

"The Regulation of Queue Size by Levying Tolls" by P. Naor (1969)

Paper presentation for BZAN605

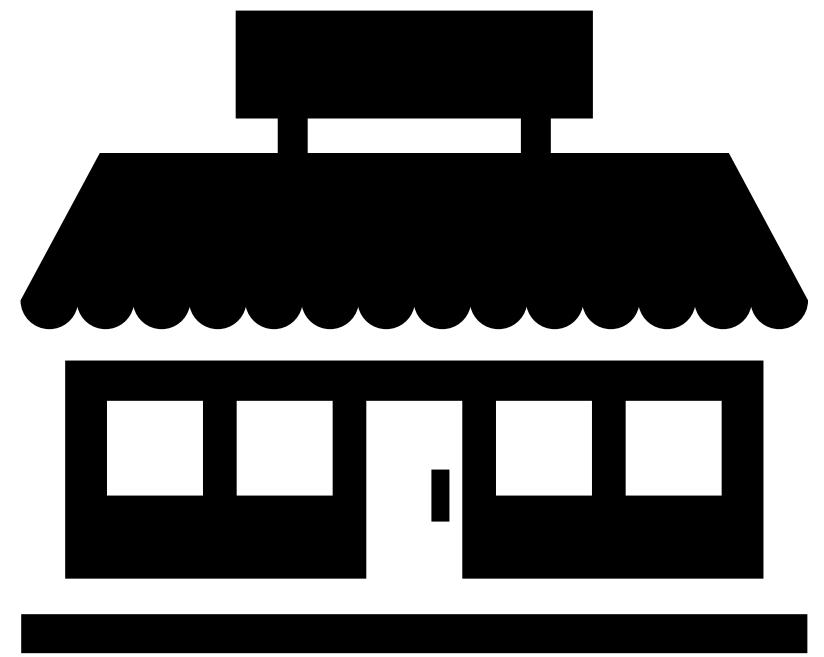
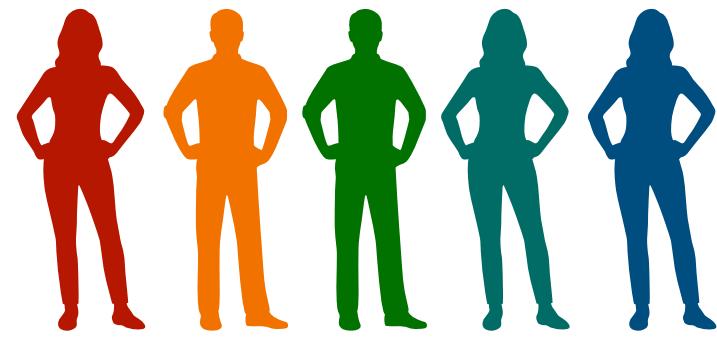
Adam Spannbauer

A little queueing theory

A little queueing theory

Lines build when things arrive faster than we process

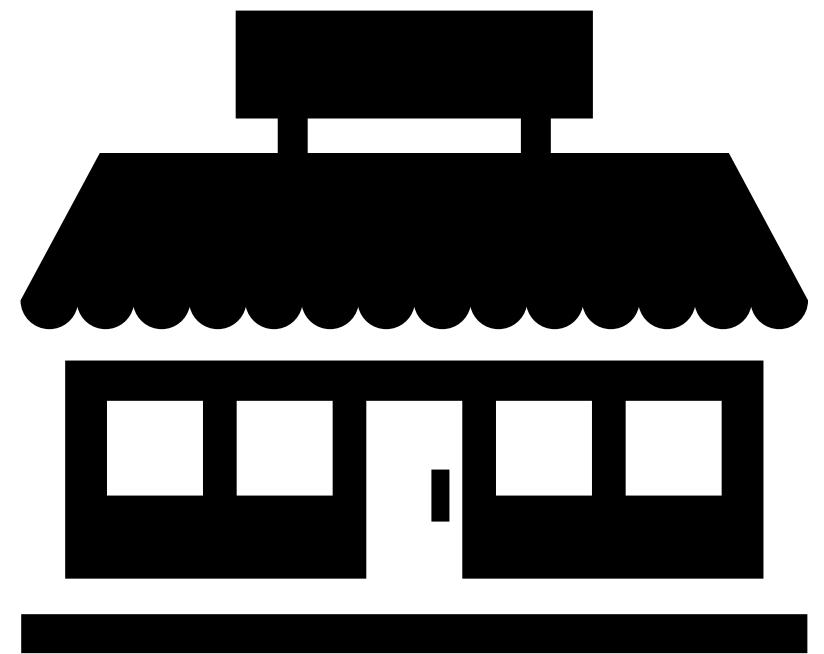
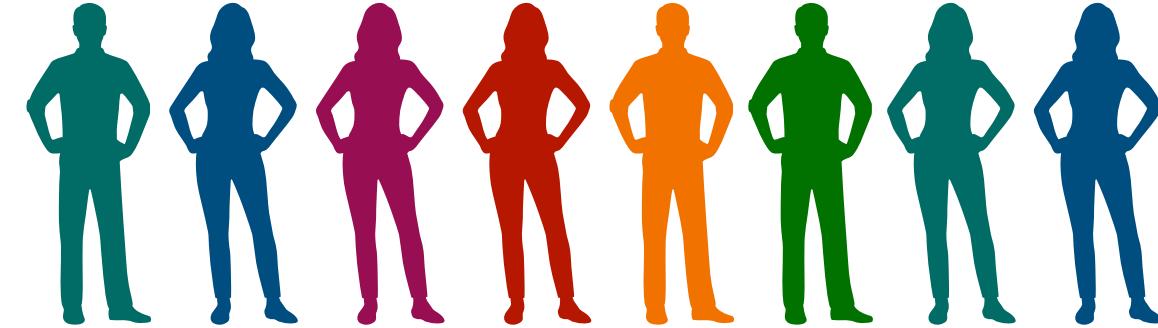
- A queue can be modeled with parameters like



A little queueing theory

Lines build when things arrive faster than we process

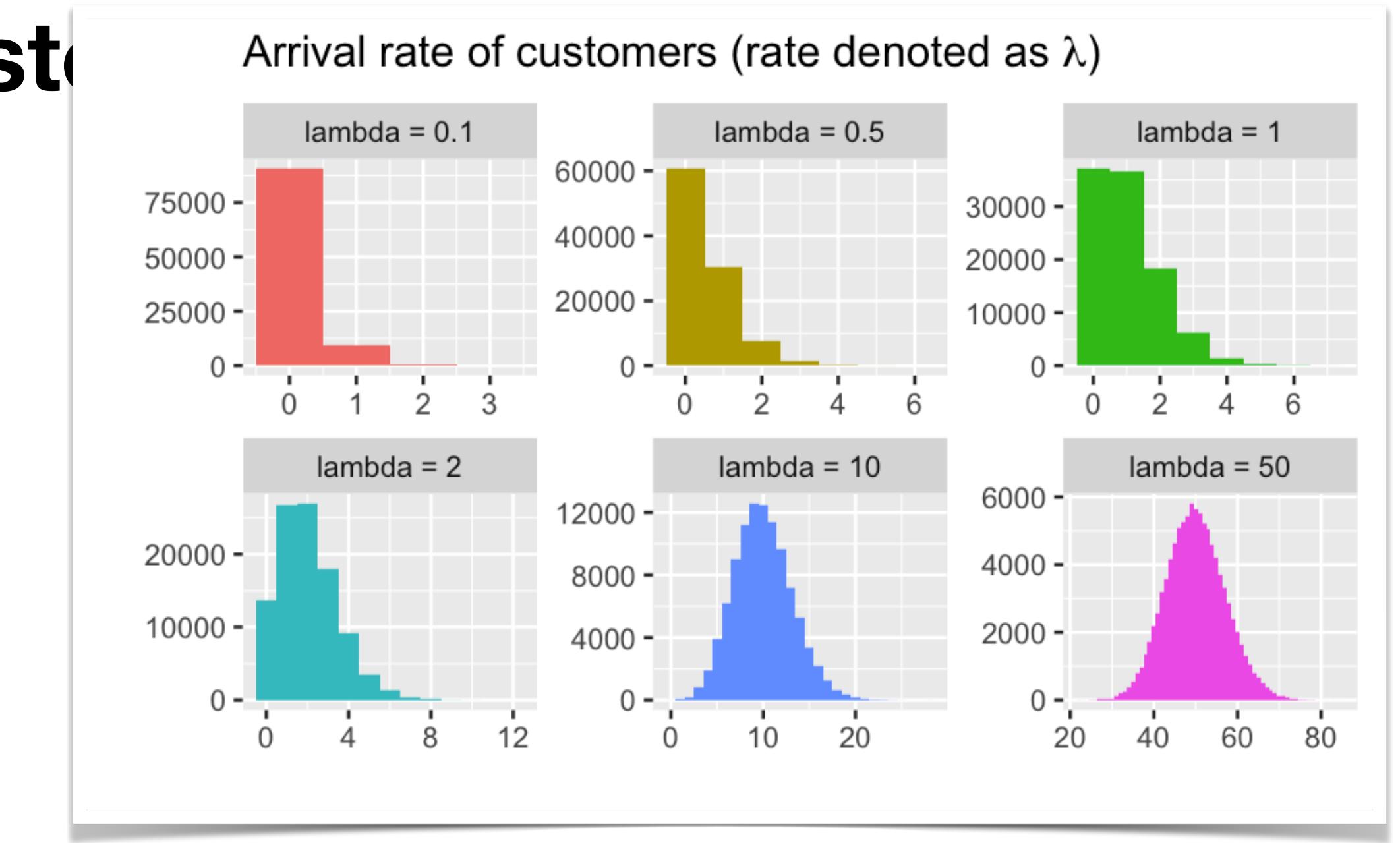
- A queue can be modeled with parameters like
 - Arrival rate



A little queueing theory

Lines build when things arrive faster

- A queue can be modeled with parameters like
 - Arrival rate
 - Ex: can be a Poisson process with arrival rate λ



