

# Patient Diversion Across Primary Health Centers Using Real Time Delay Predictors

Najiya Fatma<sup>1</sup>, Varun Ramamohan<sup>1</sup>

<sup>1</sup>Department of Mechanical Engineering, Indian Institute of Technology Delhi, Hauz Khas, New Delhi

January 2021

## Abstract

In the current work, we consider diversion of childbirth patients who arrive seeking emergency admission to public primary health centers (PHCs). PHCs are the first point of contact for an Indian patient with formal medical care, and offer medical care on an outpatient basis, and limited inpatient and childbirth care. In this context, real-time prediction of the wait time of the arriving patient becomes important in order to determine whether the patient must be diverted to another PHC or not. We study this problem using a discrete event simulation that we develop of medical care operations in two PHCs in India. We approximate the labour room service at each PHC as an M/G/1 queueing system and show how the accuracy of real-time delay predictors impacts the extent of the change in operational outcomes at each PHC. We simulate patient diversion using actual delays as well as the delay estimates generated by various delay predictors based on the state of the system such as queue-length, elapsed service time, and observed delay histories. The simulation of the diversion process also incorporates travel time between the PHCs. We also propose a new delay predictor that incorporates information regarding the system state as well as the service time distribution. We compare the operational outcomes at both PHCs without diversion and with diversion using the above delay predictors. We show numerically that more accurate delay predictors lead to more equitable distribution of resources involved in provision of childbirth care across both PHCs.

## 1 Introduction

The number of patients accessing healthcare services has increased significantly across the world, including in India [1,2], and subsequently demand at all tiers of healthcare facilities is likely to grow in the coming years. Long wait times burden the healthcare administration in addition to inconveniencing and worsening patient outcomes, and to alleviate this, referral mechanisms of various types are implemented. Inter-facility referral systems are considered to be effective mechanisms for

reducing delays in admission to hospitals [3], and strengthening the efficiency of referral system has potential to improve quality of care in the community [4]. Two types of referral policies [5] are typically practiced in the healthcare context: (a) vertical referral, when the required equipment and/or expertise are unavailable at the current healthcare facility, and therefore patients are referred to a higher level of care, and (b) horizontal referral, when patients cannot access treatment within some threshold time duration due to limited healthcare capacity and are therefore referred to other facilities typically at a similar level of care. In this work, we consider horizontal referral, and we use the term patient diversion in place of horizontal referral to be consistent with the health operations literature. We study patient diversion in the context of primary health center (PHC) operations in the Indian context. This is because previous studies [6] have shown that a significant proportion of childbirth patients are likely to not receive care at public primary healthcare facilities in the Indian context within a reasonable timeframe (e.g., two hours). Hence, we investigate whether diverting these patients to other primary healthcare facilities can reduce their estimated wait times.

While patient diversion has been proposed in many studies as a method to reduce wait times, particularly in the context of ambulance diversion to reduce emergency department delays [7], we propose real-time delay prediction as a basis for making the diversion decision. This is because estimating the delay for a given patient arriving at the facility seeking care as a function of the state of the system or delay histories on a real-time basis (i.e., at the time the patient arrives at the facility) can provide the most up to date information that can inform the diversion decision, as opposed to using steady state measures of average wait time. Note that the delay prediction must be made at all facilities to which the patient is being considered for diversion, and not only at the facility they first arrive at. In this work, we estimate real-time delays using multiple predictors and show that the extent to which diversion mechanisms affect operational outcomes at the facilities in the network depends upon the accuracy of the delay predictor used.

