

## Stable Matching

- **Definition:** A matching is stable if no unmatched man and woman both prefer each other to their current partners

- **Gale-Shapely Algorithm:**

```
Initialize each person to be free.
while (some man is free and hasn't proposed to every woman)
    Choose such a man m
    w = 1st woman on m's list to whom m has not yet proposed
    if (w is free)
        assign m and w to be engaged
    else if (w prefers m to her fiancé m')
        assign m and w to be engaged, and m' to be free
    else
        w rejects m
}
```

- Lemma: Men propose to women in decreasing order of preference
- Lemma: Women stay engaged after the first time they got engaged
- Lemma: Women's partners keep getting better
- Lemma: Upon termination, each man is engaged to a unique woman
- Lemma: Upon termination, the matching between men and women is stable
- The algorithm returns male-optimal stable matching

## Asymptotic Analysis

- **Upper:**  $T(n) = O(f(n))$   
If there  $\exists$  constants  $c > 0$ ,  $n_0 \geq 0$  s.t.  $\forall n \geq n_0$ ,  $T(n) \leq c \cdot f(n)$
- **Lower:**  $T(n) = \Omega(f(n))$   
If there  $\exists$  constants  $c > 0$ ,  $n_0 \geq 0$  s.t.  $\forall n \geq n_0$ ,  $T(n) \geq c \cdot f(n)$
- **Tight:**  $T(n) = \Theta(f(n))$   
If  $T(n) = O(f(n))$  and  $T(n) = \Omega(f(n))$

### Transitivity

- $f = O/\Omega/\Theta(g)$  and  $g = O/\Omega/\Theta(h) \rightarrow f = O/\Omega/\Theta(h)$

### Additivity

- $f = O/\Omega/\Theta(h)$  and  $g = O/\Omega/\Theta(h) \rightarrow f + g = O/\Omega/\Theta(h)$

### Rule of Thumbs

$$n^n \geq n! \geq c^{kn} \geq c^n \geq n^k \geq n \log(n) \geq n \geq \sqrt{n} \geq \log(n) \geq \log(\log(n))$$
$$O(\log(n!)) = O(n \log(n))$$

### Mathematical Properties

- $a^{\log b} = b^{\log a}$
- $1 + \frac{1}{2} + \frac{1}{3} + \frac{1}{4} + \frac{1}{5} + \dots + \frac{1}{n} = O(\log(n))$
- GP sum =  $\frac{a(r^n - 1)}{r - 1}$ ,  $|r| > 1$  or  $\frac{a}{1 - r}$ ,  $|r| < 1$
- AP sum =  $\frac{n}{2}(2a + (n - 1)d) = \frac{n}{2}(a + l)$

### Common Recurrence Relations

- $T(n) = 2T\left(\frac{n}{2}\right) + n^2 = O(n^2)$
- $T(n) = 2^k T(\sqrt[k]{n}) + c = O(\log^k(n))$
- $T(n) = 2^k T\left(n^{\frac{1}{m}}\right) + c = O\left(\log^{\frac{k}{m}}(n)\right)$

## The Master Theorem

The **Master Method** depends on the following Theorem:

**Theorem:** Let  $a \geq 1$  and  $b > 1$  be constants, let  $f(n)$  be a function, and let  $T(n)$  be defined on the nonnegative integers by the recurrence:

$$T(n) = aT\left(\frac{n}{b}\right) + f(n)$$

Then  $T(n)$  can be bounded asymptotically as follows.

1. If  $f(n) = O(n^{\log_b a - \epsilon})$  for some constant  $\epsilon > 0$ , then  $T(n) = \Theta(n^{\log_b a})$ .
2. If  $f(n) = \Theta(n^{\log_b a})$ , then  $T(n) = \Theta(n^{\log_b a} \lg n)$ .
3. If  $f(n) = \Omega(n^{\log_b a + \epsilon})$  for some constant  $\epsilon > 0$ , and if  $a \cdot f\left(\frac{n}{b}\right) \leq cf(n)$  for some constant  $c < 1$  and all sufficiently large  $n$ , then  $T(n) = \Theta(f(n))$ .

regularity condition

## Graphs

	Adjacency List	Adjacency Matrix
Check if (u,v) exists	$\Theta(\deg(u)) = O(V)$	$\Theta(1)$
Enumerate all E	$\Theta(V + E) = O(V^2)$	$\Theta(V^2)$
Use cases	Sparse graphs	Dense graphs

### Definitions and theorems

- **Trees:**
  - Def: An undirected graph is a tree if it is connected and acyclic
  - Theorem: In any undirected graph G with V nodes, any two statements below imply the third:
    1. G is connected
    2. G is acyclic
    3. G has  $V - 1$  edges
- **Simple path:** A path where all its vertices are distinct
- **Connected Graph (Undirected):** There is a path between any vertices  $u$  and  $v$
- **Strongly Connected Graph (Directed):** There is a path from  $u$  to  $v$  and from  $v$  to  $u$  for any pair of vertices  $u$  and  $v$

### Graph Traversals

- **BFS:**
  - Use cases
    1. Finding shortest paths
    2. Connectivity
    3. Testing bipartitenessA graph is not bipartite if it contains an odd-lengthed cycle  
Use extra 'Color' array and assign color whenever a node is added to L(i+1). After BFS, check for edges for which both ends have the same color
  - 4. Testing strong connectivity in directed graphs  
BFS G from S. Then, BFS  $G^{ev}$  from S. G is strongly connected if both BFS visited all vertices in G

- Properties:
  1. Edges in graph G but not in its BFS Tree T all either connect nodes in the same layer in T or connect nodes in adjacent layers
- **DFS:**
  - Use cases
    1. Finding the set of all connected components
  - Properties:
    1. For a given recursive DFS(u) call, all vertices marked "explored" between the invocation and end of this call are descendants of u in the DFS tree T
    2. If a graph G contains an edge (u, v) that is not in its DFS tree T, then one of u or v is an ancestor or of the other (i.e. diff levels)

### BFS(s):

```
visited[s] ← true
visited[v] ← false for all other v
L(0) ← list containing only s
i ← 0 // layer
T ← ∅
while L[i] is not empty do
    L(i + 1) ← empty list
    for each node u ∈ L(i)
        Consider each edge (u, v)
        if visited[v] = false then
            set visited[v] = true
            add edge (u, v) to the tree T
            add v to L(i + 1)
    endfor
    increment i by one
endwhile
```

### DFS(s):

```
visited[s] ← false for all v
S(0) ← stack containing only s
parent[] ← empty list
T ← ∅
while S is not empty do
    pop u from S
    if visited[u] = false then
        set visited[u] = true
        add (u, parent[u]) to T
        for each edge (u, v)
            push v to S
            set parent[v] to u
        endfor
    endif
endwhile
```

### DAGs

- **Properties:**
  - In every DAG, there is a node v with no incoming edges
  - A graph G is a DAG if and only if it has a topological ordering
- **Finding a topological ordering:**
  - Kahn's Algorithm [  $O(V + E)$  ]
    1. Initialise set S that contains all nodes with no incoming edges
    2. Initialise set W to count number of incoming edges for each node
    3. Repeat until S is empty:
      - 3.1. Pick any node u from S
      - 3.2. Add u to the topological order
      - 3.3. For each (u, v<sub>i</sub>), decrement W[v<sub>i</sub>]
      - 3.4. If W[v<sub>i</sub>] becomes 0, add v<sub>i</sub> to S

## Greedy

### Proving Techniques

**Exchange argument:** Show that at each step, you can exchange S's current choice with G's current choice without hurting S's quality.

