# CS 457, Fall 2019

Drexel University, Department of Computer Science

Lecture 6

# Order Statistics and the Selection Problem

The $i^{\text{th}}$ **order statistic** of a set of $n$ numbers: the $i^{\text{th}}$ smallest number in sorted sequence:

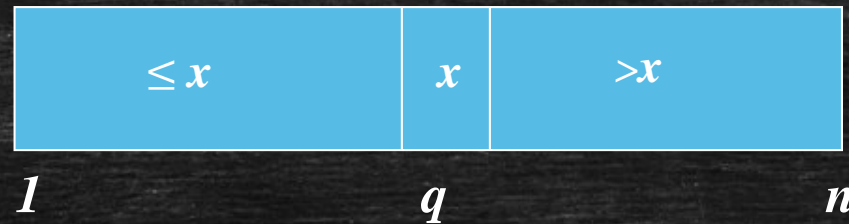| A | 4 | 1 | 3 | 2 | 16 | 9 | 10 | 14 | 8 | 7 |
|---|---|---|---|---|----|---|----|----|---|---|

- **Minimum** or **first order statistic**: **1**
- **Maximum** or **$n^{\text{th}}$ order statistic**: **16**
- **Median** or **$(n/2)^{\text{th}}$ order statistic**: **7 or 8** **(both are medians, when n is even)**

▪ **Selection Problem**
- **Input:** An array **A** of **distinct** numbers of size $n$, and a number $i$
- **Output:** The element $x$ in **A** that is larger than exactly $i-1$ other elements in **A**

- Finding *maximum* and *minimum?*
- Can be easily solved in linear time (i.e., $O(n)$. It's actually $\Theta(n)$)

# Selection Algorithms using Pivot Element

- Choose a pivot element $x$ and partition the subarray $A[1, \ldots, n]$ around it

| $\leq x$ | $x$ | $>x$ |
|---|---|---|
| **1** | **q** | **n** |

- If $q == i$, then $x$ is the $i^{th}$ order statistic

- If $q > i$, then we want the $i^{th}$ order statistic of subarray $[1, \ldots, q-1]$

- If $q < i$, then we want the $(i-q)^{th}$ order statistic of subarray $[q+1, \ldots, n]$

- But, how do we choose this pivot element?

# Simple Selection Algorithm

Select($A, p, r, i$)

1. if $p == r$
2.     return $A[p]$
3. $q$ = Partition($A, p, r$)
4. $k = q - p + 1$
5. if $i == k$
6.     return $A[q]$
7. else if $i \leq k$
8.     Select($A, \ p, \ q - 1, \ i$)
9. else
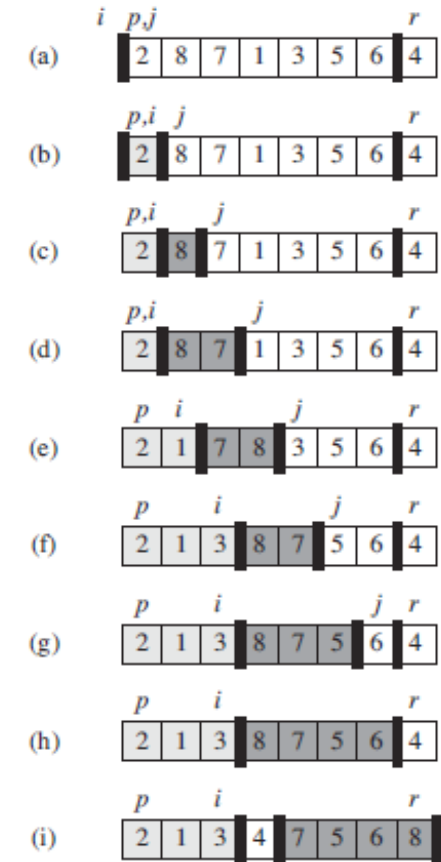10.     Select($A, \ q + 1, \ r, \ i - k$)

Partition ($A, p, r$)

1. $x = A[r]$
2. $i = p - 1$
3. for $j = p$ to $r - 1$
4.     if $A[j] \leq x$
5.         $i = i + 1$
6.         exchange $A[i]$ with $A[j]$
7. exchange $A[i + 1]$ with $A[r]$
8. return $i + 1$

