# Understanding Performance Aspects of Layered Software with Layered Resources

Murray Woodside

Department of Systems and Computer Engineering

Carleton University, Ottawa, Canada

cmw@sce.carleton.ca,  www.layered queues.org

Jan.27, 2003

**Carleton** UNIVERSITY

1

# The Challenge of Performance in Distributed and Parallel Software

- **several programs interact to complete one response**
  - clients and servers..... peers.... pipelines
- **they execute partly in sequence and partly in parallel**
- ***models are required***
  - systems are difficult to measure in the lab (too big)
- **performance is governed by many kinds of factors:**
  - congestion at different kinds of resources
  - layering of resources
  - layered system overheads
  - unbalanced parallel paths

# Advantages of the layered queueing approach

- an elegant formulation of extended queueing for layered resources
  - scales up to many resources and complex holding patterns
- layered resources have understandable patterns
  - layered resources are common...
  - client-server, parallel service, pipeline, and others
  - bottleneck patterns
- the model notation resembles software design notations such as UML
  - UML performance profile can provide parameter annotations for scenarios in UML.

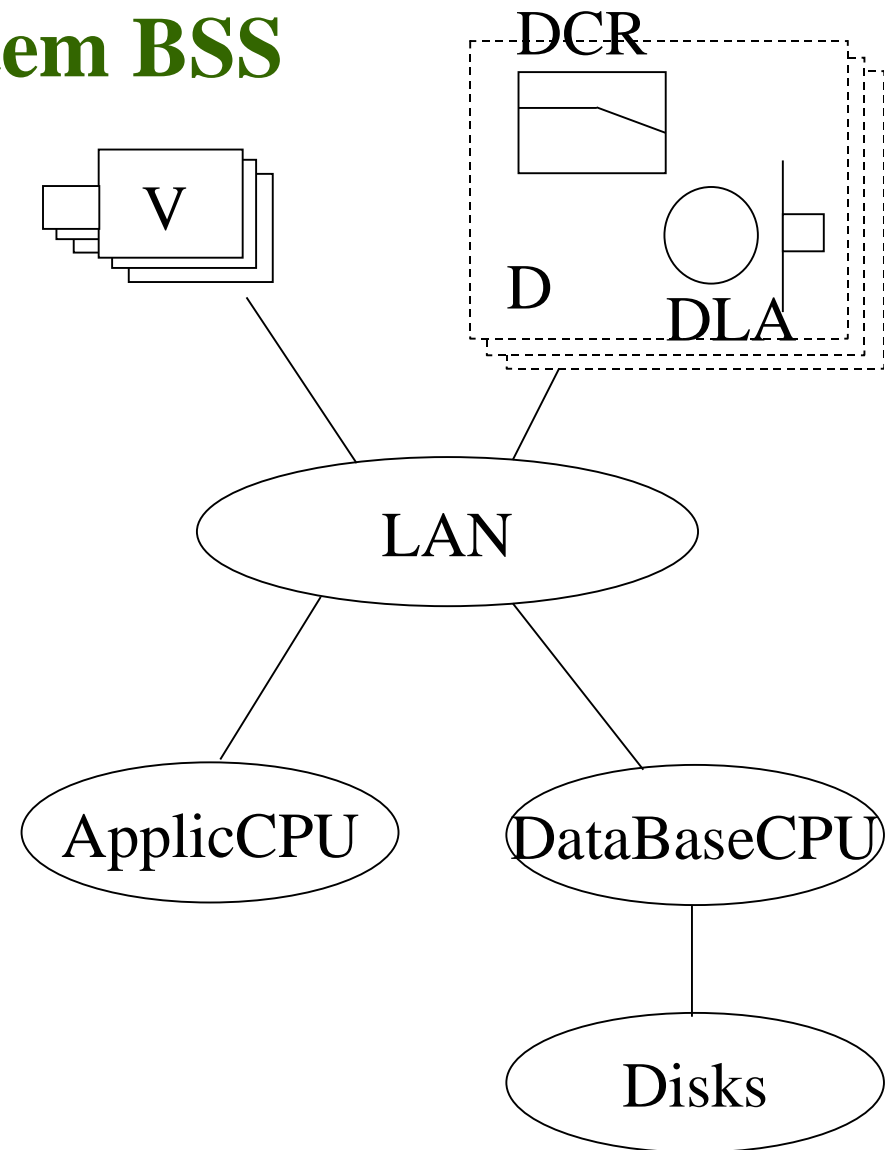# Example: layered modeling of a Building Security System BSS

.....e.g., for a hotel or a university building

- *video surveillance*:
  - poll N web cameras
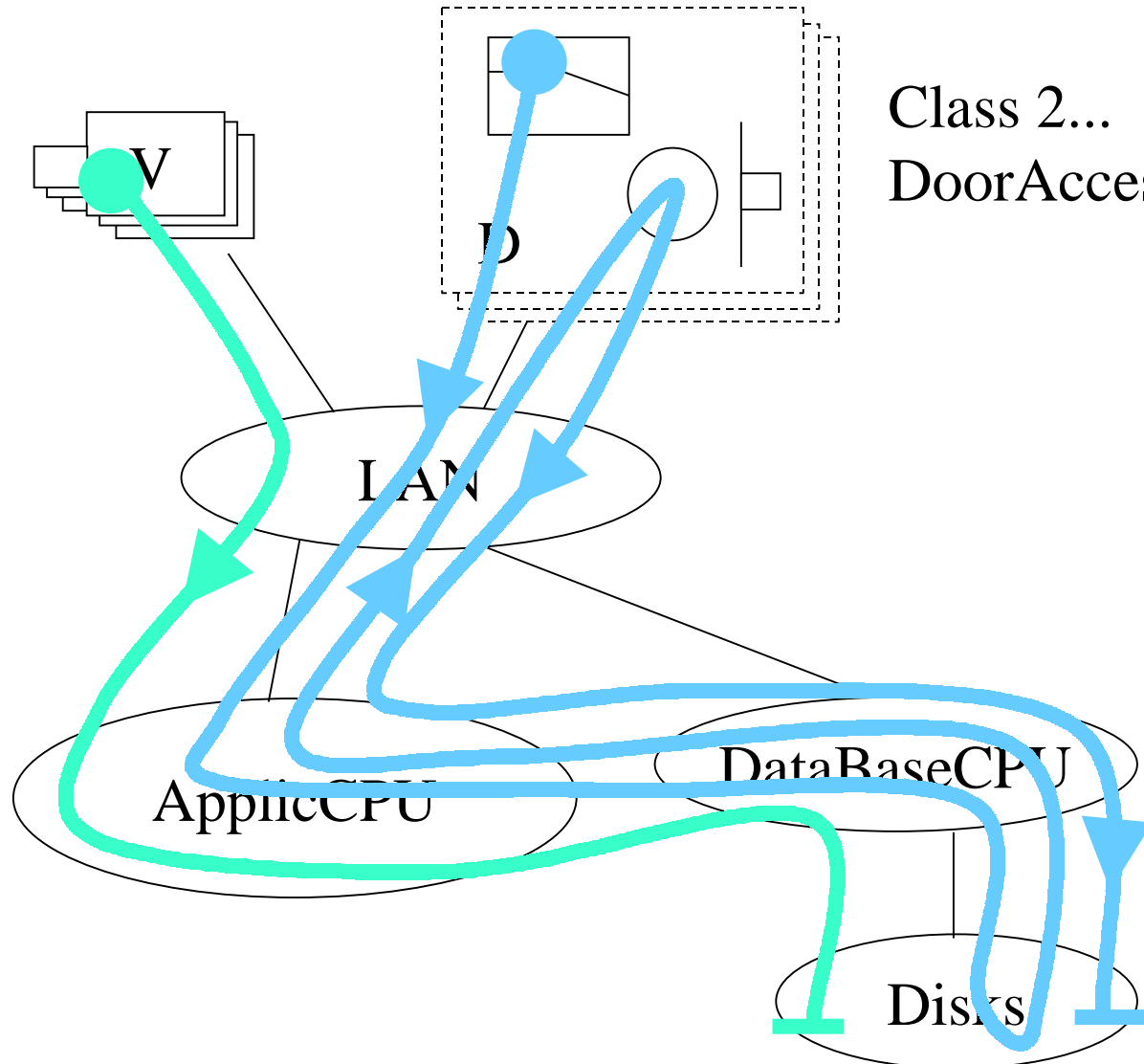  - 1 second cycle (on 95% of polling cycles)

- *door access*:
  - respond to an access card within 1 second, 95% of the time
  - card reader DCR, lock actuator DLA



4

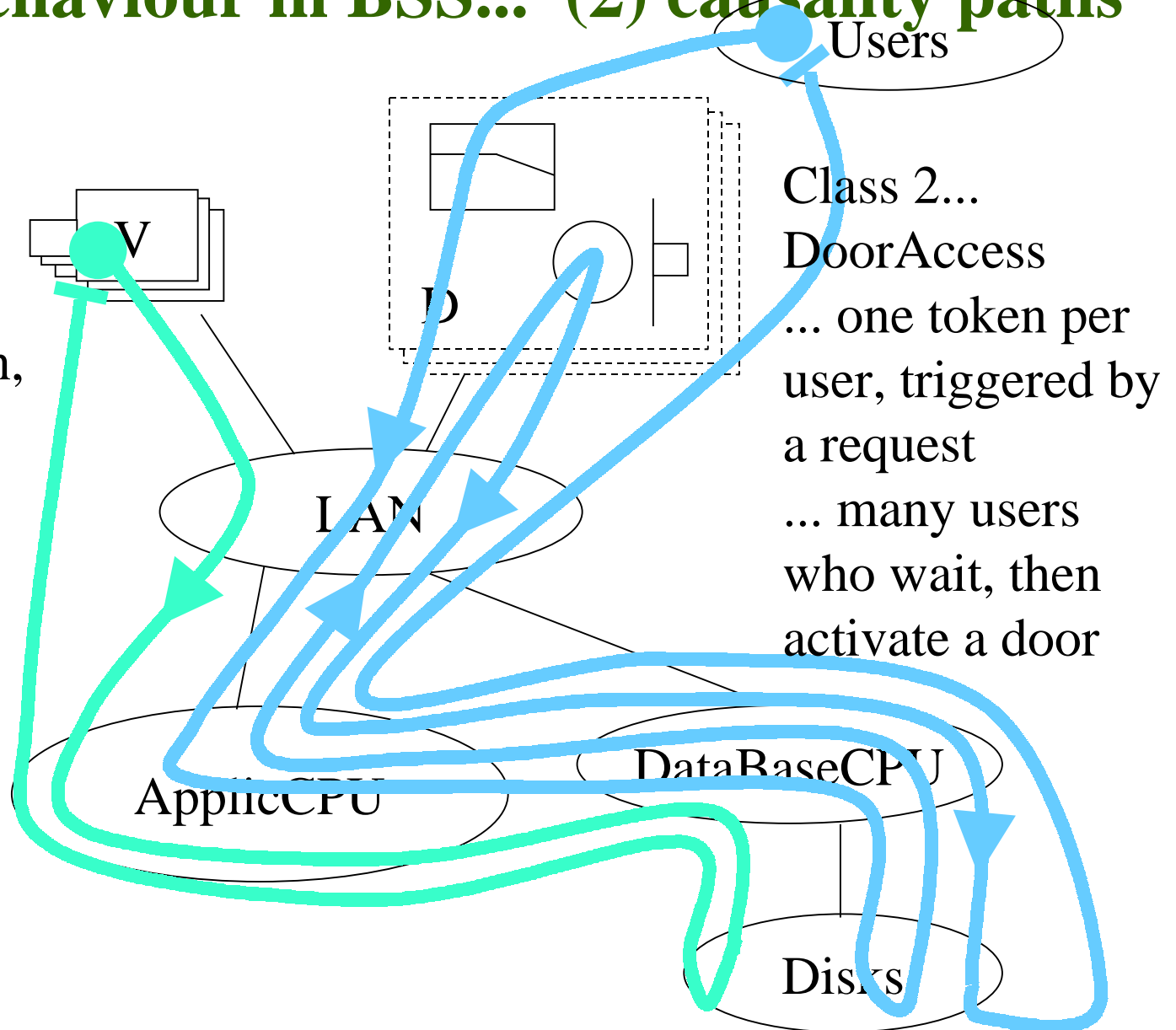# System behaviour in BSS....  (1) data paths



Class 1...
VideoScan

Class 2...
DoorAccess

LAN

ApplicCPU

DataBaseCPU

Disks

# System behaviour in BSS... (2) causality paths



Users

Class 1...
VideoScan
... one token,
cycling
through the
cameras

V

D

Class 2...
DoorAccess
... one token per
user, triggered by
a request
... many users
who wait, then
activate a door