We now briefly discuss the relevant literature. We first discuss diversion studies in the health operations literature. [8] summarized 137 articles addressing the ambulance offload delay problem and described how diversion reduced congestion and average wait times at a healthcare facility, and smoothens patient flow without increasing capacity [9]. Centralized diversion policies are found to be preferable to decentralized policies, because decisions taken by one healthcare facility affect the operational outcomes of the other healthcare facilities involved in diversion [10]. Therefore, sharing information regarding among healthcare facilities regarding their operational state becomes an important consideration during diversion. This is supported in the literature by studies that showed how diversion resulted in worse health outcomes due to lack of coordination within the diversion network, as diverted patients had to wait longer at the facility they were diverted to than at their facility of origin [11]. Previous studies [12, 13] have also proposed strategies to reduce or eliminate diversions such as increasing resource capacities at healthcare facilities. In [14], the authors presented a Markov decision process formulation for diverting ambulances in a two-facility problem,

with the objective of minimizing the wait time of patients beyond a clinically important threshold duration. They assumed that the distribution of the time to start treatment at the other facility is known. In [10], the authors formulated the AD decision as an mixed integer linear program in terms of minimizing patient wait times across the entire diversion network. The formulation is solved at discrete time intervals on a rolling horizon basis and concluded that a formulation implementing a centralized policy outperformed other models in minimizing patient tardiness.

It is thus evident that predicting delays on a real-time basis at all facilities in the diversion network can help develop a centralized diversion policy. We briefly discuss the real-time delay prediction literature in this context. Multiple studies have developed real-time delay predictors for arriving entities in different types of service systems and we refer readers to a relatively recent review [15] for a comprehensive account of the relevant literature. Three types of delay predictors have been proposed for predicting real time delays at service systems based on: (a) queue length [16], (b) delay history [17] and (c) machine learning [18]. With regard to M/G/1 queuing systems, while many delay predictors have been developed for this system, the existing delay predictors do not consider the limits or extreme quantiles of the service distributions in making their predictions, which we do in the delay predictor we develop in this study.

Previous approaches have implemented diversion without providing personalized real-time delay estimates that patients might experience at each facility in the network. In this study, we predict delays on a real-time basis at both healthcare facilities we consider using system state (queue length, elapsed service time, etc.) and delay history-based delay predictors. We simulate the diversion mechanism using real-time delay predictions across both facilities. Our main contributions are: (i) to provide a framework for real-time delay prediction based diversion in a network of queueing systems, particularly in healthcare; (ii) to propose a simple and easy to implement delay predictor based on elapsed service time, queue length, and limits of the service time distribution, and (iii) to show how the accuracy of delay predictors is related to the extent to which operational outcomes become more equitable across the diversion network. We now briefly describe PHC operations, the health facility that we consider for diversion.

## 2 Primary Health Centers

In India, PHCs are the first point of contact with a formally trained doctor and cater to outpatients, and on a limited basis to inpatient, childbirth patients and those requiring antenatal care [19]. We briefly describe the flow of patients through a PHC here. A detailed description of PHC operations, their simulation model development, including parameterization and outputs is provided in [6].

A PHC typically contains one or two doctors (typically general physicians), a staff nurse serv-

ing inpatients and childbirth patients, another nurse assisting the doctors with outpatients, four to six beds for inpatients, and a labour room with one bed for childbirth patients. PHCs also house a clinical laboratory for conducting common laboratory tests and a pharmacy that also manages patient registration. The outpatient department (OPD) operates for eight hours a day whereas the inpatient and childbirth facilities operate on a 24X7 basis, with staff nurses working in shifts to manage these departments. Note that doctors are typically available only during outpatient hours but may attend inpatients and childbirth cases outside outpatient hours. Based on our visits to these facilities, this appears to occur very infrequently, and hence we assume in the model that doctors are available only during outpatient hours.

Outpatients whose age is less than 30 years directly consult doctors upon arrival and if the doctor is busy, they join the outpatient queue. Patients whose age is greater than 30 years consult a nurse first before consulting the doctor. This nurse measures the patient's vitals, including checking for hypertension and high blood glucose levels as part of a lifestyle and non-communicable disease prevention scheme. Once the patient has finished consulting with the doctor, a certain proportion of patients are sent to an in-house laboratory if tests are required. All patients exit via the pharmacy, where in addition to obtaining pharmaceuticals as required, patients also register their visit. Inpatients typically consist of those requiring admission and care for relatively simple conditions; more complex and/or life-threatening cases are referred to secondary or tertiary care facilities. Inpatient lengths of stay are limited by the facility to 24 hours. If inpatients arrive during OPD hours, they first consult with the doctor and they are then admitted to the inpatient ward where the staff nurse monitors their condition and provides care as required. Outside OPD hours, inpatients are attended to by the staff nurse directly. Upon arriving at the PHC, childbirth patients are also first attended to by doctors during OPD hours. The doctor initiates childbirth bed request for patients so that they can be sent to the labour room. Once the labour duration is finished, these patients are shifted to an inpatient bed for between 24-48 hours before they are discharged. Outside OPD hours, childbirth patients are attended to by the staff nurse while the rest of their patient flow remains the same.