Example: Interval scheduling

1. Let  $G = i_1, i_2, \dots, i_k$  and  $S = j_1, j_2, \dots, j_m$  for an input L
2. Let P(m) be the proposition that if S returns m number of intervals, then G also returns m number of intervals

- Base case:  $P(1)$ . The optimal solution has only 1 interval. Trivially,  $G$  can pick any interval, hence  $P(1)$  is true
- Inductive hypothesis:  $P(m)$  is true
- $f(i_1) \leq f(j_1)$  since  $G$  always chooses the request with the earliest finish time.
- Therefore,  $S^* = i_1, j_2, \dots, j_{m+1}$  is also an optimal solution (explain)

- $S^{**} = j_2, j_3, \dots, j_{m+1}$  must be optimal for  $L \setminus \{i_1\}$  for  $S^*$  to be optimal for  $L$ .  $S^{**}$  outputs  $m$  intervals
- By construction,  $G$  outputs  $i_2, \dots, i_k$  for  $L \setminus \{i_1\}$ . By the inductive hypothesis,  $G$  must output  $m$  schedules. Hence,  $m = k - 1$
- Hence,  $k = m + 1$ . Therefore,  $P(m+1)$  is true.

**Structural bound:** every possible solution must adhere to some min/max and show that  $G$  produces min/max

**“Greedy stays ahead”:** Show that at each step,  $G$  is always as good as  $S$ . Show that “Greedy stays ahead” implies optimality

### Interval Scheduling

- Rule: Schedule the request with the earliest finish time

#### IntervalScheduling(R):

```

A ← []
visited[s] ← true
while R is not empty do
    choose  $r_1$  in R with earlier finish time
    add  $r_1$  to A
    delete  $r_1$  in R that are incompatible with  $r_1$ 
endwhile
return A

```

#### Pf. (Greedy stays ahead)

- Let  $i_1, \dots, i_k$  be the set of requests from  $G$  and  $j_1, \dots, j_m$  be the set of requests from  $S$
- It suffices to show that if  $f(i_r) \leq f(j_r)$  for all  $r \leq k$ , then  $k \geq m$
- Proof by induction on  $r$ 
  - Base case:  $r = 1$ .  $G$  will choose  $i_1$  which is the request with earliest finish time
  - Inductive hypothesis: Suppose  $f(i_{r-1}) \leq f(j_{r-1})$
  - $f(i_{r-1}) \leq f(j_{r-1})$ , so  $f(i_{r-1}) \leq s(j_r)$
  - Hence,  $j_r$  must be available for selection by  $G$  for the  $r^{\text{th}}$  request
  - Hence,  $f(i_r) \leq f(j_r)$

### Interval Partitioning

- Rule: Consider the resources in the order of their start time

#### IntervalPartitioning(R):

```

d ← 0
sort intervals in R in ascending order of start time
for j = 1 to n
    if interval j is compatible with some resource k
        schedule interval j to resource k
    else
        allocate new resource d + 1
        schedule interval j to resource d + 1
        d ← d + 1

```

#### Pf. (Structural Bound)

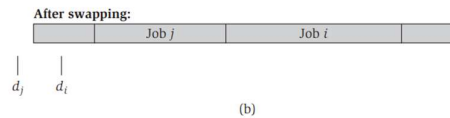
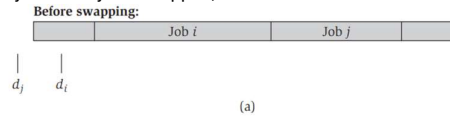
- Observation: In any instance of interval partitioning, the number of resources needed is at least the depth of the set of intervals

#### Minimising Lateness

- Rule: Schedule requests with the earliest deadlines first

#### Pf. (Exchange Argument)

- Observation: An inversion is when job  $j$  is scheduled after job  $i$  when  $d(j) < d(i)$ . If  $S$  has an inversion, then there must be a pair of jobs  $i$  and  $j$  such that  $j$  is scheduled immediately after  $i$  and has  $d(j) < d(i)$ .
- Suppose  $S$  has at least one inversion. Let job  $i$  be scheduled immediately after job  $j$  even though  $d(j) < d(i)$
  - If jobs  $i$  and  $j$  are swapped,  $S$  will have one less inversion



- Denote  $L'$  as the max lateness after swap and  $L$  as the max lateness before swap. Denote  $t(m)$  as the time taken to complete a job  $m$
- $L' = \max\{t(j) - d(j), t(j) + t(i) - d(i)\}$   
 $L = \max\{t(i) - d(i), t(i) + t(j) - d(j)\} = t(i) + t(j) - d(j)$ , since  $d(j) < d(i)$
- $t(j) - d(j) < t(i) + t(j) - d(j)$  and  $t(j) + t(i) - d(i) < t(i) + t(j) - d(j)$ . Therefore,  $L'$  must be smaller than  $L$ .
- We've shown that the lateness of  $S$  does not increase after the swap
- This shows that an optimal schedule with no inversions exists.
- All schedules with no inversions have the same maximum lateness (to proof). Hence, the schedule obtained by the greedy solution is optimal.

### Divide and Conquer

#### Counting Inversions

- Mergesort, but count the number of inversions during the merge step
- Key property: if  $L[i] > R[j]$ , then  $L[i:] > R[j:]$

#### Merge (L, R):

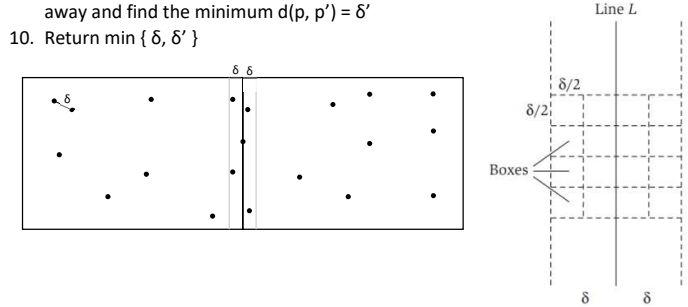
```

count ← 0
i, j ← 0
M ← empty array
while i < len(L) and j < len(R) do
    if L[i] < R[j] then
        add L[i] to back of M and increment i
    else
        add R[j] to back of M and increment j
    increment count by len(L) - i + 1
end while
... regular merge steps

```

#### Finding the closest pair of points

- Find closest pair in left half  $L$  and right half  $R$  and those in the boundary split
- Keep a sorted list  $P_x$  of points where points are sorted by  $x$ -axis. Keep a sorted list  $P_y$  of points where points are sorted by  $y$ -axis.
  - Split the points into two halves  $Q$  and  $R$  by their  $x$ -axis.
  - Recurse on both  $Q$  and  $R$  and find the minimum distance  $\delta$  between two pairs of points between the two
  - Choose a point  $q$  in  $Q$  that has highest  $x$  value  $x^*$  and draw a vertical line  $L$  through it.  $L$  essentially divides the points in  $Q$  and  $R$
  - (\*\*Pf 1) Narrow down to only points with  $x$ -values inside the  $[x^* - \delta, x^* + \delta]$  boundary.
  - Construct  $S_y$  containing only the points above, sorted by  $y$ -axis (can be done in  $O(n)$  using  $P_y$ )
  - Divide the  $[x^* - \delta, x^* + \delta]$  boundary into many  $\delta/2 \times \delta/2$  boxes
  - (\*\*Pf 2) If there exists any 2 points  $p$  and  $p'$  such that  $d(p, p') < \delta$ , then  $p$  and  $p'$  must be at most **15** positions away in  $S_y$
  - Compute each  $d(p, p')$  in  $S_y$  such that  $p$  and  $p'$  are within 15 positions away and find the minimum  $d(p, p') = \delta'$
  - Return  $\min\{\delta, \delta'\}$



- Pf 1:** If  $\exists q = (q_x, q_y) \in Q$  and  $r = (r_x, r_y) \in R$  s.t.  $d(q, r) < \delta$ , then  $q$  and  $r$  lies within distance  $\delta$  of  $L$   
 $x^* - q_x \leq r_x - q_x \leq d(q, r) \leq \delta$  and  $r_x - x^* \leq r_x - q_x \leq d(q, r) \leq \delta$   
 $\rightarrow$  Hence  $q$  and  $r$  have  $x$ -coordinate within  $\delta$  of  $L$
- Pf 2:** The max distance within a box is the length of the diagonal, which is  $\sqrt{2 \left(\frac{\delta}{2}\right)^2} = \sqrt{\frac{\delta^2}{2}} = \frac{\delta}{\sqrt{2}} \leq \delta$ . Hence, no two points can be in the same box  
 Suppose there exists points  $s$  and  $s'$  in  $S_y$  such that  $d(s, s') < \delta$  and that they are 16 positions apart. Assume WLOG that  $s_y < s'_y$ . Then,  $s$  and  $s'$  must be separated by at least 3 rows of boxes which must have a distance of at least  $\frac{3}{2}\delta \geq \delta$  – a contradiction

