

A Novel Integer Programming Approach with Separable Structure for Scheduling Job Shops

Nayanjyoti Baisya^{1*}, Shakti Prasad Jena²

^{1*} Assistant Professor, Department of Mechanical Engineering, Nalanda Institute of Technology, Bhubaneswar, Odisha, India

² Assistant Professor, Department of Mechanical Engineering, Nalanda Institute of Technology, Bhubaneswar, Odisha, India

*Corresponding author e-mail: nayanjyoti@thenalanda.com

Abstract: In this study, we examine the relevance of scheduling production lots effectively for producing medium- to high-volume goods that demand lengthy setup times. Instead of producing a full lot as is customary, lot splitting techniques separate a production lot into numerous smaller sub-lots, allowing each sub-lot to be "transferred" as soon as a step of production is finished. As a result, "transfer lots" drastically shorten lead times and decrease work-in-progress (WIP) inventory. However, it is very challenging to model, analyse, and govern transfer lots mathematically. In this research, a novel formulation of integer programming with separable structure is presented for job shop scheduling with fixed-size transfer lots. A solution approach built on the complementary concepts of Lagrangian relaxation (LR), backward dynamic programming (BDP), and heuristics is developed. Through explicit modeling of lot dynamics, transfer lots are handled on standard machines, machines with setups, and machines requiring all transfer lots within a production lot to be processed simultaneously. With the introduction of "sub-states" and the derivation of DP functional equations considering transfer lot dynamics, the standard BDP is extended to solve the lot-level sub problems. The recently developed "time step reduction technique" is also incorporated for increased efficiency. It implicitly establishes two time scales to reduce computational requirements without much loss of modeling accuracy and scheduling performance, thus enabling resolution of long horizon problems within controllable computational requirements. The method has been implemented using object-oriented programming language C++, and numerical testing results show that high quality schedules involving transfer lots are efficiently generated to achieve on-time delivery of products with low WIP inventory.

Keywords: Job Shop Scheduling; Transfer Lots; Optimization; Lagrangian Relaxation; Dynamic Programming.

1 INTRODUCTION

Scheduling of production lots is of great importance for manufacturing medium to high volume products with significant setups. In traditional production, a lot is indivisible, and individual pieces must wait for the completion of all the pieces within the lot before moving on to the next stage of operation. This often results in long lead times, low machine

utilization, and high work-in-process (WIP) inventory^[1]. In recent years, lot splitting techniques have been used to divide a production lot (briefly a 'lot') into multiple smaller sub-lots so that each sub-lot can be individually "transferred" from one stage of operation to the next as soon as that sub-lot has been completed. Lot splitting techniques thus allow individual "transfer lots" to be concurrently processed at several production stages, which

reduces manufacturing lead times, lowers WIP inventory levels, and improves product delivery times [2-4].

Most lot splitting techniques in the literature are based on heuristics. For example, two- and three- machine flow shop problems with equal-size transfer lots and a make span objective function were solved by using a modified Johnson's algorithm where each transfer lot was treated as an independent unit in Vickson and Alfredsson [5]. Heuristics for flow shop problems with three or more machines were presented in Trietsch and Baker [1]. A heuristic method for the integrated determination of transfer lot sizes and production schedules for a two-stage flow shop with a maximum flow time objective function was presented in Cetinkaya [6]. A simulation model for scheduling job shops with lot splitting using dispatching rules and a mean flow time objective function was presented in Jacobs and Bragg [7]. A heuristic algorithm is developed to minimize the make span for three-stage production processes in Glass et al. [8]. These heuristic approaches usually generate feasible schedules quickly, and demonstrate the benefits of transfer lots toward reducing lead times and lowering WIP inventory. However, it is difficult to evaluate the quality of the schedules generated, and these heuristics do not provide a systematic way for iterative improvement of the schedules. Recently, an optimization-based method has shown promise in scheduling transfer lots on standard machines without setups [9]. The method, however, cannot handle problems with long planning horizon or machine setups. There are also situations where all transfer lots within a production lot are required to be processed simultaneously. An example is the outsourcing of the entire lot, say, for heat treatment. Extension of the method is therefore needed to solve practical problems involving transfer lots.

Building on the above results of Liu and Luh [9], a novel integer programming formulation with separable structure for scheduling job shops with fixed-size transfer lots is presented in Section 2. In the formulation, transfer lots are handled on standard machines, machines with setups, and machines requiring all transfer lots within a production lot to be processed simultaneously. The formulation is "separable" in the sense that the objective function and all "coupling" machine capacity constraints are additive in terms of the basic decision variables at the lot level. A solution methodology based on a synergistic combination of Lagrangian relaxation (LR), backward dynamic programming (BDP), and heuristics is developed

in Section 3. Through the explicit modeling of lot dynamics, the introduction of "sub-states," and the derivation of dynamic programming equations considering transfer lot dynamics, the standard BDP is extended to solve lot-level sub problems within the Lagrangian relaxation framework. The recently developed "time step reduction technique" is also incorporated. It implicitly establishes two time scales to reduce computational requirements without much loss of modeling accuracy and scheduling performance, which enables the resolution of long horizon problems within reasonable computational requirements. Numerical testing results presented in Section 4 show that high quality schedules with transfer lots are generated in a timely fashion for on-time delivery with low WIP inventory.

2 PROBLEM FORMULATION

The following formulation for scheduling job shops with transfer lots is built on the previous work presented in Luh and Hoitomt [10] and Liu and Luh [9]. Instead of treating individual transfer lots as independent scheduling units as in Vickson and Alfredsson [5], only each production lot is treated as a scheduling unit, with operation beginning and completion times as decision variables. The key is to properly describe transfer lot dynamics using production lot variables only. The preliminaries that lead to the formulation are presented first, using a list of symbols given in Appendix C.

Notation and General Description

Time Step Reduction

In view of the long planning horizon under consideration relative to the time resolution required (e.g., 6 months vs. six minutes), the "time step reduction technique" originally developed in Luh et al. [1] is extended for scheduling transfer lots. The time horizon is divided into T "resolution increments" indexed by t , $0 \leq t \leq T-1$, and R consecutive resolution increments are aggregated into an "enumeration step" indexed by k , $0 \leq k \leq K$, with $T = R \times K$. For example, a 500 hour horizon can be divided into 500 one-hour resolution increments, and 50 ten-hour enumeration steps by aggregating 10 one-hour resolution increments in an enumeration step. Then an operation requiring, say, 16 hours on a machine is represented as occupying a full ten-hour enumeration step and 60% of the next enumeration step. Thus by using fractional but quantized machine utilization, multiple "short" operations are allowed to "share" a machine within

an enumeration step, and a part with several “short” operations is allowed to flow through the machines within a single enumeration step. Most input data are specified in terms of resolution increments except when stated otherwise. Since the complexity of the method depends significantly on the number of enumeration steps, this technique reduces computational effort in solving the “dual problem” through appropriate selection of the number of resolution increments R within an enumeration step.

