

CONGESTION CONTROL

Phenomenon: when too much traffic enters into system, performance degrades

→ excessive traffic can cause congestion



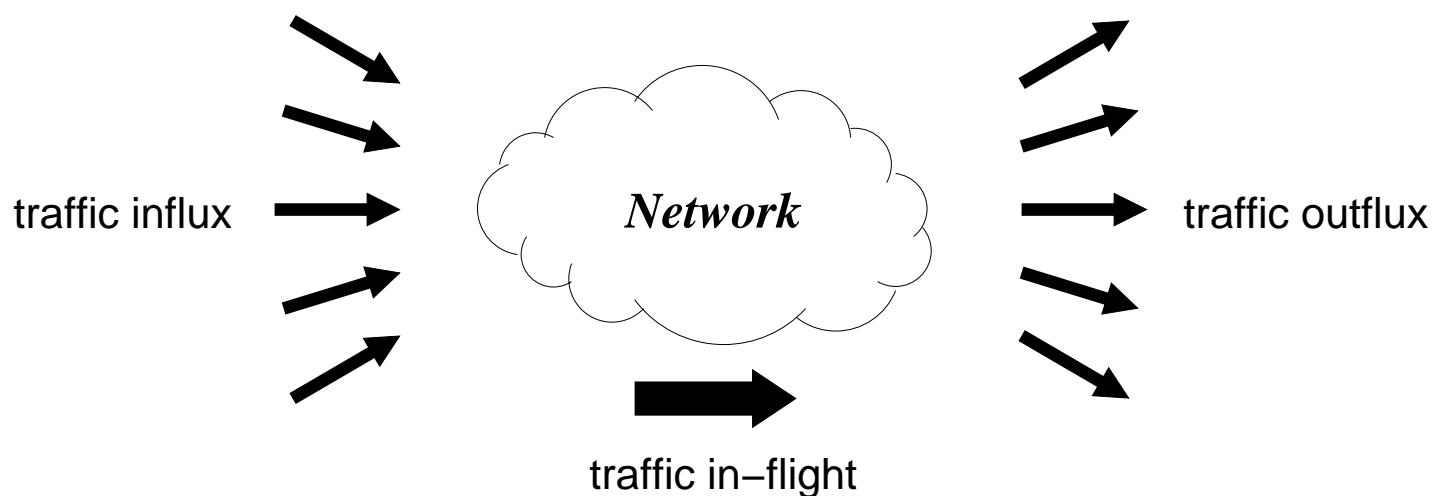
Problem: regulate traffic influx such that congestion does not occur

→ not too fast, not too slow

→ congestion control

→ first question: what is congestion?

Viewpoint: traffic coming in, in transit, going out



At time t :

- traffic influx: $\lambda(t)$ “offered load” (bps)
- traffic outflux: $\gamma(t)$ “throughput” (bps)
- traffic in-flight: $Q(t)$ “load”
→ volume: total packets in transit (no. of packets)

Examples:

Highway system:

- traffic influx: no. of cars entering highway per second
- traffic outflux: no. of cars exiting highway per second
- traffic in-flight: no. of cars traveling on highway

→ at time instance t



California Dept. of Transportation (Caltrans)

Water faucet and sink:

- traffic influx: water influx per second
- traffic outflux: water outflux per second
- traffic in-flight: water level in sink

→ “congestion?”

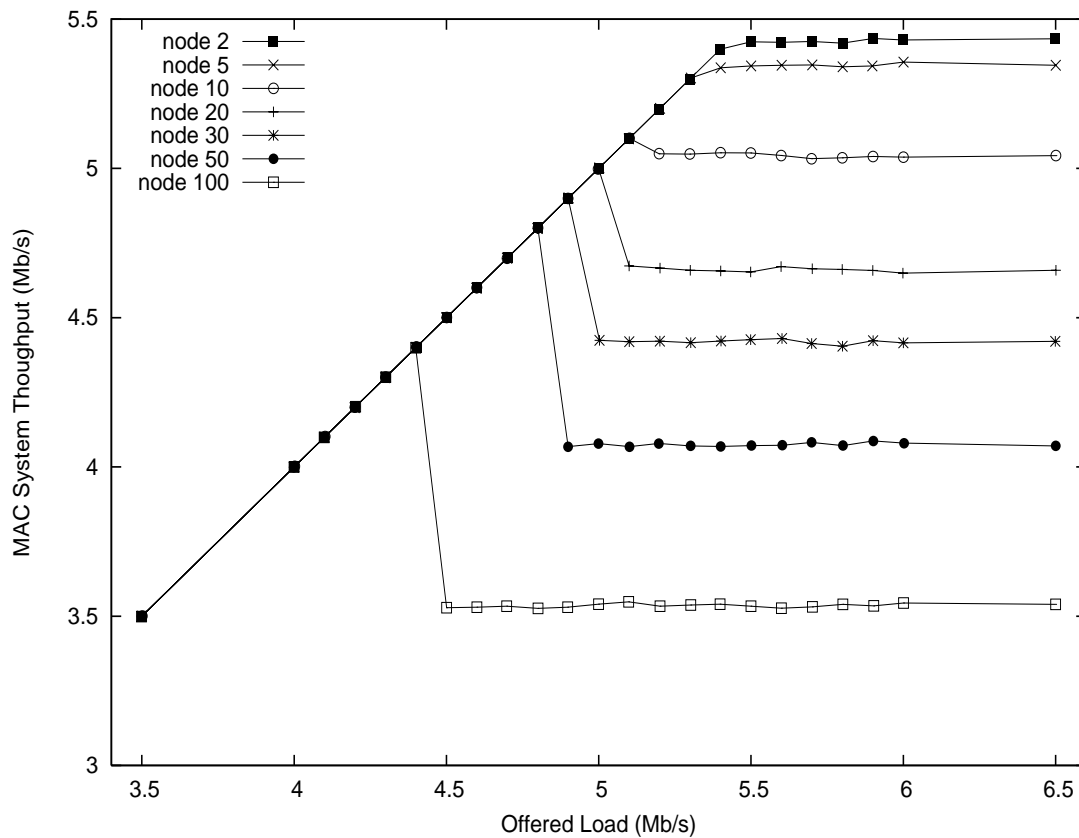


faucet.com

Thermostat ...

802.11b WLAN:

● Throughput



→ unimodal or bell-shaped

→ recall: less pronounced in real systems

What we can control:

- traffic influx rate $\lambda(t)$
- no power over anything else

Ex.:

- Faucet knob in water sink
- Temperature needle in thermostat
- Cars entering onto highway: traffic light
- Packets entering the Internet
 - from web server, P2P server, PC, laptop/handheld

How does in-flight traffic or load $Q(t)$ vary?

→ obeys simple rule

Compare two time instances t and $t + 1$.

At time $t + 1$:

$$Q(t + 1) = Q(t) + \lambda(t) - \gamma(t)$$

- $Q(t)$: what was there to begin with
- $\lambda(t)$: what newly arrived
- $\gamma(t)$: what newly exited
- $\lambda(t) - \gamma(t)$: net influx (+ or -)
- $Q(t)$ cannot be negative: no. of packets
→ $Q(t + 1) = \max\{0, Q(t) + \lambda(t) - \gamma(t)\}$
- missing item?