#### Proving Techniques

Proof by (strong) induction on the input size  $n$

- Define the proposition
- Show how the base case is fulfilled by the algorithm
- Suppose  $P(k)$  is true for all  $k < n$
- By inductive hypothesis,  $P(n/c)$  must be true
- Show how combining the subproblems causes  $P(n)$  to be true

### Dynamic Programming

## Knapsack Problem

$$DP(S, W) = \begin{cases} 0, & S = \emptyset, W \leq 0 \\ \max\{w_i + DP(S \setminus \{n_i\}, W - w_i), DP(S \setminus \{n_i\}, W)\}, & S \neq \emptyset, W > 0 \end{cases}$$

**Knapsack**( $S_n, W$ ):

```
M ← (n + 1) x (w + 1) array
M[0][W] ← 0 for all w
for i from 1 to n
    for w from 0 to W
        if w_i < w then
            M[i][w] ← max{w_i + M[i - 1][w - w_i], M[i - 1][w]}
        else
            M[i][w] ← M[i - 1][w]
return M[n][W]
```

## Network Flow

### Definitions

- s-t cut**: partition (A, B) of V such that  $s \in A$  and  $t \in B$
- cap(A, B)**: capacity of an s-t cut (A, B) =  $\sum_{e \text{ out of } A} c(e)$
- s-t flow** must satisfy the following constraints:
  - $0 \leq f(e) \leq c(e), \forall e \in E$  (capacity)
  - $f^{in}(v) = f^{out}(v), \forall v \in V \setminus \{s, t\}$  (conservation)
- Flow value**:  $v(f) = f^{out}(s)$

### Max-flow Min-cut Theorem (Ford-Fulkerson Algorithm)

- Max flow value from s to t = minimum capacity of any cut
- Flow value lemma**:  $f^{out}(A) - f^{in}(A) = v(f)$ 
  - The net flow sent across any cut in G = flow amount leaving s
- Weak duality:  $v(f) \leq \text{cap}(A, B)$  for any cut (A, B)
- Certificate of Optimality: Let f be any flow and (A, B) be any cut. If  $v(f) = \text{cap}(A, B)$ , then f is a max flow and (A, B) is a min cut
- with the smallest capacity will be the bottleneck edge b
- For every edge e along P, if e is a forward edge, we can increase  $f(e)$  by  $c(b)$ . If e is a backward edge, we decrease  $f(e)$  by  $c(b)$

### FordFulkerson(V, E):

```
f(e) ← 0 ∀ e ∈ E // initialise all flows to 0
while ∃ s-t path in Gf do
    Find simple s-t path P in Gf (BFS/DFS)
    f ← Augment(f, P, G)
    Update Gf
return f
```

### Augment(f, P, G):

```
b ← c(e) of bottleneck edge e along P
for each e in P
    if e is a forward edge then
        increment f(e) in G by b
    else
        decrement f(e) in G by b
return f + b
```

- Residual Graph**  $G_f$  of G can be constructed using the following rules:

- Vertices in  $G_f$  and G are the same
- For each edge e of G, if  $f(e) < c(e)$ , then add e to  $G_f$  but with (residual) capacity =  $c(e) - f(e)$
- For each edge e of G, if  $f(e) > 0$ , then add e to  $G_f$  but reverse the direction

- Augmenting Path** is a simple path P in  $G_f$  from s to t.

### Proof of Ford-Fulkerson

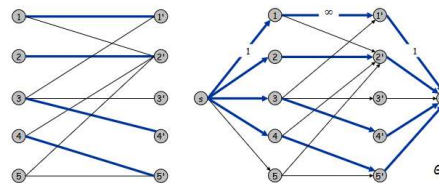
- Show TFAE
  - There exists a cut (A, B) s.t.  $v(f) = \text{cap}(A, B)$
  - Flow f is a max flow
  - There is no augmenting path relative to f
- 1 ⇒ 2: Corollary to weak duality lemma
- 2 ⇒ 3: Proof by contrapositive
  - If there exists an augmenting path, then we can improve f by sending flow along that path (assuming capacity is a non-negative integer)
- 3 ⇒ 1:
  - Let f be a flow with no augmenting paths
  - Let A be the set of vertices reachable from s in  $G_f$  and let  $B = V - A$
  - By defn of A,  $s \in A$  and  $t \notin A$  ( $t \in B$ )
  - Show that  $f^{out}(A) = \sum_{e \text{ out of } A} c(e)$ 
    - There are no directed edges from A to B to some node u in B in  $G_f$ , otherwise vertex u would have been in set A
    - Hence, if there is an edge e from A to B in G (G now, not  $G_f$ ),  $f(e) = c(e)$  for e to not remain in  $G_f$ . i.e.
  - Show that  $f^{in}(A) = 0$ 
    - Suppose  $f^{in}(A) > 0$ . Then, there exists an edge  $e(v, u)$  from B to A such that  $v \in B$  and  $u \in A$ .
    - Then, there must have been a reverse edge  $e'(u, v)$  in  $G_f$  from A to B. If that's the case, v must have been reachable from s ⇒ Contradiction
  - By Flow value lemma,  $v(f) = f^{out}(A) - f^{in}(A)$
  - Hence,  $v(f) = \sum_{e \text{ out of } A} c(e) = \text{cap}(A, B)$  as shown above

### Bipartite Matching

- Bipartite matching = set of edges s.t. no two edges in the set share the same endpoint. Maximum matching = largest of such set

### Max flow implementation

- Add source s and join s to each node in L with edge of capacity = 1
- Add sink t and join each node in R to t with edge of capacity = 1
- Edges between L and R have infinite capacity
- Max flow in this graph  $G'$  = size of maximum matching in G



**Proof** (show  $k \leq f$  and  $k \geq f$ )

Show  $k \leq f$ :

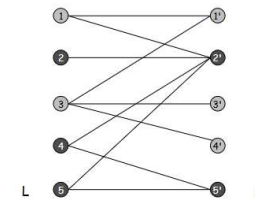
- Given bipartite graph G with max matching of k.
  - Consider a flow f that sends 1 unit along each of the k paths
  - f is a flow and has cardinality k
- Show  $k \geq f$ :
- Let f be a max flow in  $G'$  with value = k
  - By integrality theorem, flow across any edge in  $G'$  is either 0 or 1
  - Each node in L has at most one incoming edge from s, hence its incoming flow is at most 1 and outgoing flow (to R) is at most 1
  - Each node in R has at most one outgoing edge to t, hence its outgoing flow is at most 1 and its incoming flow (from L) is at most 1
  - Hence, there must be a bipartite matching between L and R with cardinality k

### Perfect Matching

- Perfect matching = bipartite matching which covers all the nodes
- Condition:  $|L| = |R|$

### Marriage (Hall's) Theorem

- A bipartite graph  $G = (L \cup R)$  has a perfect matching iff  $|N(S)| \geq |S|$  for all subsets S of L



No perfect matching:  
 $S = \{2, 4, 5\}$   
 $N(S) = \{2', 3', 4', 5'\}$

### Proof of Marriage Theorem

⇒: G has perfect matching →  $|N(S)| \geq |S|$

Each node in S needs to be mapped to a different node in  $N(S)$  for all subsets S of L

⇐:  $|N(S)| \geq |S| \rightarrow$  G has perfect matching. Proof by contrapositive

- Suppose G does not have a perfect matching
- Create  $G'$  from G using the same method in bipartite matching and let (A, B) be a min cut in  $G'$
- Max flow in  $G'$  = size of max matching in G. Since size of max matching  $< |L|$  (since G has no perfect matching),  $\text{cap}(A, B) < |L|$
- Let  $L_A = L \cap A$ ,  $L_B = L \cap B$ ,  $R_A = R \cap A$
- $\text{cap}(A, B) = |L_B| + |R_A|$ 
  - The outgoing edges from cut A must have flow = 1 because  $\text{cap}(A, B) < |L|$ . Hence, no edge in G (from L to R, with infinite capacity) will be part of the set of edges across the cut
  - The remaining edges are those from s to L or from R to t.
  - Edges from s to L across the cut are in  $L_B$ . Edges from R to t across the cut are from  $R_A$
- As established in 5a, all neighbours of nodes in  $L_A$  must be within A, so  $|N(L_A)| \leq |R_A|$
- $|N(L_A)| \leq |R_A| = \text{cap}(A, B) - |L_B| < |L| - |L_B| = |L_A|$
- Hence,  $|N(L_A)| < |L_A|$ . Choose  $S = L_A$

