

POWER & THROUGHPUT ISSUES IN NEXT-GENERATION PACKET SWITCHES

Nick Bambos
Stanford University

The Rising Problem of Power (Density)

Switching: chips ~100W+ ... systems ~ KW ... clusters/farms/centers ~ 10KW++

Computing: processors ~100W+/- ... systems ~200W+

Racks: ~ 512 CPUs+ ~100KW+ ... *high power density* ... *liquid cooled!*

Data centers: 100,000 CPUs+ ~10MW++ ... high power density

High power distribution cost... high power bills... *high cooling costs...*

Electronic infrastructure costs get amortized ... power costs don't...

Both matter: power *volume* and power *density* (perhaps more)

High/uneven power density ... high thermal stress ... fast component aging ...

reliability problems...

Throughput/Speed vs. Power

Throughput/speed can *potentially*... be adjusted by

frequency scaling... clocking up/down

voltage scaling...

component shutdown (path/structure scaling)

...

Power scales up ***super-linearly*** with throughput/speed

Where/How to Attack the Problem

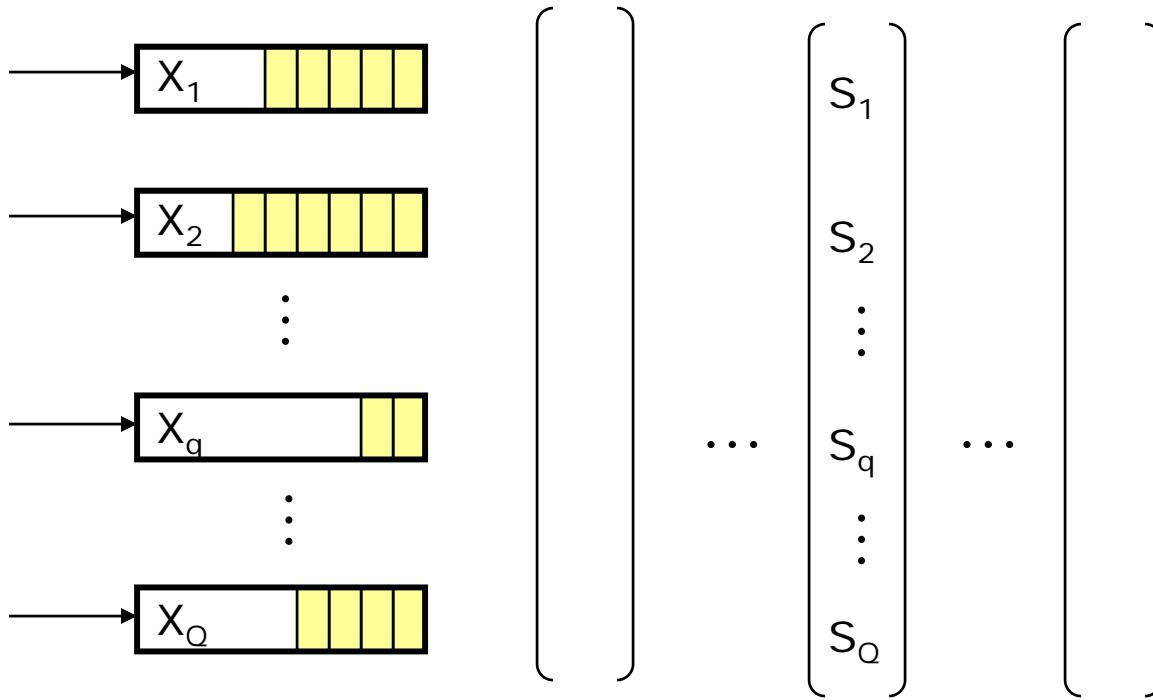
Level 3: algorithms and operations

Level 2: system architecture

Level 1: low-power circuit design

The Basic Model ...

A Canonical Model of Interacting Resource Allocation



X = backlog vector

S = service vector/mode/configuration/allocation

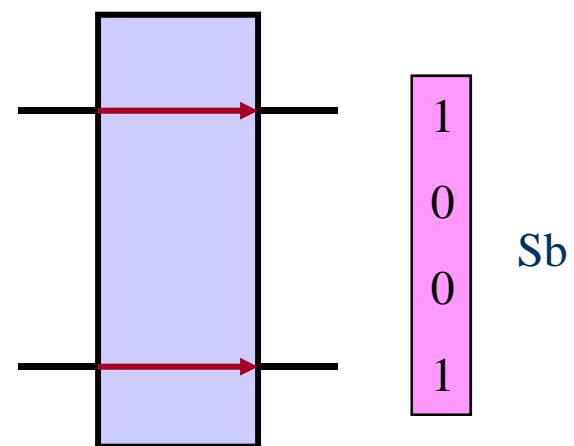
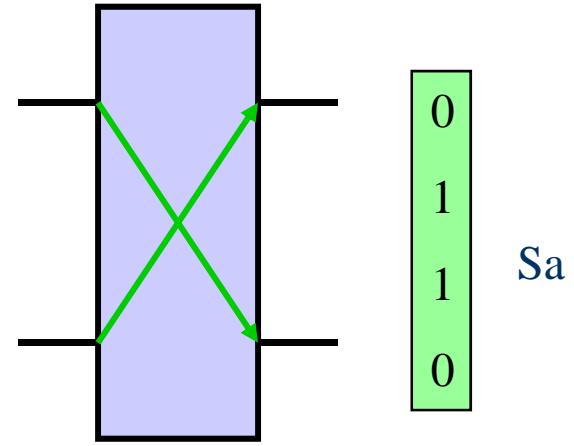
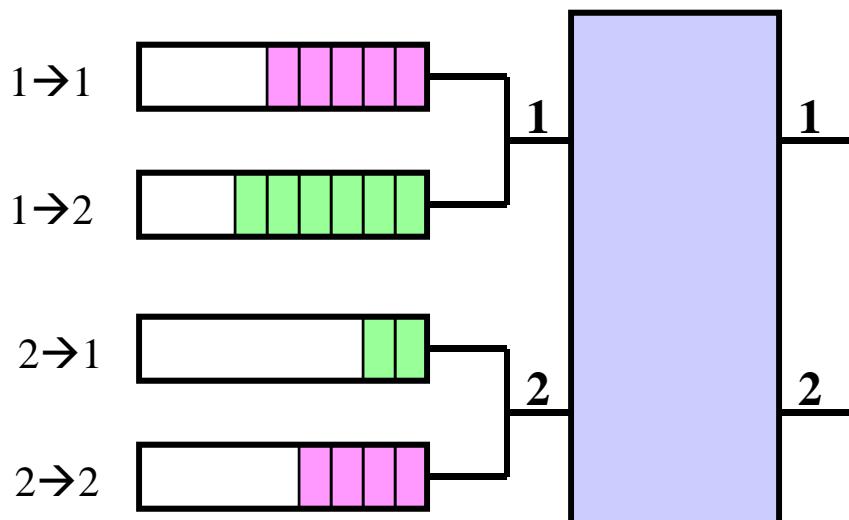
\mathcal{S} = set of all possible service vectors

P_S = cost (power...) of service configuration S

Core Issue... dynamically choose S (based on backlog/service history, etc.)
to maximize throughput, minimize (power/delay) cost...

