# Performance Guarantees for Web Server End-Systems: A Control-Theoretical Approach

### Written by Tarek F. Abdelzaher and Kang G. Shin

Presented by Chen Zhang

2006-11-15

# Performance Guarantee for Web Server

- What do we guarantee here?
  - Response time: server side delay
  - Throughput: individual hosted sites
- With …
  - Performance isolation
  - Service Differentiation
  - QoS adaptation
- Why?
  - Premium VS Basic (client view)
  - Overload/Overprovision (Server view)
- How? – Feedback control: control utilization

# The situation

- When a client visits a host…
  - Rejected … ☹
    - Always?? ☹!!! – still alive? -- robustness
  - Accepted … ☺?? Wait..
    - Not always ☹ -- availability
    - Low response, takes a long time ☹.. -- responsiveness
    - Fluctuated QoS ☹.. – stability
    - No pictures?? ☹.. -- satisfaction
- When a server receives a new client request...
  - Overloaded?
  - Is it Premium??
  - Reject or adapt?
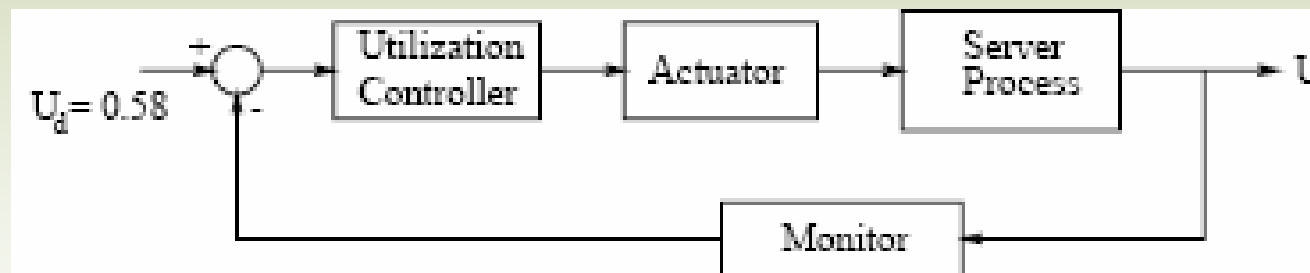
# System model and goal

- Pre-word: Guarantee ONLY for Premium!

- Basic Scenario

  - Client sends URL request to a certain host

  - Each URL refers to a host on a machine hosting various Virtual Hosts

  - Each client hold a "service level agreement"

- Server side Goals

  - Less reject when overloaded

  - Guarantee response time and throughput

| Consider for fun |
| --- |
| How many Premiums/Basics are allowed?  All Premiums are not supposed to come at the same time… |

# Control

- **<u>Target</u>**: utilization. U < 0.58. Implicit guideline for **threshold** without knowing per client **load**.
- PI controller, settings depend on
- At expense of reject, all admitted get QoS
- Control loop per class,for performance isolation



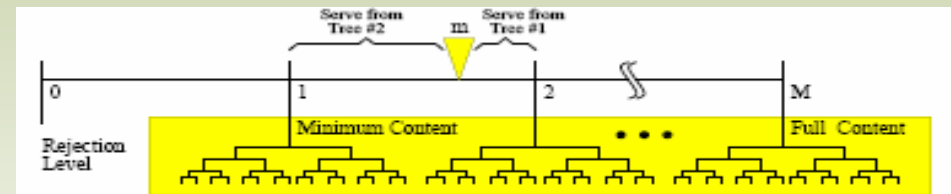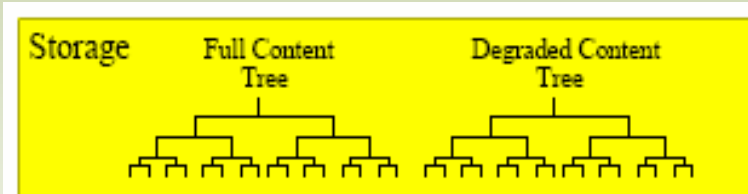| Note |
| --- |
| No detail about control function here. Just keep in mind that it is PI controller on U that both response quick to change and accurate, taking advantage of both P controller (u(t) = KPe(t)) and I controller (u(t) = u(t − 1) + KIe(t)). |

