2, 3, 4, 5, 6, 7 }
PQ = { t }



32
X 3X X

X 9

24                                                          3

0
s
            9          2                              18

            14              X 14

15                     6              2

                30         44 X X 34         45 X

                    X

                    5         11              4

                    20              16         6

        7                              19

            44                                    t

        X 15

50 5X 5X X

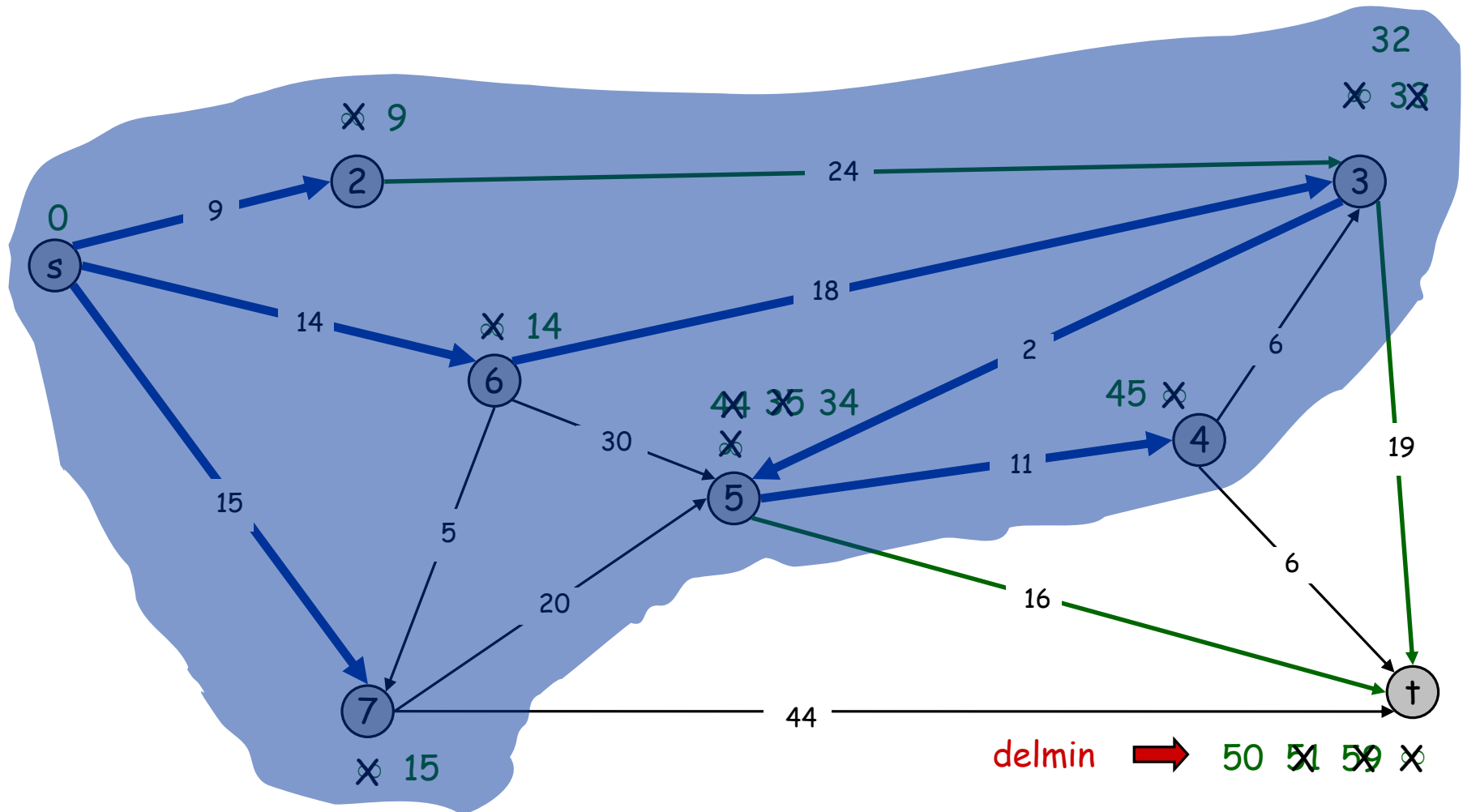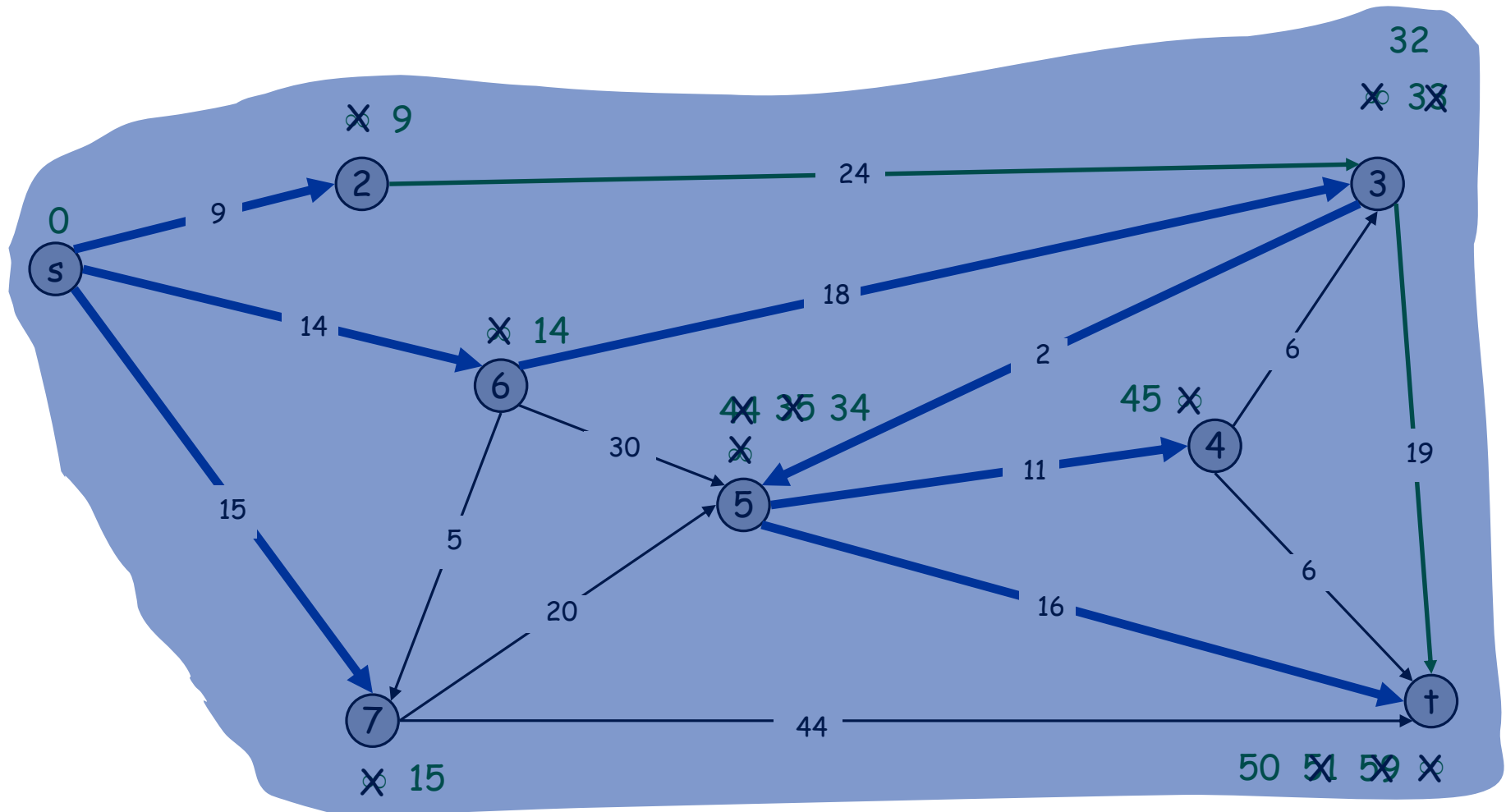# Dijkstra's Shortest Path Algorithm

S = { s, 2, 3, 4, 5, 6, 7 }
PQ = { t }

# Dijkstra's Shortest Path Algorithm

S = { s, 2, 3, 4, 5, 6, 7, † }
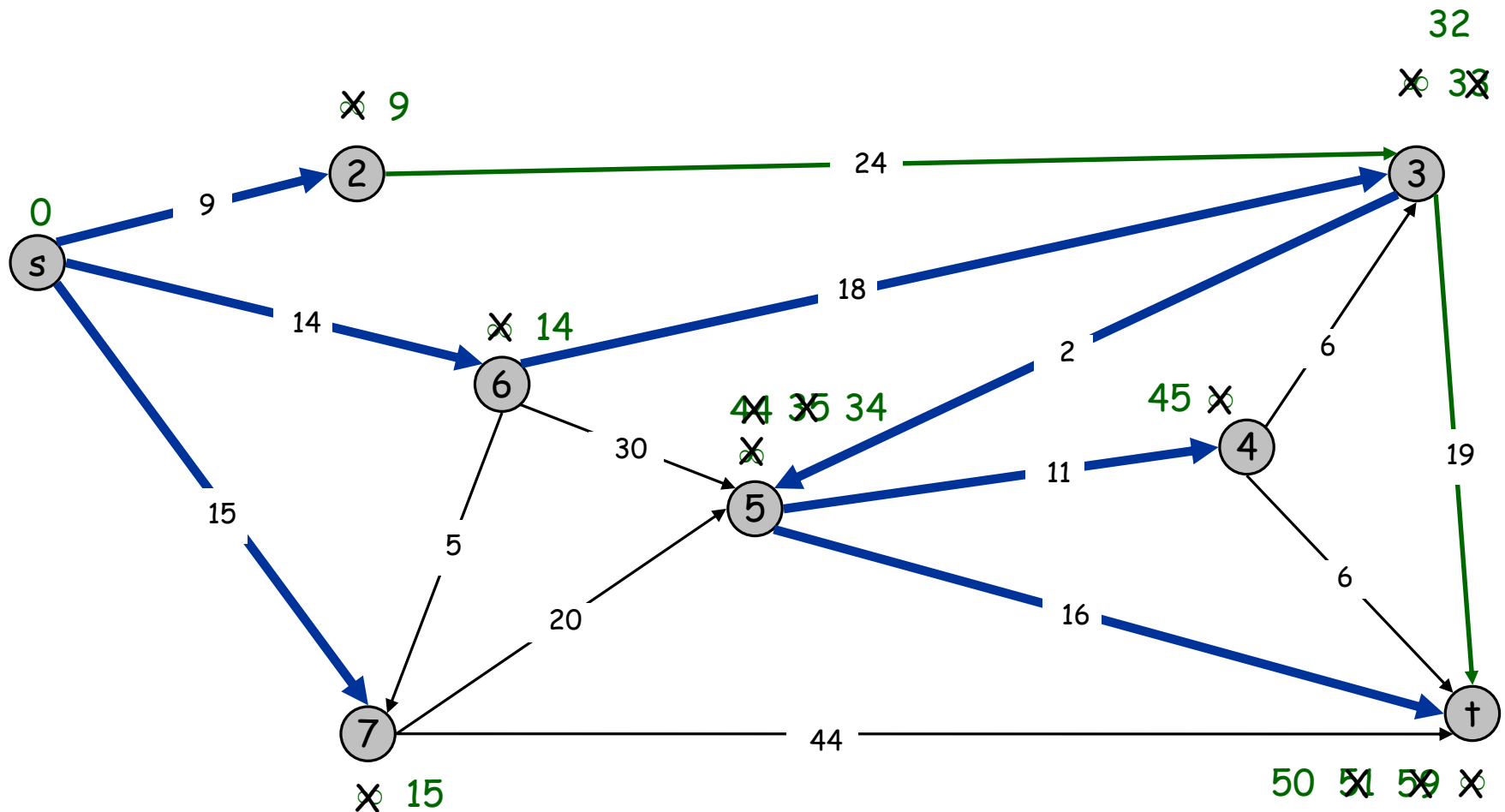PQ = { }

# Dijkstra's Shortest Path Algorithm

S = { s, 2, 3, 4, 5, 6, 7, † }
PQ = { }

# Dijkstra's Algorithm:  Implementation

For each unexplored node, explicitly maintain $\pi(v) = \min\limits_{e=(u,v)\,:\,u\in S} d(u) + \ell_e$ .