A little queueing theory

Lines build when things arrive faster than we process

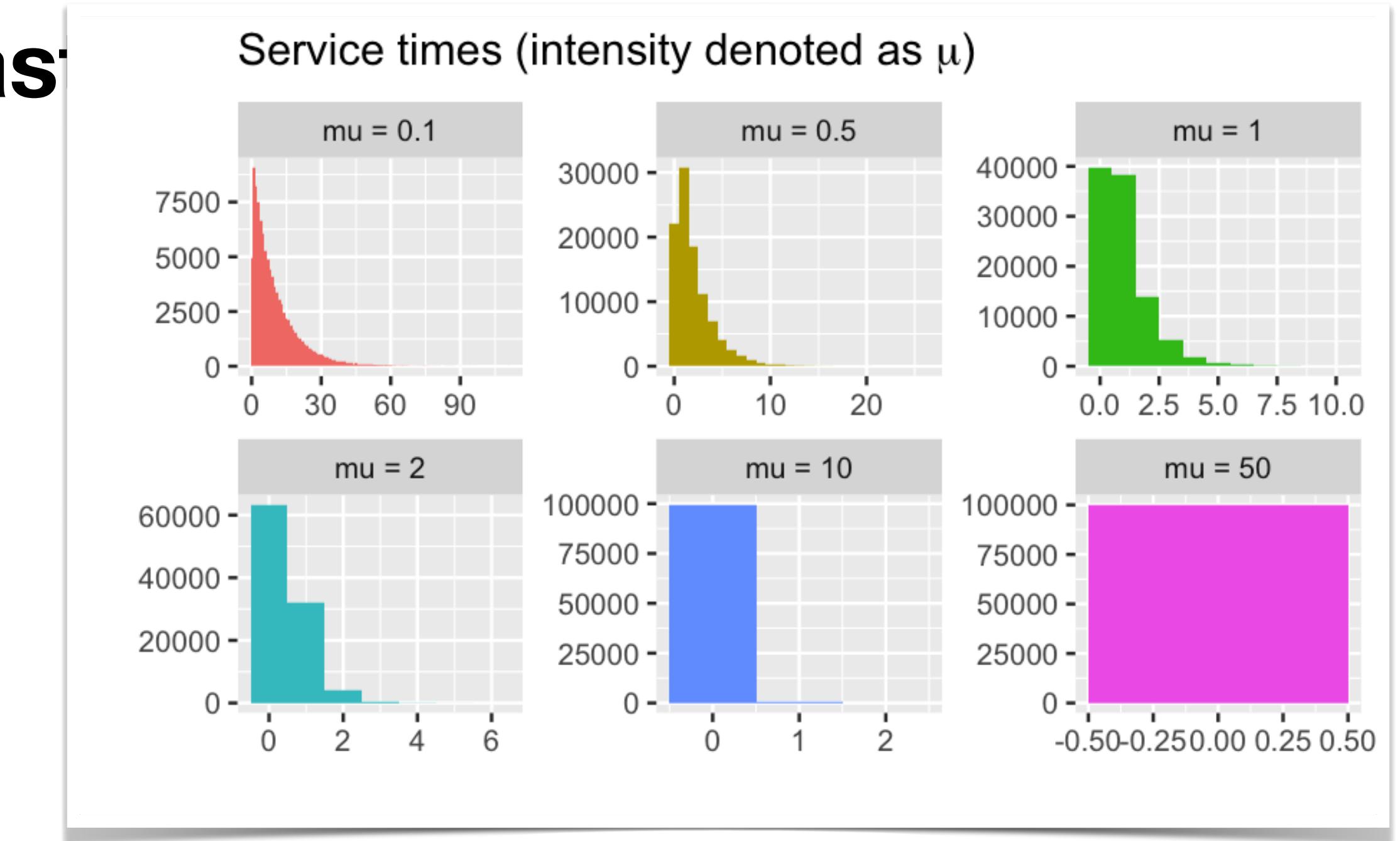
- A queue can be modeled with parameters like
 - Arrival rate
 - Ex: can be a Poisson process with arrival rate λ
 - Service rate



A little queueing theory

Lines build when things arrive fast

- A queue can be modeled with parameters like
 - Arrival rate
 - Ex: can be a Poisson process with arrival rate λ
 - Service rate
 - Ex: can be exponentially distributed with intensity μ



A little queueing theory

Metrics

- Some common ways to measure queues are
 - Expected queue length (L_q)
 - Average waiting time (W)
 - Utilization (ρ)
 - $\rho = \frac{\lambda}{\mu}$ (ratio of arrival rate to service rate)



A little queueing theory

Metrics

- Some common ways to measure queues are
 - Expected queue length (L_q)
 - Average waiting time (W)
 - Utilization (ρ)
 - $\rho = \frac{\lambda}{\mu}$ (ratio of arrival rate to service rate)

Ex:

- 5 customers arrive per hour ($\lambda = 5$)
- We can service 10 customers per hour ($\mu = 10$)
- $\rho = \frac{5}{10} = 0.5$

50% of service capacity is being utilized



A little queueing theory

Metrics

- Some common ways to measure queues are

- Expected queue length (L_q)

- Average waiting time (W)

- Utilization (ρ)

$$\bullet \quad \rho = \frac{\lambda}{\mu} \text{ (ratio of arrival rate to service rate)}$$

We're in a "steady state" when these metrics have stabilized and no longer change over time.
(i.e. viewing things "in the long run")

Our analyses take place with this steady state view.

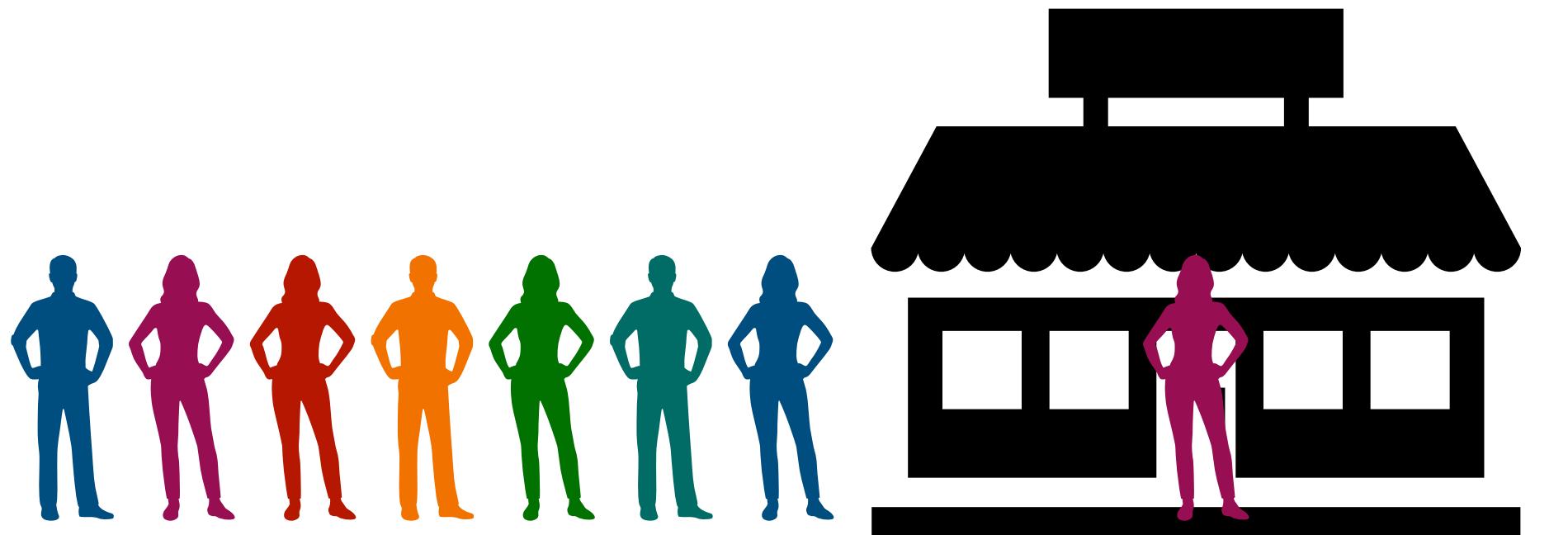
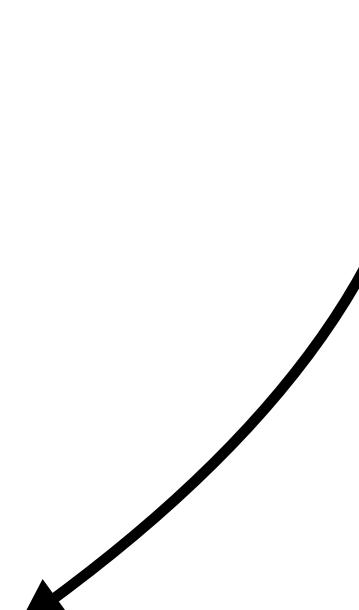


A little queueing theory

Metrics

- Some common ways to measure queues are
 - Expected queue length (L_q)
 - Average waiting time (W)
 - Utilization (ρ)
 - $\rho = \frac{\lambda}{\mu}$ (ratio of arrival rate to service rate)

If $\rho > 1$, the arrival rate is outpacing service. Queues can grow indefinitely*, and the system will not reach a steady state.

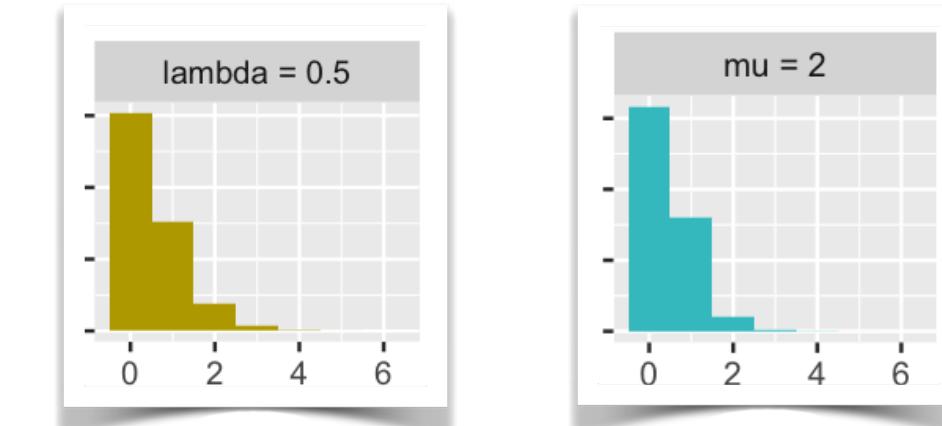


*without mechanisms like “balking”, “finite capacity”, or **tolls**

A little queueing theory

Metrics

- Some common ways to measure queues are
 - Expected queue length (L_q)
 - Average waiting time (W)
 - Utilization (ρ)
 - $\rho = \frac{\lambda}{\mu}$ (ratio of arrival rate to service rate)



“Even when the arrival rate is smaller than the service rate (so that the server can accommodate all arrivals), queues are formed due to the variability in service and inter-arrival times.”



Background to Naor's work

Pay to getting in line?