2.1.2 Machines

In a job shop, machines may have different processing capabilities; e.g., processing speed or setup requirements. Machines with identical processing capability from the scheduling point of view are grouped as a “machine type,” and all the machine types form a set denoted by H .

2.1.3 Lots and Transfer Lots

Suppose that there are L production lots, indexed by $l = 0, 1, \dots, L-1$, each consisting of a number of products of the same type. For simplicity of presentation, a production lot will be referred to as a *lot* hereafter when there is no confusion. Different lots may have different product types, due dates, or arrival dates. For feasibility, a long enough planning horizon T is selected that is sufficient for the completion of all L lots.

A production lot, say lot l , consists of N_l fixed and equal-size transfer lots. It has to go through a sequence of operations, indexed by $j = 0, 1, \dots, J_l - 1$, according to a specified process plan. Operation j of lot l , denoted as (l, j) , has to be performed by a machine belonging to an eligible machine type $h \in H_{lj}$, $h = 0, 1, \dots, |H_{lj}| - 1$. Once started, the entire lot (i.e., all the transfer lots of the production lot) must be finished on the machine before anything else can be processed by the machine. This assumption applies to various situations, e.g., when setup costs are significant or when mixed transfer lots at a machine are difficult to manage because of operator or shop-floor tracking system requirements. Let t_{jih} denote the processing time in resolution increments per transfer lot for operation (l, j) on a machine type h . Let s_{ij} represent the required “time-out” in resolution increments between (l, j) and $(l, j+1)$, representing processes not explicitly modeled in the problem formulation, such as transfer time, cooling down or curing time. It is assumed that the number of transfer lots, N_l , the transfer lot processing time t_{jih} , and “time-out” s_{ij} are given.

2.1.4 Setups

If operation (l, j) has a setup requirement on a machine, it cannot be started until the machine has been setup. Assume that the setup time t_{ijh}^s in resolution increments for operation (l, j) on machine type h is known and is sequence independent (i.e., t_{ijh}^s does not depend on what was processed earlier on that machine). Then once the machine is setup for operation (l, j) , a transfer lot can be started on the machine as soon as it has arrived from its predecessor operation $(l, j-1)$ and the machine has finished the predecessor transfer lot if it exists. Also, as stated earlier, the machine cannot process anything else until all the transfer lots in lot l are finished.

2.1.5 Decision Variables

The beginning and completion times of operation (l, j) in resolution increments are denoted by b_{lj} and c_{lj} , respectively, and are the major decision variables. To ensure schedule acceptability, b_{lj} is constrained by its given earliest beginning time b_{lj}^e and the latest beginning time b_{lj}^l , i.e., $b_{lj}^e \leq b_{lj} \leq b_{lj}^l$; similarly, c_{lj} satisfies $c_{lj}^e \leq c_{lj} \leq c_{lj}^l$. These earliest and latest beginning and completion times are determined based on factors such as the arrivals of raw materials, the desire to minimize work-in-process inventory, and due dates promised to customers.

Machine Capacity Constraints and Lot Dynamics

Machine Capacity Constraints

The number of machines available per type at each resolution increment is a given integer. The average number of type h machines available at enumeration step k , denoted as M_{kh} , is thus a quantized fraction. The machine capacity constraints state that the total number of lots being processed should not exceed the number of machines available at each time period:

$$\sum_{l=0}^{L-1} \sum_{j=0}^{J_l-1} \delta_{ljkh} \leq M_{kh} \quad (1)$$

In the above, δ_{ljkh} is the fraction of time that operation (l, j) of production lot l is assigned to machine type h at enumeration step k . It is also assumed that the capacity of machine types where all transfer lots of a lot must be simultaneously processed is large enough to accommodate the entire lot.

The dynamics of transfer lots is described through

operation precedence constraints, processing time requirements, and the setup requirements as follows.

Operation Precedence Constraints

Assume that a machine of type h has been setup for $(l, j+1)$ of lot l so that this operation can be started. The operation precedence constraints require that operation $(l, j+1)$ cannot be started until the predecessor operation (l, j) of the first transfer lot has been completed, i.e.,

$$b_{lj} + t_{ljh} + s_{lj} \leq b_{l,j+1} \quad (2)$$

Where s_{lj} is any required "time-out?" If the operation (l, j) is performed on machines that require simultaneous processing of all transfer lots within a lot, then the constraints become that $(l, j+1)$ cannot be started until the *last* transfer lot has completed operation (l, j) , i.e.,

$$c_{lj} + s_{lj} + 1 \leq b_{l,j+1} \quad (3)$$

The presence of "1" in (3) is due to our convention that when an operation begins in a period, it starts at the beginning of that period. However when it ends in a period, it finishes at the end of that period, thus occupying the entire period in both cases. This convention is often followed in practice.

2.2.3 Operation Processing Time Requirements

Each operation beginning time b_{lj} may be associated with multiple completion times c_{lj} since the completion time depends not only on the operation beginning time and the transfer lot processing time t_{ljh} , but also on the lengths of intermittent idling times. Despite the availability of a machine, intermittent idling between two transfer lots may exist since the next transfer lot may still be in processing at the previous stage if the processing time there is longer than the current one [1, 9]. The lengths of intermittent idling times depend on several factors, nevertheless, it is clear that if there is no intermittent idling between b_{lj}

and c_{lj} , then $c_{lj} = b_{lj} + N_l \times t_{ljh} - 1$; otherwise $c_{lj} = c_{l,j-1} + s_{l,j-1} + t_{ljh}$. Consequently, the operation completion time c_{lj} is given by

$$c_{lj} = \max \left(\begin{matrix} b_{lj} + N_l \times t_{ljh} - 1, \\ c_{l,j-1} + s_{l,j-1} + t_{ljh} \end{matrix} \right) \quad (4)$$

The significance of (4) is that, although the movement of individual transfer lots is not explicitly modeled, production lot beginning and completion

times can be accurately described by (2) and (4).

If operation (l, j) is performed on a machine that processes all transfer lots within a lot simultaneously, then all the transfer lots have the same beginning and completion times. The completion time thus depends only on the beginning time b_{lj} and the transfer lot processing time t_{ljh} , i.e.,

$$c_{lj} = b_{lj} + t_{ljh} - 1 \quad (5)$$

2.2.4 Setup Requirements

For an operation (l, j) requiring setups, although the first transfer lot has to wait for the completion of its predecessor operation $(l, j-1)$, the machine's setup can be started earlier. The actual setup beginning time for operation (l, j) , denoted by b'_{lj} , is

$$b'_{lj} = b_{lj} - t_{ljh}^s \quad (6)$$

Assuming that the machine is available for setup at b'_{lj} . Therefore once operation beginning time b_{lj} is known, the setup beginning time can be readily computed, and setups can be embedded within lot dynamics without introducing additional variables.

