



HAL
open science

Optimal algorithms for scheduling under time-of-use tariffs

Lin Chen, Nicole Megow, Roman Rischke, Leen Stougie, José Verschae

► **To cite this version:**

Lin Chen, Nicole Megow, Roman Rischke, Leen Stougie, José Verschae. Optimal algorithms for scheduling under time-of-use tariffs. *Annals of Operations Research*, 2021, 304 (1-2), pp.85-107. 10.1007/s10479-021-04059-3 . hal-03474019

HAL Id: hal-03474019

<https://hal.inria.fr/hal-03474019>

Submitted on 10 Dec 2021

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Optimal Algorithms for Scheduling under Time-of-Use Tariffs ^{*}

Lin Chen¹, Nicole Megow², Roman Rischke³, Leen Stougie⁴, and José Verschae⁵

¹ Department of Computer Science, University of Houston, USA. chenlin198662@gmail.com

² Department of Mathematics and Computer Science, University of Bremen, Germany.
nicole.megow@uni-bremen.de

³ Department of Video Coding & Analytics, Fraunhofer Heinrich Hertz Institute Berlin, Germany.
roman.rischke@gmail.com

⁴ CWI & Department of Econometrics and Operations Research, Vrije Universiteit Amsterdam & Erable-INRIA, The Netherlands. stougie@cwi.nl

⁵ Institute of Engineering Sciences, Universidad de O'Higgins, Chile. jose.verschae@uoh.cl

Abstract. We consider a natural generalization of classical scheduling problems in which using a time unit for processing a job causes some time-dependent cost which must be paid in addition to the standard scheduling cost. We study the scheduling objectives of minimizing the makespan and the sum of (weighted) completion times. It is not difficult to derive a polynomial-time algorithm for preemptive scheduling to minimize the makespan on unrelated machines. The problem of minimizing the total (weighted) completion time is considerably harder, even on a single machine. We present a polynomial-time algorithm that computes for any given sequence of jobs an optimal schedule, i.e., the optimal set of time-slots to be used for scheduling jobs according to the given sequence. This result is based on dynamic programming using a subtle analysis of the structure of optimal solutions and a potential function argument. With this algorithm, we solve the unweighted problem optimally in polynomial time. For the more general problem, in which jobs may have individual weights, we develop a polynomial-time approximation scheme (PTAS) based on a dual scheduling approach introduced for scheduling on a machine of varying speed. As the weighted problem is strongly NP-hard, our PTAS is the best possible approximation we can hope for.

1 Introduction

One of the classical operations research problems is the Production Planning problem. It appears in almost any introductory course in Operations Research [11, 19]. In its deterministic form a production plan at lowest total cost is required to meet known demands in the next few weeks, given holding cost for keeping inventory at the end of the week, and with unit production cost varying over the weeks. It is a very early example of a problem model in which unit cost, or tariffs, for production, service, labor, energy, etc., *vary over time*.

Nowadays, new technologies allow direct communication of a much larger variety of time-of-use tariffs to customers. E.g. in energy practice electricity prices can differ largely over the hours. Producers or providers of these resources use these variable pricing more and more to spread demand for their services, which can save enormously on the excessive costs that are usually involved to serve high peak demands. Customers are persuaded to direct their use of the scarce resources to time slots that are offered at cheaper rates. From the provider's point of view variable pricing problems have been studied quite extensively. For instance, revenue management is a well established subfield of operations research [20].

As in the Production Planning problem, in this paper we advocate models from the point of view of the user of the resources, who may take advantage from variable pricing by traveling,

^{*} A preliminary version of this paper with a subset of results appeared in the Proceedings of MFCS 2015 as [4]. This research was supported by the German Science Foundation (DFG) under contract ME 3825/1.

renting labor, using electricity, etc. at moments at which these services are offered at a lower price. This point of view forms a rich class of optimization problems in which next to classical objectives, the cost of using services needs to be taken into account.

This widely applicable framework is particularly well suited for scheduling problems, in which jobs need to be scheduled over time. Processing jobs requires labor, energy, computer power, or other resources that often exhibit variable tariffs over time. It leads to the natural generalization of scheduling problems, in which using a time slot incurs certain cost, varying over time, which we refer to as *utilization cost* that must be paid in addition to the actual scheduling cost. However natural and practicable this may seem, there appears to be very little theoretical research on such scheduling models. The only work we are aware of is by Wan and Qi [21], Kulkarni and Munagala [14], Fang et al. [9] and Chen and Zhang [3], where variable tariffs concern the cost of labor or the cost of energy.

The goal of this paper is to expedite the theoretical understanding of fundamental scheduling problems within the framework of time-varying costs or tariffs. We contribute optimal polynomial-time algorithms and best possible approximation algorithms for the fundamental scheduling objectives of minimizing the sum of weighted completion times and the makespan.

1.1 Problem definition

We first describe the underlying classical scheduling problems. We are given a set of jobs $J := \{1, \dots, n\}$ where every job $j \in J$ has a given processing time $p_j \in \mathbb{N}$ and possibly a weight $w_j \in \mathbb{Q}_{\geq 0}$. The objective is to find a preemptive schedule on a single machine such that the total (weighted) completion time, $\sum_{j \in J} w_j C_j$, is minimized; here C_j denotes the completion time of job j . Preemption means that the processing of a job may be interrupted at any time and can continue at any time later at no additional cost. In the three-field scheduling notation [10], this problem is denoted as $1 | pmtn | \sum w_j C_j$. We also consider makespan minimization on unrelated machines, $R | pmtn | C_{\max}$, where we are given a set of machines M , and each job $j \in J$ has an individual processing time $p_{ij} \in \mathbb{N}$ for running on machine $i \in M$. The objective is to find a preemptive schedule that minimizes the makespan, that is, the completion time of the latest job.

In this paper, we consider a generalization of these scheduling problems within a time-of-use tariff model. We assume that time is discretized into unit-size time slots. We are given a tariff or cost function $e : \mathbb{N} \rightarrow \mathbb{Q}_{\geq 0}$, where $e(t)$ denotes the tariff for processing job(s) at time slot $[t, t + 1)$. We assume that e is piecewise constant. I.e., we assume that the time horizon is partitioned into given intervals $I_k = [s_k, d_k)$ with $s_k, d_k \in \mathbb{N}$, $k = 1, \dots, K$, within which e has the same value e_k . To ensure feasibility, we assume that $d_K \geq \sum_{j \in J} \min_{i \in M} p_{ij}$.

Given a schedule \mathcal{S} , let $y(t)$ be a binary variable indicating if any processing is assigned to time slot $[t, t + 1)$. The *utilization cost* of \mathcal{S} is $E(\mathcal{S}) = \sum_t e(t)y(t)$. That means, for any time unit that is used in \mathcal{S} we pay the full tariff, even if the unit is only partially used. We also emphasize that in the makespan problem in which we have multiple machines, a time slot paid for can be used by all machines. This models applications in which paying for a time unit on a resource gives access to all units of the resource, e.g., all processors on a server.

The overall objective is to find a schedule that minimizes the scheduling objective, $\sum_{j \in J} w_j C_j$ resp. C_{\max} , plus the utilization cost E . We refer to the resulting problems as $1 | pmtn | \sum w_j C_j + E$ and $R | pmtn | C_{\max} + E$. We emphasize that the results in this paper also hold if we minimize any convex combination of the scheduling and utilization cost.

1.2 Related work

Scheduling with time-of-use tariffs (aka variable time slot cost) has been studied explicitly by Wan and Qi [21], Kulkarni and Munagala [14], Fang et al. [9] and Chen and Zhang [3]. In their seminal paper, Wan and Qi [21] consider several *non-preemptive* single machine problems, which are polynomial-time solvable in the classical setting, such as minimizing the total completion time, lateness, and total tardiness, or maximizing the weighted number of on-time jobs. These problems are shown to be strongly NP-hard when taking general tariffs into account, while efficient

algorithms exist for special monotone tariff functions. In particular, the problem $1 \mid \mid \sum C_j + E$ is strongly NP-hard, and it is efficiently solvable when the tariff function is increasing or convex non-increasing [21]. Practical applications, however, often require non-monotone tariff functions, which lead to wide open problems in the context of preemptive and non-preemptive scheduling. In this paper, we answer complexity and approximability questions for fundamental preemptive scheduling problems.

Kulkarni and Munagala [14] focus on a relevant but different problem in an *online setting*, namely online flow-time minimization using resource augmentation. Their main result is a scalable algorithm that obtains a constant performance guarantee when the machine speed is increased by a constant factor and there are only two distinct unit tariffs. They also show that, in this online setting, for an arbitrary number of distinct unit tariffs there is no constant speedup-factor that allows for a constant approximate solution. For the problem considered in this paper, offline scheduling without release dates, Kulkarni and Munagala [14] observed a relation to universal sequencing on a machine of varying speed [8] which implies the following results: a pseudo-polynomial 4-approximation for $1 \mid pmtn \mid \sum w_j C_j + E$, which gives an optimal solution in case that all weights are equal, and a constant approximation in quasi-polynomial time for a constant number of distinct tariffs or when using a machine that is processing jobs faster by a constant factor.

