# Decoupling QoS Control from Core Routers:
# A Novel Bandwidth Broker Architecture for Scalable Support of Guaranteed Services

Zhi-Li Zhang, Zhenhai Duan, Lixin Gao, and YiWei Thomas Hou

# Virtual Time Reference System:
# A Unifying Scheduling Framework for Scalable Support of Guaranteed Services

Zhi-Li Zhang, Zhenhai Duan, and YiWei Thomas Hou

Speaker: Wei-Ming Yin

Instructor: Ying-Dar Lin

Nov. 24th, 2000

# Agenda

- Motivation
- Virtual Time Reference System
  - Core stateless framework
  - End-to-end delay bound
- Bandwidth Broker Architecture
- Admission Control for Per-Flow Guaranteed Services
  - All rate-based vs. Mixed rate- and delay-based schedulers
- Admission Control for Class-Based Guaranteed Services
  - Dynamic flow aggregation under all rate-based schedulers
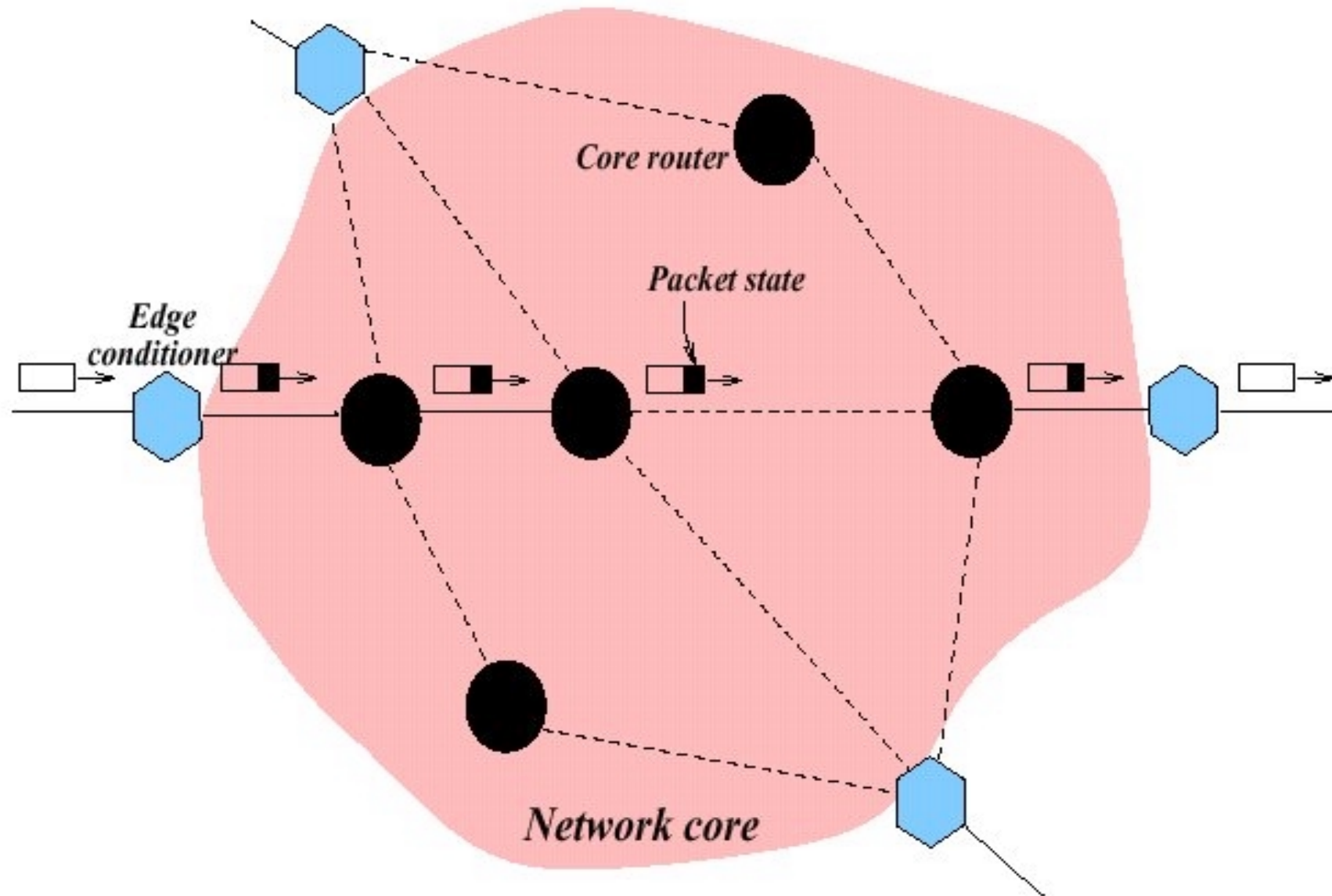- Simulation Investigation
- Conclusion and Future works

# Motivation

- Hop-by-hop admission control approach
  - Maintain per-flow or class-based QoS states at core routers
  - Perform local admission control and resource reservation
  - Maintain consistency of soft QoS states among all core routers
  - High communication overhead, less scalability, complicated design of core routers

- Path-oriented admission control approach
  - Relive core routers of QoS functions
  - Scale to both per-flow and class-based guaranteed services
  - Enable sophisticated QoS provisioning and admission control
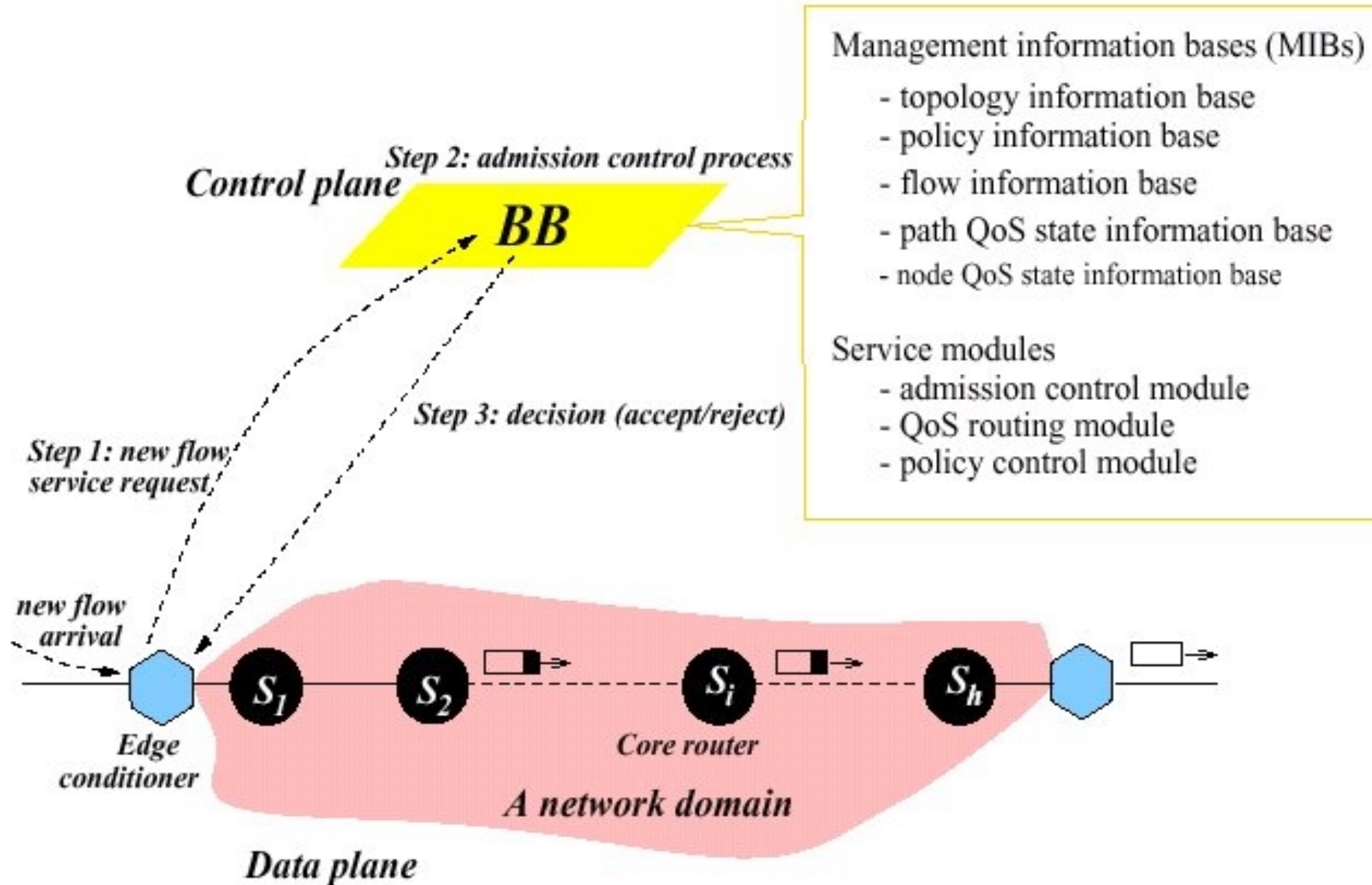  - No or minimal configuration of core routers

