

Data Structure for Disjoint Sets

21.1 Disjoint-set operations

Disjoint set data structure:

1. a data structure maintains a collection $S = \{S_1, S_2, \dots, S_k\}$ of disjoint dynamic sets.
2. Each set is identified by a representative, which is some member of the set. In ^①some applications, it doesn't matter which member is used as the representative; we only care that if we ask the representative of a set without modifying the set between the requests, we get the same answer.^② In other applications, there may be a representative rule for choosing the representative, such as choosing the smallest member in the set.
3. The following operations should be supported.

Make-Set(x): create a new set $\{x\}$.

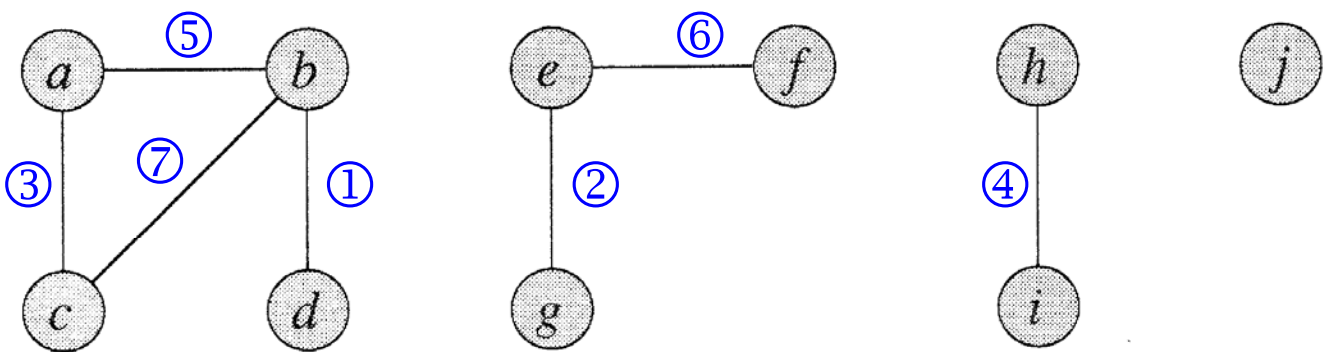
Union(x, y): unite the two sets containing x, y .

Find-Set(x): return a pointer to the representative of the set containing x .

Initially, $S = \emptyset$.

- * n : number of *Make-Set* operations
- m : total number of *Make-Set*, *Union*, and *Find-Set* operation.
- * $m \geq n$ and the number of *Union* operations is at most $n-1$.
 - * $m = n + f + u$
 - * $m \geq n; u \leq n-1$

An application of disjoint-set data structures



(a) 4 connected components

Edge processed	Collection of disjoint sets									
initial sets	{a}	<u>{b}</u>	{c}	<u>{d}</u>	{e}	{f}	{g}	{h}	{i}	{j}
① (b,d)	{a}	{b,d}	{c}		<u>{e}</u>	{f}	<u>{g}</u>	{h}	{i}	{j}
② (e,g)	<u>{a}</u>	{b,d}	<u>{c}</u>		{e,g}	{f}		{h}	{i}	{j}
③ (a,c)	{a,c}	{b,d}			{e,g}	{f}		<u>{h}</u>	<u>{i}</u>	{j}
④ (h,i)	<u>{a,c}</u>	<u>{b,d}</u>			{e,g}	{f}		{h,i}		{j}
⑤ (a,b)	{a,b,c,d}				<u>{e,g}</u>	<u>{f}</u>		{h,i}		{j}
⑥ (e,f)	{a,b,c,d}				{e,f,g}			{h,i}		{j}
★ ⑦ (b,c)	{a,b,c,d}				{e,f,g}			{h,i}		{j}

(b)

CONNECTED-COMPONENTS(G)

```

1  for each vertex  $v \in V[G]$ 
2      do MAKE-SET( $v$ )
3  for each edge  $(u, v) \in E[G]$ 
4      do if FIND-SET( $u$ )  $\neq$  FIND-SET( $v$ )
5          then UNION( $u, v$ )

```

SAME-COMPONENT(u, v)

