# Admission Control and Routing to Parallel Queues with Delayed Information via Marginal Productivity Indices

Peter Jacko* and José Niño-Mora
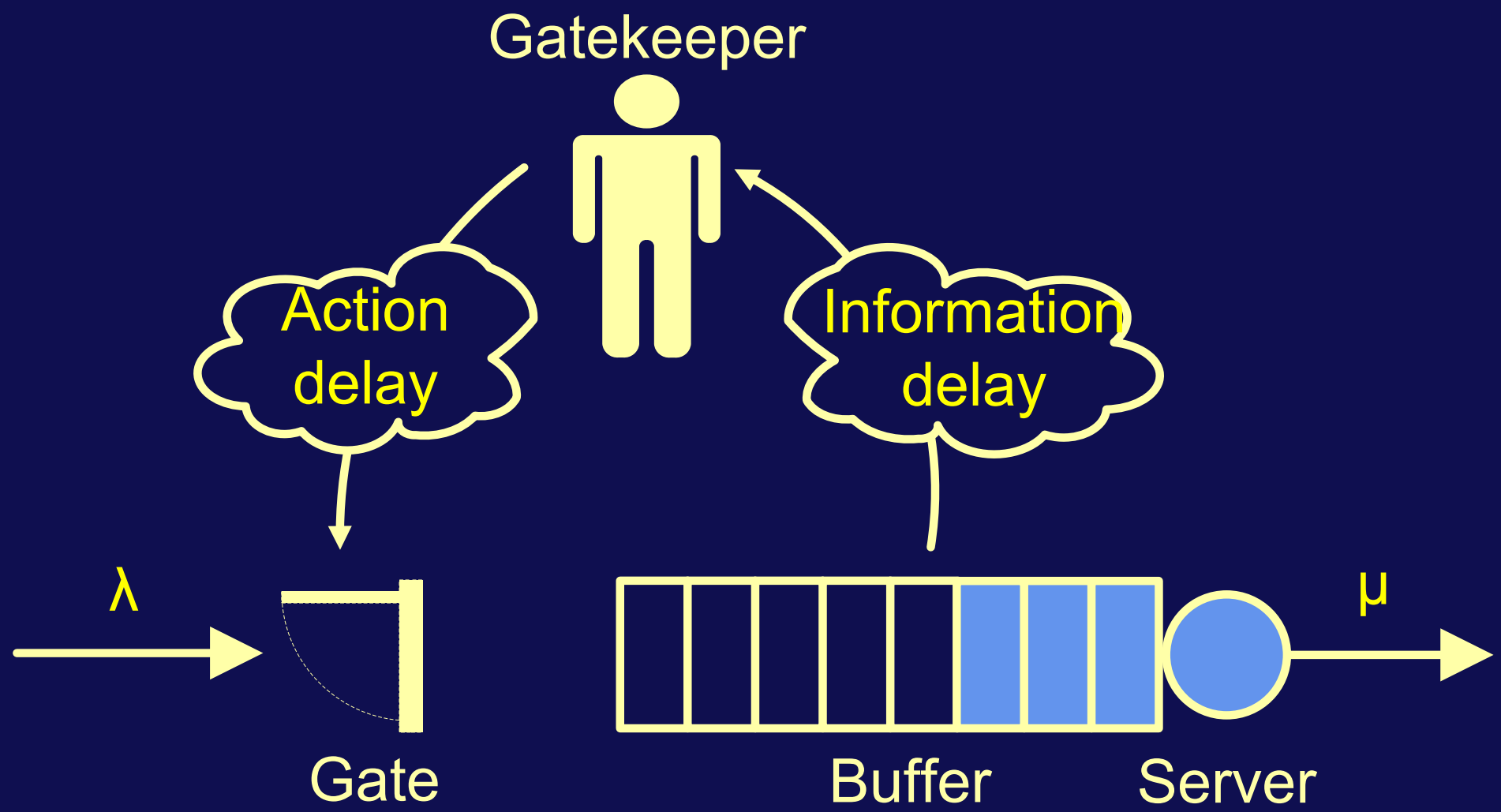
ValueTools 2008, Athens

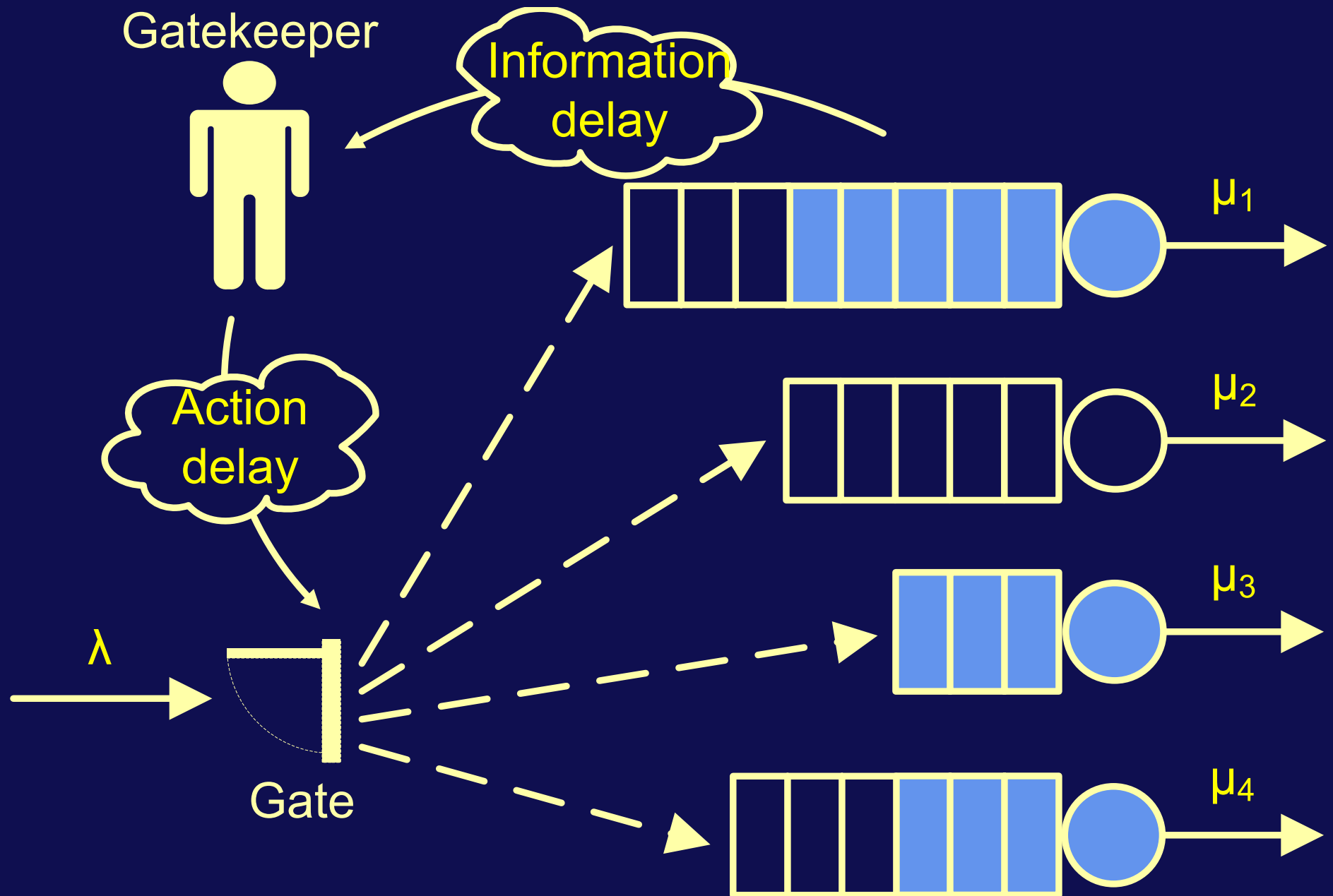*Universidad Carlos III de Madrid, Department of Statistics