- Next node to explore = node with minimum $\pi(v)$.
- When exploring v, for each incident edge e = (v, w), update

  $\pi(w) = \min \{ \ \pi(w), \ \pi(v) + \ell_e \}$.

Efficient implementation.  Maintain a priority queue of unexplored nodes, prioritized by $\pi(v)$.

| PQ Operation | Dijkstra | Binary heap | Array |
|:---:|:---:|:---:|:---:|
| Insert | n | log n | n |
| ExtractMin | n | log n | n |
| ChangeKey | m | log n | 1 |
| IsEmpty | n | 1 | 1 |
| Total | | m log n | n² |

# Dijkstra's algorithm

$d[s] \leftarrow 0$
**for** each $v \in V - \{s\}$
    **do** $d[v] \leftarrow \infty$
$S \leftarrow \varnothing$
$Q \leftarrow V$     $\triangleright$ $Q$ is a priority queue maintaining $V - S$
**while** $Q \neq \varnothing$
    **do** $u \leftarrow$ Extract-Min$(Q)$
        $S \leftarrow S \cup \{u\}$
      **for** each $v \in Adj[u]$
          **do if** $d[v] > d[u] + w(u, v)$
          **then** $d[v] \leftarrow d[u] + w(u, v)$

# Dijkstra's algorithm

$d[s] \leftarrow 0$
**for** each $v \in V - \{s\}$
    **do** $d[v] \leftarrow \infty$
$S \leftarrow \varnothing$
$Q \leftarrow V$     $\triangleright Q$ is a priority queue maintaining $V - S$
**while** $Q \neq \varnothing$
    **do** $u \leftarrow$ EXTRACT-MIN($Q$)
        $S \leftarrow S \cup \{u\}$
      **for** each $v \in Adj[u]$
          **do if** $d[v] > d[u] + w(u, v)$
          **then** $d[v] \leftarrow d[u] + w(u, v)$

Implicit DECREASE-KEY

# Analysis of Dijkstra

**while** $Q \neq \varnothing$

   **do** $u \leftarrow$ EXTRACT-MIN$(Q)$

   $S \leftarrow S \cup \{u\}$

   **for** each $v \in Adj[u]$

    **do if** $d[v] > d[u] + w(u, v)$

     **then** $d[v] \leftarrow d[u] + w(u, v)$

# Analysis of Dijkstra

$|V|$ times $\left\{\rule{0pt}{90pt}\right.$

**while** $Q \neq \varnothing$

  **do** $u \leftarrow$ EXTRACT-MIN$(Q)$

  $S \leftarrow S \cup \{u\}$

  **for** each $v \in Adj[u]$

   **do if** $d[v] > d[u] + w(u, v)$

    **then** $d[v] \leftarrow d[u] + w(u, v)$

# Analysis of Dijkstra

$|V|$
times

$degree(u)$
times

$$\textbf{while } Q \neq \varnothing$$
$$\textbf{do } u \leftarrow \text{E\small{XTRACT}-M\small{IN}}(Q)$$
$$S \leftarrow S \cup \{u\}$$
$$\textbf{for } \text{each } v \in Adj[u]$$
$$\textbf{do if } d[v] > d[u] + w(u, v)$$
$$\textbf{then } d[v] \leftarrow d[u] + w(u, v)$$

# Analysis of Dijkstra

$|V|$ times
$degree(u)$ times

$$\textbf{while } Q \neq \varnothing$$
$$\textbf{do } u \leftarrow \text{Extract-Min}(Q)$$
$$S \leftarrow S \cup \{u\}$$
$$\textbf{for } \text{each } v \in Adj[u]$$
$$\textbf{do if } d[v] > d[u] + w(u, v)$$
$$\textbf{then } d[v] \leftarrow d[u] + w(u, v)$$

$\Theta(E)$ implicit Decrease-Key's.

# Analysis of Dijkstra

$|V|$ times

$degree(u)$ times

**while** $Q \neq \varnothing$
  **do** $u \leftarrow \text{EXTRACT-MIN}(Q)$
  $S \leftarrow S \cup \{u\}$
  **for** each $v \in Adj[u]$
    **do if** $d[v] > d[u] + w(u, v)$
      **then** $d[v] \leftarrow d[u] + w(u, v)$

$\Theta(E)$ implicit DECREASE-KEY

$$\text{Time} = \Theta(V \cdot T_{\text{EXTRACT-MIN}} + E \cdot T_{\text{DECREASE-KEY}})$$

**Note:** Same formula as in the analysis of Prim's minimum spanning tree algorithm.

# Analysis of Dijkstra (continued)

$$\text{Time} = \Theta(V)\, T_{\text{EXTRACT-MIN}} + \Theta(E)\, T_{\text{DECREASE-KEY}}$$

| $Q$ | $T_{\text{EXTRACT-MIN}}$ | $T_{\text{DECREASE-KEY}}$ | Total |
| --- | --- | --- | --- |

# Analysis of Dijkstra (continued)

$$\text{Time} = \Theta(V)\, T_{\text{EXTRACT-MIN}} + \Theta(E)\, T_{\text{DECREASE-KEY}}$$

| $Q$ | $T_{\text{EXTRACT-MIN}}$ | $T_{\text{DECREASE-KEY}}$ | Total |
|---|---|---|---|
| array | $O(V)$ | $O(1)$ | $O(V^2)$ |

# Analysis of Dijkstra (continued)

$$\text{Time} = \Theta(V)\, T_{\text{EXTRACT-MIN}} + \Theta(E)\, T_{\text{DECREASE-KEY}}$$

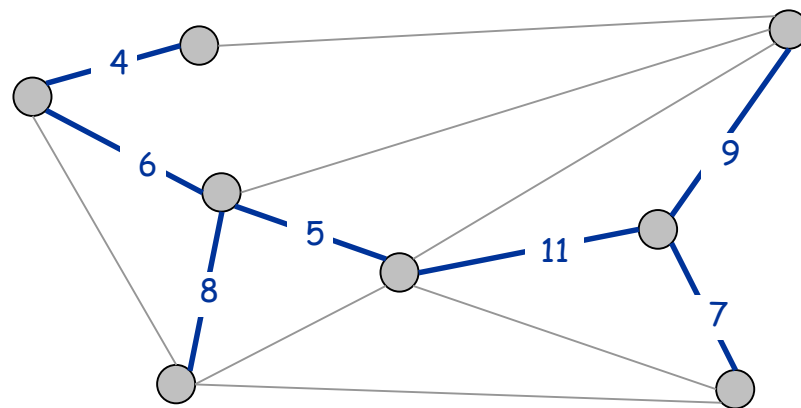| $Q$ | $T_{\text{EXTRACT-MIN}}$ | $T_{\text{DECREASE-KEY}}$ | Total |
|---|---|---|---|
| array | $O(V)$ | $O(1)$ | $O(V^2)$ |
| binary heap | $O(\lg V)$ | $O(\lg V)$ | $O(E \lg V)$ |

# Outline

- Interval scheduling

- Interval partitioning

- Scheduling to minimize lateness

- Basic elements of greedy algorithms

- Single-source shortest path: Dijkstra's algorithm

- Minimal spanning trees: Prim, Kruskal

- single-link clustering

# Minimum Spanning Tree

Minimum spanning tree.  Given a connected graph $G = (V, E)$ with real-valued edge weights $c_e$, an MST is a subset of the edges $T \subseteq E$ such that $T$ is a spanning tree whose sum of edge weights is minimized.



$G = (V, E)$

$T, \ \sum_{e \in T} c_e = 50$

Cayley's Theorem.  There are $n^{n-2}$ spanning trees of $K_n$.
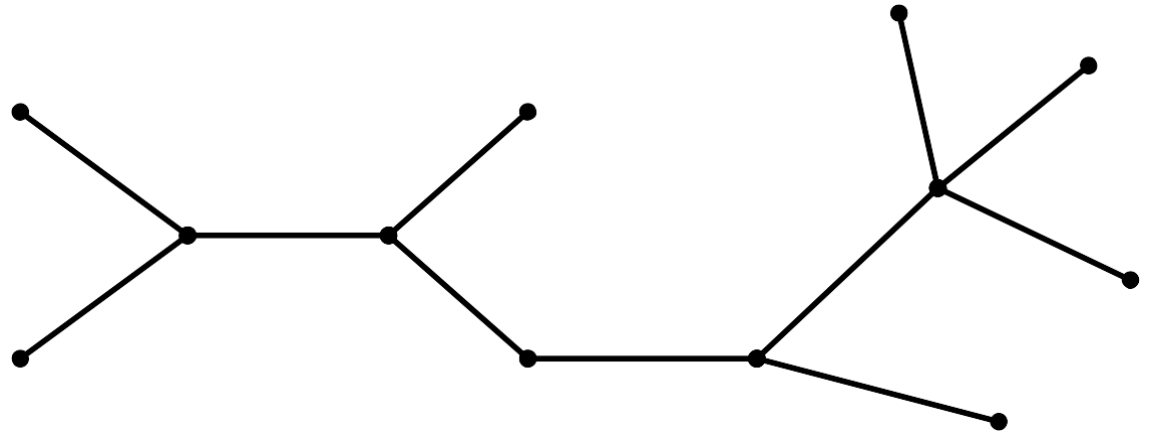
↑

can't solve by brute force

# Applications

MST is a fundamental problem with diverse applications.