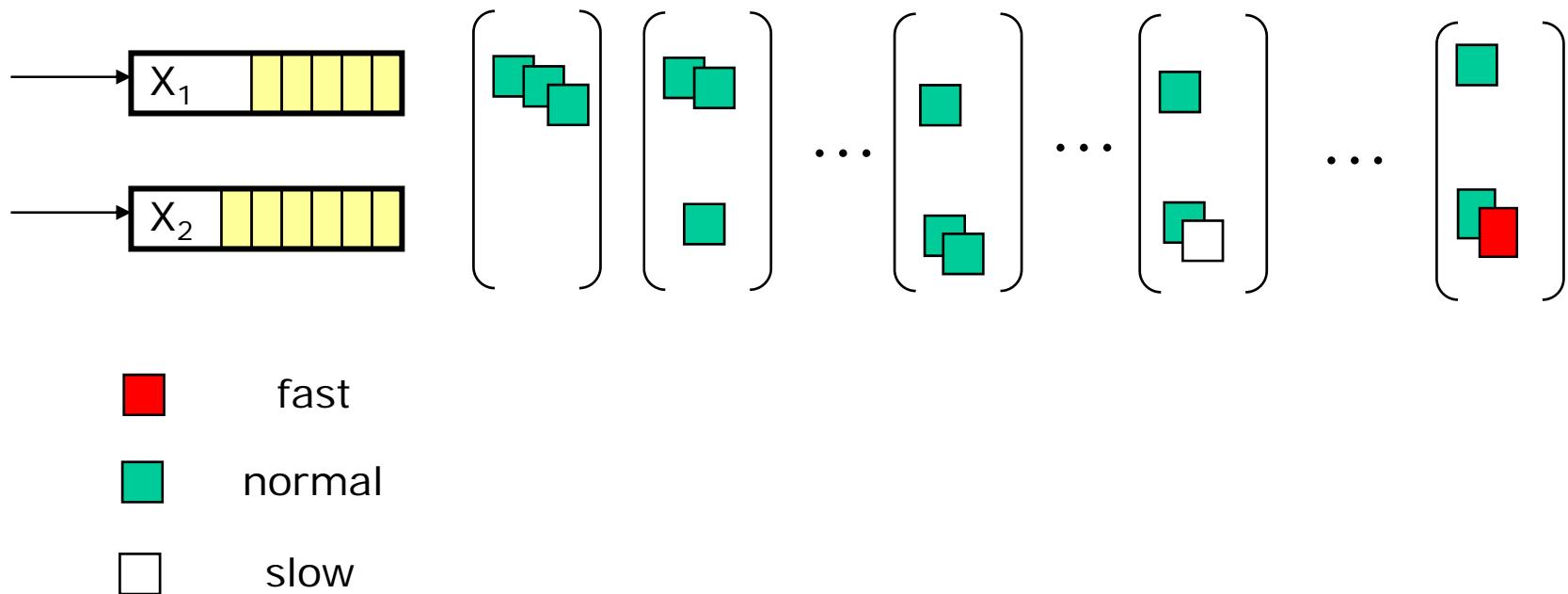
Input Queued Packet Switches ...

2x2 switch ... simplest model (...not simplistic)... scales to NxN



Multiprocessors ...

Allocating 3 processors to two job queues



... in general ... *any set* of service vectors

Throughput ...

Load and Throughput...

Arbitrary traffic traces

Load $\rho = (\rho_1, \rho_2, \dots, \rho_q, \dots, \rho_Q)$... long term avg. packet load

$$\rho_q = \frac{\text{packet load arriving to queue } q \text{ in } (0, t)}{t} \dots \text{as } t \text{ grows large}$$

Throughput job departure rate = job arrival rate

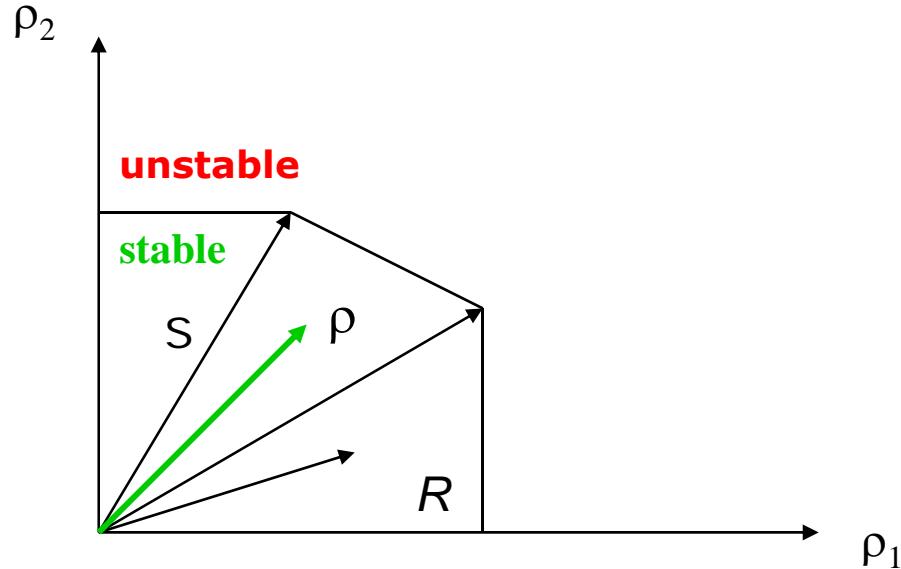
... in/out flow balance

... **flow conservation**

... stability

$$\frac{X(t)}{t} \longrightarrow 0 \dots \text{at large times } t$$

Throughput ... Capacity Region



$$R = \{ \rho : \rho \leq \sum_{S \in S} \phi_S S \dots \text{ for some } \phi_S > 0 \text{ with } \sum \phi_S = 1 \}$$