# Motivation

- Delays in information flow and action implementation

  ▷ physical distance of nodes in networks
  ▷ long-distance-controlled robots
  ▷ advanced processing of observations

- May lead to important losses if ignored

- We deal with delays in:

  ▷ admission control and routing to parallel queues
  ▷ admission control to a single queue

# Admission Control to a Single Queue with Delays

# Admission Control and Routing with Delays

# Admission Control and Routing with Delays

- Even with $2$ queues it is <span style="color:yellow">hard to analyze</span>

- Delay of one period and symmetric queues: <span style="color:yellow">JSEQ</span>

  ▷ "A large number of properties needs to be discovered and then tediously verified." (Kuri & Kumar, 1995)

- Delay of more than one period:

  ▷ ". . . the approach quickly becomes very unwieldy. . . "
  ▷ JSEQ is <span style="color:yellow">not</span> optimal (both Kuri & Kumar, 1995)
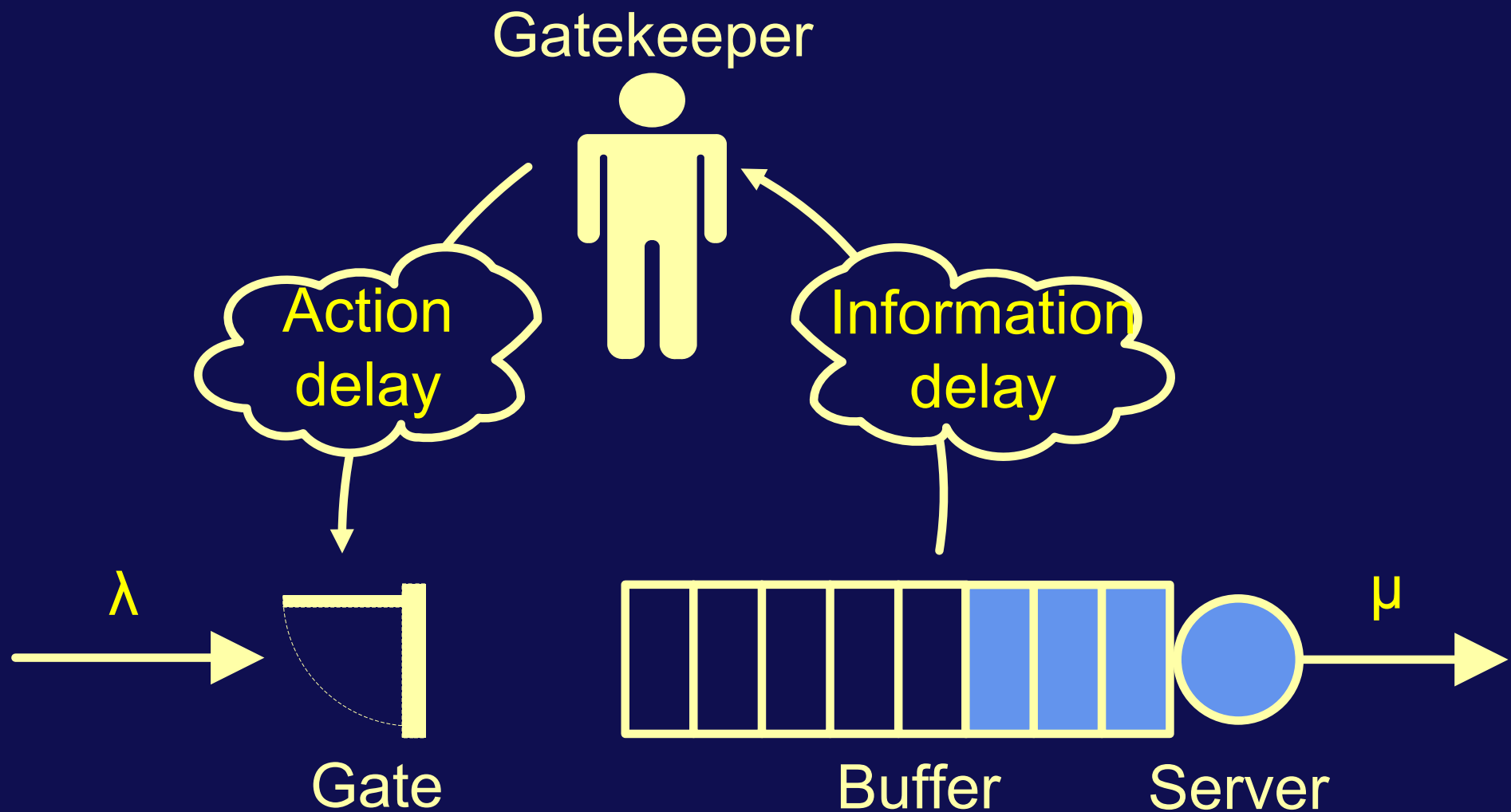  ▷ "We were not able to derive significant results. . . " (Artiges, 1993)

# Admission Control and Routing with Delays

- What if servers, buffers, holding costs and delays differ?

  ▷ curse of dimensionality

- It is a joint decision of admitting to at most 1 queue

  ▷ there must be a queue where a job is worth admitting
  ▷ if there are several such queues, route to the queue where admitting is most profitable

- Accomplished by an index policy

  ▷ via marginal productivity index (MPI)
  ▷ may be suboptimal due to ignored cross-dependence

# Outline

- Admission control to a single queue:

  ▷ MDP model with no delay
  ▷ MDP model with one period delay
  ▷ exploiting special structure via bi-threshold policies
  ▷ establishing existence of MPIs
  ▷ obtaining a fast algorithm for MPI calculation

- MPI policy properties for

  ▷ admission control and routing with one period delay
  ▷ servers assignment problem with one period delay

- Discussion of generalizations

# Admission Control to a Single Queue with Delays

# Admission Control with No Delay

- Discrete time epochs $t = 0, 1, 2, \ldots$

- Bernoulli arrivals at rate $\lambda$ per period

- Geometric server at rate $\mu$ per period

- Buffer $+$ server room: $I$

- Holding costs at rate $C_i$ per period with $i$ jobs

  $\triangleright$ convex, nondecreasing in $i$

- Loss costs at rate $\nu$ per rejected job

# MDP Model (No Delay)

- Action process $a(t) \in \mathcal{A} := \{0, 1\}$: closing the gate $(a(t) = 1)$ or opening the gate $(a(t) = 0)$

- State process $X(t) \in \mathcal{I} := \{0, 1, \ldots, I\}$

  $\triangleright$ state $I$ is uncontrollable

- At epoch $t$: $a(t)$ must be based on $X(t)$

- Transition probabilities $p_{ij}^a$

- One-period cost $C_i + \nu W_i^a$, where the work $W_i^a$ is

$$W_i^1 := \lambda \qquad\qquad W_i^0 := \begin{cases} \lambda & \text{if } i = I \\ 0 & \text{otherwise} \end{cases}$$