Based on this patient flow and parameter estimates taken from [6], we simulate patient care operations at two PHCs. We program this discrete event simulation model of the operations of two PHCs in Python using the salabim package, on an Intel i7 64-bit Microsoft Windows OS with 16 GB memory. We present estimates of operational outcomes at both PHCs from the simulation in Table 1. We only list outcomes relevant to this study – a full list is provided in [6].

We observe that a significant proportion of childbirth patients wait for more than a given threshold duration (assumed to be 2 hours) before being admitted to the labour room at the PHC they visit. It is this observation that motivated us to consider diversion for childbirth patients in this study.

Table 1: PHC Operational Outcomes

Outcome Measures	PHC 1 (4/1440/2880/1/1/1/6)*	PHC 2 (4/720/2880/1/1/1/6)*
Doctor occupancy	0.627 (0.003)	0.658 (0.004)
Staff nurse occupancy	0.307 (0.005)	0.447 (0.011)
Inpatient bed occupancy	0.114 (0.004)	0.209 (0.007)
Childbirth bed occupancy	0.470 (0.019)	0.923 (0.031)
% of childbirth cases whose wait time exceeds 2 hours	38.98 (2.77)	88.19 (5.46)

\* Outpatient load/childbirth patient load /inpatient load /number of doctors/number of staff nurse per shift/labour bed/IPD beds

### 3 Real Time Delay Prediction based Patient Diversion

In this section we present the diversion algorithm for childbirth patients and describe the delay predictors used for estimating the real-time delays of patients at both PHCs. In Figure 1, we present the centralized diversion algorithm for childbirth patients. We also consider the travel time between PHCs during diversion and divert patients to the other PHC only when the predicted delay at the other PHC (estimated at the time the patient is expected to reach the other PHC) plus the patient’s travel time is lesser. Here we estimate expected delays of patients using system state variables such as queue length and elapsed service time.

#### 3.1 Delay Predictors

We now describe the delay predictors used in our diversion model. We emphasize here that the delay prediction is made at the point in time when the patient arrives in the system - hence the term ‘real-time delay prediction’. Note that we do not use the average wait time as a delay predictor given that it is well established that average wait times are routinely outperformed by system state or delay history-based predictors for the purpose of real-time delay prediction [15]. The queueing system that we generate the delay prediction for is the labour room bed, which we approximate as an M/G/1 queueing system. The interarrival time distribution of the childbirth patients to the PHC is exponential with means for each PHC given in Table 1, and the time spent by the childbirth patient in the labour room bed (the “service time”) is uniformly distributed with parameters 360 minutes and 600 minutes.

1. Remaining service time-based delay predictor (predicted delay denoted as  $w_{rst}$ ).

$$w_{rst} = L_q E[S] + Pr \{Server \text{ is busy}\} E[remaining \text{ service time} | server \text{ is busy}]$$

Here  $L_q$  represents the length of the queue at the time the prediction is generated, and  $E[S]$  is the

---

Note: Travel time between PHCs is  $TT_{xy}$ , where  $x$  denotes the facility from where the patient is diverted, and  $y$  denotes the facility the patient is diverted to. Wait time is estimated on the basis of actual delays and predicted delays.

