



**ISE**

Industrial and  
Systems Engineering

# Stochastic Sequencing of Surgeries for a Single Surgeon Operating in Parallel Operating Rooms

Robert Storer  
Lehigh University

Camilo Mancilla  
Lehigh University

Report: 11T-001

# Stochastic Sequencing of Surgeries for a Single Surgeon Operating in Parallel Operating Rooms

Camilo Mancilla and Robert H. Storer

Lehigh University, Department of Industrial and Systems Engineering, January 5, 2011

## Abstract

We develop algorithms for a stochastic two-machine single-server sequencing problem with waiting time, idle time and overtime costs. Scheduling surgeries for a single surgeon operating in two parallel operating rooms (ORs) motivates the work. The basic idea is that OR staff perform cleanup and setup in one room while the surgeon is operating in the other. The benefit is less waiting time for the surgeon between surgeries, but may also result in added idle time for OR staff if cleanup and setup are completed prior to completion of surgery in the other OR. When surgeries are long relative to cleanup and setup times, parallel OR scheduling is not attractive as significant OR staff idle time will result. The problem we address consists of assigning surgeries to ORs and sequencing them, and is formulated as an integer stochastic program using sample average approximation. A decomposition based solution approach is developed. Computational testing based on real data shows that the proposed method solves problems to optimality in acceptable processing time. Using our solution methodology, we further provide insight as to the conditions when scheduling a single surgeon in parallel operating rooms is cost effective.

## 1 Introduction

Operating Rooms (ORs) account for more than one third of the total revenue in a hospital according to the Health Care Financial Management Association HFMA (2005). Further, many if not most hospital activities are driven by the surgery schedule, thus improvement in utilization and efficiency can have a significant cost impact. In this paper a method is developed to help the OR Manager schedule a set of surgeries for a single surgeon operating in two parallel ORs. Specifically the problem is to assign each surgical procedure to one of the two ORs and determine the sequence of surgeries in each room. It is assumed that the surgical procedures can be divided into three stages; setup, surgery, and cleanup, where surgery refers to that part of the procedure requiring the presence of the surgeon. Since the setup and cleanup stages do not require the surgeon, he or she can move to the other OR to perform surgery. If all surgeries were performed in a single OR, the surgeon would experience waiting time during the setup and cleanup stages of each procedure (excluding the setup before the first surgery and the cleanup after the last surgery). With a single OR, there will be no OR staff idle time since there is always an activity to perform. If parallel ORs are used there is an opportunity to significantly reduce the surgeon's waiting time. However, the possibility for OR staff idle time now

exists when surgery in one room extends past the time setup is complete in the other room. When surgery times are long relative to setup and cleanup, using parallel ORs is likely to be unrealistic due to the large amounts of OR staff idle time that would necessarily result. If on the other hand surgery time is comparable to setup and cleanup, the possibility exists to greatly reduce surgeon waiting time without incurring significant OR staff idle time. In this paper the setup, surgery, and cleanup times are modeled as random variables with known joint distribution. The marginal distributions of setup, surgery and cleanup durations are not assumed identical for each surgical procedure but rather differ based on the surgeon, type of surgical procedure, patient attributes, etc. The objective of our optimization problem is to minimize a weighted linear combination of expected costs that include surgeon waiting time, OR idle time, and staff overtime. The method must produce solutions fast enough to make operational decisions on a daily basis.

We formulate a two-stage stochastic integer programming model then propose a decomposition method based on the first stage variables and the mean value problem to solve it. We perform computational testing to investigate algorithm behavior, and also provide managerial insight regarding when single-surgeon, parallel OR scheduling is cost effective.

The remainder of the article is organized as follows. In the next section a brief review of the literature related to the two-machine single-server scheduling problem and OR scheduling is provided. In Section 3, the model is described and formulated. In Section 4 algorithms are proposed to solve the sequencing and assignment problems. In Section 5 computational results are presented. In Section 6 we present computational results with respect to the number of scenarios needed to solve the problem. In Section 7 we provided insight as to when parallel operating rooms is cost effective. In Section 8 conclusions and future research directions are discussed.

## 2 Literature review

Two branches of previous research are relevant to this paper, work on surgical scheduling, and work on two-machine, single-server sequencing. Batun et al. (2010) consider a model where surgical procedures are also composed of random setup, surgery and cleanup times with known distributions. They consider a problem with  $R$  operating rooms and  $N$  surgeons. Each surgeon has a known set of surgeries and a pre-specified sequence in which they will be performed. The authors find the optimal assignment of these surgeries to operating rooms and the sequence of surgeries in each room. They show computational results for problems with 4 to 11 total surgeries, 2 to 3 surgeons and 3 to 6 ORs. They used the L-Shaped method and also derived cuts based on the mean value problem to speed up the algorithm. In our work, we do not assume that the sequence of surgeries is known in advance, and have only a single surgeon and two ORs.

In the machine scheduling literature a deterministic version of the problem we address has been studied by Koullamas (1996), Abdekhodae et al. (2006), Abdekhodae and Wirth (2002). This problem involves scheduling a set of jobs on 2 machines where the setup is performed by a single server (e.g. a robot). Once the setup for a job is complete, the job is processed on one of the two machines. Problems of this type are found in manufacturing and network

computing according to Glass et al. (2000). All of these studies assume deterministic setup and processing times. The correspondence between the machine scheduling problem and our parallel OR scheduling problem can be seen if we set our setup times to zero, let our surgery times equate to their machine setup times (which require a single resource, the surgeon or robot), and their machine processing times equate to our cleanup time (which can be done in parallel). Glass et al. (2000) created an approximation algorithm with a performance guarantee and also provided a description of polynomially solvable cases. Abdekhodae and Wirth (2002) considered problems in which all setup times were greater than all processing times. Abdekhodae et al. (2006) extended the work of Abdekhodae and Wirth (2002) to the case in which the setup and processing times do not have any restrictions. They developed different heuristics that they then compared against meta-heuristic approaches. None of these studies considered stochastic processing times.

