# Midlands Decision Support Network
# Midlands Analyst Network - 21 March 2024

## The 85% bed occupancy fallacy

Dr Nathan Proudlove
nathan.proudlove@manchester.ac.uk
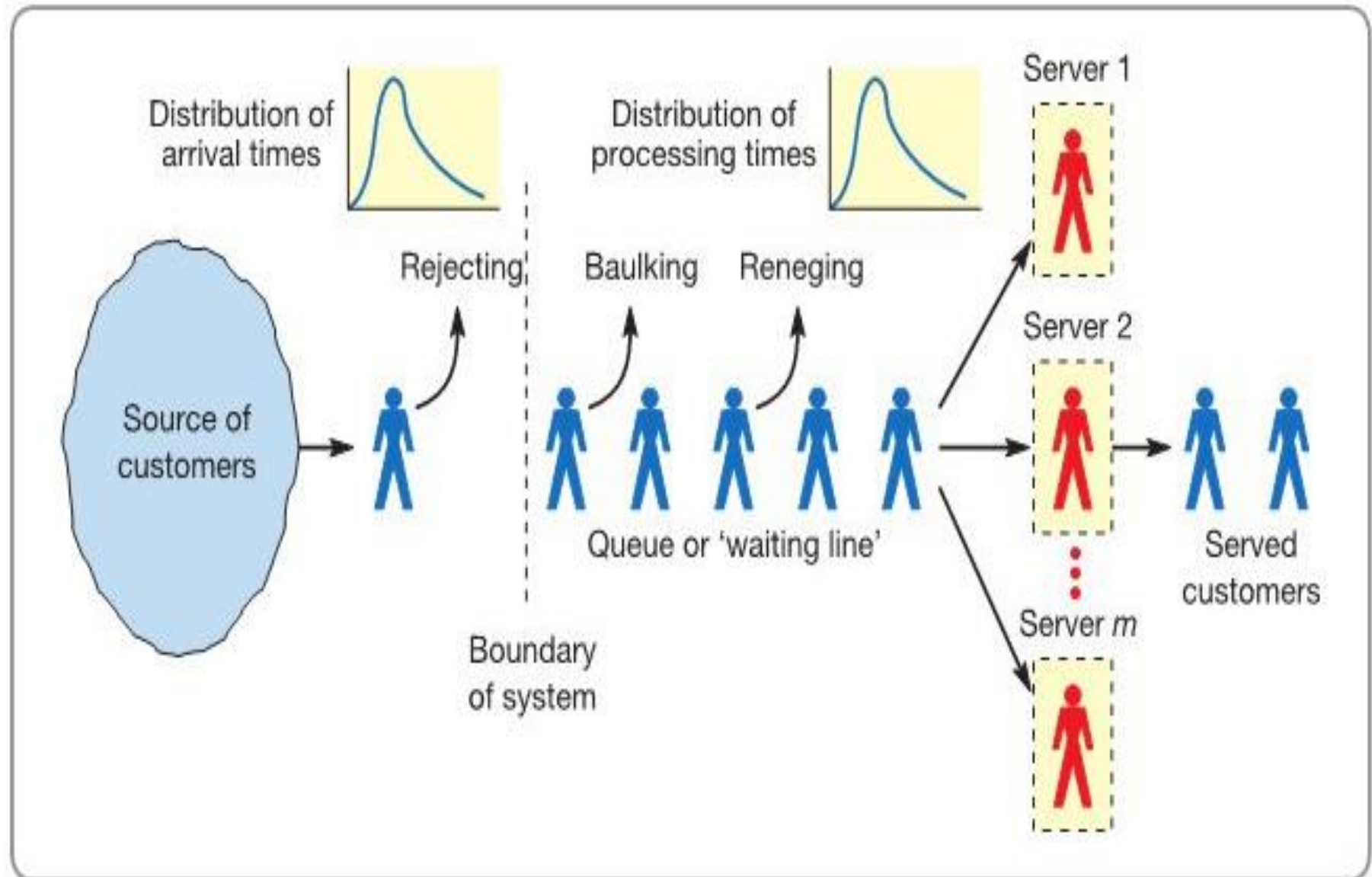
## Insights from Queuing Systems

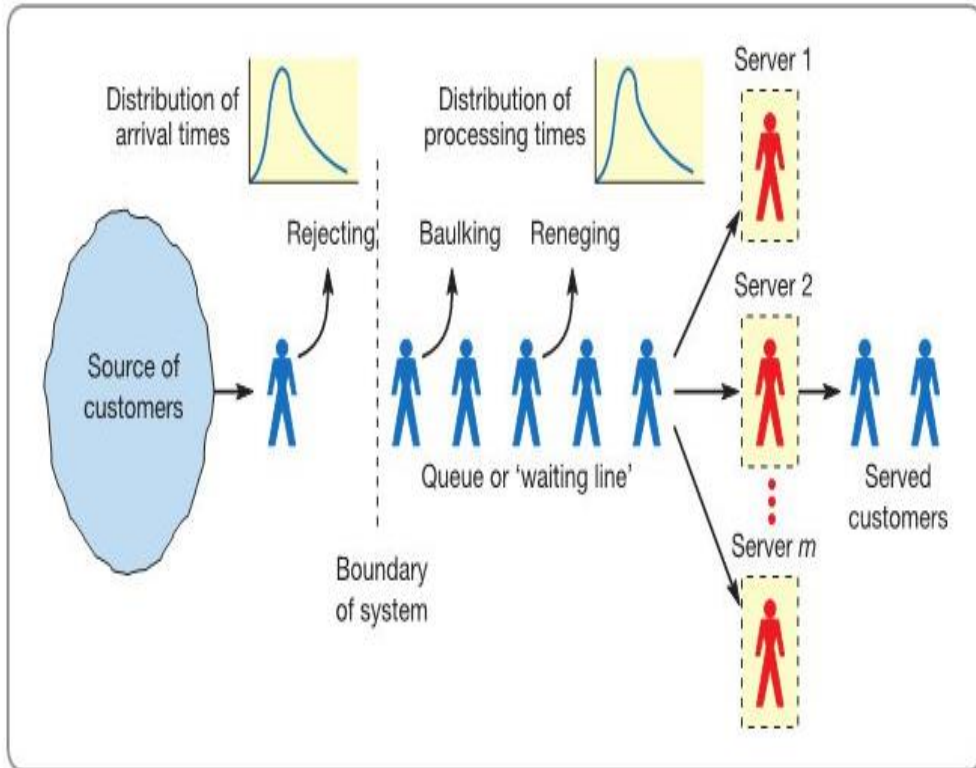Prof Neil Walton (Fong *et al.,* 2022) [Huddle 2 Nov 2023]

- www.midlandsdecisionsupport.nhs.uk/communities-of-practice/midlands-analyst-network/
- With queuing models, we usually consider the mean e.g. waiting time as a performance measure
- In healthcare we usually want some % of patients to be withing a target time
- The shape of the distribution of individual patients waiting times is (negative) exponential

**Target maximum time = 52 weeks**
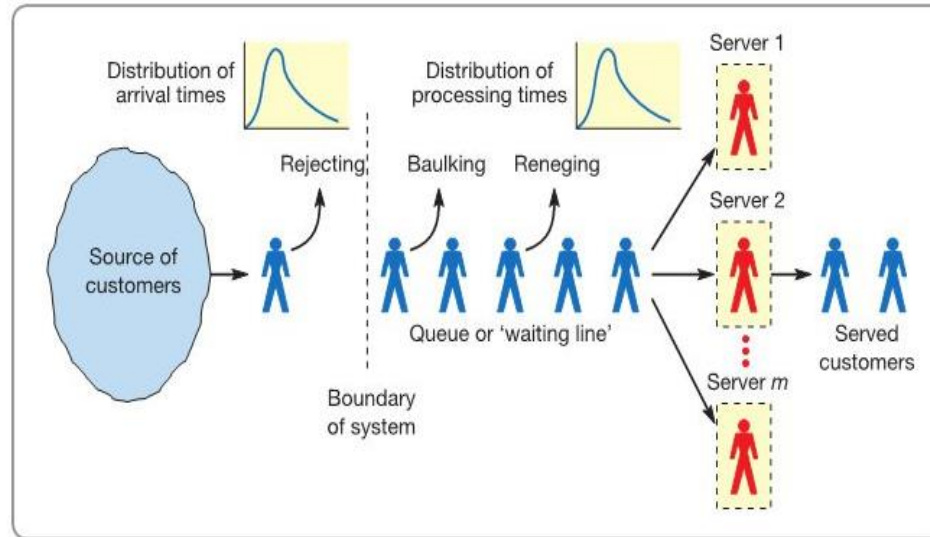
**If only 2% may exceed the target [Breaches]**