# Virtual Time Reference System

- A core stateless framework
- A unifying scheduling framework
  - Core routers only perform forwarding and scheduling
- Three logic components
  - Packet state (on packet)
  - Edge traffic conditioning (edge)
  - Virtual time reference/update mechanism (core)
- Characterize per-hop behavior and end-to-end delay bound

# Virtual Time Reference System

# System Overview



Management information bases (MIBs)

- topology information base
- policy information base
- flow information base
- path QoS state information base
- node QoS state information base

Service modules
- admission control module
- QoS routing module
- policy control module

**Control plane**

Step 2: admission control process

**BB**

Step 3: decision (accept/reject)

*Step 1: new flow, service request,*

*new flow arrival*

$S_1$   $S_2$   $S_i$   $S_h$

Edge conditioner

Core router

*A network domain*

**Data plane**

# Dynamic packet state

- State types:

> The $k$th packet of flow $j$ at core router $i$.
>
> The rate - delay parameter pair $\left(r^j, d^j\right)$. $\leftarrow$ admission control
>
> The virtual time stamp $w_i^{j,k}$. $\leftarrow$ edge
>
> The virtual time adjustment term $\delta^{j,k}$. $\leftarrow$ edge

- Carried in packet header, initialized and inserted at edge, referenced (scheduling module) and updated (forwarding module) at core.
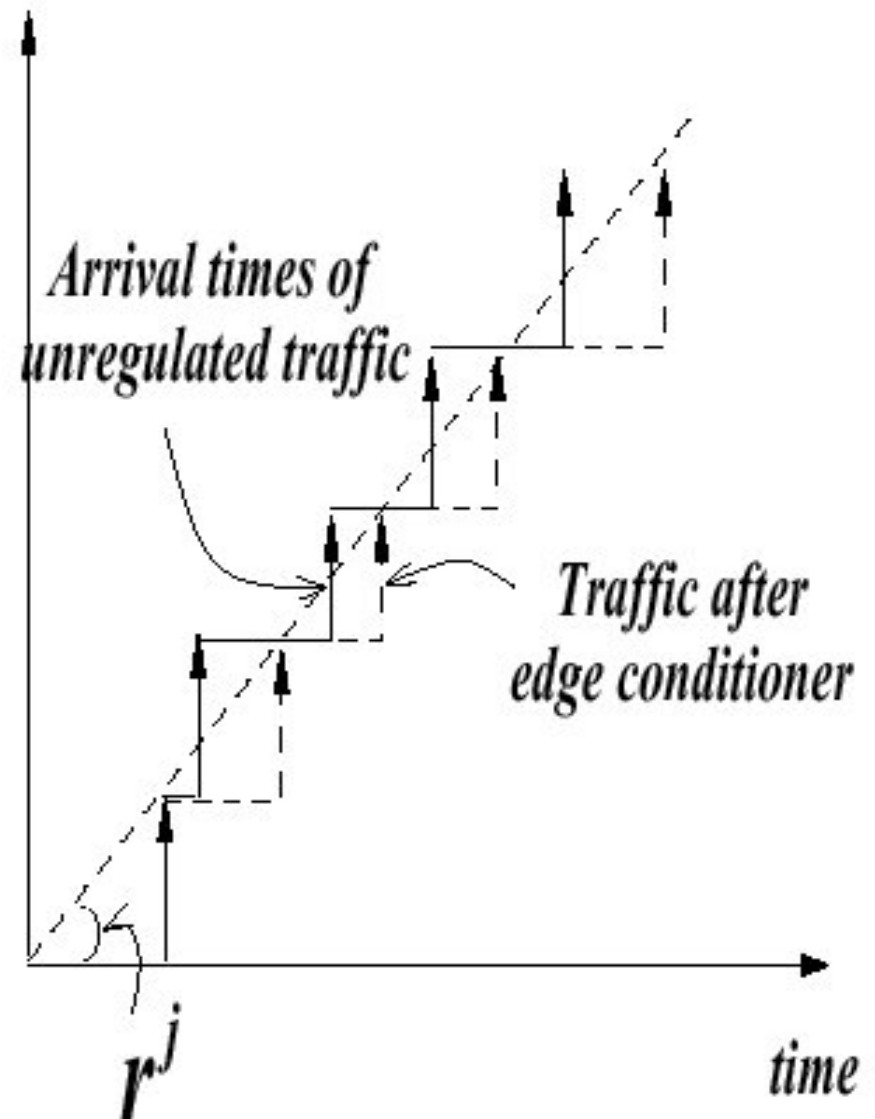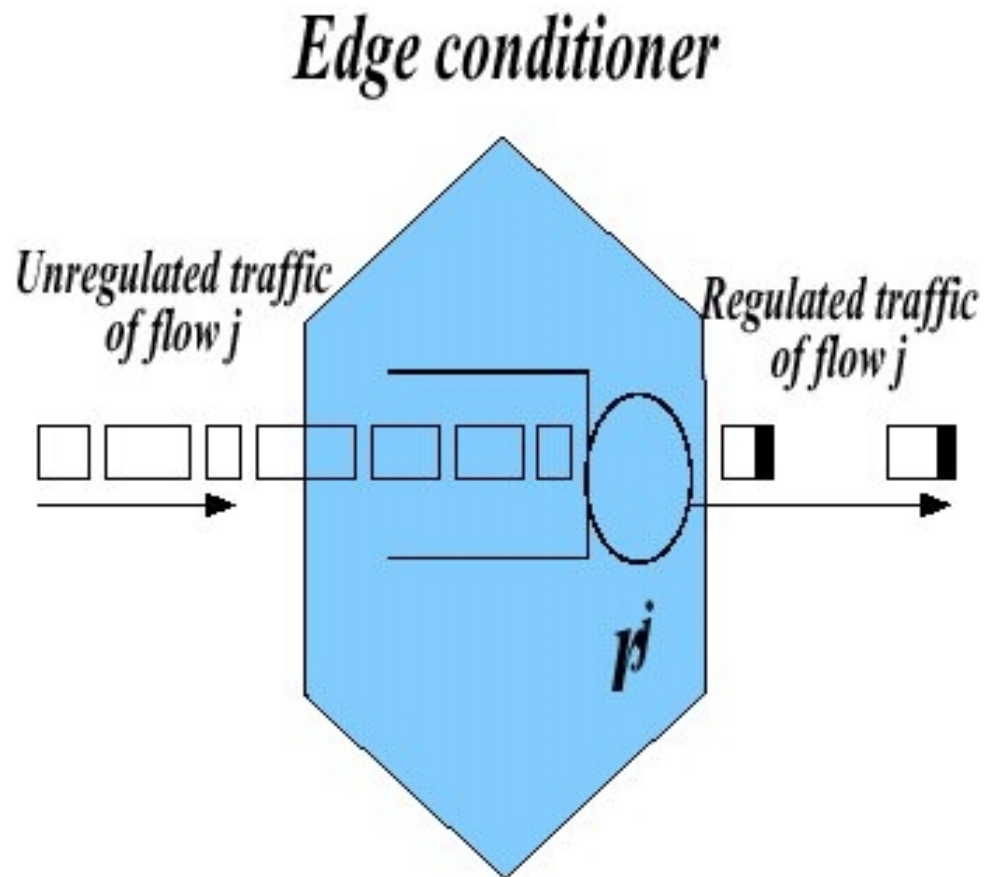
# Edge Traffic Conditioning