### **3 Mathematical formulation of the sample average parallel operating room scheduling (SAPOS) problem**

In this section we formulate a stochastic integer programming model for the problem described above. In order to account for the randomness of the setup, surgery, and cleanup times we use sample average approximation. The model we develop has pure binary first stage decision variables  $x_{ij}$  representing the sequence of surgeries, and  $m_{ir}$  representing the assignment of surgery  $i$  to one of the two ORs. The model also includes the variable  $q_{ij}$  representing the immediate predecessor surgeries in the overall surgery sequence). Note that the  $q$  variables can be calculated directly from  $\vec{x}$ . The second stage variables include, for each scenario, the start and finish time of each stage of each procedure, the surgeon waiting time, OR idle time, and overtime in each OR. Since the stage two variables are easily computed given values for the stage one binary variables  $\vec{x}$  and  $\vec{m}$ , we have “simple recourse”. Finally we can compute the average (over scenarios) cost. The problem formulation appears below.

$$\text{minimize: } \sum_{k=1}^K \frac{1}{K} \left( \sum_{i \neq j} c_w W_{ij}^k + \sum_{r \leq R} c_i I_r^k + \sum_{r \leq R} c_o O_r^k \right) \quad (1)$$

s.t.:

$$x_{ij} + x_{ji} = 1 \quad \forall i, j \quad (2)$$

$$x_{ij} + x_{je} + x_{ei} \leq 2 \quad \forall 1 \leq i \neq j \neq e \leq n \quad (3)$$

$$a_{ijr} \leq x_{ij} \quad \forall i, j, r \quad (4)$$

$$\sum_{r \leq R} m_{ir} = 1 \quad \forall i \quad (5)$$

$$a_{ijr} + a_{jir} \leq m_{ir} \quad \forall i, j > i, r \quad (6)$$

$$a_{ijr} + a_{jir} \leq m_{jr} \quad \forall i, j > i, r \quad (7)$$

$$a_{ijr} + a_{jir} \leq m_{ir} + m_{jr} - 1 \quad \forall i, j > i, r \quad (8)$$

$$\sum_{t \neq i} x_{it} + 1 = \sum_{u=1}^n u q_{iu}^1 \quad \forall i \quad (9)$$

$$\sum_{u=1}^n q_{iu}^1 = 1 \quad \forall i \quad (10)$$

$$q_{ij} = \sum_{u=1}^{n-1} q_{iju}^2 \quad \forall i, j \neq i \quad (11)$$

$$q_{iju}^2 \leq q_{iu+1}^1 \quad \forall i, j \neq i, u < n \quad (12)$$

$$q_{iju}^2 \leq q_{ju}^1 \quad \forall i, j \neq i, u < n \quad (13)$$

$$q_{iju}^2 \geq q_{ju}^1 + q_{iu+1}^1 - 1 \quad \forall i, j \neq i, u < n \quad (14)$$

$$C_i^k \leq M m_{ir} \quad \forall i, r \quad (15)$$

$$W_{ij}^k \geq \sum_{r \leq R} C_{jr}^k - \sum_{r \leq R} C_{ir}^k + post_i^k - sur_j^k - post_j^k - M(1 - q_{ij}) \quad \forall i, j \neq i, k \quad (16)$$

$$W_{ij}^k \leq \sum_{r \leq R} C_{jr}^k - \sum_{r \leq R} C_{ir}^k + post_i^k - sur_j^k - post_j^k + M(1 - q_{ij}) \quad \forall i, j \neq i, k \quad (17)$$

$$\sum_{r \leq R} C_{jr}^k \geq \sum_{r \leq R} C_{ir}^k + pre_j^k + sur_j^k + post_j^k - M(1 - a_{ijr}) \quad \forall i, j \neq i, k \quad (18)$$

$$C_{ir}^k \geq pre_i^k + sur_i^k + post_i^k - M(1 - m_{ir}) \quad \forall i, r, k \quad (19)$$

$$C_{ir}^k \leq l_r^k \quad \forall i, r, k \quad (20)$$

$$I_r^k = l_r^k - \sum_{j=1}^n m_{ir} (sur_j^k + pre_j^k + post_j^k) \quad \forall r, k \quad (21)$$

$$l_r^k - d_r = O_r^k - g_r^k \quad \forall r, k \quad (22)$$

$$C_{ir}^k, I_r^k, l_r^k, O_r^k, g_r^k, W_{ij}^k \geq 0 \quad \forall i, j, k, r$$

$$x_{ij}, m_{ir}, a_{ijr}, q_{iu}^1, q_{ij}, uq_{iju} \in \{0, 1\} \quad \forall i, j, u, r$$

### Indices and Sets

$n$  is the number of procedures to be scheduled

$K$  is the number of scenarios

$R$  is the number of ORs in this case 2.

$j, i, t$  and  $e$  index the set of procedures: from  $1, \dots, n$ .

$u$  indexes the position in the sequence of procedures:  $u=1, \dots, n$ .

$r$  indexes is the set operating rooms:  $r=1,2$

$k$  indexes scenarios:  $k=1,\dots,K$ .

#### Parameters

$c_w$  Surgeon waiting time cost per unit time..

$c_o$  overtime cost per unit time.

$c_i$  idle time cost per unit time.

$d_r$  time beyond which overtime is incurred in O.R.  $r$

$M$  is sufficiently large "big M" number.

$pre_i^k$  duration of setup for surgery  $j$  in scenario  $k$ .

$sur_i^k$  duration of surgery  $i$  in scenario  $k$ .

$post_i^k$  duration of cleanup for surgery  $i$  in scenario  $k$ .

#### Variables

$W_{ij}^k$  Surgeon waiting time between procedure  $i$  and  $j$  in scenario  $k$ .

$C_{ir}^k$  Completion time of procedure  $i$  in OR  $r$  in scenario  $k$ .

$l_r^k$  Completion time of the last procedure in OR  $r$  in scenario  $k$ .

$I_r^k$  Total OR idle time in scenario  $k$  for OR  $r$ .

$O_r^k$  Overtime time in OR  $r$  in scenario  $k$ .

$g_r^k$  a slack variable that measures the earliness with respect to time  $d_r$  in scenario  $k$  in O.R.  $r$ .

$x_{ij}$  a binary variable denoting the assignment of surgery  $i$  as a predecessor of surgery  $j$  in the overall surgery sequence.

$m_{ir}$  a binary variable denoting the assignment of surgery  $i$  to OR  $r$ .

$a_{ijr}$  a binary variable denoting the assignment of surgery  $i$  as a predecessor of surgery  $j$  in OR  $r$ .

$q_{ij}$  a binary variable denoting the assignment of surgery  $i$  as the immediate predecessor of surgery  $j$  in the overall surgery sequence.