**Then Mean Waiting Time must be = Target time / 4**

Risk of a patient waiting longer than $t$
$P(\omega_q > t)$

t

# Single Echelon Queuing Models

# Generally interested in Performance metrics such as …



Distribution of arrival times

Distribution of processing times

Server 1

Rejecting  Baulking  Reneging

Server 2

Source of customers

Queue or 'waiting line'

Served customers

Server m

Boundary of system

- Expected waiting time in queue

  o  or probability (risk) that waiting time exceeds some target time

- Expected queue length

- Probability (risk) all servers busy (so customer has wait)

  o e.g., patient waits in A&E for an inpatient bed)

- Probability (risk) a restricted queue is full (so customer is rejected)

  o  e.g., patient becomes an outlier or transferred to another hospital

- Utilisation of servers (e.g., staff, cubicles, beds)

# Modelling queuing systems



- Lots of assumptions…
- Steady state (long-run averages)

**Analytical**

**Simulation**

- Freeform

**Markov**

**Generalised**
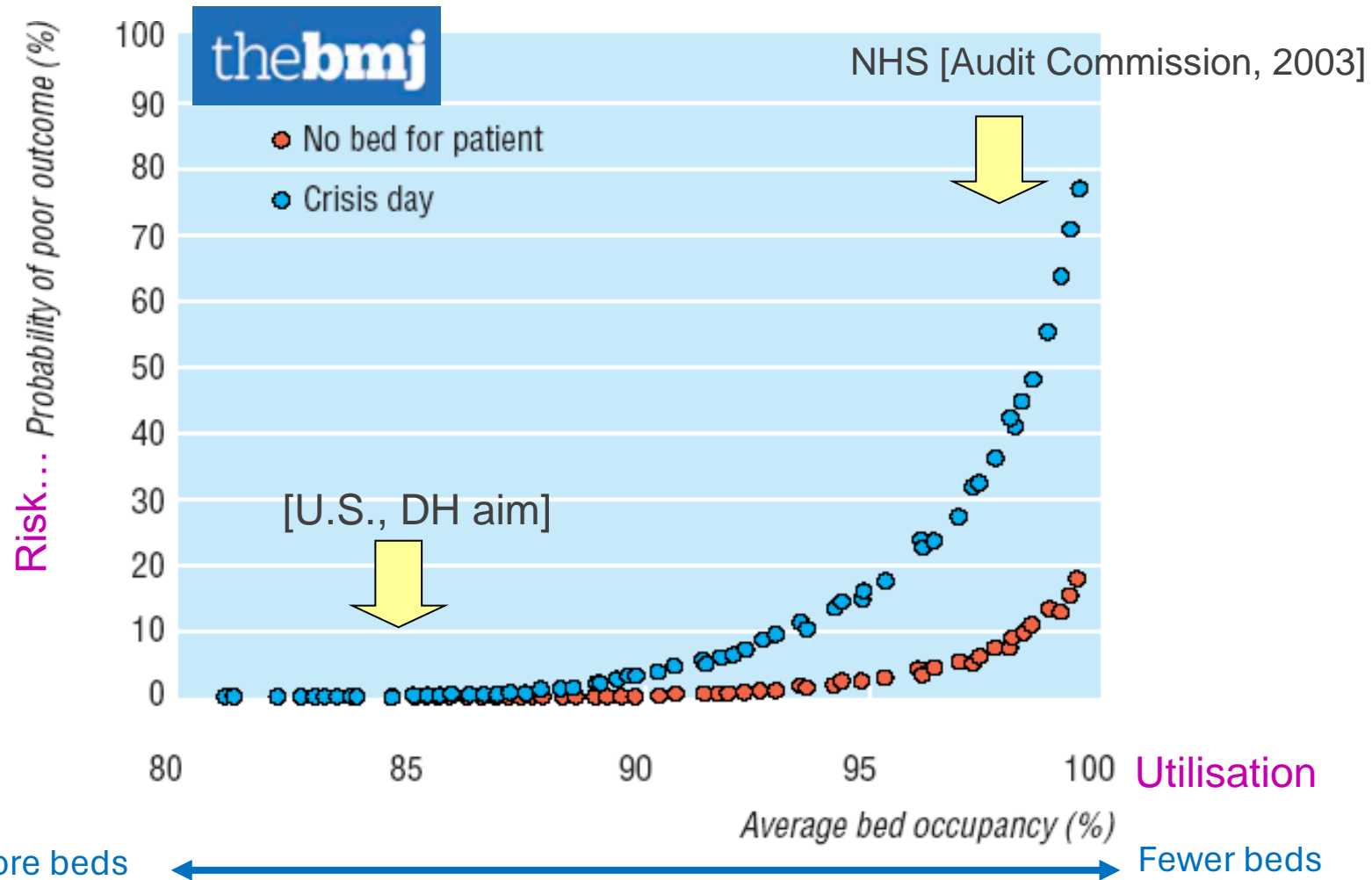
**Monte Carlo**

**Discrete Event**

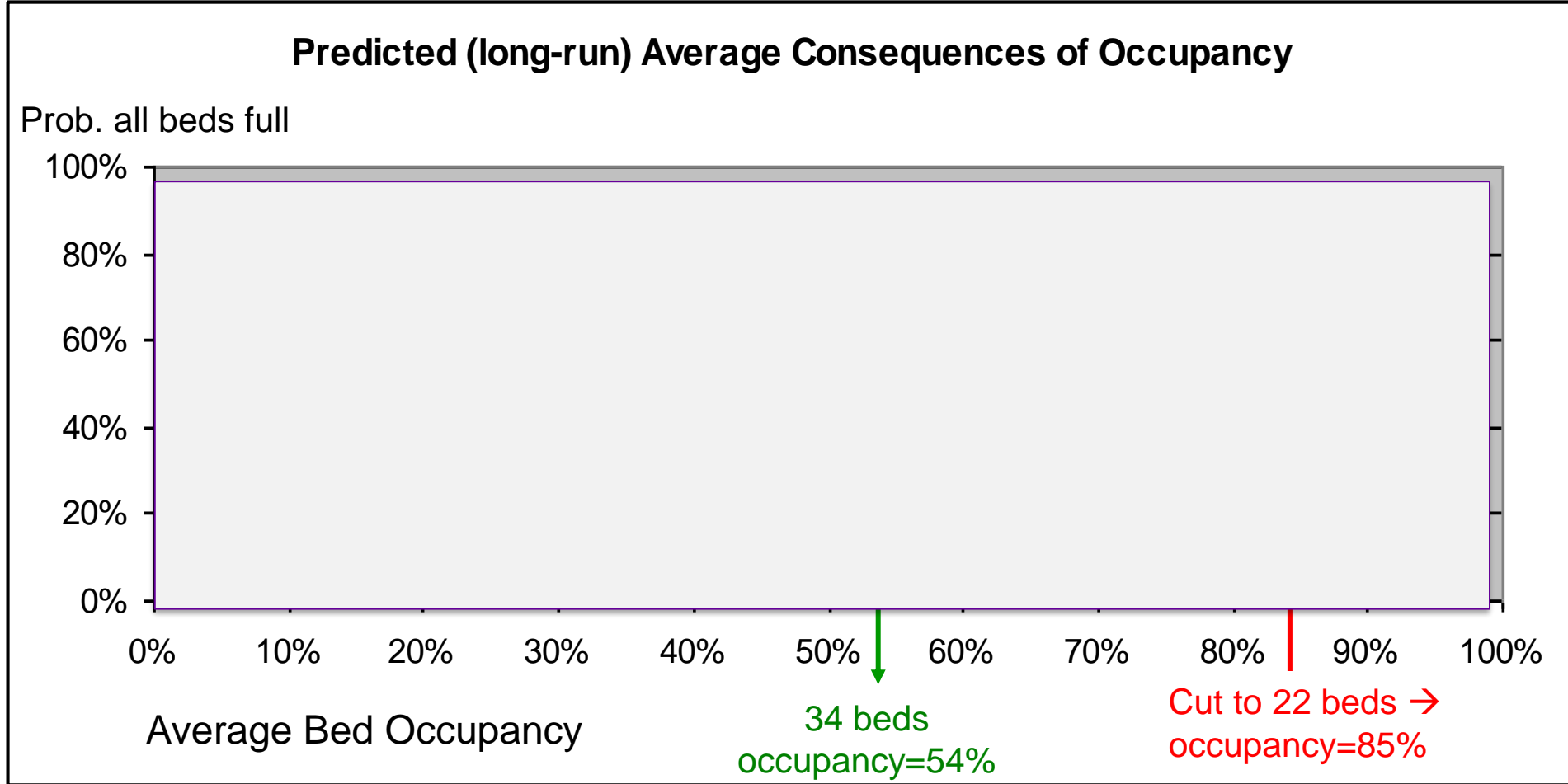| Exact equations | Approximate equations | Spreadsheets / coding | Specialist software / lots of coding |

