# The Teleology of Switched Networks

Damon Wischik, UCL

Devavrat Shah, MIT

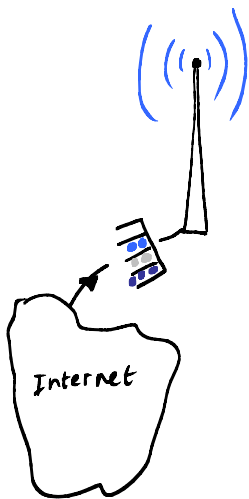# AN INPUT-QUEUED SWITCH



virtual output queues

input port 1

input port 2

input port 3

output port 1

output port 2

output port 3

Arriving packets are classified by destination.

Each timestep, the scheduling algorithm chooses a matching from inputs to outputs.
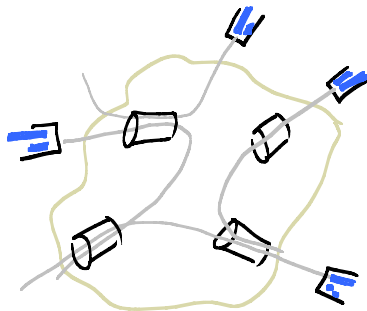
# A WIRELESS BASE STATION



Packets arrive from the Internet, and are queued up at the base station.

Each timestep, the base station chooses one queue to serve.

The service rate depends on the state of nature; the base station knows the current state of nature, and its distribution.

[31] STOLYAR, A. L. (2004). MaxWeight scheduling in a generalized switch: State space collapse and workload minimization in heavy traffic. *Annals of Applied Probability* 14, 1, 1–53.

# FLOW-LEVEL MODEL OF TCP



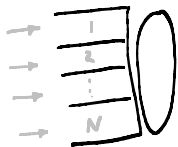Each queue represents the number of active flows on a given route.

TCP chooses transmit rates for each of the flows, a function of the number of active flows on each route.

That is, it chooses drain rates for each queue, subject to network capacity constraints.

[17] KELLY, F. P. AND WILLIAMS, R. J. (2004). Fluid model for a network operating under a fair bandwidth-sharing policy. *Annals of Applied Probability* **14**, 1055–1083.

[15] KANG, W. N., KELLY, F. P., LEE, N. H., AND WILLIAMS, R. J. (2009). State space collapse and diffusion approximation for a network operating under a fair bandwidth sharing policy. *Annals of Applied Probability*. http://www.math.ucsd.edu/~williams/bandwidth/kklw.html.

# GENERAL QUEUEING MODEL

Consider a collection of $N$ queues, $\underline{Q} = (Q_1, \cdots, Q_N)$ in slotted time.

Each timestep, work arrives, eg Bernoulli arrivals. Each queue has a dedicated arrival process; arrivals are at rate $\underline{\lambda}$

Each timestep, the switch chooses a schedule $\pi(t)$ from a finite set $S \subset \mathbb{R}_+^N$.

It picks $$\pi(t) = \underset{p \in S}{\arg\max} \; p \cdot \underline{Q}^\alpha(t)$$

$\text{i.e.} \sum_{1 \le n \le N} p_n \, Q_n^\alpha(t)$

for some fixed $\alpha > 0$.

# PRIMAL($\lambda$)

minimize $\quad \sum_{\pi \in S} \alpha_\pi$

over $\qquad \alpha \geq 0$

such that $\quad \lambda \leq \sum_{\pi \in S} \alpha_\pi \pi$

This is what an omniscient scheduler might solve: find a mixture of schedules that can serve the arriving traffic.

If PRIMAL($\lambda$) $\leq 1$, this is possible.

# DUAL($\lambda$)

maximize $\quad \xi \cdot \lambda$

over $\qquad \xi \in \mathbb{R}_+^N$

such that $\quad \xi \cdot \pi \leq 1 \quad \text{for all } \pi \in S.$

$\xi \cdot Q(t)$ represents a virtual queue; it gets $\xi_n$ tokens whenever a packet arrives at queue $n$. The maximum possible drain rate of tokens is $1$ / timeslot. On average, $\xi \cdot \lambda$ arrive per timeslot.

If DUAL($\lambda$) $> 1$, the switch must be unstable.