Projective Cone Schedules (PCS) Maximize Throughput

PCS algorithm... when backlog X , choose S to maximize projection on $\mathbf{B}X$

$$\max \langle S, \mathbf{B}X \rangle \quad \text{over } S \text{ in } S$$

maximizes throughput

for *any* fixed matrix \mathbf{B} that is

positive-definite, symmetric and has
negative/zero off-diagonal elements

MWM algorithm ... PCS with $B=I$

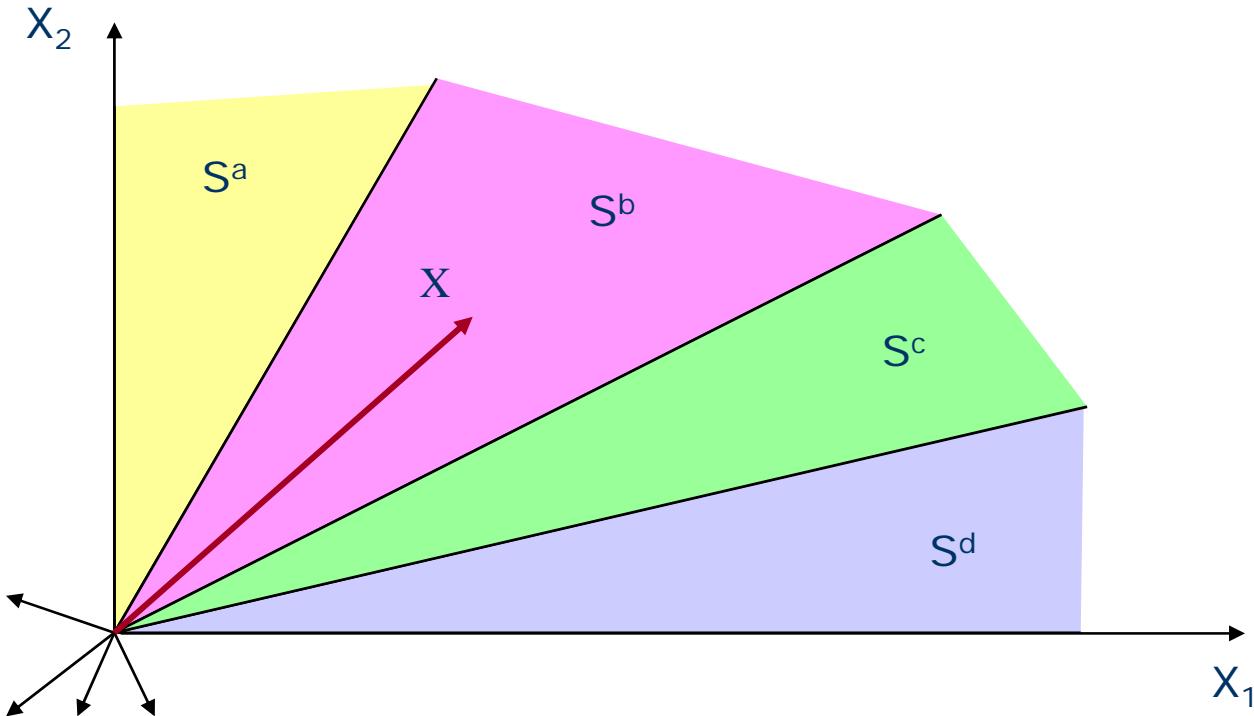
Rich family of schedules... (~ Q^2 matrix parameters to tweak and tune)

Extremely **robust** schedules

Interesting geometric intuition

Geometry of Projective Cone Schedules (PCS)

