# Ch 7 [H]

Modification Analysis; "What-If" for Closed Systems

## Preview

- asymptotic bounds
- modification analysis for closed systems
- closed vs. open networks

## Review (7.1 [H])

Little's Law for an Open Systems

$$E[N] = d \cdot E[T]$$
$$(E[N_a] = \lambda \cdot E[T_Q])$$
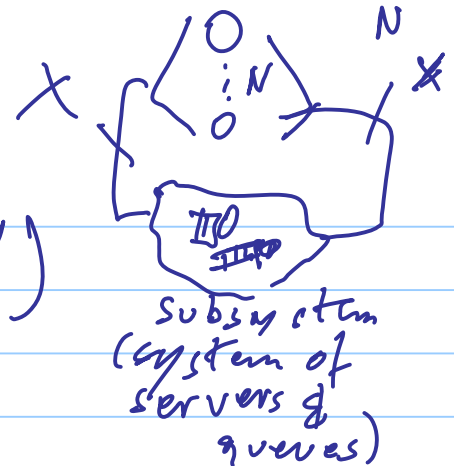$$(E[N_{red}] = d \cdot E[T_{red}])$$

Little's Law for a Closed Batch System

$$N = X \cdot E[T].$$

(This actually holds for _any_ closed system, but for interactive closed system, the following also holds.)

Little's Law for a Closed Interactive System

$$E[R] = \frac{N}{X} - E[Z] \qquad \left(\text{note}: E[T] = \frac{N}{X}\right)$$

Utilization Law

$$\rho_i = \frac{d_i}{\mu_i} = d_i \, E[S_i] = X_i \, E[S_i]$$

Forced Flow Law

$$X_i = E[V_i] \cdot X$$

Bottleneck Law

$$\rho_i = X \cdot E[D_i], \text{ where } D_i \text{ is the total service demand on device } i$$
for all visits of a single job.

## 7.2 [H] Asymptotic Bounds for Closed Systems

Let $m$ be the number of devices in the system. $E[D_i]$ is as defined before, i.e., the expected total service demand on device $i$ by a single job

Let $D = \sum_{i=1}^{m} E[D_i]$

Let $D_{max} = \max_{i \in [1..m]} \{ E[D_i] \}$

**Theorem 7.1** *For any closed interactive system with $N$ terminals,*

$$X \leq \min\left(\frac{N}{D + \mathbf{E}[Z]}, \frac{1}{D_{\max}}\right),$$

$$\mathbf{E}[R] \geq \max\left(D, N \cdot D_{\max} - \mathbf{E}[Z]\right).$$

*Importantly, the first term in each clause ($\frac{N}{D + \mathbf{E}[Z]}$ or $D$) is an asymptote for small $N$, and the second term ($\frac{1}{D_{\max}}$ or $N \cdot D_{\max} - \mathbf{E}[Z]$) is an asymptote for large $N$.*

Proof.

First, consider the large $N$ (multiprogramming level) case

1. $\forall i$  $X\left[\frac{jobs}{sec}\right] \cdot \mathbf{E}[D_i] \left[sec\right] = \rho_i \left[jobs\right]$ (utilization) $\leq 1$

   (Remember that $\rho_i$ is the expected number of jobs in server $i$.)

$$\Rightarrow \forall i \quad X = \frac{\ell_i}{E[D_i]} \leq \frac{1}{E[D_i]} \quad \Rightarrow \quad X \leq \frac{1}{D_{max}}.$$

2. $E[R] = $ (Little's Law for interactive closed systems) $= \frac{N}{X} - E[Z] \geq N \cdot D_{max} - E[Z].$

Note that for large $N$, the $D_{max}$ server is always busy ($\ell_{max} \approx 1$), so $X = \frac{1}{D_{max}}$ is an asymptote for large $N$.

Second, consider the small $N$ asymptote.

1. Let $E[R(N)]$ be the mean response time when the multiprogramming level is $N$. Then,
$$E[R(N)] \geq E[R(1)] = D \left( = \sum_{i=1}^{m} E[D_i] \right)$$

For low $N$, there is no "congestion," so the bound is tight.