Ideas cited by Naor

- Leeman proposes tolls to manage queue size
- Saaty argues Leeman's queue toll is unjust unless queueing for a luxury
 - i.e. we shouldn't toll things like healthcare lines



Leeman, W. A. (1964). The reduction of queues through the use of price. *Operations Research*, 12(5), 783-785. <https://pubsonline.informs.org/doi/pdf/10.1287/opre.12.5.783>

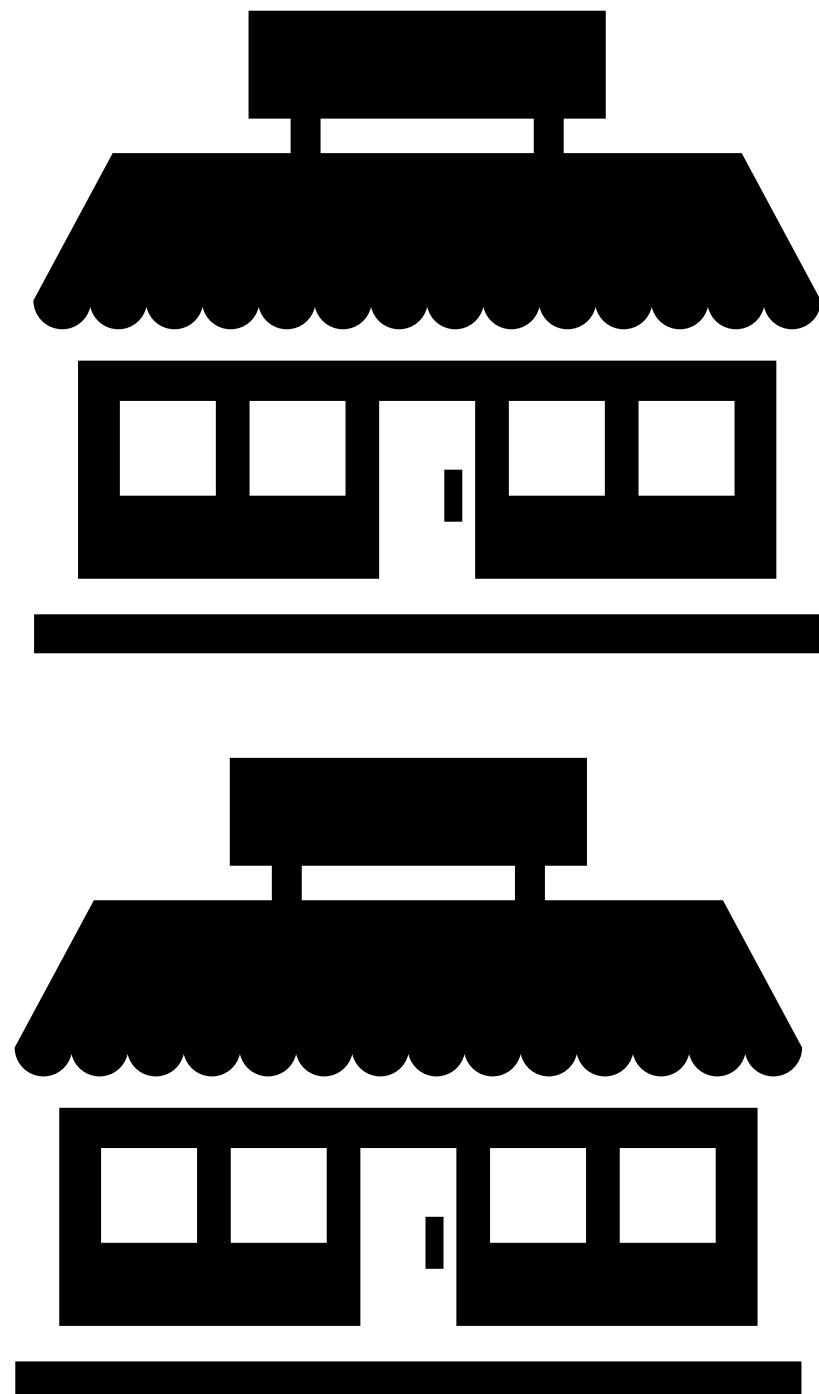
Saaty, T. L., & Leeman, W. A. (1965). The Burdens of Queuing Charges-Comments on a Letter by Leeman. *Operations Research*, 13(4), 679–681. <http://www.jstor.org/stable/167860>

Pay to getting in line?

Ideas cited by Naor

- Leeman proposes tolls to manage queue size
- Saaty argues Leeman's queue toll is unjust unless queueing for a luxury
 - i.e. we shouldn't toll things like healthcare lines

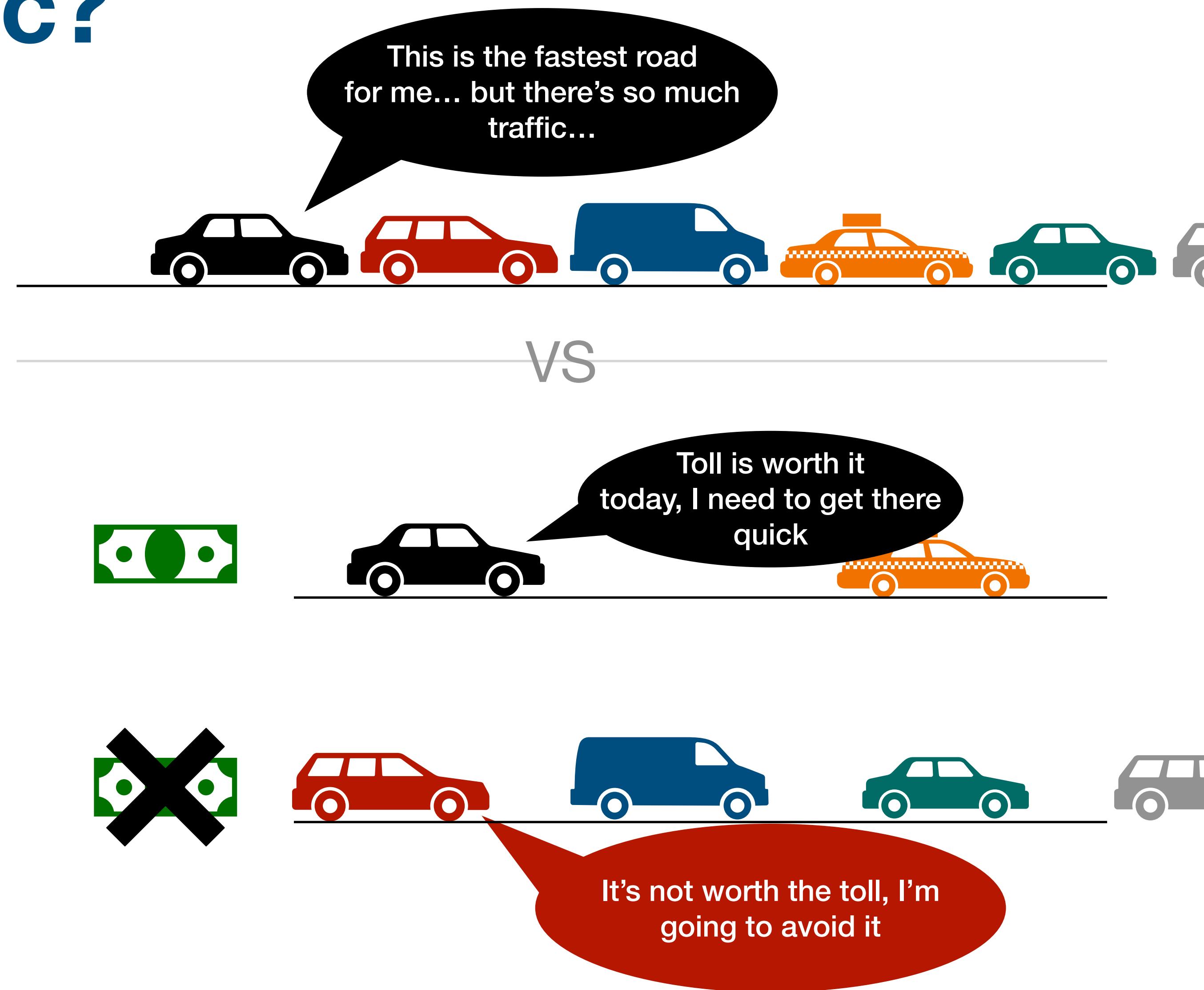
Also suggests dynamic pricing to discourage queueing in busy times and stimulate demand in less busy times



Queues are like traffic?

Ideas cited by Naor

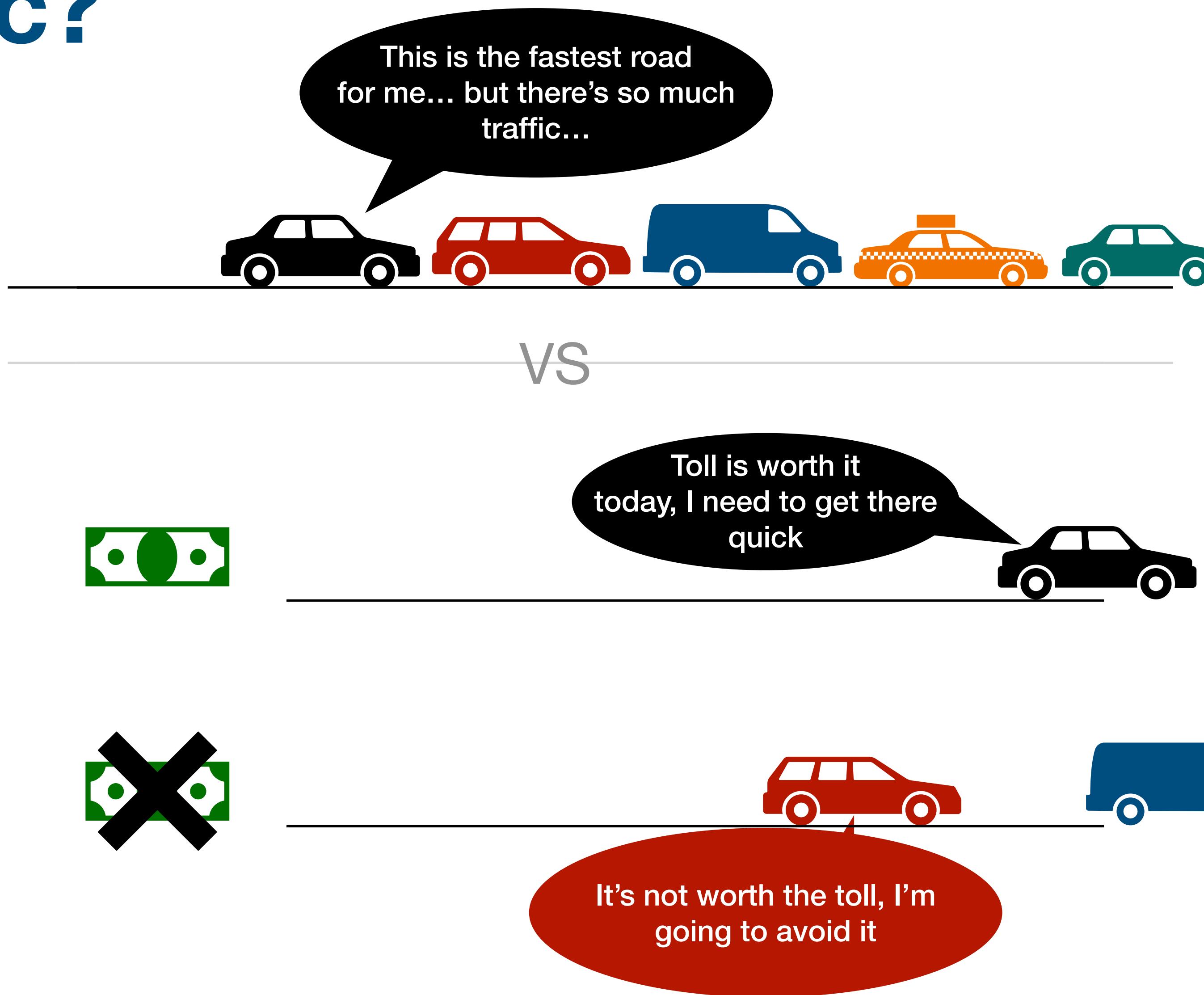
- Drivers focus on optimizing their own routes, this does not optimize the overall system
- Tolls on roads can help redistribute traffic?



Queues are like traffic?

