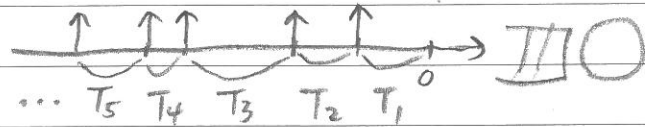A proof that the interarrival times of the Poisson
process are independent and exponentially distributed
with rate $\lambda$ :



$$\cdots \quad T_5 \quad T_4 \quad T_3 \quad T_2 \quad T_1$$

$P\{T_1 > z\} = P\{A(z) = 0\}$ (by property ① in the definition
of Poisson
process )

$$= P\{A(0+z) - A(0) = 0\}$$

$$= e^{-\lambda z} \cdot \frac{(\lambda z)^0}{0!} \quad \text{(by property ③ '' )}$$

$$= e^{-\lambda z}$$

Now, $P\{T_n > z \mid \sum_{k=1}^{n-1} T_k = s\}$ for some $s > 0$

$$= P\{A(s+z) - A(s) = 0 \mid \sum_{k=1}^{n-1} T_k = s\}$$

$$= P\{A(s+z) - A(s) = 0\} \quad \text{(by property ② '' )}$$

$$= e^{-\lambda z} \cdot \frac{(\lambda z)^0}{0!} \quad \text{(by property ③ '' )}$$
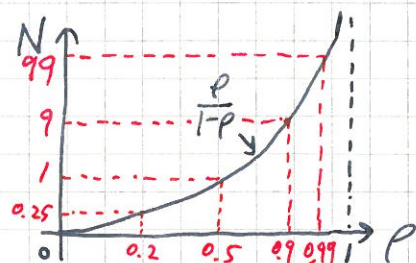
$$= e^{-\lambda z}$$

$\Rightarrow P\{T_n \leq z\} = 1 - e^{-\lambda z}$

$\Rightarrow f(T_n = z) = \lambda e^{-\lambda z}$, which is the p.d.f. of exponential
distribution ✳

In M/M/1, interpretation of $N = \frac{\rho}{1-\rho}$ :
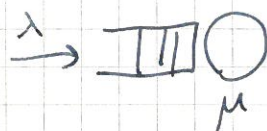
think of $\rho$ as CPU utilization ( CPU% )



since in M/M/1, $\rho = \frac{\lambda}{\mu}$

$\rightarrow \rho > 1$ : server cannot catch up with arrivals

$\rightarrow \rho < 1$ : arrivals cannot feed the server fast enough

$\Rightarrow$ # of items in a queue (i.e., N) grows much faster than the growth of CPU utilization (i.e., $\rho$) !
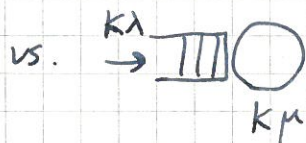
Example 3.8 :



$\rho = \frac{\lambda}{\mu}$

$N = \frac{\rho}{1-\rho} = \frac{\lambda}{\mu - \lambda}$

$T = N/\lambda$

vs.

$\rho = \frac{k\lambda}{k\mu} = \frac{\lambda}{\mu}$

$N = \frac{\rho}{1-\rho} = \frac{\lambda}{\mu - \lambda}$

$T = N/k\lambda$

$\Rightarrow$ in both cases, a new arrival will see the same # of packets ahead of it statistically.

But ~~Though~~ the packets move K times faster ~~for~~ in the configuration on the right! ※
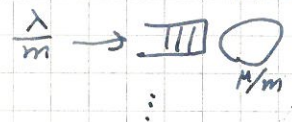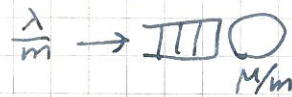
---

Example 3.9 : Statistical multiplexing vs. TDM/FDM

Consider m statistically identical and independent Poisson packet streams each with arrival rate $\frac{\lambda}{m}$. The packet length for all streams are independent and exponentially distributed. The average transmission time is $\frac{1}{\mu}$.