LAN
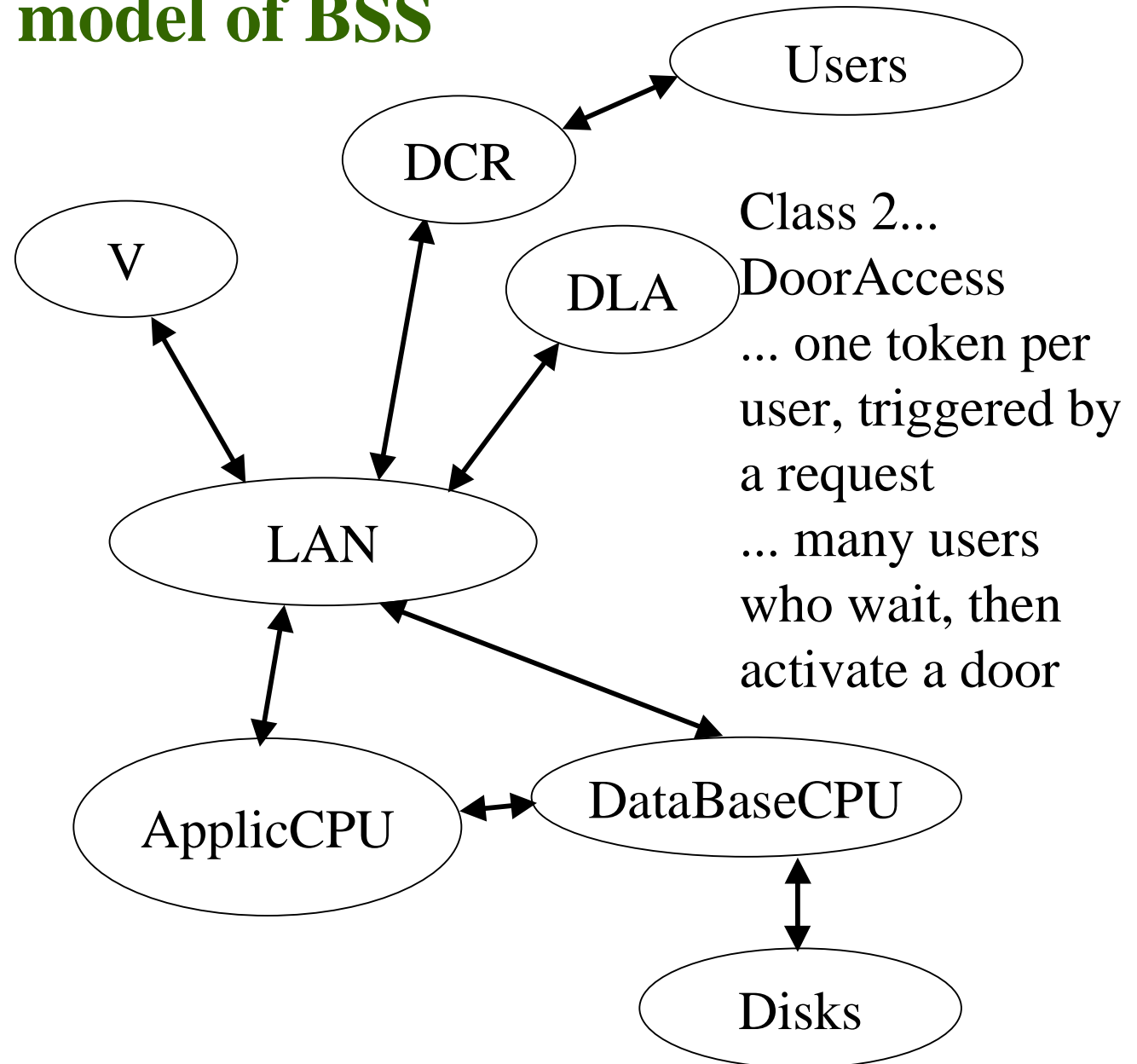
ApplicCPU

DataBaseCPU

Disks

6

# A queueing model of BSS

Class 1...
VideoScan
... one token,
cycling
through the
cameras
... or more
than one, for
double
buffering

■ *focuses on
hardware
servers only*

Users

DCR

V

DLA

Class 2...
DoorAccess
... one token per
user, triggered by
a request
... many users
who wait, then
activate a door

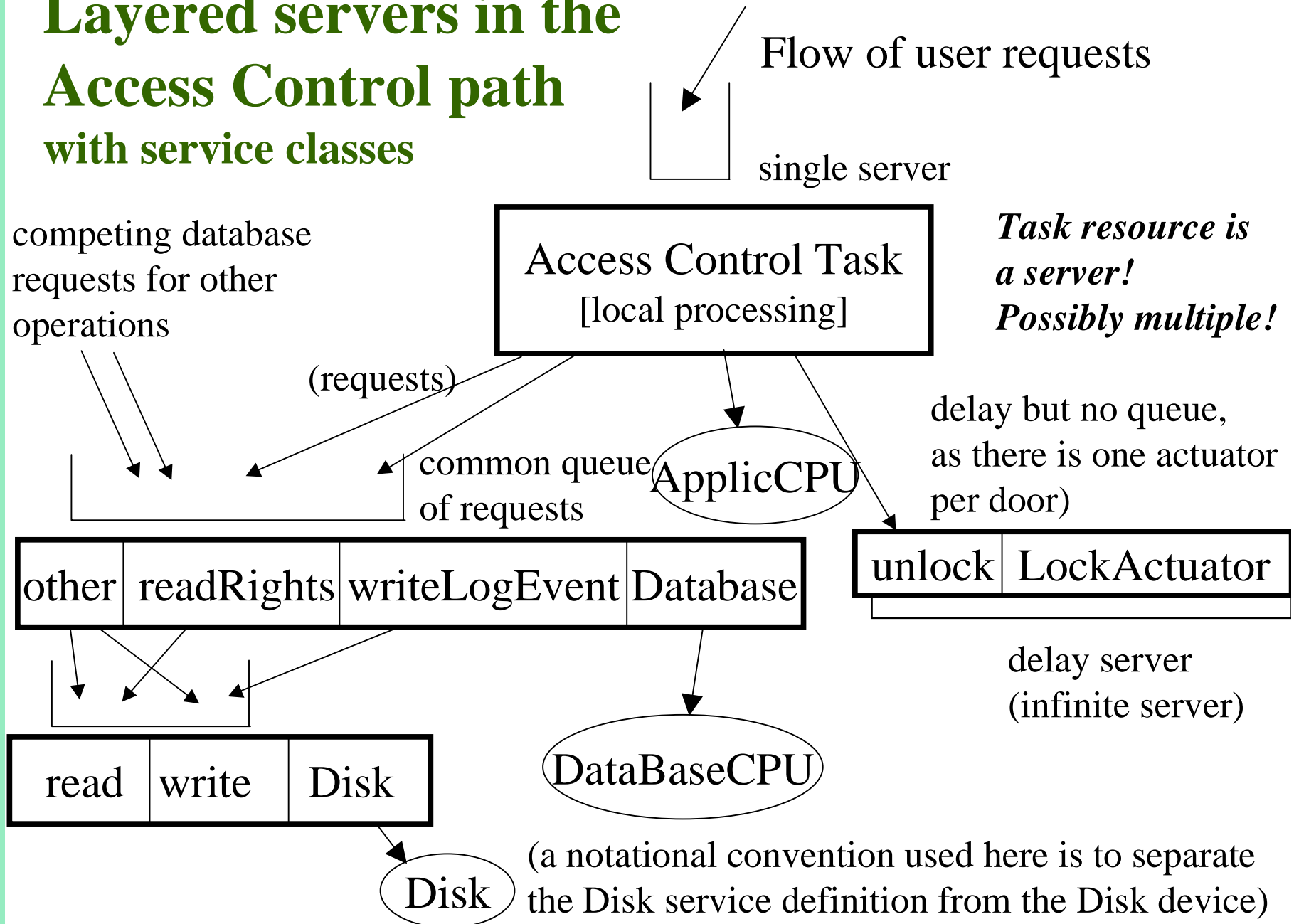LAN

ApplicCPU

DataBaseCPU

Disks

7

# Additional resources that impact performance

- multiple *buffers*, size of buffer pool
- *access control task* can have a message queue
- *database task* may have a queue of requests

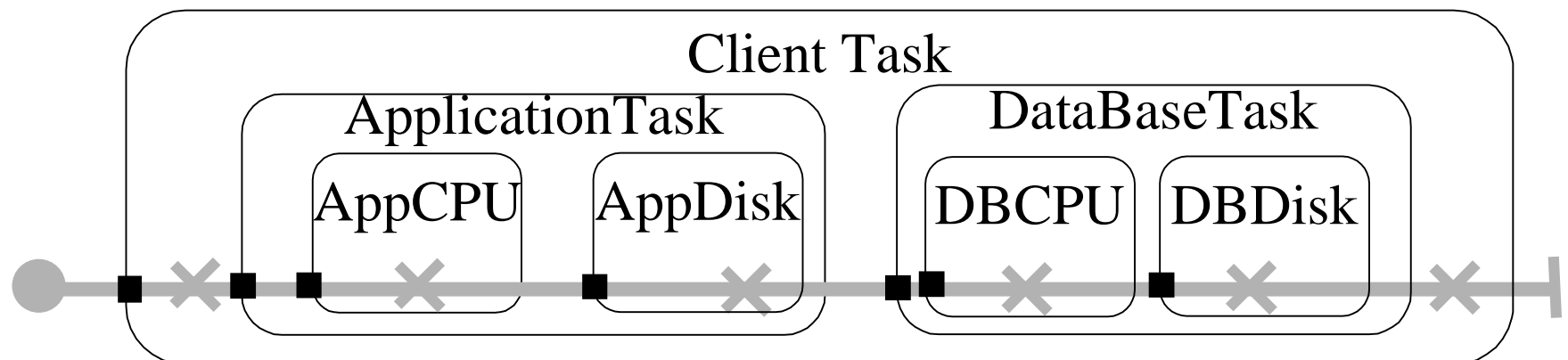- *software resources*
- *layered service*

# Layered servers in the Access Control path

## with service classes

Flow of user requests

single server

*Task resource is a server! Possibly multiple!*

Access Control Task

[local processing]

competing database requests for other operations

(requests)

common queue of requests

ApplicCPU

delay but no queue, as there is one actuator per door)

| other | readRights | writeLogEvent | Database |
|---|---|---|---|

| unlock | LockActuator |
|---|---|

delay server (infinite server)