Ideas cited by Naor

- Drivers focus on optimizing their own routes, this does not optimize the overall system
- Tolls on roads can help redistribute traffic?
 - We'd try to set tolls in a way that aligns self-interested utility with overall utility



Big impact of paper

Customers act in their own self interest

We can model this with $v_{service} = \frac{R\mu}{C}$ (reward of service per unit of cost)

- “Naor (1969) appears to be the first to incorporate customer decisions into a queueing model.”
- Naor outlines a framework for addressing *Rewards and Costs* for consumers



Customers act in their own self interest

We can model this with $v_{service} = \frac{R\mu}{C}$ (reward of service per unit of cost)

- “Naor (1969) appears to be the first to incorporate customer decisions into a queueing model.”
- Naor outlines a framework for addressing *Rewards and Costs* for consumers



Customers act in their own self interest

We can model this with $v_{service} = \frac{R\mu}{C}$ (reward of service per unit of cost)

- “Naor (1969) appears to be the first to incorporate customer decisions into a queueing model.”
- Naor outlines a framework for addressing *Rewards and Costs* for consumers

$$R = \$10$$

$$C = \$8 \text{ per minute}$$

$$\mu = 5 \text{ minutes to service per person}$$

$$n = 10 \text{ people to go including me}$$



Customers act in their own self interest

We can model this with $v_{service} = \frac{R\mu}{C}$ (reward of service per unit of cost)

- “Naor (1969) appears to be the first to incorporate customer decisions into a queueing model.”
- Naor outlines a framework for addressing *Rewards and Costs for consumers*

$$R = \$10$$

$$C = \$8 \text{ per minute}$$

$$\mu = 5 \text{ minutes to service per person}$$

$$n = 10 \text{ people to go including me}$$

$$v_s(n) = R - \frac{Cn}{\mu}$$

(utility for customer given n people wait)



Customers act in their own self interest

We can model this with $v_{service} = \frac{R\mu}{C}$ (reward of service per unit of cost)

- “Naor (1969) appears to be the first to incorporate customer decisions into a queueing model.”
- Naor outlines a framework for addressing *Rewards and Costs for consumers*
 - Sometimes the juice isn’t worth the squeeze and they won’t queue

$$R = \$10$$

$$C = \$8 \text{ per minute}$$

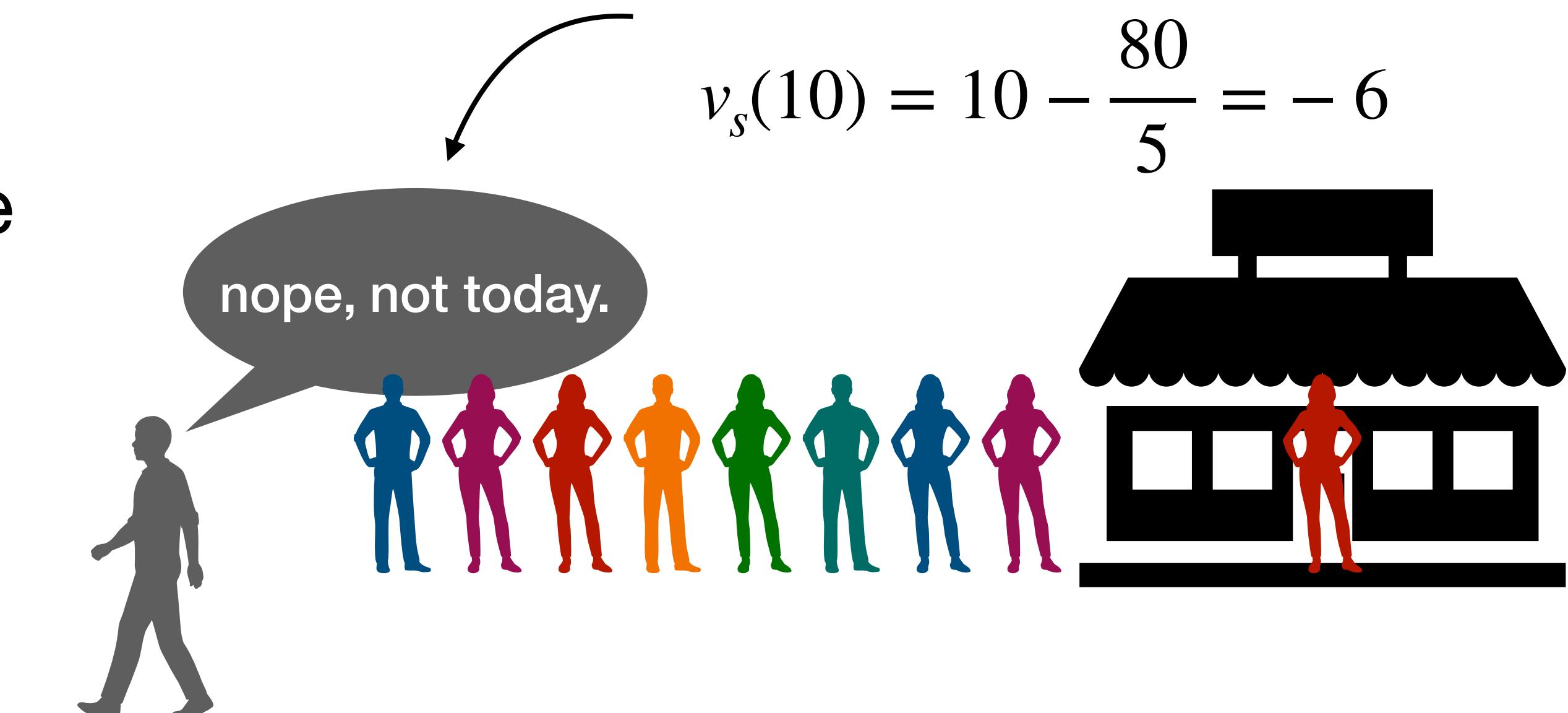
$$\mu = 5 \text{ minutes to service per person}$$

$$n = 10 \text{ people to go including me}$$

$$v_s(n) = R - \frac{Cn}{\mu}$$

$$(utility for customer given 10 people wait)$$

$$v_s(10) = 10 - \frac{80}{5} = -6$$



Customers act in their own self interest

We can model this with $v_{service} = \frac{R\mu}{C}$ (reward of service per unit of cost)

- “Naor (1969) appears to be the first to incorporate customer decisions into a queueing model.”
- Naor outlines a framework for addressing *Rewards and Costs for consumers*
 - Sometimes the juice isn’t worth the squeeze and they won’t queue

$$R = \$10$$

$$C = \$8 \text{ per minute}$$

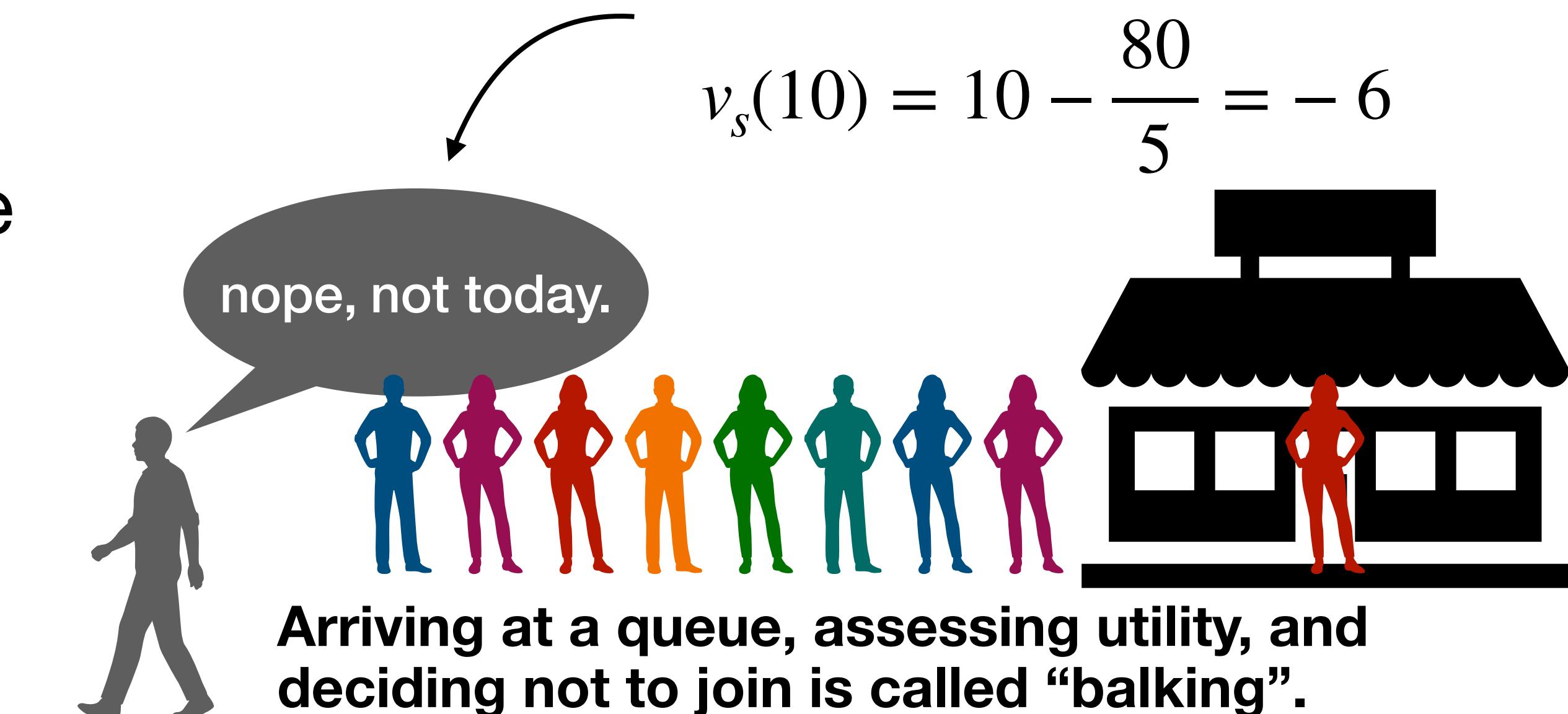
$$\mu = 5 \text{ minutes to service per person}$$

$$n = 10 \text{ people to go including me}$$

$$v_s(n) = R - \frac{Cn}{\mu}$$

$$(utility for customer given 10 people wait)$$

$$v_s(10) = 10 - \frac{80}{5} = -6$$



Customers act in their own self interest

We can model this with $v_{service} = \frac{R\mu}{C}$ (reward of service per unit of cost)