Given X , choose S to maximize $\langle S, BX \rangle$ over all S in S



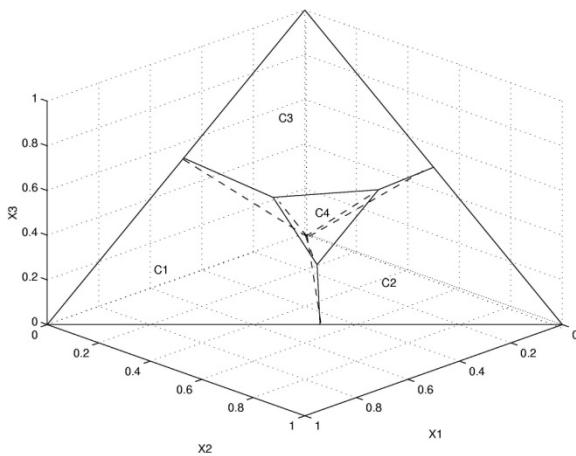
When X in cone C ,

choose $S = S(C)$ corresponding to that cone

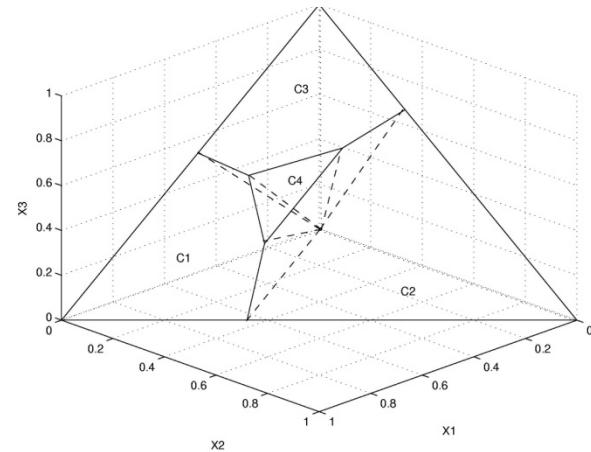
A 3-Queue Example

$$S1=(9,0,0) / S2=(0,8,0) / S3=(0,0,8) / S4=(3,4,3)$$

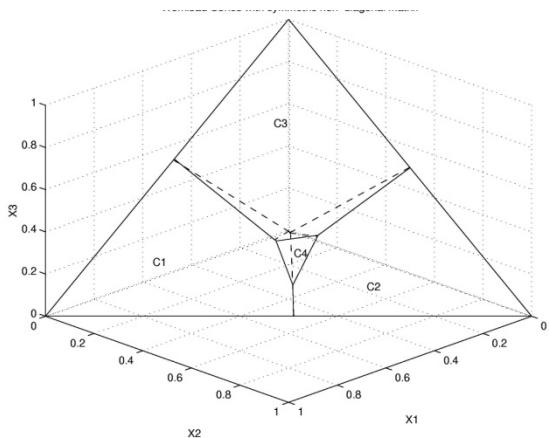
$$\mathbf{B}=[1,0,0; 0,1,0; 0,0,1]$$



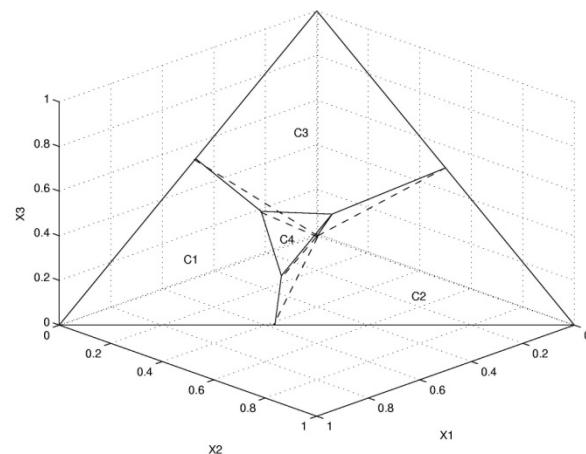
$$\mathbf{B}=[1,0,0; 0,2,0; 0,0,1]$$



$$\mathbf{B}=[1,-0.5,0; -0.5,1,0; 0,0,1]$$

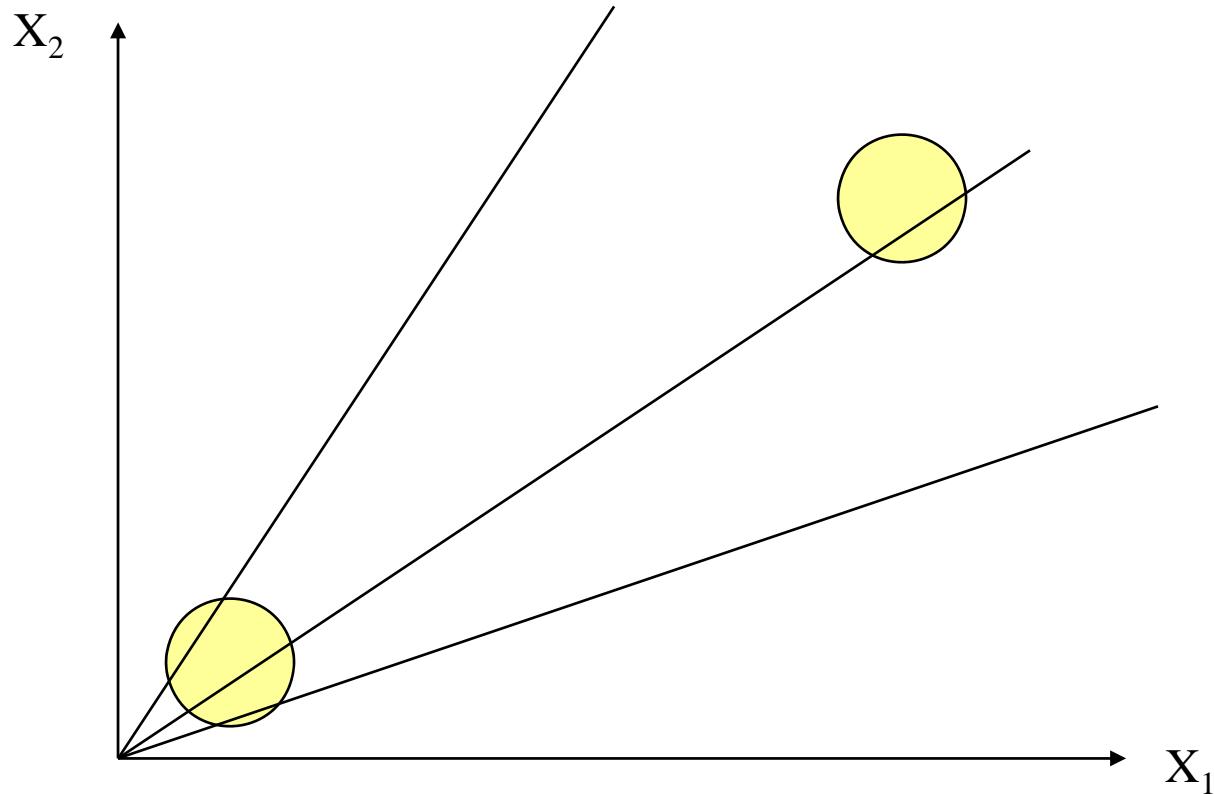


$$\mathbf{B}=[1,-0.5,0; 0,1,0; 0,0,1]$$



Scalable Projective Schedules – Local PCS

Assume bound on “workload displacement/jump”

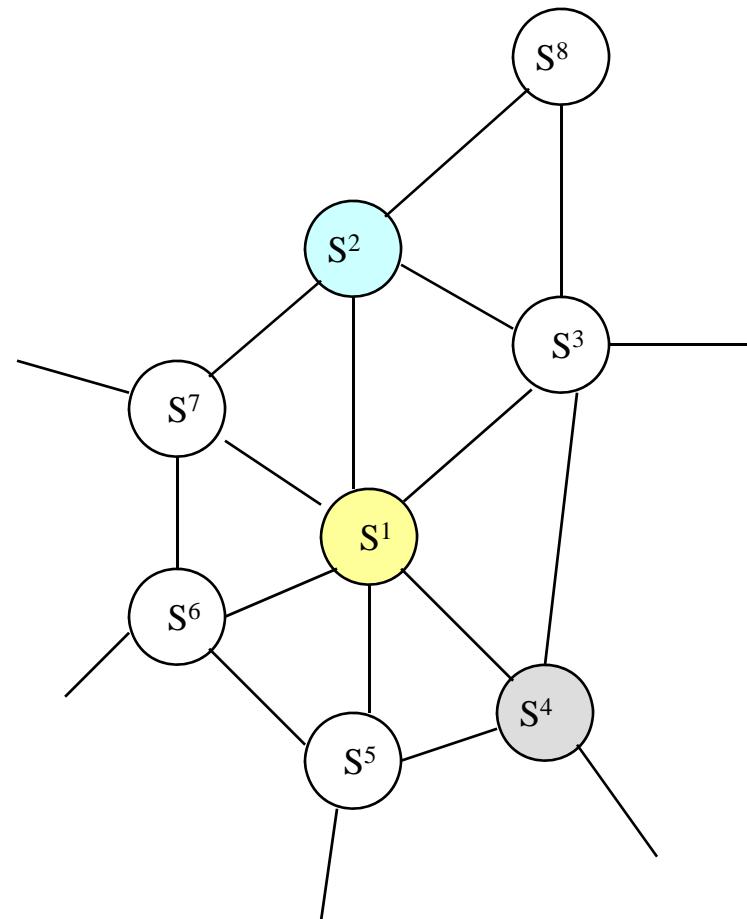
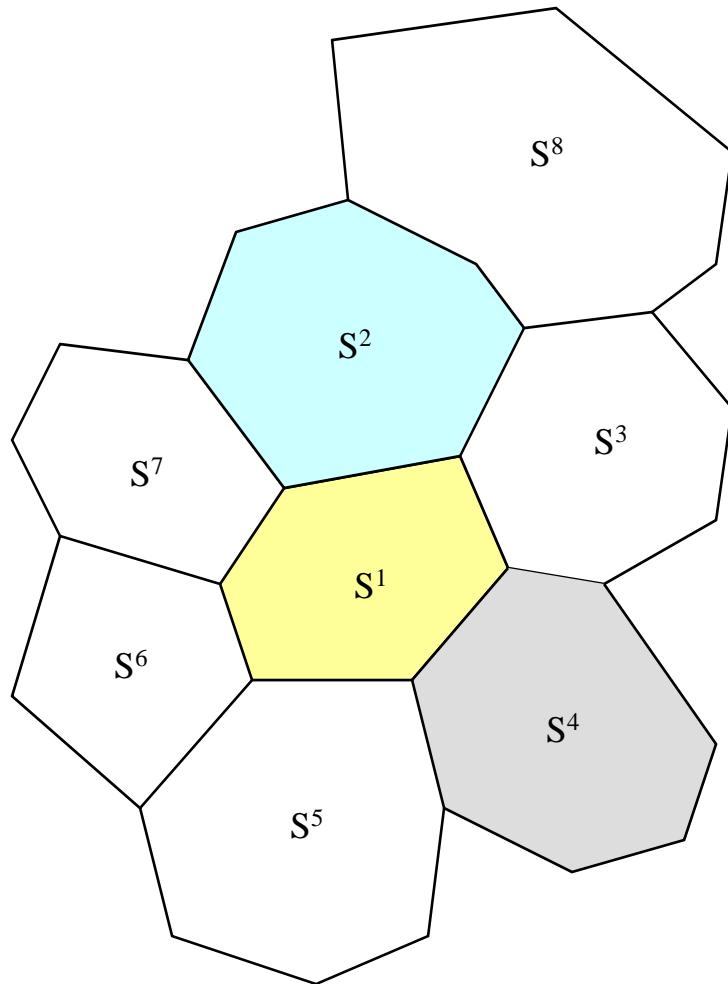