- Regulate packets injection rate not exceeding reserved rate

$$a_1^{j,k+1} - a_1^{j,k} \geq \frac{L^{j,k+1}}{r^j},$$

where $a_1^{j,k}$ denotes the arrival time of $k$th packet of flow $j$, and $L^{j,k}$ denotes the size of that packet.

# Edge Traffic Conditioning

# Virtual Time Reference/Update

- Per-hop behavior

$$\text{virtual delay}: d_i^{j,k} = \begin{cases} \frac{L^{j,k}}{r^j} + \delta^{j,k}, & \text{rate - based scheduler} \\ d^j, & \text{delay - based scheduler} \end{cases}$$
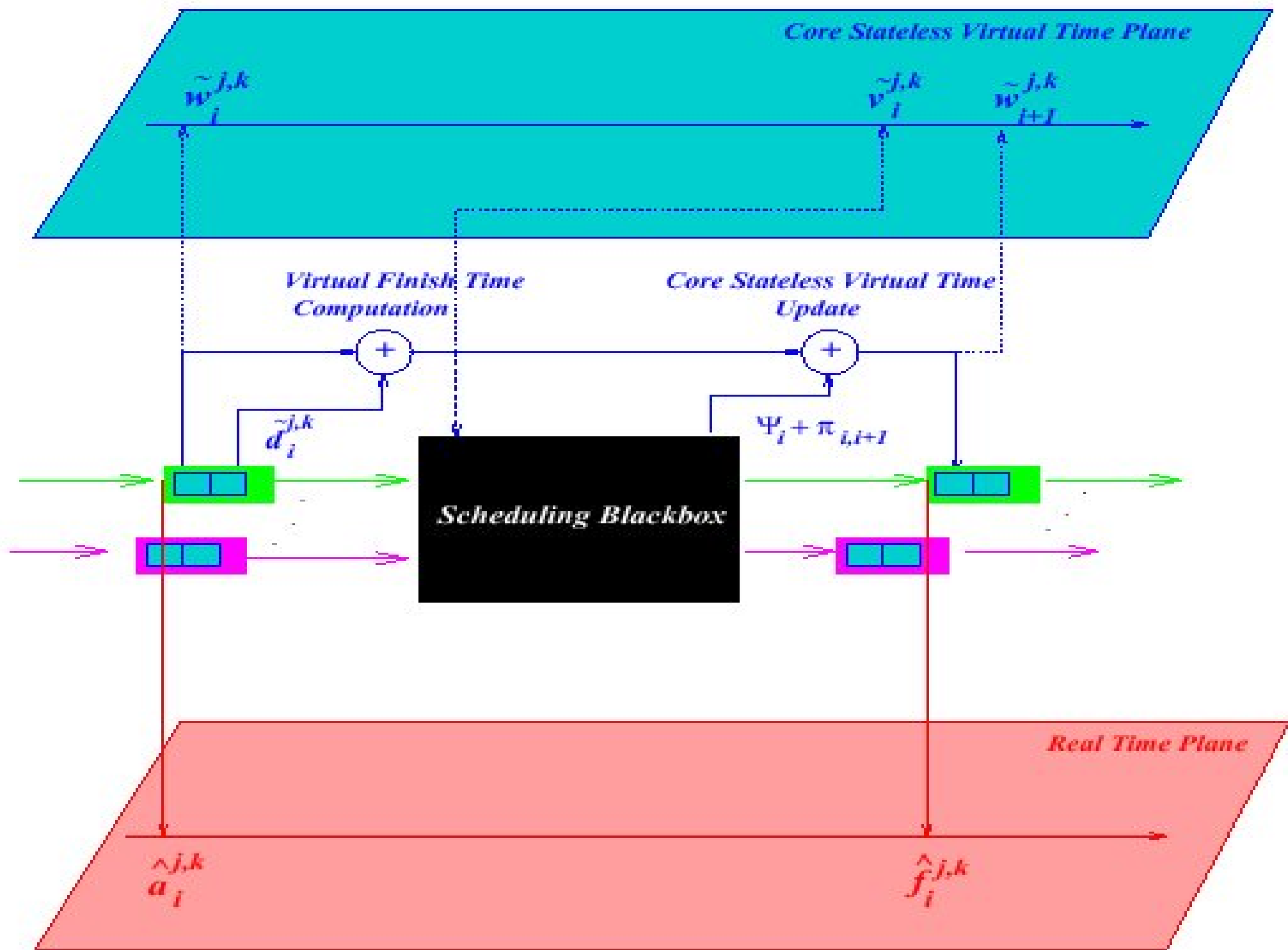
$$\text{virtual finish time}: v_i^{j,k} = w_i^{j,k} + d_i^{j,k} \quad \leftarrow \text{referenced}$$

$$\text{error term of core ronter } i : \Psi_i$$

$$\text{actual finish time}: f_i^{j,k} \leq v_i^{j,k} + \Psi_i \quad \leftarrow \text{per - hop behavior}$$

$$\text{propagation delay to next hop of core ruter } i : \pi_i$$

$$w_{i+1}^{j,k} = v_i^{j,k} + \Psi_i + \pi_i \quad \leftarrow \text{updated}$$

# Virtual Time Reference System

- Suppose total $h$ hops, of which $q$ hop are rate-based scheduler, and $h$-$q$ hops are delay-based schedulers. The traffic profile of flow $j$ is ( $s^j$, $r^j$, $P^j$, $L^{j,max}$ ).