- “Naor (1969) appears to be the first to incorporate customer decisions into a queueing model.”
- Naor outlines a framework for addressing *Rewards and Costs for consumers*
 - Sometimes the juice isn’t worth the squeeze and they won’t queue

$$R = \$10$$

$$C = \$8 \text{ per minute}$$

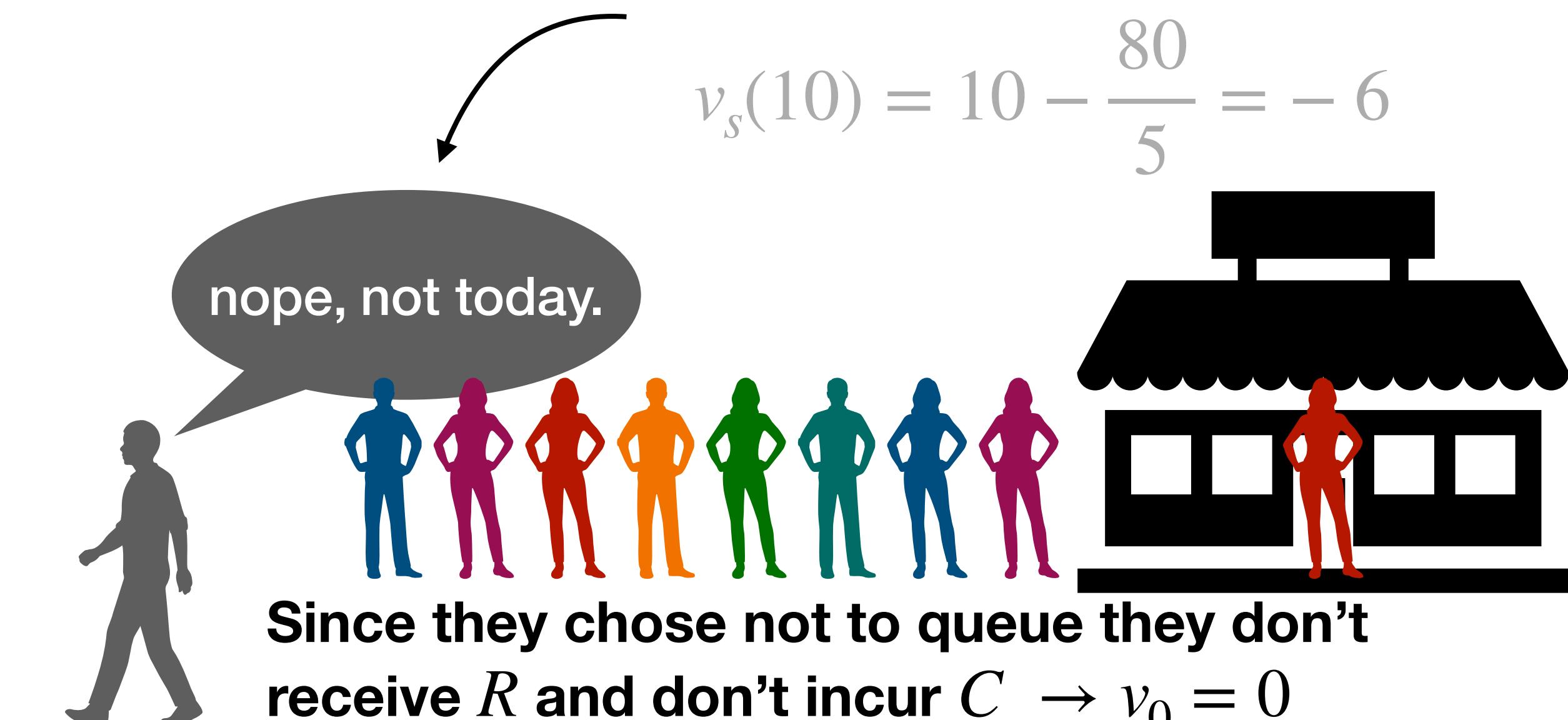
$$\mu = 5 \text{ minutes to service per person}$$

$$n = 10 \text{ people to go including me}$$

$$v_s(n) = R - \frac{Cn}{\mu}$$

$$(utility for customer given 10 people wait)$$

$$v_s(10) = 10 - \frac{80}{5} = -6$$



Naor's model in more depth

Problem statement

Can queue tolls be used to optimize the system for all?

- Excessive queue lengths cause inefficiencies
 - e.g. wasted time & reduced customer satisfaction)
- Proposal: Charge customers to enter the queue
 - just like auto traffic this can reroute individuals for the overall good of the system



Conditions for model

$R, C, \lambda, \mu, \rho, v_s, p_i, g(z), E[i], \zeta, b$



1. There exists a public good (aka well-being) that can be maximized via some objective function.
2. “Customers are liable to be diverted from the service station”

Conditions for model

$R, C, \lambda, \mu, \rho, v_s, p_i, g(z), E[i], \zeta, b$



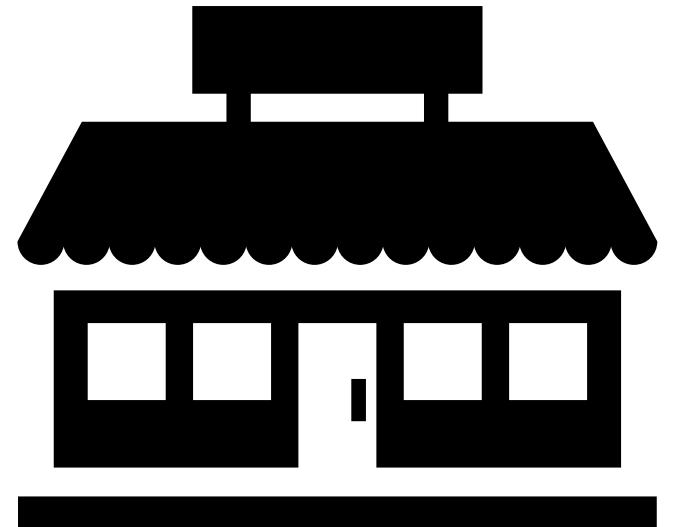
1. There exists a public good (aka well-being) that can be maximized via some objective function.

- “...the expected overall profit (in unit time) accruing to arriving customers is a proper objective function representing public good...”
- This can have a centralized decision maker (public transportation, govt healthcare, etc) or decentralized decision makers (customers act for themselves to max their own utility)

2. “Customers are liable to be diverted from the service station”

Conditions for model

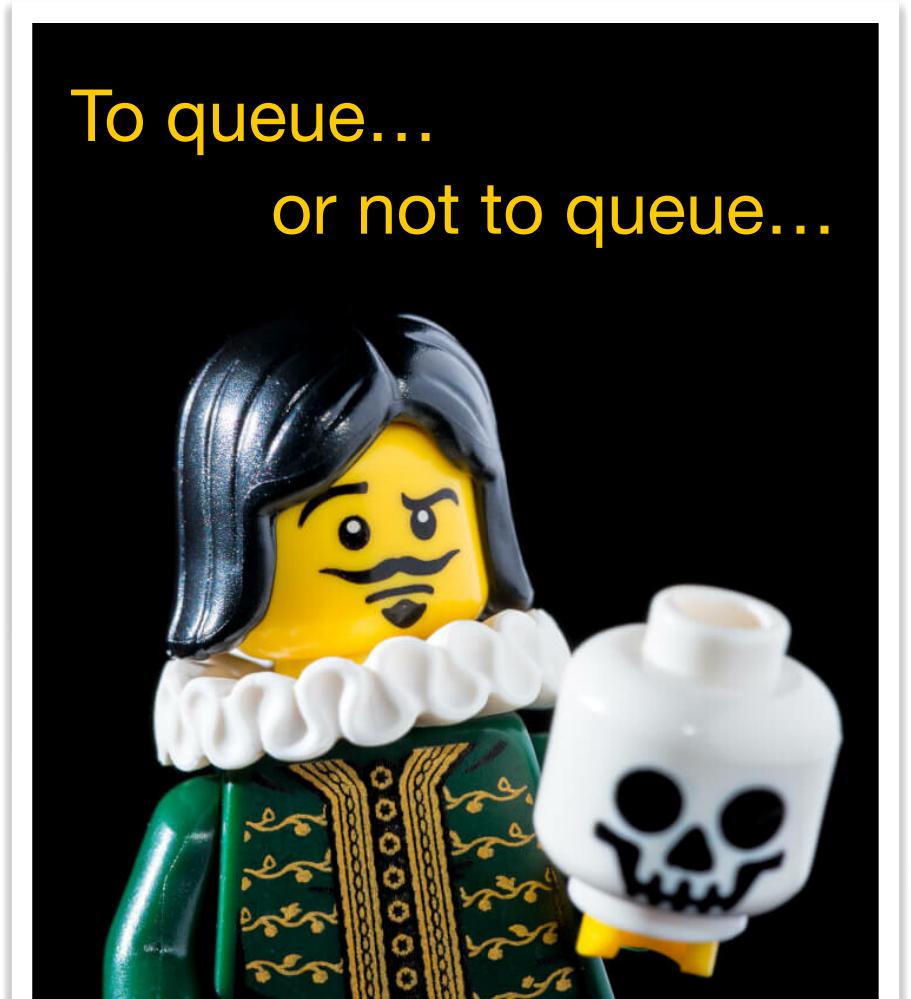
$R, C, \lambda, \mu, \rho, v_s, p_i, g(z), E[i], \zeta, b$



1. There exists a public good (aka well-being) that can be maximized via some objective function.

2. **“Customers are liable to be diverted from the service station”**

- Big contrast to prior queueing models that assumed we serve *all customers*
- Some will queue and some will not, those that don't have utility 0
- A toll can influence the decision of queuing vs leaving



<https://www.istockphoto.com/photo/lego-minifigures-series-8-figurine-the-thespian-gm458926969-22139332>

Defining the model

$R, C, \lambda, \mu, \rho, v_s, p_i, g(z), E[i], \zeta, b$



- Inputs
 - R - reward for being serviced (e.g. \$)
 - C - cost per unit of time for queueing (e.g. \$ per minute; assumed all customer C s are equal)
 - λ - arrival rate of customers (e.g. Poisson distributed)
 - μ - service rate (e.g. exponential distributed)

Defining the model

$R, C, \lambda, \mu, \rho, \nu_s, p_i, g(z), E[i], \zeta, b$



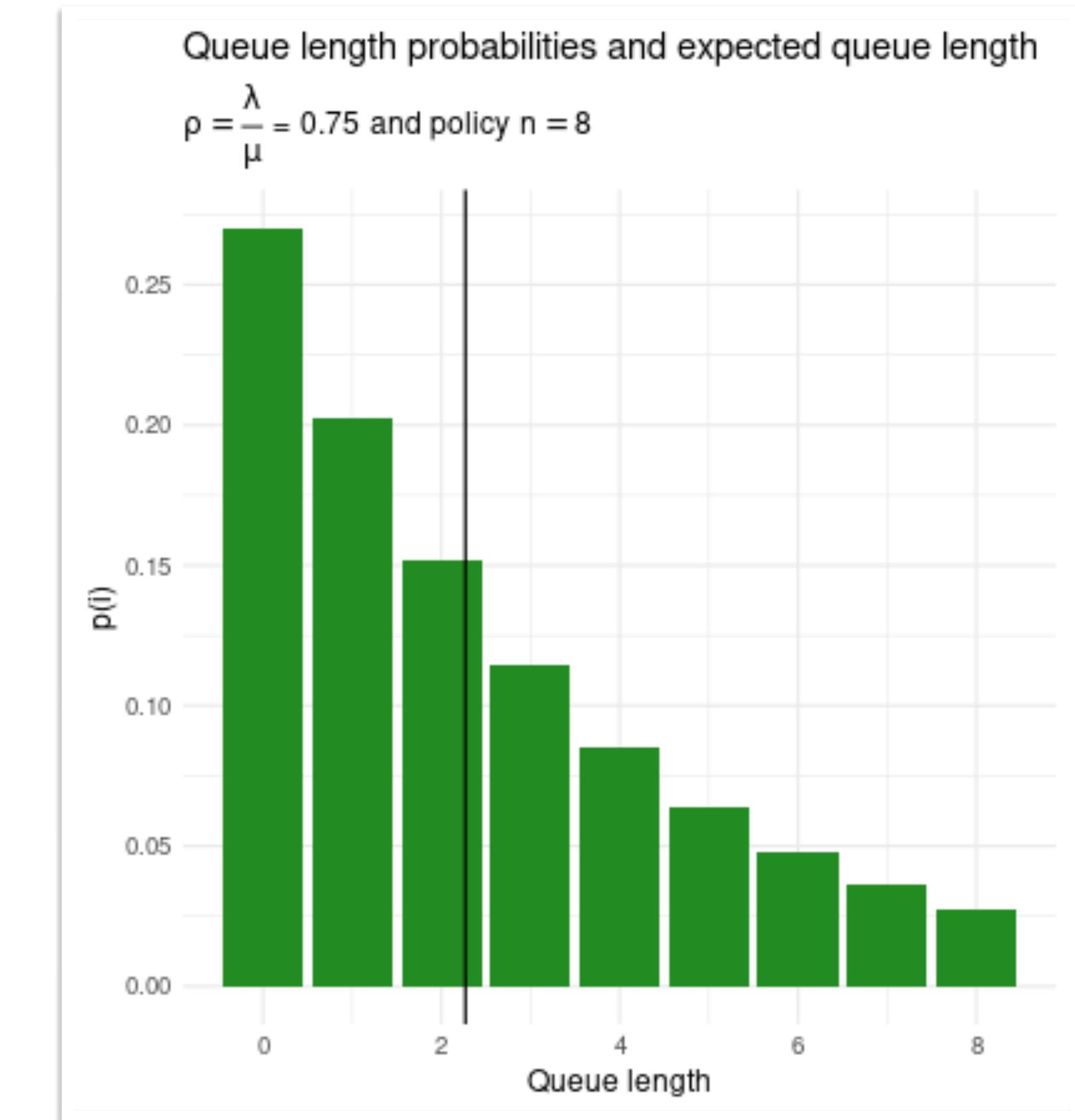
- Additional variables
 - $\rho = \frac{\lambda}{\mu}$ - classic utilization factor (ratio of arrival rate to service rate)
 - Note! A different metric will be used to calculate service's "busy fraction" to account for queue length of 0
 - i - queue size at a given time (randomly distributed)
 - n - queue capacity (if $i > n$ a customer will not queue; a customer will queue if $i \leq n$)