$q_{ij}^1, q_{ij}^2$  utility binary variables that help to compute  $q_{ij}$ .

#### Constraints

(2), (3), (4) define the possible sequences for the surgeon.

(6), (7), (8) define the sequence inside each OR.

(9), (10), (11), (12), (13), (14) compute the immediate predecessor of surgeries.

(15) set to zero the value of completion time if surgery not assigned to this room.

(16), (17) define the surgeon waiting time.

(18), (19) define the completion time for each surgery for a given assignment.

(20) calculates the completion time for every OR.

(21) calculates the OR idle time.

(22) calculates the OR overtime.

It can be seen that computing the idle time, overtime time and surgeon waiting time are simple calculations once we have the sequence of surgeries ( $x$  variables) and the respective sequence in each OR ( $m$  variables).

### 3.1 Problem complexity

The SAPOS problem can be shown to be NP-Hard through a simple extension of results from Koulamas (1996). Koulamas (1996) proved that the deterministic version of the following problem is NP-Hard. This definition is extracted as it appeared in Koulamas' paper.

“Consider a deterministic  $n$  job, two-machine one-server scheduling problem with the following assumptions.

- There are two parallel identical semiautomatic machines. The machines run unattended during job processing; however the presence of a server (robot) is required during setups.
- There is a set of jobs awaiting processing at time zero. Each job  $i$  has a processing time  $P_i$  and a sequence independent setup time  $s_i$ . All  $s_i$  and  $P_i$  are deterministic and known in advance.
- There is one server (robot) serving both machines. Overlapping requests for the server result in machine idle time (interference).
- Each machine processes one job at a time and job pre-emption is not allowed.”

The problem described in Koulamas (1996) is a special case of our problem where our setup times are zero, their setup times equate to our surgery times, and their processing times equate to our cleanup times. Therefore, the parallel OR scheduling problem is NP-Hard even in the single scenario (deterministic) case with zero setup time.

## 4 Proposed solution approach

Our first attempt at solving the problem was to simply use Cplex to solve the extended formulation. Unfortunately, Cplex was not able to achieve a reasonable relative gap (100%) in two hours with only 1 scenario and ten surgeries. Our second attempt was to use the Integer L-Shaped method (Laporte and Louveaux (1993)), with the binary ( $x$  and  $m$ ) first stage decision variables. The Integer L-Shaped method failed to close the relative gap with results similar to Cplex. During experimentation we discovered that if we fixed the surgeon-OR sequence, the resulting MIP sub-problem (which we call the “surgery-sequencing sub-problem”) was solvable in manageable processing time. To clarify, we define vector  $y$  as the “surgeon-OR sequence”. For instance, if  $y$  is (0,1,0,0,1,1) then the surgeon will perform the first surgery in room 1, then room 2, then 2 surgeries in a row in room 1, then 2 in a row in room 2, etc.

The search space of the vector  $y$  is large (there are  $2^n$  possible surgeon-OR sequences) but we can take advantage of other characteristics of the problem to make this search easier. The following is the extended problem formulation assuming we have assigned the surgeries in a “zigzag” fashion ( $y = 0, 1, 0, 1, 0, 1, \dots$ ).

$$\min \sum_{k=1}^K \frac{1}{K} \left( \sum_{i=1}^{n-1} c_w W_i^k + \sum_{i=1}^{n-2} c_i I_i^k + c_i f s^k + \sum_{r \leq R} c_o O_r^k \right) \quad (23)$$

$$\text{s.t. } : C_1^k = \sum_{j=1}^n (pre_j^k + post_j^k + sur_j^k) v_{1j} \quad k \quad (24)$$

$$W_i^k = C_{i+1}^k - C_i^k + \sum_{j=1}^n post_j^k v_{ij} - \sum_{j=1}^n (sur_j^k + post_j^k) v_{i+1,j} \quad i < n, k \quad (25)$$

$$I_i^k = C_{i+2}^k - C_i^k - \sum_{j=1}^n (pre_j^k + sur_j^k + post_j^k) v_{i+2,j} \quad i < n-1, k \quad (26)$$

$$f s^k = C_2^k - \sum_{j=1}^n (pre_j^k + sur_j^k + post_j^k) v_{2j} \quad k \quad (27)$$

$$C_n^k - d_1 = O_1^k - g_1^k \quad k \quad (28)$$

$$C_{n-1}^k - d_2 = O_2^k - g_2^k \quad k \quad (29)$$

$$\sum_{j=1}^n v_{ij} = 1 \quad i \quad (30)$$

$$\sum_{i=1}^n v_{ij} = 1 \quad j \quad (31)$$

$$C_i^k, I_i^k, f s^k, O^k, g_1^k, g_2^k, W_i^k \geq 0 \quad \forall i, k$$

$$v_{ij} \in \{0, 1\} \quad \forall i, j$$

#### Indices and Sets

$n$  procedures to be scheduled indexed by  $j$

$K$  scenarios to be considered indexed by  $k=1, \dots, K$ .

$i$  indexes positions in the surgery sequence  $i=1, \dots, n$ .

#### Variables

The variable names are the same as in the extended formulation previously described. with the following exceptions:

$C_i^k$  Completion time of procedure in position  $i$  in scenario  $k$ .

$I_i^k$  Idle time between procedure in position  $i$  and position  $i+2$  in scenario  $k$ .

$W_i^k$  Surgeon waiting time for procedure in position  $i+1$  in scenario  $k$ .

$v_{ij}$  is a binary variable denoting the assignment of surgery  $j$  to position  $i$  in the sequence of surgeries.

$f s^k$  is the idle time caused by the second surgery in scenario  $k$ .

#### Constraints

(24) Define the completion time for the first surgery.

(25),(26),(27) calculation of surgeon waiting time, OR idle time for the second surgery and OR idle time between surgeries for every scenario.

(28),(29) calculation of overtime for each OR and every scenario.

(30),(31) assignment constraints.



The surgery-sequencing sub-problem has symmetry if we consider equal overtime cost penalties in each OR, but this symmetry is easily broken by always assigning the first surgery to OR 1. We will search in the surgeon-OR sequence space ( $y$ ) that has a cardinality of  $2^{n-1}$ . Each possible surgeon-OR sequence ( $y$  vector) results in a surgery sequencing sub-problem in which the sequence of procedures must be determined.