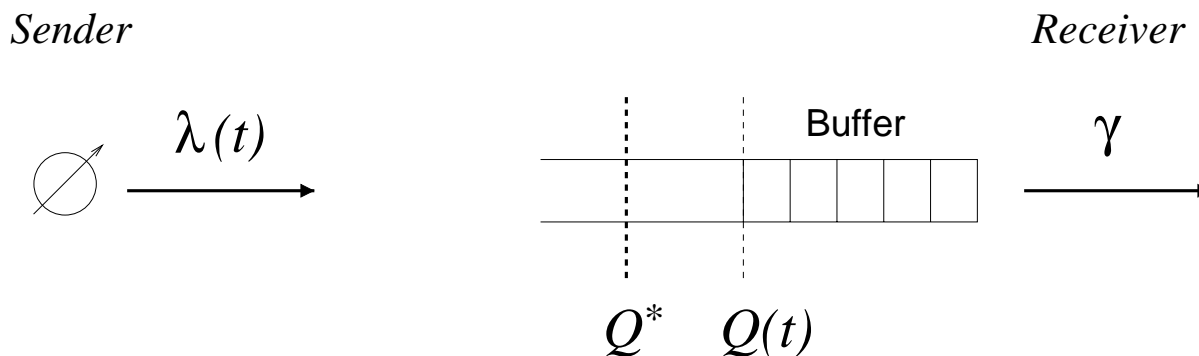
Pseudo Real-Time Multimedia Streaming

- e.g., RealPlayer, iTunes, Internet radio
- “pseudo” because of prefetching trick
- application is given headstart: few seconds
- fill buffer & prevent from becoming empty

Steps involved:

- prefetch X seconds worth of audio/video data
- causes initial delayed playback
 - e.g., couple of seconds delay after click
- keep fetching audio/video data such that X seconds worth of future data resides in receiver's buffer
 - hides spurious congestion
 - user: continuous playback experience

Pseudo real-time application architecture:



- $Q(t)$: current buffer level
- Q^* : desired buffer level
- γ : throughput, i.e., playback rate
→ e.g., for video 24 frames-per-second (fps)

Goal: vary $\lambda(t)$ such that $Q(t) \approx Q^*$

→ don't buffer too much: memory cost

→ don't buffer too little: cannot hide congestion

Other applications:

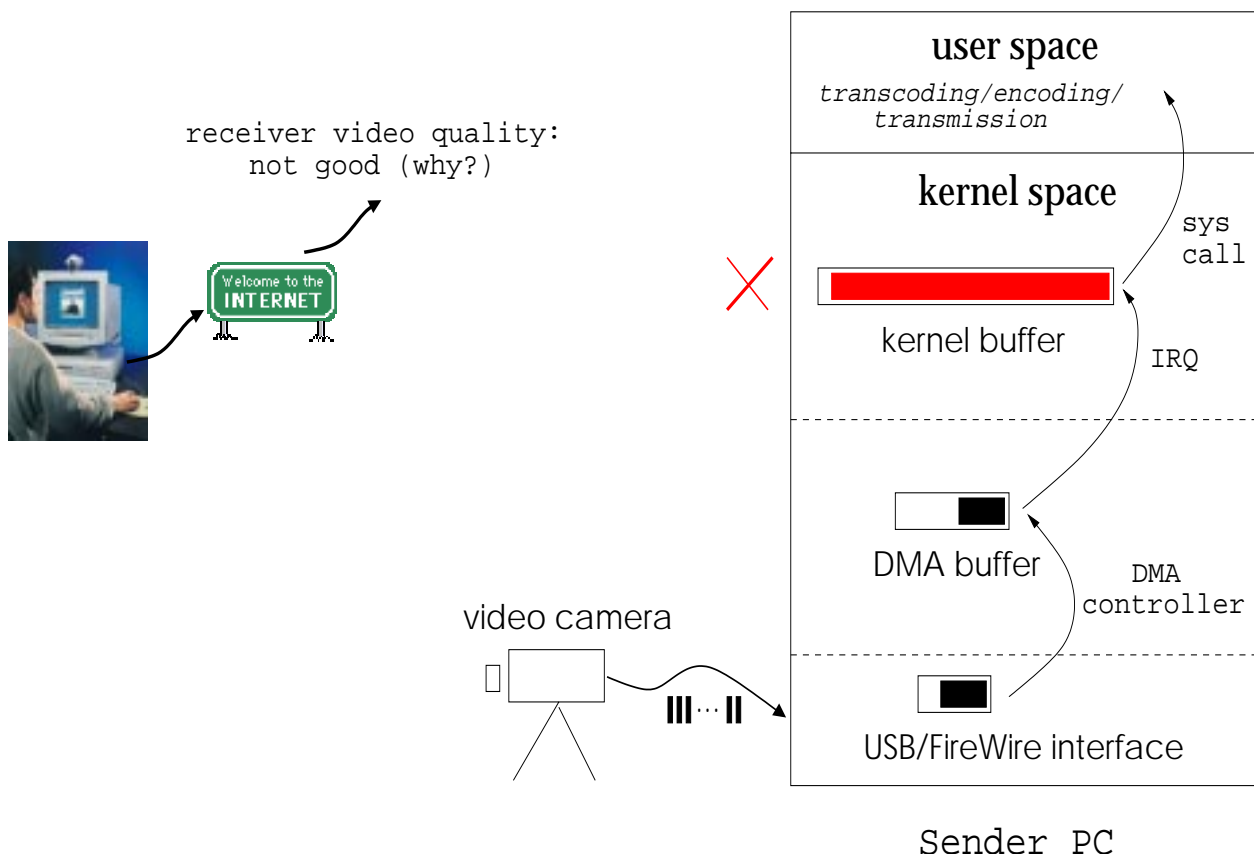
- pseudo real-time set-up is highly versatile
- captures many scenarios

Ex. 1: Router congestion control

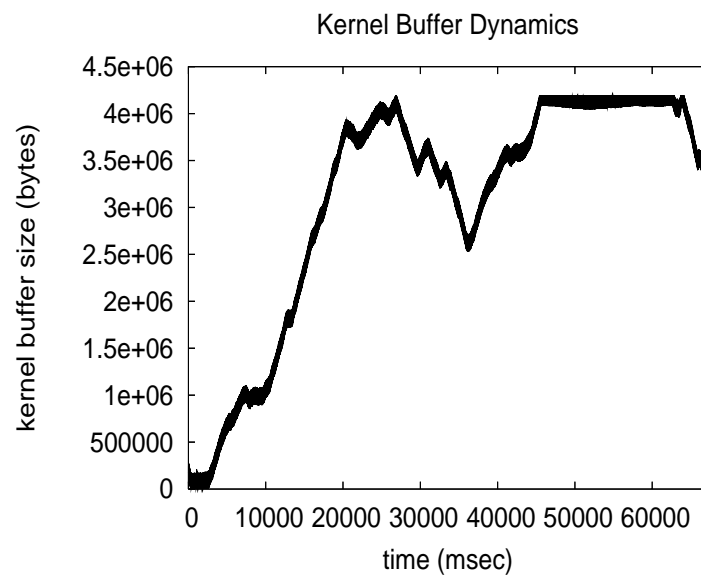
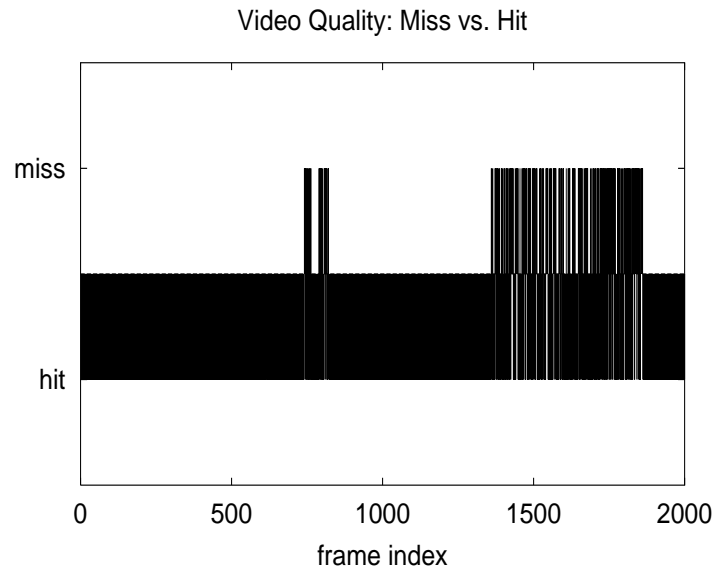
- active queue management (AQM)
- receiver is a router/switch
- Q^* is desired buffer occupancy/delay at router
- router throttles sender(s) to maintain Q^*
 - send control packets to senders
 - slow down, go faster, stay put

Ex. 2: Desktop videoconferencing

- e.g., AOL, MSN, Skype, Yahoo
- video quality is not good: why?
- misconception: network is blamed



Performance consequences:



Thus: pseudo real-time multimedia streaming application of congestion control

- producer/consumer rate mismatch problem
- also called “flow control”

Note: producer/consumer problem in OS

- focus on orderly access of shared data structure
- mutual exclusion
- e.g., use of counting semaphores
- necessary but insufficient

What is the goal?