1. Patient arrives at PHC:
    - a. Consults with doctor during OPD hours
      - I. Estimate delay to get labour bed:
        1. If delay is lesser than threshold value:
          - a. Patient is admitted to PHC  $x$
          - b. Patient is attended to by the doctor and staff nurse
          - c. Patient exits after delivery
        2. Else:
          - a. Estimate delay at PHC  $y$ 
            - i. If delay at PHC  $y + TT_{xy} <$  threshold value:
              - a. Patient is diverted to PHC  $y$
              - b. Repeat steps 1.a.I.1.b and c
            - ii. Else:
              - a. Patient is admitted to PHC where delay is lesser
              - b. Repeat steps 1.a.I.1.b and c
    - b. Consults with nurse during non OPD hours
      - I. Estimate delay to get labour bed:
        1. If delay is lesser than threshold value:
          - a. Patient is admitted to PHC  $x$
          - b. Patient is attended to by the staff nurse
          - c. Patient exits after delivery
        2. Else:
          - a. Estimate delay at PHC  $y$ 
            - i. If delay at PHC  $y + TT_{xy} <$  threshold value:
              - a. Patient is diverted to PHC  $y$
              - b. Repeat steps 1.b.I.1.b and c
            - ii. Else:
              - a. Patient is admitted to PHC where delay is lesser
              - b. Repeat steps 1.b.I.1.b and c
- 

Figure 1: Diversion algorithm for childbirth patients.

mean service time. Upon simplifying, we get

$$w_{rst} = \frac{Pr \{Server \text{ is busy}\} E[residual \text{ service time} | server \text{ is busy}]}{(1 - \rho)}$$

$\rho$  = fraction of time the server is busy (i.e., server utilization).

From [20],  $E[residual \text{ service time} | server \text{ is busy}] = \frac{E[S^2]}{2E[S]} = \frac{1 + C_B^2}{2} E[S]$ , where  $C_B^2 = \frac{var[S]}{E^2[S]}$ .

Therefore,  $w_{rst} = \left(\frac{1 + C_B^2}{2}\right) \left(\frac{\rho}{1 - \rho}\right) E[S]$ .

2. Elapsed service time and service time distribution based (proposed). The predicted delay is denoted by  $w_{est}$ .

$$w_{est} = L_q E[S] + \max((t_{avg} - t_e, \min(t_e - t_{avg}, t_{max} - t_e)))$$

Here  $t_e$  = Elapsed service time of patient on labour room bed;  $t_{min}$ ,  $t_{max}$  = length of stay distribution limits on labour room bed ( $t_{min} = 360$  minutes,  $t_{max} = 600$  minutes); and  $t_{avg} = \frac{t_{min} + t_{max}}{2}$ .

## 4 Results

In this section, we present simulation results using the delay predictors presented in Section 3.1 and compare operational outcomes for three cases: (a) no diversion, (b) actual delay-based diversion, and (c) delay prediction base diversion. We quantify accuracy of delay predictors using the mean absolute percentage error (MAPE), specified by:  $\frac{1}{N} \sum_{i=1}^N \left| \frac{A_i - P_i}{A_i} \right|$ , where  $A_i$  represents actual delay and  $P_i$  represents predicted delays estimated using different delays predictors and  $N$  is the number of patients in our sample.

As described in section 2, we observe that a significant proportion of childbirth patients experience substantial wait time before getting admitted in the healthcare facility they visit. With the no diversion case, we observe that approximately 75.92% of childbirth patients wait longer than two hours before getting admitted to the labour room and this proportion reduces significantly to 42.43% with actual delay based patient diversion. We also estimate the extent to which differences in operational outcomes between PHCs change when diversion is implemented. We note that prior to diversion, PHC 1 (see Table 1) had significantly lower resource utilization levels when compared to PHC 2. Diversion helps make the utilization levels across PHCs more equitable, and we show the extent to which this occurs for resource levels directly involved in provision of childbirth care depends upon the accuracy of delay predictors employed.