Have to search only neighbor cones ... fewer as workload grows! ... **Local Search**

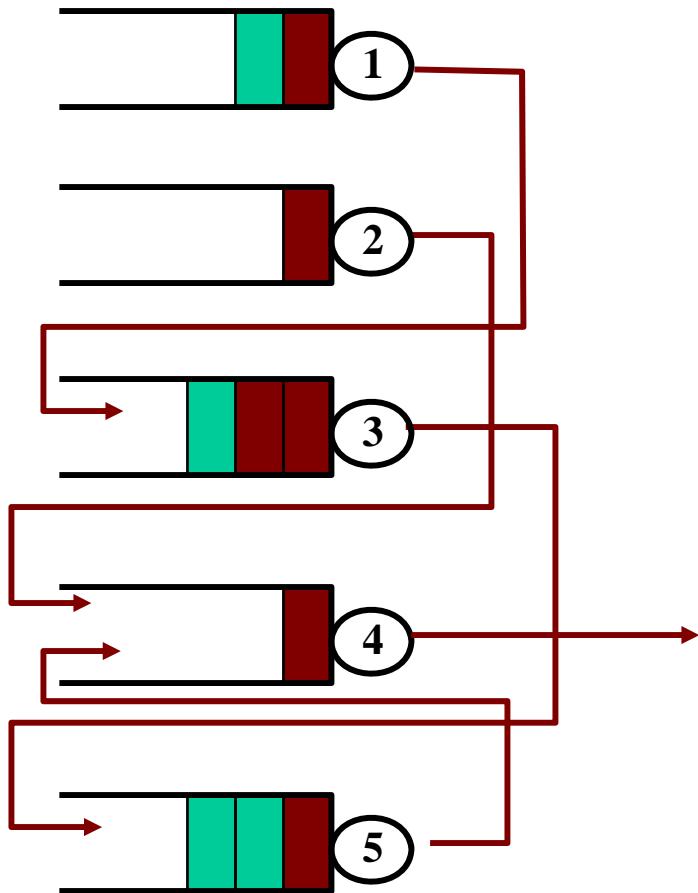
Graph Representation of PCS ... Local Search Idea ...

When at S, check out $\langle S, BX \rangle$ on **neighbor nodes** (only) ...

... chose S' with maximal $\langle S, BX \rangle$... **steepest ascent**



Networks ... and PCS

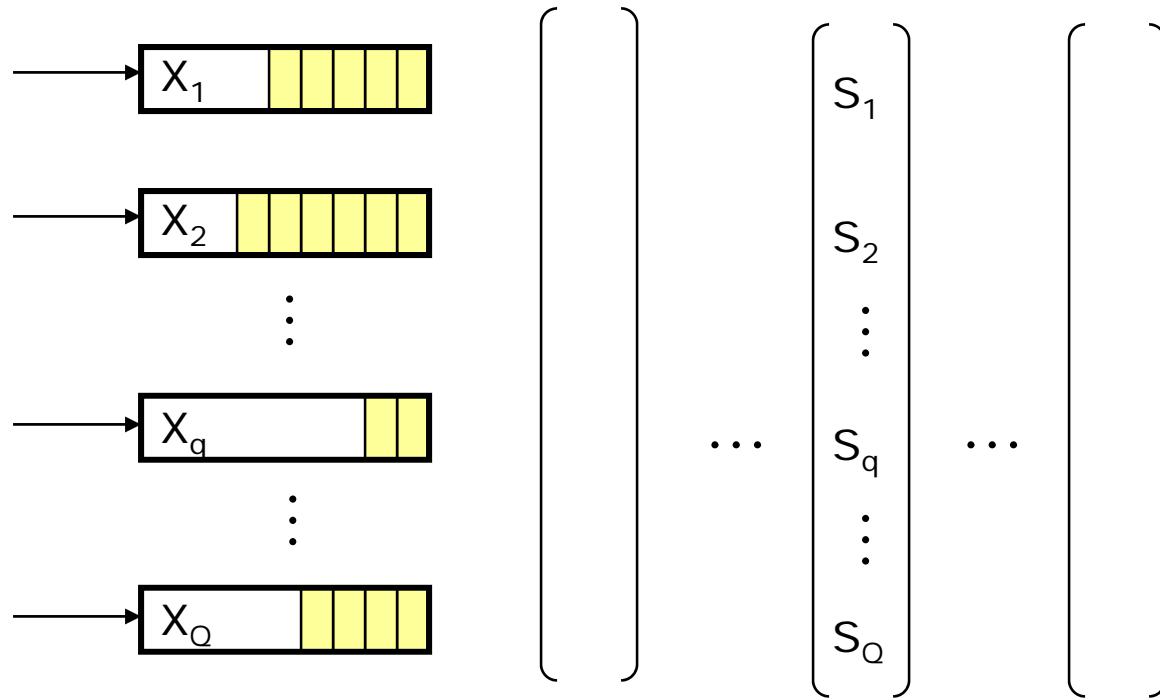


$$\begin{pmatrix} 1 \\ 1 \\ 1 \\ -1 \\ -1 \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \\ 2 \\ 1 \\ 1 \end{pmatrix} - \begin{pmatrix} 0 \\ 0 \\ 1 \\ 2 \\ 2 \end{pmatrix}$$

Service = Departures - Arrivals

Power ...