2. $X = \dfrac{N}{E[R] + E[Z]} \leq \dfrac{N}{D + E[Z]}$ . $\square$

# A Simple Example of Bounds

$E[Z] = 18 \text{ sec}$

$E[D_{cpu}] = 5 \text{ sec}$

$E[D_{disk\,a}] = 4 \text{ sec}$

$E[D_{disk\,b}] = 3 \text{ sec}$
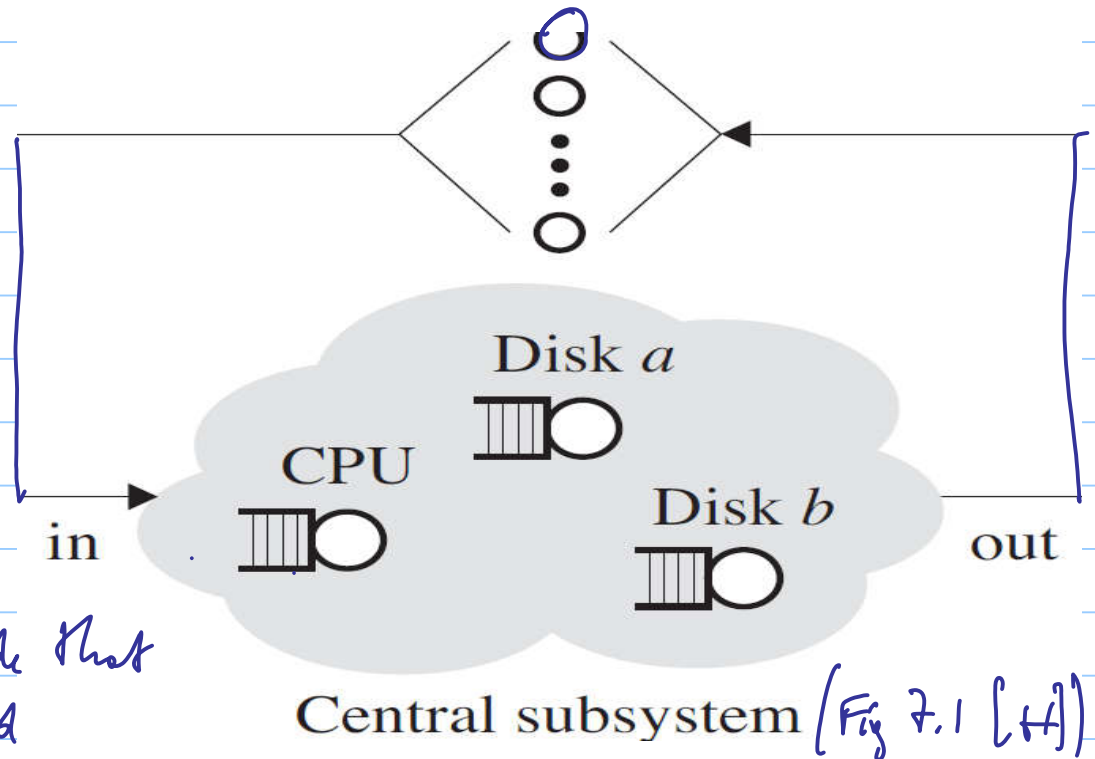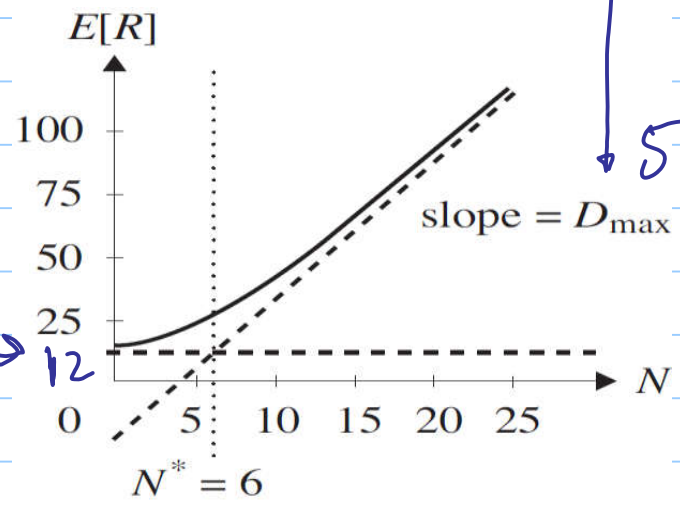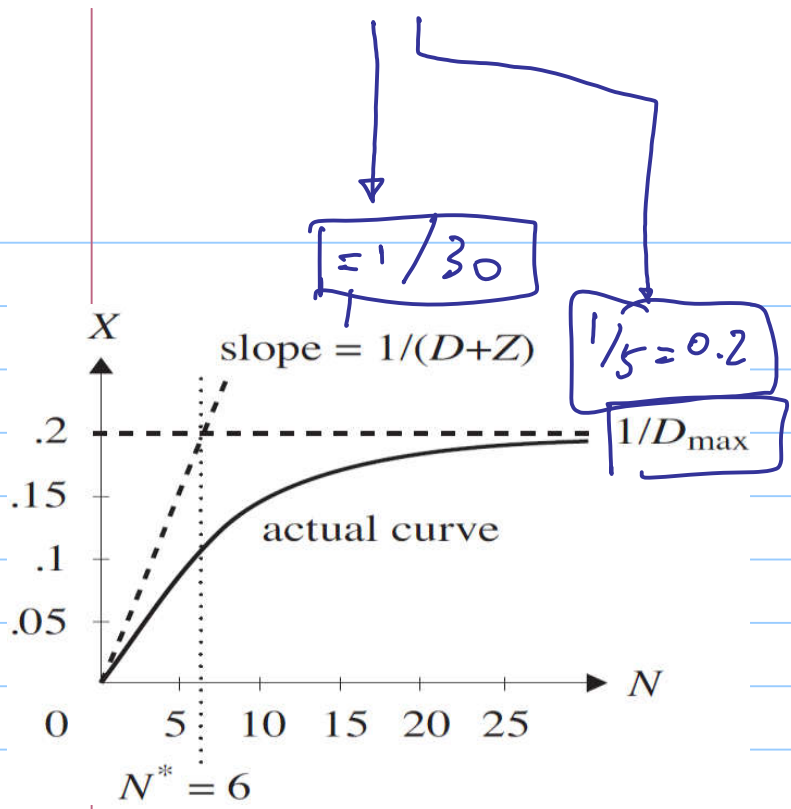
So, $D = 5 + 4 + 3 = 12 \text{ sec}$