$$f_h^{j,k} - a_1^{j,k} \leq d_{core}^j = q \frac{L^{j,\mathrm{max}}}{r^j} + (h - q)d^j + \sum_{i \in P} (\Psi_i + \pi_i)$$
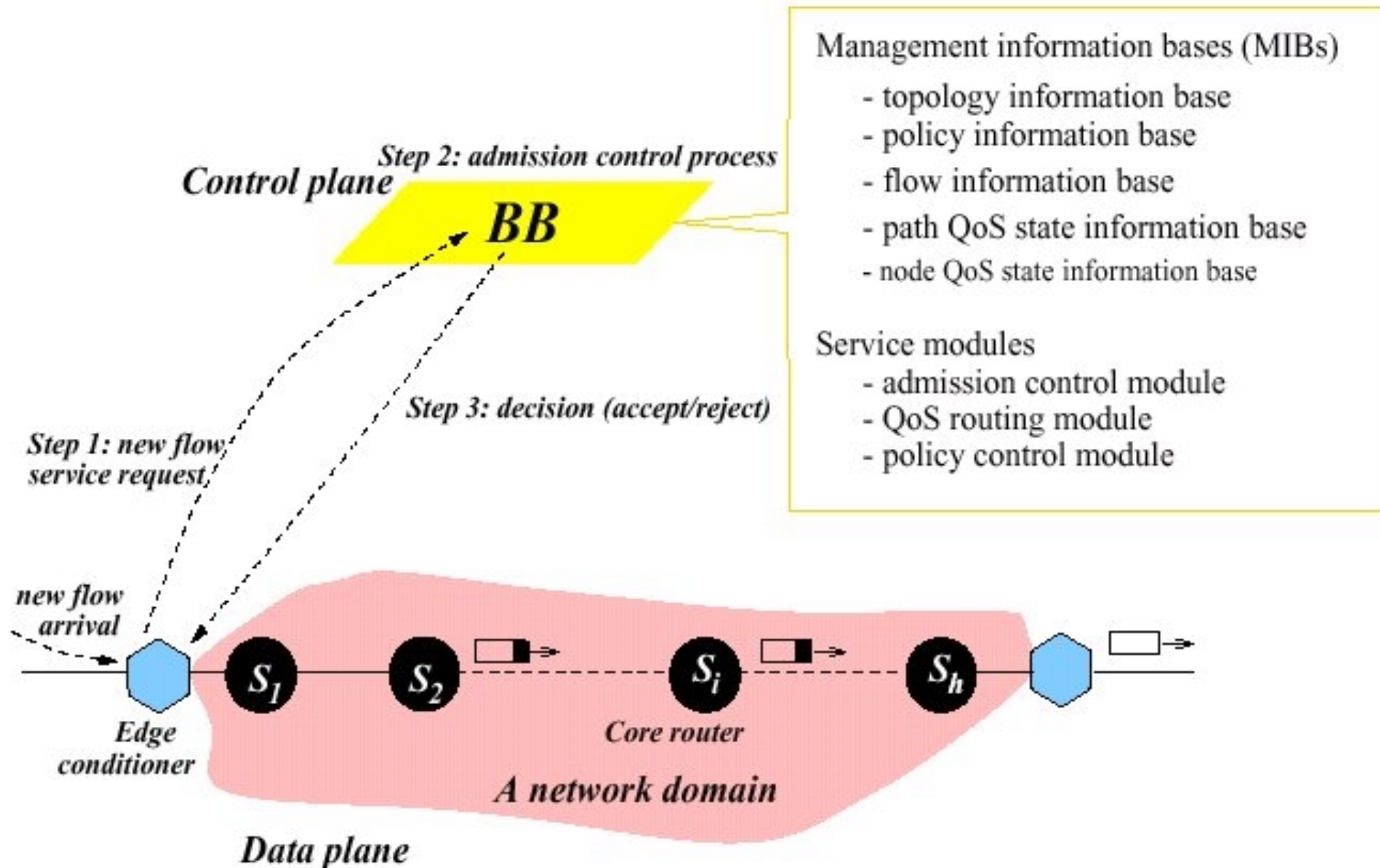
$$d_{edge}^j = \frac{\sigma^j - L^{j,\mathrm{max}}}{P^j - \rho^j} \frac{P^j - r^j}{r^j} + \frac{L^{j,\mathrm{max}}}{r^j} = T_{on}^j \frac{P^j - r^j}{r^j} + \frac{L^{j,\mathrm{max}}}{r^j}$$

end - to - end delay bound
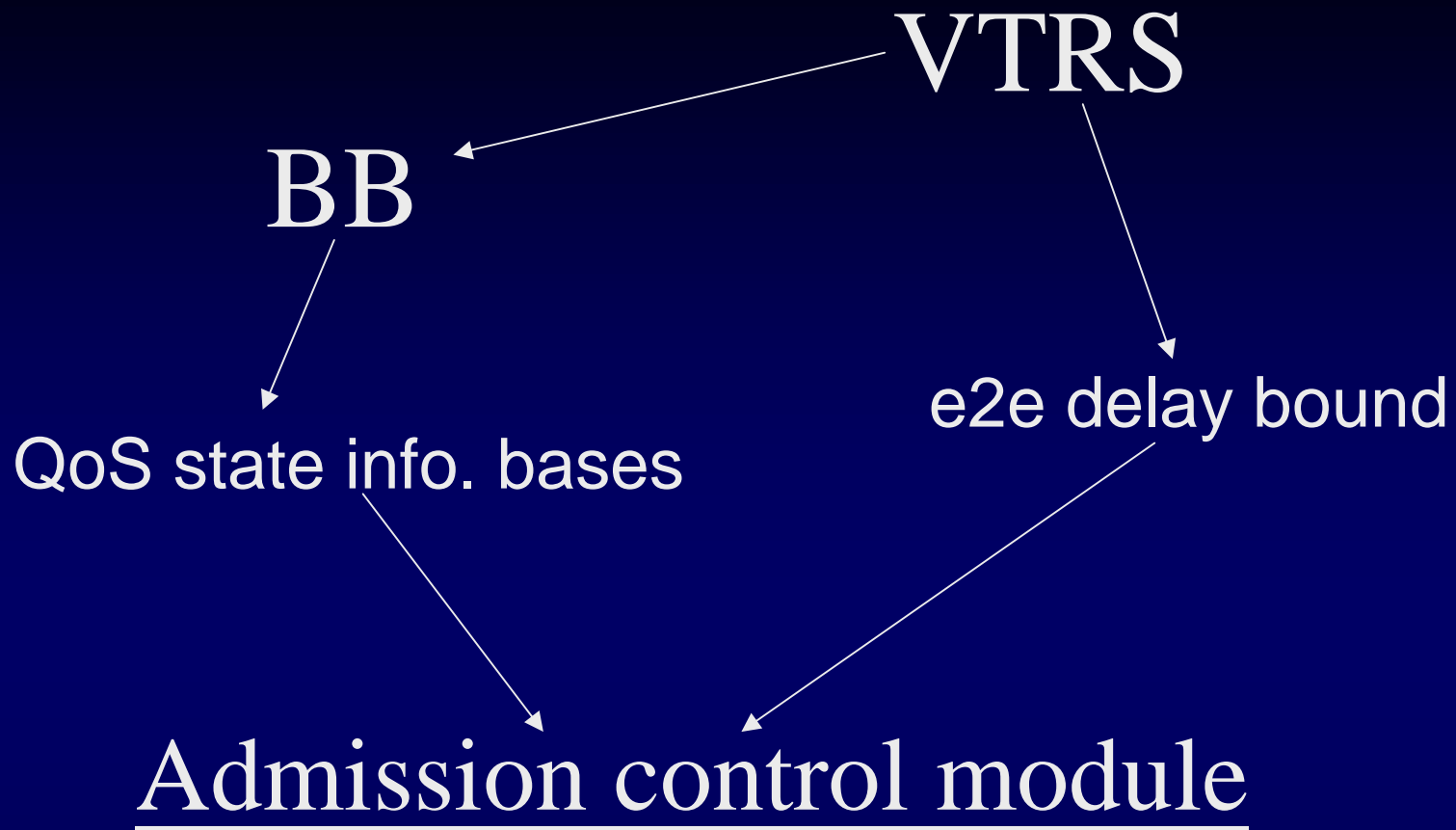
$$\Rightarrow d_{e2e}^j = d_{edge}^j + d_{core}^j$$

