# Speed Scaling under QoS Constraints with Finite Buffer
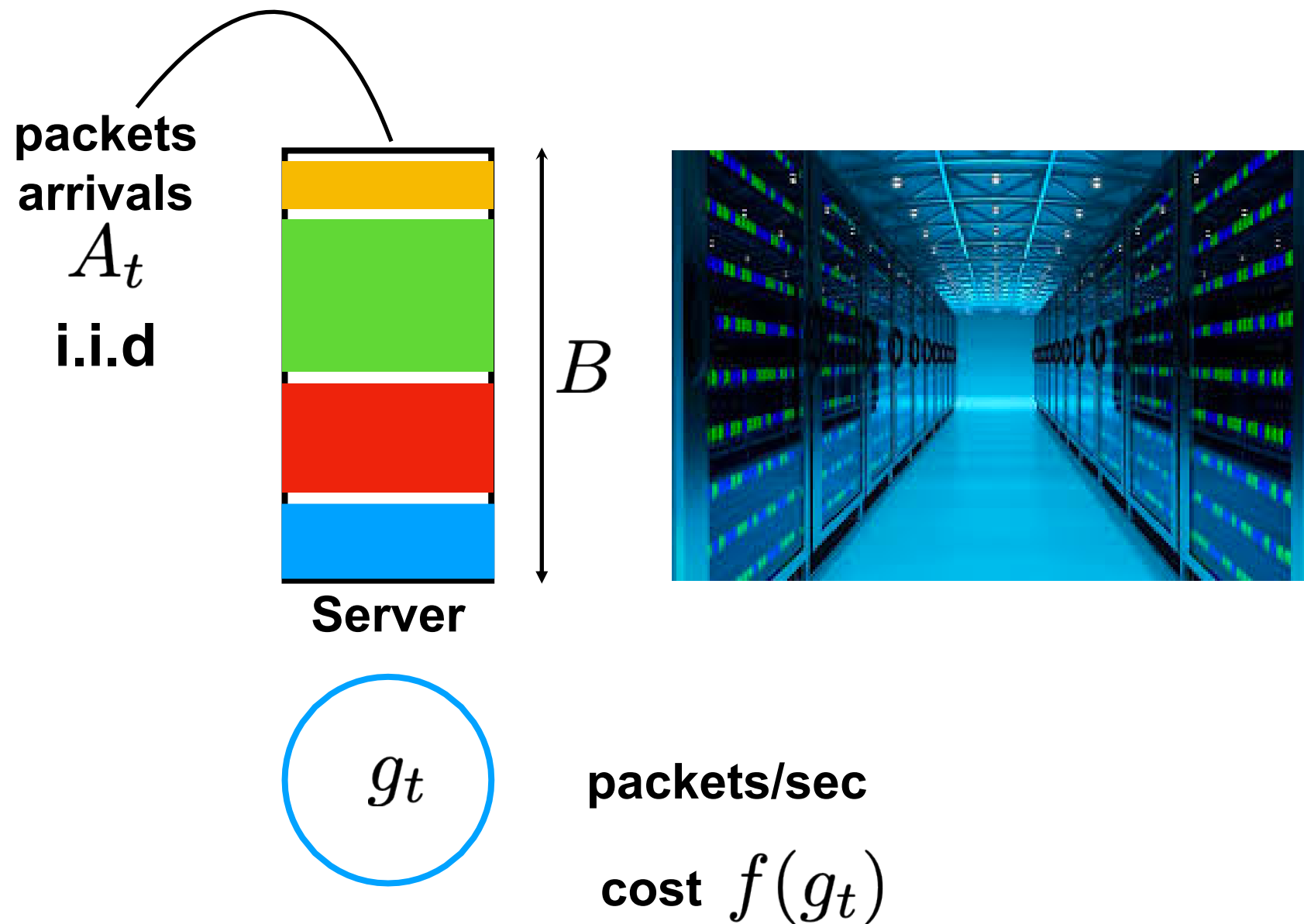
Parikshit Hegde, Akshit Kumar, and Rahul Vaze