Fang et al. [9] study scheduling on a single machine under time-of-use electricity tariffs. They do not take the scheduling cost into account, but only the energy cost. It also differs from our approach since the schedule is made by the provider and not by the user of the energy. In their model the time horizon is divided into K regions, each of which has a cost c_k per unit energy. For processing jobs the dynamic variable speed model is used; that is, the energy consumption is s^α per unit time if jobs are run at speed s , whence, within region k , the energy cost is $s^\alpha c_k$. The objective is to minimize energy cost such that all jobs are scheduled within the K regions. They prove that the non-preemptive case is NP-hard and give a non-constant approximation, and for the preemptive case, they give a polynomial-time algorithm.

Chen and Zhang [3] consider non-preemptive scheduling on a single machine so as to minimize the total utilization cost under certain scheduling feasibility constraints such as a common deadline for all jobs or a bound on the maximum lateness, maximum tardiness, maximum flow-time, or sum of completion times. They define a *valley* to be a cost interval I_k that has smaller cost than its neighboring intervals and show the following. General tariffs lead to a strongly NP-hard problem for any of the just mentioned constraints, and even very restricted tariff functions with more than one valley result in NP-hard problems that are not approximable within any constant factor. The problem with a common deadline on the job completion times is shown to admit a pseudo-polynomial time algorithm when having two valleys, a polynomial time algorithm for tariff functions with at most one valley, and an FPTAS if there are at most two valleys and $\max_k e_k / \min_k e_k$ is bounded. For the other mentioned constraints, they also present polynomial time algorithms when having no more than one valley, where the problem with a bound on the sum of completion times requires the number of cost intervals, here K , to be fixed.

The general concept of taking into consideration additional (time-dependent) cost for resource utilization when scheduling has been implemented differently in other models. We mention the area of energy-aware scheduling, where energy consumption is taken into account (see [1] for an overview). Further, the area of scheduling with generalized non-decreasing (completion-) time dependent cost functions, such as minimizing $\sum_j w_j f(C_j)$, e.g. [8, 12, 17], or even more general job-individual cost functions $\sum_j f_j(C_j)$, e.g. [2, 5, 6, 13] has received quite some attention. Our model differs fundamentally from those models since our cost function may decrease with time. In fact, delaying the processing in favor of cheaper time slots may decrease the overall cost. This is not the case in the above-mentioned models. Thus, in our framework we have the additional dimension in decision-making of selecting the time slots that shall be utilized.

Finally, we point out some similarity between our model and *scheduling on a machine of varying speed*, which (with $\sum_j w_j C_j$ as objective function) is an equivalent statement of the problem of minimizing $\sum_j w_j f(C_j)$ on a single machine with constant speed [8, 12, 17]. We do not see any mathematical reduction from one problem to the other. However, it is noteworthy that the independently studied problem of scheduling with *non-availability periods*, see e.g. the survey by

Lee [16], is a special case of both the varying-speed and the time-varying tariff model. Indeed, machine non/availability can be expressed either by 0/1-speed or equivalently by $\infty/0$ tariff. Results shown in this context imply that our problem $1 | pmtn | \sum w_j C_j + E$ is strongly NP-hard, even if there are only two distinct tariffs [22].

1.3 Our contribution

We present new optimal algorithms and best-possible approximation results, unless $P=NP$, for the generalization of basic scheduling problems to a framework with time-varying tariffs.

One of our results is a rather straightforward optimal polynomial-time algorithm for the problem $R | pmtn | C_{\max} + E$ (Section 4): We design a procedure that selects the optimal time slots to be utilized, given that we know their optimal *number*. That number can be determined by solving the scheduling problem *without* utilization cost, which can be done in polynomial time by solving a linear program [15].

Whereas minimizing makespan plus utilization cost appears to be efficiently solvable even in the most general machine model, the objective of minimizing the total weighted completion time raises significant complications. Our results on this objective concern single-machine problems (Section 2). We present an algorithm that computes for a given ordered set of jobs an optimal choice of time slots to be used. We derive this by first showing structural properties of an optimal schedule, which we then exploit together with a properly chosen potential function in a dynamic program yielding polynomial running time. Based on this algorithm, we show that the unweighted problem $1 | pmtn | \sum C_j + E$ can be solved in polynomial time and that it allows almost directly for a fully polynomial $(4+\varepsilon)$ -approximation algorithm for the weighted version $1 | pmtn | \sum w_j C_j + E$, for which a pseudo-polynomial 4-approximation was observed by Kulkarni and Munagala [14]. While pseudo-polynomial time algorithms are relatively easy to derive, it is quite remarkable that our DP's running time is polynomial in the input, in particular, independent of d_K .

In Section 3, we significantly improve the approximation result for the weighted problem by designing a polynomial-time algorithm that computes for any fixed ε a $(1+\varepsilon)$ -approximate schedule for $1 | pmtn | \sum w_j C_j + E$, that is, we give a polynomial-time approximation scheme (PTAS). Unless $P=NP$, our algorithm is best possible, since the problem is strongly NP-hard even if there are only two different tariffs [22].

Our approach is inspired by a recent PTAS for scheduling on a machine of varying speed [17] and it uses some of its properties. As discussed before, we do not see a formal mathematical relation between these two seemingly related problems which allows to apply the result from [17] directly. The key is a dual view on scheduling: instead of directly constructing a schedule in the time-dimension, we first construct a dual scheduling solution in the weight-dimension which has a one-to-one correspondence to a true schedule. We design an exponential-time dynamic programming algorithm which can be trimmed to polynomial time using techniques known for scheduling with varying speed [17].

For both the makespan and the min-sum problem, job preemption is crucial for obtaining constant worst-case performance ratios. For non-preemptive scheduling, a straightforward reduction from 2-PARTITION shows that no approximation is possible, unless $P=NP$, even if there are only two different tariffs, 0 and ∞ .

Finally, we remark that in general it is not clear that a schedule can be encoded polynomially in the input. However, for our completion-time based minimization objectives, it is easy to observe that if an algorithm utilizes p unit-size time slots in an interval of equal cost, then it utilizes the first p slots within this interval, which simplifies the structure and the output of an optimal solution in a crucial way.

We start below with presenting the more involved results for the problems with scheduling objective minimizing total (weighted) completion time. The efficient algorithm for the makespan objective is then presented in Section 4.

2 An optimal algorithm for minimizing total completion time

In this section, we show how to solve the unweighted problem $1|pmtn|\sum C_j + E$ to optimality. Our main result is as follows.

Theorem 1. *There is a polynomial-time algorithm for $1|pmtn|\sum C_j + E$.*

An algorithm for the scheduling problem with time-of-use tariffs has to make essentially two types of decisions: (i) which time slots to use and (ii) how to schedule the jobs in these slots. It is not hard to see that these two decisions can be handled separately. In fact, the following observation on the optimal sequencing of jobs holds independently of the utilization decision and follows from a standard interchange argument.

Observation 1 *In an optimal schedule \mathcal{S}^* for the problem $1|pmtn|\sum C_j + E$, jobs are processed according to the Shortest Processing Time First (SPT) rule.*

Thus, in the remainder of the section we can focus on determining which time slots to use. We design an algorithm that computes, for any given (not necessarily optimal) scheduling sequence σ , an optimal utilization decision for σ . In fact, we show our structural result even for the more general problem in which jobs have arbitrary weights.

Theorem 2. *Given an instance of $1|pmtn|\sum w_j C_j + E$ and an arbitrary processing sequence of jobs σ , we can compute an optimal utilization decision for σ in polynomial time.*

Combining the optimal choice of time slots (Theorem 2) with the optimal processing order SPT (Observation 1) immediately implies Theorem 1.

The remainder of the section is devoted to proving Theorem 2. Thus, we choose any (not necessarily optimal) order of jobs, $\sigma = (1, \dots, n)$, in which the jobs must be processed. We want to characterize an optimal schedule \mathcal{S}^* for σ , that is, the optimal choice of time slots for scheduling σ . We firstly identify structural properties of an optimal solution. Essentially, we give a full characterization which we can compute efficiently by dynamic programming.

More precisely, we establish a closed form that characterizes the relationship between the tariff of an utilized slot and job weights in an optimal solution. This relationship allows to decompose an optimal schedule into a series of sub-schedules. Our algorithm will first compute all possible sub-schedules and then use a dynamic programming approach to select and concatenate suitable sub-schedules.

In principle, an optimal schedule may preempt jobs at fractional time points. However, since time slots can only be paid for entirely, any reasonable schedule uses the utilized slots entirely as long as there are unprocessed jobs. It can be shown by a standard interchange argument that this is also true if we omit the requirement that time slots must be utilized entirely; for details, see [18]. (We remark that for the makespan problem with multiple machines considered in Section 4 this is not true.)

Lemma 1. *Allowing to pay for utilizing partial time slots, there is an optimal schedule \mathcal{S}^* for $1|pmtn|\sum w_j C_j + E$ in which all utilized time slots are entirely utilized and jobs are preempted only at integral points in time.*

Next, we split the optimal schedule \mathcal{S}^* for the given job sequence $\sigma = (1, \dots, n)$ into smaller sub-schedules. To that end we introduce the concept of a *split point*.