$D_{max} = 5$ (the CPU is the bottleneck device)

Thm. 7.1 allows us to conclude that

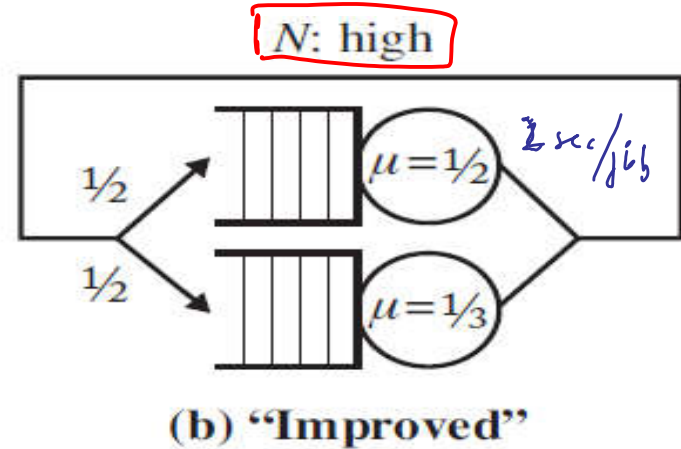$$X \leq \min\left\{ \frac{N}{12+18}, \frac{1}{5} \right\}, \text{ and }$$

$$E[R] \geq \max\{12, 5N - 18\}$$

Disk $a$

CPU

Disk $b$

in

out

Central subsystem (Fig 7.1 [H])

Left graph: $X$ axis (vertical) with values .2, .15, .1, .05 and $N$ axis (horizontal) with values 0, 5, 10, 15, 20, 25.