# MDP Model (One-Period Delay)

- $a(t)$ must be based on $\widetilde{X}(t) := (a(t-1), X(t-1))$

  $\triangleright$ because $X(t)$ is not known at $t$

- Action space $\mathcal{A}$ as before

- Augmented states $\widetilde{\mathcal{I}} := (\mathcal{A} \times \{0, 1, \ldots, I-1\}) \cup \{(*, I)\}$

  $\triangleright$ state $(*, I)$ appears by merging $(0, I)$ with $(1, I)$

- Transition probabilities $p^{a'}_{(a,i),(b,j)} := p^a_{ij} \cdot \mathbf{1}\{a' = b\}$

- One-period cost $C_{(a,i)} + \nu W_{(a,i)} := C_i + \nu W^a_i$

  $\triangleright$ note the independence of the current-period action

# Objective

- Solving the $\nu$-wage problem: $\displaystyle\min_{\pi\in\Pi} f^{\pi}_{(a,i)} + \nu g^{\pi}_{(a,i)}$

  $\triangleright$ choosing a non-anticipative control policy $\pi \in \Pi$
  $\triangleright$ expected total discounted holding cost

$$f^{\pi}_{(a,i)} := \mathbb{E}^{\pi}_{(a,i)} \left[ \sum_{t=0}^{\infty} \beta^t C_{\widetilde{X}(t)} \right]$$

  $\triangleright$ expected total discounted work (number of rejections)

$$g^{\pi}_{(a,i)} := \mathbb{E}^{\pi}_{(a,i)} \left[ \sum_{t=0}^{\infty} \beta^t W_{\widetilde{X}(t)} \right]$$

# Exploiting Special Structure

- There is an optimal policy which is stationary, deterministic, independent of the initial state

- Represent such policies as active sets $\mathcal{S} \subseteq \widetilde{\mathcal{I}}$
  ▷ the set of states in which it prescribes to shut the gate

- Bi-threshold policies are optimal (Altman & Nain, 1992)
  ▷ $\widetilde{\mathcal{I}}_{K,K}$

  ▷ $\widetilde{\mathcal{I}}_{K,K+1}$



- The family of all such active sets: $\mathcal{F}$

# Reduced Problem

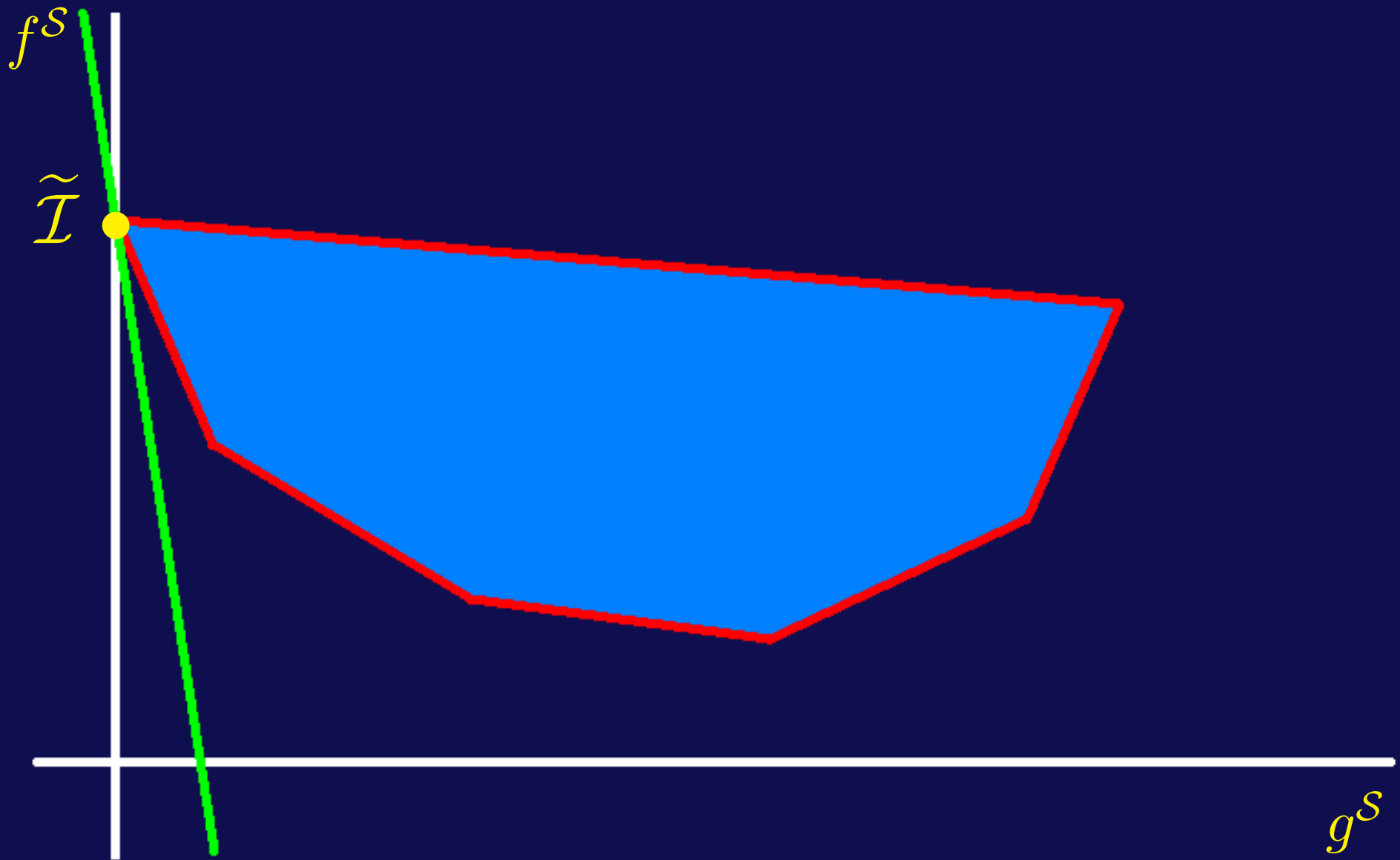- The $\nu$-wage problem can be solved by solving

$$\min_{\mathcal{S} \in \mathcal{F}} f^{\mathcal{S}}_{(a,i)} + \nu g^{\mathcal{S}}_{(a,i)}$$

- Evaluating all $\mathcal{S} \in \mathcal{F}$ requires $\mathcal{O}(I^4)$ operations