**"bed occupancy should be 85% (or 82 or 84…)"**



results from computer simulations of large, medical inpatient bed pool (Bagust *et al.,* 1999, p.156)

# A Greater Manchester hospital's Paediatric Bed Pool (simple model!)



**Predicted (long-run) Average Consequences of Occupancy**

Prob. all beds full

100%
80%
60%
40%
20%
0%

0%   10%   20%   30%   40%   50%   60%   70%   80%   90%   100%

Average Bed Occupancy

34 beds
occupancy=54%

Cut to 22 beds →
occupancy=85%

*Occupancy* is *not* a *target*!
the appropriate level is a consequence of need to *absorb variation*

"The 85% bed occupancy fallacy: The use, misuse and insights of queuing theory" (Proudlove, 2020)

The curve was drawn using results from queuing theory (Erlang equations)
- makes simplistic assumptions, but gives you a quick idea for simple situations; for more complex situations you need simulation (more laborious!)   **7**

**Predicted (long-run) Average Consequences of Occupancy**

With more resolution:

And we can also model situations where patients are transferred when busy (cf lost calls in a telephone system – Erlang)

Legend:
- Patients Wait Model $E_C$
- Patients Transferred Model $E_B$
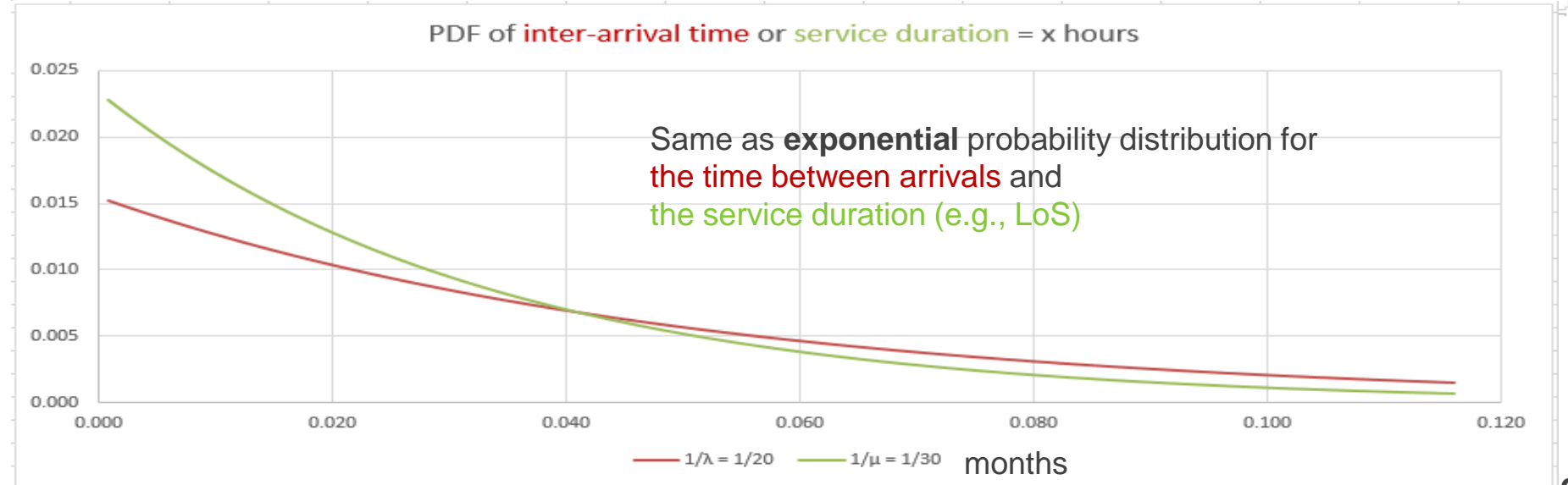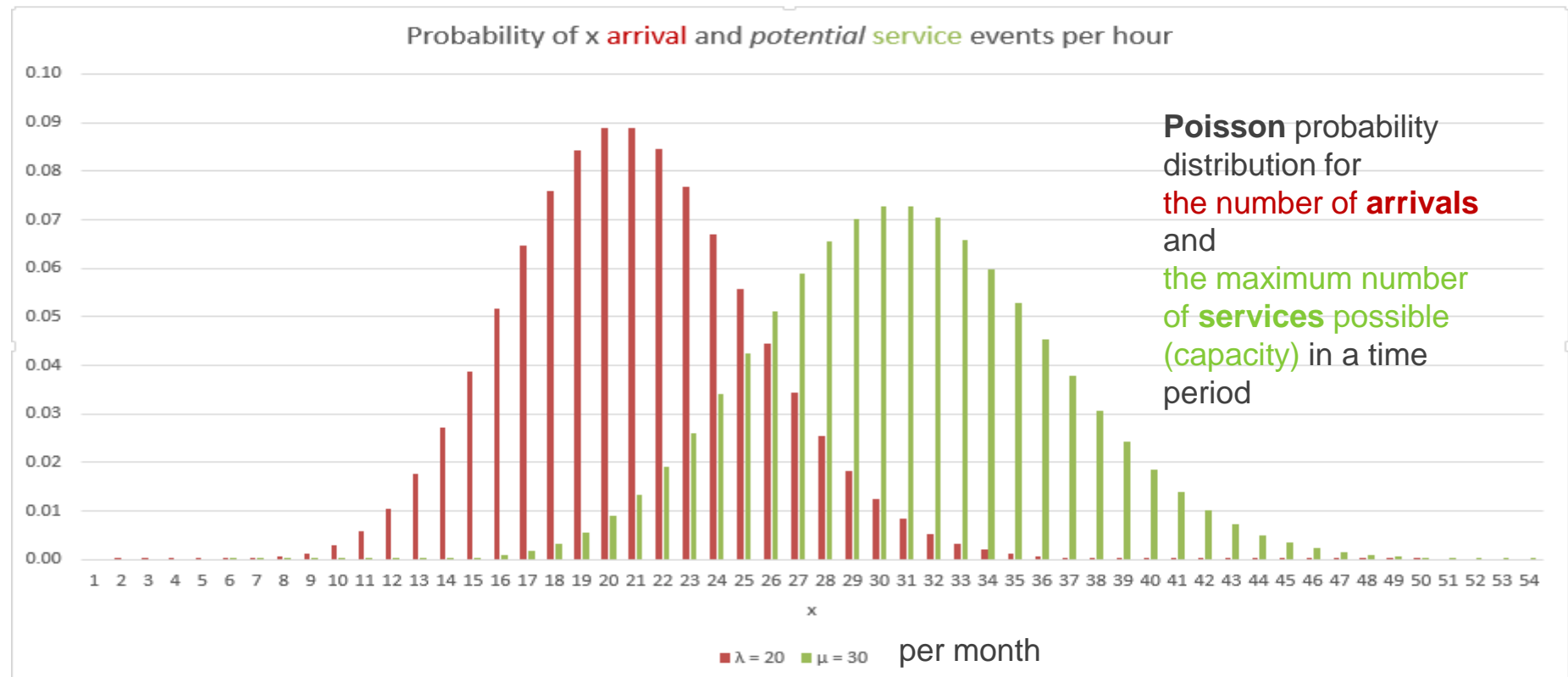
- (This graph shows the same curve as above, but for a smaller range on the X and Y scales.)
- If the starting point is that the risk of access block should be around 0.1% (i.e. a bed is available when required on 99.9% of occasions), then the Patients Wait Model suggests that the average occupancy should be around 56%, and so 33 beds would be required.
- The average occupancy should be an *outcome* of the performance required of the bed pool – it is an *output* not an input to decision making, and depends on the characteristics of the system.

