Control of Fork-Join Networks in Heavy-Traffic

Asaf Zviran

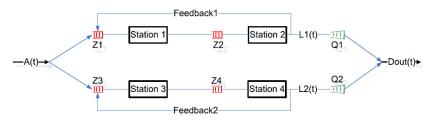
Based on MSc work under the guidance of Rami Atar (Technion) and Avishai Mandelbaum (Technion)

Industrial Engineering and Management Technion

June 2010

Network Model

We shall consider a class of Fork-Join networks that include probabilistic feedback. For example-

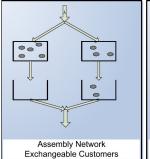


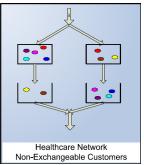
<u>Motivation</u> - Fork-Join networks are natural models for a variety of processes including communication and computer systems, manufacturing, project management and health-care.

Customers and Activities

In our model, activities are associated uniquely with customers. They are hence non-exchangeable in the sense that one can not join activities associated with different customers.

This property create dependencies between processing routes

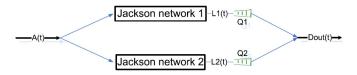




Conclusion- Customers' disorder may cause increase of Idle-time in the join nodes and hence lower throughput.

Control Problem Formulation

Consider a general Fork-Join + Jackson network. such as



We seek an optimal priority policy, in the sense of maximum throughput, under the following assumptions.

Model Assumptions

- \bullet Time-Homogeneous Poisson arrival process with average arrival rate $\lambda.$
- ullet Exponential service durations with average rate of μ_i for all servers in station j.

Policy Assumptions

- Nonanticipating.
- Work conserving.

Exact Optimality

Definition- Maximum throughput in the sense of maximum achievable departures on any finite region [0, T].

Proposition

Two equivalent conditions:

- Complementarity condition- $Q_1(T) \wedge Q_2(T) = 0$ a.s.
- Minimal buffers size- $Q_1(T) + Q_2(T) = |L_1(T) L_2(T)|$ a.s.

for any fixed T.

<u>Note-</u> These conditions are analogous to Assembly networks with exchangeable customers dynamics.

Asymptotic Optimality

Definition- Maximum achievable throughput in the Heavy-Traffic asymptotic.

<u>Notation</u>- Throughout the presentation we shall use the scaling $\hat{Q}_i^n(t) = \frac{Q_i^n(t)}{\sqrt{n}}$.

Proposition

Equivalent condition:

$$\hat{Q}_1^n(\mathit{T}) \wedge \hat{Q}_2^n(\mathit{T}) o 0$$
 in probability

for any fixed T.

One can verify that under the condition above, given any other control policy j and fixed $\mathsf{T}\colon$

$$\hat{Q}_1^n(T)+\hat{Q}_2^n(T)\leq \hat{Q}_1^{(j),n}(T)+\hat{Q}_2^{(j),n}(T)+\epsilon(n), \ \ \epsilon(n)\to 0 \ \text{in probability} \ .$$

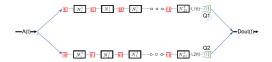
Is FCFS Optimal for Multi-Servers Network?

One can verify that FCFS is optimal for Single-Servers Fork-Join network. 1

But is this the case for Multi-Servers Fork-Join?

Obviously No, Due to the disorder effect of the multi-server stations. But the FCFS policy may be asymptotically optimal.

Consider the following system:



Heavy Traffic assumptions: ("sequence of systems" indexed by n)

- arrival rate: $\lambda^n = \lambda \cdot n + \hat{\lambda} \cdot \sqrt{n} + 1.0.t.$
- service rate: $\mu_i^{j,n} = \mu_i^j \cdot n + \hat{\mu}_i^j \cdot \sqrt{n} + \text{l.o.t.}$
- traffic intensity: $ho_i^{j,n} \equiv rac{\lambda^n}{N^j \cdot u^{j,n}}$
- Heavy Traffic Condition: $n^{1/2}(\rho_i^{j,n}-1) \longrightarrow_n \theta_i^j$, as $n \longrightarrow \infty$, $|\theta_i^j| < \infty \ \forall i,j$.

Asymptotically Optimal Control

Theorem

Given the Multi-Servers system above, with FCFS discipline, and fixed T.

$$P(Q_1^n(T) \wedge Q_2^n(T) > K) \longrightarrow_n 0;$$

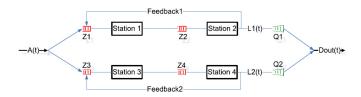
for some deterministic number K defined by the network structure.

Meaning that $\hat{Q}_1^n(T) \wedge \hat{Q}_2^n(T) \leq \frac{K}{n^{1/2}}$ w.p. converging to 1.