Definition 1 (Split Point). *Consider an optimal schedule \mathcal{S}^* and the set of potential split points $\mathcal{P} := \bigcup_{k=1}^K \{s_k, s_k + 1\} \cup \{d_K\}$. Let S_j and C_j denote the start time and completion time of job j , respectively. We call a time point $t \in \mathcal{P}$ a split point for \mathcal{S}^* if all jobs that start before t also finish their processing not later than t , i.e., if $\{j \in J : S_j < t\} = \{j \in J : C_j \leq t\}$.*

Given an optimal schedule \mathcal{S}^* , let $0 = \tau_1 < \tau_2 < \dots < \tau_\ell = d_K$ be the *maximal* sequence of split points of \mathcal{S}^* , i.e. the sequence containing all split points of \mathcal{S}^* . We denote the interval between two consecutive split points τ_x and τ_{x+1} as *region* $R_x^{\mathcal{S}^*} := [\tau_x, \tau_{x+1})$, for $x = 1, \dots, \ell - 1$.

Consider now any region $R_x^{\mathcal{S}^*}$ for an optimal schedule \mathcal{S}^* with $x \in \{1, \dots, \ell - 1\}$ and let $J_x^{\mathcal{S}^*} := \{j \in J : S_j \in R_x^{\mathcal{S}^*}\}$, the jobs that start and finish within $R_x^{\mathcal{S}^*}$. Note that $J_x^{\mathcal{S}^*}$ might be empty. Among all optimal schedules we shall consider an optimal solution \mathcal{S}^* that minimizes the value $\sum_{t=0}^{d_K-1} t \cdot y(t)$, where $y(t)$ is a binary variable that indicates if time slot $[t, t+1)$ is utilized or not.

We observe that any job j completing at the beginning of a cost interval I_k , i.e. $C_j = s_k \in R_x^{\mathcal{S}^*}$ or $C_j = s_k + 1 \in R_x^{\mathcal{S}^*}$, would make s_k resp. $s_k + 1$ a split point. Thus, no such job can exist.

Observation 2 *There is no job $j \in J_x^{\mathcal{S}^*}$ with $C_j \in R_x^{\mathcal{S}^*} \cap \mathcal{P}$.*

We say that interval I_k is *partially utilized* if at least one time slot in I_k is utilized, but not all.

Lemma 2. *There exists an optimal schedule \mathcal{S}^* in which for all $x = 1, \dots, \ell - 1$ at most one interval is partially utilized in $R_x^{\mathcal{S}^*}$.*

Proof. By contradiction, suppose that there is more than one partially utilized interval in $R_x^{\mathcal{S}^*}$. Consider any two such intervals I_h and $I_{h'}$ with $h < h'$, and all intermediate intervals utilized entirely or not at all. Let $[t_h, t_h + 1)$ and $[t_{h'}, t_{h'} + 1)$ be the last utilized time slot in I_h and $I_{h'}$, respectively. If we utilize $[t_{h'} + 1, t_{h'} + 2)$ instead of $[t_h, t_h + 1)$, then the difference in cost is $\delta_1 := e_{h'} - e_h + \sum_{j \in J'} w_j$ with $J' := \{j \in J : C_j \in \bigcup_{k=h+1}^{h'} I_k\}$ because all jobs in J' are delayed by exactly one time unit. This is true since by Observation 2 no job finishes at $d_k = s_{k+1}$ for any k . If we utilize $[t_h + 1, t_h + 2)$ instead of $[t_{h'}, t_{h'} + 1)$, then the difference in cost is $\delta_2 := e_h - e_{h'} - \sum_{j \in J'} w_j$, again using Observation 2 to assert that no job finishes at $s_k + 1$ for any $h + 1 \leq k \leq h'$. Since $\delta_1 = -\delta_2$ and \mathcal{S}^* is an optimal schedule, it must hold that $\delta_1 = \delta_2 = 0$. This, however, implies that there is another optimal schedule with earlier used time slots which contradicts our assumption that \mathcal{S}^* minimizes the value $\sum_{t=0}^{d_K-1} t \cdot y(t)$. \square

The next Lemma characterizes the time slots that are used within a region. Let e_{\max}^j be the maximum tariff spent for job j in \mathcal{S}^* . Furthermore, let $\Delta_x := \max_{j \in J_x^{\mathcal{S}^*}} (e_{\max}^j + \sum_{j' < j} w_{j'})$ and let j_x be the last job (according to sequence σ) that achieves Δ_x . Suppose, there are $b \geq 0$ jobs before and $a \geq 0$ jobs after job j_x in $J_x^{\mathcal{S}^*}$. The following lemma gives for every job $j \in J_x^{\mathcal{S}^*} \setminus \{j_x\}$ an upper bound on the tariff spent in the interval $[S_j, C_j)$.

Lemma 3. *Consider an optimal schedule \mathcal{S}^* for a given job permutation σ . For any job $j \in J_x^{\mathcal{S}^*} \setminus \{j_x\}$ a slot $[t, t+1) \in [S_j, C_j)$ is utilized if and only if the tariff $e(t)$ of $[t, t+1)$ satisfies the following upper bound:*

$$e(t) \leq \begin{cases} e_{\max}^{j_x} + \sum_{j'=j}^{j_x-1} w_{j'}, & \forall j : j_x - b \leq j < j_x \\ e_{\max}^{j_x} - \sum_{j'=j_x}^{j-1} w_{j'}, & \forall j : j_x < j \leq j_x + a. \end{cases}$$

Proof. Consider any job $j := j_x - \ell$ with $0 < \ell \leq b$. Suppose there is a job j for which a slot is utilized with cost (tariff) $e_{\max}^j > e_{\max}^{j_x} + \sum_{j'=j}^{j_x-1} w_{j'}$. Then $e_{\max}^j + \sum_{j' < j} w_{j'} > e_{\max}^{j_x} + \sum_{j' < j_x} w_{j'}$, which is a contradiction to the definition of job j_x . Thus, $e_{\max}^j \leq e_{\max}^{j_x} + \sum_{j'=j}^{j_x-1} w_{j'}$.

Now suppose that there is a slot $[t, t+1) \in [S_j, C_j)$ with cost $e(t) \leq e_{\max}^{j_x} + \sum_{j'=j}^{j_x-1} w_{j'}$ that is not utilized. There must be a slot $[t', t'+1) \in [S_{j_x}, C_{j_x})$ with cost exactly $e_{\max}^{j_x}$. If we utilize slot $[t, t+1)$ instead of $[t', t'+1)$, then the difference in cost is non-positive, because the completion times of at least ℓ jobs ($j = j_x - \ell, \dots, j_x - 1$ and maybe also j_x) decrease by one. This contradicts either the optimality of \mathcal{S}^* or our assumption that \mathcal{S}^* minimizes $\sum_{t=0}^{d_K-1} t \cdot y(t)$.

The proof of the statement for any job $j_x + \ell$ with $0 < \ell \leq a$ follows a similar argument, but now using the fact that for every job $j := j_x + \ell$ we have $e_{\max}^j < e_{\max}^{j_x} - \sum_{j'=j_x}^{j-1} w_{j'}$, because j_x was the last job with $e_{\max}^j + \sum_{j' < j} w_{j'} = \Delta_x$. \square

Corollary 1. *If the interval $[S_j, C_j]$ for processing a job $j \in J_x^{S^*} \setminus \{j_x\}$ intersects interval I_k but job j does not complete in I_k , i.e., $C_j > d_k$, then all time slots in I_k are fully utilized.*

To decide on an optimal utilization decision for the sub-schedule of the jobs in $R_x^{S^*}$, we need the following two lemmas.

Lemma 4. *If there is a partially utilized interval I_k in region $R_x^{S^*}$, then (i) I_k is the last interval of $R_x^{S^*}$, or (ii) j_x is the last job being processed in I_k and $e_k = e_{\max}^{j_x}$.*

Proof. Suppose there exists a partially utilized interval I_k in region $R_x^{S^*}$. Suppose j with $j \neq j_x$ is the last job that is processed in I_k , hence (ii) does not hold. Then either $C_j < d_k$, in which case $d_k = s_{k+1}$ is a split point and thus I_k is the last interval in the region, whence (i) is true. Or, we are in the situation of Corollary 1 and have a contradiction, because then I_k must be fully utilized.

Now suppose j_x is the last job being processed in I_k . If $C_{j_x} < d_k$, then again I_k is the last interval in the region. Otherwise $C_{j_x} \notin I_k$. If $e_k = e_{\max}^{j_x}$, then case (ii) of the lemma holds. If not, by definition of $e_{\max}^{j_x}$ we have $e_k < e_{\max}^{j_x}$. By optimality of S^* , interval I_k comes after the last utilized “expensive” interval with cost $e_{\max}^{j_x}$. Hence, job j_x is processed in an expensive interval, then in I_k and is completed in yet another interval. But then we can utilize an extra time slot in I_k instead of a time slot in the expensive interval, without increasing the completion time. This contradicts optimality, and, hence, $e_k = e_{\max}^{j_x}$, which completes the proof. \square

Lemma 5. *There exists an optimal schedule S^* for a given job permutation σ with the following property. If the last interval I_k of a region $R_x^{S^*}$ is only partially utilized then all time slots in $[S_{j_x}, C_{j_x})$ with cost at most $e_{\max}^{j_x}$ are utilized.*

Proof. Recall that $j_x + a$ is the last job being processed in the region, and hence, it is the last job processed in the partially utilized interval I_k .