Y-axis: Proportion of the time that the bed pool is full (0.0% to 1.0%)
X-axis: Average Occupancy (40% to 70%)

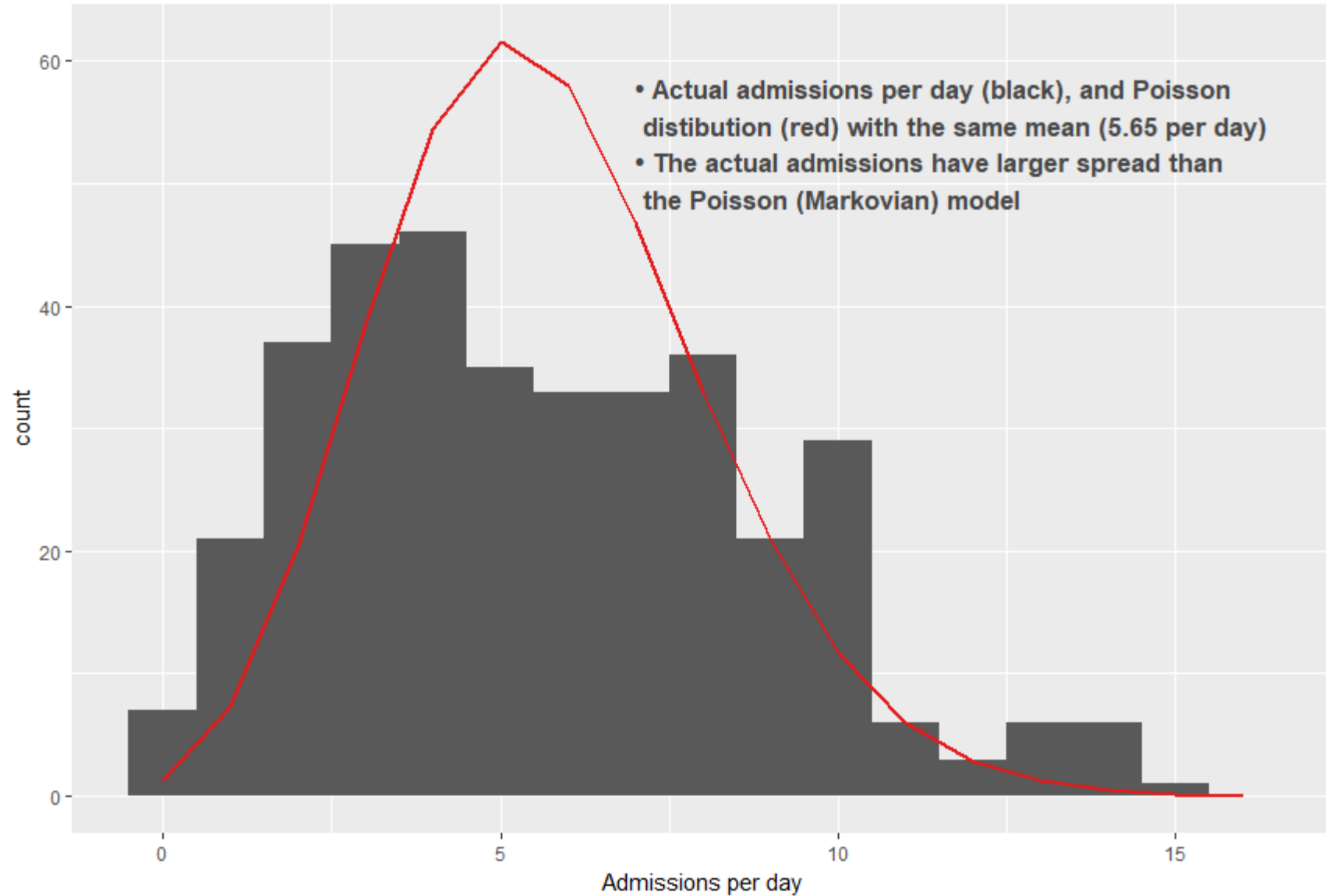| Occupancy | 46% | 47% | 49% | 50% | 51% | 53% | 54% | 56% | 58% | 60% | 62% | 64% | 66% | 69% |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Beds | 40 | 39 | 38 | 37 | 36 | 35 | 34 | 33 | 32 | 31 | 30 | 29 | 28 | 27 |

See Proudlove (2020)

8

The University of Manchester
Alliance Manchester Business School

The simplest models make lots of assumptions including:
- steady state
- 'Markovian' probability distributions

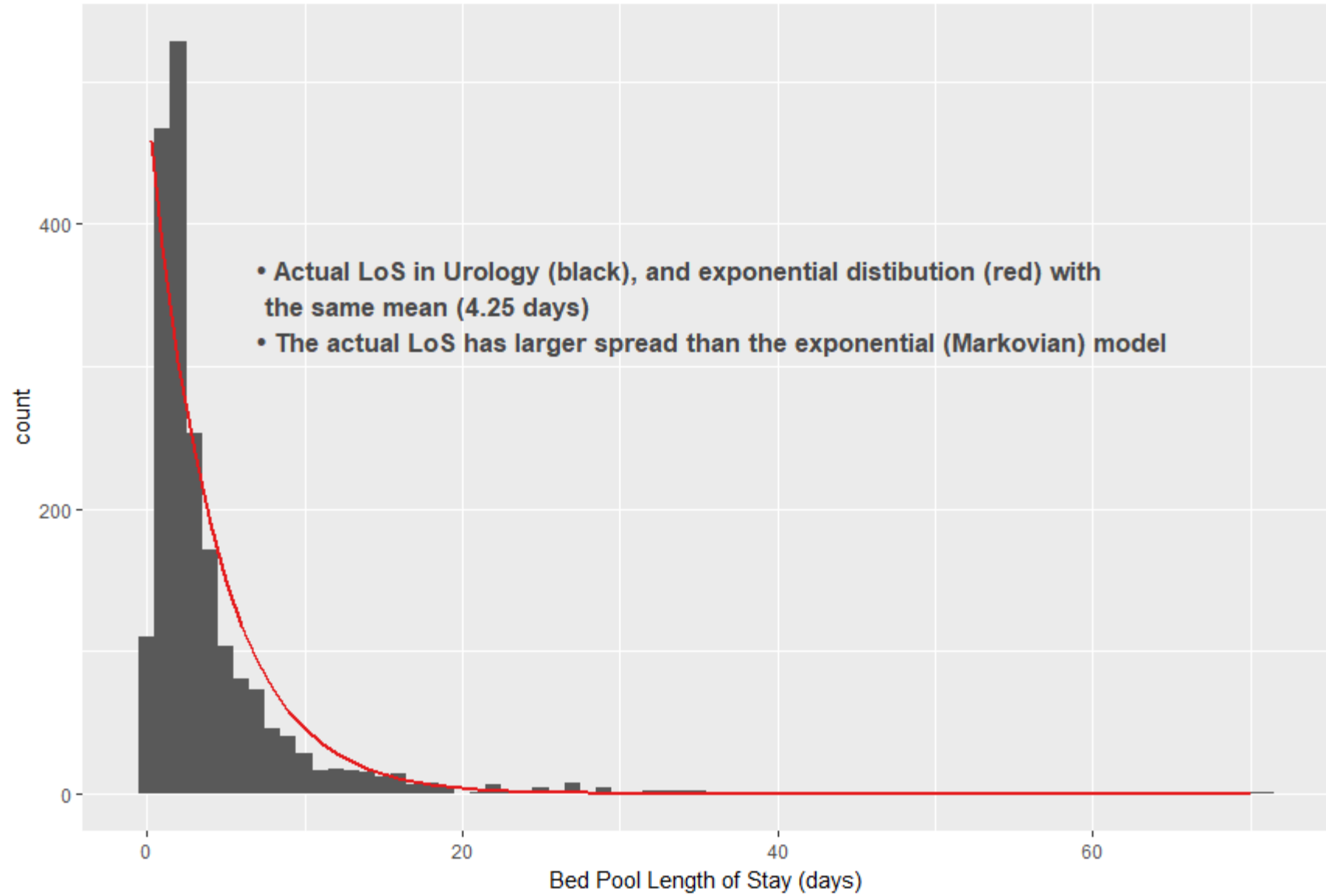## Probability of x **arrival** and *potential* service events per hour

**Poisson** probability distribution for
the number of **arrivals**
and
the maximum number of **services** possible (capacity) in a time period

λ = 20  μ = 30   per month

## PDF of **inter-arrival time** or service duration = x hours

Same as **exponential** probability distribution for
the time between arrivals and
the service duration (e.g., LoS)

1/λ = 1/20   1/μ = 1/30   months

9

A Manchester Urology Bed Pool

Do the assumptions fit?