# Partitioning

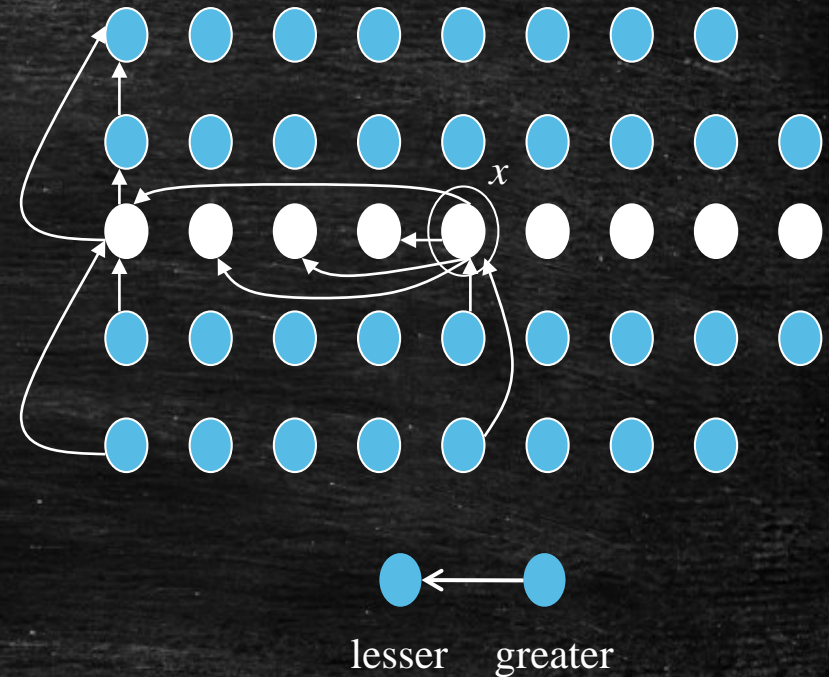Partition $(A, p, r)$

1. $x = A[r]$
2. $i = p - 1$
3. for $j = p$ to $r - 1$
4.     if $A[j] \leq x$
5.         $i = i + 1$
6.             exchange $A[i]$ with $A[j]$
7.     exchange $A[i + 1]$ with $A[r]$
8.     return $i + 1$
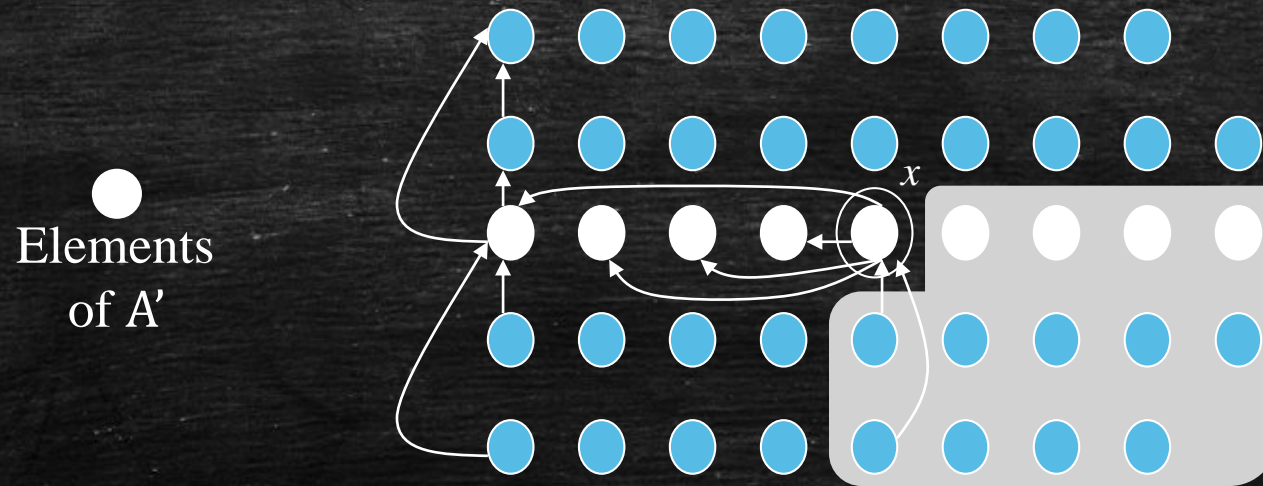
# Worst case linear time selection