$$= T_{on}^j \frac{P^j - r^j}{r^j} + (q + 1)\frac{L^{j,\mathrm{max}}}{r^j} + (h - q)d^j + \sum_{i \in P} (\Psi_i + \pi_i)$$

# Bandwidth Broker Architecture

# QoS State Information Bases

- Flow information base
  - Flow id
  - Traffic profile: ( $\sigma^j$, $\rho^j$, $P^j$, $L^{j,max}$ )
  - Service profile: $D^{j,req}$
  - QoS reservation: ( $r,d$ )

- Node QoS state information base
  - Bandwidth, buffer capacity, scheduler type, error term

- Path QoS state information base
  - Hop number, sum of error terms and propagation delays, minimal residual bandwidth along the path

VTRS

BB

QoS state info. bases

e2e delay bound

Admission control module

Whether there is a feasible rate or not,
with which
delay requirement
is less than or equals to
e2e delay bound

# Admission Control

- For per-flow guaranteed services
  - Pure rate-based schedulers
  - Mixed rate- and delay-based schedulers
  - Scalability?
- Dynamic flow aggregation
- For class-based guaranteed services
  - Pure rate-based schedulers
  - Mixed rate- and delay-based schedulers

# Per-flow: Path with Only Rate-based Schedulers

- Parameters

$P : path$

$v : flow$

$r^v$ : reserved rate of flow $v$

$d^v$ : delay parameter of flow $v$

$S_i$ : core router $i$

$F_i$ : set of flows currently traversing $S_i$

$C_i$ : total bandwidth of $S_i$

$C_{res}^{S_i}$ : residual bandwidth at $S_i$

$C_{res}^{P}$ : minimal residual bandwidth, i.e. $C_{res}^{P} = \min_{i \in P} C_{res}^{S_i}$

$(\sigma^v, \rho^v, P^v, L^{v,\max})$ : traffic profile of flow $v$

$D^{v,req}$ : end-to-end delay requirement

# Per-flow: Path with Only Rate-based Schedulers

- Fundamental inequalities

$$\rho^v \le r^v \le P^v \quad \text{and} \quad r^v \le C_{res}^P$$

$$T_{on}^v \frac{P^v - r^v}{r^v} + (h+1)\frac{L^{v,\max}}{r^v} + \sum_{i \in P}(\Psi_i + \pi_i) \le D^{v,req}$$

- Feasible rate range derivation

Let $r_{\min}^v$ be the smallest $r^v$

$$\Rightarrow r_{\min}^v = \left[T_{on}^v P^v + (h+1)L^{v,\max}\right] / \left[D^{v,req} - \sum_{i \in P}(\Psi_i + \pi_i) + T_{on}^v\right]$$

Therefore, *feasible rate range*, $R_{fea}^*$, is defined as

$$\left[r_{fea}^{low}, r_{fea}^{up}\right] = \left[\max\{\rho^v, r_{\min}^v\}, \min\{P^v, C_{res}^P\}\right]$$

- The flow is admissible if the feasible rate range is non-empty, $d^v$ is not necessary to be determined.

# Per-flow: Path with Mixed Rate- and Delay-based Schedulers

- Fundamental inequalities

$$\rho^v \le r^v \le P^v,$$

$$r^v \le C_{res}^P,$$

$$T_{on}^v \frac{P^v - r^v}{r^v} + (q+1)\frac{L^{v,\max}}{r^v} + (h-q)d^v + \sum_{i \in P}\left(\Psi_i + \pi_i\right) \le D^{v,req},$$

and

$$\sum_{\{j \in F_i : d_i^j \le d_i^k\}}\left[r^j\left(d_i^k - d_i^j\right) + L^{j,\max}\right] + \left[r^v\left(d_i^k - d^v\right) + L^{v,\max}\right] \le C_i d_i^k.$$

# Per-flow: Path with Mixed Rate- and Delay-based Schedulers

- Efficient algorithm :

