

Multi-token Markov Game with Switching Costs

Jian Li *

Daogao Liu †

Abstract

We study a general Markov game with metric switching costs: in each round, the player adaptively chooses one of several Markov chains to advance with the objective of minimizing the expected cost for at least k chains to reach their target states. If the player decides to play a different chain, an additional switching cost is incurred. The special case in which there is no switching cost was solved optimally by Dumitriu, Tetali and Winkler [DTW03] by a variant of the celebrated Gittins Index for the classical multi-armed bandit (MAB) problem with Markovian rewards [Git74, Git79]. However, for Markovian multi-armed bandit with nontrivial switching cost, even if the switching cost is a constant, the classic paper by Banks and Sundaram [BS94] showed that no index strategy can be optimal.¹

In this paper, we complement their result and show there is a simple index strategy that achieves a constant approximation factor if the switching cost is constant and $k = 1$. To the best of our knowledge, this index strategy is the first strategy that achieves a constant approximation factor for a general Markovian MAB variant with switching costs. For the general metric, we propose a more involved constant-factor approximation algorithm, via a nontrivial reduction to the stochastic k -TSP problem, in which a Markov chain is approximated by a random variable. Our analysis makes extensive use of various interesting properties of the Gittins index.

*Institute for Interdisciplinary Information Sciences, Tsinghua University. Email: lijian83@mail.tsinghua.edu.cn.

†Paul G. Allen School of Computer Science & Engineering, University of Washington. Email: dgliu@cs.washington.edu.

¹Their proof is for the discounted version of MAB, but can be extended to our setting. See Appendix D for the details.

1 Introduction

The Markovian multi-armed bandit (MAB) problem is one of the most important and well studied sequential decision problem. In this problem, at each time step, the agent knows the state of each chain and must choose to play one of n available Markov chains. The agent pays a certain cost (or receives a payoff) depending on the current state, and the chosen chain advances to the next state (according the Markovian transition rules). The goal of the agent is to optimize the expected cost (or payoff) by choosing the right sequence of actions. The infinite horizon discounted version of the problem was solved optimally by the celebrated Gittins Index theorem, first proved by Gittins and Jones [Git74]. In particular, they show that the Markov chains can be “indexed” separately, and the optimal strategy is simply choosing the chain with the smallest (or largest) index. Since then, Gittins Index has been studied and extent in a variety of ways (see [GGW11]).

A major extension to MAB is the inclusion of switching costs, that is switching to a different chain incurs a nontrivial cost [BS94, AT96, KSU08, KLM17, CGT⁺20]. This extension has found many applications in job search and labor mobility [Joh78, Vis80, Mac80, McL84, Wal84, Jov84, KW11], industrial policy [PT95, Kli04], optimal search [Wei79, BB88, BGO92, Smi99], experiment and learning [Rot74, McL84, Krä03, ASKW04] and game theory [Sch97, BV06]. The most natural problem for MAB in the presence of switching costs is to examine the extend to which the Gittins-Jones theorem remains valid, i.e., whether there is a suitably defined index strategy that is optimal. This problem was first studied in the classic paper by Banks and Sundaram [BS94], who showed that there is no index strategy that is optimal, even if the switching cost is a given nonzero constant. Motivated by this work, several authors [BGO92, AT96, VOP00, KL00, BV01, Jun01, DH03] attempted to (partially) characterize the optimal policy and present optimal solutions for several special cases. In fact, MAB with switching costs is a special case of the restless bandit problem introduced by Whittle [Whi88]: the state of the arm just abandoned changes its state to a “dummy copy” state which requires a switching cost if it is to be played. However, the restless bandit problem is known to be PSPACE-Hard, even to approximate to any non-trivial factor [PT94]. See the survey [Jun04].

We approach the problem from the perspective of approximation algorithms and focus on a finite-time version of MAB, called *multi-token Markov game*, introduced in an elegant paper by Dumitriu, Tetali and Winkler [DTW03]. In this game, we are also given n Markov chains and each chain has a target which is ultimately reachable. Each state is associated with a movement cost. In each time step, the player adaptively chooses one Markov chain to advance with the objective of minimizing the expected total cost for at least one chain to reach its target state. If there is no switching cost, they show that there is an optimal indexing strategy based on a variant of Gittins index. Even if the switching cost is a given nonzero constant, by a similar argument in [BS94], one can show that there is no indexing strategy that is optimal (see Appendix D). Hence, in this paper, we study approximation algorithms for the multi-token Markov game with switching costs.

1.1 Problem Definitions and Our Contributions

We formally define our problem as follows. We mainly follow the terminology used in [DTW03]. We first introduce the notions for a Markov system, which is simply a Markov chain with (state) movement cost.

Definition 1 (Markov System [DTW03]). *A Markov system is a tuple $\mathcal{S} = \langle V, P, C, s, t \rangle$, where V is the finite set of states, $P = \{P_{u,v}\}$ is the corresponding transition matrix (a $|V| \times |V|$ matrix), $C = \{C_u\}$ denotes a positive real movement cost for each state $u \in V$, and s (resp. t) represents the current (resp. target) state. We*

assume that the target is ultimately reachable from every state in V , and we can never exit the target state (so we can set $C_t = 0$ and $P_{t,t} = 1$). If we play \mathcal{S} in state $u \in V$, a cost of $C_u \geq 0$ is incurred and \mathcal{S} transitions from state u to state v with probability $P_{u,v}$. There is a unit reward on the target state t that we can collect.

In the following, we do not distinguish Markov systems and Markov chains, and use both terms interchangeably.

1.1.1 Unit Switching Cost

Now we define our first problem, *multi-token Markov game with unit switching cost* (MG-Unit). In this problem, we have a set of (possibly different) Markov system and switching from one Markov system to another Markov system incurs a *unit* cost. Our goal is to find a strategy that adaptively chooses the next Markov system to play until a unit of reward is collected, and the expected total cost (switching cost plus movement cost) is minimized. Formally, the problem is defined as follows.

Definition 2 (MG-Unit). *We are given a metric space $\mathcal{M} = (\mathcal{S} \cup \{\mathbf{R}\}, d)$ endowed with unit metric (i.e., $d(\mathcal{S}, \mathcal{S}') = 1$ for any $\mathcal{S} \neq \mathcal{S}' \in \mathcal{S} \cup \{\mathbf{R}\}$). Each node $\mathcal{S}_i \in \mathcal{S}$ is identified with a Markov system $\mathcal{S}_i = \langle V_i, P_i, C_i, s_i, t_i \rangle$. \mathbf{R} is the root node (the initial position) with a unit cost directed edge to every \mathcal{S}_i . If we play Markov system \mathcal{S}_i in one round and decide to play another Markov system \mathcal{S}_j in the next round, we need pay a unit switching cost $d(\mathcal{S}_i, \mathcal{S}_j) = 1$ in addition to the movement cost. The game ends when we succeed to make one Markov system reach its target state.*

For MG-Unit, we provide a simple index strategy that has a constant approximation ratio. Here, following the the definition in [BS94] (See also Definition 38 in the appendix), an index strategy means that we can define a suitable index (a real number) for each state of the Markov chains, and the strategy always chooses to play the Markov chain in which the current state has the minimum index.

Theorem 3. *There is a simple index strategy that can achieve a constant approximation ratio for the MG-Unit problem.*

Our technique: In particular, for each state u , we create a *dummy state* u' that connects to u with movement cost 1. This captures the unit switching cost. The index Γ_i for \mathcal{S}_i (at state u) is the Gittins index γ_u if \mathcal{S}_i is active, and the Gittins index of u 's dummy state u' otherwise. Our strategy is to simply choose to play the \mathcal{S}_i with the smallest index Γ_i .

Although the strategy is extremely simple to state, it is difficult to analyze it directly. Instead, we analyze an alternative strategy, via the doubling framework developed in recent works [ENS18, JLLS20] (Section 3.2). In the doubling framework, we proceeds in phases and in each phase there is an exponentially increasing cost budget. Under this framework, it suffices to show the following guarantee for the budgeted sub-problem: our strategy can succeed (i.e., collect one unit of reward from some Markov chain) with constant probability, under constraint that the total (movement plus switching) cost is below the given budget B , which is constant times the optimal cost of the original problem (Lemma 12). Solving the budgeted sub-problem this is the key technical challenge.

In order to show we can succeed with constant probability under the budget constraint, we consider the set Ω_{bad} of trajectories in which we fail (do not reach the target). A trajectory can be naturally partitioned into segments, each being a trajectory in one Markov chain and the switching cost we pay is the number

of the segments. Since the cost budget is exhausted without success, one can show that there is one segment (corresponding to one chain) that the expected cost is large but the success probability is small (conditioning on the former segments). From this, one can argue the grade (Gittins index) of the current chain is much larger than B , and the grades of all chains are also at least no smaller than it (due to the greedy rule). By the definition of the grade, it can be roughly understood as the expected movement cost one can hope for reaching the target in this chain. Hence, one can see that conditioning on Ω_{bad} , the expected movement cost to reach a target for any strategy is much larger than B (no chain is cheap). But the expected total cost of the optimal strategy is much less than B (recall B is a large constant times OPT). Therefore, one can conclude that the probability of Ω_{bad} is small.

1.1.2 General Metric Switching Cost

Next, we generalize MG-Unit to more general metric (where the switching costs form a metric without the restriction to be unit) and more general requirement that we need to collect K units of rewards for any positive integer K . We name the new problem MG-Metric.

Definition 4 (MG-Metric). *We are given a finite metric space $\mathcal{M} = (\mathbf{S} \cup \{\mathbf{R}\}, d)$ (there is no additional assumption on metric d). Each node $\mathcal{S}_i \in \mathbf{S}$ is identified with a Markov system $\mathcal{S}_i = \langle V_i, P_i, C_i, s_i, t_i \rangle$. Similarly, at the beginning of the game, the player is at the root \mathbf{R} , and needs to pay the switching cost $d(\mathbf{R}, \mathcal{S}_i)$ if he wants to play Markov system \mathcal{S}_i . Switching from \mathcal{S}_i to \mathcal{S}_j incurs a switching cost of $d(\mathcal{S}_i, \mathcal{S}_j)$. The objective is to adaptively collect at least K units of rewards (making at least K Markov system reach their targets), while minimizing the expected total cost (movement cost plus switching cost). The game ends when we succeed to make K Markov system reach their target states.*

Theorem 5. *There is a constant factor approximation algorithm for the MG-Metric problem.*

Our technique: For MG-Metric, we also adopt the doubling framework, hence only need to design an algorithm BudgetMG-Metric (Algorithm 4) for the budgeted sub-problem. BudgetMG-Metric should succeed with constant probability using a budget B when B is a constant factor of the optimal cost. At a high level, BudgetMG-Metric first transform the problem to a Stochastic- k -TSP instance \mathcal{M}_{ktsp} by reducing each Markov chain to a related random variable. Then it applies the *non-adaptive* constant factor approximation for Stochastic- k -TSP (developed in [JLLS20]) to obtain an ordering Π of vertices (chains). We pick a prefix Π_{pref} such that the switching cost for traversing Π_{pref} is no larger than a small constant proportion of the budget. One can show it is possible to collect K units of rewards from Π_{pref} such that the total movement cost is within the budget with constant probability.

Now, the key is show how to collect K units of rewards from Π_{pref} such that the movement plus switching cost is within the budget with constant probability. Obviously, ignoring the switching cost, the optimal way (optimal in terms of movement cost) of collecting K units of rewards from Π_{pref} is to play Gittins index. However, such play may switch back and forth frequently and hence leads to a high switching cost. To keep the switching cost under control, we insist visiting the chain in Π_{pref} one by one and never revisit any chain (hence the switching cost is small). However, one may not want to play a chain to the end since the current state is not economical to play and switching to the next chain is a better option. Now the Gittins index comes into rescue. We show that there is an interesting threshold γ_{j+1} (which is computed from the K -th order statistics of suitably defined random variables), such that if the Gittins index of the current state is larger than the threshold, we can give up and decide to switch to the next chain on the

Π_{pref} . It turns out such a sequential algorithm (without switching back and forth) can also succeed with constant probability without incurring a much larger movement cost than keeping playing the chain with the smallest Gittins index.

2 Related Work

The original paper by [DTW03] only deals with the $K = 1$ case (i.e., one chain reaches its target). Recent works [KWW16, GJSS19] observe that it is not difficult to extend their argument to general positive integer K without switching costs. [KWW16, Sin18] studied the problem under richer combinatorial constraints. [GJSS19] study a more general problem: there is a given packing or covering constraint $\mathcal{F} \subseteq 2^{[n]}$ (e.g., matroid, matching, knapsack) of subsets of chains. The goal is to make a subset S of chains to reach their targets ($S \in \mathcal{F}$), while minimizing the dis-utility (with upward-closed constraint) or maximize the utility (with downward-closed constraint). For semi-additive objective function, they proposed a general reduction which utilizes the “greedy” algorithms for the problem with full information and they can achieve the same approximation ratio as the “greedy” algorithm does for the full information problem.

Guha and Munagala [GM09] considered two MAX-SNP problems for bandits with switching costs: *future utilization* and *past utilization*. In their problems, each state has a reward and the rewards satisfies the *martingale property* (motivated by Bayesian considerations, see their paper for the details). Given two budgets for movement and switching, the future utilization problem aims to make the final reward of the finally chosen chain as large as possible, while the goal of the past utilization is to make the summation of rewards as large as possible. They provided $O(1)$ -approximation algorithms for both problems. They approached the problem from linear programming with Lagrangian relaxation. Their problems are very different from our problems and it is unclear how to apply their technique to our problems neither.

Our problem is also related to some problems in the stochastic probing literature, in particular the classical Pandora’s Problem defined in [Wei79]. Suppose there are n closed boxes with independent random rewards (with known distributions). The cost to open box i is c_i . When we open a box, the reward of the box is realized. At each time step, we need to decide either to pay some cost to open a new box, or stop and take the box with the maximum rewards. The goal is to maximize the expected reward minus the opening cost. Weitzman [Wei79] provided an optimal indexing strategy to this problem. In fact, one can show the problem is a special case of the Markov Game [DTW03] and Weitzman’s index can also be seen as a variant of the Gittin’s Index. Recently, the problem has been extent in various ways (see e.g., [KWW16, Dov18, BK19]).

Our problem is a stochastic combinatorial optimization problem. Designing poly-time algorithms for those problems with provable approximation guarantee has attracted significant attention in recent years (see e.g., the survey [LL16]). In this paper, we leverage the constant factor approximation algorithm for stochastic k -TSP [JLLS20] (formally defined later), which is closely related to the stochastic knapsack and stochastic orienteering problems. In stochastic knapsack, we are given a set of items with random size and profit and a knapsack with fixed capacity. We can adaptively place the items in the knapsack irrevocably, such that the expected total profit is maximized. A variant of the stochastic knapsack has been shown to be PSPACE-hard [DGV04], and several constant factor approximation algorithms have been developed [DGV04, DGV08, BGK11, GKNR12, LY13]. Stochastic orienteering [GKNR12] is a generalization of stochastic knapsack, in which there are metric switching costs between different items. If the total cost is restricted to be no more than B , Gupta et al. [GKNR12] provided an $O(\log \log B)$ upper bound of the

adaptivity gap, and Bansal and Nagarajan [BN15] showed a lower bound of $\Omega((\log \log B)^{1/2})$ even when all profits are deterministic.

In online learning literature, there is also a body of work [KSU08, Ort08, ADT12, DDKP14, KLM17] studying multi-armed bandit (MAB) with switching costs. However, here playing each arm provides i.i.d. reward, but the underlying distribution is not known. The objective is minimizing the regret. The challenges and techniques in these settings are completely different.

3 Preliminaries

In this section, we first review the notion of *grade* (a variant of Gittins index) introduced in [DTW03], then the *doubling technique* used in some previous stochastic optimization problems [ENS18, JLLS20]. The analysis requires some well known concentration inequalities such as Chernoff Bound and Freedman’s Inequality, which are presented in Appendix A.1. We define some notations which will be used throughout this paper.

Notation: For any (possibly adaptive) strategy \mathbb{P} , let $R(\mathbb{P})$ be the (random) number of units of rewards \mathbb{P} can collect (i.e. the (random) number of Markov system that reach target states). Let $C_{sw}(\mathbb{P})$ and $C_{mv}(\mathbb{P})$ be the (random) switching cost and movement cost respectively, and let $C_{tot}(\mathbb{P}) = C_{sw}(\mathbb{P}) + C_{mv}(\mathbb{P})$ be the total cost of the strategy \mathbb{P} . We say a strategy \mathbb{P} with movement budget (resp. switching budget) B means that \mathbb{P} can pay at most B to advance the Markov system (resp. to switch between different systems), and say a strategy \mathbb{P} with budget B means that if its movement cost plus switching cost is restricted to be at most B . Similarly, we say a strategy \mathbb{P} with (movement/switching) budget B in expectation means that the expectation of $(C_{mv}(\mathbb{P})/C_{sw}(\mathbb{P})) C_{tot}(\mathbb{P})$ is at most B .