$\frac{}{}$ If merging all streams into one, it is like M/M/1 :

 $m \cdot \frac{\lambda}{m}$

and thus the average delay per packet is

$$T = \frac{1}{\mu - \lambda}$$

If using TDM/FDM instead, it is like m M/M/1 :



$\frac{\lambda}{m} \rightarrow$ M/m

$\frac{\lambda}{m} \rightarrow$ M/m

Then the average delay would be m times larger

$$T' = \frac{1}{\mu/m - \lambda/m} = \frac{m}{\mu - \lambda} = m \cdot T$$

$\Rightarrow$ why would anyone ever want to use TDM/FDM ?
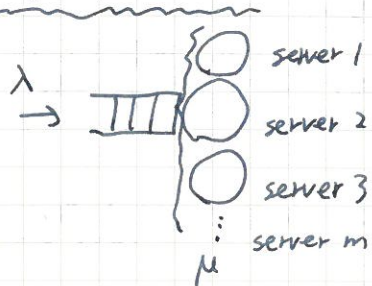
$\rightarrow$ TDM/FDM guarantees a specific service rate for each stream. Plus, statistical multiplexing may introduce^a higher variability in latency, which is undesirable for applications such as voice/video streaming.

The PASTA theorem : <u>P</u>oisson <u>A</u>rrivals <u>S</u>ee <u>T</u>ime <u>A</u>verage

( read Sections 3.3.2 – 3.3.3 )

→ " When arrivals are Poisson, [···] both an arriving and a departing customer in steady-state see a system that is statistically identical to the one seen by an observer looking at the system at an arbitrary time." We may use that to analyze the probability that an arrival will need to wait in queue.
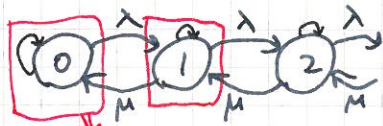
## M/M/m



let $\mu$ be the service rate,

$$\rho = \frac{\lambda}{m\mu} \text{ the utilization}$$

To analyze M/M/m, it is helpful to take a closer look at how we got $P_n \lambda = P_{n+1} \mu$ in M/M/1 :
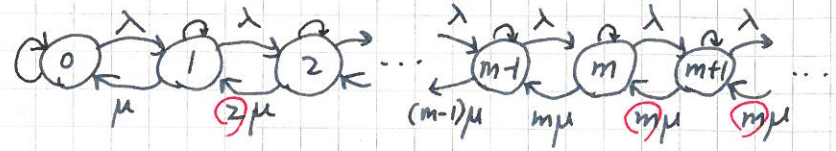
recall the Markov chain in M/M/1



$\sum \binom{\text{probability}}{\text{of entering a state}} = \sum \binom{\text{probability}}{\text{of leaving the state}}$

$P_0 \lambda = P_1 \mu$

$P_1 \lambda + P_1 \mu = P_0 \lambda + P_2 \mu$

...

$P_n \lambda + P_n \mu = P_{n-1} \lambda + P_{n+1} \mu$

$\Rightarrow P_n \lambda = P_{n+1} \mu$

On M/M/m, we have the following Markov chain



<u>for $n \leq m-1$,</u>

$P_0 \cdot \lambda = P_1 \mu$

$P_1 \lambda + P_0 \mu = P_0 \cdot \lambda + P_2 \cdot 2\mu$

$P_2 \cdot \lambda + P_1 \cdot 2\mu = P_1 \cdot \lambda + P_3 \cdot 3\mu$

...

$\Rightarrow P_n \cdot \lambda = P_{n+1} \cdot (n+1) \mu$

$\Rightarrow P_m \cdot m\mu = P_{m-1} \cdot \lambda$

$P_m = \frac{\lambda}{m\mu} P_{m-1}$

$\Rightarrow P_n = \frac{\lambda^n}{n! \, \mu^n} P_0$

from $\boxed{\rho = \frac{\lambda}{m\mu}}$

we see that $\left(\frac{\lambda}{\mu}\right)^n = (m\rho)^n$

therefore $P_n = \frac{(m\rho)^n}{n!} P_0$

<u>for $n \geq m$,</u>

$P_m \cdot \lambda + P_m (m\mu) = P_{m-1} \lambda + P_{m+1} (m\mu)$

$P_{m+1} \lambda + P_{m+1} (m\mu) = P_m \lambda + P_{m+2} (m\mu)$

...