$$
\begin{aligned}
&0. \quad t^\nu = \frac{1}{h-q}[D^{\nu,req} - D^p_{tot} + T^\nu_{on}] \\
&1. \quad \text{Let } m^* \text{ such that } d^{m^*-1} < t^\nu \le d^{m^*} \\
&2. \quad \textbf{for } m = m^*, m^*-1, \dots, 2, 1 \\
&3. \qquad \mathcal{R}^m_{fea} \leftarrow [r^{m,l}_{fea}, r^{m,r}_{fea}] \\
&4. \qquad \mathcal{R}^m_{del} \leftarrow [r^{m,l}_{del}, r^{m,r}_{del}] \\
&5. \qquad \textbf{if } (\mathcal{R}^m_{fea} \cap \mathcal{R}^m_{del} == \emptyset) \\
&6. \qquad\quad \textbf{if } (\mathcal{R}^m_{fea} == \emptyset || \mathcal{R}^m_{del} == \emptyset || r^{m,r}_{fea} < r^{m,l}_{del}) \\
&7. \qquad\qquad \textbf{break with } d^\nu = d^m \\
&8. \qquad \textbf{else } /^*\mathcal{R}^m_{fea} \cap \mathcal{R}^m_{del} \neq \emptyset^*/ \\
&9. \qquad\quad \textbf{if } (r^{m,l}_{fea} < r^{m,l}_{del}) \\
&10. \qquad\qquad r^\nu \leftarrow r^{m,l}_{del}, d^\nu \leftarrow t^\nu - \frac{\Xi^\nu}{r^\nu} \\
&11. \qquad\qquad \textbf{break with } d^\nu \\
&12. \qquad \textbf{if } (d^\nu > t^\nu) \text{ no feasible value found} \\
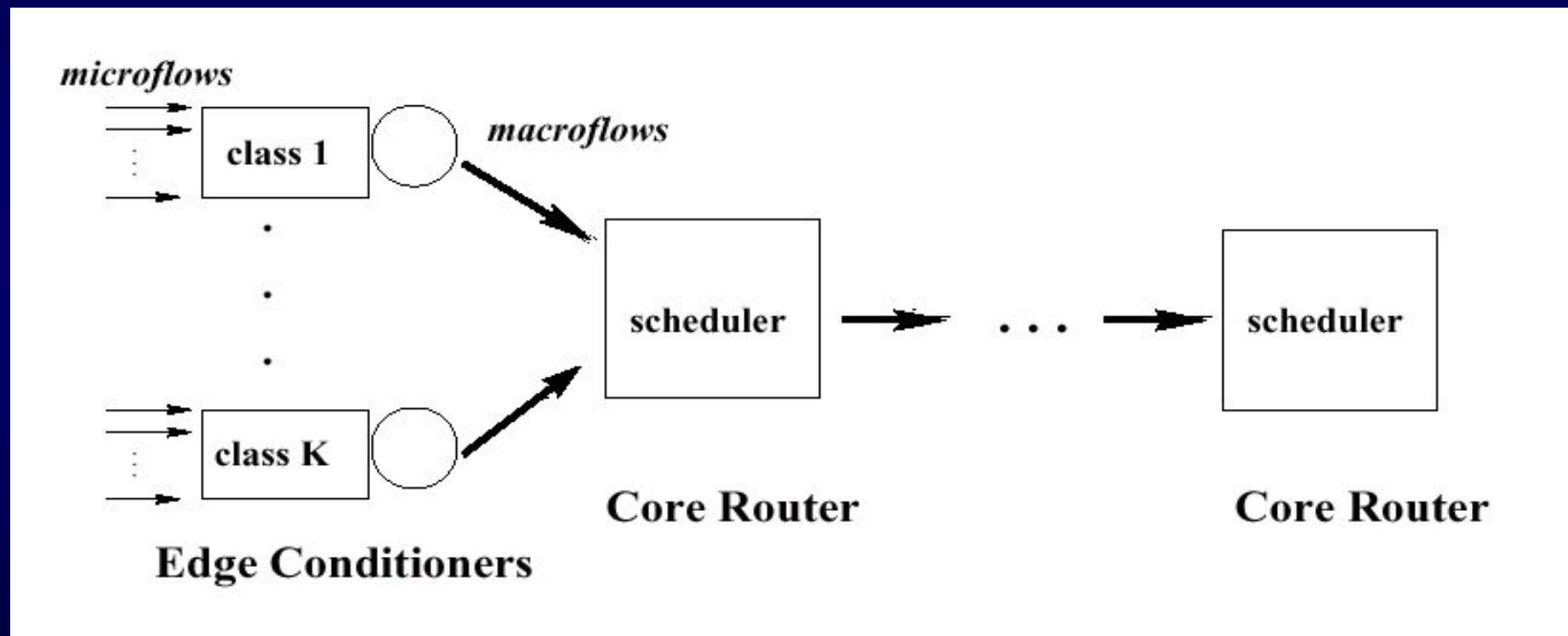&13. \qquad \textbf{else return } d^\nu
\end{aligned}
$$

- Time complexity

$$O(M), \text{ where } M \le \sum_{S_i \text{ id delay-based}} |F_i|.$$

# Class-based Guaranteed Service Model

- Enhance the scalability of proposed BB architecture
- Service Model



- Dynamic flow aggregation has not been identified nor addressed

- Impact on e2e delay ( macroflow α → α' )
  - All rate-based schedulers

worse - case delay at edge conditioner is larger than

$$d_{edge}^{\alpha'} = T_{on}^{\alpha'} \frac{P^{\alpha'} - r^{\alpha'}}{r^{\alpha'}} + \frac{L^{\alpha',\text{max}}}{r^{\alpha'}}$$
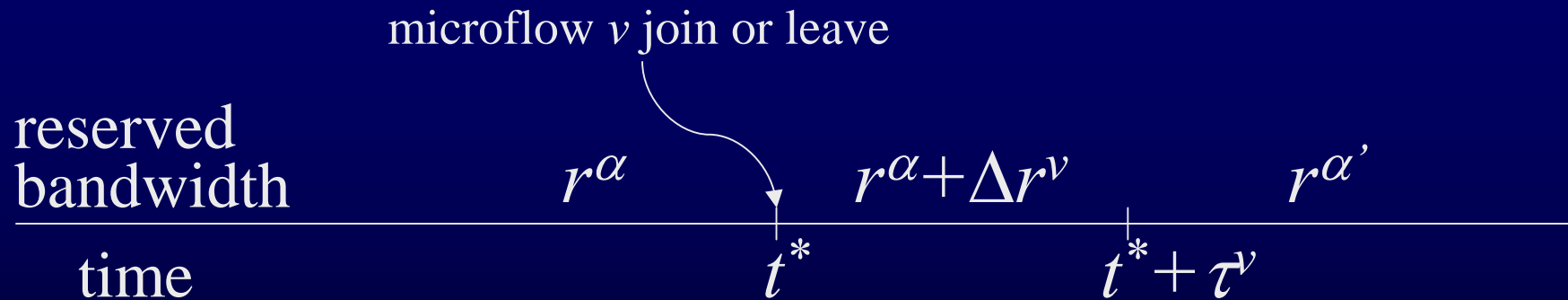
some packets from new macroflow may experience a worst - case delay in the network core

by $d_{core}^{\alpha} = h \dfrac{L^{P,\text{max}}}{r^{\alpha}} + \sum_{i \in P} \left( \Psi_i + \pi_i \right)$

instead of $d_{core}^{\alpha'} = h \dfrac{L^{P,\text{max}}}{r^{\alpha'}} + \sum_{i \in P} \left( \Psi_i + \pi_i \right)$

# Dynamic Flow Aggregation (2/5)

- ## Edge delay bound
  - Contingency bandwidth : to eliminate the lingering delay effect of the backlog packets
  - A new microflow $v$ aggregates or de-aggregates, the contingency bandwidth is $\Delta r^v$, and the contingency period is $\tau^v$.

microflow $v$ join or leave

reserved bandwidth

$r^\alpha$      $r^\alpha + \Delta r^v$      $r^{\alpha'}$

time     $t^*$         $t^* + \tau^v$

  - $\Delta r^v$ and $\tau^v$ are chosen to bound the edge delay as

$$d_{edge}^{new} \leq \max\left\{d_{edge}^{\alpha}, d_{edge}^{\alpha'}\right\}$$

# Dynamic Flow Aggregation (3/5)

- Edge delay bound
  - The microflow $v$ is with ( $\sigma^v$, $\rho^v$, $P^v$, $L^{v,max}$ ).
  - Sufficient conditions on $\Delta r^v$ and $\tau^v$ :

$$\begin{cases} \Delta r^v \geq P^v - r^v & (\text{ microflow join }) \\ \Delta r^v \geq r^v & (\text{microflow leave}) \end{cases}$$

and $\tau^v \geq \dfrac{Q(t^*)}{\Delta r^v}$,

where $Q(t^*) \leq d^{\alpha}_{edge}\left(r^{\alpha} + \Delta r^{\alpha}(t^*)\right)$ is the backlog,

where $\Delta r^{\alpha}(t^*)$ is the total contingency bandwidth allocated

to the macroflow $\alpha$ at time $t^*$.