Suppose there is a time slot $[t, t+1) \in [S_{j_x}, C_{j_x})$ with cost at most $e_{\max}^{j_x}$ that is not utilized. If we utilize $[t, t+1)$ instead of the last utilized slot in I_k , then the difference in cost is $\delta_1 := e(t) - e_k - \sum_{j=j_x}^{j_x+a} w_j$. On the other hand, if we utilize one additional time slot in I_k instead of a time slot in $[S_{j_x}, C_{j_x})$ with cost $e_{\max}^{j_x}$, then the difference in cost is $\delta_2 := e_k - e_{\max}^{j_x} + \sum_{j=j_x}^{j_x+a} w_j$. We consider an optimal schedule S^* , thus $\delta_1 \geq 0$ and $\delta_2 \geq 0$ which implies that $\delta_1 + \delta_2 = e(t) - e_{\max}^{j_x} \geq 0$. This is a contradiction if $e(t) < e_{\max}^{j_x}$. If $e(t) = e_{\max}^{j_x}$, then $\delta_1 = -\delta_2 = 0$, because we consider an optimal schedule S^* . This, however, contradicts our assumption that S^* minimizes the value $\sum_{t=0}^{d_K-1} t \cdot y(t)$. \square

We now show how to construct an optimal partial schedule for a given ordered job set in a given region in polynomial time.

Lemma 6. *Given a region R_x and an ordered job set J_x , we can find in polynomial time an optimal utilization decision for scheduling J_x within the region R_x , which does not contain any other split points than τ_x and τ_{x+1} , the boundaries of R_x .*

Proof. Given R_x and J_x , we guess the optimal combination $(j_x, e_{\max}^{j_x})$, i.e., we enumerate over all nK combinations and choose eventually the best solution.

We firstly assume that a partially utilized interval exists and it is the last one in R_x (case (i) in Lemma 4). Based on the characterization in Lemma 3 we find in polynomial time the slots to be utilized for the jobs $j_x - b, \dots, j_x - 1$. This defines $C_{j_x-b}, \dots, C_{j_x-1}$. Then starting job j_x at time C_{j_x-1} , we check intervals in the order given and utilize as much as needed of each next interval I_h if and only if $e_h \leq e_{\max}^{j_x}$, until a total of p_{j_x} time slots have been utilized for processing j_x . Lemma 5 justifies to do that. This yields a completion time C_{j_x} . Starting at C_{j_x} , we use again Lemma 3 to find in polynomial time the slots to be utilized for processing the jobs $j_x + 1, \dots, j_x + a$. This gives $C_{j_x+1}, \dots, C_{j_x+a}$.

Now we assume that there is no partially utilized interval or we are in case (ii) of Lemma 4. Similar to the case above, we find in polynomial time the slots that S^* utilizes for the jobs $j_x - b, \dots, j_x - 1$ based on Lemma 3. This defines $C_{j_x-b}, \dots, C_{j_x-1}$. To find the slots to be utilized

for the jobs $j_x + 1, \dots, j_x + a$, in this case, we start at the end of R_x and go backwards in time. We can start at the end of R_x because in this case the last interval of R_x is fully utilized. This gives $C_{j_x+1}, \dots, C_{j_x+a}$. Job j_x is thus to be scheduled in $[C_{j_x-1}, S_{j_x+1})$. In order to find the right slots for j_x we solve a makespan problem in the interval $[C_{j_x-1}, S_{j_x+1})$, which can be done in polynomial time (Theorem 4) and gives a solution that cannot be worse than what an optimal schedule \mathcal{S}^* does.

If anywhere in both cases the utilized intervals can not be made sufficient for processing the job(s) for which they are intended, or if scheduling the jobs in the utilized intervals creates any intermediate split point, then this $(j_x, e_{\max}^{j_x})$ -combination is rejected. Hence, we have computed the optimal schedules over all nK combinations of $(j_x, e_{\max}^{j_x})$ and over both cases of Lemma 4 concerning the position of the partially utilized interval. We choose the schedule with minimum total cost and return it with its value. This completes the proof. \square

Now we are ready to prove our main theorem.

Proof (Proof of Theorem 2). We give a dynamic program. Assume jobs are indexed according to the order given by σ . We define a state (j, t) , where t is a potential split point $t \in \mathcal{P}$ and j is a job from the job set J , and a dummy job 0. The value of a state, $Z(j, t)$, is the optimal scheduling cost plus utilization cost for completing jobs $1, \dots, j$ by time t . We apply the following recursion:

$$\begin{aligned} Z(j, t) &= \min \left\{ Z(j', t') + z(\{j'+1, \dots, j\}, [t', t]) \mid t', t \in \mathcal{P}, t' < t, j', j \in J, j' \leq j \right\}, \\ Z(0, t) &= 0, \quad \text{for any } t, \\ Z(j, s_1) &= \infty, \quad \text{for any } j > 0, \end{aligned}$$

where $z(\{j'+1, \dots, j\}, [t', t])$ denotes the value of an optimal partial schedule for job set $\{j'+1, j'+2, \dots, j\}$ in the region $[t', t)$, or ∞ if no such schedule exists. In case $j = j'$ there is no job to be scheduled in the interval $[t', t)$, whence we set $z(\{j'+1, \dots, j\}, [t', t]) = 0$. This models the option of leaving regions empty.

An optimal partial schedule can be computed in polynomial time as we have shown in Lemma 6. Hence, we compute $Z(j, t)$ for all $O(nK)$ states in polynomial time, which concludes the proof. \square

Remark: A simple $(4 + \epsilon)$ -approximation for the weighted problem. It is worth mentioning that the characterization of an optimal utilization decision above (Theorem 2) can be used to obtain a simple $(4 + \epsilon)$ -approximation for the *weighted* problem $1 \mid pmtn \mid \sum w_j C_j + E$.

For the weighted problem, there may not exist a job sequence that is universally optimal for *all* utilization decisions [8]. However, in the context of scheduling on an unreliable machine there has been shown a polynomial-time algorithm that computes a universal $(4 + \epsilon)$ -approximation [8]. More precisely, the algorithm constructs a sequence of jobs which approximates the scheduling cost for any utilization decision within a factor at most $4 + \epsilon$.

Consider an instance of problem $1 \mid pmtn \mid \sum w_j C_j + E$ and compute such a universally $(4 + \epsilon)$ -approximate sequence σ . Applying Theorem 2 to σ , we obtain a schedule \mathcal{S} with an optimal utilization decision for σ . Let \mathcal{S}' denote the schedule which we obtain by changing the utilization decision of \mathcal{S} to the utilization in an optimal schedule \mathcal{S}^* (but keeping the scheduling sequence σ). The schedule \mathcal{S}' has cost no less than the original cost of \mathcal{S} . Furthermore, given the utilization decision in the optimal solution \mathcal{S}^* , the sequence σ approximates the scheduling cost of \mathcal{S}^* within a factor of $4 + \epsilon$. This gives the following result.

Corollary 2. *There is a $(4 + \epsilon)$ -approximation algorithm for $1 \mid pmtn \mid \sum w_j C_j + E$.*

This result is superseded by the PTAS presented in the next section.

3 A PTAS for minimizing the total weighted completion time

The main result of this section is a polynomial time approximation scheme for minimizing the total weighted completion time with time-varying utilization cost.

Theorem 3. *There is a polynomial-time approximation scheme for $1 | pmtn | \sum w_j C_j + E$.*

In the remainder of this section we describe some preliminaries, present a dynamic programming (DP) algorithm with exponential running time, and then we argue that the running time can be reduced to polynomial time. As noted in the introduction, our approach is inspired by a PTAS for scheduling on a machine of varying speed [17], but a direct application does not seem possible.

3.1 Preliminaries and scheduling in the weight-dimension

We describe a schedule \mathcal{S} not in terms of completion times $C_j(\mathcal{S})$, but in terms of the remaining weight function $W^{\mathcal{S}}(t)$ which, for a given schedule \mathcal{S} , is defined as the total weight of all jobs not completed by time t . Notice that, by definition, $W^{\mathcal{S}}(t)$ is right-continuous. Based on the remaining weight function we can express the cost for any schedule \mathcal{S} as

$$\int_0^\infty W^{\mathcal{S}}(t) dt = \sum_{j \in J} w_j C_j(\mathcal{S}).$$

This has a natural interpretation in the standard 2D-Gantt chart, which was originally introduced in [7].

For a given utilization decision, we follow the idea of [17] and implicitly describe the completion time of a job j by the value of the function $W^{\mathcal{S}}$ at the time that j completes. This value is referred to as the *starting weight* S_j^w of job j . In analogy to the time-dimension, the value $C_j^w := S_j^w + w_j$ is called *completion weight* of job j . When we specify a schedule in terms of the remaining weight function, then we call it a *weight-schedule*, otherwise a *time-schedule*. Other terminologies, such as feasibility and idle time, also translate from the time-dimension to the weight-dimension. A weight-schedule is called *feasible* if no two jobs overlap and the machine is called *idle in weight-dimension* if there exists a point w in the weight-dimension with $w \notin [S_j^w, C_j^w]$ for all jobs $j \in J$.