# Motivation:Example
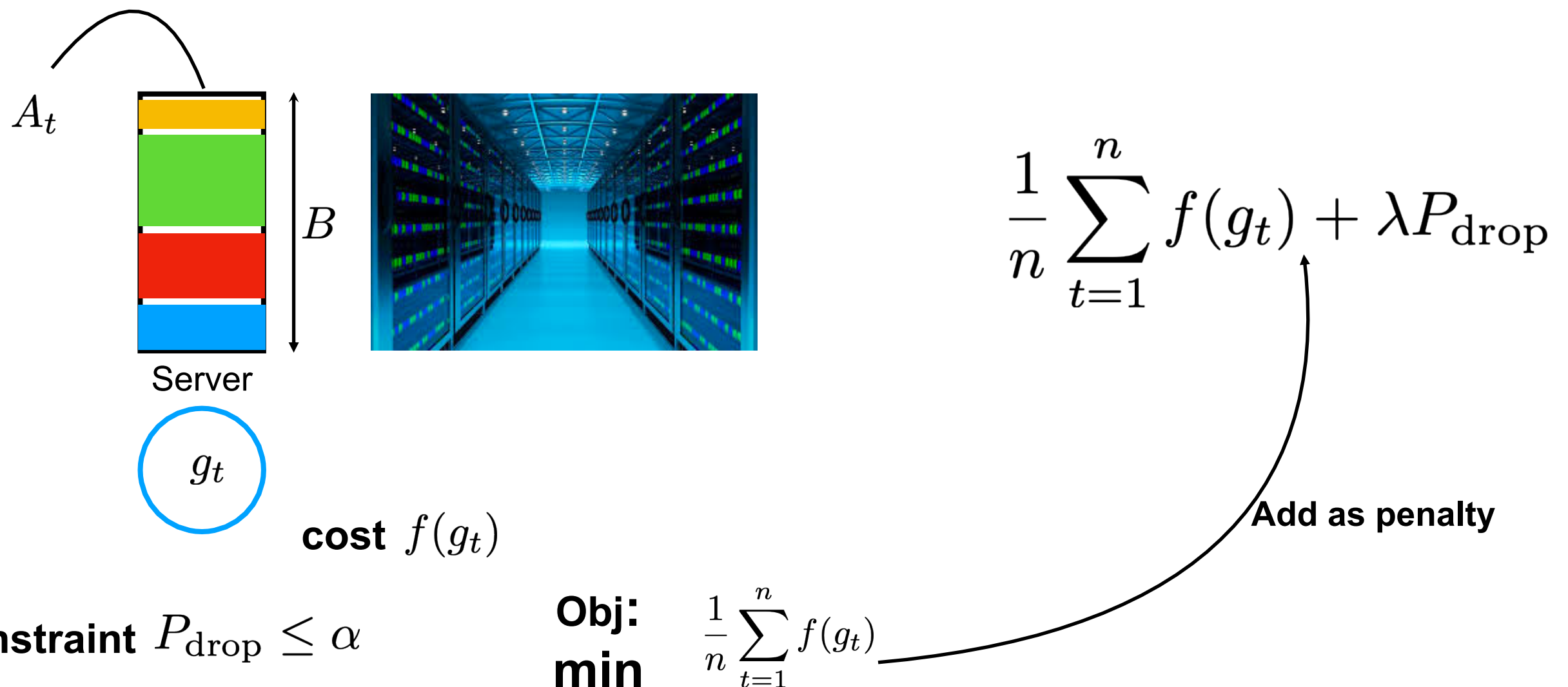


**packets arrivals**
$A_t$
**i.i.d**

$B$

**Server**

$g_t$    **packets/sec**

**cost** $f(g_t)$

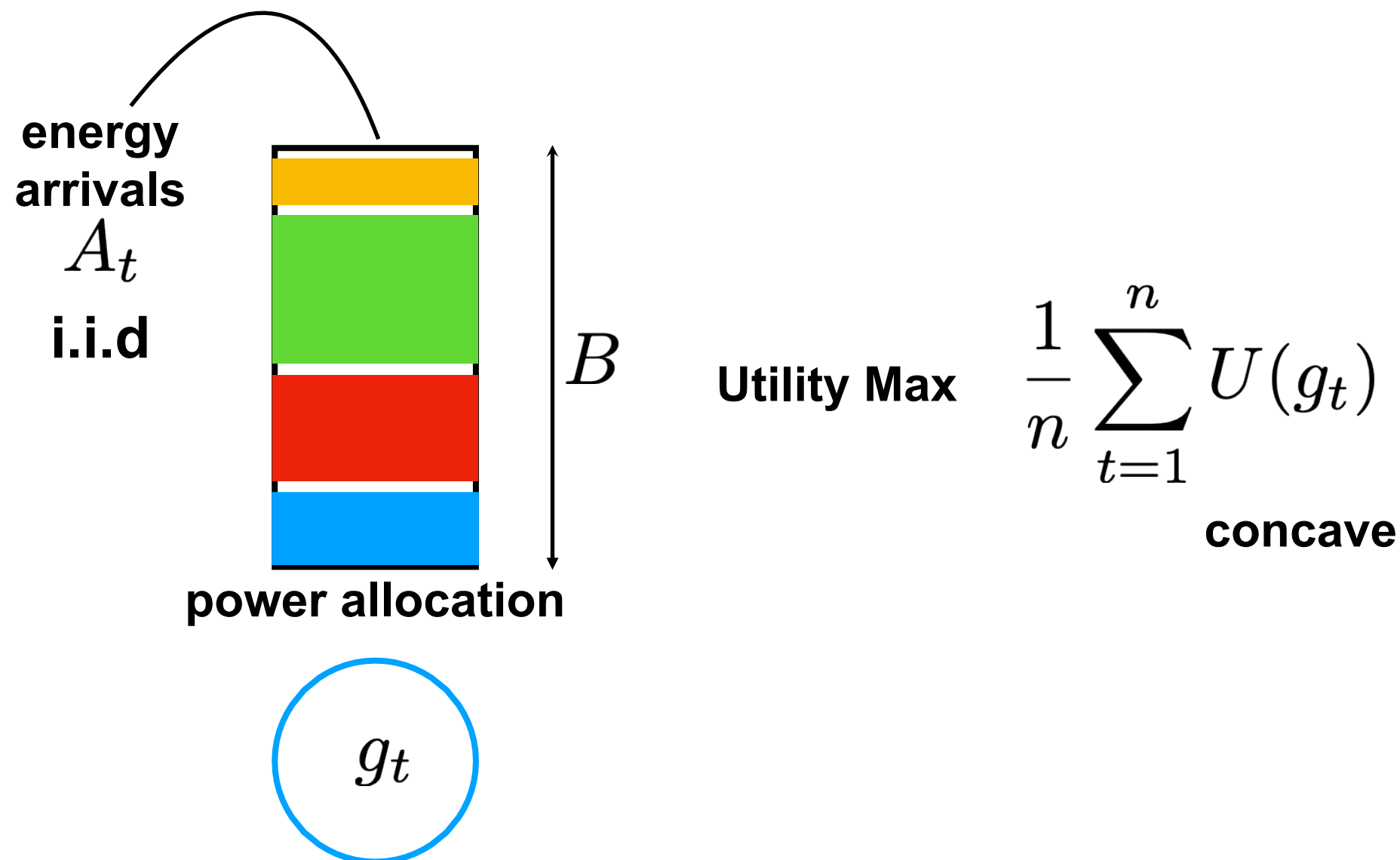**Constraint** $P_{\mathrm{drop}} \leq \alpha$     **Obj: min** $\dfrac{1}{n}\sum_{t=1}^{n} f(g_t)$

# Prior Work

- Problem formulated *Srikant &Perkins*, '99

- Computing Optimal Policies is challenging
- Simplified the problem
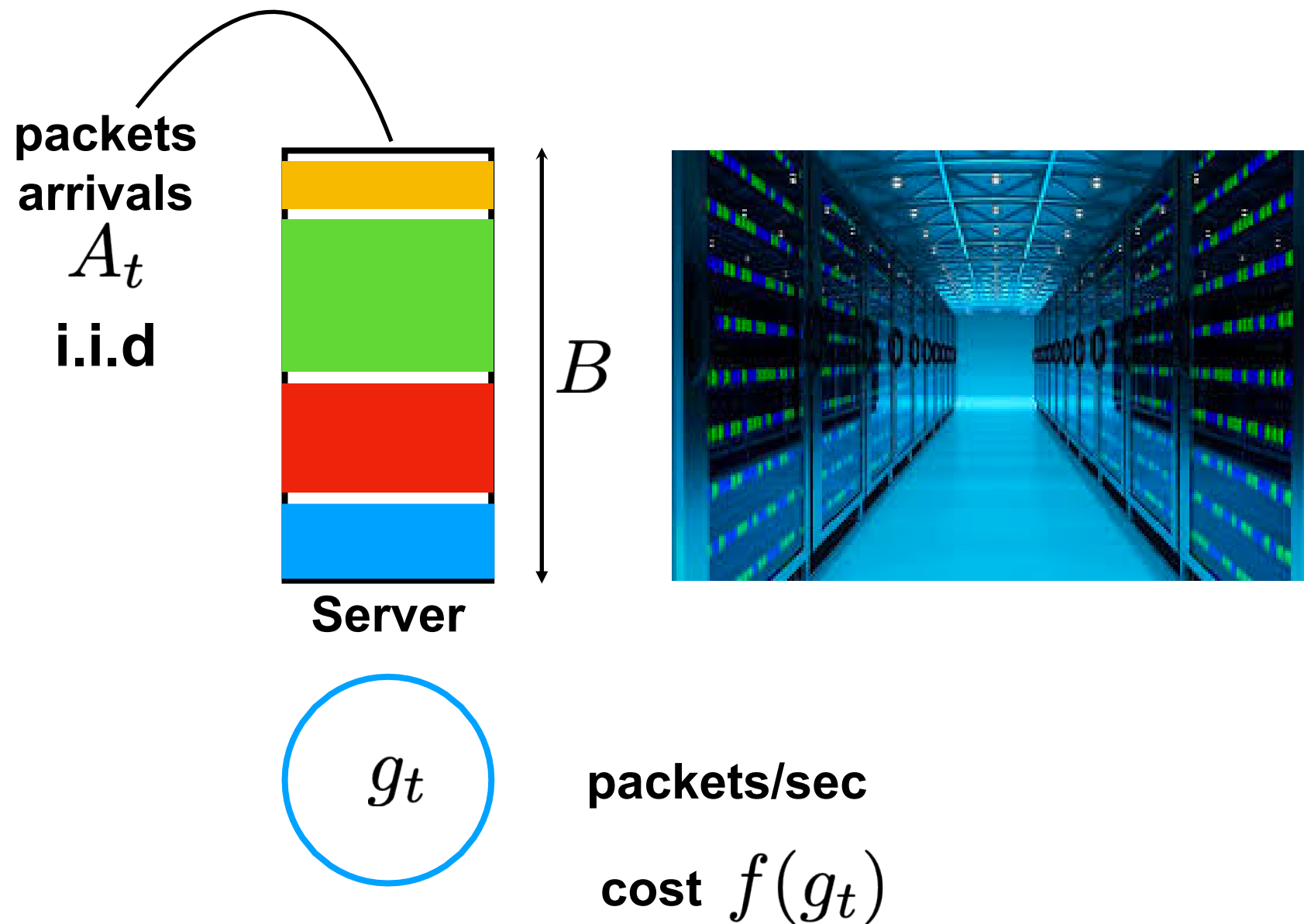  - We propose near optimal policies with provable guarantees