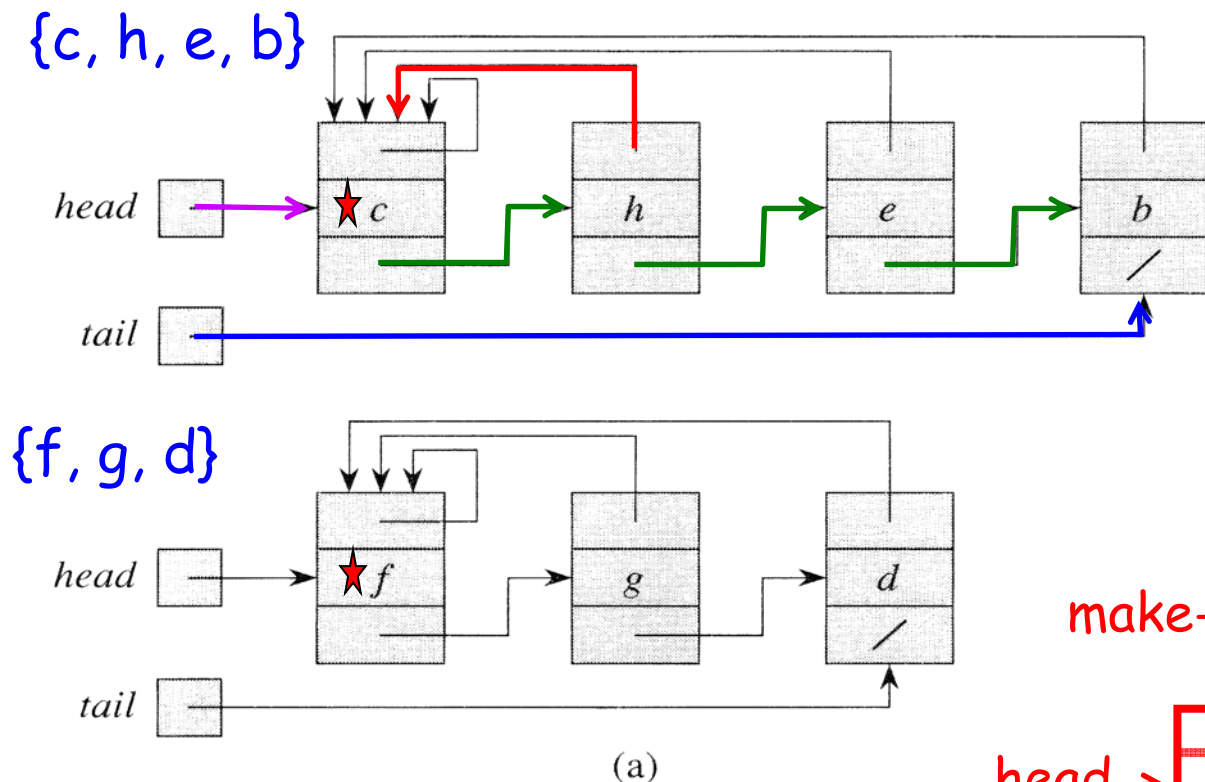
```

1  if FIND-SET( $u$ ) = FIND-SET( $v$ )
2      then return TRUE
3  else return FALSE

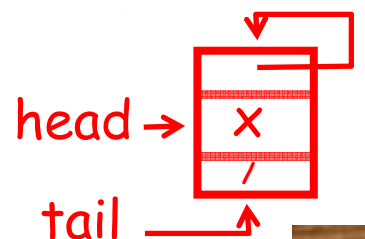
```

21.2 Linked-list representation

(First object in a list is the representative.)



make-set(x)



* Make-Set, Find-Set: $O(1)$ time.