The Canonical Model Again... Power vs. Delay



X = backlog vector

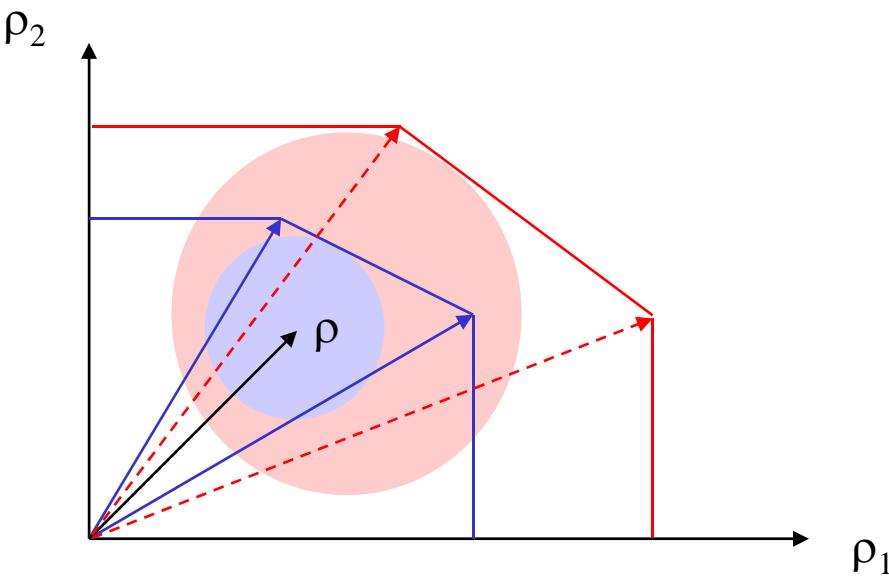
S = service vector/mode/configuration/allocation

S = set of all possible service vectors

P_S = cost (power...) of service configuration S

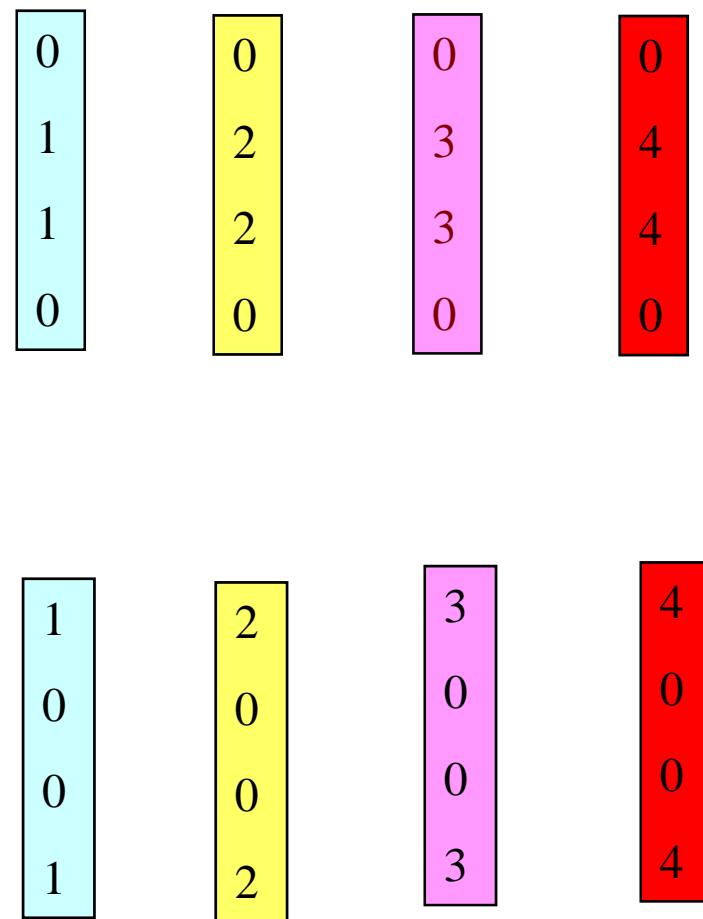
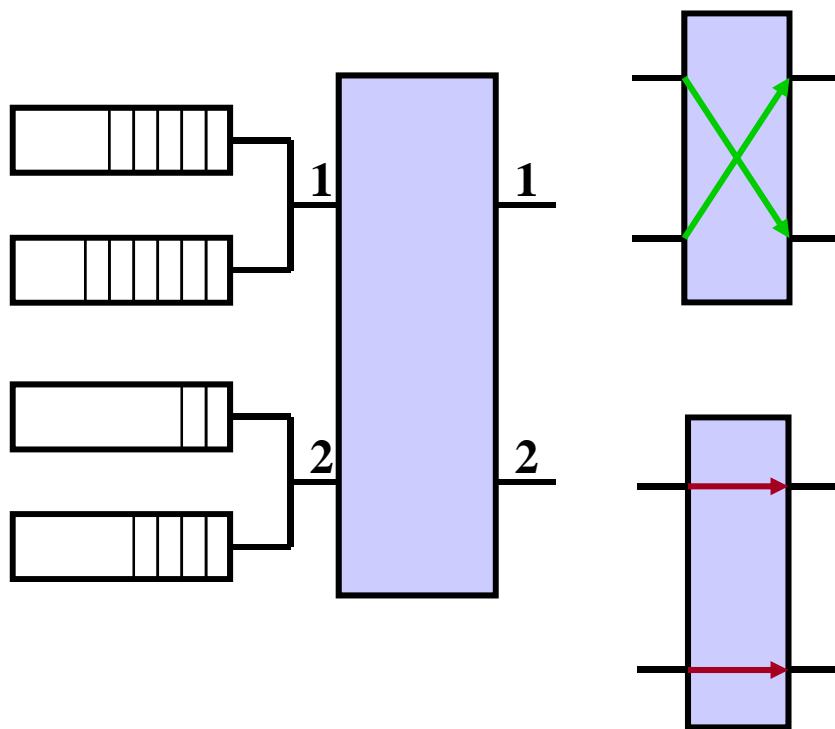
Core Issue... dynamically choose S to **minimize power + delay cost**

dynamic programming formulation ... intractable...