**Select(A,p,r,i)**

1. Divide **A** into $n/5$ groups of size **5**.

2. Find the median of each group of **5** by brute force, and store them in a set **A'** of size $n/5$.

3. Recursively use **Select**(A', **1**, $n/5, n/10$) to find the median **x** of $n/5$ medians.

4. Partition elements of **A** around **x**. Let $k$ be the order of **x** found in the partitioning.

5. if $i = k$

6.         return **x**

7. else if $i < k$

8.         **Select**(A, **p**, **q** $-$ **1**, $i$ )

9. else

10.         **Select**(A, **q** $+$ **1**, **r**, $i - k$)



$x$
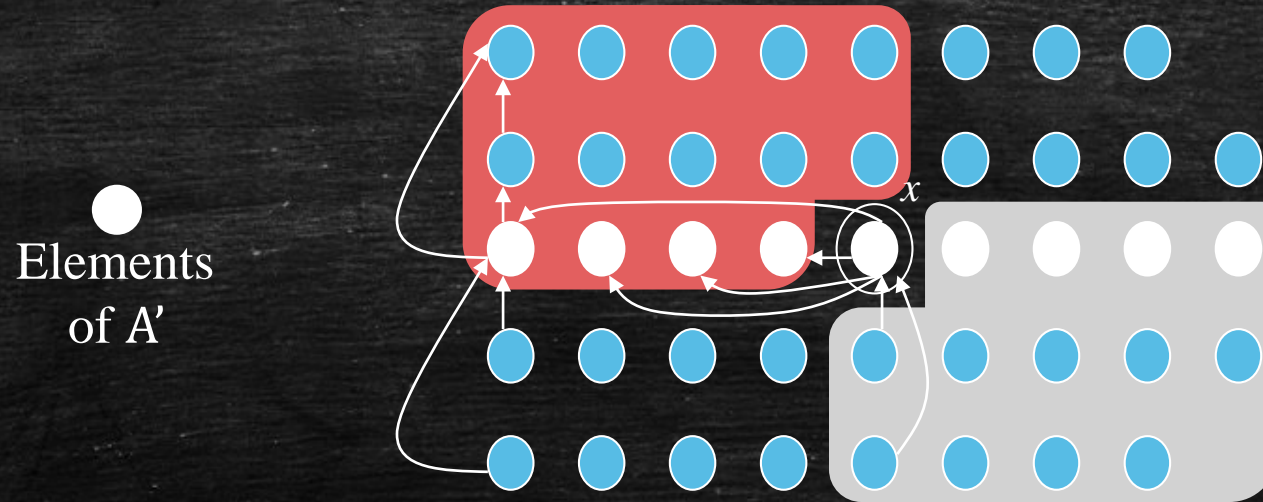
lesser    greater

# Today's Lecture

- More divide and conquer algorithms

- Probabilistic analysis and randomized algorithms

# Analysis



Elements of A'

- At least **half** of the $\lceil n/5 \rceil$ elements in **A'** are $> x$

- Groups whose median $> x$ have at least **3** elements $> x$.

- Therefore, at least $3\left(\left\lceil \frac{1}{2}\left\lceil \frac{n}{5}\right\rceil \right\rceil - 2\right) \geq \frac{3n}{10} - 6$ elements of **A** are $> x$.

# Analysis



- At least **half** of the $\lceil n/5 \rceil$ elements in **A'** are $< x$

- Groups whose median $< x$ have at least **3** elements $< x$.

- Therefore, at least $3 \left( \left\lceil \frac{1}{2} \left\lceil \frac{n}{5} \right\rceil \right\rceil - 2 \right) \geq \frac{3n}{10} - 6$ elements of $A$ are $< x$.

# Worst case linear time selection

**Select(A,p,r,i)**

1. Divide **A** into $n/5$ groups of size **5**.
2. Find the median of each group of **5** by brute force, and store them in a set **A'** of size $n/5$.
3. Recursively use **Select**$(A', \ 1, \ n/5, n/10)$ to find the median **x** of $n/5$ medians.
4. Partition elements of **A** around **x**. Let $k$ be the order of **x** found in the partitioning.
5. if $i = k$
6.      return **x**
7. else if $i < k$
8.      **Select**$(A, \ p, \ q-1, \ i)$
9. else
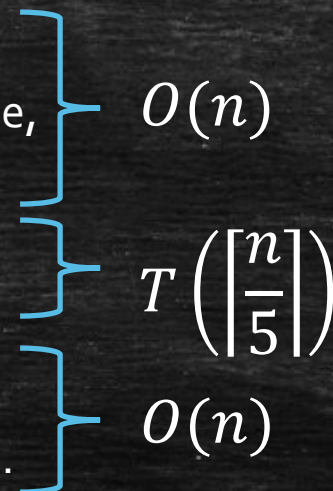10.      **Select**$(A, q+1, r, i-k)$