# Monitor

- Under the assumption of static load:
  - $U = aR + bW$
    - R: Served request rate
    - W: delivered byte bandwidth
- If only a fraction f of requests are admitted
  - $U = aRf + bW + cR(1 - f)$
    - R: Received request rate
    - W: Total byte bandwidth

| |
|---|
| **Think: Is U (R, W) a sufficient and sound measure?** |
| Cost of serving a file $C_i = a + bX_i$, a fixed overhead plus a variable overhead proportional to the length $X_i$ of the file. |

# Actuator

- **<u>Role</u>**: Translate control output to server action
- **<u>Novelty</u>**: simple admission control → degradation
- **<u>Degradation rationale</u>**: Small size VS large size
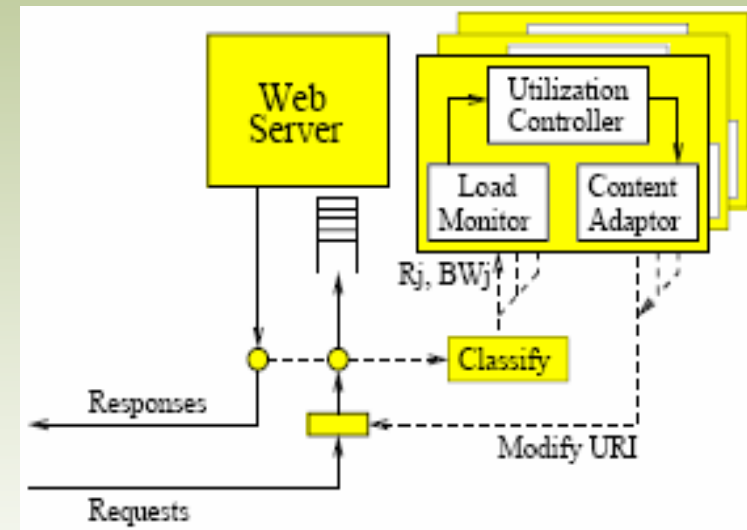- Discrete service level [0, M], mapping $m = I + F$



| Concern: | |
|---|---|
| 1. | vital VS trivial. More about criteria? Some vital information could be of large size but must not be simply degraded out when facing overload, e.g., when evacuating, people need maps, and maybe only a single copies of maps for a room of people rather than many copies verbal descriptions. |
| 2. | URL rewrite Cost? When frequent fluctuation occurs, frequent rewrite Automatically maintain consistency? Who is going to generate and maintain the "smaller" copy? |
| 3. | Who/how to decide how many service levels? How to maintain? |

# Performance Isolation

- For a single virtual server
  - Rmax, Wmax, throughput guarantee
- Capacity planning
  - $Ui^* = aRmax + bWmax$
  - $U = \sum Ui^* < 0.58$



| |
|---|
| **Think: One overload can cause Ui up and U up** |
| The importance of having Rmax and Wmax per host in order to achieve performance isolation. |

# Others

- Service differentiation
  - Per host control loop, Ui, mi
- Sharing Excess Capacity
  - Using extra system resource if
    - A certain Ui overloads
    - Overall U not overloaded
  - mb – best effort
  - mi – for a specific virtual host
  - Request is handled using max(mi, mb)

**Think about mb ..**

How to measure idle resource into mb? Will the first overload virtual server take advantage of all extra resource? Used to advertise as a Daily Bonus? :P

# Several issues

- Requires Server code modification normally
- Issues as mentioned ..
  - How many Premiums/Basics are allowed?
  - Is U (R, W) a sufficient and sound measure?
  - vital VS trivial when degrade. More about criteria?
  - URL rewrite Cost? When frequent fluctuation occurs
  - Automatically maintain consistency? Who is going to generate and maintain the "smaller" copies? How?
  - Who/how to decide how many service levels?
  - How to measure idle resource into mb?
  - Will the first overload virtual server take advantage of all extra resource?
  - ……