| read | write | Disk |
|---|---|---|

DataBaseCPU

Disk

(a notational convention used here is to separate the Disk service definition from the Disk device)
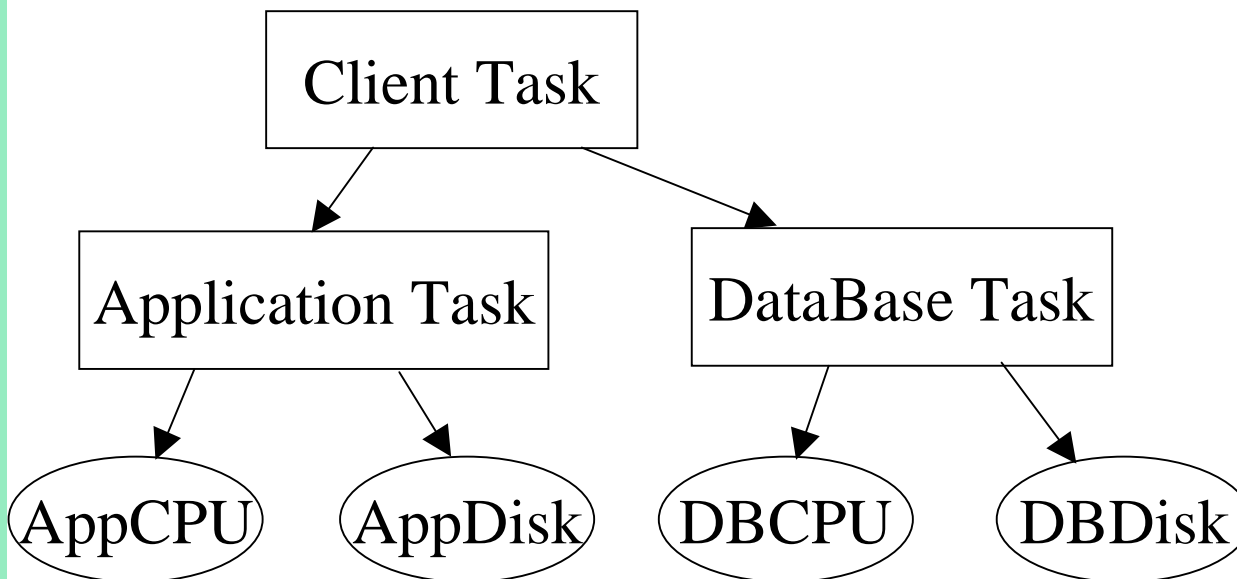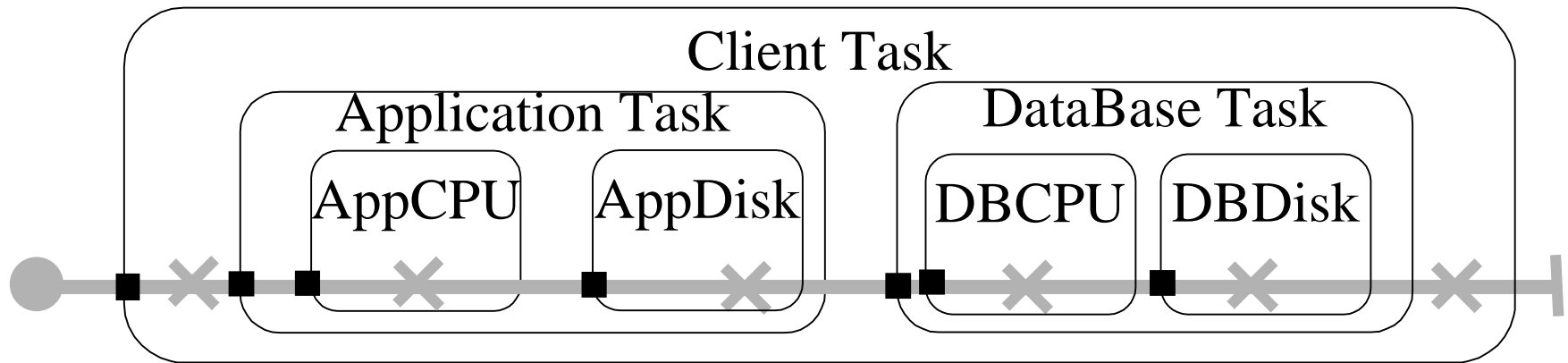
9

# Analysis must include the *software description*

- software operations

- resources that they require
  - task thread resources, buffers, devices

- a useful concept is the *resource context* of an operation:
  - its context is the set of operations during which the resource must be held
  - e.g., outside rectangle represents a client task resource
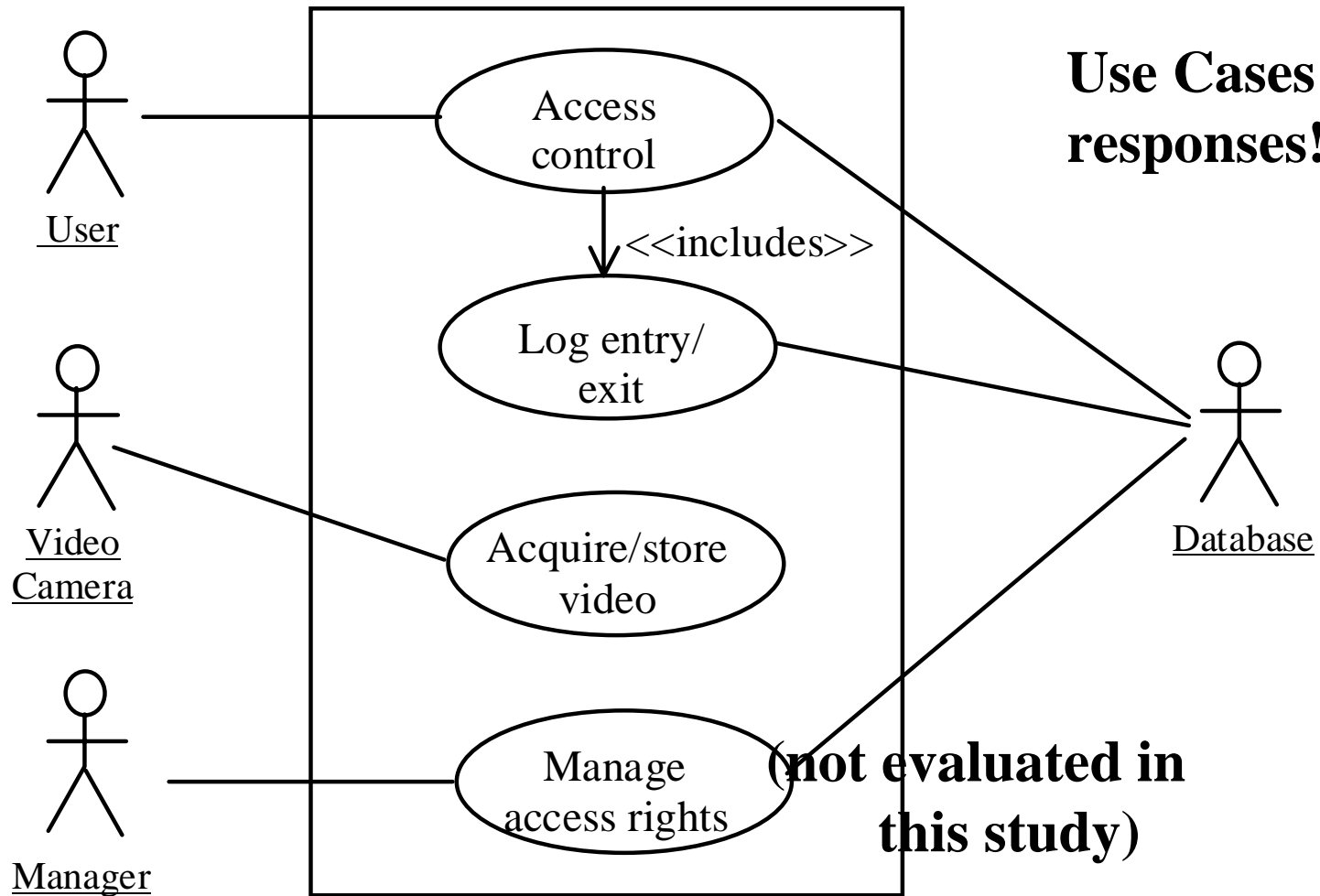    - within this context it acquires other resources

Client Task

ApplicationTask

DataBaseTask

AppCPU    AppDisk

DBCPU    DBDisk
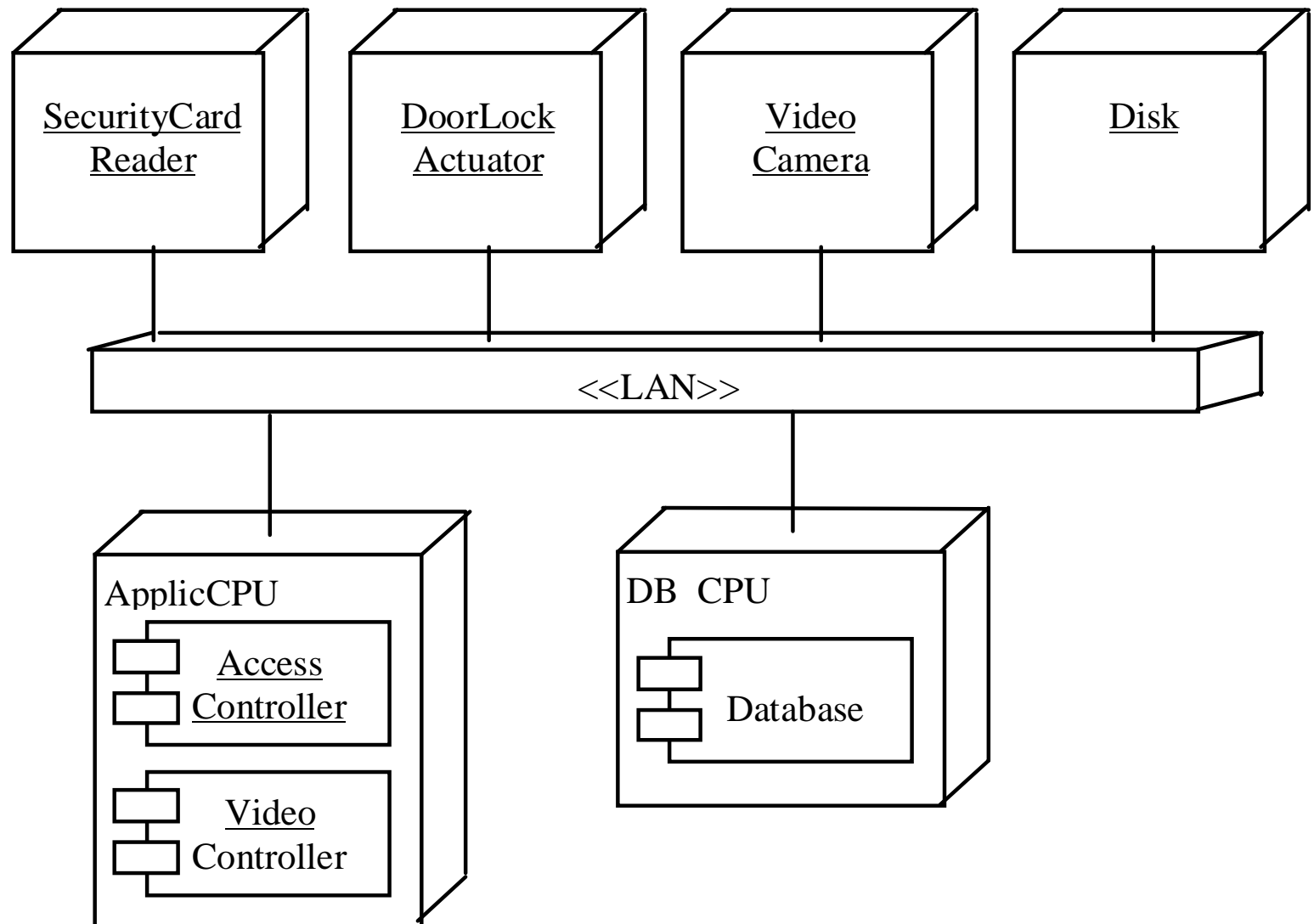
10

# Nested resource contexts gives layered resources



- like a procedure call tree

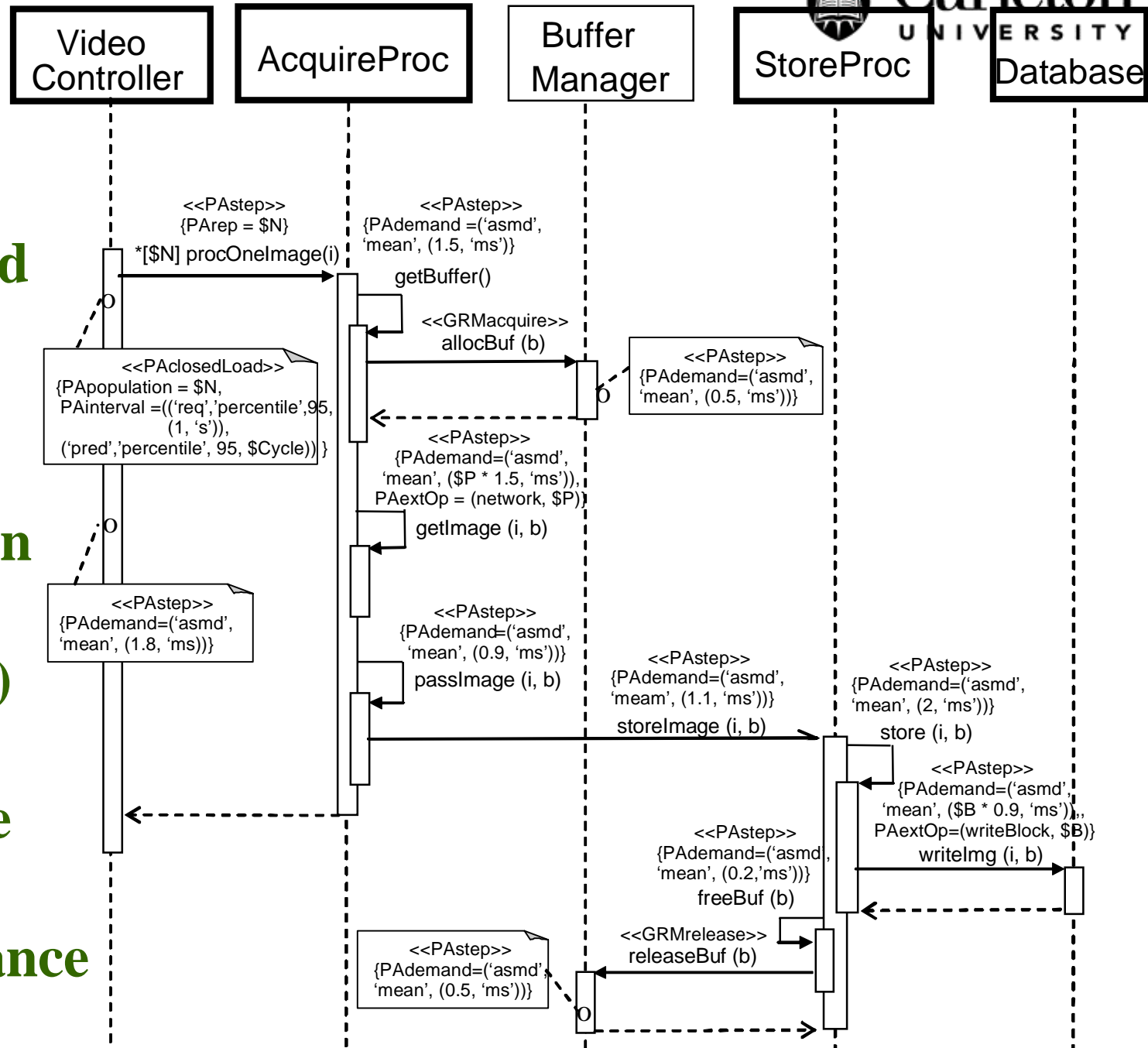# Software Specification of the Building Security System (BSS): (1) Use Cases