A weight-schedule together with a utilization decision can be translated into a time-schedule by ordering the job in decreasing order of completion weights and scheduling them in this order in the time-dimension in the utilized time slots. For a given utilization decision, consider a weight-schedule \mathcal{S} with completion weights $C_1^w > \dots > C_n^w > C_{n+1}^w := 0$ and the corresponding completion times $0 =: C_0 < C_1 < \dots < C_n$ for the jobs $j = 1, \dots, n$. We define the (*scheduling*) *cost of a weight-schedule* \mathcal{S} as $\sum_{j=1}^n (C_j^w - C_{j+1}^w) C_j$. This value equals $\sum_{j=1}^n \pi_j^{\mathcal{S}} C_j^w$, where $\pi_j^{\mathcal{S}} := C_j - S_j$, if and only if there is no idle weight. If there is idle weight, then the cost of a weight-schedule can only be greater, and we can safely remove idle weight without increasing the scheduling cost [17]. Figure 1 illustrates this fact.

Summarizing, a time-schedule implies a correspondent weight-schedule of the same cost. On the other hand, a weight-schedule plus a utilization decision implies a time-schedule with a possibly smaller cost.

3.2 Dynamic programming algorithm

Let $\varepsilon > 0$. Firstly, we scale the input parameters so that all job weights w_j , $j = 1, \dots, n$, and all tariffs e_k , $k = 1, \dots, K$, are non-negative integers. Then, we apply standard geometric rounding to the weights to gain more structure on the input, i.e, we round the weights of all jobs up to the next integer power of $(1 + \varepsilon)$, by losing at most a factor $(1 + \varepsilon)$ in the objective value. Furthermore, we discretize the weight-space into intervals of exponentially increasing size: we define intervals $WI_u := [(1 + \varepsilon)^{u-1}, (1 + \varepsilon)^u)$ for $u = 1, \dots, \nu$ with $\nu := \lceil \log_{1+\varepsilon} \sum_{j \in J} w_j \rceil$.

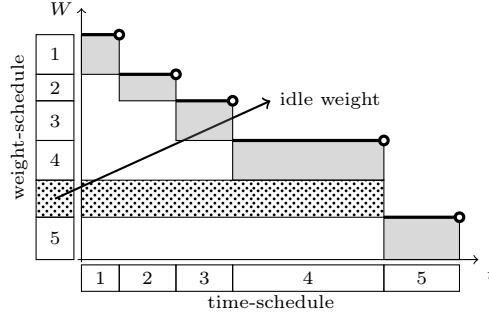


Fig. 1. 2D-Gantt chart. The x-axis shows a time-schedule, while the y-axis corresponds to $W(t) = \sum_{C_j > t} w_j$ plus the idle weight in the corresponding weight-schedule [17].

Consider a subset of jobs $J' \subseteq J$ and a partial weight-schedule of J' . In the dynamic program, the set J' represents the set of jobs at the beginning of a corresponding weight-schedule, i.e., if $j \in J'$ and $k \in J \setminus J'$, then $C_j^w < C_k^w$. However, the jobs in J' are scheduled at the end in a corresponding time-schedule. As discussed in Section 3.1, a partial weight-schedule \mathcal{S} for the jobs in $J \setminus J'$ together with a utilization decision for these jobs can be translated into a time-schedule.

Let $\mathcal{F}_u := \{J_u \subseteq J : \sum_{j \in J_u} w_j \leq (1 + \varepsilon)^u\}$ for $u = 1, \dots, \nu$. The set \mathcal{F}_u contains all the possible job sets J_u that can be scheduled in WI_u or before. Additionally, we define \mathcal{F}_0 to be the set that contains only the set of all zero-weight jobs $J_0 := \{j \in J : w_j = 0\}$. The following observation allows us to restrict to simplified completion weights.

Observation 3 Consider an optimal weight-schedule in which the set of jobs with completion weight in WI_u , $u \in \{1, \dots, \nu\}$, is exactly $J_u \setminus J_{u-1}$ for some $J_u \in \mathcal{F}_u$ and $J_{u-1} \in \mathcal{F}_{u-1}$. By losing at most a factor $(1 + \varepsilon)$ in the objective value, we can assume that for all $u \in \{1, \dots, \nu\}$ the completion weight of the jobs in $J_u \setminus J_{u-1}$ is exactly $(1 + \varepsilon)^u$.

The following observation follows from a simple interchange argument.

Observation 4 There is an optimal time-schedule in which J_0 is scheduled completely after all jobs in $J \setminus J_0$.

The dynamic program recursively constructs states $Z = [J_u, b, avg]$ and computes for every state a time point $t(Z)$ with the following meaning. A state $Z = [J_u, b, avg]$ with time point $t(Z)$ expresses that there is a feasible partial time-schedule \mathcal{S} for the jobs in $J \setminus J_u$ with $J_u \in \mathcal{F}_u$ together with a utilization decision for the time interval $[0, t(Z))$ with total utilization cost at most b and for which the *average scheduling cost*, i.e.,

$$\frac{1}{t(Z)} \cdot \int_0^{t(Z)} W^{\mathcal{S}}(t) dt,$$

is at most avg . We remark that even if \mathcal{S} only schedules jobs in $J \setminus J_u$, the remaining weight function $W^{\mathcal{S}}$ still considers jobs in $J \setminus J_u$, and thus $W^{\mathcal{S}}(t(Z)) = \sum_{j \in J \setminus J_u} w_j$. Also, \mathcal{S} implies a weight-schedule for jobs in $J \setminus J_u$ where the completion weights belong to $[\sum_{j \in J_u} w_j, \sum_{j \in J} w_j]$. Note that $avg \cdot t(Z)$ is an upper bound on the total scheduling cost of \mathcal{S} and that the average scheduling cost is non-increasing in time, because the remaining weight function $W^{\mathcal{S}}(t)$ is non-increasing in time. In the *iteration for u* , we only consider states $[J_u, b, avg]$ with $J_u \in \mathcal{F}_u$. The states in the iteration for u are created based on the states from the iteration for $u + 1$. Initially, we only have the state $Z_\nu = [J, 0, 0]$ with $t(Z_\nu) := 0$, we start the dynamic program with $u = \nu - 1$, iteratively reduce u by one, and stop the process after the iteration for $u = 0$. In the iteration for u , the states together with their time points are constructed in the following way. Consider candidate sets $J_{u+1} \in \mathcal{F}_{u+1}$

and $J_u \in \mathcal{F}_u$ with $\mathcal{F}_u \subseteq \mathcal{F}_{u+1}$, a partial time-schedule \mathcal{S} of $J \setminus J_u$, in which the set of jobs with completion weight (in the correspondent weight-schedule) in WI_{u+1} is exactly $J_{u+1} \setminus J_u$ and the set of jobs later than WI_{u+1} is exactly $J \setminus J_{u+1}$, two budgets b_1, b_2 with $b_1 \leq b_2$, and two bounds on the average scheduling cost avg_1, avg_2 . Let $Z_1 = [J_{u+1}, b_1, avg_1]$ and $Z_2 = [J_u, b_2, avg_2]$ be the corresponding states. We know that there is a feasible partial schedule for the job set $J \setminus J_{u+1}$ up to time $t(Z_1)$ having average scheduling cost at most avg_1 and utilization cost at most b_1 . By augmenting this schedule, we want to compute a minimum time point $t(Z_1, Z_2)$ that we associate with the link between Z_1 and Z_2 so that there is a feasible partial schedule for $J \setminus J_u$ that processes the jobs from $J_{u+1} \setminus J_u$ in the interval $[t(Z_1), t(Z_1, Z_2))$, has average scheduling cost at most avg_2 , and utilization cost at most b_2 . That is, $t(Z_1, Z_2)$ is the minimum makespan if we start with Z_1 and want to arrive at Z_2 . For the computation of $t(Z_1, Z_2)$, we use the following subroutine.

Using Observation 3, we approximate the area under the remaining weight function $W^S(t)$ for the jobs in $J_{u+1} \setminus J_u$ by $(1 + \varepsilon)^{u+1} \cdot (t(Z_1, Z_2) - t(Z_1))$, where $t(Z_1, Z_2)$ is the time point that we want to compute. Approximating this area gives us the flexibility to schedule the jobs in $J_{u+1} \setminus J_u$ in any order. However, we need that $avg_2 \cdot t(Z_1, Z_2)$ is an upper bound on the integral of the remaining weight function by time $t(Z_1, Z_2)$. That is, we want that

$$avg_2 \cdot t(Z_1, Z_2) \geq (1 + \varepsilon)^{u+1} \cdot t(Z_1, Z_2) + t(Z_1) \cdot (avg_1 - (1 + \varepsilon)^{u+1}).$$

Both the left-hand side and the right-hand side of this inequality are linear functions in $t(Z_1, Z_2)$. So, we can compute a smallest time point t^{LB} such that the right-hand side is greater or equal to the left-hand side for all $t(Z_1, Z_2) \geq t^{LB}$. If there is no such t^{LB} , then we set $t(Z_1, Z_2)$ to infinity and stop the subroutine. Otherwise, we know that our average scheduling cost at t^{LB} or later is at most avg_2 . Let $E(p, [t_1, t_2])$ denote the total cost of the p cheapest slots in the time-interval $[t_1, t_2]$. We compute the smallest time point $t(Z_1, Z_2) \geq t^{LB}$ so that the set of jobs $J_{u+1} \setminus J_u$ can be feasibly scheduled in $[t(Z_1), t(Z_1, Z_2))$ having utilization cost not more than $b_2 - b_1$. That is, we set

$$t(Z_1, Z_2) = \min \{t \geq \max\{t(Z_1), t^{LB}\} : E(p(J_{u+1} \setminus J_u), [t(Z_1), t]) \leq b_2 - b_1\}.$$

The time point $t(Z_1, Z_2)$ can be computed in polynomial time by applying binary search to the interval $[\max\{t(Z_1), t^{LB}\}, d_K]$, since $E(p, [t_1, t_2])$ is a monotone function in t_2 .