$A_t$

$B$

Server

$g_t$

cost $f(g_t)$

**Constraint** $P_{\mathrm{drop}} \leq \alpha$

**Obj:**
**min** $\quad \dfrac{1}{n}\sum_{t=1}^{n} f(g_t)$

$$\frac{1}{n}\sum_{t=1}^{n} f(g_t) + \lambda P_{\mathrm{drop}}$$

**Add as penalty**

# More Recent Work - EH

- [*Shaviv & Ayfer*, 2016] Competitive ratio of 2

- If B is large, [*Srivastava & Koksal*, 2013] provide near optimal policies

- In addition they show that battery overflow or battery discharge $\Theta\left(B^{-\beta}\right)$

**energy
arrivals**

$A_t$

**i.i.d**

$B$

**power allocation**

$g_t$

**Utility Max** $\frac{1}{n}\sum_{t=1}^{n}U(g_t)$

**concave**

# Recap - Problem



**packets arrivals** $A_t$ **i.i.d**

$B$

**Server**

$g_t$

**packets/sec**

**cost** $f(g_t)$

**Constraint** $P_{\mathrm{drop}} \leq \alpha$

**Obj: min** $\dfrac{1}{n}\sum_{t=1}^{n} f(g_t)$

# General Lower Bound



$$P_{\mathrm{drop}} \leq \alpha$$

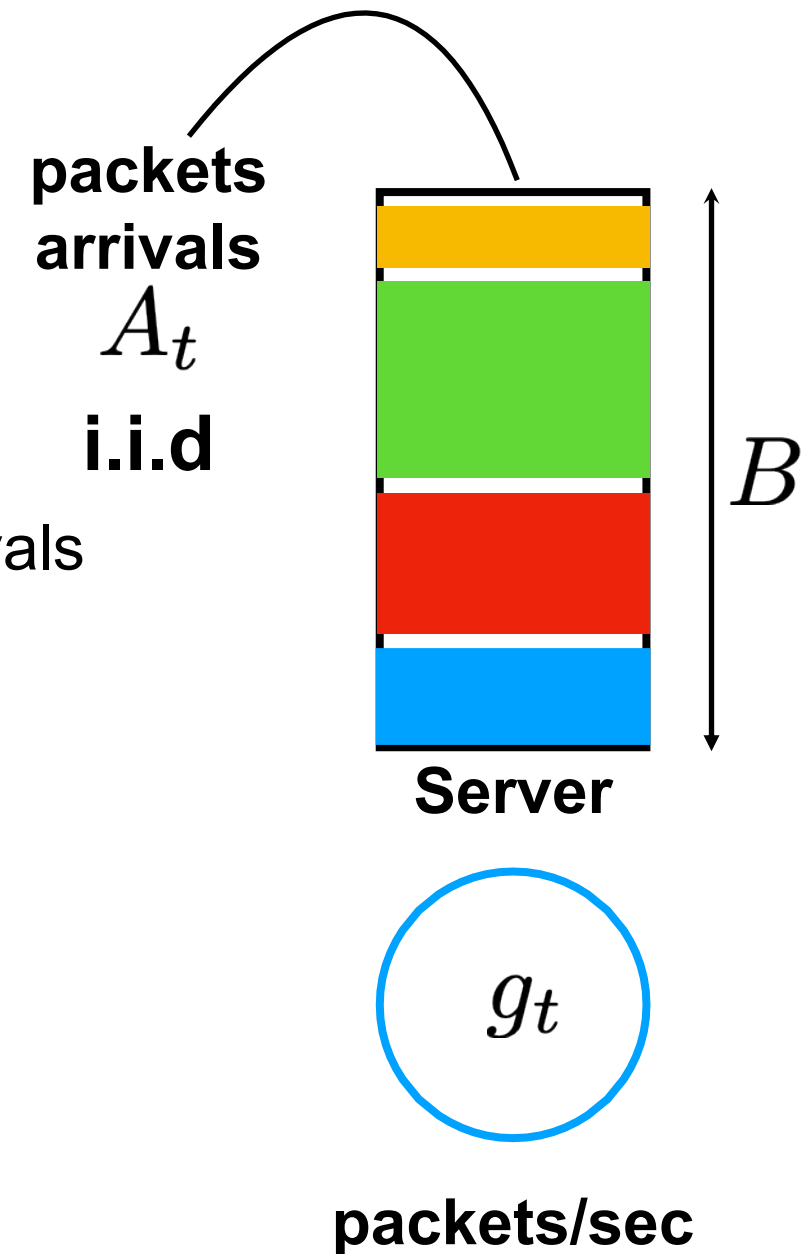Service at least $(1 - \alpha)$ fraction of total packet arrivals

$$\mathbf{E}\left[\frac{\sum_{t=1}^{n} g_t}{n}\right] \geq (1 - \alpha)\mathbf{E}\left[\frac{\sum_{t=1}^{n} A_t}{n}\right]$$

$$\lim_{n \to \infty} \mathbf{E}\left[\frac{\sum_{t=1}^{n} g_t}{n}\right] \geq (1 - \alpha)\mu$$

Invoking Jensen's

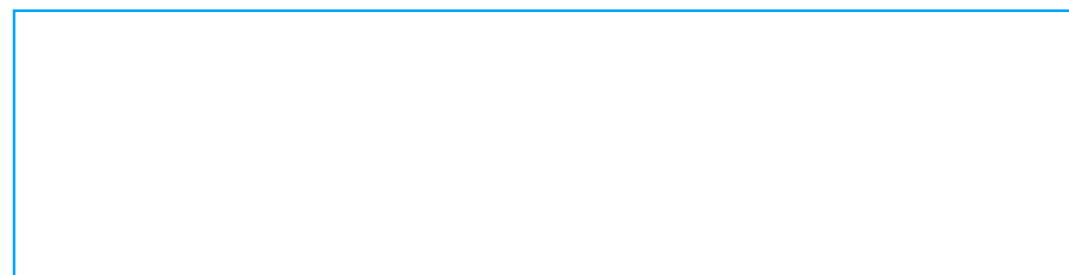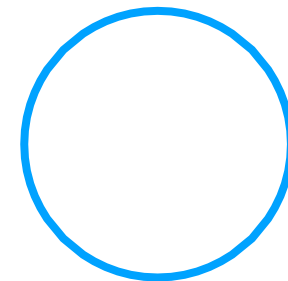$$\mathrm{Cost} \geq f((1 - \alpha)\mu)$$