# Worst case linear time selection

**Select(A,p,r,i)**

1. Divide **A** into $n/5$ groups of size **5**.
2. Find the median of each group of **5** by brute force, $\left.\vphantom{\begin{matrix}a\\b\\c\end{matrix}}\right\}$ $O(n)$
   and store them in a set **A'** of size $n/5$.
3. Recursively use **Select**$(A',\ 1,\ n/5, n/10)$
   to find the median **x** of $n/5$ medians.
4. Partition elements of **A** around **x**. $\left.\vphantom{\begin{matrix}a\\b\end{matrix}}\right\}$ $O(n)$
   Let $k$ be the order of **x** found in the partitioning.
5. if $i = k$
6.       return **x**
7. else if $i < k$
8.       **Select**$(A,\ p,\ q-1,\ i\ )$
9. else
10.       **Select**$(A, q+1, r, i-k)$

# Worst case linear time selection

**Select(A,p,r,i)**

1. Divide **A** into $n/5$ groups of size **5**.
2. Find the median of each group of **5** by brute force, and store them in a set **A'** of size $n/5$.

   $O(n)$

3. Recursively use **Select**$(A', \ 1, \ n/5, n/10)$ to find the median $x$ of $n/5$ medians.

   $T\left(\left\lceil \dfrac{n}{5}\right\rceil\right)$

4. Partition elements of **A** around $x$. Let $k$ be the order of $x$ found in the partitioning.

   $O(n)$

5. if $i = k$
6.      return $x$
7. else if $i < k$
8.      **Select**$(A, \ p, \ q - 1, \ i)$
9. else
10.      **Select**$(A, q + 1, r, i - k)$

# Worst case linear time selection

**Select(A,p,r,i)**

1. Divide **A** into $n/5$ groups of size **5**.
2. Find the median of each group of **5** by brute force, and store them in a set **A'** of size $n/5$.

   $O(n)$

3. Recursively use **Select**$(A', \ 1, \ n/5, n/10)$ to find the median **x** of $n/5$ medians.

   $T\left(\left\lceil\dfrac{n}{5}\right\rceil\right)$

4. Partition elements of **A** around **x**. Let $k$ be the order of **x** found in the partitioning.

   $O(n)$

5. if $i = k$
6.       return **x**
7. else if $i < k$
8.       **Select**$(A, \ p, \ q-1, \ i \ )$
9. else
10.       **Select**$(A, q+1, r, i-k)$

    $T\left(\dfrac{7n}{10}+6\right)$

# Worst case linear time selection

**Select(A,p,r,i)**

1. Divide **A** into $n/5$ groups of size **5**.
2. Find the median of each group of **5** by brute force, and store them in a set **A'** of size $n/5$.

   $O(n)$

3. Recursively use **Select**(A', **1**, $n/5, n/10$) to find the median **x** of $n/5$ medians.

   $T\left(\left\lceil\dfrac{n}{5}\right\rceil\right)$

4. Partition elements of **A** around **x**.
   Let $k$ be the order of **x** found in the partitioning.

   $O(n)$

5. if $i = k$
6.      return **x**
7. else if $i < k$
8.      **Select**(A, **p**, **q** $-$ **1**, $i$ )
9. else
10.      **Select**(A, **q** $+$ **1**, r, $i - k$)

   $T\left(\dfrac{7n}{10}+6\right)$

$$T\left(\left\lceil\dfrac{n}{5}\right\rceil\right) + T\left(\dfrac{7n}{10}+6\right) + O(n)$$

# Worst case linear time selection

Guess that $T(n) = O(n)$ when $T(n) \leq T\left(\left\lceil\frac{n}{5}\right\rceil\right) + T\left(\frac{7n}{10} + 6\right) + O(n)$

Assume $T(n') \leq cn'$ for some constant $c$ and every $n' < n$, and

show that this implies $T(n) \leq cn$ (for the same constant $c$)