## Intractability

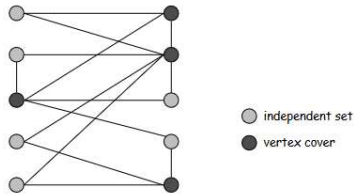
### Polynomial-Time Reduction

- X polynomial reduces to Y if arbitrary instance of X can be solved using polynomial number of calls to oracle that solved Y and polynomial number of standard computational steps (i.e.  $X \leq_p Y$ )
- If X reduces to Y, then Y is at least as hard as X because if Y can be solved in polynomial time, then X can be solved in polynomial time
- Intractability:** If  $X \leq_p Y$  and X cannot be solved in polynomial time, then Y cannot be solved in polynomial time
- Equivalence:** If  $X \leq_p Y$  and  $Y \leq_p X$ , then  $X \equiv_p Y$

### Reduction by simple equivalence

e.g. Independent set  $\equiv_p$  Vertex cover

- Independent set** = Given graph G, is there a subset of vertices S such that  $|S| \geq k$  and there is no edge between each node in S
- Vertex cover** = Given graph G, is there a subset of vertices S such that  $|S| \leq k$  and for each edge, at least one of its endpoints is in S



**Proof** (S is an independent set iff  $V - S$  is a vertex cover)

$\Rightarrow$ :

- Let S be an independent set
- Let  $(u, v)$  be an arbitrary edge
- Since S is an independent set,  $u \notin S$  or  $v \notin S \rightarrow u \in V - S$  or  $v \in V - S$
- Hence,  $V - S$  covers  $(u, v)$

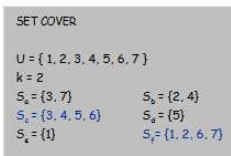
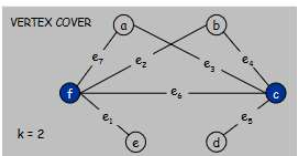
$\Leftarrow$ :

- Let  $V - S$  be a vertex cover
- Consider two nodes  $u \in S$  and  $v \in S$  (so  $u \notin V - S$  and  $v \notin V - S$ )
- If u and v are joined by an edge  $(u, v)$ , then either one of u or v must be in  $V - S$ , which contradicts (2).
- Hence, no two nodes in S are joined by an edge  $\rightarrow$  S is an independent set
- Therefore,  $\exists$  Independent set  $\geq k$  iff  $\exists$  Vertex cover  $< k$

### Reduction from special case to general case

e.g. Vertex cover  $\leq_p$  Set cover

- Set cover:** Given a set U of elements, a collection  $S_1, S_2, \dots, S_m$  of subsets of U, is there a collection of  $\leq k$  of these sets whose union is equal to U



**Proof** (Vertex cover reduces  $\leq_p$  Set cover)

- Construct specific set cover instance from graph G:
  - Define set U to be set of all edges E in G and each subset of edges  $S_v = \{e \in E: e \text{ is incident to } v\}$
- There exists a set cover of size  $\leq k$  iff there exists a vertex cover of size  $\leq k$

### Reduction via "gadgets"

e.g. 3-SAT  $\leq_p$  Independent set

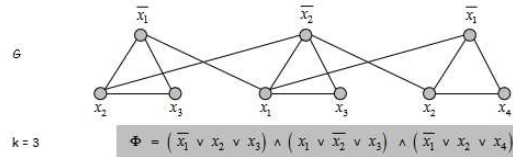
#### Definitions:

- Literal:** Boolean variable or its negation  $(x, \bar{x})$
- Clause:** Disjunction of literals  $(x_1 \vee \bar{x}_2 \vee x_3)$
- Conjunction normal form (CNF,  $\Phi$ ):** Conjunction of clauses  $(x_1 \vee \bar{x}_2 \vee x_3) \wedge (x_4 \vee x_5) \wedge (x_6 \vee \bar{x}_7)$
- 3-SAT:** Given  $\Phi$  where each clause contains exactly 3 literals, is there a satisfying truth assignment

**Proof** (3-SAT  $\leq_p$  Independent set)

Show G has an independent set S of  $|S| = k = |\Phi|$  iff  $\Phi$  of 3-SAT is satisfiable

- Construct specific graph G
  - G contains 3 vertices for each clause, one for each literal
  - Connect the 3 literals in a clause in a triangle via an edge
  - Connect literal to each of its negations via an edge
- $\Rightarrow$ :
  - Let  $S \subseteq G$  be an independent set of size k
  - S must contain at least one vertex from each triangle because  $k = |\Phi|$
  - S must contain at most one vertex from each triangle, otherwise there would be an edge between two vertices
  - Hence, S must contain exactly one vertex from each triangle
  - Furthermore, if  $u, v \in S$ , then u and v cannot be negations of each other otherwise there will be an edge  $(u, v)$  in G
  - Set these literals in S to be true
  - Then, the truth assignment will be consistent and all clauses are satisfied
- $\Leftarrow$ :
  - Given a satisfying assignment, at least one literal from each triangle must be true so that every clause is satisfied
  - Pick one literal from each triangle to form an independent set S of size k



## NP-completeness

### Definitions

- P:** Decision problem X with poly-time algorithms
- NP:** Decision problems X with a poly-time certifier
  - Certifier  $C(s, t)$  returns 'yes' for some certificate  $|t| \leq p(|s|)$
- NP-complete:**  $X \in NP$  and  $\forall Y \in NP, Y \leq_p X$ . Problem X is NP-complete if X is at least as hard as every NP problem

### Properties

**Proof:**  $P \subseteq NP$

- Given  $X \in P$ , there exists a poly-time algorithm A that returns  $A(s)$ .
- Construct a certifier  $C(s, t)$  that returns  $A(s)$  and set  $t = \emptyset$
- Then,  $C(s, t)$  runs in poly-time and  $|t| \leq p(|s|)$

**Proof:** NP-complete problem X is solvable in polynomial time iff  $P = NP$

$\Rightarrow$ : X is in NP by definition. X is solvable in poly-time. Hence,  $NP \subseteq P$

$\Leftarrow$ :  $P = NP$ . X is in NP by definition. Hence, X is solvable in poly-time

**Proof:** If  $\exists X \in NP$  s.t. X cannot be solved in polynomial time, then no NP-complete problem can be solved in polynomial time.

- Contrapositive of the above proof: If  $P = NP$ , then no NP-complete problem is solvable in polynomial time

### Polynomial Transformations

- Cook reduction:** Problem X polynomial reduces to problem Y if X can be solved using (i) polynomial # of standard computational steps and (ii) polynomial # of calls to oracle that solves Y
- Karp reduction:** Problem X polynomial transforms to problem Y if given any input x to X, we can construct an input y s.t. x is a 'yes' instance of X iff y is a 'yes' instance of Y

### Proof X is NP-complete

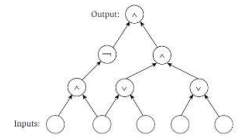
- If Y is NP-complete,  $X \in NP$  and  $Y \leq_p X$ , then X is NP-complete:
  - Take an  $Z \in NP$ , then  $Z \leq_p Y \leq_p X$ . Hence,  $Z \leq_p X$ .
  - By defn of NP-completeness, X is NP-complete
- Hence, it suffices to show:

Show that  $X \in NP$

Choose a known NP-complete problem Y

Prove  $Y \leq_p X$ : (Karp reduction)

- Construct an arbitrary instance  $s_Y$  of Y and construct instance  $s_X$  of X in polynomial time
- Show that  $s_Y = \text{'yes'}$  iff  $s_X = \text{'yes'}$