**packets arrivals** $A_t$ **i.i.d**

$B$

**Server**

$g_t$

**packets/sec**

# Greedy Policy

- Forcefully drop $\alpha$ fraction of arrivals

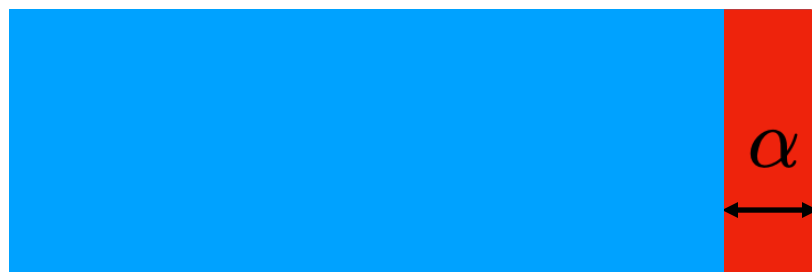- Immediately serve the $1 - \alpha$ fraction of arrivals

**Finite Buffer**

**Server**

$B$

**Serve all packets**

$\alpha$

**dropped**

$A_t$

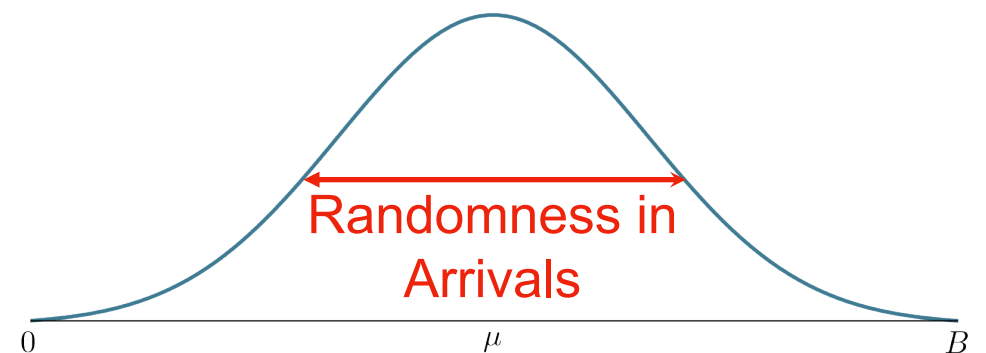# Greedy Policy:Performance

$$f(x) = x^2$$

- Immediately serve the $1 - \alpha$ fraction of arrivals

$$\text{Cost} \leq \lim_{n \to \infty} \mathbf{E}\left[\frac{1}{n}\sum_{t=1}^{n}(1-\alpha)^2 A_t^2\right]$$

$$\leq (1-\alpha)^2\left(\mu^2 + \text{var}\left(A_t\right)\right)$$

**Example Arrival Distribution**



Randomness in Arrivals

$0$     $\mu$     $B$

**Total Cost :** $(1-\alpha)^2(\mu^2 + var(A_t))$

Term from Lower Bound     Second Order Moment

sufficient in practice

# Greedy Policy:Performance

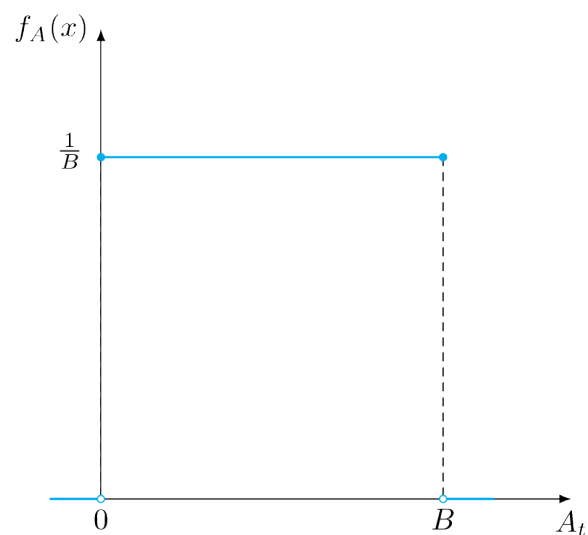**Competitive Ratio(CR) =** $1 + \text{var}(A_t)/\mu^2$

**Lower Bound**
$(1-\alpha)^2\mu^2$
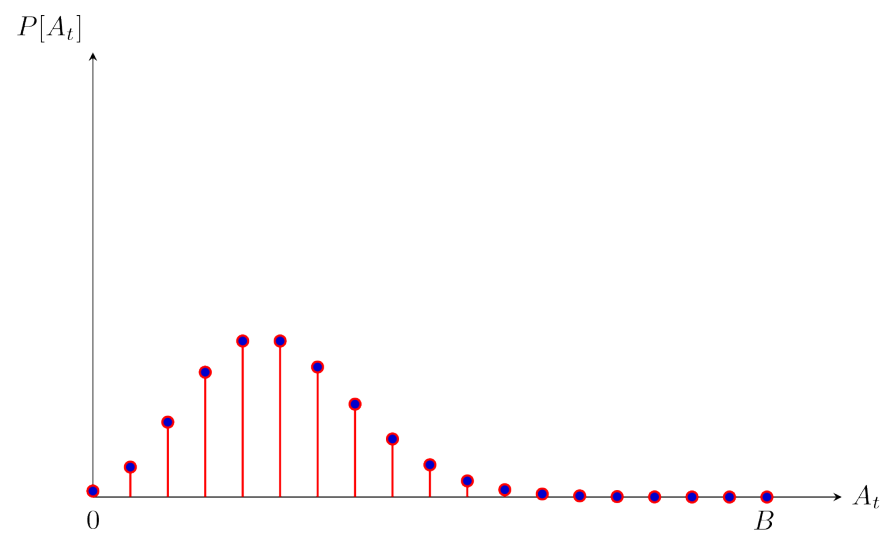
**Greedy Policy**
$(1-\alpha)^2(\mu^2 + var(A_t))$