slope = $1/(D+Z)$

$= 1/30$

$1/5 = 0.2$

$1/D_{max}$

actual curve

$N^* = 6$

Right graph: $E[R]$ axis (vertical) with values 100, 75, 50, 25 and $N$ axis (horizontal) with values 0, 5, 10, 15, 20, 25.

slope = $D_{max}$

12

5

$N^* = 6$

# 7.3 [H] Modification Analysis for Closed Systems



(a) Original

(b) "Improved"

The throughput is bounded above, for high $N$, by $1/D_{max}$, and $D_{max}$, the service demand on the slower disk ($\mu = \frac{1}{3}$) is unchanged!

# Important Observations

The bounds for $X$ and $E[R]$ meet at $N^{\#} = \dfrac{D + E[Z]}{D_{max}}$.

- $N^{\#}$ represents the point beyond which there must be some queueing in the system. ($E[R] > D$).

- For fixed $N > N^{\#}$, to increase throughput or lower response time, one must reduce $D_{max}$. Other changes will be largely ineffective.

- If $E[Z] = 0$ (batch case), $N^{*}$ decreases, i.e., the domination of $D_{max}$ occurs with fewer jobs in the system

# 7.4 (H) More Modification Analysis Examples

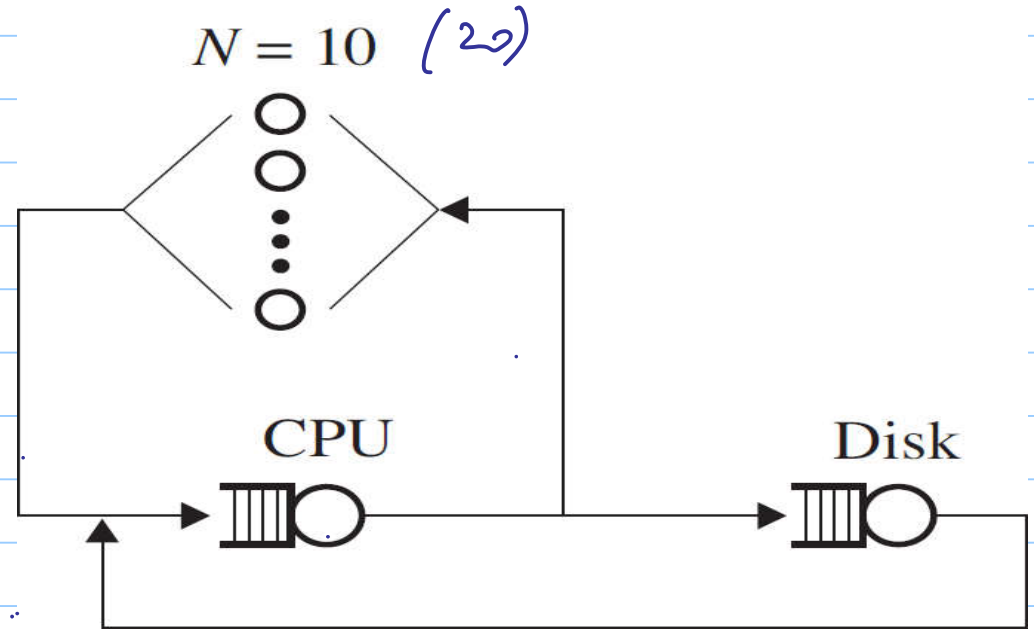Simple example

1. System A: $D_{CPU} = 4.6$, $D_{disk} = 4.0$

2. System B: $D_{CPU} = 4.9$, $D_{disk} = 1.9$
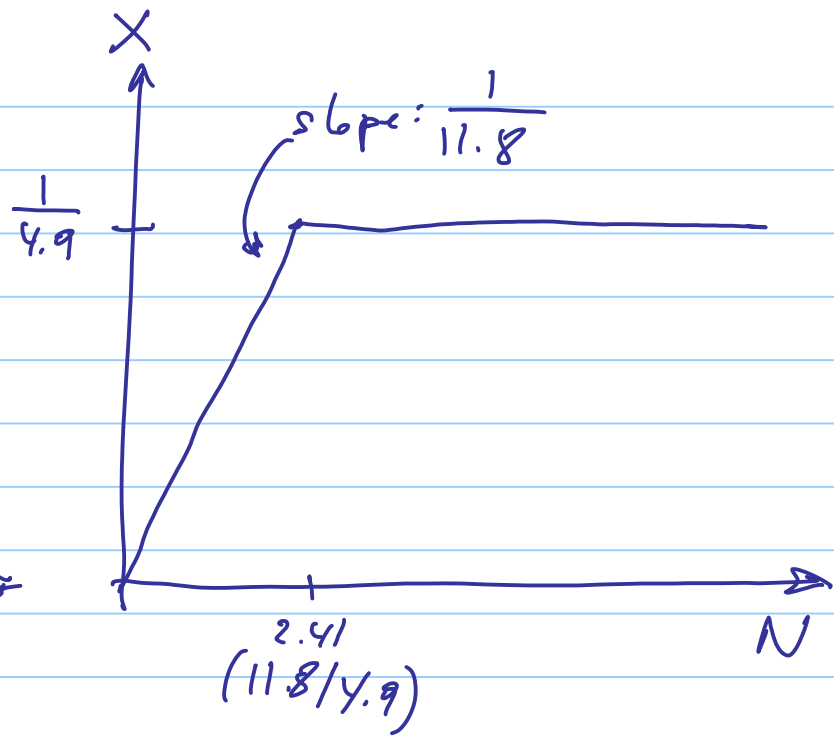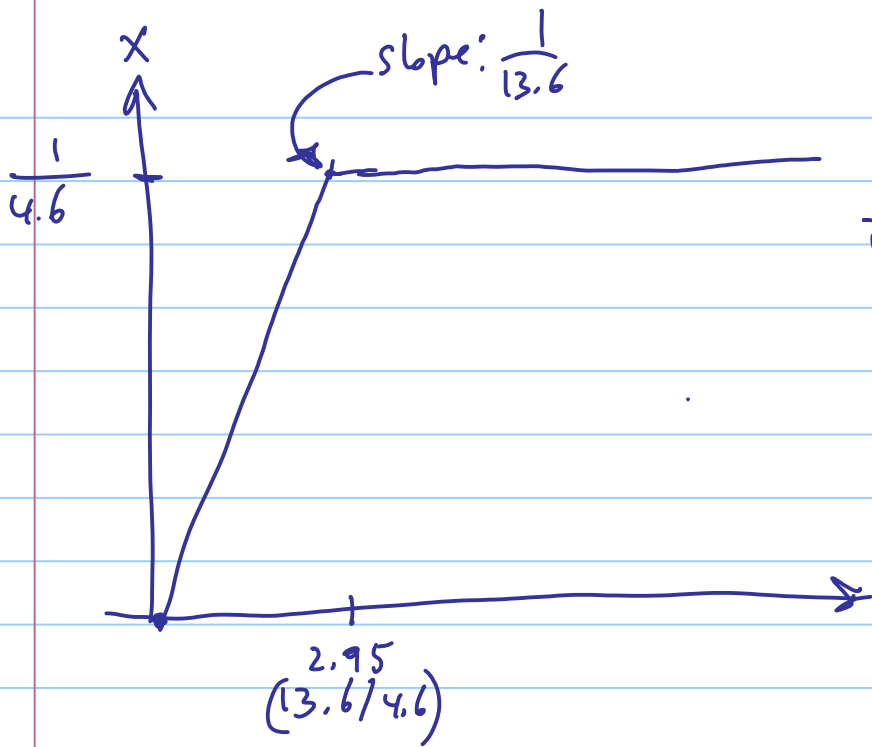
$E[Z] = 5$

$N = 10$ $(20)$

Which system has higher throughput?

$$N_A^* = \frac{D + E[Z]}{D_{max}} = \frac{8.6 + 5}{4.6} < 3$$

$$N_B^* = \frac{D + E[Z]}{D_{max}} = \frac{6.8 + 5}{4.9} < 3$$

Since $N \gg N^*$ for both systems, $X \approx \frac{1}{D_{max}}$, and therefore System A wins.