- "Dual" approach in $\mathcal{O}(I^3)$: marginal productivity indices

  ▷ but indexability (MPIs existence) must be proved
  ▷ we do by verifying PCL-indexability (Niño-Mora, 2001)
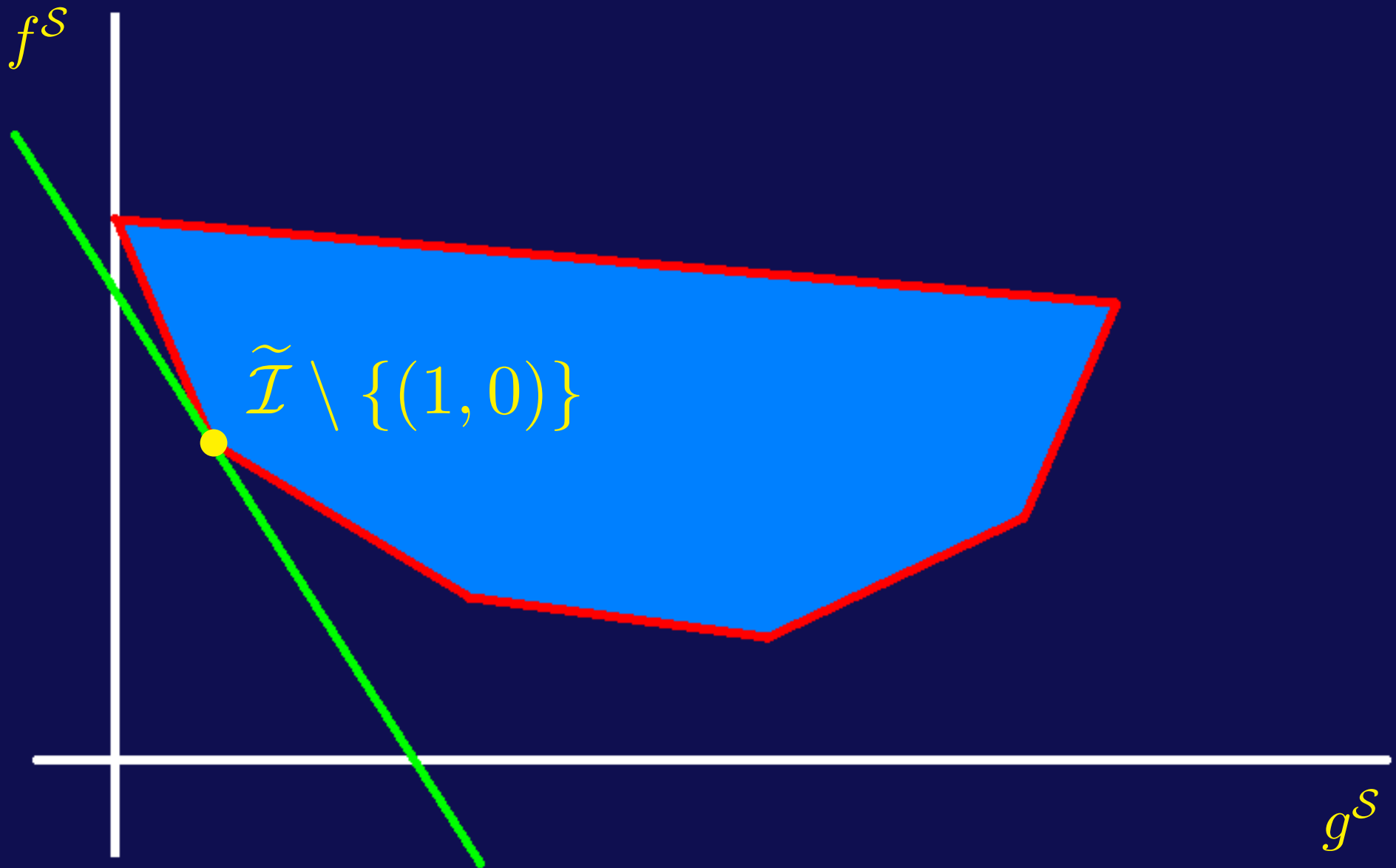  ▷ we improve algorithm to $\mathcal{O}(I)$, as in no-delay case

# Indexability

- $\nu$-wage problem is indexable, if

  ▷ the optimal active set decreases monotonically
  from $\widetilde{\mathcal{I}}$ to $\emptyset$ as $\nu$ increases from $-\infty$ to $\infty$

- Equivalently, there exist values $\nu_{(a,i)}$ such that

  ▷ it is optimal to shut the gate at state $(a, i)$ if $\nu_{(a,i)} \geq \nu$
  ▷ it is optimal to open the gate at state $(a, i)$ if $\nu_{(a,i)} \leq \nu$