# THE FLUID SCALE

From the discrete quantities

$Q(t)$

queue sizes
at time t

$A(t)$

arrivals over
time interval
$\{1, 2, \ldots, t\}$

$S_\pi(t)$

number of timesteps
in which $\pi$ was
chosen, in that interval

$Z(t)$

amount of idling
at each queue
in that interval

we define fluid-scaled quantities

$$q^r(t) = \frac{Q(rt)}{r} \qquad a^r(t) = \frac{A(rt)}{r} \qquad s^r_\pi(t) = \frac{S_\pi(rt)}{r} \qquad z^r(t) = \frac{Z(rt)}{r}$$

When $r$ is large, these scaled-quantities lie close to continuous functions $q(\cdot), a(\cdot), s_\pi(\cdot), z(\cdot)$ that satisfy differential equations.

The results in this talk are proved by analysing these diff. eq.

## Assume the arrival process satisfies

$$\mathbb{P}\left(\sup_{\tau \leq r}|\mathbf{A}(\tau) - \lambda\tau| < \varepsilon r\right) = 1 - o(R(r)) \quad \text{as } r \to \infty, \quad \text{for all } \varepsilon > 0.$$

## Specify the differential equations:

(12) $\quad \mathbf{a}(t) = \lambda t$

(13) $\quad \mathbf{q}(t) = \mathbf{q}(0) + \mathbf{a}(t) - \sum_{\pi} s_\pi(t)\pi + \mathbf{z}(t)$

(14) $\quad \sum_{\pi \in \mathcal{S}} s_\pi(t) = t$

(15) $\quad$ each $s_\pi(\cdot)$ and $z_n(\cdot)$ is increasing (not necessarily strictly increasing)

(16) $\quad$ all the components of $x(\cdot)$ are absolutely continuous—
indeed they are Lipschitz

(17) $\quad$ for almost all $t$, all $n$, $\quad \dot{z}_n(t) = 0$ if $q_n(t) > 0$

(18) $\quad$ for almost all $t$, all $\pi \in \mathcal{S}$, $\quad \dot{s}_\pi(t) = 0$ if $\pi \cdot f(\mathbf{q}(t)) < \max_{\rho \in \mathcal{S}} \rho \cdot f(\mathbf{q}(t))$

## Then, $x^r = (q^r, a^r, s^r, z^r)$ is close to a solution of those equations:

**Theorem 5.1** [3] *Make the above assumptions* (3)–(5) *and* (21)–(24). *Let FMS be the set of all processes* $x(t)$ *over* $t \in [0,T]$ *which satisfy the appropriate fluid model equations, namely*

- *equations* (12)–(17), *for any scheduling algorithm,*
- *equation* (18) *in addition if the network is running MW-f and Condition* 4.1 *holds,*
- $\mathbf{q}(0) = \mathbf{q}_0$ *in addition, if* (25) *holds.*

*Let FMS$_\varepsilon$ be the $\varepsilon$-fattening*

$$FMS_\varepsilon = \left\{x : \sup_{t \in [0,T]}|x(t) - y(t)| < \varepsilon \text{ for some } y \in FMS\right\}.$$

*Then for any $\varepsilon > 0$, $\mathbb{P}(x^r(\cdot) \in FMS_\varepsilon) = 1 - o(R(r'))$ as $r \to \infty$.*

# STABILITY

Suppose $PRIMAL(\lambda) < 1$ hence $\lambda \le \sum \alpha_\pi \pi$ where $\sum \alpha_\pi < 1$

Define the Lyapunov function $L(q) = \left(1 \cdot q^{1+\alpha}\right)^{1/1+\alpha}$

**THEOREM.** $\dfrac{d}{dt} L(q_t) \le -\eta < 0$ if $q_t \ne 0$.

**COROLLARY.** If arrivals are IID, so that the queue size process is a Markov chain, then it is +ve recurrent.

## ALGP($q$)

minimize $\quad L(r)$

over $\quad r \in \mathbb{R}_+^N$

such that $\quad r \geq q + t(\lambda - \sigma)$

for some $\sigma \in \langle \mathfrak{s} \rangle$

convex hull of $\mathfrak{s}$

## ALGD($q$)

minimize $\quad L(r)$

over $\quad r \in \mathbb{R}_+^N$

such that $\quad \mathfrak{z} \cdot r \geq \mathfrak{z} \cdot q$
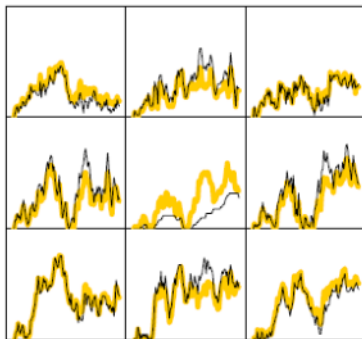
for all $\mathfrak{z} \in \mathfrak{z}^*(\lambda)$

extreme optimal solutions to DUAL($\lambda$)

These problems both ask: what is the "cheapest" way to arrange the work in the queues, given that we started at $q$?