- Actual admissions per day (black), and Poisson distibution (red) with the same mean (5.65 per day)
- The actual admissions have larger spread than the Poisson (Markovian) model

See Proudlove (2020)

- Actual LoS in Urology (black), and exponential distibution (red) with the same mean (4.25 days)
- The actual LoS has larger spread than the exponential (Markovian) model

See Proudlove (2020)

$W_q$ = expected waiting time in queue
$\lambda$ = mean arrival rate
$\mu$ = mean service rate (potential, if customers)
    $1/\mu$ = mean service duration [e.g., ALoS]
$s$ = number of servers [e.g., beds]
$\rho$ = utilisation = $\lambda/(s\mu)$
$c_a$ = coefficient of variation of arrivals
    [std dev of time between arrivals / its mean]
$c_e$ = coefficient of variation of service
    [std dev of service duration / its mean ($t_e$)]

**Simplest model** (Markovian, 1 server, customers wait)
(**M/M/1**):(GD/∞/∞)

$$W_q = \left(\frac{\rho}{1-\rho}\right)\frac{1}{\mu}$$

**Relaxing Markovian assumptions**
(to any 'General' distribution)**:**
(**G/G/1**):(GD/∞/∞)

$$W_q \approx \left(\frac{c_a^2 + c_e^2}{2}\right)\left(\frac{\rho}{1-\rho}\right)\frac{1}{\mu}$$

*The Kingman Formula*

(Markovian *c*'s are 1, so *V* term = 1)

Delay (time in queue)    Variation term    Utilisation term   Processing time

**Delay ≈ V × U × T**

**Modelling multiple servers** (from the same queue)**:**
(**G/G/s**):(GD/∞/∞)

$$W_q \approx \left(\frac{c_a^2 + c_e^2}{2}\right)\left(\frac{\rho^{\sqrt{2(s+1)}-1}}{s(1-\rho)}\right)\left(\frac{1}{\mu}\right)$$

*The VUT Relationship*

Delay (time in queue)    Variation term    Utilisation term   Processing time
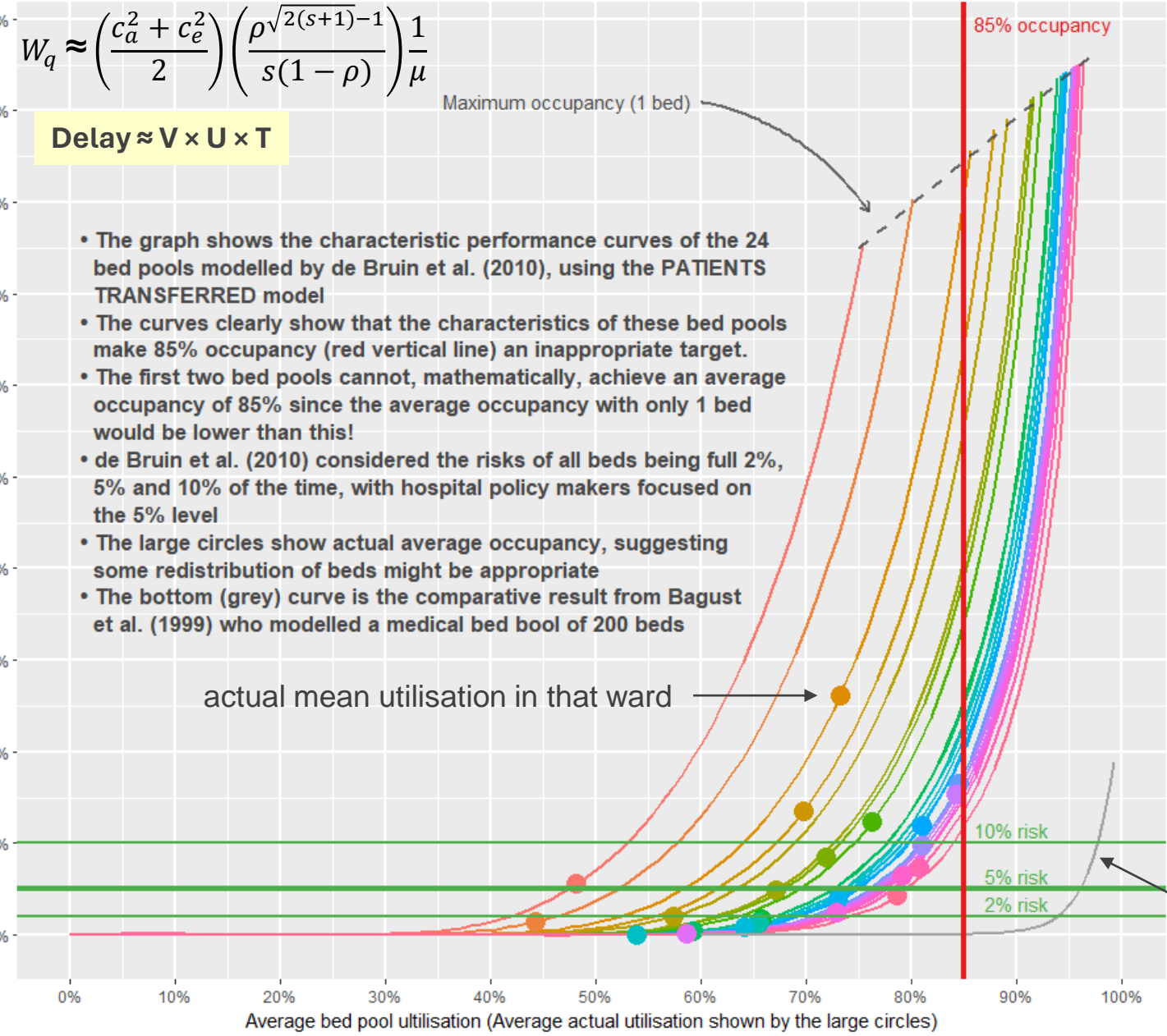
12

# Insight from the behaviour of systems:



$$W_q \approx \left(\frac{c_a^2 + c_e^2}{2}\right)\left(\frac{\rho^{\sqrt{2(s+1)}-1}}{s(1-\rho)}\right)\frac{1}{\mu}$$

**Delay ≈ V × U × T**

85% occupancy

Maximum occupancy (1 bed)

- The graph shows the characteristic performance curves of the 24 bed pools modelled by de Bruin et al. (2010), using the PATIENTS TRANSFERRED model
- The curves clearly show that the characteristics of these bed pools make 85% occupancy (red vertical line) an inappropriate target.
- The first two bed pools cannot, mathematically, achieve an average occupancy of 85% since the average occupancy with only 1 bed would be lower than this!
- de Bruin et al. (2010) considered the risks of all beds being full 2%, 5% and 10% of the time, with hospital policy makers focused on the 5% level
- The large circles show actual average occupancy, suggesting some redistribution of beds might be appropriate
- The bottom (grey) curve is the comparative result from Bagust et al. (1999) who modelled a medical bed bool of 200 beds

actual mean utilisation in that ward

10% risk
5% risk
2% risk

**Risk…** Probability all beds full, $E_B$ - so patient is transferred

Average bed pool ultilisation (Average actual utilisation shown by the large circles)

**Occupancy**

## Wards of a Dutch hospital

ward
- Special Care cardiac surgery
- Pediatric Intensive Care Unit
- Coronary Care Unit
- Medium Care
- NC Ophthalmology
- Neonatal Intensive Care Unit
- Intensive Care Unit medical
- Intensive Care Unit surgical
- NC Internal lung
- NC Pediatric unit 1
- NC Otolaryngology (ENT)
- NC Pediatric unit 2
- NC Obstetrics
- NC Internal oncology
- NC Neurology
- NC Internal medicine unit 1
- NC Internal medicine unit 2
- NC Vascular surgery
- NC Hematology
- NC Gynaecology
- NC Neuro- and orthopedic surgery
- NC Surgical oncology
- NC Cardiac surgery and cardiology
- NC Trauma surgery