The idea behind this algorithm is that often the zigzag OR sequence will be optimal. Further, the vast majority of the  $2^{n-1}$  OR sequences will be terrible. Further the LP relaxation of the mean value problem provides an easily computable lower bound that can be used to eliminate all but a few OR sequences.

#### 4.1 Basic parallel OR scheduling algorithm

The following decomposition algorithm uses bounds provided by the mean value problem.

Set the upper bound (UB) to  $\infty$ .

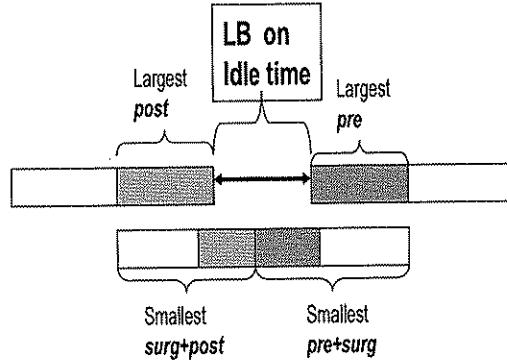
1. Initialize  $S$  as the finite set of all surgeon-OR sequences (with cardinality  $2^{n-1}$ ).  
Initialize  $s_i$  as the zigzag surgeon-OR sequence and remove from  $S$ .  
Solve the corresponding SAA MIP and initialize the UB.
2. Remove a new surgeon-OR sequence ( $s_i$ ) from  $S$  if  $S$  is not empty, otherwise end.
3. Solve the LP relaxation of the mean value problem for ( $s_i$ ). If the solution is less than the UB go to 4. Otherwise go to 2.
4. Solve the surgery sequencing sub-problem SAA for  $s_i$ . If the solution is less than the current UB update it, otherwise go to 2.

The idea behind this algorithm is that often the zigzag surgery sequencing sub-problem will be optimal, or near optimal, while the vast majority of the  $2^{n-1}$  surgeon-OR sequences will be terrible. Since the LP relaxation of the mean value problem provides an easily computable lower bound, it can be used to quickly eliminate all but a few surgeon-OR sequences. In the following sections we discuss enhancements to speed up the basic algorithm.

#### 4.2 Alternative lower bounds

We next define an alternative lower bound that can eliminate entire subsets of surgeon-OR sequence sub-problems in the algorithm. First, define a partition of the set  $S$  as follows: let  $P_i$  be the set of all possible surgeon-OR sequences in which there is a maximum of  $i$  surgeries consecutive in the same OR. For example  $(0,0,0,1,1,0,1,0,1,0,1,1,1)$  belongs to  $P_3$  since the maximum consecutive procedures in one OR is 3. This clearly partitions  $S$  since the union of  $P_i$ 's give us  $S$  and the intersection of any  $P_i$ s is empty. Next we define a lower bound based on these partitions. First, we define a lower bound for each cost component.

**Figure 1: Illustration of Lower Bound on idle time for  $P_2$**



**Idle time bound:** We use order statistic notation to simplify the reading of the bounds. For instance,  $\overline{post}_{[j]}$  means the  $j$ th smallest post value in the mean-value problem.  $L_i^I$  is the idle time lower bound for partition  $i$ . The bound is illustrated in figure 1. (In  $L_i^I$  the function  $ind_{>2}$  takes the value 1 if  $i > 2$  and takes the value 0 otherwise)

$$L_i^I = (\overline{surg} + \overline{post})_{[1]} + (\overline{surg} + \overline{pre})_{[1]} - \overline{post}_{[n]} - \overline{pre}_{[n]} + ind_{>2} \sum_{j=1}^{i-2} (\overline{pre} + \overline{surg} + \overline{post})_{[j]}$$

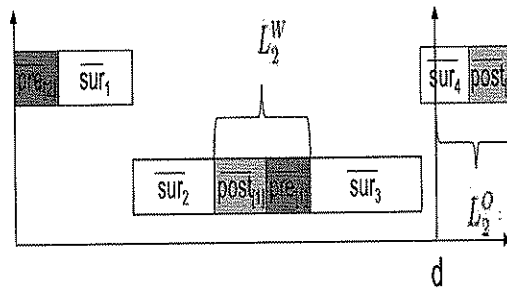
**Waiting time bound:**  $L_i^W$  is the wait time lower bound for partition  $i$ .

$$L_i^W = \sum_{j=1}^{i-1} \overline{pre}_{[j]} + \overline{post}_{[j]}$$

**Overtime time bound:**  $L_i^O$  is the overtime lower bound for partition  $i$ .

$$L_i^O = \max \left\{ \sum_{i \in J} \overline{surg}_j + \sum_{j=1}^{i-1} \overline{pre}_{[j]} + \overline{post}_{[j]} - d, 0 \right\}$$

**Figure 2: Illustration of Lower Bound on wait and over time for  $P_2$**



$$L_i = c_i L_i^I + c_w L_i^W + c_o L_i^O$$

The idea is to find a lower bound for the mean value problem by finding the best case scenario in terms of idle time, waiting time and overtime. The mean value problem is a lower bound for the expected value problem, thus the new bound is as well. Thus we can use this easily computable bound to eliminate whole sets of OR-sequences. Indeed,

when the lower bound  $L_i$  exceeds the current upper bound we can eliminate all surgeon-OR sequences belonging to  $P_i, P_{i+1},$  up to  $P_n$ .

### 4.3 A heuristic for computing an upper bound

In order to speed up the convergence of the algorithm we created a simple heuristic that seeks to find an improved initial upper bound (that is better than the solution to the zigzag problem) by spending additional computation time. This heuristic is based on Problem Search Space (Storer et al. (1992)). We perturb (using bootstrap sampling) the mean processing times in the mean value problem, solve the perturbed mean value problem to generate a new solution (sequence of surgeries and OR assignment), then evaluate the solution using the full sample average objective function. We generate 100 perturbed mean value problems and solutions, then keep the best upper bound found. Details of the perturbation algorithm are as follows:

*PSS perturbed mean value heuristic*

1. Generate  $K$  random numbers from  $\text{IntUnif}(1,K)$  and store the result in index  $i(h)$  ( $i(h) = \text{IntUnif}(1,K)$  for  $h=1, \dots, K$ ).
2. Compute the bootstrap sample average setup, cleanup and surgery times for each surgery from the bootstrap sample using the following formula

$$\mu_t = \frac{1}{K} \sum_{h=1}^K t_{i(h)} \quad t = \{\text{setup, cleanup, surgery}\}$$

3. Construct the “perturbed mean value” problem using the bootstrap averages.
4. Solve the perturbed mean value problem and then evaluate the solution using the objective function with the original set of scenarios.