Activating faster/slower, higher/lower-power, more/less expensive service vectors
... adjusting the capacity space

Packet Switches ... with Deceleration Modes



Speed Mode m: slower/cool

...

faster/hotter

The Control Problem – Riccati Recursion

Service Vector \mathbf{S} = Speed Mode \mathbf{m} x Port-Matching Configuration \mathbf{C}

In each slot, have to choose \mathbf{m} & \mathbf{C} ... $\mathbf{S}=\mathbf{m}\mathbf{C}$

Backlog Cost... $\mathbf{X}'\mathbf{E}\mathbf{X}$...per slot ... \mathbf{E} pos-def & sym

Power Cost ... $\mathbf{S}'\mathbf{F}\mathbf{S}$...per slot ... \mathbf{F} pos-def & sym

No arrivals...

$J(\mathbf{X}) = \min [\mathbf{X}'\mathbf{E}\mathbf{X} + \mathbf{S}'\mathbf{F}\mathbf{S} - J(\mathbf{X}-\mathbf{S}/\text{dept})]$ min over \mathbf{S}

Continuous Relaxation ... **Linear Quadratic Regulator (LQR)...** **Riccati Equation**

$\mathbf{S}^*(\mathbf{X}) = [(\mathbf{F}+\mathbf{P})^{-1} \mathbf{P}] \mathbf{X}$

Where \mathbf{P} solves $\mathbf{E} = \mathbf{P} (\mathbf{F}+\mathbf{P})^{-1} \mathbf{P}$

Power-Aware MWM... PA-MWM

Reduced Backlog Cost... $X'EX$ $\sim b \sum [X_q^2]$

Reduced Power Cost ... $S'FS$ $\sim f [\sum S_q]^2$

PA-MWM algorithm

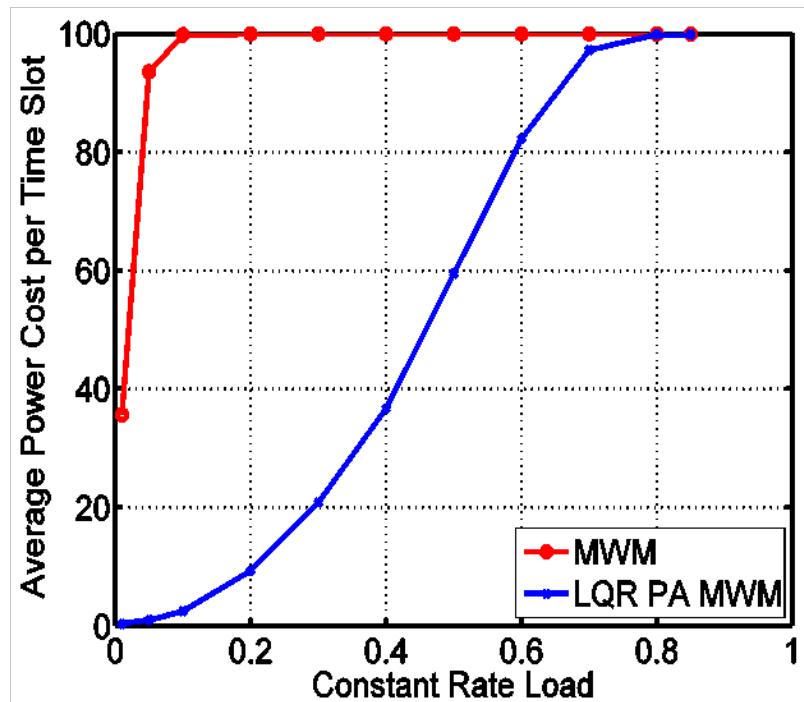
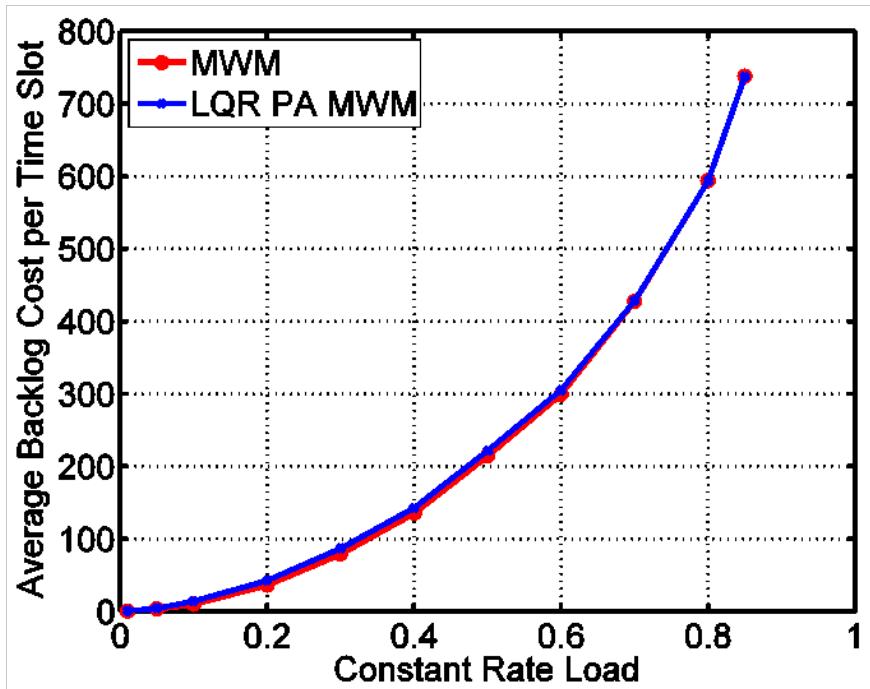
alignment