Left graph:
X (vertical axis)
$\frac{1}{4.6}$
slope: $\frac{1}{13.6}$
2.95
$(13.6/4.6)$

Right graph:
X (vertical axis)
$\frac{1}{4.9}$
slope: $\frac{1}{11.8}$
2.41
$(11.8/4.9)$
N (horizontal axis)

The following measurements were obtained for an interactive system[2]:

- $T = 650$ seconds (the length of the observation interval)
- $B_{cpu} = 400$ seconds
- $B_{slowdisk} = 100$ seconds
- $B_{fastdisk} = 600$ seconds
- $C = C_{cpu} = 200$ jobs
- $C_{slowdisk} = 2,000$ jobs
- $C_{fastdisk} = 20,000$ jobs
- $E[Z] = 15$ seconds
- $N = 20$ users

in 650 sec,
I see
200 jobs
through
here

In this example, we examine four possible improvements (modifications) – hence the name "modification analysis."

1. **Faster CPU:** Replace the CPU with one that is twice as fast.
2. **Balancing slow and fast disks:** Shift some files from the fast disk to the slow disk, balancing their demand.
3. **Second fast disk:** Buy a second fast disk to handle half the load of the busier existing fast disk.
4. **Balancing among three disks plus faster CPU:** Make all three improvements together: Buy a second fast disk, balance the load across all three disks, and also replace the CPU with a faster one.

## Derived quantities – some expectation signs are dropped.

- $D_{cpu} = B_{cpu}/C = 400 \text{ sec}/200 \text{ jobs} = 2.0 \text{ sec/job}$
- $D_{slowdisk} = B_{slowdisk}/C = 100 \text{ sec}/200 \text{ jobs} = 0.5 \text{ sec/job}$
- $D_{fastdisk} = B_{fastdisk}/C = 600 \text{ sec}/200 \text{ jobs} = 3.0 \text{ sec/job}$
- $E[V_{cpu}] = C_{cpu}/C = 200 \text{ visits}/200 \text{ jobs} = 1 \text{ visit/job}$
- $E[V_{slowdisk}] = C_{slowdisk}/C = 2{,}000 \text{ visits}/200 \text{ job} = 10 \text{ visits/job}$
- $E[V_{fastdisk}] = C_{fastdisk}/C = 20{,}000 \text{ visits}/200 \text{ job} = 100 \text{ visits/job}$
- $E[S_{cpu}] = B_{cpu}/C_{cpu} = 400 \text{ sec}/200 \text{ visits} = 2.0 \text{ sec/visit}$
- $E[S_{slowdisk}] = B_{slowdisk}/C_{slowdisk} = 100 \text{ sec}/2{,}000 \text{ visits} = .05 \text{ sec/visit}$
- $E[S_{fastdisk}] = B_{fastdisk}/C_{fastdisk} = 600 \text{ sec}/20{,}000 \text{ visits} = .03 \text{ sec/visit}$

$$D_{max} = \max\{?, 0.5, 3\} = 3$$

$$D = D_{cpu} + D_{slowdisk} + D_{fastdisk} = 2 + 0.5 + 3 = 5.5 \qquad D + k[z] = 5.5 + 15$$

**1. Faster CPU:** Originally, $D_{max} = 3$ sec/job, $D = 5.5$, $N^* = \frac{20.5}{3} \approx 7 \ll N$. $D_{cpu} \rightarrow 1$ sec/job does not change $D_{max} = 3$ sec/job. Notice that $N^*$ hardly changes at all. **The fast disk is the bottleneck.** We can never get more than 1 job done every 3 seconds on average.

**2. Balancing slow and fast disks:** Shift some files from the fast disk to the slow disk, balancing their demand. To do this we need that

$$V_{slow} + V_{fast} = 110 \text{ as originally}$$

but $S_{slow} \cdot V_{slow} = S_{fast} \cdot V_{fast}$ because we are balancing the demand.