**Use Cases define responses!**



- User
- Video Camera
- Manager

Access control

<<includes>>

Log entry/ exit

Acquire/store video

Manage access rights

Database

**(not evaluated in this study)**

# Software Specification (2) Deployment

SecurityCard Reader

DoorLock Actuator

Video Camera

Disk

<<LAN>>

ApplicCPU

Access Controller

Video Controller

DB CPU

Database
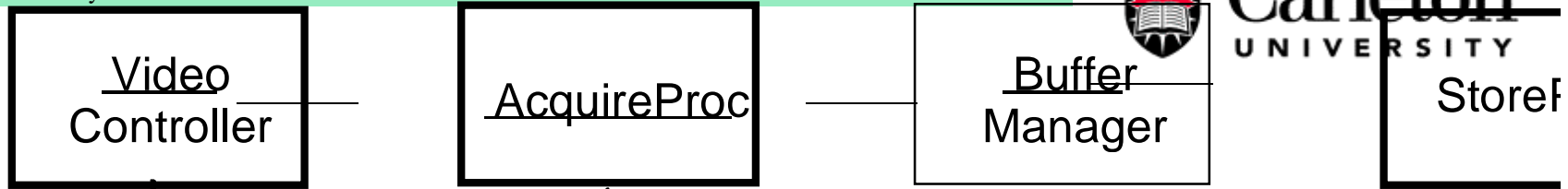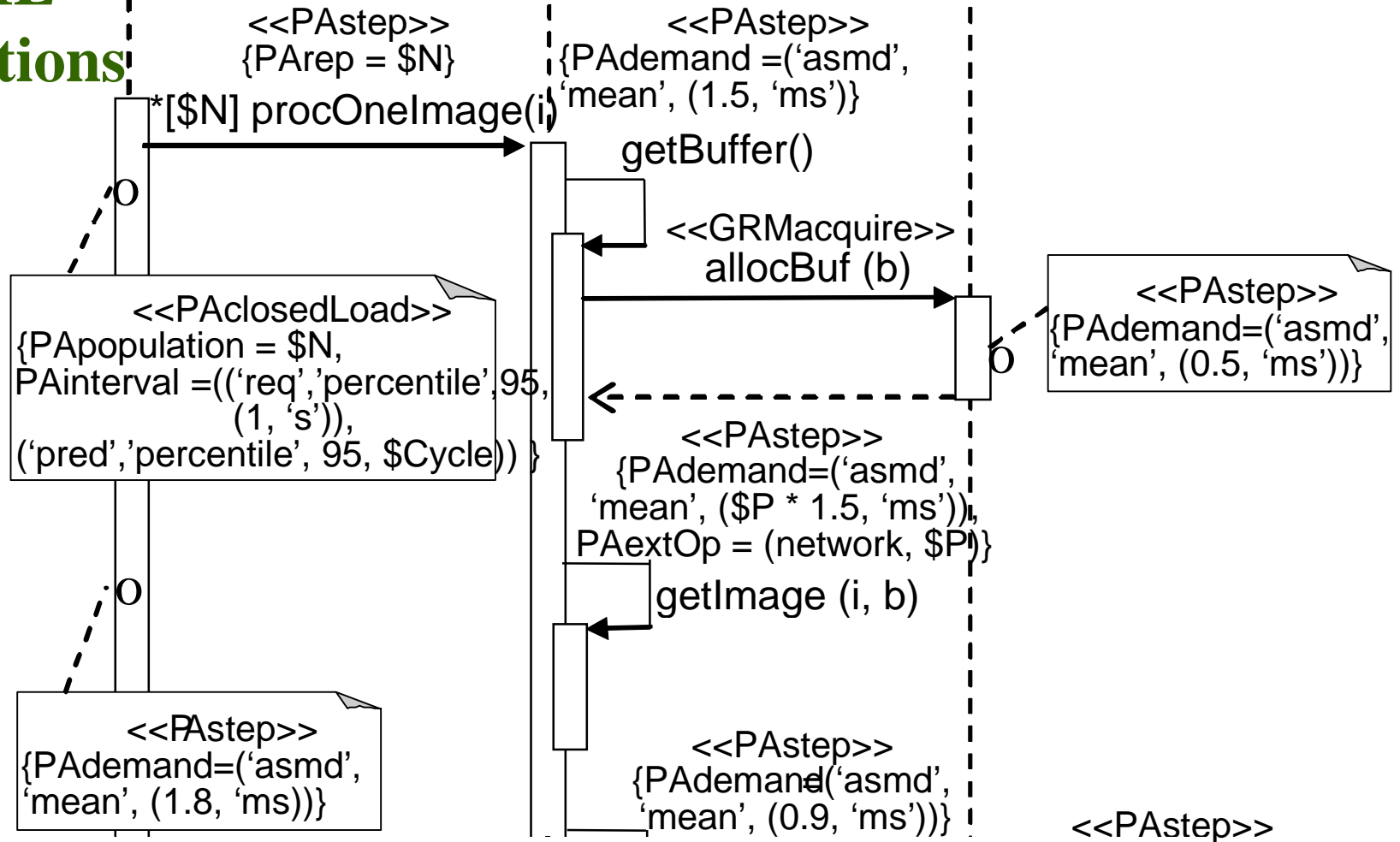
13

<<PAcontext>>

Carleton
U N I V E R S I T Y

**Spec. of BSS (3):**

**Annotated Sequence Diagram for the VideoScan Use Case (scenario)**
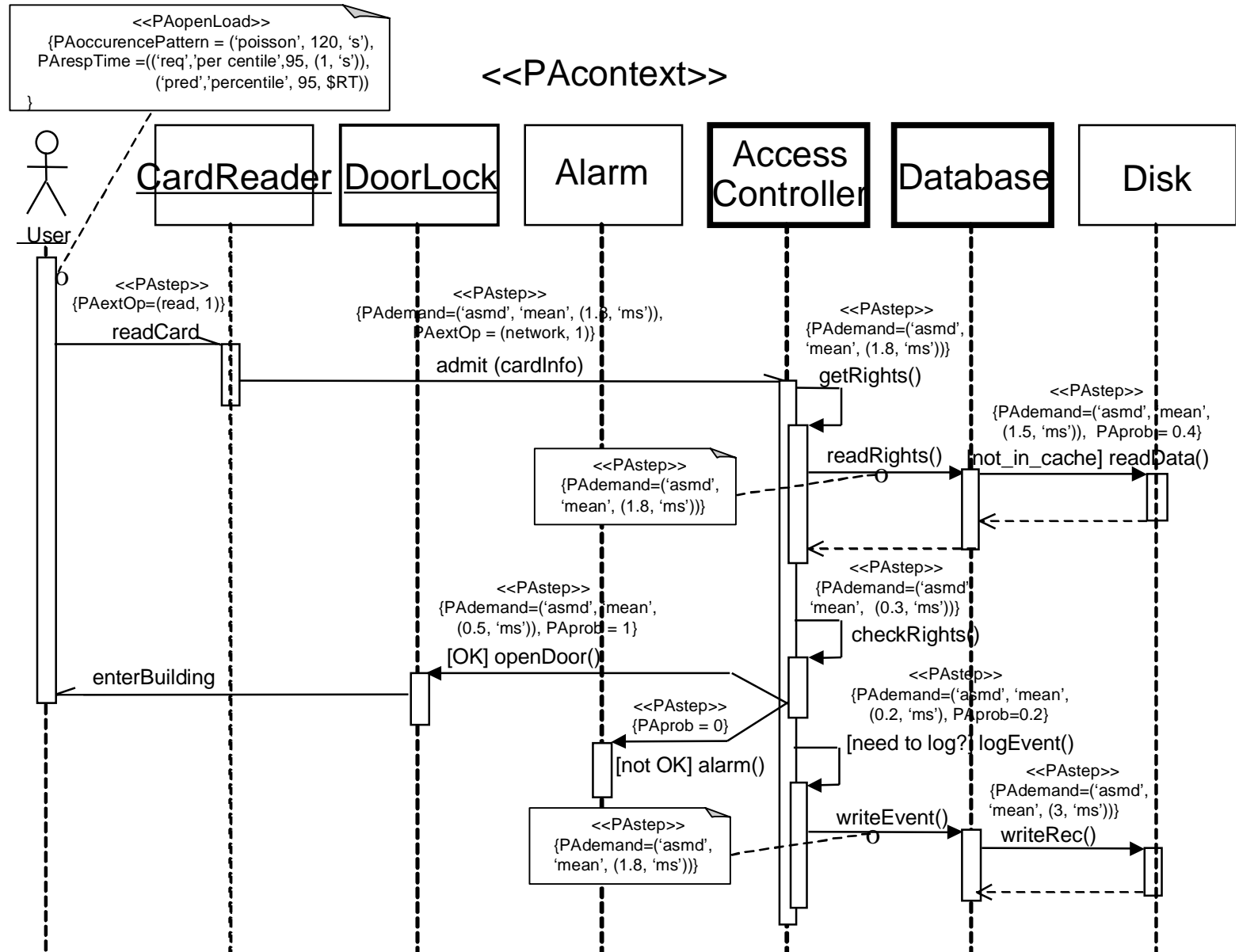
**(using the UML Performance Profile)**