Steady state equations

$R, C, \lambda, \mu, \rho, v_s, p_i, g(z), E[i], \zeta, b$



- PMF of queue length: $p_{i+1} = p_i \rho \rightarrow p_i = \frac{\rho^i(1 - \rho)}{1 - \rho^{n+1}}$
 - The probability of line being length i as a function of ρ
 $(\frac{\lambda}{\mu}$ - utilization factor) and n (queue capacity)



Kendall, D. G. (1953). Stochastic processes occurring in the theory of queues and their analysis by the method of the imbedded Markov chain. *The Annals of Mathematical Statistics*, 338-354.

Naor, P. (1969). The regulation of queue size by levying tolls. *Econometrica: journal of the Econometric Society*, 15-24. <https://doi.org/10.2307/1909200>

Steady state equations

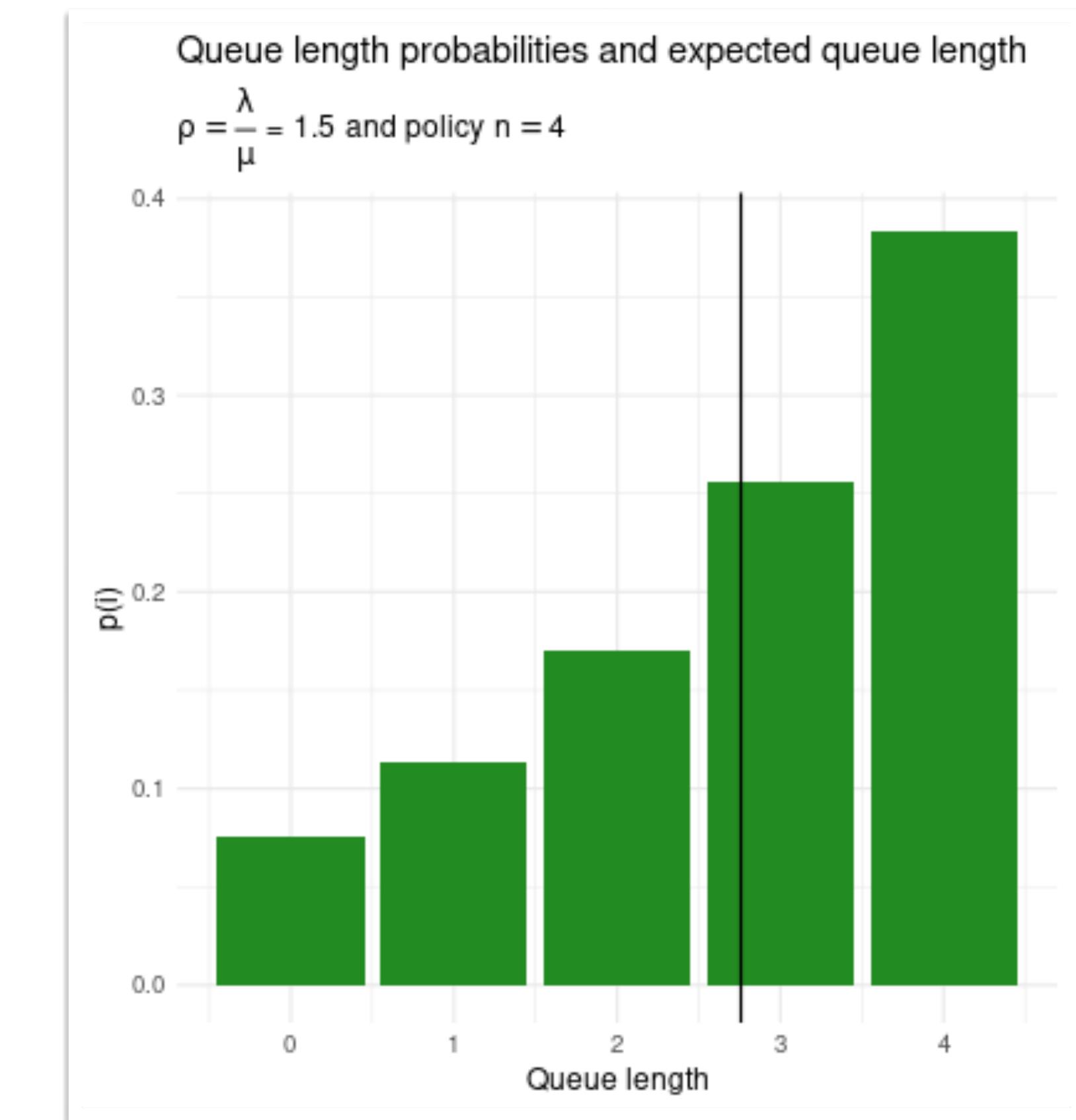
$R, C, \lambda, \mu, \rho, v_s, p_i, g(z), E[i], \zeta, b$



- PMF of queue length: $p_{i+1} = p_i \rho \rightarrow p_i = \frac{\rho^i(1 - \rho)}{1 - \rho^{n+1}}$
 - The probability of line being length i as a function of ρ ($\frac{\lambda}{\mu}$ - utilization factor) and n (queue capacity)
- Probability Generating Function:

$$g(z) = \sum_{i=0}^n p_i z^i = \frac{1 - \rho}{1 - \rho^{n+1}} \cdot \frac{1 - (\rho z)^{n+1}}{1 - \rho z}$$
- Expected length of queue:

$$E[i] = \frac{\rho[1 - (n + 1)\rho^n + n\rho^{n+1}]}{(1 - \rho)(1 - \rho^{n+1})} = \frac{\rho}{1 - \rho} - \frac{(n + 1)\rho^{n+1}}{1 - \rho^{n+1}}$$

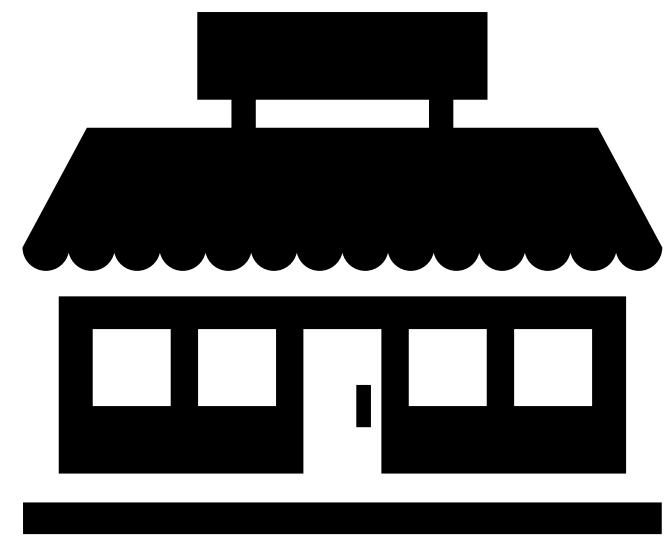


Kendall, D. G. (1953). Stochastic processes occurring in the theory of queues and their analysis by the method of the imbedded Markov chain. *The Annals of Mathematical Statistics*, 338-354.

Naor, P. (1969). The regulation of queue size by levying tolls. *Econometrica: journal of the Econometric Society*, 15-24. <https://doi.org/10.2307/1909200>

Steady state equations

$R, C, \lambda, \mu, \rho, \nu_s, p_i, g(z), E[i], \zeta, b$



- **Balking rate** (arrival rate x probability of being at capacity)

$$\bullet \quad \zeta = \lambda p_n = \frac{\lambda \rho^n (1 - \rho)}{1 - \rho^{n+1}}$$

- **Effective arrival rate** (arrival - rejection): $\lambda - \zeta = \lambda(1 - p_n)$

- **“Busy fraction”** (exclude case of no one in system)

$$\bullet \quad b = \sum_{i=1}^n p_i = 1 - p_0 = \frac{\rho(1 - \rho^n)}{1 - \rho^{n+1}}$$

- **Effective service rate** (service * busy fraction): $\mu b = \mu(1 - p_0)$

The big contribution

Modeling customers' net gain

“The newly arrived customer weighs the two alternatives - to join or not to join the queue - by the net gains associated with them.”