$\rightarrow$ *contingency period bounding*

- Core delay bound

$$d_{core}^{\alpha'} = q \max\left\{\frac{L^{P,\max}}{r^{\alpha}}, \frac{L^{P,\max}}{r^{\alpha'}}\right\} + (h - q)d^{\alpha} + \sum_{i \in P}\left(\Psi_i + \pi_i\right)$$

# Dynamic Flow Aggregation (5/5)

- Admission control: Microflow join

$$d_{e2e}^{\alpha'} = d_{edge}^{\alpha'} + \max\left\{ d_{core}^{\alpha}, d_{core}^{\alpha'} \right\} \le D^{\alpha,req}$$
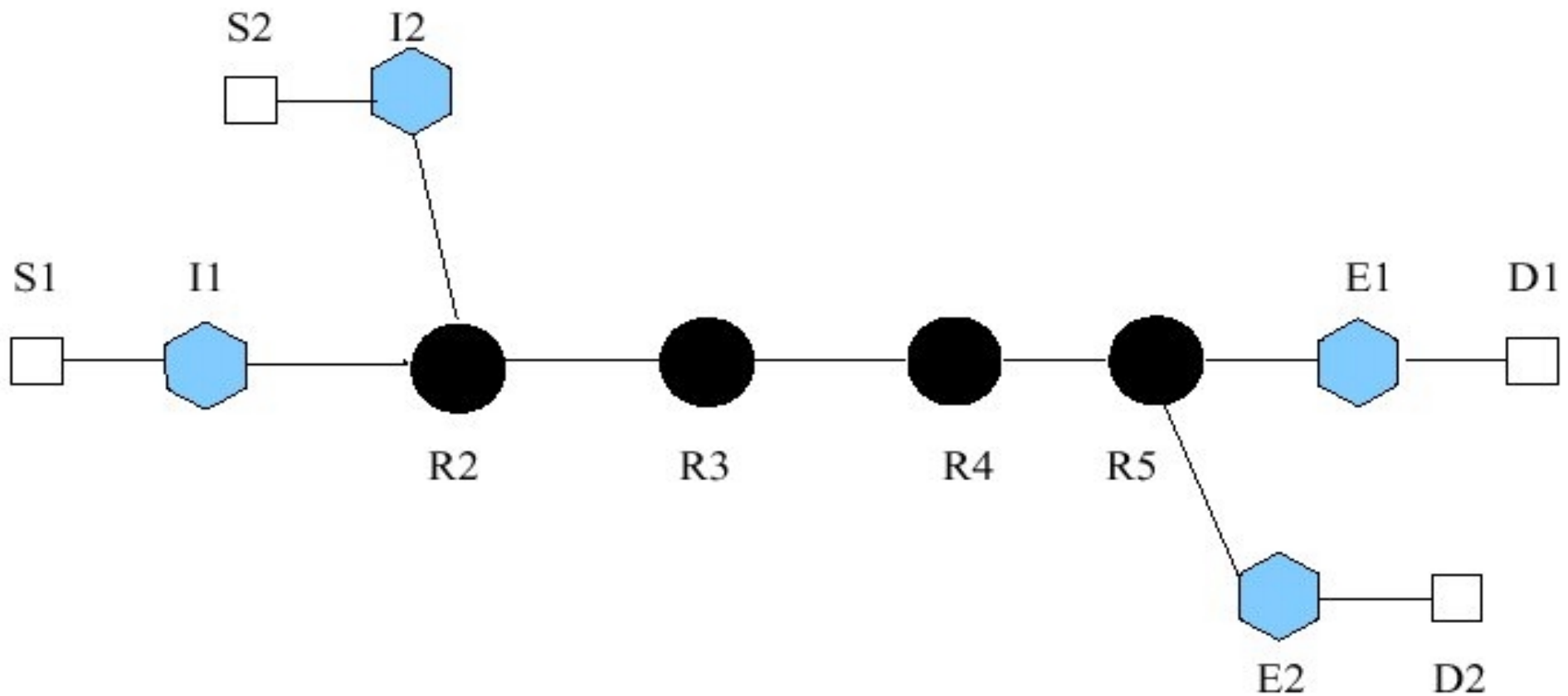
$$\text{since } r^{\alpha'} \ge r^{\alpha}, \text{ hence, } d_{core}^{\alpha} \le d_{core}^{\alpha'}.$$

$$\Rightarrow d_{edge}^{\alpha'} \le D^{\alpha,req} - d_{core}^{\alpha} \ldots\ldots(a)$$

$$\text{also, } \rho^{v} \le r^{\alpha'} - r^{\alpha} \le P^{v} \ldots\ldots(b)$$