Choose Configuration $C^*(X) = \operatorname{argmax} \langle C, X \rangle$

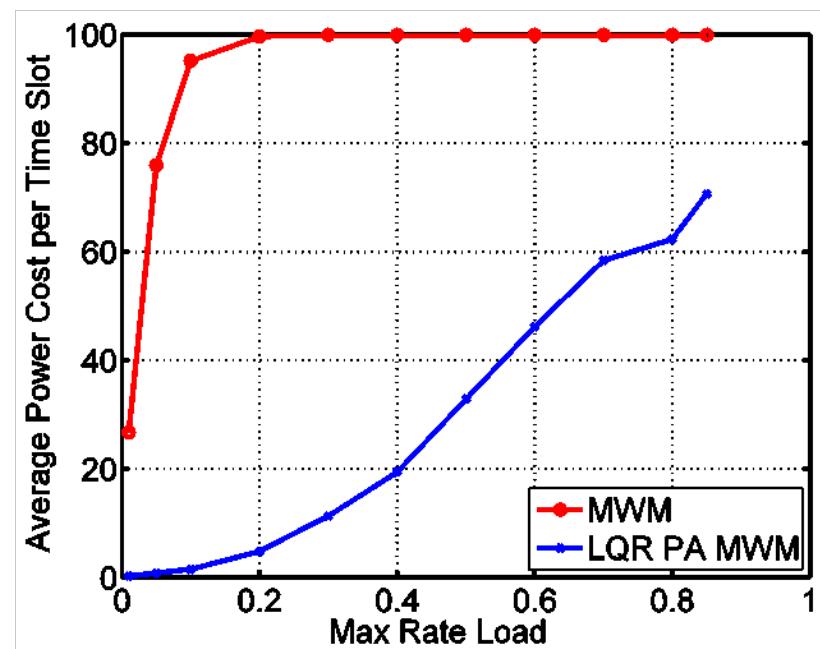
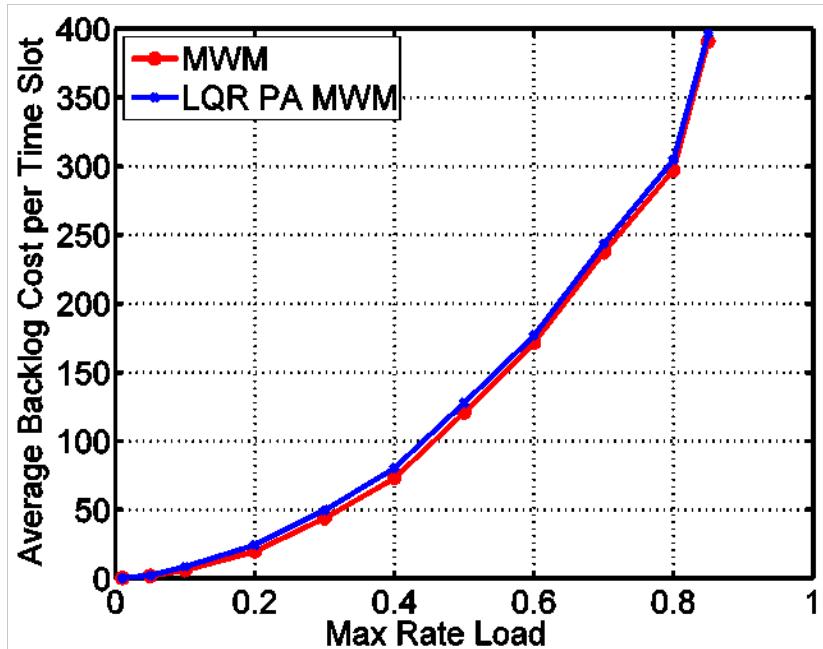
Choose Speed Mode $m^*(X) = \operatorname{argmin} |m - a \text{ TotalBacklog}|$

congestion

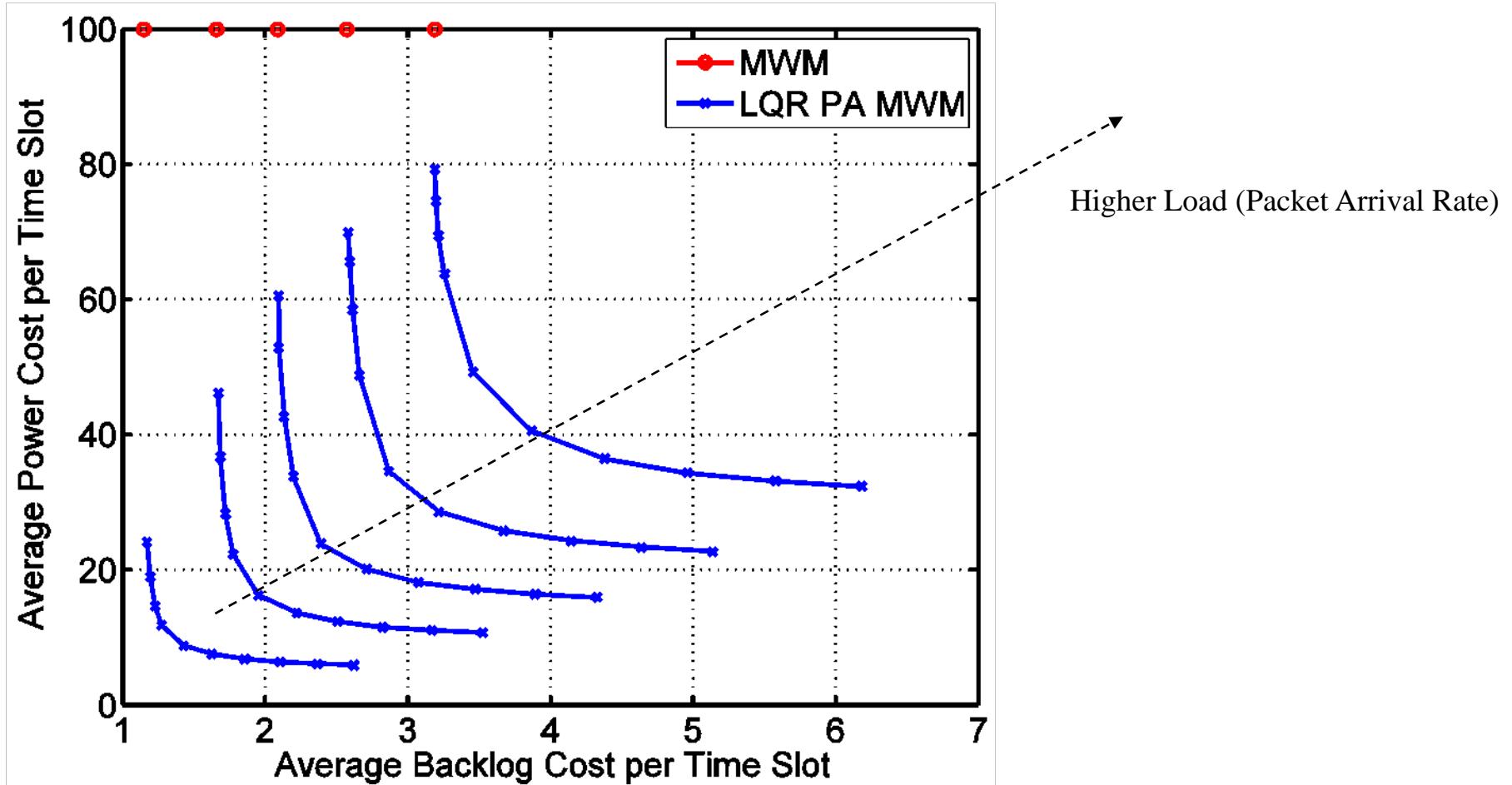
PA-MWM I



PA-MWM II



PA-MWM II



Total Cost = BacklogCost + a PowerCost ... vary 'a'

Thank You!