We apply this procedure 100 times and then we pick the best value found in step 4. In our computational experiments (section 5) we found that the gap between the optimal solution and the mean value solution (evaluated using the zigzag OR-sequence) was 70% on average. The gap between the optimal solution and the heuristic solution was 36%. Thus the heuristic ultimately saved significant computation time.

### 4.4 Solving the surgery sequence sub-problems

Two straightforward approaches for solving the surgery sequence sub-problems are branch and bound and Benders’ decomposition. We conducted a brief experiment to find which factors dictate the choice of algorithms in terms of CPU time. For the Benders’ decomposition approach we implemented the followings strategies:

#### 4.4.1 Benders' cuts

The method used to compute the Benders' cuts may have an impact on the processing time of the algorithm, so we developed ways to speed up their computation using (1) the complementary slackness theorem and (2) the structure of the basis matrix. The basis matrix is lower-triangular in the primal therefore it is upper-triangular in the dual. We can thus use simple recursive formulas to compute the dual variables.

#### 4.4.2 Computation of the recourse function

Once the first stage variables are fixed in the extended formulation, the computation of the completion time for every surgery is easily found using the following formula executed from  $i=1$  to  $n$ .

$$C_i^k = \begin{cases} C_{i-1}^k + pre_{ps(i)}^k + sur_{ps(i)}^k + post_{ps(i)}^k & \text{if } r(i-1)=r(i) \\ \max\{C_{ps(i-1)}^k - post_{ps(i-1)}^k, C_{rp(i)}^k + pre_{ps(i)}^k\} + sur_{ps(i)}^k + post_{ps(i)}^k & \text{if } r(i-1) \neq r(i) \end{cases}$$

Where  $ps(i)$  is the surgery in position  $i$ ,  $rp(i)$  is the immediate position predecessor in the same OR of surgery in position  $i$ , and  $r(i)$  is the room assigned to position  $i$  in the surgery sequence. Given the completion times for every procedure we can compute the surgeon waiting time using the following formula

$$W_i^k = C_{i+1}^k - C_i^k + post_{ps(i)}^k - sur_{ps(i+1)}^k - post_{ps(i+1)}^k$$

#### 4.4.3 Recycling Benders' cuts

Since we have to solve many surgery sequencing sub-problems, and since they all have the same integer feasible region, we can store the integer feasible solutions from the first sub-problem and then recompute Benders' cuts in succeeding sub-problems. We hope to save computing time in the Benders' decomposition scheme although there is no guarantee that the cuts are not redundant in all the new sub-problems.

#### 4.4.4 Generalized upper bound cuts

Since we have assignment constraints in the master problem we have set Cplex "generalized upper bound cuts" parameter to "intensive" in order to speed up the solution of the master problem.

#### 4.4.5 Cut strengthening

We observed significant problems with degeneracy in certain instances. We therefore used methodology from Magranti and Wong (1981) designed for such cases where they strengthen cuts using the concept of dominance. Assume that  $x^*$  is the optimal solution of an integer program. We will say that  $(\pi, \pi_0)$  dominates  $(\pi^1, \pi_0^1)$  if

$$\pi x^* + \pi_0 \geq \pi^1 x^* + \pi_0^1$$

The idea behind these cuts is to find the dual variables that most improve the lower bound. Since, we have multiple solutions in the dual problem; we solve another linear problem that approximates the search for the dual variables that most improve the lower bound in the master problem. Since this strategy is costly, we only apply it every fifth iteration.

#### 4.4.6 Generalized upper bound (GUB) branching

The idea behind this branching rule is to take advantage of the assignment constraints. Instead of branching by variables we defined weights for variables that we will use to branch on with the fractional solution that results from the LP relaxation. This will lead to improved solution times according to Nemhauser and Wolsey (1999). We implemented GUB using the SOS1 feature in Cplex 12.1. We use this feature for the master problem in Benders' and also in straightforward branch and bound.

## 5 Computational experience

The proposed algorithm based on decomposing the problem by surgeon-OR sequences will solve some problems very quickly, while others take longer. In general, when the zigzag sequence (or a sequence close to zigzag) is optimal, the algorithm solves quickly (we also note that we expect problems to behave this way in reality when parallel ORs are an attractive option). We conducted a computational experiment that was designed to discover which factors affect the run time of the algorithm. The factors investigated are shown in table 1. The experiment had 5 replicates (and thus a total of 60 instances in the experiment). The stopping criteria for the algorithm was a 1% optimality gap. In order to clarify the effect of the distribution of setup and cleanup time with respect to surgery time we created the "setup to surgery time ratio" B. This factor has two levels; "high" means that the summation of the surgery times was greater than 0.67 times the summation of the means of setup and cleanup times, while "low" means the opposite.

$$B = \frac{\sum_{j=1}^n \sum_{k=1}^K sur_j^k}{\sum_{j=1}^n \sum_{k=1}^K pre_j^k + post_j^k}$$

We also defined "cost ratio" as waiting cost divided by idle cost. The cost ratio was set to low (value=1) and high (value=2) similar to Batun et al. (2010).

**Table 1:** Experiment Design

Factor	Possible Value
Number of Scenarios	10, 50, 100
B Parameter	High , Low
Cost ratio	High , Low

In order to construct instances that reflect the complexity of real problems we gathered data from our industry partner on surgeons that performed their surgeries in parallel ORs (most commonly orthopedic and eye surgeries). Next we constructed empirical distributions for setup, surgery and cleanup times based on the type of surgery and the surgeon. Unfortunately we were not able to get sufficient data at the level of individual surgeons and procedures to directly generate (from data) sufficient scenarios to run our experiments. We next took data for all surgeons and surgeries provided by the hospital and fit a regression model between the mean and the standard deviation for setup, surgery and cleanup times. We utilized log-normal distributions to model the surgery time as suggested in May et al. (2000). We also used log-normal distributions to model setup and cleanup times as well because we observed long-tails in the histograms constructed from real data. To create a single problem instance, we began with a set of surgeries performed by a single surgeon in parallel ORs in our partner hospital. For each individual procedure in the problem instance, we set the mean setup, surgery and cleanup times equal to the average in the data set for that surgeon and procedure. We next used the regression model to get standard deviations. Finally we used lognormal distributions with the given means and standard deviations to generate scenarios for setup, surgery, and cleanup times.