The problems are identical, and have a unique solution, $\Delta(q)$.

In a critically-loaded switch, the state is nearly always in

$$\mathcal{I} = \{ q : q = \Delta(q) \}.$$

# Example: a 3x3 input-queued switch



Legend:
- $Q(t)$
- $\Delta Q(t)$

$$\lambda = \begin{pmatrix} 0.1413 & 0.4626 & 0.3910 \\ 0.4626 & 0.0055 & 0.5268 \\ 0.3910 & 0.5268 & 0.0771 \end{pmatrix}$$
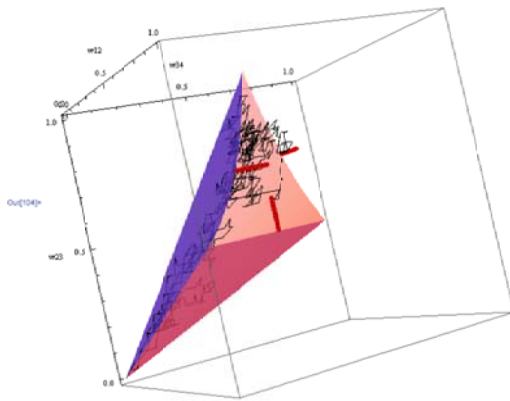
# Example: a 4-hop wireless link



let $\lambda = (\cdot 3, \quad \cdot 7, \quad \cdot 3, \quad \cdot 7)$

$$S = \left\{ \begin{array}{cccc} (1, & 0, & 1, & 0) \\ (0, & 1, & 0, & 1) \\ (1, & 0, & 0, & 0) \end{array} \right\}$$

$$S^* = \left\{ \begin{array}{cccc} (1, & 1, & 0, & 0), \\ (0, & 0, & 1, & 1), \\ (0, & 1, & 1, & 0) \end{array} \right\}$$

In all the examples I've looked at, the bigger $\alpha$ the better the performance.

Counter-Examples to the Optimality of MWS-α Policies for Scheduling in Generalized Switches

Tianxiong Ji, Eleftheria Athanasopoulou, and R. Srikant

## ALGD$^+$

minimize $\quad L(q)$

over $\qquad q \in \mathbb{R}_+^N$

such that $\quad \xi \cdot q \geq \xi \cdot \lambda - 1$

$\qquad$ for all dual-extreme $\xi$
$\qquad$ such that $\xi \cdot \lambda > 1$

## RATE$^+$

maximize $\quad \lambda \cdot q^* - \max_\pi \pi \cdot q^*$

over $\qquad q \in \mathbb{R}_+^N$

such that $\quad L(q) = 1$

let $q^+$ be the unique solution to ALGD$^+$.

**THEOREM** If the switch is <span style="color:red">overloaded</span> then $\dfrac{q(t)}{t} \longrightarrow q^+$ for any fluid model solution $q(t)$.

The ALGD$^+$ problem is: find the growth rate for $q(t)$ that incurs cost (measured by $L(q(t))$) as slowly as possible.

[8] EGOROVA, R., BORST, S., AND ZWART, B. (2007). Bandwidth-sharing networks in overload. *Performance Evaluation 64*, 978–993.

# Example: 2x2 input-queued switch with $\alpha = 1$

$$\lambda^{\text{critical}} = \begin{pmatrix} 0.3 & 0.7 \\ 0.7 & 0.3 \end{pmatrix}$$

service rate $= \begin{pmatrix} 0.3 & 0.7 \\ 0.7 & 0.3 \end{pmatrix}$

$q(t) = \begin{pmatrix} 0 & 0 \\ 0 & 0 \end{pmatrix}$

departure rate $= 2$

$$\lambda^{\text{overload}} = \begin{pmatrix} 0.3 & 1.0 \\ 0.7 & 0.3 \end{pmatrix}$$

service rate $= \begin{pmatrix} 0.2 & 0.8 \\ 0.8 & 0.2 \end{pmatrix}$

$q(t) = \begin{pmatrix} 0.1t & 0.2t \\ 0 & 0.1t \end{pmatrix}$

departure rate $= 1.9$

This switch has the paradoxical behaviour
of serving _less_ traffic the greater the arrival rate.