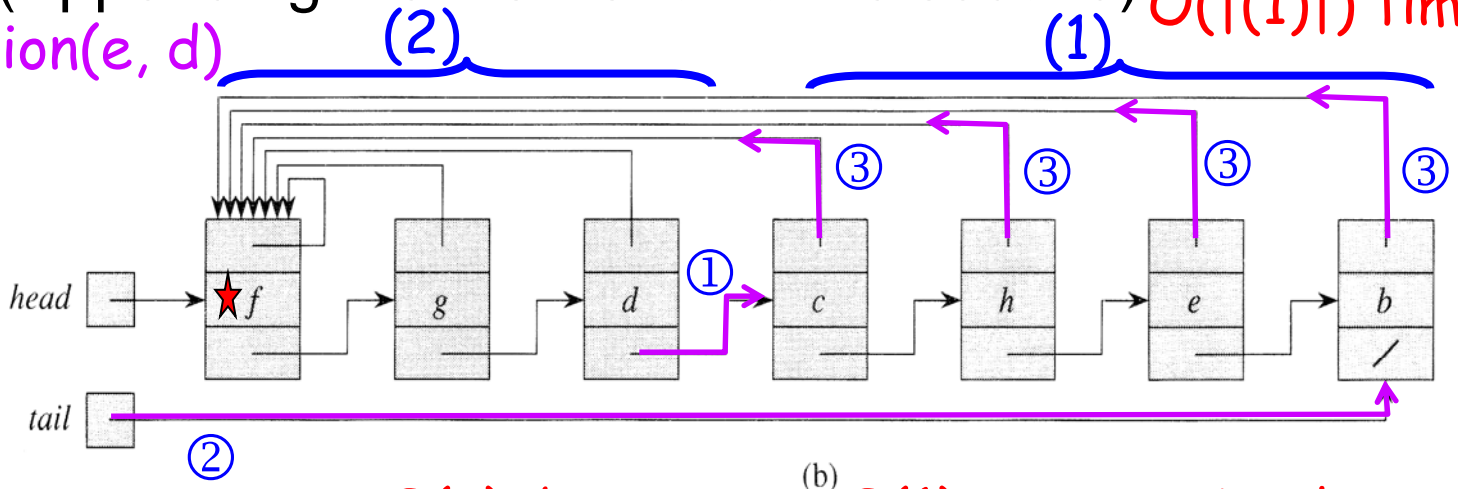
$$\text{Union}((1), (2)) \Rightarrow (2) + (1)$$

21-4

A simple implementation of Union(x, y)

(Appending the first list onto the second)

Union(e, d) $O(|(1)|)$ time



worst-case: $O(n)$, best-case: $O(1) \Rightarrow$ amortized

* $O(n^2)$ time for $m=2n-1$ operations. * $m = n + u$
* $u = n - 1$

* An example showing $\Omega(n^2)$ for $O(n)$ operations.

Operation	Number of objects updated		
MAKE-SET(x_1)		1	} n Make $O(n)$
MAKE-SET(x_2)		1	
\vdots		\vdots	
MAKE-SET(x_n)		1	
UNION(x_1, x_2)	1 + 1	1	} n-1 Union $O(\sum k)$ \parallel $O(n^2)$
UNION(x_2, x_3)	1 + 2	2	
UNION(x_3, x_4)	1 + 3	3	
\vdots		\vdots	
UNION(x_{n-1}, x_n)	1 + (n - 1)	n - 1	

* Thus, the amortized time of each operation is $O(n)$. (tight)

A weighted-union heuristic

1. Each representative stores the length of the list.
2. Append the smaller list onto the longer.

Theorem 21.1: Using the weighted-union heuristic, a sequence of m *Make-Set*, *Union*, and *Find-Set* operations takes $O(m + n \lg n)$ time, where n is the number of *Make-Set* operations.

$$\begin{aligned} * m &= n + f + u \\ * m &\geq n \end{aligned}$$

$$\left\{ \begin{array}{l} \text{Make, FIND: } O(m) \\ \text{Union: } O(n \lg n) \end{array} \right.$$