14

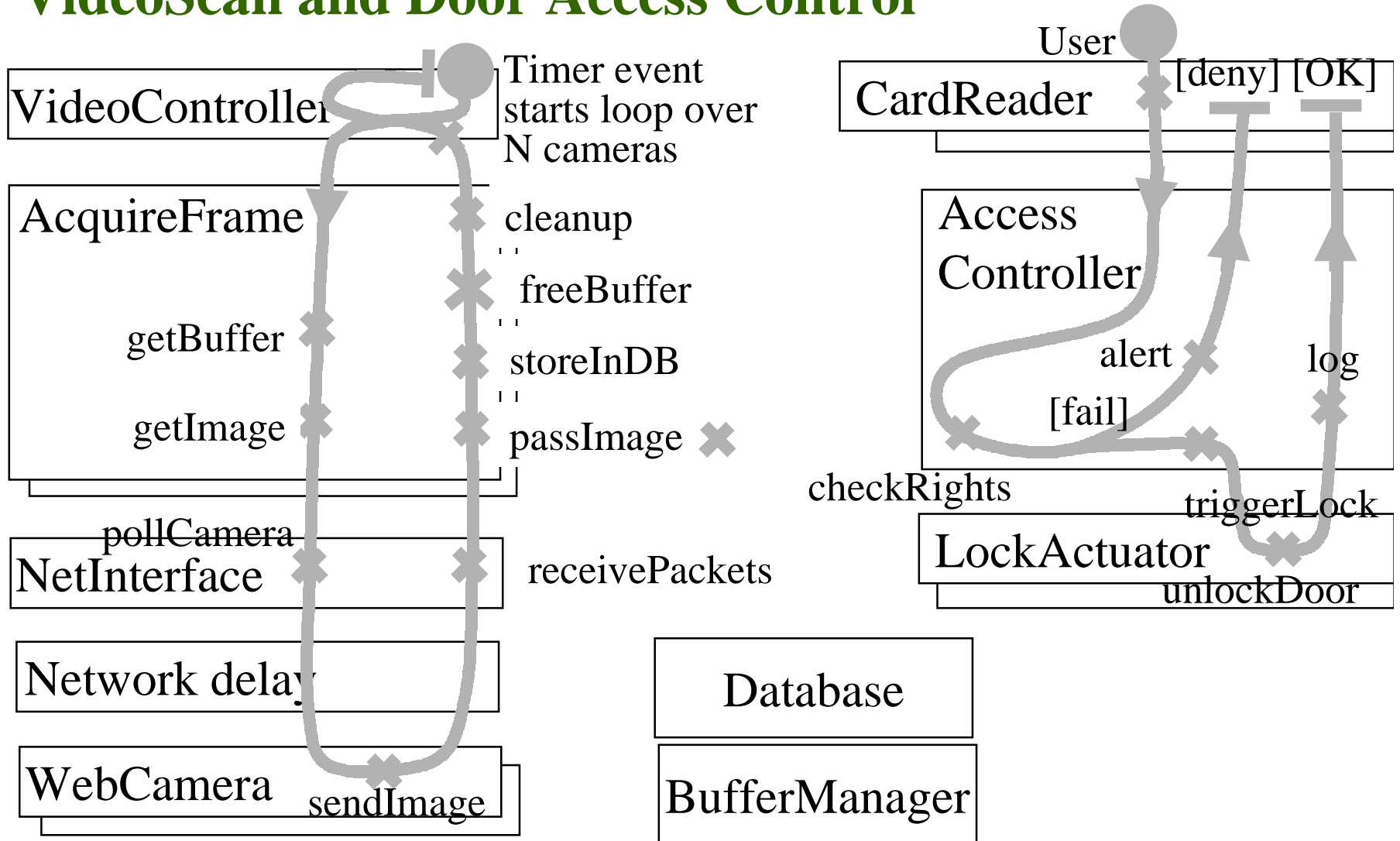| Video Controller | AcquireProc | Buffer Manager | StoreProc | Database |

<<PAstep>>
{PArep = $N}
*[$N] procOneImage(i)

<<PAstep>>
{PAdemand =('asmd',
'mean', (1.5, 'ms')}
getBuffer()

<<GRMacquire>>
allocBuf (b)

<<PAclosedLoad>>
{PApopulation = $N,
PAinterval =(('req','percentile',95,
(1, 's')),
('pred','percentile', 95, $Cycle)) }

<<PAstep>>
{PAdemand=('asmd',
'mean', (0.5, 'ms'))}

<<PAstep>>
{PAdemand=('asmd',
'mean', ($P * 1.5, 'ms')),
PAextOp = (network, $P)}
getImage (i, b)

<<PAstep>>
{PAdemand=('asmd',
'mean', (1.8, 'ms))}

<<PAstep>>
{PAdemand=('asmd',
'mean', (0.9, 'ms'))}
passImage (i, b)

<<PAstep>>
{PAdemand=('asmd',
'meam', (1.1, 'ms'))}
storeImage (i, b)

<<PAstep>>
{PAdemand=('asmd',
'mean', (2, 'ms'))}
store (i, b)

<<PAstep>>
{PAdemand=('asmd',
'mean', ($B * 0.9, 'ms')),,
PAextOp=(writeBlock, $B)}
writeImg (i, b)

<<PAstep>>
{PAdemand=('asmd',
'mean', (0.2,'ms'))}
freeBuf (b)

<<PAstep>>
{PAdemand=('asmd',
'mean', (0.5, 'ms'))}

<<GRMrelease>>
releaseBuf (b)

<<PAcontext>>

CarletoN
UNIVERSITY

## Details of the UML annotations

Video Controller

AcquireProc

Buffer Manager

StoreP

<<PAstep>>
{PArep = $N}

<<PAstep>>
{PAdemand =('asmd', 'mean', (1.5, 'ms')}

*[$N] procOneImage(i)

O

getBuffer()

<<GRMacquire>>
allocBuf (b)

<<PAclosedLoad>>
{PApopulation = $N,
PAinterval =(('req','percentile',95,
(1, 's')),
('pred','percentile', 95, $Cycle)) }

<<PAstep>>
{PAdemand=('asmd',
'mean', (0.5, 'ms'))}

O

<<PAstep>>
{PAdemand=('asmd',
'mean', ($P * 1.5, 'ms')),
PAextOp = (network, $P)}

getImage (i, b)

O

<<PAstep>>
{PAdemand=('asmd',
'mean', (1.8, 'ms))}

<<PAstep>>
{PAdemand=('asmd',
'mean', (0.9, 'ms'))}

15

<<PAstep>>

# Spec of BSS (4) Sequence Diagram
# for the Door Access Control scenario



<<PAopenLoad>>
{PAoccurencePattern = ('poisson', 120, 's'),
PArespTime =(('req','per centile',95, (1, 's')),
('pred','percentile', 95, $RT))
}

<<PAcontext>>

User

CardReader  DoorLock  Alarm  Access Controller  Database  Disk

<<PAstep>>
{PAextOp=(read, 1)}
readCard

<<PAstep>>
{PAdemand=('asmd', 'mean', (1.8, 'ms')),
PAextOp = (network, 1)}
admit (cardInfo)

<<PAstep>>
{PAdemand=('asmd', 'mean', (1.8, 'ms'))}
getRights()

<<PAstep>>
{PAdemand=('asmd', 'mean', (1.5, 'ms')), PAprob = 0.4}

<<PAstep>>
{PAdemand=('asmd', 'mean', (1.8, 'ms'))}
readRights()   [not_in_cache] readData()

<<PAstep>>
{PAdemand=('asmd', 'mean', (0.3, 'ms'))}
checkRights()

<<PAstep>>
{PAdemand=('asmd', 'mean', (0.5, 'ms')), PAprob = 1}
[OK] openDoor()

enterBuilding

<<PAstep>>
{PAprob = 0}

<<PAstep>>
{PAdemand=('asmd', 'mean', (0.2, 'ms'), PAprob=0.2}
[need to log?] logEvent()

[not OK] alarm()

<<PAstep>>
{PAdemand=('asmd', 'mean', (3, 'ms'))}
writeRec()

<<PAstep>>
{PAdemand=('asmd', 'mean', (1.8, 'ms'))}
writeEvent()

16

## Simplified view of the main scenarios: VideoScan and Door Access Control

VideoController — Timer event starts loop over N cameras

AcquireFrame — cleanup, freeBuffer, storeInDB, passImage

getBuffer

getImage

CardReader — User [deny] [OK]

Access Controller — alert, log, [fail], checkRights, triggerLock

pollCamera

NetInterface — receivePackets

LockActuator — unlockDoor

Network delay

WebCamera — sendImage

Database

BufferManager

17 *AcquireFrame can use multiple buffers to receive and store images*

# Resources are identified from annotations

- **directly...**
  - "active object" or active component == a task
  - deployment of tasks
  - task allocation of component
  - logical resource acquire/release stereotypes

... and as attributes of the activities and components

  - processor for an activity
  - database queried by an activity
  - process or task containing a component
  - network conveying messages between components
  - buffer acquired by an activity getBuffer, released

# Task resources: AcquireFrame, and AccessController

- If *AcquireFrame* is single-threaded it sequentializes the acquisition

- double buffering requires another thread, or a concurrent task, let us call it Store



19

# The service time of a task covers its resource context

- includes lower servers

- here *AccessControl* includes lower layers:
    - LockActuator,
    - Database, and
    - Disk

- service time of AccessControl can be found recursively

- it includes waiting time at the database

- *alert* includes logging

20

Carleton
UNIVERSITY

# Layered queueing applies...

Flow of user requests

single server

competing database
requests for other
operations

| Access Control Task |
| :---: |
| [local processing] |

(requests)

delay but no queue,
as there is one actuator
per door)

common queue of
requests

| other | readRights | writeLogEvent | Database |
|---|---|---|---|

| unlock | LockActuator |
|---|---|

multiclass server
(possibly a multiserver also
= multithreaded task

delay server
(infinite server)

| read | write | Disk |
|---|---|---|

21

# Service time in layered queue systems...

... is not knowable without a full analysis

- because it includes contention at lower servers

- the queueing delay is affected by competing scenarios and applications

- this is the key difficulty in understanding performance in layered systems

  - for example, bottleneck location may be unstable

# Notation for layered queueing models

- by convention, call
  - all servers and resources *"tasks"*,
  - all resource-operations *"entries"*,
  - all requests for operations by a server *"calls"*.

- all requests to a task enter a *common queue*, which can have any discipline,
  - entries define *classes* of service
  - many kinds of tasks cannot support pre-emptive disciplines

- *synchronous calls* (that block the caller) are distinguished from asynchronous (that do not).
  - sync calls always lead to a single reply
  - more complex request types can be built up

23

# LQ Model (2)

- a sub-scenario defines *entry behaviour*, by a sequence of operations and requests.
  - can use a default stochastic model for entry behaviour:
    - random "slices", with a given coefficient of variation
    - either random requests with given mean numbers, or deterministic numbers of requests, in random order
- there is a *"host" processor* server for every task, not always shown
  - host service time is divided into slices between requests for other services
  - host servers have the same semantics as tasks, all requests are synchronous, every software entry generates an entry on its host,

24

# Simplified notation for a LQ model

- with default entry behaviour (random order of calls)
- parallelograms are optional....

entries with host demand in sec.

{25}  multiple resource
or multithreaded
task (25 threads)

| entry E1 [y-host-E1] | entry E2 [y-host-E2] | Client Task |

host attachment for processor

P1

Synchronous call with
mean number
(y-serv)

| entry for a service "serv" | Server Task |

P2

25

# LQN fragment for the Door Access Control

- "alarm" has no calls or load

- doesn't show database operations to store video frames

- some parameters show "second phase" operations: (0,0.2) or [3.9, 0.2]

  - after the reply

26

# Layered Queues are a *Canonical Generalization of Queueing Networks*

- includes ordinary QN as a single layer of client and servers
  - "program" entity calling its servers

- an LQN describes any **Extended Queueing Network** with nested resource use

- advantage: it easily describes a system with hundreds of logical resources, many held at once, in many patterns.

- it has  an **economical, concise** set of parameters for the stochastic default entry behaviour

27

# Different considerations for a logical resource: the Buffer pool for Video frames

- its service time includes parts of the execution of various tasks
  - it is not identified with any particular process
- it will be modeled by a "task" which we can call a pseudo-task
  - runs on a pseudo-host
  - no execution time of its own
  - makes calls that define the operation executed with the resource
- sometimes the place of the pseudo-task can be taken by an actual "resource manager" task, if it exists

28

# Path showing the holding time of a buffer

# Modeling technique for a logical resource

■ A buffer or critical section ….

Execution path

Wait for resource, then use it

Resource Context, e.g. Critical section CS

Leave context, release resource

**program instances executing outside of resource context**

Request and wait for resource RES

RES

execution inside of resource context (zero or one instance)

CPU1    DISK1    PRINTER    CPU2    DISK2

30

# Resource pseudo-task when user tasks are distributed...

- Tasks A and B must enter a critical section (call it CS) for some work....
  - but this doesn't express what they do within CS

- So:
  - Separate out the computation within CS into *Shadow Tasks A/CS* and *B/CS*
  - to direct the call from A to A|CS, make CS a *pseudo-task* with two *pseudo-entries*



*not a suitable model*

31

# LQN fragment for the video buffer



- AcquireFrame has a fragment Acquire2 within the buffer resource context

- BufferManager also

- call to StoreImage is *second phase*

procOneImage [1.5,0] | AcquireFrame (2 threads)

(1,0)

alloc [0.5, 0] | BufferManager

bufEntry | Buffer

(1,0)  (1,0)  **(0,1)**

getImage [12,0] | passImage [0.9, 0] | Acquire2    storeImage [3.3, 0] | StoreProc

($P, 0)

(1,0)  (1,0)

network [0, 1] | Network    releaseBuf [0.5, 0] | BufMgr2    writeImg [7.2, 0] | readRights [1.8,0] | writeEvent [1.8, 0] | DataBase

($B,0)  (0.4,0)  (1,0)

DB CPU

writeBlock [1, 0] | readData [15,0] | writeRec [3, 0] | Disk

32

Carleton
UNIVERSITY

# "Second Phase Service" in software servers

- ■ Idea: often used to enhance performance
    - ▪ give a reply as early as possible
    - ▪ Do postponeable work after the reply, as "phase 2"

client

server

phase 1          phase 2, asynchronous and parallel

- ■ e.g..:  Database server update operation:
    - ▪ write to log file before returning, execute final writes later.