When the setup to surgery time ratio was high, we found that the parallel OR scheduling problem becomes much more difficult to solve. We also looked carefully at the performance of our Benders' decomposition based algorithm for solving the surgery sequence sub-problem.

Interestingly, despite the enhancements discussed in section 4, the Benders' algorithm was not faster than Cplex branch and cut. The Cplex settings we used were: (1) generalizing upper bound branching, (2) mipemphasis = emphasize moving best bound and (3) the initial starting algorithm was "barrier". The disappointing results for the Benders' algorithm appear to be caused by the fact that the technology matrix is not fixed (the coefficients of this matrix are scenarios for the different times modeled) and also because the right hand sides were usually zero. We observed that the Benders' cuts were weak. Further we saw that the number of Benders' cuts required to solve each subproblem varies with respect to the number of scenarios.

**Table 2:** Computational Results (average processing time in seconds)

Factors	Scenarios		
	10	50	100
B Parameter, Cost Ratio	10	50	100
High , Low	278	2007	6146
High , High	240	1599	4844
Low, Low	188	1724	5887
Low, High	146	1554	4403

We see in table 2 that experiments with more scenarios, low cost, and high B parameter take more time to solve.

## 5.1 Comparison of stochastic solution to a common sense heuristic.

Table 3 shows the percentage optimality gap for the solutions obtained by solving using the mean value problem assuming a zigzag OR sequence. We speculate that this “zigzag mean value problem” is precisely the approach an OR manager would likely take when sequencing parallel ORs. Thus by observing the results in table 3 we have a rough estimate of how much savings are possible over and above an adept human scheduler.

**Table 3:** Optimality gap for the mean value solution with zigzag OR Sequence.

Factors	Scenarios		
	10	50	100
B Parameter, Cost Ratio	104%	44%	33%
High , Low	104%	44%	33%
High , High	109%	55%	44%
Low , Low	124%	55%	45%
Low , High	115%	56%	50%

Table 4 shows the optimality gap for the perturbed mean value heuristic discussed previously. We see that the heuristic closes about half the optimality gap from the simple zigzag mean value solution.

**Table 4:** Optimality gap for the perturbed mean value solution with zigzag OR sequence.

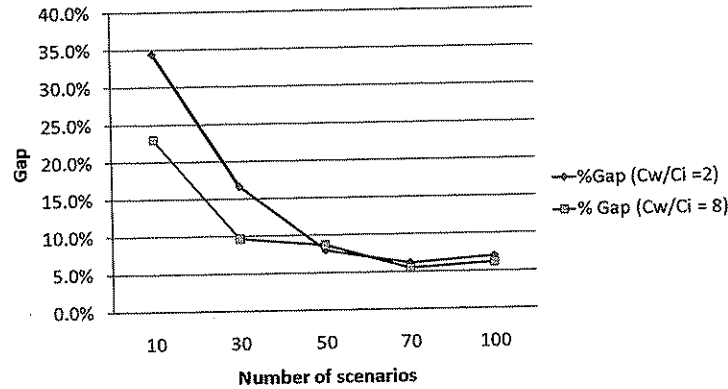
Factors	Scenarios		
	10	50	100
B Parameter, Cost Ratio	52%	28%	23%
High , Low	52%	28%	23%
High , High	56%	31%	29%
Low, Low	56%	28%	23%
Low, High	56%	29%	25%

## 6 Choosing the number of scenarios

In order to determine an appropriate number of scenarios, we applied the method proposed by Linderoth et al. (2006). The main idea is to generate several samples of scenarios that allow us to compute statistical lower and upper bounds on the true (infinite scenario) objective function value. The confidence interval on these estimates is strictly related to the number of samples used. Once we have computed these bounds for each set of scenarios, we can compute a statistical gap. The basic trade off is that as we increase the number of scenarios we will decrease the gap. The following plot shows how the gap decreases as we increase the number of scenarios. We worked with 2 instances both having the same number of surgeries and the same set of scenarios, but with different cost ratios (2 and 8). We generated 10 replicates of sets of scenarios with the same surgery, cleanup and setup duration distributions and two

different cost ratios to see if the cost ratio had an effect on the gap. In figure 3 we observe that the cost ratio did not have a large effect with 50 or more scenarios. Further we see, that for this instance, 100 scenarios is enough to achieve a 5% statistical gap.

**Figure 3:** Tradeoff between number of scenarios and statistical gap.  
**Choosing the number of scenarios**



## 7 Managerial insight

In this section we conduct experiments that provide insight about parallel OR scheduling and when it is a viable alternative. Our interest is in further exploring the tradeoff between surgeon waiting time and OR idle time when we schedule in parallel ORs (rather than a single OR). This tradeoff will depend on the relative costs of surgeon waiting time and OR idle time and on the setup to surgery time ratio  $B$ . Based on data collected from our industry partner the  $B$  factor typically ranges from 0.57 to 1.52 in cases where parallel ORs were used. We generated problem instances with varying  $B$  factor to investigate its affects. We started with a base problem with  $n=10$  procedures and  $K=100$  scenarios. The means and variances were generated as described in section 5. In order to generate problems with varying  $B$  we altered these as follows:

Let  $pre_j^k$ ,  $sur_j^k$  and  $post_j^k$  be the original setup, surgery, and cleanup times for a particular surgery  $j$  in scenario  $k$ .

Let  $D$  be the total duration ( $pre_j^k + sur_j^k + post_j^k$ )

Let  $C = \frac{pre_j^k}{post_j^k}$

Let  $B$  be the desired ratio.