Table 2 shows the results for childbirth patients when diversion is implemented using actual delay values (obtained from the simulation) and using the delay predictors  $w_{rst}$  and  $w_{est}$ . The results are benchmarked against the case when no diversion is implemented. We report the percentage differences between resource utilizations involved directly in childbirth patient care, the average wait time for the labour bed and the proportion of patients whose wait time exceeds two hours ( $\alpha$ ). It is evident that patient diversion improves outcomes in general – both  $\alpha$  and the labour bed wait time decrease with diversion. We see that the differences in resource utilization decrease the most (or in other words, become more equitable across both PHCs) as the accuracy of the delay predictor increases, with the greatest change (most equitable) observed when actual delays are used (i.e., 100% accuracy), and then decreases as the accuracy of the delay predictor decreases. The MAPEs of the delay predictors  $w_{rst}$  and  $w_{est}$  are 21.43% and 11.34% respectively. When we conducted a sensitivity analysis by increasing the arrival rate of childbirth patients, we observed MAPEs of 20.88% and 12.61% for  $w_{rst}$  and  $w_{est}$  respectively, indicating that similar trends are

observed even as patient load increases.

Table 2: Percentage difference in operational outcomes when patient diversion is implemented

Diversion case	$\Delta\rho_{doc}$	$\Delta\rho_{nurse}$	$\Delta\rho_{IPD}$	$\Delta\rho_{lb}$	Labour bed wait time (minutes)	$\alpha$
No diversion	4.73	31.3	45.73	49.07	88.31	75.92(6.43)
$w_{act}$	1.99	9.21	13.9	15.48	56.79	42.43(2.72)
$w_{rst}$	2.52	22.29	25.53	24.57	60.47	63.52(3.69)
$w_{est}$	2.2	13.97	19.74	20.1	48.37	54.59(1.24)

$\Delta\rho_{doc}$  = difference in doctor's utilization;  $\Delta\rho_{nurse}$  = difference in the staff nurse's utilization;  $\Delta\rho_{IPD}$  = difference in inpatient bed utilization;  $\Delta\rho_{lb}$  = difference in labour bed utilization;  $\alpha$  = proportion of childbirth patients with wait time > 2 hours

## 5 Conclusion

In this work, we present a framework for implementing real-time delay prediction based patient diversion across two healthcare facilities in an Indian district. Our study shows that the extent to which operational outcomes become more equitable (that is, a facility with high congestion will see a decrease in congestion, and a facility with low utilization levels will see an increase in utilization) depend upon the accuracy of delay predictors.

In addition to the fact that to the best of our knowledge, our work is the first study that considers real-time delay prediction as a basis for diversion in a health facility network, the approach that we propose may be used for diversion in general queueing systems as well – for example, in call centers. While there is a substantial body of research in the area of real-time delay prediction for queueing systems [15], we have not come across a study that implements diversion based on delay predictions. Note that articles investigating the impact of delay prediction on reneging, balking and jockeying (for multi-server systems) are present in the literature [15]; however, these assume voluntary actions on the part of the entity waiting for service, and not a policy undertaken by the queue administration.

A key assumption in our work involves the presence of a centralized administration for the network of healthcare facilities that monitors and records system state information at all facilities in the network. Such an administrator would have to generate the delay predictions whenever a new patient arrives at any facility in the network, and then make the diversion decision based on the delay estimates. This would imply availability of the requisite information technology infrastructure to facilitate deployment of this diversion mechanism.

Future avenues of research involve extension of this work to include the entire network of PHCs in a district, and potentially also secondary and tertiary levels of care. We also did not consider compliance of patients with the diversion decision, which can be investigated in further studies to determine its effect on overall network operational outcomes.