- Network design.
  - telephone, electrical, hydraulic, TV cable, computer, road

- Approximation algorithms for NP-hard problems.
  - traveling salesperson problem, Steiner tree

- Indirect applications.
  - max bottleneck paths
  - LDPC codes for error correction
  - image registration with Renyi entropy
  - learning salient features for real-time face verification
  - reducing data storage in sequencing amino acids in a protein
  - model locality of particle interactions in turbulent fluid flows
  - autoconfig protocol for Ethernet bridging to avoid cycles in a network
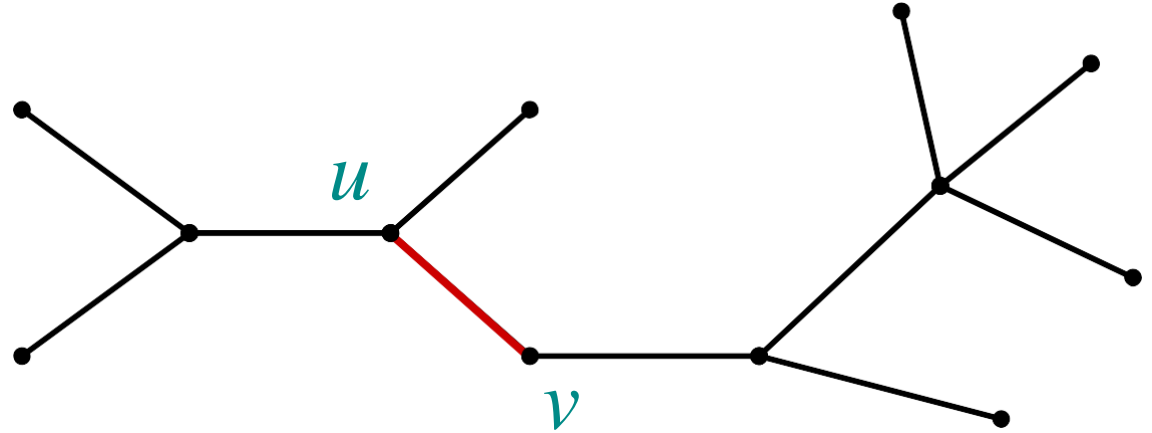
- Cluster analysis.

# Optimal substructure

MST $T$:

(Other edges of $G$
are not shown.)

# Optimal substructure

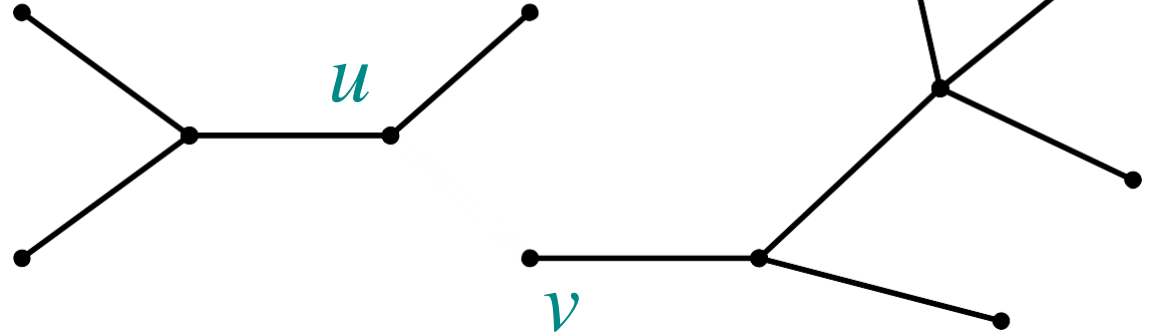MST $T$:

(Other edges of $G$ are not shown.)

Remove any edge $(u, v) \in T$.

# Optimal substructure

MST $T$:

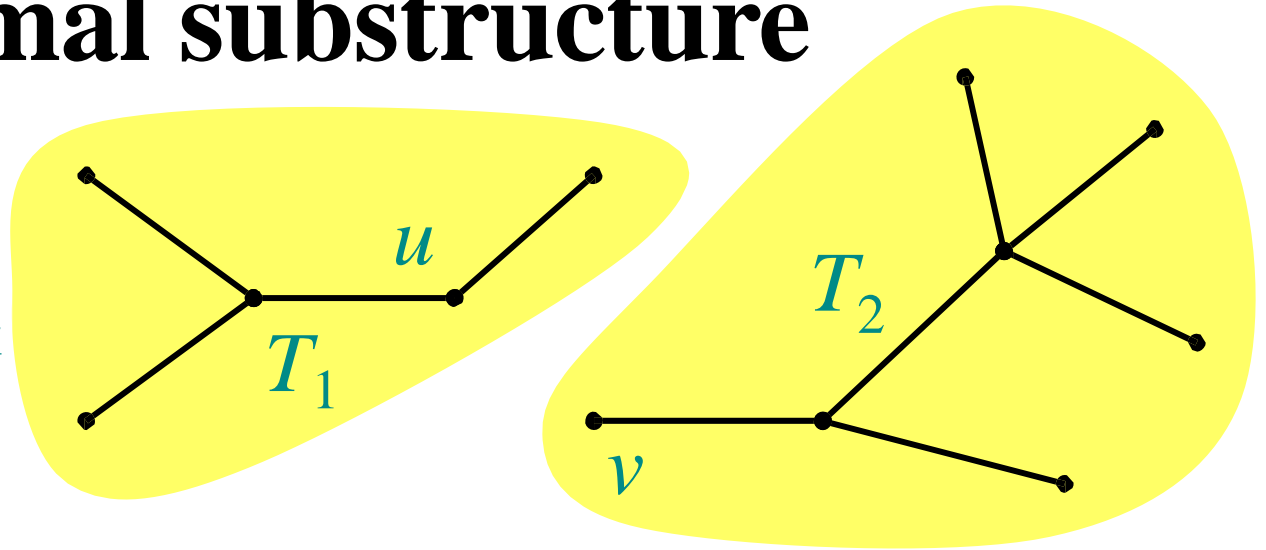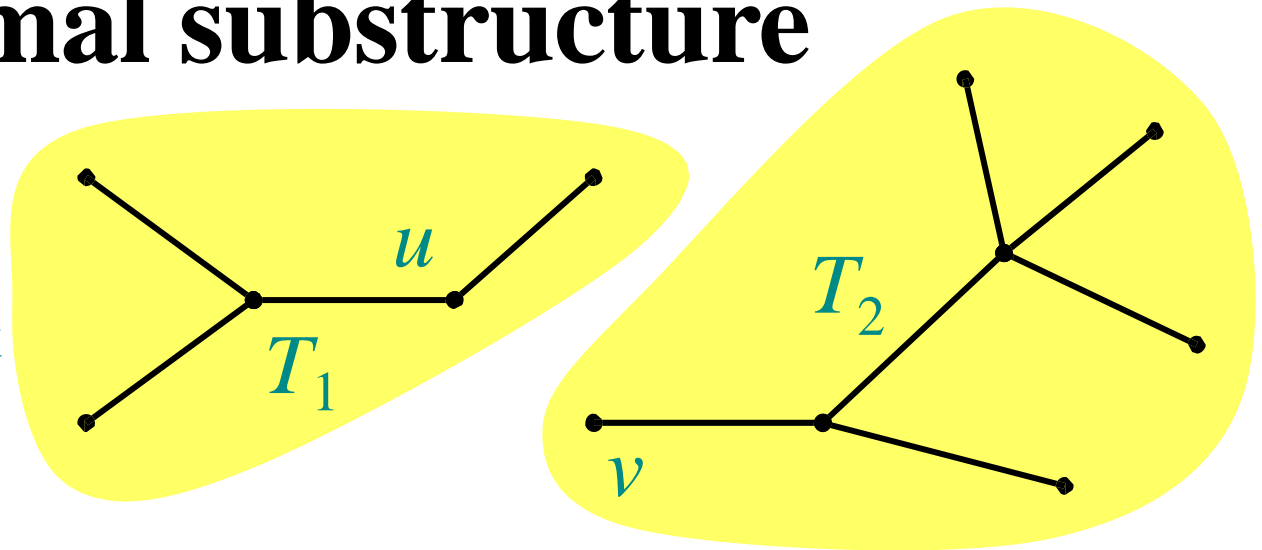(Other edges of $G$
are not shown.)

$u$

$v$

Remove any edge $(u, v) \in T$.

# Optimal substructure

MST $T$:

(Other edges of $G$ are not shown.)



Remove any edge $(u, v) \in T$. Then, $T$ is partitioned into two subtrees $T_1$ and $T_2$.

# Optimal substructure



MST $T$:

(Other edges of $G$
 are not shown.)

Remove any edge $(u, v) \in T$. Then, $T$ is partitioned into two subtrees $T_1$ and $T_2$.

**Theorem.** The subtree $T_1$ is an MST of $G_1 = (V_1, E_1)$, the subgraph of $G$ ***induced*** by the vertices of $T_1$:

$V_1$ = vertices of $T_1$,
$E_1$ = { $(x, y) \in E : x, y \in V_1$ }.

Similarly for $T_2$.

# Proof of optimal substructure

*Proof.* Cut and paste:
$$w(T) = w(u, v) + w(T_1) + w(T_2).$$
If $T_1{}'$ were a lower-weight spanning tree than $T_1$ for
$G_1$, then $T\,' = \{(u, v)\} \cup T_1{}' \cup T_2$ would be a lower-weight spanning tree than $T$ for $G$. ☐

# Greedy Algorithms

**Kruskal's algorithm.**  Start with T = $\phi$. Consider edges in ascending order of cost. Insert edge e in T unless doing so would create a cycle.
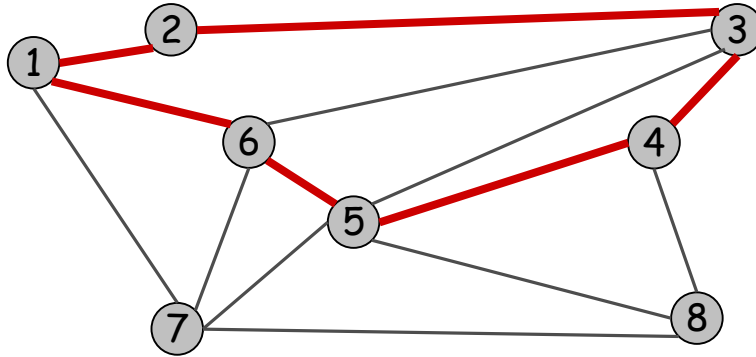
**Reverse-Delete algorithm.**  Start with T = E.  Consider edges in descending order of cost. Delete edge e from T unless doing so would disconnect T.

**Prim's algorithm.**  Start with some root node s and greedily grow a tree T from s outward.  At each step, add the cheapest edge e to T that has exactly one endpoint in T.

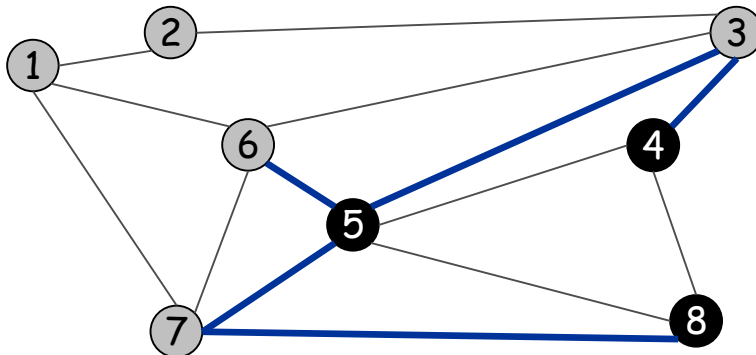**Remark.**  All three algorithms produce an MST.