Objective Function

The objective of scheduling is to ensure on-time product delivery with low WIP inventory. This is represented by minimizing the sum of weighted quadratic penalties for violating lot due dates and for releasing raw materials too early:

$$J \equiv \sum_l (w T_l^2 + \beta E_l^2) \quad (7)$$

In the above, T_l is the tardiness of lot l defined as the time the lot completion time c_l (the completion time of the *last* operation of the *last* transfer lot) exceeds the given lot due date d_l (in enumeration steps, i.e., $T_l \equiv \max \left(0, \left\lceil \frac{c_l - d_l}{\lfloor R \rfloor} \right\rceil \right)$). The lot

earliness, E_l , is similarly defined as the excess of lot's earliest beginning time \bar{b}_l over the lot's scheduled beginning time b_l (the beginning time of the *first* operation of the *first* transfer lot), i.e., $E_l \equiv \max \left(0, \left\lfloor \frac{\bar{b}_l - b_l}{\lfloor R \rfloor} \right\rfloor \right)$. The parameters w and

β are weights associated with the earliness and tardiness penalties of lot l . The above penalties define a time window in which the lot can be scheduled without penalty.

The key decision variables are operation beginning times $\{b_{ij}\}$ and completion times $\{c_{ij}\}$ for individual production lots. Once these variables are determined, transfer lot beginning times and completion times can be derived as presented in Appendix A.

3 SOLUTION METHODOLOGY

Lagrangian relaxation (LR) is a mathematical programming technique for constrained optimization. Similar to the pricing concept of a market economy, the method replaces “hard” coupling constraints (i.e., machine capacity constraints in this study) by the payment of certain “prices” (i.e., Lagrange multipliers) based on the “demand” for a machine for the use of that machine at each time unit. The original problem can thus be decomposed into many smaller and easier lot-level sub problems. Backward dynamic programming is then used to solve these lot-level sub problems where other constraints are enforced. The multipliers are then adjusted after these sub problems are solved, based on the degrees of constraint violation following again the market economy mechanism. Sub problems are then re-solved based on the new set of multipliers, and the process repeats. In mathematical terms, the “dual function” is maximized in this multiplier updating process where the values of the dual function are lower bounds to the optimal feasible cost. Since coupling constraints have been relaxed by the multipliers, the solutions of individual sub problems, when put together, may not constitute a feasible schedule. A simple heuristic is therefore used towards the end of this multiplier updating process to provide feasible schedules satisfying all constraints. The quality of the feasible schedules can be quantitatively evaluated by comparing their costs with the largest lower bound provided by the dual function. The development of the BDP technique is complicated, and will be one of the major topics of our focus in this section.

The Lagrangian Relaxation Framework

Machine capacity constraints (1) are first “relaxed” by using Lagrange multipliers $\{\pi_{kh}\}$ in enumeration steps and the Lagrangian is formed as:

$$J^+ \equiv \sum_{l=0}^{L-1} (wT_l^2 + \beta E_l^2) + \sum_k \sum_h \pi_{kh} \left(\sum_{l=0}^{L-1} \sum_{j=0}^{J_l-1} \delta_{ljkh} + m_{kh} - M_{kh} \right) \quad (8)$$

Where m_{kh} is a non-negative slack variable satisfying $0 \leq m_{kh} \leq M_{kh}$. With the multipliers given, the “relaxed problem” is to minimize the Lagrangian J^+ subject to operation precedence constraints, processing time requirements, and setup requirements (2) to (6). After regrouping relevant terms within J^+ , the problem is decomposed into individually solvable lot sub problems as follows:

$$\min_{\{b_{ij}, c_{ij}\}} L_l, \text{ with } L \equiv wT^2 + \beta E^2 + \sum_{j=0}^{J_l-1} L_{ij} (b_{ij}, c_{ij}),$$

and

$$L_{ij} (b_{ij}, c_{ij}) \sum_{t=b_{ij}-t_{ijh}}^{c_{ij}} \pi_{th} \quad (9)$$

To simplify the derivation, π_{th} is introduced to represent the multiplier for machine type h at resolution increment t , i.e., $\pi_{th} \equiv \frac{\pi_{kh}}{R}$ with $k \equiv \begin{bmatrix} t \\ R \end{bmatrix}$. The sub problems $\{L_l\}$ are subject to (2) to (6). The decision variables are operation beginning times $\{b_{ij}\}$ and completion times $\{c_{ij}\}$ of lot l .

Since $L_{ij}(b_{ij}, c_{ij})$ is the cost for using a type h machine between times b_{ij} and c_{ij} , the lot sub problem thus reflects the balance between the machine utilization cost and tardiness as well as earliness penalties.

Backward Dynamic Programming (BDP) for Lot Sub problems: Preliminaries

In view of the complications caused by the existence of multiple completion times associated with a given beginning time, the BDP algorithm developed in Luh et al. [11] must be extended to solve lot sub problems. In the following, the generic DP equations are first presented. Several key parameters are then determined, including the number of multiple operation completion times associated with each beginning time, the earliest beginning and completion times, and the latest beginning and completion times. These parameters help reduce BDP computational requirements. In addition, “sub-states” are introduced to efficiently carry out the DP procedure.

Backward Dynamic Programming Equations

Each lot sub problem has a number of DP stages, where each stage corresponds to an operation. The BDP algorithm starts with the last stage and moves backwards in time. The states for a stage correspond

to possible beginning times in Luh et al. [11]. In view of the multiple completion times associated with a given beginning time, a state in this study is represented by a pair of operation beginning and completion times (b_j, c_j) . The cumulative cost at (b_j, c_j) , denoted by $V_j(b_j, c_j)$, is obtained as the sum of the stage-wise cost $L_j(b_j, c_j)$ (plus tardiness penalty for the last stage and earliness penalty for the first stage) and the minimum cumulative cost of a reachable state $(b_{l,j+1}, c_{l,j+1})$ at the successor stage.

To be more specific, the BDP procedure starts with the last stage having the following terminal cost:

$$V_{l,j}(b_{l,j}, c_{l,j}) \equiv wT^2 + L_{l,j}(b_{l,j}, c_{l,j}) \quad (10)$$

The cumulative cost when moving backwards to the predecessor stage is then obtained recursively according to the following BDP equation subject to the constraints (2) to (6):

$$V_j(b_j, c_j) \equiv \min_{\{(b_{l,j+1}, c_{l,j+1}) \in S_{l,j+1}(b_j, c_j)\}} \{L_j(b_j, c_j) + V_{l,j+1}(b_{l,j+1}, c_{l,j+1})\} \\ = L_j(b_j, c_j) + \min_{\{(b_{l,j+1}, c_{l,j+1}) \in S_{l,j+1}(b_j, c_j)\}} \{V_{l,j+1}(b_{l,j+1}, c_{l,j+1})\} \\ 1 \leq j < J-1 \quad (11)$$

In the above, $S_{l,j+1}(b_j, c_j)$ is the set containing all allowable $\{(b_{l,j+1}, c_{l,j+1})\}$ satisfying (2) to (6) for the given (b_j, c_j) , and is determined based on possible state transitions.