## References

- [1] S. J. Poon, J. D. Schuur, and A. Mehrotra, “Trends in visits to acute care venues for treatment of low-acuity conditions in the united states from 2008 to 2015,” *JAMA internal medicine*, vol. 178, no. 10, pp. 1342–1349, 2018.
- [2] U. Shrivastava, A. Misra, V. Mohan, R. Unnikrishnan, and D. Bachani, “Obesity, diabetes and cardiovascular diseases in india: public health challenges,” *Current diabetes reviews*, vol. 13, no. 1, pp. 65–80, 2017.
- [3] P. Pouramin, C. S. Li, J. W. Busse, S. Sprague, P. Devereaux, J. Jagnoor, R. Ivers, M. Bhandari, G. Guyatt, B. Petrisor *et al.*, “Delays in hospital admissions in patients with fractures across 18 low-income and middle-income countries (inormus): a prospective observational study,” *The Lancet Global Health*, vol. 8, no. 5, pp. e711–e720, 2020.
- [4] C. Give, S. Ndima, R. Steege, H. Ormel, R. McCollum, S. Theobald, M. Taegtmeier, M. Kok, and M. Sidat, “Strengthening referral systems in community health programs: a qualitative study in two rural districts of maputo province, mozambique,” *BMC health services research*, vol. 19, no. 1, pp. 1–11, 2019.
- [5] P. W. Handayani, I. R. Saladdin, A. A. Pinem, F. Azzahro, A. N. Hidayanto, and D. Ayuningtyas, “Health referral system user acceptance model in indonesia,” *Heliyon*, vol. 4, no. 12, p. e01048, 2018.
- [6] M. Shoaib and V. Ramamohan, “Simulation modelling and analysis of primary health centre operations,” 2020, [http://web.iitd.ac.in/~varunr/PHC\\_Paper\\_Preprint.pdf](http://web.iitd.ac.in/~varunr/PHC_Paper_Preprint.pdf), accessed 23<sup>rd</sup> January 2021.
- [7] M. H. Yarmohammadian, F. Rezaei, A. Haghshenas, and N. Tavakoli, “Overcrowding in emergency departments: a review of strategies to decrease future challenges,” *Journal of research in medical sciences: the official journal of Isfahan University of Medical Sciences*, vol. 22, 2017.
- [8] M. Li, P. Vanberkel, and A. J. Carter, “A review on ambulance offload delay literature,” *Health care management science*, vol. 22, no. 4, pp. 658–675, 2019.
- [9] N. Nezamoddini and M. T. Khasawneh, “Modeling and optimization of resources in multi-emergency department settings with patient transfer,” *Operations Research for Health Care*, vol. 10, pp. 23–34, 2016.
- [10] S. Baek, Y. H. Lee, and S. H. Park, “Centralized ambulance diversion policy using rolling-horizon optimization framework to minimize patient tardiness,” in *Healthcare*, vol. 8, no. 3. Multidisciplinary Digital Publishing Institute, 2020, p. 266.
- [11] S. Deo and I. Gurvich, “Centralized vs. decentralized ambulance diversion: A network perspective,” *Management Science*, vol. 57, no. 7, pp. 1300–1319, 2011.

- [12] E. M. Castillo, G. M. Vilke, M. Williams, P. Turner, J. Boyle, and T. C. Chan, “Collaborative to decrease ambulance diversion: the california emergency department diversion project,” *The Journal of emergency medicine*, vol. 40, no. 3, pp. 300–307, 2011.
- [13] O. K. Asamoah, S. J. Weiss, A. A. Ernst, M. Richards, and D. P. Sklar, “A novel diversion protocol dramatically reduces diversion hours,” *The American journal of emergency medicine*, vol. 26, no. 6, pp. 670–675, 2008.
- [14] A. Ramirez-Nafarrate, A. B. Hafizoglu, E. S. Gel, and J. W. Fowler, “Optimal control policies for ambulance diversion,” *European Journal of Operational Research*, vol. 236, no. 1, pp. 298–312, 2014.
- [15] R. Ibrahim, “Sharing delay information in service systems: a literature survey,” *Queueing Systems*, vol. 89, no. 1, pp. 49–79, 2018.
- [16] W. Whitt, “Improving service by informing customers about anticipated delays,” *Management science*, vol. 45, no. 2, pp. 192–207, 1999.
- [17] R. Ibrahim and W. Whitt, “Real-time delay estimation based on delay history,” *Manufacturing & Service Operations Management*, vol. 11, no. 3, pp. 397–415, 2009.
- [18] V. Baldwa, S. Sehgal, V. Tandon, and V. Ramamohan, “A combined simulation and machine learning approach for real-time delay prediction for wait-listed neurosurgery candidates,” in *Proceedings of the Winter Simulation Conference, December, 2020*, pp. 13–16.
- [19] IPHS-Guidelines, *Guidelines for Primary Health Centres*, Directorate General of Health Services, New Delhi, India, 2012.
- [20] D. Gross, *Fundamentals of queueing theory*. John Wiley & Sons, 2008.