- ■ Second-phase model may
    - ▪ place this work right after the return (approx), or
    - ▪ send a message to a clean-up process that does it later

33

- ■ Queueing approximation paper in Performance 99

Carleton
UNIVERSITY

# Overall LQN: $N cameras, $R buffers



34

**BSS
LQN...
follow
the
scenarios**



35

# Results for varying buffers.

| Buffers | Cameras | Prob(miss) |
|---------|---------|------------|
| 1 | 10 | 0 |
| 1 | 20 | 0.001 |
| 1 | 30 | 0.417 |
| 3 | 10 | 0 |
| 3 | 20 | 0 |
| 3 | 30 | 0.003 |
| 3 | 40 | 0.319 |

# Building layered models

- identify resource contexts

- identify interaction patterns between them
    - synchronous, asynchronous, forwarding

- synchronous: request-reply from nested resources

# Building layered models (2)

- *forwarding* shows execution passed directly from one server to another, without returning in between
  - may traverse a route before returning



4 in sequence:

38

# Building layered models (3)

- forking to a new context is an *asynchronous* interaction

# Interaction summary

- Resource == "task"

- Resource-operation (simple): one activity....

- requests for other resources:

  - expressed via interaction types: sync, async, forward

# Activity sequence detail in an entry

- an activity has a workload (CPU demand and requests)
- sequence relationships
    - (AND/OR fork/joins)
- interaction types again: sync, async, forward



41

# An example with activity sequence detail

*multiple (25) UserT entities with default activity and its host demand*

{25}

userE [8 sec]

UserT

*host association*

UserP  inf

(1)  *1 sync call*

*one server for each client*

appE

return

compute[AppE] [45 ms]

log

cleanup

AppT

AppP

(2) *2 sync calls*

(1)

(1)  *async call*

getData [4]

cleanup [3,8]

ServerT

ServerP

logE [12]

LogT

# Example: visualize the scenario

*multiple UserT entities with default activity and its host demand*

userE
[8 sec]

UserT

*host association*

UserP

inf

*one server for each client*

(1)  *1 sync call*

appE

*return*

log

compute[AppE]
[45 ms]

cleanup

AppT

AppP

(2) *2 sync calls*

(1)

(1)  *async call*

getData
[1]

cleanup
[3,8]

ServerT

logE
[12]

LogT

ServerP

43

# LQN pattern for parallel operations

■ Server A requests services from S1 and S2 in parallel, that is it sends both requests and then waits for both replies.

■ This happens once during an execution of an entry A:

*Some resource context patterns...*

# Other patterns of resource contexts

Pipeline

Chaotic, unstructured

**Task1**   **Task2**   **Task3**   **Task4**

| RA | RA |    |    |
| RB | RD | RX |    |
| RC | RE | RF | RC |

A separate context for each activity

Resource pass-back:
... an interesting pattern that needs work...... the Buffer is released by the Agent

**Resource pipeline (above), and sliding overlapping resource contexts over time**

Program
Server
Buffer Manager
Buffer
Agent

45

# Summary of model-building from software descriptions

- ■ can be seen as a generalization of the methods defined by Connie Smith

- ■ based on tracing scenarios, detecting resources and interpreting nesting and interaction types
  - ▪ from scenarios in *Use Case Maps*: TOOLS 2002 paper
  - ▪ from tracing (*"angio traces"*): MASCOTS 95 and TSE 2000 (Hrischuk)
    - ▪ ("TLC" = trace-based load characterization)
  - ▪ now analyzing *UML scenarios*, expressed with the UML Profile on Schedulability, Performance and Time (2002)

Carleton
UNIVERSITY

# Proposed PUMA toolset architecture......
## (Performance by Unified Model Analysis)



- *general* software model input via CSM (not only UML)

- *general* performance model types via CPM (not only layered queues)

- includes heavy element of model investigation, sensitivity tools, optimization

- proposal also for component libraries for completions
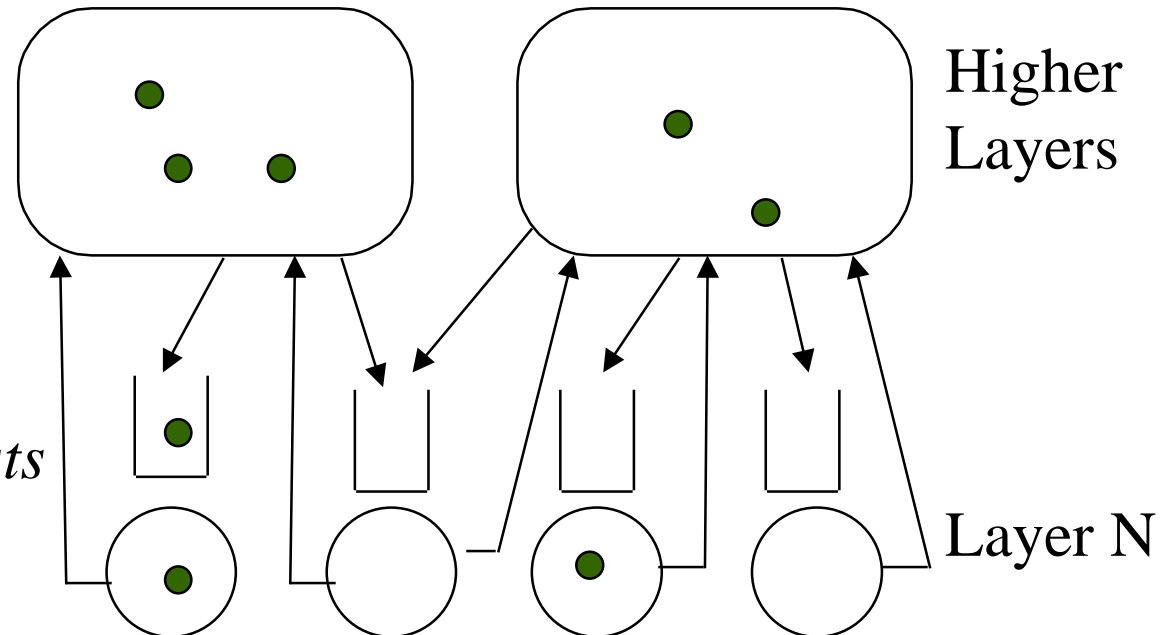
47

# LQ network *solvers*

- ***www.layeredqueues.org*** site to provide resources
- LQNS (Franks/Rolia/Petriu/Sevcik/Woodside) (84 on)
  - iterative basic MVA with lots of approximations for variance, multiservers, parallel paths, and other aspects
- Fontenot described one open layered server (Sigmetrics 1988)
- Petriu (1993) Markov model by TDA "Task Derived Aggregation"
- Ramesh/Perros (98 - 2000)
  - open systems, with close attention paid to variance effects, and a structured sequence of classes of service
- Kahkipuro (UML2000), a basic multilayer solution
- Menasce (2 layers) (2002) for critical sections, etc.

48

Carleton
UNIVERSITY

# LQNS: Iterative MVA Solver for LQN
## "Active Server" closed queueing sub-model for each layer (1984)

*"Delay Server" with tokens that represent clients or active servers from layers above*

Higher Layers

*Tokens represent  requests from upper layer servers*

Layer N

*Servers represent server tasks at layer N:*
- *Service times of the servers include delays at lower layers, including processors*
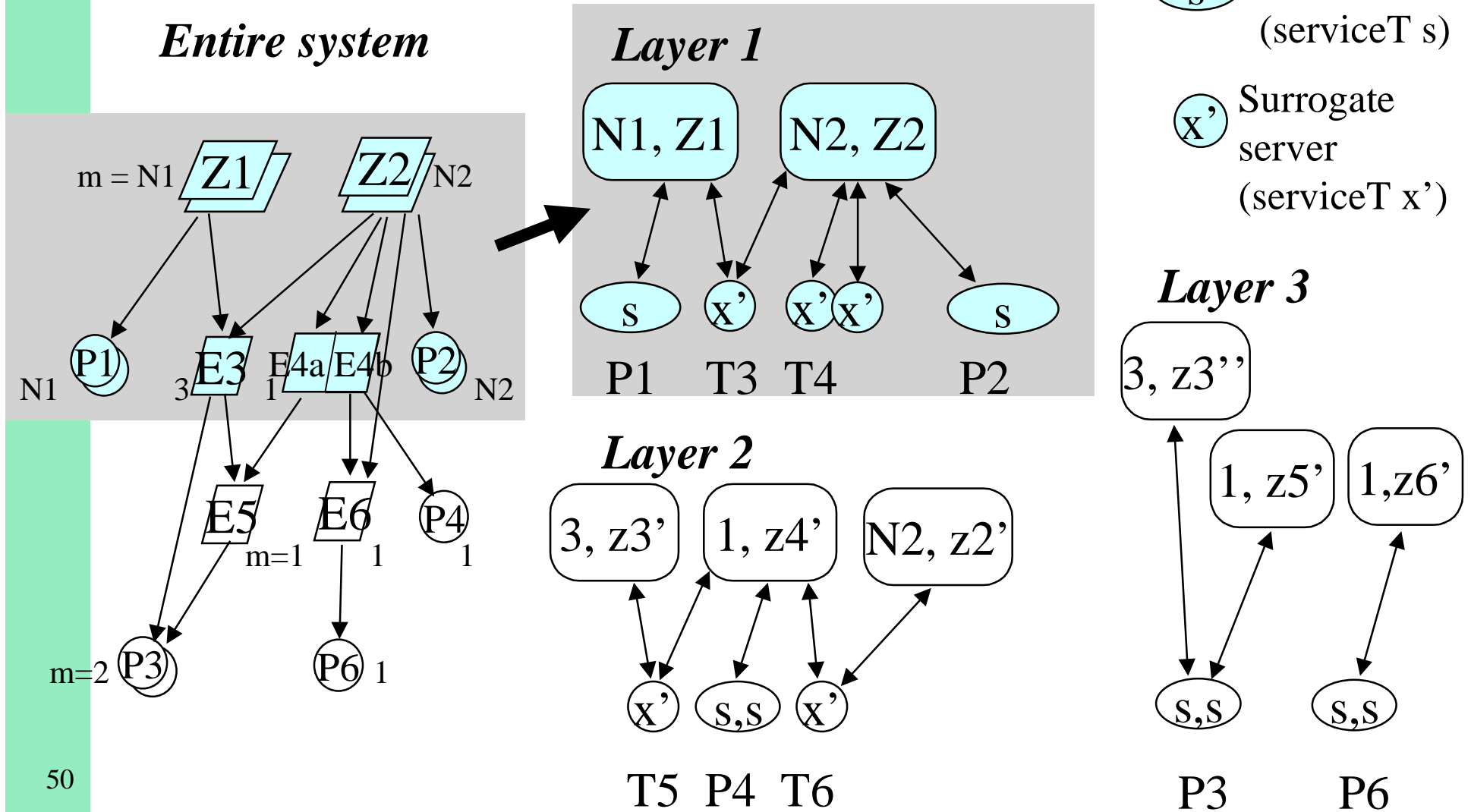- *Servers may be non-standard  (two phase!)*

49

# "Layerize": Layer 1 Submodel

Carleton UNIVERSITY

| n,z | Infinite server (tokens n, thinkT z) |

S — Processor (serviceT s)

x' — Surrogate server (serviceT x')


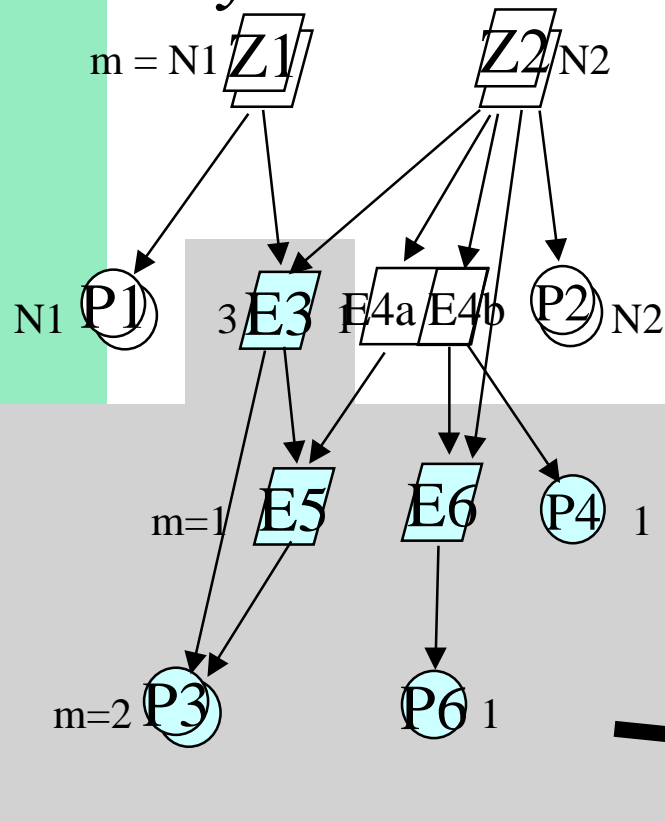
*Entire system*

*Layer 1*

*Layer 2*

*Layer 3*

# "Layerize": Layer 2 Submodel



**n,z** Infinite server (tokens n, thinkT z)

**s** Processor (serviceT s)

**x'** Surrogate server (serviceT x')

51

Carleton
UNIVERSITY

# "Layerize": Layer 3 Submodel

n,z — Infinite server (tokens n, thinkT z)