The big contribution

Modeling customers' net gain

“The newly arrived customer weighs the two alternatives - to join or **not to join** the queue - by the net gains associated with them.”

not to join: $G_i = 0$



The big contribution

Modeling customers' net gain

“The newly arrived customer weighs the two alternatives - **to join** or not to join the queue - by the net gains associated with them.”

$$\text{not to join: } G_i = 0$$

$$\text{to join: } G_i = R - (i + 1)C \frac{1}{\mu}$$



The big contribution

Modeling customers' net gain

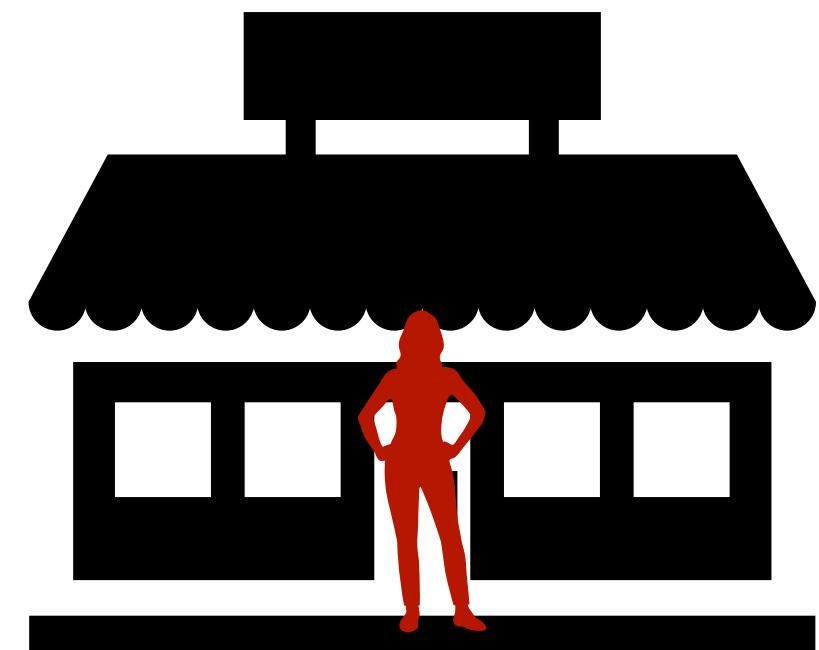
“The newly arrived customer weighs the two alternatives - to join or not to join the queue - by the net gains associated with them.”

not to join: $G_i = 0$

to join: $G_i = R - \frac{1}{\mu} (i + 1) C$

Reward for being
serviced

Cost of waiting for
everyone's service
(including the one joining)



“Self-optimization” leads to lines of length n_s

Customers optimize for their own self-interest and the system suffers

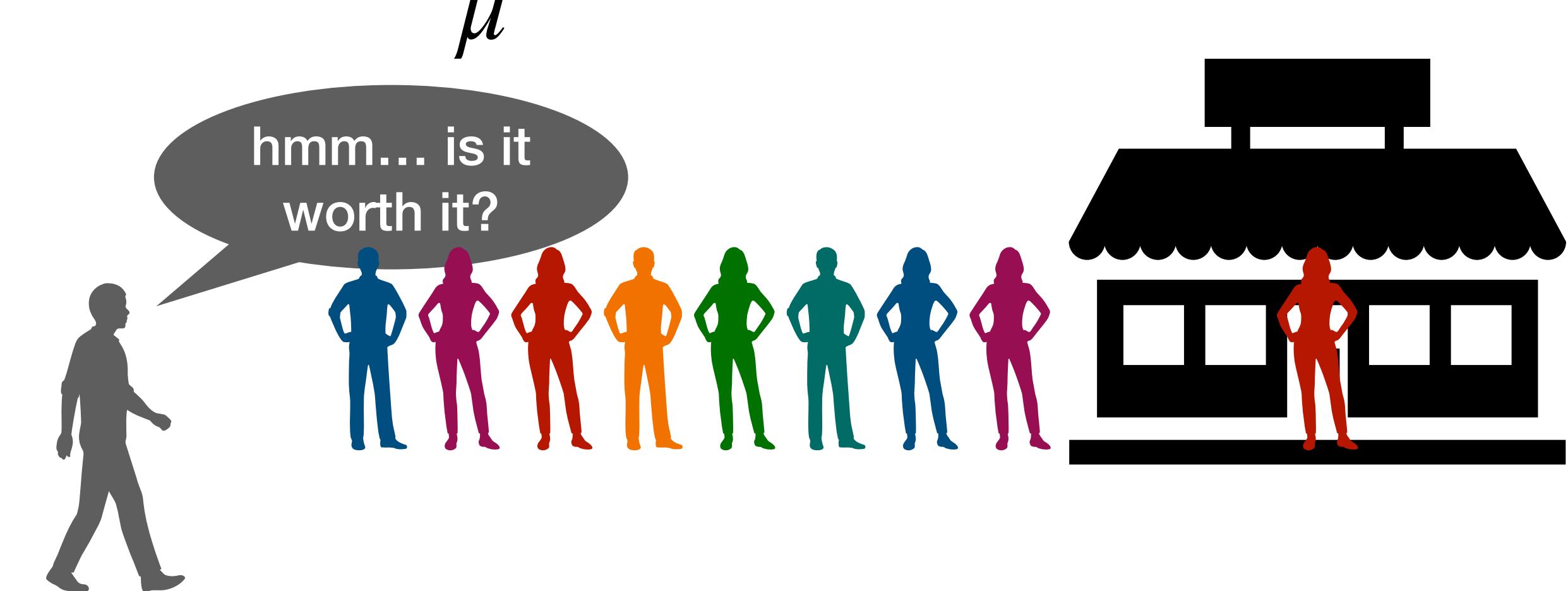
“The newly arrived customer weighs the two alternatives - to join or not to join the queue - by the net gains associated with them.”

not to join: $G_i = 0$

to join: $G_i = R - (i + 1)C \frac{1}{\mu}$

Notice this gain (G_i) does not consider the arrival rate (λ)!

If left to self-optimize, just like traffic, customers will not optimize the system overall.



“Overall optimization” leads to lines of length n_o

Maximize expected sum of customer net gains

“We note that the expected total net gain, P , under some strategy n is given by”

$$P = (\lambda - \zeta)R - CE[i]$$

Reward multiplied by expected number of customers joining queue.

This assumes once customers join the queue, they will not leave (i.e. no “reneging”)

Cost multiplied by expected number of customers in queue



“Beneficial Toll Imposition” to lower n_s to n_o

Lowering self utility via toll can raise overall utility

“a toll θ is imposed on customers joining the queue and their (individually) expected net gain is reduced in such a way that n_o is the current criterion of newly arrived customers”



“Beneficial Toll Imposition” to lower n_s to n_o

Lowering self utility via toll can raise overall utility

“a toll θ is imposed on customers joining the queue and their (individually) expected net gain is reduced in such a way that n_o is the current criterion of newly arrived customers”

$$R - \frac{C(n_o + 1)}{\mu} < \theta^* \leq R - \frac{C(n_o)}{\mu}$$

“If a toll taken from this range is levied on customers joining the queue, the combined income (in unit time) of customers and the revenue agency is maximized.”



“Revenue Maximization”

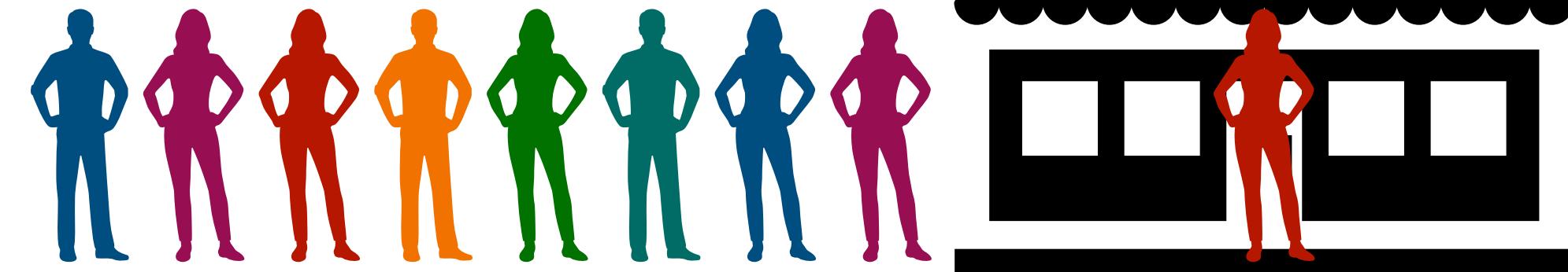
Tolls to zero out customers gain and maximize revenue

“The toll collecting agency may be completely divorced from the individual and collective economic interests of the customers. In that case the agency will seek to impose a toll, θ_r , designed to maximize its own revenue rather than to optimize the whole system.”

$$M = (\lambda - \zeta) \theta$$

Expected number of
customers joining queue.

Toll value



“Revenue Maximization” to lower n_s to n_o

Tolls to zero out customers gain and maximize revenue

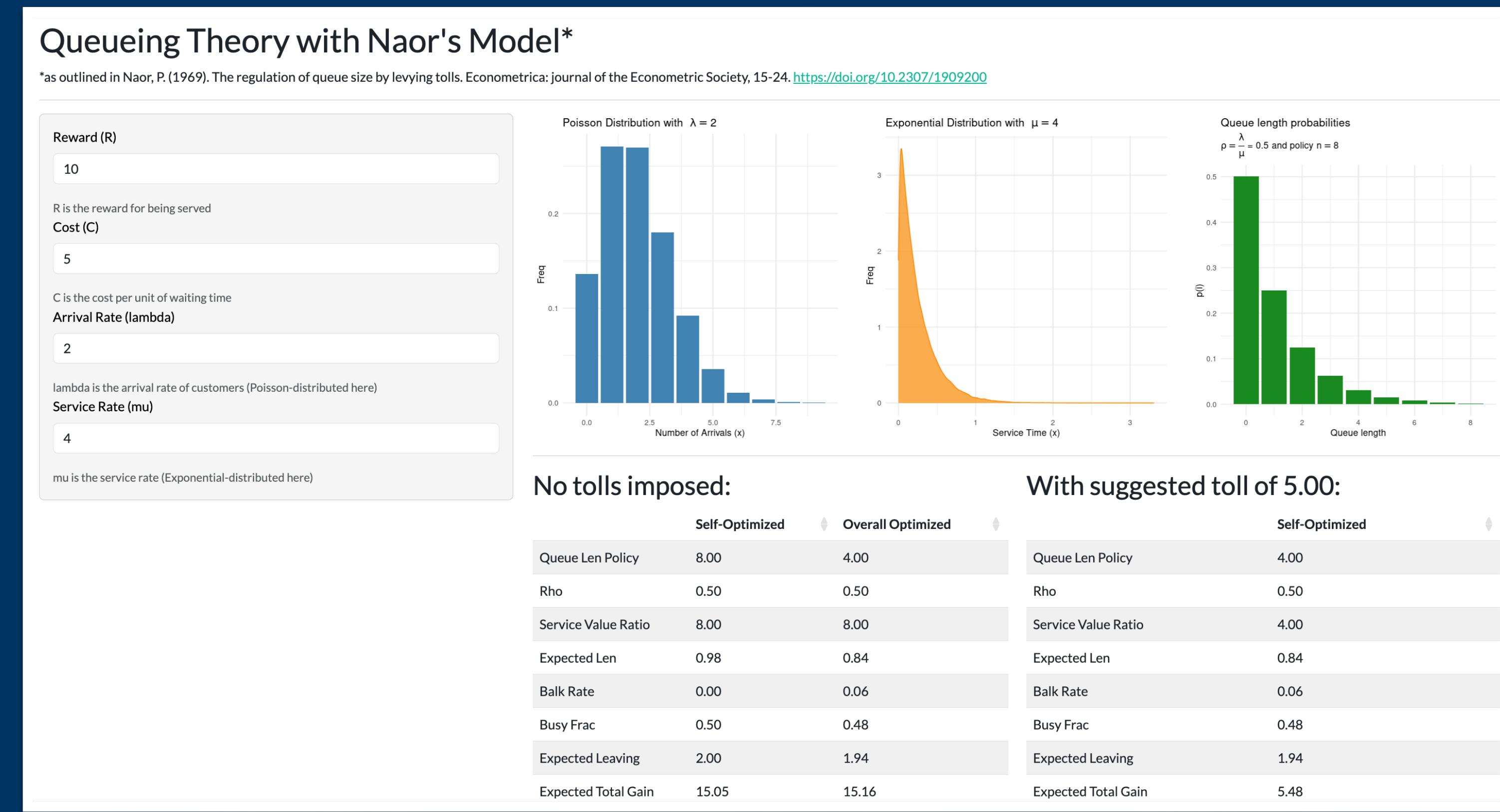
Set the toll to be equal to the customers' profit from waiting in line and being serviced