# Cycles and Cuts

Cycle. Set of edges the form a-b, b-c, c-d, …, y-z, z-a.



Cycle C = 1-2, 2-3, 3-4, 4-5, 5-6, 6-1

Cutset. A cut is a subset of nodes S. The corresponding cutset D is the subset of edges with exactly one endpoint in S.
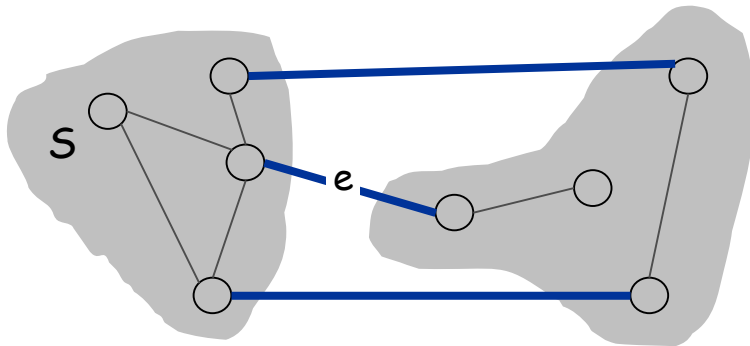


Cut S      = { 4, 5, 8 }
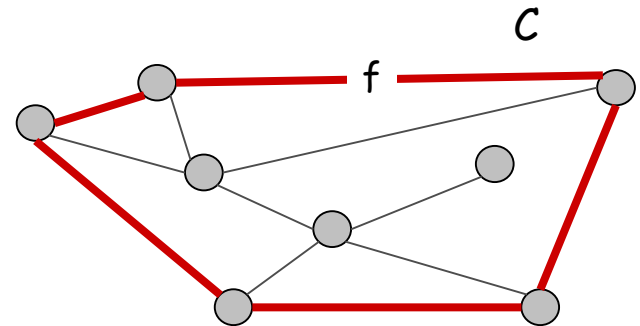Cutset  D = 5-6, 5-7, 3-4, 3-5, 7-8

# Greedy-choice property

Simplifying assumption.  All edge costs $c_e$ are distinct.

Cut property.  Let S be any subset of nodes, and let e be the min cost edge with exactly one endpoint in S.  Then the MST contains e.

Cycle property.  Let C be any cycle, and let f be the max cost edge belonging to C.  Then the MST does not contain f.
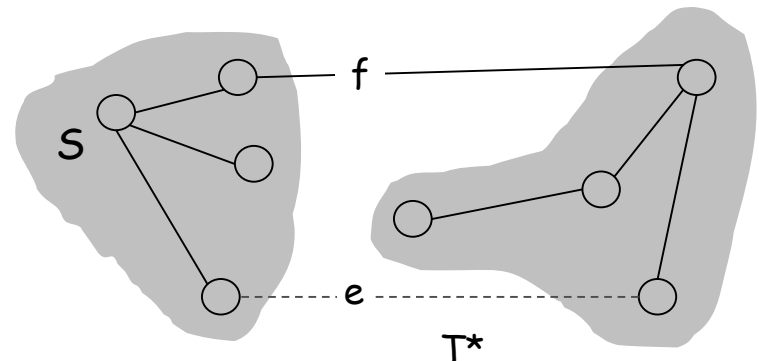


e is in the MST

f is not in the MST

# Greedy Algorithms

Simplifying assumption.  All edge costs $c_e$ are distinct.

Cut property.  Let S be any subset of nodes, and let e be the min cost edge with exactly one endpoint in S. Then the MST T* contains e.

Pf.  (exchange argument)
- Suppose e does not belong to T*, and let's see what happens.
- Adding e to T* creates a cycle C in T*.
- Edge e is both in the cycle C and in the cutset D corresponding to S
  $\Rightarrow$  there exists another edge, say f, that is in both C and D.
- T' = T* $\cup$ { e } - { f } is also a spanning tree.
- Since $c_e < c_f$, cost(T') < cost(T*).
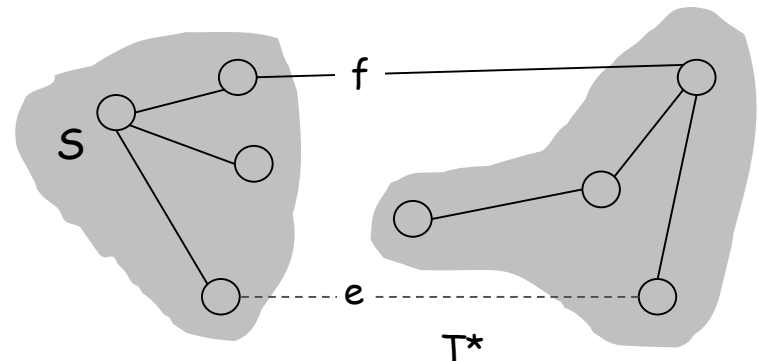- This is a contradiction.  ▪

# Greedy Algorithms
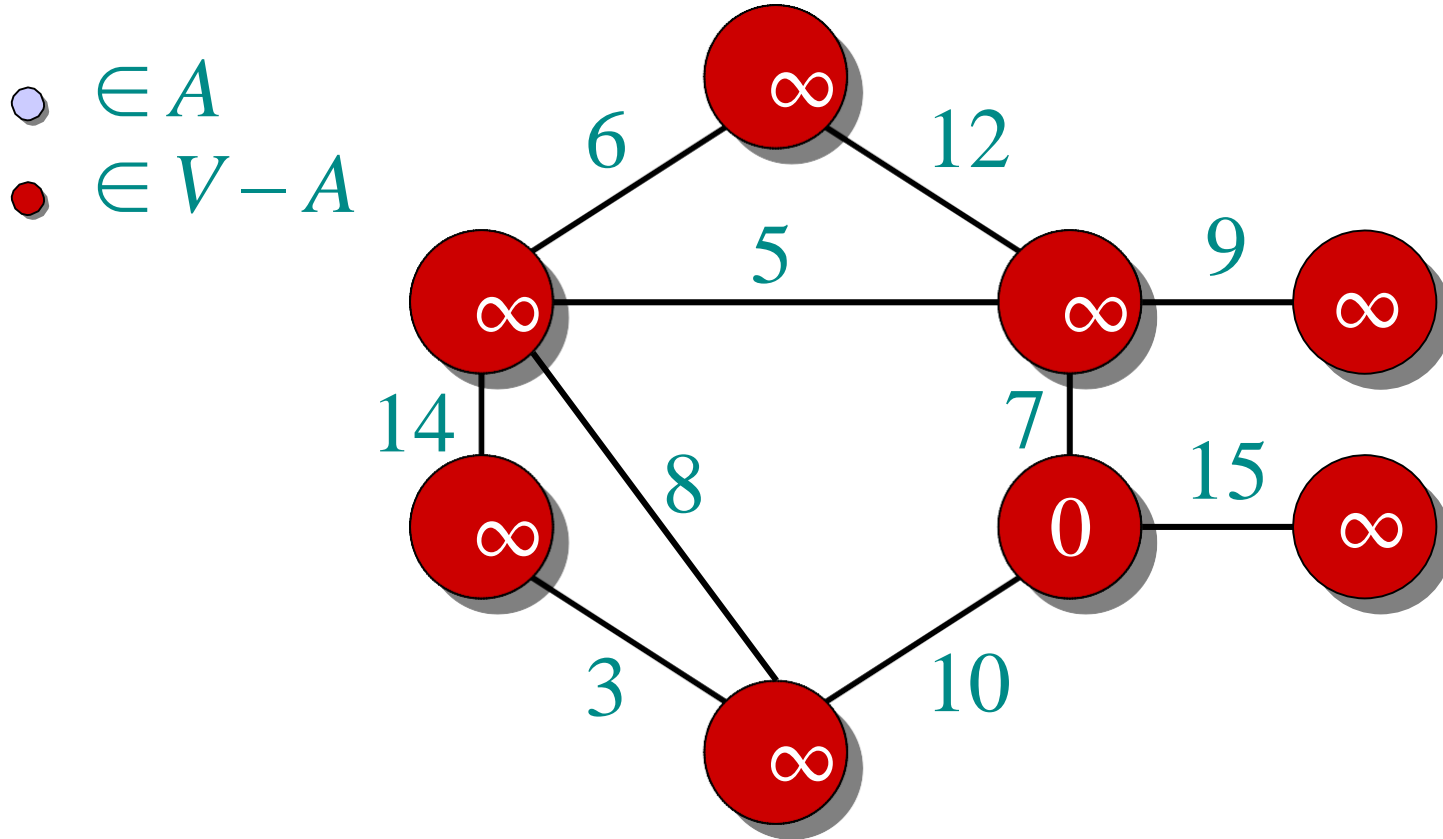
Simplifying assumption. All edge costs $c_e$ are distinct.

Cycle property. Let C be any cycle in G, and let f be the max cost edge belonging to C. Then the MST T* does not contain f.

Pf. (exchange argument)