### Circuit SAT $\leq_p$ 3-SAT

- Circuit SAT: Known NP-complete problem
- Reduce any  $X \in NP$  to circuit SAT:
  - Set first n source nodes to be the n-bit s and remaining  $p(|s|)$  source nodes to represent the bits in t
  - X outputs 'yes' iff the  $\exists t$  s.t.  $|t| \leq p(|s|)$  and the circuit is satisfiable (i.e.  $C(s, t) = \text{'yes'}$ )
- Show 3-SAT is NP-complete by reducing Circuit SAT to 3-SAT



### Construction (Circuit SAT $\leq_p$ 3-SAT)

- For each node  $v$  in Circuit SAT, let  $x_v$  be a variable in 3-SAT
- Note that  $P \Rightarrow Q \Rightarrow \neg P \vee Q$
- NOT:  $x_v \Leftrightarrow \bar{x}_u = (x_v \Rightarrow \bar{x}_u) \wedge (\bar{x}_u \Rightarrow x_v) = (\bar{x}_v \vee \bar{x}_u) \wedge (x_v \vee x_u)$
- AND:  $x_v \Leftrightarrow (x_u \wedge x_w) = \dots = (\bar{x}_v \vee x_u) \wedge (\bar{x}_v \vee x_w) \wedge (x_v \vee \bar{x}_u \vee \bar{x}_w)$
- OR:  $x_v \Leftrightarrow (x_u \vee x_w) = \dots = (x_v \vee \bar{x}_u) \wedge (x_v \vee \bar{x}_w) \wedge (\bar{x}_v \vee x_u \vee x_w)$
- SOURCE: Assign it to 0 or 1 if it is a hardcoded input
- OUTPUT: Add variable  $x_o = 1$  for output
- Some clauses have  $< 3$  variables. Create 4 new variables  $z_1, z_2, z_3, z_4$  and add the clauses  $(\bar{z}_i \vee z_3 \vee z_4), (\bar{z}_i \vee \bar{z}_3 \vee z_4), (\bar{z}_i \vee z_3 \vee \bar{z}_4), (\bar{z}_i \vee \bar{z}_3 \vee \bar{z}_4)$  for each  $i = 1, i = 2$ . This ensures that  $z_1 = z_2 = 0$
- If a clause has 1 variable  $t$ , replace it with  $(t \vee z_1 \vee z_2)$
- If a clause has 2 variables  $(s \vee t)$ , replace it with  $(s \vee t \vee z_1)$

### Proof

$\Rightarrow$  Suppose Circuit SAT is satisfiable

- The satisfying assignment to Circuit SAT will create values at all nodes of the circuit
  - This set of values will satisfy the constructed SAT instance
- $\Leftarrow$  Suppose 3-SAT is satisfiable
- The clauses in 3-SAT ensure that the values assigned to all nodes of the circuit are the same as what the circuit computes for these nodes.
  - $x_o = 1$  in 3-SAT, so the assignment is satisfiable in Circuit SAT

### 3-SAT $\leq_p$ Hamiltonian Cycle

#### Construction

- Assume a 3-SAT instance  $\Phi$  with  $n$  variables and  $k$  clauses,
- For each variable  $x_i$  in 3-SAT, construct a path  $P_i = v_{i,1}, v_{i,2}, \dots, v_{i,b}$  where  $b = 3k + 3$ . Edges in the path are bi-directional
- For each path,
  1. Add edges from  $v_{i,1}$  to  $v_{i+1,1}$  and  $v_{i+1,b}$
  2. Add edges from  $v_{i,b}$  to  $v_{i+1,1}$  and  $v_{i+1,b}$
- Add  $s$  and edges from  $s$  to  $v_{1,1}$  and  $v_{1,b}$
- Add  $t$  and edges from  $v_{n,1}$  and  $v_{n,b}$  to  $t$
- This models a Hamiltonian Cycle with  $2^n$  paths (each 'layer' can be traversed from L2R or R2L, independent of other 'layers'). Similarly, there are  $2^n$  different assignments in a 3-SAT instance
- If a path  $P_i$  is traversed from L2R, set  $x_i = 1$ . Else  $x_i = 0$
- Add an extra node  $c_j$  for each clause  $C_j$ 
  1. For each variable  $x_i$  in  $C_j$ , add the edges  $x_{i,3j} \rightarrow c_j$  and  $c_j \rightarrow x_{i,3j+1}$  if  $x_i$  is not a negation. Otherwise, add the edges  $x_{i,3j+1} \rightarrow c_j$  and  $c_j \rightarrow x_{i,3j}$
  2. e.g. if  $C_1 = x_1 \vee \bar{x}_2 \vee x_3$ ,  $P_1$  must go from L2R OR  $P_2$  must go from R2L OR  $P_3$  must go from L2R in order for  $c_j$  to be visited

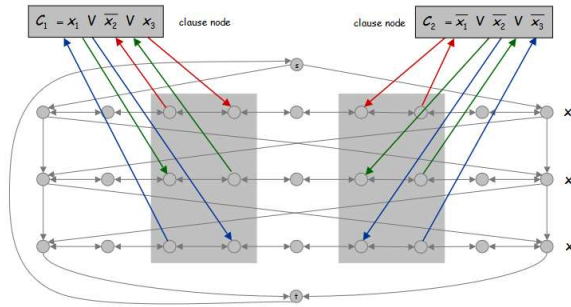
### Proof

$\Rightarrow$  Suppose  $\Phi$  is satisfiable

- If an arbitrary variable  $x_i$  is 1, then  $P_i$  traverses from L2R, otherwise it traverses from R2L.
- For each clause  $C_j$ , since it is satisfied, there must be a path  $P_i$  that traverses in the "correct" direction so  $c_j$  can be spliced into the cycle via edges incident on  $v_{i,3j}$  and  $v_{i,3j+1}$

$\Leftarrow$  Suppose there exists a Hamiltonian Cycle  $\mathcal{C}$

- Then all  $c_j$  must be visited
- If  $\mathcal{C}$  enters a node  $c_j$  from  $v_{i,3j}$ , it must immediately depart on an edge to  $v_{i,3j+1}$  otherwise  $v_{i,3j+2}$  will never be visited without breaking the Hamiltonian property
- Symmetrically, if  $\mathcal{C}$  enters a node  $c_j$  from  $v_{i,3j+1}$ , it must immediately depart on an edge to  $v_{i,3j}$  otherwise  $v_{i,3j-1}$  will never be visited without breaking the Hamiltonian property
- Hence, for each  $c_j$ , the nodes before and after  $c_j$  in  $\mathcal{C}$  are joined by an edge in  $G$ .
- Therefore, we can remove each  $c_j$  in  $\mathcal{C}$  and join  $v_{i,3j}$  and  $v_{i,3j+1}$  via an edge, forming a Hamiltonian cycle  $\mathcal{C}'$  of  $G' = G - \{c_1, \dots, c_k\}$
- Any Hamiltonian cycle in  $G'$  must traverse each  $P_i$  in only one direction. Hence, each  $P_i$  in  $\mathcal{C}'$  must traverse fully in only one direction
- If  $\mathcal{C}'$  traverses  $P_i$  from L2R, we can set  $x_i = 1$ , otherwise  $x_i = 0$
- Since the larger cycle  $\mathcal{C}$  was able to visit all  $c_j$ , at least one of the paths was traversed in the correct direction relative to  $c_j$
- Hence, the truth assignment satisfies all the clauses.