The equation for the first stage is given by

$$V_{10}(b_{10}, c_{10}) = \beta_1 E_1^2 + L_{10}(b_{10}, c_{10}) + \min_{\{(b_{11}, c_{11}) \in S_{11}(b_{10}, c_{10})\}} \{V_{11}(b_{11}, c_{11})\} \quad (12)$$

Let L_1^* represent the minimal lot sub problem cost for lot l . The minimum L_1^* can be then obtained as the minimal cumulative cost at the first stage, i.e.,

$$L_1^* \equiv \min_{\{(b_{10}, c_{10})\}} \{V_{10}(b_{10}, c_{10})\} \quad (13)$$

Finally, the optimal beginning times and completion times can be obtained by tracing forward along the

stages.

Multiple Completion Times Associated with a Beginning Time

For a given operation beginning time, the operation completion time depends on the lengths of intermittent idling times between transfer lots. As a result, there may be multiple completion times associated with a given beginning time. A detailed analysis of this complicated phenomenon can be found in Liu and Luh [9]. To deal with transfer lots effectively, a forward procedure is introduced here to determine these multiple completion times. For operation (l, j) , the number of multiple completion times associated with a beginning time, denoted by $N_j^c + 1$, is determined by considering the machine type being used, its processing time, and the processing times of its previous operations, as given in Appendix B.

The Earliest and Latest Operation Beginning and Completion Times

For a given planning horizon, in view of constraints (2) to (6), the beginning time as well as the completion time of each operation must have the earliest (least) value and the latest (largest) value, respectively. Moreover, these parameters, the earliest and latest operation beginning and completion times, are helpful in effectively limiting the computational effort in BDP. Thus the earliest beginning and completion times are determined

to ensure that every operation can be started and completed as early as possible and that these earliest times satisfy constraints (2) to (6) within the given time horizon. The earliest beginning and completion times of the current operation can be obtained recursively by proceeding forward from its predecessor operation while enforcing constraints (2) to (6).

Similarly, computed values of the latest beginning and completion times ensure that the latest completion times of their successor operations are still within the planning horizon T and they satisfy constraints (2) to (6). The latest beginning and completion time of the current operation can be calculated recursively by proceeding backward from its successor operation while enforcing constraints (2) to (6).

Backward Dynamic Programming Structure

Similar to Luh et al. [11], the DP stages correspond to operations and DP states to the possible operation beginning times. The numbers of DP states are determined by the earliest and latest beginning times. Since transfer lots lead to multiple completion times

for a beginning time, DP sub states are introduced to consider these multiple values.

Let $N_{lj}^c + 1$ denote to the number of multiple completion times associated with a beginning time for a stage (l, j) and a DP sub state be a pair of the feasible operation beginning time and one of its associated completion times. The sub state in a state is indexed by $0, 1, \dots, N_{lj}^c$; sub state 0, sub state 1,

\dots , and sub state N_{lj}^c correspond to the pairs of the beginning time and the completion time between which there is zero unit, one unit, \dots , and N_{lj}^c units intermittent idling time, respectively. For simplicity (without loss of generality), sub state 0 in a state is considered to coincide with that state itself, i.e., it is not necessary to give the corresponding completion time explicitly because there is no (zero) intermittent idling.

Fig. 1 shows a sample of stages, states, sub states, and state transitions, where a sub state in a state is represented in the format: (stage index, state index, sub-state index) and state transitions follow constraints (2) to (6). For example, (1, 0, 2) represents the second sub state in state 0 at stage 1. Based on this structure, the BDP procedure for solving a sub problem is now presented.

Backward Dynamic Programming (BDP) for Lot Sub problems

Unlike FDP, where we move from first to last

stage, in BDP the cumulative cost, given by (11), is computed from the last state b_{lj}^l to the first state b_{lj}^e (indexed by $N_{lj}^b - 1$ to 0), and for each state, from the last sub state to the first sub state (indexed by N_{lj}^c to 0). For efficiency, determining the set $S_{l,j+1}(b_{lj}, c_{lj})$ of feasible successor states is critical.

3.3.1 Set of all Feasible States at Successor Stage
 $S_{l,j+1} \left(\begin{matrix} b_{lj} \\ c_{lj} \end{matrix} \right)$

Let parameters t_b and t_c denote the earliest possible beginning and completion times of $(b_{l,j+1}, c_{l,j+1}) \in S_{l,j+1} \left(\begin{matrix} b_{lj} \\ c_{lj} \end{matrix} \right)$, respectively. In view of the constraints (2) to (6), the times t_b and t_c can be calculated as follows. If the operation ($l, j+1$) is processed on standard machines and machines with setup, then

$$t_b = b_{lj} + t_{lh} + s_{lj},$$

$$t_c = \max \left(\begin{matrix} t_b + N_1 \times t_{l,j+1,h} - \\ 1, c_{lj} + s_{lj} + t_{l,j+1,h} \end{matrix} \right) \quad (14)$$

However, if all transfer lots are required to be processed simultaneously, then

$$t_b = c_{lj} + s_{lj} + 1, \quad t_c = t_b + t_{l,j+1,h} - 1 \quad (15)$$

For the given (b_{lj}, c_{lj}) , considering the operation processing time requirements (4) and (5), the completion time $c_{l,j+1}$ can be uniquely determined

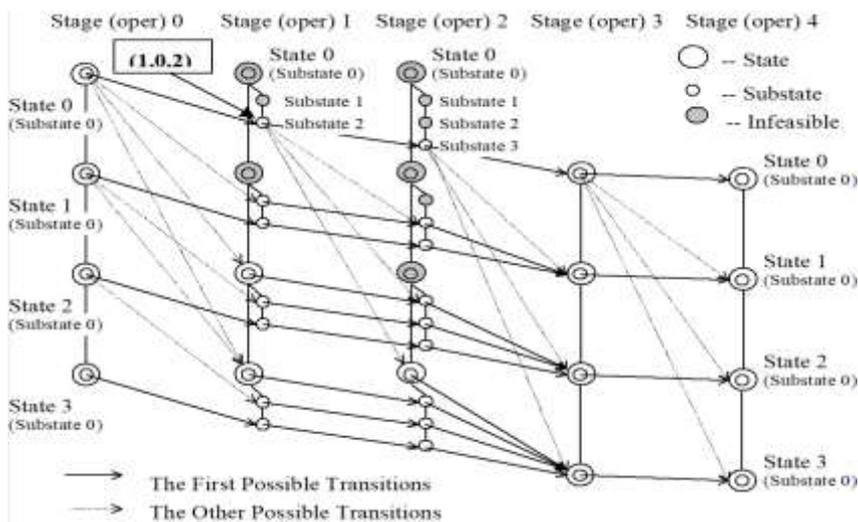


Fig. 1: DP Stages, States, Sub-states, and Transitions

from each possible beginning time $b_{l,j+1}$. Since $(b_{l,j+1}, c_{l,j+1}) \in S_{l,j+1}(b_j, c_j)$ may have intermittent idling times, the set $S_{l,j+1}(b_j, c_j)$ can be decomposed into two disjoint subsets – with and without idling times – as:

$$S_{l,j+1}^1(b_j, c_j) \equiv S_{l,j+1}^1(b_j, c_j) \cup S_{l,j+1}^2(b_j) \quad (16)$$

where

$$S_{l,j+1}^1(b_j, c_j) \equiv \left\{ \begin{array}{l} (b_{l,j+1}, c_{l,j+1}) : \\ t_b \leq b_{l,j+1} \leq t_b + N_{l,j+1}^c, \\ c_{l,j+1} = t_c \end{array} \right\} \quad (17)$$

and

$$S_{l,j+1}^2(b_j) \equiv \left\{ \begin{array}{l} b_{l,j+1} : t_b + N_{l,j+1}^c + 1 \\ 1 \leq b_{l,j+1} \leq b_{l,j+1} \end{array} \right\} \quad (18)$$