- $\nu_{(a,i)}$ is the marginal productivity index

  ▷ capturing the marginal productivity of work
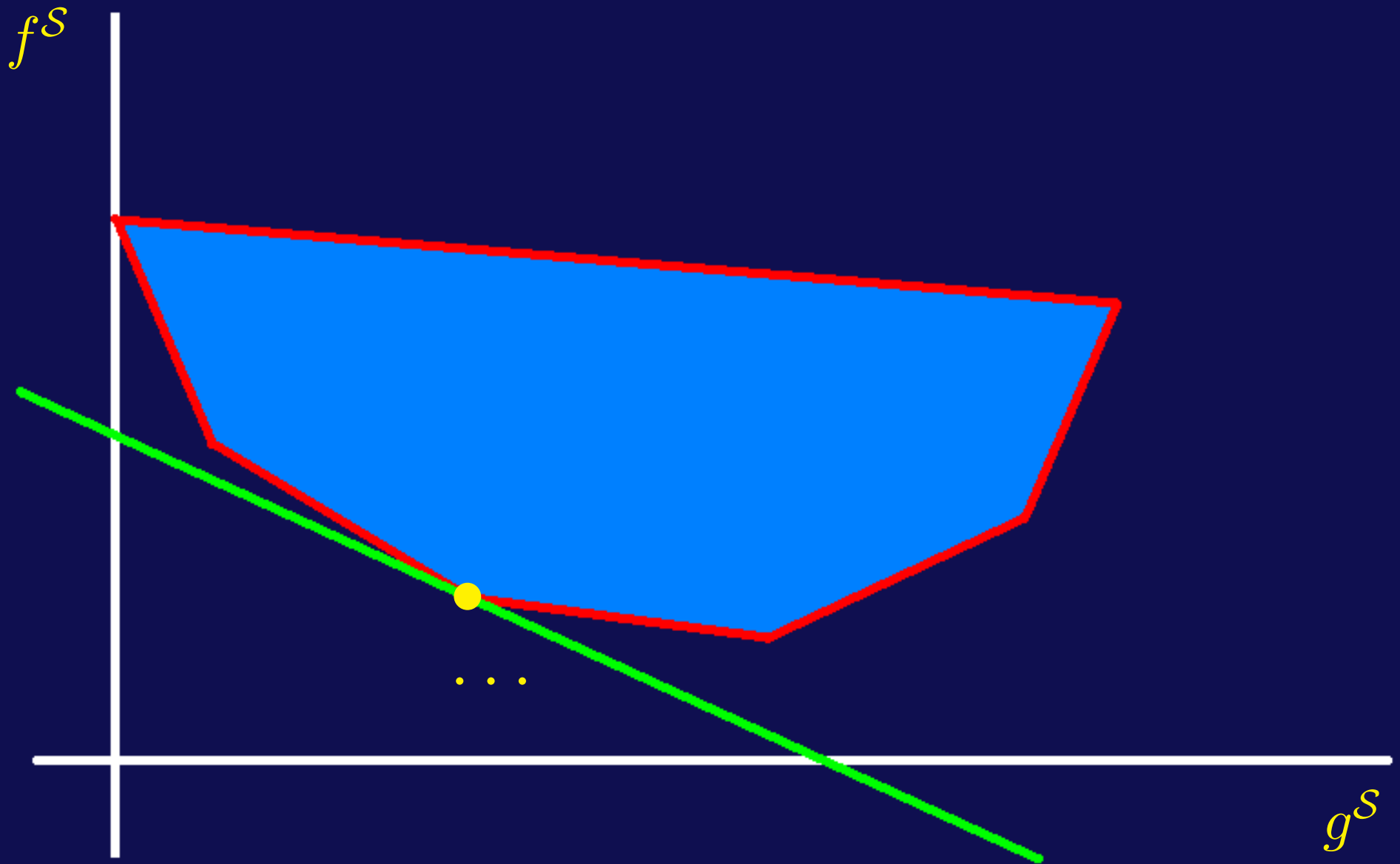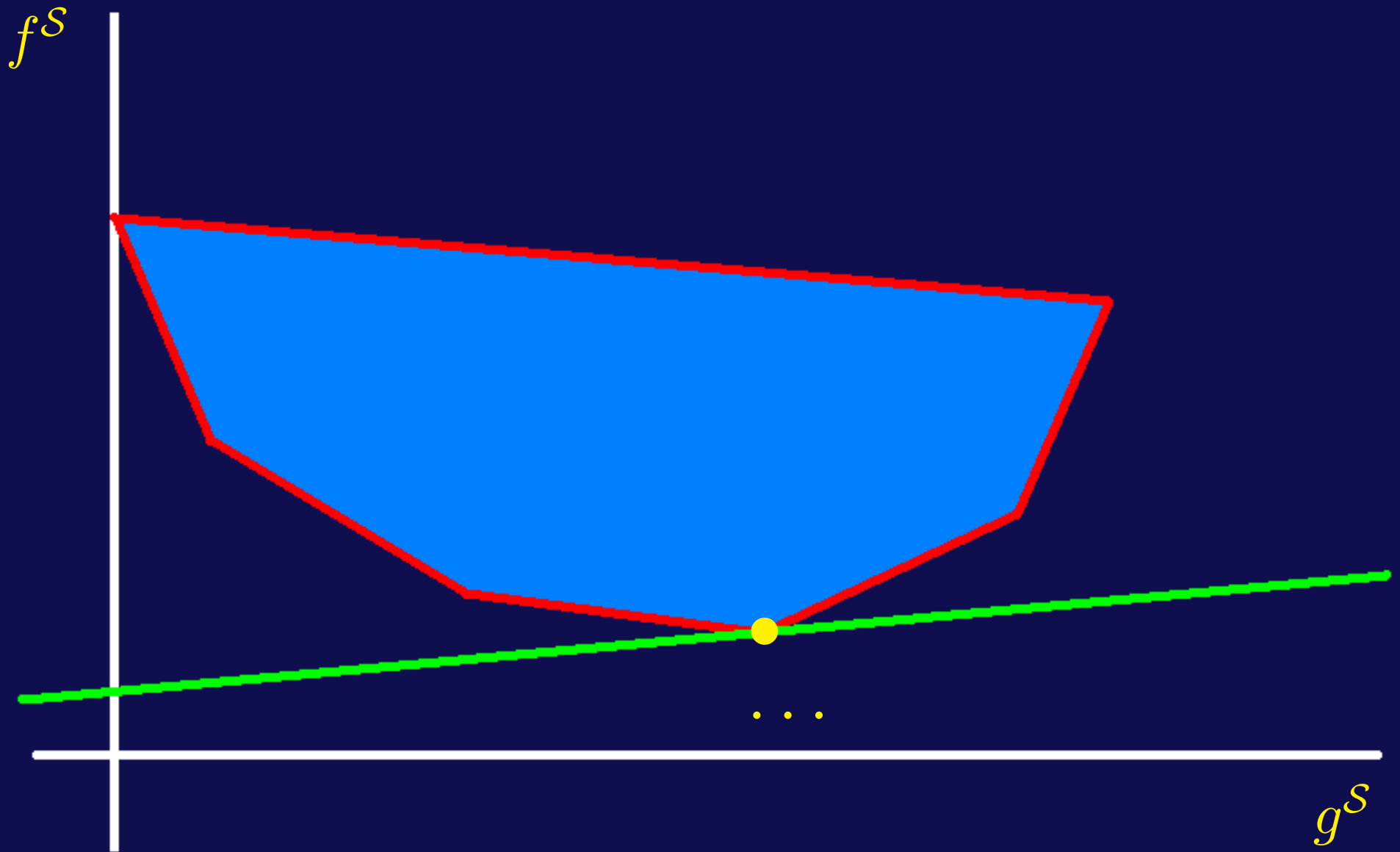  ▷ how much is worth shutting w.r.t. opening the gate
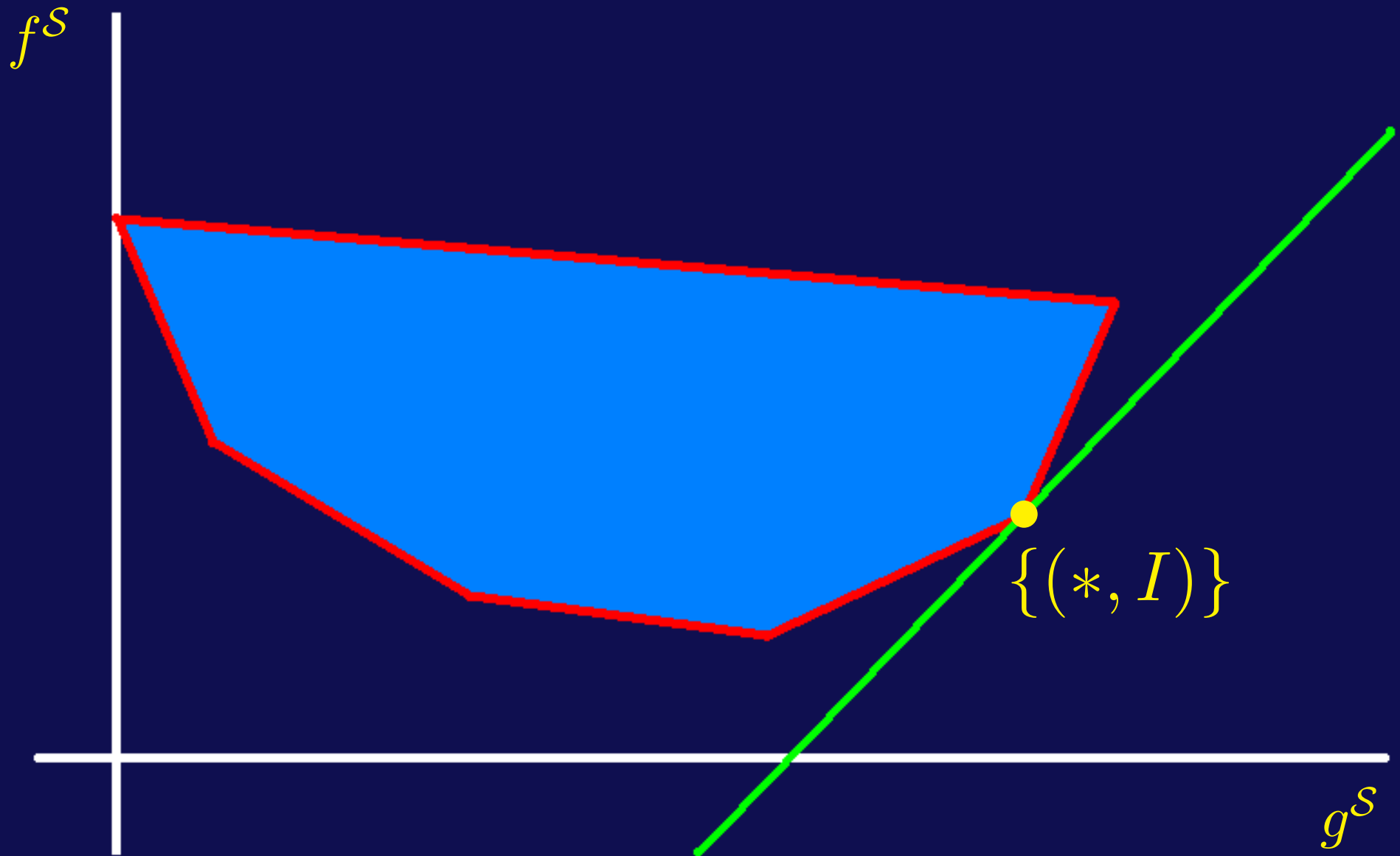
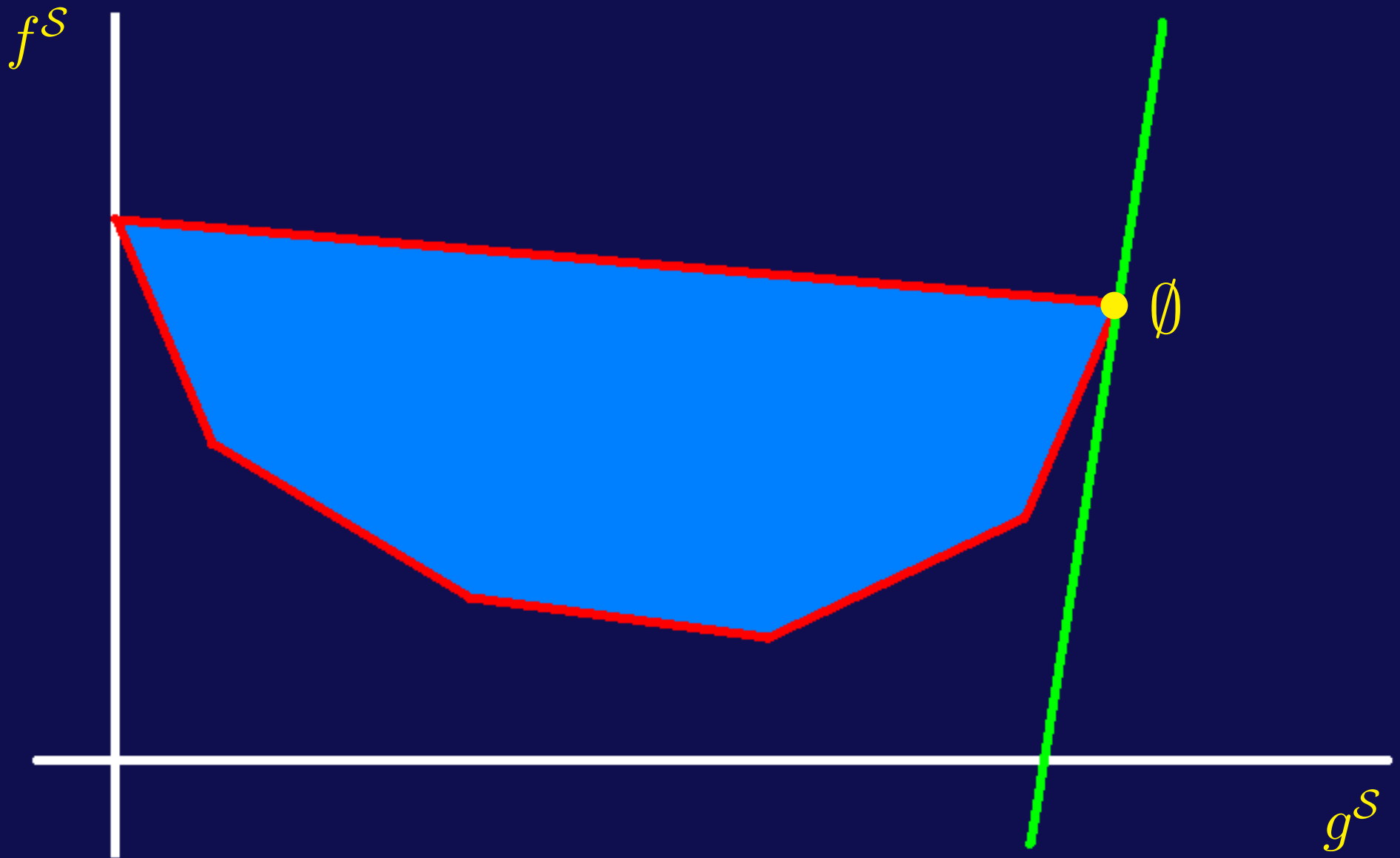# Indexability