Given all possible states $[J_{u+1}, b_1, avg_1]$ from the iteration for $u + 1$, the dynamic program enumerates for all these states all possible links to states $[J_u, b_2, avg_2]$ from the iteration for u fulfilling the above requirement on the candidate sets J_{u+1} and J_u , on the budgets b_1 and b_2 , and on the average scheduling costs avg_1 and avg_2 . For any such possible link (Z_1, Z_2) between states from the iteration for $u + 1$ and u , we apply the above subroutine and associate the time point $t(Z_1, Z_2)$ with this link. Thus, the dynamic program associates several possible time points with a state $Z_2 = [J_u, b_2, avg_2]$ from the iteration for u . However, we only keep the link with the smallest associated time point $t(Z_1, Z_2)$ (ties are broken arbitrarily) and this defines the time point $t(Z_2)$ that we associate with the state Z_2 . That is, for a state Z_2 from the iteration for u we define $t(Z_2) := \min\{t(Z_1, Z_2) \mid Z_1 \text{ is a state from the iteration for } u + 1\}$.

Let E_{\max} be an upper bound on the total utilization cost in an optimal solution, e.g., the total cost of the first $p(J)$ finite-cost time slots. The dynamic program does not enumerate all possible budgets but only a polynomial number of them, namely budgets with integer powers of $(1 + \eta_1)$ with $\eta_1 > 0$ determined later. That is, for the budget on the utilization cost, the dynamic program enumerates all values in

$$B := \{0, 1, (1 + \eta_1), (1 + \eta_1)^2, \dots, (1 + \eta_1)^{\omega_1}\} \text{ with } \omega_1 = \lceil \log_{1+\eta_1} E_{\max} \rceil.$$

The value η_1 will be chosen so that $(1 + \eta_1)^{\omega_1} \leq (1 + \varepsilon)$ and ω_1 is polynomial (see proof of Lemma 7 for the exact definition). Similarly, we observe that $(1 + \varepsilon)^\nu$ is an upper bound on the average scheduling cost. The dynamic program does also only enumerate a polynomial number of possible average scheduling costs, namely integer powers of $(1 + \eta_2)$ with $\eta_2 > 0$ also determined later. This means, for the average scheduling cost, the dynamic program enumerates all values in

$$AVG := \{0, 1, (1 + \eta_2), (1 + \eta_2)^2, \dots, (1 + \eta_2)^{\omega_2}\} \text{ with } \omega_2 = \lceil \nu \log_{1+\eta_2} (1 + \varepsilon) \rceil.$$

As before, the value η_2 will be chosen so that $(1 + \eta_2)^{\omega_2} \leq (1 + \varepsilon)$ and ω_2 is polynomial. The dynamic program stops after the iteration for $u = 0$. Now, only the set of zero-weight jobs is not scheduled yet. For any state $Z = [J_0, b, avg]$ constructed in the iteration for $u = 0$, we append the zero-weight jobs starting at time $t(Z)$ and utilizing the cheapest slots, which is justified by Observation 4. We add the additional utilization cost to b . After this, we return the state $Z = [J_0, b, avg]$ and its corresponding schedule, which can be computed by backtracking and following the established links, with minimum total cost $b + avg \cdot t(Z)$. With this, we obtain the following result.

Lemma 7. *The dynamic program computes a $(1 + O(\varepsilon))$ -approximate solution.*

Proof. Consider an arbitrary iteration u of the dynamic program and let $i = \nu - u$. We consider states $Z = [J_u, b, avg]$ with $J_u \in \mathcal{F}_u$, $b \in B$, and $avg \in AVG$ for which we construct the time points $t(Z)$. Let $Z_1^* = [J_{u+1}^*, b_1^*, avg_1^*]$ and $Z_2^* = [J_u^*, b_2^*, avg_2^*]$ with $J_{u+1}^* \in \mathcal{F}_{u+1}$ and $J_u^* \in \mathcal{F}_u$ be the states that represent an optimal solution \mathcal{S}^* for which the set of jobs with completion weight in WI_{u+1} is exactly $J_{u+1}^* \setminus J_u^*$. By Observation 3, we assume that also in \mathcal{S}^* the area under the remaining weight function $W^{\mathcal{S}^*}(t)$ for the jobs in $J_{u+1}^* \setminus J_u^*$ is approximated by $(1 + \varepsilon)^{u+1} \cdot (t(Z_2^*) - t(Z_1^*))$. We now show the following. The dynamic program constructs in iteration i a state $Z = [J_u, b, avg]$ with $J_u \in \mathcal{F}_u$, $b \in B$, and $avg \in AVG$ such that

- (i) $J_u = J_u^*$,
- (ii) $b \leq (1 + \eta_1)^i \cdot b_2^*$,
- (iii) $avg \leq (1 + \eta_2)^i \cdot avg_2^*$, and
- (iv) $t(Z) \leq t(Z_2^*)$.

We prove this statement by induction on $i = 1, \dots, \nu$. Consider the first iteration of the dynamic program, in which we consider states with job sets from $\mathcal{F}_{\nu-1}$. Let $Z^* = [J_{\nu-1}^*, b^*, avg^*]$ be the state that corresponds to the optimal solution \mathcal{S}^* . The dynamic program also considers the job set $J_{\nu-1}^*$. Suppose, we utilize the same slots that \mathcal{S}^* utilizes for the jobs in $J \setminus J_{\nu-1}^*$ in the interval $[0, t(Z^*)]$. Let b be the resulting utilization cost after rounding b^* up to the next value in B . With this, we know that $b \leq (1 + \eta_1) \cdot b^*$. Furthermore, by our assumption, we know that the average scheduling cost of \mathcal{S}^* up to time $t(Z^*)$ is $(1 + \varepsilon)^\nu$. Let avg be $(1 + \varepsilon)^\nu$ rounded up to the next value in AVG . Then we know that $avg \leq (1 + \eta_2) \cdot avg^*$. The dynamic program also considers the state $Z = [J_{\nu-1}^*, b, avg]$. However, the dynamic program computes the *minimum* time point $t(Z_\nu, Z) \geq t^{LB}$ so that the set of jobs $J \setminus J_{\nu-1}^*$ can be feasibly scheduled in $[0, t(Z_\nu, Z))$ having utilization cost not more than b . This implies that $t(Z_\nu, Z) \leq t(Z^*)$, which implies that $t(Z) \leq t(Z^*)$. Note that $t^{LB} = 0$ for the specified values in Z .

Suppose, the statement is true for the iterations $1, 2, \dots, i - 1$. We prove that it is also true for iteration i , in which we consider job sets from \mathcal{F}_u . Again, let $Z_1^* = [J_{u+1}^*, b_1^*, avg_1^*]$ and $Z_2^* = [J_u^*, b_2^*, avg_2^*]$ with $J_{u+1}^* \in \mathcal{F}_{u+1}$ and $J_u^* \in \mathcal{F}_u$ be the states that represent \mathcal{S}^* . By our hypothesis, we know that the dynamic program constructs a state $Z_1 = [J_{u+1}, b_1, avg_1]$ with

- (i) $J_{u+1} = J_{u+1}^*$,
- (ii) $b_1 \leq (1 + \eta_1)^{i-1} \cdot b_1^*$,
- (iii) $avg_1 \leq (1 + \eta_2)^{i-1} \cdot avg_1^*$, and
- (iv) $t(Z_1) \leq t(Z_1^*)$.