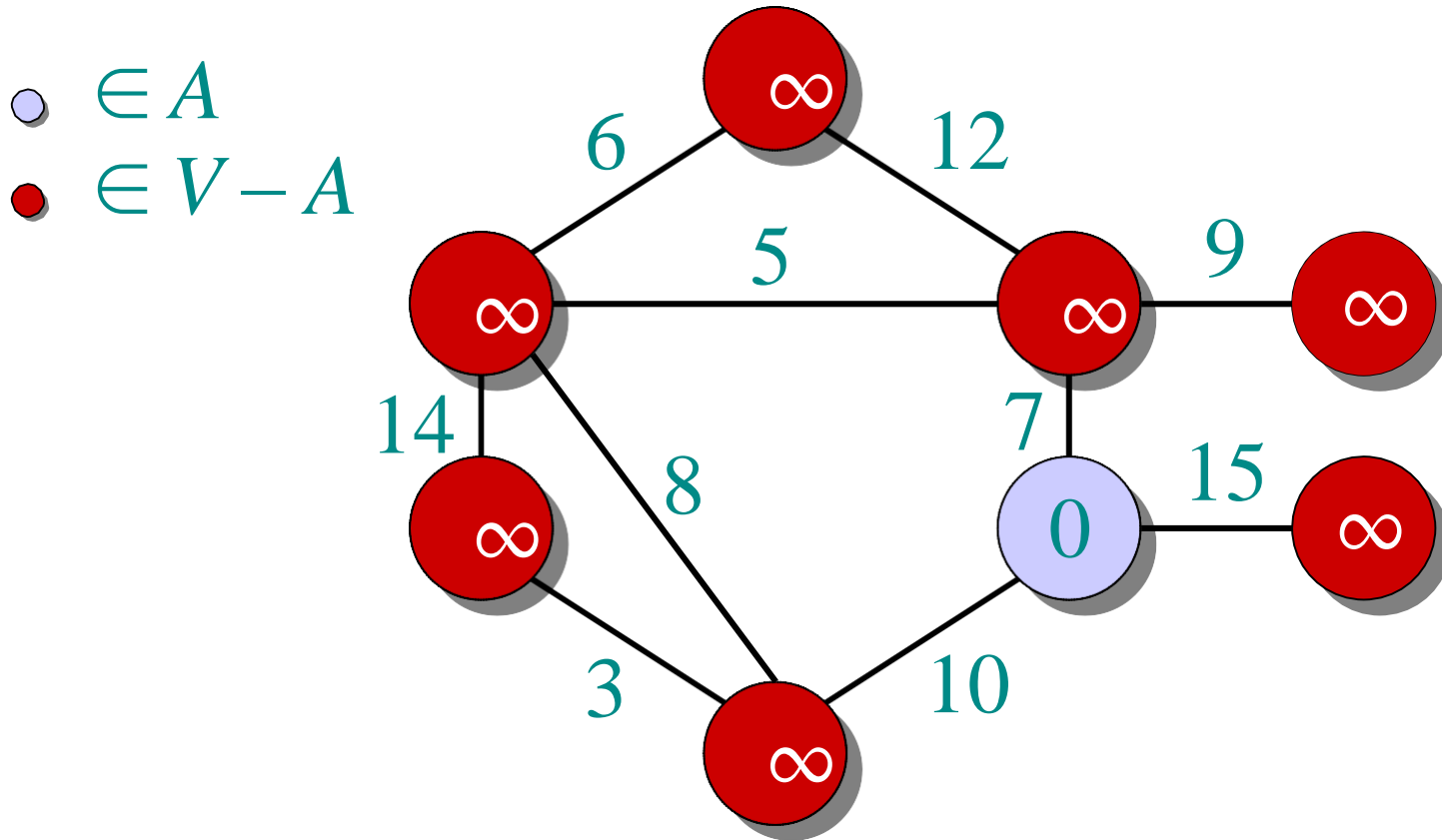- Suppose f belongs to T*, and let's see what happens.
- Deleting f from T* creates a cut S in T*.
- Edge f is both in the cycle C and in the cutset D corresponding to S
  $\Rightarrow$ there exists another edge, say e, that is in both C and D.
- T' = T* $\cup$ { e } - { f } is also a spanning tree.
- Since $c_e < c_f$, cost(T') < cost(T*).
- This is a contradiction. ▪

# Example of Prim's algorithm



$\in A$

$\in V - A$

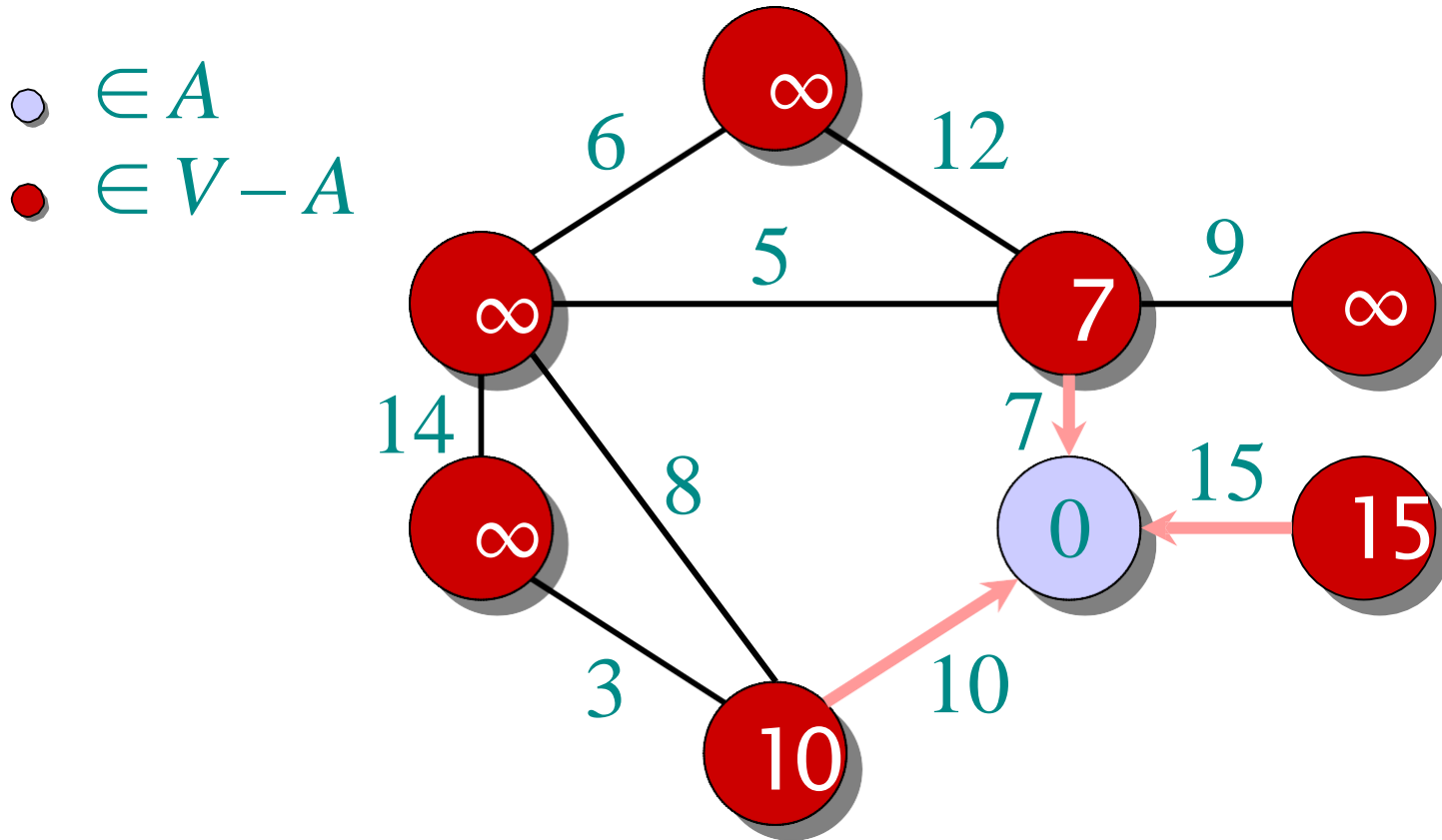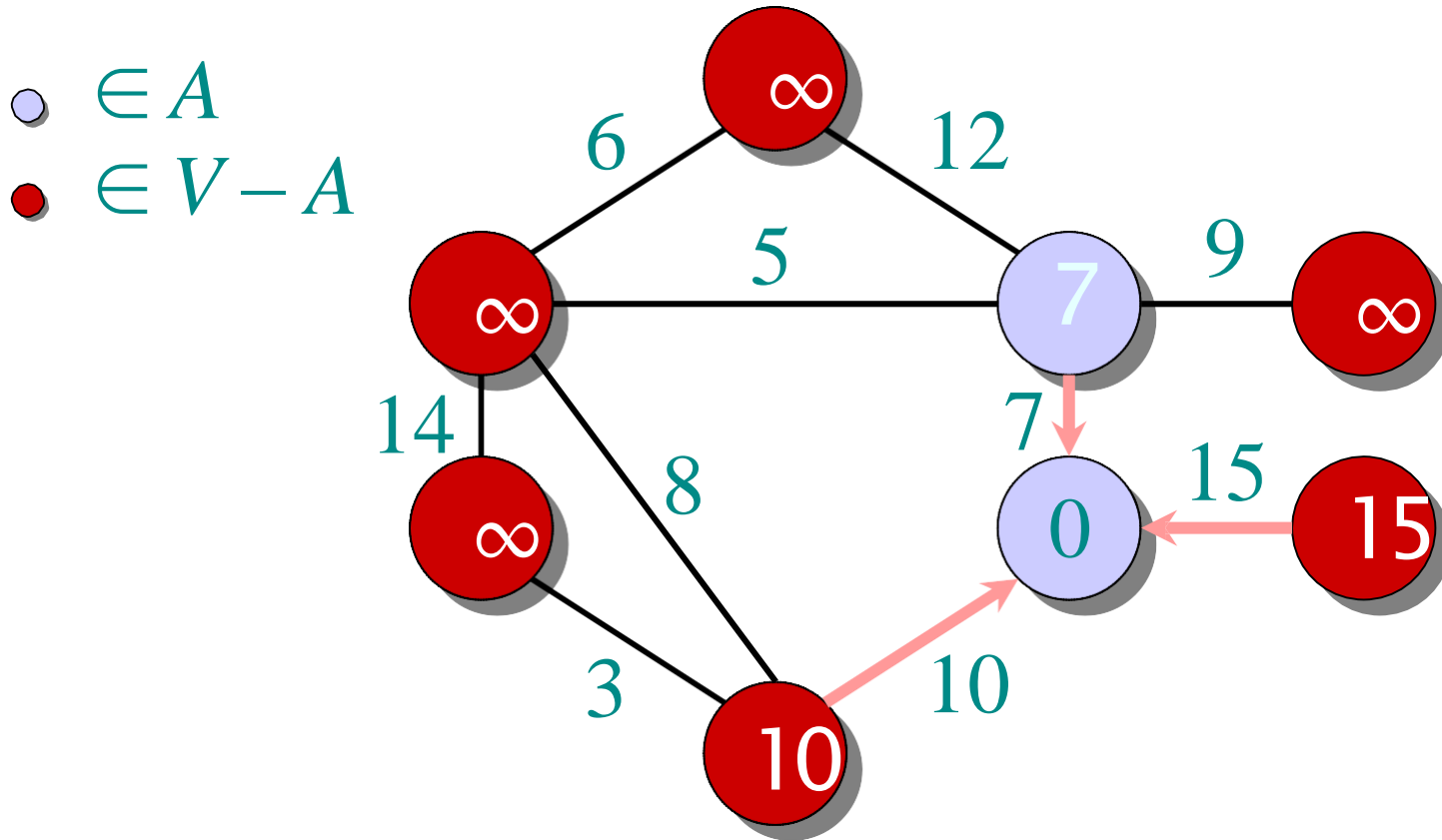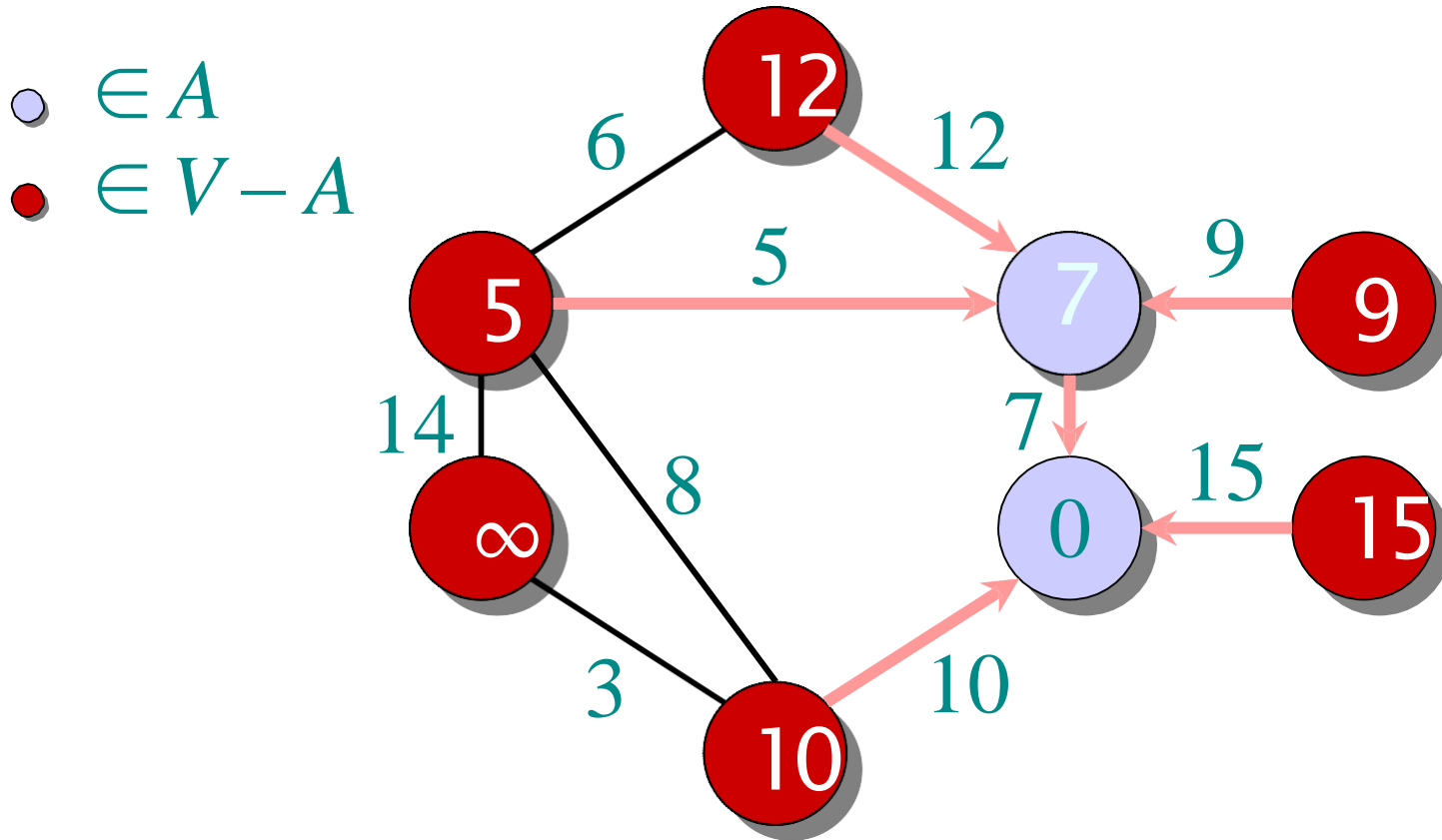# Example of Prim's algorithm