# Indexability

# Indexability

# Indexability

# Indexability

# Indexability

# PCL-Indexability

- A sufficient condition for indexability

- $\nu$-wage problem is PCL($\mathcal{F}$)-indexable, if

(i) $w^{\mathcal{S}}_{(a,i)} > 0$ for each $\mathcal{S} \in \mathcal{F}$ and $(a,i) \in \widetilde{\mathcal{I}}$
(ii) there is an optimal $\mathcal{S} \in \mathcal{F}$ for every rejection cost $\nu$

  ▷ we establish (i) by proving $\Delta_1 g^{\mathcal{S}}_{(1,i)} := g^{\mathcal{S}}_{(1,i)} - g^{\mathcal{S}}_{(0,i)} > 0$
  — because $w^{\mathcal{S}}_{(a,i)}$'s are expected values of $\Delta_1 g^{\mathcal{S}}_{(1,i)}$'s
  ▷ Altman & Nain (1992) established (ii)

- Niño-Mora (ValueTools 2007): $\mathcal{O}(I^3)$ MPI algorithm

  ▷ here, we simplify it to $\mathcal{O}(I)$ under linear holding costs

# Fast Index Algorithm (FI)

```
{Input  I, λ, μ, c, β}
{Output {ν_(a,i)}_(a,i)∈Ĩ}
{Initialization}
```

$\zeta := \lambda(1-\mu); \quad \eta := \mu(1-\lambda); \quad \varepsilon := 1 - \zeta - \eta;$

$A_0 := 0; \quad A'_0 := \beta\zeta; \quad B := \beta\mu/(1-\beta+\beta\mu); \quad B' := \beta\zeta B + \beta(\mu-\eta); \quad C := c/(1-\beta+\beta\mu); \quad D_0 := 0;$

$\nu_{(1,0)} := \beta\zeta C/\lambda;$

$$\nu_{(0,0)} := \frac{\beta\zeta C}{\lambda} \cdot \frac{(1-\beta+\beta\mu)(1+\beta\lambda+\lambda\mu) + \beta\zeta(\mu+\beta\mu+\beta\zeta)}{(1-\beta+\beta\mu)(1+\beta\zeta) + \beta\zeta(\beta\zeta - B')};$$

```
{Loop}
for K = 1 to I − 1 do
```