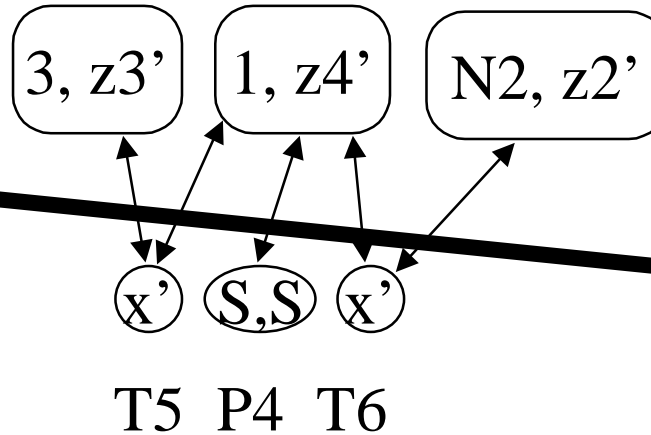S — Processor (serviceT s)

x' — Surrogate server (serviceT x')



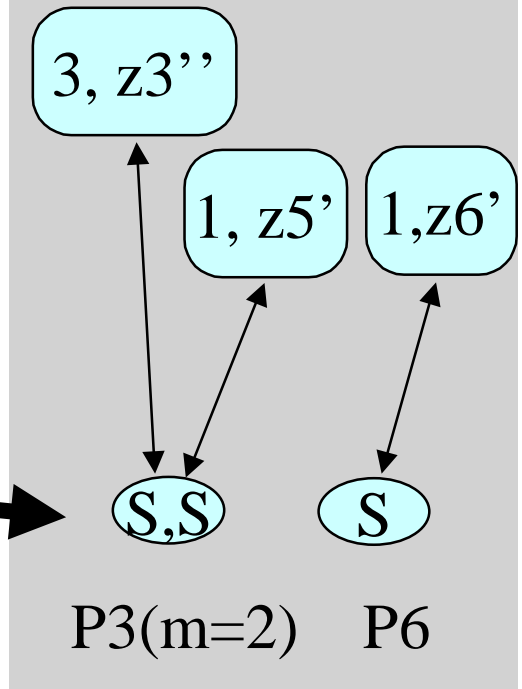*Entire system*

*Layer 1*

*Layer 2*

*Layer 3*

52

# Layerizing strategies can be interesting
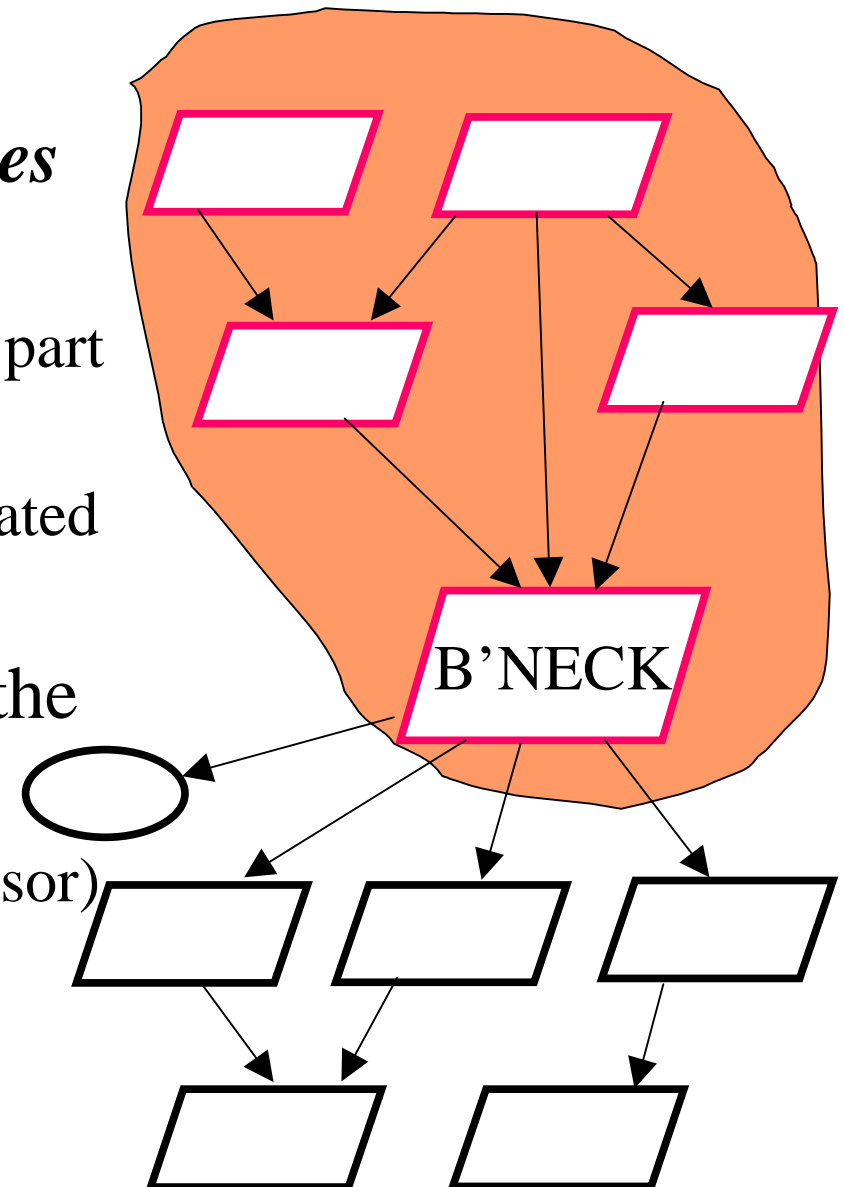
- Woodside 1984/ 88/ 95...one task per layer
  - calls can jump layers (causes a dependence effect "interlocking"
- Rolia/Sevcik 1988/92/95... wide layers, greedy from the bottom.
  - strict layering (calls cannot jump over)
  - accommodation for jumping over
- Ramesh/Perros 1998... strict layering
- Franks thesis 1999.... in LQNS... flexible choice, OR greedy from the top, OR all in one big layer (!!)
  - balanced layer sizes are best for solution time.
  - a detailed submodel for dependence effect
  - detailed study of second phases

# Exploiting the solutions to improve the software design

- sensitive parameters point to aspects with most leverage
  - sensitive execution demands can be reduced
  - sensitive interaction counts can be reduced
  - long entry service times can perhaps be parallelized
- sensitivity is highest around bottlenecks!!
  - look at resource utilization
    - task utilization includes all its nested service times when it is blocked
    - reduce service time of bottleneck server
  - increase resource multiplicity (buffers, threads)

54

# Recognizing layered bottlenecks

- a saturated server

- but.... a saturated server *pushes back* on its clients

  - the long waiting time becomes part of the client service time!!

  - result is often a cluster of saturated tasks above the bottleneck

- thus: the "real" bottleneck is the *"lowest"* saturated task

  - its servers (including its processor) are not saturated

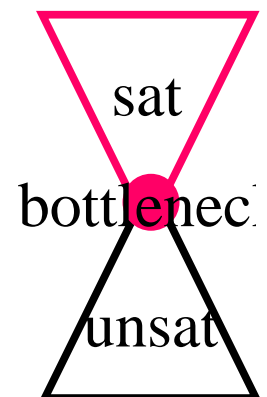  - some or all of its clients are saturated



B'NECK

# Recognizing the "real" bottleneck

- a saturated task with unsaturated servers and host
- *Strength* measure $(U_B/M_B)/[\max (U_S/M_S)]$
- IEEE TSE paper 1995

- *Notice that:*
  - if the bottleneck task has no servers, its host utilization is the same as the task (it only computes)
  - so it must have at least one additional server, a device (e.g. disk), task, or other logical server
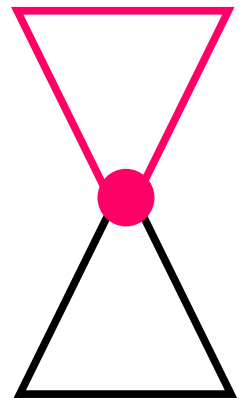  - also, it must have sufficient clients to build a queue
- thus, there is often an hourglass pattern

sat

bottleneck

unsat

56

# Bottleneck patterns and *threads* or multiplicity

- a task with M threads counts as M concurrent servers or clients
    - in identifying the "hourglass" pattern
- in the "strength" measure, a server with M threads saturates at $U = M$
- a (very rough) rule of thumb for threads, based on potential needs for concurrency:

$$M = \text{min of} \begin{cases} \text{(sum of server threads, +1)} \\ \text{(sum of client threads)} \end{cases}$$

# Curing a bottleneck
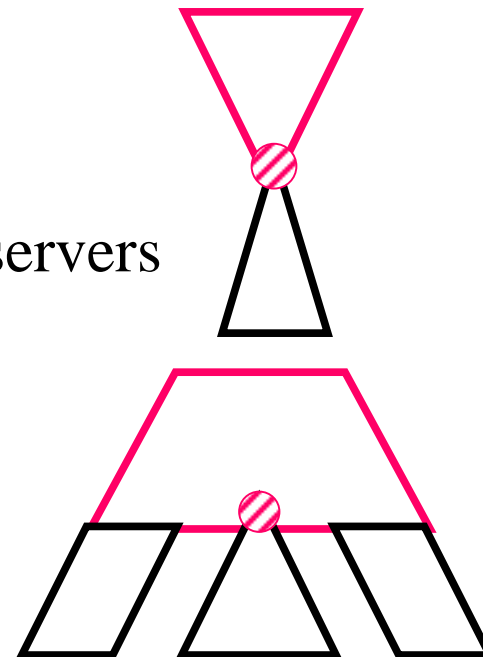
(1)  provide additional resources at the bottleneck

- for a software server, provide *multiple threads*
    - some "asynchronous server" designs provide unlimited threads
- replicated servers can split the load and distribute it
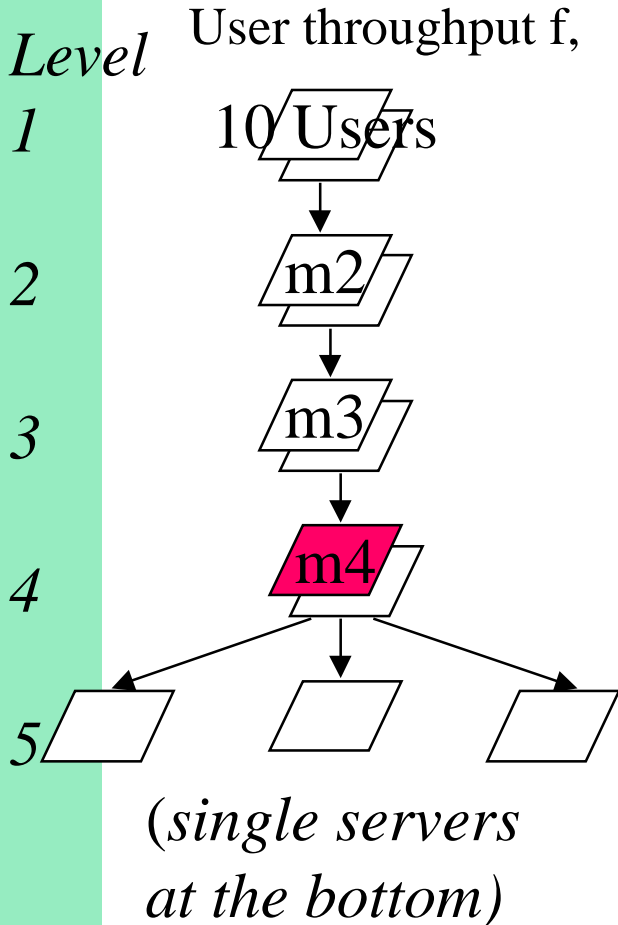- for a processor, a *multiprocessor* (or faster CPU)