$\Rightarrow P_n \cdot \lambda = P_{n+1} (m\mu).$

→ similarly, for $n \geq m$ we have

$P_n = \frac{\lambda}{m\mu} P_{n-1} = \rho^{n-m} \frac{(m\rho)^m}{m!} P_0$

$= \frac{(m\rho)^n}{m! \, m^{n-m}} P_0 = \frac{m^m \rho^n}{m!} P_0$

$= \frac{\rho^n \cdot m^n}{m! \cdot m^{n-m}} P_0$

$\Rightarrow P_n = \frac{m^m \rho^n}{m!} P_0$

Then $P_0$ can be obtained by $\sum_{n=0}^{\infty} P_n = 1.$

Let $P_Q$ be the probability that an arrival will need to wait in queue (because all m servers are busy), and we have

$P_Q = \sum_{n=m}^{\infty} P_n = \sum_{n=m}^{\infty} P_0 \frac{m^m \rho^n}{m!} = P_0 \frac{m^m \rho^m}{m!} \sum_{n=m}^{\infty} \rho^{n-m}$

$\Rightarrow \boxed{P_Q = \frac{(m\rho)^m P_0}{m! \, (1-\rho)}}$ ← The Erlang C formula, named after the pioneer of Queueing Theory, A. K. Erlang.

Now let $N_Q$ be the expected # of customers waiting in queue. We may obtain it by $\sum_{n=0}^{\infty} n \cdot P_{m+n}$.

Alternatively, we may consider M/M/m's relation to M/M/1, and get $N_Q = P_Q \cdot \frac{\rho}{1-\rho}$ ~~$\frac{\rho}{1-\rho}$~~

→ from M/M/1

Finally, $\underline{N = N_Q + N_S}$ where $N_S$ is # of customer in service

→ from the linearity of expected #.

$N_S = \sum_{n=1}^{m} n \cdot P_n = \ldots$

we can use Little's Theorem here, and $N_S = \lambda \cdot \frac{1}{\mu}$

$= m\rho$

$\Rightarrow \boxed{N = m\rho + \frac{\rho P_Q}{1-\rho}}$

$\Rightarrow \boxed{T = N/\lambda = \frac{1}{\mu} + \frac{P_Q}{m\mu - \lambda}}$ ← using Little's Theorem

for the number of customers in the queue.

Detail derivation: let $X$ be a random variable

$N_Q = P_Q \cdot E[X | queueing] + P\{no\ queueing\} \cdot E[X | no\ queueing]$

$= P_Q \cdot E[X | queueing]$

To get $E[X | queueing]$, we note that

$E[X | queueing]$ in M/M/m with service rate $\lambda$ is equal to $E[X | queueing]$ in M/M/1 with service rate $m\lambda$

let $N_{Q_1}$ be the expected # of customers waiting in this M/M/1 system, then $N_{Q_1} = \frac{\rho}{1-\rho} - \rho = \frac{\rho^2}{1-\rho}$
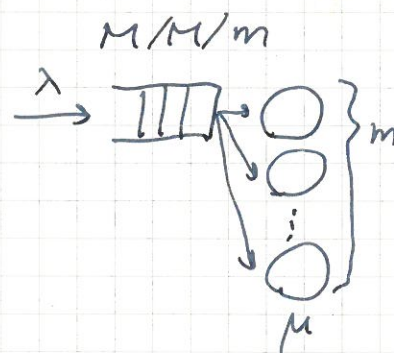
From $N_{Q_1} = P_{Q_1} \cdot E[X | queueing]$, we have

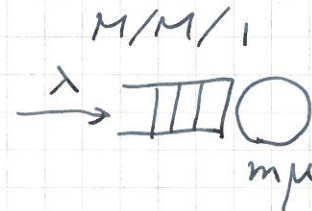$E[X | queueing] = N_{Q_1} / P_{Q_1} = \frac{\rho^2}{1-\rho} / \rho = \frac{\rho}{1-\rho}$

and therefore $N_Q = P_Q \cdot E[X | queueing] = P_Q \cdot \frac{\rho}{1-\rho}$ ✳

---

Example 3.10 (compared with examples 3.8 and 3.9)

M/M/m          v.s.          M/M/1



$T = \frac{1}{\mu} + \frac{P_Q}{m\mu - \lambda}$

$T' = \frac{1}{m\mu - \lambda}$

$\rho = \frac{\lambda}{m\mu}$

$\Rightarrow$ When $\rho \ll 1$,

$P_Q \approx 0$ and $m\mu \gg \lambda$, which imply

$\frac{T}{T'} \approx m$

$\Rightarrow$ When $\rho \to 1$,

$P_Q \approx 1$ and $\frac{1}{\mu} \ll \frac{1}{m\mu - \lambda}$, which imply

$\frac{T}{T'} \approx 1$

faster i.e., server

The above result suggests that using one (channel) for statistical multiplexing will lead to a better performance in latency, compared with one using multiple channels. i.e., servers

slower