Upon substitution, we get $T(n) \leq c\left\lceil\frac{n}{5}\right\rceil + c\left(\frac{7n}{10} + 6\right) + O(n)$

For $n > n_0$ this is at most $c\left(\frac{n}{5} + 1\right) + c\left(\frac{7n}{10} + 6\right) + an$

Therefore, $T(n) \leq cn - \left(\frac{cn}{10} - 7c - an\right)$

# Worst case linear time selection

Guess that $T(n) = O(n)$ when $T(n) \leq T\left(\left\lceil \frac{n}{5} \right\rceil\right) + T\left(\frac{7n}{10} + 6\right) + O(n)$

Assume $T(n') \leq cn'$ for some constant $c$ and every $n' < n$, and

show that this implies $T(n) \leq cn$ (for the same constant $c$)

Upon substitution, we get $T(n) \leq c\left\lceil \frac{n}{5} \right\rceil + c\left(\frac{7n}{10} + 6\right) + O(n)$

For $n > n_0$ this is at most $c\left(\frac{n}{5} + 1\right) + c\left(\frac{7n}{10} + \cdots\right)$ for $c \geq 20a$ and $n \geq 140$.

Therefore, $T(n) \leq cn - \left(\frac{cn}{10} - 7c - an\right) \leq cn$

# Simple Selection Algorithm

Select($A, p, r, i$)

1. if $p == r$
2.     return $A[p]$
3. $q$ = Partition($A, p, r$)
4. $k = q - p + 1$
5. if $i == k$
6.     return $A[q]$
7. else if $i \leq k$
8.     Select($A, p, q - 1, i$)
9. else
10.     Select($A, q + 1, r, i - k$)

Partition ($A, p, r$)

1. $x = A[r]$
2. $i = p - 1$
3. for $j = p$ to $r - 1$
4.     if $A[j] \leq x$
5.         $i = i + 1$
6.         exchange $A[i]$ with $A[j]$
7. exchange $A[i + 1]$ with $A[r]$
8. return $i + 1$

# Randomized Selection Algorithm

Randomized-Select$(A, p, r, i)$

1. if $p == r$
2.     return $A[p]$
3.   $q =$ Randomized-Partition$(A, p, r)$
4.   $k = q - p + 1$
5.   if $i == k$
6.     return $A[q]$
7.   else if $i \leq k$
8.     Randomized-Select$(A, \ p, \ q - 1, \ i)$
9.   else
10.     Randomized-Select$(A, \ q + 1, \ r, \ i - k)$

Randomized-Partition $(A, p, r)$

1. $i =$ Random$(p, r)$
2. Exchange $A[r]$ with $A[i]$
3. return Partition$(A, p, r)$

# Randomized Selection Algorithm

Randomized-Select($A, p, r, i$)

1.    if $p == r$
2.       return $A[p]$
3.    $q$ = Randomized-Partition($A, p, r$)
4.    $k = q - p + 1$
5.    if $i == k$
6.       return $A[q]$
7.    else if $i \leq k$
8.       Randomized-Select($A, \ p, \ q - 1, \ i$)
9.    else
10.      Randomized-Select($A, \ q + 1, \ r, \ i - k$)

- Worst-case running time?
  - a) $O(n^2)$ ⬅
  - b) $O(n \log n)$
  - c) $O(n)$

- Expected running time?
  - a) $O(n^2)$
  - b) $O(n \log n)$
  - c) $O(n)$ ⬅

# Running Time

- Indicator variable $\mathbb{I}\{E\}$ is 1 if event $E$ occurs and 0 o/w (see page 118)

- Consider an array $A[p, \dots, r]$ with $n$ elements

- If $X_k = \mathbb{I}\{$the subarray $A[p, \dots, q]$ has exactly $k$ elements$\}$, then

$$T(n) \le \sum_{k=1}^{n} X_k \left( T(\max\{k-1, n-k\}) + O(n) \right)$$

$$= \sum_{k=1}^{n} X_k \left( T(\max\{k-1, n-k\}) \right) + \sum_{k=1}^{n} X_k \, O(n)$$