Let  $G, H, P$  be the new setup, surgery, and cleanup times that have the desired  $B$

We have three equations:

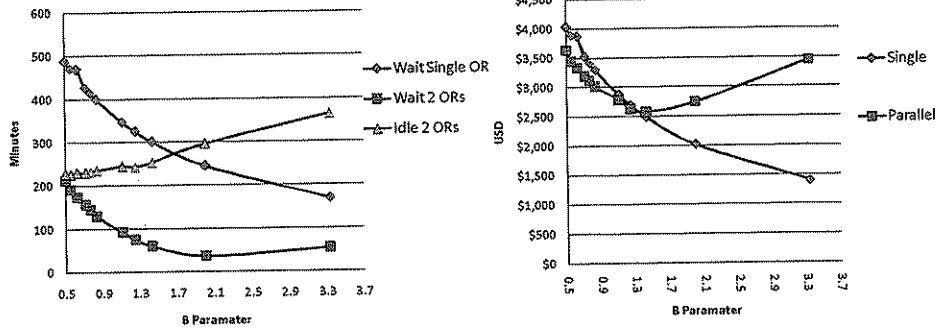
$$G+H+P=D \text{ (total duration must remain the same)}$$

$$\frac{G}{P} = C \text{ (the pre to post ratio must remain the same)}$$

$$\frac{H}{G+P} = B \text{ (to get the desired B ratio)}$$



**Figure 4: Comparison between single and parallel, cost ratio 1**  
**Parallel versus  $C_w=C_i$**  **Penalized by USD/min  $C_w=C_i$**



These equations can be easily solved to yield:

$$G = \frac{CD}{BC + C + B + 1}$$

$$H = \frac{BD(C + 1)}{BC + C + B + 1}$$

$$P = \frac{D}{BC + C + B + 1}$$

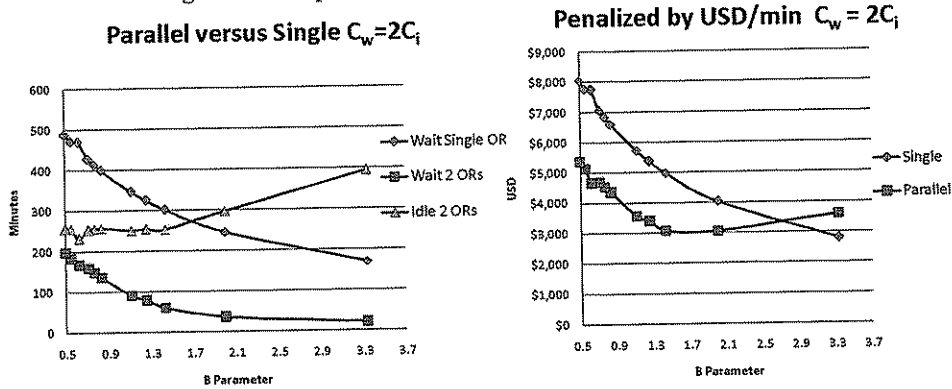
Next we solved the parallel OR scheduling problem for values of B from 0.5 to 3.32.

We assumed that any idle time in either OR throughout an 8 hour work day is penalized. The reason for this assumption is that we wanted to penalize the total idle time in both ORs. This assumption was necessary for a somewhat esoteric reason. Our first thought was to stop charging for idle time once surgeries are complete in an (either) OR. With this scheme, we frequently obtained solutions where only the second surgery was performed in OR 2 while all other surgeries were performed in OR 1. This type of solution occurred in cases where using a single OR was clearly more efficient. In these cases, scheduling the second surgery in OR 2 was cheaper than scheduling all surgeries in OR 1 since we did not charge for idle time after surgeries in OR 2 are complete. To avoid this problem, we charged for idle time up to 8 hours which is roughly when surgeries should be completed if both ORs are fully utilized throughout the schedule.

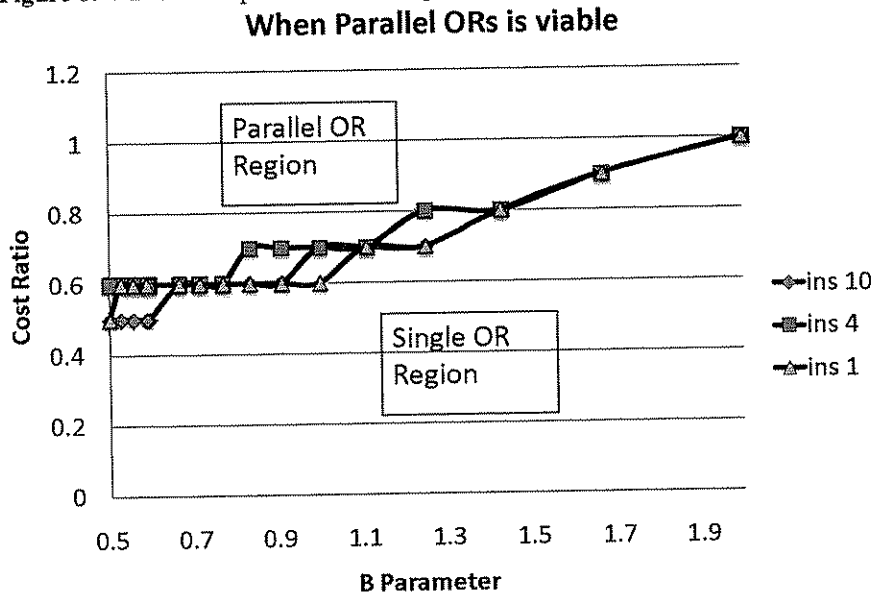
In the first experiment we assumed that surgeon waiting cost is equal to OR idle cost, and recorded the idle time and waiting time for each value of B. For purposes of comparison, we also computed surgeon waiting time in the case where all procedures are scheduled in a single OR. Note that with a single OR, the idle time is always zero. The results appear in Figure 4.

Next we repeated the experiment, with the assumption that surgeon waiting time is twice as expensive as OR idle time (cost ratio=2). The results appear in Figure 5.

**Figure 5:** Comparison between single and parallel, cost ratio 2



**Figure 6:** When to use parallel scheduling as a function of cost ratio and B parameter.



In figure 4 we observe a cutoff value for the B parameter before which parallel ORs is cost effective is about 1.4. This value changes with changing cost ratio. In figure 5 we observe that the cutoff value is around 2.7. Of course these curves were generated for a single problem instance.

Clearly the viability of parallel OR scheduling depends on the ratio B for a given set of surgeries and on the cost ratio  $c_w/c_i$ . In the next experiment we attempted to quantify when parallel OR scheduling is viable. We selected 3 test problems each with 100 scenarios and 10 surgeries from among those used in section 5. We next fixed the cost ratio  $c_w/c_i$  to a constant, then varied B until we found the value of B for which the cost of a single OR and parallel ORs was equal. This gave us one point on one of the curves in Figure 6 below. By varying the cost ratio  $c_w/c_i$  and repeating, we generated one full curve in Figure 6. Given the cost ratio for a particular hospital, one can use the graph in Figure 6 to find values of B for which parallel OR scheduling appears attractive.