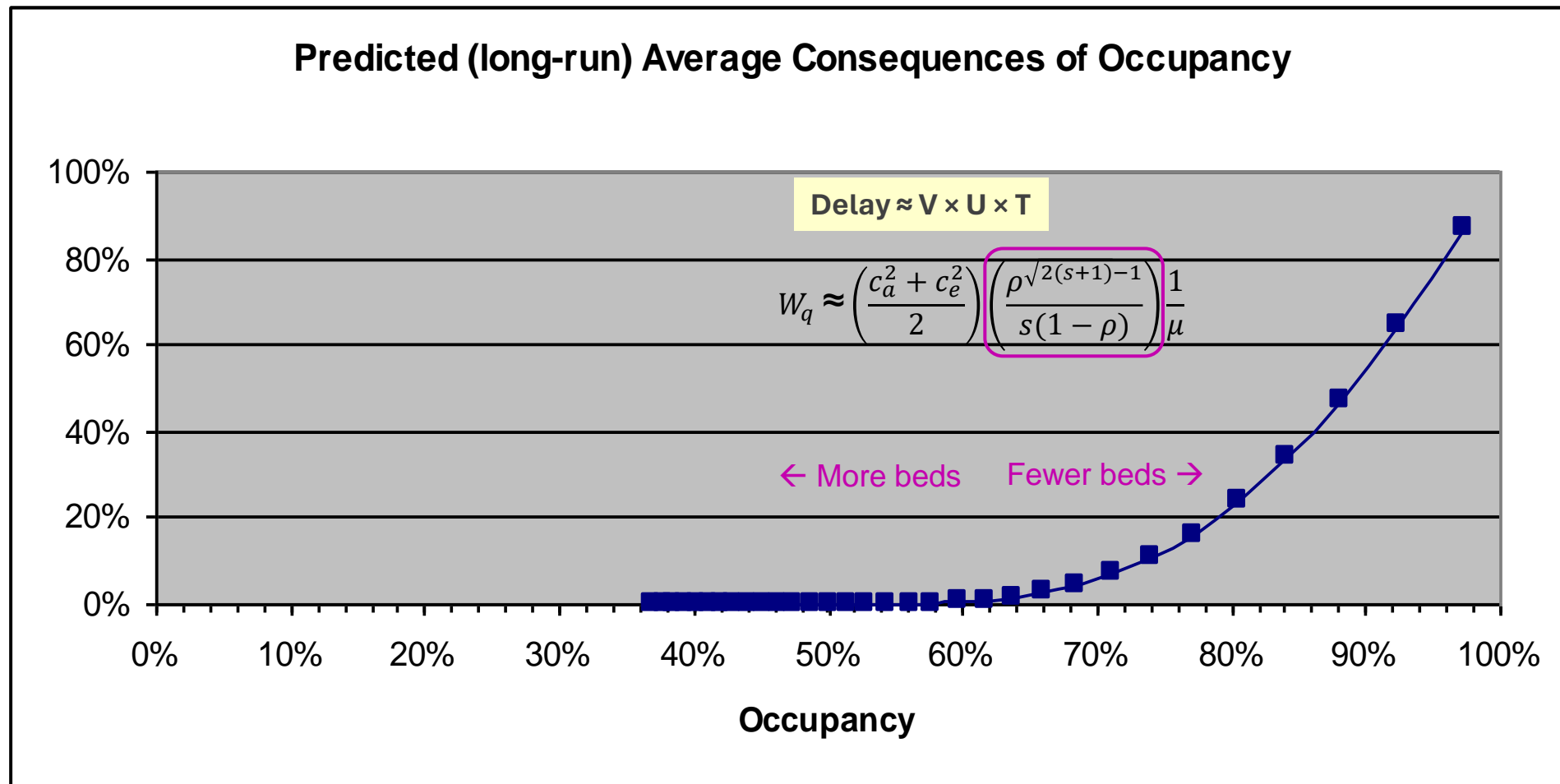Performance curve from the Bagust et al simulation (large, medical bed pool)

*Different* bed pools have *different* queuing system performance curves depending on
- Variation
- Utilisation
- (and number of servers)
- Service Duration
- ➢ *VUT curves*

- o So, the same risk of all beds being full when needed would require different average utilisations (so numbers of beds)

See Proudlove (2020)

**Accepting the systems' characteristics (variation and service times) you can slide up or down the performance curve**



**Predicted (long-run) Average Consequences of Occupancy**

Delay ≈ V × U × T

$$W_q \approx \left( \frac{c_a^2 + c_e^2}{2} \right) \left( \frac{\rho^{\sqrt{2(s+1)}-1}}{s(1-\rho)} \right) \frac{1}{\mu}$$
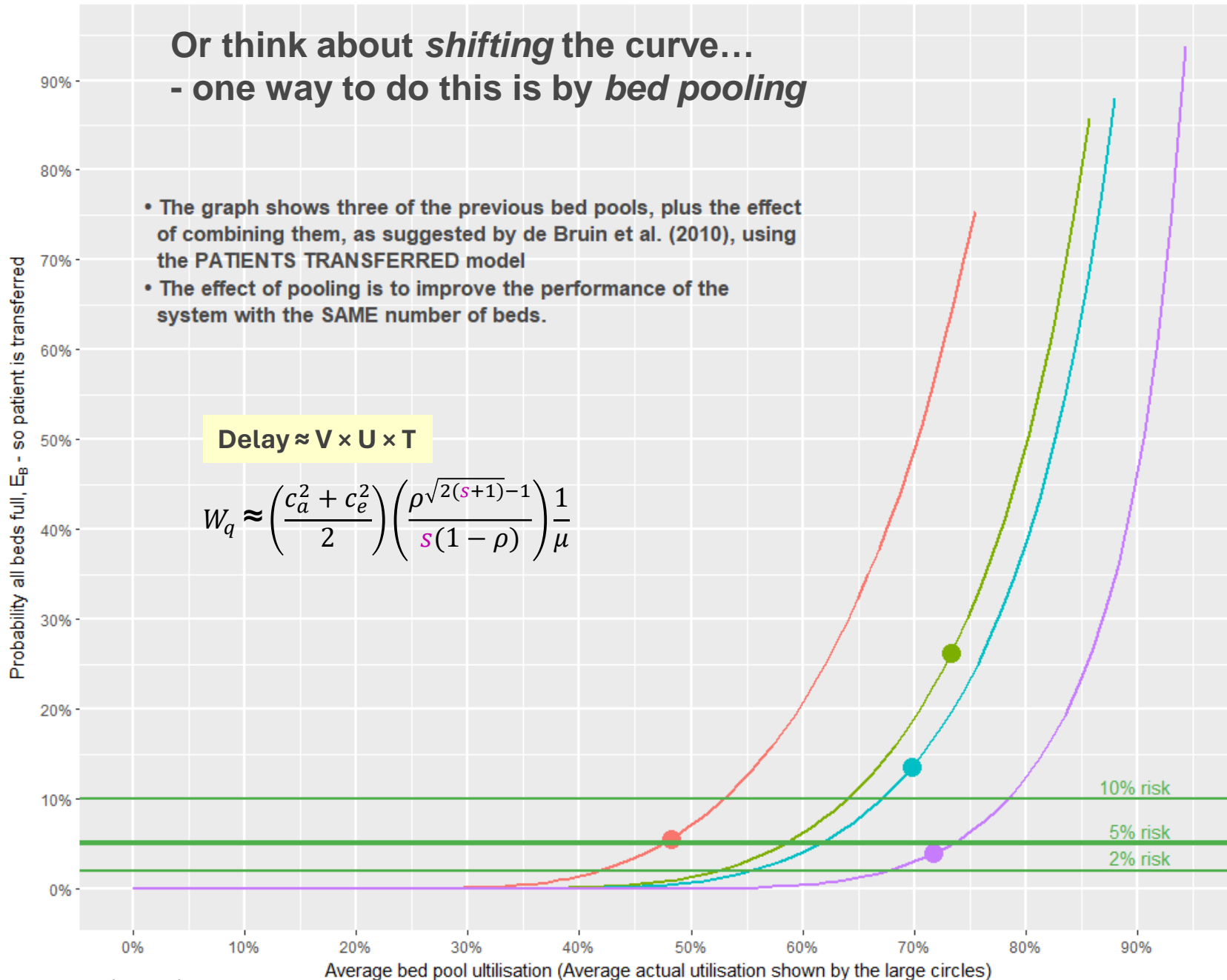
← More beds    Fewer beds →

Occupancy

Or think about *shifting* the curve…
- one way to do this is by *bed pooling*

- The graph shows three of the previous bed pools, plus the effect of combining them, as suggested by de Bruin et al. (2010), using the PATIENTS TRANSFERRED model
- The effect of pooling is to improve the performance of the system with the SAME number of beds.