3.1 Grade

Give a finite Markov game defined in Dumitriu et al. [DTW03], we can define the *grade* of each state in each Markov system. Grade is a slight variant of the original Gittins index defined for the infinite discounted game [Git74], particularly defined in a very similar way to Weber’s prevailing charge [W⁺92, FW99]. In particular, the grade of state u in Markov system \mathcal{S} depends only on \mathcal{S} , not other Markov system.

A New Game $\mathcal{S}_u(g)$: Consider a Markov system $\mathcal{S} = \langle V, P, C, s, t \rangle$. Given a non-negative real number $g \in \mathbb{R}_{\geq 0}$ and the initial state u in \mathcal{S} , we define a new game $\mathcal{S}_u(g)$: For each step, we can quit and end the game, **or** pay the movement cost to advance \mathcal{S} for one step. If we reach the target state t , we can get g units of profits² and the game halts. The objective is to maximize the objective “profit – cost” in expectation. We use $\text{val}(\mathcal{S}_u(g))$ to denote the value of this game, that is the expected objective achieved by the optimal strategy, which we denote by $\mathbb{O}(\mathcal{S}_u(g))$. Let $R(\mathbb{O}(\mathcal{S}_u(g)))$ is the (random) indicator that the strategy $\mathbb{O}(\mathcal{S}_u(g))$ reaches the target (e.g., if we quit the game before reaching the target, $R = 0$), and hence $\mathbb{E}[R(\mathbb{O}(\mathcal{S}_u(g)))]$ is the probability that we reach the target state using strategy $\mathbb{O}(\mathcal{S}_u(g))$. It is easy to see

$$\text{val}(\mathcal{S}_u(g)) = g \cdot \mathbb{E}[R(\mathbb{O}(\mathcal{S}_u(g)))] - \mathbb{E}[C_{mv}(\mathbb{O}(\mathcal{S}_u(g)))] \geq 0$$

for any $g \in \mathbb{R}_{\geq 0}$, as the strategy can always quit at the very beginning.

²We use the term *profit* here to distinguish from the term *reward* (recall we get a unit of reward by reaching a target state).

Grade $\gamma_u(\mathcal{S})$: In particular, we define the *grade* $\gamma_u(\mathcal{S})$ of state u in Markov system \mathcal{S} as the unique value of g for which an optimal player is indifferent between the two possible first moves in the game $\mathcal{S}_u(g)$, i.e. he can either play \mathcal{S} for the first step or quit at the very beginning. We also use γ_u and $\mathcal{S}(g)$ as a shorthand for $\gamma_u(\mathcal{S})$ and $\mathcal{S}_u(g)$ respectively when \mathcal{S} and its current state are clear from the context.

To gain a bit more intuition about the grade, consider a pure strategy **ALG** for the game $\mathcal{S}(g)$. Note that a pure strategy can be defined by a subset $Q \subset V$ of states: the player chooses to play \mathcal{S} , until either target t is reached, or a state in Q is reached and the player chooses to quit, which ends the game immediately. Let event Ω be that the token in \mathcal{S} reaches t and let $C_{\mathcal{S}}$ be the (random) movement cost **ALG** spends on \mathcal{S} . It is easy to see that

$$\mathbb{E}[g \cdot R(\mathbf{ALG}) - C_{mv}(\mathbf{ALG})] = g \cdot \Pr[\Omega] - \mathbb{E}[C_{\mathcal{S}} \mid \Omega] \Pr[\Omega] - \mathbb{E}[C_{\mathcal{S}} \mid \neg\Omega] \Pr[\neg\Omega].$$

One can see it is linear in g for fixed set Q . Hence, the value of the game $\text{val}(\mathcal{S}(g))$ (as a function of g) is the maximum of a set of linear functions and is therefore a piece-wise linear convex function in g (See Figure 1). When g is very small, the optimal strategy should choose to quit immediately, and both the cost and the profit are zero. When g is very large, the optimal strategy should never quit before reaching the target. Hence, as we increase g gradually from 0, there is a point at which we are indifferent between playing \mathcal{S} and quitting immediately, which is the value of the grade for state u . We set $\gamma_t = 0$ for the target state t . Readers can refer to [DTW03] for more detailed discussion. We can also define a grade for a Markov system, that is the grade of its current state u , i.e., $\gamma_u(\mathcal{S})$. As shown in [DTW03], grades can be computed in poly-time (see Section 7 in [DTW03]).

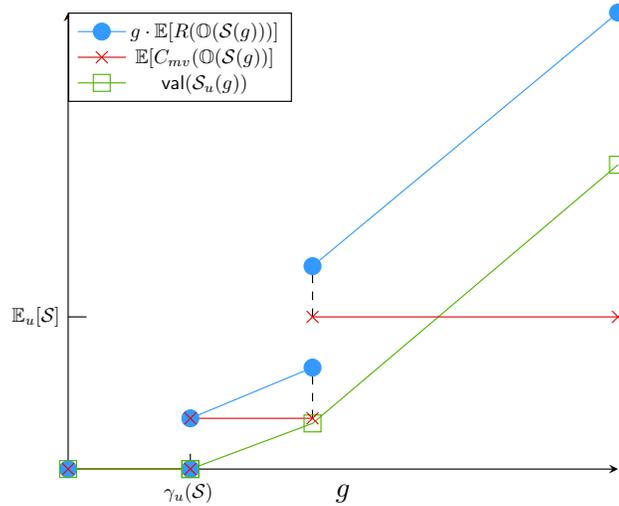


Figure 1: An illustration of $\text{val}(\mathcal{S}(g))$, which is a piecewise linear convex function in g . Each piece corresponds to a subset $Q \subset V$. The first piece corresponds to the empty set $Q = V$ (the player quits immediately), and the last piece corresponds to $Q = \emptyset$ (the player never quits before reaching the target). $\mathbb{E}_u[\mathcal{S}]$ is the expected cost of a never-quitting player. In fact, each turning point corresponds to the grade of some state, and the first one corresponds to $\gamma_u(\mathcal{S})$ where u is the current state. We also show the expected movement cost $\mathbb{E}[C_{mv}(\mathcal{O}(\mathcal{S}_u(g)))]$ and the expected profit $g \cdot \mathbb{E}[R(\mathcal{O}(\mathcal{S}_u(g)))]$.

By the definition of the grade, we can get the following lemma easily.

Lemma 6 (Optimal solution for $\mathcal{S}(g)$). *A strategy for $\mathcal{S}(g)$ is optimal if it chooses to advance \mathcal{S} whenever the grade is no more than g and it chooses to quit whenever the grade is larger than g . When the grade of \mathcal{S}*

equals to g , there is an optimal strategy that chooses to first advance \mathcal{S} and there is one that chooses to quit immediately.

We say a game is a *fair game* if the value of the game is zero. An observation that we use repeatedly is that the game $\mathcal{S}_u(\gamma_u(\mathcal{S}))$ is a fair game (from Figure 1, one can see that $\gamma_u(\mathcal{S})$ is the largest g such that the game $\mathcal{S}_u(g)$ is a fair game).

Now we define a *prevailing cost* [DTW03] and an epoch [GJSS19]. A trajectory is a sequence of states traversed by a player.

Definition 7 (Prevailing cost). *The prevailing cost of Markov system \mathcal{S} in a trajectory ω is $Y^{\max}(\omega) = \max_{u \in \omega} \gamma_u(\mathcal{S})$.*

In other word, the prevailing cost is the maximum grade at any point. In particular, the prevailing cost increases whenever the Markov system reaches a state with grade larger than each of the previously visited states. The prevailing cost can be viewed as a non-decreasing piece-wise constant function of time, which motivates the definition of epoch:

Definition 8 (Epoch). *An epoch for a trajectory ω is any maximal continuous segment of ω where the prevailing cost does not change.*

We also define an interesting teasing game introduced in [DTW03], which is useful later.

Definition 9 (Teasing game \mathcal{S}^T). *Consider the game $\mathcal{S}_s(\gamma_s)$ with initial state s . Whenever the player reaches a state u with grade $\gamma_u > \gamma_s$, we place γ_u units of profits at the target state t rather than γ_s . The objective is also to maximize the expectation of “profits - costs”. The γ_u profit provides just enough incentive for the player to continue advancing the \mathcal{S} . We denote the new teasing game by \mathcal{S}^T .*

For the new teasing game, we have the following lemma which also follows directly from the definition of the grade:

Lemma 10 (Fairness of \mathcal{S}^T , Lemma 5.3 in [DTW03]). *\mathcal{S}^T is a fair game, and a strategy for \mathcal{S}^T is optimal if and only if the player never quits in the intermediate of an epoch, and only quits at the beginning of an epoch.*

By the above lemma, one can easily see that the expected movement cost of a never-quitting player of the game \mathcal{S}^T is equal to the expected prevailing cost, where being never-quitting means that the player continues playing the system until system reaches the target state and he collects the profits.

3.2 The Doubling Technique.

In this subsection, we adopt the *doubling technique* which is similar to the ones used in related stochastic optimization problems such as [ENS18] and [JLLS20]. See the pseudo-code of the framework in Algorithm 1. Basically, the framework proceeds in phases, and in i th phase, we call a sub-procedure denoted by $\text{BudgetMG}(\mathcal{M}_{i-1}, k_{i-1}, B_i)$ in which we start with \mathcal{M}_{i-1} , the current state of all Markov system, and aim

at collecting the remaining k_{i-1} units of reward with total cost budget $B_i = O(1)\beta^i$ (in expectation).

Algorithm 1: A general algorithm Algo-MG

```

1 Input: The problem instance  $\mathcal{M}$ , objective number of rewards  $K$ 
2 Process:
3 Set  $\beta \in (1, 2)$ ,  $k_0 = K$ ,  $\mathcal{M}_0 = \mathcal{M}$ ,  $c = O(1)$ ;
4 for phase  $i = 1, 2, \dots$  do
5    $(\mathcal{M}_i, k_i) \leftarrow \text{BudgetMG}(\mathcal{M}_{i-1}, k_{i-1}, B_i = c\beta^i)$ ;
6   if  $k_i \leq 0$  then
7     Break
8   end
9 end

```

Recall that we have a unique root \mathbf{R} . In particular, in MG-Unit, we let $\text{BudgetMG}(\mathcal{M}_i, k_i, B_{i+1})$ begin to play the game at the Markov system where $\text{BudgetMG}(\mathcal{M}_{i-1}, k_{i-1}, B_i)$ halts, that's $\text{BudgetMG}(\mathcal{M}_i, k_i, B_{i+1})$ does not need pay the unit switching cost for the Markov system where $\text{BudgetMG}(\mathcal{M}_{i-1}, k_{i-1}, B_i)$ halts. In MG-Metric, we require that the strategy goes back to \mathbf{R} after $\text{BudgetMG}(\mathcal{M}_{i-1}, k_{i-1}, B_i)$ halts for simplicity of the analysis, and hence the next phase $\text{BudgetMG}(\mathcal{M}_i, k_i, B_{i+1})$ starts the game at the root \mathbf{R} . This blows up the total cost by at most a factor of 2 since switching back to the root \mathbf{R} costs at most B_i . The main reason of doing so is to avoid the case where $\text{BudgetMG}(\mathcal{M}_{i-1}, k_{i-1}, B_i)$ stops at a Markov system far-way from the other chains and $\text{BudgetMG}(\mathcal{M}_i, k_i, B_{i+1})$ has to pay a lot in the first switch (rather this switching cost is amortized to the $i - 1$ th phase).

To analyze the algorithm framework, we define $\odot(\mathcal{M}, k)$ to be the optimal strategy for the problem instance \mathcal{M} (the game starts from the root \mathbf{R}) with the target number of rewards k . Intuitively, for $j \geq i$, the expected cost of $\odot(\mathcal{M}_j, k_j)$ is no more than the one of $\odot(\mathcal{M}_i, k_i)$ as $k_j \leq k_i$, i.e. in some sense \mathcal{M}_j can get $k_i - k_j$ units of rewards for free. We can prove this formally and get the following lemma (whose proof can be found in Appendix A.2.1).

Lemma 11. *For any $j \geq i \geq 1$ and any Algorithm BudgetMG, one has*

$$\mathbb{E}[C_{tot}(\odot(\mathcal{M}_i, k_i))] \geq \mathbb{E}[C_{tot}(\odot(\mathcal{M}_j, k_j))].$$

Notice that the randomness is over the entire run of Algo-MG.

We use Lemma 11 to prove the following Lemma 12, which is used for both MG-Unit and MG-Metric.

Lemma 12 (Budgeted Subproblem). *We are given an MG-Metric (or MG-Unit) instance \mathcal{M} with positive integers k and $B \in \mathbb{R}_{\geq 0}$. For any $B > c_1 \mathbb{E}[C_{tot}(\odot(\mathcal{M}, k))]$, if there is an algorithm BudgetMG that can succeed in collecting k units of rewards with some constant probability (say more than 0.01) using expected total cost at most $c_2 B$, where c_1, c_2 are some universal constants, then there is a constant factor approximation algorithm for MG-Metric (or MG-Unit).*

The proof is similar to the previous ones [ENS18, JLLS20] with some subtle modifications and can be found in Appendix A.2.

If we can design BudgetMG which satisfies the precondition of Lemma 12, then by running BudgetMG in the framework, we get an $O(1)$ -approximation algorithm. All our effort goes into designing BudgetMG that satisfies the precondition of Lemma 12 in the following sections.

4 Markov Game with Unit Metric

In this section, we consider the MG-Unit problem (Definition 2). Our main result is Theorem 3.

Theorem 3. *There is a simple index strategy that can achieve a constant approximation ratio for the MG-Unit problem.*

Algorithm 2: Algorithm for MG-Unit

```
1 Input: The instance  $\mathcal{M}$ 
2 while we have not collected any reward do
3   | Choose to play Markov system  $\mathcal{S}_i$  with  $i = \arg \min_i \Gamma_i$ ;
4 end
```

Since our goal is to design an indexing strategy, we need to define an index which can incorporate the information of switching cost. In particular, for each Markov system $\mathcal{S}_i \in \mathcal{S}$ and every state u of \mathcal{S}_i , we create a *dummy state* u' for u with unit movement cost ($C_{u'} = 1$) and deterministic transition to u ($P_{u',u} = 1$). The dummy states are used to capture the unit switching cost. Let γ'_u denote the grade of u 's dummy state, and call γ'_u the *dummy grade* of u . We say a Markov system \mathcal{S}_i is *active* at a time step t if \mathcal{S}_i is played in the previous step (i.e., continuing to play \mathcal{S}_i in the t -th step does not incur any switching cost), and *inactive* otherwise.

Initially, we are at the root \mathbf{R} , and all of the Markov system are inactive. Now we define a grade Γ_i for \mathcal{S}_i which is at state u : if \mathcal{S}_i is active, then its Γ_i is defined to be the grade γ_u of u ; otherwise, Γ_i is defined to be the dummy grade of u , which we denote by γ'_u .

Our strategy for MG-Unit is simply choosing to play the \mathcal{S}_i with the smallest grade Γ_i , breaking ties arbitrarily. See also Algorithm 2.

Although Algorithm 2 is simple to state, directly analyzing it seems difficult. Rather, we analyze an alternative algorithm via the doubling technique framework. More specifically, we apply the framework Algo-MG (Algorithm 1), which proceeds in phases. In each phase, it calls the sub-procedure BudgetMG-Unit (Algorithm 3). In the sub-procedure, we have a movement cost budget and a switching cost budget. We repeatedly play the Markov system with the smallest grade, until one of the budgets is exhausted or all Markov system reach their target states. Note that we may not stop immediately when we reach a target state and collect one unit of reward. Instead, we should remove the present system and keep on playing (if there are still budget and available systems in this phase). Hence, the cost of the alternative algorithm is no less than that of Algorithm 2. See also Appendix B.1.

4.1 Analysis.

Recall that in order to prove the main result of this section (i.e. Theorem 3), it suffices to show that the sub-procedure BudgetMG-Unit (Algorithm 3) satisfies the precondition of Lemma 12. The key is to prove the following lemma, which is the precondition of Lemma 12 specialized for MG-Unit. Recall that we are still aiming at solving MG-Unit where Algorithm 2 stops immediately whenever it makes one Markov system reach its target state. The alternative Algorithm 3 may collect more than one unit of reward and it is only used for analysis and setting an upper bound for the expected cost of our true algorithm (Algorithm 2).

Algorithm 3: Subprocedure BudgetMG-Unit

- 1 **Input:** The instance \mathcal{M} , budget $2^8 B$
 - 2 Set movement budget $2^7 B$ and switching budget $2^7 B$;
 - 3 Set $k \leftarrow 1$;
 - 4 **while** *there are available Markov system and condition \mathcal{A} holds* **do**
 - 5 Choose to play Markov system \mathcal{S}_i with $i = \arg \min_i \Gamma_i$;
 - 6 **if** *We reach a target state t in \mathcal{S}_i* **then**
 - 7 Collect one unit reward and mark \mathcal{S}_i as *unavailable*;
 - 8 $k \leftarrow k - 1$;
 - 9 **end**
 - 10 **end**
 - 11 **Return:** The updated instance \mathcal{M} , the remaining number of target states k ;
 - 12 **Define:** Condition \mathcal{A} :
 - 13 The next move does not make the total movement cost or switching cost exceed $2^7 B$;
-

Lemma 13. *For any input \mathcal{M} , let \circledast be the optimal strategy for this instance. If $B \geq 10\mathbb{E}[C_{tot}(\circledast)]$, with probability at least $1/20$, BudgetMG-Unit (Algorithm 3) can collect at least one unit reward with budget $2^8 \cdot B$.*

Proof of Lemma 13. Due to Condition \mathcal{A} (Line 12), the total cost budget cannot be violated. Hence, it suffices to prove the success probability is at least $1/20$.