### Hamiltonian Cycle $\leq_p$ Traveling Salesman Problem

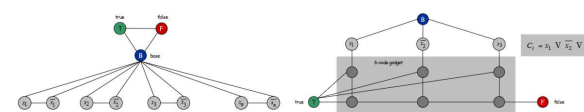
#### Construction

- Given graph  $G(V, E)$ , create  $|V|$  cities with distance function
$$d(u, v) \begin{cases} 1, & (u, v) \in E \\ 2, & (u, v) \notin E \end{cases}$$
- Then,  $G$  will have a TSP tour of length  $\leq |V|$  iff  $G$  is Hamiltonian

### 3-SAT $\leq_p$ 3-Colorability

#### Construction

- For each variable  $x_i$ , create nodes  $x_i$  and  $\bar{x}_i$  and join them via an edge
- Create 3 extra nodes, T, F and B that connect to one another. Connect B to every literal
- For each clause, create a 6-node gadget:



- $\Phi$  is satisfiable iff the constructed graph is 3-colorable

### Proof

$\Rightarrow$  Suppose  $\Phi$  is satisfiable

- Color all true literals green
- Color nodes below green nodes to be red and the node below blue
- Color remaining middle row nodes blue (first layer of nodes in gadget)
- Color remaining bottom nodes red/green as forced
- The resulting graph is 3-colorable

$\Leftarrow$  Suppose the graph is 3-colorable

- Assign each literal coloured green to be true
- Each variable must have one literal to be green and the other to be red (refer to left image). Hence, the assignment is consistent
- No literals can be blue, hence each literal is assigned T or F
- At least one literal in any clause will be true (green) to satisfy 3-colorability (refer to right graph)
- Hence,  $\Phi$  is satisfiable

### 3-SAT $\leq_p$ Subset-Sum

#### Construction

- Given  $\Phi$  with  $n$  variables and  $k$  clauses, create  $2n + 2k$  integers, each with  $n + k$  digits
- For each integer, the first  $n$  digits represent each variable and the last  $k$  digits represent each clause
- For each variable  $x_i$ , create 2 integers for  $x_i$  and  $\bar{x}_i$ . For both integers, set digit  $i$  to 1. For each clause digit, set it to 1 if the literal is part of the clause.
- For each clause  $C_j$ , create  $2k$  integers. Set first  $n$  digits to be 0. Set digit  $j$  to 1 for one of the integers and 2 for the other.
- Then,  $\Phi$  is satisfiable iff there exists a subset that sums to  $W = 11\dots144\dots4$  ( $n$  1s and  $k$  4s)

### Proof

$\Rightarrow$  Suppose  $\Phi$  is satisfiable

- For each variable  $x_i$ , choose one of the two integers representing it (depending on where  $x_i$  is T or F)
- There will be  $n$  of such integers.
- For each clause  $C_j$ , choose one or both of the two integers representing it such that the digit  $j$  of the resulting sum is 4
- Sum the  $n + k$  digits to get  $S$
- The first  $n$  digits of  $S$  will be 1 because exactly one integer is picked for each integer
- The last  $k$  digits of  $S$  will be 4; because  $\Phi$  is satisfiable, digit  $j$  has value 1, 2 or 3 from the sum of the  $n$  integers in (2). We can always pick some 1 or 2 of the dummy rows to make it up to 4
- Take the selected integers. Then, there exists a subset sum equals to  $W$

$\Leftarrow$  Suppose there exists a subset sum equals to  $W$

- Let the subset of integers that sum to  $W$  be  $S$
- Then, exactly one of the integers in  $S$  must have digit  $i \leq n$  set to 1. To achieve this, exactly one of the literals for each variable in  $\Phi$  must be true (therefore, there is consistent assignment)
- Each clause  $C_j$  is true in order for each digit  $j \geq k$  to be 4
- Hence,  $\Phi$  is satisfied

	x	y	z	C <sub>1</sub>	C <sub>2</sub>	C <sub>3</sub>
x	1	0	0	0	1	0
~x	1	0	0	1	0	1
y	0	1	0	1	0	0
~y	0	1	0	0	1	1
z	0	0	1	1	1	0
~z	0	0	1	0	0	1
	0	0	0	1	0	0
	0	0	0	2	0	0
	0	0	0	0	1	0
	0	0	0	0	0	2
	0	0	0	0	0	1
	0	0	0	0	0	2
W	1	1	1	4	4	4

dominate to get close columns to sum to 4

## Approximation Algorithms

### Load Balancing (2-approx)

- Given  $m$  machines and  $n$  jobs, each job  $j$  with processing time  $t_j$ . Each job must run contiguously on one machine and each machine processes at most one job at a time. Load  $L_i = \sum t_i$  for each machine. Assign the jobs such that maximum load on any machine (makespan) is minimised
- Greedy solution: Assign job  $j$  to machine  $i$  with the smallest load so far. Then, the makespan  $L \leq 2L^*$  where  $L^*$  is the optimal makespan
- Optimisation: Sort the jobs in descending order of processing time. Then, the algorithm is a 1.5-approximation

#### Proof (w/o sorting)

- $L^* \geq \max t_j$  since some machine must process the most time-consuming job
- $L^* \geq \text{average load} = \frac{1}{m} \sum t_j$
- Let machine  $i$  be the machine with the maximum load  $L_i$  and job  $j$  be the last job assigned to it.
- Then, machine  $i$  must have had the smallest load before this assignment. i.e.  $L_i - t_j \leq L_k$  for any machine  $L_k$
- $L_i - t_j \leq \frac{1}{m} \sum_k L_k = \frac{1}{m} \sum_k t_k \leq L^*$
- Hence,  $L_i = (L_i - t_j) + t_j \leq L^* + \frac{1}{2} L^* = \frac{3}{2} L^*$

#### Proof (with sorting)

- If there are  $\leq m$  jobs, then the greedy solution is optimal (i.e.  $L = L^*$ )
- If there are  $> m$  jobs, consider the first  $m + 1$  jobs  $t_1, t_2, \dots, t_{m+1}$
- Each job takes at least  $t_{m+1}$  time since they are sorted in desc order
- There are  $m + 1$  jobs and  $m$  machines. By pigeonhole principle, at least one machine gets 2 jobs.
- Therefore,  $L^* \geq 2t_{m+1}$
- Hence,  $L_i = (L_i - t_j) + t_j \leq L^* + \frac{1}{2} L^* = \frac{3}{2} L^*$

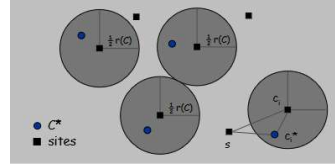
### Center Selection Problem (2-approx)

- Given  $n$  sites, select  $k$  centers  $C$  so that the max distance  $r(C)$  from a site to nearest center is minimised
- Greedy solution: Repeatedly choose the next center to be the site farthest away from any existing center (i.e.  $\text{dist}(s, C)$  is the largest). Let first center to be at any arbitrary site. Then,  $r(C) \leq r(C^*)$  where  $C^*$  is the set of optimal centers

#### Proof

- Each iteration of the greedy algorithm reduces  $\text{dist}(s, C)$  towards  $r(C)$ . Hence, upon termination, all centers in  $C$  are pairwise at least  $r(C)$  apart

- Assume (for contradiction) that  $r(C) > 2r(C^*) \Rightarrow r(C^*) < \frac{1}{2} r(C)$
- For each center  $c_i$  in  $C$ , consider a ball of radius  $\frac{1}{2} r(C)$  around it
- Because of (1), one ball will have exactly one  $c_i$ . Because of (2), one ball will have exactly one  $c_i^*$ 
  - If  $\exists$  a ball with no  $c_i^*$ , then  $\text{dist}(c_i, C^*) > \frac{1}{2} r(C)$ , a contradiction
  - If  $\exists$  a ball with multiple  $c_i^*$ , then there will  $\exists$  a ball with no  $c_i^*$



- Consider any site  $s$  and its closest center  $c_i^*$  in  $C^*$
- Then,  $\text{dist}(s, C) \leq \text{dist}(s, c_i) \leq \text{dist}(s, c_i^*) + \text{dist}(c_i^*, c_i) \leq 2r(C^*)$ , a contradiction
- Hence,  $r(C) \leq 2r(C^*)$

## Randomized Algorithms

- Monte Carlo:** Guaranteed poly-time, likely correct answer
- Las Vegas:** Guaranteed correct answer, likely poly-time

### Contention Resolution (Monte Carlo)

- Given  $n$  processes  $P_1, P_2, \dots, P_n$  and a single resource  $r$  which can only be accessed by at most one process at any given time, devise a protocol to ensure that all processes can access  $r$  on a regular basis
- Randomized Protocol: Each process requests access to  $r$  at time  $t$  with probability  $p = 1/n$