**Delay ≈ V × U × T**

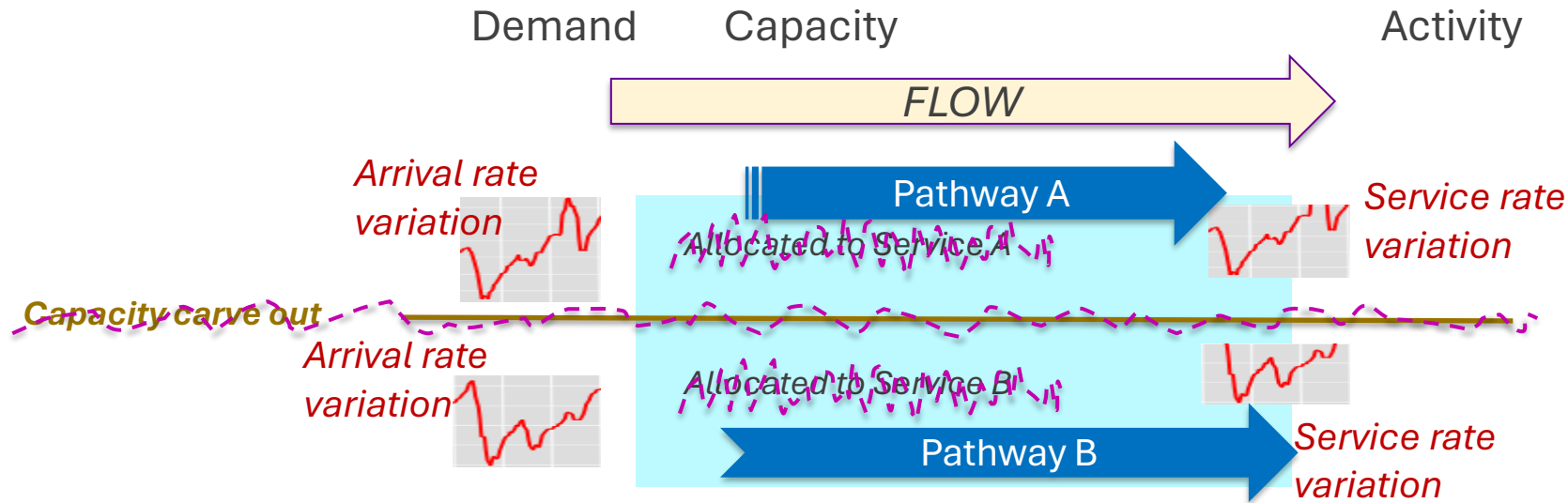$$W_q \approx \left(\frac{c_a^2 + c_e^2}{2}\right)\left(\frac{\rho^{\sqrt{2(s+1)}-1}}{s(1-\rho)}\right)\frac{1}{\mu}$$

Probability all beds full, $E_B$ - so patient is transferred

Average bed pool utilisation (Average actual utilisation shown by the large circles)

ward
- Special Care cardiac surgery
- Coronary Care Unit
- Medium Care
- Combined Pool

10% risk
5% risk
2% risk

- *Same* aggregate utilisation (work being done)
- ➤ *Much better performance*
- *Sharing the load reduces the risks from tail events*

See Proudlove (2020)

15

# Capacity *Carve-out* vs. *Pooling* or *Segmentation*

## Capacity Pooling

Demand    Capacity                                    Activity

*FLOW*

*Arrival rate variation*

Pathway A

*Allocated to Service A*

*Service rate variation*

**Capacity carve out**

*Arrival rate variation*

*Allocated to Service B*

Pathway B

*Service rate variation*



The graph shows three of the previous bed pools, plus the effect of combining them, as suggested by de Bruin et al. (2010), using the PATIENTS TRANSFERRED model
The effect of pooling is to improve the performance of the system with the SAME number of beds.



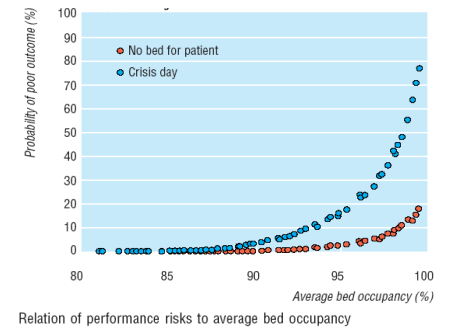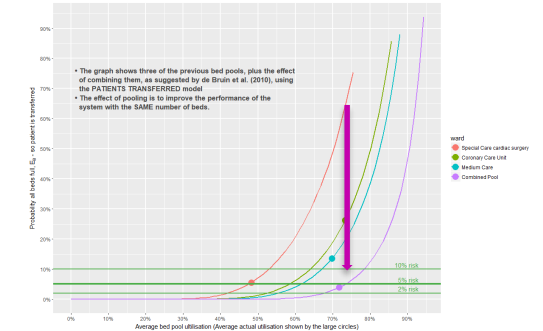Relation of performance risks to average bed occupancy

### Segmentation
- Tailoring to customer segment: faster service rates and/or lower variety of job types
- More efficient pathways outweigh carve-out
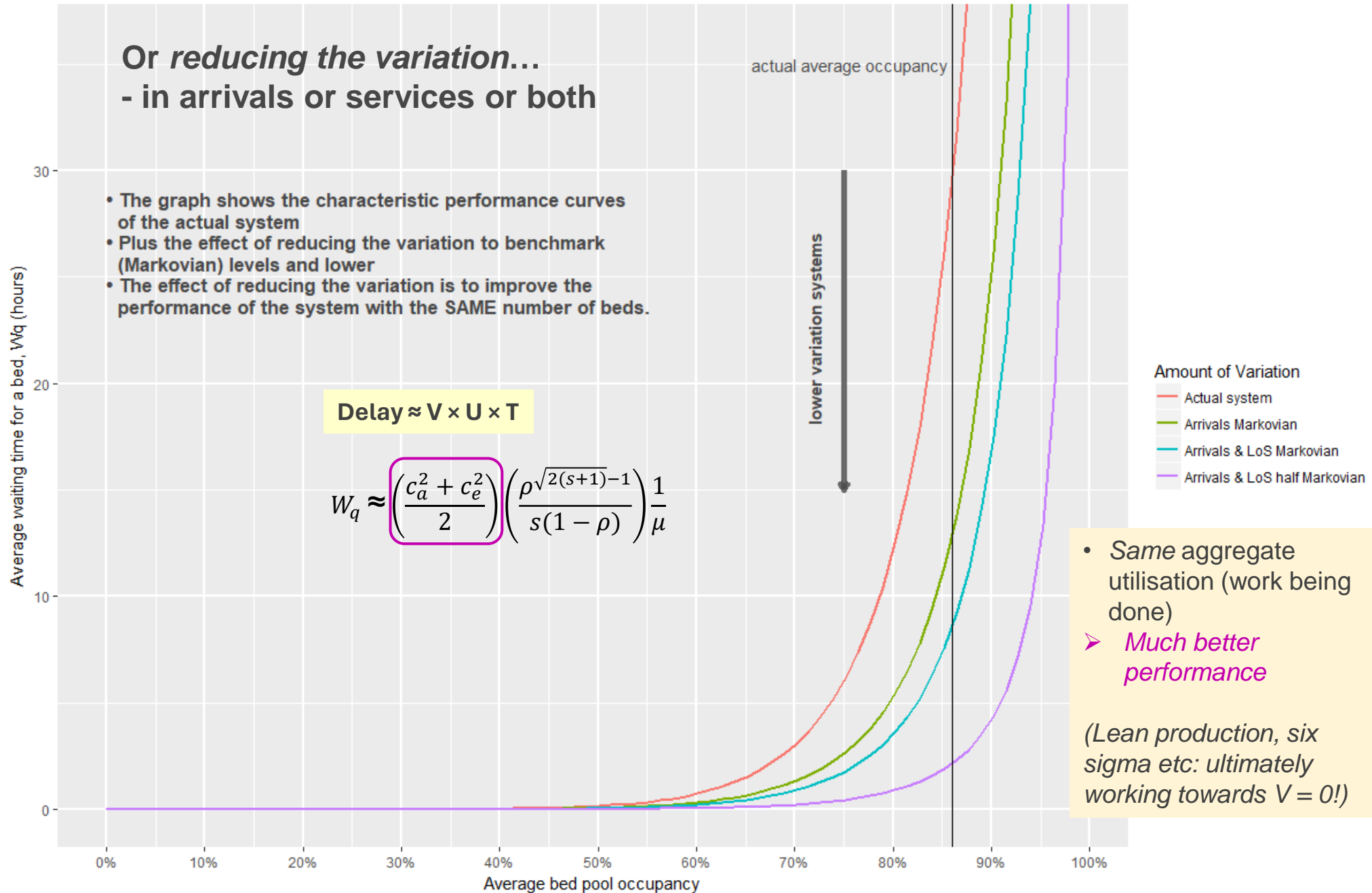
What happened in NHS trusts?!