In the first set defined by (17), every element has intermittent idling times, and there are a total of $N_{l,j+1}^c$ elements with the same $c_{l,j+1}$ (equal to t_c). In the second set defined by (18), each element is a sub state with index 0 because of no intermittent idling time, and is represented by the state itself. For example, as shown in Fig. 1, for the sub state (1, 0, 2), its (t_b, t_c) is represented by the sub state (2, 0, 3); the transition of the sub state (1, 0, 2) is represented by $S_{l_2}^1(0,2) = \{(2, 1, 2), (2, 2, 1)\}$ and $S_{l_2}^2(0) = \{(2, 2, 0), (2, 3, 0), \dots, (2, N_{l_2}^b - 1, 0)\}$, where $N_{l_2}^b$ represents the number of total states in stage 2.

3.3.2 Stage wise Minimum Cumulative Cost of the Successor Stage

In computing the cost $V_j(b_j, c_j)$ by (11), the major computational work is to find the minimum cumulative cost among all elements in $S_{l,j+1}(b_j, c_j)$. The following computation-critical equations are developed to obtain the minimum with the least amount of computation. Let $F_{l,j+1}(t_b, t_c)$ denote the minimum cumulative cost of the successor stage $(l, j+1)$ within time (t_b, t_c) among all possible $(b_{l,j+1}, c_{l,j+1}) \in S_{l,j+1}(b_j, c_j)$, i.e.,

$$F_{l,j+1}(t_b, t_c) \equiv \min_{\{b_{l,j+1} \geq t_b, c_{l,j+1} = \max(t_c, b_{l,j+1} + N_{l,j+1}^c - 1)\}} V_{l,j+1}(b_{l,j+1}, c_{l,j+1}) \quad (19)$$

In view of (15) and (16), $F_{l,j+1}(t_b, t_c)$ can be recursively obtained by comparing the cost $V_{l,j+1}(b_b, t_c)$ at time t_b and the minimum of the costs

$(F_{l,j+1}(t_b + 1, t_c), F_{l,j+1}(t_b + 1, t_c + 1))$ at time $t_b + 1$. This can be expressed as follows. If the operation $(l, j+1)$ is processed on the standard machines or machines with setups, then

$$F_{l,j+1}(t_b, t_c) = \min[V_{l,j+1}(t_b, t_c), F_{l,j+1}(t_b + 1, t_c)],$$

if

$$t_c > t_b + N_{l,j+1} \times t_{yjh} - 1 \quad (20)$$

$$F_{l,j+1}(t_b, t_c) = \min[V_{l,j+1}(t_b, t_c), F_{l,j+1}(t_b + 1, t_c + 1)],$$

if

$$t_c = t_b + N_{l,j+1} \times t_{yjh} - 1 \quad (21)$$

if all transfer lots are required to be processed simultaneously, then

$$F_{l,j+1}(t_b, t_c) = \min[V_{l,j+1}(t_b, t_c), F_{l,j+1}(t_b + 1, t_c + 1)] \quad (22)$$

The significance of (20) to (22) is that in most situations only one step comparison is needed to obtain $F_{l,j+1}(t_b, t_c)$. As a result, the cumulative cost in (11) can be rewritten as:

$$V_j(b_j, c_j) = L_j(b_j, c_j) + F_{l,j+1}(t_b, t_c) \quad (23)$$

The computational complexity of the above BDP $O(\sum_j N_j^c T)$ is the same as that of FDP analyzed in Luh and Hoiomt [10].

3.3.3 Required Setups and Simultaneous Processing of all Transfer Lots

The setup times t_s for transfer lots are not explicitly expressed in the above BDP equations; they are considered in the computation of stage wise state cost $L_j(b_j, c_j)$ by (9). Thus setup requirements have little effect on BDP computational efficiency.

However, as given by (20) to (22), BDP equations are affected by the machines requiring simultaneous

processing of all transfer lots in a lot. Simultaneous processing is required, for example, in heat treatment operation that is subcontracted out. Here all production lot units, that is, all transfer lots are sent to the subcontractor together for the heat treatment operation. But for operations processed on these machines sub states are only with index 0, not adding much computational burden in the BDP procedure.

Slack Variable Sub problems and the Dual Problem

Slack Variable Sub Problems

Little effort is needed to solve the slack variable sub problems because these sub problem solutions only require summations of the multipliers:

$$\min_{m_{kh}} L_s, \text{ with } L_s = \sum_{k,h} \pi_{kh} m_{kh} \quad (24)$$

3.4.2 The Dual Problem

With the optimal costs of lot sub problems and slack variable sub problems given by $\{L_1^*\}$ and L_s^* , respectively, the high-level dual problem, denoted by D, is obtained as

$$\max_{\{\pi_{kh}\}} D, \text{ with } D \equiv \sum_1 L_1^* + L_s^* - \sum_{k,h} \pi_{kh} M_{kh} \quad (25)$$

Since the dual function D is concave, piece-wise linear, and consists of many facets, the sub gradient method is commonly used to solve it. However this method suffers from slow convergence. To overcome this, instead of solving all sub problems for multipliers updates the Interleaved Sub-Gradient (ISG) method was suggested that update multipliers after solving *each* sub problem. Furthermore, since the dual function approaches a smooth function as the problem size increases, the Conjugate Gradient (CG) methods have more attractive convergence properties for such problems. Therefore, the newly developed Interleaved Conjugate Gradient (ICG) method that incorporates the “interleave” concept with the CG method can provide faster convergence (Zhao et al. [11]). It is used to update the multipliers in this study.

As presented in Luh et al. [10], a rough estimate of the number of multipliers when there are K enumeration steps is $\pi_{kh} : K \times |H|$. With the “time step reduction” technique, we do not need to have a multiplier for each of the resolution steps T, thus reducing the computational complexity for solving the high-level dual problems significantly.

Heuristics

The computation of sub problem solutions and multiplier updates are stopped after a fixed amount of computation time or a fixed number of iterations, where iteration consists of solving all the sub problems once. Since machine capacity constraints have been relaxed, solutions of sub problems, when put together, generally do not constitute a feasible schedule. A simple heuristic procedure is usually used to adjust sub problem solutions to form a feasible schedule. The heuristics developed for transfer lots are based on the version developed by Luh and Hoitomt [9] and Luh et al. [10].