We augment this schedule in the following way. Suppose, we utilize the same slots that \mathcal{S}^* utilizes for the jobs in $J_{u+1}^* \setminus J_u^*$ in the interval $[t(Z_1^*), t(Z_2^*)]$. Let b_2 be the resulting total utilization cost after rounding up to the next value in B . Thus, there is a feasible schedule for $J \setminus J_u^*$ having utilization cost of at most

$$\begin{aligned} b_2 &\leq (1 + \eta_1) \cdot (b_1 + b_2^* - b_1^*) \\ &\leq (1 + \eta_1)^i \cdot (b_1^* + b_2^* - b_1^*) \\ &= (1 + \eta_1)^i \cdot b_2^*. \end{aligned}$$

The new average scheduling cost after rounding to the next value in AVG is

$$\begin{aligned}
 avg_2 &\leq (1 + \eta_2) \cdot \frac{avg_1 \cdot t(Z_1) + (1 + \varepsilon)^{u+1} \cdot (t(Z_2^*) - t(Z_1))}{t(Z_2^*)} \\
 &\leq (1 + \eta_2)^i \cdot \frac{avg_1^* \cdot t(Z_1) + (1 + \varepsilon)^{u+1} \cdot (t(Z_2^*) - t(Z_1))}{t(Z_2^*)} \\
 &\leq (1 + \eta_2)^i \cdot \frac{avg_1^* \cdot t(Z_1^*) + (1 + \varepsilon)^{u+1} \cdot (t(Z_2^*) - t(Z_1^*))}{t(Z_2^*)} \\
 &= (1 + \eta_2)^i \cdot avg_2^*.
 \end{aligned}$$

The third inequality follows from the fact that $avg_1^* \geq (1 + \varepsilon)^{u+1}$. The dynamic program also considers the link between the state Z_1 and $Z_2 := [J_u^*, b_2, avg_2]$. We first observe that $t^{LB} \leq t(Z_2^*)$, since

$$avg_2 \cdot t(Z_2^*) \geq avg_1 \cdot t(Z_1) + (1 + \varepsilon)^{u+1} \cdot (t(Z_2^*) - t(Z_1))$$

by construction of avg_2 . Furthermore, we observe that $b_2 - b_1 \geq b_2^* - b_1^*$ by construction of b_2 . These two facts together with $t(Z_1) \leq t(Z_1^*)$ imply that $t(Z_1, Z_2) \leq t(Z_2^*)$, which implies that $t(Z_2) \leq t(Z_2^*)$.

To complete the proof, we need to specify the parameters η_1 and η_2 . We want that $(1 + \eta_i)^\nu \leq (1 + \varepsilon)$ for $i = 1, 2$. We claim that for a given $\nu \geq 1$ there exists an $\bar{\eta} > 0$ such that for all $\eta \in (0, \bar{\eta}]$ we have $(1 + \eta)^\nu \leq 1 + 2\nu\eta$. Consider the function $f(\eta) := (1 + \eta)^\nu - 1 - 2\nu\eta$. We have that $f(0) = 0$ and $f'(\eta) < 0$ for $\eta \in [0, 2^{1/(\nu-1)} - 1)$. This shows the claim. Hence, we choose $\eta_i = \min\{\frac{\varepsilon}{2\nu}, 2^{1/(\nu-1)} - 1\}$ for $i = 1, 2$. This shows the statement of the lemma and that the size of B as well as the size of AVG are bounded by a polynomial in the size of the input.

We remark that the given DP works for more general utilization cost functions $e : \mathbb{N} \rightarrow \mathbb{Q}_{\geq 0}$ than considered here in the paper. As argued in the proof, it is sufficient for the DP that there is a function $E(p, [t_1, t_2])$ that outputs in polynomial time for a given time interval $[t_1, t_2]$ and a given $p \in \mathbb{Z}_{\geq 0}$ the total cost of the p cheapest slots in $[t_1, t_2]$.

We also remark that the running time of the presented DP is exponential, because the size of the sets \mathcal{F}_u are exponential in the size of the input. However, in the next section we show that we can trim the sets \mathcal{F}_u down to ones of polynomial size at an arbitrarily small loss in the performance guarantee.

3.3 Trimming the state space

The set \mathcal{F}_u , containing all possible job sets J_u , is of exponential size, and so is the DP state space. In the context of scheduling with variable machine speed, it has been shown in [17] how to reduce the set \mathcal{F}_u for a similar DP (without utilization decision, though) to a set $\tilde{\mathcal{F}}_u$ of polynomial size at only a small loss in the objective value. In general, such a procedure is not necessarily applicable to our setting because of the different objective involving additional utilization cost and the different decision space. However, the compactification in [17] holds *independently of the speed of the machine* and, thus, independently of the utilization decision of the DP (interpret non/utilization as speed 0/1). Hence, we can apply it to our cost-aware scheduling framework and obtain a PTAS. We now describe the building blocks for this trimming procedure and argue why we can apply it in order to obtain the set $\tilde{\mathcal{F}}_u$ for our problem.

Light Jobs. The first building block for the trimming procedure is a classification of the jobs based on their weights.

Definition 2. *Given a weight schedule and a job $j \in J$ with starting weight $S_j^w \in WI_u$, we call job j light if $w_j \leq \varepsilon^2 |WI_u|$, otherwise j is called heavy.*

This classification enables us to structure near-optimal solutions. To impose structure on the set of light jobs, the authors in [17] describe the following routine for a given weight schedule \mathcal{S} . First, remove all light jobs from \mathcal{S} and move the remaining jobs within each interval WI_u so that the idle weight in WI_u is consecutive. Then, schedule the light jobs according to the *reverse Smith's rule*, that is, for each $u = 1, \dots, \nu$ and each idle weight $w \in WI_u$, process at w a light job j that maximizes p_j/w_j . Eventually, shift the processing of each interval WI_u to WI_{u+1} , which delays the completion of every job by at most a factor of $(1 + \varepsilon)^2$. This delay allows to completely process every light job in the weight interval where it starts processing. It can be shown that the cost of the resulting schedule is at most a factor of $1 + O(\varepsilon)$ greater than the cost of \mathcal{S} , which brings us to the following structural statement.

Lemma 8 ([17]). *At a loss of a factor of $1 + O(\varepsilon)$ in the scheduling cost, we can assume the following. For a given interval WI_u , consider any pair of light jobs j, k . If both jobs start in WI_u or later and $p_k/w_k \leq p_j/w_j$, then $C_j^w \leq C_k^w$.*

We remark, that Lemma 8 holds independently of the speed of the machine, as pointed out in [17]. This means that at a loss of a factor of $1 + O(\varepsilon)$ in the scheduling cost we can assume also for our problem that light jobs are scheduled according to *reverse Smith's rule* in the weight-dimension, which holds independently of our actual utilization decision.

Localization. We now localize jobs in the weight-dimension to gain more structure. That is, we determine for every job $j \in J$ two values r_j^w and d_j^w such that, independently of our actual utilization decision, j is scheduled completely within $[r_j^w, d_j^w]$ in some $(1 + O(\varepsilon))$ -approximate weight-schedule (in terms of the scheduling cost). We call r_j^w and d_j^w the *release-weight* and the *deadline-weight* of job j , respectively.

Lemma 9 ([17]). *We can compute in polynomial time values r_j^w and d_j^w for each $j \in J$ such that: (i) there exists a $(1 + O(\varepsilon))$ -approximate weight-schedule (in terms of the scheduling cost) that processes each job j within $[r_j^w, d_j^w]$, (ii) there exists a constant $s \in O(\log(1/\varepsilon)/\varepsilon)$ such that $d_j^w \leq r_j^w \cdot (1 + \varepsilon)^s$, (iii) r_j^w and d_j^w are integer powers of $(1 + \varepsilon)$, and (iv) the values r_j^w and d_j^w are independent of the speed of the machine.*

This lemma enables us to localize all jobs in J in polynomial time and independent of our actual utilization decision, as guaranteed by property (iv).

Compact Search Space. Based on the localization of jobs in weight space, we can cut the number of different possibilities for a candidate set J_u in iteration u of our DP down to a polynomial number. That is, we replace the set \mathcal{F}_u by a polynomially sized set $\tilde{\mathcal{F}}_u$. Instead of describing all sets $S \in \tilde{\mathcal{F}}_u$ explicitly, we give all possible complements $R = J \setminus S$ and collect them in a set \mathcal{D}_u , where a set $R \in \mathcal{D}_u$ represents a possible set of jobs having completion weights in WI_{u+1} or later. Obviously, a set $R \in \mathcal{D}_u$ must contain all jobs $j \in J$ having a release weight $r_j^w \geq (1 + \varepsilon)^u$. Furthermore, we know that $d_j^w \geq (1 + \varepsilon)^{u+1}$ is necessary for job j to be in a set $R \in \mathcal{D}_u$. Following property (ii) in Lemma 9, we thus only need to decide about the jobs having a release weight $r_j^w = (1 + \varepsilon)^i$ with $i \in \{u + 1 - s, \dots, u - 1\}$. An enumeration over basically all possible job sets for each $i \in \{u + 1 - s, \dots, u - 1\}$ gives the following desired result.

Lemma 10 ([17]). *For each u , we can construct in polynomial time a set $\tilde{\mathcal{F}}_u$ that satisfies the following: (i) there exists a $(1 + O(\varepsilon))$ -approximate weight-schedule (in terms of the scheduling cost) in which the set of jobs with completion weight at most $(1 + \varepsilon)^u$ belongs to $\tilde{\mathcal{F}}_u$, (ii) the set $\tilde{\mathcal{F}}_u$ has cardinality at most $2^{O(\log^3(1/\varepsilon)/\varepsilon^2)}$, and (iii) the set $\tilde{\mathcal{F}}_u$ is completely independent of the speed of the machine.*

Again, Property (iii) implies that we can construct the set $\tilde{\mathcal{F}}_u$ independently of our utilization decision.