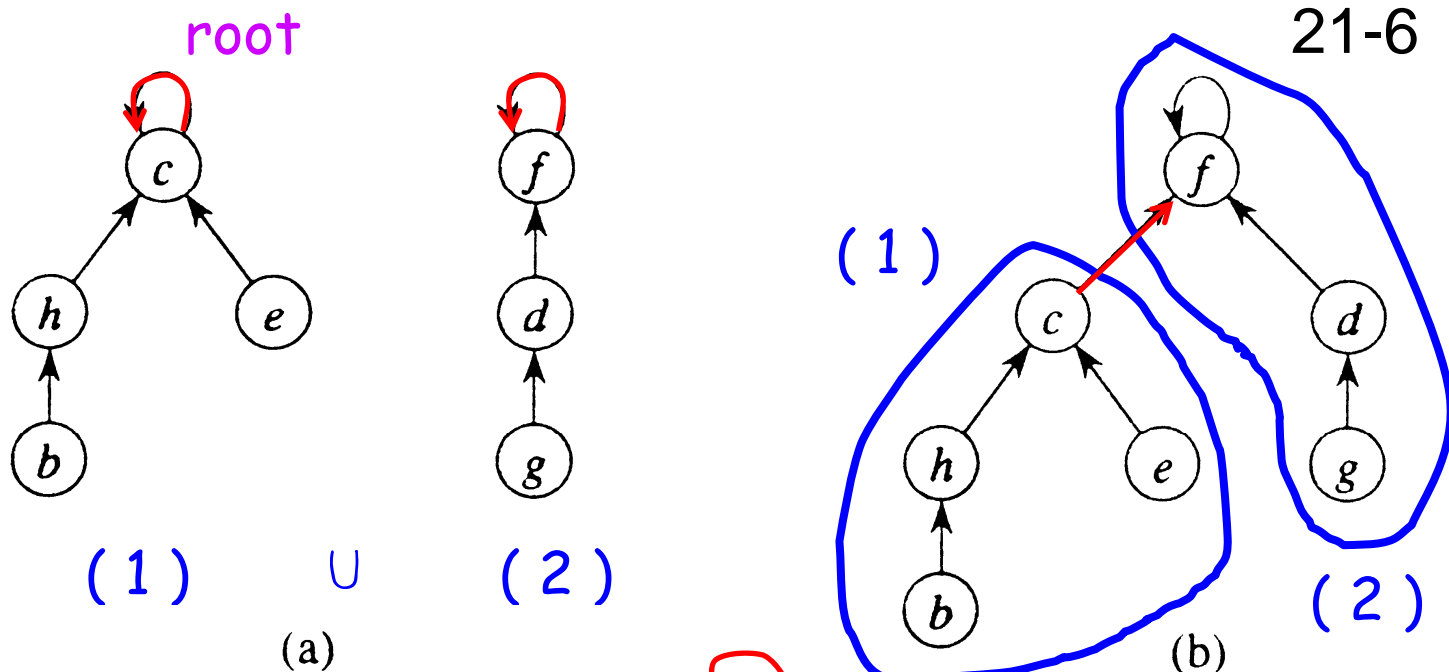
21-5a

Proof: After *Make-Set*(x) is performed, the list containing x has only one element. At the first time x 's representative pointer is updated, the list containing x has at least two elements. Continuing on, we observe that after the k -th time x 's representative is updated, the list containing x has at least 2^k elements. Since $k = O(\lg n)$, the time for all *Union* operations is at most $O(n \lg n)$. The time for each *Make-Set* and *Find-Set* operation is $O(1)$. Thus, the theorem holds. Q.E.D.

Union: (i) best: $O(1)$ (ii) worst-case: $O(n)$
 (iii) amortized: $n \lg n / (n-1) = O(\lg n)$

21.3 Disjoint-set forests

(The root of a tree is the representative.)



Make-Set(x): $O(1)$ time



Find-Set(x): $O(h)$ time, h is the height of the tree containing x . (**Find path: $x \rightarrow \text{root}$**)

Union(x, y): The root of x points to the root of y $\rightarrow O(h)$ time.