The heuristics start with the solutions of the DP sub-problems that give the actual operation beginning times and completion times. Since the setup for an operation is always immediately followed by the start of that operation on the first transfer lot, this setup is scheduled based on the availability of the desired machine and the start time of the current operation of the first transfer lot. For the first transfer lot in a production lot, its current operation can be started after the completion of its predecessor operation. For other individual transfer lots in the same lot, the current operation for a transfer lot can be processed as soon as the predecessor operation of this transfer lot and the current operation of the predecessor transfer lot are complete, together with the consideration of machine availability. For the operation processed on the machine requiring simultaneous processing of all transfer lots in the same lot, since such machine capacity is treated as large, this operation can begin after the completion of the predecessor operation of the last transfer lot.

For simplicity of presentation, we have not given the details of the situation where multiple machine types can do a given operation; the BDP can be easily extended to this situation by considering multiple stages for an operation, one for each applicable machine type. The quality of a schedule obtained is quantitatively evaluated by its relative duality gap which is the relative difference between the feasible schedule cost J and its lower bound the dual value D; i.e., Duality Gap $\equiv (J-D)/D \times 100\%$. The stopping criterion for the solution procedure may be to obtain a given duality gap within an acceptable range.

4 NUMERICAL TESTING RESULTS

The current algorithm that combines BDP and ICG

Tab. 1: Data and Results for Case 1: Transfer Lots

Lot l (N_l)	Oper. j	Machine h	t_{ljh}	d_l	w_l	Schedule	Schedule
						1.1 (b_{lj}, c_{lj})	1.2 (b_{lj}, c_{lj})
0 (5)	0	M0	2			(12, 21)	(6, 15)
	1	M1	1			(14, 22)	(16, 20)
	2	M2	2	1	1	(15, 24)	(21, 30)
1 (2)	0	M0	3			(0, 5)	(0, 5)
	1	M1	1			(3, 6)	(6, 7)
	2	M2	2	0	1	(5, 8)	(8, 11)
2 (2)	0	M1	1			(0, 1)	(0, 1)
	1	M2	2			(1, 4)	(2, 5)
	2	M0	3	1	1	(6, 11)	(16, 21)

Tab. 2: Scheduling Performance for Case 1: Transfer Lots

Schedule	Make span	Average lead-time	Average WIP inventory	Average machine utilization (%)	Average tardiness
Schedule 1.1: Transfer Lots	25	6.4	0.26	65.3	15
Schedule 1.2: Without Transfer Lots	31	16	0.52	52.7	21

within the LR framework has been implemented using the object-oriented programming language C++, and extensive initial testing has been performed on a Sun Sparc 10 workstation. Four test cases are presented below to evaluate the performance of the method developed. The first case shows that using transfer lots improves the scheduling performance greatly. Case 2 illustrates that the transfer lots can be scheduled effectively on various types of machines: standard machines, machines with setups, and machines requiring all transfer lots in the same lot to be processed simultaneously. This increases the applicability of the model to a substantially larger set of realistic environments. Case 3 and Case 4 demonstrate the capability of the method developed for scheduling real problems with different sizes (number of lots, parts, transfer lots, and time horizon). All cases assume that all machines are available throughout the planning horizon, starting from period zero. Using heuristics, feasible planning horizons are initially generated based on machine availability and lot processing time requirements. For the first three cases, the enumeration step is equal to the resolution time unit (i.e., $R = 1$). In addition to the duality gap, the following practical metrics, used by various industries, are also applied to evaluate the scheduling performance.

The metrics are: Make span = $\max_{l,n} c_{l,n}^n - \min_{l,n} b_{l,n}^n + 1$; Lead-time of the n th transfer lot in lot $l =$

$c_{l,n}^n - b_{l,n}^n + 1$; WIP inventory of the n th transfer lot in lot $l =$ (lead-time of the n th transfer lot in lot l) / make span; Machine utilization of machine $h = \frac{\sum_{l,j} t_{ljh}}{\sum_{l,j} t_{ljh} + \sum_{l,j} s_{ljh}}$ / available time in makespan; and Tardiness (delivery delay) of the n th transfer lot in lot $l = \max(0, c_{l,n}^n - d_l)$. In the above, b_{lj}^n and c_{lj}^n represent the beginning time and completion time of the operation (l, j) of the n th transfer lot, respectively.

In addition, a comparison of the current algorithm with the common dispatching rules used in practice has been suggested by the editor. The resulting schedules using these heuristic dispatching rules are being computed and will be presented for examination by the readers via the internet.

Case 1 (Schedules With and Without Transfer Lots). This case is to show by a small example that, as expected, scheduling with transfer lots can indeed improve the scheduling performance greatly. There are three lots to be scheduled on three machine types, with one machine of each type. There are 5 parts in lot 0 and 2 in lots 1 and 2. The data is given in Tab. 1. The planning horizon is 40 time units (i.e., $T = 40$).

The problem is first solved with transfer lots of size one. Each part is, thus, treated as a transfer lot. The feasible schedule, Schedule 1.1, is generated with a