(2)  reduce its service time:

- reduced host demand
- reduced requests to its servers

(3) divert load away from it

# Curing Software bottlenecks by multithreading

*Level*

*1*   User throughput f,

10 Users

*2*   m2

*3*   m3

*4*   m4

*5*

*(single servers
at the bottom)*

- bottleneck at task 4 limits the user throughput f
- f depends on the threads at all servers
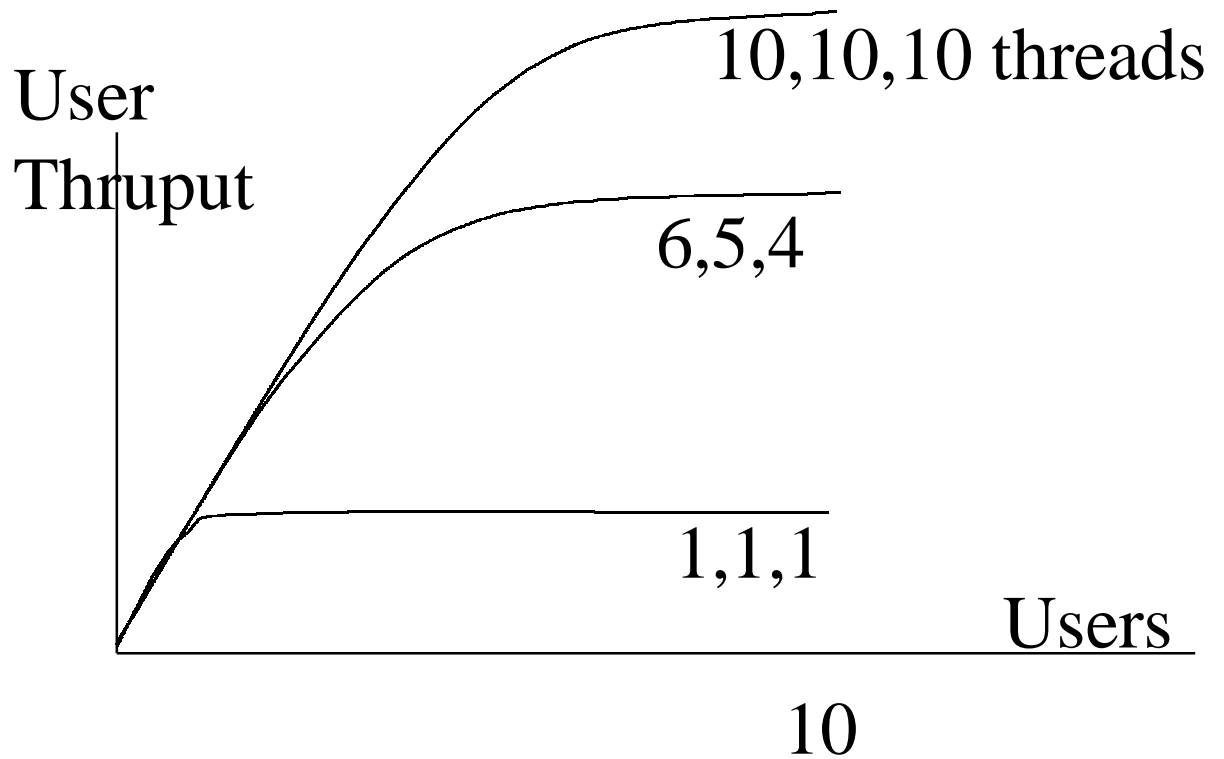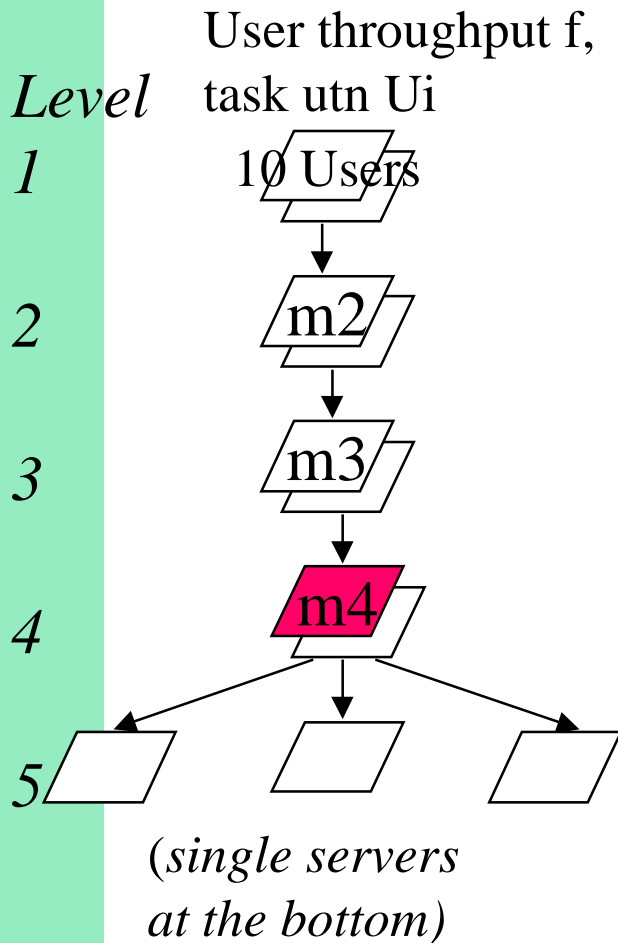  - m2 threads for task2, m3 for task 3, etc

...a multi-threaded server behaves like a multi-server; two threads can execute in parallel. *If they are sequentialized by their processor servers, that appears as waiting*

- 1 sec host demand at each server, one request to each lower task
- Ui = task utilization at level i

| (m2, m3, m4)... | f | (U2, | U3, | U4, | U5) |
|---|---|---|---|---|---|
| (1, 1, 1)..... | 0.166, | (1. | 0.83 | 0.67 | 0.167) |
| (2, 1, 1)..... | 0.200, | (0.96 | 1. | 0.8 | 0.2) |
| (3, 2, 1).... | 0.223, | (2.9 | 1.64 | 0.89 | 0.22) |
| (6, 5, 4).... | 0.475, | (5.5 | 3.9 | 2.75 | 0.475) |
| (10,10,10).. | 0.65, | (9.3 | 7.8 | 6.2 | 0.65) |

59

Carleton
UNIVERSITY

# Software bottleneck relief by multithreading

User throughput f,
task utn Ui

*Level*

*1* — 10 Users

*2* — m2

*3* — m3

*4* — m4

*5*

*(single servers at the bottom)*

*(1 sec demand at each server, one request to each lower task)*

User Thruput

10,10,10 threads

6,5,4

1,1,1

Users

10

Carleton UNIVERSITY

# Bottleneck: Results for a web server with net delay

N Users with
a thinking time
of 5 sec.

Users

Server with M threads and
holding time X

Server

0.005   0.2   0.4   1

CPU   Disk   DB   Net delay
0.015   0.01   0.5 sec

| N users | 500 | 500 | 500 | 500 | 2000 | 2000 | 2000 | 2000 |
|---|---|---|---|---|---|---|---|---|
| M threads | 10 | 30 | 100 | inf | 10 | 30 | 100 | inf |
| X server | .512 | .52 | .52 | .52 | .512 | .515 | .55 | 4.99 |
| f thruput | 19.5 | 58.2 | 90.6 | 90.6 | 19.5 | 56.7 | 180 | 200 |
| W user wait | 20.6 | 3.6 | 0.51 | 0.5 | 97.6 | 29.4 | 6.1 | 5 |
| U server | 10 | 30 | 47 | 47 | 10 | 30 | 100 | 1000 |
| U net | 9.7 | 29.1 | 45.3 | 45.3 | 9.7 | 29.1 | 90.2 | 100 |
| U CPU | .097 | .29 | .45 | .45 | .097 | .29 | .90 | 1.0 |

61

# Multiple independent bottlenecks

- there may be a web of servers and interactions

- perhaps there are multiple bottlenecks?
  - in a flat queueing network there can be as many independent bottlenecks as there are chains of customers
  - each is an independent limitation on chain throughputs
- in a layered queueing network there are only a few independent throughputs... e.g. the top-layer tasks
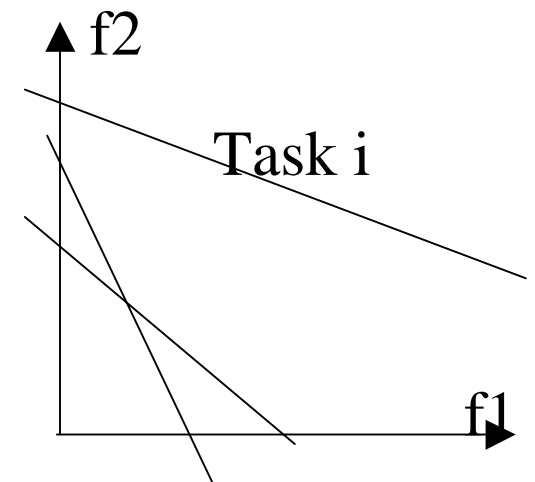
# Using bound relationships in layered queues

- reference throughputs $f_r$ at the "user" top-layer tasks
- all other throughputs (fe at entry e) are proportional, $f_e = \Sigma_r\, a_{re}\, f_r$
- no-waiting service time of entry e is $x_e$ which can be computed recursively using nested delays
- reference throughputs must satisfy the utilization constraint $U_i < M_i$ ,

$$\Sigma_{e\ in\ Task\ i}\ f_e\, x_e < M_i$$

$$\Sigma_r\ f_r\, \Sigma_{e\ in\ Task\ i}\ a_{re}\, x_e < M_i$$
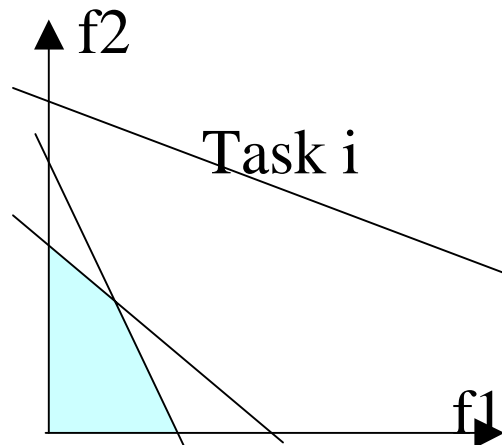
$$\Sigma_r\ f_r\ K_{ir} < 1$$

- *"rendezvous nets" paper 1995*

63

# Bound relationships are only a crude guide to bottleneck location

- ■ resources giving bounds that touch the feasible region are likely candidates
  - ▪ other bounds, e.g. for Taski, are prevented from saturating



- ■ however the bounds are not tight
  - ▪ because they ignore queueing delays at intermediate levels
  - ▪ since queueing delay can create a bottleneck... it really needs a full queueing solution

64

# Converting results to software implications

- long response or service times can be reduced by
  - parallelizing some operations
  - balancing and reduced variability in parallel paths
  - latency masking (e.g. by pre-fetching)
  - optimistic design
  - removing bottlenecks within the response
- bottlenecks can be reduced by...
  - host demands reduced, server demands reduced
  - demands made more deterministic
  - changed allocations
  - replication, threading
  - task splitting for concurrent access (servers, pipelining...)
- navigation of sensitive points (drill-down) *(Maps and Paths paper 1995)*

# Acting on the recommendations

- changes to software architecture and detailed design
- *reducing demands* is a well studied topic, e.g. Smith and Williams books
  - detailed code changes based on hot spots, locality, early binding of references
  - caching
- some of the other recommendations relate to *"performance antipatterns"* described by the same authors (WOSP 2000)
  - the "one lane bridge" is any bottleneck task
  - the "god class" is a task that can be split into smaller parts

# Summary

- the layered queueing model is a middle ground between software structure and queueing networks
  - default stochastic semantics have few parameters
  - scalable extended queueing canonical form
- fairly direct traceability of
  - software tasks to performance model objects
  - object interactions into model interactions
  - demands
  - results connected back to software observations
- similarity between model results and measurements on the software

# Where is this area going?

- solvers still pose open questions
  - improvements to accuracy and to features

- support for building models
  - models from UML
  - models created or updated from monitoring

- integration with discrete-state modeling methods
  - failure states (IPDS paper 98, others)
  - adaptation and variable configurations
  - submodels for inter-task protocols, using Petri Nets etc
  - submodels for more accurate delay distributions

- optimization (e.g. Sigmetrics 2001)