$$\Rightarrow r^{\alpha'} \text{ can be derived from (a) (b)}$$
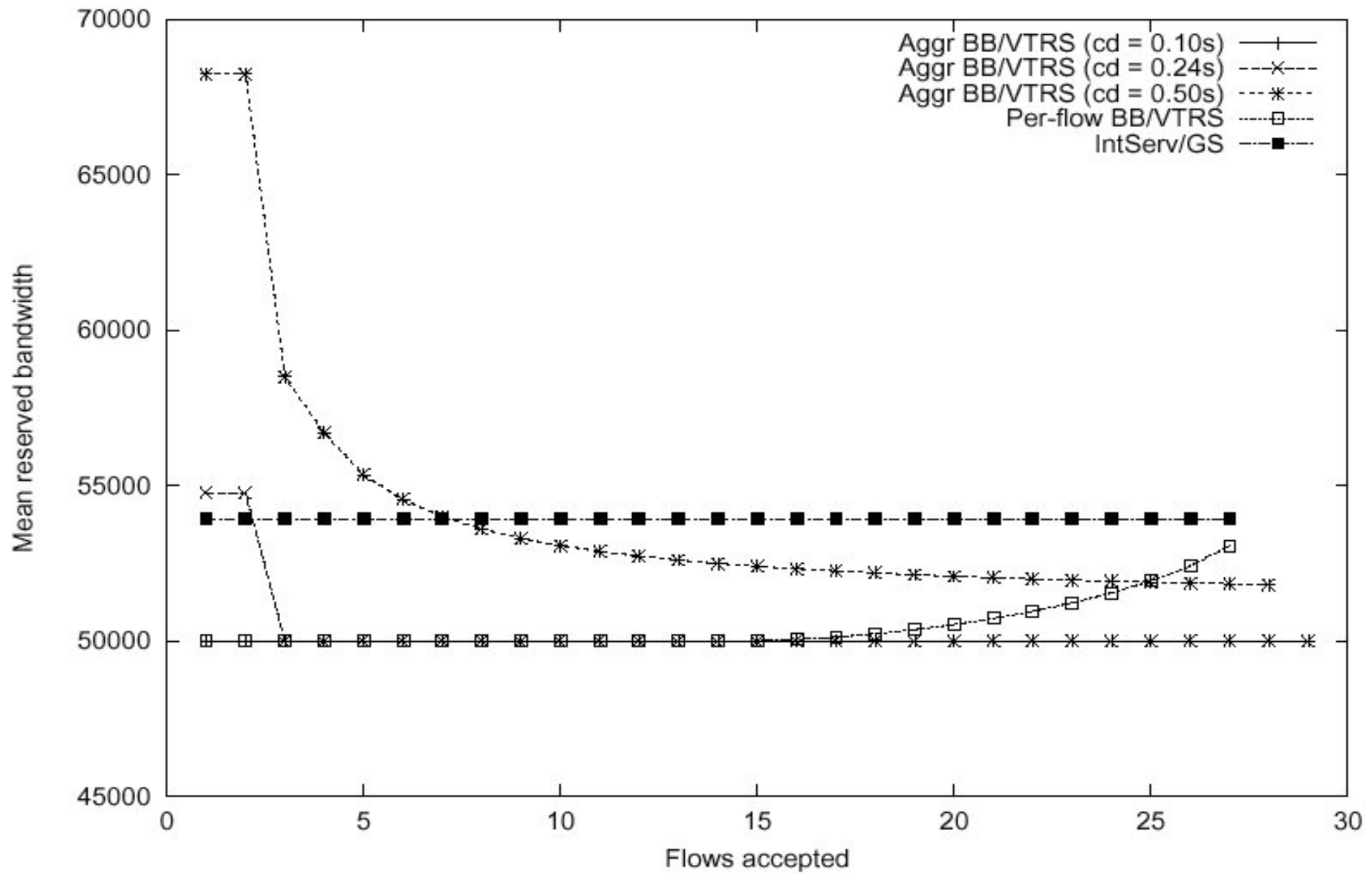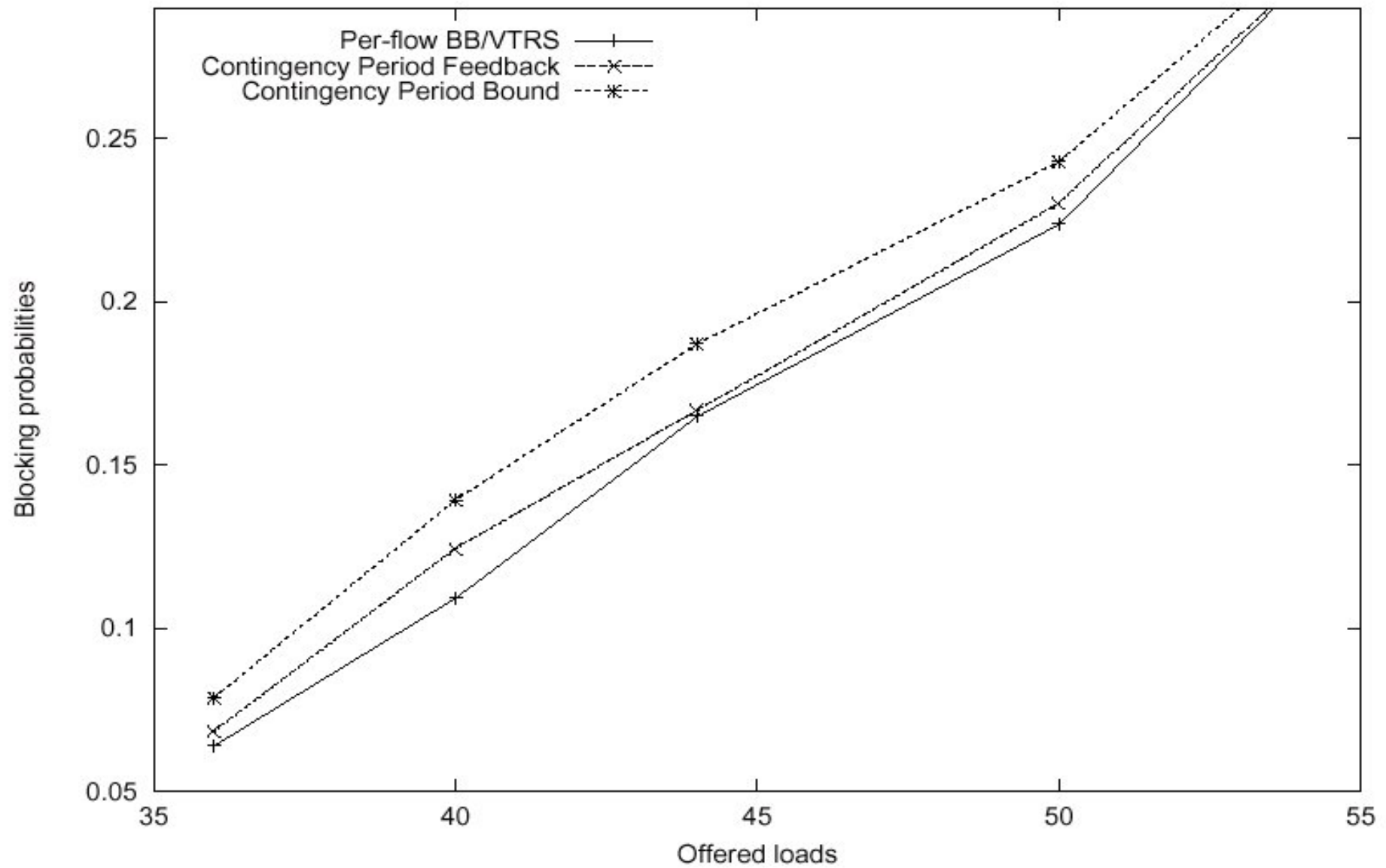
# Simulation Investigation

# Comparison

Table 2: Comparison of IntServ/GS, per-flow BB/VTRS and aggregate BB/VTRS schemes.

| | | Number of Calls admitted | | | |
|---|---|---|---|---|---|
| | | Rate-Based Only | | Mixed Rate/Delay-Based | |
| Delay bounds | | 2.44 | 2.19 | 2.44 | 2.19 |
| IntServ/GS | | 30 | 27 | 30 | 27 |
| Per-flow BB/VTRS | | 30 | 27 | 30 | 27 |
| Aggr BB/VTRS | cd = 0.10 | 29 | 29 | 29 | 29 |
| | cd = 0.24 | | | 29 | 29 |
| | cd = 0.50 | | | 29 | 28 |

# Mean Reserved Bandwidth

# Flow Blocking Rate

# Conclusion

- Present a novel BB architecture based on VTRS
- Decouple the QoS control plane from data plane
- Propose path-oriented admission control approach
- Support per-flow and class-based guaranteed services
- No or minimal configuration of core routers

# Future Works

- Distributed bandwidth broker architecture
- Inter-Domain QoS reservation and service level agreement