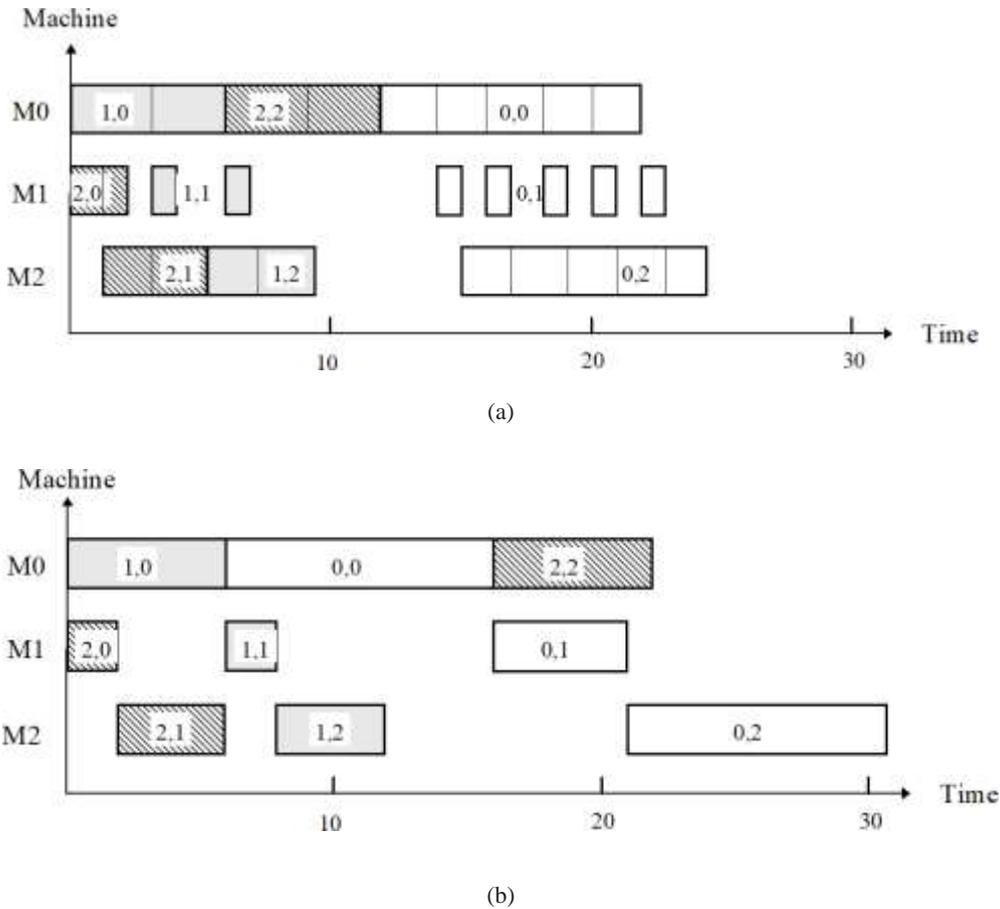


Fig. 2: (a). Gantt Chart of Schedule 1.1: Scheduling with Transfer Lots; (b). Gantt Chart of Schedule 1.2: Scheduling Without Transfer Lots

cost of 693 and a lower bound of 693, in less than CPU time of one second. Thus the solution found is the optimal solution. The operation beginning times and completion times of Schedule 1.1 are also shown in Tab. 1, with Gantt chart given in Fig. 2a.

Then the same problem is solved assuming all lots are indivisible (i.e., without using transfer lots). The feasible schedule, Schedule 1.2, has a cost of 1362 with a lower bound of 1361.87, and is obtained in less than one second CPU time. The resulting operation beginning times and completion times are also presented in Tab. 1. The Gantt chart of the Schedule 1.2 is shown in Fig. 2b.

Tab. 2 gives the metrics for the scheduling performance of both Schedule 1.1 and 1.2. Comparing with the schedule obtained without considering transfer lots (Schedule 1.2), transfer lots have significantly improved the average lead time, average WIP inventory, average machine utilization, and average delivery delay time. These

imply that the use of sub-lots of smaller size and the overlapping of consecutive operations results in less work in process and less product delivery delay.

Case 2 (Scheduling Transfer Lots with Various Machine Categories). This case is to show that the method presented here can effectively schedule transfer lots on three key machine categories: standard machines, machines with setups, and machines where all transfer lots in a lot must be processed simultaneously. There are four lots to be scheduled on five machine types, each machine type with one machine. Each lot has a different due date and Lot 0 has the highest priority (weight = 3) among the four lots while Lot 2 has a higher priority (weight = 2) than the other two. Lot 3 has an arrival time of 2 units, and each lot has four operations. Some operations processed on M0 and M2 need setups. For the operation processed on M4, all transfer lots in a lot must be processed simultaneously. The detailed data about lots and

Tab. 3: Data and Results for Case 2: Various Machine Types

Lot l (N_l)	Oper. J	Mach h	t_{ljh}	t_{ljh}^s	s_{lj}	a_l	d_l	w_l	Schedule 2 (b_{lj}, c_{lj})
0 (4)	0	M0	3	-	1	0			(6, 17)
	1	M2	2	4	-				(12, 20)
	2	M4	2	-	-				(21, 22)
	3	M1	1	-	-		0	3	(23, 26)
1 (2)	0	M1	2	-	-	0			(0, 3)
	1	M2	3	-	1				(2, 7)
	2	M0	2	2	-		1	1	(20, 23)
	3	M4	2	-	-				(24, 25)
2 (3)	0	M0	1	3	1	0			(3, 5)
	1	M1	2	-	-				(5, 10)
	2	M4	2	-	1				(11, 12)
	3	M3	5	-	-		2	2	(14, 28)
3 (3)	0	M3	2	-	-	2			(2, 7)
	1	M4	2	-	-				(8, 9)
	2	M2	1	2	-				(23, 25)
	3	M1	3	-	1		8	1	(27, 35)

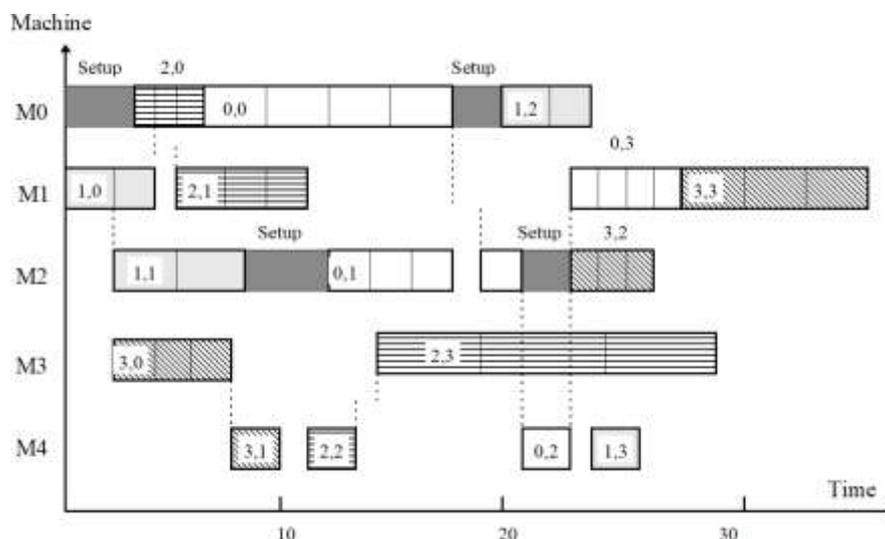


Fig. 3: Gantt Chart of Schedule 2: Scheduling Transfer Lots with Various Machine Types

operations are given in Tab. 3. “Time-outs” are also considered in this case.

The feasible schedule, Schedule 2, has a cost of 4740 with a relative duality gap 3.53%, and is obtained in 2 seconds. In Tab. 3, the resulting operation beginning times and completion times (

b_{lj}, c_{lj}) are also presented. The Gantt chart of the feasible schedule is shown in Fig. 3. The make span, average lead-time, average WIP inventory, average machine utilization, and average tardiness is 36, 21.4, 0.59, 55%, and 19.7, respectively. This once again shows that the schedules generated involving transfer lots on all three kinds of machines is of a

high quality.

The following two observations on the solution presented in Fig. 3 are made to illustrate how transfer lots are scheduled effectively on the machines with setups and machines requiring simultaneous processing of all transfer lots. First: the setup for operation (0, 1) is started at time 8 on M2 while the first transfer lot of lot 0 is still in process for its operation (0, 0) on M0. Thus once the first transfer lot completes its operation (0, 0) and is transferred to M2, the operation (0, 1) of the first transfer lot is started immediately without any delay for the machine setup. The requirement of one time unit “time-out” (transportation time, for example) between operation (0, 0) and (0, 1) can be also easily observed. Second: the operation (3, 1) is started on M4 after the operation (3, 0) of all transfer lots in lot 0 is finished. Once the operation (3, 1) is completed, the first transfer lot of lot 3 is moved to M2 and the operation (3, 2) is started as soon as the M2 is available. These two examples confirm the validity of lot dynamic modeling for different machine categories.