To complete the proof of Theorem 3 it remains to argue on the running time of the DP. The DP has ν iterations, where in each iteration for at most $2^{O(\log^3(1/\varepsilon)/\varepsilon^2)} \cdot |B| \cdot |AVG|$ previous states at most $2^{O(\log^3(1/\varepsilon)/\varepsilon^2)} \cdot |B| \cdot |AVG|$ many links to new states are considered. Therefore, the running time complexity of our DP is $\nu \cdot (2^{O(\log^3(1/\varepsilon)/\varepsilon^2)} \cdot |B| \cdot |AVG|)^2$, which is bounded by a polynomial in the size of the input.

4 Minimizing the makespan on unrelated machines

Finally we derive positive results for the problem of minimizing makespan with utilization costs on unrelated machines. The standard scheduling problem without utilization cost $R|pmtn|C_{\max}$ can be solved optimally in polynomial time by solving a linear program as was shown by Lawler and Labetoulle [15]. We show that the problem complexity does not increase significantly when taking into account time-varying utilization cost.

Consider the preemptive makespan minimization problem with utilization cost. Recall that we can use every machine in a utilized time slot and pay only once. Thus, it is sufficient to find an optimal utilization decision for solving this problem, because we can use the polynomial-time algorithm in [15] to find the optimal schedule within these slots.

Observation 5 *Given the set of time slots utilized in an optimal solution, we can compute an optimal schedule in polynomial time.*

Given an instance of our problem, let Z be the optimal makespan of the relaxed problem *without* utilization cost. Notice that Z is not necessarily integral. To determine an optimal utilization decision, we use the following observation.

Observation 6 *Given an optimal makespan C_{\max}^* for $R|pmtn|C_{\max} + E$, an optimal schedule utilizes the $\lceil Z \rceil$ cheapest slots before $\lceil C_{\max}^* \rceil$.*

Note that we must pay full tariff for a used time slot, no matter how much it is utilized, and so does an optimal solution. In particular, this holds for the last utilized slot. Hence, it remains to compute an optimal value $C^* := \lceil C_{\max}^* \rceil$ which we do by the following procedure.

We compute for every interval $I_k = [s_k, d_k)$, $k = 1, \dots, K$, an optimal point in time for C^* assuming that $C^* \in I_k$. Hereby we restrict to relevant intervals I_k which allow for a feasible schedule, i.e., $s_k \geq \lceil Z \rceil$. For a relevant interval I_k , we let $C^* = s_k$ and utilize the $\lceil Z \rceil$ cheapest time slots before C^* , which is optimal by Observation 6. Notice that any utilized time slot of cost e such that $e > e_k + 1$ can be replaced by a time slot from I_k leading to a solution of less total cost. Thus, if there is no such time slot then s_k is the best choice for C^* in I_k . Suppose there is such a time slot that could be replaced. Let $R \subseteq \{1, \dots, k-1\}$ be the index set of intervals that contain at least one utilized slot. We define I_ℓ to be the interval with $e_\ell = \max_{h \in R} e_h$ and denote by r_h the number of utilized time slots in I_h . Replace $\min\{r_\ell, d_k - s_k - r_k\}$ utilized slots from I_ℓ by slots from I_k and update R , I_ℓ and r_k . This continues until $e_\ell \leq e_k + 1$ or the interval I_k is completely utilized, i.e., $r_k = d_k - s_k$. This operation takes at most $O(K)$ computer operations per interval to compute the best C^* -value in that interval. It yields the following theorem.

Theorem 4. *The scheduling problem $R|pmtn|C_{\max} + E$ can be solved in polynomial time in the order of $O(K^2)$ plus the running time for solving $R|pmtn|C_{\max}$ without utilization cost [15].*

5 Conclusion

We investigate basic scheduling problems within the framework of time-varying costs or tariffs, where the processing of jobs causes some time-dependent cost in addition to the usual QoS measure. We presented optimal algorithms and best possible approximation algorithms for the scheduling objectives of minimizing the makespan on unrelated machines and the sum of (weighted) completion times on a single machine.

While our work closes the problems under consideration from an approximation point of view, it leaves open the approximability of multi-machine settings for the min-sum objective. Further research may also ask for the complexity status when assuming that jobs have different release dates and for other natural objective functions such as average and maximum flow-time.

Our unrelated machine model is time-slot based, that is, a utilization decision is made for a time slot and then all machines in this time slot are available. No less relevant appears to be the model with *machine-individual* tariffs, that is, a utilization decision is made for a time slot on each machine individually. It is not difficult to see that a standard LP can be adapted for optimally solving $R|pmtn, r_j|C_{\max}$ with fractional utilization cost. However, if time slots can be utilized only integrally then the integrality gap for the simple LP is unbounded and the problems seems much harder.

Time-varying cost or tariffs appear in many applications in practice but they are hardly investigated from a theoretical perspective. With our work we settle the complexity status and approximability status for very classical scheduling problems. We hope to foster further research on this framework of time-varying costs or tariffs. We emphasize that the framework is clearly not restricted to cost-aware scheduling problems. Virtually any problem in which scarce resources are to be rented from some provider lends itself to be modelled in this way, with (vehicle) routing problems as a directly appealing example.

References

1. S. Albers. Energy-efficient algorithms. *Commun. ACM*, 53(5):86–96, 2010.
2. N. Bansal and K. Pruhs. The geometry of scheduling. *SIAM J. Comput.*, 43(5):1684–1698, 2014.
3. B. Chen and X. Zhang. Scheduling with time-of-use costs. *European Journal of Operational Research*, 2018.
4. L. Chen, N. Megow, R. Rischke, L. Stougie, and J. Verschae. Optimal algorithms and a PTAS for cost-aware scheduling. In *Proceedings of the 40th International Symposium on Mathematical Foundations of Computer Science (MFCS)*, volume 9235 of *LNCS*, pages 211–222. Springer, 2015.
5. M. Cheung, J. Mestre, D. Shmoys, and J. Verschae. A primal-dual approximation algorithm for min-sum single-machine scheduling problems. *SIAM J. on Discrete Mathematics*, 31(2):825–838, 2017.
6. M. Cheung and D. B. Shmoys. A primal-dual approximation algorithm for min-sum single-machine scheduling problems. In *Proceedings of the 14th International Workshop on Approximation, Randomization, and Combinatorial Optimization (APPROX)*, volume 6845 of *LNCS*, pages 135–146, 2011.
7. W. L. Eastman, S. Even, and M. Isaac. Bounds for the optimal scheduling of n jobs on m processors. *Management Sci.*, 11(2):268–279, 1964.
8. L. Epstein, A. Levin, A. Marchetti-Spaccamela, N. Megow, J. Mestre, M. Skutella, and L. Stougie. Universal sequencing on an unreliable machine. *SIAM J. Comput.*, 41(3):565–586, 2012.
9. K. Fang, N. A. Uhan, F. Zhao, and J. W. Sutherland. Scheduling on a single machine under time-of-use electricity tariffs. *Annals of Operations Research*, 238(1):199–227, 2016.
10. R. L. Graham, E. L. Lawler, J. K. Lenstra, and A. H. G. R. Kan. Optimization and approximation in deterministic sequencing and scheduling: A survey. In *Annals of Discrete Mathematics*, volume 5, pages 287–326. Elsevier, 1979.
11. F. S. Hillier. *Introduction to operations research*. Tata McGraw-Hill Education, 2012.
12. W. Höhn and T. Jacobs. On the performance of smith’s rule in single-machine scheduling with nonlinear cost. *ACM Trans. Algorithms*, 11(4):25, 2015.
13. W. Höhn, J. Mestre, and A. Wiese. How unsplittable-flow-covering helps scheduling with job-dependent cost functions. In *Proceedings of the 41st International Colloquium on Automata, Languages, and Programming (ICALP)*, volume 8572 of *LNCS*, pages 625–636. Springer, 2014.
14. J. Kulkarni and K. Munagala. Algorithms for cost-aware scheduling. In *Proceedings of the 10th International Workshop on Approximation and Online Algorithms (WAOA)*, volume 7846 of *LNCS*, pages 201–214. Springer, 2013.
15. E. L. Lawler and J. Labetoulle. On preemptive scheduling of unrelated parallel processors by linear programming. *J. ACM*, 25(4):612–619, 1978.
16. C.-Y. Lee. Machine scheduling with availability constraints. In J. Y.-T. Leung, editor, *Handbook of Scheduling*. CRC Press, 2004.

17. N. Megow and J. Verschae. Dual techniques for scheduling on a machine with varying speed. In *Proceedings of the 40th International Colloquium on Automata, Languages, and Programming (ICALP)*, volume 7965 of *LNCS*, pages 745–756. Springer, 2013.
18. R. Rischke. *Deterministic, Stochastic, and Robust Cost-Aware Scheduling*. PhD thesis, Technical University of Munich, 2016.
19. H. A. Taha. *Operations research: An introduction*. Pearson/Prentice Hall, 2007.
20. K. T. Talluri and G. J. Van Ryzin. *The theory and practice of revenue management*, volume 68. Springer Science & Business Media, 2006.
21. G. Wan and X. Qi. Scheduling with variable time slot costs. *Naval Research Logistics*, 57:159–171, 2010.
22. G. Wang, H. Sun, and C. Chu. Preemptive scheduling with availability constraints to minimize total weighted completion times. *Ann. Oper. Res.*, 133:183–192, 2005.