# Greedy Policy:Performance

**Competitive Ratio(CR) =** $1 + \text{var}(A_t)/\mu^2$
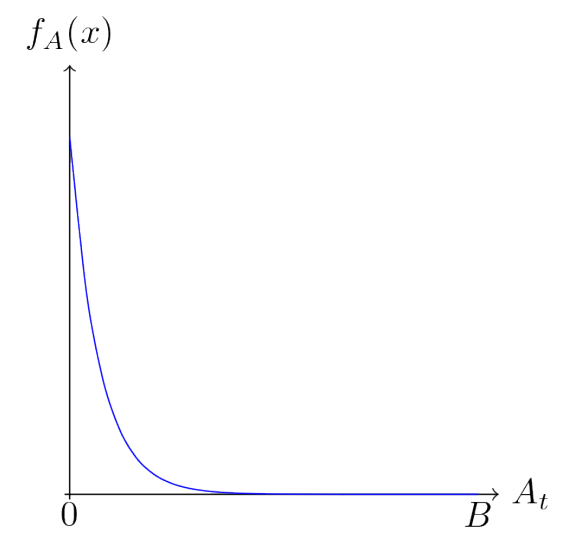


$A_t \sim \text{Unif}[0, B]$

$\text{CR} \leq \dfrac{4}{3}$

$A_t \sim \text{Poisson}(\nu), \nu \geq 1$

$\text{CR} \leq 3$

$A_t \sim \text{Exp}(\mu)$

$\text{CR} \leq 2$

# Greedy Policy:Performance

Bernoulli

$$A_t = \begin{cases} m, & \text{w.p. } \frac{\mu}{m}, \\ 0, & \text{w.p. } 1 - \frac{\mu}{m}, \end{cases} \qquad 0 < m \leq B$$
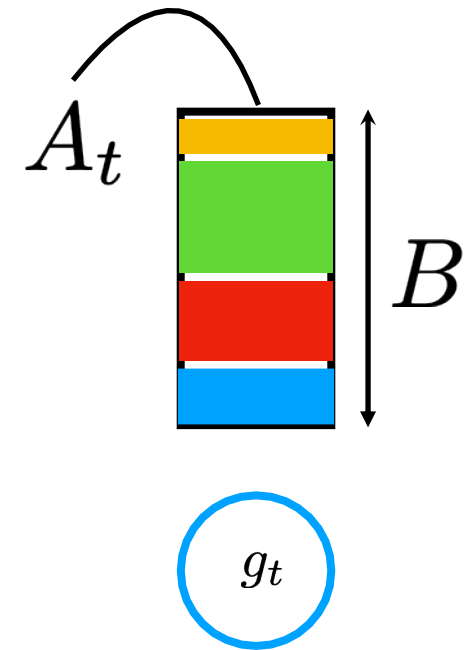
$$\text{CR}_{\text{greedy}} \leq 1 + \frac{\text{var}[X]}{\mu^2} = \frac{m}{\mu}$$

$$\mu \to 0$$

# Special Case Analysis
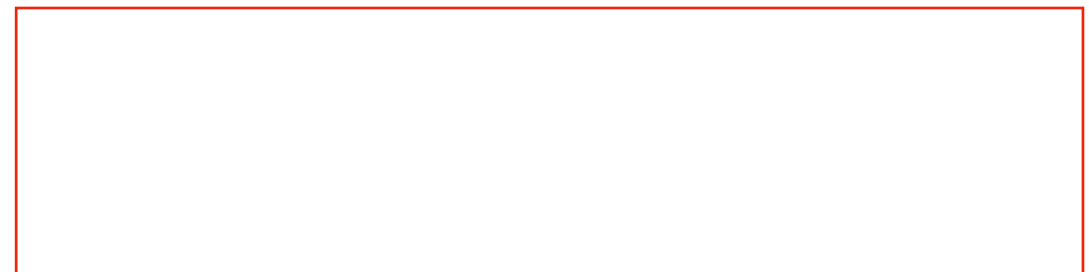
$$\mu \to 0$$

# Small $\mu$ regime



- Focus on *Extreme Bernoulli Distribution*

**With probability** $p$



$$A_t = B$$

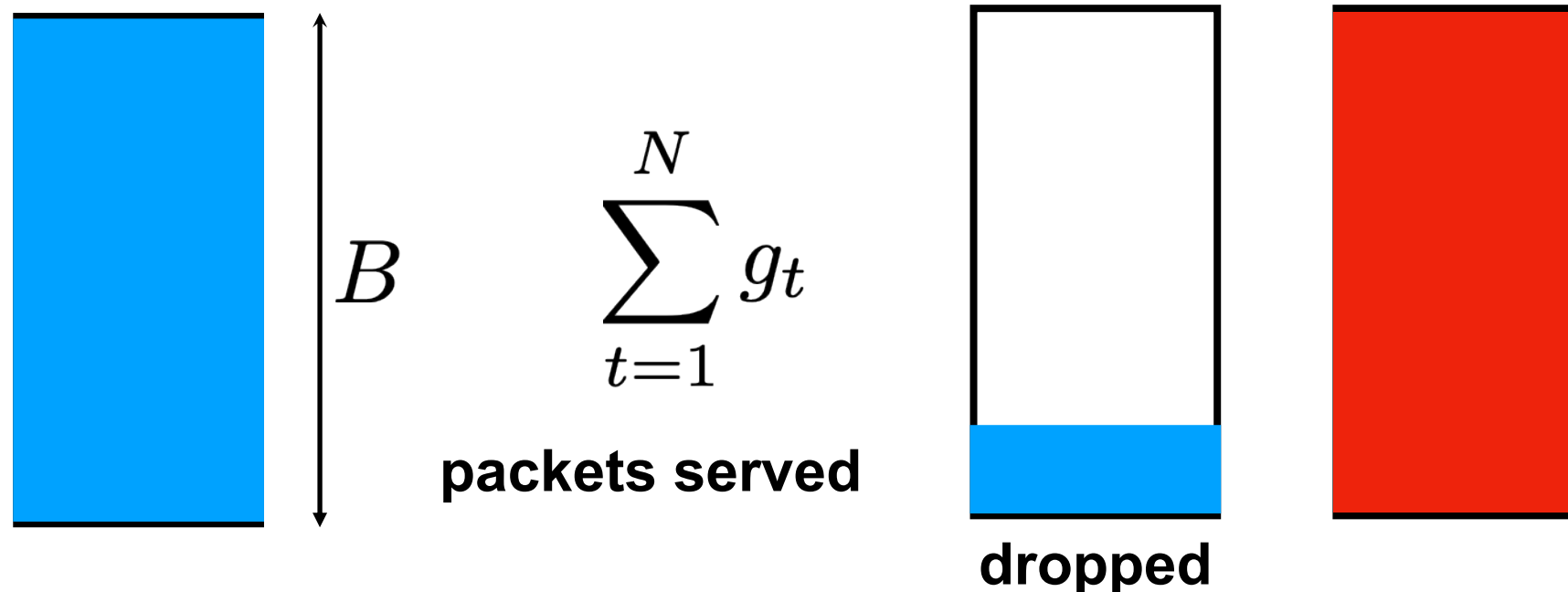**With probability** $1 - p$



$$A_t = 0$$

- Why this distribution?

- Conjecture that this is the worst case

- For similar problem in Energy Harvesting, Shaviv and Ayfer, 2016 showed this is worst case distribution