The Cases 3 and 4 draw data from a manufacturer producing aircraft/turbine-generator parts. According to the production requirements, the parts to be scheduled are grouped into a number of lots with various due dates and weights. Tab. 4 summarizes the test data for Cases 3 and 4.

Case 3 (A Real Problem Using Different Transfer-Lot Sizes). This case is to demonstrate the capability of the method developed for scheduling real problems using different transfer lot sizes. As given in Tab. 4, the planning horizon in this case is 780 hours, and therefore the number of total multipliers is 17940. First, each part is treated as a transfer lot (i.e., transfer lot size = 1). Then the total number

of transfer lots is 144 and the average number of transfer lots in a lot is 5.78. The parts include those that have operations requiring setups and can be done on alternative machines, however, these are less than about 5% of the parts. The testing results are shown in Tab. 5. It can be seen that the algorithm does not require long computational time to get high quality schedules for real problems of this size. To see the impact of transfer lot size on scheduling, every two parts are grouped into a transfer lot (i.e., transfer lot size = 2). The total number of transfer lots is 72, and the average number of transfer lots in a lot is 2.88. The testing results are also given in Tab. 5. It can be seen that the larger the transfer lots size the larger the tardiness (feasible cost), but its solution is closer to the optimum. This implies that smaller transfer lot size usually provides a better solution as expected, but it requires more computations.

Case 4 (A Larger Real Problem with and without Time Step Reduction). A larger problem is tested in this case to further illustrate the capability of the method developed for scheduling large-size practical problems. The planning horizon is 1170 hours for scheduling a total of 61 lots. First, without using time step reduction (i.e., $R = 1$), the number of total multipliers is 29150 and the results are summarized in Tab. 6. To show the effect of “time step reduction” on scheduling large problems, $R = 10$ is used, decreasing the number of total multipliers to 2915. These results are also summarized in Tab. 6. For comparison, the performance of schedules is measured at several iteration counts: 6, 24, 36, and 48. It is obvious that, at a given iteration number, $R=10$ needs much less CPU time than $R=1$ to get a good schedule because of smaller number of multipliers required by $R=10$. It can be also seen that $R=1$ gives a better schedule than $R=10$ in view of the modeling approximation caused by larger

Tab. 4: Description of Data for Cases 3 and 4: Larger Problems

Case No.	No. of Lots	No. of Parts	No. of Part Types	Ave. No. Of Operations/Lot	No. of Machines / Machine Types	Planning Horizon (hour)
3	25	144	17	8.48	39 / 25	780
4	61	339	17	10.8	44 / 27	1170

Tab. 5: Results for Case 3: Different Transfer Lot Sizes

Transfer lot size = 1, no. of transfer lots = 144				Transfer lot size = 2, total no. of transfer lots = 72			
No. of Iterations	Feasible Cost	Duality Gap	CPU Time (second)	No. of Iterations	Feasible Cost	Duality Gap	CPU Time (second)
25	86301	21.9%	302	25	98620	17.1%	305
50	86301	18.2%	593	50	96470	12.2%	602

Tab. 6: Results for Case 4: Time Step Reduction Technique

R = 1 (i.e., without time step reduction)					R = 10 (i.e., with time step reduction)				
No. of Iters.	Average lead-time	Average WIP	Average delay	CPU (s)	No. of Iter.	Average lead-time	Average WIP	Average delay	CPU (s)
6	226	0.28	38	560	6	241	0.28	41	384
24	204	0.27	34.4	2052	24	212	0.27	38.8	1364
36	199	0.266	28.5	3161	36	203	0.272	30.9	2021
48	192	0.255	28	4154	48	203	0.263	29.5	2781

values of R. Measured by practical metrics, these schedules obtained in about 20 minutes for R=10 and in 35 minutes for R=1 look quite reasonable. The duality gap (not included in the table) in this case is significantly larger than previous cases and is still under study. This case implies that for large problems good schedules can be obtained within a reasonable time by using the time step reduction technique.

5 CONCLUSION

The extended lot dynamics model and the backward dynamic programming (BDP) technique have been developed for scheduling with fixed-size transfer lots. The effective handling of transfer lots on machines with setups and machines where all transfer lots in a lot are required to be processed simultaneously is of practical significance. To apply BDP to solve lot sub problems efficiently, the number of multiple completion times associated with each operation beginning time for a lot and the earliest and latest beginning times and completion times for all operations have been determined. With sub states introduced, state transitions determined, and computation-critical DP equations derived for computing stage wise minimum cumulative costs at successor stages, the BDP algorithm is developed to efficiently solve the lot sub problems that contributes to the state-of-the-art scheduling practice. Numerical results indicate that the method can generate high quality schedules with reasonable computational effort.

REFERENCES

- [1] Trietsch D, Baker K. Basic Techniques for Lot Streaming, *Operation Research* (2015) 41(6):1065- 1076.
- [2] Kropp D, Smunt T. Optimal and Heuristic Sciences (2014) 21(1):691-709.
- [3] Wagner B, Ragatz G. The Impact of Lot Splitting on Due Date Performance. *Journal of Operations Management* (2014) 12(1):13-25.
- [4] Benjaafar S. On Production Batches, Transfer Batches, and Lead Times, *Institute of Industrial Engineering Transactions*, (2015) 28(1):357-362.
- [5] Vickson R, Alfredsson B. Two and Three- machine Flow Shop Scheduling Problems with Equal Sized Transfer Batches, *International Journal of Production Research*, (2010) 30(7):1551-1574.
- [6] Cetinkaya F. Lot Streaming in a Two-stage Flow Shop with Set-up, Processing and Removal Times Separated, *Journal of Operation Research Society*, (2013) 45(12):1445-1455.
- [7] Jacobs F, Bragg D. The Repetitive Lots: Flow- time Reduction through Sequencing and Dynamic Batch Sizing, *Decision Sciences*, (2008) 19(1):281- 294.
- [8] Glass C, Gupta J, Potts C. Lot Streaming in Three-stage Production Processes, *European Journal of Operational Research*, (2011) 75(1):378- 394.
- [9] Luh P, Hoitomt D. Scheduling of Manufacturing Systems using the Lagrangian Relaxation Technique. *IEEE Transactions on Automatic Control*, (2012) 38(7):1066-1080.
- [10] Luh P, Gou L, Zhang Y, Nagahora T, Tsuji M, Yoneda K, Hasegawa T, Kyoya Y, Kano T. Job Shop Scheduling with Group-dependent Setups, Finite Buffers, and Long Time Horizon. *Annals of Operations Research*, (2011) 76(1):233-259.
- [11] Zhao X, Luh P, Wang J. The Surrogate Gradient Algorithm for Lagrangian Relaxation Method. *Journal of Optimization Theory and Applications*, (2008) 100(3):699-712.