# HEURISTIC FOR LDP FOR UNDERLOADED SWITCH

$$\mathbb{P}\Big( L(Q(rT)) \approx r \Big) = \mathbb{P}\Big( L(q^r(T)) \approx 1 \Big) \quad \text{since } L \text{ linear}$$

$$\approx \mathbb{P}\Big( a^r(\cdot) \in \Big\{ a(\cdot): \begin{array}{l} \text{queue fed by} \\ a(\cdot) \text{ reaches} \\ L(q(T))=1 \end{array} \Big\} \Big) \quad \text{rearranging}$$

$$\approx \sup_{a(\cdot):\, L(q(T))=1} \mathbb{P}\Big( a^r(\cdot) \approx a \Big) \quad \begin{array}{l} \text{Principle of the} \\ \text{Largest Term} \end{array}$$

$$\frac{1}{r}\log\mathbb{P}\Big( L(Q(rT)) \approx r \Big) \approx \sup_{a(\cdot):\, L(q(T))=1} \frac{1}{r}\log\mathbb{P}\Big( a^r(\cdot) \approx a \Big) \quad \text{rearranging}$$

$$\approx -\inf_{a(\cdot):\, L(q(T))=1} \int_{t=0}^{T} \ell\big( \dot{a}(t) \big)\, dt \quad \begin{array}{l} \text{LDP for} \\ \text{arrival process,} \\ \text{assuming IID} \\ \text{increments.} \end{array}$$

[10] Vijay G. Subramanian, Tara Javidi, and Somsak Kittipiyakul. Many-sources large deviations for max-weight scheduling. arXiv:0902.4569v1, 2009.

[11] V. J. Venkataramanan and Xiaojun Lin. Structural properties of LDP for queue-length based wireless scheduling algorithms. In *Proceedings of Allerton*, 2007.

$$\frac{1}{r} \log \mathbb{P}\left( L(Q(rT)) \approx r \right) \quad \approx \quad - \inf_{a(\cdot)\,:\, L(q(T))=1} \int_{t=0}^{T} \ell\left( \dot{a}(t) \right) dt$$

The most likely path is:
send at mean rate $\lambda$ for time $T-U$;
then send at some overload rate $\lambda^+$ for time $U$.

Choose $U = \dfrac{1}{ALGP^+(\lambda^+)}$ ie enough time to reach $L(q(T)) = 1$.

The cost of the path is

$$= \quad - \inf_{\lambda^+ \in \mathbb{R}_+^N} \quad \frac{\ell(\lambda^+)}{ALGP^+(\lambda^+)}.$$

The proof uses $ALGP^+$ and its dual $RATE^+$.

# DESIGNING OPTIMAL ALGORITHMS

Suppose we want to keep the total queue size $Q.1$ small.
We only know how about $L(Q)$; but happily

$$\frac{1}{N^{\alpha/1 \text{im}}} \; Q.1 \; \leq \; L(Q) \; \leq \; Q.1$$

This bound translates into performance bounds.
For example, in overload,

**Theorem 12.6** *Let $\lambda \notin \Lambda$. There is some $q^{min} = q^{min}(\lambda) > 0$ such that for any scheduling algorithm, all fluid model solutions satisfy*

$$q^{min} \leq \mathbf{1} \cdot \mathbf{q}(t)/t \quad \text{for all } t \geq 0.$$

*Furthermore, for the MW-$\alpha$ scheduling algorithm,*

$$\lim_{t \to \infty} \mathbf{1} \cdot \mathbf{q}(t)/t \leq N^{\alpha/(1+\alpha)} q^{min}.$$

Hence the smaller $\alpha$ is, the closer we are to an
optimal algorithm.

# FUTURE WORK

- Multiscale phenomena

  For some scheduling algorithms, e.g. choose $\pi(t)$ to be

  $$\underset{p \in S}{\arg\max} \sum_i p_n \log(1 + a_n)$$

  it is hard to identify the fluid model. The fluid model seems to depend on higher-order statistics of the arrival process.

- Distributed scheduling algorithms

  What distributed algorithms can we analyse/design, e.g. for wireless ad-hoc networks?