→ achieve $Q(t) = Q^*$

How to: basic idea

- if $Q(t) = Q^*$ do nothing
- if $Q(t) < Q^*$ increase $\lambda(t)$
- if $Q(t) > Q^*$ decrease $\lambda(t)$

→ a rule of thumb

→ called “control law”

Network protocol implementation:

- some design options available
- control action undertaken at sender
 - smart sender/dumb receiver
 - preferred mode in many Internet protocols
 - when might the opposite be better?
- receiver informs sender of Q^* and $Q(t)$
 - feedback packet (“control signaling”)
 - or simply $+/-$ indication (binary)
 - or actual gap $Q^* - Q(t)$

Receiver sends feedback to sender; sender takes action

- called feedback control
- or closed-loop control

Key question in feedback congestion control:

- **how much** to increase/decrease $\lambda(t)$
- we already know in which direction

Desired state of the system:

$$Q(t) = Q^* \text{ and } \lambda(t) = \gamma$$

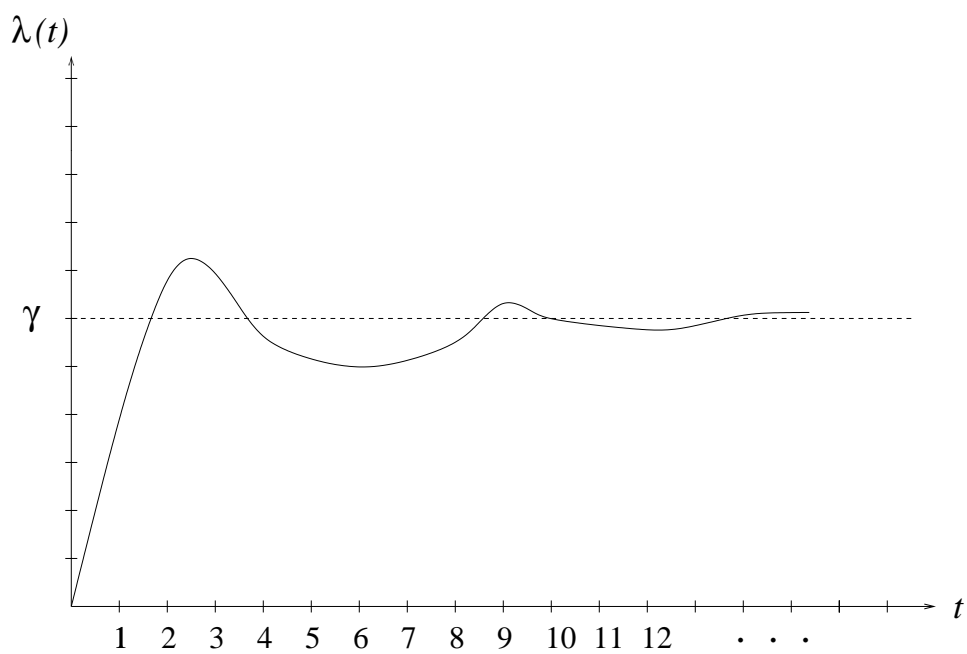
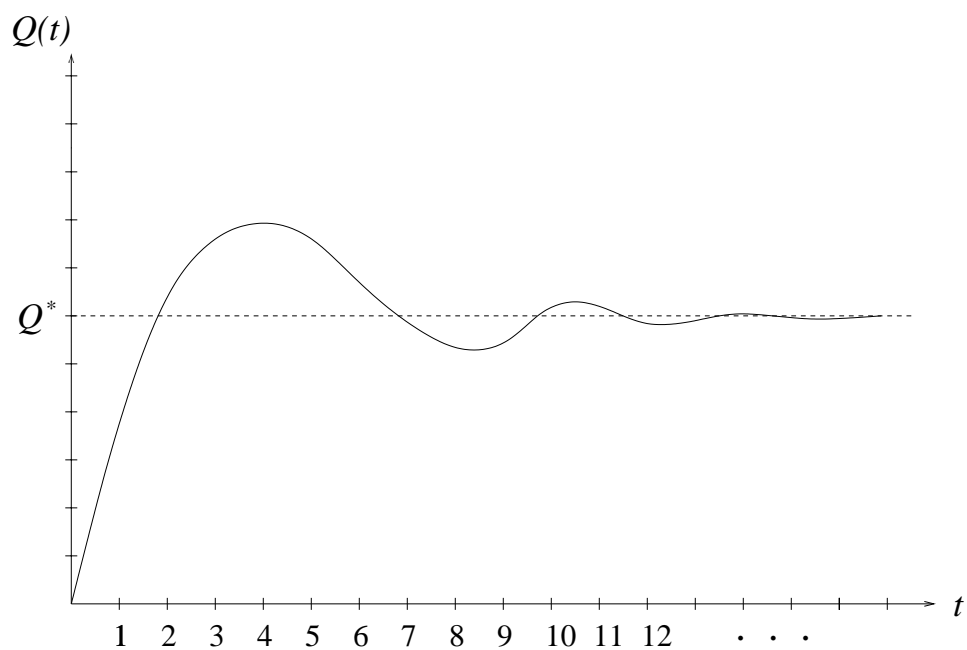
- why is “ $\lambda(t) = \gamma$ ” needed?

Starting state:

- empty buffer and nothing is being sent
- think of iTunes, Rhapsody, etc.

$$\text{i.e., } Q(t) = 0 \text{ and } \lambda(t) = 0$$

Time evolution (or dynamics): track $Q(t)$ and $\lambda(t)$



Congestion control methods: A, B, C and D

Method A:

- if $Q(t) = Q^*$ then $\lambda(t + 1) \leftarrow \lambda(t)$
- if $Q(t) < Q^*$ then $\lambda(t + 1) \leftarrow \lambda(t) + a$
- if $Q(t) > Q^*$ then $\lambda(t + 1) \leftarrow \lambda(t) - a$

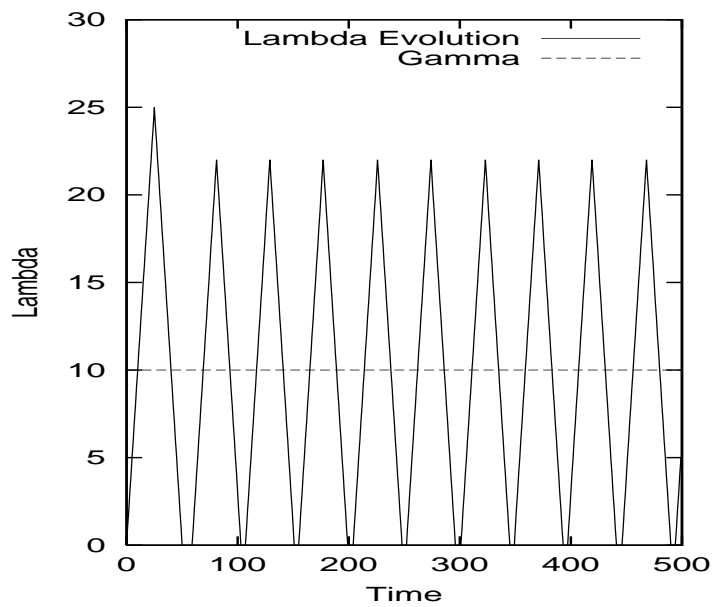
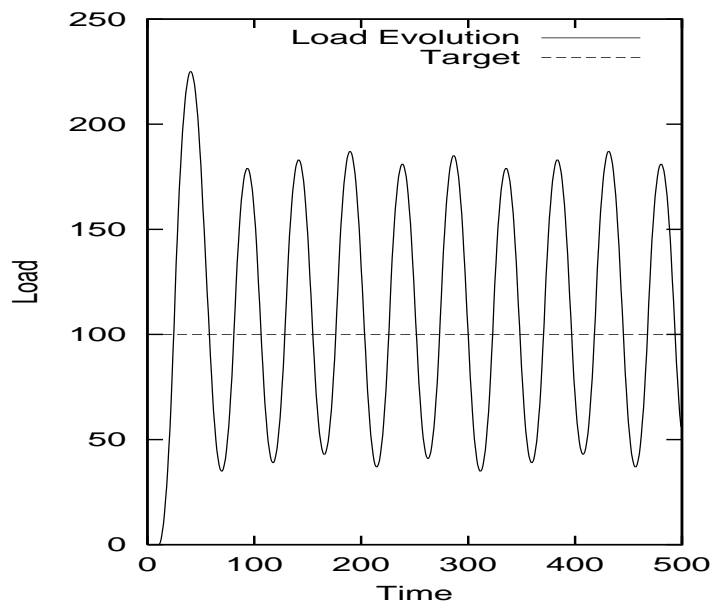
where $a > 0$ is a fixed parameter