$$\theta_r = R - \frac{Cn_r}{\mu}$$

Relies on a net neutral experience being good enough to retain customers



Examples via hands-on demo



<https://spannbaueradam.shinyapps.io/naor/>

Critiques

Adam thoughts

Before consulting the literature

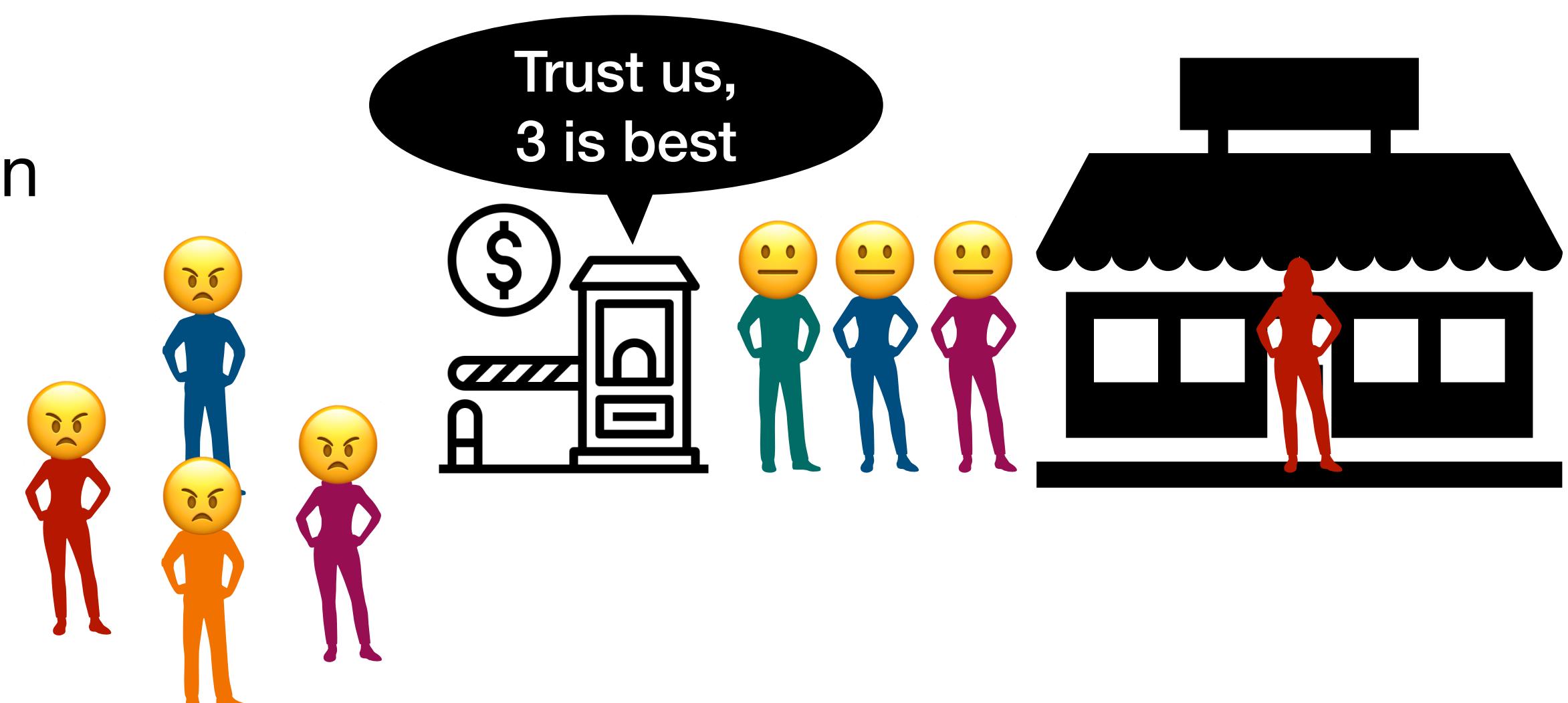
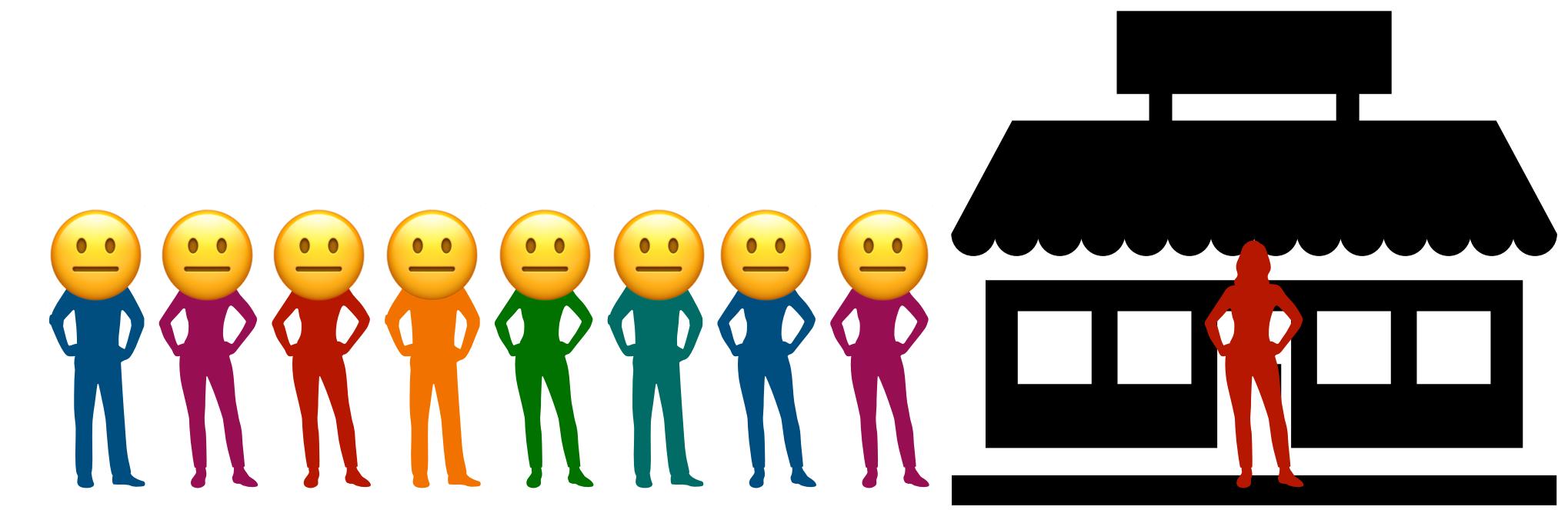
- Doesn't seem to consider how the toll affects an overall gain, just the queue size
 - We start optimizing overall gain:
 $P = (\lambda - \zeta)R - CE[i]$, but after setting a toll to bring down queue size we don't readdress how this affects P
 - P with the toll incorporated could be expressed with:
 $P = (\lambda - \zeta) \cdot (R - \theta) - CE[i]$



Adam thoughts

Before consulting the literature

- Doesn't seem to consider how the toll affects an overall gain, just the queue size
 - This could be considered net neutral because we're optimizing for firm and customers?
 - “if the toll revenue may be used for redistribution of income among the population or for socially useful purposes the proposed imposition of tolls is an optimal procedure”
 - The “if” might hint at a tradeoff in firm’s favor



Adam thoughts

Before consulting the literature

- When considering customer utility we don't consider a cost for arriving early
- One of the cited works brings up the point:
 - “Some have argued that the man ahead in a queue imposes a cost on those who wait behind him and must pay for it. But I say that he has already paid for it by coming earlier.”
- But, to be fair, this seems like a useful simplification for most applications where we focus on steady state



<https://www.jacksonville.com/story/news/2009/11/26/stub-505/15965432007/>

Critiques

Some other considerations

- Are customer's rewards and costs really all the same?
- Saaty's point against tolls on non-luxuries still stands - can be viewed as unjust to toll in some cases
- Relies on assumptions that customer's can gauge the cost-benefit accurately



Hassin, R., & Haviv, M. (2003). *To queue or not to queue: Equilibrium behavior in queueing systems* (Vol. 59). Springer Science & Business Media.

Huang, T., Allon, G., & Bassamboo, A. (2013). Bounded rationality in service systems. *Manufacturing & Service Operations Management*, 15(2), 263-279.

Saaty, T. L., & Leeman, W. A. (1965). The Burdens of Queueing Charges-Comments on a Letter by Leeman. *Operations Research*, 13(4), 679–681. <http://www.jstor.org/stable/167860>

Naor, P. (1969). The regulation of queue size by levying tolls. *Econometrica: journal of the Econometric Society*, 15-24. <https://doi.org/10.2307/1909200>

Implementations and extensions

Real life implementations

Results, not just ideas

- Dynamic pricing in ride share
 - You observe cost before entering the queue to find a driver
 - Uber/Lyft/etc can affect queues with this dynamic pricing like Naor's toll
- Fast track services in amusement parks
 - Pay toll to enter a special fast track queue



Extensions

Provided by a familiar name

Bounded Rationality in Service Systems

Tingliang Huang

Department of Management Science and Innovation, University College London,
London WC1E 6BT, United Kingdom, t.huang@ucl.ac.uk

Gad Allon, Achal Bassamboo

Kellogg School of Management, Northwestern University, Evanston, Illinois 60208
{g-allon@kellogg.northwestern.edu, a-bassamboo@kellogg.northwestern.edu}

- “Yechiali (1971, 1972) extends Naor’s model to allow for GI/M/1 queues. Knudsen (1972) extends Naor’s model to allow for a multiserver queueing system in which arriving customers’ net benefits are heterogeneous. Lippman and Stidham (1977) extend the Naor model to the finite-horizon and discounted cases, showing that, in these settings, the economic notion of an external effect has a precise quantitative interpretation. Hassin (1986) considers a revenue maximizing server who has the opportunity to suppress information on actual queue length, leaving customers to decide whether to join the queue on the basis of the known distribution of waiting times. See Van Mieghem (2000), Hassin and Haviv (2003), Afèche (2004), and Hsu et al. (2009) for other extensions and a comprehensive literature review”

Extensions

Provided by a familiar name

- “Yechiali (1971, 1972) extends Naor’s model to allow for a finite number of servers. Knudsen (1972) extends Naor’s model to allow for a system in which arriving customers’ net benefits are random variables. Hassin and Stidham (1977) extend the Naor model to the finite case and consider discounted cases, showing that, in these settings, the external effect has a precise quantitative interpretation. Hassin (1980) considers a revenue maximizing server who has the information on actual queue length, leaving customers to decide whether to join the queue on the basis of the known distribution of the queue length. Van Mieghem (2000), **Hassin and Haviv (2003)**, and **Hassin et al. (2009)** provide surveys (see also **Hassin and Haviv (2009)** for other extensions and a comprehensive literature review).

TO QUEUE OR NOT
TO QUEUE

Equilibrium Behavior in
Queueing Systems

Refael Hassin
Moshe Haviv



Kluwer's INTERNATIONAL SERIES

Qs?