Solving this system of linear equations yields the new demands $D_{slow} = D_{fast} = 2.06$. Now, $D_{max} = 2.06$ sec/job, although $D$ increases slightly because some files have been moved from the fast disk to the slow disk.

$$N^* = \frac{2 + 2.06 + 2.06 + 15}{2.06} = \frac{21.12}{2.06} \approx 10.25$$

**3. Second fast disk:** We keep $D_{slow} = 0.5$, the same as before. However, we buy a second fast disk to handle half the load of the original fast disk. So now

$$D_{fast1} = D_{fast2} = 1.5 \text{ sec/job.}$$

Thus our new $D_{max}$ is 2.0 sec/job (the CPU becomes the bottleneck).

$D_{CPU} = 2.0$

$D_{slowdisk} = 0.5$

$D_{fast_1} = D_{fast_2} = 1.5$

$D_{max} = 1.5 \text{ sec/job}$

$D = 5.5 \text{ sec/job}$

$$N^* = \frac{5.5 + 15}{1.5} = \frac{20.5}{1.5} = 13.6$$

4. **Balancing among three disks plus faster CPU:** We now make the CPU faster *and* balance load across all three disks, so

$$V_{slow} + V_{fast1} + V_{fast2} = 110.$$

$$S_{slow} \cdot V_{slow} = S_{fast1} \cdot V_{fast1} = S_{fast2} \cdot V_{fast2}.$$

Solving these simultaneous equations yields: $D_{disk1} = D_{disk2} = D_{disk3} = 1.27$. So $D_{max} = 1.27$, since we cut $D_{cpu}$ to 1 already.

A graph of the results is shown in Figure 7.5. Assuming $N$ is not too small, we conclude the following:

- Change 1 is insignificant.
- Changes 2 and 3 are about the same, which is interesting because change 2 was achieved without any hardware expense.
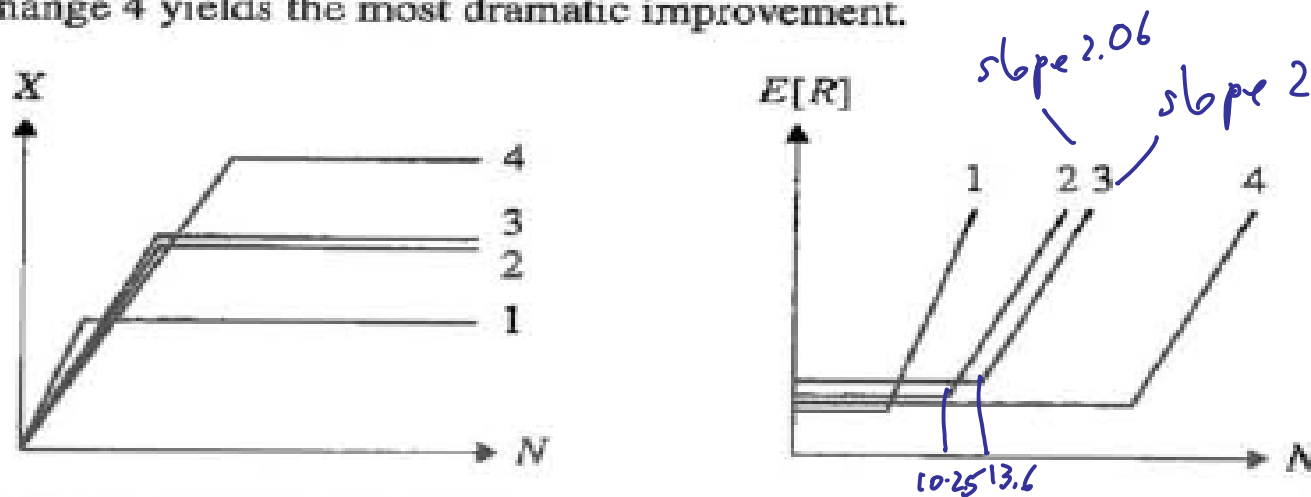- Change 4 yields the most dramatic improvement.



*slope 2.06*
*slope 2*
*10.25 13.6*

*The 'knee' is $N^*$.*

**Figure 7.5.** Throughput and response time versus $N$, showing the effects of four possible improvements from the harder example, where the improvements are labeled 1, 2, 3, and 4.

Why does modification analysis not
apply to open systems?

See 7.5 [H].