We only need to consider the case when there are more than one system, otherwise it is optimal to switch the only system and make it reach the target state. Suppose BudgetMG-Unit decides to play some Markov system \mathcal{S}_{i_j} after switching j times, i.e., the sequence of chains played by BudgetMG-Unit is $(\mathbf{R}, \mathcal{S}_{i_1}, \mathcal{S}_{i_2}, \dots, \mathcal{S}_{i_j}, \dots)$. Note that BudgetMG-Unit may revisit a chain (i.e., $\mathcal{S}_{i_t} = \mathcal{S}_{i_j}$ for some $i_t \neq i_j$). Let ω_j be the path (a sequence of states) BudgetMG-Unit traverses on \mathcal{S}_{i_j} after it switches to \mathcal{S}_{i_j} and before it switches to $\mathcal{S}_{i_{j+1}}$ (or stops due to running out budget). We define the stopping time τ as the numbers of switching of BudgetMG-Unit when it halts (equivalently, the switching cost of BudgetMG-Unit). For $\tau + 1 \leq j \leq 2^7 B$, we let $\omega_j = \emptyset$. Let $\omega = \cup_{j \geq 0}^{2^7 B} \omega_j = \cup_{j \geq 0}^{\tau} \omega_j$ be the whole trajectory traversed by BudgetMG-Unit and let Ω denote the set of all possible trajectories BudgetMG-Unit can traverse. We use the notation $\omega_{[0:j]} = \cup_{t=0}^j \omega_t$ to denote the prefix of ω .

Now we define the Boolean random variables X_j for $j \leq \tau$: if BudgetMG-Unit can get the reward in \mathcal{S}_{i_j} , then $X_j = 1$. Otherwise (i.e., BudgetMG-Unit switches out or runs out the budget) $X_j = 0$. If $j > \tau$, we let random variables $X_j = 0$. Similarly, for $j \leq \tau$, we define random variable C_j to represent the movement cost BudgetMG-Unit spends on \mathcal{S}_{i_j} (C_j does not include the unit switching cost).

Let Γ_{i_j} be the smallest grade of all systems when the j -th switch occurs, which is the dummy grade of the current state of \mathcal{S}_{i_j} . We use s_{i_j} to denote the current state of \mathcal{S}_{i_j} when the j -th switch occurs, and s'_{i_j} be its dummy state. Hence, $\Gamma_{i_j} = \gamma_{s'_{i_j}}(\mathcal{S}_{i_j})$. By the greedy process of BudgetMG-Unit, we know that

$$\Gamma_{i_1} \leq \Gamma_{i_2} \leq \dots \leq \Gamma_{i_j} \leq \dots \leq \Gamma_{i_\tau}.$$

For a trajectory ω , let $\mathcal{T}_j(\omega) = \mathbb{E}[X_j \mid \omega_{[0:j-1]}]$. Indeed, one can see that \mathcal{T}_j is a random variable with randomness from $\omega_{[0:j-1]}$. Let $\mathcal{T} = \sum_{j=1}^{\tau} \mathcal{T}_j$ and $\mathcal{T}(\omega) = \sum_{j=1}^{\tau} \mathcal{T}_j(\omega) = \sum_{j=1}^{\tau} \mathbb{E}[X_j \mid \omega_{[0:j-1]}]$.

Claim 14. We know that

$$\Gamma_{i_{j+1}} \geq \frac{\mathbb{E}[C_j + 1 \mid \omega_{[0:j-1]}]}{\mathcal{T}_j(\omega)} = \frac{\mathbb{E}[C_j + 1 \mid \omega_{[0:j-1]}]}{\mathbb{E}[X_j \mid \omega_{[0:j-1]}]} \geq \Gamma_{i_j}.$$

Proof. Recall that the algorithm plays \mathcal{S}_{i_j} whose grade is Γ_{i_j} when the j -th switch occurs. Note that we know the value of $\Gamma_{i_{j+1}}$ when the j -th switch occurs. Denote $U_j = \{u \in V_{i_j} \mid \gamma_u \leq \Gamma_{i_{j+1}}\}$. BudgetMG-Unit continues playing the system \mathcal{S}_{i_j} until it reaches the target state or reach a state outside U_j when the cost budget is not exhausted.

We know that $\mathcal{S}_{s'_{i_j}}(\Gamma_{i_j})$ is a fair game where s'_{i_j} is the dummy state of s_{i_j} . On one hand, by Lemma 6, the strategy determined by U_j may not be the optimal solution to the game $\mathcal{S}_{s'_{i_j}}(\Gamma_{i_j})$, so we have

$$\mathbb{E}[X_j \mid \omega_{[0:j-1]}] \cdot \Gamma_{i_j} - \mathbb{E}[C_j + 1 \mid \omega_{[0:j-1]}] \leq \text{val}(\mathcal{S}_{s'_{i_j}}(\Gamma_{i_j})) = 0.$$

On the other hand, as $\Gamma_{i_{j+1}} \geq \Gamma_{i_j}$, the game $\mathcal{S}_{s'_{i_j}}(\Gamma_{i_{j+1}})$ is better than fair, hence $\text{val}(\mathcal{S}_{s'_{i_j}}(\Gamma_{i_{j+1}})) \geq 0$. The strategy determined by U_j is exactly the optimal solution to the game $\mathcal{S}_{s'_{i_j}}(\Gamma_{i_{j+1}})$ according to Lemma 6. Thus we can conclude

$$\mathbb{E}[X_j \mid \omega_{[0:j-1]}] \cdot \Gamma_{i_{j+1}} - \mathbb{E}[C_j + 1 \mid \omega_{[0:j-1]}] = \text{val}(\mathcal{S}_{s'_{i_j}}(\Gamma_{i_{j+1}})) \geq 0.$$

This finishes the proof of the claim. \square

Recall that it suffices to show the success probability of BudgetMG-Unit is at least $1/20$ when $B \geq 10\mathbb{E}[C_{tot}(\odot)]$. The following key lemma states that with probability at least $3/20$ that either BudgetMG-Unit succeeds to collect at least one unit of reward, or the summation of the conditional expectation is large.

Lemma 15. (Key Lemma) With probability at least $3/20$, either $\sum_{j=1}^{\tau} \mathcal{T}_j \geq 2$ or $\sum_{j=1}^{\tau} X_j \geq 1$. Equivalently, one has

$$\Pr \left[\sum_{j=1}^{\tau} \mathcal{T}_j \geq 2 \vee \sum_{j=1}^{\tau} X_j \geq 1 \right] \geq 3/20.$$

Proof. Recall that $\mathcal{T}_j = \mathbb{E}[X_j \mid \mathcal{F}_{j-1}]$, $\mathcal{T} = \sum_{j=1}^{\tau} \mathcal{T}_j$, and Ω is the set of all possible trajectories. To prove the statement, we bound the probability of a few bad cases:

Case (1): [$\mathcal{T} < 2$, $\sum_{j=1}^{\tau} X_j = 0$ and the switching budget runs out first.]

Let $\Omega_1 = \{\omega \in \Omega \mid \mathcal{T}(\omega) = \sum_{j=1}^{\tau} \mathcal{T}_j(\omega) < 2, \sum_{j=1}^{\tau} X_j(\omega) = 0, \tau = 2^7 \cdot B\}$ be the set of trajectories corresponding to Case (1).

Conditioning on any sample path (trajectory) $\omega \in \Omega_1$, there must exist $j \in [2^7 \cdot B]$ such that $\mathcal{T}_j(\omega) \leq \frac{2}{2^7} = 2^{-6} \cdot B^{-1}$. By Claim 14, we know that

$$\Gamma_{i_{j+1}} \geq \frac{\mathbb{E}[C_j + 1 \mid \omega_{[0:j-1]}]}{\mathcal{T}_j(\omega)} \geq \frac{1}{\mathcal{T}_j(\omega)} \geq 2^6 \cdot B.$$

Recall that \mathcal{M} is the initial instance accepted by BudgetMG-Unit, and suppose the \mathcal{M}' is the resulting instance at the end of BudgetMG-Unit (their difference can be determined by ω).

Conditioning on ω , we let \mathbb{O}' be the optimal strategy for instance \mathcal{M}' . If BudgetMG-Unit does not collect any reward, we still need to collect one unit reward from \mathcal{M}' . If BudgetMG-Unit succeeds to collect at least one unit of reward, \mathbb{O}' does nothing and halts immediately. Note that $\mathbb{E}[C_{tot}(\mathbb{O}') \mid \omega]$ is a random variable with randomness from ω .

By Lemma 11, we have

$$\mathbb{E}[C_{tot}(\mathbb{O}')] \leq \mathbb{E}[C_{tot}(\mathbb{O})] \leq B/10.$$

Recall that the smallest grade of all available Markov system does not decrease and thus grades of all \mathbb{S} are at least $2^6 \cdot B$ when BudgetMG-Unit halts. We want to bound the probability $\Pr[C_{tot}(\mathbb{O}') \leq B]$. We need the following claim:

Claim 16. *For any positive real number B and any MG-Metric instance \mathcal{M} , if the grades of all Markov system are at least $\zeta \cdot B$ for some constant $\zeta \geq 1$, then for any strategy ALG with total cost budget B , one has*

$$\Pr[\text{ALG succeeds to collect at least one unit reward}] \leq 1/\zeta.$$

Proof. With the input \mathcal{M} , let $\Gamma(\mathcal{S}_i)$ be the grade of the inactive \mathcal{S}_i and u'_i be the dummy state of the initial state u_i of \mathcal{S}_i . We first simplify the game by only requiring to pay the unit switching cost once for each Markov system. Equivalently, only the initial state u_j has its corresponding dummy state u'_j in \mathcal{S}_j , and when the player decides to switch to \mathcal{S}_j again, she switches to the current normal state (instead of the dummy state).³ Consider an arbitrary strategy ALG . Obviously, the cost of ALG is only less in this simplified game (trajectory-wise).

Now, imagine that we run ALG on a composition game $\mathcal{G} = \mathcal{S}_{u'_1}(\Gamma(\mathcal{S}_1)) \circ \mathcal{S}_{u'_2}(\Gamma(\mathcal{S}_2)) \circ \dots \circ \mathcal{S}_{u'_m}(\Gamma(\mathcal{S}_m))$: The new composition game has the same set of Markov chains; hence a trajectory in the original game is also a trajectory in the new game. We know $\mathcal{S}_{u'_i}(\Gamma(\mathcal{S}_i))$ is a fair game and any combination (simultaneous, sequential, or interleaved) of independent fair games is still fair (e.g., Lemma 5.4 in [DTW03]). Hence, \mathcal{G} is a fair game.

We use p_j to denote the probability that ALG makes system \mathcal{S}_j reach its target state and C_j be the cost that ALG spends on \mathcal{S}_j in \mathcal{G} . Indeed, one can see C_j is also the cost ALG spends on \mathcal{S}_j in the simplified game. As the composition game is fair, thus we know the expected “profits-cost” of ALG is at most 0. In particular, one has $\sum_{j=1}^m \Gamma(\mathcal{S}_j)p_j - \sum_{j=1}^m \mathbb{E}[C_j] \leq \text{val}(\mathcal{G}) = 0$.

For each $j \in [m]$, one has $\Gamma(\mathcal{S}_j) \geq \zeta \cdot B$. As the budget is B , $\sum_{j=1}^m \mathbb{E}[C_j] \leq B$. Hence, we have that

$$\Pr[\text{ALG succeeds to collect at least one unit reward}] \leq \sum_{j=1}^m p_j \leq \frac{B}{\min_j \Gamma(\mathcal{S}_j)} \leq \frac{B}{\zeta \cdot B} = 1/\zeta.$$

This finishes the proof of the claim. □

Note that the randomness of \mathbb{O}' comes from ω . Applying Claim 16 to \mathbb{O}' with total cost budget B , we can see

$$\Pr[\mathbb{O}' \text{ succeeds to collect at least one unit reward} \wedge C_{tot}(\mathbb{O}') \leq B \mid \omega \in \Omega_1] \leq 2^{-6},$$

which means that $\mathbb{E}[C_{tot}(\mathbb{O}') \mid \omega \in \Omega_1] \geq B(1-2^{-6}) \geq 0.9B$. Combining the fact that $\mathbb{E}[C_{tot}(\mathbb{O}')] \leq B/10$, we know that

$$\Pr[\omega \in \Omega_1] \leq 1/9.$$

³Hence, the simplified game can be reduced to a Markov game without switching cost, hence solved optimally.

Case (2): $\mathcal{T} < 2$, $\sum_{j=1}^{\tau} X_j = 0$ and the movement budget runs out first. We define some additional variables to analyze this case. Let $\xi(\omega)$ be the movement cost of the (next) move which breaks Condition \mathcal{A} and makes BudgetMG-Unit halt. We divide Case (2) further into the following cases.

Sub-case (2.1):[Large breaking cost.] This Sub-case corresponds to the set $\Omega_{2.1} = \{\omega \in \Omega : \mathcal{T}(\omega) < 2, \sum_{j=1}^{\tau} X_j(\omega) = 0, \sum_{j=1}^{\tau} C_j(\omega) + \xi(\omega) \geq 2^7 \cdot B, \xi(\omega) \geq 2^6 \cdot B\}$. As $\xi(\omega) \geq 2^6 \cdot B$, conditioning on $\omega \in \Omega_{2.1}$, it is impossible for \odot' to make \mathcal{S}_{τ} reach the target state within budget B . Here we say \odot' does something within budget B , means the event that \odot' does something and its total cost is no more than B at the time of accomplishment.

For those $\omega \in \Omega_{2.1}$, as $\xi(\omega) \geq 2^6 \cdot B$, we know BudgetMG-Unit must have decided to play some Markov system whose grade is at least $2^6 \cdot B$. More specifically, let u be final current state of \mathcal{S}_{τ} when BudgetMG-Unit halts, and by the definition we have $C_u = \xi(\omega)$. By the definition, we know that $\gamma_u \geq C_u \geq 2^6 \cdot B$. We define $\mathcal{S}_{i_{\tau+1}}$ to be the (inactive) system with the second smallest grade when the τ -th switch occurs ($\mathcal{S}_{i_{\tau}}$ has the smallest grade then), and $\Gamma_{i_{\tau+1}}$ is the corresponding grade. One can see that $\Gamma_{i_{\tau+1}} \geq \gamma_u \geq 2^6 \cdot B$.

Similar to **Case (1)**, by Claim 16, one has

$$\Pr[\odot' \text{ succeeds to collect at least one unit reward within Budget } B \mid \omega \in \Omega_{2.1}] \leq 2^{-6},$$

which implies that $\Pr[\omega \in \Omega_{2.1}] \leq 1/9$ by the fact that $\mathbb{E}[C_{tot}(\odot')] \leq B/10$.

Sub-case (2.2):[Small expected movement cost, small breaking cost.] Particularly, this Sub-case corresponds to the set $\Omega_{2.1} = \{\omega \in \Omega : \mathcal{T}(\omega) < 2, \sum_{j=1}^{\tau} X_j(\omega) = 0, \sum_{j=1}^{\tau} \mathbb{E}[C_j \mid \omega_{[0:j-1]}] \leq 2^4 \cdot B, \sum_{j=1}^{\tau} C_j(\omega) + \xi(\omega) > 2^7 \cdot B, \xi(\omega) \leq 2^6 \cdot B\}$. Obviously, if $\sum_{j=1}^{\tau} C_j(\omega) + \xi(\omega) > 2^7 \cdot B$ and $\xi(\omega) \leq 2^6 \cdot B$, then we know $\sum_{j=1}^{\tau} C_j(\omega) \geq 2^6 \cdot B$. Thus the probability of this Sub-case can be bounded by Markov Inequality. In particular, one has

$$\begin{aligned} \Pr[\omega \in \Omega_{2.1}] &\leq \Pr \left[\omega \in \Omega : \sum_{j=1}^{\tau} C_j(\omega) \geq 2^6 \cdot B \wedge \sum_{j=1}^{\tau} \mathbb{E}[C_j \mid \omega_{[0:j-1]}] \leq 2^4 \cdot B \right] \\ &\leq \Pr \left[\sum_{j=1}^{\tau} C_j(\omega) \geq 2^6 \cdot B \mid \sum_{j=1}^{\tau} \mathbb{E}[C_j \mid \omega_{[0:j-1]}] \leq 2^4 \cdot B \right] \leq 2^{-2}, \end{aligned}$$

where the last step follows by applying Markov Inequality on the probability space $\{\omega \in \Omega : \sum_{j=1}^{\tau} \mathbb{E}[C_j \mid \omega_{[0:j-1]}] \leq 2^4 \cdot B\}$.