→ called linear increase/linear decrease

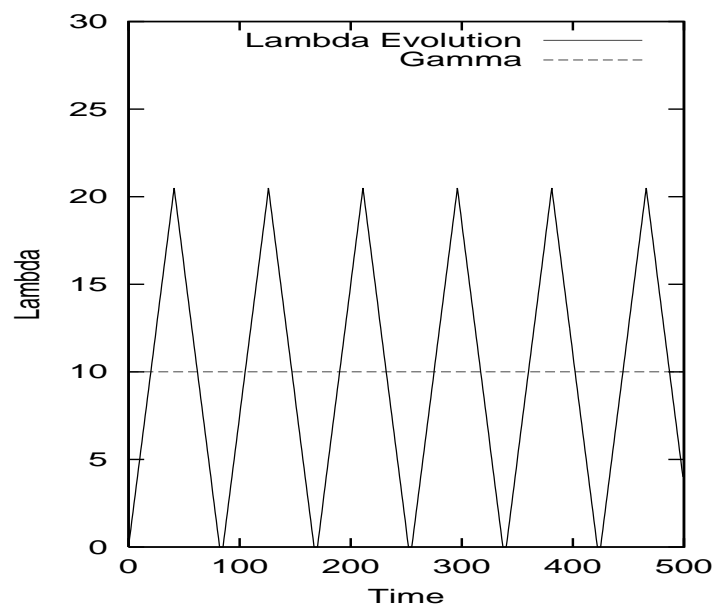
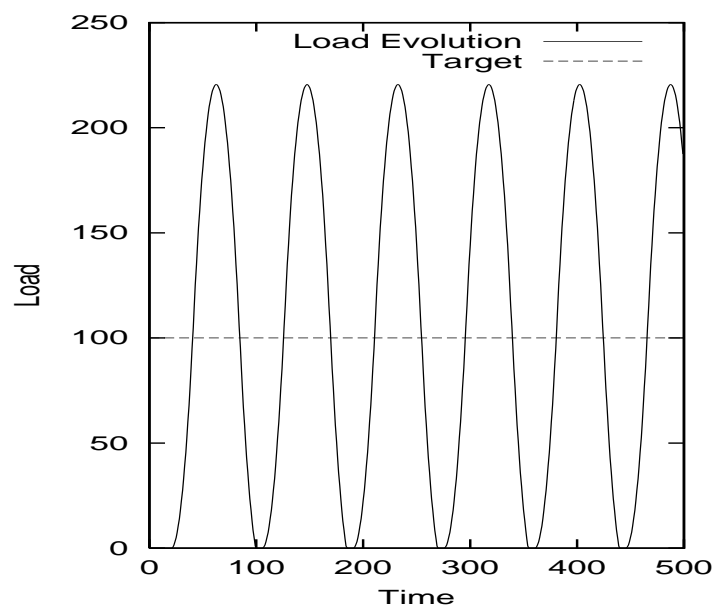
Question: how well does it work?

Example:

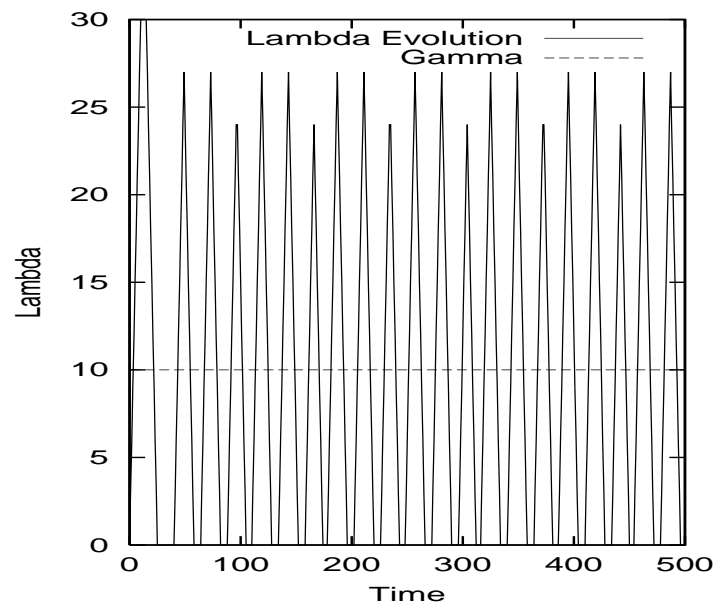
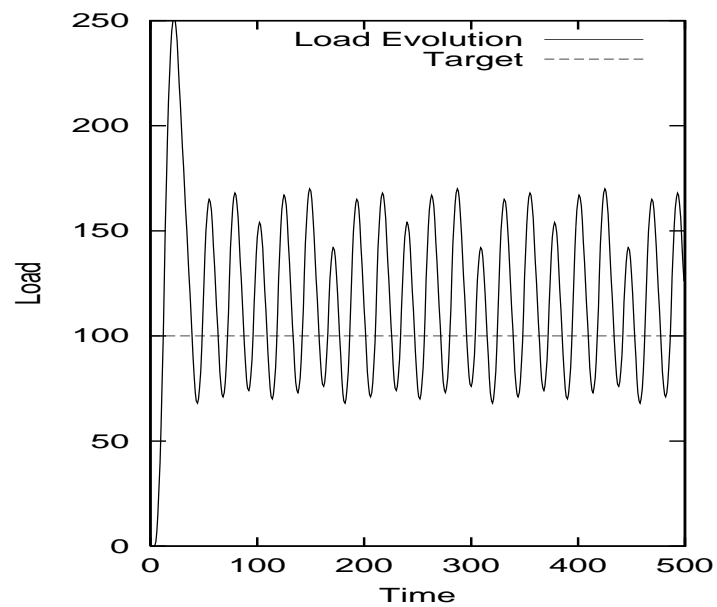
- $Q(0) = 0$
- $\lambda(0) = 0$
- $Q^* = 100$
- $\gamma = 10$
- $a = 1$



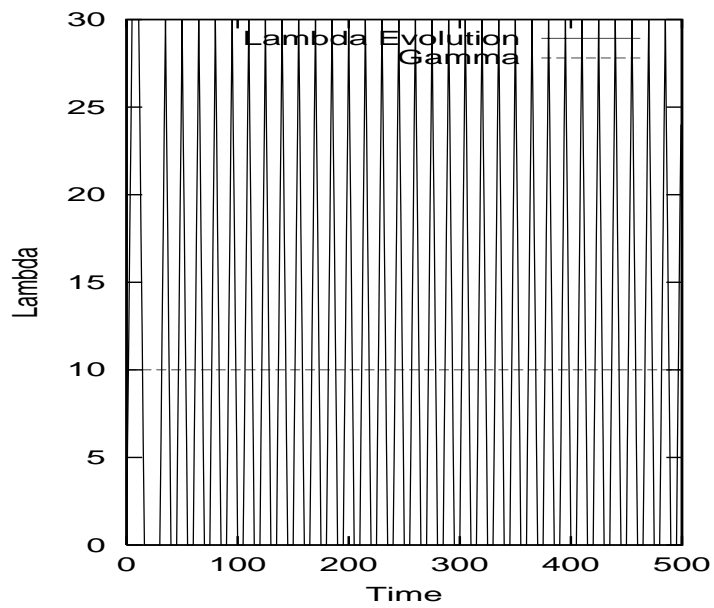
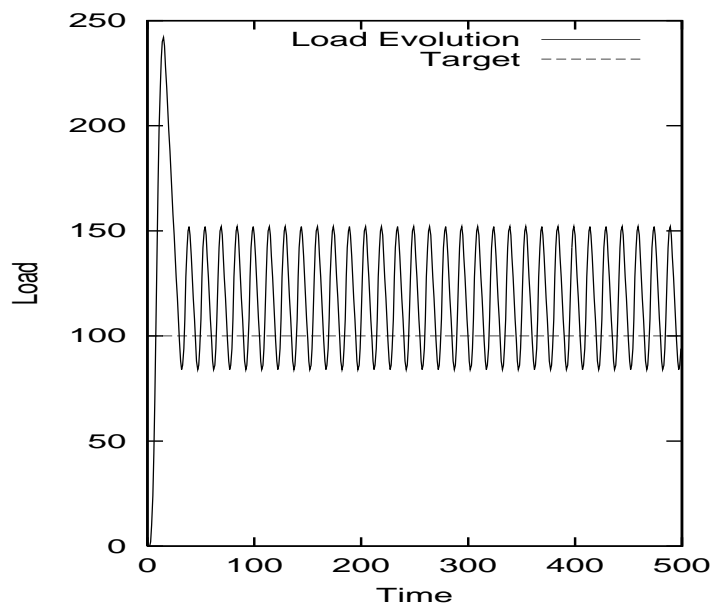
With $a = 0.5$:



With $a = 3$:



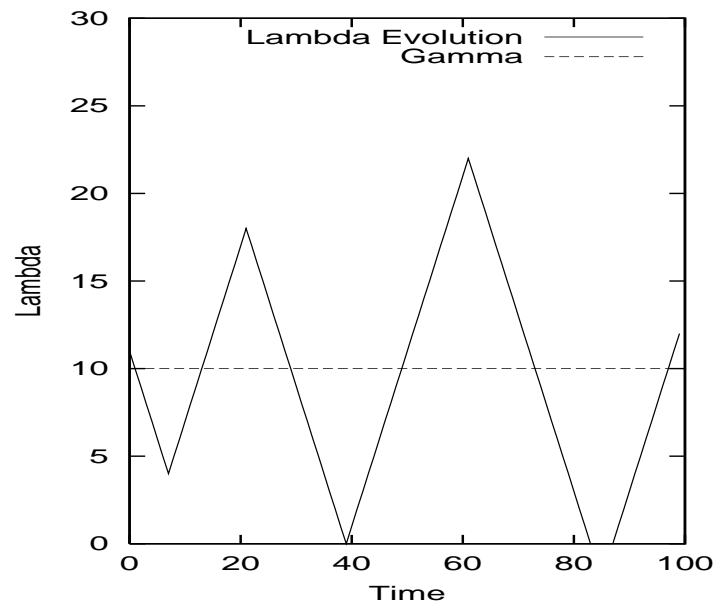
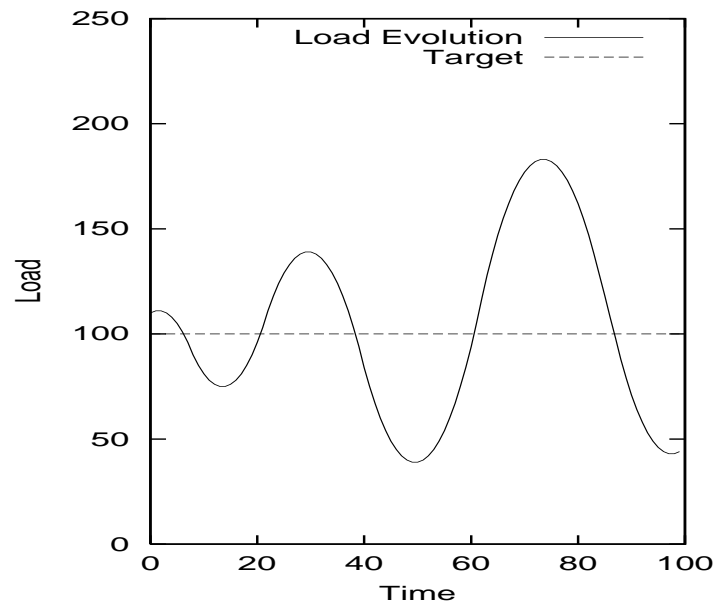
With $a = 6$:



Remarks:

- Method A isn't that great no matter what a value is used
 - keeps oscillating
- Actually: would lead to unbounded oscillation if not for physical restriction $\lambda(t) \geq 0$ and $Q(t) \geq 0$
 - i.e., bottoms out
 - easily seen: start from non-zero buffer
 - e.g., $Q(0) = 110$

With $a = 1$, $Q(0) = 110$, $\lambda(0) = 11$:



Method B:

- if $Q(t) = Q^*$ then $\lambda(t + 1) \leftarrow \lambda(t)$
- if $Q(t) < Q^*$ then $\lambda(t + 1) \leftarrow \lambda(t) + a$
- if $Q(t) > Q^*$ then $\lambda(t + 1) \leftarrow \delta \cdot \lambda(t)$

where $a > 0$ and $0 < \delta < 1$ are fixed parameters

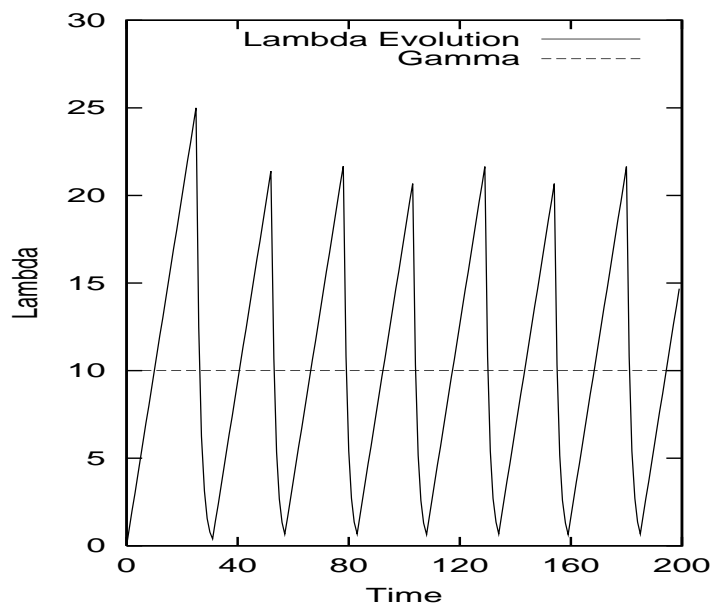
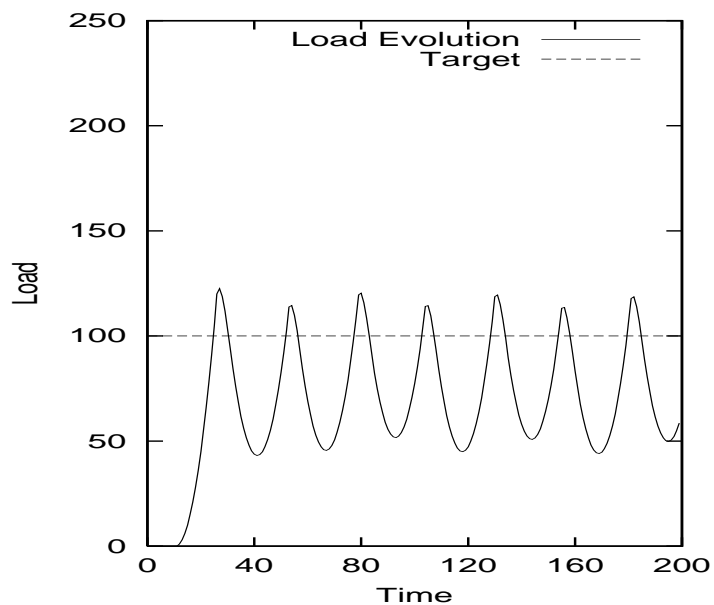
Note: only decrease part differs from **Method A**.