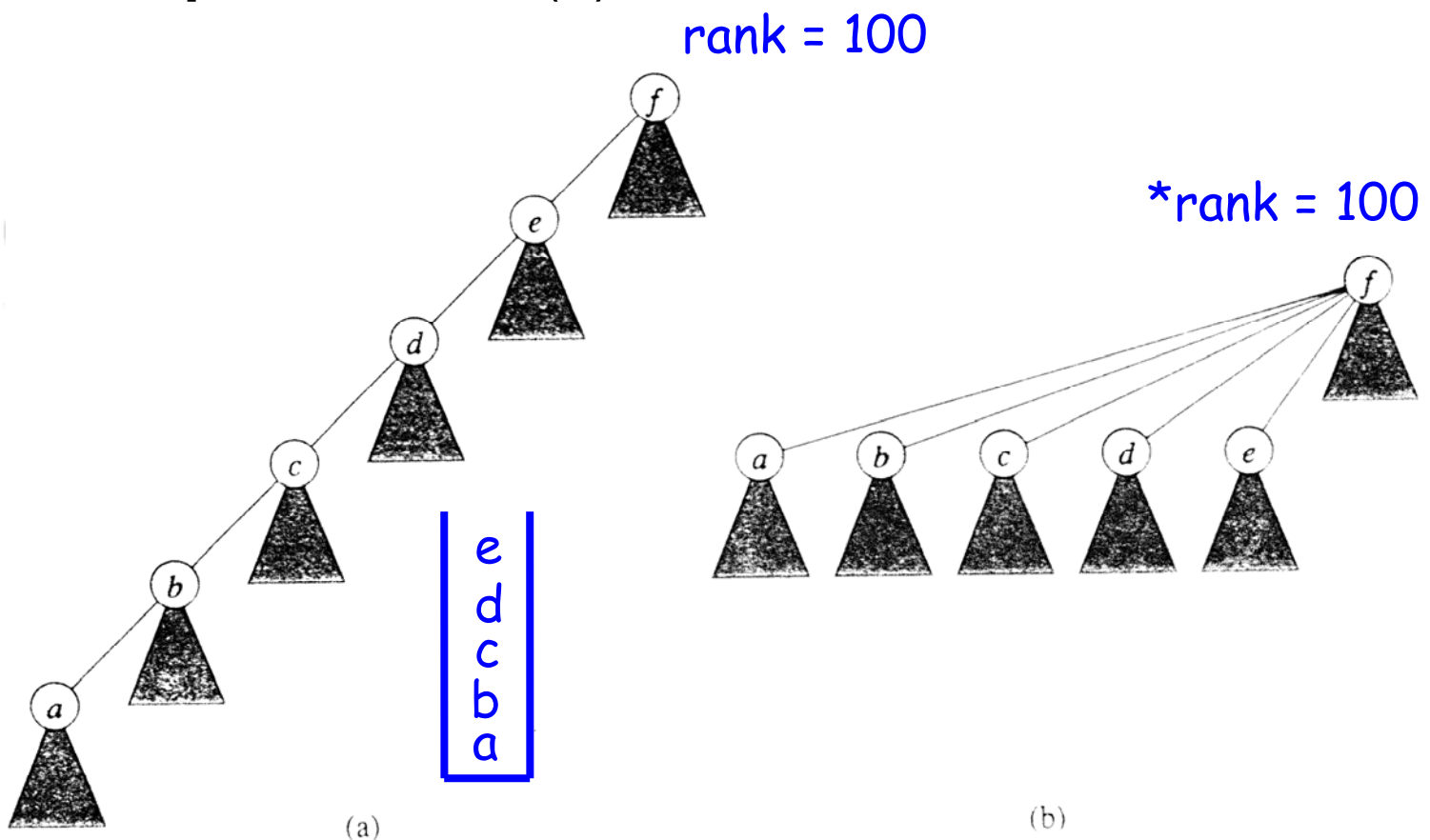
Heuristics to improve the running time

\hookrightarrow height

1. **Union by rank**: the root of the smaller tree points to the root of the larger tree (according to heights).

rank[x]: height of x (number of edges in the longest path between x and a descendant leaf)
 \hookrightarrow may be an approximation (upper bound)

2. **Path compression**: During a Find-Set(x) operation, make each node on the find path point to the root. (It will not change any rank.)

Example: Find-Set(a)**Pseudo-code for disjoint-set forests****MAKE-SET(x)**

- 1 $p[x] \leftarrow x$
- 2 $rank[x] \leftarrow 0$

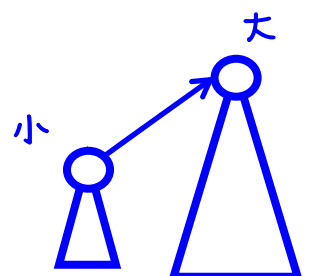
**UNION(x, y)**

- 1 $LINK(FIND-SET(x), FIND-SET(y))$

two Find: path compression

LINK(x, y)

- 1 if $rank[x] > rank[y]$ union by rank
- 2 then $p[y] \leftarrow x$
- 3 else $p[x] \leftarrow y$
- 4 if $rank[x] = rank[y]$
- 5 then $rank[y] \leftarrow rank[y] + 1$



FIND-SET(x)

```

1  if  $x \neq p[x]$     set  $x$ 's parent as the root
2      then  $p[x] \leftarrow \text{FIND-SET}(p[x])$  the root
3  return  $p[x]$  the root

```

21-8x

Effect of the heuristics on the running time

$$m = n + f + u$$

1. If only Union-by-rank is used, it can be easily shown that $O(m \lg n)$ time is required.
($h \leq \lg n$, Ex 21.4-4)
2. If only path-compression is used, it can be shown (not proved here) that the running time is

$$\Theta(n + f \cdot (1 + \log_{2+f/n} n)),$$

* better than 1

where n is the number of *Make-Set* operations and f is the number of *Find-Set* operations.