Sub-case (2.3):[Large expected movement cost, small breaking cost.] The corresponding sample path set is $\Omega_{2.3} = \{\omega \in \Omega : \mathcal{T}(\omega) < 2, \sum_{j=1}^{\tau} X_j(\omega) = 0, \sum_{j=1}^{\tau} \mathbb{E}[C_j \mid \omega_{[0:j-1]}] > 2^4 \cdot B, \sum_{j=1}^{\tau} C_j(\omega) + \xi(\omega) > 2^7 \cdot B, \xi(\omega) \leq 2^6 \cdot B\}$.

Let \mathcal{M}'' be the instance when the τ -th switch occurs, k'' be the remaining number of rewards we need to collect, and let \odot'' be the optimal strategy for \mathcal{M}'' with the objective k'' . By the same proof of Lemma 11 (trajectory-wise), we also have

$$\mathbb{E}[C_{tot}(\odot'')] \leq \mathbb{E}[C_{tot}(\odot)] \leq B/10.$$

Conditioning on any sample path $\omega \in \Omega_{2.3}$, we know that $\mathcal{T}(\omega) = \sum_{j=1}^{\tau} \mathbb{E}[X_j \mid \omega_{[0:j-1]}] \leq 2$ and $\sum_{j=1}^{\tau} \mathbb{E}[C_j + 1 \mid \omega_{[0:j-1]}] \geq 2^4 \cdot B$. Then there exists $j \leq \tau$ such that $\frac{\mathbb{E}[C_j + 1 \mid \omega_{[0:j-1]}]}{\mathbb{E}[X_j \mid \omega_{[0:j-1]}]} \geq 2^3 \cdot B$ and by Claim 14, we know $\Gamma_{i_{j+1}} \geq 2^3 \cdot B$.

Denote the subset $\Omega_{2.3.1} = \{\omega \in \Omega_{2.3} : \Gamma_{i_\tau} \geq 2^3 \cdot B\}$. We know that

$$\Pr[\textcircled{O}'' \text{ succeeds to collect at least one unit reward within Budget } B \mid \omega \in \Omega_{2.3.1}] \leq 2^{-3}.$$

by Claim 16. Hence, by the fact that $B/10 \geq \mathbb{E}[C_{tot}(\textcircled{O}'')] \geq \mathbb{E}[C_{tot}(\textcircled{O}'') \mid \omega \in \Omega_{2.3.1}] \Pr[\Omega_{2.3.1}] \geq B(1 - 2^{-3}) \Pr[\Omega_{2.3.1}]$, we can see that $\Pr[\Omega_{2.3.1}] \leq 4/35$.

Otherwise, consider those $\omega \in \Omega_{2.3} \setminus \Omega_{2.3.1}$ such that $\Gamma_{i_\tau} < 2^3 \cdot B$. By the greedy property of the algorithm, we know that $\Gamma_{i_j} \leq \Gamma_{i_\tau}$ for all $1 \leq j \leq \tau$. By Claim 14, we know that $2^3 \cdot B \geq \Gamma_{i_{j+1}} \geq \frac{\mathbb{E}[C_{j+1} \mid \omega_{[0:j-1]}]}{\mathbb{E}[X_j \mid \omega_{[0:j-1]}]}$ for $1 \leq j \leq \tau - 1$. As $\mathcal{T}(\omega) = \sum_{j=1}^{\tau} \mathbb{E}[X_j \mid \omega_{[0:j-1]}] \leq 2$, then we know that $\sum_{j=1}^{\tau-1} \mathbb{E}[C_j + 1 \mid \omega_{[0:j-1]}] \leq 2^4 \cdot B$. However, $\sum_{j=1}^{\tau} C_j(\omega) > 2^6 B$ for this case. Use the same Markov inequality argument as in Sub-case (2.2), one has

$$\begin{aligned} \Pr[\omega \in \Omega_{2.3} \setminus \Omega_{2.3.1}] &\leq \Pr \left[\omega \in \Omega : \sum_{j=1}^{\tau-1} C_j(\omega) \geq 2^6 \cdot B \wedge \sum_{j=1}^{\tau-1} \mathbb{E}[C_j \mid \omega_{[0:j-1]}] \leq 2^4 \cdot B \right] \\ &\leq \Pr \left[\sum_{j=1}^{\tau-1} C_j(\omega) \geq 2^6 \cdot B \mid \sum_{j=1}^{\tau-1} \mathbb{E}[C_j \mid \omega_{[0:j-1]}] \leq 2^4 \cdot B \right] \leq 2^{-2}. \end{aligned}$$

Union Bound: By union bound over **Case (1)** and **Case (2)**, we have that

$$\begin{aligned} \Pr[\omega \in \Omega : \text{BudgetMG-Unit fails and } T(\omega) < 2] &= \Pr[\Omega_1 \cup \Omega_{2.1} \cup \Omega_{2.2} \cup \Omega_{2.3}] \\ &\leq \Pr[\Omega_1] + \Pr[\Omega_{2.1}] + \Pr[\Omega_{2.2}] + \Pr[\Omega_{2.3}] \\ &\leq 1/9 + 1/4 + 1/9 + \Pr[\Omega_{2.3} \setminus \Omega_{2.3.1}] + \Pr[\Omega_{2.3.1}] \\ &\leq 1/9 + 1/4 + 1/9 + 1/4 + 4/35 \\ &< 17/20. \end{aligned}$$

This completes the proof. \square

The lemma below complements Lemma 15 and shows that if the summation of conditional expectation is large, then the algorithm should succeed with a constant probability.

Lemma 17. *Suppose X_1, X_2, \dots, X_n are a sequence of random variables taking values in $\{0, 1\}$, and $\mathcal{F}_j = \sigma(X_1, \dots, X_j)$ is the filtration defined by the sequence. Given any real number μ , if $\sum_{j=1}^n \mathbb{E}[X_j \mid \mathcal{F}_{j-1}] \geq \mu$, then*

$$\Pr \left[\sum_{j=1}^n X_j \geq 1 \right] \geq 1 - e^{-3\mu/8}.$$

The proof of Lemma 17 can be found in Appendix B.2. By Lemma 15, one has

$$\Pr \left[\sum_{j=1}^{2^7 B} X_j \geq 1 \right] + \Pr[T \geq 2] \geq \Pr \left[\sum_{j=1}^{2^7 B} X_j \geq 1 \vee T \geq 2 \right] \geq 3/20. \quad (1)$$

By Lemma 17, one has

$$\Pr \left[\sum_{j=1}^{2^7 B} X_j \geq 1 \mid T \geq 2 \right] \geq 1 - e^{-3/4} \geq 1/2. \quad (2)$$

By combining Equation 1 and Equation 2, one can easily show that

$$\begin{aligned} \Pr \left[\sum_{j=1}^{2^7 B} X_j \geq 1 \right] &\geq \Pr \left[\sum_{j=1}^{2^7 B} X_j \geq 1 \mid T \geq 2 \right] \cdot \Pr[T \geq 2] \geq (1 - e^{-3/4}) \cdot \Pr[T \geq 2] \\ &\geq (1 - e^{-3/4}) \cdot \left(\frac{3}{20} - \Pr \left[\sum_{j=1}^{2^7 B} X_j \geq 1 \right] \right) \geq \frac{1}{2} \cdot \left(\frac{3}{20} - \Pr \left[\sum_{j=1}^{2^7 B} X_j \geq 1 \right] \right), \end{aligned}$$

which implies that

$$\Pr \left[\sum_{j=1}^{2^7 B} X_j \geq 1 \right] \geq 1/20.$$

Thus we complete the proof of Lemma 13. \square

Theorem 3 follows from Lemma 12 and Lemma 13 directly.

The simple index-based strategy achieves a constant approximation factor in unit metric space. However, there is a simple counter-example on general metric, which shows BudgetMG-Unit may perform arbitrarily bad. See Appendix B.3 for the example. This suggests that new techniques are needed for more general metric. This is the focus of the next section.

5 Markov Game with General Metric

In this section, we consider MG-Metric (Definition 4). Note that our requirement is also more general: we need to collect K units of rewards (make K chains reach their targets). Our main result is an efficient strategy MG-Metric which can achieve a constant factor approximation.

MG-Metric (See Algorithm 6 in Appendix) adopts the same doubling framework as Algo-MG (Algorithm 1), except that we need a more complicated sub-procedure BudgetMG-Metric (which approximates the budgeted version). In particular, BudgetMG-Metric should satisfy the precondition of Lemma 12. i.e., the expected total cost of BudgetMG-Metric is bounded by $c_2 \mathbb{E}[C_{tot}(\odot(\mathcal{M}, K))]$ for some universal constant c_2 , and it can collect K units of rewards with constant probability when B is large enough. For clarity, we present the lemma below, which is the precondition of Lemma 12 specialized for MG-Metric.

Lemma 18. *For any input \mathcal{M} with the objective number of rewards K , let \odot be the optimal strategy for this instance. If $B \geq 10 \mathbb{E}[C_{tot}(\odot)]$, with probability at least 0.1, BudgetMG-Metric (Algorithm 4) can collect at least K units of rewards (i.e. make at least K system reach their targets) with budget $O(1)B$ in expectation.*

It suffices to prove that sub-procedure BudgetMG-Metric (Algorithm 4) satisfies this lemma, which is the main task of this section.

5.1 Stochastic k -TSP.

The sub-procedure BudgetMG-Metric makes use of an $O(1)$ -approximation for the Stochastic- k -TSP problem [ENS18, JLLS20], defined as follows.

Definition 19 (Stochastic- k -TSP). We are given a metric $\mathcal{M} = (\mathcal{S}, d)$ with a root $\in \mathcal{S}$ and each vertex $v \in \mathcal{S}$ has an independent stochastic selection cost $C_{sl}(v) \in \mathbb{R}_{\geq 0}$. All cost distributions are given as input but the actual cost instantiation $C_{sl}(v)$ is only known after vertex v is visited. Suppose a vertex v can only be selected if v is visited.⁴ The goal is to adaptively find a tour T originating from the root and select a set S of k visited vertices while minimizing the expected total cost, which is the sum of the length of T and the cost of the selected vertices:

$$\mathbb{E} \left[\text{Length}(T) + \sum_{v \in S} C_{sl}(v) \right].$$

Theorem 20 (Theorem 2 in [JLLS20]). There is a non-adaptive constant factor approximation algorithm ALG_{ktsp} for Stochastic- k -TSP.

Remark 21. It is important to note that the strategy in [JLLS20] is non-adaptive. Here a non-adaptive strategy is an ordering Π of all vertices and Π is independent of the realization of the costs. Note that the strategy visits the vertices according to the ordering Π and may stop and choose k visited vertices before visiting all vertices by some criterion, depending on the realization of the costs. We use $\alpha_{ktsp} = O(1)$ to be the constant approximation ratio of this algorithm.

As mentioned, ALG_{ktsp} is not only a sub-procedure which can output an ordering of the chains (vertices). In fact there is also some probing and selection process after outputting the ordering in ALG_{ktsp} . We let $C_{sl}(\text{ALG}_{ktsp})$ and $C_{sw}(\text{ALG}_{ktsp})$ be the (random) selection cost and switching cost of ALG_{ktsp} respectively.

5.2 Reduction to Stochastic- k -TSP

5.2.1 A Fair Game MG-Metric-Fair.

We define another new game which is closely related to the teasing game \mathcal{S}^T (see definition 9) and plays a key role in the following proof.

Definition 22 (MG-Metric-Fair). We are given a finite metric space $\mathcal{M} = (\mathcal{S} \cup \{\mathbf{R}\}, d)$ (there is no additional assumption on metric d). Each node $\mathcal{S}_i \in \mathcal{S}$ is identified with a Markov system $\mathcal{S}_i = \langle V_i, P_i, C_i, s_i, t_i \rangle$. Similarly, at the beginning of the game, the player is at the root \mathbf{R} , and needs to pay the switching cost $d(\mathbf{R}, \mathcal{S}_i)$ if he wants to play Markov system \mathcal{S}_i . Switching from \mathcal{S}_i to \mathcal{S}_j incurs a switching cost of $d(\mathcal{S}_i, \mathcal{S}_j)$. Let $V_i^F \subseteq V_i$ denote the subset of states \mathcal{S}_i has been transited to by the player during the game. If $t_i \in V_i^F$, he can pay the prevailing cost as the fair movement cost and get one reward from \mathcal{S}_i . The objective is to adaptively collect at least K units of rewards (making at least K Markov system reach their targets), while minimizing the expected total cost (fair movement cost plus switching cost).

Let $C_{mv}^F(\mathbb{P})$, $C_{sw}(\mathbb{P})$ and $C_{tot}^F(\mathbb{P})$ be the (random) fair movement cost, switching cost and total cost of any strategy \mathbb{P} under the rule of MG-Metric-Fair respectively. Let \mathcal{O}_{fair} be the optimal strategy to MG-Metric-Fair. We have the following claim. The proof can be found in Appendix C.4.

Claim 23. The following inequality holds: $\mathbb{E}[C_{tot}^F(\mathbb{P})] = \mathbb{E}[C_{sw}(\mathcal{O}_{fair}) + C_{mv}^F(\mathcal{O}_{fair})] \leq \mathbb{E}[C_{tot}(\mathcal{O})]$.

⁴ In the original version in [JLLS20], there is one more condition that we are currently at vertex v to select it. Up to a factor of 2, our version of the problem is equivalent to the original version.

5.2.2 The reduction.

At a high level, BudgetMG-Metric reduces an instance of MG-Metric to an instance of Stochastic- k -TSP which has the same metric. A Markov chain \mathcal{S} in MG-Metric is replaced by a random variable, i.e. the random selection cost. More specifically, for each $\mathcal{S} = \langle V, P, C, s, t \rangle$, we consider the teasing game \mathcal{S}^T . Let $C_{sl}(\mathcal{S})$ be the (random) prevailing cost of a non-quitting strategy where being non-quitting means that the player continues to play \mathcal{S}^T until it reaches the target state. Then we treat \mathcal{S} as a vertex in Stochastic- k -TSP and let its selection cost distributed as $C_{sl}(\mathcal{S})$.

We note that the distribution of $C_{sl}(\mathcal{S})$ (in particular, $\Pr[C_{sl}(\mathcal{S}) \leq B]$ for any real number B) can be computed efficiently using the technique for computing grade (see Lemma 37 in Appendix C.2).

We establish a simple relation between the optimal cost of Stochastic- k -TSP and optimal cost of our original problem MG-Metric. We let $\mathbb{E}[C_{sl}(\mathbb{P})]$, $\mathbb{E}[C_{sw}(\mathbb{P})]$, $\mathbb{E}[C_{ktsp}(\mathbb{P})]$ be the expected selection cost, expected switching cost and expected total cost of any strategy \mathbb{P} for the Stochastic- k -TSP problem, respectively. Let \mathbb{O}_{ktsp} be the optimal strategy to Stochastic- k -TSP. Note that we use the same notation $C_{sw}(\cdot)$ to represent the switching cost for all games we have defined since the switching costs are counted in the same way.

Given any instance \mathcal{M} , let \mathbb{O} be the optimal strategies for this instance \mathcal{M} for our original problem MG-Metric. Now, we claim that the optimal cost of Stochastic- k -TSP is no more than the optimal cost of our original problem. The proof of Claim 24 is via the fair game we introduced in Section 5.2.1 and can be found in Appendix C.4.

Claim 24. *With the notations defined above, it holds that*

$$\mathbb{E}[C_{ktsp}(\mathbb{O}_{ktsp})] = \mathbb{E}[C_{sw}(\mathbb{O}_{ktsp}) + C_{sl}(\mathbb{O}_{ktsp})] \leq \mathbb{E}[C_{tot}(\mathbb{O})].$$

Recall we are under the assumption that $\mathbb{E}[C_{tot}(\mathbb{O})] \leq B/10$. Then we have

$$\begin{aligned} \mathbb{E}[C_{ktsp}(\text{ALG}_{ktsp})] &\leq \alpha_{ktsp} \mathbb{E}[C_{ktsp}(\mathbb{O}_{ktsp})] \\ &\leq \alpha_{ktsp} \mathbb{E}[C_{tot}(\mathbb{O})] \\ &\leq \alpha_{ktsp} B/10, \end{aligned}$$

where the first line is by Theorem 20 and the second line follows from Claim 24.

5.3 Sub-procedure BudgetMG-Metric.

Now we provide a high level description of BudgetMG-Metric. The details can be found in Algorithm 4. We first transform the problem to a Stochastic- k -TSP instance \mathcal{M}_{ktsp} , by reducing each Markov chain to a related random variable. Then we use the constant factor approximation algorithm ALG_{ktsp} developed in [JLLS20] to obtain an ordering Π of vertices (chains). Let Π_{pref} be the prefix of Π such that the switching cost for traversing Π_{pref} is no larger than $10\alpha_{ktsp}B$. Since the expected cost of the ALG_{ktsp} of instance \mathcal{M}_{ktsp} is at most $\alpha_{ktsp}B/10$, with a large constant probability, one can collect K units of rewards from Π_{pref} , by Markov inequality. Obviously, ignoring the switching cost, the optimal way (optimal in terms of movement cost) of collecting K units of rewards from Π_{pref} is to play the chains in Π_{pref} according to grade. However, such play may switch frequently among the chains and incur a huge switching cost. To

keep the switching cost under control, we visit each chain in Π_{pref} only once (by the way we define the prefix, we can see the switching cost is certainly within the given budget).