# Extreme Bernoulli Dist

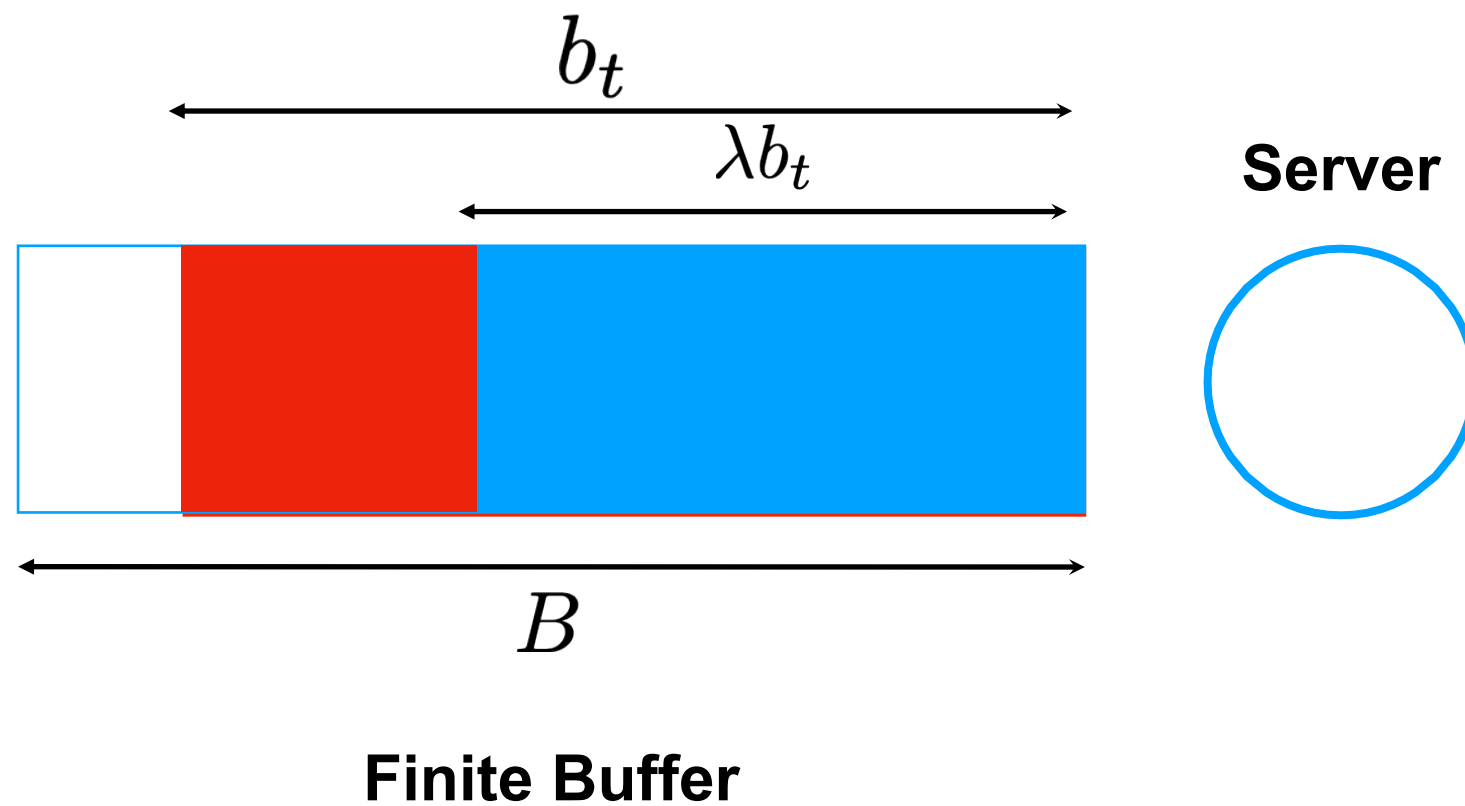$$P_{\text{drop}} = \frac{\mathbf{E}\{B - \sum_{t=1}^{N} g_t\}}{B}$$

$$P_{\text{drop}} \leq \alpha \qquad g_1 \geq B\left(1 - \frac{\alpha}{p}\right)$$

$$\text{Cost} = \frac{\mathbf{E}\left[\sum_{t=1}^{N} g_t^2\right]}{\mathbf{E}[N]} \geq \frac{g_1^2}{\mathbf{E}[N]} \geq \frac{B^2 p^2}{\alpha + p}$$

**Renewal Event**

$$N$$

$$B \qquad \sum_{t=1}^{N} g_t$$

**packets served**

**dropped**

# $\lambda$ -fraction Policy

- Serve $\lambda$ fraction of the packets in the buffer

- $\lambda$ calculated as a function of $p$ to satisfy $P_{drop}$ of $\alpha$



**Finite Buffer**

# Performance

For small $\alpha << 1, p << 1$

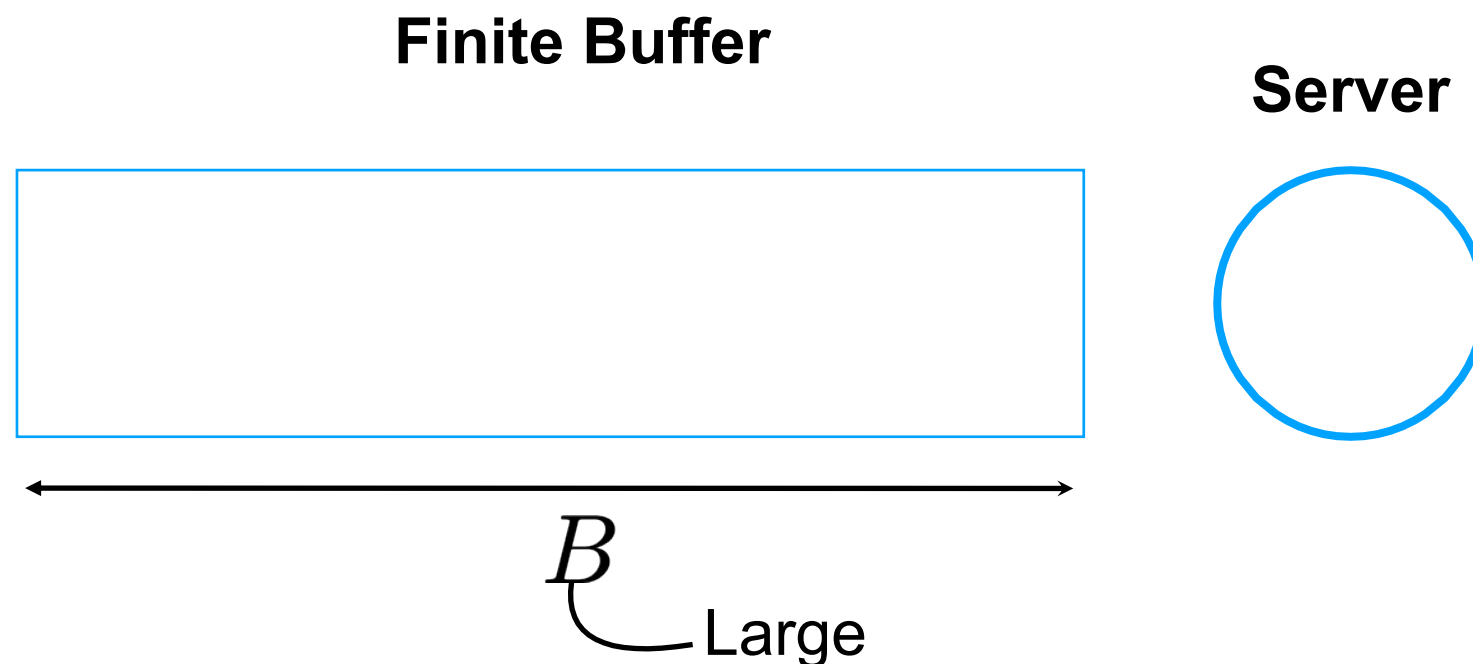$$\mathcal{J}^* = \Omega\left(\frac{\mu^2}{p + \alpha}\right)$$