- There may be good reasons for 'carve-out'? (depends on system and objectives)
  - But can you increase flexibility? [e.g., short-notice call-in to unused appointment slots carved-out for expected urgent demand?]
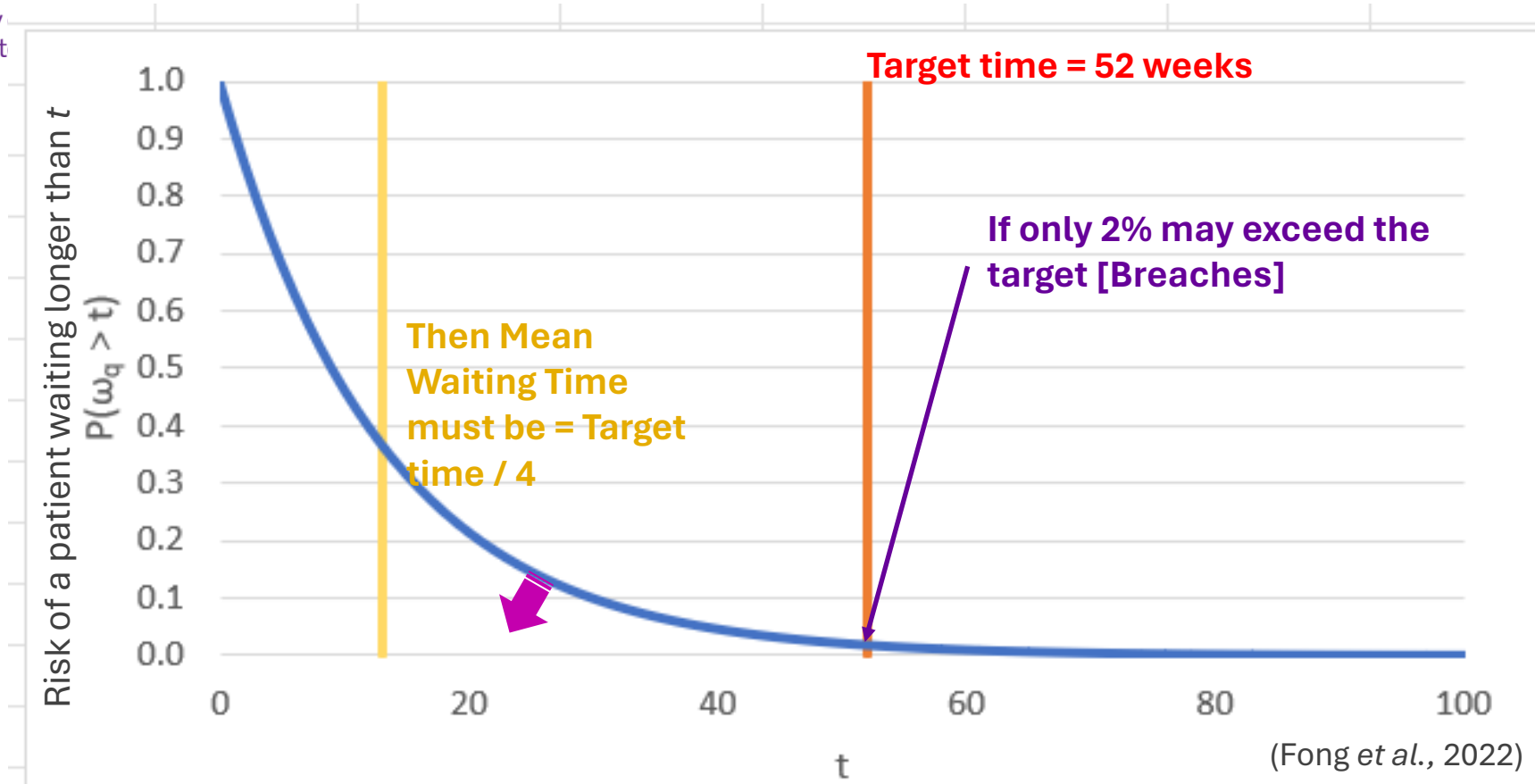
**Or *reducing the variation*…**
**- in arrivals or services or both**

- The graph shows the characteristic performance curves of the actual system
- Plus the effect of reducing the variation to benchmark (Markovian) levels and lower
- The effect of reducing the variation is to improve the performance of the system with the SAME number of beds.

**Delay ≈ V × U × T**

$$W_q \approx \left(\frac{c_a^2 + c_e^2}{2}\right)\left(\frac{\rho^{\sqrt{2(s+1)}-1}}{s(1-\rho)}\right)\frac{1}{\mu}$$

actual average occupancy

lower variation systems

Average waiting time for a bed, Wq (hours)

Average bed pool occupancy

**Amount of Variation**
— Actual system
— Arrivals Markovian
— Arrivals & LoS Markovian
— Arrivals & LoS half Markovian

- *Same* aggregate utilisation (work being done)
➢ *Much better performance*

*(Lean production, six sigma etc: ultimately working towards V = 0!)*

See Proudlove (2020)

17

Target time = 52 weeks

If only 2% may exceed the target [Breaches]

Then Mean Waiting Time must be = Target time / 4

Risk of a patient waiting longer than $t$ — $P(\omega_q > t)$

(Fong *et al.*, 2022)

- The **design of the system** and the **variation** make system performance highly non-linear
  - In particular, the long and fat tails
- Meaning low risk of poor performance [low breaches, trolley waits etc] requires very much better **average** performance
  - So lower utilisation → more resource
  - (and) or **improve the design of the system** and/or **reduce the variation**!

# Take aways

- Characteristics of system drive performance
    - "85% occupancy" does not fit all environments
    - Mean occupancy levels should be a consequence of the demand characteristics you need to absorb
        - 'empty' capacity protects the system
    - e.g., sensible utilisation levels for knee surgery vs. ICU

- Queuing theory models give quick, first-cut results
    - Make *a lot* of assumptions…
    - But give good **insights** – e.g. the VUT relationship
    - Beyond that is simulation (laborious, data-hungry, requires specialist knowledge and software)

- Can *shift* the trade-off
    - **Bed pooling**, but
        - Pooling *vs* carve-out *or* segmentation?
        - Behavioural impacts?!
    - **Reducing variation**
        - How?!



(Slack et al., 2011)

19

**References**

Bagust A, Place M and Posnett J (1999). "Dynamics of bed use in accommodating emergency admissions: stochastic simulation model". *British Medical Journal* 319, 155-158. doi: 10.1136/bmj.319.7203.155

Fong K, Mushtaq Y, House T, Gordon D, Chen Y, Griffths D, Ahmad S and Walton N (2022). "Understanding waiting lists pressures". *medRxiv*, 2022.08.23.22279117. doi: 10.1101/2022.08.23.22279117

Proudlove NC (2020). "The 85% bed occupancy fallacy: The use, misuse and insights of queuing theory". *Health Services Management Research* 33:3, 110-121. doi: 10.1177/0951484819870936

Proudlove NC (2023). "Chapter 11: Use and misuse of Queueing Theory for hospital capacity decisions" in Vissers J, Elkhuizen S and Proudlove N (Eds.) *Operations Management for Healthcare* 2nd ed. Routledge: Abingdon, UK. doi: 10.4324/9781003020011

Slack N, Brandon-Jones A and Johnston R (2016). *Operations management* 8th ed. Pearson Education: Harlow.