$\quad A_K := \beta\zeta/[1-\beta+\beta\zeta+\beta\eta(1-A_{K-1})]; \quad A'_K := \beta\zeta+\beta(\mu-\eta)A_K; \quad D_K := (c+\beta\eta D_{K-1})A_K/(\beta\zeta); \quad Z_K := A_K A'_{K-1}/A'_K;$

$$f^0 := -\frac{\frac{\beta\zeta}{A_K}D_K + \beta\zeta(c+\beta\mu B D_{K-1}) + [c-\beta(\mu-\eta)\beta D_{K-1}]B'}{\frac{A'_K}{A_K} + \beta A'_{K-1}B' + \beta\zeta\beta\mu(1-BA_{K-1})};$$

$$f^1 := -\frac{\frac{\beta\zeta}{A_K}D_K + c\beta\zeta B A_{K-1} + [\beta\mu\beta\zeta + (1-\beta)\beta(\mu-\eta)]D_{K-1} + A'_{K-1}(c-\beta\zeta\beta C)}{\frac{A'_K}{A_K} + \beta A'_{K-1}B' + \beta\zeta\beta\mu(1-BA_{K-1})};$$

$$g^0 := \frac{\beta\lambda(1+B')}{\frac{A'_K}{A_K} + \beta A'_{K-1}B' + \beta\zeta\beta\mu(1-BA_{K-1})}; \quad g^1 := \frac{1+A'_{K-1}}{1+B'}g^0;$$

```
    if K > 1 then
```

$$\nu_{(0,K-1)} := \frac{[\beta(\mu-\eta)(D_{K-1}-c) + \beta\eta\beta\zeta D_{K-1} + \beta\zeta\beta\zeta C] - [\beta\eta Z_{K-1} + \beta\varepsilon]A'_{K-1}f^0 - \beta\zeta B' f^1}{\beta\lambda - [\beta\eta Z_{K-1} + \beta\varepsilon]A'_{K-1}g^0 - \beta\zeta B' g^1};$$

```
    end {if};
```

$$\nu_{(1,K)} := \frac{[\beta(1-\mu)\beta\zeta C + \beta\mu\beta(\mu-\eta)D_{K-1}] - \beta\mu A'_{K-1}f^0 - \beta(1-\mu)B'f^1}{\beta\lambda - \beta\mu A'_{K-1}g^0 - \beta(1-\mu)B'g^1};$$

```
end {for};
{Termination}
```

$A_I := \beta\zeta/[1-\beta+\beta\zeta+\beta\eta(1-A_{I-1})]; \quad A'_I := \beta\zeta+\beta(\mu-\eta)A_I; \quad D_I := (c+\beta\eta D_{I-1})A_I/(\beta\zeta); \quad Z_I := A_I A'_{I-1}/A'_I;$

$$f^0 := -\frac{\frac{\beta\zeta}{A_I}D_I - \beta(\mu-\eta)\beta\mu D_{I-1}}{\frac{A'_I}{A_I} + \beta\mu A'_{I-1}}; \quad g^0 := \frac{\lambda(1+\beta\mu)}{\frac{A'_I}{A_I} + \beta\mu A'_{I-1}};$$

$$\nu_{(0,I-1)} := \frac{[\beta(\mu-\eta)(D_{I-1}-c) + \beta\eta\beta\zeta D_{I-1}] - [\beta\eta Z_{I-1} + \beta\varepsilon]A'_{I-1}f^0}{\beta(1-\zeta)\lambda - [\beta\eta Z_{I-1} + \beta\varepsilon]A'_{I-1}g^0};$$

$$\nu_{(*,I)} := \frac{[\beta(\mu-\eta) + \beta\eta Z_I]D_{I-1} + cZ_I}{\lambda(1-Z_I)};$$

# Optimal Bi-Threshold Policy

- Can be obtained from MPIs (nondecreasing in $i$)

  ▷ the optimal open-gate threshold is

  $$K_0 := \min\{i \in \mathcal{I} : \nu_{(0,i)} \geq \nu\}$$

  ▷ the optimal closed-gate threshold is

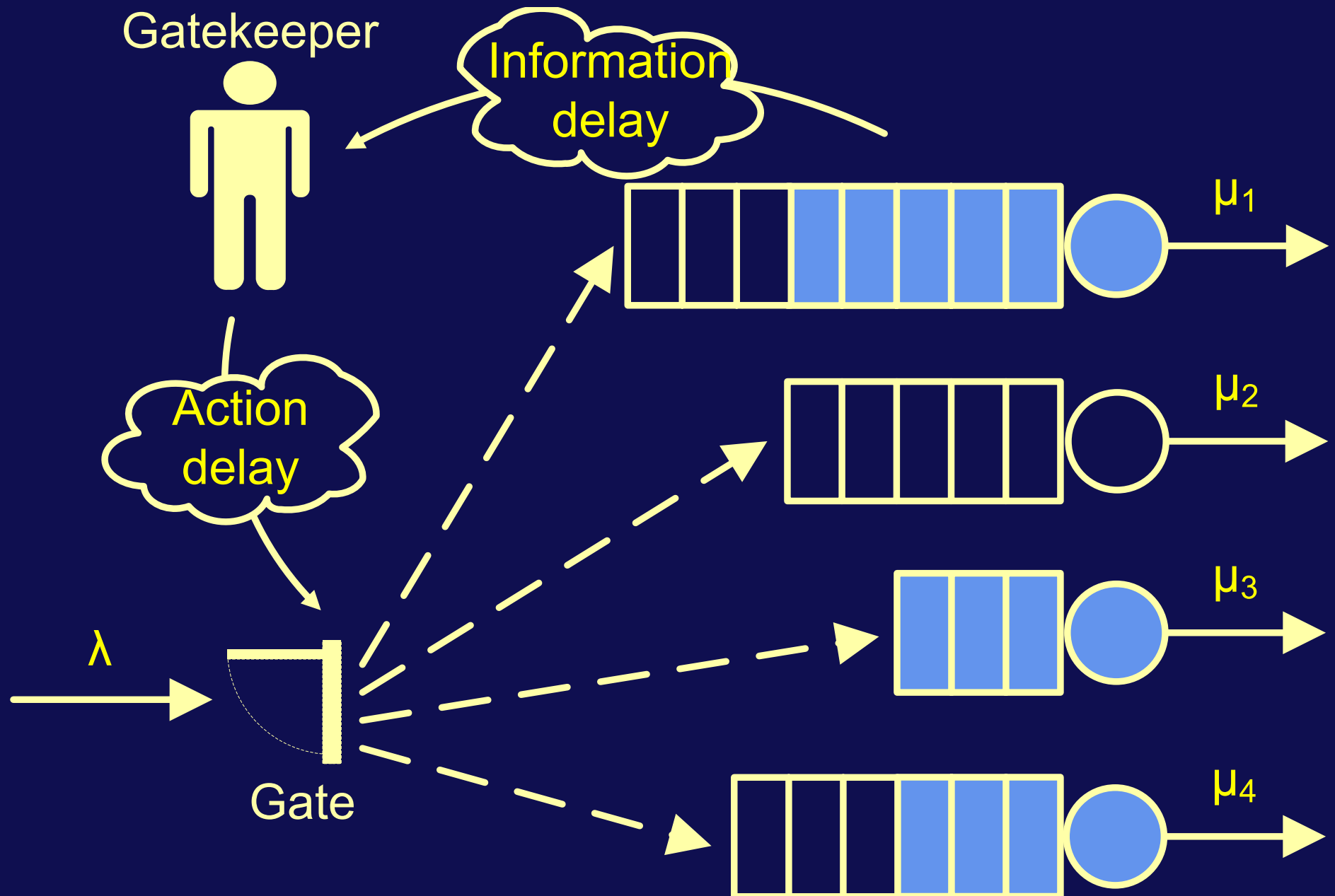  $$K_1 := \min\{i \in \mathcal{I} : \nu_{(1,i)} \geq \nu\}$$