- linear increase with slope a
- exponential decrease with backoff factor δ
- e.g., binary backoff in case $\delta = 1/2$

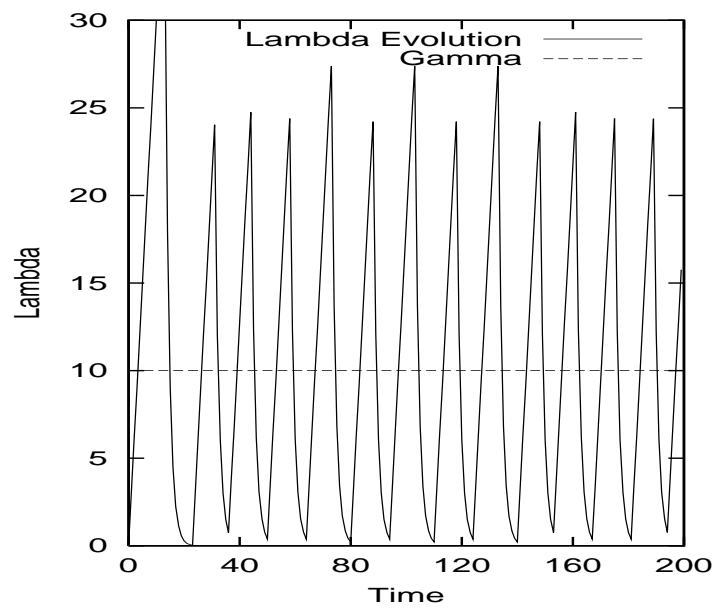
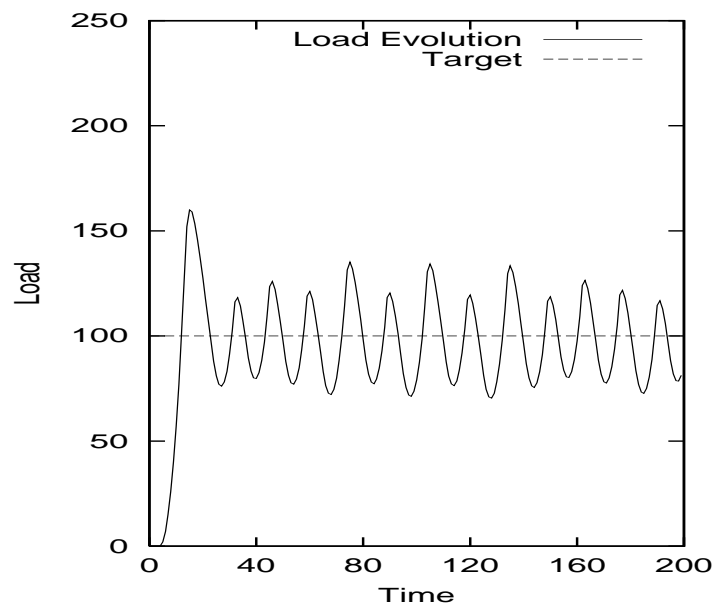
Similar to Ethernet and WLAN backoff

- question: does it work?

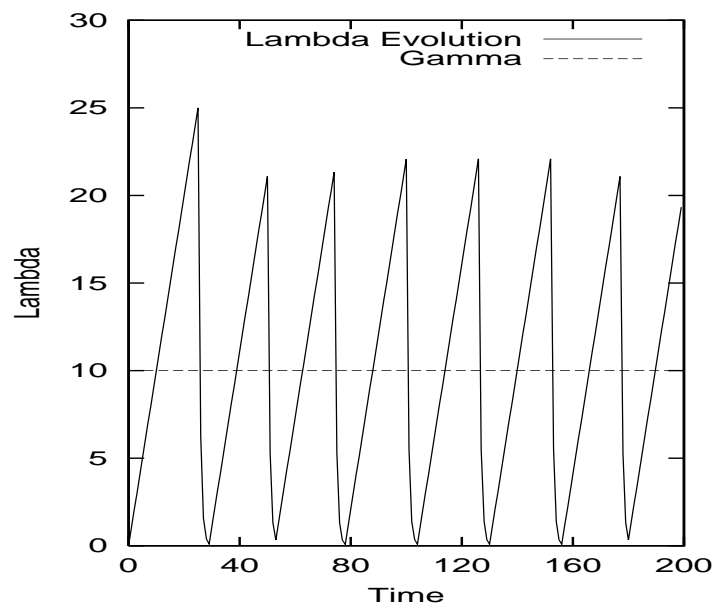
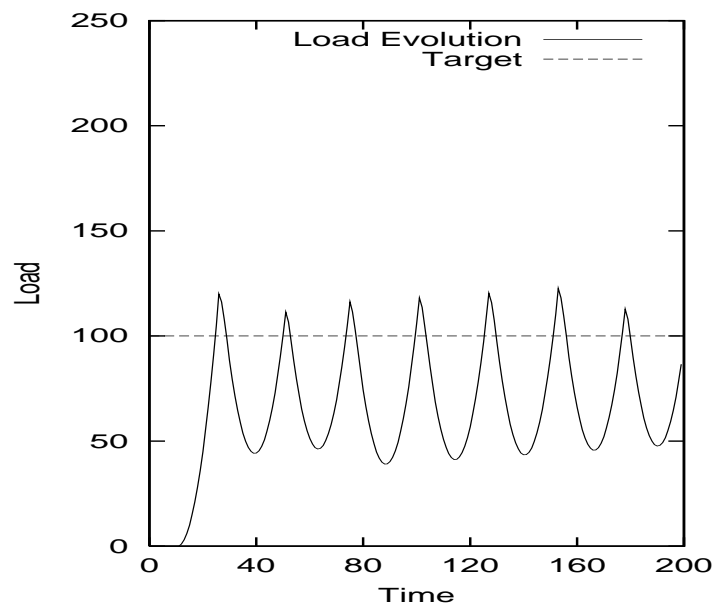
With $a = 1$, $\delta = 1/2$:



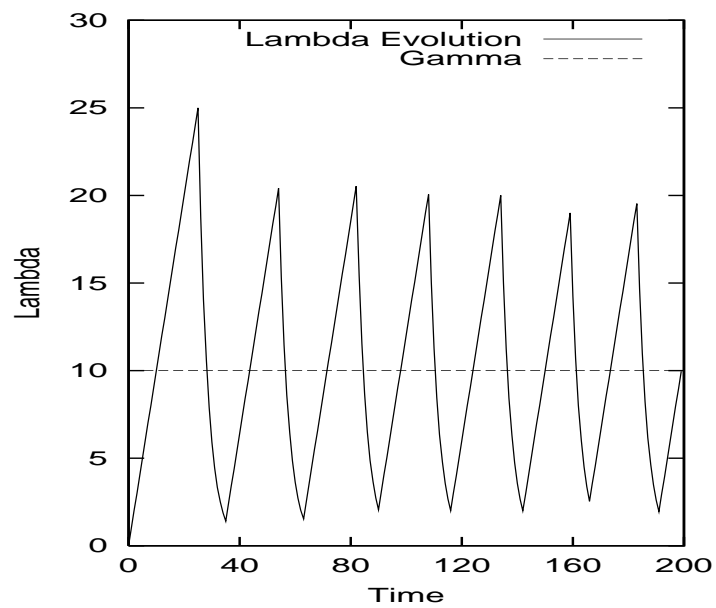
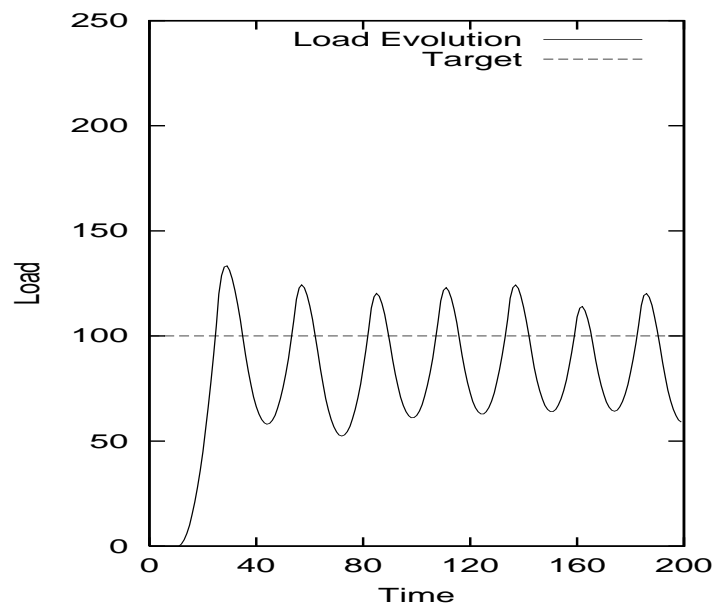
With $a = 3$, $\delta = 1/2$:



With $a = 1$, $\delta = 1/4$:



With $a = 1$, $\delta = 3/4$:



Note:

- Method B isn't that great either
- One advantage over Method A: doesn't lead to unbounded oscillation
 - note: doesn't hit "rock bottom"
 - due to asymmetry in increase vs. decrease policy
 - we observe "sawtooth" pattern
- Method B is used by TCP
 - linear increase/exponential decrease
 - additive increase/multiplicative decrease (AIMD)

Question: can we do better?