A key question now is to decide when to switch to the next if the current state is not economical to keep playing (it requires a large expected movement cost to reach the target of the chain). It turns out that we can find a threshold γ_{j+1} (which is computed from the K -th order statistics of $\{C_{sl}(\mathcal{S})\}_{\mathcal{S} \in \Pi_{pref}}$), such that if the grade of the current state is larger than the threshold, or our algorithm has spent too much on this chain, say a movement budget $100\alpha_{ktsp}B$ for each chain, then it is time to switch to the next chain on the Π_{pref} .

As mentioned, our algorithm needs to estimate the distribution of the K -th *order statistic* for a collection of random variables (that is the K -th smallest value). This can be approximated simply either by sampling (Monte Carlo) or the Bapat-Beg theorem and the fully polynomial randomized approximation scheme (FPRAS) for estimating the permanent [JSV04] (see the details in Appendix C.1). In the following, we safely ignore the estimation error and error probability for clarity purpose.

Algorithm 4: Algorithm BudgetMG-Metric

- 1 **Input:** The instance \mathcal{M} , objective number of rewards K
 - 2 **Process:**
 - 3 Compute the grades of all states and sort them in increasing order $\{\gamma_j\}_{j=1}^n$;
 - 4 Reduce the instance \mathcal{M} to a Stochastic- k -TSP instance \mathcal{M}_{ktsp} , by replacing each $\mathcal{S} \in \mathcal{M}$ with a vertex in \mathcal{M}_{ktsp} with selection cost distributed as $C_{sl}(\mathcal{S})$ (See the definition of $C_{sl}(\mathcal{S})$ in Section 5.2.2);
 - 5 $\Pi \leftarrow \text{ALG}_{ktsp}(\mathcal{M}_{ktsp})$ (Π is an ordering of vertices) ;
 - 6 Let $\Pi_{pref} = \{\mathcal{S}_0 = \mathbf{R}, \mathcal{S}_1, \mathcal{S}_2, \dots, \mathcal{S}_m\}$ be the longest prefix of Π such that $\sum_{i=0}^{m-1} d(\mathcal{S}_i, \mathcal{S}_{i+1}) \leq 10\alpha_{ktsp}B$;
 - 7 Estimate the distribution of K -th order statistic $C_{sl}^{[K]}(\mathcal{S})$ for $\{C_{sl}(\mathcal{S})\}_{\mathcal{S} \in \Pi_{pref}}$;
 - 8 Find the unique $\gamma_j, j \in [n]$ such that $\Pr[C_{sl}^{[K]}(\mathcal{S}) \leq \gamma_j] < 0.3 \wedge \Pr[C_{sl}^{[K]}(\mathcal{S}) \leq \gamma_{j+1}] \geq 0.3$;
 - 9 **for** $i = 1, \dots, m$ **do**
 - 10 Pay switching cost $d(\mathcal{S}_{i-1}, \mathcal{S}_i)$ and switches to \mathcal{S}_i ;
 - 11 Play \mathcal{S}_i until it reaches its target state or some condition(s) in \mathcal{A} (defined below) holds. ;
 - 12 If \mathcal{S}_i reaches its target state, let $K \leftarrow K - 1$;
 - 13 **end**
 - 14 **Return:** The updated instance \mathcal{M} ; the remaining number of target states K ;
 - 15 ;
 - 16 **Define:** Conditions \mathcal{A} (bad events);;
 - 17 (I) The next move on \mathcal{S}_i makes the movement cost on \mathcal{S}_i exceed $100\alpha_{ktsp}B$;
 - 18 (II) The current grade of \mathcal{S}_i is more than γ_j ;
-

5.4 Analysis.

Recall that it suffices to prove that BudgetMG-Metric satisfies Lemma 18. We only need to consider the case $B \geq 10\mathbb{E}[C_{tot}(\odot)]$, and prove two guarantees of BudgetMG-Metric:

- i) It can succeed to make K Markov system reach their target states with probability at least 0.1;

ii) The expectation of its total cost is bounded by $400\alpha_{\text{ktsp}}B$.

Recall that $C_{sl}(\mathcal{S})$ represents the (random) selection cost of \mathcal{S} under the rule of Stochastic- k -TSP. The set Π_{pref} we found in the line 5 of Algorithm 4 has some good properties stated below:

Lemma 25. *Suppose $B \geq 10\mathbb{E}[C_{tot}(\mathbb{O})]$. The set Π_{pref} found in the line 5 of Algorithm 4 satisfies that with probability at least 0.99, one has $\sum_{i=1}^K C_{sl}^{[i]}(\mathcal{S}) \leq 10\alpha_{\text{ktsp}}B$, where $C_{sl}^{[i]}(\mathcal{S})$ is the i -th order statistic for $\{C_{sl}(\mathcal{S})\}_{\mathcal{S} \in \Pi_{pref}}$.*

Proof. As $B \geq 10\mathbb{E}[C_{tot}(\mathbb{O})]$, we know that $B \geq 10\mathbb{E}[C_{sl}(\mathbb{O}_{\text{ktsp}}) + C_{sw}(\mathbb{O}_{\text{ktsp}})]$ by Claim 24.

By Theorem 20, we know that $\alpha_{\text{ktsp}}B \geq 10\mathbb{E}[C_{sl}(\text{ALG}_{\text{ktsp}}) + C_{sw}(\text{ALG}_{\text{ktsp}})]$. By Markov Inequality, we know that

$$\Pr[C_{sl}(\text{ALG}_{\text{ktsp}}) + C_{sw}(\text{ALG}_{\text{ktsp}}) \geq 10\alpha_{\text{ktsp}}B] \leq 0.01. \quad (3)$$

If $\sum_{i=1}^K C_{sl}^{[i]}(\mathcal{S}) > 10\alpha_{\text{ktsp}}B$, then it means that

$$C_{sl}(\text{ALG}_{\text{ktsp}}) + C_{sw}(\text{ALG}_{\text{ktsp}}) \geq 10\alpha_{\text{ktsp}}B. \quad (4)$$

More specifically, if $C_{sl}(\text{ALG}_{\text{ktsp}}) \leq 10\alpha_{\text{ktsp}}B$, as $\sum_{i=1}^K C_{sl}^{[i]}(\mathcal{S}) > 10\alpha_{\text{ktsp}}B$, then ALG_{ktsp} needs to visit vertices outside Π_{pref} to select some “cheap” vertices, which incurs a switching cost larger than $10\alpha_{\text{ktsp}}B$. So either switching cost or the selection cost of ALG_{ktsp} is larger than $10\alpha_{\text{ktsp}}B$.

By combining Equation 3 and Equation 4, we get

$$\Pr\left[\sum_{i=1}^K C_{sl}^{[i]}(\mathcal{S}) \leq 10\alpha_{\text{ktsp}}B\right] \geq 0.99.$$

□

We make use of the good properties of Π_{pref} via the fair game we define in Section 5.2.1. Define an important intermediate strategy \mathbb{GT} for the input instance \mathcal{M} under the rule of MG-Metric-Fair as follows: it always chooses the Markov system in Π_{pref} with the smallest grade (breaking ties arbitrarily) to play under the condition that $C_{mv}^F(\mathbb{GT}) \leq 10\alpha_{\text{ktsp}}B$, pays the *fair* movement cost *immediately* when some chain reaches its target and halts when it makes K systems reach their targets or the next step would break the fair movement cost budget.

To be more clear on the stopping condition of \mathbb{GT} on the cost budget, suppose \mathbb{GT} has not made K chains reach targets, the fair movement cost spent is X and the smallest grade among available chains (those not in the target states) is Y . Then \mathbb{GT} chooses the chain with smallest grade Y if $X + Y \leq 10\alpha_{\text{ktsp}}B$, and halts if no such chain exists.⁵ Then we have the following claim:

Claim 26. *With probability at least 0.99, \mathbb{GT} can succeed to make K systems reach their targets. Further, one has:*

$$\mathbb{E}[C_{mv}(\mathbb{GT})] \leq 10\alpha_{\text{ktsp}}B.$$

⁵In fact, if there is no switching cost and no movement cost budget, \mathbb{GT} is the optimal solution to make K chains in Π_{pref} to reach targets with the minimum expected movement cost.

Proof. The probability of success is a direct corollary of Lemma 25, as the fair movement cost of a non-quitting player on single \mathcal{S} is exactly distributed as the selection cost $C_{sl}(\mathcal{S})$. By Lemma 25, one has

$$\Pr \left[\sum_{i=1}^K C_{sl}^{[i]}(\mathcal{S}) \leq 10\alpha_{\text{ktsp}}B \right] \geq 0.99,$$

where $C_{sl}^{[i]}(\mathcal{S})$ is the i -th order statistic for $\{C_{sl}(\mathcal{S})\}_{\mathcal{S} \in \Pi_{\text{pref}}}$. Suppose there are m non-quitting players and every one plays one chain on Π_{pref} . With probability at least 0.99, the summation of the smallest k fair movement costs they pay will be no more than $10\alpha_{\text{ktsp}}B$, in which case \mathbb{GT} can succeed to make K systems reach their targets.

The bound on movement cost follows from the definition and the relation that

$$\mathbb{E}[C_{mv}(\mathbb{GT})] = \mathbb{E}[C_{mv}^F(\mathbb{GT})] \leq 10\alpha_{\text{ktsp}}B,$$

where the first equality comes from that \mathbb{GT} is a fair player under the rule of MG-Metric-Fair. More specifically, we can observe that \mathbb{GT} is playing a series of teasing game \mathcal{S}^T , and it can get the profits under the rule of \mathcal{S}^T and need to pay the fair movement cost under the rule of MG-Metric-Fair instead. Moreover, its “expected fair movement cost” equals to “expected profits”.

As it plays all of teasing games optimally by Lemma 10, the “expected profits” equals to its “expected movement cost”. Thus, we get the equality. \square

Recall again the objective is to show that BudgetMG-Metric satisfies the precondition of Lemma 18. We prove the guarantee on success probability first:

Lemma 27 (Success probability). *With probability at least 0.1, BudgetMG-Metric can make at least K Markov system reach their target states.*

Proof. Note that we use $C_{sl}(\mathcal{S})$ for system \mathcal{S} to represent the random selection cost under the rule of Stochastic- k -TSP, and now we consider the instance under the original rule, i.e. in MG-Metric.

For simplicity, we use $C_{mv}(\mathcal{S}, \gamma_{j+1})$ to represent the (random) movement cost $C_{mv}(\mathbb{O}(\mathcal{S}(\gamma_{j+1})))$, where the game $\mathcal{S}(\gamma_{j+1})$ and its optimal strategy $\mathbb{O}(\mathcal{S}(\gamma_{j+1}))$ is used to defined the grade in Subsection 3.1.

Suppose we play each $\mathcal{S} \in \Pi_{\text{pref}}$ one by one according to $\mathbb{O}(\mathcal{S}(\gamma_{j+1}))$ and denote the random subset of systems which reach their targets by $\Pi_{\text{suc}} \subseteq \Pi_{\text{pref}}$. If $|\Pi_{\text{suc}}| \geq k$, we let $C_{mv}^{[i]}$ be the i -th order statistic for $\{C_{mv}(\mathcal{S})\}_{\mathcal{S} \in \Pi_{\text{suc}}}$ and $i \in [k]$, otherwise we define $C_{mv}^{[i]}$ to be the i -th order statistic when $i \in [|\Pi_{\text{suc}}|]$ and $C_{mv}^{[i]} = 0$ for $|\Pi_{\text{suc}}| < i \leq k$. We have the following claim:

$$\Pr \left[\sum_{i=1}^K C_{mv}^{[i]} \geq 100\alpha_{\text{ktsp}}B \right] \leq 0.2. \tag{5}$$

This can be proved by contradiction. Suppose $\Pr[\sum_{i=1}^K C_{mv}^{[i]} \geq 100\alpha_{\text{ktsp}}B] > 0.2$. Since the success probability of \mathbb{GT} is at least 0.99, by union bound, one can see that

$$\Pr \left[\sum_{i=1}^K C_{mv}^{[i]} \geq 100\alpha_{\text{ktsp}}B \wedge \mathbb{GT} \text{ gets } K \text{ rewards} \right] \geq 0.19.$$

Further, we know when event that $(\sum_{i=1}^K C_{mv}^{[i]} \geq 100\alpha_{ktsp}B \wedge \mathbb{GT} \text{ gets } K \text{ rewards})$ happens, then the event $(C_{mv}(\mathbb{GT}) \geq 100\alpha_{ktsp}B \wedge \mathbb{GT} \text{ gets } K \text{ rewards})$ also happens, as $\sum_{i=1}^K C_{mv}^{[i]}$ is the minimum possible movement cost of any strategy which makes K systems on Π_{pref} reach targets. This implies that

$$\Pr [C_{mv}(\mathbb{GT}) \geq 100\alpha_{ktsp}B \wedge \mathbb{GT} \text{ gets } K \text{ rewards}] \geq \Pr \left[\sum_{i=1}^K C_{mv}^{[i]} \geq 100\alpha_{ktsp}B \wedge \mathbb{GT} \text{ gets } K \text{ rewards} \right] \geq 0.19.$$

Further, the inequality $\Pr[C_{mv}(\mathbb{GT}) \geq 100\alpha_{ktsp}B] \geq 0.19$ follows directly and the expected movement cost of \mathbb{GT} is at least $100\alpha_{ktsp}B * 0.19 = 19\alpha_{ktsp}B$. This contradicts Claim 26. Thus we complete the proof of Equation 5.

Moreover, we know that $\Pr[C_{sl}^{[k]}(\mathcal{S}) \leq \gamma_{j+1}] \geq 0.3$ by the condition in Line 8 in BudgetMG-Metric. By Union bound, we know that

$$\Pr \left[C_{sl}^{[k]}(\mathcal{S}) \leq \gamma_{j+1} \wedge \sum_{i=1}^K C_{mv}^{[i]} \leq 100\alpha_{ktsp}B \right] \geq 0.1,$$

which implies the success probability of BudgetMG-Metric. More specifically, the event $C_{sl}^{[k]}(\mathcal{S}) \leq \gamma_{j+1}$ ensures that the random subset $|\Pi_{suc}| \geq K$, in which case event $\sum_{i=1}^K C_{mv}^{[i]} \leq 100\alpha_{ktsp}B$ implies that our algorithm BudgetMG-Metric can make at least K chains reach target states. \square

Now it remains to bound the expected total cost of BudgetMG-Metric. We state a technical result at first.

Claim 28. *With probability at least 0.65, \mathbb{GT} can succeed to make K systems reach their targets with no less movement cost than BudgetMG-Metric. In other word, one has*

$$\Pr[\mathbb{GT} \text{ gets } K \text{ rewards} \wedge C_{mv}(\mathbb{GT}) \geq C_{mv}(\text{BudgetMG-Metric})] \geq 0.65.$$

Proof. On one hand, by Claim 26, we know the success probability of \mathbb{GT} is at least 0.99, i.e.

$$\Pr[\mathbb{GT} \text{ gets } K \text{ rewards}] \geq 0.99. \quad (6)$$

On the other hand, we know that $\Pr[C_{sl}^{[K]}(\mathcal{S}) > \gamma_j] \geq 0.7$, which means that

$$\Pr[C_{sl}^{[K]}(\mathcal{S}) \geq \gamma_{j+1}] \geq 0.7, \quad (7)$$

as we have assumed that the grade is distinct of each other. By Union Bound over Equation 6 and Equation 7, one has

$$\Pr[\mathbb{GT} \text{ gets } K \text{ rewards} \wedge C_{sl}^{[K]}(\mathcal{S}) \geq \gamma_{j+1}] \geq 0.69.$$

Conditioning on that $C_{sl}^{[K]}(\mathcal{S}) \geq \gamma_{j+1}$ holds and \mathbb{GT} succeeds to get K units of rewards, the grade of all *Markovsystem* $\in \Pi_{pref}$ after \mathbb{GT} halts are at least γ_{j+1} . This is true as \mathbb{GT} always chooses the *Markovsystem* in the Π_{pref} with the smallest grade to advance. As BudgetMG-Metric never plays the system when the grade of its current state exceeds γ_{j+1} , we know $C_{mv}(\mathbb{GT}) \geq C_{mv}(\text{BudgetMG-Metric})$. In particular,

$$\Pr[C_{mv}(\mathbb{GT}) \geq C_{mv}(\text{BudgetMG-Metric})] \geq \Pr[\mathbb{GT} \text{ gets } K \text{ rewards} \wedge C_{sl}^{[K]}(\mathcal{S}) \geq \gamma_{j+1}] \geq 0.69.$$

Combining these together, one has

$$\begin{aligned}
& \Pr[\text{GT gets } K \text{ rewards} \wedge C_{mv}(\text{GT}) \geq C_{mv}(\text{BudgetMG-Metric})] \\
&= \Pr[C_{mv}(\text{GT}) \geq C_{mv}(\text{BudgetMG-Metric}) \mid \text{GT gets } K \text{ rewards}] \cdot \Pr[\text{GT gets } K \text{ rewards}] \\
&\geq \Pr\left[C_{sl}^{[K]}(\mathcal{S}) \geq \gamma_{j+1} \mid \text{GT gets } K \text{ rewards}\right] \cdot \Pr[\text{GT gets } K \text{ rewards}] \\
&\geq \Pr\left[C_{sl}^{[K]}(\mathcal{S}) \geq \gamma_{j+1} \mid \text{GT gets } K \text{ rewards}\right] * 0.99,
\end{aligned}$$

where the first inequality comes from the argument above and the second inequality comes from Equation 6.