  ▷ if $\nu > \nu_{(*,I)}$, then the gate is open always

- FI can be used also for infinite buffer (never stops)

- FI also works under the time-average criterion $(\beta = 1)$

# MPI Properties

- Both $\nu_{(0,i)}$ and $\nu_{(1,i)}$ are nondecreasing in $i$, nondecreasing in $\lambda$, nonincreasing in $\mu$

- Interleaving values:

  $\triangleright$ $\nu_{(0,i)} \leq \nu_{(1,i+1)} \leq \nu_{(0,i+1)}$
  $\triangleright$ $\nu_{(1,i)} \leq \nu_{(0,i)} \leq \nu_{(1,i+1)}$

- Convergence

  $\triangleright$ $\nu_{(1,i)} \rightarrow \nu_{(0,i)}$ as $\lambda \rightarrow 0$
  $\triangleright$ $\nu_{(1,i)} \rightarrow \nu_{(0,i-1)}$ as $\lambda(1-\mu) \rightarrow 1$
  $\triangleright$ $\nu_{(0,i)}, \nu_{(1,i)} \rightarrow \beta c/(1-\beta)$ as $i \rightarrow \infty$

- $\lambda\nu_{(1,0)} =$ the expected total discounted holding cost

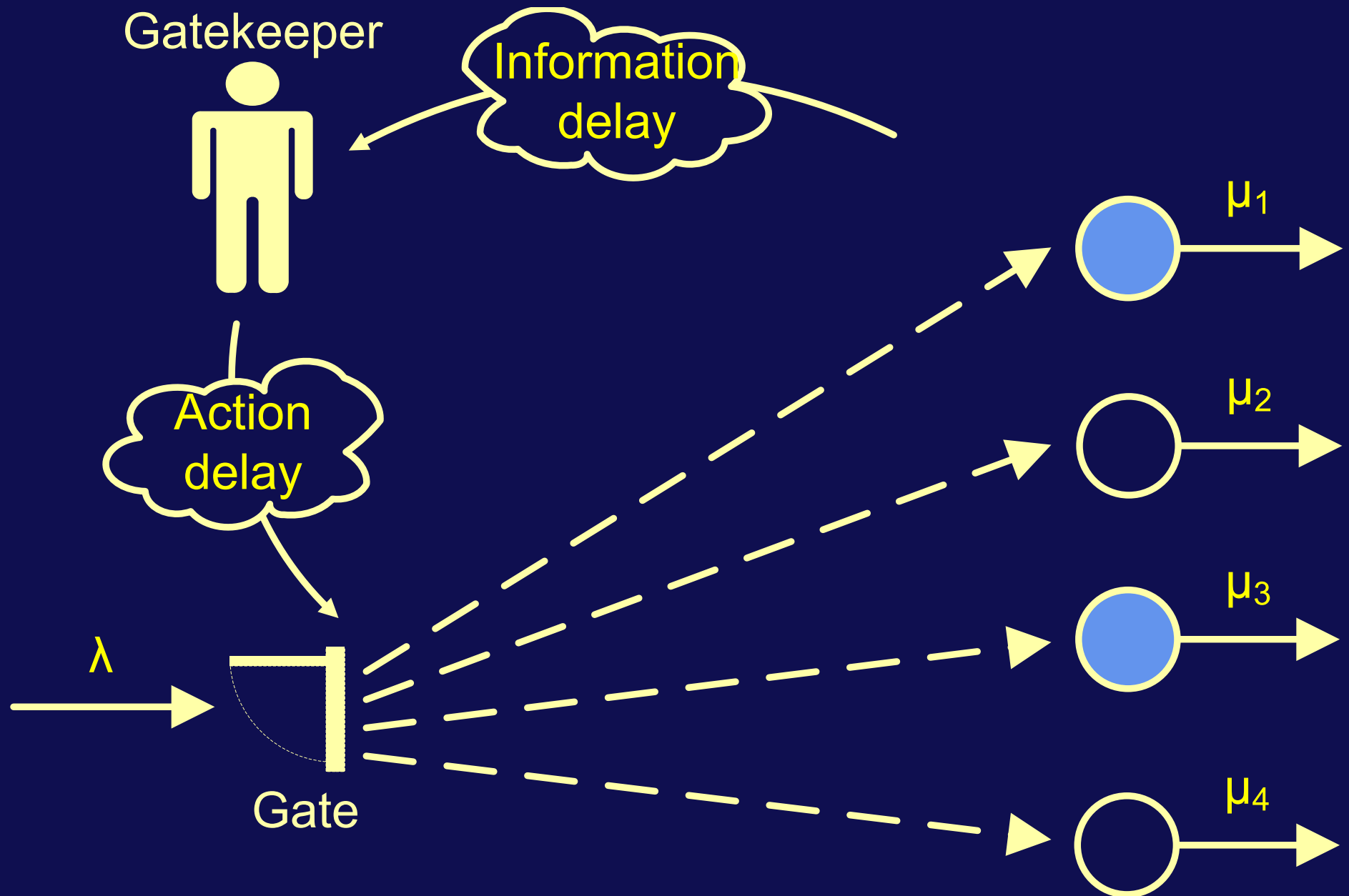# Admission Control and Routing with Delay

# Admission Control and Routing with Delay

- MPI policy for $K$ queues:

  ▷ Admit an arriving job iff
  $\nu > \nu_{\widetilde{X}_k(t)}$ for at least one queue $k$
  ▷ If admitted, route to the queue with lowest MPI

- By MPI properties, a job is routed to a queue with

  ▷ less waiting jobs
  ▷ faster server
  ▷ no job admitted in the previous period
  ▷ lower holding costs

- JSEQ is recovered in case of two symmetric queues

# Servers Assignment Problem with Delay

# Servers Assignment Problem with Delay

- The MPI is $\nu_{(a,i)} = \frac{c\beta(1-\mu)}{1-\beta(1-\mu)}$

  ▷ equal for all augmented states
  ▷ equal to the MPI with no delay
  ▷ equal under any arrival rate $\lambda$

- By MPI properties, a job is routed to a queue with

  ▷ faster server
  ▷ lower holding costs

- Jobs are routed always to the same queue

# Why MPI Policy is not Optimal?

- MPI policy is, in general, not optimal due to cross-dependence

  ▷ we do not know after-routing arrival rates for each queue; moreover, they may be time-varying
  ▷ computation of MPIs implicitly assumes that the threshold policy is the same in all periods

- MPI policy may be optimal in certain instances

- Mean behavior is nearly-optimal

# Summary of MPI Approach

- No news:

  ▷ analysis of problems with delays is hard

- Good news:

  ▷ yields tractable heuristics in heterogeneous problems
  ▷ powerful to obtain an exact algorithm of the same complexity as in the no-delay case
  ▷ some general patterns are extensible to other problems
  ▷ promising for larger delays

# Thank you for your attention!