→ what "freebie" have we not made use of?

Method C:

$$\lambda(t + 1) \leftarrow \lambda(t) + \varepsilon(Q^* - Q(t))$$

where $\varepsilon > 0$ is a fixed parameter

Tries to adjust magnitude of change as a function of the gap $Q^* - Q(t)$

→ incorporate distance from target Q^*

→ before: just the sign (above/below)

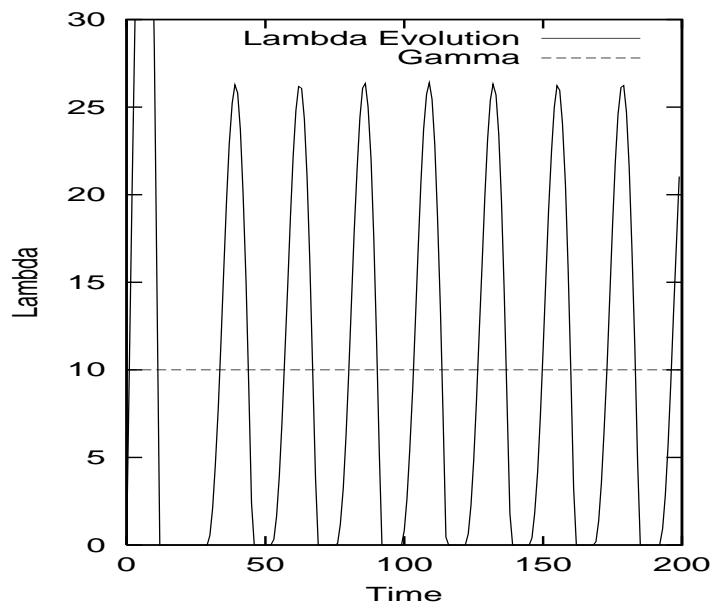
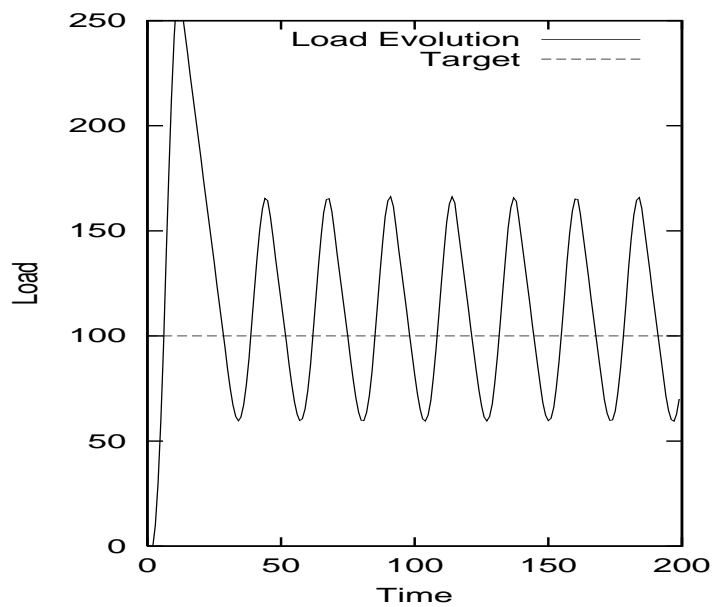
Thus:

- if $Q^* - Q(t) > 0$, increase $\lambda(t)$ proportional to gap
- if $Q^* - Q(t) < 0$, decrease $\lambda(t)$ proportional to gap

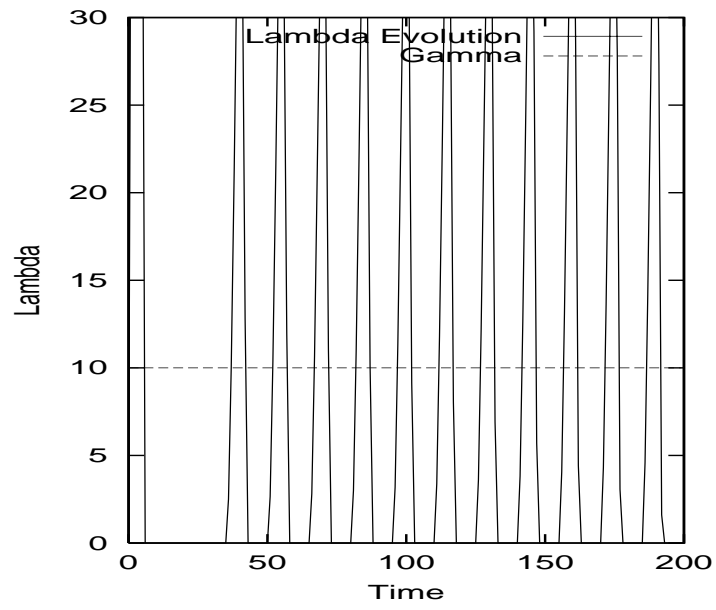
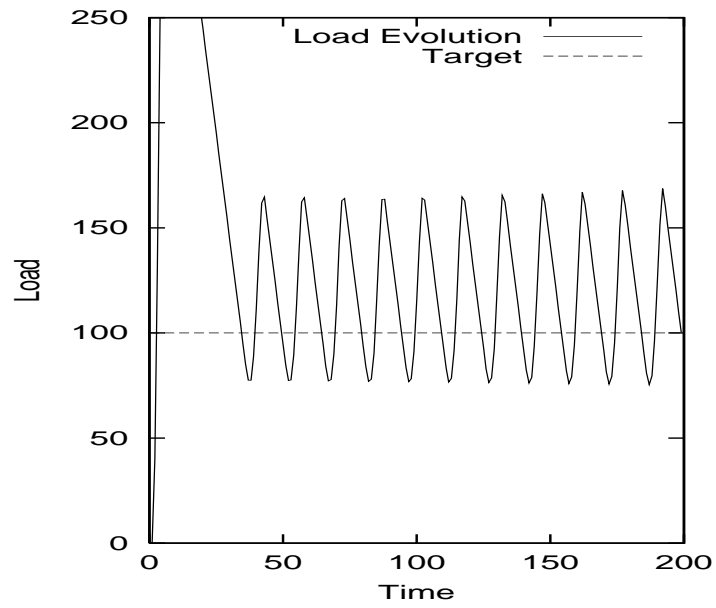
Trying to be more clever...

→ bottom line: is it any good?

With $\epsilon = 0.1$:



With $\epsilon = 0.5$:



Answer: no

→ looks good

→ but looks can be deceiving

Time to try something strange

→ any (crazy) ideas?