The remaining is to bound $\Pr\left[C_{sl}^{[K]}(\mathcal{S}) \geq \gamma_{j+1} \mid \text{GT gets } K \text{ rewards}\right]$. By conditional probability, we know that

$$\Pr\left[C_{sl}^{[K]}(\mathcal{S}) \geq \gamma_{j+1} \mid \text{GT gets } K \text{ rewards}\right] = \frac{\Pr[C_{sl}^{[K]}(\mathcal{S}) \geq \gamma_{j+1} \wedge \text{GT gets } K \text{ rewards}]}{\Pr[\text{GT gets } K \text{ rewards}]} \geq 0.69,$$

which can complete the proof with elementary calculation. \square

Now we are ready to bound the total cost.

Lemma 29 (Bounded expected cost). *The expected total cost of BudgetMG-Metric is upper bounded by $O(1)B$. More specifically, one has*

$$\mathbb{E}[C_{tot}(\text{BudgetMG-Metric})] \leq 410\alpha_{\text{ktsp}}B.$$

Proof. We prove this by contradiction. Suppose $\mathbb{E}[C_{tot}(\text{BudgetMG-Metric})] > 410\alpha_{\text{ktsp}}B$. Recall that the switching cost of BudgetMG-Metric is at most $10\alpha_{\text{ktsp}}B$. Then the assumption implies that

$$\mathbb{E}[C_{tot}(\text{BudgetMG-Metric})] > 400\alpha_{\text{ktsp}}B.$$

In fact, BudgetMG-Metric plays each system in Π_{pref} independently of other system. Recall that we let $C_{mv}(\mathcal{S})$ be the (random) movement cost $C_{mv}(\odot(\mathcal{S}(\gamma_{j+1})))$. The only difference between BudgetMG-Metric and $\odot(\mathcal{S}(\gamma_{j+1}))$ on \mathcal{S} is that BudgetMG-Metric stops playing \mathcal{S} if the next step will break the movement budget $100\alpha_{\text{ktsp}}B$. We let $C'_{mv}(\mathcal{S})$ be the (random) movement cost BudgetMG-Metric spends on \mathcal{S} and obviously, we have $\mathbb{E}[C_{mv}(\text{BudgetMG-Metric})] = \sum_{\mathcal{S} \in \Pi_{pref}} \mathbb{E}[C'_{mv}(\mathcal{S})]$.

The assumption implies that $\sum_{\mathcal{S} \in \Pi_{pref}} \mathbb{E}[C'_{mv}(\mathcal{S})] > 400\alpha_{\text{ktsp}}B$. By the Condition (I) in \mathcal{A} , $C'_{mv}(\mathcal{S}) \in [0, 100\alpha_{\text{ktsp}}B]$. By Chernoff bound (Theorem 30), one has:

$$\Pr\left[\sum_{\mathcal{S} \in \Pi_{pref}} C'_{mv}(\mathcal{S}) \leq 200\alpha_{\text{ktsp}}B\right] \leq \Pr\left[\left|\sum_{\mathcal{S} \in \Pi_{pref}} C'_{mv}(\mathcal{S}) - \mathbb{E}\left[\sum_{\mathcal{S} \in \Pi_{pref}} C'_{mv}(\mathcal{S})\right]\right| \geq 300\alpha_{\text{ktsp}}B\right] \leq 2e^{-9/2} \leq 0.1,$$

which means with probability at least 0.9, the movement cost of BudgetMG-Metric exceeds $200\alpha_{\text{ktsp}}B$, i.e.

$$\Pr[C_{mv}(\text{BudgetMG-Metric}) \geq 200\alpha_{\text{ktsp}}B] = \Pr\left[\sum_{\mathcal{S} \in \Pi_{pref}} C'_{mv}(\mathcal{S}) \geq 200\alpha_{\text{ktsp}}B\right] \geq 0.9. \quad (8)$$

In particular, by Claim 28, with probability at least 0.65, \mathbb{GT} pays more movement cost than BudgetMG-Metric. Then by union bound, one has

$$\Pr[C_{mv}(\mathbb{GT}) \geq 200\alpha_{\text{ktsp}}B] \geq 0.55,$$

which implies that $\mathbb{E}[C_{mv}(\mathbb{GT})] \geq 110\alpha_{\text{ktsp}}B$ and contradicts Claim 26.

Thus we complete the proof. \square

Hitherto we have proved that BudgetMG-Metric satisfies Lemma 18, and thus finished the proof of the main result:

Theorem 5. *There is a constant factor approximation algorithm for the MG-Metric problem.*

6 Conclusions and Open Problems

In this work, we present a simple index strategy for MG-Unit and a more involved algorithm for MG-Metric, both achieving constant approximation ratios. We did not attempt to optimize the exact constants and the constants directly implied from our analysis are quite large for both problems.⁶ Designing new algorithms or analysis with small approximation constants is a very interesting further direction. In particular, we suspect the approximation ratio of MG-Unit is a small constant.

One interesting future direction is to study the general problem proposed in [GJSS19] with switching cost. In particular, there is a given combinatorial constraint $\mathcal{F} \subseteq 2^{[n]}$, and our goal is to make a subset $F \in \mathcal{F}$ of chains reach their targets. Another interesting extension is to significantly generalize the stochastic reward k -TSP studied in [ENS18, JLLS20] as follows: we have metric graph, in which each node is associated with a Markov chain. Each target state of a Markov chain is associated with a random reward $R_v \in \mathbb{Z}^+$, which is realized when we reach the target. The goal is to collect a total reward of at least k .

7 Acknowledgments

The research is supported in part by the National Natural Science Foundation of China Grant 61822203, 61772297, 61632016, 61761146003, and Turing AI Institute of Nanjing and Xi'an Institute for Interdisciplinary Information Core Technology, NSF awards CCF-1749609, DMS-1839116, DMS-2023166, CCF-2105772, Microsoft Research Faculty Fellowship, Sloan Research Fellowship, and Packard Fellowships. Part of the work was done while Liu visited Shanghai Qi Zhi Institute.

⁶The approximation ratio implied from the current analysis of MG-Unit is around 10^5 , estimated as follows. BudgetMG-Unit uses 2^8B budget and can succeed with probability at least $1/20$ conditioning on that the expected cost of optimal solution is no more than $B/10$. Plugging these constants in Lemma 12, we know the β in the doubling framework should be set as $1/\beta^2 = 1 - 1/20 = 0.95$, and the final approximation ratio can be around $2^8 * \frac{10}{\beta-1} \approx 10^5$. The approximation ratio of MG-Metric depends on the the approximation ratio for Stochastic- k -TSP, which is already quite large.

A Appendix for Section 3

A.1 Concentration Inequalities

We need the following two well known concentration inequalities.

Theorem 30 (Chernoff-Hoeffding Bound). *Let X_1, X_2, \dots, X_n be independent random variables taking values in $[0, 1]$ and define $X := \frac{1}{n} \sum_{i \in [n]} X_i$. Then for any $\epsilon \in [0, 1]$, we have*

$$\mathbb{P}[|X - \mathbb{E}[X]| \geq \epsilon] \leq 2 \exp(-n\epsilon^2/2).$$

Theorem 31 (Freedman's Inequality, Theorem 1.6 in [Fre75]). *Consider a real-valued martingale difference sequence $\{X_t\}_{t \geq 0}$ such that $X_0 = 0$, and $\mathbb{E}[X_{t+1} | \mathcal{F}_t] = 0$ for all t , where $\{\mathcal{F}_t\}_{t \geq 0}$ is the filtration defined by the martingale. Assume that the sequence is uniformly bounded, i.e., $|X_t| \leq M$ almost surely for all t . Now define the predictable quadratic variation process of the martingale to be $W_t = \sum_{j=1}^t \mathbb{E}[X_j^2 | \mathcal{F}_{j-1}]$ for all $t \geq 1$. Then for all $\ell \geq 0$ and $\sigma^2 > 0$ and any stopping time τ , we have*

$$\mathbb{P}\left[\left|\sum_{j=0}^{\tau} X_j\right| \geq \ell \wedge W_{\tau} \leq \sigma^2 \text{ for some stopping time } \tau\right] \leq 2 \exp\left(-\frac{\ell^2/2}{\sigma^2 + M\ell/3}\right).$$

A.2 The Doubling Technique

A.2.1 Proof of Lemma 11

Lemma 11. *For any $j \geq i \geq 1$ and any Algorithm BudgetMG, one has*

$$\mathbb{E}[C_{tot}(\mathbb{O}(\mathcal{M}_i, k_i))] \geq \mathbb{E}[C_{tot}(\mathbb{O}(\mathcal{M}_j, k_j))].$$

Notice that the randomness is over an entire run of Algo-MG.

Proof. For any phases $i \geq 0$, we define:

- $\omega_{i,S}$: the trajectory of Markov system S traversed by Algo-MG in the first i phases.
- $\omega_i = \cup_{S \in \mathcal{S}} \omega_{i,S}$: the collection of the trajectories of all Markov system systems traversed by Algo-MG in the first i phases.

At a first glance, $\mathbb{O}(\mathcal{M}_j, k_j)$ is a sub-tree of $\mathbb{O}(\mathcal{M}_i, k_i)$ and this inequality holds directly. But this is not true as $\mathbb{O}(\mathcal{M}_j, k_j)$ is required to start at the root \mathbf{R} .

In particular, we know that \mathcal{M}_i and k_i can be determined by ω_i . First, for all $j \geq i \geq 0$, we have that

$$\mathbb{E}[C_{tot}(\mathbb{O}(\mathcal{M}_i, k_i))] = \mathbb{E}_{\omega_j} [\mathbb{E}[C_{tot}(\mathbb{O}(\mathcal{M}_i, k_i)) | \omega_j]].$$

Now, fixing any possible ω_j , it suffices to prove that

$$\mathbb{E}[C_{tot}(\mathbb{O}(\mathcal{M}_i, k_i)) | \omega_j] \geq \mathbb{E}[C_{tot}(\mathbb{O}(\mathcal{M}_j, k_j)) | \omega_j]. \quad (9)$$

Note that we can represent the strategy $\mathbb{O}(\mathcal{M}_i, k_i)$ as a decision tree. No matter what (\mathcal{M}_j, k_j) is for the j -th phase, we consider the algorithm $\text{ALG}(\mathcal{M}_j, k_j)$ which uses the decision tree of $\mathbb{O}(\mathcal{M}_i, k_i)$, pretending not to know $\omega_j \setminus \omega_i$. The only difference is that we do not charge the movement cost of $\text{ALG}(\mathcal{M}_j, k_j)$ for those transitions occurring at $\omega_j \setminus \omega_i$.

Thus we know that

$$\mathbb{E}[C_{tot}(\mathbb{O}(\mathcal{M}_i, k_i)) \mid \omega_j] \geq \mathbb{E}[C_{tot}(\text{ALG}(\mathcal{M}_j, k_j)) \mid \omega_j] \geq \mathbb{E}[C_{tot}(\mathbb{O}(\mathcal{M}_j, k_j)) \mid \omega_j],$$

where the second inequality comes from the optimality of $\mathbb{O}(\mathcal{M}_j, k_j)$. By integrating over all possible ω_j , we can complete the proof. \square

A.2.2 Proof of Lemma 12

Lemma 12. *We are given some MG-Metric (or MG-Unit) instance \mathcal{M} , objective number of rewards k and non-negative real number $B \in \mathbb{R}_{\geq 0}$. For any $B > c_1 \mathbb{E}[C_{tot}(\mathbb{O}(\mathcal{M}, k))]$, if we can get an algorithm BudgetMG with $\mathbb{E}[\text{BudgetMG}] \leq c_2 B$ and can succeed with probability more than 0.01 where c_1, c_2 are some universal constant, then we can get an $O(1)$ -approximation algorithm for MG-Metric (or MG-Unit).*

Proof. Assume that $\min_{\mathcal{S} \in \mathcal{S}} d(\mathcal{S}, \mathbb{R}) = 1$. If there is an algorithm BudgetMG satisfying the precondition in this lemma, we can put it into the general framework Algo-MG and argue that Algo-MG is an $O(1)$ -approximation algorithm for MG-Metric (or MG-Unit). Let \mathcal{M}, K be the original input in Algorithm 1 and for simplicity let \mathbb{O} be the optimal strategy $\mathbb{O}(\mathcal{M}, K)$.

We need some notations for any phases $i \geq 0, j \geq 1$:

- $u_j(\omega_i)$: probability that Algo-MG enters phase $j + 1$ conditioning on ω_i .
- u_j : probability that Algo-MG enters phase $j + 1$.

Notice that $u_{i-1}(\omega_{i-1})$ is the indicator variable that Algo-MG enters phase i . For i -th phase, as the budget of BudgetMG($\mathcal{M}_{i-1}, k_{i-1}, B_i$) is $B_i = c_2 \beta^i$, the total cost of Algo-MG is at most $2c_2 \beta^i$ for i th phase. Note that $\mathbb{E}[C_{tot}(\mathbb{O}(\mathcal{M}_i, k_i))]$ is a random variable with randomness from ω_i .

Let $\mathbf{1}_{\omega_{i-1}}$ be the indicator variable that $\mathbb{E}[C_{tot}(\mathbb{O}(\mathcal{M}_i, k_i)) \mid \omega_{i-1}] \geq \frac{1}{c_1} \beta^i$. If the precondition in the statement holds, we know that there exists some universal constant $\beta > 1$, for any phase $i \geq 0$, and any possible ω_{i-1} such that the algorithm Algo-MG satisfies

$$u_i(\omega_{i-1}) \leq \frac{u_{i-1}(\omega_{i-1})}{\beta^2} + \mathbf{1}_{\omega_{i-1}}. \quad (10)$$

To make it more clear, we only need to consider the case when $u_{i-1}(\omega_{i-1}) = 1$ and $\mathbf{1}_{\omega_{i-1}} = 0$, as $u_{i-1}(\omega_{i-1}) = 0$ can imply $u_i(\omega_{i-1}) = 0$ directly and the inequality holds. Besides, if $\mathbf{1}_{\omega_{i-1}} = 0$, then we know $u_i(\omega_{i-1}) \geq 1 - 0.01 = 0.99$ by the precondition of the statement. Setting $1/\beta^2 \geq 0.99$ can handle this case.

Now our objective is to show that $\mathbb{E}[C_{tot}(\text{Algo-MG})] \leq O(1) \cdot \mathbb{E}[C_{tot}(\mathbb{O})]$ by using Equation 10.