3. When both heuristics are used, the worst-case running time is $O(m\alpha(n))$, where $\alpha(n)$ is the very slowly growing inverse of Ackermann's function. Since $\alpha(n) \leq 4$ for any conceivable application, we can view the running time as linear in m in all practical situations.

21-10x

almost linear

Amortized:
 $\Rightarrow O(\alpha(n))$ per operation

21-8a

Ackermann's function and its inverse

$$\left. \begin{matrix} 2 \\ \vdots \\ 2 \end{matrix} \right\} i \quad g(0) = 2, g(1) = 2^2, g(2) = 2^{2^2}$$

* Let $g(i) = 2^{\dots^2}$ be a repeated exponentiation.
(e.g., $g(4) = 2^{2^{2^{2^2}}}$.)

* The function $\lg^* n = \min\{i \geq 0: \lg^{(i)} n \leq 1\}$ is
essentially the inverse of $g(i)$. (e.g., $\lg^* 2^{2^{2^{2^2}}} = 5$.)
 $\lg^* g(i) = i + 1$

* The Ackermann's function: for integer $i, j \geq 1$,

$$\begin{aligned} A(1, j) &= 2^j && \text{for } j \geq 1, \\ A(i, 1) &= A(i-1, 2) && \text{for } i \geq 2, \\ A(i, j) &= A(i-1, A(i, j-1)) && \text{for } i, j \geq 2, \end{aligned}$$

	$j = 1$	$j = 2$	$j = 3$	$j = 4$
$i = 1$	2^1	2^2	2^3	2^4
$i = 2$	2^2	2^{2^2}	$2^{2^{2^2}}$	$2^{2^{2^{2^2}}}$
$i = 3$	2^{2^2}	$2^{2^{\dots^2}}_{16}$	$2^{2^{\dots^2}}_{2^{2^{\dots^2}}_{16}}$	$2^{2^{\dots^2}}_{2^{2^{\dots^2}}_{2^{2^{\dots^2}}_{16}}}$
$i = 4$	$\gg 10^{80} \rightarrow A(4, 1) = A(3, 2) = g(16)$			

$$\left. \begin{matrix} 2 \\ \vdots \\ 2 \end{matrix} \right\} j$$

- * Note that $A(2, j) = 2^2 = g(j)$ for all $j \geq 1$.
Thus, $A(i, j) \geq g(j)$ for $i \geq 2$. $A(2, j) = A(1, A(2, j-1)) = 2^{A(2, j-1)}$

- * The inverse of Ackermann's function:

$$\alpha(n) = \min\{i \geq 1: A(i, 1) > \lg n\}.$$

$$\text{e.g., } \alpha(4) = 2, \alpha(32) = 3, \alpha(512) = 3, \alpha(2^{10000}) = 4$$

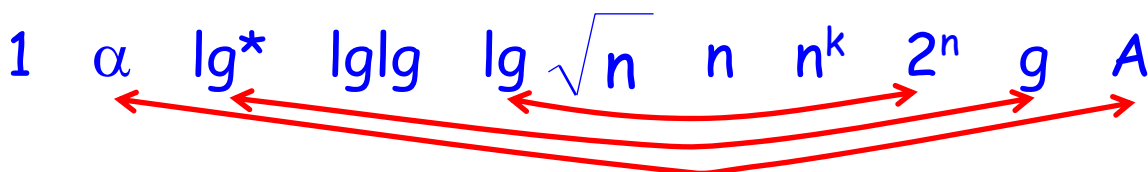
- * $A(4, 1) = A(3, 2) = g(16) \gg 10^{80}$.

- * Since $A(4, 1) \gg 10^{80}$, we have $\alpha(n) \leq 4$ for all practical cases (unless $\lg n > 10^{80}$).
or $n > 2^{10^{80}}$

- * $\lg^* n \leq 5$ for all practical cases (unless $n > 2^{65536}$).

$$= g(4) = 2^{2^{2^2}} = 2^{65536} \approx 10^{1926}$$

- * Since $A(i, 1) \geq g(i)$ for $i \geq 4$, $\alpha(n) = O(\lg^* n)$.



21-10x

Homework: Ex. 21.4-4, Prob. 21-1, 21-3.

21-10a

17-13a