We may conclude that-

- The scaled random variable $\hat{Q}_1^n(T) \wedge \hat{Q}_2^n(T) \rightarrow 0$ in probability, which is asymptotically optimal under the definition above.
- The multi-server non-exchangeable "sequence of systems" is asymptotically equivalent to an exchangeable system (Assembly system).

Fork-Join Network with Feedback



This seems to be the simplest setting of a Fork-Join network where solving for optimal scheduling is hard.

Heavy Traffic assumptions- ("sequence of systems" indexed by n)

- arrival rate- $\lambda^n = \lambda \cdot n + \hat{\lambda} \cdot \sqrt{n} + \text{l.o.t.}$
- service rate- $\mu_i^n = \mu_i \cdot n + \hat{\mu}_i \cdot \sqrt{n} + \text{l.o.t.}$
- Heavy Traffic Conditions:

Treaty Traint Conditions.
$$\begin{cases}
\left|\frac{\mu_1^n}{n} - \frac{\mu_2^n}{n}\right| \longrightarrow_n 0; & \left|\frac{\mu_3^n}{n} - \frac{\mu_4^n}{n}\right| \longrightarrow_n 0; \\
\left|\frac{\lambda^n}{n} - \frac{\mu_1^n \cdot (1 - \rho_1)}{n}\right| \longrightarrow_n 0; & \left|\frac{\lambda^n}{n} - \frac{\mu_3^n \cdot (1 - \rho_2)}{n}\right| \longrightarrow_n 0;
\end{cases}$$

Control Policy

<u>Definition</u>- At each route, assign absolute preemptive priority to customers whose service was completed in the other route.

The definition of the policy creates an artificial division of the customers into two classes:

- LP (Low Priority) Customers: Customers whose service is still incomplete in both routes.
- HP (High Priority) Customers: Customers whose service was completed in one of the routes but is still incomplete in the other.

<u>Note</u>: The set of HP customers in one route is equal to the set of customers in the synchronization queue in the other.

Asymptotically Optimal Control

Theorem

Given the system and control policy defined above, and a fixed T:

$$\textit{max}_{t \in [0,T]} \{ \hat{Z}_{1,2}^{n,H}(t) \wedge \hat{Z}_{3,4}^{n,H}(t) \} \rightarrow 0, \quad \text{in probability },$$

where
$$\left\{ \begin{array}{l} \hat{Z}_{1,2}^{n,H}(t) = \hat{Z}_{1}^{n,H}(t) + \hat{Z}_{2}^{n,H}(t); \\ \\ \hat{Z}_{3,4}^{n,H}(t) = \hat{Z}_{3}^{n,H}(t) + \hat{Z}_{4}^{n,H}(t); \end{array} \right.$$

But since

$$\hat{Z}_{1,2}^{n,H}(t) = \hat{Q}_{2}^{n}(t)$$
 and $\hat{Z}_{3,4}^{n,H}(t) = \hat{Q}_{1}^{n}(t)$.

We may conclude that-

The scaled process $\hat{Q}_1^n(t) \wedge \hat{Q}_2^n(t)$ converge uniformly to 0, in probability.

About The Proof

- We prove Heavy Traffic properties for H-P customers using the queue length of H-P customers in the processing routes.
- In the proof we focus on a time interval $[\tau, \sigma]$, defined by the two random variables:

$$\begin{split} \sigma &= \inf\{t \ : \ \hat{Z}_{1,2}^{n,H}(t) \wedge \hat{Z}_{3,4}^{n,H}(t) > \epsilon\}; \\ \tau &= \sup\{t < \sigma \ : \ \hat{Z}_{1,2}^{n,H}(t) \wedge \hat{Z}_{3,4}^{n,H}(t) \leq \frac{\epsilon}{2}\}; \end{split}$$

On this time interval, both processing routes have a queue length of H-P customers greater than $\frac{\epsilon}{2}n^{\frac{1}{2}}$, and one of the routes increases during this interval by $\frac{\epsilon}{2}n^{\frac{1}{2}}$.

• A central ingredient in the proof is the use of a down-crossings technique on the random process of H-P customers' queue length.

Summary

- We introduced a natural concept of optimality for Fork-Join networks with non-exchangeable customers.
- An optimality condition was derived via an analogy to Assembly networks and proved to be efficient for a general Fork-Join + Jackson model.
- We proposed a control policy and proved asymptotic optimality for both models - with multi-servers and with feedback.
- The proofs indicate an asymptotic equivalence between non-exchangeable and exchangeable dynamics in Heavy-Traffic.
- Simulation runs for the model with feedback show improvements of 33% in sojourn time and almost 66% in synchronization time, for the proposed policy vs. FCFS in Heavy-Traffic.

Thank You