As we have shown $\mathbb{E}[C_{tot}(\mathbb{O}) \mid \omega_{i-1}] \geq \mathbb{E}[C_{tot}(\mathbb{O}(\mathcal{M}_i, k_i)) \mid \omega_{i-1}]$ in the proof of Lemma 11 for any possible ω_{i-1} , then $\mathbb{E}[C_{tot}(\mathbb{O}(\mathcal{M}_i, k_i)) \mid \omega_{i-1}] \geq B$ means that $\mathbb{E}[C_{tot}(\mathbb{O}) \mid \omega_{i-1}] \geq B$ for any real number B , which gives us:

$$\begin{aligned} \Pr[\mathbb{E}[C_{tot}(\mathbb{O}(\mathcal{M}_i, k_i)) \geq B] &\leq \Pr[\mathbb{E}[C_{tot}(\mathbb{O})] \geq B] \\ &= \mathbf{1}_{\mathbb{E}[C_{tot}(\mathbb{O}) \geq B]}, \end{aligned}$$

where $\mathbf{1}_{\mathbb{E}[C_{tot}(\mathbb{O}) \geq B]}$ is a deterministic indicator variable that $\mathbb{E}[C_{tot}(\mathbb{O}) \geq B]$. By taking expectation of both sides of Equation 10, we have

$$u_i \leq u_{i-1}/\beta^2 + \Pr[\mathbb{E}[C_{tot}(\mathbb{O}(\mathcal{M}_i, k_i)) \geq \frac{1}{c_1}\beta^i]], \quad \forall i \geq 0.$$

Multiplying β^i on both sides and taking summation of all i , we have

$$\begin{aligned} \sum_{i \geq 1} u_i \cdot \beta^i &\leq u_0/\beta + 1/\beta \cdot \sum_{i \geq 1} u_i \cdot \beta^i + \sum_{i \geq 1} \Pr[\mathbb{E}[C_{tot}(\mathbb{O}(\mathcal{M}_i, k_i))] \geq \frac{1}{c_1}\beta^i] \cdot \beta^i \\ &\leq u_0/\beta + 1/\beta \cdot \sum_{i \geq 1} u_i \cdot \beta^i + \sum_{i \geq 1} \Pr[\mathbb{E}[C_{tot}(\mathbb{O})] \geq \frac{1}{c_1}\beta^i] \cdot \beta^i \\ &= 1/\beta + 1/\beta \cdot \sum_{i \geq 1} u_i \cdot \beta^i + \sum_{i \geq 1} \mathbf{1}_{\mathbb{E}[C_{tot}(\mathbb{O})] \geq \frac{1}{c_1}\beta^i} \cdot \beta^i \\ &\leq 1/\beta + 1/\beta \cdot \sum_{i \geq 1} u_i \cdot \beta^i + O(1) \cdot \mathbb{E}[C_{tot}(\mathbb{O})], \end{aligned}$$

which implies that

$$(1 - 1/\beta) \cdot \sum_{i \geq 1} u_i \cdot \beta^i \leq 1/\beta + O(1) \cdot \mathbb{E}[C_{tot}(\mathbb{O})].$$

Besides, recall the assumption that $\min_{\mathcal{S} \in \mathbb{S}} d(\mathcal{S}, \mathbb{R}) = 1$, one has $\mathbb{E}[C_{tot}(\mathbb{O})] \geq 1$ and that

$$\mathbb{E}[C_{tot}(\text{Algo-MG})] \leq O(1) \cdot \sum_{i \geq 0} (u_i - u_{i+1}) \cdot \beta^{i+1} = O(1)(\beta - 1) \cdot \sum_{i \geq 0} u_i \cdot \beta^i.$$

Combining these together, we get $\mathbb{E}[C_{tot}(\text{Algo-MG})] \leq O(1) \cdot \mathbb{E}[C_{tot}(\mathbb{O})]$. \square

B Appendix for Section 4

B.1 Algorithm.

Algorithm B.1 is simply an instantiation of the doubling framework Algorithm 2.

Algorithm 5: Algorithm for MG-Unit (for analysis purpose)

```

1 Set  $\beta \in (1, 2)$ ;
2 Set  $k_0 = 1$ ;
3 Set  $\mathcal{M}_0 = \mathcal{M}$ ;
4 Set  $c_1 = O(1)$ ;
5 for phase  $i = 0, 1, \dots$  do
6    $(\mathcal{M}_i, k_i) \leftarrow$  BudgetMG-Unit( $\mathcal{M}_i, c_1\beta^i$ ) (Algorithm 3);
7   if  $k_i \leq 0$  then
8     | Break
9   end
10 end

```

One has the following claim:

Claim 32. *Algorithm 5 has expected cost at least that of Algorithm 2.*

Proof. Actually, Algorithm 5 and Algorithm 2 proceed in the same manner, but whilst Algorithm 2 halts when finding a unit reward, Algorithm 5 continues until the end of the phase. Algorithm 2 only costs less. \square

By the previous claim, it suffices to prove that Algorithm 5 achieves a constant approximation factor for MG-Unit.

B.2 Proof of Lemma 17

Lemma 17. *Suppose X_1, X_2, \dots, X_n are a sequence of random variables taking values in $\{0, 1\}$, and $\mathcal{F}_j = \sigma(X_1, \dots, X_j)$ is the filtration defined by the sequence. Given any real number μ , if $\sum_{j=1}^n \mathbb{E}[X_j | \mathcal{F}_{j-1}] \geq \mu$, then*

$$\Pr\left[\sum_{j=1}^n X_j \geq 1\right] \geq 1 - e^{-3\mu/8}.$$

Proof. We construct a Martingale difference sequence (MDS) Y_1, \dots, Y_n to prove this lemma where $Y_j = X_j - \mathbb{E}[X_j | \mathcal{F}_{j-1}]$. It is easy to check that $\{Y_j\}$ is a MDS: First $\mathbb{E}[|Y_j|] \leq 1$. Second, $\mathbb{E}[Y_j | \mathcal{F}_{j-1}] = \mathbb{E}[X_j | \mathcal{F}_{j-1}] - \mathbb{E}[X_j | \mathcal{F}_{j-1}] = 0$.

We try to use Freedman's Inequality (Theorem 31). One has

$$W_t = \sum_{j=1}^t \mathbb{E}[Y_j^2 | \mathcal{F}_{j-1}] \leq \sum_{j=1}^t \mathbb{E}[X_j^2 | \mathcal{F}_{j-1}] \leq \sum_{j=1}^t \mathbb{E}[X_j | \mathcal{F}_{j-1}].$$

Hence we have

$$\begin{aligned} \Pr\left[\sum_{j=1}^n X_j = 0\right] &= \Pr\left[\sum_{j=1}^n Y_j = -\sum_{j=1}^n \mathbb{E}[X_j | \mathcal{F}_{j-1}]\right] \\ &\leq \Pr\left[\sum_{j=1}^n Y_j \leq -\mu \wedge W_n \leq \mu\right] \\ &\leq \exp(-3\mu/8). \end{aligned}$$

\square

B.3 A Counter Example

The simple index-based strategy is a constant factor approximation in unit metric space. However, there is a simple counter example with non-unit switching cost in which BudgetMG-Unit performs arbitrarily worse than the optimal strategy.

Consider the following instance: We are given a metric space $\mathcal{M} = (\mathbf{S} \cup \{\mathbf{R}\}, d)$. There are only two kinds of Markov system. The first kind of Markov system $\mathcal{S} = \langle \{s, t, x\}, P, C, s, t \rangle$, where $C_s = 0, C_x = +\infty$, and

$P_{s,t} = \epsilon, P_{s,x} = 1 - \epsilon, P_{x,t} = 1$. The second kind of Markov system \mathcal{S}' has the similar structure, except that $P'_{s,t} = \epsilon/2$ and $P'_{s,x} = 1 - \epsilon/2$. The infinite set $\mathbf{S} = \{\mathcal{S}_1, \mathcal{S}_2, \dots\} \cup \{\mathcal{S}'_1, \mathcal{S}'_2, \dots\}$.

Let $d(\mathbf{R}, \mathcal{S}_k) = 1$ for all k and $d(\mathcal{S}_k, \mathcal{S}_j) = 1$ for $k \neq j$. Let $d(\mathcal{S}'_j, \mathbf{R}) = d(\mathcal{S}'_j, \mathcal{S}_k) = \sum_{i=1}^j 2^{-i+1}$ for all k . As for the distances between the second kind of Markov system, for $k < j$, we set $d(\mathcal{S}'_k, \mathcal{S}'_j) = \sum_{i=k-1}^{j-1} 2^{-i}$.

We are at \mathbf{R} in the beginning and need to pay the switching cost to advance any one of Markov system. The optimal strategy is obviously pay unit switching cost and continue playing the second kind of Markov system until getting unit reward, and the total cost is no more than 2. But the greedy algorithm 2 always chooses one of the first kind of Markov system \mathcal{S}_i to advance, with the expected cost $1/\epsilon$, which can be large arbitrarily.

C Appendix for Section 5

C.1 Estimate Order Statistics

Definition 33 (Permanent). *The permanent of an $n \times n$ matrix A is defined by*

$$\text{per}(A) = \sum_{\sigma} A_{i,\sigma(i)},$$

where the sum is over all permutations σ of $\{1, 2, \dots, n\}$.

Theorem 34 (Bapat-Beg theorem). *Let X_1, \dots, X_n be independent real valued random variables with cumulative distribution functions respectively $F_1(x), \dots, F_n(x)$. Write $X_{(1)}, \dots, X_{(n)}$ for the order statistics. Then the joint probability distribution of the n_1, \dots, n_k order statistics (with $n_1 < n_2 < \dots < n_k$ and $x_1 < x_2 < \dots < x_k$) is*

$$\begin{aligned} F_{X_{(n_1)}, \dots, X_{(n_k)}}(x_1, \dots, x_k) &= \Pr(X_{(n_1)} \leq x_1 \wedge X_{(n_2)} \leq x_2 \wedge \dots \wedge X_{(n_k)} \leq x_k) \\ &= \sum_{i_k=n_k}^n \dots \sum_{i_2=n_2}^{i_3} \sum_{i_1=n_1}^{i_2} \frac{P_{i_1, \dots, i_k}(x_1, \dots, x_k)}{i_1! (i_2 - i_1)! \dots (n - i_k)!}, \end{aligned}$$

where

$$\text{per} \begin{pmatrix} P_{i_1, \dots, i_k}(x_1, \dots, x_k) = \\ \begin{array}{cccc} F_1(x_1) \cdots F_1(x_1) & F_1(x_2) - F_1(x_1) \cdots F_1(x_2) - F_1(x_1) & \cdots & 1 - F_1(x_k) \cdots 1 - F_1(x_k) \\ F_2(x_1) \cdots F_2(x_1) & F_2(x_2) - F_2(x_1) \cdots F_2(x_2) - F_2(x_1) & \cdots & 1 - F_2(x_k) \cdots 1 - F_1(x_k) \\ \vdots & \vdots & \vdots & \vdots \\ F_n(x_1) \cdots F_n(x_1) & F_n(x_2) - F_n(x_1) \cdots F_n(x_2) - F_n(x_1) & \cdots & 1 - F_n(x_k) \cdots 1 - F_n(x_k) \end{array} \end{pmatrix}$$

$\underbrace{\hspace{10em}}_{i_1} \quad \underbrace{\hspace{10em}}_{i_2 - i_1} \quad \underbrace{\hspace{10em}}_{n - i_k}$

Theorem 35. *There exists a fully polynomial randomized approximation scheme for the permanent of an arbitrary $n \times n$ matrix A with nonnegative entries.*

Lemma 36. *For the set of distributions $C_{sl}^1, C_{sl}^2, \dots, C_{sl}^n$ from which we can take i.i.d. samples, define $C_{sl}^{[K]}$ be the K -th order statistic. For any real number $x \in R$, we can estimate the value $\Pr[C_{sl}^{[k]} \leq x]$ within additive error ϵ with probability at least $1 - 1/\text{poly}(n)$ in polynomial time.*

Proof. $O(\log n)$ samples are enough to get an estimation of $\Pr[C_{sl}^{[k]}(\mathcal{S}) \leq x]$ within additive error $\epsilon = 0.01$ with probability at least $1 - 1/\text{poly}(n)$ by using Chernoff-Hoeffding Bound (Theorem 30).

This can also be replaced by the fully polynomial randomized approximation scheme (FPRAS) in [JSV04] to estimate the permanent of $n - k$ matrices constructed according to Bapat-Beg theorem. \square

C.2 Computation of Index and CDF

Lemma 37. *We can calculate the cumulative distribution function of $C_{sl}(\mathcal{S})$ efficiently.*

Proof of Lemma 37. Consider the Markov system $\langle V, P, C, s, t \rangle$. For simplicity, we denote the Gittins index of states in V as $\gamma_1 < \gamma_2 < \dots < \gamma_n$ where $|V| = n$ and $\gamma_1 = 0$ being the Gittins Index of the target state t .

For any $\gamma_i \leq B < \gamma_{i+1}$, we know that

$$\Pr[C_{sl}(\mathcal{S}) \leq B] = \Pr[C_{sl}(\mathcal{S}) \leq \gamma_i].$$

Let U contain those states whose Gittins Index are larger than γ_i . For state v , let x_v be the probability that the non-quitting player can get no more than γ_i units of profits under the rule of \mathcal{S}^T conditioning on he is at state v currently. Then we know that $x_v = \sum_{u \in N(v)} P_{v,u} x_u$, where $N(v) = \{u \in V \mid \gamma_u \leq \gamma_j\}$ if $v \notin U$; and $x_v = 0$ if $v \in U$.

Then it involves solving a $i \times i$ system of equations, which can be done some standard techniques like LU decomposition in $O(i^3)$ time. \square

C.3 Algorithm Framework

We present our Algorithm Framework (Algorithm 6) here.

Algorithm 6: Algorithm MG-Metric

```

1 Input: The instance  $\mathcal{M}$ , objective number of rewards  $K$ 
2 Process:
3 set  $\beta \in (1, 2)$ ;
4 set  $k_0 \leftarrow K$ ;
5 set  $\mathcal{M}_0 \leftarrow \mathcal{M}$ ;
6 for phase  $i = 1, \dots$  do
7    $(\mathcal{M}_i, k_i) \leftarrow \text{BudgetMG-Metric}(\mathcal{M}, K, 50000\beta^i)$ ;
8   if  $k_i \leq 0$  then
9     Break
10  end
11 end

```

C.4 Omitted Proof

C.4.1 Proof of Claim 23

Claim 23. *The following inequality holds: $\mathbb{E}[C_{tot}^F(\mathbb{P})] = \mathbb{E}[C_{sw}(\mathbb{O}_{fair}) + C_{mv}^F(\mathbb{O}_{fair})] \leq \mathbb{E}[C_{tot}(\mathbb{O})]$.*

Proof. From another perspective, the player is playing a series of the teasing game \mathcal{S}^T with switching cost. In the rule of \mathcal{S}^T , the player can get some amount of profits on the target state of \mathcal{S} , while in the rule MG-Metric-Fair, the player should pay the same amount of fair movement cost.

By Lemma 10, we know that \mathcal{S}^T is a fair game. In fact, if we exempt the switching cost, we can treat $C_{mv}^F(\mathbb{O})$ as the summation of profits and $C_{mv}(\mathbb{O})$ as the movement cost of a series of fair game \mathcal{S}^T for the strategy \mathbb{O} . Because of the fairness, one has $\mathbb{E}[C_{mv}^F(\mathbb{O})] \leq \mathbb{E}[C_{mv}(\mathbb{O})]$ and further

$$\mathbb{E}[C_{sw}(\mathbb{O}) + C_{mv}^F(\mathbb{O})] \leq \mathbb{E}[C_{tot}(\mathbb{O})]. \quad (11)$$

Then by the optimality of \mathbb{O}_{fair} , we know

$$\mathbb{E}[C_{sw}(\mathbb{O}_{fair}) + C_{mv}^F(\mathbb{O}_{fair})] \leq \mathbb{E}[C_{sw}(\mathbb{O}) + C_{mv}^F(\mathbb{O})]. \quad (12)$$

Combining Equation 11 and Equation 12, we complete the proof of the claim. \square

C.4.2 Proof of Claim 24

Claim 24. *With the notations defined above, it holds that*

$$\mathbb{E}[C_{ktsp}(\mathbb{O}_{ktsp})] = \mathbb{E}[C_{sw}(\mathbb{O}_{ktsp}) + C_{sl}(\mathbb{O}_{ktsp})] \leq \mathbb{E}[C_{tot}(\mathbb{O})].$$

Recall that we have defined the game MG-Metric-Fair and $\mathbb{E}[C_{mv}^F(\mathbb{P}) + C_{sw}(\mathbb{P})]$ is the expected fair movement cost plus the expected switching cost of any strategy \mathbb{P} under the rule of MG-Metric-Fair. Now we let $\mathbb{E}[C_{sl}(\mathbb{P})]$ and $\mathbb{E}[C_{sw}(\mathbb{P})]$ be the expected selected cost and expected switching cost of any strategy \mathbb{P} under the rule of Stochastic- k -TSP respectively. Let \mathbb{O}_{ktsp} be the optimal strategy to Stochastic- k -TSP. Note that we use the same notation $C_{sw}(\cdot)$ to represent the expected switching cost as these three games have the same rules about switching costs.

Given any instance \mathcal{M} , let \mathbb{O} and \mathbb{O}_{fair} be the optimal strategies for this instance \mathcal{M} under the rule of MG-Metric and MG-Metric-Fair respectively, and let \mathcal{M}_{SC} be the Stochastic- k -TSP instance transferred from \mathcal{M} by the above argument. Now we show how to prove claim 24.

Proof. In fact, with the notation defined above, one has

$$\mathbb{E}[C_{sw}(\mathbb{O}_{ktsp}) + C_{sl}(\mathbb{O}_{ktsp})] \leq \mathbb{E}[C_{sw}(\mathbb{O}_{fair}) + C_{mv}^F(\mathbb{O}_{fair})] \leq \mathbb{E}[C_{tot}(\mathbb{O})].$$

We only need to prove the first inequality as the second equality has been proved in Claim 23.

To see the first inequality, we can look at Stochastic- k -TSP from another perspective. Suppose we know the decision tree of \mathbb{O}_{fair} , and construct another strategy ALG^F based on the decision tree under a special rule equivalent to Stochastic- k -TSP: whenever \mathbb{O}_{fair} plays some system \mathcal{S}_i for the first time, ALG^F can observe the whole sample path on \mathcal{S}_i and know the largest Index of states on the sample path, and it

can take one unit reward from \mathcal{S}_i later by paying the corresponding largest Index. Then we know that $\mathbb{E}[C_{sw}(\text{ALG}^F) + C_{sl}(\text{ALG}^F)] = \mathbb{E}[C_{sw}(\text{ALG}^F) + C_{mo}^F(\text{ALG}^F)]$, as the distribution on fair movement cost of ALG^F is exactly the distribution on selected cost if we play Stochastic- k -TSP by ALG^F .