$$= \sum_{k=1}^{n} X_k \left( T(\max\{k-1, n-k\}) \right) + O(n)$$

# Expected Running Time

- $X_k = \mathbb{I}\{$the subarray $A[p, \dots, q]$ has exactly $k$ elements$\}$. Since $A[p, \dots, r]$ has $n$ elements and $q$ is chosen uniformly at random, we have $\mathbb{E}[X_k] = \frac{1}{n}$, so
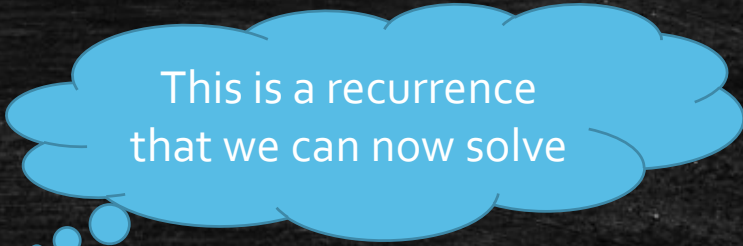
$$\mathbb{E}[T(n)] \leq \mathbb{E}\left[\sum_{k=1}^{n} X_k \left(T(\max\{k-1, n-k\})\right) + O(n)\right]$$

$$= \sum_{k=1}^{n} \mathbb{E}[X_k(T(\max\{k-1, n-k\}))] + O(n)$$

$$= \sum_{k=1}^{n} \mathbb{E}[X_k]\, \mathbb{E}[T(\max\{k-1, n-k\})] + O(n)$$

$$= \sum_{k=1}^{n} \frac{1}{n}\, \mathbb{E}[T(\max\{k-1, n-k\})] + O(n)$$

# Expected Running Time

Thus, $\mathbb{E}[T(n)] \leq \sum_{k=1}^{n} \frac{1}{n} \mathbb{E}[T(\max\{k-1, n-k\})] + O(n),$

and $\max\{k-1, n-k\}) = \begin{cases} k-1 \text{ if } k > \lceil n/2 \rceil \\ n-k \text{ if } k \leq \lceil n/2 \rceil \end{cases}$

This is a recurrence
that we can now solve

So, $\mathbb{E}[T(n)] \leq \frac{2}{n} \sum_{k=\lfloor \frac{n}{2} \rfloor}^{n-1} \mathbb{E}[T(k)] + O(n)$

# Expected Running Time

We use substitution in order to solve: $\mathbb{E}[T(n)] \leq \dfrac{2}{n} \sum_{k=\lceil \frac{n}{2} \rceil}^{n-1} \mathbb{E}[T(k)] + O(n)$

We guess that $\mathbb{E}[T(n)] = O(n)$

Hence, we assume $\mathbb{E}[T(n')] \leq cn'$ for some constant $c$ and every $n' < n$,

and we show that $\mathbb{E}[T(n)] \leq cn$ for that same constant $c$. Upon substitution:

$$\mathbb{E}[T(n)] \leq \frac{2}{n} \sum_{k=\lceil \frac{n}{2} \rceil}^{n-1} ck + an \leq \frac{2c}{n} \left( \frac{n^2 - n}{2} - \frac{n^2/4 - 3n/2 + 2}{2} \right) + an$$

which leads to $\mathbb{E}[T(n)] \leq cn - \left( \dfrac{cn}{4} - \dfrac{c}{2} - an \right)$

# Expected Running Time

We use substitution in order to solve: $\mathbb{E}[T(n)] \leq \dfrac{2}{n} \displaystyle\sum_{k=\lfloor\frac{n}{2}\rfloor}^{n-1} \mathbb{E}[T(k)] + O(n)$

We guess that $\mathbb{E}[T(n)] = O(n)$

Hence, we assume $\mathbb{E}[T(n')] \leq cn'$ for some constant $c$ and every $n' < n$,

and we show that $\mathbb{E}[T(n)] \leq cn$ for that same constant $c$. Upon substitution:

$$\mathbb{E}[T(n)] \leq \frac{2}{n} \sum_{k=\lfloor\frac{n}{2}\rfloor}^{n-1} ck + an \leq \frac{2c}{n}\left(\frac{n^2-n}{2} - \frac{n^2/4}{\phantom{x}} \phantom{xxxxxxx}\right) + an$$

for $c > 4a$,
and $n \geq \dfrac{2c}{c-4a}$

which leads to $\mathbb{E}[T(n)] \leq cn - \left(\dfrac{cn}{4} - \dfrac{c}{2} - an\right) \leq cn$