#### Proof

- Let  $S[i, t]$  = the event that  $P_i$  successfully accesses  $r$  at time  $t$
- Then,  $\Pr(S[i, t]) = p(1 - p)^{n-1} = \frac{1}{n} \left(1 - \frac{1}{n}\right)^{n-1}$
- $\frac{1}{e} \leq \left(1 - \frac{1}{n}\right)^{n-1} \leq \frac{1}{2}$ . Therefore,  $\frac{1}{en} \leq \Pr(S[i, t]) \leq \frac{1}{2n}$
- Let  $F[i, t]$  = the event that  $P_i$  fails to access  $r$  in rounds 1 through  $t$
- Then,  $\Pr(F[i, t]) = (1 - \Pr(S[i, t]))^t \leq \left(1 - \frac{1}{en}\right)^t$
- Choose  $t = \lceil en \rceil$ . Then,  $\Pr(F[i, t]) \leq \left(1 - \frac{1}{en}\right)^{\lceil en \rceil} \leq \left(1 - \frac{1}{en}\right)^{en} \leq \frac{1}{e}$
- Choose  $t = \lceil en \rceil \lceil \ln n \rceil$ . Then,  $\Pr(F[i, t]) \leq \left(\frac{1}{e}\right)^{\lceil \ln n \rceil} = n^{-c}$
- Let  $F[t]$  = the event that at least one of the  $n$  processes fails to access  $r$  in any of the rounds 1 through  $t$
- Then,  $\Pr(F[t]) = \Pr(\cup_{i=1}^n F[i, t]) \leq \sum_{i=1}^n \Pr(F[i, t]) \leq n \left(1 - \frac{1}{en}\right)^t$ 
  - Union Bound: The probability that  $\geq 1$  of the events happens is at most the sum of the probabilities of the individual events
- Choose  $t = \lceil en \rceil \lceil 2 \ln n \rceil$ . Then,  $\Pr(F[t]) \leq n \cdot n^{-2} = \frac{1}{n}$
- Therefore, the probability that all the processes succeeds to access  $r$  within rounds  $2en \lceil \ln n \rceil$  rounds is at least  $1 - \Pr(F[t]) = \frac{n-1}{n}$

### Global Minimum Cut (Monte Carlo)

- Given a connected, undirected and unweighted graph  $G(V, E)$ , find a cut  $(A, B)$  of minimum cardinality (least # edges across the cut)
- Contraction Algorithm:
  - Pick any edge  $(u, v)$  at random
  - Contract edge  $(u, v)$ : Replace vertices  $u$  and  $v$  with new super-node  $w$  and remove all edges between  $u$  and  $v$
  - Repeat until graph has just two nodes  $v_1$  and  $v_2$
  - Return the cut (all vertices that were contracted to form  $v_1$ )
  - Then, the probability that this algorithm returns a global min-cut is  $\geq 2/n^2$



#### Proof

- Let  $(A^*, B^*)$  be the optimal global min-cut of and  $F^*$  be the set of edges across the cut. Let  $|F^*| = k$
- The first iteration of the contraction algorithm picks and contracts an edge in  $F^*$  with probability  $\frac{k}{|E|}$
- Suppose  $\exists v \in V$  s.t.  $\text{degree}(v) < k$ . Take  $v$  to form  $A^*$  and the rest of the vertices to form  $B^*$ . Then, the global min-cut  $< k$ , a contradiction. Hence,  $\forall v \in V$ ,  $\text{degree}(v) \geq k$ .
- For any graph,  $|E| = \frac{1}{2} \sum_{v \in V} \text{degree}(v_i)$ . Therefore,  $|E| \geq \frac{1}{2} kn$
- Hence, by point (2), the algorithm contracts an edge in  $F^*$  in the first iteration with probability  $\leq \frac{k}{|E|} \leq \frac{k}{\frac{1}{2} kn} = \frac{2}{n}$
- After  $j$  iterations, the number of vertices remaining is  $n' = n - j$  (each iteration reduces the # vertices by one).
- Suppose no edge in  $F^*$  were contracted in these  $j$  iterations. Then, the min-cut is still  $k$ . Hence,  $|E'| \geq \frac{1}{2} kn'$  and the algorithm contracts an edge in  $F^*$  with probability  $\leq \frac{2}{n'}$
- Let  $E_j$  = event that an edge in  $F^*$  is not contracted in iteration  $j$
- Then,  $\Pr(E_1 \cap E_2 \cap \dots \cap E_{n-2}) = \Pr(E_1) \times \Pr(E_2|E_1) \times \Pr(E_3|E_1 \cap E_2) \times \dots \times \Pr(E_{n-2}|E_1 \cap E_2 \cap \dots \cap E_{n-3}) \geq \left(1 - \frac{2}{n}\right) \left(1 - \frac{2}{n-1}\right) \left(1 - \frac{2}{n-2}\right) \dots \left(1 - \frac{2}{4}\right) \left(1 - \frac{2}{3}\right) = \left(\frac{n-2}{n}\right) \left(\frac{n-3}{n-1}\right) \left(\frac{n-4}{n-2}\right) \dots \left(\frac{2}{4}\right) \left(\frac{1}{3}\right) = \frac{2}{n(n-1)} \geq \frac{2}{n^2}$
- Run the algorithm  $n^2 \ln n$  times with independent random choices. Then, the probability of failing is

$$\leq \left(1 - \frac{2}{n^2}\right)^{n^2 \ln n} = \left(\left(1 - \frac{2}{n^2}\right)^{\frac{1}{2} n^2}\right)^{2 \ln n} \leq e^{-2 \ln n} = \frac{1}{n^2}$$

$$1 + x \leq e^x \Rightarrow 1 - 1/x \leq e^{-1/x} \Rightarrow (1 - 1/x)^x \leq e^{-1}$$

#### Expectation

- $E[X] = \sum_{j=0}^{\infty} j \Pr(X = j)$
- $X = \# \text{ trials until first success} = X \sim G(p) \Rightarrow E[X] = 1/p$
- $X = \text{Bernoulli Trial} = X \sim \text{Bernoulli}(p) \Rightarrow E[X] = p$
- Linearity of expectation (LoE):  $E[X + Y] = E[X] + E[Y]$

### Max 3-SAT (Las Vegas)

- Given  $\Phi$  with  $k$  clauses, the expected # of clauses satisfied by a random assignment is  $7k/8$
- Hence, there must exist a random assignment that satisfies at least  $7k/8$  clauses with probability  $\geq 1/8k$
- Johnson's algorithm: If we repeatedly generate random assignments until one satisfies  $\geq 7k/8$  clauses, the number of trials is  $\leq 8k$

#### Proof

- Let  $X_j$  = RV that equals 1 if clause  $j$  is satisfied and 0 otherwise
- Hence,  $E[X_j] = p = 7/8$  (Bernoulli Trial)
- Let  $X$  = # of satisfied clauses i.e.  $X \sim B(k, 7/8)$ . Then,  $E[X] = E[X_1 + X_2 + \dots + X_k] = E[X_1] + E[X_2] + \dots + E[X_k] = 7k/8$  (LoE)
- Let  $p_j = \Pr[\text{\#satisfied} = j]$  and  $p = \Pr[\text{\#satisfied} \geq 7k/8]$

$$\begin{aligned} \text{Hence,} \\ 7k/8 = E[X] &= \sum_{j>0} j p_j = \sum_{j<7k/8} j p_j + \sum_{j \geq 7k/8} j p_j \\ &\leq \left(\frac{7k}{8} - \frac{1}{8}\right) \sum_{j<7k/8} p_j + k \sum_{j \geq 7k/8} p_j \\ &\leq \left(\frac{7k}{8} - \frac{1}{8}\right) (1) + kp \end{aligned}$$

- Therefore,  $p \geq 1/8k$
- Let  $Y$  = # trials until one assignment satisfies  $\geq 7k/8$  clauses. Then,  $Y \sim G(p \geq 1/8k)$
- Hence,  $E[Y] \leq 1/(1/8k) = 8k$