Besides, we know $\mathbb{E}[C_{sw}(\mathbb{O}_{ktsp}) + C_{sl}(\mathbb{O}_{ktsp})] \leq \mathbb{E}[C_{sw}(\text{ALG}^F) + C_{sl}(\text{ALG}^F)]$ by the optimality, which completes the proof of this claim. \square

D Non-optimality of Indexing Strategies

In this section, we extend the result in [BS94] to show that even when the switching cost is a constant, there is no optimal indexing strategy for our problem in general metric. More specifically, we define what is an index as follows:

Definition 38 ([BS94]). *An index in the presence of switching cost is any function γ which specifies a value $\gamma(s_i, \mathcal{S}_i, I_i)$ for any Markov system \mathcal{S}_i where s_i is the current state of Markov system \mathcal{S}_i and I_i is the indicator that \mathcal{S}_i is currently active.*

Involving the indicator is to capture the intuition that for two identical systems, the active one should be more attractive than the inactive one. See more discussion in [BS94].

We say γ is an optimal index in the presence of switching cost if it is optimal by always playing the system with the smallest index until any one reaches its target.

Now we define two kinds of Markov system with specified notations. First let $[x\delta_1 + (1-x)\delta_0] = \{V, P, C, s, t\}$ where $V = \{s, t, v_0, v_1\}$, $P_{s,v_1} = x$, $P_{s,v_0} = 1-x$, and $P_{x_i,t} = 1$ for $i \in \{0, 1\}$. As for the movement cost, we let $C_s = c$, $C_{v_0} = 0$ and $C_{v_1} = 1$. Second, let $[\delta_a] = \{\{s, t\}, P, C, s, t\}$ where $P_{s,t} = 1$ and $C_{s,t} = a$. In the following, we assume the switching cost is c , same as C_s .

In this notation, we use $\gamma([x\delta_1 + (1-x)\delta_0]; I)$ and $\gamma(\delta_a; I)$ to denote the value of the optimal index on the Markov system $[x\delta_1 + (1-x)\delta_0]$ and $[\delta_a]$ respectively. Recall that $I = 1$ denotes that the system is currently active.

Claim 39 (Similar to Claim 1 in [BS94]). *Any index γ that is optimal in the presence of switching cost c must be a strict monotone transformation of an index $\hat{\gamma}$ which satisfies*

$$\hat{\gamma}([\delta_a]; 1) = a, \tag{13}$$

$$\hat{\gamma}([\delta_a]; 0) = a + c. \tag{14}$$

for any values of a and c .

Without loss generality, we assume the optimal index γ satisfies equation 13 and equation 14. Now we use the following claim to establish the impossibility of an optimal index:

Claim 40. *There is no consistent way to define an index γ on Markov system of the form $[x\delta_1 + (1-x)\delta_0]$ if the resulting strategy is to be invariably optimal. Consequently, an optimal index cannot exist.*

Proof. Consider such a game where there are two systems, $[\delta_a]$ and $[x\delta_1 + (1-x)\delta_0]$ and suppose the optimal decision-player is on the first system $[\delta_a]$ and is indifferent from playing either one of the first step:

he can either play $[\delta_a]$ or switch to $[x\delta_1 + (1-x)\delta_0]$ and switches back if only if the state of $[x\delta_1 + (1-x)\delta_0]$ moves to v_1 .

In particular, we have $a = 2c + x(a + c)$ and elementary calculation shows that the value of a should be $\frac{2+x}{1-x}c$ under the condition that $a + c \leq 1$ in this scenario. Let $\mu(x, c) := \frac{2+x}{1-x}c$ and if $\mu(x, c) + c \leq 1$, we have

$$\gamma([x\delta_1 + (1-x)\delta_0]; 0) = \gamma([\delta_{\mu(x,c)}]; 1) = \mu(x, c) = \frac{2+x}{1-x}c. \quad (15)$$

Now consider a new game where suppose the optimal decision-player is on the second system and can either: switches to the first system directly; or plays $[x\delta_1 + (1-x)\delta_0]$ for one step and switches to $[\delta_a]$ if the state moves to v_1 .

Similarly, by simple calculation, let $v(x, c) := \frac{xc}{1-x}$ and under the condition that $v(x, c) + c \leq 1$ we have

$$\gamma([x\delta_1 + (1-x)\delta_0]; 1) = \gamma([\delta_{v(x,c)}]; 0) = v(x, c) + c = c/(1-x). \quad (16)$$

Hence we show that $\gamma([x\delta_1 + (1-x)\delta_0]; 0) = \frac{2+x}{1-x}c$ if $3c \leq 1-x$, and $\gamma([x\delta_1 + (1-x)\delta_0]; 1) = c/(1-x)$ if $c \leq 1-x$ by the above two games.

Now consider the last game to show contradiction. Suppose there are two systems $[x\delta_1 + (1-x)\delta_0]$ and $[y\delta_1 + (1-y)\delta_0]$ where $x \geq y$ but the player is current at system $[x\delta_1 + (1-x)\delta_0]$. Suppose $3c \leq 1-x \leq 1-y$. If we have

$$\min\{c+x, c+x(2c+y)\} \leq \min\{2c+y, 2c+y(2c+x)\},$$

then the optimal strategy can play the first system for the first step and hence the constraint is $2x \leq 1+2y$ and we have

$$\gamma([x\delta_1 + (1-x)\delta_0]; 1) \leq \gamma([y\delta_1 + (1-y)\delta_0]; 0).$$

To conclude, the constraints are $y \leq x$, $3c+x \leq 1$ and $2x \leq 1+2y$.

By setting $x = 4/5$, $y = 2/5$ and $c = 1/100$, all the constraints are satisfied and one has $\gamma([x\delta_1 + (1-x)\delta_0]; 1) = \frac{1}{1-x}c = 1/20 > \gamma([y\delta_1 + (1-y)\delta_0]; 0) = \frac{2+y}{1-y}c = 1/25$, which is a contradiction completing the proof. \square

References

- [ADT12] Raman Arora, Ofer Dekel, and Ambuj Tewari. Online bandit learning against an adaptive adversary: from regret to policy regret. In *Proceedings of the 29th International Conference on International Conference on Machine Learning*, pages 1747–1754, 2012.
- [ASKW04] Rina Azoulay-Schwartz, Sarit Kraus, and Jonathan Wilkenfeld. Exploitation vs. exploration: choosing a supplier in an environment of incomplete information. *Decision support systems*, 38(1):1–18, 2004.
- [AT96] Manjari Asawa and Demosthenis Teneketzis. Multi-armed bandits with switching penalties. *IEEE transactions on automatic control*, 41(3):328–348, 1996.
- [BB88] Ll Benkherouf and JA Bather. Oil exploration: sequential decisions in the face of uncertainty. *Journal of Applied Probability*, pages 529–543, 1988.
- [BGK11] Anand Bhalgat, Ashish Goel, and Sanjeev Khanna. Improved approximation results for stochastic knapsack problems. In *SODA*, pages 1647–1665, 2011.
- [BGO92] L Benkherouf, KD Glazebrook, and RW Owen. Gittins indices and oil exploration. *Journal of the Royal Statistical Society: Series B (Methodological)*, 54(1):229–241, 1992.
- [BK19] Hedyeh Beyhaghi and Robert Kleinberg. Pandora’s problem with nonobligatory inspection. In *Proceedings of the 2019 ACM Conference on Economics and Computation*, pages 131–132, 2019.
- [BN15] Nikhil Bansal and Viswanath Nagarajan. On the adaptivity gap of stochastic orienteering. *Mathematical Programming*, 154(1):145–172, 2015.
- [BS94] Jeffrey S Banks and Rangarajan K Sundaram. Switching costs and the gittins index. *Econometrica: Journal of the Econometric Society*, pages 687–694, 1994.
- [BV01] Dirk Bergemann and Juuso Välimäki. Stationary multi-choice bandit problems. *Journal of Economic dynamics and Control*, 25(10):1585–1594, 2001.
- [BV06] Thomas Brenner and Nicolaas J Vriend. On the behavior of proposers in ultimatum games. *Journal of Economic Behavior & Organization*, 61(4):617–631, 2006.
- [CGT⁺20] Shuchi Chawla, Evangelia Gergatsouli, Yifeng Teng, Christos Tzamos, and Ruimin Zhang. Pandora’s box with correlations: Learning and approximation. In *2020 IEEE 61st Annual Symposium on Foundations of Computer Science (FOCS)*, pages 1214–1225. IEEE, 2020.
- [DDKP14] Ofer Dekel, Jian Ding, Tomer Koren, and Yuval Peres. Bandits with switching costs: T 2/3 regret. In *Proceedings of the forty-sixth annual ACM symposium on Theory of computing*, pages 459–467, 2014.
- [DGV04] Brian C. Dean, Michel X. Goemans, and Jan Vondrák. Approximating the stochastic knapsack problem: The benefit of adaptivity. In *Foundations of Computer Science, 2004. Proceedings. 45th Annual IEEE Symposium on*, pages 208–217. IEEE, 2004.
- [DGV08] Brian C Dean, Michel X Goemans, and Jan Vondrák. Approximating the stochastic knapsack problem: The benefit of adaptivity. *Mathematics of Operations Research*, 33(4):945–964, 2008.

- [DH03] Fabrice Dusonchet and M-O Hongler. Optimal hysteresis for a class of deterministic deteriorating two-armed bandit problem with switching costs. *Automatica*, 39(11):1947–1955, 2003.
- [Dov18] Laura Doval. Whether or not to open pandora’s box. *Journal of Economic Theory*, 175:127–158, 2018.
- [DTW03] Ioana Dumitriu, Prasad Tetali, and Peter Winkler. On playing golf with two balls. *SIAM Journal on Discrete Mathematics*, 16(4):604–615, 2003.
- [ENS18] Alina Ene, Viswanath Nagarajan, and Rishi Saket. Approximation algorithms for stochastic k-tsp. In *37th IARCS Annual Conference on Foundations of Software Technology and Theoretical Computer Science*, 2018.
- [Fre75] David A Freedman. On tail probabilities for martingales. *the Annals of Probability*, 3(1):100–118, 1975.
- [FW99] Esther Frostig and Gideon Weiss. Four proofs of gittins’ multiarmed bandit theorem. *Applied Probability Trust*, 70:427, 1999.
- [GGW11] John Gittins, Kevin Glazebrook, and Richard Weber. *Multi-armed bandit allocation indices*. John Wiley & Sons, 2011.
- [Git74] John Gittins. A dynamic allocation index for the sequential design of experiments. *Progress in statistics*, pages 241–266, 1974.
- [Git79] John C Gittins. Bandit processes and dynamic allocation indices. *Journal of the Royal Statistical Society: Series B (Methodological)*, 41(2):148–164, 1979.
- [GJSS19] Anupam Gupta, Haotian Jiang, Ziv Scully, and Sahil Singla. The markovian price of information. In *International Conference on Integer Programming and Combinatorial Optimization*, pages 233–246. Springer, 2019.
- [GKNR12] Anupam Gupta, Ravishankar Krishnaswamy, Viswanath Nagarajan, and R Ravi. Approximation algorithms for stochastic orienteering. In *Proceedings of the twenty-third annual ACM-SIAM symposium on Discrete Algorithms*, pages 1522–1538. SIAM, 2012.
- [GM09] Sudipto Guha and Kamesh Munagala. Multi-armed bandits with metric switching costs. In *International Colloquium on Automata, Languages, and Programming*, pages 496–507. Springer, 2009.
- [JLLS20] Haotian Jiang, Jian Li, Daogao Liu, and Sahil Singla. Algorithms and adaptivity gaps for stochastic k-tsp. In *11th Innovations in Theoretical Computer Science Conference, ITCS 2020*, page 45. Schloss Dagstuhl-Leibniz-Zentrum fur Informatik GmbH, Dagstuhl Publishing, 2020.
- [Joh78] William R Johnson. A theory of job shopping. *The Quarterly Journal of Economics*, pages 261–278, 1978.
- [Jov84] Boyan Jovanovic. Matching, turnover, and unemployment. *Journal of political Economy*, 92(1):108–122, 1984.

- [JSV04] Mark Jerrum, Alistair Sinclair, and Eric Vigoda. A polynomial-time approximation algorithm for the permanent of a matrix with nonnegative entries. *Journal of the ACM (JACM)*, 51(4):671–697, 2004.
- [Jun01] Tack-seung Jun. *Essays on decision theory: effects of changes in environment on decisions*. Columbia University, 2001.
- [Jun04] Tackseung Jun. A survey on the bandit problem with switching costs. *de Economist*, 152(4):513–541, 2004.
- [KL00] Stylianos K Kavadias and Christoph H Loch. *A dynamic resource allocation policy in multiproject environments*. Number 2000-2010. INSEAD, 2000.
- [Kli04] Mikhail M Klimenko. Industrial targeting, experimentation and long-run specialization. *Journal of Development Economics*, 73(1):75–105, 2004.
- [KLM17] Tomer Koren, Roi Livni, and Yishay Mansour. Multi-armed bandits with metric movement costs. In *Advances in Neural Information Processing Systems*, pages 4119–4128, 2017.
- [Krä03] Daniel Krähmer. Entry and experimentation in oligopolistic markets for experience goods. *International Journal of Industrial Organization*, 21(8):1201–1213, 2003.
- [KSU08] Robert Kleinberg, Aleksandrs Slivkins, and Eli Upfal. Multi-armed bandits in metric spaces. In *Proceedings of the fortieth annual ACM symposium on Theory of computing*, pages 681–690, 2008.
- [KW11] John Kennan and James R Walker. The effect of expected income on individual migration decisions. *Econometrica*, 79(1):211–251, 2011.
- [KWW16] Robert Kleinberg, Bo Waggoner, and E Glen Weyl. Descending price optimally coordinates search. In *Proceedings of the 2016 ACM Conference on Economics and Computation*, pages 23–24, 2016.
- [LL16] Jian Li and Yu Liu. Approximation algorithms for stochastic combinatorial optimization problems. *Journal of the Operations Research Society of China*, 4(1):1–47, 2016.
- [LY13] Jian Li and Wen Yuan. Stochastic combinatorial optimization via poisson approximation. In *Proceedings of the forty-fifth annual ACM symposium on Theory of computing*, pages 971–980, 2013.
- [Mac80] Glenn M MacDonald. Person-specific information in the labor market. *Journal of Political Economy*, 88(3):578–597, 1980.
- [McL84] Andrew McLennan. Price dispersion and incomplete learning in the long run. *Journal of Economic dynamics and control*, 7(3):331–347, 1984.
- [Ort08] Ronald Ortner. Online regret bounds for markov decision processes with deterministic transitions. In *International Conference on Algorithmic Learning Theory*, pages 123–137. Springer, 2008.

- [PT94] Christos H Papadimitriou and John N Tsitsiklis. The complexity of optimal queueing network control. In *Proceedings of IEEE 9th Annual Conference on Structure in Complexity Theory*, pages 318–322. IEEE, 1994.
- [PT95] DG Pandalis and D Teneketzis. On the optimality of the gittins index rule in multi-armed bandits with multiple plays. In *Proceedings of 1995 34th IEEE Conference on Decision and Control*, volume 2, pages 1408–1414. IEEE, 1995.
- [Rot74] Michael Rothschild. A two-armed bandit theory of market pricing. *Journal of Economic Theory*, 9(2):185–202, 1974.
- [Sch97] KH Schlag. Why imitate. and if so, how? a bounded rationality approach to multiarmed bandits. *Journal of Economic Theory*, 78:127–159, 1997.
- [Sin18] Sahil Singla. The price of information in combinatorial optimization. In *Proceedings of the Twenty-Ninth Annual ACM-SIAM Symposium on Discrete Algorithms*. SIAM, 2018.
- [Smi99] Lones Smith. Optimal job search in a changing world. *Mathematical Social Sciences*, 38(1):1–9, 1999.
- [Vis80] W Kip Viscusi. A theory of job shopping: A bayesian perspective. *The Quarterly Journal of Economics*, 94(3):609–614, 1980.
- [VOP00] Mark P Van Oyen and Jutta Pichitlamken. Properties of optimal-weighted flowtime policies with a makespan constraint and set-up times. *Manufacturing & Service Operations Management*, 2(1):84–99, 2000.
- [W⁺92] Richard Weber et al. On the gittins index for multiarmed bandits. *The Annals of Applied Probability*, 2(4):1024–1033, 1992.
- [Wal84] Michael Waldman. Job assignments, signalling, and efficiency. *The RAND Journal of Economics*, 15(2):255–267, 1984.
- [Wei79] Martin L. Weitzman. Optimal search for the best alternative. *Econometrica: Journal of the Econometric Society*, pages 641–654, 1979.
- [Whi88] Peter Whittle. Restless bandits: Activity allocation in a changing world. *Journal of applied probability*, pages 287–298, 1988.