# Example of Prim's algorithm

# Example of Prim's algorithm

# Example of Prim's algorithm

# Example of Prim's algorithm



$\in A$

$\in V - A$

# Example of Prim's algorithm



$\in A$
$\in V - A$

# Example of Prim's algorithm

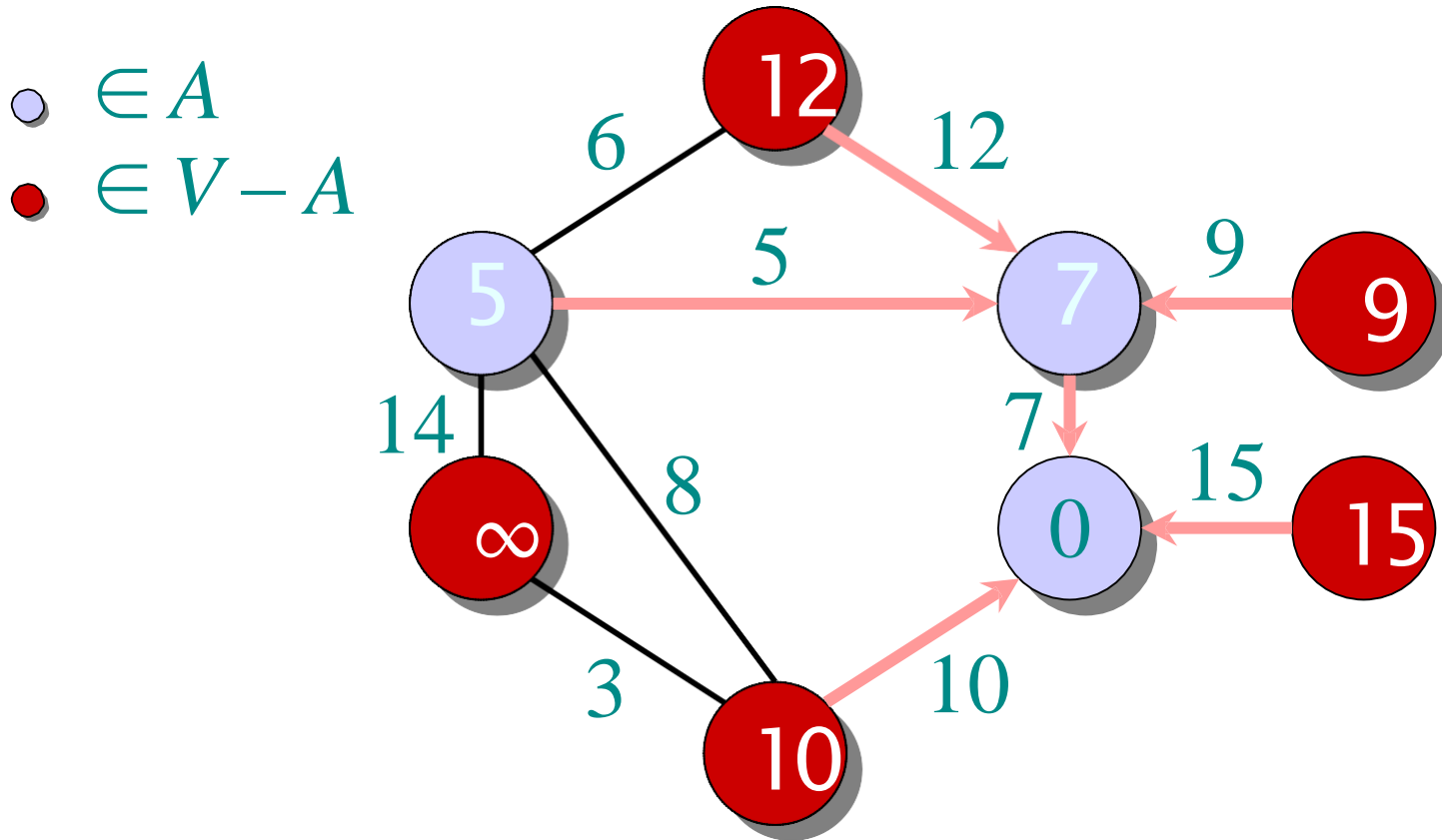# Example of Prim's algorithm

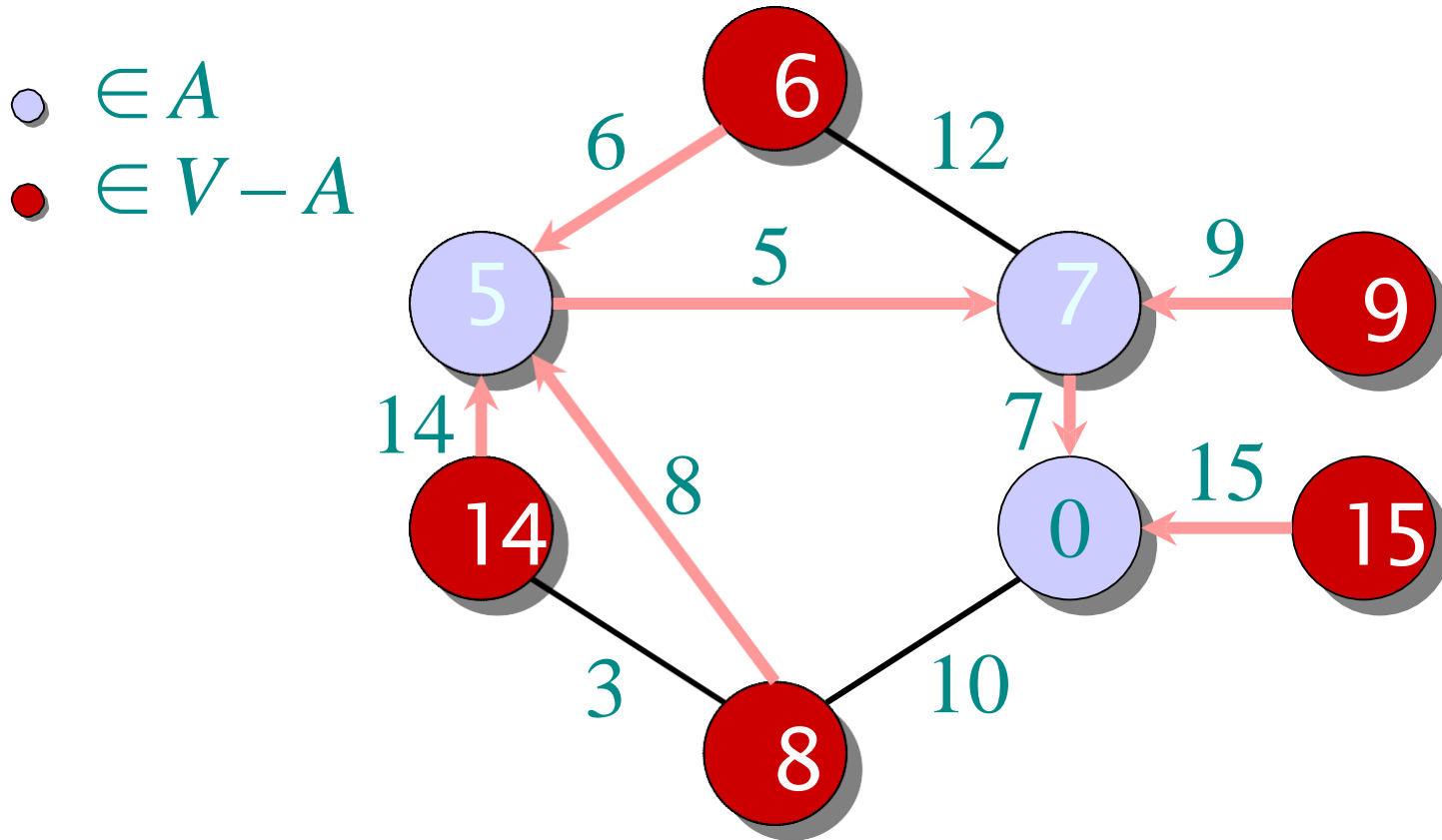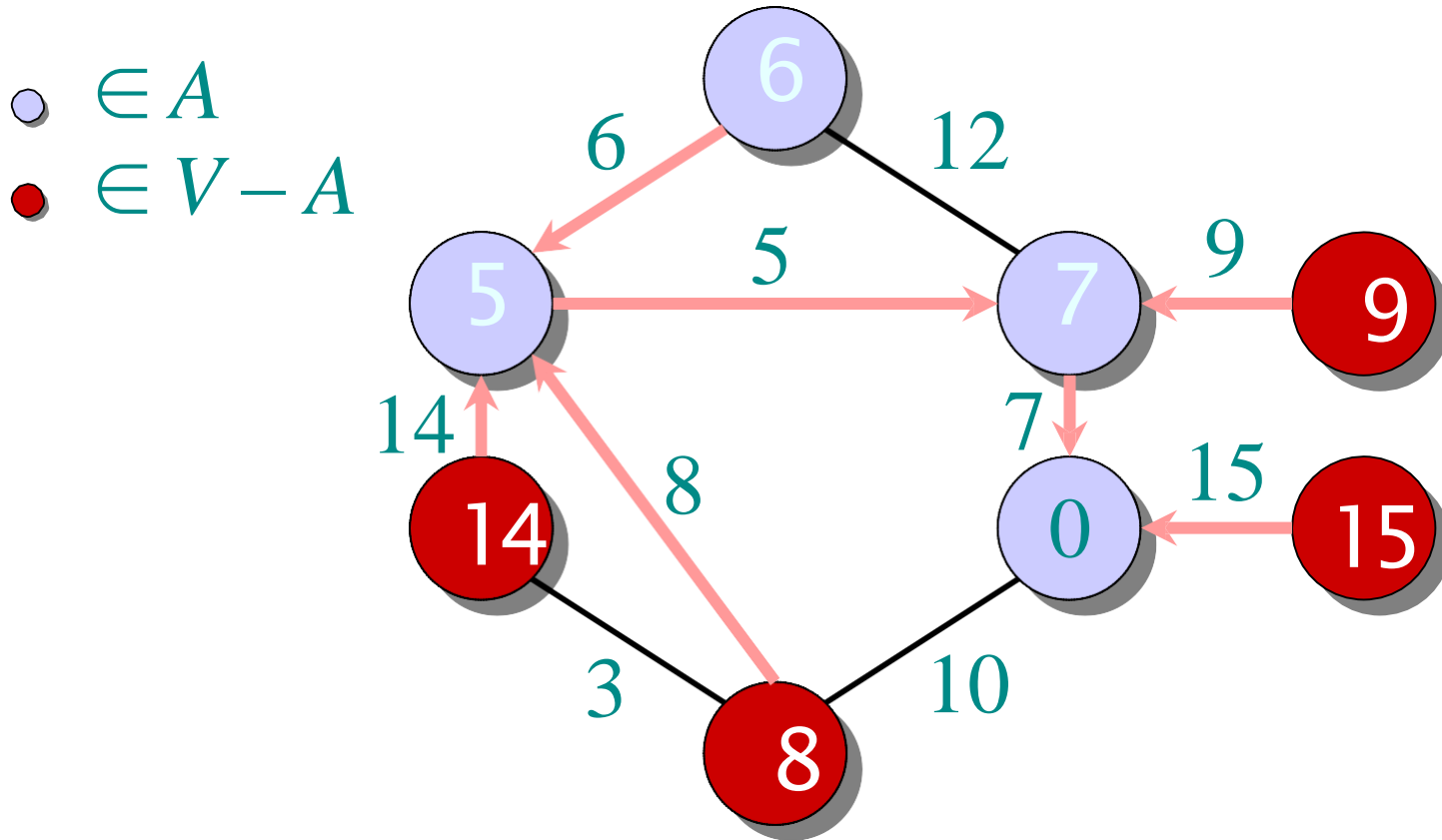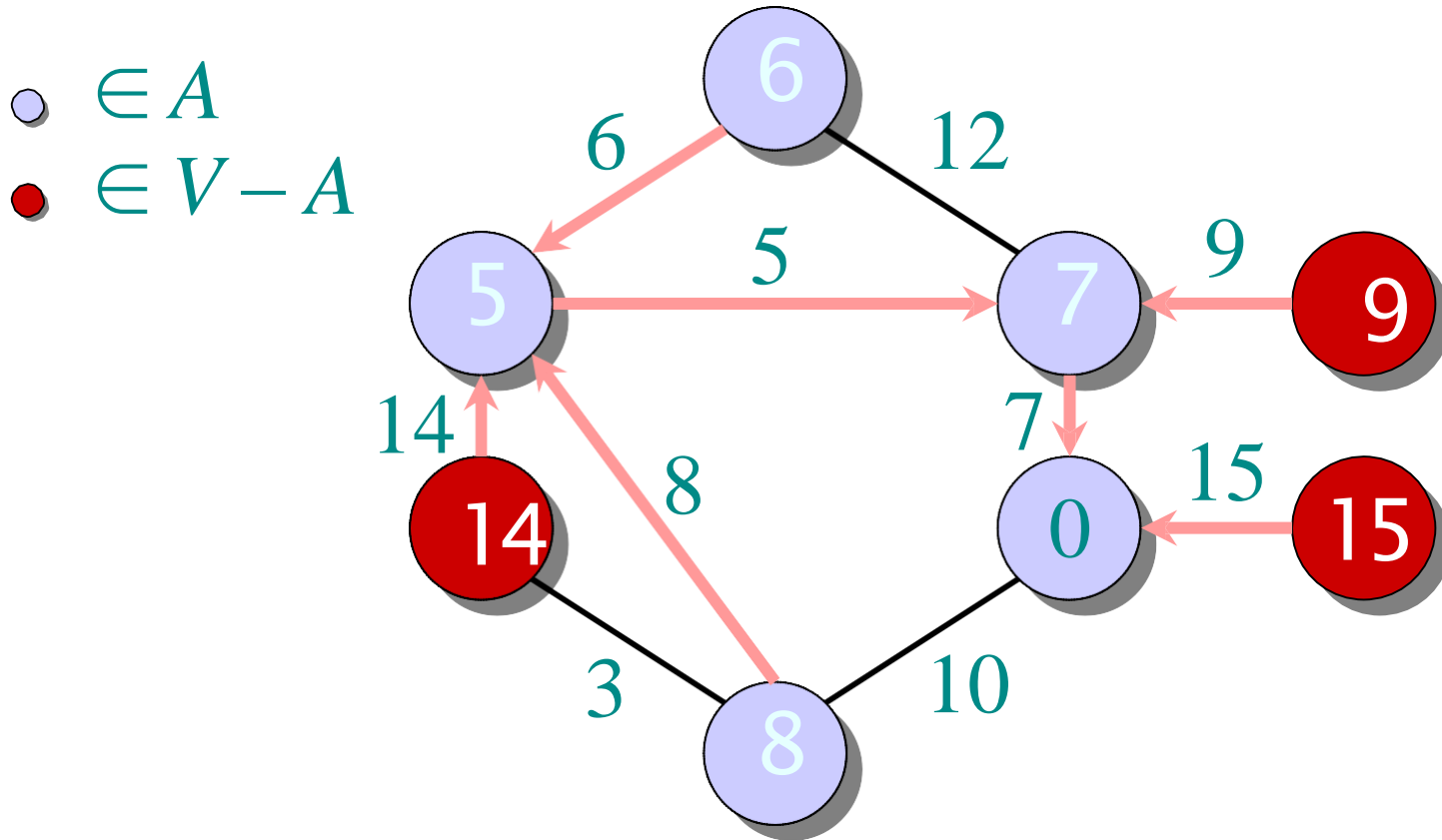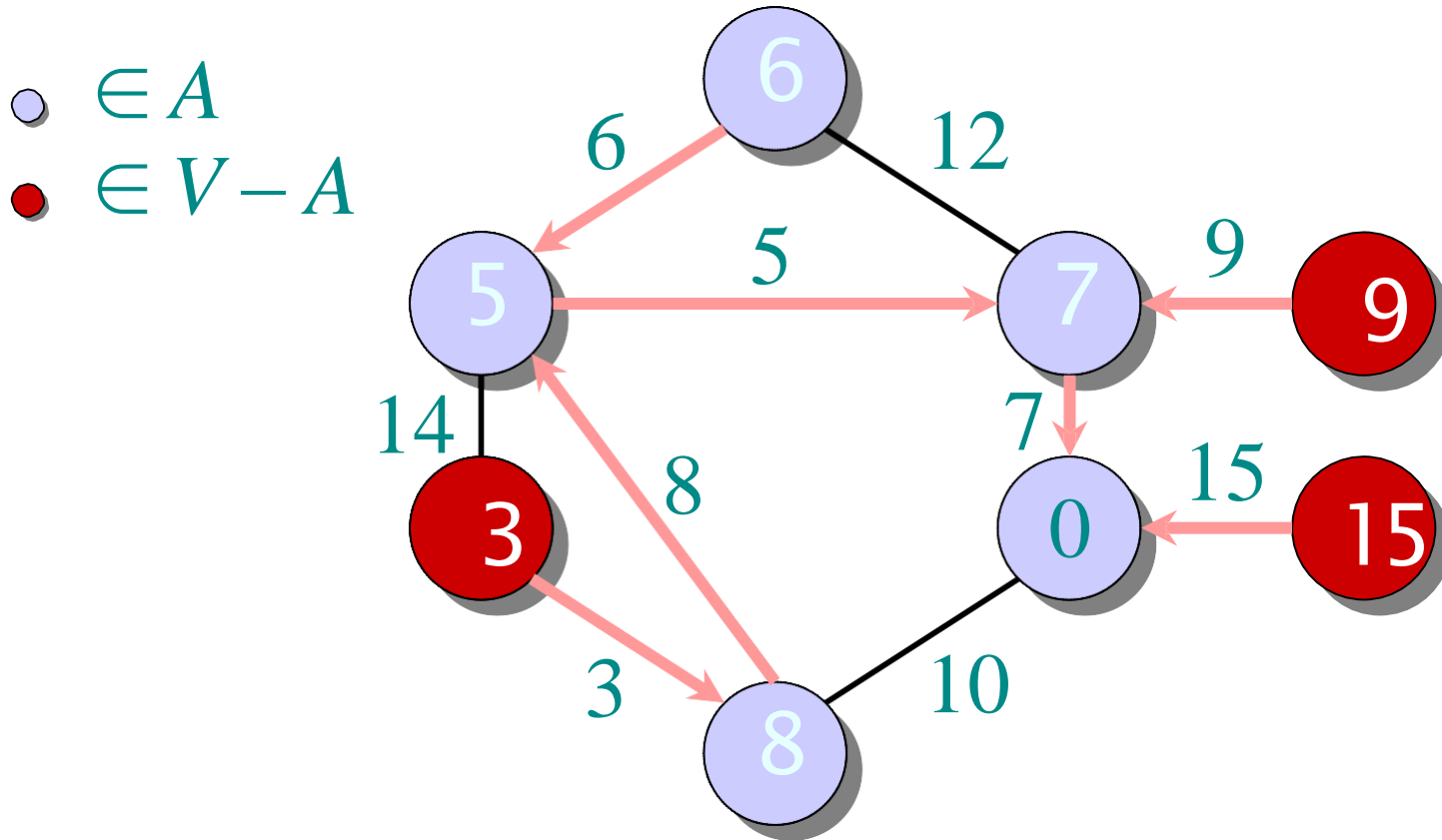# Example of Prim's algorithm

# Example of Prim's algorithm



$\in A$

$\in V - A$

# Example of Prim's algorithm



$\in A$

$\in V - A$

6

12

6

5

5

7

9

9

14

8

3

0

15

15

3

8

10

7

# Example of Prim's algorithm



$\in A$

$\in V - A$

# Prim's Algorithm:  Proof of Correctness

Prim's algorithm.  [Jarník 1930, Dijkstra 1957, Prim 1959]

- Initialize S = any node.
- Apply cut property to S.
- Add min cost edge in cutset corresponding to S to T, and add one new explored node u to S.