### Universal Hashing

- Universal class of hash functions** = a set  $H$  of hash functions  $h_i$  s.t.  $\Pr[h(u) = h(v)] \leq 1/n$ ,  $\forall u, v \in U, n = \# \text{buckets}$
- Let  $X$  = # collisions with any  $u \in U$  using  $H$ . Then, for any  $S \subseteq U$  and  $|S| \leq 1/n$ ,  $E_{h \in H}[X] \leq 1$

#### Proof

- Let  $X_v$  = RV that equals 1 if  $v$  collides with  $u$  and 0 otherwise
- Hence,  $E[X_v] = \Pr[h(u) = h(v)] \leq 1/n$
- Let  $X$  = # collisions with any  $u$  i.e.  $X \sim B(|S|, p \leq 1/n)$ . Then,  $E[X] \leq |S|/n \leq 1$

- Designing such a universal class:

- Choose a prime number  $p$ ,  $n \leq p \leq 2n$  (Chebyshev: such a  $p$  exists)
- For each  $\forall u \in U$ , identify a base- $p$  integer of  $r$  digits:  $x = (x_1, x_2, \dots, x_r)$
- Let  $A$  be the set of all possible  $r$ -digit, base- $p$  integers i.e. for each  $a \in A$ ,  $a = (a_1, a_2, \dots, a_r)$ ,  $0 \leq a_i < p$
- Then,  $H = \{ h_a \mid h_a(x) = (\sum_{i=1}^r a_i x_i) \bmod p, a \in A \}$

#### Proof

- Let  $x = (x_1, x_2, \dots, x_r)$  and  $y = (y_1, y_2, \dots, y_r)$ ,  $x \neq y$
- Then,  $\exists j$  s.t.  $x_j \neq y_j$
- $h_a(x) = h_a(y) \Leftrightarrow \sum_{i \neq j} a_i (x_i - y_i) \bmod p$
- Assume vector  $a$  is chosen randomly by first (uniformly) randomly picking each  $a_i$ ,  $i \neq j$  in the range  $[0, p]$ , then picking  $a_j$  at random

- Since  $p$  is prime,  $a_j z = m \bmod p$  has at most (exactly) 1 solution among  $p$  possibilities

- Let  $p$  be prime,  $z \neq 0 \bmod p$  ( $z$  not divisible by  $p$ ).
- Suppose  $\alpha$  and  $\beta$  are 2 different solutions (for contradiction)
- Then,  $\alpha z = m + k_1 p$  and  $\beta z = m + k_2 p$ . So,  $(\alpha - \beta)z = (k_1 - k_2)p = 0 \bmod p$
- Hence,  $(\alpha - \beta)$  is divisible by  $p$  since  $z$  not divisible by  $p$
- Therefore,  $\alpha = \beta$  since  $0 \leq \alpha, \beta < p$  (contradiction)
- Hence,  $\Pr[h_a(x) = h_a(y)] \leq 1/p \leq 1/n$  since  $n \leq p \leq 2n$

### Chernoff Bounds

- Let  $X_1, X_2, \dots, X_n$  be a Bernoulli process (i.e. they are independent 0-1 RVs) and  $X = X_1, X_2, \dots, X_n$ . Then,  $\forall \mu \geq E[X]$  and  $\forall \delta > 0$ ,

$$\Pr[X > (1 + \delta)\mu] < \left( \frac{e^\delta}{(1 + \delta)^{1+\delta}} \right)^\mu$$

and  $\forall \mu \leq E[X]$  and  $\forall \delta \in (0, 1)$ ,

$$\Pr[X < (1 + \delta)\mu] < e^{-\delta^2 \mu / 2}$$

#### Proof (above mean)

- Markov's inequality states that  $\Pr[X > a] \leq E[X]/a$
- Hence,  $\Pr[X > (1 + \delta)\mu] = \Pr[e^{tX} > e^{t(1+\delta)\mu}] \leq e^{-t(1+\delta)\mu} E[e^{tX}]$
- $E[e^{tX}] = E[e^{t \sum_i X_i}] = E[e^{tX_1} e^{tX_2} \dots e^{tX_n}] = \prod_i E[e^{tX_i}]$  (independence:  $X \perp Y \Rightarrow E(XY) = E(X)E(Y)$ )
- Let  $p_i = \Pr[X_i = 1]$ . Then,  $E[e^{tX_i}] = p_i e^t + (1 - p_i) e^0 = 1 + p_i(e^t - 1) \leq e^{p_i(e^t - 1)}$
- Therefore,  $\Pr[X > (1 + \delta)\mu] \leq e^{-t(1+\delta)\mu} \prod_i E[e^{tX_i}] \leq e^{-t(1+\delta)\mu} \prod_i e^{p_i(e^t - 1)} = e^{-t(1+\delta)\mu} e^{\mu(e^t - 1)}$  ( $\mu \geq E[X] = \sum_i p_i$ )
- Take  $t = \ln(1 + \delta)$ . Then,

$$\Pr[X > (1 + \delta)\mu] \leq e^{-\ln(1+\delta)(1+\delta)\mu} e^{\mu(e^{\ln(1+\delta)} - 1)} = \left( \frac{e^\delta}{(1 + \delta)^{1+\delta}} \right)^\mu$$

### Load Balancing

- $m$  jobs arrive in a stream and need to be processed immediately by  $n$  identical processors. Assign jobs to processors uniformly at random without centralised controller (without round-robin, where each processor will receive at most  $\lceil m/n \rceil$  jobs. How likely is it that some processor is assigned "too many" jobs?

#### Analysis

- Let  $X_i$  = # jobs assigned to processor  $i$
- Let  $Y_{ij}$  = RV that equals 1 if job  $j$  is assigned to processor  $i$
- Then,  $E[Y_{ij}] = 1/n$  ( $E(X) = p$ , if  $X \sim \text{Bernoulli}(p)$ )
- Hence,  $X_i = \sum_j Y_{ij} \Rightarrow \mu = E[X_i] = nE[Y_{ij}] = 1$
- By Chernoff bounds, with  $\delta = c - 1$ ,  $\Pr[X_i > c] < \frac{e^{c-1}}{c^c}$  for some  $c$  which will be chosen later
- Let  $\gamma(n) = x$  s.t.  $x^x = n$  and choose  $c = e\gamma(n)$
- Then,  $\Pr[X_i > c] < \frac{e^{c-1}}{c^c} < \left(\frac{e}{c}\right)^c = \left(\frac{1}{\gamma(n)}\right)^{e\gamma(n)} < \left(\frac{1}{\gamma(n)}\right)^{2\gamma(n)} = \frac{1}{n^2}$

- By Union Bound,  $\Pr[\cup_{i=1}^n X_i > c] \leq n \left(\frac{1}{n^2}\right) = \frac{1}{n}$ . That is, the probability that at least one processor received more than  $c = e\gamma(n)$  jobs is  $\frac{1}{n}$ .
- Therefore, the probability that no processors received more than  $c = e\gamma(n) = \Theta(\log n / \log \log n)$  jobs is  $1 - \frac{1}{n}$

- Suppose there are  $m = 16n \ln n$  jobs. Then,  $E[X_i] = 16 \ln n$ . With high probability, every processor will have between half and twice the average load (e.g.  $8 \ln n = \frac{1}{2}\mu \leq X_i \leq 2\mu = 16 \ln n$ )

#### Proof

- By Chernoff bounds, with  $\delta = 1$ ,

$$\Pr[X_i > 2\mu] < \left(\frac{e}{4}\right)^{16 \ln n} < \left(\frac{1}{e}\right)^{2 \ln n} = \frac{1}{n^2}$$

$$\Pr\left[X_i < \frac{1}{2}\mu\right] < e^{-\frac{1}{2}\left(\frac{1}{2}\right)^2 (16 \ln n)} = \frac{1}{n^2}$$

- By Union Bound,  $\Pr[\cup_{i=1}^n X_i > 2\mu] \leq n \left(\frac{1}{n^2}\right) = \frac{1}{n}$  and  $\Pr\left[\cup_{i=1}^n X_i < \frac{1}{2}\mu\right] \leq n \left(\frac{1}{n^2}\right) = \frac{1}{n}$ .
- Therefore, the probability that every processor has load between half and twice the average load is  $\geq 1 - \frac{1}{n} - \frac{1}{n} = 1 - \frac{2}{n}$