## 7.1 Estimating the cost ratio

Finding a realistic value for the ratio of surgeon waiting time cost to OR idle time cost seems to be problematic. Hospitals seem not to have this information available, and various authors have suggested different methods for estimating it. Batun et al. (2010) estimated the cost ratio based on the relative cost of opening a new OR in terms of staff cost. They assumed two cases: (1) 50 minutes and (2) 250 minutes of surgeon waiting time is equivalent to opening an additional OR. With either of these assumptions they can compute a penalty in USD for every minute of surgeon waiting time. Here we propose a new method to estimate this cost ratio based on an analysis of historical data. We examined actual outcomes of surgical blocks in both single OR and parallel OR cases from our industry partner. Assuming the decision to assign a block to either single or parallel ORs to be rational (i.e. that which minimizes cost), we can use this data to compute bounds on the underlying cost ratio. When the OR scheduler decides to perform surgeries in a single OR, we assume he or she did so because it is less expensive than using parallel ORs. Under this assumption the following relationship holds (where  $w_1$  and  $s_1$  refer to waiting and idle time in a single OR and  $w_2$  and  $s_2$  refer to waiting and idle time in parallel ORs).

$$c_w w_1 + c_i s_1 < c_w w_2 + c_i s_2 \Rightarrow \frac{s_2 - s_1}{w_1 - w_2} > \frac{c_w}{c_i}$$

For the cases where the block was performed in parallel ORs we obtain, by a similar argument, a lower bound in the cost ratio.

$$c_w w_1 + c_i s_1 > c_w w_2 + c_i s_2 \Rightarrow \frac{s_2 - s_1}{w_1 - w_2} < \frac{c_w}{c_i}$$

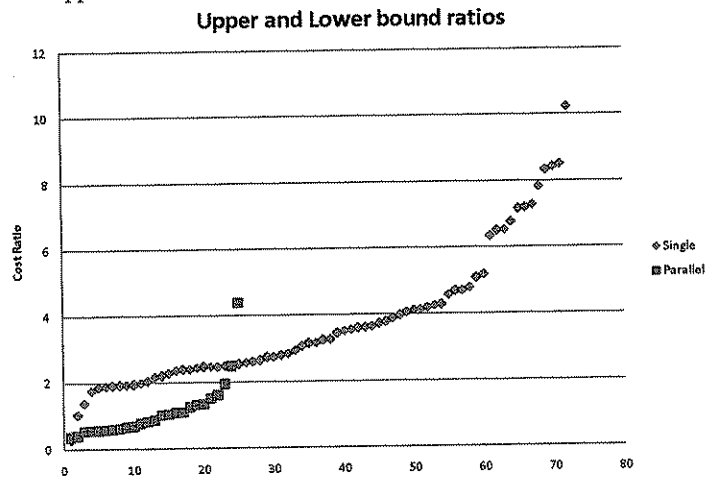
Given a single OR case, we compute  $s_2$  and  $w_2$  by assuming the same surgery sequence and the zigzag OR-sequence. Given a parallel OR case, we compute  $s_1=0$  and  $w_1$  by assuming that the surgeries with the largest setup and cleanup times are at the beginning and in the end positions in the surgery sequence respectively (note  $s_1=0$  in the single OR case). The plot in Figure 7 presents the computation of the cost ratio bounds for all surgical blocks in our data set.

While there are some overlaps that indicate all solutions may not have obeyed our rationality assumption, it appears from the plot that the implied cost ratio in our hospital is between 1 and 2. Of course this ratio may vary from hospital to hospital.

## 8 Conclusions

In this paper we developed a stochastic integer programming algorithm for sequencing surgeries for a single surgeon operating in parallel operating rooms under uncertainty. This decomposition based algorithm solves problems with sufficiently small run times for implementation on real size problems. We also provide an upper bounding heuristic that can be used to generate reasonable solutions when fast computation is required. This methodology may also be

**Figure 7:** Upper and Lower bounds in the cost ratio based on historical decisions.



useful in other application areas such as manufacturing and network computing. In particular, the lower bounds that were developed could be very useful when the number of jobs increases as one might expect in the other application areas. Using the algorithm to conduct experiments allowed us to develop insight regarding when the parallel OR approach is viable, what the basic cost tradeoffs are, and how one can measure the (typically unknown) implied cost ratio between surgeon waiting cost and OR idle time cost.

## References

- Abdekhodae, A. H., A. Wirth. 2002. Scheduling parallel machines with a single server some solvable cases and heuristics. *Computers & Operations Research* **29** 295–315.
- Abdekhodae, A. H., A. Wirth, Heng-Soon Gan. 2006. Scheduling two parallel machines with a single server: the general case. *Computers & Operations Research* **33** 994–1009.
- Batun, S., B. T. Denton, T. R. Huschka, Andrew J. Schaefer. 2010. The benefit of pooling operating rooms and parallel surgery processing under uncertainty. *To appear in Informs Journal on Computing* .
- Glass, C. A, Y M Shafransky, V A Strusevich. 2000. Scheduling for parallel dedicated machines with a single server. *Naval Research Logistics* **47**(1) 304 – 328.
- HFMA. 2005. Achieving operating room efficiency through process integration. *Health Care Financial Management Association Report* .
- Koulamas, C. P. 1996. Scheduling two parallel semiautomatic machines to minimize machine interference. *Computers & Operations Research* **23** (10) 945–956.
- Laporte, G., F.V. Louveaux. 1993. The integer l-shaped method for stochastic integer programs with complete recourse. *Operations research letters* **13** 133–142.
- Linderoth, J. T., A. Shapiro, S. J. Wright. 2006. The empirical behavior of sampling methods for stochastic programming. *Annals of Operations Research* **142** 219–245.
- Magnanti, T. L., R. T. Wong. 1981. Accelerating benders' decomposition: Algorithmic enhancement and model selection criteria. *Operations Research* **29** 464 – 484.
- May, J. H., D. P. Strum, L. G. Vargas. 2000. Fitting the lognormal distribution to surgical procedure times. *Decision Science* **31** 129–148.
- Nemhauser, G. L., L. A. Wolsey. 1999. *Integer and Combinatorial Optimization*. Wiley, New York.
- Storer, R. H., D. S. Wu, R. Vaccari. 1992. New search spaces for sequencing problems with application to job shop scheduling. *Management Science* **38**(10) 1495–1509.