Lower Bound

$$\mathcal{J}_\lambda \approx \frac{\mu^2}{2\alpha + p}$$

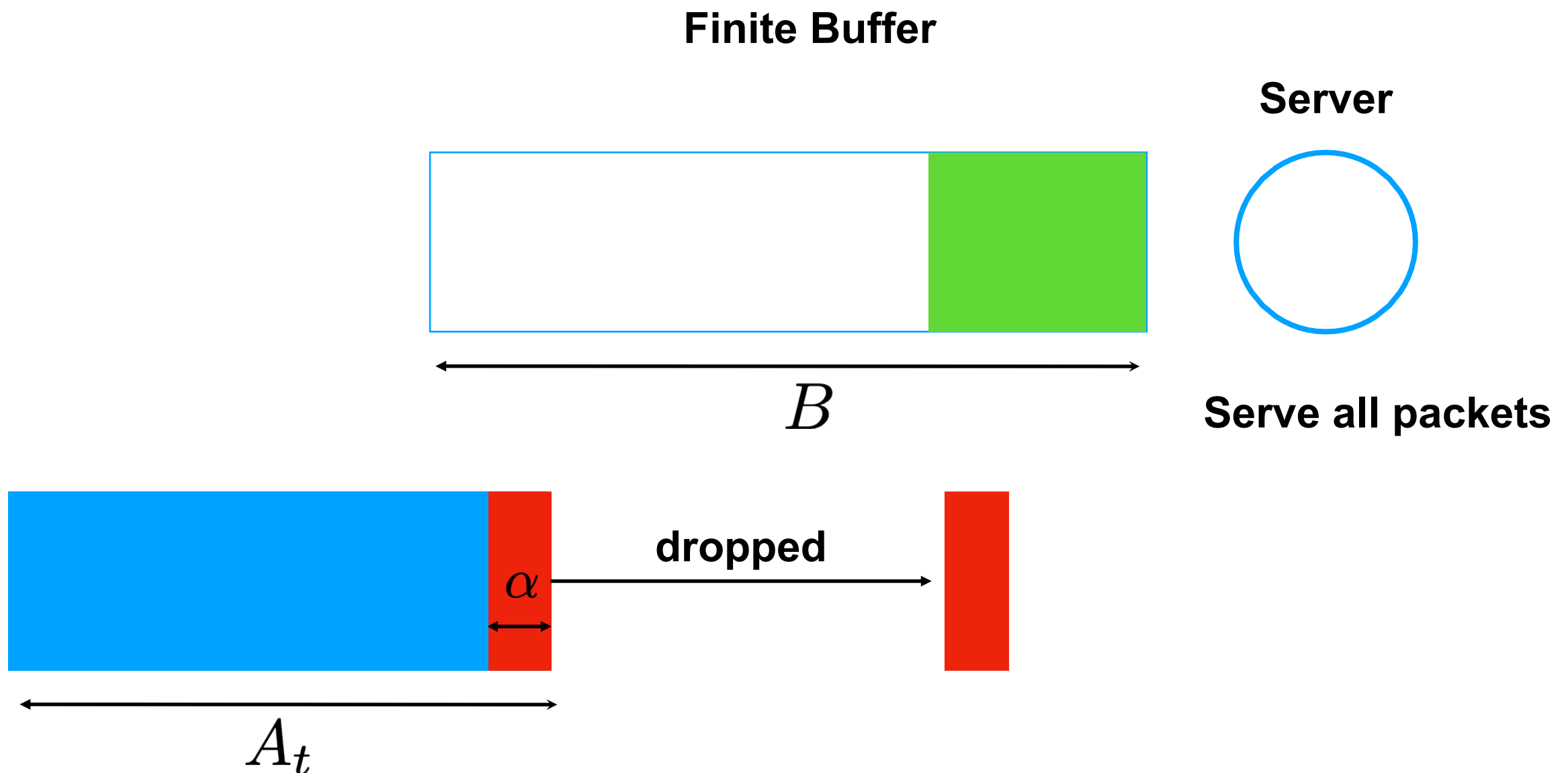$\lambda$ Fraction Policy

# Large Buffer Case

# Large Buffer Case

- Consider the case where buffer size is large

- Idea: Large buffer will reduce the "effect" of randomness in arrivals

- General Lower Bound that assumes $\mu$ arrivals every time likely to be achievable

**Finite Buffer**

**Server**

$B$

Large

# Admission Policy

- Forcefully drop $\alpha$ fraction of arrivals

**Finite Buffer**

**Server**



$B$

**Serve all packets**

$\alpha$

**dropped**

$A_t$

# Scheduling Policy

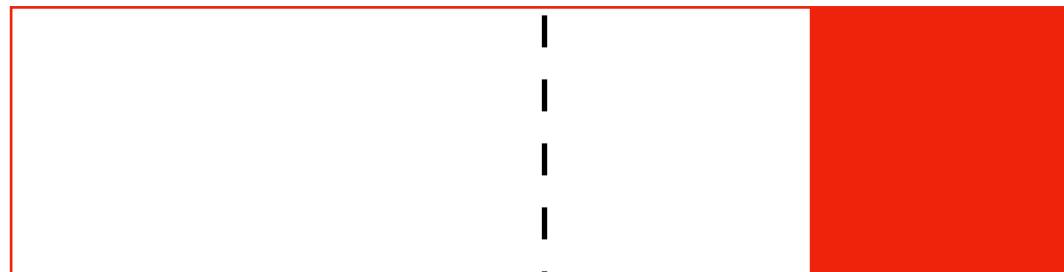**Buffer/Queue state** $b_t$     $b_t > B/2$

**Finite Buffer**

$g_t = \mu(1 - \alpha) + \delta$

$\delta$   **maintain Packet Drop constraint**   $B/2$

$\delta \quad \sim \beta \left( \dfrac{\log B}{B} \right)$

$b_t \leq B/2$

$g_t = \mu(1 - \alpha) - \delta$

# Results

$$\text{Cost} \leq \text{Opt} + \Theta(B^{-\beta}) + \Theta\left(\beta^2 \left(\frac{\log B}{B}\right)^2\right)$$

General Lower Bound

$\beta$ **tradeoffs the two terms**

$$P_{\text{drop}} \leq \alpha + \Theta(B^{-\beta})$$

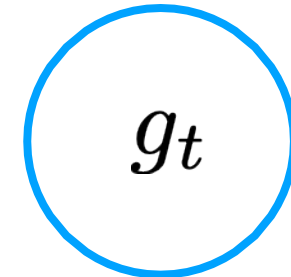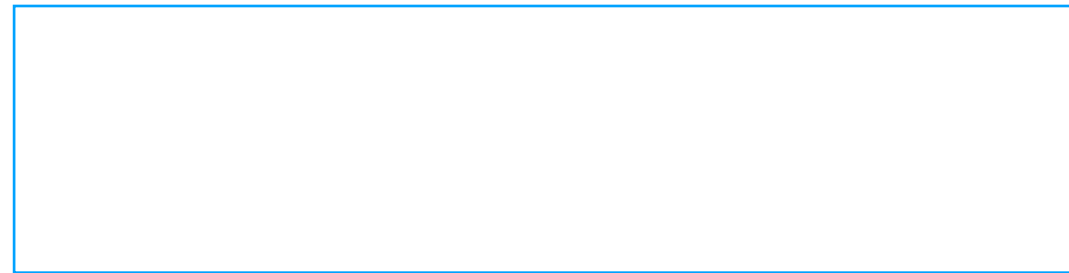Small Violation to the Packet Drop Constraint