# Implementation:  Prim's Algorithm

Implementation.  Use a priority queue ala Dijkstra.
- Maintain set of explored nodes S.
- For each unexplored node v, maintain attachment cost a[v] = cost of cheapest edge v to a node in S.
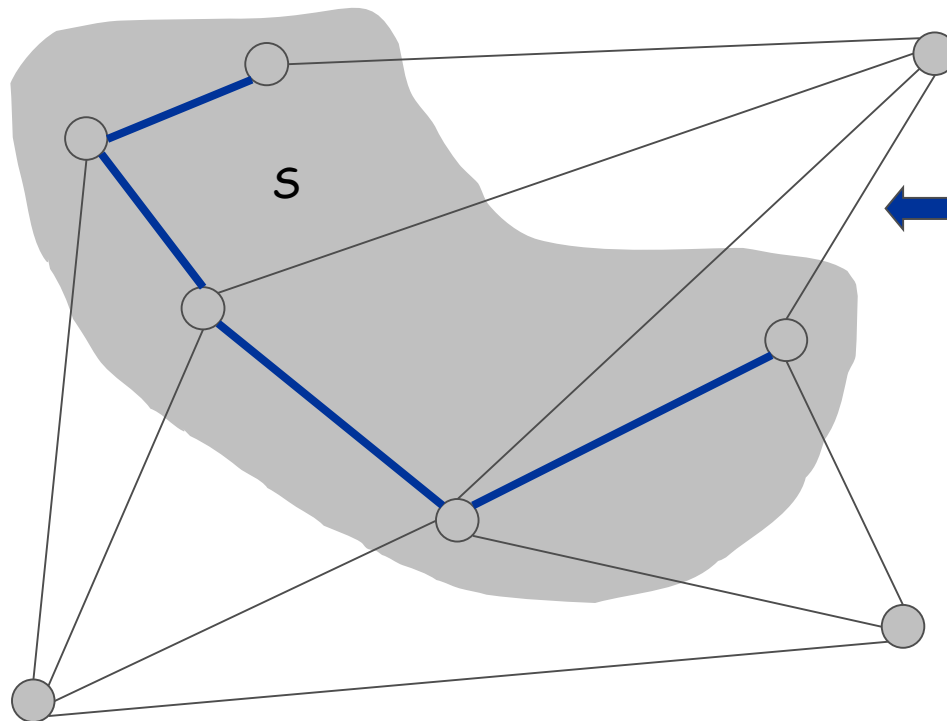- $O(n^2)$ with an array; $O(m \log n)$ with a binary heap.

```
Prim(G, c) {
    foreach (v ∈ V) a[v] ← ∞
    Initialize an empty priority queue Q
    foreach (v ∈ V) insert v onto Q
    Initialize set of explored nodes S ← φ

    while (Q is not empty) {
        u ← delete min element from Q
        S ← S ∪ { u }
        foreach (edge e = (u, v) incident to u)
            if ((v ∉ S) and (c_e < a[v]))
                decrease priority a[v] to c_e
    }
}
```

# Kruskal's Algorithm:  Proof of Correctness

Kruskal's algorithm.  [Kruskal, 1956]

- Consider edges in ascending order of weight.
- Case 1:  If adding e to T creates a cycle, discard e according to cycle property.
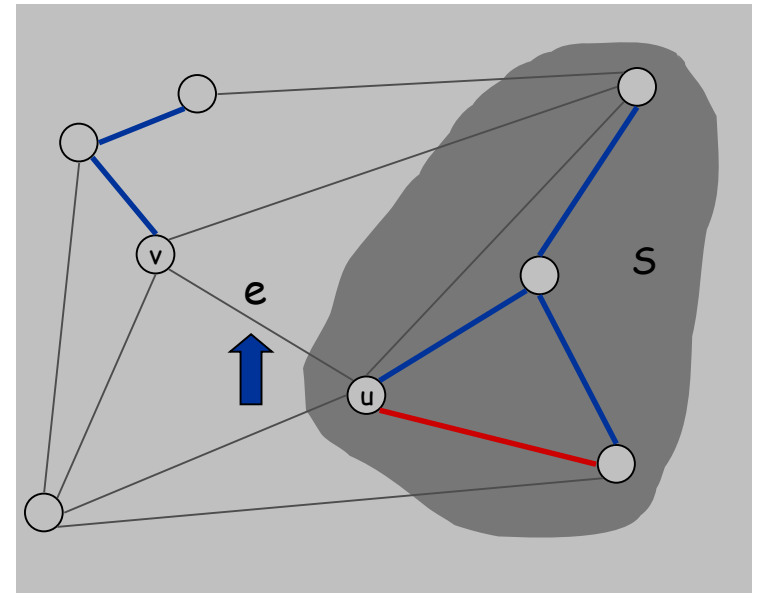- Case 2:  Otherwise, insert e = (u, v) into T according to cut property where S = set of nodes in u's connected component.



Case 1

Case 2

# Implementation:  Kruskal's Algorithm

Implementation.  Use the union-find data structure.
 - Build set T of edges in the MST.
 - Maintain set for each connected component.

Kruskal's Algorithm can be implemented on a graph with n nodes and m edges to run in O(m log n) time.

```
Kruskal(G, c) {
    Sort edges weights so that c₁ ≤ c₂ ≤ ... ≤ cₘ.
    T ← φ

    foreach (u ∈ V) make a set containing singleton u

    for i = 1 to m        are u and v in different connected components?
        (u,v) = eᵢ                    ↙
        if (u and v are in different sets) {
            T ← T ∪ {eᵢ}
            merge the sets containing u and v
        }                     ↖
    return T                      merge two components
}
```

# Lexicographic Tiebreaking

To remove the assumption that all edge costs are distinct:  perturb all edge costs by tiny amounts to break any ties.

Impact.  Kruskal and Prim only interact with costs via pairwise comparisons.  If perturbations are sufficiently small, MST with perturbed costs is MST with original costs.

↑
e.g., if all edge costs are integers,
perturbing cost of edge $e_i$ by $i / n^2$

Implementation.  Can handle arbitrarily small perturbations implicitly by breaking ties lexicographically, according to index.

```
boolean less(i, j) {
    if       (cost(e_i) < cost(e_j)) return true
    else if (cost(e_i) > cost(e_j)) return false
    else if (i < j)                  return true
    else                             return false
}
```

# Outline

- Interval scheduling

- Interval partitioning

- Scheduling to minimize lateness

- Basic elements of greedy algorithms

- Single-source shortest path: Dijkstra's algorithm

- Minimal spanning trees: Prim, Kruskal

- single-link clustering

# Clustering



Outbreak of cholera deaths  in London in 1850s.
Reference: Nina Mishra, HP Labs

# Clustering

Clustering. Given a set U of n objects labeled $p_1$, ..., $p_n$, classify into coherent groups.

↑

photos, documents. micro-organisms

Distance function. Numeric value specifying "closeness" of two objects.

↑

number of corresponding pixels whose
intensities differ by some threshold

Fundamental problem. Divide into clusters so that points in different clusters are far apart.

- Routing in mobile ad hoc networks.
- Identify patterns in gene expression.
- Document categorization for web search.
- Similarity searching in medical image databases
- Skycat: cluster $10^9$ sky objects into stars, quasars, galaxies.
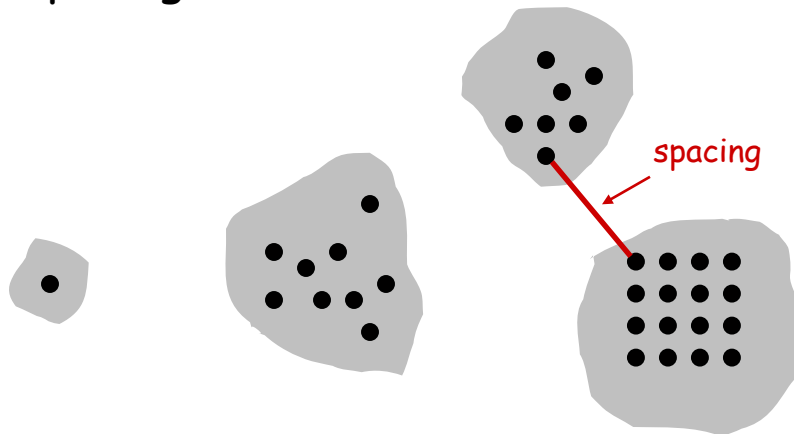
# Clustering of Maximum Spacing

k-clustering.  Divide objects into k non-empty groups.

Distance function.  Assume it satisfies several natural properties.
- $d(p_i, p_j) = 0$ iff $p_i = p_j$   (identity of indiscernibles)
- $d(p_i, p_j) \geq 0$         (nonnegativity)
- $d(p_i, p_j) = d(p_j, p_i)$    (symmetry)

Spacing.  Min distance between any pair of points in different clusters.

Clustering of maximum spacing.  Given an integer k, find a k-clustering of maximum spacing.



spacing

k = 4

# Greedy Clustering Algorithm

Single-link k-clustering algorithm.
- Form a graph on the vertex set U, corresponding to n clusters.
- Find the closest pair of objects such that each object is in a different cluster, and add an edge between them.
- Repeat n-k times until there are exactly k clusters.

Key observation.  This procedure is precisely Kruskal's algorithm (except we stop when there are k connected components).

Remark.  Equivalent to finding an MST and deleting the k-1 most expensive edges.

# Greedy Clustering Algorithm:  Analysis

**Theorem.** Let $C^*$ denote the clustering $C^*_1, ..., C^*_k$ formed by deleting the k-1 most expensive edges of a MST. $C^*$ is a k-clustering of max spacing.

**Pf.**  Let C denote some other clustering $C_1, ..., C_k$.
- The spacing of $C^*$ is the length $d^*$ of the $(k-1)^{st}$ most expensive edge.
- Let $p_i, p_j$ be in the same cluster in $C^*$, say $C^*_r$, but different clusters in C, say $C_s$ and $C_t$.
- Some edge (p, q) on $p_i$-$p_j$ path in $C^*_r$ spans two different clusters in C.
- All edges on $p_i$-$p_j$ path have length $\leq d^*$ since Kruskal chose them.
- Spacing of C is $\leq d^*$ since p and q are in different clusters.  ▪