→ good for course project

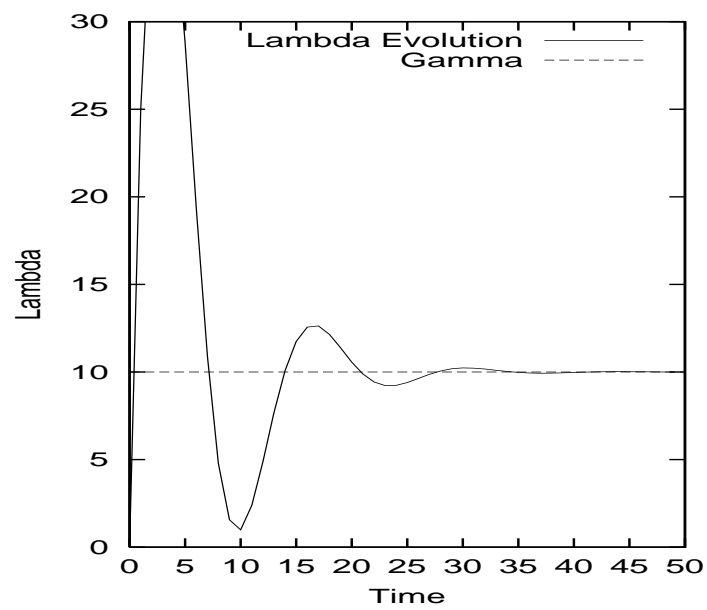
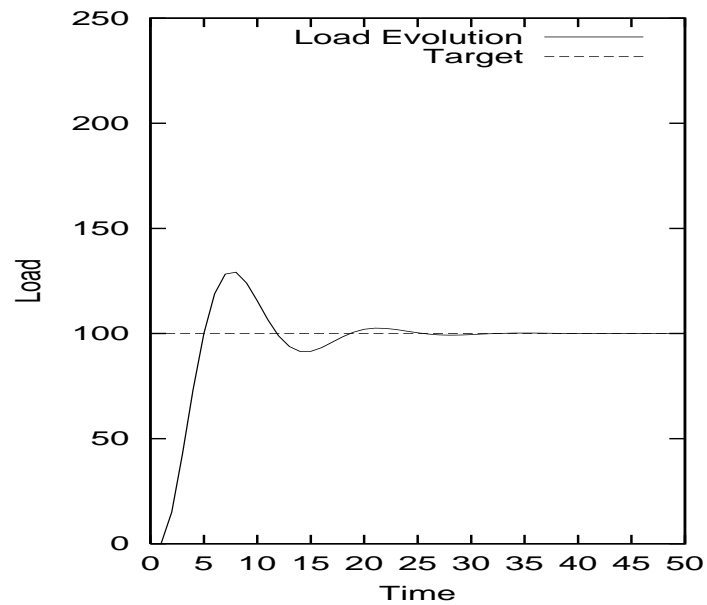
Method D:

$$\lambda(t + 1) \leftarrow \lambda(t) + \varepsilon(Q^* - Q(t)) - \beta(\lambda(t) - \gamma)$$

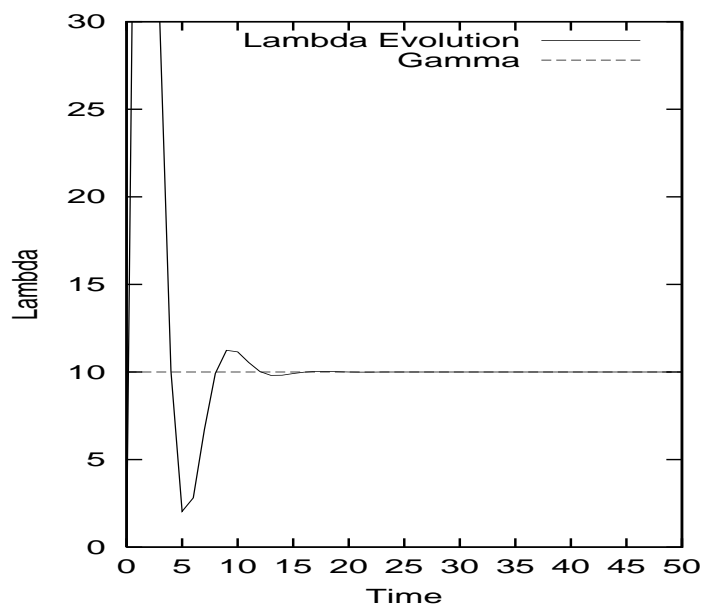
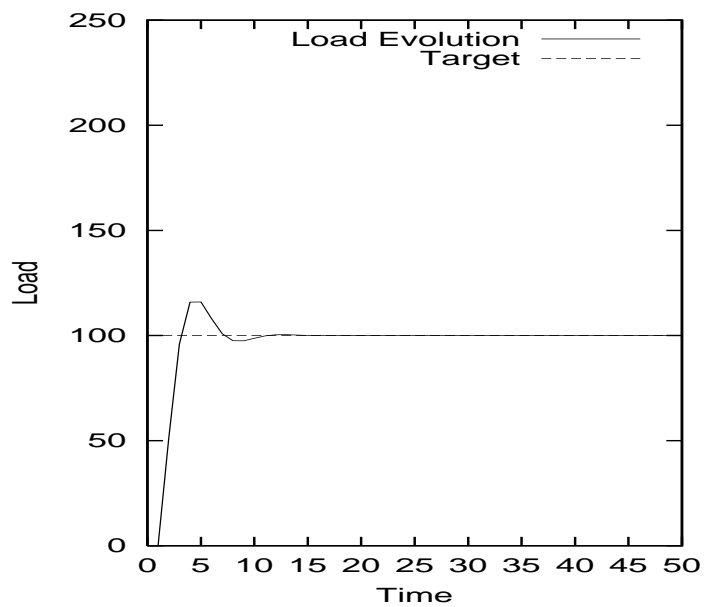
where $\varepsilon > 0$ and $\beta > 0$ are fixed parameters

- odd looking modification to **Method C**
- additional term $-\beta(\lambda(t) - \gamma)$
- what's going on?
- does it work?

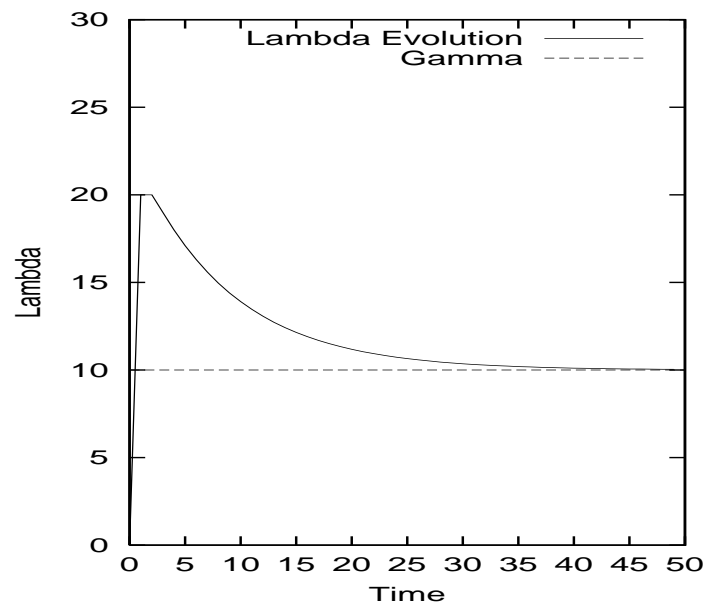
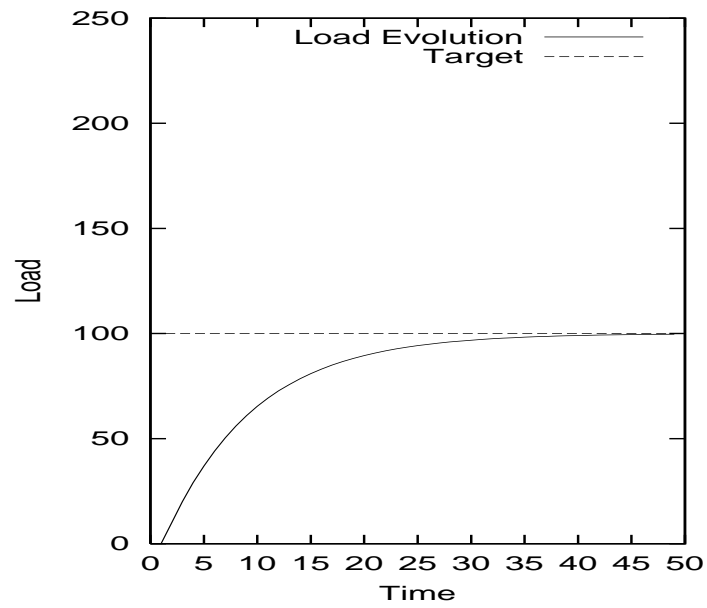
With $\varepsilon = 0.2$ and $\beta = 0.5$:



With $\varepsilon = 0.5$ and $\beta = 1.1$:



With $\varepsilon = 0.1$ and $\beta = 1.0$:



Remarks:

- Method D has desired behavior
- Is superior to Methods A, B, and C
- No unbounded oscillation
- In fact, dampening and convergence to desired state
 - converges to target operating point (Q^*, γ)
 - called asymptotically stable
 - why?