# Summary

- Problem of Speed Scaling a server

    - QoS constraint imposed by Packet Drop

    - Minimise total cost

- Greedy Policy is near optimal when **var to mean ratio** is small

- For low arrival rate

    - Focus on *extreme Bernoulli distribution*

    - $\lambda$ fraction policy is near optimal

- For more general distributions

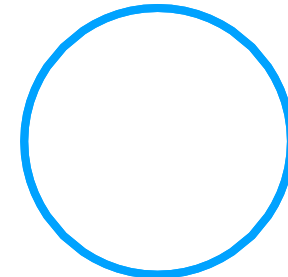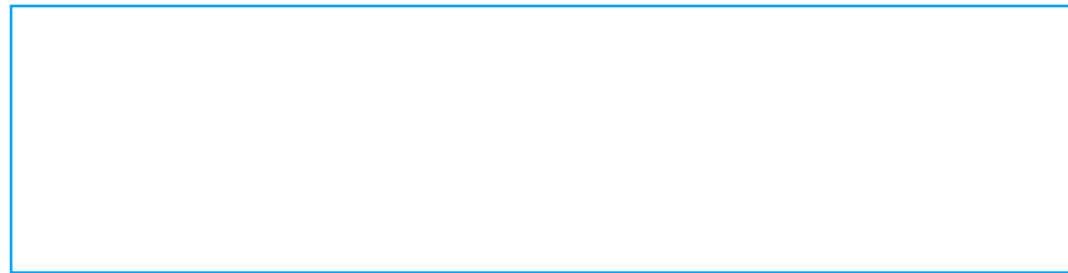    - Near Optimal policy when buffer size is large

# System Model

**Finite Buffer**

**Server**
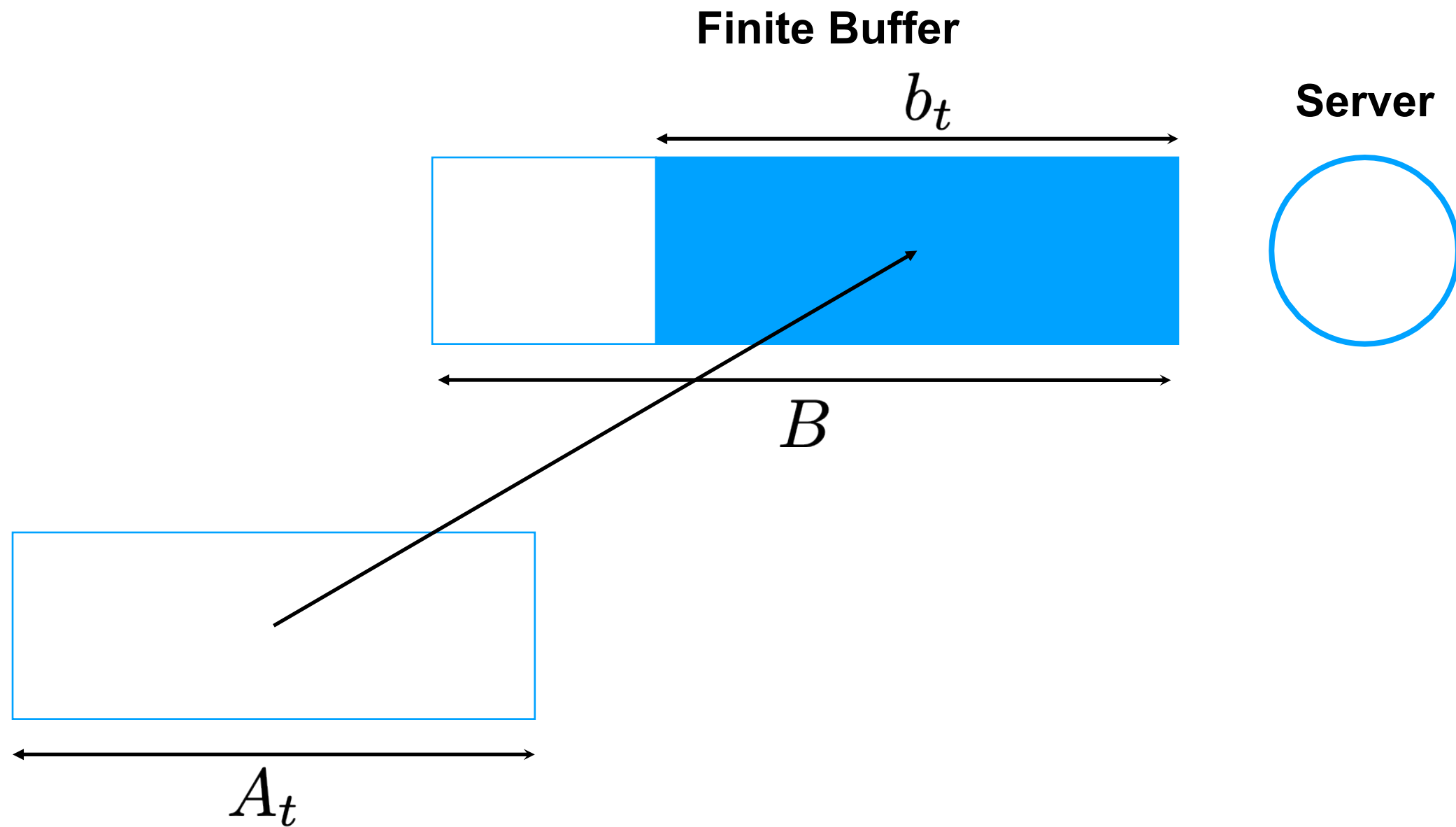
$g_t$

$B$

# System Model

**Finite Buffer**

**Server**

$$B$$

$$A_t$$

# System Model

**Finite Buffer**

**Server**

$b_t$

$B$

$A_t$

# System Model



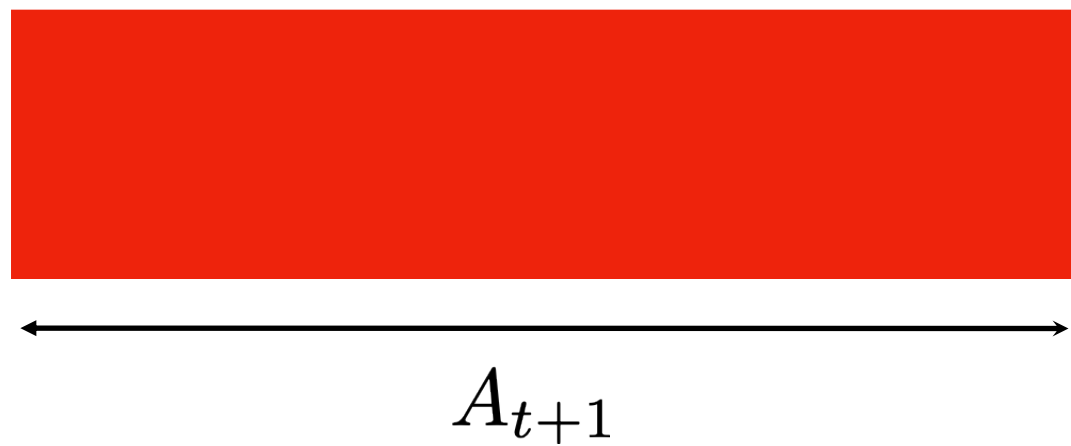**Finite Buffer**

$b_t$

**Server**

$g_t$

$B$

Cost at time t= $g_t^2$

Close approximation of microprocessor power cost (Wiermen et al. 2012)

# System Model

**Finite Buffer**

$b_t$

**Server**

$B$

$g_t$

$A_{t+1}$

$A_t, A_{t+1}$ **are**

# System Model

**Finite Buffer**

$b_{t+1}$

**Server**

$B$

$A_{t+1}$

**Packets Dropped**

# Problem Formulation



- Overall Cost = $\dfrac{1}{n}\displaystyle\sum_{t=1}^{n} g_t^2$

- P$_{\text{drop}}$ = (Total # of Packets Dropped)/(Total # of Packets Arrived)

- Constraint: $P_{\text{drop}} \leq \alpha$

- Objective: Minimise overall cost subject to constraint